

1

2 **Projections of daily mean temperature variability in**  
3 **the future: cross-validation tests with ENSEMBLES**  
4 **regional climate simulations**

5

6 Jouni Räisänen and Olle Räty

7 Department of Physics, University of Helsinki, Finland

8

9 Submitted to Climate Dynamics, 14 June 2012

10 Returned for minor revision, 22 July 2012

11 Revised, 23 August 2012

12 Accepted, 27 August 2012

13

14

15

16

17

18 **Corresponding author**

19 Jouni Räisänen

20 Department of Physics, P.O. Box 48 (Erik Palménin aukio 1),

21 FI-00014 University of Helsinki, Finland

22 Phone +358-9-191 50872; Fax + 358-9-191 48802

23 Email: [jouni.raisanen@helsinki.fi](mailto:jouni.raisanen@helsinki.fi)

24 **Abstract**

25 Because of model biases, projections of future climate need to combine model  
26 simulations of recent and future climate with information on observed climate.  
27 Here, ten methods for projecting the distribution of daily mean temperatures are  
28 compared, using six regional climate change simulations for Europe. Cross  
29 validation between the models is used to assess the potential performance of the  
30 methods in projecting future climate. Delta change and bias correction type  
31 methods show similar cross-validation performance, with methods based on the  
32 quantile mapping approach doing best in both groups due to their apparent ability  
33 to reduce the errors in the projected time mean temperature change. However, as  
34 no single method performs best under all circumstances, the optimal approach  
35 might be to use several well-behaving methods in parallel. When applying the  
36 various methods to real-world temperature projection for the late 21st century, the  
37 largest intermethod differences are found in the tails of the temperature  
38 distribution. Although the intermethod variation of the projections is generally  
39 smaller than their intermodel variation, it is not negligible. Therefore, it should be  
40 preferably included in uncertainty analysis of temperature projections, particularly  
41 in applications where the extremes of the distribution are important.

42

43 **KEYWORDS:** climate change, climate projection, temperature, daily variability,  
44 delta change, bias correction, cross validation, ENSEMBLES, Europe

45

46

## 47 **1. Introduction**

48 Despite decades of development, global and regional climate models (GCMs and  
49 RCMs) still show various kinds of biases in the simulation of the present-day  
50 climate (Randall et al. 2007, Christensen et al. 2007, van der Lindell and Mitchell  
51 2009). Therefore, model-simulated future climate as such rarely provides a  
52 plausible projection of the actual future climate. To alleviate the impact of model  
53 biases, construction of climate projections also needs to extract information from  
54 the observed and simulated climates in the recent past.

55

56 As an example, three 30-year (930-day) time series of January daily mean  
57 temperature in Jyväskylä, central Finland, are shown in Fig. 1: the first from  
58 station observations in 1971-2000, the second from an RCM simulation during the  
59 same period, and the third from the same RCM in the end of this century (2069-  
60 2098). During the years 1971-2000, the RCM simulation exhibits both a cold bias  
61 and smaller than observed variability, and the distribution of the simulated  
62 temperatures shows less negative skewness than that observed. Considering these  
63 deficiencies, the simulation for 2069-2098 is unlikely to provide a good  
64 description of the climate in this period, even if the simulated changes in mean  
65 temperature and characteristics of variability turned out to be correct.

66

67 Two commonly used approaches to account for model biases are “delta change”  
68 and “bias correction” (Fig. 2). In the former, the projection for the future is  
69 obtained by perturbing an observed time series based on the difference between  
70 the simulated future and baseline climates. In the latter, the projection is built on  
71 the future simulation by the model, after correcting this based on the differences  
72 between the simulated and observed climate during the baseline period.

73

74 If only long-term climatic means are needed, the problem is technically simple.  
75 For example, a projection for the future mean temperature is easily derived as

$$76 \quad \bar{p} = \bar{o} + (\bar{s} - \bar{c}) = \bar{s} - (\bar{c} - \bar{o}) \quad (1)$$

77 where the overline indicates temporal averaging and the four letters stand for  
78 projection ( $p$ ), baseline observations ( $o$ ), scenario simulation for the future period

79 of interest ( $s$ ) and control simulation for the baseline period ( $c$ ). In this case, the  
80 delta change and the bias correction approaches (the first and the second form in  
81 (1), respectively) give identical results. Note, however, that this result is neither  
82 unique nor necessarily optimal. As shown in recent studies (Buser et al. 2009;  
83 Boberg and Christensen 2012) and later in this paper, biases in simulated  
84 variability may also have implications for projections of the time mean climate.

85

86 If characteristics of variability are essential, the situation is more problematic. A  
87 constant delta change or constant bias correction would fail to account for either  
88 the changes or biases in the amplitude, shape of distribution, and temporal  
89 structure of the simulated variability. This can be improved by more sophisticated  
90 projection methods, but not without a potential trade-off. The more precisely a  
91 projection scheme attempts to correct for differences between simulated and  
92 observed climate or to incorporate simulated climate change, the more likely it is  
93 affected by features that are not statistically robust (e.g., random fluctuations in  
94 the tails of the distribution). The potential advantages of more sophisticated  
95 projection methods also need to be put in the context of the model- and scenario-  
96 related uncertainty in future climate change (Meehl et al. 2007, Christensen et al.  
97 2007).

98

99 A large array of methods for projecting future climate variability on daily scales  
100 have been developed particularly for precipitation (see Maraun et al. 2010, for a  
101 recent review) but also for temperature (e.g., Engen-Skaugen 2007, Piani et al.  
102 2010, Amengual et al 2012). The question thus arises, which of these different  
103 alternatives should be preferred? Although projection methods can be compared  
104 for their ability to reproduce present-day climate statistics (e.g. Themeßl et al.  
105 2011, Dosiolo and Paruolo 2011), the crucial issue is their performance in future  
106 climate.

107

108 Although future climate is unknown, some inferences on the potential  
109 performance of different methods can be drawn from intermodel cross validation  
110 (Fig. 3). If simulations for both the recent past and the future are available for  $N$   
111 models, any one of these can be left out for verification. In the same way as  
112 projections for the real future climate would be made by combining observations

113 with the baseline and future simulations from different climate models,  
114 projections for the future climate in the verifying model can be derived by  
115 replacing observations with the baseline climate in this model. Unlike in the real  
116 world, this projection is verifiable against the actual future climate in the same  
117 model. Repeating this over all choices of the verifying model, statistics can be  
118 gathered that allow comparison between different methods of projection. Such a  
119 technique has already been used in studies focusing on projection of time mean  
120 climate (e.g. Räisänen and Ylhäisi 2011, Bracegirdle and Stephenson 2012,  
121 Maraun 2012), and it is also planned to serve as one of the main tools in the  
122 recently started European Concerted Research Action ES1102 VALUE  
123 (Validating and Integrating Downscaling Methods for Climate Change Research).

124

125 In the present study, which was in part inspired by VALUE, the focus is on the  
126 projection of daily mean temperatures. Ten different projection methods, broadly  
127 similar to those used in earlier studies, are applied to a subset of six RCM  
128 simulations for Europe from the ENSEMBLES (ENSEMBLE-based Predictions  
129 of Climate Changes and their Impacts) project (van der Linden and Mitchell  
130 2009). Two main issues will be studied:

131

- 132 1. Which of the tested methods show the most promise for projection of future  
133 climate in light of their cross-validation performance? Can a single best  
134 method be identified, or would it be better to use several methods in parallel,  
135 to take into account the uncertainty in this choice (cf. Ho et al. 2012)?
- 136 2. How large is the uncertainty associated with the choice of the projection  
137 method compared with the variation of climate change between different  
138 models?

139

140 The model data and observations used are described in Section 2, and the  
141 projection methods are introduced in Section 3. The cross-validation results are  
142 presented in Section 4, whereas Section 5 studies the sensitivity of the projected  
143 future climate to the choice of the projection method. Finally, a synthesis of the  
144 main conclusions is presented in Section 6.

145

## 146 **2. Data sets**

147 Six RCM simulations from the ENSEMBLES project are used (Table 1), all using  
148 a different RCM and different driving GCM but the same (SRES A1B) emissions  
149 scenario. These data were retrieved from the ENSEMBLES Research Theme 3  
150 web page ([ensemblesrt3.dmi.dk/](http://ensemblesrt3.dmi.dk/)) in a regular  $0.25^\circ$  lon  $\times$   $0.25^\circ$  lat grid covering  
151 Europe and northernmost Africa. However, to reduce the computations, only 211  
152 land grid boxes with  $2.5^\circ \times 2.5^\circ$  spacing were used in cross validation (Fig. 4a).  
153 Here we mainly use data for a 30-year baseline period (1971-2000) and a 30-year  
154 period in the end of this century (2069-2098), but some results for an earlier  
155 projection period (2001-2030) are also shown.

156

157 For testing real-world temperature projection, a total of 139 station time series  
158 were selected from the blended European Climate Assessment & Data archive  
159 available through the Royal Netherlands Meteorological Institute Climate  
160 Explorer ([climexp.knmi.nl/](http://climexp.knmi.nl/)), one per each  $2.5^\circ$  lat  $\times$   $2.5^\circ$  lon box where a station  
161 with near-complete time series (at least 99% of valid data in 1971-2000) was  
162 available (Fig. 4b). For each station, simulated time series were chosen from that  
163 one of the nearest nine  $0.25^\circ \times 0.25^\circ$  grid boxes which has the largest land  
164 fraction, using for consistency the same land-sea mask (from SMHI-BCM) in all  
165 cases. This procedure was adopted to avoid cases in which observations from a  
166 coastal but nevertheless land-based station would be combined with simulated  
167 time series from a sea-dominated grid box.

## 168 **3. Projection methods**

169 Ten different methods for constructing projections of future temperature  
170 variability are studied (Table 2). The delta change (M1-M5) and bias correction  
171 methods (M6-M10) are technically symmetric: the same computer subroutines  
172 can be used for both by only switching the order of the three input time series  
173 (observations and the baseline and future simulations) in the argument list.  
174 Therefore, only the delta change methods are described in the following.

175

176 M1 simply adds the time mean temperature change between the baseline and  
177 scenario simulations to each daily value in the observed time series. M2 also takes

178 into account changes in the standard deviation, converting the values in the  
 179 observed time series ( $o_i$ ) to projected values ( $p_i$ ) as

$$180 \quad p_i = \bar{o} + (\bar{s} - \bar{c}) + (o_i - \bar{o}) \frac{s_s}{s_c} \quad (2)$$

181 where  $s_s$  and  $s_c$  are the simulated standard deviations during the scenario and  
 182 baseline periods. In M3, changes in the sample third-moment skewness are also  
 183 included, so that the skewness of the projected time series becomes

$$184 \quad skew_p = skew_o + (skew_s - skew_c) \quad (3)$$

185 where the subscripts  $p$ ,  $o$ ,  $s$  and  $c$  refer to the projected, observed, control  
 186 simulation and scenario simulation time series, respectively. The condition (3)  
 187 could be fulfilled by several different modifications to the data. Here, we follow  
 188 the algorithm described by Ballester et al. (2010) in their electronic supplementary  
 189 material.

190

191 M4 and M5 use the quantile mapping approach. Cumulative probability  
 192 distributions of temperature are first estimated for both the control ( $F_c$ ) and the  
 193 scenario simulations ( $F_s$ ), and each observed value  $o_i$  is then converted to

$$194 \quad p_i = F_s^{-1}(F_c(o_i)) \quad (4)$$

195 In implementing (4), two practical issues need to be solved. First, if the  
 196 conversion  $F_s^{-1}(F_c)$  is derived directly from an empirical quantile-quantile plot, it  
 197 tends to become noisy near the tails of the distribution (see the crosses in Fig. 5a).  
 198 To avoid this, some smoothing is needed. Second, if some of the observed values  
 199 fall out of the range in the control simulation, the quantile-quantile relationship  
 200 needs to be extrapolated beyond the simulated range.

201

202 M4 and M5 differ in how these practical issues are solved. In M4, the quantiles in  
 203 the model simulations are smoothed using a running average, replacing the  
 204 quantiles  $F_c^{-1}(x)$ ,  $x \in [0,1]$  with

$$205 \quad \tilde{F}_c^{-1}(x) = \frac{\int_{\max(x-D,0)}^{\min(x+D,1)} F_c^{-1}(x) dx}{\int_{\max(x-D,0)}^{\min(x+D,1)} dx} \quad (5)$$

206 and similarly for  $F_s^{-1}(x)$ . Here we use  $D = 0.05$ , which was found to be close to  
207 optimal in terms of the cross-validation statistics. This smoothed quantile-quantile  
208 relationship is illustrated by the bold red line in Fig. 5a. A disadvantage of the  
209 smoothing is that it narrows the range of the data; for example  $\tilde{F}_c^{-1}(0)$  equals the  
210 mean of the lowest 5% of the temperatures in the control simulation. We  
211 extrapolate towards low and high values assuming that the difference  $\tilde{F}_s^{-1} -$   
212  $\tilde{F}_c^{-1}$  remains constant for  $x < 0$  and  $x > 1$  (dashed red lines in Fig. 5a).

213

214 In M5, simple linear regression is used to map  $F_c^{-1}$  on  $F_s^{-1}$  (blue line in Fig. 5a).  
215 This coarse-grained implementation of quantile mapping is used to study whether  
216 the more detailed treatment in M4 has additional value.

217

218 The projected scenario period (2069-2098) time series for the case introduced in  
219 Fig. 1 are shown in the left part of Fig. 6 (also note the statistics included in the  
220 figure panels). They illustrate the following key features:

221

- 222 1. In the delta change methods (M1-M5), the structure of the projection time  
223 series follows the observations, in the bias correction methods (M6-M10)  
224 the simulation for the scenario period.
- 225 2. The projected time mean temperature is the same for all of M1-M3 and  
226 M6-M8, but not for the quantile mapping methods M4-M5 and M9-M10.  
227 The causes of this difference will be discussed later in this section.
- 228 3. M2-M3 and M7-M8 all produce the same standard deviation of  
229 temperatures.
- 230 4. M3 and M8 additionally yield the same skewness of the distribution.
- 231 5. The skewness produced by M1-M2 and M5 (-0.9) is the same as in the  
232 observed time series in Fig. 1: none of these methods modifies the shape  
233 of the temperature distribution aside from its mean and standard deviation.  
234 Similarly, the projections based on M6-M7 and M10 have the same  
235 skewness (-0.5) as the RCM simulation for 2069-2098 in Fig. 1.
- 236 6. The extremes of the temperature distribution are particularly sensitive to  
237 the choice of the method. The maximum of the projected time series varies



238 from 5.5°C to 11.6°C, the minimum from -36.8°C to -23.2°C. For M8, the  
239 minimum is actually lower than that observed in 1971-2000 (-36.1°C).

240

241 As demonstrated by the quantile representation in the right part of Fig. 6, the  
242 selection of the method is not the only uncertainty in the projection. The choice of  
243 the RCM simulation also matters, although the size of the inter-RCM variation  
244 depends on the projection method used. We will study the relative roles of inter-  
245 method and inter-RCM uncertainty in more depth in Section 5.

246

247 M4-M5 give in this case a slightly lower mean temperature than the other delta  
248 change methods, which retain the change in mean temperature exactly as  
249 simulated by the model. This is explained by the cold bias in the simulation in  
250 1971-2000 together with the decrease in variability occurring in 2069-2098 (Figs.  
251 1 and 5a). Because the observed temperatures in 1971-2000 mostly fall in the  
252 upper part of the simulated distribution, where the quantile-quantile comparison  
253 indicates a smaller difference between 1971-2000 and 2069-2098, the mean  
254 temperature change as weighted by the distribution of observations becomes  
255 smaller than that directly simulated by the model. Conversely, M9-M10 indicate a  
256 larger increase in the mean temperature than the other methods. This is due to the  
257 underestimate in variability in the simulation for 1971-2000, which in bias  
258 correction type quantile mapping implies that a larger correction should be added  
259 to higher temperatures (Fig. 5b). As the simulated temperature distribution shifts  
260 upward from 1971-2000 to 2069-2098, the average correction in the latter period  
261 becomes larger than that in the former, thus amplifying the projected change in  
262 time mean temperature. Although the details of these results are case specific, the  
263 ability of quantile mapping to modify the time mean temperature change  
264 represents a generic difference from the other six methods.

265

266 In Figs. 5-6, we used (for simplicity of interpretation) only January data when  
267 estimating the changes (M1-M5) and biases (M6-M10) of the simulated January  
268 temperature distribution. This is not necessarily optimal, because the resulting  
269 relatively small sample size may introduce substantial noise. The noise can be  
270 reduced by using a wider time window in the estimation of climate changes and  
271 model biases, although potentially at the cost of some systematic error. In Section

272 4 below, we test three choices of the window length: one, two and three months.  
 273 For the two-month window, for example, data from the second half of December  
 274 to the first half of February are used in addition to January data when estimating  
 275 the changes in January in M1-M5 and biases in M6-M10.

## 276 4. Cross validation

277 Should all of the ten methods be regarded as equally plausible, or are some of  
 278 them more likely to give useful temperature projections than others? Here, we  
 279 study this using cross validation between the six model simulations, as shown  
 280 schematically in Fig. 3. One deterministic and two probabilistic statistics are  
 281 computed, all based on a comparison between the quantiles of the projected and  
 282 verifying temperature distributions ( $T_{proj}(x)$  and  $T_{ver}(x)$ ,  $x = 0 \dots 100\%$ ) in a given  
 283 month and location. The mean square error is

$$284 \quad MSE = A[\langle T_{proj} \rangle - T_{ver}]^2 \quad (6)$$

285 where  $\langle \rangle$  denotes the ensemble mean of the five (six minus verifying model)  
 286 individual projections and  $A$  indicates averaging over the whole distribution from  
 287 0 to 100%, the 12 months, the 211 land grid boxes (weighted with the cosine of  
 288 latitude), and the six choices of the verifying model. For calculating the  
 289 continuous ranked probability score

$$290 \quad CRPS = A\left[ \int_{-\infty}^{T_{ver}} F(T_{proj})^2 dT_{proj} + \int_{T_{ver}}^{\infty} (1 - F(T_{proj}))^2 dT_{proj} \right] \quad (7)$$

291 we first form, separately for each quantile of temperature, a probabilistic forecast  
 292 for  $T_{ver}$  from the discrete cumulative distribution  $F(T_{proj})$  of the five  $T_{proj}$  values  
 293 (cf. Räisänen and Palmer 2001). Our third score, *OutOfRange*, records the  
 294 frequency of cases in which  $T_{ver}$  is below the lowest or above the highest of the  
 295 five  $T_{proj}$  values. Unlike *MSE* and *CRPS*, *OutOfRange* is not a proper validation  
 296 score in the sense that a lower value would always indicate a better forecast. By  
 297 inflating the forecast distribution sufficiently, one could ensure *OutOfRange* = 0,  
 298 while simultaneously making the forecast useless. A more useful interpretation is  
 299 as follows. If  $T_{ver}$  and the five  $T_{proj}$  values are independent samples from the same  
 300 statistical population (as they ideally should), then the probability that  $T_{proj}$  is the  
 301 lowest or highest of these six values is 1/3. If *OutOfRange* exceeds this value, this

302 indicates that the forecast obtained from the five  $T_{proj}$  values is underdispersive,  
303 thus underestimating the uncertainty in  $T_{ver}$ . Conversely,  $OutOfRange < 1/3$  would  
304 indicate an overdispersive forecast.

305

306 A summary of the cross-validation statistics is given in Fig. 7. Focusing first on  
307 the statistics for the projection period 2069-2098 (using 1971-2000 as the  
308 baseline), we can note the following:

309

310 1. *MSE* and *CRPS* are relatively insensitive to the number of months used in  
311 estimating the changes (M1-M5) or the model biases (M6-M10).  
312 However, in nearly all cases, two-month aggregation of data performs  
313 better than the use of a single month, although some exceptions are found  
314 for individual verifying models particularly for the delta change methods  
315 (not shown). Differences between two and three months are unsystematic  
316 even when considering the six-model mean statistics shown in Fig. 7.

317 2. *MSE* and *CRPS* give a similar picture of the relative performance of the  
318 ten methods (although the differences in *MSE* are larger). The two  
319 simplest methods, M1 (constant change over the whole distribution) and  
320 M6 (constant bias correction for the whole distribution) perform less well  
321 than the others. The inclusion of the standard deviation in M2 and M7  
322 gives a clear improvement, but there is little additional change in the  
323 statistics when also modifying the skewness (M3 and M8). Methods based  
324 on quantile mapping perform best within both the delta change group  
325 (M4-M5) and the bias correction group (M9-M10). M9 has both the  
326 lowest *MSE* and *CRPS* of the ten methods, although the difference from  
327 M10 is small.

328 3. *OutOfRange* is very close to the “desired” value of 1/3 (33.3%) for all of  
329 the bias correction methods. By contrast, the delta change methods tend to  
330 provide underdispersive projections, so that the verification falls more  
331 often out of the range of the five projections than it ideally should. This is  
332 understandable. Because the delta change projections from the five  
333 forecast models all modify the same underlying time series (the 1971-  
334 2000 time series in the verifying model), their differences do not fully

335 cover the effects of internal variability. This underdispersion becomes  
336 more pronounced for longer time windows in estimating the change.

337

338 The bottom row of Fig. 7 shows the corresponding verification statistics for the  
339 period 2001-2030, when climate changes from 1971-2000 are much smaller than  
340 in 2069-2098. The absolute intermodel differences in change are also smaller,  
341 making both *MSE* and *CRPS* lower during this period. Unlike in 2069-2098, M1  
342 shows in 2001-2030 nearly identical performance with the other delta change  
343 methods. At this time, changes in the width and shape of the temperature  
344 distribution still have a very low signal-to-noise ratio. Their inclusion in the  
345 projection has, therefore, little impact on the cross-validation performance. Most  
346 of the bias correction methods are also close in performance to the delta change  
347 methods at this time, M9 being again the best in the whole group. However, M6  
348 with its unrealistic assumption that model biases are constant throughout the  
349 distribution performs substantially worse than the other methods.

350

351 The *MSE* calculated for the whole temperature distribution can be written as the  
352 sum of two components: one arising from the error in the time mean temperature,  
353 and the other from errors in the differences between the individual quantiles and  
354 the mean value. The latter component (denoted as *debiased MSE* in Fig. 7) is  
355 typically much smaller than the former, particularly in 2069-2098. Therefore,  
356 most of the *MSE* and some of the intermethod differences in *MSE* actually reflect  
357 errors in the projected time mean temperature, rather than those in the width and  
358 shape of the distribution. In particular, the best performance of the quantile  
359 mapping methods (especially M9 and M10) in 2069-2098 results from the best  
360 projections of the time mean temperature. On the other hand, relatively large  
361 errors in the shape and the width of the distribution do distinguish the worst  
362 methods (M1 and M6 in 2069-2098 and M6 in 2001-2030) from the others.

363

364 In constructing Fig. 7, the same weight was given to all parts of the temperature  
365 distribution. Yet, for some applications the projection accuracy near the upper  
366 and/or lower tails of the distribution might be unproportionally important. To  
367 illustrate how the relative performance of the methods varies across the  
368 distribution, their *MSE* ranks in 2069-2098 are shown in Fig. 8 separately for all

369 percentiles of temperature (similar analysis for *CRPS* gives very similar results).  
370 For this and all the later figures in this paper, the two-month time window is used.

371

372 For most of the distribution, the ranking of the methods is broadly consistent with  
373 Fig. 7. As an exception, the simple constant bias correction M6 actually has the  
374 lowest *MSE* at 16-34%, although it performs very poorly in the upper part of the  
375 distribution. Otherwise, M9 and M10 with the lowest overall *MSEs* dominate the  
376 top rank from the extreme lower tail up to 98%. However, in the extreme high  
377 end, the performance of M9 and particularly M10 deteriorates. Thus, these  
378 projection methods might not be optimal for applications that are particularly  
379 sensitive to extremely high temperatures. Considering how M9 and M10 function,  
380 this deterioration is not surprising. Both methods attempt to estimate the  
381 temperature dependence of model bias from comparison between the observed  
382 and simulated distributions during the baseline period (Fig. 5b). However, the  
383 highest temperatures simulated in the late 21st century by far exceed those in the  
384 baseline period. This results in substantial extrapolation uncertainty in the bias  
385 correction, which probably explains the deteriorating performance of M9 and  
386 M10 in the upper end of the distribution.

387

388 The intermethod variation of cross-validation statistics also depends on the  
389 location, month of the year, and the verifying model. We do not discuss the details  
390 of this variation here, but emphasize the general implication: a method that is best  
391 in an average sense will not be the best in all individual cases. This indicates that  
392 the choice between different projection methods represents a genuine uncertainty  
393 that cannot be fully eliminated by selecting a single best method.

394

395 Given the uncertainty in the projection methods, might it not be better to use  
396 several methods simultaneously instead of just one? To test this suggestion, four  
397 combinations of methods were chosen. B2, B4 and B8 include the best two (M9  
398 and M10), four (M4-M5 and M9-M10) and eight methods (M2-M5 and M7-M10)  
399 in terms of the overall *MSE* and *CRPS* statistics for 2069-2098, while A10  
400 includes all ten methods. In each case, the same weight was given to all of the  
401 methods included. Thus, for example, the best estimate projection from B8  
402 becomes the mean of the multi-model means from all methods excluding M1 and

403 M6, while the corresponding probabilistic projection is formed by averaging the  
404 cumulative distribution functions from the same eight methods.

405

406 Figure 9 compares the cross-validation performance of these method  
407 combinations in 2069-2098 with that for the best four individual methods. When  
408 the number of methods combined increases, the range of projections included in  
409 the combination widens. Consequently, *OutOfRange* is already smaller for B2  
410 than for the individual methods, and it decreases further when more methods are  
411 added (Fig. 9c). Still, as averaged over the whole temperature distribution, 15% of  
412 verification cases falling outside the predicted range remain even for A10.

413

414 As discussed above, decreases in *OutOfRange* do not necessarily imply a better  
415 projection. Indeed, *MSE* and *CRPS* (Figs. 9a,b) tell a partly different story. While  
416 *CRPS* averaged over the whole distribution is lower for all the four tested  
417 combinations than any individual method, it is at minimum for B4 that only  
418 includes the best four methods. Similarly, the B4 combination also has the lowest  
419 *MSE*. Its superiority over the other tested methods and combinations applies to  
420 most parts of the distribution (Figs. 9d,e). In particular, the probabilistic *CRPS*  
421 measure identifies B4 as the best approach with the only exception of the absolute  
422 extremes (0 and 100%). In the period 2001-2030 as well, *CRPS* and *MSE*  
423 averaged over the whole distribution are the lowest for B4 (not shown).

424

425 These findings suggest that temperature projections might be best derived by  
426 combining the information from the two delta change (M4-M5) and the two bias  
427 correction type quantile mapping methods (M9-M10). However, future research  
428 should test whether this conclusion remains valid for other model ensembles and  
429 other parts of the world.

## 430 **5. Projections for the future**

431 Here, we apply our methodology to real-world temperature projection for the set of  
432 139 stations shown in Fig. 4b. Although the methodology provides projections for  
433 temperature in absolute units (cf. Fig. 6), we mostly show the results here as  
434 changes from the observed baseline distribution in 1971-2000.

## 435 **5.1 Intermethod differences in projections**

436 Intermethod differences in the projections are studied in Figs. 10 and 11. The first  
437 three rows of Fig. 10 summarize the projected six-RCM mean changes in the 1st,  
438 10th, 50th, 90th and 99th percentiles of temperature, averaging the month- and  
439 station-specific values over the standard three-month seasons and over northern  
440 (48 stations north of 57.5°N), central (44 stations at 47.5°N-57.5°N) and southern  
441 Europe (47 stations south of 47.5°N). M1, which applies the same delta change in  
442 all parts of the distribution, provides a reference against which to compare the  
443 projections from the other methods. In line with earlier GCM and RCM studies  
444 (Räisänen et al. 2004, Kharin et al. 2007, Kjellström et al. 2007, Nikulin et al.  
445 2011), the ENSEMBLES RCMs simulate seasonally varying changes in  
446 variability that are reflected in all of M2-M10. In winter and to some extent in  
447 autumn and spring, the simulated variability decreases particularly in northern  
448 Europe, resulting in larger changes in the lower than the upper end of the  
449 distribution. In summer, the reverse happens in central and southern Europe, with  
450 larger increases in the highest than in the lowest temperatures.

451

452 Differences also occur between the projections from M2-M10. For example, the  
453 contrast between the changes in the lowest and highest winter temperatures in  
454 northern Europe is less pronounced for M4 than the rest of M2-M10. This is at  
455 least partly due to the running averaging of the quantile-quantile relationship in  
456 M4, which contracts the range over which changes in variability can be taken into  
457 account (Fig. 5a). More strikingly, the apparent increase in the highest (lowest)  
458 summer temperatures in central and southern Europe is larger (smaller) for M6  
459 than for the other methods. This is an artifact caused by the tendency of many of  
460 the models to overestimate present-day temperature variability in summer, a bias  
461 not corrected in M6.

462

463 The differences between the methods are smaller in the middle of the distribution  
464 than in the tails. However, M9 and M10 indicate a markedly smaller increase in  
465 median (50%) temperatures in central and southern Europe in summer than any  
466 other method. In southern Europe, M10 actually projects less warming than the  
467 other methods (excluding M6 near the lower tail) throughout the distribution. The  
468 explanation is analogous to the case shown in Figs. 5b and 6, but with the sign of  
469 the difference reversed. Because the simulated variability in southern and central

470 Europe in summer is too large, that is, the temperature bias increases with  
471 increasing temperature, a more negative bias correction is applied in M9 and M10  
472 to the higher temperatures simulated in the future. This reduces the projected  
473 warming, just as recently shown for a similar bias correction method by Boberg  
474 and Christensen (2012). The results in Fig. 7 suggest that this feature may very  
475 well be an improvement: it was precisely the ability of M9 and M10 to modify the  
476 time mean temperature change that reduced the *MSEs* of these methods (and to  
477 some extent M4 and M5) in cross validation.

478

479 Another question of interest is how the choice of the method affects the  
480 intermodel variation of the projections (bottom row of Fig. 10). For M1, the  
481 intermodel standard deviation as calculated over all 139 stations is relatively small  
482 (1.1-1.2°C depending on season), being the same for all parts of the distribution.  
483 For the other methods, the standard deviation near the tails of the distribution is in  
484 most cases larger, particularly in the lower tail in winter and in the upper tail in  
485 summer. The intermodel variation tends to be the largest for M6, being amplified  
486 by uncorrected biases in variability. The standard deviation is in most cases  
487 smaller for the delta change than for the bias correction methods, because the  
488 former do not fully represent the uncertainty associated with internal variability  
489 (see the discussion of *OutOfRange* in Section 4). Note, however, the typically  
490 smaller standard deviations for M9-M10 than for the other bias correction  
491 methods.

492

493 To further compare the ten methods, pairwise intermethod root-mean-square (rms)  
494 differences in the six-RCM mean temperature projections are shown in Fig. 11  
495 (see the caption for further details). Method pairs 2-3, 4-5, 7-8, and 9-10 all stand  
496 out as closely related, with very small differences in most of the distribution. In  
497 particular, the differences between M2 and M3 and between M7 and M8 are  
498 largely negligible, except for the extreme tails where changes and bias corrections  
499 of skewness have more substantial effects. Furthermore, while the M9-M10  
500 differences are small in the lower and central parts of the distribution, they grow  
501 relatively large in the upper tail, reflecting the difficulty in the extrapolation of the  
502 bias correction beyond the range in the baseline period. As a whole, M6 and M1  
503 are the methods furthest away from the others, but the tendency of M9 and M10 to



504 differ relatively strongly from the other methods in the middle of the distribution  
505 also stands out.

## 506 **5.2 Analysis of variance**

507 In addition to the choice among the various projection methods, the projections  
508 also depend on the model simulation used. To assess the relative importance of  
509 these sources of uncertainty, fixed-effect analysis of variance was applied. The  
510 variance within each data set, consisting of all model- and method-specific  
511 projections for a given quantile of the temperature distribution at a given station  
512 and month, was decomposed as

$$513 \quad V_{tot} = V_{mod} + V_{met} + V_{int} \quad (8)$$

514 where  $V_{mod}$  is the contribution of model differences (variation of multi-method  
515 mean projections across models),  $V_{met}$  that of method differences (variation of  
516 multi-model mean projections across methods), and  $V_{int}$  that of model-method  
517 interaction (method-dependence of intermodel differences, or equivalently model-  
518 dependence of intermethod differences). The computation of these terms is  
519 analogous to Eqs. (1)-(4) of Déqué et al. (2012). We stress that this decomposition  
520 does not aim to estimate the variances that would be observed within an infinite  
521 population of independent models and methods, but is rather used to diagnose the  
522 sources of variability within our specific set of (possibly non-independent) models  
523 and (certainly non-independent) methods. Furthermore, model simulations of  
524 climate always include unforced natural variability (Räisänen 2001, Yip et al.  
525 2011). Some fraction of  $V_{mod}$  reflects this unforced variability rather than genuine  
526 intermodel differences in response to forcing, and to a lesser extent the unforced  
527 variability may also affect the other variance components.

528

529 As an illustration, the model- and method-dependence of projections for the 1st  
530 and 50th percentiles of January daily temperature in Jyväskylä, Finland in 2069-  
531 2098 is shown in Fig. 12. The projections for the 1st percentile vary widely  
532 between the models, but even more so between the methods. With all six models  
533 and all ten methods included,  $V_{tot} = 11.35 \text{ (}^\circ\text{C)}^2$ , of which 77% is attributed to  
534 method differences and only 6% to model differences, model-method interaction  
535 taking the remaining 17%. In particular, M1 gives systematically lower  
536 projections than the other methods, whereas the highest projections are generally

537 obtained from M6. However, both of these methods are suspect due to their poor  
538 performance in cross validation. Indeed, the results in Fig. 9 suggest that it might  
539 be preferable to only retain the methods 4, 5, 9 and 10 included in the B4  
540 combination. Doing this reduces the total variance by about 40%, but does not  
541 affect the relative shares of the different components in this particular case.

542

543 The projections for the 50th percentile are much less method-dependent (Fig.  
544 12b). Almost 70% of the variance is attributed to model differences when all ten  
545 methods are included, and this increases to 88% when only the best four methods  
546 are retained. Conversely, the contribution of method differences is reduced from  
547 15% in the former case to nearly zero in the latter.

548

549 Averaging over the 139 stations and 12 months confirms that intermodel  
550 differences strongly dominate the variance in the central parts of the temperature  
551 distribution (Fig. 13). Intermethod differences and model-method interaction both  
552 grow more important towards the tails of the distribution but do not become as  
553 dominant as in the case shown in Fig. 12a, especially not when only the best four  
554 methods are included. Averaging the variances over the whole distribution,  $V_{mod}$ ,  
555  $V_{int}$  and  $V_{met}$  contribute 73%, 13% and 14% in the 10-method case and 76%, 12%  
556 and 13% in the 4-method case, respectively. Therefore, the uncertainty associated  
557 with the choice of the projection method may be a secondary issue for many  
558 applications, although it clearly should not be neglected when and where the tails  
559 of the temperature distribution are particularly important. A similar conclusion –  
560 that uncertainty in bias correction is generally smaller than climate modeling  
561 uncertainty – was obtained by Chen et al. (2011), although their hydrological  
562 study addressed the bias correction uncertainty due to the choice of the baseline  
563 period rather than due to the choice of the method.

## 564 **6. Conclusions**

565 Projection of future climate cannot be generally based on model simulations alone  
566 but also requires information on the observed climate. This projection problem is  
567 often considered simple when only long-term climatic means are required, but it  
568 becomes more complicated when temporal variability is important. Here, we have  
569 focused on what is probably one of the easiest aspects of daily-scale variability for

570 both climate models and in terms of its statistical properties, distributions of daily  
571 mean temperature in a changing climate. We first studied the relative strengths  
572 and weaknesses of ten projection methods using cross validation among six RCM  
573 simulations for Europe, all made with different RCMs and different driving  
574 GCMs. The main findings from these tests include the following:

575

- 576 1. Delta change and bias correction type methods showed similar overall  
577 performance in cross validation of late 21<sup>st</sup> century (2069-2098) temperature  
578 distributions. Within both groups, quantile mapping approaches performed  
579 best, due to their smallest errors in the projected time mean temperature. The  
580 simplest approaches assuming constant change or constant bias throughout the  
581 distribution were the worst, having larger errors in the distribution of  
582 temperature around the mean value than the other methods. In projections for  
583 early 21<sup>st</sup> century (2001-2030), the intermethod differences in verification  
584 statistics were smaller, except for the poor performance of the constant-bias  
585 bias correction method.
- 586 2. The performance of different projection methods may vary across the  
587 temperature distribution. In particular, quantile mapping type bias correction  
588 methods were found to be less reliable in the extreme upper tail than in the  
589 other parts of the distribution.
- 590 3. No single method performs best under all circumstances. Thus, to some  
591 extent, the choice of the projection method represents an uncertainty  
592 analogous to the choice of the climate model used for the projection. A natural  
593 way to take this uncertainty into account is to consider a few different but  
594 well-performing projection methods instead of just one. In our cross-  
595 validation exercise, the combination of the two delta change and two bias  
596 correction quantile mapping methods generally outperformed each individual  
597 method.

598

599 Second, we assessed the sensitivity of the resulting 21<sup>st</sup> century temperature  
600 projections to the choice of the method, to find that

601

- 602 1. The choice of the projection method has typically a larger impact in the tails  
603 of the temperature distribution than in the central parts. However, the latter

604 may also be affected. In particular, our quantile mapping type bias correction  
605 methods suggest a smaller warming in southern and central Europe in summer  
606 than would be inferred directly from the model simulations. This supports the  
607 recent findings of Boberg and Christensen (2012), in particular as these  
608 methods performed well in cross validation.

609 2. The uncertainty associated with the choice of the model simulation generally  
610 exceeds that due to the choice of the projection method. However, the relative  
611 importance of the method uncertainty increases towards the tails of the  
612 distribution, indicating that this uncertainty should also be considered at least  
613 in applications where extremely low or high temperatures are important.

614

615 Our study is based on only six RCM simulations and it only covers the European  
616 area. Its conclusions, particularly regarding the relative performance of different  
617 projection methods, should therefore be verified with other data sets. New  
618 opportunities for this will be provided by the CORDEX initiative (A COordinated  
619 Regional climate Downscaling EXperiment, [cordex.dmi.dk/](http://cordex.dmi.dk/)), and to some extent  
620 also by the GCM simulations conducted in the fifth phase of the Coupled Model  
621 Intercomparison Project (Taylor et al. 2011).

622

623 We also stress that the methods studied here were only designed for, and tested  
624 for their fidelity in, changing or correcting the local frequency distribution of  
625 daily mean temperatures. Issues that we have not addressed include the temporal  
626 (Haerter et al. 2011) and spatial autocorrelation structure (Huth 2002), as well as  
627 the correlation of temperature with other variables such as precipitation (Engen-  
628 Skaugen 2007).

629

## 630 **Acknowledgments**

631 The model simulations used in this work were funded by the EU FP6 Integrated  
632 Project ENSEMBLES (Contract number 505539). This work has been supported  
633 by the Academy of Finland RECAST project (decision 140801). The two  
634 reviewers are acknowledged for their constructive comments.

635

636 **References**

- 637 Amengual A, Homar V, Romero R, Alonso S, Ramis C (2012) A statistical  
638 adjustment of regional climate model outputs to local scales: application to  
639 Platja de Palma, Spain. *J Climate* 25: 939–957, doi:  
640 <http://dx.doi.org/10.1175/JCLI-D-10-05024.1>
- 641 Ballester, J., F. Giorgi, and X. Rodó (2010) Changes in European temperature  
642 extremes can be predicted from changes in PDF central statistics. *Clim*  
643 *Change* 98: 277-284, doi:10.1007/s10584-009-9758-0
- 644 Boberg F, Christensen JH (2012) Overestimation of Mediterranean summer  
645 temperature projections due to model deficiencies. *Nature Climate Change*  
646 2: 433-436, doi: 10.1038/NCLIMATE1454
- 647 Bracegirdle TJ, Stephenson DB (2012) Higher precision estimates of regional  
648 polar warming by ensemble regression of climate model projections.  
649 *Climate Dyn*, doi: 10.1007/s00382-012-1330-3
- 650 Buser CM, Künsch HR, Lüthi D, Wild M, Schär C (2009) Bayesian multi-model  
651 projection of climate: bias assumptions and interannual variability. *Climate*  
652 *Dyn* 33: 849-868, doi: 10.1007/s00382-009-0588-6
- 653 Chen C, Haerter JO, Hagemann S, Piani C (2011) On the contribution of  
654 statistical bias correction to the uncertainty in the projected hydrological  
655 cycle. *Geophys Res Lett* 38: L20403, doi: 10.1029/2011GL049318
- 656 Christensen JH, Hewitson B, Busuioc A, Chen A, Gao X, Held I, Jones R, Kolli  
657 RK, Kwon W-T, Laprise R, Magaña Rueda V, Mearns L, Menéndez CG,  
658 Räisänen J, Rinke A, Sarr A, Whetton P (2007) Regional Climate  
659 Projections. *Climate Change 2007: the Physical Science Basis*, S Solomon  
660 et al., Eds., Cambridge University Press, pp 847-940
- 661 Déqué M, Somot S, Sanchez-Gomez E, Goodess CM, Jacob D, Lenderink G,  
662 Christensen OB (2012) The spread amongst ENSEMBLES regional  
663 scenarios: regional climate models, driving general circulation models and  
664 interannual variability. *Climate Dyn* 38: 951-964, doi: 10.1007/s00382-011-  
665 1053-x
- 666 Dosio A, Paruolo P (2011) Bias correction of the ENSEMBLES high-resolution  
667 climate change projections for use by impact models: evaluation on the  
668 present climate. *J Geophys Res* 116: D16106, doi: 10.1029/2011JD015934
- 669 Engen-Skaugen T (2007) Refinement of dynamically downscaled precipitation  
670 and temperature scenarios. *Clim Change* 84: 365-382, doi: 10.1007/s10584-  
671 007-9251-6
- 672 Haerter JO, Hagemann S, Moseley C, Piani C (2011) Climate model bias  
673 correction and the role of timescales. *Hydrol Earth Syst Sci* 15: 1065-1079,  
674 doi: 10.5194/hess-15-1065-2011
- 675 Ho CK, Stephenson DB, Collins M, Ferro CAT, Brown SJ (2012) Calibration  
676 strategies. A source of additional uncertainty in climate change projections.  
677 *Bull Am Meteorol Soc* 93: 21-26, doi: 10.1175/2011BAMS3110.1
- 678 Huth R (2002) Statistical downscaling of daily temperature in Central Europe. *J*  
679 *Climate* 15: 1731–1742, doi: <http://dx.doi.org/10.1175/1520->  
680 [0442\(2002\)015<1731:SDODTI>2.0.CO;2](http://dx.doi.org/10.1175/1520-0442(2002)015<1731:SDODTI>2.0.CO;2)

681 Kharin VV, Zwiers FW, Zhang X, Hegerl GC (2007) Changes in temperature and  
682 precipitation extremes in the IPCC ensemble of global coupled model  
683 simulations. *Journal of Climate* 20: 1419-1444, doi: 10.1175/JCLI4066.1  
684 Kjellström E, Bärring L, Jacob D, Jones R, Lenderink G, Schär C (2007)  
685 Modelling daily temperature extremes: Recent climate and future changes  
686 over Europe. *Clim Change* 81: 249-265, DOI 10.1007/s10584-006-9220-5  
687 Maraun D (2012) Nonstationarities of regional climate model biases in European  
688 seasonal mean temperature and precipitation sums. *Geophys Res Lett* 39:  
689 L06706, doi: 10.1029/2012GL051210  
690 Maraun D, Wetterhall F, Ineson AM, Chandler ER, Kendon EJ, Widmann M,  
691 Brienen S, Rust HW, Sauter T, Themeßl, Cenema VKC, Chun KP, Goodess  
692 CM, Jones RG, Onof C, Vrac M, Thiele-Eich I (2010) Precipitation  
693 downscaling under climate change: recent developments to bridge the gap  
694 between dynamical models and the end user. *Rev Geophys* 48: RG3003,  
695 doi: 10.1029/2009RG000314  
696 Meehl GA, Stocker TF, Collins W, Friedlingstein P, Gaye A, Gregory J, Kitoh A,  
697 Knutti R, Murphy J, Noda A, Raper S, Watterson I, Weaver A, Zhao Z-C  
698 (2007) Global climate projections. *Climate Change 2007: the Physical  
699 Science Basis*, S Solomon et al., Eds., Cambridge University Press, pp 747-  
700 845  
701 Nikulin G, Kjellström E, Hansson U, Strandberg G, Ullerstig A (2011) Evaluation  
702 and future projections of temperature, precipitation and wind extremes over  
703 Europe in an ensemble of regional climate simulations. *Tellus A* 63: 41–55,  
704 doi: 10.1111/j.1600-0870.2010.00466.x  
705 Piani C, Weedon GP, Best M, Gomes SM Viterbo P, Hagemann S, Haerter JO  
706 (2010) Statistical bias correction of global simulated daily precipitation and  
707 temperature for the application of hydrological models. *J Hydrology* 395:  
708 199-215, doi: 10.1016/j.jhydrol.2010.10.024  
709 Räisänen J (2001) CO<sub>2</sub>-induced climate change in CMIP2 experiments.  
710 Quantification of agreement and role of internal variability. *J Climate* 14:  
711 2088-2104, doi: [http://dx.doi.org/10.1175/1520-](http://dx.doi.org/10.1175/1520-0442(2001)014<2088:CICCIC>2.0.CO;2)  
712 [0442\(2001\)014<2088:CICCIC>2.0.CO;2](http://dx.doi.org/10.1175/1520-0442(2001)014<2088:CICCIC>2.0.CO;2)  
713 Räisänen J, Palmer TN (2001) A probability and decision-model analysis of a  
714 multi-model ensemble of climate change simulations. *J Climate* 14: 3212-  
715 3226, doi: [http://dx.doi.org/10.1175/1520-](http://dx.doi.org/10.1175/1520-0442(2001)014<3212:APADMA>2.0.CO;2)  
716 [0442\(2001\)014<3212:APADMA>2.0.CO;2](http://dx.doi.org/10.1175/1520-0442(2001)014<3212:APADMA>2.0.CO;2)  
717 Räisänen J, Ylhäisi J (2011) How much should climate model output be smoothed  
718 in space? *J Climate* 24: 867-880, doi:  
719 <http://dx.doi.org/10.1175/2010JCLI3872.1>  
720 Räisänen J, Hansson U, Ullerstig A, Döscher R, Graham LP, Jones C, Meier  
721 HEM, Samuelsson P, Willén U (2004) European climate in the late 21st  
722 century: regional simulations with two driving global models and two  
723 forcing scenarios. *Climate Dynamics* 22: 13-31, doi: 10.1007/s00382-003-  
724 0365-x  
725 Randall DA, Wood RA, Bony S, Colman R, Fichetfet T, Fyfe J, Kattsov V, Pitman  
726 A, Shukla J, Srinivasan J, Stouffer RJ, Sumi A, Taylor KE (2007) Climate

727 models and their evaluation. *Climate Change 2007: the Physical Science*  
728 *Basis*, S Solomon et al., Eds., Cambridge University Press, pp 589-662  
729 Taylor KE, Stouffer RJ, Meehl GA (2011) A summary of the CMIP5 experiment  
730 design. Available at [http://cmip-](http://cmip-pcmdi.llnl.gov/cmip5/docs/Taylor_CMIP5_design.pdf)  
731 [pcmdi.llnl.gov/cmip5/docs/Taylor\\_CMIP5\\_design.pdf](http://cmip-pcmdi.llnl.gov/cmip5/docs/Taylor_CMIP5_design.pdf)  
732 Themeßl MJ, Gobiet A, Leuprecht A (2011) Empirical-statistical downscaling and  
733 error correction of daily precipitation from regional climate models. *Int J*  
734 *Climatology* 31: 1530-1544, doi: 10.1002/joc.2168  
735 van der Linden P, Mitchell JFB, Eds. (2009) ENSEMBLES: Climate Change and  
736 its Impacts: Summary of research and results from the ENSEMBLES  
737 project. Met Office Hadley Centre, FitzRoy Road, Exeter EX1 3PB, UK,  
738 160 pp. Available at [http://ensembles-](http://ensembles-eu.metoffice.com/docs/Ensembles_final_report_Nov09.pdf)  
739 [eu.metoffice.com/docs/Ensembles\\_final\\_report\\_Nov09.pdf](http://ensembles-eu.metoffice.com/docs/Ensembles_final_report_Nov09.pdf)  
740 Yip S, Ferro CAT, Stephenson D (2011) A simple, coherent framework for  
741 partitioning uncertainty in climate change predictions. *J Climate* 24: 4634-  
742 4643, doi: 10.1175/2011JCLI4085.1  
743

744

745 **Tables**

746 **Table 1** *The RCM simulations used in this study*

Driving GCM	RCM	Institution	Shorthand
ARPEGE	ALADIN	CNRM	CNRM-A
HadCM3Q0	CLM	ETHZ	ETHZ-H0
HadCM3Q3	HadRM3Q3	Met Office	METO-H3
HadCM3Q16	HadRM3Q16	Met Office	METO-H16
ECHAM5-r3	REMO	MPI	MPI-E5
BCM	RCA3	SMHI	SMHI-BCM

747 *The first column indicates the driving global climate model, the second the regional*  
748 *climate model and the third the institution that conducted the simulations, using model*  
749 *and institution acronyms that follow the ENSEMBLES Research Theme 3 web page*  
750 *(<http://ensemblesrt3.dmi.dk/>). The last column gives the shorthand notations used in this*  
751 *article*

752

753

754

755



756 **Table 2.** *The projection methods used in this study*

---

M1	Delta change: mean
M2	Delta change: mean + standard deviation
M3	Delta change: mean + standard deviation + skewness
M4	Delta change: quantile mapping using smoothing
M5	Delta change: quantile mapping using linear regression
M6	Bias correction: mean
M7	Bias correction: mean + standard deviation
M8	Bias correction: mean + standard deviation + skewness
M9	Bias correction: quantile mapping using smoothing
M10	Bias correction: quantile mapping using linear regression

---

757

758

759

760

761

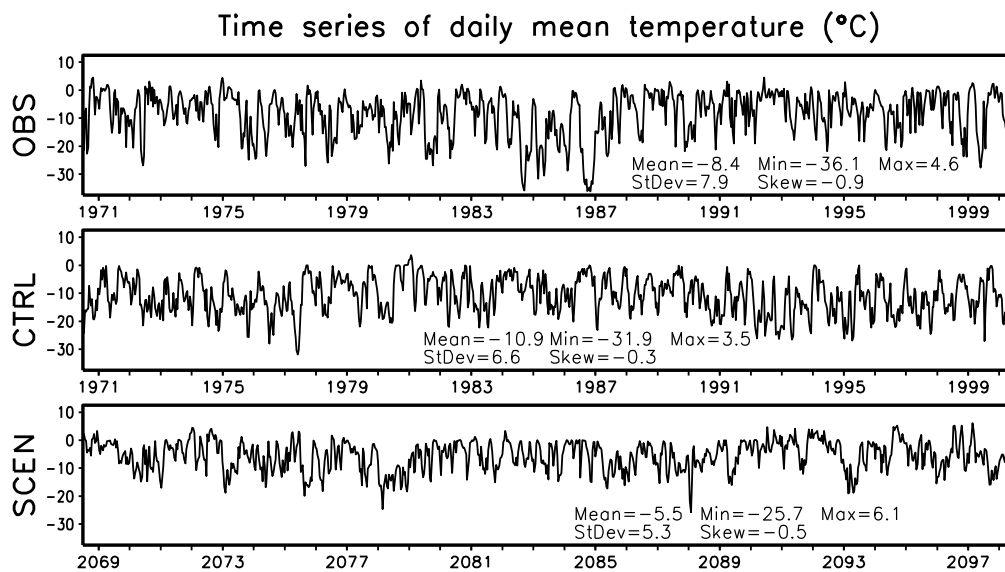
762

763

764

765

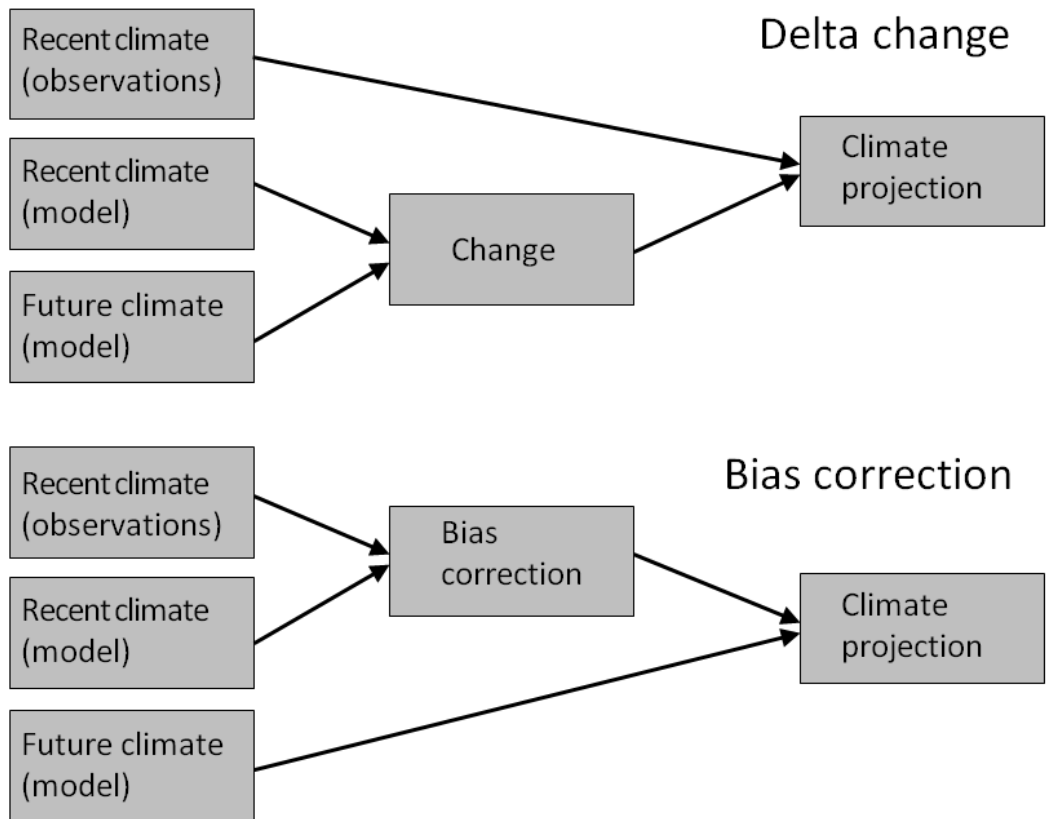
766 **Figures**



767

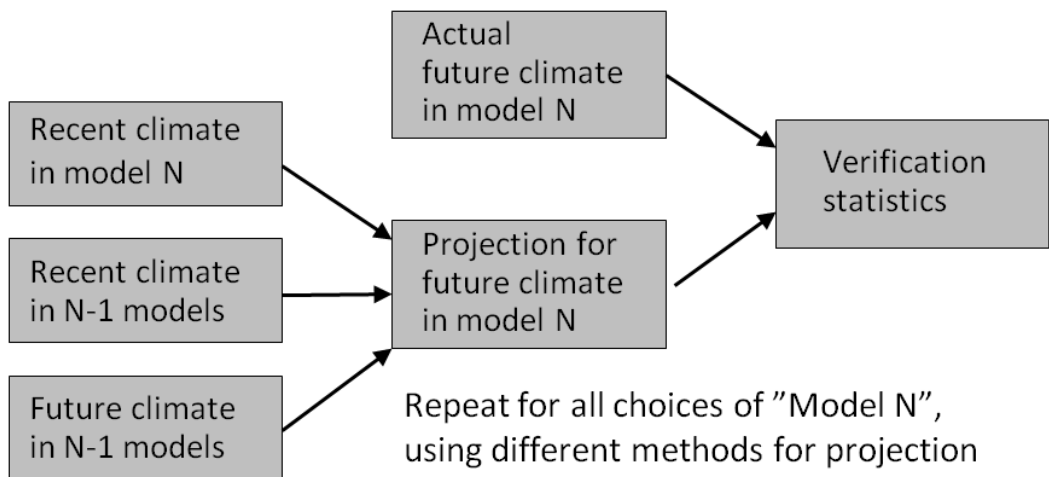
768 **Fig. 1** 30-year time series of January daily mean temperature in Jyväskylä, Finland  
769 (62.4°N, 25.7°E), as observed in 1971-2000 (top), in the ETHZ-H0 (see Table 1) RCM  
770 simulation during the same period (middle), and in the same RCM simulation in 2069-  
771 2098 (bottom)

772



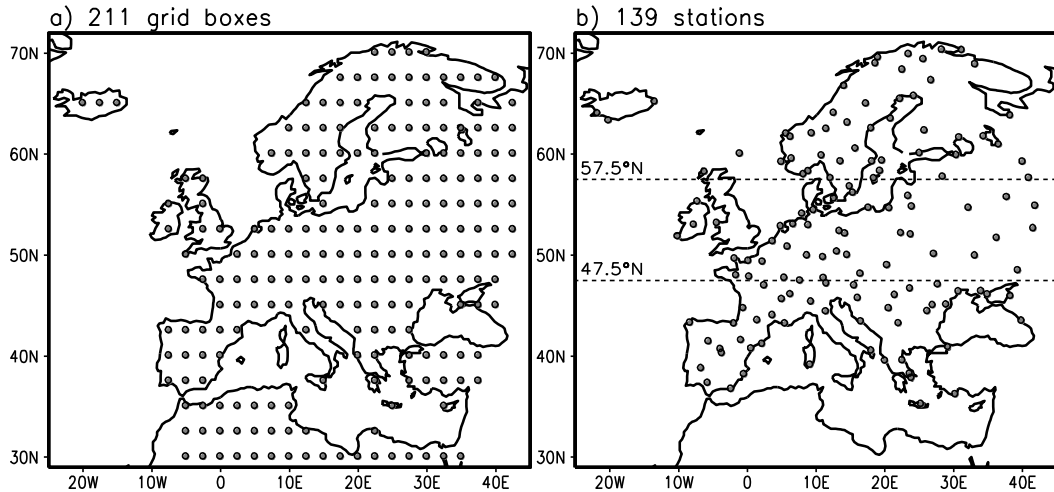
773

774 **Fig. 2** A schematic illustration of delta change and bias correction methods

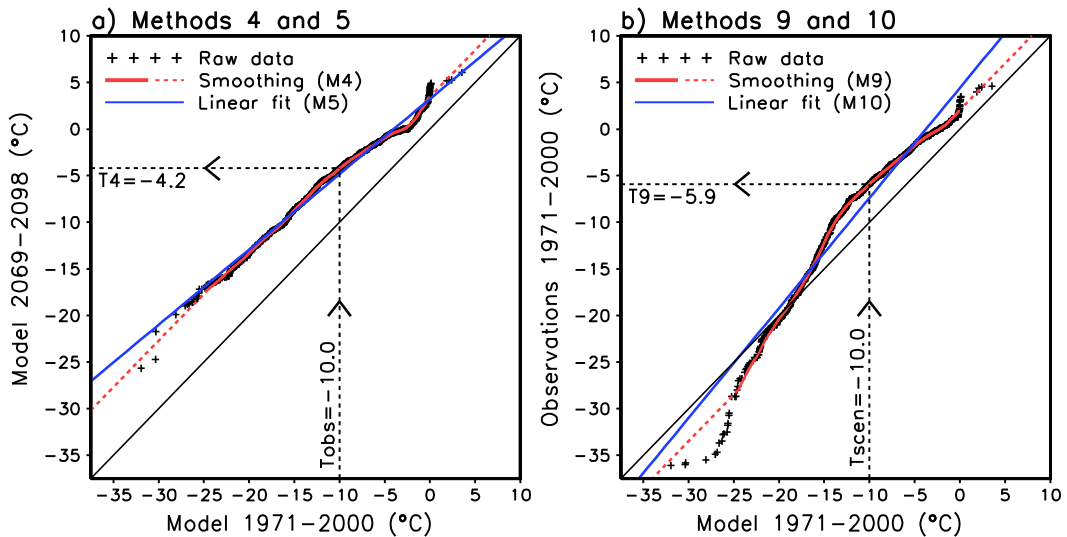


775

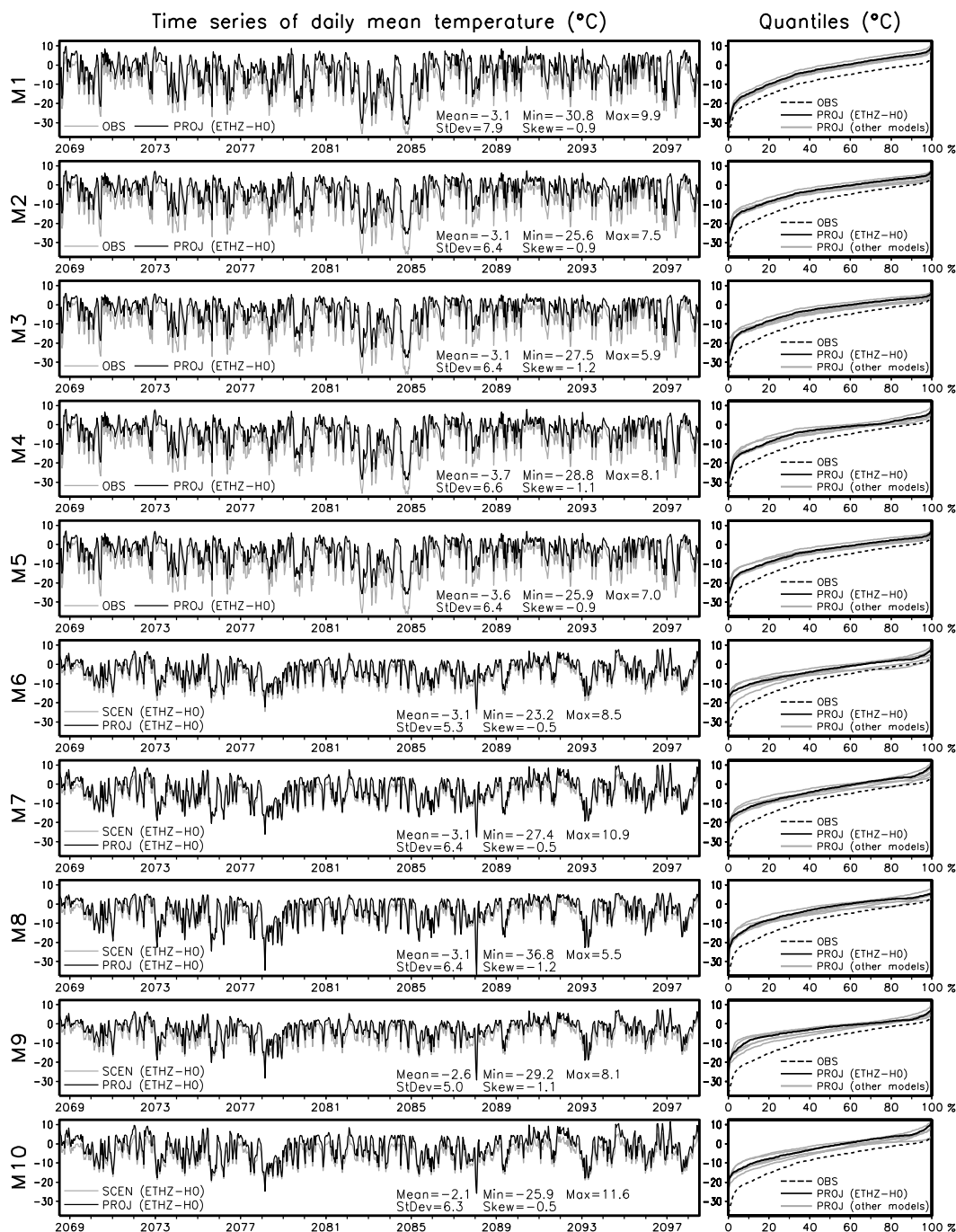
776 **Fig. 3** Principle of cross validation, as used in this study



777  
 778 **Fig. 4** Locations of (a) the 211 grid boxes used in cross validation in Section 4, and (b)  
 779 the 139 weather stations used in real-world temperature projection in Section 5. The  
 780 latitudes 47.5°N and 57.5°N used for the division in Fig. 10 are also indicated in (b)



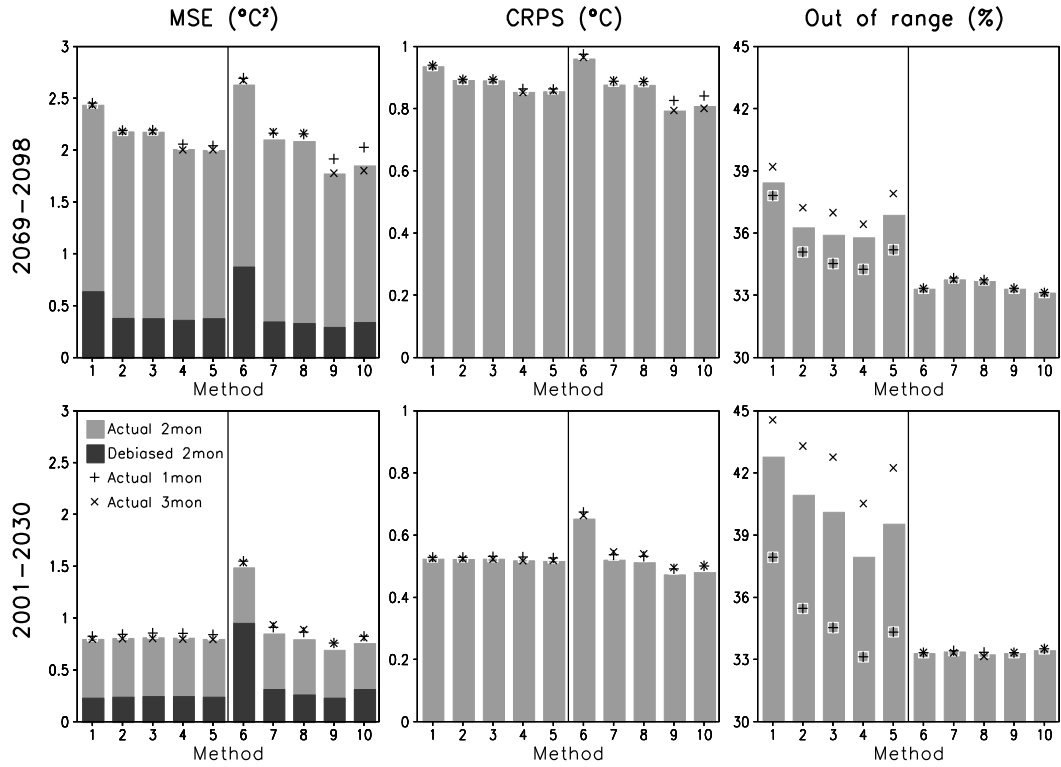
781  
 782 **Fig. 5** An illustration of quantile mapping methods for the case depicted in Figs. 1 and 6.  
 783 In (a), the crosses show the quantile-quantile plot obtained directly from the simulations  
 784 for 1971-2000 and 2069-2098, the red line gives the smoothed curve used in M4 and its  
 785 extrapolation, and the blue line depicts the linear regression used in M5. Using M4, an  
 786 observed temperature of  $-10.0^{\circ}\text{C}$  would be converted to  $-4.2^{\circ}\text{C}$  in the projection for 2069-  
 787 2098. (b) is the same for the comparison of the simulation and observations in 1971-  
 788 2000, as used in M9 and M10. In M9, a temperature of  $-10.0^{\circ}\text{C}$  in the scenario period  
 789 simulation would be converted to  $-5.9^{\circ}\text{C}$  in the projection



790

791 **Fig. 6** (left) Projected 30-year (2069-2098) time series of January daily mean  
 792 temperature in Jyväskylä, Finland, using data from the ETHZ-H0 simulation, in methods  
 793 M1-M10 (black lines). The grey lines in the top five (bottom five) panels show  
 794 observations for 1971-2000 (the ETHZ-H0 scenario simulation for 2069-2098). The  
 795 numeric values in each panel give the mean, minimum and maximum, standard deviation  
 796 and the skewness of the distribution within the projected time series. (right) Quantiles of  
 797 the observed distribution in 1971-2000 (dashed, same in all panels), and of the projected  
 798 distributions obtained using data from ETHZ-H0 (black solid line) and the other five  
 799 RCMs (grey lines)

800



801

802

**Fig. 7** Cross-validated MSE, CRPS and OutOfRange for temperature distributions in the

803

years 2069-2098 (top) and 2001-2030 (bottom). The MSE that would have been reached

804

if always predicting the correct 30-year monthly mean temperature is also shown (dark

805

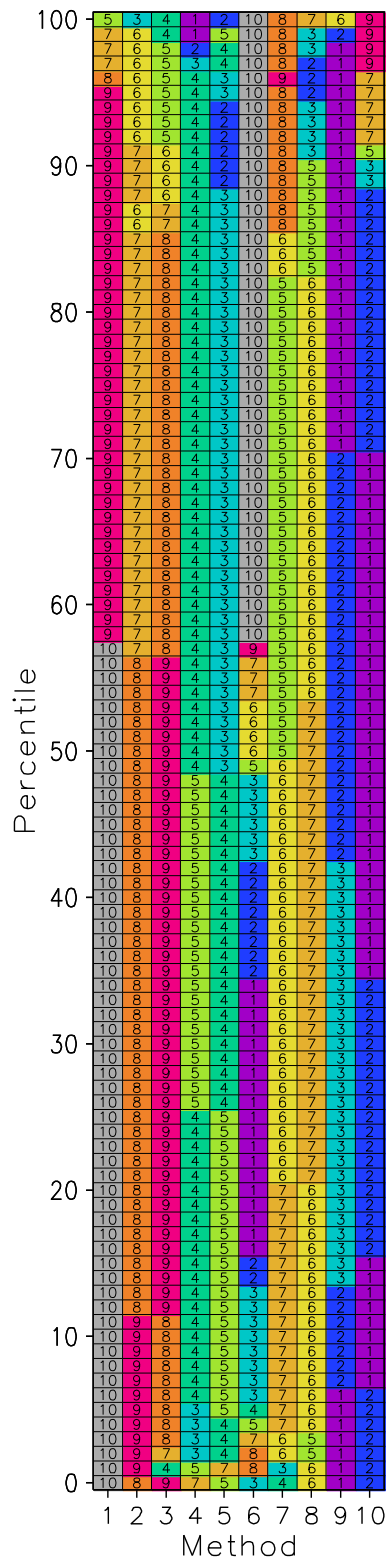
part of the bars in the left column). The bars give statistics based on two-month sampling

806

of climate changes (M1-M5) and biases (M6-M10); the plus signs (+) and crosses (×)

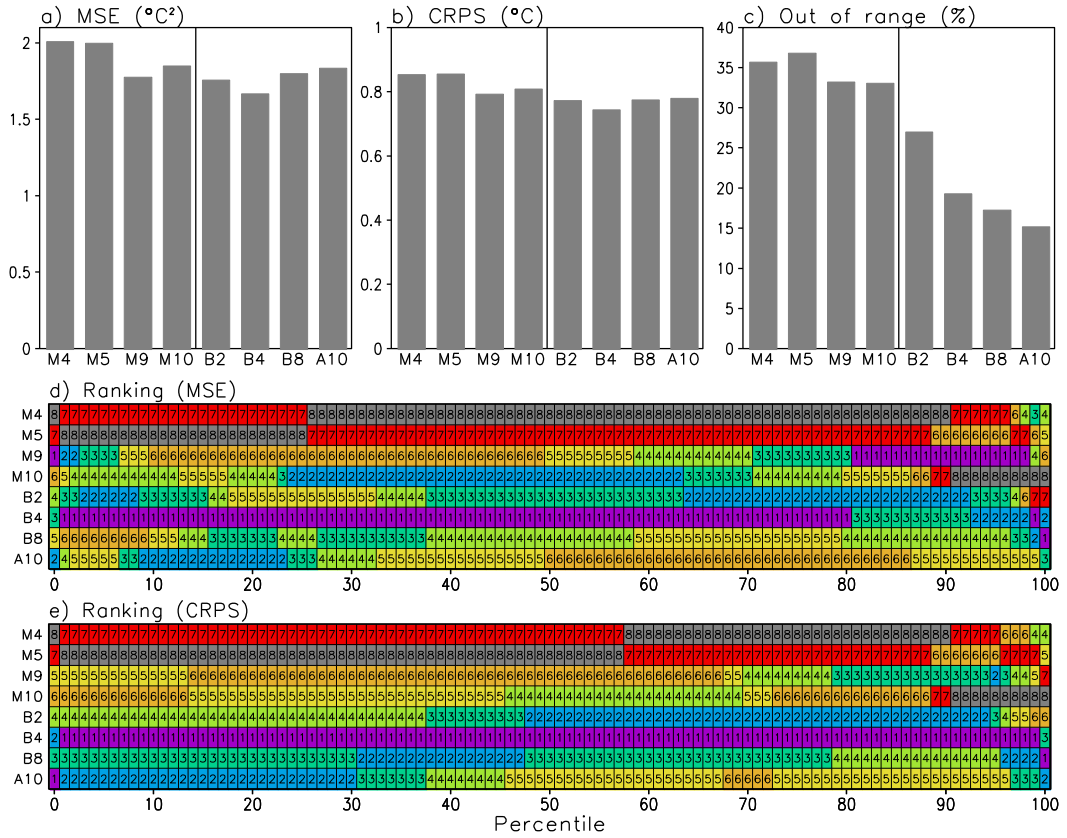
807

show the corresponding values for 1-month and 3-month sampling



808

809 **Fig. 8** Ranking of the 10 methods (1 best, 10 worst) for cross-validated MSE of different  
 810 percentiles of the temperature distribution in 2069-2098 (0% = absolute monthly minima,  
 811 100% = absolute monthly maxima)



812

813

**Fig. 9** Cross-validation statistics for temperature in the years 2069-2098. The top row

814

shows MSE, CRPS and OutOfRange separately for four individual methods (M4, M5, M9

815

and M10) and for the combinations B2, B4, B8 and A10 defined in the text. The last two

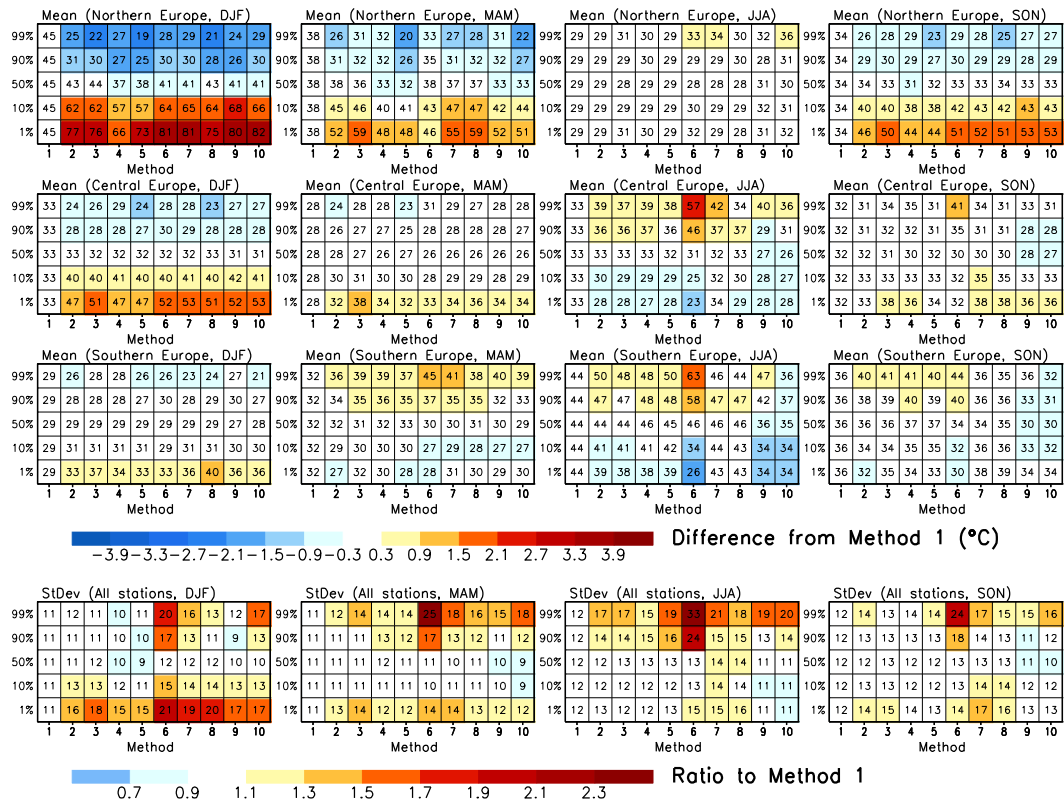
816

panels indicate the ranking (1 = best, 8 = worst) of the MSE and CRPS values within this

817

sample of methods in different parts of the temperature distribution





818

819

*Fig. 10 Summary of temperature projections for the years 2069-2098. The first three rows show the six-model mean changes in five quantiles of the temperature distribution (1% to 99%), as averaged over northern, central and southern Europe (48 stations north of 57.5°N, 44 stations at 47.5°N-57.5°N and 47 stations south of 47.5°N, respectively).*

820

821

*The quantiles were first calculated for each month and then averaged over the three-month seasons identified in the figure headers. The numeric values indicate the absolute difference from the observed value in 1971-2000 (unit: 0.1°C), and the shading gives the difference from the value for M1. The numeric values in the last row show the intermodel standard deviation of the projections (unit: 0.1°C), as calculated from the variance averaged over all 139 stations and the three months in each season. The shading indicates the ratio to the standard deviation for M1*

822

823

824

825

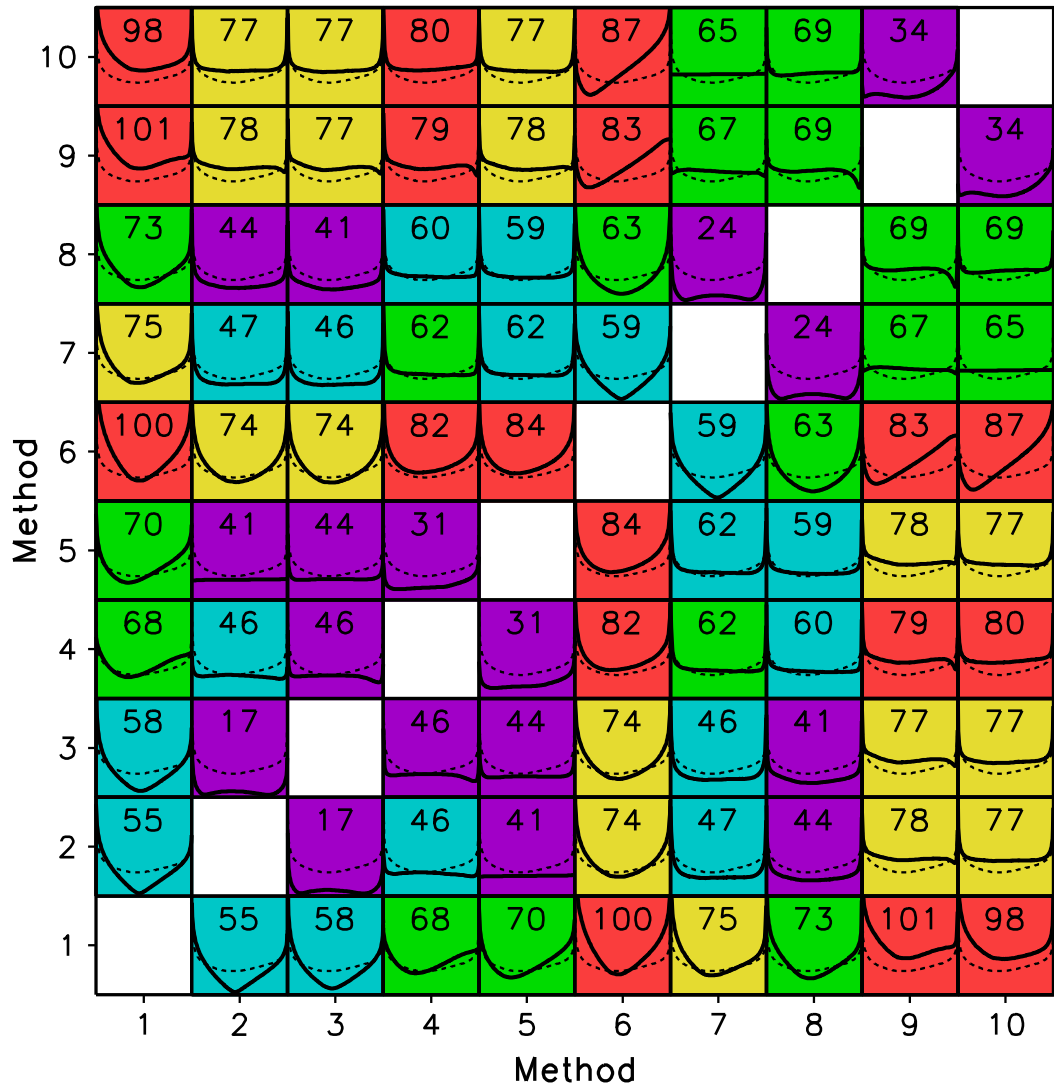
826

827

828

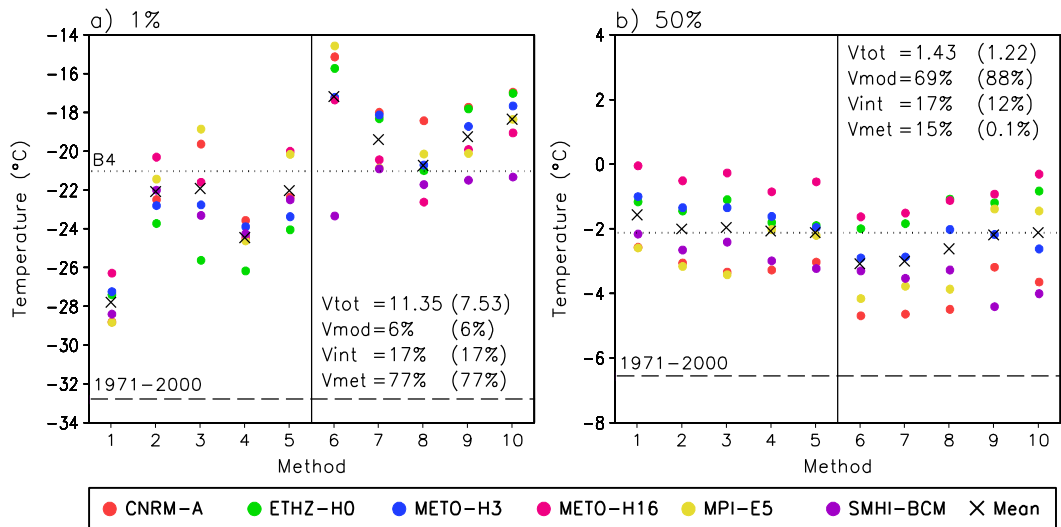
829

830



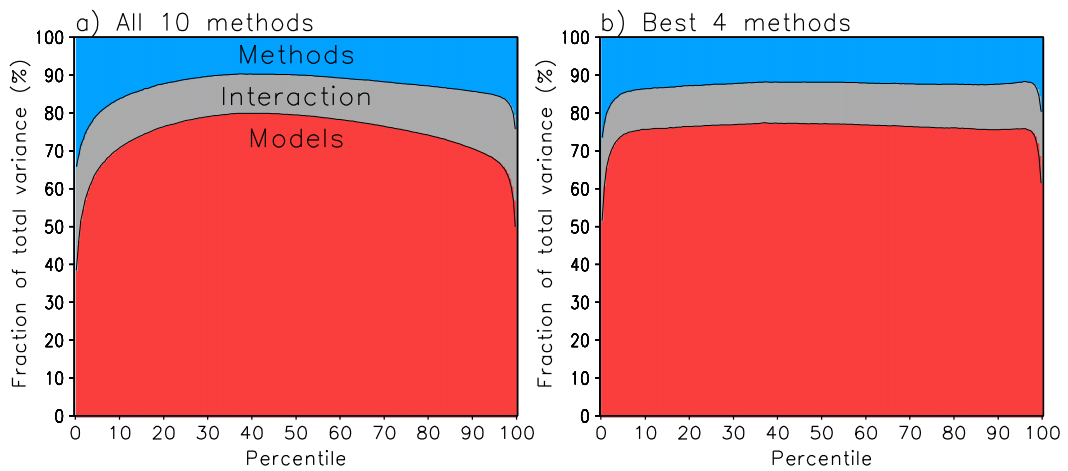
831

832 **Fig. 11** Intermethod rms differences of six-model mean temperature projections for the  
 833 years 2069-2098, using data for all 12 months and the 139 stations. In each cell, the  
 834 thick solid line shows the rms difference between the two methods compared (scale from  
 835 0 to 2°C) for quantiles ranging from 0 to 100% (left to right), while the dashed line gives  
 836 the rms difference averaged over all method pairs. The cells are coloured according to  
 837 the rms difference calculated for the whole distribution (violet: lowest 20% of cases ...  
 838 red: highest 20% of cases), which is plotted in each cell in units of 0.01°C



839

840 **Fig. 12** Projections for (a) the 1st and (b) the 50th percentile of daily mean temperature  
 841 in January in Jyväskylä, Finland in 2069-2098, based on the different methods and RCM  
 842 simulations. The dashed lines show the observed values in 1971-2000, and the dotted line  
 843 the six-model means averaged over the methods 4, 5, 9 and 10. The total variance in (°C)<sup>2</sup>  
 844 and the relative contributions of model differences, model-method interaction and method  
 845 differences are also shown, both when including all 10 methods (first numbers) and when  
 846 only including methods 4, 5, 9 and 10 (numbers in parentheses)



847

848 **Fig. 13** The relative contributions of model differences, model-method interaction and  
 849 method differences to the variance of the projected temperatures in 2069-2098 as a  
 850 function of the percentile of the distribution. The variances are averaged over the 12  
 851 months and the 139 stations. In (a), all 10 methods are included in the analysis, in (b)  
 852 only methods 4, 5, 9 and 10

853