

# Responses to Reviewer #1

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33

## Point #1

*Review to Yang et al., 2021, Bias-correcting individual inputs prior to combined calibration leads to more skillful forecasts of reference crop evapotranspiration. HESSD.*

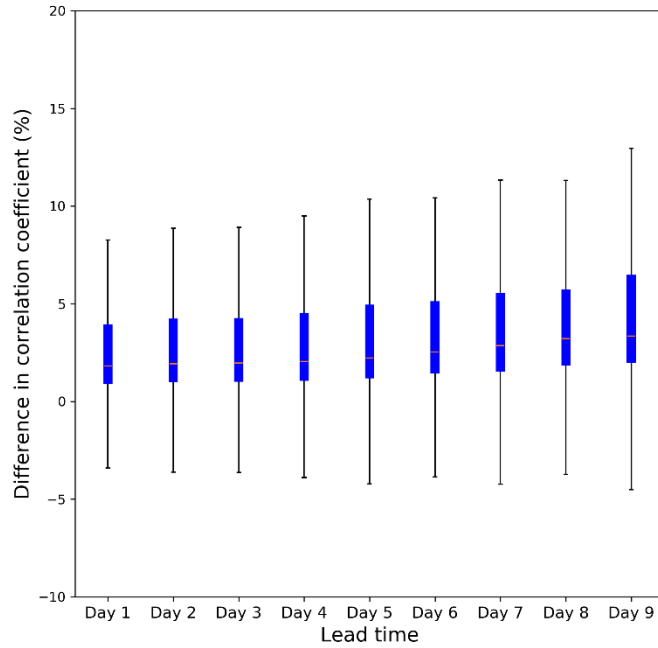
*In this study, the authors investigated a critical issue in the forecasting of short-term reference crop evapotranspiration (ET<sub>o</sub>) based on NWP outputs. It is getting popular that weather forecasts from NWP models are used to predict water loss through evapotranspiration. Such information is highly valuable for the effective management of water resources, particularly in arid/semi-arid regions. This investigation develops a new methodology that effectively corrects errors in ET<sub>o</sub> forecasts, and adds extra skills to statistical calibration. I believe this new post-processing strategy could benefit future NWP-based ET<sub>o</sub> forecasting. To improve this work, the authors should pay special attention to the following key issues:*

**Response: We appreciate the reviewer's insightful comments. We also believe the findings of this work could contribute to improving future NWP-based ET<sub>o</sub> forecasting. We address your constructive comments thoroughly and carefully and believe this work has been improved significantly. Please find more details in our point-by-point response.**

## Point #2

*1, Presentation of the results could be improved. Currently, the authors use maps to show/compare results from different model experiments. These figures could demonstrate the spatial patterns of modeling results. However, it might be more useful if the authors could summarize regional results in a different way, such as using boxplots. I believe that will better show readers the overall statistical information across the whole country than simply plotting the results as maps.*

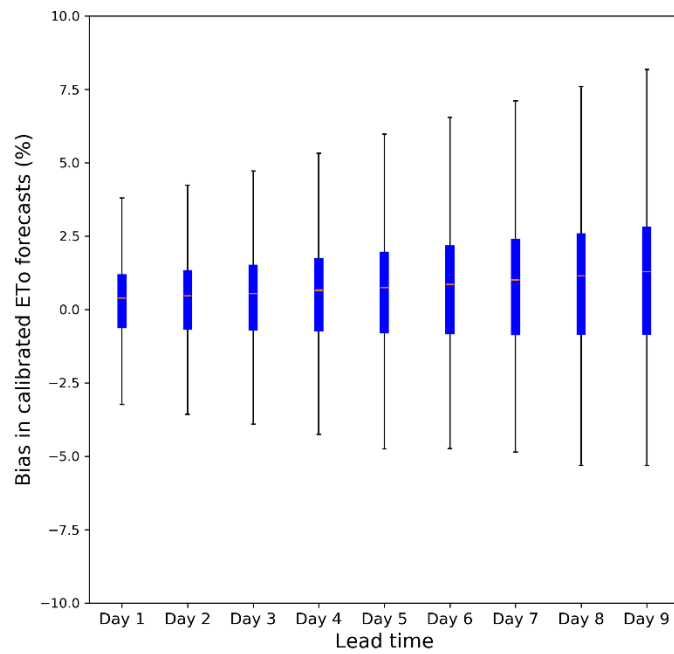
**Response: Thank you for the valuable suggestions. We create boxplots for all the maps shown in the main text. Since we already have 10 figures in the main text and 18 figures in the supplementary material, we think it is better not to add too many new figures. We combine these new boxplots with maps for Figures 2-6 and 8-9, which have extra zoom for adding new subplots. For Figures 1 and 7, which already include many subplots, we present the corresponding boxplots in the Supplementary Material. Please find the boxplots as follows:**



34

35 **Figure 2** Boxplot summarizing improvements in  $r$  in raw ETo forecasts following bias-correction to  
 36 **input variables**

37

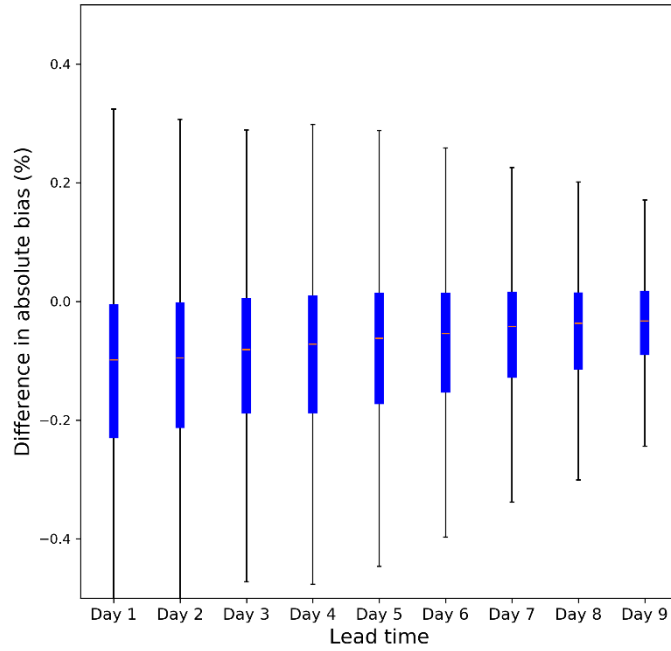


38

39 **Figure 3** Boxplot summarizing bias in calibrated ETo forecasts

40

41

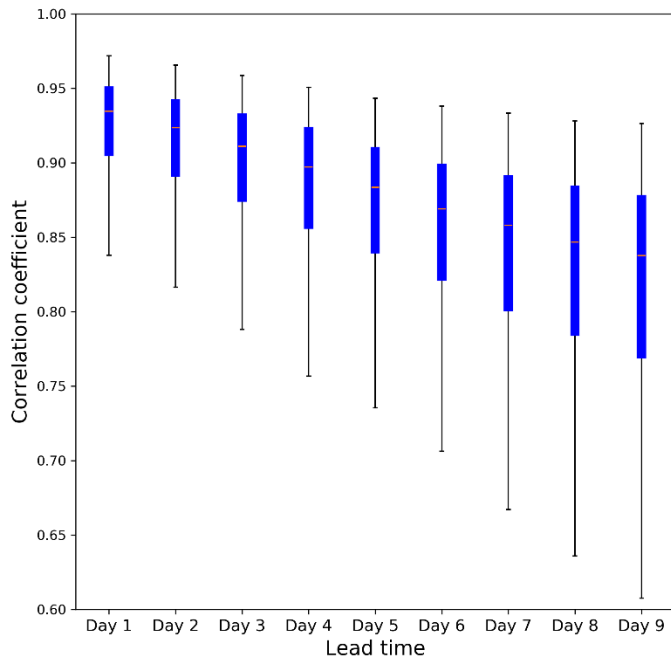


42

43 **Figure 4** Boxplot summarizing differences in absolute bias between calibrated ETo forecasts from  
 44 **Calibration 2 with Calibration 1**

45

46

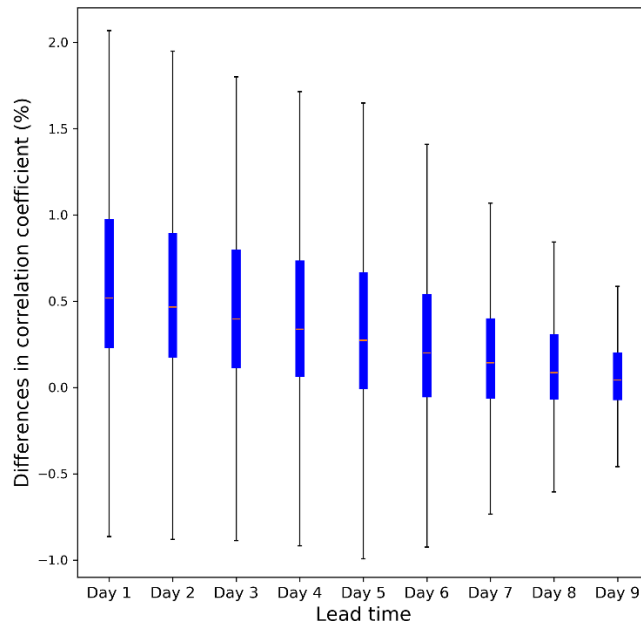


47

48 **Figure 5** Boxplot summarizing correlation coefficient between calibrated ETo forecasts from  
 49 **Calibration 2 and AWAP ETo data**

50

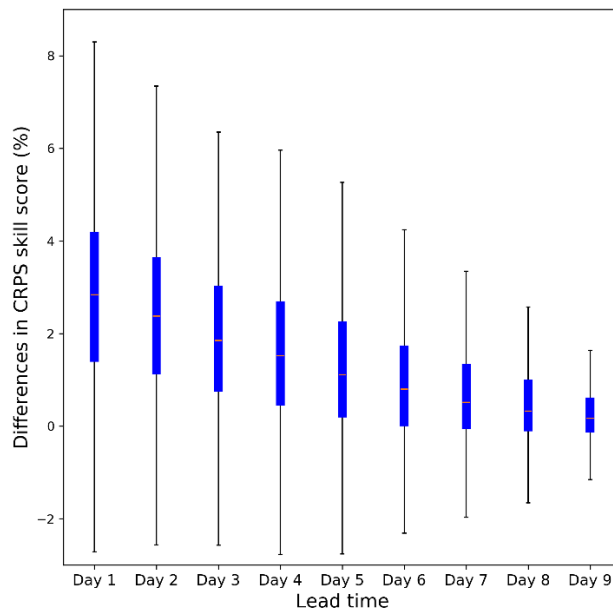
51



52

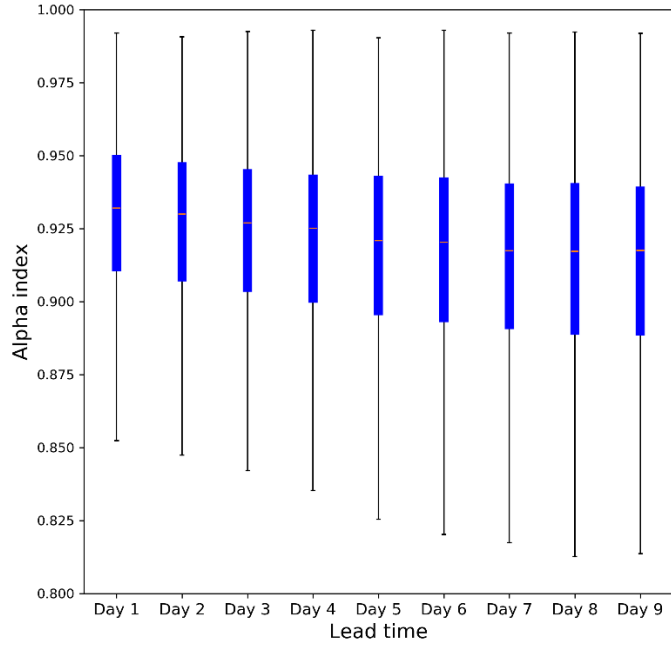
53 **Figure 6 Boxplot summarizing differences in the correlation coefficient (calibrated forecasts vs.**  
54 **AWAP ET<sub>o</sub>) between Calibrations 2 and 1**

55



56

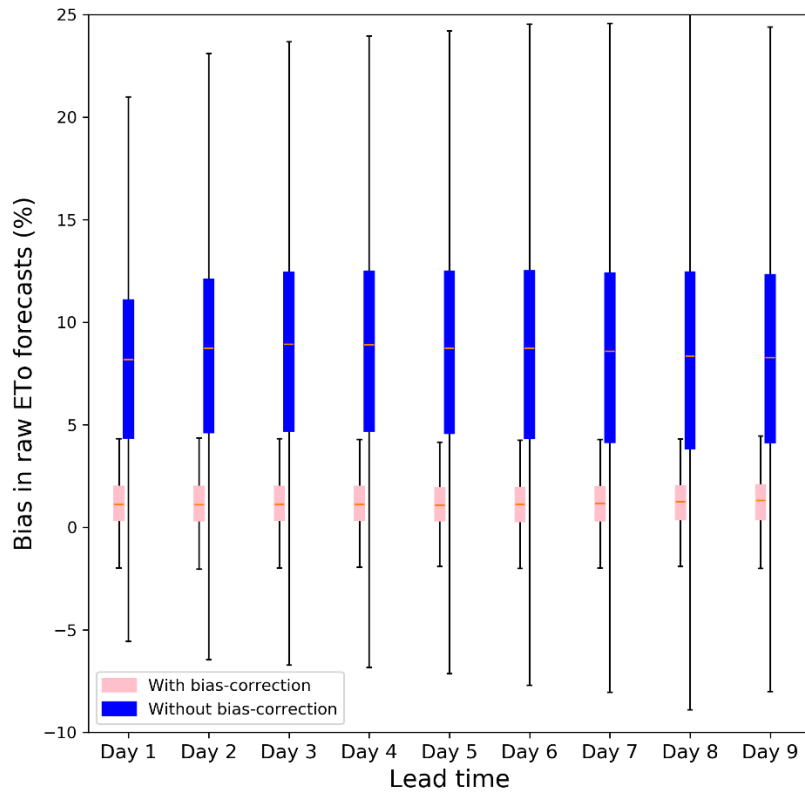
57 **Figure 8 Boxplot summarizing differences in CRPS skill scores between the calibrated forecast**  
58 **from Calibration 2 with those from Calibration 1**



59

60

**Figure 9** Boxplot summarizing the alpha index in the calibrated ETo forecasts



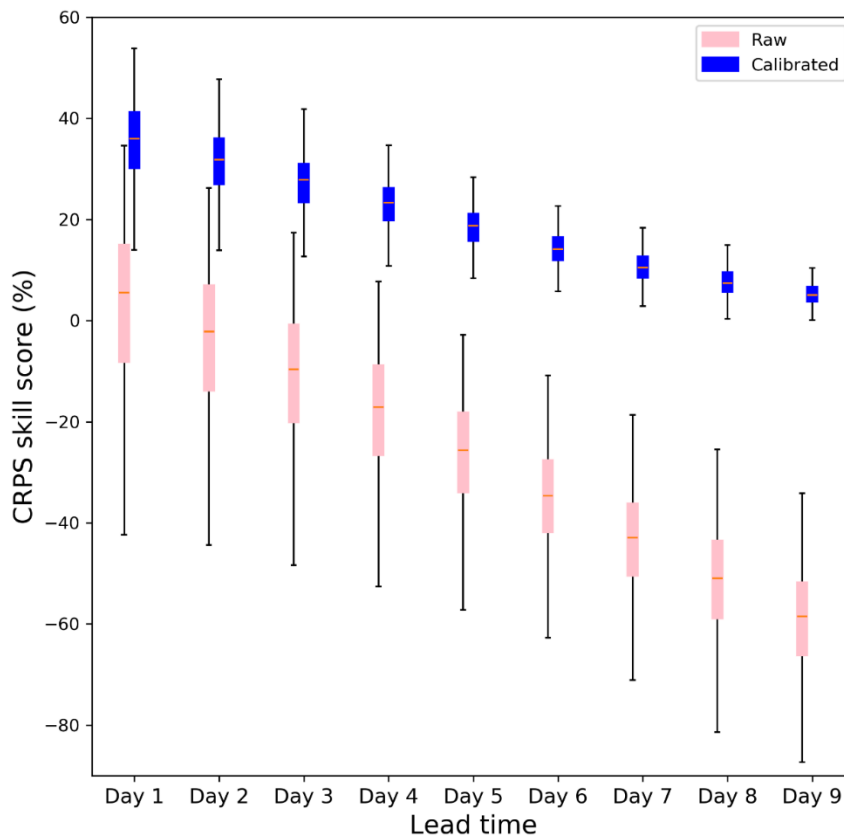
61

62

63

64

*Figure S12. Boxplot of biases in raw ETo forecasts constructed raw (blue) and bias-corrected inputs (pink)*



65

66 *Figure S13. Boxplot of CRPS skill score in raw (pink) and calibrated ETo forecast (blue) from*  
 67 *Calibration 2*

68

69 Point #3

70 *2, Implications for ETo forecasting at the monthly or seasonal scales should be further discussed. ETo*  
 71 *forecasting based on monthly or seasonal climate forecasts from GCMs is also widely performed. This*  
 72 *study develops the new strategy for short-term forecasts. The applicability of this method to ETo*  
 73 *forecasting based on GCM forecasts should be briefly discussed, to benefit a broader range of readers.*

74 **Response: We agree with the reviewer that ETo forecasting with longer forecast horizons**  
 75 **(e.g., monthly and seasonal) based on GCM forecasts is increasingly performed, and it is**  
 76 **necessary to evaluate whether the calibration strategy developed in this investigation is**  
 77 **applicable to the GCM-based seasonal ETo forecasting. As we have shown in this manuscript,**  
 78 **the reduction of error propagation from the input variables to ETo is the key reason why the**  
 79 **new strategy has better performance using raw input variables. We expect this will be the**  
 80 **case for GCM-based seasonal forecasting. However, testing this idea will be beyond the scope**

81 **of this current study. To highlight the necessity of adopting this strategy in seasonal ETo**  
82 **forecasting, we add the following paragraph to section 4.2 (Implications for forecasting of**  
83 **integrated variables and future work):**

84 "The applicability of the calibration strategy developed in this study to seasonal ETo forecasting should  
85 be further investigated. Seasonal ETo forecasting based on GCM climate forecast has been increasingly  
86 performed (Tian et al., 2014; Zhao et al., 2019b). In these investigations, raw ETo forecasts were also  
87 constructed directly with raw GCM climate forecasts. As a result, it is expected that these investigations  
88 have suffered from error propagation from input variables to seasonal ETo forecasts. Whether the  
89 calibration strategy (strategy ii) developed in this study will be applicable to seasonal ETo forecasting  
90 warrants further investigations."

91

#### 92 Point #4

93 *Specific comments:*

94 *Line 20, rewrite this sentence. Not clear*

95 **Response: we replace the original sentence:**

96 "This calibration strategy is expected to enhance future NWP-based ETo forecasting."

97 **with**

98 " We anticipate that future NWP-based ETo forecasting will benefit from adopting the calibration  
99 strategy developed in this study to produce more skillful ETo forecasts."

100

#### 101 Point #5

102 *Line 74 Calibrate->calibrate*

103 **Response: We correct the typo accordingly.**

104

#### 105 Point #6

106 *Line 80 compiled as the inputs.....*

107 **Response: We improve the sentence of:**

108 "Weather forecasts from the ACCESS-G2 model are compiled to generate ETo forecasts."

109 **with:**

110 "Weather forecasts from the Australian Community Climate and Earth System Simulator G2 version  
111 (ACCESS-G2) model are extracted as inputs for the calculation of raw ETo forecasts."

112

113 Point #7

114 *Line 95 10m -> 10 m.*

115 **Response: We add a space between the number and the unit. We also check the entire**  
116 **manuscript to correct the format of units.**

117

118 Point #8

119 *Line 107-108, need to clarify what the anomaly and climatological mean are referring to*

120 **Response: To clarify how the anomaly and climatological mean are derived, we replace the**  
121 **sentence:**

122 "Our recent investigation suggests that ETo forecast calibration based on anomaly and climatological  
123 mean produces more skillful calibrated forecasts than calibrating ETo forecasts directly."

124 **with:**

125 " Our recent investigation suggests calibrating ETo anomalies, which are calculated as departures from  
126 the climatological mean, could produce more skillful calibrated forecasts than calibrating ETo forecasts  
127 directly."

128

129 Point #9

130 *Line 165 consider rewriting this sentence. Does not read well.*

131 **Response: We replace the original sentence of**

132 "Once we obtain all the parameters for the BN distribution (equation 4), a conditional distribution is  
133 established for  $o(t)$  when a raw forecast ( $f(t)$ ) is provided."

134 **with:**

135 " With the optimized parameters (means, standard deviations, and correlations) for the BN distribution  
136 (equation 4), a conditional distribution for  $o(t)$  for a given raw forecast ( $f(t)$ ) is derived."

137

138 Point #10

139 *Line 172, what is specific month*

140 **Response: we replace "specific" with "unselected" to make the wording more specific.**

141

142 Point #11

143 *Figures in Results: shouldn't the figures be centralized?*



144 **Response: The original format following a template from HESS. After we add boxplots to**  
145 **these maps, the empty space for each figure is significantly reduced. We keep them aligned**  
146 **to the left to be consistent with the provided template.**

147

148 Point #12

149 *Line 360, not calibrate directly, should be without correcting forecasts of the inputs*

150 **Response: Thank you for the suggestion. The key message we want to present here is that**  
151 **statistical models may not be able to correct all errors in integrated variables (such as ETo).**  
152 **However, when the input variables are corrected first, error propagation from inputs to**  
153 **integrated variables, particularly for the errors which could not be corrected by calibration**  
154 **models, will be reduced. To make it clear, we improved the original sentence of:**

155 "Our investigation suggests that improving the input variables may help correct errors that could not be  
156 fixed when calibrating the integrated variables directly."

157 **with:**

158 "Our investigation suggests that improving the input variables could effectively reduce error propagation  
159 from inputs to integrated variables. This extra step is proven to be particularly useful in reducing errors in  
160 the integrated variables that could not be corrected through calibration."

161

162 Point #13

163 *Line 365, consider rewriting this sentence*

164 **Response: Thank you for the suggestion. We replace the original sentence:**

165 "As a result, using a more sophisticated calibration method to correct errors in input variables, is expected  
166 to further improve forecasts of these input variables, resulting in more significant improvements in the  
167 final calibrated ETo forecasts."

168 **with:**

169 " If a more sophisticated calibration method is employed to the input variables, error propagation from  
170 input variables to ETo forecasts will likely be further reduced. As a result, we anticipate that the  
171 calibrated ETo forecast will gain further improvements in forecast skills."

172

173 Point #14

174 *Line 377-378, two' calibration models' consider to rewrite*

175 **Response: We improve the original sentence:**

176 "Additional investigations using other calibration models will help clarify whether the improvements will  
177 hold for other calibration models."

178 **With**

179 " Additional evaluations will be needed to verify whether forecast skills will be improved using strategy ii  
180 but based on a different calibration model. "

181

182 Point #15

183 *Line 385, in the calibrated forecasts*

184 **Response: We add the missing 'in' to this sentence.**

185

186 Point #16

187 *Line 386, consider making it shorter and clearer*

188 **Response: We improve the following sentence:**

189 "Further investigation indicates that the contribution of improving input variables to the ETo forecasting  
190 tends to be independent of the calibration method applied to raw ETo forecasts."

191 **With**

192 " Further investigation indicates that the improvements tend to be independent of the calibration method  
193 applied to ETo forecasts."

194

195

196

197

198

199

200

201

202

203

204

205

206

## Responses to Reviewer #2

### 207 Point #1

208 *Comments on “Bias-correcting individual inputs prior to combined calibration leads to more skillful*  
209 *forecasts of reference crop evapotranspiration” by Yang et al. This study evaluated two calibration*  
210 *strategies for simulating reference crop evapotranspiration. The two strategies are (1) calibration*  
211 *directly applied to raw ETo forecast constructed with raw forecast of input variables; (2) bias-correcting*  
212 *input variables. The bias-correcting algorithm has been proved to be more feasible. Although this study is*  
213 *of significance, improvements and revision can make the study stronger and more compelling.*

214 **Response: We appreciate the reviewer's insightful suggestions and comments on the**  
215 **manuscript. We address comments from the reviewer carefully and improve the manuscript**  
216 **accordingly. Please see details in our point-by-point response.**

217

### 218 Point #2

219 *Core of my concerns is the results presentation and discussion, many sections are superficial; the results*  
220 *are simply described, more insightful explanation and discussion are needed. See below for my*  
221 *suggestion. A moderate revision can easily address these comments. So I suggest a moderate revision.*

222 **Response: We appreciate the reviewer's constructive comments. We improve the analysis**  
223 **and presentations by (1) creating boxplots to summarize results plotted as maps to better**  
224 **demonstrate results quantitatively, (2) performing statistical analyses (t-test) when**  
225 **comparing results from different Calibrations, (3) providing more statistical information in the**  
226 **Results section, and (4) Comparing findings of this work with published investigations. We**  
227 **further explain these improvements in detail as follows:**

#### 228 **(1) Adding boxplots to Results**

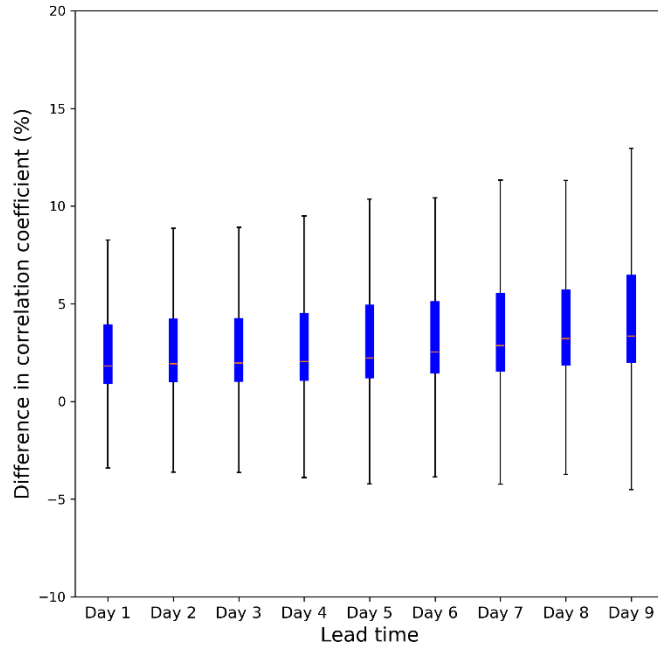
229 **We create boxplots for results shown as maps (Figures 1 to 9 in the main text). We combine**  
230 **these boxplots with maps for Figures 2-6, 8-9, which have extra zoom for adding new**  
231 **subplots. For Figures 1 and 7, which already include many subplots, we present the**  
232 **corresponding boxplots in the Supplementary Material. We also update the main text**  
233 **accordingly. Please find the boxplots as follows:**

234

235

236

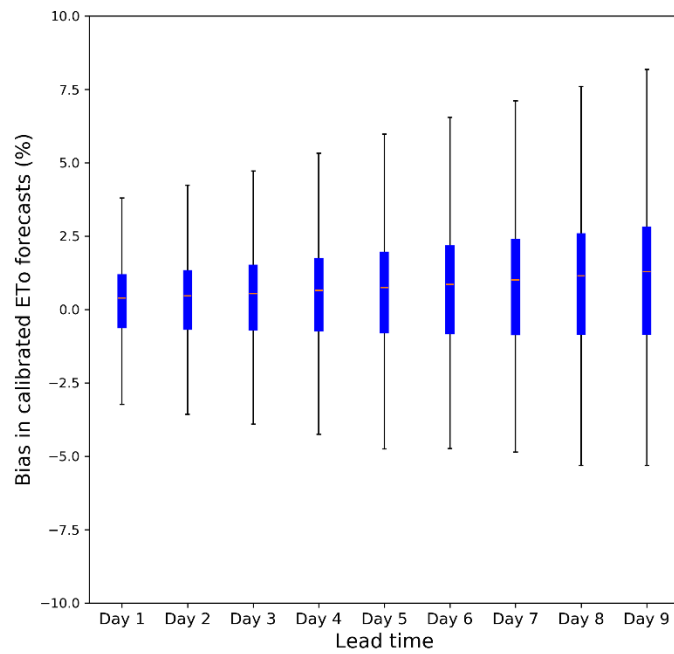
237



238

239 **Figure 2 Boxplot summarizing improvements in  $r$  in raw ETo forecasts following bias-correction to**  
 240 **input variables**

241

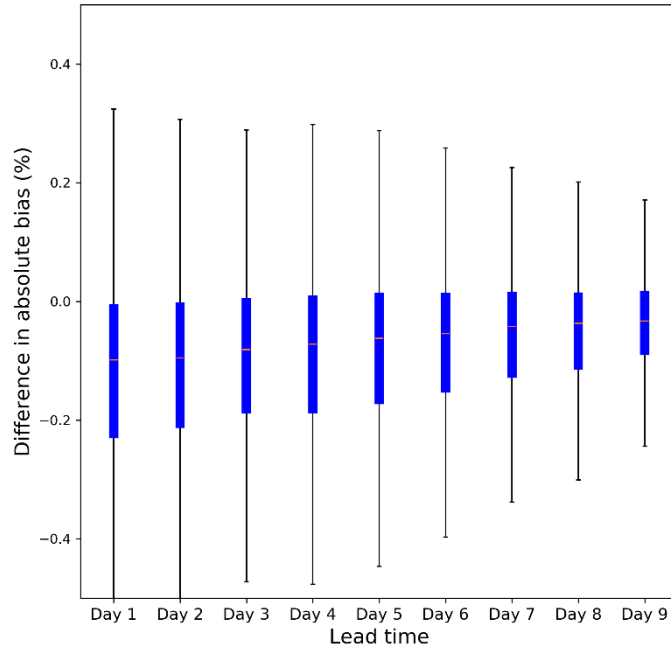


242

243 **Figure 3 Boxplot summarizing bias in calibrated ETo forecasts**

244

245



246

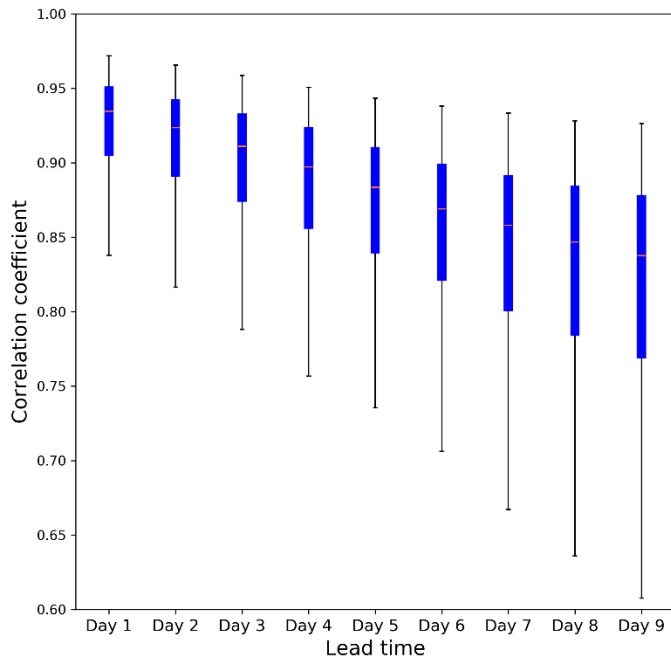
**Figure 4 Boxplot summarizing differences in absolute bias between calibrated ETo forecasts from Calibration 2 with Calibration 1**

247

248

249

250



251

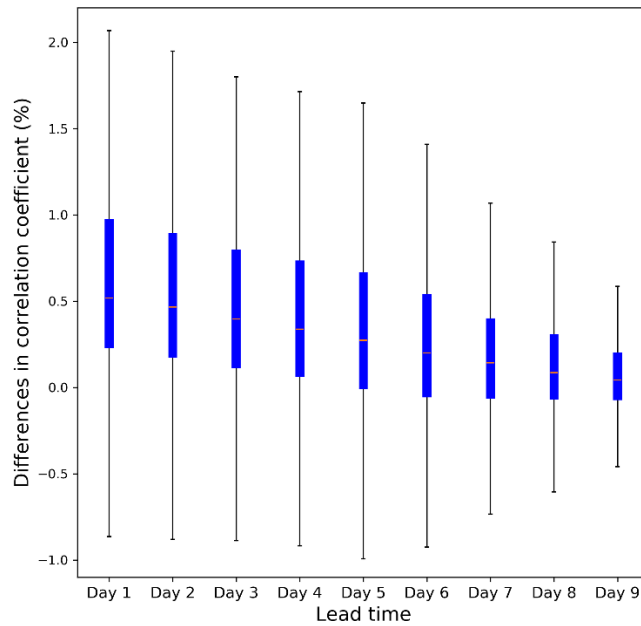
**Figure 5 Boxplot summarizing correlation coefficient between calibrated ETo forecasts from Calibration 2 and AWAP ETo data**

252

253

254

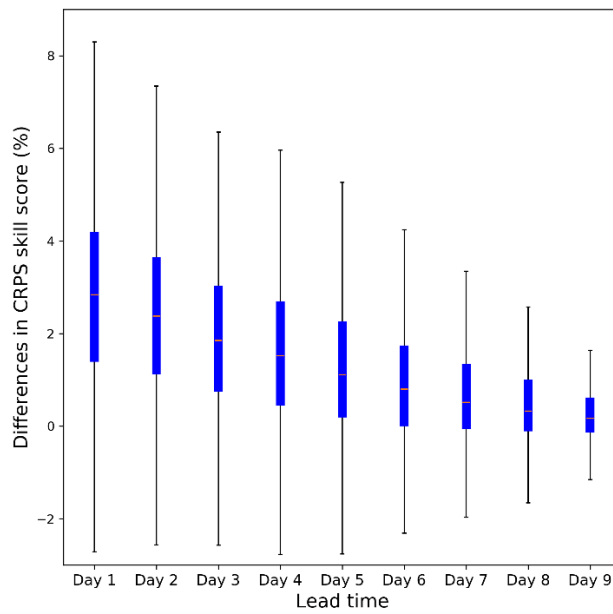
255



256

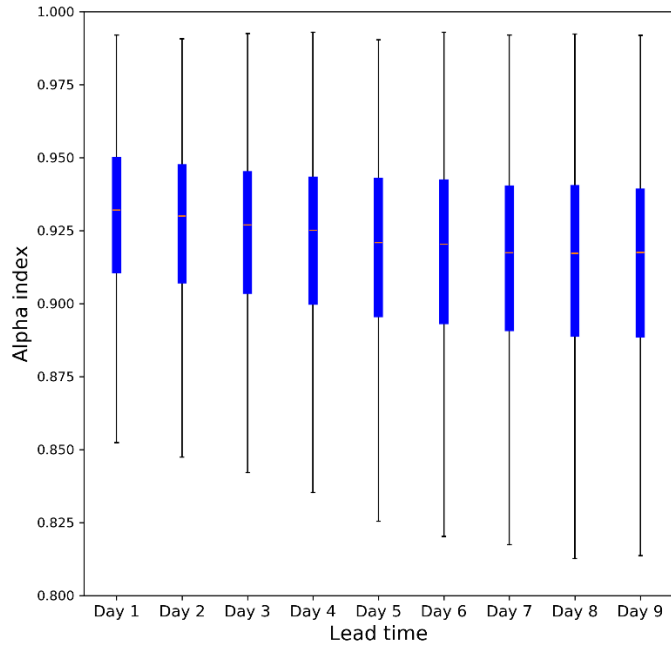
257 **Figure 6** Boxplot summarizing differences in the correlation coefficient (calibrated forecasts vs.  
258 **AWAP ETo) between Calibrations 2 and 1**

259



260

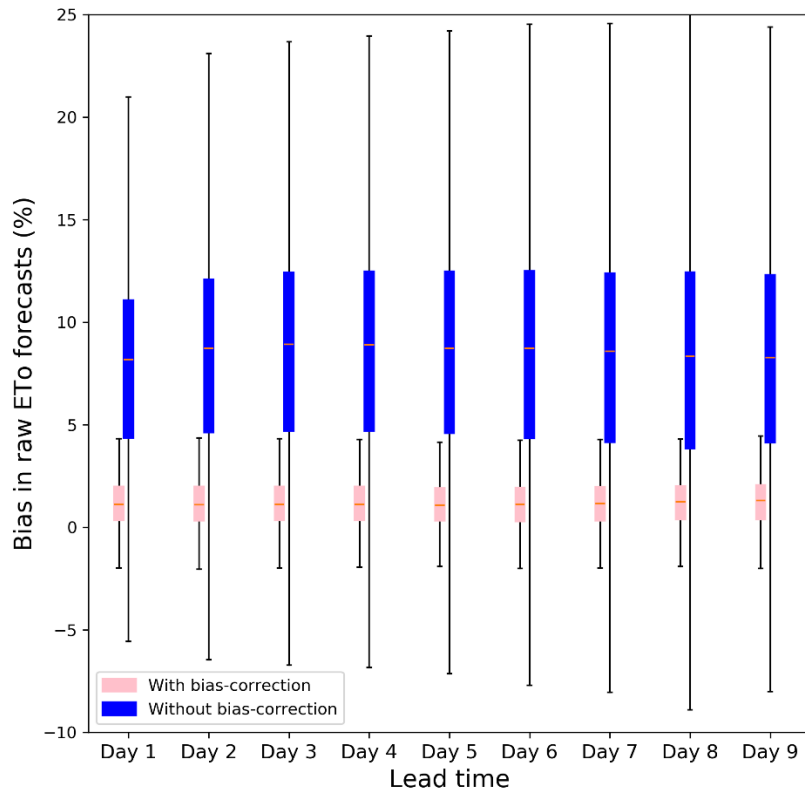
261 **Figure 8** Boxplot summarizing differences in CRPS skill scores between the calibrated forecast  
262 **from Calibration 2 with those from Calibration 1**



263

264

**Figure 9** Boxplot summarizing the alpha index in the calibrated ETo forecasts



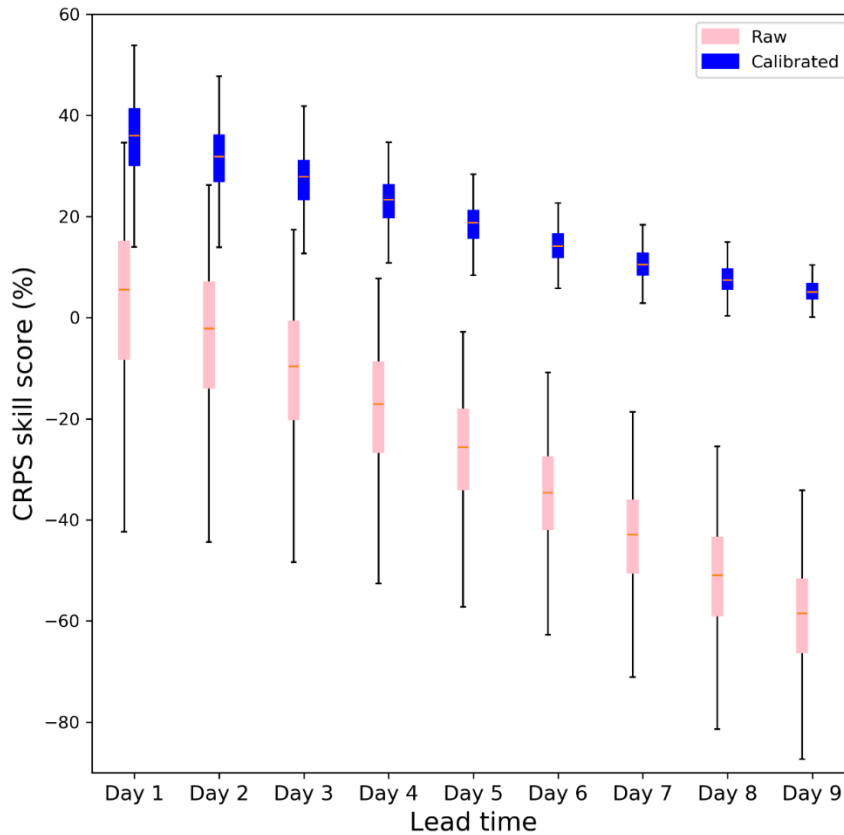
265

266

267

268

*Figure S12. Boxplot of biases in raw ETo forecasts constructed raw (blue) and bias-corrected inputs (pink)*



269

270 *Figure S13. Boxplot of CRPS skill score in raw (pink) and calibrated ETo forecast (blue) from*  
 271 *Calibration 2*

272

273 **(2) Conducting t-tests to compare results from different Calibrations.**

274 **We conduct t-tests (Table S1) to evaluate raw forecasts of the five input variables (Figures S2**  
 275 **to S6). T-tests were also conducted in the evaluation of bias, correlation coefficient, and CRPS**  
 276 **skill score (Figures 1-3, 6-9) of forecasts produced in Calibrations 1 and 2(Table S2).**

277 **In the calculation of *t* statistics, we use the Spatial Degrees of Freedom (SDOF), rather than**  
 278 **using the total grid cells in the study area, to account for the spatial correlation in the t-test.**  
 279 **The SDOF is substantially smaller than total grid cells (Toth, 1995). Wang and Shen (1999)**  
 280 **investigated SDOF of GCM outputs and reported a range of 90-120, out of 738 grid cells for**  
 281 **the southern hemisphere. In this study, we use 50 as the SDOF for our t-tests. Considering the**  
 282 **large amount of total grid cells (281,622) in this study, we believe that 50 is a conservative**  
 283 **estimate of SDOF for this investigation. We calculated the *t*-statistics and evaluate whether**



284 they are statistically significant using the SDOF of 50. Results of the t-tests (Tables S1 and S2)  
285 are added to the supplementary material.

286

287 **Reference:**

288 Toth, Z.: Degrees of freedom in Northern Hemisphere circulation data, *Tellus, Ser. A*, 47 A(4),  
289 457–472, doi:10.3402/tellusa.v47i4.11531, 1995.

290 Wang, X. and Shen, S. S.: Estimation of spatial degrees of freedom of a climate field, *J. Clim.*,  
291 12(5 I), 1280–1291, doi:10.1175/1520-0442(1999)012<1280:EOSDOF>2.0.CO;2, 1999.

292

293

294

295

296

297

298

299

300

301

302

**Table S1 Results of t-tests (*t*-statistic) for raw forecasts of input variables**

Tests Lead times	Test if bias in raw Tmax forecasts is different from zero (Figure S2)	Test if bias in raw Tmin forecasts is different from zero (Figure S3)	Test if bias in raw vapor pressure forecasts is different from zero (Figure S4)	Test if bias in raw solar radiation forecasts is different from zero (Figure S5)	Test if bias in raw wind speed forecasts is different from zero (Figure S6)
Day 1	-8.96**	1.66	-3.18**	11.83**	16.04**
Day 2	-8.16**	2.65**	-3.43**	11.39**	16.50**
Day 3	-8.19**	2.68**	-3.77**	11.81**	16.57**
Day 4	-8.12**	2.56**	-4.05**	12.17**	16.56**
Day 5	-7.87**	2.41**	-4.09**	12.45**	16.45**
Day 6	-7.70**	2.27**	-4.21**	11.88**	16.45**
Day 7	-7.73**	2.22**	-4.33**	10.81**	16.29**
Day 8	-7.70**	2.17**	-4.30**	11.41**	16.56**
Day 9	-7.44**	2.20**	-4.18**	11.95**	16.82**

303 **The Spatial Degrees of Freedom (SDOF) is 50 in the tests; \*\* indicates statistically significant differences at the 95%**  
304 **confidence interval.**

305

306

307

308

309

310

311

312

313

314

315

**Table S2 Results of t-tests (*t*-statistic) for performance evaluation**

Tests Lead times	Comparison of bias in raw ETo forecasts constructed with vs. without bias correction (Figure 1)	Test if r in raw ETo forecasts constructed with raw and bias-corrected input variables are different (Figure 2)	Test if bias in calibrated ETo forecasts from Calibration 2 (Figure 3) is different from zero	Test differences in absolute bias between calibrated ETo forecasts from Calibrations 2 and 1 (Figure 4)	Test difference in <i>r</i> between observations and calibrated ETo forecasts from Calibrations 2 and 1 (Figure 6)	Comparison of CRPS skill score between raw and calibrated ETo forecasts (Figure 7)	Test difference in CRPS skill score of calibrated ETo forecasts from Calibrations 2 and 1 (Figure 8)	Test difference in $\alpha$ -index between Calibrations 2 and 1 (Figure S14)	Test if difference in CRPS skill scores between Calibrations 3 and 4 (Figure S17)
Day 1	-9.76**	7.26**	1.80	-4.08**	5.73**	27.59**	11.53**	-0.54	11.81**
Day 2	-9.86**	7.13**	1.91	-3.93**	4.93**	29.03**	10.86**	-1.47	10.26**
Day 3	-9.86**	7.01**	2.07**	-3.68**	4.43**	31.14**	9.77**	-1.81	9.16**
Day 4	-9.81**	7.04**	2.27**	-3.54**	4.01**	33.77**	8.58**	-1.17	8.33**
Day 5	-9.71**	7.09**	2.40**	-3.36**	3.75**	38.11**	7.16**	-2.09**	7.25**
Day 6	-9.54**	7.33**	2.60**	-3.37**	3.17**	42.59**	6.44**	-1.28	6.66**
Day 7	-9.34**	7.40**	2.76**	-3.26**	2.69**	44.38**	6.15**	-1.99	6.25**
Day 8	-9.04**	7.54**	2.98**	-3.13**	2.32**	45.57**	5.85**	-1.57	5.67**
Day 9	-9.21**	7.50**	3.13**	-2.91**	1.85	51.91**	5.05**	-1.70	4.95**

316 **The Spatial Degrees of Freedom (SDOF) is 50 in the tests; \*\* indicates statistically significant differences at the 95%**  
 317 **confidence interval.**

318

319

320 **(3) Improving the Results section**

321 **We add more specific information in describing the key findings of this study and introduce**  
322 **the results of the statistical analyses (Tables S1 and S2). Since we modified many sentences,**  
323 **we decide not to list them here. Please see details in the revised manuscript.**

324 **(4) Improving the Discussion section**

325 **We further compare the findings of this investigation with existing studies in discussion:**

326 “This investigation further highlights the importance of statistical calibration in NWP-based ETo  
327 forecasting (Medina and Tian, 2020). According to an investigation across 40 sites in Australia, raw ETo  
328 forecasts constructed with NWP outputs reasonably captured the magnitude and variability of ETo, but  
329 forecast skills better than climatology were only limited to the first 6 lead times (Perera et al., 2014). Our  
330 investigation suggests that statistical calibration could substantially improve forecast skills and  
331 successfully extend the skillful forecasts to lead time 9 across Australia. Findings of this investigation  
332 agree well with the site-scale short-term ETo forecasting based on GCM outputs (Zhao et al., 2019a) in  
333 the improvements of forecast skills through statistical calibration. Calibrated forecasts from Calibration 2  
334 demonstrate similar skills as Zhao et al. (2019a) across three Australian sites. Thanks to the capability of  
335 SCC in calibrating short-archived forecasts (Wang et al., 2019), we achieve the improvements based on  
336 much shorter archived raw forecasts (3-year vs. 23-year) than Zhao et al. (2019a). Calibrated forecasts  
337 from Calibration 2 also demonstrate low biases (0.32-0.95%) comparable with calibrated ETo forecasts  
338 (0.49-0.63%) based on the Bayesian Model Averaging (BMA) model and weather forecasts from three  
339 NWP models in the U.S. during 2014-2016 (Medina and Tian, 2020).”

340

341 Point #3

342 Lines 11, fully implemented.

343 **Response: we change it to 'fully implemented '.**

344

345 Point #4

346 *Line 27, “divergent” emphasizes completely different assumption, you can just use replace it*  
347 *different to ensure a general term.*

348 **Response: We replace the word ‘divergent’ with 'different'.**

349

350 Point #5

351 *Line 38, physical processes of the atmosphere, it is unclear, atmospheric circulation or atmospheric wind*  
352 *formation, or physical processes in the atmosphere*

353 **Response: Thank you for the suggestion. We change the sentence as follows:**

354 " ETo is affected jointly by temperature, vapor pressure, solar radiation, and wind speed (Bachour et al.,  
355 2016; Luo et al., 2014). Prediction models using these weather variables as inputs allow for  
356 representations of atmospheric dynamics and often produce reasonable ETo forecasts (Torres et al.,  
357 2011)."

358

### 359 Point #6

360 *Section 3.1, 3.2, the authors described the results in the figures. However, most of those text are vague,*  
361 *please provide more specific (quantitative) information to support your statement. When you compare*  
362 *different results or method, it is better to report some statistic results (p value, r2, etc).*

363 **Response: We appreciate the constructive comments. We conduct statistical analysis to**  
364 **quantify the difference between different model runs, and update the Results section**  
365 **accordingly. Details of the t-tests could be found in our response to your comments point #2.**

366

### 367 Point #7

368 *for example, line line 223-225, you report the overprediction in Tmax, and underpredict in Tmin in*  
369 *different regions. If it is underprediction, what is the range of that underprediction, same for*  
370 *overprediction, are these different statistically significant? There are many similar issues in other*  
371 *sections.*

372 **Response: We appreciate the reviewer's valuable suggestions. We agree with the reviewer**  
373 **that more statistical information is needed. We conduct statistical analysis to quantify errors**  
374 **in raw forecasts (Table S1), and update contents in Results accordingly. Statistical analyses**  
375 **could be found in our response to your comment #2. Here is the updated description of errors**  
376 **in raw forecasts of input variables:**

377 "Raw forecasts of the five input variables demonstrate significant inconsistencies with the  
378 corresponding AWAP data (Figures S2-S6). In most parts of Australia, raw daily maximum  
379 temperature (Tmax) forecasts are lower than AWAP data by 1-2 °C. Overpredictions in Tmax  
380 are only found in coastal areas of northwestern Australia. The daily minimum temperature  
381 (Tmin) is underpredicted by more than 1.5 °C in western and central parts of Australia by the  
382 raw forecasts, but is overpredicted by ca. 1 °C in eastern and southern Australia. Vapor pressure  
383 is underpredicted in western and central regions by ca.14%, but is overpredicted by ca. 6% in  
384 coastal areas of southeastern Australia by the raw forecasts. Raw solar radiation forecasts are  
385 about 5% higher than AWAP data across Australia. Forecasted wind speed is higher than the  
386 reference data by more than 1 m s-1 (or by ca. 63%) in most parts of Australia. For each input  
387 variable, spatial patterns of biases in raw forecasts are consistent across the 9 lead times,  
388 demonstrating systematic errors in the raw NWP forecasts. According to our statistical test,  
389 overpredictions or underpredictions in raw forecasts of the input variables are statistically  
390 significant (P<0.05) for most lead times (Table S1)."

391 Point #8

392 *In the discussion section, I would be willing to see a comparison with other studies with different*  
393 *algorithms for the ETo simulation. Some quantitative comparison to elucidate the better performance of*  
394 *the new bias-correction algorithm needs to be done. I believe it will prove the reliability of the new*  
395 *algorithm.*

396 **Response: We appreciate the constructive comments. This is the first continental-scale ETo**  
397 **forecasting in Australia. Previous NWP/GCM-based ETo forecasting in Australia is conducted**  
398 **at the site scale. As a result, in the original manuscript, our evaluation was primarily focused**  
399 **on the comparison against observations. In this area of weather/climate forecasting, different**  
400 **calibration models, based on different statistical theories have been developed and**  
401 **implemented. Previous comparisons suggest that the performance of these models varied**  
402 **with study areas, NWP models, and choice of evaluation metrics (Wilks, 2018), and there is**  
403 **no conclusion regarding which group of post-processing models has the best performance.**

404 **More importantly, rather than developing a new calibration model, this investigation is to**  
405 **evaluate the necessity of including an extra step before ETo forecasts are calibrated. As we**  
406 **introduced in the main text, the objective of our investigations is to address a challenge**  
407 **commonly faced by NWP-based ETo forecasting. We expect the calibration strategy**  
408 **developed in this study will benefit ETo forecast calibrations broadly, no matter which**  
409 **statistical model is employed in ETo forecast calibration.**

410 **However, we agree with the reviewer that comparison of model performance with other**  
411 **models will help readers better understand the robustness of our calibration. We review**  
412 **previous studies and add the following content to the Discussion section (4.1):**

413 “According to an investigation across 40 sites in Australia, raw ETo forecasts constructed with NWP  
414 outputs reasonably captured the magnitude and variability of ETo, but forecast skills better than  
415 climatology were only limited to the first 6 lead times (Perera et al., 2014). Our investigation suggests  
416 that statistical calibration could substantially improve forecast skills and successfully extend the skillful  
417 forecasts to lead time 9 across Australia. Findings of this investigation agree well with the site-scale  
418 short-term ETo forecasting based on GCM outputs (Zhao et al., 2019a) in the improvements of forecast  
419 skills through statistical calibration. Calibrated forecasts from Calibration 2 demonstrate similar skills as  
420 Zhao et al. (2019a) across three Australian sites. Thanks to the capability of SCC in calibrating short-  
421 archived forecasts (Wang et al., 2019), we achieve the improvements based on much shorter archived raw  
422 forecasts (3-year vs. 23-year) than Zhao et al. (2019a). Calibrated forecasts from Calibration 2 also  
423 demonstrate low biases (0.32-0.95%) comparable with calibrated ETo forecasts (0.49-0.63%) based on  
424 the Bayesian Model Averaging (BMA) model and weather forecasts from three NWP models in the U.S.  
425 during 2014-2016 (Medina and Tian, 2020).”

426 **In addition, we also highlight the importance of testing the proposed calibration strategy**  
427 **(strategy ii) based on other calibration models in the future in section 4.2:**

428 “Third, further investigations based on other calibration models are needed to validate findings of this  
429 investigation. Our analyses based on two different methods (based on ETo anomalies vs. based on  
430 original ETo) demonstrate similar improvements in calibrated ETo forecasts with the adoption of bias-

431 correction to input variables. Additional evaluations will be needed to verify whether forecast skills will  
432 be improved using strategy ii but based on a different calibration model.”

433

434 **Reference:**

435 Wilks, D.S., 2018. Chapter 3. Univariate Ensemble Forecasting, in: Vannitsem, S., Wilks, D.S., Messner,  
436 J.W. (Eds.), Statistical Postprocessing of Ensemble Forecasts. pp. 49–89.  
437 <https://doi.org/https://doi.org/10.1016/C2016-0-03244-8>

438

439 Point #9

440 *Line 388, feasible or reliable ETo forecasting.*

441 **Response: This paragraph has been rewritten. Please see the revised contents in our response**  
442 **to your comment #10.**

443

444 Point #10

445 *Line 390, short-term ETo forecasting provides highly valuable information for real-time decision making*  
446 *on water resource management and planning farming practices. This study proved the bias-correction*  
447 *approach is a feasible method for a more robust calibration of the NWP-based ETo forecasting.*

448 **Response: We appreciate the reviewer's valuable suggestions. We remove redundant**  
449 **sentences and combine the last two paragraphs in the Conclusion section:**

450 " This investigation clearly suggests the necessity of improving input variables as part of ETo forecast  
451 calibration. With this extra step, the bias, correlation coefficient, and skills of the calibrated ETo forecasts  
452 are all improved. Further investigation indicates that the improvements tend to be independent of the  
453 calibration method applied to ETo forecasts. Forecasting the highly variable ETo is often challenging.  
454 This investigation addresses a common challenge in NWP-based ETo forecasting and develops an  
455 effective calibration strategy for adding extra skills to ETo forecasts. We anticipate that future NWP-  
456 based ETo forecasting could benefit from adopting this strategy to produce more skillful calibrated ETo  
457 forecasts. This strategy is also expected to be applicable to enhancing the forecasting of other integrated  
458 variables that are calculated using multiple NWP/GCM variables as inputs."

459

460

461

462

463

464

## Responses to Reviewer #3

465 Point #1

466 *Author(s): Qichun Yang et al.*

467 *MS No.: hess-2021-69*

468 *This paper focuses on the comparison of two calibration strategies to provide short-term reference crop*  
469 *evapotranspiration (ETo). ETo forecasting is still a relatively new area of research, in Australia and*  
470 *elsewhere, and has received more attention in the past few years. Skilful ETo forecasts in Australia would*  
471 *help support efficient water use and water management. Two strategies to calibrate ETo forecasts have*  
472 *emerged: i) the calibration of raw ETo forecasts and ii) bias-correcting input variables first before*  
473 *calibrating ETo forecasts. Little work to date compares the two approaches, it is unclear which method*  
474 *might be more advantageous or skilful. This paper therefore addresses a topical subject with a large*  
475 *audience interest.*

476 *I have some reservations regarding some methodological choices and justifications (purpose and*  
477 *inclusion of experiment 3 and 4), as well as a lack of interpretations of the results overall. I recommend*  
478 *revision to strengthen this paper.*

479 **Response: Thank you for the valuable suggestions and careful review. We revise this work**  
480 **carefully based on your constructive suggestions.**

481

482 Point #2

483 *The authors re-grid the weather forecast variables of ACCESS-G2 to match the timeframe and resolution*  
484 *of the gridded data AWAP. They perform four experiments: experiments 1) and 2) are based on the ETo*  
485 *anomaly and climatological mean, whereas experiment 3 and 4) use the ETo values directly.*  
486 *Furthermore, experiment 1) and 3) use raw inputs to calculate and calibrate ETo forecasts whereas*  
487 *experiments 2) and 4) first bias-correct inputs before ETo calibration. The SCC calibration method is used*  
488 *for ETo forecast while a quantile mapping method is used to bias-correct input forecasts. The authors*  
489 *evaluate the forecasts using three metrics for the theoretical assessment of bias, reliability and accuracy.*  
490 *Overall results suggest that the second strategy (bias-correction of inputs before ETo calibration)*  
491 *provides more skilful forecasts.*

492 **Response: We appreciate the reviewer's thorough review. The work has been substantially**  
493 **improved through addressing the valuable comments.**

494

495

496



497 Point #3

498 *Major comments:*

499 *Methodology:*

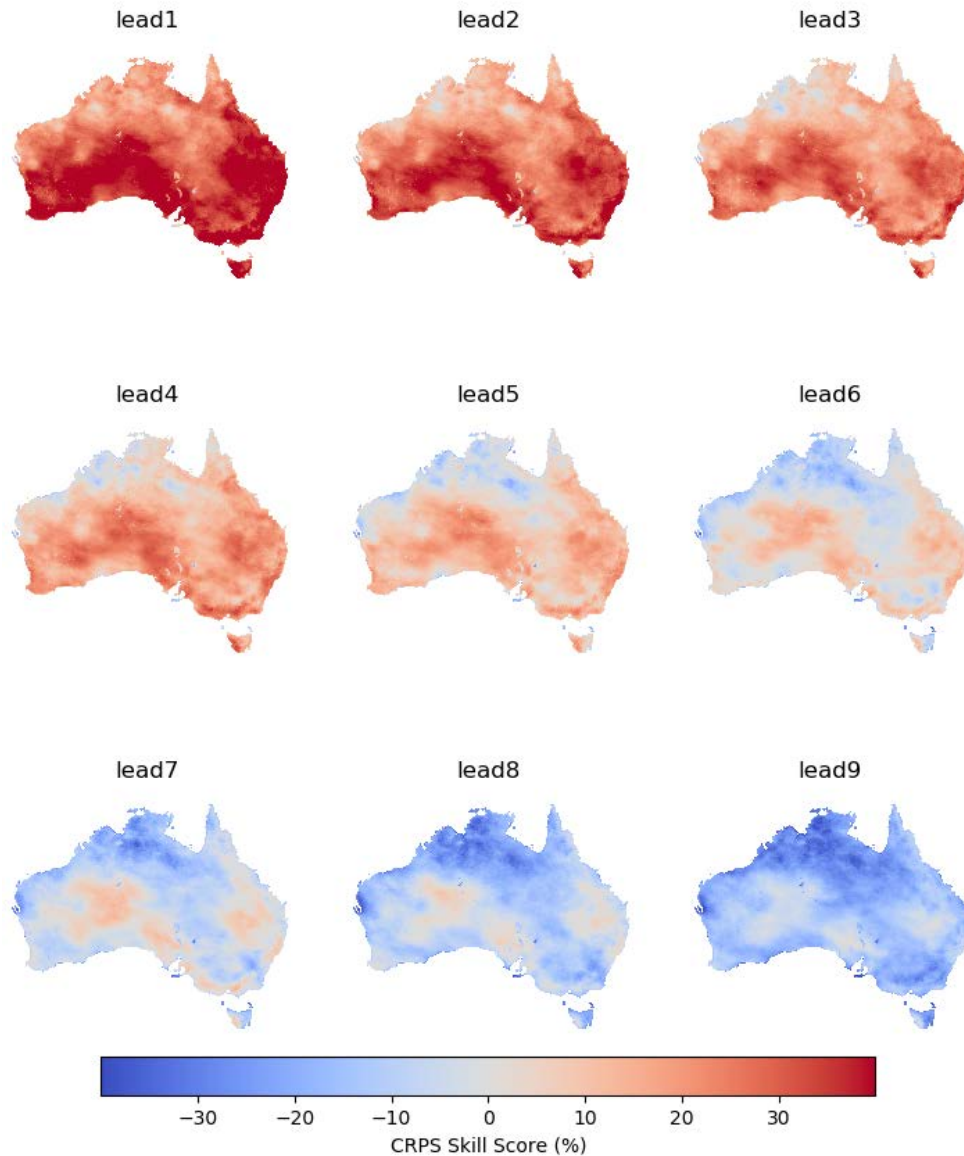
500 *P4 section 2.3: Why not compare the calibration method used SCC to other methods tested in the*  
501 *literature which would enable to place this work in context to other studies on ETo forecasting?*

502 **Response: We appreciate the constructive comments. We understand that comparing the**  
503 **performance of SCC with existing methods will help readers better understand the strengths**  
504 **of our methodology in ETo forecasting. We did not compare the calibration based on SCC**  
505 **model directly with other models in the original submission for a couple of reasons:**

506 **First, the primary objective of this investigation is to address a common challenge faced by**  
507 **NWP-based ETo forecasting, rather than to develop a new calibration model. As a result, we**  
508 **primarily focus on evaluating the necessity of correcting forecasts of input variables prior to**  
509 **calibrating ETo forecasts. As we introduced in the main text, the developed calibration**  
510 **strategy is expected to benefit ETo forecast calibrations broadly, rather than improving an**  
511 **individual calibration model. As suggested by the model experiments (Calibrations 1-4), the**  
512 **developed strategy could be applicable to other calibration models.**

513 **Second, we feel it is not necessary to compare the performance of SCC against calibration**  
514 **models, which are widely used but less sophisticated models. Simple calibration models, such**  
515 **as quantile mapping (QM), have been widely used in calibrating hydroclimate forecasts.**  
516 **These models are often readily available, or could be easily coded and implemented.**  
517 **However, the limitations of these models in forecast calibration have been reported (Zhao et**  
518 **al., 2017). When we started this investigation, we used quantile mapping to calibrate ETo**  
519 **forecasts (raw ETo forecasts constructed with raw forecasts of input variables). As**  
520 **demonstrated in the following figure, the CRPS skill score of quantile mapped ETo forecasts is**  
521 **not only lower than the SCC-calibrated forecasts for each corresponding lead time (Figure 7),**  
522 **but also becomes negative (worse than climatological forecasts) in parts of Australia starting**  
523 **from lead time 4. As a result, calibration of ETo forecasts with quantile mapping further**  
524 **confirms the limitations of this model. Therefore, using such models as a reference to**  
525 **evaluate the performance of SCC is not necessary since their limitations have been reported.**  
526 **As a result, we decide not to include a comparison with quantile mapping in this manuscript.**

527



528

529

***CRPS skill score of calibrated ETo forecasts using Quantile Mapping***

530

531 **Third, we have limited access to sophisticated calibration models. There is no global post-**  
 532 **processing software library archiving these models. We found it was hard to access the**  
 533 **source code of these models and to directly compare SCC with them. In addition, previous**  
 534 **comparisons suggest that the performance of these models varied with study areas, NWP**  
 535 **models, and choice of evaluation metrics (Wilks, 2018), and there is no conclusion regarding**  
 536 **which group of post-processing models has the best performance. Our indirect comparison**  
 537 **with other models confirms this conclusion. Details will be presented in the following**  
 538 **paragraphs.**

539 **Fourth, the short-achieved NWP forecasts (3-year) used in this study represent a challenge for**  
540 **conducting the calibration using other models. Many calibration models, particularly those**  
541 **based on models of the joint probability of forecasts and observations (Krzysztofowicz and**  
542 **Herr, 2001; Wang and Robertson, 2011), require long hindcasts (20-30 years) to establish a**  
543 **joint distribution to link observations and forecasts. Applying such models to short-archived**  
544 **forecasts such as those used in this study will substantially undermine the statistical**  
545 **assumption of these models. In contrast, the SCC model has been developed specifically to**  
546 **address the challenge associated with short-archived forecasts. The advantages of SCC in**  
547 **calibrating short-archived forecasts have been explained in our recent publications (Wang et**  
548 **al., 2019; Yang et al., 2021).**

549 **As a result, we decide not to compare SCC directly with other models. However, we totally**  
550 **agree with the reviewer that comparison of model performance with other models will help**  
551 **readers better understand the performance of our calibration. As a result, we extract our**  
552 **results at three Australia sites where ETo forecasts were also calibrated based on the**  
553 **Bayesian Joint Probability (BJP) model (Zhao et al., 2019), and compare the results of the two**  
554 **investigations. In addition, we also compare our results with site-scale investigations in other**  
555 **regions of Australia. We also compare results of this study with investigations in the U.S. We**  
556 **add the following paragraph to discuss findings of our work relative to existing investigations**  
557 **to the Discussion section (4.1):**

558 "This investigation further highlights the importance of statistical calibration in NWP-based ETo  
559 forecasting (Medina and Tian, 2020). According to an investigation across 40 sites in Australia,  
560 raw ETo forecasts constructed with NWP outputs reasonably captured the magnitude and  
561 variability of ETo, but forecast skills better than climatology were only limited to the first 6 lead  
562 times (Perera et al., 2014). Our investigation suggests that statistical calibration could  
563 substantially improve forecast skills and successfully extend the skillful forecasts to lead time 9  
564 across Australia. Findings of this investigation agree well with the site-scale short-term ETo  
565 forecasting based on GCM outputs (Zhao et al., 2019a) in the improvements of forecast skills  
566 through statistical calibration. Calibrated forecasts from Calibration 2 demonstrate similar skills  
567 as Zhao et al. (2019a) across three Australian sites. Thanks to the capability of SCC in  
568 calibrating short-archived forecasts (Wang et al., 2019), we achieve the improvements based on  
569 much shorter archived raw forecasts (3-year vs. 23-year) than Zhao et al. (2019a). Calibrated  
570 forecasts from Calibration 2 also demonstrate low biases (0.32-0.95%) comparable with  
571 calibrated ETo forecasts (0.49-0.63%) based on the Bayesian Model Averaging (BMA) model  
572 and weather forecasts from three NWP models in the U.S. during 2014-2016 (Medina and Tian,  
573 2020)."

574 **In addition, we also highlight the importance of further testing the proposed calibration**  
575 **strategy (strategy ii) based on other calibration models. We add the following contents to**  
576 **section 4.2:**

577 “Third, further investigations based on other calibration models are needed to validate findings  
578 of this investigation. Our analyses based on two different methods (based on ETo anomalies vs.  
579 based on original ETo) demonstrate similar improvements in calibrated ETo forecasts with the  
580 adoption of bias-correction to input variables. Additional evaluations will be needed to verify  
581 whether forecast skills will be improved using strategy ii but based on a different calibration  
582 model.”

583 **Reference:**

- 584 Medina, H. and Tian, D.: Comparison of probabilistic post-processing approaches for improving  
585 numerical weather prediction-based daily and weekly reference evapotranspiration forecasts,  
586 Hydrol. Earth Syst. Sci., 24, 1011–1030, 2020.
- 587 Perera, K. C., Western, A. W., Nawarathna, B. and George, B.: Forecasting daily reference  
588 evapotranspiration for Australia using numerical weather prediction outputs, Agric. For. Meteorol.,  
589 194, 50–63, doi:10.1016/j.agrformet.2014.03.014, 2014.
- 590 Wilks, D.S., 2018. Chapter 3. Univariate Ensemble Forecasting, in: Vannitsem, S., Wilks, D.S., Messner,  
591 J.W. (Eds.), Statistical Postprocessing of Ensemble Forecasts. pp. 49–89.  
592 <https://doi.org/https://doi.org/10.1016/C2016-0-03244-8>
- 593 Krzysztofowicz, R., Herr, H.D., 2001. Hydrologic uncertainty processor for probabilistic river stage  
594 forecasting: precipitation-dependent model. J. Hydrol. 249, 46–68.
- 595 Wang, Q.J., Robertson, D.E., 2011. Multisite probabilistic forecasting of seasonal flows for streams with  
596 zero value occurrences. Water Resour. Res. 47, 1–19. <https://doi.org/10.1029/2010WR009333>
- 597 Wang, Q.J., Zhao, T., Yang, Q., Robertson, D., 2019. A Seasonally Coherent Calibration ( SCC ) Model for  
598 Postprocessing Numerical Weather Predictions. Mon. Weather Rev. 147, 3633–3647.  
599 <https://doi.org/10.1175/MWR-D-19-0108.1>
- 600 Yang, Q., Wang, Q.J., Hakala, K., 2021. Achieving effective calibration of precipitation forecasts over a  
601 continental scale. J. Hydrol. Reg. Stud. 35, 100818. <https://doi.org/10.1016/j.ejrh.2021.100818>
- 602 Zhao, T., Wang, Q.J., Schepen, A., 2019. A Bayesian modelling approach to forecasting short-term  
603 reference crop evapotranspiration from GCM outputs. Agric. For. Meteorol. 269–270, 88–101.  
604 <https://doi.org/10.1016/j.agrformet.2019.02.003>

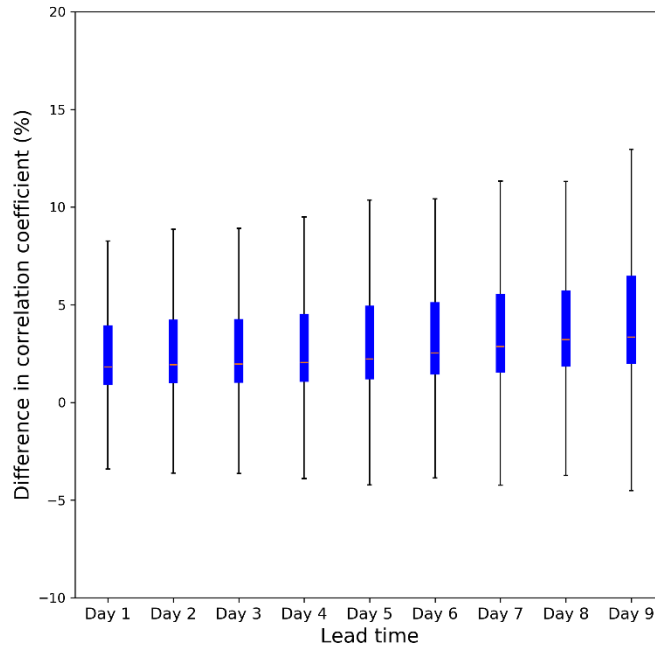
605

606 **Point #4**

607 *Presentation of summary statistics. Why not use boxplots to present overall statistics and across lead*  
608 *times (for example next to figure 4 and so on)? Reliability diagrams for particular ETo thresholds would*  
609 *be helpful to communicate when the forecasts are reliable.*

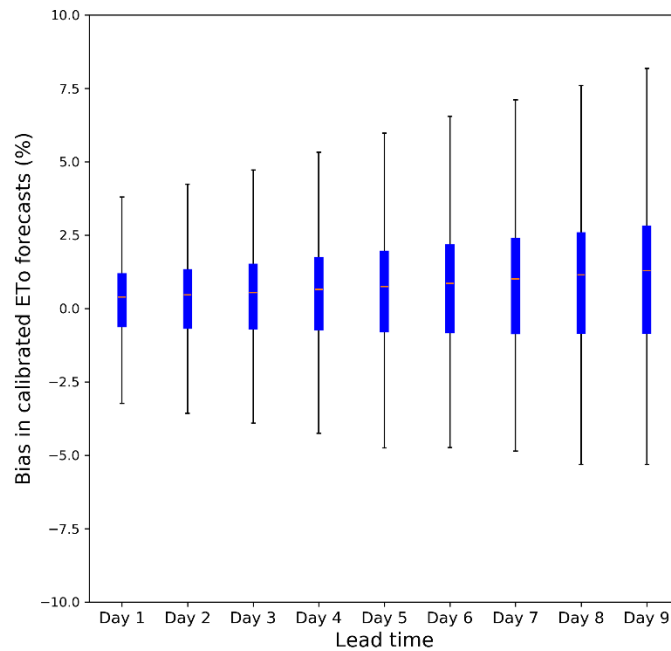
610 **Response: Thank you for the constructive suggestions. We created boxplots for results**  
611 **shown as maps (Figures 1 to 9 in the main text). For Figures 1 and 7, which already include**  
612 **many subplots, we present the corresponding boxplots in the Supplementary Material. For**  
613 **other map figures (Figures 2-6, and 8-9), which have extra zoom for adding new subplots, we**

614 combine boxplots with the maps. We also update the main text accordingly. Please find the  
615 boxplots as follows:



616

617 **Figure 2** Boxplot summarizing improvements in  $r$  in raw ETo forecasts following bias-correction to  
618 input variables

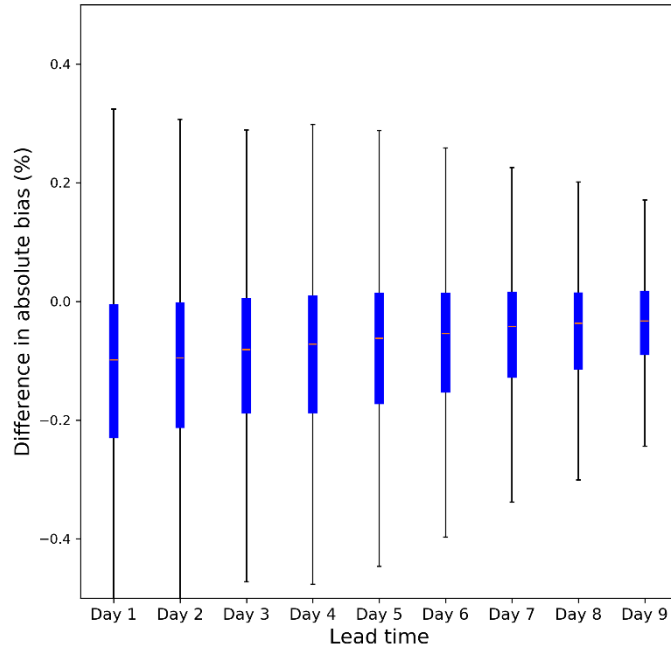


619

620 **Figure 3** Boxplot summarizing bias in calibrated ETo forecasts

621

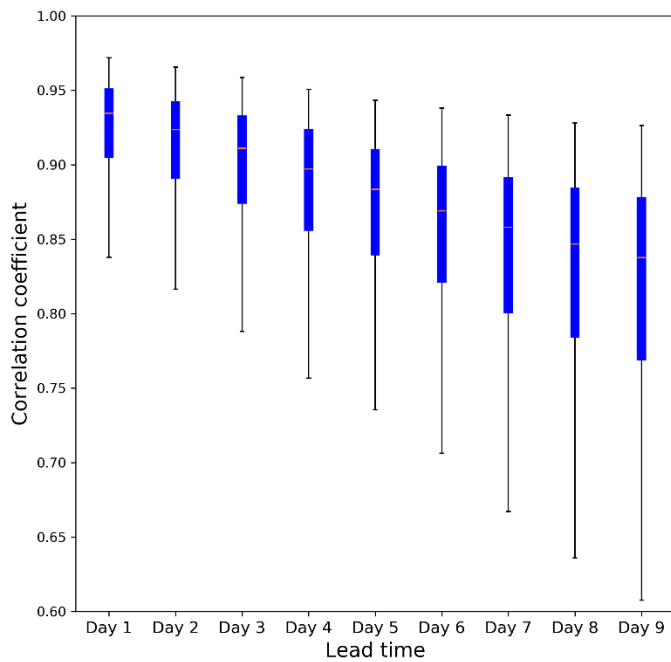
622



623

624 **Figure 4** Boxplot summarizing differences in absolute bias between calibrated ETo forecasts from  
 625 **Calibration 2 with Calibration 1**

626

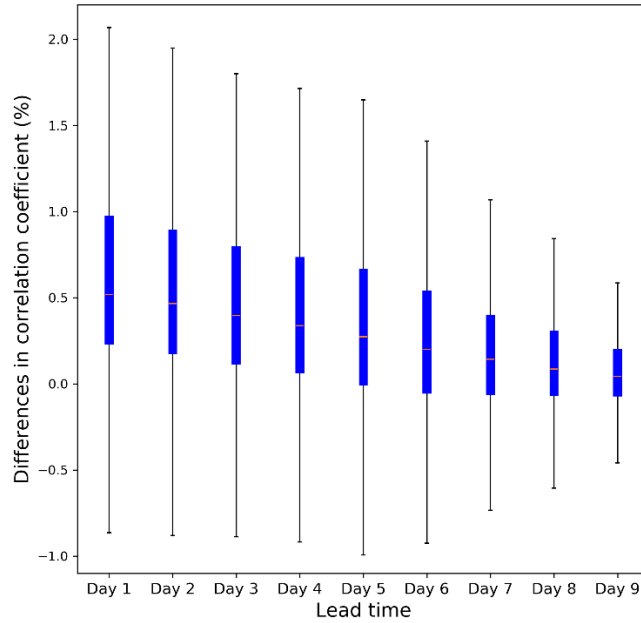


627

628 **Figure 5** Boxplot summarizing correlation coefficient between calibrated ETo forecasts from  
 629 **Calibration 2 and AWAP ETo data**

630

631



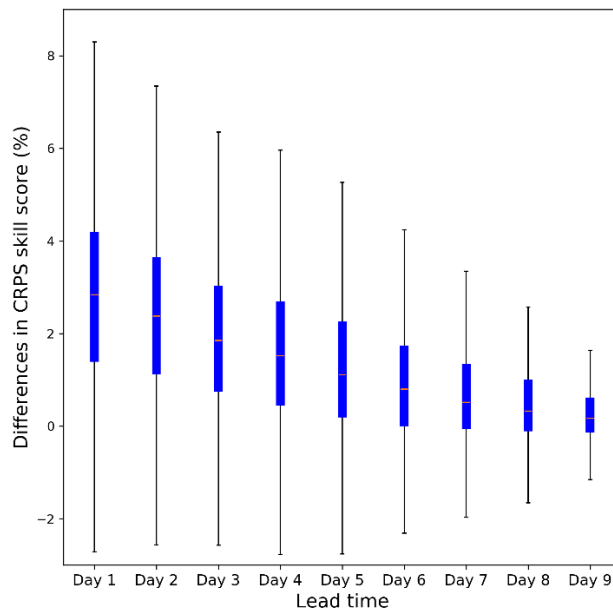
632

**Figure 6 Boxplot summarizing differences in the correlation coefficient (calibrated forecasts vs. AWAP ET<sub>o</sub>) between Calibrations 2 and 1**

633

634

635

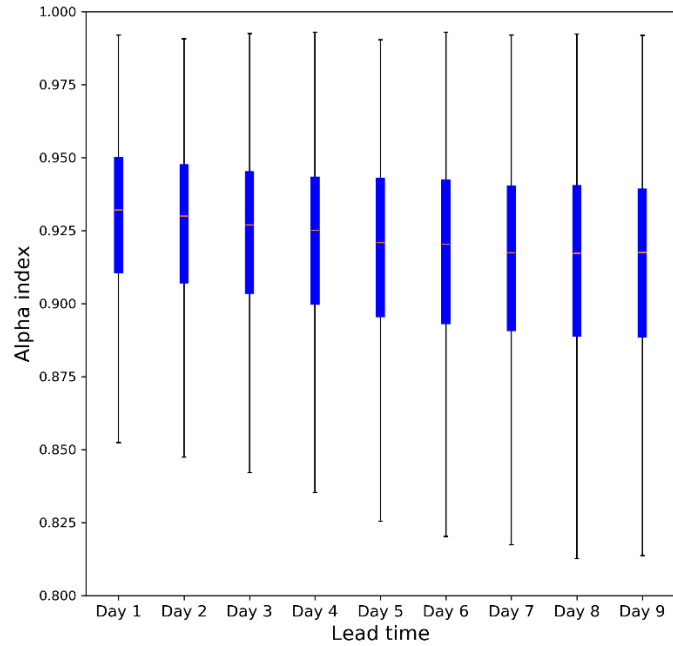


636

**Figure 8 Boxplot summarizing differences in CRPS skill scores between the calibrated forecast from Calibration 2 with those from Calibration 1**

637

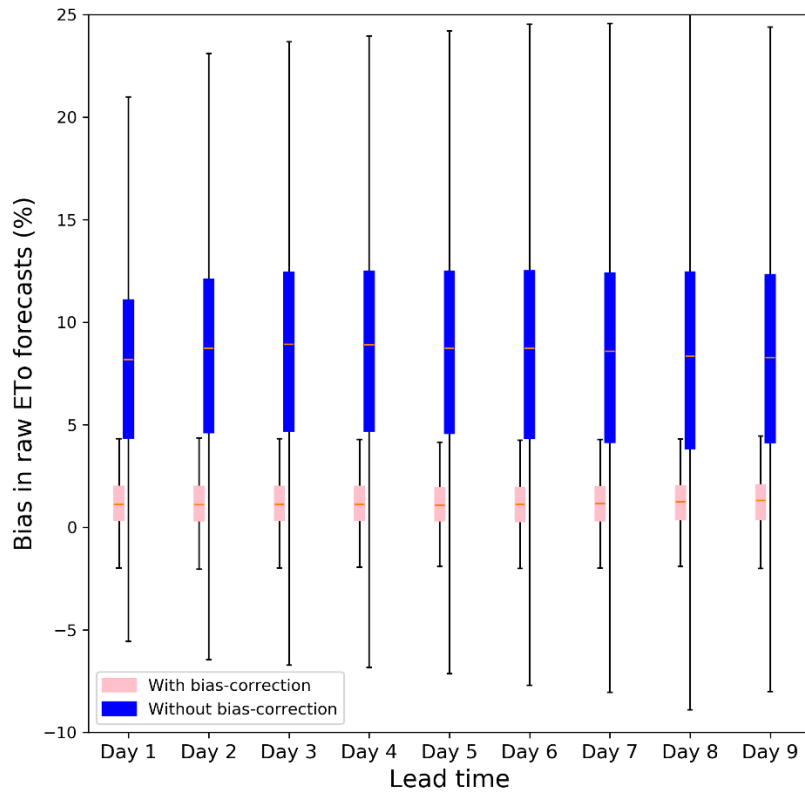
638



639

640

**Figure 9** Boxplot summarizing the alpha index in the calibrated ETo forecasts



641

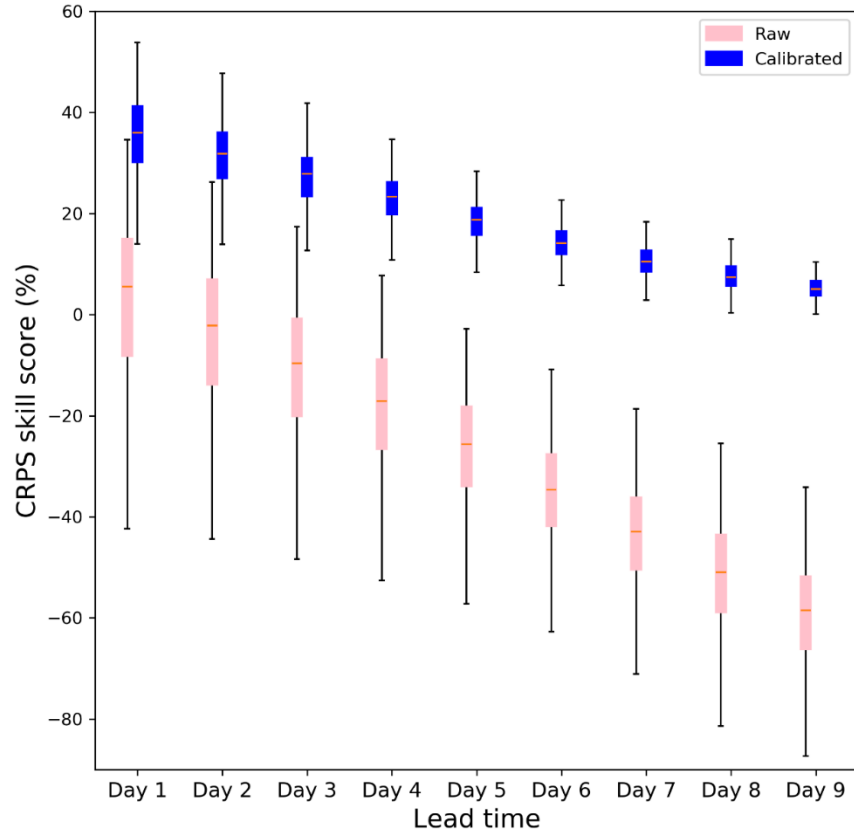
642

643

*Figure S12. Boxplot of biases in raw ETo forecasts constructed raw (blue) and bias-corrected inputs (pink)*

644

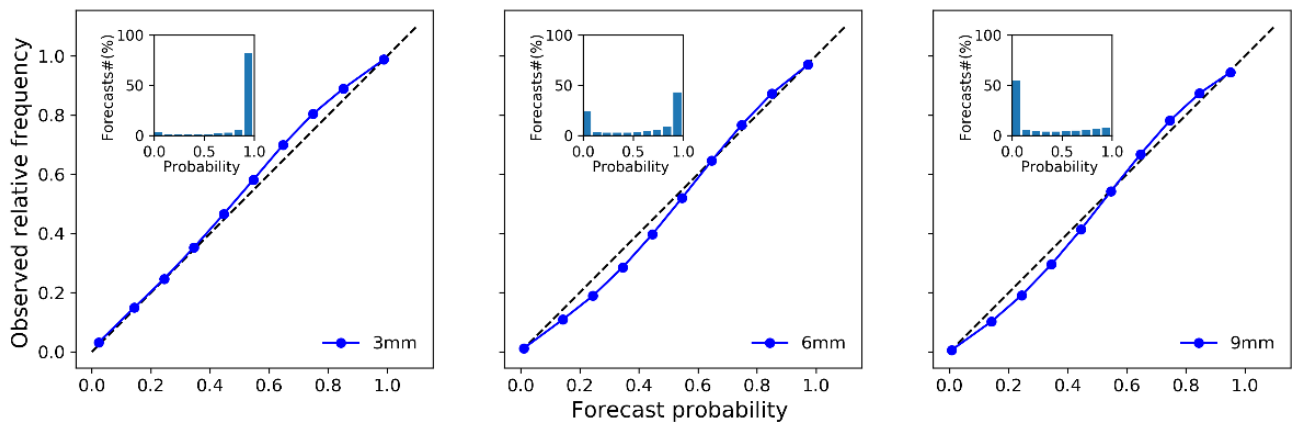




645

646 *Figure S13. Boxplot of CRPS skill score in raw (pink) and calibrated ETo forecast (blue) from*  
 647 *Calibration 2*

648 **We also create reliability diagrams to summarize to evaluate the calibrated ensemble**  
 649 **forecasts from Calibration 2. The three thresholds used to generate the reliability diagram are**  
 650 **3 mm/day, 6mm/day, and 9 mm/day:**



651

652 **Figure 10: Reliability diagrams of calibrated ETo forecasts during 4/2016-3/2019 with thresholds of**  
 653 **3, 6, and 9 mm day<sup>-1</sup>.**

654

655 **We update the Method section to introduce how the reliability diagram is created and how to**  
656 **understand the diagram:**

657 “We further evaluate the reliability of calibrated ETo forecasts from calibration 2 using the  
658 reliability diagram (Hartmann et al., 2002), which assesses how well the predicted probabilities  
659 of forecasts match observed frequencies. We convert the calibrated ensemble ETo forecasts to  
660 forecast probabilities exceeding three thresholds, including 3, 6, and 9 mm day-1. We pool  
661 forecasts of different grid cells, days, and lead times together in the calculation of forecast  
662 probability. In the reliability diagram, perfectly reliable forecasts would demonstrate a curve  
663 along the diagonal. A plotted curve above the diagonal indicates underestimations and vice  
664 versa.”

665

666 **We add the following sentence to section 3.5 (Reliability of calibrated ETo forecasts) to**  
667 **introduce the reliability diagram.**

668 “The reliability diagram further confirms the consistency between forecast probabilities and  
669 observed frequencies (Figure 10). The plotted curves based on three thresholds (3, 6, and 9 mm  
670 day-1) are mainly distributed along the 1:1 line, further indicating the high reliability of  
671 calibrated ETo forecasts.”

672

### 673 Point #5

674 *Authors present experiments 1-4 in the method but then only present some results one experiment 3)*  
675 *and 4) in the last section of results (CRPSS in 3.5). No explanation are provided of why calibration 3) and*  
676 *4) are only briefly introduced. Why is there a big gap with no results on calibration 3) and 4) on the bias*  
677 *and reliability results? Could the authors please expand on the purpose of including these at all in? At*  
678 *p17 l350-354, 'a further evaluation based on a different way of implementing the calibration*  
679 *demonstrate similar improvements in calibrated ETo forecasts with the adoption of bias-correction to*  
680 *input variables'. Is the purpose of including experiment 3) and 4) to test the generalisation of the*  
681 *method? If so, it needs to be clearly stated and justified earlier.*

682 **Response: Thank you for the valuable comments. The reviewer is correct that adding**  
683 **calibrations 3 and 4 is to further evaluate whether the developed calibration strategy could**  
684 **be generally applied to future NWP-based ETo forecasting, and will the strategy be**  
685 **independent of calibration models. We further clarify why we include Calibrations 3 and 4 in**  
686 **Method (section 2.3):**

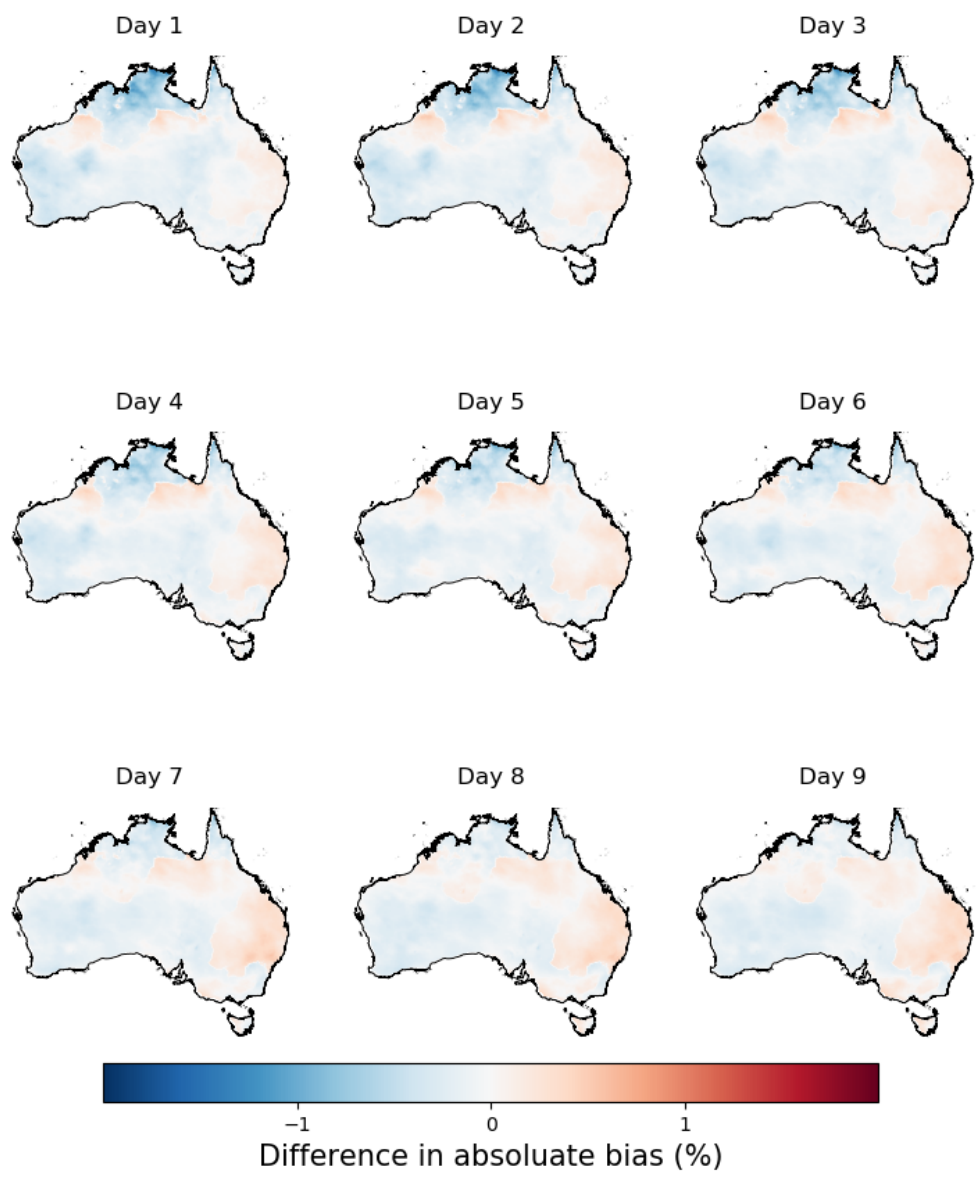
687 “The comparison between Calibrations 1 and 2 is to investigate whether the bias-correction of input  
688 variables would further improve ETo forecasts when the calibration is conducted based on ETo anomalies  
689 and climatological mean. We also conduct additional calibrations which post-process ETo forecasts  
690 directly (Calibrations 3 and 4), to test whether the contribution of improving input variables to ETo  
691 forecast calibration, if there is any, will depend on how ETo forecasts are calibrated (based on anomalies  
692 vs. based on ETo). Calibrations 3 and 4 will help evaluate the general applicability of strategy ii to  
693 enhance NWP/GCM-based ETo forecasting. Key steps of the four calibrations could be found in the

694 schematic diagram introducing how raw ETo forecasts are constructed and how calibrations are  
695 conducted (Figure S1). In the main text, we primarily analyze results from Calibrations 1 and 2.  
696 Improvements with the adoption of bias-correction to input variables in Calibrations 3 and 4 are very  
697 similar to Calibrations 1 and 2 (see the Supplementary Material). To avoid redundancy, we mainly  
698 present results from Calibrations 3 and 4 in the Supplementary Material.”

699

700 **In the original submission, we did not present all results from Calibrations 3 and 4 because**  
701 **these two calibrations were complementary for supporting findings from Calibrations 1 and**  
702 **2. In addition, differences in bias, reliability, and correlation coefficient between Calibrations**  
703 **3 and 4 are very similar to those between Calibrations 1 and 2. We thought it might be a bit**  
704 **redundant and may confuse readers if we present all results from Calibrations 3 and 4 in the**  
705 **main text. However, we agree with the reviewer that it is necessary to present results from**  
706 **Calibrations 3 and 4 in case readers are interested in them. In the revised manuscript, we**  
707 **present them in the supplementary material (See the figures below), in order not to distract**  
708 **readers from understanding key objectives (e.g., the necessity of bias-correcting input**  
709 **variables prior to ETo calibration) of this investigation. Specifically, in addition to the figure**  
710 **showing improvements in CRPS skill score, we also add figures demonstrating differences in**  
711 **absolute bias (Figure S15), correlation coefficients (Figure S16), and alpha index (Figure S18)**  
712 **between Calibrations 3 and 4 in the Supplementary Material:**

713



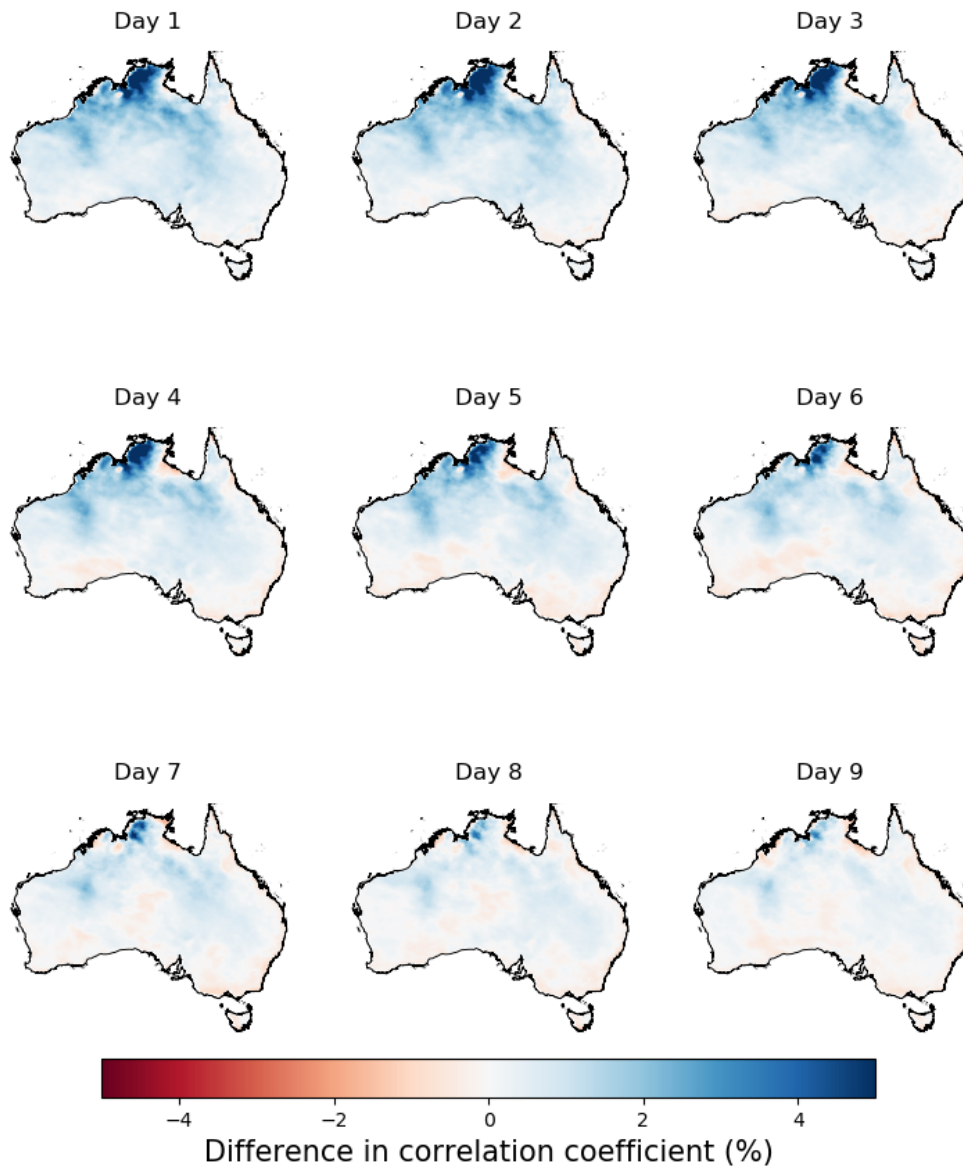
714

715

716

717

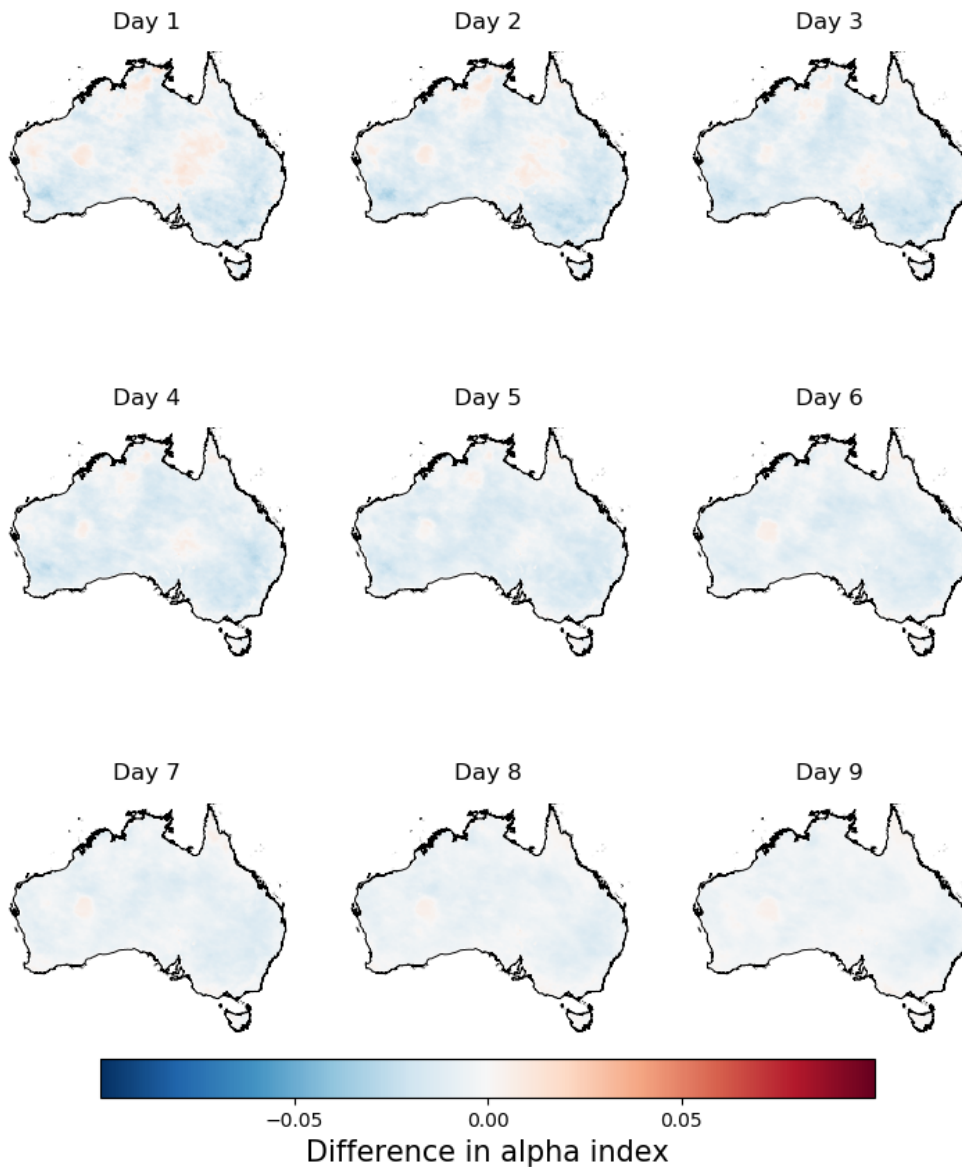
*Figure S15. Differences in absolute bias between Calibrations 3 and 4*



718

719

*Figure S16. Differences in correlation coefficient between Calibrations 3 and 4*



720  
721 *Figure S18. Differences in alpha index between Calibrations 3 and 4*

722

723

724 **We add one new subsection in Results to introduce results from Calibrations 3 and 4**

725 **3.7 Results from Calibrations 3 and 4**

726 “We also compare the bias, correlation coefficient, CRPS skill score, and reliability of calibrated forecasts  
727 from Calibrations 3 and 4, to evaluate whether we can obtain similar improvements through the bias-  
728 correction of input variables if we conduct the ETo forecast calibration in a different way (without using  
729 ETo climatological mean and anomalies). Results show that the adoption of bias-correction also leads to  
730 lower bias, higher correlation coefficient, and higher CRPS skill score in terms of magnitude, spatial

731 patterns, and trend along the lead times, when ETo forecasts are calibrated directly (Figure S15-S17). In  
732 addition, the alpha index was only slightly different between Calibrations 3 and 4 (Figure S18). This  
733 additional comparison further confirms the general applicability of strategy ii for enhancing NWP-based  
734 ETo forecasting.”

735

## 736 Point #6

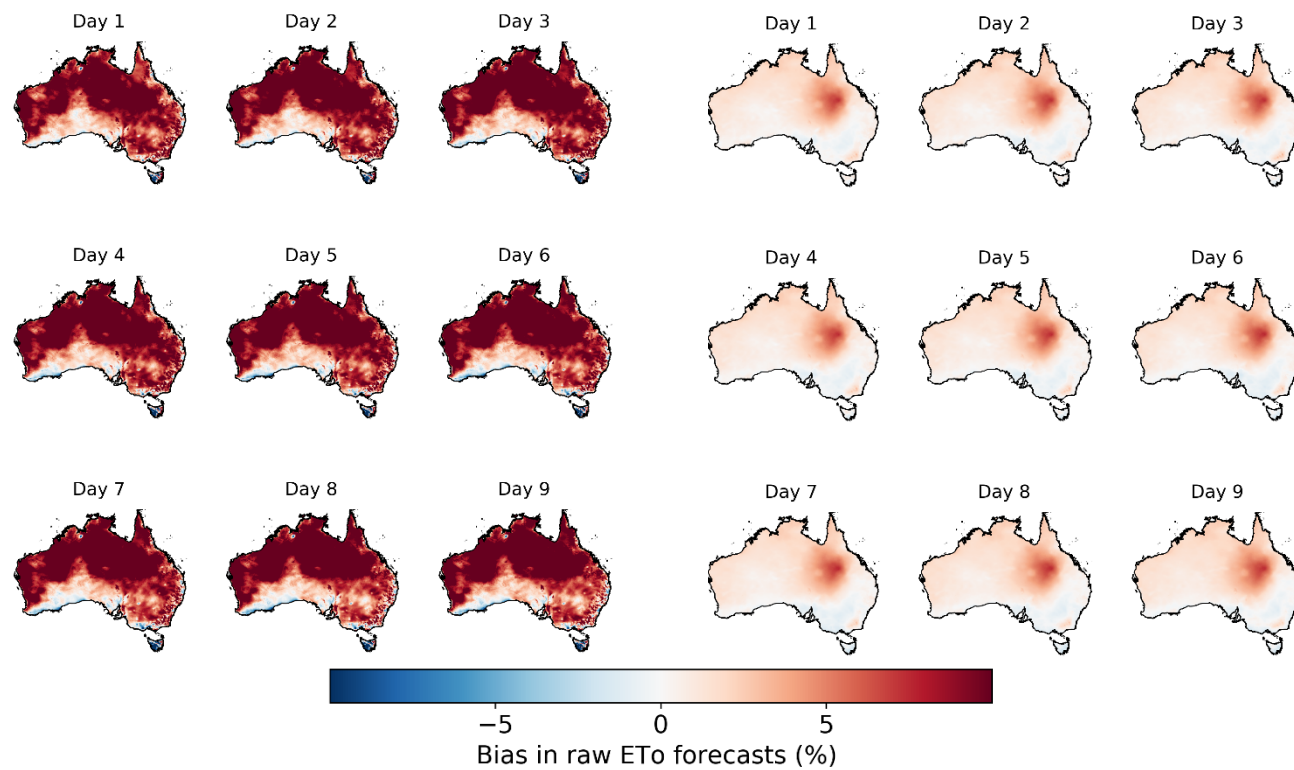
737 *Methodological choices for evaluation:*

738 *P7 | 180-185 : why choosing the absolute bias and over a relative measure e.g. percentage bias? This*  
739 *choice makes it difficult to compare the magnitude of the errors in the results across different variables*  
740 *and studies. For example, figure 1 shows a bias between -2 to 2mm/day which does not seem like much*  
741 *compared to other input variables such as precipitation. Figure 3 with a range of -0.1 to 0.1 seems very*  
742 *small. Conversely, percentages are used for the correlation coefficient in Figure 6 so why not use it for*  
743 *the bias?*

744 **Response: We appreciate the reviewer's valuable comments. Bias shows differences of the**  
745 **mean between forecasts and observations, and could be either positive (overestimation) or**  
746 **negative (underestimation). Larger departures from the observed mean, no matter the bias is**  
747 **positive or negative, suggest more significant inconsistencies with observations. Absolute**  
748 **bias is a good indicator measuring the departure from the observed mean. As a result, using**  
749 **absolute bias, we can compare results from two different calibrations, with smaller absolute**  
750 **bias indicating closer to the observed mean, and thus suggesting better performance.**

751 **We agree with the reviewer that using percentages will make the results more comparable**  
752 **with other variables, or with other studies. As a result, we change the unit of bias in figures 1,**  
753 **S12, 3, 4 to percentage:**

754



755

756

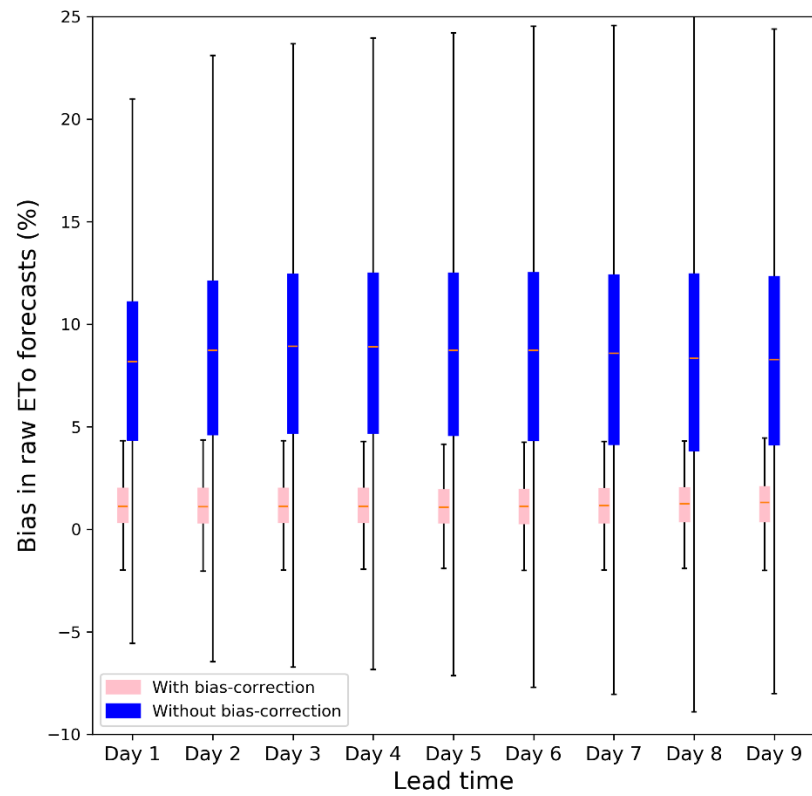
757 **Figure 1: Bias in (three panels on the left) raw ETo forecasts constructed with raw forecasts of input variables and (three panels on the right)**  
 758 **raw ETo forecasts constructed with bias-corrected input variables.**

759

760

761





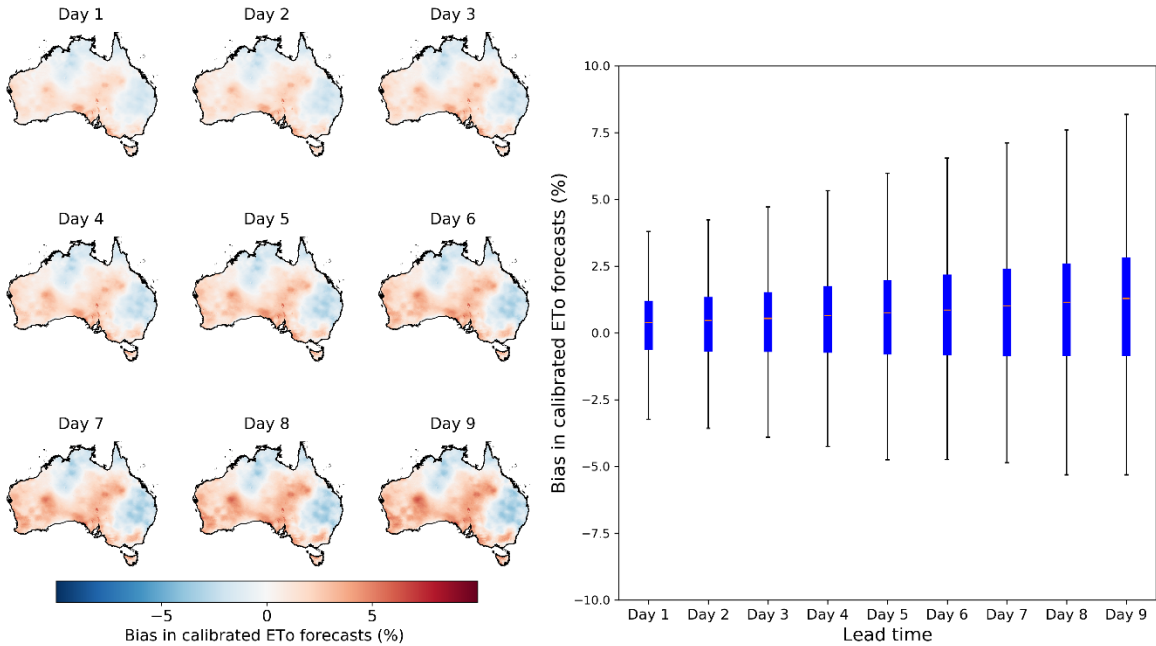
762

763

*Figure S12. Boxplot of biases in raw ETo forecasts constructed raw (blue) and bias-corrected inputs (pink)*

764

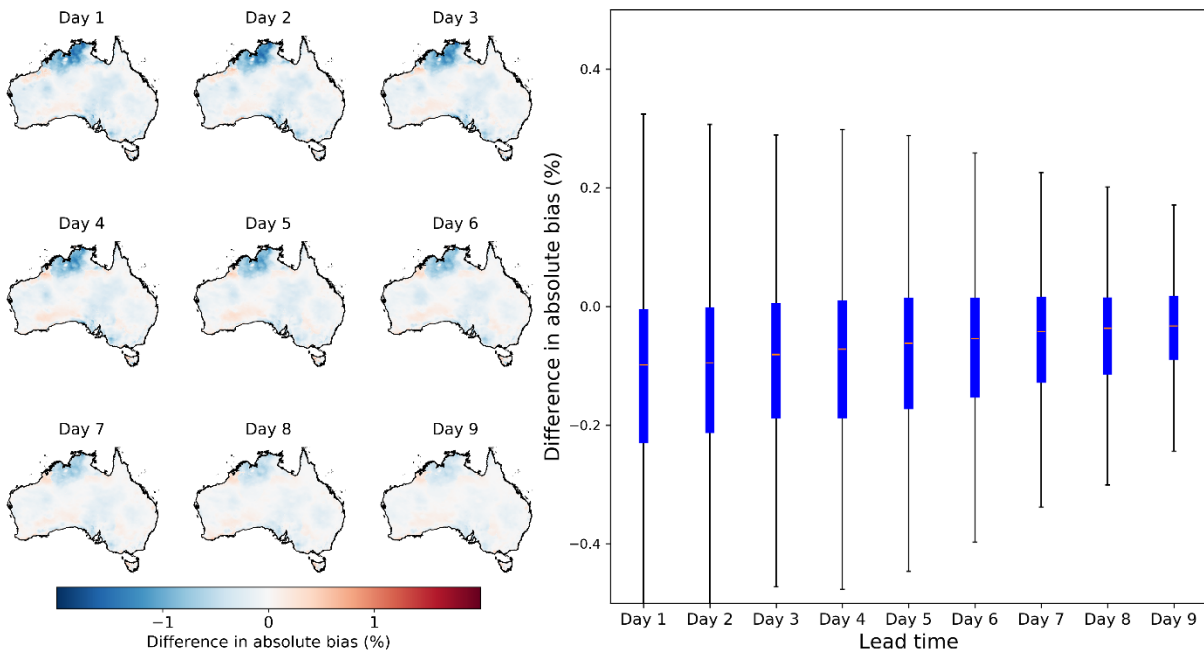
765



766

767 **Figure 3: Bias in calibrated ETo forecasts of 9 lead times from Calibration 2, in which raw ETo forecasts**  
 768 **are constructed with bias-corrected input variables. Maps on the left show the spatial patterns of**  
 769 **bias, and the boxplot on the right summarizes results for all grid cells.**

770



771

772 **Figure 4: Differences in absolute bias between calibrated ETo forecasts from Calibration 2 with**  
 773 **Calibration 1. Maps on the left show the spatial patterns of difference in absolute bias, and the**  
 774 **boxplot on the right summarizes results for all grid cells.**

775

776 Point #7

777 *P8 I205-2015: why is climatology used as reference forecast for the skill score? In hydrological*  
778 *forecasting persistence is typically used for short lead times, whereas climatology would be used for*  
779 *longer lead times, see fore example (Pappenberger, Ramos et al. 2015). Could you please expand and*  
780 *justify the choice of reference forecast used and implication of interpretation of results?*

781 **Response: We really appreciate the reviewer’s valuable suggestion and the introduction of**  
782 **this classic paper. We choose the climatology forecasts as the reference rather than using**  
783 **persistence for several reasons:**

784 **1, Climatology forecasts have been widely used as the reference in the calculation of CRPS**  
785 **skill score for short-term hydroclimate forecasts. Since climatology forecasts have similar**  
786 **errors across all lead times (Bennett et al., 2014), they have been used as the reference to**  
787 **compare forecast skills among different lead times (Academies, 2014; Zhao et al., 2019).**

788 **2, Persistence is also a good reference, but it's been mainly used for the first two lead times.**  
789 **As demonstrated in figure 5 of Bennett et al. (2014), errors in persistence could increase**  
790 **quickly with lead time. As a result, multiple studies suggested that persistence is good for skill**  
791 **discrimination for short lead times (Pappenberger et al., 2015; Thiemig et al., 2015).**

792 **Since we investigate 9 lead times in this study, errors in persistency are expected to be**  
793 **significant at long lead times. Using persistence as the reference may artificially exemplify**  
794 **forecast skills at long lead times. As a result, we think the use of climatology forecasts as the**  
795 **reference for the calculation of the CRPS skill score is acceptable.**

796 **We add the following sentences to section 2.4.3 (Skills of the raw and calibrated forecasts) to**  
797 **explain the use of climatology forecasts as the reference for the calculation of CRPS skill score**

798 “In the calculation of CRPS skill score, both climatology forecasts or the last observations  
799 (persistence) have been used as reference forecasts (Pappenberger et al., 2015; Thiemig et al.,  
800 2015). However, reference forecasts based on persistence are more suitable for evaluating the  
801 performance of forecasts shorter than two days. As a result, we choose climatology forecasts as  
802 the reference, since errors in climate forecasts are similar among all lead times and thus could be  
803 used to demonstrate the increasing errors in raw and calibrated forecasts as lead time advances.”

804

805 **Reference:**

806 Academies, N.: The science of NOAA’S Operational Hydrologic Ensemble Forecast Service, Bull.  
807 Am. Meteorol. Soc., (January), 79–98, doi:10.1175/BAMS-D-12-00081.1, 2014.

808 Bennett, J. C., Robertson, D. E., Lal, D., Wang, Q. J., Enever, D., Hapuarachchi, P. and Tuteja, N.  
809 K.: A System for Continuous Hydrological Ensemble Forecasting (SCHEF) to lead times of 9 days,  
810 J. Hydrol., 519, 2832–2846, doi:10.1016/j.jhydrol.2014.08.010, 2014.

811 Pappenberger, F., Ramos, M. H., Cloke, H. L., Wetterhall, F., Alfieri, L., Bogner, K., Mueller, A.  
812 and Salamon, P.: How do I know if my forecasts are better? Using benchmarks in hydrological  
813 ensemble prediction, *J. Hydrol.*, 522, 697–713, doi:10.1016/j.jhydrol.2015.01.024, 2015.

814 Thiemig, V., Bisselink, B., Pappenberger, F. and Thielen, J.: A pan-African medium-range  
815 ensemble flood forecast system, *Hydrol. Earth Syst. Sci.*, 19, 3365–3385, doi:10.5194/hess-19-  
816 3365-2015, 2015.

817 Zhao, T., Wang, Q. J. and Schepen, A.: A Bayesian modelling approach to forecasting short-term  
818 reference crop evapotranspiration from GCM outputs, *Agric. For. Meteorol.*, 269–270(January),  
819 88–101, doi:10.1016/j.agrformet.2019.02.003, 2019.

820

### 821 Point #8

822 *P8 I214. Why is the definition of CRPSS using percentage? As far as I am aware, most studies do not*  
823 *present the CRPSS in terms of percentage, could you please comment on the reason of this choice with*  
824 *references that also use percentages and if there is any advantages?*

825 **Response: Thank you for the comments. We agree with the reviewer that many studies use**  
826 **ratios when presenting the CRPS skill score. Meanwhile, we also notice that some studies (see**  
827 **the reference list at the bottom of our response to this comment) use percentage as the unit**  
828 **of CRPS skill score.**

829 **As shown in Figure 7, skills of calibrated forecasts decreased quickly with lead time. As a**  
830 **result, the CRPS skill score approaches zero at lead time 9. One advantage of using the**  
831 **percentage as the unit of CRPS skill score is that small decimals of low skills will be converted**  
832 **to more readable percent.**

833 **We add the following sentence to explain why the percentage is used as the unit of CRPS skill**  
834 **score:**

835 *“We use percentage as the unit of CRPS skill score so low skill scores at long lead times will be*  
836 *converted from small decimals to more readable percent.”*

837 **We believe the choice of percentage as the unit of CRPS skill score will not affect the**  
838 **conclusions of this study. Here are some investigations using % as the unit of CRPS skill score:**

839 Brown, J. D. and Seo, D. J.: A nonparametric postprocessor for bias correction of hydrometeorological  
840 and hydrologic ensemble forecasts, *J. Hydrometeorol.*, 11(3), 642–665, doi:10.1175/2009JHM1188.1,  
841 2010.

842 Kumar, L. G. A., Smith, A. S. D., Gonzalez, G. B. P., Merryfield, V. K. W. and Newman, A. S. Á. M.: A  
843 verification framework for interannual-to-decadal predictions experiments, *Clim. Dyn.*, 40, 245–272,  
844 doi:10.1007/s00382-012-1481-2, 2013.

845 Munkhammar, J., van der Meer, D. and Widén, J.: Probabilistic forecasting of high-resolution clear-sky  
846 index time-series using a Markov-chain mixture distribution model, *Sol. Energy*, 184(January), 688–695,  
847 doi:10.1016/j.solener.2019.04.014, 2019.

848 Robertson, D. E. and Wang, Q. J.: Seasonal Forecasts of Unregulated Inflows into the Murray River ,  
849 Australia, *Water Resour. Manag.*, 27, 2747–2769, doi:10.1007/s11269-013-0313-4, 2013.

850 Schepen, A., Wang, Q. J. and Robertson, D. E.: Seasonal Forecasts of Australian Rainfall through  
851 Calibration and Bridging of Coupled GCM Outputs, *Mon. Weather Rev.*, 142, 1758–1770,  
852 doi:10.1175/MWR-D-13-00248.1, 2014.

853

#### 854 Point #9

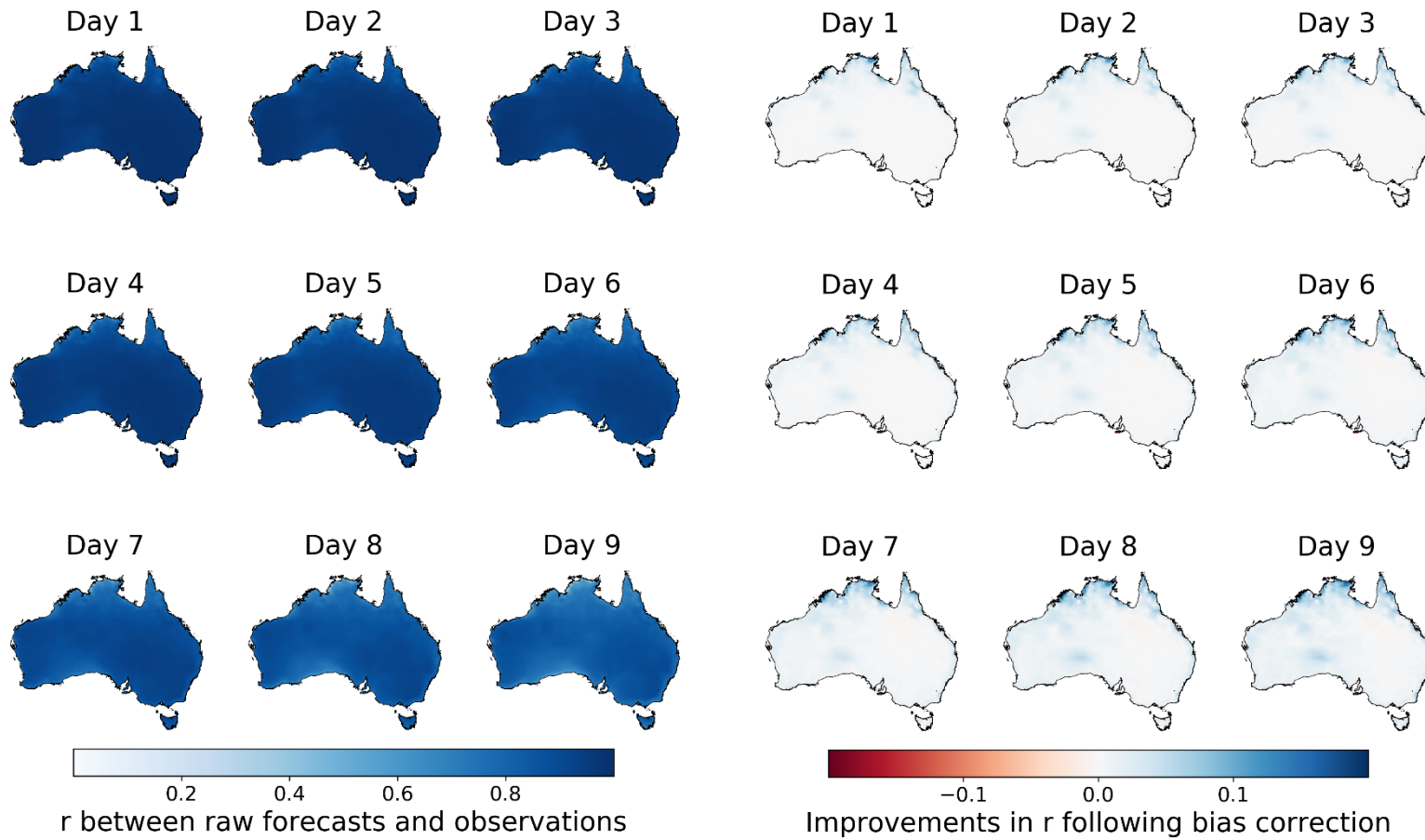
855 Analysis and interpretation of results:

856 *P11 I259-261: why the higher difference in bias in approaches for the Northern Territory? How does this*  
857 *relate to the biases, errors and assumptions of the NWP? Is it correlated to the biases of specific input*  
858 *variables? How is it correlated to the nonlinear relationship in calculating ETo? Why are the biases most*  
859 *pronounced for shorter lead times? Please comment.*

860 **Response: Thank you for the valuable comments. To answer these questions, we present**  
861 **more results to explain how quantile mapping to input variables contributes to improving**  
862 **calibrated ETo forecasts. Specifically, we (1) calculate the correlation coefficients ( $r$ ) between**  
863 **raw/bias-corrected forecasts of the five input variables and AWAP data to further analyze**  
864 **how quantile mapping has improved input variables, in addition to correcting bias (shown in**  
865 **figure 1); (2) investigate the improvements in correlation coefficients between raw ETo**  
866 **forecasts following the bias-correction to input variables and AWAP ETo, to examine how**  
867 **improvements in each variable are translated into the resultant raw ETo forecasts; (3) explain**  
868 **how improvements in raw ETo forecasts through bias-correcting input variables lead to**  
869 **improvements in calibrated ETo forecasts. Please find more details as follows:**

870 **1, In addition to correcting bias (Figures S2 to S6), quantile mapping also generally improves**  
871 **the temporal patterns of raw forecasts of the input variables. Following figures shows  $r$**   
872 **between raw forecasts of the input variables and their corresponding AWAP data (three**  
873 **columns on the left), and improvements in  $r$  by quantile mapping (three columns on the**  
874 **right):**

875



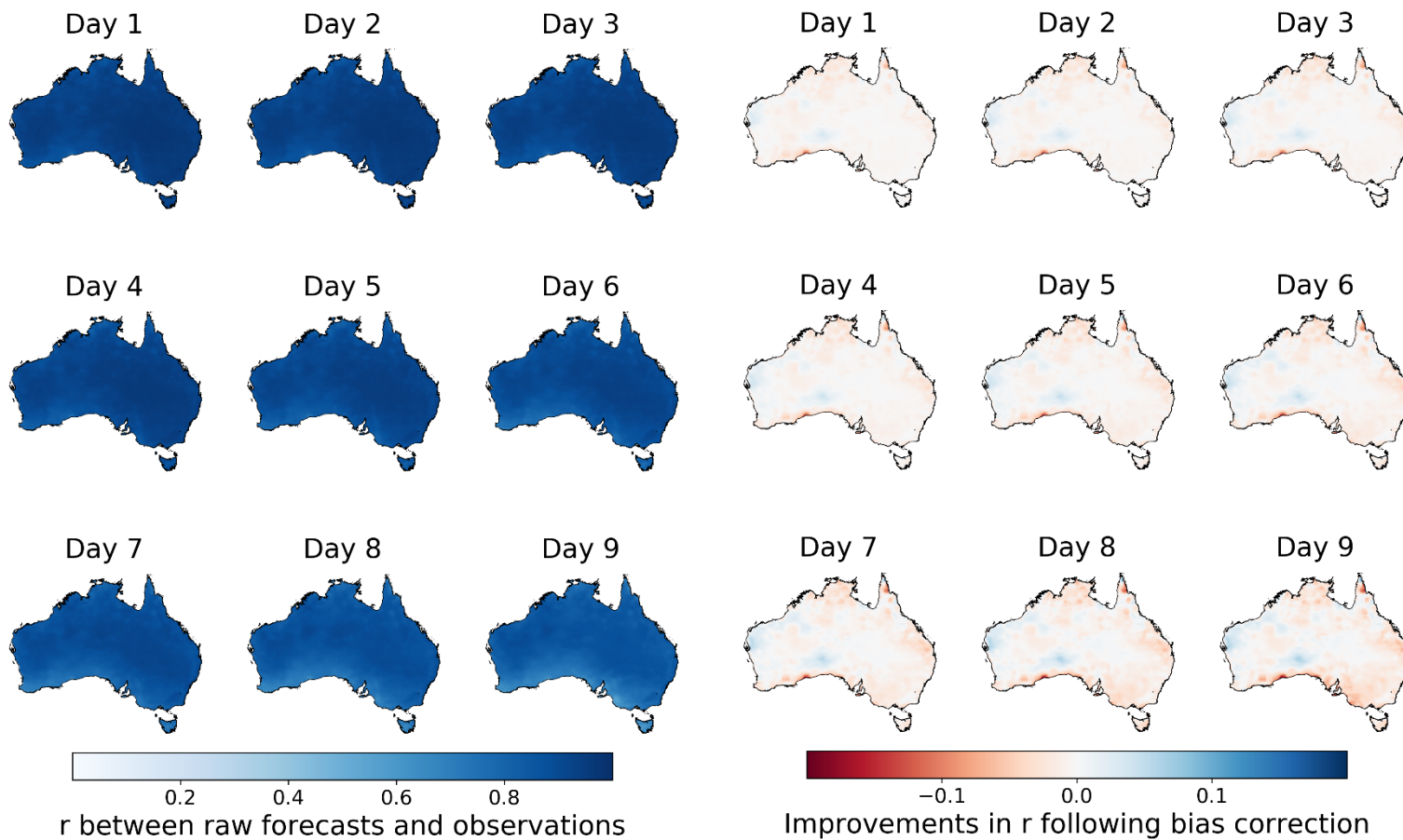
876

877

878

879

*Figure S7. Correlation coefficients ( $r$ ) between raw  $T_{max}$  forecasts and AWAP data (three panels on the left), and improvements in  $r$  (three panels on the right) through quantile mapping*



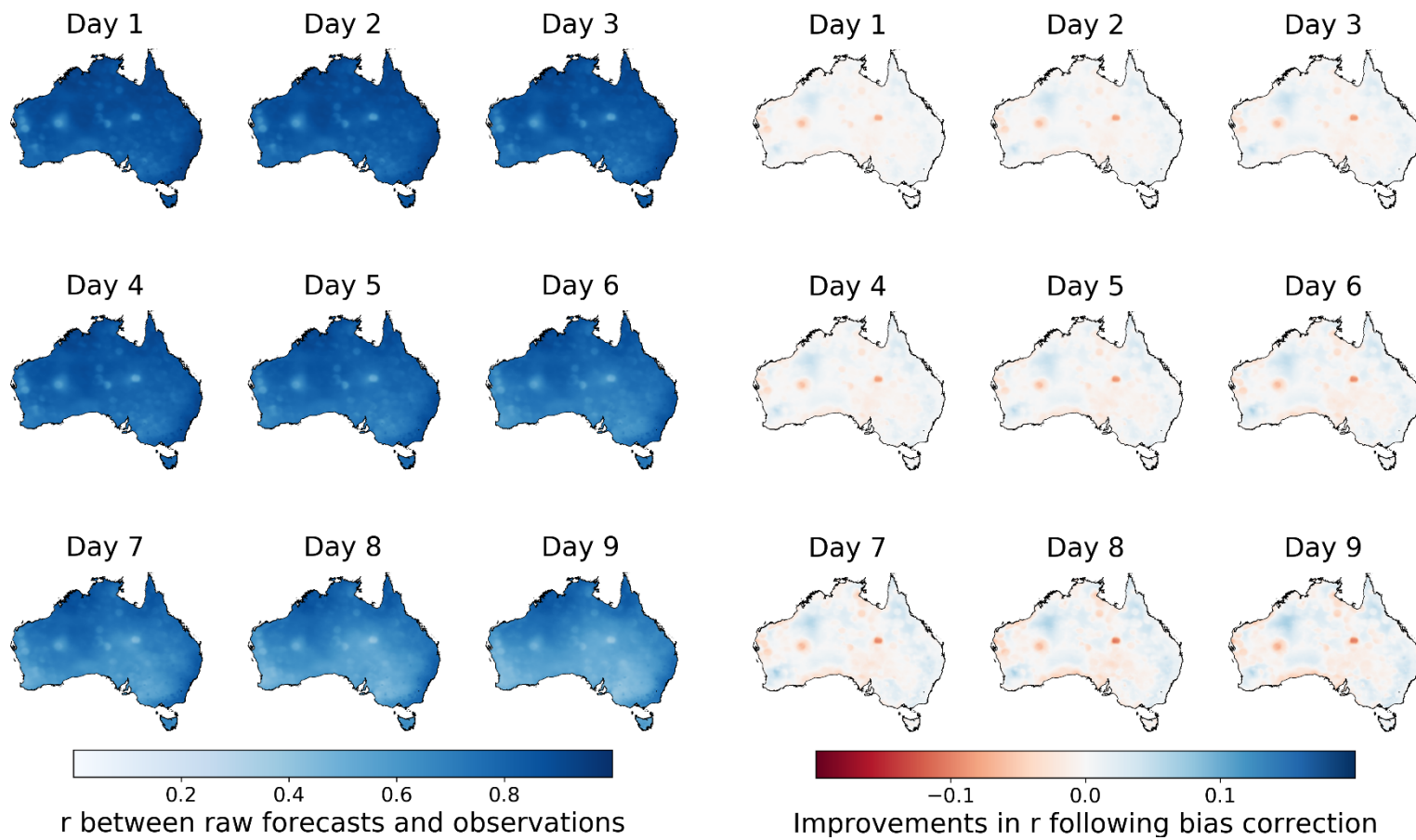
880

881

882

883

*Figure S8. Correlation coefficients ( $r$ ) between raw  $T_{min}$  forecasts and AWAP data (three panels on the left), and improvements in  $r$  (three panels on the right) through quantile mapping*



884

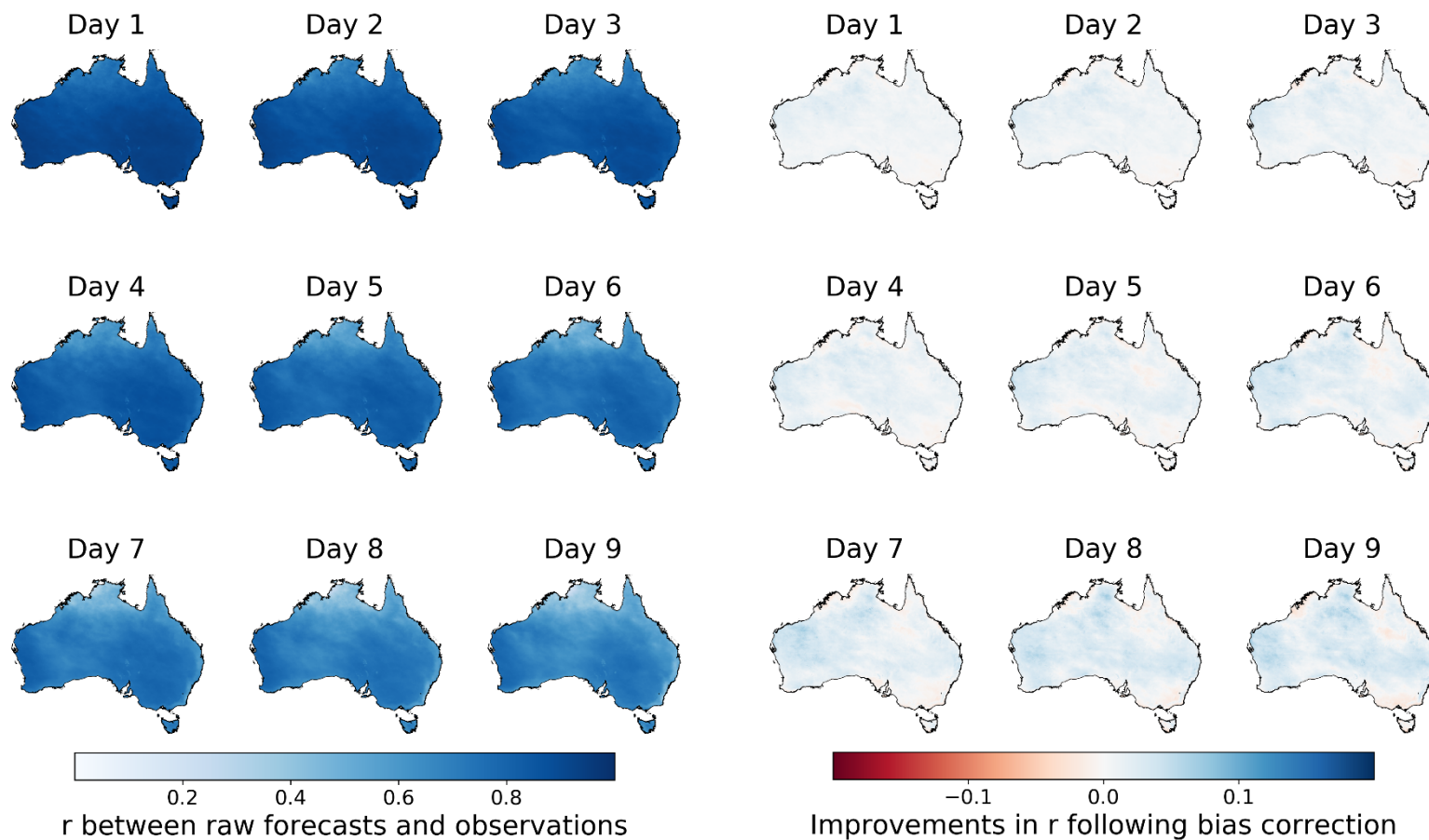
885

886

887

*Figure S9. Correlation coefficients ( $r$ ) between raw vapor pressure forecasts and AWAP data (three panels on the left), and improvements in  $r$  (three panels on the right) through quantile mapping*



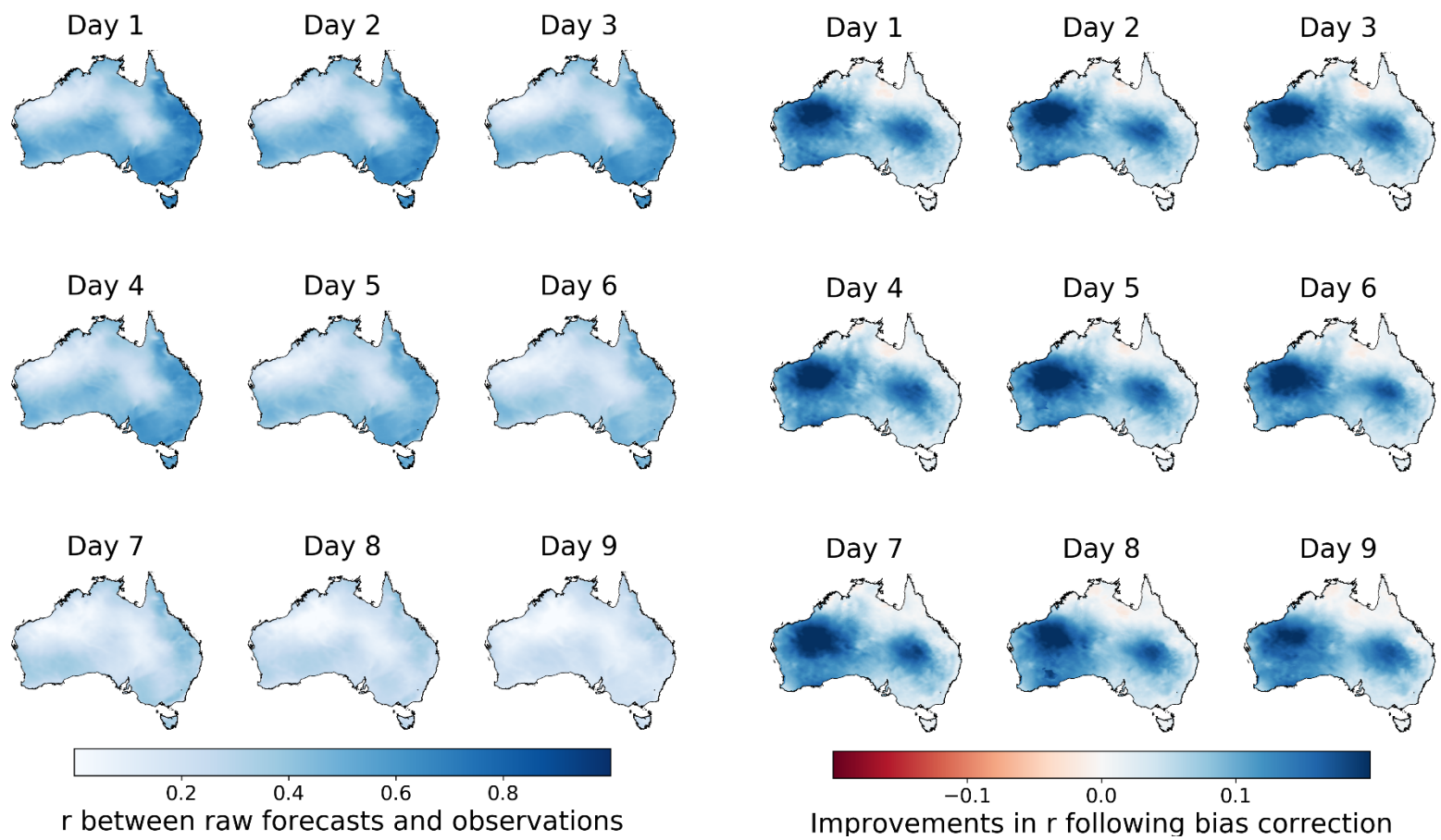


888

889

890

*Figure S10. Correlation coefficients ( $r$ ) between raw solar radiation forecasts and AWAP data (three panels on the left), and improvements in  $r$  (three panels on the right) through quantile mapping*



891

892

893

894

895

896

*Figure S11. Correlation coefficients ( $r$ ) between raw wind speed forecasts and AWAP data (three panels on the left), and improvements in  $r$  (three panels on the right) through quantile mapping*

897 **As shown in the above figures,  $r$  between raw forecasts of the input variables and AWAP data**  
898 **varies with the input variables. The two temperature variables have higher  $r$  values than the**  
899 **other three variables, and wind speed forecasts demonstrate the lowest correlation with**  
900 **AWAP data. For all variables, the  $r$  decreases with lead time, indicating higher uncertainties in**  
901 **raw forecasts at longer lead times.**

902 **Quantile mapping generally improves the correlation between forecasts of the input**  
903 **variables and AWAP data. The above figures show that bias-corrected forecasts demonstrate**  
904 **higher  $r$  for Tmax, solar radiation, and wind speed across most parts of Australia; for Tmin**  
905 **and vapor pressure, changes in  $r$  are less significant, and both increases and slight decreases**  
906 **in  $r$  are observed.**

907 **We add the above figures to the supplementary. We also add following descriptions to**  
908 **section 3.1:**

909 “Raw forecasts of the input variables generally agree with the AWAP data in temporal patterns during the  
910 study period, but the  $r$  varies with variables (Figures S7-S11). The two temperature variables (Tmax and  
911 Tmin) have higher  $r$  ( $>0.9$ ) than the other three variables, and wind speed forecasts demonstrate the  
912 lowest correlations with AWAP data. For all variables, the  $r$  decreases with lead time, indicating higher  
913 uncertainties at long lead times in raw forecasts.”

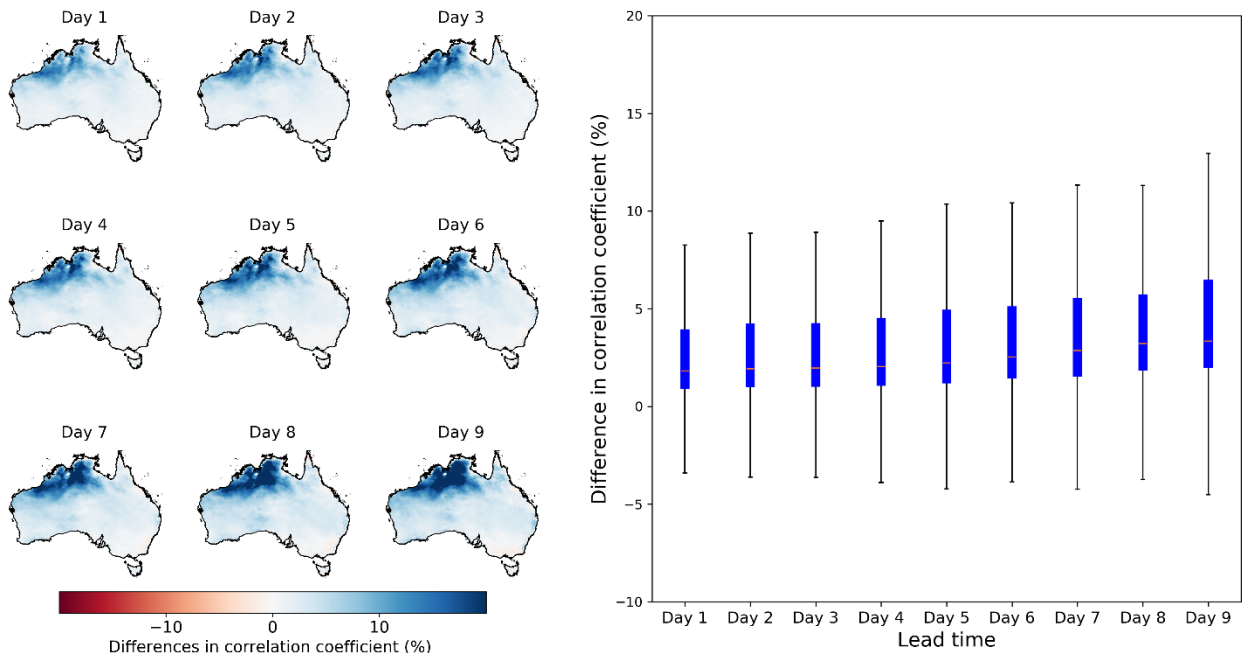
914 “In addition, quantile mapping also improves the correlation between forecasts of input variables and  
915 AWAP data (Figures S7-S11). The most significant improvements are found in wind speed forecasts, in  
916 which the  $r$  is improved by up to 0.2 in central and southern parts of Australia. Forecasts of Tmax and  
917 solar radiation also demonstrate higher  $r$  with the adoption of quantile mapping. Both increases and slight  
918 decreases were found for vapor pressure and Tmin, indicating less significant improvements in temporal  
919 patterns than other variables. ”

920 **2, With the adoption of quantile mapping to raw forecasts of individual variables, raw ETo**  
921 **forecasts (Calibrations 2 or 4) also show higher  $r$  with observations, than the raw ETo**  
922 **forecasts constructed with the raw forecasts of input variables (Calibrations 1 or 3):**

923

924

925



926

927 **Figure 2: The comparison between the correlation coefficient of AWAP ETo and raw ETo forecasts**  
 928 **constructed with the bias-corrected inputs vs. the correlation coefficient of AWAP ETo and raw ETo**  
 929 **forecasts constructed with the uncorrected inputs. The boxplot on the right summarizes results for all**  
 930 **grid cells.**

931 **As is shown in the above figure, the quantile mapping also improves the temporal patterns of raw ETo**  
 932 **forecasts, for the lead times. More significant improvements are found in northern Australia. Due to**  
 933 **the nonlinearity in the calculation of ETo using the input variables, spatial patterns of improvements**  
 934 **in  $r$  (Figure 2) do not resemble improvements of any individual input variables. Although both T<sub>max</sub>**  
 935 **and wind speed show more significant improvements in northern Australia, where the  $r$**   
 936 **improvements are greater than other regions (Figure 2), improvements in the two variables do not**  
 937 **lead to higher  $r$  in other parts of Australia. As a result, we believe that improvements in  $r$  of raw ETo**  
 938 **forecasts are contributed jointly by these input variables and their interactions.**

939 **We add the above figure (Figure 2) and the following contents to the manuscript:**

940 “The adoption of quantile mapping to input variables also improves the temporal patterns of raw ETo  
 941 forecasts (Figure 2). Compared with the raw ETo forecasts constructed with raw input variables, the raw  
 942 ETo forecasts based on bias-corrected inputs generally shows higher correlations with AWAP ETo,  
 943 particularly in northern Australia, where  $r$  is improved by more than 10%. However, due to the  
 944 nonlinearity in the calculation of ETo using the input variables, spatial patterns of improvements in  $r$   
 945 (Figure 2) does not resemble improvements in any individual input variables (Figures S7 to S11). The  
 946 improvements in  $r$  of raw ETo forecasts seem to be contributed jointly by these input variables and their  
 947 interactions.”

948 **3, We add the following contents to section 3.3 to explain the spatial patterns of changes in  $r$**   
 949 **and absolute bias:**

950 “Larger reductions in absolute bias in northern Australia coincide with the improvements in the  
 951 correlation between raw ETo forecasts and AWAP ETo (Figure 2). However, unlike the improvements in

952 *r* for all lead times in raw ETo forecasts, the improvements in absolute bias are more pronounced at short  
953 lead times (Days 1-3) than long lead times (Days 7-9). The uneven improvements across different lead  
954 times may be caused by the significant intrinsic uncertainties in forecasts, which have hindered the  
955 manifestation of improvements to raw ETo forecasts at long lead times in calibrated forecasts.”

956 **Based on the above analyses, we can then answer the questions the reviewer raised in this**  
957 **comment.**

958 **More significant reductions in absolute bias in northern Australia show similar spatial**  
959 **patterns with that of the improvements in *r* between raw ETo forecasts and AWAP ETo. As we**  
960 **further explained in our response to your next comment (#10), deficiencies in NWP models in**  
961 **simulating weather dynamics in tropical regions have been reported. Bias-correction**  
962 **effectively corrects errors in these areas. However, improvements to raw ETo forecasts in *r***  
963 **with the application of quantile mapping could not be explained by any individual variable.**  
964 **The nonlinearity in calculating ETo based on the individual variables may have combined**  
965 **improvements in each variable and lead to more significant improvements in northern**  
966 **Australia. Less significant improvements in calibrated ETo forecasts at longer lead times may**  
967 **be caused by the more significant intrinsic uncertainties in forecasts than short lead times.**  
968 **These uncertainties have inhibited the translation of improvements in raw ETo forecasts to**  
969 **calibrated forecasts.**

970

#### 971 Point #10

972 *P13 l282-285: Why lowest score of correlation coefficient in northern Territory? Is it linked to the NWP*  
973 *(and if so how?) or is it linked to observations? E.g. differences in observations compared to rest of*  
974 *country?*

975 **Response: Thank you for the comments. We believe the low correlation results from the**  
976 **NWP forecasts rather than from observations for several reasons:**

977 **1, Evaluation of the observations (AWAP data) did not show larger errors in northern**  
978 **Australia than other areas of Australia (Jones et al., 2009). As a result, we do not have**  
979 **evidence that the quality of observations in this region is lower than in other regions**

980 **2, Deficiencies of NWP forecasts in tropical regions in Australia have been well documented.**  
981 **Due to its highly dynamic nature, forecasts for tropical regions often demonstrate larger**  
982 **uncertainties than other climate zones. In the evaluation of NWP forecasts in Australia,**  
983 **tropical zones show lower skills than other regions (Ebert and McBride, 2000; McBride and**  
984 **Ebert, 2000; Roux et al., 2010). According to Huang et al. (2018), ACCESS models have been**  
985 **suffering from low skills in simulating the convective processes in tropical zones of Australia.**

986 **3, Raw ETo forecasts constructed with outputs of an early version of the ACCESS model in**  
987 **another study showed higher RMSE in Northern Territory than other regions (Perera et al.,**

988 **2014), further confirms that lower correlation coefficient is mainly caused by the NWP**  
989 **forecasts.**

990 **We add the following sentences to section 3.3:**

991 “Deficiencies in ACCESS models in simulating dynamics of tropical climate systems may have  
992 resulted in the low  $r$  in northern Australia.”

993

994 **Reference:**

995 Ebert, E. E. and McBride, J. L.: Verification of precipitation in weather systems : determination  
996 of systematic errors, *J. Hydrol.*, 239, 179–202, 2000.

997 Huang, J., Rikus, L. J., Qin, Y. and Katzfey, J.: Assessing model performance of daily solar  
998 irradiance forecasts over Australia, *Sol. Energy*, 176(November), 615–626,  
999 doi:10.1016/j.solener.2018.10.080, 2018.

1000 Jones, D. A., Wang, W. and Fawcett, R.: High-quality spatial climate data-sets for Australia, *Aust.*  
1001 *Meteorol. Oceanogr. J.*, 58, 233–248, 2009.

1002 McBride, J. L. and Ebert, E. E.: Verification of quantitative precipitation forecasts from  
1003 operational numerical weather prediction models over Australia, *Weather Forecast.*, 15(1),  
1004 103–121, doi:10.1175/1520-0434(2000)015<0103:VOQPFF>2.0.CO;2, 2000.

1005 Perera, K. C., Western, A. W., Nawarathna, B. and George, B.: Forecasting daily reference  
1006 evapotranspiration for Australia using numerical weather prediction outputs, *Agric. For.*  
1007 *Meteorol.*, 194, 50–63, doi:10.1016/j.agrformet.2014.03.014, 2014.

1008 Roux, B., Seed, A., Pagano, T. and Roux, B.: Improved use of precipitation forecasts in short-  
1009 term water forecasting – progress report, The Centre for Australian Weather and Climate  
1010 Research A partnership between CSIRO and the Bureau of Meteorology Improved., 2010.

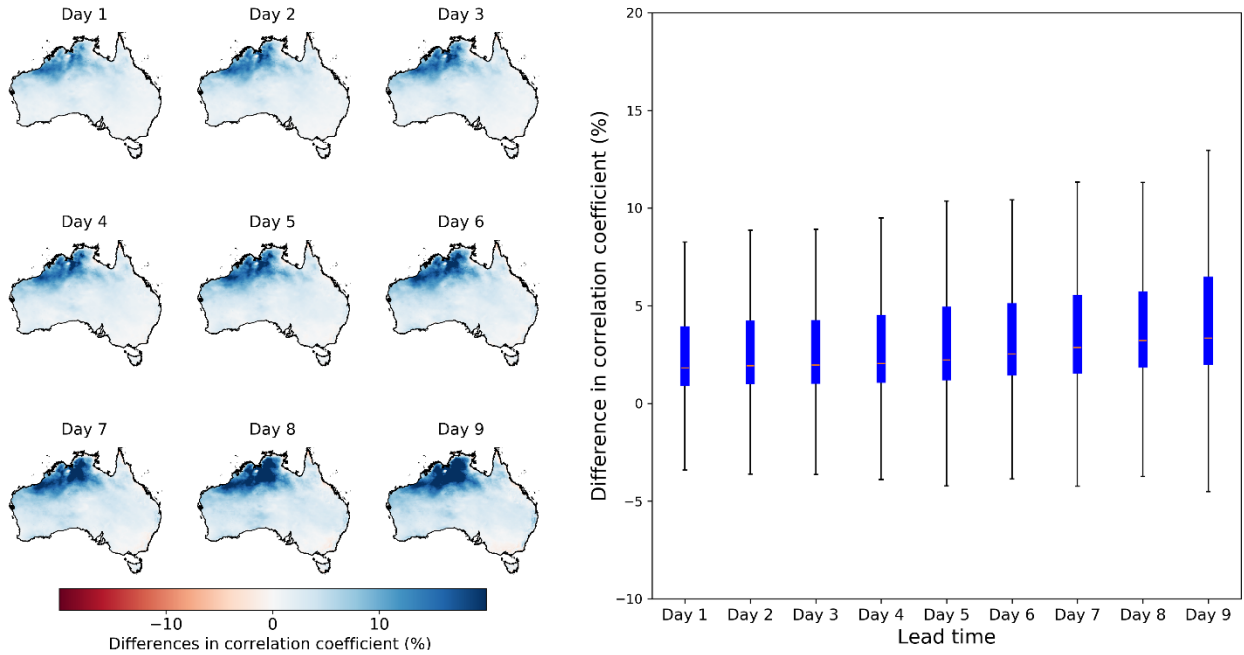
1011

1012

1013 Point #11

1014 *P14 I294-297: The geographical patterns of the correlation performance is very similar to the patterns of*  
1015 *the bias performance. Could you please comment why and if the reasons are the same? Are these related*  
1016 *to either the NWP or observations?*

1017 **Response: Thank you for the valuable comments. We add the following figure to the**  
1018 **manuscript to demonstrate how bias-correction of input variables improves correlations**  
1019 **between raw ETo forecasts and AWAP ETo:**



1020

1021 **Figure 2: The comparison between the correlation coefficient of AWAP ETo and raw ETo forecasts**  
 1022 **constructed with the bias-corrected inputs vs. the correlation coefficient of AWAP ETo and raw ETo**  
 1023 **forecasts constructed with the uncorrected inputs. The boxplot on the right summarizes results for all**  
 1024 **grid cells.**

1025 **The above figure shows that when input variables are bias-corrected, the resultant raw ETo**  
 1026 **forecasts show higher correlation coefficients, than raw ETo forecasts constructed with raw**  
 1027 **inputs. Spatial patterns of the improvements in  $r$  in raw forecasts for short lead times are**  
 1028 **consistent with the improvements in  $r$  in calibrated forecasts (Figure 6). As a result, we**  
 1029 **believe this is how the new calibration strategy improves the calibration of ETo forecasts.**  
 1030 **Less significant improvements in ETo forecasts at longer lead times may be caused by the**  
 1031 **more significant intrinsic uncertainties in raw forecasts than short lead times. These**  
 1032 **uncertainties have inhibited the translation of improvements to raw ETo forecasts in**  
 1033 **calibrated forecasts. We have explained the connections between improvements in raw ETo**  
 1034 **forecasts and calibrated ETo forecasts in response to your comment #9.**

1035 **As we introduced in the manuscript, when we calibrate the raw ETo forecasts ( $f(t)$ ), we built a**  
 1036 **conditional distribution ( $\delta(m(t))$ ) for observations ( $o(t)$ ), and 100 values will be drawn from**  
 1037 **this conditional distribution to generate the calibrated ensemble forecasts:**

1038

$$\delta(m(t)) \sim N\left(\mu_o(m(t)) + r \frac{\sigma_o(m(t))}{\sigma_f(m(t))} (f(t) - \mu_f(m(t))), (1 - r^2)\sigma_o^2\right)$$

1039 **in which where  $m(t)$  returns the month  $k$  ( $k=1$  to  $12$ ) of daily forecasts or observations of day  $t$ ;**  
 1040  **$\mu_f(m(t))$  and  $\sigma_f(m(t))$  refer to the marginal distribution's mean and standard deviation of  $f(t)$  in**  
 1041 **month  $m(t)$ , respectively;  $\mu_o(m(t))$  and  $\sigma_o(m(t))$  are the mean and standard deviation of the**

1042 marginal distribution of  $o(t)$  in month  $m(t)$ ;  $r$  is the correlation between  $f(t)$  and  $o(t)$  in the  
1043 transformed space.

1044 **As a result, when the correlation is improved, it will help improve the estimation of the mean  
1045 and standard deviation of the above conditional distributions. As a result, bias in calibrated  
1046 forecasts will be further reduced. That is why improvements in bias demonstrate a similar  
1047 spatial pattern as those of the correlation coefficient.**

1048 **To explain improvements in  $r$  in calibrated forecasts, we add the following sentence to  
1049 section 3.4:**

1050 “Spatial patterns of improvements in  $r$  in calibrated ETo forecasts (Figure 6) are consistent with  
1051 the improvements in  $r$  of raw ETo forecasts with the adoption of bias-correction (Figure 2),  
1052 particularly for the short lead times. The improvements in  $r$  of calibrated ETo forecasts (Figure  
1053 6) may also lead to more reasonable conditional distributions for a given raw forecast (equation  
1054 4). As a result, regions showing improvements in  $r$  in calibrated ETo forecasts (Figure 6) often  
1055 demonstrate reductions in absolute bias (Figure 4).”

1056

#### 1057 Point #12

1058 *P16 I320-328. Please comment on why the accuracy has larger differences in terms of geographical  
1059 patterns than for the bias and PIT performance which had very strong localised performance.*

1060 **Response: Thank you for the comments. We believe there are four reasons for the differences  
1061 in spatial patterns of CRPS skill score (Figure 8) with changes in bias (Figure 4), correlation  
1062 coefficient (Figure 6), and alpha index (Figure S13):**

1063 **1, The metrics measure different features of the quality of forecasts, and may have different  
1064 sensitivities to changes in calibrated forecasts. As a result, it is not unexpected that their  
1065 spatial patterns show differences. Bias measures average differences; correlation coefficient  
1066 shows consistency between observations and forecasts; the CRPS skill score measures the  
1067 performance of calibrated forecasts relative to the climatology forecast; the alpha index is an  
1068 indicator showing whether the distribution of calibrated forecasts is overconfident or  
1069 underconfident. As a result, improvements indicated by these metrics do not necessarily  
1070 show exactly the same spatial patterns.**

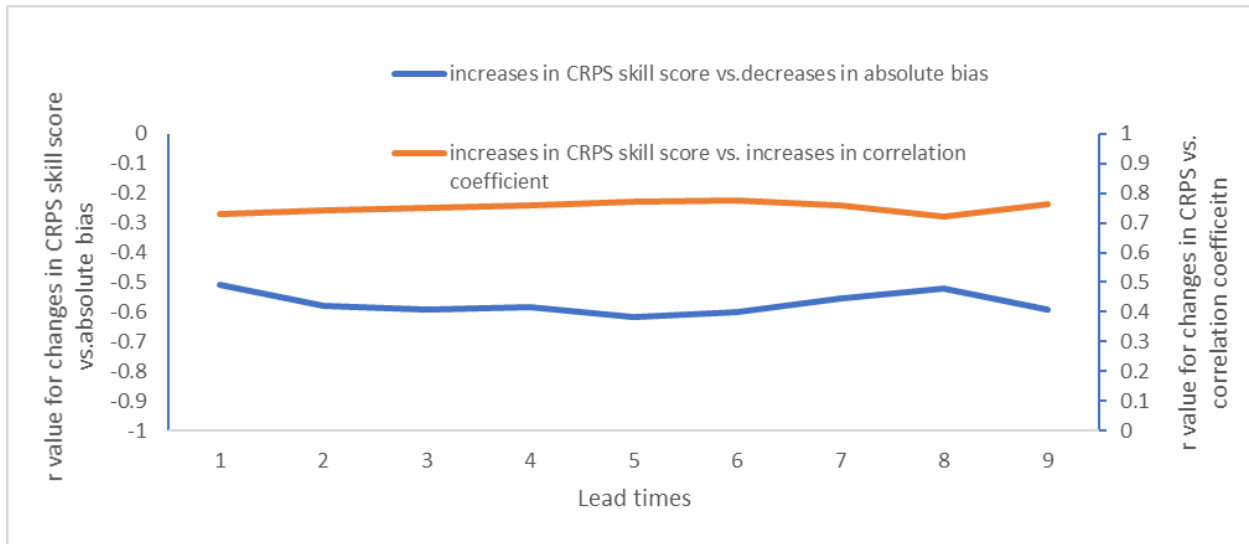
1071 **2, The alpha index is less sensitive to changes in forecasts than other metrics. It is well known  
1072 that the quality of forecasts often declines with lead time, even for calibrated forecasts. This  
1073 tendency can be seen from the correlation coefficient (Figure 5) and CRPS skill score (Figure  
1074 7). However, the same trend is not shown in the alpha index (Figure 9). As demonstrated by  
1075 figure 9, the alpha index demonstrates similar magnitudes and spatial patterns among the 9  
1076 lead times. As was introduced in equations 13 and 14, PIT value and alpha index are mainly  
1077 used to measure the consistency between distributions of forecasts and observations.**

1078 **Improvements achieved through the adoption of calibration strategy ii (e.g., Calibrations 2**



1079 and 4) may not significantly change the statistical distributions of the calibrated forecasts, as  
1080 evidenced by the *t-test* (Table S2). As a result, differences in the alpha index (Figure S13)  
1081 between Calibrations 2 and 1 do not show spatial patterns resembling absolute bias (Figure  
1082 4), correlation coefficient (Figure 6), and CRPS skill score (Figure 8).

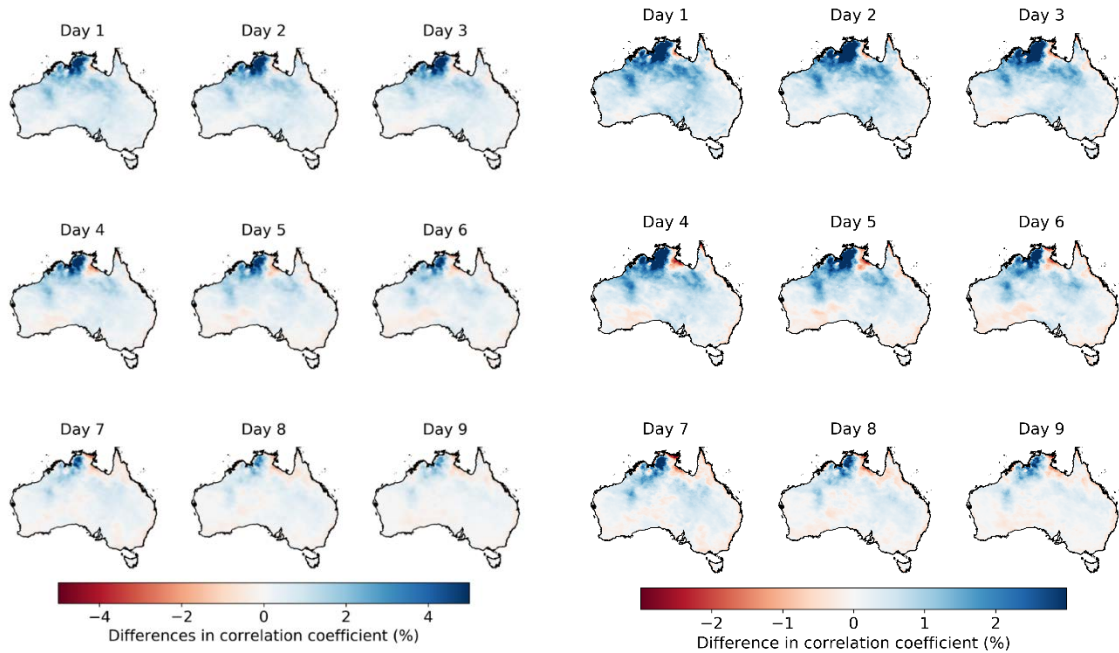
1083 **3, The spatial patterns of improvements in absolute bias, correlation coefficient, and CRPS**  
1084 **skill score are generally consistent. We calculate the spatial correlation for changes in CRPS**  
1085 **skill score vs. changes in absolute bias (Figure 8 vs. Figure 4), and the spatial correlation for**  
1086 **changes in CRPS skill score vs. changes in correlation coefficients (Figure 8 vs. Figure 6). As is**  
1087 **shown in the following figure, the metrics show high spatial correlation.**



1088 **4, The upper and lower limits used for the maps may have affected our understanding of the**  
1089 **spatial patterns of the evaluation metrics. Following comparison shows that when using**  
1090 **narrower limits (-3% to 3%, rather than -5% to 5%) for the color bar of the maps showing**  
1091 **improvements in correlation coefficients (the subplot on the right), the spatial pattern looks**  
1092 **more consistent with the maps showing increases in CRPS skill score (Figure 8). In the revised**  
1093 **manuscript, we use the plot with narrower color bar limits.**

1095

1096



1097

1098

1099

1100 **To explain spatial patterns of the evaluation metrics, we add a new subsection to the Results section**  
 1101 **(3.8 Summary of results):**

1102 “Although the selected metrics measure different aspects of forecast quality, they generally agree with  
 1103 each other in demonstrating improvements in calibrated ETo forecasts with the adoption of the Strategy ii.  
 1104 As introduced in the Method section, bias measures average differences; correlation coefficient shows  
 1105 consistency between observations and forecasts in temporal variability; the CRPS skill score measures the  
 1106 performance of the calibrated forecasts relative to climatology forecast; the  $\alpha$  index is an indicator  
 1107 showing whether the distribution of calibrated forecasts is overconfident or underconfident. As a result,  
 1108 these metrics may differ from each other in magnitude when used to evaluate different calibrations  
 1109 (Figures 4, 6, 8, and S14). However, improvements in bias, correlation, and skills with the adoption of  
 1110 bias-correction to input variables are generally consistent in spatial patterns. Compared with the other  
 1111 three metrics, the  $\alpha$  index demonstrates less significant changes when input variables are bias-corrected  
 1112 first (Table S2 and Figure S14), mainly because this index is less sensitive to changes in calibrated  
 1113 forecasts than other metrics.”

1114

1115 **Point #13**

1116 *P16 I329: Results on calibration 2 and 4: what is the comparison between 2 and 4? Why are these only*  
 1117 *addressed in the evaluation of forecast accuracy section? Why is there no mention of these for the bias*  
 1118 *and reliability evaluation? I suggest changing the section order and moving this section first. Then, add a*  
 1119 *sentence in the bias and reliability section to explicitly communicate what results of experiment 3) and 4)*  
 1120 *are not presented and why.*

1121 **Response: Thank you for the valuable suggestions. We check the original submission and**  
1122 **believe your comments refer to Calibrations 3 and 4 here.**

1123 **As we explain in our response to your comment #5, calibrations 3 and 4 are to further confirm**  
1124 **that whether our strategy is suitable for general application. We further explain the reason of**  
1125 **by adding the following sentences to clarify why Calibrations 3 and 4 are included in this**  
1126 **study in Method:**

1127 “The comparison between Calibrations 1 and 2 is to investigate whether the bias-correction of input  
1128 variables would further improve ETo forecasts when the calibration is conducted based on ETo anomalies  
1129 and climatological mean. We also conduct additional calibrations which post-process ETo forecasts  
1130 directly (Calibrations 3 and 4), to test whether the contribution of improving input variables to ETo  
1131 forecast calibration, if there is any, will depend on how ETo forecasts are calibrated (based on anomalies  
1132 vs. based on ETo). Calibrations 3 and 4 will help evaluate the general applicability of strategy ii to  
1133 enhance NWP/GCM-based ETo forecasting. Key steps of the four calibrations could be found in the  
1134 schematic diagram introducing how raw ETo forecasts are constructed and how calibrations are  
1135 conducted (Figure S1). In the main text, we primarily analyze results from Calibrations 1 and 2.  
1136 Improvements with the adoption of bias-correction to input variables in Calibrations 3 and 4 are very  
1137 similar to Calibrations 1 and 2 (see the Supplementary Material). To avoid redundancy, we mainly  
1138 present results from Calibrations 3 and 4 in the Supplementary Material.”

1139 **As we introduced in our response to your comment point #5, we add more results (bias,**  
1140 **correlation, and alpha-index) from Calibrations 3 and 4 to the Supplementary Material. We**  
1141 **also add one new subsection (3.7) to briefly introduce the results shown in these figures**  
1142 **(Figures S15-S18).**

1143

#### 1144 Point #14

1145 *Discussion:*

1146 *There are little to no direct comparison of results and calibration work presented here to any previous*  
1147 *methods or studies (which were mentioned in the introduction). To address a research closure, please put*  
1148 *the work presented in this paper in context with other studies applying strategy 1 and strategy 2.*

1149 **Response: We appreciate the reviewer’s valuable suggestion. We explain in detail why we do**  
1150 **not compare our calibration directly with calibrations using other models in our response to**  
1151 **your comment #3. However, we totally agree with the reviewer that it is necessary to**  
1152 **compare our results with previous investigations in ETo forecasting to help the audience**  
1153 **better understand the performance of our calibration. Therefore, we add the following**  
1154 **contents to the Discussion:**

1155 “This investigation further highlights the importance of statistical calibration in NWP-based ETo  
1156 forecasting (Medina and Tian, 2020). According to an investigation across 40 sites in Australia,  
1157 raw ETo forecasts constructed with NWP outputs reasonably captured the magnitude and  
1158 variability of ETo, but forecast skills better than climatology were only limited to the first 6 lead

1159 times (Perera et al., 2014). Our investigation suggests that statistical calibration could  
1160 substantially improve forecast skills and successfully extend the skillful forecasts to lead time 9  
1161 across Australia. Findings of this investigation agree well with the site-scale short-term ETo  
1162 forecasting based on GCM outputs (Zhao et al., 2019a) in the improvements of forecast skills  
1163 through statistical calibration. Calibrated forecasts from Calibration 2 demonstrate similar skills  
1164 as Zhao et al. (2019a) across three Australian sites. Thanks to the capability of SCC in  
1165 calibrating short-archived forecasts (Wang et al., 2019), we achieve the improvements based on  
1166 much shorter archived raw forecasts (3-year vs. 23-year) than Zhao et al. (2019a). Calibrated  
1167 forecasts from Calibration 2 also demonstrate low biases (0.32-0.95%) comparable with  
1168 calibrated ETo forecasts (0.49-0.63%) based on the Bayesian Model Averaging (BMA) model  
1169 and weather forecasts from three NWP models in the U.S. during 2014-2016 (Medina and Tian,  
1170 2020)."

1171 **In addition, we also highlight the importance of testing the proposed calibration strategy**  
1172 **(strategy ii) in the future in section 4.2:**

1173 "Third, further investigations based on other calibration models are needed to validate findings  
1174 of this investigation. Our analyses based on two different methods (based on ETo anomalies vs.  
1175 based on original ETo) demonstrate similar improvements in calibrated ETo forecasts with the  
1176 adoption of bias-correction to input variables. Additional evaluations will be needed to verify  
1177 whether forecast skills will be improved using strategy ii but based on a different calibration  
1178 model."

1179

#### 1180 Point #15

1181 *It is unclear whether authors recommend the use of experiment 2) or 4), when and why. In that sense, I*  
1182 *question again the inclusion of these experiments without further elaborating and discussing these*  
1183 *results.*

1184 **Response: Thank you for the valuable suggestion. As we explain in our response to your**  
1185 **comments #5 and #13, the objective of this study is to evaluate the necessity of correcting the**  
1186 **input variables prior to ETo forecast calibration. We also further explain that including**  
1187 **Calibrations 3 and 4 was to further evaluate whether the strategy could be generally applied**  
1188 **to other calibration models. In addition, we add results from Calibrations 3 and 4 and**  
1189 **discussed implications from these two calibrations (section 3.7):**

1190 "We also compare the bias, correlation coefficient, CRPS skill score, and reliability of calibrated forecasts  
1191 from Calibrations 3 and 4, to evaluate whether we can obtain similar improvements through the bias-  
1192 correction of input variables if we conduct the ETo forecast calibration in a different way (without using  
1193 ETo climatological mean and anomalies). Results show that the adoption of bias-correction also leads to  
1194 lower bias, higher correlation coefficient, and higher CRPS skill score in terms of magnitude, spatial  
1195 patterns, and trend along the lead times, when ETo forecasts are calibrated directly (Figure S15-S17). In  
1196 addition, the alpha index was only slightly different between Calibrations 3 and 4 (Figure S18). This

1197 additional comparison further confirms the general applicability of strategy ii for enhancing NWP-based  
1198 ETo forecasting.”

1199 Point #16

1200 Structure:

1201 *The introduction is well structured and appropriately present previous work studies and existing*  
1202 *strategies.*

1203 **Response: We appreciate your constructive comments.**

1204

1205 Point #17

1206 *The title is a bit lengthy, authors could consider shortening it.*

1207 **Response: We change the title from:**

1208 “Bias-correcting input variables prior to combined calibration leads to more skillful forecasts of  
1209 reference crop evapotranspiration”

1210 **to:**

1211 "Bias-correcting input variables enhances forecasting of reference crop evapotranspiration."  
1212

1213 Point #18

1214 *As noted above, I suggest authors consider the order of results presented in the context of results from*  
1215 *experiment 3) and 4).*

1216 **Response: As we explained in our response to your comments #5, #13, and #15, we add a**  
1217 **new subsection (3.7) to present results from calibrations 3 and 4 and discuss the implications**  
1218 **of these two Calibrations.**

1219

1220 Point #19

1221 *Minor comments:*

1222 *P4 l106: I suggest adding a diagram clearly explaining steps and differences of procedure between the*  
1223 *calibration experiments.*

1224 **Response: We appreciate the valuable suggestions and create a diagram to show the key**  
1225 **steps of the four calibrations**

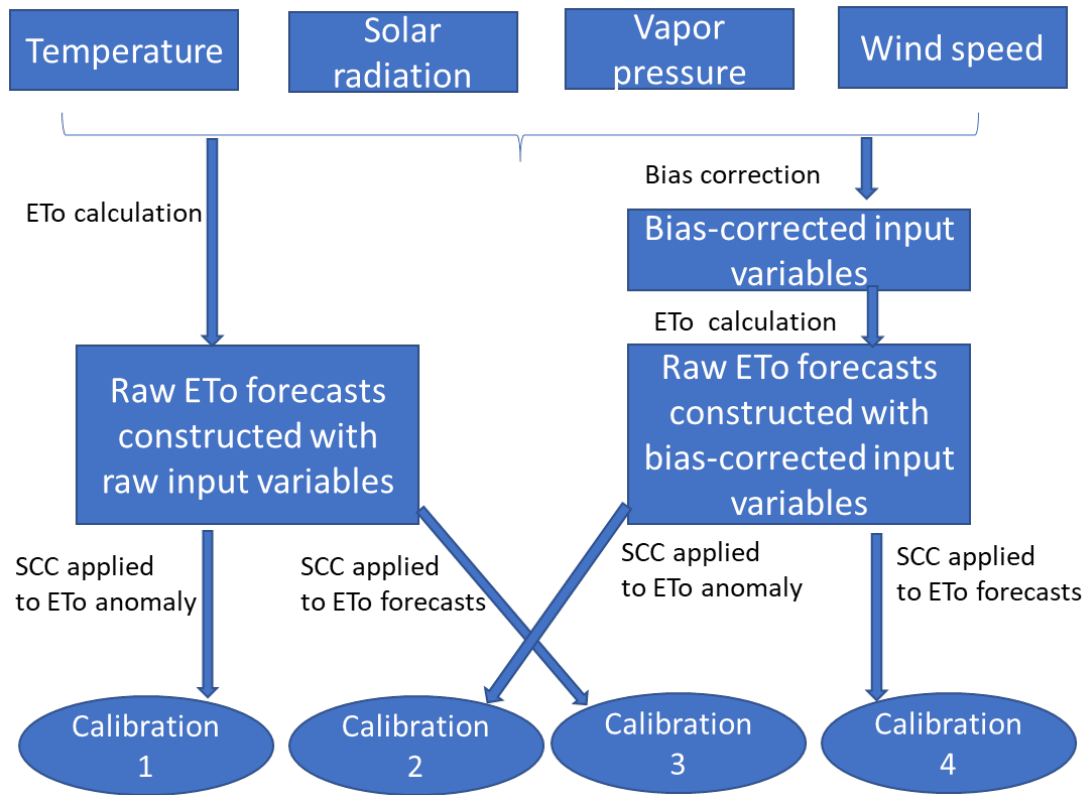


Figure S1. Schematic of the four calibrations

1226

1227

1228

1229 Point #20

1230 P3 I68: '...pressing need to investigate.' Please expand why it is pressing?

1231 **Response: Thank you for the comments. ETo forecasts have been increasingly used in the**  
 1232 **planning of farming activities (e.g., amount and timing of irrigation) in Australia. We improve**  
 1233 **this sentence as follows:**

1234 “Since NWP/GCM-based ETo forecasting is increasingly conducted to support water resource  
 1235 management, there is a need to investigate the necessity of correcting raw forecasts of the input variables  
 1236 in ETo forecast calibration.”

1237

1238 Point #21

1239 P3 I74: Calibrate should be calibrate with small cap letter.

1240 **Response: Thank you for the careful review. We correct this typo.**

1241

1242 Point #22

1243 *P3 l80-84: There are many efforts to develop downscaling methods, please comment on what was been*  
1244 *done here to downscale ACCESS-G2 to the AWAP grid. Why not scaling AWAP to the match the forecast*  
1245 *grid?*

1246 **Response: Thank you for the valuable suggestions. In the revised manuscript, we further**  
1247 **introduce that we used bilinear interpolation to remap ACCESS-G2 forecasts. Meanwhile, we**  
1248 **agree with the reviewer that sophisticated methods have been developed to downscale**  
1249 **coarse resolution forecasts to match observations.**

1250 **In this study, the purpose of the regridding is to connect forecasts with the corresponding**  
1251 **observations so we can calibrate the forecasts, rather than trying to reconstruct the spatial**  
1252 **patterns of forecasts at a finer scale.**

1253 **We conducted a literature review on the remapping methods used in forecasts post-**  
1254 **processing. It is common that raw forecasts and references data have different spatial**  
1255 **resolutions. We found that bilinear interpolation of forecasts from a coarser resolution to a**  
1256 **finer resolution has been widely used in forecast post-processing and verification. For**  
1257 **example, Hamill et al. (2015) used bilinear interpolation to downscale the resolution of**  
1258 **Global Ensemble Forecast System (GEFS) forecasts from 1° to 1/8° to match observations**  
1259 **before post-processing with an analogy-based model. Yuan et al. (2014) used bilinear**  
1260 **interpolation to remap the Global Ensemble Forecast System (GEFS, with resolutions of**  
1261 **~0.469° and ~0.625°) to match the North-American Land Data Assimilation System (NLDAS,**  
1262 **with the resolution of 1/8°), before the forecasts were post-processed with a quantile**  
1263 **mapping method. Zeng and Yuan (2018) used bilinear interpolation to remap sub-seasonal to**  
1264 **seasonal forecasts from ECMWF (0.25°X0.25°to 0.5°X0.5° for different lead times), NCEP**  
1265 **(1°X1°), China Meteorological Administration (CMA, 1°X1°), Hydrometeorological Centre of**  
1266 **Russia (HMCR, 1.1°X1.4°), and Australian Bureau of Meteorology (BoM, 2°X2°) to a common**  
1267 **resolution of 0.7°, in order to match the reanalysis data. James et al. (2017) regridded the**  
1268 **wind forecasts with bilinear interpolation from the 3-km High-Resolution Rapid Refresh**  
1269 **(HRRR) NWP model to an observation tower in Colorado to evaluate forecast quality. Bowler**  
1270 **et al. (2008) interpolated the ECMWF forecasts with a grid spacing of 1.5° bilinearly to the site**  
1271 **scale for forecast verification. Yuan and Wood (2012) used bilinear interpolation to match**  
1272 **forecasts from the Euro- Mediterranean Centre for Climate Change (CMCC-INGV), the**  
1273 **European Centre for Medium-Range Weather Forecasts (ECMWF), the Leibniz Institute of**  
1274 **Marine Sciences at Kiel University (IFM-GEOMAR), Météo France, and UK Met Office (UKMO),**  
1275 **which have a spatial resolution of 2.5° to match the observation of 1°.**

1276 **As a result, previous investigations suggested that downscaling with a sophisticated method**  
1277 **could potentially be useful, but that is not necessarily essential in forecast post-processing,**  
1278 **and bilinear interpolation is acceptable.**

1279 **However, we agree with the reviewer that whether a better remapping method will further**  
1280 **improve the forecast calibration should be investigated in the future. Therefore, we add the**  
1281 **following sentence to section 4.2 (Implications for forecasting of integrated variables and**  
1282 **future work):**

1283 “More sophisticated remapping methods should be evaluated to understand the impacts of forecast  
1284 regridding on statistical calibration.”

1285

1286 **Reference:**

1287 Bowler, N.E., Arribas, A., Mylne, K.R., Robertson, K.B., Beare, S.E., 2008. The  
1288 MOGREPS short-range ensemble prediction system. *Q. J. R. Meteorol. Soc.* 722,  
1289 703–722. <https://doi.org/10.1002/qj>

1290 Hamill, T., Scheuerer, M., Bates, G., 2015. Analog Probabilistic Precipitation  
1291 Forecasts Using GEFS Reforecasts and Climatology-Calibrated Precipitation  
1292 Analyses. *Mon. Weather Rev.* 143, 3300–3309. [https://doi.org/10.1175/MWR-](https://doi.org/10.1175/MWR-D-15-0004.1)  
1293 [D-15-0004.1](https://doi.org/10.1175/MWR-D-15-0004.1)

1294 James, E.P., Benjamin, S.G., Marquis, M., 2017. A unified high-resolution wind and  
1295 solar dataset from a rapidly updating numerical weather prediction model.  
1296 *Renew. Energy* 102, 390–405. <https://doi.org/10.1016/j.renene.2016.10.059>

1297 Monteiro, J.A.F., Strauch, M., Srinivasan, R., Abbaspour, K., Gucker, B., 2016.  
1298 Accuracy of grid precipitation data for Brazil : application in river discharge  
1299 modelling of the Tocantins catchment. *Hydrol. Process.* 30, 1419–1430.  
1300 <https://doi.org/10.1002/hyp.10708>

1301 Yuan, X., Wood, E.F., 2012. On the clustering of climate models in ensemble  
1302 seasonal forecasting. *Geophys. Res. Lett.* 39, 1–7.  
1303 <https://doi.org/10.1029/2012GL052735>

1304 Yuan, X., Wood, E.F., Liang, M., 2014. Integrating weather and climate prediction:  
1305 Toward seamless hydrologic forecasting. *Geophys. Res. Lett.* 5891–5896.  
1306 <https://doi.org/10.1002/2014GL061076>.Received

1307 Zeng, D., Yuan, X., 2018. Multiscale Land – Atmosphere Coupling and Its  
1308 Application in Assessing Subseasonal Forecasts over East Asia. *J.*  
1309 *Hydrometeorology* 19, 745–760. <https://doi.org/10.1175/JHM-D-17-0215.1>

1310



1311 Point #23

1312 *P4 I100: please add a comment that SCC model will be described in section 2.3.2*

1313 **Response: We added the following sentence to this section:**

1314 “The calibration model used in this study is the Seasonally Coherent Calibration (SCC) model, which is  
1315 introduced in detail in section 2.3.2.”

1316

1317 Point #24

1318 *P5 I134 climatological means or mean? Please rephrase and clarify this sentence.*

1319 **Response: Thank you, and we change it to ‘climatological mean’.**

1320

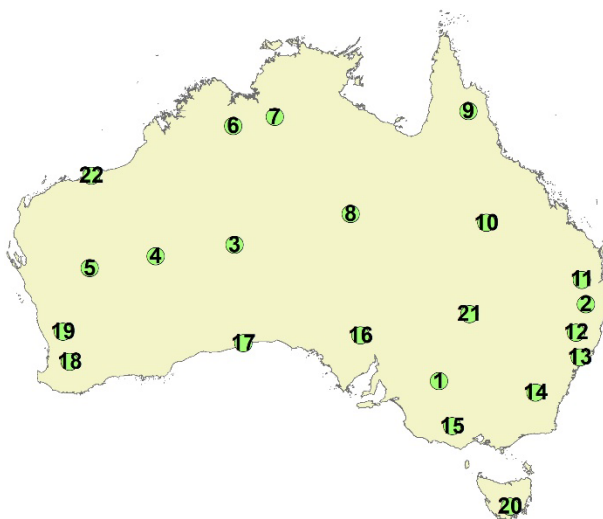
1321 Point #25

1322 *P6 I165: Why are only 100 members drawn, is there any difference with a varying number of ensemble  
1323 members for forecast reliability?*

1324 **Response: Thank you for the comments. We use 100 members because the computation cost  
1325 is more affordable than using a larger ensemble size.**

1326 **In order to evaluate how different ensemble sizes would affect the reliability and skills of  
1327 forecasts, we choose 22 sites randomly across Australia and compare the alpha index and  
1328 CRPS skill score across these sites using 100, 500, and 1000 ensemble members. The  
1329 following map shows the locations of the 22 sites.**

1330



1331

1332

1333

1334 The following figure shows the alpha index is almost identical across the selected sites for the  
1335 three ensemble sizes:

1336

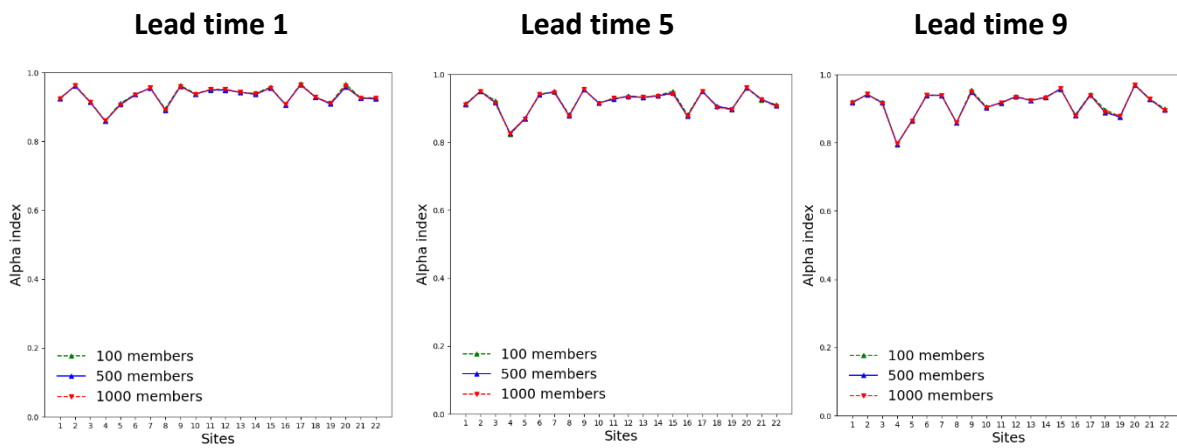
1337

1338

1339

1340

1341

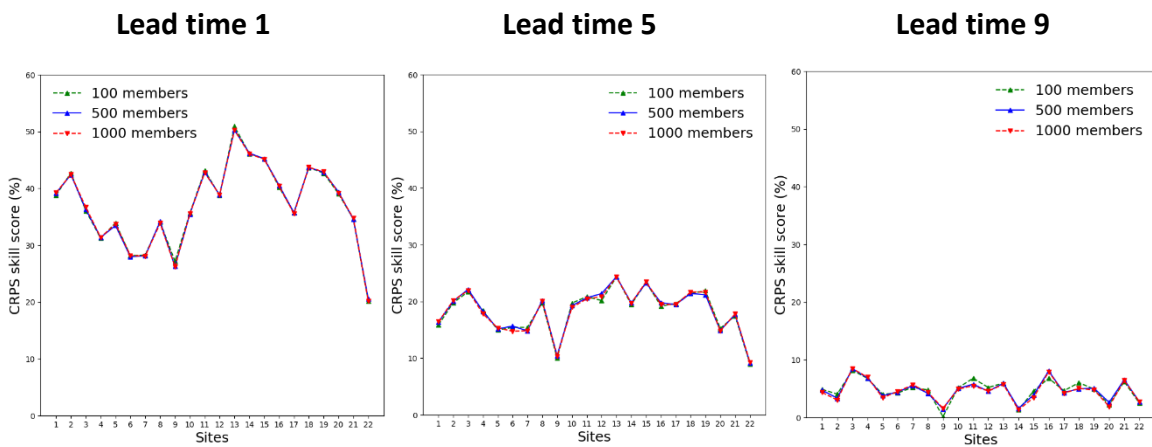


1342

1343

1344 Comparison of CRPS skill score shows that different ensemble sizes have negligible impacts  
1345 on the score:

1346



1347

1348

1349 As a result, we conclude that the ensemble size used in this study is reasonable.

1350

1351 Point #26

1352 *Is there a need or a reason to verify accumulated Eto forecast values across lead times (as is often the*  
1353 *case for streamflow forecasting)? Please comment.*

1354 **Response: Thank you for the comments. For short-term weather forecasts, which are issued**  
1355 **on a daily basis, users are often interested in the short-lead-time forecasts (e.g., lead times 1**  
1356 **to 3). Accumulated forecasts across all lead times will not provide the information that users**  
1357 **are particularly interested.**

1358 **In addition, the evaluation by lead time shows that improvements with the adoption of the**  
1359 **new calibration strategy (Calibrations 2 and 4) decrease with lead time, but still show better**  
1360 **performance than the calibrations (Calibrations 1 and 3) based on raw forecasts of input**  
1361 **variables, event at lead time 9. As a result, we are confident that evaluation based on**  
1362 **accumulated ETo will not change the conclusion of this study.**

1363 Point #27

1364 *P8 I225: 'wind speed is higher than 1m/s than the reference in Australia'. Could you please translate that*  
1365 *in terms of percentage so that this statement can be more easily compared to other locations.*

1366 **Response: We add more quantitative information in the evaluation of raw forecasts of input**  
1367 **variables and use percentage to measure the changes:**

1368 “The daily minimum temperature (Tmin) is underpredicted by more than 1.5 °C in western and  
1369 central parts of Australia by the raw forecasts, but is overpredicted by ca. 1 °C in eastern and  
1370 southern Australia. Vapor pressure is underpredicted in western and central regions by ca.14%,  
1371 but is overpredicted by ca. 6% in coastal areas of southeastern Australia by the raw forecasts.  
1372 Raw solar radiation forecasts are about 5% higher than AWAP data across Australia. Forecasted  
1373 wind speed is higher than the reference data by more than 1 m s-1 (or by ca. 63%) in most parts  
1374 of Australia. For each input variable, spatial patterns of biases in raw forecasts are consistent  
1375 across the 9 lead times, demonstrating systematic errors in the raw NWP forecasts.”

1376

1377 Point #28

1378 *P18 I380' NWP outputs have been increasingly used for ETo forecasting.' For which applications? Please*  
1379 *finish the sentence.*

1380 **Response: We modify this sentence as follows:**

1381 " NWP outputs have been increasingly used for ETo forecasting to support water resource management."

1382

1383 Point #29

1384 *P18 I385 Addition 'of' in ... skill 'of' the calibrated ETo forecasts.*

1385 **Response: We add the missing 'of' to this sentence:**

1386 "With this extra step, the bias, correlation coefficient, and skills of the calibrated ETo forecasts  
1387 are all improved."

1388

1389 Point #30

1390 *References:*

1391 *Pappenberger, F., M. H. Ramos, H. L. Cloke, F. Wetterhall, L. Alfieri, K. Bogner, A. Mueller and P. Salamon*  
1392 *(2015). "How do I know if my forecasts are better? Using benchmarks in hydrological ensemble*  
1393 *prediction." Journal of Hydrology 522: 697-713.*

1394 **Response: We cite this paper in the revised manuscript in introducing the CRPS skill score.**

1395

1396

1397