

# Responses to Reviewer #2

1

## 2 Point #1

### 3 1. General comments

4 *This paper presents a method to improve monthly seasonal forecasts of potential evapotranspiration*  
5 *using a trend-aware statistical model (BJP-tj). This model builds on previous work by the authors on the*  
6 *BJP model combining a data transform with a multivariate normal distribution.*

7 *The topic of trend-aware forecasts is fundamental in a changing climate where the use of long historical*  
8 *time series to calibrate statistical forecast and post-processing models becomes questionable. This paper*  
9 *provides a valuable contribution to the field by showing how an existing statistical model can be*  
10 *extended to include trends with limited additional complexity. The model performance is thoroughly*  
11 *analysed using well established metrics that target a wide range of forecast attributes. Finally, the paper*  
12 *is well written, clear and to the point concise, with figures that provide strong visual evidence to support*  
13 *the authors' analysis.*

14 *We do not see any major issues with the paper and recommend it to be published with minor revisions.*  
15 *The two main items that could be improved by the authors aside of the detailed points raised in the*  
16 *following section relate to:*

17 **Response: We appreciate the excellent summary and assessment. We addressed the valuable**  
18 **comments carefully and provided point-by-point responses. Please see details as follows.**

19

## 20 Point #2

21 *[cross validation scheme] The authors used a traditional leave-one-out cross validation scheme where a*  
22 *single month is left aside for validation and the model is calibrated against the remaining data points.*  
23 *This an optimistic cross-validation scheme because the validation month is likely to show a similar trend*  
24 *and is not completely independent from the calibration data. A more conservative approach would be to*  
25 *split the data set in two parts, although this would not solve the problem completely. This an important*  
26 *issue but would require complex theoretical developments that are probably beyond the scope of this*  
27 *paper. However, we recommend a bit of discussion around this point.*

28 **Response: Thank you for the valuable comments and suggestions. We agree with the**  
29 **reviewer that the current leave-one-out cross-validation strategy is not perfect for inferring**  
30 **the trend parameters. As we introduced in the Method section (equation 5), the two trend**  
31 **parameters are inferred together with parameters (mean vector and covariance matrix)**  
32 **defining the bivariate distribution. The current strategy (leave-one-out) has been proven**  
33 **effective for the inferencing mean vector and covariance matrix, but may not be good enough**  
34 **for the inference of trend parameters, since the left-out month may not be fully independent**  
35 **of the remaining 29 months. Leaving out longer years, such as splitting the 30-year data**

36 **equally into two parts, as the reviewer suggested, could alleviate this problem to some**  
37 **extent. However, we have another concern about data splitting. This strategy will**  
38 **substantially reduce samples used for parameter inference from 29 to 15, and thus may lead**  
39 **to significant sampling problems. We feel solving this problem may need additional efforts**  
40 **and more sophisticated solutions. As a result, we highlight this as a challenge that should be**  
41 **addressed in our future work in section 4.3 (Future work)**

42 “First of all, more sophisticated cross-validation methods should be developed for the inference  
43 of trend parameters. The current leave-one-out method has been proven to be effective in the  
44 inference of the mean vector and covariance matrix (Shao et al., 2020). However, this strategy  
45 may not guarantee the independence between the left-out data and data used for the inference of  
46 trend parameters. We decided not to implement the data-splitting method for cross-validation  
47 because of the risk of introducing sampling errors. Future investigations should take this  
48 challenge into consideration and develop better cross-validation methods for the inference of  
49 trend parameters.”

50

### 51 Point #3

52 *[risk of overfitting when there is no observed trend] The authors demonstrate that the BJP-ti model*  
53 *outperforms BJP and raw forecasts when there is a trend in observed data. However, some of the results*  
54 *shown by the authors suggest that its performance is worse than BJP when the trends are not significant.*  
55 *This result is to be expected because of the higher number of parameters of BJP-ti which may increase*  
56 *the risk of overfitting and counter-performance over validation data. We recommend highlighting this*  
57 *point in the manuscript to better identify the strengths and weaknesses of BJP-ti.*

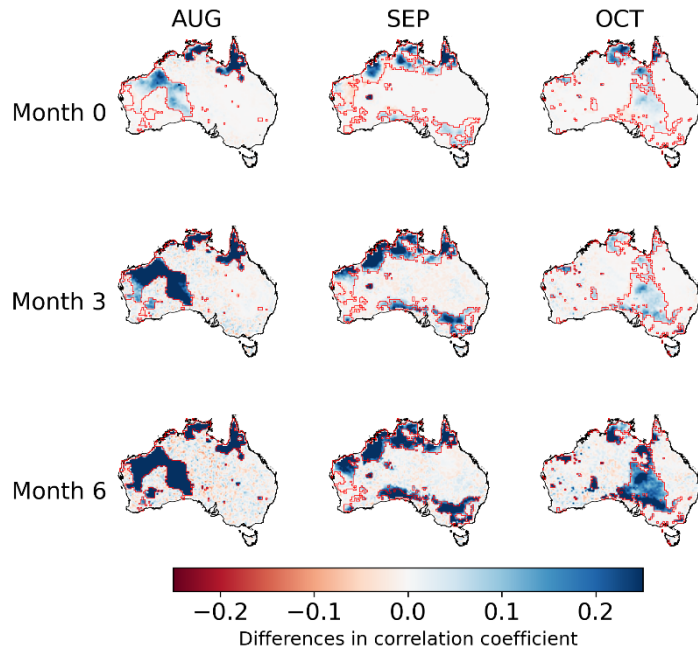
58 **Response: We agree with the reviewer that the original BJP-ti parameterization suffered from**  
59 **parameter overfitting and resulted in degradation in performance when compared with the**  
60 **BJP model.**

61 **To solve this problem, we accept your valuable suggestions and add limits to inferred trends**  
62 **in trend reconstruction. Specifically, we use  $P < 0.05$  as the threshold to define statistically**  
63 **significant trends. For trends that are statistically insignificant ( $P > 0.05$ ), we set the inferred**  
64 **trends to zero to avoid overfitting:**

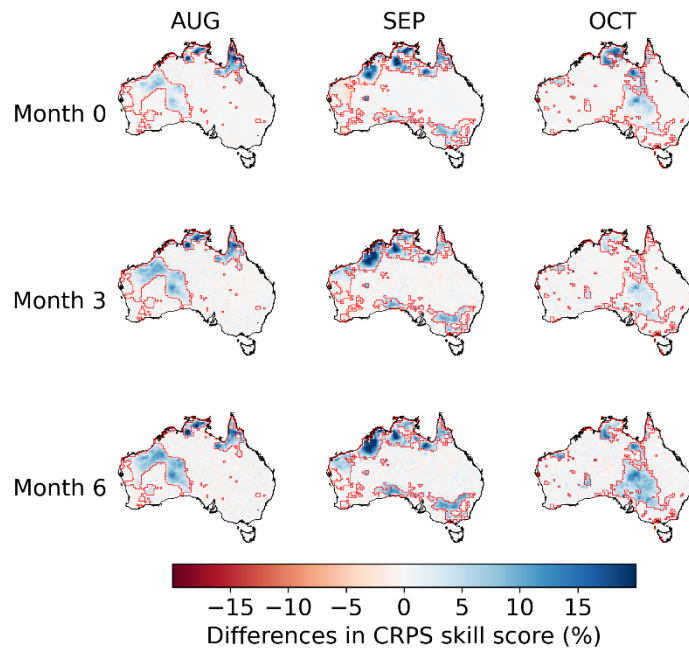
65 “For trends that are insignificant ( $P > 0.05$ ), we set  $m_i$  to 0 to avoid overfitting trends in calibrated  
66 forecasts. For significant trends, we set the  $m_i$  value based on observations and raw forecasts  
67 during 1981-2019”

68 **This new strategy is not only effective in limiting the trend reconstruction to regions with**  
69 **significant observed trends (Figure 1), but also avoids the reductions in correlation**  
70 **coefficients (Figure 2) and CRPS skill score (Figure 3) following trend reconstruction.**

71 **We have updated all results based on the new calibration. We present Figures 2 and 3 here to**  
72 **show the advantage and effectiveness of the new strategy.**



74 **Figure 2. Differences in the correlation coefficient ( $r$ ) between BJP-ti calibrated forecasts and**  
 75 **observations with that between BJP calibrated forecasts and observations for three selected months**  
 76 **(AUG, SEP, OCT) and three lead times (Months 0, 3, and 6). Red polygons show regions with significant**  
 77 **trends.**



79  
 80 **Figure 3. Differences in CRPS skill score between BJP-ti calibrated forecasts and the BJP calibrated**  
 81 **forecasts for three selected months (AUG, SEP, OCT) and three lead times (Months 0, 3, and 6). Red**  
 82 **polygons show regions with significant observed trends.**

83 Point #4

84 [Line 26] “Reference crop evapotranspiration (ET<sub>o</sub>) measures the evaporative demand of the  
85 atmosphere”: Please provide additional details regarding the definition of ET<sub>o</sub>. We suggest the following:  
86 “Reference crop evapotranspiration (ET<sub>o</sub>) measures the evaporative demand of the atmosphere for a  
87 hypothetical crop of given height, with defined surface resistance factor and albedo. It is generally  
88 computed using the Penman-Monteith equation following Allen et al. (1998, see section 2.1), which is  
89 known as FAO56. McMahon et al. (2013) provides additional information about the process. ”

90 **Response: Thank you for your valuable suggestions. We add the suggested introduction of ET<sub>o</sub>**  
91 **and the reference to the manuscript.**

92 **Reference:**

93 McMahon T.A., Peel, M. C., Lowe, L., Srikanthan, R. and McVicar, T.R.: Estimating actual, potential,  
94 reference crop and pan evaporation using standard meteorological data: A pragmatic synthesis. Hydrol.  
95 Earth Syst. Sci., 17, 1331–1363, doi: /10.5194/hess-17-1331-2013, 2013

96

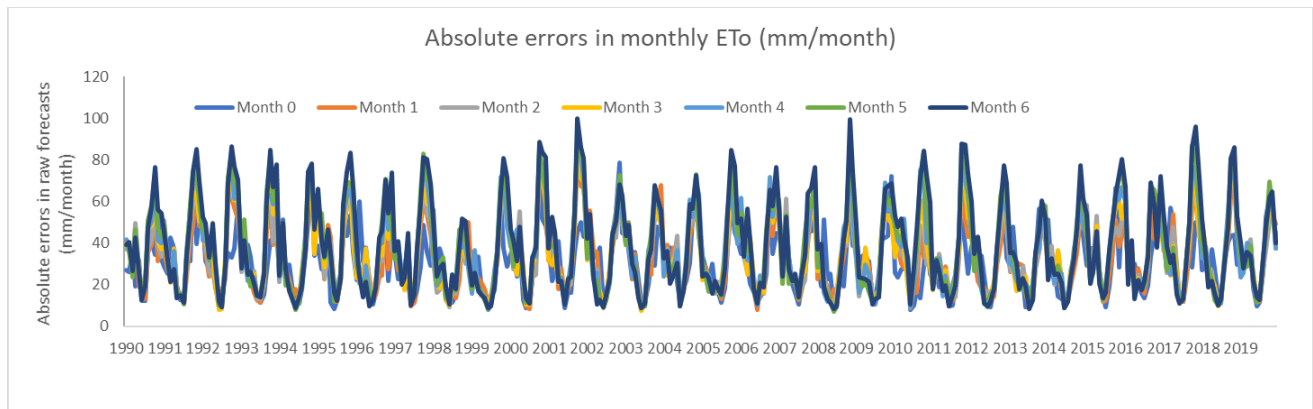
97 Point #5

98 [Line 94] “we combine the archived re-forecasts and operational forecasts”: Please comment briefly on  
99 the potential differences in skill between the re-forecast and operational data aside of the number of  
100 ensembles generated.

101 **Response: Thank you for the suggestions. According to the ECMWF SEAS5 documentations**  
102 **(Stockdale et al., 2017; Johnson et al., 2019), SEAS5 runs for the re-forecast and operational**  
103 **forecasts periods were configured as similar as possible to maintain consistencies. However,**  
104 **there are some slight differences. In addition to ensemble size, initial conditions for the two**  
105 **sets of runs are from different data sources. As a result, performance during the two periods**  
106 **may vary for some weather variables. For example, according to the ECMWF user guide**  
107 **(ECMWF 2021), because of the different initializations, ‘the real-time forecasts of Lake**  
108 **Superior (including the Great Lakes and the Caspian Sea) are cooler in the summer than the re-**  
109 **forecasts were’. In addition, according to the latest evaluation of the SEAS5 forecasts (Figure**  
110 **40 in Haiden et al., 2021), forecasts of accumulated cyclone energy for the Atlantic tropical**  
111 **storm demonstrate larger errors during 2016-2021 than the re-forecasts.**

112 **However, we feel it is hard to draw a conclusion on the relative performance of the re-**  
113 **forecasts and operational forecasts, because they have different lengths and cover different**  
114 **years, and their performances may vary with the ECMWF output variables.**

115 **In addition, we did not see significant differences in absolute errors in raw ET<sub>o</sub> forecasts**  
116 **during the re-forecast period (1990-2016) vs. operational forecasts (2017-2019). We**  
117 **calculated the average absolute errors in raw ET<sub>o</sub> forecasts across Australia during the study**  
118 **period (1990-2019). The absolute errors during the re-forecasts and real-time periods seem to**  
119 **be comparable. We added the following figure to the Supplementary material.**



121 Figure S1. Absolute errors in raw ECMWF ETo forecasts.

122

123 **Based on these investigations, we modified the introduction of the re-forecast and**  
 124 **operational forecasts as follows:**

125 “To match  $ET_0$  observations, we combine the archived re-forecasts and operational forecasts to derive  
 126 raw  $ET_0$  forecasts for the period of 1990-2019. ECMWF runs for the two sets of forecasts are configured  
 127 in a similar way, except for differences in initialization (Johnson et al., 2019). Absolute errors in raw  $ET_0$   
 128 forecasts during the two periods are comparable (Figure S1). We choose the first 25 ensemble members  
 129 of the real-time forecasts (2017-2019) to match the ensemble size of the re-forecasts (1990-2016).”

130

131 **Reference:**

132 Stockdale, T., Johnson, S., Ferranti, L., Balmaseda, M. and Briceag, S.: ECMWF’s new long-range  
 133 forecasting system SEAS5. Meteorology section of ECMWF Newsletter No. 154., 2017.

134 Johnson, S. J., Stockdale, T. N., Ferranti, L., Balmaseda, M. A., Molteni, F., Magnusson, L., Tietsche, S.,  
 135 Decremer, D., Weisheimer, A., Balsamo, G., Keeley, S. P. E., Mogensen, K., Zuo, H. and Monge-sanz, B.  
 136 M.: SEAS5 : the new ECMWF seasonal forecast system, Geosci. Model Dev., 12, 1087–1117, 2019.

137 ECMWF. SEAS5 user guide. Version 1.2, March 2021.  
 138 [https://www.ecmwf.int/sites/default/files/medialibrary/2017-10/System5\\_guide.pdf](https://www.ecmwf.int/sites/default/files/medialibrary/2017-10/System5_guide.pdf)

139 Haiden, T., Janousek, M., Vitart, F., Ben-Bouallegue Z., Ferranti, L. and Prates, F.: Evaluation of  
 140 ECMWF forecasts, including the 2021 upgrade. Technical Memo 884. 2021.  
 141 [https://www.ecmwf.int/sites/default/files/elibrary/2021/20142-evaluation-ecmwf-forecasts-including-](https://www.ecmwf.int/sites/default/files/elibrary/2021/20142-evaluation-ecmwf-forecasts-including-2021-upgrade.pdf)  
 142 [2021-upgrade.pdf](https://www.ecmwf.int/sites/default/files/elibrary/2021/20142-evaluation-ecmwf-forecasts-including-2021-upgrade.pdf)

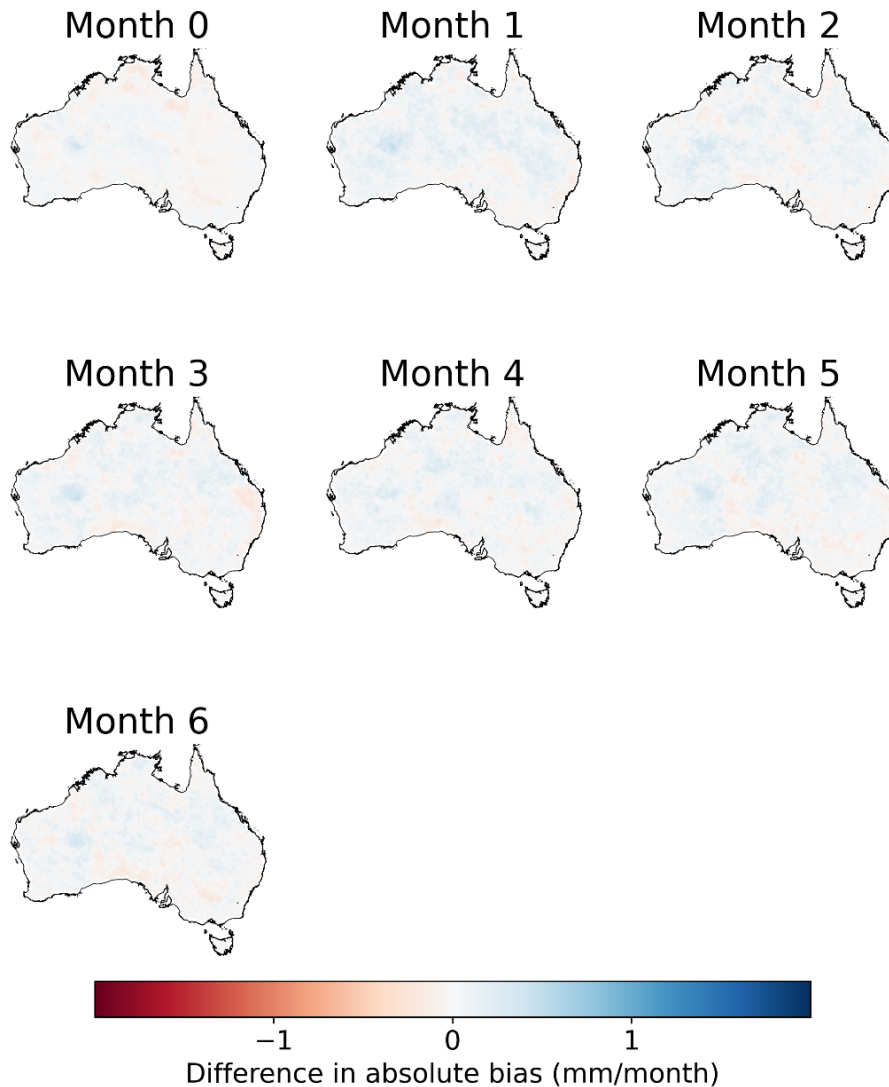
143

144 **Point #6**

145 [Line 125] “trends in transformed forecasts and observations are removed to produce detrended data”:  
 146 This is quite an aggressive process because removing trend linearly in transform space, as described in  
 147 equations 3 and 4, can lead to substantial reduction in un-transformed space after a certain time. When

148 *trends parameters in BJP-Tri are significant (which seems frequent as suggested by Figure 1), we are a bit*  
149 *concerned that this could lead to forecasts becoming unrealistically large or systematically zero if left*  
150 *unchecked.*

151 **Response: We appreciate the reviewer’s valuable comments. We further evaluated our**  
152 **methodology results and confirmed that parameter inference in the transformed space did**  
153 **not result in extreme values in calibrated forecasts. First of all, the removed trend will be**  
154 **added back to transformed forecasts/observation through the retrending process (step 5 in**  
155 **section 2.3). As a result, even a large trend is removed from transformed data in the**  
156 **detrending process, it will be added back to the transformed data before calibrated forecasts**  
157 **are transformed back to their original space. Second, as we introduced in section**  
158 **2.3(equations 7 and 8), we’ve set limits to inferred trends to avoid extreme values. Third, we**  
159 **further compared the absolute errors in calibrated forecasts produced using the BJP-ti model**  
160 **vs. those using the BJP model (See the following figure), and did not see significant increases**  
161 **in errors after trend reconstruction:**



163 Figure S2. Differences in absolute bias between BJP-ti and BJP calibrated forecasts

164 **The above figure indicates that differences in the two sets of calibrated forecasts (with vs.**  
 165 **without trend reconstruction) are almost negligible. We added the above figure to the**  
 166 **Supplementary Material, and explained findings in the comparison in section 2.3:**

167 “Our analysis indicated that our trend-reconstruction strategy (detrending and retrending in the  
 168 transformed space, and setting limits to inferred trends) would not introduce unwanted bias in the  
 169 calibrated forecasts (Figure S2).”

170 **As a result, we can reassure the reviewer that our trend reconstruction strategy will not lead**  
 171 **to extreme values in calibrated forecasts.**

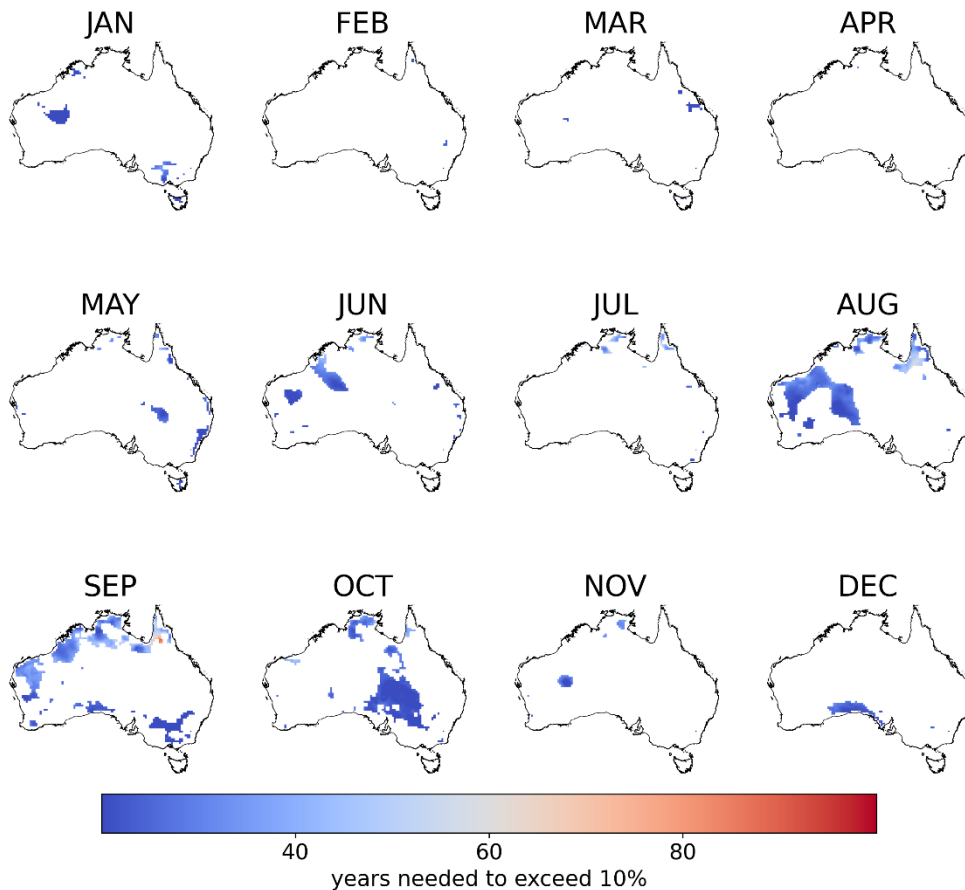
172

173 Point #7

174 *We suggest commenting briefly on the time needed for the mean unconditional forecast (i.e. considering*  
175 *zo only in Equation 5) to depart from the unconditional forecast mean obtained at t=tm by more than,*  
176 *say, 50% in untransformed space. Perhaps consider showing the distribution of this time across the*  
177 *gridded domain and provide guidance on how frequently BJP-tri should be reviewed to monitor the*  
178 *accuracy.*

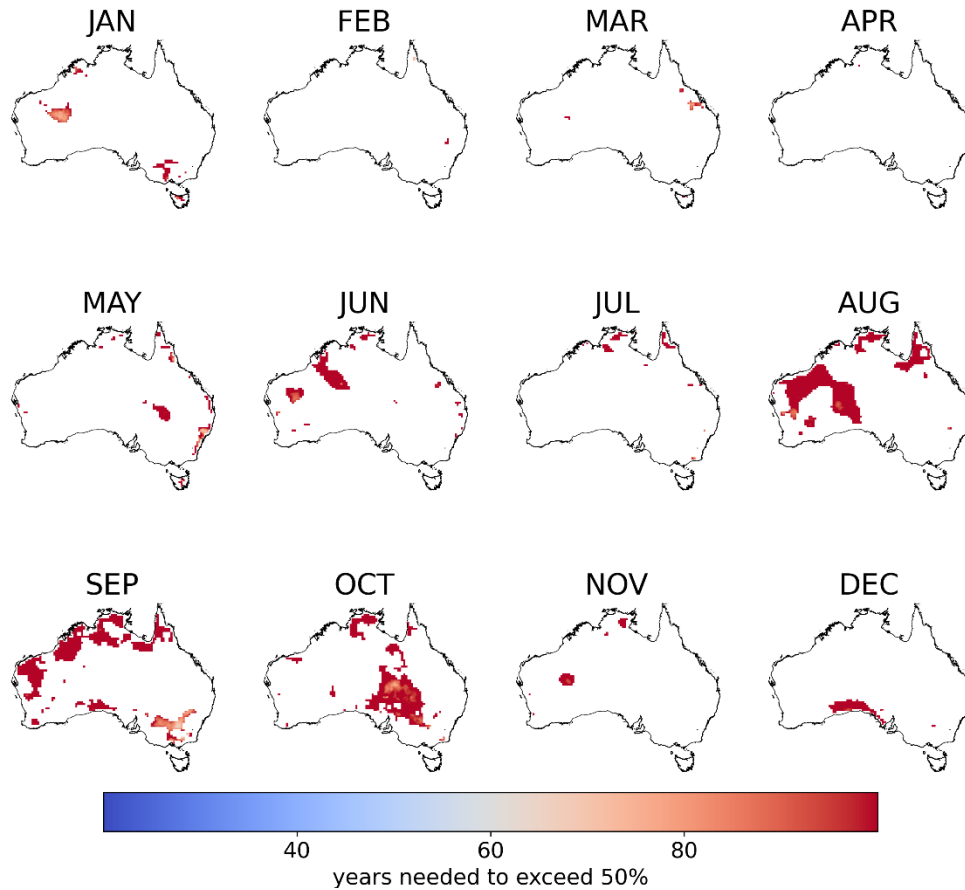
179 **Response: Thank you for the comments. We create figures to show the time needed for the**  
180 **departure of climatology forecasts which does not consider temporal trends from the**  
181 **calibrated forecasts with reconstructed trends. Here we considered both 10% and 50%**  
182 **departure. As we explained in our response to your comment #3, we adopted a new strategy**  
183 **that only allows trend reconstruction in regions with significant observed trends. As a result,**  
184 **we only focus on these regions when investigating the departures.**

185



187 Figure S11. Years needed for the departures of climatology forecasts from the calibrated  
188 forecasts with reconstructed trends to exceed 10%





190 Figure S12. Years needed for the departures of climatology forecasts from the calibrated  
 191 forecasts with reconstructed trends to exceed 50%

192

193 **As suggested by the above plots, it will take about 20-30 years for the departure to reach**  
 194 **10%, and more than 100 years to reach 50%. However, we believe correcting time-dependent**  
 195 **errors is still necessary, since increasing extreme weather conditions across the globe in**  
 196 **recent years indicate that climate change is intensifying. We add the following discussions to**  
 197 **section 4.1:**

198 “Although it may take decades for climate change to substantially alter the magnitude of  $ET_o$   
 199 (Figures S11 and 12), we recommend that future GCM-based  $ET_o$  forecasting should still correct  
 200 time-dependent errors. More frequent extreme weather events in recent years support model  
 201 projections that climate change will intensify in the future (Kharin et al., 2013). It is expected  
 202 that future climate change may induce more significant temporal trends in  $ET_o$ .”

203

204 Point #8

205 [Line 132] “ $t_m$  is approximately the middle year”: does moving  $t_m$  has an impact on generated  
206 forecasts? I believe not because it is compensated by the value of the mean parameter  $\mu$ . Please  
207 confirm. If this the case, please highlight that the position of  $t_m$  is arbitrary and does not affect the  
208 forecasts.

209 **Response: Thank you for the valuable comments. The reviewer is right that using different**  
210 **years as the reference for trend removal will impact the magnitude of the resultant**  
211 **detrended data (both forecasts and observations), but will not affect the trend**  
212 **reconstruction. When using a different year other than 2004 as a reference year, all**  
213 **detrended data point will be larger (or smaller) by the same value than data using the**  
214 **middle year as the reference. These differences will be lead to different mean and standard**  
215 **deviation parameters. However, after we add the trend back (retrending) to data, the**  
216 **difference will be canceled out. As a result, choosing a different reference year will not affect**  
217 **the trend reconstruction and forecast calibration.**

218 **We clarify this point by adding the following explanations:**

219 “The position of  $t_m$  is empirically selected, but it will not affect the calibration if we choose a different  
220 year as  $t_m$ ”

221

222 Point #9

223 “Equation 8 shows the conditional posterior distribution of parameter  $\delta$  conditional on  $\delta$ .”: We suggest  
224 “Equation 8 shows the posterior distribution of parameter  $\delta$  conditional on  $\delta$ ”.

225 **Response: We changed the wording accordingly.**

226

227 Point #10

228 “In equation 8,  $\delta$  is the mean and  $\delta$  is the standard deviation for predictors or  
229 predictands.”: Please move this sentence just after Equation 8. In addition, we suggest the following  
230 clarification: “ $\delta$  is the standard deviation for predictors or predictands extracted from the  
231 diagonal of covariance matrix  $S$  (see equation 5)”.

232 **Response: We moved this sentence to the beginning of this paragraph to better introduce**  
233 **Equation 8. We also improved the descriptions of parameters based on your suggestions.**

234

235 Point #11

236 [Line 160] “we adopt a leave-one-year-out cross-validation strategy”: for a trend-aware model, this is an  
237 optimistic approach to model validation because the model has seen both past and future data during  
238 calibration. A more challenging validation would be to split the data in two parts, infer the trend from

239 *one part and validate on the other. We understand that this is challenging with a heavily parameterised*  
240 *model such a BJP, consequently it is probably beyond the scope of this paper to solve this question here.*  
241 *However, it is important to flag the potential issue of using traditional leave-out validation for trend*  
242 *analysis.*

243 **Response: We agree with the reviewer about the issue in the leave-one-out cross-validation.**  
244 **Please see our response to the same point in your comment #2.**

245

## 246 Point #12

247 *[Line 166] “The comparison is conducted for months with large areas of statistically significant (at the*  
248 *95% confidence interval) temporal trends in observed ETo.”: this approach is problematic because it does*  
249 *not check the performance of the BJP-ti model when there is no observed trend. BJP-ti is more*  
250 *parameterised than BJP, consequently it is always exposed to the risk of overfitting the data when there*  
251 *is no trend, i.e. when trend parameters cannot be calibrated reliably. Please comment on this point and*  
252 *justify why performance assessment excluded month with no significant observed trend.*

253 **Response: Thank you for the valuable comments. As we explained in our response to your**  
254 **comment #3, we adopted a new strategy to deal with the overfitting problem. In the latest**  
255 **calibration with this strategy, the degradations in CRPS skill score and correlation coefficients**  
256 **caused by trend overfitting have been effectively corrected.**

257 **We add the evaluation results for the remaining 9 months to the supplementary material. As**  
258 **we can see in the following figures, improvements in the two metrics mainly occurred to**  
259 **regions with significant observed trends. For regions with insignificant observed trends,**  
260 **changes in the metrics are generally negligible. We introduced how results are presented in**  
261 **section 2.4 as follows:**

262

263 *“We present results of the comparison in the main text for months (August, September, and October) with*  
264 *large areas of statistically significant (at the 95% confidence interval) temporal trends in observed ETo;*  
265 *results for the remaining nine months are presented in the Supplementary Material.”*

266

267

268

269

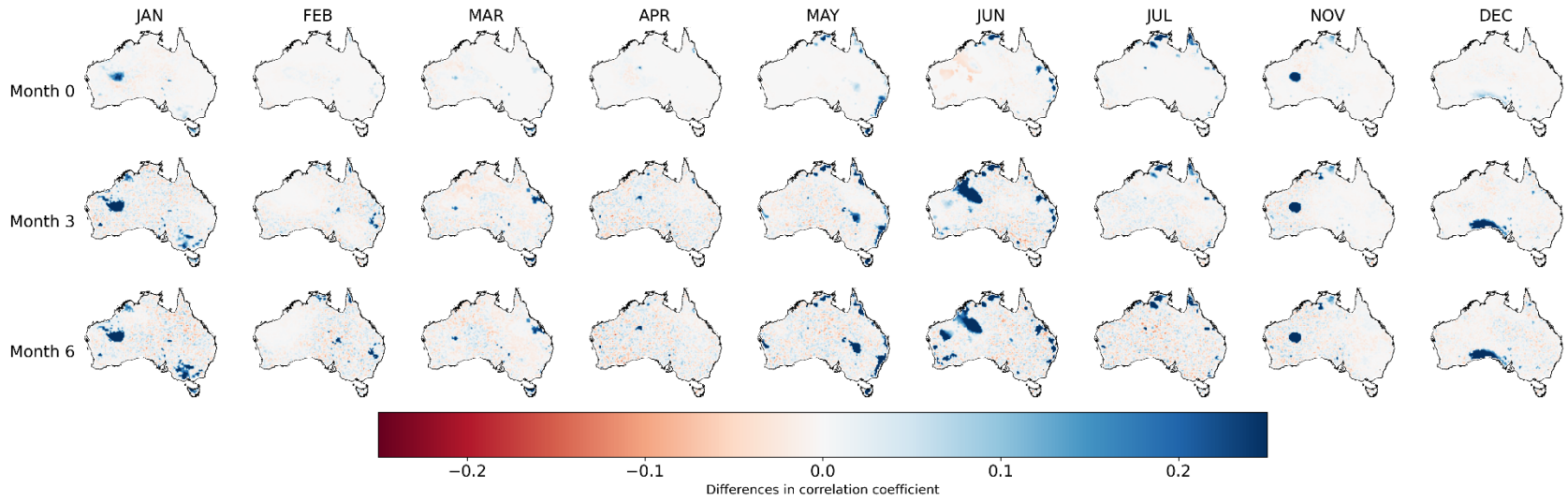
270

271

272

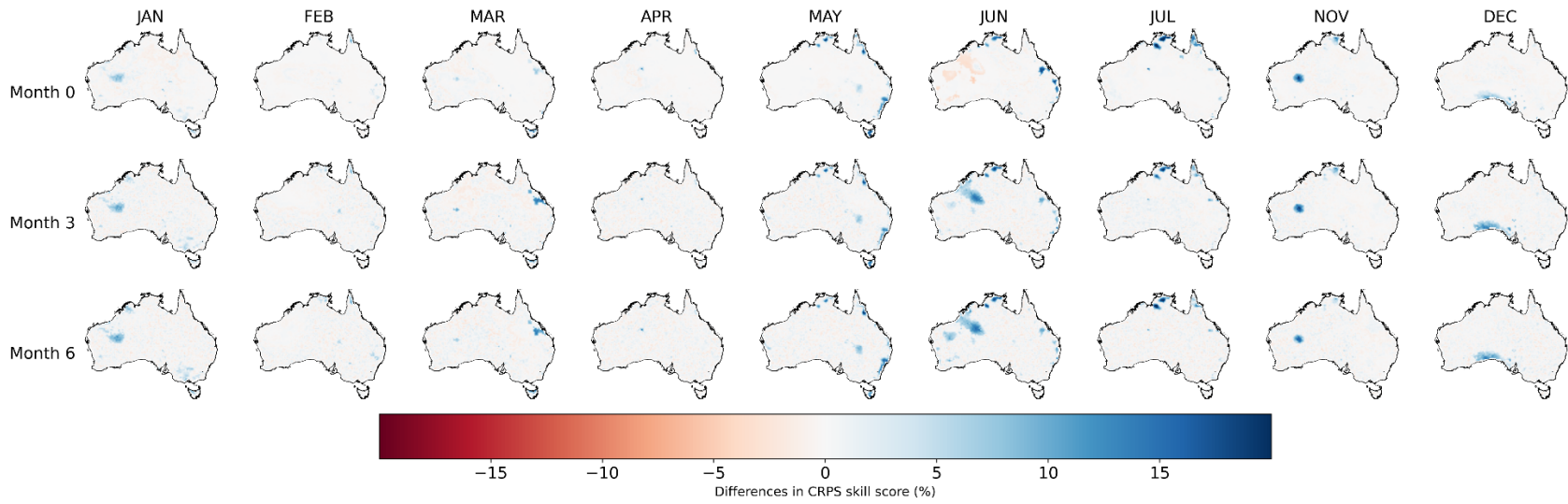
273

274



275

278



279

282

283

284

285

286

287

288

289 Point #13

290 [Line 197] “*Raw or calibrated forecasts of ETo (mm month<sup>-1</sup>)*”: This is a deterministic  
291 metric, so we believe that  $x(t)$  is the mean of raw or calibrated forecast. Please clarify.

292 **Response: Thank you for the suggestion. The reviewer is correct that for raw forecasts, they**  
293 **are calculated with the ensemble mean of each input variable (temperature, solar radiation,**  
294 **and vapor pressure), so they are deterministic; for calibrated forecasts, we used ensemble**  
295 **mean here to calculate the bias. We further explained the differences as follows:**

296 “Raw forecasts are deterministic since they are calculated based on the ensemble mean of each input  
297 variable. For calibrated forecasts, we use the ensemble mean to calculate bias. ”

298

299 Point #14

300 “*Observed ETo shows increasing trends in many parts of Australia in the three selected months*”: There is  
301 a significant body of literature related to trends in evapotranspiration related to climate change  
302 (McVicar et al., 2012). Please comment briefly on how this statement relates to current research in the  
303 field.

304 **Response: Thank you for the valuable suggestions. We reviewed a few classic publications on**  
305 **temporal trends of ETo based on the reviewer’s suggestions (Donohue et al., 2010; McVicar et**  
306 **al., 2012). Because these investigations focus on a period (1981-2006) earlier than our**  
307 **investigation (1990-2019), the negative trends across Australia from their research were not**  
308 **observed in our study. We add the following contents to briefly introduce analyses of**  
309 **temporal trends in ETo in Australia.**

310 “Compared with findings from previous investigations, observed trends identified in this study  
311 also demonstrate significant spatial variability and varying magnitudes in different months  
312 (Donohue et al., 2010; McVicar et al., 2012). We found more positive trends in our study period  
313 (1990-2019) than the period of 1981-2006 (Donohue et al., 2010) ”

314

315 **Reference:**

316 Donohue, R.J., McVicar, T.R. and Roderick, M.L.: Assessing the ability of potential evaporation  
317 formulations to capture the dynamics in evaporative demand within a changing climate, J.  
318 Hydrol., 386 (1–4), 186-197, doi: 10.1016/j.jhydrol.2010.03.020, 2010

319 McVicar, T.R., Roderick, M.L., Donohue, R.J., Li, L.T., Van Niel, T.G., Thomas, A., Grieser, J.,  
320 Jhajharia, D., Himri, Y., Mahowald, N.M., Mescherskaya, A.V., Kruger, A.C., Rehman, S. and  
321 Dinpashoh, Y.: Global review and synthesis of trends in observed terrestrial near-surface wind speeds:  
322 Implications for evaporation, J. Hydrol., 416–417, 182-205, doi: 10.1016/j.jhydrol.2011.10.024, 2012

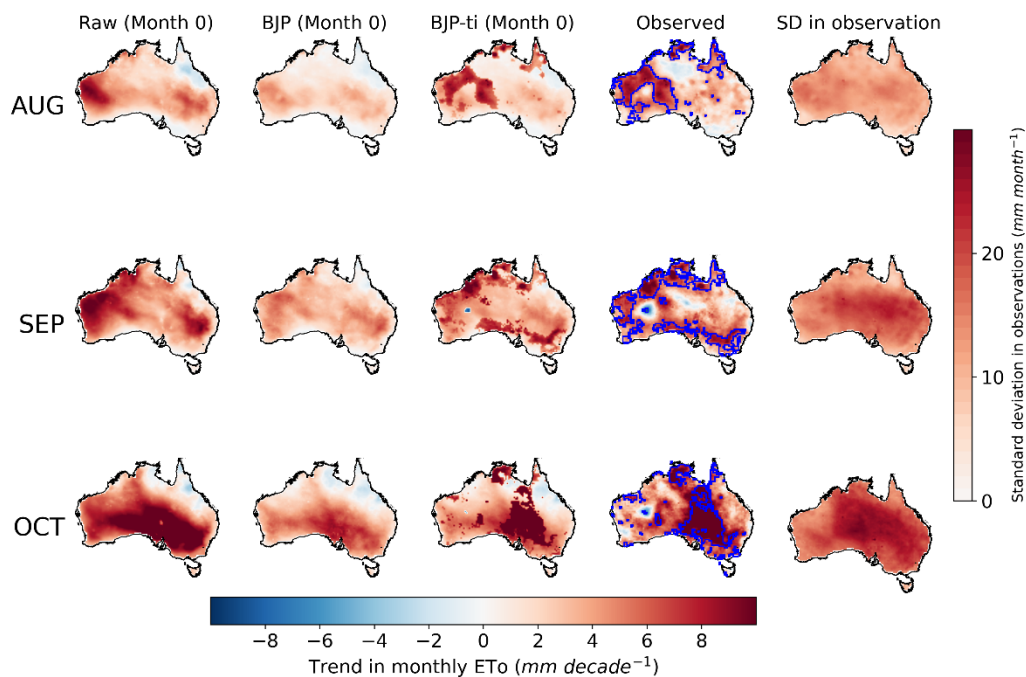
323

324 Point #15

325 [Figure 1.] We suggest adding the standard deviation of annual  $ET_o$  in the first column of figure 1 to  
326 highlight the significance of trend values. It is important to understand if the observed trends of 6 to 8  
327 mm/decade reported below are large compared to climatological variance.

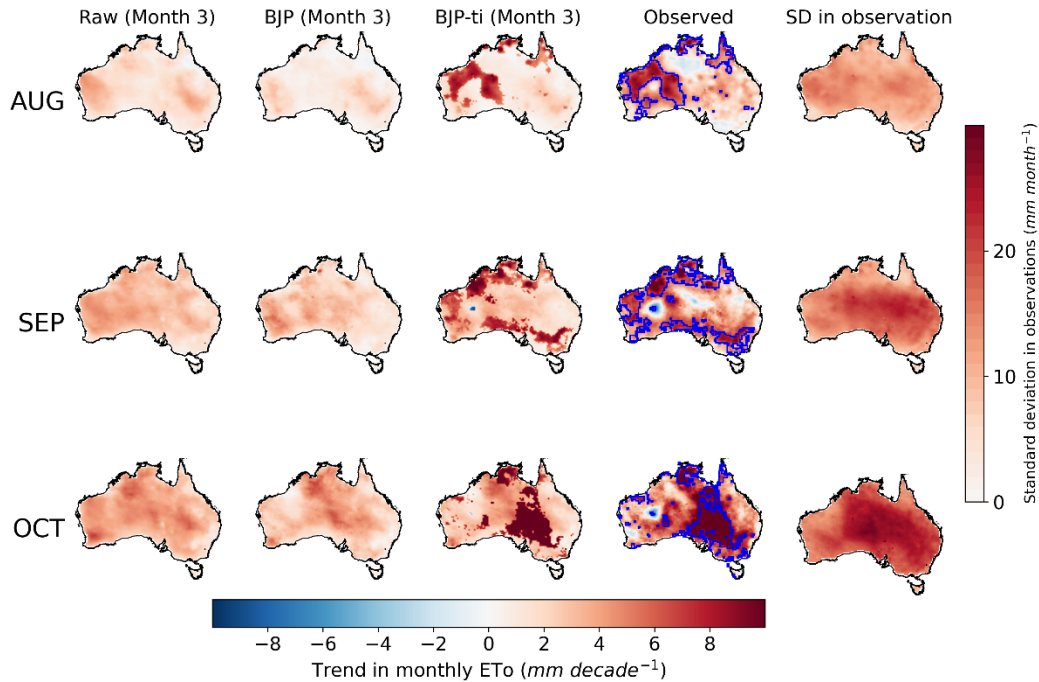
328 **Response:** Thank you for the valuable comments. We add the standard deviation to the  
329 figure. We present the standard deviation in the last column because it is easier to show the  
330 legend. In response to your comment #17, we also add contour lines to show regions with  
331 significant observed trends. Figure 1 (Month 0) and results for other lead times (Month 3 and  
332 6) in the Supplementary Material were all updated:

333



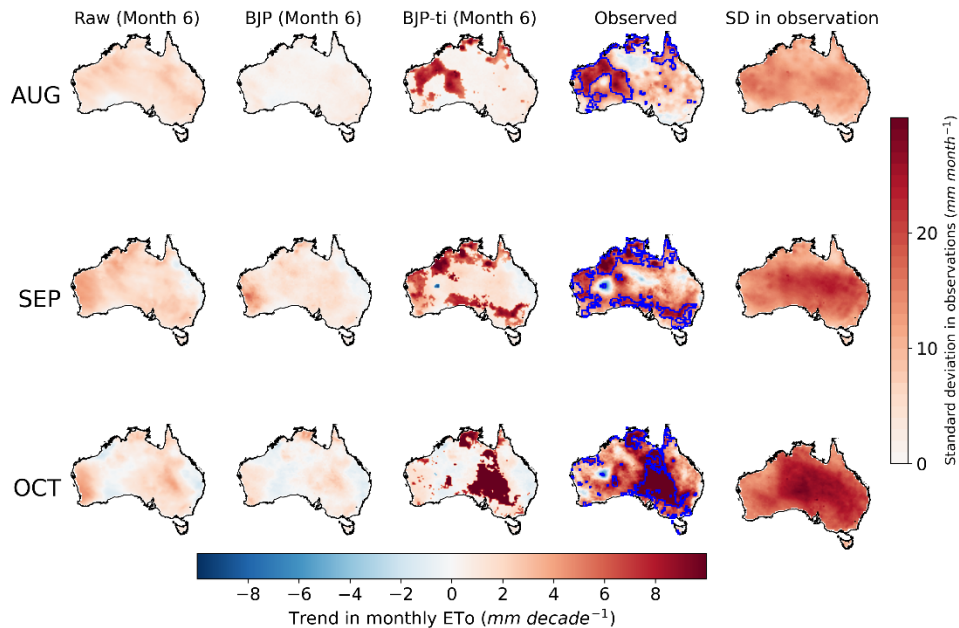
334

335 **Figure 1.** Trends in raw forecasts, BJP calibrated forecasts, and BJP-ti calibrated forecasts at the lead  
336 time of month 0, and observed  $ET_o$  in August, September, and October. Blue polygons show regions  
337 where observed trends are statistically significant. SD refers to standard deviation.



339

340 Figure S2. Trends in raw forecasts, BJP calibrated forecasts, BJP-ti calibrated forecasts for  
 341 Month 3, and observed  $ET_0$  for three selected months. Blue polygons show regions where  
 342 observed trends are statistically significant. SD refers to standard deviation.



344

345 Figure S3. Trends in raw forecasts, BJP calibrated forecasts, BJP-ti calibrated forecasts  
 346 for Month 6, and observed  $ET_0$  for three selected months. Blue polygons show regions  
 347 where observed trends are statistically significant. SD refers to standard deviation.



349 Point #16

350 *“Slight decreases in r are also found in regions where the observed trends are not statistically*  
351 *significant.”: This statement seems to support the comment made against line 166 suggesting that BJP-ti*  
352 *might suffer from over-parameterisation when observed trends are not significant. If confirmed, this is*  
353 *an important limitation of the model that should be highlighted more clearly.*

354 **Response: We agree with the reviewer on the overfitting issue. We have explained how we**  
355 **address this challenge in our response to your comment #3. Specifically, we have set fitted**  
356 **trends for regions where observed trends are statistically insignificant to zero. This new**  
357 **strategy successfully resolved the overfitting problem, and degradation in performance of**  
358 **calibration following trend reconstruction (BJP-ti vs. BJP) was also corrected. We have**  
359 **updated the manuscript based on the new calibration.**

360

361 Point #17

362 *[Figure 2.] We suggest adding in this figure a contour line showing the area where observed trend is not*  
363 *significant. This could help understand better the strength and weaknesses of BJP-ti.*

364 **Response: Thank you for the valuable suggestion. After we adopted a new calibration**  
365 **strategy, as we explained in our response to your comments #3 and #16, degradation in the**  
366 **performance of the calibration was removed. We use contour lines to show the boundaries of**  
367 **regions with significant observed trends in Figures 1, 2, and 3.**

368 **Please see details in our response to your comments #3 and #15.**

369

370 Point #18

371 *Please also report the proportion of the study area where CRPS of BJP-ti is greater than the one of BJP.*  
372 *From Figure 3, it seems that BJP-ti underperforms in large parts of the domain, even if the decrease*  
373 *remains limited.*

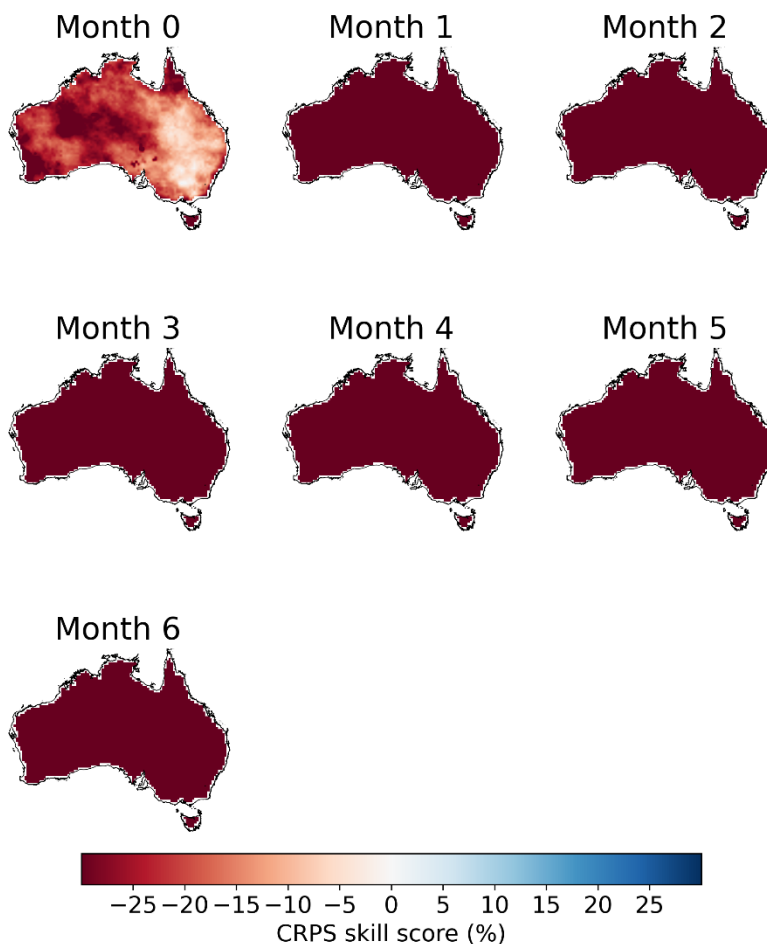
374 **Response: Thank you for the comments. After we resolve the overfitting issues, degradation**  
375 **in forecast skills is removed. Please see details in our response to your comment #3.**

376

377 Point #19

378 *“with CRPS skill scores lower than -25% in all grid cells”: this comparison is informative, but a little bit*  
379 *biased because raw operational forecasts are generally post-processed using techniques such as*  
380 *quantile-quantile mapping. We believe it is useful to show that raw forecasts have serious deficiency to*  
381 *reproduce on-ground observations, but it is also important to highlight that these forecasts would not*  
382 *normally be used for direct estimation of ET0.*

383 **Response: Thank you for the valuable suggestion. We agree with the reviewer that simple**  
384 **bias correction is often applied to raw seasonal climate forecasts. We adopted quantile**  
385 **mapping to raw ETo forecasts before the calibration with the BJP-ti model. However, we**  
386 **found that bias-corrected ETo forecasts still demonstrate low skills for lead times beyond the**  
387 **Month 0:**



388

389

Figure S13, CRPS skill score of bias-corrected ETo forecasts

395 “We need to point out that simple bias-correction is often applied to raw ECMWF forecasts  
396 before they are used. We applied quantile mapping to the raw ETo forecasts and were able to  
397 improve skills in ETo forecasts (Figure S13). However, the bias-corrected forecasts still  
398 demonstrate skills much worse than climatology forecasts, particularly at long lead times.”

399

400 In addition, since the primary objective of this investigation is to understand how trend  
401 reconstruction would affect forecast calibration, we decided to use the raw ET<sub>o</sub> forecasts for  
402 this current investigation because we are not clear how would the quantile mapping affect  
403 trends in ECMWF forecasts.

404 However, we totally agree with the reviewer that improving the raw forecasts of ECMWDF  
405 forecasts will be a very interesting point which needs further investigation. Trends in  
406 individual input variables (e.g., temperature, vapor pressure, and solar radiation) needed for  
407 ET<sub>o</sub> calculation have been reported by Donohue et al. (2010) and McVicar et al. (2012). It is  
408 not clear whether correcting bias and reconstructing trends in each of the input variables  
409 first, prior to calculating the raw ET<sub>o</sub> forecasts, will further enhance the ET<sub>o</sub> forecasts  
410 calibration. We highlight this point in our Future work section (4.2):

411 “In this study, we directly use the raw forecasts of individual input variables (e.g., temperature,  
412 solar radiation, and vapor pressure) to construct the raw ET<sub>o</sub> forecasts. However, trends in these  
413 variables have been reported in previous investigations. Whether correcting errors including  
414 time-dependent errors in the raw forecasts of each input variable, will lead to more skillful  
415 calibrated ET<sub>o</sub> forecasts, warrants further investigation in the future”

416

#### 417 **Reference:**

418 Donohue, R.J., McVicar, T.R. and Roderick, M.L.: Assessing the ability of potential evaporation  
419 formulations to capture the dynamics in evaporative demand within a changing climate, J.  
420 Hydrol., 386 (1–4), 186-197, doi: 10.1016/j.jhydrol.2010.03.020, 2010

421 McVicar, T.R., Roderick, M.L., Donohue, R.J., Li, L.T., Van Niel, T.G., Thomas, A., Grieser, J.,  
422 Jhajharia, D., Himri, Y., Mahowald, N.M., Mescherskaya, A.V., Kruger, A.C., Rehman, S. and  
423 Dinpashoh, Y.: Global review and synthesis of trends in observed terrestrial near-surface wind speeds:  
424 Implications for evaporation, J. Hydrol., 416–417, 182-205, doi: 10.1016/j.jhydrol.2011.10.024, 2012

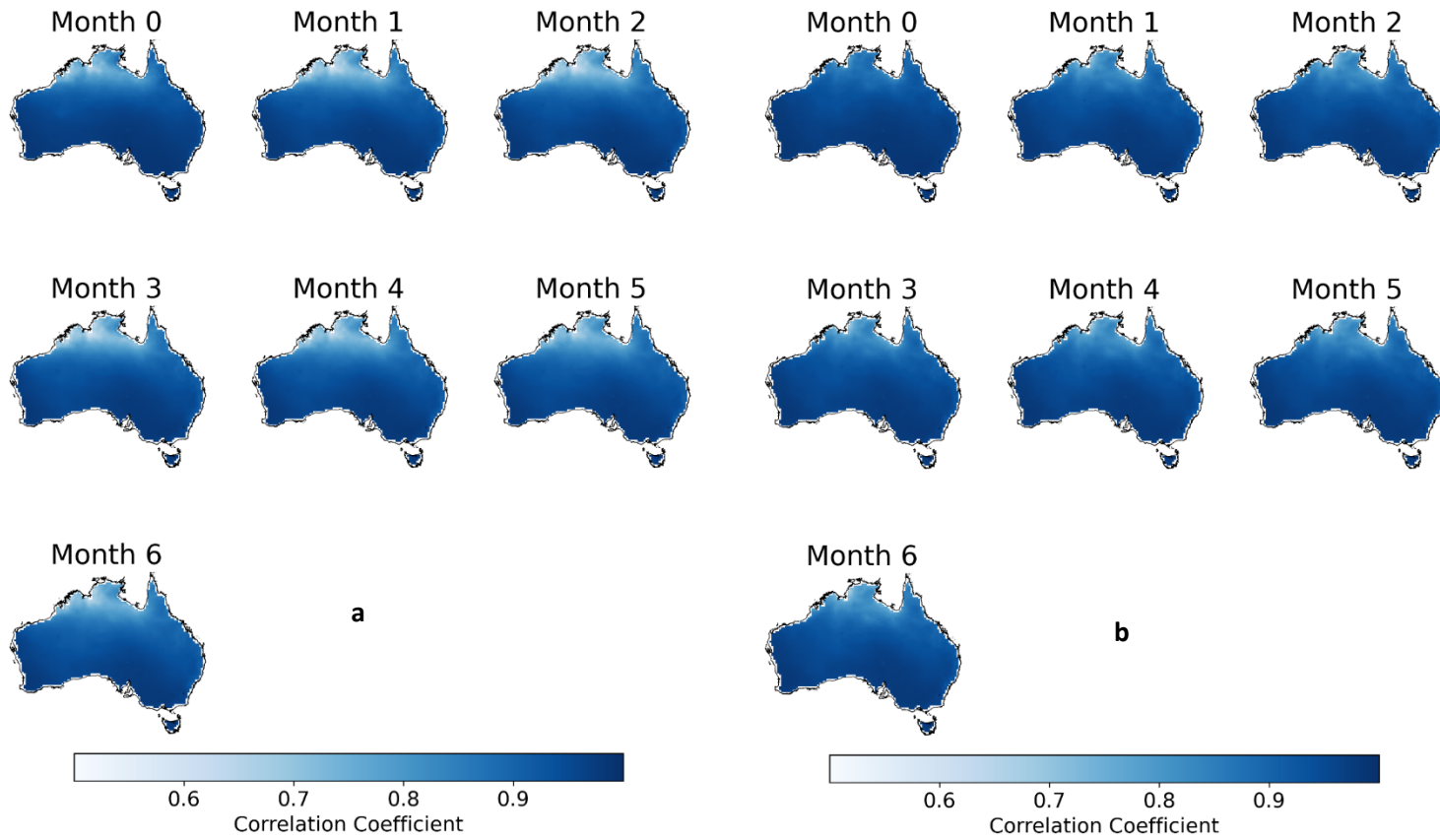
425

#### 426 Point #20

427 *It would be perhaps more interesting to compare the correlation score between raw and BJP-ti forecasts,*  
428 *which discards some the known deficiencies of raw forecasts.*

429 **Response: Thank you for the valuable suggestions. We agree with the reviewer that the**  
430 **correlation coefficient could be less impacted by the systematic errors in raw ECMWF**  
431 **forecasts than other metrics. We calculated the correlation coefficients between raw/BJP-ti**  
432 **calibrated forecasts and observations. Because of the high seasonality in ET<sub>o</sub>, both raw and**  
433 **calibrated forecasts demonstrate high correlations with observations:**

434



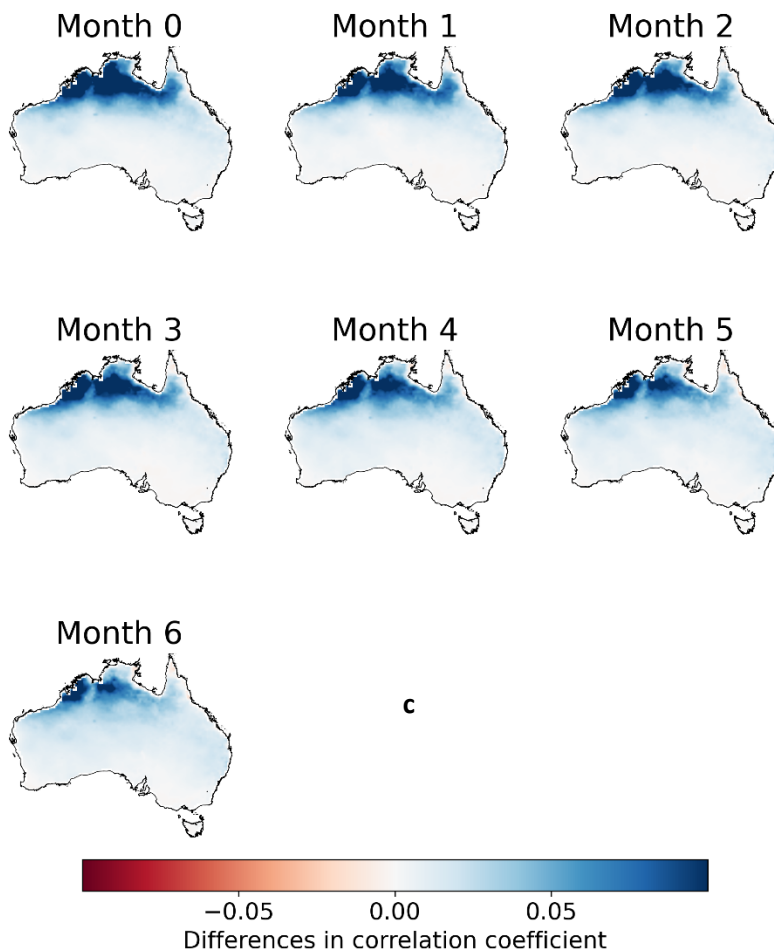
435

436

**Correlation coefficients between (a) raw forecasts/(b) calibrated forecasts and observations.**

437 To demonstrate the improvements in correlation through the calibration with the BJP-ti  
438 model, we compared the correlation coefficients between calibrated forecasts and  
439 observation with those between raw forecasts and observation:

440



441

442 **(c) improvements in correlation coefficient through the calibration with the BJP-ti model**

443 Results show improvements in correlation coefficients for all lead times, particularly in  
444 northern Australia, where raw forecasts demonstrate low correlations with observations.

445 Since the correlation plots for (a) raw and (b) calibrated forecasts are very similar, we decided  
446 to keep (b) and (c) in the main text (Figure 8 in the revised manuscript) and present (a) in the  
447 Supplementary Material (Figure S10).

448 We add the new section in the main text to demonstrate the evaluation of the performance  
449 of calibration in improving correlation coefficients:

450 **“3.5 Correlation between raw/calibrated forecasts and observations**

451 The calibration based on the BJP-ti model also improves the correlation coefficients between forecasts  
452 and observations. Raw forecasts are able to capture the high seasonality in  $ET_o$  and thus demonstrate high  
453 correlation coefficients with observations (Figure S10). The  $r$  values are generally over 0.9 across most  
454 parts of central and southern Australia. Lower  $r$  values are mainly distributed in coastal regions of  
455 northern Australia. Calibration with the BJP-ti model further improved the representation of  $ET_o$  temporal  
456 dynamics (Figure 8). The  $r$  values for calibrated forecasts are over 0.9 in most parts of Australia.  
457 Improvements in  $r$  are more pronounced in northern Australia, where raw forecasts show lower  
458 correlations with observations. ”

459

#### 460 Point #21

461 *Same comment than for Line 290.*

462 **Response: We understand the reviewer’s concern about how we evaluate the raw forecasts.**  
463 **As we explained in our response to your comments #19 and #20, we further 1) apply bias-**  
464 **correction to raw forecasts, 2) highlight the necessity of improving individual input variables**  
465 **prior to the calculation of raw  $ET_o$  forecasts, and 3) use the correlation coefficients as another**  
466 **evaluation metrics to show the performance of raw forecasts. Please see details in our**  
467 **response to your comments #19 and #20.**

468

#### 469 Point #22

470 *“We recommend that future GCM-based  $ET_o$  forecasting should correct time-dependent errors”: this*  
471 *comment should be toned down to include the risk of model overfitting discussed previously in relation to*  
472 *lines 166 and 271.*

473 **Response: Thank you for the comments. First, as we explained in our response to your**  
474 **comment #3, the overfitting problem has been resolved by setting the trend to zero in**  
475 **calibration for grid cells where observations do not demonstrate statistically significant**  
476 **trends. Second, we agree with the reviewer that it is necessary to remind the audience of the**  
477 **importance of avoiding overfitting in forecast trend reconstruction.**

478 **We feel it is better to highlight the necessity of dealing with overfitting in the discussion of**  
479 **BJP-ti model’s strengths. As a result, we add the following discussions to the second**  
480 **paragraph of section 4.2 (Implications for improving statistical calibration models):**

481 *“The successful application to  $ET_o$  forecasts confirms the robustness of trend reconstruction algorithms*  
482 *based on the data transformation, Bayesian inference, and using statistical significance of observed trends*  
483 *to deal with overfitting of trend parameters in the BJP-ti model. This study further demonstrates the*  
484 *feasibility for the general application of BJP-ti to different hydroclimate variables showing temporal*  
485 *trends. We also anticipate that the BJP-ti algorithms for trend reconstruction could be adopted by other*  
486 *calibration models to enhance seasonal forecast calibration.”*

487 Point #23

488 *“Future work for seasonal ETo forecasting”:* We suggest adding the two challenges of model overfitting  
489 *when there is no observed trend and validation of trend-aware forecast beyond leave-one-out approach.*

490 **Response:** Since the overfitting issue has been resolved (response to comment #3), and we  
491 already highlighted the importance of dealing with this issue in section 4.2 (response to  
492 comment #22), we decided to emphasize the challenge in cross-validation only here. Our  
493 discussion on the limitations of the leave-one-month out strategy and future work needed to  
494 address this challenge are presented in our response to your comment #2.

495

496

497

498