

Responses to Reviewer #1

Point #1

In the manuscript “Reconstructing climate trends adds skills to seasonal reference crop evapotranspiration forecasting”, Yang et al adopted a new method to improve the prediction of evaporative water loss based on seasonal climate forecasts from the ECMWF model. This method is capable of dealing with the impacts of the changing climate on the prediction of future evapotranspiration (Reference crop evapotranspiration, ETo), and could lead to more realistic predictions. The changing climate has substantially altered the water cycle, representing one of the most critical challenges in hydrological modelling and water resource management. This work is innovative in taking this impact into account and addressing the challenges associated with climate change in the prediction of future evapotranspiration. The developed method is expected to be applicable to other models and thus benefit both forecasters (weather/climate centers) and forecast users (irrigators, hydrological modelers).

The manuscript is generally well written. Introduction clearly explains the background, challenges, motivation, and objective of this work; Method provides detailed information of the model, how the model runs are conducted, and evaluation metrics; Results generally are clear and readable; Discussion provides valuable insights and important implications for future improvements of climatology-based models in hydrological modeling and forecasting.

I encourage the authors to address the following issues before publishing this work.

Response: We appreciate the reviewer’s nice summary and constructive comments.

Point #2

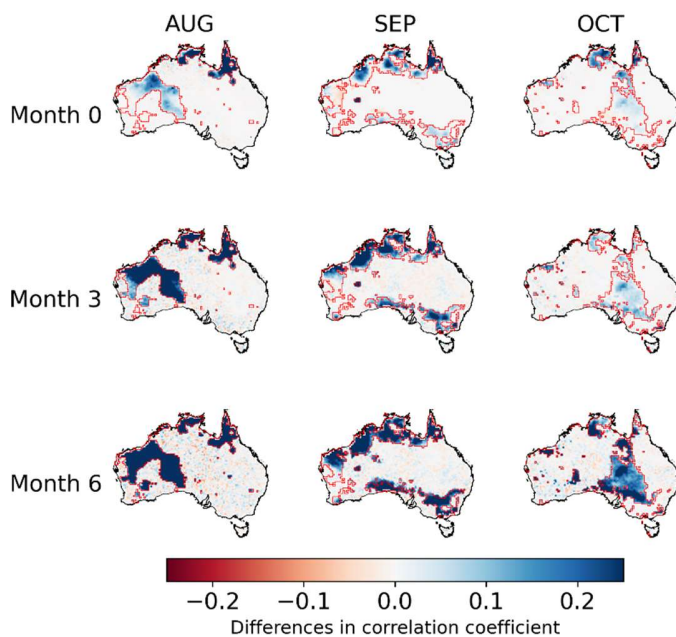
1. For time-series data, in addition to the magnitude of trend, another important feature is the statistical significance. I noticed the authors had taken this into consideration in selecting the months (8,9,10) for evaluating the performance of trend construction. In constructing the observed trends in calibrated forecasts, you empirically set limits of the trends in equation 8. I understand this is to avoid extremely large trend values. In addition to this adjustment, I think you should limit trends to zero, in grid cells where observed trends are insignificant ($P < 0.05$). Otherwise, the trend reconstruction may overestimate climate trends. I see decreases in the correlation coefficients and skill scores when compared with the calibration without trend reconstruction (Figures 2 and 3). I think limiting the insignificant trends could avoid these unwanted decreases. I suggest the authors rerun the trend-reconstruction calibration and take statistical significance into account. If you see improvements in the new runs, update the results accordingly.

Response: We agree with the reviewer that the statistical significance of trends in observations should be tested and used to limit the reconstructed trends. We accepted your valuable suggestions and redid the calibration and analysis by setting limits in trend

37 reconstruction. Specifically, we used $P < 0.05$ as the threshold to define statistically significant
38 trends. For grid cells with insignificant observed trends ($P > 0.05$), we set inferred trends to
39 zero to avoid overfitting. We introduced this new strategy in section 2.3 as follows:

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

43 **New results show that this strategy is not only effective in limiting the trend reconstruction to**
44 **regions where observed trends are significant, but also helps avoid the reductions in**
45 **correlation coefficient and CRPS skill score caused by overfitting (Figures 2 and 3):**

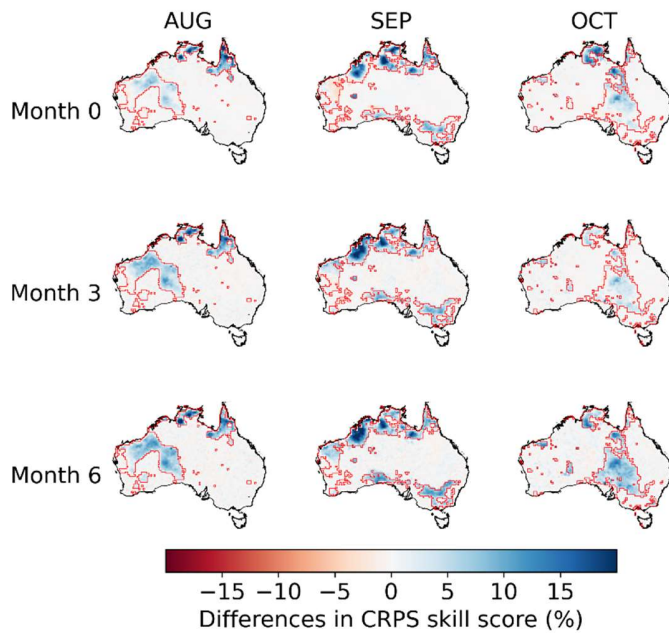


46

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

51

52



53

54

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

58

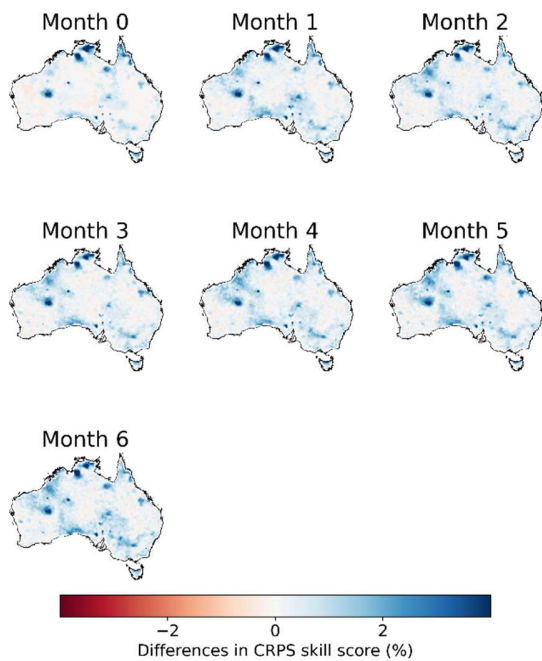
59 **We updated all results in the manuscript based on the new calibration.**

60

61 Point #3

62 *2. In addition to the improvements in the 3 selected months, whether trend construction improve the*
 63 *calibration over the whole study period?*

64 **Response: Thank you for the valuable suggestions. We added a new figure (Figure 4) to show**
 65 **the overall improvements in CRPS skill score and updated section 3.3 accordingly:**



66

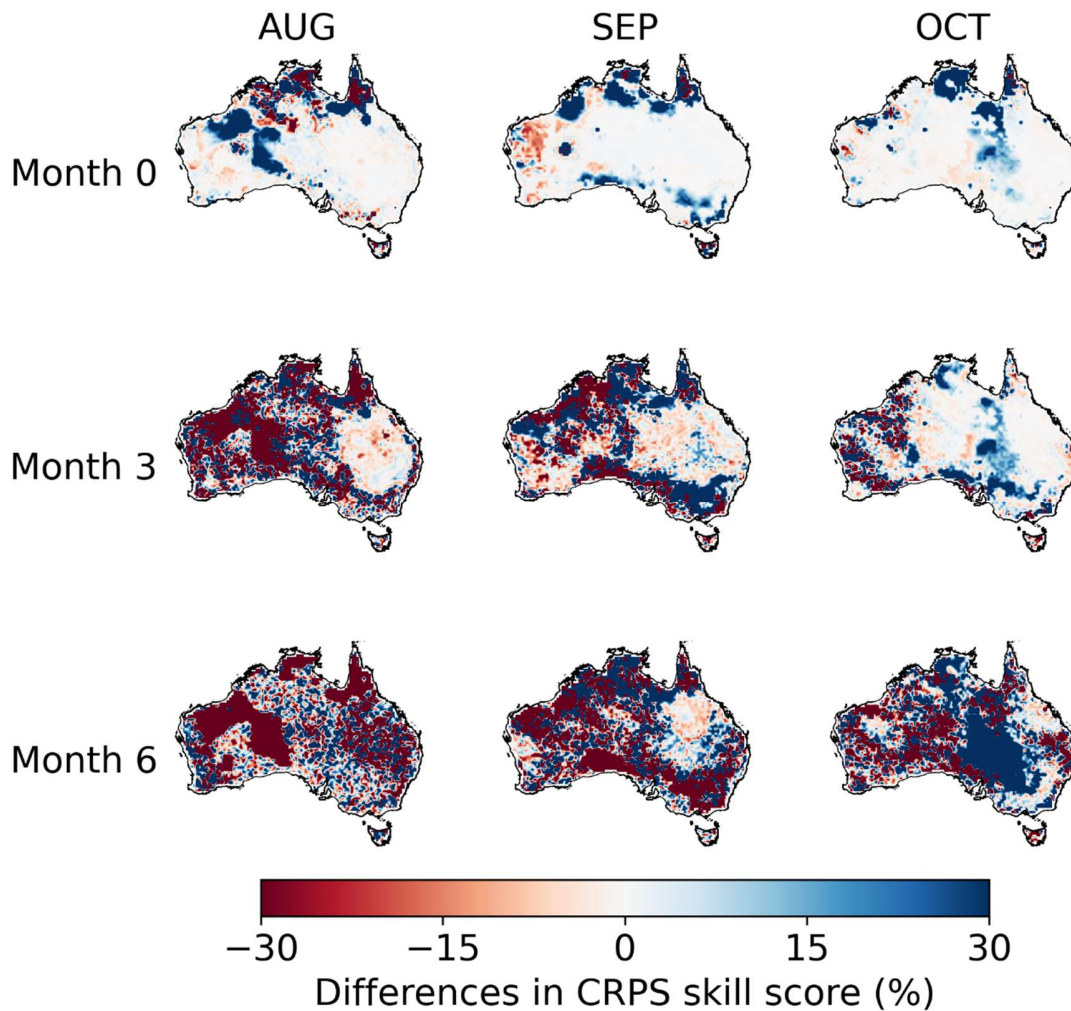
67 **Figure 4. Differences in CRPS skill score between BJP-ti calibrated forecasts and the BJP calibrated**
 68 **forecasts over 1990-2019**

69

70 Point #4

71 *3. Presentation of the improvements in figures 2 and 3. I suggest the authors use the percentage of*
 72 *changes to demonstrate the differences. Since correlation and skill score vary largely from short to long*
 73 *lead times, using percentages could better demonstrate the more significant improvements at long lead*
 74 *times.*

75 **Response: Thank you for the valuable suggestions. We did not use percentage as the unit**
 76 **because we found that at long lead times, CRPS skill score in calibrated forecasts based on**
 77 **the BJP model could be slightly negative, and thus make the plot based on percentage**
 78 **confusing:**



79

80 **As a result, we decided to use their original unit. Actually, after fixing the problems in**
 81 **overfitting, figure 2 and 3 could better demonstrate how trend reconstruction improve the**
 82 **correlation and skill scores, particularly at long lead times. Please see details in our response**
 83 **to your comment #2.**

84 Point #5

85 *Specific comments:*

86 *Page 1. line 22, forecast should be forecasting*

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

88

89 Point #6

90 *Page 3. line 92-93. This study is performed across Australia only*

91 **Response: We added the following sentence to clarify the spatial extent of this investigation:**

92 “While SEAS5 produces climate forecasts across the globe, the calibration in this study is performed
93 across Australia only.”

94

95 Point #7

96 *Page 4. line 100, Calculation of ETo observations and forecasts*

97 **Response: We changed the subtitle accordingly.**

98

99 Point #8

100 *Page 6. line 160-165. Please italicize k in this paragraph and throughout the manuscript to be consistent*
101 *with the equations.*

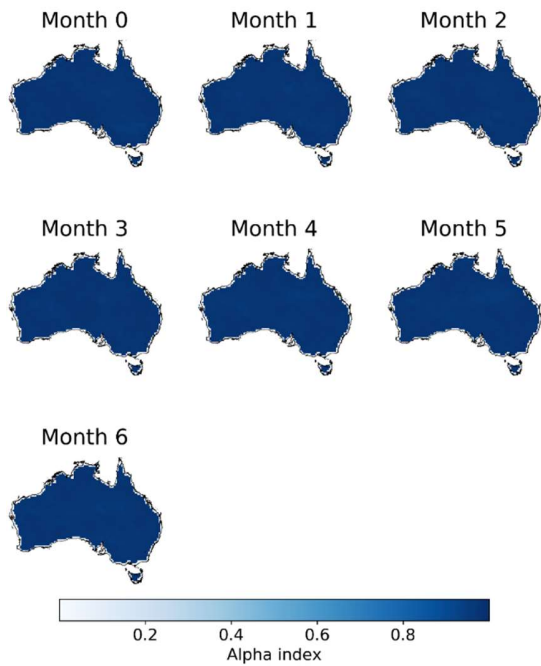
102 **Response: We italicized k in the manuscript.**

103

104 Point #9

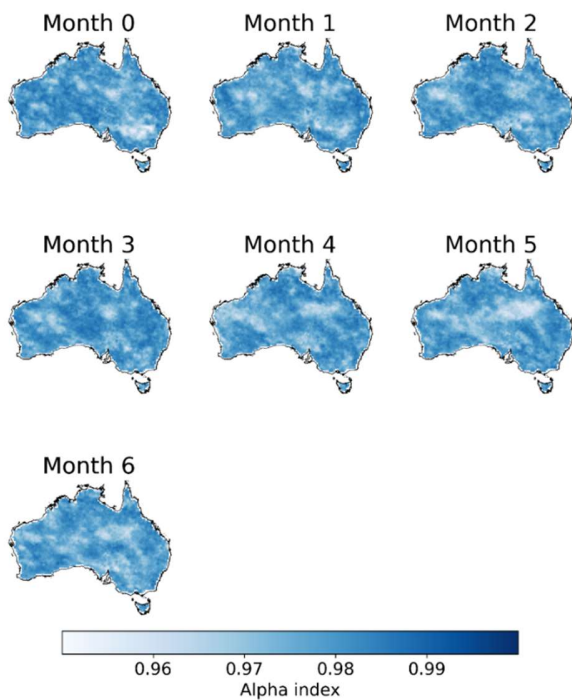
105 *Page 15. Figure 7, It is hard to read the alpha index values in the figure. Please consider changing the*
106 *limits of the color bar, and use narrower limits (e.g.,0.8-1), to make the alpha index maps more readable.*

107 **Response: We replotted the figure with a new color bar of 0.95-1 and replaced the original**
108 **figure:**



109

110 **With the following one:**



111

112 Point #10

113 *Page 17. line 378. To change with time?*

114 **Response: We changed the wording based on your suggestions.**

115

Responses to Reviewer #2

116 Point #1

117 1. General comments

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

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

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

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

133

134 Point #2

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

142 **Response: Thank you for the valuable comments and suggestions. We agree with the**
143 **reviewer that the current leave-one-out cross-validation strategy is not perfect for inferring**
144 **the trend parameters. As we introduced in the Method section (equation 5), the two trend**
145 **parameters are inferred together with parameters (mean vector and covariance matrix)**
146 **defining the bivariate distribution. The current strategy (leave-one-out) has been proven**
147 **effective for the inferencing of mean vector and covariance matrix, but may not be good**
148 **enough for the inference of trend parameters, since the left-out month may not be fully**
149 **independent of the remaining 29 months. Leaving out longer years, such as splitting the 30-**
150 **year data equally into two parts, as the reviewer suggested, could alleviate this problem to**

151 **some extent. However, we have another concern about data splitting. This strategy will**
152 **substantially reduce samples for parameter inference from 29 to 15, and thus may lead to**
153 **significant sampling errors. We feel solving this problem may need additional efforts and**
154 **more sophisticated solutions. As a result, we highlight this as a challenge that should be**
155 **addressed in our future work in section 4.3 (Future work)**

156 “First of all, more sophisticated cross-validation methods should be developed for the inference
157 of trend parameters. The current leave-one-out method has been proven to be effective in the
158 inference of the mean vector and covariance matrix (Shao et al., 2020). However, this strategy
159 may not guarantee the independence between the left-out data and data used for the inference of
160 trend parameters. We decided not to implement the data-splitting method for cross-validation
161 because of the risk of introducing sampling errors. Future investigations should take this
162 challenge into consideration and develop more robust cross-validation methods for the inference
163 of trend parameters.”

164

165 Point #3

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

172 **Response: We agree with the reviewer that the original BJP-ti parameterization suffered from**
173 **parameter overfitting and resulted in degradations in performance when compared with the**
174 **BJP model.**

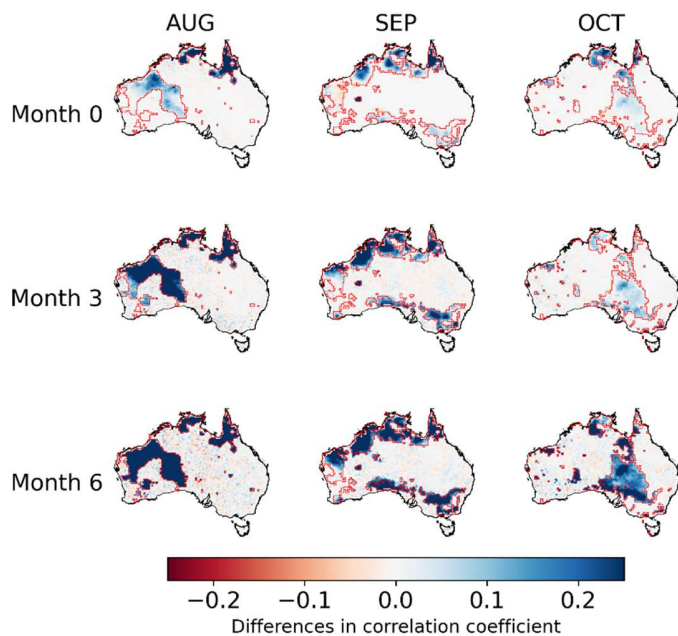
175 **To solve this problem, we accept your valuable suggestions and add limits to inferred trends**
176 **in trend reconstruction. Specifically, we use $P < 0.05$ as the threshold to define statistically**
177 **significant trends. For trends that are statistically insignificant ($P > 0.05$), we set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

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

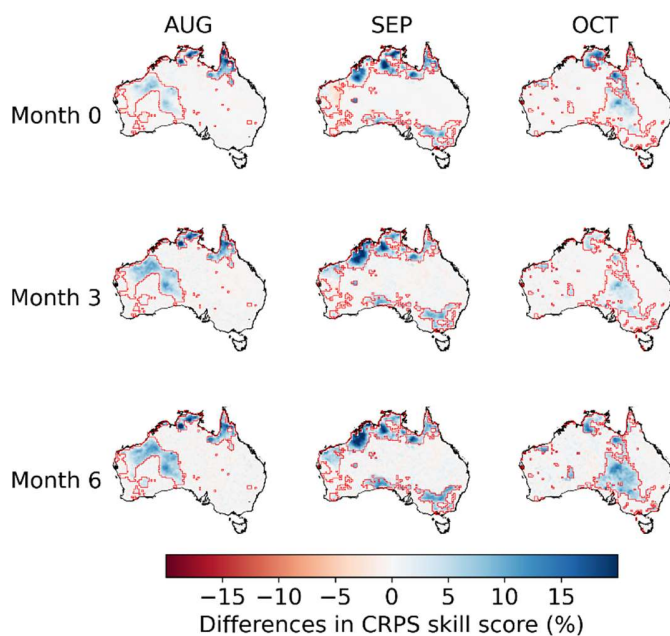
182 **This new strategy is not only effective in limiting the trend reconstruction to regions with**
183 **significant observed trends (Figure 1), but also avoids the reductions in correlation**
184 **coefficients (Figure 2) and CRPS skill score (Figure 3) following trend reconstruction.**

185 **We have updated the manuscript based on the new calibration. We present Figures 2 and 3**
186 **here to show the advantage and effectiveness of the new strategy. As you can see, the**

187 decreases in correlation and CRPS skill score were removed. For regions with statistically
188 insignificant trends, changes in the two metrics are negligible.



189
190 **Figure 2. Differences in the correlation coefficient (r) between BJP-ti calibrated forecasts and**
191 **observations with that between BJP calibrated forecasts and observations for three selected months**
192 **(AUG, SEP, OCT) and three lead times (Months 0, 3, and 6). Red polygons show regions with significant**
193 **trends.**



194
195

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

199 Point #4

200 *[Line 26] “Reference crop evapotranspiration (ET₀) measures the evaporative demand of the*
201 *atmosphere”: Please provide additional details regarding the definition of ET₀. We suggest the following:*
202 *“Reference crop evapotranspiration (ET₀) measures the evaporative demand of the atmosphere for a*
203 *hypothetical crop of given height, with defined surface resistance factor and albedo. It is generally*
204 *computed using the Penman-Monteith equation following Allen et al. (1998, see section 2.1), which is*
205 *known as FAO56. McMahon et al. (2013) provides additional information about the process.”*

206 **Response: Thank you for your valuable suggestions. We add the suggested introduction of ET₀**
207 **and the suggested reference to the manuscript.**

208 **Reference:**

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

212

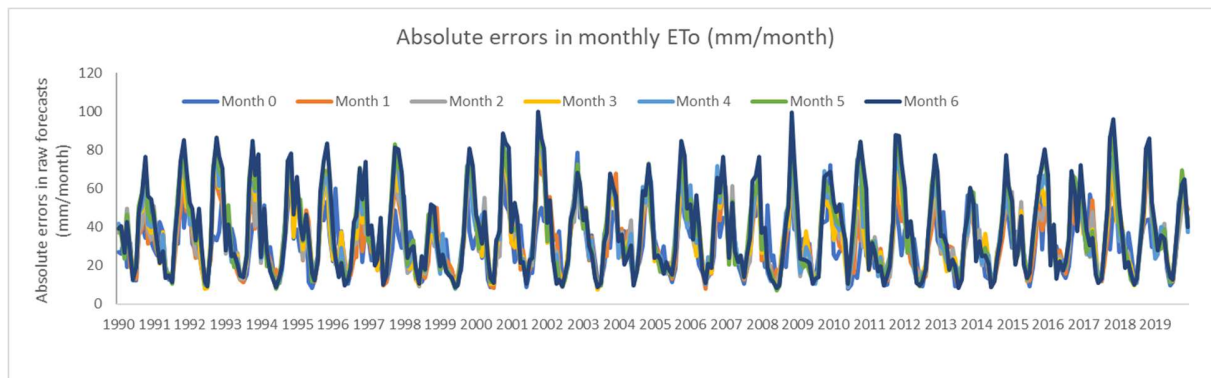
213 Point #5

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

217 **Response: Thank you for the suggestions. According to the ECMWF SEAS5 documentations**
218 **(Stockdale et al., 2017; Johnson et al., 2019), SEAS5 runs for the re-forecast and operational**
219 **forecasts periods were configured as similar as possible to maintain consistencies. However,**
220 **there are some slight differences. In addition to ensemble size, initial conditions for the two**
221 **sets of runs are from different data sources. As a result, performance during the two periods**
222 **may vary for some weather variables. For example, according to the ECMWF user guide**
223 **(ECMWF 2021), because of the different initializations, ‘the real-time forecasts of Lake**
224 **Superior (including the Great Lakes and the Caspian Sea) are cooler in the summer than the re-**
225 **forecasts were’. In addition, according to the latest evaluation of the SEAS5 forecasts (Figure**
226 **40 in Haiden et al., 2021), forecasts of accumulated cyclone energy for the Atlantic tropical**
227 **storm demonstrate larger errors during 2016-2021 than the re-forecasts.**

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

231 **In addition, we did not see significant differences in absolute errors in raw ET_0 forecasts**
232 **during the re-forecast period (1990-2016) vs. operational forecasts (2017-2019). As shown in**
233 **the following figure, the absolute errors during the re-forecasts and real-time periods seem to**
234 **be comparable. We added this figure to the Supplementary Material.**



235
236 Figure S1. Absolute errors in raw ECMWF ET_0 forecasts.

237

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

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

245

246 Reference:

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

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

252 ECMWF. SEAS5 user guide. Version 1.2, March 2021.

253 https://www.ecmwf.int/sites/default/files/medialibrary/2017-10/System5_guide.pdf

254 Haiden, T., Janousek, M., Vitart, F., Ben-Bouallegue Z., Ferranti, L. and Prates, F.: Evaluation of
255 ECMWF forecasts, including the 2021 upgrade. Technical Memo 884. 2021.

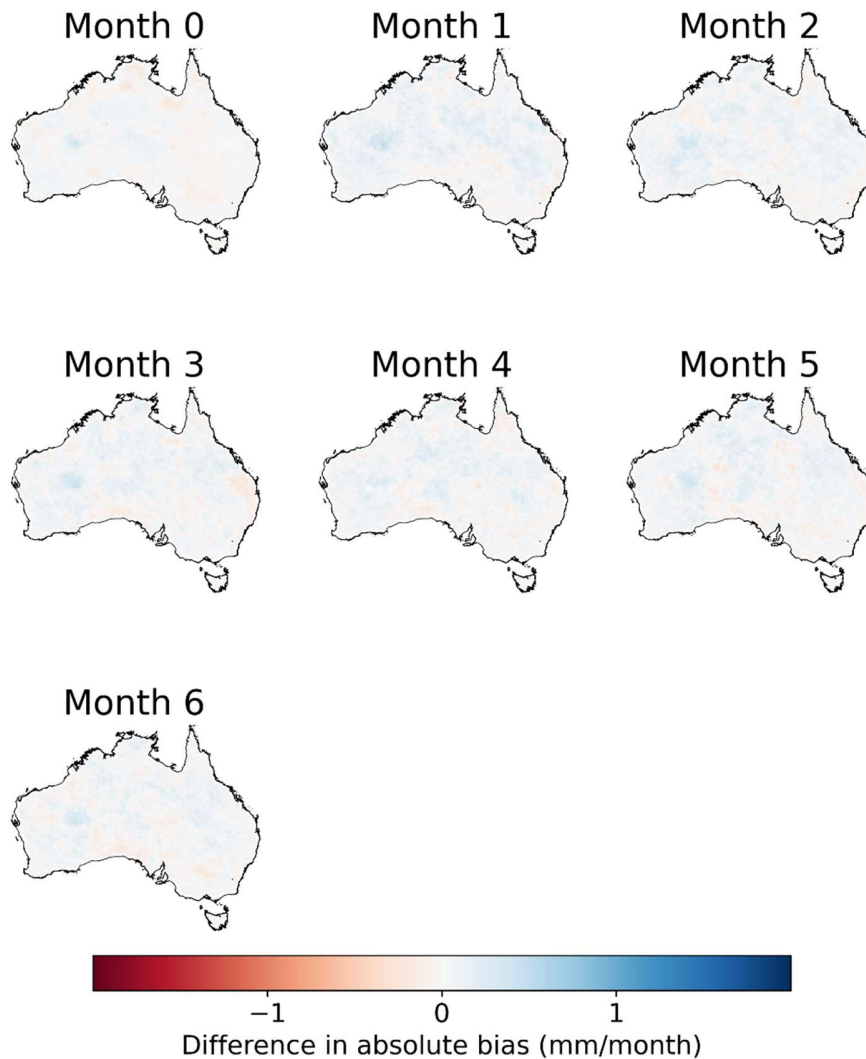
256 [https://www.ecmwf.int/sites/default/files/elibrary/2021/20142-evaluation-ecmwf-forecasts-including-](https://www.ecmwf.int/sites/default/files/elibrary/2021/20142-evaluation-ecmwf-forecasts-including-2021-upgrade.pdf)
257 [2021-upgrade.pdf](https://www.ecmwf.int/sites/default/files/elibrary/2021/20142-evaluation-ecmwf-forecasts-including-2021-upgrade.pdf)

258

259 Point #6

260 *[Line 125] “trends in transformed forecasts and observations are removed to produce detrended data”:*
261 *This is quite an aggressive process because removing trend linearly in transform space, as described in*
262 *equations 3 and 4, can lead to substantial reduction in un-transformed space after a certain time. When*
263 *trends parameters in BJP-Tri are significant (which seems frequent as suggested by Figure 1), we are a bit*
264 *concerned that this could lead to forecasts becoming unrealistically large or systematically zero if left*
265 *unchecked.*

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



277

278

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

279

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

280

281

282

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

283

284

285

As a result, we can reassure the reviewer that our trend reconstruction strategy is reliable.

286

287

Point #7

288

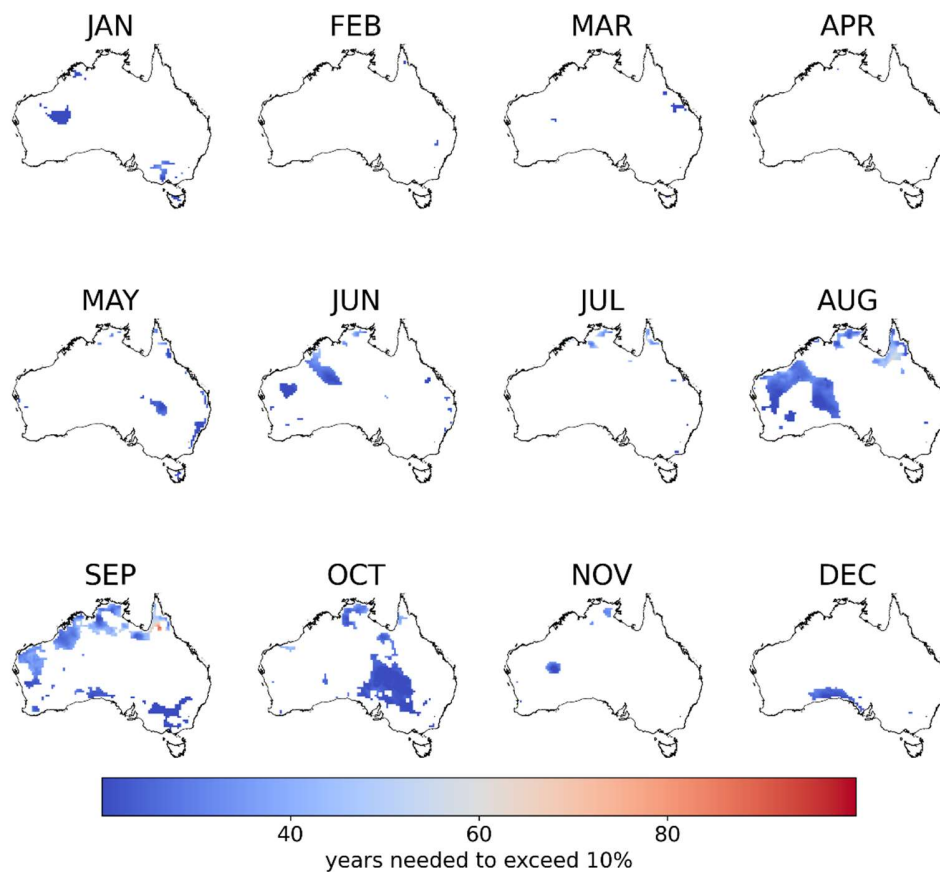
We suggest commenting briefly on the time needed for the mean unconditional forecast (i.e. considering zo only in Equation 5) to depart from the unconditional forecast mean obtained at $t=t_m$ by more than,

289

290 say, 50% in untransformed space. Perhaps consider showing the distribution of this time across the
291 gridded domain and provide guidance on how frequently BJP-tri should be reviewed to monitor the
292 accuracy.

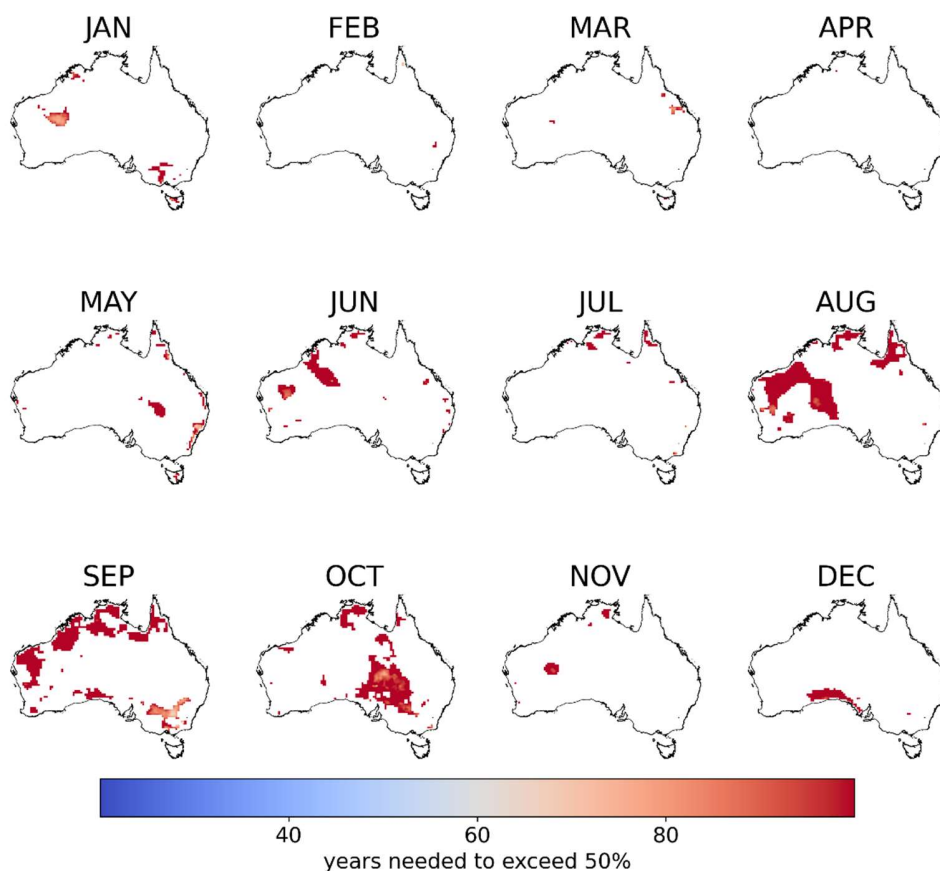
293 **Response:** Thank you for the comments. We create figures to show the time needed for the
294 departure of climatology forecasts which does not consider temporal trends from the
295 calibrated forecasts with reconstructed trends. Here we considered both 10% and 50%
296 departure. As we explained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

299



300

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



303

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

306

307 **As suggested by the above plots, it will take about 20-30 years for the departure to reach**
 308 **10%, and more than 100 years to reach 50%. However, we believe correcting time-dependent**
 309 **errors is still necessary, since increasing extreme weather conditions across the globe in**
 310 **recent years indicate that climate change is intensifying. We add the following discussions to**
 311 **section 4.1:**

312 “Although it may take decades for climate change to substantially alter the magnitude of ETo
 313 (Figures S11 and 12), we recommend that future GCM-based ETo forecasting should still correct
 314 time-dependent errors. More frequent extreme weather events in recent years support model
 315 projections that climate change will intensify in the future (Kharin et al., 201), and may induce
 316 more significant temporal trends in ETo.”

317

318 Point #8

319 [Line 132] " t_m is approximately the middle year": does moving t_m has an impact on generated
320 forecasts? I believe not because it is compensated by the value of the mean parameter μ . Please
321 confirm. If this the case, please highlight that the position of t_m is arbitrary and does not affect the
322 forecasts.

323 **Response: Thank you for the valuable comments. The reviewer is correct that using different**
324 **years as the reference for trend removal will impact the magnitude of the resultant**
325 **detrended data (both forecasts and observations), but will not affect the trend**
326 **reconstruction. When using a different year other than 2004 as a reference year, all**
327 **detrended data points will be larger (or smaller) by the same value than data using the**
328 **middle year as the reference. These differences will be lead to different mean and standard**
329 **deviation parameters. However, after we add the trend back (retrending) to data, the**
330 **difference will be canceled out. As a result, choosing a different reference year will not affect**
331 **the trend reconstruction and forecast calibration.**

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

333 "The position of t_m is empirically selected, but it will not affect the calibration if we choose a different
334 year as t_m "

335

336 Point #9

337 "Equation 8 shows the conditional posterior distribution of parameter $\sigma^2 | \delta^2$ ": We suggest
338 "Equation 8 shows the posterior distribution of parameter $\sigma^2 | \delta^2$ conditional on δ^2 ".

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

340

341 Point #10

342 "In equation 8, δ is the mean and σ is the standard deviation for predictors or
343 predictands.": Please move this sentence just after Equation 8. In addition, we suggest the following
344 clarification: " σ is the standard deviation for predictors or predictands extracted from the
345 diagonal of covariance matrix S (see equation 5)".

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

348

349 Point #11

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

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

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

359

360 Point #12

361 *[Line 166] “The comparison is conducted for months with large areas of statistically significant (at the*
362 *95% confidence interval) temporal trends in observed ETo.”: this approach is problematic because it does*
363 *not check the performance of the BJP-ti model when there is no observed trend. BJP-ti is more*
364 *parameterised than BJP, consequently it is always exposed to the risk of overfitting the data when there*
365 *is no trend, i.e. when trend parameters cannot be calibrated reliably. Please comment on this point and*
366 *justify why performance assessment excluded month with no significant observed trend.*

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

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

376

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

380

381

382

383

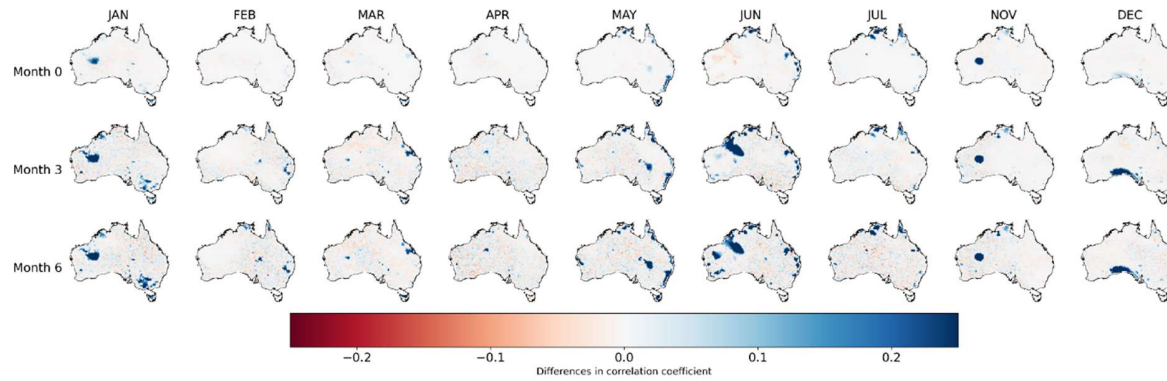
384

385

386

387

388



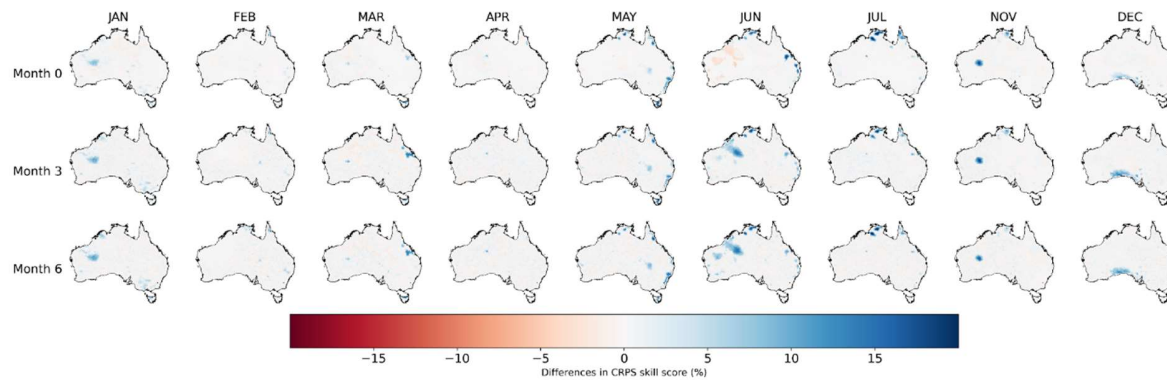
389

390

391

Figure S6. Differences in correlation coefficient between BJP-ti calibrated forecasts and observations with that between BJP calibrated forecasts and observations for nine selected months and three lead times (months 0, 3, and 6)

392



393

394

395

Figure S7. Differences in CRPS skill score between BJP-ti calibrated forecasts and observations with that between BJP calibrated forecasts and observations for nine selected months and three lead times (months 0, 3, and 6)

396 **Point #13**

397 [Line 197] “ $\delta \neq (\delta - i)$ is raw or calibrated forecasts of ETo (mm month-1)”: This is a deterministic
398 metric, so we believe that $x(t)$ is the mean of raw or calibrated forecast. Please clarify.

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

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

405

406 **Point #14**

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

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

417 “Compared with findings from previous investigations, observed trends identified in this study
418 also demonstrate significant spatial variability and varying magnitudes in different months
419 (Donohue et al., 2010; McVicar et al., 2012). We found more positive trends in our study period
420 (1990-2019) than the period of 1981-2006 (Donohue et al., 2010) ”

421

422 **Reference:**

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

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

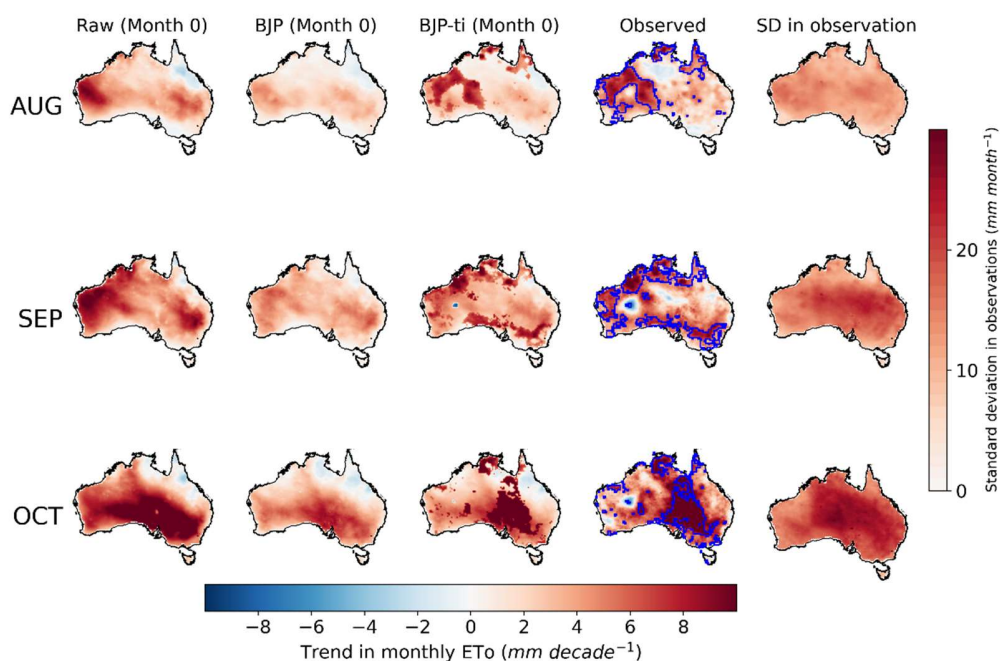
430

431 Point #15

432 [Figure 1.] We suggest adding the standard deviation of annual E_{T0} in the first column of figure 1 to
433 highlight the significance of trend values. It is important to understand if the observed trends of 6 to 8
434 mm/decade reported below are large compared to climatological variance.

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

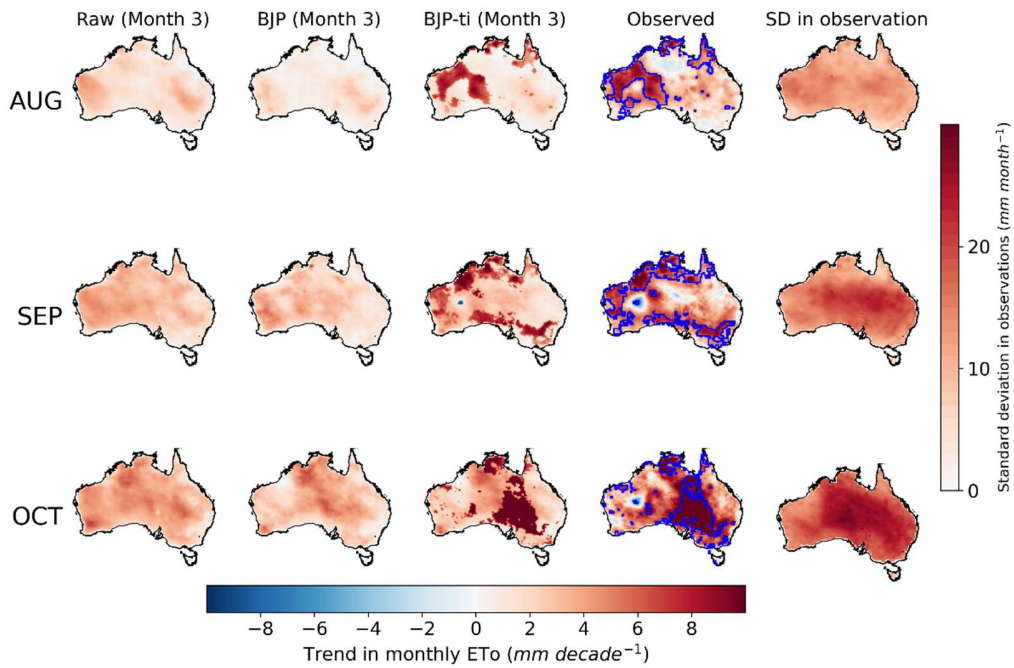
440



441

442 **Figure 1.** Trends in raw forecasts, BJP calibrated forecasts, and BJP-ti calibrated forecasts at the lead
443 time of month 0, and observed E_{T0} in August, September, and October. Blue polygons show regions
444 where observed trends are statistically significant. SD refers to standard deviation.

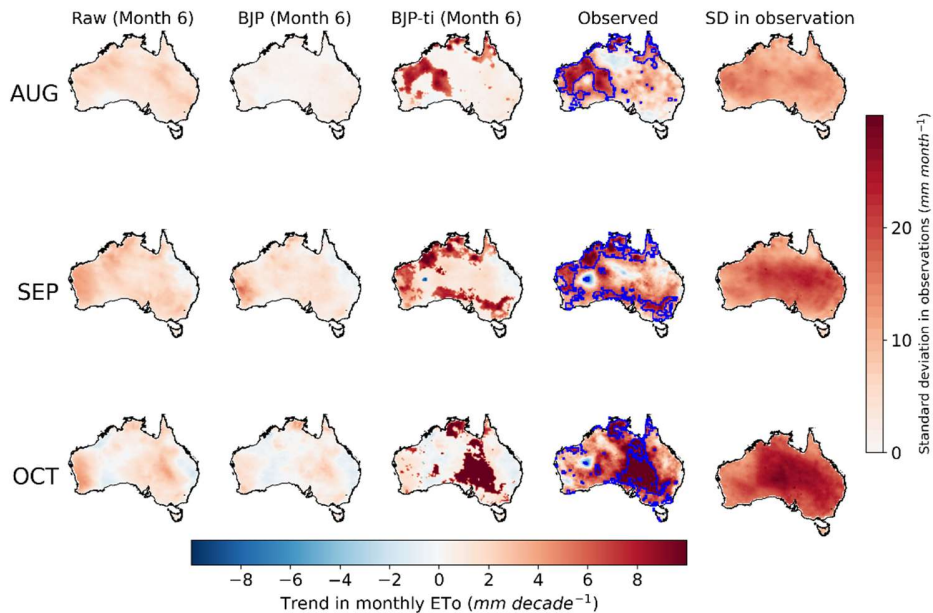
445



446

447 Figure S2. Trends in raw forecasts, BJP calibrated forecasts, BJP-ti calibrated forecasts for
 448 Month 3, and observed ETo for three selected months. Blue polygons show regions where
 449 observed trends are statistically significant. SD refers to standard deviation.

450



451

452

453

454

451 Figure S3. Trends in raw forecasts, BJP calibrated forecasts, BJP-ti calibrated forecasts
 452 for Month 6, and observed ETo for three selected months. Blue polygons show regions
 453 where observed trends are statistically significant. SD refers to standard deviation.
 454

455

456 Point #16

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

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

467

468 Point #17

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

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

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

476

477 Point #18

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

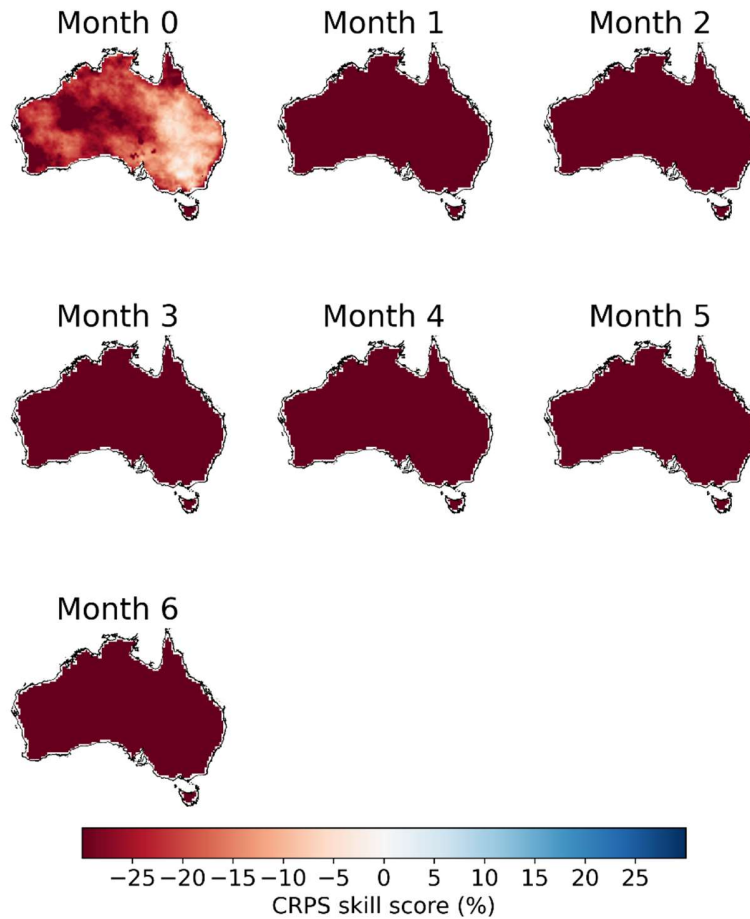
481 **Response: Thank you for the comments. After we resolve the overfitting issues, degradation**
482 **in forecast skills is removed. Please see details in our response to your comment #3.**

483

484 Point #19

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

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



495

496

497

498

499

500

501

502

503

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

504 **As a result, we feel simple bias-correction methods may not be sophisticated enough for**
505 **calibrating seasonal ETo forecasts. We add the above figure to the Supplementary Material to**
506 **show that we are aware that simple bias correction is often used to post-process raw ECMWF**
507 **forecasts. We also highlighted that simple bias correction is not sophisticated enough to**
508 **produce skillful ETo forecasts:**

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

513

514 **In addition, since the primary objective of this investigation is to understand how trend**
515 **reconstruction would affect forecast calibration, we decided to use the raw ETo forecasts for**
516 **this current investigation because we are not clear how would the quantile mapping affect**
517 **trends in ECMWF forecasts.**

518 **However, we totally agree with the reviewer that improving the raw forecasts of ECMWF**
519 **forecasts will be a very interesting point which needs further investigation. Trends in**
520 **individual input variables (e.g., temperature, vapor pressure, and solar radiation) needed for**
521 **ETo calculation have been reported by Donohue et al. (2010) and McVicar et al. (2012). It is**
522 **not clear whether correcting bias and reconstructing trends in each of the input variables**
523 **first, prior to calculating the raw ETo forecasts, will further enhance the ETo forecasts**
524 **calibration. We highlight this point in our Future work section (4.3):**

525 “In this study, we directly use the raw forecasts of individual input variables (e.g., temperature,
526 solar radiation, and vapor pressure) to construct the raw ETo forecasts. However, trends in these
527 variables have been reported in previous investigations. Whether correcting errors including
528 time-dependent errors in the raw forecasts of each input variable, will lead to more skillful
529 calibrated ETo forecasts, warrants further investigation.”

530

531 **Reference:**

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

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

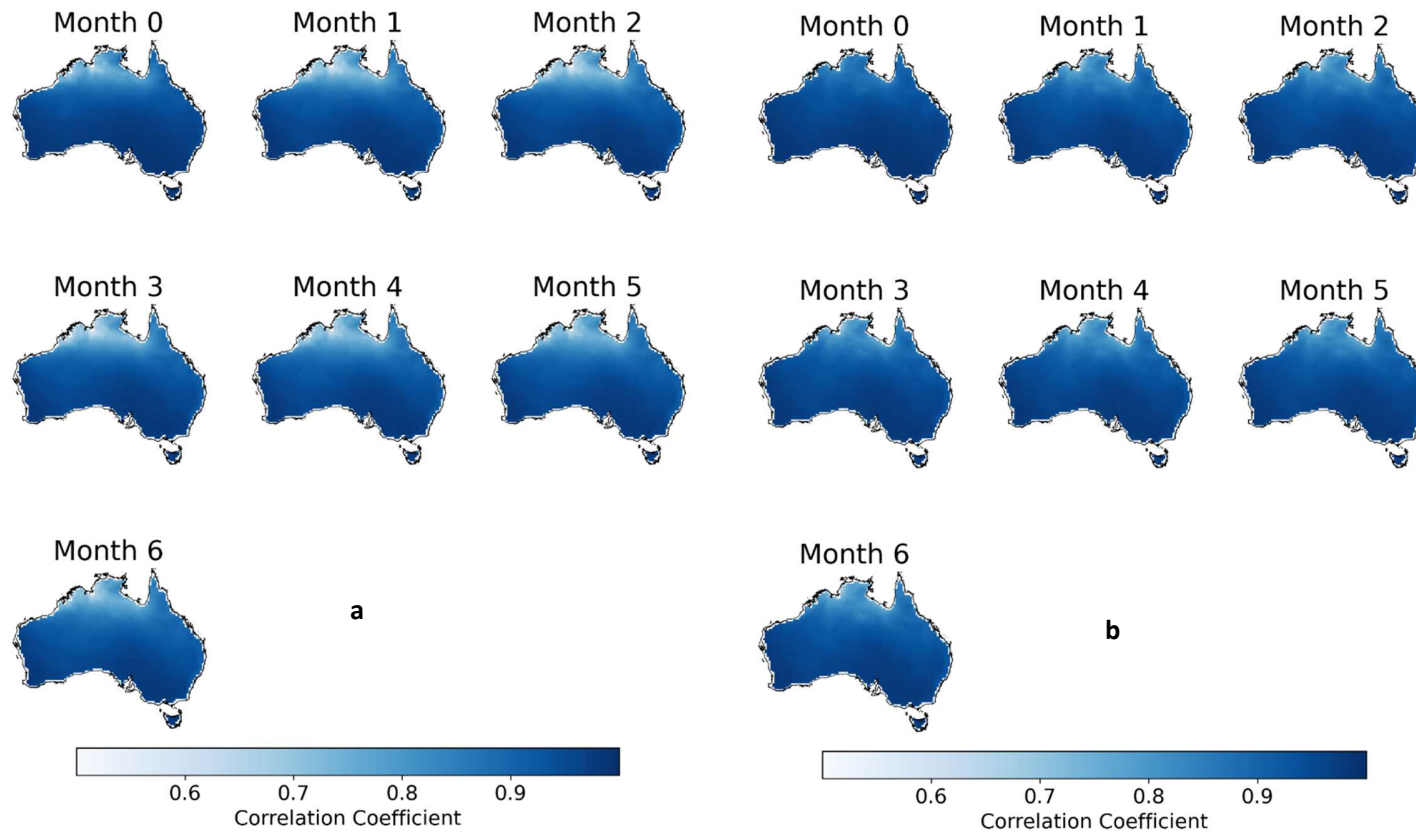
539

540 Point #20

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

543 **Response: Thank you for the valuable suggestions. We agree with the reviewer that the**
544 **correlation coefficient could be less impacted by the systematic errors in raw ECMWF**
545 **forecasts than other metrics. We calculated the correlation coefficients between raw/BJP-ti**
546 **calibrated forecasts and observations. Because of the high seasonality in ET_o , both raw and**
547 **calibrated forecasts demonstrate high correlations with observations:**

548



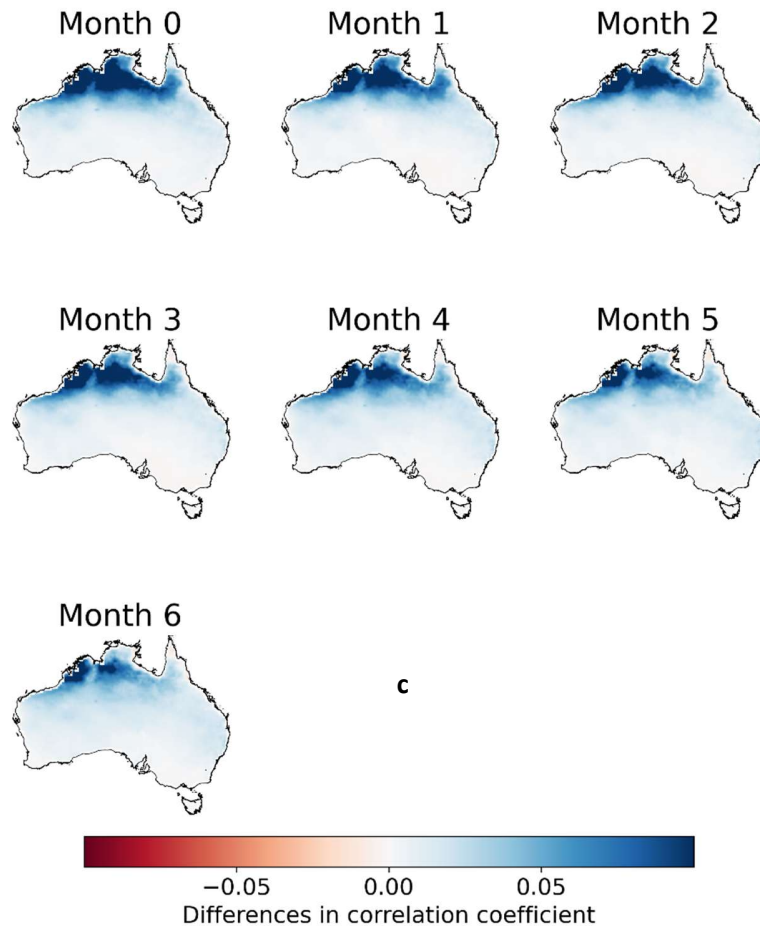
549

550

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

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

554



555

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

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

559 **Since the correlation plots for (a) raw and (b) calibrated forecasts are very similar, making it**
560 **hard to tell the difference, we decided to keep (b) and (c) in the main text (Figure 8 in the**
561 **revised manuscript) and present (a) in the Supplementary Material (Figure S10).**

562 **We add the new section in the main text to demonstrate the evaluation of the performance**
563 **of calibration in improving correlation coefficients:**

564 **“3.5 Correlation between raw/calibrated forecasts and observations**

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

573

574 Point #21

575 *Same comment than for Line 290.*

576 **Response: We understand the reviewer’s concern about how we evaluate the raw forecasts.**
577 **As we explained in our response to your comments #19 and #20, we further 1) applied bias-**
578 **correction to raw forecasts, 2) highlighted the necessity of improving individual input**
579 **variables prior to the calculation of raw ET_o forecasts, and 3) used the correlation coefficients**
580 **as another evaluation metrics to show the performance of raw forecasts. Please see details in**
581 **our response to your comments #19 and #20.**

582

583 Point #22

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

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

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

595 *“This study further demonstrates the feasibility for the general application of BJP-ti to different*
596 *hydroclimate variables showing temporal trends (Shao et al., 2021b, 2021c). The successful application to*
597 *ET_o forecasts confirms the robustness of trend reconstruction algorithms based on the data*
598 *transformation, Bayesian inference, and using statistical significance of observed trends to deal with*
599 *overfitting of trend parameters in the BJP-ti model. We also anticipate that the BJP-ti algorithms for trend*
600 *reconstruction could be adopted by other calibration models to enhance seasonal forecast calibration.”*

601 Point #23

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

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

609

610

611

612

613

614