

1 **Manuscript #hess-2020-595**

2 **Interactive comment on “Robust historical evapotranspiration trends across**
3 **climate regimes” by Sanaa Hobeichi et al.**
4

5 We would like to thank the referees for their constructive comments on our manuscript. This document
6 outlines our responses to their comments. We provide a track changed version of the manuscript to
7 highlight the changes made to the manuscript and the supplementary material.

8 In addition to the suggested changes by the two referees, we have further improved the analysis by
9 introducing a parallel, complementary dataset version to DOLCE V2.1, DOLCE V3, that has fewer parent
10 datasets than V2.1, reducing the number of temporal tiers and temporal discontinuities found in DOLCE
11 V2.1, mostly over the tropics. DOLCE V2.1 remains a more optimal dataset in many senses as it
12 minimises bias and maximises correlation with in-situ observation, whereas V3 prioritises temporal
13 continuity. Similar to DOLCE V2.1, the superiority of DOLCE V3 over its parents is demonstrated using an
14 out-of-sample testing approach.

15 DOLCE V3 is presented alongside DOLCE V2.1 throughout the manuscript and has not resulted in any
16 new sections or qualitative change to the manuscript. The main change is in section ‘3.5 Changes in ET
17 since 1980’, in which DOLCE V3 was used instead of DOLCE V2.1 to carry out the analysis of trends. The
18 new results show that trends in DOLCE V3 ET are mostly within the range of trends in available ET
19 datasets, unlike DOLCE V2.1 whose temporal inconsistencies resulted in higher trends than the available
20 datasets mostly over the wet ET regimes. We have amended related text, figures and tables accordingly.
21 These updated results also help to address the concerns of the referees, as outlined below.

22

23 **Response to Referee 1, Jasper Denissen**

24 **General Comments:**

25

26 1. Because of the efforts made to improve on DOLCE v1, a big part of the paper is about verifying DOLCE
27 v2 against the parent datasets and in-situ observations. I think the title should reflect that.

28 We agree with the referee in that the technical side of the paper which includes improving DOLCE,
29 comparing it with its parents, and verifying it against in-situ observations constitute a big part of the
30 paper. However, given that we are not publishing this work as a data paper, we chose a title that
31 highlights the scientific side of the this work, which is the robust analysis of trends in ET. Furthermore,
32 most of the figures that show the performance of DOLCE against in-situ observations and highlight its
33 superiority over its parent datasets have now moved to the supplementary material as suggested by
34 both referees, and as a result of this, 4 out of the 7 figures included in the main manuscript are now
35 focused on the assessment of trends.

36 2. The authors put a lot of effort on trying to find ways to improve DOLCE v2 in comparison to DOLCE v1
37 by i) weighting groups (Figure 1; right column) and ii) Bias correction strategies (Figure 2;

38 Supplementary Figure 3) as described in 2.2.4, 2.2.5 and 3.1. Despite the authors efforts to
39 significantly improve the ET estimates, I think that the added complexity does not justify the little
40 improvement gained. I would suggest moving sections 2.2.4, 2.2.5 and 3.1 and corresponding figures
41 to the Supplementary materials. In the main text, the authors can shortly motivate the weighting
42 group and bias correction strategy by referring to the Supplementary materials. This would make the
43 derivation of DOLVE v2 ET more straight-forward and benefit the readability of the manuscript.

44 We agree with the referee that the clustering methods added little improvement to the weighting,
45 however a little improvement is better than nothing. Technically, grouped weighting requires
46 aggregating the time and/or space domains prior to applying the weighting technique and adds a small
47 amount of work and a few additional seconds of processing time. It is therefore worth investigating the
48 efficacy of grouping approaches, and choosing the methods that adds the most improvement to the
49 weighting even if the improvement is marginal. Furthermore, as detailed in the text, most of the
50 weighting strategies have been previously suggested as ways to improve merging but have not been
51 tested. Therefore, testing them here provides valuable information to the science community.

52 To address the referee's concern and increase the readability of the manuscript, we have now moved
53 sections '2.2.5 Bias correction strategies' and '3.1 Selection of a grouping strategy' and associated
54 figures to the supplementary material.

55 **Specific comments:**

56 3. Throughout the paper, please properly introduce i) table contents, ii) figure axes labels and color
57 codes and iii) statistics of box-plots.

58 We thank the referee for his comment. We have now explained the figures and the tables further
59 throughout the text and in the captions.

60 4. lines 20: I found the notion that these climatology clusters / climate regimes are able to summarize
61 or even replace the Köppen-Geiger climate regimes quite interesting. That would be worth a mention
62 in the abstract, also putting this sentence into context.

63 We thank the referee for his suggestion, we have now mentioned the agreement of these classes with
64 the Köppen-Geiger climate regimes in the abstract

65 *The new clusters include three wet and three dry regimes and provide an approximation of Köppen-*
66 *Geiger climate classes.*

67 5. lines 23-25: "We find that despite robust . . . ET clusters". This is the only time this is mentioned in
68 the entire manuscript. I don't see the relevance of it for the abstract.

69 Good point. We have now shown this finding in section 3.5.3 Global annual trends across the ET regimes

70 ... *Our results indicate that decreasing ET trends observed in some regions oppose the consistent positive*
71 *trends across the majority of ET clusters.*

72 6. lines 82-84: FLUXCOM also belongs in this summation of gridded datasets that successfully exploit in-
73 situ measurements.

74 Here we are listing the studies that applied fusion techniques attempting to match a global dataset that
75 is deemed more reliable than the original datasets. FLUXCOM was listed earlier in the introduction with
76 the machine-learning based datasets:

77 *.... techniques including machine-learning algorithms (Jung et al. 2010; Hamed Alemohammad et*
78 *al. 2017; Jung et al. 2019), typically incorporating a range of remote sensing inputs.*

79 7. lines 198-200: I assume that the ‘very large spatiotemporal domains’ are equal to the spatial and
80 temporal resolutions of the ET data sets? Do the authors mean that through time and varying wind
81 directions, you might actually get closer to the grid cell mean than looking at individual days?

82 We thank the referee for his comment. ‘Very large spatiotemporal domains’ means over many sites and
83 time steps. This paragraph is trying to say that weights are derived by assessing the agreement between
84 flux tower measurements and the value of underlying grid cells over many locations and time steps. This
85 however assumes that the point scale of the flux towers can represent the grid scale. We don’t expect
86 this representativeness to be true at each site, however the ensemble of flux tower observations as a
87 whole do represent the underlying grid cells. This has been thoroughly tested and validated in previous
88 work. We have now made the paragraph clearer:

89 *First, weights for each product are constructed over very large spatiotemporal domains, i.e.*
90 *more than 13000 space-time records as described below, so that the (assumed stochastic) biases*
91 *of individual sites relative to grid cell values are unlikely to influence weights over a large*
92 *sample. In fact representativeness of point-scale measurement for the grid scale does exist*
93 *across all the flux tower sites as a whole, this has been verified by (Hobeichi et al. 2018).*

94 8. lines 223-230: Just out of interest: What is the total amount of days initially available from all sites?
95 After the filtering based on data availability, how many days are left?

96 The original sites data was a mix of half-hourly, daily and monthly data. The majority of daily data come
97 from Ameriflux sites. The raw Ameriflux data consisted of 147 sites with daily data and a total of 191,583
98 daily records. After quality control and filtering, the number of sites and daily records dropped to 56 and
99 81,142 respectively.

100 9. line 241: Could you elaborate on these conditions? As I’m not a flux tower measurements expert: Are
101 these 20 and 30 W m⁻² thresholds usually applied? Is there a paper where this methodology is also
102 applied?

103 Good point, our response to this comment is now included in the text

104 *A study by Paca et al. (2019) examined the changes to flux tower LE by three means of*
105 *corrections, and found that these on average differ by around 20 Wm⁻² from one another. On*
106 *this basis, we expect that typically, the correction of flux tower LE should not exceed 20 Wm⁻²,*
107 *unless errors in other components of the budgets are propagating in the corrected ET. The rule*
108 *for correcting small fluxes and the condition in which each rule is applied (i.e. LE= 30 Wm⁻²) are*
109 *in part subjective and in another part based on a case by case assessment of changes induced to*
110 *ET by the correction techniques, and achieve a reasonable trade-off between data quality and*
111 *availability.*

112 10. lines 246-249: How do you justify using LE values without any correction – any value is better than
113 nothing? Did you verify the differences between i) only LE data with correction and ii) LE data with
114 and without correction? Are there any biases there?

115 All ET measurements are prone to systematic errors, here we are using the physical constraints of the
116 energy balance to minimize these errors. Constraining ET this way is a ‘plus’ rather than a ‘must’. ET
117 measurements, despite systematic errors in them, provide the most reliable information on ET and can
118 still be used for ground-truthing gridded ET estimates. There is certainly a bias between i) and ii),
119 however the majority of the sites, this bias is small relative to the values in gridded estimates.

120 11. lines 250-255: Why not just take an average of the different towers within one grid cell, weighted by
121 fractional cover of biome within that grid cell? I would assume that you want to have the possibility
122 of retaining as much data as possible, as there could be data gaps in the data records in the tower
123 measurements which could be filled with values from towers in the vicinity. And an additional
124 question: Do you somehow account for varying footprints for the in-situ ET measurements?
125 Depending on the location of the flux tower within the grid cell, what the tower measures could be
126 from a neighboring grid cell.

127 Weighting ET measurements from two neighbouring towers located in two different biomes, with one
128 dominant and another less dominant is expected to provide a better ET than only taking ET from the
129 dominant biome. However, in all these cases, flux towers have different years of coverage. Therefore,
130 considering data from both sites will retain more data but at the same time will lead to temporal
131 inconsistency in the timeseries of the combined ET.

132 No, we haven’t accounted for varying footprints for the in-situ measurements. However, we have
133 visually assessed the position of each flux tower in the grid cell using ArcGIS, and we have found that of
134 the 260 sites used in this study, only two sites are close to the edge of the grid which means that their
135 footprint might extend to that neighbouring grid cell. However, at both sites, the land cover and the
136 climate within the expected footprint of the tower (< 2000 m) across the underlying and neighbouring
137 grid cells is the same.

138 12. lines 272-275: Are w_k are the weights, which sum to 1, based on the error correlation in
139 Supplementary Figure 2? If so, please state that here for clarity.

140 We thank the referee for his suggestion. We have now added a reference to Figure S2 in the text

141 *The analytical solution to this problem accounts for both the performance differences between*
142 *the parent datasets and their error covariance (Fig. S2), a proxy for dependence.*

143 13. lines 296-301: “the discrepancy between DOLCE and actual ET at any spatial scale greater than that
144 of a tower footprint should be less than this uncertainty estimate” I am slightly confused here: This
145 would only be true for spatial scales greater than the tower footprint and smaller than 0.25x0.25
146 degree grid cell res- olution (which is used in the study), right? If the spatial scale would be greater
147 than 0.25x0.25 degrees, the discrepancy should be even larger?

148 The referee is right. We have now made this clear in the text

149 *... the discrepancy between DOLCE and actual ET at any spatial scale greater than that of a tower*

150 *footprint and smaller than that of DOLCE should be less than this uncertainty estimate*

151 14. lines 351-353: This might be a naïve question, but is that not just because of the small number of flux
152 towers on the SH? I could imagine that due to the higher availability of in-situ measurements and
153 abundance of other measurements net- works ET estimates are much more well constrained in the
154 NH than the SH.

155 This is true to a large degree for the data driven approaches such as machine learning ET estimates,
156 which will have their reliability degraded in areas with less observations. However, we cannot assert
157 that the lower availability of in-situ measurements in the Southern Hemisphere is driving disagreement
158 among modelled and remote sensing ET estimates given that these typically do not rely on flux
159 measurements.

160 15. 357-359: I was confused here, as I assumed the authors were referring to boreal summer. Please
161 clarify.

162 We thank the referee for spotting this. Incorrect months were originally attributed to summer-fall and
163 winter-spring seasons in the Northern and Southern hemispheres. We have now corrected the text by
164 listing the correct months in each season

165 *We consider two combined seasons i.e. summer-fall and winter-spring. In the summer-fall*
166 *season, we constrain the weighting with (1) monthly observations from sites located in the*
167 *Northern hemisphere during the period June–November, and (2) monthly observations from sites*
168 *located in the Southern hemispheres during the December–May. The remaining observational*
169 *data is used to constrain the weighting during the winter-spring combined season.*

170 16. lines 380-282: I do not understand what is meant with ‘extrapolating the bias field’. Please explain
171 more clearly.

172 The Bias field is the ET bias values described in lines 378 – 380. We have spatially extrapolated the ET
173 bias values from the grid cells containing the sites to the entire global land. We have now made this
174 clearer in the text.

175 *We then assign those ET bias values, or bias field, to the grid cells containing the sites. Finally,*
176 *using the bias values at these grid cells, we extrapolate the bias field spatially to the entire global*
177 *land domain within the tier using several different extrapolation strategies,...*

178 17. Figure 3: Does zonal ET follow a Gaussian distribution? If not, it would make more sense to define
179 the grey ribbon as the interquartile range instead of the standard deviation. Also, I would suggest
180 making a mask of the grid cells where all of the parent data sets have values, so that a comparison is
181 fairer. The figure without the mask could then be moved to the supplementary material.

182 We haven’t examined the distribution of zonal ET, and the uncertainty standard deviation term is not
183 describing the latitudinal spread of ET around the mean (i.e. zonal standard deviation), but rather it
184 shows the range of DOLCE ET values bounded by its uncertainty, i.e. $DOLCE \pm \text{uncertainty}$. The
185 uncertainty of DOLCE is explained in Section 2.2.2, and is described as the ‘standard deviation of
186 uncertainty’.

187 We have now omitted 'standard deviation' from the caption to ensure that the reader does not
188 misinterpret the grey ribbon and clarified in the text that the grey ribbon is 'DOLCE \pm uncertainty'. We
189 have also replaced this plot with a new version that applies the same spatial mask to all datasets. In the
190 new figure, we have now added a similar plot for DOLCE V3 and its parents.

191 18. lines 500-501: Mentioning the seasonal cycle of DOLCE v2 ET but not elaborating on it feels out of
192 place. Either elaborate on differences between seasonal cycles between DOLCE v2 ET and others or
193 remove.

194 We thank the referee for his suggestions. We have now removed the plot of seasonal climatology

195 19. lines 525-528: Please clarify this sentence; I found it confusing as written.

196 We thank the referee for his comment. We have now clarified the sentence.

197 *The uncertainty estimate of DOLCE, however, is firmly grounded in the discrepancy between the*
198 *gridded DOLCE product and in-situ tower data. The variance of this discrepancy is used to*
199 *recalibrate the variance of the parent datasets, which are then used to estimate uncertainty,*
200 *allowing spatiotemporally varying uncertainty estimate that is both consistent with the*
201 *discrepancy between DOLCE and surface observations while at the same time being spatially and*
202 *temporally complete. This process is detailed by Hobeichi et al (2018).*

203 20. lines 540-542: Either put into context by comparing with all the other literature references or remove.
204 These two sentences seem lost.

205 We thank the referee for his suggestion. We have now removed lines 540 - 542 from the text

206 21. lines 546-550: In the first sentence the authors state the RMSE is not computed because the means
207 between DOLCE and sites are not equal. In the sentence after you explain that all data has been
208 normalized and therefore all have a zero mean. So, by normalizing you could in principle calculate the
209 RMSE?

210 We used the Taylor diagram to show how DOLCE performs at all sites, rather than a single site. To do
211 this it was necessary to normalise the data because there will be different values of the observed
212 variable across different sites. Since the data is now normalised RMSE is not helpful here. To avoid any
213 confusion, we have now removed this sentence from the paragraph.

214 22. lines 558-560: Would the signal of land cover be clearer when the authors would aggregate the land
215 cover types to short/tall vegetation? Next to that, in Supplementary Figure 6, the color legend is
216 blocking some of the extreme values.

217 We thank the referee for suggesting this. We have now combined 6 classes of broadleaved and needle
218 leaved tree covers in one group, but we still could not find any clear signal. We have made the fill color
219 of the legend transparent to make sure no values are covered.

220 23. lines 571-572: Is there a specific conclusion drawn from these figures? Otherwise mentioning them
221 seems unnecessary to me.

222 We thank the referee for his suggestion. We have now removed these plots.

223 24. Have the authors also looked at ET trends from flux tower measurements? As trends across all KS-
224 clustering defined climate regimes are positive, the flux tower observations could corroborate that if
225 they are also generally positive, right?

226 Good point. Flux tower measurements have only started around 2002, and the longest observational
227 records we had at any site was only 17 years. Therefore, there is no in-situ observations that cover a
228 long enough period to examine trends in annual ET from observations. This was already acknowledged
229 in the text in section 3.5.1 : *Unfortunately, given the absence of adequate in-situ observations that cover*
230 *a long enough period to establish trends analysis, it is difficult to validate the identified trends directly.*

231 25. line 678-680: I don't know how the fact that the global ET trends are different than the other ET
232 products reflects usefulness; the fact that the DOLCE trends are different does not necessarily mean
233 they are correct. However, it would be really interesting to see whether flux tower observations find
234 similar long-term ET trends.

235 We agree with the referee in that trends in DOLCE V2.1 are much stronger than those in other datasets.
236 After some additional investigation we have made significant changes to this section, as explained
237 below. However we do note that Table 5 shows that the global ET trends in the other datasets show
238 differences in sign, magnitude and statistical significance of the trends across all ET regimes.

239 We have now introduced DOLCE V3, a complementary dataset to DOLCE V2.1, which we now use to
240 carry out the analysis of trends instead of DOLCE V2.1. The main difference between the two datasets is
241 the number of contributing parent datasets. The focus for DOLCE V3 was reducing the number of
242 temporal tiers to reduce temporal discontinuities in DOLCE V2.1 (these were revealed in a separate
243 analysis not related to this manuscript), mostly over the tropics. We believe these discontinuities and
244 inhomogeneity lead to misrepresentative trends in some cases. We present V3 as a parallel version
245 (rather than replacement) as V2.1 still remains a better performing data set in many of the out of
246 sample tests.

247 We have now repeated the analysis of trends in Section 3.5 and have shown new trends results
248 incorporating DOLCE V3 (instead of DOLCE V2.1). The new results show that trends in DOLCE V3 agree
249 with some products more than the others, and its trends' slopes are within the range of slopes of trends
250 in the available products across all the ET regimes. This is now shown in the text and in the updated
251 Table 5

252 *We repeat the same analysis for all the participating parent datasets that span at least 30 years.*
253 *Sen's slope of the trends and their confidence interval (computed at the 95% confidence level)*
254 *are presented in Table 5. As noted earlier, trends' behaviour is deemed inconclusive when the CI*
255 *encompasses negative and positive values. These are presented with regular (as opposed to*
256 *bold) typeface and are exhibited by FLUXCOM-MET in all regimes except the driest. ERA5-land*
257 *shows downward trends in the 'M.H.ET, M.variability' and 'H.ET, L.variability' regimes. Both*
258 *GLEAM 3.5A and PLSH show upward ET trends in all regimes, with the exception of GLEAM which*
259 *shows no reliable trends in the driest and wettest ET regimes. Differences exist in the magnitude*
260 *of trends across the majority the products and the regimes. As in DOCLE V3, the strongest trends*
261 *in GLEAM 3.5A occur in the 'M.H.ET, M.variability' regime at a rate 0.5 mm year^{-1} . Finally, the*
262 *slopes of DOLCE V3 trends are within the range of slopes of trends in available ET products.*

263 Also, none of the available datasets incorporate the same degree of observational constraint in either
264 their mean field or uncertainty estimates, which makes us believe that trends exhibited by DOLCE V3 are
265 more reliable than those observed in other datasets.

266 As we mentioned in our response to the previous comment above, there are not enough in-situ
267 observations that cover a long enough period to examine trends in annual ET directly from observations.

268

269 **Technical corrections**

270

271 26. line 21: 'at each location'. Do you mean globally?

272

273 Yes, we have now made it clear in the text by adding 'on land'. The sentence now reads:

274

275 *... we derive novel ET climatology clusters for the land surface, based on the magnitude*
276 *and variability of ET at each location **on land**.*

277

278 27. line 42: remove the comma after 'approaches'

279 We thank the referee for spotting this. We have now removed the coma after 'approaches'.

280 line 113: replace 'trends (5) behavioural' with 'trends and (5) behavioural'

281 Thank you for suggesting this. We have now made the change in the text.

282 28. lines 236: to avoid confusion, maybe rephrase "latent heat measurements are used directly" to
283 "latent heat measurements are used without any corrections".

284 We thank the referee for his suggestion. We have now replaced 'directly' with 'without any corrections'.

285 29. line 435: Replace 'Fig. S3' with 'Fig S4.'

286 We thank the referee for spotting this. We have now made the change in the text.

287 30. line 590: replace 'intensified' with 'increased'

288 We thank the referee for his suggestion. We have now made the change.

289 31. line 624: replace 'modified' with 'modified'

290 We thank the referee for spotting this. We have now made the correction.

291 32. line 743: replace 'Figure2' with 'Figure 2'

292 We thank the referee for spotting this. We have now added a space after 'Figure'.

293 **Response to Referee 2**

294 **General Comments:**

295 Overall this work seems like a useful addition to the literature. I have some detailed comments for
296 clarification. The results and discussion section of the paper often for large parts mostly just list what is
297 shown in the figures, but it would make it a lot more interesting to more read about what the figures
298 teach us. In addition, please check if small things table contents, figure axes, etc. are introduced. Often
299 this seems somewhat lacking.

300 We thank the referee for his comment. We have now explained the figures and the tables further
301 throughout the text and in the captions.

302 Detailed comments

303 1. L13: why “gridded”?

304 The scale of ET observations from Eddy Covariance towers is typically less than 1000 m which does not
305 allow to study ET at the regional scale. Therefore, gridded ET datasets are needed to understand ET at
306 the regional scale.

307 2. L19: “After successful evaluation of the efficacy”: a “successful evaluation” does not say
308 anything about the efficacy, so please rephrase.

309 We thank the referee for their suggestion. We have now changed “after successful evaluation of
310 the efficacy of these uncertainty estimates out-of-sample” to:

311 *after demonstrating the efficacy of these uncertainty estimates out-of-sample*

312

313 3. L19: coverage, rather than reach?

314 We thank the referee for their suggestion. We have now changed ‘reach’ to ‘coverage’

315 4. L33: “with different scopes” is unclear in its meaning to me.

316 We thank the referee for their comment. We have now clarified ‘with different scopes’ in the text.

317 with different scopes (e.g. addressing key questions in ecology, hydrology, or other
318 disciplines),

319 5. L35: “typically incorporating a range of remote sensing inputs” would benefit from some
320 citations.

321 We thank the referee for their comment. Citation was already included before “typically incorporating a
322 range of remote sensing inputs”, we have now moved it to the end of the sentence.

323 6. L36: “have been recognised for their potential to outperform single source datasets” can the
324 strength of these methods be made in a more explicit statement that is more specific?

325 We thank the referee for their suggestion. We have now specified the strength of merging methods.

326 ... have been recognised for their potential to outperform single source datasets in reducing bias
327 against in reducing bias against tower-based eddy-covariance ET measurements

328

329 7. L40: time resolution (rather than step)?

330 We thank the referee for their suggestion. We believe that both ‘time resolution’ and ‘time step’ can be
331 used here interchangeably.

332 8. L43: chemical seems redundant?

333 We thank the referee for their comment. However we couldn’t find any redundancy.

334 9. L70: physically-based

335 We thank the referee for spotting this, we have now made the correction in the text

336 10. L70: which ET trends did Pan look at?

337 Pan looked at ET trends during 1982-2011. We have now specified this in the text.

338 11. L142-147: it seems some references could be added here?

339 References of these products were given in the describing paragraphs that follow. We have now added
340 those references in L142-147.

341 12. Section 2.2.4. I do not suggest to redo the analysis, but why aren’t weighing groups
342 considered based on their physical similarity linked to ET (e.g. landcover) rather than these
343 currently somewhat oddly chosen groups?

344 Good point. We agree with the referee in that the most obvious weighting groups are land cover and
345 climate zones. We have tried both grouping methods in a paper describing the first version of DOLCE,
346 but this did not improve the final hybrid product. This was explained in L336 – L337

347 *Hobeichi et al. (2018) tried to group flux tower sites based on their land cover type and*
348 *computed weights for each land cover type. However, this approach did not improve the results,*
349 *whether grouping by climate zone or aridity index, with the main reason being attributed to the*
350 *small number of sites in many groups.*

351 We disagree with the referee that the grouping approaches are ‘somewhat oddly chosen’. We have
352 clarified in ‘Section 2.2.4 Weighting groups’ the motivation behind each grouping method.

353 We have now added a new weighting group that considers both physical similarities linked to ET, and
354 seasons. We have applied this grouping to derive the new version 3 of DOLCE which we now use to
355 examine ET trends. We have now explained the new grouping method in Section 2.2.4.

356 • *Grouping by ET regime and months: Land was classified into three distinct broad ET regimes*
357 *(Fig. S4) according to two aspects of ET, mean annual total ET and within-year relative variability*
358 *throughout 1980 – 2018, derived from GLEAM V3.5a, and using K-means unsupervised*
359 *classification (MacQueen, 1967). We explain the classification method further in section 3.5.2.*
360 *Different sets of weights were computed at each ET regime during June–November and*
361 *December–May. Implementing weighting this way ensured that we account for performance*
362 *differences across different physical aspects of the land and seasons. Despite that observational*

363 *data was divided into six distinct groups, the observational data available in each group was still*
364 *appropriate to merge the four parent datasets of DOLCE V3. However, we found this grouped*
365 *weighting strategy not appropriate for merging 11 parent datasets of DOLCE V2.*

366 13. Table 1: indicate what a (lack of) marker means. It's somewhat obvious but it's still good to
367 specify. . .

368 We are not sure what the referee is asking us to specify, but it seems from the referee's comment, the
369 suggested change is not that important....?

370 14. Table 3: why are uncertainties this large for DOLCE V2?

371 DOLCE's uncertainty gives an accurate upper bound estimate of the likely discrepancy between the
372 product and unseen ET measurements. Uncertainties in DOLCE V2 are large compared to uncertainties
373 of hybrid estimates derived by different merging approaches which typically consider the spread the
374 parent datasets. This has been clarified in Section 2.2.2

375
376 *This process ensures that the computed uncertainty provides a better uncertainty estimate of the*
377 *hybrid ET than simply using the spread of the parent datasets.*

378 *One additional advantage of defining uncertainty in this way is that it should give an accurate*
379 *upper bound estimate of the likely discrepancy between the product and unseen ET*
380 *measurements at a range of spatial scales. That is, since it is based on the discrepancy of the*
381 *final hybrid product and point-based flux tower estimates, which are essentially at the extremes*
382 *of spatial discrepancy, the discrepancy between DOLCE and actual ET at any spatial scale greater*
383 *than that of a tower footprint should be less than this uncertainty estimate (noting however that*
384 *this is the estimated standard deviation of uncertainty, rather than a hard upper limit)*

385

386 15. Table 5: specify unit of the trends.

387 We thank the referee for their suggestion. We have now specified the unit of the trend as mm year^{-1} .

388 16. Figure 1: is this necessary to include in the main paper, or could it be supplementary
389 materials?

390 We thank the referee for their suggestion. We have now moved this Figure to the supplementary
391 material.

392 17. Figure 2: idem

393 We thank the referee for their suggestion. We have now moved this Figure to the supplementary
394 material.

395 18. Figure 3: can more distinguishable lines styles (i.e. color, thickness etc) be used better allow
396 interpreting this figure?

397 We thank the referee for their suggestions. We have now improved this figure and made the lines more
398 distinguishable.

399 19. L759: reliable or robust?

400 We have now rephrased the caption:

401 *Spatial pattern of ET climate trends in DOLCE V3 over 1980 – 2018 derived using Mann-Kendall*
402 *and Sen’s slope methods. Grid cells in white correspond to unreliable ET trends because (i) the*
403 *confidence interval of the slope encompasses a mix of negative and positive values; or (ii) trends’*
404 *slopes computed for multiple different random samples of ET within the interval $ET \pm$ uncertainty*
405 *do not agree in sign.*