# Robust historical evapotranspiration trends across climate regimes

3 Sanaa Hobeichi<sup>1,2</sup>, Gab Abramowitz<sup>1,2</sup>, and Jason Evans<sup>1,2</sup>

<sup>5</sup> <sup>1</sup>Climate Change Research Centre, UNSW Sydney, NSW 2052, Australia.

<sup>6</sup> <sup>2</sup>ARC Centre of Excellence for Climate Extremes, UNSW Sydney, NSW 2052, Australia.

7

4

8 Correspondence: Sanaa Hobeichi (s.hobeichi@unsw.edu.au)

## 9 10 Abstract

11 Evapotranspiration (ET) links the hydrological, energy, and carbon cycle on the land surface. Quantifying 12 ET and its spatiotemporal changes is also key to understanding climate extremes such as droughts, 13 heatwaves and flooding. Regional ET estimates require reliable observationally-based gridded ET 14 datasets, and while many have been developed using physically-based, empirically-based and hybrid 15 techniques, their efficacy, and particularly the efficacy of their uncertainty estimates, is difficult to 16 verify. In this work, we extend the methodology used in Hobeichi et al. (2018) to derive two new 17 versions of the Derived Optimal Linear Combination Evapotranspiration (DOLCE) product, with 18 observationally constrained spatiotemporally varying uncertainty estimates, higher spatial resolution, 19 more constituent products and extended temporal coverage (1980-2018). After demonstrating the 20 efficacy of these uncertainty estimates out-of-sample, we derive novel ET climatology clusters for the 21 land surface, based on the magnitude and variability of ET at each location on land. The new clusters 22 include three wet and three dry regimes and provide an approximation of Köppen-Geiger climate 23 classes. The verified uncertainty estimates and extended time period then allow us to examine the 24 robustness of historical trends spatially and in each of these six ET climatology clusters. We find that 25 despite robust decreasing ET trends in some regions, these do not correlate with behavioural ET 26 clusters. Each cluster, and the majority of the Earth's surface, show clear robust increases in ET over the 27 recent historical period. The new datasets DOLCE V2.1 and DOLCE V3 can be used for benchmarking 28 global ET estimates and for examining ET trends respectively.

29

# 30 **1. Introduction**

31 Understanding the spatiotemporal variability of evapotranspiration (ET) is a critical part of

32 understanding the processes that lead to high impact weather phenomena, such as droughts (Han et al.,

33 2018; Montano et al., 2015; Sheffield et al., 2012; Teuling et al., 2013), heatwaves (Teuling, 2018; Ukkola

et al., 2018) and flooding (Dawdy et al., 1972; Sharma et al., 2018). Several global gridded ET datasets

35 have been developed, using physical schemes with different scopes (e.g. addressing key questions in

36 ecology, hydrology, or other disciplines), and complexity (see Fisher and Koven, 2020), and empirical

37 techniques including machine-learning algorithms, typically incorporating a range of remote sensing

inputs (Hamed Alemohammad et al., 2017; Jung et al., 2010, 2019). Recently, ET datasets derived with a

- 39 hybrid approach have been recognised for their potential to outperform single source datasets in
- 40 reducing bias against tower-based eddy-covariance ET measurements (Ershadi et al., 2014; Feng et al.,
- 41 2016; Hobeichi et al., 2018; Jiménez et al., 2018; McCabe et al., 2016).
- 42

While most observational products are global (or near global) in their spatial extent, and typically
available with a monthly time step, different products are constrained by very different types of
observations, and vary significantly in their treatment of uncertainty. As detailed below when describing
the datasets we use here, 'physically-based' approaches use equations that represent different physical,
chemical, and biological processes and incorporate satellite-based atmospheric forcing, and
parameterization of land surface characteristics, while 'empirical' approaches integrate ground-based

- 49 measurements of ET together with satellite data and ground-based measurements of vegetation
- 50 characteristics and land surface parameters. These differences result in a diverse group of products and
- 51 estimates, but it is their approach to deriving uncertainty estimates that is arguably more important.
- 52

53 Very few datasets provide uncertainty estimates associated with the ET flux, these include datasets

54 described in Bodesheim et al. (2018) and Jung et al. (2019). In Bodesheim et al. (2018), monthly

55 uncertainty estimates are computed from the standard deviation of the half-hourly ET values that were

used to derive monthly ET averages. Jung et al. (2019) provide an ensemble of global ET estimates,

57 deviations from the ensemble median are used to derive ET uncertainties. In both cases, uncertainties

58 do not reflect the actual deviation from the measured ET at site locations. Without well calibrated

59 uncertainty estimates we are unable to tell whether an identified property of any given data set, such as

a trend or a proportion of the surface energy or water budget, is robust, rather than a result of bias or

- 61 stochastic uncertainty.
- 62

63 ET trends computed from different approaches (i.e. physical and empirical) show general agreement at

64 the global scale, and indicate that ET has increased since early 1980s (Miralles et al., 2014; Pan et al.,

65 2020; Zhang et al., 2016). However, different ET products exhibit considerable disparities in regional and

66 continental ET trends. For instance, Miralles et al. (2014) detected upward ET trends in GLEAM (Global

- 67 Land Evaporation Amsterdam Model; Miralles et al. 2011) in the northern latitudes caused by vegetation
- 68 greening. In water limited regions, they found that ET is characterised by a multidecadal variability that
- 69 follows ENSO dynamics, mainly in eastern and central Australia, southern Africa and eastern South

70 America. In comparison, ET trends estimated from the observation-driven Penman-Monteith-Leuning

- 71 (PML; Zhang et al. 2016) model show increasing ET since 1980 in the northern latitudes, arid regions in
- northern Africa, and northern and eastern Amazon. On the other hand, PML exhibits negative trends in

73 southern South America and western United States. More recently, Pan et al. (2020) found that ET

74 trends exhibited during 1982-2011 by a range of empirical and physically-based estimates disagree in

75 the direction of trend in the Amazon basin and many arid and semi-arid regions. Without incorporating

- uncertainties in ET estimates in the analysis of trends, it becomes difficult to assess the reliability of the
   established trends.
- 78

79 The gridded ET product derivation technique implemented by Hobeichi et al. (2018) offers the potential

80 for robust out-of-sample testing of its uncertainty estimates, as well as several other advantages over

81 other techniques. Like other merging approaches, it offers the potential to minimise the eccentricities or 82 biases of any one product, by averaging them (in this case using weights). However, unlike several other 83 merging techniques (Mueller et al., 2013; Paca et al., 2019; Rodell et al., 2015; Stephens et al., 2012) it 84 accounts for performance differences between parent estimates using in-situ data as the observational 85 constraint, rather than assigning weights based on the ability to match another gridded dataset that is 86 deemed more reliable, or the ensemble mean of a selection of datasets (Munier et al., 2014; Sahoo et 87 al., 2011; Wan et al., 2015; Zhang et al., 2018). The efficacy of using in-situ measurements for 88 constraining much larger scale gridded estimates has also been shown explicitly (Hobeichi et al., 2018, 89 2020b). Next, most available merging techniques do not account for dependence between parent 90 estimates, where redundant information in different parent products is likely to bias the hybrid estimate 91 (Abramowitz et al., 2019; Herger et al., 2018). Finally, and perhaps most important for this work, the 92 technique calculates global spatially and temporally varying uncertainty estimates that are 93 observationally-based, in that they are based on the discrepancy between the hybrid ET estimate and in-94 situ data. Aside from being more defensible than simply taking the spread of the parent products 95 around their mean (e.g. Pan et al., 2012, Zhang et al., 2018), this approach also allows for out-of-sample 96 testing, by leaving some sites out of the derivation of the hybrid product and its uncertainty, and then 97 using them to test its accuracy. 98

99 Despite these advantages, out-of-sample testing of uncertainty estimates was not explored by Hobeichi 100 et al (2018), and the short temporal availability of the DOLCE product (2000 – 2009) limited its 101 application, particularly in examining historical trends. While different subsets of parent products were

- 102 used over different regions to expand the spatial coverage of DOLCE, the possibility of different product
- 103 subsets in different time periods to extend its temporal reach was not explored. Additionally, since the
- 104 development of DOLCE, four of its six parent datasets (Jung et al., 2010; Martens et al., 2016; Miralles et
- 105 al., 2011; Mu et al., 2011; Zhang et al., 2016) have been improved and several new global ET datasets
- 106 have been developed (Balsamo et al., 2015; Bodesheim et al., 2018; Jung et al., 2019). Most of these are
- 107 available at a higher spatial resolution than the original 0.5° in DOLCE and cover different subsets of the
- 108 period 1980 – 2018, with at least two available every year during this period (Table1).
- 109
- 110 In this paper we amend these shortcomings and explore some of the insights that the new versions of 111 DOLCE offer, in particular focusing on the temporal trends in ET in different regions, and the assessment 112 of robustness of trends that well calibrated uncertainty estimates afford. Roughly in order, we detail 113 below: (1) how we update the DOLCE product with new parent datasets and extend its temporal 114 coverage; (2) how the improved products compare to their previous version and other existing ET 115 estimates from the literature; (3) the efficacy of uncertainty estimates, in particular whether or not they 116 are overconfident; (4) an exploration of historical trends in ET using the extended temporal coverage, 117 and how the uncertainty estimates allow us to examine the robustness of these trends; and (5) 118 behavioural ET clusters that describe ET based climate regimes, as a mean to understand the spatial
- 119 distribution of trends we find.
- 120
- 121
- 2. Data and Methods 122
- 123

124 To derive two new versions of DOLCE, one suitable for benchmarking ET dataset and another for trends

- analysis, we combine 11 and 4 available global gridded ET datasets respectively using the same merging
- 126 technique as in DOLCE V1. This technique derives a linear combination of the participating ET datasets
- based on their ability to match in-situ observations while also accounting for their error dependency.
- 128 While we acknowledge the obvious spatial mismatch between gridded and in-situ data, we refer readers
- 129 to Hobeichi et al (2018) where it was shown that in-situ observations do contain useful information
- about grid scale fluxes, using out-of-sample testing in a similar framework to the one we present here.
- 131

132 Our aim is to increase the time coverage and spatial resolution of DOLCE V1, as well as examine

- 133 strategies to improve the effectiveness of the weighting strategy. Below we detail newly available global
- datasets that allow us to derive DOLCE V2 and DOLCE V3 at 0.25° spatial resolution, and an improved
- 135 collection of in-situ constraining data. We then briefly revisit the weighting and uncertainty estimation
- approach, before describing our tiering approach to extending the temporal reach of DOLCE V2 and
   DOLCE V3. Finally, we examine alternative clustering and bias-correction approaches to improve the
- 138 out-of-sample performance of the weighting technique.
- 139

140 Throughout the paper, we use the two terms evapotranspiration (ET) and latent heat (LE)

141 interchangeably, and the unit  $W m^{-2}$  for heat fluxes and  $mm y ear^{-1}$  for the water flux equivalent. For

reference:  $1 W m^{-2} = 12.86 mm y ear^{-1}$ . As above, we refer to the product from Hobeichi et al (2018)

- as DOLCE V1 and the new products we are deriving as DOLCE V2 or DOLCE V2.1 and DOLCE V3.
- 144
- 145

146 **2.1** Data

# 147 2.1.1 Global ET datasets:

148 DOLCE V1 was derived from 6 global ET datasets: MPIBGC (Jung et al., 2010), GLEAM v2a, GLEAM v2b

149 (Miralles et al., 2011), GLEAM v3a (Martens et al., 2016, 2017), MOD16 (Mu et al., 2011) and PML

- 150 (Zhang et al., 2016). In DOLCE V2, we keep both MOD16 and PML datasets, substitute the GLEAM
- 151 products with their improved versions GLEAM3.3A and GLEAM3.3B (Martens et al., 2016, 2017), and
- replace MPIBGC with newly developed empirical ET datasets from the Max Planck Institute for
- 153 Biogeochemistry: BACI (Bodesheim et al., 2018) and two ET estimates from the FLUXCOM project (Jung
- et al., 2019). Additionally, we incorporate a recently published dataset ERA5-Land (Muñoz Sabater,
- 155 2019) and three newly available ET datasets PLSH (Zhang et al., 2015), SEBS (Chen et al., 2019; Su, 2002)
- and SRB-GEWEX (Vinukollu et al., 2011). In comparison, DOLCE V3 was derived from 4 global ET
- datasets. These are: ERA5-Land, an ET dataset from the FLUXOM project, and the two latest versions of
- 158 the GLEAM products, GLEAM V3.5A and GLEAM V3.5B. We provide a brief description of these datasets
- 159 below, with URLs and download dates shown in supplementary Table S2.
- 160 Biosphere Atmosphere Change Index (BACI; Bodesheim et al., 2018): The dataset is derived by upscaling
- diurnal cycles of ET and other land-Atmosphere fluxes from a large set of FLUXNET sites based on a
- 162 random forest regression framework. It uses seasonal vegetation variables and indices from MODIS
- satellites, and meteorological data either measured at the flux tower sites or retrieved from the ERA-
- 164 Interim data.

165 ERA5-Land(Muñoz Sabater, 2019): A global land surface reanalysis dataset that has been developed by 166 rerunning the land component of the ECMWF ERA5 climate reanalysis with a series of improvements 167 (mainly higher temporal frequency and spatial resolution) that makes it more reliable for land 168 applications. ERA5-Land is produced under a single simulation that uses adjusted atmospheric inputs from 169 ERA5 atmospheric variables without being coupled to the atmospheric module of ERA5.

170 FLUXCOM (Jung et al., 2019): An empirical upscaling of observations from 224 flux tower sites using 171 machine learning methods. The full FLUXCOM product includes 63 global ET datasets that have been 172 produced using two different setups, a remote sensing (RS) setup and a remote sensing + meteorological 173 (MET) setup. The development of the global datasets incorporates 9 machine learning techniques, 4 global 174 meteorological datasets (used only with the MET setup), 3 correction methods for energy imbalance at 175 the flux tower sites and MODIS remote sensing input. In DOLCE V2, we include one dataset from each 176 setup, that we refer to as FLUXCOM-RS (from the RS setup) and FLUXCOM-MET (from the MET setup). To 177 choose the two datasets we analysed the pair-wise error correlations of all the products against in-situ 178 flux tower and selected the two that had the lowest pair-wise error correlation (and so were deemed least 179 dependent). In DOLCE V3, we include a dataset from the MET setup only.

180 Process-based Land Surface Evapotranspiration/Heat Fluxes algorithm (PLSH; Zhang et al., 2015):

181 Terrestrial ET is derived using an improved NDVI-based Penman-Monteith algorithm originally developed in

182 (Zhang et al., 2010). ET is regulated by a set of geophysical data from GIMMS and Vegetation Index and

183 Phenology along with radiative data from World Climate Research Programme/Global Energy and

184 Water-Cycle Experiment (WCRP/GEWEX) Surface Radiation Budget (SRB) and CERES along with other

185 meteorological observations data from the NCEP/DOE AMIP-II Reanalysis (NCEP2; Kanamitsu et al.,

186 2002).

187 Surface Energy Balance System (SEBS; Chen et al., 2019; Su, 2002): ET estimates are produced with the

188 revised Surface Energy Balance System (SEBS) algorithm in Chen et al. (2013; 2019). It uses

189 meteorological observations, ground heat flux, net radiation and canopy measurements collected from

190 flux tower sites, and NDVI and emissivity data from MODIS.

191 Surface Radiation Budget (SRB)-GEWEX (Vinukollu et al., 2011): ET is estimated based on the Penman-

192 Monteith equation. Input data sets include remote sensing data from AVHRR and MODIS,

193 meteorological data derived from the Variable Infiltration Capacity (VIC; Liang et al., 1994) land surface

194 model forced by PGF and radiative data from the NASA Global Energy and Water Exchanges (GEWEX)

195 Surface Radiation Budget Project (Stackhouse Jr et al., 2011).

196

197 It is clear that different parent datasets share forcing, parameterisations, and physical and empirical

assumptions. Therefore, they do not constitute entirely independent estimates. Furthermore, their error

199 correlation (when compared with data from 254 sites – details on these below), which can be used as a

200 measure of their dependence (Bishop and Abramowitz, 2013) is high (Fig. S2, correlation > 0.5),

reinforcing the potential for benefit using a weighting approach that can account for this redundancy.

202

Part of the high correlation is of course due to spatial heterogeneity and the scale mismatch between insitu and gridded datasets since individual site locations within a grid cell are likely biased with respect to

- the (unknown) true grid cell averaged flux. While it might appear that a weighting approach that
- accounts for error correlations between parent datasets might be in danger of overfitting to error
- 207 correlation resulting from spatial heterogeneity, we have two mechanisms that ensure this is not a
- 208 concern for our final product. First, weights for each product are constructed over very large
- 209 spatiotemporal domains, i.e. more than 13000 space-time records as described below, so that the
- 210 (assumed stochastic) biases of individual sites relative to grid cell values are unlikely to influence weights
- over a large sample. In fact, representativeness of point-scale measurement for the grid scale does exist
   across all the flux tower sites as a whole, this has been verified by Hobeichi et al., (2018). Second, and
- more categorical, all results here are presented out-of-sample, so that any overfitting will degrade,
- rather than improve the results we present. More detail on this is presented below.
- 215

Given that most of the parent datasets provide ET information at a 0.25° or finer spatial resolution

- 217 (Table 1), it is possible to enhance the resolution of DOLCE from 0.5° to 0.25°. All the parent datasets are
- resampled from their original spatial resolution to a common 0.25° grid using the nearest-neighbour
- resampling method, and aggregated to monthly temporal scale before implementing the weighting
- 220 technique.
- 221
- 222

# 223 2.1.2 Flux tower data

224 We use flux tower observations from a range of networks including Ameriflux (ameriflux.lbl.gov), the 225 Atmospheric Radiation Measurement (ARM; arm.gov), AsiaFlux (asiaflux.net), European Fluxes Database 226 (europe-fluxdata.eu), Fluxnet 2015, LaThuile Free Fair Use (fluxnet.fluxdata.org), Oak Ridge data 227 repository (daac.ornl.gov), OzFlux (ozflux.org.au), and data acquired through communication with 228 individual site principal investigators (PI). Particular efforts were made to establish connections with PIs 229 in regions where ET observations are scarce, including all areas outside North America, Europe and 230 Australia, particularly the MENA regions, Siberia, Central Africa and the Amazon basin. Our efforts and 231 communications with many PIs unfortunately failed to incorporate flux data from some of these regions

- 232 (excepting those that are already available from the cited networks). Before the quality control process
- 233 detailed below, we had obtained data from 366 flux tower sites.
- 234

235 The raw data consists of a composite of half hourly, daily and monthly records. We compute daily 236 averages from half-hourly records for days where at least 80% of half-hourly LE records are available. 237 Subsequently, we compute monthly averages from daily records for months where at least 80% of daily 238 LE records are available. In DOLCE V1 we applied a less strict guality control on the observational data in 239 which up to 50% of gap filling was allowed. The reason was that DOLCE V1 incorporated much fewer 240 observational data - sourced from Fluxnet 2015 and LaThuile Free Fair use only. In order to retain 241 enough observational data to constrain the weighting, it was necessary to make a trade-off between the 242 quality and the quantity of the data.

- 243
- We also apply energy balance corrections to the monthly LE at all sites where monthly averages of the other variables of the surface energy budget - net radiation  $(R_n)$ , ground heat flux (G), and sensible

- heat flux (H) are available with the same high quality (quality flag > 80%). Corrections are carried out
- independently for every monthly record. Where any of the other components of the energy budgets are
- absent, latent heat measurements are used without any corrections. The energy balance correction is
- applied as a Bowen Ratio (BR) based correction that distributes the energy budget residuals among *H*
- and LE in such a way that their ratio is conserved. This is done under pre-defined constraints that
- disallow large changes to be applied to LE. As a result of this, we accept the BR correction and use the
- 252 corrected LE (*LE*<sub>cor</sub>) values if the original monthly LE and *LE*<sub>cor</sub> satisfy:  $\begin{pmatrix} LE_{cor} \in \begin{bmatrix} 1 & 2 \end{bmatrix}$ where LE < 20 W/m<sup>-2</sup>

253 
$$\begin{cases} \overline{LE} \in \left[\frac{1}{2} - 2\right], & \text{where } LE \leq 30 \text{ W m}^{-2} \\ LE_{cor} - LE \leq 20 \text{ W m}^{-2}, & \text{where } LE \geq 30 \text{ W m}^{-2} \end{cases}$$

254 In DOLCE V1, we did not set a threshold for LE adjustments, which resulted in LE being changed 255 drastically in a few sites to offset errors in the other energy balance components. If the BR correction 256 does not meet the above criterion, we reject the correction and try using a residual correction, which 257 simply calculates LE as the residual term in the energy balance equation, i.e.  $LE_{cor} = R_n - H - G$ . 258 Similarly, we reject the residual correction if the relation between LE and  $LE_{cor}$  above is not satisfied. In 259 this case, we use the original monthly LE values without correction. A simplified flowchart of these steps 260 is displayed in Fig. S3 in the supplementary material. A study by Paca et al. (2019) examined the changes 261 to flux tower LE by three means of correction, and found that these on average differ by around 20 Wm-262 2 from one another. On this basis, we expect that typically, the correction of flux tower LE should not 263 exceed 20 Wm-2, unless errors in other components of the budgets are propagating in the corrected ET. 264 The rule for correcting small fluxes and the condition in which each rule is applied (i.e. LE= 30 Wm-2) are in part subjective and in another part based on a case by case assessment of changes induced to ET 265 266 by the correction techniques, and achieve a reasonable trade-off between data quality and availability.

267 In a further pre-processing step, if a site is located in close proximity to other sites such that they all sit 268 on the same 0.25° grid-cell, we use observational data from the site that is more representative of the 269 underlying grid-cell. Selecting the most representative site among these sites involves 1) identifying the 270 biome cover at each site; 2) computing the fraction of the grid area covered by each biome; the most 271 representative site is the one whose biome is more abundant in the underlying grid-cell (i.e. scores the 272 highest fraction of the total area). If all sites are equally representative of the underlying grid-cell, we 273 consider them as one site and we combine monthly LE from the sites by taking the average. We use the 274 high resolution 300 m - land cover maps from the European Space Agency (ESA; http://www.esa.int/) 275 downloaded from <u>https://cds.climate.copernicus.eu/</u> to determine the biome types of neighbouring 276 sites and the corresponding grid-cells. This step has ensured that we are not matching a grid-cell with 277 inappropriate observational data. All the excluded sites are in Europe and North America. This filtering 278 along with the quality control measures described earlier reduced the number of employed sites in this 279 study from 366 sites to 260 sites (Fig. S1). Furthermore, we exclude 6 sites from the weighting, located 280 on flooded land area, wetlands or intensively irrigated land. As a result of this, the constraining 281 observational dataset used to derive DOLCE V2 includes 254 sites with a total of 13641 monthly records. 282

#### 283 2.2. Methods

#### 284 2.2.1 Weighting approach

The weighting technique is the same as that used in DOLCE V1 and was originally presented by Bishop and 285 Abramowitz (2013) and implemented for merging observational estimates by Hobeichi et al. (2018, 2019, 286 287 2020a). It consists of building a linear combination,  $\mu$ , of the parent datasets that minimise  $\sum_{j=1}^{J} (\mu^j - y^j)^2$ , where  $j \in [1, J]$  are the monthly time-site records,  $y^j$  is the observed ET at the  $j^{\text{th}}$  time-288 site record. The linear combination  $\mu^{j} = \sum_{k=1}^{K} w_{k} x_{k}^{j}$  is subject to the constraint that  $\sum_{k=1}^{K} w_{k} = 1$ , 289 where  $k \in [1, K]$  represents the parent datasets and  $x_k^j$  is the value of the  $k^{th}$  bias-corrected parent 290 dataset (i.e. after subtracting its mean bias relative to the all-site observational dataset) corresponding to 291 292 the  $i^{th}$  time-site record. The analytical solution to this problem accounts for both the performance 293 differences between the parent datasets and their error covariance (Fig. S2), a proxy for dependence. Further details on the merging technique can be found in Abramowitz and Bishop (2015) and Bishop and 294 295 Abramowitz (2013). The weighting approach is used to combine the global parent datasets separately on 296 different spatiotemporal subsets of the entire period and globe, using a tiered approach detailed in 297 section 2.2.3.

298

#### 299 2.2.2 Computing uncertainty in ET

300 The ensemble dependence transformation process developed by Bishop and Abramowitz (2013) is used 301 to calculate the spatiotemporal uncertainty of DOLCE V2 and DOLCE V3. The process transforms the 302 global parent datasets to a new ensemble so that the variance of the transformed ensemble about the derived hybrid ET estimate,  $\mu$ , is constrained to be equal to the error variance of  $\mu$  with respect to the 303 flux tower data, averaged over time and space (i.e. across all *J* records). We use the spread  $\sqrt{\sigma^2}$  of the 304 transformed ensemble as the spatially and temporally varying estimate of uncertainty standard 305 306 deviation, which we will refer to as uncertainty. We refer the reader to Bishop and Abramowitz (2013) 307 for the derivation of this approach, and Hobeichi et al (2018) for its implementation in this context. The spread  $\sqrt{\sigma^2}$  of the transformed ensemble accurately reflects the uncertainty of  $\mu$  in those grid-cells 308 where flux tower observations are available. This process ensures that the computed uncertainty 309 310 provides a better uncertainty estimate of the hybrid ET than simply using the spread of the parent 311 datasets. One additional advantage of defining uncertainty in this way is that it should give an accurate upper 312 313 bound estimate of the likely discrepancy between the product and unseen ET measurements at a range 314 of spatial scales. That is, since it is based on the discrepancy of the final hybrid product and point-based 315 flux tower estimates, which are essentially at the extremes of spatial discrepancy, the discrepancy between DOLCE and actual ET at any spatial scale greater than that of a tower footprint and smaller 316 317 than that of DOLCE should be less than this uncertainty estimate (noting however that this is the 318 estimated standard deviation of uncertainty, rather than a hard upper limit). In Section 2.2.5 below, we 319 detail the out-of-sample testing of this uncertainty estimate at the point scale. 320

#### 321 2.2.3 Tiering of data set subsets in time and space to maximise coverage

322 To derive DOLCE V1 over the global land, we applied spatial tiering (using different subsets of parent 323 products in different regions to maximise spatial coverage). We now expand this approach to include 324 temporal tiering to improve the temporal reach of DOLCE. Collectively, the incorporated parent datasets 325 have a temporal cover over 1980 – 2018, but only a short common overlap during 2003-2007 in DOLCE 326 V2, and during 2003 – 2016 in DOLCE V3, and their spatial intersection does not cover the global land. 327 Therefore, to achieve a global land coverage from 1980 through 2018 without excluding any of their parent products, it was necessary to build DOLCE V2 and DOLCE V3 from different subsets of parent 328 329 datasets in time periods and land regions depending on the availability of the parent datasets as shown 330 in Table 1. To this end, we consider 14 and 4 distinct temporal tiers in DOLCE V2 and DOLCE V3 331 respectively. For example, in DOLCE V2, tier 9 covers 2008 - 2012 and incorporates all datasets except 332 SRB-GEWEX. Tier 1 incorporates the least parent datasets, for the year 1980 (i.e. FLUXCOM-MET and 333 GLEAM3.3A), while tier 8 uses all the parent datasets and covers 2003 – 2007. Furthermore, within each 334 temporal tier, we consider three spatial sub-tiers, with each spatial sub-tier covering a part of the land. 335 These consist of (a) all land except Antarctica, Greenland and North Africa, (b) only Antarctica and 336 Greenland, (c) only North Africa. A similar spatial tiering approach was also applied in DOLCE V1. Other 337 spatial tiers, each consisting of a small number of grid cells were also considered where necessary to ensure that no grid cell in DOLCE V2 or DOLCE V3 is missing ET data if a single parent is missing ET data 338 339 for that grid cell. As a result of the tiering approach, weighting is computed separately using a different 340 subset of parent data sets and site data in each tier, resulting in distinct spatiotemporal subsets of the 341 entire period. Collectively, the hybrid estimates developed throughout the temporal tiers and their spatial 342 sub-tiers form DOLCE V2 and DOLCE V3 over the global land throughout 1980 – 2018. The reduced number 343 of temporal tiers in DOLCE V3 is to ensure that no temporal discontinuities occur throughout the covered 344 period, which otherwise would have reduced the suitability of DOLCE V3 for trend analysis. In comparison, 345 the incorporation of a larger ensemble of parent products in DOLCE V2 is to derive an optimal ET product 346 that minimises discrepancy with in-situ observations.

347

#### 348 2.2.4 Weighting groups

349 Previous studies have found that the performance of a global product can vary with different climatic 350 circumstances, suggesting that separating the weighting into separate regions or other groupings might 351 well improve the results of the weighting overall (Ershadi et al., 2014; Hobeichi et al., 2018; Michel et al., 2016). Grouped weighting simply involves dividing the time and/or space covered by a particular tier 352 353 into different subsets or groups (e.g., with different climatic conditions), and then applying the 354 weighting technique separately for each group (within a single tier). We expect that grouped weighting 355 has the potential to improve weighting by accounting for the variation in performance of the parent 356 datasets over different climate or land conditions and can hopefully improve biases detected in DOLCE 357 V1. Hobeichi et al. (2018) tried to group flux tower sites based on their land cover type and computed 358 weights for each land cover type. However, this approach did not improve the results, whether grouping 359 by climate zone or aridity index, with the main reason being attributed to the small number of sites in 360 many groups. Despite the availability of 100 additional sites to constrain the weighting here compared

to Hobeichi et al., (2018), the ratio of the observational data to the number of parents has not improved
 across several climate or land cover types for DOLCE V2. We therefore investigate new approaches to
 grouped weighting that allow sufficiently low group numbers to keep a reasonable sample size in each
 of them, including:

- Grouping by latitudinal zone: this is a simplification of grouping by climate type in which
   climates are aggregated into three latitudinal zones: (i) high latitudes (±60° poleward), (ii) mid latitudes (±60° towards the subtropics ±40°), and (iii) tropics and sub-tropics (between -40° and
   40°). In each zone we apply a separate weighting using the corresponding group of sites.
- Grouping by continents: Sites are naturally separated by continental boundaries and we might suspect that a particular ET product performs differently across continents. For instance, precipitation is involved in the derivation of many of the parent datasets, and has been found to have different fidelity over different continents (Hobeichi et al., 2020b).
- Grouping by hemisphere: Pan et al. (2020) found that ET estimates agree more in the Northern
   hemisphere than in the Southern hemisphere. Therefore, performing separate weighting in each
   hemisphere could be better than weighting across all global land.
- 376 Grouping by seasons: Several studies have shown that the skill of ET datasets vary by seasons • 377 (Jiménez et al., 2018; Long et al., 2014; Mueller et al., 2011). To capture these differences, we 378 implement grouping by seasons, and grouping by month (detailed below). We consider two 379 combined seasons i.e., summer-fall and winter-spring. In the summer-fall season, we constrain 380 the weighting with (1) monthly observations from sites located in the Northern hemisphere 381 during the period June–November, and (2) monthly observations from sites located in the 382 Southern hemispheres during the period December–May. The remaining observational data is 383 used to constrain the weighting during the winter-spring combined season.
- Grouping by months: This is similar to grouping by seasons, the only difference is that the two
   groups are June–November and December–May, without accounting for the different seasonal
   phase between hemispheres.
- 387 Grouping by ET regime and months: Land was classified into three distinct broad ET regimes (Fig. • 388 S4) according to two aspects of ET, mean annual total ET and within-year relative variability 389 throughout 1980 – 2018, derived from GLEAM V3.5a, and using K-means unsupervised 390 classification (MacQueen, 1967). We explain the classification method further in section 3.5.2. 391 Different sets of weights were computed at each ET regime during June–November and 392 December–May. Implementing weighting this way ensured that we account for performance 393 differences across different physical aspects of the land and seasons. Despite that observational 394 data was divided into six distinct groups, the observational data available in each group was still 395 appropriate to merge the four parent datasets of DOLCE V3. However, we found this grouped 396 weighting strategy not appropriate for merging 11 parent datasets of DOLCE V2.
- 397

As an alternative to the grouping strategies, we also investigate if deriving a spatially varying bias
 correction within each tier could further improve the weighting. We describe the examined bias
 correction approaches and their effectiveness in the Supplementary Material.

#### 402 2.2.5 Out-of-sample testing approach

403

404 To test the effectiveness of different weighting groups or bias-correction approaches, and assess which 405 strategy offers the best performance, we use out-of-sample tests. To do this, we first divide the flux 406 tower sites between the in-sample and out-of-sample groups by randomly selecting 25% of the sites as 407 out of sample. The remaining sites form the in-sample training set are used to compute bias correction 408 terms and weights for the parent datasets in each tier using the weighting technique without weighting 409 groups (as adopted in DOLCE V1), and with each of the groups and bias correction strategies detailed in 410 section 2.2.5 and S4 (supplementary material). In each case, these bias correction terms and weights are 411 then applied to the parent datasets and compared to the out-of-sample sites to test efficacy of the 412 clustering or bias correction approach employed. The process is repeated for each grouping or bias 413 correction strategy to derive several hybrid ET datasets for each out of sample group of sites. 414 415 For each strategy, the test was repeated 1000 times with a different random selection of sites being out 416 of sample. The performance of each hybrid ET estimate was evaluated across five statistical metrics.

417 These were root mean squared error (RMSE), absolute standard deviation difference  $|\sigma_{dataset} - \sigma_{observation}|$ , correlation, mean absolute deviation (i.e. mean(|dataset - observation|)) and median

419 absolute deviation (i.e. median(|dataset – observation|)). DOLCE V1 has not been included in this test

420 because its coarser spatial resolution (i.e. 0.5°) excludes many coastal sites and so significantly reduces

421 the observational data we could use in this analysis. The out-of-sample test is carried out over the

422 common period of availability of all the parent datasets i.e. 2003 – 2007 and 2003 – 2016 to enable

423 comparison of the out-of-sample performance of each approach with all of the 11 and 4 parent datasets

- 424 of DOLCE V2 and DOLCE V3 respectively.
- 425

426 We perform another out-of-sample experiment to test if the uncertainty estimate derived by the 427 successful grouping/bias correction strategy performs well out of sample. In this test, we first select a 428 site S, but instead of constraining the weighting using observed ET from this site, we compute the weights and bias correction terms of the parent datasets by using all the sites except S (i.e. just one site 429 430 is out of sample). We then calculate the MSE of the derived hybrid ET against observations from all the sites except S. We denote this value by uncertainty<sub>in-sample</sub>, since it represents the uncertainty 431 432 estimate computed using the same observational dataset that we used to train the weighting. We also 433 calculate the MSE of the hybrid ET against the out-of-sample observations from S, and we denote this as uncertainty<sub>out-sample</sub>, since we perform the comparison against ET observations that have not been 434 used to train the weighting. We repeat this test for all the sites, and each time we calculate the ratio 435 uncertainty<sub>in-sample</sub>. In an ideal case, this ratio should equal to unity. 436 uncertainty<sub>out-sample</sub>, 437

- 439 **3. Results and Discussion**
- 440

# 441 3.1 Out-of-Sample Performance of DOLCE V2 and DOLCE V3

We derive DOLCE V2.1 (Hobeichi, 2020) from 11 parent datasets by applying a grouped weighting by
months. As detailed in the Supplementary material (section S5), this approach achieves slightly better

- out-of-sample performance than the other grouped weighting approaches in estimating ET (Fig. S6) and
- in deriving more robust uncertainty estimates (Fig. S7). We recall that in grouped weighting by months,
- the observational and gridded ET data are split into two groups, one covering the period June –
- 447 November and the other covering December May. Weighting and bias correction is then implemented
- in each group separately for each tier to create the subsets from which the hybrid ET product is derived.
- 449 We derive DOLCE V3 (Hobeichi, 2021) from 4 parent datasets by applying a grouped weighting by ET
- 450 regimes and months. Both DOLCE V3 and DOLCE V2.1 outperform their parent datasets in the out-of-
- 451 sample tests across all performance metrics (Fig. S6 and Fig. S8). DOLCE V2.1 performs better than
- 452 DOLCE V3 across all performance metrics except Standard deviation difference as illustrated in Fig. S8.
- 453 The overall better performance of DOLCE V2.1 is expected given that more ET estimates contributing to
- 454 the weighting. On the other hand, DOLCE V2.1 has proven worse performance than DOLCE V3 in
- 455 capturing variation in ET observations since variability in ET should have decreased when the variations
- in individual products are not temporally coincident.
- 457

# 458 3.2 Comparison of DOLCE V2 and DOLCE V3 with their parent datasets

459 Figure 1 displays the latitudinal means of each of DOLCE V2 and DOLCE V3 and their parent datasets 460 computed over a common spatial mask and common periods of 2003 – 2007 and 2003 – 2016 in the 461 case of DOLCE V2 and DOLCE V3 respectively. The grey ribbon represents the uncertainty of DOLCE V2 462 and DOLCE V3 in Fig. 1a and Fig. 1b respectively, defined by the ± uncertainty standard deviation 463 interval. The uncertainty standard deviation of the two DOLCE products mostly contain the latitudinal 464 variations of their parent datasets with the exception of FLUXCOM-RS which exhibits larger ET over the 465 tropics and subtropics of the southern hemisphere relative to DOLCE V2 (Fig. 1b). This containment 466 should not be surprising since uncertainty estimates should be robust for point-scale estimates. Figure 467 1a shows that DOLCE V1 exhibits a slightly lower ET than DOLCE V2 in the tropics and sub-tropics. DOLCE 468 V2 appears in the lower end of the range of the other datasets from 60° poleward. All the datasets 469 exhibit considerable disparities over the mid-latitude south of -50°, where the contribution of the 470 terrestrial ET comes mostly from the lower Andes. The difference between DOLCE V2 and DOLCE V3 is 471 smallest over the mid latitudes of the Northern hemisphere where most of the flux tower sites are 472 located, and is largest over the tropics where very few observations are available. Also, both the number 473 and the spread of parent datasets is larger in DOLCE V2 which explains its larger uncertainty compared 474 to DOLCE V3. The parent datasets of DOLCE V3 are in general in the upper range of ET across all the 475 different participating products, which also explains why DOLCE V3 exhibits larger ET than DOLCE V2 476 throughout the land and mostly over the tropics.

478 Figure 2 shows the spatial distribution of differences in the ET mean between DOLCE V2 and each of its
479 parent datasets. We apply different spatiotemporal masks for each comparison based on parent dataset

- 480 coverage (Table 1). We also compute the climatological difference of DOLCE V2 with its predecessor
- 481 DOLCE V1 over 2000 2009. A similar plot showing the spatial distribution of differences in the ET mean
- 482 between DOLCE V3 and each of its parent datasets is provided in Fig. S9. Fig. S9 shows that DOLCE V3
- 483 exhibits higher ET than DOLCE V2.1 and DOLCE V1 over most of the land, particularly over the tropics
- and the high latitudes. On the other hand, the climatological difference between DOLCE V3 and its
   parent datasets show different spatial patterns, and the least climatological difference is between
- 486 DOLCE V3 and GLEAM V3.5B.
- 487
- 488 Over the temperate regions of the northern hemisphere, DOLCE V2 exhibits lower mean ET than all its
- parents except SEBS. We have computed the mean bias of all these datasets relative to the
- 490 observational data available from sites located in these temperate latitudes. DOLCE V2 has a negligible
- 491 bias of 0.2  $W m^{-2}$  relative to the observational data. This bias results from a positive bias of 0.4  $W m^{-2}$
- 492 during June November and a negative bias of -0.2  $W m^{-2}$  during December–May. All the parent
- 493 datasets except SEBS exhibit a positive bias that ranges between 2.7 and 11.4  $W m^{-2}$  and SEBS has a
- 494 negative mean bias of -3.4  $W m^{-2}$ , that varies between -0.2  $Wm^{-2}$  during December May and -6.3 495  $Wm^{-2}$  during June – November. We note that the bias relative to the in-situ observational datasets is 496 only indicative of the performance of the gridded datasets at the sites and do not necessarily represent 497 the actual mean bias over these regions. The discrepancy between DOLCE V2 and DOLCE V1 is relatively
- 498 small across all land.
- 499

500 Large differences between DOLCE V2 and FLUXCOM-RS are seen over the Congo and the Amazon basins, 501 southern Africa, and the Brazilian highlands. The mean climatological bias of FLUXCOM-RS relative to observational data from these regions is 30  $W m^{-2}$ . This large bias likely results from the lack of 502 503 sufficient data available to train the machine learning algorithm over climatically distinct biomes, which made ET prediction less constrained. This bias did not appear in FLUXCOM-MET possibly because ET 504 prediction is based on a larger set of predictor variables. DOLCE V2 exhibits a relatively small bias 505 ranging between 2.6  $W m^{-2}$  during June–November and 6.4  $W m^{-2}$  during December–May. In 506 comparison, DOLCE V3 exhibits no significant bias during June–November and a bias of 12.2  $W~m^{-2}$ 507 508 during December–May which is similar to the bias in GLEAM V3.5B over these latitudes and seasons, 509 and is less than the bias in the remaining parent datasets (i.e. GLEAM V3.5A, FLUXCOM-MET, and ERA5 510 In general, there are apparent disparities in the patterns of climatological differences in the tropics 511 across all the maps. This results from the fact that global ET datasets exhibit large differences over the tropics which has been highlighted previously (Paca et al., 2019; Pan et al., 2020), particularly over the 512 513 Amazon basin.

514

# 3.3 Comparison of basin and continental ET with existing literature

516 We now compare DOLCE V2 and DOLCE V3 with annual mean ET aggregates over a range of river basins 517 documented in a recent study (Table 4 of Zhang et al., 2018). ET in this study - which we'll refer to as 518 CDR-ET- is derived by merging 10 available ET datasets into a hybrid ET which then receives corrections, 519 so that the surface water budget - established by derived hybrid estimates of the other hydrological 520 variables - is closed. Table 2 displays the mean annual ET aggregates in  $mm \ year^{-1}$  across 20 river 521 basins calculated for DOLCE V2, DOLCE V3 and CDR-ET over the common period 1984 – 2010. Our 522 results show that there is an overall agreement between DOLCE V2 and CDR-ET across all the non-523 Siberian rivers where the difference in ET estimates is mostly around 10%. The agreement worsens over 524 the Arctic basins Indigirka, Kolyma, Lena, Northern Dvina, Yenisei and particularly over Olenik and 525 Pechora where the differences in ET estimates exceed 20%. Previous studies have reported large 526 uncertainties in the water fluxes over the Siberian basins (Lorenz et al., 2015) most likely due to the 527 absence of a proper representation of snow and permafrost dynamics (Candogan Yossef et al., 2012). 528 Interestingly, over the north American arctic basins Mackenzie and Yukon, DOLCE V2 and CDR-ET exhibit 529 much smaller relative differences than at their Siberian counterparts. DOLCE V3 exhibits higher ET than 530 DOLCE V2 and CDR-ET across the majority of the river basins, and particularly over the Arctic basins. 531 DOLCE V3 is within the range of its recently developed parent datasets which exhibit higher ET than the 532 old generation products such as SRB-GEWEX and SEBS incorporated in DOLCE V2.

533

534 We also compare DOLCE V2 and DOLCE V3 with continental annual means of ET shown by L'Ecuyer et al. 535 (2015). In their study, they derive a hybrid ET by merging three global datasets. Then, they adjust the 536 hybrid ET and its associated uncertainty by enforcing the physical constraints of the surface and 537 atmospheric water and energy budgets using a data assimilation technique (DAT). Our results show that 538 DOLCE V2 has smaller ET with larger associated uncertainties compared to those derived in L'Ecuyer et 539 al. (2015) (Table 3). The range of their ET estimate overlaps with the range of DOLCE V2 and DOLCE V3 throughout all continents. In L'Ecuyer et al. (2015), the uncertainty estimates are originally taken from 540 541 the literature and are deemed constant across time and space, then these are reduced by the DAT. The 542 uncertainty estimate of DOLCE, however, is firmly grounded in the discrepancy between the gridded 543 DOLCE product and in-situ tower data. The variance of this discrepancy is used to recalibrate the 544 variance of the parent datasets, which are then used to estimate uncertainty, allowing spatiotemporally 545 varying uncertainty estimate that is both consistent with the discrepancy between DOLCE and surface 546 observations while at the same time being spatially and temporally complete. This process is detailed by 547 Hobeichi et al (2018).

548

549 Finally, we compare DOLCE V2 with the ET component of Conserving Land Atmosphere Synthesis Suite 550 (CLASS; Hobeichi, 2019; Hobeichi et al., 2020a) which we denote as CLASS-ET. CLASS dataset comprises 551 coherent estimates of the surface water and energy budgets at the gridded monthly scale. CLASS-ET has 552 been derived by adjusting DOLCE V1 by enforcing the simultaneous closure of the surface water and 553 energy budgets using the same DAT as in L'Ecuyer et al. (2015), and can be therefore considered an 554 improved version of DOLCE V1. Table S3 displays the continental area weighted averages of DOLCE V2, 555 DOLCE V1 and CLASS-ET and the mean differences DOLCE V2 – DOLCE V1 and DOLCE V2 – CLASS 556 computed over a common time period 2003-2009, and using a common spatial mask. We find that, in 557 general, DOLCE-V2 is closer to CLASS-ET (i.e. the improved version of DOLCE V1), than DOLCE V1. 558

# 3.4 Performance of DOLCE V2 at flux sites

560 We now compare DOLCE V2 with ET measured at the 260 sites used in this study (Table S1). We display 561 two performance metrics - correlation and standard deviation - on a Taylor Diagram (Fig. 3). All data has

- been normalised before computing the statistical metrics so that the observational data at each site has
- a mean of zero and a standard deviation of 1. Each coloured point summarises the performance
- statistics of DOLCE V2 at a single site. The observational data is represented by a single "reference"
- point i.e., the hollow point at one on the horizontal axis. The plot in Fig. 3 shows that most of the
- coloured points lie close to the reference point, indicating that DOLCE V2 is highly correlated with most
- of the observational data. Overall, Fig. 3 shows good agreement with the observational datasets. Poor
- 568 performance is seen over a small number of sites. These are represented by points located outside the
- 569 Taylor diagram area. Most of these sites have less than one year of monthly records with several gaps,
- 570 perhaps raising questions about observational quality.
- 571 In a further analysis, we investigate whether the performance of DOLCE V2 is reduced over a particular 572 land cover type. For this purpose, we repeat Fig. 3, but this time we colour-code the statistics points by
- 573 the land cover type of the sites they represent as shown in Fig. S10. The new plot does not reveal clear
- 574 links between the performance of DOLCE V2 and the biome types of the sites. Similarly, we could not
- 575 find performance links with the degree of representativeness of the site to the underlying grid-cell. This
- 576 is shown in Fig. S11 where colours represent the degree of agreement between the land cover type at
- the footprint of the tower site and the dominant land cover of the grid-cell containing the site. As shown
- 578 in Fig. S11, we carry out this analysis on the basis of three levels of agreement. These include blue points 579 representing sites whose land types match the dominant land types of the underlying grid-cells; green
- 580 points representing sites whose land types cover more than 25% of the underlying grid-cells without 581 being the dominant land cover at these grid-cells; and pink points representing sites whose land types
- 561 being the dominant land cover at these grid-cells; and pink points representi
   582 covers less than 25% of the underlying grid-cells.
- 583
- 584

# 585 3.5 Changes in ET since 1980

## 586 3.5.1 Annual ET trends over the global land

587 We use DOLCE V3 to produce a long-term (1980 – 2018) map of trends in annual ET totals (Fig. 4) as 588 proposed by Mann-Kendall (Kendall, 1948; Mann, 1945) using the Sen's slope method (Sen, 1968). We 589 use the uncertainty estimates associated with the ET fields and the confidence interval of the slope as 590 two confidence measures to filter out spurious trends. These confidence measures consider trends' 591 behaviour as reliable only if (i) the confidence interval of the slope does not encompasses a mix of 592 negative and positive values; and (ii) trends' slopes computed for multiple different random samples of 593 ET within the interval ET ± uncertainty standard eviation agree in sign at least 90% of the time.

594 Unreliable trends occur in regions where ET uncertainty is relatively high, such as in north Africa and 595 Sahel, and in the high latitudes where ET observations are sparse or do not exist. Inconsistent trend 596 behaviour (CI includes positive and negative values) is found in regions that experienced long phases of 597 droughts and non-droughts during 1980-2018, mainly in Australia, or a succession of drought and wet 598 events, mainly in southern United States and the Amazons basin (Marengo et al., 2018). As a result of 599 this, a general long trend in ET is not identified in these regions. Miralles et al. (2014) report that these 600 changes in ET over these regions reflect El-Niño-La-Niña cycle. Similarly, we have not detected clear 601 long trends in southern South America and eastern and southern Africa. This partially agrees with the 602 study of Pan et al. (2020) where their figure 8 shows no ET trend in eastern Africa, and no agreement on the sign of trend between the participating datasets has been found in southern South America. Figure 4

604 indicates that ET has increased over most of the northern latitudes which has been highlighted in many

- studies (e.g. Miralles et al. 2014; Pan et al. 2020; Zhang et al. 2016), and declined in western United
- 506 States, central Africa and South Amercia. Unfortunately, given the absence of adequate in-situ
- 607 observations that cover a long enough period to establish trends analysis, it is difficult to validate the
- 608 identified trends directly.

In further analysis, we verify that the spatiotemporal tiering adopted in DOLCE V3 has not resulted in
 temporal discontinuities. Figure 5 illustrates the annual average line plot of the area weighted mean of

611 continental ET exhibited by DOLCE V3. The vertical dashed lines mark the beginning of a new tier, i.e. in

612 1981, 2003 and 2017. While the line plot does shows some marked changes, these do not coincide with

613 changes in tiers, and rather coincide with extreme events, and are specific to the continents where

614 these events occurred. For instance, in Australia, ET shows high mean annual total in three very wet

615 years 2000, 2010 and 2011, and low levels throughout 2001 – 2009 during the millennium drought.

- Additionally, the decline in ET since 2017 is caused by severe droughts that developed across most of
- 617 Australia.
- 618

## 619 3.5.2 ET regimes

620 To understand changes in ET across wet and dry regions, we classify land into 6 distinct dry and wet ET 621 regimes according to two aspects of ET: annual averages and within-year relative variability derived from DOLCE V3. We apply K-means clustering (MacQueen, 1967) - an unsupervised machine learning 622 623 algorithm known for its outstanding efficiency in clustering data – by implementing the K-Means 624 function and the least squares quantisation method (Lloyd, 1982) using R software. K-Means identifies K 625 centroids (i.e. imaginary values representing the centre of the clusters) and assigns each data point to 626 the cluster of the nearest centroid using – in this paper - the least squares quantisation method. For 627 each grid cell, we compute 1) the average of the annual total ET across 39 years (1980-2018); and 2) 628 within-year relative variability climatology by temporally averaging the relative standard deviation of 629 monthly ET calculated over a year and across all years. These have been used as input features for the 630 unsupervised classification. After trial and error, we find that the global land can be adequately classified into six distinct regimes that include three dry and three wet regimes. According to centroids 631 632 values (Table S4), we label the six regimes from driest to wettest and we list the proportion of the land 633 covered by each regime : (i) very low ET with high variability (16%), (ii) low ET with high variability (34%), 634 (iii) mild low ET with medium variability (22%), (iv) mild high ET with medium variability (13%), (v) high 635 ET with low variability (8%), and (vi) very high ET with low variability (7%). Figure 6 displays the spatial

- 636 distribution of the 6 ET regimes.
- 637 We compare the derived ET regimes map with the modified Köppen climate (KC) classification map by 638 Chen and Chen (2013). We find that each KC class overlaps with only one ET regime with only two 639 exceptions (Table 4): i) Land characterised by a 'Dry Steppe Hot arid' (coded BSh in KC) climate belongs 640 the 'Mild low ET with medium variablity regime', but in two regions, the Indian Deccan plateau and 641 Argentinean Gran Chaco low forests, where the climate is BSh, the ET regime is 'Mild high ET with 642 medium variability'; ii) Regions with a 'Mild temperate Fully humid Hot summer' climate (coded Cfa in 643 KC) overlaps with the 'Mild high ET with medium variability' regime in coastal regions, and to the 'Very

high ET with low variability' regime in inland regions. These two KC classes (i.e. BSh and Cfa) are shown

- 645 in bold in Table 4. Overall, ET-regimes defined in this paper provide an efficient way to aggregate the KC
- classes in less varied classes. This is not surprising knowing that KC classes are developed based on the
- 647 empirical relationship between climate and vegetation, and that ET links the water, energy (climate) and
- 648 carbon (vegetation) budgets.
- 649
- 650 3.5.3 Global annual trends across the ET regimes

651 We now explore annual trends in mean ET exhibited in each ET regime during 1980-2018. First, we 652 calculate the annual ET total climatology and ET relative variability climatology spatially averaged across 653 each regime separately, then we compute the trends in yearly ET as above (i.e. using Mann-Kendall and 654 the Sen's slope methods). Figure 7 illustrates trends' results for the dry regimes (V.L.ET, H.variability, 655 L.ET, H.variability and M.L.ET, M.Variability) and the wet regimes (M.H.ET, M.variability, H.ET, 656 L.variability and V.H.ET, L.variability). Across all regimes except the wettest one, trends in yearly ET total 657 are upward as indicated by the positive signs of both the slopes and their complete confidence intervals. The strongest trends occur in the 'M.H.ET, M.variability' regime at a rate 0.6 mm year<sup>-1</sup>, while the 658 659 slowest trend occurs in the 'V.L.ET H.variability' regime where ET is in general low. In the wettest ET 660 regime 'V.H. ET, L. variability', while the slope of the trend is positive, its confidence interval contains 661 mixed positive and negative values. This suggest that the tendency for increasing ET in the wettest ET 662 regime is not robust. Our results indicate that decreasing ET trends observed in some regions oppose 663 the consistent positive trends across the majority of ET clusters.

664 We repeat the same analysis for all the participating parent datasets that span at least 30 years. Sen's 665 slope of the trends over the period 1982 – 2012 and their confidence interval (computed at the 95% 666 confidence level) are presented in Table 5. As noted earlier, trends' behaviour is deemed inconclusive 667 when the CI encompasses negative and positive values. These are presented with regular (as opposed to 668 bold) typeface and are exhibited by FLUXCOM-MET in all regimes except the driest. In contrast, PLSH 669 shows reliable upward trends in all regimes. ERA5-land shows downward trends in the 'M.H.ET, 670 M.variability' and 'H.ET, L.variability' regimes. Both GLEAM 3.5A and DOLCE V3 show reliable upward ET 671 trends in the two middle regimes. Differences exist in the magnitude of trends across the majority the 672 products and the regimes. In DOCLE V3, the strongest trend occur in the 'M.H.ET, M.variability' regime 673 at a rate 0.56  $mm y ear^{-1}$ . Finally, the slopes of DOLCE V3 trends are within the range of slopes of 674 trends in available ET products.

675

676 There are of course some notable limitations to the approach we have taken here, some of which were 677 previously discussed in Hobeichi et al. (2018). First, the weighting approach adopted here relies heavily 678 on flux tower observations, which can suffer from a range of technical issues (Burba and Anderson, 679 2010; Fratini et al., 2019), as well as temporal gaps during particular weather conditions such as 680 extremes (Van Der Horst et al., 2019), which can affect our results. Next, unresolved land surface 681 processes in the parent datasets due for example to the absence of a proper representation of snow and 682 permafrost dynamics, or the heterogeneity of the land surface are likely to lead to uncertain ET 683 estimation in DOLCE V2 and DOLCE V3, since each of these is only a combination of its parent datasets. 684 This applies particularly in regions where observations are scarce or do not exist.

# 687 **4.** Conclusions

This work presents two new hybrid ET datasets DOLCE V2.1 and DOLCE V3. The new datasets are the 688 689 result of several key improvements over their predecessor, incorporating more parent products in 690 DOLCE V2.1, more in-situ data, testing a range of alternative implementations of its weighting and bias 691 correction approach, increased spatial resolution, and covering a longer time period. The incorporation 692 of a large ensemble of parent datasets in DOLCE V2.1 allowed us to derive a more optimal ET product 693 that can be used to benchmark global ET estimates. In comparison, the reduced number of parent 694 datasets in DOLCE V3 minimised temporal tiering and ensured that no temporal discontinuities occur 695 throughout the covered period. This allowed us to examine historical trends in ET and their robustness 696 to observational uncertainty. Despite the observationally constrained approach to defining uncertainty, 697 we found robust ET trends across most areas of the land surface, enough to present a clear signal in 698 most of the ET climate regimes we examined. These trends indicate a global increase in land derived ET 699 between 1980 and 2018. This contrasts with other gridded ET products that did not incorporate the 700 same degree of observational constraint in either their mean field or uncertainty estimates, and 701 demonstrates the usefulness of this long-term hybrid ET dataset.

702

5. Data Availability

703 704

DOLCE V2.1 dataset (Hobeichi, 2020) is publicly available in NetCDF-4 format and can be freely
downloaded from the NCI data catalogue at <a href="http://dx.doi.org/10.25914/5f1664837ef06">http://dx.doi.org/10.25914/5f1664837ef06</a>.

DOLCE V3 dataset (Hobeichi, 2021) is publicly available in NetCDF-4 format and can be freely
 downloaded from the NCI data catalogue at <a href="https://doi.org/10.25914/606e9120c5ebe">https://doi.org/10.25914/606e9120c5ebe</a>.

- 6. Competing interests.
- 711 712

710

713 The authors declare that they have no competing interests.

714

# 7. Acknowledgment

715 716

717 The authors acknowledge the support of the Australian Research Council Centre of Excellence for 718 Climate Extremes (CE170100023). This research was undertaken with the assistance of resources and 719 services from the National Computational Infrastructure (NCI), which is supported by the Australian 720 Government. We thank Franklin (Pete) Robertson (NASA Marshall Space Flight Center) for his valuable 721 contribution to DOLCE V3. This work used eddy covariance data acquired and shared by the FLUXNET 722 community, including these networks: AmeriFlux, AfriFlux, AsiaFlux, CarboAfrica, CarboEuropeIP, 723 Carboltaly, CarboMont, ChinaFlux, Fluxnet-Canada, GreenGrass, ICOS, KoFlux, LBA, NECC, OzFlux-TERN, 724 TCOS-Siberia, and USCCC. The FLUXNET eddy covariance data processing and harmonization was carried

- 725 out by the ICOS Ecosystem Thematic Center, AmeriFlux Management Project and Fluxdata project of
- 726 FLUXNET, with the support of CDIAC, and the OzFlux, ChinaFlux and AsiaFlux offices. Data were also
- 727 obtained from the Atmospheric Radiation Measurement (ARM) Program sponsored by the U.S.
- 728 Department of Energy, Office of Science, Office of Biological and Environmental Research, Climate and
- 729 Environmental Sciences Division. This works used data sourced from Terrestrial Ecosystem Research
- 730 Network (TERN) infrastructure, an Australian Government NCRIS enabled project; the Oak Ridge
- 731 National Laboratory Distributed Active Archive Center (ORNL DAAC); the Land Cover project of the ESA
- 732 Climate Change Initiative. We would like to thank all the principal investigators that authorised us to
- 733 download site data from the European Fluxes Database, and all the research institutes that made
- publicly available and/or hosted the gridded ET datasets used in this study. 734
- 735

# 8. Tables

736 737

738 Table 1: Spatial and temporal coverage and original resolution of the global ET datasets (at the time of analysis) used to develop 739 DOLCE V2.1 and DOLCE V3. DOLCE V2.1 was derived from 11 datasets and 14 temporal tiers. DOLCE V3 was derived from 4 740 datasets and 4 temporal tiers i.e. (1) 1980, (2) 1981 – 2002, (3) 2003 – 2016, (4) 2017 - 2018.

	Time period	BACI	ERA5- land	FLUXCO M-MET	FLUXCO M - RS	GLEAM 3.3A	GLEAM 3.3B	MOD16	PML	PLSH	SEBS	SRB- GEWEX
DOLCE Version		V2.1	V2.1 V3	V2.1 V3 (1980- 2016)	V2.1	V2 V3 (GLEA E	2.1 MV3.5A & 3)	V2.1	V2.1	V2.1	V2.1	V2.1
Tier	Excluded Land domain	Antarctica Greenland North Africa		Antarctica Greenland North Africa	Antarctica Greenland North Africa			Antarctica Greenland North Africa	Antarctica Greenland			
	Original resolution	0.5° half hourly	0.1° hourly	$\frac{1^{\circ}}{12}$ monthly	$\frac{1^{\circ}}{12}$ monthly	0.25° monthly	0.25° monthly	0.05° monthly	0.5° monthly	$\frac{1^{\circ}}{12}$ monthly	0.05° monthly	0.1° 3- hourly
1	1980			•		•						
2	1981		•	•		•			•			
3	1982 –		•	•		•			•	•		
	1983											
4	1984 –		•	•		•			•	•		•
	1999											
5	2000		•	•		•		•	•	•		•
	1,2&3											
6	2000		•	•		•		•	•	•	•	•
	(4 – 12)											
7	2001 –	•	•	•	•	•		•	•	•	•	•
	2002											
8	2003 -	•	•	•	•	•	•	•	•	•	•	•
	2007											
9	2008 -	•	•	•	•	•	•	•	•	•	•	
10	2012	-		-				-			-	
10	2013	•	•	•	•	•	•	•		•	•	
11	2014	•	•	•	•	•	•	•			•	
12	2015		•		•	•	٠	•			٠	

13	2016 –	•		•	•		•	
	2017 (1							
	- 6)							
14	2017	•		•	•			
	(7 – 12)							
	- 2018							

742 Table 2: Mean annual ET aggregates in mm year<sup>-1</sup> across 20 river basins calculated for DOLCE V2, DOLCE V3 and CDR-ET

743 (Table 4 of Zhang et al., 2018) over a common period 1984 – 2010. CDR-ET is derived by merging 10 available ET datasets into a

hybrid ET which then receives corrections, so that the surface water budget established by derived hybrid estimates of the other
 hydrological variables is closed.

746

Basin	CDR-ET	DOLCE V2	DOLCE V3
	1984 – 2010	1984 – 2010	1984 – 2010
Amazon	1153	1167	1314
Amur	295	309	421
Columbia	331	340	436
Congo	1045	1084	1160
Danube	503	451	550
Indigirka	138	107	231
Indus	277	323	365
Kolyma	167	132	243
Lena	245	185	283
Mackenzie	241	214	333
Mississippi	577	513	555
Murray-Darling	411	419	445
Niger	401	456	427
Northern Dvina	324	232	376
Ob	323	245	357
Olenek	174	108	237
Paraná	892	854	856
Pechora	244	166	276
Yenisei	265	216	325
Yukon	175	158	261

747

Table 3: Annual continental averages of ET ( $W m^{-2}$ ) and its standard deviation uncertainty calculated for DOLCE V2, DOLCE V3

and developed in (L'Ecuyer et al., 2015) over a common period 2000 – 2009. In (L'Ecuyer et al., 2015), ET is derived by merging
 three global datasets, and then adjusted by enforcing the physical constraints of the surface and atmospheric water and energy

751 budgets.

continent	ET± uncertainty	ET± uncertainty	ET± uncertainty
	(L'Ecuyer et al., 2015)	DOLCE V2	DOLCE V3

Africa	45 ± 3	40 ± 17	39 ± 13
Australia	27 ± 3	28 ± 16	28 ± 13
Eurasia	33 ± 3	30 ± 13	34 ± 13
North	22 + 6	28 + 12	22 + 12
America	55 ± 0	20112	52 ± 12
South	77 + 4	72 + 22	76 + 10
America	77±4	/3±23	70 ± 19

753 Table 4: correspondence between ET-regimes derived here and Köppen climate classes derived in (Chen and Chen, 2013. Text in

bold fontface indicates that the Köppen climate is associated with more than one ET regime.

755

ET regimes	Köppen climate classes (Chen and Chen, 2013)
Very low ET with high	Polar (Tundra/Frost)
variability	Dry Desert (Hot/Cold) arid
Low ET with high variability	Snow Fully humid Cold summer/Cool summer
	Snow Dry summer Cool summer
	Snow Dry winter Cold summer
	Dry Steppe Cold arid
	Dry Desert Hot arid/Cold arid
	Mild temperate Dry summer Cool summer
	Mild temperate Dry summer Warm summer
Mild low ET with medium	Snow Fully humid (Hot/Warm summer)
variability	Snow Dry winter (Hot/Warm/Summer)
	Dry Steppe Hot arid
	Mild temperate Dry summer Hot summer
	Mild temperate Fully humid Warm summer
Mild high ET with medium	Dry Steppe Hot arid (observed only in the Indian Deccan plateau and
variability	Argentinean Gran Chaco low forests)
	Mild temperate Fully humid Hot summer (observed in inland regions)
	Mild temperate Dry winter (Hot/Warm summer)
	Tropical Dry summer
High ET with low variability	Mild temperate Fully humid Hot summer/Warm summer (observed in coastal
	regions)
	Tropical Dry winter
Very high ET with low	Tropical Fully humid
variability	Topical Monsoon

756

757 Table 5: Trends in yearly ET total (mm year<sup>1</sup>) spatially averaged across each ET regime calculated for DOLCE V3 and five

758 participating parent datasets available during 1982 – 2012. The text shows slopes of the trend line and their confidence interval

759 calculated at the 95% confidence level, bold text indicates that the trend is reliable since the confidence interval is strictly

760 *positive or negative.* 

Dataset and	V.L.ET,	L.ET,	M.L.ET,	M.H.ET,	H.ET,	V.H.ET,
time span	H.variability	H.variability	M.variability	M.variability	L.variability	L.variability
DOLCE V3	-0.04 [-0.23, 0.16]	0.26 [-0.11, 0.63]	0.44 [0.1, 0.76]	0.56 [0.2, 0.87]	0.07 [-0.27, 0.4]	0.34 [-0.1, 0.9]
ERA5-land	-0.18 [-0.36, 0.04]	0.02 [-0.42, 0.47]	0.14 [-0.38, 0.6]	-0.65 [-1.14, -0.22]	-0.89 [-1.28, -0.51]	0.11 [-0.2, 0.5]
FLUXCOM- MET	-0.02 [-0.04, 0]	0.04 [-0.11, 0.23]	0.05 [-0.07, 0.2]	-0.11 [-0.27, 0.04]	-0.003 [-0.18, 0.17]	0.25 [-0.04, 0.57]

GLEAM 3.5A -0	0.08 [-0.28, 0.16]	0.35 [-0.04, 0.76]	0.59 [0.34, 0.95]	0.43 [0.1, 0.77]	0.05 [-0.33, 0.44]	0.62 [0.12, 1.31]
<b>PML</b> -0	0.1 [-0.28, 0.15]	0.42 [0.11, 0.75]	1 [0.64, 1.45]	0.21 [-0.19, 0.64]	0.28 [-0.38, 0.81]	-0.32 [-1.24, 0.62]
PLSH 0.:	.17 [0.1, 0.24]	0.39 [0.16, 0.66]	1.3 [0.8, 1.77]	1.41 [0.85, 1.89]	1.53 [0.75, 2.17]	0.82 [0.36, 1.35]

# 

# 9. Figures





Figure 1: (a) Latitudinal means of DOLCE V2 and its parent datasets computed over a common period 2003–2007, and a
 common spatial mask. (b) Latitudinal means of DOLCE V3 and its parent datasets computed over a common period 2003–2016,
 and common spatial mask. The grey ribbon represents the values of DOLCE ± uncertainty. DOLCE V1 and DOLCE V2 are included
 in (a) and (b) respectively for comparison. FLUXCOM-METa and FLUXCOM-METb are two different datasets from the FLUXCOM-

*MET setup*.





774 775 776 777 Figure 2: Spatial distribution of differences in ET climatology between DOLCE V2 and each of its parent datasets and DOLCE V1. Different spatiotemporal masks are applied for each comparison based on the spatiotemporal coverage of DOLCE V2 and the

other datasets.



779 Figure 3: Taylor Diagram displaying two performance metrics i.e. correlation and standard deviation of DOLCE V2 relative to 780 normalised observational data presented by a hollow point (reference point) at one unit on the x-axis. Pink points represent 781 performance statistics scored at sites located on wetlands, flooded plain or intensively irrigated areas.



783 784

Figure 4: Spatial pattern of ET climate trends in DOLCE V3 over 1980 – 2018 derived using Mann-Kendall and Sen's slope 785 methods. Grid cells in white correspond to unreliable ET trends because (i) the confidence interval of the slope encompasses a 786 787 mix of negative and positive values; or (ii) trends' slopes computed for multiple different random samples of ET within the

- *interval ET ± uncertainty do not agree in sign.*
- 788
- 789





Figure 5: Annual average line plot of the area weighted mean of continental ET exhibited by DOLCE V3. The vertical dashed lines 792 mark the beginning of a new tier in 1981, 2003 and 2017





V.H.ET, L.variability
H.ET, L.variability
M.H.ET, M.variability
M.L.ET, M.variability
L.ET, H.variability
V.L.ET, H.variability

795 100°W 0° 100°E
Figure 6: Classification of the land into 6 distinct dry and wet ET regimes using K-means unsupervised classification based on
796 DOLCE V3 annual ET mean and within-year relative variability both computed for 1980-2018. The six ET regimes are labelled
798 from driest to wettest as very low ET with high variability (V.L.ET, H. variability), (ii) low ET with high variability (L.ET, H.
799 variability), (iii) mild low ET with medium variability (M.L.ET, M. variability), (iv) mild high ET with medium variability (M.H.ET, M. variability), (v) high ET with low variability (H.ET, L. variability), and (vi) very high ET with low variability (V.H.ET, L.
801 variability).
802



803

804

Figure 7: Trends in mean annual ET total computed for the dry and wet ET regimes during 1980-2018. Slopes and confidence
 intervals are computed using Mann-Kendall and the Sen's slope methods. The spatial distribution of the ET regimes is illustrated
 in Fig. 6.

808

809

## 810 **10.** References

- 811 Abramowitz, G. and Bishop, C. H.: Climate model dependence and the ensemble dependence
- transformation of CMIP projections, J. Clim., 28(6), 2332–2348, doi:10.1175/JCLI-D-14-00364.1, 2015.

- Abramowitz, G. ., Herger, N., Gutmann, Ethan; Hammerling, D. and Knutti, Reto; Leduc, Martin; Lorenz,
- 814 Ruth; Pincus, Robert; Schmidt, G. A.: ESD Reviews: Model dependence in multi-model climate
- ensembles: weighting, sub-selection and out-of-sample testing, Earth Syst. Dynam, 10(1), 91–105,

816 doi:10.5194/esd-10-91-2019, 2019.

- 817 Balsamo, G., Albergel, C., Beljaars, A., Boussetta, S., Brun, E., Cloke, H., Dee, D., Dutra, E., Muñoz-
- Sabater, J., Pappenberger, F., De Rosnay, P., Stockdale, T. and Vitart, F.: ERA-Interim/Land: a global land
- surface reanalysis data set, Hydrol. Earth Syst. Sci, 19, 389–407, doi:10.5194/hess-19-389-2015, 2015.
- Bishop, C. H. and Abramowitz, G.: Climate model dependence and the replicate Earth paradigm, Clim.
- 821 Dyn., 41(3–4), 885–900, doi:10.1007/s00382-012-1610-y, 2013.
- Bodesheim, P., Jung, M., Gans, F., Mahecha, M. D. and Reichstein, M.: Upscaled diurnal cycles of land-
- Atmosphere fluxes: A new global half-hourly data product, Earth Syst. Sci. Data, 10(3), 1327–1365,
  doi:10.5194/essd-10-1327-2018, 2018.
- 825 Burba, G. B. P. G. to E. C. F. M. P. and W. E. for S. and I. A. and Anderson, D.: A Brief Practical Guide to
- 826 Eddy Covariance Flux Measurements: Principles and Workflow Examples for Scientific and Industrial
- 827 Applications, LI-COR Biosciences., 2010.
- 828 Candogan Yossef, N., Van Beek, L. P. H., Kwadijk, J. C. J. and Bierkens, M. F. P.: Assessment of the
- potential forecasting skill of a global hydrological model in reproducing the occurrence of monthly flow
- 830 extremes, Hydrol. Earth Syst. Sci., 16(11), 4233–4246, doi:10.5194/hess-16-4233-2012, 2012.
- 831 Chen, D. and Chen, H. W.: Using the Köppen classification to quantify climate variation and change: An
- example for 1901–2010, Environ. Dev., 6, 69–79, 2013.
- 833 Chen, X., Su, Z., Ma, Y., Yang, K., Wen, J. and Zhang, Y.: An improvement of roughness height
- parameterization of the Surface Energy Balance System (SEBS) over the Tibetan plateau, J. Appl.
- 835 Meteorol. Climatol., 52(3), 607–622, doi:10.1175/JAMC-D-12-056.1, 2013.
- 836 Chen, X., Massman, W. J. and Su, Z.: A column canopy-air turbulent diffusion method for different
- 837 canopy structures, J. Geophys. Res. Atmos., 124(2), 488–506, 2019.
- B38 Dawdy, D. R., Lichty, R. W. and Bergmann, J. M.: A rainfall-runoff simulation model for estimation of
- 839 flood peaks for small drainage basins, US Government Printing Office., 1972.
- 840 Erfanian, A., Wang, G. and Fomenko, L.: Unprecedented drought over tropical South America in 2016:
- Significantly under-predicted by tropical SST, Sci. Rep., 7(1), doi:10.1038/s41598-017-05373-2, 2017.
- 842 Ershadi, A., McCabe, M. F., Evans, J. P., Chaney, N. W. and Wood, E. F.: Multi-site evaluation of
- terrestrial evaporation models using FLUXNET data, Agric. For. Meteorol., 187, 46–61,
- doi:10.1016/j.agrformet.2013.11.008, 2014.
- 845 Feng, F., Li, X., Yao, Y., Liang, S., Chen, J., Zhao, X., Jia, K., Pintér, K. and McCaughey, J. H.: An Empirical
- 846 Orthogonal Function-Based Algorithm for Estimating Terrestrial Latent Heat Flux from Eddy Covariance,
- 847 Meteorological and Satellite Observations, PLoS One, 11(7), e0160150,
- 848 doi:10.1371/journal.pone.0160150, 2016.
- 849 Fisher, R. A. and Koven, C. D.: Perspectives on the future of Land Surface Models and the challenges of
- 850 representing complex terrestrial systems, J. Adv. Model. Earth Syst., 12, 1–24,
- doi:10.1029/2018ms001453, 2020.
- 852 Fratini, G., Sabbatini, S., Ediger, K., Riensche, B., Burba, G., Nicolini, G., Vitale, D. and Papale, D.:
- 853 Characterization of Eddy Covariance flux errors due to data synchronization issues during data
- acquisition, in Geophysical Research Abstracts, vol. 21., 2019.
- Hamed Alemohammad, S., Fang, B., Konings, A. G., Aires, F., Green, J. K., Kolassa, J., Miralles, D., Prigent,
- 856 C. and Gentine, P.: Water, Energy, and Carbon with Artificial Neural Networks (WECANN): A statistically

- 857 based estimate of global surface turbulent fluxes and gross primary productivity using solar-induced
- 858 fluorescence, Biogeosciences, 14(18), 4101–4124, doi:10.5194/bg-14-4101-2017, 2017.
- Han, D., Wang, G., Liu, T., Xue, B.-L., Kuczera, G. and Xu, X.: Hydroclimatic response of
- 860 evapotranspiration partitioning to prolonged droughts in semiarid grassland, J. Hydrol., 563, 766–777,861 2018.
- 862 Herger, N., Abramowitz, G., Knutti, R., Angélil, O., Lehmann, K. and Sanderson, B. M.: Selecting a climate
- 863 model subset to optimise key ensemble properties, Earth Syst. Dyn., 9(1), 135–151, doi:10.5194/esd-9-
- 864 135-2018, 2018.
- Hobeichi, S.: Conserving Land-Atmosphere Synthesis Suite (CLASS) v 1.1, , doi:10.25914/5c872258dc183,
  2019.
- 867 Hobeichi, S.: Derived Optimal Linear Combination Evapotranspiration DOLCE v2.1, ,
- 868 doi:10.25914/5f1664837ef06, 2020.
- 869 Hobeichi, S., Abramowitz, G., Evans, J. and Ukkola, A.: Derived Optimal Linear Combination
- 870 Evapotranspiration (DOLCE): A global gridded synthesis et estimate, Hydrol. Earth Syst. Sci., 22(2), 1317–
- 871 1336, doi:10.5194/hess-22-1317-2018, 2018.
- Hobeichi, S., Abramowitz, G., Evans, J. and Beck, H. E.: Linear Optimal Runoff Aggregate (LORA): A global
- 873 gridded synthesis runoff product, Hydrol. Earth Syst. Sci., 23, 851–870, doi:10.5194/hess-23-851-2019,
  874 2019.
- 875 Hobeichi, S., Abramowitz, G. and Evans, J. P.: Conserving Land Atmosphere Synthesis Suite (CLASS), J.
- 876 Clim., 33, 1821–1844, doi:10.1175/JCLI-D-19-0036.1, 2020a.
- 877 Hobeichi, S., Abramowitz, G., Contractor, S. and Evans, J.: Evaluating precipitation datasets using surface
- water and energy budget closure, J. Hydrometeorol., 989–1009, doi:10.1175/jhm-d-19-0255.1, 2020b.
- 879 Van Der Horst, S. V. J., Pitman, A. J., De Kauwe, M. G., Ukkola, A., Abramowitz, G. and Isaac, P.: How
- 880 representative are FLUXNET measurements of surface fluxes during temperature extremes?,
- 881 Biogeosciences, 16(8), 1829–1844, doi:10.5194/bg-16-1829-2019, 2019.
- Jiménez, C., Martens, B., Miralles, D. M., Fisher, J. B., Beck, H. E. and Fernández-prieto, D.: Exploring the
- 883 merging of the global land evaporation WACMOS-ET products based on local tower measurements,
- 884 Hydrol. Earth Syst. Sci, 22, 4513–4533, doi:https://doi.org/10.5194/hess-22-4513-2018, 2018.
- Jung, M., Reichstein, M., Ciais, P., Seneviratne, S. I., Sheffield, J., Goulden, M. L., Bonan, G., Cescatti, A.,
- 886 Chen, J., De Jeu, R., Dolman, A. J., Eugster, W., Gerten, D., Gianelle, D., Gobron, N., Heinke, J., Kimball, J.,
- Law, B. E., Montagnani, L., Mu, Q., Mueller, B., Oleson, K., Papale, D., Richardson, A. D., Roupsard, O.,
- 888 Running, S., Tomelleri, E., Viovy, N., Weber, U., Williams, C., Wood, E., Zaehle, S. and Zhang, K.: Recent
- decline in the global land evapotranspiration trend due to limited moisture supply, Nature, 467(7318),
- 890 951–954, doi:10.1038/nature09396, 2010.
- Jung, M., Koirala, S., Weber, U., Ichii, K., Gans, F., Gustau-Camps-Valls, Papale, D., Schwalm, C.,
- 892 Tramontana, G. and Reichstein, M.: The FLUXCOM ensemble of global land-atmosphere energy fluxes, ,
- 893 6:74, 1–14, doi:10.1038/s41597-019-0076-8, 2019.
- Kanamitsu, M., Ebisuzaki, W., Woollen, J., Yang, S.-K., Hnilo, J. J., Fiorino, M. and Potter, G. L.: NCEP-DOE
- AMIP-II Reanalysis (R-2), Bull. Am. Meteorol. Soc., 83(11), 1631–1644, doi:10.1175/BAMS-83-11-1631,
- 896 2002.
- 897 Kendall, M. G.: Rank correlation methods., 1948.
- 898 L'Ecuyer, T. S., Beaudoing, H. K., Rodell, M., Olson, W., Lin, B., Kato, S., Clayson, C. A., Wood, E.,
- Sheffield, J., Adler, R., Huffman, G., Bosilovich, M., Gu, G., Robertson, F., Houser, P. R., Chambers, D.,
- 900 Famiglietti, J. S., Fetzer, E., Liu, W. T., Gao, X., Schlosser, C. A., Clark, E., Lettenmaier, D. P. and Hilburn,

- 901 K.: The observed state of the energy budget in the early twenty-first century, J. Clim., 28(21), 8319–
- 902 8346, doi:10.1175/JCLI-D-14-00556.1, 2015.
- Liang, X., Lettenmaier, D. P., Wood, E. F. and Burges, S. J.: A simple hydrologically based model of land
- surface water and energy fluxes for general circulation models, J. Geophys. Res. Atmos., 99(D7), 14415-14428, doi:10.1029/94JD00483, 1994.
- Lloyd, S.: Least squares quantization in PCM, IEEE Trans. Inf. theory, 28(2), 129–137, 1982.
- 907 Long, D., Longuevergne, L. and Scanlon, B. R.: Uncertainty in evapotranspiration fromland
- 908 surfacemodeling, remote sensing, and GRACE satellites, Water Resour. Res., 50(2), 1131–1151,
- 909 doi:10.1002/2013WR014581.Received, 2014.
- 910 Lorenz, C., Tourian, M. J., Devaraju, B., Sneeuw, N. and Kunstmann, H.: Basin-scale runoff prediction: An
- 911 Ensemble Kalman Filter framework based on global hydrometeorological data sets, Water Resour. Res.,
- 912 51, 8450–8475, doi:10.1002/2014WR016794, 2015.
- 913 MacQueen, J.: Some methods for classification and analysis of multivariate observations, in Proceedings
- of the fifth Berkeley symposium on mathematical statistics and probability, vol. 1, pp. 281–297, Oakland,
- 915 CA, USA., 1967.
- 916 Mann, H. B.: Nonparametric tests against trend, Econom. J. Econom. Soc., 245–259, 1945.
- 917 Marengo, J. A., Souza, C. M., Thonicke, K., Burton, C., Halladay, K., Betts, R. A., Alves, L. M. and Soares,
- 918 W. R.: Changes in Climate and Land Use Over the Amazon Region: Current and Future Variability and
- 919 Trends, Front. Earth Sci., 6, doi:10.3389/feart.2018.00228, 2018.
- 920 Martens, B., Miralles, D., Lievens, H., Van Der Schalie, R., De Jeu, R., Fernández-Prieto, D. and Verhoest,
- N.: GLEAM v3: updated land evaporation and root-zone soil moisture datasets, Geophys. Res. Abstr. EGU
  Gen. Assem., 18, 2016–4253, 2016.
- 923 Martens, B., Miralles, D. G., Lievens, H., Van Der Schalie, R., De Jeu, R. A. M., Fernández-Prieto, D., Beck,
- H. E., Dorigo, W. A. and Verhoest, N. E. C.: GLEAM v3: Satellite-based land evaporation and root-zone
- soil moisture, Geosci. Model Dev., 10(5), 1903–1925, doi:10.5194/gmd-10-1903-2017, 2017.
- 926 McCabe, M. F., Ershadi, A., Jimenez, C., Miralles, D. G., Michel, D. and Wood, E. F.: The GEWEX LandFlux
- 927 project: Evaluation of model evaporation using tower-based and globally gridded forcing data, Geosci.
- 928 Model Dev., 9, 283–305, doi:10.5194/gmd-9-283-2016, 2016.
- 929 Michel, D., Jiménez, C., Miralles, D. G., Jung, M., Hirschi, M., Ershadi, A., Martens, B., McCabe, M. F.,
- 930 Fisher, J. B., Mu, Q., Seneviratne, S. I., Wood, E. F. and Fernández-Prieto, D.: The WACMOS-ET project
- 931 Part 1: Tower-scale evaluation of four remote-sensing-based evapotranspiration algorithms, Hydrol.
- 932 Earth Syst. Sci., 20(2), 803–822, doi:10.5194/hess-20-803-2016, 2016.
- 933 Miralles, D. G., Holmes, T. R. H., De Jeu, R. A. M., Gash, J. H., Meesters, A. G. C. A. and Dolman, A. J.:
- 934 Global land-surface evaporation estimated from satellite-based observations, Hydrol. Earth Syst. Sci.,
- 935 15(2), 453–469, doi:10.5194/hess-15-453-2011, 2011a.
- 936 Miralles, D. G., De Jeu, R. A. M., Gash, J. H., Holmes, T. R. H. and Dolman, A. J.: Magnitude and variability
- of land evaporation and its components at the global scale, Hydrol. Earth Syst. Sci., 15(3), 967–981,
- 938 doi:10.5194/hess-15-967-2011, 2011b.
- 939 Miralles, D. G., Van Den Berg, M. J., Gash, J. H., Parinussa, R. M., De Jeu, R. A. M., Beck, H. E., Holmes, T.
- 940 R. H., Jiménez, C., Verhoest, N. E. C., Dorigo, W. A., Teuling, A. J., & Johannes Dolman, A. (2014). El Niño-
- La Niña cycle and recent trends in continental evaporation. In Nature Climate Change (Vol. 4, Issue 2, pp.
- 942 122–126). https://doi.org/10.1038/nclimate2068.
- 943 Montano, B. Q., Westerberg, I., Wetterhall, F., Hidalgo, H. G. and Halldin, S.: Characterising droughts in
- 944 Central America with uncertain hydro-meteorological data, in 2015 AGU Fall Meeting, AGU., 2015.

- 945 Mu, Q., Zhao, M. and Running, S. W.: Improvements to a MODIS global terrestrial evapotranspiration
- 946 algorithm, Remote Sens. Environ., 115(8), 1781–1800, doi:10.1016/j.rse.2011.02.019, 2011.
- 947 Mueller, B., Seneviratne, S. I., Jimenez, C., Corti, T., Hirschi, M., Balsamo, G., Ciais, P., Dirmeyer, P.,
- 948 Fisher, J. B., Guo, Z., Jung, M., Maignan, F., McCabe, M. F., Reichle, R., Reichstein, M., Rodell, M.,
- 949 Sheffield, J., Teuling, A. J., Wang, K., Wood, E. F. and Zhang, Y.: Evaluation of global observations-based
- 950 evapotranspiration datasets and IPCC AR4 simulations, Geophys. Res. Lett., 38(6), 3–10,
- 951 doi:10.1029/2010GL046230, 2011.
- 952 Mueller, B., Hirschi, M., Jimenez, C., Ciais, P., Dirmeyer, P. A., Dolman, A. J., Fisher, J. B., Jung, M.,
- 953 Ludwig, F., Maignan, F., Miralles, D. G., McCabe, M. F., Reichstein, M., Sheffield, J., Wang, K., Wood, E.
- 954 F., Zhang, Y. and Seneviratne, S. I.: Benchmark products for land evapotranspiration: LandFlux-EVAL
- 955 multi-data set synthesis, Hydrol. Earth Syst. Sci., 17, 3707–3720, doi:10.5194/hess-17-3707-2013, 2013.
- 956 Munier, S., Aires, F., Schlaffer, S., Prigent, C., Papa, F., Maisongrande, P. and Pan, M.: Combining
- 957 datasets of satellite retrieved products for basin-scale water balance study. Part II: Evaluation on the
- 958 Mississippi Basin and closure correction model, J. Geophys. Res. Atmos., 119, 12,100-12,116,
- 959 doi:10.1002/2014JD021953, 2014.
- 960 Paca, V. H. da M., Espinoza-Dávalos, G. E., Hessels, T. M., Moreira, D. M., Comair, G. F. and Bastiaanssen,
- 961 W. G. M.: The spatial variability of actual evapotranspiration across the Amazon River Basin based on
- remote sensing products validated with flux towers, Ecol. Process., 8(1), doi:10.1186/s13717-019-01588, 2019.
- 964 Pan, S., Pan, N., Tian, H., Friedlingstein, P., Sitch, S., Shi, H., Arora, V. K., Haverd, V., Jain, A. K., Kato, E.,
- Lienert, S., Lombardozzi, D., Ottle, C., Poulter, B. and Zaehle, S.: Evaluation of global terrestrial
- 966 evapotranspiration by state-of-the-art approaches in remote sensing, machine learning, and land
- 967 surface models, Hydrol. Earth Syst. Sci., 24, 1485–1509, doi:10.5194/hess-24-1485-2020, 2020.
- 968 Rodell, M., Beaudoing, H. K., L'Ecuyer, T. S., Olson, W. S., Famiglietti, J. S., Houser, P. R., Adler, R.,
- 969 Bosilovich, M. G., Clayson, C. A., Chambers, D., Clark, E., Fetzer, E. J., Gao, X., Gu, G., Hilburn, K.,
- 970 Huffman, G. J., Lettenmaier, D. P., Liu, W. T., Robertson, F. R., Schlosser, C. A., Sheffield, J. and Wood, E.
- 971 F.: The observed state of the water cycle in the early twenty-first century, J. Clim., 28, 8289–8318,
- 972 doi:10.1175/JCLI-D-14-00555.1, 2015.
- 973 Sahoo, A. K., Pan, M., Troy, T. J., Vinukollu, R. K., Sheffield, J. and Wood, E. F.: Reconciling the global
- 974 terrestrial water budget using satellite remote sensing, Remote Sens. Environ., 115, 1850–1865,
- 975 doi:10.1016/j.rse.2011.03.009, 2011.
- 976 Sen, P. K.: Estimates of the regression coefficient based on Kendall's tau, J. Am. Stat. Assoc., 63(324),
- 977 1379–1389, 1968.
- 978 Sharma, A., Wasko, C. and Lettenmaier, D. P.: If precipitation extremes are increasing, why aren't
- 979 floods?, Water Resour. Res., 54(11), 8545–8551, 2018.
- 980 Sheffield, J., Wood, E. F. and Roderick, M. L.: Little change in global drought over the past 60 years,
- 981 Nature, 491(7424), 435–438, doi:10.1038/nature11575, 2012.
- 982 Stackhouse Jr, P. W., Gupta, S. K., Cox, S. J., Zhang, T., Mikovitz, J. C. and Hinkelman, L. M.: 24.5-Year
- 983 surface radiation budget data set released, Glob. Energy Water Cycle Exp. News, 21(1), 1–20, 2011.
- 984 Stephens, G. L., Li, J., Wild, M., Clayson, C. A., Loeb, N., Kato, S., L'Ecuyer, T., Stackhouse, P. W., Lebsock,
- 985 M. and Andrews, T.: An update on Earth's energy balance in light of the latest global observations, Nat.
- 986 Geosci., 5, 691–696, doi:10.1038/ngeo1580, 2012.
- 987 Su, Z.: The Surface Energy Balance System (SEBS) for estimation of turbulent heat fluxes, Hydrol. Earth
- 988 Syst. Sci., 6(1), 85–100, doi:10.5194/hess-6-85-2002, 2002.

- Teuling, A. J.: A hot future for European droughts, Nat. Clim. Chang., 8(5), 364–365, 2018.
- 990 Teuling, A. J., Van Loon, A. F., Seneviratne, S. I., Lehner, I., Aubinet, M., Heinesch, B., Bernhofer, C.,
- 991 Grünwald, T., Prasse, H. and Spank, U.: Evapotranspiration amplifies European summer drought,
- 992 Geophys. Res. Lett., 40(10), 2071–2075, 2013.
- 993 Ukkola, A. M., Pitman, A. J., Donat, M. G., De Kauwe, M. G. and Angélil, O.: Evaluating the Contribution
- of Land-Atmosphere Coupling to Heat Extremes in CMIP5 Models, Geophys. Res. Lett., 45(17), 9003–
- 995 9012, doi:10.1029/2018GL079102, 2018.
- 996 Vinukollu, R. K., Wood, E. F., Ferguson, C. R. and Fisher, J. B.: Global estimates of evapotranspiration for
- climate studies using multi-sensor remote sensing data: Evaluation of three process-based approaches,
  Remote Sens. Environ., 115, 801–823, doi:10.1016/j.rse.2010.11.006, 2011.
- 999 Wan, Z., Zhang, K., Xue, X., Hong, Z., Hong, Y. and J. Gourley, J.: Water balance-based actual
- 1000 evapotranspiration reconstruction fromground and satellite observations over the conterminous United
- 1001 States Zhanming, Water Resour. Res., 51, 6485–6499, doi:10.1002/2015WR017311, 2015.
- 1002 Zhang, K., Kimball, J. S., Nemani, R. R. and Running, S. W.: A continuous satellite-derived global record of
- 1003 land surface evapotranspiration from 1983 to 2006, Water Resour. Res., 46(9),
- 1004 doi:10.1029/2009WR008800, 2010.
- 1005 Zhang, K., Kimball, J. S., Nemani, R. R., Running, S. W., Hong, Y., Gourley, J. J. and Yu, Z.: Vegetation
- Greening and Climate Change Promote Multidecadal Rises of Global Land Evapotranspiration, Sci. Rep.,
  5(June), 1–9, doi:10.1038/srep15956, 2015.
- 1008 Zhang, Y., Peña-Arancibia, J. L., McVicar, T. R., Chiew, F. H. S., Vaze, J., Liu, C., Lu, X., Zheng, H., Wang, Y.,
- 1009 Liu, Y. Y., Miralles, D. G. and Pan, M.: Multi-decadal trends in global terrestrial evapotranspiration and its
- 1010 components, Sci. Rep., 6, 19124, doi:10.1038/srep19124, 2016.
- 1011 Zhang, Y., Pan, M., Sheffield, J., Siemann, A. L., Fisher, C. K., Liang, M., Beck, H. E., Wanders, N.,
- 1012 Maccracken, R. F., Houser, P. R., Zhou, T., Lettenmaier, D. P., Pinker, R. T., Bytheway, J., Kummerow, C.
- 1013 D. and Wood, E. F.: A Climate Data Record (CDR) for the global terrestrial water, Earth Syst. Sci, 22(1),
- 1014 241–263, doi:10.5194/hess-22-241-2018, 2018.
- 1015
- 1016