

1 **Global assessment of how averaging over spatial heterogeneity in precipitation and potential evapotranspiration**  
2 **affects modeled evapotranspiration rates**

3

4 Elham Rouholahnejad Freund<sup>1,2</sup>, Ying Fan<sup>3</sup>, James W. Kirchner<sup>2,4,5</sup>

5

6 <sup>1</sup>Laboratory of Hydrology and Water Management, Ghent University, Ghent, Belgium

7 <sup>2</sup>Department of Environmental Systems Science, ETH Zurich, 8092, Zurich, Switzerland

8 <sup>3</sup>Department of Earth and Planetary Sciences, Rutgers University, New Brunswick, NJ, United States

9 <sup>4</sup>Swiss Federal Research Institute WSL, Birmensdorf, 8903, Switzerland

10 <sup>5</sup>Dept. of Earth and Planetary Science, University of California, Berkeley, CA 94720, United States

11

12 *Correspondence to:* Elham Rouholahnejad Freund, elham.rouholahnejad@gmail.com

13

14 **Short summary**

15 Evapotranspiration (ET) rates and the properties that regulate them are spatially heterogeneous. Averaging over  
16 spatial heterogeneity in precipitation and potential evapotranspiration as main drivers of ET may lead to biased  
17 estimates of energy and water fluxes from the land surface to the atmosphere. Here we show that this bias will be  
18 largest in mountainous terrain, in regions with temperate climates and dry summers, and in landscapes where  
19 spatial variations in precipitation and potential evapotranspiration are inversely correlated.

20

21 **Abstract**

22 Accurately estimating large-scale evapotranspiration rates is essential to understanding and predicting global  
23 change. Evapotranspiration models that are applied at continental scale typically operate on relatively large spatial  
24 grids, with the result that the heterogeneity in land surface properties and processes at smaller spatial scales  
25 cannot be explicitly represented. Averaging over this spatial heterogeneity may lead to biased estimates of energy  
26 and water fluxes. Here we estimate how averaging over spatial heterogeneity in precipitation (P) and potential  
27 evapotranspiration (PET) may affect grid-cell-averaged evapotranspiration (ET) rates, as seen from the atmosphere  
28 over heterogeneous landscapes across the globe. Our goal is to identify where, under what conditions, and at what  
29 scales this "heterogeneity bias" could be most important, but not to quantify its absolute magnitude. We use  
30 Budyko curves as simple functions that relate ET to precipitation (P) and potential evapotranspiration (PET).  
31 Because the relationships driving ET are nonlinear, averaging over sub-grid heterogeneity in P and PET will lead to  
32 biased estimates of average ET. We examine the global distribution of this bias, its scale dependence, and its  
33 sensitivity to variations in P versus PET. Our analysis shows that this heterogeneity bias is more pronounced in  
34 mountainous terrain, in landscapes where spatial variations in P and PET are inversely correlated, and in regions  
35 with temperate climates and dry summers. We also show that this heterogeneity bias increases on average, and  
36 expands over larger areas, as the grid cell size increases.

37

38

## 39 1. Introduction

40 Estimates of evapotranspiration (ET) fluxes have significant implications for future temperature predictions. Smaller  
41 ET fluxes imply greater sensible heat fluxes and, therefore, drier and warmer conditions in the context of climate  
42 change (Seneviratne et al., 2010). Surface evaporative fluxes (and thus energy partitioning over land surfaces) are  
43 nonlinear functions of available water and energy, and thus are coupled to spatially heterogeneous surface  
44 characteristics (e.g., soil type, vegetation, topography) and meteorological inputs (e.g., radiative flux, wind, and  
45 precipitation; Kalma et al., 2008; Shahraeeni and Or, 2010; Holland et al., 2013). These characteristics are spatially  
46 variable on length scales of <1 m to many kilometers. Even the highest-resolution continental-scale  
47 evapotranspiration models, such as those that are embedded in Earth System Models (ESMs), typically cannot  
48 explicitly represent the spatial heterogeneity of land surface hydrological properties at scales that are important to  
49 atmospheric fluxes. Instead, these models usually calculate grid-averaged evapotranspiration fluxes based on grid-  
50 averaged properties of the land surface (Sato et al., 1989; Koster et al., 2006; Santanello and Peters-Lidard, 2011).  
51 Thus, ET estimates that are derived from spatially-averaged land surface properties do not capture ET variations  
52 driven by the underlying surface heterogeneity (McCabe and Wood, 2006). These spatially averaged ET estimates  
53 may differ from the average of the actual spatially heterogeneous ET flux, because the relationships driving ET are  
54 nonlinear (Rouholahnejad Freund and Kirchner, 2017).

55  
56 Several studies have quantified the effects of land surface heterogeneity on potential evapotranspiration (PET) and  
57 latent heat (LH) fluxes, and have found that averaging over land surface heterogeneity can potentially bias ET  
58 estimates either positively or negatively. For example, Boone and Wetzel (1998) studied the effects of soil texture  
59 variability within each pixel in the Land-Atmosphere-Cloud Exchange (PLACE) model, which has a spatial resolution  
60 of approximately 100 by 100 km. They reported that accounting for sub-grid variability in soil texture reduced  
61 global ET by 17%, increased total runoff by 48%, and increased soil wetness by 19%, compared to using a  
62 homogenous soil texture to describe the entire grid cell. Kollet (2009) found that heterogeneity in soil hydraulic  
63 conductivity had a strong influence on evapotranspiration during the dry months of the year, but not during  
64 months with sufficient moisture availability. Hong et al. (2009) reported that aggregating radiance data from 30 m  
65 to 60, 120, 250, 500, and 1000 m resolution (input upscaling) and then calculating ET from these aggregated inputs  
66 at these grid scales using Surface Energy Balance Algorithm for Land (SEBAL, Bastiaanssen et al., 1998a) yields  
67 slightly larger ET estimates as compared to ET calculated with finer resolution inputs and then aggregated at the  
68 desired grid scales (output upscaling). The discrepancy between ET estimated with the output upscaling method  
69 and the input upscaling method grows as the size of the grid cell increases (the difference between ET calculated  
70 from the input and output upscaling methods is ~20% more at a grid scale of 1 km by 1 km compared to a grid scale  
71 of 120 m by 120 m). Aminzadeh et al. (2017) investigated the effects of averaging surface heterogeneity and soil  
72 moisture availability on potential evaporation from a heterogeneous land surface including bare soil and vegetation  
73 patches. They found that if the heterogeneity length scale is smaller than the convective atmospheric boundary  
74 layer (ABL) thickness, averaging over heterogeneous land surfaces has only a small effect on average potential

75 evaporation rates. Averaging over larger-scale heterogeneities, however, led to overestimates of potential  
76 evaporation.

77  
78 Heterogeneity biases have also been identified in ET calculation algorithms that use remote sensing data as inputs.  
79 McCabe and Wood (2006) found that remote sensing retrievals of ET are larger than the corresponding in-situ flux  
80 estimates and characterized the roles of land surface heterogeneity and remote sensing resolution in the retrieval  
81 of evaporative flux. McCabe and Wood (2006) used Landsat (60 m), Advanced Space borne Thermal Emission and  
82 Reflection Radiometer (ASTER) (90 m), and MODIS (1020 m) independently to estimate ET over the Walnut Creek  
83 watershed in Iowa. They compared these remote sensing estimates to eddy covariance flux measurements and  
84 reported that Landsat and ASTER ET estimates had a higher degree of consistency with one another and correlated  
85 better to the ground measurements ( $r=0.87$  and  $r=0.81$ , respectively) than MODIS- based ET estimates did. All three  
86 remote sensing products overestimated ET as compared to ground measurements (at 12 out of 14 tower sites).  
87 Upon aggregation of Landsat and ASTER retrievals to MODIS scale (1 km), the correlation with the ground  
88 measurements decreased to  $r=0.75$  and  $r=0.63$  for Landsat and ASTER, respectively.

89  
90 Contrary to overestimation bias, many remotely sensed ET estimates that include parameters related to  
91 aerodynamic resistance are significantly affected by heterogeneity, and underestimate ET as the scale increases  
92 (Ershadi et al., 2013). Because aerodynamic resistance is significantly affected by land surface properties (e.g.,  
93 vegetation height, roughness length, and displacement height), decreases in aerodynamic resistance at coarser  
94 resolutions could lead to smaller estimates of evapotranspiration. Ershadi et al. (2013) showed that input  
95 aggregation from 120m to 960 m in Surface Energy Balance System (SEBS, Su, 2002) leads to up to 15 %  
96 underestimation of ET at the larger grid resolution in a study area in the south-east of Australia.

97  
98 Rouholahnejad Freund and Kirchner (2017) quantified the impact of sub-grid heterogeneity on grid-average ET  
99 using a simple Budyko curve (Turc, 1954; Mezentsev, 1955) in which long-term average ET is a non-linear function  
100 of long-term averages of precipitation (P) and potential evaporation (PET). They showed mathematically that  
101 averaging over spatially heterogeneous P and PET results in overestimation of ET within the Budyko framework (Fig.  
102 1). Their analysis implies that large-scale ESMs that overlook land surface heterogeneity will also yield biased  
103 evapotranspiration estimates due to the inherent nonlinearity in ET processes. They did not, however, determine  
104 where around the globe, and under what conditions, this heterogeneity bias is likely to be most important.

105  
106 The recognition that spatial averaging can potentially lead to biased flux estimates has prompted methods for  
107 representing sub-grid-scale heterogeneities and processes within large scale land surface models and ESMs.  
108 Accounting for land surface heterogeneity in large-scale ESMs is not merely constrained by limitations in both  
109 computational power (Baker et al. 2017) and the availability of high-resolution forcing data, but also by the fact  
110 that the atmospheric and land surface components of some ESMs operate at different resolutions. There have been

111 several attempts to integrate sub-grid heterogeneity in ESMs while keeping the computational costs affordable. In  
112 “mosaic” approaches, the model is run separately for each surface type in a grid cell, and then the surface-specific  
113 fluxes are area-weighted to calculate the grid-cell average fluxes (e.g., Avissar and Pielke, 1989; Koster and Suarez,  
114 1992). The “effective parameter” approach (e.g., Wood and Mason, 1991; Mahrt et al., 1992), by contrast, seeks to  
115 estimate effective parameter values at the grid cell scale that subsume the effects of sub-grid heterogeneity.  
116 Estimating these effective parameters can be challenging because the relevant land-surface processes typically  
117 depend nonlinearly on multiple interacting parameters, and land-surface signals at different scales are propagated  
118 and diffused differently in the atmosphere. Alternatively, the “correction factor” approach (e.g., Maayar and Chen,  
119 2006) uses sub-grid information on spatially heterogeneous land-surface processes and properties to estimate  
120 multiplicative correction factors for fluxes that are originally calculated from spatially averaged inputs at the grid-  
121 cell scale. All three approaches try to reduce the heterogeneous problem to a homogeneous one that has  
122 equivalent effects on the atmosphere at the grid-cell scale.

123  
124 There is a growing need to understand how sub-grid heterogeneity (and the atmosphere’s integration of it) affect  
125 grid-scale water and energy fluxes, and to develop effective methods to incorporate these effects in ESMs (Clark et  
126 al., 2015, Fan et al., 2019). In a previous study, we proposed a general framework for quantifying systematic biases  
127 in ET estimates due to averaging over heterogeneities (Rouholahnejad Freund and Kirchner, 2017). We used the  
128 Budyko framework as a simple estimator of ET, and demonstrated theoretically how averaging over heterogeneous  
129 precipitation and potential evapotranspiration can lead to systematic overestimation of long-term average ET  
130 fluxes from heterogeneous landscapes. In the present study, we apply this analysis across the globe and highlight  
131 the locations where the resulting heterogeneity bias is largest. Our hypotheses, derived from the Budyko  
132 framework as summarized in Eq. (4) below, are that (1) strongly heterogeneous landscapes, such as mountainous  
133 terrain, will exhibit greater heterogeneity bias, (2) this bias will be larger in climates where P and PET are inversely  
134 correlated in space, and (3) heterogeneity bias will decrease as the spatial scales of averaging decrease.

135

## 136 **2. Effects of sub-grid heterogeneity on ET estimates in the Budyko framework**

137 Budyko (1974) showed that long-term annual average evapotranspiration is a function of both the supply of water  
138 (precipitation, P) and the evaporative demand (potential evapotranspiration, PET) under steady-state conditions  
139 and in catchments with negligible changes in storage (Eq. 1; Turc, 1954; Mezentsev, 1955):

$$140 \quad ET = f(P, PET) = \frac{P}{\left(\left(\frac{P}{PET}\right)^n + 1\right)^{1/n}}. \quad (1)$$

141 where ET is actual evapotranspiration, P is precipitation, PET is potential evaporation, and  $n$  (dimensionless) is a  
142 catchment-specific parameter that modifies the partitioning of P between ET and discharge.

143

144 Evapotranspiration rates are inherently bounded by energy and water limits. Under arid conditions ET is limited by  
 145 the available supply of water (the water limit line in Fig. 1b), while under humid conditions ET is limited by  
 146 atmospheric demand (PET) and converges toward PET (the energy limit line in Fig. 1b). Budyko showed that over a  
 147 long period and under steady-state conditions, hydrological systems function close to their energy or water limits.  
 148 These intrinsic water and energy constraints make the Budyko curve downward-curving.

149  
 150 In a heterogeneous landscape, like the simple example of two model columns in Fig. 1a, P and PET vary spatially.  
 151 The two columns with heterogeneous P and PET are represented by the two solid black circles on the Budyko curve  
 152 in Fig. 1b. In this hypothetical two-column example, the true average of ET values calculated from individual  
 153 heterogeneous inputs (the solid black circles) lies below the curve (the grey circle, labeled "true average").  
 154 However, if we aggregate the two columns and consider the system as one column with average properties, the  
 155 function of average inputs (averaged P and PET over the two columns) lies on the Budyko curve (the open circle)  
 156 which is larger than the true average of the two columns. In short, in any downward curving function, the function  
 157 of the average inputs (the open circle) will always be larger than the average of the individual function values (the  
 158 true average; grey circle). The difference between the two can be termed the "heterogeneity bias".

159  
 160 In a previous study (Rouholahnejad Freund and Kirchner, 2017) we showed that when nonlinear underlying  
 161 relationships are used to predict average behaviour from averaged properties, the magnitude of the resulting  
 162 heterogeneity bias can be estimated from the degree of the curvature in the underlying function and the range  
 163 spanned by the individual data being averaged. Here we summarize these findings as building blocks of the current  
 164 study. The second-order, second-moment Taylor expansion of the ET function  $f(P, PET)$  (Eq. 1) around its mean  
 165 directly yields:

$$166 \quad \bar{f}(P, PET) = \overline{ET} \approx f(\bar{P}, \overline{PET}) + \frac{1}{2} \frac{\partial^2 f}{\partial P^2} var(P) + \frac{1}{2} \frac{\partial^2 f}{\partial PET^2} var(PET) + \frac{\partial^2 f}{\partial P \partial PET} cov(P, PET) \quad , \quad (2)$$

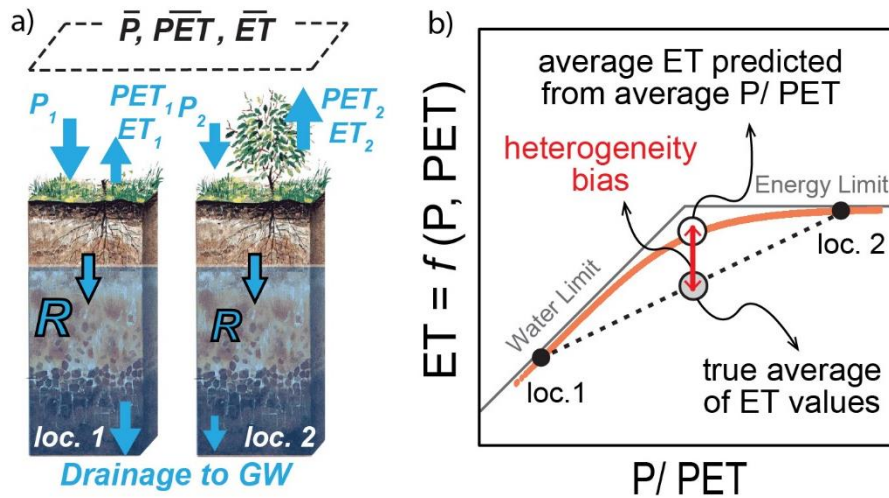
167 where  $\bar{f}(P, PET)$  is the true average of the spatially heterogeneous ET function,  $f(\bar{P}, \overline{PET})$  is the ET function  
 168 evaluated at its average inputs  $\bar{P}$  and  $\overline{PET}$ , and the derivatives are calculated at  $\bar{P}$  and  $\overline{PET}$ . Evaluating the  
 169 derivatives using Eq. (1) and reshuffling the terms, Rouholahnejad Freund and Kirchner (2017) obtained the  
 170 following expression for the heterogeneity bias, the difference between the average ET,  $\bar{f}(P, PET)$ , and the ET  
 171 function evaluated at the mean of its inputs,  $f(\bar{P}, \overline{PET})$ :

$$172 \quad f(\bar{P}, \overline{PET}) - \bar{f}(P, PET) \approx (n + 1) \frac{\bar{P}^{n+1} \overline{PET}^{n+1}}{(\bar{P}^n + \overline{PET}^n)^{2+1/n}} \left[ \frac{1}{2} \frac{var(P)}{\bar{P}^2} + \frac{1}{2} \frac{var(PET)}{\overline{PET}^2} - \frac{cov(P, PET)}{\bar{P} \overline{PET}} \right]. \quad (3)$$

173 To more clearly show the effects of variations in P and PET, Eq. (3) can be reformulated as follows:

$$174 \quad (n + 1) \frac{\bar{P}^{n+1} \overline{PET}^{n+1}}{(\bar{P}^n + \overline{PET}^n)^{2+1/n}} \left[ \frac{1}{2} \left( \frac{SD(P)}{\bar{P}} \right)^2 + \frac{1}{2} \left( \frac{SD(PET)}{\overline{PET}} \right)^2 - r_{P, PET} \left( \frac{SD(P)}{\bar{P}} \right) \left( \frac{SD(PET)}{\overline{PET}} \right) \right]. \quad (4)$$

175 Equation (4) shows that the heterogeneity bias depends on only four quantities: the fractional variation (i.e., the  
 176 coefficient of variation) in precipitation ( $\frac{SD(P)}{\bar{P}}$ ) and in potential ET ( $\frac{SD(PET)}{\bar{PET}}$ ), the correlation between precipitation  
 177 and potential ET ( $r_{P,PET}$ ), and the function  $(n + 1) \frac{\bar{P}^{n+1} \bar{PET}^{n+1}}{(\bar{P}^n + \bar{PET}^n)^{2+1/n}}$ , which quantifies the curvature in the ET function  
 178 in Budyko space. As shown by Fig. 1b and Eq. (2), the discrepancy between average of the ET function and the ET  
 179 function of the average inputs (the heterogeneity bias) is proportional to both the degree of nonlinearity in the  
 180 function, as defined by its second derivatives, and the variability of P and PET. Equation (4) allows one to estimate  
 181 how much the curvature of the ET function and the fractional variability (standard deviation divided by mean) of P  
 182 and PET will affect estimates of ET. However, to the best of our knowledge, the consequences of these  
 183 nonlinearities for global evaporative flux estimates have not previously been quantified.  
 184



185  
 186 Figure 1. Heterogeneity bias in a hypothetical two-column model in the Budyko framework. The true average ET of  
 187 the columns (gray circle) lies below the curve and is less than the average ET estimated from the average P/PET of  
 188 the two columns (open circle). The heterogeneity bias depends on the curvature of the function and the spread of  
 189 its inputs. Both panels are adapted from Rouholahnejad Freund and Kirchner (2017).  
 190

### 191 3. Effects of sub-grid heterogeneity on ET estimates at 1° by 1° grid scale across the globe

192 Across a landscape of similar size to a typical ESM grid cell (1° by 1°), soil moisture, atmospheric demand (PET) and  
 193 precipitation (P) will vary with topographic position; hillslopes will typically be drier, and riparian regions will be  
 194 wetter. To map the spatial pattern in the heterogeneity bias that could result from averaging over this land surface  
 195 heterogeneity, we applied the approach outlined in section 2 to the global land surface area at 1° by 1° grid scale.  
 196 Within each 1° by 1° grid cell, we used 30 arc-second values of P (WorldClim; Hijmans et al., 2005) and PET  
 197 (WorldClim; Hijmans et al., 2005) to examine the variations in small-scale climatic drivers of ET. Because 30 arc-  
 198 seconds is nearly 1 km, hereafter we refer to the 30 arc-second data as 1km values for simplicity. The spatial  
 199 distribution of long-term annual averages (1960-1990) of P and PET values at 1 km resolution, along with 1km

200 values of the aridity index ( $AI=P/PET$ ), are shown in Fig 2a-c. ET values calculated from these 1km P and PET values  
201 using Eq. (1) are then averaged at  $1^\circ$  by  $1^\circ$  scale (“true average”, Fig. 2e). We also averaged the 1km values of P and  
202 PET within each grid cell and then modeled ET using the Budyko curve (Eq. 1) applied to these averaged input  
203 values. The difference between these two ET estimates is the heterogeneity bias.

204

205 We also calculated the heterogeneity bias using Eq. (4), which describes how the nonlinearity in the governing  
206 equation and the heterogeneity in P and PET jointly contribute to the heterogeneity bias. The heterogeneity bias  
207 estimates obtained by Eq. (4) were functionally equivalent ( $R^2=0.97$ , root mean square error of 0.17%) to those  
208 obtained by direct calculation using Eq. (1) as described above.

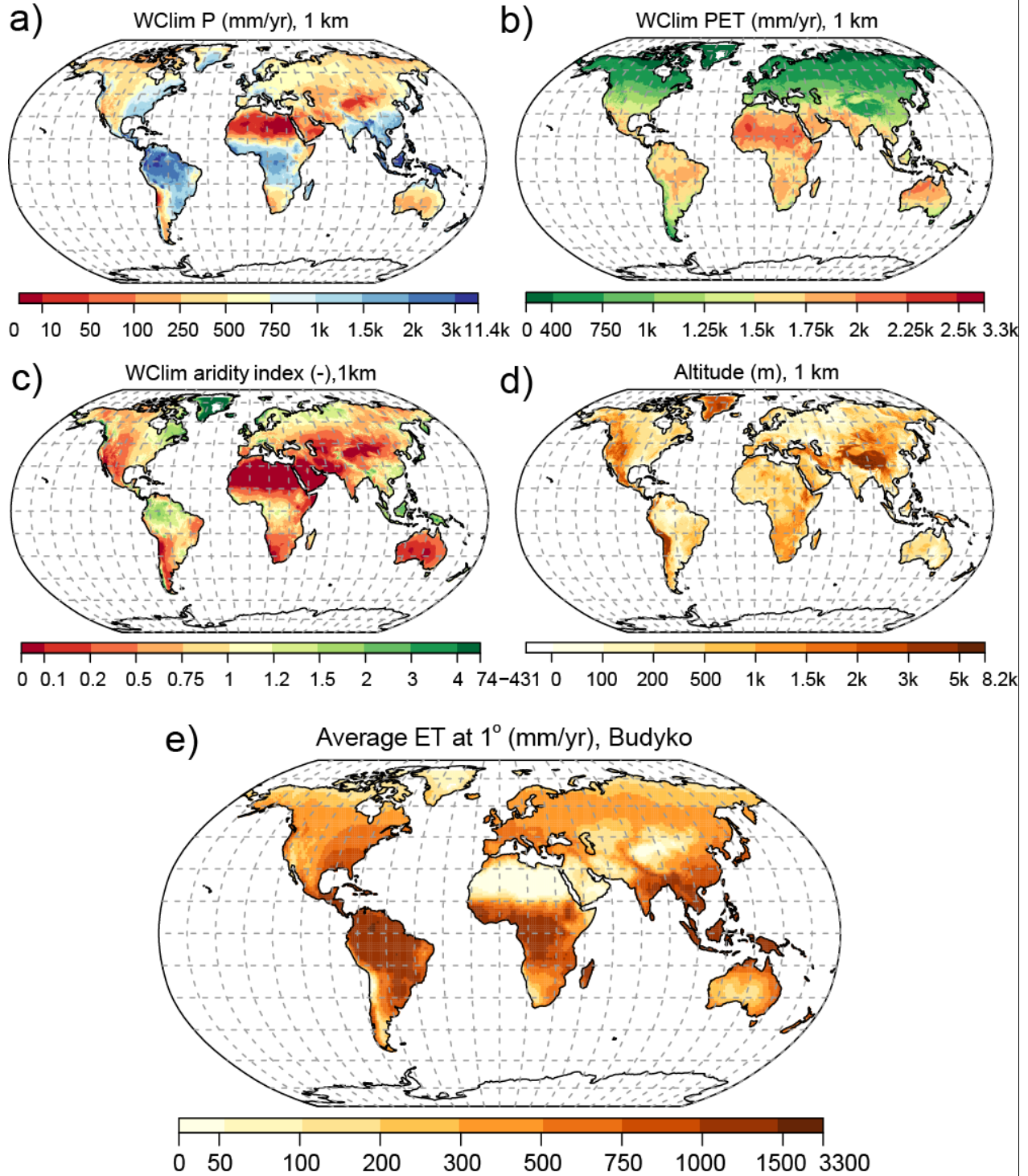
209

210 Fig. 3a-d illustrates the variability (quantified by standard deviation) of 1km values of P, PET, aridity index, and  
211 altitude at the  $1^\circ$  by  $1^\circ$  grid scale. The heterogeneity bias in long-term average ET fluxes at the  $1^\circ$  by  $1^\circ$  grid scale  
212 (Fig. 3e) highlights regions around the globe where ET fluxes are likely to be systematically overestimated. The  
213 spatial distribution of the heterogeneity bias calculated using Eq. 4 (Fig. 3e) closely coincides with locations where  
214 the aridity index is highly variable (Fig. 3c), which is driven in turn by topographic variability (Fig. 3d). Strongly  
215 heterogeneous landscapes exhibit larger estimated heterogeneity biases in long-term average ET fluxes. Although  
216 the global average of our Budyko-based heterogeneity bias estimates is small (<1%), physically based ET  
217 calculations may exhibit larger heterogeneity biases than the modest values we calculate here, because the Budyko  
218 approach already subsumes spatial heterogeneity effects at the catchment scale (and also temporal heterogeneity  
219 effects due to its steady-state assumptions). The heterogeneity biases in ET estimates shown in Fig. 3e correspond  
220 to long-term average ET estimates. Given the fact that P and PET can vary temporally (i.e., seasonality), the actual  
221 bias could be much larger, particularly where P and PET are inversely correlated (see the last term of Eq. 4).

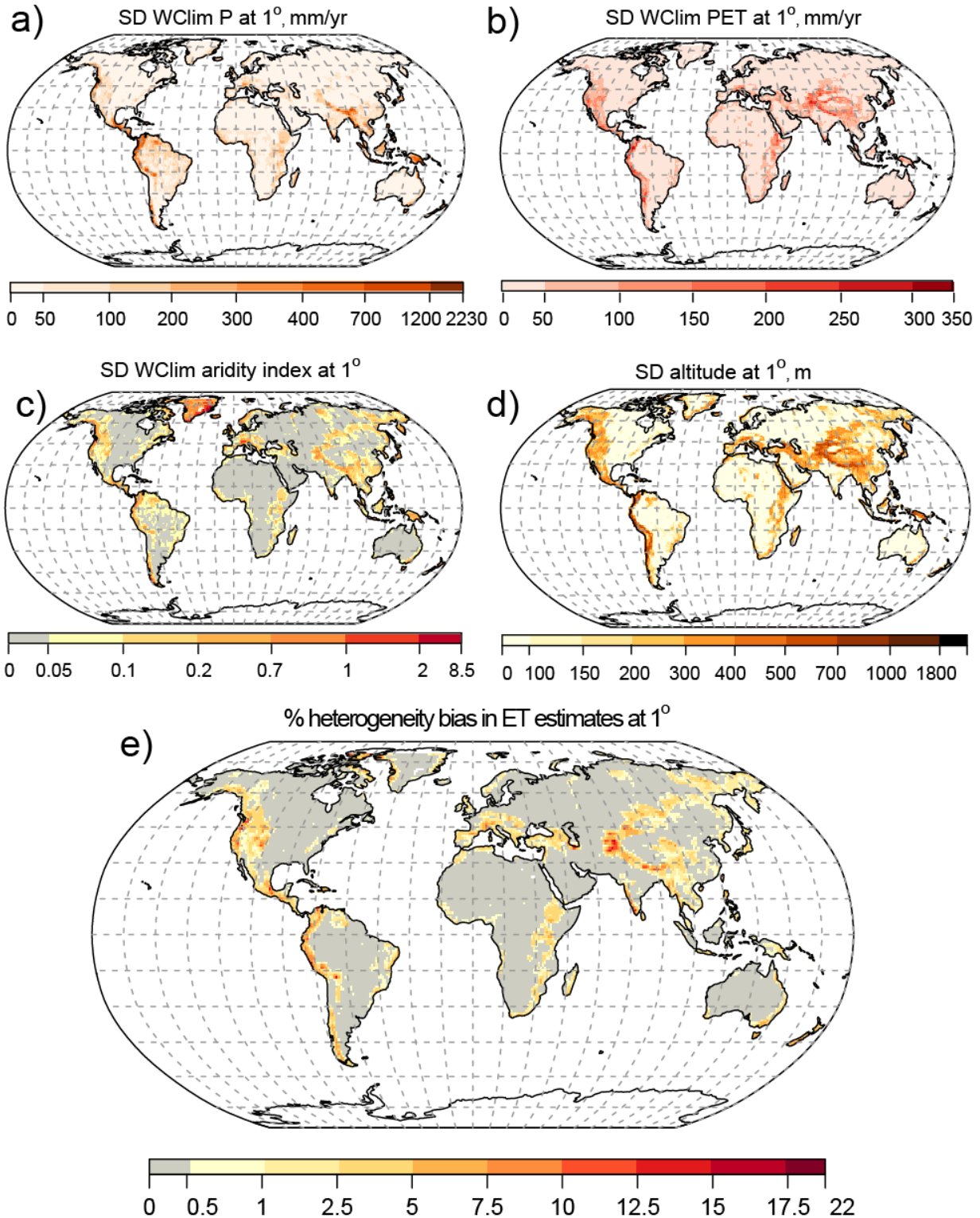
222

223 Our results show that the topographic gradient, and hence the variability in the aridity index across a given grid  
224 scale, drives consistent, predictable patterns of heterogeneity bias in evapotranspiration estimates at that scale.  
225 Equation 4 shows that this bias is equally sensitive to fractional variability in P and PET (standard deviation divided  
226 by mean). However, because P is typically more variable (in percentage terms) than PET across landscapes, the  
227 variability in P will usually make a larger contribution to the estimated heterogeneity bias.





229  
 230 Figure 2. Global distribution of one-kilometer resolution annual mean precipitation (a: P; WorldClim; Hijmans et al.,  
 231 2005), potential evapotranspiration (b: PET; WorldClim; Hijmans et al., 2005), aridity index (c: AI=P/PET; WorldClim;  
 232 Hijmans et al., 2005), and topography (d: SRTM; Jarvis et al., 2008), along with (e) evapotranspiration (ET) at 1° by  
 233 1° scale by averaging 1km values of ET calculated using the Budyko function (Eq. 1).  
 234



235  
 236 Figure 3. Global spatial distribution of variability (standard deviation) of one-kilometer values of a) precipitation (P),  
 237 b) potential evapotranspiration (PET), c) aridity index (AI=P/PET), and d) altitude at 1° by 1° grid cell. The  
 238 heterogeneity bias in ET estimates (e) is calculated using Eq. (4). Grid cells with larger standard deviation in altitude  
 239 and aridity index have larger heterogeneity bias.

#### 240 **4. Variation in heterogeneity bias across climate zones, data sources, and grid scales**

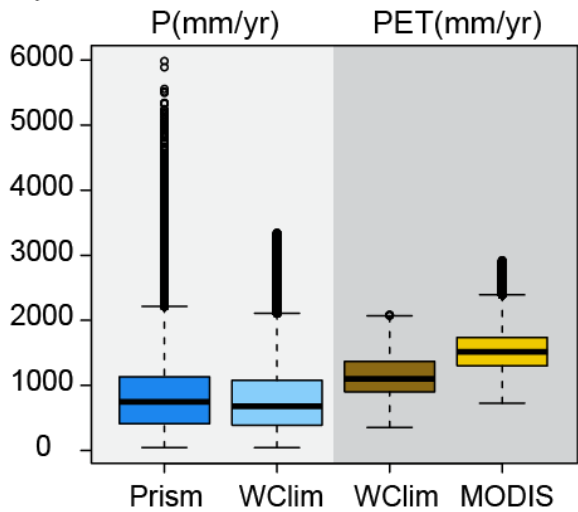
241 With increased availability of spatial data, it is becoming standard practice to assess input data uncertainties and  
242 their propagated impacts on water and energy flux estimates in land surface models. To quantify how choices  
243 among alternative input data products could affect the heterogeneity bias in ET estimates, we calculated the  
244 heterogeneity bias at 1° by 1° grid cell resolution across the contiguous US using four different pairs of P and PET  
245 data products. Two precipitation data sets, Prism (<http://prism.oregonstate.edu>) and WorldClim (Hijmans et al.,  
246 2005), along with two PET data sets, MODIS (Mu et al., 2007) and WorldClim (Hijmans et al., 2005). As Prism  
247 precipitation data is available at 4 km resolution, all other data sets were aggregated to 4 km. Two P products and  
248 two PET products were combined in all possible pairs. The WorldClim PET dataset (Hijmans et al., 2005) is based on  
249 the Hargreaves method (Hargreaves and Samani 1985) while the MODIS PET product (Mu et al, 2007) is based on  
250 the Penman–Monteith equation (Monteith, 1965). The heterogeneity bias in ET estimates (Eq. 4), as outlined in  
251 Sect. 2, was evaluated from 4km values of P, PET, and the estimated average ET using the Budyko relationship (Eq.  
252 1) for each of the four input data pairs. Figure 4a-e compares the spatial distributions of heterogeneity bias across  
253 the contiguous US for the four pairs of P and PET data products. The heterogeneity bias in ET estimates reached as  
254 high as 36 % in the western US using Prism P and WorldClim PET as input to the ET model (Fig. 4b). A visual  
255 comparison of Figs. 4b and Fig. 4d shows that the choice of P data source (Prism vs. WorldClim) had a bigger effect  
256 on the heterogeneity bias than the choice of PET data source (MODIS vs. WorldClim), meaning that the fractional  
257 variability in P is the dominant variable. In all cases, data sources that were more variable in relation to their means  
258 (Prism for P and WorldClim for PET; Fig. 4b) led to larger estimates of heterogeneity bias, as expected from Eq. (4).  
259 Thus we infer that we would have obtained larger heterogeneity biases if we had conducted our global analysis  
260 (Fig. 3) with Prism P and either WorldClim or MODIS PET, but we cannot show that result explicitly at global scale  
261 because Prism P is not freely available globally.

262  
263 If we separate the heterogeneity biases shown in Fig. 4 according to Köppen-Geiger climate zones (Peel et al., 2007;  
264 Fig. 5a), we see that they are distinctly higher in particular climate-terrain combinations. Estimated heterogeneity  
265 biases are higher in regions with temperate climates and dry summers (climate zone Cs) and in regions with cold,  
266 dry summers (climate zone Ds), most likely due to the sharp spatial gradient in their water and energy sources for  
267 evapotranspiration (Fig. 5b). These areas typically have high topographic relief, combined with seasonal climate.  
268 The heterogeneity effects on ET estimates in these regions are expected to be even larger when a mechanistic  
269 model of ET is used. We expect that averaging over temporal variations of drivers of ET, especially in places with  
270 strong seasonality, could substantially bias the ET estimates, but this cannot be quantified in the Budyko framework  
271 due to its underlying steady-state assumptions. Figure 5b also illustrates the relative magnitudes of the  
272 heterogeneity biases obtained with the four pairs of P and PET data sources. The estimated heterogeneity bias is  
273 highest when the Prism P and WorldClim PET datasets are used, followed by the combination of Prism P and MODIS  
274 PET, which resulted in the second-highest heterogeneity bias across different climate zones. Wilcoxon signed-rank  
275 tests was performed to evaluate the statistical significance of the differences between heterogeneity bias in ET

276 estimates using all pairs of climate zones and data sources that are shown in Fig. 5b (Table S1). These analysis show  
277 that while the difference between heterogeneity biases estimated in Cs and Ds climate zones are not statistically  
278 significant across all four combinations of datasets, the difference between estimated heterogeneity bias in Cs  
279 versus Cf, Ds versus Cf, as well as Cs versus Bs climate zones are significant across all four data combinations  
280 (highlighted in Table S1 of the supplementary material).

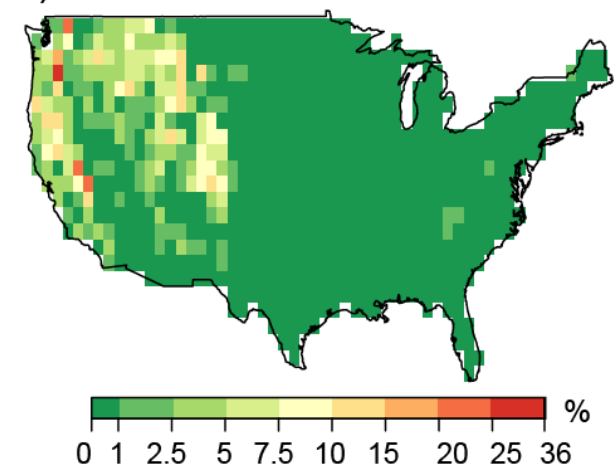
281  
282 Equation 4 shows that heterogeneity biases in Budyko estimates of ET are equally sensitive to the same percentage  
283 variability in P and PET. Thus the degree of sensitivity, per se, to P and PET variations expressed in percentage terms  
284 is the same. Although Figs. 5c and 5d give the visual impression that PET is more variable than P across climate  
285 zones and between data sources, Fig. 5e shows that the fractional variability in P is systematically higher than PET,  
286 and it also varies more across the climate zones and between the two data sets. Because P is typically more  
287 variable than PET (in percentage terms) across landscapes, the variability in P will make a larger contribution to the  
288 heterogeneity bias (Fig. 5e) estimated using the Budyko approach. Whether this is true for more physically based ET  
289 estimates remains to be seen. Analysis of percent variability of P and PET products shows that percent variabilities  
290 of precipitation products are in general larger than PET products and hence contribute more to heterogeneity (Fig  
291 5e). While the percent variabilities of the two PET products are in the same range, the percent variability in Prism  
292 precipitation is slightly larger than in WorldClim precipitation, in regions with dry summers (Cs and Ds climate zones  
293 in Fig. 5a).  
294

a) Distribution of P and PET in the four datasets



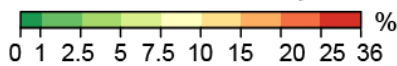
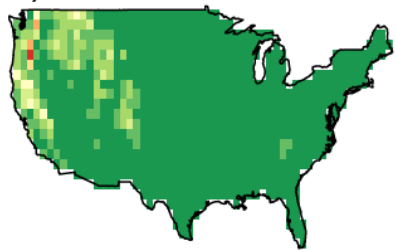
% heterogeneity bias in ET estimates at 1°

b) Prism P, Wclim PET as inputs

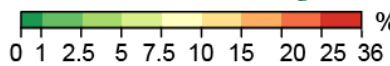
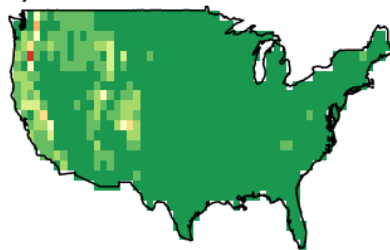


% heterogeneity bias in ET estimates at 1°

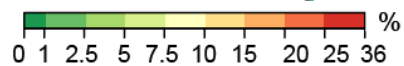
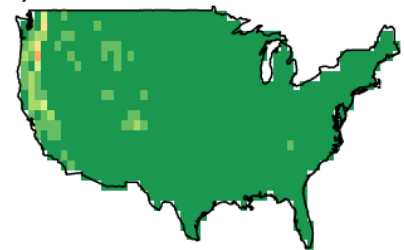
c) Prism P, MODIS PET as inputs



d) Wclim P, Wclim PET as inputs



e) Wclim P, MODIS PET as inputs



295

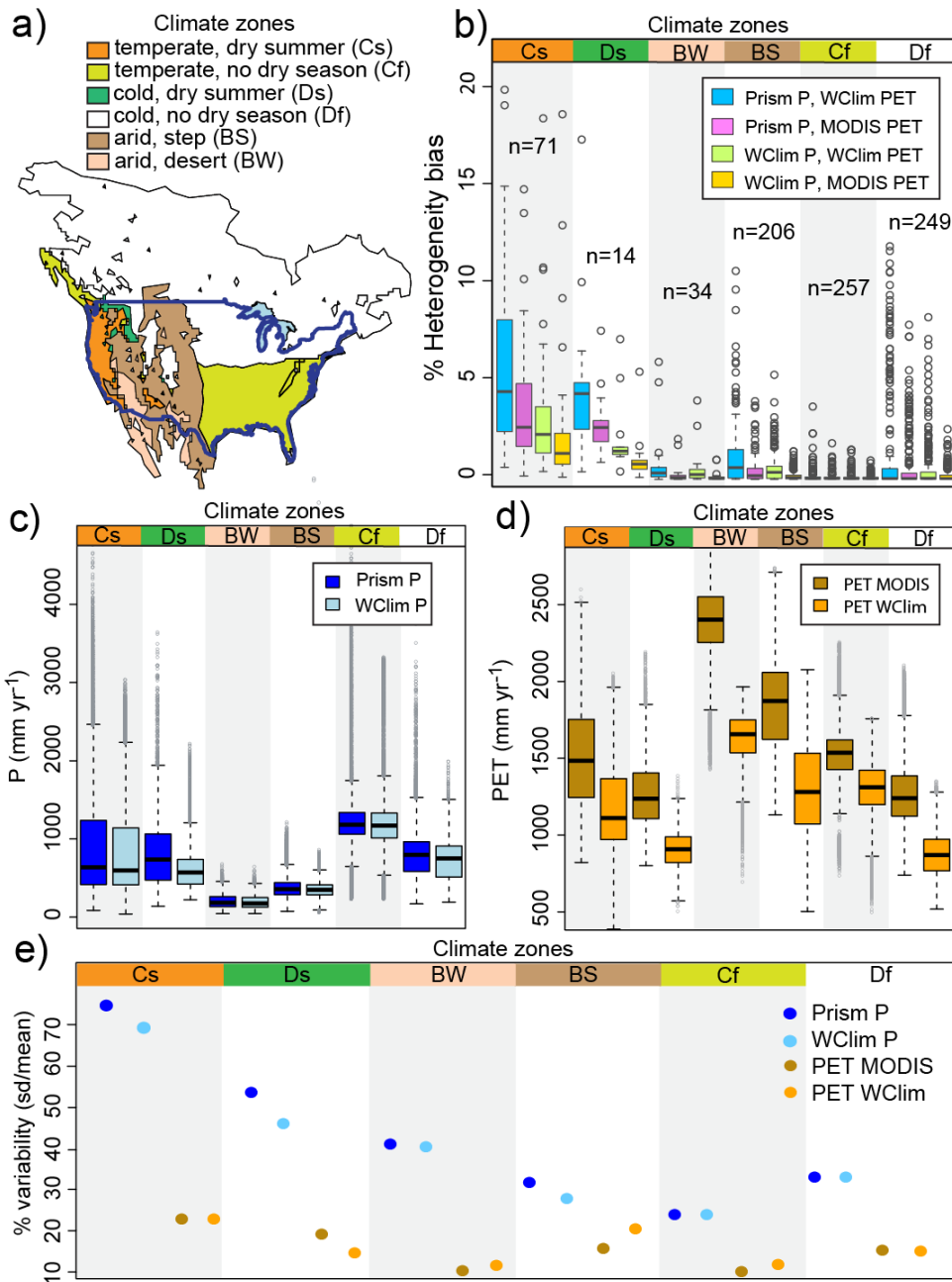
296

297

298

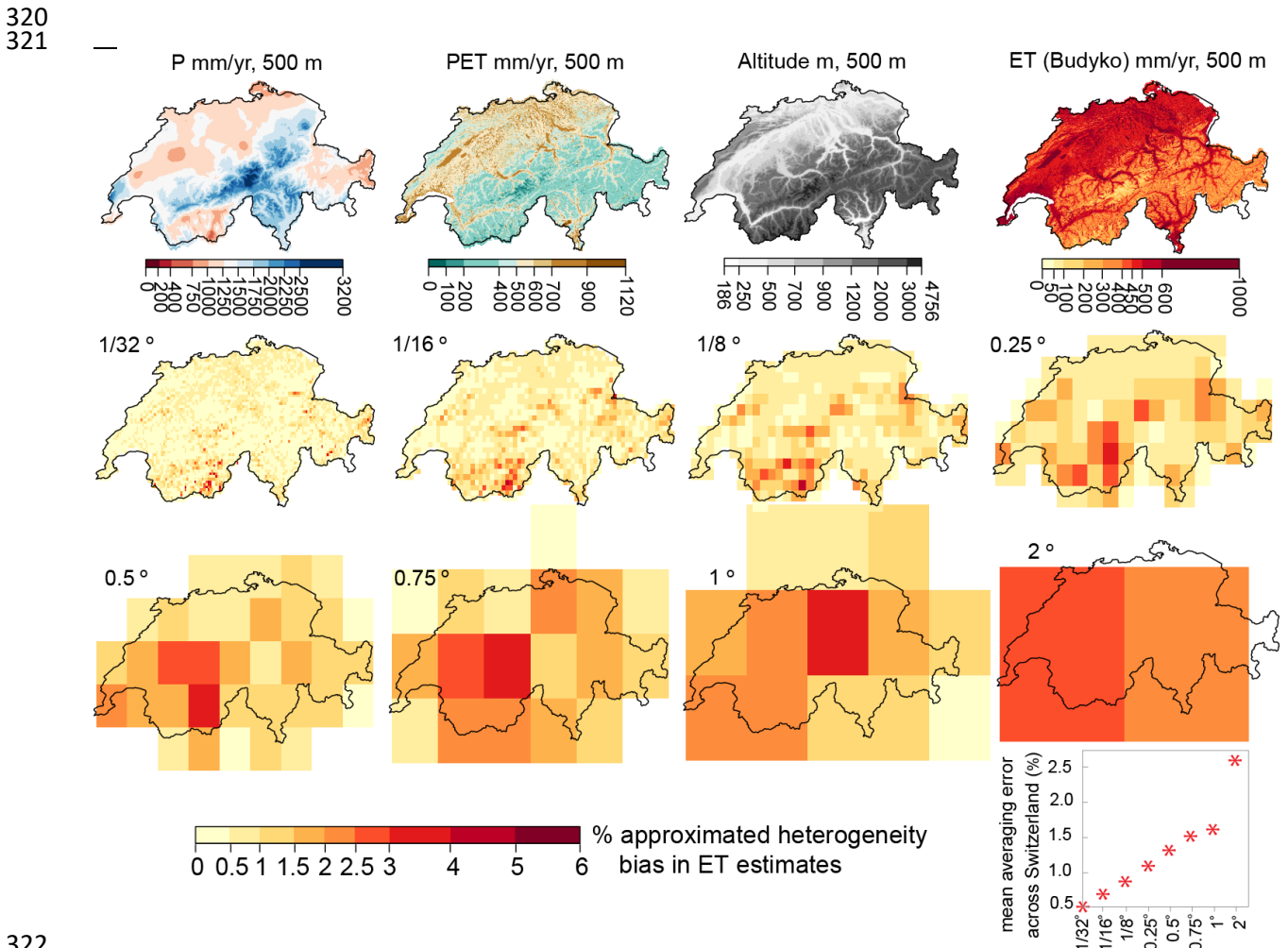
299

Figure 4. The distribution of P and PET in the four datasets is shown in a). Estimated heterogeneity bias (Eq. 4) across the contiguous US using four-kilometer values of b) Prism P and WorldClim PET c) Prism P and MODIS PET d) WorldClim P and WorldClim PET, and e) WorldClim P and MODIS PET as inputs.



300  
 301 Figure 5. a) Köppen-Geiger climate classification (Peel et al., 2007 in Beck et al. 2013) across the contiguous US, b)  
 302 the distribution of calculated heterogeneity bias in ET estimates (Eq. 4) at 1° by 1° grid cell in individual climate  
 303 zones, shown by boxplot (three data points with heterogeneity biases of over 20% are off-scale). The significance of  
 304 differences between the pairs are presented in Table S1. Panels c and d show the distribution of precipitation  
 305 products (Prism and WorldClim) and potential evaporation products (MODIS and WorldClim) at individual climate  
 306 zones, respectively. The color-coded climate zones at the tops of panels b, c, and d correspond to the climate zones  
 307 mapped in panel a. Panel e compares the percentage variability of the two P and PET data products across climate  
 308 zones, showing that the percentage variability in P is markedly higher than in PET, and the percentage variability in  
 309 Prism P is somewhat higher than in WorldClim P, particularly in climate zones with dry summers.

310 Because future increases in computing power will lead to ESMs with smaller grid cells, it is useful to ask how  
 311 changes in grid resolution affect the heterogeneity biases that we have estimated in this paper. To quantify the  
 312 heterogeneity bias in ET estimates as a function of grid scale, we repeated our analysis at various grid resolutions  
 313 using Switzerland as a test case. We started with high-resolution (500m) maps of long-term average annual  
 314 precipitation and PET across the Swiss landscape (Fig. 6), and then used Eq. 4 to estimate the heterogeneity bias at  
 315 grid scales ranging from  $1/32^\circ$  to  $2^\circ$  ( $\sim 3$  km to  $\sim 200$  km). As Fig. 6 shows, aggregating P and PET over larger scales  
 316 leads to larger, and more widespread, overestimates in ET. Conversely, at finer grid resolutions, the average  
 317 heterogeneity bias is smaller, and the locations with large biases are more localized. On average, the heterogeneity  
 318 bias across Switzerland as a whole grows exponentially as the inputs are averaged over larger grids (as shown in the  
 319 lower-right panel in Fig. 6).



322 Figure 6. Heterogeneity bias in ET estimates at various scales across Switzerland, estimated from 500m climate  
 323 data. ET is calculated using the Budyko relationship (Eq. 1). Heterogeneity bias was estimated from 500m  
 324 precipitation (P) and potential evapotranspiration (PET), and their variances at each grid scale, using Eq. 4. At finer  
 325 grid resolutions, the heterogeneity bias is more localized, and smaller on average.  
 326  
 327

328 **5. Summary and discussion**

329 Because evapotranspiration (ET) processes are inherently bounded by water and energy constraints, over the long  
330 term, ET is always a nonlinear function of available water and PET, whether this function is expressed as a Budyko  
331 curve or another ET model. These nonlinearities imply that spatial heterogeneity will not simply average out in  
332 predictions of land surface water and energy fluxes. Overlooking sub-grid spatial heterogeneity in large-scale land  
333 surface models could lead to biases in estimating these fluxes. Here we have shown that, across several scales,  
334 averaging over spatially heterogeneous land surface properties and processes leads to biases in evapotranspiration  
335 estimates. We examined the global distribution of this bias, its scale dependence, and its sensitivity to variations in  
336 P versus PET, and showed under what conditions this heterogeneity bias is likely to be most important. Our analysis  
337 does not quantify the heterogeneity biases in ESMs, owing to the many differences between these mechanistic  
338 models and the simple empirical Budyko curve. But if the heterogeneity biases in ESMs can be quantified, they can  
339 be used as correction factors to improve ESM estimates of surface-atmosphere water and energy fluxes across  
340 landscapes. Our paper highlights a general methodology that can be used to estimate heterogeneity biases and to  
341 map their spatial patterns, but not to calculate their absolute magnitudes because those will change significantly  
342 depending on the ET formulation that is used.

343  
344 In this study, we used Budyko curves as simple models of ET, in which long-term average ET rates are functionally  
345 related to long-term averages of P and PET. We used an approach outlined by Rouholahnejad Freund and Kirchner  
346 (2017) to estimate the heterogeneity bias in modeled ET at 1-degree grid scale across the globe (Fig. 3), and also at  
347 multiple grid scales across Switzerland (Fig. 6), using finer-resolution P and PET values as drivers of ET. We showed  
348 how the heterogeneity effects on ET estimates vary with the nonlinearity in the governing equations and with the  
349 variability in land surface properties. Our analysis shows that heterogeneity effects on ET fluxes matter the most in  
350 areas with sharp gradients in the aridity index, which are in turn controlled by topographic gradients, and not  
351 merely in areas that are either arid or humid (e.g., compare Fig. 3e with Fig. 2c).

352  
353 According to our analysis, regions within the U.S. that have temperate climates and dry summers exhibit greater  
354 heterogeneity bias in ET estimates (Fig. 5). We show that the estimated heterogeneity bias at each grid scale  
355 depends on the variance in the drivers of ET at that scale (Fig. 4), and on the choice of data sources used to  
356 estimate ET. Heterogeneity bias estimates were significantly larger across the contiguous United States when P and  
357 PET data sources with larger variances were used (Fig. 4).

358  
359 We also explored how heterogeneity biases and their spatial distribution vary with the scale at which the climatic  
360 drivers of ET are averaged. We found that as heterogeneous climatic variables are aggregated to larger scales, the  
361 heterogeneity biases in ET estimates become greater on average, and extend over larger areas (Fig. 6). At smaller  
362 grid scales, estimated heterogeneity biases do not completely disappear, but instead become more localized  
363 around areas with sharp topographic gradients. Finding an effective scale at which one can average over the



364 heterogeneity of land surface properties and processes has been a longstanding problem in Earth science. Our  
365 analysis shows that at smaller resolutions the average heterogeneity bias as seen from the atmosphere becomes  
366 smaller, but there is no characteristic scale at which it vanishes entirely (Fig. 6). The magnitude and spatial  
367 distribution of this bias depend strongly on the scale of the averaging and degree of the nonlinearity in the  
368 underlying processes. The heterogeneity bias concept is general and extendable to any convex or concave function  
369 (Rouholahnejad Freund and Kirchner 2017), meaning that in any nonlinear process, averaging over spatial and  
370 temporal heterogeneity can potentially lead to bias.

371  
372 In the analysis presented here, we have assumed a value of 2 for the Budyko parameter  $n$ , which approximates the  
373 variation of ET/PET with respect to P/PET in MODIS and WorldClim data across continental Europe (Mu et al. 2007;  
374 Hijmans et al. 2005; Rouholahnejad Freund & Kirchner, 2017). Although there are suggestions in the literature that  
375  $n$  can vary with land use and other landscape properties (e.g., Teuling et al., 2019), here we have assumed that  $n$  is  
376 spatially and temporally constant in order to focus on the effects of heterogeneity in P and PET. In the supplement  
377 we present a sensitivity analysis with values of  $n$  ranging from 2 to 5 (Fig. S1). That analysis shows that, as expected  
378 from Eqs. 3 and 4, higher values of  $n$  lead to larger heterogeneity biases, because higher values of  $n$  localize the  
379 curvature of the Budyko function more strongly at the transition between the energy and water limits (Fig. 1b),  
380 increasing the heterogeneity bias for P/PET values near this transition. Nonetheless, the spatial pattern shown in  
381 Fig. 3e remains largely unchanged over the full range of  $n$  values that we analyzed, and the Taylor approximation in  
382 Eqs. 3 and 4 yields realistic estimates of the heterogeneity bias for all values of  $n$  that were tested (Fig. S2). Thus  
383 while our numerical estimates of heterogeneity bias depend somewhat on the value of  $n$ , our conclusions do not.

384  
385 One should keep in mind that the true mechanistic equations that determine point-scale ET as a function of point-  
386 scale water availability and PET (if such data were available) may be much more nonlinear than Budyko's empirical  
387 curves, because these curves already average over significant spatial and temporal heterogeneity. Thus, we expect  
388 that the real-world effects of sub-grid heterogeneity are probably larger than those we have estimated in Sects. 3  
389 and 4 of this study. In addition, the 1km P and PET values that are used in our global analysis might be still too  
390 coarse to represent small-scale heterogeneity that is important to evapotranspiration processes.

391  
392 Budyko curves are empirical relationships that functionally relate evaporation processes to the supply of water and  
393 energy under steady-state conditions in closed catchments with no changes in storage. Our analysis likewise  
394 assumes no changes in storage, nor any lateral transfer between the model grid cells, although both lateral  
395 transfers and changes in storage may be important, both in the real world and in models. Unlike the Budyko  
396 framework, ET fluxes in most ESMs are often physically based (not merely functions of P and PET) and are  
397 calculated at much smaller time steps (seconds to minutes). These models often represent more processes that are  
398 important to evapotranspiration (such as storage variations) and include their dynamics to the extent that is  
399 computationally feasible. Because these relationships may be much more nonlinear than Budyko curves, much

400 larger heterogeneity biases could result when complex physically based models are used to estimate ET from  
401 spatially aggregated data. Therefore, we are now working to quantify heterogeneity bias in ET fluxes using a more  
402 mechanistic land surface model.

403

#### 404 **Acknowledgements**

405 E.R.F. acknowledges support from the Swiss National Science Foundation (SNSF) under Grant No. P2EZP2\_162279.  
406 The authors thank Massimiliano Zappa of the Swiss Federal Research Institute WSL for providing the 500m  
407 resolution data that enabled the analysis shown in Fig. 6.

408

#### 409 **References**

410 Aminzadeh M., and D. Or: The complementary relationship between actual and potential evaporation for spatially  
411 heterogeneous surfaces, *Water Resour. Res.*, 53, 580–601, doi:10.1002/2016WR019759, 2017.

412 Avissar, R., R. A. Pielke: A Parameterization of Heterogeneous Land Surfaces for Atmospheric Numerical Models and  
413 Its Impact on Regional Meteorology, *Monthly Weather Review*, vol. 117, issue 10, p. 2113, doi:10.1175/1520-  
414 0493(1989)117<2113:APOHLS>2.0.CO;2, 1989.

415 Baker I. T. , P. J. Sellers , A. S. Denning, I. Medina , P. Kraus, K. D. Haynes , and S. C. Biraud: Closing the scale gap  
416 between land surface parameterizations and GCMs with a new scheme, SiB3-Bins, *Journal of Advances in Modeling  
417 Earth Systems*, *J. Adv. Model. Earth Syst.*, 9, 691–711, doi:10.1002/2016MS000764, 2017.

418 Bastiaanssen, W. G. M., M. Menenti, R. A. Feddes, and A. A. M. Holtslag: A remote sensing surface energy balance  
419 algorithm for land (SEBAL): 1. Formulation, *Journal of Hydrology*, 212-213, 198–212, 1998.

420 Beck H. E., A. I. J. M. van Dijk, D. G. Miralles, R. A. M. de Jeu, L. A. Bruijnzeel, T. R. McVicar, and J. Schellekens:  
421 Global patterns in base flow index and recession based on streamflow observations from 3394 catchments, *Water  
422 Resour. Res.*, 49, 7843–7863, doi:10.1002/2013WR013918, 2013.

423 Boone, A., and O. J. Wetzel: A simple scheme for modeling sub-grid soil texture variability for use in an atmospheric  
424 climate model. *Journal of the Meteorological Society of Japan*, 77(1), 317–333, 1998.

425 Budyko, M. I.: *Climate and life*, Academic, New York, 1974.

426 Clark, M. P., Y. Fan, D. M. Lawrence, J. C. Adam, D. Bolster, D. J. Gochis, R. P. Hooper, M. Kumar, L. R. Leung, D. S.  
427 Mackay, R. M. Maxwell, C. Shen, S. C. Swenson, and X. Zeng: Improving the representation of hydrologic processes  
428 in Earth System Models, *Water Resour. Res.*, 51, 5929–5956, doi:10.1002/2015WR017096, 2015.

429 Ershadi A., M. F. McCabe, J. P. Evans, J. P. Walker: Effects of spatial aggregation on the multi-scale estimation of  
430 evapotranspiration, *Remote Sensing of Environment* 131, 51–62, <http://dx.doi.org/10.1016/j.rse.2012.12.007>,  
431 2013.

432 Fan, Y., M. Clark, D. M. Lawrence, S. Swenson, L. E. Band, S. L. Brantley, P. D. Brooks, W. E. Dietrich, A. Flores, G.  
433 Grant, J. W. Kirchner, D. S. Mackay, J. J. McDonnell, P. C. D. Milly, P. L. Sullivan, C. Tague, H. Ajami, N. Chaney, A.  
434 Hartmann, P. Hazenberg, J. McNamara, J. Pelletier, J. Perket, E. Rouholahnejad-Freund, T. Wagener, X. Zeng, E.  
435 Beighley, J. Buzan, M. Huang, B. Livneh, B. P. Mohanty, B. Nijssen, M. Safeeq, C. Shen, W. van Verseveld, J. Volk, D.  
436 Yamazaki: Hillslope hydrology in global change research and Earth system modeling, *Water Resources Research*, 55,  
437 doi:10.1029/2018WR023903, 2019.

438 Hargreaves, G. H., and Z. A. Samani: Reference crop evaporation from temperature, *Appl. Eng. Agric.*, 1(2), 96-99,  
439 1985.

440 Hijmans, R. J., S. E. Cameron, J. L. Parra, P. G. Jones, and A. Jarvis: Very high resolution interpolated climate surfaces  
441 for global land areas, *Int. J. Climatol.*, 25, 1965–1978, doi:10.1002/joc.1276, 2005.

442 Holland, S., J. L. Heitman, A. Howard, T. J. Sauer, W. Giese, A. Ben-Gal, N. Agam, D. Kool, and J. Havlin: Micro Bowen  
443 ratio system for measuring evapotranspiration in a vineyard interrow, *Agric. For. Meteorol.*, 177, 93–100, 2013.

444 Hong, S. H., J. M. H. Hendrickx, and B. Borchers: Up-scaling of SEBAL derived evapotranspiration maps from Landsat  
445 (30 m) to MODIS (250 m) scale, *Journal of Hydrology*, 370, 122–138, 2009.

446 Jarvis, A., Reuter, H. I., Nelson, A., and Guevara, E.: Hole-filled SRTM for the globe Version 4, available from the  
447 CGIARCSI SRTM 90m Database, <http://srtm.csi.cgiar.org> (last access: 26 February 2016), 2008.

448 Kalma, J. D., T. R. McVicar, and M. F. McCabe: Estimating land surface evaporation: A review of methods using  
449 remotely sensed surface temperature data, *Surv. Geophys.*, 29, 421–469, doi:10.1007/s10712-008-9037-z, 2008.

450 Kollet S. J.: Influence of soil heterogeneity on evapotranspiration under shallow water table conditions: transient,  
451 stochastic simulations, *Environmental Research Letters*, 4, 35007, doi:10.1088/1748-9326/4/3/035007, 2009.

452 Koster R. D. et al.: GLACE: The Global Land– Atmosphere Coupling Experiment. Part I: Overview. *J. Hydrometeorol.*, 7,  
453 590–610, 2006.

454 Koster R. D., and M. Suarez: Modeling the land surface boundary in climate models as a composite of independent  
455 vegetation stands, *J. Geophysical Research*, 97 (D3), 26-97-2715, 1992.

456 Maayar, M. E., J. M. Chen: Spatial scaling of evapotranspiration as affected by heterogeneities in vegetation,  
457 topography, and soil texture, *Remote Sensing of Environment*, 102, 33–51, 2006.

458 Mahrt, L., J. Sun, D. Vickers, J. I. MacPherson, J. R. Perderson, and R. L. Desjardins: Observations of fluxes and inland  
459 breezes over a heterogeneous surface, *J. Atmos. Sci.* 51, 2165e2178, 1992.

460 McCabe M., and E. Wood: Scale influences on the remote estimation of evapotranspiration using multiple satellite  
461 sensors, *Remote Sensing of Environment* 105 (2006) 271–285, 2006.

462 Mezentsev, V. S.: More on the calculation of average total evaporation, *Meteorol. Gidrol.*, 5, 24–26, 1955.

463 Montheith, J. L.: Evaporation and environment, the state of and movement of water in living organisms, *Proceeding*  
464 *of Soc. for Exp. Biol.*, 19, 205–234, doi:10.1002/qj.49710745102, 1965.

465 Mu, Q., F. A. Heinsch, M. Zhao, and S. W. Running: Development of a global evapotranspiration algorithm based on  
466 MODIS and global meteorology data, *Remote Sens. Environ.*, 111, 519–536, doi:10.1016/j.rse.2007.04.015, 2007.

467 Peel, M. C., B. L. Finlayson, and T. A. McMahon: Updated world map of the Köppen-Geiger climate classification,  
468 *Hydrol. Earth Syst. Sci.*, 11, 1633–1644, <https://doi.org/10.5194/hess-11-1633-2007>, 2007.

469 PRISM Climate Group, Oregon State University, <http://prism.oregonstate.edu>, created 22 Feb 2017.

470 Rouholahnejad Freund, E., and J. W. Kirchner: A Budyko framework for estimating how spatial heterogeneity and  
471 lateral moisture redistribution affect average evapotranspiration rates as seen from the atmosphere, *Hydrology*  
472 *and Earth System Sciences*, 21(1), 217–233, 2017.

473 Santanello J. R., and C. D. Peters-Lidard: Diagnosing the Sensitivity of Local Land–Atmosphere Coupling via the Soil  
474 Moisture–Boundary Layer Interaction, *J. Hydrometeorology*, 12, 766–786, doi: 10.1175/JHM-D-10-05014.1, 2011.

475 Sato N., P. J. Sellers, D. A. Randall, E. K. Schneider, J. Shukla, J. L. Kinter III, Y. T. Hou, and E. Albertazzi: Effects of  
476 Implementing the Simple Biosphere Model in a General Circulation Model, *J. Atmospheric Sciences*, 46(18), 2757–  
477 2782, 1989.

478 Seneviratne, S. I., T. Corti, E. L. Davin, M. Hirschi, E. B. Jaeger, I. Lehner, B. Orlowsky, and A. J. Teuling: Investigating  
479 soil moisture–climate interactions in a changing climate: A review, *Earth-Science Reviews*, 99(3–4), 125–161, 2010.

480 Shahraeeni, E., and D. Or: Thermo-evaporative fluxes from heterogeneous porous surfaces resolved by infrared  
481 thermography, *Water Resour. Res.*, 46, W09511, doi:10.1029/2009WR008455, 2010.

482 Su, Z.: The Surface Energy Balance System (SEBS) for estimation of turbulent heat fluxes. *Hydrology and Earth*  
483 *System Sciences*, 6, 85–100, 2002.

484 Teuling, A. J. and de Badts, E. A. G. and Jansen, F. A. and Fuchs, R. and Buitink, J. and Hoek van Dijke, A. J. and  
485 Sterling, S. M., Climate change, reforestation/afforestation, and urbanization impacts on evapotranspiration and  
486 streamflow in Europe, *Hydrology and Earth System Sciences*, 23, 3631–3652, DOI = {10.5194/hess-23-3631-2019,  
487 2019.Turc, L.: Le bilan d'eau des sols: relation entre la precipitations, l'évaporation et l'écoulement, *Ann. Agron. A*,  
488 5, 491–569, 1954.

489 Wood, N., and P. J. Mason: The influence of static stability on the effective roughness length for momentum and  
490 heat transfer, *Quart. J. Roy. Meteor. Soc.* 117, 1025e1056, 1991.