

**Response to editor's comments – report #2**

Dear Editor,

Thank you very much for your comments. Following please find our point-by-point response to your questions and suggestions. The editor's comments are in regular font and our response is in bold. The page and line numbers refer to the revised manuscript that will be submitted with this response (with "all mark up" display for review).

As can be seen by the two review reports, both reviewers still have major concerns. Especially, the concerns about the parameter choice  $n$  are important. Since you study scale dependency, having a constant  $n$  is dangerous. Please provide information on:

- 1) the value of  $n$
- 2) was  $n$  constant?
- 3) how did you determine it?
- 4) was it spatially variable?
- 5) how sensitive is the approach to  $n$ ?

**This is an important point and we should have mentioned the reason to hold  $n$  constant.**

**We did not vary  $n$  either in space or in time, because doing so would create artifacts that would confound the effects of spatial heterogeneity in  $P$  and  $PET$ . For example, if we vary  $n$  from place to place, then how do we separate the effects of spatial heterogeneity from the effects of the imposed variation in  $n$ ? For similar reasons, we do not agree that having a constant  $n$  is "dangerous" for a study of scale-dependency. Instead, it is essential for  $n$  to be held constant because otherwise one cannot separate the effects of scale-dependent variation in  $P$  and  $PET$  from the effects of scale-dependent variation in  $n$ . Recall that in our analysis we use Budyko curves as an analytical framework (or a simple "see-through" function) for exploring the consequences of spatial heterogeneity in landscape properties. We do not quantify the heterogeneity bias in ESMs (which are not based on Budyko curves), nor do we use Budyko curves as a proxy for what ESM ET estimates would be.**

**In the revised version we present a sensitivity analysis (for  $n$  values ranging from 2 to 5) in the supplement (Figs S1, S2) and discuss its main results in the manuscript (under Summary and Discussion). Those results are: 1) the spatial patterns of aggregation bias are similar, 2) the absolute magnitude of aggregation bias increases somewhat for higher values of  $n$ , as predicted by Eqs. 3 and 4, and 3) the Taylor approximation in Eqs. 3 and 4 yield realistic estimates of the aggregation bias for all values of  $n$  that were tested.**

Furthermore, spatial heterogeneity in  $P$  and  $PET$  should be better explained as indicated in report #2, reviewer #1 (first point). As well as the temporal dynamics (report #1, reviewer#3).

**Please see our response to reviewer #1 first point.**

Lastly, I am doubting whether this manuscript should be transferred into a technical note as suggested

by reviewer #1. In principle I agree, since the manuscript presents a method to define scaling issues and the 'research component' is a bit on the back. I leave it up to the authors to decide whether they prefer a technical note or a research article. However, if you chose the latter, please emphasize the research component (what can we learn from it, e.g., process understanding).

We still think that it is appropriate for this manuscript to be published as a research article. We explain the added value of the manuscript in response to the second last comment by reviewer#1, report#2.

**Response to Referee #1, report#2**

**We thank Reviewer #1 for her/his comments on the manuscript, and present our responses below. The Reviewer's comments are in regular font and our responses are in bold.**

**Response to Referee #1**

I now read the revised manuscript version and rebuttal letter of the manuscript entitled "Global assessment of how averaging over spatial heterogeneity in precipitation and potential evapotranspiration affects modeled evapotranspiration rates". I appreciate authors effort in improving the manuscript and clarifying several aspects of their work. Having said that, I must also admit that I am still not fully convinced by some of the arguments the authors used in the reply letter to present (and defend) their work. Therefore, I highlight below those important points that need to be addressed:

- Authors treat P and PET as model input data defined at the same scale (see lines 292-294). However, this does not reflect how typically ESMs dynamical cores are designed and/or have been evolving. Forcing terms (e.g., precipitation, temperature, humidity) are defined at the atmospheric model grid (usually coarser) while PET is calculated at the PFT-level using land surface features (e.g., LAI, aerodynamic resistance). In light of this, variability in modelled P and PET occurs at different spatial scales. This makes, in my opinion, the calculation of a correction factor for ET less straightforward. Could authors elaborate on this point?

**The introduction describes "mosaic" approaches in which PET (and ET) are calculated for individual PFTs and then aggregated. In any case, our purpose is not to mimic the way that ESMs actually calculate ET (obviously so, since ESMs do not employ Budyko curves). Our purpose is instead to illustrate how variability in P and PET would be translated into biased ET estimates, using Budyko curves as a simple "see-through" function for illustration purposes. In the case that the reviewer mentions (P and PET calculated at different scales), the magnitude of the bias would depend on whether the PET estimates were first averaged at the atmospheric grid cell scale (Case 1), or whether P and PET were jointly used to estimate ET for each PFT within each grid cell, and then these ET estimates were averaged (Case 2). The aggregation bias would be greater in Case 1 than in Case 2, but we don't think we would be justified in going into these details in the present paper, because, again, the Budyko calculations presented here are in any case not the way that ESMs actually calculate ET. We make sure we emphasize this point in the manuscript.**

- Authors affirm that their work highlights under which climate conditions averaging in P and PET has an effect on ET estimates. Without giving clear explanations I suspect that the largest differences found for Cs and Ds climate zones are still mainly driven by topography. Note also that this analysis is limited to the CONUS domain where there are not many sampling points for certain climate zones (e.g., Ds). This information should be provided in the plots and the statistical significance of the differences should be tested. An analysis at the global scale would be certainly more convincing.

**We agree that the large aggregation biases found in the Cs and Ds climate zones are mainly driven by topography We highlight this in the second paragraph of Section 4: "Heterogeneity biases are higher in regions with temperate climates and dry summers (climate zone Cs) and in regions with cold, dry summers (climate zone Ds), most likely due to the sharp spatial gradient in their water**

and energy sources for evapotranspiration (Fig. 5b). These areas typically have high topographic relief, combined with seasonal climate." We focused our analysis on the CONUS domain because we wanted to compare Prism and WorldClim as precipitation data sources, and fine resolution Prism data are only publicly available for CONUS. Thus while a global analysis would arguably be more comprehensive, it is not possible without acquiring (and paying for) proprietary Prism data.

The number of 1-degree by 1-degree grid cells (sampling points) at which heterogeneity biases are calculated per climate zones are now added to Figure 5b.

We present a table in the supplement in which we report statistical significance of differences between heterogeneity bias estimated at 1-degree by 1-degree grid cell across the contiguous US using 4 sets of P and PET data. The difference between heterogeneity bias estimated at the two climate zones that are raised by the reviewer (Cs and Ds) is not statistically significant across all 4 combinations of datasets (highlighted in yellow in Table S1 of the supplementary material). However, the difference between estimated heterogeneity bias in Cs vs Cf climate zones, and Ds versus Cf climate zones, as well as Cs versus Bs climate zones are significant across all four data combinations (highlighted in grey, blue, and green in Table S1 of the supplementary material). We discuss the main results of the statistical difference analysis in the manuscript (Section 4., second paragraph).

- The grid-scale dependence is tested for Switzerland and I don't think we can "extrapolate" this exponential relationship everywhere around the globe. If you want to convince the reader you need to repeat this assessment for all regions (identified for instance in Fig. 3) where averaging effects are not negligible. Juxtaposing the different curves will (or will not) support the existence of a general "scaling" relationship with the grid resolution.

**A global analysis would indeed be more comprehensive, but it would require high-resolution global data that we simply do not have. The graph of average heterogeneity bias versus grid resolution was added to figure 6 upon the reviewer's request in the previous round. In the manuscript, we report that "On average, the heterogeneity bias across Switzerland as a whole grows exponentially as the inputs are averaged over larger grids" and do not generalize it to any other region or the globe.**

- In the first iteration I asked authors to provide more information about the value of "n" parameter. This information is still missing in the manuscript. Can the authors provide some concrete numbers on the sensitivity of their global estimates with respect to different "n" values?

**This is an important point, thank you for raising it. We now present a sensitivity analysis for n values ranging from 2 to 5 in the supplement (Figures S1 and S2) and discuss its main results in the manuscript (Section 5. Summary and discussions, 5<sup>th</sup> paragraph). Those results are: 1) the spatial patterns of aggregation bias are similar, 2) the absolute magnitude of aggregation bias increases somewhat for higher values of n, as predicted by Eqs. 3 and 4, and 3) the Taylor approximation in Eqs. 3 and 4 yield realistic estimates of the aggregation bias for all values of n that were tested.**

Other comments:

- Please do not include any discussion in the captions of the figures. See Figure 3-5-6.

**We include concise statements of the main takeaway messages that the figures convey. We think that this is very helpful to readers – particularly those who scan the figures of a paper to get a first**

**impression of its main points. As this is a matter of style, we prefer to keep the captions as they are.**

- Lines 36-37 in the abstract. I do not see how the results of this paper can be used for guiding a more detailed mechanistic modelling. Note also that averaging- or grid-scale effects have been largely reported also when using mechanistic models. So please remove this sentence.

**This sentence is removed from the abstract.**

- You do not want to quantify the absolute magnitude of the averaging effects and at the same you claim that your methodology is potentially a way for correcting such bias. The second statement imply a sort of quantification, in my opinion.

**This is correct, and we said this in the first paragraph of Section 5. Obviously, to correct for aggregation biases one needs to quantify them. As the first paragraph of Section 5 explains, the general approach outlined here could be used to quantify and correct for aggregation biases – but it would need to be applied to the mechanistic ET equations that are actually used in ESMs, rather than the simple Budyko curves that we have used here for purposes of illustration.**

- I found an imbalance between the emphasis you put on the introduction and the actual findings of the manuscript. As I said in the first iteration, this is an applied study of a previously described methodology that does not contain general insights.

**The entire introduction, except for the last paragraph, sets up the general problem of heterogeneity bias and how it is typically handled in ESMs. This introduction is essential for readers who do not already know this material. Although we do indeed apply a previously described methodology (and the introduction is quite explicit about this), the present paper presents a series of new insights, including:**

- **Our previous work showed mathematically that averaging over spatially heterogeneous P and PET results in overestimation of ET within the Budyko framework. We did not, however, determine where around the globe, and under what conditions, this heterogeneity bias is likely to be most important. In this work, we examine the global distribution of this bias, its scale dependence, and its sensitivity to variations in P versus PET.**
- **Our goal is to identify where, under what conditions, and at what spatial scales averaging over heterogeneities in P and PET could be most important to estimates of evapotranspiration, but not to quantify the absolute magnitude of these averaging effects.**
- **Our work outlines a strategy for quantifying heterogeneity biases and potentially correcting for them, and highlights regions where more detailed mechanistic modeling is needed.**
- **Our analysis of percent variability of P and PET products shows that percent variabilities of precipitation products are in general larger than PET products and hence contribute more to heterogeneity bias.**
- **Our analyses show that mountainous terrain, regions with temperate climates and dry summers, and landscapes where spatial variations in precipitation and potential evapotranspiration are inversely correlated exhibit greater heterogeneity bias in ET estimates.**

- **Our analysis of scale dependence (using Switzerland as a test case) shows that heterogeneity bias in Switzerland increases almost exponentially as grid cell sizes increase.**

- Lines 317-319 (“most likely due to the sharp spatial gradient...”). This is a quite generic statement.

**In this sentence and the next one, we are making exactly the statement that the reviewer said needed to be made (that the high aggregation biases in the Cs and Ds climatic zones are largely attributable to topography).**

Dear Reviewer #3,

Thank you for your review and the detailed comments. Following please find our point by point response to your suggestions and questions. The Reviewer's comments are in regular font and our response is in bold.

### **Response to Referee #3 #report1**

Review of "Global assessment of how averaging over spatial heterogeneity in precipitation and potential evapotranspiration affects modeled evapotranspiration rates" by Elham Rouholahnejad Freund et al.

This is my first review of this work. The manuscript by Elham Rouholahnejad Freund et al. addresses the interesting issue of scaling of water and energy exchange at the land surface. While the manuscript focusses on a novel issue that could provide an interesting new addition to decades of literature on scaling, I do not find the manuscript in its current version to be convincing. This has to do with: a) a poor link between the main methodology and motivation as outlined in the Introduction, and b) a complete neglect of surface heterogeneity and scale-dependency in the Budyko n-parameter. In my view, the work could be a valuable contribution to HESS only if these deficiencies are addressed.

The work is motivated by potential scaling issues in ESMs. In my view, these relate mainly to effects of land surface heterogeneity (land use, soil type, groundwater tables, soil moisture, all of which can show large variability on small scales). To the degree these depend on forcing, this will to a large extent be caused by spatio-temporal variability of rainfall (i.e. convective storms leading to temporary wetting of part of a water-limited region only) and not just spatial variability. This is a big simplification where most of the scaling problem already is solved, and which is inherent to the choice for Budyko. A possible solution would be to solely focus on scaling issues within the Budyko framework due to P and PET. This would be a fairly novel approach, and it would avoid (artificially) linking too much to ESMs.

**Our analysis uses Budyko curves as a simple analytical framework (or a "see-through" function) to demonstrate our analysis. Our purpose is not to mimic the way that ESMs actually calculate ET (obviously so, since ESMs do not employ Budyko curves). Our purpose is instead to illustrate how variability in P and PET would be translated into biased ET estimates, using Budyko curves as a simple ET function for illustration purposes. Thus our purpose is not to highlight scaling issues within the Budyko framework per se.**

**We agree that Budyko curves already average over temporal heterogeneity (and the manuscript says this explicitly in the third paragraph of Section 3). It is unclear whether this temporal heterogeneity would lead to significant aggregation bias in ESMs, because they are usually solved on relatively short time steps.**

My second concern deals with the choice for a single Budyko n-parameter. Effectively, the authors show that at larger scales due to forcing heterogeneity, the Budyko curve tends to become more linear.

**This is indeed a consequence of our analysis, but it is not the point that we are trying to make. Our point is that any nonlinear function will yield averages that lie "inside" the curve, and for ET functions this will always be below the curve. We are just using Budyko curves to illustrate this point.**

But what is the motivation for the baseline choice of  $n$ ?

**For the calculations in the main paper, we used  $n=2$  because this is a commonly used value in the existing literature. We will add this detail to the manuscript.**

Where is it shown that this value corresponds better to observations (i.e. is a more valid model) at finer scales (1 km) than at coarser scales (1 degree)?

**We do not show this (and indeed we are not aware of any literature that does show it). This would require long-term catchment mass balances at the 1 km scale and at 1-degree scale, which are not widely available. In any case, our analysis is mainly concerned with the spatial pattern of aggregation bias (where is it larger? where is it smaller?) and this will not be particularly sensitive to the choice of  $n$ . We will add figures to the supplement where we compare aggregation bias calculations for different values of  $n$ .**

Would the results not strongly depend on the choice of  $n$ ? In reality,  $n$  will also very strongly depend on land use (see for example Fig. 1 in <https://www.hydrol-earth-syst-sci-discuss.net/hess-2018-634/#discussion> or any of the many other studies on this subject). I think any analysis of scaling should focus on the main nonlinearities to avoid becoming a purely academic exercise (nothing wrong with the latter, but then it should be presented as such). As a minimum, I would expect a sensitivity analysis on how the results depend on the value of  $n$ , accompanied by a discussion on how  $n$  might vary locally and with scale.

**We present this sensitivity analysis (for  $n$  values ranging from 2 to 5) in the supplement (Figs. S1 and s2) and discuss its main results in the manuscript. Those results are: 1) the spatial patterns of aggregation bias are similar, 2) the absolute magnitude of aggregation bias increases somewhat for higher values of  $n$ , as predicted by Eqs. 3 and 4, and 3) the Taylor approximation in Eqs. 3 and 4 yield realistic estimates of the aggregation bias for all values of  $n$  that were tested.**

Ideally, I would see the manuscript being restructured towards a more theoretical analysis of effects of forcing heterogeneity on the Budyko model, which would result in a clearly testable hypothesis that Budyko curves should become more linear with increasing scale as a result of heterogeneity, and based on the maps regions can be identified where this effect should be largest and best observable. It should also be noted that the value of  $n$  used in the analysis is not reported, at least I could not find it.

**As explained above, our paper is not intended as an analysis of scaling effects in Budyko curves. We do not think that it is a particularly interesting hypothesis that Budyko curves should become more linear with increasing scale, since this is generally true of all curved functions (it is basically a mathematical theorem rather than an empirical hypothesis). We will report the value of  $n$  that we used, in addition to the results of the sensitivity analysis covering a range of  $n$  values.**



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 The major goal of large-scale Earth System Models (ESMs) is to understand and predict global change. However,  
23 computational constraints require ESMs to operate on relatively large spatial grids (typically ~1 degree or ~100 km  
24 in size), 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 Budyko  
30 curves as simple functions that relate ET to precipitation (P) and potential evapotranspiration (PET). Because the  
31 relationships driving ET are nonlinear, averaging over sub-grid heterogeneity in P and PET will lead to biased  
32 estimates of average ET. We examine the global distribution of this bias, its scale dependence, and its sensitivity to  
33 variations in P versus PET. Our analysis shows that this "heterogeneity bias" is more pronounced in mountainous  
34 terrain, in landscapes where spatial variations in P and PET are inversely correlated, and in regions with temperate  
35 climates and dry summers. We also show that this heterogeneity bias increases on average, and expands over  
36 larger areas, as the grid cell size increases. ~~Our work outlines a strategy for quantifying heterogeneity biases and  
37 potentially correcting for them, and highlights regions where more detailed mechanistic modeling is needed.~~

38

39

## 40 1. Introduction

41 Earth System Models (ESMs) are designed to understand interactions between the land surface, atmosphere, and  
42 oceans and to predict global environmental changes. However, the Earth system and its underlying physical  
43 processes are highly heterogeneous across orders of magnitude in scale below the scale of typical ESM grids (e.g.,  
44 1° by 1°). Despite increasing recognition of the need to mechanistically represent physical processes in ESMs,  
45 currently even the most disaggregated large-scale ESMs cannot explicitly represent the spatial heterogeneity of  
46 land surface hydrological properties at scales that are important to atmospheric fluxes. Averaging over land surface  
47 properties at the scale of ESM model grid cells may have important implications for water and energy flux estimates  
48 (Avisar and Pielke, 1989; Giorgi and Avisar, 1997; Ershadi et al., 2013; Lu et al., 2014).

49  
50 Estimates of evapotranspiration (ET) fluxes have significant implications for future temperature predictions. Smaller  
51 ET fluxes imply greater sensible heat fluxes and, therefore, drier and warmer conditions in the context of climate  
52 change (Seneviratne et al., 2010). Surface evaporative fluxes (and thus energy partitioning over land surfaces) are  
53 nonlinear functions of available water and energy, and thus are coupled to spatially heterogeneous surface  
54 characteristics (e.g., soil type, vegetation, topography) and meteorological inputs (e.g., radiative flux, wind, and  
55 precipitation; Kalma et al., 2008; Shahraeeni and Or, 2010; Holland et al., 2013). These characteristics are spatially  
56 variable on length scales of <1 m to many kilometers, well below typical ESM grid scales of ~100 km. ESMs calculate  
57 grid-averaged surface and atmospheric fluxes using parameterizations that correspond to grid-averaged properties  
58 of the land surface (Sato et al., 1989; Koster et al., 2006; Santanello and Peters-Lidard, 2011). Thus ET estimates  
59 that are derived from spatially-averaged land surface properties do not capture ET variations driven by the  
60 underlying surface heterogeneity (McCabe and Wood, 2006). Because the relationships driving ET are nonlinear,  
61 the average ET flux from a heterogeneous landscape may be different from an ET estimate calculated from spatially  
62 averaged inputs (Rouholahnejad Freund and Kirchner, 2017).

63  
64 Several studies have quantified the effects of land surface heterogeneity on potential evapotranspiration (PET) and  
65 latent heat (LH) fluxes, and have found that averaging over land surface heterogeneity can potentially bias ET  
66 estimates either positively or negatively. For example, Boone and Wetzel (1998) studied the effects of soil texture  
67 variability within each pixel in the Land-Atmosphere-Cloud Exchange (PLACE) model, which has a spatial resolution  
68 of approximately 100 by 100 km. They reported that accounting for sub-grid variability in soil texture reduced  
69 global ET by 17%, increased total runoff by 48%, and increased soil wetness by 19%, compared to using a  
70 homogenous soil texture to describe the entire grid cell. Kollet (2009) found that heterogeneity in soil hydraulic  
71 conductivity had a strong influence on evapotranspiration during the dry months of the year, but not during  
72 months with sufficient moisture availability. Hong et al. (2009) reported that aggregating radiance data from 30 m  
73 to 60, 120, 250, 500, and 1000 m resolution (input upscaling) and then calculating ET from these aggregated inputs  
74 at these grid scales using Surface Energy Balance Algorithm for Land (SEBAL, Bastiaanssen et al., 1998a) yields  
75 slightly larger ET estimates as compared to ET calculated with finer resolution inputs and then aggregated at the

76 desired grid scales (output upscaling). The discrepancy between ET estimated with the output upscaling method  
77 and the input upscaling method grows as the size of the grid cell increases (the difference between ET calculated  
78 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  
79 of 120 m by 120 m). Aminzadeh et al. (2017) investigated the effects of averaging surface heterogeneity and soil  
80 moisture availability on potential evaporation from a heterogeneous land surface including bare soil and vegetation  
81 patches. They found that if the heterogeneity length scale is smaller than the convective atmospheric boundary  
82 layer (ABL) thickness, averaging over heterogeneous land surfaces has only a small effect on average potential  
83 evaporation rates. Averaging over larger-scale heterogeneities, however, led to overestimates of potential  
84 evaporation.

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

97  
98 Contrary to overestimation bias, many remotely sensed ET estimates that include parameters related to  
99 aerodynamic resistance are significantly affected by heterogeneity, and underestimate ET as the scale increases  
100 (Ershadi et al., 2013). Because aerodynamic resistance is significantly affected by land surface properties (e.g.,  
101 vegetation height, roughness length, and displacement height), decreases in aerodynamic resistance at coarser  
102 resolutions could lead to smaller estimates of evapotranspiration. Ershadi et al. (2013) showed that input  
103 aggregation from 120m to 960 m in Surface Energy Balance System (SEBS, Su, 2002) leads to up to 15 %  
104 underestimation of ET at the larger grid resolution in a study area in the south-east of Australia.

105  
106 Rouholahnejad Freund and Kirchner (2017) quantified the impact of sub-grid heterogeneity on grid-average ET  
107 using a simple Budyko curve (Turc, 1954; Mezentsev, 1955) in which long-term average ET is a non-linear function  
108 of long-term averages of precipitation (P) and potential evaporation (PET). They showed mathematically that  
109 averaging over spatially heterogeneous P and PET results in overestimation of ET within the Budyko framework (Fig.  
110 1). Their analysis implies that large-scale ESMs that overlook land surface heterogeneity will also yield biased

111 evapotranspiration estimates due to the inherent nonlinearity in ET processes. They did not, however, determine  
112 where around the globe, and under what conditions, this heterogeneity bias is likely to be most important.

113  
114 The recognition that spatial averaging can potentially lead to biased flux estimates has prompted methods for  
115 representing sub-grid-scale heterogeneities and processes within [large scale land surface models and](#) ESMs.  
116 Accounting for land surface heterogeneity in large-scale ESMs is not merely constrained by limitations in both  
117 computational power (Baker et al. 2017) and the availability of high-resolution forcing data, but also by the fact  
118 that the atmospheric and land surface components of some ESMs operate at different resolutions. There have been  
119 several attempts to integrate sub-grid heterogeneity in ESMs while keeping the computational costs affordable. In  
120 “mosaic” approaches, the model is run separately for each surface type in a grid cell, and then the surface-specific  
121 fluxes are area-weighted to calculate the grid-cell average fluxes (e.g., Avissar and Pielke, 1989; Koster and Suarez,  
122 1992). The “effective parameter” approach (e.g., Wood and Mason, 1991; Mahrt et al., 1992), by contrast, seeks to  
123 estimate effective parameter values at the grid cell scale that subsume the effects of sub-grid heterogeneity.  
124 Estimating these effective parameters can be challenging because the relevant land-surface processes typically  
125 depend nonlinearly on multiple interacting parameters, and land-surface signals at different scales are propagated  
126 and diffused differently in the atmosphere. Alternatively, the “correction factor” approach (e.g., Maayar and Chen,  
127 2006) uses sub-grid information on spatially heterogeneous land-surface processes and properties to estimate  
128 multiplicative correction factors for fluxes that are originally calculated from spatially averaged inputs at the grid-  
129 cell scale. All three approaches try to reduce the heterogeneous problem to a homogeneous one that has  
130 equivalent effects on the atmosphere at the grid-cell scale.

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

143

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

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

$$148 \quad ET = f(P, PET) = \frac{P}{\left(\left(\frac{P}{PET}\right)^n + 1\right)^{1/n}} \quad (1)$$

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

151  
 152 Evapotranspiration rates are inherently bounded by energy and water limits. Under arid conditions ET is limited by  
 153 the available supply of water (the water limit line in Fig. 1b), while under humid conditions ET is limited by  
 154 atmospheric demand (PET) and converges toward PET (the energy limit line in Fig. 1b). Budyko showed that over a  
 155 long period and under steady-state conditions, hydrological systems function close to their energy or water limits.  
 156 These intrinsic water and energy constraints make the Budyko curve downward-curving.

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

167  
 168 In a previous study (Rouholahnejad Freund and Kirchner, 2017) we showed that when nonlinear underlying  
 169 relationships are used to predict average behaviour from averaged properties, the magnitude of the resulting  
 170 heterogeneity bias can be estimated from the degree of the curvature in the underlying function and the range  
 171 spanned by the individual data being averaged. Here we summarize these findings as building blocks of the current  
 172 study. The second-order, second-moment Taylor expansion of the ET function  $f(P, PET)$  (Eq. 1) around its mean  
 173 directly yields:

$$174 \quad \bar{f}(P, PET) = \overline{ET} \approx f(\bar{P}, \overline{PET}) + \frac{1}{2} \frac{\partial^2 f}{\partial P^2} var(P) + \frac{1}{2} \frac{\partial^2 f}{\partial PET^2} var(PET) + \frac{\partial^2 f}{\partial P \partial PET} cov(P, PET) \quad , \quad (2)$$

175 where  $\bar{f}(P, PET)$  is the true average of the spatially heterogeneous ET function,  $f(\bar{P}, \overline{PET})$  is the ET function  
 176 evaluated at its average inputs  $\bar{P}$  and  $\overline{PET}$ , and the derivatives are calculated at  $\bar{P}$  and  $\overline{PET}$ . Evaluating the  
 177 derivatives using Eq. (1) and reshuffling the terms, Rouholahnejad Freund and Kirchner (2017) obtained the

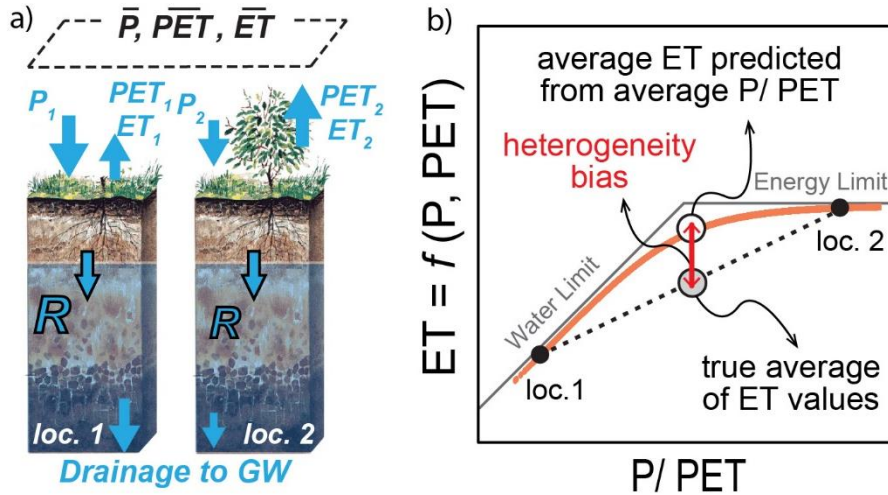
178 following expression for the heterogeneity bias, the difference between the average ET,  $\bar{f}(P, PET)$ , and the ET  
 179 function evaluated at the mean of its inputs,  $f(\bar{P}, \bar{PET})$ :

$$180 \quad f(\bar{P}, \bar{PET}) - \bar{f}(P, PET) \approx (n + 1) \frac{\bar{P}^{n+1} \bar{PET}^{n+1}}{(\bar{P}^n + \bar{PET}^n)^{2+1/n}} \left[ \frac{1}{2} \frac{var(P)}{\bar{P}^2} + \frac{1}{2} \frac{var(PET)}{\bar{PET}^2} - \frac{cov(P, PET)}{\bar{P} \bar{PET}} \right]. \quad (3)$$

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

$$182 \quad (n + 1) \frac{\bar{P}^{n+1} \bar{PET}^{n+1}}{(\bar{P}^n + \bar{PET}^n)^{2+1/n}} \left[ \frac{1}{2} \left( \frac{SD(P)}{\bar{P}} \right)^2 + \frac{1}{2} \left( \frac{SD(PET)}{\bar{PET}} \right)^2 - r_{P, PET} \left( \frac{SD(P)}{\bar{P}} \right) \left( \frac{SD(PET)}{\bar{PET}} \right) \right]. \quad (4)$$

183 Equation (4) shows that the heterogeneity bias depends on only four quantities: the fractional variation (i.e., the  
 184 coefficient of variation) in precipitation  $\left( \frac{SD(P)}{\bar{P}} \right)$  and in potential ET  $\left( \frac{SD(PET)}{\bar{PET}} \right)$ , the correlation between precipitation  
 185 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  
 186 in Budyko space. As shown by Fig. 1b and Eq. (2), the discrepancy between average of the ET function and the ET  
 187 function of the average inputs (the heterogeneity bias) is proportional to both the degree of nonlinearity in the  
 188 function, as defined by its second derivatives, and the variability of P and PET. Equation (4) allows one to estimate  
 189 how much the curvature of the ET function and the fractional variability (standard deviation divided by mean) of P  
 190 and PET will affect estimates of ET. However, to the best of our knowledge, the consequences of these  
 191 nonlinearities for global evaporative flux estimates have not previously been quantified.  
 192



193  
 194 Figure 1. Heterogeneity bias in a hypothetical two-column model in the Budyko framework. The true average ET of  
 195 the columns (gray circle) lies below the curve and is less than the average ET estimated from the average P/PET of  
 196 the two columns (open circle). The heterogeneity bias depends on the curvature of the function and the spread of  
 197 its inputs. Both panels are adapted from Rouholahnejad Freund and Kirchner (2017).

198

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

200 Across a landscape of similar size to a typical ESM grid cell (1° by 1°), soil moisture, atmospheric demand (PET) and  
201 precipitation (P) will vary with topographic position; hillslopes will typically be drier, and riparian regions will be  
202 wetter. To map the spatial pattern in the heterogeneity bias that results from averaging over this land surface  
203 heterogeneity, we applied the approach outlined in section 2 to the global land surface area at 1° by 1° grid scale.  
204 Within each 1° by 1° grid cell, we used 30 arc-second values of P (WorldClim; Hijmans et al., 2005) and PET  
205 (WorldClim; Hijmans et al., 2005) to examine the variations in small-scale climatic drivers of ET. Because 30 arc-  
206 seconds is nearly 1 km, hereafter we refer to the 30 arc-second data as 1km values for simplicity. The spatial  
207 distribution of long-term annual averages (1960-1990) of P and PET values at 1 km resolution, along with 1km  
208 values of the aridity index (AI=P/PET), are shown in Fig 2a-c. ET values calculated from these 1km P and PET values  
209 using Eq. (1) are then averaged at 1° by 1° scale (“true average”, Fig. 2e). We also averaged the 1km values of P and  
210 PET within each grid cell and then modeled ET using the Budyko curve (Eq. 1) applied to these averaged input  
211 values. The difference between these two ET estimates is the heterogeneity bias.

212

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

217

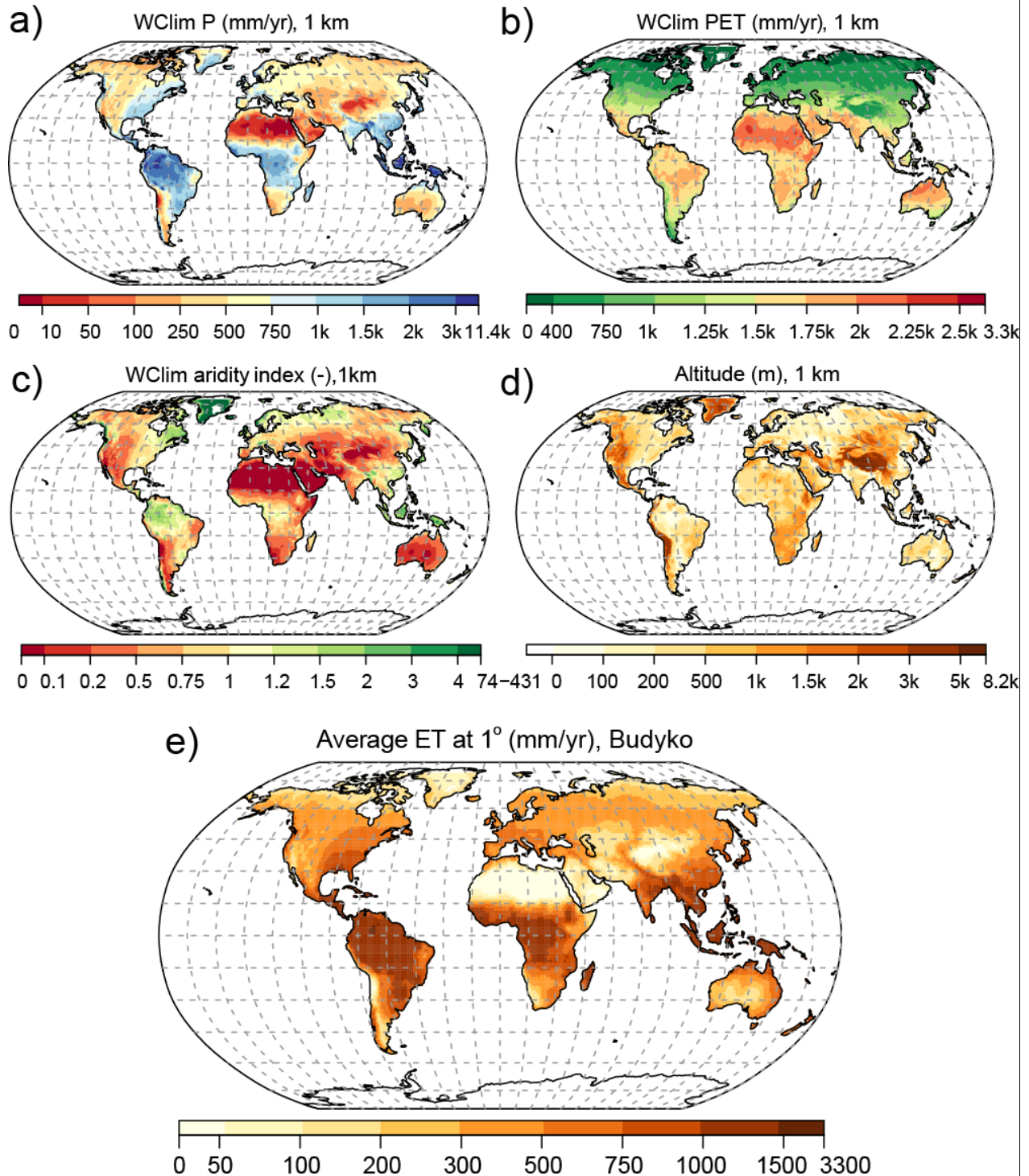
218 Fig. 3a-d illustrates the variability (quantified by standard deviation) of 1km values of P, PET, aridity index, and  
219 altitude at the 1° by 1° grid scale. The heterogeneity bias in long-term average ET fluxes at the 1° by 1° grid scale  
220 (Fig. 3e) highlights regions around the globe where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

230

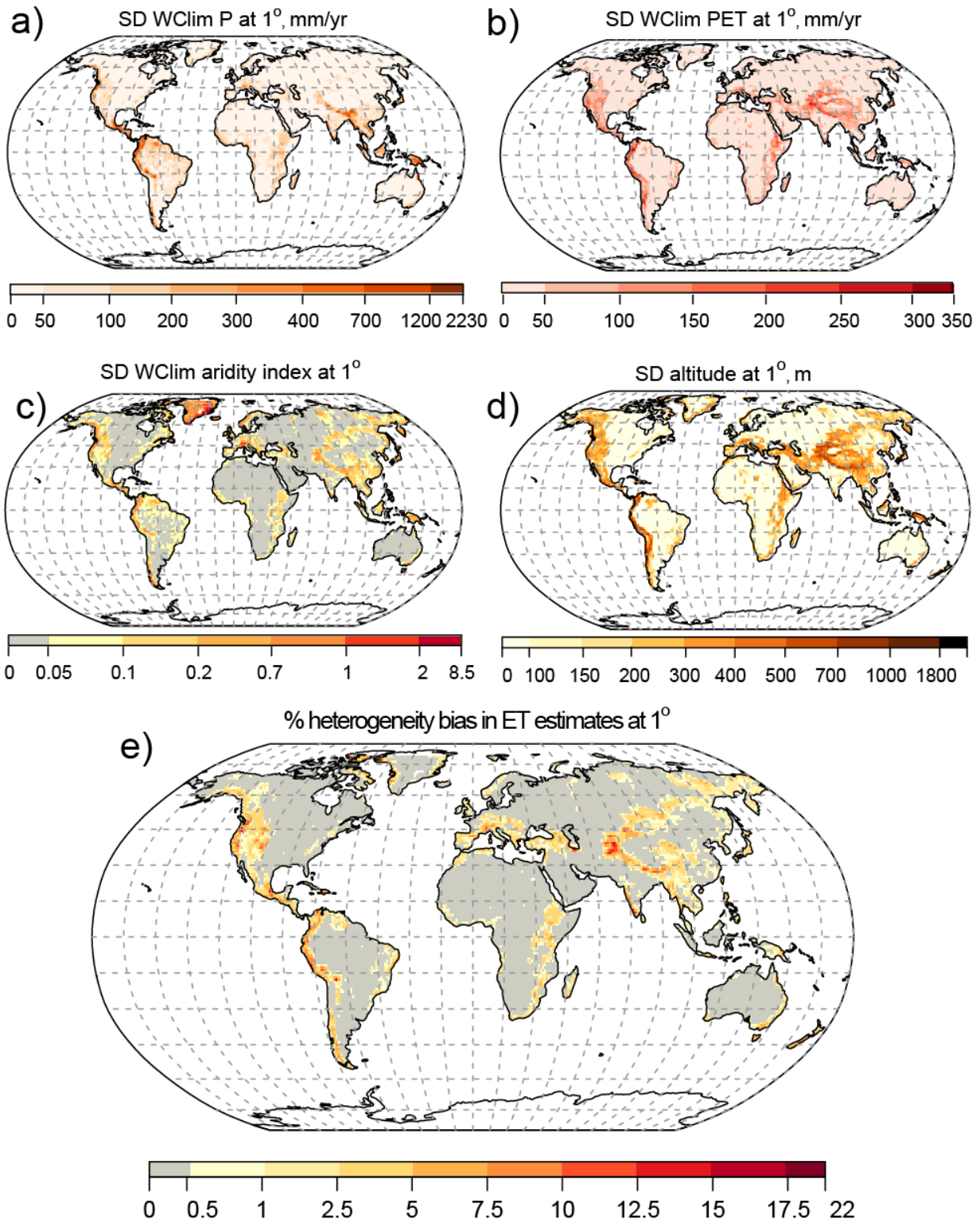
231 Our results show that the topographic gradient, and hence the variability in the aridity index across a given grid  
232 scale, drives consistent, predictable patterns of heterogeneity bias in evapotranspiration estimates at that scale.  
233 Equation 4 shows that this bias is equally sensitive to fractional variability in P and PET (standard deviation divided



234 by mean). However, because P is typically more variable (in percentage terms) than PET across landscapes, the  
 235 variability in P will usually make a larger contribution to the heterogeneity bias.  
 236



237  
 238 Figure 2. Global distribution of one-kilometer resolution annual mean precipitation (a: P; WorldClim; Hijmans et al.,  
 239 2005), potential evapotranspiration (b: PET; WorldClim; Hijmans et al., 2005), aridity index (c: AI=P/PET; WorldClim;  
 240 Hijmans et al., 2005), and topography (d: SRTM; Jarvis et al., 2008), along with (e) evapotranspiration (ET) at 1° by  
 241 1° scale by averaging 1km values of ET calculated using the Budyko function (Eq. 1).



243

244 Figure 3. Global spatial distribution of variability (standard deviation) of one-kilometer values of a) precipitation (P),

245 b) potential evapotranspiration (PET), c) aridity index (AI=P/PET), and d) altitude at 1° by 1° grid cell. The

246 heterogeneity bias in ET estimates (e) is calculated using Eq. (4). Grid cells with larger standard deviation in altitude  
247 and aridity index have larger heterogeneity bias.

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

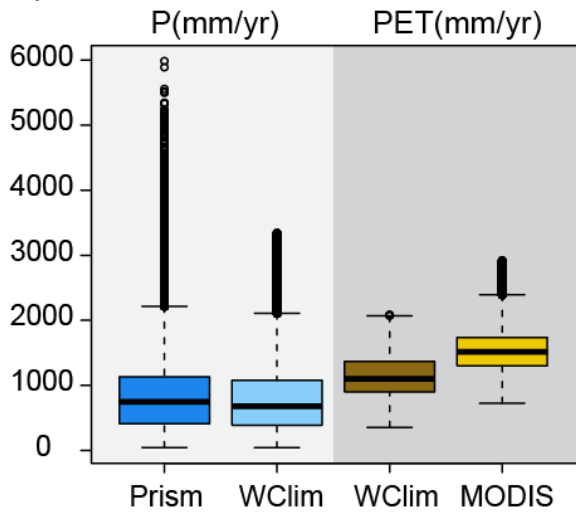
249 With increased availability of spatial data, it is becoming standard practice to assess input data uncertainties and  
250 their propagated impacts on water and energy flux estimates in land surface models. To quantify how choices  
251 among alternative input data products could affect the heterogeneity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

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

282 which resulted in the second-highest heterogeneity bias across different climate zones. [Wilcoxon signed-rank tests](#)  
283 [was performed to evaluate the statistical significance of the differences between heterogeneity bias in ET estimates](#)  
284 [using all pairs of climate zones and data sources that are shown in Fig. 5b \(Table S1\)](#). These analysis show that while  
285 [the difference between heterogeneity biases estimated in Cs and Ds climate zones are not statistically significant](#)  
286 [across all four combinations of datasets, the difference between estimated heterogeneity bias in Cs versus Cf, Ds](#)  
287 [versus Cf, as well as Cs versus Bs climate zones are significant across all four data combinations \(highlighted in Table](#)  
288 [S1 of the supplementary material\)](#).

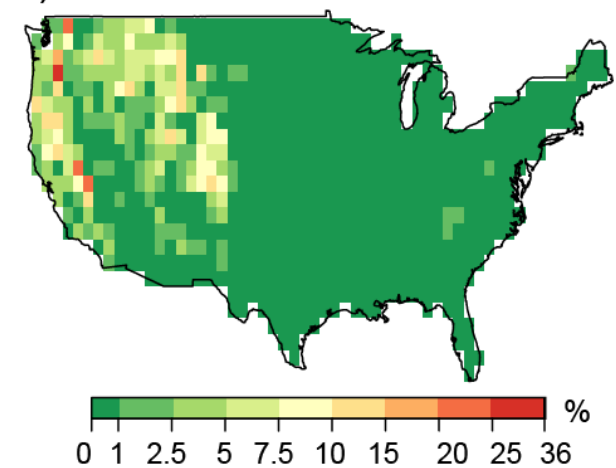
289  
290 Equation 4 shows that heterogeneity biases in Budyko estimates of ET are equally sensitive to the same percentage  
291 variability in P and PET. Thus the degree of sensitivity, per se, to P and PET variations expressed in percentage terms  
292 is the same. Although Figs. 5c and 5d give the visual impression that PET is more variable than P across climate  
293 zones and between data sources, Fig. 5e shows that the fractional variability in P is systematically higher than PET,  
294 and it also varies more across the climate zones and between the two data sets. Because P is typically more  
295 variable than PET (in percentage terms) across landscapes, the variability in P will make a larger contribution to the  
296 heterogeneity bias (Fig. 5e) in the Budyko approach. Whether this is true for more physically based ET estimates  
297 remains to be seen. Analysis of percent variability of P and PET products shows that percent variabilities of  
298 precipitation products are in general larger than PET products and hence contribute more to heterogeneity (Fig 5e).  
299 While the percent variabilities of the two PET products are in the same range, the percent variability in Prism  
300 precipitation is slightly larger than in WorldClim precipitation, in regions with dry summers (Cs and Ds climate zones  
301 in Fig. 5a).  
302

a) Distribution of P and PET in the four datasets



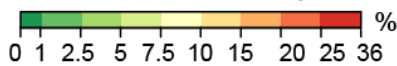
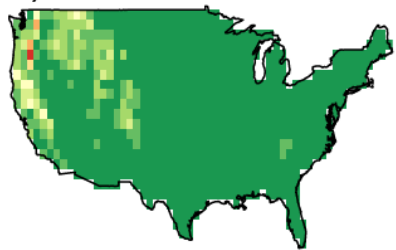
% heterogeneity bias in ET estimates at 1°

b) Prism P, Wclim PET as inputs

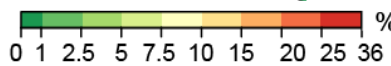
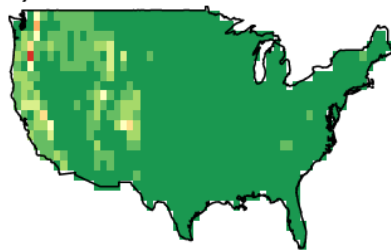


% heterogeneity bias in ET estimates at 1°

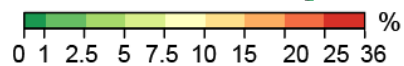
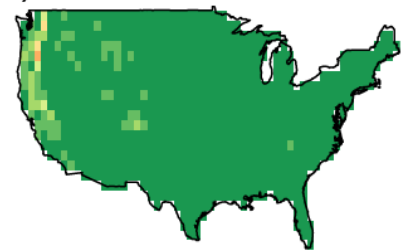
c) Prism P, MODIS PET as inputs



d) Wclim P, Wclim PET as inputs



e) Wclim P, MODIS PET as inputs



303

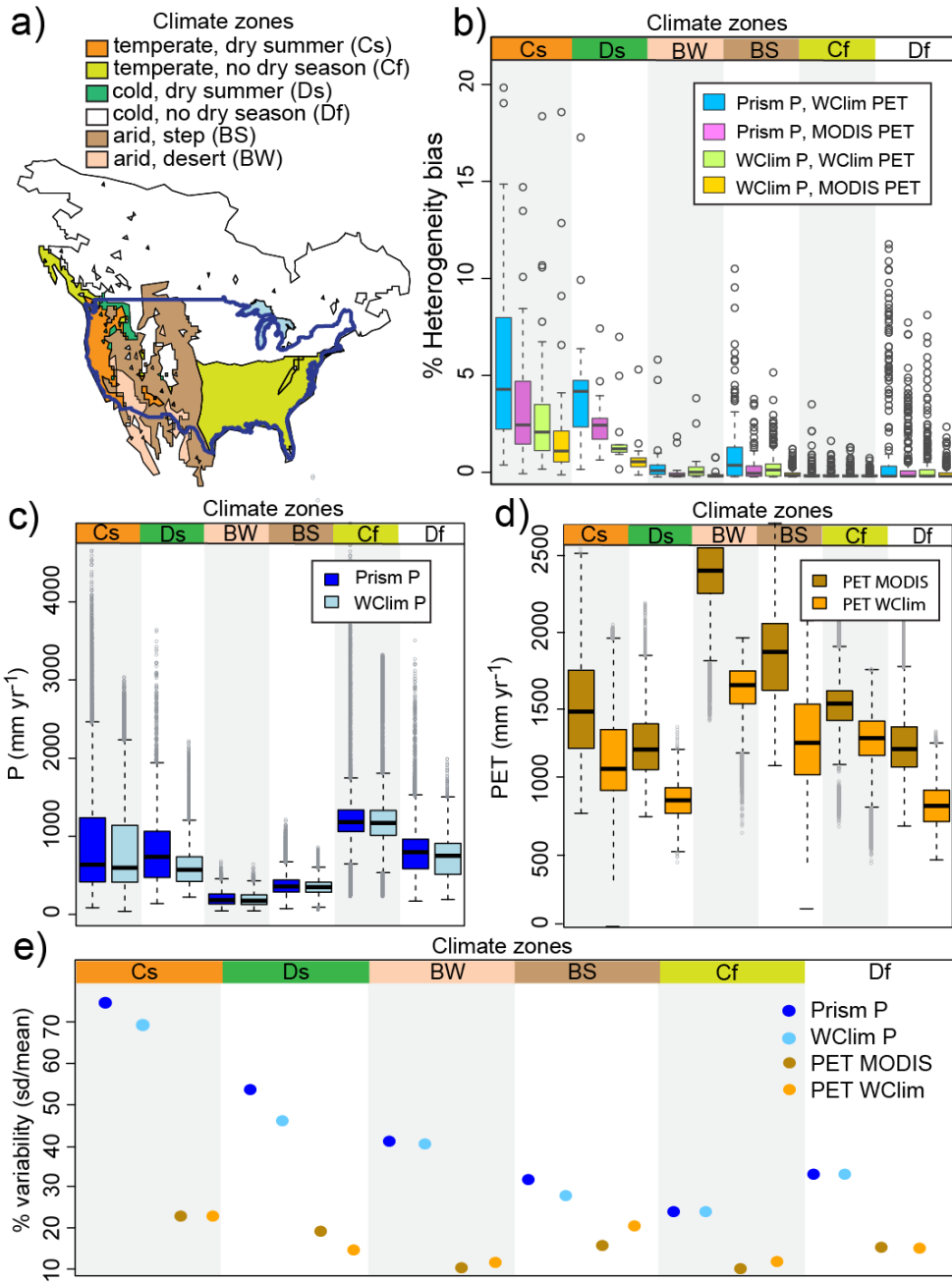
304

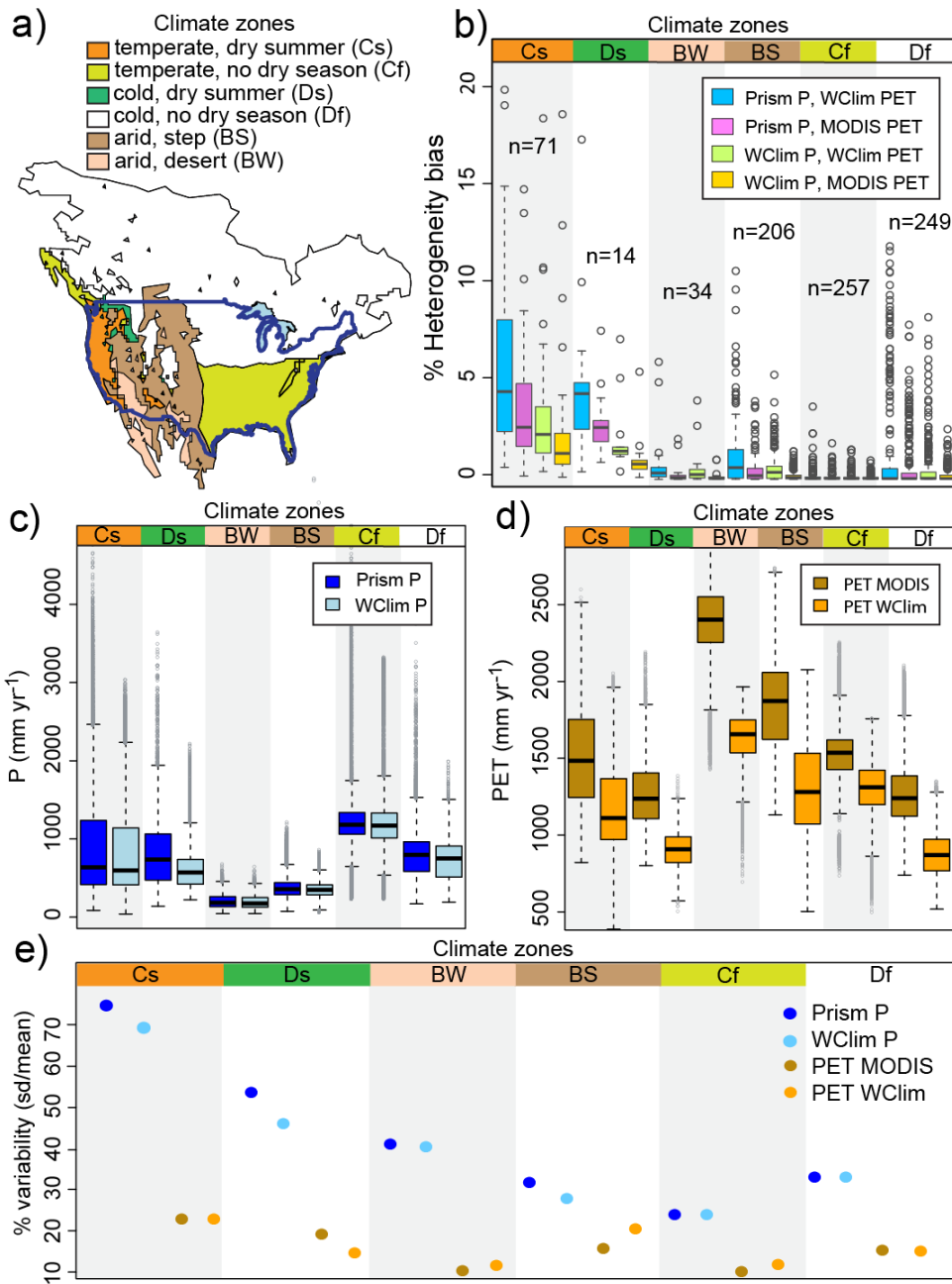
305

306

307

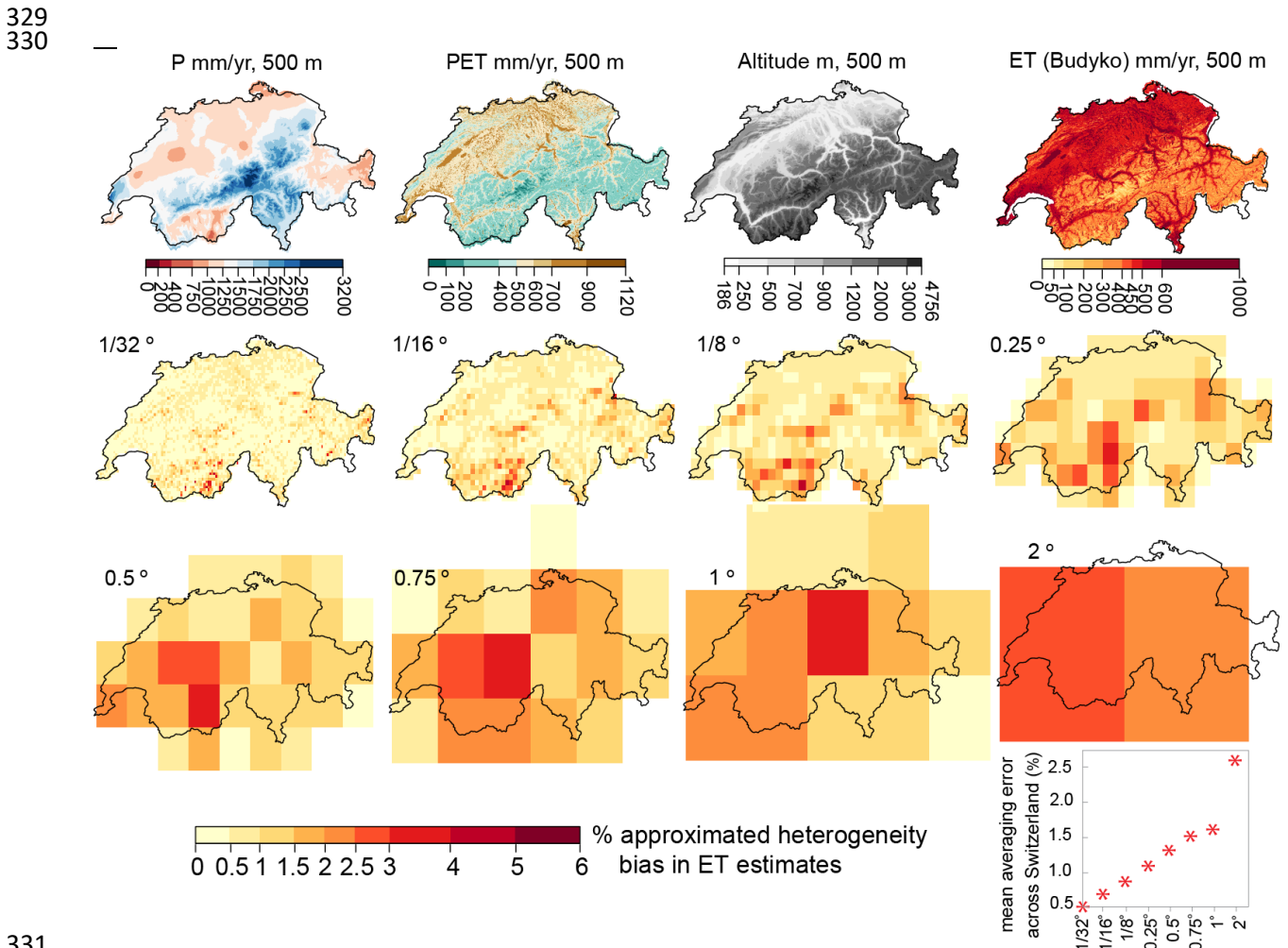
Figure 4. The distribution of P and PET in the four datasets is shown in a). Estimated heterogeneity bias (Eq. 4) across the contiguous US using [onefour](#)-kilometer values of b) Prism P and WorldClim PET c) Prism P and MODIS PET d) WorldClim P and WorldClim PET, and e) WorldClim P and MODIS PET as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





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

319 Because future increases in computing power will lead to ESMs with smaller grid cells, it is useful to ask how  
 320 changes in grid resolution affect the heterogeneity biases that we have estimated in this paper. To quantify the  
 321 heterogeneity bias in ET estimates as a function of grid scale, we repeated our analysis at various grid resolutions  
 322 using Switzerland as a test case. We started with high-resolution (500m) maps of long-term average annual  
 323 precipitation and PET across the Swiss landscape (Fig. 6), and then used Eq. 4 to estimate the heterogeneity bias at  
 324 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  
 325 leads to larger, and more widespread, overestimates in ET. Conversely, at finer grid resolutions, the average  
 326 heterogeneity bias is smaller, and the locations with large biases are more localized. On average, the heterogeneity  
 327 bias across Switzerland as a whole grows exponentially as the inputs are averaged over larger grids (as shown in the  
 328 lower-right panel in Fig. 6).



332 Figure 6. Heterogeneity bias in ET estimates at various scales across Switzerland, estimated from 500m climate  
 333 data. ET is calculated using the Budyko relationship (Eq. 1). Heterogeneity bias was estimated from 500m  
 334 precipitation (P) and potential evapotranspiration (PET), and their variances at each grid scale, using Eq. 4. At finer  
 335 grid resolutions, the heterogeneity bias is more localized, and smaller on average.  
 336



## 337 5. Summary and discussion

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

352  
353 In this study, we used Budyko curves as simple models of ET, in which long-term average ET rates are functionally  
354 related to long-term averages of P and PET. We used an approach outlined by Rouholahnejad Freund and Kirchner  
355 (2017) to estimate the heterogeneity bias in modeled ET at 1-degree grid scale across the globe (Fig. 3), and also at  
356 multiple grid scales across Switzerland (Fig. 6), using finer-resolution P and PET values as drivers of ET. We showed  
357 how the heterogeneity effects on ET estimates vary with the nonlinearity in the governing equations and with the  
358 variability in land surface properties. Our analysis shows that heterogeneity effects on ET fluxes matter the most in  
359 areas with sharp gradients in the aridity index, which are in turn controlled by topographic gradients, and not  
360 merely in areas that are either arid or humid (e.g., compare Fig. 3e with Fig. 2c).

361  
362 According to our analysis, regions within the U.S. that have temperate climates and dry summers exhibit greater  
363 heterogeneity bias in ET estimates (Fig. 5). We show that the heterogeneity bias in ET estimates at each grid scale  
364 depends on the variance in the drivers of ET at that scale (Fig. 4), and on the choice of data sources used to  
365 estimate ET. Heterogeneity bias was significantly larger across the contiguous United States when P and PET data  
366 sources with larger variances were used (Fig. 4).

367  
368 We also explored the magnitude and spatial distribution of heterogeneity bias in ET estimates as a function of the  
369 scale at which the climatic drivers of ET are averaged. We found that as heterogeneous climatic variables are  
370 aggregated to larger scales, the heterogeneity biases in ET estimates become greater on average, and extend over  
371 larger areas (Fig. 6). At smaller grid scales, the heterogeneity bias does not completely disappear, but instead  
372 becomes more localized around areas with sharp topographic gradients. Finding an effective scale at which one can

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

380  
381 In the analysis presented here, we have assumed a value of 2 for the Budyko parameter  $n$ , which approximates the  
382 variation of ET/PET with respect to P/PET in MODIS and WorldClim data across continental Europe (Mu et al. 2007;  
383 Hijmans et al. 2005; Rouholahnejad Freund & Kirchner, 2017). Although there are suggestions in the literature that  
384  $n$  can vary with land use and other landscape properties (e.g., Teuling et al., 2019), here we have assumed that  $n$  is  
385 spatially and temporally constant in order to focus on the effects of heterogeneity in P and PET. In the supplement  
386 we present a sensitivity analysis with values of  $n$  ranging from 2 to 5 (Fig. S1). That analysis shows that, as expected  
387 from Eqs. 3 and 4, higher values of  $n$  lead to larger heterogeneity biases, but the spatial pattern shown in Fig. 3e  
388 remains largely unchanged. The Taylor approximation in Eqs. 3 and 4 yields realistic estimates of the heterogeneity  
389 bias for all values of  $n$  that were tested (Fig. S2). Thus while our numerical estimates of heterogeneity bias depend  
390 somewhat on the value of  $n$ , our conclusions do not.

391  
392 One should keep in mind that the true mechanistic equations that determine point-scale ET as a function of point-  
393 scale water availability and PET (if such data were available) may be much more nonlinear than Budyko's empirical  
394 curves, because these curves already average over significant spatial and temporal heterogeneity. Thus, we expect  
395 that the real-world effects of sub-grid heterogeneity are probably larger than those we have estimated in Sects. 3  
396 and 4 of this study. In addition, the 1km P and PET values that are used in our global analysis might be still too  
397 coarse to represent small-scale heterogeneity that is important to evapotranspiration processes.

398  
399 Budyko curves are empirical relationships that functionally relate evaporation processes to the supply of water and  
400 energy under steady-state conditions in closed catchments with no changes in storage. Our analysis likewise  
401 assumes no changes in storage, nor any lateral transfer between the model grid cells, although both lateral  
402 transfers and changes in storage may be important, both in the real world and in models. Unlike the Budyko  
403 framework, ET fluxes in most ESMs are often physically based (not merely functions of P and PET) and are  
404 calculated at much smaller time steps (seconds to minutes). These models often represent more processes that are  
405 important to evapotranspiration (such as storage variations) and include their dynamics to the extent that is  
406 computationally feasible. Because these relationships may be much more nonlinear than Budyko curves, there may  
407 also be significant heterogeneity biases when complex physically based models are used to estimate ET from

408 spatially aggregated data. Therefore, we are now working to quantify heterogeneity bias in ET fluxes using a more  
409 mechanistic land surface model.

410

#### 411 **Acknowledgements**

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

415

#### 416 **References**

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

419 Avissar, R., R. A. Pielke: A Parameterization of Heterogeneous Land Surfaces for Atmospheric Numerical Models and  
420 Its Impact on Regional Meteorology, *Monthly Weather Review*, vol. 117, issue 10, p. 2113, doi:10.1175/1520-  
421 0493(1989)117<2113:APOHLS>2.0.CO;2, 1989.

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

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

427 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:  
428 Global patterns in base flow index and recession based on streamflow observations from 3394 catchments, *Water  
429 Resour. Res.*, 49, 7843–7863, doi:10.1002/2013WR013918, 2013.

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

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

433 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.  
434 Mackay, R. M. Maxwell, C. Shen, S. C. Swenson, and X. Zeng: Improving the representation of hydrologic processes  
435 in Earth System Models, *Water Resour. Res.*, 51, 5929–5956, doi:10.1002/2015WR017096, 2015.

436 Ershadi A., M. F. McCabe, J. P. Evans, J. P. Walker: Effects of spatial aggregation on the multi-scale estimation of  
437 evapotranspiration, *Remote Sensing of Environment* 131, 51–62, <http://dx.doi.org/10.1016/j.rse.2012.12.007>,  
438 2013.

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

445 Giorgi, F., and R. Avissar: Representation of heterogeneity effects in Earth system modeling: Experience from land  
446 surface modeling, *Rev. Geophys.*, 35, 413–437, doi:10.1029/97RG01754, 1997.

447 Hargreaves, G. H., and Z. A. Samani: Reference crop evaporation from temperature, *Appl. Eng. Agric.*, 1(2), 96-99,  
448 1985.

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

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

453 Hong, S. H., J. M. H. Hendrickx, and B. Borchers: Up-scaling of SEBAL derived evapotranspiration maps from Landsat  
454 (30 m) to MODIS (250 m) scale, *Journal of Hydrology*, 370, 122–138, 2009.

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

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

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

461 Koster R. D. et al.: GLACE: The Global Land– Atmosphere Coupling Experiment. Part I: Overview. *J. Hydrometeorol.*, 7,  
462 590–610, 2006.

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

465 Lu, H., T., Liu, Y. Yang, D. Yao: A hybrid dual-Source model of estimating evapotranspiration over different  
466 ecosystems and implications for satellite-based approaches, *Remote Sens.* 6, 8359–8386, 2014.

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

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

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

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

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

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

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

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

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

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

486 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  
487 Implementing the Simple Biosphere Model in a General Circulation Model, *J. Atmospheric Sciences*, 46(18), 2757-  
488 2782, 1989.

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

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

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

495 [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](#)  
496 [Sterling, S. M., Climate change, reforestation/afforestation, and urbanization impacts on evapotranspiration and](#)

497 [streamflow in Europe, Hydrology and Earth System Sciences, 23, 3631—3652, DOI = {10.5194/hess-23-3631-2019,](#)  
498 [2019.](#)

499 Turc, L.: Le bilan d'eau des sols: relation entre la precipitations, l'évaporation et l'écoulement, Ann. Agron. A, 5,  
500 491–569, 1954.

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