

Manuscript Number: hess-2018-430

Title: Bayesian performance evaluation of evapotranspiration models: a case study based on eddy covariance system of a maize field in northwestern China

Corresponding Author: Xiaoying Zhang

Authors: GuoxiaoWei, Xiaoying Zhang, Ming Ye , Ning Yue, Fei Kan

Dear Editor,

On behalf of my co-authors, we thank you very much for giving us an opportunity to revise our manuscript. We are grateful to the editors and reviewers for their positive and constructive comments and suggestions on our manuscript (hess-2018-430) entitled “Bayesian performance evaluation of evapotranspiration models for an arid region in northwestern China”.

We have studied reviewer’s comments carefully, and revised the manuscript thoroughly to address the comments. The revision is marked in red in the revised manuscript. We have tried our best to revise our manuscript according to the comments. Attached please find the revised version, which we would like to submit for your kind consideration.

We would like to express our great appreciation to you and reviewers for comments on our paper. Looking forward to hearing from you.

Thank you and best regards.

Yours sincerely,

Xiaoying Zhang

Corresponding author:

Name Xiaoying Zhang

E-mail: xiaoyingzh@jlu.edu.cn

List of Responses

Dear Editors and Reviewers:

Thank you for your letter and for the reviewers' comments concerning our manuscript entitled "Bayesian performance evaluation of evapotranspiration models for an arid region in northwestern China" (hess-2018-430). These comments are all valuable and very helpful for revising and improving our paper, as well as the important guiding significance to our researches. We have studied comments carefully and have made corrections which we hope meet with approval. Revised portion are marked in red in the paper. The main corrections in the paper and the responds to the reviewer's comments are followed.

Thank you and best regards.

Yours sincerely,

Xiaoying Zhang

Corresponding author:

Name: Xiaoying Zhang

E-mail: xiaoyingzh@jlu.edu.cn

Responds to the comments:

Referee #1:

General comments:

1. Comment: *For the objective (1), what is the purpose of selecting the best model using BME, to improve model prediction? If the purpose is to improve model prediction, did the authors try Bayesian model averaging? Based on the results, some models are underestimate, some models are overestimate, it is possible that model averaging could give a better prediction performance.*

Response: We believe this comment is very important for considerably improving our manuscript. Our original idea was to identify which model is optimal for ET prediction, and then to improve the model with the Bayesian model averaging. However, our result showed that model SW obtained a weight of 100%. This means that BME assigned a weight of 100% to the competing model and the weights of 0% to other three alternative models, and thus, the Bayesian model averaging prediction is also the SW prediction. Other studies on hydrological model selection have yielded similar results in that one model obtained an weight of close to 100% (e.g., Meyer et al., 2007; Lu et al., 2013; Schöniger et al., 2007). Therefore, Bayesian model averaging was not used in our study.

2. Comment: *For the objective (2), theoretically we know these statistics only measure model fit without considering model complexity, so they are not as robust as BME. And we know these statistics can be efficiently calculated, so there is no need for testing. Please justify the objective (2).*

Response: We have changed the original objective (1) and (2) to “(1) to calibrate ET model parameters using the differential evolution adaptive metropolis (DREAM) algorithm; (2) to identify which parameters had a greater impact on the model performance and to explain why the selected optimal model performed best”. These changes can be seen at P5, L120-121.

3. Comment: *objective (3) is very meaningful. I would like to see more analysis on the model-data mismatch to improve model development and model performance.*

Response: Considering the Reviewer’s suggestion, we added some sentences in 4.1 “Parameter uncertainty analysis”, which can be seen at P18-P19, L439-453, and in 4.3 “Analysis of model-data mismatch”, which can be seen at P24, L597-601.

4. Comment: *In several places of the manuscript, the logic is not very clear. The English writing needs improvement.*

Response: It is true as you suggested that our manuscript needs the revision of English sentence. After the revision according to the reviewer’s comments, the manuscript have been edited

by the professional translation services.

Specific comments:

1. Comment: *Abstract, I think including some insights obtained from the numerical experiments in the abstract would attract more audience and make this work more meaningful.*

Response: Considering the Reviewer's suggestion, we added and revised the content in abstract. These changes can be seen at P1, L19-23.

2. Comment: *2. Line 35-36, the SW model performs best in this study area, but may not be the best in other areas. For example, in Li et al., (2013)'s study, PM performed better than SW in estimation of maize. Please justify the statement that SW should be the first choice for evaluating ET of spring maize in arid desert oasis areas.*

Response: We have changed this statement to "The mismatch analysis indicated that explicit treatment of energy imbalance and energy interaction will be the primary way to further improve ET model performance." The change can be seen at P2, L38-39.

3. Comment: *Line 93-94, BME can be used to compare and select the best-performing model. This is well-known and not a hypothesis that needs to be determined. What do the authors mean by saying "an unbiased view"?*

Response: We have deleted this sentence at P4, L111.

4. Comment: *Line 95-97, the first part of the sentence says Bayesian applications have focused on comparison of alternative models, but the second part of the sentence says that little attention has been given to the Bayesian model comparison. The sentence is self-contradictory. Please clarify.*

Response: Our original intention is to say that Bayesian applications have focused calibration of individual models but the model comparison is still conducted using traditional statistical criteria. Our expression was not very clear. We changed this statement as "Most applications of Bayesian methods have focused on the calibration of individual models, while the comparison of alternative models continues to be performed using traditional error metrics." This changes can be seen at P4, L107-108.

5. Comment: *Line 277, for each chain? I thought you total have 40,000 samples from all chains. In addition, Line 848, from one chain? Please clarify.*

Response: This referred the 40,000 samples from one chain. Total have 40,000 multiplied by N (chain number) samples from all chains. Please see P13, L298-299; P32, L874; P35, L905.

6. Comment: Line 280-282, based on Figure 1, DREAM needs far less than 8000 generations to make the GR statistic smaller than 1.2. Also, based on Figure 1's x-axis scale, it is hard to tell "obviously" the chain converged after about 620 and 450 generations.

Response: Thanks for the comment. We changed the sentence "DREAM needs far less than 8000 generations to make the GR statistic smaller than 1.2 for the both models" to "The algorithm needs about 8,000 generations to make the G-R statistic close to 1.0 for the both models." In addition, we deleted statement "Obviously, the complete mixing of the different chains and convergence of DREAM were attained after about 620 and 450 generations for PM and SW models, respectively". The change can be seen at P13-14, L302-304.

7. Comment: Figure 1. In Figure 1(b) the position of the dash line is not at 1.2. The position of the label (b) is not aligned well with the label (a).

Response: Thanks for the comment. We have redrawn the Figure 1. This can be seen Figure 1.

8. Comment: Figure 2. If the authors cannot get more information from the CDFs than the histograms, I suggest deleting the CDFs which make Figure 2 busy and confusing. Also, I would like to see more discussion about Figure 2; what insights the authors can obtain from these plots?

Response: We have redrawn the Figure 2. This can be seen Figure 2.

The main insights summarized as following: g_s and modeled ET in PM model are relatively insensitive to Q_{50} , D_{50} and K_q . Hence, these parameters could not be well constrained. The calculation of g_s^c in SW model is the same as in PM model, and thus, g_s^c and modeled ET for SW model are also insensitive to parameters of Q_{50} , D_{50} , K_q . Therefore, these three parameters were also not be well constrained in SW model. In addition, for edge-hitting parameters, their uncertainties may be also the outcome of model biases or EC-measured ET data errors. Although the ecophysiological parameter g_{max} is a variable in the g_s^c equation in both PM and SW models, but this parameter is sensitive to g_s^c and has large influences on the evaluated ET. Its effect is relatively independent compared to the other meteorological parameters in the models, and therefore this parameter was well specified in SW model. The parameter K_a is insensitive to g_s and modeled ET. In contrast, K_a is contained in equation of net radiation flux into the substrate (Eq.12) in SW model. From the above analysis, we could see that K_a not only involved the distribution of energy between the canopy and the soil surface but also the energy imbalance. Therefore, parameter K_a has a great influence on the performance of the SW model.

The revisions and changes can be seen P14, L314-321; P14, L326-330 and P17, L397-412.

9. Comment: *Line 294-297, I found the discussion of the figure 2 is confusing. I think, the figure 2 says the histograms tend to concentrate in the upper bounds, not the lower bounds. Also, the authors should increase the upper limits of these parameters not decrease, because the histograms are concentrated in the upper bounds.*

Response: Thanks for the point. We have revised this sentence. The change can be seen at P14, L315-321.

11. Comment: *Line 355-356, what do the author mean by saying “to sample groups of variable in turn”?*

Response: The sentence should be “to sample one or groups of variable in turn”. We have removed some sentence and this can be seen P17, L388-393.

12. Technical corrections: *Line 29, obstained → obtained. Line 92, beed → been*

Response: We have corrected the words already.

Once again, thank you very much for your comments and suggestions.

Referee #2:

Comments:

1. Comment: *The grammar of this paper needs some improvements, some grammar errors and ambiguous sentences can be found.*

Response: It is really true as you suggested that our manuscript needs a language improvement. After the revision according to the reviewer's comments, the manuscript have been edited by the professional translation services.

2. Comment: *The universality of this study and its conclusions need to be clarified since the study area and methodology are both very spatial and temporal specific. PS, are the conclusions valid under other conditions or not?*

Response: The ET models and BME model selection can be applied to other conditions as long as the required data can be obtained. Although there are many studies on ET model evaluation, their conclusions about model ranking are all based on traditional error metrics. Just as you said, the conclusion about whether SW model is optimal selected by BME method under other conditions still needs further confirmation. We have added relevant contents at P24, L594-596.

3. Comment: *Following the last comment, is it possible to provide results for other study areas or using other time scales? This will provide strong evidences to support the conclusions.*

Response: It is really true that providing results for other study areas or using other time scales would be very useful for providing strong evidences to support the conclusions. We've been looking for reliable data from other study area or from other crops for BME model selection to confirm whether the SW model is the optimal model under other conditions. However, it is difficult to obtain the required data by ET models, especially the soil water contents. So far, we haven't got the requited data yet. And thus, Thanks for the comment that we are not able to provide results for other study areas or using other time scales for BME model selection by now.

4. Comment: *I am not sure I can agree with some conclusions, for example, the one in lines 531-532, the authors suggest prioritizing BME over other measurements, but BME can also provide inaccurate results.*

Response: We think this is true, and we deleted statement "and that BME should be used instead", and reorganized the original sentences. Please see P25, L609-615.

Once again, thank you very much for your comments and suggestions.

Referee #3:

General Comments:

1. Comment: *Language issues should be fully checked throughout the entire text before publication in HESS.*

Response: It is really true as you suggested that our manuscript needs the revision of English sentence. After the revision according to the reviewer's comments, the manuscript have been edited by the professional translation services.

2. Comment: *Novelty of the paper should be better emphasized rather than "BME has not been used for evaluating the ET models".*

Response: Thanks for the comment. We have changed the original statement to "Currently, ET model selection and comparison have been still conducted using traditional error metrics. It is known that error metrics are not adequate to provide reasonable result of model ranking for disregarding model complexity (Marshall et al., 2005; Samani et al., 2018). The focus of this study is to use a Bayesian approach to evaluate the performance of the PM, SW, PT-FC, and AA models, which is a novelty contribution of this study." These changes can be seen at P3, L81-85.

3. Comment: *Model complexity for each model should be better described. For example, authors can directly introduce number of parameters with uncertainties in their experiment?*

Response: Considering the comments, we have described the number of parameters of each model at P13, L289-291 and added a sentences "The results illustrate that with the addition of parameters, the model complexity and the model performance are both increased." at P16, L382-383.

Specific comments:

1. Comment: *The Abstract is out of organization. It seems to me that you never mention the model complexity but always write "underestimation" or "overestimate" to explain why the SW is the best one.*

Response: We have reorganized the Abstract. These can be seen P1-P2, L19-39. We also have added sentence "Although the SW model with seven parameters is sophisticated, it's good fitting to observations can counterbalance its higher complexity." These can be seen P2, L35-36.

2. Comment: *Lines 25-27: unclear, please rephrase this sentence.*

Response: Considering the comments, we have reorganized the abstract again, and changed the original sentence to "The parameters in each model were first calibrated using Differential Evolution Adaptive Metropolis (DREAM) algorithm, and then were analyzed to identify their

impacts on the model performance. The Bayesian model evidence (BME) approach, was further adopted to select the optimal model by incorporating the mathematically rigorous thermodynamic integration algorithm.” These modifications can be seen at P1, L19-23.

3. Comment: *It is unclear for me why ‘SW’ is best one from the abstract.*

Response: Considering the comments, we have added the sentences “Our results revealed that the extinction coefficient was the most significant parameter in the ET models. It was not merely partitioning the total available energy into the canopy and surface, but also including the energy imbalance correction. The extinction coefficient is well constrained in the SW model and poorly constrained in the PM model, but not considered in PT-FC and AA models” to explain why SW is the best one in abstract. Please see P2, L31-35.

4. Comment: *Line37: please check the symbol.*

Response: We have checked the symbol at P2, L40.

5. Comment: Simulate ET or estimate ET? Please be very sure of this word.

Response: We have changed some “estimate ET” to “simulate ET”.

6. Comment: *Line 41: add a reference.*

Response: Thanks for the comments. We have added corresponding reference “(Brutsaert, 2005)”. Please see at P2, L45.

7. Comment: *Lines 55-56: unclear, please rephrase this sentence.*

Response: We rephased this sentence as “These ET models are generally complex, because of for the coupling of the land surface and atmospheric processes, and high-dimensional with a large number of parameters”. Please see at P3, L59-60.

8. Comment: *Lines 62-63: ‘These quantitative criteria’ refer to what?*

Response: We have reorganized this paragraph, and deleted the original statement. Please see at P3, L67-68.

9. Comment: *Line 70: performances.*

Response: Thanks. We have reorganized this paragraph and corrected the statement. Please see at P3, L75.

10. Comment: *Line 71: remove ‘the’ from ‘the SW model’*

Response: We have removed ‘the’ from ‘the SW model’.

11. Comment: *Lines 71-72: please rephrase this sentence.*

Response: We have changed the sentence “Ershadi et al. (2014) evaluated the surface energy balance system (SEBS), PM, PT-JPL (a modified Priestley-Taylor model), and AA models.” This can be seen at P3, L76-77.

12. Comment: *Line 73: should be model ranking? Please check the terminology.*

Response: Thanks. We have corrected this mistake at P3, L69.

13. Comment: *Lines 75-76: unclear, significant variability of model performances?*

Response: Considering the comments, we have changed “significant” to “considerable”. This can be seen at P3, L80.

14. Comment: *Lines 92-93: been?*

Response: We have corrected the sentence.

15. Comment: *Lines 102-103: add a reference*

Response: We have added the reference “(Vrugt et al., 2008, 2009)” at P5, L119-120.

Once again, thank you very much for your comments and suggestions.

Bayesian performance evaluation of evapotranspiration models: a case study for an arid region based on eddy covariance system of a maize field in northwestern China

Guoxiao Wei^{1,2,*}, Xiaoying Zhang^{3,*}, Ming Ye⁴, Ning Yue^{1,2}, Fei Kan^{1,2}

¹ Key Laboratory of Western China's Environmental System (Ministry of Education), Lanzhou University, China, 730000

² School of Earth and Environmental Sciences, Lanzhou University, China, 730000

³ Construct Engineering College, Jilin University, China, 130400

⁴ Department of Earth, Ocean, and Atmospheric Science, Florida State University, USA, 32306

* Corresponding author: xiaoyingzh@jlu.edu.cn.

Abstract

Evapotranspiration (ET) is a major component of the land surface process involved in energy fluxes and energy balance, especially in the hydrological cycle of agricultural ecosystems. While many models have been developed as powerful tools to estimate ET, there has been no agreement on which model has the best describing the loss of water to the atmosphere. In this study, we present a solid study to evaluate four widely used ET models and their parameter contributions (i.e., the Shuttleworth-Wallace (SW) model, Penman-Monteith (PM) model, Priestley-Taylor and Flint-Childs (PT-FC) model, and Advection-Aridity (AA) model) by using half-hourly ET observations obtained at a spring maize field in an arid region. The four tested models are the Shuttleworth-Wallace (SW) model, Penman-Monteith (PM) model, Priestley-Taylor and Flint-Childs (PT-FC) model, and Advection-Aridity (AA) model. The parameters in each model were first calibrated using Differential Evolution Adaptive Metropolis (DREAM) algorithm, and then were analyzed to identify their impacts on the model performance. The Bayesian model evidence (BME) approach, was further adopted to select the optimal model by incorporating the mathematically rigorous thermodynamic integration algorithm. The BME-based model ranking (from the best to the worst) is SW, PM, PT-FC, and AA. The residuals between observations and corresponding model simulations are also analyzed, and the same model ranking is also obtained by using residual-based statistics, i.e., the coefficient of determination (R^2), index of agreement (IA), root mean square error (RMSE) and model efficiency (EF). The PM and SW models overestimate ET, whereas the PT-FC and AA models underestimate ET in the study period. The four models also underestimate ET during the periods of partial crop cover. Especially during the late maturity stage, the PT-FC and AA models consistently produce an underestimation, and provide the worst simulated ET. As a result, at the half-hourly time scale, the SW model is the best model and recommend as the first choice for evaluating ET of spring maize in arid-desert-oasis areas. Our results revealed that the extinction coefficient was the most significant parameter in the ET models. It was not merely partitioning the total available energy into the canopy and surface, but also including the energy imbalance correction. The extinction coefficient is well constrained in the SW model and poorly constrained in the PM model, but not considered in PT-FC and AA models. The extinction coefficient is well constrained in the SW model and poorly constrained in the PM model, but not considered in PT-FC and AA models. This is the main reason that the SW model outperforming the other models. Although the SW model with seven parameters is sophisticated, its good fitting to observations can counterbalance its higher

38 complexity. In addition, the discrepancies between observations and model simulations were evaluated using
39 traditional error metrics. The mismatch analysis indicated that explicit treatment of energy imbalance and
40 energy interaction will be the primary way to further improve ET model performance.

41 **Keywords:** Bayesian analysis; ET models; Eddy covariance; ~~Penman-Monteith; Shuttleworth-Wallace; Model~~
42 ~~performance; Extinction coefficient; Maize~~

43 1. Introduction

44 Surface energy fluxes are an important component of Earth's global energy budget and a primary
45 determinant of surface climate. Evapotranspiration (ET), as a major energy flux process for energy balance,
46 accounts for about 60--65% of the average precipitation over the surface of the Earth (Brutsaert, 2005). In
47 agricultural ecosystems, more than 90% of the total water losses are due to ET (BrutsaertMorison et al.,
48 20052008). Therefore, robust ET estimation is crucial to a wide range of problems in hydrology (Xu and Singh,
49 1998), ecology, and global climate change (Xu and Singh, 1998Morison et al., 2008). In practice, much of our
50 understanding of how land surface processes and vegetation affect weather and climate is based on numerical
51 modeling of surface energy fluxes and the atmospherically-coupled hydrological cycle (Bonan, 2008). Several
52 models are commonly used in agricultural systems to evaluate ET. The Penman-Monteith (PM) and
53 Shuttleworth-Wallace (SW) models are physically sound and rigorous (Zhu et al., 2013), and thus widely used
54 to estimate ET for seasonally varied vegetations. The models consider the relationships among-between net
55 radiation, all kinds of heat flux (such as latent heat, sensible heat, and heat from soil and canopy), and surface
56 temperature. The Priestley-Taylor and Flint-Childs (PT-FC) model (based on radiation) and the
57 advection-aridity (AA) model (based on meteorological variables) have also been widely used because they only
58 require a small amountnumber of ground-based measurements for settingto set up the models (Ershadi et al.,
59 2014).

60 These ET models are generally complex, because of for-the coupling of the land surface and atmospheric
61 processes, and high-dimensional- with a large number of parameters. Modelers are challenged by how to
62 compare-Comparing the performance of competing models and howand to-evaluating eand understanding the
63 mismatch-discrepancies between simulations of the -model-simulations and corresponding observed
64 surface-atmosphere water flux are remain challenging problems (Legates, 1999). Moreover, how to choose a
65 criterion to reliably evaluate model performance is another crucial issue.-Both non-Bayesian analysis (Szilagyi
66 and Jozsa, 2008; Vinukollu et al., 2011; Li et al., 2013; Ershadi et al., 2015) and Bayesian analysis have been
67 used for to evaluating-evaluate the model-performance of ET models (Zhu et al., 2014; Chen et al., 2015; Liu et
68 al., 2016; Zhang et al., 2017; Elshall et al., 2018; Samani et al., 2018; Zeng et al., 2018) . These quantitative-
69 criteria used for model evaluation and selection include residual-based measures (e.g., regression line slope and
70 mean bias error, MBE), squared residual based measures (e.g., coefficient of determination, R^2), root mean-
71 square error (RMSE), model efficiency (EF), and index of agreement (IA).-Li et al. (2013) compared the maize-
72 ET estimates-simulations of the PM, SW and adjusted SW models under film-mulching conditions of maize
73 growth in an arid region of China. They found that the half-hourly ET was overestimated by 17% by the SW
74 model. ,with relatively high MBE, RMSE, and lower R^2 and IA.-In contrast, the PM and MSW-adjusted SW
75 models underestimated the daily ET by 6% and 2%, respectively,-during the entire experimental period of 116-

76 ~~days~~. Therefore, the performances of PM and adjusted SW models are better than that of the SW model in their
77 case ~~study~~. Ershadi et al. (2014) evaluated the surface energy balance system (SEBS), PM, PT-JPL (a modified
78 Priestley–Taylor model, ~~similar to the PT-FC~~) and AA models. Based on the average value of EF and RMSE,
79 the model ranking from ~~the~~-worst to ~~the~~-best was AA, PM, SEBS, and PT-JPL. Ershadi et al. (2015) also
80 ~~compared the evaluated model~~-response of the models to ~~the~~-different formulations of aerodynamic and surface
81 resistances ~~against with~~ global FLUXNET data. Their results showed ~~significant-considerable~~ variability in
82 model performance among and within biome types. ~~Currently, ET model selection and comparison have been~~
83 ~~still conducted using traditional error metrics. It is known that error metrics are not adequate to provide~~
84 ~~reasonable result of model ranking for disregarding model complexity (Marshall et al., 2005; Samani et al.,~~
85 ~~2018). The focus of this study is to use a Bayesian approach to evaluate the performance of the PM, SW, PT-FC,~~
86 ~~and AA models, which is a novelty contribution of this study.~~

87 In ET models, the land surface energy system is governed by presumably infinite-dimensional physics.
88 However, considering the ET models as finite-dimensional can be more precisely by covering all relevant
89 relations. Therefore, employing consistent criteria for model selection might be justified when the aim is to
90 better understand the processes involved (Höge et al., 2018). When using consistent model selection, ~~The~~
91 Bayesian model evidence (BME), also known as marginal likelihood, measures the average fit of ~~a~~-model
92 simulations to their corresponding observations ~~to the data~~ over a model's prior parameter space. This feature
93 enables BME to consider model complexity (in terms of number of model parameters) for model performance
94 evaluation. When comparing several alternative conceptual models, the model with the largest marginal
95 likelihood is selected as the best model (Lartillot and Philippe, 2006). BME can thus be used for evaluating the
96 model fit (over the parameter space) and for comparing alternative models. In previous studies, ~~the~~ Bayesian
97 information criterion (BIC; ~~KashyapSchwarz, 1982~~1978) ~~or and the~~ Kashyap information criterion (KIC;
98 ~~Schwarz 1978~~Kashyap, 1982) ~~were have been~~ used to approximate BME ~~for reducing~~ by using maximum
99 likelihood theories to reduce computational cost of evaluating BME (Ye et al., 2004). However, these
100 approximations have theoretical and computational limitations (Ye et al., 2008; Xie et al., 2011; Schöniger et
101 al., 2014), and a numerical evaluation (not a likelihood approximation) of BME is necessary, especially for
102 complex models (Lartillot and Philippe, 2006). Lartillot and Philippe (2006) advocated the use of
103 thermodynamic integration (TI) for estimating BME, ~~which is~~-also known as path sampling (Gelman and Meng,
104 1998; Neal, 2000), in order to avoid sampling solely in the prior or posterior parameter space. TI uses samples
105 that are systematically generated from the prior to the posterior parameter space by conducting path sampling
106 with several discrete power coefficient values (Liu et al., 2016). It is ~~both mathematically rigorous and~~
107 ~~more numerically~~ accurate than the generally used harmonic mean method (Xie et al., 2011).

108 Most applications of Bayesian methods have focused on the calibration of individual models, while the
109 comparison of alternative models continues to be performed using traditional error metrics. ~~While many~~
110 ~~statistical criteria have been used to evaluate different ET models, BME has not been used for evaluating the ET~~
111 ~~models. It remains to be determined whether BME can be used to compare and select the best model and~~
112 ~~whether BME can provide an unbiased view of the performance of the models. Furthermore, most Bayesian~~
113 ~~applications have focused on the calibration of individual models and comparison of alternative models using~~

114 ~~these statistical measures, with little attention given to the Bayesian model comparison.~~ More generally,
115 Bayesian approaches to Mmodel calibration, comparison, and analysis ~~underlying the Bayesian paradigm has~~
116 ~~have been used much far~~ less used in the evaluation of ET models than in other areas of environmental science.
117 In this study, the Bayesian approach ~~was is~~ used to calibrate and evaluate the four ET models (PM, SW, PT-FC,
118 and AA) based on an experiment over a spring maize field in an arid area of northwest China, from 3 June to 27
119 September 2014. The objectives of the study are as follows: (1) ~~The to calibrate ET model parameters were~~
120 ~~calibrated~~ using the DiffeRential Evolution Adaptive Metropolis (DREAM) algorithm (Vrugt et al., 2008,
121 2009)-; ~~The objectives of the study are as follows:-~~(12) to identify which parameters had a greater impact on the
122 model performance and to explain why the selected optimal model performed best; ~~compare the four models~~
123 ~~and select the best one using BME~~; (23) to evaluate the performance of the models using traditional error
124 metrics and BME; and ~~various general statistics such as correlation based measures (R^2), relative error measures~~
125 ~~(IA and EF), and absolute error measures (such as RMSE and MBE) and to determine whether these methods~~
126 ~~are efficient and reasonable for evaluating the ET models~~; (34) to analyze discrepancies between model
127 simulations and observation data in order to ~~model data mismatch for~~ better understanding model performance
128 and identify ways to improve these models. ~~Using BME for evaluating the ET models has not been reported in~~
129 ~~the literature.~~ We expect that the study will not only boost the development of model parameterization and
130 model selection but also contribute to the improvement of the ET models.

131 2. Data and methodology

132 2.1. Description of the study area

133 The experiment of ~~of maize growth~~ was conducted at Daman Superstation, located in Zhangye City, Gansu
134 ~~province~~, northwest China. Daman Oasis is located in the middle Heihe River basin, which is the second largest
135 inland river basin in the arid region of northwest China. The midstream area of the Heihe River basin is
136 characterized by oases with irrigated agriculture, and is a ~~major zone of water region that consumes large amount~~
137 ~~of water consumption~~ for both domestic and agricultural uses. The annual average precipitation and temperature
138 are 125 mm and 7.2 °C (1960–2000), respectively. The annual accumulated temperature (>10 °C) is 3,234 °C,
139 and the annual average potential evaporation is about 2,290 mm. The average annual duration of sunshine is
140 3,106 h with 148 frost-free days. The predominant soil type is silty-clay loam and the depth of the frozen layer
141 is about 143 mm. The study area is a typical irrigated ~~agriculture agricultural~~ region, and the major ~~source of~~
142 water ~~resources are is the~~ snowmelt from the Qilian Mountains. ~~The maize-Maize~~ and spring wheat are the
143 principal crops ~~grown in the region.~~; Maize is generally sown in late April and harvested in mid-September, and
144 is planted with a row spacing of 40 cm and a plant spacing of 30 cm. The plant density is about 66,000 plants
145 per hectare ~~in the study area.~~

146 2.2. Measurements and data processing

147 Our ~~observation~~ data were collected from the field observation systems of the Heihe Watershed Allied
148 Telemetry Experimental Research (HiWATER) project as described in Li et al (2013). The observation period
149 was from DOY (day of the year) 154 to DOY 270 in 2014. An open-path eddy covariance (EC) system was
150 installed in a maize field, with the sensors at ~~thea~~ height of 4.5 m. Maize is the main crop in the study region,

151 and thus covers ~~which can supply~~ sufficient planting area to set the EC measurements. The EC data was logged
 152 at a frequency of 10 Hz and then processed with an average time interval of 30 min. Sensible and latent heat
 153 fluxes were computed by the EC approach of Baldocchi (2003). Flux data measured by EC were controlled by
 154 traditional ~~routes~~ methods, including three-dimensional rotation (Aubinet et al., 2000), WPL
 155 (Webb-Penman-Leuning) density fluctuation correction (Webb et al., 1980), frequency response correction (Xu
 156 et al., 2014), and spurious data removal caused by rainfall, water condensation, and system failure. About 85%
 157 of the energy balance closure was observed in the EC data (Liu et al., 2011).

158 Standard hydro-meteorological variables, including rainfall, air temperature, wind speed, and wind
 159 direction, were continuously measured at the heights of 3, 5, 10, 15, 20, 30 and 40 m above the ground. Soil
 160 temperature and moisture were measured at heights of 2, 4, 10, 20, 40, 80, 120 and 160 cm. Photosynthetically
 161 active radiation was measured at a height of 12 m. Net radiation, including downward, ~~and~~ upward and
 162 longwave radiation, was measured by a four-component net radiometer. An infrared thermometer was installed
 163 at a height of 12 m. Leaf Area Index (LAI) was measured approximately every 10 days during the growing
 164 season.

165 2.3. Model description

166 In this section, we summarize the mathematical definitions forming the basis of each of the four models.
 167 Appendix A contains a summary of the names and physical meanings of the model parameters.

168 2.3.1 Penman-Monteith (PM) model

169 The PM model can be formulated in the ~~as~~ following way (Monteith, 1965) ~~and most of the parameters are~~
 170 ~~explained in Appendix A:~~

$$171 \quad \lambda E = \frac{\varepsilon A + (\rho C_p / \gamma) D_a g_a}{\varepsilon + 1 + g_a / g_s} \quad (1)$$

172 where $\varepsilon = \Delta / \gamma$; and A is defined to be $A = R_n - G$.

173 In the present study, g_a is parameterized ~~as in the way~~ suggested by Leuning (2008) and g_s is defined as:

$$174 \quad g_s = g_s^c \left[\frac{1 + \frac{\tau g_a}{(\varepsilon + 1) g_s^c} \left[f - \frac{(\varepsilon + 1)(1 - f) g_s^c}{g_a} \right] + \frac{g_a}{\varepsilon g_i}}{1 - \tau \left[f - \frac{(\varepsilon + 1)(1 - f) g_s^c}{g_a} \right] + \frac{g_a}{\varepsilon g_i}} \right] \quad (32)$$

175 where $1 - \tau$ and τ ~~is~~ are the fraction of the total available energy absorbed by the canopy and by the soil, and $\tau =$
 176 $\exp(-K_d LAI)$; and g_i and g_s^c ~~is~~ are defined ~~as in~~ equations (3) and (4), respectively (Monteith, 1965):

$$177 \quad g_i = \frac{A}{(\rho C_p / \gamma) D_a} \quad (43)$$

178 ~~(Monteith, 1965);~~ g_s^c is expressed as:

$$179 \quad g_s^c = \frac{g_{\max}}{K_q} \ln \left[\frac{Q_h + Q_{50}}{Q_h \exp(-K_q \text{LAI}) + Q_{50}} \right] \left[\frac{1}{1 + D_a / D_{50}} \right] f(\theta) \quad (54)$$

180 where $f(\theta)$ is the factor considersrepresents water stress and is expressed as:

$$181 \quad f(\theta) = \begin{cases} 1 & \theta > \theta_a \\ \frac{\theta - \theta_b}{\theta_a - \theta_b} & \theta_b < \theta < \theta_a \\ 0 & \theta < \theta_b \end{cases} \quad (65)$$

182 ~~where and~~ θ_a ~~was is~~ set as $\theta_a = 0.75 \theta_b$. Aerodynamic conductance g_a is calculated as:

$$183 \quad g_a = \frac{k^2 u_m}{\ln[(z_m - d)/z_{0m}] \ln[(z_m - d)/z_{0v}]} \quad (76)$$

184 where the quantities d , z_{0m} and z_{0v} are calculated using $d = 2h/3$, $z_{0m} = 0.123h$ and $z_{0v} = 0.1z_{0m}$ (Allen 1998).

185 2.3.2. Shuttleworth-Wallace (SW) model

186 The SW model comprises a one-dimensional model of plant transpiration and a one-dimensional model of
187 soil evaporation. The two terms are calculated by the following equations:

$$188 \quad \lambda E T = \lambda E + \lambda T = C_s E T_s + C_c E \quad (7)$$

$$189 \quad ET_s = \frac{\Delta A + \{\rho C_p (e_s - e_a) - \Delta r_a^s (A - A_s)\} / (r_a^a + r_a^s)}{\Delta + \gamma \{1 + r_s^s / (r_a^a + r_a^s)\}} \quad (8)$$

$$190 \quad ET_c = \frac{\Delta A + \{\rho C_p (e_s - e_d) - \Delta r_a^c A\} / (r_a^a + r_a^d)}{\Delta + \gamma \{1 + r_s^c / (r_a^a + r_a^d)\}} \quad (9)$$

191 where the available energy input above the soil surface is defined as $A_s = R_{ns} - G$.

192 R_{ns} can be calculated using the Beer's law relationship:

$$193 \quad R_{ns} = R_n \exp(-K_d L A) \quad (10)$$

194 The coefficients C_s and C_c are obtained as follows:

195
$$C_s = \left\{1 + R_s R_a / R_c (R_s + R_a)\right\}^{-1} \quad (11)$$

196
$$C_c = \left\{1 + R_c R_a / R_s (R_c + R_a)\right\}^{-1} \quad (12)$$

197 where

198
$$R_a = (\Delta + \gamma) r_a^a \quad (13)$$

199
$$R_s = (\Delta + \gamma) r_a^s + \gamma r_s^s \quad (14)$$

200
$$R_c = (\Delta + \gamma) r_a^c + \gamma r_s^c \quad (15)$$

201 Soil surface resistance is expressed as:

202
$$r_s^s = \exp\left(-b_2 \frac{\theta}{\theta_s}\right) \quad (16)$$

203 In this study, we consider the reciprocal of bulk stomatal resistance, known as canopy conductance. The
 204 calculation of g_s^c is the same as in the PM model. The two aerodynamic resistances (r_a^a and r_a^s) and the
 205 boundary layer resistance (r_a^c) are modeled following the approach proposed by Shuttleworth and Gurney
 206 (1990).

207 **2.3.3. Priestley–Taylor and Flint-Childs (PT-FC) model**

208 The Priestley–Taylor model (Priestley and Taylor, 1972) model was introduced to estimate evaporation
 209 from an extensive wet surface under conditions of minimum advection (Stannard, 1993; Sumner and Jacobs,
 210 2005). The ET is expressed as:

211
$$\lambda ET = \alpha_{PT} \frac{\Delta}{\Delta + \gamma} (R_n - G) \quad (17)$$

212 where α_{PT} is a unitless coefficient. The Priestley-Taylor model was modified by Flint and Childs (1991) in order
 213 to scale the Priestley-Taylor potential ET to actual ET for nonpotential conditions (hereafter the PT-FC model):

214
$$\lambda ET = \alpha \frac{\Delta}{\Delta + \gamma} (R_n - G) \quad (18)$$

215 where α is as a function of the environmental variables, which could be related to any process that limits ET
 216 (e.g., soil hydraulic resistance, aerodynamic resistance, stomatal resistance); however, only soil moisture status
 217 was considered to simplify ET estimation in the PT-FC model (Flint and Childs, 1991). In this model, α is
 218 defined as:

$$219 \quad \alpha = \beta_1 [1 - \exp(-\beta_2 \Theta)] \quad (19)$$

220 where ~~Θ is calculated as~~ $\Theta = \frac{\theta - \theta_r}{\theta_s - \theta_r}$.

221 2.3.4. Advection-aridity (AA) model

222 The AA model was first proposed by Brutsaert and Stricker (1979) and further improved by Parlange and
 223 Katul (1992). The model relies on the feedback between actual (λET) and potential ET , which assumes that
 224 actual potential ET should converge to wet surface ET at wet surface conditions. Its general form is:

$$225 \quad \lambda ET = (2\alpha_{PT} - 1) \frac{\Delta}{\Delta + \gamma} (R_n - G) - \frac{\gamma}{\Delta + \gamma} \frac{\rho(q^* - q)}{r_a} \quad (20)$$

226 where α_{PT} is the Priestley–Taylor coefficient, usually taken as 1.26 (Priestley and Taylor, 1972); and r_a is
 227 similar to that used for the Penman-Monteith model (Brutsaert and Stricker, 1979; Brutsaert, 2005; Ershadi et al.,
 228 2014). This model is based mainly on meteorological variables and does not require any information related to
 229 soil moisture, canopy resistance or other measures of aridity (Ershadi et al., 2014). In this study, ~~similar to as for~~
 230 the PT-FC model, we ~~modified-changed~~ α_{PT} to α , which is calculated using the same equation as in the PT-FC
 231 model. ~~The detailed list of symbols and physical characteristics in ET models are stated in Appendix A.~~

232 2.4 BME Estimation

233 The Bayesian model evidence (BME) of a model, M , is defined as (Schöniger et al., 2014):

$$234 \quad \text{BME} = p(\mathbf{D}|M) = \int p(\mathbf{D}|\boldsymbol{\theta}, M) p(\boldsymbol{\theta}|M) d\boldsymbol{\theta} \quad (21)$$

235 where \mathbf{D} is observed or estimated data, $\boldsymbol{\theta}$ is the vector of parameters associated with model M , $p(\boldsymbol{\theta}|M)$ is
 236 the prior density of $\boldsymbol{\theta}$ under model M , $p(\mathbf{D}|\boldsymbol{\theta}, M)$ is the joint likelihood of model M and its parameters $\boldsymbol{\theta}$.
 237 Estimating BME using power posterior estimators such as thermodynamic integration (TI) (Lartillot and
 238 Philippe, 2006) depends mainly on the calculation of the marginal likelihood $p(\mathbf{D}|M)$. The main idea of power
 239 posterior sampling is to define a path that links the prior to the unnormalized posterior. Thus, using an
 240 unnormalized power posterior density

$$241 \quad q_\beta(\boldsymbol{\theta}) = p(\mathbf{D}|\boldsymbol{\theta}, M)^\beta p(\boldsymbol{\theta}|M) \quad (22)$$

242 the power coefficient $\beta \in [0, 1]$ is a scalar parameter for discretizing a continuous and differentiable path
 243 linking two unnormalized power posterior densities. The unnormalized power posterior density $q_\beta(\boldsymbol{\theta})$ in
 244 Equation (22) uses the normalizing constant Z_β to yield the normalized power posterior density:

245
$$p_{\beta}(\boldsymbol{\theta}) = \frac{q_{\beta}(\boldsymbol{\theta})}{Z_{\beta}} \quad (23)$$

246 such that

247
$$Z_{\beta} = \int q_{\beta}(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (24)$$

248 The above integral takes a simplified form by the potential:

249
$$U(\boldsymbol{\theta}) = \frac{\partial \ln q_{\beta}(\boldsymbol{\theta})}{\partial \beta} \quad (25)$$

250 thus, the integral can be directly estimated by the following way:

251
$$p(\mathbf{D}|M) = \frac{Z_1}{Z_0} = \exp\left\{\int_0^1 E_{\theta}[\ln p(\mathbf{D}|\boldsymbol{\theta}, M)] d\beta\right\} \quad (26)$$

252 The one-dimensional integral with respect to β is evaluated by using numerical methods by discretizing β into a
 253 set of β_k . Since there is no theoretical method for selecting β_k values (Liu et al., 2016), we determined these
 254 values using an empirical but straightforward method. Following Xie et al. (2011), a schedule of the power
 255 posterior coefficients β_k is generated by

256
$$\beta_k = (k / K)^{1/\varepsilon} \quad (27)$$

257 for $k=0, 1, 2, \dots, K$. Using $\varepsilon = 0.3$ and $K = 20$ is a reasonable initial choice. By using the trapezoidal rule of
 258 numerical inregration, equation (26) is evaluated via

259
$$p(\mathbf{D}|M) = \exp\left(\int_0^1 y_{\beta} d\beta\right) = \exp\left(\sum_{k=0}^K r_{TI,k}\right) \quad (28)$$

260 such that

261
$$r_{TI,k} = (\beta_k - \beta_{k-1}) \left[\frac{y_k - y_{k-1}}{2} \right] \quad (29)$$

262 and

263
$$y_k = E_{\beta}[\log p(\mathbf{D}|\boldsymbol{\theta}_k, M)] = \frac{1}{n} \sum_{i=1}^n \log p(\mathbf{D}|\boldsymbol{\theta}_{k,i}, M) \quad (30)$$

264 where n is the number of random samples of $\boldsymbol{\theta}_k$ corresponding to β_k , and $\boldsymbol{\theta}_{k,i}$ is the i -th sample.

265 The random samples, $\boldsymbol{\theta}_{k,i}$, are drawn by using the MCMC method implemented in the DREAM code. See
 266 Appendix B for further details on Bayesian inference and the DREAM algorithm. In the DREAM-based

276 calculation, the Metropolis acceptance ratio is $\alpha_k = \min\left(1, [\alpha_{k, power-posterior} \alpha_{k, prior}]\right)$ with the power
 277 posterior ratio given by $\alpha_{k, power-posterior} = \left(\alpha_{k, posterior}\right)^{\beta_k}$. The prior probability ratio
 278 $\alpha_{k, prior} = \Pr(\theta_{k, new} | M) / \Pr(\theta_{k, old} | M)$ is the ratio of the probability of the newly proposed sample
 279 $\theta_{k, new}$ and the probability of the previously accepted sample $\theta_{k, old}$. The posterior probability ratio
 280 $\alpha_{k, posterior} = L(\mathbf{D} | \theta_{k, new}, M) / L(\mathbf{D} | \theta_{k, old}, M)$ is the likelihood ratio of samples $\theta_{k, new}$ and $\theta_{k, old}$, and
 281 β_k is the power posterior coefficient. Thus, to use the DREAM algorithm to sample any power posterior
 282 distribution, the regular Metropolis acceptance ratio $\alpha = \min\left(1, [\alpha_{posterior} \alpha_{prior}]\right)$ is changed to
 283 $\alpha_k = \min\left(1, [\alpha_{k, power-posterior} \alpha_{k, prior}]\right)$ in DREAM.

275 2.5 Traditional statistical metrics of evaluating model performance

276 The traditional **error metrics** for evaluating model performance include ~~correlation-based measures of~~
 277 R^2 and slope (**correlation-based measures**), ~~relative error measures of~~ index of agreement (IA) and model
 278 **efficiency (EF)** (**relative error measures**), and ~~absolute the root mean square error measures of~~ (RMSE) and
 279 mean bias error (MBE) (Poblete-Echeverria and Ortega-Farias, 2009). ~~Their~~ **definitions of the listed metrics are**
 280 **as follows:**

$$281 \quad IA = 1 - \frac{\sum_{t=1}^n [O(t) - M(t)]^2}{\sum_{t=1}^n \left[\left| O(t) - \overline{O(t)} \right| + \left| O(t) - \overline{M(t)} \right| \right]^2} \quad (33301)$$

$$282 \quad EF = 1 - \frac{\sum_{t=1}^n [O(t) - M(t)]^2}{\sum_{t=1}^n [O(t) - \overline{O(t)}]^2} \quad (34312)$$

$$283 \quad RMSE = \sqrt{\frac{1}{n} \sum_{t=1}^n [O(t) - M(t)]^2} \quad (35323)$$

$$284 \quad MBE = \frac{1}{n} \sum_{t=1}^n [O(t) - M(t)] \quad (36334)$$

285 where $O(t)$ is the observations and $\overline{O(t)}$ is the mean observation at time t ; $M(t)$ is the modeled value and
 286 $\overline{M(t)}$ is the mean modeled value estimated by the posterior median parameter values; and n is the total number
 287 of the observed values.

288 3. Results

289 3.1 Parameter estimation

290 The ~~re-were~~PM model has five parameters g_{max} , D_{50} , Q_{50} , K_q and K_a . ~~in the PM model, and two the SW~~
291 ~~model has seven parameters – the five used in the PM model and additional parameters, b_1 and b_2 .~~ ~~in the SW~~
292 ~~model.~~—The PT-FC and AA models each include two parameters, denoted ~~as~~ by β_1 and β_2 (Table 1). The prior
293 probability density of each parameter ~~was~~ is specified as a uniform distribution with the ranges listed in Table 1.
294 A total of 50,000 realizations were generated with the DREAM algorithm, which was used to estimate the
295 posterior probability density function of each parameter with the calibration period data from DOY 154 to DOY
296 202. In the calculations, the chain number, N , was equal to the ~~number of parameter numbers in the associated~~
297 ~~model. Therefore, i.e., N is equal to 5, 7, 2 and 2 for the PM, SW, PT-FC and AA models, respectively.~~ For
298 each model, the first 10,000 samples were discarded as burn-in data, and the remaining 40,000 samples were
299 ~~used for calibration. to set up posterior density functions for each chain. In total, $40,000 \times N$ realizations were~~
300 ~~used to set up posterior density functions for each model. To understand~~ illustrate the efficiency and
301 convergence of DREAM for the ET models, Figure 1 shows the trace plots of the G-R statistic for each of the
302 different parameters in the PM and SW models ~~with the using a different color (PT-FC and AA models not~~
303 ~~shown).~~ The algorithm required about 8,000 generations to make the G-R statistic close to 1.0 ~~1.2~~ for the
304 ~~both two models. Obviously, the complete mixing of the different chains and convergence of DREAM were~~
305 ~~attained after about 620 and 450 generations for PM and SW models, respectively.~~ The acceptance rates for the
306 PM and SW models were about 15.3% and 18.9%, respectively.

307 Histograms ~~and cumulative distribution functions (CDFs)~~ of the DREAM-derived marginal distributions of
308 the parameters are presented in Figure 2 and summarized in Table 2 by Maximum Likelihood Estimates (MLEs),
309 posterior medians and 95% probability intervals. Figures 2a-2e, 2f-2l, 2m-2n, and 2o-2p show histograms of the
310 PM, SW, PT-FC and AA models, respectively. Parameter g_{max} (Fig. 2a) in the PM model, parameters g_{max} , K_a , b_2
311 (Fig. 2f, 2j, 2l) in the SW model, and parameter β_1 (Fig. 2m) in the PT-FC model and AA model (Fig. 2o) were
312 well constrained and occupied a relatively small range. ~~Parameter g_{max} (Fig. 2A) in the PM model, parameters~~
313 ~~g_{max} , K_a , b_1 , b_2 (Fig. 2I, 2M, 2N, 2O) in the SW model, and parameter β_1 (Fig. 2F) were well constrained and~~
314 ~~occupied a relatively small range.~~ These parameters displayed a unimodal distribution and appeared
315 approximately Gaussian. In contrast, the distributions of the other parameters differed significantly from a
316 Gaussian distribution, as shown by the corresponding histograms. The distributions of all but one of these
317 parameters concentrated most of the probability mass at their upper limits. ~~Parameters Q_{50} , D_{50} , K_q and Q_A (Fig.~~
318 ~~2B-2E) in the PM model and parameters D_{50} , K_q in the SW model (Fig. 2K-2L) exhibited relatively large~~
319 ~~uncertainty reductions. However, the histograms~~ The exception was parameter b_1 for the SW model (Fig. 2k),
320 which clearly does not follow a normal distribution with most of the mass ~~obviously deviated from normality~~
321 ~~and tended to~~ concentrated in the lower bounds. ~~When the upper limits of these parameters were decreased,~~
322 ~~similar histograms were reached (not shown) and still did not show statistically meaningful distributions.~~ In
323 contrast, Q_{50} was not only poorly constrained (Fig. 2jg) but was also the upper edge-hitting parameter in the SW
324 model. ~~In addition~~ Moreover, the corresponding distributions of the same parameter in different models were
325 slightly different; ~~for~~ For example, the mean of g_{max} in the PM model (0.04 mm s^{-1}) was less than that in the
326 SW model (0.01 mm s^{-1}) (Fig. 2A-2a and 2If, Table 2), except that D_{50} in the PM and SW models and β_2 in

the PT-FC and AA models exhibited similar regions. It is interesting to observe that the distribution of K_a in PM model (Figure 2e) has a truncated distribution with highest probability mass at the upper bound, whereas the distribution of K_a in the SW model (Figure 2j) tends to become approximately normal. Overall, the marginal posterior probability density function of most of the individual parameters occupied only a relatively small region compared with the uniform prior distributions, and exhibited relatively large uncertainty reduction.

3.2 Performance of the models

The performances of each of the four evaporation-ET models were evaluated during-over the course of the whole season in 2014. The calibrated parameters of the four models were used and individual ET models were run to estimate the half-hourly λ ET values. Table 3 summarizes the statistical results for the performance of the models using regression line slope, R^2 , RMSE, MBE, IA, and EF as shown in Table 3. The regressions between measured and modeled λ ET values and MBE are shown in Figures 3 and 4, respectively.

In general, the four models produced slightly better fits to the measured λ ET for all the seasons with R^2 larger than 0.75 (Fig. 3). However, obvious discrepancies in the predictions made by among the models were detected by comparing measured and modeled λ ET. According to the regression line slope and MBE, the PM model overestimated ET by 1% with a MBE of -9.52 W m^{-2} , and the SW model overestimated ET by 5% with a relatively higher MBE of -19.07 W m^{-2} compared to the PM model. The PT-FC and AA models tended to underestimate λ ET by 9% and 8% with an MBE of 25.42 and 23.29 W m^{-2} , respectively. From a comparison between the slope and MBE, the PM model performance was higher than that of the SW, PT-FC and AA other three models, with a slope almost equal to 1 and with relatively lower MBE. The SW model was ranked second, while performance of the AA model was slightly higher comparable to that of the PT-FC model, but slightly higher, and was ranked third. However, if R^2 , RMSE, IA, and EF were used to evaluate the model performances, the SW model had the best overall performance with $R^2=0.83$, $\text{RMSE}= 76.34 \text{ W m}^{-2}$, $\text{IA} = 0.95$ and $\text{EF} = 0.79$. The second-best model was the PM model, with $R^2=0.76$, $\text{RMSE}= 85.38 \text{ W m}^{-2}$, $\text{IA}=0.93$ and $\text{EF}=0.74$, and the PT-FC performance was ranked third with $R^2=0.75$, $\text{RMSE}= 94.39 \text{ W m}^{-2}$, $\text{IA}=0.92$ and $\text{EF}=0.68$, while the AA model ranked fourth with $R^2=0.75$, $\text{RMSE}= 95.09 \text{ W m}^{-2}$, $\text{IA}=0.92$ and $\text{EF}=0.67$. Based on the analysis of these traditional statistical criteria error metrics, the performances of the PT-FC and AA models yielded similar results. The observed and modeled λ ET for the four ET models were tightly grouped along the regression lines (Figure 3), and the PT-FC and AA models had similar modeled ET values with a similar degree of point scattering along the regression lines (Figure 3c-3d).

Figure 4 shows that large seasonal variations were exist arise in MBE for the four ET models. From the variations of in the MBE, the estimated λ ET values for all models were generally lower than the measured values before the early jointing stage-of maize growth (DOY 154-177, left dashed line) and after the late maturity stage (DOY 256-265, right dash line) with the corresponding $\text{LAI} < 2.5 \text{ m}^2 \text{ m}^{-2}$. More positive MBE values for the PT-FC and AA models after the late maturity stage indicated their underestimated performances; however, these estimations appeared even more consistent with a symmetrical scattering of points along the 0-0 line (Figure 4c, 4d) during DOY 177-256 with $\text{LAI} > 2.5 \text{ m}^2 \text{ m}^{-2}$.

364 3.3 Comparison of the models using BME

365 Since there ~~was is currently~~ no theoretical method ~~so far~~ for selecting power posterior β values, we
366 determined these values using empirical but straightforward methods. For any ~~different~~ power coefficient of
367 $\beta \in [0, 1]$, a sample was drawn from the distribution p_β (Eq. 25) through running DREAM. ~~Although adding~~
368 ~~more β_k values might improve the BME estimation, this was not done because of the computational cost. For~~
369 ~~each β_k value, at least 150,000 DREAM simulations were large enough to ensure convergence. Figure 5 showed~~
370 ~~shows the evolution of $\ln p(D|\theta, M)$ for the four models as a function of β for a dataset covering the entire~~
371 ~~period. The potential values of the PM model increased from -6533.02 (the logarithm of the prior likelihood) to~~
372 ~~-6290.71 and the potential values increased from -6544.49 to -6016.17 for the SW model. In addition, the~~
373 ~~potential values increased from -6708.02 to -6361.76 for the PT-FC model and from -7732.98 to -7033.32 for~~
374 ~~the AA model. Table 3 showed that the estimated BME is -6300.5 natural log units (nits) for the PM model,~~
375 ~~-6025.1 nits for the SW model, -6366.8 nits for the PT-FC model, and -7042.8 nits for the AA model. The BME~~
376 ~~for the SW model was substantially larger than that for the other three models, and the BME for the AA model~~
377 ~~was the smallest. The BME-based model ranking (from the best to the worst) is SW, PM, PT-FC, and AA.~~
378 ~~Although the parameters of the PM model were less than for the SW model, the potential evolution of the SW-~~
379 ~~model was substantially different to that of the PM model. In summary, the PT-FC and AA models, which~~
380 ~~consisting of the same number of parameters, had similar potential patterns of evolution with the respect to the~~
381 ~~coefficient β_k . Although adding more β_k values may improve the BME estimation, it was not undertaken because~~
382 ~~of the computational cost. For each β_k value, 150,000 DREAM simulations were large enough to ensure~~
383 ~~convergence. The results illustrate that with the addition of parameters, the model complexity and the model~~
384 ~~performance are both increased.~~

385 4. Discussion

386 4.1 Parameter uncertainty analysis

387 With regard to the efficiency of the DREAM algorithm, the acceptance rates of the PM (15.33%) and SW
388 (18.94%) models were much higher than those obtained by some Markov Chain Monte Carlo (MCMC)
389 algorithms ~~which that have been used in the previous studies. , like 0.01902% in the population Monte Carlo~~
390 ~~sampling algorithm~~ (Sadegh et al., 2014). ~~This was is a large improvement in search efficiency, which in large~~
391 ~~part resulted resultes from its ability to sample one or groups of variable in turn. Furthermore, this the DREAM~~
392 ~~algorithm method ran runs multiple chains in parallel and adaptively updateds the scale and orientation of the~~
393 ~~proposal proposed distribution~~ (Vrugt et al., 2008). ~~Therefore, the DREAM scheme substantially improved~~
394 ~~improves not only the convergence, but also its sampling efficiency for ET models. The posterior parameter~~
395 ~~bounds exhibit a larger reduction using the DREAM algorithm compared with other studies using the~~
396 ~~Metropolis-Hasting algorithm. This demonstrates that DREAM could efficiently handle problems involving~~
397 ~~high-dimensionality, multimodality, nonlinearity.~~

398 The results showed that the assumed prior uncertainty ranges from most parameters in the four models
399 were significantly reduced. This indicates that the observed ET data contained sufficient information to estimate

400 these parameters. Surface conductance g_s and modeled ET in the PM model are relatively insensitive to Q_{50} , D_{50}
401 and K_q . Hence, these parameters could not be well constrained, and further relaxing the ranges for these
402 parameters could not result in physically realistic behavior of the model. The calculation of g_s^c —in the SW
403 model is the same as in the PM model, and thus, g_s^c and modeled ET in the SW model are also insensitive to
404 parameters of Q_{50} , D_{50} , K_q . Therefore, these three parameters were also not well constrained in the SW model. In
405 addition, the uncertainties present in the edge-hitting parameters, may be the outcome of model biases or
406 EC-measured ET data errors, or the characteristic time scale of parameters governing the processes affecting ET
407 is not exactly on the order of half-hours (Braswell et al., 2005). For example, Q_{50} and D_{50} govern changes in
408 visible radiation flux and the humidity deficit at which stomatal conductance is half its maximum value,
409 respectively, and these parameters may change over a shorter or longer time scale than half-hours.

410 The ecophysiological parameter g_{max} is a variable in the g_s^c —equation in both the PM and SW models, but
411 this parameter is sensitive to g_s^c and has a significant impact on the evaluated ET. ~~—in both the PM and SW~~
412 ~~models, the maximum stomatal conductance of leaves, and the soil surface resistance parameters b_1 and b_2 in the~~
413 ~~SW models, all had large influences on the evaluated ET. Thus, their effects were~~ is relatively independent
414 compared to the other meteorological parameters in the models, and therefore this parameter was well specified
415 in the PM and SW models. The posterior mean value of g_{max} (0.04 m s^{-1}) in the PM model from our study was
416 close to that (0.05 m s^{-1}) reported in northwestern China (Li et al., 2013; Zhu et al., 2014), but g_{max} (0.01 m s^{-1})
417 in the SW model was less than the reported value. ~~The estimated posterior means for b_1 and b_2 were different (b_1~~
418 ~~$= 9.3$, $b_2 = 6.2$) from those for maize suggested by Zhu et al. (2014) using the same equation of soil surface~~
419 ~~resistance (r_s^s). Though Zhu et al. (2014) concluded that the responses of g_s^c to VPD and LAI calculated using~~
420 ~~the modified Leuning model were close to those using Jarvis model (Jarvis, 1976), Li et al. (2015) showed that~~
421 ~~the performance of PM model was different using the two canopy resistance formula. Therefore, the different~~
422 ~~results of parameters b_1 and b_2 between our study and the previous study by Zhu et al. (2014) were mainly due~~
423 ~~to the usage of different canopy conductance models.—Parameter β_1 was well constrained in the PT-FC and AA~~
424 ~~models because it was relatively independent and did not directly relate to other observed variables.~~

425 For edge hitting parameters, their uncertainties may be the outcome of model biases or EC measured ET
426 data, or the characteristic time scale of parameters govern processes that was not exactly on the order of
427 half hours (Braswell et al., 2005). For example, Q_{50} and D_{50} govern changes in visible radiation flux and
428 humidity deficit at which stomatal conductance at its half maximum value, which may change over a shorter or
429 longer time scale rather than half hours. K_q was another parameter that cannot be well constrained, and this may
430 be resulted from either the estimated ET was insensitive to these parameters, or there were correlations between
431 the parameters. We expected a complementary correlative relationship between the visible radiation flux and
432 extinction coefficient for shortwave radiation, which indicated that the information in EC-measured ET data was
433 insufficient to separate these parameters, and therefore the parameters cannot be constrained separately.—

434 The sensitive parameters (such as g_{max} , b_1 and b_2) were just corresponding to the well constrained
435 parameters. Therefore, the major parameters in PM and SW models were well optimized, except that several

436 parameters (Q_{so} and K_e) appeared to be not well constrained. In addition, the posterior parameter bounds
437 exhibited a larger reduction using the DREAM algorithm compared with other studies using the Metropolis-
438 Hasting algorithm (Zhu et al., 2014). This further demonstrated that DREAM can efficiently handle problems
439 involving high dimensionality, multimodality, nonlinearity, and local optima.

440 Parameter K_a implicitly appears in the surface conductance equation (Eq.2) in PM model and K_a is
441 insensitive to g_s and modeled ET (Leuning et al., 2008). In contrast, K_a is contained in the equation of net
442 radiation flux into the substrate (Eq.10) in the SW model. This parameter can explicitly partition the total
443 available energy into that absorbed by the canopy and by the soil in the SW model. An analysis of equation (10),
444 found that the variation of K_a could not only account for the extinction effect but also correct the energy forcing
445 data errors. This also meant that the estimated value of K_a using calibration data was actually not just the true
446 extinction coefficient, but also included the energy imbalance correction in the SW model. From this analysis,
447 we could see that K_a not only involved the distribution of energy between the canopy and the soil surface but
448 also the energy imbalance. Therefore, parameter K_a has a great influence on the performance of the SW model.
449 This is why K_a is poorly constrained in the PM model but well constrained in the SW model. To further illustrate
450 the insights regarding the influence of parameter K_a on the performance of the SW model, we calibrated the SW
451 model again and reran the model with a constant value of K_a . The results showed a significant reduction in
452 model performance when K_a was held constant. This implied that the main reason for the SW model
453 outperforming the PM model in our study was not only the more physically rigorous structure of the SW model
454 but also the key parameter K_a being well constrained in the SW model.

455 In general, parameters related to soil surface resistance in the SW model were well evaluated, while
456 parameters related to canopy surface resistance in PM and SW models were poorly estimated. Therefore, using a
457 reliable canopy surface resistance equation in the ET model was crucial for improving its performance. In
458 addition, in our study, the traditional approach was used to quantify the uncertainty, which assumed that the
459 uncertainty mainly arose because of ~~same from~~ the parameter uncertainty. However, this method ~~did not~~ cannot
460 explicitly consider errors in the input data and model structural inadequacies. This is unrealistic for real
461 applications, and it is desirable to develop a more reliable inference method to treat all sources of uncertainty
462 separately and appropriately (Vrugt et al., 2008). Moreover, simultaneous direct measurement by
463 micro-lysimeter of sap flow and daily soil evaporation will further help to constrain the model parameters.

464 4.2 Evaluation and selection of the models

465 In this study, the traditional statistical measures and BME were chosen to evaluate and compare the
466 performance of four ET models. From the respective composition of these measures, the statistical measures can
467 be divided into residual-based ~~measures-metrics~~ (such as regression slope and MBE) and squared-residual-based
468 measures (such as R^2 , RMSE, IA, and EF). The rankings of the models obtained using the same type of metric
469 (residual-based or squared-residual-based) are similar. ~~Table 3 shows the values evaluated by BME method,~~
470 ~~residual based and squared residual based measures. By comparison, the estimates obtained within the same~~
471 ~~measure (residual based or squared residual based) were congruent. For example, slope-Slope and MBE, for~~
472 ~~example, which have similar results in the~~ are both residual-based measures, produce identical rankings.

473 However, the rankings produced by metrics of different types are not the same. ~~results from different kind of~~
474 ~~measures were incongruent;~~ For example, the PM model ~~outperformed~~ ~~outperforms~~ the SW model according
475 to the residual-based ~~measures~~ ~~metrics~~, but the performance of the PM model ~~was~~ ~~is~~ worse than SW model based
476 on the squared-residual-based measures. The comparative analysis shows a consistency between BME and the
477 squared-residual-based ~~statistics~~ ~~metrics~~ (hence the residual-based metrics disagreed with the BME measures). ~~–~~
478 ~~whereas residual based criteria were obvious disagreement with the BME measures. It~~ This ~~revealed~~ ~~reveals~~ that
479 the more complex SW model ~~was~~ ~~is~~ the best model based on ~~the~~ BME and squared-residual-based statistics. The
480 rank order of overall performance of the models from best to worst ~~was~~ ~~is~~: SW, PM, PT-FC, and AA model.

481 Previous studies had shown that BME evaluated by ~~the~~ TI provided estimates similar to the true values, and
482 selected the true model if the true model was included within the candidate models (Marshall et al., 2005;
483 Lartillot and Philippe, 2006). Meanwhile, some ~~have~~ argued that Bayesian analysis would choose the simplest
484 model (Jefferys and Berger, 1992; Xie et al., 2011) because of the best trade-off between good fit with ~~the~~ data
485 and model complexity (Schöniger et al., 2014). In this case, the most complex SW model had the highest BME
486 and was chosen as the ~~model with the best performance~~ ~~behaved model~~. This ~~likely~~ ~~probably~~ resulted from the
487 fact that the complex SW model ~~was~~ ~~is~~ indeed the most reliable model among the alternative ET models and ~~can~~
488 ~~provide a good fit to justify its higher complexity~~. The SW model is a two-layer model, and simulates soil
489 evaporation and plant transpiration separately, whereas the PM model is a single-layer model in which the plant
490 transpiration and soil evaporation cannot be separated (Monteith, 1965). The PT-FC model ~~was~~ ~~is~~ a simplified
491 ~~model version~~ of the PM model, and ~~it~~ only ~~required~~ ~~requires~~ meteorological and radiation information
492 (Priestley and Taylor, 1972), whereas ~~the~~ AA model only ~~relied~~ ~~relies~~ on the feedback between actual ET and
493 potential ET (Brutsaert and Stricker, 1979). Based on these physical mechanisms and processes ~~that each offer~~
494 these ET models ~~take into account~~, the rank order of the models ~~was~~ ~~is~~ reasonable.

495 ~~The estimates showed that the maximum values of R^2 , IA and EF, and the minimum value of RMSE, all~~
496 ~~selected the most complex SW model as the best performing model.~~ The results indicated that the SW model
497 ~~was~~ ~~is~~ the best performing model ~~–evaluated by~~ in terms of squared-residual-based ~~measures~~ ~~metrics~~, which
498 ~~resulted~~ ~~resultes~~ from the ability of the model to fit the measured data, irrespective of model complexity. It was
499 interesting to note that both the squared-residual-based measures and the BME consistently yielded the same
500 rank order. Although the squared-residual-based ~~measures~~ ~~metrics~~ seemed to identify a reasonable rank order,
501 this ~~had~~ ~~has~~ ~~often~~ not been the case, since the simple traditional statistical measures were known to usually
502 provide a biased view of the efficacy of a model (Kessler and Neas, 1994; Legates and McCabe, 1999). In
503 addition, sensitivity to outliers ~~was~~ ~~is~~ associated with these ~~measures~~ ~~metrics~~ and leads to relatively high values
504 due to the squaring of the residual terms (Willmott, 1981). Furthermore, these traditional statistical ~~measures~~
505 ~~metrics~~ ignore ~~edges~~ the priors, without penalizing model complexity, which ~~was~~ ~~is~~ in fact used in a Bayesian
506 analysis. ~~The dimensionality (model's parameter space) not only affected model evaluation by BME (Schöniger~~
507 ~~et al., 2014) but it may also affect the evaluation using traditional statistical measures. Here, two dimensional~~
508 ~~models of~~ PT-FC and AA, provided identical estimates of R^2 and IA. This ~~was~~ ~~is~~ most likely because both ~~the~~
509 ~~PT-FC and AA~~ models had the same dimensions and a similar model structure, ~~whereas BME estimates remain~~
510 ~~well behaved for the two ET models~~. Marshall et al. (2005) argued that EF would provide an incorrect

511 conclusion, and Samani et al. (2018) suggested that RMSE ~~also would selected~~ the complex model as the best
512 performing model. ~~Thus, we deduced that SRB measures are also problematic.~~ As for slope and MBE, the
513 ~~rankings produced by these~~ residual-based ~~measures-metrics~~ were in obvious disagreement with the ~~one based~~
514 ~~on BME-measure~~. Part of the lower values of slope and MBE may be counter-balanced by the higher values of
515 slope and MBE, thus these criteria ~~s-provided~~ an erroneous and unreliable ~~evaluation of the models-evaluation~~.
516 Therefore, the squared-residual-based and residual-based measures were not certain to provide reasonable
517 results in terms of model ranking.

518 BME is a consistent model selection which tries to identify which of the models produced the observed
519 data. Conversely, nonconsistent model selection uses the available data to estimate which of the models might
520 be best in predicting the future data. In fact, the error metrics are essentially nonparsimonious model selection,
521 which is a special case of nonconsistent model selection, where only the goodness of fit is used for rating
522 models without penalizing the model complexity and thus lacking consistency for the selected model (Höge et
523 al., 2018). The consistency between BME and the squared-residual-based metrics only indicates that the optimal
524 model evaluated by BME would also provide the best predictions, and thus consistent model selection should
525 also be asymptotically efficient (Leeb & Pötscher, 2009; Shao, 1997).

526 4.3 Analysis of model-data mismatch

527 Conceptual and structural inadequacies of the hydrological model ~~andtogether with~~ measurement errors of
528 the model input (forcing) and output (calibration) data introduced errors in the estimated parameters and model
529 simulations (Laloy, 20142015). Hydrological systems ~~were-are~~ indeed heavily input-driven and errors in
530 forcing data can dramatically impair the quality of calibration results and model output (Bardossy and Das, 2008;
531 Giudice, 2015). Measurement errors ~~were-raiseoccured~~ for a variety of reasons, including unreasonable
532 gap-filling in rainy days; dew and fog; inadequate areal coverage of point-scale soil water measurement;
533 mechanical limitations of the EC system; and inaccurate measurements of wind-speed, soil water, radiation and
534 vapor pressure deficit. ET processe ~~was-is~~ described using equations that can only capture parts of the complex
535 natural processes and ~~the-any ET model structures-wereis~~ an inherent simplification of the real system. These
536 inadequacies can thus lead to biased parameters and implausible predictions.

537 In our study, the results indicated that the PM and SW models overestimated the half-hourly ET compared
538 to the measured ET. Several studies also indicated that ~~the-ET values-were-was~~ overestimated by the PM model
539 (Fisher et al., 2005; Ortega-Farias et al., 2006; Li et al., 2015) and the SW model (Li et al., 2013; Li et al., 2015;
540 Zhang et al., 2008). Possible reasons for the inaccurate estimates included the following: (1) Anisotropic
541 turbulence with weak vertical and strong horizontal fluctuation leads to energy imbalance. The total turbulent
542 heat flux was lower by ~10–30% compared to the available energy in many land surface experiments (Tsvang et
543 al., 1991; Beyrich et al., 2002; Oncley et al., 2007; Foken et al., 2010) and influx networks (Franssen et al.,
544 2010). Liang et al. (2017) also showed an energy imbalance result in the semiarid area in China, and indicated
545 that the energy balance closure ratio ranged from 0.52 to 0.90 during the ~~daytime~~, whereas it was about 0.25
546 ~~during night-timeat night~~. However, the measured ET only included vertical flux and not horizontal flux,
547 leading to the measured ET being lower than that of ~~modeled-ET predicted~~ by the PM and SW models using the

548 available energy. (2) The absence of a mechanistic representation of the physiological response to plant
549 hydrodynamics ~~cause-makes~~ it difficult for the available ET models to resolve the dynamics of intradaily
550 hysteresis, producing patterns of diurnal error, while the imbalance or lack of between-leaf water demand and
551 soil water supply imposes hydrodynamic limitations on stomatal conductance (Thomsen et al., 2013; Zhang et
552 al., 2014; Matheny et al., 2014). Li et al. (2015) also concluded that neglecting the restrictive effect of the soil
553 on water transport in empirical canopy resistance equations can result in large errors in the partial canopy stage.
554 However, these equations can estimate ET accurately under the full canopy stage (Alves and Pereira, 2000;
555 Katerji and Rana, 2006; Katerji et al., 2011; Rana et al., 2011). Li et al. (2015) showed that the PM model
556 combined with canopy resistance overestimated maize ET during the partial and dense canopy stages by 16%
557 and 13%, respectively. Moreover, in a study of ET in vineyards, Leuning (2008) found that the PM model
558 coupled with canopy resistance overestimated ET during the entire growth stage by 29%.

559 The ~~estimated-estimates for~~ ET ~~for-produced by~~ the PT-FC and AA models ~~was-were~~ generally lower than
560 the measured values during the entire season. In addition, the four models also underestimated ~~the~~-ET during
561 periods of partial cover ($LAI < 2.5 \text{ m}^2 \text{ m}^{-2}$). ~~Especially during the late maturity stage, †~~The PT-FC and AA
562 models consistently underestimated ET, ~~especially during the late maturity stage, and provided the worst-~~
563 ~~simulated ET.~~ The underestimation probably resulted from the following: (1) Non classical situations, such as
564 the oasis effect, may occur in the study area. Strong evaporation from the moist ground and plants results in
565 latent heat cooling. However, this upward latent heat flux was opposed by a downward sensible heat flux from
566 the warm air to the cool ground, and thus the latent heat flux was positive while the sensible heat flux is
567 negative. Therefore, the latent heat flux can be greater in magnitude than the solar heating, because of the
568 additional energy extracted from the warm air by evaporation (Stull, 1988). (2) ~~The~~ Lack of mechanistic
569 representation of rainfall interception in ET models probably ~~also-~~led to inaccurate simulation ~~on-for periods~~
570 ~~soon~~~~shortly~~ after rainy days. Bohn and Vivoni (2016) found that evaporation of canopy interception accounted
571 for 8% of the annual ET across the North American monsoon region.

572 Comparing the AA and PT-FC models, the former includes forcing data of available radiation, soil water
573 content and relative humidity, but the PT-FC model only requires available radiation and soil water content and
574 is independent of relative humidity. However, the similar statistical results and similar degrees of MBE scatter
575 indicate that relative humidity has little influence on the AA model simulation. The consistent and consecutive
576 underestimation of ET by the PT-FC and AA models during the late maturity stage show that the model-data
577 disagreement is not caused by regional advection and rainfall interception, because atmospheric processes and
578 thermally-induced circulation can only occur at certain times and during certain days. Therefore, we think that
579 the consistent underestimation of ET by the PT-FC and AA models results primarily from conceptual and
580 structural inadequacies, energy imbalance, and soil water stress. Although the PM and SW models share a
581 common theoretical basis and the PT-FC model is a simplification of the PM model, these models perform
582 significantly differently. Part of the overestimation of ET by the PM and SW models, caused by coupling with
583 the canopy resistance, may be offset by underestimation caused by energy imbalance and soil water stress.
584 However, underestimation of ET by the PT-FC and AA models cannot be counterbalanced by overestimation
585 during the later maturity stage because the PT-FC and AA models are independent of the canopy resistance.

586 Consequently, the half-hourly patterns of errors in the estimates of ET by the PM and SW models are
587 characterized by symmetry and a low degree of scatter, but the PT-FC and AA models exhibit consistently
588 asymmetrical error patterns. By contrast, other studies showed that the PM model (Kato et al., 2004) and the SW
589 model (Chen et al., 2015) underestimated half-hourly ET. As for the PT-FC and AA models, some studies
590 reported that the PT-JPL (Zhang et al., 2017) and the AA model showed an overall poor performance (Zhang et
591 al., 2017). While other studies have indicated that the AA method performed well for both maize and canola
592 crops (Liu et al., 2012). Therefore, the performance of the four ET models appears to vary not only for different
593 crops and locations but also for different meteorological, physiological and soil conditions. Moreover, the
594 performance is also related to the stage of crop growth. Note that these conclusions about the ET models
595 evaluation are derived from traditional error metrics rather than those based on BME model selection. It would
596 be desirable to use available data from other study areas or from other crops for BME-based model selection to
597 confirm whether the SW model is the optimal model under other conditions.

598 Overall, combined with the parameter uncertainty analysis described in Section 4.1, we conclude that
599 energy imbalance and energy interaction between canopy and soil surface have a greater impact on the model
600 performance. And thus, explicitly treating of energy error, and incorporating the elements of existing hydrologic
601 theory about energy interaction between canopy and surface or conceptually correcting the energy interaction
602 are a practicable option for model improvement and application.

603 5. Conclusions

604 This study illustrated the application of the Bayesian approach for on the statistical analysis and model
605 selection of four widely used ET models. The results showed that the DREAM algorithm successfully reduced
606 the assumed prior uncertainties for most of the parameters in the four models. In the model calibration, the key
607 parameters which had a significant influence on ET simulations were well constrained. The main reasons for the
608 outperforming of SW model were its physically rigorous structure and the extinction coefficient parameter,
609 which is sensitive and has a significant impact on the performance of the model, being well constrained. –

610 ~~BME can be used to rank the alternative models in our study, although numerical evaluation of BME is~~
611 ~~computationally expensive particularly for high dimensional models.~~ BME is a consistent model selection
612 ~~to which tries to identify the best fitting to the observed data.~~ Although the squared-residual-based metrics,
613 including R^2 , IA, RMSE, and EF, produced a ranking identical to that of BME, it must be noted that these
614 squared-residual-based metrics do not allow using prior information and do not penalize the model complexity
615 when comparing the models. Therefore, some cautions are needed when using these statistical methods to
616 compare different models.

617 The model–data mismatches discrepancies were analyzed to facilitate model improvement after using
618 Bayesian model calibration and comparison. The results indicated that ~~model–data mismatches~~ the discrepancies
619 ~~arose are~~ mainly as a result of energy imbalance caused by anisotropic turbulence, additional energy induced by
620 advection processes, the absence of a mechanistic representation of the physiological response to plant
621 hydrodynamics and the energy interaction between canopy and surface. Among these causes, energy imbalance
622 and additional energy are related to forcing data errors rather than to an unreasonable model structure. Thus,
623 understanding the process of the physiological response to plant hydrodynamics and the interaction between

624 **canopy and surface** is essential for improving the performance of evapotranspiration models. Overall, the
625 applications of Bayesian calibration, Bayesian model evaluation and analysis of model–data discrepancies in our
626 study, provide a promising framework for reducing uncertainty and improving the performance of ET models. It
627 would be desirable to confirm whether the SW is the optimal model using data of other crops.or other climate
628 regions.

629 **Author contribution**

630 Guoxiao Wei and Xiaoying Zhang designed the experiments. Ning Yue and Fei Kan carried them out.
631 Ming Ye developed the model selection scheme. Guoxiao Wei performed the simulations. Guoxiao Wei and
632 Xiaoying Zhang prepared the manuscript with contributions from all co-authors.

633 **Competing interests**

634 The authors declare that they have no conflict of interest.

635 **Acknowledgments**

636 We thank Ying Guo, Huihui Dang, Jun Dong for the data collection and analysis. This work was funded by
637 the National Natural Science Foundation of China Nos. 41471023 and Nos.41702244. The third author was
638 supported in part by DOE Early Career Award DE-SC0008272 and National Science Foundation-Division of
639 Earth Science Grant 1552329. All observed data used in this study are from Heihe Watershed Allied Telemetry
640 Experimental Research (HiWATER). We thank all the staff who participated in HiWATER field campaigns.
641 Considerate and helpful comments by anonymous reviewers have considerably improved the manuscript.

642 **References**

- 643 Allen, R. G., Perista, L. S., Raes, D., and Smith, M.: Crop Evapotranspiration-Guidelines for Computing Crop
644 Water Requirements; FAO Irrigation and Drainage apers-56, FAO-Food and Agriculture Organization
645 of the United Nations, Rome, 1998.
- 646 Alves, I. and Pereira, L. S.: Modeling surface resistance from climatic variables? *Agric. Water Manage.*, 42,
647 371–385, 2000.
- 648 Aubinet, M., Grelle, A., Ibrom, A., Rannik, Ü., Moncrieff, J., and Foken, T.: Estimates of the annual net carbon
649 and water exchange of forests: the euroflux methodology, *Adv. Ecol. Res.*, 30(1), 113-175, 2000.
- 650 Baldocchi, D. D.: Assessing the eddy covariance technique for evaluating carbon dioxide exchange rates of
651 ecosystems: past, present and future, *Glob. Change. Biol.*, 9, 479–492, 2003.
- 652 Bardossy, A., and Das, T.: Influence of rainfall observation network on model calibration and application,
653 *Hydrol. Earth Syst. Sci.*, 12(1), 77-89, doi:10.5194/hess-12-77-2008, 2008.
- 654 Barton, I. J.: A Parameterization of the Evaporation from Nonsaturated Surfaces, *J. Appl Meteorol.*, 18(1), 43-47,
655 1979.
- 656 Beyrich, F., Richter, S. H., Weisensee, U., Kohsiek, W., Lohse, H., de Bruin, H. A. R., Foken, T., Göckede, M.,
657 Berger, F., Vogt, R., and Batchvarova, E.: Experimental determination of turbulent fluxes over the
658 heterogeneous litfass area: selected results from the litfass-98 experiment, *Theor. Appl. Climatol.*,
659 73(1–2), 19–34, doi:10.1007/s00704-002-0691-7, 2002.
- 660 Bohn, T. J., and Vivoni, E. R.: Process-based characterization of evapotranspiration sources over the North
661 American monsoon region, *Water Resour. Res.*, 52, 358–384, doi:10.1002/2015WR017934, 2016.
- 662 Bonan, G.: *Ecological climatology: concepts and applications*, Cambridge University Press, 2008.
- 663 Braswell, B. H., Sacks, W. J., Linder, E., and Schimel, D. S.: Estimating diurnal to annual ecosystem parameters
664 by synthesis of a carbon flux model with eddy covariance net ecosystem exchange observations, *Global.*
665 *Change. Biol.*, 11, 335–355, 2005,

- 666 Brutsaert, W., and Han, S.: An advection-aridity approach to estimate actual regional evapotranspiration, *Water*
667 *Resour. Res.*, 15(2), 443-450, 1979.
- 668 Brutsaert, W.: *Hydrology: An Introduction*. Cambridge University Press, Cambridge, 2005.
- 669 Chen, D. Y., Wang, X., Liu, S. Y., Wang, Y. K., Gao, Z. Y., Zhang, L. L., Wei, X. G., and Wei, X. D.: Using
670 Bayesian analysis to compare the performance of three evapotranspiration models for rainfed jujube
671 (*Ziziphus jujuba* Mill.) plantations in the Loess Plateau, *Agr. Water Manage.*, 159, 341–357, 2015.
- 672 Elshall, A. S., Ye, M., Pei, Y., Zhang, F., Niu, G. Y., and Barron-Gafford, G. A.: Relative model score: A scoring
673 rule for evaluating ensemble simulations with application to microbial soil respiration modeling, *Stoch.*
674 *Env. Res. A.*, 1-11, DOI: 10.1007/s00477-018-1592-3, 2018.
- 675 Ershadi, A., McCabe, M. F., Evans, J. P., Chaney, N. W., and Wood, E. F.: Multi-site evaluation of terrestrial
676 evaporation models using fluxnet data, *Agric. For. Meteorol.*, 187(8), 46-61, 2014.
- 677 Ershadi, A., McCabe, M. F., Evans, J. P., and Wood, E. F.: Impact of model structure and parameterization on
678 Penman–Monteith type evaporation models, *J. Hydrol.*, 525, 521–535, 2015.
- 679 Fisher, J. B., DeBiase, T. A., Qi, Y., Xu, M., and Goldstein, A. H.: Evapotranspiration models compared on a
680 Sierra Nevada forest ecosystem, *Environ. Model. Softw.*, 20 (6), 783–796, 2005.
- 681 Flint A. L., Childs, S. W.: Use of the Priestley–Taylor evaporation equation for soil water limited conditions in a
682 small forest clearcut, *Agric. For. Meteorol.*, 56(3–4), 247–260, 1991.
- 683 Foken, T., Mauder, M., Liebethal, C., Wimmer, F., Beyrich, F., Leps, J. P., Raasch, S., DeBruin, H. A. R.,
684 Meijninger, W. M. L., and Bange, J.: Energy balance closure for the LITFASS-2003 experiment, *Theor.*
685 *Appl. Climatol.*, 101(1–2), 149–160, doi:10.1007/s00704-009-0216-8, 2010.
- 686 Franssen, H. J. H., Stöckli, R., Lehner, I., Rotenberg, E., and Seneviratne S. I.: Energy balance closure of
687 eddy-covariance data: A multisite analysis for European FLUXNET stations, *Agric. For. Meteorol.*,
688 150(12), 1553–1567, doi:10.1016/j.agrformet.2010.08.005, 2010.
- 689 Gelman, A., and Rubin, D. B.: Inference from iterative simulation using multiple sequences, *Stat. Sci.*, 7, 457–
690 472, 1992.
- 691 Gelman, A.: Simulating normalizing constants: From importance sampling to bridge sampling to path sampling,
692 *Stat. Sci.* 13, 163–185, 1998.
- 693 Giudice, D., Albert, C., Rieckermann, J., and Reichert, P.: Describing the catchment-averaged precipitation as a
694 stochastic process improves parameter and input estimation, *Water Resour. Res.*, 52, 3162–3186,
695 doi:10.1002/2015WR017871, 2016.
- 696 Höge, M., Wöhling, T., and Nowak, W.: A primer for model selection: The decisive role of model complexity.
697 *Water Resour. Res.*, 54, 1688–1715, doi.org/10.1002/2017WR021902, 2018.
- 698 Jefferys, W. H., and Berger, J. O.: Sharpening Ockham’s razor on a Bayesian strop, *Am. Sci.*, 89, 64-72, 1992.
- 699 Kashyap, R. L.: Optimal choice of AR and MA parts in autoregressive moving average models, *IEEE Trans.*
700 *Pattern Anal. Mach. Intell.*, 4(2), 99–104, 1982.
- 701 Katerji, N., and Rana, G.: Modelling evapotranspiration of six irrigated crops under Mediterranean climate
702 conditions, *Agric. For. Meteorol.*, 138, 142–155, 2006.
- 703 Katerji, N., Rana, G., Fahed, S.: Parameterizing canopy resistance using mechanistic and semi-empirical
704 estimates of hourly evapotranspiration: critical evaluation for irrigated crops in the Mediterranean,
705 *Hydrol. Process.*, 25, 117–129, 2011.
- 706 Kato, T., Kimura, R., and Kamichika, M.: Estimation of evapotranspiration, transpiration ratio and water-use
707 efficiency from a sparse canopy using a compartment model, *Agric. Water Manage.*, 65, 173–191,
708 2004.
- 709 Kessler, E., and Neas, B.: On correlation, with applications to the radar and raingage measurement of rainfall,
710 *Atmos. Res.*, 34, 217-229, 1994.
- 711 Laloy, E., Linde, N., Jacques, D., and Vuigt, J. A.: Probabilistic inference of multi-Gaussian fields from indirect
712 hydrological data using circulant embedding and dimensionality reduction, *Water Resour. Res.*, 51,
713 4224–4243, doi:10.1002/2014WR016395, 2015.
- 714 Lartillot, N., and Philippe, H.: Computing Bayes factors using thermodynamic integration, *Syst. Biol.*, 55(2),
715 195-207, 2006.
- 716 Leeb, H., and Pötscher, B. M.: *Model selection*, Berlin, Germany: Springer., pp, 889–925,
717 doi.org/10.1007/978-3-540-71297-839, 2009.
- 718 Legates, D. R., and McCabe, G. J.: Evaluating the use of “goodness-of-fit” measures in hydrologic and
719 hydroclimatic model validation, *Water Resour. Res.*, 35, 233–241, 1999.
- 720 Leuning, R., Zhang, Y. Q., Rajaud, A., Cleugh, H., and Tu, K.: A simple surface conductance model to estimate

- 721 regional evaporation using MODIS leaf area index and the Penman–Monteith equation, *Water Resour.*
722 *Res.*, 44, W10419, doi.org/10.1029/2007WR006562, 2008.
- 723 Liang, J., Zhang, L., Cao, X., Wen, J., Wang, J., and Wang, G.: Energy balance in the semiarid area of the Loess
724 Plateau, China, *J. Geophys. Res. Atmos.*, 122, 2155–2168, doi:10.1002/2015JD024572, 2017.
- 725 Li, S., Kang, S., Zhang, L., Ortega-Farias, S., Li, F., Du, T., Tong, L., Wang, S., Ingman, M., and Guo, W.:
726 Measuring and modeling maize evapotranspiration under plastic film-mulching condition, *J. Hydrol.*,
727 503, 153–168, 2013.
- 728 Li, S., Zhang, L., Kang, S., Tong, L., Du, T., Hao, X., Zhao, P.: Comparison of several surface resistance models
729 for estimating crop evapotranspiration over the entire growing season in arid regions. *Agric. For.*
730 *Meteorol.* 208, 1-15, 2015.
- 731 Li, X., Cheng, G. D., Liu, S. M., Xiao, Q., Ma, M. G., Jin, R., Che, T., Liu, Q. H., Wang, W. Z., Qi, Y., Wen, J.
732 G., Li, H. Y., Zhu, G. F., Guo, J. W., Ran, Y. H., Wang, S. G., Zhu, Z. L., Zhou, J., Hu, X. L., and Xu, Z.
733 W.: Heihe Watershed Allied Telemetry Experimental Research (HiWATER): Scientific objectives and
734 experimental design, *B. Am. Meteorol. Soc.*, 94, 1145–1160, 2013.
- 735 Liu, S. M., Xu, Z. W., Wang, W. Z., Jia, Z. Z., Zhu, M. J., Bai, J., and Wang, J. M.: A comparison of
736 eddy-covariance and large aperture scintillometer measurements with respect to the energy
737 balance closure problem, *Hydrol. Earth Syst. Sci.*, 15, 1291–1306, doi:10.5194/hess-15-1291-2011,
738 2011.
- 739 Liu, G., Liu, Y., Hafeez, M., Xu, D., Vote, C.: Comparison of two methods to derive time series of actual
740 evapotranspiration using eddy covariance measurements in the southeastern Australia, *J. Hydrol.*, 454–
741 455 (4), 1–6, 2012.
- 742 Liu, P., Elshall, A. S., Ye, M., Beerli, P., Zeng, X., Lu, D., and Tao, Y.: Evaluating marginal likelihood with
743 thermodynamic integration method and comparison with several other numerical methods. *Water*
744 *Resour. Res.*, 52(2), 734-758, doi:10.1002/2014WR016718, 2016.
- 745 Marshall, L., Nott, D., and Sharma, A.: Hydrological model selection: A Bayesian alternative, *Water Resour.*
746 *Res.*, 41(10), 3092-3100, doi: 10.1029/2004WR003719, 2005.
- 747 Matheny, A. M., Bohrer, G., Stoy, P. C., Baker, I. T., Black, A. T., Desai, A. R., Dietze, M. C., Gough, C. M.,
748 Ivanov, V. Y., Jassal, R. S., Novick, K. A., Schäfer, K. V. R., and Verbeeck, H.: Characterizing the
749 diurnal patterns of errors in the prediction of evapotranspiration by several land-surface models: An
750 NACP analysis, *J. Geophys. Res. Biogeosci.*, 119(7), 1458-1473, 2014.
- 751 Monteith, J. L.: Evaporation and environment, *Symp. Soc. Exp. Biol.*, 19, 205–234, 1965.
- 752 Morison, J. I. L., Baker, N. R., Mullineaux, P. M., and Davies, W. J.: Improving water use in crop production,
753 *Philos. T. Roy. Soc. B.*, 363, 639–658, 2008.
- 754 Neal, R. M.: Markov chain sampling methods for Dirichlet process mixture models, *J. Comput. Graph. Stat.*, 9,
755 249–265, 2000.
- 756 Oncley, S. P., Foken, T., Vogt, R., Kohsiek, W., DeBruin, H., Bernhofer, C., Christen, A., Van Gorsel, E., Grantz,
757 D., and Feigenwinter, C.: The energy balance experiment EBEX-2000. Part I: Overview and energy
758 balance, *Boundary Layer Meteorol.*, 123(1), 1–28, doi:10.1007/s10546-007-9161-1, 2007.
- 759 Ortega-Farias, S., Olioso, A., Fuentes, S., and Valdes, H.: Latent heat flux over a furrow-irrigated tomato crop
760 using Penman–Monteith equation with a variable surface canopy resistance, *Agric. Water Manage.*, 82,
761 421–432, 2006.
- 762 Ortega-Farias, S., Poblete-Echeverria, C., and Brisson, N.: Parameterization of a two-layer model for estimating
763 vineyard evapotranspiration using meteorological measurements, *Agr. For. Meteorol.*, 150, 276–286,
764 2010.
- 765 Parlange, M. B., and Katul, G. G.: An advection-aridity evaporation model, *Water Resour. Res.*, 28 (1), 127-132,
766 1992.
- 767 Poblete-Echeverria, C., and Ortega-Farias, S.: Estimation of actual evapotranspiration for a drip-irrigated Merlot
768 vineyard using a three-source model, *Irrig. Sci.*, 28, 65–78, 2009.
- 769 Priestley, C. H. B., and Taylor, R. J.: On the assessment of surface heat flux and evaporation using large-scale
770 parameters, *Mon. Weather Rev.*, 100 (2), 81-92, 1972.
- 771 Rana, G., Katerji, N., Ferrara, R.M., and Martinelli, N.: An operational model to estimate hourly and daily crop
772 evapotranspiration in hilly terrain: validation on wheat and oat crops, *Theory Appl. Climatol.*, 103,
773 413–426, 2011.
- 774 Sadegh, M., and Vrugt J. A.: Approximate Bayesian Computation using Markov Chain Monte Carlo simulation:
775 DREAM(ABC), *Water Resour. Res.*, 50, 6767–6787, doi:10.1002/2014WR015386, 2014.

- 776 Samani, S., Ye, M., Zhang, F., Pei, Y. Z., Tang, G. P., Elshall, A. S., and Moghaddam, A. A.: Impacts of prior
777 parameter distributions on bayesian evaluation of groundwater model complexity, *Water Science &*
778 *Engineering.*, 11(2), 89-100, doi.org/10.1016/j.wse.2018.06.001, 2018.
- 779 Schöniger, A., Wohling, T., Samaniego, L., and Nowak, W.: Model selection on solid ground: Rigorous
780 comparison of nine ways to evaluate Bayesian model evidence, *Water Resour. Res.*, 50, 9484–9513,
781 doi:10.1002/2014WR016062, 2014.
- 782 Schwarz, G.: Estimating the dimension of a model, *Ann. Stat.*, 6(2), 461–464, doi:10.1214/aos/1176344136,
783 1978.
- 784 Shao, J.: An asymptotic theory for linear model selection, *Statistica Sinica*, 7(2), 221–242, 1997.
- 785 Shuttleworth, W. J., Gurney, R. J.: The theoretical relationship between foliage temperature and canopy
786 resistance in sparse crops, *Q. J. Roy. Meteorol. Soc.*, 116, 497–519, 1990.
- 787 Stannard, D. I.: Comparison of Penman-Monteith, Shuttleworth-Wallace, and modified Priestley-Taylor
788 evapotranspiration models for wildland vegetation in semiarid rangeland, *Water Resour. Res.*, 29 (5),
789 1379-1392, 1993.
- 790 Stull, R. B.: An introduction to boundary layer meteorology, Kluwer Academic Publ., 255pp, 1988.
- 791 Sumner, D. M., and Jacobs, J. M.: Utility of Penman–Monteith Priestley–Taylor reference evapotranspiration,
792 and pan evaporation methods to estimate pasture evapotranspiration, *J. Hydrol.*, 308 (1-4), 81-104,
793 2005.
- 794 Szilagyi, J., and Jozsa, J.: New findings about the complementary relationship based evaporation estimation
795 methods, *J. Hydrol.*, 354: 171–186, 2008.
- 796 Thomsen, J., Bohrer, G., Matheny, M. V., Ivanov, Y., He, L., Renninger, H., and Schäfer, K.: Contrasting
797 hydraulic strategies during dry soil conditions in *Quercus rubra* and *Acer rubrum* in a sandy site in
798 Michigan, *Forests.*, 4(4), 1106–1120, 2013.
- 799 Tsvang, L., Fedorov, M., Kader, B., Zubkovskii, S., Foken, T., Richter, S., and Zeleny, Y.: Turbulent exchange
800 over a surface with chessboardtype inhomogeneities, *Boundary Layer Meteorol.*, 55(1–2), 141–160,
801 1991.
- 802 Vinukollu R, K., Wood, E. F., Ferguson, C. R., and Fisher, J. B.: Global estimates of evapotranspiration for
803 climate studies using multi-sensor remote sensing data: evaluation of three process-based approaches,
804 *Remote Sens. Environ.*, 115(3), 801–823, 2011.
- 805 Vrugt, J. A., ter Braak, C. J. F., Clark, M. P. J., Hyman, M., and Robinson, B. A.: Treatment of input uncertainty
806 in hydrologic modeling: Doing hydrology backward with Markov chain Monte Carlo simulation, *Water*
807 *Resour. Res.*, 44, W00B09, doi:10.1029/2007WR006720, 2008.
- 808 Vrugt, J. A., ter Braak, C. J. F., Diks, C. G. H., Higdon, D., Robinson, B. A., and Hyman, J. M.: Accelerating
809 Markov chain Monte Carlo simulation by differential evolution with self-adaptive randomized
810 subspace sampling, *Int. J. Nonlinear Sci. Numer. Simul.*, 10(3), 273-290, 2009.
- 811 Webb, E. K., Pearman, G. I., and Leuning, R.: Correction of flux measurements for density effects due to heat
812 and water-vapor transfer, *Q. J. R. Meteorol. Soc.*, 106(447), 85–100, 1980.
- 813 Willmott, C. J.: On the validation of models, *Phys. Geogr.*, 2, 184-194, 1981.
- 814 Xie, W., Lewis, P. O., Fan, Y., Kuo, L., and Chen, M. H.: Improving marginal likelihood estimaton for Bayesian
815 phylogenetic model selection, *Syst. Biol.*, 60(2), 150-160, 2011.
- 816 Xu, C. Y., and Singh, V. P.: A review on monthly water balance models for water resources investigations, *Water*
817 *Resour. Manage.*, 12, 31-50, 1998.
- 818 Xu, Z. W., Liu, S. M., Li, X., Shi, S. J., Wang, J. M., Zhu, Z. L., Xu, T. R., Wang, W. Z., and Ma, M. G.:
819 Intercomparison of surface energy flux measurement systems used during the HiWATERUSOEXE, *J.*
820 *Geophys. Res.*, 118, 13140–13157, 2014.
- 821 Ye, M., Neuman, S. P., and Meyer, P. D.: Maximum likelihood Bayesian averaging of spatial variability models
822 in unsaturated fractured tuff, *Water Resour. Res.*, 40, W05113, doi:10.1029/2003WR002557, 2004.
- 823 Ye, M., Meyer, P. D., and Neuman, S. P.: On model selection criteria in multimodel analysis, *Water Resour. Res.*,
824 44, W03428, doi:10.1029/2008WR006803, 2008.
- 825 Zhang, B., Kang, S., Li, F., and Zhang, L.: Comparison of three evapotranspiration models to Bowen
826 ratio-energy balance method for vineyard in an arid desert region of northwest China, *Agr. Forest*
827 *Meteorol.*, 148: 1629–1640, 2008.
- 828 Zhang, X. Y., Liu, C. X., Hu, B. X., and Zhang, G. N.: Uncertainty analysis of multi-rate kinetics of uranium
829 desorption from sediments, *J. Contam. Hydrol.*, 156(1), 1-15, 2014.
- 830 Zhang, K., Ma, J., Zhu, G., Ma, T., Han, T., and Feng, L. L.: Parameter sensitivity analysis and optimization for

831 a satellite-based evapotranspiration model across multiple sites using moderate resolution imaging
832 spectroradiometer and flux data. *Journal of Geophysical Research: Atmospheres*, 122(1), 230-245,
833 2017.

834 Zhu, G. F., Su, Y. H., Li, X., Zhang, K., and Li, C. B.: Estimating actual evapotranspiration from an alpine
835 grassland on Qinghai–Tibetan plateau using a two-source model and parameter uncertainty analysis by
836 Bayesian approach, *J. Hydrol.*, 476, 42–51, 2013.

837 Zhu, G. F., Li, X., Su, Y. H., Zhang, K., Bai, Y., Ma, J. Z., Li, C. B., Hu, X. L., and He, J. H.: Simultaneously
838 assimilating multivariate data sets into the two-source evapotranspiration model by Bayesian approach:
839 Application to spring maize in an arid region of northwestern China, *Geosci. Model. Dev.*, 7(4), 1467–
840 1482, 2014.

841 **Appendix A: List of symbols and physical characteristics in ET models**

A	Available energy for the whole canopy (W m^{-2})
A_s	Available energy (W m^{-2})
R_n	Net radiation fluxes into the canopy (W m^{-2})
R_{ns}	Net radiation flux into the substrate (W m^{-2})
G	Soil heat flux (W m^{-2})
λET	Sum of the latent heat flux from the crop (λT) and soil (λE) (W m^{-2})
ET_c	Canopy transpiration (W m^{-2})
ET_s	Soil evaporation (W m^{-2})
C_c	Canopy resistance coefficient (dimensionless)
C_s	Soil surface resistance coefficient (dimensionless)
LAI	Leaf area index
Q_{50}	Visible radiation flux (W m^{-2})
D_{50}	Vapor pressure deficit (kPa)
D_a	Vapor pressure deficit at the reference height ($D_a=e_s-e_a$) (kPa)
Q_h	Flux density of visible radiation at the top of the canopy (W m^{-2})
K_q	Extinction coefficient
K_a	Extinction coefficient
f	Fraction of evaporation soil and total evaporation
λ	Latent heat of water evaporation (MJ kg^{-1})
Δ	Slope of the saturated vapour pressure curve (Pa K^{-1})
γ	Psychrometric constant (kPa K^{-1})
ρ	Density of air (kg m^{-3})
k	Karman constant (0.41)
e_s	Saturated vapor pressure (kPa)
e_a	Actual vapor pressure (kPa)
q^*	Saturation-specific humidity at air temperature (kg kg^{-1})
q	Specific humidity of the atmosphere (kg kg^{-1})
b_1	Empirical constant (s m^{-1})
b_2	Empirical constant (s m^{-1})
β_1	empirical constant
β_2	empirical constant
θ	Soil water content ($\text{m}^3 \text{m}^{-3}$)
θ_a	Critical water content at which plant stress starts ($\text{m}^3 \text{m}^{-3}$)
θ_b	Water content at the wilting point ($\text{m}^3 \text{m}^{-3}$)
θ_r	Residual soil water content ($\text{m}^3 \text{m}^{-3}$)
θ_s	Saturated water content ($\text{m}^3 \text{m}^{-3}$)
Θ	Relative water saturation

d	Zero plane displacement height (m)
z_m	Height of the wind speed and humidity measurements (3 m)
z_{0m}	Roughness length governing the transfer of momentum (m)
z_{0v}	Roughness length governing the transfer of water vapor (m)
h	Canopy height (m)
u_z	Wind speed at height z_m (m s^{-1})
g_a	Aerodynamic conductance (m s^{-1})
g_s	Surface conductance (m s^{-1})
g_{max}	Maximum stomatal conductance of leaves at the top of the canopy (m s^{-1})
g_s^c	Canopy conductance (m s^{-1})
r_a	Aerodynamic resistance (s m^{-1})
r_a^a	Aerodynamic resistance between canopy source height and a reference level (s m^{-1})
r_a^s	Aerodynamic resistance between the substrate and the canopy source height (s m^{-1})
r_a^c	Bulk boundary layer resistance of the vegetation element in the canopy (s m^{-1})
r_s^s	Surface resistance of the canopy (s m^{-1});
r_s^c	Bulk stomatal resistance of the canopy (s m^{-1})

842

843 **Appendix B: Bayesian inference and the DREAM algorithm**

844 The posterior probability distribution of the parameter is calculated by Bayes' theorem:

$$845 \quad \pi(\boldsymbol{\theta} | D, M) \propto \frac{\pi(\boldsymbol{\theta} / M) p(D | \boldsymbol{\theta}, M)}{p(D | M)} \quad (\text{A1})$$

846 where $\pi(\boldsymbol{\theta} / M)$ represents the prior density of $\boldsymbol{\theta}$ under model M ; $p(D | \boldsymbol{\theta}, M)$ is the joint likelihood of
847 model M and its parameters $\boldsymbol{\theta}$; and

$$848 \quad p(D | M) = \int p(D | \boldsymbol{\theta}, M) \pi(\boldsymbol{\theta} | M) d\boldsymbol{\theta} \quad (\text{A2})$$

849 is the marginal likelihood, or Bayesian model evidence (BME).

850 The likelihood function, $p(D | \boldsymbol{\theta}, M)$, used for parameter estimation, is specified according to the
851 distributions of observation errors. Error $e(t)$ in each observation $D(t)$ at time t is expressed by

$$852 \quad e(t) = D(t) - f(t) \quad (\text{A3})$$

853 . Assuming $e(t)$ follows a Gaussian distribution with a zero mean, and the likelihood function can be
854 expressed as

$$p(D|\theta) = \prod_{t=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{[e(t)]^2}{2\sigma^2}} \quad (\text{A4})$$

856 where n is the number of observations and σ represents the error variances.

857 In this study, we used the DREAM algorithm (Vrugt et al., 2008, 2009) to explore the ET models'
 858 parameter space and to estimate BME. The DREAM sampling scheme is an adaptation of the global
 859 optimization algorithm of a shuffled complex evolution metropolis (SCEM-UA). This algorithm was
 860 described in more detail in Vrugt et al. (2008, 2009).

861 List of Tables

862 **Table 1.** Prior distributions and parameter limits for the PM, SW, PT-FC and AA models. The values are
 863 derived from the literature.

864 **Table 2.** Maximum Likelihood Estimates (MLEs), Mean Estimates, 95% High-Probability Intervals
 865 (Lower Limit, Upper Limit).

866 **Table 3.** Slope and coefficient of determination (R^2) of regression between measured and modeled
 867 half-hourly evapotranspiration values, and statistics of root mean square error (RMSE), mean bias error
 868 (MBE), index of agreement (IA), model efficiency (EF) and Logarithm of BME for the four ET models.

869 List of Figures

870 **Figure 1.** Trace plots of the G-R statistic of Gelman and Rubin (Gelman and Rubin, 1992) using DREAM
 871 for the PM model (a) and (b) the SW model. Different parameters are coded with different colors. The
 872 dashed line denotes the default threshold used to diagnose convergence to a limiting distribution.

873 **Figure 2.** (a)-(e), (f)-(l), (m)-(n), and (o)-(p) show histograms for the PM (black), SW (cyan), PT-FC
 874 (magenta) and AA (orange) models, respectively. These histograms are constructed from all chains for each
 875 model and a total of 40,000 $\times N$ realizations are simulated using DREAM. The x axes represent the
 876 prespecified limits of the parameters.

877 **Figure 3.** Regressions between measured and modeled half-hourly ET values produced by different
 878 models from DOY 154 to DOY 270: (a) PM, (b) SW, (c) PT-FC and (d) AA. The regressions are: $Y =$
 879 $0.99X$ ($R^2 = 0.76$), $Y = 1.05X$ ($R^2 = 0.82$), $Y = 0.91X$ ($R^2 = 0.75$), and $Y = 0.92X$ ($R^2 = 0.75$) for the PM,
 880 SW, PT-FC and AA models, respectively.

881 **Figure 4.** Mean bias error (MBE) of predicted and observed ET values for (a) PM, (b) SW, (c) PT-FC and
 882 (d) AA models from DOY 154 to DOY 270. Parameters used for prediction are estimated by DREAM with
 883 the dataset for the calibration period from DOY 154 to DOY 202.

884 **Figure 5.** Variation of the mean posterior expectation of the potential y_k with β_k for the PM, SW, PT-FC
 885 and AA models.

886

887 **Table 1** Prior distributions and parameter limits for the PM, SW, PT-FC and AA models. The values are
 888 derived from the literature.

Parameter	Description	Prior range	PM	Prior for SW	Prior for PT and AA	References
-----------	-------------	-------------	----	--------------	---------------------	------------

		Lower	upper	Lower	upper	Lower	upper	
g_{max} (mm s ⁻¹)	maximum stomatal conductance	0	50	0	50			Kelliher et al. (1995)
Q_{50} (W m ⁻²)	visible radiation flux	10	50	10	50			Leuning et al. (2008)
D_{50} (kPa)	vapor pressure deficit	0.5	3	0.5	3			Leuning et al. (2008)
K_q	extinction coefficient	0	1	0	1			Leuning et al. (2008)
K_a	extinction coefficient	0	1	0	1			Leuning et al. (2008)
b_1 (s m ⁻¹)	empirical constant			4.5	11.3			Sellers et al. (1992)
b_2 (s m ⁻¹)	empirical constant			0	8			Sellers et al. (1992)
β_1	empirical constant					0.5	1.5	Flint et al. (1991); Barton. (1979)
β_2	empirical constant					0.1	10	

889

890 **Table 2** Maximum Likelihood Estimates (MLEs), Mean Estimates, 95% High-Probability Intervals
891 (Lower Limit, Upper Limit).

Parameter	Posterior for PM			Posterior for SW			Posterior for PT and AA		
	MLE	Mean	CI	MLE	Mean	CI	MLE	Mean	CI
g_{max} (mm s ⁻¹)	0.04	0.04	(0.03, 0.04)	0.01	0.01	(0.005, 0.012)			
Q_{50} (W m ⁻²)	49.96	48.52	(39.73, 49.74)	47.49	40.32	(11.02, 48.99)			
D_{50} (kPa)	3.00	2.87	(1.92, 2.97)	2.98	2.88	(2.26, 2.98)			
K_q	1.00	0.99	(0.911, 0.998)	0.99	0.88	(0.06, 0.98)			
K_a	1.00	0.98	(0.822, 0.995)	0.12	0.12	(0.074, 0.184)			
b_1 (s m ⁻¹)				4.51	4.57	(4.52, 4.96)			
b_2 (s m ⁻¹)				0.39	0.57	(0.07, 1.38)			
β_1							1.1 ^a 1.5 ^b	1.098 ^a 1.499 ^b	(1.06, 1.16) ^a (1.492, 1.499) ^b
β_2							10.00 ^a 10.00 ^b	9.75 ^a 9.94 ^b	(7.97, 9.95) ^a (9.44, 9.99) ^b

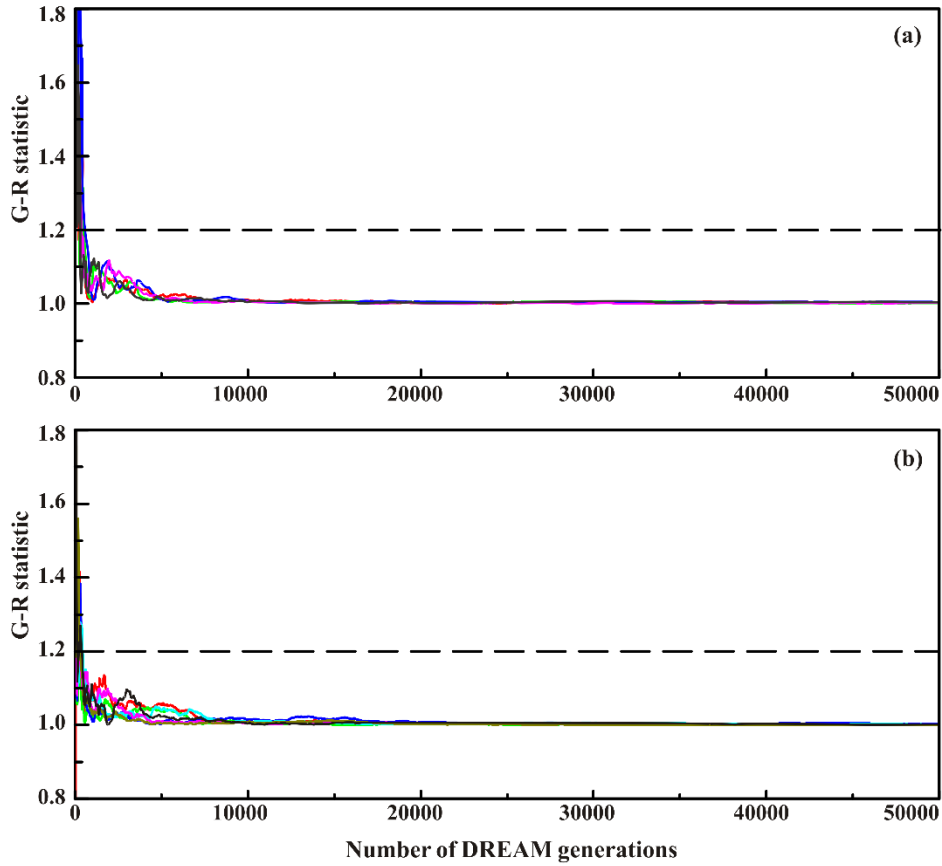
892 ^a PT-FC model; ^b AA model.

893 **Table 3** Slope and coefficient of determination (R^2) of regression between measured and modeled
894 half-hourly evapotranspiration values, and statistics of root mean square error (RMSE), mean bias error
895 (MBE), index of agreement (IA), model efficiency (EF) and Logarithm of BME for the four ET models.

Model	Slope	R^2	RMSE	MBE	IA	EF	BME
PM	1.01	0.76	85.38	-9.52	0.93	0.74	-6300.5

SW	1.05	0.82	76.34	-19.07	0.95	0.79	-6025.1
PT-FC	0.91	0.75	94.39	25.42	0.92	0.68	-6366.8
AA	0.92	0.75	95.09	23.29	0.92	0.67	-6390.3

896

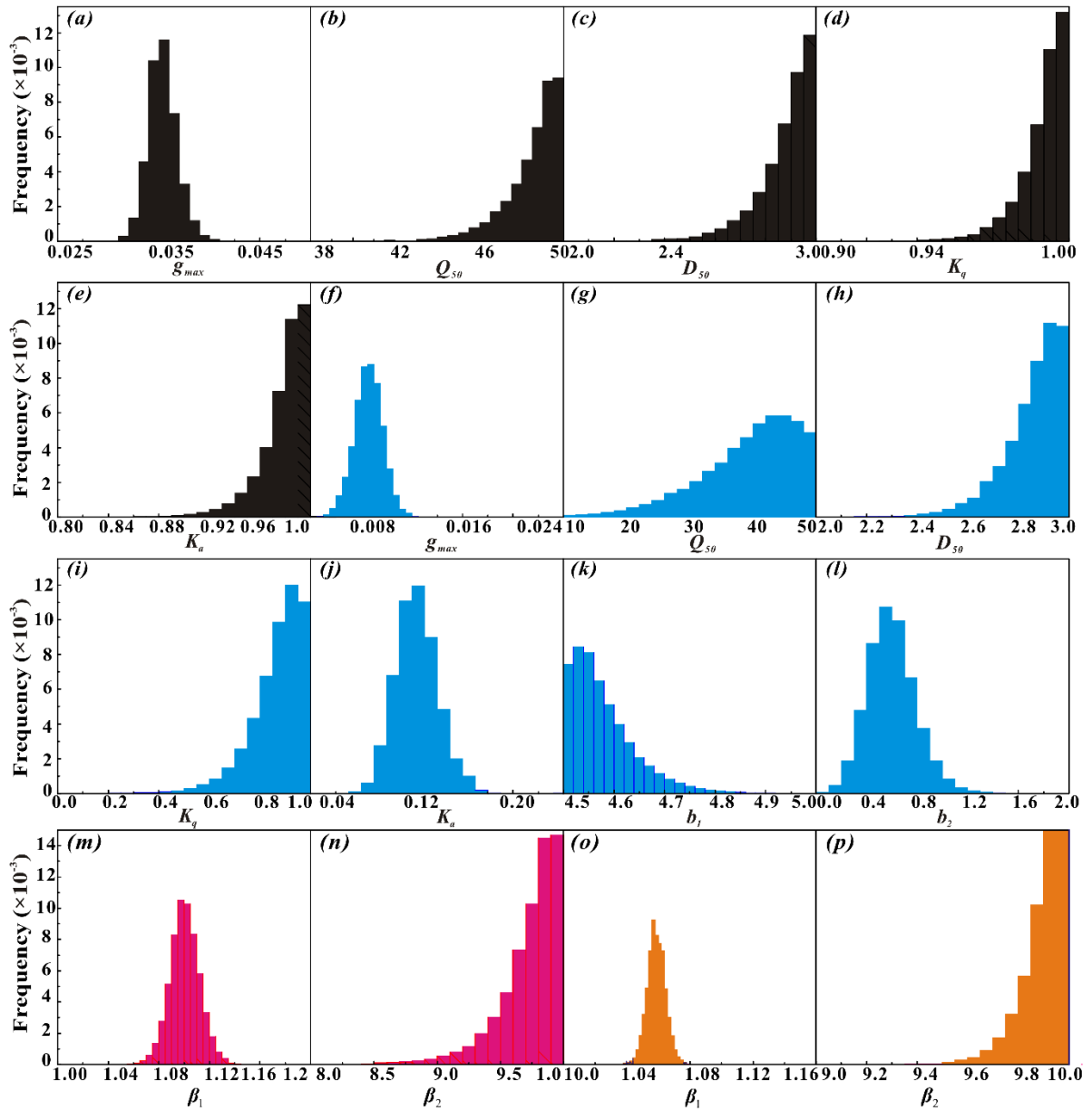


897

898 **Figure 1** Trace plots of the G-R statistic of Gelman and Rubin (Gelman and Rubin, 1992) using DREAM
 899 for the PM model (a) and (b) the SW model. Different parameters are coded with different colors. The
 900 dashed line denotes the default threshold used to diagnose convergence to a limiting distribution.

901

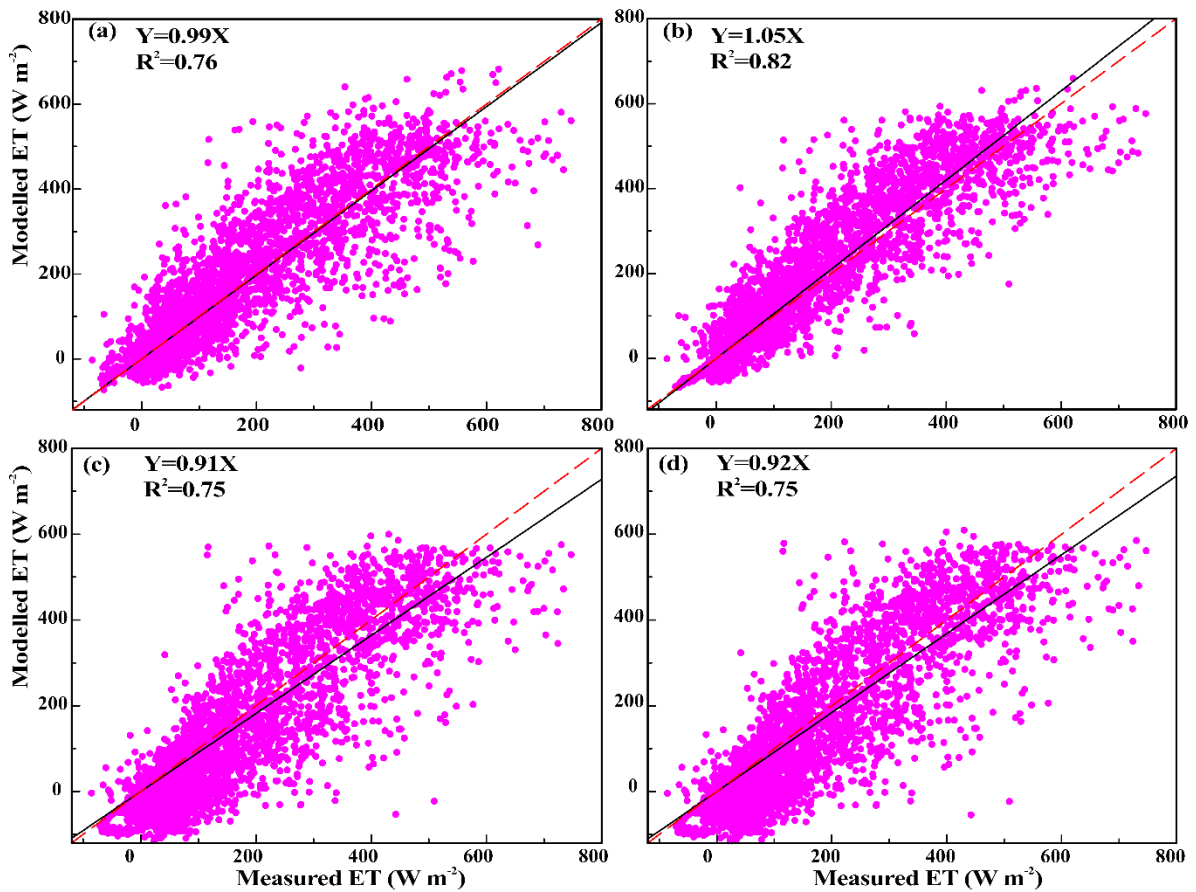
902



903

904 **Figure 2** (a)-(e), (f)-(l), (m)-(n), and (o)-(p) show histograms for the PM (black), SW (cyan), PT-FC
 905 (magenta) and AA (orange) models, respectively. These histograms are constructed from all chains for each
 906 model and a total of $40,000 \times N$ realizations are simulated using DREAM. The x axes represent the
 907 prespecified limits of the parameters.

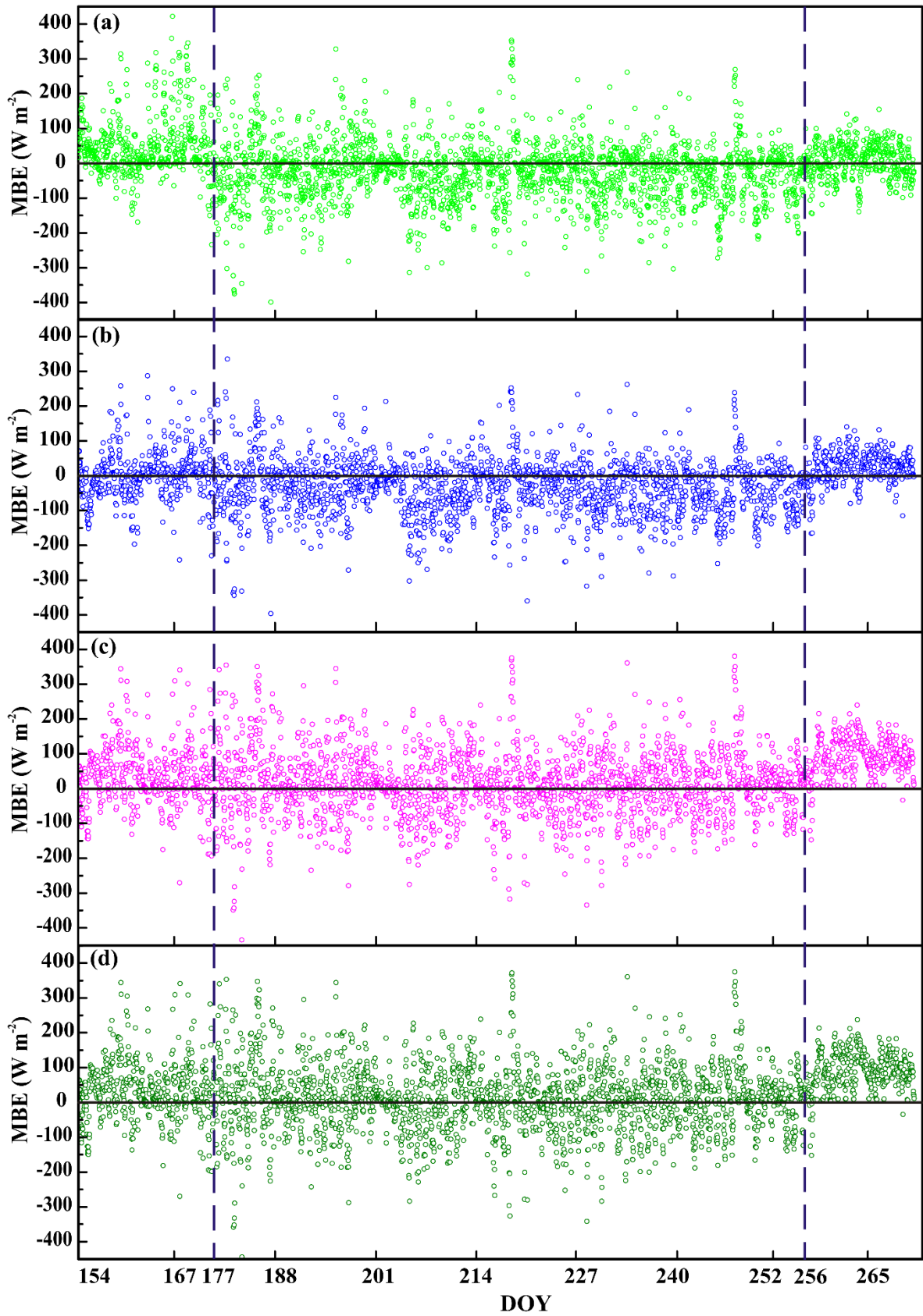
908



909

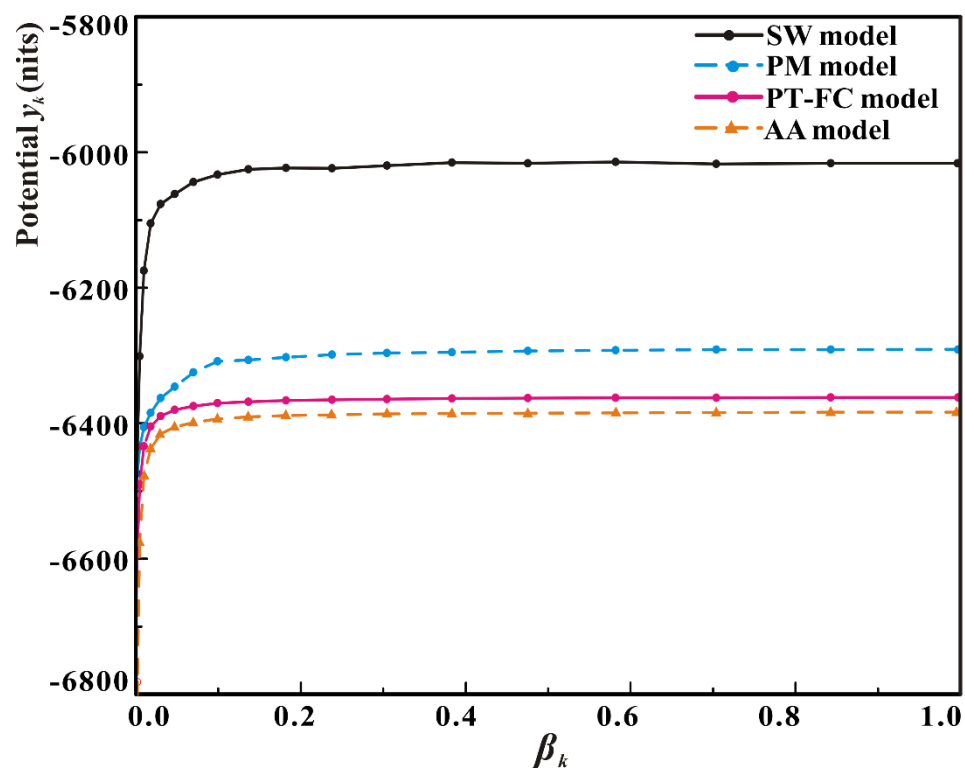
910 **Figure 3** Regressions between measured and modeled half-hourly ET values produced by different models
 911 from DOY 154 to DOY 270: (a) PM, (b) SW, (c) PT-FC and (d) AA. The regressions are: $Y = 0.99X$ ($R^2 =$
 912 0.76), $Y = 1.05X$ ($R^2 = 0.82$), $Y = 0.91X$ ($R^2 = 0.75$), and $Y = 0.92X$ ($R^2 = 0.75$) for the PM, SW, PT-FC
 913 and AA models, respectively.

914



915

916 **Figure 4** Mean bias error (MBE) of predicted and observed ET values for (a) PM, (b) SW, (c) PT-FC and
 917 (d) AA models from DOY 154 to DOY 270. Parameters used for prediction are estimated by DREAM with
 918 the dataset for the calibration period from DOY 154 to DOY 202.



920

921 **Figure 5** Variation of the mean posterior expectation of the potential y_k with β_k for the PM, SW, PT-FC and
 922 AA models.

923