

**Manuscript Number:** hess-2018-430

**Title:** Bayesian performance evaluation of evapotranspiration models based on eddy covariance systems in the arid region

**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 a positive comment on our revised manuscript. We are grateful to the editors and reviewers for their constructive suggestions on our manuscript (hess-2018-430).

Based on the reviewer's comment, we have briefed our manuscript title to "Bayesian performance evaluation of evapotranspiration models based on eddy covariance systems in the arid region". Meanwhile, we have corrected the technical issues through the manuscript. 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](mailto:xiaoyingzh@jlu.edu.cn)

**Responds to the comments:**

**Referee #1:**

**General comments:**

**1. Comment:** *The title of this manuscript should be more general.*

Response: Thanks for the comment, we have modified the manuscript title to “Bayesian performance evaluation of evapotranspiration models based on eddy covariance systems in the arid region”.

# Bayesian performance evaluation of evapotranspiration models: ~~a case study based on eddy covariance systems~~ in the arid region ~~of a maize field in northwestern China~~

Guoxiao Wei <sup>1,2</sup>, Xiaoying Zhang <sup>3,\*</sup>, Ming Ye <sup>4</sup>, Ning Yue <sup>1,2</sup>, Fei Kan <sup>1,2</sup>

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

<sup>2</sup> School of Earth and Environmental Sciences, Lanzhou University, China, 730000

<sup>3</sup> Construct Engineering College, Jilin University, China, 130400

<sup>4</sup> 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 ~~simulate~~estimate ET, there is no agreement on which model best describing the loss of water to the atmosphere. This study focuses on two aspects, evaluating the performance of four widely used ET models, and identifying parameters, as well the physical mechanisms that have significant impacts on the model performance. The four tested models are 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.~~ By incorporating the mathematically rigorous thermodynamic integration algorithm, the Bayesian model evidence (BME) approach is adopted to select the optimal model with half-hourly ET observations obtained at a spring maize field in an arid region. Our results reveal that SW has the best performance and the extinction coefficient is ~~the most significant factor in the ET models.~~ It is 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. This is one of the main reasons that the SW model outperforming the other models. Meanwhile, the good fitting of SW model to observations can counterbalance its higher complexity. In addition, the detailed analysis of the discrepancies between observations and model simulations during the crop growth season indicate that explicit treatment of energy imbalance and energy interaction will be the primary way to further improve ET model performance.

30 **Keywords:** Bayesian analysis; ET models; Eddy covariance; Model performance; Extinction coefficient

## 31 1. Introduction

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

47 ~~These ET models are generally complex, because of the coupling of the land surface and atmospheric~~  
48 ~~processes, and high-dimensional with a large number of parameters.~~ Comparing the performance of the  
49 competing ET models and evaluating and understanding the discrepancies between simulations of the models  
50 and corresponding observed surface-atmosphere water flux are remain challenging problems (Legates, 1999).  
51 Both non-Bayesian analysis (Szilagyi and Jozsa, 2008; Vinukollu et al., 2011; Li et al., 2013; Ershadi et al.,  
52 2015) and Bayesian analysis have been used to evaluate the performance of ET models (Zhu et al., 2014; Chen  
53 et al., 2015; Liu et al., 2016; Zhang et al., 2017; Elshall et al., 2018; Samani et al., 2018; Zeng et al., 2018) . Li  
54 et al. (2013) compared the ET simulations of the PM, SW and adjusted SW models under film-mulching  
55 conditions of maize growth in an arid region of China. They found that the half-hourly ET was overestimated by  
56 17% by the SW model. In contrast, the PM and adjusted SW models underestimated the daily ET by 6% and 2%,  
57 respectively. Therefore, the performances of PM and adjusted SW models are better than that of the SW model

58 in their case study. Ershadi et al. (2014) evaluated the surface energy balance system (SEBS), PM, PT-JPL (a  
59 modified Priestley–Taylor model) and AA models. Based on the average value of EF and RMSE, the model  
60 ranking from worst to best was AA, PM, SEBS, and PT-JPL. Ershadi et al. (2015) also compared the response  
61 of the models to different formulations of aerodynamic and surface resistances with global FLUXNET data.  
62 Their results showed considerable variability in model performance among and within biome types. Currently,  
63 ET model selection and comparison have been still conducted using traditional error metrics. It is known that  
64 error metrics are not adequate to provide a reasonable result of model ranking for disregarding model  
65 complexity (Marshall et al., 2005; Samani et al., 2018). The focus of this study is to use a Bayesian approach to  
66 evaluate the performance of the PM, SW, PT-FC, and AA models, which is a novelty contribution of this study.  
67 In ET models, the land surface energy system is governed by presumably infinite-dimensional physics. However,  
68 considering the ET models as finite-dimensional can be more precisely by covering all relevant relations.  
69 Therefore, employing consistent criteria for model selection might be justified when the aim is to better  
70 understand the processes involved (Höge et al., 2018). When using consistent model selection, Bayesian model  
71 evidence (BME), also known as marginal likelihood, measures the average fit of model simulations to their  
72 corresponding observations over a model’s prior parameter space. This feature enables BME to consider model  
73 complexity (in terms of the number of model parameters) for model performance evaluation. When comparing  
74 several alternative conceptual models, the model with the largest marginal likelihood is selected as the best  
75 model (Lartillot and Philippe, 2006). BME can thus be used for evaluating the model fit (over the parameter  
76 space) and for comparing alternative models. In previous studies, the Bayesian information criterion (BIC;  
77 Schwarz, 1978) and the Kashyap information criterion (KIC; Kashyap, 1982) have been used to approximate  
78 BME by using maximum likelihood theories to reduce the computational cost of evaluating BME (Ye et al.,  
79 2004). However, these approximations have theoretical and computational limitations (Ye et al., 2008; Xie et al.,  
80 2011; Schöniger et al., 2014), and a numerical evaluation (not a likelihood approximation) of BME is necessary,  
81 especially for complex models (Lartillot and Philippe, 2006). Lartillot and Philippe (2006) advocated the use of  
82 thermodynamic integration (TI) for estimating BME, also known as path sampling (Gelman and Meng, 1998;  
83 Neal, 2000), in order to avoid sampling solely in the prior or posterior parameter space. TI uses samples that are  
84 systematically generated from the prior to the posterior parameter space by conducting path sampling with  
85 several discrete power coefficient values (Liu et al., 2016). It is numerically accurate than the generally used  
86 harmonic mean method (Xie et al., 2011).

87 Most applications of Bayesian methods have focused on the calibration of individual models, while the  
88 comparison of alternative models continues to be performed using traditional error metrics. More generally,  
89 Bayesian approaches to model calibration, comparison, and analysis have been used far less used in the  
90 evaluation of ET models than in other areas of environmental science. In this study, the Bayesian approach is  
91 used to calibrate and evaluate the four ET models (PM, SW, PT-FC, and AA) based on an experiment over a  
92 spring maize field in an arid area of northwest China, from 3 June to 27 September 2014. The objectives of the  
93 study are as follows: (1) to calibrate ET model parameters using the Differential Evolution Adaptive Metropolis  
94 (DREAM) algorithm (Vrugt et al., 2008, 2009); (2) to identify which parameters had a greater impact on the  
95 model performance and to explain why the selected optimal model performed best; (3) to evaluate the  
96 performance of the models using traditional error metrics and BME; and (4) to analyze discrepancies between  
97 model simulations and observation data in order to better understand model performance and identify ways to  
98 improve these models. We expect that the study will not only boost the development of model parameterization  
99 and model selection but also contribute to the improvement of the ET models.

## 100 **2. Data and methodology**

### 101 **2.1. Description of the study area**

102 The experiment of maize growth was conducted at Daman Superstation, located in Zhangye City, Gansu  
103 province, northwest China. Daman Oasis is located in the middle Heihe River basin, which is the second largest  
104 inland river basin in the arid region of northwest China. The midstream area of the Heihe River basin is  
105 characterized by oases with irrigated agriculture, and is a region that consumes a large amount of water for both  
106 domestic and agricultural uses. The annual average precipitation and temperature are 125 mm and 7.2 °C  
107 (1960–2000), respectively. The annual accumulated temperature (>10 °C) is 3,234 °C, and the annual average  
108 potential evaporation is about 2,290 mm. The average annual duration of sunshine is 3,106 h with 148 frost-free  
109 days. The predominant soil type is silty-clay loam and the depth of the frozen layer is about 143 mm. The study  
110 area is a typical irrigated agricultural region, and the major source of water is snowmelt from the Qilian  
111 Mountains. Maize and spring wheat are the principal crops grown in the region. Maize is generally sown in late  
112 April and harvested in mid-September, and is planted with a row spacing of 40 cm and a plant spacing of 30 cm.  
113 The plant density is about 66,000 plants per hectare in the study area.

## 114 2.2. Measurements and data processing

115 Our data were collected from the field observation systems of the Heihe Watershed Allied Telemetry  
116 Experimental Research (HiWATER) project as described in Li et al (2013). The observation period was from  
117 DOY (day of the year) 154 to DOY 270 in 2014. An open-path eddy covariance (EC) system was installed in a  
118 maize field, with the sensors at a height of 4.5 m. Maize is the main crop in the study region, and thus covers  
119 sufficient planting area to set the EC measurements. The EC data was logged at a frequency of 10 Hz and then  
120 processed with an average time interval of 30 min. Sensible and latent heat fluxes were computed by the EC  
121 approach of Baldocchi (2003). Flux data measured by EC were controlled by traditional methods, including  
122 three-dimensional rotation (Aubinet et al., 2000), [Webb-Penman-Leuning WPL](#) (~~WPL Webb Penman Leuning~~)  
123 density fluctuation correction (Webb et al., 1980), frequency response correction (Xu et al., 2014), and spurious  
124 data removal caused by rainfall, water condensation, and system failure. About 85% of the energy balance  
125 closure was observed in the EC data (Liu et al., 2011).

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

## 132 2.3. Model description

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

### 135 2.3.1 Penman-Monteith (PM) model

136 The PM model can be formulated in the following way (Monteith, 1965):

$$137 \lambda E = \frac{\varepsilon A + (\rho C_p / \gamma) D_a g_a}{\varepsilon + 1 + g_a / g_s} \quad (1)$$

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

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

$$140 \quad g_s = g_s^c \left[ \frac{1 + \frac{\tau g_a}{(\varepsilon + 1) g_s^c} \left[ f - \frac{(\varepsilon + 1)(1 - f) g_s^c}{g_a} \right] + \frac{g_a}{\varepsilon g_i}}{1 - \tau \left[ f - \frac{(\varepsilon + 1)(1 - f) g_s^c}{g_a} \right] + \frac{g_a}{\varepsilon g_i}} \right] \quad (2)$$

141 where  $1 - \tau$  and  $\tau$  are the fractions of the total available energy absorbed by the canopy and by the soil, and  $\tau =$   
 142  $\exp(-K_a LAI)$ , and  $g_i$  and  $g_s^c$  are defined in equations (3) and (4), respectively (Monteith, 1965):

$$143 \quad g_i = \frac{A}{(\rho C_p / \gamma) D_a} \quad (3)$$

$$144 \quad g_s^c = \frac{g_{\max}}{K_q} \ln \left[ \frac{Q_h + Q_{50}}{Q_h \exp(-K_q LAI) + Q_{50}} \right] \left[ \frac{1}{1 + D_a / D_{50}} \right] f(\theta) \quad (4)$$

145 where  $f(\theta)$  represents water stress and is expressed as:

$$146 \quad f(\theta) = \begin{cases} 1 & \theta > \theta_a \\ \frac{\theta - \theta_b}{\theta_a - \theta_b} & \theta_b < \theta < \theta_a \\ 0 & \theta < \theta_b \end{cases} \quad (5)$$

147 and  $\theta_a$  is set as  $\theta_a = 0.75 \theta_b$ . Aerodynamic conductance  $g_a$  is calculated as:

$$148 \quad g_a = \frac{k^2 u_m}{\ln[(z_m - d)/z_{0m}] \ln[(z_m - d)/z_{0v}]} \quad (6)$$

149 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).

### 150 2.3.2. Shuttleworth-Wallace (SW) model

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



153  $\lambda ET = \lambda E + \lambda T = C_s ET_s + C_c ET_c$  (7)

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

155  $ET_c = \frac{\Delta A + \{\rho C_p (e_s - e_a) - \Delta r_a^c A_s\}}{\Delta + \gamma \left\{1 + r_s^c / (r_a^a + r_a^c)\right\}}$  (9)

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

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

158  $R_{ns} = R_n \exp(-K_a LAI)$  (10)

159 The coefficients  $C_s$  and  $C_c$  are obtained as follows:

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

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

162 where

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

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

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

166 Soil surface resistance is expressed as:

167  $r_s^s = \exp\left(b_1 - b_2 \frac{\theta}{\theta_s}\right)$  (16)

168 In this study, we consider the reciprocal of bulk stomatal resistance, known as canopy conductance. The

169 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

170 boundary layer resistance ( $r_a^c$ ) are modeled following the approach proposed by Shuttleworth and Gurney

171 (1990).

### 172 2.3.3. Priestley–Taylor and Flint-Childs (PT-FC) model

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

$$176 \quad \lambda ET = \alpha_{PT} \frac{\Delta}{\Delta + \gamma} (R_n - G) \quad (17)$$

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

$$179 \quad \lambda ET = \alpha \frac{\Delta}{\Delta + \gamma} (R_n - G) \quad (18)$$

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

$$184 \quad \alpha = \beta_1 [1 - \exp(-\beta_2 \Theta)] \quad (19)$$

$$185 \quad \text{where } \Theta = \frac{\theta - \theta_r}{\theta_s - \theta_r}.$$

### 186 2.3.4. Advection-aridity (AA) model

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

$$190 \quad \lambda ET = (2\alpha_{PT} - 1) \frac{\Delta}{\Delta + \gamma} (R_n - G) - \frac{\gamma}{\Delta + \gamma} \frac{\rho(q^* - q)}{r_a} \quad (20)$$

191 where  $\alpha_{PT}$  is the Priestley-Taylor coefficient, usually taken as 1.26 (Priestley and Taylor, 1972); and  $r_a$  is similar  
 192 to that used for the Penman-Monteith model (Brutsaert and Stricker, 1979; Brutsaert, 2005; Ershadi et al., 2014).  
 193 This model is based mainly on meteorological variables and does not require any information related to soil  
 194 moisture, canopy resistance or other measures of aridity (Ershadi et al., 2014). In this study, ~~as for the PT-FC~~  
 195 ~~model~~, we changed  $\alpha_{PT}$  to  $\alpha$ , which is calculated using the same equation as in the PT-FC model.

## 196 2.4 BME Estimation

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

$$198 \quad \text{BME} = p(\mathbf{D}|M) = \int p(\mathbf{D}|\boldsymbol{\theta}, M) p(\boldsymbol{\theta}|M) d\boldsymbol{\theta} \quad (21)$$

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

$$205 \quad q_\beta(\boldsymbol{\theta}) = p(\mathbf{D}|\boldsymbol{\theta}, M)^\beta p(\boldsymbol{\theta}|M) \quad (22)$$

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

$$209 \quad p_\beta(\boldsymbol{\theta}) = \frac{q_\beta(\boldsymbol{\theta})}{Z_\beta} \quad (23)$$

210 such that

$$211 \quad Z_\beta = \int q_\beta(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (24)$$

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

$$213 \quad U(\boldsymbol{\theta}) = \frac{\partial \ln q_{\beta}(\boldsymbol{\theta})}{\partial \beta} \quad (25)$$

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

$$215 \quad p(\mathbf{D}|M) = \frac{Z_1}{Z_0} = \exp\left\{\int_0^1 E_{\theta}[\ln p(\mathbf{D}|\boldsymbol{\theta}, M)] d\beta\right\} \quad (26)$$

216 The one-dimensional integral with respect to  $\beta$  is evaluated by using numerical methods by discretizing  $\beta$  into a  
 217 set of  $\beta_k$ . Since there is no theoretical method for selecting  $\beta_k$  values (Liu et al., 2016), we determined these  
 218 values using an empirical but straightforward method. Following Xie et al. (2011), a schedule of the power  
 219 posterior coefficients  $\beta_k$  is generated by

$$220 \quad \beta_k = (k / K)^{1/\varepsilon} \quad (27)$$

221 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  
 222 numerical inregration, equation (26) is evaluated via

$$223 \quad p(\mathbf{D}|M) = \exp\left(\int_0^1 y_{\beta} d\beta\right) = \exp\left(\sum_{k=0}^K r_{PI,k}\right) \quad (28)$$

224 such that

$$225 \quad r_{PI,k} = (\beta_k - \beta_{k-1}) \left[ \frac{y_k - y_{k-1}}{2} \right] \quad (29)$$

226 and

$$227 \quad y_k = E_{\beta}[\log p(\mathbf{D}|\boldsymbol{\theta}_k, M)] = \frac{1}{n} \sum_{i=1}^n \log p(\mathbf{D}|\boldsymbol{\theta}_{k,i}, M) \quad (30)$$

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

229 The random samples,  $\boldsymbol{\theta}_{k,i}$ , are drawn by using the MCMC method implemented in the DREAM code. See  
 230 Appendix B for further details on Bayesian inference and the DREAM algorithm. In the DREAM-based  
 231 calculation, the Metropolis acceptance ratio is  $\alpha_k = \min\left(1, [\alpha_{k, power-posterior} \alpha_{k, prior}]\right)$  with the power

232 posterior ratio given by  $\alpha_{k,power-posterior} = (\alpha_{k,prior})^{\beta_k}$ . The prior probability ratio  
 233  $\alpha_{k,prior} = \Pr(\theta_{k,new} | M) / \Pr(\theta_{k,old} | M)$  is the ratio of the probability of the newly proposed sample  
 234  $\theta_{k,new}$  and the probability of the previously accepted sample  $\theta_{k,old}$ . The posterior probability ratio  
 235  $\alpha_{k,prior} = L(\mathbf{D} | \theta_{k,new}, M) / L(\mathbf{D} | \theta_{k,old}, M)$  is the likelihood ratio of samples  $\theta_{k,new}$  and  $\theta_{k,old}$ , and  
 236  $\beta_k$  is the power posterior coefficient. Thus, to use the DREAM algorithm to sample any power posterior  
 237 distribution, the regular Metropolis acceptance ratio  $\alpha = \min(1, [\alpha_{posterior} \alpha_{prior}])$  is changed to  
 238  $\alpha_k = \min(1, [\alpha_{k,power-posterior} \alpha_{k,prior}])$  in DREAM.

## 239 2.5 Traditional statistical metrics of evaluating model performance

240 The traditional error metrics for evaluating model performance include—  $R^2$  and slope (correlation-based  
 241 measures), index of agreement (IA) and model efficiency (EF) (relative error measures), and the root mean  
 242 square error (RMSE) and mean bias error (MBE) (Poblete-Echeverria and Ortega-Farias, 2009). The definitions  
 243 of the listed metrics are:

$$244 \quad \text{IA} = 1 - \frac{\sum_{t=1}^n [O(t) - M(t)]^2}{\sum_{t=1}^n [ |O(t) - \overline{O(t)}| + |O(t) - \overline{M(t)}| ]^2} \quad (31)$$

$$245 \quad \text{EF} = 1 - \frac{\sum_{t=1}^n [O(t) - M(t)]^2}{\sum_{t=1}^n [O(t) - \overline{O(t)}]^2} \quad (32)$$

$$246 \quad \text{RMSE} = \sqrt{\frac{1}{n} \sum_{t=1}^n [O(t) - M(t)]^2} \quad (33)$$

$$247 \quad \text{MBE} = \frac{1}{n} \sum_{t=1}^n [O(t) - M(t)] \quad (34)$$

248 where  $O(t)$  is the observation and  $\overline{O(t)}$  is the mean observation at time  $t$ ;  $M(t)$  is the modeled value and  
 249  $\overline{M(t)}$  is the mean modeled value estimated by the posterior median parameter values; and  $n$  is the total number

250 of the observed values.

### 251 3. Results

#### 252 3.1 Parameter estimation

253 The PM model has five parameters  $g_{max}$ ,  $D_{50}$ ,  $Q_{50}$ ,  $K_q$  and  $K_a$ ; the SW model has seven parameters – the five  
254 used in the PM model and parameters  $b_1$  and  $b_2$ . The PT-FC and AA models each include two parameters,  
255 denoted by  $\beta_1$  and  $\beta_2$  (Table 1). The prior probability density of each parameter is specified as an uniform  
256 distribution with the ranges listed in Table 1. A total of 50,000 realizations were generated with the DREAM  
257 algorithm, which was used to estimate the posterior probability density function of each parameter with the  
258 calibration period data from DOY 154 to DOY 202. In the calculations, the chain number,  $N$ , was equal to the  
259 number of parameters in the associated model. Therefore,  $N$  is equal to 5, 7, 2 and 2 for the PM, SW, PT-FC and  
260 AA models, respectively. For each model, the first 10,000 samples were discarded as burn-in data, and the  
261 remaining 40,000 samples were used for calibration. In total,  $40,000 \times N$  realizations were used to set up  
262 posterior density functions for each model. To illustrate the efficiency and convergence of DREAM for the ET  
263 models, Figure 1 shows the trace plots of the G-R statistic for each of the different parameters in the PM and  
264 SW models using a different color. The algorithm required about 8,000 generations to make the G-R statistic  
265 close to 1.0 for the two models. The acceptance rates for the PM and SW models were about 15.3% and 18.9%,  
266 respectively.

267 Histograms of the DREAM-derived marginal distributions of the parameters are presented in Figure 2 and  
268 summarized in Table 2 by Maximum Likelihood Estimates (MLEs), posterior medians and 95% probability  
269 intervals. Figures 2a-2e, 2f-2l, 2m-2n, and 2o-2p show histograms of the PM, SW, PT-FC and AA models,  
270 respectively. Parameter  $g_{max}$  (Fig. 2a) in the PM model, parameters  $g_{max}$ ,  $K_a$ ,  $b_2$  (Fig. 2f, 2j, 2l) in the SW model,  
271 and parameter  $\beta_1$  (Fig. 2m) in the PT-FC model and AA model (Fig. 2o) were well constrained and occupied a  
272 relatively small range. These parameters displayed a unimodal distribution and appeared approximately  
273 Gaussian. In contrast, the distributions of the other parameters differed significantly from a Gaussian  
274 distribution, as shown by the corresponding histograms. The distributions of all but one of these parameters  
275 concentrated most of the probability mass at their upper limits. The exception was parameter  $b_1$  for the SW  
276 model (Fig. 2k), which clearly does not follow a normal distribution with most of the mass concentrated in the

277 lower bounds. In contrast,  $Q_{50}$  was not only poorly constrained (Fig. 2g) but was also the upper edge-hitting  
278 parameter in the SW model. Moreover, the corresponding distributions of the same parameter in different  
279 models were slightly different. For example, the mean of  $g_{max}$  in the PM model ( $0.04 \text{ mm s}^{-1}$ ) was less than that  
280 in the SW model ( $0.01 \text{ mm s}^{-1}$ ) (Fig. 2a and 2f, Table 2), except that  $D_{50}$  in the PM and SW models and  $\beta_2$  in the  
281 PT-FC and AA models exhibited similar regions. It is interesting to observe that the distribution of  $K_a$  in PM  
282 model (Figure 2e) has a truncated distribution with highest probability mass at the upper bound, whereas the  
283 distribution of  $K_a$  in the SW model (Figure 2j) tends to become approximately normal. Overall, the marginal  
284 posterior probability density function of most of the individual parameters occupied only a relatively small  
285 region compared with the uniform prior distributions, and exhibited relatively large uncertainty reduction.

### 286 3.2 Performance of the models

287 The performance of each of the four ET models was evaluated over the course of the whole season in 2014.  
288 The calibrated parameters of the four models were used and individual ET models were run to estimate the  
289 half-hourly  $\lambda\text{ET}$  values. Table 3 summarizes the statistical results for the performance of the models using  
290 regression line slope,  $R^2$ , RMSE, MBE, IA, and EF. The regressions between measured and modeled  $\lambda\text{ET}$  values  
291 and MBE are shown in Figures 3 and 4, respectively.

292 In general, the four models produced slightly better fits to the measured  $\lambda\text{ET}$  for all the seasons with  $R^2$   
293 larger than 0.75 (Fig. 3). However, obvious discrepancies in the predictions made by the models were detected  
294 by comparing measured and modeled  $\lambda\text{ET}$ . According to the regression line slope and MBE, the PM model  
295 overestimated ET by 1% with a MBE of  $-9.52 \text{ W m}^{-2}$ , and the SW model overestimated ET by 5% with a  
296 relatively higher MBE of  $-19.07 \text{ W m}^{-2}$  compared to the PM model. The PT-FC and AA models tended to  
297 underestimate  $\lambda\text{ET}$  by 9% and 8% with an MBE of 25.42 and 23.29  $\text{W m}^{-2}$ , respectively. From a comparison  
298 between the slope and MBE, the PM model performance was higher than that of the other three models, with a  
299 slope almost equal to 1 and relatively lower MBE. The SW model was ranked second, while the performance of  
300 the AA model was slightly higher ~~comparable to~~ than that of the PT-FC model. However, if  $R^2$ , RMSE, IA, and  
301 EF were used to evaluate performance, the SW model had the best overall performance with  $R^2=0.83$ , RMSE=  
302  $76.34 \text{ W m}^{-2}$ , IA = 0.95 and EF = 0.79. The second-best model was the PM model, and the PT-FC was ranked  
303 third, while the AA model ranked fourth. Based on the analysis of these traditional error metrics, the PT-FC and  
304 AA models yielded similar results. The observed and modeled  $\lambda\text{ET}$  for the four ET models were tightly grouped

305 along the regression lines (Figure 3), and the PT-FC and AA models had similar<sup>ly</sup> modeled ET values with a  
306 similar degree of point scattering along the regression lines (Figure 3c-3d).

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

### 314 3.3 Comparison of the models using BME

315 Since there is currently no theoretical method for selecting power posterior  $\beta$  values, we determined these  
316 values using empirical but straightforward methods. For any power coefficient of  $\beta \in [0, 1]$ , a sample was drawn  
317 from the distribution  $p_\beta$  (Eq. 2523) through running DREAM. Although adding more  $\beta_k$  values might improve  
318 the BME estimation, this was not done because of the computational cost. For each  $\beta_k$  value, at least 150,000  
319 DREAM simulations were large enough to ensure convergence. Figure 5 shows the evolution of  $\ln p(D|\theta, M)$   
320 for the four models as a function of  $\beta$  for a dataset covering the entire period. The BME for the SW model was  
321 substantially larger than that for the other three models, and the BME for the AA model was the smallest. The  
322 BME-based model ranking (from the best to the worst) is SW, PM, PT-FC, and AA. The PT-FC and AA models,  
323 which consist<sup>ing</sup> of the same number of parameters, had similar potential patterns of evolution with respect to  
324 the coefficient  $\beta_k$ . The results illustrate that with the addition of parameters, the model complexity and the model  
325 performance are both increased.

## 326 4. Discussion

### 327 4.1 Parameter uncertainty analysis

328 With regard to the efficiency of the DREAM algorithm, the acceptance rates of the PM (15.3%) and SW  
329 (18.9%) models were much higher than those obtained by some Markov Chain Monte Carlo (MCMC) algorithms  
330 that have been used in previous studies— (Sadegh et al., 2014). The posterior parameter bounds exhibit a larger



331 reduction using the DREAM algorithm compared with other studies using the Metropolis–Hasting algorithm.  
332 This demonstrates that DREAM could efficiently handle problems involving high-dimensionality, multimodality  
333 and, nonlinearity.

334 The results showed that the assumed prior uncertainty ranges from most parameters in the four models  
335 were significantly reduced. This indicates that the observed ET data contained sufficient information to estimate  
336 these parameters. Surface conductance  $g_s$  and modeled ET in the PM model are relatively insensitive to  $Q_{50}$ ,  $D_{50}$   
337 and  $K_q$ . Hence, these parameters could not be well constrained, and further relaxing the ranges for these  
338 parameters could not result in physically realistic behavior of the model. The calculation of  $-g_s^c$  in the SW model  
339 is the same as in the PM model, and thus,  $g_s^c$  and modeled ET in the SW model are also insensitive to  
340 parameters of  $Q_{50}$ ,  $D_{50}$ ,  $K_q$ . Therefore, these three parameters were also not well constrained in the SW model. In  
341 addition, the uncertainties present in the edge-hitting parameters, may be the outcome of model biases or  
342 EC-measured ET data errors, or the characteristic time scale of parameters governing the processes affecting ET  
343 is not exactly on the order of half-hours (Braswell et al., 2005). For example,  $Q_{50}$  and  $D_{50}$  govern changes in  
344 visible radiation flux and the humidity deficit at which stomatal conductance is half its maximum value,  
345 respectively, and these parameters may change over a shorter or longer time scale than half-hours.

346 The ecophysiological parameter  $g_{max}$  is a variable in the  $g_s^c$  equation in both the PM and SW models, but  
347 this parameter is sensitive to  $-g_s^c$  and has a significant impact on the evaluated ET. Its effects is relatively  
348 independent compared to the other meteorological parameters in the models, and therefore this parameter was  
349 well specified in the PM and SW models. The posterior mean value of  $g_{max}$  ( $0.04 \text{ m s}^{-1}$ ) in the PM model from  
350 our study was close to that ( $0.05 \text{ m s}^{-1}$ ) reported in northwestern China (Li et al., 2013; Zhu et al., 2014), but  
351  $g_{max}$  ( $0.01 \text{ m s}^{-1}$ ) in the SW model was less than the reported value. Parameter  $\beta_l$  was well constrained in the  
352 PT-FC and AA models because it was relatively independent and did not directly relate to other observed  
353 variables.

354 Parameter  $K_a$  implicitly appears in the surface conductance equation (Eq.2) in PM model and  $K_a$  is  
355 insensitive to  $g_s$  and modeled ET (Leuning et al., 2008). In contrast,  $K_a$  is contained in the equation of net  
356 radiation flux into the substrate (Eq.10) in the SW model. This parameter can explicitly partition the total  
357 available energy into that absorbed by the canopy and by the soil in the SW model. An analysis of equation (10),

358 found that the variation of  $K_a$  could not only account for the extinction effect but also correct the energy forcing  
359 data errors. This also meant that the estimated value of  $K_a$  using calibration data was actually not just the true  
360 extinction coefficient, but also included the energy imbalance correction in the SW model. From this analysis,  
361 we could see that  $K_a$  not only involved the distribution of energy between the canopy and the soil surface but  
362 also the energy imbalance. Therefore, parameter  $K_a$  has a great influence on the performance of the SW model.  
363 This is why  $K_a$  is poorly constrained in the PM model but well constrained in the SW model. To further illustrate  
364 the insights regarding the influence of parameter  $K_a$  on the performance of the SW model, we calibrated the SW  
365 model again and reran the model with a constant value of  $K_a$ . The results showed a significant reduction in  
366 model performance when  $K_a$  was held constant. This implied that the main reason for the SW model  
367 outperforming the PM model in our study was not only the more physically rigorous structure of the SW model  
368 but also the key parameter  $K_a$  being well constrained in the SW model.

369 In general, parameters related to soil surface resistance in the SW model were well evaluated, while  
370 parameters related to canopy surface resistance in PM and SW models were poorly estimated. Therefore, using a  
371 reliable canopy surface resistance equation in the ET model was crucial for improving its performance. In  
372 addition, in our study, the traditional approach was used to quantify the uncertainty,  $\sigma_{\text{total}}$ —which assumed that the  
373 uncertainty mainly arose because of the parameter uncertainty. However, this method cannot explicitly consider  
374 errors in the input data and model structural inadequacies. This is unrealistic for real applications, and it is  
375 desirable to develop a more reliable inference method to treat all sources of uncertainty separately and  
376 appropriately (Vrugt et al., 2008). Moreover, simultaneous direct measurement by micro-lysimeter of sap flow  
377 and daily soil evaporation will further help to constrain the model parameters.

## 378 **4.2 Evaluation and selection of the models**

379 In this study, the traditional statistical measures and BME were chosen to evaluate and compare the  
380 performance of four ET models. From the respective composition of these measures, the statistical measures can  
381 be divided into residual-based metrics (such as regression slope and MBE) and squared-residual-based measures  
382 (such as  $R^2$ , RMSE, IA, and EF). The rankings of the models obtained using the same type of metric  
383 (residual-based or squared-residual-based) are similar. Slope and MBE, for example, which are both  
384 residual-based measures, produce identical rankings. However, the rankings produced by metrics of different  
385 types are not the same. For example, the PM model outperforms the SW model according to the residual-based

386 metrics, but the performance of the PM model is worse than the SW model based on the squared-residual-based  
387 measures. The comparative analysis shows ~~a~~ consistency between BME and the squared-residual-based metrics  
388 (hence the residual-based metrics disagreed with the BME measures). This reveals that the more complex SW  
389 model is the best model based on BME and squared-residual-based statistics. The rank order of overall  
390 performance of the models from best to worst is: SW, PM, PT-FC, and AA model.

391 Previous studies had shown that BME evaluated by TI provided estimates similar to the true values, and  
392 selected the true model if the true model was included within the candidate models (Marshall et al., 2005;  
393 Lartillot and Philippe, 2006). Meanwhile, some have argued that Bayesian analysis would choose the simplest  
394 model (Jefferys and Berger, 1992; Xie et al., 2011) because of the best trade-off between good fit with the data  
395 and model complexity (Schöniger et al., 2014). In this case, the most complex SW model had the highest BME  
396 and was chosen as the model with the best performance. This probably resulted from the fact that the complex  
397 SW model is indeed the most reliable model among the alternative ET models and can provide a good fit to  
398 justify its higher complexity. The SW model is a two-layer model, and simulates soil evaporation and plant  
399 transpiration separately, whereas the PM model is a single-layer model in which the plant transpiration and soil  
400 evaporation cannot be separated (Monteith, 1965). The PT-FC model is a simplified version of the PM model,  
401 and only requires meteorological and radiation information (Priestley and Taylor, 1972), whereas the AA model  
402 only relies on the feedback between actual ET and potential ET (Brutsaert and Stricker, 1979).

403 The results indicate that the squared-residual-based measures yielded the same rank order as the BME  
404 consistently, which makes the squared-residual-based metrics seemed to identify a reasonable rank order.

405 However, this has not been the general case, since the error metrics and BME belong to different types of model  
406 selection and there are differences in the behavior and optimality of the two types of model selection. BME is a  
407 consistent model selection which tries to identify which of the models produced the observed data. Conversely,  
408 nonconsistent model selection uses the available data to estimate which of the models might be best in  
409 predicting future data. In fact, the error metrics are essentially nonparsimonious model selection, which is a  
410 special case of nonconsistent model selection. The simple traditional statistical measures were known to usually  
411 provide a biased view of the efficacy of a model (Kessler and Neas, 1994; Legates and McCabe, 1999), where  
412 only the goodness of fit is used for rating models without penalizing the model complexity and thus lacking  
413 consistency for the selected model (Höge et al., 2018). In addition, sensitivity to outliers is associated with these

414 metrics and leads to relatively high values due to the squaring of the residual terms (Willmott, 1981).  
415 Furthermore, these traditional statistical metrics ignore the priors, ~~without penalizing model complexity~~, which  
416 is in fact used in a Bayesian analysis. PT-FC and AA, provide identical estimates of  $R^2$  and IA. This is most  
417 likely because both models had the same dimension and a similar model structure. Marshall et al. (2005) argued  
418 that EF would provide an incorrect conclusion, and Samani et al. (2018) suggested that RMSE would select the  
419 complex model as the best performing model. As for the slope and MBE, the rankings produced by these  
420 residual-based metrics were in obvious disagreement with the one based on BME. Part of the lower simulation  
421 values ~~of slope and MBE may could~~ be counter-balanced by the higher values ~~of that~~ of in the slope and MBE  
422 methods, thus these criteria provide an erroneous and unreliable evaluation of the models. Therefore, the  
423 squared-residual-based and residual-based measures were not certain to provide reasonable results in terms of  
424 model ranking. The consistency between BME and the squared-residual-based metrics only indicates that the  
425 optimal model evaluated by BME would also provide the best predictions, and thus consistent model selection  
426 should also be asymptotically efficient (Leeb & Pötscher, 2009; Shao, 1997).

#### 427 **4.3 Analysis of model-data mismatch**

428 Conceptual and structural inadequacies of the hydrological model together with measurement errors of the  
429 model input (forcing) and output (calibration) data introduce errors in the estimated parameters and model  
430 simulations (Laloy, 2015). Hydrological systems are indeed heavily input-driven and errors in forcing data can  
431 dramatically impair the quality of calibration results and model output (Bardossy and Das, 2008; Giudice, 2015).  
432 Measurement errors occur for a variety of reasons, including unreasonable gap-filling in rainy days; dew and fog;  
433 inadequate areal coverage of point-scale soil water measurement; mechanical limitations of the EC system; and  
434 inaccurate measurements of wind-speed, soil water, radiation and vapor pressure deficit. ET processes is  
435 described using equations that can only capture parts of the complex natural processes and any ET model is an  
436 inherent simplification of the real system. These inadequacies can thus lead to biased parameters and  
437 implausible predictions.

438 In our study, the results indicated that the PM and SW models overestimated the half-hourly ET compared  
439 to the measured ET. Several studies also indicated that ET was overestimated by the PM model (Fisher et al.,  
440 2005; Ortega-Farias et al., 2006; Li et al., 2015) and the SW model (Li et al., 2013; Li et al., 2015; Zhang et al.,  
441 2008). Possible reasons for the inaccurate estimates included the following: (1) Anisotropic turbulence with

442 weak vertical and strong horizontal fluctuation leads to energy imbalance. The total turbulent heat flux was  
443 lower by ~10–30% compared to the available energy in many land surface experiments (Tsvang et al., 1991;  
444 Beyrich et al., 2002; Oncley et al., 2007; Foken et al., 2010) and influx networks (Franssen et al., 2010). Liang  
445 et al. (2017) also showed an energy imbalance result in the semiarid area in China, and indicated that the energy  
446 balance closure ratio ranged from 0.52 to 0.90 during the day, whereas it was about 0.25 at night. However, the  
447 measured ET only included vertical flux and not horizontal flux, leading to the measured ET being lower than  
448 that of ET predicted by the PM and SW models using the available energy. (2) The absence of a mechanistic  
449 representation of the physiological response to plant hydrodynamics makes it difficult for the available ET  
450 models to resolve the dynamics of intradaily hysteresis, producing patterns of diurnal error, while the imbalance  
451 or lack of between-leaf water demand and soil water supply imposes hydrodynamic limitations on stomatal  
452 conductance (Thomsen et al., 2013; Zhang et al., 2014; Matheny et al., 2014). Li et al. (2015) also concluded  
453 that neglecting the restrictive effect of the soil on water transport in empirical canopy resistance equations can  
454 result in large errors in the partial canopy stage. However, these equations can estimate ET accurately under the  
455 full canopy stage (Alves and Pereira, 2000; Katerji and Rana, 2006; Katerji et al., 2011; Rana et al., 2011). Li et  
456 al. (2015) showed that the PM model combined with canopy resistance overestimated maize ET during the  
457 partial and dense canopy stages by 16% and 13%, respectively. Moreover, in a study of ET in vineyards,  
458 Leuning (2008) found that the PM model coupled with canopy resistance overestimated ET during the entire  
459 growth stage by 29%.

460 The estimates for ET produced by the PT-FC and AA models were generally lower than the measured  
461 values during the entire season. In addition, the four models also underestimated ET during periods of partial  
462 cover ( $LAI < 2.5 \text{ m}^2 \text{ m}^{-2}$ ). The PT-FC and AA models consistently underestimated ET, especially during the late  
463 maturity stage. The underestimation probably resulted from the following: (1) Non classical situations, such as  
464 the oasis effect, may occur in the study area. Strong evaporation from the moist ground and plants results in  
465 latent heat cooling. However, this upward latent heat flux was opposed by a downward sensible heat flux from  
466 the warm air to the cool ground, and thus the latent heat flux was positive while the sensible heat flux is  
467 negative. Therefore, the latent heat flux can be greater in magnitude than the solar heating, because of the  
468 additional energy extracted from the warm air by evaporation (Stull, 1988). (2) The lack of mechanistic  
469 representation of rainfall interception in ET models probably led to inaccurate simulation for periods soon after  
470 rainy days. Bohn and Vivoni (2016) found that evaporation of canopy interception accounted for 8% of the

471 annual ET across the North American monsoon region. Comparing the AA and PT-FC models, the former  
472 includes forcing data of available radiation, soil water content and relative humidity, but the PT-FC model only  
473 requires available radiation and soil water content and is independent of relative humidity. However, the similar  
474 statistical results and similar degrees of MBE scatter indicate that relative humidity has little influence on the  
475 AA model simulation. The consistent and consecutive underestimation of ET by the PT-FC and AA models  
476 during the late maturity stage show that the model-data disagreement is not caused by regional advection and  
477 rainfall interception, because atmospheric processes and thermally-induced circulation can only occur at certain  
478 times and during certain days. Therefore, we think that the consistent underestimation of ET by the PT-FC and  
479 AA models results primarily from conceptual and structural inadequacies, energy imbalance, and soil water  
480 stress. Although the PM and SW models share a common theoretical basis and the PT-FC model is a  
481 simplification of the PM model, these models perform significantly differently. Part of the overestimation of ET  
482 by the PM and SW models, caused by coupling with the canopy resistance, may be offset by underestimation  
483 caused by energy imbalance and soil water stress. However, underestimation of ET by the PT-FC and AA  
484 models cannot be counterbalanced by overestimation during the later maturity stage because the PT-FC and AA  
485 models are independent of the canopy resistance. Consequently, the half-hourly patterns of errors in the  
486 estimates of ET by the PM and SW models are characterized by symmetry and a low degree of scatter, but the  
487 PT-FC and AA models exhibit consistently asymmetrical error patterns. By contrast, other studies showed that  
488 the PM model (Kato et al., 2004) and the SW model (Chen et al., 2015) underestimated half-hourly ET. As for  
489 the PT-FC and AA models, some studies reported that the PT-JPL (Zhang et al., 2017) and the AA model showed  
490 an overall poor performance (Zhang et al., 2017). While other studies have indicated that the AA method  
491 performed well for both maize and canola crops (Liu et al., 2012). Therefore, the performance of the four ET  
492 models appears to vary not only for different crops and locations but also for different meteorological,  
493 physiological and soil conditions. Moreover, the performance is also related to the stage of crop growth. Note  
494 that these conclusions about the ET models evaluation are derived from traditional error metrics rather than  
495 those based on BME model selection. It would be desirable to use available data from other study areas or from  
496 other crops for BME-based model selection to confirm whether the SW model is the optimal model under other  
497 conditions. Overall, combined with the parameter uncertainty analysis described in Section 4.1, we conclude that  
498 energy imbalance and energy interaction between canopy and soil surface have a greater impact on the model  
499 performance. And thus, explicitly treating of energy error, and incorporating the elements of existing hydrologic

500 theory about energy interaction between canopy and surface or conceptually correcting the energy interaction  
501 are a practicable option for model improvement and application.

## 502 **5. Conclusions**

503 This study illustrated the application of the Bayesian approach on the statistical analysis and model  
504 selection of four widely used ET models. The results showed that the DREAM algorithm successfully reduced  
505 the assumed prior uncertainties for most of the parameters in the four models. In the model calibration, the key  
506 parameters which had a significant influence on ET simulations were well constrained. The main reasons for the  
507 outperforming of SW model were its physically rigorous structure and the extinction coefficient parameter,  
508 which is sensitive and has a significant impact on the performance of the model, being well constrained. BME is  
509 a consistent model selection to identify the best fitting to the observed data. Although the squared-residual-based  
510 metrics, including  $R^2$ , IA, RMSE, and EF, produced a ranking identical to that of BME, it must be noted that  
511 these squared-residual-based metrics do not allow using prior information and do not penalize the model  
512 complexity when comparing the models. Therefore, some cautions are needed when using these statistical  
513 methods to compare different models.

514 The model–data discrepancies were analyzed to facilitate model improvement after Bayesian model  
515 calibration and comparison. The results indicate that the discrepancies arose mainly as a result of energy  
516 imbalance caused by anisotropic turbulence, additional energy induced by advection processes, the absence of a  
517 mechanistic representation of the physiological response to plant hydrodynamics and the energy interaction  
518 between canopy and surface. Among these causes, energy imbalance and additional energy are related to forcing  
519 data errors rather than to an unreasonable model structure. Thus, understanding the process of the physiological  
520 response to plant hydrodynamics and the interaction between canopy and surface is essential for improving the  
521 performance of evapotranspiration models. Overall, the applications of Bayesian calibration, Bayesian model  
522 evaluation and analysis of model–data discrepancies in our study, provide a promising framework for reducing  
523 uncertainty and improving the performance of ET models. It would be desirable to confirm whether the SW is  
524 the optimal model using data of other crops\_ or other climate regions.

## 525 **Author contribution**

526 Guoxiao Wei and Xiaoying Zhang designed the experiments. Ning Yue and Fei Kan carried them out.  
527 Ming Ye developed the model selection scheme. Guoxiao Wei performed the simulations. Guoxiao Wei and

528 Xiaoying Zhang prepared the manuscript with contributions from all co-authors.

## 529 **Competing interests**

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

## 531 **Acknowledgments**

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

## 538 **References**

- 539 Allen, R. G., Perista, L. S., Raes, D., and Smith, M.: Crop Evapotranspiration-Guidelines for Computing Crop  
540 Water Requirements; FAO Irrigation and Drainage apers-56, FAO-Food and Agriculture Organization  
541 of the United Nations, Rome, 1998.
- 542 Alves, I. and Pereira, L. S.: Modeling surface resistance from climatic variables? *Agric. Water Manage.*, 42,  
543 371–385, 2000.
- 544 Aubinet, M., Grelle, A., Ibrom, A., Rannik, Ü., Moncrieff, J., and Foken, T.: Estimates of the annual net carbon  
545 and water exchange of forests: the euroflux methodology, *Adv. Ecol. Res.*, 30(1), 113-175, 2000.
- 546 Baldocchi, D. D.: Assessing the eddy covariance technique for evaluating carbon dioxide exchange rates of  
547 ecosystems: past, present and future, *Glob. Change. Biol.*, 9, 479–492, 2003.
- 548 Bardossy, A., and Das, T.: Influence of rainfall observation network on model calibration and application,  
549 *Hydrol. Earth Syst. Sci.*, 12(1), 77-89, doi:10.5194/hess-12-77-2008, 2008.
- 550 Barton. I. J.: A Parameterization of the Evaporation from Nonsaturated Surfaces, *J. Appl Meteorol.*, 18(1), 43-47,  
551 1979.
- 552 Beyrich, F., Richter, S. H., Weisensee, U., Kohsiek, W., Lohse, H., de Bruin, H. A. R., Foken, T., Göckede, M.,  
553 Berger, F., Vogt, R., and Batchvarova, E.: Experimental determination of turbulent fluxes over the  
554 heterogeneous litfass area: selected results from the litfass-98 experiment, *Theor. Appl. Climatol.*,  
555 73(1–2), 19–34, doi:10.1007/s00704-002-0691-7, 2002.
- 556 Bohn, T. J., and Vivoni, E. R.: Process-based characterization of evapotranspiration sources over the North  
557 American monsoon region, *Water Resour. Res.*, 52, 358–384, doi:10.1002/2015WR017934, 2016.
- 558 Bonan, G.: *Ecological climatology: concepts and applications*, Cambridge University Press, 2008.
- 559 Braswell, B. H., Sacks, W. J., Linder, E., and Schimel, D. S.: Estimating diurnal to annual ecosystem parameters  
560 by synthesis of a carbon flux model with eddy covariance net ecosystem exchange observations, *Global*  
561 *Change. Biol.*, 11, 335–355, 2005,
- 562 Brutsaert, W., and Han, S.: An advection-aridity approach to estimate actual regional evapotranspiration, *Water*  
563 *Resour. Res.*, 15(2), 443-450, 1979.
- 564 Brutsaert, W.: *Hydrology: An Introduction*. Cambridge University Press, Cambridge, 2005.
- 565 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  
566 Bayesian analysis to compare the performance of three evapotranspiration models for rainfed jujube  
567 (*Ziziphus jujuba* Mill.) plantations in the Loess Plateau, *Agr. Water. Manage.*, 159, 341–357, 2015.
- 568 Elshall, A. S., Ye, M., Pei, Y., Zhang, F., Niu, G. Y., and Barron-Gafford, G. A.: Relative model score: A scoring



569 rule for evaluating ensemble simulations with application to microbial soil respiration modeling, *Stoch.*  
570 *Env. Res. A.*, 1-11, DOI: 10.1007/s00477-018-1592-3, 2018.

571 Ershadi, A., McCabe, M. F., Evans, J. P., Chaney, N. W., and Wood, E. F.: Multi-site evaluation of terrestrial  
572 evaporation models using fluxnet data, *Agric. For. Meteorol.*, 187(8), 46-61, 2014.

573 Ershadi, A., McCabe, M. F., Evans, J. P., and Wood, E. F.: Impact of model structure and parameterization on  
574 Penman–Monteith type evaporation models, *J. Hydrol.*, 525, 521–535, 2015.

575 Fisher, J. B., DeBiase, T. A., Qi, Y., Xu, M., and Goldstein, A. H.: Evapotranspiration models compared on a  
576 Sierra Nevada forest ecosystem, *Environ. Model. Softw.*, 20 (6), 783–796, 2005.

577 Flint A. L., Childs, S. W.: Use of the Priestley–Taylor evaporation equation for soil water limited conditions in a  
578 small forest clearcut, *Agric. For. Meteorol.*, 56(3–4), 247–260, 1991.

579 Foken, T., Mauder, M., Liebethal, C., Wimmer, F., Beyrich, F., Leps, J. P., Raasch, S., DeBruin, H. A. R.,  
580 Meijninger, W. M. L., and Bange, J.: Energy balance closure for the LITFASS-2003 experiment, *Theor.*  
581 *Appl. Climatol.*, 101(1–2), 149–160, doi:10.1007/s00704-009-0216-8, 2010.

582 Franssen, H. J. H., Stöckli, R., Lehner, I., Rotenberg, E., and Seneviratne S. I.: Energy balance closure of  
583 eddy-covariance data: A multisite analysis for European FLUXNET stations, *Agric. For. Meteorol.*,  
584 150(12), 1553–1567, doi:10.1016/j.agrformet.2010.08.005, 2010.

585 Gelman, A., and Rubin, D. B.: Inference from iterative simulation using multiple sequences, *Stat. Sci.*, 7, 457–  
586 472, 1992.

587 Gelman, A.: Simulating normalizing constants: From importance sampling to bridge sampling to path sampling,  
588 *Stat. Sci.* 13, 163–185, 1998.

589 Giudice, D., Albert, C., Rieckermann, J., and Reichert, P.: Describing the catchment-averaged precipitation as a  
590 stochastic process improves parameter and input estimation, *Water Resour. Res.*, 52, 3162–3186,  
591 doi:10.1002/2015WR017871, 2016.

592 Höge, M., Wöhling, T., and Nowak, W.: A primer for model selection: The decisive role of model complexity.  
593 *Water Resour. Res.*, 54, 1688–1715, doi.org/10.1002/2017WR021902, 2018.

594 Jefferys, W. H., and Berger, J. O.: Sharpening Ockham’s razor on a Bayesian strop, *Am. Sci.*, 89, 64-72, 1992.

595 Kashyap, R. L.: Optimal choice of AR and MA parts in autoregressive moving average models, *IEEE Trans.*  
596 *Pattern Anal. Mach. Intell.*, 4(2), 99–104, 1982.

597 Katerji, N., and Rana, G.: Modelling evapotranspiration of six irrigated crops under Mediterranean climate  
598 conditions, *Agric. For. Meteorol.*, 138, 142–155, 2006.

599 Katerji, N., Rana, G., Fahed, S.: Parameterizing canopy resistance using mechanistic and semi-empirical  
600 estimates of hourly evapotranspiration: critical evaluation for irrigated crops in the Mediterranean,  
601 *Hydrol. Process.*, 25, 117–129, 2011.

602 Kato, T., Kimura, R., and Kamichika, M.: Estimation of evapotranspiration, transpiration ratio and water-use  
603 efficiency from a sparse canopy using a compartment model, *Agric. Water Manage.*, 65, 173–191,  
604 2004.

605 Kessler, E., and Neas, B.: On correlation, with applications to the radar and raingage measurement of rainfall,  
606 *Atmos. Res.*, 34, 217-229, 1994.

607 Laloy, E., Linde, N., Jacques, D., and Vrugt, J. A.: Probabilistic inference of multi-Gaussian fields from indirect  
608 hydrological data using circulant embedding and dimensionality reduction, *Water Resour. Res.*, 51,  
609 4224–4243, doi:10.1002/2014WR016395, 2015.

610 Lartillot, N., and Philippe, H.: Computing Bayes factors using thermodynamic integration, *Syst. Biol.*, 55(2),  
611 195-207, 2006.

612 Leeb, H., and Pötscher, B. M.: Model selection, Berlin, Germany: Springer., pp, 889–925,  
613 doi.org/10.1007/978-3-540-71297-839, 2009.

614 Legates, D. R., and McCabe, G. J.: Evaluating the use of “goodness-of-fit” measures in hydrologic and  
615 hydroclimatic model validation, *Water Resour. Res.*, 35, 233–241, 1999.

616 Leuning, R., Zhang, Y. Q., Rajaud, A., Cleugh, H., and Tu, K.: A simple surface conductance model to estimate  
617 regional evaporation using MODIS leaf area index and the Penman–Monteith equation, *Water Resour.*  
618 *Res.*, 44, W10419, doi.org/10.1029/2007WR006562, 2008.

619 Liang, J., Zhang, L., Cao, X., Wen, J., Wang, J., and Wang, G.: Energy balance in the semiarid area of the Loess  
620 Plateau, China, *J. Geophys. Res. Atmos.*, 122, 2155–2168, doi:10.1002/2015JD024572, 2017.

621 Li, S., Kang, S., Zhang, L., Ortega-Farias, S., Li, F., Du, T., Tong, L., Wang, S., Ingman, M., and Guo, W.:  
622 Measuring and modeling maize evapotranspiration under plastic film-mulching condition, *J. Hydrol.*,  
623 503, 153–168, 2013.

- 624 Li, S., Zhang, L., Kang, S., Tong, L., Du, T., Hao, X., Zhao, P.: Comparison of several surface resistance models  
625 for estimating crop evapotranspiration over the entire growing season in arid regions. *Agric. For.*  
626 *Meteorol.* 208, 1-15, 2015.
- 627 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.  
628 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.  
629 W.: Heihe Watershed Allied Telemetry Experimental Research (HiWATER): Scientific objectives and  
630 experimental design, *B. Am. Meteorol. Soc.*, 94, 1145–1160, 2013.
- 631 Liu, S. M., Xu, Z. W., Wang, W. Z., Jia, Z. Z., Zhu, M. J., Bai, J., and Wang, J. M.: A comparison of  
632 eddy-covariance and large aperture scintillometer measurements with respect to the energy  
633 balance closure problem, *Hydrol. Earth Syst. Sci.*, 15, 1291–1306, doi:10.5194/hess-15-1291-2011,  
634 2011.
- 635 Liu, G., Liu, Y., Hafeez, M., Xu, D., Vote, C.: Comparison of two methods to derive time series of actual  
636 evapotranspiration using eddy covariance measurements in the southeastern Australia, *J. Hydrol.*,  
637 454–455 (4), 1–6, 2012.
- 638 Liu, P., Elshall, A. S., Ye, M., Beerli, P., Zeng, X., Lu, D., and Tao, Y.: Evaluating marginal likelihood with  
639 thermodynamic integration method and comparison with several other numerical methods. *Water*  
640 *Resour. Res.*, 52(2), 734-758, doi:10.1002/2014WR016718, 2016.
- 641 Marshall, L., Nott, D., and Sharma, A.: Hydrological model selection: A Bayesian alternative, *Water Resour.*  
642 *Res.*, 41(10), 3092-3100, doi: 10.1029/2004WR003719, 2005.
- 643 Matheny, A. M., Bohrer, G., Stoy, P. C., Baker, I. T., Black, A. T., Desai, A. R., Dietze, M. C., Gough, C. M.,  
644 Ivanov, V. Y., Jassal, R. S., Novick, K. A., Schäfer, K. V. R., and Verbeeck, H.: Characterizing the  
645 diurnal patterns of errors in the prediction of evapotranspiration by several land-surface models: An  
646 NACP analysis, *J. Geophys. Res. Biogeosci.*, 119(7), 1458-1473, 2014.
- 647 Monteith, J. L.: Evaporation and environment, *Symp. Soc. Exp. Biol.*, 19, 205–234, 1965.
- 648 Morison, J. I. L., Baker, N. R., Mullineaux, P. M., and Davies, W. J.: Improving water use in crop production,  
649 *Philos. T. Roy. Soc. B.*, 363, 639–658, 2008.
- 650 Neal, R. M.: Markov chain sampling methods for Dirichlet process mixture models, *J. Comput. Graph. Stat.*, 9,  
651 249–265, 2000.
- 652 Oncley, S. P., Foken, T., Vogt, R., Kohsiek, W., DeBruin, H., Bernhofer, C., Christen, A., Van Gorsel, E., Grantz,  
653 D., and Feigenwinter, C.: The energy balance experiment EBEX-2000. Part I: Overview and energy  
654 balance, *Boundary Layer Meteorol.*, 123(1), 1–28, doi:10.1007/s10546-007-9161-1, 2007.
- 655 Ortega-Farias, S., Olioso, A., Fuentes, S., and Valdes, H.: Latent heat flux over a furrow-irrigated tomato crop  
656 using Penman–Monteith equation with a variable surface canopy resistance, *Agric. Water Manage.*, 82,  
657 421–432, 2006.
- 658 Ortega-Farias, S., Poblete-Echeverria, C., and Brisson, N.: Parameterization of a two-layer model for estimating  
659 vineyard evapotranspiration using meteorological measurements, *Agr. For. Meteorol.*, 150, 276–286,  
660 2010.
- 661 Parlange, M. B., and Katul, G. G.: An advection-aridity evaporation model, *Water Resour. Res.*, 28 (1), 127-132,  
662 1992.
- 663 Poblete-Echeverria, C., and Ortega-Farias, S.: Estimation of actual evapotranspiration for a drip-irrigated Merlot  
664 vineyard using a three-source model, *Irrig. Sci.*, 28, 65–78, 2009.
- 665 Priestley, C. H. B., and Taylor, R. J.: On the assessment of surface heat flux and evaporation using large-scale  
666 parameters, *Mon. Weather Rev.*, 100 (2), 81-92, 1972.
- 667 Rana, G., Katerji, N., Ferrara, R.M., and Martinelli, N.: An operational model to estimate hourly and daily crop  
668 evapotranspiration in hilly terrain: validation on wheat and oat crops, *Theory Appl. Climatol.*, 103,  
669 413–426, 2011.
- 670 Sadegh, M., and Vrugt J. A.: Approximate Bayesian Computation using Markov Chain Monte Carlo simulation:  
671 DREAM(ABC), *Water Resour. Res.*, 50, 6767–6787, doi:10.1002/2014WR015386, 2014.
- 672 Samani, S., Ye, M., Zhang, F., Pei, Y. Z., Tang, G. P., Elshall, A. S., and Moghaddam, A. A.: Impacts of prior  
673 parameter distributions on bayesian evaluation of groundwater model complexity, *Water Science &*  
674 *Engineering.*, 11(2), 89-100, doi.org/10.1016/j.wse.2018.06.001, 2018.
- 675 Schöniger, A., Wohling, T., Samaniego, L., and Nowak, W.: Model selection on solid ground: Rigorous  
676 comparison of nine ways to evaluate Bayesian model evidence, *Water Resour. Res.*, 50, 9484–9513,  
677 doi:10.1002/2014WR016062, 2014.
- 678 Schwarz, G.: Estimating the dimension of a model, *Ann. Stat.*, 6(2), 461–464, doi:10.1214/aos/1176344136,

- 679 1978.
- 680 Shao, J.: An asymptotic theory for linear model selection, *Statistica Sinica*, 7(2), 221–242, 1997.
- 681 Shuttleworth, W. J., Gurney, R. J.: The theoretical relationship between foliage temperature and canopy  
682 resistance in sparse crops, *Q. J. Roy. Meteorol. Soc.*, 116, 497–519, 1990.
- 683 Stannard, D. I.: Comparison of Penman-Monteith, Shuttleworth-Wallace, and modified Priestley-Taylor  
684 evapotranspiration models for wildland vegetation in semiarid rangeland, *Water Resour. Res.*, 29 (5),  
685 1379-1392, 1993.
- 686 Stull, R. B.: An introduction to boundary layer meteorology, Kluwer Academic Publ., 255pp, 1988.
- 687 Sumner, D. M., and Jacobs, J. M.: Utility of Penman–Monteith Priestley–Taylor reference evapotranspiration,  
688 and pan evaporation methods to estimate pasture evapotranspiration, *J. Hydrol.*, 308 (1-4), 81-104,  
689 2005.
- 690 Szilagyi, J., and Jozsa, J.: New findings about the complementary relationship based evaporation estimation  
691 methods, *J. Hydrol.*, 354: 171–186, 2008.
- 692 Thomsen, J., Bohrer, G., Matheny, M. V., Ivanov, Y., He, L., Renninger, H., and Schäfer, K.: Contrasting  
693 hydraulic strategies during dry soil conditions in *Quercus rubra* and *Acer rubrum* in a sandy site in  
694 Michigan, *Forests.*, 4(4), 1106–1120, 2013.
- 695 Tsvang, L., Fedorov, M., Kader, B., Zubkovskii, S., Foken, T., Richter, S., and Zeleny, Y.: Turbulent exchange  
696 over a surface with chessboardtype inhomogeneities, *Boundary Layer Meteorol.*, 55(1–2), 141–160,  
697 1991.
- 698 Vinukollu R, K., Wood, E. F., Ferguson, C. R., and Fisher, J. B.: Global estimates of evapotranspiration for  
699 climate studies using multi-sensor remote sensing data: evaluation of three process-based approaches,  
700 *Remote Sens. Environ.*, 115(3), 801–823, 2011.
- 701 Vrugt, J. A., ter Braak, C. J. F., Clark, M. P. J., Hyman, M., and Robinson, B. A.: Treatment of input uncertainty  
702 in hydrologic modeling: Doing hydrology backward with Markov chain Monte Carlo simulation, *Water  
703 Resour. Res.*, 44, W00B09, doi:10.1029/2007WR006720, 2008.
- 704 Vrugt, J. A., ter Braak, C. J. F., Diks, C. G. H., Higdon, D., Robinson, B. A., and Hyman, J. M.: Accelerating  
705 Markov chain Monte Carlo simulation by differential evolution with self-adaptive randomized  
706 subspace sampling, *Int. J. Nonlinear Sci. Numer. Simul.*, 10(3), 273-290, 2009.
- 707 Webb, E. K., Pearman, G. I., and Leuning, R.: Correction of flux measurements for density effects due to heat  
708 and water-vapor transfer, *Q. J. R. Meteorol. Soc.*, 106(447), 85–100, 1980.
- 709 Willmott, C. J.: On the validation of models, *Phys. Geogr.*, 2, 184-194, 1981.
- 710 Xie, W., Lewis, P. O., Fan, Y., Kuo, L., and Chen, M. H.: Improving marginal likelihood estimaton for Bayesian  
711 phylogenetic model selection, *Syst. Biol.*, 60(2), 150-160, 2011.
- 712 Xu, C. Y., and Singh, V. P.: A review on monthly water balance models for water resources investigations, *Water  
713 Resour. Manage.*, 12, 31-50, 1998.
- 714 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.:  
715 Intercomparison of surface energy flux measurement systems used during the HiWATERUSOEXE, *J.  
716 Geophys. Res.*, 118, 13140–13157, 2014.
- 717 Ye, M., Neuman, S. P., and Meyer, P. D.: Maximum likelihood Bayesian averaging of spatial variability models  
718 in unsaturated fractured tuff, *Water Resour. Res.*, 40, W05113, doi:10.1029/2003WR002557, 2004.
- 719 Ye, M., Meyer, P. D., and Neuman, S. P.: On model selection criteria in multimodel analysis, *Water Resour. Res.*,  
720 44, W03428, doi:10.1029/2008WR006803, 2008.
- 721 Zhang, B., Kang, S., Li, F., and Zhang, L.: Comparison of three evapotranspiration models to Bowen  
722 ratio-energy balance method for vineyard in an arid desert region of northwest China, *Agr. Forest  
723 Meteorol.*, 148: 1629–1640, 2008.
- 724 Zhang, X. Y., Liu, C. X., Hu, B. X., and Zhang, G. N.: Uncertainty analysis of multi-rate kinetics of uranium  
725 desorption from sediments, *J. Contam. Hydrol.*, 156(1), 1-15, 2014.
- 726 Zhang, K., Ma, J., Zhu, G., Ma, T., Han, T., and Feng, L. L.: Parameter sensitivity analysis and optimization for  
727 a satellite-based evapotranspiration model across multiple sites using moderate resolution imaging  
728 spectroradiometer and flux data. *Journal of Geophysical Research: Atmospheres*, 122(1), 230-245,  
729 2017.
- 730 Zhu, G. F., Su, Y. H., Li, X., Zhang, K., and Li, C. B.: Estimating actual evapotranspiration from an alpine  
731 grassland on Qinghai–Tibetan plateau using a two-source model and parameter uncertainty analysis by  
732 Bayesian approach, *J. Hydrol.*, 476, 42–51, 2013.
- 733 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

734 assimilating multivariate data sets into the two-source evapotranspiration model by Bayesian approach:  
 735 Application to spring maize in an arid region of northwestern China, Geosci. Model. Dev., 7(4),  
 736 1467–1482, 2014.

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

$A$	Available energy for the whole canopy ( $\text{W m}^{-2}$ )
$A_s$	Available energy <a href="#">for the soil surface</a> ( $\text{W m}^{-2}$ )
$R_n$	Net radiation fluxes into the canopy ( $\text{W m}^{-2}$ )
$R_{ns}$	Net radiation flux into the substrate ( $\text{W m}^{-2}$ )
$G$	Soil heat flux ( $\text{W m}^{-2}$ )
$\lambda ET$	Sum of the latent heat flux from the crop ( $\lambda T$ ) and soil ( $\lambda E$ ) ( $\text{W m}^{-2}$ )
$ET_c$	Canopy transpiration ( $\text{W m}^{-2}$ )
$ET_s$	Soil evaporation ( $\text{W m}^{-2}$ )
$C_c$	Canopy resistance coefficient (dimensionless)
$C_s$	Soil surface resistance coefficient (dimensionless)
$LAI$	Leaf area index
$Q_{50}$	Visible radiation flux <a href="#">when stomatal conductance is half its maximum value</a> ( $\text{W m}^{-2}$ )
$D_{50}$	Vapor pressure deficit <a href="#">at which stomatal conductance is half its maximum value</a> (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 ( $\text{W m}^{-2}$ )
$K_q$	Extinction coefficient
$K_a$	Extinction coefficient
$f$	Fraction of evaporation soil and total evaporation
$\lambda$	Latent heat of water evaporation ( $\text{MJ kg}^{-1}$ )
$\Delta$	Slope of the saturated vapour pressure curve ( $\text{Pa K}^{-1}$ )
$\gamma$	Psychrometric constant ( $\text{kPa K}^{-1}$ )
$\rho$	Density of air ( $\text{kg m}^{-3}$ )
$k$	Karman constant (0.41)
$e_s$	Saturated vapor pressure (kPa)
$e_a$	Actual vapor pressure (kPa)
$q^*$	Saturation-specific humidity at air temperatur ( $\text{kg kg}^{-1}$ )
$q$	Specific humidity of the atmosphere ( $\text{kg kg}^{-1}$ )
$b_1$	Empirical constant ( $\text{s m}^{-1}$ )
$b_2$	Empirical constant ( $\text{s m}^{-1}$ )
$\beta_1$	empirical constant
$\beta_2$	empirical constant
$\theta$	Soil water content ( $\text{m}^3 \text{m}^{-3}$ )
$\theta_a$	Critical water content at which plant stress starts ( $\text{m}^3 \text{m}^{-3}$ )
$\theta_b$	Water content at the wilting point ( $\text{m}^3 \text{m}^{-3}$ )
$\theta_r$	Residual soil water content ( $\text{m}^3 \text{m}^{-3}$ )
$\theta_s$	Saturated water content ( $\text{m}^3 \text{m}^{-3}$ )
$\Theta$	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 $z_m$ ( $\text{m s}^{-1}$ )
$g_a$	Aerodynamic conductance ( $\text{m s}^{-1}$ )
$g_s$	Surface conductance ( $\text{m s}^{-1}$ )
$g_{max}$	Maximum stomatal conductance of leaves at the top of the canopy ( $\text{m s}^{-1}$ )
$g_s^c$	Canopy conductance ( $\text{m s}^{-1}$ )
$r_a$	Aerodynamic resistance ( $\text{s m}^{-1}$ )
$r_a^a$	Aerodynamic resistance between canopy source height and a reference level ( $\text{s m}^{-1}$ )
$r_a^s$	Aerodynamic resistance between the substrate and the canopy source height ( $\text{s m}^{-1}$ )
$r_a^c$	Bulk boundary layer resistance of the vegetation element in the canopy ( $\text{s m}^{-1}$ )
$r_s^s$	Surface resistance of the canopy ( $\text{s m}^{-1}$ );
$r_s^c$	Bulk stomatal resistance of the canopy ( $\text{s m}^{-1}$ )

738

## 739 **Appendix B: Bayesian inference and the DREAM algorithm**

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

$$741 \quad \pi(\boldsymbol{\theta} | D, M) = \frac{\pi(\boldsymbol{\theta} / M) p(D | \boldsymbol{\theta}, M)}{p(D | M)} \quad (\text{A1})$$

742 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  
743 model  $M$  and its parameters  $\boldsymbol{\theta}$ ; and

$$744 \quad p(D | M) = \int p(D | \boldsymbol{\theta}, M) p(\boldsymbol{\theta} | M) d\boldsymbol{\theta} \quad (\text{A2})$$

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

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

$$748 \quad e(t) = D(t) - f(t) \quad (\text{A3})$$

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

751 
$$p(D|\theta) = \prod_{t=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{[e(t)]^2}{2\sigma^2}} \quad (\text{A4})$$

752 where  $n$  is the number of observations and  $\sigma$  represents the error variances.

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

757 **List of Tables**

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

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

762 **Table 3.** Slope and coefficient of determination ( $R^2$ ) of regression between measured and modeled  
 763 half-hourly evapotranspiration values, and statistics of root mean square error (RMSE:  $\text{W m}^{-2}$ ), mean bias  
 764 error (MBE:  $\text{W m}^{-2}$ ), index of agreement (IA), model efficiency (EF) and Logarithm of BME for the four  
 765 ET models.

766 **List of Figures**

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

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

774 **Figure 3.** Regressions between measured and modeled half-hourly ET values produced by different  
 775 models from DOY 154 to DOY 270: (a) PM, (b) SW, (c) PT-FC and (d) AA. The regressions are:  $Y =$   
 776  $0.99X$  ( $R^2 = 0.76$ ),  $Y = 1.05X$  ( $R^2 = 0.82$ ),  $Y = 0.91X$  ( $R^2 = 0.75$ ), and  $Y = 0.92X$  ( $R^2 = 0.75$ ) for the PM,  
 777 SW, PT-FC and AA models, respectively.

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

780 the dataset for the calibration period from DOY 154 to DOY 202.

781 **Figure 5.** Variation of the mean posterior expectation of the potential  $y_k$  with  $\beta_k$  for the PM, SW, PT-FC  
 782 and AA models. [The nits denote natural log units.](#)

783

784 **Table 1** Prior distributions and parameter limits for the PM, SW, PT-FC and AA models. The values are  
 785 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 <sup>-1</sup> )	maximum stomatal conductance	0	50	0	50			Kelliher et al. (1995)
$Q_{50}$ (W m <sup>-2</sup> )	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 <sup>-1</sup> )	empirical constant			4.5	11.3			Sellers et al. (1992)
$b_2$ (s m <sup>-1</sup> )	empirical constant			0	8			Sellers et al. (1992)
$\beta_1$	empirical constant					0.5	1.5	Flint et al. (1991);
$\beta_2$	empirical constant					0.1	10	Barton. (1979)

786

787 **Table 2** Maximum Likelihood Estimates (MLEs), Mean Estimates, 95% High-Probability Intervals  
 788 (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 <sup>-1</sup> )	0.04	0.04	(0.03, 0.04)	0.01	0.01	(0.005, 0.012)			
$Q_{50}$ (W m <sup>-2</sup> )	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 <sup>-1</sup> )				4.51	4.57	(4.52, 4.96)		
$b_2$ (s m <sup>-1</sup> )				0.39	0.57	(0.07, 1.38)		
$\beta_1$							1.1 <sup>a</sup>	1.098 <sup>a</sup> (1.06, 1.16) <sup>a</sup>
							1.5 <sup>b</sup>	1.499 <sup>b</sup> (1.492, 1.499) <sup>b</sup>
$\beta_2$							10.00 <sup>a</sup>	9.75 <sup>a</sup> (7.97, 9.95) <sup>a</sup>
							10.00 <sup>b</sup>	9.94 <sup>b</sup> (9.44, 9.99) <sup>b</sup>

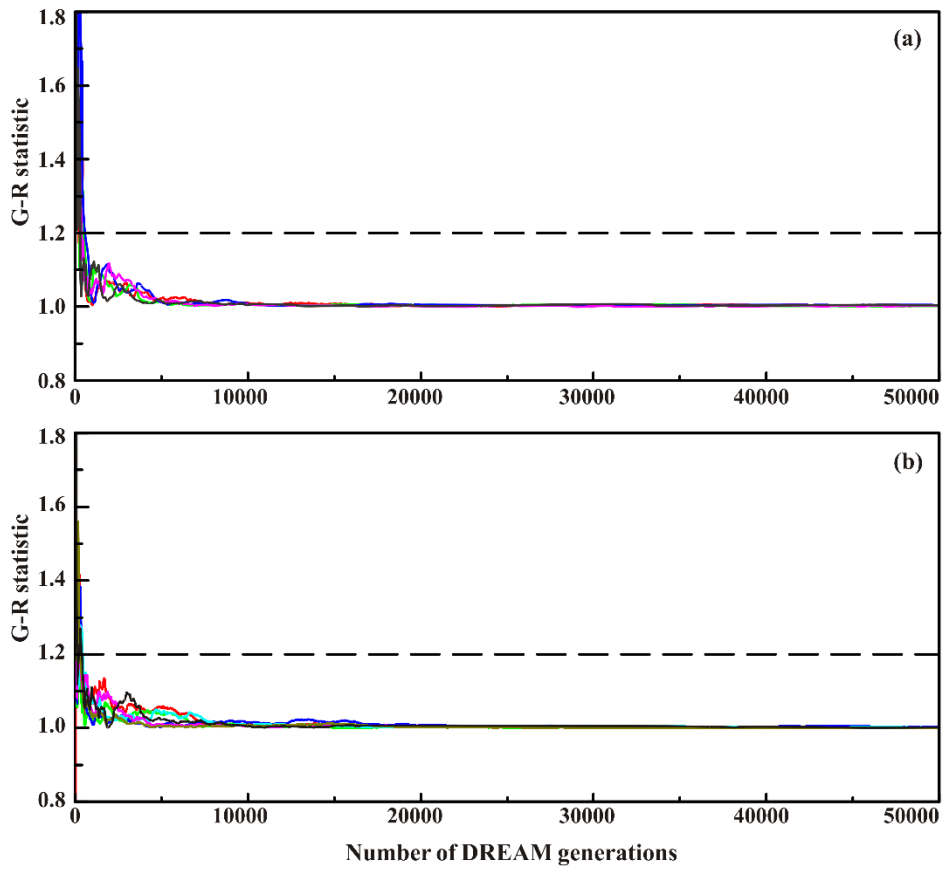
789 <sup>a</sup> PT-FC model; <sup>b</sup> AA model.

790 **Table 3** Slope and coefficient of determination ( $R^2$ ) of regression between measured and modeled  
791 half-hourly evapotranspiration values, and statistics of root mean square error (RMSE:  $W m^{-2}$ ), mean bias  
792 error (MBE:  $W m^{-2}$ ), index of agreement (IA), model efficiency (EF) and Logarithm of BME for the four  
793 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

794



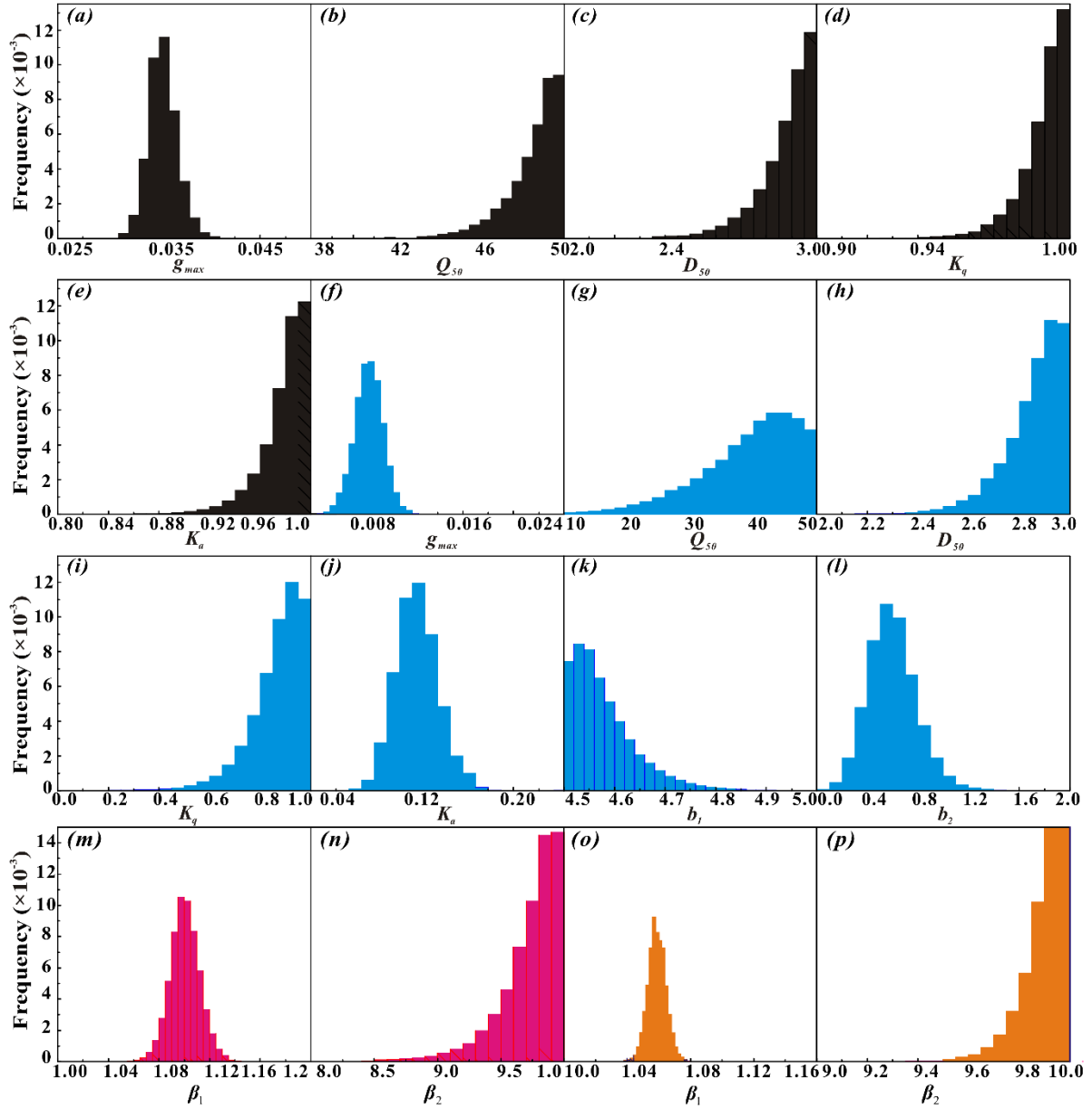


795

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

799

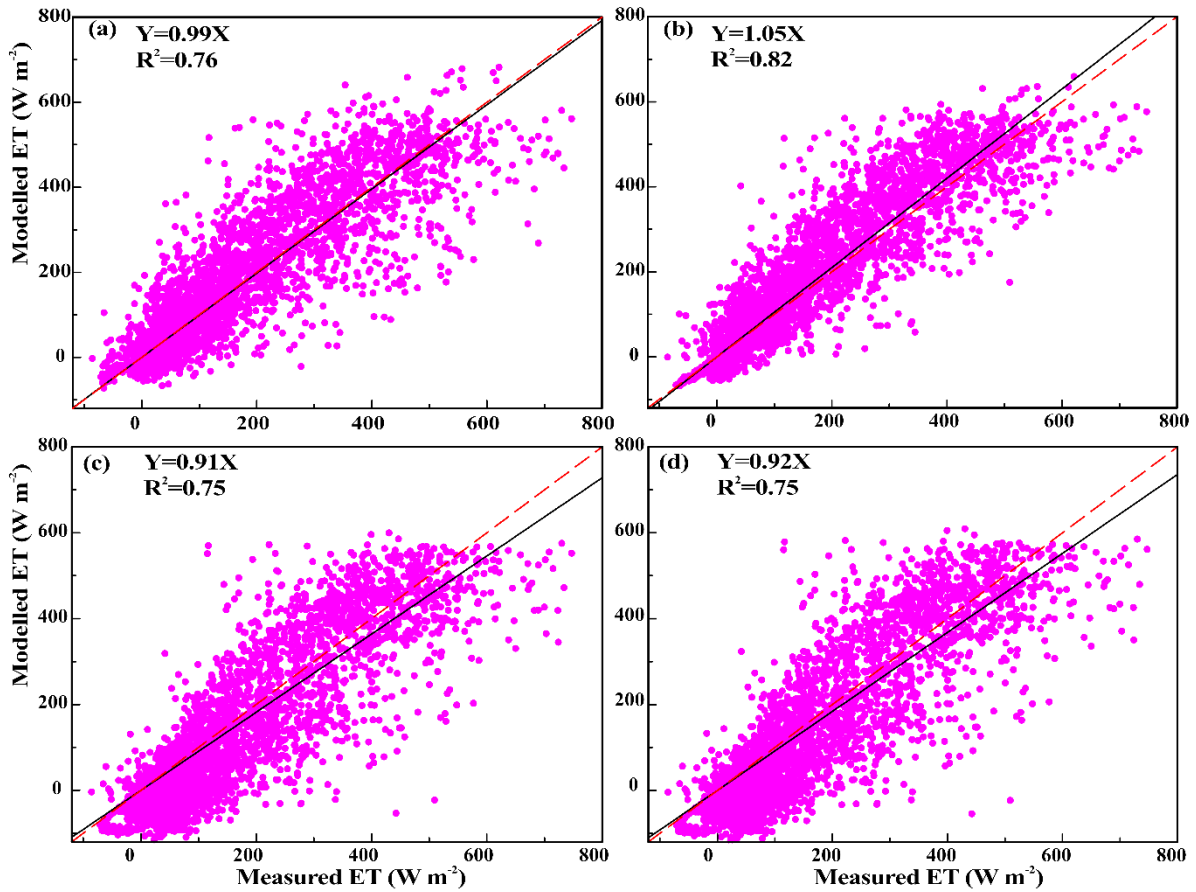
800



801

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

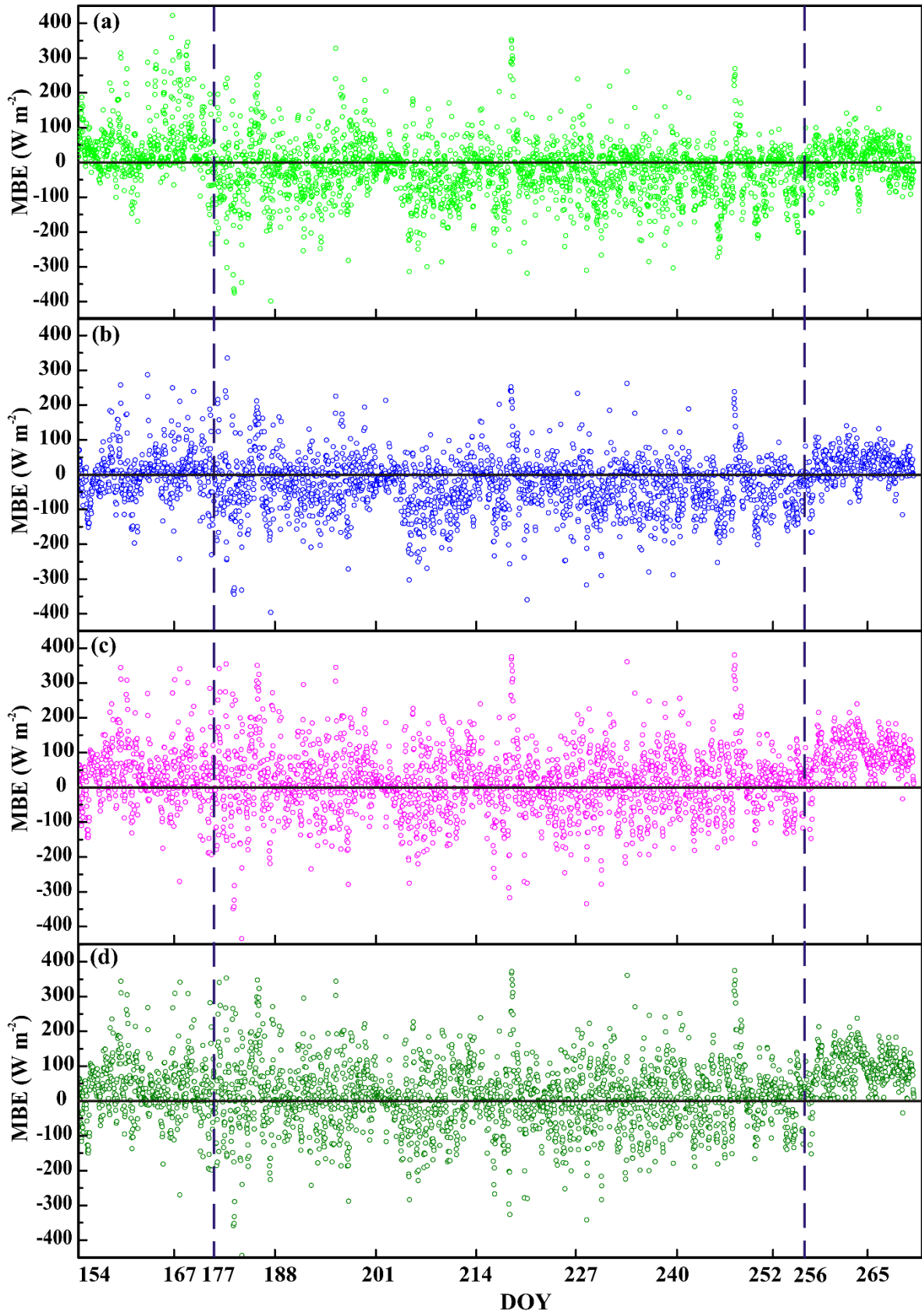
806



807

808 **Figure 3** Regressions between measured and modeled half-hourly ET values produced by different models  
 809 from DOY 154 to DOY 270: (a) PM, (b) SW, (c) PT-FC and (d) AA. The regressions are:  $Y = 0.99X$  ( $R^2 =$   
 810  $0.76$ ),  $Y = 1.05X$  ( $R^2 = 0.82$ ),  $Y = 0.91X$  ( $R^2 = 0.75$ ), and  $Y = 0.92X$  ( $R^2 = 0.75$ ) for the PM, SW, PT-FC  
 811 and AA models, respectively.

812

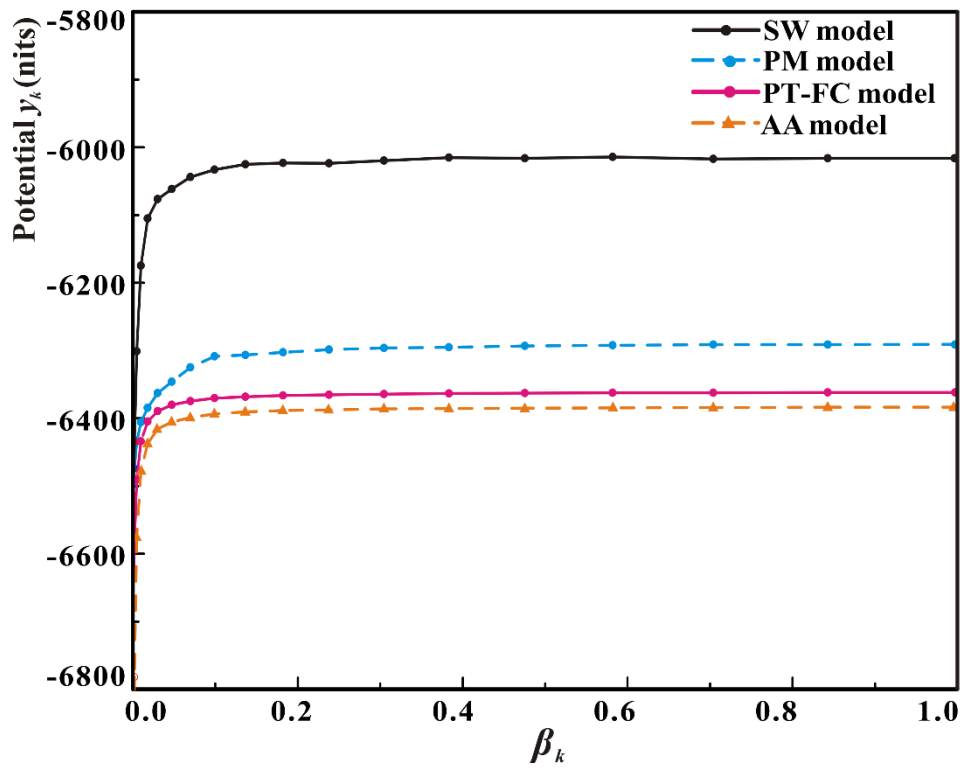


813

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

816 the dataset for the calibration period from DOY 154 to DOY 202.

817



818

819 **Figure 5** Variation of the mean posterior expectation of the potential  $y_k$  with  $\beta_k$  for the PM, SW, PT-FC and  
820 AA models. [The nits denote natural log units.](#)

821