



1 **Bayesian performance evaluation of evapotranspiration models for an arid region in**
2 **northwestern China**

3 Guoxiao Wei ^{1,2,*}, Xiaoying Zhang ^{3,4,*}, Ming Ye ⁵, Ning Yue ^{1,2}, Fei Kan ^{1,2}

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

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

7 ³Institute of Groundwater and Earth Sciences, Jinan University, China, 510632

8 ⁴Construct Engineering College, Jilin University, China, 130400

9 ⁵Department of Earth, Ocean, and Atmospheric Science, Florida State University, USA, 32306

10

11

12

13

14

15

16

17 * Corresponding author: xiaoyingzh@jlu.edu.cn; gxwei@lzu.edu.cn.

18



19 **Abstract**

20 Evapotranspiration (ET) is a major component of the land surface process involved in energy fluxes
21 and balance, especially in the hydrological cycle of agricultural ecosystems. While many models have been
22 developed to estimate ET, there has been no agreement on which model has the best performance. In this
23 study, we evaluate four widely used ET models (i.e., the Shuttleworth Wallace (SW) model,
24 Penman-Monteith (PM) model, Priestley-Taylor and Flint-Childs (PT-FC) model, and Advection-Aridity
25 (AA) model) by using half-hourly ET observations obtained at a spring maize field in an arid region. The
26 model evaluation is based on Bayesian model comparison and ranking using the Bayesian model evidence
27 (BME), which balances between goodness-of-fit to data and model complexity. The BME-based model
28 ranking (from the best to the worst) is SW, PM, PT-FC, and AA. The residuals between observations and
29 corresponding model simulations are also analyzed, and the same model ranking is also obtained by using
30 residual-based statistics, i.e., the coefficient of determination (R^2), index of agreement (IA), root mean
31 square error (RMSE) and model efficiency (EF). The PM and SW models overestimate ET, whereas the
32 PT-FC and AA models underestimate ET in the study period. The four models also underestimate ET
33 during the periods of partial crop cover. Especially during the late maturity stage, the PT-FC and AA
34 models consistently produce an underestimation, and provide the worst simulated ET. As a result, at the
35 half-hourly time scale, the SW model is the best model and recommend as the first choice for evaluating
36 ET of spring maize in arid desert oasis areas.

37 **Keywords:** Evapotranspiration; Bayesian analysis; Penman–Monteith; Shuttleworth-Wallace; Maize

38 **1. Introduction**

39 Surface energy fluxes are an important component of Earth's global energy budget and a primary



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

55 These ET models are generally complex for the coupling of the land surface and atmospheric
56 processes, and high-dimensional with a large number of parameters. Modelers are challenged by how to
57 compare the competing models and how to evaluate the mismatch between model simulations and
58 corresponding observed surface-atmosphere water flux (Legates, 1999). Moreover, how to choose a
59 criterion to reliably evaluate model performance is another crucial issue. Both non-Bayesian analysis
60 (Szilagyi and Jozsa, 2008; Vinukollu et al., 2011; Li et al., 2013; Ershadi et al., 2015) and Bayesian



61 analysis have been used for evaluating model performance (Zhu et al., 2014; Chen et al., 2015; Liu et al.,
62 2016; Zhang et al., 2017; Elshall et al., 2018; Samani et al., 2018; Zeng et al., 2018) . These quantitative
63 criteria used for model evaluation and selection include residual-based measures (e.g., regression line slope
64 and mean bias error, MBE), squared residual-based measures (e.g., coefficient of determination, R^2), root
65 mean square error (RMSE), model efficiency (EF), and index of agreement (IA). Li et al. (2013) compared
66 the maize ET estimates given by PM, SW and adjusted SW models under film-mulching conditions in an
67 arid region of China. They found that the half-hourly ET was overestimated by 17% by the SW model,
68 with relatively high MBE, RMSE, and lower R^2 and IA. In contrast, the PM and MSW models
69 underestimated the daily ET by 6% and 2%, respectively, during the entire experimental period of 116
70 days. Therefore, the performance of PM and MSW models are better than that of the SW model in their
71 case. Ershadi et al. (2014) evaluated the surface energy balance system (SEBS), PM, PT-JPL (modified
72 Priestley–Taylor model, similar to the PT-FC) and AA models. Based on the average value of EF and
73 RMSE, the model rank from the worst to the best was AA, PM, SEBS, and PT-JPL. Ershadi et al. (2015)
74 also evaluated model response to the different formulations of aerodynamic and surface resistances against
75 global FLUXNET data. The results showed significant variability in model performance among and within
76 biome types.

77 The Bayesian model evidence (BME), also known as marginal likelihood, measures the average fit
78 of a model to the data over a model's parameter space. When comparing several alternative conceptual
79 models, the model with the largest marginal likelihood is selected as the best model (Lartillot and Philippe,
80 2006). BME can thus be used for evaluating the model fit (over the parameter space) and for comparing
81 alternative models. In previous studies, Bayesian information criterion (BIC; Kashyap, 1982) or Kashyap



82 information criterion (KIC; Schwarz 1978) were used to approximate BME for reducing computational
83 cost of evaluating BME (Ye et al., 2004). However, these approximations have theoretical and
84 computational limitations (Ye et al., 2008; Xie, 2011; Schöniger et al., 2014), and a numerical evaluation
85 (not an approximation) of BME is necessary, especially for complex models (Lartillot and Philippe, 2006).
86 Lartillot and Philippe (2006) advocated the use of thermodynamic integration (TI) for estimating BME,
87 which is also known as path sampling (Gelman and Meng, 1998; Neal, 2000), to avoid sampling solely in
88 the prior or posterior parameter space. TI uses samples that are systematically generated from the prior to
89 the posterior parameter space by conducting path sampling with several discrete power coefficient values
90 (Liu et al., 2016). It is both mathematically rigorous and more accurate than the generally used harmonic
91 mean method (Xie et al., 2011).

92 While many statistical criteria have been used to evaluate different ET models, BME has not been
93 used for evaluating the ET models. It remains to be determined whether BME can be used to compare and
94 select the best model and whether BME can provide an unbiased view of the performance of the models.
95 Furthermore, most Bayesian applications have focused on the calibration of individual models and
96 comparison of alternative models using these statistical measures, with little attention given to the
97 Bayesian model comparison. Model calibration, comparison, and analysis underlying the Bayesian
98 paradigm has been much less used in the evaluation of ET models than in other areas of environmental
99 science.

100 In this study, the Bayesian approach was used to calibrate and evaluate the four ET models (PM, SW,
101 PT-FC, and AA) based on an experiment over a spring maize field in an arid area of northwest China, from
102 3 June to 27 September 2014. The ET models were calibrated using the Differential Evolution Adaptive



103 Metropolis (DREAM) algorithm. The objectives of the study are as follows: (1) to compare the four
104 models and select the best one using BME; (2) to evaluate various general statistics such as
105 correlation-based measures (R^2), relative error measures (IA and EF), and absolute error measures (such as
106 RMSE and MBE) and to determine whether these methods are efficient and reasonable for evaluating the
107 ET models; (3) to analyze model-data mismatch for better understanding model performance. Using BME
108 for evaluating the ET models has not been reported in the literature. We expect that the study will not only
109 boost the development of model parameterization and model selection but also contribute to the
110 improvement of the ET models.

111

112 **2. Data and methodology**

113 ***2.1. Description of the study area***

114 The experiment was conducted at Daman Superstation, located in Zhangye, Gansu, northwest China.
115 Daman Oasis is located in the middle Heihe River basin, which is the second largest inland river basin in
116 the arid region of northwest China. The midstream area of the Heihe River basin is characterized by oases
117 with irrigated agriculture, and is a major zone of water consumption for domestic and agricultural uses.
118 The annual average precipitation and temperature are 125 mm and 7.2 °C (1960–2000), respectively. The
119 annual accumulated temperature (>10 °C) is 3,234 °C, and the annual average potential evaporation is
120 about 2,290 mm. The average annual duration of sunshine is 3,106 h with 148 frost-free days. The
121 predominant soil type is silty-clay loam and the depth of the frozen layer is about 143 mm. The study area
122 is a typical irrigated agriculture region, and the major water resources are the snowmelt from the Qilian
123 Mountains. The maize and spring wheat are the principal crops, Maize is generally sown in late April and



124 harvested in mid-September and is planted with a row spacing of 40 cm and a plant spacing of 30 cm. The
125 plant density is about 66,000 plants per hectare.

126 ***2.2. Measurements and data processing***

127 Our observation data were collected from the field observation systems of the Heihe Watershed Allied
128 Telemetry Experimental Research (HiWATER) project as described in Li et al (2013). The observation
129 period was from DOY (day of the year) 154 to DOY 270 in 2014. An open-path eddy covariance (EC)
130 system was installed in a maize field, with the sensors at the height of 4.5 m. Maize is the main crop in the
131 study region, which can supply sufficient planting area to set the EC measurements. The EC data was
132 logged at a frequency of 10 Hz and then processed with an average time interval of 30 min. Sensible and
133 latent heat fluxes were computed by the EC approach of Baldocchi (2003). Flux data measured by EC
134 were controlled by traditional routes, including three-dimensional rotation (Aubinet et al., 2000), WPL
135 (Webb-Penman-Leuning) density fluctuation correction (Webb et al., 1980), frequency response correction
136 (Xu et al., 2014), and spurious data removal caused by rainfall, water condensation, and system failure.
137 About 85% of the energy balance closure was observed in the EC data (Liu et al., 2011).

138 Standard hydro-meteorological variables, including rainfall, air temperature, wind speed, and wind
139 direction, were continuously measured at the heights of 3, 5, 10, 15, 20, 30 and 40 m above the ground.
140 Soil temperature and moisture were measured at heights of 2, 4, 10, 20, 40, 80, 120 and 160 cm.
141 Photosynthetically active radiation was measured at a height of 12 m. Net radiation, including downward
142 and upward and longwave radiation, was measured by a four-component net radiometer. An infrared
143 thermometer was installed at a height of 12 m. LAI was measured approximately every 10 days during the
144 growing season.



145 **2.3. Model description**

146 **2.3.1 Penman-Monteith (PM) model**

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

$$149 \quad \lambda E = \frac{\varepsilon A + (\rho C_p / \gamma) D_a g_a}{\varepsilon + 1 + g_a / g_s} \quad (1)$$

150 where $\varepsilon = A/\gamma$; and A is defined as:

$$151 \quad A = R_n - G \quad (2)$$

152 In the present study, g_a is parameterized as suggested by Leuning (2008) and g_s is defined as:

$$153 \quad g_s = g_s^c \frac{\left[1 + \frac{\tau g_a}{(\varepsilon + 1) g_s^c} \left[f - \frac{(\varepsilon + 1)(1 - f) g_s^c}{g_a} \right] + \frac{g_a}{\varepsilon g_i} \right]}{1 - \tau \left[f - \frac{(\varepsilon + 1)(1 - f) g_s^c}{g_a} \right] + \frac{g_a}{\varepsilon g_i}} \quad (3)$$

154 where $1 - \tau$ is the fraction of the total available energy absorbed by the canopy and by the soil, and $\tau = \exp(-$

155 $k_A LAI)$; g_i is defined as:

$$156 \quad g_i = \frac{A}{(\rho C_p / \gamma) D_a} \quad (4)$$

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

$$158 \quad g_s^c = \frac{g_{\max}}{K_Q} \ln \left[\frac{Q_h + Q_{50}}{Q_h \exp(-K_Q LAI) + Q_{50}} \right] \left[\frac{1}{1 + D_a / D_{50}} \right] f(\theta) \quad (5)$$

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



$$f(\theta) = \begin{cases} 1 & \theta > \theta_a \\ \frac{\theta - \theta_b}{\theta_a - \theta_b} & \theta_b < \theta < \theta_a \\ 0 & \theta < \theta_b \end{cases} \quad (6)$$

161 where θ_a was set as $\theta_a = 0.75 \theta_b$. Aerodynamic conductance is calculated as:

$$g_a = \frac{k^2 u_m}{\ln[(z_m - d)/z_{0m}] \ln[(z_m - d)/z_{0v}]} \quad (7)$$

163 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).

164 2.3.2. Shuttleworth-Wallace (SW) model

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

$$167 \quad \lambda ET = \lambda E + \lambda T = C_s ET_s + C_c ET_c \quad (8)$$

$$168 \quad ET_s = \frac{\Delta A + \{\rho C_p (e_s - e_a) - \Delta r_a^s (A - A_s)\} / (r_a^a + r_a^s)}{\Delta + \gamma \{1 + r_s^s / (r_a^a + r_a^s)\}} \quad (9)$$

$$169 \quad ET_c = \frac{\Delta A + \{\rho C_p (e_s - e_a) - \Delta r_a^c A_s\} / (r_a^a + r_a^c)}{\Delta + \gamma \{1 + r_s^c / (r_a^a + r_a^c)\}} \quad (10)$$

170 where the available energy input above the soil surface is defined as:

$$171 \quad A_s = R_{ns} - G \quad (11)$$

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

$$173 \quad R_{ns} = R_n \exp(-K_A LAI) \quad (12)$$

174 The two coefficients C_s and C_c are obtained as follows:



$$175 \quad C_s = \left\{1 + R_s R_a / R_c (R_s + R_a)\right\}^{-1} \quad (13)$$

$$176 \quad C_c = \left\{1 + R_c R_a / R_s (R_c + R_a)\right\}^{-1} \quad (14)$$

177 where R_c , R_a , and R_s are given as:

$$178 \quad R_a = (\Delta + \gamma) r_a^a \quad (15)$$

$$179 \quad R_s = (\Delta + \gamma) r_a^s + \gamma r_s^s \quad (16)$$

$$180 \quad R_c = (\Delta + \gamma) r_a^c + \gamma r_s^c \quad (17)$$

181 Soil surface resistance is expressed as:

$$182 \quad r_s^s = \exp\left(b_1 - b_2 \frac{\theta}{\theta_s}\right) \quad (18)$$

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

184 The 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

185 the boundary layer resistance (r_{ac}) are modeled following the approach proposed by Shuttleworth and

186 Gurney (1990).

187 2.3.3. Priestley–Taylor and Flint–Childs (PT-FC) model

188 The Priestley–Taylor (Priestley and Taylor, 1972) model was introduced to estimate evaporation from

189 an extensive wet surface under conditions of minimum advection (Stannard, 1993; Sumner and Jacobs,

190 2005). It is expressed as:

$$191 \quad \lambda ET = \alpha_{PT} \frac{\Delta}{\Delta + \gamma} (R_n - G) \quad (19)$$



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

$$195 \quad \lambda ET = \alpha \frac{\Delta}{\Delta + \gamma} (R_n - G) \quad (20)$$

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

$$200 \quad \alpha = \beta_1 [1 - \exp(-\beta_2 \Theta)] \quad (21)$$

201 where Θ is calculated as

$$202 \quad \Theta = \frac{\theta - \theta_r}{\theta_s - \theta_r} \quad (22)$$

203 **2.3.4. Advection-aridity (AA) model**

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

$$207 \quad \lambda ET = (2\alpha_{PT} - 1) \frac{\Delta}{\Delta + \gamma} (R_n - G) - \frac{\gamma}{\Delta + \gamma} \frac{\rho(q^* - q)}{r_a} \quad (23)$$

208 where α_{PT} is the Priestley–Taylor coefficient, usually taken as 1.26 (Priestley and Taylor, 1972); and r_a is
 209 similar to that used for the Penman-Monteith model (Brutsaert and Stricker, 1979; Brutsaert, 2005; Ershadi



210 et al., 2014). This model is based mainly on meteorological variables and does not require any information
 211 related to soil moisture, canopy resistance or other measures of aridity (Ershadi et al., 2014). In this study,
 212 similar to the PT-FC model, we modified α_{PT} to α , which is calculated using the same equation as in the
 213 PT-FC model. The detailed list of symbols and physical characteristics in ET models are stated in
 214 Appendix A.

215 **2.4 BME Estimation**

216 **2.4.1 Thermodynamic Integration Estimator**

217 Estimating the BME using power posterior estimators such as thermodynamic integration (TI)
 218 (Lartillot and Philippe, 2006) depends mainly on the likelihood $p(\mathbf{D}|\boldsymbol{\theta}, M)$ calculation. The main idea of
 219 power posterior sampling is to define a path that links the prior to the unnormalized posterior. Thus, using
 220 an unnormalized power posterior density

$$221 \quad q_{\beta}(\boldsymbol{\theta}) = p(\mathbf{D}|\boldsymbol{\theta}, M)^{\beta} p(\boldsymbol{\theta}|M) \quad (24)$$

222 the power coefficient $\beta \in [0, 1]$ is a scalar parameter for discretizing a continuous and differentiable path
 223 linking two unnormalized power posterior densities. The unnormalized power posterior density $q_{\beta}(\boldsymbol{\theta})$ in
 224 Equation (24) uses the normalizing constant Z_{β} to yield the normalized power posterior density:

$$225 \quad p_{\beta}(\boldsymbol{\theta}) = \frac{q_{\beta}(\boldsymbol{\theta})}{Z_{\beta}} \quad (25)$$

226 such that

$$227 \quad Z_{\beta} = \int q_{\beta}(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (26)$$



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

$$229 \quad U(\theta) = \frac{\partial \ln q_{\beta}(\theta)}{\partial \beta} \quad (27)$$

230 thus, the integration can be directly estimated by following:

$$231 \quad p(\mathbf{D}|M) = \frac{Z_1}{Z_0} = \exp\left\{\int_0^1 E_{\theta}[\ln p(\mathbf{D}|\theta, M)] d\beta\right\} \quad (28)$$

232 2.4.2 Power posterior sampling

233 The Metropolis acceptance ratio is $\alpha_k = \min(1, [\alpha_{k, power-posterior} \alpha_{k, prior}])$ with the power posterior

234 ratio given by $\alpha_{k, power-posterior} = (\alpha_{k, posterior})^{\beta_k}$. The prior probability ratio

235 $\alpha_{k, prior} = \Pr(\theta_{k, new} | M) / \Pr(\theta_{k, old} | M)$ is the ratio of the probability of the newly proposed sample

236 $\theta_{k, new}$ and the probability of the previously accepted sample $\theta_{k, old}$. The posterior probability ratio

237 $\alpha_{k, posterior} = L(\mathbf{D} | \theta_{k, new}, M) / L(\mathbf{D} | \theta_{k, old}, M)$ is the likelihood ratio of samples $\theta_{k, new}$ and $\theta_{k, old}$, and β_k

238 is the power posterior coefficient. Thus, to use the DREAM to sample any power posterior distributions

239 (Bayesian inference and the DREAM algorithm please see the Appendix B), the regular Metropolis

240 acceptance ratio $\alpha = \min(1, [\alpha_{posterior} \alpha_{prior}])$ is changed to $\alpha_k = \min(1, [\alpha_{k, power-posterior} \alpha_{k, prior}])$ in DREAM.

241 Since there has been no theoretical method so far for selecting β values (Liu et al., 2016), we

242 determined these values using an empirical but straightforward method. Following Xie et al. (2011), a

243 schedule of the power posterior coefficients β_k is generated by

$$244 \quad \beta_k = (k / K)^{1/\alpha} \quad (29)$$

245 for $k=0, 1, 2, \dots, K$. Using $\alpha = 0.3$ and $K = 20$ is a reasonable initial choice.



246 2.4.3 Estimation of BME using TI

247 After obtaining the power posterior samples, the corresponding likelihoods are used to estimate the
 248 BME. This step simply requires the log-likelihoods $\log p(\mathbf{D}|\boldsymbol{\theta}_{k,i}, M)$ to be inserted into the following
 249 equation:

$$250 \quad \hat{r}_{TI} \equiv p(\mathbf{D} | M) = \exp\left(\int_0^1 y_{\beta} d\beta\right) = \exp\left(\sum_{k=0}^K r_{TI,k}\right) \quad (30)$$

251 such that

$$252 \quad r_{TI,k} = (\beta_k - \beta_{k-1}) \left[\frac{y_k - y_{k-1}}{2} \right] \quad (31)$$

253 and

$$254 \quad y_k = E_{\beta}[\log p(\mathbf{D} | \boldsymbol{\theta}_k, M)] = \frac{1}{n} \sum_{i=1}^n \log p(\mathbf{D} | \boldsymbol{\theta}_{k,i}, M) \quad (32)$$

255 Each panel in this one-dimensional integral is given by $r_{TI,k}^*$ for the case of the trapezoidal rule (Eq.
 256 31), and the summation of these panels gives the natural logarithm of BME.

257 2.5 Traditional statistical evaluations

258 The traditional statistics for evaluating model performance include correlation-based measures of
 259 R^2 and slope, relative error measures of IA and EF, and absolute error measures of RMSE and mean bias
 260 error (MBE) (Poblete-Echeverria and Ortega-Farias, 2009). Their definitions are as follows:

$$261 \quad \text{IA} = 1 - \frac{\sum_{t=1}^n [O(t) - M(t)]^2}{\sum_{t=1}^n \left[|O(t) - \overline{O(t)}| + |O(t) - \overline{M(t)}| \right]^2} \quad (33)$$



$$262 \quad \text{EF} = 1 - \frac{\sum_{t=1}^n [O(t) - M(t)]^2}{\sum_{t=1}^n [O(t) - \overline{O(t)}]^2} \quad (34)$$

$$263 \quad \text{RMSE} = \sqrt{\frac{1}{n} \sum_{t=1}^n [O(t) - M(t)]^2} \quad (35)$$

$$264 \quad \text{MBE} = \frac{1}{n} \sum_{t=1}^n [O(t) - M(t)] \quad (36)$$

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

267

268 **3. Results**

269 **3.1 Parameter estimation**

270 There were five parameters g_{max} , D_{50} , Q_{50} , K_Q and Q_A in the PM model, and two additional parameters,
271 b_1 and b_2 , in the SW model. The PT-FC and AA models include two parameters, denoted as β_1 and β_2
272 (Table 1). The prior probability density of each parameter was specified as a uniform distribution with the
273 ranges listed in Table 1. A total of 50,000 realizations were generated with DREAM and the calibration
274 period data were from DOY 154 to DOY 202. In the calculations, the chain number, N , was equal to the
275 parameter number, i.e., $N = 5, 7, 2$ and 2 for the PM, SW, PT-FC and AA models, respectively. For each
276 model, the first 10,000 samples were discarded as burn-in data, and the remaining 40,000 samples were
277 used to set up posterior density functions for each chain.

278 To understand the efficiency and convergence of DREAM for the ET models, Figure 1 shows the



279 trace plots of the G-R statistic for each of the different parameters in the PM and SW models with the
280 different color (PT-FC and AA models not shown). The algorithm requires about 8,000 generations to
281 make the G-R statistic smaller than 1.2 for the both models. Obviously, the complete mixing of the
282 different chains and convergence of DREAM were attained after about 620 and 450 generations for PM
283 and SW models, respectively. The acceptance rates for the PM and SW models were about 15.3% and
284 18.9%, respectively.

285 Histograms and cumulative distribution functions (CDFs) of the DREAM-derived marginal
286 distributions of the parameters are presented in Figure 2 and summarized in Table 2 by Maximum
287 Likelihood Estimates (MLEs), posterior medians and 95% probability intervals. The uppercase in Figures
288 2A-2E, 2I-2O, 2F-2G, and 2H and 2P showed histograms, and the corresponding lowercase of 2a-2e, 2i-2o,
289 2f-2g, and 2h and 2p show CDFs, for the PM, SW, PT-FC and AA models, respectively.

290 Parameter g_{max} (Fig. 2A) in the PM model, parameters g_{max} , K_A , b_1 , b_2 (Fig. 2I, 2M, 2N, 2O) in the SW
291 model, and parameter β_1 (Fig. 2F) were well constrained and occupied a relatively small range. These
292 parameter sample displayed a unimodal distribution, and became approximately symmetric. Parameters
293 Q_{50} , D_{50} , K_Q and Q_A (Fig. 2B-2E) in the PM model and parameters D_{50} , K_Q in the SW model (Fig. 2K-2L)
294 exhibited relatively large uncertainty reductions. However, the histograms obviously deviated from
295 normality and tended to concentrate in the lower bounds. When the upper limits of these parameters were
296 decreased, similar histograms were reached (not shown) and still did not show statistically meaningful
297 distributions. In contrast, Q_{50} was not only poorly constrained (Fig. 2J) but was also the edge-hitting
298 parameter in the SW model. In addition, the corresponding distributions of the same parameter in different
299 models were slightly different; for example, the mean of g_{max} in the PM model (0.04 mm s^{-1}) was less than



300 that in the SW model (0.01 mm s^{-1}) (Fig. 2A and 2I; Table 2), except that D_{50} in the PM and SW models
301 and β_2 in the PT-FC and AA models exhibited similar regions (Fig. 2C and 2K, 2G and 2P; Table 2).

302 **3.2 Performance of the models**

303 The performances of the four evaporation models were evaluated during the whole season in 2014.
304 The calibrated parameters of the four models were used and individual ET models were run to estimate the
305 half-hourly λET values. Statistical results for the performance of the models were summarized in tables as
306 the regression line slope, R^2 , RMSE, MBE, IA, and EF as shown in Table 3. The regressions between
307 measured and modeled λET values and MBE are shown in Figures 3 and 4, respectively.

308 In general, the four models produced slightly better fits to the measured λET for all the seasons with
309 R^2 larger than 0.75 (Fig. 3). However, obvious discrepancies among models were detected by comparing
310 measured and modeled λET . According to the regression line slope and MBE, the PM model
311 overestimated ET by 1% with a MBE of -9.52 W m^{-2} , and the SW overestimates ET by 5% with a
312 relatively higher MBE of -19.07 W m^{-2} compared to the PM model. The PT-FC and AA models tended to
313 underestimate λET by 9% and 8% with an MBE of 25.42 and 23.29 W m^{-2} , respectively. From a
314 comparison between the slope and MBE, the PM model performance was higher than the SW, PT-FC and
315 AA models, with a slope almost equal to 1 and with relatively lower MBE. The SW model was ranked
316 second, while the AA model was comparable to the PT-FC, but slightly higher, and was ranked third.
317 However, if R^2 , RMSE, IA, and EF were used to evaluate the model performances, the SW model had the
318 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
319 model was the PM model with $R^2 = 0.76$, $\text{RMSE} = 85.38 \text{ W m}^{-2}$, $\text{IA} = 0.93$ and $\text{EF} = 0.74$. The PT-FC
320 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



321 AA model ranked fourth with $R^2 = 0.75$, $RMSE = 95.09 \text{ W m}^{-2}$, $IA = 0.92$ and $EF = 0.67$. Based on the
322 analysis of these traditional statistical criteria, the performances of the PT-FC and AA models yielded
323 similar results. The observed and modeled λET for the four ET models were tightly grouped along the
324 regression lines (Figure 3), and the PT-FC and AA models had similar modeled ET values with a similar
325 degree of point scattering along the regression lines (Figure 3c-3d).

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

333 **3.3 Comparison of the models using BME**

334 Since there was no theoretical method so far for selecting β values, we determined these values using
335 empirical but straightforward methods. For any different power coefficient of $\beta \in [0, 1]$, a sample was
336 drawn from the distribution p_β (Eq. 25) through running DREAM. Figure 5 showed the evolution of \ln
337 $p(D|\theta, M)$ for the four models as a function of β for a dataset covering the entire period. The potential
338 values of the PM model increased from -6533.02 (the logarithm of the prior likelihood) to -6290.71, and
339 the potential values increased from -6544.49 to -6016.17 for the SW model. In addition, the potential
340 values increased from -6708.02 to -6361.76 for the PT-FC model and from -7732.98 to -7033.32 for the
341 AA model. Table 3 showed that the estimated BME is -6300.5 natural log units (nits) for the PM model,



342 -6025.1 nits for the SW model, -6366.8 nits for the PT-FC model, and -7042.8 nits for the AA model. The
343 BME for the SW model was substantially larger than that for the other three models, and the BME for the
344 AA model was the smallest. Although the parameters of the PM model were less than for the SW model,
345 the potential evolution of the SW model was substantially different to that of the PM model. In summary,
346 the PT-FC and AA models, consisting of the same number of parameters, had similar potential patterns of
347 evolution with the coefficient β_k . Although adding more β_k values may improve the BME estimation, it was
348 not undertaken because of the computational cost. For each β_k value, 150,000 DREAM simulations were
349 large enough to ensure convergence.

350 **4. Discussion**

351 ***4.1 Parameter uncertainty analysis***

352 With regard to the efficiency of the DREAM algorithm, the acceptance rates of the PM (15.33%) and
353 SW (18.94%) models were much higher than some MCMC algorithms which used in the previous studies,
354 like 0.019% in the population Monte Carlo sampling algorithm (Sadegh et al., 2014). This was a large
355 improvement in search efficiency, which in large part resulted from its ability to sample groups of variable
356 in turn. Furthermore, this method ran multiple chains in parallel and adaptively updated the scale and
357 orientation of the proposal distribution (Vrugt et al., 2008). Therefore, the DREAM scheme substantially
358 improved not only the convergence, but also its sampling efficiency for ET models.

359 The results showed that the DREAM algorithm successfully reduced the assumed prior uncertainties
360 from the large number of parameters in the four models. The well-constrained parameters were those that
361 had significant contribution. For example, the ecophysiological parameter g_{max} , in both the PM and SW



362 models, the maximum stomatal conductance of leaves, and the soil surface resistance parameters b_1 and b_2
363 in the SW models, all had large influences on the evaluated ET. Thus, their effects were relatively
364 independent compared to the other meteorological parameters in the models. The posterior mean value of
365 g_{max} (0.04 m s^{-1}) in the PM model from our study was close to that (0.05 m s^{-1}) reported in northwestern
366 China (Li et al., 2013; Zhu et al., 2014), but g_{max} (0.01 m s^{-1}) in the SW model was less than the reported
367 value. The estimated posterior means for b_1 and b_2 were different ($b_1 = 9.3$, $b_2 = 6.2$) from those for maize
368 suggested by Zhu et al. (2014) using the same equation of soil surface resistance (r_s^s). Though Zhu et al.
369 (2014) concluded that the responses of g_s^c to VPD and LAI calculated using the modified Leuning model
370 were close to those using Jarvis model (Jarvis, 1976), Li et al. (2015) showed that the performance of PM
371 model was different using the two canopy resistance formula. Therefore, the different results of parameters
372 b_1 and b_2 between our study and the previous study by Zhu et al. (2014) were mainly due to the usage
373 of different canopy conductance models.

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



383 parameters cannot be constrained separately.

384 The sensitive parameters (such as g_{max} , b_1 and b_2) were just corresponding to the well-constrained
385 parameters. Therefore, the major parameters in PM and SW models were well optimized, except that
386 several parameters (Q_{50} and K_D) appeared to be not well constrained. In addition, the posterior parameter
387 bounds exhibited a larger reduction using the DREAM algorithm compared with other studies using the
388 Metropolis–Hasting algorithm (Zhu et al., 2014). This further demonstrated that DREAM can efficiently
389 handle problems involving high-dimensionality, multimodality, nonlinearity, and local optima.

390 In general, parameters related to soil surface resistance in the SW model were well evaluated, while
391 parameters related to canopy surface resistance in PM and SW models were poorly estimated. Therefore,
392 using a reliable canopy surface resistance equation in the ET model was crucial for improving its
393 performance. In addition, in our study, the traditional approach was used to quantify the uncertainty which
394 assumed that the uncertainty mainly came from the parameter uncertainty. However, this method did not
395 explicitly consider errors in the input data and model structural inadequacies. This was unrealistic for real
396 applications, and it was desirable to develop a more reliable inference method to treat all sources of
397 uncertainty separately and appropriately (Vrugt et al., 2008). Moreover, simultaneous direct measurement
398 by micro-lysimeter of sap flow and daily soil evaporation will further help to constrain the model
399 parameters.

400 ***4.2 Evaluation and selection of the models***

401 In this study, the traditional statistical measures and BME were chosen to evaluate and compare the
402 performance of four ET models. From the respective composition of these measures, the statistical



403 measures can be divided into residual-based measures (such as regression slope and MBE) and
404 squared-residual-based measures (such as R^2 , RMSE, IA, and EF). Table 3 shows the values evaluated by
405 BME method, residual-based and squared-residual-based measures. By comparison, the estimates obtained
406 within the same measure (residual-based or squared-residual-based) were congruent. For example, slope
407 and MBE have similar results in the residual-based measures. However, the results from different kind of
408 measures were incongruent; for example, PM model outperformed SW model according to the
409 residual-based measures, but PM model was worse than SW model based on the squared-residual-based
410 measures. The comparative analysis showed a consistency between BME and the squared-residual-based
411 statistics, whereas residual-based criteria were obvious disagreement with the BME measures. It revealed
412 that the more complex SW model was the best model based on the BME and squared-residual-based
413 statistics. The rank order of overall performance of the models from best to worst was: SW, PM, PT-FC,
414 and AA model.

415 Previous studies had shown that BME evaluated by the TI provided estimates similar to the true
416 values and selected the true model if the true model was included within the candidate models (Marshall et
417 al., 2005; Lartillot and Philippe, 2006). Meanwhile, some argued that Bayesian analysis would choose the
418 simplest model (Jefferys and Berger, 1992; Xie et al., 2011) because of the best trade-off between good fit
419 with data and model complexity (Schöniger et al., 2014). In this case, the most complex SW model had the
420 highest BME and was chosen as the best-behaved model. This likely resulted from the fact that the
421 complex SW model was indeed the most reliable model among the alternative ET models. SW model was
422 a two-layer model, and estimated soil evaporation and plant transpiration separately, but PM model was a
423 single-layer model while the plant transpiration and soil evaporation cannot be separated (Monteith, 1965).



424 The PT-FC model was a simplified model of PM, and it only required meteorological and radiation
425 information (Priestley and Taylor, 1972), whereas AA model only relied on the feedback between actual
426 ET and potential ET (Brutsaert and Stricker, 1979). Based on these physical mechanisms and processes for
427 these ET models, the rank order of the models was reasonable.

428 The estimates showed that the maximum values of R^2 , IA and EF, and the minimum value of RMSE,
429 all selected the most complex SW model as the best performing model. The results indicated that the SW
430 model was the best performing model evaluated by squared-residual-based measures, which resulted from
431 the ability of the model to fit the measured data, irrespective of model complexity. It was interesting to
432 note that both the squared-residual-based measures and the BME consistently yielded the same rank order.
433 Although the squared-residual-based measures seemed to identify a reasonable rank order, this had often
434 not been the case, since the simple traditional statistical measures were known to usually provide a biased
435 view of the efficacy of a model (Kessler and Neas, 1994; Legates and McCabe, 1999). In addition,
436 sensitivity to outliers was associated with these measures and leads to relatively high values due to the
437 squaring of the residual terms (Willmott, 1981). Furthermore, these traditional statistical measures ignored
438 the priors, without penalizing model complexity, which was in fact used in a Bayesian analysis. The
439 dimensionality (model's parameter space) not only affected model evaluation by BME (Schöniger et al.,
440 2014) but it may also affect the evaluation using traditional statistical measures. Here, two-dimensional
441 models of PT-FC and AA provided identical estimates of R^2 and IA. This was most likely because both the
442 PT-FC and AA models had the same dimensions and a similar model structure, whereas BME estimates
443 remain well-behaved for the two ET models. Marshall et al. (2005) argued that EF would provide an
444 incorrect conclusion, and Samani et al. (2018) suggested that RMSE also selected the complex model as



445 the best performing model. Thus, we deduced that SRB measures are also problematic. As for slope and
446 MBE, these residual-based measures were obvious disagreement with the BME measure. Part of the lower
447 values of slope and MBE may be counter balanced by the higher values of slope and MBE, thus these
448 criterias provided an erroneous and unreliable model evaluation. Therefore, the squared-residual-based and
449 residual-based measures were not certain to provide reasonable results in terms of model ranking.

450 ***4.3 Analysis of model-data mismatch***

451 Conceptual and structural inadequacies of the hydrological model and measurement errors of the
452 model input (forcing) and output (calibration) data introduced errors in the estimated parameters and model
453 simulations (Laloy, 2014). Hydrological systems were indeed heavily input driven and errors in forcing
454 data can dramatically impair the quality of calibration results and model output (Bardossy and Das, 2008;
455 Giudice, 2015). Measurement errors were raised for a variety of reasons, including unreasonable
456 gap-filling in rainy days; dew and fog; inadequate areal coverage of point-scale soil water measurement;
457 mechanical limitations of the EC system; and inaccurate measurements of wind-speed, soil water, radiation
458 and vapor pressure deficit. ET processe was described using equations that can only capture parts of the
459 complex natural processes and the model structures were an inherent simplification of the real system.
460 These inadequacies can thus lead to biased parameters and implausible predictions.

461 In our study, the results indicated that the PM and SW models overestimated the half-hourly ET
462 compared to the measured ET. Several studies also indicated that the ET values were overestimated by the
463 PM model (Fisher et al., 2005; Ortega-Farias et al., 2006; Li et al., 2015) and the SW model (Li et al., 2013;
464 Li et al., 2015; Zhang et al., 2008). Possible reasons for the inaccurate estimates included the following: (1)
465 Anisotropic turbulence with weak vertical and strong horizontal fluctuation leads to energy imbalance. The



466 total turbulent heat flux was lower by ~10–30% compared to the available energy in many land surface
467 experiments (Tsvang et al., 1991; Beyrich et al., 2002; Oncley et al., 2007; Foken et al., 2010) and influx
468 networks (Franssen et al., 2010). Liang et al. (2017) also showed an energy imbalance result in the
469 semiarid area in China, and indicated that the energy balance closure ratio ranged from 0.52 to 0.90 during
470 the daytime, whereas it was about 0.25 during night-time. However, the measured ET only included
471 vertical flux and not horizontal flux, leading to the measured ET being lower than that of modeled ET by
472 the PM and SW models using the available energy. (2) The absence of a mechanistic representation of the
473 physiological response to plant hydrodynamics cause it difficult for the available ET models to resolve the
474 dynamics of intradaily hysteresis, producing patterns of diurnal error, while the imbalance or lack of
475 between-leaf water demand and soil water supply imposes hydrodynamic limitations on stomatal
476 conductance (Thomsen et al., 2013; Zhang et al., 2014; Matheny et al., 2014). Li et al. (2015) also
477 concluded that neglecting the restrictive effect of the soil on water transport in empirical canopy resistance
478 equations can result in large errors in the partial canopy stage. However, these equations can simulate ET
479 accurately under the full canopy stage (Alves and Pereira, 2000; Katerji and Rana, 2006; Katerji et al.,
480 2011; Rana et al., 2011). Li et al. (2015) showed the PM model combined with the canopy resistance
481 overestimated maize ET during the partial and dense canopy stages by 16% and 13%, respectively
482 (Leuning, 2008). Moreover, the PM model coupled with the canopy resistance overestimated vineyard ET
483 during the entire growth stage by 29% (Leuning, 2008).

484 The estimated ET for the PT-FC and AA models was generally lower than the measured values during
485 the entire season. In addition, the four models also underestimated the ET during periods of partial cover
486 ($LAI < 2.5 \text{ m}^2 \text{ m}^{-2}$). Especially during the late maturity stage, the PT-FC and AA models consistently



487 underestimated ET and provided the worst simulated ET. The underestimation probably resulted from the
488 following: (1) Non classical situations, such as the oasis effect, may occur in the study area. Strong
489 evaporation from the moist ground and plants results in latent heat cooling. However, this upward latent
490 heat flux was opposed by a downward sensible heat flux from the warm air to the cool ground, and thus the
491 latent heat flux was positive while the sensible heat flux is negative. Therefore, the latent heat flux can be
492 greater in magnitude than the solar heating, because of the additional energy extracted from the warm air
493 by evaporation (Stull, 1988). (2) Lack of mechanistic representation of rainfall interception in ET models
494 probably also led to inaccurate simulation on shortly after rainy days. Bohn and Vivoni (2016) found that
495 evaporation of canopy interception accounted for 8% of the annual ET across the North American
496 monsoon region.

497 Comparing the AA and PT-FC models, the former included forcing data of available radiation, soil
498 water content and relative humidity, but the PT-FC model only requires available radiation and soil water
499 content and was independent of relative humidity. However, the similar statistical results and similar
500 degrees of MBE scatter indicated that relative humidity has little influence on the AA model simulation.
501 The consistent and consecutive underestimation of ET by the PT-FC and AA models during the late
502 maturity stage showed that the model-data disagreement is caused mainly by regional advection and
503 rainfall interception, because atmospheric processes and thermally-induced circulation can only occur at
504 certain times and during certain days. Therefore, we suggested that the consistent underestimation of ET
505 by the PT-FC and AA models primarily results from conceptual and structural inadequacies, energy
506 imbalance, and soil water stress. Although the PM and SW models shared a common theoretical basis and
507 the PT-FC model was the simplification of the PM model, these models performed significantly differently.



508 Part of the overestimation of ET by the PM and SW models, caused by coupling with the canopy resistance,
509 may be offset by underestimation caused by energy imbalance and soil water stress. However,
510 underestimation of ET by the PT-FC and AA models cannot be counterbalanced by overestimation during
511 the later maturity stage because the PT-FC and AA models are independent of the canopy resistance.
512 Consequently, the half-hourly patterns of errors in the prediction of ET by the PM and SW models were
513 characterized by symmetry and a low degree of scatter, but the PT-FC and AA models exhibited consistent
514 and asymmetrical error patterns.

515 By contrast, other studies showed that the PM model (Kato et al., 2004) and the SW model (Chen et
516 al., 2015) underestimated half-hourly ET. As for the PT-FC and AA models, while some studies reported
517 that the PT-JPL (Zhang et al., 2017) and AA model showed an overall poor performance, however, other
518 studies have indicated that the AA method performed well for both maize and canola crops (Liu et al.,
519 2012). Therefore, the performance of the four ET models appears to vary not only for different crops and
520 locations (Zhu et al., 2014) but also for different meteorological, physiological and soil conditions.
521 Moreover, the performance was also related to the stage of crop growth.

522

523 **5. Conclusions**

524 This study illustrated the use of the Bayesian approach for the statistical analysis and model selection
525 of four widely used ET models. BME can be used to rank the alternative models in our study, although
526 numerical evaluation of BME is computationally expensive particularly for high-dimensional models.
527 Bayesian model comparison identified the SW model as the best ET model. Although the
528 squared-residual-based measures, including R^2 , IA, RMSE, and EF, provide a congruent model ranking



529 with that of BME, it must be noted that these squared-residual-based measures do not allow using prior
530 information for comparing the models. We advocated that caution is needed when using these statistical
531 methods, and that BME should be used instead. In contrast, residual-based measures disagree with the
532 BME measure, and thus these measures can not be used for evaluating model performance.

533 The model–data mismatches were analyzed to facilitate model improvement after using Bayesian
534 model calibration and comparison. The results indicated that model–data mismatches are mainly resulted
535 from energy imbalance caused by anisotropic turbulence, the absence of a mechanistic representation of
536 the physiological response to plant hydrodynamics, and additional energy induced by advection processes.
537 Among them, energy imbalances and additional energy were related to forcing data error rather than to an
538 unreasonable model structure. Thus, understanding the process of the physiological response to plant
539 hydrodynamics, such as developing or selecting more reasonable and process-based canopy resistance
540 models, was essential for improving the performance of evapotranspiration models. Overall, in our study,
541 the applications of Bayesian calibration, Bayesian model evaluation and analysis of model–data
542 mismatches, provided a promising framework for reducing uncertainty and improving the performance of
543 ET models.

544

545 **Author contribution**

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



549

550 **Competing interests**

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

552

553 **Acknowledgments**

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

561

562 **References**

- 563 Akaike, H.: Information theory and an extension of the maximum likelihood principle, in: Breakthroughs
564 in Statistics, vol. 1, Foundations and Basic Theory, edited by: Kotz, S. and John-15 son, N. L.,
565 Springer-Verlag, New York, USA, 610–624, 1973.
- 566 Allen, R. G., Perista, L. S., Raes, D., and Smith, M.: Crop Evapotranspiration-Guidelines for Computing
567 Crop Water Requirements; FAO Irrigation and Drainage apers-56, FAO-Food and Agriculture
568 Organization of the United Nations, Rome, 1998.
- 569 Alves, I. and Pereira, L. S.: Modeling surface resistance from climatic variables? Agric. Water Manage., 42,



- 570 371–385, 2000.
- 571 Aubinet, M., Grelle, A., Ibrom, A., Rannik, Ü., Moncrieff, J., and Foken, T.: Estimates of the annual net
572 carbon and water exchange of forests: the euroflux methodology, *Adv. Ecol. Res.*, 30(1), 113-175,
573 1999.
- 574 Baldocchi, D. D.: Assessing the eddy covariance technique for evaluating carbon dioxide exchange rates of
575 ecosystems: past, present and future, *Glob. Change. Biol.*, 9, 479–492, 2003.
- 576 Bardossy, A., and Das, T.: Influence of rainfall observation network on model calibration and application,
577 *Hydrol. Earth Syst. Sci.*, 12(1), 77-89, [https://doi:10.5194/hess-12-77-2008](https://doi.org/10.5194/hess-12-77-2008), 2008.
- 578 Barton, I. J.: A Parameterization of the Evaporation from Nonsaturated Surfaces, *J. Appl Meteorol.*, 18(1),
579 43-47, 1979.
- 580 Beyrich, F., Richter, S. H., Weisensee, U., Kohsiek, W., Lohse, H., de Bruin, H. A. R., Foken, T., Göckede,
581 M., Berger, F., Vogt, R., and Batchvarova, E.: Experimental determination of turbulent fluxes over
582 the heterogeneous litfass area: selected results from the litfass-98 experiment, *Theor. Appl.*
583 *Climatol.*, 73(1–2), 19–34, [https://doi:10.1007/s00704-002-0691-7](https://doi.org/10.1007/s00704-002-0691-7), 2002.
- 584 Bohn, T. J., and Vivoni, E. R.: Process-based characterization of evapotranspiration sources over the North
585 American monsoon region, *Water Resour. Res.*, 52, 358–384, [https://doi:10.1002/2015WR017934](https://doi.org/10.1002/2015WR017934),
586 2016.
- 587 Bonan, G.: *Ecological climatology: concepts and applications*, Cambridge University Press, 2008.
- 588 Braswell, B. H., Sacks, W. J., Linder, E., and Schimel, D. S.: Estimating diurnal to annual ecosystem
589 parameters by synthesis of a carbon flux model with eddy covariance net ecosystem exchange
590 observations, *Global. Change. Biol.*, 11, 335–355, 2005,
- 591 Brutsaert, W., and Han, S.: An advection-aridity approach to estimate actual regional evapotranspiration,
592 *Water Resour. Res.*, 15(2), 443-450, 1979.
- 593 Brutsaert, W.: *Hydrology: An Introduction*. Cambridge University Press, Cambridge, 2005.
- 594 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



- 595 Bayesian analysis to compare the performance of three evapotranspiration models for rainfed
596 jujube (*Ziziphus jujuba* Mill.) plantations in the Loess Plateau, *Agr. Water. Manage.*, 159, 341–357,
597 2015.
- 598 Elshall, A. S., Ye, M., Pei, Y., Zhang, F., Niu, G. Y., and Barron-Gafford, G. A.: Relative model score: A
599 scoring rule for evaluating ensemble simulations with application to microbial soil respiration
600 modeling, *Stoch. Env. Res. A.*, 1-11, DOI: 10.1007/s00477-018-1592-3, 2018.
- 601 Ershadi, A., McCabe, M. F., Evans, J. P., Chaney, N. W., and Wood, E. F.: Multi-site evaluation of
602 terrestrial evaporation models using fluxnet data, *Agric. For. Meteorol.*, 187(8), 46-61, 2014.
- 603 Ershadi, A., McCabe, M. F., Evans, J. P., and Wood, E. F.: Impact of model structure and parameterization
604 on Penman–Monteith type evaporation models, *J. Hydrol.*, 525, 521–535, 2015.
- 605 Fisher, J. B., DeBiase, T. A., Qi, Y., Xu, M., and Goldstein, A. H.: Evapotranspiration models compared on
606 a Sierra Nevada forest ecosystem, *Environ. Model. Softw.*, 20 (6), 783–796, 2005.
- 607 Flint A. L., Childs, S. W.: Use of the Priestley–Taylor evaporation equation for soil water limited
608 conditions in a small forest clearcut, *Agric. For. Meteorol.*, 56(3–4), 247–260, 1991.
- 609 Foken, T., Mauder, M., Liebethal, C., Wimmer, F., Beyrich, F., Leps, J. P., Raasch, S., DeBruin, H. A. R.,
610 Meijninger, W. M. L., and Bange, J.: Energy balance closure for the LITFASS-2003 experiment,
611 *Theor. Appl. Climatol.*, 101(1–2), 149–160, <https://doi:10.1007/s00704-009-0216-8>, 2010.
- 612 Franssen, H. J. H., Stöckli, R., Lehner, I., Rotenberg, E., and Seneviratne S. I.: Energy balance closure of
613 eddy-covariance data: A multisite analysis for European FLUXNET stations, *Agric. For. Meteorol.*,
614 150(12), 1553–1567, <https://doi:10.1016/j.agrformet.2010.08.005>, 2010.
- 615 Gelman, A., and Rubin, D. B.: Inference from iterative simulation using multiple sequences, *Stat. Sci.*, 7,
616 457– 472, 1992.
- 617 Gelman, A.: Simulating normalizing constants: From importance sampling to bridge sampling to path
618 sampling, *Stat. Sci.* 13, 163–185, 1998.
- 619 Giudice, D., Albert, C., Rieckermann, J., and Reichert, P.: Describing the catchment-averaged precipitation



- 620 as a stochastic process improves parameter and input estimation, *Water Resour. Res.*, 52,
621 3162–3186, <https://doi.org/10.1002/2015WR017871>, 2016.
- 622 Jarvis, P. G.: The interpretation of the variations in leaf water potential and stomatal conductance found in
623 canopies in the field. *Philos. T. Roy. Soc. B.*, 273, 563–610, 1976.
- 624 Jefferys, W. H., and Berger, J. O.: Sharpening Ockham’s razor on a Bayesian strop, *Am. Sci.*, 89, 64-72,
625 1992.
- 626 Kashyap, R. L.: Optimal choice of AR and MA parts in autoregressive moving average models, *IEEE*
627 *Trans. Pattern Anal. Mach. Intell.*, 4(2), 99–104, 1982.
- 628 Katerji, N., and Rana, G.: Modelling evapotranspiration of six irrigated crops under Mediterranean climate
629 conditions, *Agric. For. Meteorol.*, 138, 142–155, 2006.
- 630 Katerji, N., Rana, G., Fahed, S.: Parameterizing canopy resistance using mechanistic and semi-empirical
631 estimates of hourly evapotranspiration: critical evaluation for irrigated crops in the Mediterranean,
632 *Hydrol. Process.*, 25, 117–129, 2011.
- 633 Kato, T., Kimura, R., and Kamichika, M.: Estimation of evapotranspiration, transpiration ratio and
634 water-use efficiency from a sparse canopy using a compartment model, *Agric. Water Manage.*, 65,
635 173–191, 2004.
- 636 Kessler, E., and Neas, B.: On correlation, with applications to the radar and raingage measurement of
637 rainfall, *Atmos. Res.*, 34, 217-229, 1994.
- 638 Laloy, E., Linde, N., Jacques, D., and Vrugt, J. A.: Probabilistic inference of multi-Gaussian fields from
639 indirect hydrological data using circulant embedding and dimensionality reduction, *Water Resour.*
640 *Res.*, 51, 4224–4243, <https://doi.org/10.1002/2014WR016395>, 2015.
- 641 Lartillot, N., and Philippe, H.: Computing Bayes factors using thermodynamic integration, *Syst. Biol.*,
642 55(2), 195-207, 2006.
- 643 Legates, D. R., and Davis, R. E.: The continuing search for an anthropogenic climate change signal:
644 Limitations of correlation-based approaches, *Geophys. Res. Lett.*, 24, 2319-2322, 1997.



- 645 Legates, D. R., and McCabe, G. J.: Evaluating the use of “goodness-of-fit” measures in hydrologic and
646 hydroclimatic model validation, *Water Resour. Res.*, 35, 233–241, 1999.
- 647 Leuning, R., Zhang, Y. Q., Rajaud, A., Cleugh, H., and Tu, K.: A simple surface conductance model to
648 estimate regional evaporation using MODIS leaf area index and the Penman–Monteith equation,
649 *Water Resour. Res.*, 44, W10419, <http://dx.doi.org/10.1029/2007WR006562>, 2008.
- 650 Liang, J., Zhang, L., Cao, X., Wen, J., Wang, J., and Wang, G.: Energy balance in the semiarid area of the
651 Loess Plateau, China, *J. Geophys. Res. Atmos.*, 122, 2155–2168,
652 <https://doi:10.1002/2015JD024572>, 2017.
- 653 Li, S., Kang, S., Zhang, L., Ortega-Farias, S., Li, F., Du, T., Tong, L., Wang, S., Ingman, M., and Guo, W.:
654 Measuring and modeling maize evapotranspiration under plastic film-mulching condition, *J.*
655 *Hydrol.*, 503, 153–168, 2013.
- 656 Li, S., Zhang, L., Kang, S., Tong, L., Du, T., Hao, X., Zhao, P., 2015. Comparison of several surface
657 resistance models for estimating crop evapotranspiration over the entire growing season in arid
658 regions. *Agric. For. Meteorol.* 208, 1-15
- 659 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.,
660 Wen, J. G., Li, H. Y., Zhu, G. F., Guo, J. W., Ran, Y. H., Wang, S. G., Zhu, Z. L., Zhou, J., Hu, X.
661 L., and Xu, Z. W.: Heihe Watershed Allied Telemetry Experimental Research (HiWATER):
662 Scientific objectives and experimental design, *B. Am. Meteorol. Soc.*, 94, 1145–1160, 2013.
- 663 Liu, S. M., Xu, Z. W., Wang, W. Z., Jia, Z. Z., Zhu, M. J., Bai, J., and Wang, J. M.: A comparison of
664 eddy-covariance and large aperture scintillometer measurements with respect to the energy
665 balance closure problem, *Hydrol. Earth Syst. Sci.*, 15, 1291–1306,
666 <https://doi:10.5194/hess-15-1291-2011>, 2011.
- 667 Liu, G., Liu, Y., Hafeez, M., Xu, D., Vote, C.: Comparison of two methods to derive time series of actual
668 evapotranspiration using eddy covariance measurements in the southeastern Australia, *J. Hydrol.*,
669 454–455 (4), 1–6, 2012.
- 670 Liu, P., Elshall, A. S., Ye, M., Beerli, P., Zeng, X., Lu, D., and Tao, Y.: Evaluating marginal likelihood with



- 671 thermodynamic integration method and comparison with several other numerical methods. *Water*
672 *Resour. Res.*, 52(2), 734-758, [https://doi:10.1002/2014WR016718](https://doi.org/10.1002/2014WR016718), 2016.
- 673 Marshall, L., Nott, D., and Sharma, A.: Hydrological model selection: A Bayesian alternative, *Water*
674 *Resour. Res.*, 41(10), 3092-3100, [https://doi: 10.1029/2004WR003719](https://doi.org/10.1029/2004WR003719), 2005.
- 675 Matheny, A. M., Bohrer, G., Stoy, P. C., Baker, I. T., Black, A. T., Desai, A. R., Dietze, M. C., Gough, C.
676 M., Ivanov, V. Y., Jassal, R. S., Novick, K. A., Schäfer, K. V. R., and Verbeeck, H.:
677 Characterizing the diurnal patterns of errors in the prediction of evapotranspiration by several
678 land-surface models: An NACP analysis, *J. Geophys. Res. Biogeosci.*, 119(7), 1458-1473, 2014.
- 679 Monteith, J. L.: Evaporation and environment, *Symp. Soc. Exp. Biol.*, 19, 205–234, 1965.
- 680 Morison, J. I. L., Baker, N. R., Mullineaux, P. M., and Davies, W. J.: Improving water use in crop
681 production, *Philos. T. Roy. Soc. B.*, 363, 639–658, 2008.
- 682 Neal, R. M.: Markov chain sampling methods for Dirichlet process mixture models, *J. Comput. Graph.*
683 *Stat.*, 9, 249–265, 2000.
- 684 Oncley, S. P., Foken, T., Vogt, R., Kohsiek, W., DeBruin, H., Bernhofer, C., Christen, A., Van Gorsel, E.,
685 Grantz, D., and Feigenwinter, C.: The energy balance experiment EBEX-2000. Part I: Overview
686 and energy balance, *Boundary Layer Meteorol.*, 123(1), 1–28,
687 [https://doi:10.1007/s10546-007-9161-1](https://doi.org/10.1007/s10546-007-9161-1), 2007.
- 688 Ortega-Farias, S., Olioso, A., Fuentes, S., and Valdes, H.: Latent heat flux over a furrow-irrigated tomato
689 crop using Penman–Monteith equation with a variable surface canopy resistance, *Agric. Water*
690 *Manage.*, 82, 421–432, 2006.
- 691 Ortega-Farias, S., Carrasco, M., Olioso, A., Acevedo, C., and Poblete, C.: Latent heat flux over a Cabernet
692 Sauvignon vineyard using the Shuttleworth and Wallace model, *Irrig. Sci.*, 25, 161–170, 2007.
- 693 Ortega-Farias, S., Poblete-Echeverria, C., and Brisson, N.: Parameterization of a two-layer model for
694 estimating vineyard evapotranspiration using meteorological measurements, *Agr. For. Meteorol.*,
695 150, 276–286, 2010.



- 696 Parlange, M. B., and Katul, G. G.: An advection-aridity evaporation model, *Water Resour. Res.*, 28 (1),
697 127-132, 1992.
- 698 Poblete-Echeverria, C., and Ortega-Farias, S.: Estimation of actual evapotranspiration for a drip-irrigated
699 Merlot vineyard using a three-source model, *Irrig. Sci.*, 28, 65–78, 2009.
- 700 Priestley, C. H. B., and Taylor, R. J.: On the assessment of surface heat flux and evaporation using
701 large-scale parameters, *Mon. Weather Rev.*, 100 (2), 81-92, 1972.
- 702 Rana, G., Katerji, N., Ferrara, R.M., and Martinelli, N.: An operational model to estimate hourly and daily
703 crop evapotranspiration in hilly terrain: validation on wheat and oat crops, *Theory Appl. Climatol.*,
704 103, 413–426, 2011.
- 705 Sadegh, M., and Vrugt J. A.: Approximate Bayesian Computation using Markov Chain Monte Carlo
706 simulation: DREAM(ABC), *Water Resour. Res.*, 50, 6767–6787,
707 <https://doi.org/10.1002/2014WR015386>, 2014.
- 708 Samani, S., Ye, M., Zhang, F., Pei, Y. Z., Tang, G. P., Elshall, A. S., and Moghaddam, A. A.: Impacts of
709 prior parameter distributions on bayesian evaluation of groundwater model complexity, *Water
710 Science & Engineering.*, 11(2), 89-100, <https://doi.org/10.1016/j.wse.2018.06.001>, 2018.
- 711 Schöniger, A., Wohling, T., Samaniego, L., and Nowak, W.: Model selection on solid ground: Rigorous
712 comparison of nine ways to evaluate Bayesian model evidence, *Water Resour. Res.*, 50,
713 9484–9513, <https://doi.org/10.1002/2014WR016062>, 2014.
- 714 Schwarz, G.: Estimating the dimension of a model, *Ann. Stat.*, 6(2), 461–464,
715 <https://doi.org/10.1214/aos/1176344136>, 1978.
- 716 Shuttleworth, W. J., Gurney, R. J.: The theoretical relationship between foliage temperature and canopy
717 resistance in sparse crops, *Q. J. Roy. Meteorol. Soc.*, 116, 497–519, 1990.
- 718 Stannard, D. I.: Comparison of Penman-Monteith, Shuttleworth-Wallace, and modified Priestley-Taylor
719 evapotranspiration models for wildland vegetation in semiarid rangeland, *Water Resour. Res.*, 29
720 (5), 1379-1392, 1993.



- 721 Stull, R. B.: An introduction to boundary layer meteorology, Kluwer Academic Publ., 255pp, 1988.
- 722 Sumner, D. M., and Jacobs, J. M.: Utility of Penman–Monteith Priestley–Taylor reference
723 evapotranspiration, and pan evaporation methods to estimate pasture evapotranspiration, J.
724 Hydrol., 308 (1-4), 81-104, 2005.
- 725 Szilagyi, J., and Jozsa, J.: New findings about the complementary relationship based evaporation
726 estimation methods, J. Hydrol., 354: 171–186, 2008.
- 727 Thomsen, J., Bohrer, G., Matheny, M. V., Ivanov, Y., He, L., Renninger, H., and Schäfer, K.: Contrasting
728 hydraulic strategies during dry soil conditions in *Quercus rubra* and *Acer rubrum* in a sandy site in
729 Michigan, Forests., 4(4), 1106–1120, 2013.
- 730 Tsvang, L., Fedorov, M., Kader, B., Zubkovskii, S., Foken, T., Richter, S., and Zeleny, Y.: Turbulent
731 exchange over a surface with chessboardtype inhomogeneities, Boundary Layer Meteorol.,
732 55(1–2), 141–160, 1991.
- 733 Vinukollu R, K., Wood, E. F., Ferguson, C. R., and Fisher, J. B.: Global estimates of evapotranspiration for
734 climate studies using multi-sensor remote sensing data: evaluation of three process-based
735 approaches, Remote Sens. Environ., 115(3), 801–823, 2011.
- 736 Vrugt, J. A., ter Braak, C. J. F., Clark, M. P. J., Hyman, M., and Robinson, B. A.: Treatment of input
737 uncertainty in hydrologic modeling: Doing hydrology backward with Markov chain Monte Carlo
738 simulation, Water Resour. Res., 44, W00B09, [https://doi:10.1029/2007WR006720](https://doi.org/10.1029/2007WR006720), 2008.
- 739 Vrugt, J. A., ter Braak, C. J. F., Diks, C. G. H., Higdon, D., Robinson, B. A., and Hyman, J. M.:
740 Accelerating Markov chain Monte Carlo simulation by differential evolution with self-adaptive
741 randomized subspace sampling, Int. J. Nonlinear Sci. Numer. Simul., 10(3), 273-290, 2009.
- 742 Webb, E. K., Pearman, G. I., and Leuning, R.: Correction of flux measurements for density effects due to
743 heat and water-vapor transfer, Q. J. R. Meteorol. Soc., 106(447), 85–100, 1980.
- 744 Willmott, C. J.: On the validation of models, Phys. Geogr., 2, 184-194, 1981.
- 745 Xie, W., Lewis, P. O., Fan, Y., Kuo, L., and Chen, M. H.: Improving marginal likelihood estimator for



- 746 Bayesian phylogenetic model selection, *Syst. Biol.*, 60(2), 150-160, 2011.
- 747 Xu, C. Y., and Singh, V. P.: A review on monthly water balance models for water resources investigations,
748 *Water Resour. Manage.*, 12, 31-50, 1998.
- 749 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.:
750 Intercomparison of surface energy flux measurement systems used during the HiWATERUSOEXE,
751 *J. Geophys. Res.*, 118, 13140–13157, 2014.
- 752 Ye, M., Neuman, S. P., and Meyer, P. D.: Maximum likelihood Bayesian averaging of spatial variability
753 models in unsaturated fractured tuff, *Water Resour. Res.*, 40, W05113,
754 doi:10.1029/2003WR002557, 2004.
- 755 Ye, M., Meyer, P. D., and Neuman, S. P.: On model selection criteria in multimodel analysis, *Water Resour.*
756 *Res.*, 44, W03428, doi:10.1029/2008WR006803, 2008.
- 757 Zeng, X., Ye, M., Wu, J., Wang, D., & Zhu, X.: Improved nested sampling and surrogate - enabled
758 comparison with other marginal likelihood estimators, *Water Resour. Res.*, 4, DOI:
759 10.1002/2017WR020782, 2018.
- 760 Zhang, K., Ma, J., Zhu, G., Ma, T., Han, T., and Feng, L. L.: Parameter sensitivity analysis and
761 optimization for a satellite-based evapotranspiration model across multiple sites using Moderate
762 Resolution Imaging Spectroradiometer and flux data, *J. Geophys. Res. Atmos.*, 122, 230–245,
763 [https://doi:10.1002/2016JD025768](https://doi.org/10.1002/2016JD025768), 2017.
- 764 Zhang, B., Kang, S., Li, F., and Zhang, L.: Comparison of three evapotranspiration models to Bowen
765 ratio-energy balance method for vineyard in an arid desert region of northwest China, *Agr. Forest
766 Meteorol.*, 148: 1629–1640, 2008.
- 767 Zhu, G. F., Su, Y. H., Li, X., Zhang, K., and Li, C. B.: Estimating actual evapotranspiration from an alpine
768 grassland on Qinghai–Tibetan plateau using a two-source model and parameter uncertainty
769 analysis by Bayesian approach, *J. Hydrol.*, 476, 42–51, 2013.
- 770 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.:



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

774 Zhang, X. Y., Liu, C. X., Hu, B. X., and Zhang, G. N.: Uncertainty analysis of multi-rate kinetics of
 775 uranium desorption from sediments, J. Contam. Hydrol., 156(1), 1-15, 2014.

776 **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)
k_A	Extinction coefficient for available energy
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})

777

778 Appendix B: Bayesian inference and the DREAM algorithm

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

$$780 \quad \pi(\theta | D, M) = \frac{\pi(\theta | M) p(D | \theta, M)}{p(D | M)} \quad (\text{A1})$$

781 where $\pi(\theta | M)$ represents the prior density of θ under model M ; $p(D | \theta, M)$ is the joint likelihood of

782 model M and its parameters θ ; and

$$783 \quad p(D | M) = \int p(D | \theta, M) p(\theta | M) d\theta \quad (\text{A2})$$

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

785 The likelihood function, $p(D | \theta, M)$, used for parameter estimation, is specified according to the

786 distributions of observation errors. Error $\varepsilon(t)$ in each observation $D(t)$ at time t is expressed by



787
$$e(t) = D(t) - f(t) \quad (\text{A3})$$

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

789 expressed as

790
$$p(D|\theta) = \prod_{t=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{[e(t)]^2}{2\sigma^2}} \quad (\text{A4})$$

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

792 In this study, we used the DREAM algorithm (Vrugt et al., 2008, 2009) to explore the ET models'

793 parameter space and to estimate BME. The DREAM sampling scheme is an adaptation of the global

794 optimization algorithm of a shuffled complex evolution metropolis (SCEM-UA). This algorithm was

795 described in more detail in Vrugt et al. (2008, 2009).

796

797 List of Tables

798 **Table 1** Prior distributions and parameter limits for the PM, SW, PT-FC and AA models. The values are

799 derived from the literature.

800 **Table 2** Maximum Likelihood Estimates (MLEs), Mean Estimates, 95% High-Probability Intervals

801 (Lower Limit, Upper Limit).

802 **Table 3** Slope and coefficient of determination (R^2) of regression between measured and modeled

803 half-hourly evapotranspiration values, and statistics of root mean square error (RMSE), mean bias error

804 (MBE), index of agreement (IA), model efficiency (EF) and Logarithm of BME for the four ET models.

805

806 **List of Figures**807 **Figure 1** Trace plots of the G-R statistic of Gelman and Rubin (Gelman and Rubin, 1992) using DREAM

808 for the PM model (a) and (b) the SW model. Different parameters are coded with different colors. The

809 dashed line denotes the default threshold used to diagnose convergence to a limiting distribution.

810 **Figure 2** Uppercase of 2A-2E, 2I-2O, 2F-2G, and 2H and 2P show histograms, and corresponding

811 lowercases of 2a-2e, 2i-2o, 2f-2g, and 2h and 2p show CDFs for the PM, SW, PT-FC and AA models,

812 respectively. The histograms and the CDFs are constructed from the one chain and 40000 generations

813 simulated using DREAM. The y axes represent the prespecified limits of the parameters.

814 **Figure 3.** Regression between measured and modeled half-hourly evapotranspiration values produced by815 different models: (a) PM, (b) SW, (c) PT-FC and (d) AA. The regressions are: $Y = 0.99X$ ($R^2 = 0.76$), $Y =$ 816 $1.05X$ ($R^2 = 0.82$), $Y = 0.91X$ ($R^2 = 0.75$), and $Y = 0.92X$ ($R^2 = 0.75$) for the PM, SW, PT-FC and AA

817 models, respectively.

818 **Figure 4.** Mean bias error (MBE) of predicted and observed ET (W m^{-2}) values for (a) PM, (b) SW, (c)

819 PT-FC and (d) AA models from DOY 154 to DOY 270. Parameters used for prediction are estimated by

820 DREAM with the dataset for the calibration period from DOY 154 to DOY 202.

821 **Figure 5.** Variation of the mean posterior expectation of the potential y_k (equation (36)) with β_k (power

822 coefficient in equation (33)) for the PM, SW, PT-FC and AA models.

823

824 **Table 1** Prior distributions and parameter limits for the PM, SW, PT-FC and AA models. The values are

825 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);
B_2	empirical constant					0.1	10	Barton. (1979)

826

827

828

829 **Table 2** Maximum Likelihood Estimates (MLEs), Mean Estimates, 95% High-Probability Intervals
 830 (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	1.098 ^a	(1.06, 1.16) ^a
	1.5 ^b	1.499 ^b	(1.492, 1.499) ^b
β_2	10.00 ^a		
	10.00	9.75 ^a	(7.97, 9.95) ^a
		9.94 ^b	(9.44, 9.99) ^b

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

832

833

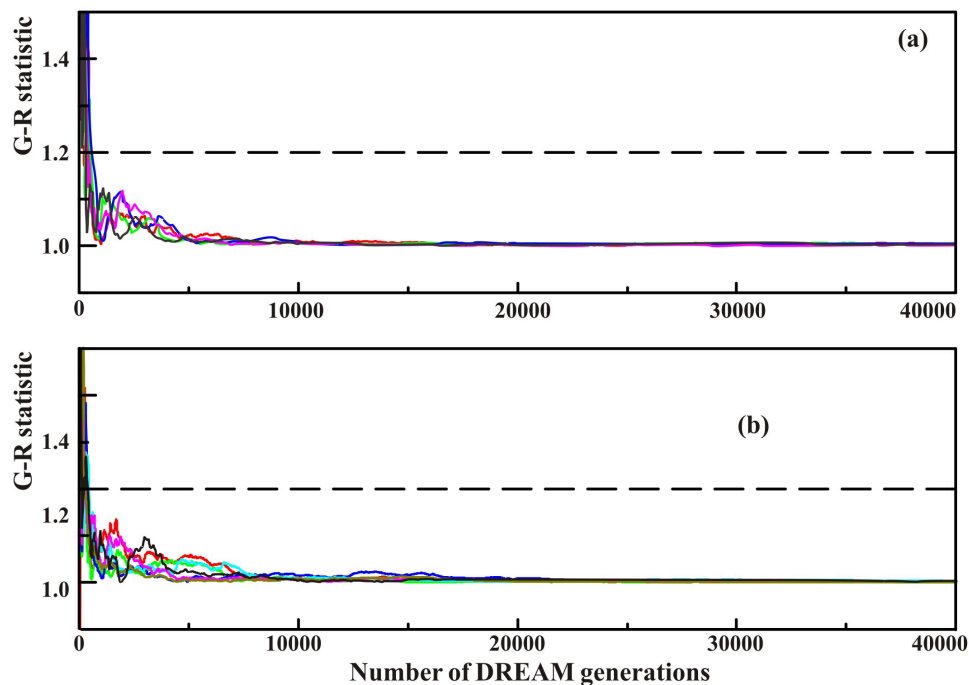
834 **Table 3** Slope and coefficient of determination (R^2) of regression between measured and modeled
 835 half-hourly evapotranspiration values, and statistics of root mean square error (RMSE), mean bias error
 836 (MBE), index of agreement (IA), model efficiency (EF) and Logarithm of BME for the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

837

838

839



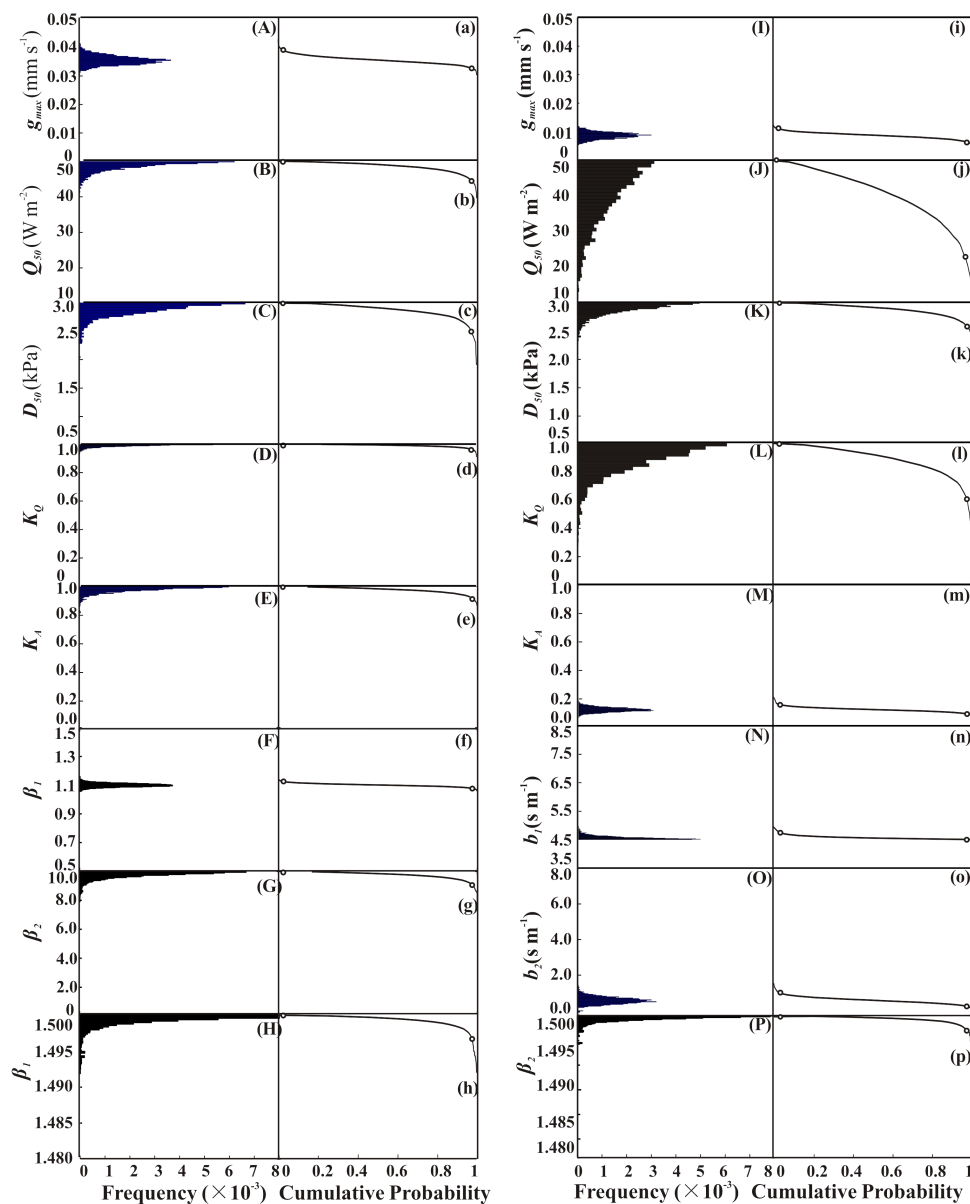
840

841 **Figure 1** Trace plots of the G-R statistic of Gelman and Rubin (Gelman and Rubin, 1992) using DREAM

842 for the PM model (a) and (b) the SW model. Different parameters are coded with different colors. The

843 dashed line denotes the default threshold used to diagnose convergence to a limiting distribution.

844



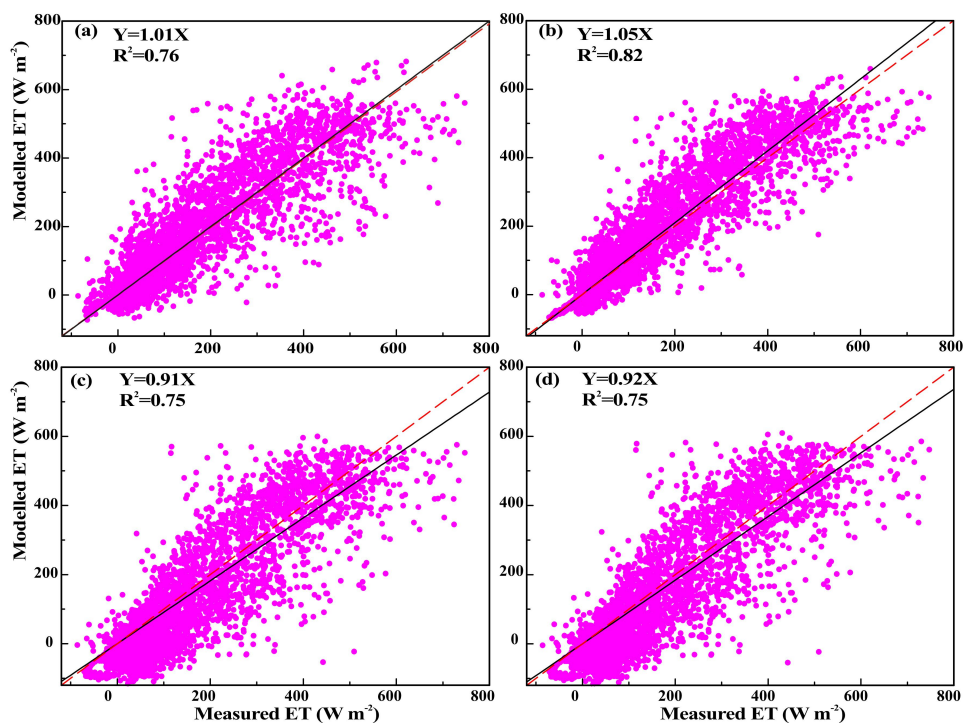
845

846 **Figure 2** Uppercase of 2A-2E, 2I-2O, 2F-2G, and 2H and 2P show histograms, and corresponding

847 lowercases of 2a-2e, 2i-2o, 2f-2g, and 2h and 2p show CDFs for the PM, SW, PT-FC and AA models,

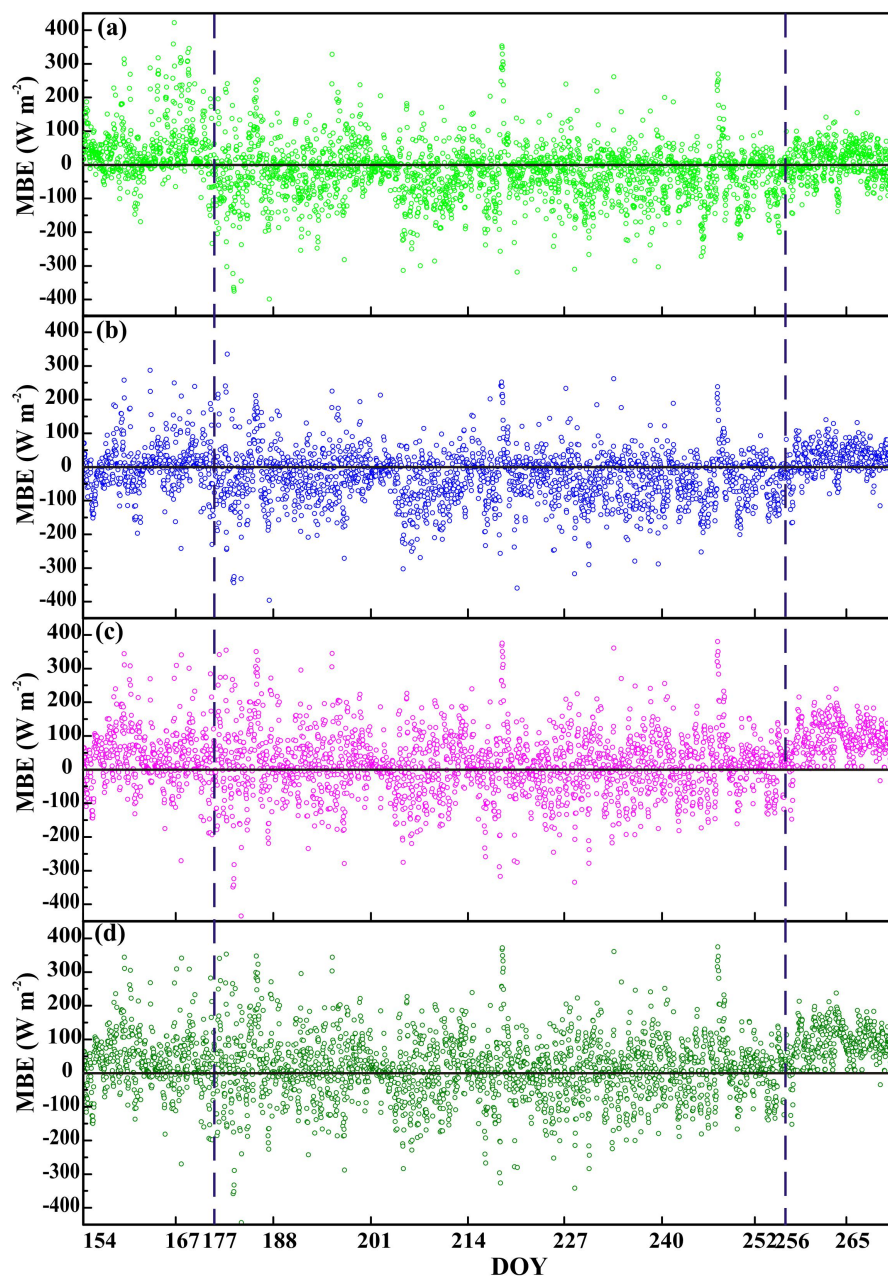
848 respectively. The histograms and the CDFs are constructed from the one chain and 40000 generations

849 simulated using DREAM. The y axes represent the prespecified limits of the parameters.



850

851 **Figure 3.** Regression between measured and modeled half-hourly evapotranspiration values produced by
852 different models: (a) PM, (b) SW, (c) PT-FC and (d) AA. The regressions are: $Y = 0.99X$ ($R^2 = 0.76$), $Y =$
853 $1.05X$ ($R^2 = 0.82$), $Y = 0.91X$ ($R^2 = 0.75$), and $Y = 0.92X$ ($R^2 = 0.75$) for the PM, SW, PT-FC and AA
854 models, respectively.



855

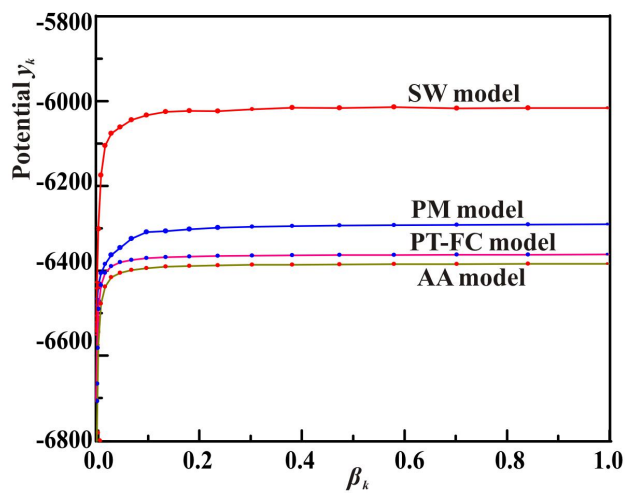
Fig

856 **ure 4.** Mean bias error (MBE) of predicted and observed ET (W m^{-2}) values for (a) PM, (b) SW, (c) PT-FC

857 and (d) AA models from DOY 154 to DOY 270. Parameters used for prediction are estimated by DREAM



858 with the dataset for the calibration period from DOY 154 to DOY 202.



859

860 **Figure 5.** Variation of the mean posterior expectation of the potential y_k (equation (36)) with β_k (power

861 coefficient in equation (33)) for the PM, SW, PT-FC and AA models.

862