

1 **Evaluation of water flux predictive models developed using eddy**
2 **covariance observations and machine learning: a meta-analysis**

3 Haiyang Shi^{1,2,4,5}, Geping Luo^{1,2,3,5}, Olaf Hellwich⁶, Mingjuan Xie^{1,2,4,5}, Chen Zhang^{1,2}, Yu Zhang^{1,2}, Yuangang
4 Wang^{1,2}, Xiuliang Yuan¹, Xiaofei Ma¹, Wenqiang Zhang^{1,2,4,5}, Alishir Kurban^{1,2,3,5}, Philippe De Maeyer^{1,2,4,5} and
5 Tim Van de Voorde^{4,5}

6
7 ¹State Key Laboratory of Desert and Oasis Ecology, Xinjiang Institute of Ecology and Geography, Chinese
8 Academy of Sciences, Urumqi, Xinjiang, 830011, China.

9 ²University of Chinese Academy of Sciences, 19 (A) Yuquan Road, Beijing, 100049, China.

10 ³Research Centre for Ecology and Environment of Central Asia, Chinese Academy of Sciences, Urumqi, China.

11 ⁴Department of Geography, Ghent University, Ghent 9000, Belgium.

12 ⁵Sino-Belgian Joint Laboratory of Geo-Information, Ghent, Belgium and Urumqi, China.

13 ⁶Department of Computer Vision & Remote Sensing, Technische Universität Berlin, 10587 Berlin, Germany.

14

15 **Correspondence to: Geping Luo (luogp@ms.xjb.ac.cn) and Olaf Hellwich (olaf.hellwich@tu-berlin.de)**

16 Submitted to *Hydrology and Earth System Sciences*

17 **Abstract.**

18 With the rapid accumulation of water flux observations from global eddy-covariance flux sites, many studies
19 have used data-driven approaches to model site-scale water fluxes with various predictors and machine learning
20 algorithms used. However, systematic evaluation of such models is still limited. We therefore performed a meta-
21 analysis of 32 such studies, derived 139 model records, and evaluated the impact of various features on model
22 accuracy throughout the modeling flow. SVM (average R-squared = 0.82) and RF (average R-squared = 0.81)
23 outperformed over evaluated algorithms with sufficient sample size in both cross-study and intra-study (with the
24 same data) comparisons. The average accuracy of the model applied to arid regions is higher than in other
25 climate types. The average accuracy of the model was slightly lower for forest sites (average R-squared = 0.76)
26 than for croplands and grasslands (average R-squared = 0.8 and 0.79), but higher than for shrubland sites
27 (average R-squared = 0.67). Using Rn/Rs, precipitation, Ta, and FAPAR improved the model accuracy. The
28 combined use of Ta and Rn/Rs is very effective especially in forests, while in grasslands the combination of Ws
29 and Rn/Rs is also effective. Random cross-validation showed higher model accuracy than spatial cross-
30 validation and temporal cross-validation, but spatial cross-validation is more important in spatial extrapolation.
31 The findings of this study are promising to guide future research on such machine learning-based modeling.

32 **1 Introduction**

33 Evapotranspiration (ET) is one of the most important components of the water cycle in terrestrial ecosystems. It
34 also represents the key variable in linking ecosystem functioning, carbon and climate feedbacks, agricultural
35 management, and water resources (Fisher et al., 2017). The quantification of ET for regional, continents, or the
36 globe can improve our understanding of the water, heat, and carbon interactions, which is important for global
37 change research (Xu et al., 2018). Information on ET has been used in many fields, including, but not limited to,
38 droughts and heatwaves (Miralles et al., 2014), regional water balance closures (Chen et al., 2014; Sahoo et al.,
39 2011), agricultural management (Allen et al., 2011), water resources management (Anderson et al., 2012),
40 biodiversity patterns (Gaston, 2000). In addition, accurate large-scale and long-time series ET prediction at high
41 spatial and temporal resolution has been of great interest (Fisher et al., 2017).

42
43 Currently, there are three main approaches for simulation and spatial and temporal prediction of ET: (i) physical
44 models based on remote sensing such as surface energy balance models (Minacapilli et al., 2009; Wagle et al.,
45 2017), Penman-Monteith equation (Mu et al., 2011; Zhang et al., 2010), Priestley-Taylor equation (Miralles et
46 al., 2011); (ii) process-based land surface models, biogeochemical models and hydrological models (Barman et
47 al., 2014; Pan et al., 2015; Sándor et al., 2016; Chen et al., 2019); and (iii) the observation-based machine
48 learning modeling approach with in situ eddy covariance (EC) observations of water flux (Jung et al., 2011; Li
49 et al., 2018; Van Wijk and Bouten, 1999; Xie et al., 2021; Xu et al., 2018; Yang et al., 2006; Zhang et al., 2021).
50 For remote sensing-based physical models and process-based land surface models, some physical processes
51 have not been well characterized due to the lack of understanding of the detailed mechanisms influencing ET
52 under different environmental conditions. For example, the inaccurate representation and estimation of stomatal
53 conductance (Li et al., 2019) and the linearization (McColl, 2020) of the Clausius-Clapeyron relation in the
54 Penman-Monteith equation may introduce both empirical and conceptual errors into estimates of ET. Limited by

55 complicated assumptions and model parametrizations, these process-based models face challenges in the
56 accuracy of their ET estimations over heterogeneous landscapes (Pan et al., 2020; Zhang et al., 2021).
57 Therefore, many researchers have used data-driven approaches for the simulation and prediction of ET with the
58 accumulation of a large volume of measured site-scale observational data of water fluxes in the past decades.
59 Various machine learning models have been developed to simulate water fluxes at the flux site scale. Besides,
60 various predictor variables (e.g., meteorological factors, vegetation conditions, and moisture supply conditions)
61 have been incorporated into such models for upscaling (Fang et al., 2020; Jung et al., 2009) of water flux to a
62 larger scale or understanding the driving mechanisms with the variable importance analysis performed in such
63 models.

64

65 However, to date, the systematic assessment of the uncertainty in the processes of water flux prediction models
66 constructed using the machine learning approach is limited. Although considerable effort has been invested in
67 improving the accuracy of such prediction models, our understanding of the expected accuracy of such models
68 under different conditions is still limited. It is still not easy for us to give the general guidelines for selecting
69 appropriate predictor variables and models. Questions such as ‘Which predictor variables are the best in water
70 flux simulations?’ and ‘How to improve the prediction accuracy of water flux effectively?’ etc. still confuse the
71 researchers in the field. Therefore, we should synthesize the findings from published studies to determine which
72 predictor variables, machine learning models, and other features can significantly improve the prediction
73 accuracy of water flux. Also, we are interested in understanding under which specific conditions they are more
74 effective.

75

76 A variety of features control the accuracy of such models, including the predictor variables used, the inherent
77 heterogeneity within the dataset, the plant functional type (PFT) of the flux sites, the method of model
78 construction and validation, and the algorithm chosen:

79 a) Predictor variables used: Compared to process-based models, the data used may have a more significant
80 impact on the final model performance in data-driven models. Various biophysical covariates and other
81 environmental factors have been used for the simulation and prediction of water fluxes. The most
82 commonly used factors include mainly precipitation (Prec), air temperature (T_a), wind speed (W_s), net/sun
83 radiation (R_n/R_s), soil temperature (T_s), soil texture, vapor-pressure deficit (VPD), the fraction of absorbed
84 photosynthetically active radiation (FAPAR), vegetation index (e.g., Normalized Difference Vegetation
85 Index (NDVI), Enhanced Vegetation Index (EVI)), Leaf area index (LAI), and carbon fluxes (e.g., Gross
86 Primary Productivity (GPP)). These used predictor variables and their complex interactions drive the
87 fluctuations and variability of water fluxes. They affect the accuracy of water flux simulations in two ways:
88 their actual impact on water fluxes at the process-based level and their spatio-temporal resolution and
89 inherent accuracy. The relationship between water fluxes and these variables at the process-based driving
90 mechanism level is very different under different PFTs, different climate types, and different
91 hydrometeorological conditions. For example, in irrigated croplands in arid regions, water fluxes may be
92 highly correlated with irrigation practices, and thus soil moisture may be a very important predictor
93 variable, and its importance may be significantly higher than in other PFTs. And in models that incorporate
94 data from multiple PFTs, some variables that play important roles in multiple PFTs may have higher

95 importance. In terms of data spatial and temporal resolution, the data for these predictor variables may have
96 different scales. In terms of spatial resolution, meteorological observations such as precipitation and air
97 temperature are at the flux site scale, while factors extracted from satellite remote sensing and reanalysis
98 climate datasets cover a much larger spatial scale (i.e. the grid-scale). This leads to considerable differences
99 in the degree of spatial match between different variables and the site scale EC observations (approximately
100 100 m x 100 m). It is therefore difficult for some variables to be fairly compared in the subsequent
101 importance analysis of driving factors. In terms of temporal resolution, the importance of predictor
102 variables with different temporal resolutions may be variable for models with different time scales (e.g.,
103 half-hourly, daily, and monthly models). For example, the daily or 8-day NDVI data based on MODIS
104 satellite imagery may better capture the temporal dynamics of water fluxes concerning vegetation growth
105 than the 16-daily NDVI data derived from Landsat images. Besides, data on non-temporal dynamic
106 variables such as soil texture cannot explain temporal variability in water fluxes in the data-driven
107 simulations, although soil texture may be important in the interpretation of the actual driving mechanisms
108 of ET (which may need to be quantified in detail in ET simulations by process-based models). In addition,
109 some inherent accuracy issues (e.g., remote sensing-based NDVI may not be effective at high values) of the
110 predictors may propagate into the consequent machine learning models, thus affecting the modeling and our
111 understanding of its importance. Therefore, it is necessary to consider the spatial and temporal resolution of
112 the data and their inherent accuracy for the predictors used in different studies in the systematic evaluation
113 of data-driven water flux simulations.

- 114 b) The heterogeneity of the dataset and model validation: the volume and inherent spatiotemporal
115 heterogeneity of the training dataset (with more variability and extremes incorporated) may affect model
116 accuracy. Typically, training data with larger regions, multiple sites, multiple PFTs, and longer year spans
117 may have a higher degree of imbalance (Kaur et al., 2019; Van Hulse et al., 2007; Virkkala et al., 2021;
118 Zeng et al., 2020). And in machine learning, in general, modeling with unbalanced data (with significant
119 differences in the distribution between the training and validation sets) may result in lower model accuracy.
120 Currently, the most common ways of model validation include spatial, temporal, and random cross-
121 validation. Spatial validation is mainly to evaluate the ability of the model to be applied in different regions
122 or flux sites with different PFT types, and one of the common methods is 'leave one site out' (Fang et al.,
123 2020; Papale et al., 2015; Zhang et al., 2021). If the data of the site left out for validation differs
124 significantly from the distribution of the training data set, the expected accuracy of the model applied at that
125 site may be low because the trained model may not capture the specific and local relationships between the
126 water flux and the various predictor variables at that site. For temporal validation, to assess the ability of the
127 models to adapt to the interannual variability, typically some years of data are used for training and the
128 remaining years for model validation (Lu and Zhuang, 2010). If a year with extreme climate is used for
129 validation, the accuracy may be low because the training dataset may not contain such extreme climate
130 conditions. In the case of PFTs that are significantly affected by human activities, such as cropland, the
131 possible different crops grown and different land use practices (e.g., irrigation) across years can also lead to
132 low accuracy in temporal validation.
- 133 c) Various machine learning algorithms: Some machine learning algorithms may have specific advantages
134 when applied to model the relationships between water fluxes and covariates. For example, neural networks

135 may have an advantage in nonlinear fitting, while random forests can avoid serious overfitting problems.
136 However, which algorithm is better overall in different situations (i.e. applied to different data sets)? Which
137 algorithm is generally more accurate than the others when using the same data set? A comprehensive
138 evaluation is important.

139

140 Therefore, to systematically and comprehensively assess the impact of various features in such modeling, we
141 perform a meta-analysis of published water flux simulation studies that combine the flux site water flux
142 observations, various predictors, and machine learning. The accuracy of model records collected from the
143 literature was linked with various model features to assess the impacts of predictor data types, algorithms, and
144 other features on model accuracy. The findings of this study may be promising to improve our understanding of
145 the impact of various features of the models to guide future research on such machine learning-based modeling.

146 **2 Methodology**

147 **2.1 Protocol for selecting the sample of articles**

148 We applied a general query (on December 1st, 2021) on title, abstract, and keywords to include articles with the
149 “OR” operator applied among expressions (Table 1) in the Scopus database. Preferred Reporting Items for
150 Systematic Reviews and Meta-Analyses (PRISMA) (Moher et al., 2009) are followed when filtering the papers.
151 We first excluded articles that obviously did not fit the topic of this study based on the abstract, and then
152 performed the article screening with the full-text reading.

153

154 The inclusion of articles follows the following criteria:

- 155 a) Articles were filtered for those with water fluxes (or latent heat) simulated.
- 156 b) The water flux or latent heat observations used in the prediction models should be from the eddy-
157 covariance flux measurements.
- 158 c) Articles focusing only on gap-filling (Hui et al., 2004) techniques (i.e., the objective was not simulation
159 and extrapolation of water fluxes using machine learning) were excluded.
- 160 d) Only articles that used multivariate regression (with the number of covariates greater than or equal to 3)
161 were included.
- 162 e) The determination coefficient (R-squared) of the validation step should be reported as the metric of model
163 performance (Shi et al., 2021; Tramontana et al., 2016; Zeng et al., 2020) in the articles.
- 164 f) The articles should be published in English-language journals.

165

166 Although RMSE is also often used for model accuracy assessment, its dependence on the magnitude of water
167 flux values makes it difficult to use for fair comparisons between studies. For example, due to the difference in
168 the range of ET values, models developed from flux stations in dry grasslands will typically have lower RMSE
169 than models developed by flux stations based on forests in humid regions. Therefore, RMSE may not be a good
170 metric for cross-study comparisons in this meta-analysis.

171

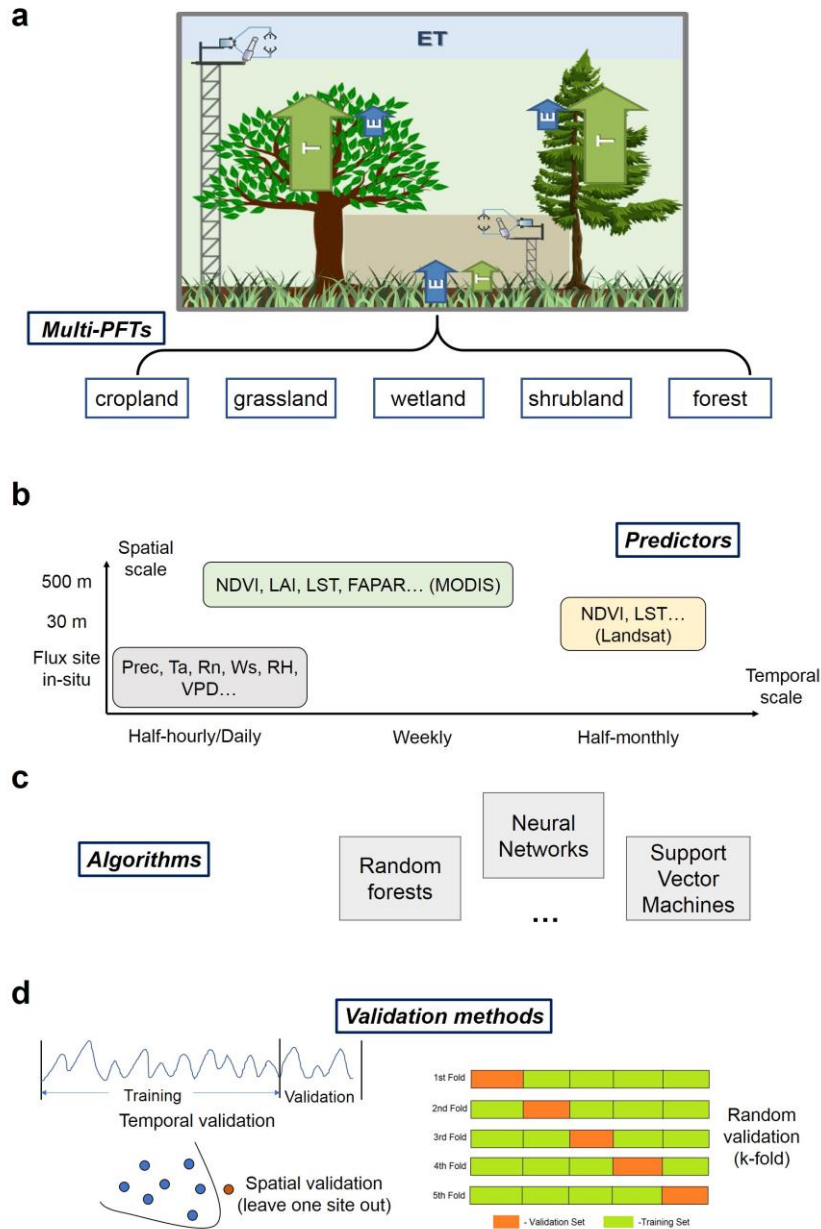
172 Table 1. Article search: ‘[A1 OR A2 OR A3...] AND [B1 OR B2 OR B3...] AND [C1 OR C2 OR C3 OR C4...]’

ID	A	B	C
1	Water flux	Eddy covariance	Machine learning
2	Evapotranspiration	Flux tower	Support Vector
3	Latent heat	Flux site	Neural Network
4			Random Forest

173

174 **2.2 Features of the prediction processes evaluated**

175 The various features (Table 2) involved in the water flux modeling framework (Fig. 1) include the PFTs of the
 176 sites, the predictors used, the machine learning algorithms, the validation methods, and other features. Each
 177 model for which R-squared is reported is treated as a data record. If multiple algorithms were applied to the
 178 same dataset, then multiple records were extracted. Models using different data or features are also recorded as
 179 multiple records.



180

181 Figure 1. Features of the machine learning-based water flux prediction process. (a) the eddy-covariance-based

182 water flux observations of various plant function types (PFTs), modified from Paul-Limoges et al., 2020. ET,

183 evapotranspiration. E, evaporation. T, transpiration. (b) Predictors and their spatial and temporal resolution. (c)

184 The machine learning algorithms used for the modeling, such as neural networks, random forests, etc. (d) The

185 model validation methods used including the spatial, temporal, and random cross-validations.

186

187

Table 2. Description of information extracted from the included papers.

Field	Definition & Categories adopted	Harmonization
Climate	Climate zones of the study location derived from the Köppen climate classification (Peel et al., 2007)	
Plant functional type (PFT)	PFT of the flux sites: 1-forest, 2-grassland, 3-cropland, 4-wetland, 5-shrubland, 6-savannah, and multi-PFTs	The categorization is based on the descriptions in the article. For example, cropland for various crops is classified

		as ‘cropland’, and both woody savannah and savannah are classified as ‘savannah’.
Location	More precise location (with the latitude and longitude of the center of the studied sites): latitude, longitude	
Algorithms	Random Forests (RF), Multiple Linear Regressions (MLR), Artificial Neural Networks (ANN), Support Vector Machines (SVM), Cubist, model tree ensembles (MTE), K-nearest neighbors (KNN), long short-term memory (LSTM), gradient boosting regression tree (GBRT), extra tree regressor (ETR), Gaussian process regression (GPR), Bayesian model averaging (BMA), extreme learning machine (ELM), and deep belief network (DBN)	Various model algorithms with parameter optimization or other improvements are categorized as their algorithm family. For example, various improved models of RF algorithms are classified as RF, rather than as another algorithm family.
Sites number	Number of the flux sites used	
Spatial scale	Area representatively covered by the flux sites: local (less than 100 x 100 km), regional, global (continent-scale and global scale)	The spatial scale is roughly categorized based on the area covered by the site. The model is classified as ‘global’ only when the spatial extent reaches the continental scale.
Temporal scale	The temporal scale of the model: half-hourly, hourly, daily, 4-daily, 8-daily, monthly, seasonally (i.e., 0.02, 0.04, 1, 4, 8, 30, 90 days)	Models with a temporal scale greater than one month and less than one year are classified as seasonal scale models.
Year span	The span of years of the flux data used	Year span is calculated as the span from the earliest to the latest year of available flux data.
Site year	Describe the volume of total flux data with the number of sites and years aggregated.	
Cross-validation	Describe the chosen method of cross-validation: Spatial (e.g., ‘leave one site out’), temporal (e.g., ‘leave one year out’), random (e.g., ‘k-fold’)	
Training/validation	Describe the ratio of the data volume in the training and validation sets.	In spatial validation, this ratio is represented by the ratio of the number of sites used for training to the number of sites used for validation. In temporal validation, this is represented by the ratio of the span of time periods used for training to the span of time periods used for validation.
Satellite images	Describe the source of satellite images used to derive NDVI, EVI, LAI, LST, etc: Landsat, MODIS, AVHRR	
Biophysical predictors	LAI, NDVI/EVI, the fraction of absorbed photosynthetically active radiation/photosynthetically active radiation (FAPAR/PAR), leaf area index (LAI), Carbon fluxes (CF) including NEE/GPP, etc.	The predictor variables of different measurement methods are categorized according to their definitions. For example, both using the NDVI calculated based on satellite remote

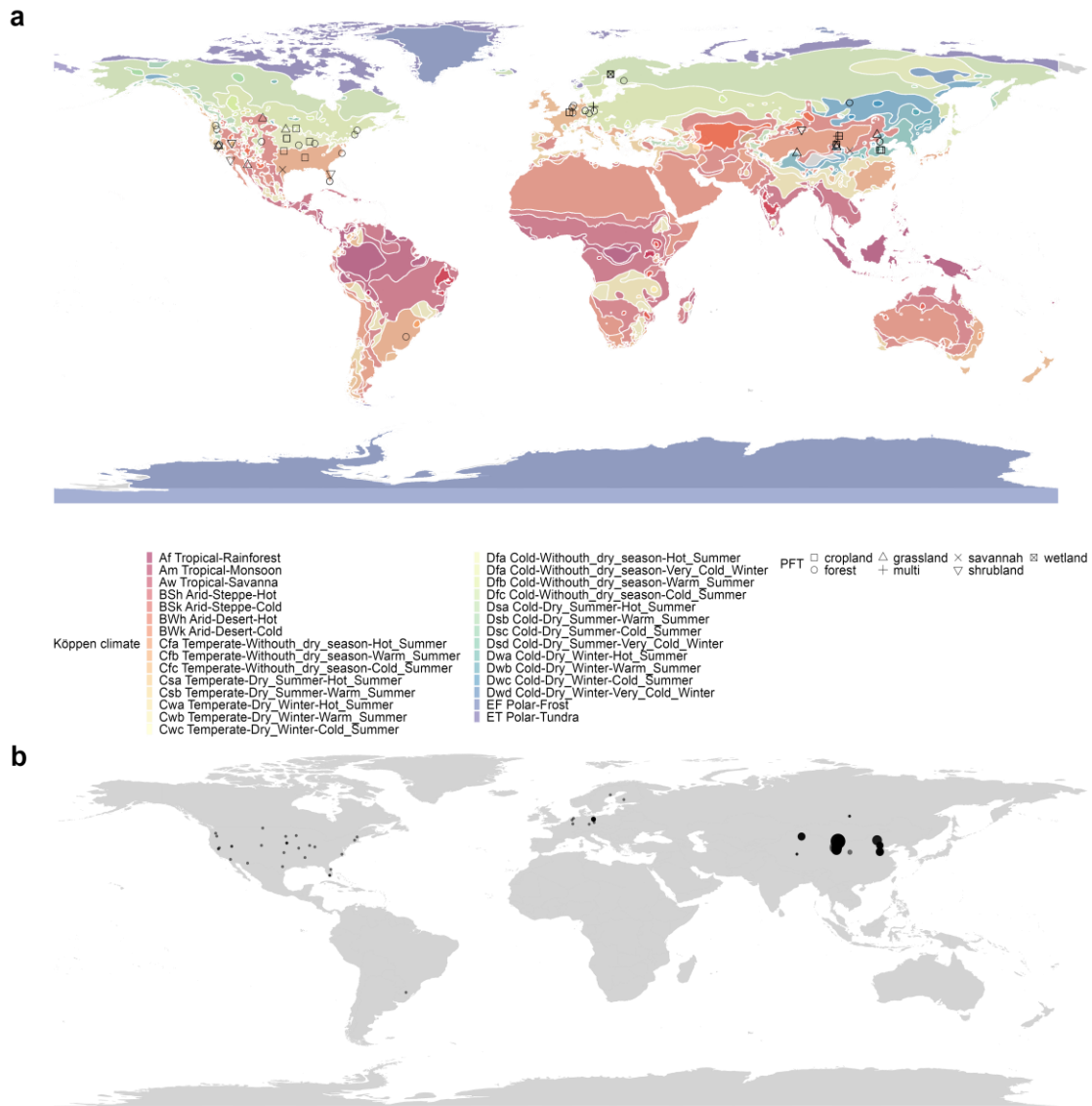
		sensing bands and in situ measurements were classified as the use of 'NDVI'.
Meteorological variables	precipitation (Prec), net radiation/solar radiation (Rn/Rs), air temperature (Ta), vapour-pressure deficit (VPD), relative humidity (RH) , etc.	The way meteorological data are measured is not differentiated. For example, both using Ta from reanalysis data and Ta measured at flux sites were classified as the use of Ta.
Ancillary data	Describe the ancillary variables used: soil texture, terrain (DEM), soil moisture/land surface water index (SM/LSWI), etc.	Both the use of in situ measured soil moisture and the use of remote sensing-based LSWI was classified as using surface moisture-related indicators SM/LSWI.
Accuracy measure	Accuracy measure used to assess the model performance: R-squared (in the validation phase)	

188

189 3 Results

190 3.1 Articles included in the meta-analysis

191 A total of 32 articles (Table S1) containing a total of 139 model records were included. The geographical scope
 192 of these articles was mainly Europe, North America, and China (Fig. 2).



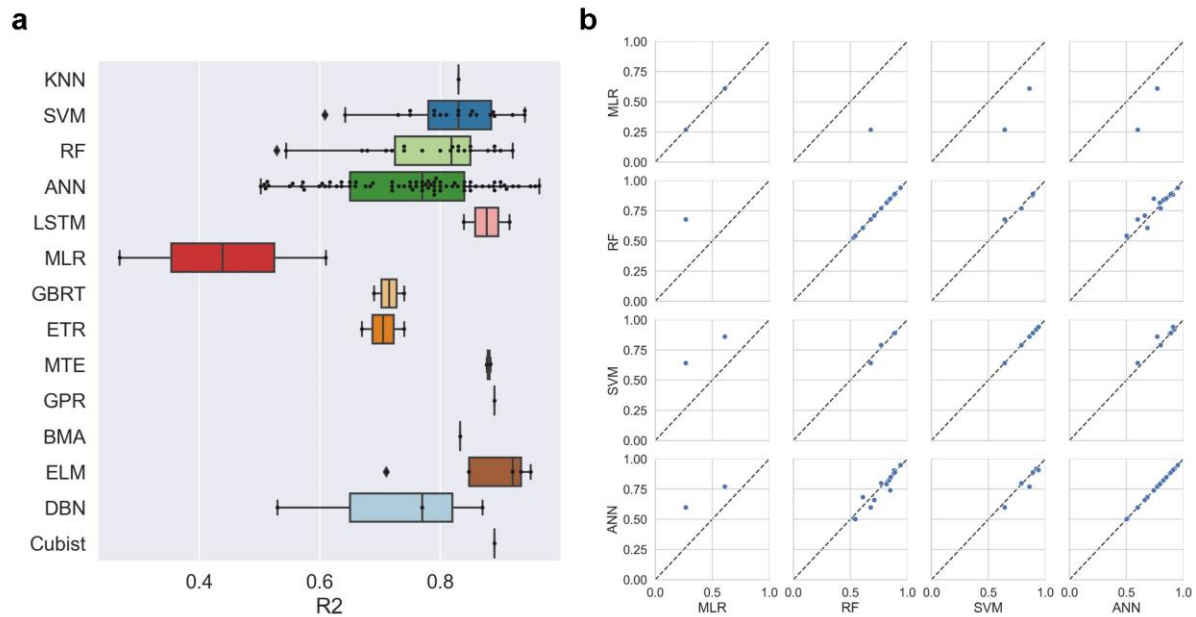
193
 194 Figure 2. Location of the included studies in the meta-analysis. (a) PFTs and the climate zones (from Köppen
 195 climate classification) of these studies and (b) the number of flux sites included in each study. Global and
 196 continental-scale studies (e.g., models developed based on FLUXNET of the global scale) are not shown on the
 197 map due to the difficulty of identifying specific locations.

198 3.2 The formal Meta-analysis

199 3.2.1 Algorithms

200 SVM and RF outperformed (Fig. 3a) across studies (better than other algorithms with sufficient sample size in
 201 Fig. 3a such as ANN). These three machine learning algorithms (i.e., ANN, SVM, RF) were significantly more
 202 accurate than the traditional MLR. Other algorithms such as MTE, ELM, Cubist, etc. also correspond to high
 203 accuracy, but with limited evidence sample size (Fig. 3a). In the internal comparison (different algorithms
 204 applied to the same data set) in single studies, we also find that SVM and RF were slightly more accurate than
 205 ANN (Fig. 3b), and all these three (i.e., ANN, SVM, RF) are considerably more accurate than MLR. Overall,

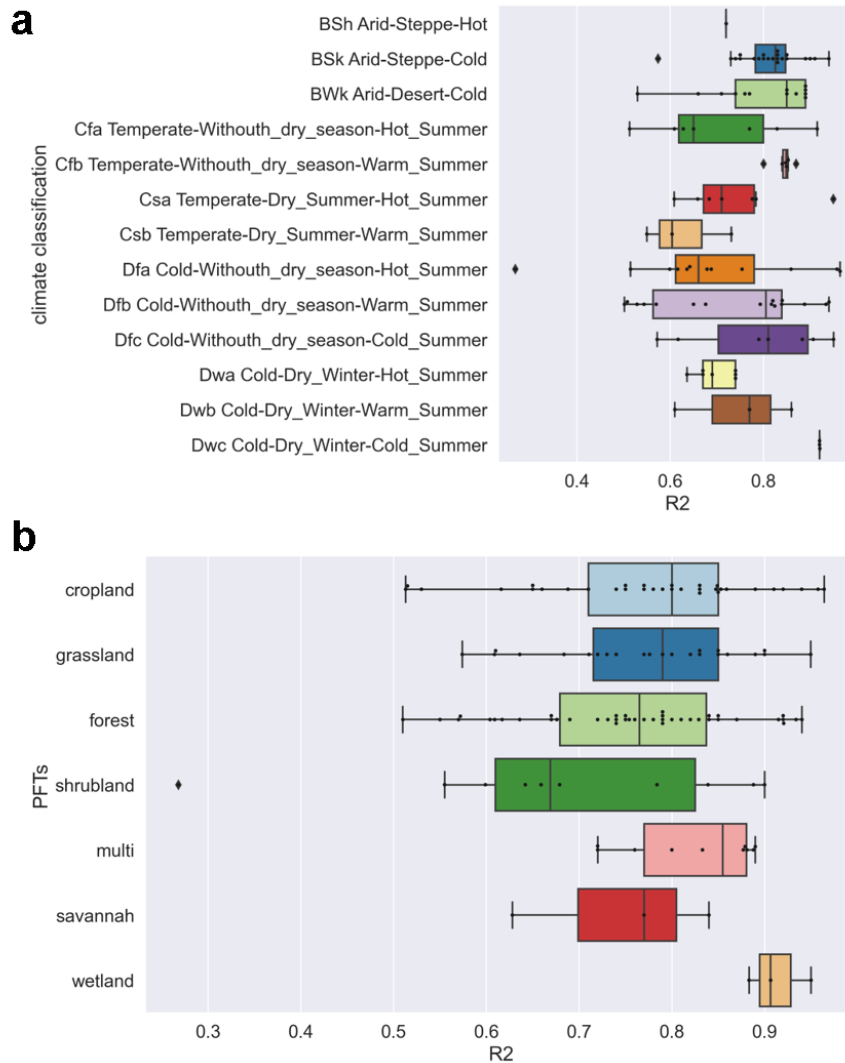
206 SVM and RF have shown higher accuracy in water flux simulations in both inter and intra-study comparisons
 207 with sufficient sample size as evidence.



208
 209 Figure 3. Model accuracy (R-squared) using various algorithms across studies (a) and internal comparisons of
 210 selected pairs of algorithms within studies (b). Algorithms: Random Forests (RF), Multiple Linear Regressions
 211 (MLR), Artificial Neural Networks (ANN), Support Vector Machines (SVM), Bayesian model averaging
 212 (BMA), Cubist, model tree ensembles (MTE), gradient boosting regression tree (GBRT), extra tree regressor
 213 (ETR), K-nearest neighbors (KNN), long short-term memory (LSTM), Gaussian process regression (GPR),
 214 extreme learning machine (ELM), and deep belief network (DBN).

215 3.2.2 Climate types and PFTs

216 We found higher average model accuracy in arid climate zones (Fig. 4a), such as the Cold semi-arid (steppe)
 217 climate (BSk) and Cold desert climate (BWk). Most of these studies were located in northwest China and the
 218 western USA. It may be caused by the simpler relationship between water fluxes and biophysical covariates in
 219 arid regions. In arid zones, due to the high potential ET, the variability in the actual ET may be largely explained
 220 by water availability (moisture supply) and vegetation change with the effect of variability in thermal conditions
 221 reduced. As for the various PFTs, the average model accuracy was slightly lower for forest types than for
 222 cropland and grassland types (Fig. 4b). The lowest average accuracy was found for shrub sites, which may be
 223 related to the difficulty of the remote sensing-based vegetation index (e.g., NDVI) to quantify the physiological
 224 and ecological conditions of shrubs (Zeng et al., 2022), and the heterogeneity of the spatial distribution of
 225 shrubs within the EC observation area may also cause difficulties in capturing their relationships with
 226 biophysical variables. We also found high model accuracy for the wetland type, although records as evidence to
 227 support this finding may be limited. Compared to other PFTs, the more steady and adequate water availability in
 228 the wetland type may make the variations of water fluxes less explained by other biophysical covariates.



229

230 Figure 4. Differences in model accuracy (R-squared) of (a) various climate zones (classified by Köppen climate
 231 classification) across studies and (b) PFTs. BSh, Hot semi-arid (steppe) climate. BSk, Cold semi-arid (steppe)
 232 climate. BWk, Cold desert climate. Cfa, Humid subtropical climate. Cfb, Temperate oceanic climate. Csa, Hot-
 233 summer Mediterranean climate. Csb, Warm-summer Mediterranean climate. Dfa, Hot-summer humid
 234 continental climate. Dfb, Warm-summer humid continental climate. Dfc, Subarctic climate. Dwa, Monsoon-
 235 influenced hot-summer humid continental climate. Dwb, Monsoon-influenced warm-summer humid continental
 236 climate. Dwc, Monsoon-influenced subarctic climate.

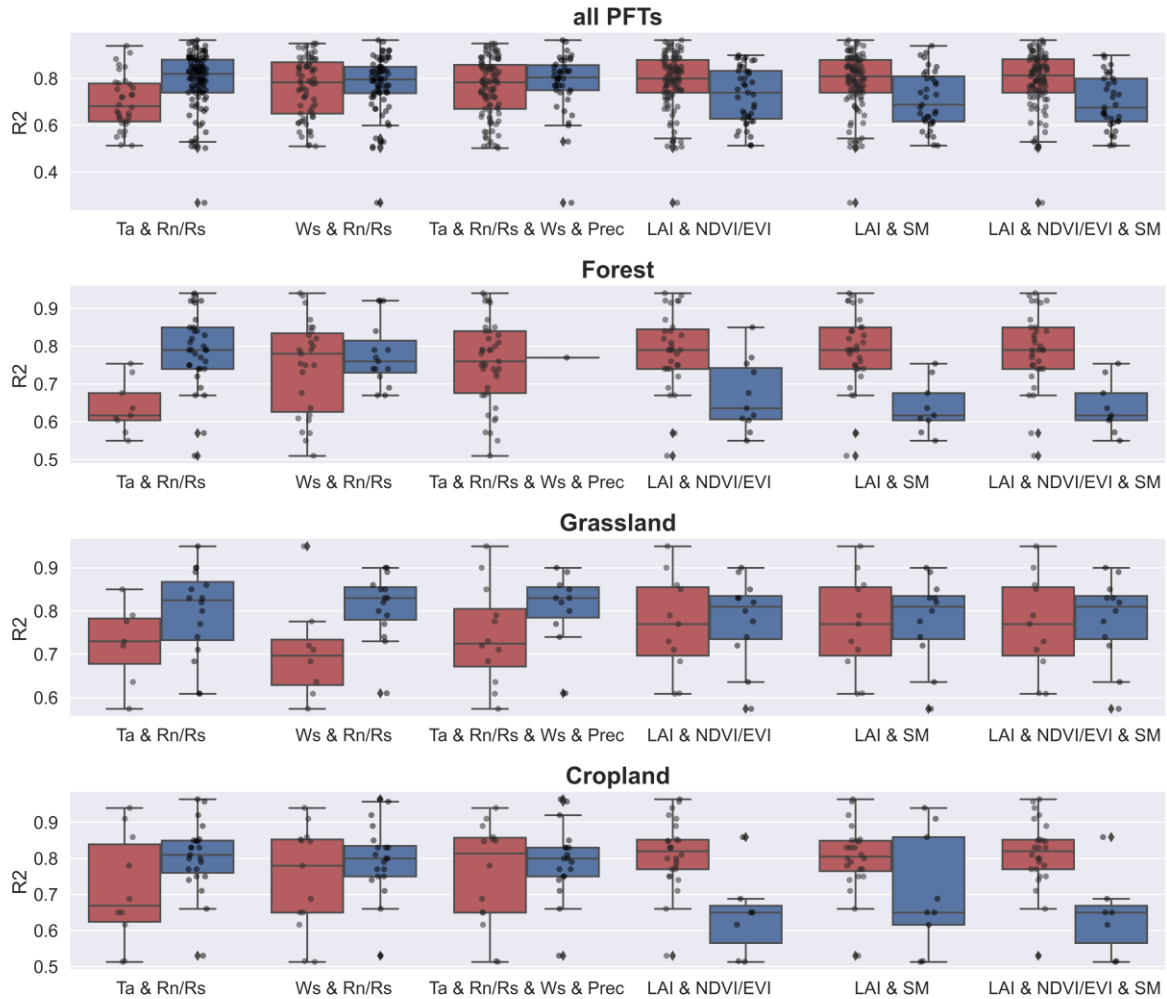
237 **3.3.3 Predictors and their combinations**

238 On one hand, for the effects of individual predictors, the use of Rn/Rs, Prec, Ta, and FAPAR improved the
 239 accuracy of the model (Fig. S1). This pattern partially changed in the different PFTs. In the forest sites, the
 240 accuracy of the models with Rn/Rs and Ta used was higher than that of the models with Rn/Rs and Ta not used.
 241 For the grassland sites, the use of Ws, FAPAR, Prec, and Rn/Rs improved the model accuracy. For the cropland
 242 sites, Ta and FAPAR were more important for improving the model accuracy.

243

244 On the other hand, the evaluation of the effect of individual predictors on model accuracy is not necessarily
245 reliable because some predictor variables are used together (e.g., the high model accuracy corresponding to a
246 particular variable may be because it is often used together with another variable that plays the dominant role in
247 improving accuracy). Therefore, we tested for independence between the use of variables and assessed the effect
248 of the combination of variables on model accuracy. We calculated the correlation matrix (Fig. S2) between the
249 use of various predictors (not used is set as 0 and used is set as 1). We found there was a dependence between
250 the use of some predictors, the use of NDVI/EVI, LAI, and SM was significantly negatively correlated with the
251 use of Rn/Rs and Ta (Fig. S2). It indicated that many of the models that used Rn/Rs and Ta did not use
252 NDVI/EVI, LAI, and SM, and the models that used NDVI/EVI, LAI, and SM also happened to not use Rn/Rs
253 and Ta. Given this dependence, we evaluated the effect of the combination of variables on the model accuracy
254 (Fig. 5). In Fig. 5, the three variable combinations on the left side are mainly meteorological variables while the
255 three variable combinations on the right side are mainly vegetation-related variables based on remote sensing
256 (e.g., NDVI, EVI, LAI, LSWI). We found that, overall, the accuracy of the models using only meteorological
257 variable combinations was higher than that of the models using only remote sensing-based vegetation-related
258 variables. It demonstrated the importance of using meteorological variables in machine learning-based ET
259 prediction (probably especially for models with small time scales such as hourly scale, and daily scale). For
260 example, in the forest type, the combination of Ta and Rn/Rs is very effective compared to using only remote
261 sensing-based vegetation index variable combinations. The combination of Ta and Rn/Rs is also effective in the
262 grassland and cropland types. The combination of Ws and Rn/Rs played an important role in the grassland type
263 for improving model accuracy. Despite this, it does not negate the positive role of remote sensing-based
264 vegetation-related variables in ET prediction. This effectiveness can be dependent on the time scale of the model
265 as well as the PFTs. In models with large time scales (monthly scale, seasonal scale) and PFTs in which ET is
266 sensitive to vegetation dynamics, remote sensing-based vegetation-related variables may also be of high
267 importance.

268



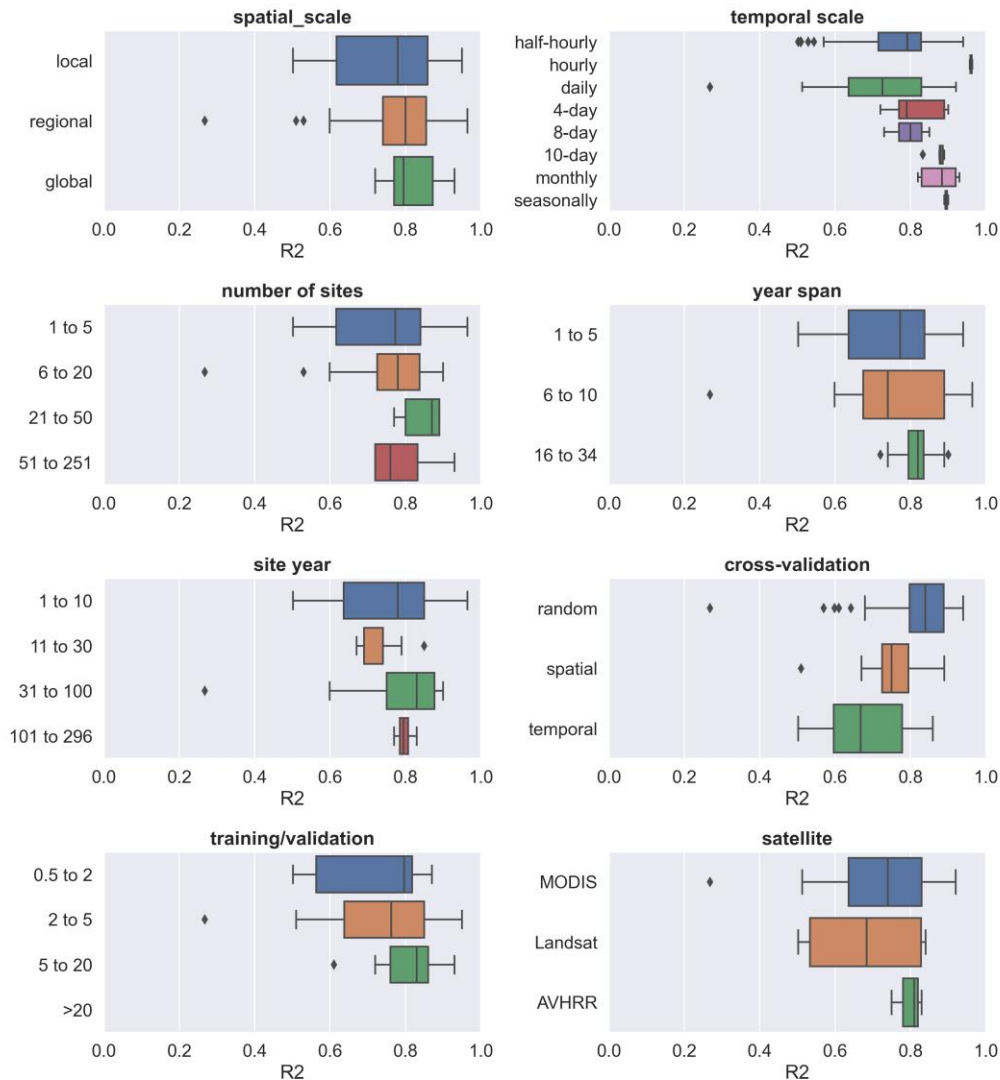
269

270 Figure 5. Effects of combinations of predictor variables on model accuracy in various PFTs (all data, forest,
 271 grassland, and cropland). Dark blue boxes indicate that the predictors were together used in the model (e.g., for
 272 ‘Ta & Rn/Rs’, the dark blue box represents Ta and Rn/Rs were together used in the model), while dark red
 273 boxes indicate the other conditions (i.e., the combination was not used). Predictors: precipitation (Prec), soil
 274 moisture/remote sensing-based land surface water index (SM), net radiation/solar radiation (Rn/Rs), enhanced
 275 vegetation index (EVI), air temperature (Ta), leaf area index (LAI), Normalized Difference Vegetation
 276 Index/Enhanced Vegetation Index (NDVI/EVI).

277 **3.3.4 Other model features**

278 We also evaluated the impact of some other features on accuracy. The differences in accuracy of models with
 279 different spatial scales, year spans, number of sites, and volume of data (Fig. 6) appear to be insignificant. This
 280 seems to be related to the fact that in large-scale water flux simulations, the sites of similar PFTs are selected
 281 such as for modeling multiple forest sites across Europe (Van Wijk and Bouten, 1999) which focus on ‘forest’
 282 and multiple grassland sites across arid northern China (Xie et al., 2021; Zhang et al., 2021) which focus on
 283 ‘grassland’, rather than mixing different PFT types to train models as the way in machine learning modeling of
 284 carbon fluxes (Zeng et al., 2020). In terms of the time scales of the models, the 4-day, 8-day, and monthly scales
 285 appear to correspond to higher accuracy compared to the half-hourly and daily scales. The higher the ratio of the
 286 volume of data in the training and validation sets, the higher the model accuracy. Compared to the models using

287 Landsat data, the models using MODIS data showed slightly higher accuracy probably due to the advantage of
 288 MODIS data in capturing the temporal dynamics of biophysical covariates. There were significant differences in
 289 the accuracy of the models using different cross-validation methods, with the models using random cross-
 290 validation showing higher accuracy than those using temporal cross-validation. This suggests that interannual
 291 variability may have a high impact on the models in water flux simulations. The driving mechanism of ET may
 292 vary significantly across years, and the inclusion of some extreme climatic conditions in the training set may be
 293 important for model accuracy and robustness.
 294

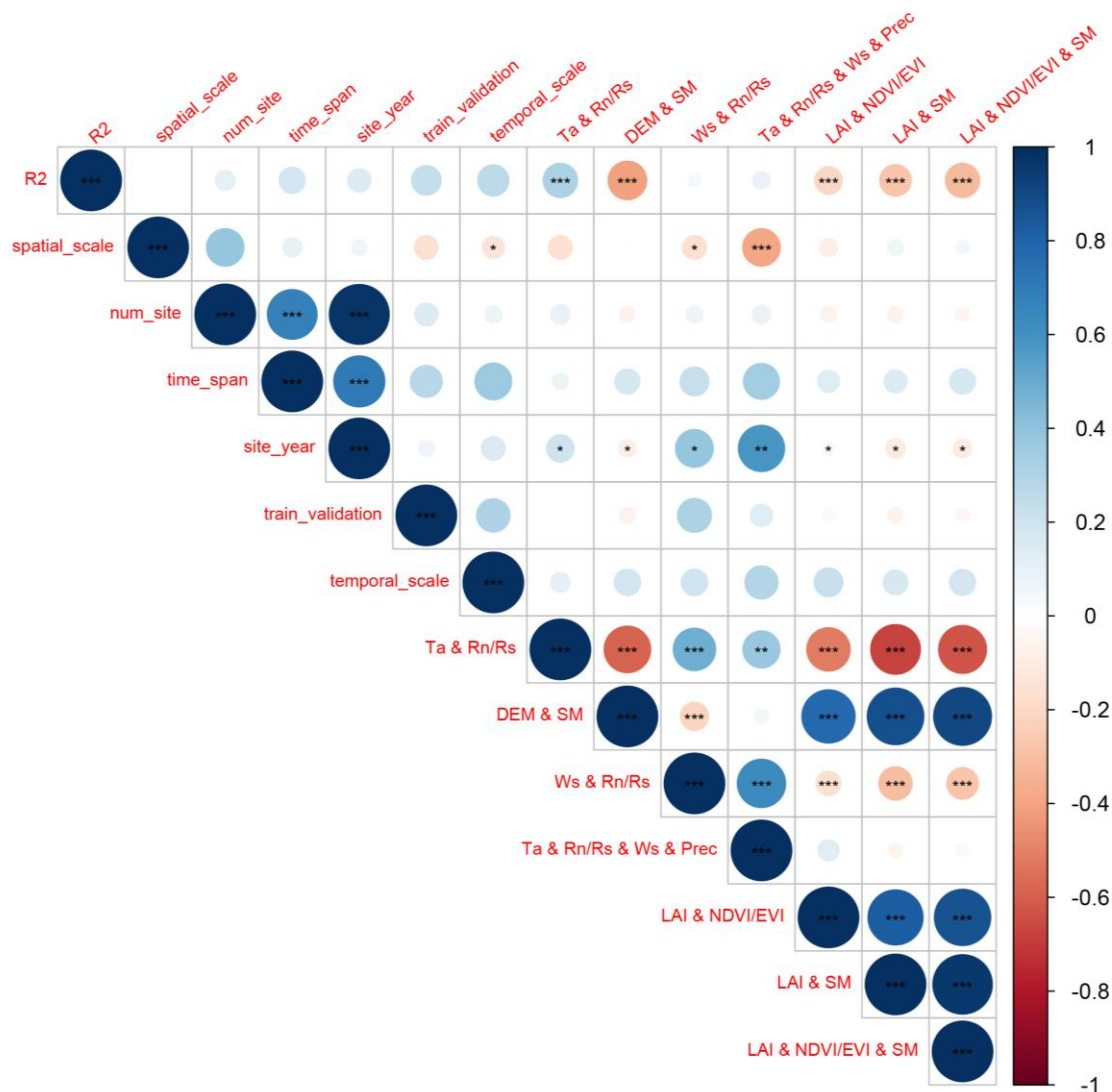


295
 296 Figure 6. The effects of other model features (i.e. spatial scale, number of sites, temporal scale, year span, site
 297 year, validation method, training/validation ratio, and satellite imagery used) on the R-squared.

298 3.3.5 Linear correlation of quantitative features and R-squared

299 We also analyzed the linear correlation (Fig. 7) between multiple quantitative features and the R-squared. We
 300 found that the magnitude of the linear correlation coefficients between the use of predictor combinations and the
 301 R-squared was higher than other features. The use of the predictor combination ‘Ta and Rn/Rs’ significantly
 302 improved the model accuracy. ‘Temporal scale’, ‘time span’, ‘training/validation ratio’, and ‘number of sites’

303 showed weak positive correlations with R-squared (not significant, p-value > 0.1). The positive correlation
 304 between 'temporal scale' and R-squared is higher among these features, although not significant. It should also
 305 be paid more attention to in future studies. The feature 'training/validation ratio' and 'time span' are also
 306 positively correlated (although not significantly) with the R-squared, suggesting the importance of the volume of
 307 data in the training set in a data-driven machine learning model. Larger 'training/validation ratio' and 'time span'
 308 may correspond to greater proportional coverage of the scenarios/conditions in the training set over the
 309 validation set, and thus correspond to higher accuracy.
 310



311
 312 Figure 7. Evaluation of linear correlations between multiple features and the R-squared records with the
 313 statistical significance test. For the feature 'spatial scale', the 'local' scale was set to 1, the 'regional' scale was
 314 set to 2, and the 'global' scale was set to 3 in the analysis of linear correlation. For the use of various predictor
 315 combinations with '&', the value for 'together used' is set as 1 and other conditions are set as 0 (e.g., for the
 316 feature 'Ta & Rn/Rs & Ws & Prec', if Ta, Rn/Rs, Ws, and Prec were used together in the model, the value is set
 317 as 1). Significance: the p-value < 0.01 (***), 0.05 (**), and 0.1 (*).

318 **4 Discussions**

319 With the accumulation of in situ EC observations around the world, the study of ET simulations based on data-
320 driven approaches has received more attention from researchers in the last decade. Many studies have combined
321 EC observations, various predictors, and machine learning algorithms to improve the prediction accuracy of
322 site-scale water fluxes. To date, the results of these studies have not been comprehensively evaluated to provide
323 clear guidance for feature selection in water flux prediction models. To better understand the approach and
324 guide future research, we performed a meta-analysis of such studies. Machine learning-based water flux
325 simulations and predictions still suffer from high uncertainty. By investigating the expected improvements that
326 can be achieved by incorporating different features, we can avoid practices that may reduce model accuracy in
327 future research.

328 **4.1 Opportunities and challenges in the site-scale water flux simulation**

329 In the above meta-analysis of the models, we found that water flux simulations based on EC observations can
330 achieve high accuracy but also have high uncertainty through the modeling workflow. The R-squared of many
331 water flux simulation models exceeds 0.8, possibly higher than some remote sensing-based and process-based
332 models, and possibly higher than carbon flux simulations such as the net ecosystem exchange (Shi et al., 2022)
333 in the same modeling framework.

334
335 Biophysical and meteorological variables are considered both important in ET simulations. This study found
336 that models using a combination of meteorological variables had higher accuracy than models using only
337 remotely sensed vegetation dynamic information. However, due to the high proportion of models with small
338 temporal scales (e.g., half-hourly scale, hourly scale, and daily scale) in this study, this advantage of the
339 combination of meteorological variables may be more suitable for small temporal scales. A possible explanation
340 is that vegetation-related variables such as NDVI and LAI at the daily scale, 8-day scale, and 16-day scale have
341 limited explanatory ability for hourly or daily-scale variability in ET. At a small temporal scale, the use of
342 combinations of meteorological variables can capture moisture and energy conditions that control the rapid
343 fluctuations of ET and thus has a dominant role in hourly or daily-scale ET prediction. This also corroborates
344 the high accuracy of some physic-based ET estimation models (Rigden and Salvucci, 2015) that use only
345 meteorological variables and not vegetation-related variables such NDVI (only an estimate of vegetation height
346 derived from land cover maps is used to represent vegetation conditions (Rigden and Salvucci, 2015)).

347
348 There are differences in model accuracy among different PFTs. For example, in forest sites, limitations in data
349 accuracy of factors were possible because some remote sensing-based predictors such as NDVI, FAPAR, and
350 LAI have limited accuracy when applied to forest types (Liu et al., 2018b; Zeng et al., 2022). In addition, factors
351 such as crown density, which may significantly affect the proportion of soil evaporation, transpiration, and
352 evaporation of canopy interception, were not considered in these models, which may also lead to low model
353 accuracy. This suggests that in water flux simulation, the driving mechanisms of water fluxes in different PFTs
354 do affect the accuracy of machine learning models, and we need to consider more the actual and specific
355 influencing factors in specific PFTs. More variables that can quantify the ratio of evaporation and transpiration

356 should be considered for inclusion, which also appears to improve the mechanistic interpretability of such
357 machine learning models. A previous study (Zhao et al., 2019) has combined the physics-based approach (e.g.,
358 Penman-Monteith equation) and machine learning to build hybrid models to improve interpretability. We should
359 make full use of empirical knowledge and experiences from process-based models to improve the accuracy and
360 interpretability of the machine learning approach.

361

362 Among the validation methods, random cross-validation has higher accuracy than spatial cross-validation and
363 temporal cross-validation. However, spatial cross-validation and temporal cross-validation may be able to better
364 help us recognize the robustness of the model when extrapolated (i.e., applied to new stations and new years).
365 The lower accuracy in the temporal cross-validation approach implies that we need to focus on interannual
366 hydrological and meteorological variability in the water flux simulations. In cropland sites, we may also need to
367 pay more attention to the effects of interannual variability in anthropogenic cropping patterns. If some extreme
368 weather years are not included, the robustness of the model when extrapolated to other years may be challenged,
369 especially in the context of the various extreme weather events of recent years. This can also inform the siting of
370 future flux stations. Regions where climate extremes may occur and biogeographic types not covered by
371 existing flux observation networks should be given more attention to achieve global-scale, accurate and robust
372 machine learning-based spatio-temporal prediction of water fluxes.

373 **4.2 Uncertainties and limitations of this meta-analysis**

374 **4.2.1 The limited number of available literature and model records**

375 Despite many articles and model records collected through our efforts to perform this meta-analysis, there still
376 appears to be a long way to go to finally and completely understand the various mechanisms involved in water
377 flux simulation with machine learning. Some of the insights provided by this study can be not robust (due to the
378 limited sample size available when the goal is to assess the effects of multiple features), but this does not negate
379 the fact that this study does obtain some meaningful findings. Therefore, researchers should treat the results of
380 this study with caution, as they were obtained only statistically. Overall, it is still positive to conduct a meta-
381 analysis of such studies, considering their rapid growth in the number and lack of guiding directions.

382 **4.2.2 Publication bias and weighting**

383 Publication bias and weighting: Due to the relatively limited number of articles that could be included in the
384 meta-analysis, this study did not focus much on publication bias. Meta-analytic studies in other fields typically
385 measure the quality of journals and the public availability of research data (Borenstein et al., 2011; Field and
386 Gillett, 2010) to determine the weighting of the literature in a comprehensive assessment. However, most of the
387 articles did not publicly provide flux observations or share developed models. Meta-analysis studies in other
388 fields typically measure the impact of included studies based on sample size and variance of experimental
389 results (Adams et al., 1997; Don et al., 2011; Liu et al., 2018a). In this study, due to the lack of a convincing
390 manner to determine weights among articles, we assigned the same weight to the results for all the literature.

391 **4.2.3 Uncertainties in the information of the extracted features**

392 At the information extraction level, the following issues may also introduce uncertainties:

- 393 a) Uncertainties caused by data quality control (e.g. gap-filling (Hui et al., 2004)) are difficult to assess
394 effectively. Gap-filling is a commonly used technique to fill in low-quality data in flux observations (Chen
395 et al., 2012; Hui et al., 2004). However, the impact of this practice on machine learning-based ET
396 prediction models is unclear, due to the difficulty of directly assessing how this technique is performed in
397 various studies by this meta-analysis. Typically, models with small time scales (e.g., hourly scale, daily
398 scale) can exclude low-quality observations and use only high-quality data. However, for models with
399 large time scales (e.g., monthly scales), gap-filling (e.g., based on meteorological data) may be
400 unavoidable. This may lead to decrease in training data purity and introduce uncertainty in the subsequent
401 prediction model development.
- 402 b) Systematic uncertainties caused by the energy balance closure (EBC) issue in eddy-covariance flux
403 measurements are also difficult to assess by this meta-analysis. EBC is a common problem (Eshonkulov et
404 al., 2019) in eddy-covariance flux observations. For that reason, the latent heat flux measured potentially
405 underestimates ET. Some prediction models corrected EBC (e.g., using Bowen ratio preserving (Mauder et
406 al., 2013, 2018) and energy balance residuals (Charuchittipan et al., 2014; Mauder et al., 2018)) in the
407 processing of training data, but some did not. How this will affect the accuracy of the prediction model is
408 not clear due to multiple factors that need to be evaluated that influence EBC (Foken, 2008), including
409 measurement errors of the energy balance components, incorrect sensor configurations, influences of
410 heterogeneous canopy height, unconsidered energy storage terms in the soil-plant-atmosphere system,
411 inadequate time averaging intervals, and long-wave eddies (Jacobs et al., 2008; Foken, 2008; Eshonkulov
412 et al., 2019). To reduce this uncertainty, more attention to flux site characteristics (Eshonkulov et al., 2019)
413 related to PFT, topography, flux footprint area, etc., to select the appropriate correction method is
414 necessary for future studies.
- 415 c) As most studies used far more water flux observation records than the number of covariates in their
416 regression models, we did not adjust the R-squared in this study to an adjusted R-squared.
- 417 d) The various specific ways in which the parameters of the model are optimized are not differentiated. They
418 are broadly categorized into different families or kinds of algorithms, which may also introduce uncertainty
419 into the assessment.
- 420 e) The assessment of some features is not detailed due to the limitations of the available model records. For
421 example, the classification of PFT could be more detailed. ‘Forest’ could be further classified as broadleaf
422 forest, coniferous forest, etc. while ‘cropland’ could be further classified as rainfed and irrigated cropland
423 based on differences in their response mechanisms of water fluxes to environmental factors.

424 **5 Conclusion**

425 We performed a meta-analysis of the site-scale water flux simulations combining in situ flux observations,
426 meteorological, biophysical, and ancillary predictors, and machine learning. The main conclusions are as
427 follows:

- 428 1. SVM (average R-squared = 0.82) and RF (average R-squared = 0.81) outperformed over evaluated
429 algorithms with sufficient sample size in both cross-study and intra-study (with the same training dataset)
430 comparisons.

- 431 2. The average accuracy of the model applied to arid regions is higher than in other climate types.
- 432 3. The average accuracy of the model was slightly lower for forest sites (average R-squared = 0.76) than for
- 433 cropland and grassland sites (average R-squared = 0.8 and 0.79), but higher than for shrub sites (average R-
- 434 squared = 0.67).
- 435 4. Among various predictor variables, the use of Rn/Rs, Prec, Ta, and FAPAR improved the model accuracy.
- 436 The combination of Ta and Rn/Rs is very effective especially in the forest type, while in the grassland type
- 437 the combination of Ws and Rn/Rs is also effective.
- 438 5. Among the different validation methods, random cross-validation shows higher model accuracy than spatial
- 439 cross-validation and temporal cross-validation.
- 440
- 441

442 **Acknowledgements**

443 We thank the editor and two anonymous reviewers for their insightful comments which contributed substantially
444 to the improvement of this manuscript.

445 **Financial support**

446 This research was supported by the National Natural Science Foundation of China (Grant No. U1803243), the
447 Key projects of the Natural Science Foundation of Xinjiang Autonomous Region (Grant No. 2022D01D01), the
448 Strategic Priority Research Program of the Chinese Academy of Sciences (Grant No. XDA20060302), and High-
449 End Foreign Experts Project.

450 **Author Contributions**

451 HS and GL were responsible for the conceptualization, methodology, formal analysis, investigation, visualization,
452 and writing. OH contributed to the investigation. XM, XY, YW, WZ, MX, CZ and YZ processed the data. AK,
453 TVDV and PDM provided supervision.

454 **Competing interests**

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

456 **Code availability**

457 The codes that were used for all analyses are available from the first author (shihaiyang16@mails.ucas.ac.cn)
458 upon request.

459 **Data availability**

460 The data used in this study can be accessed by contacting the first author (shihaiyang16@mails.ucas.ac.cn) upon
461 request.

462

463 **References**

- 464 Adams, D. C., Gurevitch, J., and Rosenberg, M. S.: Resampling tests for meta - analysis of ecological
465 data, *Ecology*, 78, 1277–1283, 1997.
- 466 Allen, R. G., Pereira, L. S., Howell, T. A., and Jensen, M. E.: Evapotranspiration information
467 reporting: I. Factors governing measurement accuracy, *Agricultural Water Management*, 98, 899–920,
468 <https://doi.org/10.1016/j.agwat.2010.12.015>, 2011.
- 469 Anderson, M. C., Allen, R. G., Morse, A., and Kustas, W. P.: Use of Landsat thermal imagery in
470 monitoring evapotranspiration and managing water resources, *Remote Sensing of Environment*, 122,
471 50–65, <https://doi.org/10.1016/j.rse.2011.08.025>, 2012.
- 472 Barman, R., Jain, A. K., and Liang, M.: Climate-driven uncertainties in modeling terrestrial energy
473 and water fluxes: a site-level to global-scale analysis, 20, 1885–1900,
474 <https://doi.org/10.1111/gcb.12473>, 2014.
- 475 Borenstein, M., Hedges, L. V., Higgins, J. P., and Rothstein, H. R.: *Introduction to meta-analysis*,
476 John Wiley & Sons, 2011.
- 477 Charuchittipan, D., Babel, W., Mauder, M., Leps, J.-P., and Foken, T.: Extension of the Averaging
478 Time in Eddy-Covariance Measurements and Its Effect on the Energy Balance Closure, *Boundary-
479 Layer Meteorol*, 152, 303–327, <https://doi.org/10.1007/s10546-014-9922-6>, 2014.
- 480 Chen, Y., Xia, J., Liang, S., Feng, J., Fisher, J. B., Li, X., Li, X., Liu, S., Ma, Z., Miyata, A., Mu, Q.,
481 Sun, L., Tang, J., Wang, K., Wen, J., Xue, Y., Yu, G., Zha, T., Zhang, L., Zhang, Q., Zhao, T., Zhao,
482 L., and Yuan, W.: Comparison of satellite-based evapotranspiration models over terrestrial
483 ecosystems in China, *Remote Sensing of Environment*, 140, 279–293,
484 <https://doi.org/10.1016/j.rse.2013.08.045>, 2014.
- 485 Chen, Y., Wang, S., Ren, Z., Huang, J., Wang, X., Liu, S., Deng, H., and Lin, W.: Increased
486 evapotranspiration from land cover changes intensified water crisis in an arid river basin in northwest
487 China, *Journal of Hydrology*, 574, 383–397, <https://doi.org/10.1016/j.jhydrol.2019.04.045>, 2019.
- 488 Chen, Y.-Y., Chu, C.-R., and Li, M.-H.: A gap-filling model for eddy covariance latent heat flux:
489 Estimating evapotranspiration of a subtropical seasonal evergreen broad-leaved forest as an example,
490 468–469, 101–110, <https://doi.org/10.1016/j.jhydrol.2012.08.026>, 2012.
- 491 Don, A., Schumacher, J., and Freibauer, A.: Impact of tropical land-use change on soil organic carbon
492 stocks – a meta-analysis, 17, 1658–1670, <https://doi.org/10.1111/j.1365-2486.2010.02336.x>, 2011.
- 493 Eshonkulov, R., Poyda, A., Ingwersen, J., Wize mann, H.-D., Weber, T. K. D., Kremer, P., Högy, P.,
494 Pulatov, A., and Streck, T.: Evaluating multi-year, multi-site data on the energy balance closure of
495 eddy-covariance flux measurements at cropland sites in southwestern Germany, 16, 521–540,
496 <https://doi.org/10.5194/bg-16-521-2019>, 2019.
- 497 Fang, B., Lei, H., Zhang, Y., Quan, Q., and Yang, D.: Spatio-temporal patterns of evapotranspiration
498 based on upscaling eddy covariance measurements in the dryland of the North China Plain, 281,
499 <https://doi.org/10.1016/j.agrformet.2019.107844>, 2020.
- 500 Field, A. P. and Gillett, R.: How to do a meta - analysis, *British Journal of Mathematical and
501 Statistical Psychology*, 63, 665–694, 2010.
- 502 Fisher, J. B., Melton, F., Middleton, E., Hain, C., Anderson, M., Allen, R., McCabe, M. F., Hook, S.,
503 Baldocchi, D., Townsend, P. A., Kilic, A., Tu, K., Miralles, D. D., Perret, J., Lagouarde, J.-P.,

504 Waliser, D., Purdy, A. J., French, A., Schimel, D., Famiglietti, J. S., Stephens, G., and Wood, E. F.:
505 The future of evapotranspiration: Global requirements for ecosystem functioning, carbon and climate
506 feedbacks, agricultural management, and water resources, 53, 2618–2626,
507 <https://doi.org/10.1002/2016WR020175>, 2017.

508 Foken, T.: The energy balance closure problem: An overview, *Ecological Applications*, 18, 1351–
509 1367, 2008.

510 Gaston, K. J.: Global patterns in biodiversity, 405, 220–227, <https://doi.org/10.1038/35012228>, 2000.

511 Hui, D., Wan, S., Su, B., Katul, G., Monson, R., and Luo, Y.: Gap-filling missing data in eddy
512 covariance measurements using multiple imputation (MI) for annual estimations, 121, 93–111,
513 [https://doi.org/10.1016/S0168-1923\(03\)00158-8](https://doi.org/10.1016/S0168-1923(03)00158-8), 2004.

514 Jacobs, A. F. G., Heusinkveld, B. G., and Holtslag, A. A. M.: Towards Closing the Surface Energy
515 Budget of a Mid-latitude Grassland, *Boundary-Layer Meteorol*, 126, 125–136,
516 <https://doi.org/10.1007/s10546-007-9209-2>, 2008.

517 Jung, M., Reichstein, M., and Bondeau, A.: Towards global empirical upscaling of FLUXNET eddy
518 covariance observations: Validation of a model tree ensemble approach using a biosphere model, 6,
519 2001–2013, <https://doi.org/10.5194/bg-6-2001-2009>, 2009.

520 Jung, M., Reichstein, M., Margolis, H. A., Cescatti, A., Richardson, A. D., Arain, M. A., Arneth, A.,
521 Bernhofer, C., Bonal, D., Chen, J., Gianelle, D., Gobron, N., Kiely, G., Kutsch, W., Lasslop, G., Law,
522 B. E., Lindroth, A., Merbold, L., Montagnani, L., Moors, E. J., Papale, D., Sottocornola, M., Vaccari,
523 F., and Williams, C.: Global patterns of land-atmosphere fluxes of carbon dioxide, latent heat, and
524 sensible heat derived from eddy covariance, satellite, and meteorological observations, 116,
525 <https://doi.org/10.1029/2010JG001566>, 2011.

526 Kaur, H., Pannu, H. S., and Malhi, A. K.: A Systematic Review on Imbalanced Data Challenges in
527 Machine Learning: Applications and Solutions, *ACM Comput. Surv.*, 52, 79:1-79:36,
528 <https://doi.org/10.1145/3343440>, 2019.

529 Li, X., He, Y., Zeng, Z., Lian, X., Wang, X., Du, M., Jia, G., Li, Y., Ma, Y., Tang, Y., Wang, W., Wu,
530 Z., Yan, J., Yao, Y., Ciais, P., Zhang, X., Zhang, Y., Zhang, Y., Zhou, G., and Piao, S.:
531 Spatiotemporal pattern of terrestrial evapotranspiration in China during the past thirty years, 259,
532 131–140, <https://doi.org/10.1016/j.agrformet.2018.04.020>, 2018.

533 Li, X., Kang, S., Niu, J., Huo, Z., and Liu, J.: Improving the representation of stomatal responses to
534 CO₂ within the Penman–Monteith model to better estimate evapotranspiration responses to climate
535 change, *Journal of Hydrology*, 572, 692–705, <https://doi.org/10.1016/j.jhydrol.2019.03.029>, 2019.

536 Liu, Q., Zhang, Y., Liu, B., Amonette, J. E., Lin, Z., Liu, G., Ambus, P., and Xie, Z.: How does
537 biochar influence soil N cycle? A meta-analysis, *Plant and soil*, 426, 211–225, 2018a.

538 Liu, Y., Xiao, J., Ju, W., Zhu, G., Wu, X., Fan, W., Li, D., and Zhou, Y.: Satellite-derived LAI
539 products exhibit large discrepancies and can lead to substantial uncertainty in simulated carbon and
540 water fluxes, *Remote Sensing of Environment*, 206, 174–188,
541 <https://doi.org/10.1016/j.rse.2017.12.024>, 2018b.

542 Lu, X. and Zhuang, Q.: Evaluating evapotranspiration and water-use efficiency of terrestrial
543 ecosystems in the conterminous United States using MODIS and AmeriFlux data,
544 <https://doi.org/10.1016/j.rse.2010.04.001>, 2010.

- 545 Mauder, M., Cuntz, M., Drüe, C., Graf, A., Rebmann, C., Schmid, H. P., Schmidt, M., and
546 Steinbrecher, R.: A strategy for quality and uncertainty assessment of long-term eddy-covariance
547 measurements, *Agricultural and Forest Meteorology*, 169, 122–135,
548 <https://doi.org/10.1016/j.agrformet.2012.09.006>, 2013.
- 549 Mauder, M., Genzel, S., Fu, J., Kiese, R., Soltani, M., Steinbrecher, R., Zeeman, M., Banerjee, T., De
550 Roo, F., and Kunstmann, H.: Evaluation of energy balance closure adjustment methods by
551 independent evapotranspiration estimates from lysimeters and hydrological simulations, 32, 39–50,
552 <https://doi.org/10.1002/hyp.11397>, 2018.
- 553 McColl, K. A.: Practical and Theoretical Benefits of an Alternative to the Penman-Monteith
554 Evapotranspiration Equation, 56, e2020WR027106, <https://doi.org/10.1029/2020WR027106>, 2020.
- 555 Minacapilli, M., Agnese, C., Blanda, F., Cammalleri, C., Ciraolo, G., D’Urso, G., Iovino, M., Pumo,
556 D., Provenzano, G., and Rallo, G.: Estimation of actual evapotranspiration of Mediterranean perennial
557 crops by means of remote-sensing based surface energy balance models, 13, 1061–1074,
558 <https://doi.org/10.5194/hess-13-1061-2009>, 2009.
- 559 Miralles, D. G., Holmes, T. R. H., De Jeu, R. a. M., Gash, J. H., Meesters, A. G. C. A., and Dolman,
560 A. J.: Global land-surface evaporation estimated from satellite-based observations, 15, 453–469,
561 <https://doi.org/10.5194/hess-15-453-2011>, 2011.
- 562 Miralles, D. G., Teuling, A. J., van Heerwaarden, C. C., and Vilà-Guerau de Arellano, J.: Mega-
563 heatwave temperatures due to combined soil desiccation and atmospheric heat accumulation, *Nature*
564 *Geosci*, 7, 345–349, <https://doi.org/10.1038/ngeo2141>, 2014.
- 565 Moher, D., Liberati, A., Tetzlaff, J., Altman, D. G., and Prisma Group: Preferred reporting items for
566 systematic reviews and meta-analyses: the PRISMA statement, *PLoS medicine*, 6, e1000097, 2009.
- 567 Mu, Q., Zhao, M., and Running, S. W.: Improvements to a MODIS global terrestrial
568 evapotranspiration algorithm, *Remote Sensing of Environment*, 115, 1781–1800,
569 <https://doi.org/10.1016/j.rse.2011.02.019>, 2011.
- 570 Pan, S., Tian, H., Dangal, S. R. S., Yang, Q., Yang, J., Lu, C., Tao, B., Ren, W., and Ouyang, Z.:
571 Responses of global terrestrial evapotranspiration to climate change and increasing atmospheric CO₂
572 in the 21st century, 3, 15–35, <https://doi.org/10.1002/2014EF000263>, 2015.
- 573 Pan, S., Pan, N., Tian, H., Friedlingstein, P., Sitch, S., Shi, H., Arora, V. K., Haverd, V., Jain, A. K.,
574 Kato, E., Lienert, S., Lombardozzi, D., Nabel, J. E. M. S., Ottlé, C., Poulter, B., Zaehle, S., and
575 Running, S. W.: Evaluation of global terrestrial evapotranspiration using state-of-the-art approaches
576 in remote sensing, machine learning and land surface modeling, 24, 1485–1509,
577 <https://doi.org/10.5194/hess-24-1485-2020>, 2020.
- 578 Papale, D., Black, T. A., Carvalhais, N., Cescatti, A., Chen, J., Jung, M., Kiely, G., Lasslop, G.,
579 Mahecha, M. D., Margolis, H., Merbold, L., Montagnani, L., Moors, E., Olesen, Jø. E., Reichstein,
580 M., Tramontana, G., Van Gorsel, E., Wohlfahrt, G., and Ráduly, B.: Effect of spatial sampling from
581 European flux towers for estimating carbon and water fluxes with artificial neural networks, 120,
582 1941–1957, <https://doi.org/10.1002/2015JG002997>, 2015.
- 583 Paul-Limoges, E., Wolf, S., Schneider, F. D., Longo, M., Moorcroft, P., Gharun, M., and Damm, A.:
584 Partitioning evapotranspiration with concurrent eddy covariance measurements in a mixed forest,
585 *Agricultural and Forest Meteorology*, 280, 107786, <https://doi.org/10.1016/j.agrformet.2019.107786>,
586 2020.

- 587 Peel, M. C., Finlayson, B. L., and McMahon, T. A.: Updated world map of the Köppen-Geiger
588 climate classification, 11, 1633–1644, <https://doi.org/10.5194/hess-11-1633-2007>, 2007.
- 589 Rigden, A. J. and Salvucci, G. D.: Evapotranspiration based on equilibrated relative humidity
590 (ETRHEQ): Evaluation over the continental U.S., 51, 2951–2973,
591 <https://doi.org/10.1002/2014WR016072>, 2015.
- 592 Sahoo, A. K., Pan, M., Troy, T. J., Vinukollu, R. K., Sheffield, J., and Wood, E. F.: Reconciling the
593 global terrestrial water budget using satellite remote sensing, *Remote Sensing of Environment*, 115,
594 1850–1865, <https://doi.org/10.1016/j.rse.2011.03.009>, 2011.
- 595 Sándor, R., Barcza, Z., Hidy, D., Lellei-Kovács, E., Ma, S., and Bellocchi, G.: Modelling of grassland
596 fluxes in Europe: Evaluation of two biogeochemical models, *Agriculture, Ecosystems &
597 Environment*, 215, 1–19, <https://doi.org/10.1016/j.agee.2015.09.001>, 2016.
- 598 Shi, H., Hellwich, O., Luo, G., Chen, C., He, H., Ochege, F. U., Van de Voorde, T., Kurban, A., and
599 de Maeyer, P.: A global meta-analysis of soil salinity prediction integrating satellite remote sensing,
600 soil sampling, and machine learning, 1–15, <https://doi.org/10.1109/TGRS.2021.3109819>, 2021.
- 601 Shi, H., Luo, G., Hellwich, O., Xie, M., Zhang, C., Zhang, Y., Wang, Y., Yuan, X., Ma, X., and
602 Zhang, W.: Variability and Uncertainty in Flux-Site Scale Net Ecosystem Exchange Simulations
603 Based on Machine Learning and Remote Sensing: A Systematic Evaluation, *Biogeosciences
604 Discussions*, 1–25, 2022.
- 605 Tramontana, G., Jung, M., Schwalm, C. R., Ichii, K., Camps-Valls, G., Ráduly, B., Reichstein, M.,
606 Arain, M. A., Cescatti, A., Kiely, G., Merbold, L., Serrano-Ortiz, P., Sickert, S., Wolf, S., and Papale,
607 D.: Predicting carbon dioxide and energy fluxes across global FLUXNET sites with regression
608 algorithms, *Biogeosciences*, 13, 4291–4313, <https://doi.org/10.5194/bg-13-4291-2016>, 2016.
- 609 Van Hulse, J., Khoshgoftaar, T. M., and Napolitano, A.: Experimental perspectives on learning from
610 imbalanced data, in: *Proceedings of the 24th international conference on Machine learning*, New
611 York, NY, USA, 935–942, <https://doi.org/10.1145/1273496.1273614>, 2007.
- 612 Van Wijk, M. T. and Bouten, W.: Water and carbon fluxes above European coniferous forests
613 modelled with artificial neural networks, [https://doi.org/10.1016/S0304-3800\(99\)00101-5](https://doi.org/10.1016/S0304-3800(99)00101-5), 1999.
- 614 Virkkala, A.-M., Aalto, J., Rogers, B. M., Tagesson, T., Treat, C. C., Natali, S. M., Watts, J. D.,
615 Potter, S., Lehtonen, A., Mauritz, M., Schuur, E. A. G., Kochendorfer, J., Zona, D., Oechel, W.,
616 Kobayashi, H., Humphreys, E., Goeckede, M., Iwata, H., Lafleur, P. M., Euskirchen, E. S., Bokhorst,
617 S., Marushchak, M., Martikainen, P. J., Elberling, B., Voigt, C., Biasi, C., Sonnentag, O., Parmentier,
618 F.-J. W., Ueyama, M., Celis, G., St.Louis, V. L., Emmerton, C. A., Peichl, M., Chi, J., Järveoja, J.,
619 Nilsson, M. B., Oberbauer, S. F., Torn, M. S., Park, S.-J., Dolman, H., Mammarella, I., Chae, N.,
620 Poyatos, R., López-Blanco, E., Christensen, T. R., Kwon, M. J., Sachs, T., Holl, D., and Luoto, M.:
621 Statistical upscaling of ecosystem CO₂ fluxes across the terrestrial tundra and boreal domain:
622 Regional patterns and uncertainties, *Global Change Biology*, 27, 4040–4059,
623 <https://doi.org/10.1111/gcb.15659>, 2021.
- 624 Wagle, P., Bhattarai, N., Gowda, P. H., and Kakani, V. G.: Performance of five surface energy
625 balance models for estimating daily evapotranspiration in high biomass sorghum, *ISPRS Journal of
626 Photogrammetry and Remote Sensing*, 128, 192–203, <https://doi.org/10.1016/j.isprsjprs.2017.03.022>,
627 2017.
- 628 Xie, M., Luo, G., Hellwich, O., Frankl, A., Zhang, W., Chen, C., Zhang, C., and De Maeyer, P.:
629 Simulation of site-scale water fluxes in desert and natural oasis ecosystems of the arid region in
630 Northwest China, 35, e14444, <https://doi.org/10.1002/hyp.14444>, 2021.

631 Xu, T., Guo, Z., Liu, S., He, X., Meng, Y., Xu, Z., Xia, Y., Xiao, J., Zhang, Y., Ma, Y., and Song, L.:
632 Evaluating Different Machine Learning Methods for Upscaling Evapotranspiration from Flux Towers
633 to the Regional Scale, 123, 8674–8690, <https://doi.org/10.1029/2018JD028447>, 2018.

634 Yang, F., White, M. A., Michaelis, A. R., Ichii, K., Hashimoto, H., Votava, P., Zhu, A.-X., and
635 Nemani, R. R.: Prediction of Continental-Scale Evapotranspiration by Combining MODIS and
636 AmeriFlux Data Through Support Vector Machine, 44, 3452–3461,
637 <https://doi.org/10.1109/TGRS.2006.876297>, 2006.

638 Zeng, J., Matsunaga, T., Tan, Z.-H., Saigusa, N., Shirai, T., Tang, Y., Peng, S., and Fukuda, Y.:
639 Global terrestrial carbon fluxes of 1999–2019 estimated by upscaling eddy covariance data with a
640 random forest, 7, <https://doi.org/10.1038/s41597-020-00653-5>, 2020.

641 Zeng, Y., Hao, D., Huete, A., Dechant, B., Berry, J., Chen, J. M., Joiner, J., Frankenberg, C., Bond-
642 Lamberty, B., Ryu, Y., Xiao, J., Asrar, G. R., and Chen, M.: Optical vegetation indices for monitoring
643 terrestrial ecosystems globally, *Nat Rev Earth Environ*, 1–17, <https://doi.org/10.1038/s43017-022-00298-5>, 2022.

645 Zhang, C., Luo, G., Hellwich, O., Chen, C., Zhang, W., Xie, M., He, H., Shi, H., and Wang, Y.: A
646 framework for estimating actual evapotranspiration at weather stations without flux observations by
647 combining data from MODIS and flux towers through a machine learning approach, *Journal of*
648 *Hydrology*, 603, 127047, <https://doi.org/10.1016/j.jhydrol.2021.127047>, 2021.

649 Zhang, K., Kimball, J. S., Nemani, R. R., and Running, S. W.: A continuous satellite-derived global
650 record of land surface evapotranspiration from 1983 to 2006, 46,
651 <https://doi.org/10.1029/2009WR008800>, 2010.

652 Zhao, W. L., Gentine, P., Reichstein, M., Zhang, Y., Zhou, S., Wen, Y., Lin, C., Li, X., and Qiu, G.
653 Y.: Physics-Constrained Machine Learning of Evapotranspiration, 46, 14496–14507,
654 <https://doi.org/10.1029/2019GL085291>, 2019.

655
656