



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

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

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

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

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

13 ⁶ Department 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



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 in both cross-study and intra-study (with the same training dataset)
24 comparisons. The average accuracy of the model applied to arid regions is higher than other climate classes. The
25 average accuracy of the model was slightly lower for forest sites (average R-squared = 0.76) than for cropland
26 and grassland sites (average R-squared = 0.8 and 0.79), but higher than for shrub sites (average R-squared =
27 0.67). Among various predictor variables, the use of net/sun radiation, precipitation, air temperature, and
28 the fraction of absorbed photosynthetically active radiation improved the model accuracy. Among the different
29 validation methods, random cross-validation shows higher model accuracy than spatial cross-validation and
30 temporal cross-validation, but spatial cross-validation is more important for the application for water flux
31 predictive models when used for spatial extrapolation. The findings of this study are promising to guide future
32 research on such machine learning-based modeling.

33 **1 Introduction**

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

43
44 Currently, there are three main approaches for simulation and spatial and temporal prediction of ET: (i) physical
45 models based on remote sensing such as surface energy balance models (Minacapilli et al., 2009; Wagle et al.,
46 2017), Penman-Monteith equation (Mu et al., 2011; Zhang et al., 2010), Priestley-Taylor equation (Miralles et
47 al., 2011); (ii) process-based land surface models, biogeochemical models and hydrological models (Barman et
48 al., 2014; Pan et al., 2015; Sándor et al., 2016; Chen et al., 2019); and (iii) the observation-based machine
49 learning modeling approach with in situ eddy covariance (EC) observations of water flux (Jung et al., 2011; Li
50 et al., 2018; Van Wijk and Bouten, 1999; Xie et al., 2021; Xu et al., 2018; Yang et al., 2006; Zhang et al., 2021).
51 For remote sensing-based physical models and process-based land surface models, some physical processes
52 have not been well characterized due to the lack of understanding of the detailed mechanisms influencing ET
53 under different environmental conditions. Limited by complicated assumptions and model parametrizations,
54 these process-based models face challenges in the accuracy of their ET estimations over heterogeneous



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

62

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

73

74 A variety of features may affect the accuracy of such models, including the predictor variables used, the inherent
75 heterogeneity within the dataset, the plant functional type (PFT) of the flux sites, the method of model
76 construction and validation, and the machine learning algorithm chosen:

77 a) Predictor variables used: Compared to process-based models, data used may have a more significant impact
78 on the final model performance in data-driven models. Various biophysical covariates and other
79 environmental factors have been used for the simulation and prediction of water fluxes. The most
80 commonly used factors include mainly precipitation (Prec), air temperature (Ta), wind speed (Ws), net/sun
81 radiation (Rn/Rs), soil temperature (Ts), soil texture, vapor-pressure deficit (VPD), the fraction of absorbed
82 photosynthetically active radiation (FAPAR), vegetation index (e.g., NDVI, EVI), LAI, and carbon fluxes
83 (e.g., GPP). These used predictor variables and their complex interactions drive the fluctuations and
84 variability of water fluxes. They affect the accuracy of water flux simulations in two ways: their actual
85 impact on water fluxes at the process-based level and their spatio-temporal resolution and inherent accuracy.
86 The relationship between water fluxes and these variables at the process-based driving mechanism level is
87 very different under different PFTs, different climate types, and different hydrometeorological conditions.
88 For example, in irrigated croplands in arid regions, water fluxes may be highly correlated with irrigation
89 practices, and thus soil moisture may be a very important predictor variable, and its importance may be
90 significantly higher than in other PFTs. And in models that incorporate data from multiple PFTs, some
91 variables that play important roles in multiple PFTs may have higher importance. In terms of data spatial
92 and temporal resolution, the data for these predictor variables may have different scales. In terms of spatial
93 resolution, meteorological observations such as precipitation and air temperature are at the flux site scale,
94 while factors extracted from satellite remote sensing and reanalysis climate datasets cover a much larger



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

110 b) The volume of the dataset, inherent heterogeneity of the dataset, and how the model is validated: the
111 volume and inherent spatiotemporal heterogeneity of the training dataset (with more variability and
112 extremes incorporated) may affect model accuracy. Typically, training data with larger regions, multiple
113 sites, multiple PFTs, and longer year spans may have a higher degree of imbalance (Kaur et al., 2019; Van
114 Hulse et al., 2007; Virkkala et al., 2021; Zeng et al., 2020). And in machine learning, in general, modeling
115 with unbalanced data (with significant differences in the distribution between the training and validation
116 sets) may result in lower model accuracy. Currently, the most common ways of model validation include
117 spatial, temporal, and random cross-validation. Spatial validation is mainly to evaluate the ability of the
118 model to be applied in different regions or flux sites with different PFT types, and one of the common
119 methods is 'leave one site out' (Fang et al., 2020; Papale et al., 2015; Zhang et al., 2021). If the data of the
120 site left out for validation differs significantly from the distribution of the training data set, the expected
121 accuracy of the model applied at that site may be low because the trained model may not capture the
122 specific and local relationships between the water flux and the various predictor variables at that site. For
123 temporal validation, to assess the ability of the models to adapt to the interannual variability, typically some
124 years of data are used for training and the remaining years for model validation (Lu and Zhuang, 2010). If a
125 year with extreme climate is used for validation, the accuracy may be low because the training dataset may
126 not contain such extreme climate conditions. In the case of PFTs that are significantly affected by human
127 activities, such as cropland, the possible different crops grown and different land use practices (e.g.,
128 irrigation) across years can also lead to low accuracy in temporal validation. K-fold cross-validation is
129 commonly used in random cross-validation to assess the fitness of the model to the spatio-temporal
130 variability. In this case, different values of K may also affect the model accuracy. For example, for an
131 unbalanced dataset, the average model accuracy obtained from a 10-fold ($K = 10$) validation approach is
132 likely to be higher than that of a 3-fold ($K = 3$) validation approach.

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 may avoid overfitting due to the
136 introduction of randomness. However, which algorithm is better overall in different situations (i.e. applied
137 to different data sets)? Which algorithm is generally more accurate than the others when using the same
138 data set? A comprehensive evaluation of this is necessary.

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 title, abstract, and keywords to include articles with the “OR” operator applied
149 among expressions (Table 1) in the Scopus database. Preferred Reporting Items for Systematic Reviews and
150 Meta-Analyses (PRISMA) (Moher et al., 2009) is followed when filtering the papers. Articles were filtered for
151 those with water fluxes (or latent heat) simulated, with multi-variable regression used, with the determination
152 coefficient (R-squared) of the validation step reported as the metric of model performance (Shi et al., 2021;
153 Tramontana et al., 2016; Zeng et al., 2020), and published in English journals. Although RMSE is also often
154 used for model accuracy assessment, its dependence on the magnitude of water flux values makes it difficult to
155 use for fair comparisons between studies.

156 Table 1. Article search: ‘[A1 OR A2 OR A3...] AND [B1 OR B2...] AND [C1 OR C2...]

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

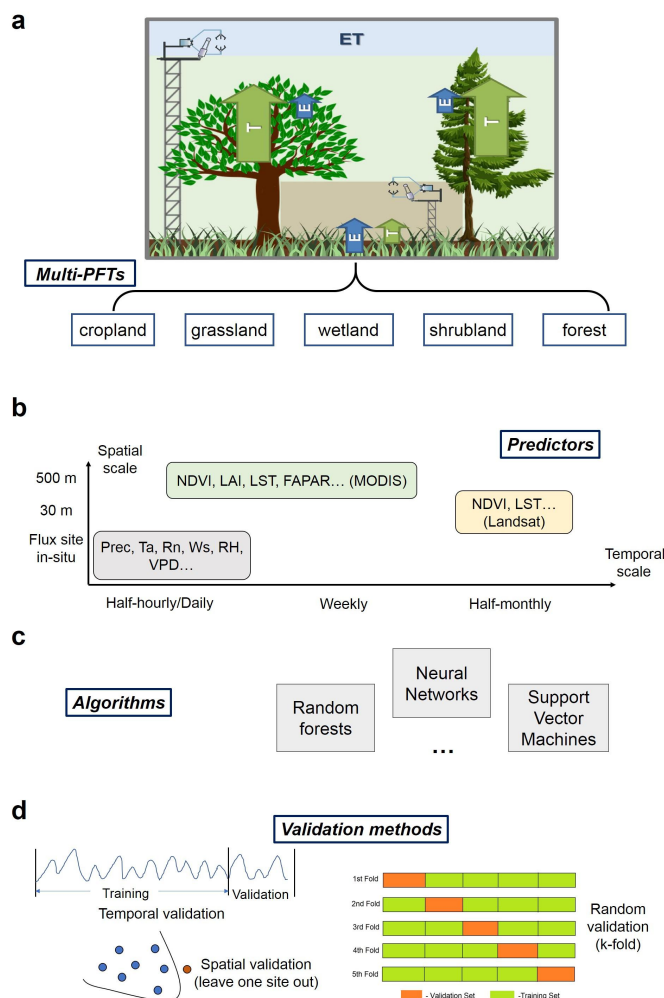
157

158 2.2 Features of the prediction processes evaluated

159 The various features (Table 2) involved in the water flux modeling framework (Fig. 1) include the PFTs of the
160 sites, the predictors used, the machine learning algorithms, the validation methods, and other features. Each
161 model for which R-squared is reported is treated as a data record. If multiple algorithms were applied to the



162 same dataset, then multiple records were extracted. Models using different data or features are also recorded as
 163 multiple records.



164
 165 Figure 1. Features of the machine learning-based water flux prediction process. (a) the eddy-covariance-based
 166 water flux observations of various plant function types (PFTs), modified from Paul-Limoges et al., 2020. ET,
 167 evapotranspiration. E, evaporation. T, transpiration. (b) Predictors and their spatial and temporal resolution. (c)
 168 The machine learning algorithms used for the modeling, such as neural networks, random forests, etc. (d) The
 169 model validation methods used including the spatial, temporal, and random cross-validations.

170

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

Field/Feature	Definition	Categories adopted
Climate	Climate zone 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
Location	More precise location (with the latitude and longitude of the center of the studied sites).	latitude, longitude
Algorithms	Algorithm families used	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)
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)
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)
Year span	The span of years of the flux data used	
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.	
Satellite images	Describe the source of satellite images used to derive NDVI, EVI, LAI, LST, etc.	Landsat, MODIS, AVHRR



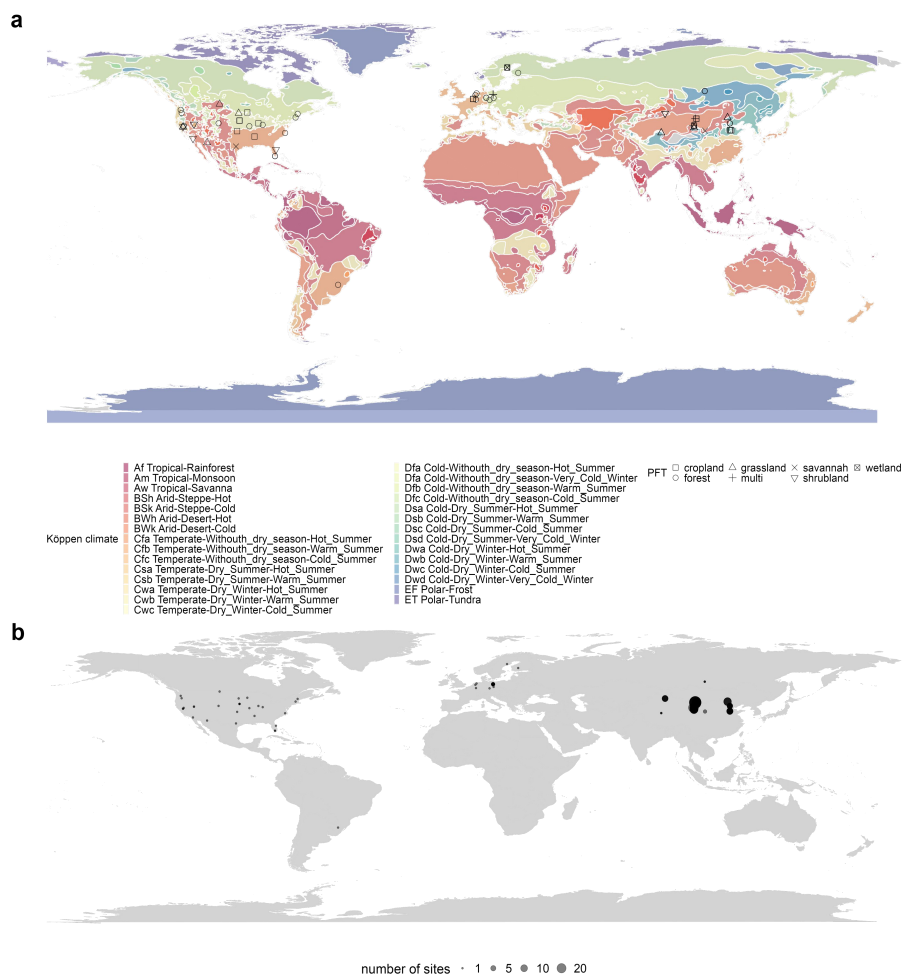
Biophysical predictors	LAI, NDVI/EVI, enhanced vegetation index (EVI), the fraction of absorbed photosynthetically active radiation/photosynthetically active radiation (FAPAR/PAR), leaf area index (LAI), Carbon fluxes (CF) including NEE/GPP, etc.	Used (recorded as '1') or not used (recorded as '0')
Meteorological variables	precipitation (Prec), net radiation/solar radiation (Rn/Rs), air temperature (Ta), vapour-pressure deficit (VPD), relative humidity (RH), etc.	Used (recorded as '1') or not used (recorded as '0')
Ancillary data	Describe the ancillary variables used: soil texture, terrain (DEM), soil moisture/land surface water index (SM/LSWI), etc.	Used (recorded as '1') or not used (recorded as '0')
Top three variables in the ranking of importance of predictors	Describe the interpretation of the importance of variables reported in the machine learning models.	
Accuracy measure	Accuracy measure used to assess the model performance	R-squared (in the validation phase)

172

173 **3 Results**

174 **3.1 Articles included in the meta-analysis**

175 A total of 32 articles (see Supplement Information) containing a total of 139 model records were included. The
 176 geographical scope of these articles was mainly Europe, North America, and China (Fig. 2).



177
 178 Figure 2. Location of the included studies in the meta-analysis. (a) PFTs and the climate zones (from Köppen
 179 climate classification) of these studies and (b) the number of flux sites included in each study. Global and
 180 continental-scale studies (e.g., models developed based on FLUXNET of the global scale) are not shown on the
 181 map due to the difficulty of identifying specific locations.

182 3.2 The formal Meta-analysis

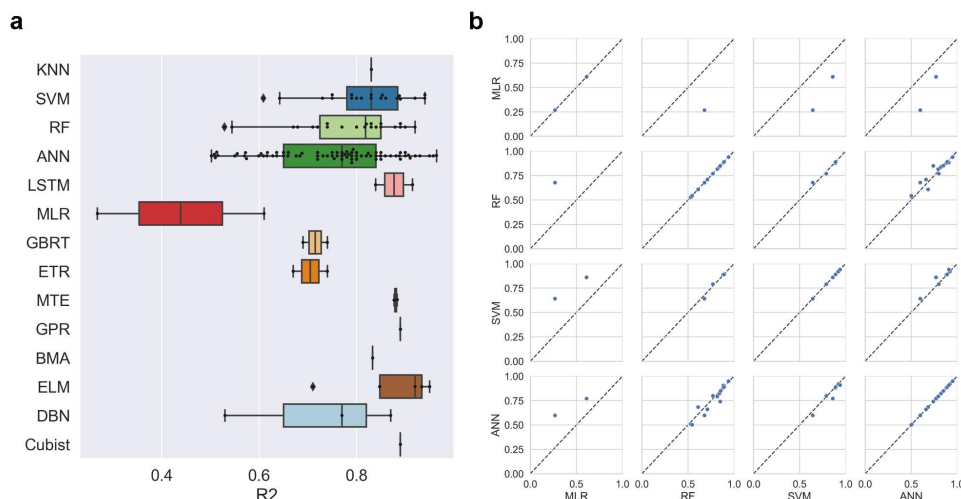
183 We formally assessed the impact of the features (e.g., algorithms, study area, PFTs, the volume of data used,
 184 validation methods, predictor variables, etc.) used in the different models based on differences of R-squared.

185

186 SVM and RF outperformed (Fig. 3a) across studies (lightly better than ANN). These three machine learning
 187 algorithms (i.e., ANN, SVM, RF) were significantly more accurate than the traditional MLR. Other algorithms
 188 such as MTE, ELM, Cubist, etc. also correspond to high accuracy, but with limited evidence sample size. In the
 189 internal comparison (different algorithms applied to the same data set) in single studies, we also find that SVM



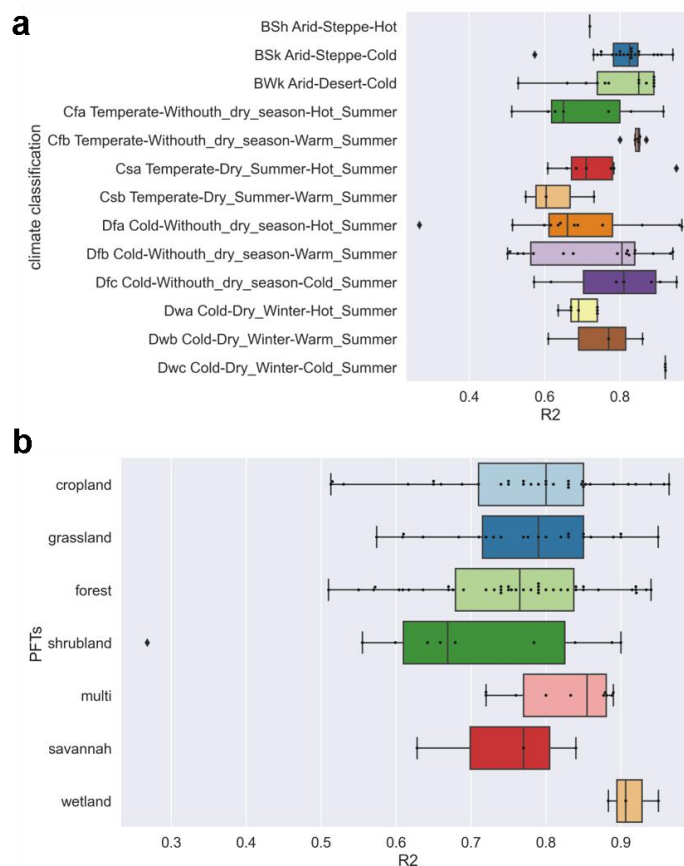
190 and RF were significantly more accurate than ANN (Fig. 3b), and all these three (i.e., ANN, SVM, RF) are
191 significantly more accurate than MLR. Overall, SVM and RF have shown higher accuracy in water flux
192 simulations.



193
194 Figure 3. Differences in model accuracy (R-squared) using different algorithms across studies (a) and internal
195 comparisons of the model accuracy (R-squared) of selected pairs of algorithms within individual studies (b).
196 Regression algorithms: Random Forests (RF), Multiple Linear Regressions (MLR), Artificial Neural Networks
197 (ANN), Support Vector Machines (SVM), Cubist, model tree ensembles (MTE), K-nearest neighbors (KNN),
198 long short-term memory (LSTM), gradient boosting regression tree (GBRT), extra tree regressor (ETR),
199 Gaussian process regression (GPR), Bayesian model averaging (BMA), extreme learning machine (ELM), and
200 deep belief network (DBN).

201

202 We found higher average model accuracy in arid climate zones (Fig. 4a), such as BSk and BWk. Most of these
203 studies were located in northwest China and the western USA. It may be caused by the simpler relationship
204 between water fluxes and biophysical covariates in arid regions. In arid zones, due to the high potential ET, the
205 variability in the actual ET may be largely explained by water availability (moisture supply) and vegetation
206 change with the effect of variability in thermal conditions reduced. As for the various PFTs, the average model
207 accuracy was slightly lower for forest types than for cropland and grassland types (Fig. 4b) possibly because
208 some remote sensing-based predictors such as FAPAR and LAI have limited accuracy when applied to forest
209 types (Fig. 5). The lowest average accuracy was found for shrub sites, which may be related to the difficulty of
210 remote sensing-based NDVI, etc., to quantify the physiological and ecological conditions of shrubs, and the
211 heterogeneity of the spatial distribution of shrubs within the EC observation area may also cause difficulties in
212 capturing their relationships with biophysical variables. We also found high model accuracy for the wetland
213 type, although records as evidence to support this finding may be limited. Compared to other PFTs, the more
214 steady and adequate water availability in the wetland type may make the variations of water fluxes less
215 explained by other biophysical covariates.

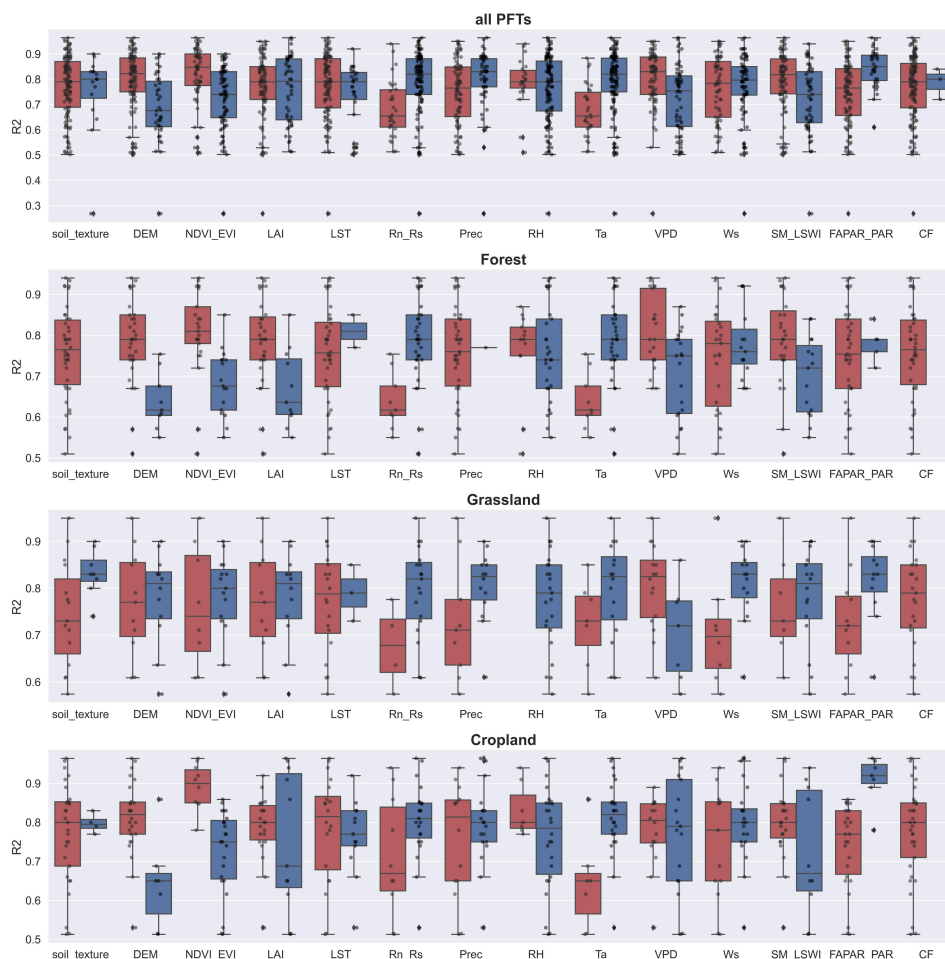


216

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

224

225 Among the various predictors, the use of Rn/Rs, Prec, Ta, and FAPAR significantly improved the accuracy of
 226 the model (Fig. 5). This pattern partially changed in the different PFTs. In the forest sites, the accuracy of the
 227 models with Rn/Rs and Ta used was significantly higher than that of the models with Rn/Rs and Ta not used.
 228 For the grassland sites, the use of Ws, FAPAR, Prec, and Rn/Rs significantly improved the model accuracy. For
 229 the cropland sites, Ta and FAPAR were more important for improving the model accuracy.
 230



231

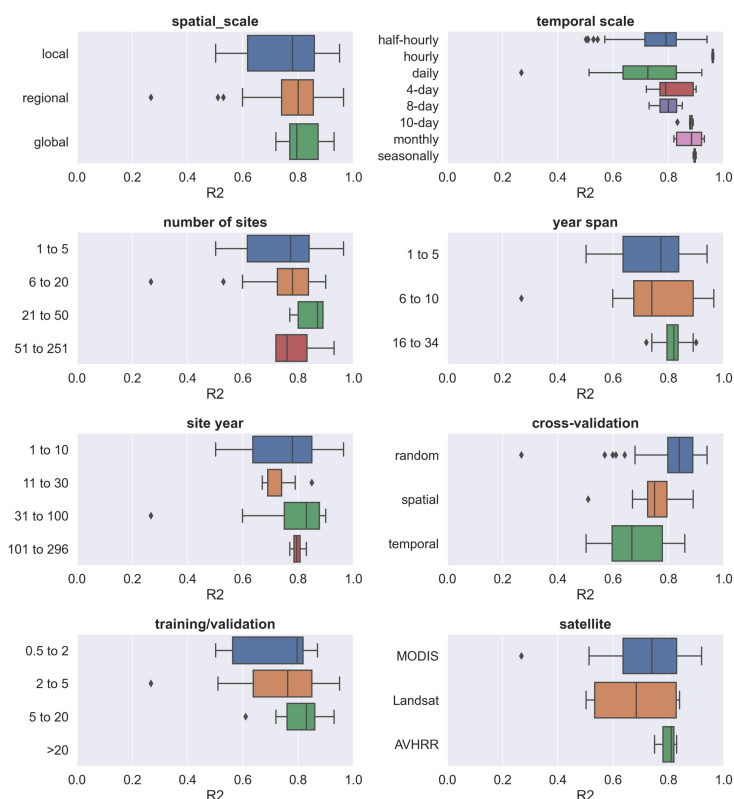
232 Figure 5. The impact of the various predictors used in models of different PFTs (all data, forest, grassland, and
233 cropland) on R-squared. Dark blue boxes indicate that the predictor was used in the model, while dark red boxes
234 indicate that the predictor was not used. Predictors: precipitation (Prec), soil moisture/land surface water index
235 (SM_LSWI), net radiation/solar radiation (Rn_Rs), enhanced vegetation index (EVI), air temperature (Ta),
236 vapor-pressure deficit (VPD), the fraction of absorbed photosynthetically active radiation/photosynthetically
237 active radiation (FAPAR_PAR), relative humidity (RH), carbon flux (CF), leaf area index (LAI).

238

239 We also evaluated the impact of some other features on accuracy. The differences in accuracy of models with
240 different spatial scales, year spans, number of sites, and volume of data (Fig. 6) appear to be insignificant. This
241 seems to be related to the fact that in large scale water flux simulations, the sites of similar PFTs are selected
242 such as for modeling multiple forest sites across Europe (Van Wijk and Bouten, 1999) which focus on ‘forest’
243 and multiple grassland sites across arid northern China (Xie et al., 2021; Zhang et al., 2021) which focus on
244 ‘grassland’, rather than mixing different PFT types to train models as the way in machine learning modeling of
245 carbon fluxes (Zeng et al., 2020). In terms of the time scales of the models, the 4-day, 8-day, and monthly scales



246 appear to correspond to higher accuracy compared to the half-hourly and daily scales. Also, the variability of the
247 accuracy of the half-hourly and daily scale models is higher. The higher the ratio of the volume of data in the
248 training and validation sets, the higher the model accuracy. Compared to the models using Landsat data, the
249 models using MODIS data showed slightly higher accuracy probably due to the advantage of MODIS data in
250 capturing the temporal dynamics of biophysical covariates. There were significant differences in the accuracy of
251 the models using different cross-validation methods, with the models using random cross-validation showing
252 higher accuracy than those using temporal cross-validation. This suggests that interannual variability may have a
253 high impact on the models in water flux simulations. The driving mechanism of ET may vary significantly
254 across years, and the inclusion of some extreme climatic conditions in the training set may be important for
255 model accuracy and robustness.
256



257
258 Figure 6. The impacts of other features (i.e. spatial scale, temporal scale, number of sites, year span, site year,
259 cross-validation method, training/validation, and satellite imagery) on the model performance.

260 4 Discussions

261 With the accumulation of in situ EC observations around the world, compared to remote sensing or process
262 model-based approaches, the study of ET simulations based on data-driven approaches has received more



263 attention from researchers in the last decade. Many studies have combined EC observations, various predictors,
264 and machine learning algorithms to improve the prediction accuracy of site-scale water fluxes. To date, the
265 results of these studies have not been comprehensively evaluated to provide clear guidance for feature selection
266 in water flux prediction models. To better understand the approach and guide future research, we performed a
267 meta-analysis of such studies. Machine learning-based water flux simulations and predictions still suffer from
268 high uncertainty. By investigating the expected improvements that can be achieved by incorporating different
269 features, we can avoid practices that may reduce model accuracy in future research.

270 **4.1 Opportunities and challenges in the site-scale water flux simulation**

271 In the above meta-analysis of the models, we found that water flux simulations based on EC observations can
272 achieve high accuracy but also have high uncertainty through the modeling workflow. The R-squared of many
273 water flux simulation models exceeds 0.8, possibly higher than some remote sensing-based and process-based
274 models, and possibly higher than carbon flux simulations in the same modeling framework. This suggests that in
275 general, these currently used biophysical and meteorological variables are closely related to water fluxes.

276
277 There are differences in model accuracy among different PFTs. For example, in forest sites, limitations in data
278 accuracy of factors were possible because some remote sensing-based predictors such as FAPAR and LAI have
279 limited accuracy when applied to forest types (Liu et al., 2018b). In addition, factors such as crown density,
280 which may significantly affect the proportion of soil evaporation, transpiration, and evaporation of canopy
281 interception, were not considered in these models, which may also lead to low model accuracy. This suggests
282 that in water flux simulation, the driving mechanisms of water fluxes in different PFTs do affect the accuracy of
283 machine learning models, and we need to consider more the actual and specific influencing factors in specific
284 PFTs. More variables that can quantify the ratio of evaporation and transpiration should be considered for
285 inclusion, which also appears to improve the mechanistic interpretability of such machine learning models.
286 Several studies (Zhao et al., 2019) have combined the physics-based approach (e.g., Penman-Monteith equation)
287 and machine learning to build hybrid models to improve interpretability. We should make full use of empirical
288 knowledge and experiences from process-based models to improve the accuracy and interpretability of the
289 machine learning approach.

290
291 The impact of differences in different satellite images on model accuracy and performance may be limited since
292 most studies used windows of 2 km x 2 km or 3 km x 3km when extracting covariates based on satellite remote
293 sensing (Walther et al., 2021) and the effects of differences in image resolution were smoothed out (i.e., the
294 differences in values averaged over a 2 km window may not be significant at 30m and 500m resolutions).
295 However, the coarse resolution of MODIS images may not be effective when the extraction window is smaller
296 (e.g., 200 m) to reduce the inconsistency of the flux footprint extent and the extracted covariates from remote
297 sensing images due to the non-homogeneity of the underlying conditions (Chu et al., 2021). Compared to the
298 16-daily temporal scale of Landsat data, the daily or 8-daily temporal scale of MODIS data may improve the
299 accuracy slightly possibly because more temporal dynamic information is explained. The inclusion of some
300 ancillary variables that do not have the temporal dimension (e.g., soil texture, topographic variables) may be of



301 more limited use unless the model includes many flux sites for which the spatial variability of the ancillary
302 variables is large enough and does affect water fluxes.
303
304 Among the different validation methods, random cross-validation has higher accuracy than spatial cross-
305 validation and temporal cross-validation. However, spatial cross-validation and temporal cross-validation may
306 be able to better help us recognize the robustness of the model when extrapolated (i.e., applied to new stations
307 and new years). The lower accuracy in the temporal cross-validation approach implies that we need to focus on
308 interannual hydrological and meteorological variability in the water flux simulations. In cropland sites, we may
309 also need to pay more attention to the effects of interannual variability in anthropogenic cropping patterns. If
310 some extreme weather years are not included, the robustness of the model when extrapolated to other years may
311 be challenged, especially in the context of the various extreme weather events of recent years. This can also
312 inform the siting of future flux stations. Regions where climate extremes may occur and biogeographic types not
313 covered by existing flux observation networks should be given more attention to achieve global-scale, accurate
314 and robust machine learning-based spatio-temporal prediction of water fluxes.

315 **4.2 Uncertainties and limitations of this meta-analysis**

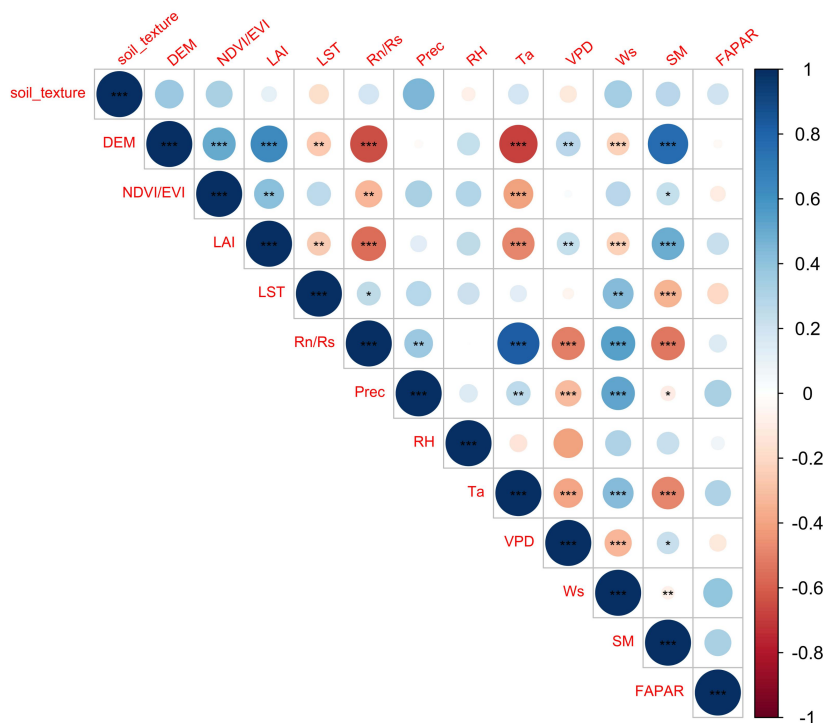
316 The potential uncertainties and limitations of the results of this meta-analysis are as follows:

- 317 a) The number of available literature and model records that can be collected: Despite many articles and
318 model records collected through our efforts to perform this meta-analysis, there still appears to be a long
319 way to go to finally and completely understand the various mechanisms involved in water flux simulation
320 with machine learning. Some of the insights provided by this study can be not robust (due to the limited
321 sample size available when the goal is to assess the effects of multiple features), but this does not negate the
322 fact that this study does obtain some meaningful findings. Therefore, researchers should treat the results of
323 this study with caution, as they were obtained only statistically. Overall, it is still positive to conduct a
324 meta-analysis of such studies, considering their rapid growth in the number and lack of guiding directions.
- 325 b) Publication bias and weighting: Due to the relatively limited number of articles that could be included in the
326 meta-analysis, this study did not focus much on publication bias. Meta-analytic studies in other fields
327 typically measure the quality of journals and the public availability of research data (Borenstein et al., 2011;
328 Field and Gillett, 2010) to determine the weighting in the literature in a comprehensive assessment.
329 However, most of the articles did not publicly provide flux observations or share developed models. Meta-
330 analysis studies in other fields typically measure the impact of included studies based on sample size and
331 variance of experimental results (Adams et al., 1997; Don et al., 2011; Liu et al., 2018a). In this study, due
332 to the lack of a convincing manner to determine weights among articles, we assigned the same weight to the
333 results for all the literature.
- 334 c) Uncertainties in the information of the extracted features: First, as most studies used far more water flux
335 observation records than the number of covariates in their regression models, we did not adjust the R-
336 squared in this study to an adjusted R-squared. Secondly, uncertainties caused by data quality control (e.g.
337 gap-filling (Hui et al., 2004)) and differences in the eddy covariance observation instruments used to
338 observe water fluxes, etc., are difficult to assess effectively. Thirdly, the various specific ways in which the



339 parameters of the model are optimized are not differentiated. They are broadly categorized into different
 340 families or kinds of algorithms, which may also introduce uncertainty into the assessment. Fourth, the
 341 assessment of some features is not detailed due to the limitations of the available model records. For
 342 example, the classification of PFT could be more detailed. ‘Forest’ could be further classified as broadleaf
 343 forest, coniferous forest, etc. while ‘cropland’ could be further classified as rainfed and irrigated cropland
 344 based on differences in their response mechanisms of water fluxes to environmental factors.

345 d) Independence between features: There is dependence between some of the features being evaluated, which
 346 may affect the assessment of the impact of single features on the accuracy of the model. We found that the
 347 use of NDVI/EVI, LAI, VPD, and SM was significantly negatively correlated with the use of Rn/Rs and Ta
 348 (Fig. 7) when unused was set to 0 and used was set to 1. It means that many of the models that used Rn/Rs
 349 and Ta did not use NDVI/EVI, LAI, VPD, and SM, and the models that used NDVI/EVI, LAI, VPD, and
 350 SM also happened to not use Rn/Rs and Ta. It can indirectly explain the fact that the accuracy of the models
 351 with NDVI/EVI, LAI, VPD, and SM is even lower than that of the models without NDVI/EVI, LAI, VPD,
 352 and SM in the above analysis (Fig. 5) because of the disturbance from the use of Rn/Rs and Ta.



353
 354 Fig. 7. Correlation matrix between the use of various predictors (not used is set as 0 and used is set as 1) which
 355 may introduce uncertainty in the assessment of the impact of an individual predictor on model performance.
 356 Significance: the p-value < 0.01 (***), 0.05 (**), and 0.1 (*).



357 **5 Conclusion**

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

- 361 a) SVM (average R-squared = 0.82) and RF (average R-squared = 0.81) outperformed over evaluated
362 algorithms in both cross-study and intra-study (with the same training dataset) comparisons.
- 363 b) The average accuracy of the model applied to arid regions is higher than other climate classes.
- 364 c) The average accuracy of the model was slightly lower for forest sites (average R-squared = 0.76) than for
365 cropland and grassland sites (average R-squared = 0.8 and 0.79), but higher than for shrub sites (average R-
366 squared = 0.67).
- 367 d) Among various predictor variables, the use of Rn/Rs, Prec, Ta, and FAPAR improved the model accuracy.
- 368 e) Among the different validation methods, random cross-validation shows higher model accuracy than spatial
369 cross-validation and temporal cross-validation.

370

371 **Financial support**

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

376 **Author Contributions**

377 Haiyang Shi: Conceptualization, Methodology, Data, Writing. Geping Luo: Conceptualization, Supervision,
378 Revision. Olaf Hellwich: Methodology. Alishir Kurban: Supervision. Tim Van De Voorde: Supervision.
379 Philippe De Maeyer: Supervision, Revision. Xiaofei Ma, Xiuliang Yuan, Yuangang Wang, Wenqiang Zhang,
380 Mingjuan Xie, Chen Zhang, Yu Zhang: Data.

381 **Competing interests**

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

383 **Data availability**

384 The data used in this study can be accessed by contacting the first author (shihaiyang16@mails.ucas.ac.cn)
385 based on reasonable request.

386



387 **References**

- 388 Adams, D. C., Gurevitch, J., and Rosenberg, M. S.: Resampling tests for meta - analysis of ecological data,
389 Ecology, 78, 1277–1283, 1997.
- 390 Allen, R. G., Pereira, L. S., Howell, T. A., and Jensen, M. E.: Evapotranspiration information reporting: I.
391 Factors governing measurement accuracy, Agricultural Water Management, 98, 899–920,
392 <https://doi.org/10.1016/j.agwat.2010.12.015>, 2011.
- 393 Anderson, M. C., Allen, R. G., Morse, A., and Kustas, W. P.: Use of Landsat thermal imagery in monitoring
394 evapotranspiration and managing water resources, Remote Sensing of Environment, 122, 50–65,
395 <https://doi.org/10.1016/j.rse.2011.08.025>, 2012.
- 396 Barman, R., Jain, A. K., and Liang, M.: Climate-driven uncertainties in modeling terrestrial energy and water
397 fluxes: a site-level to global-scale analysis, 20, 1885–1900, <https://doi.org/10.1111/gcb.12473>, 2014.
- 398 Borenstein, M., Hedges, L. V., Higgins, J. P., and Rothstein, H. R.: Introduction to meta-analysis, John Wiley &
399 Sons, 2011.
- 400 Chen, Y., Xia, J., Liang, S., Feng, J., Fisher, J. B., Li, X., Li, X., Liu, S., Ma, Z., Miyata, A., Mu, Q., Sun, L.,
401 Tang, J., Wang, K., Wen, J., Xue, Y., Yu, G., Zha, T., Zhang, L., Zhang, Q., Zhao, T., Zhao, L., and Yuan, W.:
402 Comparison of satellite-based evapotranspiration models over terrestrial ecosystems in China, Remote Sensing
403 of Environment, 140, 279–293, <https://doi.org/10.1016/j.rse.2013.08.045>, 2014.
- 404 Chen, Y., Wang, S., Ren, Z., Huang, J., Wang, X., Liu, S., Deng, H., and Lin, W.: Increased evapotranspiration
405 from land cover changes intensified water crisis in an arid river basin in northwest China, Journal of Hydrology,
406 574, 383–397, <https://doi.org/10.1016/j.jhydrol.2019.04.045>, 2019.
- 407 Chu, H., Luo, X., Ouyang, Z., Chan, W. S., Dengel, S., Biraud, S. C., Torn, M. S., Metzger, S., Kumar, J., Arain,
408 M. A., Arkebauer, T. J., Baldocchi, D., Bernacchi, C., Billesbach, D., Black, T. A., Blanken, P. D., Bohrer, G.,
409 Bracho, R., Brown, S., Brunsell, N. A., Chen, J., Chen, X., Clark, K., Desai, A. R., Duman, T., Durden, D.,
410 Fares, S., Forbrich, I., Gamon, J. A., Gough, C. M., Griffis, T., Helbig, M., Hollinger, D., Humphreys, E., Ikawa,
411 H., Iwata, H., Ju, Y., Knowles, J. F., Knox, S. H., Kobayashi, H., Kolb, T., Law, B., Lee, X., Litvak, M., Liu, H.,
412 Munger, J. W., Noormets, A., Novick, K., Oberbauer, S. F., Oechel, W., Oikawa, P., Papuga, S. A., Pendall, E.,
413 Prajapati, P., Prueger, J., Quinton, W. L., Richardson, A. D., Russell, E. S., Scott, R. L., Starr, G., Staebler, R.,
414 Stoy, P. C., Stuart-Haëntjens, E., Sonnentag, O., Sullivan, R. C., Suyker, A., Ueyama, M., Vargas, R., Wood, J.
415 D., and Zona, D.: Representativeness of Eddy-Covariance flux footprints for areas surrounding AmeriFlux sites,
416 Agricultural and Forest Meteorology, 301–302, 108350, <https://doi.org/10.1016/j.agrformet.2021.108350>, 2021.
- 417 Don, A., Schumacher, J., and Freibauer, A.: Impact of tropical land-use change on soil organic carbon stocks – a
418 meta-analysis, 17, 1658–1670, <https://doi.org/10.1111/j.1365-2486.2010.02336.x>, 2011.
- 419 Fang, B., Lei, H., Zhang, Y., Quan, Q., and Yang, D.: Spatio-temporal patterns of evapotranspiration based on
420 upscaling eddy covariance measurements in the dryland of the North China Plain, 281,
421 <https://doi.org/10.1016/j.agrformet.2019.107844>, 2020.
- 422 Field, A. P. and Gillett, R.: How to do a meta - analysis, British Journal of Mathematical and Statistical
423 Psychology, 63, 665–694, 2010.
- 424 Fisher, J. B., Melton, F., Middleton, E., Hain, C., Anderson, M., Allen, R., McCabe, M. F., Hook, S., Baldocchi,
425 D., Townsend, P. A., Kilic, A., Tu, K., Miralles, D. D., Perret, J., Lagouarde, J.-P., Waliser, D., Purdy, A. J.,
426 French, A., Schimel, D., Famiglietti, J. S., Stephens, G., and Wood, E. F.: The future of evapotranspiration:
427 Global requirements for ecosystem functioning, carbon and climate feedbacks, agricultural management, and
428 water resources, 53, 2618–2626, <https://doi.org/10.1002/2016WR020175>, 2017.
- 429 Gaston, K. J.: Global patterns in biodiversity, 405, 220–227, <https://doi.org/10.1038/35012228>, 2000.



- 430 Hui, D., Wan, S., Su, B., Katul, G., Monson, R., and Luo, Y.: Gap-filling missing data in eddy covariance
431 measurements using multiple imputation (MI) for annual estimations, 121, 93–111,
432 [https://doi.org/10.1016/S0168-1923\(03\)00158-8](https://doi.org/10.1016/S0168-1923(03)00158-8), 2004.
- 433 Jung, M., Reichstein, M., and Bondeau, A.: Towards global empirical upscaling of FLUXNET eddy covariance
434 observations: Validation of a model tree ensemble approach using a biosphere model, 6, 2001–2013,
435 <https://doi.org/10.5194/bg-6-2001-2009>, 2009.
- 436 Jung, M., Reichstein, M., Margolis, H. A., Cescatti, A., Richardson, A. D., Arain, M. A., Arneeth, A., Bernhofer,
437 C., Bonal, D., Chen, J., Gianelle, D., Gobron, N., Kiely, G., Kutsch, W., Lasslop, G., Law, B. E., Lindroth, A.,
438 Merbold, L., Montagnani, L., Moors, E. J., Papale, D., Sottocornola, M., Vaccari, F., and Williams, C.: Global
439 patterns of land-atmosphere fluxes of carbon dioxide, latent heat, and sensible heat derived from eddy
440 covariance, satellite, and meteorological observations, 116, <https://doi.org/10.1029/2010JG001566>, 2011.
- 441 Kaur, H., Pannu, H. S., and Malhi, A. K.: A Systematic Review on Imbalanced Data Challenges in Machine
442 Learning: Applications and Solutions, *ACM Comput. Surv.*, 52, 79:1-79:36, <https://doi.org/10.1145/3343440>,
443 2019.
- 444 Li, X., He, Y., Zeng, Z., Lian, X., Wang, X., Du, M., Jia, G., Li, Y., Ma, Y., Tang, Y., Wang, W., Wu, Z., Yan,
445 J., Yao, Y., Ciais, P., Zhang, X., Zhang, Y., Zhang, Y., Zhou, G., and Piao, S.: Spatiotemporal pattern of
446 terrestrial evapotranspiration in China during the past thirty years, 259, 131–140,
447 <https://doi.org/10.1016/j.agrformet.2018.04.020>, 2018.
- 448 Liu, Q., Zhang, Y., Liu, B., Amonette, J. E., Lin, Z., Liu, G., Ambus, P., and Xie, Z.: How does biochar
449 influence soil N cycle? A meta-analysis, *Plant and soil*, 426, 211–225, 2018a.
- 450 Liu, Y., Xiao, J., Ju, W., Zhu, G., Wu, X., Fan, W., Li, D., and Zhou, Y.: Satellite-derived LAI products exhibit
451 large discrepancies and can lead to substantial uncertainty in simulated carbon and water fluxes, *Remote
452 Sensing of Environment*, 206, 174–188, <https://doi.org/10.1016/j.rse.2017.12.024>, 2018b.
- 453 Lu, X. and Zhuang, Q.: Evaluating evapotranspiration and water-use efficiency of terrestrial ecosystems in the
454 conterminous United States using MODIS and AmeriFlux data, <https://doi.org/10.1016/j.rse.2010.04.001>, 2010.
- 455 Minacapilli, M., Agnese, C., Blanda, F., Cammalleri, C., Ciraolo, G., D’Urso, G., Iovino, M., Pumo, D.,
456 Provenzano, G., and Rallo, G.: Estimation of actual evapotranspiration of Mediterranean perennial crops by
457 means of remote-sensing based surface energy balance models, 13, 1061–1074, <https://doi.org/10.5194/hess-13-1061-2009>, 2009.
- 459 Miralles, D. G., Holmes, T. R. H., De Jeu, R. a. M., Gash, J. H., Meesters, A. G. C. A., and Dolman, A. J.:
460 Global land-surface evaporation estimated from satellite-based observations, 15, 453–469,
461 <https://doi.org/10.5194/hess-15-453-2011>, 2011.
- 462 Miralles, D. G., Teuling, A. J., van Heerwaarden, C. C., and Vilà-Guerau de Arellano, J.: Mega-heatwave
463 temperatures due to combined soil desiccation and atmospheric heat accumulation, *Nature Geosci*, 7, 345–349,
464 <https://doi.org/10.1038/ngeo2141>, 2014.
- 465 Moher, D., Liberati, A., Tetzlaff, J., Altman, D. G., and Prisma Group: Preferred reporting items for systematic
466 reviews and meta-analyses: the PRISMA statement, *PLoS medicine*, 6, e1000097, 2009.
- 467 Mu, Q., Zhao, M., and Running, S. W.: Improvements to a MODIS global terrestrial evapotranspiration
468 algorithm, *Remote Sensing of Environment*, 115, 1781–1800, <https://doi.org/10.1016/j.rse.2011.02.019>, 2011.
- 469 Pan, S., Tian, H., Dangal, S. R. S., Yang, Q., Yang, J., Lu, C., Tao, B., Ren, W., and Ouyang, Z.: Responses of
470 global terrestrial evapotranspiration to climate change and increasing atmospheric CO₂ in the 21st century, 3,
471 15–35, <https://doi.org/10.1002/2014EF000263>, 2015.
- 472 Pan, S., Pan, N., Tian, H., Friedlingstein, P., Sitch, S., Shi, H., Arora, V. K., Haverd, V., Jain, A. K., Kato, E.,
473 Lienert, S., Lombardozzi, D., Nabel, J. E. M. S., Ottlé, C., Poulter, B., Zaehle, S., and Running, S. W.:
474 Evaluation of global terrestrial evapotranspiration using state-of-the-art approaches in remote sensing, machine
475 learning and land surface modeling, 24, 1485–1509, <https://doi.org/10.5194/hess-24-1485-2020>, 2020.



- 476 Papale, D., Black, T. A., Carvalhais, N., Cescatti, A., Chen, J., Jung, M., Kiely, G., Lasslop, G., Mahecha, M. D.,
477 Margolis, H., Merbold, L., Montagnani, L., Moors, E., Olesen, Jø. E., Reichstein, M., Tramontana, G., Van
478 Gorsel, E., Wohlfahrt, G., and Ráduly, B.: Effect of spatial sampling from European flux towers for estimating
479 carbon and water fluxes with artificial neural networks, 120, 1941–1957, <https://doi.org/10.1002/2015JG002997>,
480 2015.
- 481 Paul-Limoges, E., Wolf, S., Schneider, F. D., Longo, M., Moorcroft, P., Gharun, M., and Damm, A.:
482 Partitioning evapotranspiration with concurrent eddy covariance measurements in a mixed forest, *Agricultural
483 and Forest Meteorology*, 280, 107786, <https://doi.org/10.1016/j.agrformet.2019.107786>, 2020.
- 484 Peel, M. C., Finlayson, B. L., and McMahon, T. A.: Updated world map of the Köppen-Geiger climate
485 classification, 11, 1633–1644, <https://doi.org/10.5194/hess-11-1633-2007>, 2007.
- 486 Sahoo, A. K., Pan, M., Troy, T. J., Vinukollu, R. K., Sheffield, J., and Wood, E. F.: Reconciling the global
487 terrestrial water budget using satellite remote sensing, *Remote Sensing of Environment*, 115, 1850–1865,
488 <https://doi.org/10.1016/j.rse.2011.03.009>, 2011.
- 489 Sándor, R., Barcza, Z., Hidy, D., Lellei-Kovács, E., Ma, S., and Bellocchi, G.: Modelling of grassland fluxes in
490 Europe: Evaluation of two biogeochemical models, *Agriculture, Ecosystems & Environment*, 215, 1–19,
491 <https://doi.org/10.1016/j.agee.2015.09.001>, 2016.
- 492 Shi, H., Hellwich, O., Luo, G., Chen, C., He, H., Ochege, F. U., Van de Voorde, T., Kurban, A., and de Maeyer,
493 P.: A global meta-analysis of soil salinity prediction integrating satellite remote sensing, soil sampling, and
494 machine learning, 1–15, <https://doi.org/10.1109/TGRS.2021.3109819>, 2021.
- 495 Tramontana, G., Jung, M., Schwalm, C. R., Ichii, K., Camps-Valls, G., Ráduly, B., Reichstein, M., Arain, M. A.,
496 Cescatti, A., Kiely, G., Merbold, L., Serrano-Ortiz, P., Sickert, S., Wolf, S., and Papale, D.: Predicting carbon
497 dioxide and energy fluxes across global FLUXNET sites with regression algorithms, 13, 4291–4313,
498 <https://doi.org/10.5194/bg-13-4291-2016>, 2016.
- 499 Van Hulse, J., Khoshgoftaar, T. M., and Napolitano, A.: Experimental perspectives on learning from imbalanced
500 data, in: *Proceedings of the 24th international conference on Machine learning*, New York, NY, USA, 935–942,
501 <https://doi.org/10.1145/1273496.1273614>, 2007.
- 502 Van Wijk, M. T. and Bouten, W.: Water and carbon fluxes above European coniferous forests modelled with
503 artificial neural networks, [https://doi.org/10.1016/S0304-3800\(99\)00101-5](https://doi.org/10.1016/S0304-3800(99)00101-5), 1999.
- 504 Virkkala, A.-M., Aalto, J., Rogers, B. M., Tagesson, T., Treat, C. C., Natali, S. M., Watts, J. D., Potter, S.,
505 Lehtonen, A., Mauritz, M., Schuur, E. A. G., Kochendorfer, J., Zona, D., Oechel, W., Kobayashi, H.,
506 Humphreys, E., Goeckede, M., Iwata, H., Lafleur, P. M., Euskirchen, E. S., Bokhorst, S., Marushchak, M.,
507 Martikainen, P. J., Elberling, B., Voigt, C., Biasi, C., Sonntag, O., Parmentier, F.-J. W., Ueyama, M., Celis,
508 G., St. Louis, V. L., Emmerton, C. A., Peichl, M., Chi, J., Järveoja, J., Nilsson, M. B., Oberbauer, S. F., Torn, M.
509 S., Park, S.-J., Dolman, H., Mammarella, I., Chae, N., Poyatos, R., López-Blanco, E., Christensen, T. R., Kwon,
510 M. J., Sachs, T., Holl, D., and Luoto, M.: Statistical upscaling of ecosystem CO₂ fluxes across the terrestrial
511 tundra and boreal domain: Regional patterns and uncertainties, 27, 4040–4059,
512 <https://doi.org/10.1111/gcb.15659>, 2021.
- 513 Wagle, P., Bhattarai, N., Gowda, P. H., and Kakani, V. G.: Performance of five surface energy balance models
514 for estimating daily evapotranspiration in high biomass sorghum, *ISPRS Journal of Photogrammetry and
515 Remote Sensing*, 128, 192–203, <https://doi.org/10.1016/j.isprsjprs.2017.03.022>, 2017.
- 516 Walther, S., Besnard, S., Nelson, J. A., El-Madany, T. S., Migliavacca, M., Weber, U., Ermida, S. L., Brümmer,
517 C., Schrader, F., Prokushkin, A. S., Panov, A. V., and Jung, M.: Technical note: A view from space on global
518 flux towers by MODIS and Landsat: The FluxnetEO dataset, 1–40, <https://doi.org/10.5194/bg-2021-314>, 2021.
- 519 Xie, M., Luo, G., Hellwich, O., Frankl, A., Zhang, W., Chen, C., Zhang, C., and De Maeyer, P.: Simulation of
520 site-scale water fluxes in desert and natural oasis ecosystems of the arid region in Northwest China, 35, e14444,
521 <https://doi.org/10.1002/hyp.14444>, 2021.



- 522 Xu, T., Guo, Z., Liu, S., He, X., Meng, Y., Xu, Z., Xia, Y., Xiao, J., Zhang, Y., Ma, Y., and Song, L.:
523 Evaluating Different Machine Learning Methods for Upscaling Evapotranspiration from Flux Towers to the
524 Regional Scale, 123, 8674–8690, <https://doi.org/10.1029/2018JD028447>, 2018.
- 525 Yang, F., White, M. A., Michaelis, A. R., Ichii, K., Hashimoto, H., Votava, P., Zhu, A.-X., and Nemani, R. R.:
526 Prediction of Continental-Scale Evapotranspiration by Combining MODIS and AmeriFlux Data Through
527 Support Vector Machine, 44, 3452–3461, <https://doi.org/10.1109/TGRS.2006.876297>, 2006.
- 528 Zeng, J., Matsunaga, T., Tan, Z.-H., Saigusa, N., Shirai, T., Tang, Y., Peng, S., and Fukuda, Y.: Global
529 terrestrial carbon fluxes of 1999–2019 estimated by upscaling eddy covariance data with a random forest, 7,
530 <https://doi.org/10.1038/s41597-020-00653-5>, 2020.
- 531 Zhang, C., Luo, G., Hellwich, O., Chen, C., Zhang, W., Xie, M., He, H., Shi, H., and Wang, Y.: A framework
532 for estimating actual evapotranspiration at weather stations without flux observations by combining data from
533 MODIS and flux towers through a machine learning approach, *Journal of Hydrology*, 603, 127047,
534 <https://doi.org/10.1016/j.jhydrol.2021.127047>, 2021.
- 535 Zhang, K., Kimball, J. S., Nemani, R. R., and Running, S. W.: A continuous satellite-derived global record of
536 land surface evapotranspiration from 1983 to 2006, 46, <https://doi.org/10.1029/2009WR008800>, 2010.
- 537 Zhao, W. L., Gentine, P., Reichstein, M., Zhang, Y., Zhou, S., Wen, Y., Lin, C., Li, X., and Qiu, G. Y.: Physics-
538 Constrained Machine Learning of Evapotranspiration, 46, 14496–14507,
539 <https://doi.org/10.1029/2019GL085291>, 2019.
- 540
541