

1 **Variability and Uncertainty in Flux-Site Scale Net Ecosystem**
2 **Exchange Simulations Based on Machine Learning and**
3 **Remote Sensing: A Systematic Evaluation**

4 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
5 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
6 Tim Van de Voorde^{4,5}

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

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

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

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

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

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

15

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

17 Submitted to *Biogeosciences*

18 **Abstract.** Net ecosystem exchange (NEE) is an important indicator of carbon cycling in terrestrial ecosystems.
19 Many previous studies have combined flux observations, meteorological, biophysical, and ancillary predictors
20 using machine learning to simulate the site-scale NEE. However, systematic evaluation of the performance of
21 such models is limited. Therefore, we performed a meta-analysis of these NEE simulations. A total of 40 such
22 studies and 178 model records were included. The impacts of various features throughout the modeling process
23 on the accuracy of the model were evaluated. Random Forests and Support Vector Machines performed better
24 than other algorithms. Models with larger time scales have lower average R-squared, especially when the time
25 scale exceeds the monthly scale. Half-hourly models (average R-squared = 0.73) were significantly more
26 accurate than daily models (average R-squared = 0.5). There are significant differences in the predictors used
27 and their impacts on model accuracy for different plant functional types (PFTs). Studies at continental and
28 global scales (average R-squared = 0.37) with multiple PFTs, more sites, and a large span of years correspond to
29 lower R-squared than studies at local (average R-squared = 0.69) and regional scales (average R-squared = 0.7).
30 Also, the site-scale NEE predictions need more focus on the internal heterogeneity of the NEE dataset and the
31 matching of the training set and validation set.

32 **1 Introduction**

33 Net ecosystem exchange (NEE) of CO₂ is an important indicator of carbon cycling in terrestrial ecosystems (Fu
34 et al., 2019), and accurate estimation of NEE is important for the development of global carbon neutral policies.
35 Although process-based models have been used for NEE simulations (Mitchell et al., 2009), their accuracy and
36 spatial resolutions of the model outputs are limited probably due to the lack of understanding and quantification
37 of complex processes. Many researchers have tried to use a data-driven approach as an alternative (Fu et al.,
38 2014; Tian et al., 2017; Tramontana et al., 2016; Jung et al., 2011). On the one hand, it was made possible by
39 the increase in the growth of global carbon flux observations and the large amount of flux observation data
40 being accumulated. Since the 1990s, the use of the eddy covariance technique to monitor NEE has been rapidly
41 promoted (Baldocchi, 2003). Several regional and global flux measurement networks have been established for
42 the big data management of the flux sites, including CarboEuro-flux (Europe), AmeriFlux (North America),
43 OzFlux (Australia), ChinaFlux (China), FLUXNET (global), etc. On the other hand, machine learning
44 approaches are increasingly used to extract patterns and insights from the ever-increasing stream of geospatial
45 data (Reichstein et al., 2019). The rapid development of various algorithms and high public availability of model
46 tools in the field of machine learning have made these techniques easily available to more researchers in the
47 field of geography and ecology (Reichstein et al., 2019). Since the above two major advances (i.e., increasing
48 availability of flux data and machine learning techniques) in the last two decades, various machine learning
49 algorithms have been used to simulate NEE at the flux station scale with various predictor variables (e.g.,
50 meteorological variables, biophysical variables) incorporated for spatial and temporal mapping of NEE or
51 understanding the driving mechanisms of NEE.

52
53 To date, studies on using machine learning to predict NEE have a high diversity in terms of modeling
54 approaches. To obtain a comprehensive understanding of machine learning-based NEE prediction, a synthesis
55 evaluation of these machine learning models is necessary. Since the beginning of this century, when machine

56 learning approaches were still rarely used in geography and ecology research, neural networks were already
57 used to perform simulations and mapping of NEE in European forests (Papale and Valentini, 2003).
58 Subsequently, considerable efforts have been made by researchers to improve such predictive models. Many
59 studies have demonstrated the effectiveness of their proposed improvements (i.e., using predictors with a higher
60 spatial resolution (Reitz et al., 2021) and using data from the local flux site network (Cho et al., 2021)) by
61 comparing with previous studies. However, the improvements achieved in these studies may be limited to
62 smaller areas and specific conditions and may not be generalizable (Cleverly et al., 2020; Reed et al., 2021; Cho
63 et al., 2021). We are more interested in guidelines with universal applicability that improve the model accuracy,
64 such as the selection of appropriate predictors and algorithms under different conditions. Therefore, we should
65 synthesize the results of models applied to different conditions and regions to obtain general insights.

66

67 Many factors may affect the performance of these NEE prediction models, such as the predictor variables, the
68 spatial and temporal span of the observed flux data, the plant functional type (PFT) of the flux sites, the model
69 validation method, the machine learning algorithm used, as described below:

70 a) Predictors: Various biophysical variables (Zeng et al., 2020; Cui et al., 2021; Huemmrich et al., 2019) and
71 other meteorological and environmental factors have been used in the simulation of NEE. The most
72 commonly used predictor variables include precipitation (Prec), air temperature (Ta), wind speed (Ws),
73 net/sun radiation (Rn/Rs), soil temperature (Ts), soil texture, soil moisture (SM) (Zhou et al., 2020), vapor-
74 pressure deficit (VPD) (Moffat et al., 2010; Park et al., 2018), the fraction of absorbed photosynthetically
75 active radiation (FAPAR) (Park et al., 2018; Tian et al., 2017), vegetation index (e.g., NDVI, EVI), LAI,
76 and evapotranspiration (ET) (Berryman et al., 2018). The predictor variables used vary with the natural
77 conditions and vegetation functional types of the study area. In contrast, in models that include multiple
78 PFTs, some variables that play a significant role in the prediction of each of the multiple PFTs may have
79 higher importance. For example, growing degree days (GDD) may be a more effective variable for NEE of
80 tundra in the northern hemisphere high latitudes (Virkkala et al., 2021), while measured groundwater levels
81 may be important for wetlands (Zhang et al., 2021). Some of these predictor variables are measured at flux
82 stations (e.g., meteorological factors such as precipitation and temperature), while others are extracted
83 from reanalyzed meteorological datasets and satellite remote sensing image data (e.g., vegetation indices).
84 The spatial and temporal resolution of predictors can lead to differences in their relevance to NEE
85 observations. Most measured in situ meteorological factors have a good spatio-temporal match to the
86 observed NEE (site scale, half-hourly scale). However, the proportion of NEE explained by remotely
87 sensed biophysical covariates may depend on their spatial and temporal scales. For example, the MODIS-
88 based 8-daily NDVI data may better capture temporal variation in the relationship between NEE and
89 vegetation growth than the Landsat-based 16-daily NDVI data. In contrast, the interpretation of NEE by
90 variables such as soil texture and soil organic content (SOC), which do not have temporal dynamic
91 information, may be limited to the interpretation of spatial variability, although they are considered to be
92 important drivers of NEE. Therefore, the importance of variables obtained from NEE simulations based on
93 a data-driven approach may differ from that in process-based models as well as in the actual driving
94 mechanisms. This may be related to the spatial and temporal resolution of the predictors used and the
95 quality of the data. It is necessary to consider the spatio-temporal resolution of the data for the actual

96 biophysical variables used in the different studies in the systematic evaluation of data-driven NEE
97 simulations.

98 b) The spatio-temporal heterogeneity of data sets, and validation method: The spatio-temporal heterogeneity
99 of the dataset may affect model accuracy. Typically, training data with larger regions, multiple sites,
100 multiple PFTs, and longer spans of years may have a higher degree of imbalance (Kaur et al., 2019; Van
101 Hulse et al., 2007; Virkkala et al., 2021; Zeng et al., 2020). Modeling with unbalanced data (where the
102 difference between the distribution of the training and validation sets is significant even if selected at
103 random) may result in lower model accuracy. To date, the most commonly used methods for validating
104 such models include spatial (Virkkala et al., 2021), temporal (Reed et al., 2021), and random (Cui et al.,
105 2021) cross-validation. The imbalance of data between the training and validation sets may affect the
106 accuracy of the models when using these validation methods. Spatial validation is used to assess the ability
107 of the model to adapt to different regions or flux sites of different PFTs, and a common method is 'leave
108 one site out' cross-validation (Virkkala et al., 2021; Zeng et al., 2020). If the data from the site left out is
109 not covered (or partially covered) by the distribution of the training dataset, the model's prediction
110 performance at that site may be poor due to the absence of a similar type in the training set. Temporal
111 validation typically uses some years of data as training and the remaining years as validation to assess the
112 model's fitness for interannual variability. For a year that is left out (e.g. a special extreme drought year
113 which does not occur in the training set), the accuracy of the model may be limited if there are no similar
114 years (extreme drought years) in the training dataset. K-fold cross-validation is commonly used in random
115 cross-validation to assess the fitness of the model to the spatio-temporal variability. In this case, different
116 values of K may also have a significant impact on the model accuracy. For example, for an unbalanced
117 dataset, the average model accuracy obtained from a 10-fold ($K = 10$) validation approach is likely to be
118 higher than that of a 3-fold ($K = 3$) validation approach (Marcot and Hanea, 2021).

119 c) Machine learning algorithms used: Simulating NEE using different machine learning algorithms may
120 influence the model accuracy, which may be induced by the characteristics of these algorithms themselves
121 and the specific data distribution of the NEE training set. For example, Neural Networks can be used
122 effectively to deal with nonlinearities, while as an ensemble learning method, Random Forests can avoid
123 overfitting due to the introduction of randomness. Therefore, a comprehensive evaluation of this is
124 necessary.

125

126 In this study, to evaluate the impacts of predictors use, algorithms, spatial/temporal scale, and validation
127 methods on model accuracy, we performed a meta-analysis of papers with prediction models that combine NEE
128 observations from flux towers, various predictors, and machine learning for the data-driven NEE simulations. In
129 addition, we also analyzed the causality of multiple features in NEE simulations and the joint effects of multiple
130 features on model accuracy using the Bayesian Network (BN) (a multivariate statistical analysis approach
131 (Pearl, 1985)). The findings of this study can provide some general guidance for future NEE simulations.

132 **2 Methodology**

133 **2.1 Criteria for including articles**

134 In the Scopus database, a literature query was applied to titles, abstracts, and keywords (Table 1) according to
135 Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) (Moher et al., 2009) (Fig. 1):

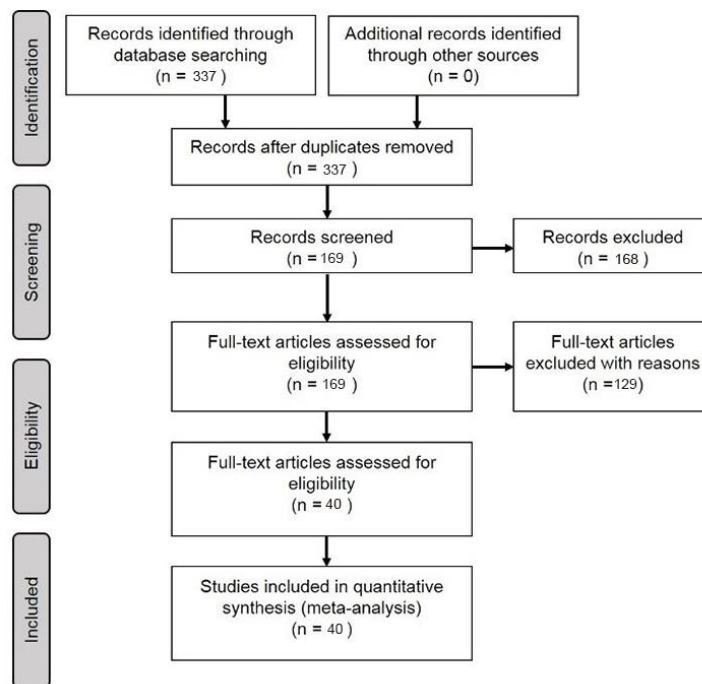
- 136 a) Articles were filtered for those that modeled NEE. Articles that modeled other carbon fluxes such as
137 methane flux were not included.
- 138 b) Articles that used only univariate regression rather than multiple regression were screened out.
- 139 c) Articles reported the determination coefficient (R-squared) of the validation step (Shi et al., 2021;
140 Tramontana et al., 2016; Zeng et al., 2020) as the measure of model performance. Although RMSE is also
141 often used for model accuracy assessment, its dependence on the magnitude of water flux values makes it
142 difficult to use for fair comparisons between studies.
- 143 d) Articles were published in journals with language limited to English.
- 144 e) Articles were filtered for those that were published in the specific journals (Table S1) for research quality
145 control because the data, model implements, and peer review in these journals are often more reliable.

146

147 Table 1. Article search query design: '[A1 OR A2 OR A3...] AND [B1 OR B2...] AND [C1 OR C2...]'

ID	A	B	C
1	Carbon flux	"Eddy covariance"	"machine learning"
2	CO ₂ flux	"Flux tower"	regress*
3	"net ecosystem exchange"		"Support Vector"
4	net ecosystem produc		"Neural Network"
5	gross primary produc		"Random Forest"
6	Carbon exchange		

148



149

150 Figure 1. PRISMA-based paper filtering flowchart.

151 **2.2 Features of prediction models**

152 The information of R-squared (at the validation phase) and the associated model features reported in the article
 153 are considered as one data record for the formal meta-analysis (i.e., each R-squared record corresponding to a
 154 prediction model). From the included papers, R-squared records and various features (Table 2) involved in the
 155 NEE modeling framework (Fig. 2) were extracted (including the used algorithms, modeling/validation methods,
 156 remote sensing data, meteorological data, biophysical data, and ancillary data). In some studies, multiple
 157 algorithms were applied to the same dataset, or models with different features were developed. In these cases,
 158 multiple data records will be documented.

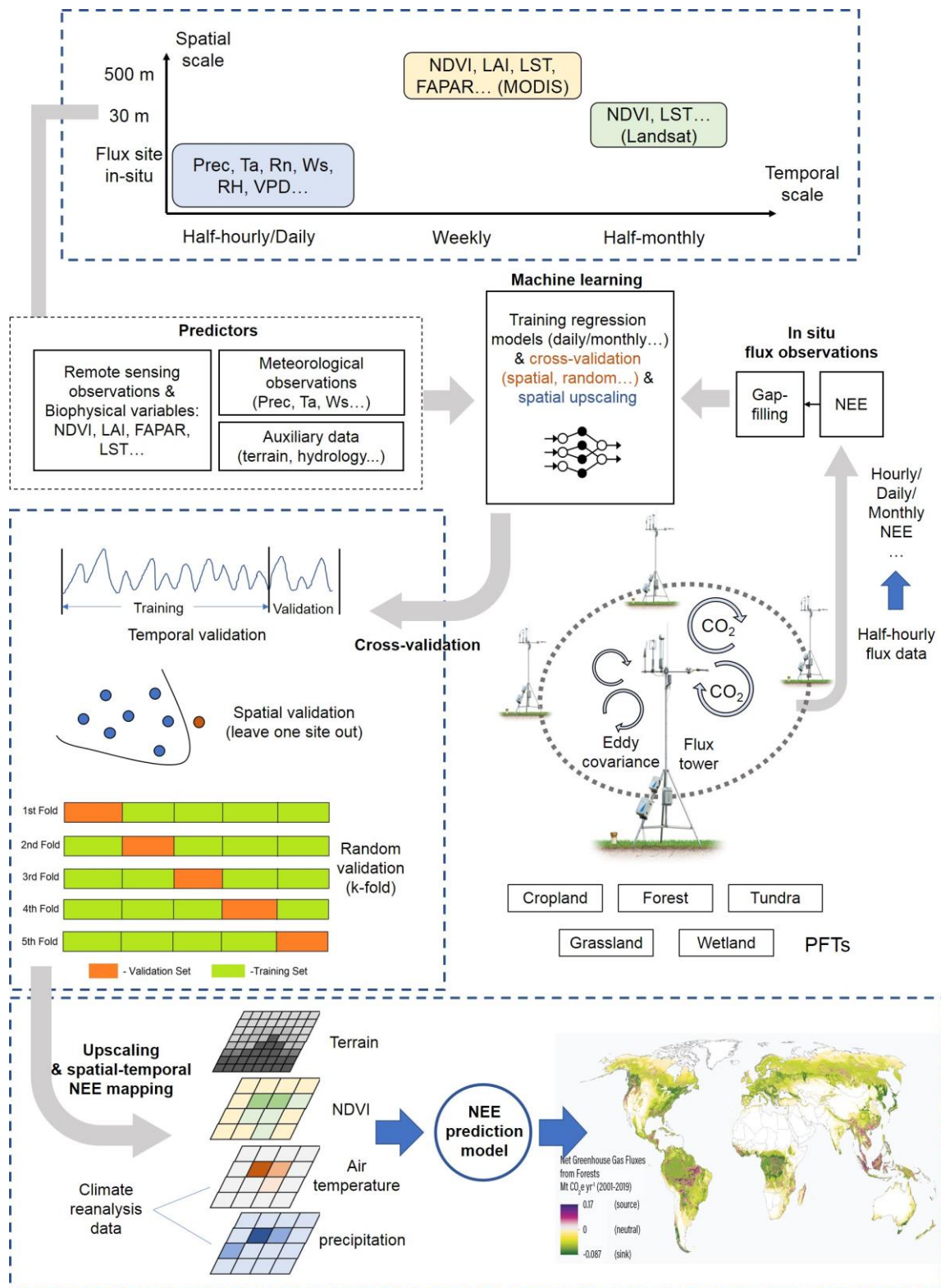
159

160 In the practical information extracting step, we categorized such features in a comparable manner. First, we
 161 categorized the various algorithms used in these papers, although the same algorithm may also have a variant
 162 form or an optimized parameter scheme. They are categorized into the following families of algorithms:
 163 Random Forests (RF), Multiple Linear Regressions (MLR), Artificial Neural Networks (ANN), Support Vector
 164 Machines (SVM), Partial Least Squares Regression (PLSR), Generalized additive model (GAM), Boosted
 165 Regression Tree (BRT), Bayesian Additive Regression Trees (BART), Cubist, model tree ensembles (MTE).
 166 Second, we classified the spatial scales of these studies. Models with study areas (spatial extent covered by flux
 167 stations) smaller than 100x100 km were classified as ‘local’ scale models, those with study area sizes exceeding
 168 continental scale were classified as ‘global’ scale, and those with study area sizes in between were classified as
 169 ‘regional’ scale. Third, for various predictors, we only recorded whether the predictors were used or not without
 170 distinguishing the detailed data sources and categories (e.g., grid meteorological data from various reanalysis
 171 datasets and in-situ meteorological observations from flux stations), measurement methods (e.g., soil moisture

172 measured/estimated by remote sensing or in situ sensors), etc. Fourth, we documented PFTs for the prediction
173 models from the description of study areas or sites in these papers. They are classified into the following types:
174 forest, grassland, cropland, wetland, savannah, tundra, and multi-PFTs (models containing a mixture of multiple
175 PFTs). Models not belonging to the above PFTs were not given a PFT field and were not included in the
176 subsequent analysis of the PFT differences. Other features (Table 2) are extracted directly from the
177 corresponding descriptions in the papers in an explicit manner.

178

179



180

181

182

183

184

185

186

Figure 2. Features of the machine learning-based NEE prediction process. The flux tower photo is from <https://www.licor.com/env/support/Eddy-Covariance/videos/ec-method-02.html> (last accessed: 23rd March 2022). The map in the lower part is from Harris et al., 2021. Prec, Ta, Rn, Ws, RH, and VPD represent precipitation, air temperature, net surface radiation, wind speed, relative humidity, and vapour-pressure deficit respectively. FAPAR is the fraction of absorbed photosynthetically active radiation. LST is the land surface temperature. LAI is the leaf area index.

Field/Feature	Definition	Categories adopted
Id paper	Identification number of the paper (internal)	
Paper	Paper metadata	
Author/s	Name/s of author/s	
Title	Title of the paper	
Year	Year of publication	
Publication title	Name of the journal where the paper was published	
Plant functional type (PFT)	PFTs for the flux sites used	1-forest, 2-grassland, 3-cropland, 4-wetland, 5-savannah, 6-tundra and multi-PFTs
Location	More precise location (with the latitude and longitude of the center of the studied sites). Global (mainly based on FluxNet (Tramontana et al., 2016)) and continental-scale studies are not shown on the map due to the difficulty of identifying specific locations.	latitude, longitude
Algorithms	Algorithm families used in the multivariate regression	Random Forests (RF), Multiple Linear Regressions (MLR), Artificial Neural Networks (ANN), Support Vector Machines (SVM), Partial Least Squares Regression (PLSR), Generalized additive model (GAM), Boosted Regression Tree (BRT), Bayesian Additive Regression Trees (BART), Cubist, model tree ensembles (MTE).
Sites number	Number of the flux sites used	
Study area/Spatial scale	Area representatively covered by the flux sites	local (less than 100×100 km), regional, global (continent-scale and global scale)
Temporal scale	The temporal scale of the model	half-hourly, hourly, daily, weekly, 8-daily, monthly, seasonally, yearly
Study period	The period of the data used in the model	year, growing season, daytime, spring, summer, autumn, winter
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 in training and validation sets.	
Satellite images	Describe the source of satellite images used to derive NDVI, EVI, LAI, LST, etc.	Landsat, MODIS, Hyperion (EO-1), AVHRR, IKONOS
Biophysical predictors	LAI, NDVI/EVI, evapotranspiration (ET) (i.e., the latent heat observed by the flux station), enhanced vegetation index (EVI), the fraction of absorbed photosynthetically active radiation/photosynthetically active radiation (FAPAR/PAR), leaf area index (LAI), 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 source of ancillary variables including terrain variables derived from DEM, soil texture, or hydrology-related data: soil organic content (SOC), soil texture, terrain, 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 in machine learning models.	
Accuracy measure	Accuracy measure used to assess the performance of the estimation/prediction	R-squared (in the validation phase)

189

190 2.3 Bayesian Network for analyzing joint effects

191 Based on the Bayesian network (BN), the joint impacts of multiple model features on the R-squared are
192 analyzed. A BN can be represented by nodes (X_1, \dots, X_n) and the joint distribution (Pearl, 1985):

$$193 \quad P(X) = P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i | pa(X_i)) \quad (1)$$

194 where $pa(X_i)$ is the probability of the parent node X_i . Expectation-maximization (EM) approach (Moon, 1996) is
195 used to incorporate the collected model records and compile the BN.

196

197 Sensitivity analysis is used for the evaluation of node influence based on mutual information (MI) which is
198 calculated as the entropy reduction of the child node resulting from changes at the parent node (Shi et al., 2020):

199
$$MI = H(Q) - H(Q|F) = \sum_q \sum_f P(q, f) \log_2 \left(\frac{P(q, f)}{P(q)P(f)} \right) \quad (2)$$

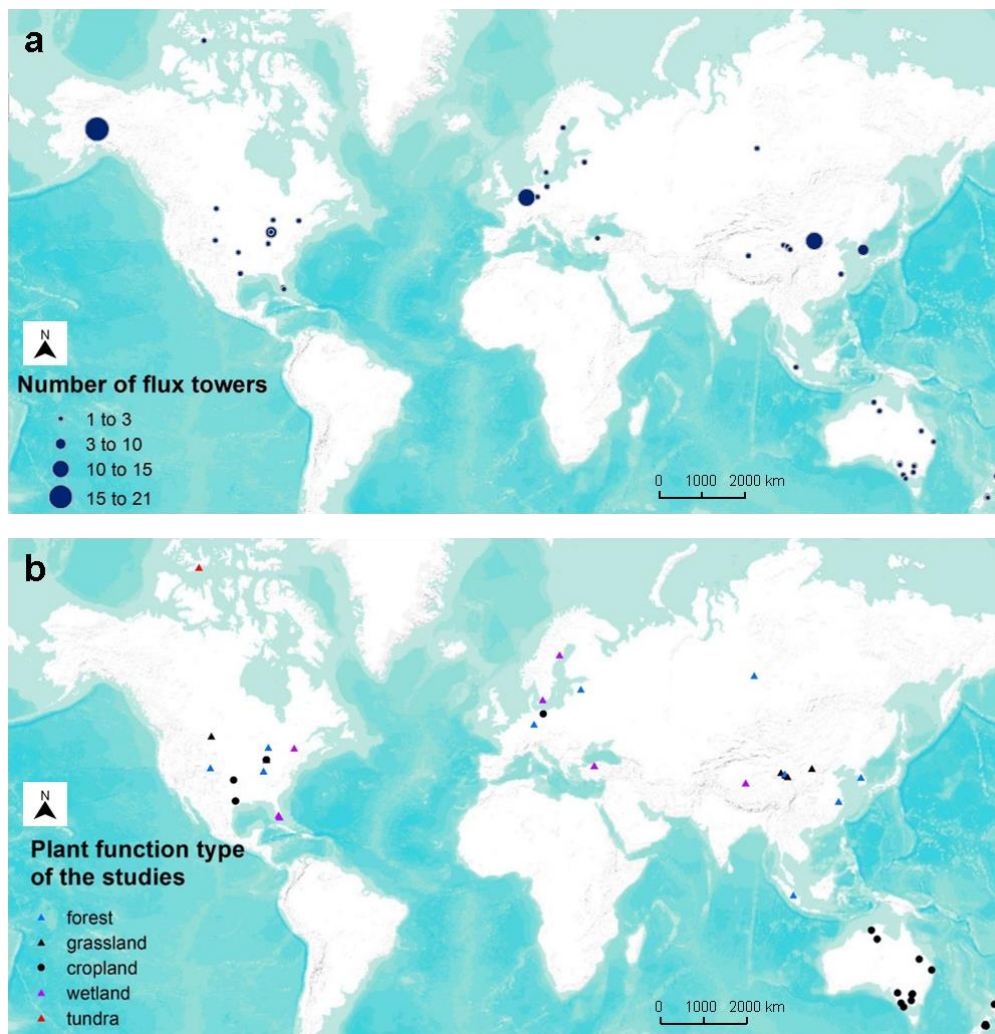
200 where H represents the entropy, Q represents the target node, F represents the set of other nodes and q and f
 201 represent the status of Q and F.

202 **3 Results**

203 **3.1 Articles included in the meta-analysis**

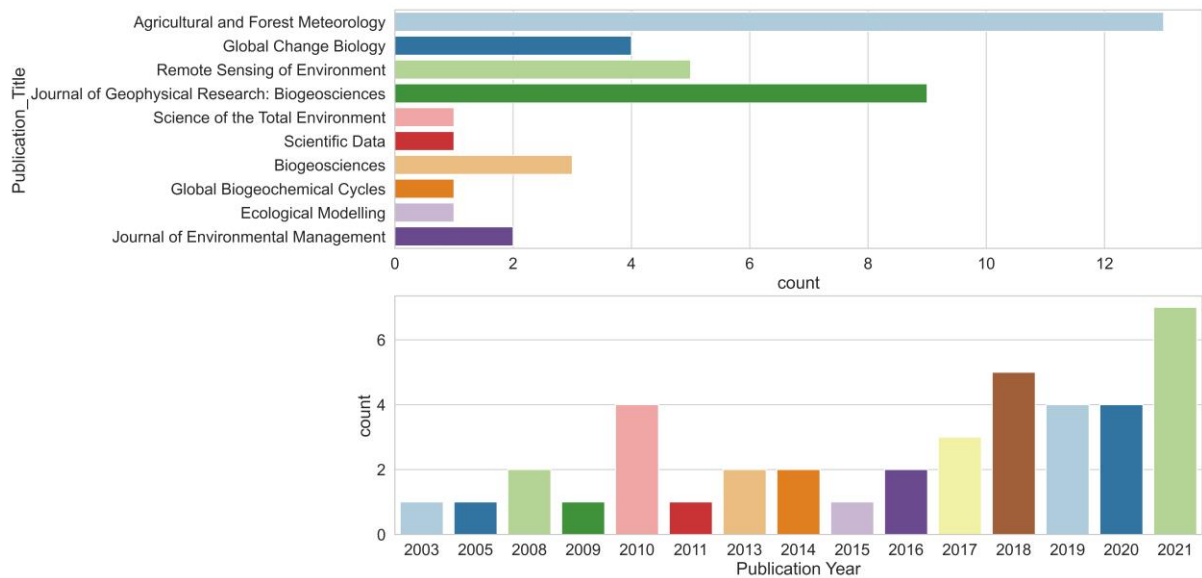
204 We included 40 articles (Table S2) and extracted 178 model records for the formal meta-analysis (Fig. 1). Most
 205 studies were implemented in Europe, North America, Oceania, and China (Fig. 3). The number of such papers is
 206 increasing recently (Fig. 4) and it shows the machine learning approach for NEE prediction has been of interest
 207 to more researchers. The main journals in which these articles have been published (Fig. 4) include Remote
 208 Sensing of Environment, Global Change Biology, Agricultural and Forest Meteorology, Biogeosciences, and
 209 Journal of Geophysical Research: Biogeosciences, etc.

210



211
 212 Figure 3. Location of studies (a) included with the number of flux sites included and (b) their PFTs in the meta-
 213 analysis (total of 40 studies and 178 model records). Global (mainly based on FluxNet (Tramontana et al.,

214 2016)) and continental-scale studies are not shown on the map due to the difficulty of identifying specific
 215 locations.
 216



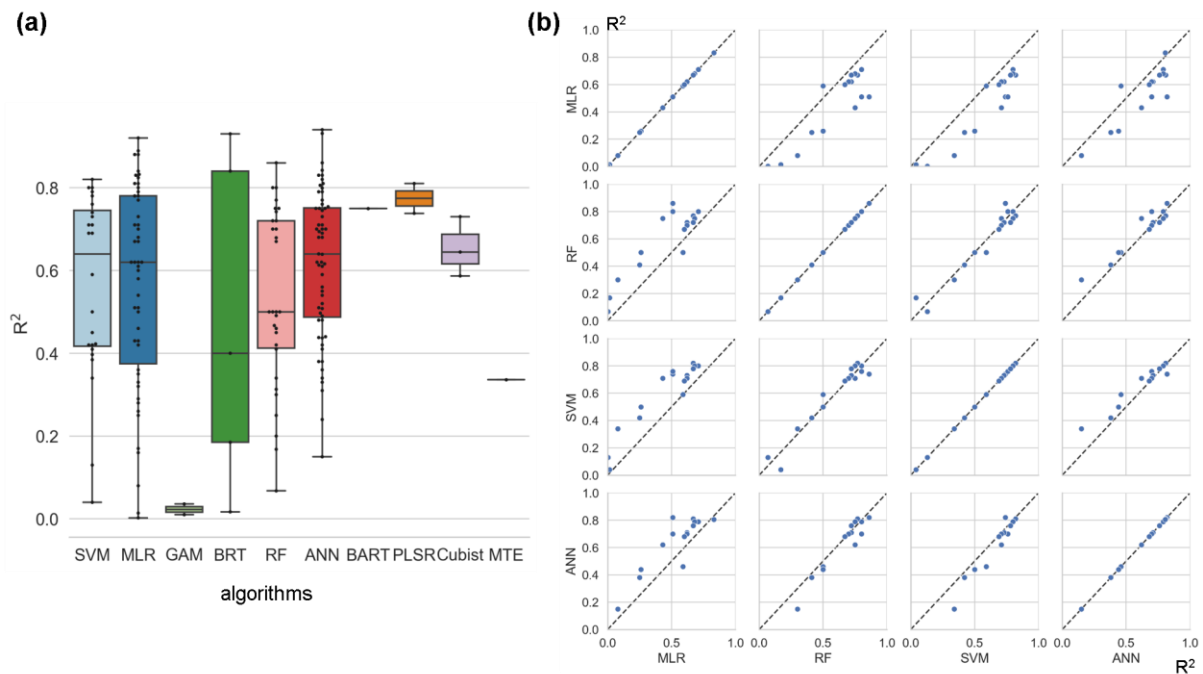
217
 218 Figure 4. The number of studies published across journals and the total number of publications per year.

219 **3.2 The formal Meta-analysis**

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

222 **3.2.1 Algorithms**

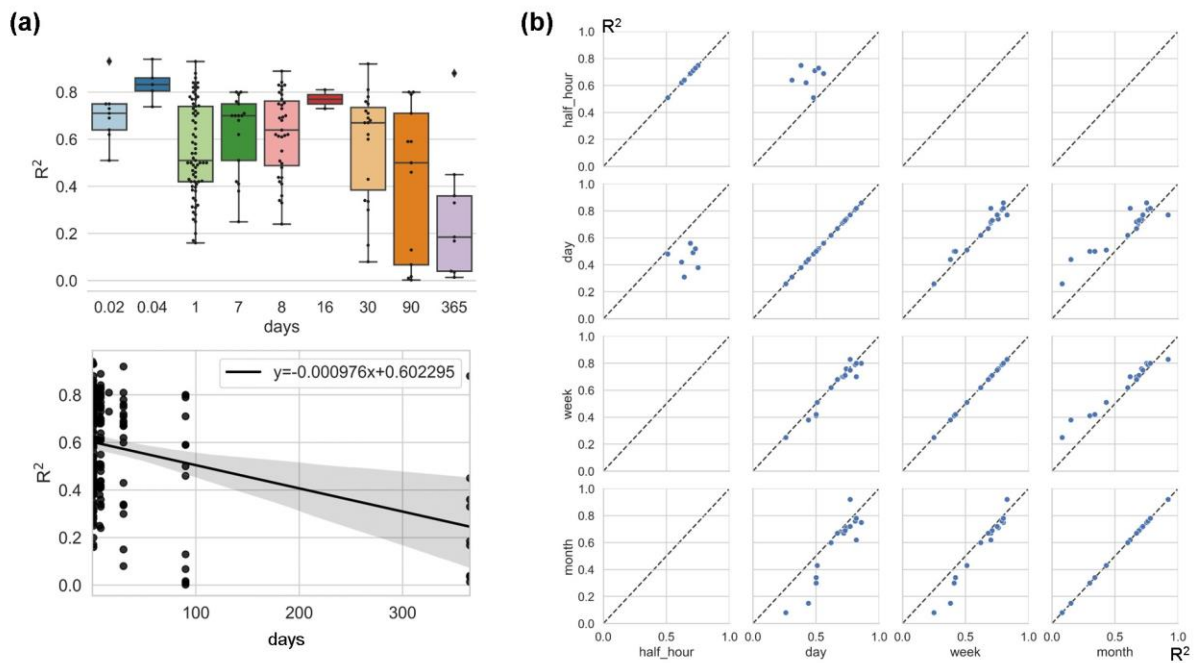
223 Among the more frequently used algorithms, ANN and SVM performed better (Fig. 5a) on average across
 224 studies (lightly better than RF). On the other hand, since cross-study comparisons of algorithm accuracy include
 225 differences in data used in model construction, we performed a pairwise comparison (Fig. 5b) of these four
 226 algorithms (i.e., ANN, SVM, RF, and MLR). In these studies, multiple models are developed for consistent
 227 training data with the interference of training data differences removed. It shows that RF and SVM perform best
 228 in the inter-study comparison (Fig. 5b). Whereas ANN performed slightly worse than RF and SVM, all three of
 229 them were stronger than MLR. Overall, the performance of RF and SVM may be good and similar in the NEE
 230 simulations.



231
 232 Figure 5. Differences in model accuracy (R-squared) using different algorithms across studies (a) and internal
 233 comparisons of the model accuracy (R-squared) of selected pairs of algorithms within individual studies (b).
 234 Regression algorithms: Random Forests (RF), Multiple Linear Regressions (MLR), Artificial Neural Networks
 235 (ANN), Support Vector Machines (SVM), Partial Least Squares Regression (PLSR), Generalized additive
 236 model (GAM), Boosted Regression Tree (BRT), Bayesian Additive Regression Trees (BART), Cubist, model
 237 tree ensembles (MTE). In panel (a), the horizontal line in the box indicates the medians. The top and bottom
 238 border lines of the box indicate the 75% and 25% percentiles, respectively.

239 3.2.2 Time scales

240 The impact of time scale on R-squared is considerable (Fig. 6), with models with larger time scales having
 241 lower average R-squared, especially when the time scale exceeds the monthly scale. The most frequently used
 242 scales were the daily, 8-day, and monthly scales. In studies where multiple time scales were used with other
 243 characteristics being the same, we found that models with half-hourly scales were significantly more accurate
 244 than models with daily scales (Fig. 6). However, the difference in accuracy between the day-scale and week-
 245 scale models is small. The accuracy of models with a monthly scale is the lowest.



246

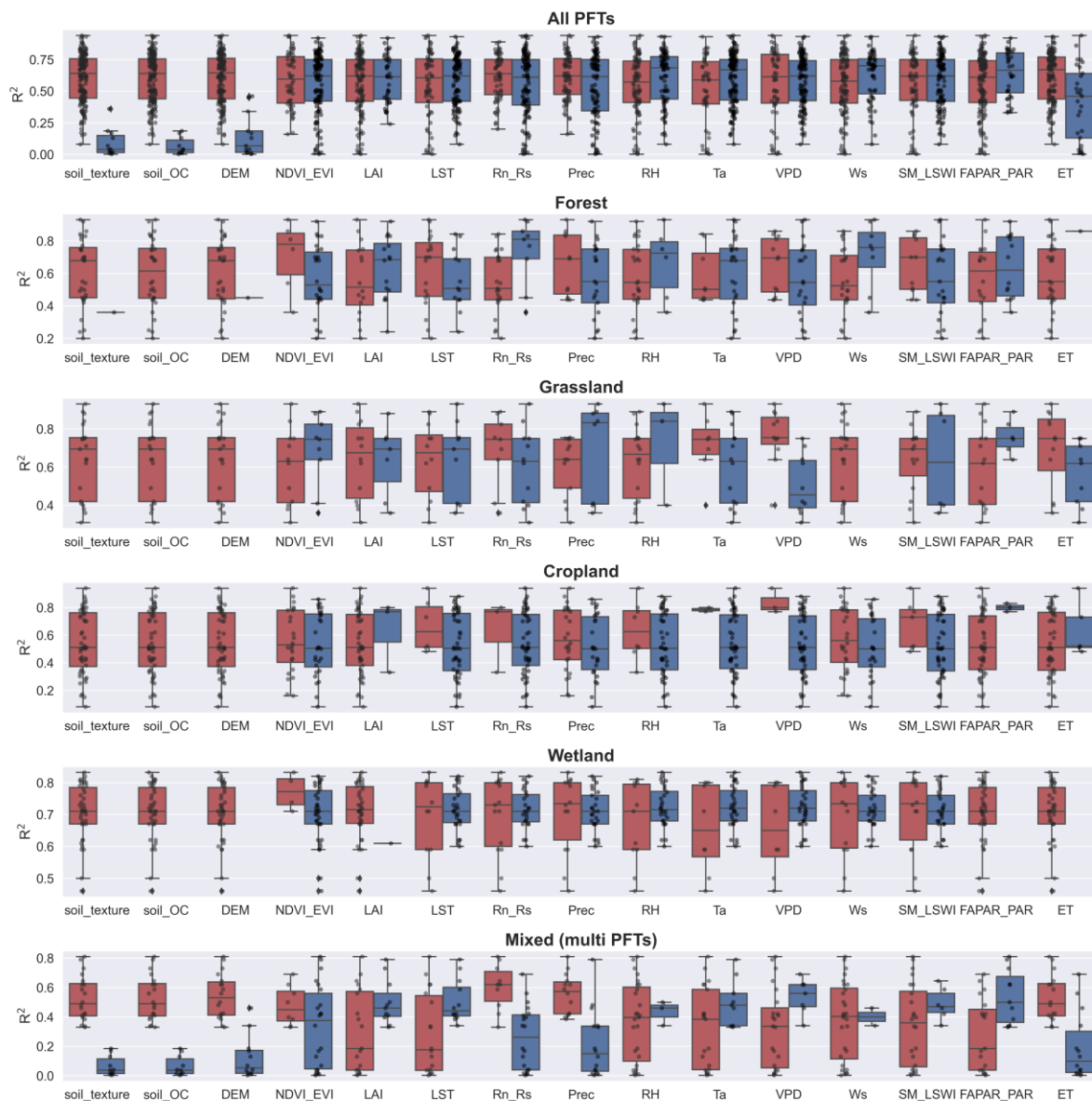
247 Figure 6. Differences in model accuracy (R-squared) at different time scales across studies with the
 248 linear regression between R-squared and time scales (a), and comparison of the model accuracy (R-
 249 squared) of selected pairs of time scales within individual studies (b). All model records were
 250 included in panel (a), while studies that used multiple time scales (with other model characteristics
 251 unchanged) were included in panel (b). Time scales: 0.02 days (half-hourly), 0.04 days (hourly), 30
 252 days (monthly), and 90 days (quarterly).

253 3.2.3 Various predictors

254 Among the commonly used predictors for NEE, there are significant differences in the predictors used and their
 255 impacts on model accuracy for different PFTs (Fig. 7). Ancillary data (e.g. soil texture, soil organic content,
 256 topography) that do not have temporal variability are used less frequently because they can only explain spatial
 257 heterogeneity. In contrast, the biophysical variables LAI, FAPAR, and ET were used significantly less
 258 frequently than NDVI/EVI, especially in the cropland and wetland types. The meteorological variables Ta,
 259 Rn/Rs, and VPD were used most frequently. For forest sites, Rn/Rs and Ws appear to be the variables that
 260 improve model accuracy. For grassland sites, we found that NDVI/EVI appears to be the most effective, despite
 261 the small sample size. For sites in croplands and wetlands, we did not find predictor variables that had a
 262 significant impact on model accuracy.

263

264 For different PFTs, the top three variables in the ranking of model importance differed (Fig. S1). SM, Rn/Rs,
 265 Ta, Ts, and VPD all showed high importance across PFTs. This suggests that the variability of measured site-
 266 scale moisture and temperature conditions is important for the simulation of NEE for all PFTs. In contrast, in the
 267 importance ranking, other variables such as precipitation and NDVI/EVI may not lead because of the lag in their
 268 effect on NEE (Hao et al., 2010; Cranko Page et al., 2022). And some other variables may improve model
 269 accuracy for specific PFTs such as groundwater table depth (GWT) for wetland sites and growing degree days
 270 (GDD) for tundra sites.



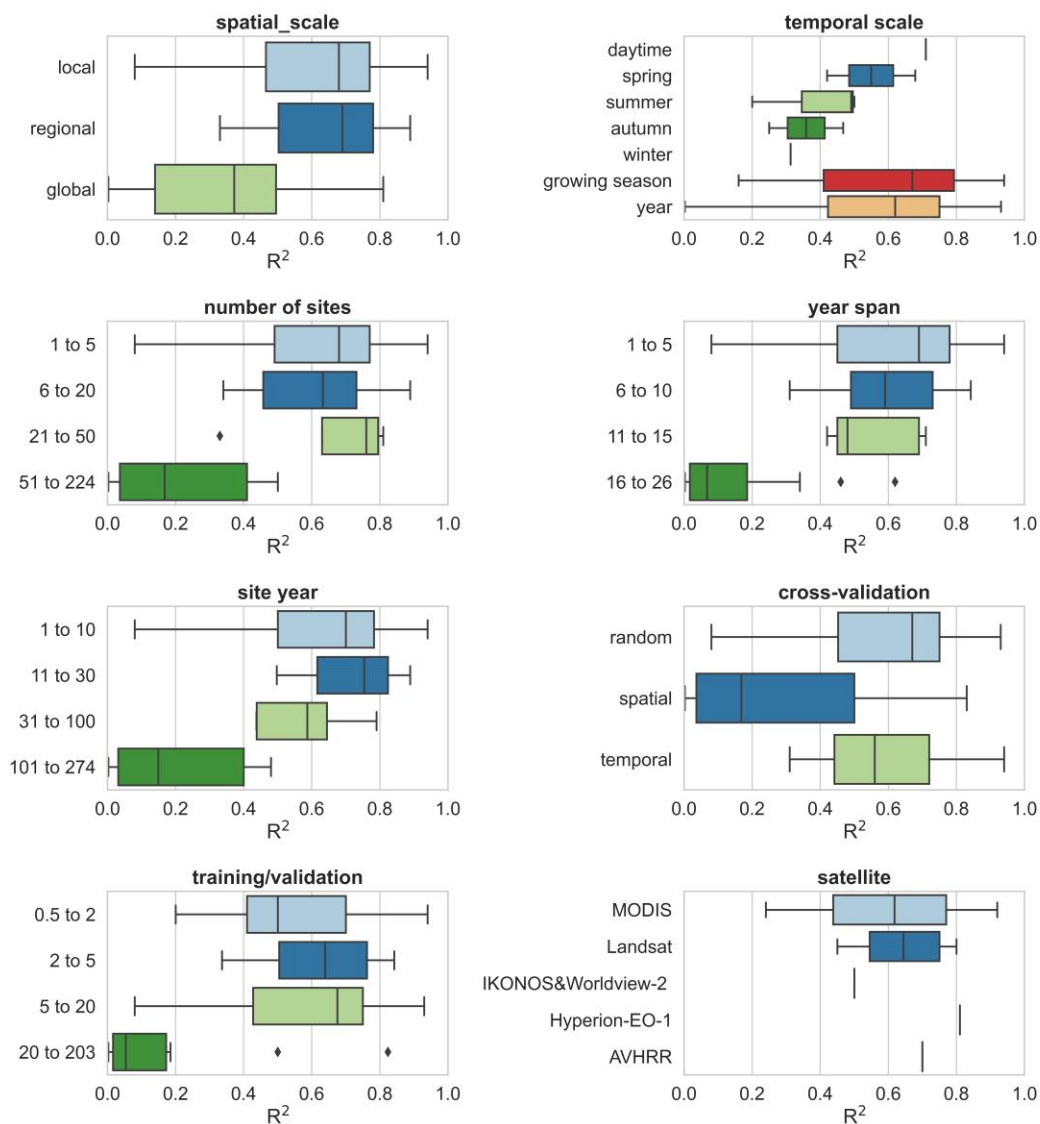
272

273 Figure 7. The impact of the various predictors incorporated in models of different PFTs (1-forest, 2-grassland, 3-
 274 cropland, 4-wetland, 6-tundra) on R-squared. Dark blue boxes indicate that the predictor was used in the model,
 275 while dark red boxes indicate that the predictor was not used. Predictors: soil organic content (Soil_OC),
 276 precipitation (Prec), soil moisture/land surface water index (SM_LSWI), net radiation/solar radiation (Rn_Rs),
 277 enhanced vegetation index (EVI), air temperature (Ta), vapor-pressure deficit (VPD), the fraction of absorbed
 278 photosynthetically active radiation/photosynthetically active radiation (FAPAR_PAR), relative humidity (RH),
 279 evapotranspiration (ET), leaf area index (LAI).

280 3.2.4 Other features

281 In addition, we evaluated other features of the model construction that may contribute to differences in model
 282 accuracy (Fig. 8). Studies at continental and global scales with a large number of sites and a large span of years
 283 correspond to lower R-squared than studies at local and regional scales, suggesting that studies with a large
 284 number of sites across large regions are likely to have high variability in the relationship between NEE and

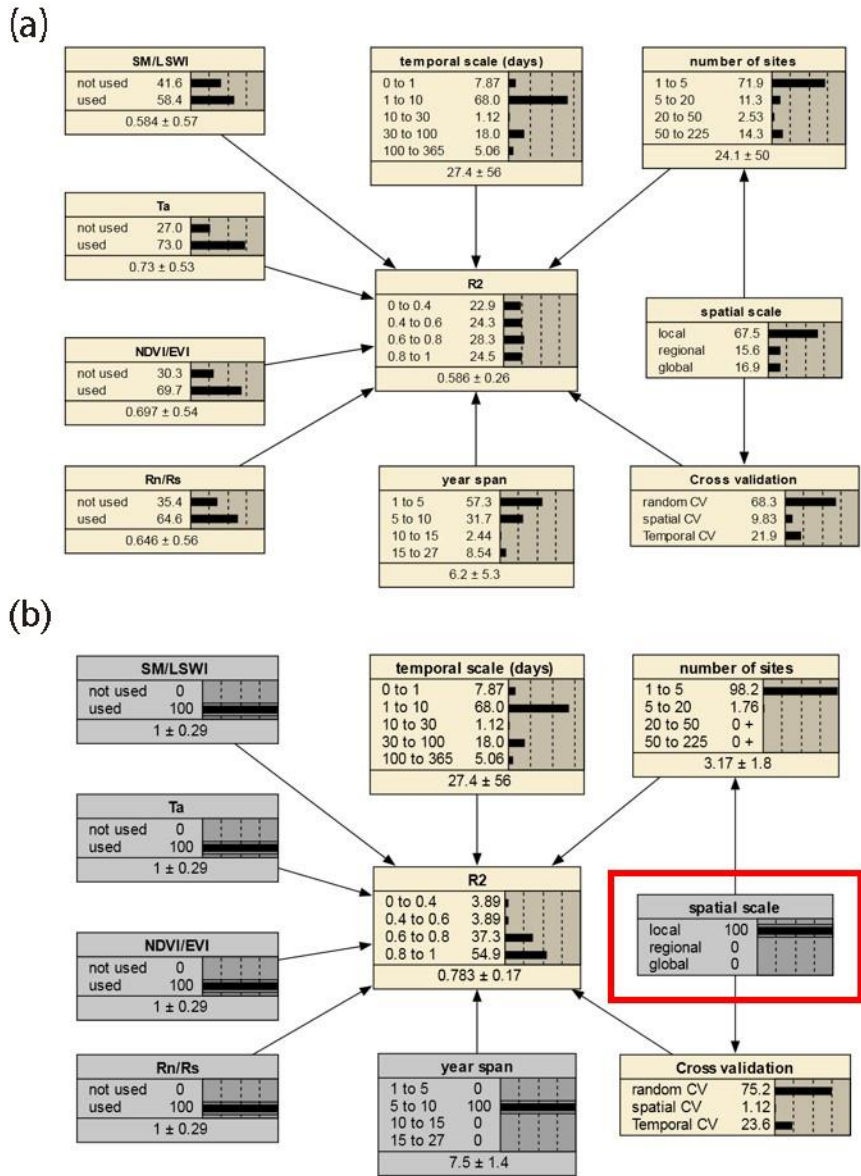
285 covariates and that studies at small scales are more likely to have higher model accuracy. Spatial validation
 286 (usually 'leave one site out') corresponds to lower model accuracy compared to random and temporal validation.
 287 This again confirms the dominant role of heterogeneity in the relationship between NEE and covariates across
 288 sites in explaining model accuracy. This seems to be indirectly supported by the fact that a high ratio of training
 289 to validation sets corresponds to a low R-squared, as this high ratio tends to be accompanied by the use of the
 290 'leave one site out' validation approach. The accuracy of the models with a growing season period was slightly
 291 higher than that of the models with an annual period. For the satellite remote sensing data used, the models
 292 based on MODIS data with biophysical variables extracted were slightly less accurate than those based on
 293 Landsat data. For the daily scale models, Landsat data performed a little better than MODIS (Fig. S2). This
 294 suggests that the higher temporal resolution of MODIS compared to Landsat may not play a dominant role in
 295 improving model accuracy. This may also be partially attributed to studies using MODIS-based explanatory data
 296 that tend to include too large surrounding areas around the site (e.g., 2x2 km), which can lead to a scale
 297 mismatch between the flux footprint and the explanatory variables.



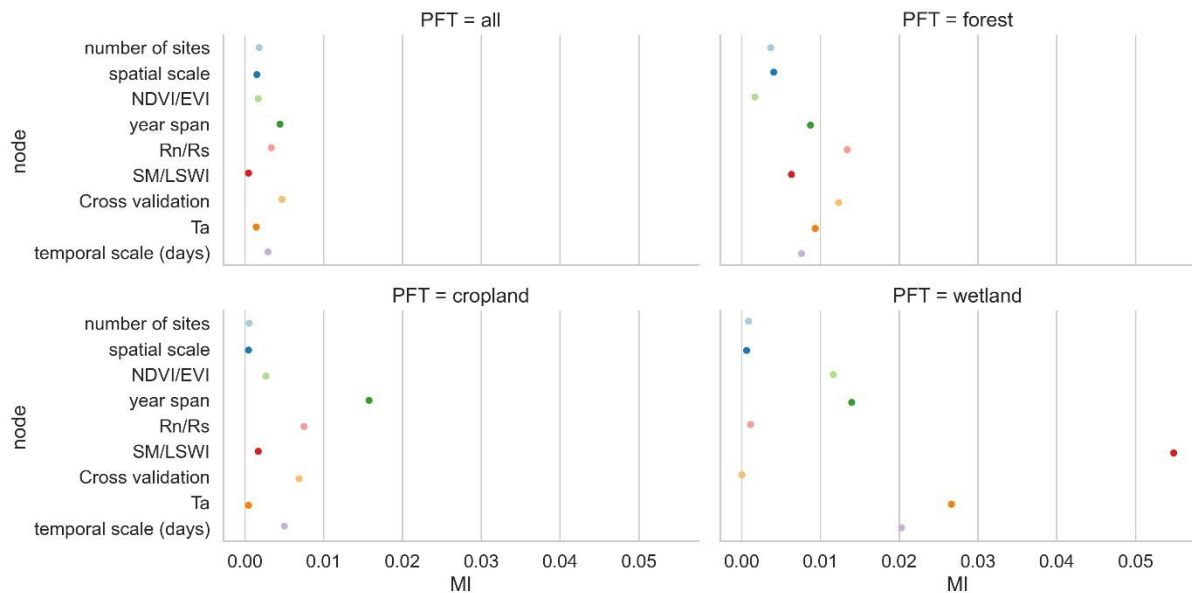
298
 299 Figure 8. The impacts of other features (i.e. spatial scale, study period, number of sites, year span, site year,
 300 cross-validation method, training/validation, and satellite imagery) on the model performance.

301 **3.3. The joint causal impacts of multi-features based on the BN**

302 We selected the features that had a more significant impact on model accuracy in the above assessment and
303 further incorporated them into the BN-based multivariate assessment to understand the joint impact of multiple
304 features on R-squared. The features incorporated included the spatial scale, the number of sites, the temporal
305 scale, the span of years, the cross-validation method, and whether some specific predictors were used. We
306 discretized the distribution of individual nodes and compiled the BN (Fig. 9.a) using records from different
307 PFTs as input. Sensitivity analysis of the R-squared node (Fig. 10) showed that R-squared was most sensitive to
308 'year span', cross-validation method, Rn/Rs, and time scale under multi-feature control. In the forest and
309 cropland types, R-squared is more sensitive to Rn/Rs, while in the wetland type it is more sensitive to SM/LSWI
310 and Ta. The sensitivity of R-squared to 'year span' was much higher in the cropland type compared to the other
311 PFTs, which may suggest that the interannual variability in the NEE simulations of the cropland type is higher
312 due to potential interannual variability of the planting structure and irrigation practices. For the cropland type,
313 differences in the phenology, harvesting, and irrigation (water volume and frequency) in different years can lead
314 to significant inter-annual differences in NEE simulations. Subsequently, using the constructed BN (with the
315 empirical information in previous studies incorporated), for new studies we can instructively infer the
316 probability distribution of the possible R-squared (Fig. 9.b) with some model features predetermined. In
317 previous studies, spatio-temporal mapping of NEE based on statistical models has often lacked accuracy
318 assessment since there are no grid-scale NEE observations, and this BN may have the potential to be used to
319 validate the accuracy (R-squared) of the NEE time series output of the grid-scale (i.e. inferring possible R-
320 squared from model features, where the output of the grid-scale is considered to be of the form 'leave one site
321 out').



322
 323 Figure 9. The joint effects of multiple features on the R-squared based on the BN with all records input (a) and
 324 the inference on the probability distribution of R-squared based on the BN with the status of some nodes
 325 determined (b). The values before and after the “±” indicate the mean and standard deviation of the distribution,
 326 respectively. The gray boxes indicate that the status of the nodes has been determined. In panel (b), specific
 327 values of parent nodes such as ‘spatial scale’ are determined (shown in the red box), leading to an increase in the
 328 expected R-squared compared to the average scenario of the panel (a) (as inferred from the posterior conditional
 329 probabilities with the status of the node ‘spatial scale’ are determined as ‘local’).
 330



331

332 Figure 10. The sensitivity analysis of the R-squared node to other nodes based on the mutual information (MI)
 333 across PFTs. ‘Cross-validation’ is the cross-validation method including spatial, temporal, and random cross-
 334 validation.

335 4 Discussions

336 Many studies have evaluated the incorporation of various predictors and model features using machine learning
 337 for improving the site-scale NEE predictions (Tramontana et al., 2016; Zeng et al., 2020; Jung et al., 2011). A
 338 comprehensive evaluation of these studies to provide definitive guidance on the selection of features in NEE
 339 prediction modeling is limited. This study fills the research gap with a meta-analysis of the literature through
 340 statistics on the accuracy and performance of models. Machine learning-based NEE simulations and predictions
 341 still suffer from high uncertainty. By better understanding the expected improvements that can be achieved
 342 through the inclusion of different features, we can identify priorities for the consideration of different features in
 343 modeling efforts and avoid operations decreasing model accuracy.

344

345 Compared to previous comparisons of machine learning-based NEE prediction models, this study is more
 346 comprehensive. Previous studies (Abbasian et al., 2022) have also found advantages of RF over other
 347 algorithms in NEE prediction. This study consolidated this finding using a larger amount of evidence. Previous
 348 studies (Tramontana et al., 2016) have also compared the impact of different practices in NEE prediction models
 349 based on the R-squared, such as comparing the difference in accuracy between the two predictor combinations
 350 (i.e., using only remotely sensed data and using remotely sensed data and meteorological data together). In
 351 contrast, since this study incorporated more detailed factors influencing model accuracy, the understanding of
 352 such issues was deepened. However, there are still many uncertainties and challenges in NEE prediction not
 353 clarified in this study.

354 **4.1 Challenges in the site-scale NEE simulation and implications for other carbon flux simulations**

355 **4.1.1 Variations in time scales**

356 In the above analysis, we found that the effect of the time scale of the model is considerable. This suggests that
357 we should be careful in determining the time scale of the model to consider whether the predictor variables used
358 will work at this time scale. Previous studies have reported the dependence of the NEE variability and
359 mechanism on the time scales. On the one hand, the importance of variables affecting NEE varies at different
360 time scales. For example, in tropical and subtropical forests in southern China (Yan et al., 2013), seasonal NEE
361 variability is predominantly controlled by soil temperature and moisture, while interannual NEE variability is
362 controlled by the annual precipitation variation. A study (Jung et al., 2017) showed that for annual-scale NEE
363 variability, water availability and temperature were the dominant drivers at the local and global scales,
364 respectively. This indicates the need to recognize the temporal and spatial driving mechanisms of NEE in
365 advance in the development of NEE prediction models. On the other hand, dependence may exist between NEE
366 anomalies at various time scales. For example, previous studies (Luyssaert et al., 2007) showed that short-term
367 temperature anomalies may interpret both the daily and seasonal NEE anomalies. This implies that the models at
368 different time scales may not be independent. In the previous studies, the relationship between prediction
369 models at different scales has not been well investigated, and it may be valuable to compare the relations
370 between data and models at different scales in depth. Larger time scales correspond to lower model accuracy,
371 possibly related to the fact that some small-time-scale relations between NEE and covariates (especially
372 meteorological variables) are smoothed. In particular, for models with time scales smaller than one day (e.g.
373 half-hourly models), the 8-daily and 16-daily biophysical variable data obtained from satellite remote sensing
374 are difficult to explain the temporal variation in the sub-daily NEE. Therefore, for models at small time scales
375 (i.e. half-hourly, hourly, daily scale models), in situ meteorological variables may be more important. The
376 inclusion of some ancillary variables (e.g. soil texture, topographic variables) with no temporal dynamic
377 information may be ineffective unless many sites are included in the model and the spatial variability of the
378 ancillary variables for these sites is sufficiently large (Virkkala et al., 2021).

379

380 In terms of completeness and purity of training data, hourly and daily models can be better compared to monthly
381 and yearly models. Hourly and daily models can usually preclude those low-quality data and gaps in the flux
382 observations. However, for monthly and yearly scale models, gap-filling (Ruppert et al., 2006; Moffat et al.,
383 2007; Zhu et al., 2022) is necessary because there are few complete and continuous fluxes observations without
384 data gaps on the monthly to yearly scales. Since various gap-filling techniques rely on environmental factors
385 (Moffat et al., 2007) such as meteorological observations, this may introduce uncertainty in the predictive
386 models (i.e., a small fraction of the observed information of NEE is estimated from a combination of
387 independent variables). How it would affect the accuracy of prediction models at various time scales remains
388 uncertain, although various gap-filling techniques have been widely used in the pre-processing of training data.

389

390 In addition, the impacts of lagged effects (Hao et al., 2010; Cranko Page et al., 2022) of covariates are not
391 considered in most models, which may underestimate the degree of explanation of NEE for some predictor
392 variables (e.g. precipitation). Most of the machine learning-based models use only the average Ta and do not
393 take into account the maximum temperature, minimum temperature, daily difference in temperature, etc., as in

394 the process-based ecological models (Mitchell et al., 2009). This suggests that the inclusion of different
395 temporal characteristics of individual variables in machine learning-based NEE prediction models may be
396 insufficient.

397 **4.1.2 Scale mismatch of explanatory predictors and flux footprints**

398 An excessively large extraction area of remote sensing data (e.g., 2x2 km) may be inappropriate. In the non-
399 homogeneous underlying conditions, the agreement of the area of flux footprints with the scale of the predictors
400 should be considered in the extraction of the predictor variables in various PFTs (Chu et al., 2021).

401

402 The effects of this mismatch between explanatory variables and flux footprints may be diverse for different
403 PFTs. For example, for cropland types, the NEE is monitored at a range of several hundred meters around the
404 flux towers, but remote sensing variables such as FAPAR, NDVI, LAI, etc. can be extracted at coarse scales
405 (e.g., 2x2 km), some effects outside the extent of the flux footprint (Chu et al., 2021; Walther et al., 2021) are
406 incorporated (e.g., planting structures with high spatial heterogeneity, agricultural practices such as irrigation).
407 And for more homogeneous types such as grasslands, coarse-scale meteorological data may still cause spatial
408 mismatches, even though the differences in land cover types within the 2x2 km and 200x200 m extent around
409 the flux stations in grasslands may not be considerable. For example, precipitation with high spatial
410 heterogeneity can dominate the spatial variability of soil moisture and thus affect the spatial variability of
411 grassland NEE (Wu et al., 2011; Jongen et al., 2011). However, using 0.25°x0.25° reanalysis precipitation data
412 (Zeng et al., 2020) may make it difficult for predictive models to capture this spatial heterogeneity around the
413 flux station.

414

415 Since few of the studies included in this meta-analysis considered the effect of variation in flux footprint, this
416 feature was difficult to consider in this study. However, its influence should still be further investigated in future
417 studies. With flux footprints calculated (Kljun et al., 2015) and the factors around the flux site (Walther et al.,
418 2021) that affect the flux footprint incorporated, it is promising to clarify this issue.

419 **4.1.3 Possible unbalance of training and validation sets**

420 In addition to the time scale of the models, the most significant differences in model accuracy and performance
421 were found in the heterogeneity within the NEE dataset and the match of the training set and validation set.

422 Often NEE simulations can achieve high accuracy in local studies, where the main factor negatively affecting
423 model accuracy may be the interannual variability in the relationship between NEE and covariates. However,
424 the complexity may increase when the dataset contains a large study area, many sites, PFTs, and year spans.

425 Under this condition, the accuracy of the model in the 'leave one site out' validation may be more dependent on
426 the correlation and match between the training and validation sets (Jung et al., 2020). When the model is applied
427 to an outlier site (of which the NEE, covariates, and their relationship are very different compared with the
428 remaining sites), it appears to be difficult to achieve a high prediction accuracy (Jung et al., 2020). If we further
429 upscale the prediction model to large spatial and temporal scales, the uncertainties involved may be difficult to
430 assess (Zeng et al., 2020). We can only infer the possible model accuracy based on the similarity of the
431 distribution of predictors in the predicted grid to that of the existing sites in the model. In the upscaling process,

432 reanalysis data with the coarse spatial resolution are often used as an alternative for site-scale meteorological
433 predictors. However, most studies did not assess in detail the possible errors associated with spatial mismatches
434 in this operation.

435

436 In summary, the site-scale NEE predictions may require more focus on the internal heterogeneity of the NEE
437 dataset and the matching of the training set and validation set, and also require a better understanding of the
438 influence of different scales of the same variable (e.g. site-scale precipitation and grid-scale precipitation in the
439 reanalysis meteorological data) across modeling and upscaling steps. For the prediction of other carbon fluxes
440 such as methane fluxes (in the same framework as the NEE predictions), the results of this study may also be
441 partially applicable, although there may be significant differences in the use of specific predictors (Peltola et al.,
442 2019).

443 **4.2 Uncertainties**

444 The uncertainties in this analysis may include:

- 445 a) **Publication bias and weighting:** Publication bias is not refined due to the limitations of the number of
446 articles that can be included. Meta-analyses often measure the quality of journals and the data availability
447 (Borenstein et al., 2011; Field and Gillett, 2010) to determine the weighting of the literature in a
448 comprehensive assessment. However, a high proportion of the articles in this study did not make flux
449 observations publicly available or share the NEE prediction models developed. Furthermore, meta-analysis
450 studies in other fields typically measure the impact of papers by evidence/data volume, and the variance of
451 the evaluated effects (Adams et al., 1997; Don et al., 2011; Liu et al., 2018). However, in this study,
452 because no convincing method is found to quantify the weights of results from included articles, some
453 features (e.g. the number of flux sites, the span of years) were directly assessed rather than used to
454 determine the weights of the articles.
- 455 b) **Limitations of the criteria for inclusion in the literature:** in the model accuracy-based evaluation, we
456 selected only literature that developed multiple regression models. Potentially valuable information from
457 univariate regression models was not included. In addition, only papers in high-quality English journals
458 were included in this study to control for possible errors due to publication bias. However, many studies
459 that fit this theme may have been published in other languages or other journals.
- 460 c) **Independence between features:** There is dependence between the evaluated features (e.g. the dependency
461 between the spatial extent and the number of sites). It may negatively affect the assessment of the impact
462 of individual features on the accuracy of the model, although the BN-based analysis of joint effects can
463 reduce the impact of this dependence between variables by specifying causal relationships between
464 features. The interference of unknown dependencies between features may still not be eliminated when we
465 focus on the effects of an individual feature on the model performance. The sample size collected in this
466 study (178 records in total) is not very large. This also suggests that more future efforts should be devoted
467 to the comprehensive evaluation and summarization of NEE simulations.

468

469 Additionally, there are still other potential factors not considered by this study such as the uncertainty of climate
470 data (site vs reanalysis), footprint matching between site and satellite images, etc. Overall, although the

471 quantitative results of this study should be used with caution, they still have positive implications for guiding
472 future such studies.

473 **5 Conclusion**

474 We performed a meta-analysis of the site-scale NEE simulations combining in situ flux observations,
475 meteorological, biophysical, and ancillary predictors, and machine learning. The impacts of various features
476 throughout the modeling process on the accuracy of the model were evaluated. The main findings of this study
477 include:

- 478 1. RF and SVM performed better than other evaluated algorithms.
- 479 2. The impact of time scale on model performance is significant. Models with larger time scales have lower
480 average R-squared, especially when the time scale exceeds the monthly scale. Models with half-hourly
481 scales (average R-squared = 0.73) were significantly more accurate than models with daily scales (average
482 R-squared = 0.5).
- 483 3. Among the commonly used predictors for NEE, there are significant differences in the predictors used and
484 their impacts on model accuracy for different PFTs.
- 485 4. It is necessary to focus on the potential imbalance between the training and validation sets in NEE
486 simulations. Studies at continental and global scales (average R-squared = 0.37) with multiple PFTs, more
487 sites, and a large span of years correspond to lower R-squared than studies at local (average R-squared =
488 0.69) and regional scales (average R-squared = 0.7).

489

490 **Acknowledgments**

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

495 **Contributions**

496 H.S and G.L initiated this research and were responsible for the integrity of the work as a whole. H.S performed
497 formal analysis, and calculations and drafted the manuscript. H.S, G.L, X.M, X.Y, Y.W, W.Z, M.X, C.Z, and
498 Y.Z were responsible for the data collection and analysis. G.L, P.D.M, T.V.D.V, O.H, and A.K contributed
499 resources and financial support.

500 **Competing interests**

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

502 **Data availability**

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

505 **Code availability**

506 The code used in this study can be accessed by contacting the first author (shihaiyang16@mails.ucas.ac.cn)
507 based on a reasonable request.

508

509

510 **References**

- 511 Abbasian, H., Solgi, E., Mohsen Hosseini, S., and Hossein Kia, S.: Modeling terrestrial net ecosystem
512 exchange using machine learning techniques based on flux tower measurements, *Ecological*
513 *Modelling*, 466, 109901, <https://doi.org/10.1016/j.ecolmodel.2022.109901>, 2022.
- 514 Adams, D. C., Gurevitch, J., and Rosenberg, M. S.: Resampling tests for meta - analysis of ecological
515 data, *Ecology*, 78, 1277–1283, 1997.
- 516 Baldocchi, D. D.: Assessing the eddy covariance technique for evaluating carbon dioxide exchange
517 rates of ecosystems: past, present and future, 9, 479–492, [https://doi.org/10.1046/j.1365-](https://doi.org/10.1046/j.1365-2486.2003.00629.x)
518 [2486.2003.00629.x](https://doi.org/10.1046/j.1365-2486.2003.00629.x), 2003.
- 519 Berryman, E. M., Vanderhoof, M. K., Bradford, J. B., Hawbaker, T. J., Henne, P. D., Burns, S. P.,
520 Frank, J. M., Birdsey, R. A., and Ryan, M. G.: Estimating soil respiration in a subalpine landscape
521 using point, terrain, climate, and greenness data, *Journal of Geophysical Research: Biogeosciences*,
522 123, 3231–3249, 2018.
- 523 Borenstein, M., Hedges, L. V., Higgins, J. P., and Rothstein, H. R.: *Introduction to meta-analysis*,
524 John Wiley & Sons, 2011.
- 525 Cho, S., Kang, M., Ichii, K., Kim, J., Lim, J.-H., Chun, J.-H., Park, C.-W., Kim, H. S., Choi, S.-W.,
526 and Lee, S.-H.: Evaluation of forest carbon uptake in South Korea using the national flux tower
527 network, remote sensing, and data-driven technology, *Agricultural and Forest Meteorology*, 311,
528 108653, 2021.
- 529 Chu, H., Luo, X., Ouyang, Z., Chan, W. S., Dengel, S., Biraud, S. C., Torn, M. S., Metzger, S.,
530 Kumar, J., Arain, M. A., Arkebauer, T. J., Baldocchi, D., Bernacchi, C., Billesbach, D., Black, T. A.,
531 Blanken, P. D., Bohrer, G., Bracho, R., Brown, S., Brunzell, N. A., Chen, J., Chen, X., Clark, K.,
532 Desai, A. R., Duman, T., Durden, D., Fares, S., Forbrich, I., Gamon, J. A., Gough, C. M., Griffis, T.,
533 Helbig, M., Hollinger, D., Humphreys, E., Ikawa, H., Iwata, H., Ju, Y., Knowles, J. F., Knox, S. H.,
534 Kobayashi, H., Kolb, T., Law, B., Lee, X., Litvak, M., Liu, H., Munger, J. W., Noormets, A., Novick,
535 K., Oberbauer, S. F., Oechel, W., Oikawa, P., Papuga, S. A., Pendall, E., Prajapati, P., Prueger, J.,
536 Quinton, W. L., Richardson, A. D., Russell, E. S., Scott, R. L., Starr, G., Staebler, R., Stoy, P. C.,
537 Stuart-Haëntjens, E., Sonnentag, O., Sullivan, R. C., Suyker, A., Ueyama, M., Vargas, R., Wood, J.
538 D., and Zona, D.: Representativeness of Eddy-Covariance flux footprints for areas surrounding
539 AmeriFlux sites, *Agricultural and Forest Meteorology*, 301–302, 108350,
540 <https://doi.org/10.1016/j.agrformet.2021.108350>, 2021.
- 541 Cleverly, J., Vote, C., Isaac, P., Ewenz, C., Harahap, M., Beringer, J., Campbell, D. I., Daly, E.,
542 Eamus, D., He, L., Hunt, J., Grace, P., Hutley, L. B., Laubach, J., McCaskill, M., Rowlings, D.,
543 Rutledge Jonker, S., Schipper, L. A., Schroder, I., Teodosio, B., Yu, Q., Ward, P. R., Walker, J. P.,
544 Webb, J. A., and Grover, S. P. P.: Carbon, water and energy fluxes in agricultural systems of
545 Australia and New Zealand, 287, <https://doi.org/10.1016/j.agrformet.2020.107934>, 2020.
- 546 Cranko Page, J., De Kauwe, M. G., Abramowitz, G., Cleverly, J., Hinko-Najera, N., Hovenden, M. J.,
547 Liu, Y., Pitman, A. J., and Ogle, K.: Examining the role of environmental memory in the
548 predictability of carbon and water fluxes across Australian ecosystems, *Biogeosciences*, 19, 1913–
549 1932, 2022.
- 550 Cui, X., Goff, T., Cui, S., Menefee, D., Wu, Q., Rajan, N., Nair, S., Phillips, N., and Walker, F.:
551 Predicting carbon and water vapor fluxes using machine learning and novel feature ranking
552 algorithms, *Science of The Total Environment*, 775, 145130, 2021.

553 Don, A., Schumacher, J., and Freibauer, A.: Impact of tropical land-use change on soil organic carbon
554 stocks – a meta-analysis, 17, 1658–1670, <https://doi.org/10.1111/j.1365-2486.2010.02336.x>, 2011.

555 Field, A. P. and Gillett, R.: How to do a meta - analysis, *British Journal of Mathematical and*
556 *Statistical Psychology*, 63, 665–694, 2010.

557 Fu, D., Chen, B., Zhang, H., Wang, J., Black, T. A., Amiro, B. D., Bohrer, G., Bolstad, P., Coulter,
558 R., and Rahman, A. F.: Estimating landscape net ecosystem exchange at high spatial–temporal
559 resolution based on Landsat data, an improved upscaling model framework, and eddy covariance flux
560 measurements, *Remote Sensing of Environment*, 141, 90–104, 2014.

561 Fu, Z., Stoy, P. C., Poulter, B., Gerken, T., Zhang, Z., Wakbulcho, G., and Niu, S.: Maximum carbon
562 uptake rate dominates the interannual variability of global net ecosystem exchange, *Global Change*
563 *Biology*, 25, 3381–3394, 2019.

564 Hao, Y., Wang, Y., Mei, X., and Cui, X.: The response of ecosystem CO₂ exchange to small
565 precipitation pulses over a temperate steppe, *Plant Ecol*, 209, 335–347,
566 <https://doi.org/10.1007/s11258-010-9766-1>, 2010.

567 Harris, N. L., Gibbs, D. A., Baccini, A., Birdsey, R. A., de Bruin, S., Farina, M., Fatoyinbo, L.,
568 Hansen, M. C., Herold, M., Houghton, R. A., Potapov, P. V., Suarez, D. R., Roman-Cuesta, R. M.,
569 Saatchi, S. S., Slay, C. M., Turubanova, S. A., and Tyukavina, A.: Global maps of twenty-first
570 century forest carbon fluxes, *Nat. Clim. Chang.*, 11, 234–240, [https://doi.org/10.1038/s41558-020-](https://doi.org/10.1038/s41558-020-00976-6)
571 [00976-6](https://doi.org/10.1038/s41558-020-00976-6), 2021.

572 Huemmrich, K. F., Campbell, P., Landis, D., and Middleton, E.: Developing a common globally
573 applicable method for optical remote sensing of ecosystem light use efficiency, *Remote Sensing of*
574 *Environment*, 230, 111190, 2019.

575 Jongen, M., Pereira, J. S., Aires, L. M. I., and Pio, C. A.: The effects of drought and timing of
576 precipitation on the inter-annual variation in ecosystem-atmosphere exchange in a Mediterranean
577 grassland, *Agricultural and Forest Meteorology*, 151, 595–606,
578 <https://doi.org/10.1016/j.agrformet.2011.01.008>, 2011.

579 Jung, M., Reichstein, M., Margolis, H. A., Cescatti, A., Richardson, A. D., Arain, M. A., Arneth, A.,
580 Bernhofer, C., Bonal, D., and Chen, J.: Global patterns of land - atmosphere fluxes of carbon dioxide,
581 latent heat, and sensible heat derived from eddy covariance, satellite, and meteorological observations,
582 *Journal of Geophysical Research: Biogeosciences*, 116, 2011.

583 Jung, M., Reichstein, M., Schwalm, C. R., Huntingford, C., Sitch, S., Ahlström, A., Arneth, A.,
584 Camps-Valls, G., Ciais, P., Friedlingstein, P., Gans, F., Ichii, K., Jain, A. K., Kato, E., Papale, D.,
585 Poulter, B., Raduly, B., Rödenbeck, C., Tramontana, G., Viovy, N., Wang, Y.-P., Weber, U., Zaehle,
586 S., and Zeng, N.: Compensatory water effects link yearly global land CO₂ sink changes to
587 temperature, 541, 516–520, <https://doi.org/10.1038/nature20780>, 2017.

588 Jung, M., Schwalm, C., Migliavacca, M., Walther, S., Camps-Valls, G., Koirala, S., Anthoni, P.,
589 Besnard, S., Bodesheim, P., Carvalhais, N., Chevallier, F., Gans, F., S Goll, D., Haverd, V., Köhler,
590 P., Ichii, K., K Jain, A., Liu, J., Lombardozzi, D., E M S Nabel, J., A Nelson, J., O’Sullivan, M.,
591 Pallandt, M., Papale, D., Peters, W., Pongratz, J., Rödenbeck, C., Sitch, S., Tramontana, G., Walker,
592 A., Weber, U., and Reichstein, M.: Scaling carbon fluxes from eddy covariance sites to globe:
593 Synthesis and evaluation of the FLUXCOM approach, 17, 1343–1365, [https://doi.org/10.5194/bg-17-](https://doi.org/10.5194/bg-17-1343-2020)
594 [1343-2020](https://doi.org/10.5194/bg-17-1343-2020), 2020.

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

598 Kljun, N., Calanca, P., Rotach, M., and Schmid, H. P.: A simple two-dimensional parameterisation for
599 Flux Footprint Prediction (FFP), *Geoscientific Model Development*, 8, 3695–3713, 2015.

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

602 Luyssaert, S., Janssens, I. A., Sulkava, M., Papale, D., Dolman, A. J., Reichstein, M., Hollmén, J.,
603 Martin, J. G., Suni, T., Vesala, T., Loustau, D., Law, B. E., and Moors, E. J.: Photosynthesis drives
604 anomalies in net carbon-exchange of pine forests at different latitudes, 13, 2110–2127,
605 <https://doi.org/10.1111/j.1365-2486.2007.01432.x>, 2007.

606 Marcot, B. G. and Hanea, A. M.: What is an optimal value of k in k-fold cross-validation in discrete
607 Bayesian network analysis?, *Comput Stat*, 36, 2009–2031, [https://doi.org/10.1007/s00180-020-00999-](https://doi.org/10.1007/s00180-020-00999-9)
608 9, 2021.

609 Mitchell, S., Beven, K., and Freer, J.: Multiple sources of predictive uncertainty in modeled estimates
610 of net ecosystem CO₂ exchange, *Ecological Modelling*, 220, 3259–3270,
611 <https://doi.org/10.1016/j.ecolmodel.2009.08.021>, 2009.

612 Moffat, A. M., Papale, D., Reichstein, M., Hollinger, D. Y., Richardson, A. D., Barr, A. G.,
613 Beckstein, C., Braswell, B. H., Churkina, G., Desai, A. R., Falge, E., Gove, J. H., Heimann, M., Hui,
614 D., Jarvis, A. J., Kattge, J., Noormets, A., and Stauch, V. J.: Comprehensive comparison of gap-filling
615 techniques for eddy covariance net carbon fluxes, 147, 209–232,
616 <https://doi.org/10.1016/j.agrformet.2007.08.011>, 2007.

617 Moffat, A. M., Beckstein, C., Churkina, G., Mund, M., and Heimann, M.: Characterization of
618 ecosystem responses to climatic controls using artificial neural networks, 16, 2737–2749,
619 <https://doi.org/10.1111/j.1365-2486.2010.02171.x>, 2010.

620 Moher, D., Liberati, A., Tetzlaff, J., Altman, D. G., and Prisma Group: Preferred reporting items for
621 systematic reviews and meta-analyses: the PRISMA statement, *PLoS medicine*, 6, e1000097, 2009.

622 Moon, T. K.: The expectation-maximization algorithm, 13, 47–60, 1996.

623 Papale, D. and Valentini, R.: A new assessment of European forests carbon exchanges by eddy fluxes
624 and artificial neural network spatialization, 9, 525–535, [https://doi.org/10.1046/j.1365-](https://doi.org/10.1046/j.1365-2486.2003.00609.x)
625 2486.2003.00609.x, 2003.

626 Park, S.-B., Knohl, A., Lucas-Moffat, A. M., Migliavacca, M., Gerbig, C., Vesala, T., Peltola, O.,
627 Mammarella, I., Kolle, O., Lavrič, J. V., Prokushkin, A., and Heimann, M.: Strong radiative effect
628 induced by clouds and smoke on forest net ecosystem productivity in central Siberia, *Agricultural and*
629 *Forest Meteorology*, 250–251, 376–387, <https://doi.org/10.1016/j.agrformet.2017.09.009>, 2018.

630 Pearl, J.: Bayesian networks: A model of self-activated memory for evidential reasoning, in:
631 *Proceedings of the 7th Conference of the Cognitive Science Society*, University of California, Irvine,
632 CA, USA, 15–17, 1985.

633 Peltola, O., Vesala, T., Gao, Y., Rätty, O., Alekseychik, P., Aurela, M., Chojnicki, B., Desai, A. R.,
634 Dolman, A. J., Euskirchen, E. S., Friborg, T., Göckede, M., Helbig, M., Humphreys, E., Jackson, R.
635 B., Jocher, G., Joos, F., Klatt, J., Knox, S. H., Kowalska, N., Kutzbach, L., Lienert, S., Lohila, A.,
636 Mammarella, I., Nadeau, D. F., Nilsson, M. B., Oechel, W. C., Peichl, M., Pypker, T., Quanton, W.,

637 Rinne, J., Sachs, T., Samson, M., Schmid, H. P., Sonnentag, O., Wille, C., Zona, D., and Aalto, T.:
638 Monthly gridded data product of northern wetland methane emissions based on upscaling eddy
639 covariance observations, *Earth System Science Data*, 11, 1263–1289, [https://doi.org/10.5194/essd-11-](https://doi.org/10.5194/essd-11-1263-2019)
640 1263-2019, 2019.

641 Reed, D. E., Poe, J., Abraha, M., Dahlin, K. M., and Chen, J.: Modeled Surface-Atmosphere Fluxes
642 From Paired Sites in the Upper Great Lakes Region Using Neural Networks, *Journal of Geophysical*
643 *Research: Biogeosciences*, 126, <https://doi.org/10.1029/2021JG006363>, 2021.

644 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., and Prabhat:
645 Deep learning and process understanding for data-driven Earth system science, 566, 195–204,
646 <https://doi.org/10.1038/s41586-019-0912-1>, 2019.

647 Reitz, O., Graf, A., Schmidt, M., Ketzler, G., and Leuchner, M.: Upscaling Net Ecosystem Exchange
648 Over Heterogeneous Landscapes With Machine Learning, 126, e2020JG005814,
649 <https://doi.org/10.1029/2020JG005814>, 2021.

650 Ruppert, J., Mauder, M., Thomas, C., and Lüers, J.: Innovative gap-filling strategy for annual sums of
651 CO₂ net ecosystem exchange, 138, 5–18, <https://doi.org/10.1016/j.agrformet.2006.03.003>, 2006.

652 Shi, H., Luo, G., Zheng, H., Chen, C., Bai, J., Liu, T., Ochege, F. U., and De Maeyer, P.: Coupling the
653 water-energy-food-ecology nexus into a Bayesian network for water resources analysis and
654 management in the Syr Darya River basin, *Journal of Hydrology*, 581, 124387,
655 <https://doi.org/10.1016/j.jhydrol.2019.124387>, 2020.

656 Shi, H., Hellwich, O., Luo, G., Chen, C., He, H., Ochege, F. U., Van de Voorde, T., Kurban, A., and
657 de Maeyer, P.: A global meta-analysis of soil salinity prediction integrating satellite remote sensing,
658 soil sampling, and machine learning, 1–15, <https://doi.org/10.1109/TGRS.2021.3109819>, 2021.

659 Tian, X., Yan, M., van der Tol, C., Li, Z., Su, Z., Chen, E., Li, X., Li, L., Wang, X., Pan, X., Gao, L.,
660 and Han, Z.: Modeling forest above-ground biomass dynamics using multi-source data and
661 incorporated models: A case study over the qilian mountains, *Agricultural and Forest Meteorology*,
662 246, 1–14, <https://doi.org/10.1016/j.agrformet.2017.05.026>, 2017.

663 Tramontana, G., Jung, M., Schwalm, C. R., Ichii, K., Camps-Valls, G., Ráduly, B., Reichstein, M.,
664 Arain, M. A., Cescatti, A., Kiely, G., Merbold, L., Serrano-Ortiz, P., Sickert, S., Wolf, S., and Papale,
665 D.: Predicting carbon dioxide and energy fluxes across global FLUXNET sites with regression
666 algorithms, *Biogeosciences*, 13, 4291–4313, <https://doi.org/10.5194/bg-13-4291-2016>, 2016.

667 Van Hulse, J., Khoshgoftaar, T. M., and Napolitano, A.: Experimental perspectives on learning from
668 imbalanced data, in: *Proceedings of the 24th international conference on Machine learning*, New
669 York, NY, USA, 935–942, <https://doi.org/10.1145/1273496.1273614>, 2007.

670 Virkkala, A.-M., Aalto, J., Rogers, B. M., Tagesson, T., Treat, C. C., Natali, S. M., Watts, J. D.,
671 Potter, S., Lehtonen, A., Mauritz, M., Schuur, E. A. G., Kochendorfer, J., Zona, D., Oechel, W.,
672 Kobayashi, H., Humphreys, E., Goeckede, M., Iwata, H., Lafleur, P. M., Euskirchen, E. S., Bokhorst,
673 S., Marushchak, M., Martikainen, P. J., Elberling, B., Voigt, C., Biasi, C., Sonnentag, O., Parmentier,
674 F.-J. W., Ueyama, M., Celis, G., St.Louis, V. L., Emmerton, C. A., Peichl, M., Chi, J., Järveoja, J.,
675 Nilsson, M. B., Oberbauer, S. F., Torn, M. S., Park, S.-J., Dolman, H., Mammarella, I., Chae, N.,
676 Poyatos, R., López-Blanco, E., Christensen, T. R., Kwon, M. J., Sachs, T., Holl, D., and Luoto, M.:
677 Statistical upscaling of ecosystem CO₂ fluxes across the terrestrial tundra and boreal domain:
678 Regional patterns and uncertainties, *Global Change Biology*, 27, 4040–4059,
679 <https://doi.org/10.1111/gcb.15659>, 2021.

680 Walther, S., Besnard, S., Nelson, J. A., El-Madany, T. S., Migliavacca, M., Weber, U., Ermida, S. L.,
681 Brümmer, C., Schrader, F., Prokushkin, A. S., Panov, A. V., and Jung, M.: Technical note: A view
682 from space on global flux towers by MODIS and Landsat: The FluxnetEO dataset, *Biogeosciences*
683 *Discussions*, 1–40, <https://doi.org/10.5194/bg-2021-314>, 2021.

684 Wu, Z., Dijkstra, P., Koch, G. W., Peñuelas, J., and Hungate, B. A.: Responses of terrestrial
685 ecosystems to temperature and precipitation change: a meta-analysis of experimental manipulation,
686 *17*, 927–942, <https://doi.org/10.1111/j.1365-2486.2010.02302.x>, 2011.

687 Yan, J., Zhang, Y., Yu, G., Zhou, G., Zhang, L., Li, K., Tan, Z., and Sha, L.: Seasonal and inter-
688 annual variations in net ecosystem exchange of two old-growth forests in southern China, *Agricultural*
689 *and Forest Meteorology*, 182–183, 257–265, <https://doi.org/10.1016/j.agrformet.2013.03.002>, 2013.

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

693 Zhang, C., Brodylo, D., Sirianni, M. J., Li, T., Comas, X., Douglas, T. A., and Starr, G.: Mapping
694 CO₂ fluxes of cypress swamp and marshes in the Greater Everglades using eddy covariance
695 measurements and Landsat data, *Remote Sensing of Environment*, 262,
696 <https://doi.org/10.1016/j.rse.2021.112523>, 2021.

697 Zhou, Y., Li, X., Gao, Y., He, M., Wang, M., Wang, Y., Zhao, L., and Li, Y.: Carbon fluxes response
698 of an artificial sand-binding vegetation system to rainfall variation during the growing season in the
699 Tengger Desert, *Journal of Environmental Management*, 266,
700 <https://doi.org/10.1016/j.jenvman.2020.110556>, 2020.

701 Zhu, S., Clement, R., McCalmont, J., Davies, C. A., and Hill, T.: Stable gap-filling for longer eddy
702 covariance data gaps: A globally validated machine-learning approach for carbon dioxide, water, and
703 energy fluxes, *Agricultural and Forest Meteorology*, 314, 108777,
704 <https://doi.org/10.1016/j.agrformet.2021.108777>, 2022.

705