

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 time scales. For example, the MODIS-based
88 8-daily NDVI data may better capture temporal variation in the relationship between NEE and vegetation
89 growth than the Landsat-based 16-daily NDVI data. In contrast, the interpretation of NEE by variables
90 such as soil texture and soil organic content (SOC), which do not have temporal dynamic information, may
91 be limited to the interpretation of spatial variability, although they are considered to be important drivers of
92 NEE. Therefore, the importance of variables obtained from NEE simulations based on a data-driven
93 approach may differ from that in process-based models as well as in the actual driving mechanisms. This
94 may be related to the spatial and temporal resolution of the predictors used and the quality of the data. It is

95 necessary to consider the spatio-temporal resolution of the data for the actual biophysical variables used in
96 the different studies in the systematic evaluation of data-driven NEE simulations.

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

124
125 In this study, to evaluate the impacts of predictors use, algorithms, spatial/time scale, and validation methods on
126 model accuracy, we performed a meta-analysis of papers with prediction models that combine NEE
127 observations from flux towers, various predictors, and machine learning for the data-driven NEE simulations. In
128 addition, we also analyzed the causality of multiple features in NEE simulations and the joint effects of multiple
129 features on model accuracy using the Bayesian Network (BN) (a multivariate statistical analysis approach
130 (Pearl, 1985)). The findings of this study can provide some general guidance for future NEE simulations.

131 **2 Methodology**

132 **2.1 Criteria for including articles**

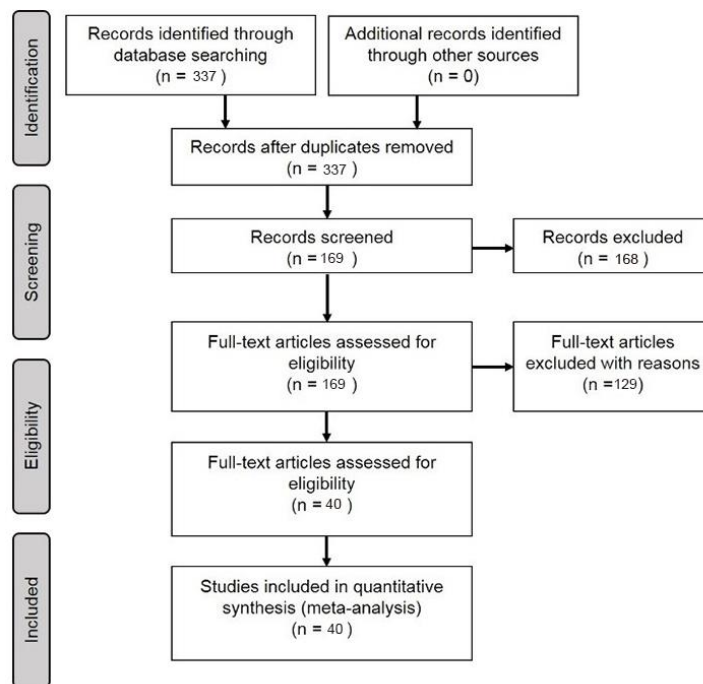
133 In the Scopus database, a literature query was applied to titles, abstracts, and keywords (Table 1) according to
134 Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) (Moher et al., 2009) (Fig. 1):

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

145
146 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		

147



148

149 Figure 1. PRISMA-based paper filtering flowchart.

150 **2.2 Features of prediction models**

151 Typically, the flow of the NEE prediction modeling framework (Fig. 2) based on flux observations and machine
 152 learning is as follows: first, half-hourly scale NEE flux observations are aggregated into various time scale NEE
 153 data, and gap-filling techniques (Moffat et al., 2007) are often used in this step to obtain complete NEE series
 154 when data are missing. Various predictors including meteorological variables, remote sensing-based biophysical
 155 variables, etc. are extracted to match site-scale NEE series to generate a training dataset containing the target
 156 variable NEE and various covariates. Subsequently, various algorithms are used for the NEE prediction model
 157 construction and validated in different ways (e.g., leave-one-site-out validation (Zeng et al., 2020)). Finally, in
 158 some studies, prediction models were applied to gridded covariate data to map the regional or global-scale NEE
 159 spatial and temporal variations (Zeng et al., 2020; Papale and Valentini, 2003; Jung et al., 2020). The
 160 information of R-squared (at the validation phase) and the associated model features reported in the article are
 161 considered as one data record for the formal meta-analysis (i.e., each R-squared record corresponding to a
 162 prediction model). From the included papers, R-squared records and various features (Table 2) involved in the
 163 NEE modeling framework (Fig. 2) were extracted (including the used algorithms, modeling/validation methods,
 164 remote sensing data, meteorological data, biophysical data, and ancillary data). In some studies, multiple
 165 algorithms were applied to the same dataset, or models with different features were developed (Virkkala et al.,
 166 2021; Zhang et al., 2021; Cleverly et al., 2020; Tramontana et al., 2016). In these cases, multiple data records
 167 will be documented.

168

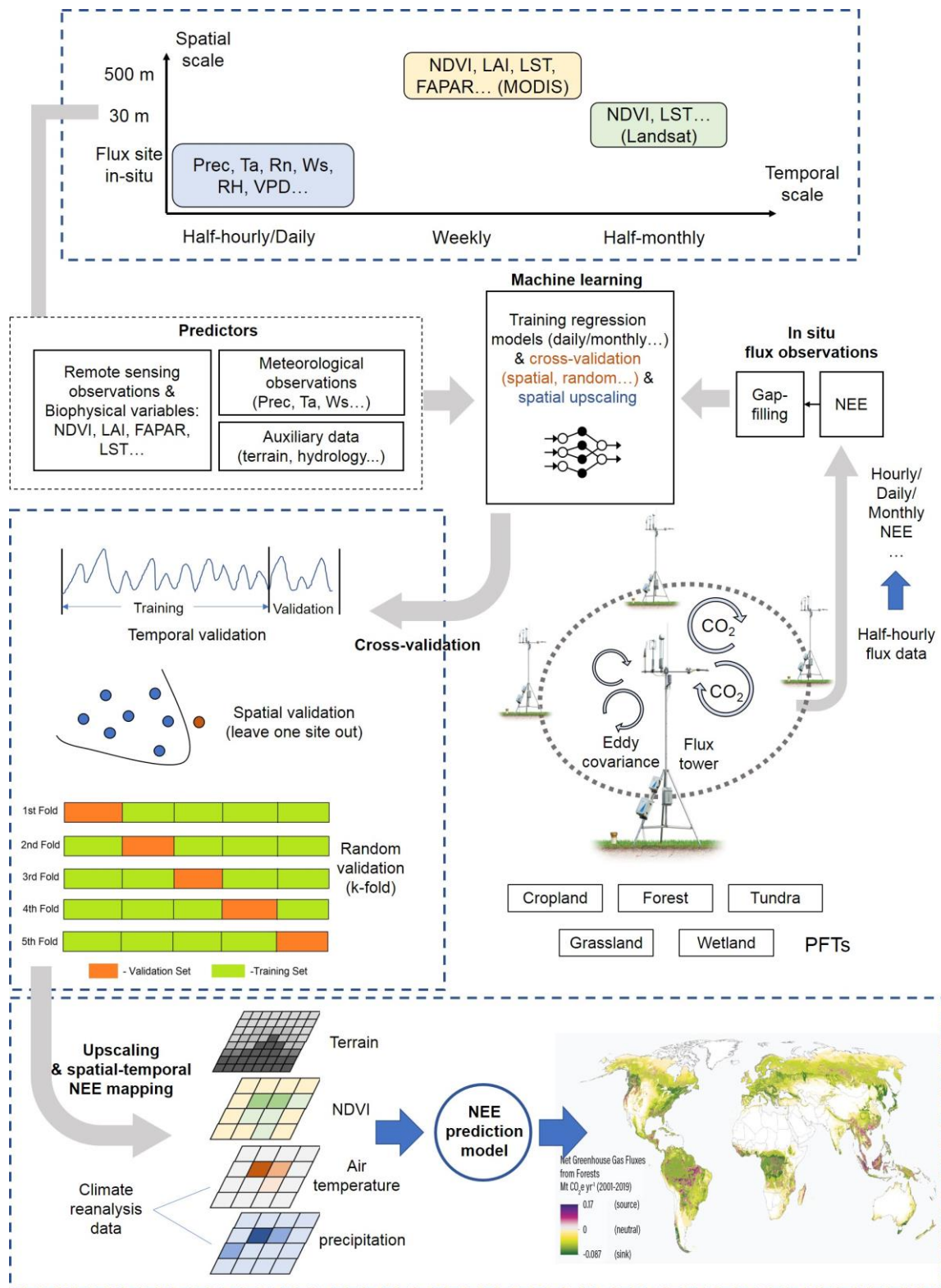
169 In the practical information extracting step, we categorized such features in a comparable manner. First, we
 170 categorized the various algorithms used in these papers, although the same algorithm may also have a variant

171 form or an optimized parameter scheme. They are categorized into the following families of algorithms:
172 Random Forests (RF), Multiple Linear Regressions (MLR), Artificial Neural Networks (ANN), Support Vector
173 Machines (SVM), Partial Least Squares Regression (PLSR), Generalized additive model (GAM), Boosted
174 Regression Tree (BRT), Bayesian Additive Regression Trees (BART), Cubist, model tree ensembles (MTE).
175 Second, we classified the spatial scales of these studies. Models with study areas (spatial extent covered by flux
176 stations) smaller than 100x100 km were classified as ‘local’ scale models, those with study area sizes exceeding
177 continental scale were classified as ‘global’ scale, and those with study area sizes in between were classified as
178 ‘regional’ scale. Third, for various predictors, we only recorded whether the predictors were used or not without
179 distinguishing the detailed data sources and categories (e.g., grid meteorological data from various reanalysis
180 datasets and in-situ meteorological observations from flux stations), measurement methods (e.g., soil moisture
181 measured/estimated by remote sensing or in situ sensors), etc. Fourth, we documented PFTs for the prediction
182 models from the description of study areas or sites in these papers. They are classified into the following types:
183 forest, grassland, cropland, wetland, savannah, tundra, and multi-PFTs (models containing a mixture of multiple
184 PFTs). Models not belonging to the above PFTs were not given a PFT field and were not included in the
185 subsequent analysis of the PFT differences. Other features (Table 2) are extracted directly from the
186 corresponding descriptions in the papers in an explicit manner.

187

188 Subsequently, the model accuracies corresponding to different levels of various features are compared in a
189 cross-study fashion. In the evaluation of algorithms and time scales, we also implement comparisons within
190 individual studies. For example, in the evaluation of the effects of the algorithms, we compare the accuracy of
191 models using the same training data and keeping other features as constants in individual studies. In this intra-
192 study comparison step, only algorithms with relatively large sample sizes in the cross-study comparisons were
193 selected. In this study, algorithms with less than 10 available model records are not considered to have a
194 sufficient sample size and we do not give further conclusive opinions on the accuracy of these algorithms due to
195 their small samples (e.g., PLSR and BART with high R-squared but very few records as evidence). MLR, RF,
196 SVM, and ANN were found to have large sample sizes (Fig. 5a), and thus their accuracies can be comparable.
197 Based on this, in the intra-study comparison step, we only compare the accuracy differences between MLR, RF,
198 SVM, and ANN in the context of using the same data and the same other model features (Fig. 5b).

199



200

201

202

203

204

205

206

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.

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

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)
Time scale	The time 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)

209

210 **2.3 Bayesian Network for analyzing joint effects**

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

$$213 \quad P(X) = P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i | pa(X_i)) \quad (1)$$

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

216

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

219
$$MI = H(Q) - H(Q|F) = \sum_q \sum_f P(q, f) \log_2 \left(\frac{P(q, f)}{P(q)P(f)} \right) \quad (2)$$

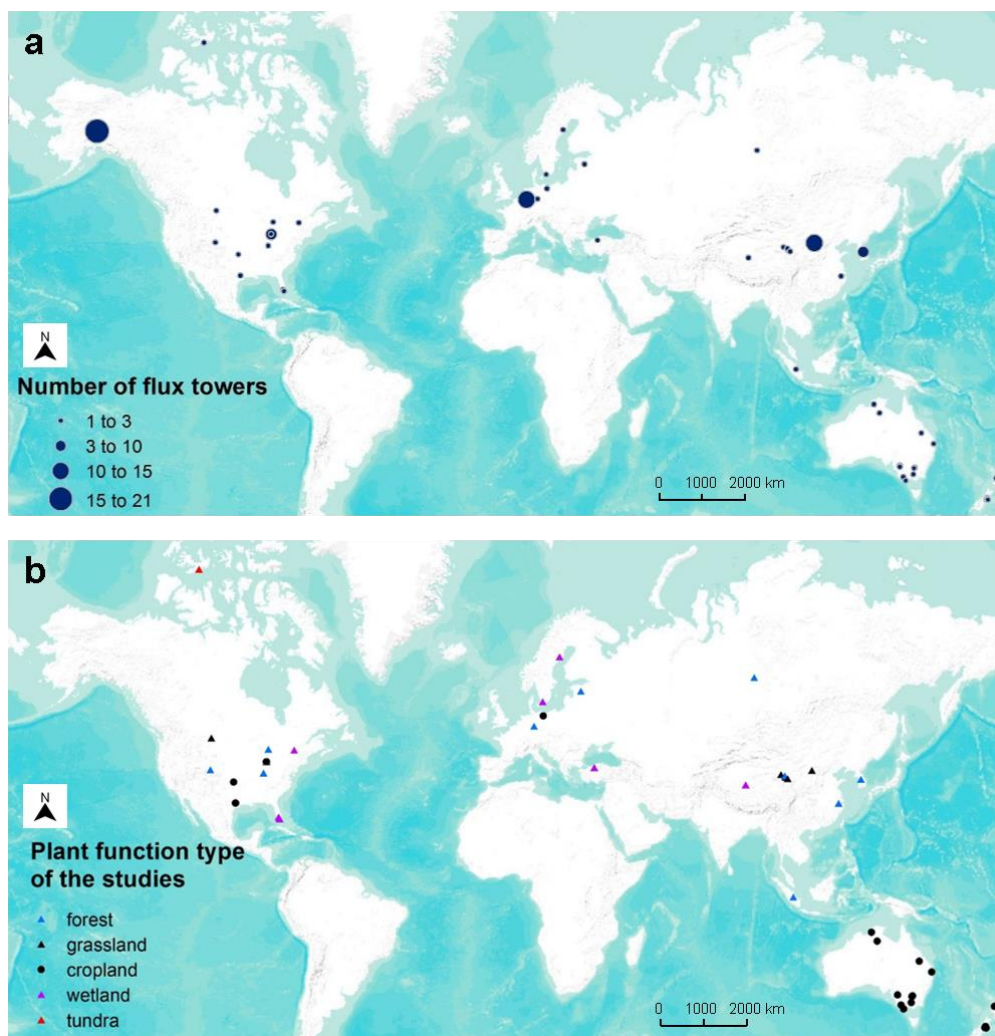
220 where H represents the entropy, Q represents the target node, F represents the set of other nodes and q and f
 221 represent the status of Q and F.

222 **3 Results**

223 **3.1 Articles included in the meta-analysis**

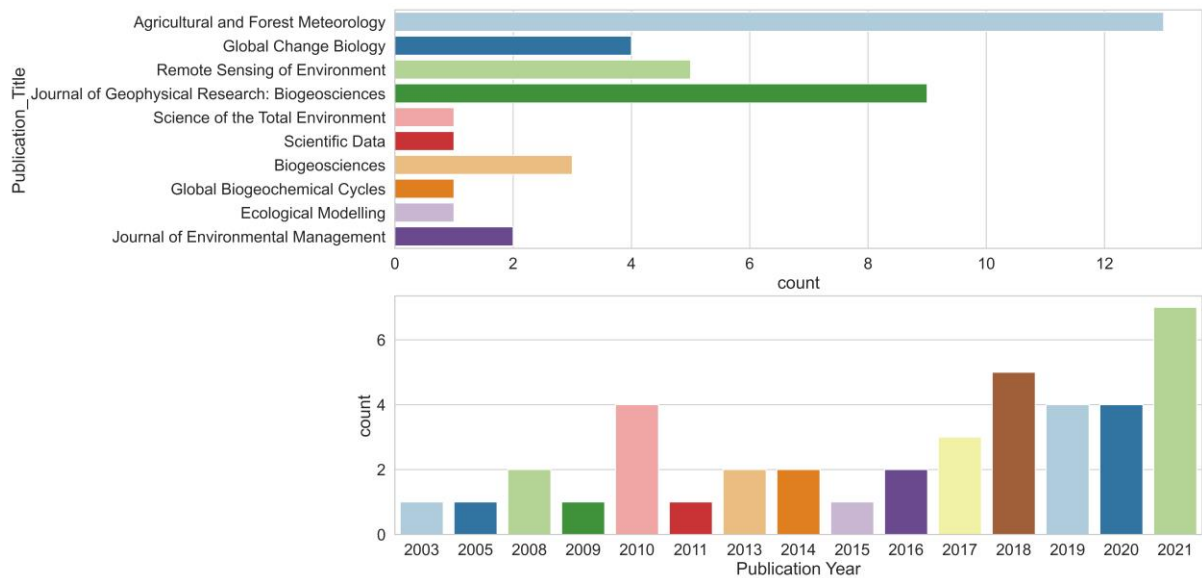
224 We included 40 articles (Table S2) and extracted 178 model records for the formal meta-analysis (Fig. 1). Most
 225 studies were implemented in Europe, North America, Oceania, and China (Fig. 3). The number of such papers is
 226 increasing recently (Fig. 4) and it shows the machine learning approach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

230



231
 232 Figure 3. Location of studies (a) included with the number of flux sites included and (b) their PFTs in the meta-
 233 analysis (total of 40 studies and 178 model records). Global (mainly based on FluxNet (Tramontana et al.,

234 2016)) and continental-scale studies are not shown on the map due to the difficulty of identifying specific
 235 locations.
 236



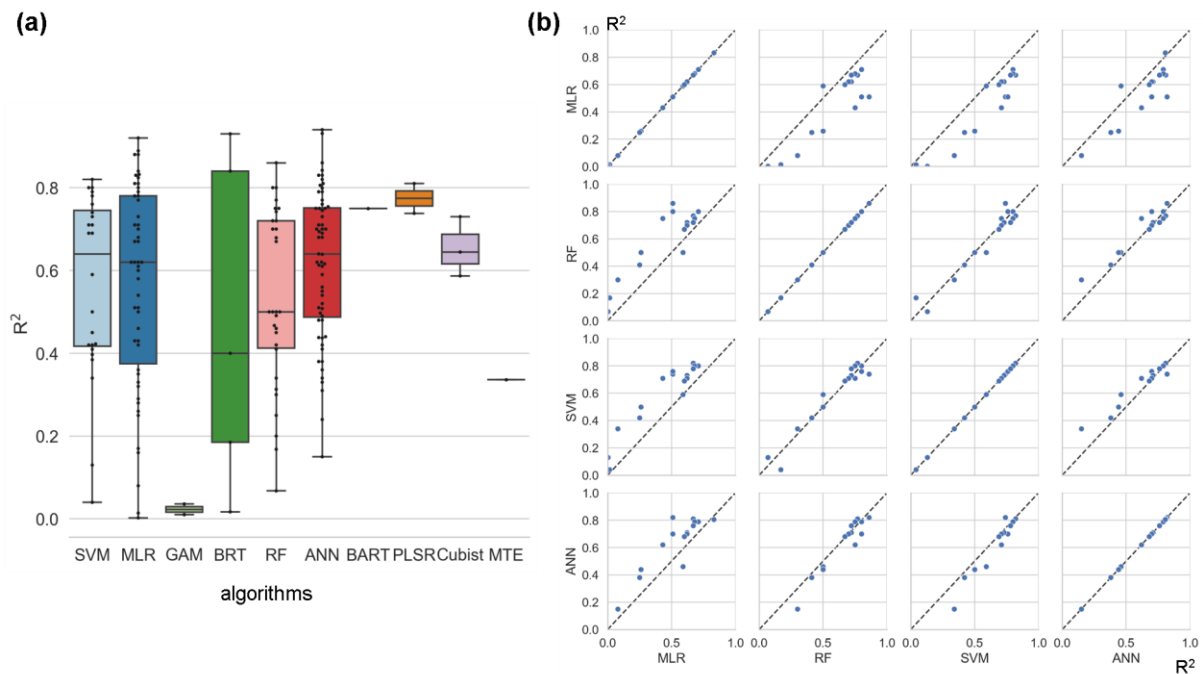
237
 238 Figure 4. The number of studies published across journals and the total number of publications per year.

239 **3.2 The formal Meta-analysis**

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

242 **3.2.1 Algorithms**

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



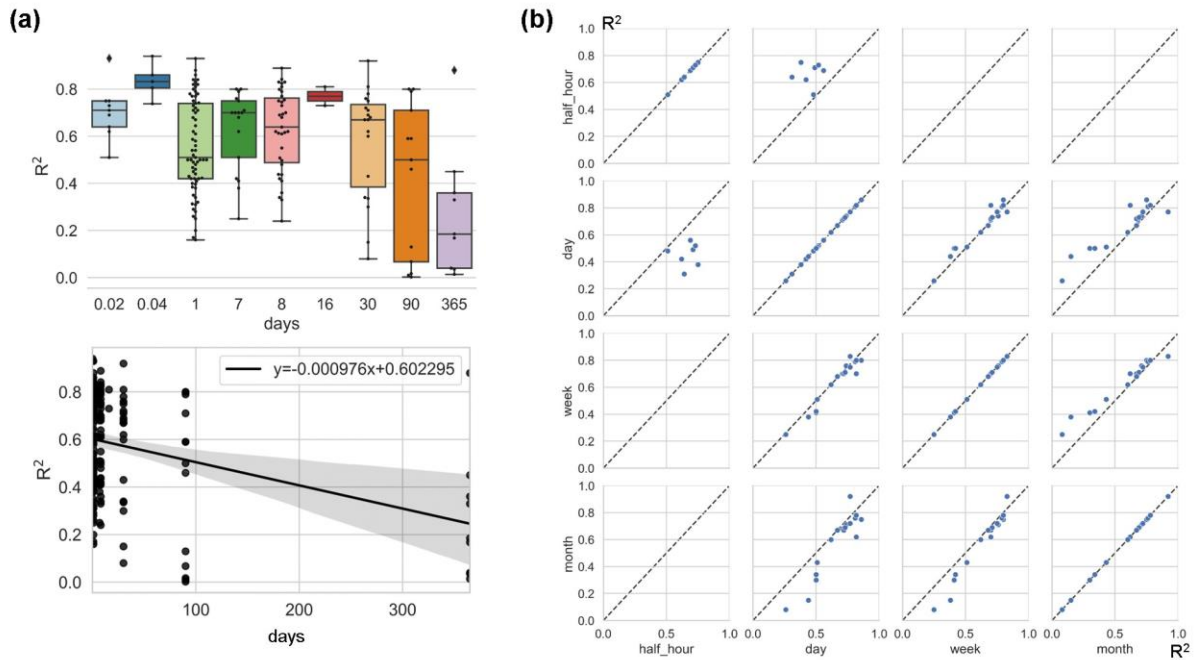
251

252 Figure 5. Differences in model accuracy (R-squared) using different algorithms across studies (a) and internal
 253 comparisons of the model accuracy (R-squared) of selected pairs of algorithms within individual studies (b).

254 Regression algorithms: Random Forests (RF), Multiple Linear Regressions (MLR), Artificial Neural Networks
 255 (ANN), Support Vector Machines (SVM), Partial Least Squares Regression (PLSR), Generalized additive
 256 model (GAM), Boosted Regression Tree (BRT), Bayesian Additive Regression Trees (BART), Cubist, model
 257 tree ensembles (MTE). In panel (a), the horizontal line in the box indicates the medians. The top and bottom
 258 border lines of the box indicate the 75% and 25% percentiles, respectively.

259 3.2.2 Time scales

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



266
 267
 268
 269
 270
 271
 272

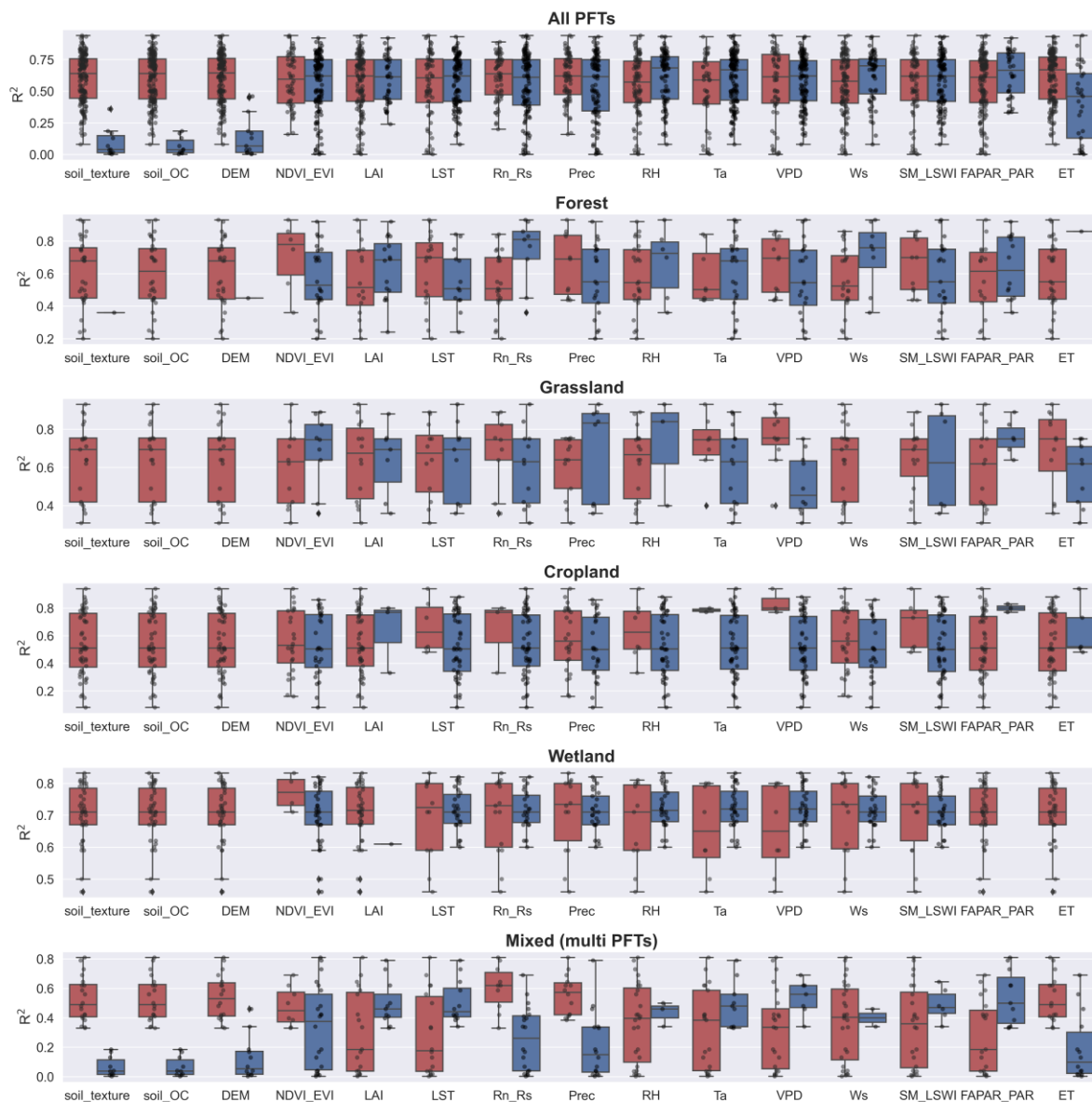
Figure 6. Differences in model accuracy (R-squared) at different time scales across studies with the linear regression between R-squared and time scales (a), and comparison of the model accuracy (R-squared) of selected pairs of time scales within individual studies (b). All model records were included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

273 3.2.3 Various predictors

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

283

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



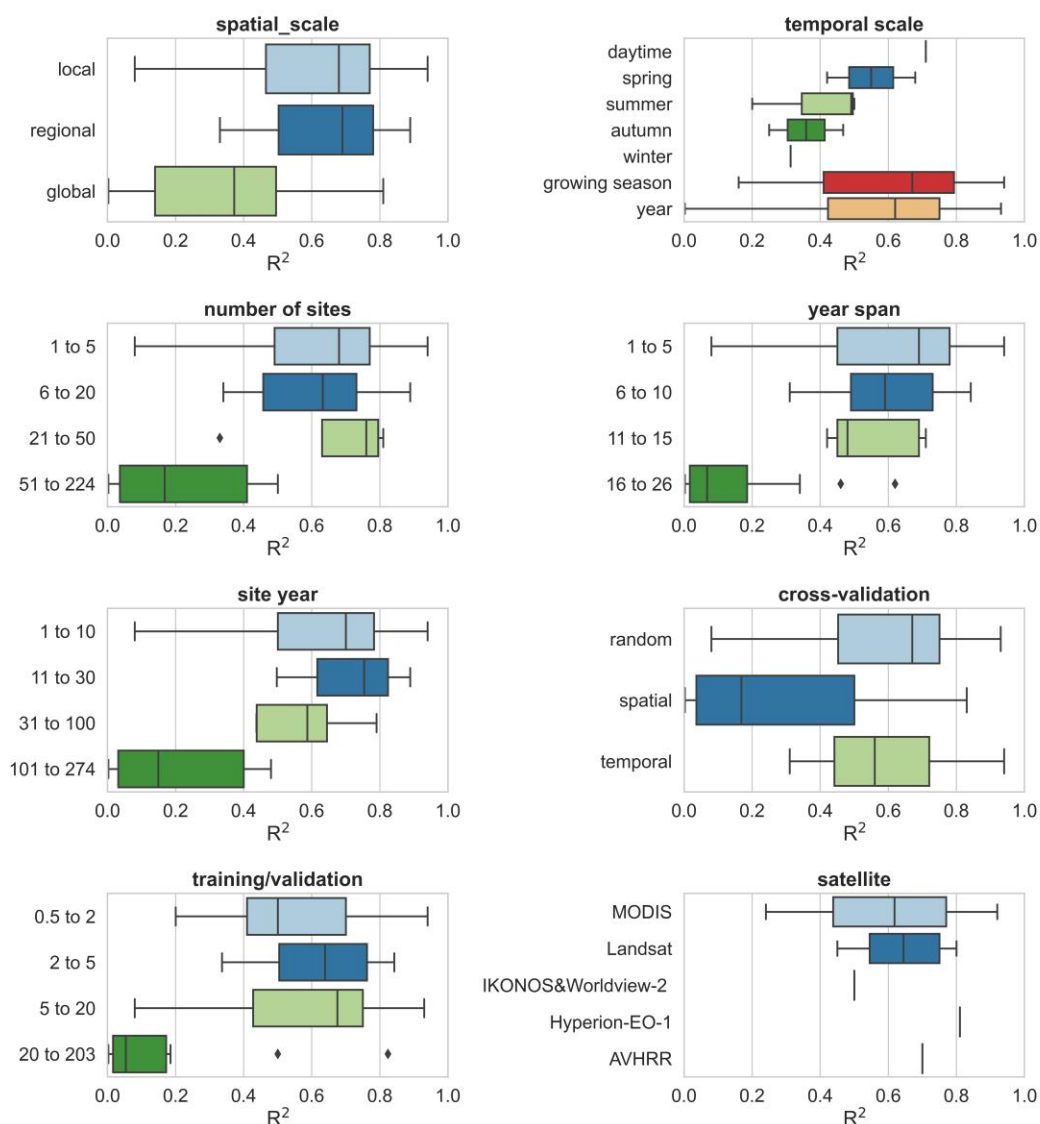
292

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

300 **3.2.4 Other features**

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

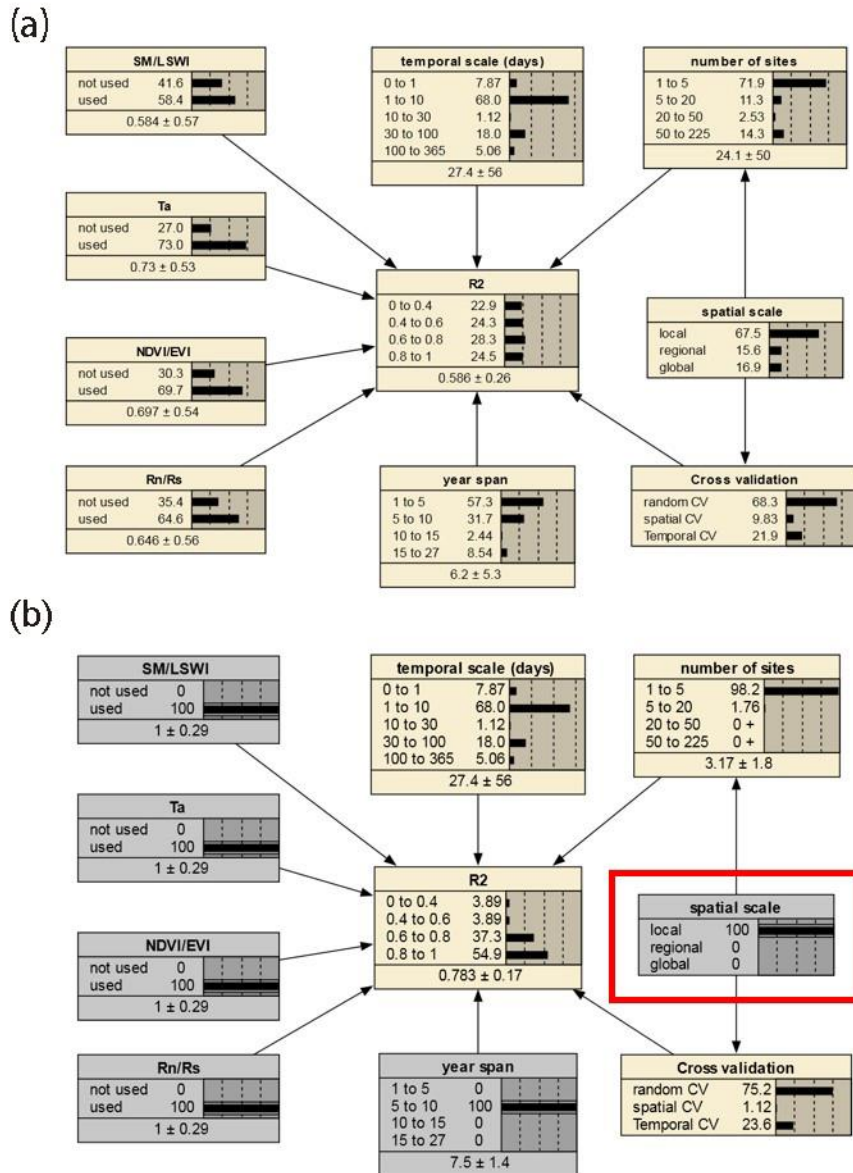
305 covariates and that studies at small scales are more likely to have higher model accuracy. Spatial validation
 306 (usually 'leave one site out') corresponds to lower model accuracy compared to random and temporal validation.
 307 This again confirms the dominant role of heterogeneity in the relationship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



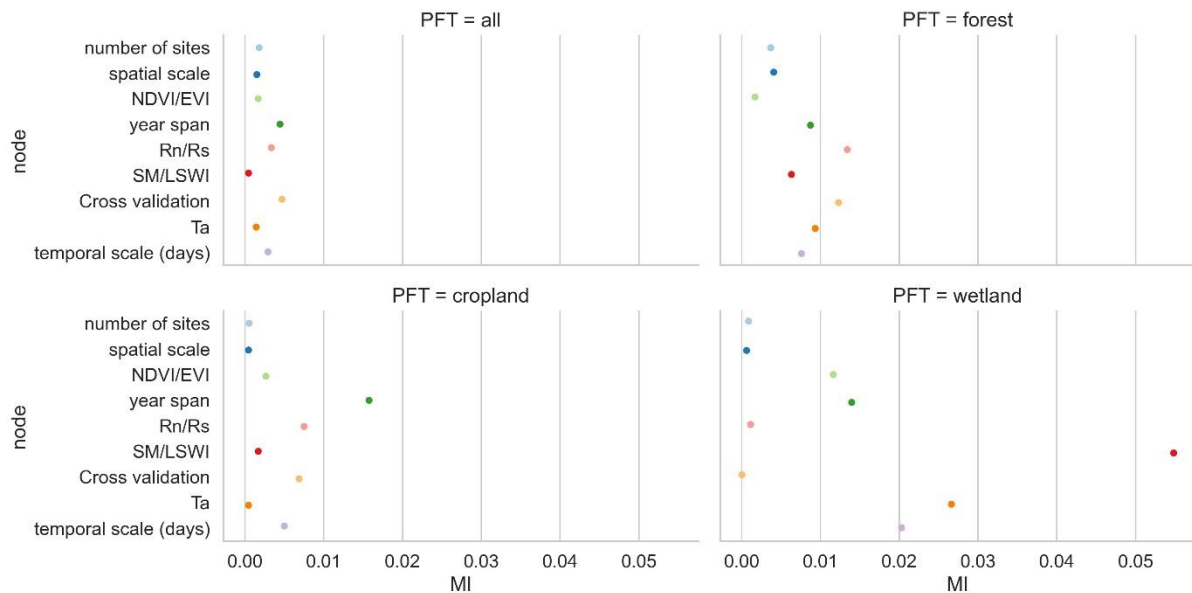
318
 319 Figure 8. The impacts of other features (i.e. spatial scale, study period, number of sites, year span, site year,
 320 cross-validation method, training/validation, and satellite imagery) on the model performance.

321 **3.3. The joint causal impacts of multi-features based on the BN**

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



341
 342 Figure 9. The joint effects of multiple features on the R-squared based on the BN with all records input (a) and
 343 the inference on the probability distribution of R-squared based on the BN with the status of some nodes
 344 determined (b). The values before and after the “±” indicate the mean and standard deviation of the distribution,
 345 respectively. The gray boxes indicate that the status of the nodes has been determined. In panel (b), specific
 346 values of parent nodes such as ‘spatial scale’ are determined (shown in the red box), leading to an increase in the
 347 expected R-squared compared to the average scenario of panel (a) (as inferred from the posterior conditional
 348 probabilities with the status of the node ‘spatial scale’ are determined as ‘local’).
 349



350

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

354 4 Discussions

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

363

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

373 **4.1 Challenges in the site-scale NEE simulation and implications for other carbon flux simulations**

374 **4.1.1 Variations in time scales**

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

398

399 In terms of completeness and purity of training data, hourly and daily models can be better compared to monthly
400 and yearly models. Hourly and daily models can usually preclude those low-quality data and gaps in the flux
401 observations. However, for monthly and yearly scale models, gap-filling (Ruppert et al., 2006; Moffat et al.,
402 2007; Zhu et al., 2022) is necessary because there are few complete and continuous fluxes observations without
403 data gaps on the monthly to yearly scales. Since various gap-filling techniques rely on environmental factors
404 (Moffat et al., 2007) such as meteorological observations, this may introduce uncertainty in the predictive
405 models (i.e., a small fraction of the observed information of NEE is estimated from a combination of
406 independent variables). How it would affect the accuracy of prediction models at various time scales remains
407 uncertain, although various gap-filling techniques have been widely used in the pre-processing of training data.

408

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

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

416 **4.1.2 Scale mismatch of explanatory predictors and flux footprints**

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

420

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

433

434 Since few of the studies included in this meta-analysis considered the effect of variation in flux footprint, this
435 feature was difficult to consider in this study. However, its influence should still be further investigated in future
436 studies. With flux footprints calculated (Kljun et al., 2015) and the factors around the flux site (Walther et al.,
437 2021) that affect the flux footprint incorporated, it is promising to clarify this issue.

438 **4.1.3 Possible unbalance of training and validation sets**

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

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

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

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

454

455 In summary, the site-scale NEE predictions may require more focus on the internal heterogeneity of the NEE
456 dataset and the matching of the training set and validation set, and also require a better understanding of the
457 influence of different scales of the same variable (e.g. site-scale precipitation and grid-scale precipitation in the
458 reanalysis meteorological data) across modeling and upscaling steps. For the prediction of other carbon fluxes
459 such as methane fluxes (in the same framework as the NEE predictions), the results of this study may also be
460 partially applicable, although there may be significant differences in the use of specific predictors (Peltola et al.,
461 2019).

462 **4.2 Uncertainties**

463 The uncertainties in this analysis may include:

- 464 a) **Publication bias and weighting:** Publication bias is not refined due to the limitations of the number of
465 articles that can be included. Meta-analyses often measure the quality of journals and the data availability
466 (Borenstein et al., 2011; Field and Gillett, 2010) to determine the weighting of the literature in a
467 comprehensive assessment. However, a high proportion of the articles in this study did not make flux
468 observations publicly available or share the NEE prediction models developed. Furthermore, meta-analysis
469 studies in other fields typically measure the impact of papers by evidence/data volume, and the variance of
470 the evaluated effects (Adams et al., 1997; Don et al., 2011; Liu et al., 2018). However, in this study,
471 because no convincing method is found to quantify the weights of results from included articles, some
472 features (e.g. the number of flux sites, the span of years) were directly assessed rather than used to
473 determine the weights of the articles.
- 474 b) **Limitations of the criteria for inclusion in the literature:** in the model accuracy-based evaluation, we
475 selected only literature that developed multiple regression models. Potentially valuable information from
476 univariate regression models was not included. In addition, only papers in high-quality English journals
477 were included in this study to control for possible errors due to publication bias. However, many studies
478 that fit this theme may have been published in other languages or other journals.
- 479 c) **Independence between features:** There is dependence between the evaluated features (e.g. the dependency
480 between the spatial extent and the number of sites). It may negatively affect the assessment of the impact
481 of individual features on the accuracy of the model, although the BN-based analysis of joint effects can
482 reduce the impact of this dependence between variables by specifying causal relationships between
483 features. The interference of unknown dependencies between features may still not be eliminated when we
484 focus on the effects of an individual feature on the model performance. We should pay more attention to
485 the effect of features on model accuracy individually in future studies, and it may be valuable to keep other
486 features as constants while changing the level of only one feature and assessing the difference. It may help
487 us to understand the real sensitivity of model accuracy to different features in specific conditions. The
488 sample size collected in this study (178 records in total) is not very large. This also suggests that more
489 future efforts should be devoted to the comprehensive evaluation and summarization of NEE simulations.

490
491
492
493
494

495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511

Additionally, there are still other potential factors not considered by this study such as the uncertainty of climate data (site vs reanalysis), footprint matching between site and satellite images, etc. Overall, although the quantitative results of this study should be used with caution, they still have positive implications for guiding future such studies.

5 Conclusion

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

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

512 **Acknowledgments**

513 We thank the editors and three anonymous referees for their insightful comments on this paper which
514 substantially improved.

515 **Financial support**

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

520 **Author contributions**

521 H.S and G.L initiated this research and were responsible for the integrity of the work as a whole. H.S performed
522 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
523 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
524 resources and financial support.

525 **Competing interests**

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

527 **Data availability**

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

530 **Code availability**

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

533

534

535 **References**

- 536 Abbasian, H., Solgi, E., Mohsen Hosseini, S., and Hossein Kia, S.: Modeling terrestrial net ecosystem
537 exchange using machine learning techniques based on flux tower measurements, *Ecological*
538 *Modelling*, 466, 109901, <https://doi.org/10.1016/j.ecolmodel.2022.109901>, 2022.
- 539 Adams, D. C., Gurevitch, J., and Rosenberg, M. S.: Resampling tests for meta - analysis of ecological
540 data, *Ecology*, 78, 1277–1283, 1997.
- 541 Baldocchi, D. D.: Assessing the eddy covariance technique for evaluating carbon dioxide exchange
542 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)
543 [2486.2003.00629.x](https://doi.org/10.1046/j.1365-2486.2003.00629.x), 2003.
- 544 Berryman, E. M., Vanderhoof, M. K., Bradford, J. B., Hawbaker, T. J., Henne, P. D., Burns, S. P.,
545 Frank, J. M., Birdsey, R. A., and Ryan, M. G.: Estimating soil respiration in a subalpine landscape
546 using point, terrain, climate, and greenness data, *Journal of Geophysical Research: Biogeosciences*,
547 123, 3231–3249, 2018.
- 548 Borenstein, M., Hedges, L. V., Higgins, J. P., and Rothstein, H. R.: *Introduction to meta-analysis*,
549 John Wiley & Sons, 2011.
- 550 Cho, S., Kang, M., Ichii, K., Kim, J., Lim, J.-H., Chun, J.-H., Park, C.-W., Kim, H. S., Choi, S.-W.,
551 and Lee, S.-H.: Evaluation of forest carbon uptake in South Korea using the national flux tower
552 network, remote sensing, and data-driven technology, *Agricultural and Forest Meteorology*, 311,
553 108653, 2021.
- 554 Chu, H., Luo, X., Ouyang, Z., Chan, W. S., Dengel, S., Biraud, S. C., Torn, M. S., Metzger, S.,
555 Kumar, J., Arain, M. A., Arkebauer, T. J., Baldocchi, D., Bernacchi, C., Billesbach, D., Black, T. A.,
556 Blanken, P. D., Bohrer, G., Bracho, R., Brown, S., Brunzell, N. A., Chen, J., Chen, X., Clark, K.,
557 Desai, A. R., Duman, T., Durden, D., Fares, S., Forbrich, I., Gamon, J. A., Gough, C. M., Griffis, T.,
558 Helbig, M., Hollinger, D., Humphreys, E., Ikawa, H., Iwata, H., Ju, Y., Knowles, J. F., Knox, S. H.,
559 Kobayashi, H., Kolb, T., Law, B., Lee, X., Litvak, M., Liu, H., Munger, J. W., Noormets, A., Novick,
560 K., Oberbauer, S. F., Oechel, W., Oikawa, P., Papuga, S. A., Pendall, E., Prajapati, P., Prueger, J.,
561 Quinton, W. L., Richardson, A. D., Russell, E. S., Scott, R. L., Starr, G., Staebler, R., Stoy, P. C.,
562 Stuart-Haëntjens, E., Sonnentag, O., Sullivan, R. C., Suyker, A., Ueyama, M., Vargas, R., Wood, J.
563 D., and Zona, D.: Representativeness of Eddy-Covariance flux footprints for areas surrounding
564 AmeriFlux sites, *Agricultural and Forest Meteorology*, 301–302, 108350,
565 <https://doi.org/10.1016/j.agrformet.2021.108350>, 2021.
- 566 Cleverly, J., Vote, C., Isaac, P., Ewenz, C., Harahap, M., Beringer, J., Campbell, D. I., Daly, E.,
567 Eamus, D., He, L., Hunt, J., Grace, P., Hutley, L. B., Laubach, J., McCaskill, M., Rowlings, D.,
568 Rutledge Jonker, S., Schipper, L. A., Schroder, I., Teodosio, B., Yu, Q., Ward, P. R., Walker, J. P.,
569 Webb, J. A., and Grover, S. P. P.: Carbon, water and energy fluxes in agricultural systems of
570 Australia and New Zealand, 287, <https://doi.org/10.1016/j.agrformet.2020.107934>, 2020.
- 571 Cranko Page, J., De Kauwe, M. G., Abramowitz, G., Cleverly, J., Hinko-Najera, N., Hovenden, M. J.,
572 Liu, Y., Pitman, A. J., and Ogle, K.: Examining the role of environmental memory in the
573 predictability of carbon and water fluxes across Australian ecosystems, *Biogeosciences*, 19, 1913–
574 1932, 2022.
- 575 Cui, X., Goff, T., Cui, S., Menefee, D., Wu, Q., Rajan, N., Nair, S., Phillips, N., and Walker, F.:
576 Predicting carbon and water vapor fluxes using machine learning and novel feature ranking
577 algorithms, *Science of The Total Environment*, 775, 145130, 2021.

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

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

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

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

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

592 Harris, N. L., Gibbs, D. A., Baccini, A., Birdsey, R. A., de Bruin, S., Farina, M., Fatoyinbo, L.,
593 Hansen, M. C., Herold, M., Houghton, R. A., Potapov, P. V., Suarez, D. R., Roman-Cuesta, R. M.,
594 Saatchi, S. S., Slay, C. M., Turubanova, S. A., and Tyukavina, A.: Global maps of twenty-first
595 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)
596 [00976-6](https://doi.org/10.1038/s41558-020-00976-6), 2021.

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

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

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

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

613 Jung, M., Schwalm, C., Migliavacca, M., Walther, S., Camps-Valls, G., Koirala, S., Anthoni, P.,
614 Besnard, S., Bodesheim, P., Carvalhais, N., Chevallier, F., Gans, F., S Goll, D., Haverd, V., Köhler,
615 P., Ichii, K., K Jain, A., Liu, J., Lombardozzi, D., E M S Nabel, J., A Nelson, J., O’Sullivan, M.,
616 Pallandt, M., Papale, D., Peters, W., Pongratz, J., Rödenbeck, C., Sitch, S., Tramontana, G., Walker,
617 A., Weber, U., and Reichstein, M.: Scaling carbon fluxes from eddy covariance sites to globe:
618 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)
619 [1343-2020](https://doi.org/10.5194/bg-17-1343-2020), 2020.

- 620 Kaur, H., Pannu, H. S., and Malhi, A. K.: A Systematic Review on Imbalanced Data Challenges in
621 Machine Learning: Applications and Solutions, *ACM Comput. Surv.*, 52, 79:1-79:36,
622 <https://doi.org/10.1145/3343440>, 2019.
- 623 Kljun, N., Calanca, P., Rotach, M., and Schmid, H. P.: A simple two-dimensional parameterisation for
624 Flux Footprint Prediction (FFP), *Geoscientific Model Development*, 8, 3695–3713, 2015.
- 625 Liu, Q., Zhang, Y., Liu, B., Amonette, J. E., Lin, Z., Liu, G., Ambus, P., and Xie, Z.: How does
626 biochar influence soil N cycle? A meta-analysis, *Plant and soil*, 426, 211–225, 2018.
- 627 Luyssaert, S., Janssens, I. A., Sulkava, M., Papale, D., Dolman, A. J., Reichstein, M., Hollmén, J.,
628 Martin, J. G., Suni, T., Vesala, T., Loustau, D., Law, B. E., and Moors, E. J.: Photosynthesis drives
629 anomalies in net carbon-exchange of pine forests at different latitudes, 13, 2110–2127,
630 <https://doi.org/10.1111/j.1365-2486.2007.01432.x>, 2007.
- 631 Marcot, B. G. and Hanea, A. M.: What is an optimal value of k in k-fold cross-validation in discrete
632 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)
633 9, 2021.
- 634 Mitchell, S., Beven, K., and Freer, J.: Multiple sources of predictive uncertainty in modeled estimates
635 of net ecosystem CO₂ exchange, *Ecological Modelling*, 220, 3259–3270,
636 <https://doi.org/10.1016/j.ecolmodel.2009.08.021>, 2009.
- 637 Moffat, A. M., Papale, D., Reichstein, M., Hollinger, D. Y., Richardson, A. D., Barr, A. G.,
638 Beckstein, C., Braswell, B. H., Churkina, G., Desai, A. R., Falge, E., Gove, J. H., Heimann, M., Hui,
639 D., Jarvis, A. J., Kattge, J., Noormets, A., and Stauch, V. J.: Comprehensive comparison of gap-filling
640 techniques for eddy covariance net carbon fluxes, 147, 209–232,
641 <https://doi.org/10.1016/j.agrformet.2007.08.011>, 2007.
- 642 Moffat, A. M., Beckstein, C., Churkina, G., Mund, M., and Heimann, M.: Characterization of
643 ecosystem responses to climatic controls using artificial neural networks, 16, 2737–2749,
644 <https://doi.org/10.1111/j.1365-2486.2010.02171.x>, 2010.
- 645 Moher, D., Liberati, A., Tetzlaff, J., Altman, D. G., and Prisma Group: Preferred reporting items for
646 systematic reviews and meta-analyses: the PRISMA statement, *PLoS medicine*, 6, e1000097, 2009.
- 647 Moon, T. K.: The expectation-maximization algorithm, 13, 47–60, 1996.
- 648 Papale, D. and Valentini, R.: A new assessment of European forests carbon exchanges by eddy fluxes
649 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)
650 2486.2003.00609.x, 2003.
- 651 Park, S.-B., Knohl, A., Lucas-Moffat, A. M., Migliavacca, M., Gerbig, C., Vesala, T., Peltola, O.,
652 Mammarella, I., Kolle, O., Lavrič, J. V., Prokushkin, A., and Heimann, M.: Strong radiative effect
653 induced by clouds and smoke on forest net ecosystem productivity in central Siberia, *Agricultural and*
654 *Forest Meteorology*, 250–251, 376–387, <https://doi.org/10.1016/j.agrformet.2017.09.009>, 2018.
- 655 Pearl, J.: Bayesian networks: A model of self-activated memory for evidential reasoning, in:
656 *Proceedings of the 7th Conference of the Cognitive Science Society*, University of California, Irvine,
657 CA, USA, 15–17, 1985.
- 658 Peltola, O., Vesala, T., Gao, Y., Rätty, O., Alekseychik, P., Aurela, M., Chojnicki, B., Desai, A. R.,
659 Dolman, A. J., Euskirchen, E. S., Friborg, T., Göckede, M., Helbig, M., Humphreys, E., Jackson, R.
660 B., Jocher, G., Joos, F., Klatt, J., Knox, S. H., Kowalska, N., Kutzbach, L., Lienert, S., Lohila, A.,
661 Mammarella, I., Nadeau, D. F., Nilsson, M. B., Oechel, W. C., Peichl, M., Pypker, T., Quanton, W.,

662 Rinne, J., Sachs, T., Samson, M., Schmid, H. P., Sonnentag, O., Wille, C., Zona, D., and Aalto, T.:
663 Monthly gridded data product of northern wetland methane emissions based on upscaling eddy
664 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)
665 1263-2019, 2019.

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

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

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

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

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

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

684 Tian, X., Yan, M., van der Tol, C., Li, Z., Su, Z., Chen, E., Li, X., Li, L., Wang, X., Pan, X., Gao, L.,
685 and Han, Z.: Modeling forest above-ground biomass dynamics using multi-source data and
686 incorporated models: A case study over the qilian mountains, *Agricultural and Forest Meteorology*,
687 246, 1–14, <https://doi.org/10.1016/j.agrformet.2017.05.026>, 2017.

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

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

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

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

709 Wu, Z., Dijkstra, P., Koch, G. W., Peñuelas, J., and Hungate, B. A.: Responses of terrestrial
710 ecosystems to temperature and precipitation change: a meta-analysis of experimental manipulation,
711 *17*, 927–942, <https://doi.org/10.1111/j.1365-2486.2010.02302.x>, 2011.

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

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

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

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

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

730