



## Meta-modelling of carbon fluxes from crop and grassland multi-model outputs

5 Roland Hollós<sup>1,2,3</sup>, Nándor Zrinyi<sup>2,4</sup>, Zoltán Barcza<sup>2,3</sup>, Gianni Bellocchi<sup>5</sup>, Renáta Sándor<sup>1</sup>, János Ruff<sup>6</sup>, Nándor Fodor<sup>1,7</sup>

<sup>1</sup> Centre for Agricultural Research HUN-REN, Agricultural Institute, Martonvásár, 2462, Hungary

<sup>2</sup> Eötvös Loránd University, Department of Meteorology, Budapest, 1117, Hungary

10 <sup>3</sup> Czech Academy of Sciences, Global Change Research Institute, Brno, 603 00, Czech Republic

<sup>4</sup> Eötvös Loránd University, Doctoral School of Earth Sciences, , Budapest, 1117, Hungary

<sup>5</sup> INRAE, VetAgro Sup, Unité Mixte de Recherche sur l'Ecosystème Prairial (UREP), Clermont-Ferrand, 63000, France

<sup>6</sup> University of Pécs, Institute of Mathematics and Informatics, Pécs, 7624, Hungary

15 <sup>7</sup> University of Debrecen, Faculty of Agricultural and Food Sciences and Environmental Management, Debrecen, 4032, Hungary

\* Correspondence to: Nándor Fodor ([fodor.nandor@atk.hu](mailto:fodor.nandor@atk.hu))

20 **Abstract.** We evaluated four stacking-based meta-models - Multiple Linear Regression, Random Forest, XGBoost, and XGBoost with environmental covariates (XGB+) - against the multi-model median (MMM) and best individual process-based models for gross primary production (GPP), ecosystem respiration (RECO) and net ecosystem exchange (NEE) at two cropland and two grassland sites. All meta-models were associated with improved RMSE, bias and correlation, with explained variance gains of ~10–38.5% over MMM, largest for RECO in croplands and  
25 smallest for NEE in grasslands. Bias was nearly eliminated except at one cropland site. SHAP analysis showed that diverse individual models, not always the top performers, contributed most, and that temperature - especially for RECO in croplands and NEE in grasslands - was the dominant environmental driver, while precipitation had minor effects. These findings highlight the predictive and diagnostic advantages of stacking-based approaches over equal-weight MMM, with potential applications across agroecosystem, Earth system and environmental model ensembles.

### 30 1 Introduction

Biogeochemical processes are central to agricultural planning, underpinning concepts such as “climate-smart agriculture”, “low-carbon agriculture” and “greenhouse gas-mitigating farming practices” within narratives/storylines (e.g. Thornton et al., 2018; Hou and Hou, 2019; Anuga et al., 2020). As governments and societies reshape economies, particularly in the context of decarbonisation (e.g. Sroufe and Watts, 2022),  
35 biogeochemical modelling has emerged as an essential tool to support agricultural policies (Bellocchi, 2023). These models simulate the complex interactions between agriculture and ecosystem services, such as C sequestration, biodiversity conservation and water quality regulation, thereby empowering policymakers to integrate ecosystem value into agricultural decision-making and land-use strategies (e.g. Lambin and Meyfroidt, 2011). Moreover, by simulating a wide variety of ecosystem processes, the models facilitate the assessment of agricultural emissions and  
40 mitigation strategies, forming a scientific basis for climate policy and agricultural resilience (Li et al., 2006; Valin et al., 2013; Sándor et al., 2018; Lembaid et al., 2021, 2022; Gascuel-Odoux et al., 2022).



Many state-of-the-art biogeochemical models have a long development history with some of them spanning 50 years (Hidy et al., 2022). Still, inherent uncertainties in model structures, parameterisations and assumptions pose challenges for reliable predictions under diverse environmental and management conditions (e.g. Riccio et al., 2007; 45 Bellocchi et al., 2010; Therond et al., 2011; Harrison et al., 2012; Brilli et al., 2017; Bilotto et al., 2021). This calls for improved mathematical tools and innovative solutions to strengthen the trust in the models.

Ensemble modelling has gained prominence as a robust method to address these uncertainties in biogeochemical models (e.g. Challinor et al., 2013; Snow et al., 2014; Calanca et al., 2016; Jones et al., 2017; Knutti et al., 2019). Each model adopts unique assumptions about processes like soil nutrient cycling, photosynthesis, allocation and 50 crop phenology, leading to significant variability in predictions under similar scenarios. Such structural uncertainties (e.g. divergent representations of water stress or heat tolerance), and also parameter uncertainty can lead to inconsistent predictions. Ensemble approaches mitigate these discrepancies by synthesising outputs from multiple models, compensating biases and highlighting consistent trends (e.g. Bassu et al., 2014; Rosenzweig et al., 2014; Kollas et al., 2015; Li et al., 2015; Ruane et al., 2016, 2017; Sándor et al., 2017, 2020; Ehrhardt et al., 2018; Wallach 55 et al., 2018). This increases the reliability of predictions and reduces the risks associated with over-reliance on individual models.

Although ensemble techniques are proven to improve modelling accuracy and increase the trust in biogeochemical models, problems still exist. For example, the ensemble models are hard to explain, and currently only the simplest methods (averaging or median calculation; Sándor et al., 2016) are the most widely used for ecosystem models. 60 These simple model combination algorithms usually have many implicit presumptions that models often cannot meet (for example, same variance and independence). In the field of machine learning (ML), which generates models directly from data, instead of following physical, biological and biogeochemical principles, much more developed techniques exist, which can be potentially exploited to improve the predictive power of the ensemble output. These techniques usually have less and also more explicit conditions (Scowen et al., 2023).

As an advanced, potentially promising approach, ML based meta-modelling extends the concept of ensemble technique by synthesising outputs even across multiple calibration scenarios. Recent advancements in ML offer innovative tools to refine data-intensive modelling in Earth sciences (e.g. Li et al., 2015; Jackson et al., 2017; Keskin et al., 2019; Reichstein et al., 2019; Bai et al., 2021; Chandel et al., 2024; Wang et al., 2024). By using ensemble learning – a subset of ML methods – predictive frameworks can integrate diverse model outputs with improved 70 accuracy and interpretability (e.g. Hansen and Salamon, 1990; Opitz and Maclin, 1999; Dietterich, 2000; Hagedorn et al., 2005; Palmer, 2019). Stacking, as a flexible ensemble learning technique, combines outputs from heterogeneous models into a meta-model, assigning weights to highlight their relative importance (e.g. Breiman, 1996; Van der Laan et al., 2007; Sagi and Rokach, 2018). This approach has potential for biogeochemical modelling, offering simplicity and adaptability with Multiple Linear Regression (e.g. Kutner et al., 2005), and enhanced prediction 75 accuracy through decision-tree-based methods like Random Forest (e.g. Breiman, 2001a,b; Liaw and Wiener, 2002) and XGBoost (e.g. Chen and Guestrin, 2016). In spite of the potential in stacking methods, up to the knowledge of the authors it was not used so far as a stacking ensemble method, but instead some papers used it for predicting the output from the driving input data directly.

The integration of diverse environmental drivers with sophisticated analytical methods is central to advancing 80 predictive modelling in environmental science. For instance, Hengl et al. (2017) demonstrated that spatial predictions of soil properties improved significantly when large stacks of remote-sensing and environmental



covariates were incorporated into ensemble machine-learning frameworks. This approach is further supported by research such as Luo et al. (2009), who found that high-frequency driver data reduces parameter equifinality in ecosystem data assimilation, and Pappas et al. (2014), who developed efficient methods for gap-filling  
85 hydrometeorological observations. More recently, Sándor et al. (2023) showed that analysing residual correlations among crop and grassland model ensembles can reveal structural model deficiencies and improve simulation of C-N fluxes when combined with ensemble averaging. Collectively, these studies highlight the substantial value of leveraging diverse variables, advanced techniques, and residual-based ensemble diagnostics to build more robust and accurate predictive models. Building on these advancements, the present study introduces a meta-modelling  
90 framework that integrates biogeochemical model outputs and environmental variables, guided by residual correlation insights, to address structural model error and enhance predictive performance.

This study aims to improve the predictive accuracy of biogeochemical flux calculations, particularly C fluxes, in diverse agricultural systems (crop rotations and grassland systems). Building on an already established multi-model framework, a novel calibration methodology is used to develop a meta-model that balances scientific rigour and  
95 practical feasibility. The results of this meta-model are compared with the multi-model ensemble median and other meta-models to demonstrate the improvement, interpretability and reliability of ensemble predictions. This dual focus aims to improve methodologies for sustainable crop and grassland management and to inform agricultural policy to better address the challenges of climate resilience and sustainability.

## 2 Materials and methods

### 100 2.1 Meta-modelling framework

The methods employed in this study arise from the contemporary landscape of crop and grassland modelling, specifically focusing on ensemble modelling. We adopt and extend the concept of meta-modelling by comparing it to the multi-model median (MMM) estimator (e.g. Sándor et al., 2018).

While the term *meta-modelling* can be conveyed ambiguously, some clarification is needed. Widely used in  
105 mathematics, computer science and engineering, a meta-model is understood as a model of one or more models. Here, we define it as the use of model outputs from a multi-model ensemble - for one particular method supplemented by environmental variables - as inputs to a higher-level statistical model. Unlike traditional surrogate models that approximate the structure of a process-based model, our meta-model approach complements, rather than replaces, process-based models by exploiting the information embedded in multiple model outputs. The enhanced  
110 predictive power and the resulting trust in model outcomes are particularly valuable for informing policy decisions. This method goes beyond Bayesian model averaging, non-democratic model selection and some machine learning based output combination techniques, allowing for more flexibility in capturing nonlinear relationships.

The foundation of the method is the so-called “No Free Lunch” (NFL) theorem (Wolpert and Macready, 2005), which formally states that when evaluating performance over the entire space of possible tasks (assuming a uniform  
115 distribution), all machine learning models achieve the same average result. This implies that no single algorithm is universally superior. Crucially, however, real-world problems deviate from this uniform distribution, displaying inherent structure and regularities that models with well-suited inductive biases can effectively exploit.



We adopt *stacked generalization* (stacking), as an ensemble method wherein a meta-model integrates predictions generated by multiple base models. This method relies on the observation that diverse models inherently capture  
120 different aspects of real-world systems like the plant-soil system. Consequently, training a metamodel to optimally integrate these varied predictions allows stacking to leverage the individual strengths of the base models, ultimately leading to improved generalisation across various tasks focusing on the plant-soil systems.

Linear models are the most prevalent meta-models in ensemble learning, largely due to their simplicity, computational efficiency, and ease of interpretation. Their additive structure makes it straightforward to quantify  
125 the relative influence of individual models within an ensemble, which is particularly valuable for diagnostic purposes. Notably, a precursor to formal meta-modelling can be seen in the Granger–Ramanathan averaging method (Granger and Ramanathan, 1984; Nand et al., 2025), which introduced the idea of optimally weighting model outputs through constrained linear regression. This approach laid the foundation for modern stacking techniques by demonstrating how combining forecasts in a regression framework could systematically improve predictive  
130 accuracy.

However, when base models exhibit substantial structural differences and higher structural errors, the relationship between predicted and observed outcomes may no longer be well-approximated by a linear model. In such scenarios, simple linear approaches like the Granger–Ramanathan averaging may prove insufficient. Furthermore, the often-implicit assumption of conditional independence among base-model predictions in linear meta-modelling is  
135 frequently violated, rendering multiple linear regression a suboptimal choice.

To address this limitation, ensemble-based machine-learning meta-models offer a compelling alternative. These models can effectively capture complex, non-linear relationships and do not rely on the conditional independence of base-model predictions. Random Forests, for example, have been among the earliest and most effective ensemble methods used as meta-models, particularly in contexts where linear approaches underperform (Zhao and Cheng,  
140 2022).

While effective, Random Forests can be computationally demanding during both training and inference. They also tend to overfit the training data, potentially compromising generalization on independent data. XGBoost offers greater flexibility, often delivering superior predictive accuracy with generally lower computational requirements. However, XGBoost typically requires extensive hyperparameter tuning to achieve optimal performance.

145 A key limitation of all these approaches is their static nature: the functional relationship learned between the base-models' outputs and the target variable remains fixed. This static feature limits the long-term usability of meta-models, as evolving environmental conditions can shift the relevance of different structural components within the base-models. Consequently, the optimal meta-model may vary over time, necessitating continuous retraining to maintain accuracy, which represents a major drawback.

150 With advancements in machine learning, today we can significantly extend the approach described above. By incorporating differentiating environmental factors — such as precipitation — as additional features, the meta-model can learn not only the functional relationship between the base-model outputs and the measurements, but also how this relationship varies with environmental conditions.

This approach is only possible with flexible meta-models. Simpler models, such as generalized linear models, rely  
155 on assumptions like conditional independence and uncorrelated features — assumptions that are often violated when additional environmental factors are introduced. As a result, these simpler models may struggle to capture the conditional dependencies required for enhanced predictive performance.



Such an *environment-aware meta-modelling* framework (i.e. XGBoost with additional meteorological drivers) improves both predictive accuracy and interpretability by revealing the specific conditions under which different models succeed or fail. Examining this meta-model allows researchers and modellers to understand their model's limitations in varying contexts. This knowledge guides further model development and facilitates the selection of the most appropriate model for given circumstances, promoting valuable synergy between modelling approaches.

## 2.2 Source of the model ensemble

The model ensemble is based on data from international initiatives, primarily the Agricultural Model Intercomparison and Improvement Project (<https://agmip.org>) and the Integrative Research Group of the Global Research Alliance on agricultural GHGs (<https://globalresearchalliance.org/research/integrative>). These exercises have shown that ensembles of process-based biogeochemical models can reliably estimate agricultural productivity, as well as C and N emissions (and stocks) of agricultural systems (Ehrhardt et al., 2018; Sándor et al., 2018, 2020, 2024). These studies, mostly funded by national agencies, also contribute to the assessment of C storage potential (e.g. Farina et al., 2021), aligning with the '4 per mille Soils for Food Security and Climate' initiative established at the 2015 United Nations Climate Change Conference (COP21) by the French Ministry of Agriculture (<https://www.4p1000.org>).

Here, we refer to a multi-model scheme, using daily model outputs as inputs for the meta-model. The multi-modelling exercise discussed here is not the primary content but serves as the foundational basis for the subsequent meta-modelling exercise. We are revisiting salient elements from previously published studies to build the framework for our meta-modelling analysis. Specifically, we delve into comparisons with the multi-model median from the study by Sándor et al. (2020) on C fluxes, which, in turn, was based on the protocol of Ehrhardt et al. (2018).

**Table 1:** C-flux outputs provided by different models (Sándor et al., 2024), denoted by symbols such as ✓ (present) and NA (not available). Models are marked as M01-M26, and they are kept anonymous as in the Sándor et al. (2020) study. In cropland sites, we had gross primary production (GPP) from six models, net ecosystem exchange (NEE) from seven models and ecosystem respiration (RECO) from 15 models. At grassland sites, we had GPP from 10 models, NEE from 10 models and RECO from 12 models.

Model type	Crop models												Grassland models								Both systems			
	Model id	M01	M02	M04	M09	M12	M13	M18	M19	M20	M25	M26	M03	M06	M16	M21	M22	M23	M24	M28	M05	M07	M08	M14
Outputs	GPP	✓	NA	NA	✓	NA	NA	NA	✓	NA	NA	NA	NA	✓	✓	✓	✓	✓	✓	✓	✓	✓	NA	✓
	RECO	✓	✓	✓	✓	✓	✓	✓	✓	NA	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
	NEE	✓	NA	NA	✓	NA	NA	NA	✓	✓	NA	NA	NA	✓	✓	✓	✓	✓	✓	✓	✓	✓	NA	✓

We use daily outputs from 23 crop and grassland simulation models/versions (Table 1; for details see Sándor et al., 2020; 2024). Model names (M01–M28) were anonymised for consistency, with M11 excluded due to missing C fluxes. These models were applied to four long-term field sites: two grazed grasslands (G3, G4) and two croplands (C1, C2) across the UK, France (two sites) and Canada (Table 2). In the original Sándor et al. (2020) paper there



was one additional cropland site (India) that was rejected from this study due to the relatively poor temporal coverage of the observation data.

**Table 2:** Cropland and grassland sites, and temporal coverage of available data used for the analysis. Different crop rotations were used in the cropland sites, including cereals (spring and winter wheat [W], triticale [T], maize [M] and rice [R]), legumes (soybean [S]), rapeseeds (canola and mustard [C]), borages (phacelia, P) and fallow intercropping periods [I].

Sites, country (latitude, longitude, altitude)	Years of available data (simulation period)	Land use	Annual mean temperature (°C)	Annual mean precipitation (mm)
C1: Ottawa, Canada (45.29, -75.77, 94 m a.s.l.)	2007-2012	W/S/C/M/W/C	7.2±0.9	936±121
C2: Grignon, France (48.85, 1.95, 125 m a.s.l.)	2008-2012	C/M/W/T/P/M /W/I	10.9±0.7	571±35
G3: Laqueuille, France (45.64, 2.74, 1040 m a.s.l.)	2003-2012	Permanent grassland (cattle grazing)	13.7±0.4	910±96
G4: Easter Bush, United Kingdom (55.52, -3.33, 190 m a.s.l.)	2002-2010	Permanent grassland (ewe grazing)	7.8±0.8	1078±205

High quality data covering climate, soil, agricultural practices and C fluxes were gathered from Sándor et al. (2024). Observations at these sites include eddy covariance and chamber measurements of net ecosystem CO<sub>2</sub> exchange (NEE) data, further divided into two main fluxes (e.g. Reichstein et al., 2005; Raj et al., 2016): gross primary production (GPP), representing photosynthetic production from atmospheric CO<sub>2</sub>, and total ecosystem respiration (RECO), encompassing the total C respired by plants, soil organisms and, in the case of grasslands, grazing animals. Initialisation and calibration procedures aligned models with vegetation, soil and atmospheric fluxes from the study sites, following the protocol described in Ehrhardt et al (2018). This comprehensive exercise involved a multi-stage approach, granting modellers access to increasingly detailed data for running and evaluating their models, progressing from uncalibrated to fully calibrated simulations. In summary, the calibration included five ascending levels, incorporating additional data for refinement: blind test (S1), utilising only site specific weather and management data for the simulation periods; initialisation (S2), incorporating additional historical climate and management data, extending to years preceding simulation periods, along with regional productivity for initialisation purposes; partial calibration (S3), including biomass production and phenology data; intermediate calibration (S4), integrating soil temperature, soil water content and mineral N data; full calibration (S5), using N<sub>2</sub>O emissions, NEE, GPP, RECO and soil organic C and N data.

For the purposes of this study, emphasis was placed on the outputs from the partial calibration stage (S3) stage. S3 involves calibration using plant data exclusively, which enhances the practical implementability of models for end users and beneficiaries. This approach recognises the importance of ensuring models are not only scientifically

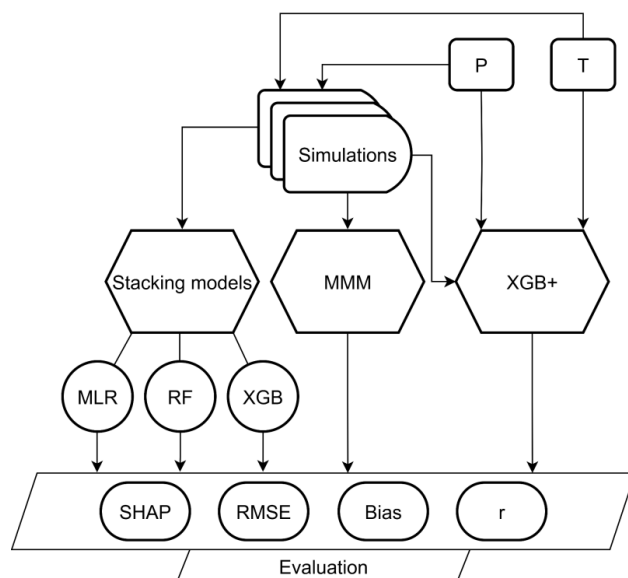


robust but also practically useful. Its validity is supported by findings that additional calibration stages beyond S3 yielded only modest improvements (Sándor et al., 2023).

### 2.3 Meta-modelling approaches

220 The meta-model construction was done at the daily time resolution for GPP, RECO and NEE. Two main approaches  
were tested in this study for the meta-model construction. The first one focuses on the classic stacking ensemble  
method (referred here as CSEM), where the meta-models were constructed using multiple linear regression (MLR),  
random forest (RF) method, and XGBoost (XGB) using the model outputs described in Section 2.2. The second  
225 approach, that is introduced in this study as a *novel method*, uses the same stacking methods but extends the input  
dataset so that besides the individual, process-based model outputs, we use two additional environmental factors  
(temperature and precipitation) to predict the final outcome (i.e. the meta-model). This approach (i.e. inclusion of  
the meteorological drivers) is referred to here as extended generic stacking ensemble method (XGB+). The inclusion  
of only temperature and precipitation is grounded in their central roles as primary climatic drivers of terrestrial C  
230 fluxes. Temperature governs essential biological processes such as photosynthesis and respiration (Lloyd and  
Taylor, 1994; Xu and Baldocchi, 2004), while precipitation affects soil moisture and plant water stress, both critical  
for GPP and RECO (Reichstein et al., 2005; Schwalm et al., 2010). Extensive research has shown that fluctuations  
in these two variables explain much of the interannual and seasonal variability in ecosystem C exchange (Jung et  
al., 2007; Beer et al., 2010). Moreover, temperature and precipitation are consistently measured across sites, widely  
available and commonly integrated into ecological and climate models. Their inclusion enhances predictive  
235 performance without adding complexity or compromising generalisability (e.g. Hijmans et al., 2005). Supporting  
this approach, Dorman et al. (2013) advocate for the use of core climatic variables over an excess of redundant  
predictors to maintain model robustness and transferability.  
All meta-models were constructed using 70% of the observations for training, and the remaining 30% was used for  
validation.

240 For the MLR method we used the *lm* function of base R. The RF method (*randomForest* package from CRAN;  
Liaw and Wiener, 2002) was used with 1000 trees. For each tree the number of predictors were 1/3 of the number  
of input data streams (GPP, RECO, NEE and meteorological drivers if applicable) in order to prevent overfitting.  
The subset of predictors were used randomly. The minimum size of the terminal nodes was five. We have split the  
dataset into training and testing dataset randomly, using 70% and 30% of the data for training and validation,  
245 respectively. The XGB method was implemented using the *XGBoost* package from CRAN (Chen and Guestrin,  
2016). For hyper-parameter optimization we used grid-search (in other words systematic search) technique for  
XGBoost. Figure 1 shows the overview of the stacking methods used in the study.



**Figure 1:** Schematic illustration of the applied stacking ensemble learning methodology. The workflow combines  
250 outputs from multiple simulations with multi-model median (MMM) and stacking models (MLR, RF, XGB). The  
XGB+ model is an extended version that incorporates additional input features: temperature (T) and precipitation  
(P). The performance of these models and the ensembles are evaluated using four metrics: SHAP (for  
interpretability), RMSE (root mean square error), Bias and the correlation coefficient (r). The final output is the site-  
level carbon flux, which is the same variable type as the initial model outputs.

255

To facilitate the interpretation of the linear meta-model, we examined the normalised feature weights, which provide  
a direct measure of the relative importance of predictors in a linear framework. For the non-linear ensemble models  
(RF and XGBoost), feature contributions were quantified using SHAP (SHapley Additive exPlanations) values  
260 (Shapley, 1953; Lundberg and Lee, 2017). SHAP values are grounded in cooperative game theory, representing the  
mean marginal contribution of a feature across all possible feature coalitions. This formulation provides a  
theoretically consistent and locally accurate decomposition of model predictions, enabling a granular understanding  
of feature effects. In contrast to traditional importance measures such as the mean decrease in impurity (often  
referred to as the Gini importance; Breiman, 2001a,b), SHAP values address several critical limitations. The Gini  
265 metric, while computationally efficient, is known to exhibit biases toward features with a higher number of  
categories or greater variance (Strobl et al., 2007). Moreover, it does not provide insight into the directionality or  
context-dependent contribution of features to individual predictions. SHAP overcomes these limitations by  
providing additive, model-agnostic attributions that remain consistent across different models and capture complex,  
non-linear interactions between features (Lundberg et al., 2020). These properties make SHAP particularly well-  
270 suited for interpreting ensemble methods like RF and gradient-boosted decision trees, where feature interactions and  
non-linearities are prominent.

All evaluations were performed separately for each study site, and for the studied carbon fluxes (GPP, RECO, NEE).  
No cross-site meta-model was developed in this study. To assess how well models approximate site-level



275 observations, we compared each meta-model and the MMM using three common performance metrics (Richter et al., 2012) across 26 site-stage-output combinations: the root mean square error (RMSE), square of the Pearson's correlation coefficient ( $R^2$ ), and the bias.

The results are organized according to the carbon flux type. First GPP is discussed, as it represents plant photosynthesis and is the most important driver of the plant production and yield. Most likely GPP is the process that is the simplest to simulate by process-based models. It is due to the fact that the mechanism of photosynthesis is well-discovered, and practically it is separated from other autotrophic or heterotrophic processes (since this is the only flux that comes into the ecosystem). GPP is a relatively large flux (compared e.g. to NEE or net biome production), so model optimization is perhaps simpler and more interpretable. Of course GPP can be biased for some cases due to improper representation of phenology, and lack of model optimization.

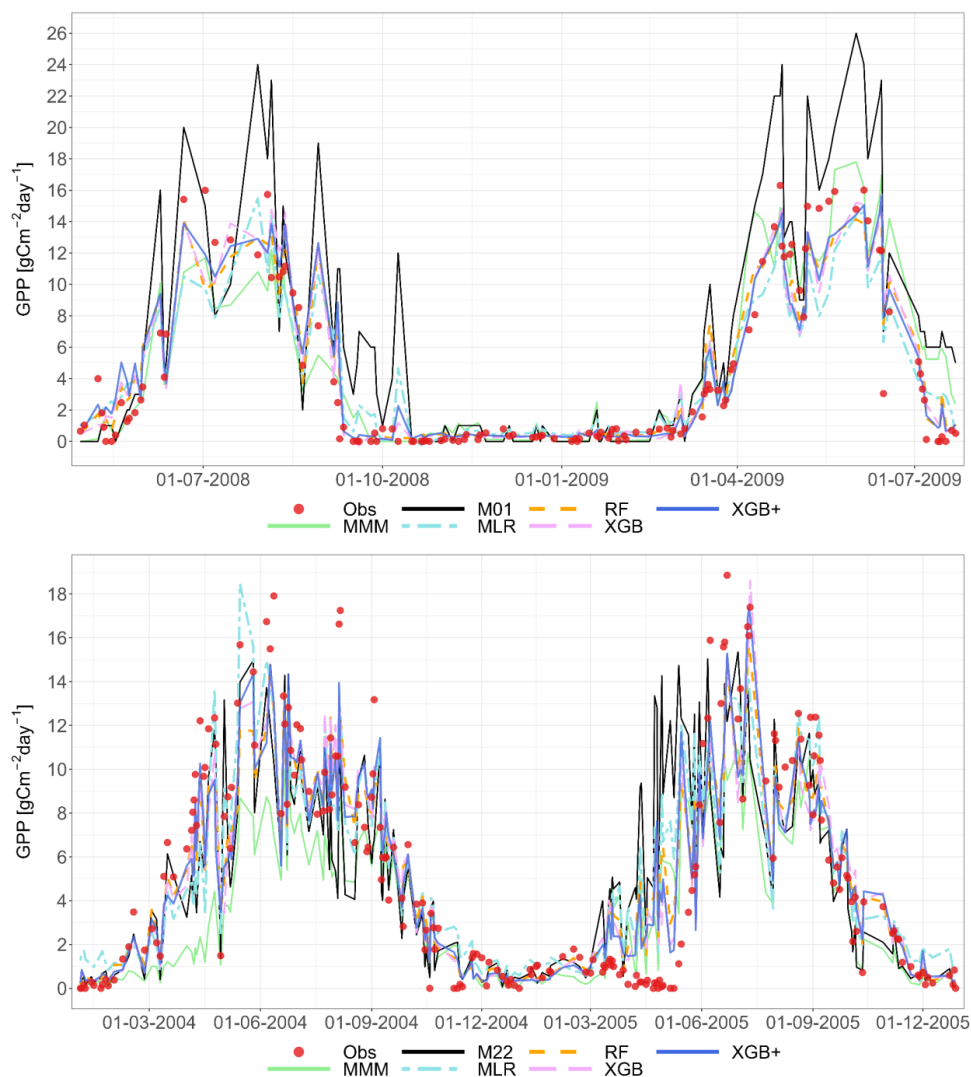
280 RECO follows the discussion that is typically harder to simulate than GPP. It is due to the fact that RECO is the sum of autotrophic and heterotrophic respiration, where both components have their own intrinsic uncertainties. For example, heterotrophic respiration is typically problematic due to the complexity of the decomposition processes in the soil including the rhizosphere and the presence, complexity and activity of the microbial/fungal life forms. RECO can be erroneously simulated if one of the two major components are misrepresented (or even worse if, due to the compensation of errors, RECO looks good due to wrong reasons). Nevertheless, RECO is still a larger carbon flux, so its representation can be expected to be relatively well-adjusted.

285 Finally, we present NEE related results. This flux is probably the hardest to simulate, since it is the sensitive balance between RECO and GPP (by definition  $NEE = RECO - GPP$ ). Capturing the magnitude and variability of NEE requires the proper representation of plant phenology (start and end of season), and accurate representation of GPP and RECO. If the meta-model captures NEE with greater precision than the more traditional model ensembles, this represents a significant accomplishment. Although NEE is defined as  $RECO - GPP$ , in this study NEE is not calculated with this definition but rather it is modelled independently from GPP and RECO (due to methodological reasons). This might introduce some inconsistency, but the aim was to check the ability of the metamodeling approaches for any given variable, independent of other observations. Successful and unbiased simulation of NEE is a major step forward supporting e.g. atmospheric inversions or any other bottom-up carbon flux estimations.

## 300 3. Results

### 3.1. GPP

305 Fig. 2 shows a representative example of the observed and simulated time series of GPP at the Grignon cropland site (C2) for the 2008-2009 growing seasons (maize and after that winter wheat), and for two consecutive years at the Easter Bush permanent grassland site (G4) (2004-2005). All meta-modelling approaches are plotted alongside the best-performing individual model (M01 for Grignon and M22 for Easter Bush). Appendix A contains the complete simulated dataset for all sites, and for all years.



**Figure 2:** Performance of the multimodel median (MMM, green), the constructed meta-models and the best-performing individual models for simulating GPP. The top panel shows results for the Grignon cropland site (C2), which includes maize and winter wheat. The bottom panel shows two years of data for the Easter Bush grassland site (G4). Observations are marked by red dots. The meta-models include Multiple Linear Regression (MLR, hatched light blue), Random Forest (RF, hatched orange), XGBoost (XGB, hatched purple), and XGBoost+ (XGB+, blue). The best-performing individual models (M01 at C2, M22 at G4) are shown in black. Dates are provided in the format of dd-mm-yyyy.

The plots indicate a relatively large scatter in model results compared to observations, especially for croplands. Notably, the best-performing individual model for the cropland time series exhibited considerable biases, generally overestimating GPP. In contrast, all meta-model approaches indicate reduced bias. However, to provide objective

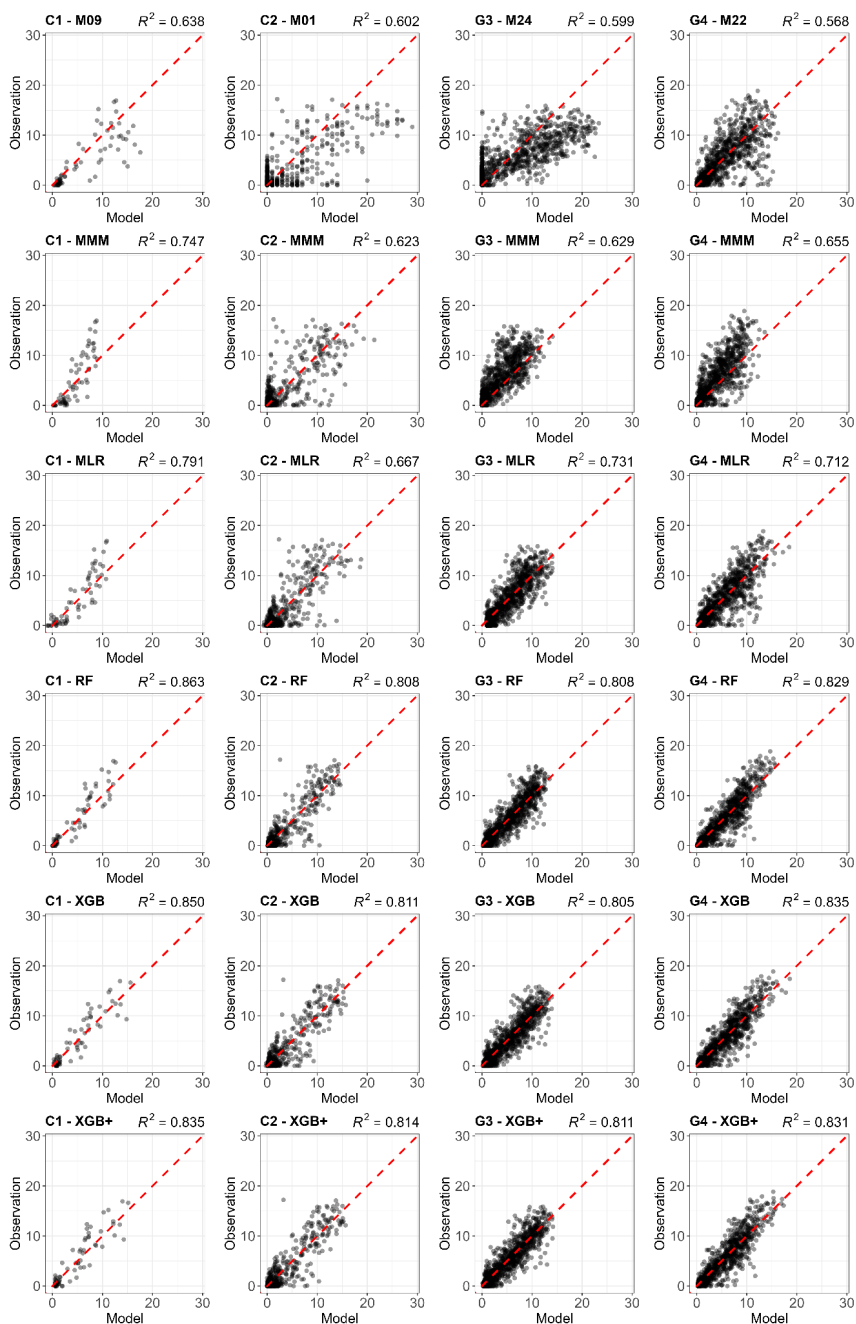


320 quality indicators, further quantitative analysis was performed. Fig. 3 compares observed and modelled GPP for all  
four experimental sites and all metamodel construction methods used in the study, including the MMM. The figure  
includes a scatterplot showing the performance of the best individual model, selected based on its RMSE score.  
RMSE was chosen as the primary selection metric because it correlates well with other performance metrics  
(Kobayashi and Salam, 2000; Robeson and Willmott, 2023). The figure was constructed using all data from the  
325 sites, not only training data.

The figure shows a gradual improvement in performance from top to bottom, indicated by increases in  $R^2$  and tighter  
alignment of data points along the 1:1 line. At every site, MMM outperforms the best individual model in terms of  
explained variance. For sites C2 and G3, MMM corrects much of the bias seen in the individual simulations.

At site C1, improvements across MLR, RF, XGB and XGB+ do not always correspond to continuous increases in  
330  $R^2$ , but the alignment of the observation-model data pairs improves relative to the 1:1 line. For C2, explained  
variance steadily increases with a marked improvement upon introduction of RF. Similar trends are observed for the  
grassland sites G3 and G4, where the longer time series yield a larger number of data points and a more pronounced  
improvement. Overall, across all sites, explained variance increases by ~20% for the best performing metamodels  
(RF, XGB and XGB+).

335



**Figure 3:** Comparison of the best individual model, the constructed meta-models and the traditional Multi-Model Median with the observations for all sites (from left to right C1, C2, G3 and G4) and for the entire time series for GPP. From top to bottom: best individual model with ID, MMM, MLR, RF, XGB and XGB+. All units are in  $\text{g C m}^{-2} \text{ day}^{-1}$ . Red dashed line represents the 1:1 relationship.

340  $\text{m}^{-2} \text{ day}^{-1}$ .



Table 3 presents model statistical measures for all sites and modelling approaches evaluated, alongside the best-performing individual model, using data only from the validation subset.

345

**Table 3:** Statistical evaluation of the best performing individual model, the multi-model median (MMM) and the applied meta-models (MLR, RF, XGB and XGB+) with respect to three performance metrics for GPP: root mean square error (RMSE), bias and Pearson’s correlation coefficient (r). Only validation data were used for the calculation of the statistics. RMSE and BIAS are provided in  $\text{g C m}^{-2} \text{day}^{-1}$  units.

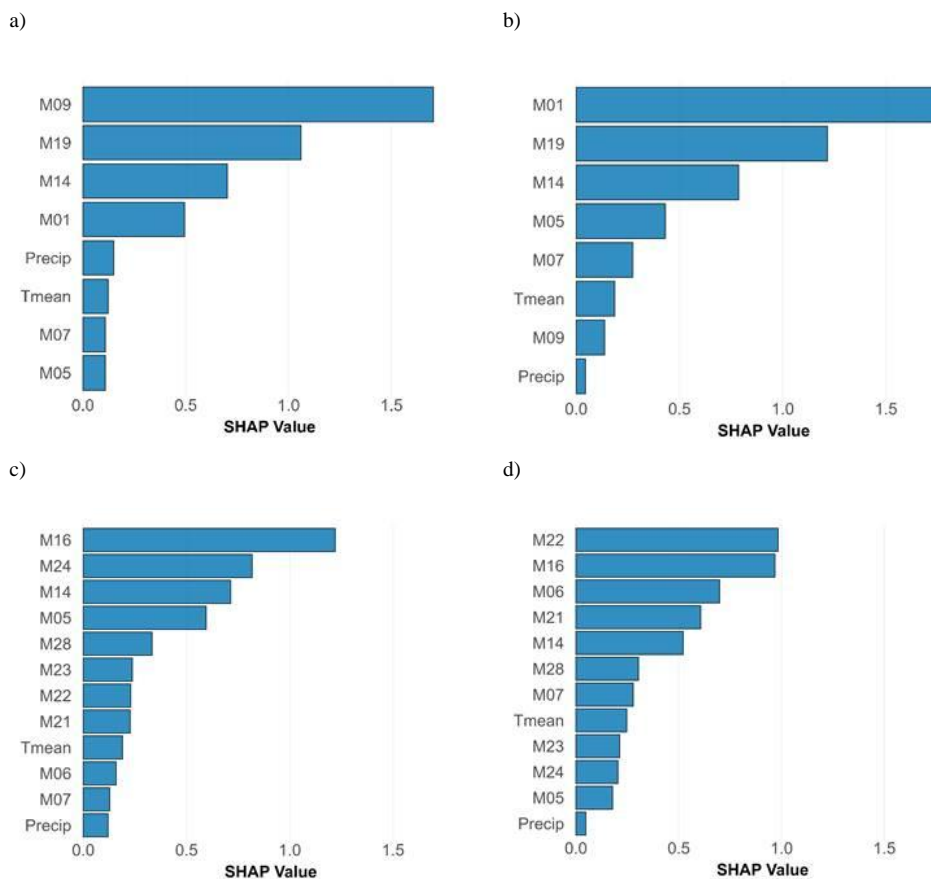
Site	Metric	best individual model	MMM	MLR	RF	XGB	XGB+
C1	RMSE	3.65	3.03	2.41	1.98	1.93	2.05
	BIAS	1.52	-0.61	-0.56	-0.51	-0.19	-0.28
	r	0.799	0.864	0.889	0.929	0.922	0.914
C2	RMSE	4.22	2.88	2.55	1.93	1.92	1.90
	BIAS	1.02	-0.05	0.16	0.1	0.12	0.15
	r	0.776	0.790	0.817	0.899	0.901	0.902
G3	RMSE	4.48	2.96	2.12	1.79	1.80	1.78
	BIAS	1.22	-1.6	0.02	0.03	0.06	0.05
	r	0.774	0.793	0.855	0.899	0.897	0.901
G4	RMSE	3.02	2.96	2.37	1.84	1.80	1.82
	BIAS	0.11	-1.36	0.04	0.02	0.02	0.03
	r	0.754	0.809	0.844	0.911	0.914	0.912

350

The table highlights the improved performance of meta-models relative to both the best individual model and the MMM, which was previously published in Sándor et al. (2020). Notably, moving to more advanced meta-models results in a substantial reduction in both RMSE and absolute bias. For the two grassland sites the meta-models provide almost bias-free estimations, while for the crop sites bias is low but not approaching zero. In addition, explained variance typically increases by about 10% for the best-performing metamodel compared to MMM (but only 6% at C1).

355

360



**Figure 4:** The SHAP values of the XGB+ meta-model for GPP. Larger values mean stronger contribution to the resulting GPP. a) C1 site; b) C2 site; c) G3 site; d) G4 site. Tmean stands for daily mean temperature, and Precip is daily precipitation.

365 Given the novelty of the XGB+ approach and its occasional superior performance over all other approaches, here  
we analyse its functioning using SHAP values of the input data streams (individual models, plus temperature and  
precipitation). SHAP value analysis (Fig. 4) shows that certain models (e.g. M09, M19, M14 at C1; M01, M19 at  
C2; M16, M24 at G3; M22, M16 at G4) dominate contributions to GPP prediction across sites. Temperature is also  
identified as a notable predictor, especially at sites C1 and C2. Precipitation consistently exhibits minimal impact  
370 on GPP prediction across all sites, with minor contributions only at C1 and G4.

### 3.2 RECO

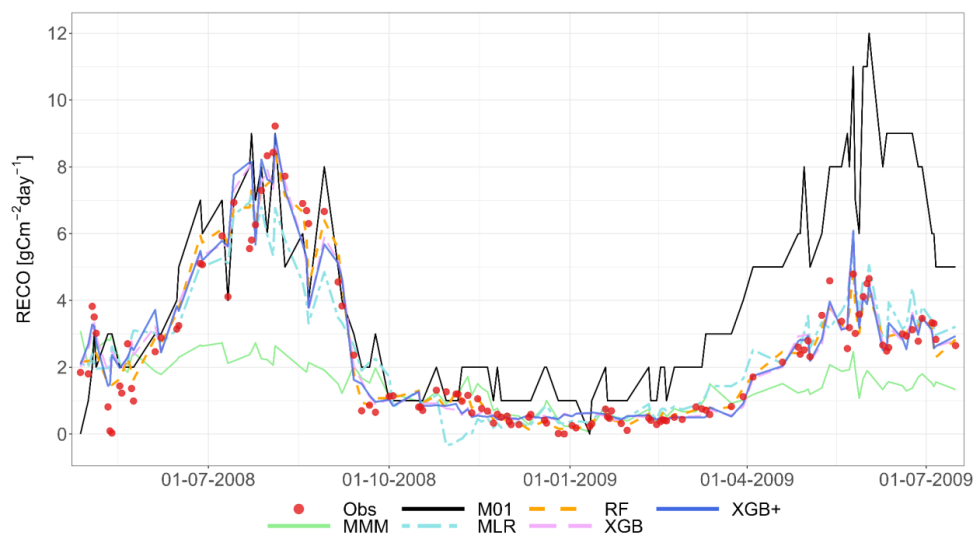


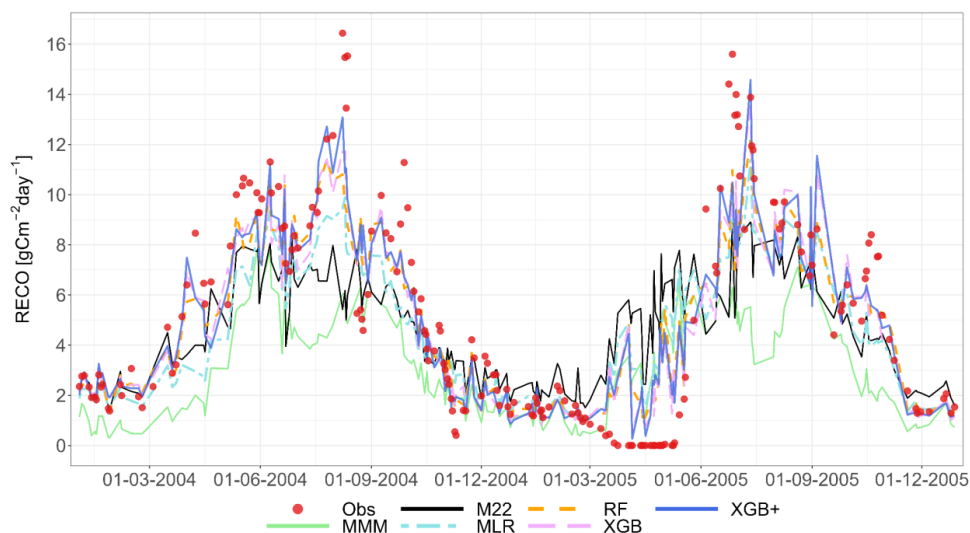
Fig. 5 presents a two-year comparison of simulated and observed RECO for Grignon (C2) and Easter Bush (G4). The graphs display the results of the meta-models, the MMM and the best-performing individual models (M01 and M22). Appendix A contains the complete simulated dataset for all sites, and for all years.

375 At both sites, the MMM consistently underestimates respiration, particularly during peak periods. This bias is especially pronounced at G4, where the MMM produces a negative peak around June-July, resulting in a poor fit to the observations. In contrast, the meta-models and the best individual model successfully follow the positive peak of the measured values at this site. Visual inspection suggests that the XGB and XGB+ meta-models were the most accurate in simulating the maximum values for G4. They were also among the best performers for C2, alongside RF and MLR. However, all meta-models at G4 exhibited a negative bias during the observation peaks. Regarding the individual models, the best-performing model at C2 (M01) accurately simulates respiration for maize but strongly overestimates it for winter wheat. At G4, the best individual model (M22) greatly underestimates the maximum values for both years, though its performance is still better than the MMM. This visual analysis suggests that the

380

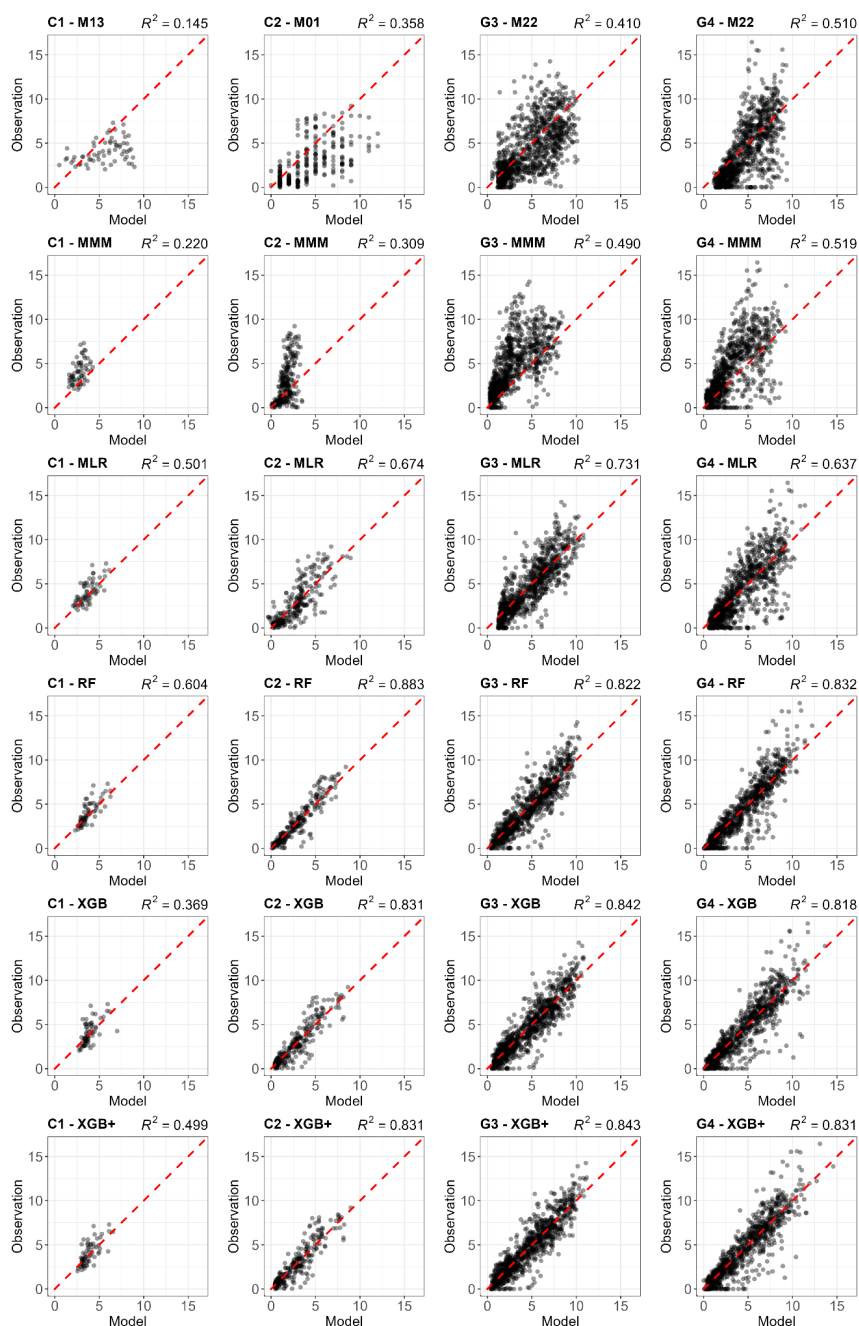
385





390 **Figure 5:** Performance of the multimodel median (MMM, green), the constructed meta-models and the best-  
performing individual models for simulating RECO. The top panel shows results for the Grignon cropland site (C2),  
which includes maize and winter wheat. The bottom panel shows two years of data for the Easter Bush grassland  
site (G4). Observations are marked by red circles. The meta-models include Multiple Linear Regression (MLR,  
395 hatched light blue), Random Forest (RF, hatched orange), XGBoost (XGB, hatched purple), and XGBoost+ (XGB+,  
blue). The best-performing individual models (M01 at C2, M22 at G4) are shown in black. Dates are provided in  
the format of dd-mm-yyyy.

A more quantitative analysis can be performed by inspecting Fig. 6, which presents scatterplots comparing simulated  
400 and observed RECO. In terms of  $R^2$  values (as shown in Fig. 6), the MMM generally outperforms the best-  
performing individual models, with the exception of C2 (a trend also visible in Fig. 5). However, both the MMM  
and the individual models exhibit low  $R^2$  values, especially at the crop sites, indicating that they explain a limited  
part of the variance. Among the meta-models, RF shows the strongest correlation with observed data at the crop  
sites, while at the grassland sites, its performance is comparable to XGB and XGB+. The MLR consistently provides  
405 a “middle-ground” performance. The scatter plots in Fig. 6 distinctly illustrate that the MMM has a strong negative  
bias, consistently underestimating the observed RECO. This bias is eliminated even by the simplest meta-model  
(MLR), and as more sophisticated meta-model techniques are employed, the alignment with the 1:1 line becomes  
significantly better.



410 **Figure 6:** Comparison of the best individual model, the multi-model median and the constructed meta-models with observations for RECO. Each row represents a different model type, and columns correspond to the sites (from left to right: C1, C2, G3, and G4). The top row shows the best individual models with their identifiers (M13 at C1, M01 at C2, M22 at G3 and G4). The remaining rows show the MMM, MLR, RF, XGB and XGB+. All units are in  $\text{g C m}^{-2} \text{ day}^{-1}$ . The red dashed line represents the 1:1 relationship.



415

**Table 4:** Statistical evaluation of the best-performing individual model, the multi-model median (MMM) and the applied meta-models (MLR, RF, XGB and XGB+) for RECO. Three performance metrics are used: root mean square error (RMSE), bias and Pearson’s correlation coefficient (r). Only validation data were used for the calculation of the statistics. RMSE and BIAS are provided in  $\text{gC m}^{-2} \text{day}^{-1}$  units.

Site	Metric	Best individual model	MMM	MLR	RF	XGB	XGB+
C1	RMSE	2.61	1.72	0.96	0.84	1.03	0.92
	BIAS	1.56	-1.29	-0.32	-0.23	-0.14	-0.17
	r	0.381	0.469	0.707	0.777	0.607	0.707
C2	RMSE	2.72	2.24	1.30	0.79	0.94	0.93
	BIAS	1.48	-1.09	0.01	0.01	0.03	0.05
	r	0.599	0.556	0.821	0.940	0.912	0.912
G3	RMSE	2.42	2.72	1.56	1.27	1.19	1.19
	BIAS	0.44	-1.69	0.08	0.06	0.05	0.02
	r	0.641	0.700	0.855	0.907	0.918	0.918
G4	RMSE	2.36	2.52	1.97	1.36	1.40	1.35
	BIAS	0.51	-1.08	0.04	0.01	0.06	0.05
	r	0.714	0.720	0.798	0.912	0.905	0.912

420

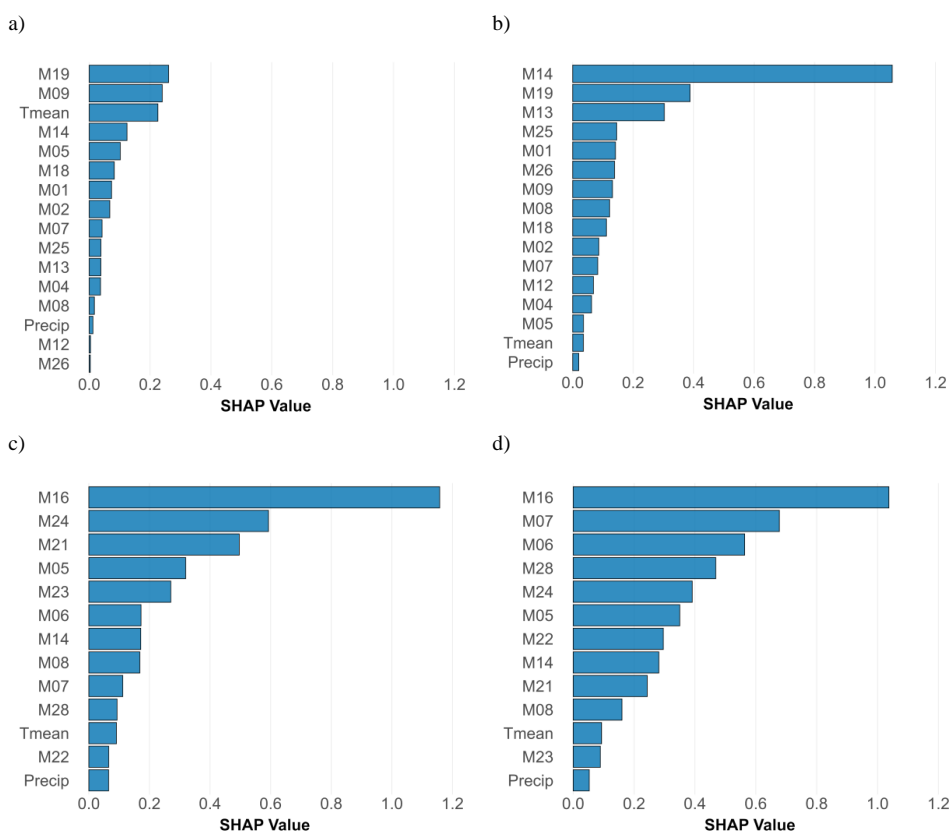
Table 4 provides statistics calculated based on the validation dataset. A comparison of the MMM and the best individual models reveals some key differences. At both crop sites (C1 and C2), the MMM exhibits lower RMSE values, indicating better accuracy. Its correlation coefficient (r) is higher than the best individual model at C1 (0.469 vs. 0.381) but slightly lower at C2 (0.556 vs. 0.599). At the grassland sites (G3 and G4), the MMM shows slightly higher RMSE values but a better correlation coefficient than the best individual models. Overall, the meta-models demonstrate a distinct improvement over both the MMM and the best individual models. Among the meta-models, RF, XGB and XGB+ are the top performers. At the crop sites, RF generally outperforms XGB+, particularly in the correlation coefficient at C1 (0.777 vs. 0.707). Notably at this site, the MLR meta-model also yielded better metric results than XGB. At the grassland sites, the performance differences between RF, XGB, and XGB+ were minimal, with all three models demonstrating strong performance across all metrics. Overall, explained variance typically

425

430



increases by ~19-38.5% for the best-performing metamodel compared to MMM (at C4 this was the lowest, and largest at C2). The meta-models typically show almost unbiased estimates, with the exception of C1.



435 **Figure 7:** The SHAP values of the XGB+ meta-model for RECO. Larger values mean stronger contribution to the  
 resulting RECO. a) C1 site; b) C2 site; c) G3 site; d) G4 site. Tmean stands for daily mean temperature, and Precip  
 is daily precipitation.

440 `SHAP values (Fig. 7) indicate that model outputs are the primary contributors of RECO predictions. Among  
 environmental predictors, mean temperature emerges as an important contributor, especially at site C1, where its  
 influence nearly matches top models M19 and M09. Precipitation influence on RECO is negligible across all sites.

### 3.3 NEE

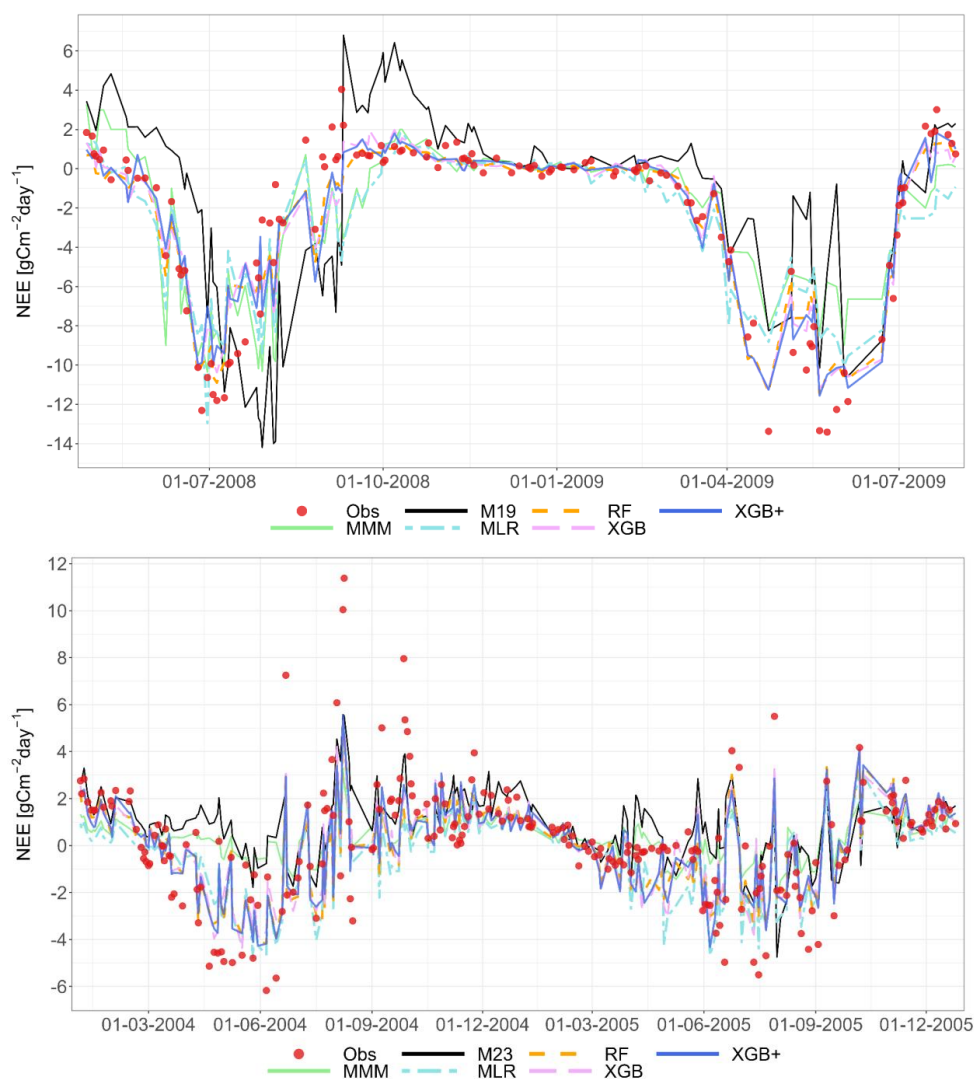
445

Fig. 8 illustrates the simulated and observed NEE for Grignon (C2) and Easter Bush (G4) over two years. Appendix  
 A contains the complete simulated dataset for all sites, and for all years.



At G4, none of the models fully capture the magnitude of the observed peaks (i.e. highest and lowest NEE). While the performance appears better during the first year at C2, where MLR closely tracks the observed C uptake, none of the models accurately represent the winter wheat (second year) minimum. Among the meta-models, RF, XGB and XGB+ provide the best representation of this minimum, while the best-performing individual model (M19) exhibits a large bias. The MMM shows greater consistency with other meta-models at C2 compared to its performance for RECO in Fig. 5. For both sites, M19 tends to diverge most from the other models during observation peaks, whereas RF, XGB and XGB+ are consistently following a similar trace, suggesting similar performance.

455



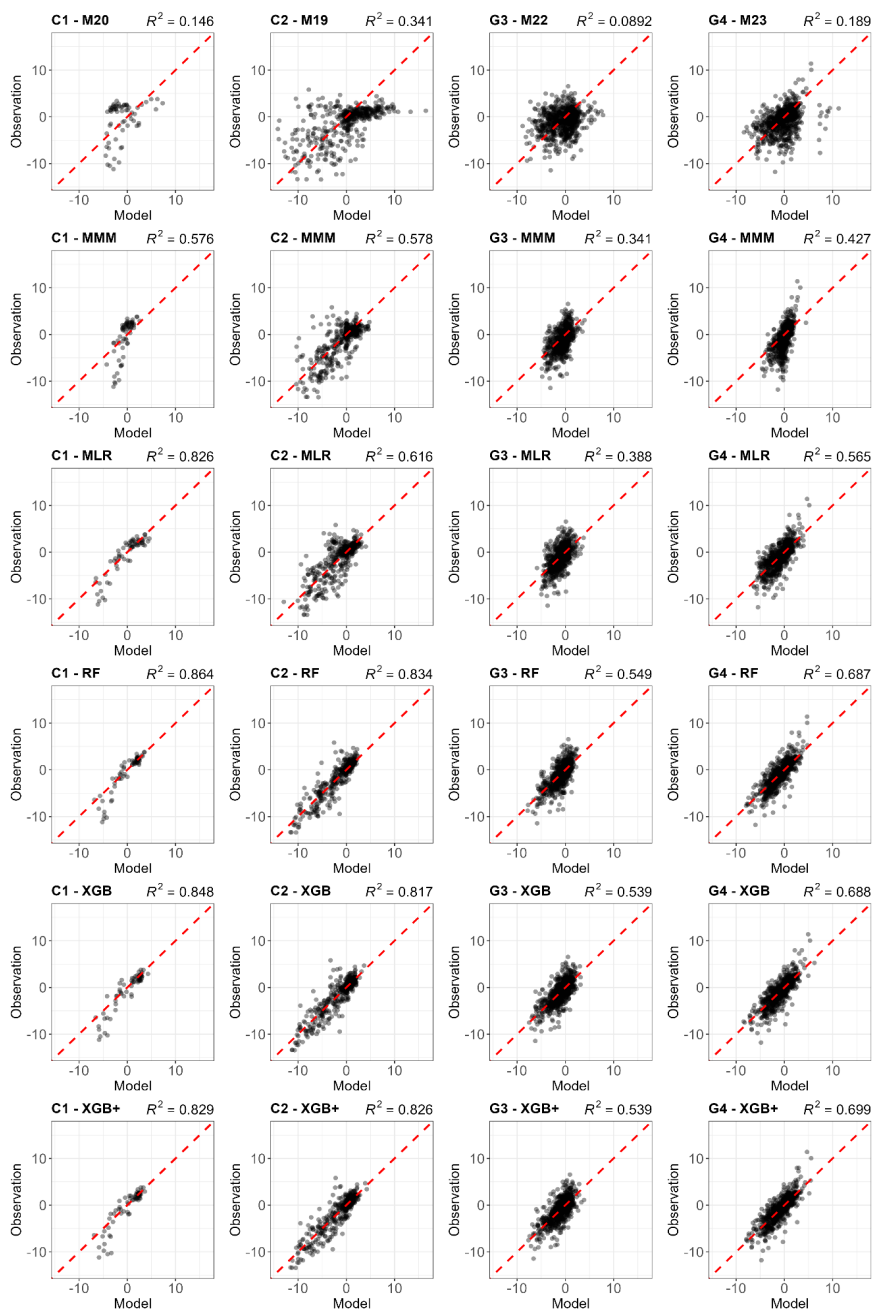
**Figure 8:** Performance of the multimodel median (MMM, green), the constructed meta-models and the best-performing individual models for simulating NEE. The top panel shows results for the Grignon cropland site (C2),

460



which includes maize and winter wheat. The bottom panel shows two years of data for the Easter Bush grassland site (G4). Observations are marked by red circles. The meta-models include Multiple Linear Regression (MLR, hatched light blue), Random Forest (RF, hatched orange), XGBoost (XGB, hatched purple), and XGBoost+ (XGB+, blue). The best-performing individual models (M19 at C2, M23 at G4) are shown in black. Dates are provided in the format of dd-mm-yyyy.

To provide a more objective assessment, Fig. 9 presents scatterplots for the entire dataset, with the  $R^2$  values shown. For all sites, the best individual models perform poorly, practically, capturing only a small portion of the observed variability. The MMM, however, aligns more closely with the meta-models (unlike the case with RECO in Fig. 5), and its  $R^2$  values are significantly higher than those of the best individual models, although they still indicate a relatively low correlation. The NEE tendencies at G3 were not well-captured by any model, as the best  $R^2$  value, produced by RF, is only 0.549. While this holds true for G4 as well, the correlations for the meta-models were 20-40% higher than at G3. At the crop sites (C1 and C2), the meta-models, particularly RF, XGB and XGB+, show a strong linear relationship with the observed data, explaining a large portion of the variance.



**Figure 9:** Comparison of the best individual model, the multi-model median and the constructed meta-models with observations for NEE. Each row represents a different model type, and columns correspond to the sites (from left to right: C1, C2, G3, and G4). The top row shows the best individual models with their identifiers (M20 at C1, M19 at C2, M22 at G3, M23 at G4). The remaining rows show the MMM, MLR, RF, XGB and XGB+. All units are in  $\text{g C m}^{-2} \text{ day}^{-1}$ . The red dashed line represents the 1:1 relationship.

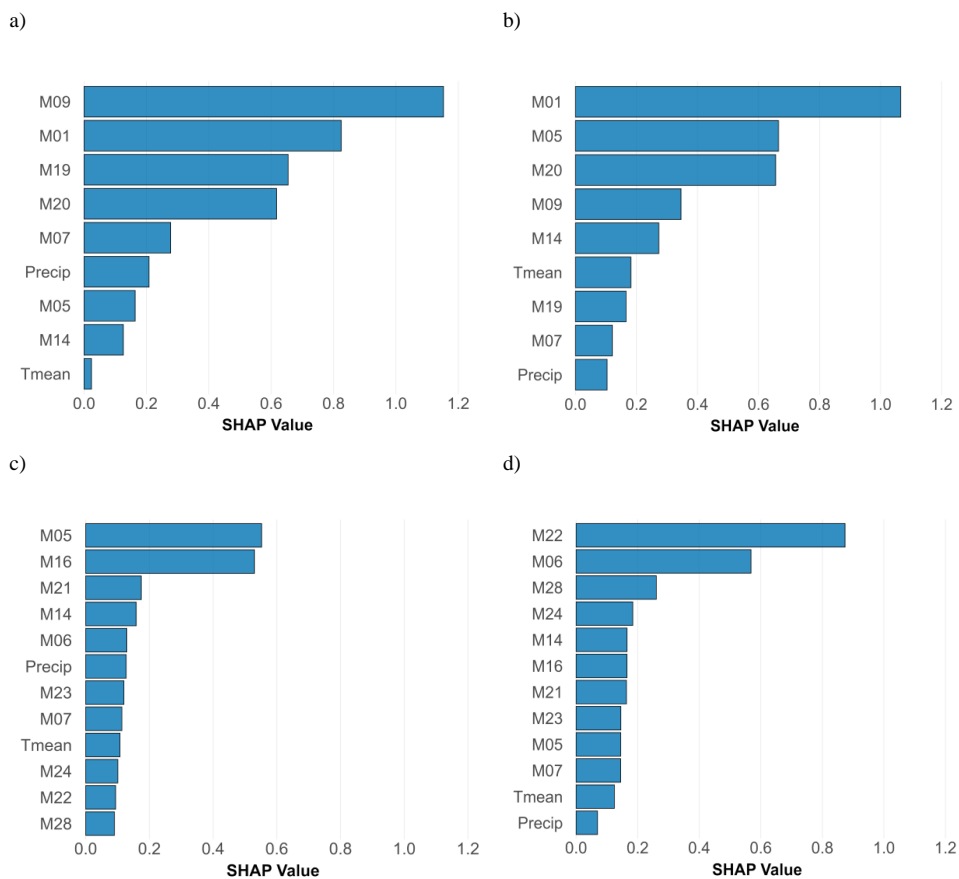
480



**Table 5:** Statistical evaluation of the best-performing individual model, the multi-model median (MMM) and the applied meta-models (MLR, RF, XGB and XGB+) for NEE. Three performance metrics are used: root mean square error (RMSE), bias and Pearson’s correlation coefficient (r). Only validation data were used for the calculation of the statistics. RMSE and BIAS are provided in  $\text{g C m}^{-2} \text{d}^{-1}$  units.

Site	Metric	best individual model	MMM	MLR	RF	XGB	XGB+
C1	RMSE	3.87	3.06	1.77	1.78	1.81	1.89
	BIAS	-0.48	0.1	0.38	0.46	0.66	0.52
	r	0.382	0.759	0.909	0.929	0.921	0.910
C2	RMSE	3.98	2.29	2.15	1.42	1.48	1.44
	BIAS	1.03	0.4	-0.1	-0.07	-0.02	0.02
	r	0.584	0.761	0.785	0.913	0.904	0.909
G3	RMSE	2.82	1.87	1.78	1.53	1.55	1.55
	BIAS	0.44	0.31	-0.03	0.01	0.02	-0.01
	r	0.299	0.584	0.623	0.741	0.734	0.734
G4	RMSE	2.61	2.06	1.64	1.39	1.39	1.36
	BIAS	0.48	0.67	-0.03	-0.01	0.03	0.02
	r	0.435	0.653	0.751	0.829	0.829	0.836

490 Table 5 shows statistics for all the models using three metrics based on the validation subset. Across all the four sites, the RF, XGB and XGB+ values differ little from each other in all three metrics, resulting in a similar performance, which is also visible in Fig. 8. The MLR meta-model is on par with these three meta-models, except at C2, where its performance is worse across all three metrics, making it more comparable to the MMM. The best-performing individual models have poor statistics compared to all other models. At the crop sites, they produce the worst values, with RMSE values that are, on average, 50% higher and a r value that is 36% lower than that of the MMM. At the grassland sites, the performance gap is slightly smaller, with the MMM’s RMSE values being, on average, 38% better and its correlation coefficient being 41% higher than that of the best individual models. Overall, the explained variance generally increased by ~15-18% for the best-performing meta-model compared to MMM, with the most pronounced improvement observed at G4. Bias shows similarity to those presented for GPP and RECO, meaning that with the exception of C1 the meta-models provide results that are almost bias-free.



505 **Figure 10:** The SHAP values of the XGB+ meta-model for NEE. Larger values mean stronger contribution to the resulting NEE. a) C1 site; b) C2 site; c) G3 site; d) G4 site. Tmean stands for daily mean temperature, and Precip is daily precipitation.

510 SHAP analysis presented in Fig. 10 reveals that a few key model outputs primarily determine NEE predictions. Meteorological variables, particularly temperature and precipitation, contribute modestly but importantly, with a stronger effect noted at grassland sites G3 and G4.



## 515 4 Discussion

### 4.1 Synthesis of the meta-model approaches

Across all sites and flux components, the stacking-based meta-models (MLR, RF, XGB, XGB+) consistently outperformed both the best individual process-based models and the traditional MMM. By combining ensemble  
520 learning with key environmental covariates, XGB+ emerged as a good candidate as well, in some cases achieving lower errors and higher correlations than all other approaches. Thus, while the MMM retained its role as a robust benchmark, its predictive power was consistently surpassed by the stacking-based methods.

This aligns with findings from Shahhosseini et al. (2020, 2021), who demonstrated that ensemble learning approaches can surpass the performance of individual models for crop yield prediction, and with Zhang et al. (2022),  
525 who reported similar benefits of ensemble learning in winter wheat yield modelling. The added value of explicitly incorporating drivers such as temperature into the meta-modelling framework also supports the conclusions of Mathieu and Aires (2018), who emphasised that agro-climatic indices like temperature and precipitation can significantly enhance model accuracy. Our results go beyond these earlier studies by demonstrating that such improvements are really possible (even without incorporating meteorology) across multiple flux components (GPP,  
530 RECO, NEE) and agroecosystem types, including both croplands and grasslands.

### 4.2 Cross-comparison of model contributions and environmental drivers for XGB+

Across all three fluxes (GPP, RECO, NEE), XGB+ also heavily relies on individual model outputs, confirming that the meta-model is successfully leveraging model consensus to improve predictive performance. However, the  
535 consistent appearance of temperature - particularly at the C1 site for RECO and at grassland sites for NEE - highlights that certain temperature-driven processes are not fully captured by the base models and need to be explicitly considered. Precipitation, while generally having less impact, may still play a role in certain site-specific cases (e.g. GPP at C1), suggesting localised interactions between water availability and C fluxes. Importantly, the climate characteristics of the study sites provide important context for interpreting these results. These locations lie  
540 within mid-latitude temperate zones with well-defined seasonality, and some (e.g. Ottawa) experience high continentality, marked by large annual temperature ranges and pronounced growing season transitions. In such climates, temperature becomes the primary constraint on biological activity during both dormancy and active periods, while precipitation is often more evenly distributed or limiting only during episodic droughts (Baldocchi, 2008; Koster et al., 2004; Sun et al., 2025). This climatic backdrop helps explain why temperature consistently  
545 emerges as a key predictor across sites and flux types, while precipitation plays a more site-specific and secondary role. These patterns align with broader findings from temperate ecosystems, where thermal constraints dominate respiration and photosynthesis processes, particularly outside of water-limited systems.

The dominant role of specific models in GPP predictions, supplemented by temperature's influence (Fig. 4), suggests that while ensemble model outputs capture much of the variability, temperature-dependent processes remain only  
550 partially resolved by the base models. This aligns with global studies such as Zhu et al. (2016) and Bellocchi et al.



(2023), who identified temperature as a key driver of GPP dynamics, especially in temperate and high-latitude ecosystems. The minimal role of precipitation in GPP prediction indicates that water availability was not a limiting factor at the examined sites, consistent with ecosystem-specific findings reported by Reichstein et al. (2013), which showed temperature or radiation often dominate in temperate regions.

555 For RECO, the strong temperature effect (Fig. 7) corroborates well-established biological principles, such as the temperature sensitivity of respiration processes described by Lloyd and Taylor (1994). The negligible role of precipitation further suggests these sites are not subject to drought stress, echoing findings from Xu et al. (2025) who highlighted thermal thresholds as primary controls over ecosystem carbon fluxes.

560 NEE's XGB+ related sensitivity to meteorological variables at grassland sites (Fig. 10) confirms the known climate sensitivity of these systems (Baldocchi et al., 2018). The significant influence of temperature and soil water content on GPP and RECO in grasslands, as noted by Xia et al. (2024), explains the prominence of these variables in NEE predictions. These findings underscore the importance of climate drivers in modulating carbon fluxes beyond what model ensembles alone capture.

	C1	C2	G3	G4
GPP	1/8	1/8	2/12	1/12
RECO	11/16	5/16	12/13	7/13
NEE	4/9	7/9	11/12	8/12

565

**Figure 11:** Ranking of the top-performing individual models (selected based on RMSE) across each site and model output based on their contribution to the XGB+ model. The first number in each cell indicates the ranking among the other models according to the SHAP value, while the second number is the total number of models. Red shades are just visual representations of the first number relative to the second.

570

Fig. 11 provides deeper insight into the functioning of XGB+ by ranking the influence of the top-performing process-based models (using SHAP values) for each flux and site. Several distinct patterns emerge. For RECO predictions, the influence of high-performing models was minimal, particularly at site G3, where the best model ranked only 12<sup>th</sup> out of 13. This indicates that, for RECO, the ensemble relied less on the best individual models. In contrast, for GPP, the top-performing models had a substantial impact across all sites, ranking first at both C1 and C2. This suggests that XGB+ leveraged structurally diverse yet individually strong models to enhance predictive accuracy. NEE exhibited moderate to low contributions from top-performing models - especially at grassland sites pointing to a more balanced integration of individual models and ensemble diversity. These findings reinforce a key principle of ensemble learning: the models with the highest standalone accuracy are not always the most influential

580



within the ensemble. By contributing unique and non-redundant information, even structurally weaker models can significantly improve an ensemble's predictive ability. This finding, based on SHAP analysis, is consistent with the principle of ensemble diversity highlighted by Chergui and Kechadi (2022), which posits that combining a variety of models with different strengths and weaknesses leads to more robust and accurate predictions. However, our SHAP analysis adds a novel layer of interpretability by quantifying the relative influence of top-performing versus structurally diverse models within the ensemble, providing diagnostic insights that go beyond earlier ensemble studies (Shahhosseini et al., 2020).

Gains in performance were more pronounced for croplands than for grasslands, especially for GPP and NEE. This aligns with previous research (Bansal et al., 2024; Nand et al., 2025) which found that management-intensive crop systems, like maize and winter wheat, benefit more from dynamic environmental and management data due to their seasonal and management-driven variability. For instance, Nand et al. (2025) demonstrated that weighted multi-model averaging, specifically the Granger-Ramanathan B method for actual evapotranspiration (ETa), reduced RRMSE by 4–8.5% in croplands. The Granger-Ramanathan B method is a multi-model averaging approach that requires non-negative weights that sum to one, and it provides the closest match to measured values for daily ETa in maize simulations. Conversely, grasslands, characterised by more stable phenological cycles, show smaller improvements. Their flux variability is often more strongly influenced by biotic controls rather than environmental drivers alone (Reichstein et al., 2007; Stoy et al., 2013). This is because grasslands exhibit more stable, biologically mediated carbon fluxes, while the sharp seasonal transitions and management-driven dynamics of crop systems make them particularly responsive to external data.

Overall, the XGB+ meta-model achieved good predictive performance across diverse agroenvironmental contexts, while also providing enhanced interpretability by quantifying the contributions of individual models and environmental drivers. These dual benefits - improved performance and diagnostic insight - position XGB+ as a powerful framework for advancing predictive modelling and guiding process-model refinement.

#### 4.3 Implications and applications

The demonstrated improvements in predictive performance of stacking-based meta-models have direct implications for crop and grassland biogeochemical modelling. By consistently outperforming both the best-performing individual models and the MMM, these approaches provide a pathway for enhancing C-flux predictions across diverse agro-environmental contexts. This is particularly relevant for applications such as greenhouse gas inventorying, site-specific management planning, and scenario analysis under climate change, where reliability is critical.

Beyond performance gains, the interpretability of the meta-models offers valuable diagnostic insights for process modellers. The analysis of regression coefficients and feature contributions highlights the value of retaining structurally diverse models in ensembles, rather than limiting selection to the top-performing candidates. Moreover, identifying temperature as a dominant external contributor points to opportunities for refining process-based models, such as improving their representation of phenological and temperature response mechanisms. These insights can also inform calibration practices: by quantifying the relative influence of models and drivers, the meta-modelling framework can guide targeted recalibration and prioritisation of model improvement efforts. In this way, meta-model outputs can serve as both a predictive tool and a diagnostic instrument for iterative process model development.



Together, these findings have broader implications for how model ensembles should be constructed and applied.

620 Our results highlight that the traditional practice of equal-weight multi-model averaging (like with the MMM) may not be sufficient for delivering credible predictions. As noted in Section 4.4, the approach depends on long, high-quality time series for training, which may limit application in some contexts. Consistent with the arguments of Mathieu and Aires (2018), Eyring et al. (2019), Shahhosseini et al. (2020) and Chergui and Kechadi (2022), who emphasised that structurally dependent and unevenly performing models require differentiated treatment, our

625 findings demonstrate that stacking-based meta-models - combining adaptive weighting with key environmental drivers - consistently outperform both the best individual models and the MMM. This evidence calls for a paradigm shift in ensemble design: future model intercomparison and synthesis efforts should move toward performance-informed, diagnostic-driven weighting strategies that explicitly incorporate relevant covariates. Such approaches not only improve predictive accuracy but also generate actionable insights for refining underlying

630 process-based models, ultimately accelerating progress toward more reliable and process-rich representations of agroecosystem C dynamics. This calls for coordinated action by the modelling community to adopt adaptive, performance-informed ensemble frameworks in future intercomparison efforts. Furthermore, our framework could be applied as a post-processing step to archived ensemble outputs from major multi-model initiatives (discussed in Section 5), maximising the value of existing model intercomparison investments without requiring new simulations.

#### 635 **4.4 Limitations and future research**

While the proposed meta-modelling framework offers notable advantages, several limitations warrant acknowledgment. First, the approach depends on sufficiently long and consistent time series for training and evaluation, which were not uniformly available across all sites and flux components, limiting generalisability in data-scarce regions. This data dependence is highlighted in Section 4.3 as a key consideration for ensemble design.

640 Second, incorporating external covariates improves performance but increases complexity and introduces dependencies on auxiliary datasets that may not always be accessible or reliable. Our analysis also focused on a limited set of covariates (temperature and precipitation), leaving unexplored the potential benefits of integrating hydrological or soil-related variables.

The interpretability of regression coefficients also comes with caveats: while they provide insight into model influence within the ensemble, they do not directly reveal mechanistic underpinnings, and the relationships captured remain largely empirical. Consequently, care must be taken when using these findings to inform process-level changes.

645

Future research should expand the set of environmental covariates to include soil moisture, global radiation, management practices, and hydrological variables, which are likely to improve representation of key drivers in both

650 crops and grasslands. Testing the framework across a broader range of management systems, crop types, and climatic zones would also help assess scalability and robustness. Moreover, integrating advanced interpretability tools (e.g., causal inference methods) could move beyond purely statistical associations toward more mechanistically grounded insights. Finally, co-developing meta-model frameworks with process-based modellers could establish a feedback loop in which ensemble diagnostics directly inform iterative improvements to individual models, accelerating

655 convergence toward more reliable, process-rich representations.



## 5 Conclusions

International initiatives that foster collaborations between researchers working on agricultural and grassland models are creating new opportunities in the field of process-oriented modelling. By gathering outputs from several models with different representations of plant and soil processes and using standardised protocols, exploitation of the potential of the ensembles becomes possible. However, those multi-model ensemble techniques are still at their infancy, as typically simple multi-model means or medians are constructed as robust estimations that typically overperform the individual models. Some studies attempted to use more sophisticated methods like skill-based model selection and even machine learning, but the potential of the multi-model frameworks is still being explored. In this study, building on a previous multi-model exercise performed under the umbrella of international initiatives, new combinations of models are tested that we call here as meta-models.

The introduced meta-models significantly improved the accuracy of C-flux estimates in crop and grassland ecosystems compared to individual process-based models and traditional multi-model medians. By integrating structurally diverse models and incorporating key environmental variables - particularly temperature -, these meta-models deliver not only more reliable predictions but also diagnostic insights into the relative contributions of models and environmental drivers. This underscores the importance of maintaining model diversity in ensembles and highlights opportunities to refine process-based models, especially regarding temperature responses and phenological processes.

Performance gains were more pronounced in crop systems than in grasslands, likely reflecting the stronger influence of management and pronounced seasonal dynamics. Nevertheless, the approach relies on long, high-quality datasets and auxiliary covariates, and its empirical nature limits direct mechanistic interpretation. These limitations were detailed in Section 4.4.

Importantly, this framework opens the door to re-analyzing outputs from major multi-model initiatives like AgMIP (Agricultural Model Intercomparison and Improvement Project; <https://agmip.org>) and MACSUR (Modelling European Agriculture with Climate Change for Food Security; <https://www.facejpi.net/en/facejpi/actions/core-theme-1/knowledge-hub-macsur.htm>). Rather than requiring new simulations, archived ensemble datasets from these projects could be post-processed using stacking meta-models to enhance predictive skill and extract new diagnostic insights - maximising the value of past investments in large-scale model intercomparison. We strongly encourage international modelling communities to pilot such stacking-based re-analyses, which offer a low-cost, high-impact opportunity to unlock new insights and improve ensemble predictions.

While promising, this approach requires long, high-quality datasets and auxiliary inputs and remains empirical in nature, calling for caution when inferring mechanistic causation. Future research should expand the set of environmental drivers (e.g., soil and hydrological variables), test scalability across broader agroecosystems, and apply advanced interpretability tools. Collaborative development with process-based modellers could translate these statistical gains into mechanistic improvements, ultimately leading to a new generation of hybrid ensemble frameworks for agricultural and grassland biogeochemistry. We encourage the international modelling communities to pilot such stacking-based re-analyses, leveraging their rich archives to unlock new insights and improve ensemble predictions without additional simulation costs.



### Code and Data availability

695 The exact versions of the R scripts used to produce the results presented in this paper are available from the GitHub repository: [https://github.com/hollorol/metamodeling\\_of\\_c](https://github.com/hollorol/metamodeling_of_c) under the GPL-3 licence. The input data used to produce the results presented in this paper is archived on the Harvard Dataverse repository under doi:10.7910/DVN/5TO4HE.

### Author contribution

700 Roland Hollós: Methodology, Conceptualization, Visualization, Writing (original draft preparation)  
Nándor Zrinyi: Software, Visualization, Writing (original draft preparation)  
Zoltán Barcza: Methodology, Writing (original draft preparation)  
Gianni Bellocchi: Methodology, Writing (original draft preparation), Supervision  
Renáta Sándor: Data curation, Validation  
705 János Ruff: Conceptualization, Formal analysis  
Nándor Fodor: Funding acquisition, Resources, Project administration

### Acknowledgements

The present article was published under the auspices of the MACSUR (Modelling European Agriculture with Climate Change for Food Security) Science-Policy Knowledge Forum (MACSUR SciPol Pilot, June 2021-  
710 December 2022, and MACSUR SciPolNet, May 2024-April 2026), with the support of the INRAE metaprogramme “Climate change in agriculture and forests: Adaptation and mitigation” (CLIMAE) and INRAE’s Public Policy-Support Directorate (DAPP). It falls within the thematic area of the French government IDEX-ISITE initiative (reference: 16-IDEX-0001; project CAP 20-25). This work has been partly implemented by the National Multidisciplinary Laboratory for Climate Change (RRF-2.3.1-21-2022-00014) project within the framework of  
715 Hungary’s National Recovery and Resilience Plan supported by the Recovery and Resilience Facility of the European Union. This project was also supported by the FK 131813 project, implemented with support provided by the National Research, Development and Innovation Fund of Hungary, financed under the FK\_19 funding scheme. RH and ZB was supported by the “Advanced methods of greenhouse gases emission reduction and sequestration in agriculture and forest landscape for climate change mitigation” (CZ.02.01.01/00/22\_008/0004635) project. RS, RH  
720 and GB received mobility funding from the French-Hungarian bilateral partnership through the BALATON (N° 44703TF)/TÉT (2019-2.1.11-TÉT-2019-00031) programme. Support was also provided by the TKP2021-NKTA-06 project that has been implemented with the support provided by the Ministry of Innovation and Technology of Hungary from the National Research, Development and Innovation Fund, financed under the [TKP2021-NKTA] funding scheme.

725



## References

- 730 Anuga, S.W., Chirinda, N., Nukpezah, D., Ahenkan, A., Andrieu, N., Gordon, C.: Towards low carbon agriculture: Systematic-narratives of climate-smart agriculture mitigation potential in Africa. *Current Research in Environmental Sustainability* 2, 100015. <https://doi.org/10.1016/j.crsust.2020.100015>, 2020.
- 735 Asseng, S., Ewert, F., Rosenzweig, C., Jones, J.W., Hatfield, J.L., Ruane, A.C., Boote, K.J., Thorburn, P.J., Rötter, R.P., Cammarano, D., Brisson, N., Basso, B., Martre, P., Aggarwal, P.K., Angulo, C., Bertuzzi, P., Biernath, C., Challinor, A.J., Doltra, J., Gayler, S., Goldberg, R., Grant, R., Heng, L., Hooker, J., Hunt, L.A., Ingwersen, J., Izaurralde, R.C., Kersebaum, K.C., Müller, C., Naresh Kumar, S., Nendel, C., O'Leary, G., Olesen, J.E., Osborne, T.M., Palosuo, T., Priesack, E., Ripoche, D., Semenov, M.A., Shcherbak, I., Steduto, P., Stöckle, C., Stratonovitch, P., Streck, T., Supit, I., Tao, F., Travasso, M., Waha, K., Wallach, D., White, J.W., Williams, J.R., Wolf, J.: Uncertainty in simulating wheat yields under climate change. *Nature Climate Change* 3, 827-832. <https://doi.org/10.1038/nclimate1916>, 2013.
- 740 Bai, Y., Zhang, S., Bhattarai, N., Mallick, K., Liu, Q., Tang, L., Im, J., Guo, L., Zhang, J. : On the use of machine learning based ensemble approaches to improve evapotranspiration estimates from croplands across a wide environmental gradient. *Agricultural and Forest Meteorology* 298-299, 108308. <https://doi.org/10.1016/j.agrformet.2020.108308>, 2021.
- 745 Baldocchi, D.: 'Breathing' of the terrestrial biosphere: Lessons learned from a global network of carbon dioxide flux measurement systems. *Australian Journal of Botany* 56, 1-26. <https://doi.org/10.1071/BT07151>, 2008.
- 750 Baldocchi, D., Chu, H., Reichstein, M.: Intercomparison of ten eddy covariance flux partitioning methods with a global dataset of CO<sub>2</sub> fluxes. *Agricultural and Forest Meteorology* 256-257, 223-233. <https://doi.org/10.1016/j.agrformet.2017.05.015>, 2018.
- 755 Bansal, Y., Lillis, D., Kechadi, M.T.: A neural meta model for predicting winter wheat crop yield. *Machine Learning* 113: 3771-3788. <https://doi.org/10.1007/s10994-023-06455-1>, 2024.
- 760 Bassu, S., Brisson, N., Durand, J.L., Boote, K.J., Lizaso, J., Jones, J.W., Rosenzweig, C., Adam, M., Basso, B., Baron, C., Basso, B., Biernath, C., Boogaard, H., Conijn, S., Corbeels, M., Deryng, D., De Sanctis, G., Gayler, S., Grassini, P., Hatfield, J., Hoek, S., Izaurralde, C., Jongschaap, R., Kemanian, A.R., Kersebaum, K.C., Kim, S.-H., Kumar, N.S., Makowski, D., Müller, C., Nendel, C., Priesack, E., Pravia, M.V., Sau, F., Shcherbak, I., Tao, F., Teixeira, E., Timlin, D., Waha, K.: How do various maize crop models vary in their responses to climate change factors? *Global Change Biology* 20, 2301-2320. <https://doi.org/10.1111/gcb.12520>, 2014.
- 765 Beer, C., Reichstein, M., Tomelleri, E., Ciais, P., Jung, M., Carvalhais, N., Rödenbeck, C., Arain, M. A., Baldocchi, D., Bonan, G. B., Bondeau, A., Cescatti, A., Lasslop, G., Lindroth, A., Lomas, M., Luysaert, S., Margolis, H., Oleson, K.W., Rouspard, O., Veenendaal, E., Viovy, N., Williams, C., Woodward, F.I., Papale, D.: Terrestrial gross



- carbon dioxide uptake: Global distribution and covariation with climate. *Science* 329, 834–838. <https://doi.org/10.1126/science.1184984>, 2010.
- 770 Bellocchi, G.: MACSUR SciPol Policy brief 5: Assessing Emissions and Mitigation Practices in Agriculture, towards an effective use of models for policy. Zenodo. <https://doi.org/10.5281/zenodo.8038881>, 2023.
- 775 Bellocchi, G., Barcza, Z., Hollós, R., Acutis, M., Bottyán, E., Doro, L., Hidy, D., Lellei-Kovács, E., Ma, S., Minet, J., Pacskó, V., Perego, A., Ruget, F., Seddaiu, G., Wu, L., Sándor, R.: Sensitivity of simulated soil water content, evapotranspiration, gross primary production and biomass to climate change factors in Euro-Mediterranean grasslands. *Agricultural and Forest Meteorology* 343, 109778. <https://doi.org/10.1016/j.agrformet.2023.109778>, 2023.
- 780 Bellocchi, G., Rivington, M., Donatelli, M., Matthews, K.: Validation of biophysical models: issues and methodologies. A review. *Agronomy for Sustainable Development* 30, 109–113. <https://doi.org/10.1051/agro/2009001>, 2010.
- 785 Bilotto, F., Harrison, M.T., Migliorati, M.D.A., Christie, K.M., Rowlings, D.W., Grace, P.R., Smith, A.P., Rawnsley, R.P., Thorburn, P.J., Eckard, R.J.: Can seasonal soil N mineralisation trends be leveraged to enhance pasture growth? *Science of the Total Environment* 772:145031. <https://doi.org/10.1016/j.scitotenv.2021.145031>, 2021.
- Breiman, L.: Stacked regressions. *Machine Learning* 24, 123–140. <https://doi.org/10.1007/BF00058655>, 1996.
- 790 Breiman, L.: Stacked regressions. *Machine Learning* 24, 49–64. <https://doi.org/10.1007/BF00117832>, 2001a.
- Breiman, L.: Random forests. *Machine Learning* 45, 5–32. <https://doi.org/10.1023/A:1010933404324>, 2001b.
- 795 Brill, L., Bechini, L., Bindi, M., Carozzi, M., Cavalli, D., Conant, R., Dorich, C.D., Doro, L., Ehrhardt, F., Farina, R., Ferrise, R., Fitton, N., Francaviglia, R., Grace, P., Iocola, I., Klumpp, K., Léonard, J., Martin, R., Massad, R.S., Recous, S., Seddaiu, G., Sharp, J., Smith, P., Smith, W.N., Soussana, J-F., Bellocchi, G.: Review and analysis of strengths and weaknesses of agro-ecosystem models for simulating C and N fluxes. *Science of the Total Environ.* 598, 445-470. <https://doi.org/10.1016/j.scitotenv.2017.03.208>, 2017.
- 800 Calanca, P., Deléglise, C., Martin, R., Carrère, P., Mosimann, E.: Testing the ability of a simple grassland model to simulate the seasonal effects of drought on herbage growth. *Field Crops Research* 187, 12-23. <https://doi.org/10.1016/j.fcr.2015.12.008>, 2016.
- 805 Challinor, A.J., Smith, M.S., Thornton, P.: Use of agro-climate ensembles for quantifying uncertainty and informing adaptation. *Agricultural and Forest Meteorology* 170, 2-7. <https://doi.org/10.1016/j.agrformet.2012.09.007>, 2013.



- Chandel, S., Kleber, M., Jahn, R., Vogel, C.: Soil science-informed machine learning. *Geoderma* 452, 117094. <https://doi.org/10.1016/j.geoderma.2024.117094>, 2024.
- 810 Chen, T., Guestrin, C.: XGBoost: A scalable tree boosting system. In: *KDD '16: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 785-794. <https://doi.org/10.1145/2939672.2939785>, 2016.
- Chergui, N., Kechadi, M.T.: Data analytics for crop management: a big data view. *Journal of Big Data* 9, 1-37. <https://doi.org/10.1186/s40537-022-00668-2>, 2022.
- 815 Dieterich, T.G.: Ensemble methods in machine learning. In: *Multiple classifier systems. MCS 2000. Lecture Notes in Computer Science*, vol 1857. Springer, Berlin, Heidelberg. [https://doi.org/10.1007/3-540-45014-9\\_1](https://doi.org/10.1007/3-540-45014-9_1), 2000.
- 820 Dormann, C.F., Elith, J., Bacher, S., Buchmann, C., Carl, G., Carré, G., Marquéz, J.R.G., Gruber, B., Lafourcade, B., Leitão, P.J., Münkemüller, T., McClean, C., Osborne, P.E., Reineking, B., Schröder, B., Skidmore, A.K., Zurell, D., Lautenbach, S.: Collinearity: A review of methods to deal with it and a simulation study evaluating their performance. *Ecography* 36, 27–46. <https://doi.org/10.1111/j.1600-0587.2012.07348.x>, 2013.
- 825 Ehrhardt, F., Soussana, J.-F., Bellocchi, G., Grace, P., McAuliffe, R., Recous, S., Sándor, R., Smith, P., Snow, V., Migliorati, M.D.A., Basso, B., Bhatia, A., Brilli, L., Doltra, J., Dorich, C.D., Doro, L., Fitton, N., Giacomini, S.J., Grant, B., Harrison, M.T., Jones, S.K., Kirschbaum, M.U.F., Klumpp, K., Laville, P., Léonard, J., Liebig, M., Lieffering, M., Martin, R., Massad, R.S., Meier, E., Merbold, L., Moore, A.D., Myrgiotis, V., Newton, P., Pattey, E., Rolinski, S., Sharp, J., Smith, W.N., Wu, L., Zhang, Q.: Assessing uncertainties in crop and pasture ensemble model simulations of productivity and N2O emissions. *Global Change Biology* 24, e603-e616. <https://doi.org/10.1111/gcb.13965>, 2018.
- 830 Eyring, V., Cox, P.M., Flato, G.M., Friedlingstein, P., Hall, A., Hawkins, E., Hewitt, H.T., Joshi, M., Klein, S.A., Knutti, R., Meehl, G.A., O'Neill, B.C., Piani, C., Raper, S.C.B., Riahi, K., Roeckner, E., Sanderson, B.M., Wenzel, S.: Taking climate model evaluation to the next level. *Nature Climate Change* 9, 102–110. <https://doi.org/10.1038/s41558-018-0355-y>, 2019.
- 835 Farina, R., Sándor, R., Abdalla, M., Álvaro-Fuentes, J., Bechini, L., Bolinder, M.A., Brilli, L., Chenu, C., Clivot, H., De Antoni Migliorati, M., Di Bene, C., Dorich, C.D., Ehrhardt, F., Ferchaud, F., Fitton, N., Francaviglia, R., Franko, U., Giltrap, D.L., Grant, B.B., Guenet, B., Harrison, M.T., Kirschbaum, M.U.F., Kuka, K., Kulmala, L., Liski, J., McGrath, M.J., Meier, E., Menichetti, L., Moyano, F., Nendel, C., Recous, S., Reibold, N., Shepherd, A., Smith, W.N., Smith, P., Soussana, J.-F., Stella, T., Taghizadeh-Toosi, A., Tsutsikh, E., Bellocchi, G.: Ensemble modelling, uncertainty and robust predictions of organic carbon in long-term bare-fallow soils. *Global Change Biology* 27, 904-928. <https://doi.org/10.1111/gcb.15441>, 2021.



- 845 Faticchi, S., Pappas, C., Zscheischler, J., Leuzinger, S.: Modelling carbon sources and sinks in terrestrial vegetation. *New Phytologist* 221, 652-668. <https://doi.org/10.1111/nph.15451>, 2019.
- Gascuel-Oudou, C., Lescourret, F., Dedieu, B., Detang-Dessendre, C., Faverdin, P., Hazard, L., Litrico-Chiarelli, I., Petit, S., Roques, L., Reboud, X., Tixier-Boichard, M., de Vries, H., Caquet, T.: A research agenda for scaling up agroecology in European countries. *Agronomy for Sustainable Development* 42, 53. <https://doi.org/10.1007/s13593-022-00786-4>, 2022.
- 850
- Granger, C.W.J., Ramanathan R.: Improved methods of combining forecasts. *Journal of Forecasting* 3:197-204. <https://doi.org/10.1002/for.3980030207>, 1984.
- 855
- Hagedorn, R., Doblas-Reyes, F.J., Palmer, T.N.: The rationale behind the success of multi-model ensembles in seasonal forecasting – I. Basic concept. *Tellus A: Dynamic Meteorology and Oceanography* 57, 219-233. <https://doi.org/10.1111/j.1600-0870.2005.00103.x>, 2005.
- 860
- Hansen, L.K., Salamon, P.: Neural network ensembles. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 12, 993-1001. <https://doi.org/10.1109/34.58871>, 1990.
- Harrison, M.T., Evans, J.R., Moore, A.D.: Using a mathematical framework to examine physiological changes in winter wheat after livestock grazing: 1. Model derivation and coefficient calibration. *Field Crops Research* 136, 116-126. <https://doi.org/10.1016/j.fcr.2012.06.015>, 2012.
- 865
- Hengl, T., de Jesus, J.M., Heuvelink, G.B., Gonzalez, M.R., Kilibarda, M., Blagotić, A., Shangguan, W., Wright, M.N., Geng, X., Bauer-Marschallinger, B., Guevara, M.A., Vargas, R., MacMillan, R.A., Batjes, N.H., Leenaars, J.G.B., Ribeiro, E., Wheeler, I., Mantel, S., Kempen, B.: SoilGrids250m: Global gridded soil information based on machine learning. *PLoS ONE* 12, e0169748. <https://doi.org/10.1371/journal.pone.0169748>, 2018.
- 870
- Hijmans, R.J., Cameron, S.E., Parra, J.L., Jones, P.G., Jarvis, A.: Very high-resolution interpolated climate surfaces for global land areas. *International Journal of Climatology* 25, 1965-1978. <https://doi.org/10.1002/joc.1276>, 2005.
- 875
- Hou, J., Hou, B.: Farmers' adoption of low-carbon agriculture in China: An extended theory of the planned behavior model. *Sustainability* 11, 1399. <https://doi.org/10.3390/su11051399>, 2019.
- Jackson, R. B., Lajtha, K., Crow, S. E., Hugelius, G., Kramer, M. G., Rhodes, C.: The ecology of soil carbon: pools, vulnerabilities, and biotic and abiotic controls. *Annual Review of Ecology, Evolution, and Systematics* 48, 419-445. <https://doi.org/10.1146/annurev-ecolsys-112414-054234>
- 880
- Janes-Bassett, V., Davies, J., Rowe, Ed C., Tipping, E., 2020. Simulating long-term carbon nitrogen and phosphorus biogeochemical cycling in agricultural environments. *Science of the Total Environment* 714, 136599. <https://doi.org/10.1016/j.scitotenv.2020.136599>, 2017.



885

Jones, J.W., Antle, J.M., Basso, B.O., Boote, K.J., Conant, R.T., Foster, I., Godfray, H.C.J., Herrero, M., Howitt, R.E., Janssen, S., Keating, B.A., Muñoz-Carpena, R., Porter, C.H., Rosenzweig, C., Wheeler, T.R.: Brief history of agricultural systems modelling. *Agricultural Systems* 155, 240-254. <https://doi.org/10.1016/j.agsy.2016.05.014>, 2017.

890

Jung, M., Vetter, M., Herold, M., Churkina, G., Reichstein, M., Zaehle, S., Ciais, P., Viovy, N., Bondeau, A., Chen, Y., Trusilova, K., Feser, F., Heimann, M.: Uncertainties of modelling gross primary productivity over Europe: A systematic study on the effects of using different drivers and terrestrial biosphere models. *Global Biogeochemical Cycles* 21, GB4021. <https://doi.org/10.1029/2006GB002915>, 2007.

895

Keskin, H., Grunwald, S., Basso, B.: Machine learning advances to predict crop yield for agricultural systems. *Soil Science Society of America Journal* 83, 1521-1531. <https://doi.org/10.2136/sssaj2019.06.0203>, 2019.

900

Knutti, R., Baumberger, C., Hirsch Hadorn, G.: Uncertainty quantification using multiple models - prospects and challenges. In: Beisbart C., Saam N.J. (eds.) *Computer simulation validation: fundamental concepts, methodological frameworks, and philosophical perspectives*. Springer: Cham, pp. 835–855, 2019.

905

Kobayashi, K., Salam, M.U.: Comparing simulated and measured values using mean squared deviation and its components. *Agronomy Journal* 92, 345-352. <https://doi.org/10.2134/agronj2000.922345x>, 2000.

910

Kollas, C., Kersebaum, K.C., Nendel, C., Manevski, K., Müller, C., Palosuo, T., Armas-Herrera, C.M., Beaudoin, N., Bindi, M., Charfeddine, M., Conradt, T., Constantin, J., Eitzinger, J., Ewert, F., Ferrise, R., Gaiser, T., Garcia de Cortazar-Atauri, I., Giglio, L., Hlavinka, P., Hoffmann, H., Hoffmann, M.P., Launay, M., Manderscheid, R., Mary, B., Mirschel, W., Moriondo, M., Olesen, J.E. Öztürk, I., Pacholski, A., Ripoché-Wachter, D., Roggero, P.P., Roncossek, S., Rötter, R.P., Ruget, F., Sharif, B., Trnka, M., Ventrella, D., Waha, K., Wegehenkel, M., Weigel, H.-J., Wu, L.: Crop rotation modelling - A European model intercomparison. *European Journal of Agronomy* 70, 98–111. <https://doi.org/10.1016/j.eja.2015.06.007>, 2015.

915

Koster, R.D., Dirmeyer, P.A., Guo, Z., Bonan, G., Chan, E., Cox, P., Gordon, C.T., Kanae, S., Kowalczyk, E., Lawrence, D., Liu, P., Lu, C.-H., Malyshev, S., McAvaney, B., Mitchell, K., Mocko, D., Oki, T., Oleson, K., Pitman, A., Sud, Y. C., Taylor, C. M., Verseghy, D., Vasic, R., Xue, Y., Yamada, T.: Regions of strong coupling between soil moisture and precipitation. *Science* 305, 1138–1140. <https://doi.org/10.1126/science.1100217>, 2004.

920

Kutner, M.H., Nachtsheim, C.J., Neter, J. and Li, W.: *Applied linear statistical models*. 5<sup>th</sup> Edition, McGraw-Hill, Irwin, New York, 2005.

Lambin, E.F., Meyfroidt, P.: Global land use change, economic globalization, and the looming land scarcity. *Proceedings of the National Academy of Sciences* 108, 3465-3472. <https://doi.org/10.1073/pnas.1100480108>, 2011.



- 925 Lembaid, I., Moussadek, R., Mrabet, R., Bouhaouss, A.: Modeling soil organic carbon changes under alternative climatic scenarios and soil properties using DNDC model at a semi-arid Mediterranean environment. *Climate* 10, 23. <https://doi.org/10.3390/cli10020023>, 2022.
- Lembaid, I., Moussadek, R., Mrabet, R., Doauik, A., Bouhaouss, A.: Modeling the effects of farming management practices on soil organic carbon stock under two tillage practices in a semi-arid region, Morocco. *Heliyon* 7, e05889. <https://doi.org/10.1016/j.heliyon.2020.e05889>, 2021.
- 930 Li, T., Cui, L., Kuhnert, M., McLaren, T.I., Pandey, R., Liu, H., Wang, W., Xu, Z., Xia, A., Dalal, R.C., Dang, Y.P.: A comprehensive review of soil organic carbon estimates: Integrating remote sensing and machine learning technologies. *Journal of Soils and Sediments* 24, 3556-3571. <https://doi.org/10.1007/s11368-024-03913-8>, 2015.
- 935 Li, C., Farahbakhshazad, N., Jaynes, D.B., Dinnes, D.L., Salas, W., McLaughlin, D.: Modeling nitrate leaching with a biogeochemical model modified based on observations in a row-crop field in Iowa. *Ecological Modelling* 196, 116-130. <https://doi.org/10.1016/j.ecolmodel.2006.02.007>, 2006.
- 940 Li, T., Hasegawa, T., Yin, X., Zhu, Y., Boote, K., Adam, M., Bregalgio, S., Buis, S., Confalonieri, R., Fumoto T., Gaydon, D., Marcaida III, M., Nakagawa, H., Oriol, P., Ruane, A.C., Ruget, F., Balwinder-Singh, B., Singh, U., Tang, L., Tao, F., Wilkens, P., Yoshida, H., Zhang, Z., Bouman, B.: Uncertainties in predicting rice yield by current crop models under a wide range of climatic conditions. *Global Change Biology* 21, 1328–1341. <https://doi.org/10.1111/gcb.12758>, 2015.
- 945 Liaw, A., Wiener, M.: Classification and regression by randomForest. *R News* 2, 18-22. <http://CRAN.R-project.org/doc/Rnews>, 2002.
- 950 Lobell, D.B., Schlenker, W., Costa-Roberts, J.: Climate trends and global crop production since 1980. *Science* 333, 616-620. <https://doi.org/10.1126/science.1204531>, 2011.
- Lloyd, J., Taylor, J.A.: On the temperature dependence of soil respiration. *Functional Ecology* 8, 315–323. <https://doi.org/10.2307/2389824>, 1994.
- 955 Lundberg, S.M., Erion, G., Lee, S.-I.: Consistent individualized feature attribution for tree ensembles. arXiv:1802.03888. <https://doi.org/10.48550/arXiv.1802.03888>, 2020.
- Lundberg, S.M., Lee, S.-I., 2017. A unified approach to interpreting model predictions. Proceedings of the 31<sup>st</sup> International Conference on Neural Information Processing Systems, Long Beach, 4-9 December, pp. 4766-4777.
- 960 Luo, Y., Weng, E., Wu, X., Gao, C., Zhou, X., Zhang, L.: Parameter identifiability, constraint, and equifinality in data assimilation with ecosystem models. *Ecological Applications* 19, 571-574. <https://doi.org/10.1890/08-0561.1>, 2009.



- 965 Mangani, R., Tesfamariam, E., Engelbrecht, C.J., Bellocchi, G., Hassen, A., Mangani, T.: Potential impacts of extreme weather events in main maize (*Zea mays* L.) producing areas of South Africa under rainfed conditions. *Regional Environmental Change* 19, 1441-1452. <https://doi.org/10.1007/s10113-019-01486-8>, 2019.
- Martre, P., Wallach, D., Asseng, S., Ewert, F., Jones, J.W., Rotter, R.P., Boote, K.J., Ruane, A.C., Thorburn, P.J.,  
970 Cammarano, D., Hatfield, J.L., Rosenzweig, C., Aggarwal, P.K., Angulo, C., Basso, B., Bertuzzi, P., Biernath, C.,  
Brisson, N., Challinor, A.J., Doltra, J., Gayler, S., Goldberg, R., Grant, R.F., Heng, L., Hooker, J., Hunt, L.A.,  
Ingwersen, J., Izaurralde, R.C., Kersebaum, K.C., Müller, C., Kumar, S.N., Nendel, C., O'leary, G., Olesen, J.E.,  
Osborne, T.M., Palosuo, T., Priesack, E., Ripoche, D., Semenov, M.A., Shcherback, I., Steduto, P., Stöckle, C.O.,  
Stratonovitch, P., Streck, T., Supit, I., Tao, F., Travasso, M., Waha, K., White, J.W., Wolf, J., 2015. Multimodel  
975 ensembles of wheat growth: many models are better than one. *Global Change Biology* 21, 911-925.  
Mathieu, J.A., Aires, F.: Assessment of the agro-climatic indices to improve crop yield forecasting. *Agricultural and  
Forest Meteorology* 253-254:15-30. <https://doi.org/10.1016/j.agrformet.2018.01.031>, 2018.
- Nand, V., Qi, Z., Ma, L., Helmers, M.J., Madramootoo, C.A., Smith, W.N., Zhang, T., Weber, T.K.D., Pattey, E.,  
980 Li, Z., Wang, J., Jin, V.L., Jiang, Q., Tenuta, M., Trout, T.J., Cheng, H., Harmel, R.D., Kimball, B.A., Thorp, K.R.,  
Boote, K.J., Stockle, C., Suyker, A.E., Evett, S.R., Brauer, D.K., Coyle, G.G., Copeland, K.S., Marek, G.W.,  
Colaizzi, P.D., Acutis, M., Alimagham, S.M., Archontoulis, S., Babacar, F., Barcza, Z., Basso, B., Bertuzzi, P.,  
Constantin, J., De Antoni Migliorati, M., Dumont, B., Durand, J.L., Fodor, N., Gaiser, T., Garofalo, P., Gayler, S.,  
985 Giglio, L., Grant, R., Guan, K., Hoogenboom, G., Kim, S.H., Kisekka, I., Lizaso, J., Masia, S., Meng, H., Mereu,  
V., Mukhtar, A., Perego, A., Peng, B., Priesack, E., Shelia, V., Snyder, R., Soltani, A., Spano, D., Srivastava, A.,  
Thomson, A., Timlin, D., Trabucco, A., Webber, H., Willaume, M., Williams, K., van der Laan, M., Ventrella, D.,  
Viswanathan, M., Xu, X., Zhou, W.: Evaluation of multimodel averaging approaches for ensembling  
evapotranspiration and yield simulations from maize models. *Journal of Hydrology* 661: 133631.  
<https://doi.org/10.1016/j.jhydrol.2025.133631>, 2025.
- 990 Olesen, J.E., Bindi, M.: Consequences of climate change for European agricultural productivity, land use and policy.  
*European Journal of Agronomy* 16, 239-262. [https://doi.org/10.1016/S1161-0301\(02\)00004-7](https://doi.org/10.1016/S1161-0301(02)00004-7), 2002.
- Opitz, D., Maclin, R.: Popular ensemble methods: An empirical study. *Journal of Artificial Intelligence Research*  
995 11, 169-198. <https://doi.org/10.1613/jair.614>, 1999.
- Ostle, N.J., Smith, P., Fisher, R., Woodward, F.I., Fisher, J.B., Smith, J.U., Galbraith, D., Levy, P., Meir, P.,  
McNamara, N.P., Bardgett, R.D.: Integrating plant-soil interactions into global carbon cycle models. *Journal of  
Ecology* 97, 851-863. <https://doi.org/10.1111/j.1365-2745.2009.01547.x>, 2009.
- 1000 Pappas, C., Papalexiou, S.M., Koutsoyiannis D.: A quick gap filling of missing hydrometeorological data. *Journal  
of Geophysical Research Atmosphere* 119, 9290-9300. <https://doi.org/10.1002/2014JD021633>, 2014.



- 1005 Raj, R., Hamm, N.A.S., van de Tol, C., Stein, A.: Uncertainty analysis of gross primary production partitioned from net ecosystem exchange measurements. *Biogeosciences* 13, 1409-1422. <https://doi.org/10.5194/bg-13-1409-2016>, 2006.
- 1010 Reichstein, M., Bahn, M., Ciais, P., Frank, D., Mahecha, M. D., Seneviratne, S. I., Zscheischler, J., Beer, C., Buchmann, N., Frank, D. C., Papale, D., Rammig, A., Smith, P., Thonicke, K., van der Velde, M., Vicca, S., Walz, A., Wattenbach, M.: Climate extremes and the carbon cycle. *Nature* 500, 287–295. <https://doi.org/10.1038/nature12350>, 2013.
- 1015 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., Prabhat: Deep learning and process understanding for data-driven Earth system science. *Nature* 566, 195–204. <https://doi.org/10.1038/s41586-019-0912-1>, 2019.
- 1020 Riccio, G., Giunta, G., Galmarini, S.: Seeking for the rational basis of the Median Model: the optimal combination of multi-model ensemble results. *Atmospheric Chemistry and Physics* 7, 6085-6098. <https://doi.org/10.5194/acp-7-6085-2007>, 2007.
- Richter, K., Atzberger, C., Hank, T.B., Mauser, W.: Derivation of biophysical variables from Earth observation data: validation and statistical measures. *Journal of Applied Remote Sensing* 6, 063557. <https://doi.org/10.1117/1.JRS.6.063557>, 2012.
- 1025 Reichstein, M., Ciais, P., Papale, D., Valentini, R., Running, S., Viovy, N., Cramer, W., Granier, A., Ogee, J., Allard, V., Aubinet, M., Bernhofer, C., Buchmann, N., Carrara, A., Grünwald, T., Heimann, M., Heinesch, B., Knohl, A., Kutsch, W., Loustau, D., Manca, G., Matteucci, G., Miglietta, F., Ourcival, J.M., Pilegaard, K., Pumpanen, J., Rambal, S., Schaphoff, S., Seufert, G., Soussana, J.-F., Sanz, M.-J., Vesala, T., Zhao, M.: Reduction of ecosystem productivity and respiration during the European summer 2003 climate anomaly: A joint flux tower, remote sensing and modelling analysis. *Global Change Biology* 13, 634-651. <https://doi.org/10.1111/j.1365-2486.2006.01224.x>, 2007.
- 1030 Reichstein, M., Falge, E., Baldocchi, D., Papale, D., Aubinet, M., Berbigier, P., Bernhofer, C., Buchmann, N., Gilmanov, T., Granier, A., Grünwald, T., Havránková, K., Ilvesniemi, H., Janous, D., Knohl, A., Laurila, T., Lohila, A., Loustau, D., Matteucci, G., Meyers, T., Miglietta, F., Ourcival, J.-M., Pumpanen, J., Rambal, S., Rotenberg, E., Sanz, M., Tenhunen, J., Seufert, G., Vaccari, F., Vesala, T., Yakir, D., Valentini, R.: On the separation of net ecosystem exchange into assimilation and ecosystem respiration: Review and improved algorithm. *Global Change Biology* 11, 1424–1439. <https://doi.org/10.1111/j.1365-2486.2005.001002.x>, 2005.
- 1040 Robeson, S.M., Willmott, C.J.: Decomposition of the mean absolute error (MAE) into systematic and unsystematic components. *PLoS ONE* 18, e0279774. <https://doi.org/10.1371/journal.pone.0279774>, 2023.



- Rosenzweig, C., Elliott, J., Deryng, D., Ruane, A.C., Müller, C., Arneth, A., Boote, K.J., Folberth, C., Glotter, M., Khabarov, N., Neumann, K., Piontek, F., Pugh, T.A.M., Schmid, E., Stehfest, E., Yang, H., Jones, J.W.: Assessing agricultural risks of climate change in the 21st century in a global gridded crop model intercomparison. Proceedings of the National Academy of Sciences of the United States of America 111, 3268-3273. <https://doi.org/10.1073/pnas.1222463110>, 2014.
- 1045
- Ruane, A.C., Hudson, N.I., Asseng, S., Camarrano, D., Ewert, F., Martre, P., Boote, K.J., Thorburn, P.J., Aggarwal, P.K., Angulo, C., Basso, B., Bertuzzi, P., Biernath, C., Brisson, N., Challinor, A.J., Doltra, J., Gayler, S., Goldberg, R., Grant, R.F., Heng, L., Hooker, J., Hunt, L.A., Ingwersen, J., Izaurralde, R.C., Kersebaum, K.C., Kumar, S.N., Müller, C., Nendel, C., O'Leary, G., Olesen, J.E., Osborne, T.M., Palosuo, T., Priesack, E., Ripoche, D., Rötter, R.P., Semenov, M.A., Shcherbak, I., Steduto, P., Stöckle, C.O., Stratonovitch, P., Streck, T., Supit, I., Tao, F., Travasso, M., Waha, K., Wallach, D., White, J.W., Wolf, J.: Multi-wheat-model ensemble responses to interannual climate variability. Environmental Modelling & Software 81, 86-101. <https://doi.org/10.1016/j.envsoft.2016.03.008>, 2016.
- 1050
- 1055
- Ruane, A.C., Rosenzweig, C., Asseng, S., Boote, K.J., Elliott, J., Ewert, F., Jones, J.W., Martre, P., McDermid, S.P., Müller, C., Snyder, A., Thorburn, P.J.: An AgMIP framework for improved agricultural representation in integrated assessment models. Environmental Research Letters 12, 125003. <https://doi.org/10.1088/1748-9326/aa8da6>, 2017.
- 1060
- Sagi, O., Rokach, L.: Ensemble learning: A survey. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery 8, e1249. <https://doi.org/10.1002/widm.1249>, 2018.
- 1065
- Sándor, R., Barcza, Z., Acutis, M., Doro, L., Hidy, D., Köchy, M., Minet, J., Lellei-Kovács, E., Ma, S., Perego, A., Rolinski, S., Ruget, F., Sanna, M., Seddaiu, G., Wu, L., Bellocchi, G.: Multi-model simulation of soil temperature, soil water content and biomass in Euro-Mediterranean grasslands: Uncertainties and ensemble performance. European Journal of Agronomy 88, 22-40. <https://doi.org/10.1016/j.eja.2016.06.006>, 2017.
- 1070
- Sándor, R., Ehrhardt, F., Basso, B., Bellocchi, G., Bhatia, A., Brilli, L., Migliorati, M.D., Doltra, J., Dorich, C., Doro, L., Fitton, N., Giacomini, S.J., Grace, P., Grant, B., Harrison, M.T., Jones, S., Kirschbaum, M.U.F., Klumpp, K., Laville, P., Léonard, J., Liebig, M., Liewering, M., Martin, R., McAuliffe, R., Meier, E., Merbold, L., Moore, A., Myrgiotis, V., Newton, P., Pattey, E., Recous, S., Rolinski, S., Sharp, J., Massad, R.S., Smith, P., Smith, W., Snow, V., Wu, L., Zhang, Q., Soussana, J.-F.: C and N models Intercomparison – benchmark and ensemble model estimates for grassland production. Advances in Animal Biosciences 7, 245-247. <https://doi.org/10.1017/S2040470016000297>, 2016.
- 1075
- Sándor, R., Ehrhardt, F., Brilli, L., Carozzi, M., Recous, S., Smith, P., Snow, V., Soussana, J.F., Dorich, C.D., Fuchs, K., Fitton, N., Gongadze, K., Klumpp, K., Liebig, M., Martin, R., Merbold, L., Newton, P.C.D., Rees, R.M., Rolinski, S., Bellocchi, G.: The use of biogeochemical models to evaluate mitigation of greenhouse gas emissions from managed grasslands. Science of the Total Environment 15, 292-306. <https://doi.org/10.1016/j.scitotenv.2018.06.020>, 2018.
- 1080



- 1085 Sándor, R., Ehrhardt, F., Grace, P., Recous, S., Smith, P., Snow, V., Soussana, J.-F., Basso, B., Bhatia, A., Brilli, L., Doltra, J., Dorich, C.D., Doro, L., Fitton, N., Grant, B., Harrison, M.T., Kirschbaum, M.U.F., Klumpp, K., Laville, P., Léonard, J., Martin, R., Massad, R.S., Moore, A., Myrگیotis, V., Pattey, E., Rolinski, R., Sharp, J., Skiba, U., Smith, W., Wu, L., Zhang, Q., Bellocchi, G.: Ensemble modelling of carbon fluxes in grasslands and croplands. *Field Crops Research* 252, 107791. <https://doi.org/10.1016/j.fcr.2020.107791>, 2020.
- 1090 Sándor, R., Ehrhardt, F., Grace, P., Recous, S., Smith, P., Snow, V., Soussana, J.-F., Basso, B., Bhatia, A., Brilli, L., Doltra, J., Dorich, C.D., Doro, L., Fitton, N., Grant, B., Harrison, M.T., Kirschbaum, M.U.F., Klumpp, K., Laville, P., Léonard, J., Martin, R., Massad, R.S., Moore, A., Myrگیotis, V., Pattey, E., Rolinski, R., Sharp, J., Skiba, U., Smith, W., Wu, L., Zhang, Q., Bellocchi, G.: Experimental and simulated data for crop and grassland production and carbon-nitrogen fluxes. *Open Data Journal for Agricultural Research* 24, 22-27. Available at:  
1095 <https://odjar.org/article/view/18594>, 2024.
- Sándor, R., Ehrhardt, F., Grace, P., Recous, S., Smith, P., Snow, V., Soussana, J.-F., Basso, B., Bhatia, A., Brilli, L., Doltra, J., Dorich, C.D., Doro, L., Fitton, N., Grant, B., Harrison, M.T., Skiba, U., Kirschbaum, M.U.F., Klumpp, K., Laville, P., Léonard, J., Martin, R., Massad, R.S., Moore, A., Myrگیotis, V., Pattey, E., Rolinski, R., Sharp, J.,  
1100 Smith, W., Wu, L., Zhang, Q., Bellocchi, G.: Residual correlation and ensemble modelling to improve crop and grassland models. *Environmental Modelling & Software* 161, 105625. <https://doi.org/10.1016/j.envsoft.2023.105625>, 2023.
- Schwalm, C.R., Williams, C.A., Schaefer, K., Arneث, A., Bonal, D., Buchmann, N., Chen, J., Law, B.E., Lindroth, A., Luysaert, S., Reichstein, M., Richardson, A.D.: Assimilation exceeds respiration sensitivity to drought: A  
1105 FLUXNET synthesis. *Global Change Biology* 16, 657–670. <https://doi.org/10.1111/j.1365-2486.2009.01991.x>, 2010.
- Scowen, M., Athanasiadis, I. N., Bullock, J. M., Eigenbrod, F., Willcock, S.: The current and future uses of machine  
1110 learning in ecosystem service research. *Science of the Total Environment* 799, 149263. <https://doi.org/10.1016/j.scitotenv.2021.149263>, 2021.
- Shahhosseini, M., Hu, G., Archontoulis, S.V.: Forecasting corn yield with machine learning ensembles. *Frontiers in  
1115 Plant Science* 11, 1120. <https://doi.org/10.3389/fpls.2020.01120>, 2020.
- Shahhosseini, M., Hu, G., Huber, I., Archontoulis, S.V.: Coupling machine learning and crop modeling improves  
crop yield prediction in the US Corn Belt. *Scientific Reports* 11, 1606. <https://doi.org/10.1038/s41598-020-80820-1>, 2021.
- 1120 Shapley, L.S.: A value for n-person games. In: Kuhn H.W., Tucker A.W. (eds.), *Contributions to the theory of games II* (Vol. 28, pp. 307–317). Princeton University Press. 1953.



- Smith, P., House, J.I., Bustamante, M., Sobocká, J., Harper, R., Pan, G., West, P.C., Clark, J.M., Adhya, T., Rumpel, C., Paustian, K., Kuikman, P., Cotrufo, M.F., Elliott, J.A., McDowell, R., Griffiths, R.I., Asakawa, S., Bondeau, A.,  
1125 Jain, A.K., Meersmans, J., Pugh, T.A.M.: Global change pressures on soils from land use and management. *Global Change Biology* 22, 1008-1028. <https://doi.org/10.1111/gcb.13068>, 2016.
- Snow, V., Rotz, C.A., Moore, A.D., Martin-Clouaire, R., Johnson, I.R., Hutchings, N.J., Eckard, R.J.: The challenges  
- and some solutions - to process-based modelling of grazed agricultural systems. *Environmental Modelling &*  
1130 *Software* 62, 420-436. <https://doi.org/10.1016/j.envsoft.2014.03.009>, 2014.
- Sroufe, R., Watts, A., 2022. Pathways to agricultural decarbonization: Climate change obstacles and opportunities  
in the US. *Resources, Conservation and Recycling* 182, 106276. <https://doi.org/10.1016/j.resconrec.2022.106276>
- 1135 Stoy, P.C., Mauder, M., Foken, T., Marcolla, B., Boegh, E., Ibrom, A., Arain, M.A., Arneth, A., Aurela, M.,  
Bernhofer, C., Cescatti, A., Dellwik, E., Duce, P., Gianelle, D., van Gorsel, E., Kiely, G., Knohl, A., Margolis, H.,  
McCaughy, H., Merbold, L., Montagnani, L., Papale, D., Reichstein, M., Saunders, M., Serrano-Ortiz, P.,  
Sottocornola, M., Spano, D., Vaccari, F., Varlagin, A.: A data-driven analysis of energy balance closure across  
FLUXNET research sites: The role of landscape scale heterogeneity. *Agricultural and Forest Meteorology* 171-172,  
1140 137-152. <https://doi.org/10.1016/j.agrformet.2012.11.004>, 2013.
- Strobl, C., Boulesteix, A. L., Zeileis, A., Hothorn, T.: Bias in random forest variable importance measures:  
Illustrations, sources and a solution. *BMC Bioinformatics* 8, 25. <https://doi.org/10.1186/1471-2105-8-25>, 2007.
- 1145 Sun, W., Zhou, S., Yu, B., Zhang, Y., Keenan, T.F., Fu, B.: Soil moisture-atmosphere interactions drive terrestrial  
carbon-water trade-offs. *Communications Earth & Environment* 6, 169. [https://doi.org/10.1038/s43247-025-02145-](https://doi.org/10.1038/s43247-025-02145-z)  
[z](https://doi.org/10.1038/s43247-025-02145-z), 2025.
- Therond, O., Hengsdijk, H., Casellas, E., Wallach, D., Adam, M., Belhouchette, H., Oomen, R., Russell, G., Ewert,  
1150 F., Bergez, J.-E., Janssen, S., Wery, J., van Ittersum, M.K.: Using a cropping system model at regional scale: low-  
data approaches for crop management information and model calibration. *Agriculture, Ecosystems & Environment*  
142, 85-94. <https://doi.org/10.1016/j.agee.2010.05.007>, 2011.
- Thornton, P.K., Whitbread, A., Baedeker, T., Cairns, J., Claessens, L., Beethgen, W., Bunn, C., Friedmann, M.,  
1155 Giller, K.E., Herrero, M., Howden, M., Kilcline, K., Nangia, V., Ramirez-Villegas, J., Kumar, S., West, P.C.,  
Keating, B.: A framework for priority-setting in climate smart agriculture research. *Agricultural Systems* 167, 161-  
175. <https://doi.org/10.1016/j.agsy.2018.09.009>, 2018.
- Valin, H., Havlík, P., Mosnier, A., Herrero, M., Schmid, E., Obersteiner, M.: Agricultural productivity and  
1160 greenhouse gas emissions: trade-offs or synergies between mitigation and food security? *Environmental Research*  
*Letters* 8, 035019. <https://doi.org/10.1088/1748-9326/8/3/035019>, 2013



- 1165 Van der Laan, M.J., Polley, E.C., Hubbard, A.E.: Super learner. *Statistical Applications in Genetics and Molecular Biology* 6, 25. <https://doi.org/10.2202/1544-6115.1309>, 2007.
- 1170 Van der Velde, M., Tubiello, F.N., Vrieling, A., Bouraoui, F.: Impacts of extreme weather on wheat and maize in France: Evaluating regional crop simulations against observed data. *Climatic Change* 123, 699–711. <https://doi.org/10.1007/s10584-011-0368-2>, 2014.
- 1175 Wallach, D., Martre, P., Liu, B., Asseng, S., Ewert, F., Thonburn, P.J., van Ittersum, M., Aggarwal, P.K., Ahmed, M., Basso, B., Biernath, C., Cammarano, D., Challinor, A.J., De Sanctis, G., Dumont, B., Rezaei, E.E., Fereres, E., Fitzgerald, G.J., Gao, Y., Garcia-Vila, M., Gayler, S., Girousse, C., Hoogenboom, G., Horan, H., Izaurralde, R.C., Jones, C.D., Kassie, B.T., Kersebaum, K.C., Klein, C., Koehler, A.-K., Maiorano, A., Minoli, S., Müller, C., Kumar, S.N., Nendel, C., O’Leary, G.J., Palosuo, T., Priesack, E., Ripoche, D., Rötten, R.P., Semenov, M.A., Stöckle, C., Stratonovitch, P., Streck, T., Supit, I., Fao, F., Wolf, J., Zhang, Z.: Multi-model ensembles improve predictions of crop-environment-management interactions. *Global Change Biology* 24, 5072–5083. <https://doi.org/10.1111/gcb.14411>, 2018.
- 1180 Wang, Z., Liu, Z., Huang, M.: NDVI joint process-based models drive a learning ensemble model for accurately estimating cropland net primary productivity (NPP). *Frontiers in Environmental Science* 11, 1304400. <https://doi.org/10.3389/fenvs.2023.1304400>, 2024.
- 1185 Wolpert, D.H., Macready, W.G.: Coevolutionary free lunches. *IEEE Transactions on Evolutionary Computation* 9, 721–35. <https://doi.org/10.1109/TEVC.2005.856205>, 2005.
- 1190 Xia, J., Chen, T., Zhang, K., Wang, Y., Chen, G.: Impacts of climate extremes on carbon fluxes and their underlying mechanisms in a typical temperate grassland ecosystem. *Science of the Total Environment* 907, 167755. <https://doi.org/10.1016/j.scitotenv.2023.167755>, 2024.
- 1195 Xu, L., Baldocchi, D.D.: Seasonal variation in carbon dioxide exchange over a Mediterranean annual grassland in California. *Agricultural and Forest Meteorology* 123, 79–96. <https://doi.org/10.1016/j.agrformet.2003.10.004>, 2004.
- 1200 Xu, X., Xu, J., Li, B., Li, J., Nie, M.: Ecosystem carbon fluxes exhibit thermal response thresholds at which carbon–climate feedback changes. *Global Ecology and Biogeography* 34, e70030. <https://doi.org/10.1111/geb.70030>, 2025.
- Zhang, J., Tian, H., Wang, P., Tansey, K., Zhang, S., Li, H.: Improving wheat yield estimates using data augmentation models and remotely sensed biophysical indices within deep neural networks in the Guanzhong Plain, PR China. *Computers and Electronics in Agriculture* 192, 106616. <https://doi.org/10.1016/j.compag.2021.106616>, 2022.



Zhu, Z., Piao, S., Myneni, R.B., Huang, M., Zeng, Z., Canadell, J. G., Ciais, P., Sitch, S., Friedlingstein, P., Arneeth, A., Cao, C., Cheng, L., Kato, E., Koven, C., Li, Y., Lian, X., Liu, Y., Liu, R., Mao, J., Pan, Y., Peng, S., Peñuelas, J., Poulter, B., Pugh, T.A.M., Stocker, B. D., Viovy, N., Wang, X., Wang, Y., Xiao, Z., Yang, H., Zaehle, S., Zeng, N.: Greening of the Earth and its drivers. *Nature Climate Change* 6, 791–795. <https://doi.org/10.1038/nclimate3004>, 2016.

1205

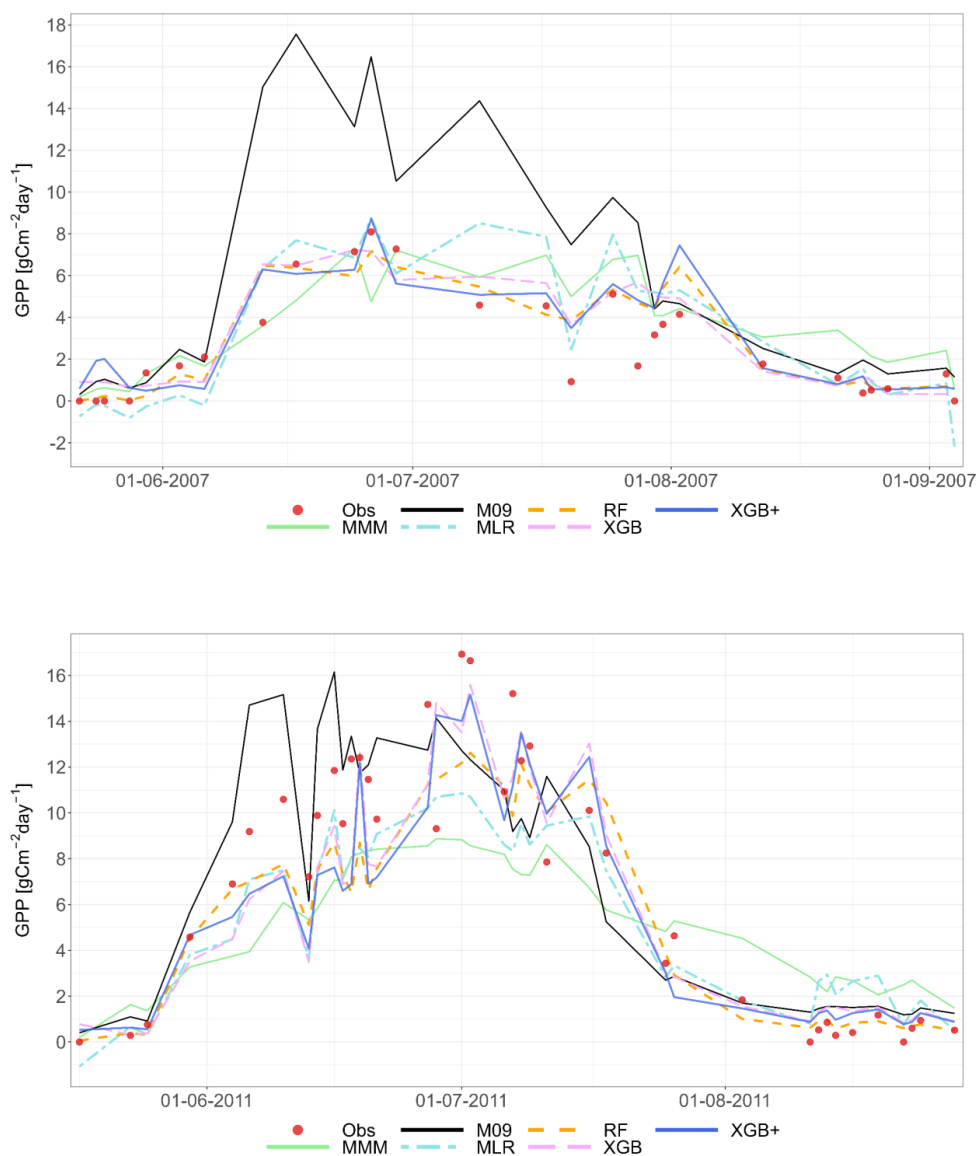
1210



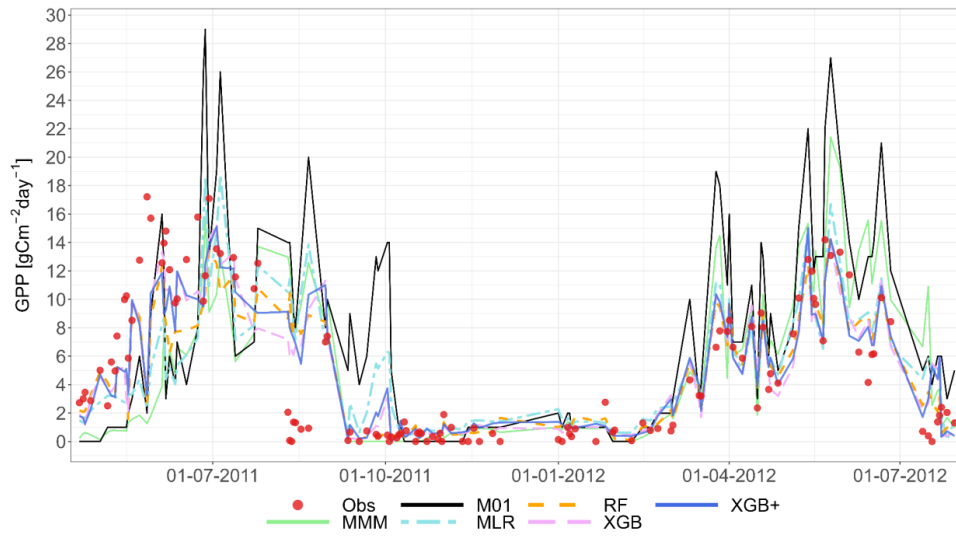
1215 **Appendix A**

The following figures demonstrate the simulated fluxes of GPP, RECO and NEE for all sites and years included in the study, except those that are presented in the main text. Graphs are presented with a very short caption only containing: flux type, site name and type, covered years. See the main text for explanations.

1220



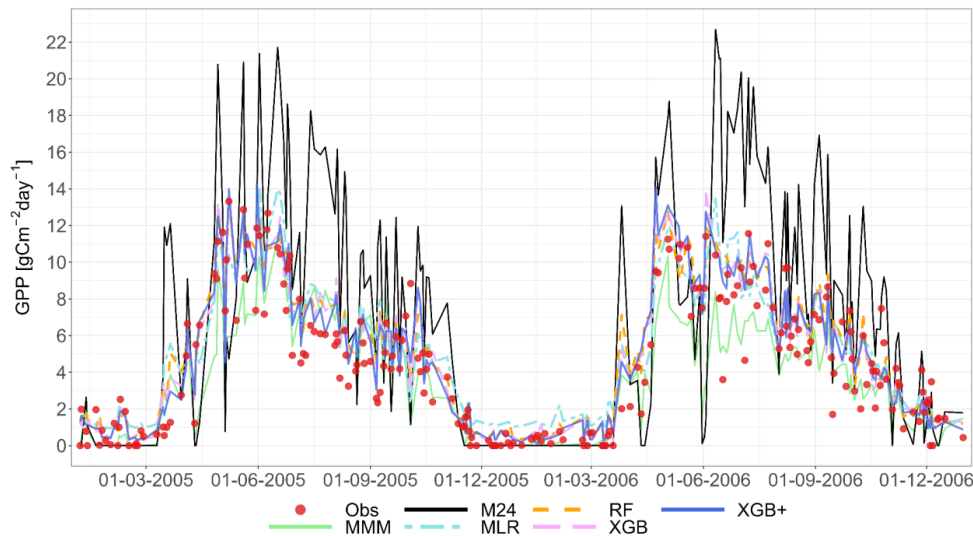
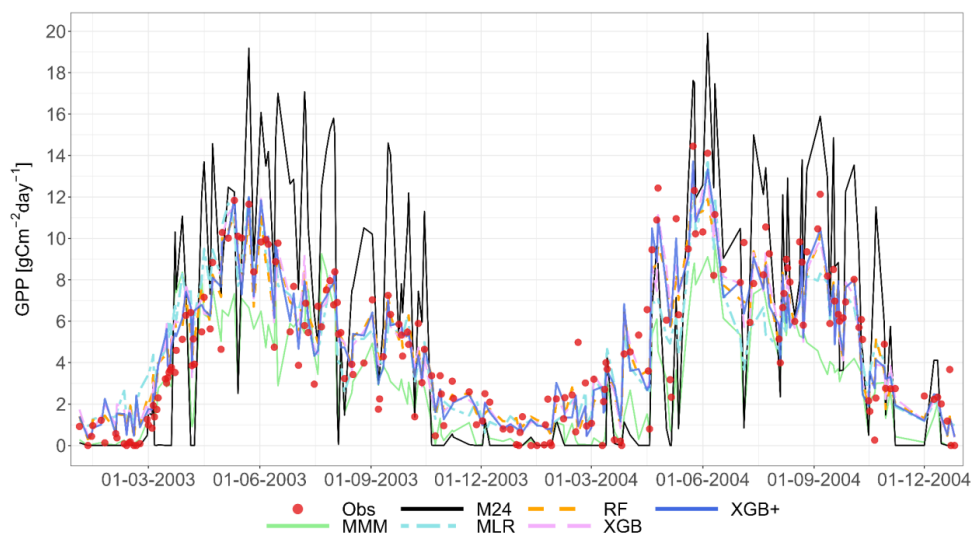
**Figure A1:** GPP, C1 - Ottawa (CA), cropland, 2007 and 2011.



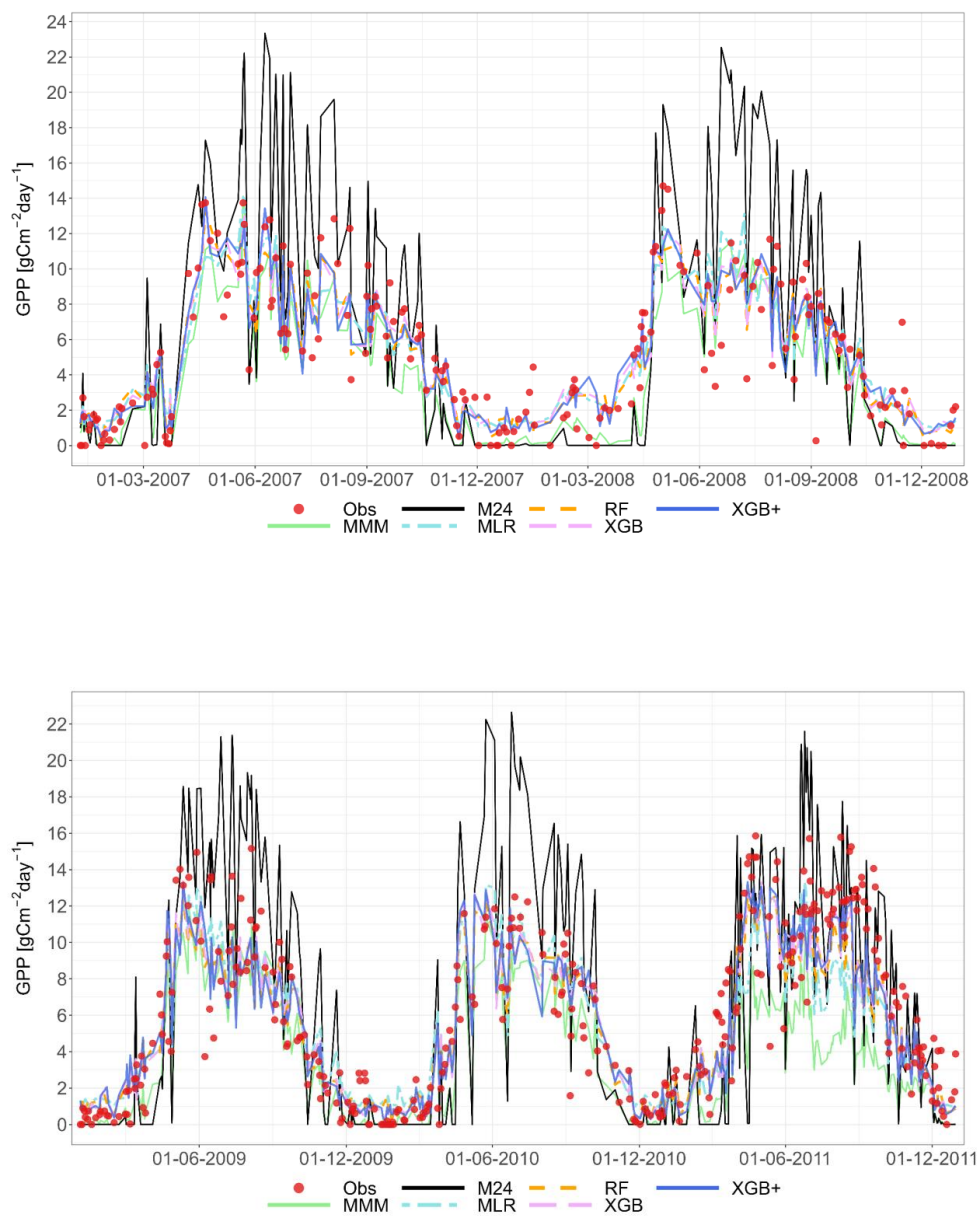
1225

**Figure A2:** GPP, C2 - Grignon, (FR), cropland, 2011 and 2012.

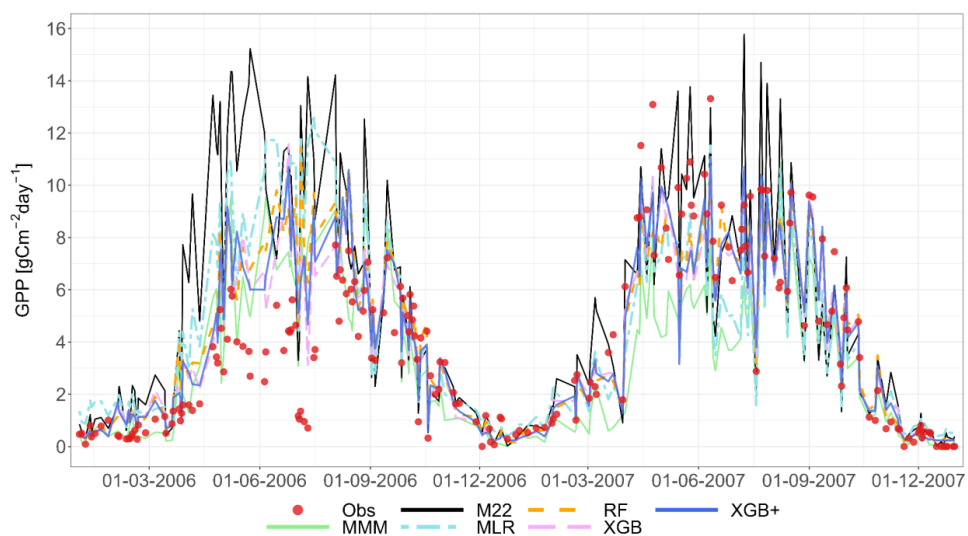
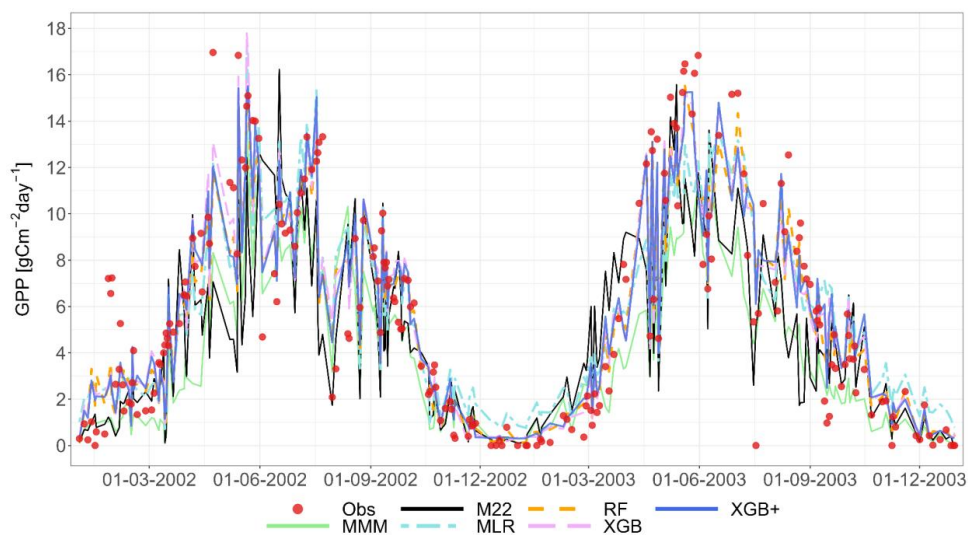
1230



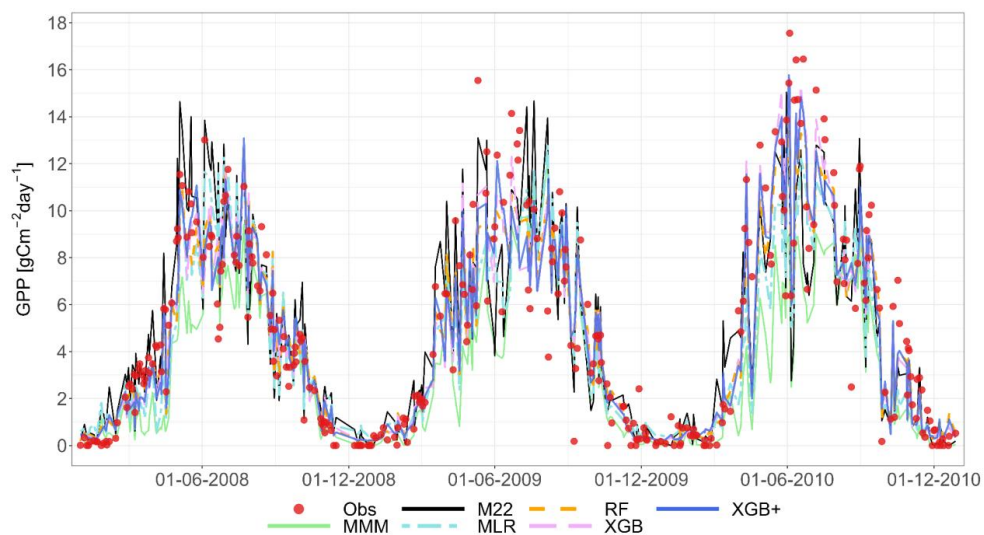
1235



**Figure A3:** GPP, G3 - Laqueuille (FR), grassland, 2003-2011.

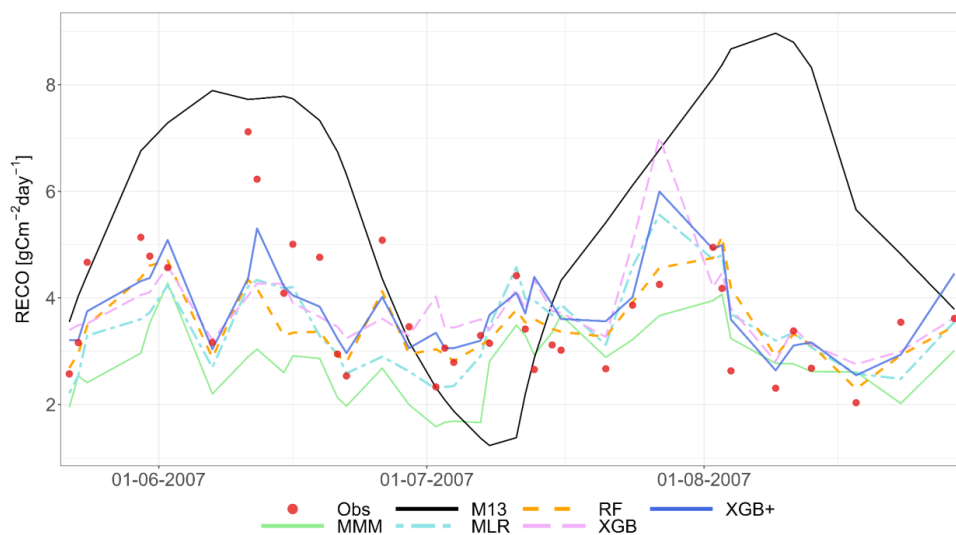


1250

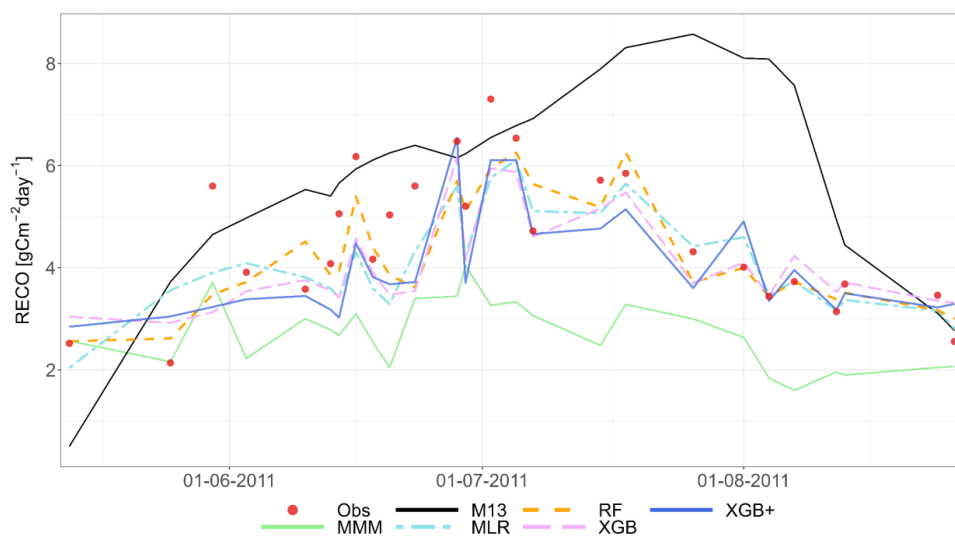


**Figure A4:** GPP, G4 - Easter Bush (UK), grassland

1255

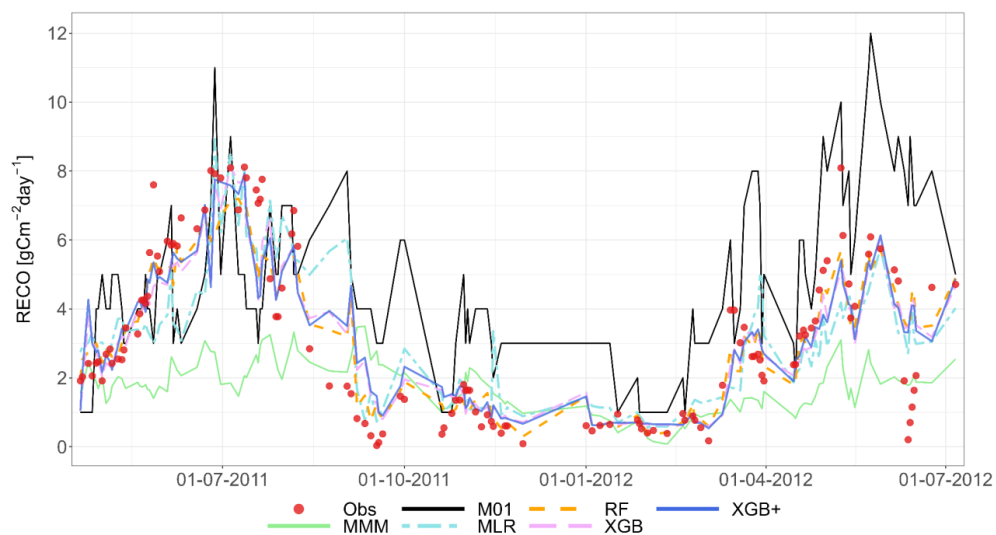


1260



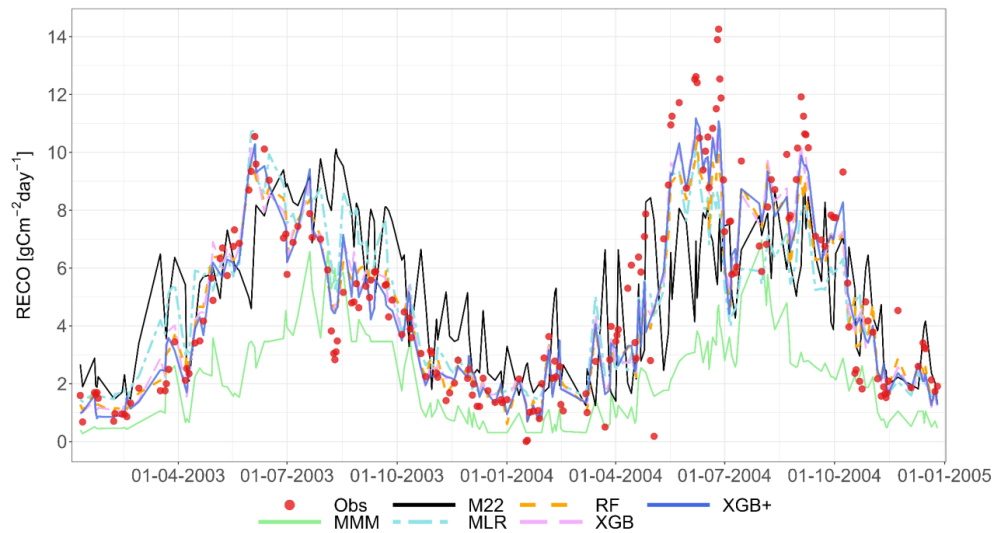
**Figure A5:** RECO, C1 - Ottawa (CA), cropland, 2007 and 2011.

1265

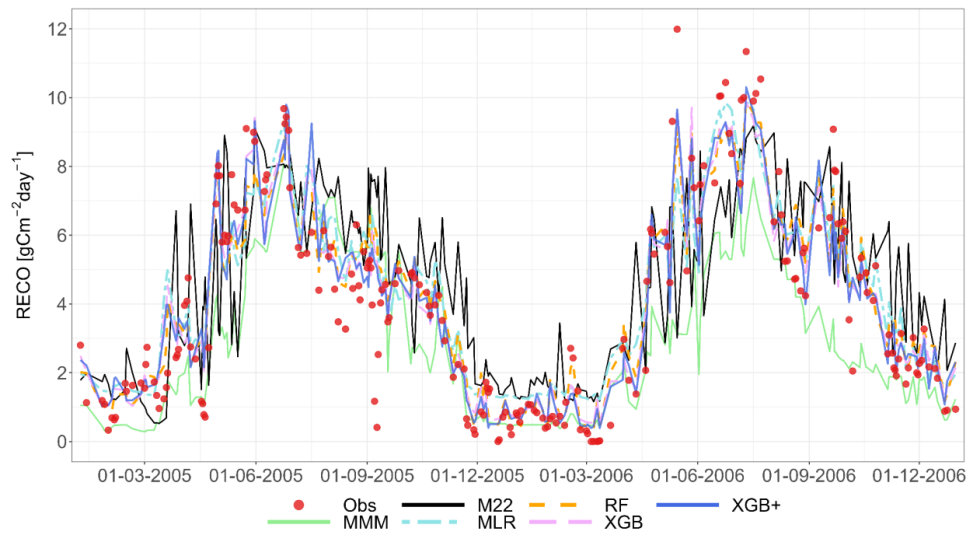


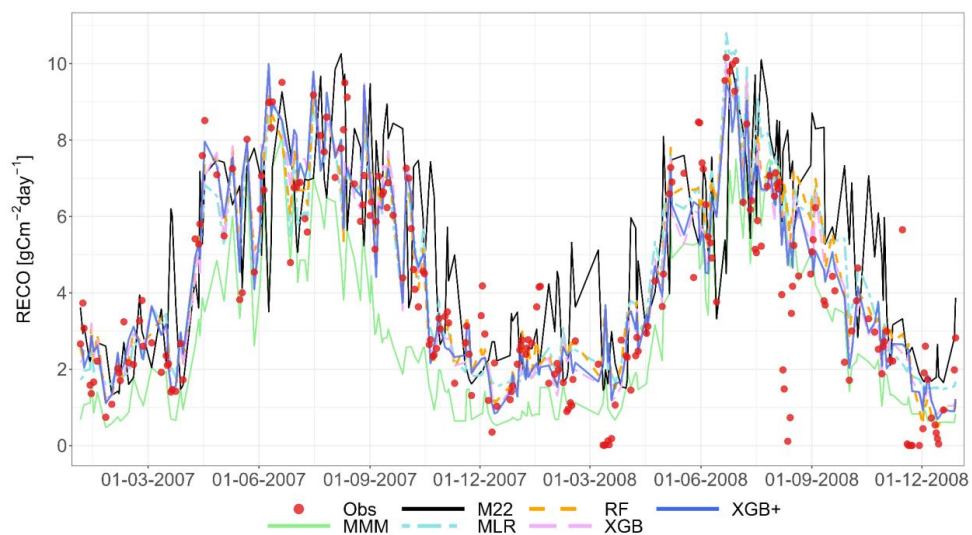
**Figure A6:** RECO, C2 - Grignon, (FR), cropland, 2011 and 2012.

1270



1275





1280

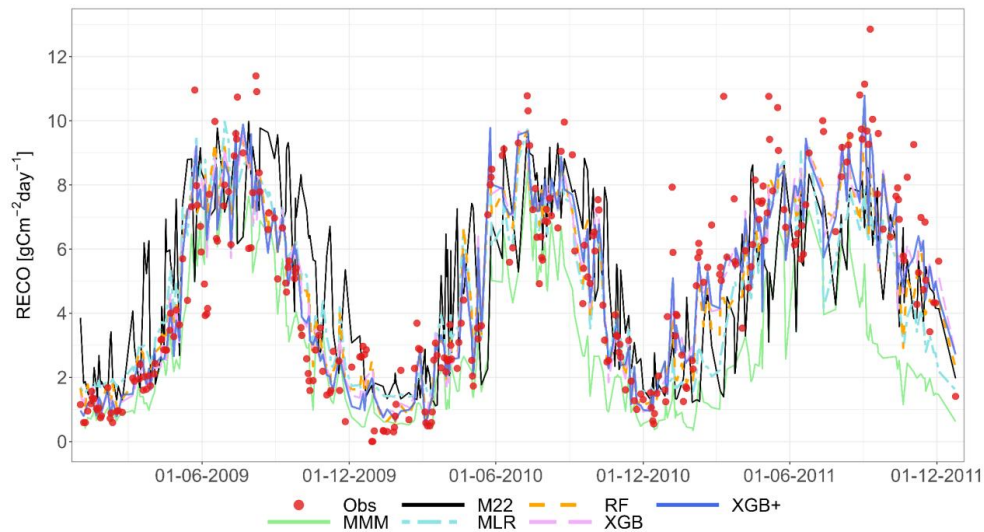
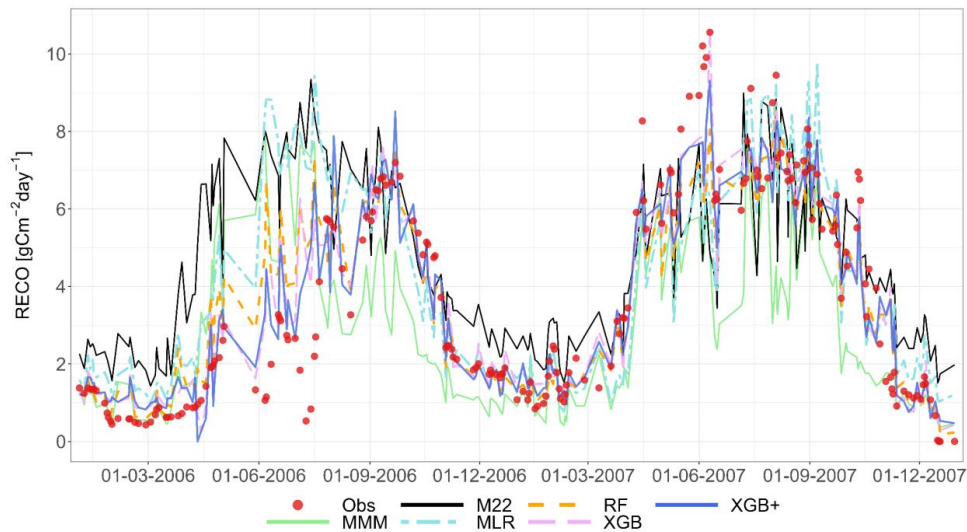
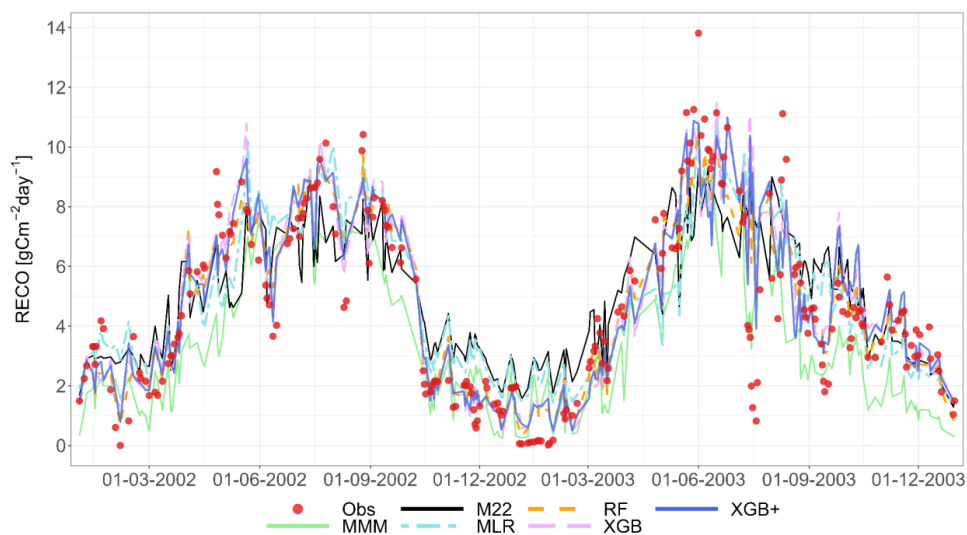
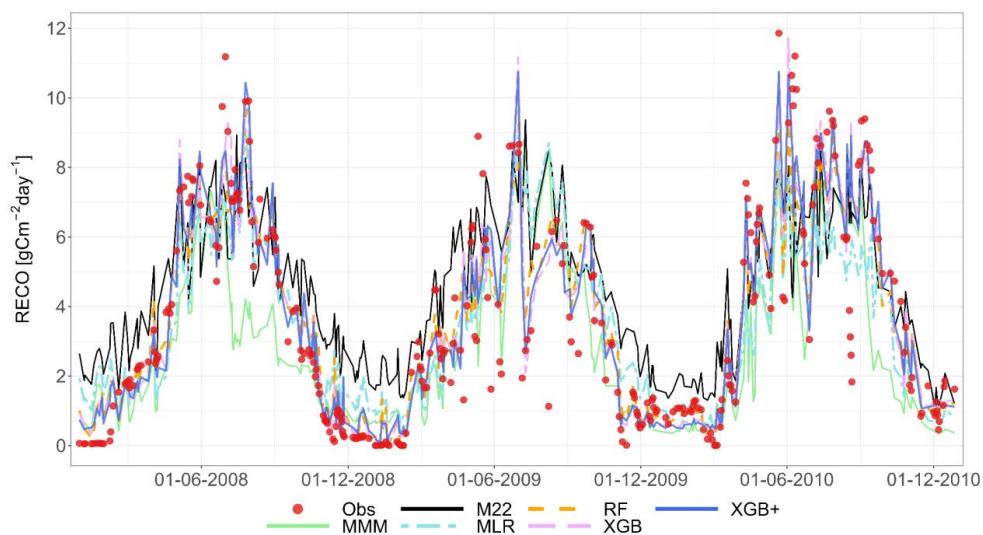


Figure A7: RECO, G3 - Laqueuille (FR), grassland, 2003-2011.

1285



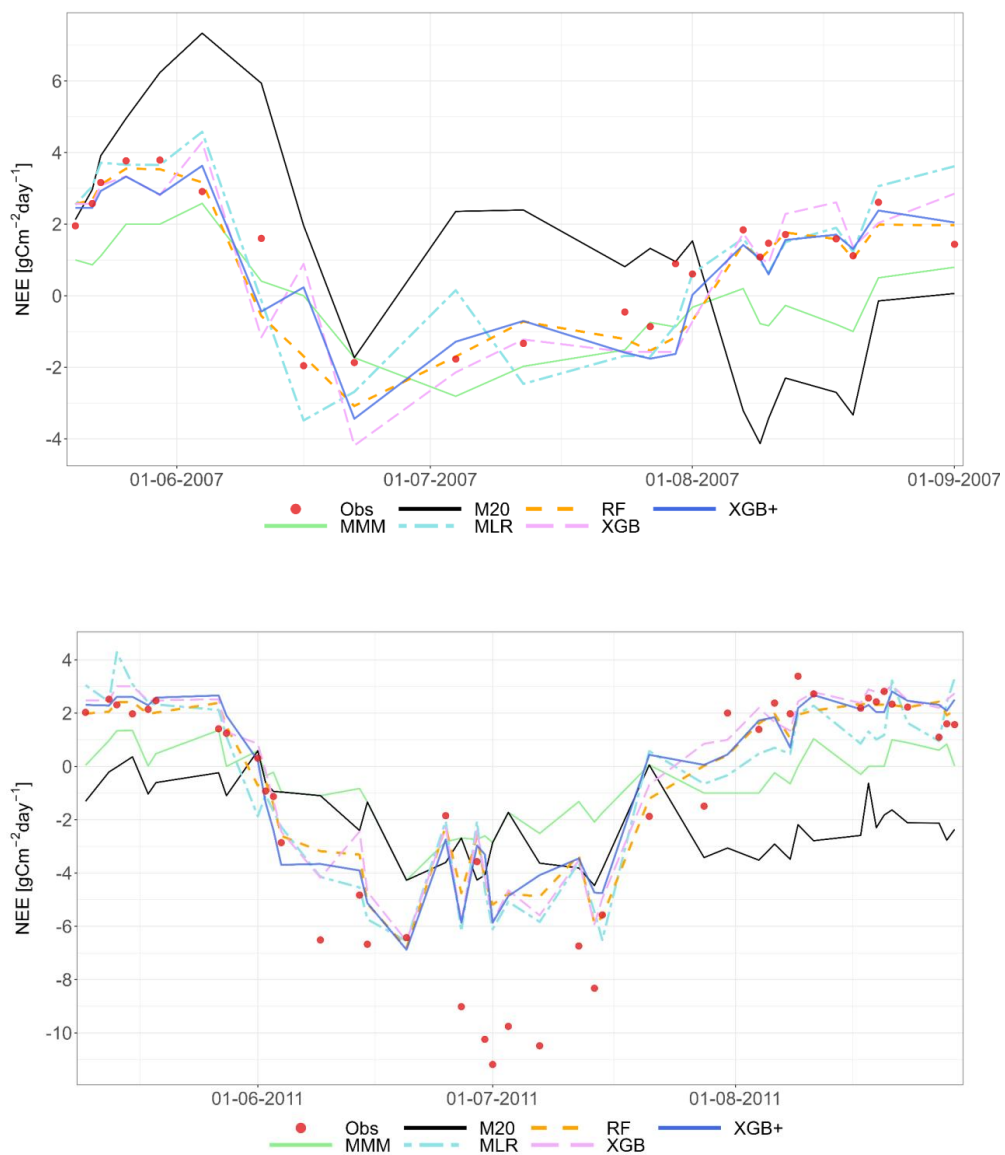
1290



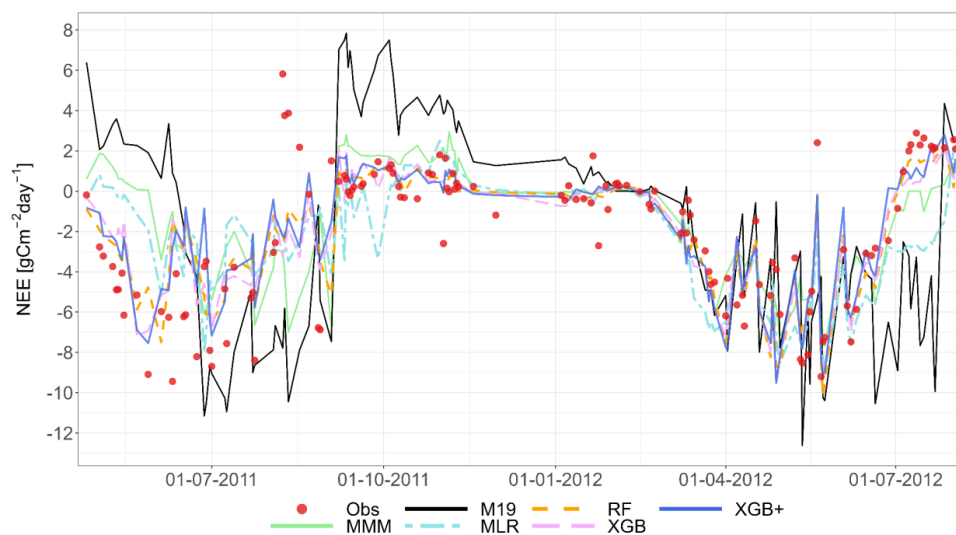
**Figure A8:** RECO, G4 - Easter Bush (UK), grassland

1295

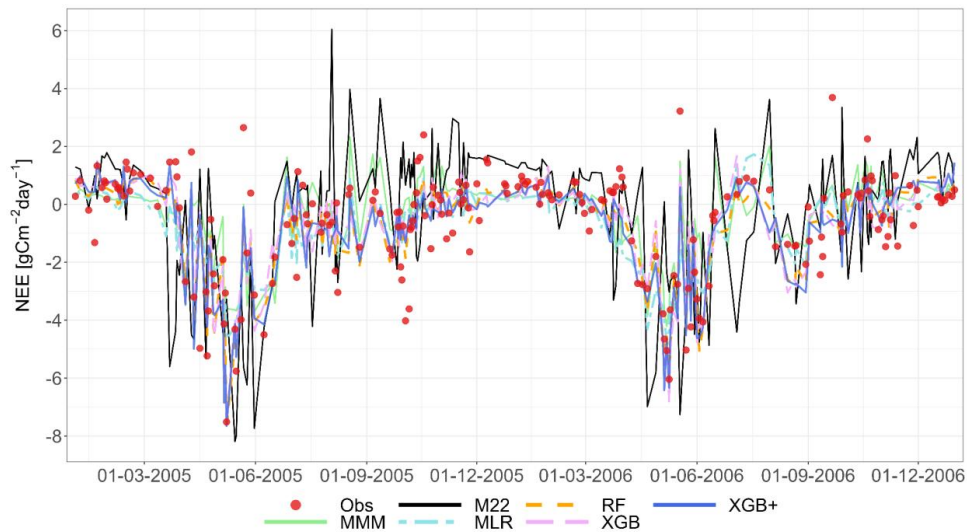
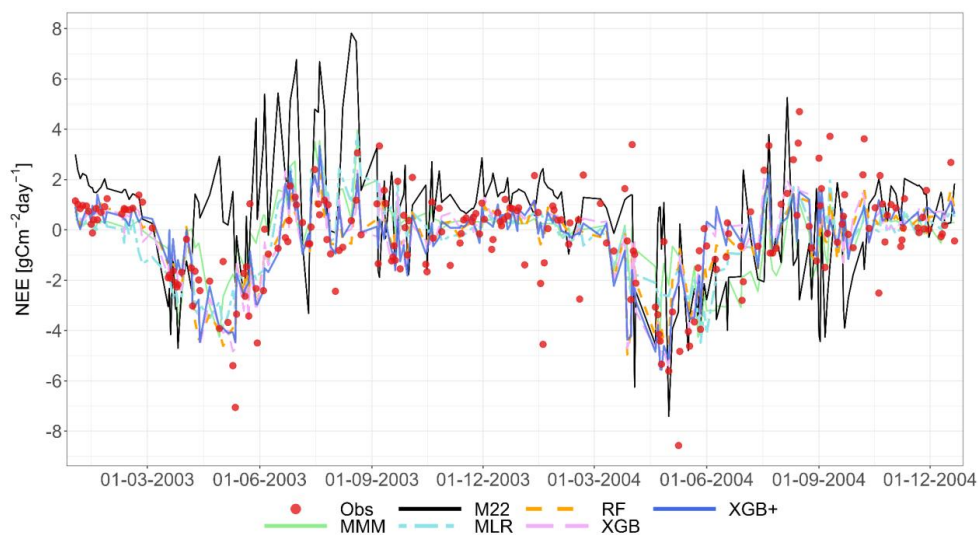
1300



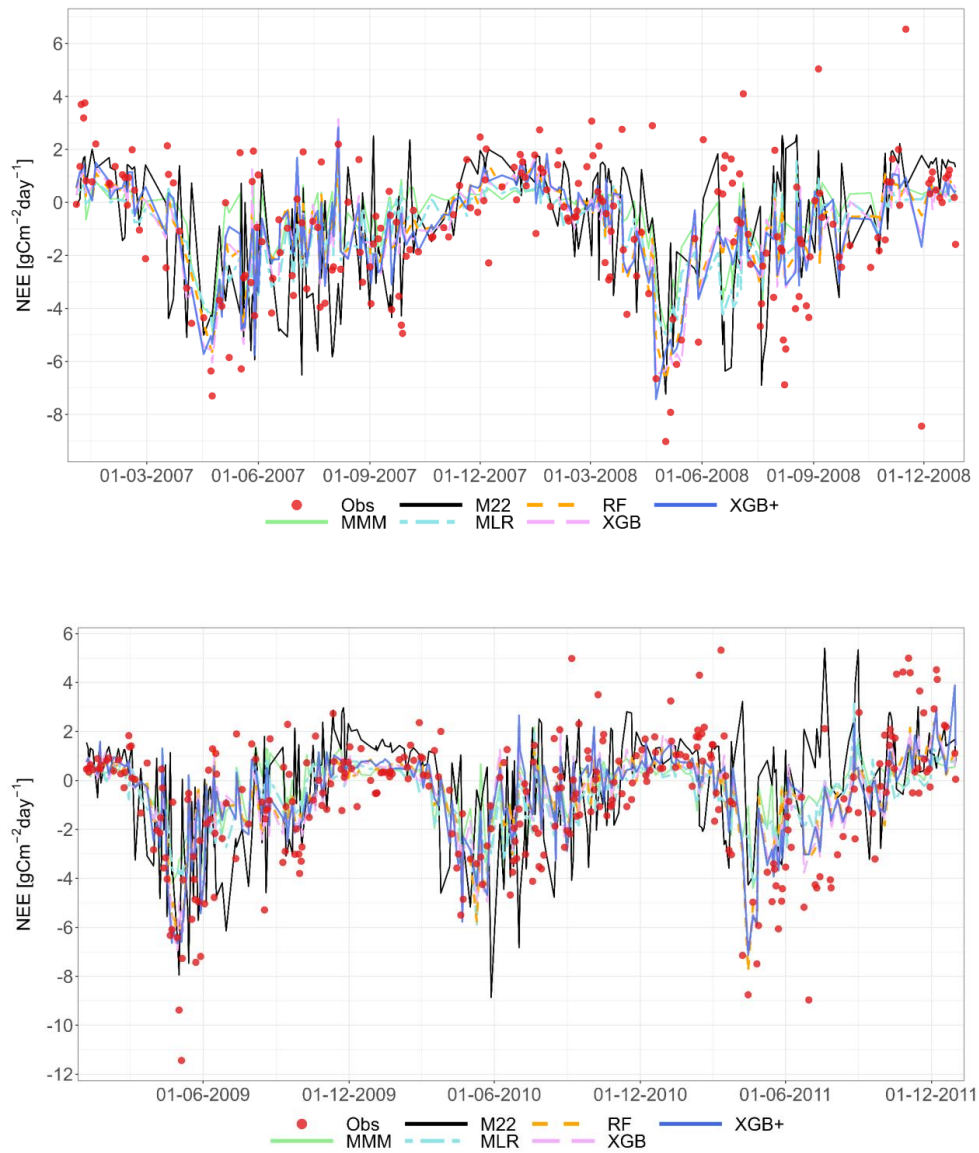
**Figure A9:** NEE, C1 - Ottawa (CA), cropland, 2007 and 2011.



**Figure A10:** NEE, C2 - Grignon, (FR), cropland, 2011 and 2012.



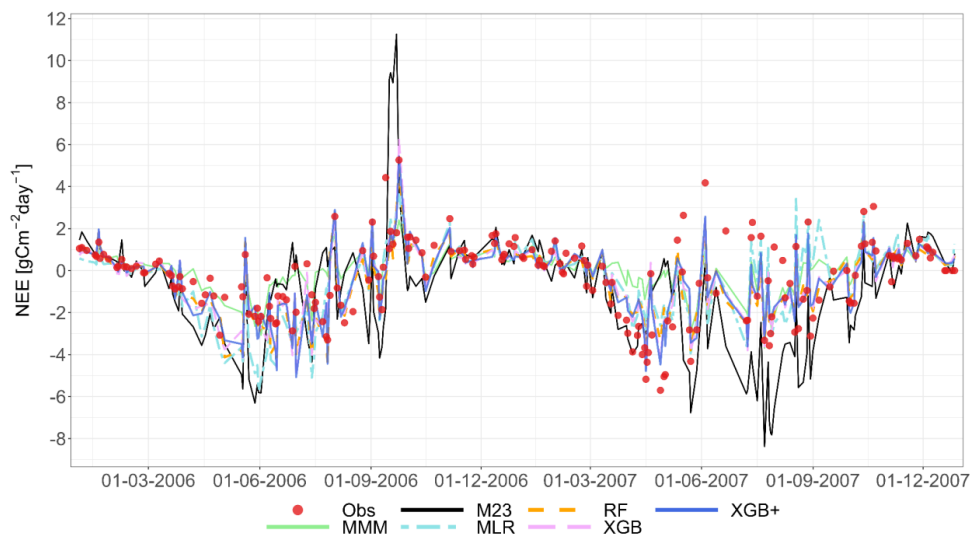
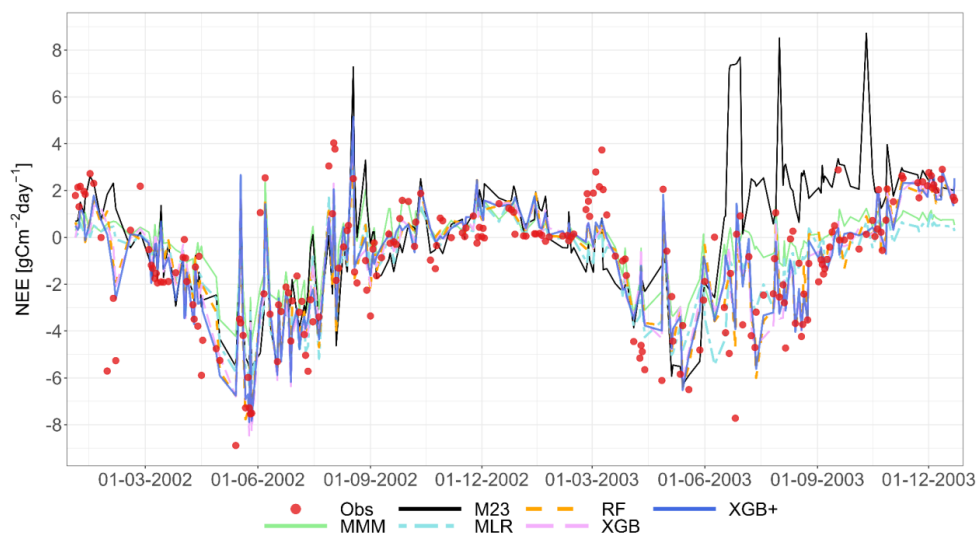
1315

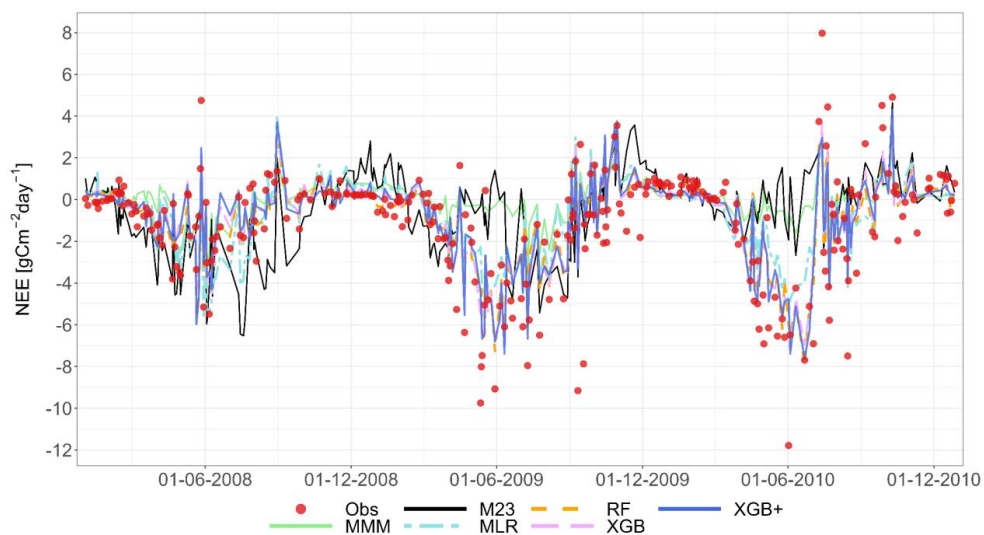


1320

**Figure A11:** NEE, G3 - Laqueuille (FR), grassland, 2003-2011.

1325





1335

**Figure A12:** NEE, G4 - Easter Bush (UK), grassland