

KGML-ag: A Modeling Framework of Knowledge-Guided Machine Learning to Simulate Agroecosystems: A Case Study of Estimating N₂O Emission using Data from Mesocosm Experiments

Licheng Liu¹, Shaoming Xu², Jinyun Tang³, Kaiyu Guan^{4,5,6}, Timothy J. Griffis⁷, Matthew D. Erickson⁷, Alexander L. Frie⁷, Xiaowei Jia⁸, Taegon Kim^{1,9}, Lee T. Miller⁷, Bin Peng^{4,5,6}, Shaowei Wu¹⁰, Yufeng Yang¹, Wang Zhou^{4,5}, Vipin Kumar², Zhenong Jin^{1*}

¹Department of Bioproducts and Biosystems Engineering, University of Minnesota, Saint Paul, MN, 55108, USA

²Department of Computer Science and Engineering, University of Minnesota, Minneapolis, MN, 55455, USA

³Climate and Ecosystem Sciences Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA

⁴Agroecosystem Sustainability Center, Institute for Sustainability, Energy, and Environment, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA

⁵Department of Natural Resources and Environmental Sciences, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA

⁶National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA

⁷Department of Soil, Water, and Climate, University of Minnesota, Saint Paul, MN 55108, USA

⁸Department of Computer Science, University of Pittsburgh, Pittsburgh, PA, 15260, USA

⁹Department of Smart Farm, Jeonbuk National University, Jeonju, Jeollabuk-do, 54896, Republic of Korea

¹⁰School of Physics and Astronomy, University of Minnesota, Minneapolis, MN, 55455, USA

Correspondence to: Zhenong Jin (jinzn@umn.edu)

Abstract.

Agricultural nitrous oxide (N₂O) emission accounts for a non-trivial fraction of global greenhouse gases (GHGs) budget. To date, estimating N₂O fluxes from cropland remains a challenging task because the related microbial processes (e.g., nitrification and denitrification) are controlled by complex interactions among climate, soil, plant and human activities. Existing approaches such as process-based (PB) models have well-known limitations due to insufficient representations of the processes or uncertainties of model parameters, and to leverage recent advances in machine learning (ML) a new method is needed to unlock the “black box” to overcome its limitations such as low interpretability, out-of-sample failure and massive data demand. In this study, we developed a first-of-kind knowledge-guided machine learning model for agroecosystems (KGML-ag), by incorporating biogeophysical/chemical domain knowledge from an advanced PB model, *ecosys*, and tested it by comparing simulating daily N₂O fluxes with real observed data from mesocosm experiments. The Gated Recurrent Unit (GRU) was used as the basis to build the model structure. To optimize the model performance, we have investigated a range of ideas, including: 1) Using initial values of intermediate variables (IMVs) instead of time series as model input to reduce data demand; 2) Building hierarchical structures to explicitly estimate IMVs for further N₂O prediction; 3) Using multitask learning to balance the simultaneous training on multiple variables; and 4) Pretraining with millions of synthetic data generated from *ecosys* and fine tuning with mesocosm observations. Six other pure ML models were developed using the same mesocosm data to serve as the benchmark for the KGML-ag model. Results show that KGML-ag did an excellent job in reproducing the mesocosm

37 N₂O fluxes (overall $r^2 = 0.81$, and RMSE = 3.6 mg N m⁻² day⁻¹ from cross-validation). Importantly KGML-ag always
38 outperforms the PB model and ML models in predicting N₂O fluxes, especially for complex temporal dynamics and emission
39 peaks. Besides, KGML-ag goes beyond the pure ML models by providing more interpretable predictions as well as pinpointing
40 desired new knowledge and data to further empower the current KGML-ag. We believe the KGML-ag development in this
41 study will stimulate a new body of research on interpretable ML for biogeochemistry and other related geoscience processes.

42 **1 Introduction**

43 Nitrous oxide (N₂O), with its global warming potential 273 ± 118 times greater than that of carbon dioxide (CO₂) for a 100-
44 year time horizon, is one of the major greenhouse gases (IPCC6; Forster et al., 2021). The increasing rate of atmospheric N₂O
45 concentration during the period 2010-2015 is 44% higher than during 2000-2005, mainly driven by increased anthropogenic
46 sources that have increased total global N₂O emissions to ~ 17 Tg N yr⁻¹ (Syakila and Kroeze, 2011; Thompson et al., 2019).
47 It is estimated that approximately 60% of the contemporary N₂O emission increases are from agriculture management at global
48 scale (Pachauri et al., 2014; Robertson et al., 2014; Tian et al., 2020), but the estimation uncertainty can exceed 300% (Barton
49 et al., 2015; Solazzo et al., 2021). Quantifying N₂O emissions from agricultural soils is extremely challenging, partly because
50 the related microbial processes, mainly about incomplete denitrification and nitrification, are controlled by many environment
51 and management factors such as temperature/water conditions, soil/crop properties, and N fertilization rate, all of which
52 together have collectively led to large temporal and spatial variabilities of N₂O emissions (Butterbach-Bahl et al., 2013; Grant
53 et al., 2016).

54
55 Process-based (PB) models are often used for simulating N₂O fluxes from agroecosystems, but they have some inherent
56 limitations, including incomplete knowledge of the processes, low accuracy due to the under-constrained parameters,
57 expensive computing cost, and rigid structure for further improvements, that we could not resolve by using PB model itself.
58 For example, an advanced agroecosystem model, *ecosys* (Grant et al., 2003, 2006, 2016), simulates N₂O production rates
59 through nitrification and denitrification processes when oxygen (O₂) is limited, with equations considering the influence from
60 related substrate concentrations (e.g., NO₂⁻, N₂O, and CO₂), nitrifier and denitrifier populations, and soil thermal, hydrological
61 physical and chemical conditions. The produced N₂O accumulates, transfers in gaseous phase, aqueous phase, over different
62 soil layers, and eventually exchanges with atmosphere at the soil surface. Other PB models, including DNDC (Zhang et al.,
63 2002; Zhang and Niu, 2016), DAYCENT (Del Grosso et al., 2000; Nécipálová et al., 2015), and APSIM (Keating et al., 2003;
64 Holzworth et al., 2014), have also included processes to simulate N₂O production, but adopt different parameterizations using
65 static partition parameters to estimate N₂O emission from nitrification, and other empirical parameters to control the influence
66 on nitrification from soil water content, pH, temperature and substrate concentrations. Besides, N₂O is intimately connected
67 with the soil organic carbon (SOC) dynamics, because soil nitrifiers and denitrifiers interact strongly with aerobic and
68 anaerobic heterotrophs that process SOC evolution, and all of these microbes are driven by shared environmental variables

69 including soil temperature, moisture, redox status, and physical and chemical properties (Thornley et al., 2007). As expected,
70 these connections make it difficult for PB models, even the most advanced ones like *ecosys*, to find sufficient representations
71 of the physical and biogeochemical processes or obtain enough data to calibrate a large number of model parameters with
72 strong spatio-temporal variations. Thus, novel approaches are needed for addressing the big challenge of agricultural N₂O flux
73 simulations.

74
75 Machine learning (ML) models can automatically learn patterns and relationships from data. Recent studies have investigated
76 the potential to predict agricultural N₂O emission with ML models, including random forest (RF, Saha et al., 2021),
77 metamodelling with extreme gradient boosting (XGBoost) (Kim et al., 2021), and deep learning neural network (DNN)
78 (Hamrani et al., 2020). Notably, Hamrani et al. (2020) compared nine widely used ML models for predicting agricultural N₂O.
79 That study pointed out that the long short term memory (LSTM) model with recurrent networks containing memory cells as
80 building blocks will be most suitable for N₂O predictions, but the challenge remains with respect to the ability of capturing the
81 sharp peak of N₂O fluxes and lag time between N fertilizer application and the emission peak. Although there is an increasing
82 interest in leveraging recent advances in machine learning, capturing this opportunity requires going beyond the ML
83 limitations, including limited generalizability to out-of-sample scenarios, demand for massive training data, and low
84 interpretability due to the “black-box” use of ML (Karpatne et al., 2017). PB models with their transparent structures built by
85 representations of physical and biogeochemical processes, seem to be exact complementary to ML models. Thus, combining
86 the power of ML model and PB model understanding innovatively is likely a path forward.

87
88 The above need to integrate ML and PB models can be potentially addressed by the newly proposed framework of Knowledge-
89 guided Machine Learning (KGML) models. In the review by Willard et al. (2021), five research frontiers have been identified
90 regarding the development of KGML for diverse disciplines including earth system science, they are: 1) Loss function design
91 according to physical or chemical laws (Jia et al., 2019, 2021; Read et al., 2019); 2) Knowledge-guided initialization through
92 pretraining ML models with synthetic data generated from PB models (Jia et al., 2019, 2021; Read et al., 2019); 3) Architecture
93 design according to causal relations or adding dense layers containing domain knowledge (Khandelwal et al., 2020; Beucler
94 et al., 2019, 2021); 4) Residual modeling with ML models to reduce the bias between PB model outputs and observations
95 (Hanson et al., 2020); and 5) Other hybrid modeling approaches combining PB and ML models (Kraft et al., 2021). These
96 recent advances in KGML pave the pathway to a more efficient, accurate and interpretable solution for estimating N₂O fluxes
97 from the agroecosystem.

98
99 In this study, we present a first-of-its-kind attempt of developing a KGML for agricultural GHG fluxes prediction (KGML-ag)
100 with knowledge-guided initialization and architecture design, and demonstrate the potential of KGML-ag with a case study on
101 quantifying N₂O flux observed by a multi-year mesocosm experiments. We designed the KGML-ag structure based on the
102 causal relations of related N₂O processes informed by an advanced agroecosystem model, *ecosys* (Grant et al., 2003, 2006,

103 2016). We used the synthetic data generated from *ecosys* to design the KGML-ag input/output, and to pre-train the KGML-ag
104 model to learn the basic patterns of each variable. Observations from multi-season controlled-environment mesocosm
105 chambers (Miller, 2021, thesis; Miller et al., 2021, in review) were used to refine the pretrained KGML-ag and evaluate the
106 model performance. Since there is limited literature that guides the development of KGML-ag and not a one that directly
107 addressed GHG fluxes, we investigated a range of ideas to optimize the model performance, including: 1) Using initial values
108 of intermediate variables (IMVs) instead of sequences as model input to reduce data demand; 2) Building hierarchical
109 structures to explicitly estimate IMVs for further N₂O prediction; 3) Using multitask learning to balance the simultaneous
110 training on multiple variables; and 4) Pretraining with millions of synthetic data generated from *ecosys* and fine tuning with
111 mesocosm observations. Although we evaluated the KGML-ag models with real measurements only from a mesocosm
112 experiment, the lessons learned from the development process and various KGML-ag structures can be transferred to other
113 data, other variables and large scale simulations, therefore have broader implications on further KGML related research in
114 agriculture. We believe this study will stimulate a new body of research on interpretable machine learning for biogeochemistry
115 and other related topics in geoscience.

116 **2 Methods**

117 **2.1 Experimental design overview**

118 To develop and evaluate the KGML-ag models and compare their performance with pure ML models, we designed the
119 following experiments:

- 120 1) With the synthetic data, we developed and pretrained multiple KGML-ag models to learn general patterns and
121 interactions among variables, and evaluated their model performance (Fig. S2, Table 1);
- 122 2) With the observed data, we finetuned multiple KGML-ag models to adapt real-world situations, and evaluated their
123 model performance (Fig. 2-3; Fig. S3-5; Table 2-3);
- 124 3) We further benchmarked KGML-ag models and uncertainties with other pure ML models without considering
125 temporal dependence, including Decision Tree (DT), Random Forest (RF), Gradient Boosting (GB) from the sklearn
126 package (<https://scikit-learn.org/stable/>), Extreme Gradient Boosting (XGB) from the XGBoost package
127 (<https://xgboost.readthedocs.io/en/latest/>) and a 6-linear-layer artificial neural network (ANN) with the mesocosm
128 experiment data by 10 times ensemble experiments (Fig. 4-5; Fig. S6-8);
- 129 4) We conducted a few small experiments to further investigate how various model configurations, such as the
130 pretraining process, data augmentation and IMV initial values would influence KGML-ag model performance (Table
131 3).

132 2.2 KGML-ag structure development

133 2.2.1 Generating synthetic data with *ecosys*

134 We generated synthetic data using a PB model, *ecosys*. The *ecosys* model is an advanced agroecosystem model constructed
135 from detailed biophysical and biogeochemical rules instead of using empirical relations (Grant et al., 2001). It represents N₂O
136 evolution in the microbe-engaged processes of nitrification-denitrification using substrate kinetics that are sensitive to soil
137 nitrogen availability, soil temperature, soil moisture, and soil oxygen status (Grant and Pattey 2008). Two groups of microbial
138 populations, autotrophic nitrifiers and heterotrophic denitrifiers, produce N₂O with specific competitive or cooperative
139 relations in *ecosys* when O₂ availability fails to meet O₂ demand for their respiration, and NO₂⁻ become alternative electron
140 acceptors. N₂O transfer within soil layers and from soil to the atmosphere is driven by concentration gradient using diffusion-
141 convection-dispersion equations, in the forms of gaseous and aqueous N₂O under control of volatilization-dissolution (Grant
142 et al., 2016). Unlike the pipeline model described by Davidson et al. (2000), which mainly considers the correlations of N₂O
143 production with nitrogen availability and of N₂O emissions with soil water content, *ecosys* enables integrative effects of energy,
144 water, nitrogen availability on N₂O production and N₂O transfer via the microbial population dynamics and their interactions
145 with soil, plant, and atmospheric dynamics, under diverse meteorological and anthropogenic disturbances (e.g. runoff,
146 drainage, tillage, irrigation, soil erosion). Many previous studies have demonstrated its robustness in simulating agricultural
147 carbon and nitrogen cyclings at different spatial/temporal scales, and under different management practices (Grant et al., 2003,
148 2006, 2016; Metivier et al., 2009; Zhou et al., 2021). For the agricultural ecosystems in the US Midwest, whose simulations
149 are used for synthetic data in this study, the performance of *ecosys* on CO₂ have been extensively benchmarked, including CO₂
150 exchange (daily Reco R² = 0.80-0.86; daily NEE, R² = 0.75-0.89) and leaf area index (LAI, R² = 0.78) from six flux towers,
151 USDA census reported corn yield (R² = 0.83) and soybean yield (R² = 0.80), satellite-derived GPP for corn (R² = 0.83) and
152 soybean (R² = 0.85) in the US Midwest (Zhou et al., 2021). In addition, *ecosys* model can capture the dynamics and magnitude
153 of N₂O flux in hourly frequency (R² = 0.2-0.4 and RMSE = 0.1-0.2 mg N m⁻² h⁻¹ in Grant et al., 2008; R² = 0.28-0.37 and
154 RMSE = 0.2-0.28 mg N m⁻² h⁻¹ in Grant et al., 2003), and in various ecosystems (e.g. agriculture soil in Grant et al., 2006,
155 2008; forest in Grant et al., 2010; and grassland in Grant et al., 2016). Therefore, *ecosys* is an appropriate choice of
156 domain knowledge provider and synthetic data generator in the development of KGML models. We generated daily synthetic
157 data including N₂O flux and 76 IMVs (e.g. CO₂ flux from soil, layerwise soil NO₃⁻ concentration, layerwise soil temperature,
158 and layerwise soil moisture; detailed in Table S1) from *ecosys* simulations for 2000-2018 over 99 randomly selected counties
159 in Iowa, Illinois, and Indiana, USA. We used hourly meteorological inputs (downward shortwave radiation, air temperature,
160 precipitation, relative humidity, and wind speed) from the phase 2 of North American Land Data Assimilation System
161 (NLDAS-2, Xia et al., 2012) and layerwise soil properties (e.g. bulk density, texture, pH, SOC concentration) from the
162 SSURGO database (Soil Survey Staff, 2020) as inputs to *ecosys*. Crop management except N fertilization rates were configured
163 to the same settings as mesocosm experiments (described in Sec 2.2.2). To increase the variability in synthetic data, we

164 implemented 20 different N fertilization rates ranging from 0 to 33.6 g N m⁻² (i.e. 0 to 300 lb N ac⁻¹) in each simulation of 99
165 counties, and more detailed information for model setup refers to Zhou et al. (2021).

166
167 The generated synthetic data were then processed for further use by KGML-ag development. Meanwhile, the hourly weather
168 forcings were converted to seven daily variables, including the maximum air temperature (TMAX_AIR, °C), difference
169 between the maximum and the minimum air temperature (TDIF_AIR, °C), the maximum humidity (HMAX_AIR, fraction),
170 difference between the maximum and the minimum humidity (HDIF_AIR, fraction), surface downward shortwave radiation
171 (RADN, W m⁻²), precipitation (PREC, mm day⁻¹), and wind speed (WIND, m s⁻¹). Six soil properties were retrieved from the
172 SSURGO database, including total averaged (depth weighted averaged for all layers) bulk density (TBKDS, Mg m⁻³), sand
173 content (TCSAND, g kg⁻¹), silt content (TCSILT, g kg⁻¹), pH (TPH), cation exchange capacity (TCEC, cmol⁺ kg⁻¹) and soil
174 organic carbon (TSOC, g C kg⁻¹); and two crop properties were retrieved, including planting day of the year (PDOY) and crop
175 type (CROPT, 1 for corn and 0 for soybean). Finally, each synthetic data sample has daily N₂O flux, 76 selected IMVs, 7
176 weather forcings (W), 1 N fertilization rate (FN, g N m⁻²) and 8 soil/crop properties (SCP) (Fig. 1.a; Table S1). The periods
177 from April 1st to July 31st (122 days) were selected to cover the mesocosm observations (around 30 days before and 90 days
178 after N fertilizer date). The total amount of synthetic data sample is 122 days x 18 years x 99 counties x 20 N fertilizer rates
179 (about 4.3 million data points). We randomly selected the samples from 70 counties for training, 10 counties for validation,
180 and 19 counties for testing.

181 **2.2.2 Mesocosm experiments for KGML-ag model fine-tuning and evaluation**

182 Observations were acquired from a controlled-environment mesocosm facility on the St. Paul campus of the University of
183 Minnesota. Soil samples were sourced in 2015 from a farm in Goodhue County, MN (44.2339° N and 92.8976° W), which had
184 been under corn-soybean rotation for 25 years. Six chambers with a soil surface area of 2 m² and column depth of 1.1 m were
185 used to plant continuous corn during 2015-2018 and monitor the N₂O flux response to different precipitation treatments. The
186 experiment also measured other environmental variables including air temperature and photosynthetically active radiation
187 (PAR), which were controlled to mimic the outdoor ambient environment. Granular urea fertilizer was hand broadcasted and
188 incorporated to a depth of 0.05 m to each chamber at a rate of 22.4 g N m⁻² (200 lb N ac⁻¹) on May 1st of 2015, May 4th of
189 2016 and May 3rd of 2017, and 10.3 g N m⁻² (92 lb N ac⁻¹) on May 8th of 2018. Corn hybrid (DKC-53-56RIB) were hand
190 planted to a depth of 0.05 m in two rows spaced 0.76 m apart 3-5 days after fertilizer application, at a seeding rate of 35,000
191 seeds ac⁻¹ in 2015 to 2017, and 70,000 seeds ac⁻¹ in 2018 but thinned upon emergence to ensure 100 percent emergence at
192 35,000 seeds ac⁻¹. Crops were harvested at the end of September by cutting the stover five inches above the soil. Hourly N₂O
193 fluxes (mg N m⁻² h⁻¹) and CO₂ fluxes (g C m⁻² h⁻¹) were measured using non-steady-state flux chambers with a CO₂ analyzer
194 (LI-10820 for 2016 and LI-7000 for 2017 and 2018, LI-COR Biosciences, Lincoln, NE) and a N₂O analyzer (Teledyne
195 M320EU, Teledyne Technologies International Corp, Thousand Oaks, CA) (Detail method can be retrieved from Fassbinder
196 et al., 2012, 2013). We also collected soil moisture at 15 cm depth (VWC as abbreviation of volumetric water content, m³ m⁻³)

197 ³), weekly 0-15 cm depth soil $\text{NO}_3^- + \text{NO}_2^-$ concentration (NO_3^- for short in the following text, g N Mg^{-1}), soil NH_4^+
198 concentration (NH_4^+ , g N Mg^{-1}), and related environment variables including air temperature, radiation, humidity and soil/crop
199 properties from three growing seasons during 2016-2018 and six mesocosm chambers (Fig. S1). The magnitude of N_2O flux
200 and NO_3^- soil concentration and their responses following fertilizer application from this mesocosm experiment are slightly
201 higher than several field studies of agricultural soils (Fassbinder et al., 2013; Grant et al., 1999, 2006, 2008, 2016; Hamrani et
202 al., 2020; Venterea et al., 2011). More details about the mesocosm facility and experimental design can be found in the thesis
203 of Miller L. (2021).

204

205 The observed data were then processed to fine-tune and evaluate the KGML-ag models. The N_2O flux and four IMVs and
206 weather variables were collected from the measurements in the selected period (i.e., April 1st to July 31st). Weekly NO_3^- (short
207 for soil NO_3^- within 0-15 cm depth), and NH_4^+ (short for soil NH_4^+ within 0-15 cm) were linearly interpolated to the daily time
208 scale on days containing VWC (short for soil VWC in 15 cm) data. Hourly air temperature, net radiation, N_2O (short for N_2O
209 fluxes from soil), CO_2 (short for CO_2 fluxes from soil) and VWC were resampled to daily scale. All SCP were derived from
210 mesocosm measurements except that TCEC was derived from the SSURGO database according to the soil origin. We used the
211 leave-one-out cross-validation (LOOCV) method for the evaluation process. Each time we used five chambers' data for model
212 finetuning and another one chamber data for validation. For example, if we used chamber 1-5 to train the model, then chamber
213 6 would serve as the out-of-sample data to validate the results. Only the validation results would be presented in our study.

214

215 To reduce overfitting and increase the generalization of the trained model based on the small amount of mesocosm data, we
216 applied the following method to augment the experimental measurements and weather forcings to 1000 times larger by
217 sampling hourly data and averaging them to daily scale. In this method, 16 hours (or maximum valid hours) of data are
218 randomly selected from 24 hours of data to compute their mean as the daily value. Since up to 2/3 of the day is covered by the
219 selected data (16 hours /24 hours), the augmented daily values should be representative enough for the source day with slight
220 variations from each other. Furthermore, the observation ratio, (24 hours - missing hours) / 24 hours, can be used as the weights
221 in loss function to inject the data quality information in model optimization. If the day has more than 16 hours missing values,
222 we consider the observations in that day as not trustworthy and drop the day by setting the weight to 0. This method can not
223 only augment the data to 1000 times larger but also deal with the missing values in observed data inherently. The total amount
224 of observed mesocosm data and related weather forcings are augmented to 122 days x 3 years x 6 chambers x 1000 data
225 samples in this study.

226

227 **2.2.3 Gated Recurrent Unit (GRU) as the basis of KGML-ag**

228 Hamrani et al. (2020) compared different models and reported that LSTM provided the highest accuracy in predicting N_2O
229 fluxes, because N_2O flux is time dependent by its production/consumption nature and LSTM simulates target variables by

230 considering both current and historical states. The LSTM model, proposed by Hochreiter and Schmidhuber (1997), uses a cell
231 state as an internal memory to preserve the historical information. At each time step, it creates a set of gating variables to filter
232 the input and historical information and then uses the processed data to update the cell state. Similar to LSTM, GRU is a gated
233 recurrent neural network but only keeps one hidden state (Cho et al., 2014). Though simpler than LSTM, GRU is proved to
234 have similar performance (Chung et al., 2014). Our preliminary test on synthetic data for N₂O prediction showed that GRU
235 indeed provided similar or higher accuracy and model efficiency under different model settings than LSTM (Table S2). This
236 is possible because simpler models with fewer weights and hyperparameters are more robust in combating the overfitting
237 problem. Therefore, we choose GRU as the basis of KGML-ag development.

238 **2.2.4 Incorporating domain knowledge to the development of KGML-ag**

239 To quantitatively reveal the correlations between N₂O fluxes and IMVs and guide the KGML-ag development, we conducted
240 feature importance analysis by a customized 4-layer GRU ML model (Fig. 1b). Each layer of the model has a GRU cell with
241 64 hidden units. The 4-layer structure makes the model deeper and capable of capturing complex interactions. Between each
242 GRU cell, 20% of the output hidden states are randomly dropped by replacing them with zero values (so called 20% dropout)
243 to avoid overfitting. A linear dense layer is used to map the final output to N₂O. We first trained GRU models using synthetic
244 data with different combinations of IMVs as inputs to predict the N₂O fluxes (original-test, Table S2). The feature importance
245 analysis of well-trained models was then implemented by replacing one input feature with a Gaussian noise with mean $\mu=0$
246 and standard deviation $\sigma=0.01$, while keeping others untouched (new-test). The importance score was calculated by the new-
247 test's root mean square error (RMSE) (replacing one feature) minus the original-test's RMSE (no replacing). RMSE was

248 calculated by $\sqrt{\frac{\sum_1^N (y_i - y_i')^2}{N}}$ where N is the total number of observations across time and space, y_i is i -th measurement from
249 synthetic data or observed data and y_i' is its corresponding prediction.

250
251 To find important variables for N₂O flux prediction in an ideal situation where all variables are available, we conducted a
252 feature importance analysis for GRU models with all IMVs and basic inputs including FN, 7 W and 8 SCP (Fig. S2a). Results
253 indicated that flux variables including NH₃, H₂, N₂, O₂, CH₄, evapotranspiration (ET) and CO₂ had significant influence on the
254 model performance. Variables ranked high in feature importance analysis are considered with priority during model
255 development. To develop a functionable KGML-ag, we further investigated the feature importance of four IMVs that are
256 available from mesocosm observations including CO₂, NO₃⁻, VWC and NH₄⁺, which were ranked 7th, 20th, 58th, 60th
257 respectively in 92 input features of synthetic data (Fig. S2a). We used these four available IMVs to create two input
258 combinations: 1) CO₂ flux, NO₃⁻, VWC and NH₄⁺ (IMVcb1), and 2) NO₃⁻, VWC and NH₄⁺ (IMVcb2). The objective of
259 building IMVcb2 was to investigate the importance of the highly ranked variable CO₂ flux (by removing it from the inputs),
260 and the impact of mixing-up flux and non-flux variables on model performance. We tested the feature importance of the GRU
261 models built with IMVcb1 and IMVcb2 to check whether they would help in N₂O prediction (Fig. S2b-c). All the feature

262 importance results above indicated the correlation intensity between N₂O and many other variables, which would help the
263 KGML-ag model development and interpretation in this study (rest of this section and Sec. 3.1), and would guide future N₂O
264 related measurements and KGML model development (discussed in Sec. 4.3).

265
266 Next we used the knowledge learned from synthetic data to develop the structure of KGML-ag (Fig. 1c-d). Previous studies
267 for KGML models have used physical laws, e.g., conservation of mass or energy, to design the loss function for constraining
268 the ML model to produce physically consistent results (Read et al., 2019; Khandelwal et al., 2020). However, for complex
269 systems like agroecosystems, it is challenging to incorporate physical laws, such as mass balance for N₂O, into the loss function
270 due to the incomplete understanding of the processes and the lack of mass balance related data for validation. An alternative
271 solution is to incorporate such information in the design of the neural network (Willard et al., 2021). Effectiveness of such an
272 approach was demonstrated by Khandelwal et al. (2020) in the context of modeling stream flow in a river basin using Soil &
273 Water Assessment Tool (SWAT). They used a hierarchical neural network to explicitly model IMVs (e.g., soil moisture, snow
274 cover) and their relationships with the target variable (streamflow) and showed that this model is much more effective than a
275 neural network that attempts to directly learn the relationship between input drivers and the target variables. Following this
276 idea, we identified four desired features of an effective KGML-ag model, including: 1) We used initial values instead of
277 sequence of the IMVs from synthetic data or observed data to provide a solid starting state for the ML system and reduce the
278 IMV data demand, and then used the rest of the data to further constrain the prediction of IMVs; 2) We built a hierarchical
279 structure based on the structure of process representation in *ecosys* to first predict IMVs and then simulate N₂O with predicted
280 IMVs; 3) We trained all variables together using multitask learning to reach the best prediction scores, which generalized the
281 model and incorporated interactions between IMVs and N₂O; 4) We initialized the KGML-ag model by pretraining with
282 synthetic data before using real observed data to transfer physical knowledge, which further reduced the demand on large
283 training samples and aided in faster convergence for fine-tuning.

284
285 To meet these desired features, we proposed two KGML-ag models (Fig. 1c-d). The first model, KGML-ag1, is a hierarchical
286 structure containing two modules to simulate IMVs and N₂O sequentially. Each module is a 2-layer 64 units GRU ML model.
287 The inputs to the module of the KGML-ag1 model for IMV predictions (KGML-ag1-IMV module) are FN, 7W and 8SCP
288 together with the initial values of IMVs, and the outputs are IMV predictions. The inputs to the module of the KGML-ag1
289 model for N₂O predictions (KGML-ag1-N₂O module) are FN, 7W, 8SCP and predicted IMVs from KGML-ag1-IMV, and the
290 output is the target variable N₂O. Linear dense layers were coded for both modules to map output states to IMVs or N₂O. The
291 dropout method was applied to drop 20% of the state output between GRU cells and dense layers. The second model, KGML-
292 ag2, is also a hierarchical structure similar to KGML-ag1, but has multiple KGML-ag2-IMV modules to explicitly simulate
293 IMVs by tuning them separately in the fine-tuning process (discussed in Sec. 2.2.5). Each KGML-ag2-IMV module in KGML-
294 ag2 is a 2-layer 64 units GRU cell with the inputs of FN+7W+8SCP and one IMV initial value, and the output of one IMV

295 prediction. The KGML-ag2-N₂O module collects the IMV predictions from KGML-ag2-IMV modules and predicts the N₂O
296 with inputs of FN+7W+8SCP and predicted IMVs.

297 **2.2.5 Strategies for pretraining and fine-tuning processes**

298 To increase the efficiency of the training process, we used the Z-normalization ($\frac{X-\mu}{\sigma}$, where X is the vector of a particular
299 variable over all the data samples in the data set; μ is the mean value of X ; σ is the standard deviation of X) method to normalize
300 each variable separately on synthetic data. Then the scaling factors (μ , σ) derived from *ecosys* synthetic data for each variable
301 were used to Z-normalize observed data into the same ranges as synthetic data. As mentioned in Sec. 2.2.1, the TDIF_AIR,
302 HDIF_AIR were used instead of absolute min temperature (TMIN_AIR) and humidity (HMIN_AIR). This is done because
303 TMIN_AIR and HMIN_AIR follow similar trends as TMAX_AIR and HMAX_AIR, making Z-normalization numerically
304 poorly defined. Using the difference between maximum and minimum can provide a clearer information of daily air
305 temperature/humidity variation.

306
307 During the pretraining process, we initialized the IMV of KGML-ag using the first day value of synthetic IMV time series.
308 Adam optimizer with a start learning rate of 0.0001 was used for the training process. The learning rate would decay by 0.5
309 times after every 600 training epochs. At each epoch, synthetic data samples were randomly shuffled before being input to the
310 model to predict N₂O (and IMVs if any). The mean square error (MSE) loss (calculation was equal to the square of RMSE) or
311 sum of MSE loss (if multitask learning) between predictions and *ecosys* synthetic observations were calculated to optimize the
312 weights of GRU cells. After the training process updated the model's weights, the validation process was performed to evaluate
313 the model performance based on untouched samples with RMSE and the square of Pearson correlation coefficient (r^2). r^2 was

314 calculated as $\frac{(\sum_i (y_i' - \underline{y_i}') (y_i - \underline{y_i}))^2}{\sum_i (y_i' - \underline{y_i}')^2 (y_i - \underline{y_i})^2}$, where y_i is the i -th measurement from synthetic data or observed data, y_i' is its
315 corresponding prediction, $\underline{y_i}$ is the mean of the measurement y in diagnosing space and $\underline{y_i}'$ is the mean of the predicted y' in
316 diagnosing space. If both validated r^2 and RMSE were better than the best values in previous epochs, the updated model in this
317 epoch would be saved. Normalized RMSE (NRMSE, calculated by RMSE/(max-min) of each variable observation) was
318 introduced to evaluate IMV predictions between variables with different value ranges.

319
320 During the fine-tuning process, we used estimated IMV initial values of 1.0 g C m⁻², 0.2 m³ m⁻³, 0.0 g N Mg⁻¹, and 20.0 g N
321 Mg⁻¹ for CO₂, VWC, NH₄⁺, and NO₃⁻ respectively, from starting day (April 1st) to the day before the first day of real
322 observations, as input to KGML-ag models. Then the first-day values of observed IMVs were input into KGML-ag during the
323 rest days of the period as IMV initial values. In addition, as described in Sec. 2.2.2, we used a data augmentation method to
324 augment the total amount of data 1000 times larger for the fine-tuning process. The purpose of this data augmentation method
325 was to increase the generalization of the fine-tuned model and to overcome the overfitting due to small sample size. The mask

326 matrix was elementarily multiplied to the output matrix to calculate the MSE, r^2 and RMSE only for days with observations.
327 The similar optimizer was used with an initial learning rate of 0.00005 and decay fraction of 0.5 per 200 epochs. Other
328 training/validation methods in each epoch were similar to the pretraining process. Specifically, in the KGML-ag1 model
329 finetuning process, we first froze the KGML-ag1-N₂O module and only trained the KGML-ag1-IMV module for IMVs. After
330 finishing the KGML-ag1-IMV module training, we froze the KGML-ag1-IMV module and trained the KGML-ag1-N₂O
331 module for N₂O. In the KGML-ag2 fine-tuning process, the similar freezing method was used but different KGML-ag2-IMV
332 modules were trained separately one by one.

333 2.3 Development environment description

334 We used the Pytorch 1.6.0 (<https://pytorch.org/get-started/previous-versions/>) and python 3.7.9
335 (<https://www.python.org/downloads/release/python-379/>) as the programming environment for the model development. In order
336 to use the GPU to speed-up the training process, we installed cudatoolkit 10.2.89 (<https://developer.nvidia.com/cuda-toolkit>).
337 A desktop with Nvidia 2080 super GPU was used for code development and testing. The Mangi cluster
338 (<https://www.msi.umn.edu/mangi>) from High Performance Computing of Minnesota Supercomputing Institute (HPC-MSI,
339 <https://www.msi.umn.edu/content/hpc>) with 2-way Nvidia Tesla V100 GPU was used in training processes which consumed
340 longer time and bigger memories.

341 3 Results

342 3.1 Pretraining experiments using synthetic data from *ecosys*

343 In the pretraining stage, the GRU model with 76 IMVs achieved the best performance in predicting N₂O fluxes ($r^2=0.98$, RMSE
344 =0.54 mg N m⁻² day⁻¹ and normalized RMSE (NRMSE) = 0.01) on the test set of synthetic data generated from *ecosys* (Table
345 1). The high performance was due to some flux IMVs such as NH₃, H₂, O₂, CO₂ and ET, which are highly correlated to N₂O
346 (Fig. S2a), were used as input to the model. The good performance of GRU with all IMVs indicates that ML models are able
347 to perfectly mimic *ecosys* when sufficient information about IMVs is available. The GRU model with only basic input of N
348 fertilizer rate, 7 weather forcings, and 8 soil/crop properties (FN+7W+8SCP) had the accuracy of $r^2=0.89$ and RMSE = 1.37
349 mg N m⁻² day⁻¹ (Table 1). The relatively low performance is likely because this model failed to capture several highly nonlinear
350 pathways that are employed by *ecosys* to predict N₂O (e.g., one influence pathway from precipitation to N₂O can be:
351 Precipitation → soil moisture → N components solubility/concentration → nitrification/denitrification rate/amount → soil
352 N₂O concentration → gas N₂O flux). When adding sequences of IMV combinations (i.e., IMVcb1 of CO₂ flux, NO₃⁻, NH₄⁺
353 and VWC, and IMVcb2 of NO₃⁻, NH₄⁺ and VWC), the GRU models performed slightly better than the GRU model using only
354 basic inputs, achieving r^2 of 0.92 and 0.90, respectively (Table 1). The KGML-ag1 with IMVcb1 and IMVcb2 initial values
355 provided better performance (both $r^2 = 0.90$) than GRU with basic input and comparable performance to the GRU with inputs

356 of IMVcb1 and IMVcb2 sequence. Besides, KGML-ag1 provided predicted IMVs of CO₂, NO₃⁻, NH₄⁺, and VWC with r² over
357 0.91, and NRMSE below 0.06 (Table 1). KGML-ag2 also provided comparable N₂O performance but relatively better IMVs
358 performance of r² over 0.92 and NRMSE below 0.05. Results indicated that KGML-ag models with IMV initial values as extra
359 input performed similar or better than pure ML models in synthetic data.

360 **3.2 KGML-ag evaluation using observed data from mesocosm**

361 After being fine-tuned with observed data, KGML-ag1 had N₂O prediction overall accuracy of r²=0.81 and RMSE=3.6 mg N
362 m⁻² day⁻¹, while non-pretrained GRU model provided r²=0.78 and RMSE=4.0 mg N m⁻² day⁻¹, and pretrained GRU model
363 provided r²=0.80 and RMSE=3.77 mg N m⁻² day⁻¹ (Table 3). The time series of N₂O predictions from KGML-ag1 and the non-
364 pretrained GRU model were further compared (Fig. 2), from which we found at least two advantages of using KGML-ag1 for
365 N₂O predictions: 1) For the region without observation data (normally before day 25), KGML-ag1 predicted stable N₂O fluxes
366 close to 0 mg N m⁻² day⁻¹ (which is close to the reality in the experiment setting) while GRU caused anomalous peaks of fluxes.
367 This is because KGML-ag1 has learned knowledge for the whole period from the pretraining process with *ecosys* model
368 generated synthetic data, but GRU model has no prior knowledge for the period without any data in observations; 2) Although
369 KGML-ag1 had a lower accuracy than GRU in some chambers, KGML-ag1 can better capture the temporal dynamics of N₂O
370 fluxes compare to GRU, especially when the fluxes are highly variable (e.g. Fig 2 chamber 2).

371
372 To validate KGML-ag1 robustness, we further investigated the KGML-ag1 and GRU model performance in different temporal
373 windows, shrinking from the whole period to the N₂O peak occurrence time (days 1-122, day 30-80, day 40-65 and day 45-60
374 for year 2016-2018), and performance in N₂O flux, first order gradient of N₂O (slope) and second order gradient of the N₂O
375 (curvature) (Table 2). Slope represents the speed of N₂O flux changes through time and curvature represents the acceleration.
376 Assessing prediction performance with these two metrics will reveal the model robustness on capture variable dynamics, which
377 is critical when predicting fast-change variables with hot moments (a short period of time with rare events like flux increasing
378 quickly) like N₂O. First of all, the overall r² and RMSE of KGML-ag1 for values, slope and curvature were always better than
379 GRU. In particular, KGML-ag1 captured the peak region (e.g., days 45-60) much better than GRU in both magnitude and
380 dynamics (Table 2, Fig 2). Even for chamber 2 and 5 in which KGML-ag1 made worse N₂O predictions than GRU (Δr^2 ranging
381 from -0.07 to -0.03), it better captured temporal dynamics than GRU in terms of slope (Δr^2 ranging from 0.08 to 0.16) and
382 curvature (Δr^2 from 0.11 to 0.23) (Table 2). For other chambers, KGML-ag1 outperformed GRU consistently. For chamber 1,
383 KGML-ag1 had worse N₂O predictions RMSE than GRU but the Δr^2 increased as the window shrinks to the peak emission
384 time (0.07 → 0.13). The slope and curvature for chamber 1 also indicated that KGML-ag1 captured the dynamics much
385 better than GRU. For chamber 3, KGML-ag1 predicted better N₂O but presented worse slope and curvature RMSE than
386 GRU (Table 2). However, when explicitly investigating the time series of N₂O flux, slope and curvature in each year, KGML-
387 ag1 outperformed GRU more significantly in 2017, the year with more complex temporal dynamics of N₂O fluxes, than in

388 2016 and 2018, especially for chamber 3 (Fig. 2; Fig. S3-4). This investigation supported that KGML-ag1 was more capable
389 for complex dynamics predictions.

390

391 Interestingly, the fine-tuned KGML-ag1 model predicted reasonable IMVs including CO₂, NO₃⁻, NH₄⁺, and VWC with overall
392 r² of 0.37, 0.39, 0.60, and 0.33 and NRMSE of 0.14, 0.21, 0.09 and 0.18, respectively (Table 3). The time series comparisons
393 between IMV predictions and observations further indicated that KGML-ag1 could reasonably capture both magnitude and
394 dynamics (Fig. 3). KGML-ag2 presented better IMVs predictions than KGML-ag1, with overall r² of CO₂, NO₃⁻, NH₄⁺, and
395 VWC increasing by 0.37, 0.17, 0.06 and 0.51, and NRMSE decreasing by 0.05, 0.03, 0.01 and 0.10, respectively, but a slightly
396 lower r² (decreasing 0.02) of N₂O (Table 3; Fig. S5). This indicated that explicitly simulating each IMV with separated KGML-
397 ag2-IMV modules did not benefit the N₂O flux prediction accuracy, likely due to increasing model complexity which resulted
398 in reduced stability and ignoring the IMV interactions. In addition, we also found all KGML-ag models would perform better
399 by using IMVcb1 (with CO₂) than using IMVcb2 (without CO₂) in real data tests, indicating feature importance analysis based
400 on synthetic data can be a reasonable substitute for analysis with the often limited real-world data.

401 3.3 KGML-ag comparing with other pure ML models

402 The results from eight different models showed that KGML-ag1 comparing with other pure ML models consistently provided
403 the lowest RMSE (3.59-3.94 mg N m⁻² day⁻¹, 1.14-1.23 mg N m⁻² day⁻², and 0.84-0.89 mg N m⁻² day⁻³) and highest r² (0.78-
404 0.81, 0.48-0.56, and 0.23-0.31) for N₂O fluxes, slope and curvature, respectively (Fig. 4). This indicated that KGML-ag1
405 outperformed other pure ML models in capturing both the magnitude and dynamics of N₂O flux. Meanwhile, we have
406 calculated the uncertainty of mesocosm measurement due to converting hourly data to daily data during 30-80 days by using
407 augmented values minus the mean of the augmented values with lower and upper limits being -10.2 and 10.4 mg N m⁻² day⁻¹,
408 respectively (standard deviation = 1.4 mg N m⁻² day⁻¹). KGML-ag1 during the same period has comparable uncertainties based
409 on ensemble simulations with lower and upper limits being -14.4 and 15.2 mg N m⁻² day⁻¹, respectively (calculated by ensemble
410 values minus the mean of ensemble values; standard deviation = 1.3 mg N m⁻² day⁻¹). KGML-ag2 presented slightly better
411 mean scores for N₂O flux predictions than KGML-ag1, but worse scores for slope and curvature and larger uncertainties. This
412 proved the hypothesis discussed in section 3.2 that KGML-ag2 didn't benefit the magnitude and dynamics predictions of N₂O
413 flux with its more complex structure and less connections between IMVs.

414

415 Within the tree-based models (DT, RF, GB and XGB), the simplest model DT provided the worst predictions for N₂O flux,
416 slope and curvature. The XGB model provided the highest N₂O flux accuracy with r² of 0.61-0.63 and RMSE of 5.07-5.17 mg
417 N m⁻² day⁻¹, while the GB model provided best slope and curvature predictions with r² of 0.38-0.40 and 0.23-0.26, and RMSE
418 of 1.34-1.37 mg N m⁻² day⁻² and 0.91-0.95 mg N m⁻² day⁻³, respectively. The highest N₂O flux accuracy and relatively low
419 slope and curvature accuracy of the XGB model implied that there is a trade-off between the abilities of capturing dynamics
420 and magnitude.

421

422 In the group of deep learning models including ANN, GRU and KGML-ag1, ANN provided the worst predictions. Even with
423 the better N₂O flux predictions than most tree-based models (except XGB), the slope and curvature predictions of ANN were
424 the worst among all eight models. This implied that the trade-off between accurately capturing N₂O dynamics to magnitude in
425 ANN was significant. But when considering the temporal dependence, deep learning model GRU and KGML-ag1
426 outperformed all other models in flux, slope and curvature predictions. This indicated that without considering temporal
427 dependence the improvement in N₂O flux prediction accuracy could be risky by causing the performance drop in capturing
428 dynamics.

429

430 The detailed model comparisons in each chamber are shown in Fig. 5 (N₂O flux) and Fig. S6-7 (N₂O slope and curvature),
431 where the results are found to follow the same pattern as described above. In addition, time series comparisons of chamber 3
432 and 4 in 2017 between different models are presented in Fig. S8 as two examples. For periods without any observed data, we
433 assumed that the good model predictions should be stable, consistent with the nearest period and close to the reality in the
434 experiment setting (e.g. no erratic peak and N₂O flux near 0 mg N m⁻² day⁻¹ before day 25). From these comparisons, we infer
435 that without considering temporal dependence and pretraining process, the tree-based model including DT, RF, GB and XGB
436 and deep learning model ANN predicted erratic peaks in almost every missing data point, while the GRU model was stable in
437 short missing period (1-2 days of missing data) and only presented poor performance in long missing period (before day 25).
438 This improvement by the GRU model may be attributed to the structure of GRU that naturally keeps the historical information
439 using hidden states, which enables GRU to consider the temporal dependence and make consistent predictions over time.

440 **3.4 Influence of pretraining process, data augmentation and using IMV initial values as input feature**

441 After we pretrained the GRU model with synthetic data, the overall r² of N₂O flux predictions in observed data increased by
442 0.02, 0.12 and 0.14, and RMSE decreased by 0.23 mg N m⁻² day⁻¹, 0.15 mg N m⁻² day⁻² and 0.02 mg N m⁻² day⁻³ for flux, slope
443 and curvature predictions, respectively, compared to non-pretrained GRU (No.1-6 in Table 3). The gap between the GRU
444 model with pretrain and KGML-ag1 in N₂O value prediction shows the improvement resulting from architecture change (r²
445 increases by 0.01 and RMSE decreases by 0.17 mg N m⁻² day⁻¹). Although pretrained GRU had higher slope and curvature
446 prediction accuracy than KGML-ag models, it still couldn't achieve the current N₂O value prediction accuracy of KGML-ag1.
447 Besides, the KGML-ag models had relatively shallow N₂O prediction modules (2-layer GRU KGML-ag-N₂O module of
448 KGML-ag models vs 4-layer GRU) but included modules for IMV predictions, which therefore increased the model
449 interpretability.

450

451 It's worth noting that prediction accuracy of all KGML-ag models dropped without augmenting the training dataset in the fine-
452 tuning process (No.7-10 in Table 3). Moreover, the maximum training epochs increased from 800 to 20000, which resulted in

453 overfitting on the small data set. This indicated that the data augmentation indeed helped the models become more
454 generalizable and gain better accuracy.

455

456 Experiments using zero initial values presented a significant drop in every variable's prediction accuracy (No.11-14 in Table
457 3). This indicated that the IMV initial values input into the KGML-ag-IMV modules of KGML-ag models influenced not only
458 the IMV prediction but also the N₂O prediction of the KGML-ag-N₂O module. This shows that there is useful information
459 transferred from IMVs in the KGML-ag-IMV module to the KGML-ag-N₂O module.

460 **4 Discussion**

461 In the previous section, we showed that KGML-ag models can outperform ML models, by invoking architectural constraints
462 and PB model synthetic data initialization. Compared to traditional PB models such as *ecosys*, KGML-ag models provide
463 computationally more accurate and efficient predictions (KGML-ag few seconds vs *ecosys* half hour), which is similar to
464 traditional ML surrogate models (Fig. S9). But KGML-ag goes beyond that by providing more interpretable predictions than
465 pure ML models.

466 **4.1 Interpretability of KGML-ag**

467 The proposed KGML-ag models incorporate causal relations among N₂O related variables/processes as shown in Fig. S10.
468 Managements, weather forcings and initial values of IMVs influence soil water, soil temperature and soil properties, which
469 influence the availability of O₂ and N as well as the microbe populations in soil, and further influence the nitrification and
470 denitrification rates. N₂O is produced during both nitrification and denitrification when soil O₂ concentration is limited. Our
471 KGML-ag follows this hierarchical structure by designing KGML-ag-IMV modules representing the soil processes for IMVs
472 predictions (Fig. 1c-d).

473

474 To better explain the time series predictions of N₂O flux (Fig. S1; Fig. 2-3), we separated the observations of each year into
475 three periods: leading period (before N₂O increasing), increasing period (increasing to the peak) and decreasing period (peak
476 decreasing to near zero). During the leading period, both NH₄⁺ and CO₂ were increasing immediately in the following few days
477 following urea N fertilizer application, indicating that urea was decomposing into NH₄⁺ and CO₂ in soil water. With
478 accumulating NH₄⁺ in soil, nitrification started producing NO₃⁻ and consuming O₂. N₂O didn't respond to the fertilizer
479 immediately due to enough O₂ in soil. Then when the soil became sufficiently hypoxic, N₂O fluxes entered an increasing
480 period with N₂O being produced by nitrification and denitrification processes. CO₂ fluxes were relatively low and NH₄⁺ kept
481 decreasing during this period. Finally, when soil NH₄⁺ was exhausted and NO₃⁻ started decreasing due to denitrification, N₂O
482 fluxes then entered the decreasing period. CO₂ flux was related to urea decomposition during the leading period, and was more
483 closely related to O₂ demand in other periods. The KGML-ag predictions of N₂O and IMV captured the three periods and

484 transition points, demonstrating the connections between those variables following the description as above (Fig. 3; Fig. S5).
485 Although KGML-ag1 obtained lower IMVs prediction accuracy compared to KGML-ag2, it captured the general trends and
486 was doing better for transitions, especially in NH_4^+ predictions. KGML-ag2 overfitted on the observations and ignored the
487 correlations between IMVs, which resulted in loss in pretrain knowledge, poorer performance in the leading period, and erratic
488 predictions in the period with missing observations (before day 25).

489 **4.2 Lessons for KGML-ag development**

490 The development of KGML-ag in our study is suitable to predict not only N_2O but also other variables, such as CO_2 , CH_4 and
491 ET, with complicated generation processes relying on the historical states. To develop a capable KGML model, we need to
492 carefully address three questions:

493
494 What kind of ML model is suitable for developing KGML? The answer could be determined by the dominant variation type
495 of the target variable in the data. If the dominant type is temporal variance, like flux variables in high temporal resolution (e.g.,
496 daily, or hourly), we should consider ML models with temporal dependency. RNN models such as GRU used in this study,
497 and CNN models such as casual CNN (Oord et al., 2016) can be good starting ML models. If the dominant type is spatial
498 variation, like variables in coarse temporal resolution (e.g., monthly or annually) but with high diversity due to soil property,
499 land cover and climate, we should consider ML models with the ability to deal with edges, hotpoints and categories, such as
500 CNN;

501
502 What physical/chemical constraints can be used to build KGML models? Although physical rules such as mass balance or
503 energy balance are conceptually straightforward and were proved capable of constraining KGML in predicting lake phosphorus
504 and temperature dynamics (Hanson et al., 2020; Read et al., 2019), they were excluded in this study according to our
505 preliminary analysis. The reason is that the mass balance equation of N in the agriculture ecosystem includes too many
506 unknown and unobservable components such as N_2 flux, NH_3 flux, N leaching, microbial N, plant N and soil/plant exchange,
507 which collectively introduce large uncertainties in balance equations and make them hard to be directly applied in the KGML-
508 ag framework. Other related physical (e.g., diffusion, solution) or chemical (e.g., nitrification, denitrification) processes cannot
509 be easily added into the KGML-ag structure as rules due to lack of understanding of the process. Instead, as mentioned in Sect.
510 2.2.4, we used hierarchical structure to enforce an architectural constraint and causal relations among variables, and pretraining
511 processes to infuse knowledge from *ecosys* to KGML-ag models.

512
513 How to involve PB models in the KGML development? An advanced PB model like *ecosys* built upon biophysical and
514 biochemical rules instead of empirical relations will be a good basis to learn the process, guide the structure and provide the
515 constraints for KGML. The generated synthetic data in this study helped us get some knowledge about variables such as their
516 general trends, dynamics and correlations. Such knowledge can be transferred to KGML models from synthetic data in the

517 pretraining process, which can reduce the efforts to collect large numbers of real-world observation data. Moreover, while
518 KGML shows great potential beyond PB models, we reckon that equally important for improving N₂O modeling is to continue
519 improving our understanding of the related processes and mechanisms. Novel data collection and incorporating new
520 understanding into PB models (e.g., *ecosys*) could provide foundation to further empower KGML (see further discussion in
521 Sect. 4.3).

522

523 **4.3 Limitation and possible improvement**

524 First, the KGML-ag models in this study are limited by the available observed data. The mesocosm measurements of N₂O
525 fluxes (16.9±11.7 mg N m⁻² day⁻¹ during days of 45-60; Highest value is 71 mg N m⁻² day⁻¹) and NO₃⁻ soil concentrations
526 (59.3±20.7 g N Mg⁻¹ during days of 45-60; Highest value is 95.2 g N Mg⁻¹) are at the high end of the range that has been
527 observed by field studies (Fassbinder et al., 2013; Grant et al., 1999, 2006, 2008, 2016; Hamrani et al., 2020; Venterea et al.,
528 2011). Some IMVs with high feature importance scores (e.g., O₂ flux, N₂ flux) or at different depths (e.g., soil NO₃⁻ at 5 cm
529 depth, VWC at 5 cm depth), and data out of growing seasons are not included. The direct consequences are that some important
530 processes cannot be well represented by the current KGML-ag (e.g., O₂ demand and N availability for nitrification and
531 denitrification). Further improvement of KGML should consider three categories of data: target variable N₂O flux, IMVs and
532 basic inputs (Fig. 1a). For N₂O flux observation, we lack sub-hourly to sub-daily observations to capture the hot moment of
533 emission during 0-30 days after N fertilizer applications. Besides, the non-growing season can provide 35-65% of the annual
534 direct N₂O emissions from seasonally frozen croplands and lead to a 17–28 % underestimate of the global agricultural N₂O
535 budget if ignoring its contribution (Wagner-Riddle et al., 2017), but we can barely find observations from non-growing
536 seasons. For IMVs, we found oxygen demand indicator (e.g., O₂ concentration or flux, CO₂ flux, CH₄ flux), N mass balance
537 related variables (e.g., N₂ flux, soil NO₃⁻, soil NH₄⁺, N leaching) and soil water and temperature, can be used to better constrain
538 the processes and therefore improve the KGML performance. Rohe et al. (2021) also indicated the importance of O₂, CO₂ and
539 N₂ soil fluxes for N₂O predictions. In addition, the layerwise soil observations (e.g., soil NO₃⁻, soil VWC) at 0-30 cm depth
540 can be used to significantly improve the KGML model quality, according to our feature importance analysis (Fig. S2a).
541 Moreover, continuous monitoring on these variables during the whole year is preferred rather than only during the growing
542 season, since N₂O flux is largely influenced by previous states. To apply the KGML-ag to large scale, other observational data
543 including basic inputs of soil/crop properties (e.g., soil bulk density, pH, crop type), management information (e.g., fertilizer,
544 irrigation, tillage) and weather forcings along with N₂O flux observations are critical for fine-tuning and validating the
545 developed KGML-ag and therefore explicitly simulating the N₂O or IMVs dynamics under specific conditions. Recent
546 advances in remote sensing and machine learning have enabled estimating these variables with high-resolution at a large scale
547 (Peng et al., 2020)

548

549 Second, the physical/chemical constraints can be more comprehensive in KGML-ag models. Although current KGML-ag
550 models are well-initialized with *ecosys* synthetic data and constrained by causal relations of processes with hierarchical
551 structure, the predicted N₂O flux and IMVs can still violate some basic physical rules like mass balance. As we discussed in
552 Sec. 4.2, it will be challenging to add physical rules like mass balance equation for N in a complicated agriculture ecosystem
553 due to data limitations such as missing observations on certain key variables. Using inequalities instead of equations for mass
554 balance may be one alternative solution. For example, we could use ReLU to add in a limitation for N mass balance residues
555 which are calculated from known terms not larger than an empirical static value. Besides, better understanding of processes in
556 the N cycle from fieldworks and lab experiments can also help us design new constraints. This limitation is also partially
557 related to the data limitation and can be overcome by involving more complete N₂O data to introduce more powerful
558 constraints to KGML-ag.

559

560 Third, the KGML-ag currently are suffering from dealing with physical/chemical boundary transitions. Boundary transitions
561 are common in the real world, such as phase change, volume solubility, and soil porosity etc. A detailed PB model generally
562 coded plenty of “if/else/switch” statements inside to deal with the boundaries. But KGML-ag models based on the GRU are
563 better at capturing continuous changes, rather than discrete changes. One solution is to include data with boundary information.
564 In this study, involving IMVs like O₂, CO₂ and N₂, which already have boundary information like water freezing point, N pool
565 volumes and other complicated boundaries related to soil/crop properties, can significantly improve the model performance.
566 The data with boundary information could be continuous observation or estimated value from existing data. By using initial
567 values to predict IMVs, KGML-ag in this study can partially solve the boundary transition problem when observation data is
568 limited. Another solution is designing new structures of KGML-ag, such as combining ReLU function or including CNN
569 model which are robust for discrete situations to the RNN models, or designing new constraints to limit the model working
570 within the thresholds.

571

572 Finally, at the current stage we can not claim to have completely opened the black box of KGML-ag, but this framework is a
573 significant step towards this goal. For example, some ideas implemented in our study, such as using pretraining to transfer
574 knowledge from a PB model to a ML model, incorporating causal relations by hierarchical structure, predicting IMVs for
575 tracking middle changes and using initial values as input to reduce data demand, would shed light on the future KGML-ag
576 framework improvement. Besides, we acknowledge the importance of further testing the KGML-ag over completely
577 independent datasets, but results presented in this manuscript are sufficient to justify the power of KGML as a framework. The
578 mesocosm experiment data we used in this study has provided a comprehensive set of inputs and intermediate variables in
579 addition to the output of N₂O fluxes, thus serving as a unique testbed. We expect to further validate and refine our KGML-ag
580 model once more gold standard data of N₂O fluxes along with other relevant inputs and intermediate variables become publicly
581 available. Moreover, incorporating more and more domain knowledge into KGML-ag will be possible for further
582 improvement, but we don't think KGML-ag will become inefficient as it becomes more like the PB model. In fact, to efficiently

583 emulate components of PB models has been proposed as a research frontier in hybrid modeling for earth system science
584 (Reichstein et al., 2019; Irrgang et al., 2021), with latest advances occurring in weather forecasts (Bauer et al., 2021). By using
585 a hybrid model, computationally inefficient components of PB can be identified one by one, and be replaced with more efficient
586 ML-based surrogates to eventually obtain the most efficient model. Further KGML-ag model development will also need to
587 balance efficiency, accuracy and interpretability.

588 **5 Conclusions**

589 In this study, two KGML-ag models have been developed, validated, and tested for agricultural soil N₂O flux prediction using
590 synthetic data generated by the PB model *ecosys* and observational data from a mesocosm facility. The results show that
591 KGML-ag models can outperform PB and pure ML models in N₂O prediction in not only magnitude (KGML-ag1 $r^2 = 0.81$ vs
592 best ML model GRU $r^2 = 0.78$) but also dynamics (KGML-ag1 accuracy minus GRU accuracy, slope $\Delta r^2 = 0.06$ and curvature
593 $\Delta r^2 = 0.08$). KGML-ag can also defeat the PB model *ecosys* in efficiency by completing *ecosys*'s half-hour job within a few
594 seconds. Compared to ML models, KGML-ag models can better represent complex dynamics and high peaks of N₂O flux.
595 Moreover, with IMV predictions and hierarchical structures, KGML-ag models can provide biogeophysical/chemical
596 information about key processes controlling N₂O fluxes, which will be useful for interpretable forecasting and developing
597 mitigation strategies. Data demand for the KGML-ag models is significantly reduced due to involving IMV initial values and
598 pretrain processes with synthetic data. This study demonstrated that the potential of KGML-ag application in the complex
599 agriculture ecosystem is high and illustrates possible pathways of KGML-ag development for similar tasks. Further
600 improvement of our KGML-ag models can involve general principles to further constrain the predictions through loss functions
601 or architectures, but call for more detailed, high temporal resolution N₂O observation data from field measurements.

602 **Code and Data Availability**

603 The code and data used in this study can be found at <https://doi.org/10.5281/zenodo.5504533>.

604 **Author contributions**

605 LL and ZJ conceived the study. WZ and YY conducted *ecosys* simulations and provided synthetic data. LL and SX processed
606 the data and developed the KGML-ag model. LL, SX and SW carried the experiments out with supervisions from ZJ, JT, KG,
607 and VK. TJG, MDE, ALF and LTM shared mesocosm observations and interpreted the data. LL wrote the first draft of the
608 manuscript with further editing from TK on figures and tables. ZJ, SX, JT, KG, XJ, BP, YY, WZ and VK further edited the
609 manuscript.

610 **Competing interests**

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

612 **References**

613 Barton, L., Wolf, B., Rowlings, D., Scheer, C., Kiese, R., Grace, P., ... & Butterbach-Bahl, K.: Sampling frequency affects
614 estimates of annual nitrous oxide fluxes, *Scientific reports*, 5(1), 1-9, 2015.

615 Bauer, P., Dueben, P. D., Hoefler, T., Quintino, T., Schulthess, T. C., & Wedi, N. P.: The digital revolution of Earth-system
616 science. *Nature Computational Science*, 1(2), 104-113, 2021.

617 Beucler, T., Pritchard, M., Rasp, S., Ott, J., Baldi, P., & Gentine, P.: Enforcing analytic constraints in neural networks
618 emulating physical systems, *Physical Review Letters*, 126(9), 098302, 2021.

619 Beucler, T., Rasp, S., Pritchard, M., & Gentine, P.: Achieving conservation of energy in neural network emulators for climate
620 modeling, arXiv preprint arXiv:1906.06622, 2019.

621 Butterbach-Bahl, K., Baggs, E. M., Dannenmann, M., Kiese, R., & Zechmeister-Boltenstern, S.: Nitrous oxide emissions from
622 soils: how well do we understand the processes and their controls? *Philosophical Transactions of the Royal Society B:
623 Biological Sciences*, 368(1621), 20130122, 2013.

624 Cho, K., Van Merriënboer, B., Bahdanau, D., & Bengio, Y.: On the properties of neural machine translation: Encoder-decoder
625 approaches, arXiv preprint arXiv:1409.1259, 2014.

626 Chung, Junyoung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio.: Empirical evaluation of gated recurrent neural
627 networks on sequence modeling, arXiv preprint arXiv:1412.3555, 2014.

628 Daw, A., Thomas, R. Q., Carey, C. C., Read, J. S., Appling, A. P., & Karpatne, A.: Physics-guided architecture (pga) of neural
629 networks for quantifying uncertainty in lake temperature modeling. In *Proceedings of the 2020 siam international conference
630 on data mining* (pp. 532-540), Society for Industrial and Applied Mathematics, 2020.

631 Del Grosso, S. J., Parton, W. J., Mosier, A. R., Ojima, D. S., Kulmala, A. E., & Phongpan, S.: General model for N₂O and N₂
632 gas emissions from soils due to denitrification, *Global biogeochemical cycles*, 14(4), 1045-1060, 2020.

633 Fassbinder, J. J., Schultz, N. M., Baker, J. M., & Griffis, T. J.: Automated, Low-Power Chamber System for Measuring Nitrous
634 Oxide Emissions, *Journal of environmental quality*, 42, 606. doi: 10.2134/jeq2012.0283, 2013.

635 Fassbinder, J. J., Griffis, T. J., & Baker, J. M.: Evaluation of carbon isotope flux partitioning theory under simplified and
636 controlled environmental conditions, *Agricultural and forest meteorology*, 153, 154-164, 2012.

637 Forster, P., Storelvmo, T., Armour, K. , Collins, W., ... & Zhang, H.: The Earth's Energy Budget, Climate Feedbacks, and
638 Climate Sensitivity. In: *Climate Change 2021: The Physical Science Basis. Contribution of Working Group I to the Sixth
639 Assessment Report of the Intergovernmental Panel on Climate Change*, Cambridge University Press. In Press, 2021.

640 Gilhespy, S. L., Anthony, S., Cardenas, L., Chadwick, D., del Prado, A., Li, C., ... & Yeluripati, J. B.: First 20 years of DNDC
641 (DeNitrification DeComposition): model evolution, *Ecological modelling*, 292, 51-62, 2014.

642 Grant, R. F.: Modeling Carbon and Nitrogen Dynamics for Soil Management, (Boca Raton, FL: CRC Press) A review of the
643 Canadian ecosystem model ecosys 173–264, 2021.

644 Grant, R. F., & Pattey, E.: Mathematical modeling of nitrous oxide emissions from an agricultural field during spring thaw.
645 Global Biogeochemical Cycles, 13(2), 679-694, 1999.

646 Grant, R. F., & Pattey, E.: Modelling variability in N₂O emissions from fertilized agricultural fields, Soil Biology and
647 Biochemistry, 35(2), 225-243, 2003.

648 Grant, R. F., & Pattey, E.: Temperature sensitivity of N₂O emissions from fertilized agricultural soils: Mathematical modeling
649 in ecosys. Global biogeochemical cycles, 22(4), 2008.

650 Grant, R. F., Neftel, A., & Calanca, P.: Ecological controls on N₂O emission in surface litter and near-surface soil of a managed
651 grassland: modelling and measurements, Biogeosciences, 13(12), 3549-3571, 2016.

652 Grant, R. F., Pattey, E., Goddard, T. W., Kryzanowski, L. M., & Puurveen, H.: Modeling the effects of fertilizer application
653 rate on nitrous oxide emissions, Soil Science Society of America Journal, 70(1), 235-248, 2006.

654 Hamrani, A., Akbarzadeh, A., & Madramootoo, C. A.: Machine learning for predicting greenhouse gas emissions from
655 agricultural soils, Science of The Total Environment, 741, 140338, 2020.

656 Hanson, P. C., Stillman, A. B., Jia, X., Karpatne, A., Dugan, H. A., Carey, C. C., ... & Kumar, V.: Predicting lake surface
657 water phosphorus dynamics using process-guided machine learning, Ecological Modelling, 430, 109136, 2020.

658 Holzworth, D. P., Huth, N. I., deVoil, P. G., Zurcher, E. J., Herrmann, N. I., McLean, G., ... & Keating, B. A.: APSIM–
659 evolution towards a new generation of agricultural systems simulation, Environmental Modelling & Software, 62, 327-350,
660 2014.

661 Irrgang, C., Boers, N., Sonnewald, M., Barnes, E. A., Kadow, C., Staneva, J., & Saynisch-Wagner, J.: Towards neural Earth
662 system modelling by integrating artificial intelligence in Earth system science. Nature Machine Intelligence, 3(8), 667-674,
663 2021.

664 Jia, X., Willard, J., Karpatne, A., Read, J. S., Zwart, J. A., Steinbach, M., & Kumar, V.: Physics-guided machine learning for
665 scientific discovery: An application in simulating lake temperature profiles, ACM/IMS Transactions on Data Science, 2(3), 1-
666 26, 2021.

667 Jia, X., Willard, J., Karpatne, A., Read, J., Zwart, J., Steinbach, M., & Kumar, V.: Physics guided RNNs for modeling
668 dynamical systems: A case study in simulating lake temperature profiles, In Proceedings of the 2019 SIAM International
669 Conference on Data Mining (pp. 558-566), Society for Industrial and Applied Mathematics, 2019.

670 Karpatne, A., Atluri, G., Faghmous, J. H., Steinbach, M., Banerjee, A., Ganguly, A., ... & Kumar, V.: Theory-guided data
671 science: A new paradigm for scientific discovery from data, IEEE Transactions on knowledge and data engineering, 29(10),
672 2318-2331, 2017.

673 Keating, B. A., Carberry, P. S., Hammer, G. L., Probert, M. E., Robertson, M. J., Holzworth, D., ... & Smith, C. J.: An overview
674 of APSIM, a model designed for farming systems simulation, European journal of agronomy, 18(3-4), 267-288, 2003.

675 Khandelwal, A., Xu, S., Li, X., Jia, X., Stienbach, M., Duffy, C., ... & Kumar, V., Physics guided machine learning methods
676 for hydrology, arXiv preprint arXiv:2012.02854, 2020.

677 Kim, T., Jin, Z., Smith, T., Liu, L., Yang, Y., Yang, Y., ... & Zhou, W.: Quantifying nitrogen loss hotspots and mitigation
678 potential for individual fields in the US Corn Belt with a metamodeling approach, *Environmental Research Letters*, 2021.

679 Kraft, B., Jung, M., Körner, M., Koirala, S., & Reichstein, M.: Towards hybrid modeling of the global hydrological cycle,
680 *Hydrology and Earth System Sciences Discussions*, 1-40, 2021.

681 Meyer, D., Nagler, T., & Hogan, R. J.: Copula-based synthetic data augmentation for machine-learning emulators.
682 *Geoscientific Model Development*, 14(8), 5205-5215, 2021.

683 Miller, L. T. , Griffis, T. J., Erickson, M. D., Turner, P. A., Deventer, M. J., Chen, Z., Yu, Z., Venterea, R.T., Baker, J. M.,
684 and Frie, A. L.: Response of nitrous oxide emissions to future changes in precipitation and individual rain events, *Journal of*
685 *Environmental Quality*, In review, 2021

686 Miller, L. T., *Assessing Agricultural Nitrous Oxide Emissions and Hot Moments Using Mesocosm Simulations*, (Master
687 Thesis, University of Minnesota) Retrieved from the University of Minnesota Digital Conservancy,
688 <https://hdl.handle.net/11299/219276>, 2021

689 Necpálová, M., Anex, R. P., Fienen, M. N., Del Grosso, S. J., Castellano, M. J., Sawyer, J. E., ... & Barker, D. W.:
690 Understanding the DayCent model: Calibration, sensitivity, and identifiability through inverse modeling, *Environmental*
691 *Modelling & Software*, 66, 110-130, 2015.

692 Oord, A. V. D., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., ... & Kavukcuoglu, K.: Wavenet: A generative
693 model for raw audio, arXiv preprint arXiv:1609.03499, 2016.

694 Pachauri, R. K., Allen, M. R., Barros, V. R., Broome, J., Cramer, W., Christ, R., ... & van Ypserle, J. P.: *Climate change 2014:*
695 *synthesis report. Contribution of Working Groups I, II and III to the fifth assessment report of the Intergovernmental Panel on*
696 *Climate Change* (p. 151). *Ippc*, 2014.

697 Read, J. S., Jia, X., Willard, J., Appling, A. P., Zwart, J. A., Oliver, S. K., ... & Kumar, V.: Process-guided deep learning
698 predictions of lake water temperature, *Water Resources Research*, 55(11), 9173-9190, 2019.

699 Peng, B., Guan, K., Tang, J., Ainsworth, E. A., Asseng, S., Bernacchi, C. J., ... & Zhou, W.: Towards a multiscale crop
700 modelling framework for climate change adaptation assessment, *Nature plants*, 6(4), 338-348, 2020.

701 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., & Carvalhais, N.: Deep learning and process understanding
702 for data-driven Earth system science. *Nature*, 566(7743), 195-204, 2019.

703 Robertson, M., BenDor, T. K., Lave, R., Riggsbee, A., Ruhl, J. B., & Doyle, M.: Stacking ecosystem services, *Frontiers in*
704 *Ecology and the Environment*, 12(3), 186-193, 2014.

705 Rohe, L., Apelt, B., Vogel, H. J., Well, R., Wu, G. M., & Schlüter, S.: Denitrification in soil as a function of oxygen availability
706 at the microscale, *Biogeosciences*, 18(3), 1185-1201, 2021.

707 Saha, D., Basso, B., & Robertson, G. P.: Machine learning improves predictions of agricultural nitrous oxide (N₂O) emissions
708 from intensively managed cropping systems, *Environmental Research Letters*, 16(2), 024004, 2021.

709 Solazzo, E., Crippa, M., Guizzardi, D., Muntean, M., Choulga, M., & Janssens-Maenhout, G.: Uncertainties in the Emissions
710 Database for Global Atmospheric Research (EDGAR) emission inventory of greenhouse gases, *Atmospheric Chemistry and*
711 *Physics*, 21(7), 5655-5683, 2021.

712 Solazzo, E., Crippa, M., Guizzardi, D., Muntean, M., Choulga, M., & Janssens-Maenhout, G.: Uncertainties in the Emissions
713 Database for Global Atmospheric Research (EDGAR) emission inventory of greenhouse gases, *Atmospheric Chemistry and*
714 *Physics*, 21(7), 5655-5683, 2021.

715 Syakila, A., & Kroeze, C.: The global nitrous oxide budget revisited, *Greenhouse gas measurement and management*, 1(1),
716 17-26, 2011.

717 Thompson, R. L., Lassaletta, L., Patra, P. K., Wilson, C., Wells, K. C., Gressent, A., ... & Canadell, J. G.: Acceleration of
718 global N₂O emissions seen from two decades of atmospheric inversion, *Nature Climate Change*, 9(12), 993-998, 2019.

719 Thornley, J. H., & France, J.: *Mathematical models in agriculture: quantitative methods for the plant, animal and ecological*
720 *sciences*, Cabi, 2007.

721 Tian, H., Xu, R., Canadell, J. G., Thompson, R. L., Winiwarter, W., Suntharalingam, P., ... & Yao, Y.: A comprehensive
722 quantification of global nitrous oxide sources and sinks, *Nature*, 586(7828), 248-256, 2020.

723 Venterea, R. T., Maharjan, B., & Dolan, M. S.: Fertilizer source and tillage effects on yield-scaled nitrous oxide emissions in
724 a corn cropping system. *Journal of Environmental Quality*, 40(5), 1521-1531, 2011.

725 Wagner-Riddle, C., Congreves, K. A., Abalos, D., Berg, A. A., Brown, S. E., Ambadan, J. T., ... & Tenuta, M.: Globally
726 important nitrous oxide emissions from croplands induced by freeze–thaw cycles, *Nature Geoscience*, 10(4), 279-283, 2017.

727 Willard, J., Jia, X., Xu, S., Steinbach, M., & Kumar, V.: Integrating Scientific Knowledge with Machine Learning for
728 Engineering and Environmental Systems, arXiv preprint arXiv:2003.04919, 2020.

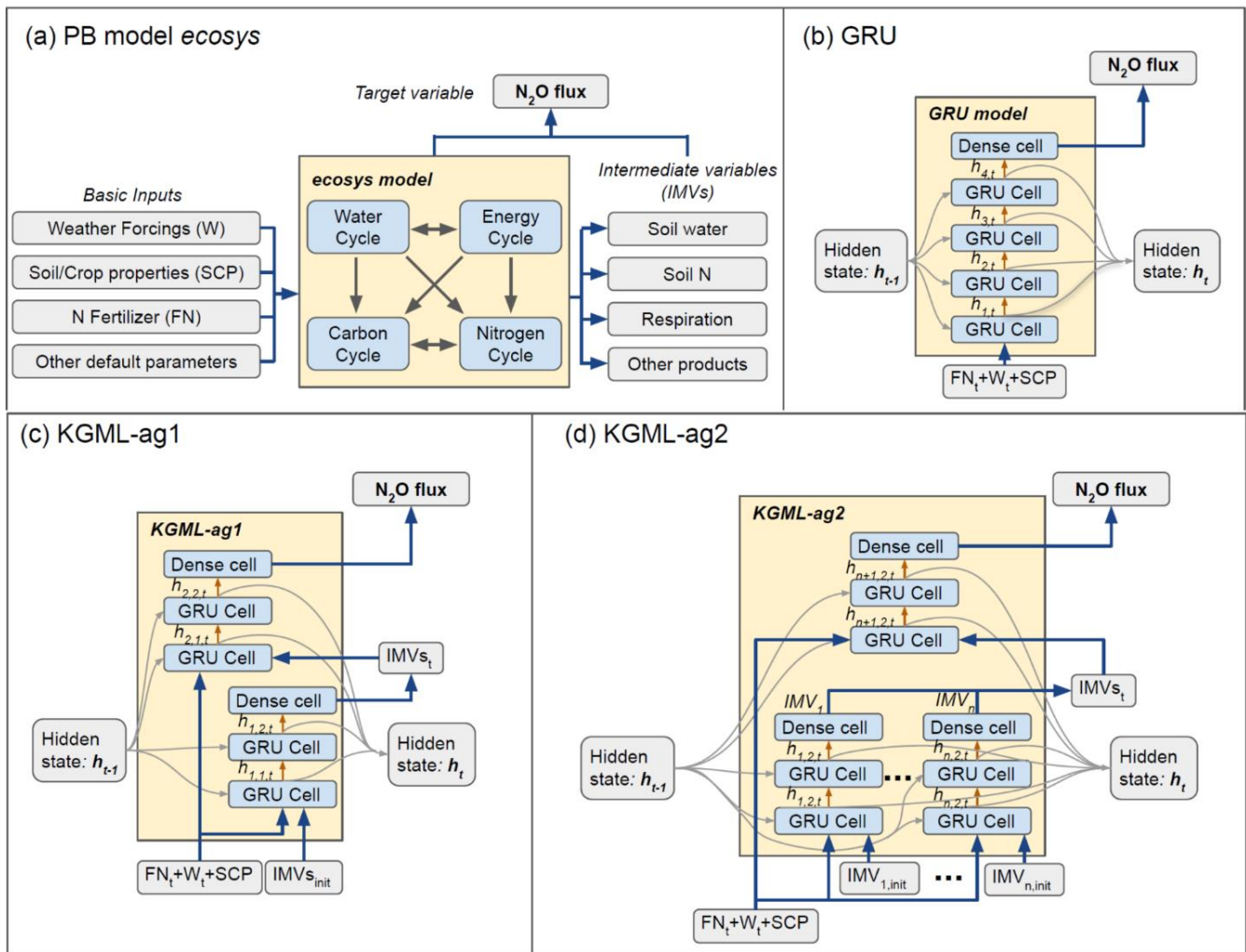
729 Zhang, Y., & Niu, H.: The development of the DNDC plant growth sub-model and the application of DNDC in agriculture: a
730 review, *Agriculture, Ecosystems & Environment*, 230, 271-282, 2016.

731 Zhang, Y., Li, C., Zhou, X., & Moore III, B.: A simulation model linking crop growth and soil biogeochemistry for sustainable
732 agriculture, *Ecological modelling*, 151(1), 75-108, 2002.

733 Zhou, W., Guan, K., Peng, B., Tang, J., Jin, Z., Jiang, C., ... & Mezbahuddin, S.: Quantifying carbon budget, crop yields and
734 their responses to environmental variability using the ecosys model for US Midwestern agroecosystems. *Agricultural and*
735 *Forest Meteorology*, 307, 108521, 2021.

736

737



Input/output
 PB Module/ML cell
 Model

→ Input/output transfer
 → PB module connection
 → h transfer between ML cells with 20% dropout
 → h input and output

738

739

740

741

742

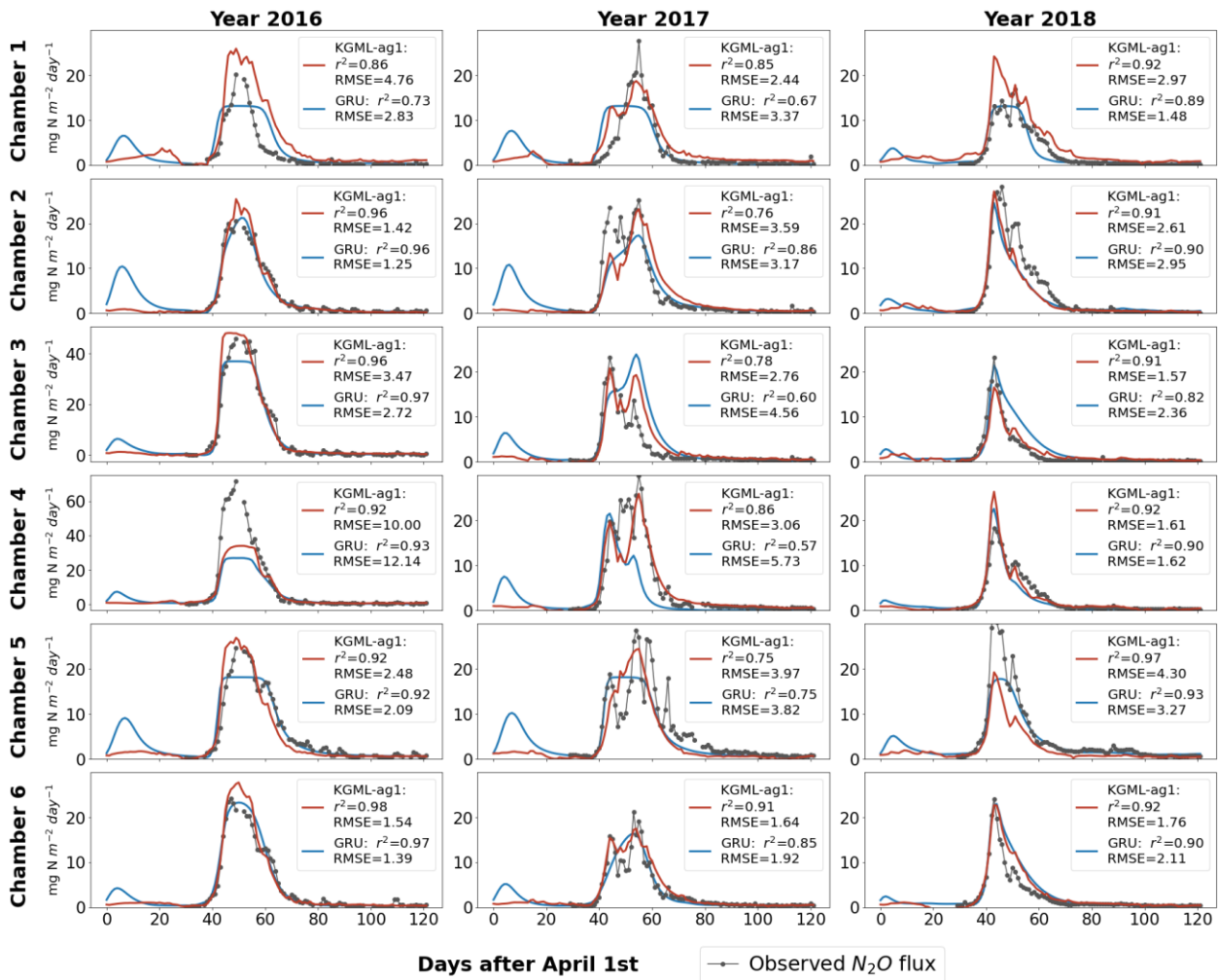
743

744

745

746

Figure 1: The model structures. a) The *ecosys* model; b) Gated recurrent unit (GRU) model; c) KGML-ag1 model with a hierarchical structure; d) KGML-ag2 model with a hierarchical structure with separated GRU modules for IMV predictions. Specifically, in our KGML model design, weather forcings (W) include temperature (TMAX, TDIF), precipitation (PRECN), radiation (RADN), humidity (HMAX and HDIF) and wind speed (WIND); soil/crop properties (SCP) include bulk density (TBKDS), sand content (TCSAND), silt content (TCSILT), pH (TPH), cation exchange capacity (TCEC), soil organic carbon (TSOC), planting day of the year (PDOY) and crop type (CROPT); IMVs include CO₂ flux, soil NO₃⁻ concentration, soil NH₄⁺ concentration, and soil volumetric water content (VWC).



747

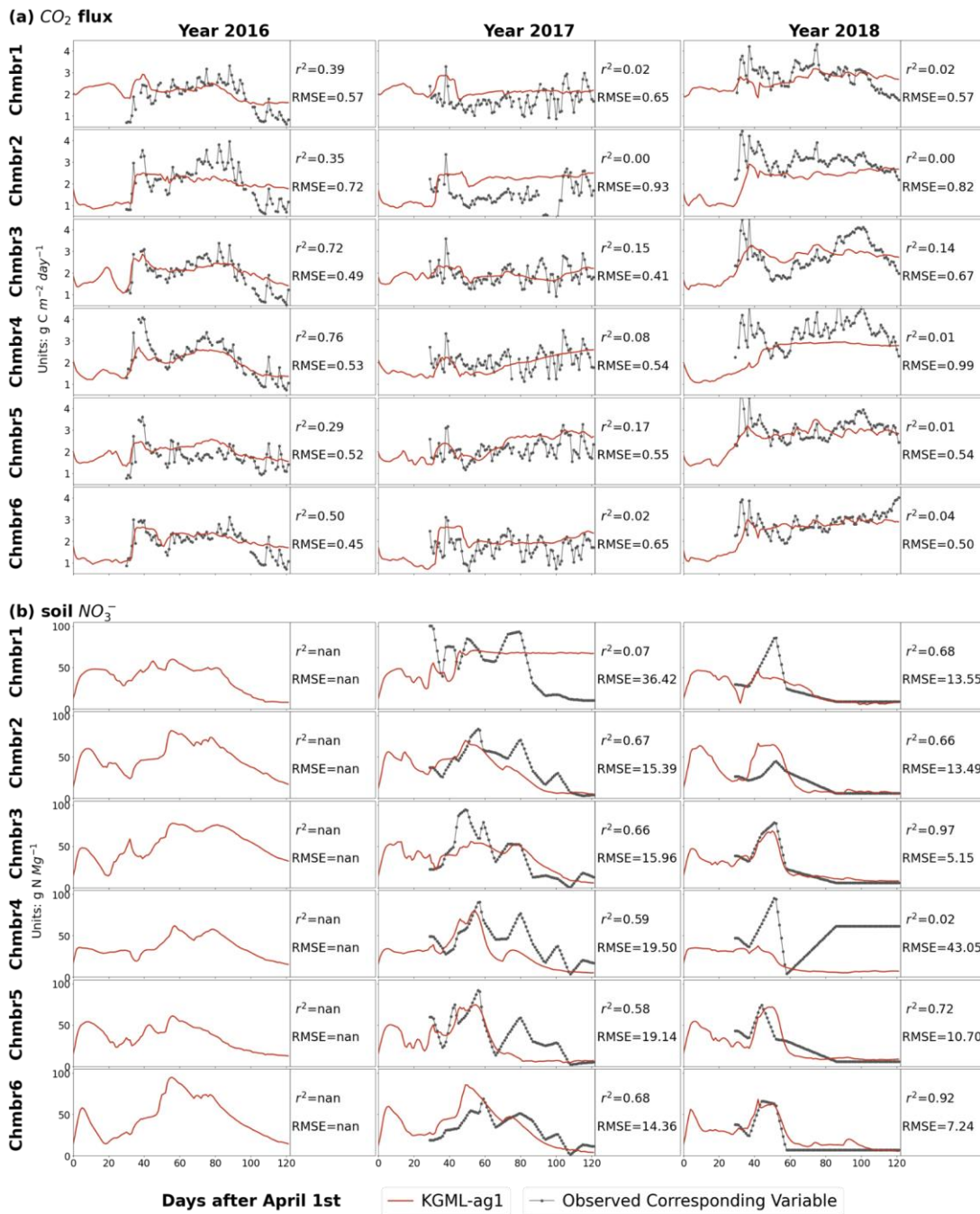
748

749

750

Figure 2: Leave-one-out cross validation of time series of N_2O flux ($mg\ N\ m^{-2}\ day^{-1}$) predicted by the pure non-pretrained GRU model (blue line) and KGML-ag1 model (red line). Observations are shown as black line-dots. Validation results for each chamber were based on out-of-sample predictions by models trained by other five chambers.

751



752

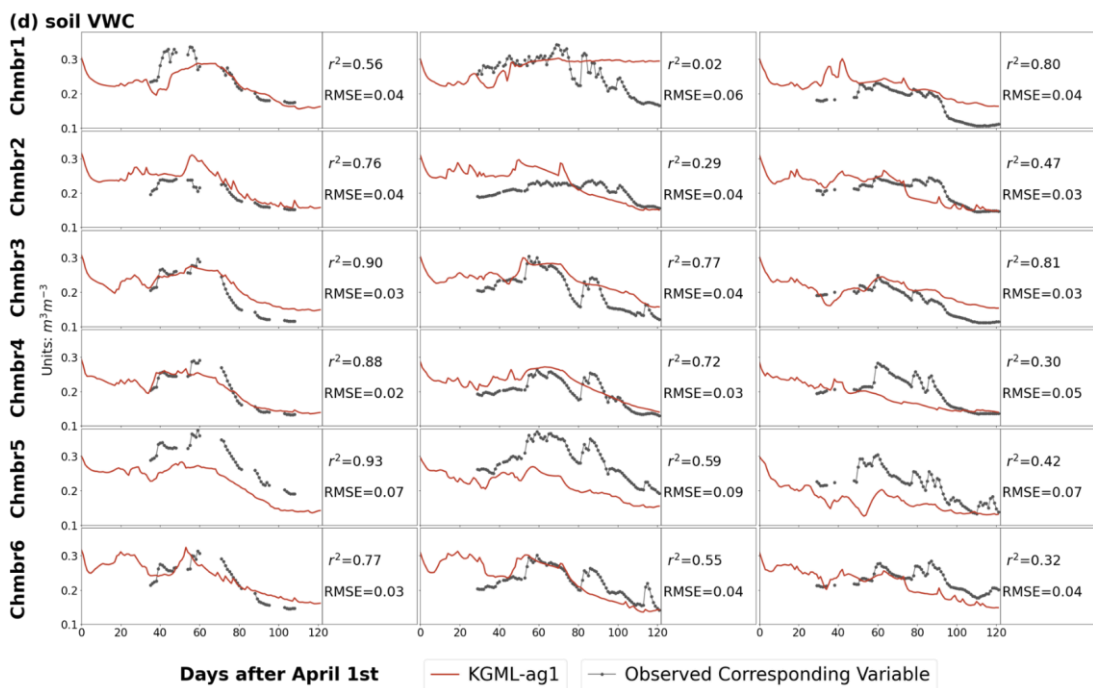
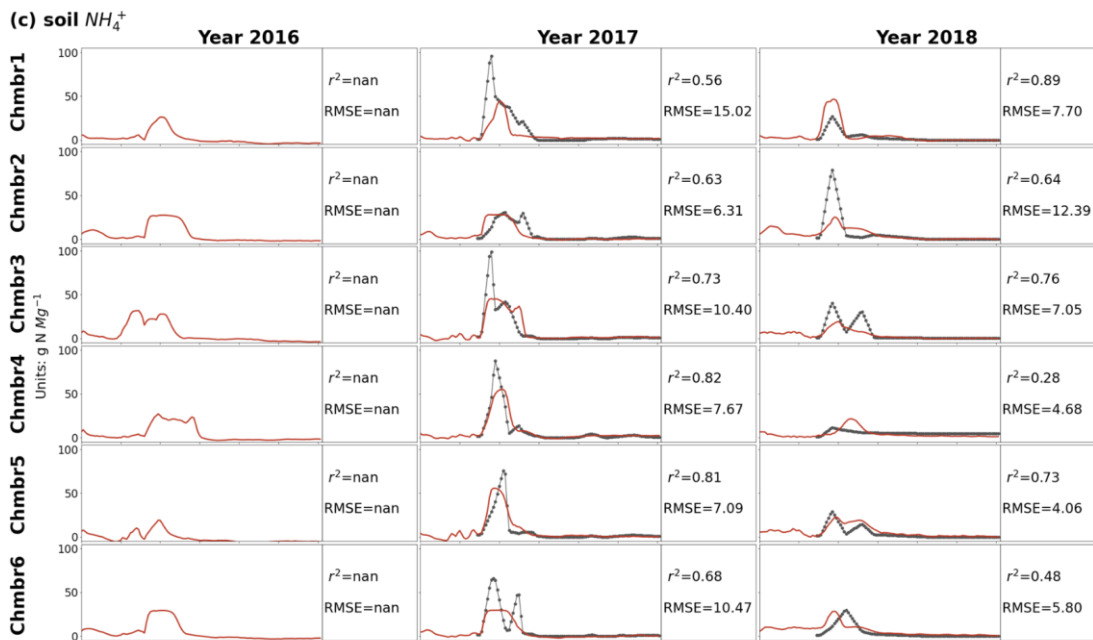
753

754

755

756

Figure 3: Leave-one-out cross validation of time series of IMVs predicted by KGML-ag1 model (red line). Observations are shown as black line-dots. Validation results for each chamber were based on out-of-sample predictions by models trained by other five chambers. Chmb is the abbreviation for chamber. r^2 and RMSE are calculated and present in each year and chamber. The CO₂ flux and soil NO₃⁻ concentration units are g C m⁻² day⁻¹ and g N Mg⁻¹, respectively.



757

758

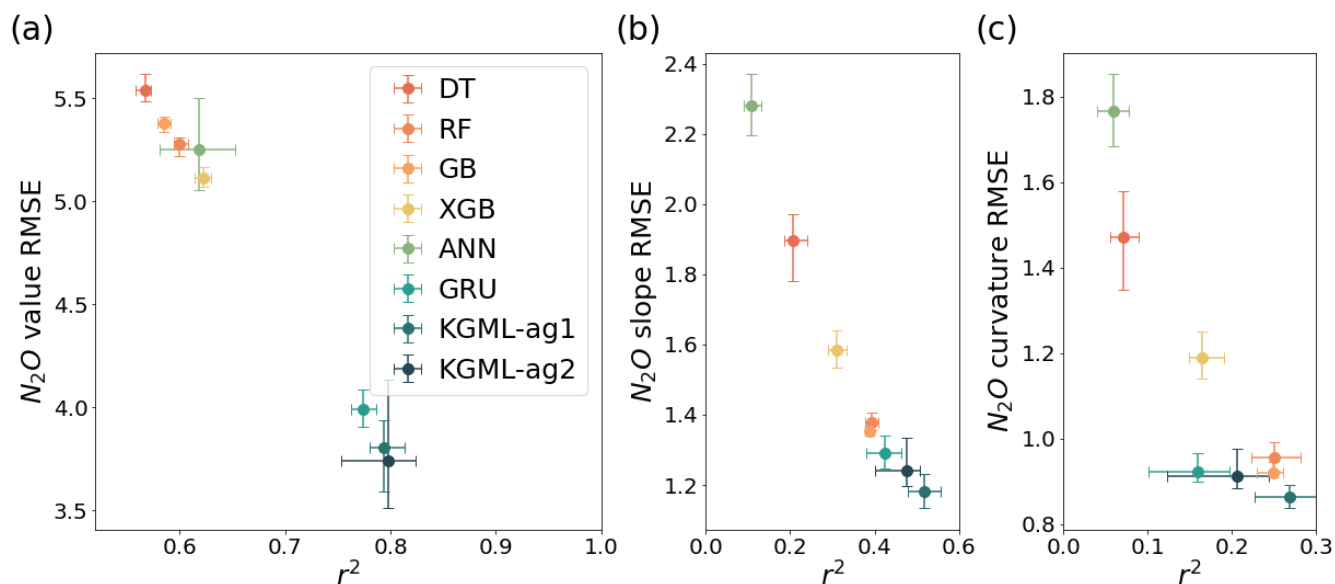
759

760

761

Figure 3 Contd.: Leave-one-out cross validation of time series of IMVs predicted by KGML-ag1 model (red line). Observations are shown as black line-dots. Validation results for each chamber were based on out-of-sample predictions by models trained by other five chambers. Chmb is the abbreviation for chamber. r^2 and RMSE are calculated and present in each year and chamber. The soil NH_4^+ concentration and soil VWC units are $g\ N\ Mg^{-1}$ and $m^3\ m^{-3}$, respectively.

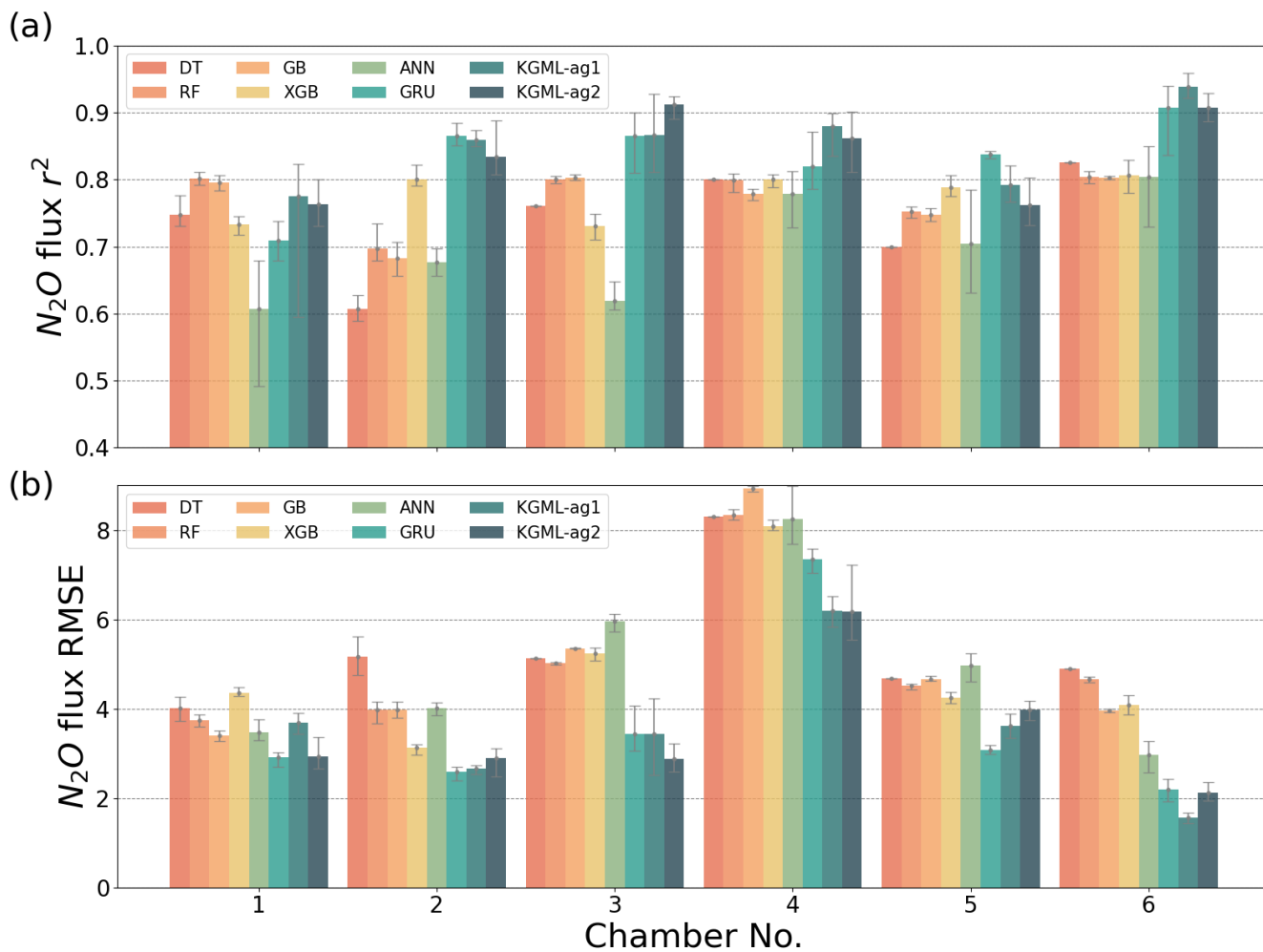
762



764

765 **Figure 4: The comparisons of overall prediction accuracy from leave-one-out cross validation for N_2O value (a), 1st order gradient**
 766 **(slope, b) and 2nd order gradient (curvature, c) between four tree-based ML models (DT, RF, GB and XGB), two deep learning**
 767 **models (ANN and GRU) and KGML-ag models. The overall performances were calculated by comparing out-of-sample predictions**
 768 **(each chamber's predictions were from models trained by other five chambers) from all validated chambers with observations.**
 769 **Different color symbols represent the different models. The x- and y-error bars are coming from the maximum and minimum scores**
 770 **of ensemble experiments. The dot represents the mean score of the ensemble experiments.**

771



772
 773 **Figure 5: The comparisons of N_2O flux prediction accuracy r^2 (a) and (b) RMSE from leave-one-out cross validation, between four**
 774 **tree-based ML models (DT, RF, GB and XGB), two deep learning models (ANN and GRU) and KGML-ag models in six chambers.**
 775 **Validation results for each chamber were based on out-of-sample predictions by models trained by other five chambers. The gray**
 776 **error bars are coming from the maximum and minimum scores of ensemble experiments.**

777

778
779**Table 1: Pretrain results for different model and IMV combinations using *ecosys* synthetic data. Only performances from testing data sets (synthetic data from 19 counties) were presented.**

No.	Pretrain Model	Input Feature N	N ₂ O		CO ₂		NO ₃ ⁻		NH ₄ ⁺		VWC	
			r ²	RMSE	r ²	NRMSE	r ²	NRMSE	r ²	NRMSE	r ²	NRMSE
1	GRU+76IMVs	76 IMVs+FN+7Ws+8SCP	0.98	0.54	-- ^a	--	--	--	--	--	--	--
2	GRU+IMVcb1	4 IMVs+FN+7Ws+8SCP	0.92	1.15	--	--	--	--	--	--	--	--
3	GRU+IMVcb2	3 IMVs+FN+7Ws+8SCP	0.90	1.26	--	--	--	--	--	--	--	--
4	GRU	FN+7Ws+8SCP	0.89	1.37	--	--	--	--	--	--	--	--
5	KGML-ag1+IMVcb1_ini	FN+7Ws+8SCP+4IMV_ini	0.90	1.24	0.91	0.06	0.95	0.03	0.98	0.03	0.95	0.04
6	KGML-ag1+IMVcb2_ini	FN+7Ws+8SCP+3IMV_ini	0.90	1.26	--	--	0.94	0.03	0.97	0.03	0.95	0.04
7	KGML-ag2+IMVcb1_ini	FN+7Ws+8SCP+4IMV_ini	0.90	1.27	0.92	0.05	0.95	0.02	0.98	0.03	0.96	0.04
8	KGML-ag2+IMVcb2_ini	FN+7Ws+8SCP+3IMV_ini	0.91	1.19	--	--	0.95	0.00	0.99	0.02	0.95	0.04

780 ^aThe empty slot indicates that the model does not predict that variable.

781

782

Table 2: Prediction accuracy comparisons between non-pretrained GRU model and KGML-ag1.

	No.	N ₂ O, KGML-ag1 minus GRU				N ₂ O 1st order gradient, KGML-ag1 minus GRU				N ₂ O 2nd order gradient, KGML-ag1 minus GRU			
		All time ^b	Day 30-80	Day 40-65	Day 45-60	All time	Day 30-80	Day 40-65	Day 45-60	All time	Day 30-80	Day 40-65	Day 45-60
Δr^{2a}	All data	0.03 ^c	0.04	0.07	0.10	0.07	0.07	0.07	0.15	0.08	0.08	0.09	0.11
	Chamber1	0.07	0.10	0.20	0.13	0.18	0.18	0.19	0.14	0.08	0.09	0.09	0.02
	Chamber2	-0.04	-0.05	-0.07	-0.05	0.08	0.09	0.09	0.16	0.20	0.20	0.20	0.23
	Chamber3	0.06	0.06	0.08	0.06	0.04	0.04	0.04	0.13	-0.01	-0.01	-0.01	0.07
	Chamber4	0.06	0.08	0.12	0.07	0.05	0.05	0.05	0.14	0.07	0.07	0.08	0.12
	Chamber5	-0.05	-0.06	-0.07	-0.03	0.09	0.09	0.10	0.16	0.13	0.13	0.15	0.11
	Chamber6	0.03	0.04	0.08	0.17	0.14	0.14	0.15	0.22	0.12	0.13	0.14	0.23
$\Delta RMSE^a$	All data	-0.41	-0.56	-0.84	-1.19	-0.07	-0.10	-0.14	-0.20	-0.03	-0.05	-0.07	-0.08
	Chamber1	0.80	1.06	1.21	1.70	0.00	0.00	-0.02	0.00	0.05	0.07	0.10	0.18
	Chamber2	0.08	0.11	0.07	-0.04	-0.10	-0.13	-0.18	-0.14	-0.10	-0.14	-0.19	-0.22
	Chamber3	-0.71	-0.96	-1.30	-2.09	0.03	0.04	0.07	-0.25	0.09	0.13	0.17	0.08
	Chamber4	-1.68	-2.27	-3.09	-3.81	-0.11	-0.15	-0.21	-0.26	-0.05	-0.07	-0.09	-0.16
	Chamber5	0.53	0.69	0.86	0.99	-0.10	-0.14	-0.20	-0.23	-0.09	-0.12	-0.18	-0.14
	Chamber6	-0.20	-0.27	-0.37	-0.61	-0.14	-0.20	-0.29	-0.33	-0.07	-0.10	-0.15	-0.19

783 ^aLeave-one-out cross validation results for each chamber were based on out-of-sample predictions by models trained by other five chambers.
784 The “All data” performances were calculated by comparing out-of-sample predictions from all validated chambers with observations. The
785 difference of r² (Δr^2), and difference of RMSE ($\Delta RMSE$, units are mg N m⁻² day⁻¹, mg N m⁻² day⁻², mg N m⁻² day⁻³ for N₂O value, 1st order
786 gradient and 2nd order gradient, respectively) were calculated by values from KGML-ag1 minus values from GRU.

787 ^bResults from different time windows of different chambers during the period of April 1st-July31st (Days1-122) were detected.

788 ^cNone bold values mean KGML-ag1 outperforms GRU, while **bold** values mean the opposite.

789

790

791 **Table 3: Experiments for measuring GRU and KGML-ag models performance, and influence of pretraining process, training data**
 792 **augmentation and IMV initial values.**

No.	Retrain Model	Experiment	N ₂ O		N ₂ O 1st order gradient		N ₂ O 2nd order gradient		CO ₂		NO ₃ ⁻		NH ₄ ⁺		VWC	
			r ^{2c}	RMSE ^c	r ²	RMSE	r ²	RMSE	r ²	NRMSE	r ²	NRMSE	r ²	NRMSE	r ²	NRMSE
1	GRU, baseline^a	No Pretrain	0.78	4.00	0.45	1.27	0.20	0.90	-- ^b	--	--	--	--	--	--	--
2	GRU	Pretrain	0.80	3.77	0.57	1.12	0.34	0.82	--	--	--	--	--	--	--	--
3	KGML-ag1+ IMVcb1_ini	Original setting	0.81	3.60	0.51	1.20	0.28	0.87	0.37	0.14	0.39	0.21	0.60	0.09	0.33	0.18
4	KGML-ag1+ IMVcb2_ini	Original setting	0.80	3.71	0.49	1.22	0.21	0.91	--	--	0.37	0.22	0.53	0.10	0.33	0.19
5	KGML-ag2+ IMVcb1_ini	Original setting	0.79	3.77	0.48	1.23	0.22	0.90	0.74	0.09	0.46	0.18	0.66	0.08	0.84	0.08
6	KGML-ag2+ IMVcb2_ini	Original setting	0.78	3.91	0.47	1.24	0.20	0.91	--	--	0.49	0.18	0.69	0.08	0.84	0.08
7	KGML-ag1+ IMVcb1_ini	No augmentation	0.80	3.73	0.49	1.22	0.22	0.90	0.38	0.14	0.38	0.21	0.61	0.09	0.37	0.17
8	KGML-ag1+ IMVcb2_ini	No augmentation	0.77	4.04	0.41	1.31	0.13	0.95	--	--	0.38	0.21	0.53	0.10	0.35	0.18
9	KGML-ag2+ IMVcb1_ini	No augmentation	0.76	4.06	0.45	1.27	0.16	0.95	0.69	0.10	0.21	0.25	0.60	0.09	0.80	0.09
10	KGML-ag2+ IMVcb2_ini	No augmentation	0.74	4.27	0.48	1.23	0.21	0.90	--	--	0.40	0.21	0.60	0.09	0.81	0.09
11	KGML-ag1+ IMVcb1_ini	Zero initial values	0.48	6.27	0.26	1.49	0.08	1.00	0.19	0.16	0.25	0.25	0.47	0.12	0.14	0.25
12	KGML-ag1+ IMVcb2_ini	Zero initial values	0.49	5.94	0.31	1.41	0.13	0.95	--	--	0.31	0.25	0.38	0.13	0.24	0.25
13	KGML-ag2+ IMVcb1_ini	Zero initial values	0.48	6.05	0.12	1.66	0.01	1.09	0.58	0.12	0.34	0.25	0.21	0.13	0.56	0.31
14	KGML-ag2+ IMVcb2_ini	Zero initial values	0.39	6.60	0.15	1.59	0.04	1.01	--	--	0.16	0.27	0.27	0.12	0.53	0.31

793 ^aNo.1-6 includes the experiments with original simulation settings as described in Sec. 2 and bold values refers to the baseline GRU
 794 simulation; No.7-10 include the experiments without data augmentation during the finetuning process; And No. 11-14 includes the
 795 experiments of replacing original IMV initial values with zeros.

796 ^bThe empty slot indicates that the model does not predict that variable.

797 ^cThe leave-one-out cross validation overall performances were calculated by comparing out-of-sample predictions (each chamber's
 798 predictions were from models trained by other five chambers) from all validated chambers with observations.