# 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 Liu[1], Shaoming Xu[2], Jinyun Tang[3], Kaiyu Guan[4,5,6], Timothy J. Griffis[7], Matthew D. Erickson[7], Alexander L. Frie[7], Xiaowei Jia[8], Taegon Kim[1, 9], Lee T. Miller[7], Bin Peng[4,5,6], Shaowei Wu[10], Yufeng Yang[1], Wang Zhou[4,5], Vipin Kumar[2], Zhenong Jin[1,11]*

[1]Department of Bioproducts and Biosystems Engineering, University of Minnesota, Saint Paul, MN, 55108, USA
[2]Department of Computer Science and Engineering, University of Minnesota, Minneapolis, MN, 55455, USA
[3]Climate and Ecosystem Sciences Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA
[4]Agroecosystem Sustainability Center, Institute for Sustainability, Energy, and Environment, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA
[5]Department of Natural Resources and Environmental Sciences, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA
[6]National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA
[7]Department of Soil, Water, and Climate, University of Minnesota, Saint Paul, MN 55108, USA
[8]Department of Computer Science, University of Pittsburgh, Pittsburgh, PA, 15260, USA
[9]Department of Smart Farm, Jeonbuk National University, Jeonju, Jeollabuk-do, 54896, Republic of Korea
[10]School of Physics and Astronomy, University of Minnesota, Minneapolis, MN, 55455, USA
[11]Institute on the Environment, University of Minnesota, Saint Paul, MN, 55108, 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

38    $N_2O$ fluxes (overall $r^2 = 0.81$, and RMSE = 3.6 mg N $m^{-2}$ $day^{-1}$ from cross-validation). Importantly KGML-ag always

39    outperforms the PB model and ML models in predicting $N_2O$ fluxes, especially for complex temporal dynamics and emission

40    peaks. Besides, KGML-ag goes beyond the pure ML models by providing more interpretable predictions as well as pinpointing

41    desired new knowledge and data to further empower the current KGML-ag. We believe the KGML-ag development in this

42    study will stimulate a new body of research on interpretable ML for biogeochemistry and other related geoscience processes.

43    **1 Introduction**

44    Nitrous oxide ($N_2O$), with its global warming potential $273 \pm 118$ times greater than that of carbon dioxide ($CO_2$) for a 100-

45    year time horizon, is one of the major greenhouse gases (IPCC6; Forster et al., 2021). The increasing rate of atmospheric $N_2O$

46    concentration during the period 2010-2015 is 44% higher than during 2000-2005, mainly driven by increased anthropogenic

47    sources that have increased total global $N_2O$ emissions to ~17 Tg N $yr^{-1}$ (Syakila and Kroeze, 2011; Thompson et al., 2019).

48    It is estimated that approximately 60% of the contemporary $N_2O$ emission increases are from agriculture management at global

49    scale (Pachauri et al., 2014; Robertson et al., 2014; Tian et al., 2020), but the estimation uncertainty can exceed 300% (Barton

50    et al., 2015; Solazzo et al., 2021). Quantifying $N_2O$ emissions from agricultural soils is extremely challenging, partly because

51    the related microbial processes, mainly about incomplete denitrification and nitrification, are controlled by many environment

52    and management factors such as temperature/water conditions, soil/crop properties, and N fertilization rate, all of which

53    together have collectively led to large temporal and spatial variabilities of $N_2O$ emissions (Butterbach-Bahl et al., 2013; Grant

54    et al., 2016).

55

56    Process-based (PB) models are often used for simulating $N_2O$ fluxes from agroecosystems, but they have some inherent

57    limitations, including incomplete knowledge of the processes, low accuracy due to the under-constrained parameters,

58    expensive computing cost, and rigid structure for further improvements, that we could not resolve by using PB model itself.

59    For example, an advanced agroecosystem model, *ecosys* (Grant et al., 2003, 2006, 2016), simulates $N_2O$ production rates

60    through nitrification and denitrification processes when oxygen ($O_2$) is limited, with equations considering the influence from

61    related substrate concentrations (e.g., $NO_2^-$, $N_2O$, and $CO_2$), nitrifier and denitrifier populations, and soil thermal, hydrological

62    physical and chemical conditions. The produced $N_2O$ accumulates, transfers in gaseous phase, aqueous phase, over different

63    soil layers, and eventually exchanges with atmosphere at the soil surface. Other PB models, including DNDC (Zhang et al.,

64    2002; Zhang and Niu, 2016), DAYCENT (Del Grosso et al., 2000; Necpálová et al., 2015), and APSIM (Keating et al., 2003;

65    Holzworth et al., 2014), have also included processes to simulate $N_2O$ production, but adopt different parameterizations using

66    static partition parameters to estimate $N_2O$ emission from nitrification, and other empirical parameters to control the influence

67    on nitrification from soil water content, pH, temperature and substrate concentrations. Besides, $N_2O$ is intimately connected

68    with the soil organic carbon (SOC) dynamics, because soil nitrifiers and denitrifiers interact strongly with aerobic and

69    anaerobic heterotrophs that process SOC evolution, and all of these microbes are driven by shared environmental variables

70  including soil temperature, moisture, redox status, and physical and chemical properties (Thornley et al., 2007). As expected,

71  these connections make it difficult for PB models, even the most advanced ones like *ecosys*, to find sufficient representations

72  of the physical and biogeochemical processes or obtain enough data to calibrate a large number of model parameters with

73  strong spatio-temporal variations. Thus, novel approaches are needed for addressing the big challenge of agricultural $N_2O$ flux

74  simulations.

75

76  Machine learning (ML) models can automatically learn patterns and relationships from data. Recent studies have investigated

77  the potential to predict agricultural $N_2O$ emission with ML models, including random forest (RF, Saha et al., 2021),

78  metamodelling with extreme gradient boosting (XGBoost) (Kim et al., 2021), and deep learning neural network (DNN)

79  (Hamrani et al., 2020). Notably, Hamrani et al. (2020) compared nine widely used ML models for predicting agricultural $N_2O$.

80  That study pointed out that the long short term memory (LSTM) model with recurrent networks containing memory cells as

81  building blocks will be most suitable for $N_2O$ predictions, but the challenge remains with respect to the ability of capturing the

82  sharp peak of $N_2O$ fluxes and lag time between N fertilizer application and the emission peak. Although there is an increasing

83  interest in leveraging recent advances in machine learning, capturing this opportunity requires going beyond the ML

84  limitations, including limited generalizability to out-of-sample scenarios, demand for massive training data, and low

85  interpretability due to the "black-box" use of ML (Karpatne et al., 2017). PB models with their transparent structures built by

86  representations of physical and biogeochemical processes, seem to be exact complementary to ML models. Thus, combining

87  the power of ML model and PB model understanding innovatively is likely a path forward.

88

89  The above need to integrate ML and PB models can be potentially addressed by the newly proposed framework of Knowledge-

90  guided Machine Learning (KGML) models. In the review by Willard et al. (2021), five research frontiers have been identified

91  regarding the development of KGML for diverse disciplines including earth system science, they are: 1) Loss function design

92  according to physical or chemical laws (Jia et al., 2019, 2021; Read et al., 2019); 2) Knowledge-guided initialization through

93  pretraining ML models with synthetic data generated from PB models (Jia et al., 2019, 2021; Read et al., 2019); 3) Architecture

94  design according to causal relations or adding dense layers containing domain knowledge (Khandelwal et al., 2020; Beucler

95  et al., 2019, 2021); 4) Residual modeling with ML models to reduce the bias between PB model outputs and observations

96  (Hanson et al., 2020); and 5) Other hybrid modeling approaches combining PB and ML models (Kraft et al., 2021). These

97  recent advances in KGML pave the pathway to a more efficient, accurate and interpretable solution for estimating $N_2O$ fluxes

98  from the agroecosystem.

99

100  In this study, we present a first-of-its-kind attempt of developing a KGML for agricultural GHG fluxes prediction (KGML-ag)

101  with knowledge-guided initialization and architecture design, and demonstrate the potential of KGML-ag with a case study on

102  quantifying $N_2O$ flux observed by a multi-year mesocosm experiments. We designed the KGML-ag structure based on the

103  causal relations of related $N_2O$ processes informed by an advanced agroecosystem model, *ecosys* (Grant et al., 2003, 2006,

3

104    2016). We used the synthetic data generated from *ecosys* to design the KGML-ag input/output, and to pre-train the KGML-ag

105    model to learn the basic patterns of each variable. Observations from multi-season controlled-environment mesocosm

106    chambers (Miller, 2021, thesis; Miller et al., 2021, in review) were used to refine the pretrained KGML-ag and evaluate the

107    model performance. Since there is limited literature that guides the development of KGML-ag and not a one that directly

108    addressed GHG fluxes, we investigated a range of ideas to optimize the model performance, including: 1) Using initial values

109    of intermediate variables (IMVs) instead of sequences as model input to reduce data demand; 2) Building hierarchical

110    structures to explicitly estimate IMVs for further $N_2O$ prediction; 3) Using multitask learning to balance the simultaneous

111    training on multiple variables; and 4) Pretraining with millions of synthetic data generated from *ecosys* and fine tuning with

112    mesocosm observations. Although we evaluated the KGML-ag models with real measurements only from a mesocosm

113    experiment, the lessons learned from the development process and various KGML-ag structures can be transferred to other

114    data, other variables and large scale simulations, therefore have broader implications on further KGML related research in

115    agriculture. We believe this study will stimulate a new body of research on interpretable machine learning for biogeochemistry

116    and other related topics in geoscience.

117    **2 Methods**

118    **2.1 Experimental design overview**

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

120    following experiments:

121        1)   With the synthetic data, we developed and pretrained multiple KGML-ag models to learn general patterns and

122            interactions among variables, and evaluated their model performance (Fig. S2, Table 1);

123        2)   With the observed data, we finetuned multiple KGML-ag models to adapt real-world situations, and evaluated their

124            model performance (Fig. 2-3; Fig. S3-5; Table 2-3);

125        3)   We further benchmarked KGML-ag models and uncertainties with other pure ML models without considering

126            temporal dependence, including Decision Tree (DT), Random Forest (RF), Gradient Boosting (GB) from the sklearn

127            package (https://scikit-learn.org/stable/), Extreme Gradient Boosting (XGB) from the XGBoost package

128            (https://xgboost.readthedocs.io/en/latest/) and a 6-linear-layer artificial neural network (ANN) with the mesocosm

129            experiment data by 10 times ensemble experiments (Fig. 4-5; Fig. S6-8);

130        4)   We conducted a few small experiments to further investigate how various model configurations, such as the

131            pretraining process, data augmentation and IMV initial values would influence KGML-ag model performance (Table

132            3).

## 2.2 KGML-ag structure development

### 2.2.1 Generating synthetic data with *ecosys*

We generated synthetic data using a PB model, *ecosys*. The *ecosys* model is an advanced agroecosystem model constructed from detailed biophysical and biogeochemical rules instead of using empirical relations (Grant et al., 2001). It represents $N_2O$ evolution in the microbe-engaged processes of nitrification-denitrification using substrate kinetics that are sensitive to soil nitrogen availability, soil temperature, soil moisture, and soil oxygen status (Grant and Pattey 2008). Two groups of microbial populations, autotrophic nitrifiers and heterotrophic denitrifiers, produce $N_2O$ with specific competitive or cooperative relations in *ecosys* when $O_2$ availability fails to meet $O_2$ demand for their respirations, and $NO_2^-$ become alternative electron acceptors. $N_2O$ transfer within soil layers and from soil to the atmosphere is driven by concentration gradient using diffusion-convection-dispersion equations, in the forms of gaseous and aqueous $N_2O$ under control of volatilization-dissolution (Grant et al., 2016). Unlike the pipeline model described by Davidson et al. (2000) , which mainly considers the correlations of $N_2O$ production with nitrogen availability and of $N_2O$ emissions with soil water content, *ecosys* enables integrative effects of energy, water, nitrogen availability on $N_2O$ production and $N_2O$ transfer via the microbial population dynamics and their interactions with soil, plant, and atmospheric dynamics, under diverse meteorological and anthropogenic disturbances (e.g. runoff, drainage, tillage, irrigation, soil erosion). Many previous studies have demonstrated its robustness in simulating agricultural carbon and nitrogen cyclings at different spatial/temporal scales, and under different management practices (Grant et al., 2003, 2006, 2016; Metivier et al., 2009; Zhou et al., 2021). For the agricultural ecosystems in the US Midwest, whose simulations are used for synthetic data in this study, the performance of *ecosys* on $CO_2$ and $N_2O$ fluxes have been extensively benchmarked, including $CO_2$ exchange (daily Reco $R^2$ = 0.80-0.86; daily NEE, $R^2$ = 0.75-0.897) and leaf area index (LAI, $R^2$ = 0.78) from six flux towers, USDA census reported corn yield ($R^2$ = 0.83) and soybean yield ($R^2$ = 0.80), satellite-derived GPP for corn ($R^2$ = 0.83) and soybean ($R^2$ = 0.85) in the US Midwestfrom Illinois, Iowa and Indiana, and hourlycumulative $N_2O$ fluxesemissions ($R^2$ = 0.36) across eight Midwestern states (ZhouWang et al., 2021; Yang et al., 2022). In addition, *ecosys* model can capture the dynamics and magnitude of $N_2O$ flux in hourly frequency ($R^2$ = 0.2-0.4 and RMSE = 0.1-0.2 mg N m$^{-2}$ h$^{-1}$ in Grant et al., 2008; $R^2$ = 0.28-0.37 and 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, 2008; forest in Grant et al., 2010; and grassland in Grant et al., 2016). Therefore, *ecosys* is an appropriate choice of domain knowledge provider and synthetic data generator in the development of KGML models. We generated daily synthetic data including $N_2O$ flux and 76 IMVs (e.g. $CO_2$ flux from soil, layerwise soil $NO_3^-$ concentration, layerwise soil temperature, and layerwise soil moisture; detailed in Table S1) from *ecosys* simulations for 2000-2018 over 99 randomly selected counties in Iowa, Illinois, and Indiana, USA. We used hourly meteorological inputs (downward shortwave radiation, air temperature, precipitation, relative humidity, and wind speed) from the phase 2 of North American Land Data Assimilation System (NLDAS-2, Xia et al., 2012) and layerwise soil properties (e.g. bulk density, texture, pH, SOC concentration) from the SSURGO database (Soil Survey Staff, 2020) as inputs to *ecosys*. Crop management except N fertilization rates were configured to the same settings as mesocosm experiments (described in Sec 2.2.2). To increase

166      the variability in synthetic data, we implemented 20 different N fertilization rates ranging from 0 to 33.6 g N m$^{-2}$ (i.e. 0 to 300

167      lb N ac$^{-1}$) in each simulation of 99 counties, and more detailed information for model setup refers to Zhou et al. (2021).

168

169      The generated synthetic data were then processed for further use by KGML-ag development. Meanwhile, the hourly weather

170      forcings were converted to seven daily variables, including the maximum air temperature (TMAX_AIR, ℃), difference

171      between the maximum and the minimum air temperature (TDIF_AIR, ℃), the maximum humidity (HMAX_AIR, fraction),

172      difference between the maximum and the minimum humidity (HDIF_AIR, fraction), surface downward shortwave radiation

173      (RADN, W m$^{-2}$), precipitation (PREC, mm day$^{-1}$), and wind speed (WIND, m s$^{-1}$). Six soil properties were retrieved from the

174      SSURGO database, including total averaged (depth weighted averaged for all layers) bulk density (TBKDS, Mg m$^{-3}$), sand

175      content (TCSAND, g kg$^{-1}$), silt content (TCSILT, g kg$^{-1}$), pH (TPH), cation exchange capacity (TCEC, cmol$^{+}$ kg$^{-1}$) and soil

176      organic carbon (TSOC, g C kg$^{-1}$); and two crop properties were retrieved, including planting day of the year (PDOY) and crop

177      type (CROPT, 1 for corn and 0 for soybean). Finally, each synthetic data sample has daily N$_2$O flux, 76 selected IMVs, 7

178      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

179      from April 1st to July 31st (122 days) were selected to cover the mesocosm observations (around 30 days before and 90 days

180      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

181      (about 4.3 million data points). We randomly selected the samples from 70 counties for training, 10 counties for validation,

182      and 19 counties for testing.

183      **2.2.2 Mesocosm experiments for KGML-ag model fine-tuning and evaluation**

184      Observations were acquired from a controlled-environment mesocosm facility on the St. Paul campus of the University of

185      Minnesota. Soil samples were sourced in 2015 from a farm in Goodhue County, MN (44.2339° N and 92.8976° W), which had

186      been under corn-soybean rotation for 25 years. Six chambers with a soil surface area of 2 m$^2$ and column depth of 1.1 m were

187      used to plant continuous corn during 2015-2018 and monitor the N$_2$O flux response to different precipitation treatments. The

188      experiment also measured other environmental variables including air temperature and photosynthetically active radiation

189      (PAR), which were controlled to mimic the outdoor ambient environment. Granular urea fertilizer was hand broadcasted and

190      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

191      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

192      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

193      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

194      35,000 seeds ac$^{-1}$. Crops were harvested at the end of September by cutting the stover five inches above the soil. Hourly N$_2$O

195      fluxes (mg N m$^{-2}$ h$^{-1}$) and CO$_2$ fluxes (g C m$^{-2}$ h$^{-1}$) were measured using non-steady-state flux chambers with a CO$_2$ analyzer

196      (LI-10820 for 2016 and LI-7000 for 2017 and 2018, LI-COR Biosciences, Lincoln, NE) and a N$_2$O analyzer (Teledyne

197      M320EU, Teledyne Technologies International Corp, Thousand Oaks, CA) (Detail method can be retrieved from Fassbinder

198      et al., 2012, 2013). We also collected soil moisture at 15 cm depth (VWC as abbreviation of volumetric water content, m$^3$ m$^{-}$

199   $^3$), weekly 0-15 cm depth soil $NO_3^-$ + $NO_2^-$ concentration ($NO_3^-$ for short in the following text, g N $Mg^{-1}$), soil $NH_4^+$
200   concentration ($NH_4^+$, g N $Mg^{-1}$), and related environment variables including air temperature, radiation, humidity and soil/crop
201   properties from three growing seasons during 2016-2018 and six mesocosm chambers (Fig. S1). The magnitude of $N_2O$ flux
202   and $NO_3^-$ soil concentration and their responses following fertilizer application from this mesocosm experiment are slightly
203   higher than~~consistent with~~ several field studies of agricultural soils (Fassbinder et al., 2013; Grant et al., 1999, 2006, 2008,
204   2016; Hamrani et al., 2020; Venterea et al., 2011). More details about the mesocosm facility and experimental design can be
205   found in the thesis of Miller L. (2021).

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

218   To reduce overfitting and increase the generalization of the trained model based on the small amount of mesocosm data, we
219   applied the following method to augment the experimental measurements and weather forcings to 1000 times larger by
220   sampling hourly data and averaging them to daily scale. In this method, 16 hours (or maximum valid hours) of data are
221   randomly selected from 24 hours of data to compute their mean as the daily value. Since up to 2~~3~~/3~~4~~ of the day is~~are~~ covered
222   by the selected data (16 hours /24 hours), the augmented daily values should be representative enough for the source day with
223   ~~and meanwhile present~~ slight variations from each other. Furthermore, the observation ratio, (24 hours - missing hours) / 24
224   hours, can be used as the weights in loss function to inject the data quality information in model optimization. If the day has
225   more than 16 hours missing values, we consider the observations in that day as not trustworthy and drop the day by setting the
226   weight to 0. This method can not only augment the data to 1000 times larger but also deal with the missing values in observed
227   data inherently. The total amount of observed mesocosm data and related weather forcings are augmented to 122 days x 3
228   years x 6 chambers x 1000 data samples in this study.

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

Hamrani et al. (2020) compared different models and reported that LSTM provided the highest accuracy in predicting $N_2O$ fluxes, because $N_2O$ flux is time dependent by its production/consumption nature and LSTM simulates target variables by considering both current and historical states. The LSTM model, proposed by Hochreiter and Schmidhuber (1997), uses a cell state as an internal memory to preserve the historical information. At each time step, it creates a set of gating variables to filter the input and historical information and then uses the processed data to update the cell state. Similar to LSTM, GRU is a gated recurrent neural network but only keeps one hidden state (Cho et al., 2014). Though simpler than LSTM, GRU is proved to have similar performance (Chung et al., 2014). Our preliminary test on synthetic data for $N_2O$ prediction showed that GRU indeed provided similar or higher accuracy and model efficiency under different model settings than LSTM (Table S2). This is possible because simpler models with fewer weights and hyperparameters are more robust in combating the overfitting problem. Therefore, we choose GRU as the basis of KGML-ag development.

### 2.2.4 Incorporating domain knowledge to the development of KGML-ag

To quantitatively reveal the correlations between $N_2O$ fluxes and IMVs and guide the KGML-ag development, we conducted feature importance analysis by a customized 4-layer GRU ML model (Fig. 1b). Each layer of the model has a GRU cell with 64 hidden units. The 4-layer structure makes the model deeper and capable of capturing complex interactions. Between each GRU cell, 20% of the output hidden states are randomly dropped by replacing them with zero values (so called 20% dropout) to avoid overfitting. A linear dense layer is used to map the final output to $N_2O$. We first trained GRU models using synthetic data with different combinations of IMVs as inputs to predict the $N_2O$ fluxes (original-test, Table S2). The feature importance analysis of well-trained models was then implemented by replacing one input feature with a Gaussian noise with mean μ=0 and standard deviation σ=0.01, while keeping others untouched (new-test). The importance score was calculated by the new-test's root mean square error (RMSE) (replacing one feature) minus the original-test's RMSE (no replacing). RMSE was calculated by $\frac{\sqrt{\sum_1^N (y_i - y_{i'})^2}}{N}$ where $N$ is the total number of observations across time and space, $y_i$ is i-th measurement from synthetic data or observed data and $y_i'$ is its corresponding prediction.

To find important variables for $N_2O$ flux prediction in an ideal situation where all variables are available, we conducted a feature importance analysis for GRU models with all IMVs and basic inputs including FN, 7 W and 8 SCP (Fig. S2a). Results indicated that flux variables including $NH_3$, $H_2$, $N_2$, $O_2$, $CH_4$, evapotranspiration (ET) and $CO_2$ had significant influence on the model performance. Variables ranked high in feature importance analysis are considered with priority during model development. To develop a functionable KGML-ag, we further investigated the feature importance of four IMVs that are available from mesocosm observations including $CO_2$, $NO_3^-$, VWC and $NH_4^+$, which were ranked 7th, 20th, 58th, 60th respectively in 92 input features of synthetic data (Fig. S2a). We used these four available IMVs to create two input

261  combinations: 1) $CO_2$ flux, $NO_3^-$, VWC and $NH_4^+$ (IMVcb1), and 2) $NO_3^-$, VWC and $NH_4^+$ (IMVcb2). The objective of

262  building IMVcb2 was to investigate the importance of the highly ranked variable $CO_2$ flux (by removing it from the inputs),

263  and the impact of mixing-up flux and non-flux variables on model performance. We tested the feature importance of the GRU

264  models built with IMVcb1 and IMVcb2 to check whether they would help in $N_2O$ prediction (Fig. S2b-c). All the feature

265  importance results above indicated the correlation intensity between $N_2O$ and many other variables, which would help the

266  KGML-ag model development and interpretation in this study (rest of this section and Sec. 3.1), and would guide future $N_2O$

267  related measurements and KGML model development (discussed in Sec. 4.3).

268

269  Next we used the knowledge learned from synthetic data to develop the structure of KGML-ag (Fig. 1c-d). Previous studies

270  for KGML models have used physical laws, e.g., conservation of mass or energy, to design the loss function for constraining

271  the ML model to produce physically consistent results (Read et al., 2019; Khandelwal et al., 2020). However, for complex

272  systems like agroecosystems, it is challenging to incorporate physical laws, such as mass balance for $N_2O$, into the loss function

273  due to the incomplete understanding of the processes and the lack of mass balance related data for validation. An alternative

274  solution is to incorporate such information in the design of the neural network (Willard et al., 2021). Effectiveness of such an

275  approach was demonstrated by Khandelwal et al. (2020) in the context of modeling stream flow in a river basin using Soil &

276  Water Assessment Tool (SWAT). They used a hierarchical neural network to explicitly model IMVs (e.g., soil moisture, snow

277  cover) and their relationships with the target variable (streamflow) and showed that this model is much more effective than a

278  neural network that attempts to directly learn the relationship between input drivers and the target variables. Following this

279  idea, we identified four desired features of an effective KGML-ag model, including: 1) We used initial values instead of

280  sequence of the IMVs from synthetic data or observed data to provide a solid starting state for the ML system and reduce the

281  IMV data demand, and then used the rest of the data to further constrain the prediction of IMVs; 2) We built a hierarchical

282  structure based on the structure of process representation in *ecosys* to first predict IMVs and then simulate $N_2O$ with predicted

283  IMVs; 3) We trained all variables together using multitask learning to reach the best prediction scores, which generalized the

284  model and incorporated interactions between IMVs and $N_2O$; 4) We initialized the KGML-ag model by pretraining with

285  synthetic data before using real observed data to transfer physical knowledge, which further reduced the demand on large

286  training samples and aided in faster convergence for fine-tuning.

287

288  To meet these desired features, we proposed two KGML-ag models (Fig. 1c-d). The first model, KGML-ag1, is a hierarchical

289  structure containing two modules to simulate IMVs and $N_2O$ sequentially. Each module is a 2-layer 64 units GRU ML model.

290  The inputs to the module of the KGML-ag1 model for IMV predictions (KGML-ag1-IMV module) are FN, 7W and 8SCP

291  together with the initial values of IMVs, and the outputs are IMV predictions. The inputs to the module of the KGML-ag1

292  model for $N_2O$ predictions (KGML-ag1-$N_2O$ module) are FN, 7W, 8SCP and predicted IMVs from KGML-ag1-IMV, and the

293  output is the target variable $N_2O$. Linear dense layers were coded for both modules to map output states to IMVs or $N_2O$. The

294  dropout method was applied to drop 20% of the state output between GRU cells and dense layers. The second model, KGML-

9

295 ag2, is also a hierarchical structure similar to KGML-ag1, but has multiple KGML-ag2-IMV modules to explicitly simulate
296 IMVs by tuning them separately in the fine-tuning process (discussed in Sec. 2.2.5). Each KGML-ag2-IMV module in KGML-
297 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
298 prediction. The KGML-ag2-N$_2$O module collects the IMV predictions from KGML-ag2-IMV modules and predicts the N$_2$O
299 with inputs of FN+7W+8SCP and predicted IMVs.

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

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

310 During the pretraining process, we initialized the IMV of KGML-ag using the first day value of synthetic IMV time series.
311 Adam optimizer with a start learning rate of 0.0001 was used for the training process. The learning rate would decay by 0.5
312 times after every 600 training epochs. At each epoch, synthetic data samples were randomly shuffled before being input to the
313 model to predict N$_2$O (and IMVs if any). The mean square error (MSE) loss (calculation was equal to the square of RMSE) or
314 sum of MSE loss (if multitask learning) between predictions and *ecosys* synthetic observations were calculated to optimize the
315 weights of GRU cells. After the training process updated the model's weights, the validation process was performed to evaluate
316 the model performance based on untouched samples with RMSE and the square of Pearson correlation coefficient (r$^2$). r$^2$ was
317 calculated as $\frac{(\sum_i (y_i{'} - \overline{y_i{'}})(y_i - \overline{y_i}))^2}{\sum_i (y_i{'} - \overline{y_i{'}})^2 (y_i - \overline{y_i})^2}$, where $y_i$ is the i-th measurement from synthetic data or observed data, $y_i{'}$ is its
318 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
319 diagnosing space. If both validated r$^2$ and RMSE were better than the best values in previous epochs, the updated model in this
320 epoch would be saved. Normalized RMSE (NRMSE, calculated by RMSE/(max-min) of each variable observation) was
321 introduced to evaluate IMV predictions between variables with different value ranges.

323 During the fine-tuning process, we used estimated IMV initial values of 1.0 g C m$^{-2}$, 0.2 m$^3$ m$^{-3}$, 0.0 g N Mg$^{-1}$, and 20.0 g N
324 Mg$^{-1}$ for CO$_2$, VWC, NH$_4{}^+$, and NO$_3{}^-$ respectively, from starting day (April 1st) to the day before the first day of real
325 observations, as input to KGML-ag models. Then the first-day values of observed IMVs were input into KGML-ag during the

326　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

327　augment the total amount of data 1000 times larger for the fine-tuning process. The purpose of this data augmentation method

328　was to increase the generalization of the fine-tuned model and to overcome the overfitting due to small sample size. The mask

329　matrix was elementarily multiplied to the output matrix to calculate the MSE, $r^2$ and RMSE only for days with observations.

330　The similar optimizer was used with an initial learning rate of 0.00005 and decay fraction of 0.5 per 200 epochs. Other

331　training/validation methods in each epoch were similar to the pretraining process. Specifically, in the KGML-ag1 model

332　finetuning process, we first froze the KGML-ag1-$N_2O$ module and only trained the KGML-ag1-IMV module for IMVs. After

333　finishing the KGML-ag1-IMV module training, we froze the KGML-ag1-IMV module and trained the KGML-ag1-$N_2O$

334　module for $N_2O$. In the KGML-ag2 fine-tuning process, the similar freezing method was used but different KGML-ag2-IMV

335　modules were trained separately one by one.

336　**2.3 Development environment description**

337　We used the Pytorch 1.6.0 (https://pytorch.org/get-started/previous-versions/) and python 3.7.9

338　(https://www.python.org/downloads/release/python-379/) as the programing environment for the model development. In order

339　to use the GPU to speed-up the training process, we installed cudatoolkit 10.2.89 (https://developer.nvidia.com/cuda-toolkit).

340　A desktop with Nvidia 2080 super GPU was used for code development and testing. The Mangi cluster

341　(https://www.msi.umn.edu/mangi) from High Performance Computing of Minnesota Supercomputing Institute (HPC-MSI,

342　https://www.msi.umn.edu/content/hpc) with 2-way Nvidia Tesla V100 GPU was used in training processes which consumed

343　longer time and bigger memories.

344　**3 Results**

345　**3.1 Pretraining experiments using synthetic data from *ecosys***

346　In the pretraining stage, the GRU model with 76 IMVs achieved the best performance in predicting $N_2O$ fluxes ($r^2$=0.98, RMSE

347　=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

348　1). The high performance was due to some flux IMVs such as $NH_3$, $H_2$, $O_2$, $CO_2$ and ET, which are highly correlated to $N_2O$

349　(Fig. S2a), were used as input to the model. The good performance of GRU with all IMVs indicates that ML models are able

350　to perfectly mimic *ecosys* when sufficient information about IMVs is available. The GRU model with only basic input of N

351　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

352　mg N m$^{-2}$ day$^{-1}$ (Table 1). The relatively low performance is likely because this model failed to capture several highly nonlinear

353　pathways that are employed by ecosys to predict $N_2O$ (e.g., one influence pathway from precipitation to $N_2O$ can be:

354　Precipitation $\rightarrow$ soil moisture $\rightarrow$ N components solubility/concentration $\rightarrow$ nitrification/denitrification rate/amount $\rightarrow$ soil

355　$N_2O$ concentration $\rightarrow$ gas $N_2O$ flux). When adding sequences of IMV combinations (i.e., IMVcb1 of $CO_2$ flux, $NO_3^-$, $NH_4^+$

356 and VWC, and IMVcb2 of $NO_3^-$, $NH_4^+$ and VWC), the GRU models performed slightly better than the GRU model using only

357 basic inputs, achieving $r^2$ of 0.92 and 0.90, respectively (Table 1). The KGML-ag1 with IMVcb1 and IMVcb2 initial values

358 provided better performance (both $r^2$ = 0.90) than GRU with basic input and comparable performance to the GRU with inputs

359 of IMVcb1 and IMVcb2 sequence. Besides, KGML-ag1 provided predicted IMVs of $CO_2$, $NO_3^-$, $NH_4^+$, and VWC with $r^2$ over

360 0.91, and NRMSE below 0.06 (Table 1). KGML-ag2 also provided comparable $N_2O$ performance but relatively better IMVs

361 performance of $r^2$ over 0.92 and NRMSE below 0.05. Results indicated that KGML-ag models with IMV initial values as extra

362 input performed similar or better than pure ML models in synthetic data.

363 **3.2 KGML-ag evaluation using observed data from mesocosm**

364 After being fine-tuned with observed data, KGML-ag1 had $N_2O$ prediction overall accuracy of $r^2$=0.81 and RMSE=3.6 mg N

365 $m^{-2}$ $day^{-1}$, while non-pretrained GRU model provided $r^2$=0.78 and RMSE=4.0 mg N $m^{-2}$ $day^{-1}$, and pretrained GRU model

366 provided $r^2$=0.80 and RMSE=3.77 mg N $m^{-2}$ $day^{-1}$ (Table 3). The time series of $N_2O$ predictions from KGML-ag1 and the non-

367 pretrained GRU model were further compared (Fig. 2), from which we found at least two advantages of using KGML-ag1 for

368 $N_2O$ predictions: 1) For the region without observation data (normally before day 25), KGML-ag1 predicted stable $N_2O$ fluxes

369 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.

370 This is because KGML-ag1 has learned knowledge for the whole period from the pretraining process with *ecosys* model

371 generated synthetic data, but GRU model has no prior knowledge for the period without any data in observations; 2) Although

372 KGML-ag1 had a lower accuracy than GRU in some chambers, KGML-ag1 can better capture the temporal dynamics of $N_2O$

373 fluxes compare to GRU, especially when the fluxes are highly variable (e.g. Fig 2 chamber 2).

374

375 To validate KGML-ag1 robustness, we further investigated the KGML-ag1 and GRU model performance in different temporal

376 windows, shrinking from the whole period to the $N_2O$ peak occurrence time (days 1-122, day 30-80, day 40-65 and day 45-60

377 for year 2016-2018), and performance in $N_2O$ flux, first order gradient of $N_2O$ (slope) and second order gradient of the $N_2O$

378 (curvature) (Table 2). Slope represents the speed of $N_2O$ flux changes through time and curvature represents the acceleration.

379 Assessing prediction performance with these two metrics will reveal the model robustness on capture variable dynamics, which

380 is critical when predicting fast-change variables with hot moments (a short period of time with rare events like flux increasing

381 quickly) like $N_2O$. First of all, the overall $r^2$ and RMSE of KGML-ag1 for values, slope and curvature were always better than

382 GRU. In particular, KGML-ag1 captured the peak region (e.g., days 45-60) much better than GRU in both magnitude and

383 dynamics (Table 2, Fig 2). Even for chamber 2 and 5 in which KGML-ag1 made worse $N_2O$ predictions than GRU ($\Delta r^2$ ranging

384 from -0.07 to -0.03), it better captured temporal dynamics than GRU in terms of slope ($\Delta r^2$ ranging from 0.08 to 0.16) and

385 curvature ($\Delta r^2$ from 011 to 0.23) (Table 2). For other chambers, KGML-ag1 outperformed GRU consistently. For chamber 1,

386 KGML-ag1 had worse $N_2O$ predictions RMSE than GRU but the $\Delta r^2$ increased as the window shrinks to the peak emission

387 time ($0.07 \rightarrow 0.13$). The slope and curvature for chamber 1 also indicated that KGML–ag1 captured the dynamics much

12

388  better than GRU. For chamber 3, KGML–ag1 predicted better $N_2O$ but presented worse slope and curvature RMSE than

389  GRU (Table 2). However, when explicitly investigating the time series of $N_2O$ flux, slope and curvature in each year, KGML-

390  ag1 outperformed GRU more significantly in 2017, the year with more complex temporal dynamics of $N_2O$ fluxes, than in

391  2016 and 2018, especially for chamber 3 (Fig. 2; Fig. S3-4). This investigation supported that KGML-ag1 was more capable

392  for complex dynamics predictions.

393

394  Interestingly, the fine-tuned KGML-ag1 model predicted reasonable IMVs including $CO_2$, $NO_3^-$, $NH_4^+$, and VWC with overall

395  $r^2$ 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

396  between IMV predictions and observations further indicated that KGML-ag1 could reasonably capture both magnitude and

397  dynamics (Fig. 3). KGML-ag2 presented better IMVs predictions than KGML-ag1, with overall $r^2$ of $CO_2$, $NO_3^-$, $NH_4^+$, and

398  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

399  lower $r^2$ (decreasing 0.02) of $N_2O$ (Table 3; Fig. S5). This indicated that explicitly simulating each IMV with separated KGML-

400  ag2-IMV modules did not benefit the $N_2O$ flux prediction accuracy, likely due to increasing model complexity which resulted

401  in reduced stability and ignoring the IMV interactions. In addition, we also found all KGML-ag models would perform better

402  by using IMVcb1 (with $CO_2$) than using IMVcb2 (without $CO_2$) in real data tests, indicating feature importance analysis based

403  on synthetic data can be a reasonable substitute for analysis with the often limited real-world data.

404  **3.3 KGML-ag comparing with other pure ML models**

405  The results from eight different models showed that KGML-ag1 comparing with other pure ML models consistently provided

406  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 $r^2$ (0.78-

407  0.81, 0.48-0.56, and 0.23-0.31) for $N_2O$ fluxes, slope and curvature, respectively (Fig. 4). This indicated that KGML-ag1

408  outperformed other pure ML models in capturing both the magnitude and dynamics of $N_2O$ flux. Meanwhile, we have

409  calculated the uncertainty of mesocosm measurement due to converting hourly data to daily data during 30-80 days by using

410  augmented value minus mean of the augmented values (-10.2 to 10.4 mg N m$^{-2}$ day$^{-1}$, and standard deviation =1.4 mg N m$^{-2}$

411  day$^{-1}$). KGML-ag1 during the same period has comparable uncertainties based on ensemble simulations (calculated by

412  ensemble value minus mean of ensemble values; -14.4 to 15.2 mg N m$^{-2}$ day$^{-1}$, with standard deviation = 1.3 mg N m$^{-2}$ day$^{-1}$).

413  KGML-ag2 presented slightly better mean scores for $N_2O$ flux predictions than KGML-ag1, but worse scores for slope and

414  curvature and larger uncertainties. This proved the hypothesis discussed in section 3.2 that KGML-ag2 didn't benefit the

415  magnitude and dynamics predictions of $N_2O$ flux with its more complex structure and less connections between IMVs.

416

417  Within the tree-based models (DT, RF, GB and XGB), the simplest model DT provided the worst predictions for $N_2O$ flux,

418  slope and curvature. The XGB model provided the highest $N_2O$ flux accuracy with $r^2$ of 0.61-0.63 and RMSE of 5.07-5.17 mg

419  N m$^{-2}$ day$^{-1}$, while the GB model provided best slope and curvature predictions with $r^2$ of 0.38-0.40 and 0.23-0.26, and RMSE

420  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 $N_2O$ flux accuracy and relatively low

421 slope and curvature accuracy of the XGB model implied that there is a trade-off between the abilities of capturing dynamics
422 and magnitude.
423
424 In the group of deep learning models including ANN, GRU and KGML-ag1, ANN provided the worst predictions. Even with
425 the better $N_2O$ flux predictions than most tree-based models (except XGB), the slope and curvature predictions of ANN were
426 the worst among all eight models. This implied that the trade-off between accurately capturing $N_2O$ dynamics to magnitude in
427 ANN was significant. But when considering the temporal dependence, deep learning model GRU and KGML-ag1
428 outperformed all other models in flux, slope and curvature predictions. This indicated that without considering temporal
429 dependence the improvement in $N_2O$ flux prediction accuracy could be risky by causing the performance drop in capturing
430 dynamics.
431
432 The detailed model comparisons in each chamber are shown in Fig. 5 ($N_2O$ flux) and Fig. S6-7 ($N_2O$ slope and curvature),
433 where the results are found to follow the same pattern as described above. In addition, time series comparisons of chamber 3
434 and 4 in 2017 between different models are presented in Fig. S8 as two examples. For periods without any observed data, we
435 assumed that the good model predictions should be stable, consistent with the nearest period and close to the reality in the
436 experiment setting (e.g. no erratic peak and $N_2O$ flux near 0 mg N m$^{-2}$ day$^{-1}$ before day 25). From these comparisons, we infer
437 that without considering temporal dependence and pretraining process, the tree-based model including DT, RF, GB and XGB
438 and deep learning model ANN predicted erratic peaks in almost every missing data point, while the GRU model was stable in
439 small gaps short missing period (1-2 days of missing data) and only presented poor performance in long missing period (before
440 25 day 25). This improvement by the GRU model maycan be attributed to the structure of GRU that naturally keeps the
441 historical information using hidden states, which enables GRU to consider the temporal dependence and make consistent
442 predictions over time.

443 **3.4 Influence of pretraining process, data augmentation and using IMV initial values as input feature**

444 After we pretrained the GRU model with synthetic data, the overall $r^2$ of $N_2O$ flux predictions in observed data increased by
445 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
446 and curvature predictions, respectively, compared to non-pretrained GRU (No.1-6 in Table 3 gray region). The gap between
447 the GRU model with pretrain and KGML-ag1 in $N_2O$ value prediction shows the improvement resulting from architecture
448 change ($r^2$ increases by 0.01 and RMSE decreases by 0.17 mg N m$^{-2}$ day$^{-1}$). Although pretrained GRU had higher slope and
449 curvature prediction accuracy than KGML-ag models, it still couldn't achieve the current $N_2O$ value prediction accuracy of
450 KGML-ag1. Besides, the KGML-ag models had relatively shallow $N_2O$ prediction modules (2-layer GRU KGML-ag-$N_2O$
451 module of KGML-ag models vs 4-layer GRU) but included modules for IMV predictions, which therefore increased the model
452 interpretability.
453

454 It's worth noting that prediction accuracy of all KGML-ag models dropped without augmenting the training dataset in the fine-
455 tuning process (No.7-10 in Table 3 ~~blue region~~). Moreover, the maximum training epochs increased from 800 to 20000, which
456 resulted in overfitting on the small data set. This indicated that the data augmentation indeed helped the models become more
457 generalizable and gain better accuracy.
458
459 Experiments using zero initial values presented a significant drop in every variable's prediction accuracy (No.11-14 in Table
460 3 ~~yellow region~~). This indicated that the IMV initial values input into the KGML-ag-IMV modules of KGML-ag models
461 influenced not only the IMV prediction but also the $N_2O$ prediction of the KGML-ag-$N_2O$ module. This shows that there is
462 useful information transferred from IMVs in the KGML-ag-IMV module to the KGML-ag-$N_2O$ module.

463 **4 Discussion**

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

469 **4.1 Interpretability of KGML-ag**

470 The proposed KGML-ag models incorporate causal relations among $N_2O$ related variables/processes as shown in Fig. S10.
471 Managements, weather forcings and initial values of IMVs influence soil water, soil temperature and soil properties, which
472 influence the availability of $O_2$ and N as well as the microbe populations in soil, and further influence the nitrification and
473 denitrification rates. $N_2O$ is produced during both nitrification and denitrification when soil $O_2$ concentration is limited. Our
474 KGML-ag follows this hierarchical structure by designing KGML-ag-IMV modules representing the soil processes for IMVs
475 predictions (Fig. 1c-d).
476
477 To better explain the time series predictions of $N_2O$ flux (Fig. S1; Fig. 2-3), we separated the observations of each year into
478 three periods: leading period (before $N_2O$ increasing), increasing period (increasing to the peak) and decreasing period (peak
479 decreasing to near zero). During the leading period, both $NH_4^+$ and $CO_2$ were increasing immediately in the following few days
480 following urea N fertilizer application, indicating that urea was decomposing into $NH_4^+$ and $CO_2$ in soil water. With
481 accumulating $NH_4^+$ in soil, nitrification started producing $NO_3^-$ and consuming $O_2$. $N_2O$ didn't respond to the fertilizer
482 immediately due to enough $O_2$ in soil. Then when the soil became sufficiently hypoxic, $N_2O$ fluxes entered an increasing
483 period with $N_2O$ being produced by nitrification and denitrification processes. $CO_2$ fluxes were relatively low and $NH_4^+$ kept
484 decreasing during this period. Finally, when soil $NH_4^+$ was exhausted and $NO_3^-$ started decreasing due to denitrification, $N_2O$

485     fluxes then entered the decreasing period. $CO_2$ flux was related to urea decomposition during the leading period, and was more

486     closely related to $O_2$ demand in other periods. The KGML-ag predictions of $N_2O$ and IMV captured the three periods and

487     transition points, demonstrating the connections between those variables following the description as above (Fig. 3; Fig. S5).

488     Although KGML-ag1 obtained lower IMVs prediction accuracy compared to KGML-ag2, it captured the general trends and

489     was doing better for transitions, especially in $NH_4^+$ predictions. KGML-ag2 overfitted on the observations and ignored the

490     correlations between IMVs, which resulted in loss in pretrain knowledge, poorer performance in the leading period, and erratic

491     predictions in the period with missing observations (before day 25).

492     **4.2 Lessons for KGML-ag development**

493     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

494     ET, with complicated generation processes relying on the historical states. To develop a capable KGML model, we need to

495     carefully address three questions:

496

497     What kind of ML model is suitable for developing KGML? The answer could be determined by the dominant variation type

498     of the target variable in the data. If the dominant type is temporal variance, like flux variables in high temporal resolution (e.g.,

499     daily, or hourly), we should consider ML models with temporal dependency. RNN models such as GRU used in this study,

500     and CNN models such as casual CNN (Oord et al., 2016) can be good starting ML models. If the dominant type is spatial

501     variation, like variables in coarse temporal resolution (e.g., monthly or annually) but with high diversity due to soil property,

502     land cover and climate, we should consider ML models with the ability to deal with edges, hotpoints and categories, such as

503     CNN;

504

505     What physical/chemical constraints can be used to build KGML models? Although physical rules such as mass balance or

506     energy balance are conceptually straightforward and were proved capable of constraining KGML in predicting lake phosphorus

507     and temperature dynamics (Hanson et al., 2020; Read et al., 2019), they were excluded in this study according to our

508     preliminary analysis. The reason is that the mass balance equation of N in the agriculture ecosystem includes too many

509     unknown and unobservable components such as $N_2$ flux, $NH_3$ flux, N leaching, microbial N, plant N and soil/plant exchange,

510     which collectively introduce large uncertainties in balance equations and make them hard to be directly applied in the KGML-

511     ag framework. Other related physical (e.g., diffusion, solution) or chemical (e.g., nitrification, denitrification) processes cannot

512     be easily added into the KGML-ag structure as rules due to lack of understanding of the process. Instead, as mentioned in Sect.

513     2.2.4, we used hierarchical structure to enforce an architectural constraint and causal relations among variables, and pretraining

514     processes to infuse knowledge from *ecosys* to KGML-ag models.

515

516     How to involve PB models in the KGML development? An advanced PB model like *ecosys* built upon biophysical and

517     biochemical rules instead of empirical relations will be a good basis to learn the process, guide the structure and provide the

518 constraints for KGML. The generated synthetic data in this study helped us get some knowledge about variables such as their
519 general trends, dynamics and correlations. Such knowledge can be transferred to KGML models from synthetic data in the
520 pretraining process, which can reduce the efforts to collect large numbers of real-world observation data. Moreover, while
521 KGML shows great potential beyond PB models, we reckon that equally important for improving $N_2O$ modeling is to continue
522 improving our understanding of the related processes and mechanisms. Novel data collection and incorporating new
523 understanding into PB models (e.g., *ecosys*) could provide foundation to further empower KGML (see further discussion in
524 Sect. 4.3).

525

**4.3 Limitation and possible improvement**

527 First, the KGML-ag models in this study are limited by the available observed data. The mesocosm measurements of $N_2O$
528 fluxes ($16.9\pm11.7$ mg N m$^{-2}$ day$^{-1}$ during days of 45-60; Highest value is 71 mg N m$^{-2}$ day$^{-1}$) and $NO_3^-$ soil concentrations
529 ($59.3\pm20.7$ 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
530 observed by field studies (Fassbinder et al., 2013; Grant et al., 1999, 2006, 2008, 2016; Hamrani et al., 2020; Venterea et al.,
531 2011). Some IMVs with high feature importance scores (e.g., $O_2$ flux, $N_2$ flux) or at different depths (e.g., soil $NO_3^-$ at 5 cm
532 depth, VWC at 5 cm depth), and data out of growing seasons are not included. The direct consequences are that some important
533 processes cannot be well represented by the current KGML-ag (e.g., $O_2$ demand and N availability for nitrification and
534 denitrification). Further improvement of KGML should consider three categories of data: target variable $N_2O$ flux, IMVs and
535 basic inputs (Fig. 1a). For $N_2O$ flux observation, we lack sub-hourly to sub-daily observations to capture the hot moment of
536 emission during 0-30 days after N fertilizer applications. Besides, the non-growing season can provide 35-65% of the annual
537 direct $N_2O$ emissions from seasonally frozen croplands and lead to a 17–28 % underestimate of the global agricultural $N_2O$
538 budget if ignoring its contribution (Wagner-Riddle et al., 2017), but we can barely find observations from non-growing
539 seasons. For IMVs, we found oxygen demand indicator (e.g., $O_2$ concentration or flux, $CO_2$ flux, $CH_4$ flux), N mass balance
540 related variables (e.g., $N_2$ flux, soil $NO_3^-$, soil $NH_4^+$, N leaching) and soil water and temperature, can be used to better constrain
541 the processes and therefore improve the KGML performance. Rohe et al. (2021) also indicated the importance of $O_2$, $CO_2$ and
542 $N_2$ soil fluxes for $N_2O$ predictions. In addition, the layerwise soil observations (e.g., soil $NO_3^-$, soil VWC) at 0-30 cm depth
543 can be used to significantly improve the KGML model quality, according to our feature importance analysis (Fig. S2a).
544 Moreover, continuous monitoring on these variables during the whole year is preferred rather than only during the growing
545 season, since $N_2O$ flux is largely influenced by previous states. To apply the KGML-ag to large scale, other observational data
546 including basic inputs of soil/crop properties (e.g., soil bulk density, pH, crop type), management information (e.g., fertilizer,
547 irrigation, tillage) and weather forcings along with $N_2O$ flux observations are critical for fine-tuning and validating the
548 developed KGML-ag and therefore explicitly simulating the $N_2O$ or IMVs dynamics under specific conditions. Recent
549 advances in remote sensing and machine learning have enabled estimating these variables with high-resolution at a large scale
550 (Peng et al., 2020)

17

551

552  Second, the physical/chemical constraints can be more comprehensive in KGML-ag models. Although current KGML-ag
553  models are well-initialized with *ecosys* synthetic data and constrained by causal relations of processes with hierarchical
554  structure, the predicted $N_2O$ flux and IMVs can still violate some basic physical rules like mass balance. As we discussed in
555  Sec. 4.2, it will be challenging to add physical rules like mass balance equation for N in a complicated agriculture ecosystem
556  due to data limitations such as missing observations on certain key variables. Using inequalities instead of equations for mass
557  balance may be one alternative solution. For example, we could use ReLU to add in a limitation for N mass balance residues
558  which are calculated from known terms not larger than an empirical static value. Besides, better understanding of processes in
559  the N cycle from fieldworks and lab experiments can also help us design new constraints. This limitation is also partially
560  related to the data limitation and can be overcomed by involving more complete $N_2O$ data to introduce more powerful
561  constraints to KGML-ag.

562

563  Third, the KGML-ag currently are suffering from dealing with physical/chemical boundary transitions. Boundary transitions
564  are common in the real world, such as phase change, volume solubility, and soil porosity etc. A detailed PB model generally
565  coded plenty of "if/else/switch" statements inside to deal with the boundaries. But KGML-ag models based on the GRU are
566  better at capturing continuous changes, rather than discrete changes. One solution is to include data with boundary information.
567  In this study, involving IMVs like $O_2$, $CO_2$ and $N_2$, which already have boundary information like water freezing point, N pool
568  volumes and other complicated boundaries related to soil/crop properties, can significantly improve the model performance.
569  The data with boundary information could be continuous observation or estimated value from existing data. By using initial
570  values to predict IMVs, KGML-ag in this study can partially solve the boundary transition problem when observation data is
571  limited. Another solution is designing new structures of KGML-ag, such as combining ReLU function or including CNN
572  model which are robust for discrete situations to the RNN models, or designing new constraints to limit the model working
573  within the thresholds.

574

575  Finally, at the current stage we can not claim to have completely opened the black box of KGML-ag, but this framework is a
576  significant step towards this goal. For example, some ideas implemented in our study, such as using pretraining to transfer
577  knowledge from a PB model to a ML model, incorporating causal relations by hierarchical structure, predicting IMVs for
578  tracking middle changes and using initial values as input to reduce data demand, would shed light on the future KGML-ag
579  framework improvement. Besides, we acknowledge the importance of further testing the KGML-ag over completely
580  independent datasets, but results presented in this manuscript are sufficient to justify the power of KGML as a framework. The
581  mesocosm experiment data we used in this study has provided a comprehensive set of inputs and intermediate variables in
582  addition to the output of $N_2O$ fluxes, thus serving as a unique testbed. We expect to further validate and refine our KGML-ag
583  model our validation results will be more solid once more gold standard data of $N_2O$ fluxes along with other relevant inputs
584  and intermediate variables become publicly available. Moreover, incorporating more and more domain knowledge into

18

585 KGML-ag will be possible~~inevitable~~ for~~in~~ further improvement, but we don't think KGML-ag will become inefficient as it
586 becomes more like the PB model. In fact, to efficiently emulate~~surrogate~~ components of PB models has been proposed as a
587 research frontier in hybrid modeling for earth system science (Reichstein et al., 2019; Irrgang et al., 2021), with latest advances
588 occurring in weather forecasts (Bauer et al., 2021). By using a hybrid model, computationally inefficient components of PB
589 can be identified one by one, and be replaced with more efficient ML-based surrogates to eventually obtain the most efficient
590 model. Further KGML-ag model development will also need to balance efficiency, accuracy and interpretability.

## 5 Conclusions

592 In this study, two KGML-ag models have been developed, validated, and tested for agricultural soil $N_2O$ flux prediction using
593 synthetic data generated by the PB model *ecosys* and observational data from a mesocosm facility. The results show that
594 KGML-ag models can outperform PB and pure ML models in $N_2O$ prediction in not only magnitude (KGML-ag1 $r^2 = 0.81$ vs
595 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
596 $\Delta r^2 = 0.08$). KGML-ag can also defeat the PB model *ecosys* in efficiency by completing *ecosys*'s half-hour job within a few
597 seconds. Compared to ML models, KGML-ag models can better represent complex dynamics and high peaks of $N_2O$ flux.
598 Moreover, with IMV predictions and hierarchical structures, KGML-ag models can provide biogeophysical/chemical
599 information about key processes controlling $N_2O$ fluxes, which will be useful for interpretable forecasting and developing
600 mitigation strategies. Data demand for the KGML-ag models is significantly reduced due to involving IMV initial values and
601 pretrain processes with synthetic data. This study demonstrated that the potential of KGML-ag application in the complex
602 agriculture ecosystem is high and illustrates possible pathways of KGML-ag development for similar tasks. Further
603 improvement of our KGML-ag models can involve general principles to further constrain the predictions through loss functions
604 or architectures, but call for more detailed, high temporal resolution $N_2O$ observation data from field measurements.

## Code and Data Availability

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

## Author contributions

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

613 **Competing interests**

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

615 **References**

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

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

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

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

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

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

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

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

634 Del Grosso, S. J., Parton, W. J., Mosier, A. R., Ojima, D. S., Kulmala, A. E., & Phongpan, S.: General model for $N_2O$ and N2
635 gas emissions from soils due to dentrification, Global biogeochemical cycles, 14(4), 1045-1060, 2020.

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

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

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

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

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

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

649 Grant, R. F., & Pattey, E.: Modelling variability in $N_2O$ emissions from fertilized agricultural fields, Soil Biology and
650 Biochemistry, 35(2), 225-243, 2003.

651 Grant, R. F., & Pattey, E.: Temperature sensitivity of $N_2O$ emissions from fertilized agricultural soils: Mathematical modeling
652 in ecosys. Global biogeochemical cycles, 22(4), 2008.

653 Grant, R. F., Neftel, A., & Calanca, P.: Ecological controls on $N_2O$ emission in surface litter and near-surface soil of a managed
654 grassland: modelling and measurements, Biogeosciences, 13(12), 3549-3571, 2016.

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

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

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

661 Holzworth, D. P., Huth, N. I., deVoil, P. G., Zurcher, E. J., Herrmann, N. I., McLean, G., ... & Keating, B. A.: APSIM–
662 evolution towards a new generation of agricultural systems simulation, Environmental Modelling & Software, 62, 327-350,
663 2014.

664 Irrgang, C., Boers, N., Sonnewald, M., Barnes, E. A., Kadow, C., Staneva, J., & Saynisch-Wagner, J.: Towards neural Earth
665 system modelling by integrating artificial intelligence in Earth system science. Nature Machine Intelligence, 3(8), 667-674,
666 2021.

667 Jia, X., Willard, J., Karpatne, A., Read, J. S., Zwart, J. A., Steinbach, M., & Kumar, V.: Physics-guided machine learning for
668 scientific discovery: An application in simulating lake temperature profiles, ACM/IMS Transactions on Data Science, 2(3), 1-
669 26, 2021.

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

673 Karpatne, A., Atluri, G., Faghmous, J. H., Steinbach, M., Banerjee, A., Ganguly, A., ... & Kumar, V.: Theory-guided data
674 science: A new paradigm for scientific discovery from data, IEEE Transactions on knowledge and data engineering, 29(10),
675 2318-2331, 2017.

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

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

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

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

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

686    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.,
687    and Frie, A. L.: Response of nitrous oxide emissions to future changes in precipitation and individual rain events, Journal of
688    Environmental Quality, In review, 2021

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

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

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

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

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

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

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

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

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

710    Saha, D., Basso, B., & Robertson, G. P.: Machine learning improves predictions of agricultural nitrous oxide ($N_2O$) emissions
711    from intensively managed cropping systems, Environmental Research Letters, 16(2), 024004, 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  Solazzo, E., Crippa, M., Guizzardi, D., Muntean, M., Choulga, M., & Janssens-Maenhout, G.: Uncertainties in the Emissions
716  Database for Global Atmospheric Research (EDGAR) emission inventory of greenhouse gases, Atmospheric Chemistry and
717  Physics, 21(7), 5655-5683, 2021.
718  Syakila, A., & Kroeze, C.: The global nitrous oxide budget revisited, Greenhouse gas measurement and management, 1(1),
719  17-26, 2011.
720  Thompson, R. L., Lassaletta, L., Patra, P. K., Wilson, C., Wells, K. C., Gressent, A., ... & Canadell, J. G.: Acceleration of
721  global $N_2O$ emissions seen from two decades of atmospheric inversion, Nature Climate Change, 9(12), 993-998, 2019.
722  Thornley, J. H., & France, J.: Mathematical models in agriculture: quantitative methods for the plant, animal and ecological
723  sciences, Cabi, 2007.
724  Tian, H., Xu, R., Canadell, J. G., Thompson, R. L., Winiwarter, W., Suntharalingam, P., ... & Yao, Y.: A comprehensive
725  quantification of global nitrous oxide sources and sinks, Nature, 586(7828), 248-256, 2020.
726  Venterea, R. T., Maharjan, B., & Dolan, M. S.: Fertilizer source and tillage effects on yield-scaled nitrous oxide emissions in
727  a corn cropping system. Journal of Environmental Quality, 40(5), 1521-1531, 2011.
728  Wagner-Riddle, C., Congreves, K. A., Abalos, D., Berg, A. A., Brown, S. E., Ambadan, J. T., ... & Tenuta, M.: Globally
729  important nitrous oxide emissions from croplands induced by freeze–thaw cycles, Nature Geoscience, 10(4), 279-283, 2017.
730  Willard, J., Jia, X., Xu, S., Steinbach, M., & Kumar, V.: Integrating Scientific Knowledge with Machine Learning for
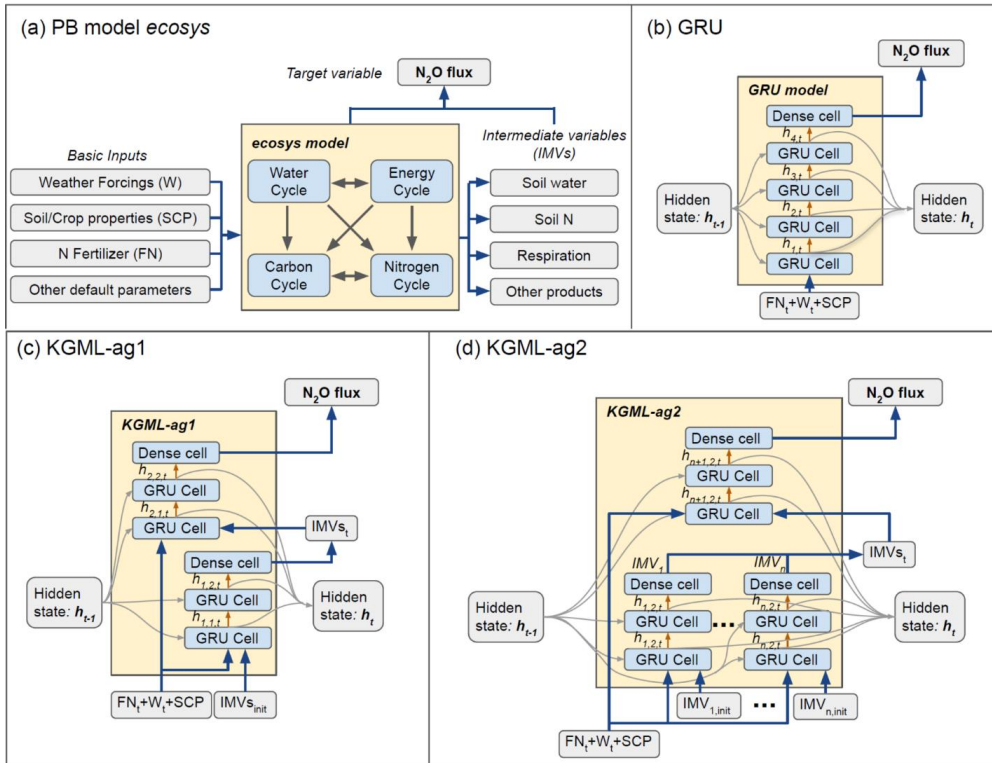731  Engineering and Environmental Systems, arXiv preprint arXiv:2003.04919, 2020.
732  Yang, Y., Liu, L., Zhou, W., Guan, K., Kim, T., Tang, J., Peng, B., Zhu, P., Grant, R. F., Griffis, T. J., Jin, Z.: Distinct driving
733  mechanisms of non growing season N2O emissions call for spatial-specific mitigation strategies in the US Midwest. One
734  Earth. Submitted, 2022.
735  Zhang, Y., & Niu, H.: The development of the DNDC plant growth sub-model and the application of DNDC in agriculture: a
736  review, Agriculture, Ecosystems & Environment, 230, 271-282, 2016.
737  Zhang, Y., Li, C., Zhou, X., & Moore III, B.: A simulation model linking crop growth and soil biogeochemistry for sustainable
738  agriculture, Ecological modelling, 151(1), 75-108, 2002.
739  Zhou, W., Guan, K., Peng, B., Tang, J., Jin, Z., Jiang, C., ... & Mezbahuddin, S.: Quantifying carbon budget, crop yields and
740  their responses to environmental variability using the ecosys model for US Midwestern agroecosystems. Agricultural and
741  Forest Meteorology, 307, 108521, 2021.
742
743

**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 $CO_2$ flux, soil $NO_3^-$ concentration, soil $NH_4^+$ concentration, and soil volumetric water content (VWC).**

Figure 2: Leave-one-out cross validation of time series of N$_2$O flux (mg N m$^{-2}$ day$^{-1}$) time series predicted by thecomparisons among pure non-pretrained GRU modelpredictions (blue line) and, KGML-ag1 modelpredictions (red line). and oObservations are shown as (black line-dots,) from cross-validation. The N$_2$O flux unit is mg N m$^{-2}$ day$^{-1}$. Validation results for each chamber were based on out-of-sample predictions by models trained by other five chambers.Leave-one-out cross validation (LOOCV) method was used to train/validate the models. Only validation results were presented and each chamber validation results were from models trained by other five chambers.

Formatted: Font color: Auto

**(a)** $CO_2$ **flux**
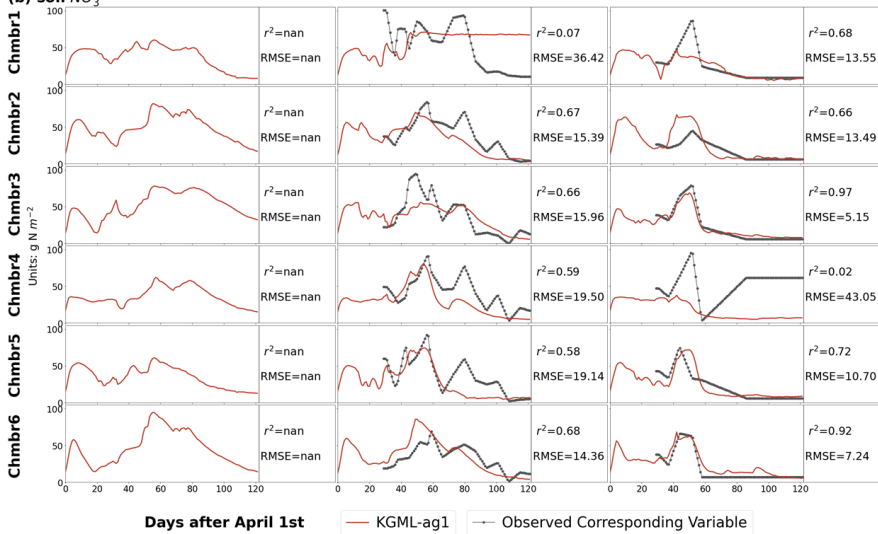
**(b)** soil $NO_3^-$

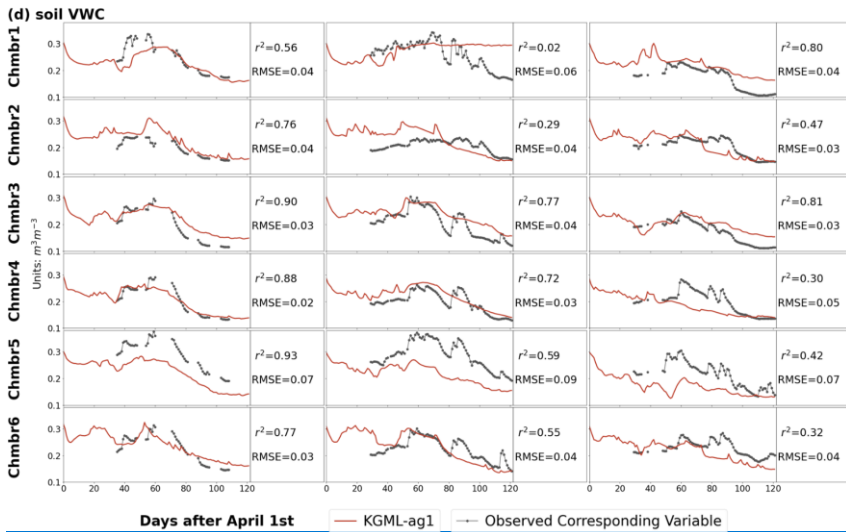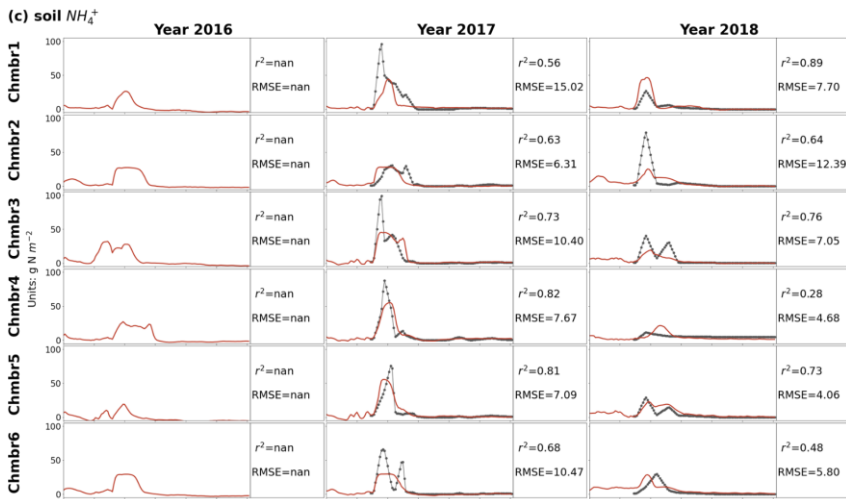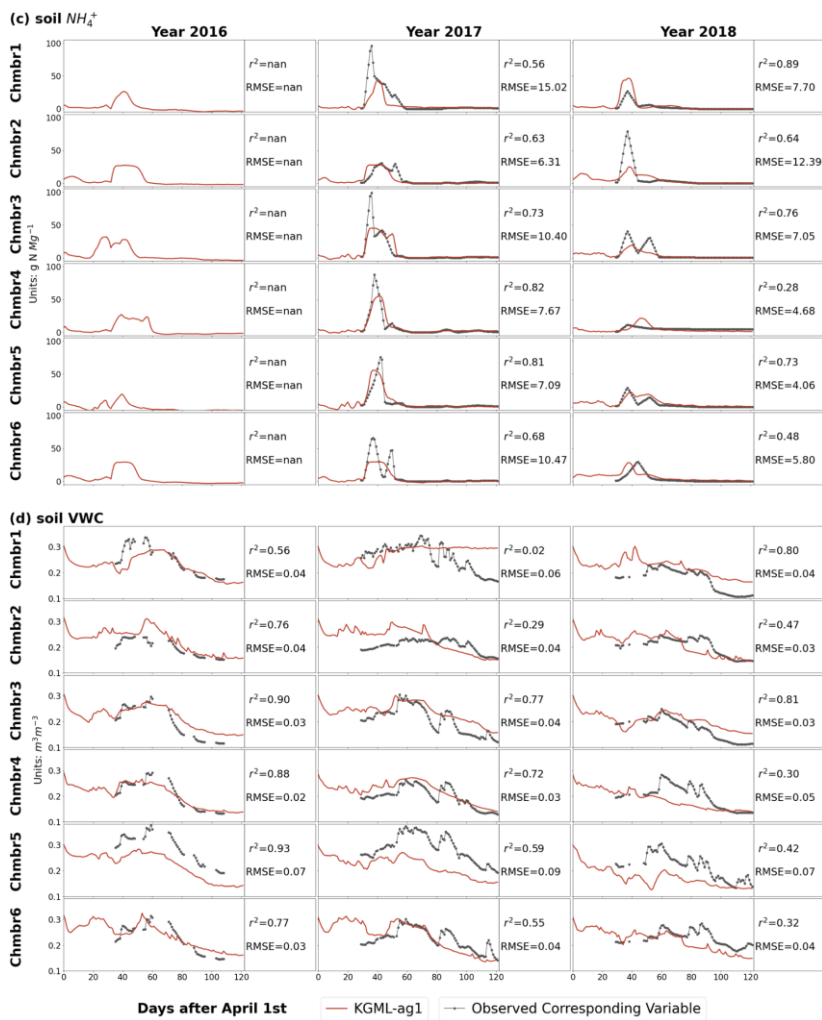Days after April 1st · ─── KGML-ag1 · ─┼─ Observed Corresponding Variable
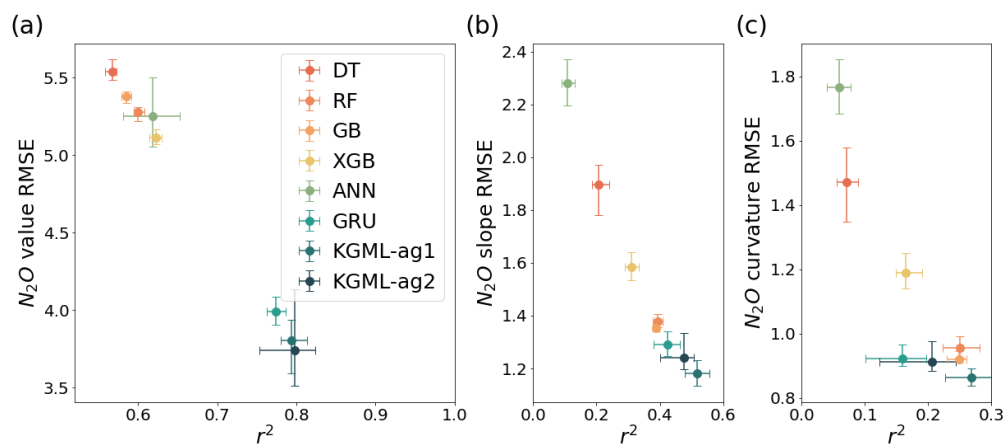
761

**Figure 3: Leave-one-out cross validation of time series of ~~I~~MVs predict~~edion byfrom~~ KGML-ag1 model (red line). Observations are shown as black line-dots.~~The black-dot line represents observations and the red line represents the results from KGML-ag1.~~ 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 $CO_2$ flux and soil $NO_3^-$ concentration units are g C m$^{-2}$ day$^{-1}$ and g N ~~Mgm$^{-1}$~~, respectively.

(c) soil $NH_4^+$



(d) soil VWC



Days after April 1st — KGML-ag1 — Observed Corresponding Variable

**(c) soil** $NH_4^+$

**(d) soil VWC**

Days after April 1st    — KGML-ag1    -·- Observed Corresponding Variable

**Figure 3 Contd.:** ~~Leave-one-out cross validation of time series of~~ IMVs predict~~edion~~ by~~from~~ KGML-ag1 model (red line). Observations are shown as black line-dots. ~~, The black-dot line represents observations and the red line represents the results from KGML-ag1.~~ 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 ~~Mgm~~$^{-1.2}$ and $m^3$ $m^{-3}$, respectively.
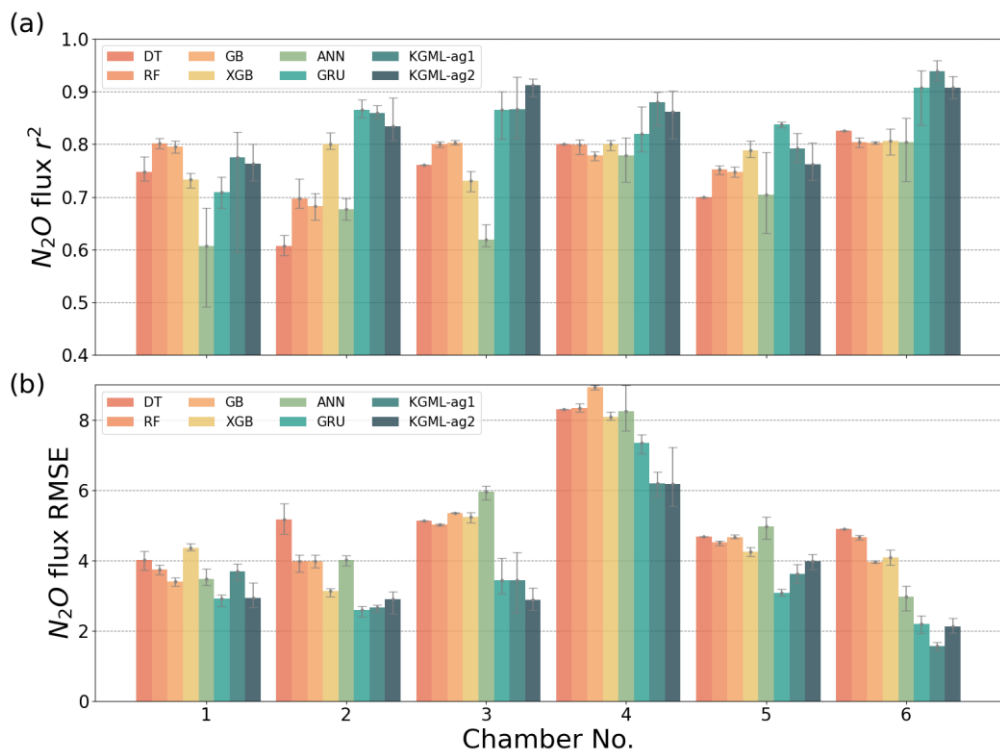
777

**Figure 4: The comparisons of overall prediction accuracy** <u>from leave-one-out cross validation</u> **for N₂O value (a), 1st order gradient (slope, b) and 2nd order gradient (curvature, c) between four tree-based ML models (DT, RF, GB and XGB), two deep learning models (ANN and GRU) and KGML-ag models.** <u>The overall performances were calculated by comparing out-of-sample predictions (each chamber's predictions were from models trained by other five chambers) from all validated chambers with observations.</u> **Different color symbols represent the different models. The x- and y-error bars are coming from the maximum and minimum scores of ensemble experiments. The dot represents the mean score of the ensemble experiments.**

784

Figure 5: The comparisons of N₂O flux prediction accuracy r² (a) and (b) RMSE from Leave-one-out cross validation, between four tree-based ML models (DT, RF, GB and XGB), two deep learning models (ANN and GRU) and KGML-ag models in six6 chambers. Validation results for each chamber were based on out-of-sample predictions by models trained by other five chambers. The gray error bars are coming from the maximum and minimum scores of ensemble experiments.

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.**

| | | | N$_2$O | | CO$_2$ | | NO$_3^-$ | | NH$_4^+$ | | VWC | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| No. | Pretrain Model | Input Feature N | 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 |

[a]The empty slot indicates that the model does not predict that variable.


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$^2$ [a]** | 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 |
| **ΔRMSE[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 |

[a]Leave-one-out cross validation results for each chamber were based on out-of-sample predictions by models trained by other five chambers. The "All data" performances were calculated by comparing out-of-sample predictions from all validated chambers with observations. The difference of r$^2$ (Δr$^2$), and difference of RMSE (ΔRMSE, units are mg N m$^{-2}$ day$^{-1}$, mg N m$^{-2}$ day$^{-2}$, mg N m$^{-2}$ day$^{-3}$ for N$_2$O value, 1st order gradient and 2nd order gradient, respectively) were calculated by values from KGML-ag1 minus values from GRU.

[b]Results from different time windows of different chambers during the period of April 1st-July31st (Days1-122) were detected.

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

| No. | Retrain Model | Experiment | N₂O $r^{2,c}$ | N₂O RMSE[c] | N₂O 1st order gradient $r^2$ | N₂O 1st order gradient RMSE | N₂O 2nd order gradient $r^2$ | N₂O 2nd order gradient RMSE | $CO_2$ $r^2$ | $CO_2$ NRMSE | $NO_3^-$ $r^2$ | $NO_3^-$ NRMSE | $NH_4^+$ $r^2$ | $NH_4^+$ NRMSE | VWC $r^2$ | VWC 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 |

[a][No.1-6]Gray region includes the experiments with original simulation settings as described in Sec. 2 and bold valuesdark gray refers to the baseline GRU simulation; No.7-10Blue region includes the experiments without data augmentation during the finetuning process; And No. 11-14yellow region includes the experiments of replacing original IMV initial values with zeros.

[b]The empty slot indicates that the model does not predict that variable.

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