

# Simulating N<sub>2</sub>O Emission from Fertilized Mesocosm Using Knowledge Guided Machine Learning

Licheng LIU<sup>1</sup>, Shaoming Xu<sup>1</sup>, Zhenong Jin<sup>2</sup>, Jinyun Tang<sup>3</sup>, Timothy Griffis<sup>4</sup>, Kaiyu Guan<sup>5</sup>, Alexander Frie<sup>6</sup>, Taegon Kim<sup>1</sup>, Bin Peng<sup>5</sup>, Yufeng Yang<sup>1</sup>, Wang Zhou<sup>5</sup>, and Vipin Kumar<sup>1</sup>

<sup>1</sup>University of Minnesota Twin Cities

<sup>2</sup>University of Minnesota-Twin Cities

<sup>3</sup>Lawrence Berkeley Natl Lab

<sup>4</sup>Univ Minnesota

<sup>5</sup>University of Illinois at Urbana Champaign

<sup>6</sup>University of California Riverside

November 24, 2022

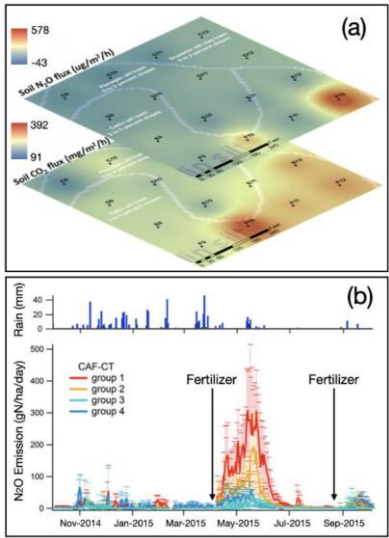
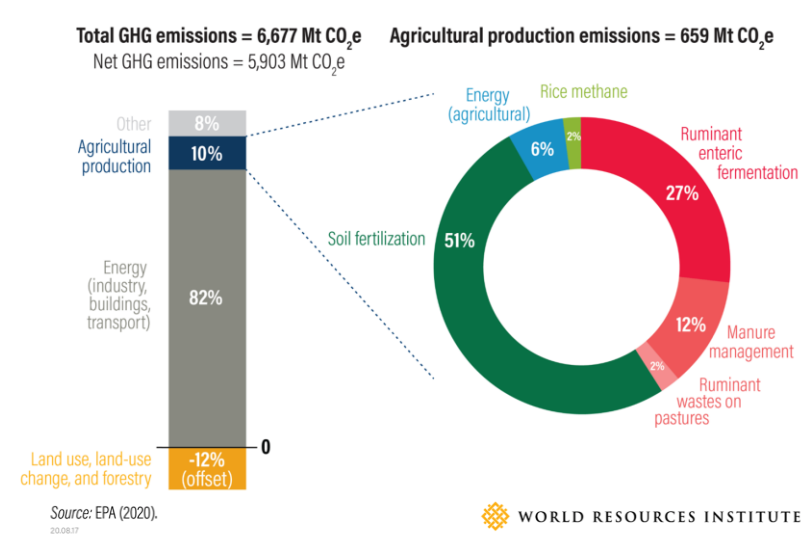
## Abstract

Nitrous oxide (N<sub>2</sub>O) is one of the important greenhouse gases (GHGs), with its global warming potential 265 times greater than that of carbon dioxide (CO<sub>2</sub>). About 60% of the anthropogenic N<sub>2</sub>O emission is from agriculture production. To date, estimating N<sub>2</sub>O emissions from cropland remains a challenging task because the related microbial origin processes (e.g. incomplete nitrification and denitrification) are controlled by a diverse factors of climate, soil, plant and human activities. In this study, we developed a ML model with physical/biogeochemical domain knowledge, namely knowledge guided machine learning (KGML), for simulating daily N<sub>2</sub>O fluxes from the agriculture ecosystem. The Gated Recurrent Unit (GRU) was used as the basis to build the model structure. A range of ideas have been implemented to optimize the model performance, including 1) hierarchical structure based on variable causal relations, 2) intermediate variable (IMV) prediction and transfer, 3) inputting IMV initials for constraints, 4) model pretrain/retrain, and 5) multitask learning. The developed KGML was pre-trained by millions of synthetic data generated by an advanced PB model, *ecosys*, and then re-trained by observations from six mesocosm chambers during three growing seasons. Six other pure ML models were developed using the same data from mesocosm chambers to serve as the benchmark for the KGML model. The results show that KGML can always outperform the PB model in efficiency and ML models in prediction accuracy of capturing N<sub>2</sub>O flux magnitude and dynamics. Besides, the reasonable predictions of IMVs increase the interpretability of KGML. We believe the footprint of KGML development in this study will stimulate a new body of research on interpretable machine learning for biogeochemistry and other related geoscience processes.

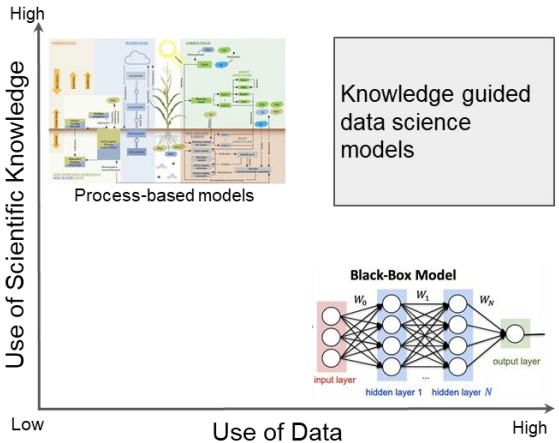
# Knowledge Guided Machine Learning for Simulating Agricultural N<sub>2</sub>O Emission

Presenter: Licheng Liu; [lichengl@umn.edu](mailto:lichengl@umn.edu)

## Motivation:



Waldo et al. 2019



❑ Fertilizer use accounts for 51% of the ag emissions, largely in forms of **N<sub>2</sub>O**, **265x more powerful than CO<sub>2</sub>** as a GHG

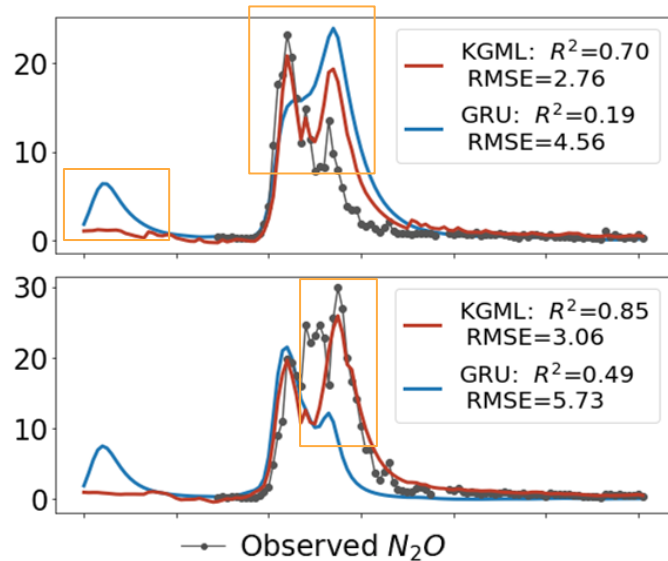
❑ Hard to estimate due to **hot spots, hot moment** of N<sub>2</sub>O fluxes

❑ KGML model can take full advantage of data without ignoring the treasure of accumulated scientific knowledge

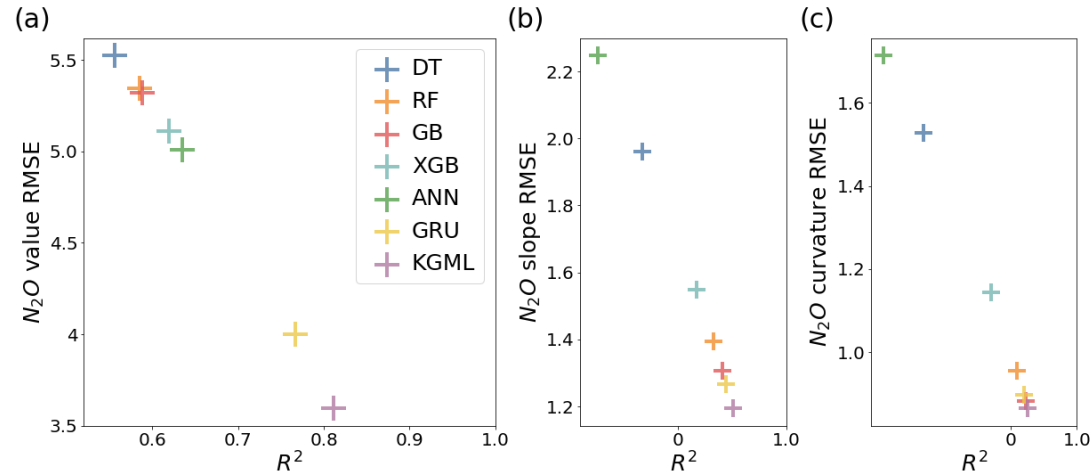
# Knowledge Guided Machine Learning for Simulating Agricultural N<sub>2</sub>O Emission

Presenter: Licheng Liu; [lichengl@umn.edu](mailto:lichengl@umn.edu)

## Key results:



One example of KGML model comparing to GRU model in mesocosm experiment data

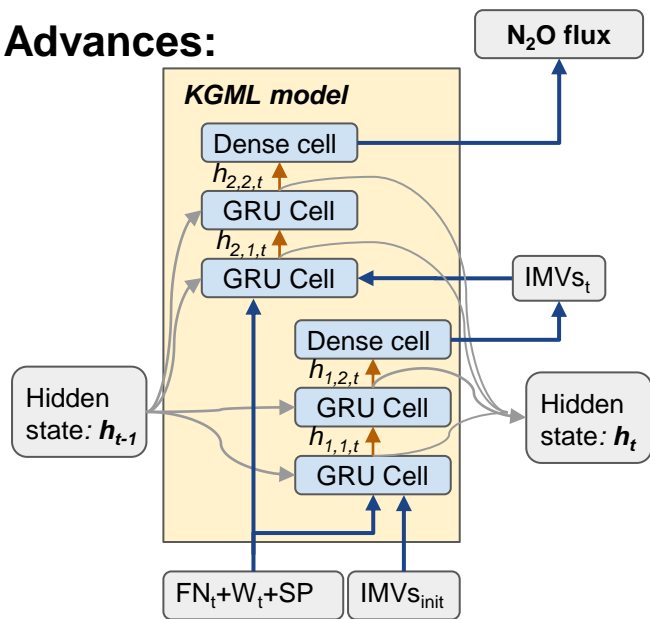


- ❑ KGML (purple) outperform all other ML models
- ❑ This is mainly because (1) pre-training using synthetic data, (2) knowledge guided architecture, (3) knowledge guided initial values

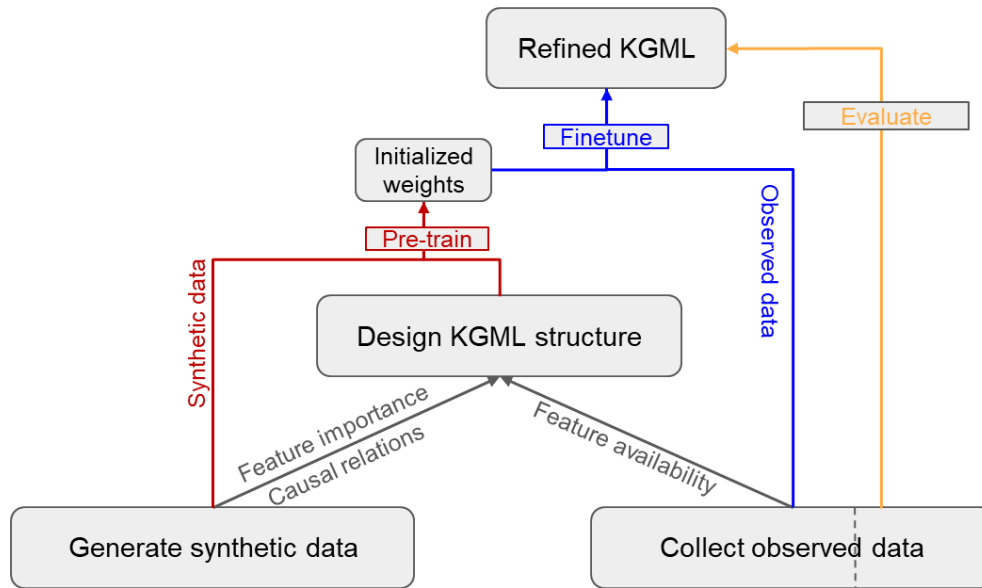
# Knowledge Guided Machine Learning for Simulating Agricultural N<sub>2</sub>O Emission

Presenter: Licheng Liu; [lichengl@umn.edu](mailto:lichengl@umn.edu)

## Advances:



KGML model structure



KGML model development workflow

- ❑ High performance
- ❑ Low data demand
- ❑ Flexible structure
- ❑ Structure and workflow can be easily transfer to other similar geoscience tasks!



## Knowledge Guided Machine Learning for Simulating Agricultural N<sub>2</sub>O Emission

**Licheng Liu**<sup>1</sup>, Shaoming Xu<sup>2</sup>, Zhenong Jin<sup>1,3</sup>, Jinyun Tang<sup>4</sup>, Timothy J Griffis<sup>5</sup>, Kaiyu Guan<sup>6,7</sup>, AlexanderLee Frie<sup>5</sup>, Taegon Kim<sup>1</sup>, Bin Peng<sup>6,7</sup>, Xiaowei Jia<sup>8</sup>, Yufeng Yang<sup>1</sup>, Wang Zhou<sup>6</sup> and Vipin Kumar<sup>2</sup>

<sup>1</sup>Department of Bioproducts & Biosystems Engineering, University of Minnesota

<sup>2</sup>Department of Computer Science, University of Minnesota

<sup>3</sup>Institute on the Environment, University of Minnesota

<sup>4</sup>Earth and Environmental Sciences Area, Lawrence Berkeley National Laboratory

<sup>5</sup>Department of Soil, Water & Climate, University of Minnesota

<sup>6</sup>Department of Natural Resources and Environmental Sciences, University of Illinois at Urbana-Champaign

<sup>7</sup>National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign

<sup>8</sup>Department of Computer Science, University of Pittsburgh

Contact: [lichengl@umn.edu](mailto:lichengl@umn.edu)

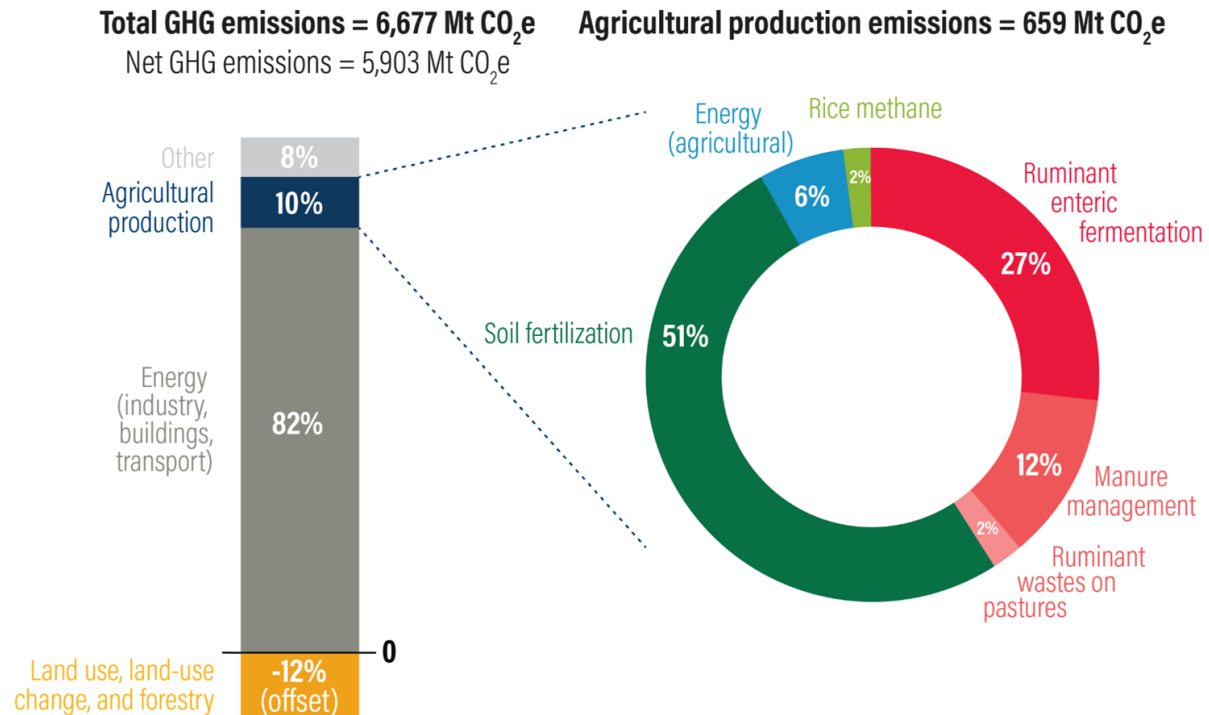
*AGU Fall meeting 2021*



UNIVERSITY OF MINNESOTA  
Driven to Discover®

# Agriculture and Climate Change

Agriculture contributes a quarter of global greenhouse gas (GHG) emissions that are causing climate change: **~14%** directly from agricultural activities and **~10%** through land use change.



Source: EPA (2020),  
20.08.17



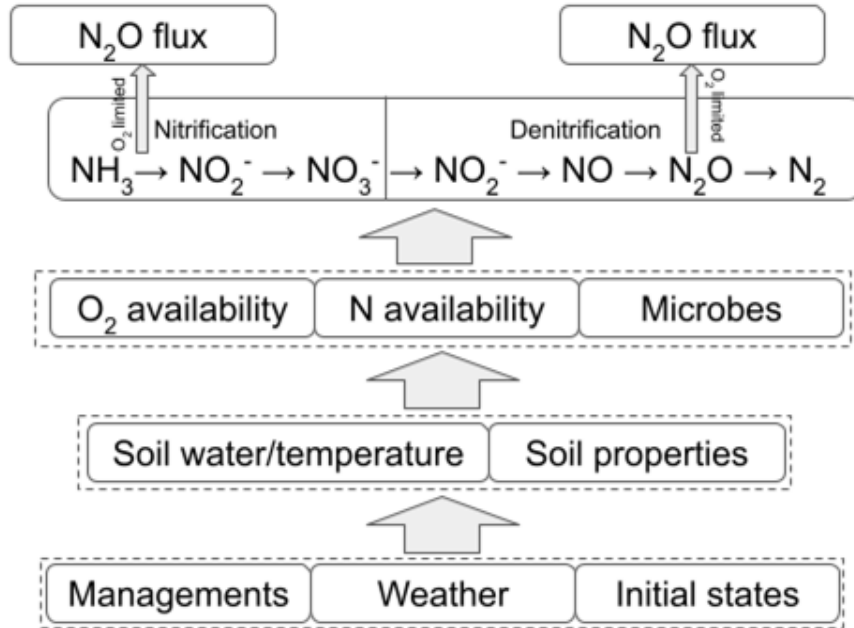
WORLD RESOURCES INSTITUTE

- ❑ Agricultural production accounts for **10%** of U.S. GHG emissions in 2018
- ❑ Fertilizer use accounts for 51% of the ag emissions, largely in forms of **N<sub>2</sub>O**, **265x more powerful than CO<sub>2</sub>** as a GHG
- ❑ Over applying fertilizer also causes **water** & **air** pollution and **land** degradation

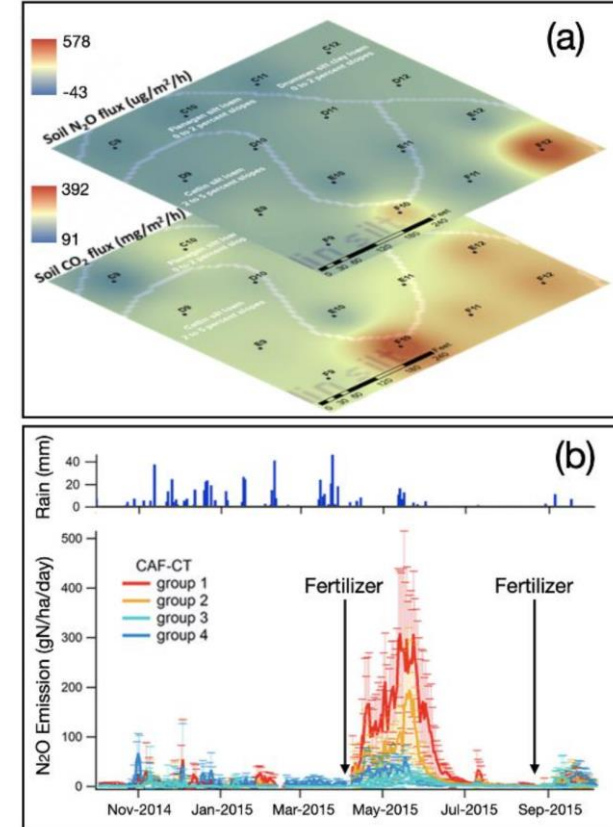


# Why estimating $N_2O$ is so hard?

Soil nitrous oxide ( $N_2O$ ) emissions are highly variable in space and time due to dynamic controls by a range of biotic and abiotic factors.



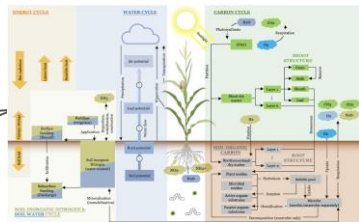
Hot spots, hot moment of  $N_2O$  fluxes



# Opportunities with knowledge guided machine learning

Contain more knowledges but also many empirical parameters

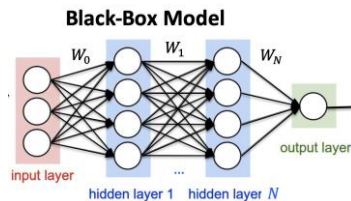
Use of Scientific Knowledge



Process-based models

Knowledge guided data science models

Take full advantage of data without ignoring the treasure of accumulated scientific knowledge



Require large number of data to train, often work as black-box

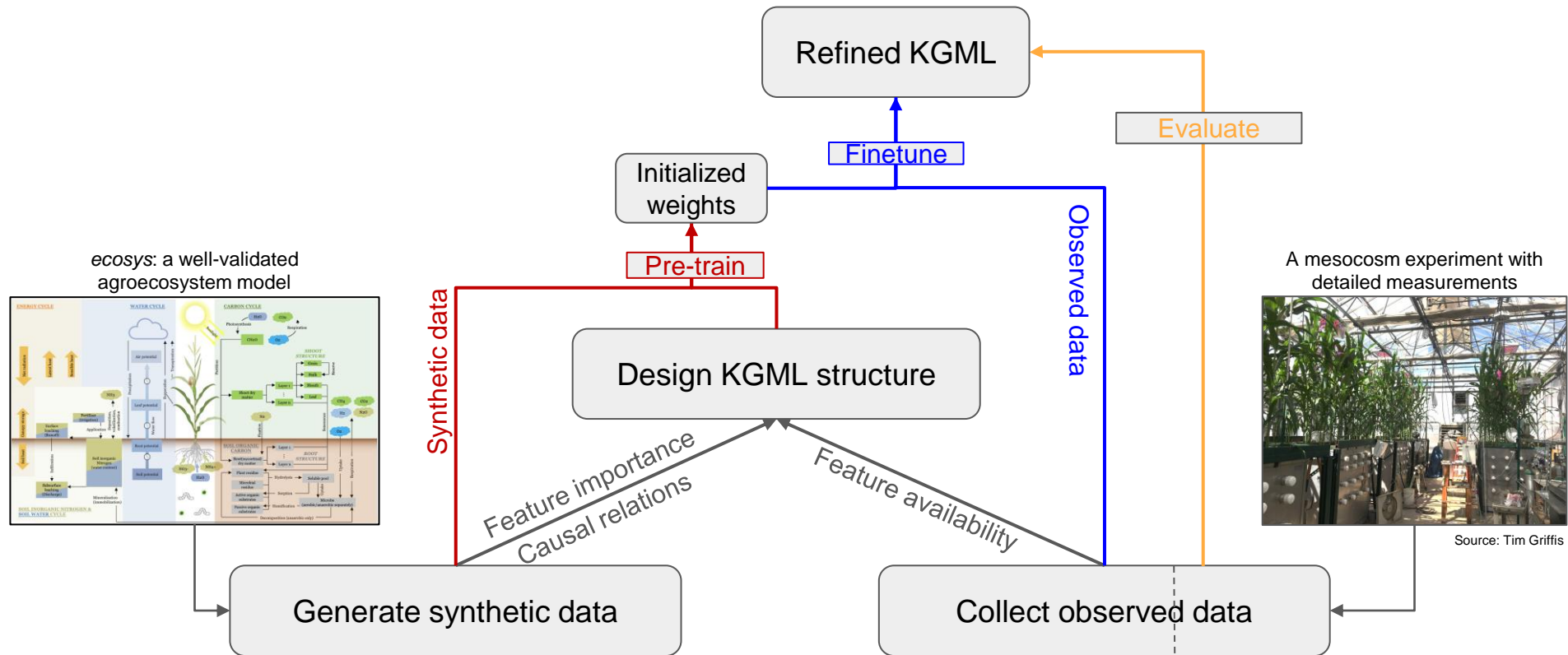
Use of Data

Low

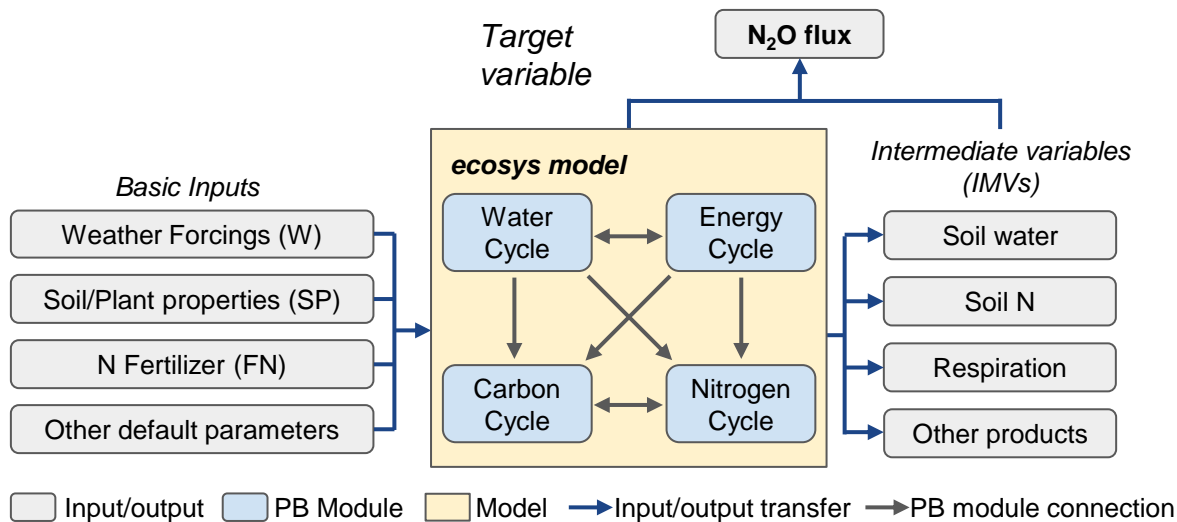
High



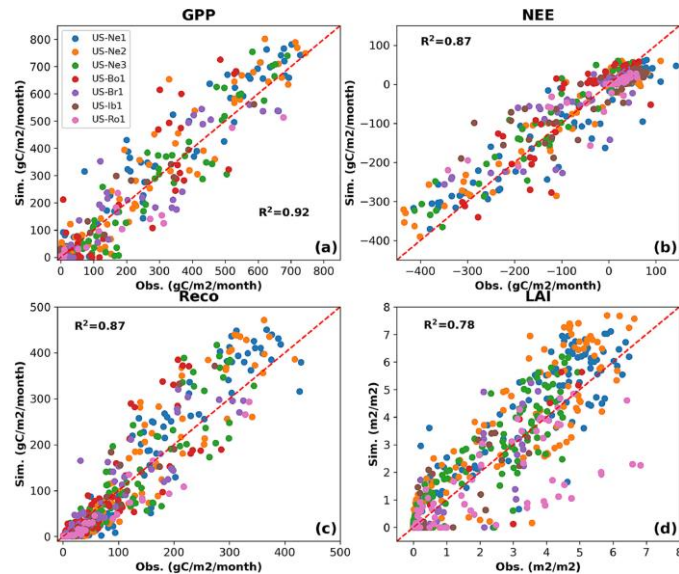
# Experiments overview



# Generate synthetic data from an advanced agroecosystem model, *Ecosys*



## Validation at 7 US FluxNet sites for ag



## Ecosys simulation:

- 99 random sites from Illinois, Indiana and Iowa
- 18 years simulations
- Over 4 million synthetic data samples

Zhou et al. (2021)

# Observed N<sub>2</sub>O fluxes from mesocosm experiments

## Experiment setup:

- Growing seasons during 2016-2018
- 6 chambers with different precipitation treatment
- N<sub>2</sub>O flux was measured by Teledyne M320EU Analyzer in automatic chamber

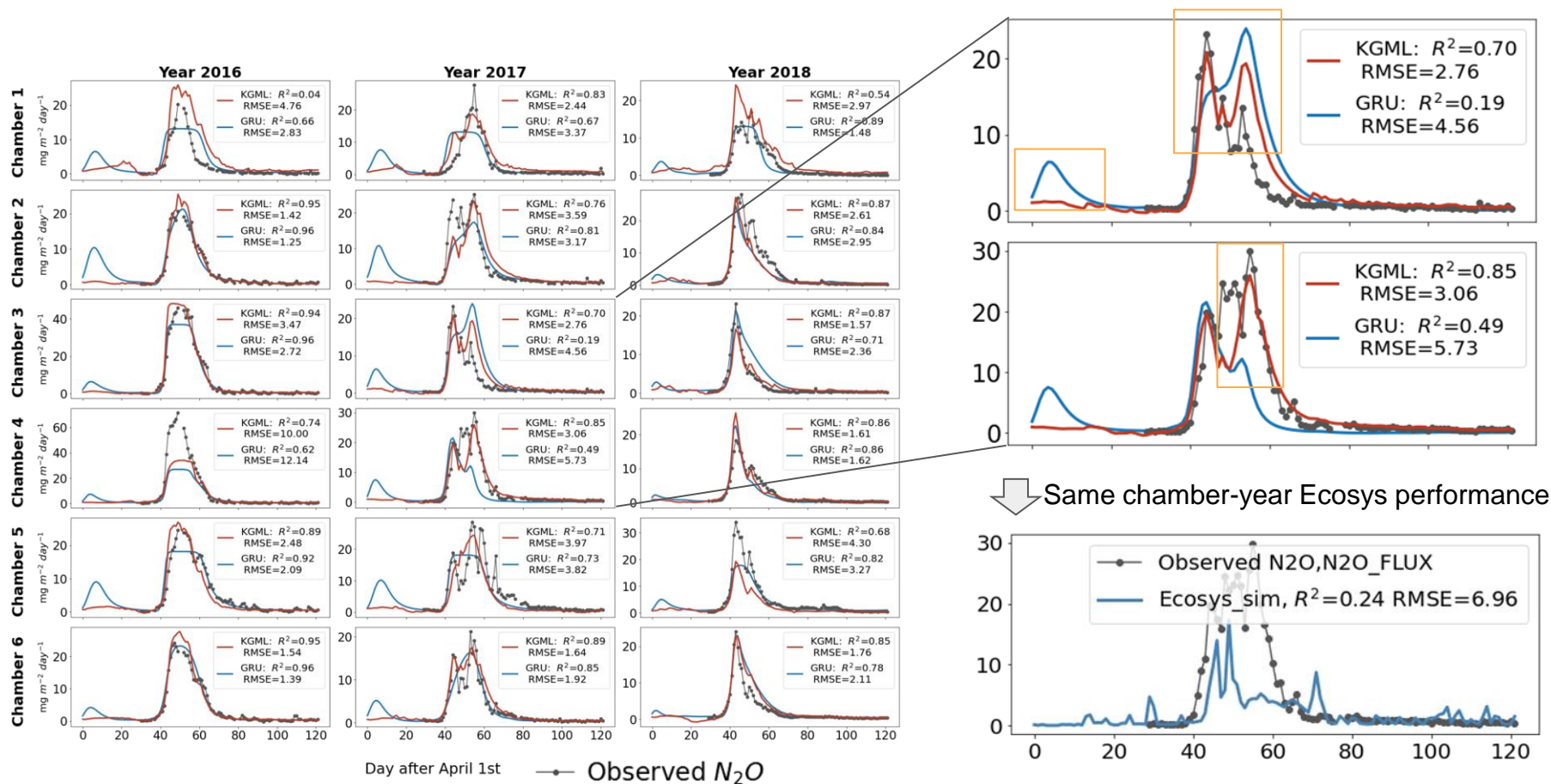
## Available variables:

- Controlled weather conditions
- Hourly N<sub>2</sub>O fluxes, CO<sub>2</sub> fluxes, and soil moisture at 15 cm depth
- Weekly soil [NO<sub>3</sub>-], [NH<sub>4</sub>+] at 15cm depth
- Management info: planting/harvesting dates, fertilizer application timing and rate



Source: Tim Griffis

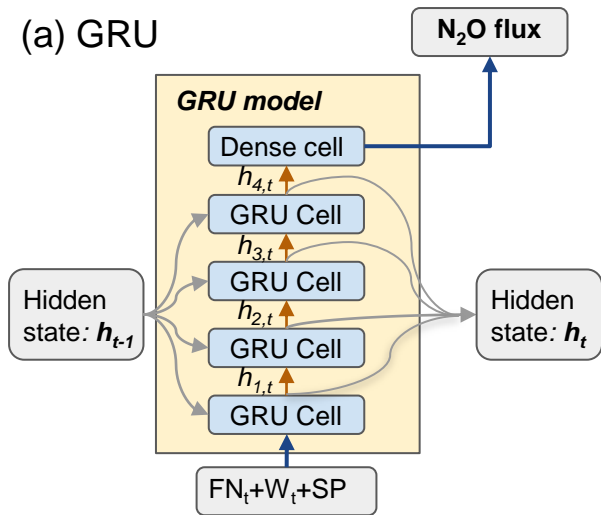
# The KGML model outperformed the pure ML model



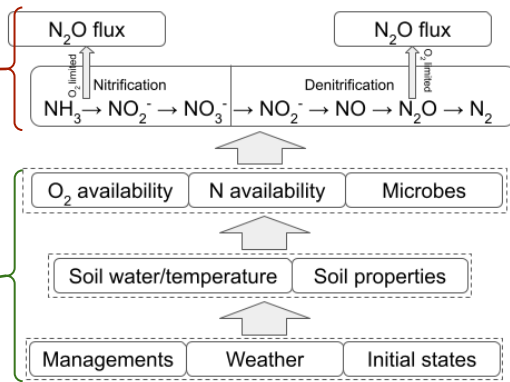
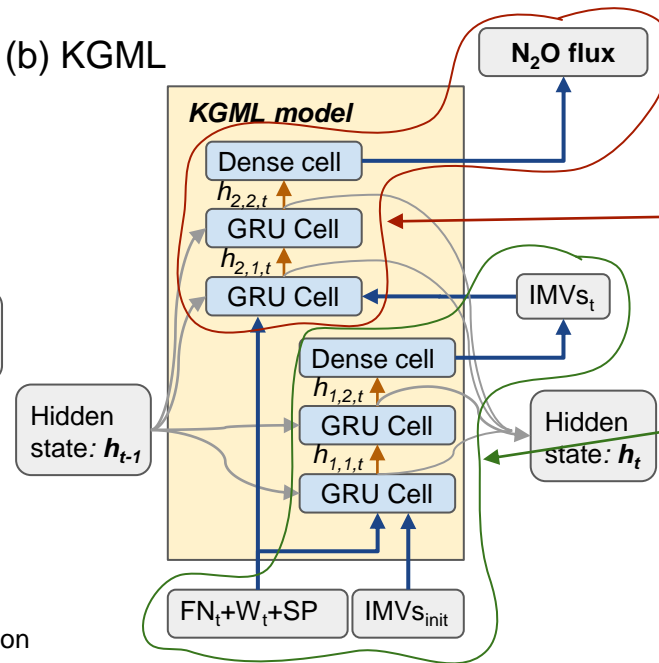


# Develop KGML model based on causal relations and feature importance

(a) GRU



(b) KGML

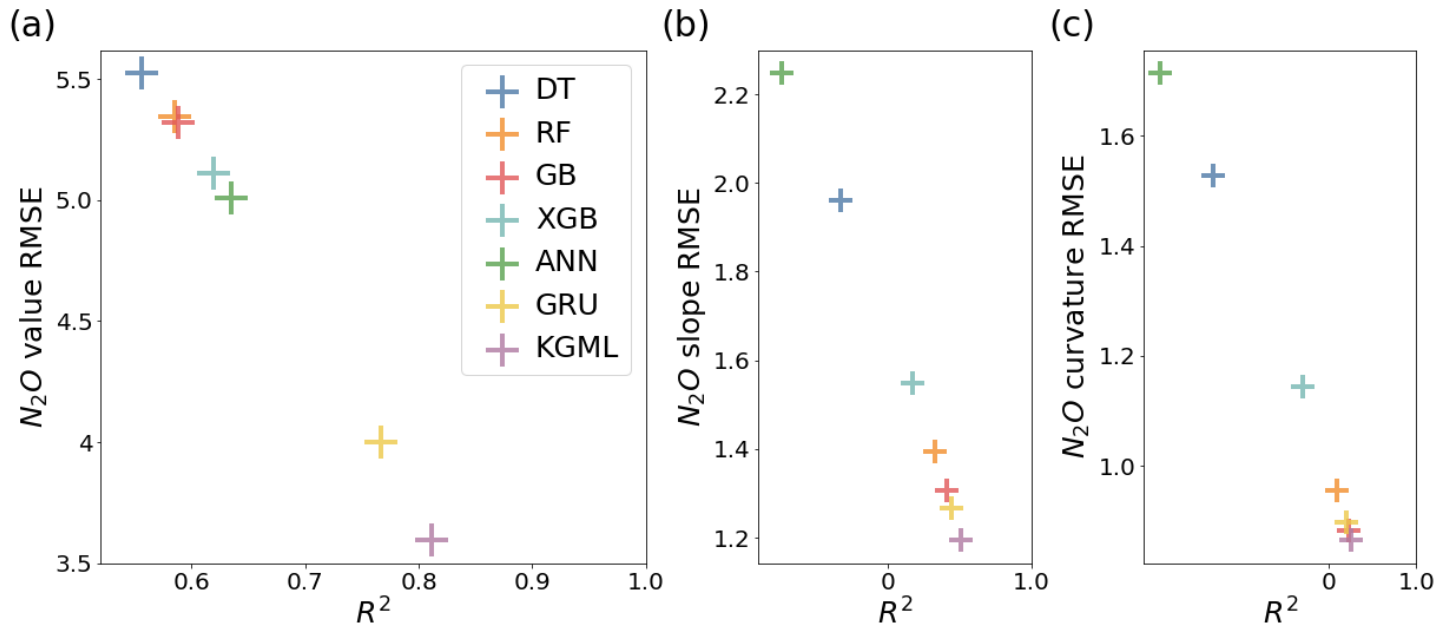


Detected key processes and causal relations for N<sub>2</sub>O fluxes

- Input/output (grey box)
- ML cell (blue box)
- Model (yellow box)
- Input/output transfer (blue arrow)
- PB module connection (grey arrow)
- $h$  transfer between ML cells with 20% dropout (orange arrow)
- $h$  input and output (grey arrow)

- GRU outperformed LSTM with its simpler structure in N<sub>2</sub>O simulation
- GRU model was used to do feature importance tests
- Knowledge guided initialization and architecture constraints were applied

# The KGML model outperformed the pure ML model



□ KGML (purple) outperform all other ML models

□ This is mainly because (1) pre-training using synthetic data, (2) knowledge guided architecture, (3) knowledge guided initial values



## Conclusion

---

- ❑ We used 1) knowledge-guided initialization, 2) hierarchical architecture, and 3) initial values of intermediate variables to develop the **KGML-ag structure for N<sub>2</sub>O prediction**
- ❑ The KGML-ag model has been tested on mesocosm experiment observations and can **outperform all other pure ML models**
- ❑ KGML **reduced data demand** significantly comparing to ML
- ❑ More N<sub>2</sub>O flux data are needed for further improvement of KGML
- ❑ The **structures are flexible** and can be easily revised or transferred to other data/study (**Another study of KGML-ag for Carbon** is presented in AGU poster named: *Estimating the Autotrophic and Heterotrophic Respiration in the US Crop Fields using Knowledge Guided Machine Learning*)

Don't have a good day, have a great day!

---

Thank you so much for your interest!

If you have any questions, please feel free to contact  
**Licheng Liu**: [lichengl@umn.edu](mailto:lichengl@umn.edu)