



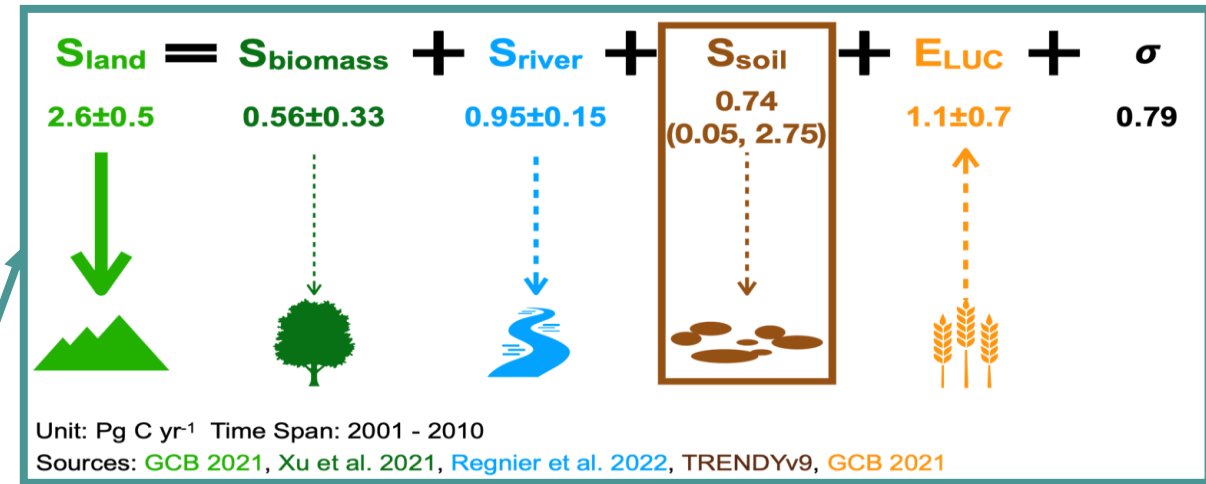
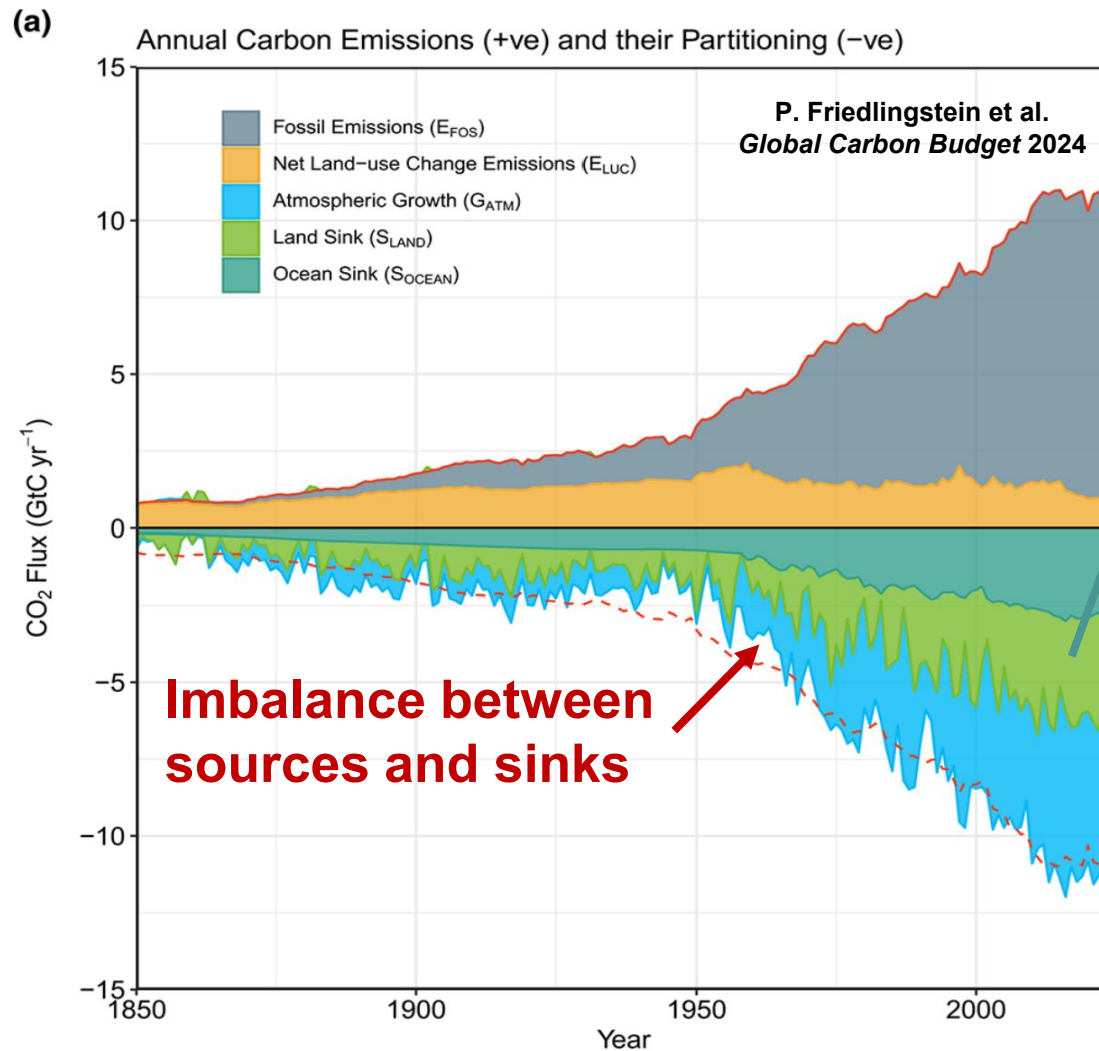
Artificial Intelligence to Advance Modeling and Understanding of Soil Organic Carbon

Haodi Xu¹, Joshua Fan², Feng Tao³, Md Nasim², Carla P. Gomes², Yiqi Luo¹

Cornell University, Department of ¹Soil & Crop Science, ²Computer Science, ³Ecology & Evolutionary Biology

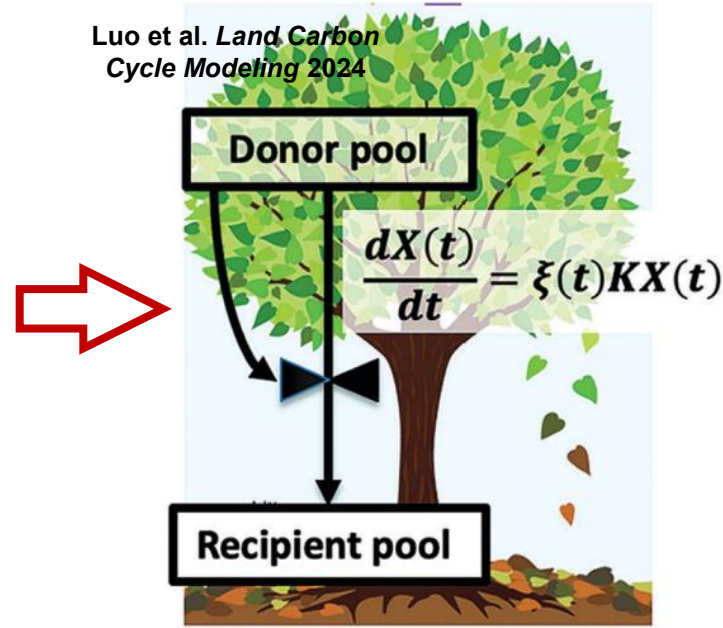
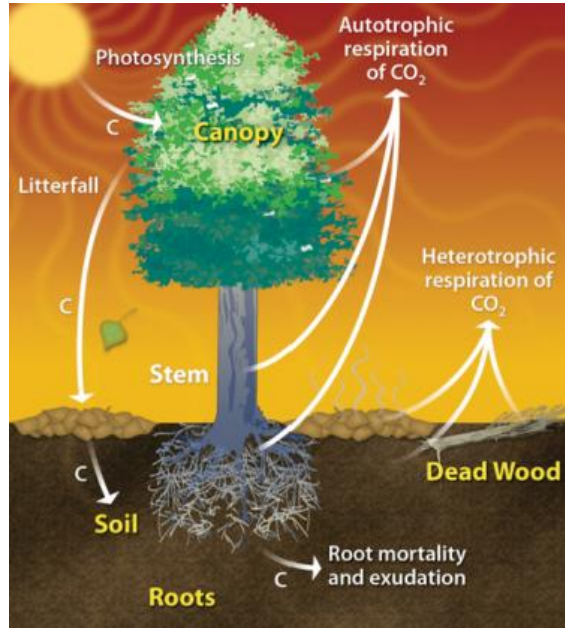
Cornell University

Soil is Critical in Understanding Global Carbon Cycle

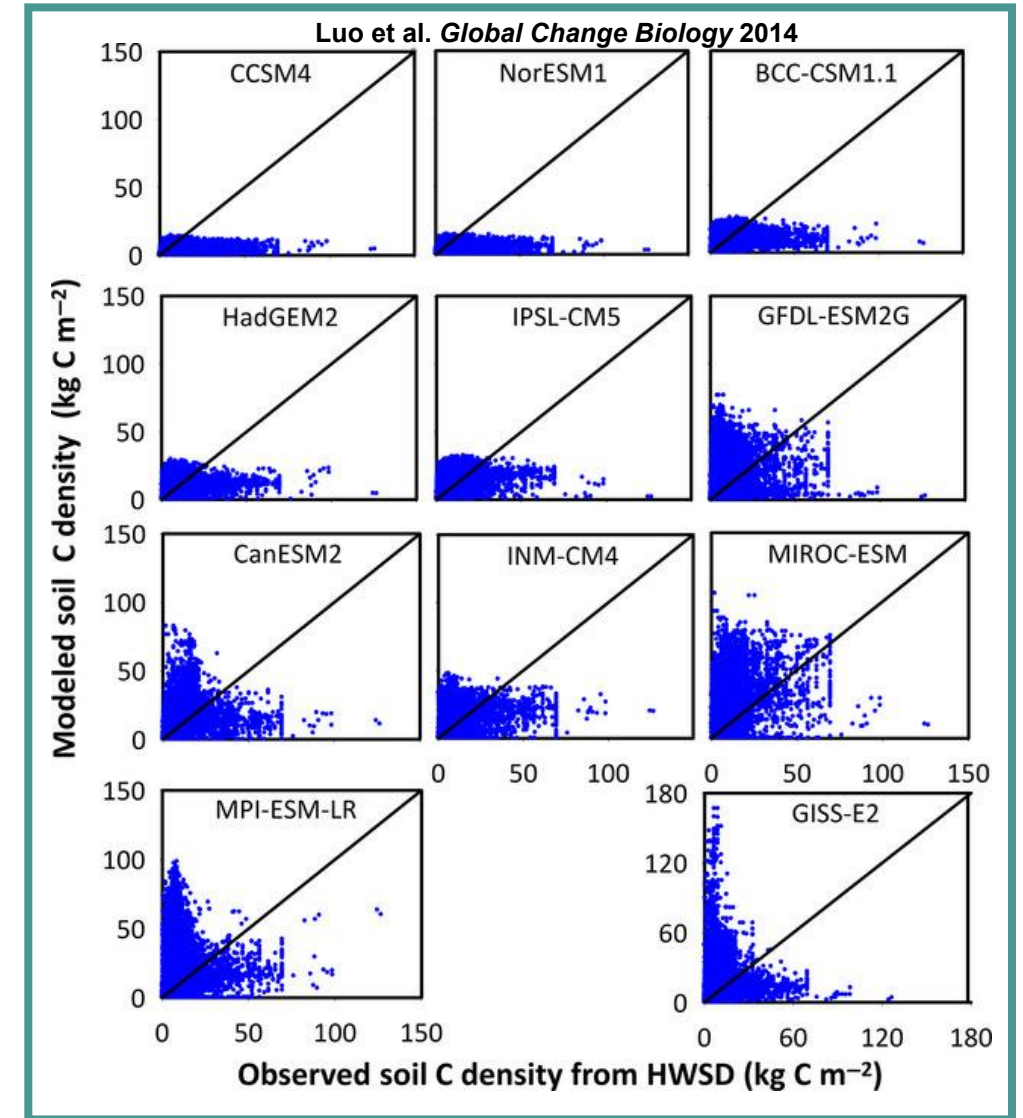


- **Imbalance** within global carbon cycle can partially explained by the uncertainties in land carbon sinks
- **Soils** hold a large part of uncertainties in land carbon sinks

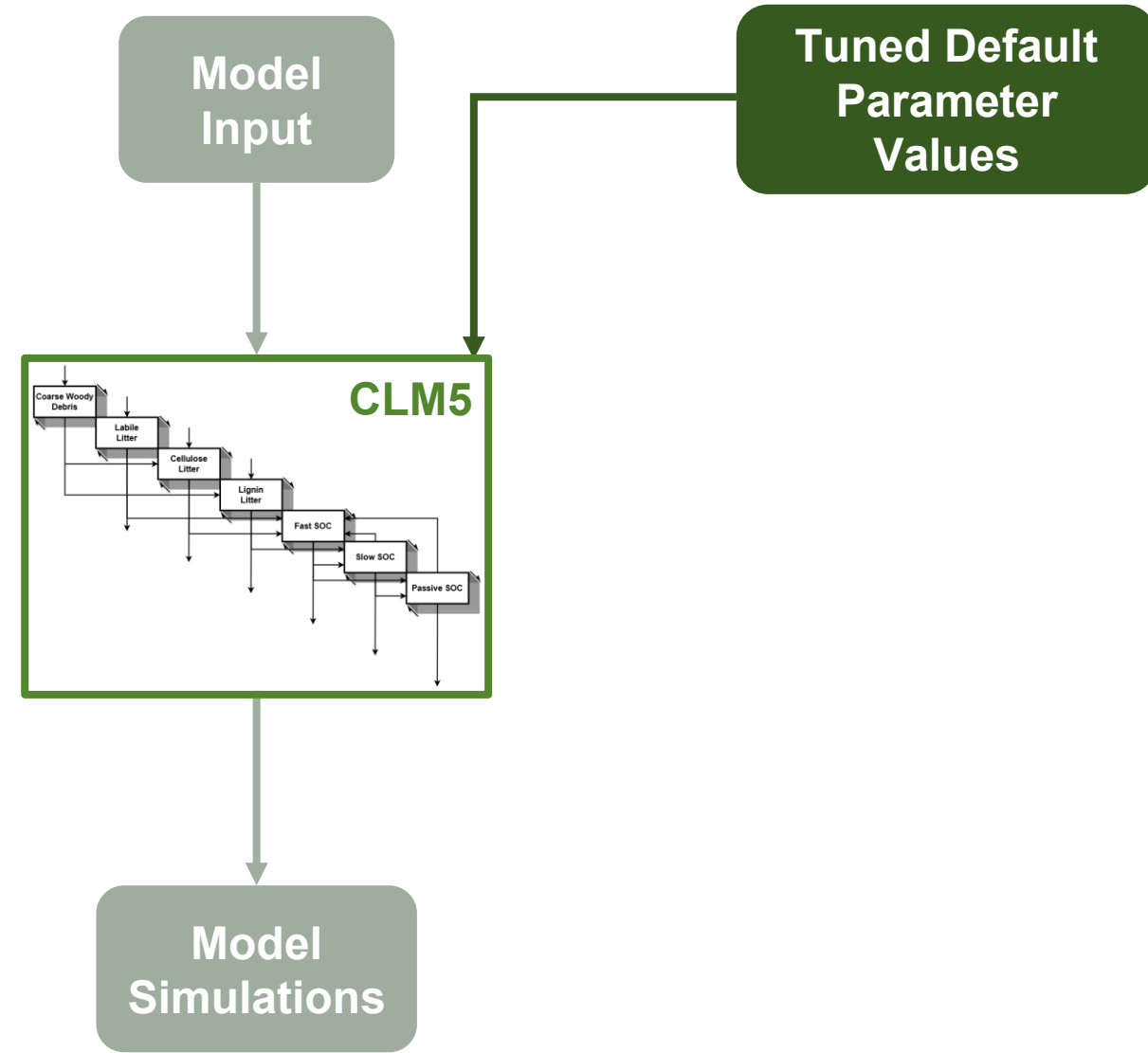
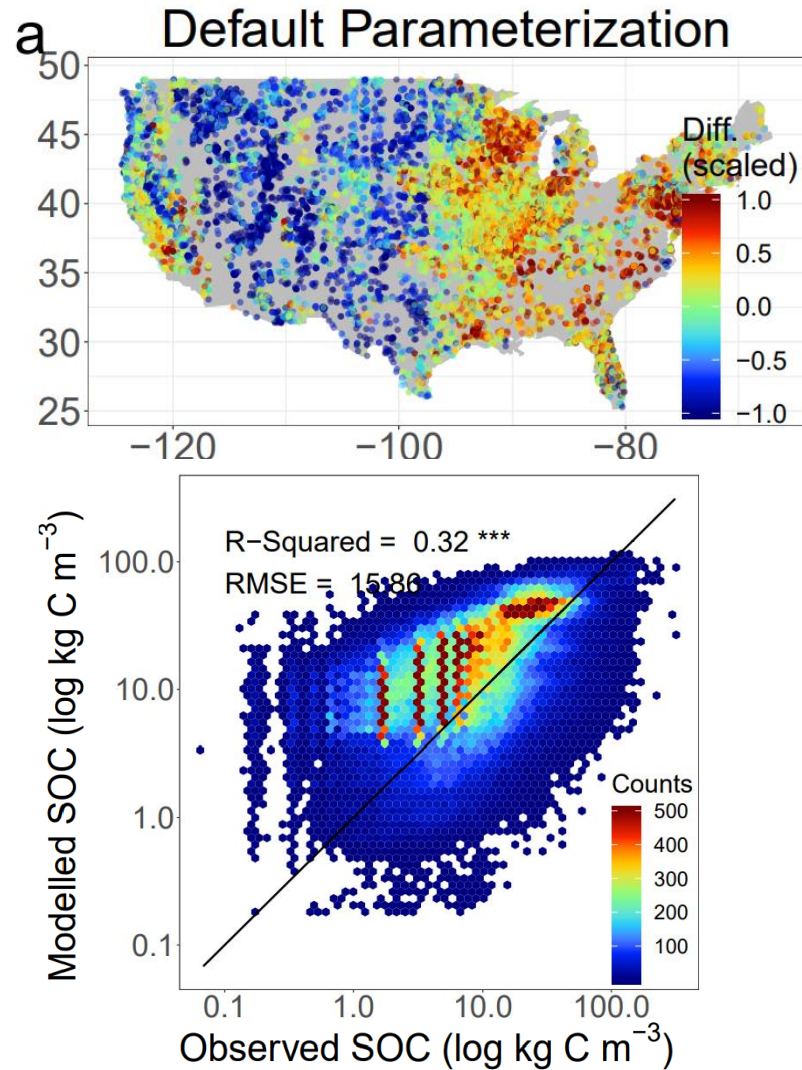
Modeling SOC Dynamics



- Our knowledge of SOC processes can be incorporated into process-based models
- Yet **high uncertainties** remain in capturing the global SOC cycle

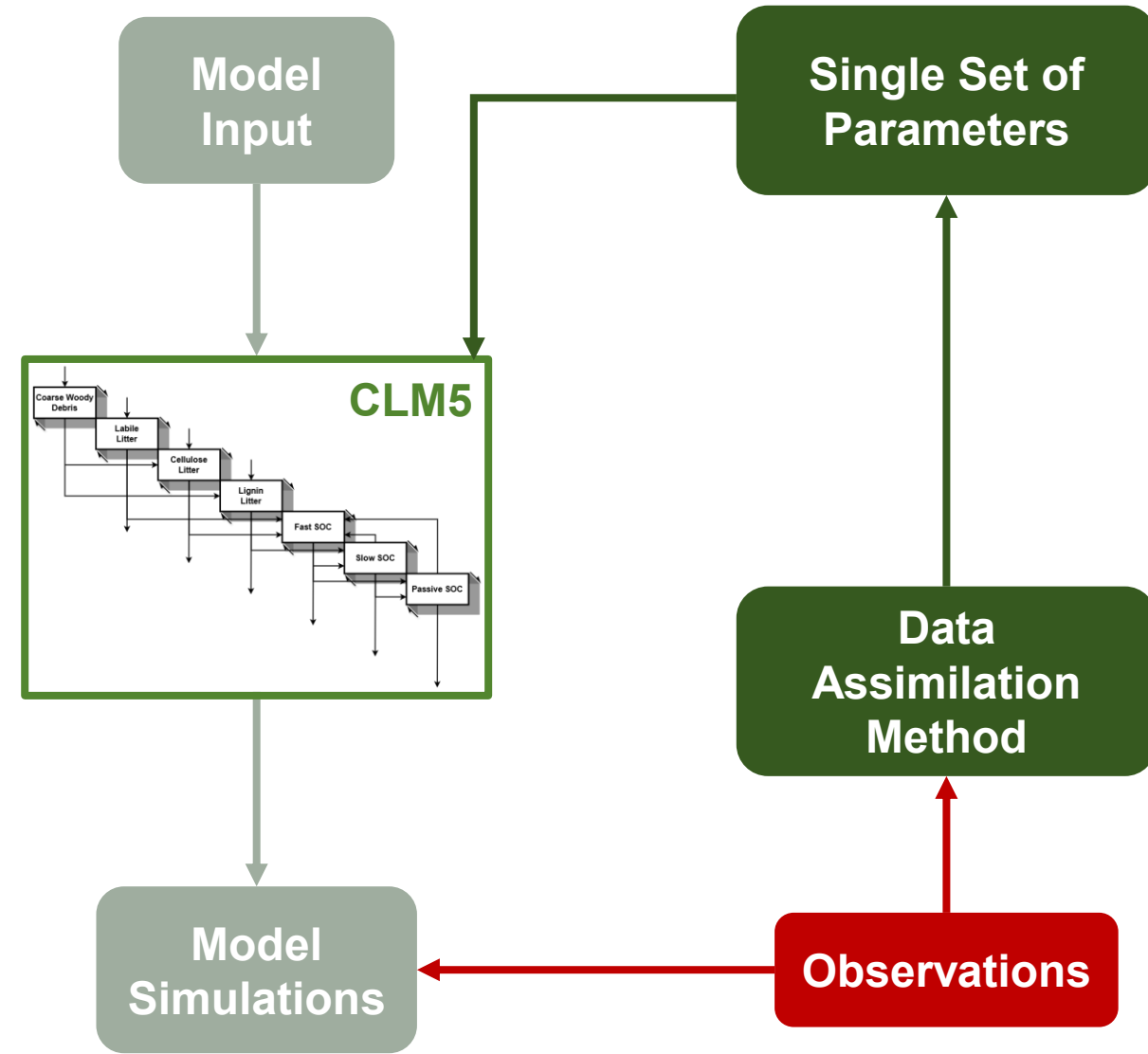
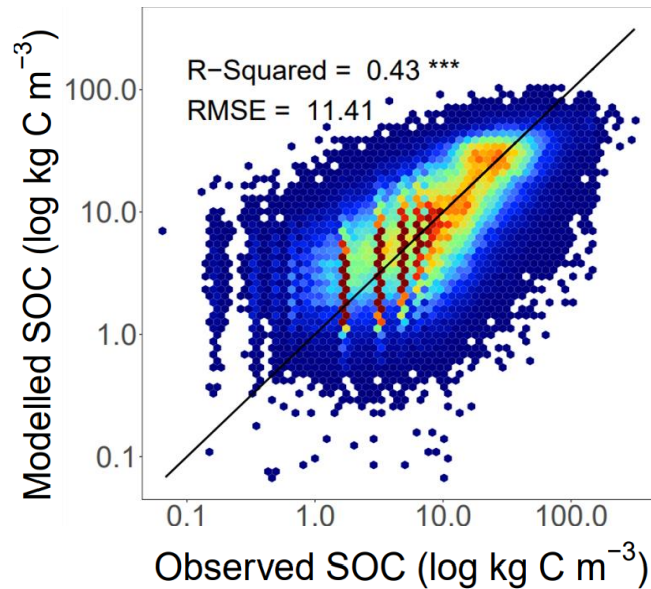
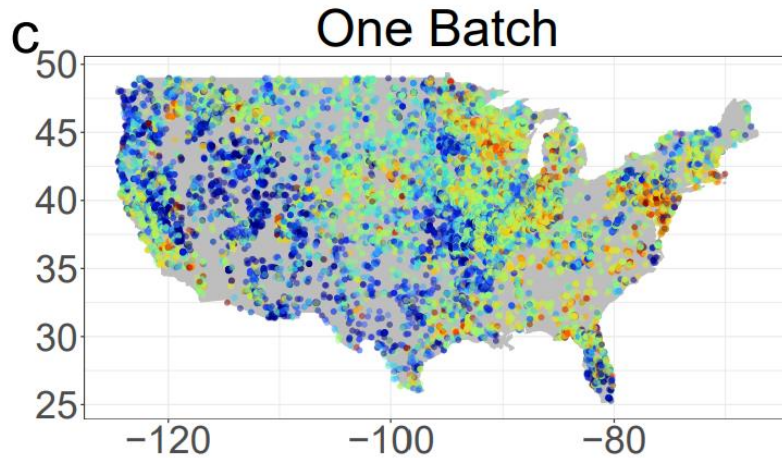


Modeled SOC across Conterminous US



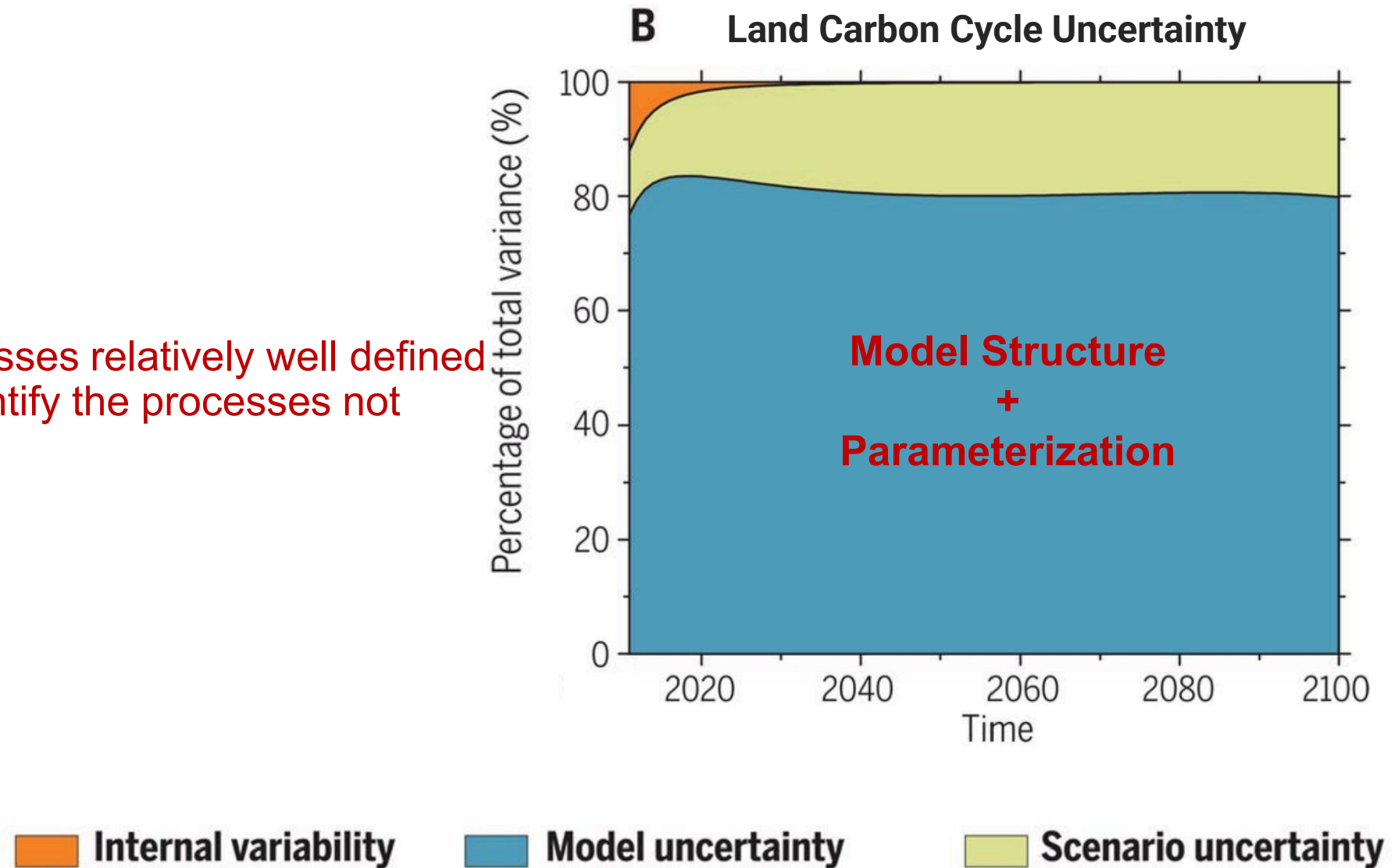
Tao et al. *Frontiers in Big Data* 2020

Data Assimilation with One Set of Parameters



Tao et al. *Frontiers in Big Data* 2020

Model Structure: processes relatively well defined
Parameterization: quantify the processes not explicitly represented

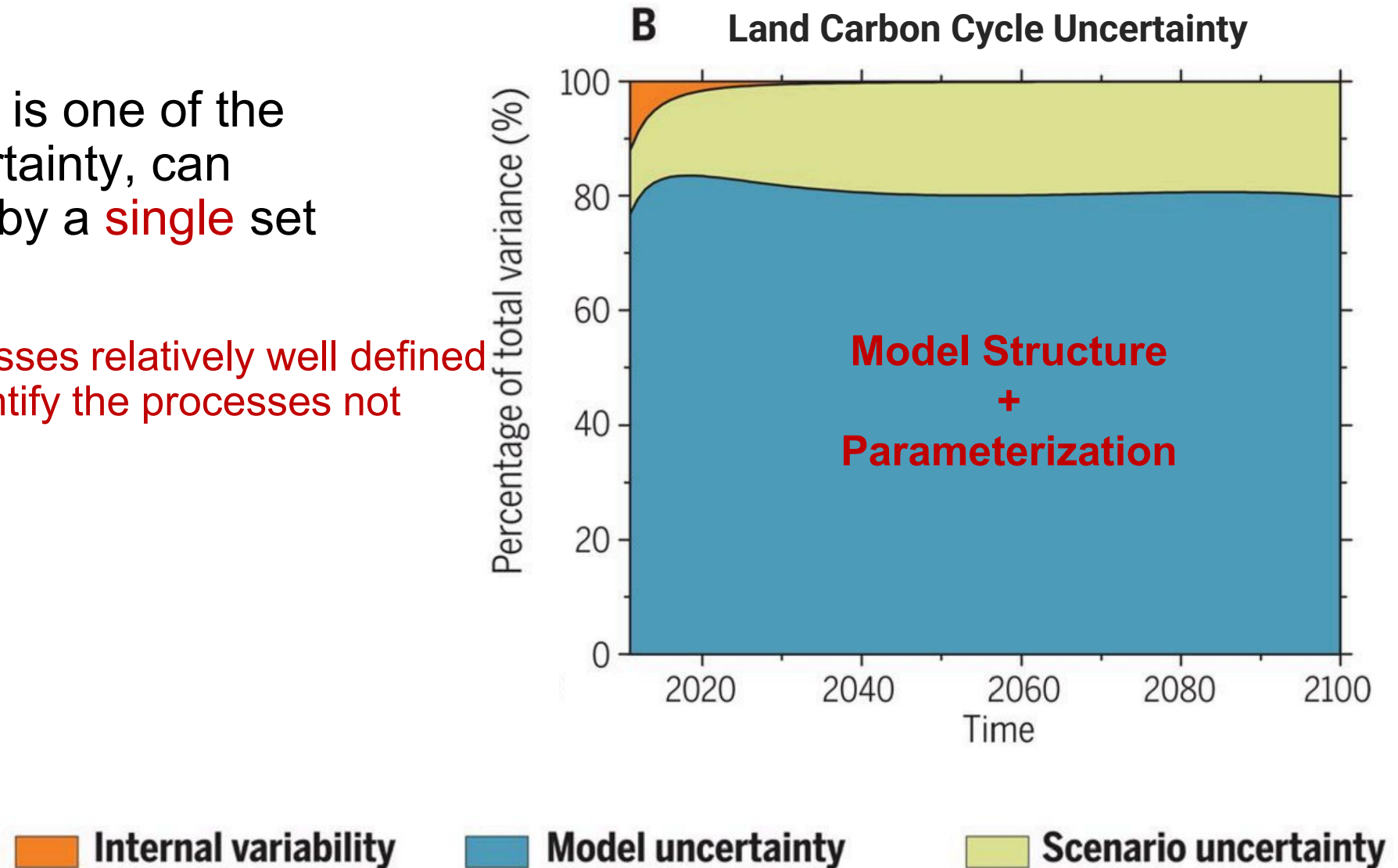


Bonan et al., *Science*, 2018

Heterogeneity of soils is one of the major sources of uncertainty, can hardly be represented by a **single** set of parameters

Model Structure: processes relatively well defined

Parameterization: quantify the processes not explicitly represented



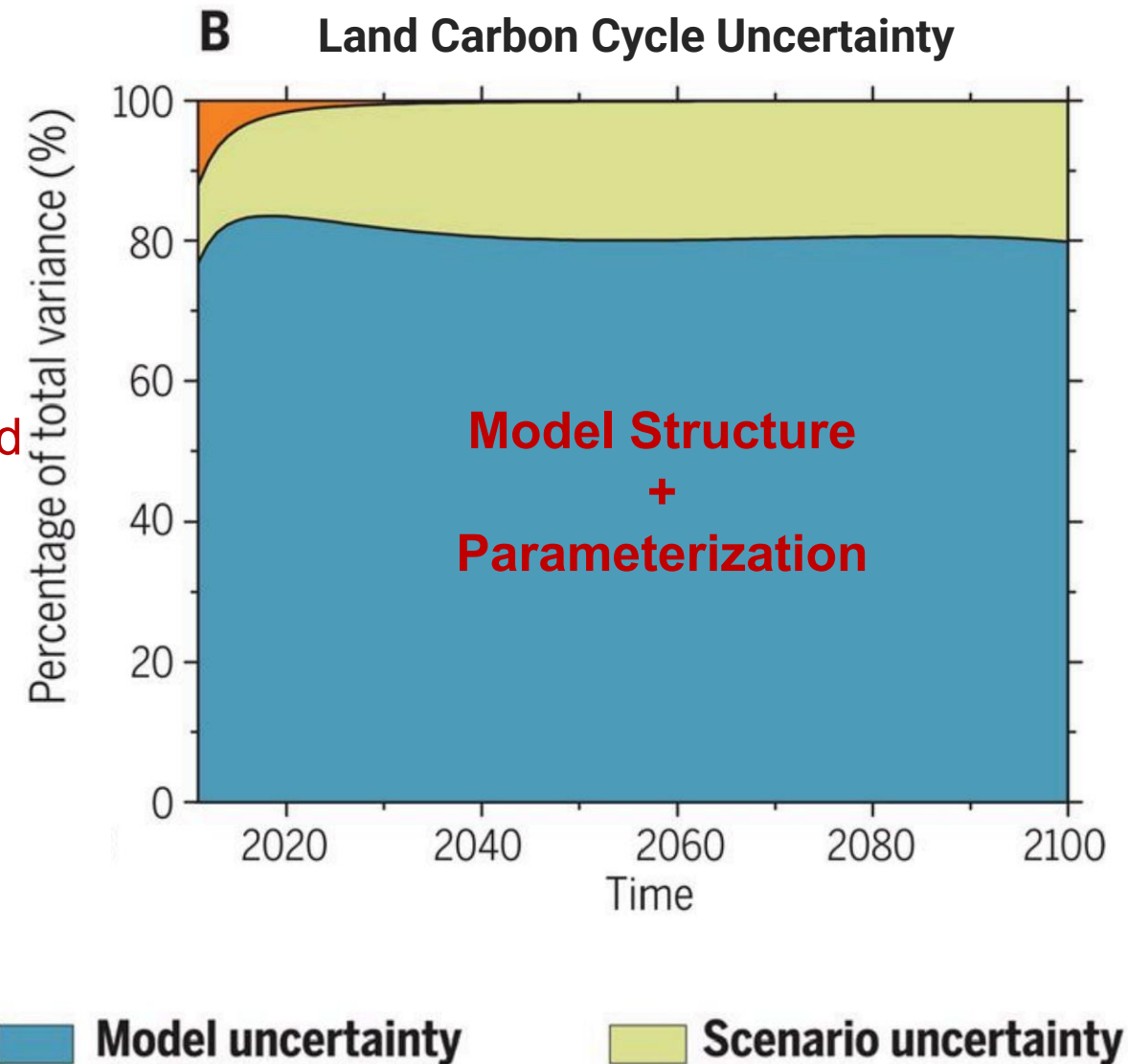
Bonan et al., *Science*, 2018

Heterogeneity of soils is one of the major sources of uncertainty, can hardly be represented by a **single** set of parameters

Model Structure: processes relatively well defined

Parameterization: quantify the processes not explicitly represented

Identify parameter sets that accurately simulate SOC storage, allowing the model to more effectively represent the complexity of real soil systems.



Bonan et al., *Science*, 2018

Environmental Information

Climate

Soil pH

Lon, Lat, Elevation

Soil Texture

Bulk Density

Land Cover

⋮

**Environmental
Information**



**Process-based
Model Parameters**

Climate

Soil pH

Lon, Lat, Elevation

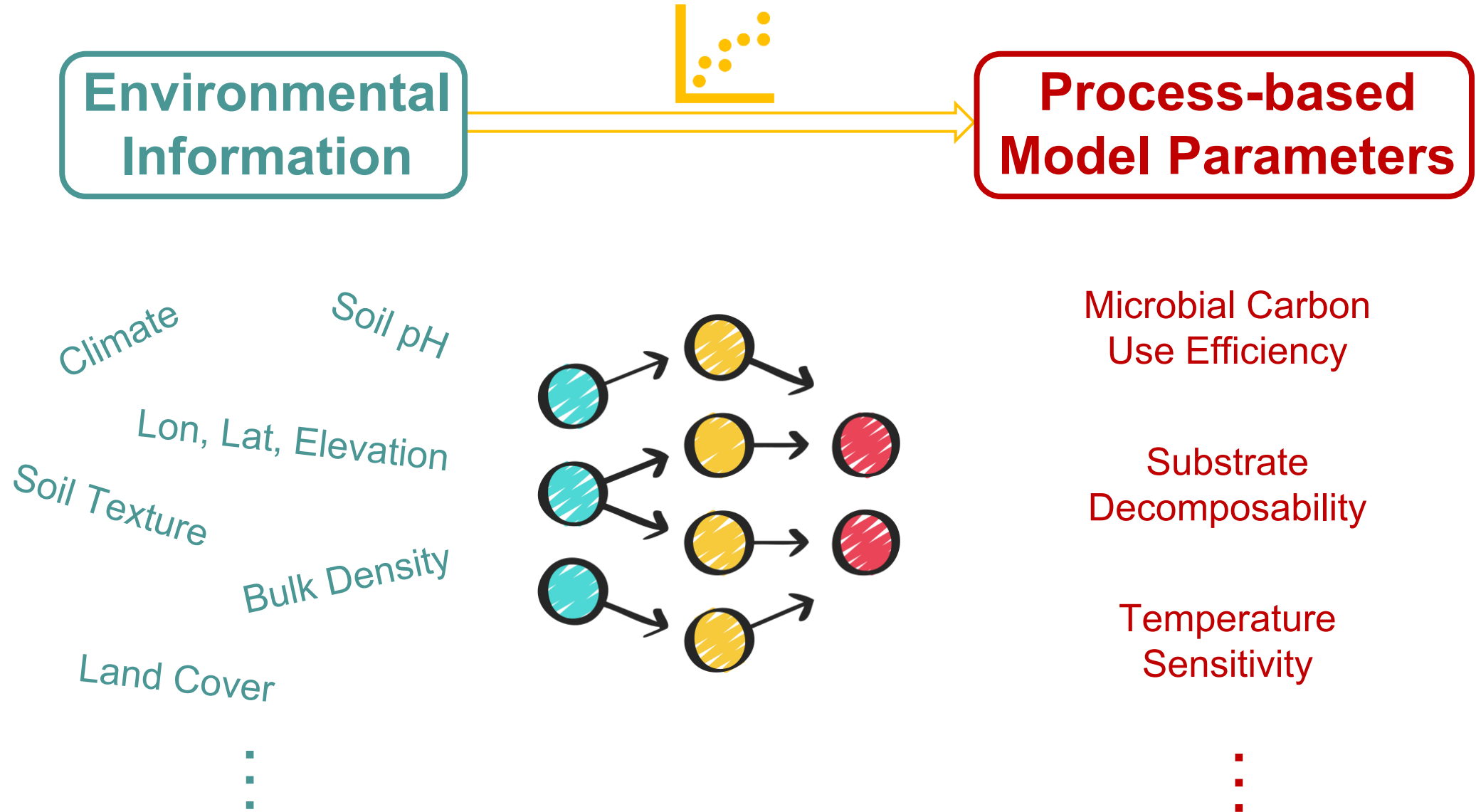
Soil Texture

Bulk Density

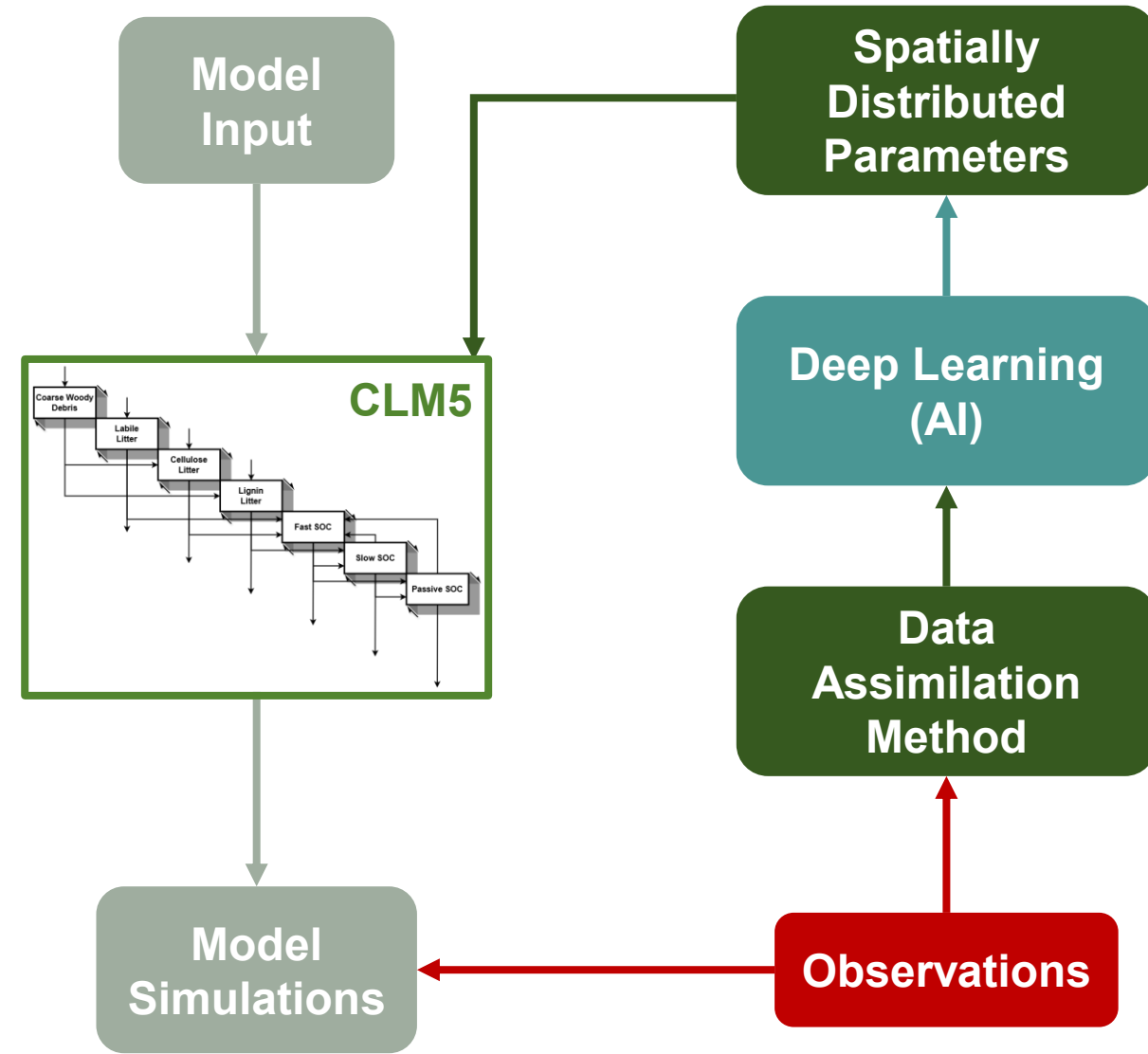
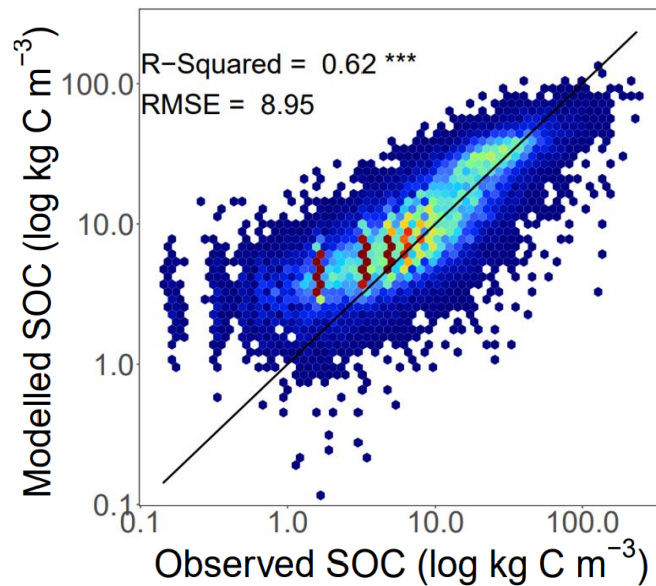
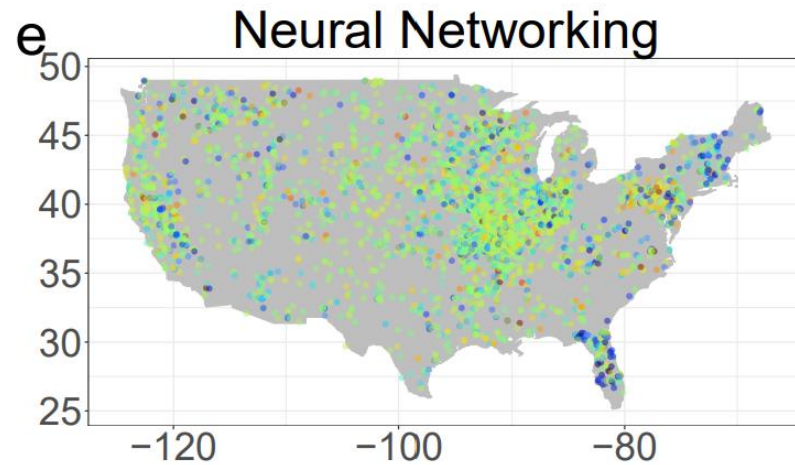
Land Cover



Artificial Intelligence (Deep Learning)

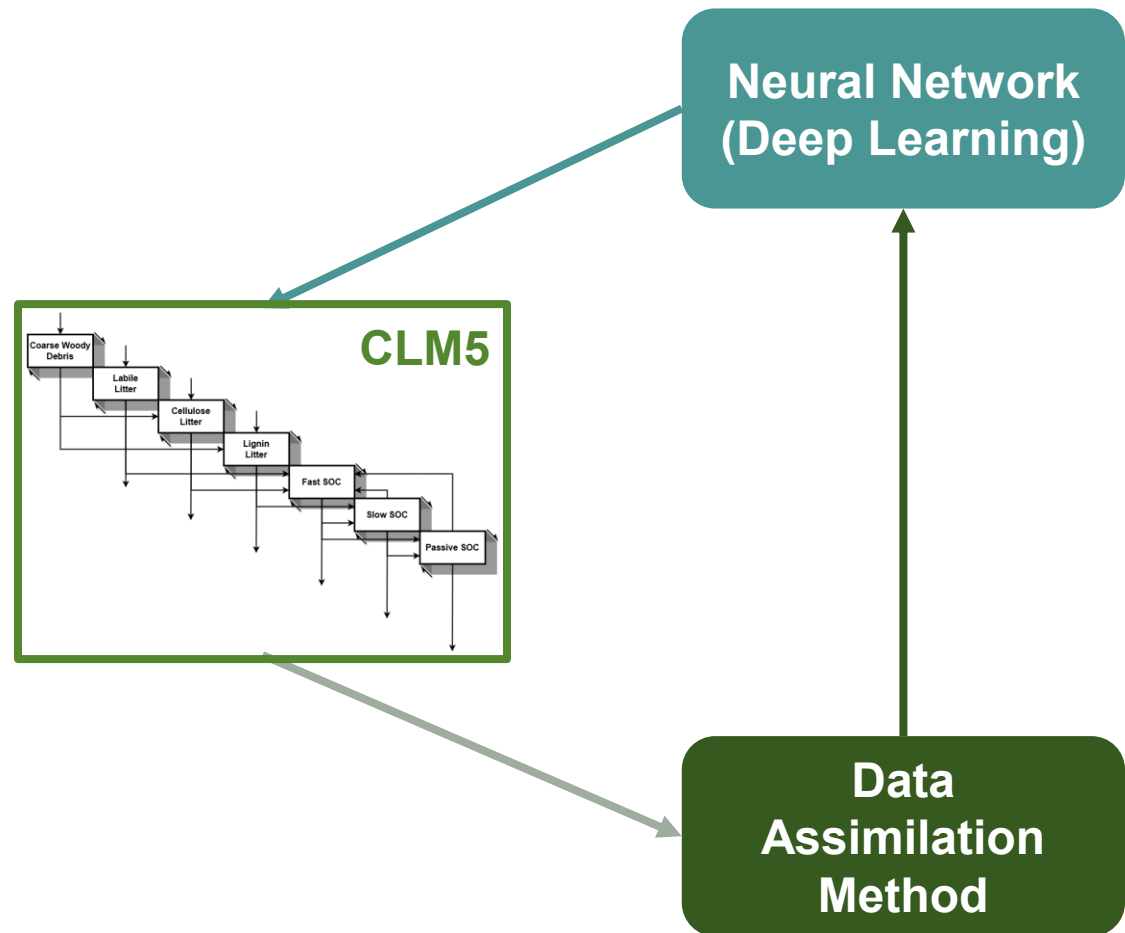


PROcess guided deep learning and DATA driven modeling (PRODA)



Tao et al. *Frontiers in Big Data* 2020

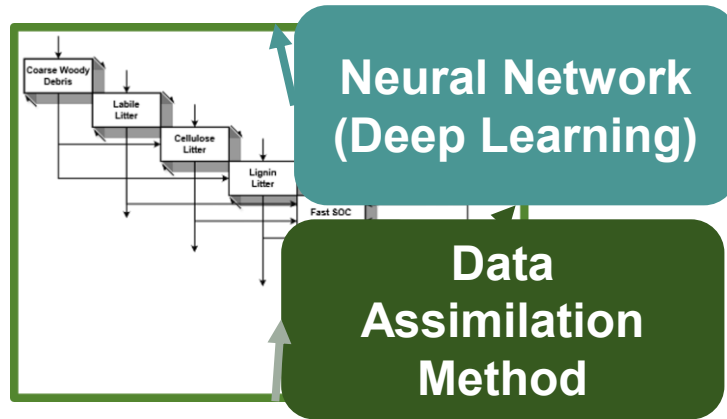
The need to further improve PRODA



PRODA includes three components:

- High computational cost due to **Bayesian-based** data assimilation method (e.g. MCMC)
- Limited flexibility with observed data
- Hard to understand and use

The need to further improve PRODA



PRODA includes three components:

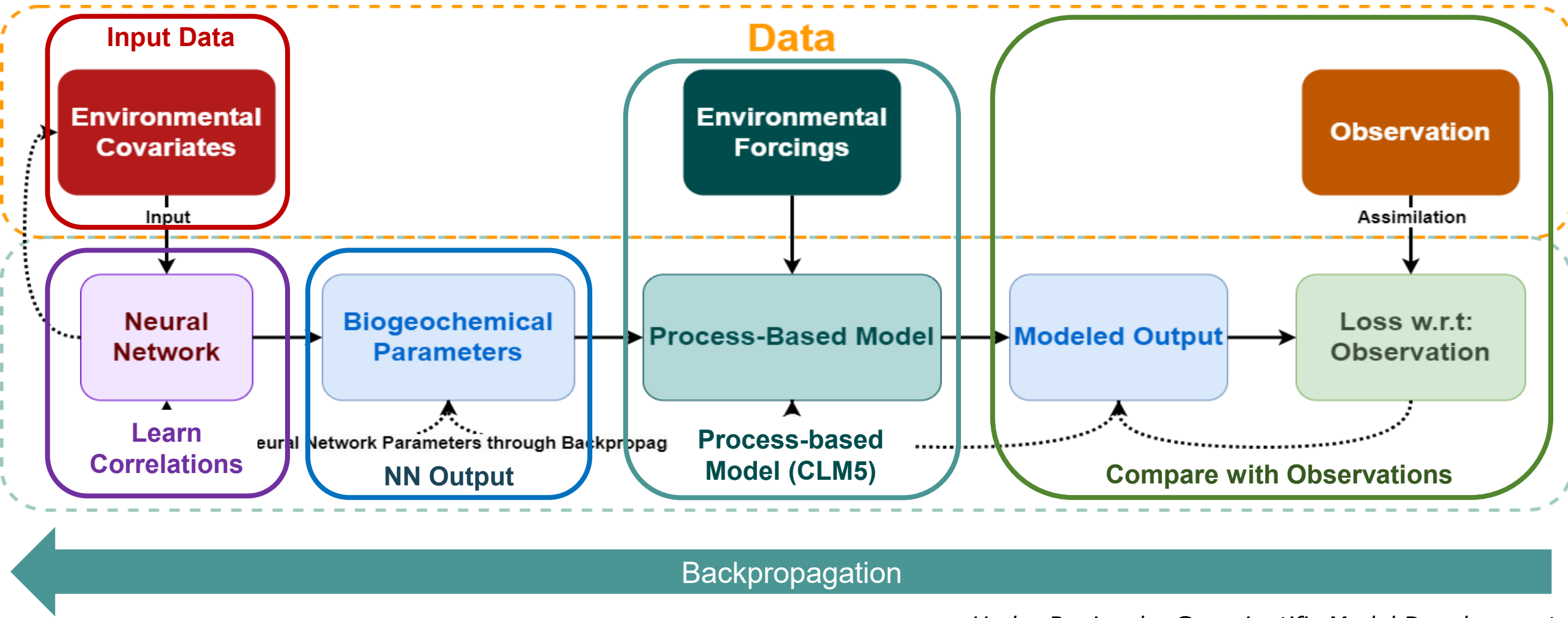
- High computational cost due to **Bayesian-based** data assimilation method (e.g. MCMC)
- Limited flexibility with observed data
- Hard to understand and use

Combine these altogether?



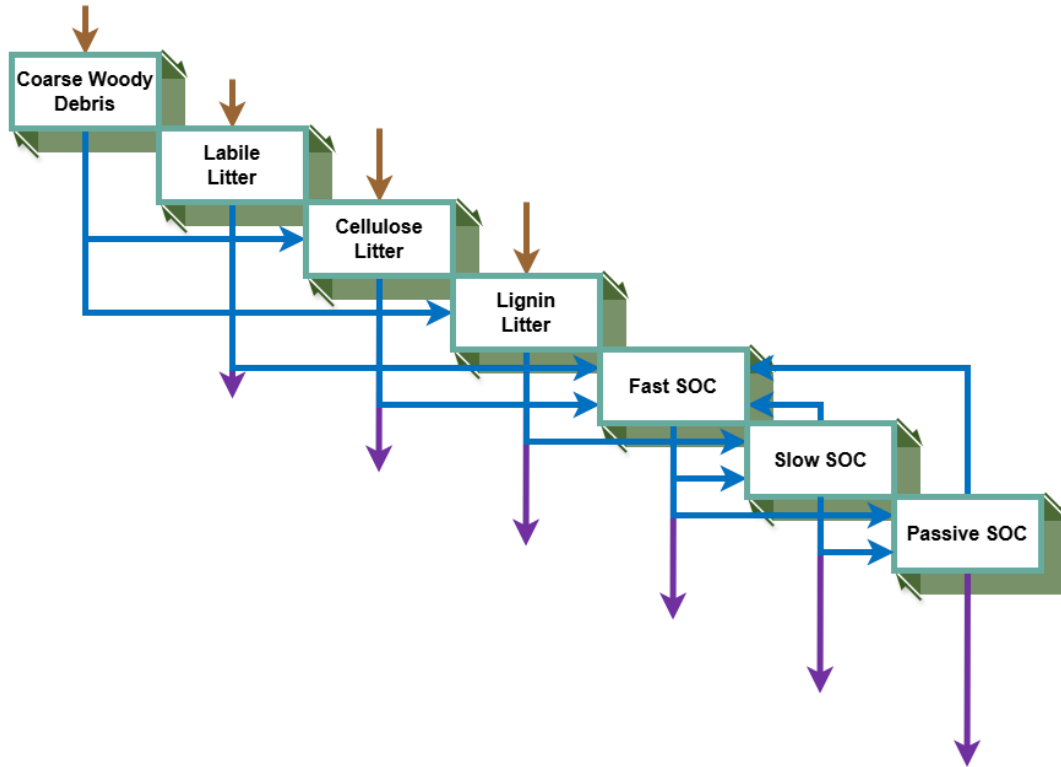
A Biogeochemistry-Informed Neural Network (BINN)

Workflow of BINN



Under Review by Geoscientific Model Development

Matrix Representation of CLM5



Changes in Carbon in each pool through out the time

C movement among different pools

Carbon Dynamics through Diffusion

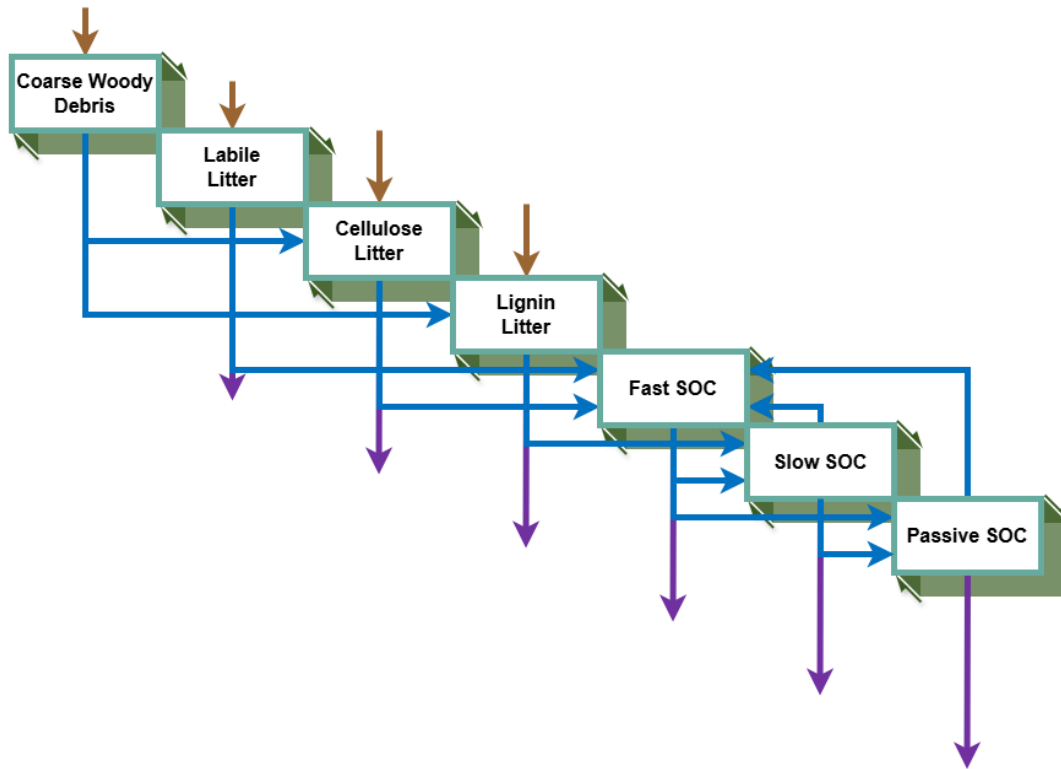
$$\frac{dX(t)}{dt} = B(t)I(t) - A \xi(t) K X(t) - V(t)X(t)$$

Carbon Input

The intrinsic decomposition rate of each C pool, modified by the environmental scalar

Each matrix is constructed from one or more unique parameters (21 parameters in total)

Matrix Representation of CLM5



Changes in Carbon in each pool through out the time

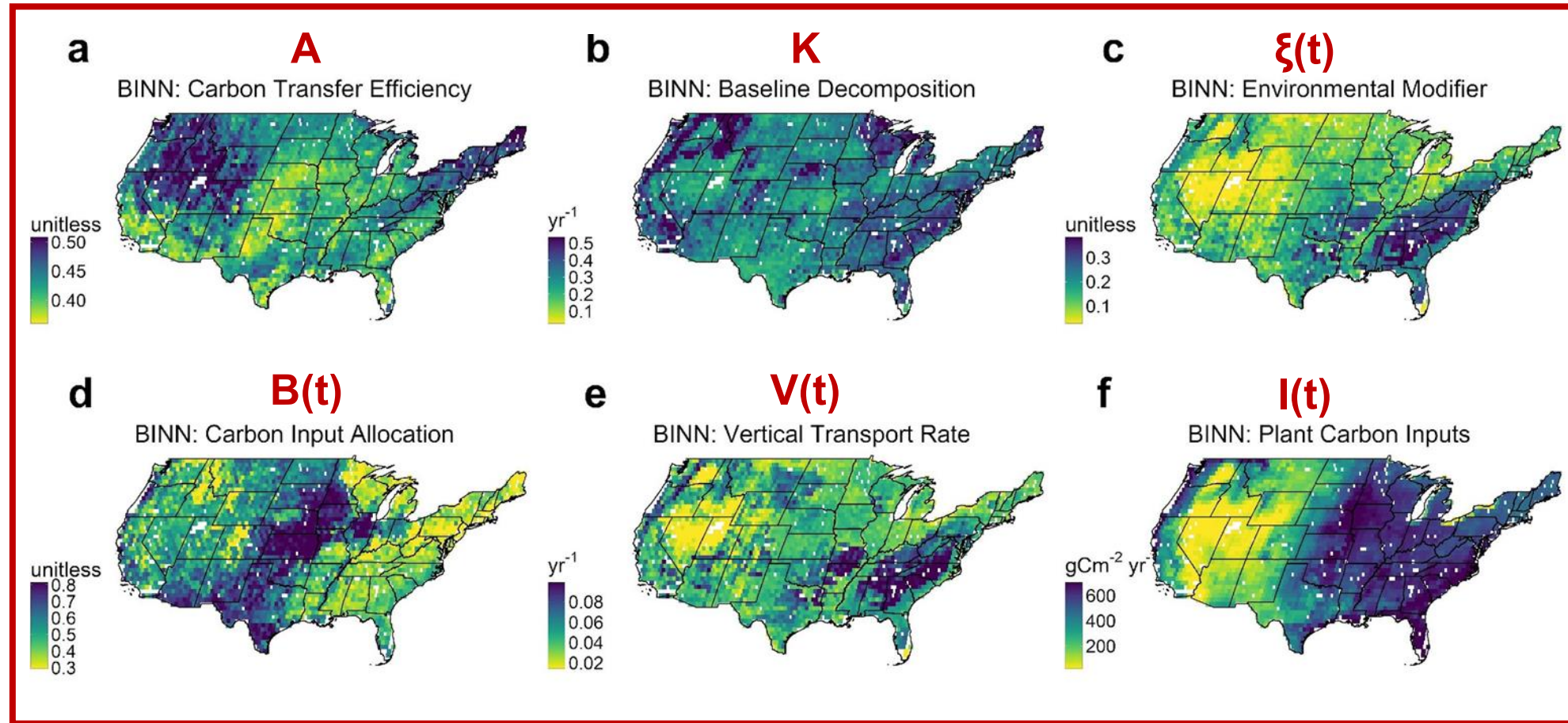
C movement among different pools

Carbon Dynamics through Diffusion

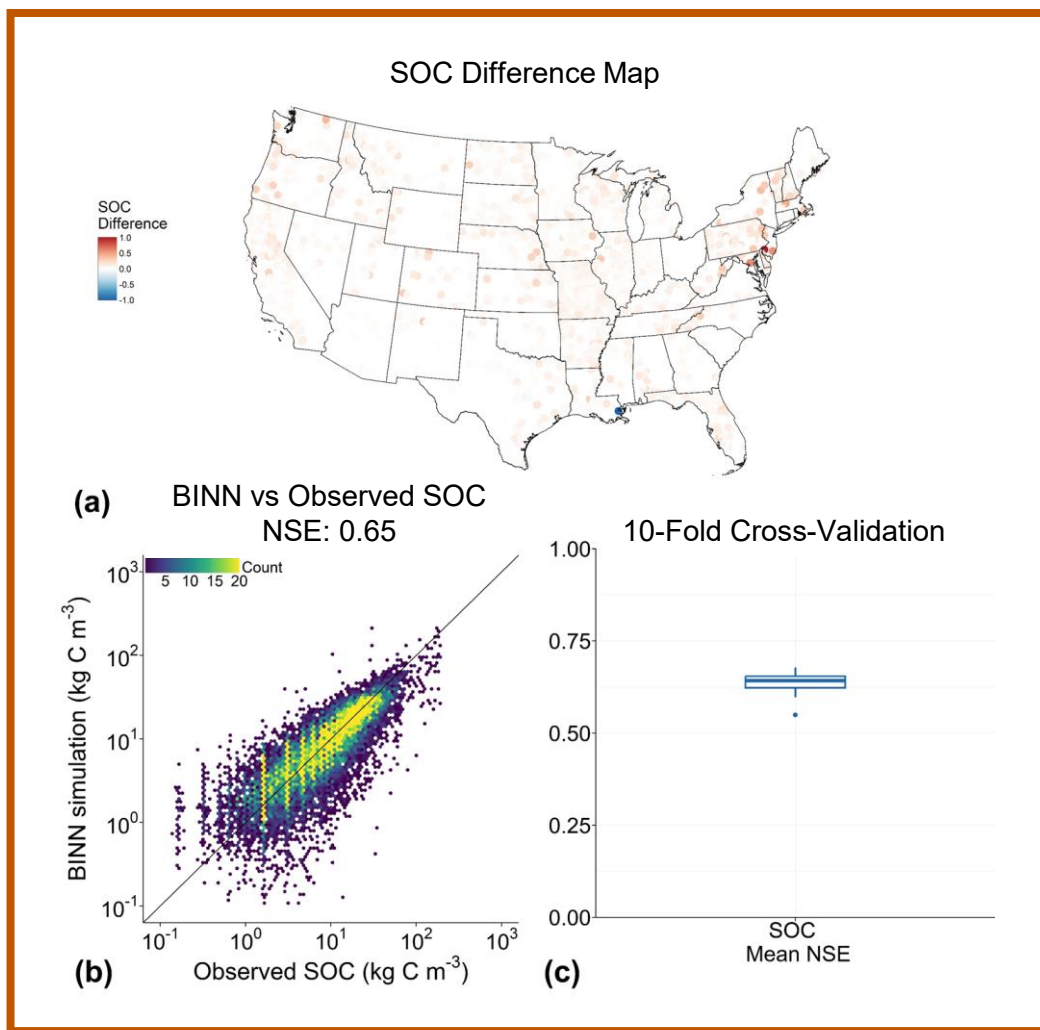
$$\frac{dX(t)}{dt} = \underbrace{B(t)I(t)}_{\text{Carbon Input}} - \underbrace{A}_{\text{C movement among different pools}} \underbrace{\xi(t)K}_{\text{The intrinsic decomposition rate of each C pool, modified by the environmental scalar}} X(t) - \underbrace{V(t)}_{\text{Carbon Dynamics through Diffusion}} X(t)$$

I : Model Components, quantified by 21 model parameters

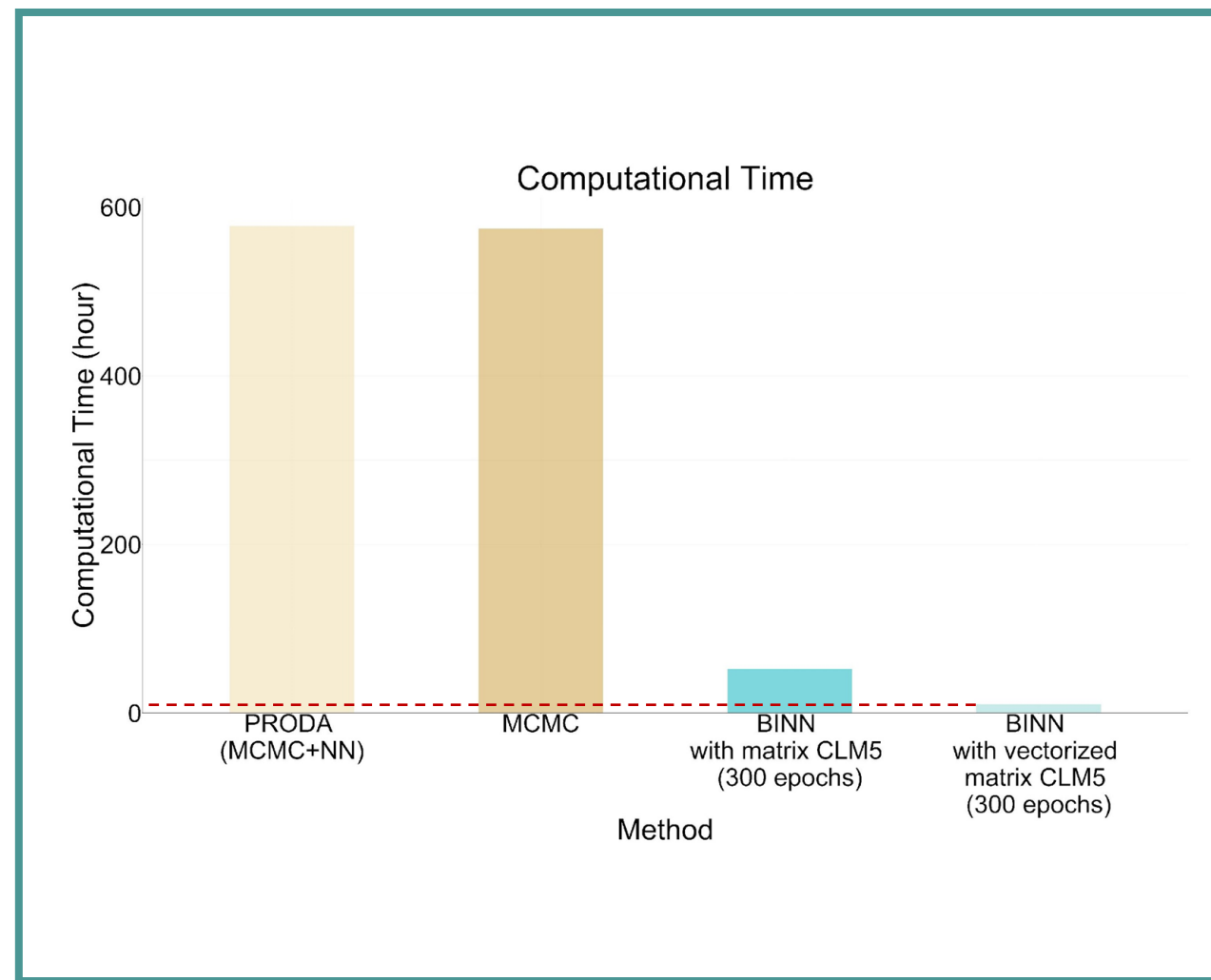
BINN discovers the spatial distribution of underlying mechanisms from SOC observations and environmental covariates (without any data for these mechanisms)

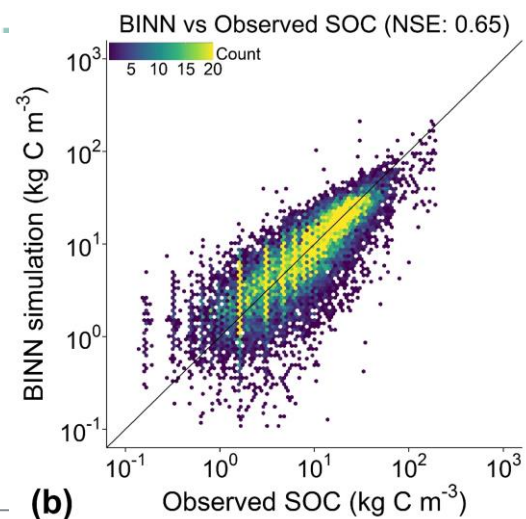
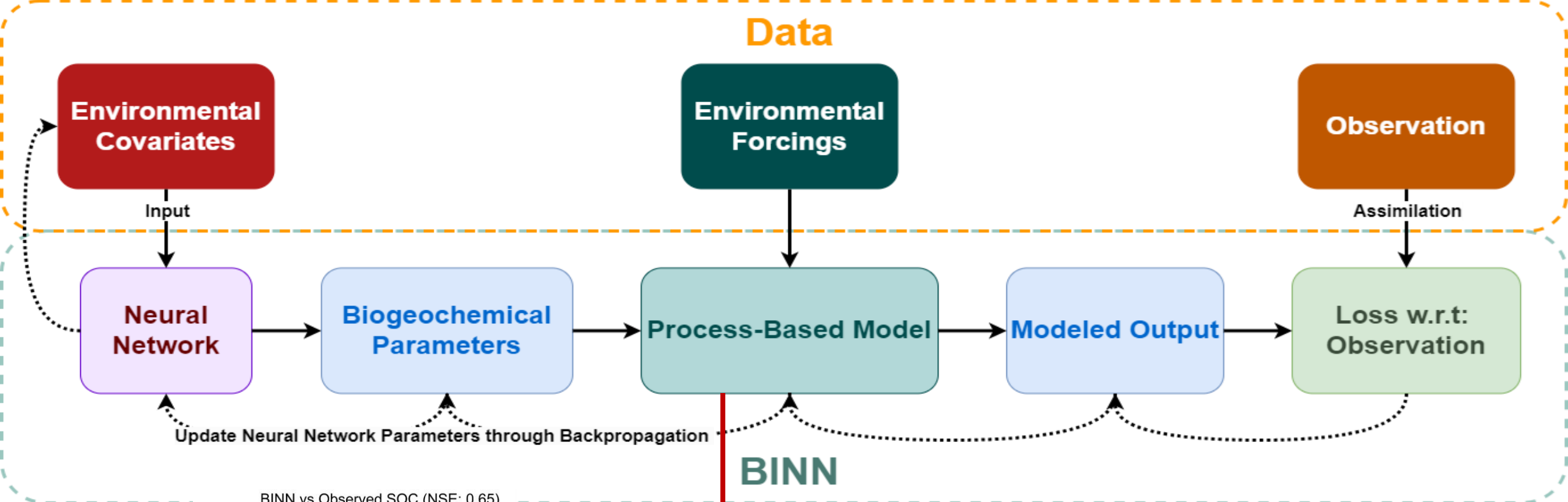


More Realistic Representations of SOC Distribution

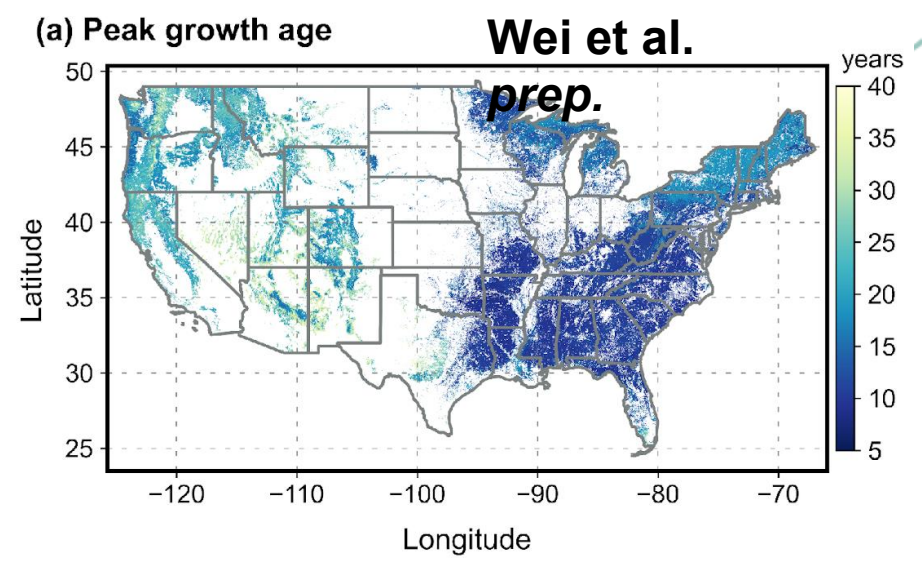
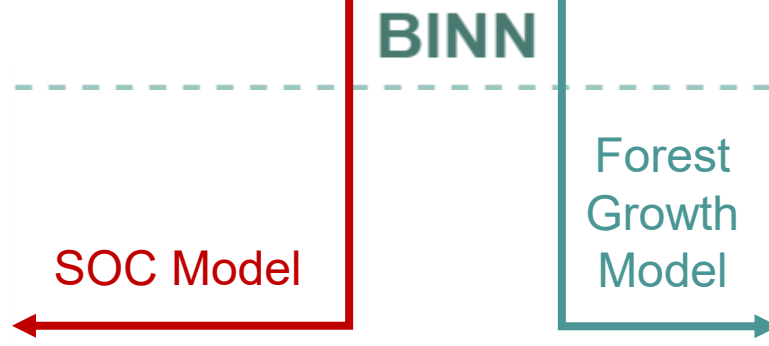
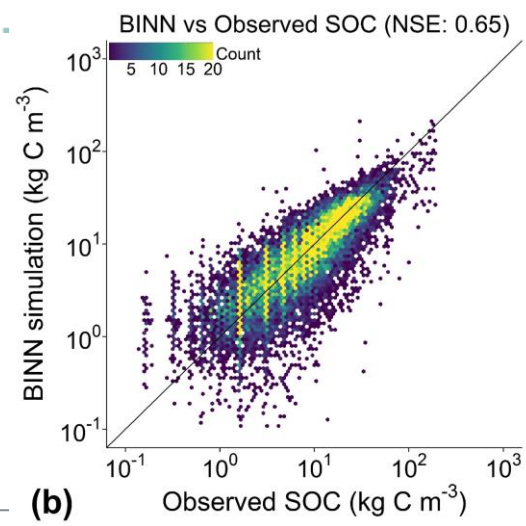
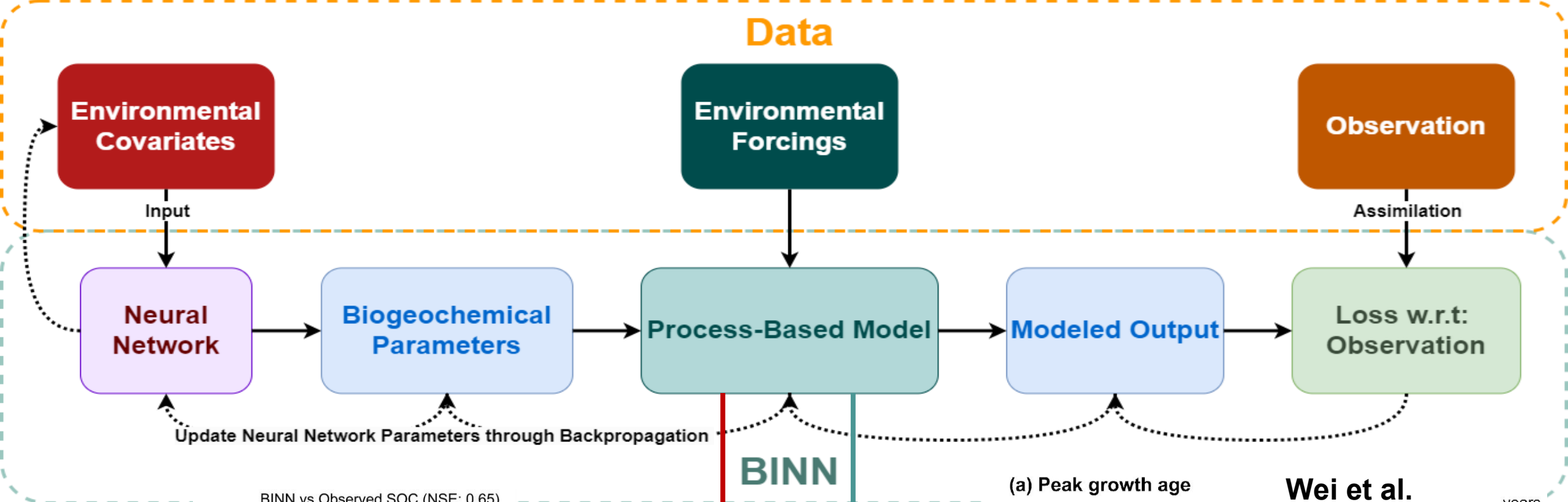


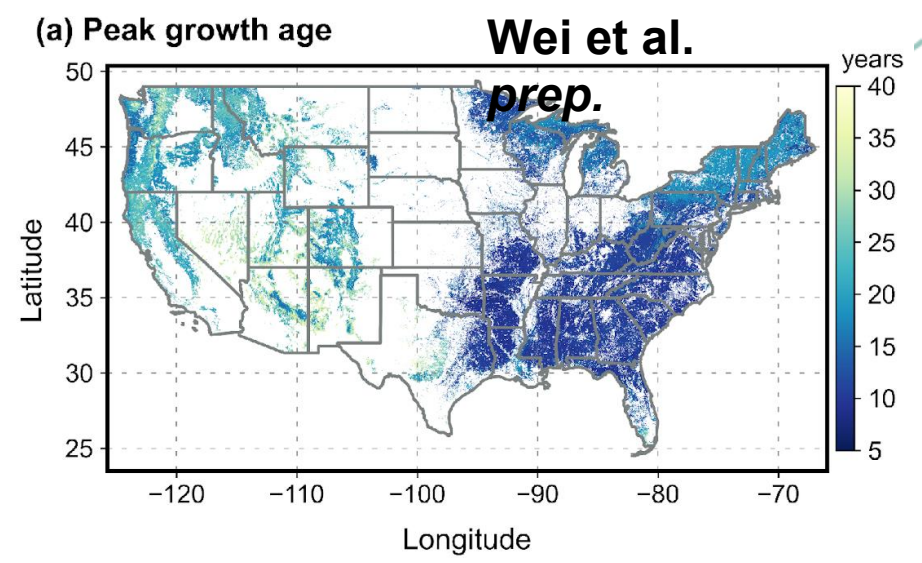
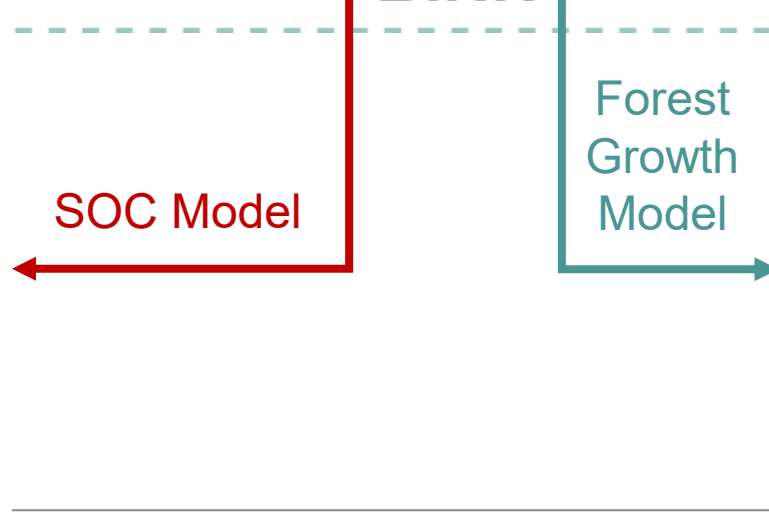
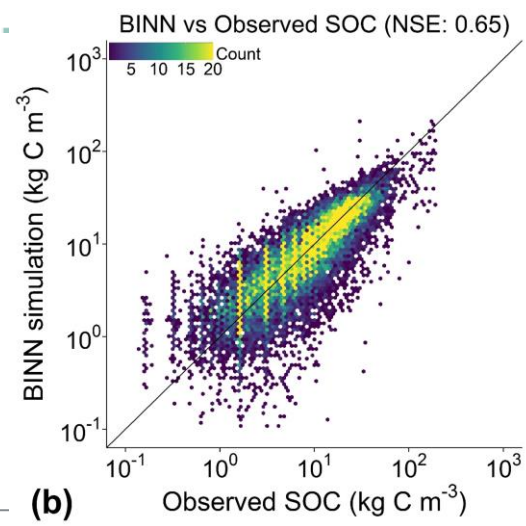
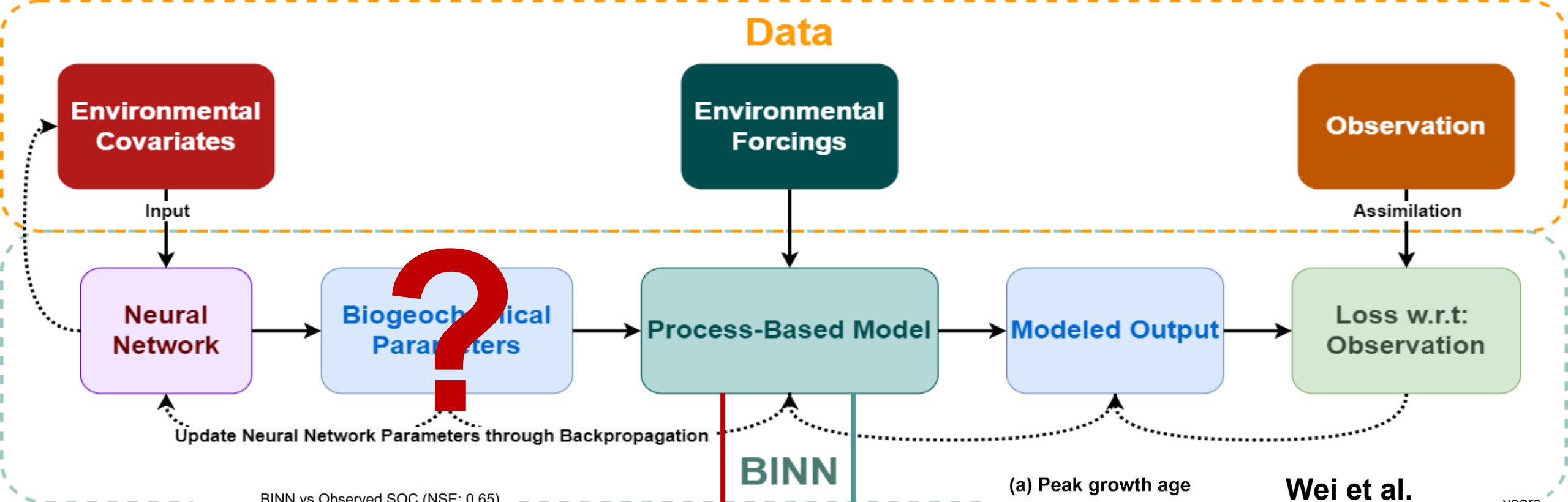
Way Less Computational Requirements (More than 50 times)



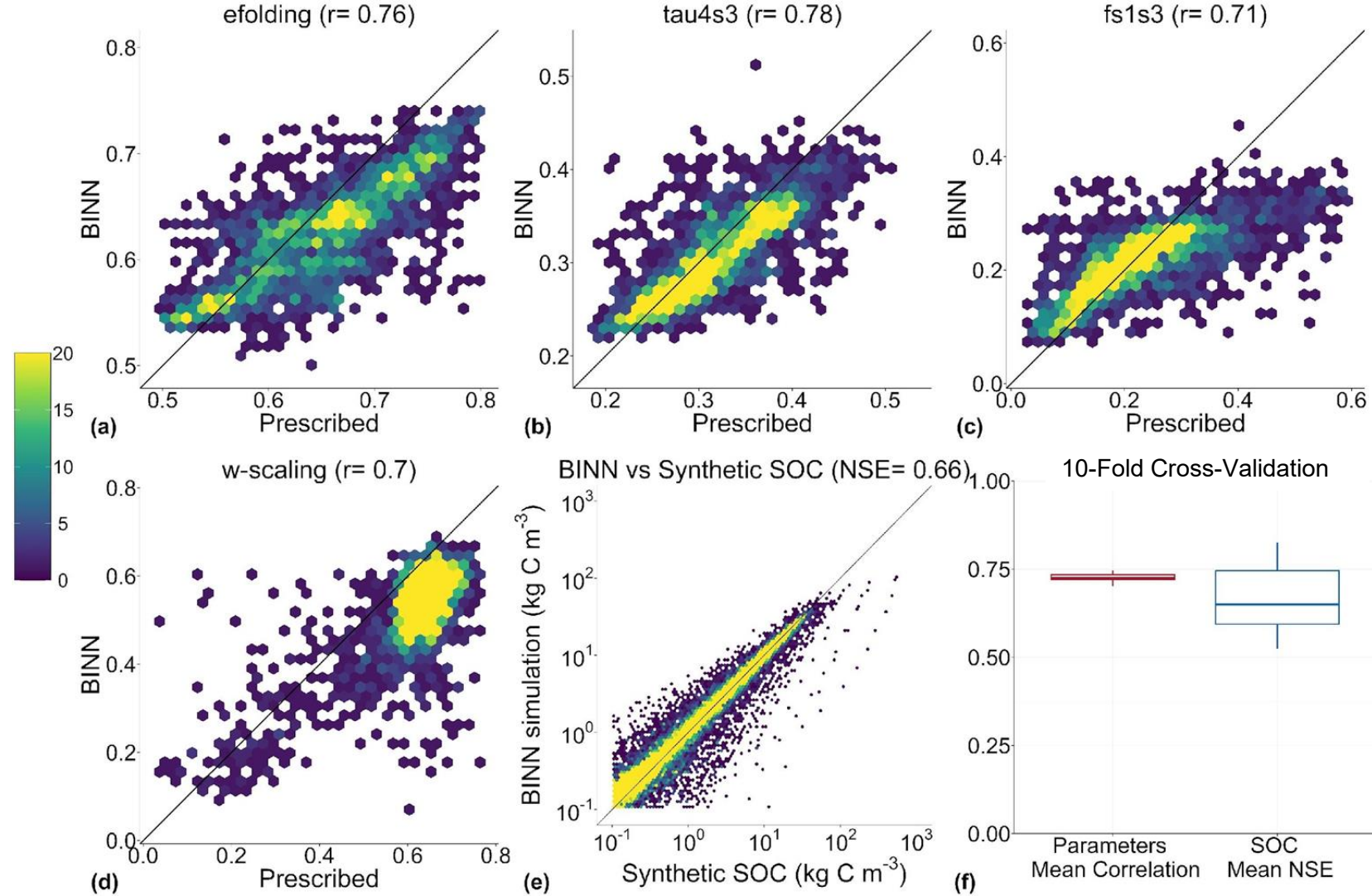
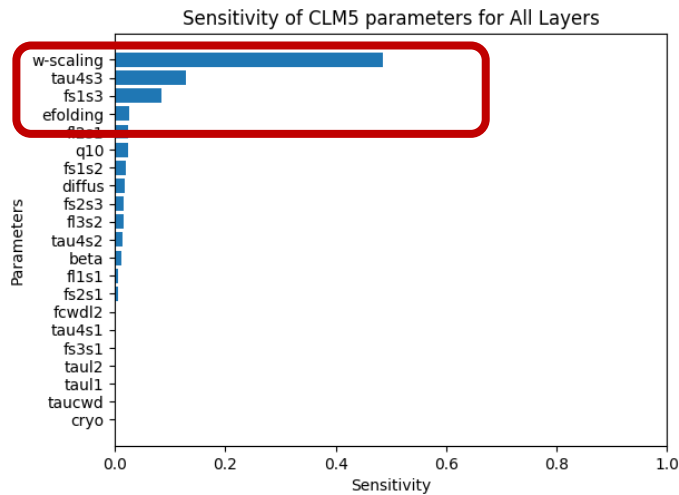
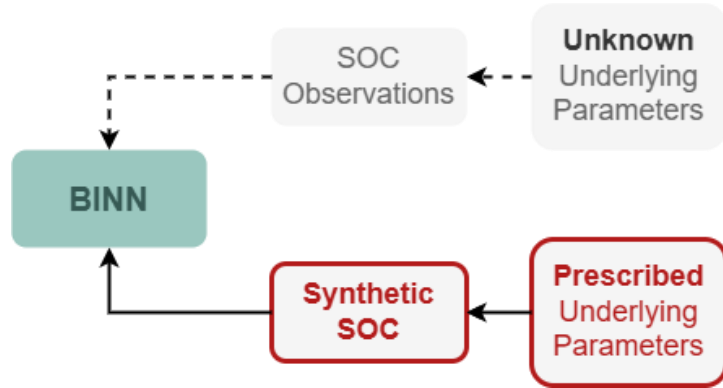


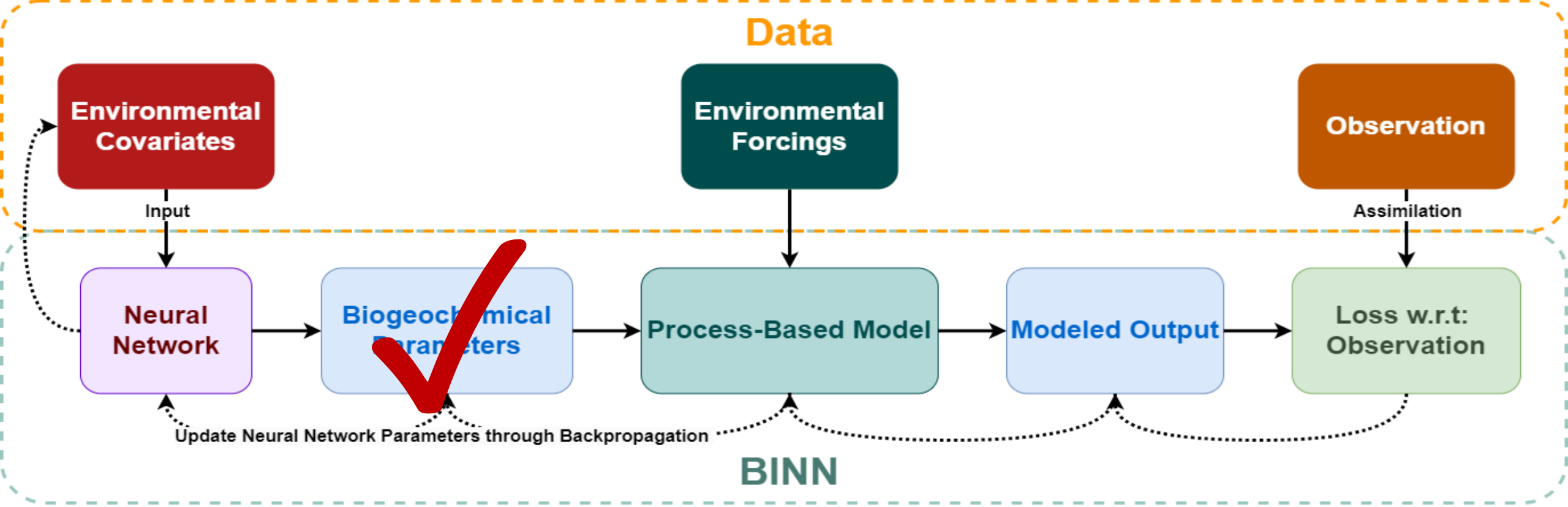
SOC Model





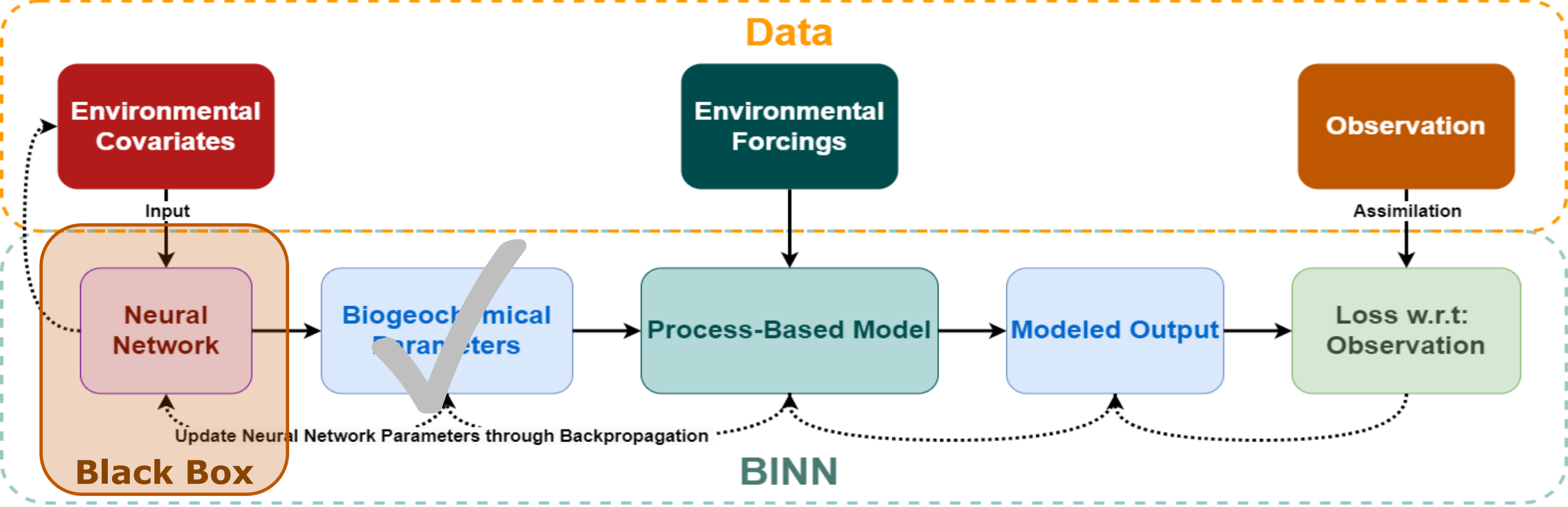
Recovery Experiment with 4 Most Influential Parameters

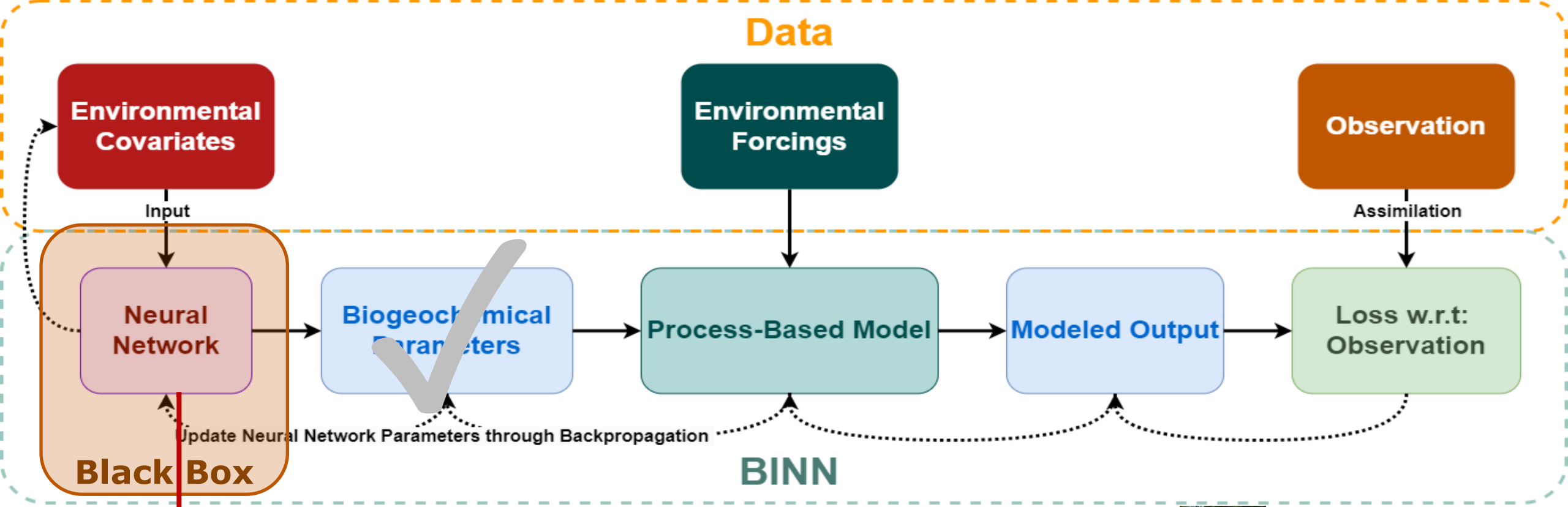




Recover causes from data
(biogeochemical parameters)

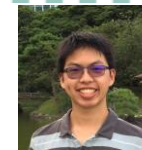
Discover emerging mechanisms governing
SOC from big data (observations)





Kolmogorov Arnold Network (KAN)

An alternative to neural networks that is easier to interpret



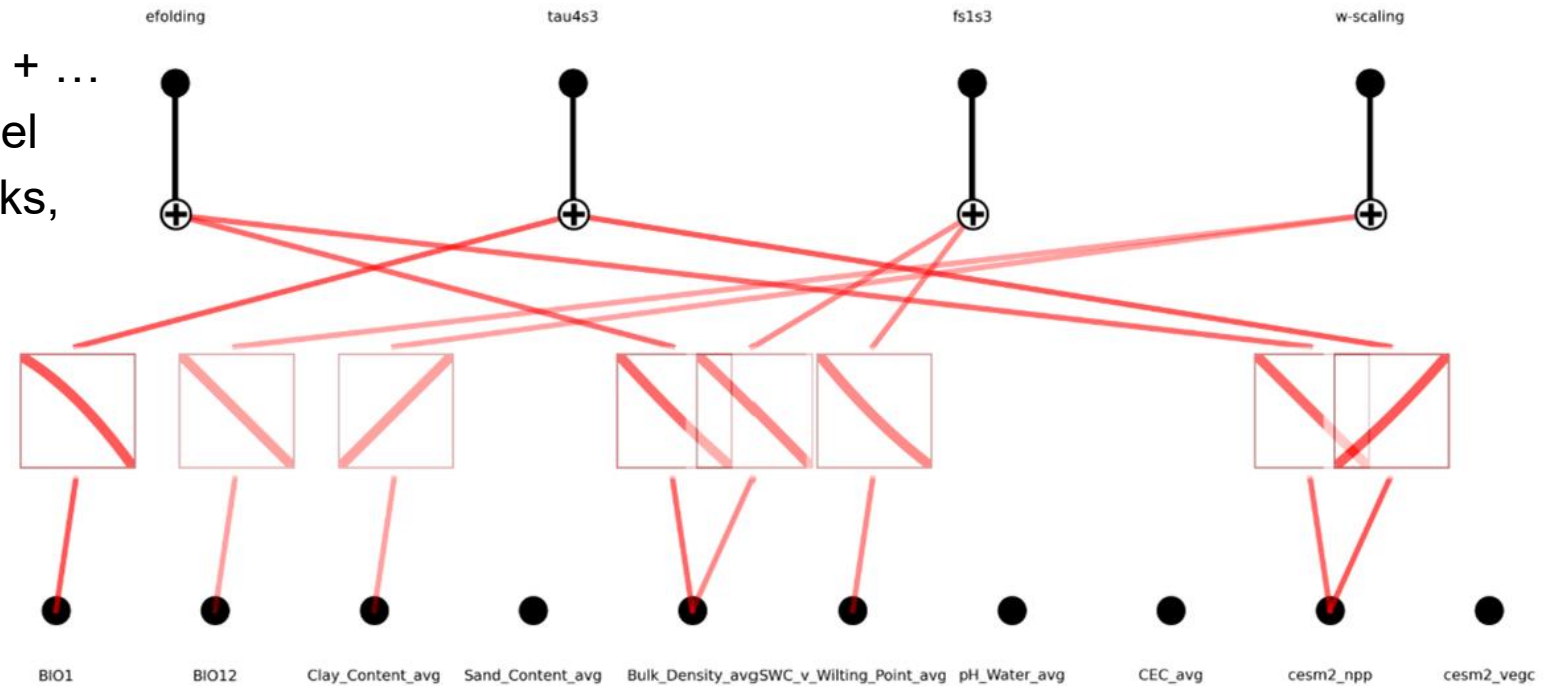
Scientifically-Interpretable Reasoning Network(SciReN)

Submitted to NeurIPS 2025

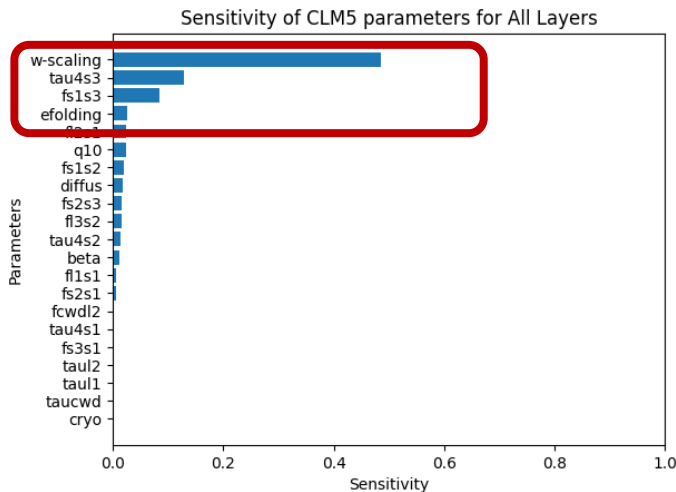
Structure of KAN

- Curves show how changing one environmental variable will affect a biogeochemical parameter
- Each output = sum of interpretable contributions from each input, e.g.
 - $\text{parameter} = (\text{input1})^2 + \exp(\text{input6}) + \dots$
- 1-layer KAN = generalized additive model
 - Not as expressive as neural networks, but similar accuracy in our case
 - Can add more layers

Outputs (Biogeochemical Parameters) - not observed

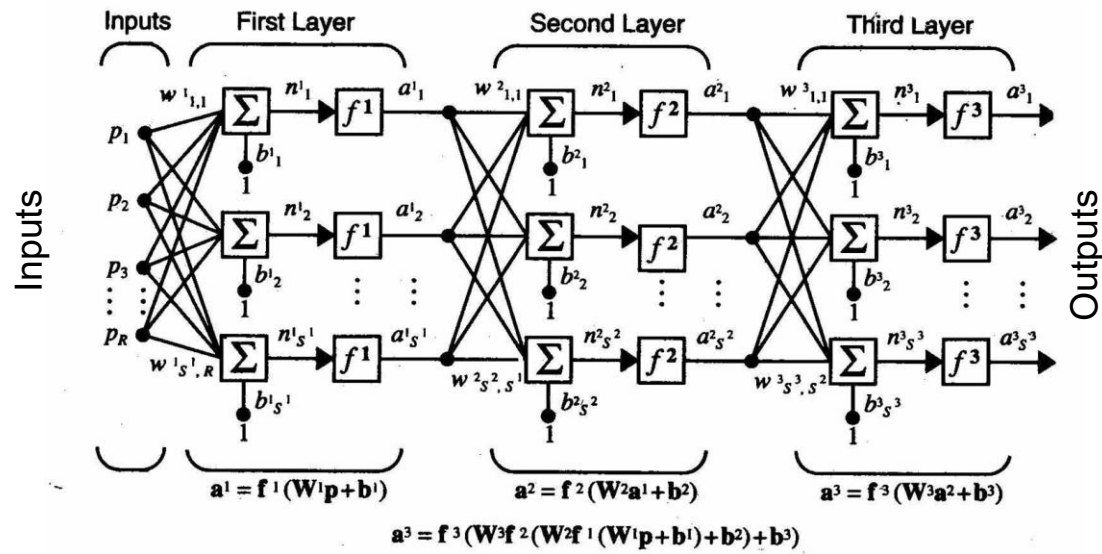


Inputs (Environmental Covariates) - observed

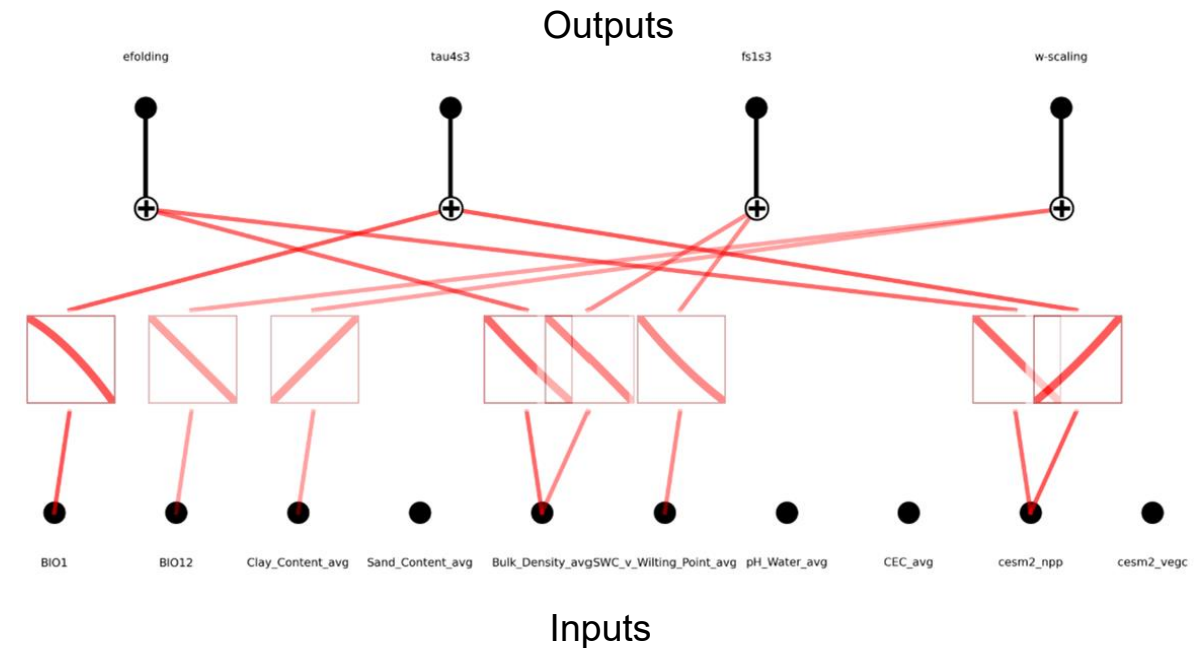


Difference between NN and KAN

- **Standard neural networks:** every input variable affects every output
 - Inputs are mixed together and passed through nonlinearities - hard to understand
 - Unclear how each input (environmental variable) influences each output (biogeochemical parameters)

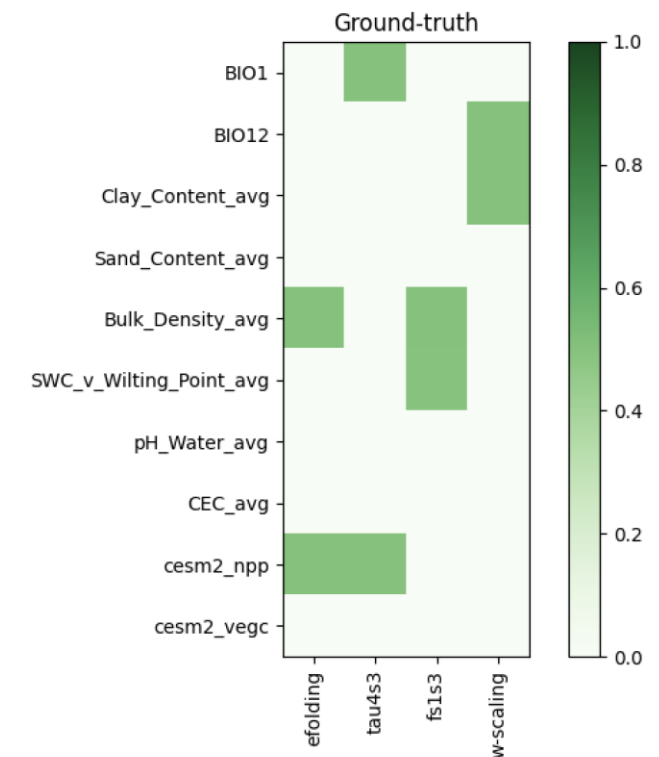
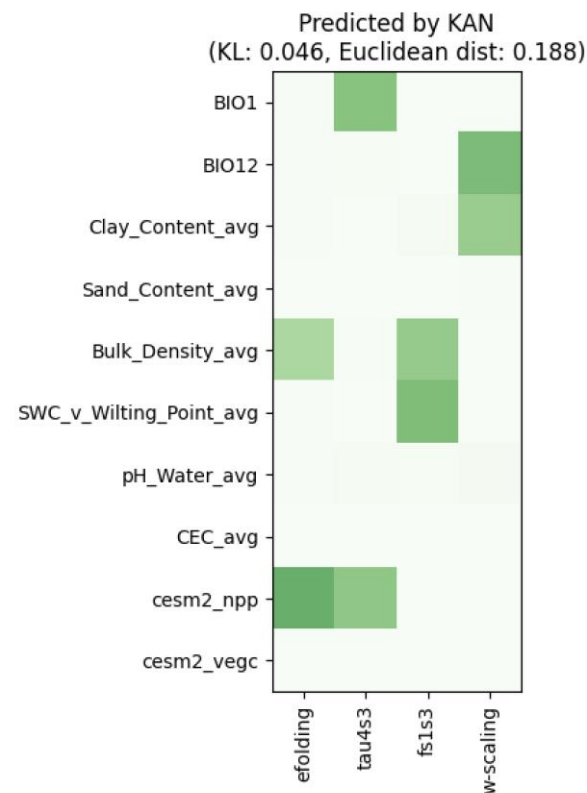
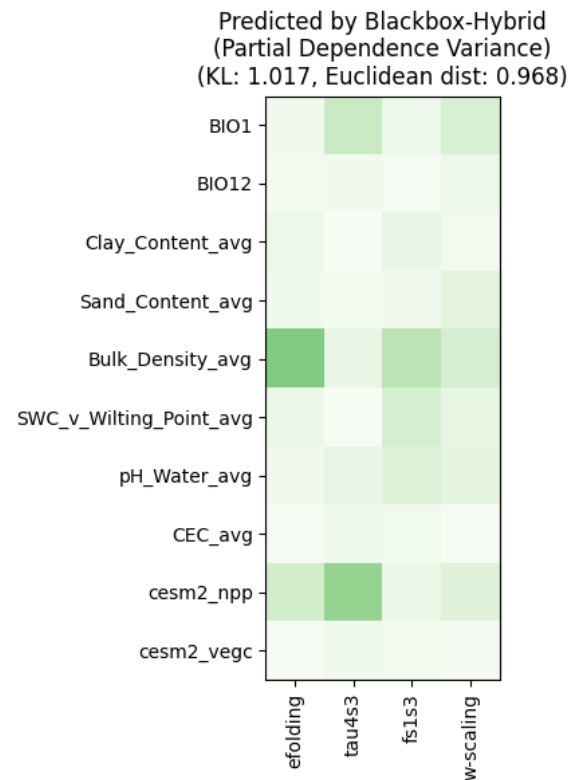


- **Kolmogorov-Arnold Networks:** learn **interpretable** functional relationships between inputs and outputs
 - Apply a 1D function on each edge (link), then add contributions from each input
 - **Sparsity:** focus on a few crucial relationships; remove unnecessary connections



Functional Relationship Retrieval Test

- Prescribed functional relationships between 10 environmental inputs and 4 most sensitive biogeochemical parameters in CLM5
- BINN does not reveal functional relationships. Even after applying a post-hoc interpretation method (Partial Dependence Variance), we find it did not even implicitly learn the correct relationships.
- SciReN (center) revealed the relationships very accurately.



In Conclusion, BINN/SciReN can ...

Recover causes from data
(biogeochemical parameters)

Discover emerging mechanisms governing
SOC from big data (observations)

Uncover novel relationships between environmental
conditions and underlying mechanisms



Thanks for Listening

