



# AI-enabled ModEx: A scalable AI workflow to efficiently calibrate PFLOTRAN process models

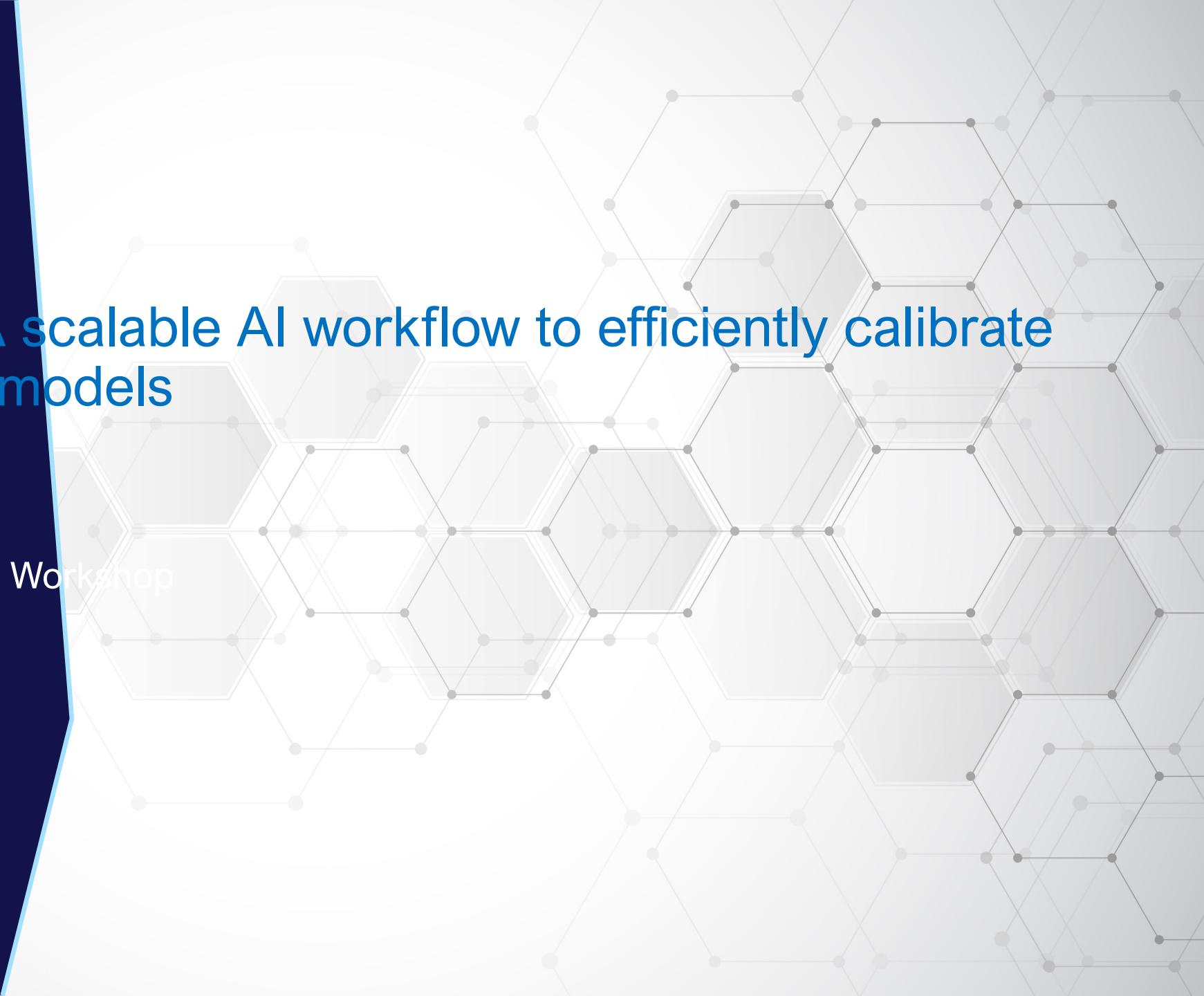
EMSL UP # 60592

2023 RemPlex Summit Training Workshop

**Maruti Mudunuru**  
Earth Scientist, PNNL

Date: Nov-17-2023

PNNL-SA-184045



# Background and relevance

- **Model** – The “PFLOTRAN reaction sandbox” provides users with a venue for implementing user-defined reactions with minimal code
- **Data** – Relevance -- Reaction networks models developed within the reaction sandbox can be leveraged to better understand carbon cycling
- **AI-enabled ModEx** – The PFLOTRAN reaction sandbox models and experimental/field data integration can be performed using AI/ML

PFLOTRAN / code / pflotran  
reaction\_sandbox\_flexbiohill.F90

Pull requests Check out

PFLOTRAN is an open source, state-of-the-art massively parallel subsurface flow and reactive transport code.

Source master ⚡ 7642e2e Full commit

pflotran / src / pflotran / reaction\_sandbox\_flexbiohill.F90

```
1 module Reaction_Sand_FlexBioHill_class
2
3 #include "petsc/finclude/petsccsys.h"
4 use petsccsys
5
6 use Reaction_Sandbox_BioHill_class
7 use PFLOTRAN_Constants_module
8
9 implicit none
10
11 private
12
13 type, public, &
14 extends(reaction_sandbox_biohill_type) :: reaction_sandbox_flexbiohill_type
15 PetscReal :: k_max
16 PetscReal :: K_Aaq_n
17 PetscReal :: K_Baq
18 PetscReal :: I_Caq
19 PetscReal :: yield
20 PetscReal :: k_decay
21 PetscReal :: n
22 PetscBool :: molarity_units
23 PetscReal, pointer :: stoich(:)
24 contains
25 procedure, public :: ReadInput => FlexBioHillReadInput
26 procedure, public :: Setup => FlexBioHillSetup
27 procedure, public :: Evaluate => FlexBioHillEvaluate
28 procedure, public :: Destroy => FlexBioHillDestroy
29 end type reaction_sandbox_flexbiohill_type
30
31 public :: FlexBioHillCreate
32
33 contains
34
35 ! *****
36
37 function FlexBioHillCreate()
38 !
39 ! Allocates flexible biodegradation reaction object.
40 !
41 implicit none
42
43 class(reaction_sandbox_flexbiohill_type), pointer :: FlexBioHillCreate
```

Example – Reaction Sandbox  
F90 code for implementing  
user-defined reactions

Example – PFLOTRAN input deck for numerical simulations

```
1 # Description: Reaction Sandbox Simple (A + B --> C + D)
2 # Number of reaction network parameters: D = 7
3 # Sobol-based GSA: N * (2D+2) samples for first and second-order indices
4
5 =====
6 SIMULATION
7   SIMULATION_TYPE SUBSURFACE
8   PROCESS_MODELS
9     SUBSURFACE_TRANSPORT transport
10    MODE GIRL
11  /
12  /
13 END
14 =====
15 SUBSURFACE
16
17 ===== constraints =====
18 # modify these initial concentration
19 CONSTRAINT initial
20 CONCENTRATIONS # [mol/L]
21   Aaq 1.d-3 F
22   Baq 5.d-4 F
23   Caq 1.d-10 F
24   Daq 1.d-10 F
25   Eq 1.d-10 F
26  /
27 IMMOBILE # [mol/m^3 bulk]
28   Xim 1.d-4
29   Yim 1.d-10
30  /
31 END
32
33 ===== chemistry =====
34 CHEMISTRY
35 PRIMARY_SPECIES
36   Aaq
37   Baq
38   Eq
39   Caq
40   Daq
41  /
42 IMMOBILE_SPECIES
43   Yim
44   Xim
45  /
```

# Why AI-enabled Model-Experiment-Data integration (ModEx)?

## ▪ Background

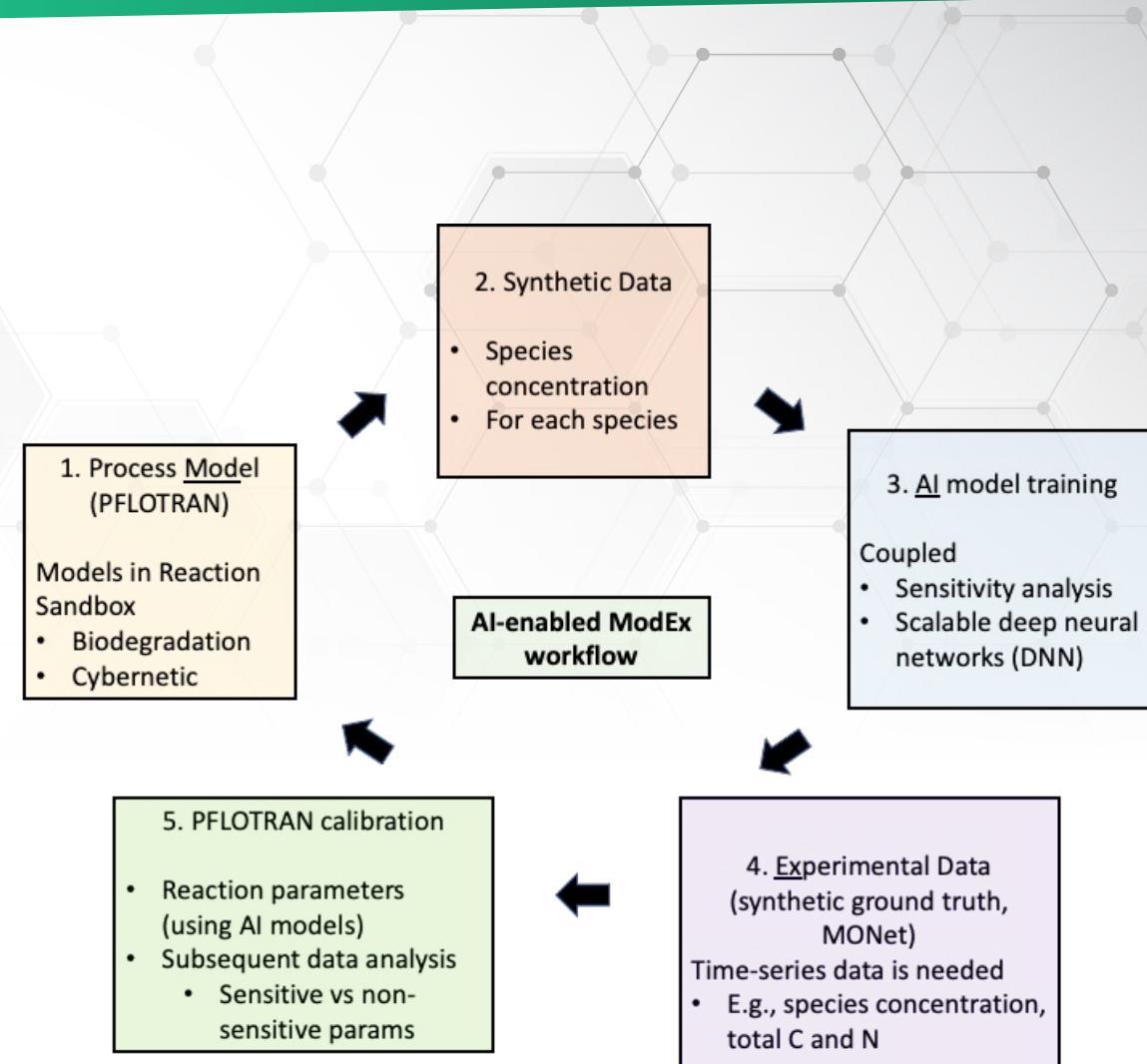
- Provide calibrated PFLOTRAN reaction process models that can work with field or experimental data (e.g., carbon degradation)
- Sensitivity analysis to identify important parameters
- Model-Data integration using AI/ML

## ▪ End-to-end workflow outcomes

- An open-source and user-friendly AI/ML workflow available for users to analyze field/experimental data
- A skeleton that can be easily adapted for future workflows
  - Other process models relevant to coupled flow, transport, and ERT field data
  - AI/ML modeling for PFLOTRAN users

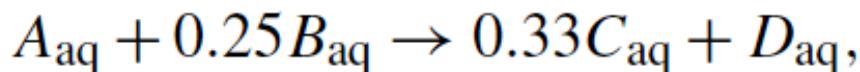
## ▪ Technical readiness plan (proof of concept, development, deployment)

- Proof-of-concept (TRL-3)
- AI/ML workflows for PFLOTRAN users (TRL-3)



# PFLOTRAN reaction sandbox – Biomass degradation process model

- **Carbon degradation** – Microbially mediated carbon biodegradation reaction with biomass growth and decay over time



Electron donor  
(e.g., carbon compounds)

Electron acceptor  
(e.g., oxygen)

- **Model** – The rate of biomass growth and decay over time can be modeled as

$$\frac{dX}{dt} = yield_{X_{\text{im}}} I_r - k_{\text{decay}} X$$

Change in concentration

Yield

Biomass decay

- **Reaction rate** – Monod expression + Hill function

$$I_r = k_{\max} X_{\text{im}} \frac{A_{\text{aq}}^n}{K_{A_{\text{aq}}}^n + A_{\text{aq}}^n} \times \frac{B_{\text{aq}}}{K_{B_{\text{aq}}} + B_{\text{aq}}} \times \frac{I_{C_{\text{aq}}}}{I_{C_{\text{aq}}} + C_{\text{aq}}}$$

- Process model parameters in this reaction network that are varied

1. Max Specific Utilization Rate –  $k_{\max}$ 
    - [mole mole<sup>-1</sup> biomass s<sup>-1</sup>]
  2. Aaq Half Saturation Constant –  $K_{A_{\text{aq}}}$  [M]
  3. Baq Half Saturation Constant –  $K_{B_{\text{aq}}}$  [M]
  4. Caq Monod Inhibition Constant –  $K_{C_{\text{aq}}}$  [M]
  5. Yield –  $yield_{X_{\text{im}}}$  [mole<sub>biomass</sub> mole<sup>-1</sup>]
  6. Biomass Decay Rate Constant –  $k_{\text{decay}}$  [s<sup>-1</sup>]
  7. Hill Exponent –  $n$  [-]
- Inhibitor concentration –  $I_{C_{\text{aq}}}$  [M] is not varied
  - Immobile biomass species –  $X_{\text{im}}$  [mole<sub>biomass</sub> m<sup>-3</sup><sub>bulk</sub>] is not varied

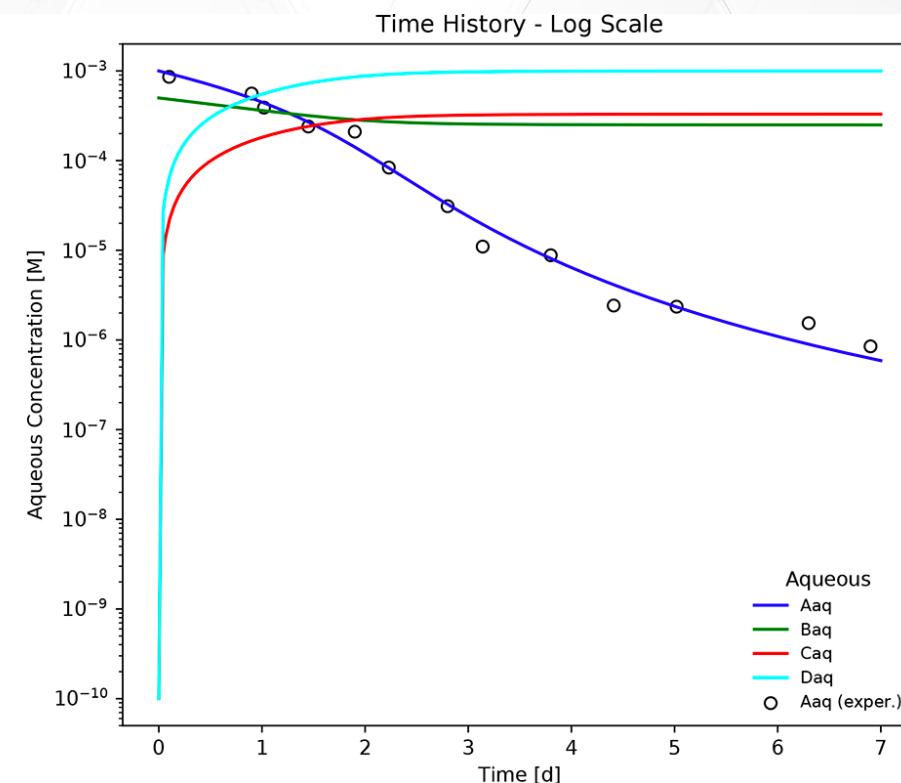
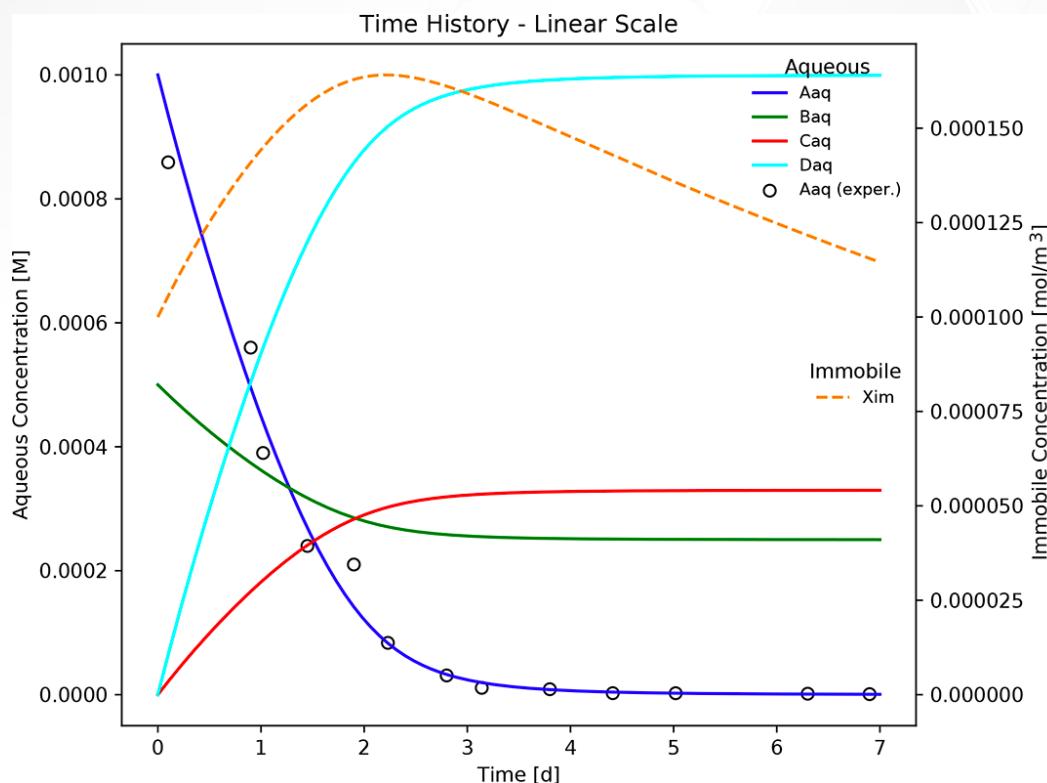
Hammond et al., PFLOTRAN reaction sandbox,  
2022

# Data requirements – What is needed to calibrate PFLOTRAN reaction sandbox? (1/2)

- Data – What are the data requirements for calibrating such reaction network parameters?

- We need time-series data on species concentration (i.e., organic carbon and major cation/anion concentrations)
- We will use synthetic time-series to develop and test the workflow

For example, incase of the biomass degradation process model



## Data requirements – What is needed to calibrate PFLOTRAN reaction sandbox? (2/2)

- Data – What are the data requirements for calibrating such reaction network parameters?

- We need time-series data on species concentration (i.e., organic carbon and major cation/anion concentrations)
- We will use synthetic time-series to develop and test the workflow

**Table 2.** Microbially mediated reaction parameters for the batch biodegradation experiment.  $n$  only applies to the reaction incorporating the Hill function (i.e., Eq. 15). M signifies molarity or mole per liter of water.

Parameter	Value	Units
$k_{\max}$	$9 \times 10^{-2}$	mole mole $^{-1}_{\text{biomass}}$ s $^{-1}$
$K_{A_{\text{aq}}}$	$2 \times 10^{-4}$	M
$K_{B_{\text{aq}}}$	$1.25 \times 10^{-5}$	M
$I_{C_{\text{aq}}}$	$2.5 \times 10^{-4}$	M
$\text{yield}_{X_{\text{im}}}$	$1 \times 10^{-4}$	mole biomass mole $^{-1}$
$k_{\text{decay}}$	$1 \times 10^{-6}$	1 s $^{-1}$
$n$	1.2	–

**Table 3.** Initial concentrations for the batch biodegradation experiment.

Species	Concentration	Units
$A_{\text{aq}}$	$1 \times 10^{-3}$	M
$B_{\text{aq}}$	$5 \times 10^{-4}$	M
$C_{\text{aq}}$	$1 \times 10^{-10}$	M
$D_{\text{aq}}$	$1 \times 10^{-10}$	M
$X_{\text{im}}$	$1 \times 10^{-4}$	mole m $^{-3}_{\text{bulk}}$

Microbially mediated biodegradation reaction with biomass growth and decay over time: Hammond et al., PFLOTRAN reaction sandbox, 2022

# AI-enabled ModEx: Overall summary

AI-enabled ModEx pipeline that will be available to PFLOTRAN users has five steps

## 1. Microbial reaction specifics with synthetic data

- Examples from PFLOTRAN reaction sandbox

## 2. PFLOTRAN simulation data generation

- Sobol-based sampling for realizations
  - $N * (2D + 2)$  realizations = 16000

## 3. Sensitivity analysis (Local, Global, Obs-based)

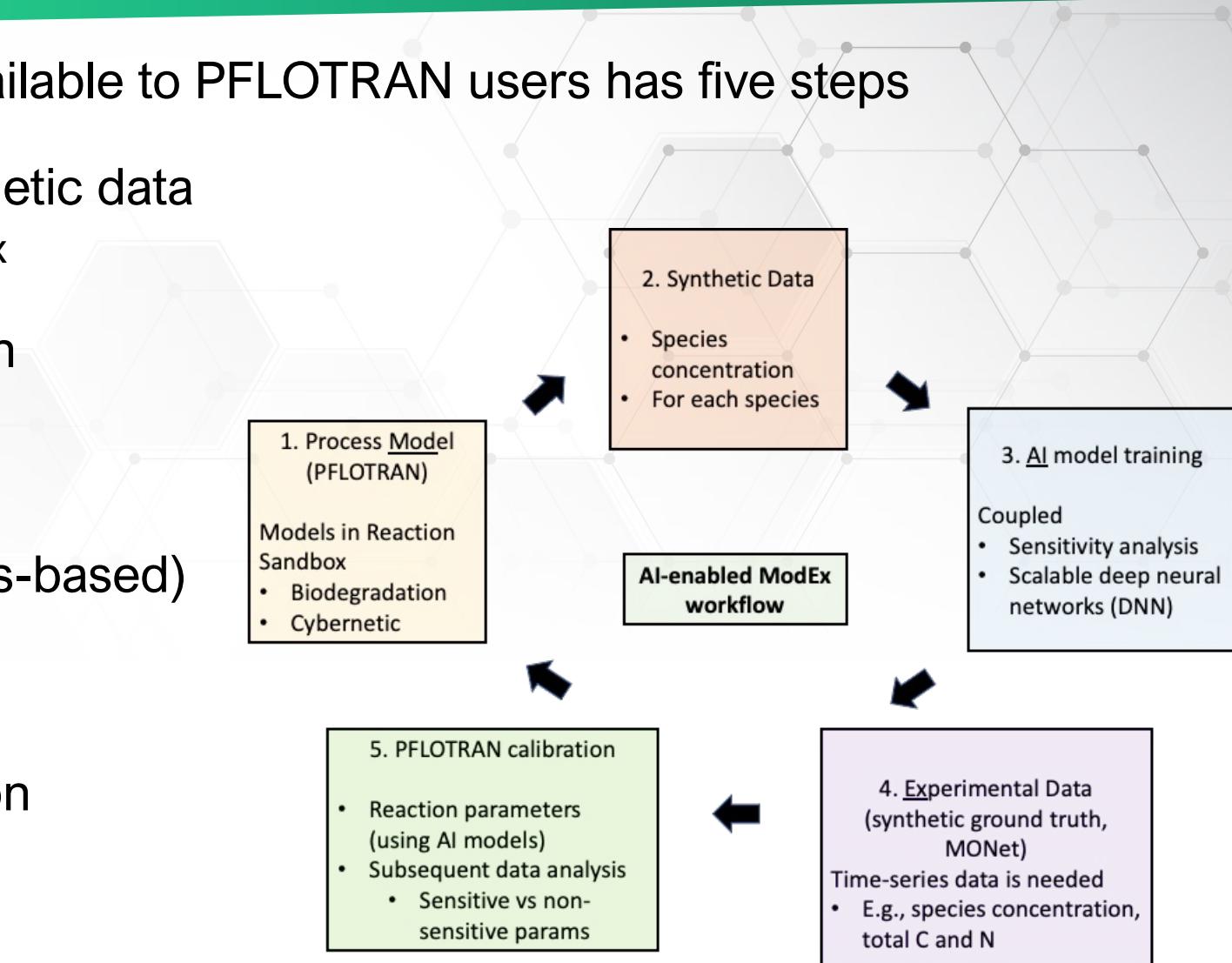
- F-test, Mutual Information
- Random Forests, SHAPley values

## 4. Reaction network parameter estimation

- AI model training
- With and without noise in data

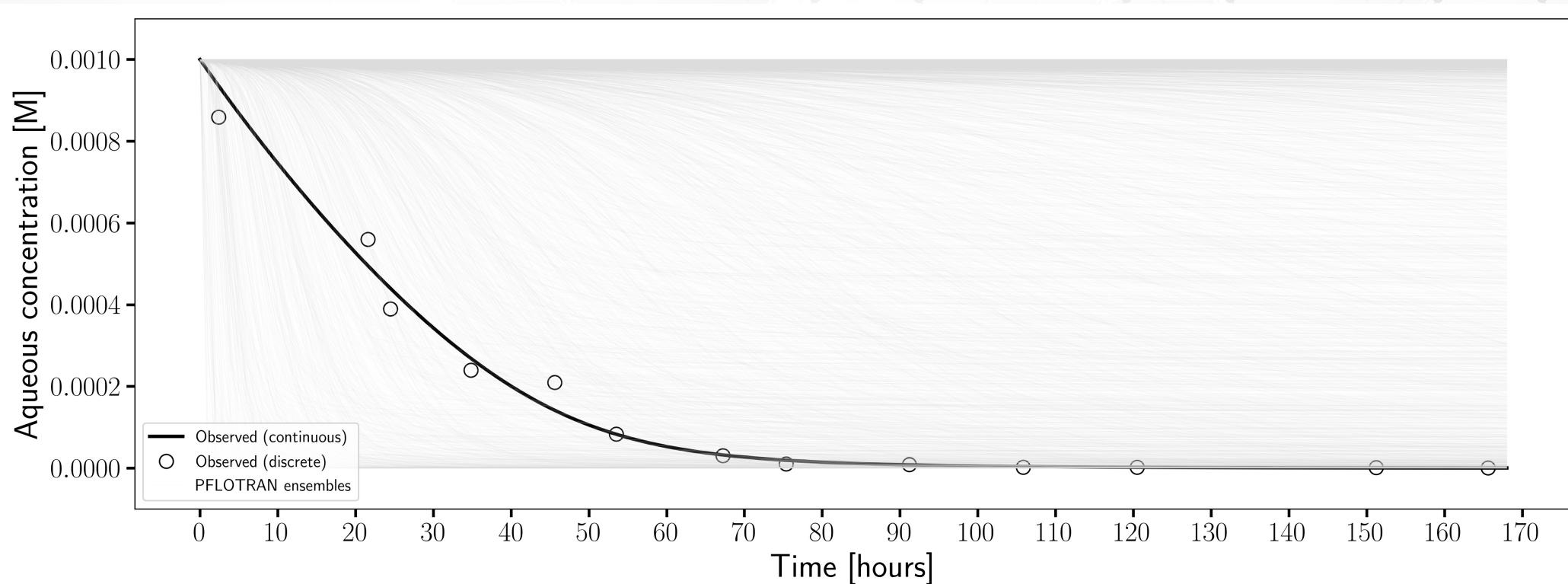
## 5. Comparison with synthetic data

- Performance metrics (e.g., R<sup>2</sup>-score, MSE)
- Predictive uncertainty



# Step-1 and Step-2 – PFLOTRAN simulations

- **Carbon degradation** – Microbially mediated carbon biodegradation reaction with biomass growth and decay over time
- **Simulation data generation**
  - 16000 Sobol-sampling realizations
  - 15539 realizations ran to completion



Process model parameters in this reaction network

1. Max Specific Utilization Rate
2. Aaq Half Saturation Constant
3. Baq Half Saturation Constant
4. Caq Monod Inhibition Constant
5. Yield
6. Biomass Decay Rate Constant
7. Hill Exponent

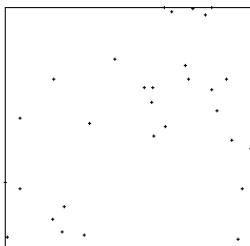
# What is a Sobol Sequence?

Sobol sequence is a quasi-random numbers generated by Sobol's algorithm in place of (pseudo-)random numbers.

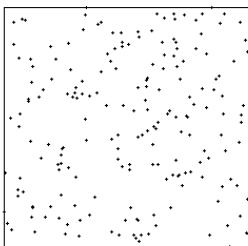
Quasi-random numbers offer a lower discrepancy (they fill the space of possibilities more evenly) resulting in a faster convergence and more stable estimates of quantities of interest and sensitivities

SALib: <https://salib.readthedocs.io/en/latest/>

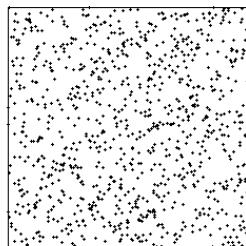
Random - 32 points



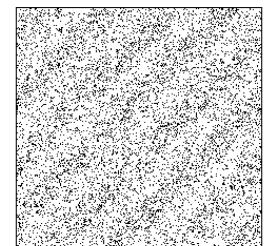
Random - 200 points



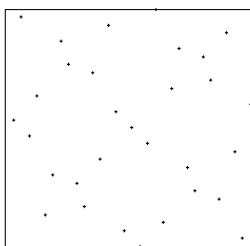
Random - 1000 points



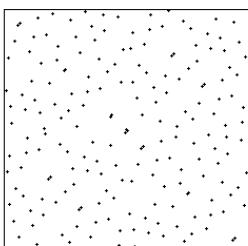
Random - 10.000 points



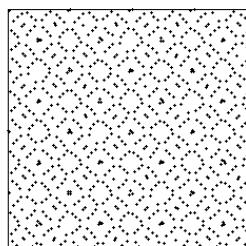
Sobol - 32 points



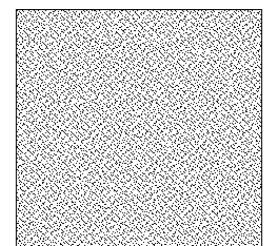
Sobol - 200 points



Sobol - 1000 points



Sobol - 10.000 points



A screenshot of the SALib documentation website. The header features the SALib logo and navigation links for Getting started, Basics, SALib Interface, Advanced, Wrappers, and More. A search bar and a version indicator (1.4.7 (stable)) are also present. The main content area displays the title "SALib - Sensitivity Analysis Library in Python" and a brief description of the library's purpose and applications. It includes a "Supported Methods" section listing various sensitivity analysis techniques. On the right side, there are links for "On this page", "Supported Methods", "Getting Started", "For Developers", "Other Info", and a "Show Source" button.

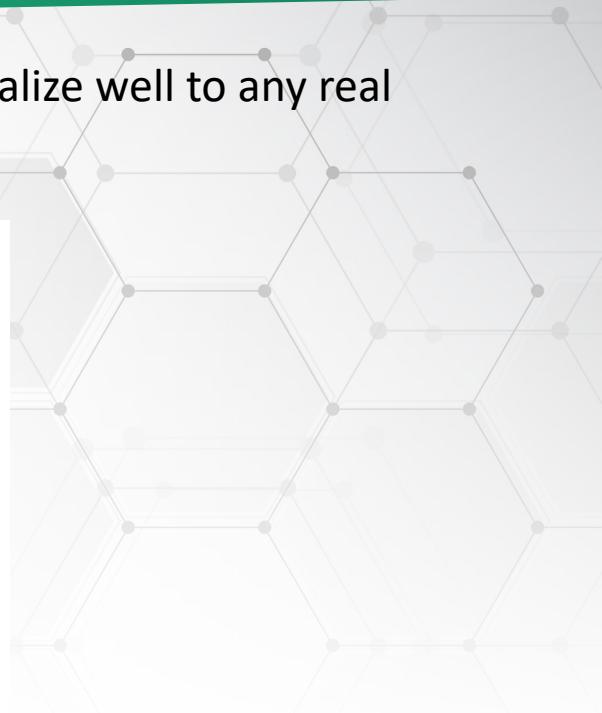
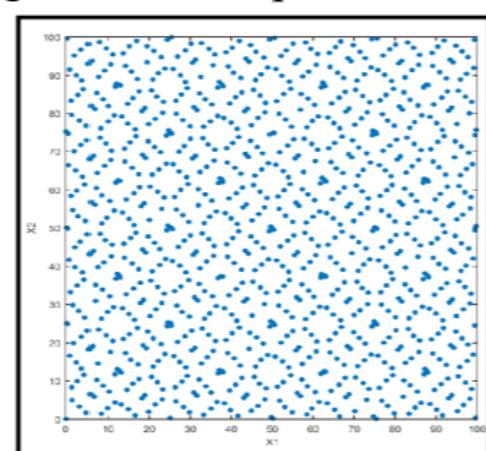
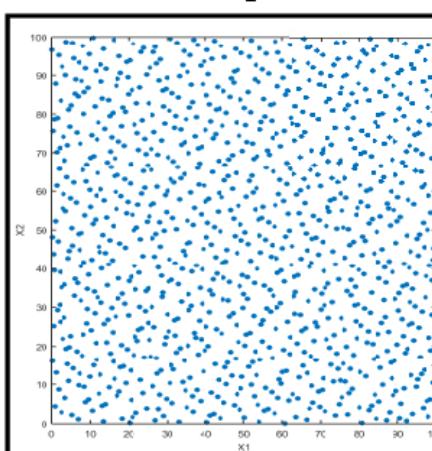
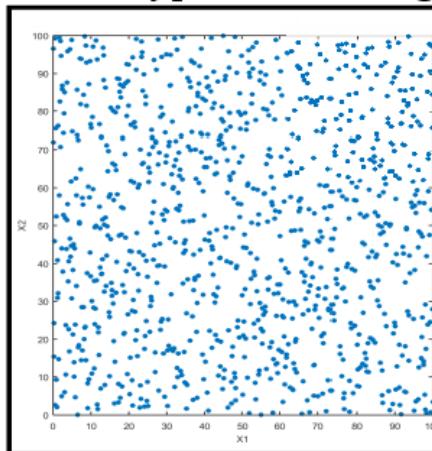
## Supported Methods

- Sobol Sensitivity Analysis ([Sobol 2001](#), [Saltelli 2002](#), [Saltelli et al. 2010](#))
- Method of Morris, including groups and optimal trajectories ([Morris 1991](#), [Campolongo et al. 2007](#))
- Fourier Amplitude Sensitivity Test (FAST) ([Cukier et al. 1973](#), [Saltelli et al. 1999](#))
- Random Balance Designs - Fourier Amplitude Sensitivity Test (RBD-FAST) ([Tarantola et al. 2006](#), [Elmar Plischke 2010](#), [Tissot et al. 2012](#))
- Delta Moment-Independent Measure ([Borgonovo 2007](#), [Plischke et al. 2013](#))
- Derivative-based Global Sensitivity Measure (DGSM) ([Sobol and Kucherenko 2009](#))
- Fractional Factorial Sensitivity Analysis ([Saltelli et al. 2008](#))
- High Dimensional Model Representation ([Li et al. 2010](#))
- PAWN ([Pianosi and Wagener 2018](#), [Pianosi and Wagener 2015](#))
- Regional Sensitivity Analysis (based on [Hornberger and Spear, 1981](#), [Saltelli et al. 2008](#), [Pianosi et al., 2016](#))

# Alternate sampling strategies and low-discrepancy sequences

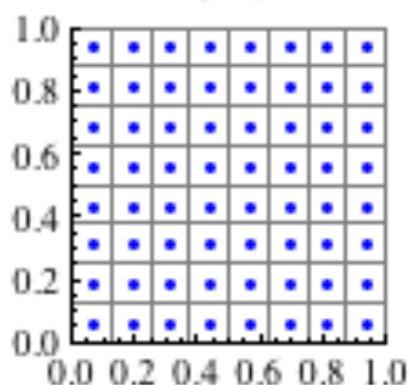
Low-discrepancy sequences are an excellent tool for evenly sampling a sample space. They generalize well to any real space, if you use an appropriate mapping.

○ Latin hypercube design ○ Halton sequence design ○ Sobol sequence design

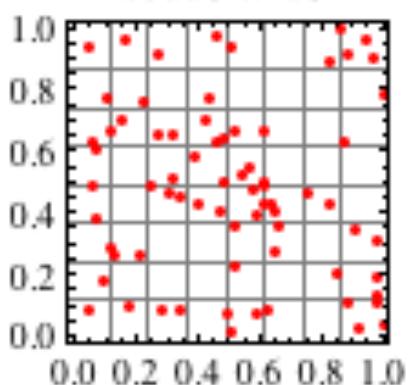


Alternative 2D sequences

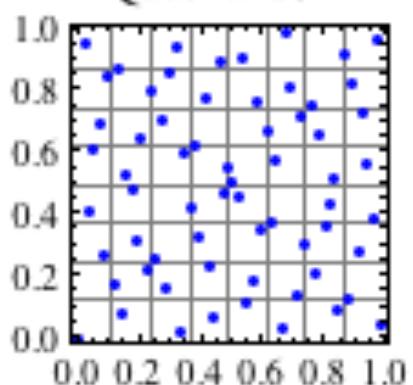
Grid



Pseudorandom



Quasirandom



# Step-3: Comprehensive sensitivity analysis (1/3)

## ▪ Local sensitivity analysis

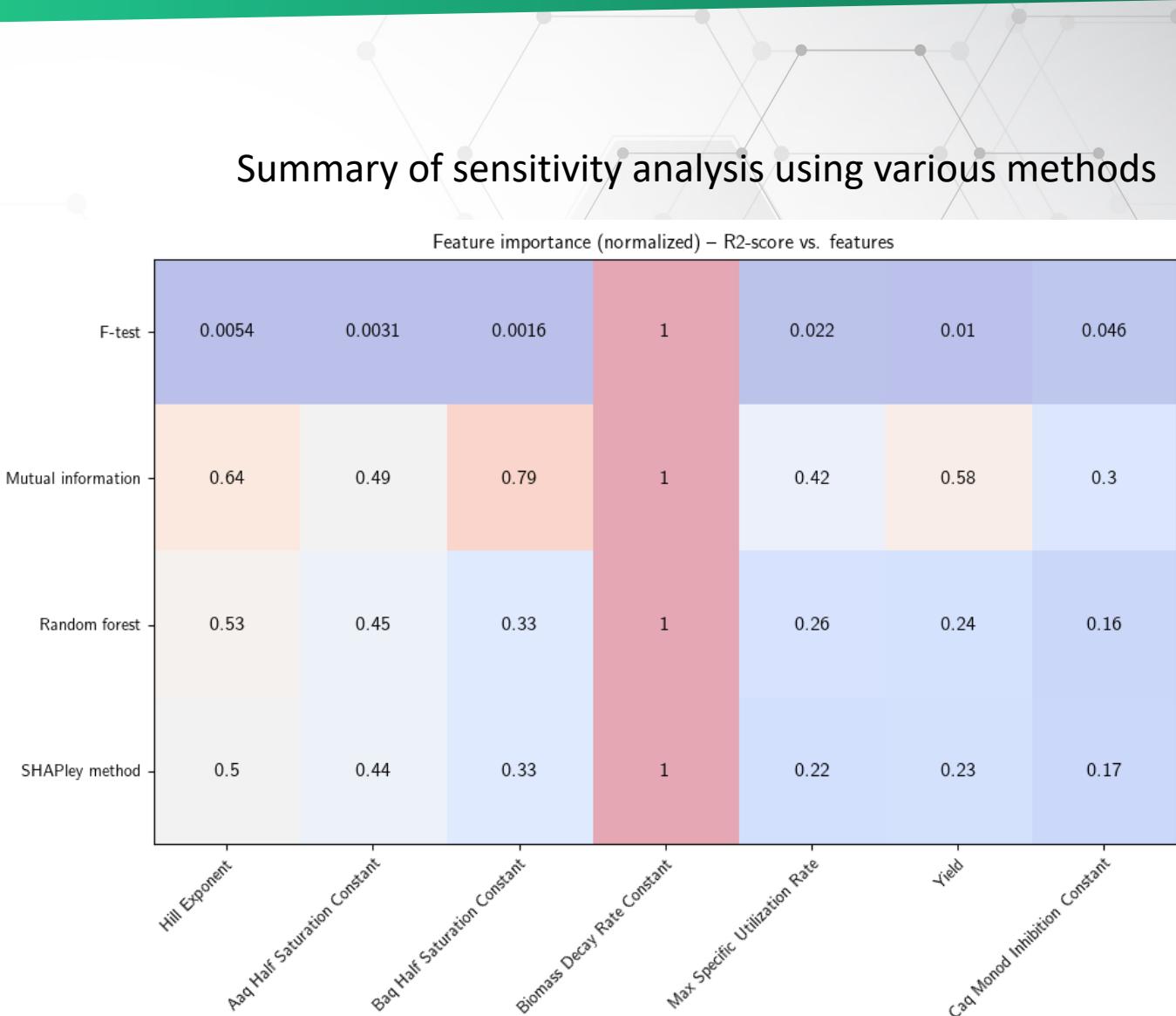
- F-test
- Mutual information (MI)

## ▪ Global sensitivity analysis

- Averaged local sensitivities

## ▪ Sensitivity analysis representative of observations

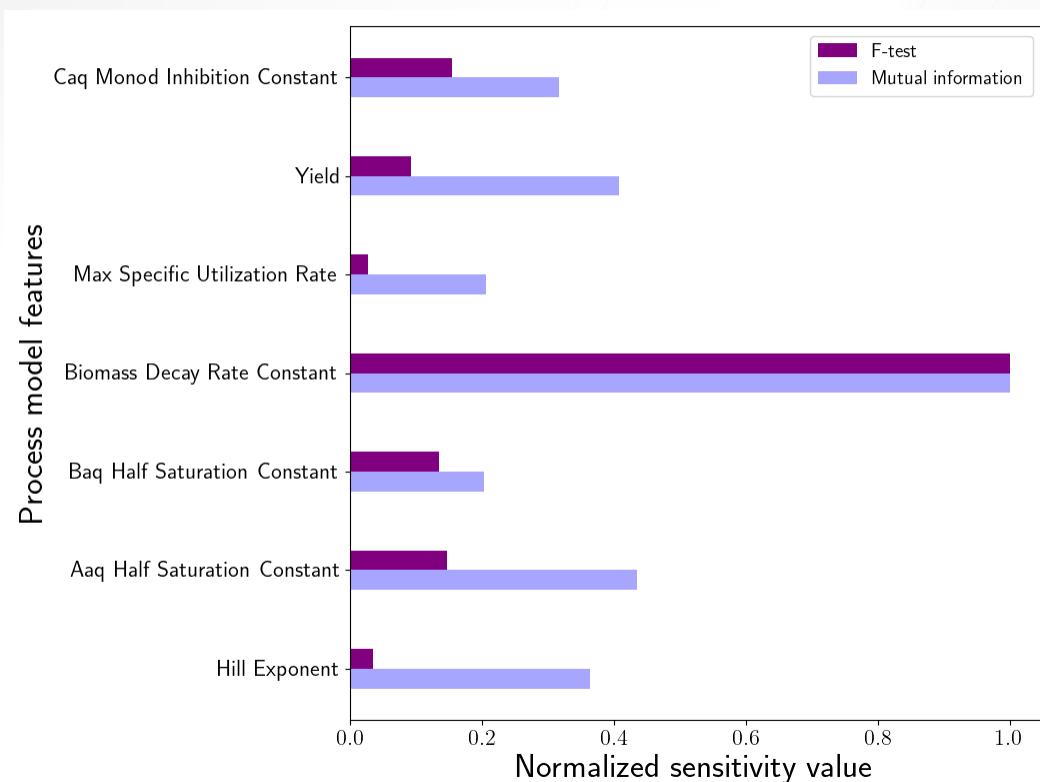
- F-test
- MI
- Random Forest
- SHAPley values
- Performance metrics
  - $R^2$ , NSE, logNSE, KGE, mKGE, npKGE



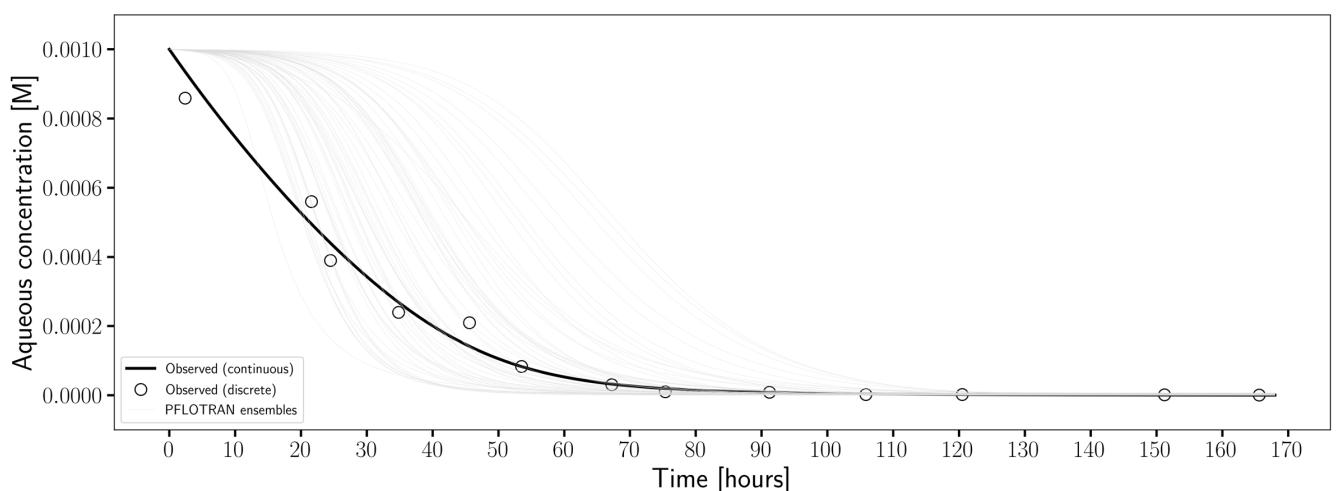
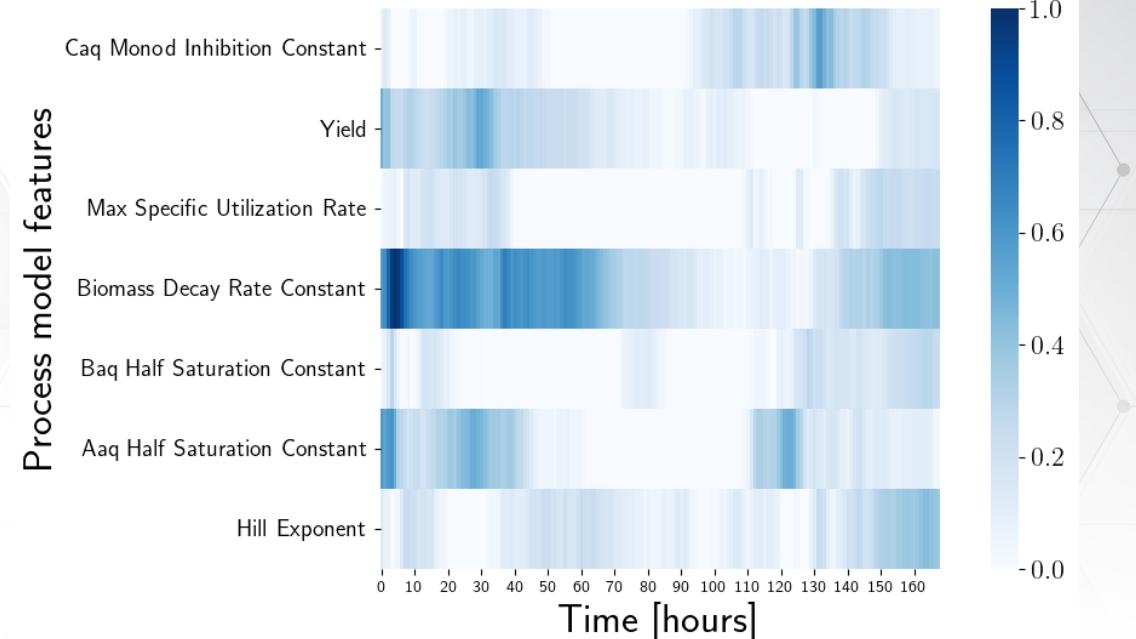
## Step-3: Comprehensive sensitivity analysis (2/3)

### Sensitivity analysis representative of observations

- F-test, MI, Random Forest, SHAPley values
- Performance metrics
  - $R^2$ , NSE, logNSE, KGE, mKGE, npKGE



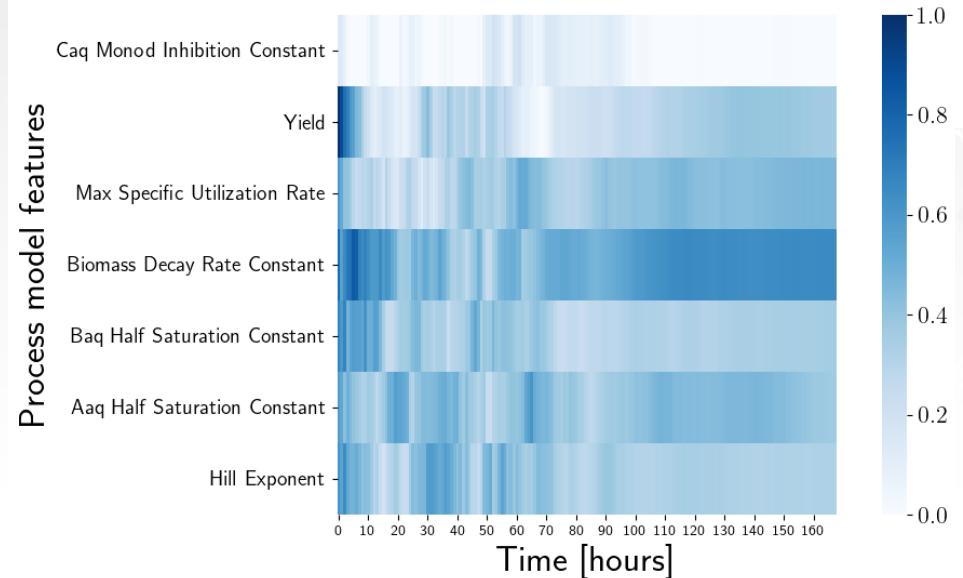
Sensitivity (MI) for all realizations with  $\log\text{NSE}$ -score  $> 0.5$



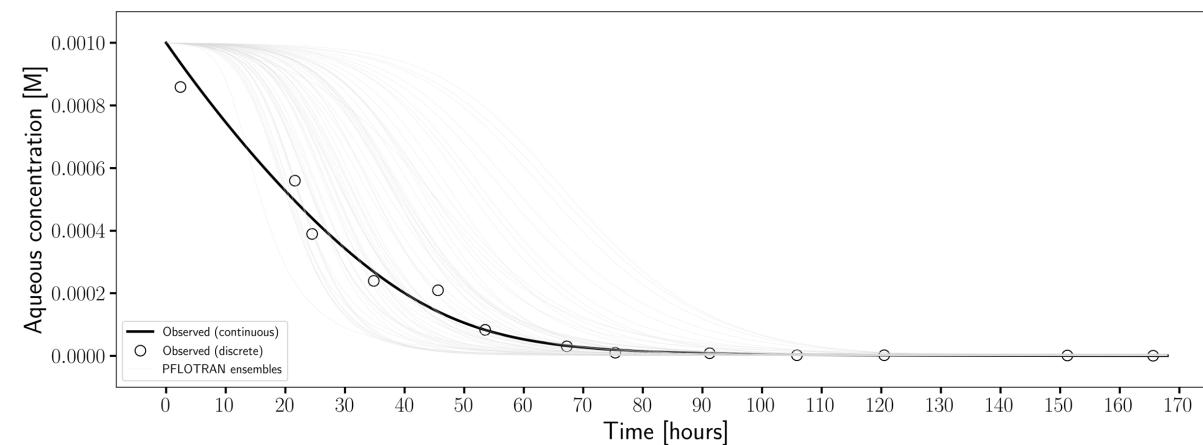
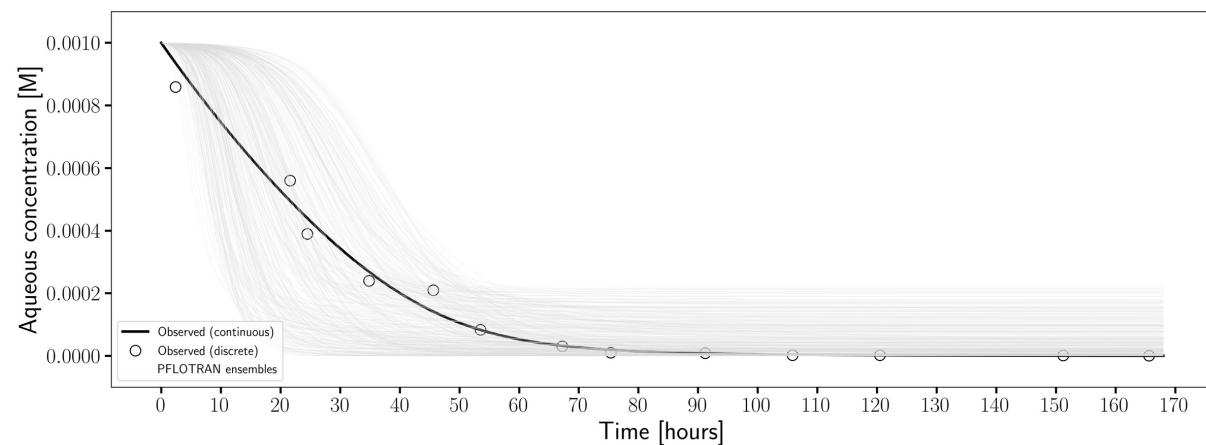
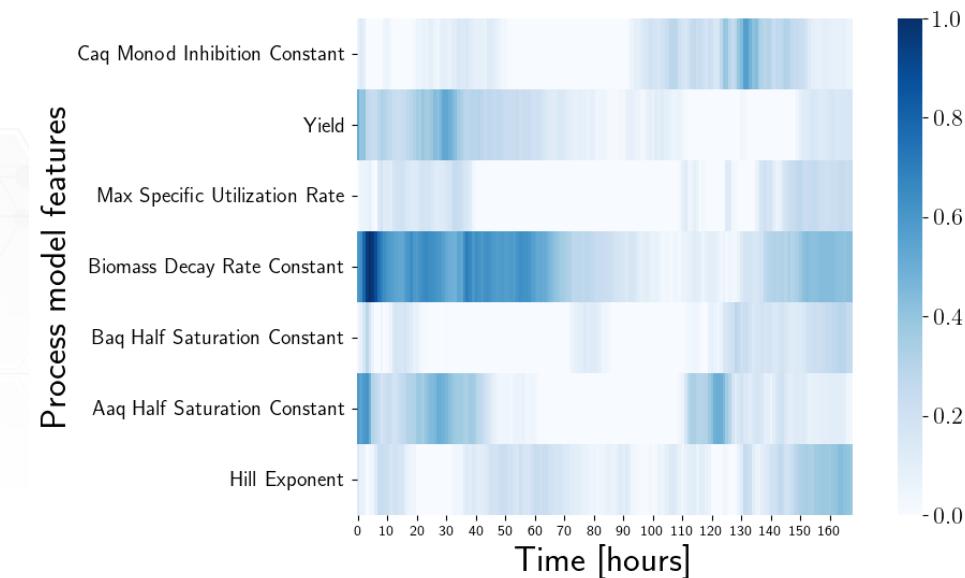
## Step-3: Comprehensive sensitivity analysis (3/3)

- Local sensitivity analysis across realizations that are closer to the observational data based on  $R^2$ -score or  $\logNSE > 0.5$

Sensitivity (MI) for all realz with  $R^2$ -score > 0.5



Sensitivity (MI) for all realz with  $\logNSE$ -score > 0.5



# Step-4 – ML model training and estimation of PFLOTRAN reaction network parameters

1. Species time-series data

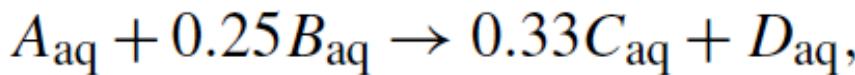


2. Data standardization



## 3. AI model – ML or DL + Hyperparameter tuning (on Tahoma)

- Single task learning (STL)
  - Each parameter → One model
  - Seven different models
- Multi-task learning (MTL)
  - All parameters → One model



Electron  
donor

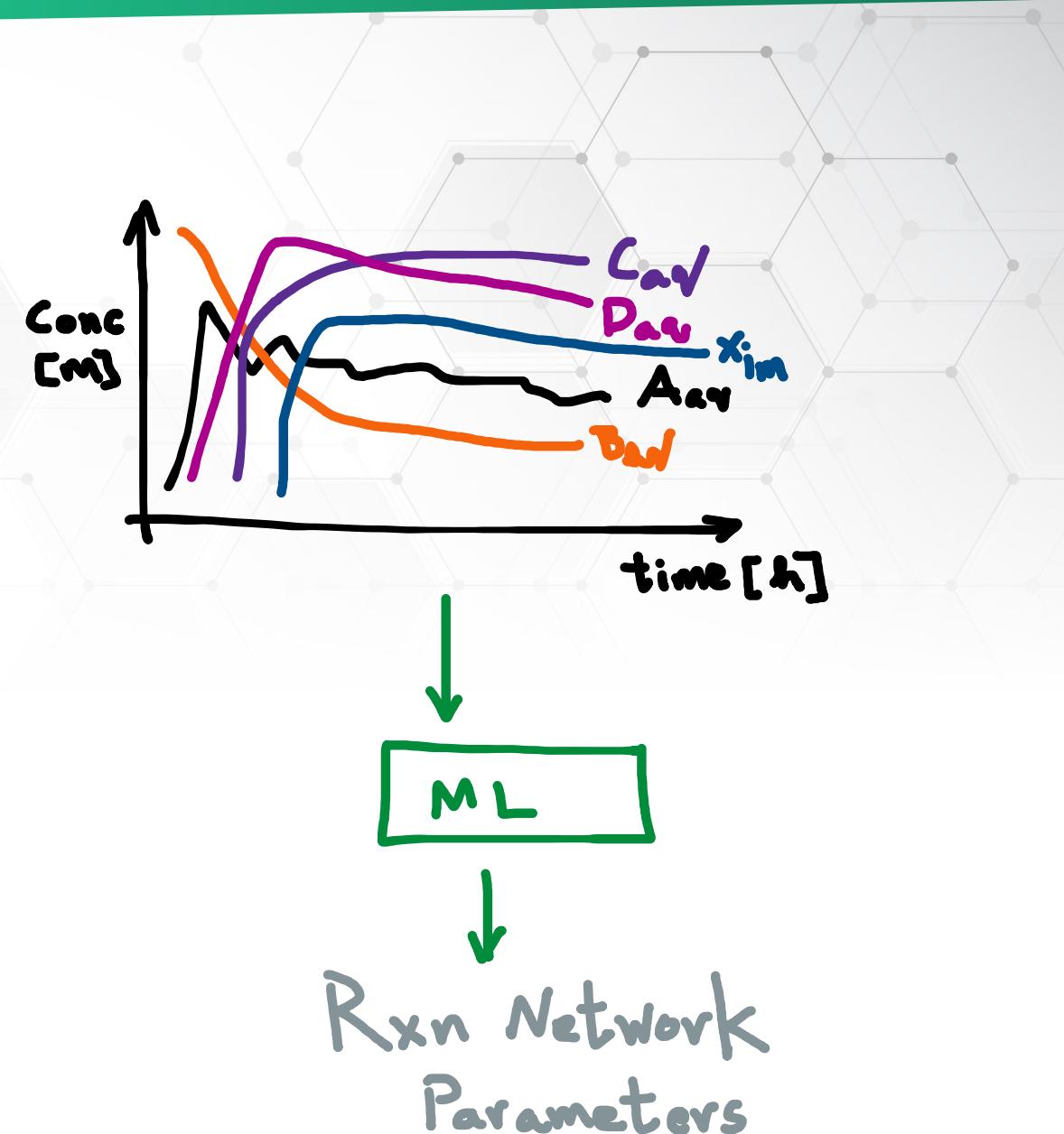
Electron  
acceptor

4. PFLOTRAN  
reaction  
network  
parameters

## Step-4 – ML model training and estimation of PFLOTRAN reaction network parameters

Inputs and outputs of ML/DL model training

- Input to ML/DL model
  - Time-series data of chemical species
- Output to ML/DL model
  - Reaction network parameters



## Step-4 – Component specifics

1. Species time-series data



- Aaq time-series
- Aaq + Baq + Caq + Daq + Xim
- Standard, MinMax, MaxAbs, Robust, Quantile Transformer, PowerTransformer Scalers
- Log10 transformation
- Traditional ML – Random Forests, SVM
- Deep learning (DL) –
  - Deep neural networks (DNN)
  - Convolutional neural networks (CNN)
- Reaction sandbox parameters
  1. Max Specific Utilization Rate
  2. Aaq Half Saturation Constant
  3. Baq Half Saturation Constant
  4. Caq Monod Inhibition Constant
  5. Yield
  6. **Biomass Decay Rate Constant**
  7. Hill Exponent

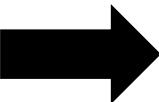
2. Data standardization



3. AI model – ML or DL + Hyperparameter tuning  
**(on Tahoma)**



4. PFLOTRAN reaction network parameters



- Scenarios
  - Scenario-1: Full data (Aaq) + Multi-task learning
  - Scenario-2: GLUE data (Aaq) + Multi-task learning
  - Scenario-3: GLUE data (Aaq) + Single-task learning
  - Scenario-4: GLUE data (Aaq) + Single-task learning + reduced time-series steps
  - Scenario-5: GLUE data (Aaq) + Single-task learning + reduced time-series steps
  - **Scenario-6: GLUE data (Aaq, Baq, Caq, Daq, Xim) + Single-task learning + reduced time-series steps**

GLUE → Generalized Likelihood Uncertainty Estimation

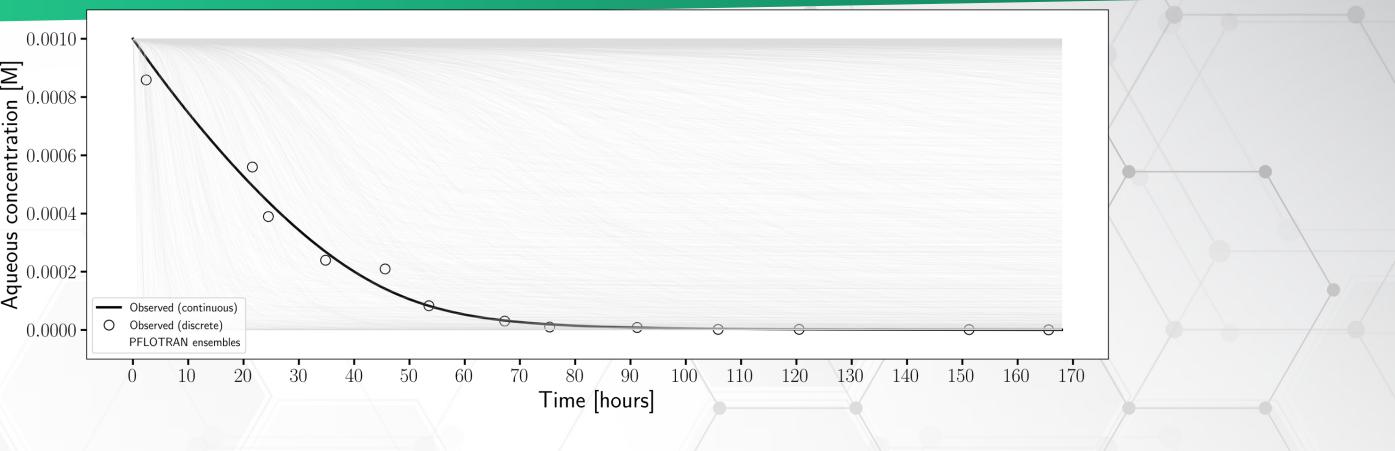
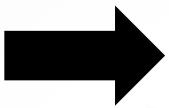
- A methodology to distill full data into actionable information for ML model training
- **A way to perform data worth analysis**

STL or Single-task learning → Train one ML model to predict each PFLOTRAN reaction network parameter

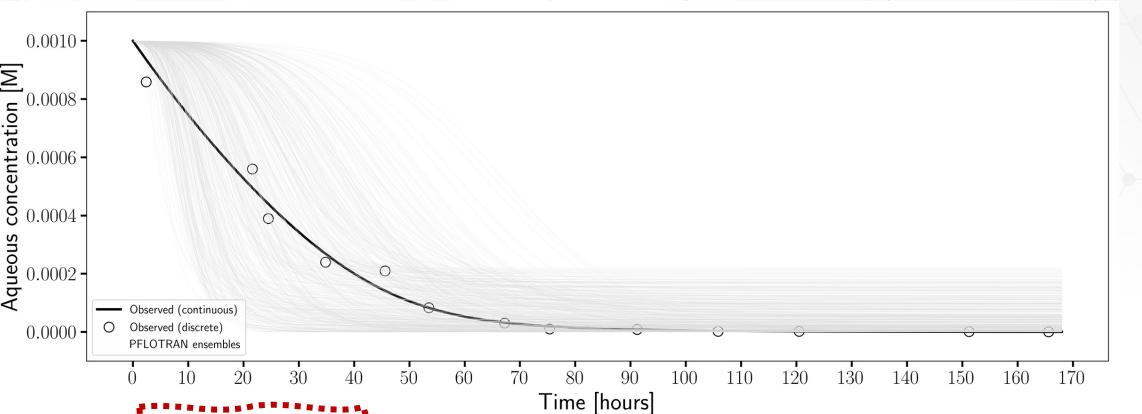
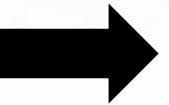
MTL or Multi-task learning → Train a single ML model that predicts all the PFLOTRAN reaction network parameters

## Step-4 – Data worth analysis for ML model training

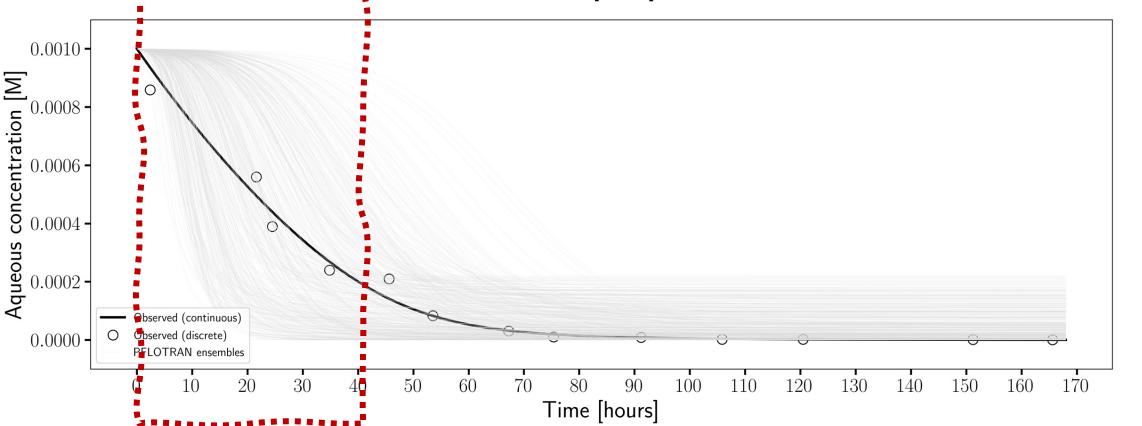
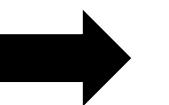
Full data – 15539 realz with 169 time-steps (7 days)



GLUE-based data – 489 realz with 169 time-steps (7 days)



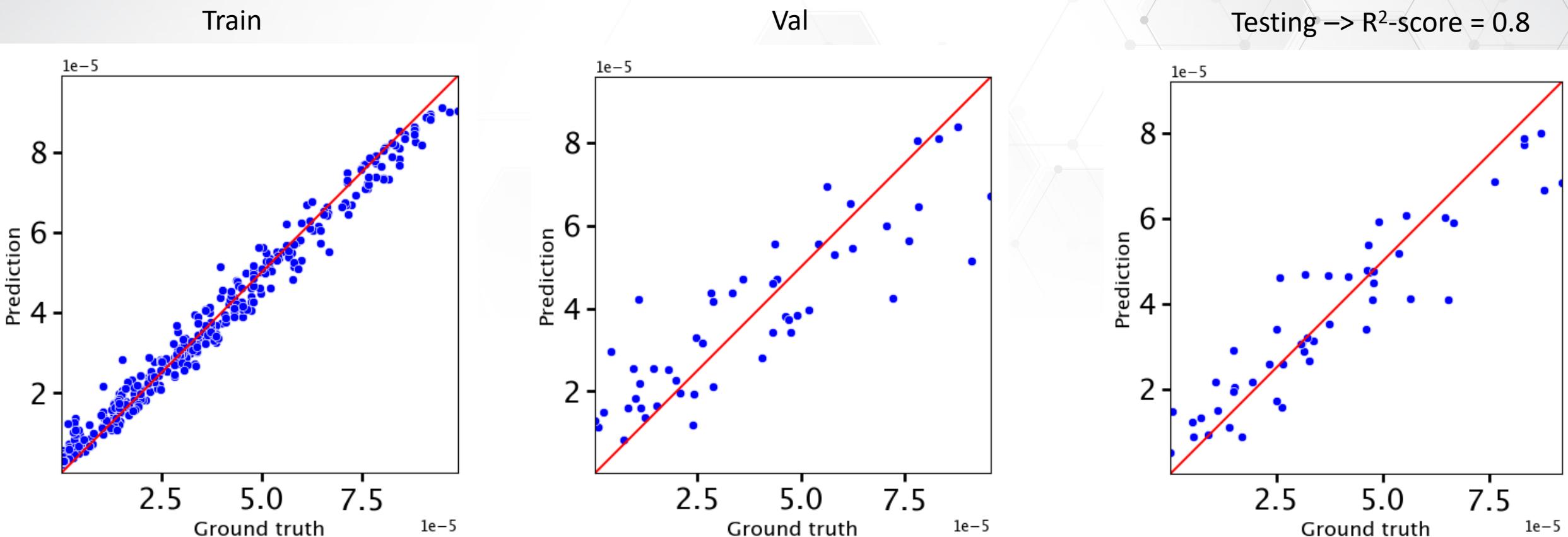
GLUE-based data with reduced time-steps – 489 realz with 40 time-steps (1.5 days)



## Step-5 – Preliminary results and main inference

Scenario-6: ML + Single-task learning + reduced time-series steps (Inputs = Aaq + Baq + Caq + Daq + Xim)

- Biomass decay rate prediction

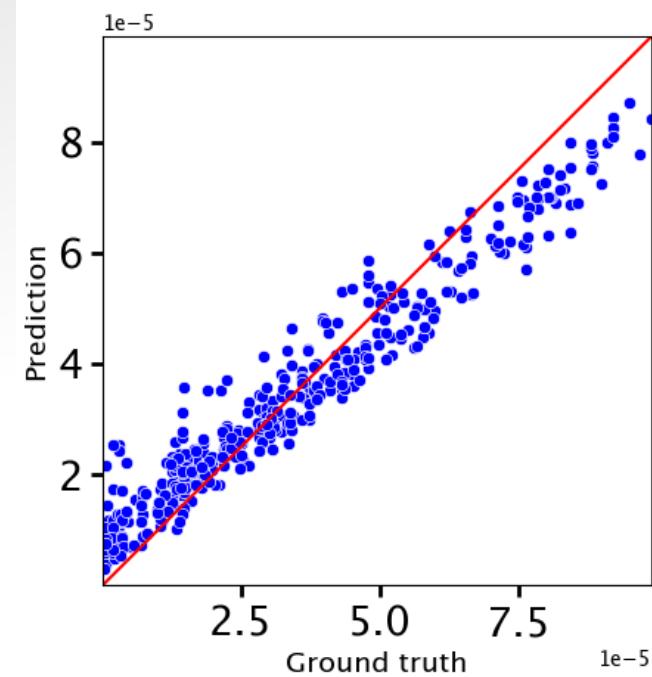


## Step-5 – Results with only Aaq data

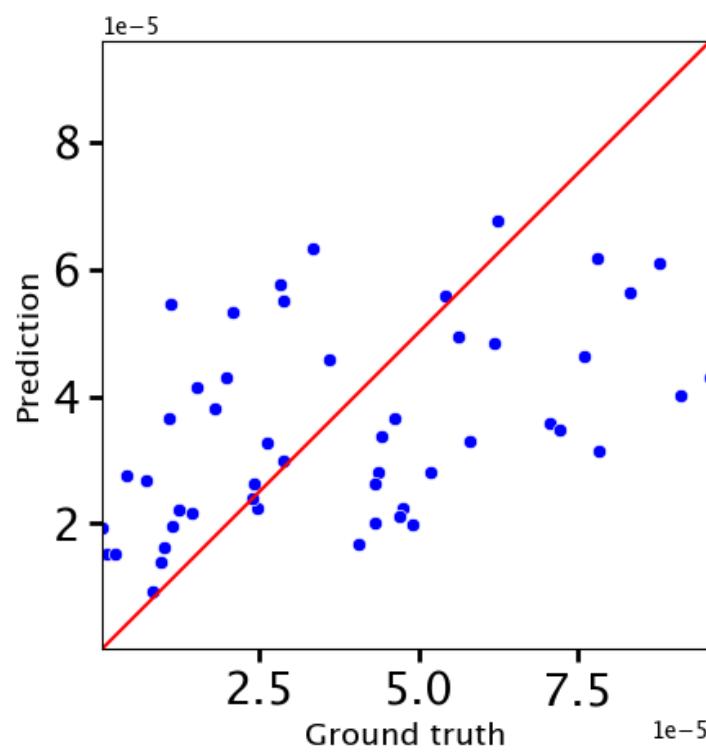
Scenario-6: ML + Single-task learning + reduced time-series steps (Inputs = Aaq only)

- Biomass decay rate prediction

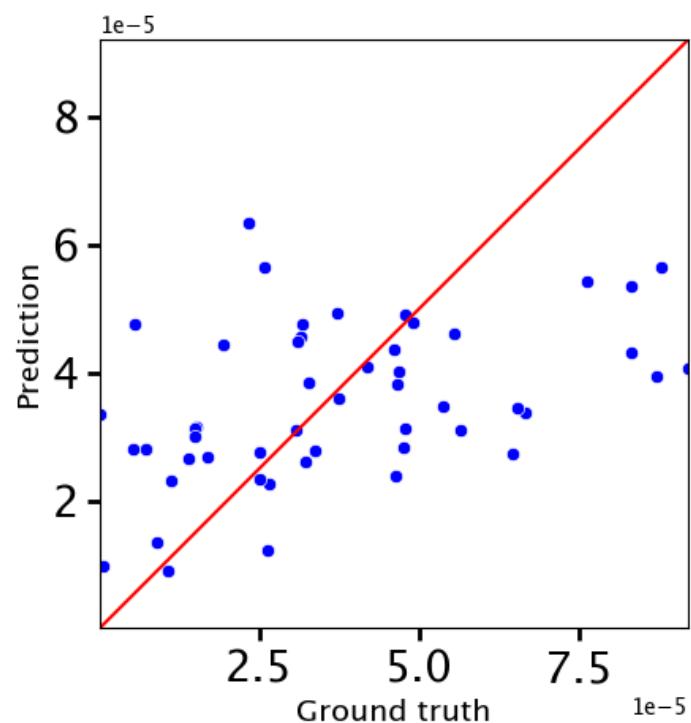
Train



Val



Testing →  $R^2$ -score = 0.6



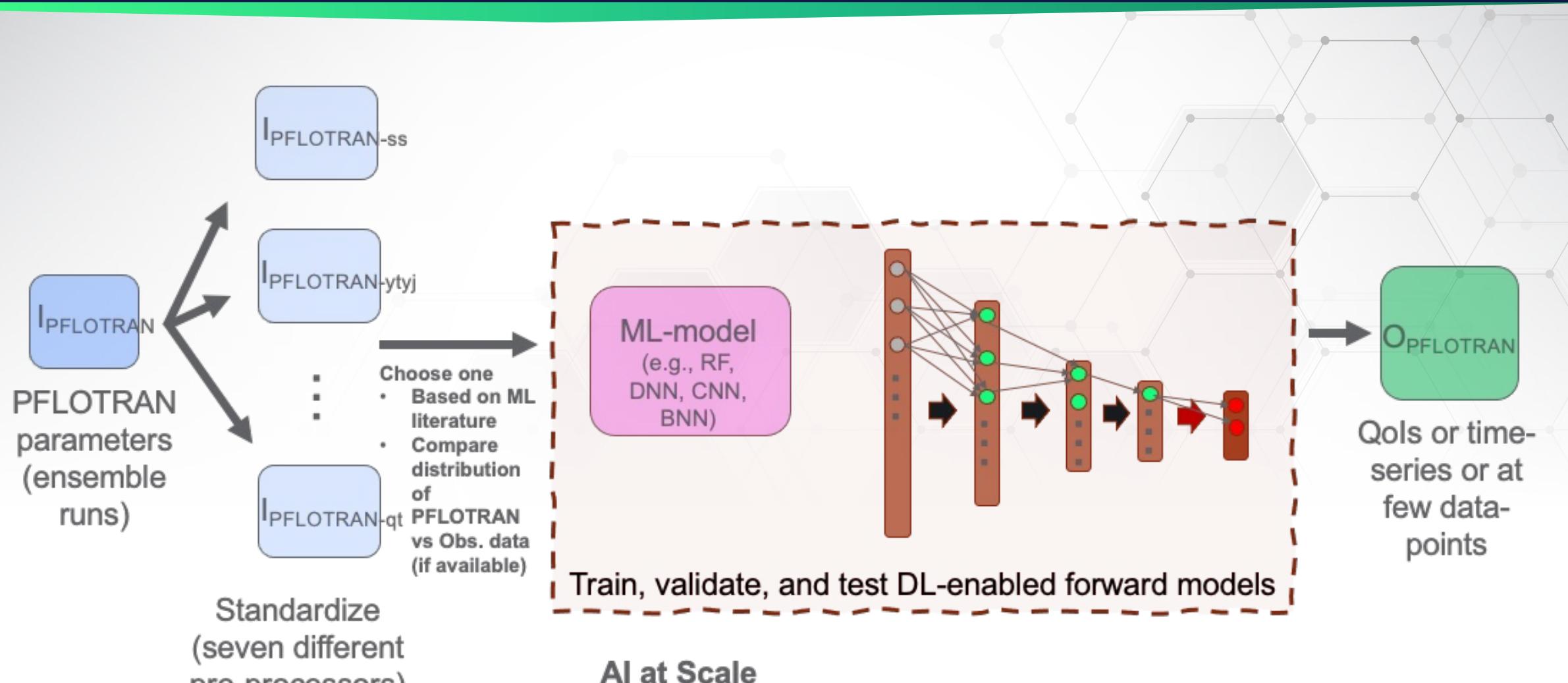
- DNNs – 21875 models trained for each scenario
  - Number of layers -- [1, 2, 3]
  - Neurons – [100, 50, 25]
  - Dropout values for reducing overfitting – [0.0, 0.1, 0.2, 0.3, 0.4]
  - Activation function (LeakyReLU) values -- [0.0, 0.1, 0.2, 0.3, 0.4]
  - Learning rate -- [1e-6, 1e-5, 1e-4, 1e-3, 1e-2]
  - Epochs -- [500, 1000, 1500, 2000, 2500]
  - Batch size -- [4, 8, 16, 32, 64]
- CNNs – 18750 models for each scenario
  - Number of layers -- [1]
  - Filter sizes – [256, 128, 64, 32, 16]
  - Kernel sizes -- [32, 16, 8, 4, 2]
  - Dropout values for reducing overfitting – [0.0, 0.1, 0.2, 0.3, 0.4]
  - Learning rate -- [1e-6, 1e-5, 1e-4, 1e-3, 1e-2]
  - Epochs -- [50, 100, 200, 300, 400, 500]
  - Batch size -- [4, 8, 16, 32, 64]

- Standard scaler
- Max-absolute scaler
- Min-Max Scaler
- Robust Scaler
- Quantile Transformer
  - Normal distribution
  - Uniform distribution
- Power Transformer

- Why need them?

- Each pre-processor has its own advantage
- Pre-processor has considerable effect on ML-model training/predictions/inference
- We need pre-processor that do well on data with outliers
  - (due to process model structural errors)

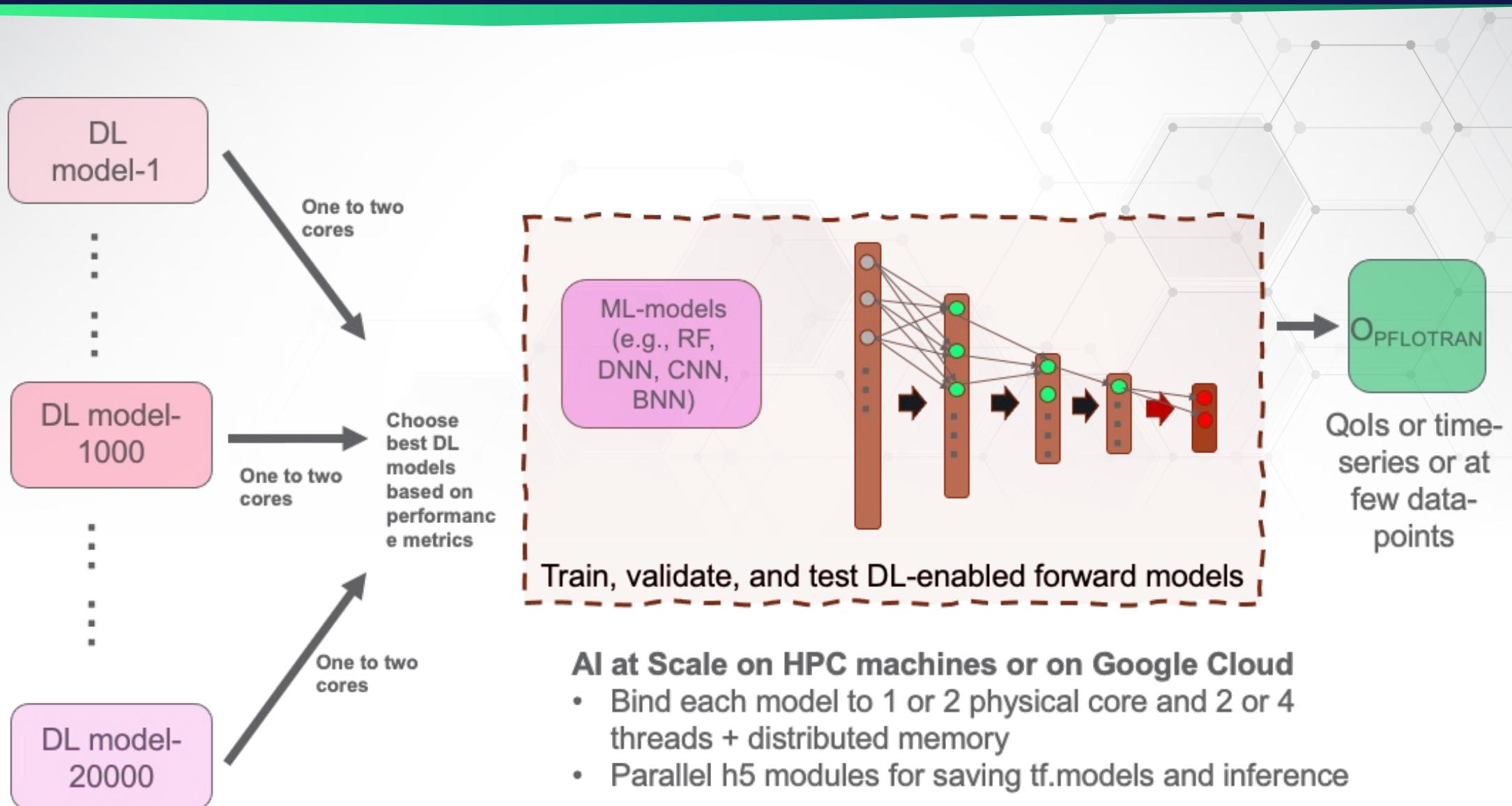
# In-progress: Train, validate, and test ML models at scale (1/2)



## AI at Scale

- Model training is time-consuming and expensive
- Model training is embarrassingly parallel
- Hyperparameter tuning using grid search

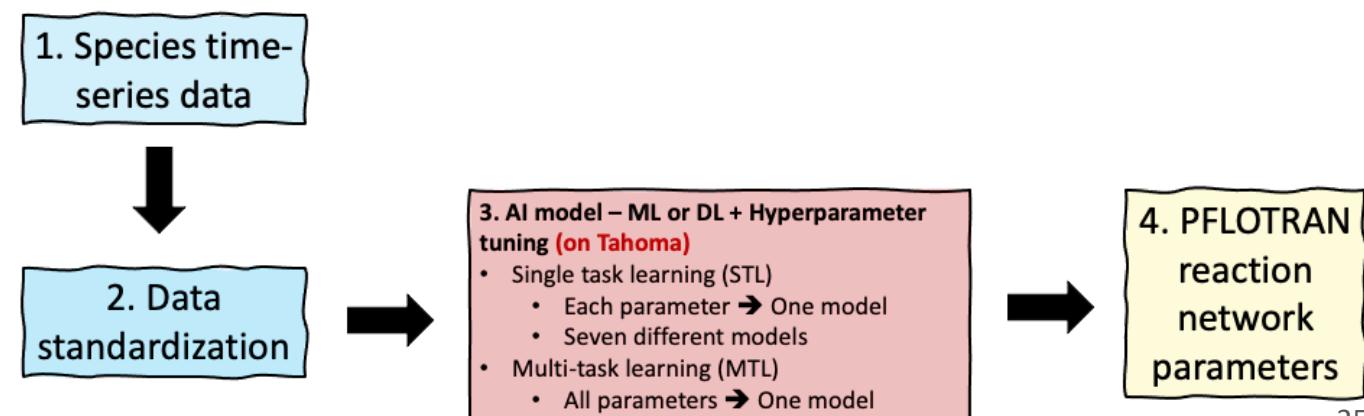
## In-progress: Train, validate, and test ML models at scale (2/2)



# Conclusions, outcomes, and Impact

What have we achieved? How the results impact PFLOTRAN users?

- Before this project there was no AI-enabled ModEx workflow for PFLOTRAN users.
- AI shows promise in parameter estimation and inverse modeling based on performance metrics (good R<sup>2</sup>-scores)
- Informative data (e.g., GLUE, adding other species) can easily improve AI model performance
- PFLOTRAN users can use this workflow or components of it in their research (e.g., data analysis, sensitivity analysis)



## Our next steps: Improving the AI/ML workflows

### Key takeaways:

- Preliminary results show promise in AI/ML models but extensive hyperparameter tuning is necessary
- AI-enabled ModEx: A scalable AI workflow to efficiently calibrate PFLOTRAN process models
  - Manuscript in development with team members and collaborators – Glenn Hammond, Cheng Shi, Satish Karra, and Emily Graham

Once the workflow is fully tested and validated, Python scripts will be made available and open-sourced at: [https://github.com/maruti-iitm/AI\\_ModEx\\_PFLOTRAN.git](https://github.com/maruti-iitm/AI_ModEx_PFLOTRAN.git)

"This research was performed on a project award from the Environmental Molecular Sciences Laboratory, a DOE Office of Science User Facility sponsored by the Biological and Environmental Research program under Contract No. DE-AC05-76RL01830."

Questions?

