# DATA-DRIVEN MODELING IN THE WATER SECTOR: A PARADIGM SHIFT IN PREDICTIVE METHODS

**Mostafa Khalil, PhD**

*Data Scientist | Innovation Engineer, Stantec*

*Views are my own*

# MODELING IN THE WATER SECTOR

**First-Principles Models**

Initial State
(t = 0)

Disturbances
(e.g., influent
flow)

Parameter
Values (e.g.,
reaction
kinetics)

Physical Laws
(e.g., mass balance)

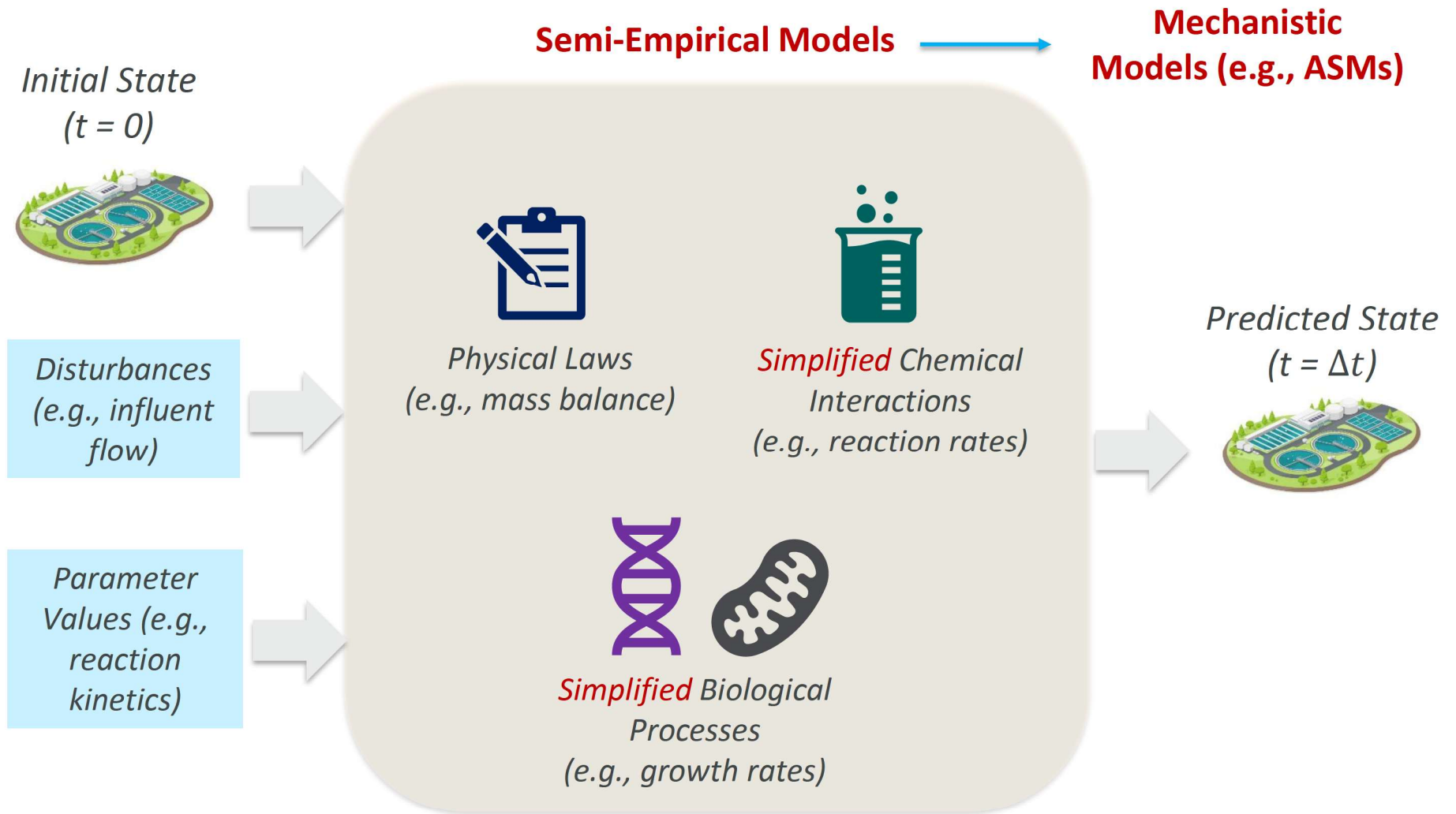*All* Chemical Interactions
(e.g., reaction rates)

*All* Biological Processes
(e.g., growth rates)

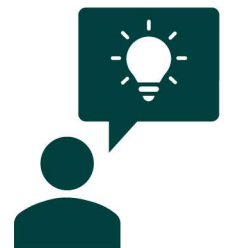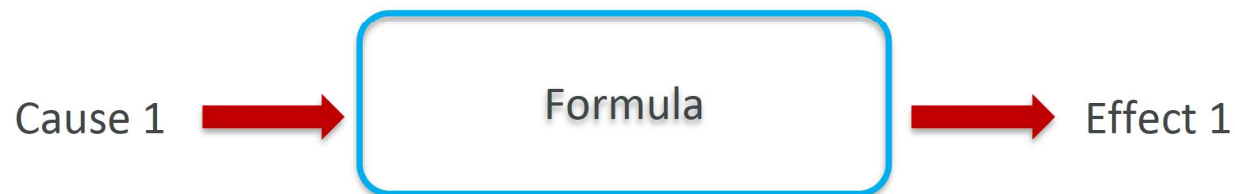Predicted State
(t = $\Delta t$)

# MODELING IN THE WATER SECTOR

**Semi-Empirical Models** → **Mechanistic Models (e.g., ASMs)**

*Initial State (t = 0)*

*Disturbances (e.g., influent flow)*

*Parameter Values (e.g., reaction kinetics)*

*Physical Laws (e.g., mass balance)*

*Simplified Chemical Interactions (e.g., reaction rates)*

*Simplified Biological Processes (e.g., growth rates)*

*Predicted State (t = Δt)*

# MECHANISTIC MODELS

- Gold standards in our field

- Model cause – effect relationships

- Can answer "what-if" questions

- Works outside historical range

$$\text{biomass growth} = \mu_h X_{bh} \left(\frac{S_O}{K_O + S_O}\right)\left(\frac{S_S}{K_S + S_S}\right)\left(\frac{S_{NH}}{K_{NH} + S_{NH}}\right)\left(\frac{S_{ALK}}{K_{ALK} + S_{ALK}}\right)$$

Max. Specific Growth Rate → $\mu_h$

Half-Saturation Coefficient → $S_S$

Biomass — Oxygen — Substrate — Nutrient — Alkalinity

Cause 1 ⟶ **Formula** ⟶ Effect 1

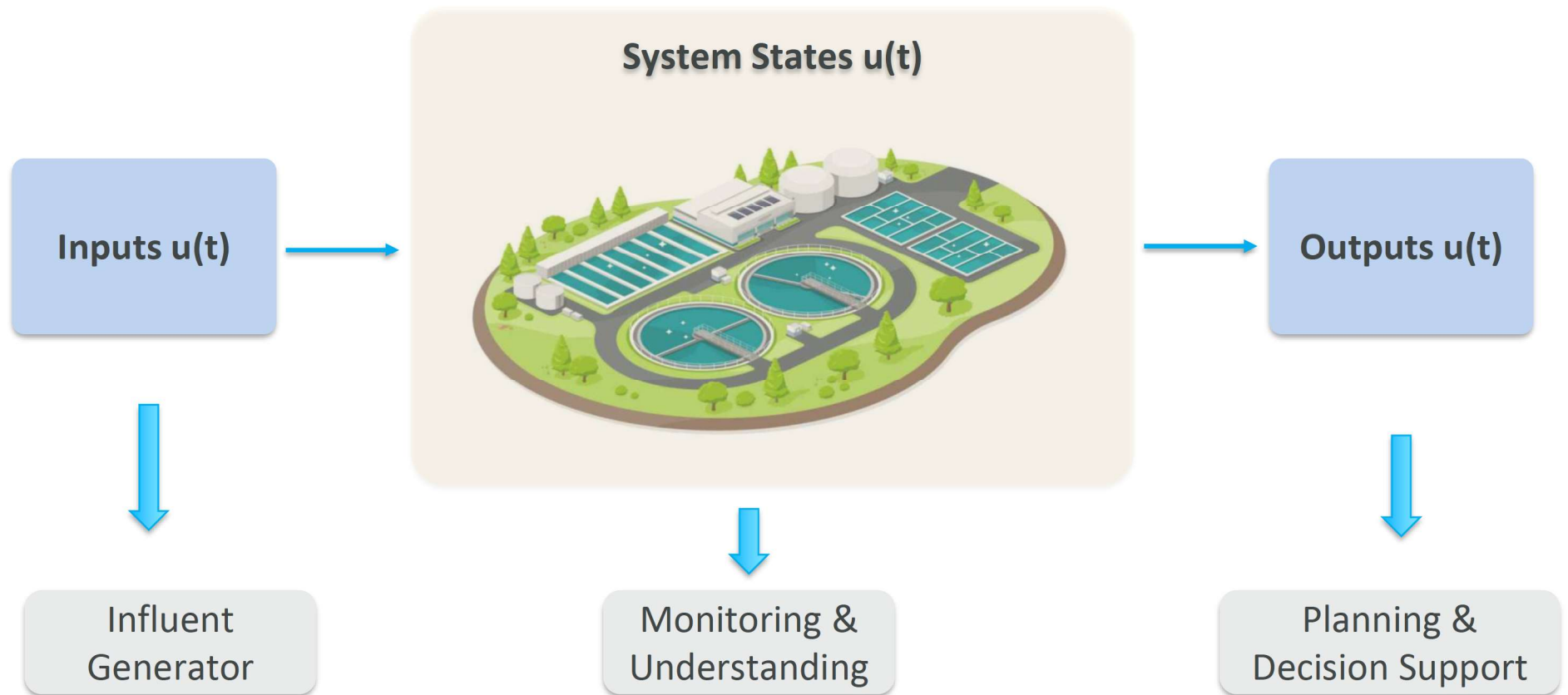# MECHANISTIC MODELS

**Conditional trust and confidence:**

- No Unknown or poorly understood relationships

- Parameters are accurate and (somewhat) fixed over time

- Same formulas will hold at all conditions
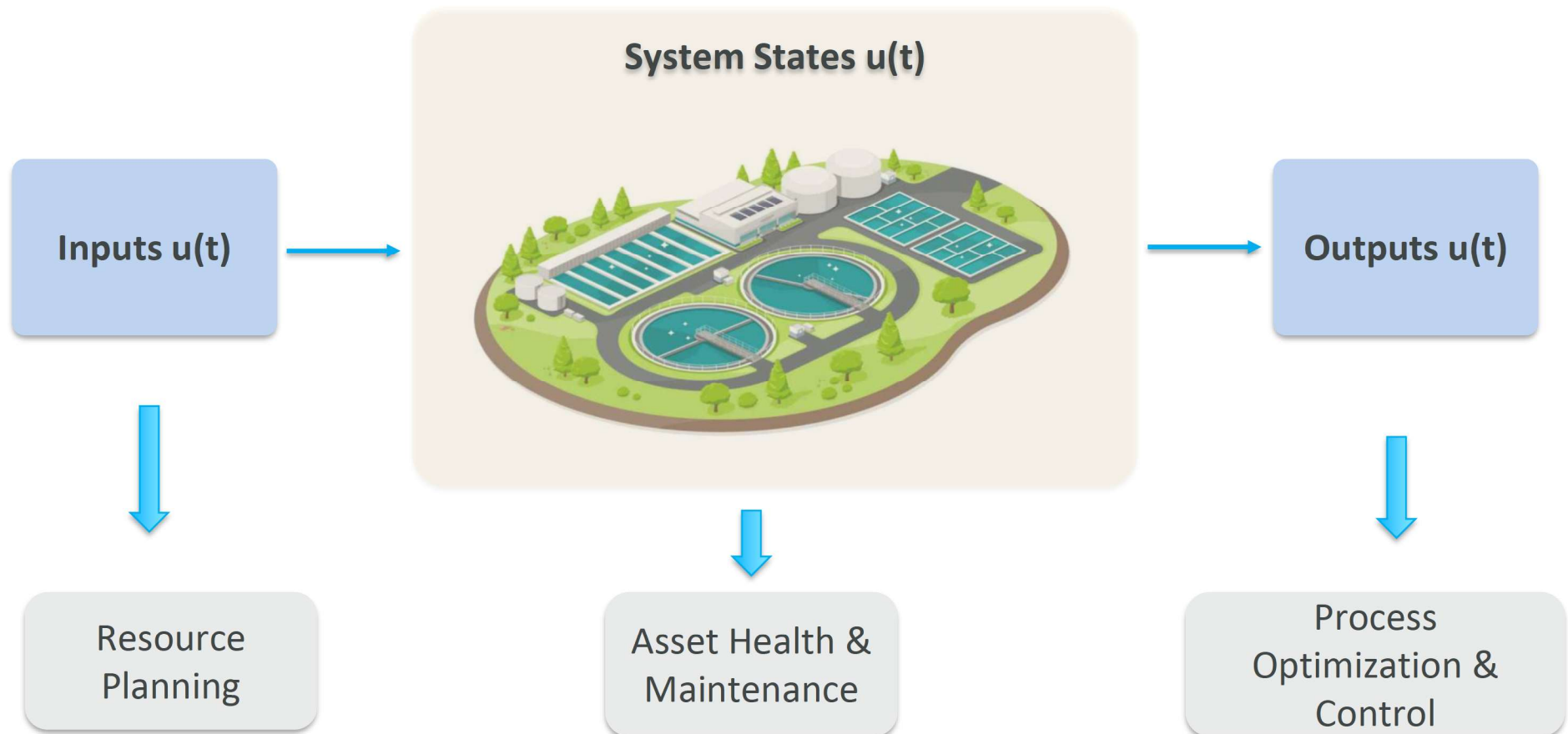
# WHY DO WE BUILD MODELS?

**System States u(t)**



**Inputs u(t)** → **Outputs u(t)**

Influent Generator

Monitoring & Understanding

Planning & Decision Support

# WHY DO WE BUILD MODELS?

**System States u(t)**



**Inputs u(t)** → → **Outputs u(t)**

Resource Planning

Asset Health & Maintenance

Process Optimization & Control

# WHY DO WE BUILD MODELS?

Process Monitoring & Understanding

Process Optimization & Control

Planning and Decision Support

Asset Health & Maintenance

Forecasting

# UTILIZING THE POWER OF DATA

Artificial Intelligence

Machine Learning

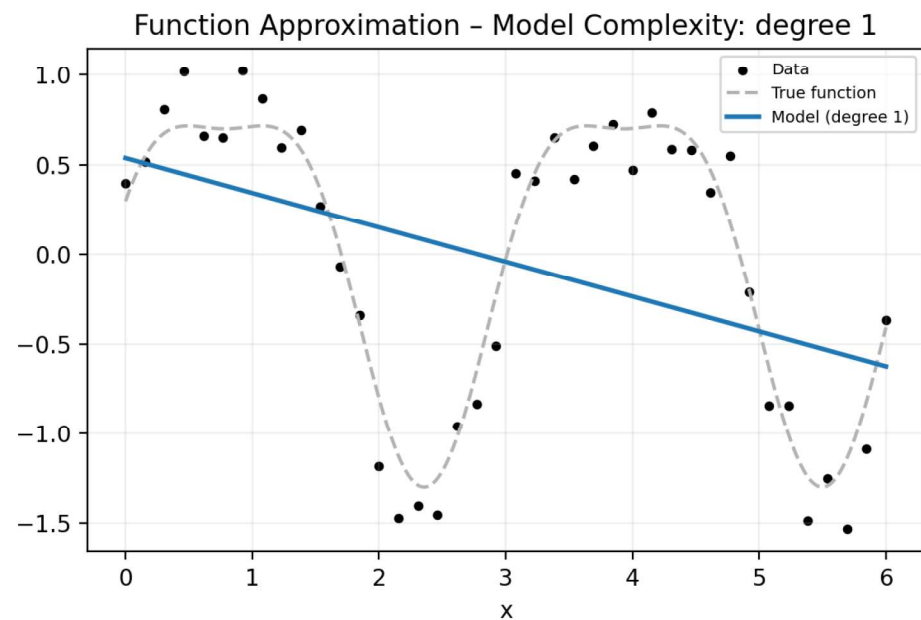| Natural Language Processing | Computer Vision | Robotics | Speech Recognition |

# UTILIZING THE POWER OF DATA: MACHINE LEARNING

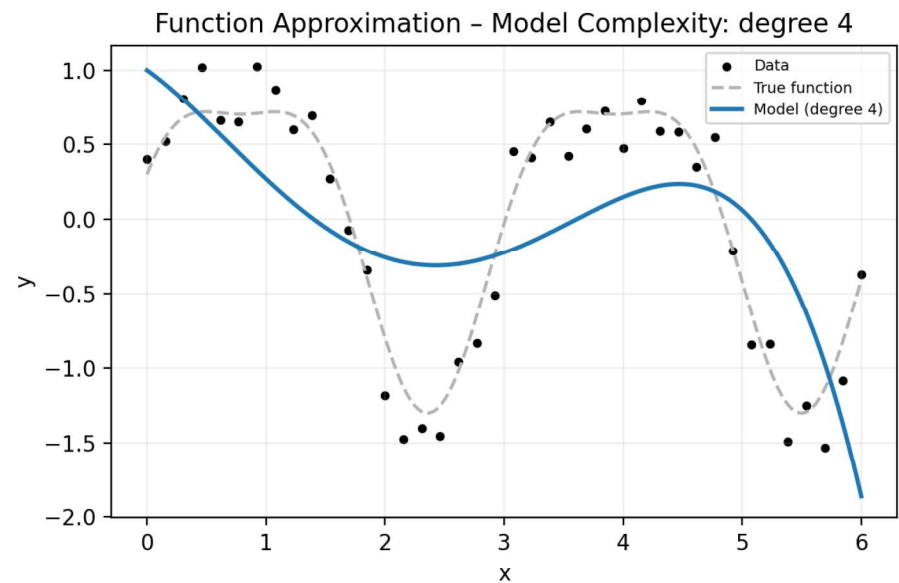Input (Predictors)        Output
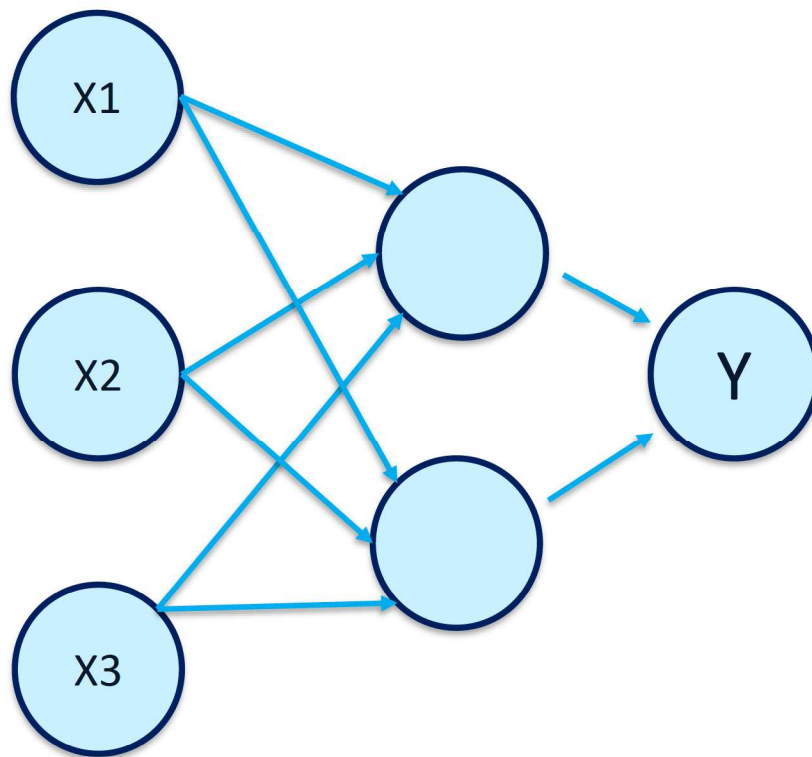
| X1 | X2 | X3 | Xn | Y |
|----|----|----|----|---|
|    |    |    |    |   |
|    |    |    |    |   |
|    |    |    |    |   |
|    |    |    |    |   |
|    |    |    |    |   |

Observations

X1

X2 → Y

X3

Data Points
Linear Regression (y = 2.12x + 4.55)

Y (Target)

X (Feature)

# UTILIZING THE POWER OF DATA



Function Approximation – Model Complexity: degree 1

# UTILIZING THE POWER OF DATA

X1

X2

X3

Y

Output

Input

Function Approximation – Model Complexity: degree 4

# UTILIZING THE POWER OF DATA



Function Approximation – Model Complexity: degree 8

Input

# UTILIZING THE POWER OF DATA
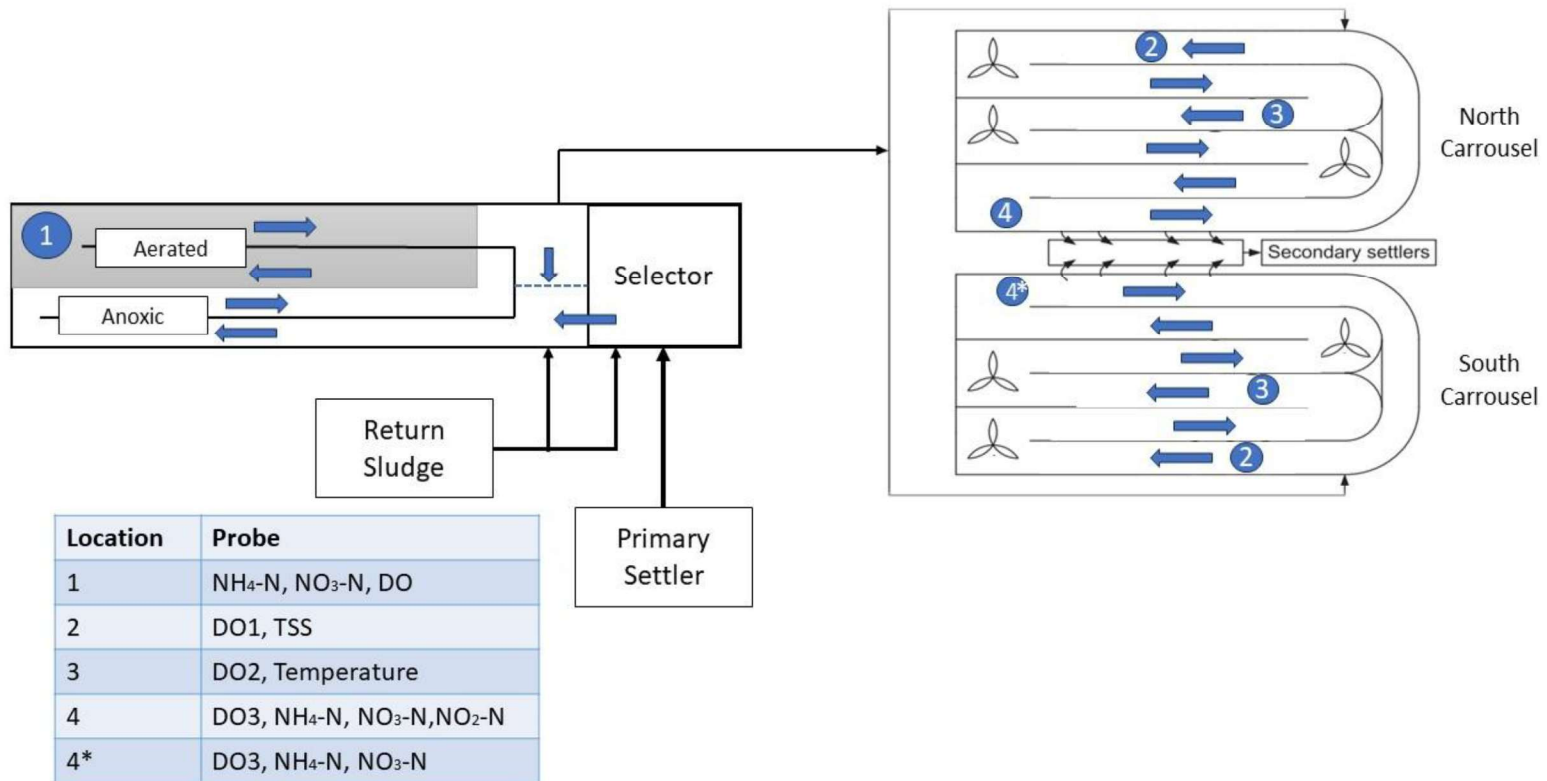


Input



Function Approximation – Model Complexity: degree 16

# INTRODUCTION TO CHALLENGES

# EXAMPLE: N₂O EMISSIONS MODELING



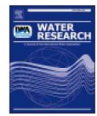| Location | Probe |
|----------|-------|
| 1 | NH₄-N, NO₃-N, DO |
| 2 | DO1, TSS |
| 3 | DO2, Temperature |
| 4 | DO3, NH₄-N, NO₃-N,NO₂-N |
| 4* | DO3, NH₄-N, NO₃-N |

# EXAMPLE: N₂O EMISSIONS MODELING

# ALIGNMENT WITH DOMAIN KNOWLEDGE



Model Complexity (number of trees)

800 Before optimization

231 After optimization

# ALIGNMENT WITH DOMAIN KNOWLEDGE

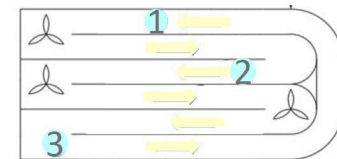DO Measurement Locations

What features are the most important
for the model to make prediction?

Remove NO₂ from input
features



Permutation feature importance

Permutation feature importance
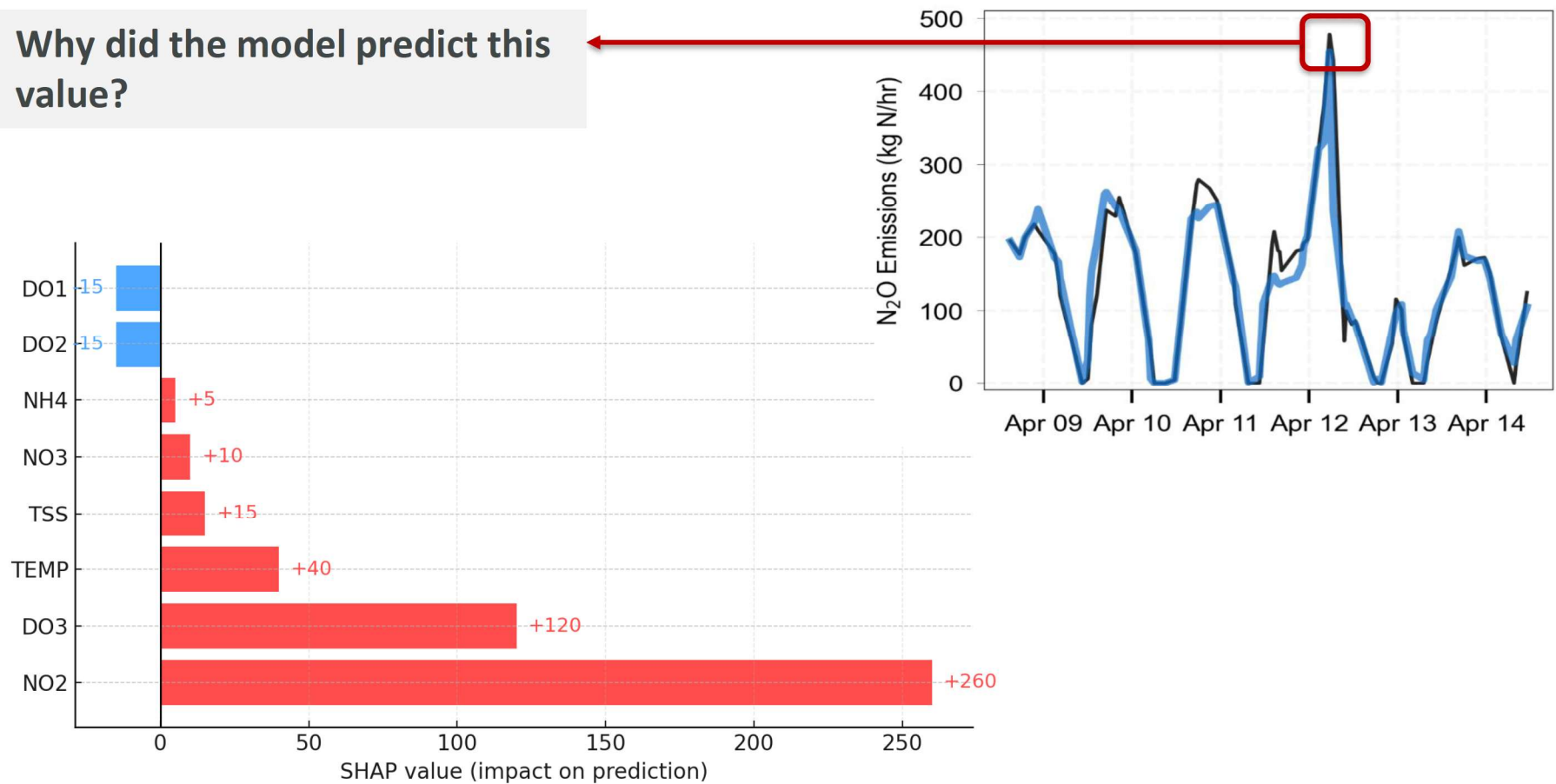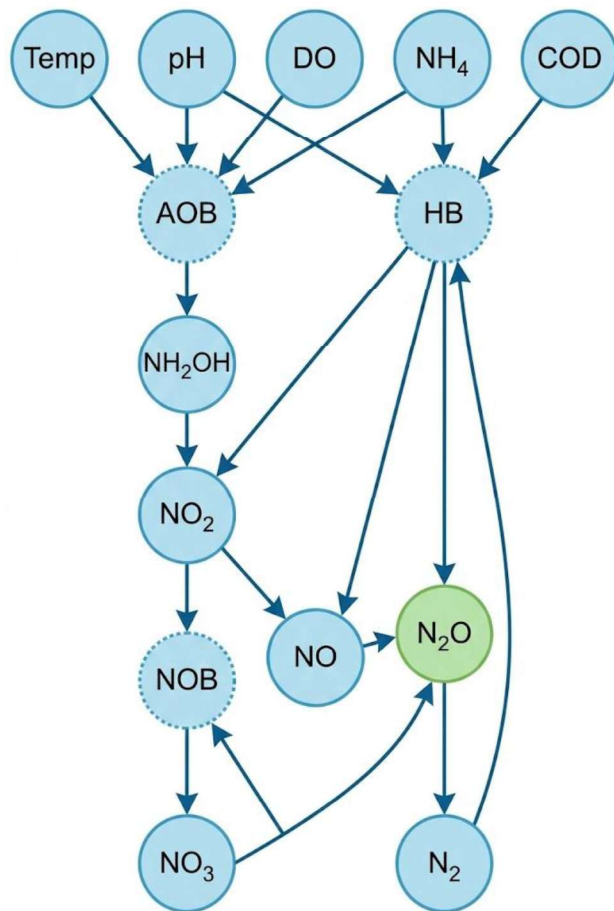
# INTERPRETABILITY OF MODEL PREDICTIONS
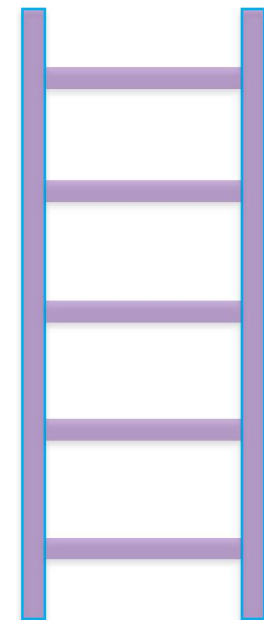
Why did the model predict this value?



**This is not a causal effect!**

# ADVANCING TOWARDS DECISION SUPPORT



Interpretability ≠ Causality

○ Counterfactuals

○ Interventions

✔ Local Interpretability

✔ Global Interpretability

# THE CAUSALITY PROBLEM

- Correlation can be misleading when an unmeasured factor (confounder) influences both variables

  - P ($N_2O$ | DO): What we observe in the data (correlation)

  - P ($N_2O$ | do(DO) = d): What would happen if we *intervene* on DO (causal effect)



ML models can easily learn correlations
but extracting *causation* requires extra work

# EXAMPLE: RO MEMBRANE FOULING FORECASTING
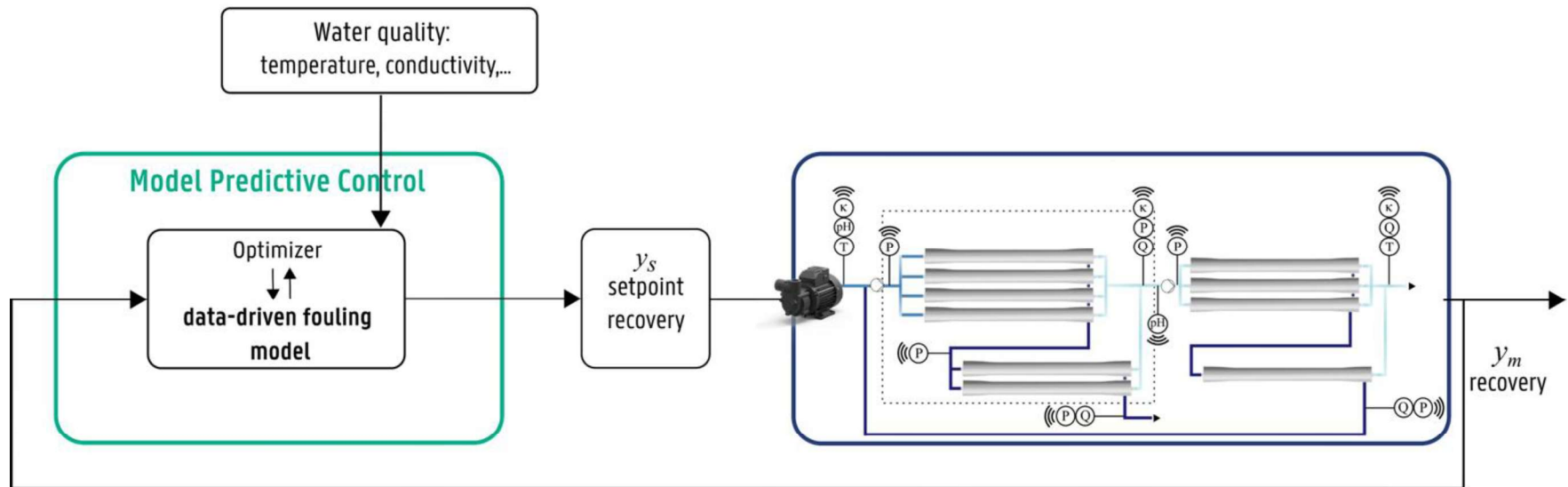
# EXAMPLE: RO MEMBRANE FOULING FORECASTING

Last 6 days

| Conductivity |
| Temperature |
| Recovery |
| Pressure |

→ LSTM →

4 hours into the future

| Fouling |

Rm forecasting on valid range



- valid Rm data
- LSTM forecast Rm

2023-04-28 13:15:00    2023-08-04 17:15:00

# EXAMPLE: RO MEMBRANE FOULING FORECASTING



Rm forecasting on zoomed range

Lag in forecasted fouling (Rm) was caused by frequent OFF periods in the installation

# EXAMPLE: RO MEMBRANE FOULING FORECASTING
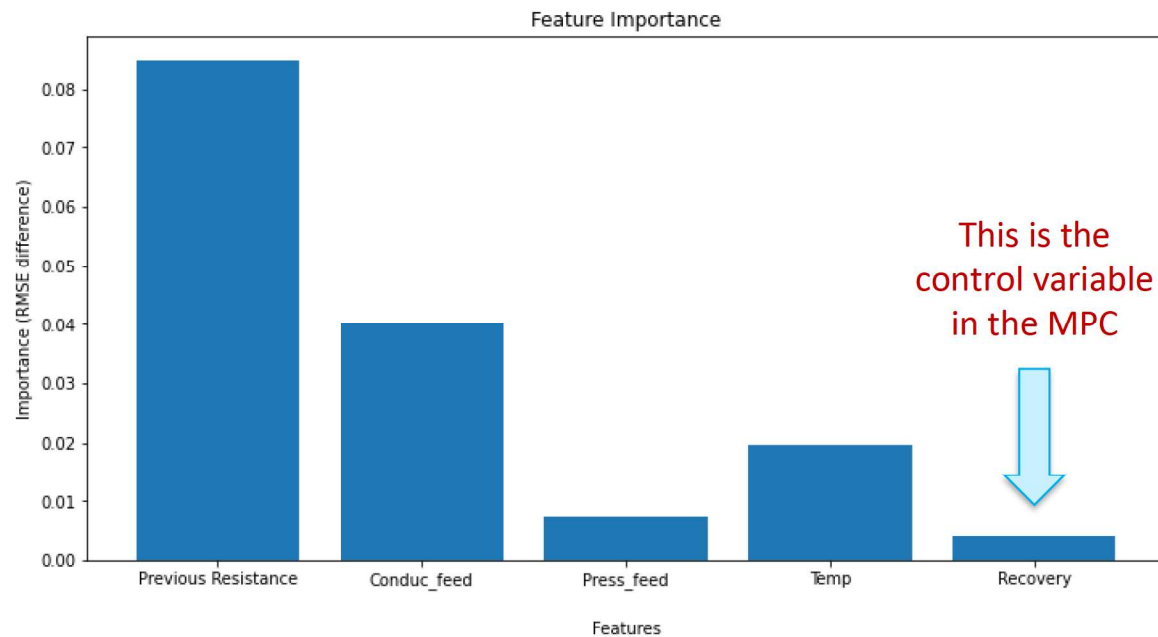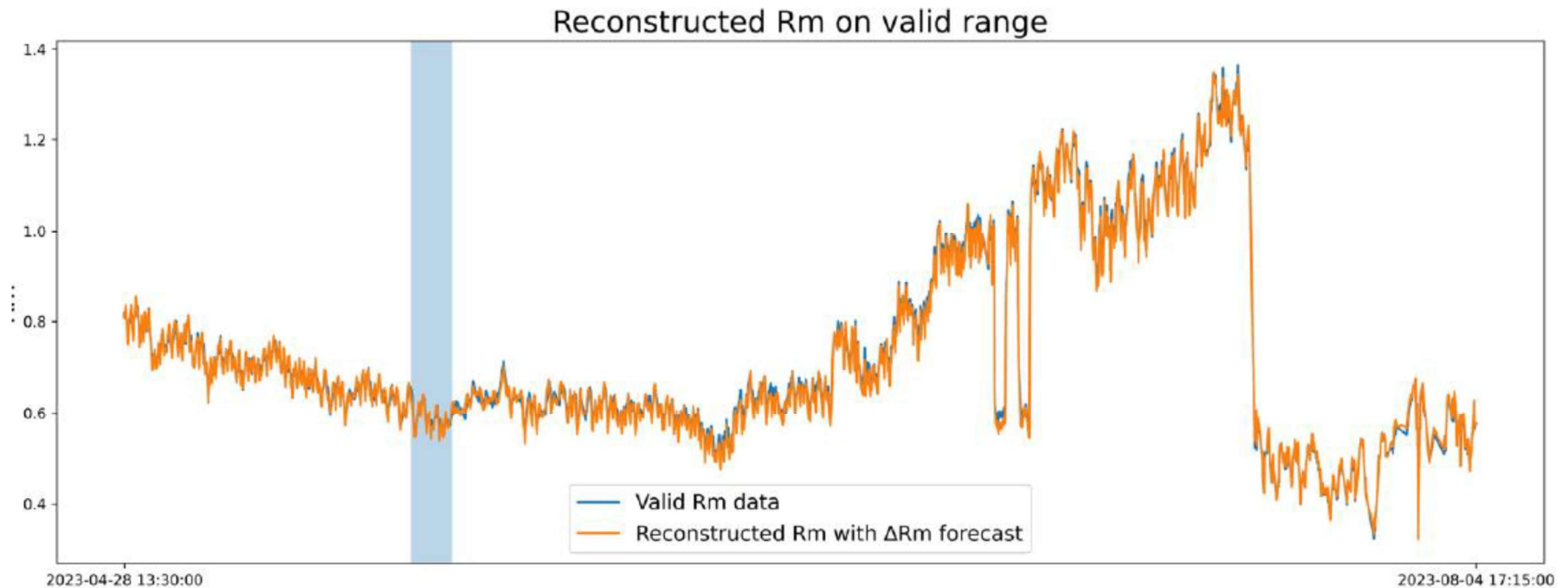
# EXAMPLE: RO MEMBRANE FOULING FORECASTING



Reconstructed Rm on valid range

# EXAMPLE: RO MEMBRANE FOULING FORECASTING



Predicting the *change* in fouling makes the model learn system dynamics directly, eliminating lag and improving responsiveness.

# EXAMPLE: RO MEMBRANE FOULING FORECASTING

- Even for a powerful model, forecasting Rm directly made it slower to react to sudden changes (e.g., frequent on – off)

- Predictions were lagged during fast transitions

- $\Delta Rm$ prediction is simpler as it removed slow trends and noise (stationary signal)

- Final Rm is reconstructed by adding predicted $\Delta Rm,$ leading to better accuracy

# KEY TAKEAWAYS!

# Thank you!

- Data-driven methods expand our modeling toolbox — they don't replace physics or expertise

- ML is powerful but could be fragile: performance depends more on data and context than on algorithms

- Accuracy can be misleading: a good fit does not mean the model is correct

- ML is not magic. It introduces new challenges (drift, retraining, explainability) that must be managed deliberately

- A model can look right and be wrong. Accuracy is not the whole story