# KGML for Aquatic Sciences
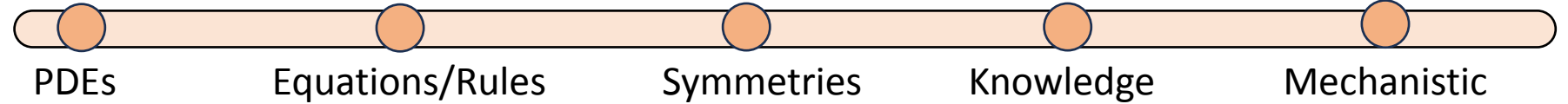
# Organizing KGML Research: A Multi-Dimensional View

**Format Used for Representing Knowledge**

PDEs — Equations/Rules — Symmetries — Knowledge Graphs — Mechanistic Models
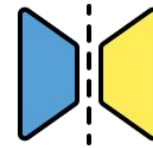
$$\frac{\partial u(x,t)}{\partial t} + \mathcal{N}\big(\lambda, u(x,t)\big) = 0$$

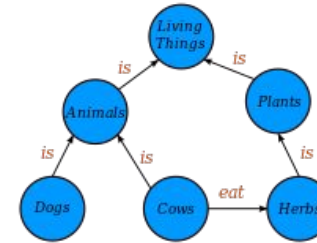Navier Stokes Equation,
Wave Equation,
Schrodinger Equation, …

$$\frac{\partial u}{\partial x} \propto \frac{\partial^2 u}{\partial x^2}$$

$$a < \frac{\partial u}{\partial x} < b$$

**PINNs:** Raissi et al. 2019
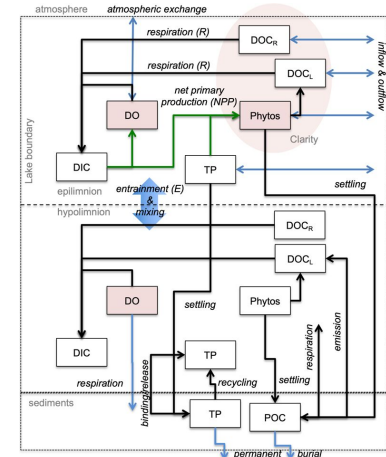
**PGNN:** Karpatne et al. 2017
**PGRNN:** Jia et al. 2019
**PGA-LSTM:** Daw et al. 2020

**NequIP:** Batzner et al. 2022
**Cormorant:** Anderson et al. 2019
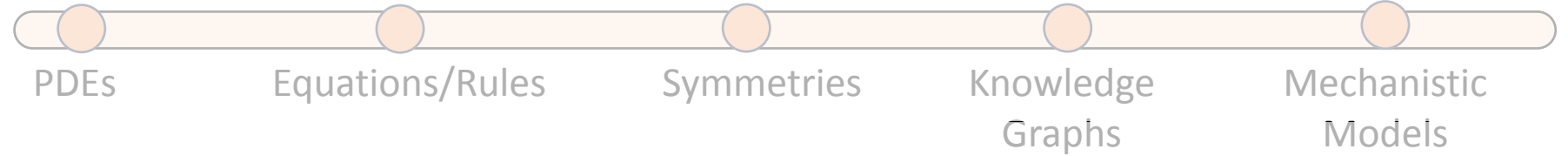**Equivariant-Net:** Wang et al. 2021
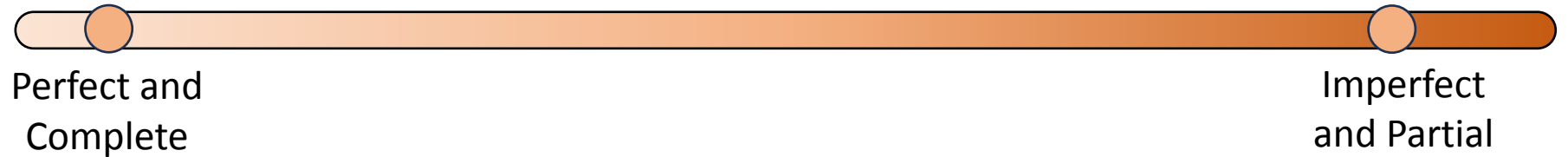
Zareian et al. 2020

**MCL:** Ladwig et al. 2024
**dPL:** Shen et al. 2023

Karpatne, Jia, and Kumar. "Knowledge-guided Machine Learning: Current Trends and Future Prospects." *arXiv:2403.15989* (2024).

# Organizing KGML Research: A Multi-Dimensional View

**Format Used for Representing Knowledge**

PDEs — Equations/Rules — Symmetries — Knowledge Graphs — Mechanistic Models

**Type of Scientific Knowledge**

Perfect and Complete

Imperfect and Partial

Example: Solving *known* PDEs



Initial Vorticity | t=15 | t=20

Navier Stokes Eq., Heat Eq., Wave Eq., Schrodinger Eq., …

Primary Objective: Improve Computational Efficiency

**PINNs:** Raissi et al. 2019, **DeepONets**: Lu et al. 2021, **FNOs**: Li et al. 2021

Example: Modeling complex dynamical systems with missing/imperfect physics

Modeling Turbulence, Multi-phase Flow, Cloud Physics, Aerosols, …



Additional Objective: Improve Modeling Accuracy

**PGNN:** Karpatne et al. 2017, **PGRNN**: Jia et al. 2019, **PGA-LSTM**: Daw et al. 2020

Karpatne, Jia, and Kumar. "Knowledge-guided Machine Learning: Current Trends and Future Prospects." *arXiv:2403.15989* (2024).

# Organizing KGML Research: A Multi-Dimensional View



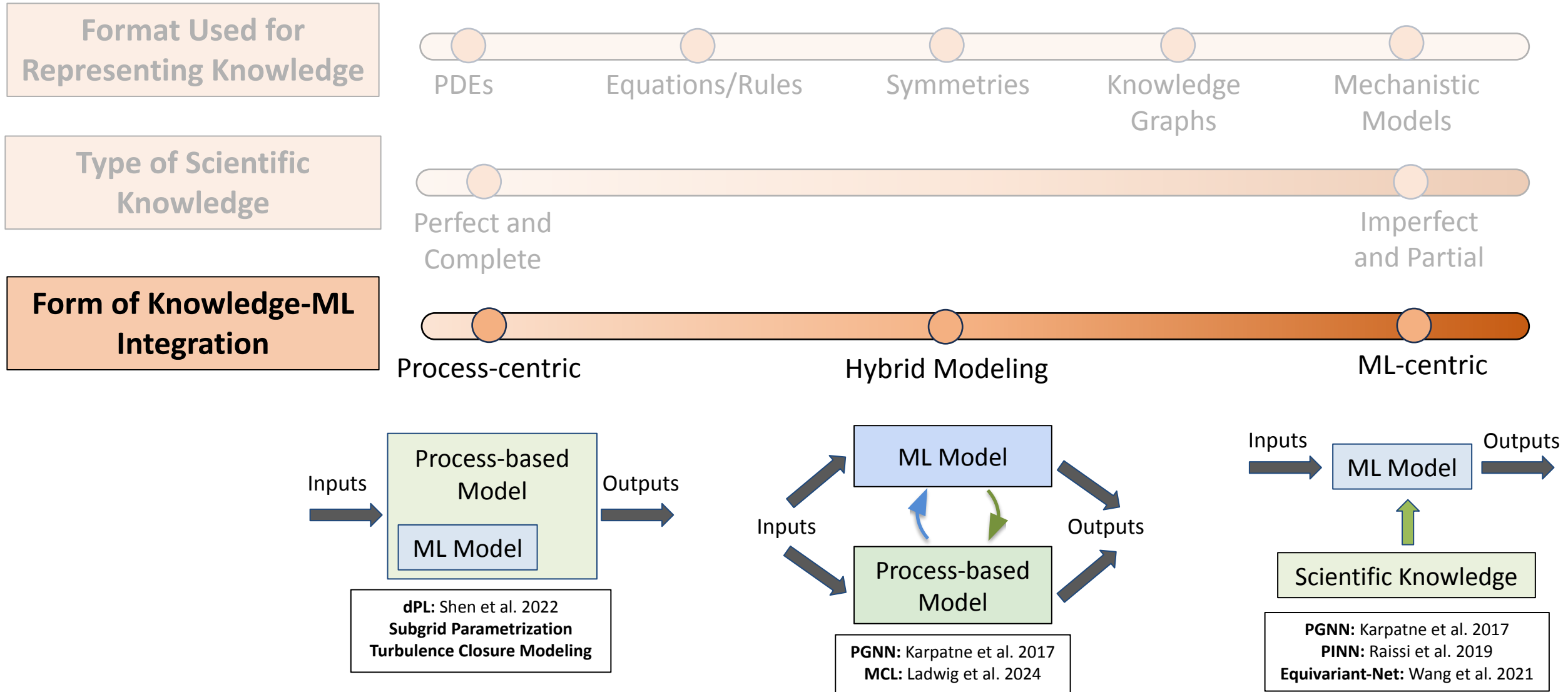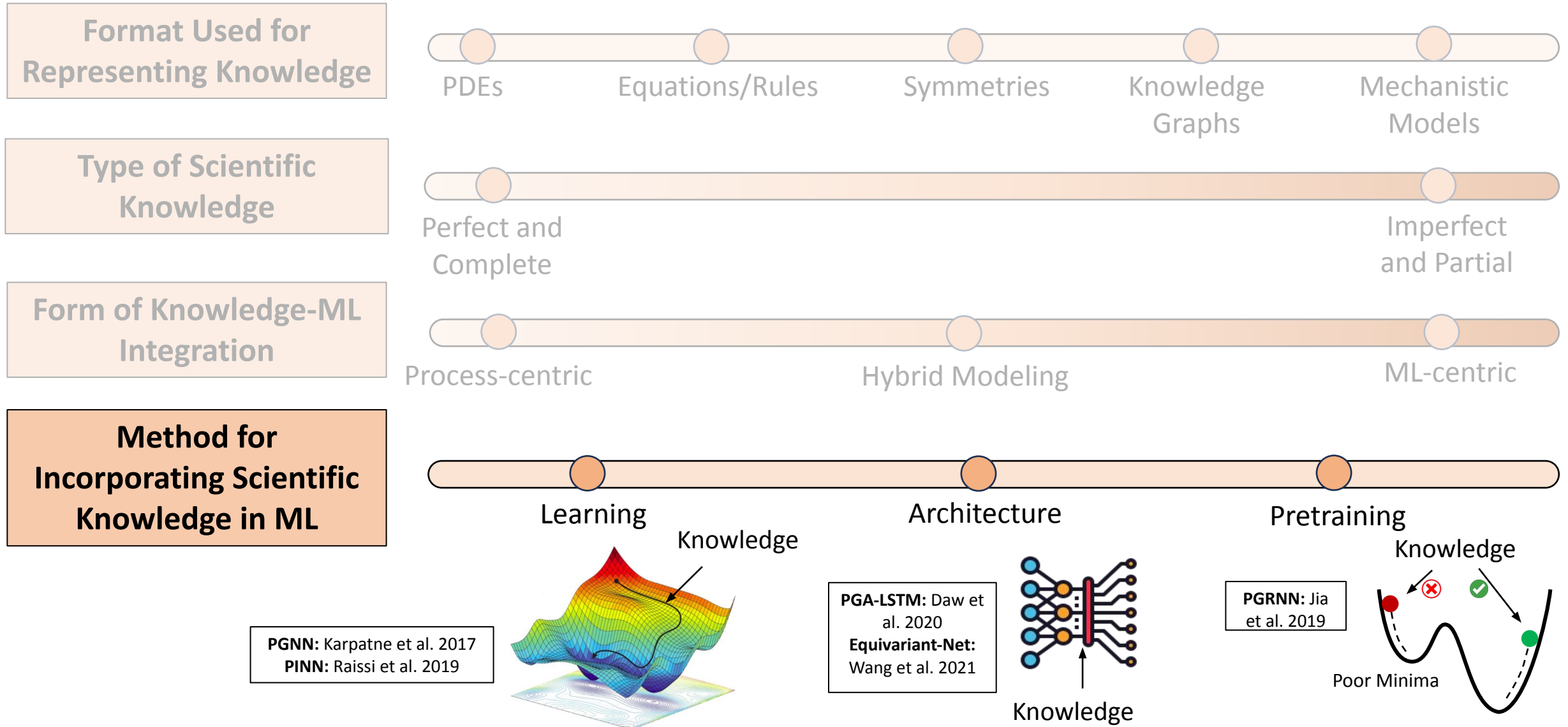**Format Used for Representing Knowledge**

PDEs — Equations/Rules — Symmetries — Knowledge Graphs — Mechanistic Models

**Type of Scientific Knowledge**

Perfect and Complete — Imperfect and Partial

**Form of Knowledge-ML Integration**

Process-centric — Hybrid Modeling — ML-centric

Inputs → Process-based Model [ML Model] → Outputs

**dPL:** Shen et al. 2022
**Subgrid Parametrization**
**Turbulence Closure Modeling**

Inputs → ML Model ⇄ Process-based Model → Outputs

**PGNN:** Karpatne et al. 2017
**MCL:** Ladwig et al. 2024

Inputs → ML Model → Outputs
Scientific Knowledge

**PGNN:** Karpatne et al. 2017
**PINN:** Raissi et al. 2019
**Equivariant-Net:** Wang et al. 2021

Karpatne, Jia, and Kumar. "Knowledge-guided Machine Learning: Current Trends and Future Prospects." *arXiv:2403.15989* (2024).

# Organizing KGML Research: A Multi-Dimensional View



Karpatne, Jia, and Kumar. "Knowledge-guided Machine Learning: Current Trends and Future Prospects." *arXiv:2403.15989* (2024).

# Organizing KGML Research: A Multi-Dimensional View



Karpatne, Jia, and Kumar. "Knowledge-guided Machine Learning: Current Trends and Future Prospects." *arXiv:2403.15989* (2024).

# Organizing KGML Research: A Multi-Dimensional View



Karpatne, Jia, and Kumar. "Knowledge-guided Machine Learning: Current Trends and Future Prospects." *arXiv:2403.15989* (2024).
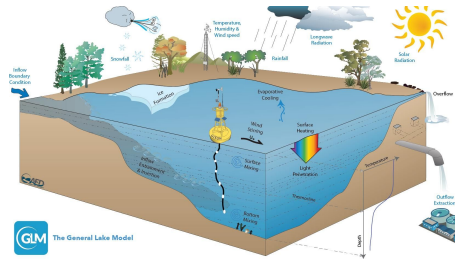
# KGML Use Cases

## Lake Temperature Modeling

**Goal:** Predicting the temperature of the lake.
- Use *imperfect* and *partial* knowledge as loss functions
- Use *simulation data* for pre-training and *observational data* for finetuning
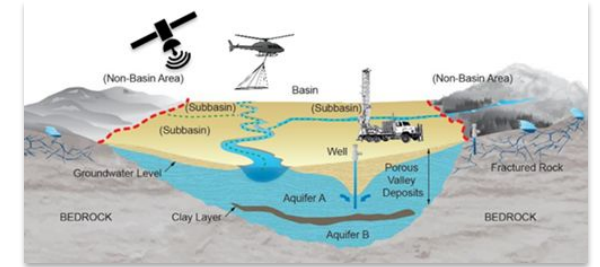
Physics-guided NNs **(PGNNs):** Daw et al. 2017

Physics-guided RNNs **(PGRNNs):** Jia et al. 2019



## River-basin Characterization

**Goal:** Predict basin characteristics of rivers.
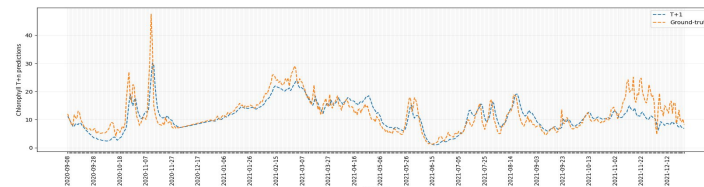- Extract system characteristics from driver and response data.



Knowledge-guided Self-supervised **(KGSSL):** Ghosh et al. 2022

Uncertainty Quantification **(UQ-KGSSL):** Sharma et al. 2022

## Chlorophyll-a Prediction



**Goal:** Predicting the chlorophyll-a content of water bodies.
- Sparse observed data for chlorophyll
- Interested in predicting the blooms.

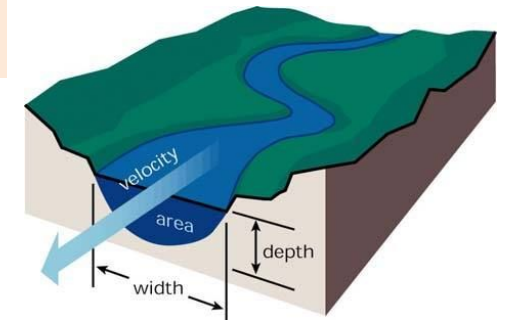LSTM based Chl-a Prediction**:** Cen et al. 2022



## Streamflow Forecasting

**Goal:** Predict the stream flow of rivers.
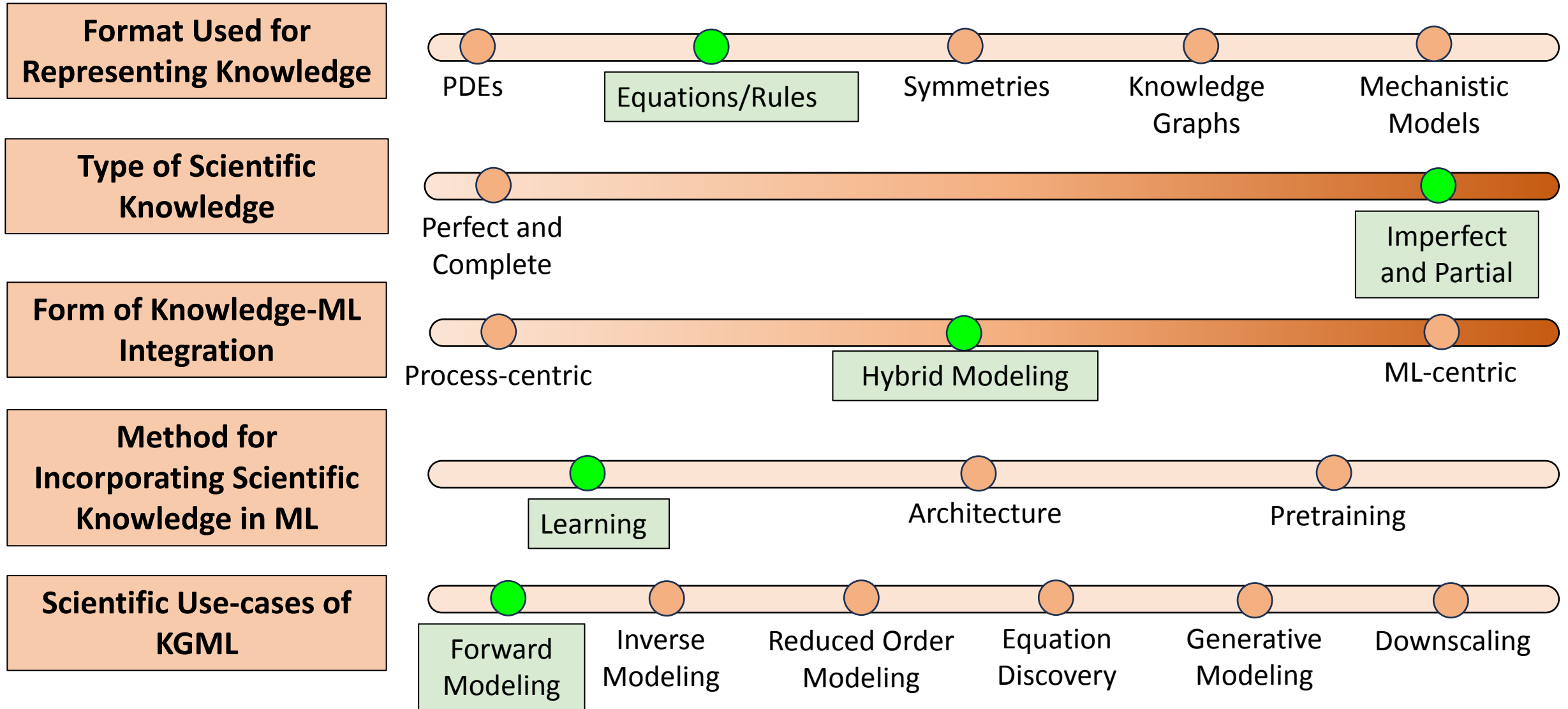- Use river-network data (graph) and the knowledge of thermodynamics to improve predictions.



Physics-guided Recurrent Graph Model **(PGRGnN):** Jia et al. 2020

KGML for Multi-scale Process and Data Assimilation**:** Kumar et al. 2023

# Use Case 1:
# Lake Temperature Modeling

# Organizing KGML Research: A Multi-Dimensional View



**Format Used for Representing Knowledge**

PDEs · Equations/Rules · Symmetries · Knowledge Graphs · Mechanistic Models

**Type of Scientific Knowledge**

Perfect and Complete · Imperfect and Partial

**Form of Knowledge-ML Integration**

Process-centric · Hybrid Modeling · ML-centric

**Method for Incorporating Scientific Knowledge in ML**

Learning · Architecture · Pretraining

**Scientific Use-cases of KGML**

Forward Modeling · Inverse Modeling · Reduced Order Modeling · Equation Discovery · Generative Modeling · Downscaling

Karpatne, Jia, and Kumar. "Knowledge-guided Machine Learning: Current Trends and Future Prospects." *arXiv:2403.15989* (2024).
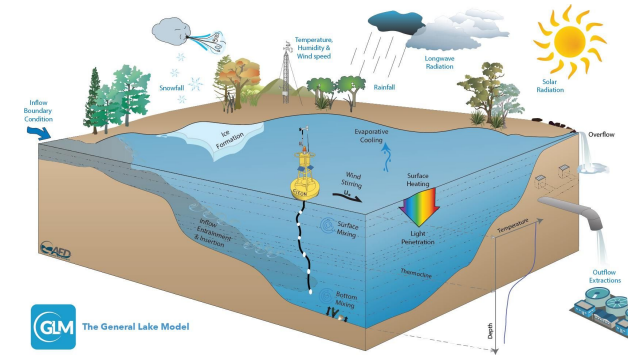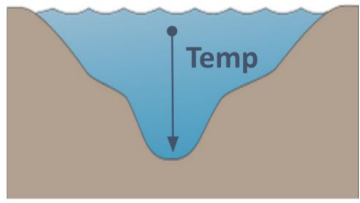
# Lake Temperature Modeling

## 1D Model of Temperature

**Meteorological Input Drivers**
E.g., longwave/shortwave radiation, air temperature, humidity, wind speed, etc.



**Target**
Temperature of water at every depth of the lake

Temp



GLM The General Lake Model

## Motivation



Growth and survival of fisheries
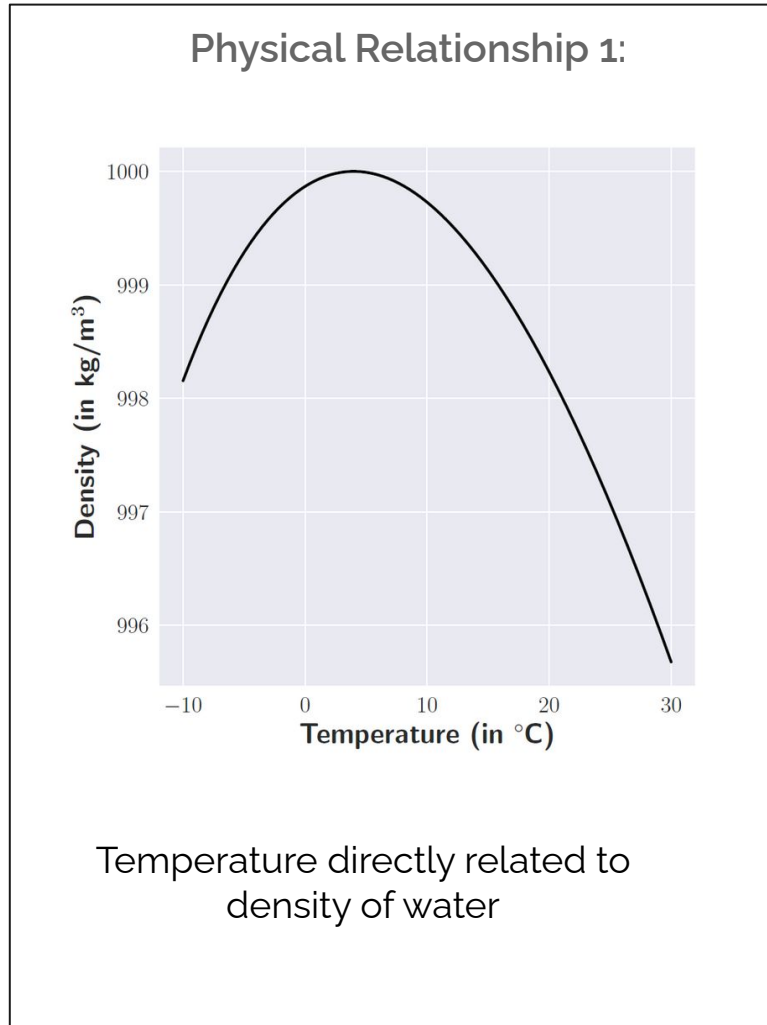


Harmful Algal Blooms



Chemical Constituents: $O_2$, C, N

# Physical Relationships of Temperature



Physical Relationship 1:

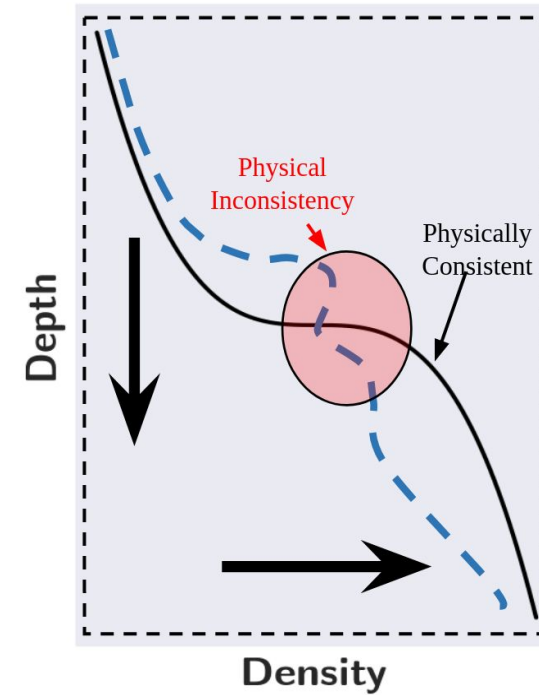Temperature directly related to density of water

# Physical Relationships of Temperature



Physical Relationship 1:

Temperature directly related to density of water



Physical Relationship 2:

Density of water **monotonically increases** with depth

# Physics-guided Neural Networks (PGNN)

The physics supervision is enforced as a soft constraint where the model is penalized when the predictions of the model violate the physics constraint.

$$\min_{\theta} \quad \boxed{L(y, \widehat{y})} \quad + \quad \boxed{\lambda_{PHY} \, L_{PHY}(\widehat{y})}$$

Empirical Error            Physics-Loss

😊 ✓ **PROS**

- **Easy to use:** Constraints can be easily incorporated as physics loss functions.

- **Unsupervised:** Physics loss functions can be evaluated on unlabeled data.

Daw, Arka, Anuj Karpatne, William D. Watkins, Jordan S. Read, and Vipin Kumar. "Physics-guided neural networks (pgnn): An application in lake temperature modeling." In *Knowledge Guided Machine Learning*, pp. 353-372. Chapman and Hall/CRC, 2022.

# PGNN shows improved generalization

Results on two different lakes: Lake Mille Lacs and Lake Mendota



(a) Results on Mille Lacs Lake

(b) Results on Lake Mendota

PGNN consistently outperforms the other baselines for both lakes showing better Test RMSE and Physics Consistency.

Daw, Arka, Anuj Karpatne, William D. Watkins, Jordan S. Read, and Vipin Kumar. "Physics-guided neural networks (pgnn): An application in lake temperature modeling." In *Knowledge Guided Machine Learning*, pp. 353-372. Chapman and Hall/CRC, 2022.

# Pretraining on Simulation Lakes

Simulation Data from the different lakes can be used to pretrain the RNN model. This will serve as a "better" initialization.

Jia, Xiaowei, Jared Willard, Anuj Karpatne, Jordan Read, Jacob Zwart, Michael Steinbach, and Vipin Kumar. "Physics guided RNNs for modeling dynamical systems: A case study in simulating lake temperature profiles." In *Proceedings of the 2019 SIAM international conference on data mining*, pp. 558-566. Society for Industrial and Applied Mathematics, 2019.

# Use Case 2:
# KGML with Uncertainty Quantification

# Organizing KGML Research: A Multi-Dimensional View

Karpatne, Jia, and Kumar. "Knowledge-guided Machine Learning: Current Trends and Future Prospects." *arXiv:2403.15989* (2024).

# Uncertainty Quantification

Generate a distribution over the predictions rather than point estimates.

- Regression: Predict the variance along with the output mean.

- Classification: Predict the confidence along with the output labels.



Aims to quantify the **robustness** of the ML models by assessing prediction reliability.

Daw, Arka, R. Quinn Thomas, Cayelan C. Carey, Jordan S. Read, Alison P. Appling, and Anuj Karpatne. "Physics-guided architecture (pga) of neural networks for quantifying uncertainty in lake temperature modeling." In *Proceedings of the 2020 siam international conference on data mining*, pp. 532-540. Society for Industrial and Applied Mathematics, 2020.

# Uncertainty Quantification with MC Dropout

A schematic representation of using Dropouts to estimate uncertainty.



The **red** nodes are dropped while the **green** nodes contribute to the output of the neural network.

**Repeat this "N" times**

Daw, Arka, Anuj Karpatne, William D. Watkins, Jordan S. Read, and Vipin Kumar. "Physics-guided neural networks (pgnn): An application in lake temperature modeling." In *Knowledge Guided Machine Learning*, pp. 353-372. Chapman and Hall/CRC, 2022.

# Approach 1: Dropouts with Physics-based Loss



Train with physics-loss functions

$$\min_{\theta} L(y, \hat{y}) + \lambda_{PHY} L_{PHY}(\hat{y})$$

Layer 1

Layer 2

Inputs

Output

The **red** nodes are dropped while the **green** nodes contribute to the output of the neural network.

Evaluation

1st MC Sample

Physically Inconsistent

nth MC Sample

$\mu_y$

😠 ✗ **Limitations**

- Trained PGL-models are physically consistent.

- Each of the perturbed dropout networks are no longer physically consistent.

Daw, Arka, Anuj Karpatne, William D. Watkins, Jordan S. Read, and Vipin Kumar. "Physics-guided neural networks (pgnn): An application in lake temperature modeling." In *Knowledge Guided Machine Learning*, pp. 353-372. Chapman and Hall/CRC, 2022.
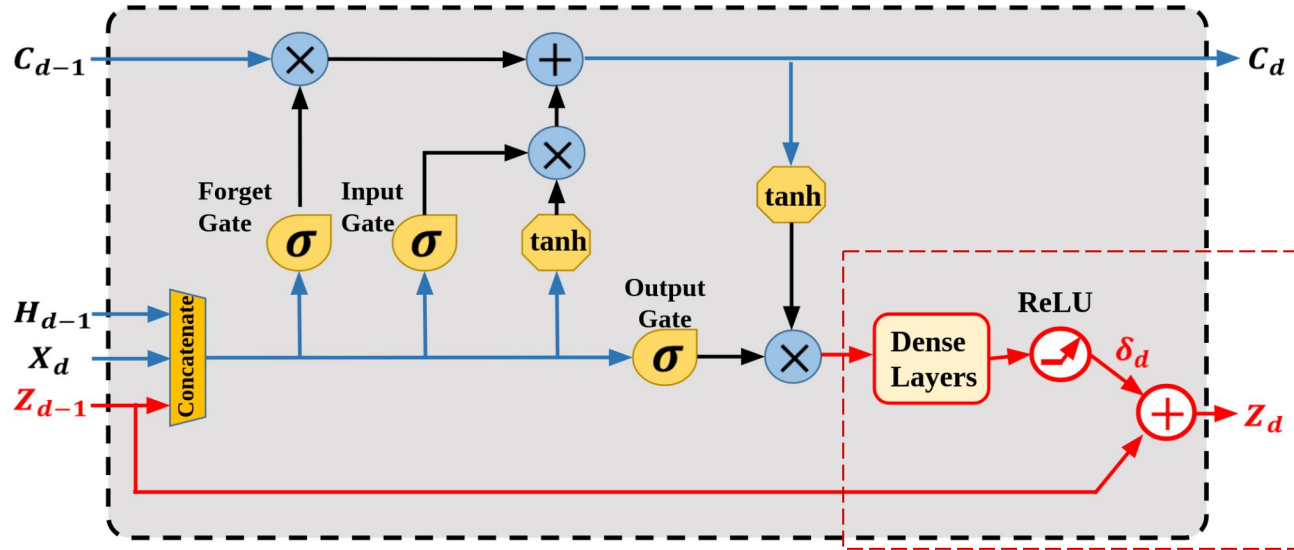
# Proposed PGA-LSTM Framework

- **Temporal Autoencoder:** Encodes the input time series to obtain a temporal embedding.

- **Monotonicity Preserving LSTM:** Enforces the monotonicity constraint on the density predictions.

- **Dense Layers:** Takes the density estimates and the input drivers to predict temperature.

Daw, Arka, Anuj Karpatne, William D. Watkins, Jordan S. Read, and Vipin Kumar. "Physics-guided neural networks (pgnn): An application in lake temperature modeling." In *Knowledge Guided Machine Learning*, pp. 353-372. Chapman and Hall/CRC, 2022.

# Monotonicity Preserving LSTM



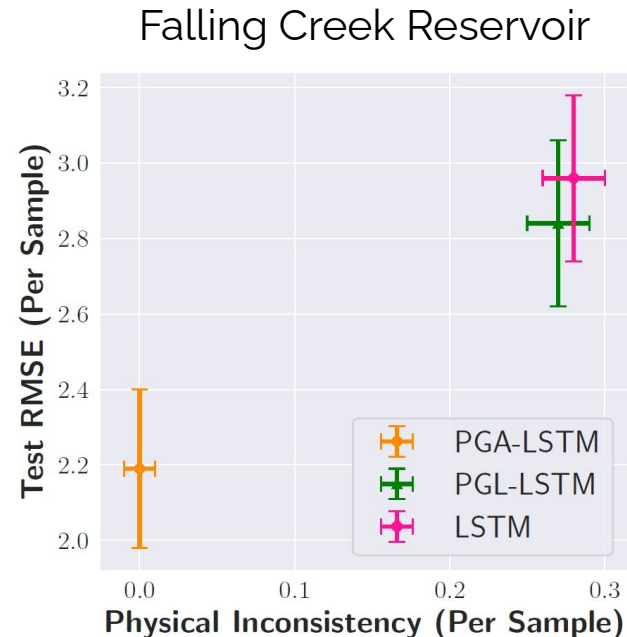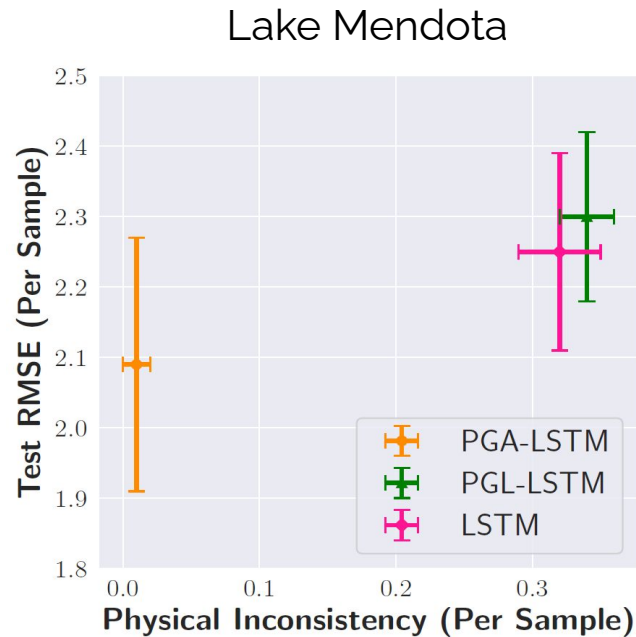Components in **red** represent the novel physics-informed innovations in LSTM

**Key Idea**

The **ReLU function** ensures that the residual outputs are non-negative, thus enforcing the monotonicity constraint.

The monotonicity preserving LSTM:
1. Adds a layer of **interpretability** into the model outputs,
2. Makes it more **robust** to small perturbations in the model weights
3. Ensures physics-**generalization** on unseen test set.

Daw, Arka, Anuj Karpatne, William D. Watkins, Jordan S. Read, and Vipin Kumar. "Physics-guided neural networks (pgnn): An application in lake temperature modeling." In *Knowledge Guided Machine Learning*, pp. 353-372. Chapman and Hall/CRC, 2022.

# Impact on predictive performance and physical consistency



PGA-LSTM improves the Test RMSE while always being physically consistent across both lakes.

Daw, Arka, Anuj Karpatne, William D. Watkins, Jordan S. Read, and Vipin Kumar. "Physics-guided neural networks (pgnn): An application in lake temperature modeling." In *Knowledge Guided Machine Learning*, pp. 353-372. Chapman and Hall/CRC, 2022.

# Monotonicity Preserving LSTM



The mean and the variance of the three models are computed from **100 MC-Samples**.

Daw, Arka, Anuj Karpatne, William D. Watkins, Jordan S. Read, and Vipin Kumar. "Physics-guided neural networks (pgnn): An application in lake temperature modeling." In *Knowledge Guided Machine Learning*, pp. 353-372. Chapman and Hall/CRC, 2022.
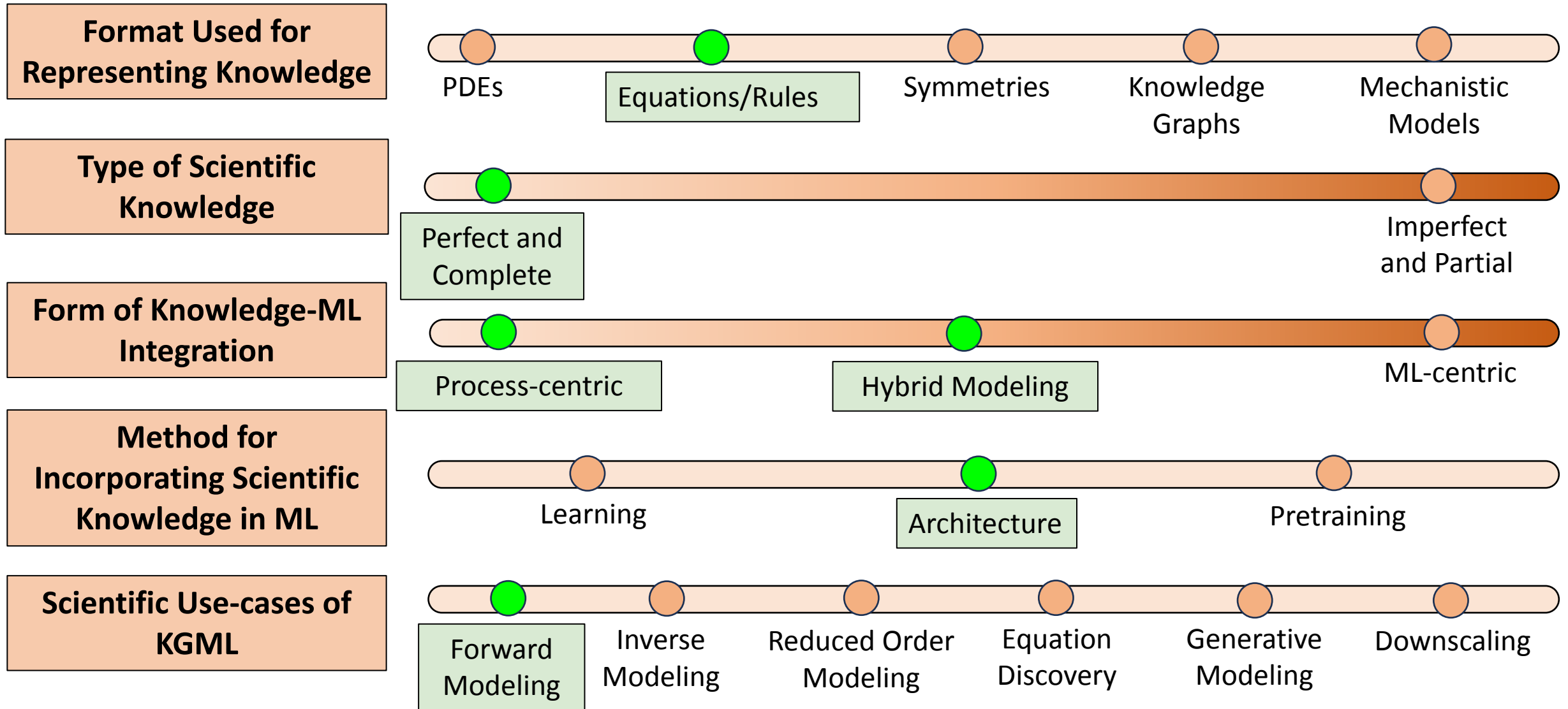
# Monotonicity Preserving LSTM



The PGA-LSTM samples are always physically consistent while PGL-LSTM and LSTM samples are very much physically inconsistent.

✅ Predictions are **more robust to minor perturbations** in model weights!

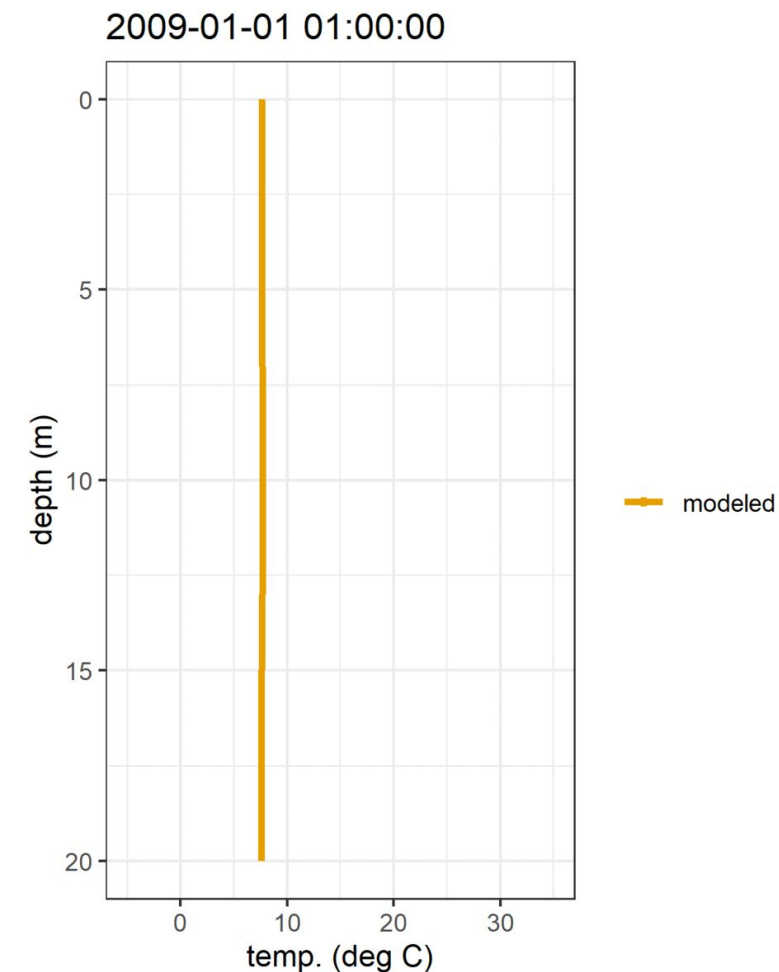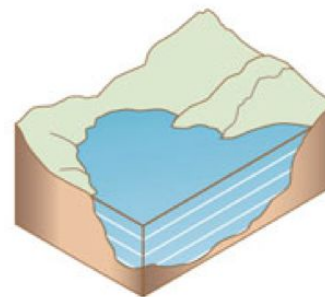Daw, Arka, Anuj Karpatne, William D. Watkins, Jordan S. Read, and Vipin Kumar. "Physics-guided neural networks (pgnn): An application in lake temperature modeling." In *Knowledge Guided Machine Learning*, pp. 353-372. Chapman and Hall/CRC, 2022.

# Use Case 3:
# Hybrid Modeling

# Organizing KGML Research: A Multi-Dimensional View



Karpatne, Jia, and Kumar. "Knowledge-guided Machine Learning: Current Trends and Future Prospects." *arXiv:2403.15989* (2024).

# Process-based Modeling

- plethora of model approaches:
    - **energy-balance** models: mixing depth by external energy
    - **turbulence-based** models: advanced turbulence-closure

**1D: One-dimensional**



2009-01-01 01:00:00



Ladwig, Robert, Arka Daw, Ellen A. Albright, Cal Buelo, Anuj Karpatne, Michael Frederick Meyer, Abhilash Neog, Paul C. Hanson, and Hilary A. Dugan. "Modular Compositional Learning Improves 1D Hydrodynamic Lake Model Performance by Merging Process-Based Modeling With Deep Learning." Journal of Advances in Modeling Earth Systems 16, no. 1 (2024)

# Process-based Modeling

- plethora of model approaches:
  - **energy-balance** models: mixing depth by external energy
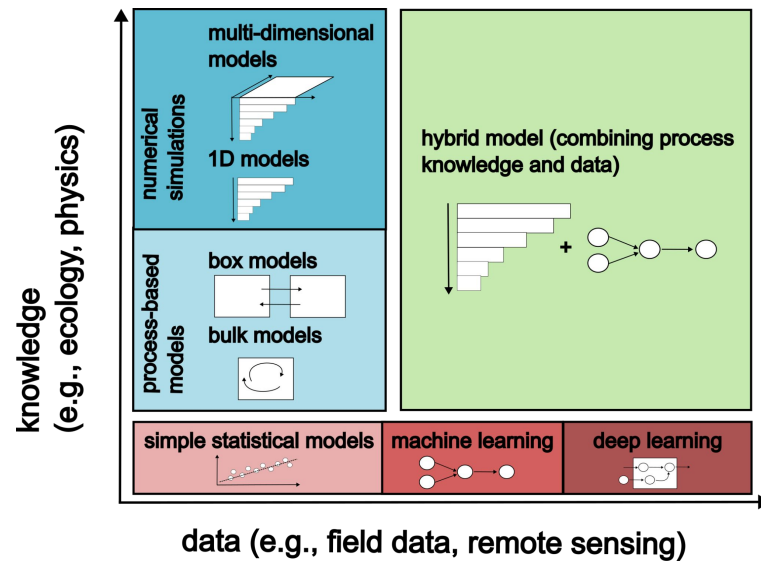  - **turbulence-based** models: advanced turbulence-closure

1D: One-dimensional

2009-01-01 01:00:00
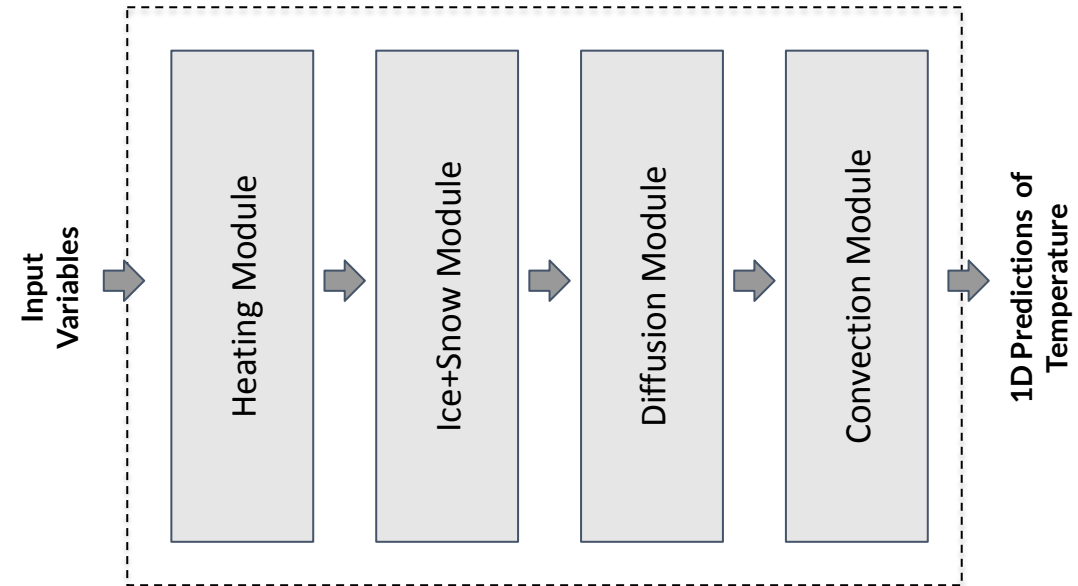


**Can we combine these process models with data?**



Ladwig, Robert, Arka Daw, Ellen A. Albright, Cal Buelo, Anuj Karpatne, Michael Frederick Meyer, Abhilash Neog, Paul C. Hanson, and Hilary A. Dugan. "Modular Compositional Learning Improves 1D Hydrodynamic Lake Model Performance by Merging Process-Based Modeling With Deep Learning." Journal of Advances in Modeling Earth Systems 16, no. 1 (2024)

# Modularized 1D Model

**Modularized Process Models:**
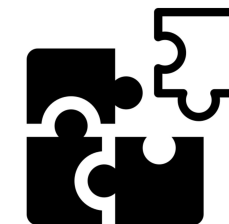
a)   heating (atmosphere and geothermal)
b)   ice, snow and snow-ice formation
c)   vertical diffusion
d)   convective overturn

😠 ❌ **CONS**

- **Imperfect Module:** All of the physics modules are not perfect, i.e., some of the physical phenomena are more complex.
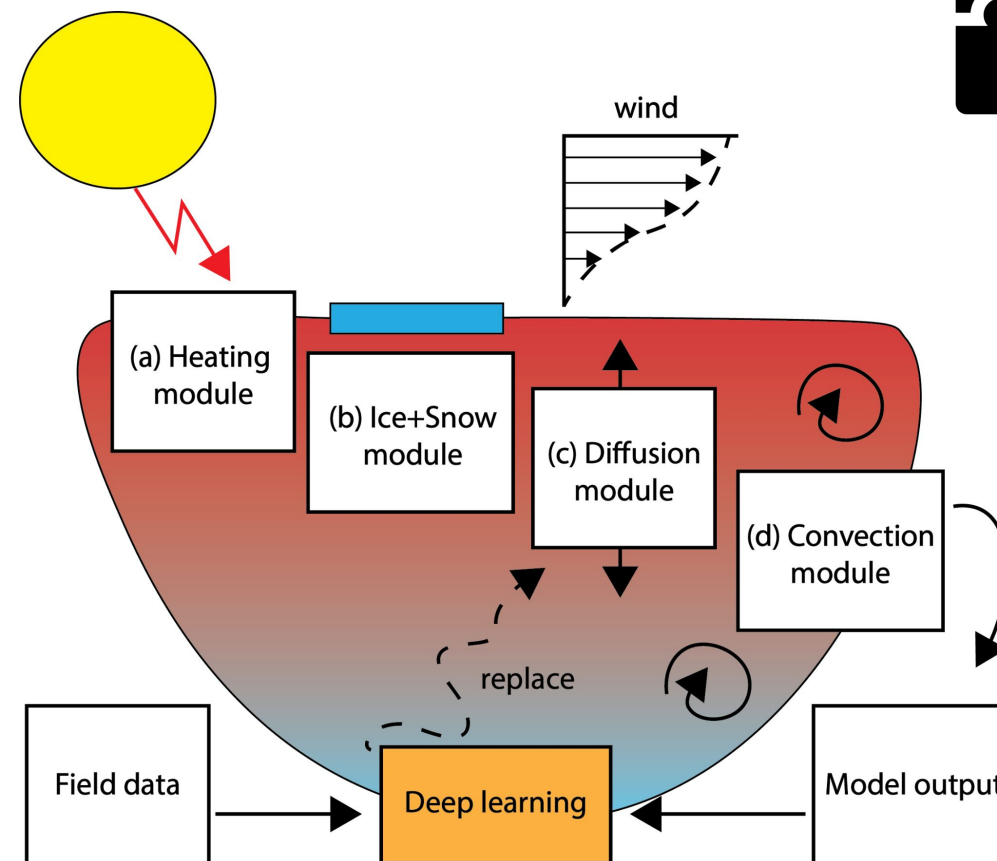


Ladwig, Robert, Arka Daw, Ellen A. Albright, Cal Buelo, Anuj Karpatne, Michael Frederick Meyer, Abhilash Neog, Paul C. Hanson, and Hilary A. Dugan. "Modular Compositional Learning Improves 1D Hydrodynamic Lake Model Performance by Merging Process-Based Modeling With Deep Learning." Journal of Advances in Modeling Earth Systems 16, no. 1 (2024)

# Modular Compositional Learning (MCL)

**Imperfect Modules:** Diffusion Module

**Idea:** Replace the imperfect modules with deep learning based models.

## ✅ PROS

- **Richer Physics knowledge:** We retain the interpretability and knowledge of the modular process based modules.

- **Hybrid modeling:** Deep learning modules learns to dynamics of the necessary "missing" module (in this case diffusion module) to learn a more accurate model.



Ladwig, Robert, Arka Daw, Ellen A. Albright, Cal Buelo, Anuj Karpatne, Michael Frederick Meyer, Abhilash Neog, Paul C. Hanson, and Hilary A. Dugan. "Modular Compositional Learning Improves 1D Hydrodynamic Lake Model Performance by Merging Process-Based Modeling With Deep Learning." Journal of Advances in Modeling Earth Systems 16, no. 1 (2024)

# Modular Compositional Learning (MCL)

**Process-based framework**



Robert Ladwig

Ladwig, Robert, Arka Daw, Ellen A. Albright, Cal Buelo, Anuj Karpatne, Michael Frederick Meyer, Abhilash Neog, Paul C. Hanson, and Hilary A. Dugan. "Modular Compositional Learning Improves 1D Hydrodynamic Lake Model Performance by Merging Process-Based Modeling With Deep Learning." Journal of Advances in Modeling Earth Systems 16, no. 1 (2024)

# Modular Compositional Learning (MCL)

**Process-based framework**

| (a)   Heat | (b) Ice | (c) Diffusion | (d) Convection |

**Pretraining**

| (a)   Heat | (b) Ice | (c) Diffusion | (d) Convection |

Process-based, pretrained deep learning, finetuned deep learning

Robert Ladwig

Ladwig, Robert, Arka Daw, Ellen A. Albright, Cal Buelo, Anuj Karpatne, Michael Frederick Meyer, Abhilash Neog, Paul C. Hanson, and Hilary A. Dugan. "Modular Compositional Learning Improves 1D Hydrodynamic Lake Model Performance by Merging Process-Based Modeling With Deep Learning." Journal of Advances in Modeling Earth Systems 16, no. 1 (2024)

# Modular Compositional Learning (MCL)

**Process-based framework**

| (a) | Heat | (b) Ice | (c) Diffusion | (d) Convection |

**Pretraining**

| (a) | Heat | (b) Ice | (c) Diffusion | (d) Convection |

**Finetuning**

Observed data

| (a) | Heat | (b) Ice | (c) Diffusion | (d) Convection |

Process-based, pretrained deep learning, finetuned deep learning

Robert Ladwig

Ladwig, Robert, Arka Daw, Ellen A. Albright, Cal Buelo, Anuj Karpatne, Michael Frederick Meyer, Abhilash Neog, Paul C. Hanson, and Hilary A. Dugan. "Modular Compositional Learning Improves 1D Hydrodynamic Lake Model Performance by Merging Process-Based Modeling With Deep Learning." Journal of Advances in Modeling Earth Systems 16, no. 1 (2024)

# Modular Compositional Learning (MCL)



Ladwig, Robert, Arka Daw, Ellen A. Albright, Cal Buelo, Anuj Karpatne, Michael Frederick Meyer, Abhilash Neog, Paul C. Hanson, and Hilary A. Dugan. "Modular Compositional Learning Improves 1D Hydrodynamic Lake Model Performance by Merging Process-Based Modeling With Deep Learning." Journal of Advances in Modeling Earth Systems 16, no. 1 (2024)

# Empirical Evaluation (Test Period 2015-17)

Comparing Observed Data and Processed-based model

**A** Observed data

**D** Process-based framework

Test RMSE: 4.46

**1** Process-based model framework

a) Process-based
b) Process-based
c) Process-based
d) Process-based

Ladwig, Robert, Arka Daw, Ellen A. Albright, Cal Buelo, Anuj Karpatne, Michael Frederick Meyer, Abhilash Neog, Paul C. Hanson, and Hilary A. Dugan. "Modular Compositional Learning Improves 1D Hydrodynamic Lake Model Performance by Merging Process-Based Modeling With Deep Learning." Journal of Advances in Modeling Earth Systems 16, no. 1 (2024)

# Empirical Evaluation (Test Period 2015-17)

Comparing the models:

1. After pretraining each of the deep learning models on simulation data.
2. Finetuning the entire deep learning pipeline on observed data.



A Observed data

D Process-based framework
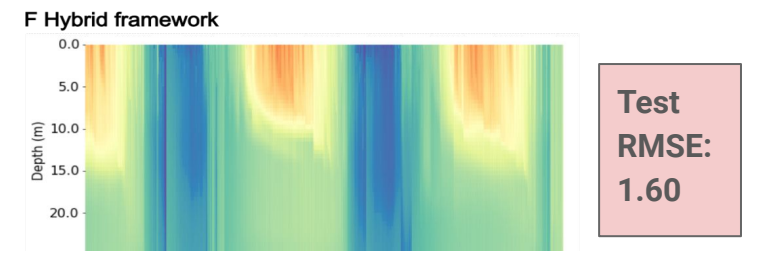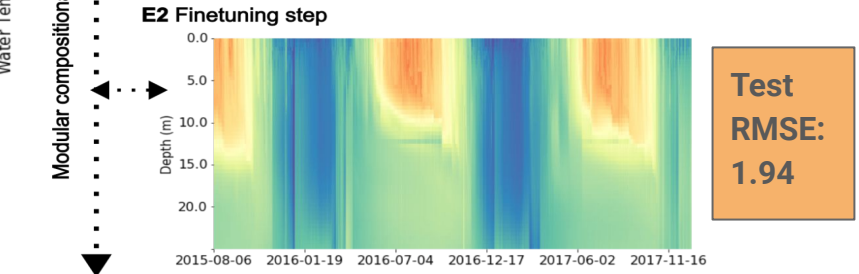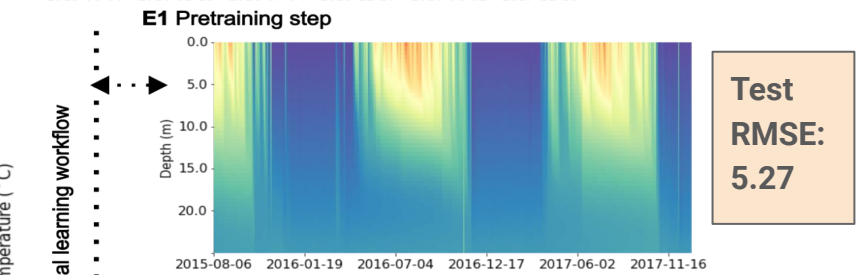
Test RMSE: 4.46

E1 Pretraining step

Test RMSE: 5.27

E2 Finetuning step

Test RMSE: 1.94

2 Pretrained deep learning framework

a) Deep learning
b) Deep learning
c) Deep learning
d) Deep learning

3 Finetuned deep learning framework

a) Deep learning
b) Deep learning
c) Deep learning
d) Deep learning

Observed data

Ladwig, Robert, Arka Daw, Ellen A. Albright, Cal Buelo, Anuj Karpatne, Michael Frederick Meyer, Abhilash Neog, Paul C. Hanson, and Hilary A. Dugan. "Modular Compositional Learning Improves 1D Hydrodynamic Lake Model Performance by Merging Process-Based Modeling With Deep Learning." Journal of Advances in Modeling Earth Systems 16, no. 1 (2024)
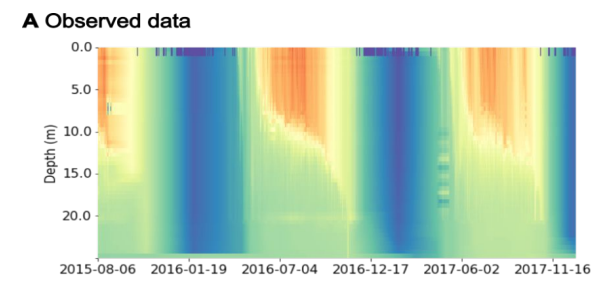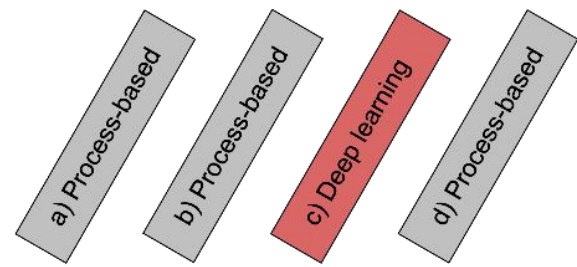
# Empirical Evaluation (Test Period 2015-17)

Plugging the deep-learning module into the process-based module pipeline.



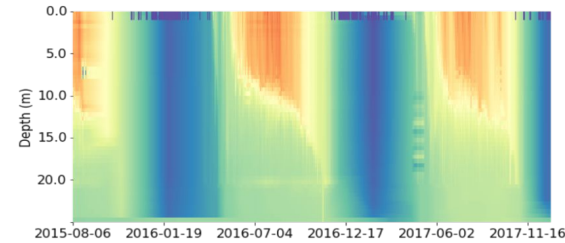Ladwig, Robert, Arka Daw, Ellen A. Albright, Cal Buelo, Anuj Karpatne, Michael Frederick Meyer, Abhilash Neog, Paul C. Hanson, and Hilary A. Dugan. "Modular Compositional Learning Improves 1D Hydrodynamic Lake Model Performance by Merging Process-Based Modeling With Deep Learning." Journal of Advances in Modeling Earth Systems 16, no. 1 (2024)

# Empirical Evaluation (Test Period 2015-17)



**A** Deep learning model (no process)

Observed data

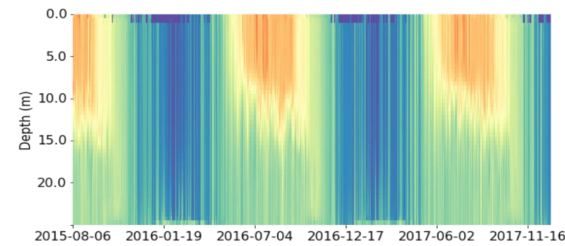Boundary conditions → Deep learning model → $T_{HY1}(z)$
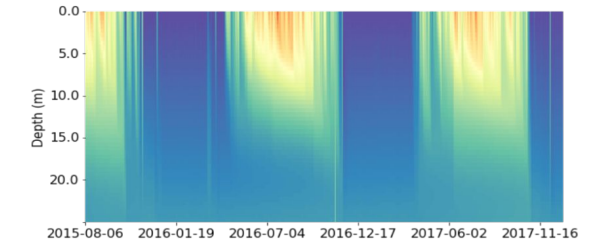
Loss

Model output

**A** Observed data

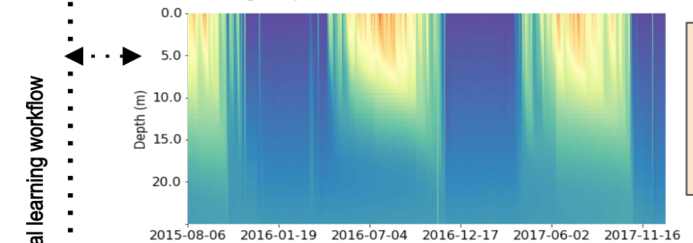**B** Deep learning model (no process)

**D** Process-based framework

**Test RMSE: 4.46**

**E1** Pretraining step

**Test RMSE: 5.27**

**E2** Finetuning step

**Test RMSE: 1.94**

**Test RMSE: 2.10**

Modular compositional learning workflow

**F** Hybrid framework

**Test RMSE: 1.60**

Ladwig, Robert, Arka Daw, Ellen A. Albright, Cal Buelo, Anuj Karpatne, Michael Frederick Meyer, Abhilash Neog, Paul C. Hanson, and Hilary A. Dugan. "Modular Compositional Learning Improves 1D Hydrodynamic Lake Model Performance by Merging Process-Based Modeling With Deep Learning." Journal of Advances in Modeling Earth Systems 16, no. 1 (2024)

# Current Work in MCL



## 4 Hybrid model framework

a) Process-based   b) Process-based   c) Deep learning   d) Process-based

## 1D Lake Physics with MCL

## 1D Water Quality with MCL

Model Structure:

T   WQ

f) Atmospheric gas exchange
g) Net primary production
a) Heat generation
b) Ice and snow formation
h) Ecosystem respiration
c) Diffusive transport          eddy diffusivity   i) Diffusive transport
d) Convective wind mixing                                          WQ
light extinction
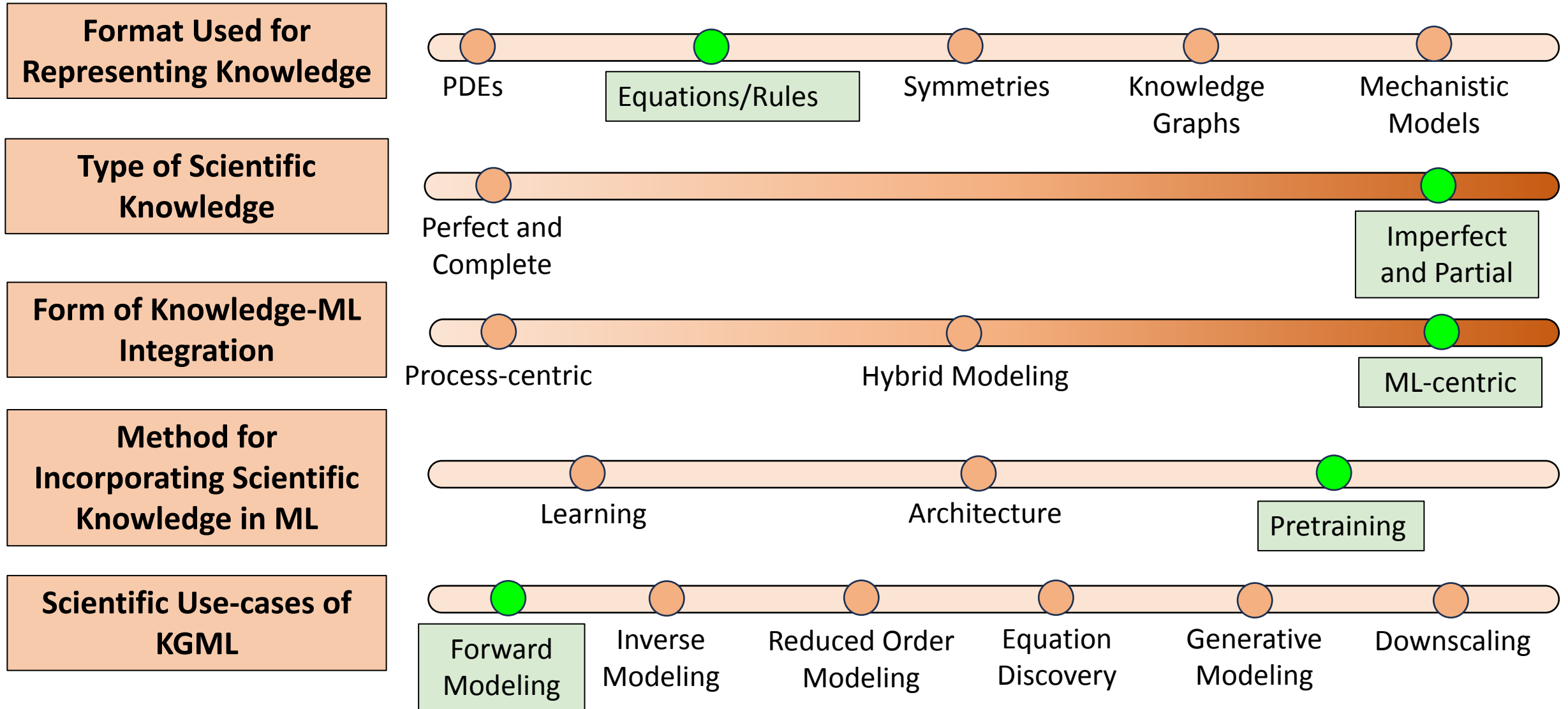e) Density instabilities                                          T

1D-AEMpy

## 1D Lake Physics with MCL: memory for multiple lakes

# Use Case 5:
# Lake Chlorophyll-a Prediction

# Organizing KGML Research: A Multi-Dimensional View



Karpatne, Jia, and Kumar. "Knowledge-guided Machine Learning: Current Trends and Future Prospects." *arXiv:2403.15989* (2024).
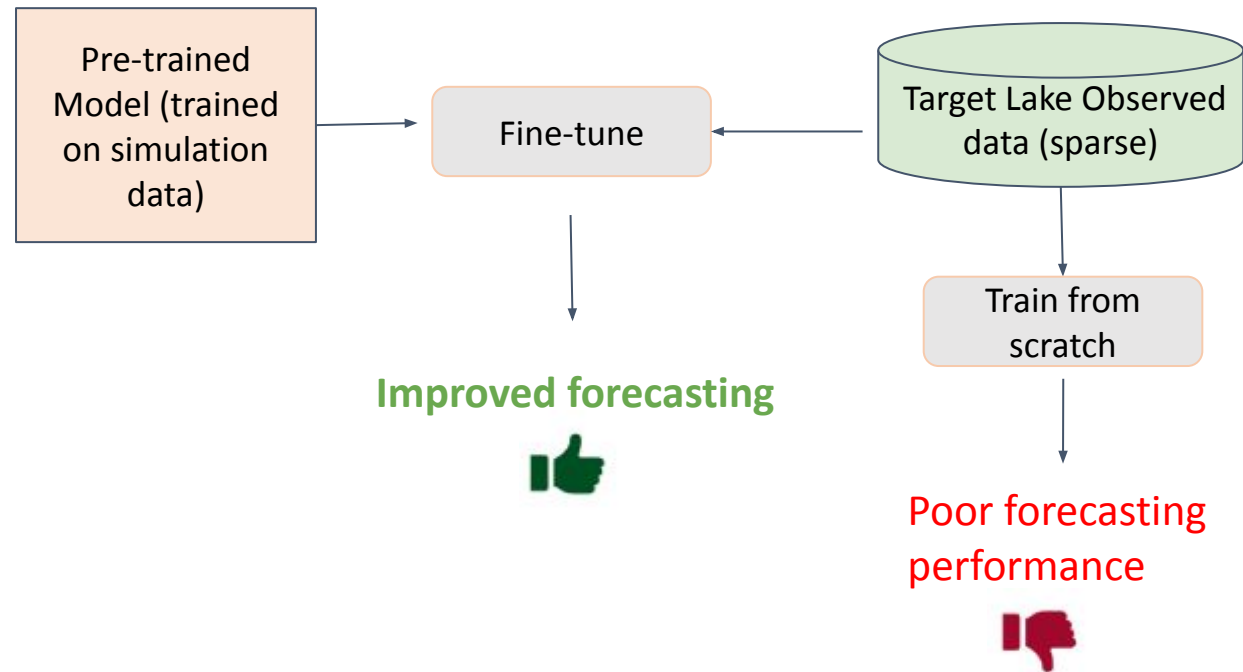
# Transfer Learning for Chlorophyll-a Prediction
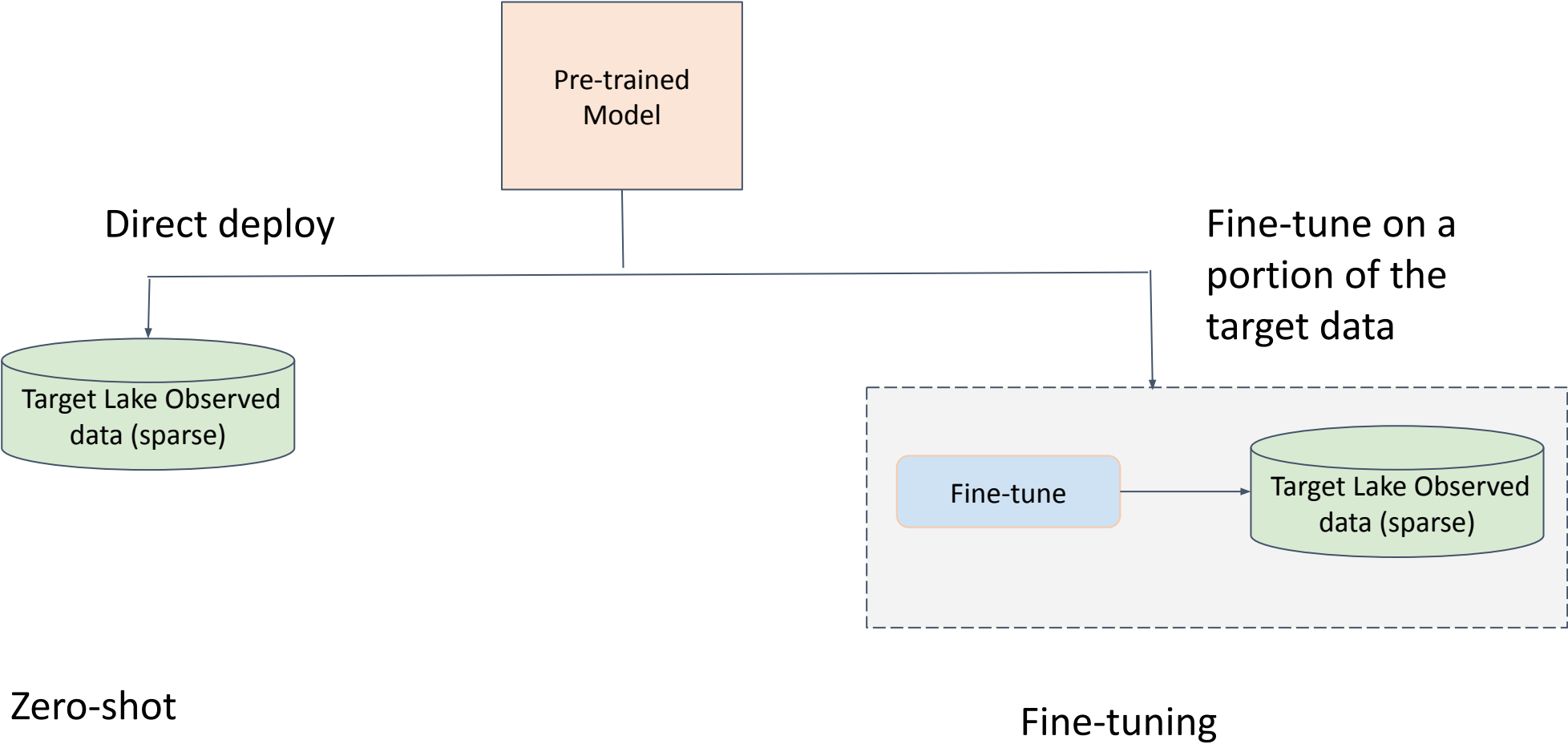
**Problem Context:**
- Observations of chlorophyll-a vary across lakes, some being well-observed, others less-observed.
- Deep learning models are data-hungry, show poor forecasting performance on target lakes with sparse data.

**Research Question:** *How can we improve forecasting performance of chlorophyll-a on lakes with few observations?*

**Approach:** Instead of "training from scratch" *transfer Learning* enables us to transfer knowledge learned from data-rich source lakes (in the form of pre-trained models) to target lakes.

# Types of Transfer Learning methods



Pre-trained Model

Direct deploy

Fine-tune on a portion of the target data

Target Lake Observed data (sparse)

Fine-tune

Target Lake Observed data (sparse)

Zero-shot

Fine-tuning

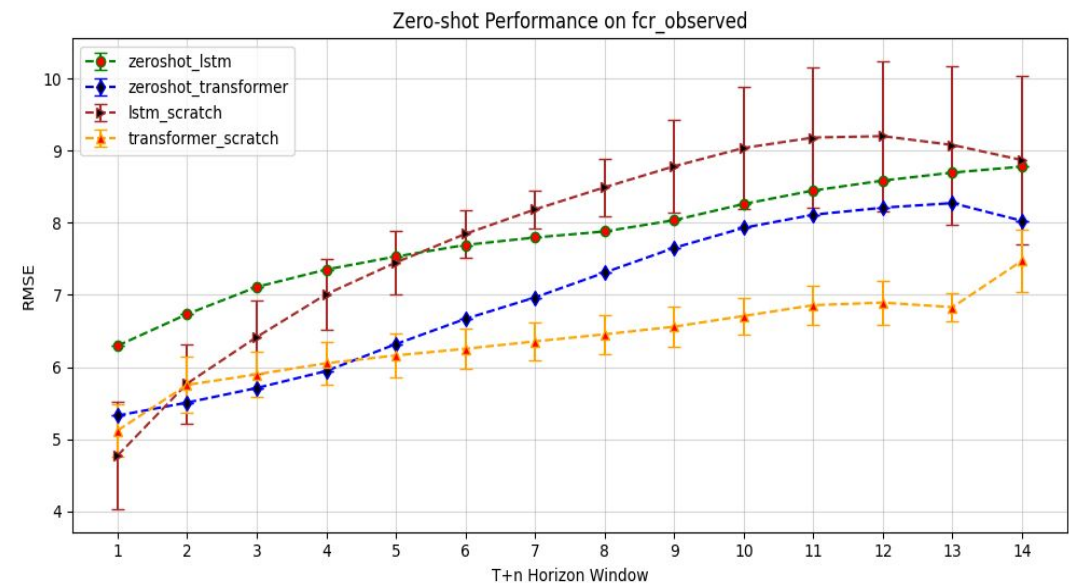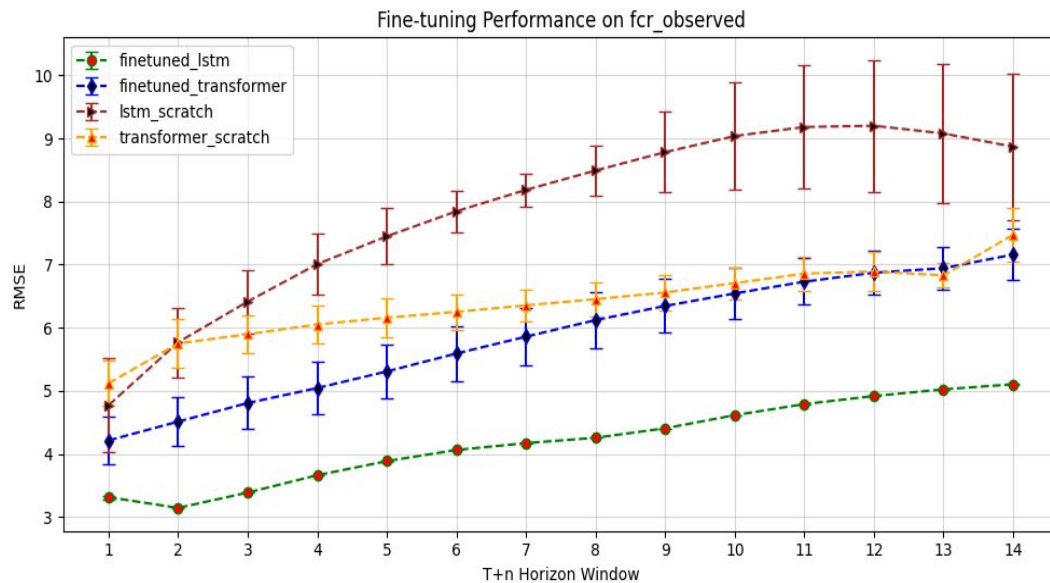# Transfer Learning for Chlorophyll-a Prediction

**Problem Setup**
Pre-training: Model pre-trained on simulation data of lakes Mendota, Sunapee, FCR.
Models: LSTM [1], Transformer [2]

Data split in target lake = 70:30
Model trained/fine-tuned on the 70% and tested on the 30% data.

>Following results are on the test set (i.e. 30% of data)



1. Hochreiter, S., & Schmidhuber, J"urgen. (1997). Long short-term memory. *Neural Computation*, *9*(8), 1735–1780.
2. Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).
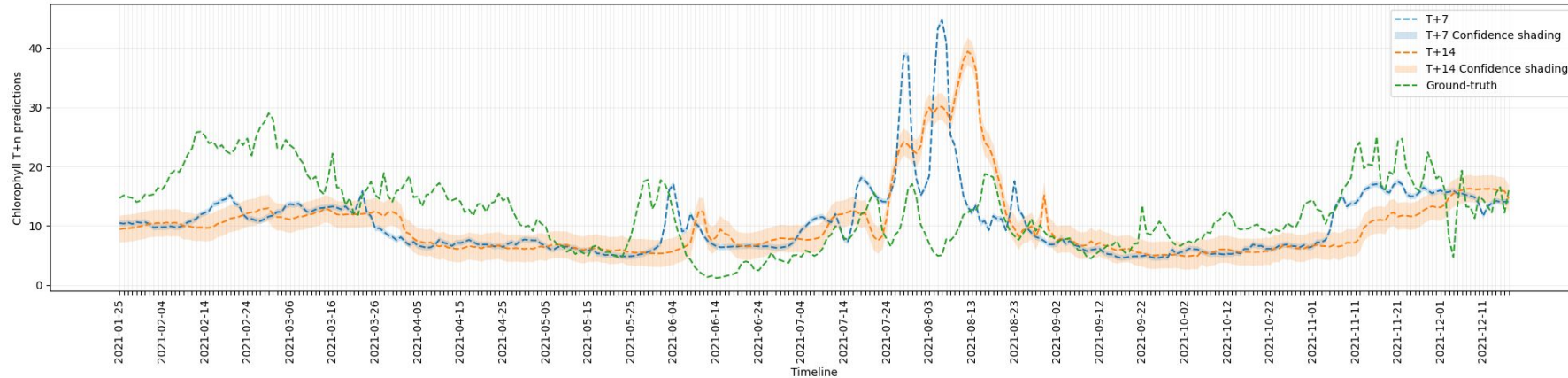
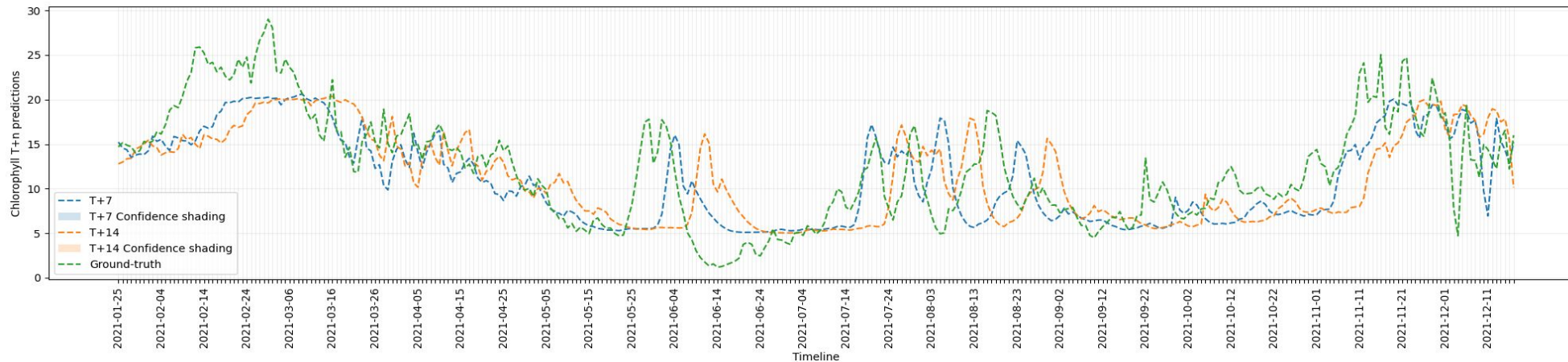# Transfer Learning for Chlorophyll-a Prediction



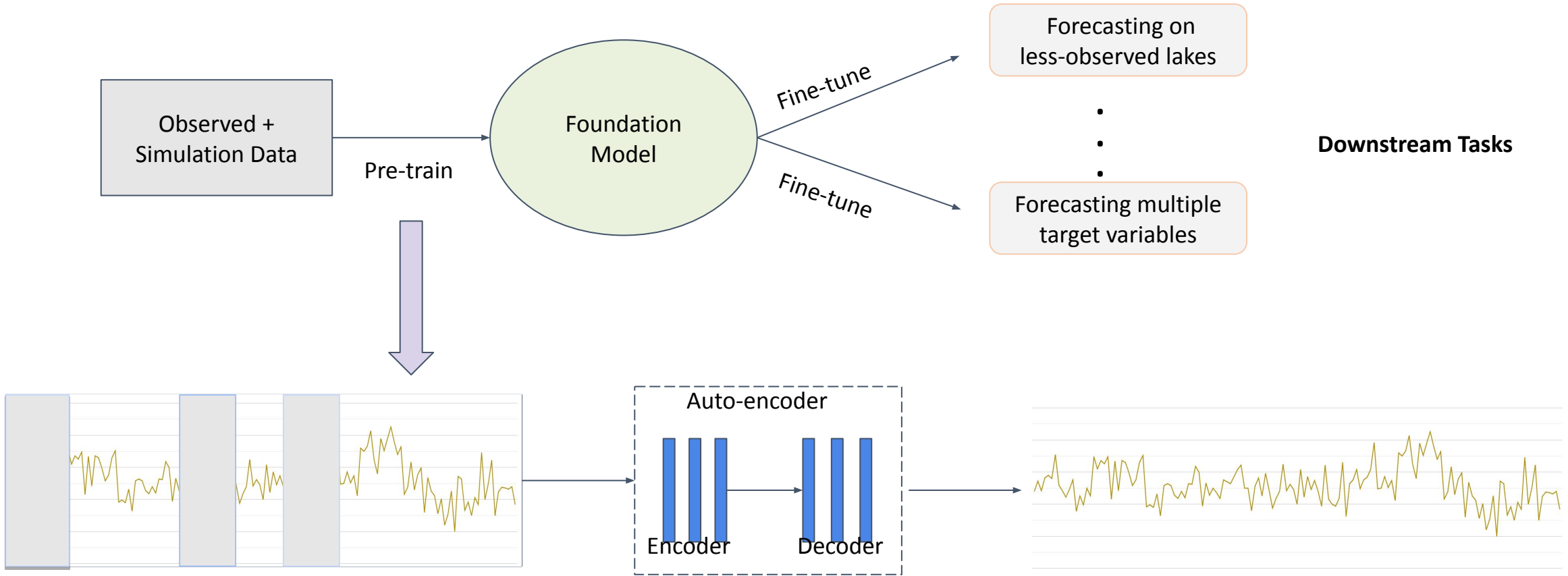Fig. 1 Predictions on FCR observed Test portion - Model trained from scratch

- Fine-tuned model aligns with the ground-truth scale of chlorophyll data
- Fine-tuned model shows relatively more confident predictions



LSTM model

Fig. 2 Predictions on FCR observed Test portion - Model fine-tuned on FCR observed

# Towards a Foundation Model



Observed + Simulation Data — Pre-train → Foundation Model

Fine-tune → Forecasting on less-observed lakes

Fine-tune → Forecasting multiple target variables

**Downstream Tasks**

Auto-encoder

Encoder — Decoder

Learning Time-series representation

NAIRR Pilot 240161

**LakeGPT:** Building A Foundation Model for Aquatic Sciences

48