# Toward Scientific Foundation Models for Aquatic Ecosystems

## Abhilash Neog

PhD Computer Science, Virginia Tech

**H33D: Advancing Water Science Through Artificial Intelligence: Lessons, Strategies, and New Frontiers**
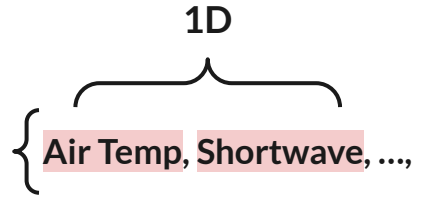
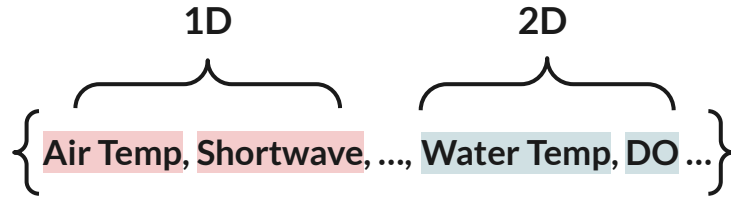*AGU 2025, New Orleans, LA*

# Motivation
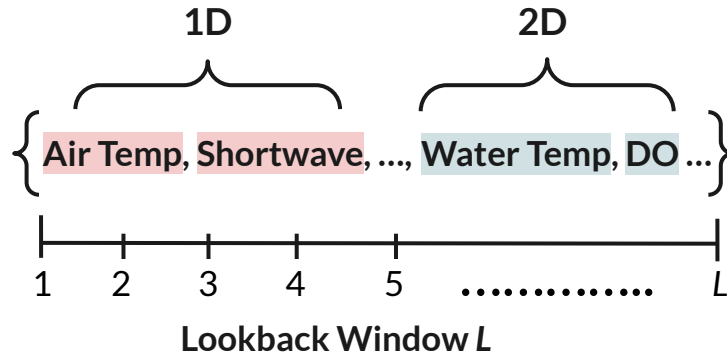
Aquatic Ecosystem 1

# Motivation

**Aquatic Ecosystem 1**

**1D**

**Air Temp**, **Shortwave**, ...,

# Motivation

Aquatic Ecosystem 1

1D          2D

{ Air Temp, Shortwave, ..., Water Temp, DO ... }

# Motivation



Aquatic Ecosystem 1

1D · 2D

{ Air Temp, Shortwave, ..., Water Temp, DO ... }

1 2 3 4 5 ............. $L$

Lookback Window $L$

# Motivation

Aquatic Ecosystem 1



1D    2D

{ Air Temp, Shortwave, ..., Water Temp, DO ... }  →  { Water Temp, DO, .... }

1   2   3   4   5   ............   $L$        $L+1$  $L+2$  ........  $L+H$

**Lookback Window** $L$          **Horizon Window** $H$

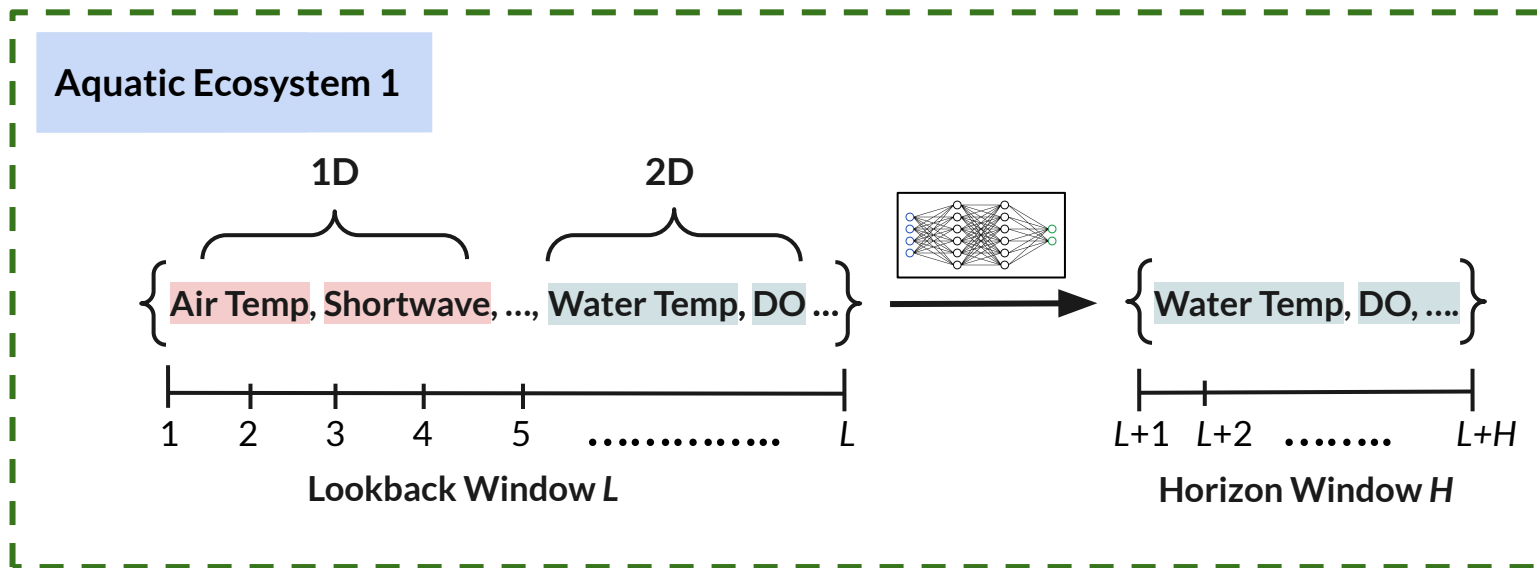# Motivation

# Motivation

# Motivation



⊗ **Different subset of variables available in different ecosystems**

⊗ **Large amounts of missing data (*e.g., Falling Creeks Reservoir, VA, has 70% missing data, 2017-04 to 2022-10* )**

Aquatic Ecosystem 1

1D 2D

Air Temp, Shortwave, ..., Water Temp, DO ...  Water Temp, DO, ....
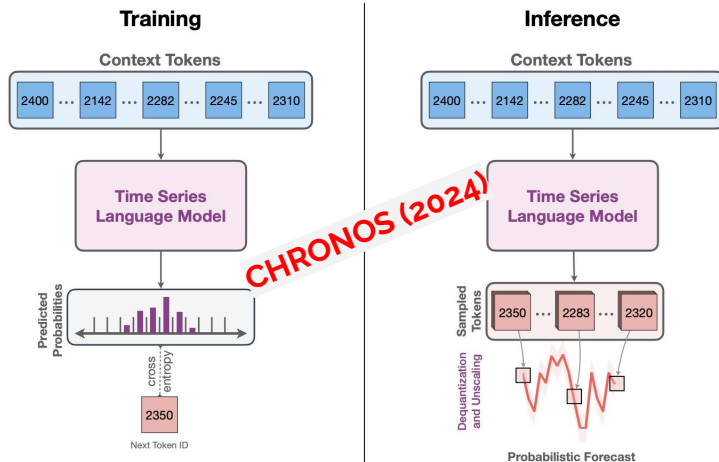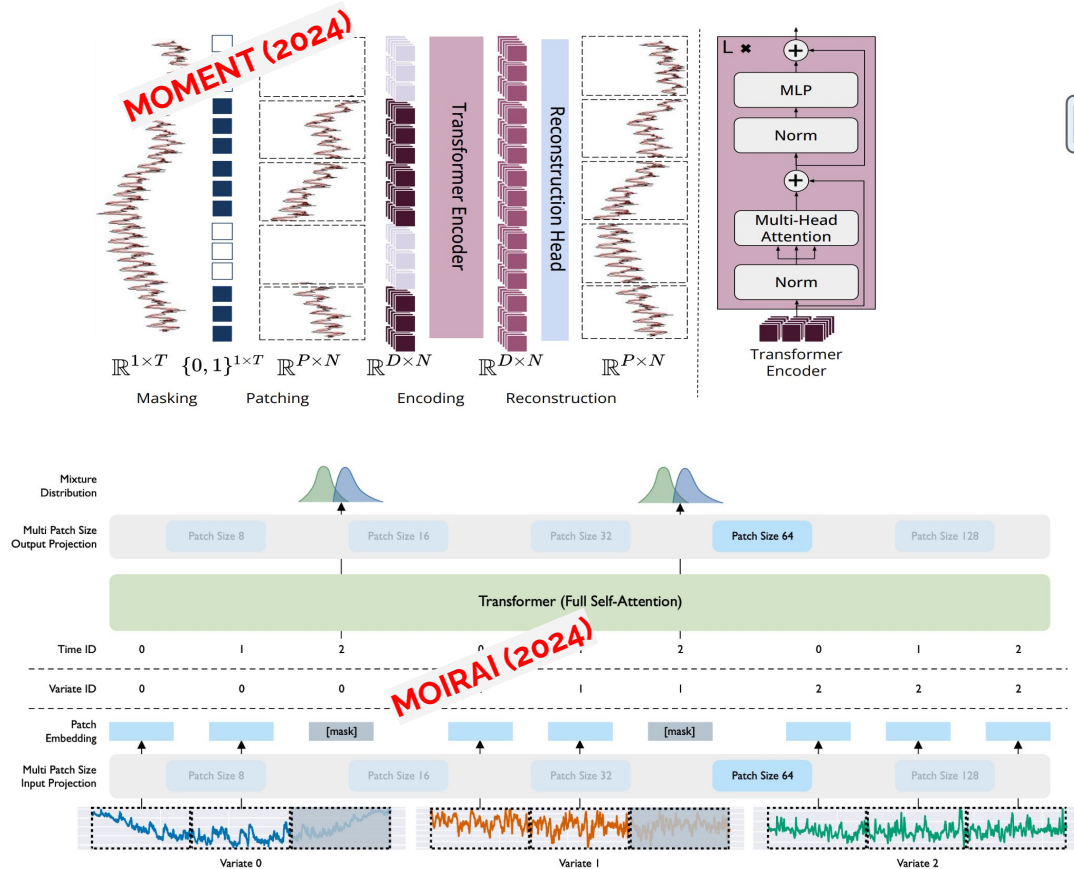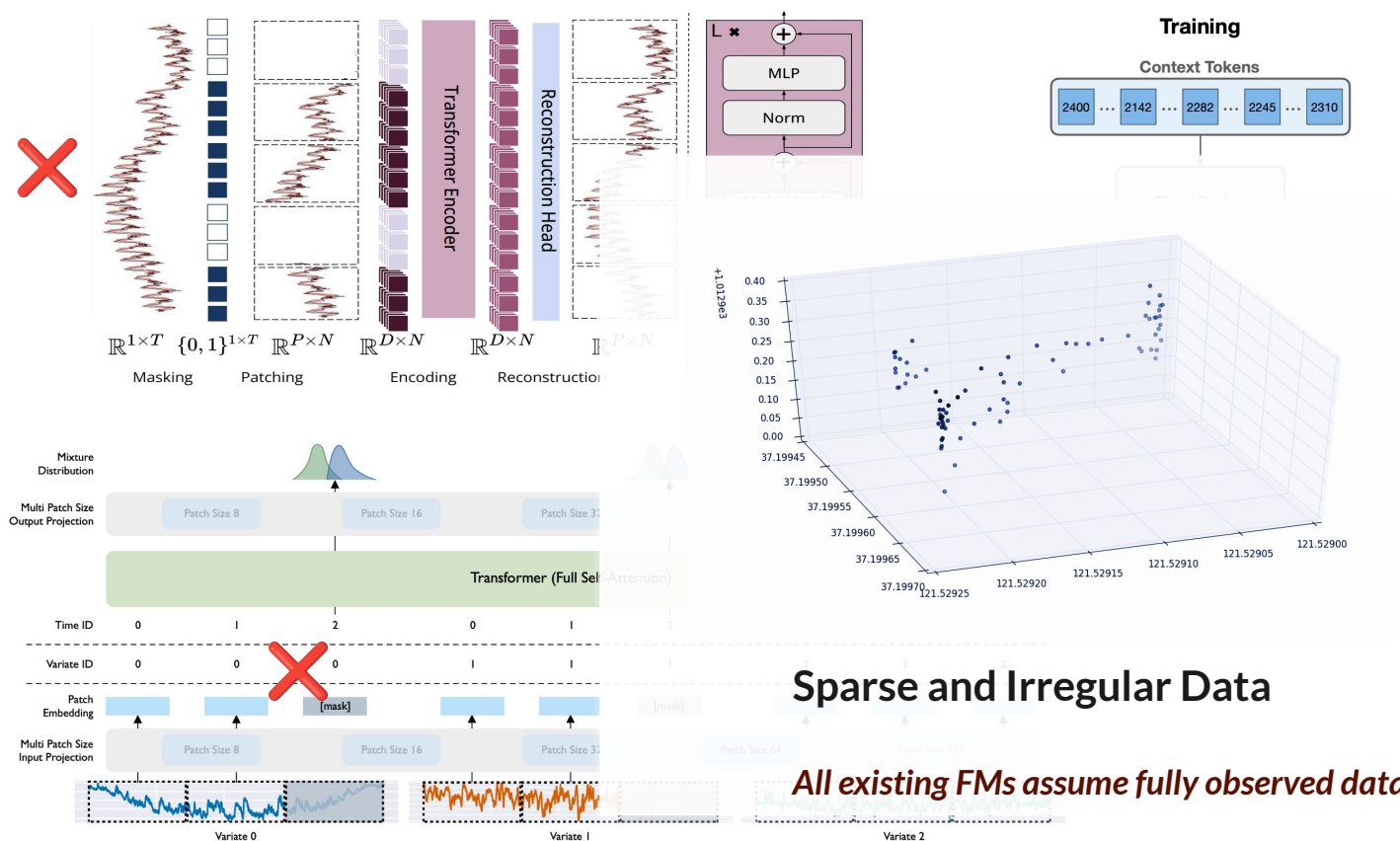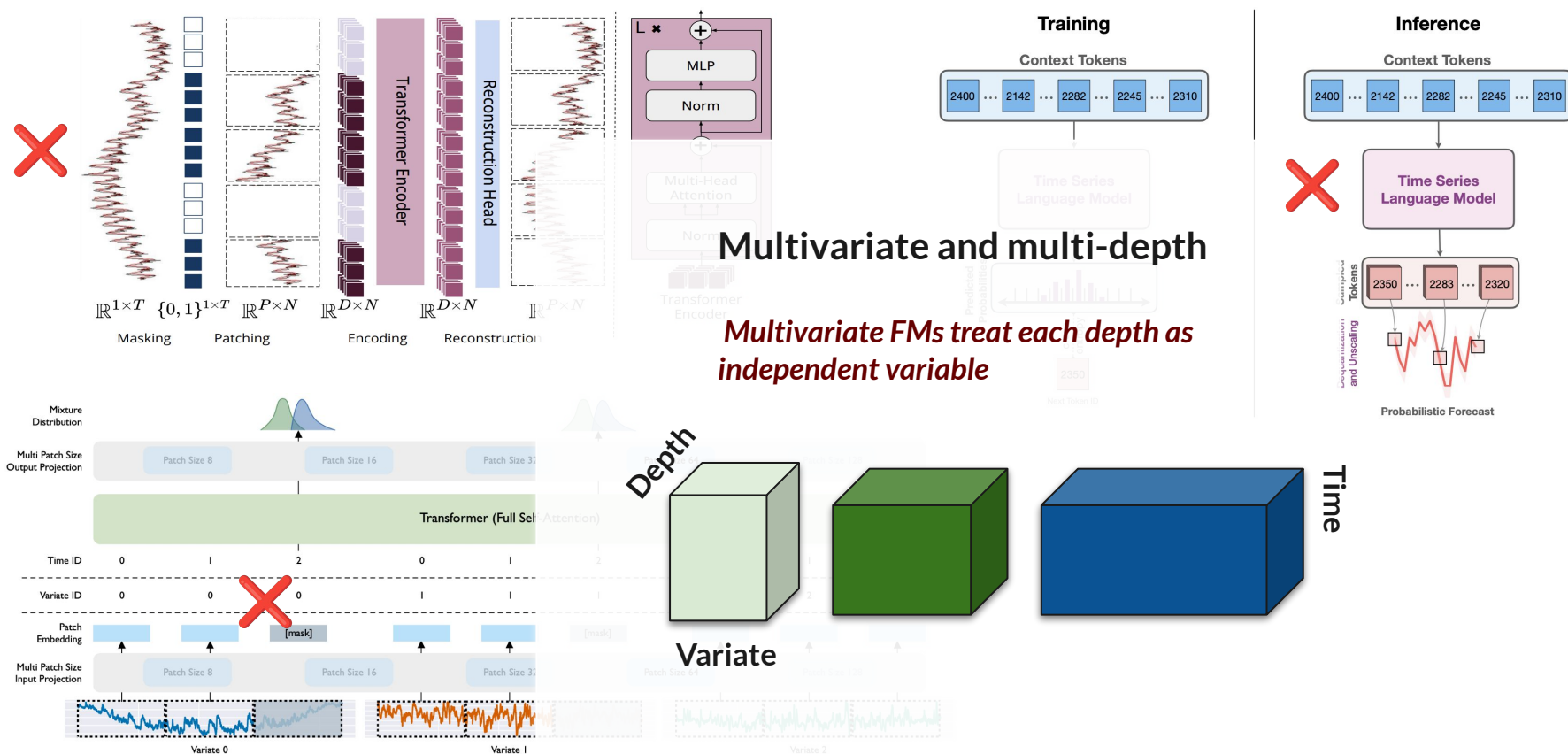
1  2  3  4  5  .............  L  L+1  L+2  ........  L+H

## Can we build a Foundation Model that can generalize across different lake ecosystems with different variables and missing values?

Transfer Learning

Aquatic Ecosystem 2
*Well observed*

⊗ Different subset of variables available in different ecosystems

⊗ **Large amounts of missing data (*e.g., Falling Creeks Reservoir, VA, has 70% missing data, 2017-04 to 2022-10* )**

# Existing Foundation Models

$\mathbb{R}^{1 \times T}$  $\{0,1\}^{1 \times T}$  $\mathbb{R}^{P \times N}$  $\mathbb{R}^{D \times N}$  $\mathbb{R}^{D \times N}$  $\mathbb{R}^{P \times N}$

Masking   Patching   Encoding   Reconstruction

**Training**

Context Tokens

| 2400 | ... | 2142 | ... | 2282 | ... | 2245 | ... | 2310 |

**Inference**

Context Tokens

| 2400 | ... | 2142 | ... | 2282 | ... | 2245 | ... | 2310 |

**Time Series Language Model**

| 2350 | ... | 2283 | ... | 2320 |

Probabilistic Forecast

## Sparse and Irregular Data

*All existing FMs assume fully observed data*

Training

Inference

Context Tokens

| 2400 | ... | 2142 | ... | 2282 | ... | 2245 | ... | 2310 |

Context Tokens

| 2400 | ... | 2142 | ... | 2282 | ... | 2245 | ... | 2310 |

Time Series Language Model

| 2350 | ... | 2283 | ... | 2320 |

Probabilistic Forecast

## Multivariate and multi-depth

*Multivariate FMs treat each depth as independent variable*

$\mathbb{R}^{1\times T}$  $\{0,1\}^{1\times T}$  $\mathbb{R}^{P\times N}$  $\mathbb{R}^{D\times N}$  $\mathbb{R}^{D\times N}$

Masking    Patching    Encoding    Reconstruction

L ✖

MLP

Norm

Multi-Head Attention

Norm

Transformer Encoder

Transformer Encoder

Reconstruction Head

Mixture Distribution

Multi Patch Size Output Projection — Patch Size 8 | Patch Size 16 | Patch Size 32

Transformer (Full Self-Attention)

Time ID    0    1    2    0    1

Variate ID    0    0    0    1    1

Patch Embedding    [mask]

Multi Patch Size Input Projection — Patch Size 8 | Patch Size 16 | Patch Size 32

Variate 0    Variate 1    Variate 2

Depth

Variate

Time

13

# Lake Foundation Model (LakeFM) - An Overview

✓ **Variable Context and Prediction Length**

✓ **Variate-wise Distribution**

✓ **Any-variate any-depth prediction**

✓ **Variable Context and Prediction Length**

✓ **Variate-wise Distribution**

✓ **Any-variate any-depth prediction**

# Lake Foundation Model (LakeFM) - Tokenization & Embedding

$[ (T_k, D_k, V_k, Value_k), ... ]$

✓ **Irregular Grid in depth, time & variates with missing values**

# Lake Foundation Model (LakeFM) - Findings (I) - Increasing Horizon Length

# Lake Foundation Model (LakeFM) - Findings (I) - Increasing Horizon Length



*Key Observation : Overall, LakeFM maintains stable performance across increasing horizon lengths*

# Lake Foundation Model (LakeFM) - Findings (II) - Incomplete Data (Variables)



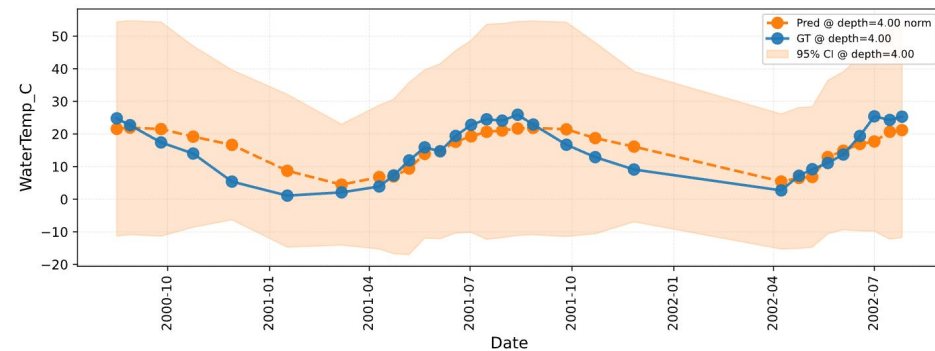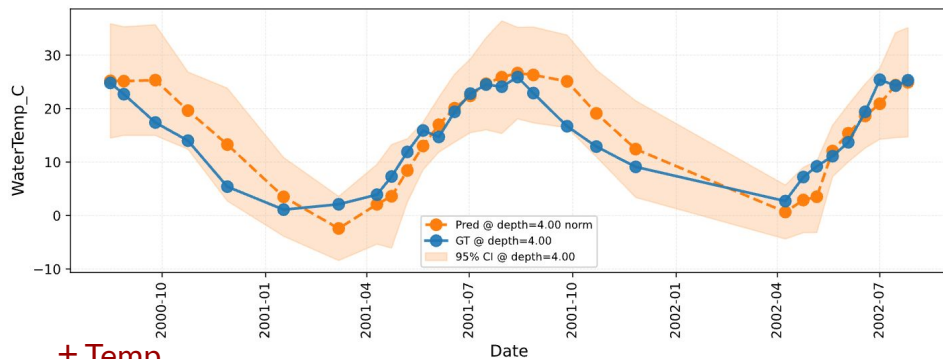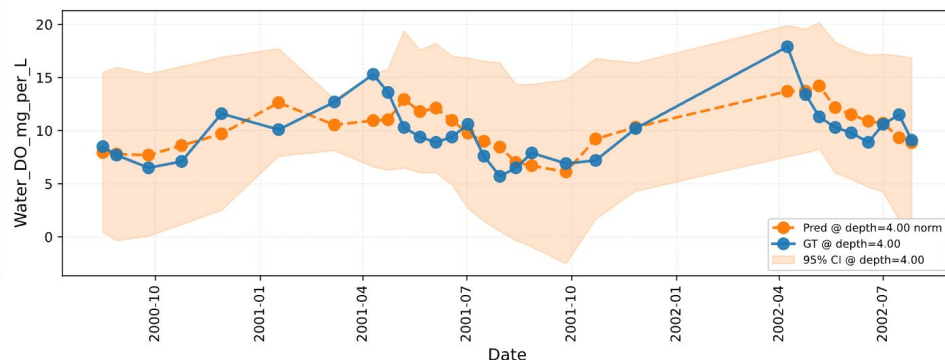ME : 30 timesteps ahead forecast, at Depth 4.00m (shaded = 95% CI)

+ DO

- DO

*Key Observation : Removing DO from inputs increases uncertainty in predictions of water temperature*

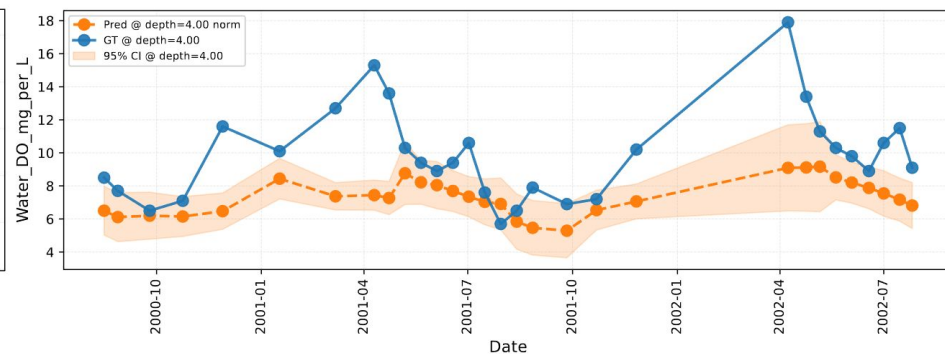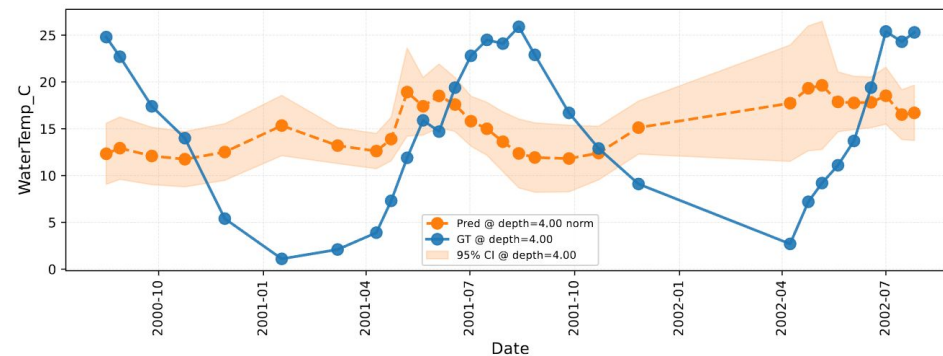# Lake Foundation Model (LakeFM) - Findings (II) - Incomplete Data (Variables)



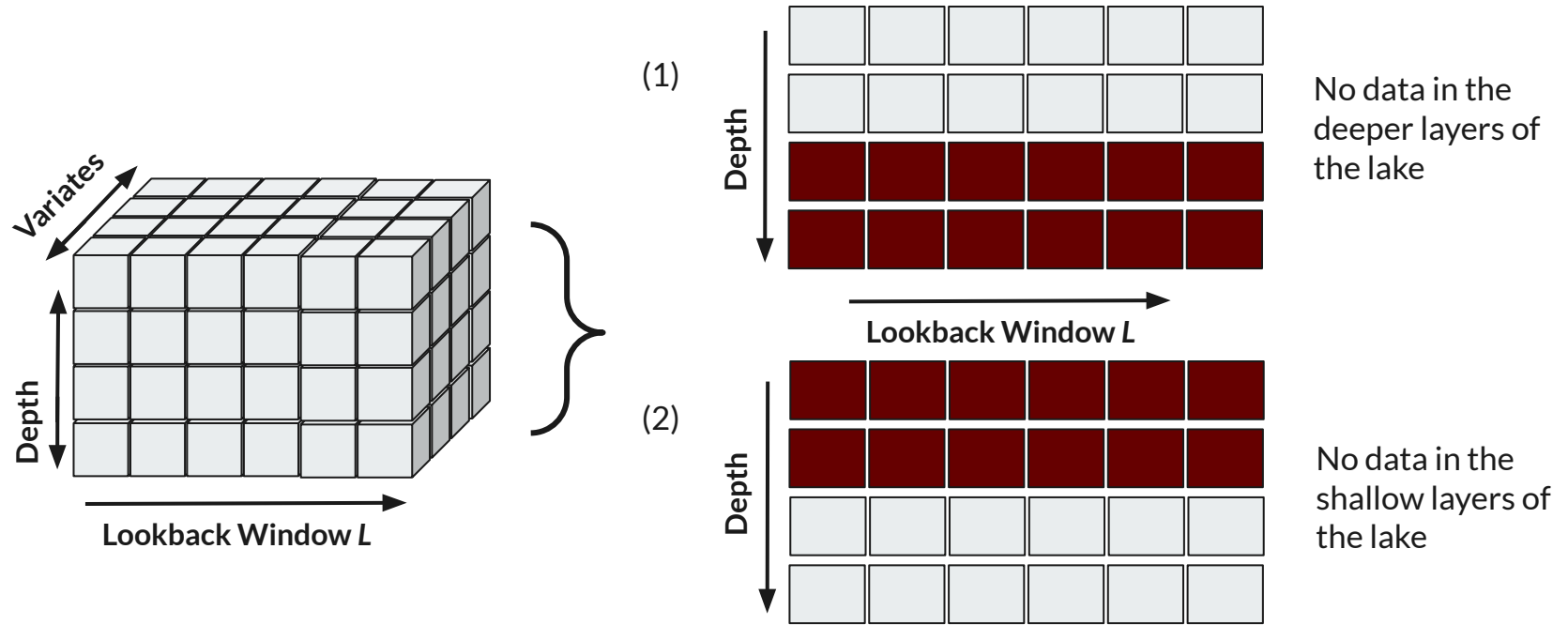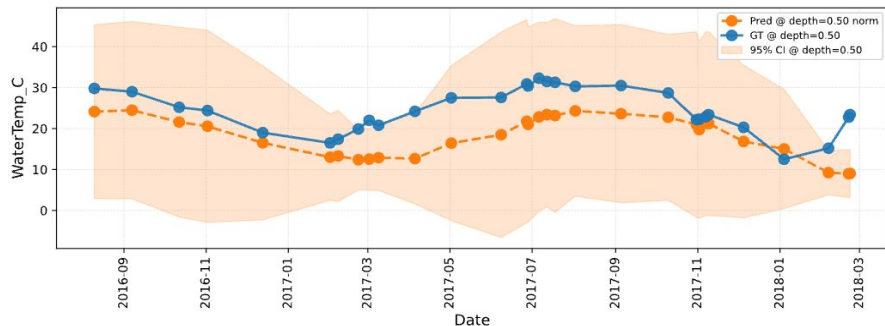ME : 30 timesteps ahead forecast, at Depth 4.00m  (shaded = 95% CI)

+ Temp

- Temp

*Key Observation : Water temperature is a critical variable. Removing it degrades all predictions.*

(1) No data in the deeper layers of the lake

(2) No data in the shallow layers of the lake

**Water Temp Predictions @ 0.5 m using full-depth history**

**Water Temp Predictions @ 0.5 m using only deeper-depth history**

**Water Temp Predictions @ 2.0 m using full-depth history**

**Water Temp Predictions @ 2.0 m using only deeper-depth history**

*Key Observation : Water temperature predictions remain stable even without shallow-layer variables*

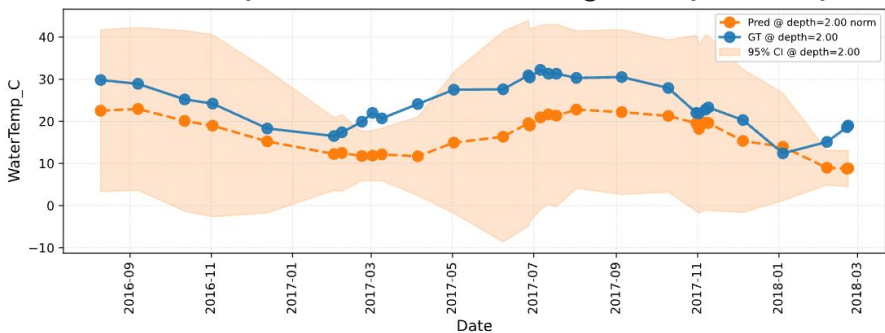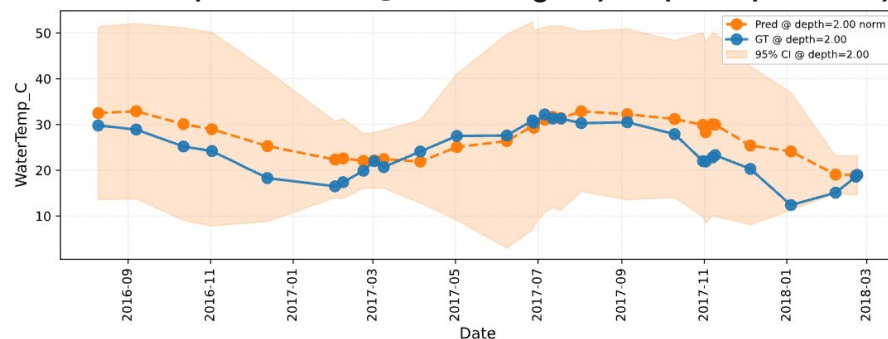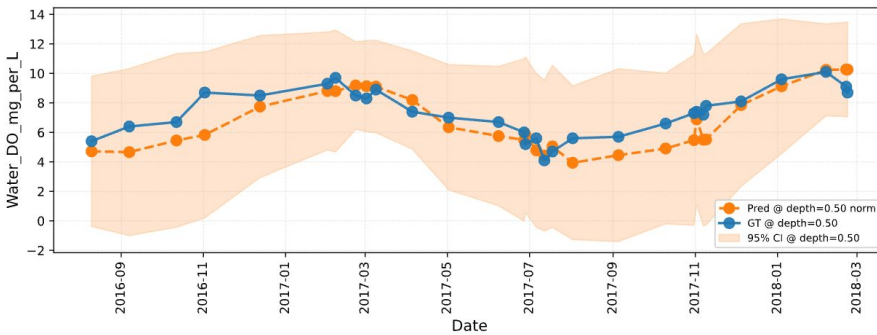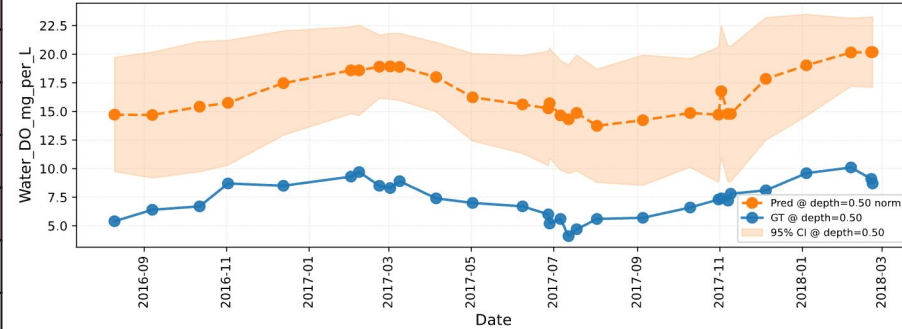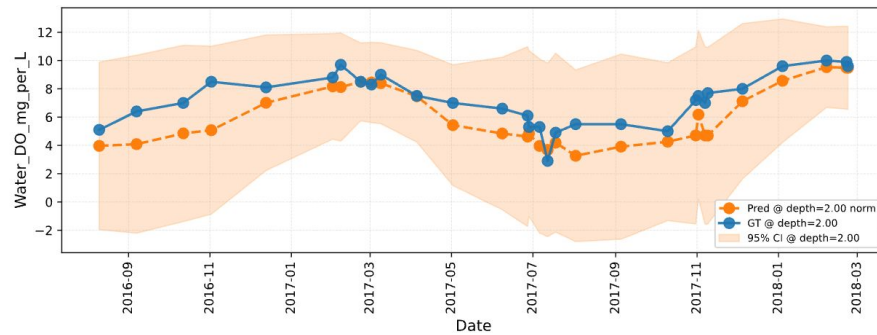# Lake Foundation Model (LakeFM) - Findings (III) - Incomplete Data (Depth)



*Key Observation : In contrast, DO predictions cannot rely on deeper-layer variables, indicating stronger vertical variability along the water column*

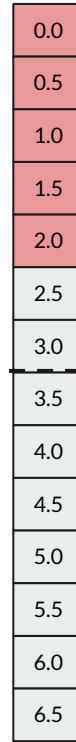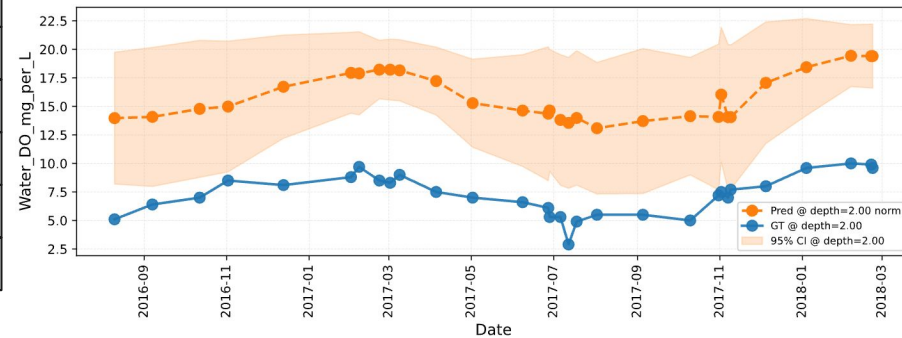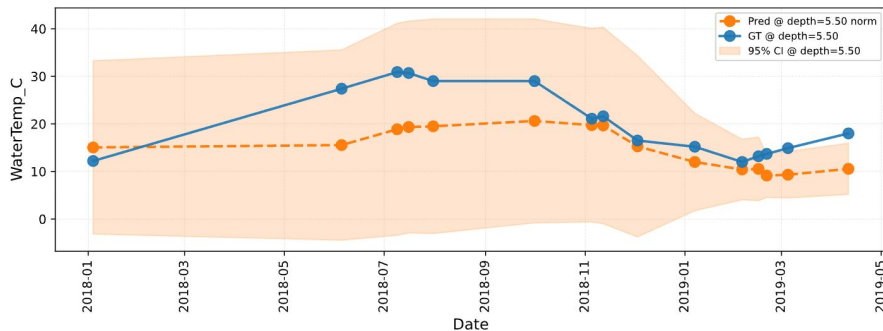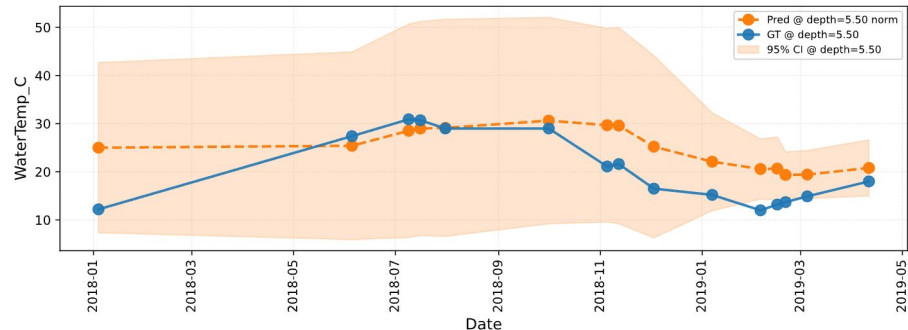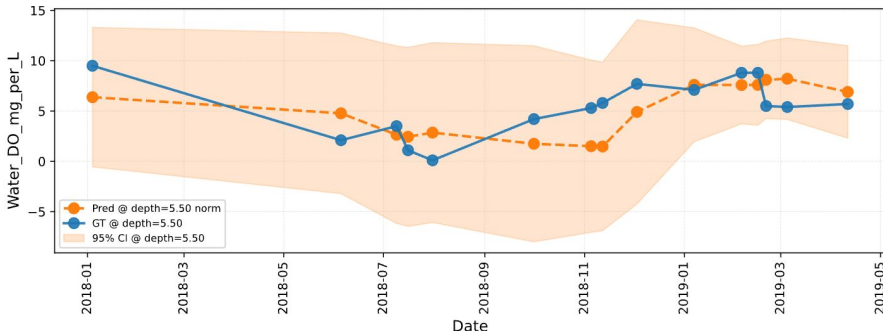Water Temp Predictions @ 5.5 m using full-depth history
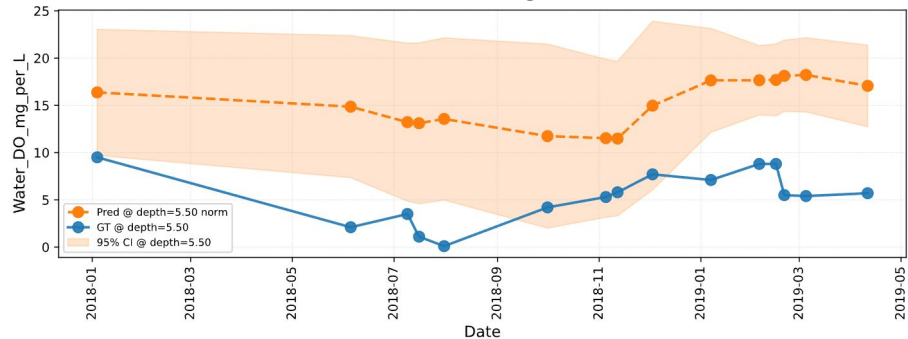
Water Temp Predictions @ 5.5 m using only shallow-depth history

Water DO Predictions @ 5.5 m using full-depth history
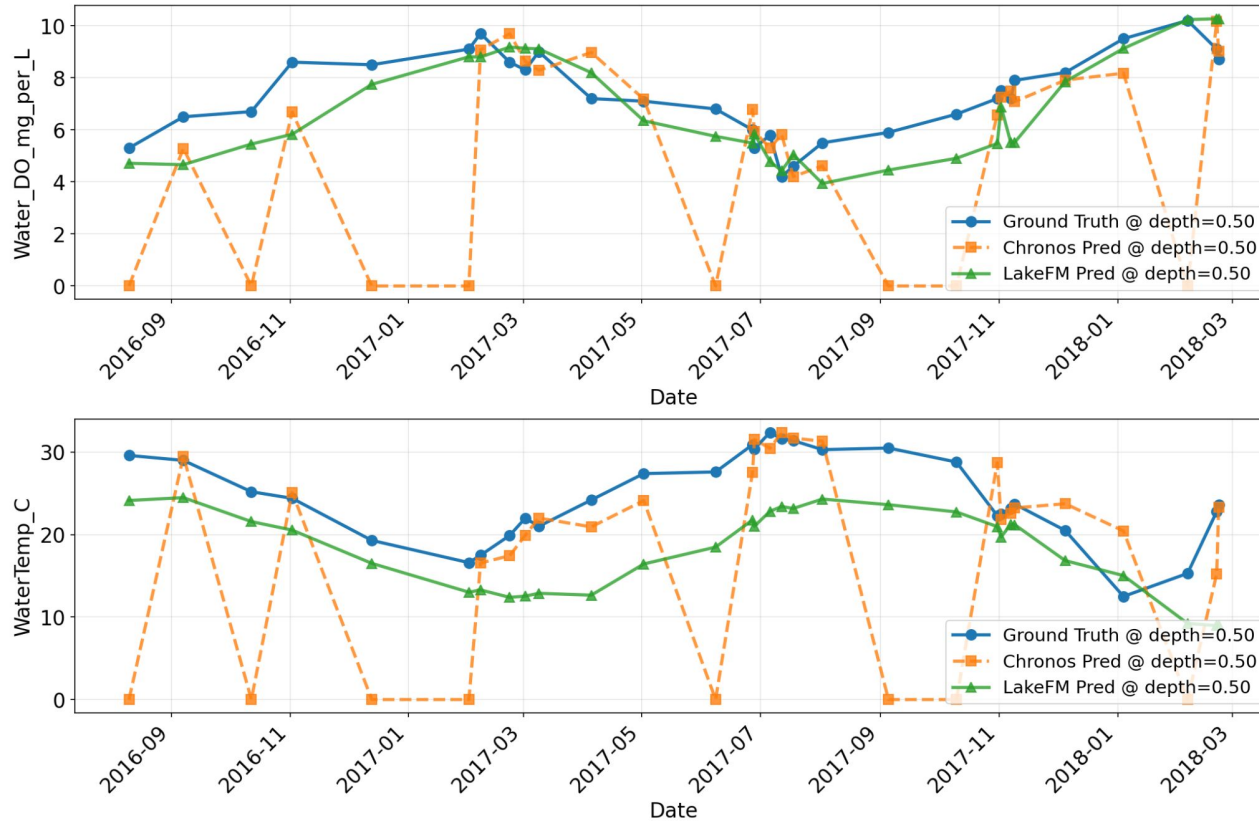
Water DO Predictions @ 5.5 m using only shallow-depth history

*Key Observation : In general, water temperature remains stable and is predictable using either shallow or deeper layers, while DO dynamics are tightly coupled to the local depth.*

31

# Lake Foundation Model (LakeFM) - Findings (IV) - Performance Comparison

**Comparing LakeFM predictions (on a horizon window of 30 timesteps) with Chronos Foundation Model on Lake BARC at Depth 0.5m**



*Key Observation :*

*Chronos struggles with missing data, leading to context-dependent instability, whereas LakeFM remains stable under the same conditions*

# Lake Foundation Model (LakeFM) - Ongoing Work

**Lake representation analysis**

- ❖ Visualization of learned lake embeddings

- ❖ Analyzing seasonal clustering patterns in lake representations

- ❖ Analyzing temporal trajectories of lake representations over time

- ❖ *Geographic structure emerging in embedding space*

**Variable representation analysis**

- ❖ Visualization of learned variable embeddings

- ❖ Analyzing variable similarities inferred from embedding clusters

- ❖ Empirical verification of embedding-based variable similarities

Thank you