

# Some Preliminary Transfer Learning Results for Dynamical Systems

**Joshua Hammond, Brian Korgel, Michael Baldea**

McKetta Department of Chemical Engineering  
The University of Texas at Austin

**Acknowledgments: Tyler Soderstrom, ExxonMobil**

TWCCC Fall Meeting 2025

# Chemical Processes are “Copies” of Dynamical Systems



- Several plants are built off the same blueprint
- Similar but not identical
- May use a module (or more) in multiple designs

How can \*data driven\* modeling knowledge developed for one system be leveraged in the operation and control of another “copy”

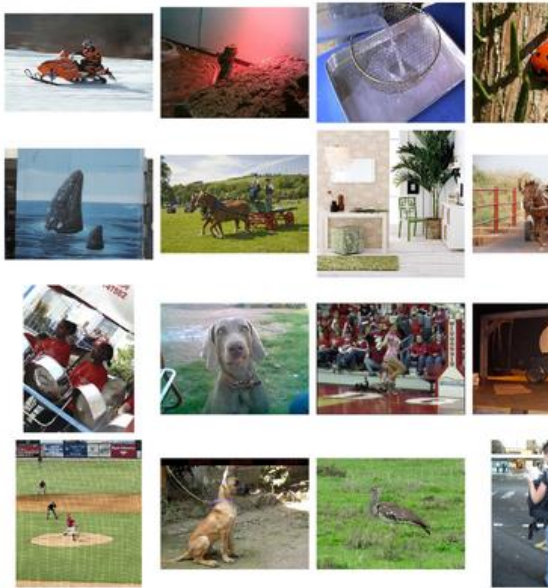
>> Transfer learning

M. Baldea, T.F. Edgar, B. Stanley, A.A. Kiss, Modular Manufacturing: Status, Challenges and Opportunities, AICHE Journal, 63(10), 4262-4272, 2017

# Most existing TL applications focus on image classification

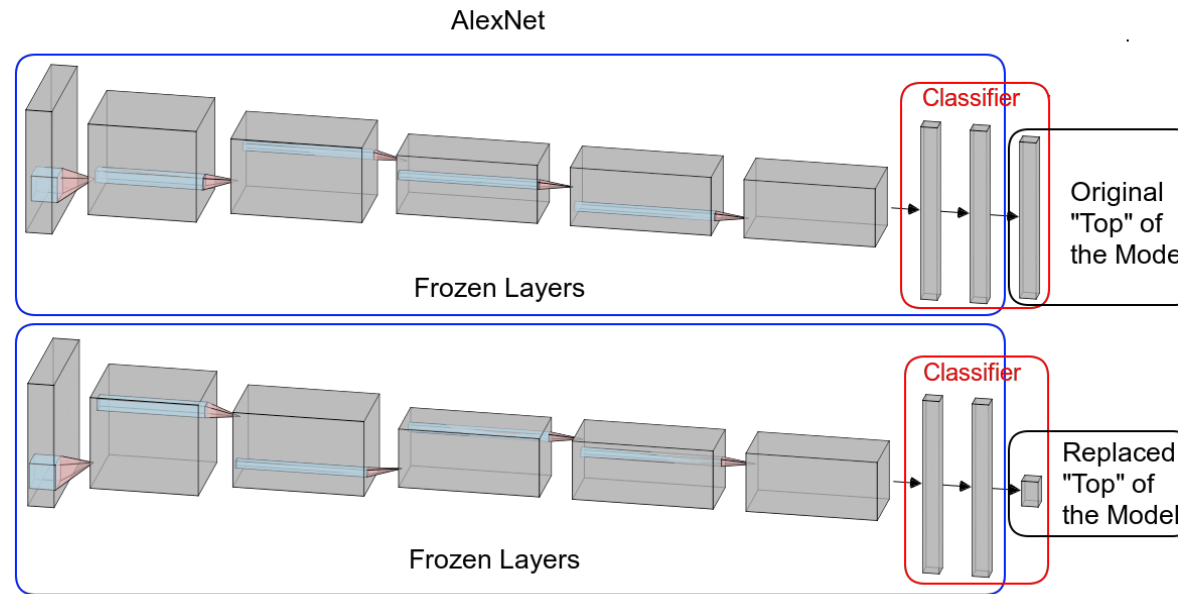
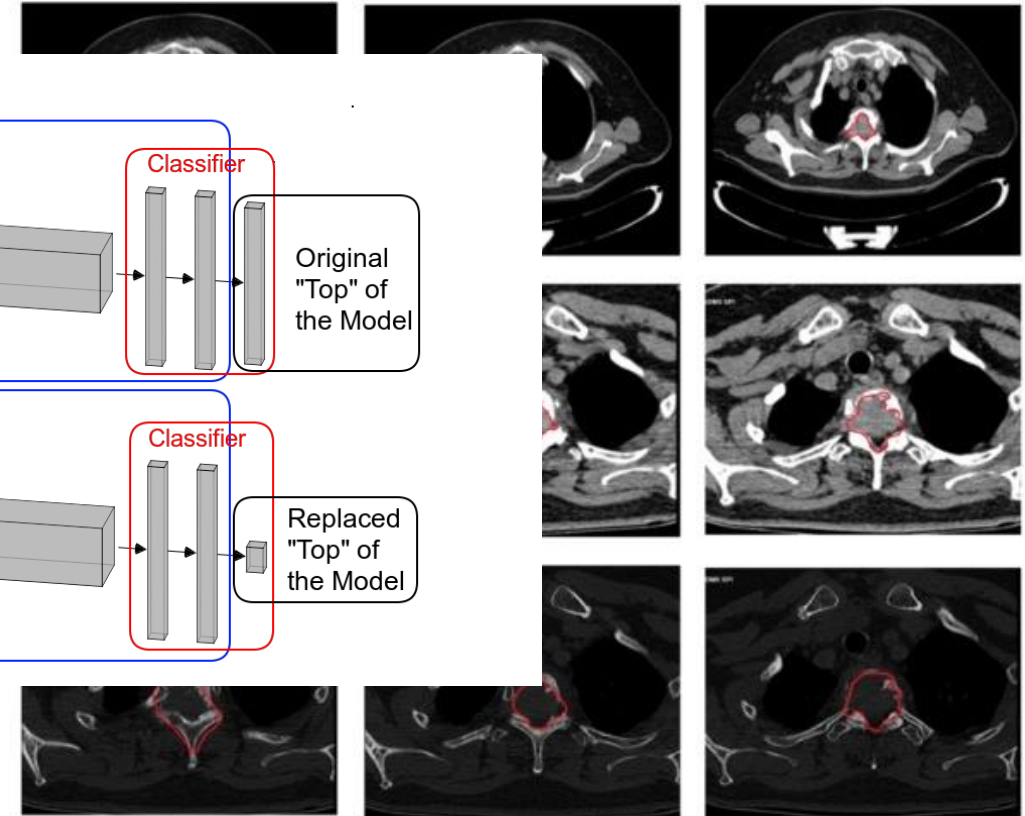
## Source: ImageNet Classification<sup>[1]</sup>

Image classification from 1,000 common objects



## Target: Medical Diagnostics<sup>[2]</sup>

Classification of tumor type



Malignant

**Not dynamical systems – no governing equation constraints**

[1] Russakovsky et. Al. "ImageNet Large Scale Visual Recognition Challenge" arXiv: 1409.0575 [2] Guo et. Al. "Radiographic imaging and diagnosis of spinal bone tumors: AlexNet and ResNet for the classification of tumor malignancy" 2024 Journal of Bone Oncology [3] AlexNet Visualization: Daniel Voight Godoy <https://github.com/dvgodoy/dl-visuals/>

# Machine learning models of dynamical systems

---

We consider a general physical system, whose *true* dynamics are given by

$$\frac{d}{dt}\mathbf{x}(t) = \mathbf{f}^S(\mathbf{x}, \mathbf{u}, \mathbf{p})$$

where the function  $\mathbf{f}^S: \mathbb{R}^{n_x \times n_u \times n_p} \rightarrow \mathbb{R}^{n_x}$  is the *governing equation(s)*, and is a function of the current state  $\mathbf{x}$ , inputs  $\mathbf{u}$ , and parameters  $\mathbf{p}$ .

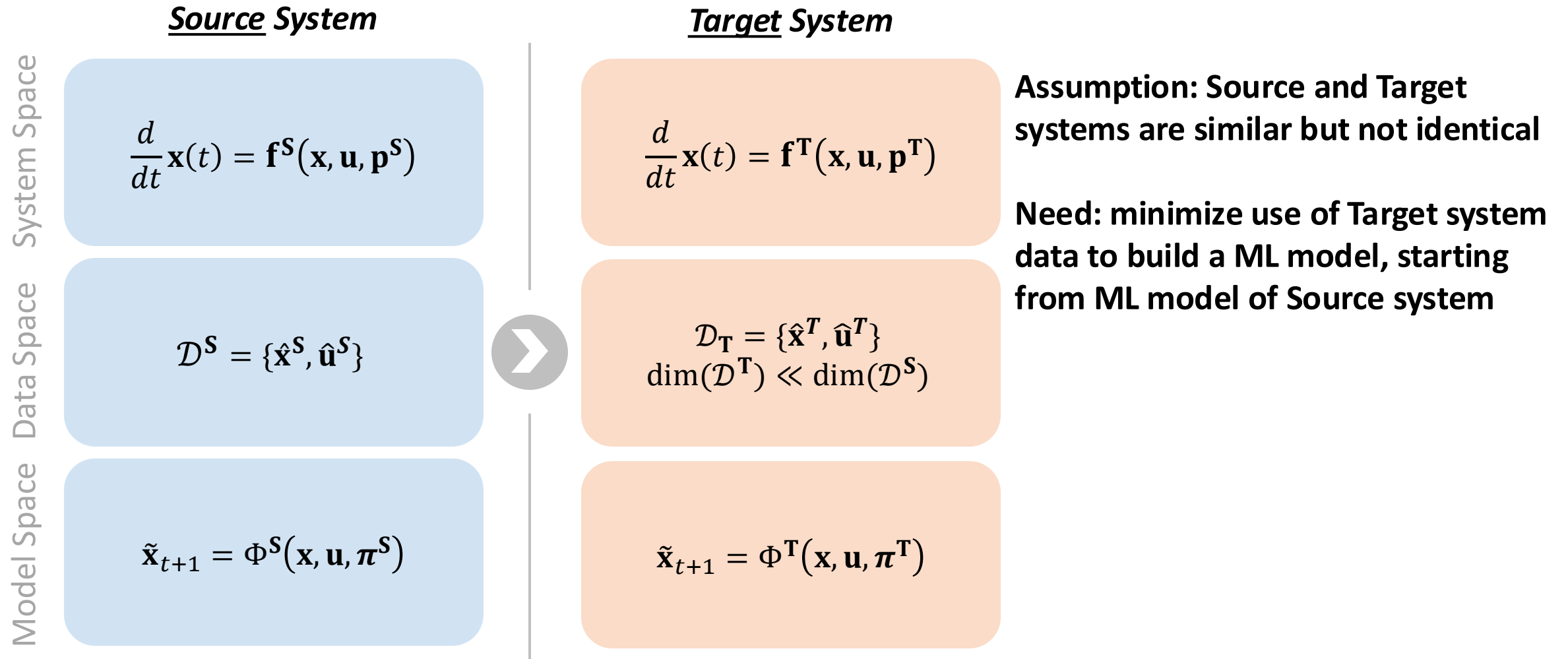
**NOTE:** The true form of the function  $\mathbf{f}$  is *unknown*, but data  $\mathcal{D}^S \triangleq \{\hat{\mathbf{x}}, \hat{\mathbf{u}}\}$  are available with noise that is assumed to be Gaussian.

A NN model produces discrete predictions of future states

$$\tilde{\mathbf{x}}_{t+1} = \Phi^S(\mathbf{x}_t, \mathbf{u}_t, \boldsymbol{\pi})$$

where  $\boldsymbol{\pi}$  is a vector of the NN parameters that are fit using the data

# Research Problem Formulation



# Example: Damped Spring

$$m_s \ddot{x} + c_s \dot{x} + k_s x = u_s$$

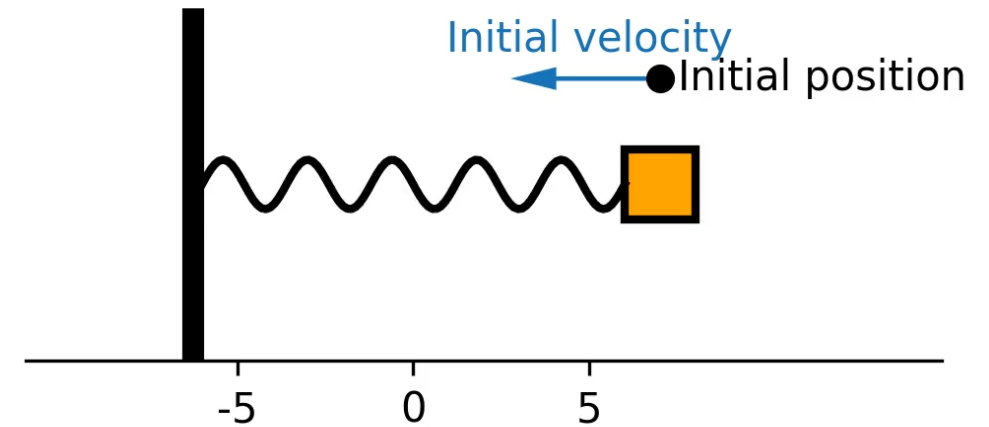
**Given:**  $x_0, \dot{x}_0$

**Predict:**  $x_1, x_2, \dots, x_{20}$

**Source System**  $m_s \ddot{x} + c_s \dot{x} + k_s x = u_s$

Measurement Noise:  $\mathcal{N}(0, 0.1)$

$\dim(D_s) = 100,000$



Use a 2-layer MLP

- 32 neurons each layer
- Sigmoid activation
- 2,868 trainable parameters

# Example: Damped Spring

$$m_s \ddot{x} + c_s \dot{x} + k_s x = u_s$$

**Given:**  $x_0, \dot{x}_0$

**Predict:**  $x_1, x_2, \dots, x_{20}$

**Source System**  $m_s \ddot{x} + c_s \dot{x} + k_s x = u_s$

Measurement Noise:  $\mathcal{N}(0, 0.1)$

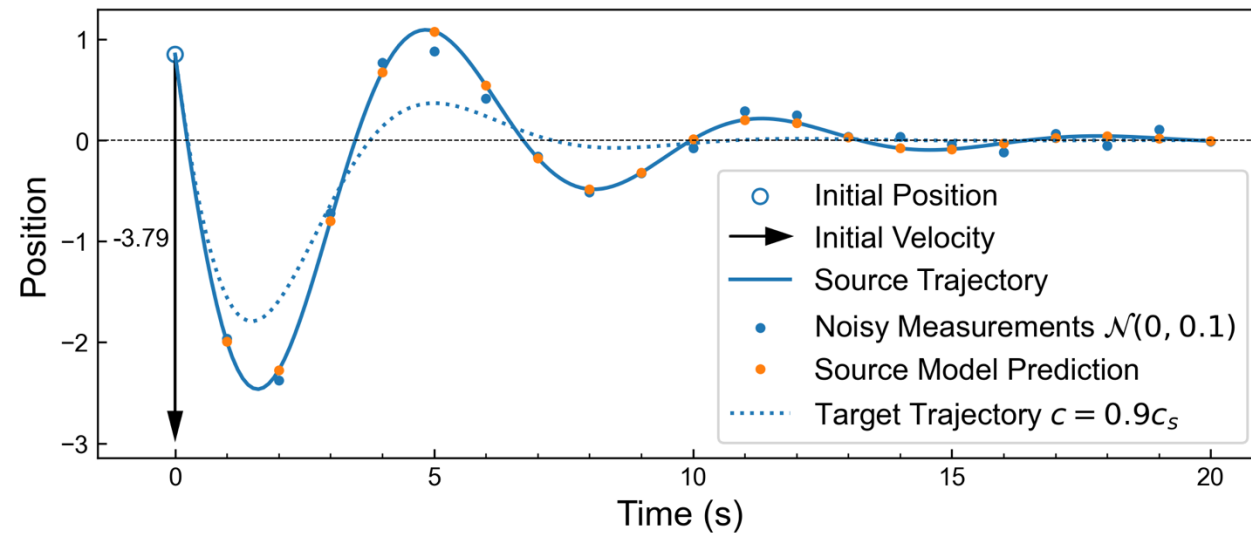
$\dim(D_S) = 100,000$

**Target System (set of systems)**

Change  $m_s, c_s, k_s, u_s, \pm \Delta 10\%$

$\dim(D_T) = 1,000$

An additional 99,000 data points are used for evaluation

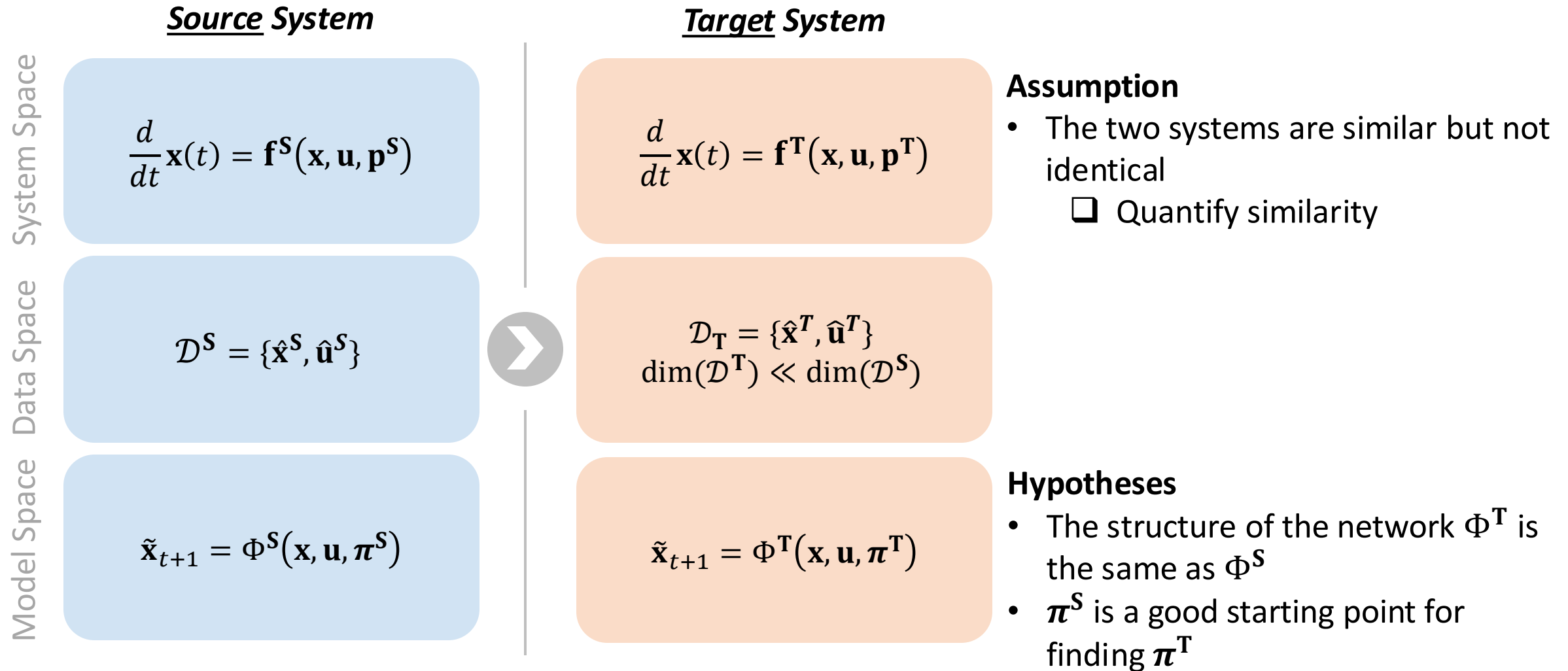


Use a 2-layer MLP

- 32 neurons each layer
- Sigmoid activation
- 2,868 trainable parameters

Can the NN model of the Source oscillator be easily adapted to the Target oscillator?

# Transfer Learning in Dynamical Systems





# Cosine Similarity

Given two vectors  $a$  and  $b$ :

$$\theta(a, b) \triangleq \frac{a \cdot b}{|a||b|}$$

$\theta = 1$  Perfectly correlated

$\theta = 0$  Uncorrelated

$\theta = -1$  Negatively correlated

Some benefits:

- Scale, translation independent
- Robust to different lengths, high dimensionality
- Interpretability

## Data Similarity

Done on a variable by variable basis

## Functional Similarity

For a given set of inputs  $(\mathbf{x}, \mathbf{u})$  to two dynamical systems

$$\begin{aligned} f^S(\mathbf{x}, \mathbf{u}, \mathbf{p}^S) \\ f^T(\mathbf{x}, \mathbf{u}, \mathbf{p}^T) \end{aligned}$$

Simulate/sample in discrete time the input and output trajectories and then evaluate similarity using the data

# NN Training

## Prediction

$$\tilde{\mathbf{y}}_k = \Phi(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k) + \mathbf{v}_k$$

## Gradient Descent

$$\mathbf{e}_k = \mathbf{y}_k - \tilde{\mathbf{y}}_k$$

$$\mathcal{L}_k = \text{Loss Function}(\mathbf{e}_k)$$

$$\frac{d\mathcal{L}}{d\hat{\boldsymbol{\pi}}_k} \leftarrow \text{Backpropagation}(\mathcal{L}, \hat{\boldsymbol{\pi}}_k)$$

$$\hat{\boldsymbol{\pi}}_{k+1} = \hat{\boldsymbol{\pi}}_k - \eta \frac{d\mathcal{L}}{d\hat{\boldsymbol{\pi}}_k}$$

Repeat across all datapoints 1 ...  $k$  until convergence

$k$  – Discrete time index

$\boldsymbol{\pi}_k$  – Neural network parameters ( $N_p, 1$ )

$\hat{\boldsymbol{\pi}}_k$  – NN parameter estimate ( $N_p, 1$ )

$\mathbf{y}_k$  – Discrete measurement ( $N_{out}, 1$ )

$\tilde{\mathbf{y}}_k$  – NN output ( $N_{out}, 1$ )

$\Phi$  – NN model

$\mathbf{w}_k$  – Model uncertainty

$\mathbf{v}_k$  – Measurement uncertainty

$\mathbf{j}$  – Selected parameter indices

$\mathbf{e}_k$  – Error vector ( $N_{out}, 1$ )

$\mathbf{H}_k$  – Gradients of network parameters

w.r.t. outputs (Jacobian) ( $N_{params}, N_{out}$ )

$\mathbf{P}_k$  – Covariance matrix ( $N_{params}, N_{params}$ )

$\mathbf{Q}$  – Process noise matrix ( $N_{params},$

$N_{params}$ )  $E[\mathbf{w}_k \mathbf{w}_{k+1}] = \delta_{k,k+1} \mathbf{Q}_k$

$\mathbf{R}_k$  – Measurement Noise ( $N_{out}, N_{out}$ )

$E[\mathbf{v}_k \mathbf{v}_{k+1}] = \delta_{k,k+1} \mathbf{R}_k$

# Strategies for updating NN models

**Idea:** If models are similar, and their parameters will be close, we can use model-updating techniques to evolve model parameters to reflect the target system

## The Subset Extended Kalman Filter (SEKF)

- Uses well-established Kalman Filtering techniques to update a subset of model parameters
- Probabilistic methods for state-estimation in the presence of noise ideally limit spurious correlations and overfitting

$$\mathbf{e}_k = \mathbf{y}_k - \tilde{\mathbf{y}}_k$$

$$\mathbf{H}_k = \nabla_{\pi} \mathbf{y}_k$$

$$\mathbf{j} = \{\text{Prop. } q, \text{Mag. } q\}$$

$$\mathbf{H}'_k = \mathbf{H}_k[:, \mathbf{j}]; \mathbf{P}'_k = \mathbf{P}_k[\mathbf{j}, \mathbf{j}]$$

**Parameter selection**

$$\mathbf{A}_k = [\mathbf{R}_k + \mathbf{H}'_k{}^T \mathbf{P}'_k \mathbf{H}'_k]^{-1}$$

$$\mathbf{K}_k = \mathbf{P}'_k \mathbf{H}'_k \mathbf{A}_k$$

$$\hat{\pi}[\mathbf{j}]_{k+1} = \hat{\pi}[\mathbf{j}]_k + \mathbf{K}_k \mathbf{e}_k$$

$$\mathbf{P}[\mathbf{j}, \mathbf{j}]_{k+1} = \mathbf{P}'_k - \mathbf{K}_k \mathbf{H}'_k{}^T \mathbf{P}'_k + \mathbf{Q}'$$

$k$  – Discrete time index

$\pi_k$  – Neural network parameters ( $N_{p,1}$ )

$\hat{\pi}_k$  – NN parameter estimate ( $N_{p,1}$ )

$\mathbf{y}_k$  – Discrete measurement ( $N_{out,1}$ )

$\tilde{\mathbf{y}}_k$  – NN output ( $N_{out,1}$ )

$\Phi$  – NN model

$\mathbf{w}_k$  – Model uncertainty

$\mathbf{v}_k$  – Measurement uncertainty

$\mathbf{j}$  – Selected parameter indices

$\mathbf{e}_k$  – Error vector ( $N_{out,1}$ )

$\mathbf{H}_k$  – Gradients of network parameters w.r.t. outputs (Jacobian) ( $N_{params}, N_{out}$ )

$\mathbf{P}_k$  – Covariance matrix ( $N_{params}, N_{params}$ )

$\mathbf{Q}$  – Process noise matrix ( $N_{params}, N_{params}$ )

$E[\mathbf{w}_k \mathbf{w}_{k+1}] = \delta_{k,k+1} \mathbf{Q}_k$

$\mathbf{R}_k$  – Measurement Noise ( $N_{out}, N_{out}$ )  $E[\mathbf{v}_k \mathbf{v}_{k+1}] =$

$\delta_{k,k+1} \mathbf{R}_k$

# Transfer Learning for Damped Spring

$$m_s \ddot{x} + c_s \dot{x} + k_s x = u_s$$

Given:  $x_0, \dot{x}_0$

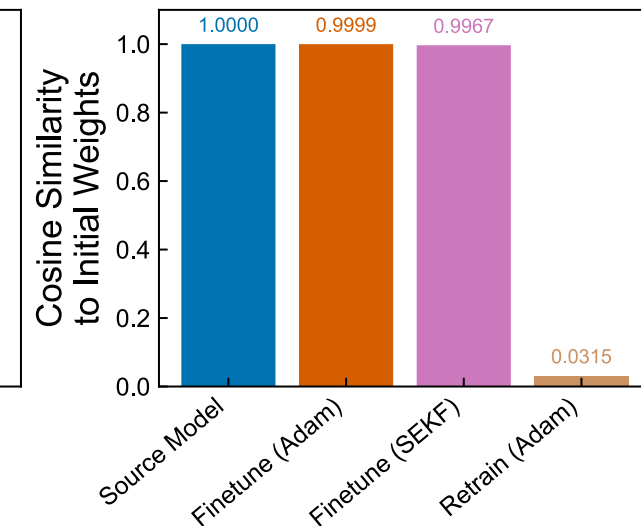
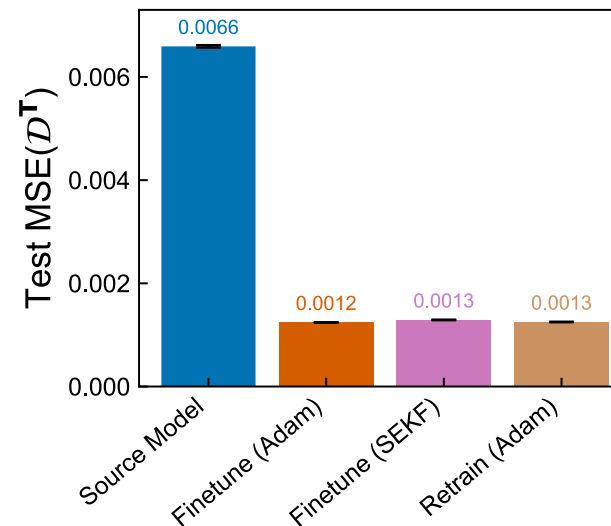
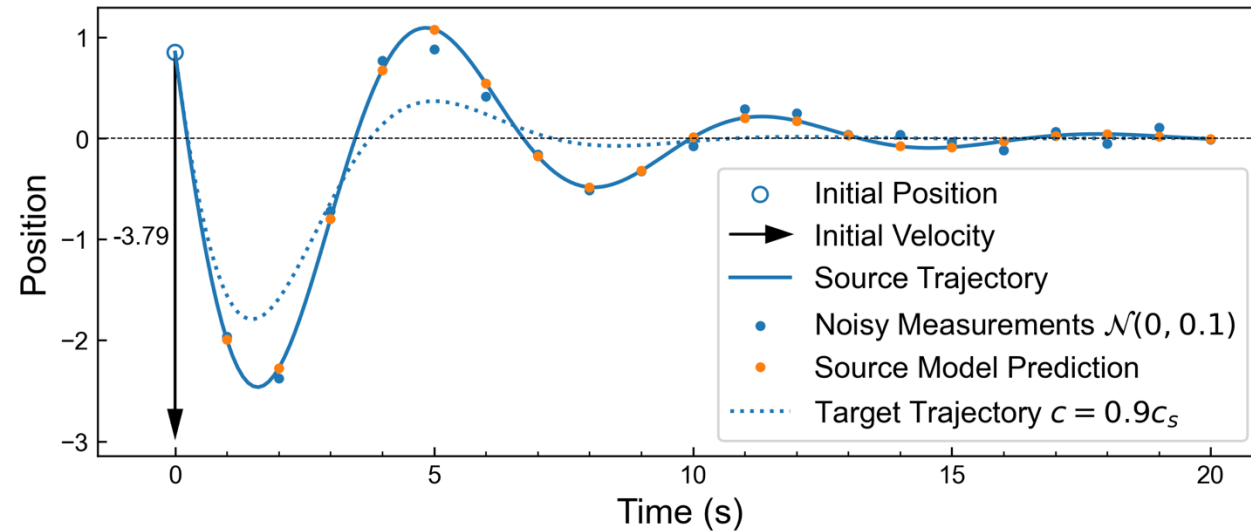
Predict:  $x_1, x_2, \dots, x_{20}$

## Competing Techniques:

- Apply source model
- Finetune (train NN starting from source model parameters) using Adam
- Finetune using SEKF
- Retrain (reinitialize NN weights) using ADAM

Use a 2-layer MLP

- 32 neurons each layer
- Sigmoid activation
- 2,868 trainable parameters



# Transfer Learning for Damped Spring

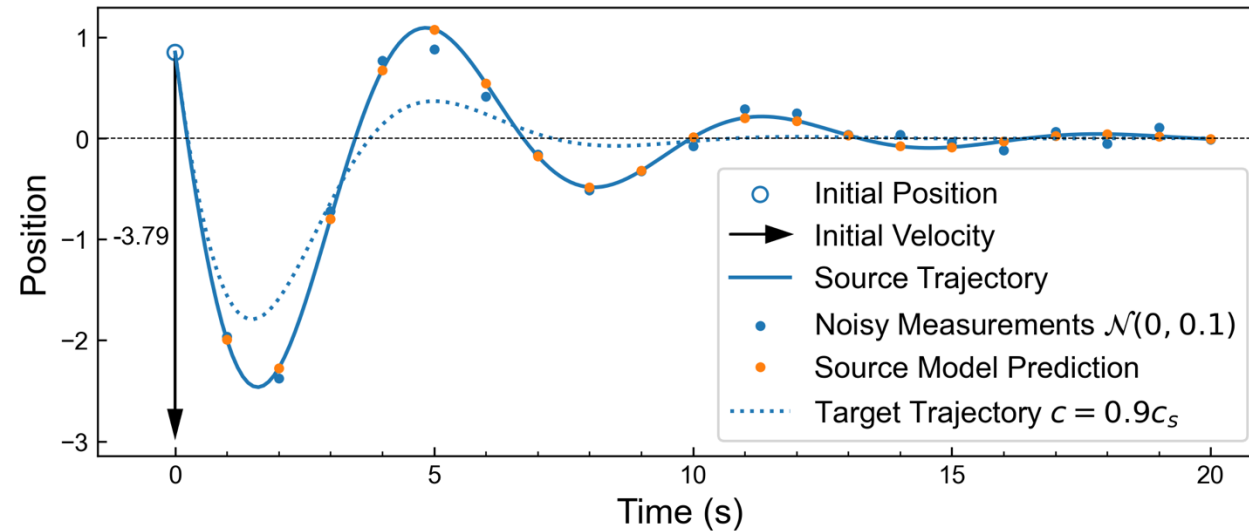
$$m_s \ddot{x} + c_s \dot{x} + k_s x = u_s$$

**Given:**  $x_0, \dot{x}_0$

**Predict:**  $x_1, x_2, \dots, x_{20}$

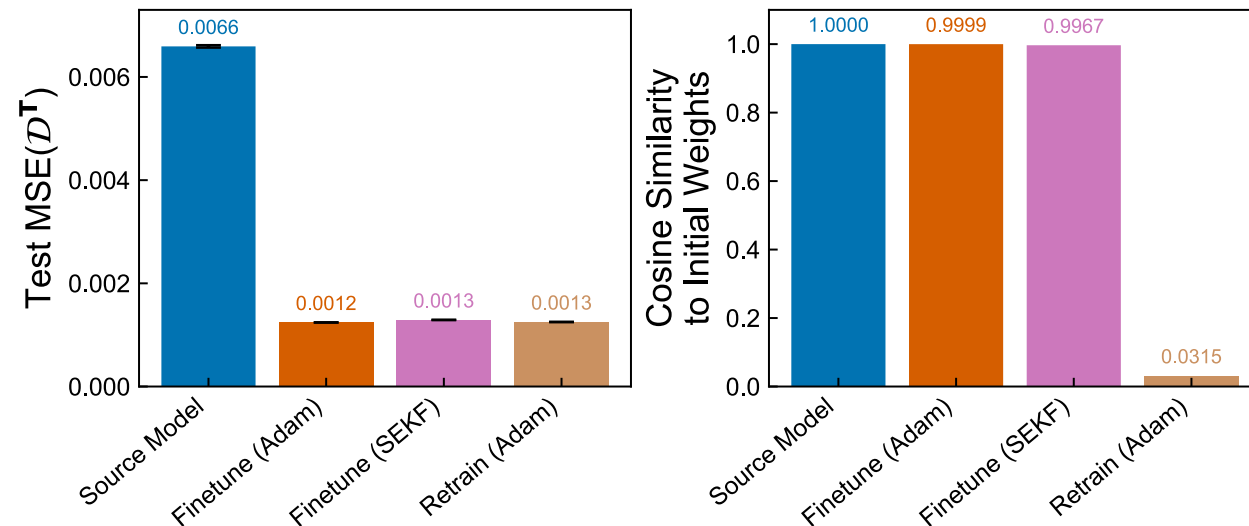
## Competing Techniques:

- Apply source model
- Finetune (train NN starting from source model parameters) using Adam
- Finetune using SEKF
- Retrain (reinitialize NN weights) using ADAM

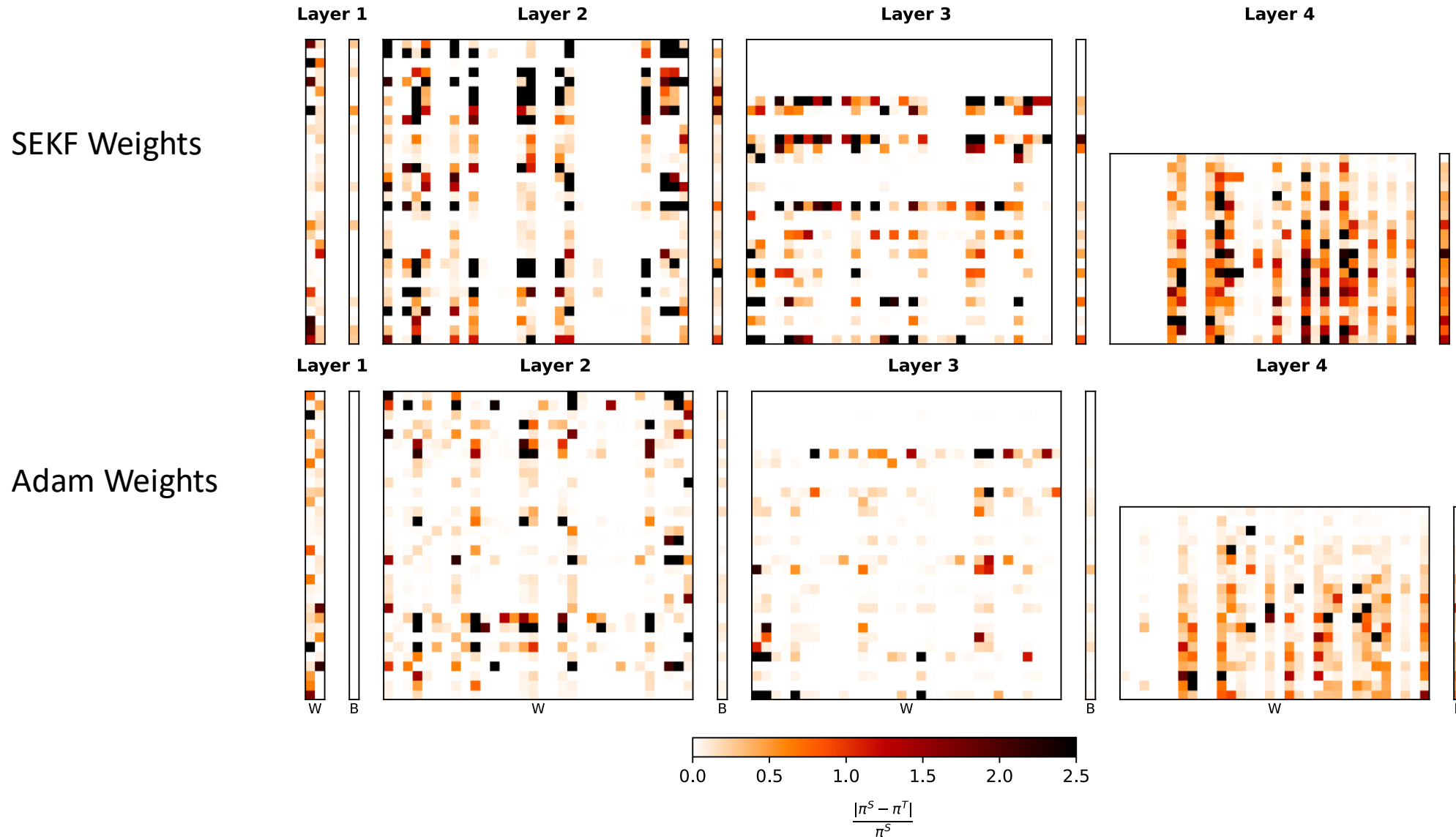


## Preliminary findings

- Finetuning results in similar model parameters between source and target
- Retraining may change NN model parameters substantially
- Loss is comparable



# Changes to the NN parameters are distributed through the model

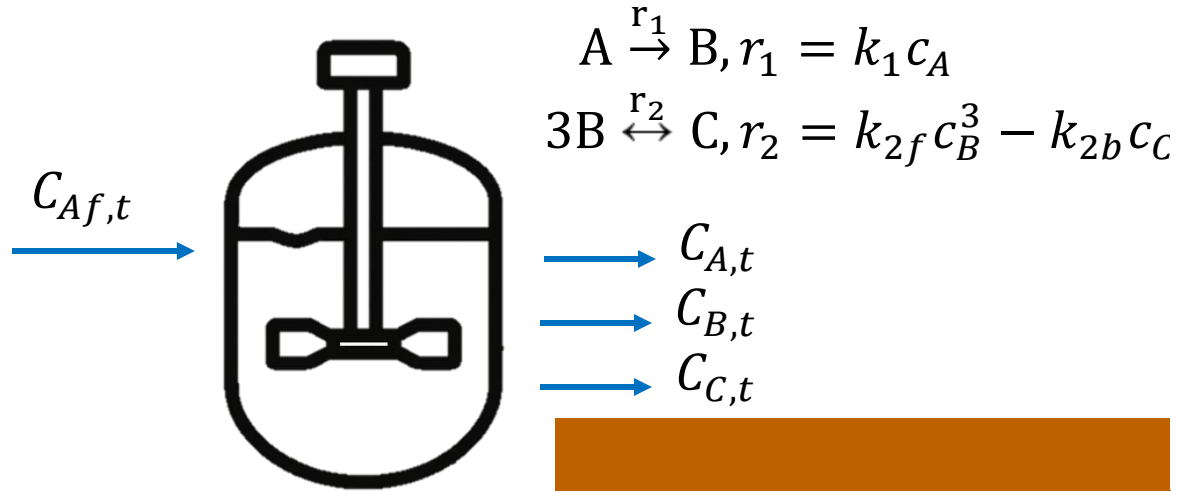


# Example: Official TWCCC CSTR 😊

Given:  $C_{A0}, C_{B0}, C_{C0}, U_{0,...,59}$

Predict:  $C_{i1}, \dots, C_{i,60}$

Source → Target: Reactor Volume  $\Delta 20\%$

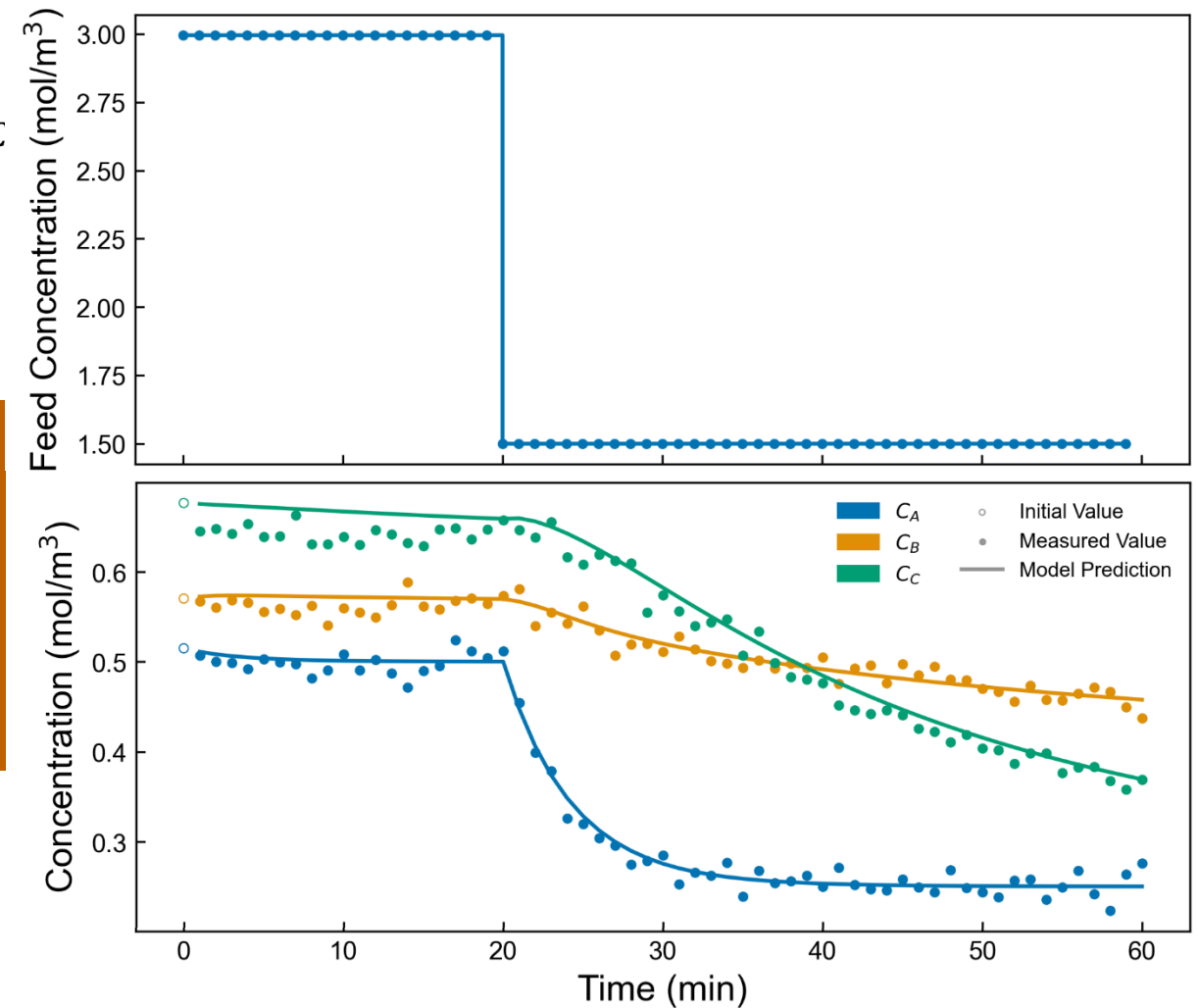


## Source System

$\dim(D_S) = 1$  year of data

Use a 2-layer Neural ODE

- 16 neurons each layer
- Sigmoid activation
- 419 trainable params



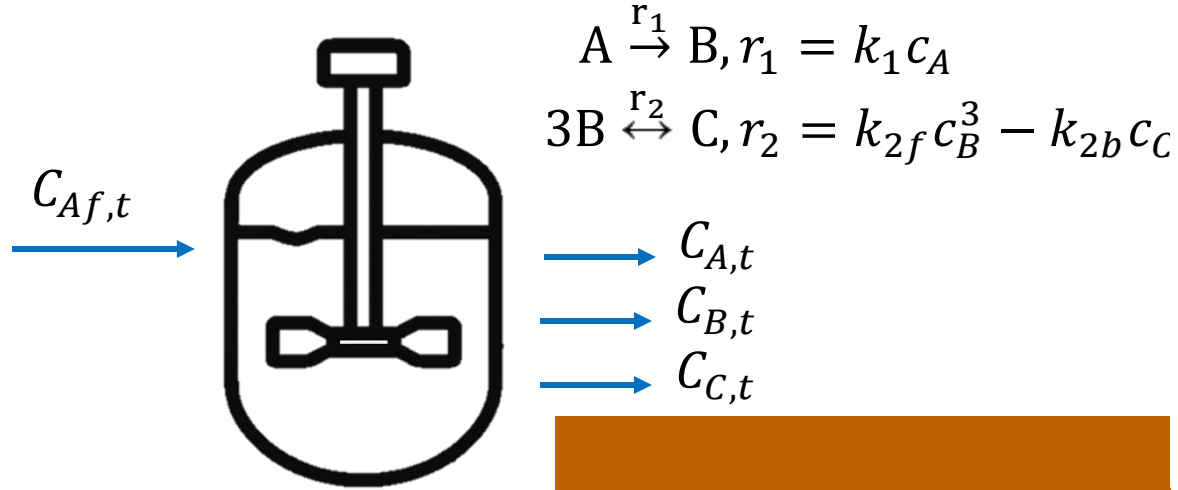
[1] Kumar, P., & Rawlings, J. B. (2023). Structured nonlinear process modeling using neural networks and application to economic optimization. *Computers and Chemical Engineering*, 177. <https://doi.org/10.1016/j.compchemeng.2023.108314>

# Example: Official TWCCC CSTR 😊

Given:  $C_{A0}, C_{B0}, C_{C0}, U_{0,...,59}$

Predict:  $C_{i1}, \dots, C_{i,60}$

Source → Target: Reactor Volume  $\Delta 20\%$



## Source System

$\dim(D_S) = 1$  year of data

## Target System

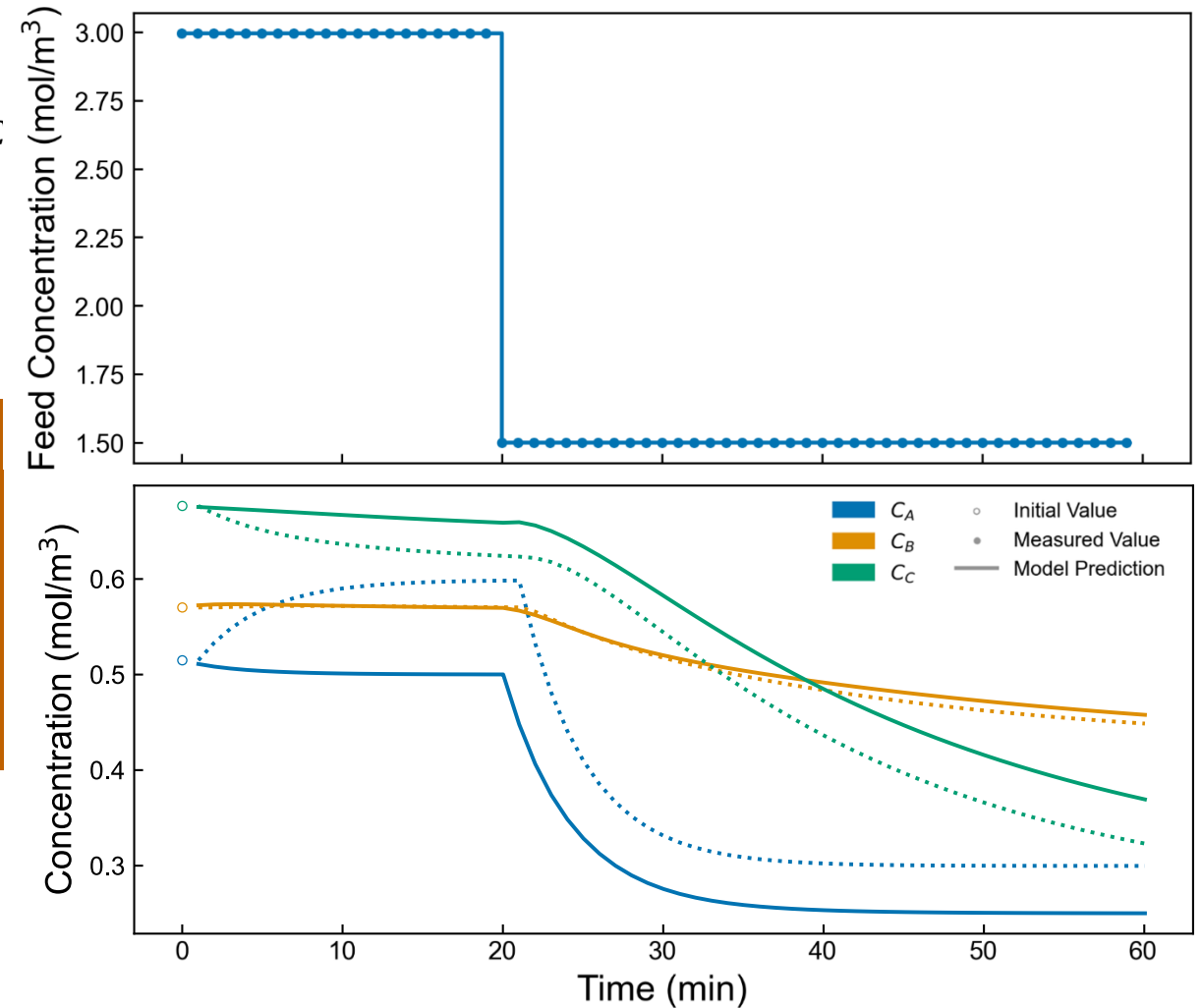
Reactor Volume  $\Delta 20\%$

$\dim(D_T) =$  One day of training data

Evaluated on remaining 363 days

Use a 2-layer Neural ODE

- 16 neurons each layer
- Sigmoid activation
- 419 trainable params



[1] Kumar, P., & Rawlings, J. B. (2023). Structured nonlinear process modeling using neural networks and application to economic optimization. *Computers and Chemical Engineering*, 177. <https://doi.org/10.1016/j.compchemeng.2023.108314>

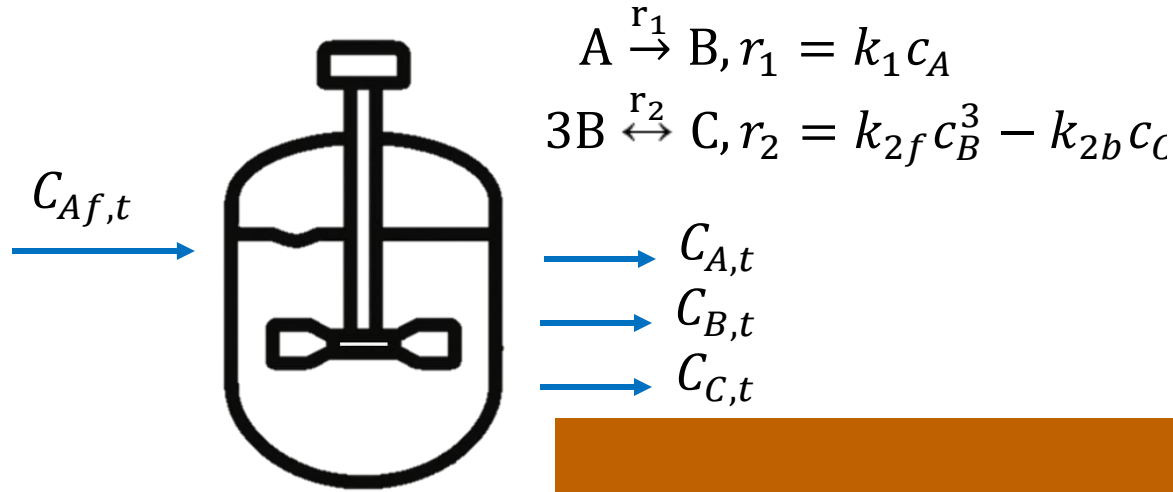


# Example: Official TWCCC CSTR 😊

Given:  $C_{A0}, C_{B0}, C_{C0}, U_{0,...,59}$

Predict:  $C_{i1}, \dots, C_{i,60}$

Source → Target: Reactor Volume  $\Delta 20\%$



## Source System

$\dim(D_S) = 1$  year of data

## Target System

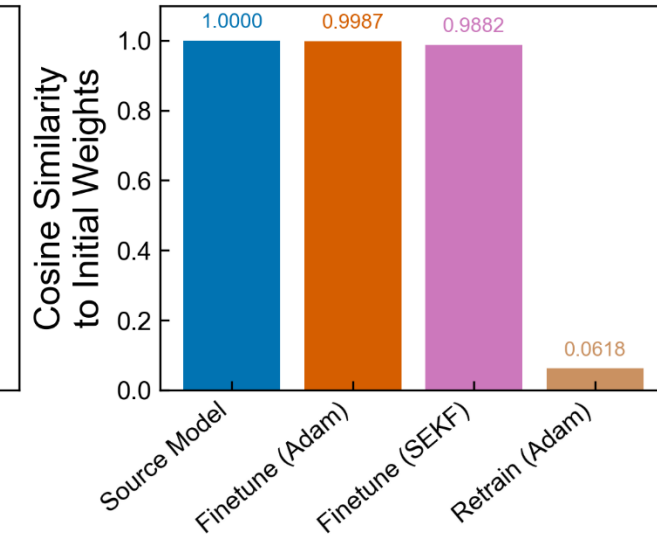
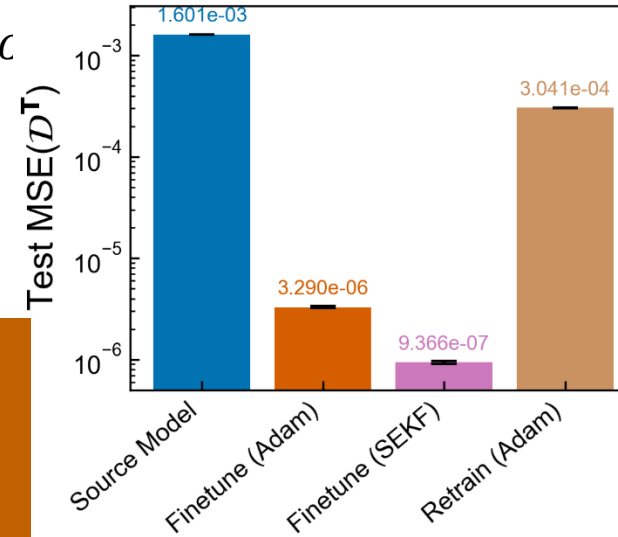
Reactor Volume  $\Delta 20\%$

$\dim(D_T) =$  One day of training data

Evaluated on remaining 363 days

Use a 2-layer Neural ODE

- 16 neurons each layer
- Sigmoid activation
- 419 trainable params



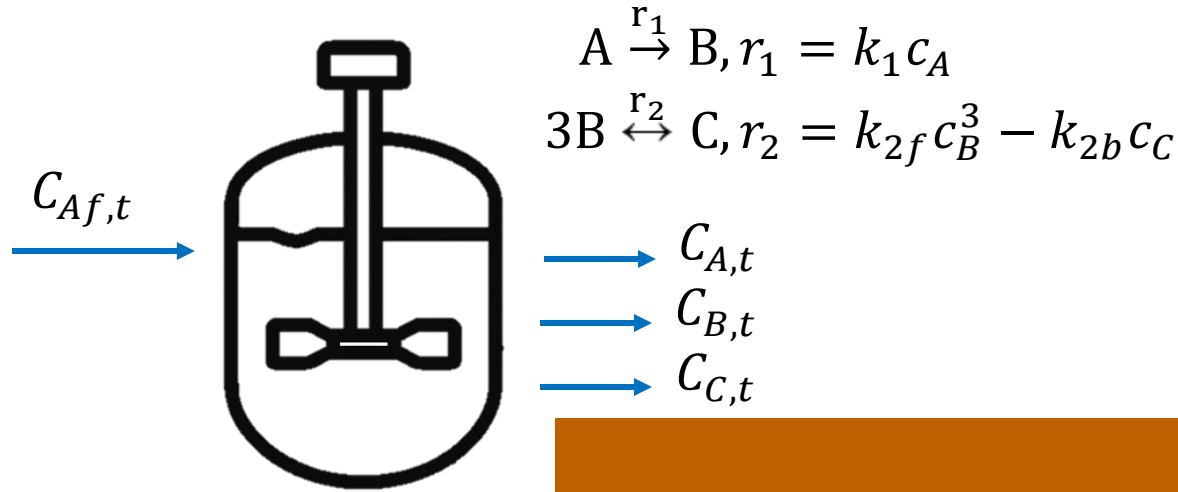
[1] Kumar, P., & Rawlings, J. B. (2023). Structured nonlinear process modeling using neural networks and application to economic optimization. *Computers and Chemical Engineering*, 177. <https://doi.org/10.1016/j.compchemeng.2023.108314>

# Example: Official TWCCC CSTR 😊

Given:  $C_{A0}, C_{B0}, C_{C0}, U_{0,...,59}$

Predict:  $C_{i1}, \dots, C_{i,60}$

Source → Target: Reactor Volume  $\Delta 20\%$



## Source System

$\dim(D_S) = 1$  year of data

## Target System

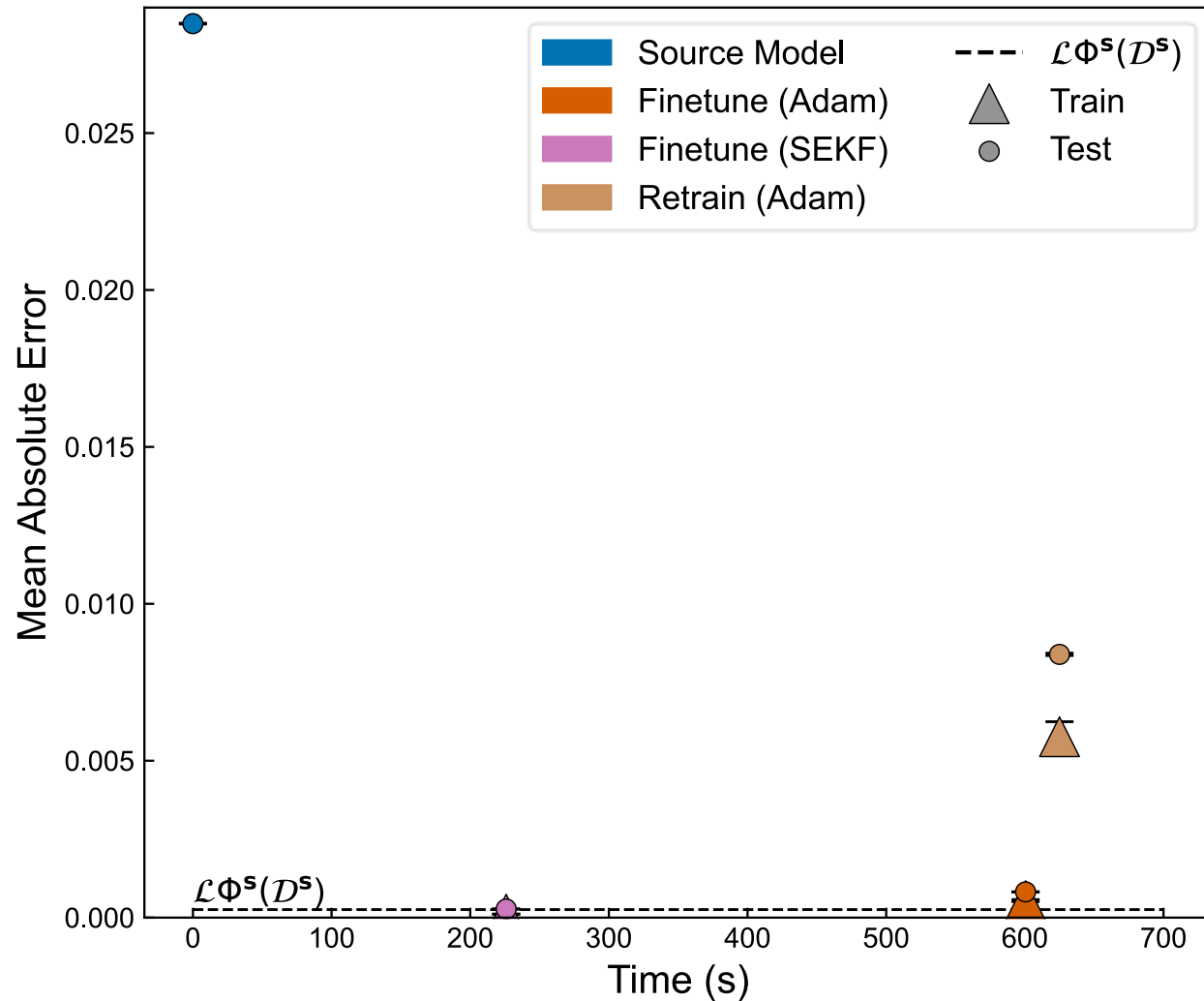
Reactor Volume  $\Delta 20\%$

$\dim(D_T) =$  One day of training data

Evaluated on remaining 363 days

Use a 2-layer Neural ODE

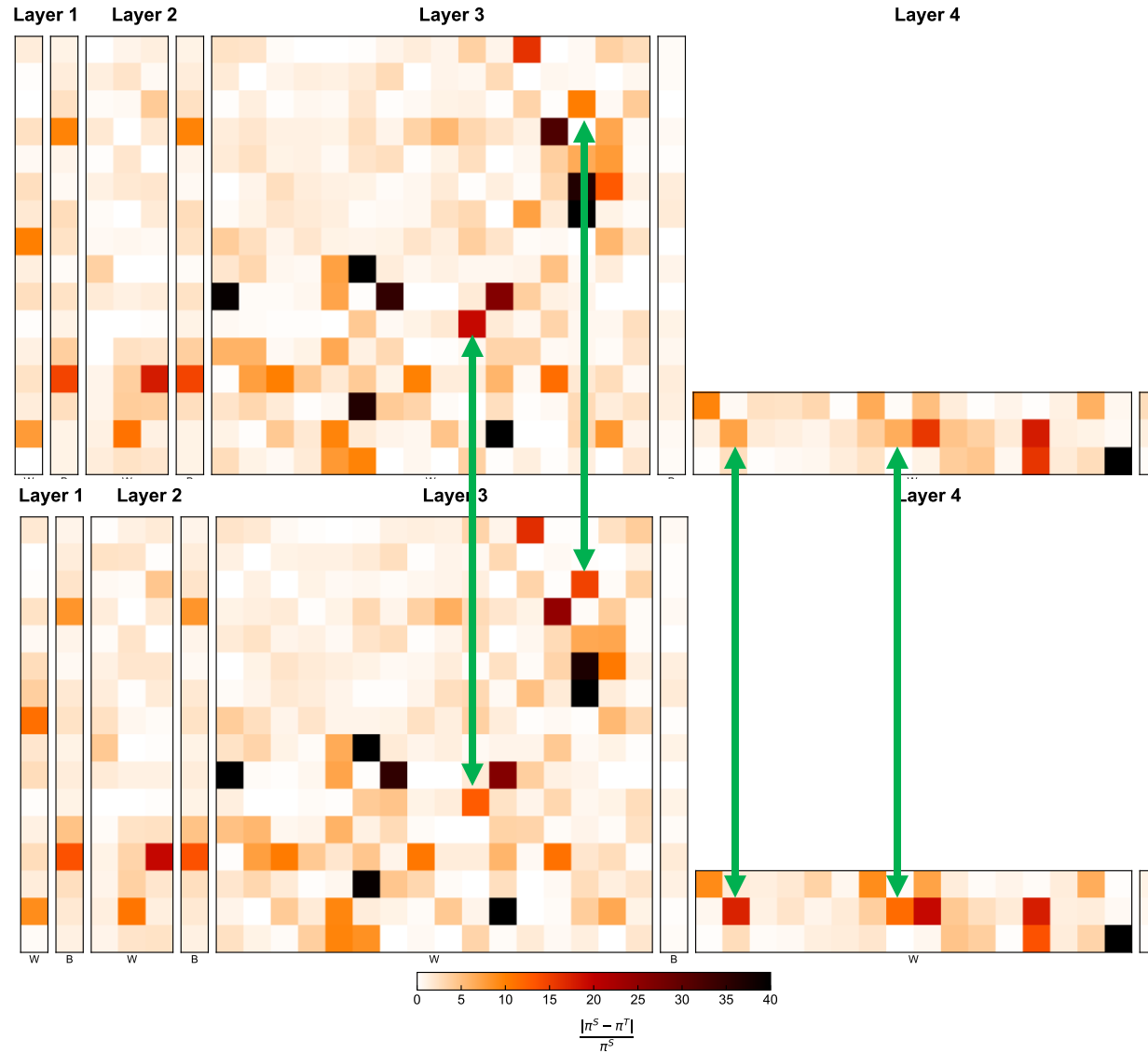
- 16 neurons each layer
- Sigmoid activation
- 419 trainable params



[1] Kumar, P., & Rawlings, J. B. (2023). Structured nonlinear process modeling using neural networks and application to economic optimization. *Computers and Chemical Engineering*, 177. <https://doi.org/10.1016/j.compchemeng.2023.108314>

# Changes to the NN parameters are distributed through the model

SEKF Weights



Adam Weights

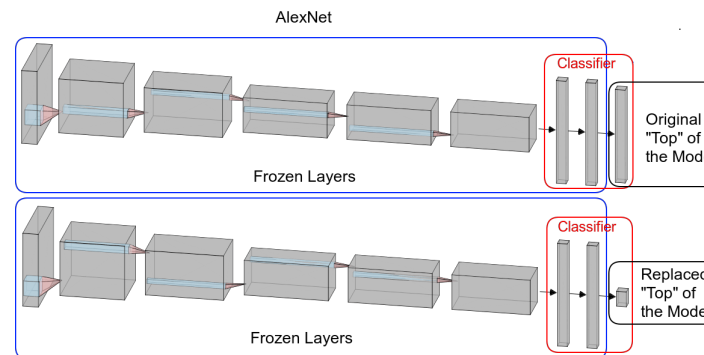
Similar but not identical

# Small Changes in NN Parameters Adapt the Model (Source → Target)

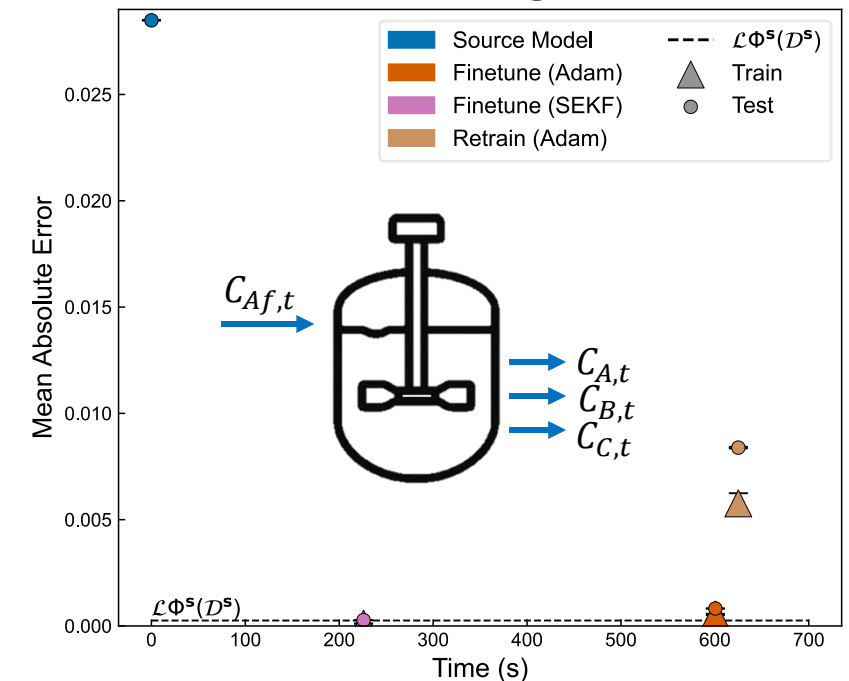
Identical processes vary in behavior.

Transfer Learning is most common in image processing. Determining **which parameters to modify** for dynamical regression is an open question.

Small changes in the source model parameters can adapt the model to the target system.

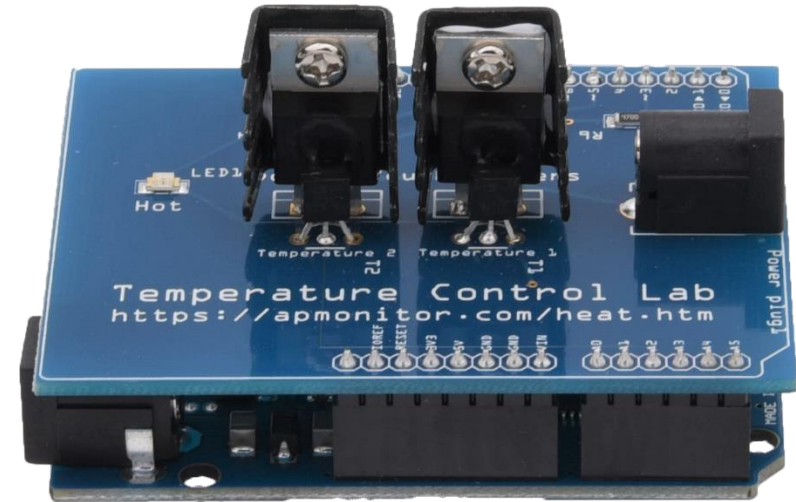
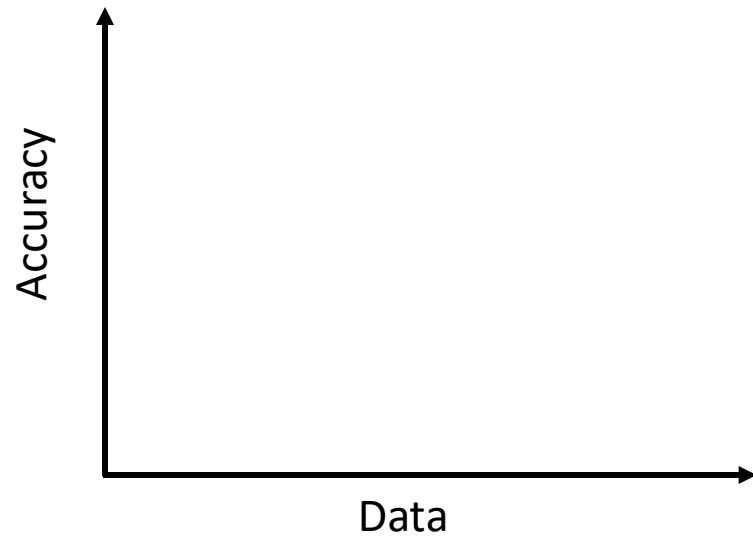


The SEKF may reduce the generalization error associated with the limited target dataset



# Future Work

---





# Thank you

---

Baldea Group

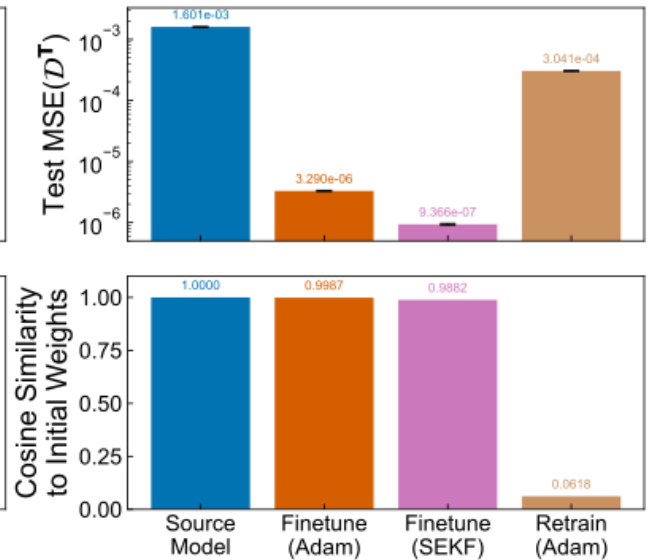
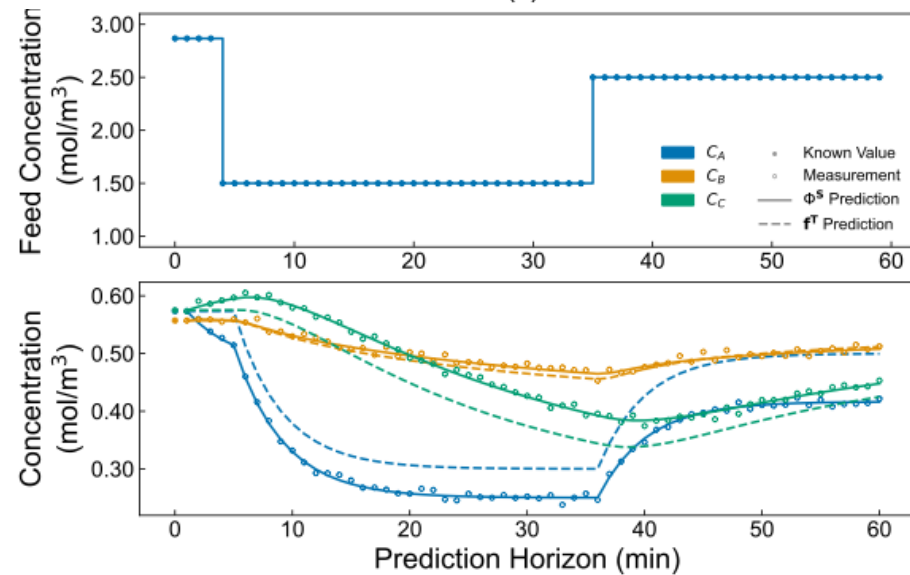
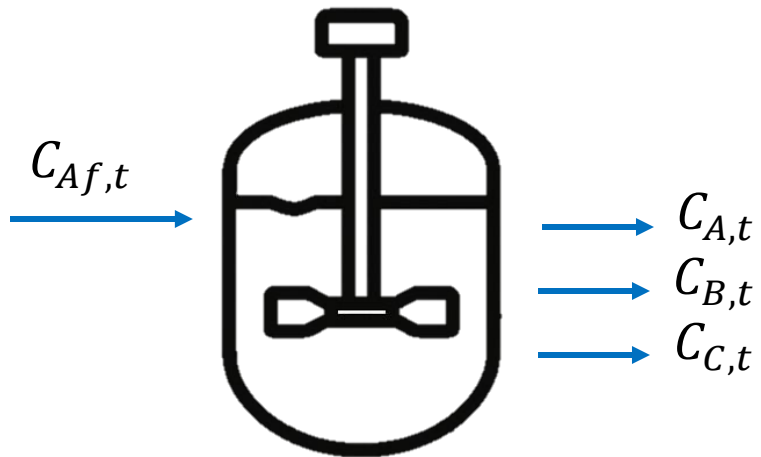
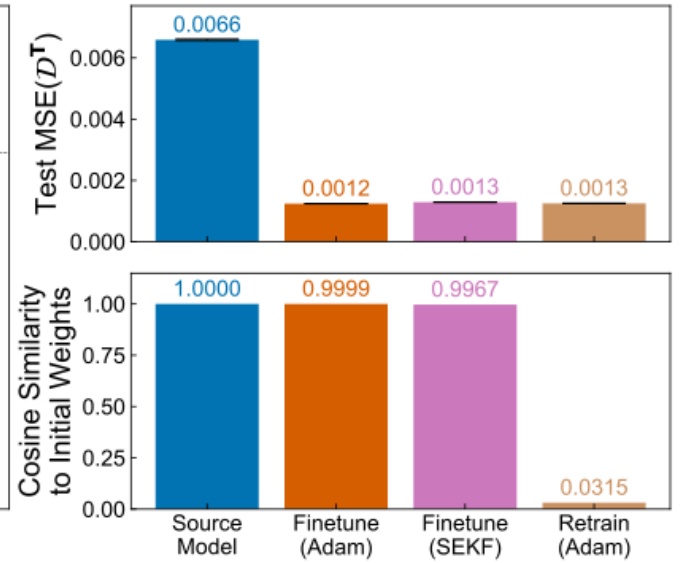
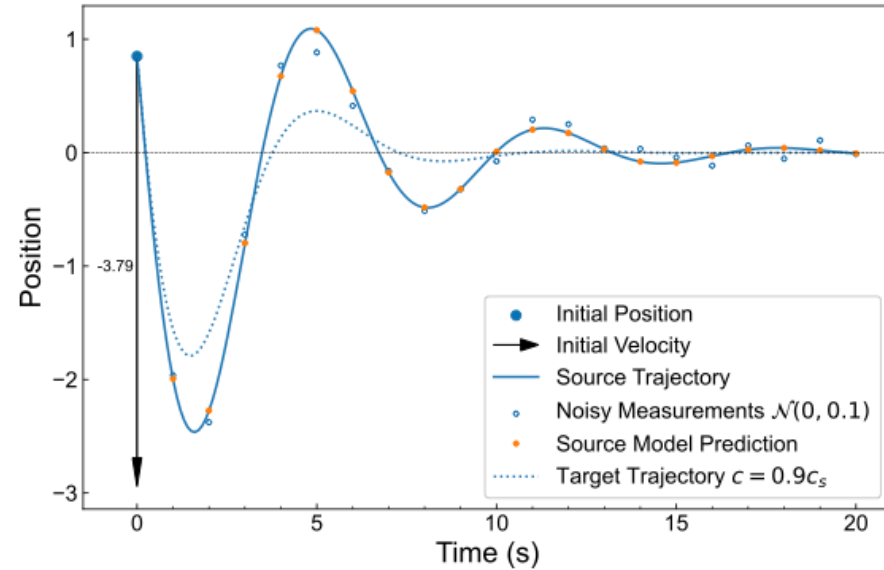
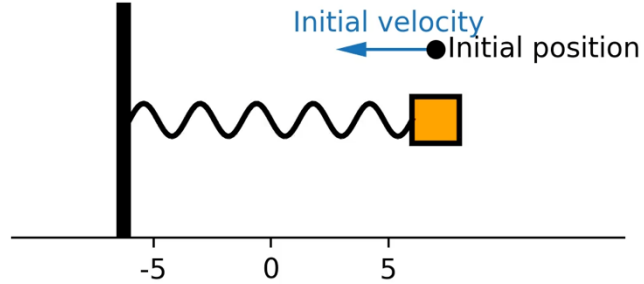


Korgel Group



Thank you to Tyler Soderstrom and ExxonMobil for the support

# Questions?



# Supplemental Slides



# Damped Spring System: Source Model Training Details

## Training Data

- 100,000 Examples:
  - 10,000 for training
  - 10,000 for validation
  - 80,000 for testing
- $x_0, \dot{x}_0$  for each examples drawn from uniform distribution  $[-5,5]$
- AHSA Hyperparameter Tuning w/ 100 samples
- Max 200 Epochs
- Hyperparameters
  - Batch size (16, 32, 64, 128, 256, 1028, 4096)
  - Decrease LR on Plateau Scheduler
    - Initial learning rate ( $1 \times 10^{-6}$  -  $1 \times 10^{-1}$ )
    - Learning rate patience(10, 20, 30, 40, 50)
    - Learning rate factor (0.1 - 0.9)

### Best Model:

- Batch Size: 16
- Initial learning rate: 0.03182
- Patience: 10 epochs
- LR factor: 0.3442

# Damped Spring System: Transfer Learning Details

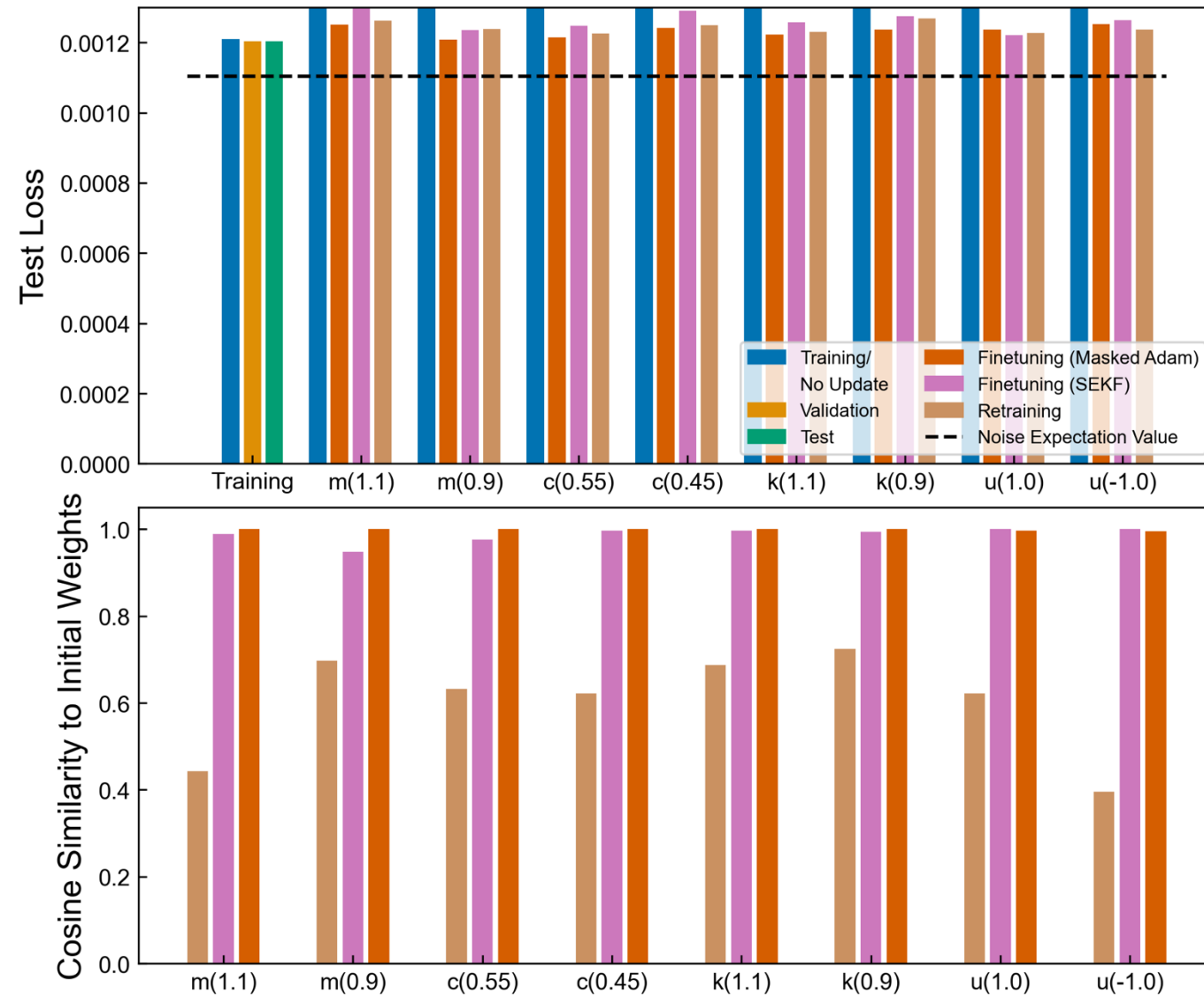
## Training Data

- 100,000 Examples:
  - 10,000 for training
  - 10,000 for validation
  - 80,000 for testing
- $x_0, \dot{x}_0$  for each examples drawn from uniform distribution  $[-5,5]$
- AHSA Hyperparameter Tuning w/ 100 samples
- Max 200 Epochs
- Hyperparameters
  - Batch size (16, 32, 64, 128, 256, 1028, 4096)
  - Decrease LR on Plateau Scheduler
    - Initial learning rate ( $1 \times 10^{-6}$  -  $1 \times 10^{-1}$ )
    - Learning rate patience(10, 20, 30, 40, 50)
    - Learning rate factor (0.1 - 0.9)

### Best SEKF Parameters:

- Batch Size: 32
- Initial learning rate: 100
- Patience: 100 epochs
- LR factor: -
- $Q = 1 \times 10^{-6} \text{ I}$
- Initial P = 100 I
- % parameters updated: 54.5 %

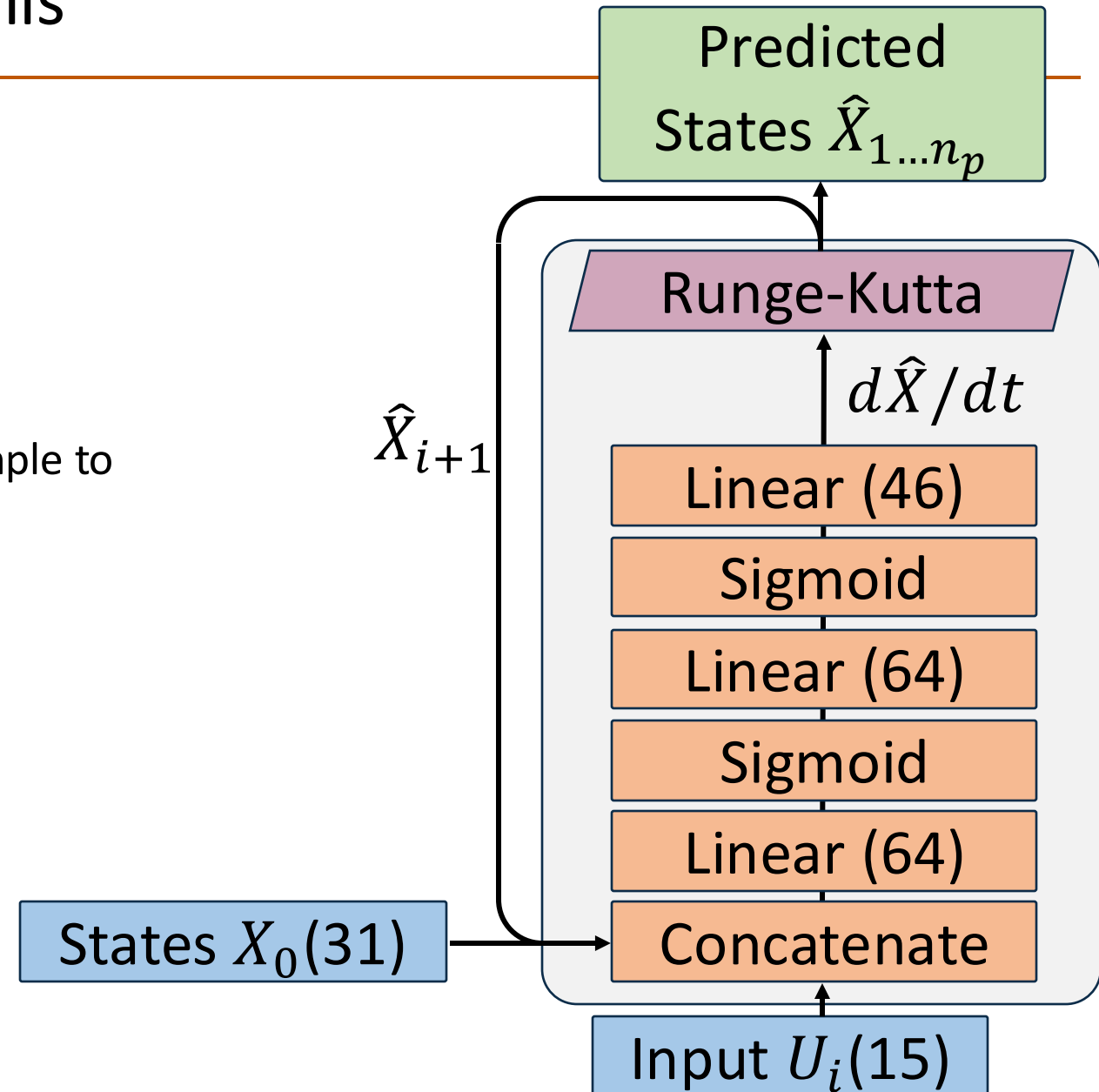
# Damped Spring System: Full Results



# CSTR System: Initial Training Details

## Training Data

- 1 Year of data
  - 80% for training
  - 10% for validation
  - 10% for testing
- 5 minute interval between each input-output example to maximize throughput while minimizing redundant information
- Hyperparameters
  - Epochs: 10
  - Learning Rate: 0.1
  - Batch Size: 256



# Multi-stream SEKF

**Problem:** The SEKF performs updates on single datapoints. When performing this across large datasets, this becomes time-consuming to iterate across all datapoints individually.

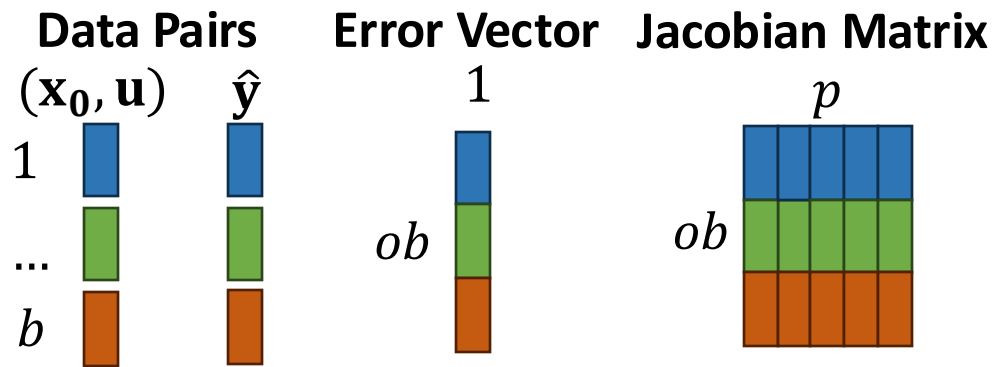
For a **batch size**  $b$ , where each prediction includes  $o$  **outputs** produced from a model with  $p$  **parameters**:

## Standard Batched Operations

- Individual losses averaged into a single scalar value
- Update calculated with respect to that single averaged loss

## Multi-Stream Operations

- Individual losses and Jacobians concatenated
- Update is the sum of each individual update calculated in parallel



$$\mathbf{e}_{ob,1} = \mathbf{y}_{ob,1} - \tilde{\mathbf{y}}_{ob,1}$$

$$\mathbf{H}_{ob,p} = \nabla_{\pi} \mathbf{y}_{ob,1}$$

$$\mathbf{A}_{ob,ob} = [\mathbf{R}_{ob,ob} + \mathbf{H}_{ob,p} \mathbf{P}_{p,p} \mathbf{H}_{p,ob}^T]^{-1}$$

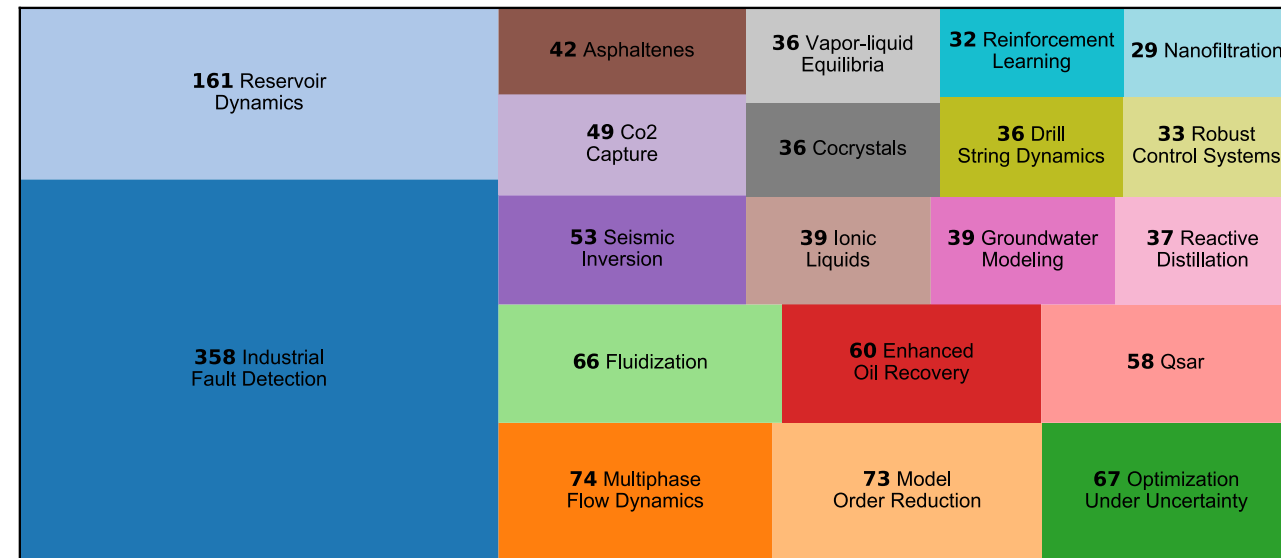
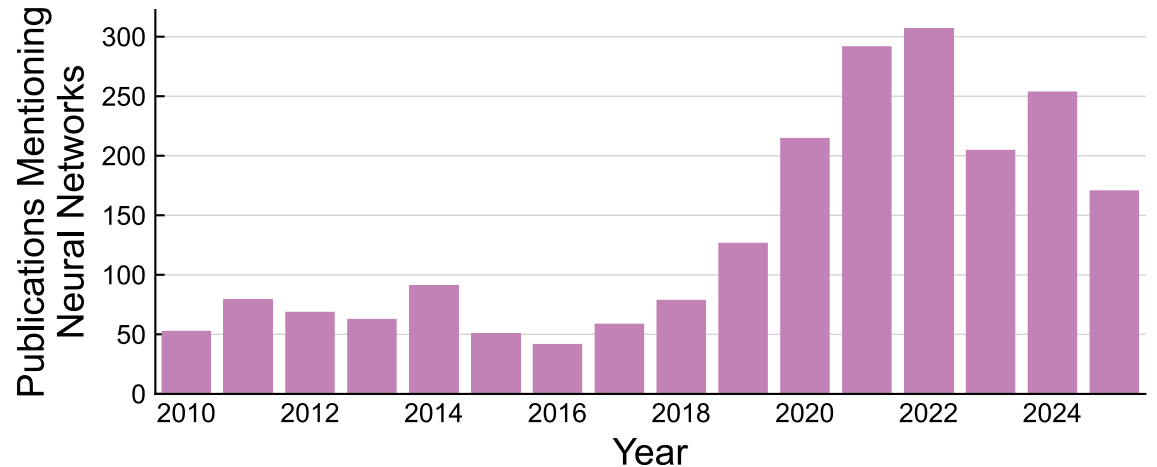
$$\mathbf{K}_{p,ob} = \mathbf{P}_{p,p} \mathbf{H}_{p,ob}^T \mathbf{A}_{ob,ob}$$

$$\hat{\pi}_{p,1} = \hat{\pi}_{p,1} + \mathbf{K}_{p,ob} \mathbf{e}_{ob,1}$$

$$\mathbf{P}_{p,p} = \mathbf{P}_{p,p} - \mathbf{K}_{p,ob} \mathbf{H}_{ob,p} \mathbf{P}_{p,p} + \mathbf{Q}_{p,p}$$

# Hot topic: neural networks in Process Systems Engineering

- Data science, machine learning, and artificial intelligence research and tools have been increasing at an exponential rate.
  - **Neural networks** have been the “hot topic”
- Most research assumes that there is sufficient **data coverage and quantity**, that the training data is drawn from the **same system** as the application, and that the **system is stationary**
- Data-driven models have no generalization or extrapolation guarantees



Web of Science Query Topic: “Neural Network” Dates: 2010 – 2025 Journals: Computers & Chemical Engineering, Industrial & Engineering Chemistry, Scientific Reports, Renewable Sustainable Energy Reviews, Journal of Process Control, Chemical Engineering Science, Chemical Engineering Research Design, AIChE Journal