

# Some time, some space, and some equations: Machine-learning of model error in dynamical systems

Matthew E. Levine<sup>1</sup>    Andrew M. Stuart<sup>1</sup>

<sup>1</sup>Department of Computing and Mathematical Sciences  
California Institute of Technology

Third Symposium on Machine Learning and Dynamical Systems  
Fields Institute, University of Toronto  
September 29, 2022

# Introduction

- Machine learning works (with enough data)!

# Introduction

- Machine learning works (with enough data)!
- Mechanistic models based on physics work (with enough knowledge and compute)!

# Introduction

- Machine learning works (with enough data)!
- Mechanistic models based on physics work (with enough knowledge and compute)!
- In most open prediction problems, we have SOME data and SOME prior knowledge.

# Introduction

- Machine learning works (with enough data)!
- Mechanistic models based on physics work (with enough knowledge and compute)!
- In most open prediction problems, we have SOME data and SOME prior knowledge.
- The next generation of high-performing prediction models will **hybridize physics-based and data-driven modeling techniques**
- How can we help lay the groundwork for this future?

# Our problem

**True system (ODE):**

$$\begin{aligned}\dot{x} &= f^\dagger(x, y) \\ \dot{y} &= \frac{1}{\varepsilon} g^\dagger(x, y)\end{aligned}\tag{1}$$

- **Relevance:** across disciplines (climatology, physiology, celestial mechanics, etc.).

# Our problem

**True system (ODE):**

$$\begin{aligned}\dot{x} &= f^\dagger(x, y) \\ \dot{y} &= \frac{1}{\varepsilon} g^\dagger(x, y)\end{aligned}\tag{1}$$

- **Relevance:** across disciplines (climatology, physiology, celestial mechanics, etc.).
- **Goal:** Given noisy observations of  $x$ , learn predictive model for future  $x$  dynamics.

# Our problem

**True system (ODE):**

$$\begin{aligned}\dot{x} &= f^\dagger(x, y) \\ \dot{y} &= \frac{1}{\varepsilon} g^\dagger(x, y)\end{aligned}\tag{1}$$

- **Relevance:** across disciplines (climatology, physiology, celestial mechanics, etc.).
- **Goal:** Given noisy observations of  $x$ , learn predictive model for future  $x$  dynamics.
- **Methodological constraints:**
  - Partial, noisy observations (e.g. observe  $x$ , but not  $y$ )



# Our problem

**True system (ODE):**

$$\begin{aligned}\dot{x} &= f^\dagger(x, y) \\ \dot{y} &= \frac{1}{\varepsilon} g^\dagger(x, y)\end{aligned}\tag{1}$$

- **Relevance:** across disciplines (climatology, physiology, celestial mechanics, etc.).
- **Goal:** Given noisy observations of  $x$ , learn predictive model for future  $x$  dynamics.
- **Methodological constraints:**
  - Partial, noisy observations (e.g. observe  $x$ , but not  $y$ )
  - No knowledge of  $y$ ,  $g^\dagger$ ,  $\varepsilon$ , nor  $\dim(y)$

# Our problem

**True system (ODE):**

$$\begin{aligned}\dot{x} &= f^\dagger(x, y) \\ \dot{y} &= \frac{1}{\varepsilon} g^\dagger(x, y)\end{aligned}\tag{1}$$

- **Relevance:** across disciplines (climatology, physiology, celestial mechanics, etc.).
- **Goal:** Given noisy observations of  $x$ , learn predictive model for future  $x$  dynamics.
- **Methodological constraints:**
  - Partial, noisy observations (e.g. observe  $x$ , but not  $y$ )
  - No knowledge of  $y$ ,  $g^\dagger$ ,  $\varepsilon$ , nor  $\dim(y)$
  - Observations may be irregularly spaced and noisy

# Our problem

**True system (ODE):**

$$\begin{aligned}\dot{x} &= f^\dagger(x, y) \\ \dot{y} &= \frac{1}{\varepsilon} g^\dagger(x, y)\end{aligned}\tag{1}$$

- **Relevance:** across disciplines (climatology, physiology, celestial mechanics, etc.).
- **Goal:** Given noisy observations of  $x$ , learn predictive model for future  $x$  dynamics.
- **Methodological constraints:**
  - Partial, noisy observations (e.g. observe  $x$ , but not  $y$ )
  - No knowledge of  $y$ ,  $g^\dagger$ ,  $\varepsilon$ , nor  $\dim(y)$
  - Observations may be irregularly spaced and noisy
  - Ability to leverage partial knowledge of  $f^\dagger$

**Big picture: Learning dynamics using w/ partial observations  $\{x_k\}_{k=0}^K$**

	Introduce delays	Introduce augmented states
Discrete-time	<div><math display="block">x_{k+1} = \Psi(x_k, x_{k-1}, \dots, x_{k-\kappa}), \quad \kappa \in \mathbb{Z}^+</math><p>Takens</p><ul style="list-style-type: none"><li><math>\Psi \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul></div>	<div><math display="block">x_{k+1} = \Psi_1(x_k, r_k)</math><math display="block">r_{k+1} = \Psi_2(x_k, r_k), \quad r_k \in \mathbb{R}^d</math><p>Tries to represent the missing states!</p><ul style="list-style-type: none"><li>RNNs / Reservoir Computers</li><li><math>\Psi_1, \Psi_2 \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul></div>
Continuous-time	<div><math display="block">\dot{x} = f(\{x(t-s)\}_{s=\tau}^t)</math><p>Takens / Mori-Zwanzig</p><ul style="list-style-type: none"><li><math>f \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul></div>	<div><math display="block">\dot{x} = f_1(x, r)</math><math display="block">\dot{r} = f_2(x, r), \quad r \in \mathbb{R}^d</math><p>Tries to represent the missing states!</p><ul style="list-style-type: none"><li>“Continuous-time” RNNs</li><li>“Continuous-time” Reservoir Computers</li><li><math>f_1, f_2 \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul></div>

# Big picture: Learning dynamics using w/ partial observations $\{x_k\}_{k=0}^K$

	Introduce delays	Introduce augmented states
Discrete-time	$x_{k+1} = \Psi(x_k, x_{k-1}, \dots, x_{k-\kappa}), \quad \kappa \in \mathbb{Z}^+$ <p>Takens</p> <ul style="list-style-type: none"><li><math>\Psi \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul>	$x_{k+1} = \Psi_1(x_k, r_k)$ $r_{k+1} = \Psi_2(x_k, r_k), \quad r_k \in \mathbb{R}^d$ <p>Tries to represent the missing states!</p> <ul style="list-style-type: none"><li>RNNs / Reservoir Computers</li><li><math>\Psi_1, \Psi_2 \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul>
Continuous-time	$\dot{x} = f(\{x(t-s)\}_{s=\tau}^t)$ <p>Takens / Mori-Zwanzig</p> <ul style="list-style-type: none"><li><math>f \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul>	$\dot{x} = f_1(x, r)$ $\dot{r} = f_2(x, r), \quad r \in \mathbb{R}^d$ <p>Tries to represent the missing states!</p> <ul style="list-style-type: none"><li>“Continuous-time” RNNs</li><li>“Continuous-time” Reservoir Computers</li><li><math>f_1, f_2 \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul>

**Big picture: Learning dynamics using w/ partial observations  $\{x_k\}_{k=0}^K$**

	Introduce delays	Introduce augmented states
Discrete-time	$x_{k+1} = \Psi(x_k, x_{k-1}, \dots, x_{k-\kappa}), \quad \kappa \in \mathbb{Z}^+$ <p>Takens</p> <ul style="list-style-type: none"><li><math>\Psi \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul>	$x_{k+1} = \Psi_1(x_k, r_k)$ $r_{k+1} = \Psi_2(x_k, r_k), \quad r_k \in \mathbb{R}^d$ <p>Tries to represent the missing states!</p> <ul style="list-style-type: none"><li>RNNs / Reservoir Computers</li><li><math>\Psi_1, \Psi_2 \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul>
Continuous-time	$\dot{x} = f(\{x(t-s)\}_{s=\tau}^t)$ <p>Takens / Mori-Zwanzig</p> <ul style="list-style-type: none"><li><math>f \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul>	$\dot{x} = f_1(x, r)$ $\dot{r} = f_2(x, r), \quad r \in \mathbb{R}^d$ <p>Tries to represent the missing states!</p> <ul style="list-style-type: none"><li>“Continuous-time” RNNs</li><li>“Continuous-time” Reservoir Computers</li><li><math>f_1, f_2 \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul>

# Big picture: Learning dynamics using w/ partial observations $\{x_k\}_{k=0}^K$

	Introduce delays	Introduce augmented states
Discrete-time	$x_{k+1} = \Psi(x_k, x_{k-1}, \dots, x_{k-\kappa}), \quad \kappa \in \mathbb{Z}^+$ <p>Takens</p> <ul style="list-style-type: none"><li><math>\Psi \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul>	$x_{k+1} = \Psi_1(x_k, r_k)$ $r_{k+1} = \Psi_2(x_k, r_k), \quad r_k \in \mathbb{R}^d$ <p>Tries to represent the missing states!</p> <ul style="list-style-type: none"><li>RNNs / Reservoir Computers</li><li><math>\Psi_1, \Psi_2 \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul>
Continuous-time	$\dot{x} = f(\{x(t-s)\}_{s=\tau}^t)$ <p>Takens / Mori-Zwanzig</p> <ul style="list-style-type: none"><li><math>f \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul>	$\dot{x} = f_1(x, r)$ $\dot{r} = f_2(x, r), \quad r \in \mathbb{R}^d$ <p>Tries to represent the missing states!</p> <ul style="list-style-type: none"><li>“Continuous-time” RNNs</li><li>“Continuous-time” Reservoir Computers</li><li><math>f_1, f_2 \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul>

# Big picture: Learning dynamics using w/ partial observations $\{x_k\}_{k=0}^K$

	Introduce delays	Introduce augmented states
Discrete-time	$x_{k+1} = \Psi(x_k, x_{k-1}, \dots, x_{k-\kappa}), \quad \kappa \in \mathbb{Z}^+$ <p>Takens</p> <ul style="list-style-type: none"><li><math>\Psi \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul>	$x_{k+1} = \Psi_1(x_k, r_k)$ $r_{k+1} = \Psi_2(x_k, r_k), \quad r_k \in \mathbb{R}^d$ <p>Tries to represent the missing states!</p> <ul style="list-style-type: none"><li>RNNs / Reservoir Computers</li><li><math>\Psi_1, \Psi_2 \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul>
Continuous-time	$\dot{x} = f(\{x(t-s)\}_{s=\tau}^t)$ <p>Takens / Mori-Zwanzig</p> <ul style="list-style-type: none"><li><math>f \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul>	$\dot{x} = f_1(x, r)$ $\dot{r} = f_2(x, r), \quad r \in \mathbb{R}^d$ <p>Tries to represent the missing states!</p> <ul style="list-style-type: none"><li>“Continuous-time” RNNs</li><li>“Continuous-time” Reservoir Computers</li><li><math>f_1, f_2 \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li></ul>



# Big picture: Learning dynamics using w/ partial observations $\{x_k\}_{k=0}^K$

	Introduce delays	Introduce augmented states
Discrete-time	$x_{k+1} = \Psi(x_k, x_{k-1}, \dots, x_{k-\kappa}), \quad \kappa \in \mathbb{Z}^+$ <p>Takens</p> <ul style="list-style-type: none"> <li><math>\Psi \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li> </ul>	$x_{k+1} = \Psi_1(x_k, r_k)$ $r_{k+1} = \Psi_2(x_k, r_k), \quad r_k \in \mathbb{R}^d$ <p>Tries to represent the missing states!</p> <ul style="list-style-type: none"> <li>RNNs / Reservoir Computers</li> <li><math>\Psi_1, \Psi_2 \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li> </ul>
Continuous-time	$\dot{x} = f(\{x(t-s)\}_{s=\tau}^t)$ <p>Takens / Mori-Zwanzig</p> <ul style="list-style-type: none"> <li><math>f \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li> </ul>	$\dot{x} = f_1(x, r)$ $\dot{r} = f_2(x, r), \quad r \in \mathbb{R}^d$ <p>Tries to represent the missing states!</p> <ul style="list-style-type: none"> <li>“Continuous-time” RNNs</li> <li>“Continuous-time” Reservoir Computers</li> <li><math>f_1, f_2 \in \{\text{NNs, GPs, Polynomials, RFs, ...}\}</math></li> </ul>

# Leveraging partial knowledge of the dynamics

For any  $f_0$  (regardless of its fidelity), there exists an  $m^\dagger(x, y)$  such that (1) can be re-written as

$$\dot{x} = f_0(x) + m^\dagger(x, y) \quad (2a)$$

$$\dot{y} = \frac{1}{\varepsilon} g^\dagger(x, y). \quad (2b)$$

Qianxiao Li, MLDS 2022 slides!!

- Approximation Rates. Let  $\mathcal{H} = \cup_m \mathcal{H}_m$ , where  $\mathcal{H}_m \subset \mathcal{H}_{m+1}$ ,  $m$  measures size of hypothesis space (approximation budget)

$$\inf_{\hat{F} \in \mathcal{H}_m} \|F^* - \hat{F}\| \leq \text{Complexity}(F^*) \text{rate}(m), \quad \text{rate}(m) \rightarrow 0$$

# Learning latent dynamics in continuous-time

$$\begin{aligned} \dot{x} &= f_0(x) + m(x, r; \theta) \\ \dot{r} &= g(x, r; \theta) \end{aligned} \quad \Longleftrightarrow \quad \begin{aligned} \dot{u} &= f(u; \theta), \quad u = [x, r]^T \\ Hu &= x \end{aligned}$$

# Learning latent dynamics in continuous-time

$$\begin{aligned} \dot{x} &= f_0(x) + m(x, r; \theta) \\ \dot{r} &= g(x, r; \theta) \end{aligned} \quad \Longleftrightarrow \quad \begin{aligned} \dot{u} &= f(u; \theta), \quad u = [x, r]^T \\ Hu &= x \end{aligned}$$

Assume noisy observations  $z = Hu + \eta$ .

# Learning latent dynamics in continuous-time

$$\begin{aligned} \dot{x} &= f_0(x) + m(x, r; \theta) & \iff & \dot{u} = f(u; \theta), \quad u = [x, r]^T \\ \dot{r} &= g(x, r; \theta) & & Hu = x \end{aligned}$$

Assume noisy observations  $z = Hu + \eta$ .  
Let  $u(t; v, \theta)$  solve  $\dot{u} = f(u; \theta)$ ,  $u(0) = v$ .

# Learning latent dynamics in continuous-time

$$\begin{aligned} \dot{x} &= f_0(x) + m(x, r; \theta) \\ \dot{r} &= g(x, r; \theta) \end{aligned} \quad \Longleftrightarrow \quad \begin{aligned} \dot{u} &= f(u; \theta), \quad u = [x, r]^T \\ Hu &= x \end{aligned}$$

Assume noisy observations  $z = Hu + \eta$ .  
Let  $u(t; v, \theta)$  solve  $\dot{u} = f(u; \theta)$ ,  $u(0) = v$ .

**Hard Constraint Idea 1:** Infer init. cond. and parameters (Rubanova *et al.* 2019)

$$\operatorname{argmin}_{\theta, u_0} \int_0^T \|z(t) - Hu(t; u_0, \theta)\|^2 dt.$$

- $\times$  Poorly-posed with larger  $T$  for chaotic systems with sensitivity to  $u_0$ .

# Learning latent dynamics in continuous-time

$$\begin{aligned} \dot{x} &= f_0(x) + m(x, r; \theta) & \iff & \dot{u} = f(u; \theta), \quad u = [x, r]^T \\ \dot{r} &= g(x, r; \theta) & & Hu = x \end{aligned}$$

Assume noisy observations  $z = Hu + \eta$ .  
Let  $u(t; v, \theta)$  solve  $\dot{u} = f(u; \theta)$ ,  $u(0) = v$ .

**Hard Constraint Idea 2:** Break data into chunks to cope with sensitivities

$$\operatorname{argmin}_{\theta, \{u_0^{(k)}\}_{k=1}^K} \sum_{k=1}^K \int_0^T \|z^{(k)}(t) - Hu(t; u_0^{(k)}, \theta)\|^2 dt.$$

- ✗ Dimensionality of inference grows with  $T$ .
- ✗ When chunks are small, can overfit by "cherry-picking"  $u_0^{(k)}$ .



# Learning latent dynamics in continuous-time

$$\begin{aligned} \dot{x} &= f_0(x) + m(x, r; \theta) & \iff & \dot{u} = f(u; \theta), \quad u = [x, r]^T \\ \dot{r} &= g(x, r; \theta) & & Hu = x \end{aligned}$$

Assume noisy observations  $z = Hu + \eta$ .  
Let  $u(t; v, \theta)$  solve  $\dot{u} = f(u; \theta)$ ,  $u(0) = v$ .

**Data Assimilation for inferring missing dynamics:**

$$\operatorname{argmin}_{\theta, \{u_0^{(k)}\}_{k=1}^K} \sum_{k=1}^K \int_0^T \|z^{(k)}(t) - Hu(t; u_0^{(k)}, \theta)\|^2 dt.$$

- Initial conditions  $u_0^{(k)}$  can be inferred using a sequence of warmup data (and assumption on  $\theta$ ) using standard *Data Assimilation* techniques.



# Learning latent dynamics in continuous-time

$$\begin{aligned} \dot{x} &= f_0(x) + m(x, r; \theta) & \iff & \dot{u} = f(u; \theta), \quad u = [x, r]^T \\ \dot{r} &= g(x, r; \theta) & & Hu = x \end{aligned}$$

Assume noisy observations  $z = Hu + \eta$ .

Let  $u(t; v, \theta)$  solve  $\dot{u} = f(u; \theta)$ ,  $u(0) = v$ .

Let  $\hat{m}(t, \tau, \theta_{\text{DYN}}, \theta_{\text{DA}})$  be an estimate of  $u(t) \mid \{z(t-s)\}_{s=0}^{\tau}, \theta_{\text{DYN}}, u(t-\tau) = 0$ .

# Learning latent dynamics in continuous-time

$$\begin{aligned} \dot{x} &= f_0(x) + m(x, r; \theta) & \iff & \dot{u} = f(u; \theta), \quad u = [x, r]^T \\ \dot{r} &= g(x, r; \theta) & & Hu = x \end{aligned}$$

Assume noisy observations  $z = Hu + \eta$ .

Let  $u(t; v, \theta)$  solve  $\dot{u} = f(u; \theta)$ ,  $u(0) = v$ .

Let  $\hat{m}(t, \tau, \theta_{\text{DYN}}, \theta_{\text{DA}})$  be an estimate of  $u(t) \mid \{z(t-s)\}_{s=0}^\tau$ ,  $\theta_{\text{DYN}}$ ,  $u(t-\tau) = 0$ .

**DA-based inference:** Initial conditions can be estimated jointly with parameters

$$\operatorname{argmin}_{\theta_{\text{DYN}}, \theta_{\text{DA}}} \sum_{k=1}^K \int_0^T \|z^{(k)}(t) - Hu(t; \hat{m}(t_k, \tau, \theta_{\text{DYN}}, \theta_{\text{DA}}), \theta_{\text{DYN}})\|^2 dt.$$

# Learning latent dynamics in continuous-time

$$\begin{aligned} \dot{x} &= f_0(x) + m(x, r; \theta) & \iff & \dot{u} = f(u; \theta), \quad u = [x, r]^T \\ \dot{r} &= g(x, r; \theta) & & Hu = x \end{aligned}$$

Assume noisy observations  $z = Hu + \eta$ .

Let  $u(t; v, \theta)$  solve  $\dot{u} = f(u; \theta)$ ,  $u(0) = v$ .

Let  $\hat{m}(t, \tau, \theta_{\text{DYN}}, \theta_{\text{DA}})$  be an estimate of  $u(t) \mid \{z(t-s)\}_{s=0}^{\tau}$ ,  $\theta_{\text{DYN}}$ ,  $u(t-\tau) = 0$ .

**DA-based inference:** Initial conditions can be estimated jointly with parameters

$$\operatorname{argmin}_{\theta_{\text{DYN}}, \theta_{\text{DA}}} \sum_{k=1}^K \int_0^T \|z^{(k)}(t) - Hu(t; \hat{m}(t_k, \tau, \theta_{\text{DYN}}, \theta_{\text{DA}}), \theta_{\text{DYN}})\|^2 dt.$$

- Here, we perform joint estimation with auto-differentiable 3DVAR
- Chen *et al.* 2021 perform joint estimation with auto-differentiable Ensemble Kalman Filter
- Carassi *et al.* 2021 apply alternating descent (EnKF for  $\hat{m}$ , supervised SGD for  $\theta$ )

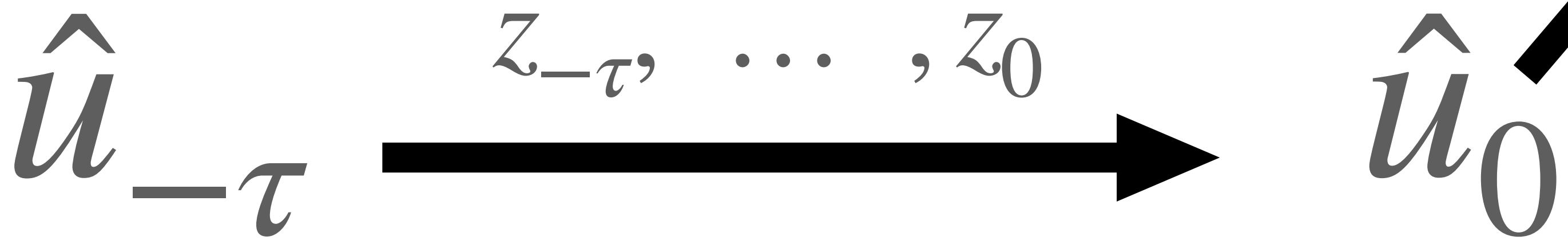
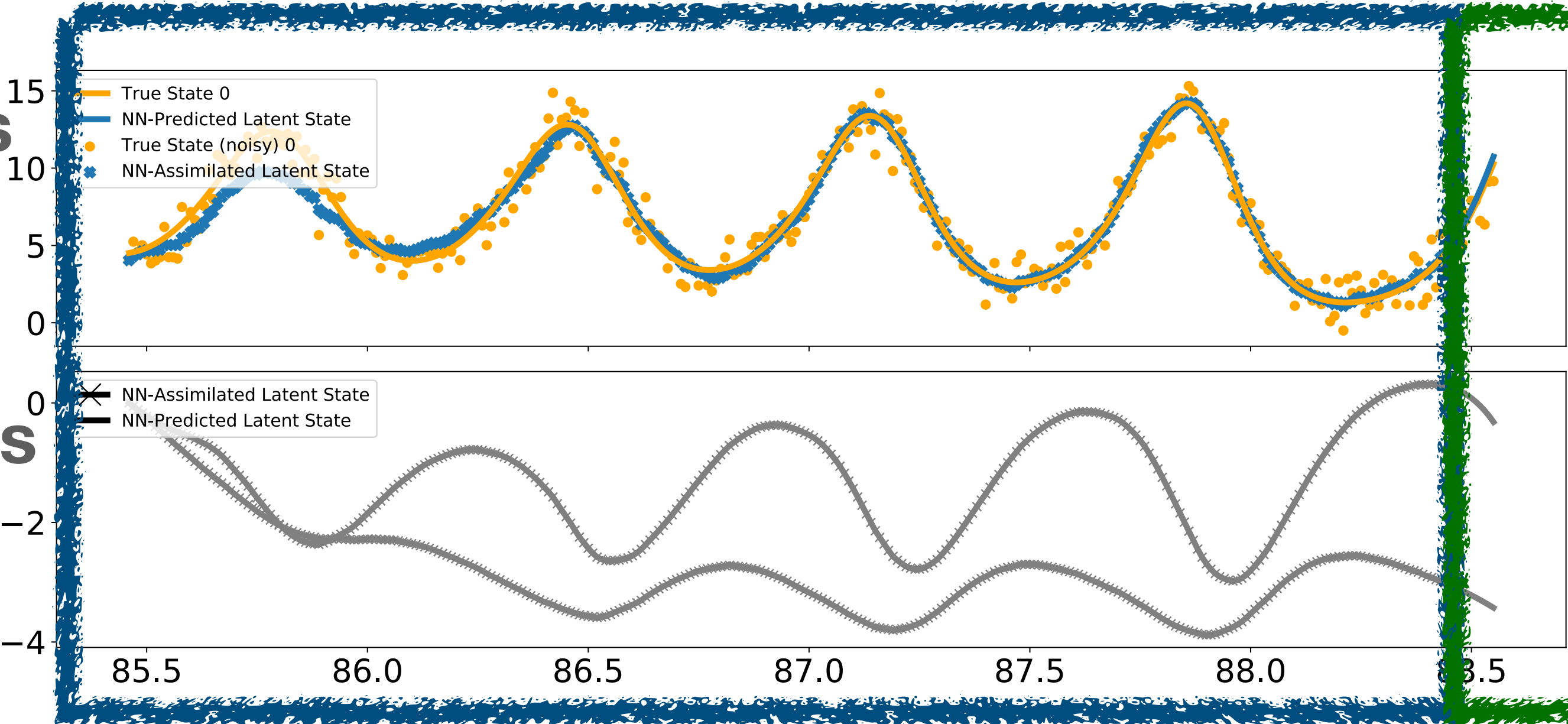
Let  $\Psi(v; \theta) := u(\Delta t; v, \theta)$  denote a integrator of our RHS that maps us to the next data element

Auto-differentiable WARMUP via Data Assimilation

Auto-differentiable FORECAST

Observed states  
 $x_1(t), z(t)$

Unobserved states  
 $x_2(t), x_3(t)$

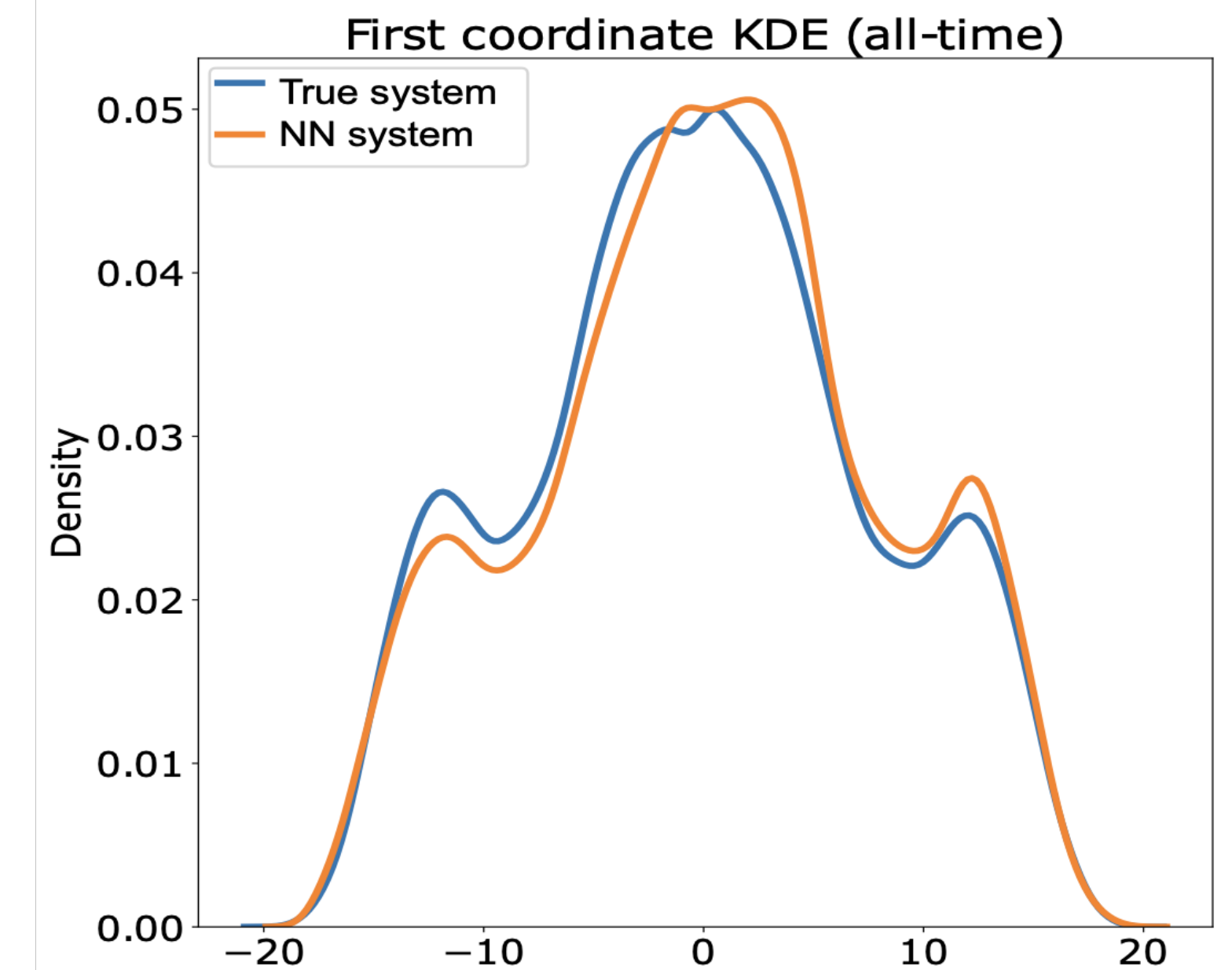
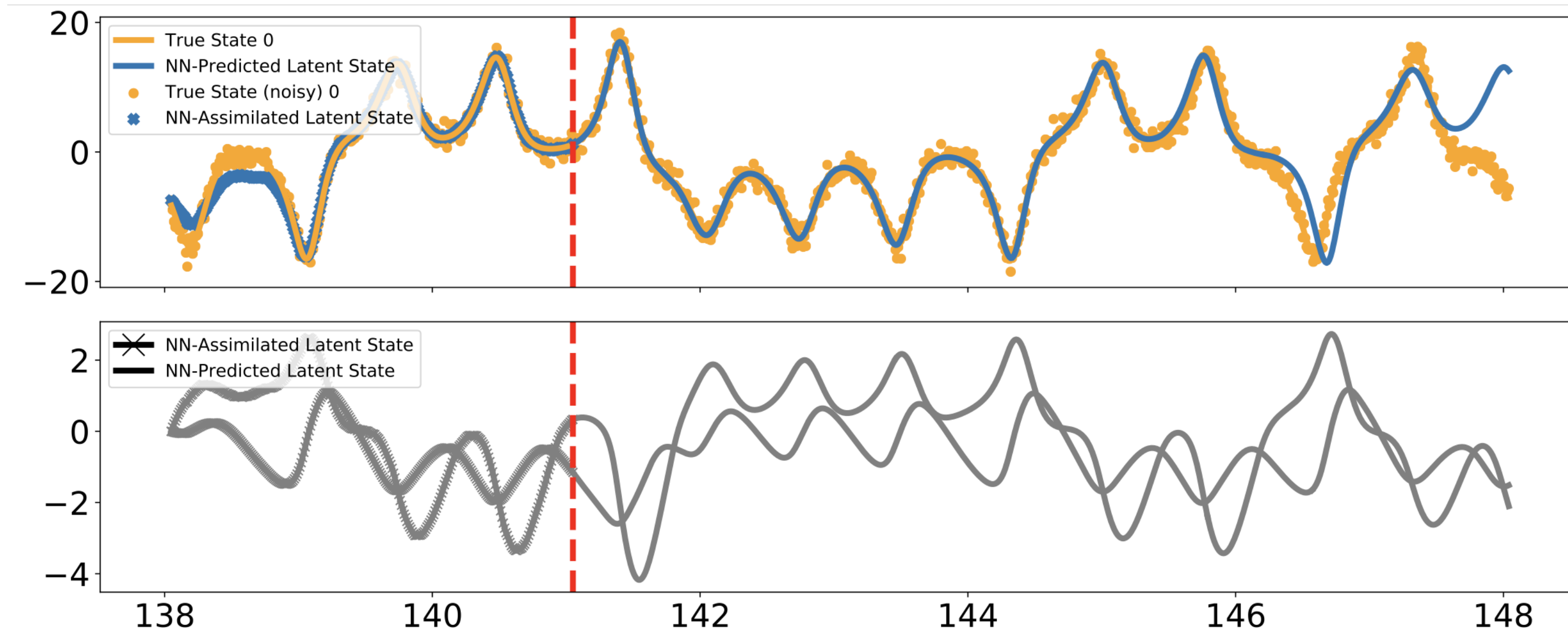


$$\hat{u}_{k+1} = \Psi(\hat{u}_k; \Delta t) + K \left( z_{k+1} - H \Psi(\hat{u}_k; \Delta t) \right)$$

Sequential updates w/ constant gain

## Accurate short-term forecasts and long-term statistics for first component of L63

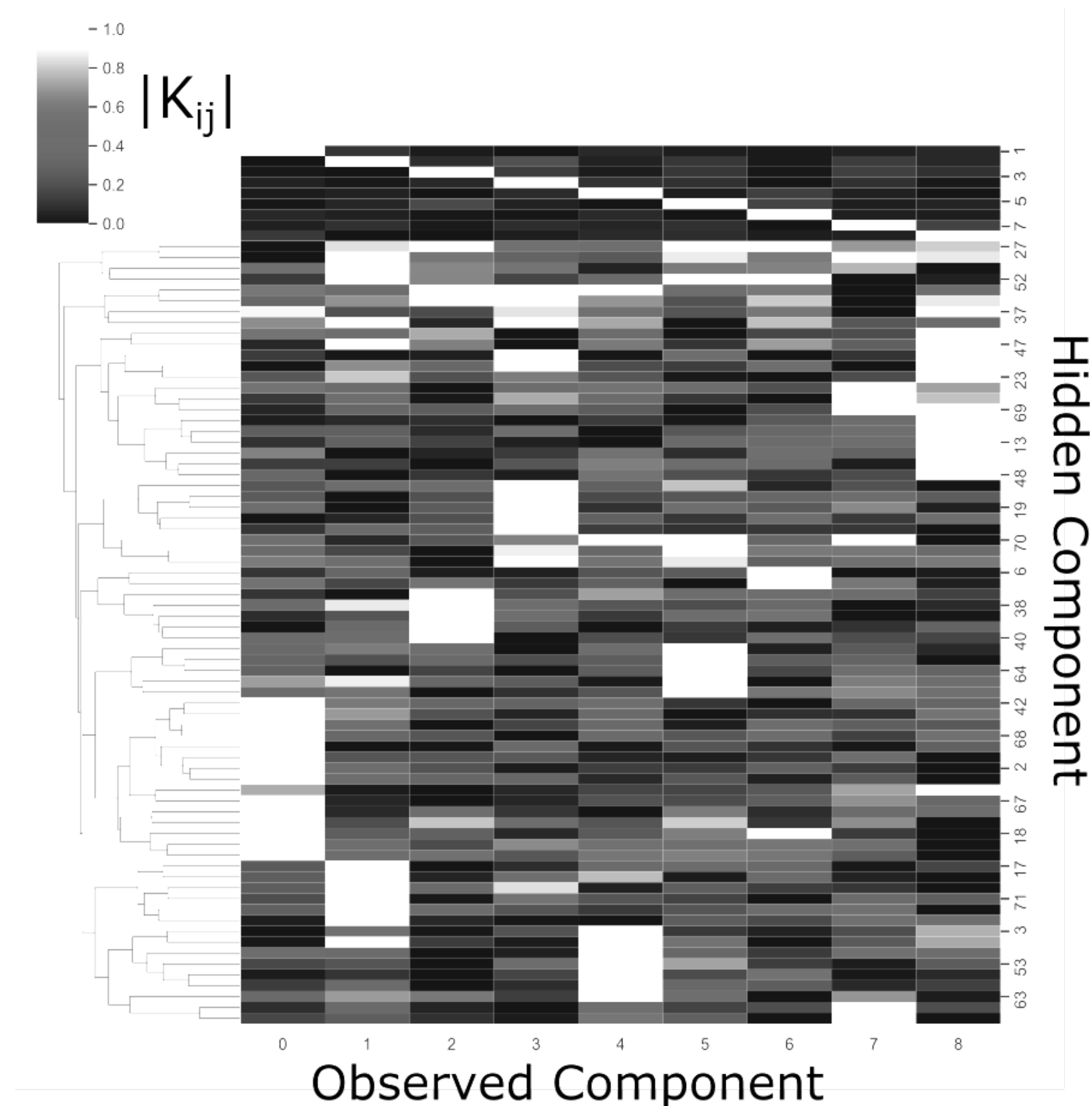
**Lorenz '63** with partial, noisy observations—noisily observe only the first component, and model its dynamics using our augmented-state model.





# Learning L96MS memory-based closure

- The true L96MS system has a clustered subgrouping of fast variables—our model has re-discovered this structure, and the DA gain  $K$  has learnt to exploit these correlations for improved filtering.



# Reservoir computing with connections to random features

$$x_{k+1} = Cr_{k+1} \tag{5a}$$

$$r_{k+1} = \sigma(W_1 r_k + W_2 x_k + b) \tag{5b}$$

- Randomize and fix  $(W_1, W_2, b, r_0)$ .
- Given  $\{x_k\}_{k=0}^K, r_0, W_1, W_2, b$ , we can determine  $\{r_k\}_{k=1}^K$
- This is great, because now we just need to do a linear regression!  $C: r_k \mapsto x_k$ .

**Debate:** For a fully observed system, what choice of  $W_1$  gives the best RC?

# Reservoir computing with connections to random features

$$x_{k+1} = Cr_{k+1} \tag{5a}$$

$$r_{k+1} = \sigma(W_1 r_k + W_2 x_k + b) \tag{5b}$$

- Randomize and fix  $(W_1, W_2, b, r_0)$ .
- Given  $\{x_k\}_{k=0}^K, r_0, W_1, W_2, b$ , we can determine  $\{r_k\}_{k=1}^K$
- This is great, because now we just need to do a linear regression!  $C: r_k \mapsto x_k$ .

**Debate:** For a fully observed system, what choice of  $W_1$  gives the best RC?

**My answer:** Easy, just  $W_1 = 0$ . Then  $r_{k+1} = \sigma(W_2 x_k + b)$ , and

$$x_{k+1} = \sum_j C^{(j)} \sigma(W_2^{(j)} x_k + b^{(j)}), \quad \text{a random feature model!}$$

For example, choosing  $\sigma := \cos(\cdot)$ ,  $W_2^{(j)} \sim \mathcal{N}(0, \Sigma)$ , and  $b^{(j)} \sim \mathcal{U}[-2\pi, 2\pi]$ , approximates Gaussian Process with RBF kernel in large feature limit.



# Empirical evaluations of learnt memory

$$x_{k+1} = Cr_{k+1} \tag{6a}$$

$$r_{k+1} = \sigma(W_1 r_k + W_2 x_k + b) \tag{6b}$$

- Often we apply RCs to partially-observed systems, which have substantial Markovian properties on the observables.
- Do you ever worry that your RC performance is not *really* learning *memory*?
- Easy sanity check: Set  $W_1 = 0$  and re-tune hyperparameters.
- The difference in performance between this and our RC represents the amount of memory that the RC has accounted for.

# Expressivity of RCs/RNNs: latent dims vs function complexity

$$x_{k+1} = Cr_{k+1}$$

$$r_{k+1} = \sigma(W_1 r_k + W_2 x_k + b)$$

Versus

$$x_{k+1} = C\sigma(Ar_{k+1} + a)$$

$$r_{k+1} = B\sigma(W_1 r_k + W_2 x_k + b)$$

- # latent variables  $\equiv \dim(r)$
- Expressivity  $\propto \dim(r)$
- Approximator limit leads to  $\dim(r) = \infty$ , even for a finite dimensional system!

- # latent variables  $\equiv \dim(r)$
- Expressivity  $\propto \dim(A, W_1, W_2)$
- Decouples dimension from expressivity;  
Can use infinitely many parameters in a finite dimensional space.

# Expressivity of RCs/RNNs: latent dims vs function complexity

$$x_{k+1} = Cr_{k+1}$$

$$r_{k+1} = \sigma(W_1 r_k + W_2 x_k + b)$$

Versus

$$x_{k+1} = C\sigma(Ar_{k+1} + a)$$

$$r_{k+1} = B\sigma(W_1 r_k + W_2 x_k + b)$$

- # latent variables  $\equiv \dim(r)$
- Expressivity  $\propto \dim(r)$
- Approximator limit leads to  $\dim(r) = \infty$ , even for a finite dimensional system!

- # latent variables  $\equiv \dim(r)$
- Expressivity  $\propto \dim(A, W_1, W_2)$
- Decouples dimension from expressivity;  
Can use infinitely many parameters in a finite dimensional space.

# Continuous-time RCs and Random Features on Banach Spaces

$$\dot{x} = Cr(t) \tag{7a}$$

$$\dot{r} = \sigma(W_1 r + W_2 x + b) \tag{7b}$$

**Question:** How to get RCs to work well in continuous time? What distributions should we choose for  $W_1, W_2$ ?

**My answer:** Shrug? I haven't gotten it to work well. I would love to hear from you all on this!

# Conclusions

- We can learn ODEs from partially observed, noisy data by embedding state-estimation techniques within the optimization
  - Can be used for tuning DA parameters
  - Can move this to a derivative-free optimization (if your models are too huge to differentiate through)
  - Currently using for modeling endocrine dynamics in patients with diabetes (joint work with Emily Fox)
- Reservoir Computers and random feature methods are quite connected
  - With  $W_r = 0$ , RC approximates a markovian GP.
  - With  $W_r \neq 0$ , something deeper is going on! If you have ideas on this, come see me! We (Ollie Dunbar, Nick Nelsen, and I) are organizing a small ICIAM minisymposium around this topic.
- I'll graduate in 2023 and need a job!
  - I want to deploy this work in biomedical applications and improve the methods until they work with real data and solve real problems!

# Conclusions

- We can learn ODEs from partially observed, noisy data by embedding state-estimation techniques within the optimization
  - Can be used for tuning DA parameters
  - Can move this to a derivative-free optimization (if your models are too huge to differentiate through)
  - Currently using for modeling endocrine dynamics in patients with diabetes (joint work with Emily Fox)
- Reservoir Computers and random feature methods are quite connected
  - With  $W_r = 0$ , RC approximates a markovian GP.
  - With  $W_r \neq 0$ , something deeper is going on! If you have ideas on this, come see me! We (Ollie Dunbar, Nick Nelsen, and I) are organizing a small ICIAM minisymposium around this topic.
- I'll graduate in 2023 and need a job!
  - I want to deploy this work in biomedical applications and improve the methods until they work with real data and solve real problems!

# Conclusions

- We can learn ODEs from partially observed, noisy data by embedding state-estimation techniques within the optimization
  - Can be used for tuning DA parameters
  - Can move this to a derivative-free optimization (if your models are too huge to differentiate through)
  - Currently using for modeling endocrine dynamics in patients with diabetes (joint work with Emily Fox)
- Reservoir Computers and random feature methods are quite connected
  - With  $W_r = 0$ , RC approximates a markovian GP.
  - With  $W_r \neq 0$ , something deeper is going on! If you have ideas on this, come see me! We (Ollie Dunbar, Nick Nelsen, and I) are organizing a small ICIAM minisymposium around this topic.
- I'll graduate in 2023 and need a job!
  - I want to deploy this work in biomedical applications and improve the methods until they work with real data and solve real problems!



# Related Work

- Kaheman, Kadierdan, Eurika Kaiser, Benjamin Strom, J. Nathan Kutz, and Steven L. Brunton. “Learning Discrepancy Models From Experimental Data.” ArXiv:1909.08574 [Cs, Eess, Stat], September 18, 2019. <http://arxiv.org/abs/1909.08574>.
- Tipireddy, Ramakrishna, Paris Perdikaris, Panos Stinis, and Alexandre Tartakovsky. “A Comparative Study of Physics-Informed Neural Network Models for Learning Unknown Dynamics and Constitutive Relations.” ArXiv:1904.04058 [Physics], April 2, 2019. <http://arxiv.org/abs/1904.04058>.
- Lovelett, Robert J., Jose L. Avalos, and Ioannis G. Kevrekidis. “Partial Observations and Conservation Laws: Grey-Box Modeling in Biotechnology and Optogenetics.” ArXiv:1909.04234 [Math], September 9, 2019. <http://arxiv.org/abs/1909.04234>.
- Pathak, Jaideep, Alexander Wikner, Rebeckah Fussell, Sarthak Chandra, Brian R. Hunt, Michelle Girvan, and Edward Ott. “Hybrid Forecasting of Chaotic Processes: Using Machine Learning in Conjunction with a Knowledge-Based Model.” Chaos: An Interdisciplinary Journal of Nonlinear Science 28, no. 4 (April 1, 2018): 041101. <https://doi.org/10.1063/1.5028373>.



# Related Work

- Rico-Martines, R., I. G. Kevrekidis, M. C. Kube, and J. L. Hudson. “Discrete- vs. Continuous-Time Nonlinear Signal Processing: Attractors, Transitions and Parallel Implementation Issues.” In 1993 American Control Conference, 1475–79. San Francisco, CA, USA: IEEE, 1993. <https://doi.org/10.23919/ACC.1993.4793116>.
- Chang, Bo, Minmin Chen, Eldad Haber, and Ed H. Chi. “AntisymmetricRNN: A Dynamical System View on Recurrent Neural Networks.” ArXiv:1902.09689 [Cs, Stat], February 25, 2019. <http://arxiv.org/abs/1902.09689>.
- Funahashi, Ken-ichi, and Yuichi Nakamura. “Approximation of Dynamical Systems by Continuous Time Recurrent Neural Networks.” Neural Networks 6, no. 6 (January 1, 1993): 801–6. [https://doi.org/10.1016/S0893-6080\(05\)80125-X](https://doi.org/10.1016/S0893-6080(05)80125-X).
- Schäfer, Anton Maximilian. “RECURRENT NEURAL NETWORKS ARE UNIVERSAL APPROXIMATORS,” 2007, 11.
- Wan, Zhong Yi, Pantelis Vlachas, Petros Koumoutsakos, and Themistoklis Sapsis. “Data-Assisted Reduced-Order Modeling of Extreme Events in Complex Dynamical Systems.” PLOS ONE 13, no. 5 (May 24, 2018): e0197704.

# Related Work

- Lei, Youming, Jian Hu, and Jianpeng Ding. “A Hybrid Model Based on Deep LSTM for Predicting High-Dimensional Chaotic Systems.” ArXiv:2002.00799 [Cs, Eess], January 21, 2020. <http://arxiv.org/abs/2002.00799>.
- Sherstinsky, Alex. “Fundamentals of Recurrent Neural Network (RNN) and Long Short-Term Memory (LSTM) Network.” Physica D: Nonlinear Phenomena 404 (March 1, 2020): 132306. <https://doi.org/10.1016/j.physd.2019.132306>.

# Thank you!

- Many thanks to my adviser **Andrew Stuart** and the graduate students and postdocs of the Caltech CMS department.
- Many more thanks to **Boumediene** for the kind invitation and to all of the **MLDS organizers** and **members of the Fields Institute** for welcoming us!
- Thanks to the **participants and speakers** who have made this an enriching week!
- For more info, see:
  - **My website:** `mlevine@netlify.com`
  - **Our paper:** Levine and Stuart, A Framework for Machine Learning of Dynamical Systems, To appear in Communications of the AMS: Volume 2. 2022.  
<https://arxiv.org/abs/2107.06658>