

1           **DATA-DRIVEN MODEL REDUCTION USING WELDNET:  
2           WINDOWED AUTOENCODERS FOR LEARNING DYNAMICS**

3           BIRAJ DAHAL\* AND OTHERS

4       **Abstract.** Many problems in science and engineering involve time-dependent, high dimensional  
5       datasets arising from complex physical processes, which are costly to simulate. In this work, we  
6       propose **WeldNet**: Windowed Encoders for Learning Dynamics, a data-driven nonlinear model re-  
7       duction framework to build a low-dimensional surrogate model for complex evolution systems. Given  
8       time-dependent training data, we split the time domain into multiple overlapping windows, within  
9       which nonlinear dimension reduction is performed by auto-encoders to capture latent codes. Once a  
10      low-dimensional representation of the data is learned, a propagator network is trained to capture the  
11      evolution of the latent codes in each window, and a transcoder is trained to connect the latent codes  
12      between adjacent windows. The proposed windowed decomposition significantly simplifies propa-  
13      gator training by breaking long-horizon dynamics into multiple short, manageable segments, while  
14      the encoders ensure consistency across windows. In addition to the algorithmic framework, we  
15      develop a mathematical theory establishing the representation power of WeldNet under the manifold  
16      hypothesis, justifying the success of nonlinear model reduction via deep autoencoder-based archi-  
17      tectures. Our numerical experiments on various differential equations indicate that WeldNet can  
18      capture nonlinear latent structures and their underlying dynamics, outperforming both traditional  
19      projection-based approaches and recently developed nonlinear model reduction methods.

20       **Key words.** model reduction, manifold learning, scientific machine learning, dynamical systems

21       **MSC codes.** 68Q25, 68R10, 68U05

22       **1. Introduction.** Many real-world applications in science and engineering in-  
23       volve large-scale, complex, and costly data simulations or inversions of physical pro-  
24       cesses. However, the high-dimensional nature of these models often creates over-  
25       whelming demands on computational resources. Model reduction plays a crucial role  
26       in addressing this challenge, which helps to reduce the data dimension and problem  
27       size [21, 6, 5, 52, 46, 9, 42].

28       Linear model reduction techniques for differential equations and dynamical sys-  
29       tems have been well-established in literature [5, 6, 52]. Classical projection-based  
30       model reduction approaches have demonstrated great success when the underlying  
31       model is linear and low-dimensional. Representative projection-based model reduction  
32       methods include Proper Orthogonal Decomposition (POD) [7, 40, 52], reduced-basis  
33       techniques [41, 46], the Principal Component Analysis (PCA) approach [33], rational  
34       interpolation [19, 3], Galerkin projection [23, 45], etc. Most projection-based model  
35       reduction methods rely on projecting high-dimensional models onto a low-dimensional  
36       linear subspace.

37       In real-world applications, many objects exhibit low-dimensional nonlinear struc-  
38       tures [49, 44]. Simple transformations, such as translations or rotations, place these  
39       objects on low-dimensional nonlinear manifolds, which cannot be efficiently captured  
40       by linear subspaces. Meanwhile, many physical processes are inherently nonlinear,  
41       including fluid dynamics, nonlinear optical phenomena, and shallow water wave prop-  
42       agation. When linear model reduction methods are applied to reduce evolutionary  
43       equations, their optimal performance is quantified by the Kolmogorov  $n$ -width [39].  
44       For diffusion-dominated problems, the Kolmogorov  $n$ -width decays rapidly as the  
45       subspace dimension increases, enabling linear methods to achieve significant success  
46       with well-established justification [2, 8]. Unfortunately, many differential equations,

---

\*School of Mathematics, Georgia Institute of Technology, Atlanta, GA, USA bdahal6@gatech.edu.

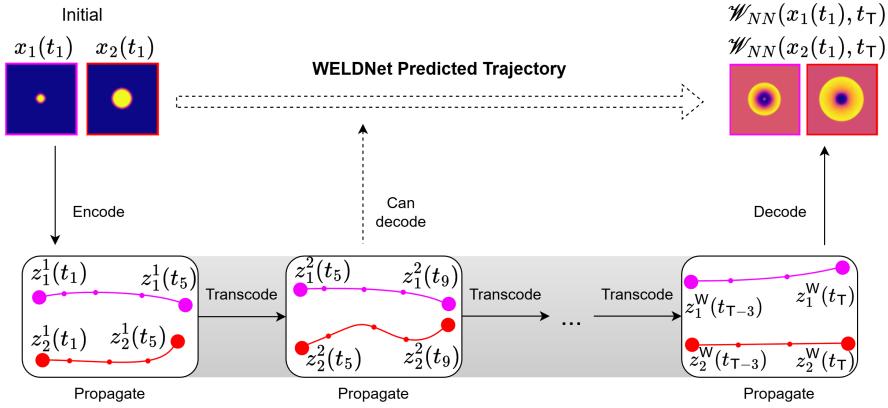


FIG. 1. *WeldNet*: Initial conditions are encoded, propagated within windows, transcoded between windows, and decoded. In this example, there are  $W$  windows and  $T$  time steps.  $z_j^i(t_k)$  denotes the latent space representation of  $x_j(t_k)$  according to window  $i$ .

47 particularly advection-dominated problems, display a slow decay of the Kolmogorov  
 48  $n$ -width. In these cases, linear model reduction methods require a sufficiently large  
 49 reduced dimensionality to achieve acceptable accuracy [34], which demonstrates the  
 50 limitations of linear model reduction methods in handling nonlinear structures.

51 The limitations of linear methods drive the development of nonlinear model reduction  
 52 techniques, which aim to leverage low-dimensional nonlinear structures in data  
 53 and models [11, 12, 55]. As deep learning has gained popularity in recent years, deep  
 54 learning methods have been increasingly applied to address nonlinear model reduction  
 55 [36, 54, 26, 17, 18, 17]. Many works utilize autoencoder [4] for dimension reduction,  
 56 and then learn the unknown physical process on latent variables by a neural network  
 57 [22, 37, 53, 36, 54, 26, 17, 18, 17, 16, 29]. However, when the solution manifold of  
 58 evolutionary equations has evolved for a long time, a global dimension reduction may  
 59 not well capture the effective latent parameters at all time points.

60 In this work, we propose Windowed Encoders for Learning Dynamics with Neural  
 61 Networks (WeldNet) for nonlinear model reduction. WeldNet uses encoder networks  
 62 to perform nonlinear dimension reduction and learn the evolution operation in the  
 63 encoded domain. The time domain is divided into several sequentially overlapping  
 64 subintervals called “windows”. In each window, we train an autoencoder to learn  
 65 the latent codes and a propagator to evolve the latent codes in time. Transcoder  
 66 networks are trained to connect two adjacent windows on their overlap, allowing an  
 67 initial condition to be encoded, evolved to the terminal time, and then decoded. After  
 68 training, WeldNet can be used to approximate the trajectory at any time given an  
 69 initial condition. Figure 1 diagrams our WeldNet method of trajectory learning for  
 70 time-dependent data.

71 The separation between dimension reduction and trajectory learning allows the  
 72 exploitation of existing nonlinear low-dimensional evolutionary structures in the data.  
 73 Figure 2 indicates the notation for each component of WeldNet. We consider the  
 74 trajectory manifold of an evolutionary process, denoted by  $\mathcal{M}([0, T])$ , which collects  
 75 all evolution trajectories in the time interval  $[0, T]$ . We split this time interval into  
 76  $W$  windows. For each window  $i \in \{1, 2, \dots, W\}$ ,  $\mathcal{E}^i$  and  $\mathcal{D}^i$  denote the encoder and  
 77 decoder networks,  $\mathcal{P}^i$  denotes the propagator network for the evolution of latent

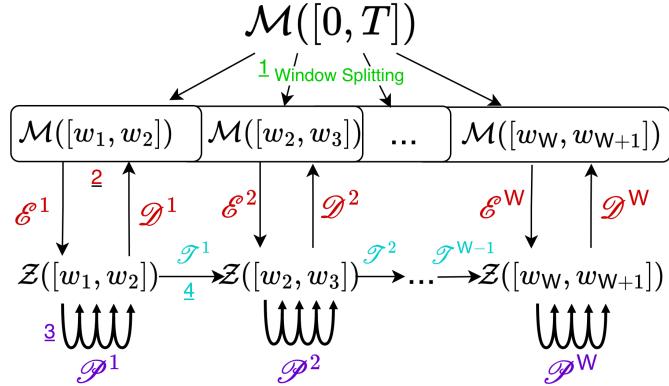


FIG. 2. Components of WeldNet: 1) Window splitting, 2) Autoencoder, 3) Propagator, 4) Transcoder.

78 codes, and  $\mathcal{T}_i$  (for  $i < W$ ) denotes the transcoder network from the  $i$ th window to the  
79 ( $i + 1$ )th window.

80 We validate WeldNet on a few examples, including the Burgers' equation, trans-  
81 port equation, Korteweg–De Vries (KdV) equation and 2D shallow-water equation.  
82 Furthermore, this paper provides a mathematical theory on the representation power  
83 of WeldNet to justify the success of nonlinear model reduction by deep learning under  
84 manifold hypothesis. Suppose the evolution operator satisfies a regularity assumption  
85 and all trajectories of this evolutionary process lie on a low-dimensional manifold.  
86 We prove that, WeldNet can approximate the trajectories of this evolutionary process  
87 up to arbitrary accuracy, if the encoder, decoder and propagator network architec-  
88 tures are properly set up. Our theory justifies the representation power of WeldNet  
89 for a large class of evolutionary processes, and provides a theoretical foundation for  
90 nonlinear model reduction using auto-encoder-based methods.

91 The main contributions of this paper are:

- 92 1. **WeldNet framework.** We propose a windowed autoencoder–propagator  
93 architecture for nonlinear model reduction. The method performs dimension  
94 reduction and latent-space trajectory learning by dividing the time domain  
95 into overlapping windows connected via transcoder networks.
- 96 2. **Efficient trajectory prediction.** WeldNet enables end-to-end approxima-  
97 tion of system trajectories by encoding an initial condition, evolving the latent  
98 representations across windows, and decoding the result at any desired time.
- 99 3. **Theoretical foundation.** Under the manifold hypothesis, we establish a  
100 representation theory to show that WeldNet can express evolutionary pro-  
101 cesses with low-dimensional structures to arbitrarily high accuracy when the  
102 encoder, decoder, and propagator networks are appropriately designed.
- 103 4. **Empirical validation.** The effectiveness of WeldNet is demonstrated on  
104 several nonlinear PDEs, including the Burgers', transport, KdV, and 2D  
105 shallow-water equations.

106 **Organization.** We first introduce our WeldNet model in Section 2. A repre-  
107 sentation theory of WeldNet is presented in Section 3, and comprehensive numerical  
108 experiments are given in Section 4. Finally, we conclude in Section 5.

109     **Notations.** For any  $n \in \mathbb{N}$ , we denote  $[n] = \{1, \dots, n\}$ . For  $\mathbf{x} \in \mathbb{R}^n$ , we  
 110   denote  $\|\mathbf{x}\|_{\mathbb{R}^n} = \|\mathbf{x}\|_2$  and  $\|\mathbf{x}\|_\infty = \max_{i \in [n]} x_i$  where  $x_i$  is the  $i$ th component of  
 111    $\mathbf{x}$ . For any  $\mathbf{x} \in \mathbb{R}^n$ , we denote the ReLU function  $\sigma : \mathbb{R}^n \rightarrow \mathbb{R}^n$  by  $\sigma(\mathbf{x}) =$   
 112    $(\max(x_i, 0))_{i=1}^n$ . For any function  $f : A \rightarrow B$  defined on sets  $A, B$ , we denote  
 113    $\|f\|_{L^\infty(A; B)} = \sup_{\mathbf{x} \in A} \|f(\mathbf{x})\|_B$  where  $\|\cdot\|_B$  is a norm on  $B$ . Given a finite set  $A$ ,  
 114   we denote the cardinality of  $A$  by  $|A|$ . We use  $\bigcirc_{k=1}^K f_k$  to denote the composition  
 115    $f_K \circ f_{K-1} \circ \dots \circ f_1$ . In particular,  $\bigcirc_{k=1}^K f$  denotes the composition of  $f$  for  $K$  times.

116     **2. WeldNet for Model Reduction.** In this section, we present our WeldNet  
 117   model, which operates on time-dependent trajectory data collected from an evolu-  
 118   tionary process. In science and engineering applications, the solution trajectory in an  
 119   evolutionary process often depends on few parameters in the initial condition or in  
 120   the evolution equation [56, 26]. When the trajectories of this evolutionary process lie  
 121   on a low-dimensional manifold, WeldNet parametrizes the trajectory manifold by a  
 122   low-dimensional latent code and builds a surrogate evolutionary model in the latent  
 123   space.

We consider an evolutionary process in  $\mathbb{R}^D$ , whose initial states can be (locally) parameterized by a small number of parameters. Suppose the initial states are supported on a set  $\mathcal{M}(0) \subseteq \mathbb{R}^D$ . We are interested in learning an evolutionary process driven by an unknown continuous **time-evolution operator**  $\mathcal{F} : \mathcal{M}(0) \times [0, T] \rightarrow \mathbb{R}^D$ . For simplicity, we denote

$$\mathbf{x}(t) = \mathcal{F}(\mathbf{x}(0), t).$$

124     We denote  $\mathcal{M}(t) = \mathcal{F}(\mathcal{M}(0), t) = \{\mathcal{F}(\mathbf{x}(0), t) : \mathbf{x}(0) \in \mathcal{M}(0)\}$  as the initial data  
 125   evolved by  $t$  time units, and  $\mathcal{M}([a, b]) = \{\mathcal{M}(t) : t \in [a, b]\}$  denotes the collection of  
 126   state manifold in the time interval  $[a, b]$ . We call the set  $\mathcal{M}(0)$  the initial manifold  
 127   and the set  $\mathcal{M}([0, T])$  the **trajectory manifold**. We will assume that  $\mathcal{M}([0, T])$   
 128   is a  $d$ -dimensional smooth manifold smoothly isometrically embedded in  $\mathbb{R}^D$ . For  
 129   convenience, we extend the domain of the time-evolution operator  $\mathcal{F}$  to cover all  
 130   of  $\mathcal{M}([0, T])$  and all relevant times, such that for all  $\mathbf{x}(s) = \mathcal{F}(\mathbf{x}(0), s) \in \mathcal{M}(s)$ ,  
 131   we assume that the dynamic satisfies  $\mathcal{F}(\mathbf{x}(s), t) = \mathcal{F}(\mathbf{x}(0), s + t) = \mathbf{x}(s + t)$  for  
 132    $t \in [0, T - s]$ .

133     In addition, we extend the domain of  $\mathcal{F}$  to all of  $\mathcal{M}([0, T]) \times [0, T]$  in the following  
 134   way: if  $\mathbf{x} \in \mathcal{M}(s)$  then define  $\mathcal{F}(\mathbf{x}, t) = \mathcal{F}(\mathbf{x}, \max(T - s, t))$ . Note that this extension  
 135   preserves the Lipschitz constant.

136     **2.1. Data Collection.** In applications, one can measure the trajectory of the  
 137   above evolutionary process and collect trajectory data for multiple initial states. Sup-  
 138   pose  $N$  initial states are randomly sampled from a probability measure  $\rho(0)$  supported  
 139   on  $\mathcal{M}(0)$ :

$$\{\mathbf{x}_1(0), \dots, \mathbf{x}_N(0)\} \stackrel{\text{iid}}{\sim} \rho(0).$$

140     Let  $0 = t_1 < t_2 < \dots < t_T = T$  be a time grid, denoted  $\mathbb{T} = \{t_1, \dots, t_T\}$ .

$$141 \quad \{\mathbf{x}_n\}_{n=1}^N := \{(\mathbf{x}_n(t_k))_{k=1}^T\}_{n=1}^N \in \mathbb{R}^{N \times T \times D}, \text{ where } \mathbf{x}_n(t_k) = \mathcal{F}(\mathbf{x}_n(0), t_k).$$

142     Given this dataset, our goal is to construct a low-dimensional surrogate model  $\mathcal{W}_{\text{model}} :$   
 143    $\mathcal{M}(0) \times \mathbb{T} \rightarrow \mathbb{R}^D$  such that

$$144 \quad \mathcal{W}_{\text{model}}(\mathbf{x}(0), t_k) \approx \mathcal{F}(\mathbf{x}(0), t_k), \forall \mathbf{x}(0) \in \mathcal{M}(0), t_k \in \mathbb{T}.$$

145     After training, one can predict the solution trajectory for a new initial condition  
 146   sampled from  $\rho(0)$  by evolving the surrogate model in the low-dimensional latent  
 147   space.

Notation	Explanation
$[w_i, w_{i+1}]$	Time range in window $i$
$\mathbb{T}$	Time grid: $\{t_1, \dots, t_T\}$
$\mathbb{T}^i$	$\mathbb{T} \cap [w_i, w_{i+1}]$
$T_i$	Size of $\mathbb{T} \cap (w_i, w_{i+1})$
$\mathcal{F}$	Evolution operator
$\mathcal{E}_*^i / \mathcal{D}_*^i$	Window $i$ 's oracle encoder/decoder
$\mathcal{P}_*^i$	Window $i$ 's oracle propagator $= \mathcal{E}_*^i \circ \mathcal{F}(\cdot, \Delta t) \circ \mathcal{D}_*^i$ .
$\mathcal{T}_*^i$	Window $i$ 's oracle transcoder $= \mathcal{E}_*^{i+1} \circ \mathcal{E}_*^i$ .
$\mathcal{E}_{\text{NN}}^i / \mathcal{D}_{\text{NN}}^i / \mathcal{P}_{\text{NN}}^i / \mathcal{T}_{\text{NN}}^i$	Window $i$ 's oracle encoder/decoder/propagator/transcoder
$\mathcal{M}([a, b])$	Trajectory manifold from time $a$ to $b$
$\mathcal{Z}([a, b])$	Latent space from time $a$ to $b$

TABLE 1  
Notation overview.

149    **2.2. WeldNet Training.** Training a WeldNet model involves four stages, shown  
 150 in Figure 2: 1) Window splitting, 2) Autoencoder training, 3) Propagator training,  
 151 4) Transcoder training. We provide details of each stages below. We also provide an  
 152 overview of notations is provided in Table 1.

153    **1) Window splitting.** A WeldNet model with  $W$  windows divides the time  
 154 domain into  $W$  sequentially overlapping windows, with the  $i$ -th window denoted by  
 155  $[w_i, w_{i+1}]$  and  $\bigcup_{i=1}^W [w_i, w_{i+1}] = [0, T]$ . Suppose  $\{w_i\}_{i=1}^{W+1} \subset \{t_k\}_{k=1}^T$ , i.e., the end-  
 156 points of each window are on the time grid.

157    **2) Autoencoder and Propagator training.** WeldNet first trains  $W$  autoen-  
 158 coders and  $W$  propagators, where the  $i$ th autoencoder aims to learn a representation  
 159 of the manifold  $\mathcal{M}([w_i, w_{i+1}])$ , and the  $i$ th propagator will predict the displacement  
 160 required to evolve the latent code from one time discretization point to the subsequent  
 161 one in latent manifold.

162    Fixing a latent space dimension  $d \in \mathbb{N}$ , We denote the domain of the decoders  
 163 and propagators, also known as the **latent space**, by  $\mathcal{Z}([0, T]) = [0, 1]^d \times [0, T]$ . We  
 164 also use the notation  $\mathcal{Z}(t) = [0, 1]^d \times \{t\}$  and notation  $\mathcal{Z}([a, b]) = [0, 1]^d \times [a, b]$ . Note  
 165 that we always incorporate time in our latent codes to track the time information.

166    The autoencoder reconstruction loss for the  $i$ th window is

$$(2.1) \quad \mathbf{L}_{\text{ae}}^i(\mathcal{E}^i, \mathcal{D}^i) = \frac{1}{N|\mathbb{T} \cap [w_i, w_{i+1}]|} \sum_{n=1}^N \sum_{t_k \in \mathbb{T} \cap [w_i, w_{i+1}]} \left\| \mathcal{D}^i(\mathcal{E}^i(\mathbf{x}_n(t_k))) - \mathbf{x}_n(t_k) \right\|_{\mathbb{R}^d}^2,$$

168    where  $(\mathcal{E}^i, \mathcal{D}^i)$  is an encoder-decoder pair in the  $i$ th window. The propagator loss for  
 169 the  $i$ th window (denoted  $\mathbf{L}_{\text{prop}}^i(\mathcal{E}^i, \mathcal{P}^i)$ ) is

$$(2.2) \quad \frac{1}{N|\mathbb{T} \cap [w_i, w_{i+1}]|} \sum_{n=1}^N \sum_{t_k \in \mathbb{T} \cap [w_i, w_{i+1}]} \left| \mathcal{P}^i(\mathcal{E}^i(\mathbf{x}_n(t_k))) - \mathcal{E}^i(\mathbf{x}_n(t_{k+1})) \right|_{\mathbb{R}^{d+1}}^2,$$

171    where  $\mathcal{P}^i$  denotes an one-step propagator for the latent code in the  $i$ th window.

172    We train the encoder, decoder, and propagator together using the following ob-  
 173 jective function:

$$(2.3) \quad (\mathcal{E}_{\text{NN}}^i, \mathcal{D}_{\text{NN}}^i, \mathcal{P}_{\text{NN}}^i) \in \arg \min_{\mathcal{E} \in \mathcal{F}_{\text{NN}}^{\mathcal{E}}, \mathcal{D} \in \mathcal{F}_{\text{NN}}^{\mathcal{D}}, \mathcal{P} \in \mathcal{F}_{\text{NN}}^{\mathcal{P}}} \mathbf{L}_{\text{ae}}^i(\mathcal{E}^i, \mathcal{D}^i) + \lambda \mathbf{L}_{\text{prop}}^i(\mathcal{E}^i, \mathcal{P}^i),$$

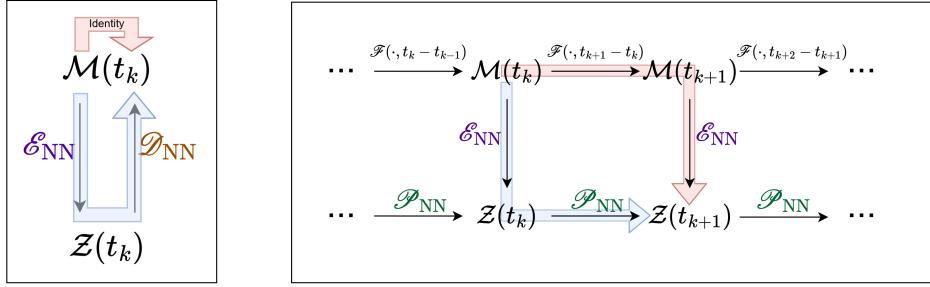


FIG. 3. Diagram of the autoencoder loss (left) and propagator loss (right), where the loss is the MSE between the output of the red arrow and the blue arrow.

175 where  $\mathcal{F}_{\text{NN}}^{\mathcal{D}}$ ,  $\mathcal{F}_{\text{NN}}^{\mathcal{E}}$ , and  $\mathcal{F}_{\text{NN}}^{\mathcal{P}}$  are network classes for decoder, encoder, and propagator,  
176 respectively, and  $\lambda > 0$  is a hyperparameter. In this work, we use  $\lambda = 0.1$  for all  
177 experiments. Figure 3 diagrams the autoencoder and propagator losses (the loss is  
178 the mean squared error between applying the functions in the red path and applying  
179 the functions in the blue path).

180 **3) Propagator finetuning.** After we train the autoencoder and propagator  
181 together, we then finetune the propagator in order to reduce the accumulation of  
182 error that occurs when applying  $\mathcal{P}^i$  to propagate a latent code over multiple time  
183 steps. Specifically, we freeze the encoder and train the propagator with the objective  
184 (where  $k(w_i)$  is the time index of  $w_i$ , i.e. the first time step in window  $i$ ):

$$(2.4) \quad \mathcal{P}_{\text{NN}}^i = \arg \min_{\mathcal{P}^i \in \mathcal{F}_{\text{NN}}^{\mathcal{P}}} \frac{1}{N \mathbf{T}_i} \sum_{n=1}^N \sum_{s=1}^{\mathbf{T}_i} \| (\bigcirc_{k=1}^s \mathcal{P}^i) (\mathcal{E}^i (\mathbf{x}_n(t_{k(w_i)}))) - \mathcal{E}^i (\mathbf{x}_n(t_{k(w_i)+s})) \|_{\mathbb{R}^{d+1}}^2,$$

186 where we denote  $\mathbf{T}_i = |\mathbb{T} \cap (w_i, w_{i+1}]|$ .

187 The loss for the finetuning of propagator network is illustrated in the Figure 4.

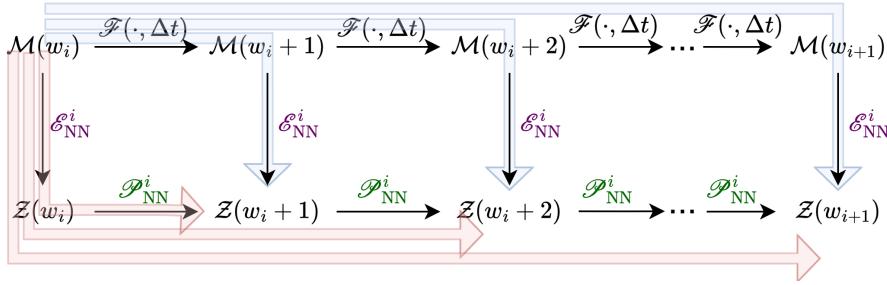


FIG. 4. Diagram for propagator accumulation loss. The loss is the average of the mean squared error between the blue and red arrows that point to the same symbol.

188 **4) Transcoder training.** We have trained autoencoders and propagators on  
189 each time window separately. In order to connect the windows, we train transcoder  
190 networks on the overlap between the windows. The goal of the  $i$ th transcoder is to  
191 connect the codes at the end of window  $i$  with the codes at the beginning of window  
192  $i + 1$ . To train the transcoder, we will propagate codes from the beginning of window

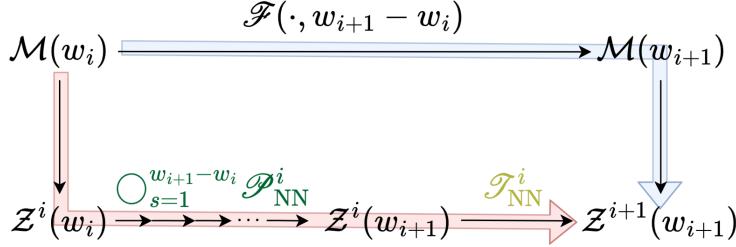


FIG. 5. Diagram of the transcoder loss, where the loss is the MSE between the output of the red arrow and the blue arrow.

193  $i$  to the end. Specifically, we use the objective  $\mathcal{T}_{\text{NN}}^i =$   
(2.5)

$$194 \quad \arg \min_{\mathcal{T}^i \in \mathcal{F}_{\text{NN}}^{\mathcal{T}}} \frac{1}{N} \sum_{n=1}^N \left\| \mathcal{T}^i \left( \left( \bigcirc_{s=1}^{|T \cap (w_i, w_{i+1})|} \mathcal{P}^i \right) (\mathcal{E}_{\text{NN}}^i(\mathbf{x}_n(w_i))) \right) - \mathcal{E}_{\text{NN}}^{i+1}(\mathbf{x}_n(w_{i+1})) \right\|_{\mathbb{R}^{d+1}}^2$$

195 where  $\mathcal{F}_{\text{NN}}^{\mathcal{T}}$  denotes the network class for transcoders. Figure 5 diagrams the loss  
196 function for the transcoder.

197 **2.3. WeldNet Inference.** After WeldNet is trained, one can predict the tra-  
198 jectory given a new initial state  $\mathbf{x}(0)$ . The inference involves encoding the initial state  
199 to the latent code, evolving the latent code in time, and decoding at the final time.  
200 This inference procedure can be represented by our WeldNet model, denoted by  $\mathcal{W}_{\text{NN}}$ ,  
201 such that, for any  $\mathbf{x}(0) \in \mathcal{M}(0)$  and  $t_k \in T \cap (w_i, w_{i+1}]$ , WeldNet gives rise to (where  
202  $k(w_i)$  is the time grid index of the time  $w_i$ )

$$203 \quad (2.6) \quad \begin{aligned} \mathcal{W}_{\text{NN}}(\mathbf{x}(0), t_k) := \\ \mathcal{D}_{\text{NN}} \circ \bigcirc_{\ell=1}^{k-k(w_i)} \mathcal{P}_{\text{NN}}^i \circ \bigcirc_{j=1}^{i-1} (\mathcal{T}_{\text{NN}}^j \circ \bigcirc_{\ell=1}^{|T \cap (w_j, w_{j+1})|} \mathcal{P}_{\text{NN}}^j) \circ \mathcal{E}_{\text{NN}}^1(\mathbf{x}(0)). \end{aligned}$$

204 Figure 1 shows the process of encoding, propagating, and decoding for inference  
205 with WeldNet.

206 **3. Approximation Theory.** In this section, we prove the approximation ability  
207 of WeldNet for a large class of evolutionary operators. The low-dimensional struc-  
208 ture in the trajectory manifold results in the existence of oracle maps for charts,  
209 time-evolution operator, and transition maps which can be approximated by the au-  
210 toencoder, propagator, and transcoder networks respectively.

### 211    3.1. Preliminaries.

#### 212     3.1.1. Neural Networks.

In this work, we consider ReLU networks as follows:

DEFINITION 3.1 (Feedforward Neural Network (FNN)). Let  $L \in \mathbb{N}$  and suppose  
 $\mathfrak{W}_1, \dots, \mathfrak{W}_{L+1}$  are weight matrices and  $\mathfrak{b}_1, \dots, \mathfrak{b}_{L+1}$  are bias vectors with  $\mathfrak{W}_\ell \in \mathbb{R}^{d_{\ell-1} \times d_\ell}$  and  $\mathfrak{b}_\ell \in \mathbb{R}^{d_\ell}$  for all  $\ell \in [L+1]$ . Let  $\theta = (\mathfrak{W}_1, \mathfrak{b}_1, \dots, \mathfrak{W}_{L+1}, \mathfrak{b}_{L+1})$ . We  
define a **feedforward neural network (FNN)** with weights  $\theta$ , denoted  $f_{\text{NN}}$ , as a  
function of the form:

$$f_{\text{NN}}(x) = \mathfrak{W}_{L+1} \sigma(\mathfrak{W}_L \sigma(\dots \mathfrak{W}_1 x + \mathfrak{b}_1) \dots + \mathfrak{b}_L) + \mathfrak{b}_{L+1}.$$

213 We say the **depth** and **width** of  $f_{\text{NN}}$  are  $L$  and  $W = \max_\ell d_\ell$  respectively.

214 A FNN with depth  $L$  and width  $W$  is therefore a composition of  $L + 1$  functions, all  
 215 but the last one of which involves an affine transform with output dimension at most  
 216  $W$ . With the (maximum) depth and width specified, we can define a function class  
 217 of feedforward neural networks.

218 **DEFINITION 3.2** (FNN Class). *Let  $d_{in}, d_{out}, L, W \in \mathbb{N}$ . The class of feedforward  
 219 neural networks (**FNN class**) with  $d_{in}$  inputs,  $d_{out}$  outputs, depth at most  $L$ , and  
 220 width at most  $W$  is denoted  $\mathcal{F}_{\text{NN}}(d_{in}, d_{out}, L, W)$ .*

221 **3.1.2. Manifolds.**

222 **DEFINITION 3.3** (Lipschitz Function). *Let  $A, B$  be metric spaces with metrics  $d_A$   
 223 and  $d_B$  respectively, and let  $L > 0$ . A function  $f : A \rightarrow B$  is called  **$L$ -Lipschitz** if*

$$224 \quad \sup_{x \neq y \in A} \frac{\|f(x) - f(y)\|_{d_B}}{\|x - y\|_{d_A}} \leq L,$$

225 and we define the **Lipschitz constant** of  $f$ , denoted  $\text{Lip}_A(f)$ , as the smallest  $L$  for  
 226 which  $f : A \rightarrow B$  is  $L$ -Lipschitz.

227 **DEFINITION 3.4** (Manifold). *A  $d$ -dimensional **manifold**  $\mathcal{M}$  is a topological space  
 228 endowed with a collection of **charts**  $(\phi_i : U_i \subseteq \mathcal{M} \rightarrow \tilde{U}_i \subseteq \mathbb{R}^d)_{i \in \mathcal{I}}$  (for some index  
 229 set  $\mathcal{I}$ ) such that for all  $i \in \mathcal{I}$ ,  $\phi_i$  is a homeomorphism between the open sets  $U_i$  and  
 230  $\tilde{U}_i$ , and we have  $\bigcup_{i \in \mathcal{I}} U_i = \mathcal{M}$ . We say  $\mathcal{M}$  is a **smooth** manifold if for all charts  $\phi_i$   
 231 and  $\phi_j$  on  $\mathcal{M}$ , the function*

$$232 \quad \phi_j \circ \phi_i^{-1} : \tilde{U}_i \cap \tilde{U}_j \subseteq \mathbb{R}^d \rightarrow \tilde{U}_i \cap \tilde{U}_j \subseteq \mathbb{R}^d$$

233 is smooth in the usual Euclidean sense. Finally, for any  $f : \mathcal{M} \rightarrow \mathbb{R}^k$ , we say that  $f$   
 234 has **injective derivative** if for each chart  $\phi_i : U_i \rightarrow \tilde{U}_i$ , we have the function  $f \circ \phi_i^{-1}$   
 235 is differentiable in the Euclidean sense, and its derivative is nonzero on its domain.

236 **DEFINITION 3.5** (Smooth Embedding). *We say that a manifold is **smoothly em-**  
 237 **bedded** in  $\mathbb{R}^D$  if  $\mathcal{M} \subseteq \mathbb{R}^D$  and the inclusion function  $\iota : \mathcal{M} \rightarrow \mathbb{R}^D$  is a smooth  
 238 homeomorphism onto its image with injective derivative.*

239 More details on manifolds can be found in standard texts on differential geometry  
 240 [25]. Next, we introduce a regularity/curvature parameter for embedded submanifolds  
 241 of Euclidean space known as reach. Informally speaking, the reach of a manifold  
 242 describes the curvature of the embedded manifold in the ambient Euclidean space. It  
 243 is the radius of the largest “ball” (interior of a hypersphere) you can roll around the  
 244 manifold without crossing it.

245 **DEFINITION 3.6** (Reach [15]). *Let  $\mathcal{M} \subseteq \mathbb{R}^D$  be an isometrically embedded mani-  
 246 fold. The reach of  $\mathcal{M}$ , denoted  $\tau(\mathcal{M})$ , is*

$$247 \quad \inf\{r > 0 : \exists x \neq y \in \mathcal{M}, v \in \mathbb{R}^D \\ 248 \quad \text{such that } r = \|x - v\| = \|y - v\| = \text{dist}(v, \mathcal{M})\}.$$

249 Linear subspaces have reach  $\infty$ , and the reach of a hypersphere equals to its radius. It  
 250 is well-known that compact embedded manifolds have positive reach [50]. The larger  
 251 the reach, the easier the manifold is to be represented with neural networks, and that  
 252 will affect the bounds in neural network approximation theory [10, 30].

253     **3.2. Main Theorems.** According to Section 2.1, the time-evolution operator  
 254      $\mathcal{F}$  is defined on the set  $\mathcal{M}([0, T]) \times [0, T]$ . We will assume that  $\mathcal{F}$  is Lipschitz for all  
 255     times:

256     *Assumption 1.* Assume  $\mathcal{F}$  is Lipschitz, i.e.  $\text{Lip}(\mathcal{F}) :=$

$$257 \quad \sup_{\substack{t \in [0, T], \\ s \in [0, T-t]}} \text{Lip}_{\mathcal{M}(t)}(\mathcal{F}(\cdot, s)) = \sup_{\substack{t \in [0, T], \\ s \in [0, T-t]}} \sup_{\substack{\mathbf{x}(t) \neq \mathbf{y}(t) \\ \in \mathcal{M}(t)}} \frac{\|\mathcal{F}(\mathbf{x}(t), s) - \mathcal{F}(\mathbf{y}(t), s)\|}{\|\mathbf{x}(t) - \mathbf{y}(t)\|} < \infty.$$

258 In this work, we will consider evolutionary processes whose initial conditions are  
 259 sampled from an embedded  $d$ -dimensional manifold. This means that each initial  
 260 condition can be described (locally) by  $d$  parameters. More specifically, we will assume  
 261 that the collection of trajectories at various segmented times form a  $(d+1)$ -dimensional  
 262 manifold, with  $d$ -dimensions in parameter space and one dimension given by time. For  
 263 convenience, we will use the notation  $\mathcal{Z}(t) = [0, 1]^d \times \{t\}$  and  $\mathcal{Z}([a, b]) = [0, 1]^d \times [a, b]$ .

264     *Assumption 2* (Segmented Manifold). Suppose there exist  $0 = s_1 < s_2 < \dots <$   
 265  $s_{S+1} = T$  such that for all  $i \in [S]$ , the subset  $\mathcal{M}([s_i, s_{i+1}])$  is a compact  $(d+1)$ -  
 266 dimensional smooth manifold smoothly embedded in  $\mathbb{R}^D$  with reach  $\tau(\mathcal{M}([s_i, s_{i+1}])) > 0$ .  
 267 In addition, for all  $i \in [S]$  we assume there is a smooth function  $\mathcal{D}_*^i$  that maps from  
 268 a Euclidean space to the manifold with smooth inverse  $\mathcal{E}_*^i = (\mathcal{D}_*^i)^{-1}$ :

$$269 \quad \mathcal{E}_*^i : \mathcal{M}([s_i, s_{i+1}]) \rightarrow \mathcal{Z}([s_i, s_{i+1}]), \quad \mathcal{D}_*^i : \mathcal{Z}([s_i, s_{i+1}]) \rightarrow \mathcal{M}([s_i, s_{i+1}]),$$

270 such that for any  $t \in [s_i, s_{i+1}]$ ,  $\mathcal{M}(t) = \mathcal{D}_*^i(\mathcal{Z}(t))$ .

271     In Assumption 2,  $\mathcal{E}_*^i$  and  $\mathcal{D}_*^i$  serve as the oracle encoder and decoder for the  
 272 trajectory manifold for each segment  $i \in [S]$ . We denote the Lipschitz constant of the  
 273 oracle encoder on each time segment by  $\text{Lip}_{\mathcal{E}^i} = \sup_{t \in [s_i, s_{i+1}]} \text{Lip}_{\mathcal{M}(t)}(\mathcal{E}^i)$  and the  
 274 oracle decoder in each time segment by  $\text{Lip}_{\mathcal{D}^i} = \sup_{t \in [s_i, s_{i+1}]} \text{Lip}_{\mathcal{Z}(t)}(\mathcal{D}^i)$ .

275     We will first establish the result where we set the windows to be exactly aligned  
 276 with the manifold segments in Assumption 2. In other words, we will choose  $W = S$   
 277 and define  $0 < w_1 < \dots < w_{W+1} = T$  such that  $w_i = s_i$  for all  $i \in [W]$ . We  
 278 construct encoder, decoder, and propagator networks over each window (which is  
 279 the same as a segment). In order to translate between adjacent windows, we use  
 280 transcoder networks. The inference by WeldNet for any initial state  $\mathbf{x}(0) \in \mathcal{M}(0)$  can  
 281 be represented in (2.6).

282     We introduce a latent time-evolution map, defined on the domain  $\mathcal{Z}([w_i, w_{i+1}])$   
 283 denoted as  $\mathcal{P}_*^i$  such that  $\mathcal{P}_*^i(\mathbf{z}, t) = \mathcal{E}_*^i(\mathcal{F}(\mathcal{D}_*^i(\mathbf{z}), t))$  for  $t \leq w_{i+1} - w_i$ , which gives  
 284 the evolution of latent code for  $\mathcal{F}$  within window  $i$ :  $\mathcal{D}_*^i \circ \mathcal{P}_*^i \circ \mathcal{E}_*^i = \mathcal{F}$ . The evolution  
 285 of  $\mathcal{P}_*^i$  operates in latent space, yet it matches the dynamics of the time-evolution  
 286 operator  $\mathcal{F}$  on a section of the trajectory manifold. We will call  $\mathcal{F}$  the time-evolution  
 287 operator, and  $\mathcal{P}_*^i$  the **oracle propagator** for the  $i$ th window.

288     We are primarily interested in the case where the latent codes follow an ODE  
 289 within each segment:

290     *Assumption 3* (Latent Dynamics). Using the notation in Assumption 2, suppose  
 291 for all  $i \in [S]$ , there is a Lipschitz function  $g^i : \mathcal{Z}([s_i, s_{i+1}]) \rightarrow \mathbb{R}^{d+1}$  such that for  
 292 all  $\mathbf{z}(s_i) \in \mathcal{Z}(s_i)$ , the latent code  $\mathbf{z}(t) = \mathcal{P}_*^i(\mathbf{z}(s_i), t - s_i)$  for  $t \in [s_i, s_{i+1}]$  satisfies  
 293  $\frac{\partial \mathbf{z}}{\partial t}(t) = g^i(\mathbf{z}(t))$ .

294     The ODE in Assumption 3 models the dynamics of the latent code. In particular,  
 295 the  $(d+1)$ th coordinate of the latent code  $\mathbf{z}(t)$  is time, i.e.  $\mathbf{z}(t)_{d+1} = t$  and therefore

296 the  $(d+1)$ th coordinate of  $g^i$  is 1, i.e.  $g_{d+1}^i = 1$ . We denote the Lipschitz constant  
 297 of  $g^i$  (for the  $i$ th segment) as  $\text{Lip}(g^i) := \sup_{t \in [s_i, s_{i+1}]} \text{Lip}_{\mathcal{Z}(t)}(g^i)$  and the Lipschitz  
 298 constant of  $g = \{g^i\}_{i=1}^S$  (among all segments) as  $\text{Lip}(g) = \sup_{i \in [W]} \text{Lip}(g^i)$ .

299 Our first result (presented in Lemma 3.7 below) gives an approximation guarantee  
 300 for WeldNet in the latent dynamics setting, assuming that the windows are set to be  
 301 equal to the manifold segments.

302 **LEMMA 3.7.** *Suppose Assumptions 1, 2, and 3 hold. Assume the time grid  $\mathbb{T}$  is  
 303 uniform such that  $t_{k+1} - t_k = \Delta t$  for all  $k \in [\mathsf{T}-1]$ . Let  $\epsilon > 0$ , and suppose there are  
 304  $W = S$  windows such that for all  $i \in [W]$ , we have  $w_i = s_i$  and each window has number  
 305 of time steps  $|\mathbb{T}^i| > 1 + \text{Lip}(g)\|g\|_{L^\infty}(w_{i+1} - w_i)^2 e^{\text{Lip}(g)(w_{i+1} - w_i)} \epsilon^{-1}$ . Then for each  
 306  $i \in [W]$  there exist an encoder network  $\mathcal{E}_N^i \in \mathcal{F}_{\text{NN}}(D, d+1, L_{\mathcal{E}^i}, W_{\mathcal{E}^i})$ , a decoder  
 307 network  $\mathcal{D}_N^i \in \mathcal{F}_{\text{NN}}(d+1, D, L_{\mathcal{D}^i}, W_{\mathcal{D}^i})$ , a propagator network  $\mathcal{P}_N^i \in \mathcal{F}_{\text{NN}}(d+  
 308 1, d+1, L_{\mathcal{P}^i}, W_{\mathcal{P}^i})$ , and a transcoder network (for  $i < W$ )  $\mathcal{T}_N^i \in \mathcal{F}_{\text{NN}}(d+1, d+  
 309 1, L_{\mathcal{T}^i}, W_{\mathcal{T}^i})$  such that for any  $k \in [\mathsf{T}]$ , the WeldNet  $\mathcal{W}_{\text{NN}}$  given in (2.6) guarantees*

$$310 \quad \sup_{\mathbf{x}(0) \in \mathcal{M}(0)} \|\mathcal{W}_{\text{NN}}(\mathbf{x}(0), t_k) - \mathcal{F}(\mathbf{x}(0), t_k)\|_{\mathbb{R}^D} < \epsilon.$$

311 The network parameters are

$$\begin{aligned} 312 \quad L_{\mathcal{E}^i} &= O(\log^2(S/\epsilon)), \quad W_{\mathcal{E}^i} = O(D(S/\epsilon)^{d+1}) \\ 313 \quad L_{\mathcal{D}^i} &= O(\log(S/\epsilon)), \quad W_{\mathcal{D}^i} = O(D(S/\epsilon)^{d+1}) \\ 314 \quad L_{\mathcal{P}^i} &= O(\log(S/\epsilon)), \quad W_{\mathcal{P}^i} = O((S/\epsilon)^{d+1}) \\ 315 \quad L_{\mathcal{T}^i} &= O(\log(S/\epsilon)), \quad W_{\mathcal{T}^i} = O((S/\epsilon)^d) \end{aligned}$$

316 where we hide constants depending on (for all  $i \in [W]$ )  $\|g_i\|_{L^\infty}$ ,  $\text{Lip}(g_i)$ ,  $\log(D)$ ,  $d$ ,  
 317  $\tau(\mathcal{M}([s_i, s_{i+1}]))$ ,  $\max_i \text{Lip}(\mathcal{E}_*^i)$ ,  $\max_i \text{Lip}(\mathcal{D}_*^i)$ ,  $\text{Lip}(\mathcal{F})$ , the volume of  $\mathcal{M}([s_i, s_{i+1}])$ ,  
 318 and  $\sup_{\mathbf{x} \in \mathcal{M}([s_i, s_{i+1}])} \|\mathbf{x}\|_{\mathbb{R}^D}$ .

319 Lemma 3.7 provides a representation guarantee of WeldNet, and the network size  
 320 scales crucially with the intrinsic dimension  $d$  instead of the data dimension  $D$ , which  
 321 demonstrates the efficiency of model reduction. The proof is deferred to Appendix  
 322 B.1.

323 Lemma 3.7 requires the number of windows to be equal to the number of segments.  
 324 One major advantage of WeldNet is the ability to train with a number of windows  
 325 potentially higher than the number of segments. Fortunately, the approximation  
 326 guarantee of Lemma 3.7 can be extended to this scenario, as long as the windows  
 327 partition each segment. This is the subject of the next theorem, which follows from  
 328 Lemma 3.7.

329 **DEFINITION 3.8.** *We say a sequence  $\{w_i\}_{i=1}^{W+1}$  with  $0 = w_1 < w_2 < \dots < w_{W+1} =$   
 330 **T subdivides the segments**  $\{s_j\}_{j=1}^{S+1}$  if for all  $i \in [W]$ , there is a unique  $j \in [S]$  such  
 331 that  $[w_i, w_{i+1}] \subseteq [s_j, s_{j+1}]$ . We denote by  $\pi : [W] \rightarrow [S]$  that assigns each window  
 332 index to the unique segment index it lies within.*

333 **THEOREM 3.9.** *Suppose Assumptions 1, 2, and 3 hold. Assume the time grid  $\mathbb{T}$  is  
 334 uniform such that  $t_{k+1} - t_k = \Delta t$  for all  $k \in [\mathsf{T}-1]$ . Let  $\epsilon > 0$  and suppose there are  
 335  $W \geq S$  windows so that  $0 = w_1 < w_2 < \dots < w_{W+1}$  subdivides the segments  $\{s_j\}_{j=1}^{S+1}$ .  
 336 Additionally, suppose for all  $i \in [W]$ , each window has number of time steps  $|\mathbb{T}^i| >$   
 337  $1 + \text{Lip}(g)\|g\|_{L^\infty}(w_{i+1} - w_i)^2 e^{\text{Lip}(g)(w_{i+1} - w_i)} \epsilon^{-1}$ . Then for each  $i \in [W]$  there exist an*

338 encoder  $\mathcal{E}_{\text{NN}}^i \in \mathcal{F}_{\text{NN}}(D, d+1, L_{\mathcal{E}^i}, W_{\mathcal{E}^i})$ , a decoder  $\mathcal{D}_{\text{NN}}^i \in \mathcal{F}_{\text{NN}}(d+1, D, L_{\mathcal{D}^i}, W_{\mathcal{D}^i})$ ,  
 339 a propagator  $\mathcal{P}_{\text{NN}}^i \in \mathcal{F}_{\text{NN}}(d+1, d+1, L_{\mathcal{P}^i}, W_{\mathcal{P}^i})$ , and a transcoder (for  $i < W$ )  
 340  $\mathcal{T}_{\text{NN}}^i \in \mathcal{F}_{\text{NN}}(d+1, d+1, L_{\mathcal{T}^i}, W_{\mathcal{T}^i})$  such that for any  $k \in [\mathsf{T}]$ , the WeldNet  $\mathcal{W}_{\text{NN}}$   
 341 given in (2.6) guarantees

342 
$$\sup_{\mathbf{x}(0) \in \mathcal{M}(0)} \|\mathcal{W}_{\text{NN}}(\mathbf{x}(0), t_k) - \mathcal{F}(\mathbf{x}(0), t_k)\|_{\mathbb{R}^D} < \epsilon.$$

343 The network parameters are

344 
$$L_{\mathcal{E}^i} = O(\log^2(S/\epsilon)), \quad W_{\mathcal{E}^i} = O(D(S/\epsilon)^{d+1})$$
  
 345 
$$L_{\mathcal{D}^i} = O(\log(S/\epsilon)), \quad W_{\mathcal{D}^i} = O(D(S/\epsilon)^{d+1})$$
  
 346 
$$L_{\mathcal{P}^i} = O(\log(S/\epsilon)), \quad W_{\mathcal{P}^i} = O((S/\epsilon)^{d+1})$$
  
 347 
$$L_{\mathcal{T}^i} = O(\log(S/\epsilon)), \quad W_{\mathcal{T}^i} = O((S/\epsilon)^d)$$

348 where we hide constants depending on (for all  $i \in [W]$ )  $\|g_{\pi(i)}\|_{L^\infty}$ ,  $\text{Lip}(g_{\pi(i)})$ ,  $\log(D)$ ,  
 349  $d$ ,  $\tau(\mathcal{M}([s_{\pi(i)}, s_{\pi(i)+1}]))$ ,  $\max_i \text{Lip}(\mathcal{E}_*^{\pi(i)})$ ,  $\max_i \text{Lip}(\mathcal{D}_*^{\pi(i)})$ ,  $\text{Lip}(\mathcal{F})$ , the volume of  
 350  $\mathcal{M}([s_{\pi(i)}, s_{i+1}])$ , and  $\sup_{\mathbf{x} \in \mathcal{M}([w_i, w_{i+1}])} \|\mathbf{x}\|_{\mathbb{R}^D}$ .

351 Theorem 3.9 is a generalization of Lemma 3.7 (they both rely on Assumptions  
 352 1, 2, and 3), and its proof is deferred to Appendix B.2. Without Assumption 3,  
 353 we can still establish an approximation error guarantee for WeldNet, but with a  
 354 different construction for the propagator network. The proof of the following theorem  
 355 is deferred to Appendix B.3.

356 THEOREM 3.10. Suppose Assumption 1 and 2 hold. Let  $\epsilon > 0$ , and suppose there  
 357 are  $W \geq S$  windows so that  $0 = w_1 < w_2 < \dots < w_{W+1}$  subdivides the segments.  
 358 Suppose the time grid  $\mathsf{T}$  is uniform such that  $t_{k+1} - t_k = \Delta t$  for all  $k \in [\mathsf{T}-1]$ . For all  
 359  $i \in [S]$ , let  $\overline{T}_s = |\mathsf{T} \cap [s_i, s_{i+1}]|$  denote the number of time steps in segment  $i$ . Then for  
 360 each  $i \in [W]$ , there exist an encoder network  $\mathcal{E}_{\text{NN}}^i \in \mathcal{F}_{\text{NN}}(D, d+1, L_{\mathcal{E}^i}, W_{\mathcal{E}^i})$ , a decoder  
 361 network  $\mathcal{D}_{\text{NN}}^i \in \mathcal{F}_{\text{NN}}(d+1, D, L_{\mathcal{D}^i}, W_{\mathcal{D}^i})$ , a propagator network  $\mathcal{P}_{\text{NN}}^i \in \mathcal{F}_{\text{NN}}(d +$   
 362  $1, d+1, L_{\mathcal{P}^i}, W_{\mathcal{P}^i})$ , and a transcoder network (for  $i < W$ )  $\mathcal{T}_{\text{NN}}^i \in \mathcal{F}_{\text{NN}}(d+1, d +$   
 363  $1, L_{\mathcal{T}^i}, W_{\mathcal{T}^i})$  such that for any  $k \in [\mathsf{T}]$ , the WeldNet  $\mathcal{W}_{\text{NN}}$  given in (2.6) guarantees

364 
$$\sup_{\mathbf{x}(0) \in \mathcal{M}(0)} \|\mathcal{W}_{\text{NN}}(\mathbf{x}(0), t_k) - \mathcal{F}(\mathbf{x}(0), t_k)\|_{\mathbb{R}^D} < \epsilon.$$

365 The network parameters are

366 
$$L_{\mathcal{E}^i} = O(\log^2(S/\epsilon)), \quad W_{\mathcal{E}^i} = O(D(S/\epsilon)^{d+1})$$
  
 367 
$$L_{\mathcal{D}^i} = O(\log(S/\epsilon)), \quad W_{\mathcal{D}^i} = O(D(S/\epsilon)^{d+1})$$
  
 368 
$$L_{\mathcal{P}^i} = O(\overline{T}_{\pi(i)} \log(S/\epsilon)), \quad W_{\mathcal{P}^i} = O((\overline{T}_{\pi(i)}/\epsilon)^d)$$
  
 369 
$$L_{\mathcal{T}^i} = O(\log(S/\epsilon)), \quad W_{\mathcal{T}^i} = O((S/\epsilon)^{d+1})$$

370 where we hide constants in  $O$  depending on (for all  $i \in [W]$ )  $d$ ,  $\log(D)$ ,  $\tau(\mathcal{M}([0, T]))$ ,  
 371  $\max_i \text{Lip}(\mathcal{E}_*^i)$ ,  $\max_i \text{Lip}(\mathcal{D}_*^i)$ ,  $\text{Lip}(\mathcal{F})$ , volume of  $\mathcal{M}([0, T])$ , and  $\sup_{\mathbf{x} \in \mathcal{M}([0, T])} \|\mathbf{x}\|_{\mathbb{R}^D}$ .

372     **Remark.** The main difference between Theorem 3.10 and Theorem 3.9 is that  
 373 the former does not assume Assumption 3, which makes the approximation of the  
 374 propagator network more challenging, resulting in the large propagator size. We  
 375 briefly remark that the idea of the proof of Theorem 3.10 does not require a  
 376 uniform time grid size, so it can be extended to the case of non-uniform time grids  
 377 (in which case the size of the propagator will depend on the minimum and maximum  
 378 time grid spacing).

379     **4. Experiments.**

380     **4.1. Data Generation.** For simplicity, we start from an evolutionary ordinary  
 381 differential equation (ODE)

$$382 \quad (4.1) \quad \begin{aligned} \partial_t \mathbf{x} &= \mathfrak{F}(\mathbf{x}(t), t), \\ \mathbf{x}(0) &= \mathbf{x}_0 \end{aligned}$$

383 with solution  $\mathbf{x}(t) : [0, T] \rightarrow \mathbb{R}^D$ , where  $T > 0$  is the end time. We will measure the  
 384 solution of the ODE in (4.1) at discrete time points; let  $0 = t_1 < t_2 < \dots < t_T = T$  be  
 385 the discretized time locations for some fixed  $T \in \mathbb{N}$ . The data for the ODE represents  
 386 a trajectory of (4.1) with initial condition  $\mathbf{x}(t_1) = \mathbf{x}_0$ :  $\{\mathbf{x}(t_1), \mathbf{x}(t_2), \dots, \mathbf{x}(t_T)\} \subseteq \mathbb{R}^D$ .

387 In the data-driven framework, we collect trajectory data from multiple initial con-  
 388 ditions. Let  $\{\mathbf{x}_n(t_1)\}_{n=1}^N$  be  $N$  sets of initial conditions. The collection of discretized  
 389 trajectory data with these initial conditions is denoted by

$$390 \quad (4.2) \quad \{\mathbf{x}_n(t_1), \mathbf{x}_n(t_2), \dots, \mathbf{x}_n(t_T)\}_{n=1}^N,$$

391 which serves as the training data (and can be thought of as  $N$  elements of  $\mathbb{R}^{T \times D}$ ).

392 We next consider a partial differential equation (PDE)

$$393 \quad (4.3) \quad \begin{aligned} \partial_t \mathbf{u} &= \mathfrak{F}(\mathbf{u}(\mathbf{x}, t), t), \\ \mathbf{u}(\mathbf{x}, 0) &= f(\mathbf{x}), \quad \forall \mathbf{x} \in \Omega, \end{aligned}$$

394 with solution  $\mathbf{u} : \Omega \times [0, T] \rightarrow \mathbb{R}$  where  $\Omega \subset \mathbb{R}^{d_\Omega}$  is the spatial domain. We will  
 395 discretize the solution of the PDE in the spatial and temporal domain. Let  $\mathbb{X} =$   
 396  $(X_1, \dots, X_D) \subseteq \Omega$  be a discretization set in the spatial domain and consider the  
 397 same temporal discretization as before. The discretized data for the PDE solution  
 398  $\mathbf{u}$  is denoted  $\{\mathbf{u}(\mathbb{X}, t_1), \mathbf{u}(\mathbb{X}, t_2), \dots, \mathbf{u}(\mathbb{X}, t_T)\} \subseteq \mathbb{R}^D$ , which represents a discretized  
 399 trajectory of the PDE in (4.3) with the initial condition  $f$ .

400 Given initial conditions  $\{f_n\}_{n=1}^N$  and solution trajectories  $\{\mathbf{u}_n\}_{n=1}^N$  (for all  $n \in [N]$ ),  
 401  $\mathbf{u}_n$  solves (4.3) with initial condition  $f = f_n$ , the trajectory data in this case has the  
 402 same form as (4.2) if we use the notation  $\mathbf{x}_n(t_k) = \mathbf{u}_n(\mathbb{X}, t_k)$ , that is,  $\mathbf{x}_n(t_k) \in \mathbb{R}^D$   
 403 with  $i$ th component given by  $\mathbf{u}_n(X_i, t_k)$ .

404 Many differential equations in applications can be complex to simulate - the dis-  
 405cretization set size grows exponentially in the spatial dimension  $d_\Omega$ . When the spatial  
 406 dimension  $d_\Omega$  is higher than one, a large number of discretization points are needed.  
 407 Our goal is to build a low-dimensional surrogate model based on training data of the  
 408 form (4.2) for ODEs, PDEs, and other evolutionary processes. After training, one can  
 409 predict the solution trajectory for a new initial condition by evolving the surrogate  
 410 model in a low-dimensional space.

411     **4.2. Benefit of Windowed Approach.** In this section, we illustrate the ben-  
 412 efit of using a windowed autoencoder approach on a simple but instructive example.

413 Consider the one-dimensional transport equation for  $T = 0.3$ , given by

414 (4.4) 
$$u_t = -u_x; \quad u(x, 0) = g(x), \quad x \in (0, 1),$$

415 with zero Dirichlet boundary conditions. We consider the weak version of this PDE,  
416 so the initial condition  $g$  does not have to be differentiable everywhere.

We consider initial conditions containing two hats. Let  $\sigma(x) = \max(x, 0)$ , and fix  $\epsilon = 0.05$ . For any  $\alpha$  and  $t$ , we define the “hat” function with height  $\alpha$ , location  $s$ , and width  $\epsilon$  as

$$H_{\alpha,s,\epsilon}(x) = 2\alpha/\epsilon \left( \sigma(x-s) - 2\sigma\left(x-s-\frac{\epsilon}{2}\right) + \sigma(x-s-\epsilon) \right).$$

417 We consider the following one-parameter set of initial conditions:

418 (4.5) 
$$\hat{g}_{\text{tscale}} = \{H_{a,0.1,0.05}(x) + H_{2.5,0.2,0.05}(x) : a \in [1, 4]\}$$

419 We discretize the spatial domain  $[0, 1]$  with 512 equally spaced points. We con-  
420 sider a probability measure on the set  $\hat{g}_{\text{tscale}}$  obtained by sampling the parameter  $a$   
421 uniformly from  $[1, 4]$  and then outputting the discretized initial condition correspond-  
422 ing to that value of  $a$ . Formally, we define a function  $G_{\text{tscale}} : [1, 4] \rightarrow \mathbb{R}^{512}$  that  
423 maps from parameter ( $a$ ) to initial condition. This allows us to construct the initial  
424 measure  $\rho_{\text{tscale}}(0) = (G_{\text{tscale}})_\sharp \text{Unif}([1, 4])$ . We sample by  $\rho_{\text{tscale}}(0)$  by first choosing  
425  $1 \leq a \leq 4$  uniformly and outputting the corresponding (discretized) initial condition  
426 for  $a$  as defined in (4.5).

427 The trajectory manifold  $\mathcal{M}_{\text{tscale}}$  is obtained by collecting solutions from the PDE  
428 in (4.4) with initial conditions from  $\rho_{\text{tscale}}(0)$ . For simplicity, we omit the time interval  
429 notation for the manifolds in this section. We collect 500 trajectories with 301 time  
430 steps until  $T = 0.3$ , so our data is of the form  $((\mathbf{x}_n(t_k)))_{k=1}^{51}{}_{n=1}^{500}$ , as detailed in Section  
431 2.1.

432 We will train WeldNet models on this dataset with one window, two windows,  
433 and four windows, and we will attempt to roughly equalize the number of trainable  
434 parameters and number of training epochs. Specifically, we will implement each  
435 component of WeldNet using 3 layer neural networks so that we implement the one-  
436 window, two-window, and four-window models with networks of width 1000, 500, and  
437 250 (respectively) and train each component for 1200, 600, and 300 epochs (respec-  
438 tively). The total training time for each model in minutes is 42.28, 15.73, and 8.31  
439 (respectively).

440 In this section, we will use a latent space dimension of 2 for the autoencoders -  
441 equal to the intrinsic dimension of the data (see Table 2). This results are in a lower  
442 performance than a latent space dimension of 4 that we use later in the paper. We use  
443 this reduced latent space dimension to show that using windows with the estimated  
444 intrinsic dimension (instead of using a higher latent space dimension) can overcome  
445 the error accumulation issues of a single-window autoencoder model, and due to the  
446 ease of visualizing two dimensional latent spaces.

447 We train autoencoders and propagators together for all models. We show the  
448 latent space of each model in Figures 6, 7, and 8 for the one-window, two-window,  
449 and four-window models respectively. We show each window individually, and we  
450 show the points clouds colored by their parameter in the left plot and time in the  
451 right plot. The latent space for the one window model in Figure 6 involves pinches  
452 and twists, while the latent spaces (in each window) for the two-window and four-  
453 window are much more uniform. We note that even though the latent space in the

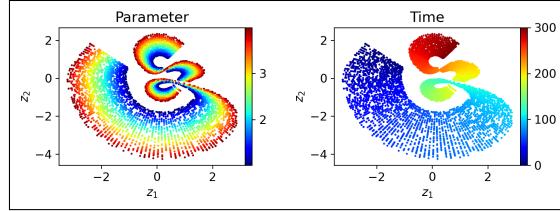


FIG. 6. Latent space of illustrative example with one-window model.

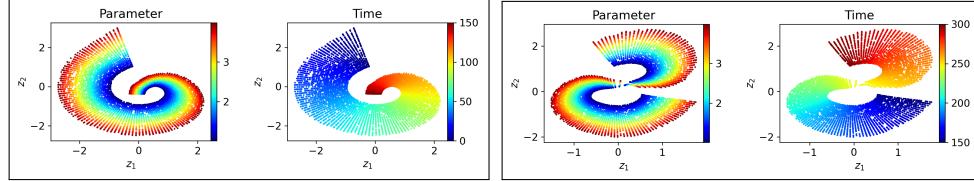


FIG. 7. Latent space of illustrative example with two-window model.

454 second window of the two-window model has a compressed area in Figure 7, the  
 455 propagator is still able to learn these latent dynamics.

456 Next, we train propagators for the one-window model and the propagators and  
 457 transcoders for the two-window and four-window models. We compare the projection  
 458 and operator error of the models over all times in Figure 9. The projection error is  
 459 the error of the autoencoder reconstruction on the test data (i.e. the representation  
 460 or approximation error of the autoencoder in representing the solution manifold).  
 461 The operator error is the prediction error of a solution given an initial condition,  
 462 which depends on both the projection error and the error of the propagator network  
 463 in learning the dynamics. We observe severe error accumulation in the one-window  
 464 model. This problem is resolved by using more windows.

465 **4.3. Implementation Details.** The datasets for the Burgers' equation and  
 466 KdV equation trajectory manifolds were generated using the chebfun package [13]  
 467 in Matlab. The datasets for the transport equation trajectory manifolds were gener-  
 468 ated using the analytic solution of the equation.

469 **4.3.1. Model Descriptions.** In this section, we detail the architectures of the  
 470 components of WeldNet and comparison models. Suppose the data is  $D$ -dimensional.

471 FF-WeldNet and Conv-WeldNet uses ReLU networks with 3-hidden layers and  
 472 width 200 for the propagators and transcoders. FF-WeldNet uses feedforward ReLU  
 473 networks with 3-hidden layers and width 500 for the encoder and decoder. Conv-  
 474 WeldNet uses a 4 layer convolutional encoder with the convolution layers having  
 475 channel size 8, 16, 32, 32, and kernel size 8, 8, 8, 4, respectively, with stride 2, and  
 476 symmetric zero padding of 1, and it uses a symmetric architecture of convolution  
 477 transpose layers for the decoder.

478 Note that a one-window WeldNet model is a normal latent dynamics model which  
 479 has already been considered in the literature (e.g. see [26]), so we will use the label  
 480 "AENet" to refer to that. The goal of our work is to show the benefit of multiple  
 481 windows, so the performance of a one-window WeldNet model or an AENet should  
 482 only be considered as a baseline.

483 To illustrate the advantage of dimension reduction, we train a model consisting  
 484 of a propagator operating directly on the  $D$ -dimensional data. We call this approach

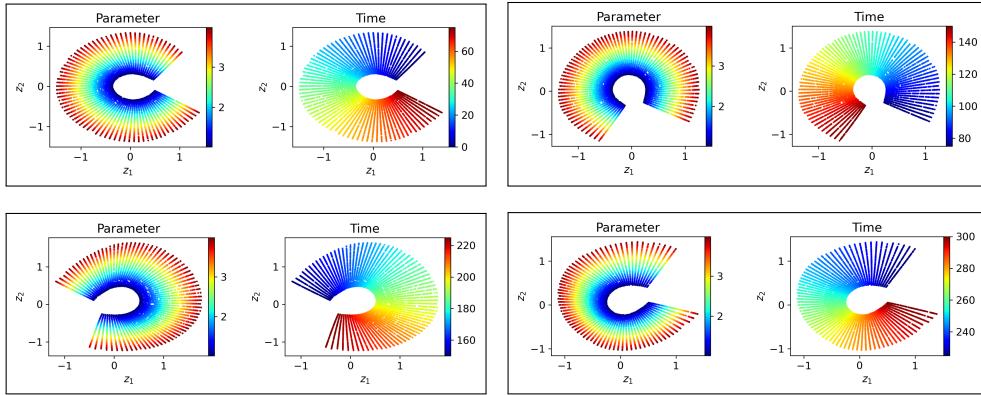


FIG. 8. Latent space of illustrative example with four-window model.

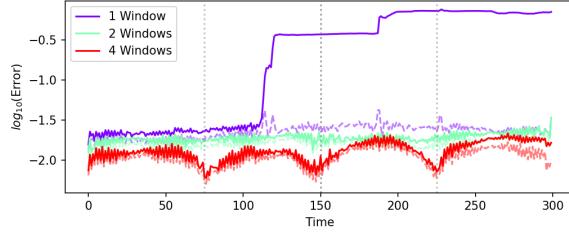


FIG. 9. Projection (dashed) and Operator (solid) error for the example in Section 4.2.

485 High Dimensional Propagator (HDP). Specifically, we train (depth 3) ReLU networks  
 486 with input and output dimension equal to  $D$  ( $D = 512$  for most of our examples).  
 487 We use a width of 1000 since HDP only consists of a single network (while WeldNet  
 488 and other models have several networks). We found that using residual/displacement  
 489 based training as done for the the WeldNet propagator (as described in Section 4.3.2)  
 490 leads to training instability. Thus, we will train the HDP to predict the high dimen-  
 491 sional state at the next time directly, i.e. given an input of  $x(t) \in \mathcal{M}(t) \subseteq \mathbb{R}^D$ , the  
 492 HDP should predict  $x(t + \Delta t) \in \mathbb{R}^D$ .

493 The Time-Input feedforward networks are three-hidden layer ReLU networks with  
 494 width 1000. If the data is  $D$ -dimensional, the Time-Input feedforward network has  
 495  $D+1$  input neurons and  $D$  output neurons. Given an input  $u(0)$  (representing a data  
 496 point initial condition) and a time  $t$ , the Time-Input network outputs the evolution  
 497 of that input by  $t$  time units, i.e. it should approximately be  $u(t)$ .

498 We also compare to a model called Latent Deep Operator Network (Latent-DON)  
 499 [24]. Similar to a one-window WeldNet, Latent-DON trains a single autoencoder over  
 500 all times and then trains a latent prediction model that uses a Deep Operator Network  
 501 (DON) architecture. DON is an operator learning method that takes time as an input  
 502 and tries to predict the evolution of an initial condition by that given time, similar  
 503 to time-input models, but in latent space. Given an input  $u(0) \in \mathbb{R}^D$  and a time  $t$ ,  
 504 Latent-DON first projects to a  $k$ -dimensional latent code  $z(0) \in \mathbb{R}^k$  using the encoder,  
 505 and then uses the latent code and time as inputs to the DON to try to predict the  
 506 evolution of that latent code by  $t$  time units. We then decode the predicted evolved

507 latent code, and this value should approximately be  $u(t)$ . See Appendix C.2 for  
 508 specific details on the DON architecture used in Latent-DON.

509 We remark that the Time-Input and Latent-DON models have an easier learning  
 510 task than WeldNet, because they are only tested to predict the evolution of the initial  
 511 conditions by a given time, while WeldNet is capable of evolving data from any start  
 512 time to any end time on the time grid.

513 LDNet is a model reduction based method for learning dynamical systems. The  
 514 method trains a dynamics network that evolves latent codes and a reconstruction  
 515 network that maps from latent code to function predictions on a grid. Specifically,  
 516 each initial condition is assigned a latent code consisting of  $f$  fixed dimensions and  
 517  $k$  dynamic dimensions that start at 0. The fixed dimensions are given by the codes  
 518 assigned to the initial conditions (which need to be known ahead of time). The  
 519 dynamics network is trained to evolve the dynamic portion of each latent code, with  
 520  $f + k$  input neurons and  $k$  output neurons. The latent code at a given time can be  
 521 evolved by the dynamics network iteratively. The output prediction is then obtained  
 522 by passing the evolved latent code through the reconstruction network. This means  
 523 that the reconstruction network has  $f + k$  input neurons and  $D$  output neurons. We  
 524 use a width of 500 for the dynamics and reconstruction networks, and we use the  
 525 ReLU activation function for both networks.

526 **Remark.** This description of LDNet is reliant on a fixed size  $D$  for the data and  
 527 reconstruction network, while the original implementation of LDNet in [43] uses a  
 528 grid-independent reconstruction network. We use a modified grid based version for a  
 529 more direct comparison to our grid based WeldNet models, and due to training issues  
 530 with the grid independent version. See Appendix C.1 for a comparison between the  
 531 original LDNet and the grid based LDNet that we will subsequently use in this paper.

532 Weak-form Latent Space Dynamics Identification, WLaSDI, is a reduced-order  
 533 modeling technique recently proposed in [51]. WLaSDI uses a one-hidden layer au-  
 534 toencoder with 1024 hidden neurons in the encoder and 6168 hidden neurons in the  
 535 decoder with the *silu* activation function,  $\text{silu}(x) = \frac{x}{1+e^{-x}}$ . The latent dynamics are  
 536 then learned by solving a least-squares problem involving a dictionary of trial func-  
 537 tions; specifically we apply the region-based local dynamics identification algorithm  
 538 discussed in Section 3.3.2 of [51].

539 All WeldNet, Time-Input, HDP, Latent-DON, and LDNet models are trained  
 540 with a batch size of 32, learning rate of 1e-4, with an adaptive learning schedule  
 541 that decay by a factor of 0.3 with a patience of 15 epochs, to a minimum learning  
 542 rate of 1e-6. We use PyTorch [38] to implement the neural networks, and we use  
 543 PyTorch’s imlementation of the AdamW optimization algorithm with default settings  
 544 for training the models.

545 For WLaSDI, the autoencoder is trained with a batch size of 20, learning rate of  
 546 1e-4, with an adaptive learning rate schedule that decays by a factor of 0.1 with a  
 547 patience of 10 epochs, using the implementation in [51].

548 **4.3.2. Displacement Networks.** When training the propagators, we noticed  
 549 that using a “displacement” or residual architecture for the propagator and transcoder  
 550 significantly aids in training. Specifically, we implement the propagator network  $\mathcal{P}_{\text{NN}}$   
 551 as a function of the form:

$$552 \quad \mathcal{P}_{\text{NN}}(\mathbf{z}) = \mathbf{z} + f_{\text{NN}}(\mathbf{z}),$$

553 where  $f_{\text{NN}}$  is a trainable neural network. We implement the transcoder in the same  
 554 way. While changing the propagator and transcoder to this form does not affect

555 the approximation theory, we found experimentally that this change resulted in a  
 556 much lower loss in the training. On the other hand, the loss explodes even for low  
 557 learning rate when this method is used for the HDP model, so we only employ this  
 558 displacement network strategy for WeldNet. Note that LDNet uses a displacement  
 559 form for its dynamics network as part of its original specification, so we also use this  
 560 strategy there.

561 **4.4. Numerical Results.** We test WeldNet on several numerical examples, sim-  
 562 ilar to the illustrative example in Section 4.2, demonstrating the efficacy of our non-  
 563 linear dimension reduction and trajectory learning model. For all examples, we collect  
 564 500 trajectories with 301 time steps, so our data is of the form  $((\mathbf{x}_n(t_k)))_{k=1}^{301}_{n=1}$ , as  
 565 detailed in Section 2.1.

566 We consider three versions of WeldNet based on the architecture for the autoen-  
 567 coder: FF-Weld (feedforward networks), Conv-Weld (convolutional networks), and  
 568 PCA-WeldNet (Principal Component Analysis for linear dimension reduction). Note  
 569 that PCA-WeldNet is analogous to the classic POD or linear model reduction meth-  
 570 ods.

571 We compare with Feedforward Neural Networks, LDNets in [31] and Weak-form  
 572 Latent Space Dynamics Identification (WLaSDI) as described in [51]. We use a latent  
 573 dimension of 4 for WeldNet and WLaSDI, and we use 1 fixed dimension and 3 dynamic  
 574 dimensions for LDNet.

575 The numerical test error of all the trained models at various times can be found  
 576 in Appendix D.

577 **4.4.1. Burgers' Equation.** We consider the viscous Burgers' equation with  
 578  $\nu = 1/1000$  and with periodic boundary conditions for  $t \in [0, 1]$ .

579 (4.6) 
$$u_t = \nu u_{xx} - uu_x; \quad u(x, 0) = g(x), \quad x \in (0, 1).$$

580 We consider initial conditions generated by combining two complex base waves.  
 581 Let  $w_0, w_1$  be functions sampled from the Gaussian Random Field  $N(0, 7^4(-\frac{d^2}{dx^2} +$   
 582  $7^2 I)^{-2.5})$  on  $[0, 1]$ . For any  $a \in [-0.9, 0.9]$ , consider the two sets of initial conditions:

583 (4.7) 
$$\hat{g}_{\text{bscale}} = \{aw_0(x) + \sqrt{1-a^2}w_1(x) : a \in [-0.9, 0.9]\},$$

584 (4.8) 
$$\hat{g}_{\text{bshift}} = \{0.5w_0(x-h) + \sqrt{0.75}w_1(x-h) : h \in [0, 1]\}.$$

585 The trajectory manifolds  $\mathcal{M}_{\text{bscale}}$  and  $\mathcal{M}_{\text{bshift}}$  are obtained by collecting solutions  
 586 from the PDE in (4.6) with initial conditions from  $\hat{g}_{\text{bscale}}$  and  $\hat{g}_{\text{bshift}}$  respectively. We  
 587 collect 500 data trajectories with 301 time steps. We display the estimated intrinsic  
 588 dimensionality of these datasets by Maximum Likelihood Estimation (MLE) [28] at  
 589 various times and as a whole in Table 2 in Section 4.5 for details on this computation).  
 590 The intrinsic dimension for each time is approximately 1, and the dimension of the  
 591 datasets for all times is approximately 2, matching our expectation of 1 parameter  
 592 dimension plus 1 time dimension.

593 We compare the performance of our WeldNet models (for  $W = 4$ ) with other  
 594 models in Figures ?? and 11 in terms of final time Prediction error. We use a latent  
 595 space dimension of **four** for this example. The  $x$ -axis is the parameter value for  $a$  or  $h$   
 596 (respectively), and the  $y$ -axis shows the prediction error at the final time for the initial  
 597 condition corresponding to that parameter value. In Figures ?? and 11, FF-WeldNet  
 598 significantly outperforms HDP, highlighting the advantages of model reduction: low-  
 599 dimensional latent propagators are substantially easier to learn than high-dimensional

600 ones. FF-WeldNet also significantly outperforms PCA-WeldNet, which demonstrates  
 601 the advantage of nonlinear dimension reduction. On this problem, FF-WeldNet and  
 602 LDNet achieve the best result.

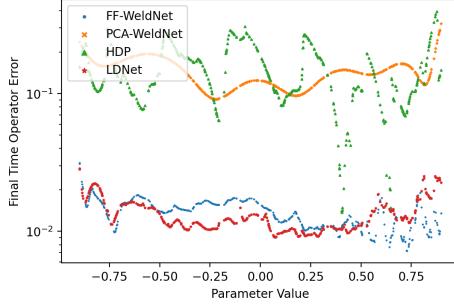


FIG. 10. Error vs Parameter value for  $\mathcal{M}_{bscale}$  about the trajectory manifold of the Burgers' equation (4.6) with initial conditions in (4.7).

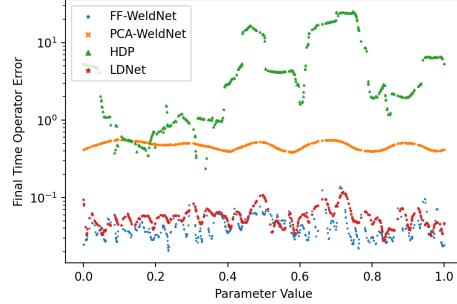


FIG. 11. Error vs Parameter value for  $\mathcal{M}_{bshift}$  about the trajectory manifold of the Burgers' equation (4.6) with initial conditions in (4.8).

603 For both of the trajectory manifolds  $\mathcal{M}_{bscale}$  and  $\mathcal{M}_{bshift}$ , we can compare the  
 604 reconstruction error and the operator/prediction error at each time. Figures 12 and 13  
 605 are line plots of error versus time for FF-WeldNet, Conv-WeldNet, and PCA-WeldNet  
 606 for  $\mathcal{M}_{bscale}$  and  $\mathcal{M}_{bshift}$ , respectively. For this example, FF-WeldNet and Conv-  
 607 WeldNet with latent dimension 4 outperforms PCA-WeldNet with latent dimension  
 608 4, since the solution manifold are nonlinear in these examples.

609 **4.4.2. Transport Equation.** We next consider the transport equation given by  
 610 (4.4). We consider initial conditions given by  $\hat{g}_{tscale}$  in (4.5) and

$$611 \quad (4.9) \quad \hat{g}_{tshift}(x) = \{H_{1,0.1}(x) + H_{2.5,0.2+0.1h}(x) : h \in [0, 3]\}$$

612 and collect the trajectory manifold  $\mathcal{M}_{tscale}$  and  $\mathcal{M}_{tshift}$  from the corresponding condition sets. We use a latent space dimension of **four** for this example. We compare the  
 613 performance of four window WeldNet models with other models in terms of final time  
 614 prediction error in Figures 14 and 15 respectively. WLaSDI has a very high relative  
 615 operator error, so it does not appear in the figure to allow for clarity of scaling. The  
 616 transport equation is a representative example of advection-dominated PDEs, which  
 617 generates a highly nonlinear solution manifold. In Figures 14 and 15, FF-WeldNet  
 618 with latent dimension 4 outperforms PCA-WeldNet with latent dimension 4, HDP  
 619 and LDNet.

620 For both of the trajectory manifolds  $\mathcal{M}_{tscale}$  and  $\mathcal{M}_{tshift}$ , we can compare the  
 621 reconstruction error and the operator error at each time. Figures 16 and 17, are  
 622 line plots of error versus time for FF-WeldNet, Conv-WeldNet, and PCA-WeldNet  
 623 for  $\mathcal{M}_{tscale}$  and  $\mathcal{M}_{tshift}$ , respectively. In Figures 16 and 17, Conv-WeldNet achieves  
 624 the lowest error, followed by FF-WeldNet. Both models substantially outperform  
 625 PCA-WeldNet for the same latent space dimension.

626 **4.4.3. KdV Equation.** We also consider the Korteweg–De Vries (KdV) equa-  
 627 tion for  $T = 0.01$ , given by

$$629 \quad (4.10) \quad u_t = -u_{xxx} - uu_x, \quad u(x, 0) = g(x), \quad x \in (0, 6).$$

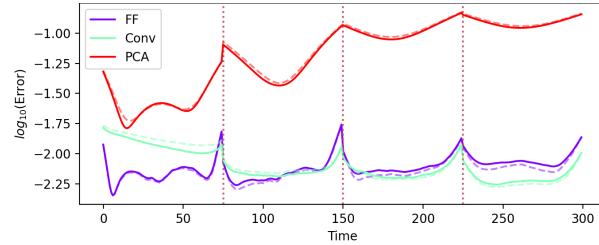


FIG. 12. Autoencoder error (dashed) and Test error for 4-window WeldNet models for  $\mathcal{M}_{\text{kscale}}$ . FF refers to FF-WeldNet, Conv refers to Conv-WeldNet, and PCA refers to PCA-WeldNet. The latent space dimension is 4.

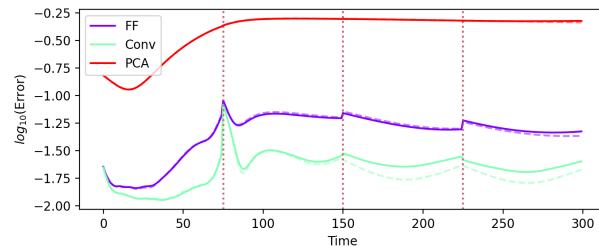


FIG. 13. Autoencoder error (dashed) and Test error for 4-window WeldNet models for  $\mathcal{M}_{\text{bshift}}$ .

630 We describe our sets of initial conditions. Given a constant  $c$ , we define the soliton of  
 631 size  $c$  as  $\phi_c(x) = c/2 \operatorname{sech}(\sqrt{c}x/2)^2$ . We consider initial conditions:

$$(4.11) \quad \hat{g}_{\text{kscale}} = \{\phi_{a^2}(x-1) + \phi_{36}(x-2) : a \in [6, 18]\},$$

$$(4.12) \quad \hat{g}_{\text{kshift}} = \{\phi_{36}(x-1) + \phi_{36}(x-2-h) : h \in [0, 0.4]\},$$

634 and collect the trajectory manifold  $\mathcal{M}_{\text{kscale}}$  and  $\mathcal{M}_{\text{kshift}}$  from the corresponding condition sets. We compare the performance of four-window WeldNet models (with latent 635 space dimension four) with other models in terms of final time operator error in Figure 18 and 19 respectively. WLSDI has a very high relative operator error, so most 636 of it is cut off from the figure to allow for clarity of scaling. For both examples, 637 FF-WeldNet is the best model, outperforming PCA-WeldNet, HDP, and LDNet for 638 almost all parameter values.  
 639

640 For both of the trajectory manifolds  $\mathcal{M}_{\text{kscale}}$  and  $\mathcal{M}_{\text{kshift}}$ , we can compare the 641 reconstruction error and the operator error at each time. Figures 20 and 21, are 642 line plots of error versus time for FF-WeldNet, Conv-WeldNet, and PCA-WeldNet 643 for  $\mathcal{M}_{\text{kscale}}$  and  $\mathcal{M}_{\text{kshift}}$ , respectively. FF-WeldNet and Conv-WeldNet perform 644 similarly on this problem, and they both greatly outperform PCA-WeldNet, showing the 645 advantage of nonlinear dimension reduction.  
 646

647 **4.4.4. Shallow-Water Equations.** We now consider an example with two  
 648 spatial dimensions. Consider the shallow-water equations over the spatial domain  
 649  $[-2.5, 2.5]^2$  and time domain  $[0, 1]$ :

$$(4.13) \quad \partial_t h + \partial_x(hu) + \partial_y(hv) = 0,$$

$$(4.14) \quad \partial_t(hu) + \partial_x(u^2h + \frac{1}{2}g_r h^2) + \partial_y(uvh) = -g_r h \partial_x b,$$

$$(4.15) \quad \partial_t(hv) + \partial_y(v^2h + \frac{1}{2}g_r h^2) + \partial_x(uvh) = -g_r h \partial_y b,$$

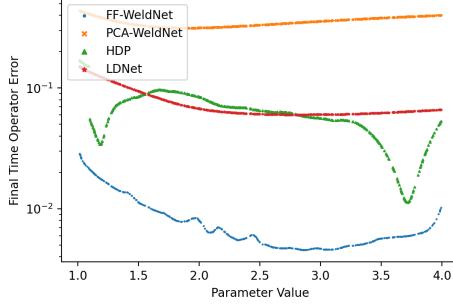


FIG. 14. Error vs Parameter value for  $\mathcal{M}_{tscale}$  about the trajectory manifold of the transport equation (4.4) with initial conditions in (4.7).

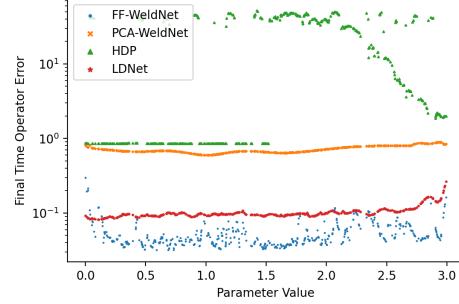


FIG. 15. Error vs Parameter value for  $\mathcal{M}_{tshift}$  about the trajectory manifold of the transport equation (4.4) with initial conditions in (4.8).

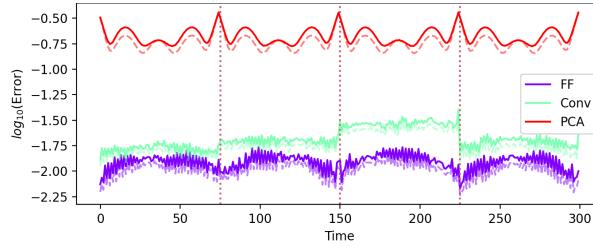


FIG. 16. Autoencoder error (dashed) and Test error for 4-window WeldNet models for  $\mathcal{M}_{tscale}$ .

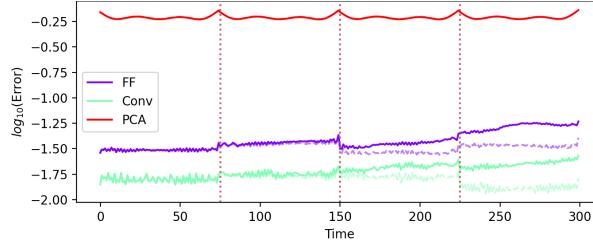


FIG. 17. Autoencoder error (dashed) and Test error for 4-window WeldNet models for  $\mathcal{M}_{tshift}$ .

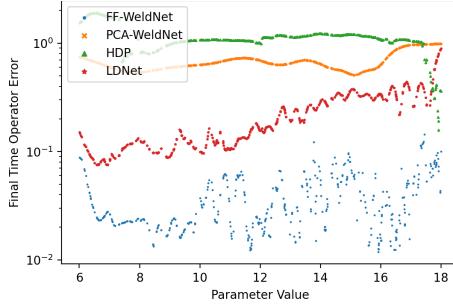


FIG. 18. Error vs Parameter value for  $\mathcal{M}_{kscale}$  about the trajectory manifold of the KdV equation (4.10) with initial conditions in (4.11).

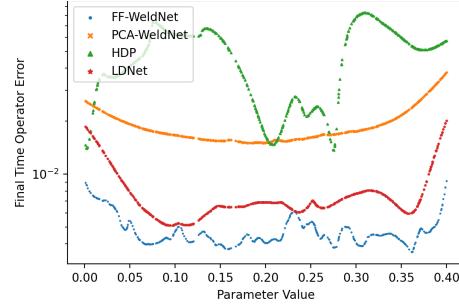
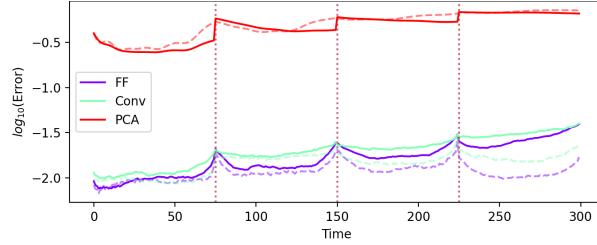
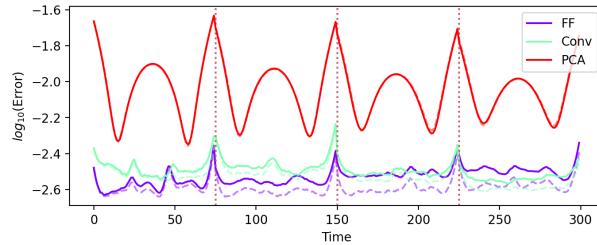


FIG. 19. Error vs Parameter value for  $\mathcal{M}_{kshift}$  about the trajectory manifold of the KdV equation (4.10) with initial conditions in (4.12).

FIG. 20. Autoencoder error (dashed) and Test error for WeldNet models for  $\mathcal{M}_{kscale}$ .FIG. 21. Autoencoder error (dashed) and Test error for WeldNet models for  $\mathcal{M}_{kshift}$ .

653 where  $u$  is horizontal velocity,  $v$  is vertical velocity,  $h$  is the water depth,  $g_r$  is the  
 654 gravitational acceleration, and  $b$  is a scalar field known as bathymetry. We impose  
 655 zero dirichlet boundary conditions for  $u$  and  $v$ .

656 We consider the following set of initial conditions representing centered bumps of  
 657 varying radii (provided by the PDEBench dataset [48, 47]):

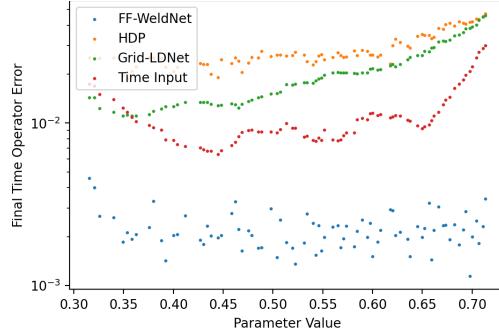
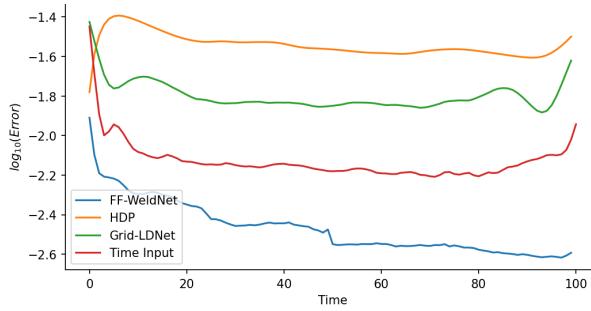
$$658 \quad (4.16) \quad \hat{g}_{\text{shallow}} = \left\{ h(x, y) = 1 + \mathbb{1}[\sqrt{x^2 + y^2} < r] : r \in [0.3, 0.7] \right\}$$

659 and collect the trajectory manifold  $\mathcal{M}_{\text{shallow}}$ . Note that the initial condition is dis-  
 660 continuous, and the location of the discontinuity depends on the sample. For this  
 661 example only, we use a  $128 \times 128$  grid for each sample, and we use 101 time steps.  
 662 We collect 81 examples.

663 Note that the inputs are sampled on a size 16384 grid. This makes the size  
 664 of feedforward networks prohibitively high for training. While more sophisticated  
 665 architectures such as two-dimensional convolutions exist, we will instead reduce the  
 666 dimension from 16384 to 128 using PCA, and then train our models to predict the  
 667 evolution of the PCA modes. We compute the top 128 principal components (over all  
 668 times) using the training data only, and this results in a relative projection error of  
 669 0.11%.

670 Now we have 81 samples with 101 time steps and 128 features each. We train each  
 671 of our models on this dataset. Even though the compression to 128 features is very  
 672 accurate, it can be challenging to learn the dynamics. We compare the performance of  
 673 the four-window WeldNet model (with latent space dimension of 4) with other models  
 674 in terms of final time operator error in Figure 22. WeldNet greatly outperforms all  
 675 other models by almost an entire order of magnitude.

676 The lower error for WeldNet stems from the fact that the trajectory manifold is  
 677 harder to represent at earlier times than later times, and WeldNet is able to separate  
 678 the latent dynamics learning in each of its segments. The transcoder is able to lower  
 679 the error as a code is evolved between windows. On the other hand, there is no

FIG. 22. *Error vs Parameter value for  $\mathcal{M}_{shallow}$ .*FIG. 23. *Error vs Parameter value for  $\mathcal{M}_{shallow}$ .*

680 error reduction mechanism in the other dynamics learning models, so the error only  
 681 increases with time. We can see this clearly in Figure 23, a plot of the operator  
 682 error versus time for various models. This plot shows that FF-WeldNet has improved  
 683 performance for later times, while the error stagnates or increases with time on the  
 684 other models. This shows the advantage of having multiple windows; a decoupled  
 685 representation of this trajectory manifold outperforms a fully coupled one.

686 **4.5. Intrinsic Dimension Estimation.** We estimate the intrinsic dimensionality  
 687 of our datasets using the Scikit-dimension [1] Python package and display it in  
 688 Table 2. Specifically, we use the MLE algorithm [28] with the “Haro” integral approxi-  
 689 mation [20] and the TwoNN algorithm [14] with the parameter `discard_fraction=0`.  
 690 We use a random size 50 000 subset for computing the intrinsic dimension of all times  
 691 together. By construction, the initial conditions for all considered datasets are gener-  
 692 ated by sampling one random scalar, so we expect the intrinsic dimensionality of each  
 693 time slice to be approximately 1 and the intrinsic dimensionality of all times together  
 694 to be approximately 2. This justifies the use of dimension reduction methods for these  
 695 datasets and indicates that a small latent space dimension can effectively capture the  
 696 data.

697 For the shallow water dataset, the relatively high intrinsic dimensionality for  
 698 the initial time and all times together likely comes from the discontinuity of the  
 699 initial condition in (4.16). The estimated intrinsic dimensionality for a continuous  
 700 approximation to the initial conditions can be computed, and we empirically found  
 701 that it was about 1 regardless of the specific construction used (as long as the functions  
 702 are continuous).

Dataset	t=0	t=60	t=120	t=180	t=240	t=300	All t
bscale	0.973	0.974	0.974	0.976	0.978	0.977	1.89
bshift	0.998	1.019	1.128	1.091	1.065	1.048	1.685
tscale	0.976	0.976	0.976	0.976	0.976	0.976	1.118
tshift	1.013	1.013	1.012	1.013	1.013	1.012	1.973
kscale	0.945	0.945	0.945	0.945	0.945	0.945	1.803
kshift	0.995	0.995	0.995	0.995	0.995	0.995	1.901
shallow	2.184	1.282	1.257	1.264	1.244	1.231	2.593

TABLE 2

*Estimated intrinsic dimensionality of all datasets for various times and as a whole. MLE is used for all datasets except for kscale, for which TwoNN was used (see Appendix 4.5 for details).*

703     **4.6. Propagator Ablation Study.** WeldNet first trains the autoencoder and  
 704 propagator together using the propagator displacement loss in (2.2) (where the name  
 705 displacement comes from being one step of applying propagator) and then finetunes  
 706 the propagator using the accumulation loss in (2.4). Alternatively, we could train the  
 707 autoencoder and propagator separately, or we could use only a single loss function  
 708 for the propagator. We compare our Weldnet model training algorithm against three  
 709 variants.

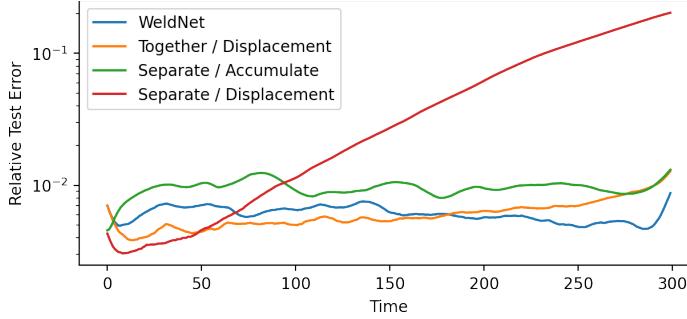
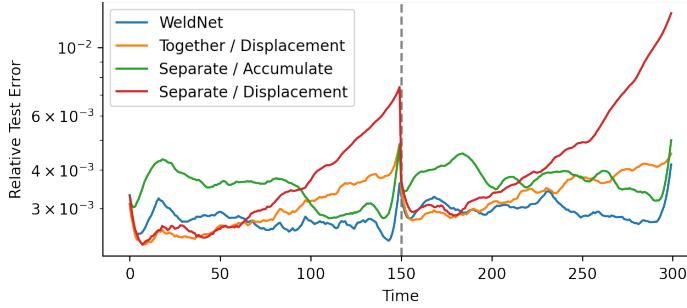
- 710       (i) “WeldNet” (original) - train autoencoder and propagator together with dis-  
 711       placement loss, then finetune the propagator with accumulation loss.
- 712       (ii) “Together / Displacement” - same as (i) but finetune with displacement loss.
- 713       (iii) “Separate / Accumulate” - train autoencoder first, then train propagator with  
 714       accumulation loss.
- 715       (iv) “Separate / Displacement” - same as (iii) but with the displacement loss.

716     We train autoencoder, propagator, and transcoder models using each of the above  
 717 four training algorithms. For algorithms (i) and (ii), we jointly train the autoencoder  
 718 and propagator for 300 epochs and then finetune the propagator for 150 epochs.  
 719 For algorithms (iii) and (iv), we train the autoencoder for 300 epochs and then the  
 720 propagator for 450 epochs. We train any transcoders for 300 epochs. All other settings  
 721 are identical to the main paper. We use the *kshift* dataset (for the KdV equation)  
 722 described in Section 4.4.3.

723     Figures 24 and 25 compare the relative test error of one-window and two-window  
 724 models (respectively) trained with each algorithm. The final time test error for each  
 725 one-window model is (in order): 0.87%, 1.28%, 1.31%, 20.4%, and the average test  
 726 error for each two-window model is (in order): 0.42%, 0.45%, 0.50%, 1.3%. While  
 727 increasing the number of windows reduces the gap between the models, the best model  
 728 is (i) which is the training algorithm used in the main paper.

729     **5. Conclusion.** This paper introduces WeldNet, a data-driven framework for  
 730 nonlinear model reduction that constructs low-dimensional surrogate models for com-  
 731 plex evolutionary systems. The architecture consists of three key components: autoen-  
 732 coders, propagators, and transcoders. These components operate together to model  
 733 dynamics in reduced latent spaces. By transferring the evolution of high-dimensional  
 734 systems into these latent spaces and decomposing the time domain into overlapping  
 735 windows, WeldNet effectively captures intricate nonlinear structures while simplifying  
 736 long-time propagation into a sequence of short, tractable segments.

737     In addition to the algorithmic design, a representation theory is developed to  
 738 establish the approximation capability of WeldNet under the manifold hypothesis,

FIG. 24. *Test error vs time for variant training algorithms on  $\mathcal{M}_{kshift}$  (1 window models).*FIG. 25. *Test error vs time for variant training algorithms on  $\mathcal{M}_{kshift}$  (2 window models).*

739 thereby providing a mathematical foundation for nonlinear model reduction via deep  
 740 learning. Extensive numerical experiments on a variety of differential equations  
 741 demonstrate the robustness and accuracy of the proposed approach. Across all tested  
 742 scenarios, WeldNet consistently achieves smaller prediction errors than both classical  
 743 projection-based techniques and recently developed nonlinear reduced-order models.  
 744 In addition, WeldNet is more computationally efficient to train than other methods.

745 Overall, this work shows that windowed latent-space learning offers a powerful and  
 746 principled strategy for modeling complex dynamical systems, and it opens promising  
 747 ways for future research in nonlinear model reduction in scientific machine learning.

748

## REFERENCES

- 749 [1] J. BAC, E. M. MIRKES, A. N. GORBAN, I. TYUKIN, AND A. ZINOVYEV, *Scikit-dimension: a*  
 750 *python package for intrinsic dimension estimation*, Entropy, 23 (2021), p. 1368.
- 751 [2] M. BACHMAYR AND A. COHEN, *Kolmogorov widths and low-rank approximations of parametric*  
 752 *elliptic pdes*, Mathematics of Computation, 86 (2017), pp. 701–724.
- 753 [3] U. BAUR, C. BEATTIE, P. BENNER, AND S. GUGERCIN, *Interpolatory projection methods for pa-*  
 754 *rameterized model reduction*, SIAM Journal on Scientific Computing, 33 (2011), pp. 2489–  
 755 2518.
- 756 [4] M. BELKIN AND P. NIYOGI, *Laplacian eigenmaps for dimensionality reduction and data repre-*  
 757 *sentation*, Neural computation, 15 (2003), pp. 1373–1396.
- 758 [5] P. BENNER, S. GUGERCIN, AND K. WILLCOX, *A survey of projection-based model reduction*  
 759 *methods for parametric dynamical systems*, SIAM review, 57 (2015), pp. 483–531.
- 760 [6] P. BENNER, M. OHLEBERGER, A. COHEN, AND K. WILLCOX, *Model reduction and approximation:*  
 761 *theory and algorithms*, SIAM, 2017.
- 762 [7] G. BERKOZ, P. HOLMES, AND J. L. LUMLEY, *The proper orthogonal decomposition in the*  
 763 *analysis of turbulent flows*, Annual review of fluid mechanics, 25 (1993), pp. 539–575.

- [8] P. BINEV, A. COHEN, W. DAHMEN, R. DEVORE, G. PETROVA, AND P. WOJTASZCZYK, *Convergence rates for greedy algorithms in reduced basis methods*, SIAM journal on mathematical analysis, 43 (2011), pp. 1457–1472.
- [9] S. L. BRUNTON AND J. N. KUTZ, *Data-driven science and engineering: Machine learning, dynamical systems, and control*, Cambridge University Press, 2022.
- [10] M. CHEN, H. JIANG, W. LIAO, AND T. ZHAO, *Efficient approximation of deep relu networks for functions on low dimensional manifolds*, Advances in neural information processing systems, 32 (2019), pp. 8174–8184.
- [11] R. R. COIFMAN, I. G. KEVREKIDIS, S. LAFON, M. MAGGIONI, AND B. NADLER, *Diffusion maps, reduction coordinates, and low dimensional representation of stochastic systems*, Multiscale Modeling & Simulation, 7 (2008), pp. 842–864.
- [12] M. CROSSKEY AND M. MAGGIONI, *Atlas: a geometric approach to learning high-dimensional stochastic systems near manifolds*, Multiscale Modeling & Simulation, 15 (2017), pp. 110–156.
- [13] T. A. DRISCOLL, N. HALE, AND L. N. TREFETHEN, *Chebfun guide*, 2014.
- [14] E. FACCO, M. d'ERRICO, A. RODRIGUEZ, AND A. LAIO, *Estimating the intrinsic dimension of datasets by a minimal neighborhood information*, Scientific reports, 7 (2017), p. 12140.
- [15] H. FEDERER, *Curvature measures*, Transactions of the American Mathematical Society, 93 (1959), pp. 418–491.
- [16] N. FRANCO, A. MANZONI, AND P. ZUNINO, *A deep learning approach to reduced order modelling of parameter dependent partial differential equations*, Mathematics of Computation, 92 (2023), pp. 483–524.
- [17] S. FRESCA, L. DEDE, AND A. MANZONI, *A comprehensive deep learning-based approach to reduced order modeling of nonlinear time-dependent parametrized pdes*, Journal of Scientific Computing, 87 (2021), pp. 1–36.
- [18] F. J. GONZALEZ AND M. BALAJEWICZ, *Deep convolutional recurrent autoencoders for learning low-dimensional feature dynamics of fluid systems*, arXiv preprint arXiv:1808.01346, (2018).
- [19] S. GUGERCIN, A. C. ANTOULAS, AND C. BEATTIE,  *$H_2$  model reduction for large-scale linear dynamical systems*, SIAM journal on matrix analysis and applications, 30 (2008), pp. 609–638.
- [20] G. HARO, G. RANDALL, AND G. SAPIRO, *Translated poisson mixture model for stratification learning*, International Journal of Computer Vision, 80 (2008), pp. 358–374.
- [21] J. S. HESTHAVEN, G. ROZZA, B. STAMM, ET AL., *Certified reduced basis methods for parametrized partial differential equations*, vol. 590, Springer, 2016.
- [22] J. S. HESTHAVEN AND S. UBBIALI, *Non-intrusive reduced order modeling of nonlinear problems using neural networks*, Journal of Computational Physics, 363 (2018), pp. 55–78.
- [23] P. HOLMES, J. L. LUMLEY, G. BERKOOZ, AND C. W. ROWLEY, *Turbulence, coherent structures, dynamical systems and symmetry*, Cambridge university press, 2012.
- [24] K. KONTOLATI, S. GOSWAMI, G. EM KARNIADAKIS, AND M. D. SHIELDS, *Learning nonlinear operators in latent spaces for real-time predictions of complex dynamics in physical systems*, Nature Communications, 15 (2024), p. 5101.
- [25] J. M. LEE, *Introduction to Riemannian manifolds*, vol. 2, Springer, 2018.
- [26] K. LEE AND K. T. CARLBERG, *Model reduction of dynamical systems on nonlinear manifolds using deep convolutional autoencoders*, Journal of Computational Physics, 404 (2020), p. 108973.
- [27] R. J. LEVEQUE, *Finite difference methods for ordinary and partial differential equations: steady-state and time-dependent problems*, SIAM, 2007.
- [28] E. LEVINA AND P. BICKEL, *Maximum likelihood estimation of intrinsic dimension*, Advances in neural information processing systems, 17 (2004).
- [29] H. LIU, B. DAHAL, R. LAI, AND W. LIAO, *Generalization error guaranteed auto-encoder-based nonlinear model reduction for operator learning*, Applied and Computational Harmonic Analysis, 74 (2025), p. 101717.
- [30] H. LIU, A. HAVRILLA, R. LAI, AND W. LIAO, *Deep nonparametric estimation of intrinsic data structures by chart autoencoders: Generalization error and robustness*, Applied and Computational Harmonic Analysis, 68 (2024), p. 101602.
- [31] L. LU, P. JIN, G. PANG, Z. ZHANG, AND G. E. KARNIADAKIS, *Learning nonlinear operators via deeponet based on the universal approximation theorem of operators*, Nature Machine Intelligence, 3 (2021), pp. 218–229.
- [32] L. LU, P. JIN, G. PANG, Z. ZHANG, AND G. E. KARNIADAKIS, *Learning nonlinear operators via deeponet based on the universal approximation theorem of operators*, Nature Machine Intelligence, 3 (2021), pp. 218–229, <https://doi.org/10.1038/s42256-021-00302-5>.

- [33] B. MOORE, *Principal component analysis in linear systems: Controllability, observability, and model reduction*, IEEE transactions on automatic control, 26 (1981), pp. 17–32.
- [34] M. OHLBERGER AND S. RAVE, *Reduced basis methods: Success, limitations and future challenges*, arXiv preprint arXiv:1511.02021, (2015).
- [35] K. OONO AND T. SUZUKI, *Approximation and non-parametric estimation of resnet-type convolutional neural networks*, in International conference on machine learning, PMLR, 2019, pp. 4922–4931.
- [36] S. E. OTTO AND C. W. ROWLEY, *Linearly recurrent autoencoder networks for learning dynamics*, SIAM Journal on Applied Dynamical Systems, 18 (2019), pp. 558–593.
- [37] T. O'LEARY-ROSEBERRY, U. VILLA, P. CHEN, AND O. GHATTAS, *Derivative-informed projected neural networks for high-dimensional parametric maps governed by pdes*, Computer Methods in Applied Mechanics and Engineering, 388 (2022), p. 114199.
- [38] A. PASZKE, S. GROSS, F. MASSA, A. LERER, J. BRADBURY, G. CHANAN, T. KILLEEN, Z. LIN, N. GIMELSHAIN, L. ANTIGA, ET AL., *Pytorch: An imperative style, high-performance deep learning library*, Advances in neural information processing systems, 32 (2019).
- [39] A. PINKUS, *N-widths in Approximation Theory*, vol. 7, Springer Science & Business Media, 2012.
- [40] R. PINNAU, *Model reduction via proper orthogonal decomposition*, Model order reduction: theory, research aspects and applications, (2008), pp. 95–109.
- [41] C. PRUD'HOMME, D. V. ROVAS, K. VEROY, L. MACHIELS, Y. MADAY, A. T. PATERA, AND G. TURINICI, *Reliable real-time solution of parametrized partial differential equations: Reduced-basis output bound methods*, J. Fluids Eng., 124 (2002), pp. 70–80.
- [42] A. QUARTERONI, A. MANZONI, AND F. NEGRI, *Reduced basis methods for partial differential equations: an introduction*, vol. 92, Springer, 2015.
- [43] F. REGAZZONI, S. PAGANI, M. SALVADOR, L. DEDE', AND A. QUARTERONI, *Learning the intrinsic dynamics of spatio-temporal processes through latent dynamics networks*, Nature Communications, 15 (2024), p. 1834.
- [44] S. T. ROWEIS AND L. K. SAUL, *Nonlinear dimensionality reduction by locally linear embedding*, science, 290 (2000), pp. 2323–2326.
- [45] C. W. ROWLEY, T. COLONIUS, AND R. M. MURRAY, *Model reduction for compressible flows using pod and galerkin projection*, Physica D: Nonlinear Phenomena, 189 (2004), pp. 115–129.
- [46] G. ROZZA, D. B. P. HUYNH, AND A. T. PATERA, *Reduced basis approximation and a posteriori error estimation for affinely parametrized elliptic coercive partial differential equations: application to transport and continuum mechanics*, Archives of Computational Methods in Engineering, 15 (2008), p. 229.
- [47] M. TAKAMOTO, T. PRADITIA, R. LEITERITZ, D. MACKINLAY, F. ALESIANI, D. PFLÜGER, AND M. NIEPERT, *Pdebench: An extensive benchmark for scientific machine learning*, Advances in Neural Information Processing Systems, 35 (2022), pp. 1596–1611.
- [48] M. TAKAMOTO, T. PRADITIA, R. LEITERITZ, D. MACKINLAY, F. ALESIANI, D. PFLÜGER, AND M. NIEPERT, *PDEBench Datasets*, 2022, <https://doi.org/10.18419/DARUS-2986>, <https://doi.org/10.18419/DARUS-2986>.
- [49] J. B. TENENBAUM, V. d. SILVA, AND J. C. LANGFORD, *A global geometric framework for nonlinear dimensionality reduction*, science, 290 (2000), pp. 2319–2323.
- [50] C. THÄLE, *50 years sets with positive reach—a survey.*, Surveys in Mathematics and its Applications, 3 (2008), pp. 123–165.
- [51] A. TRAN, X. HE, D. A. MESSENGER, Y. CHOI, AND D. M. BORTZ, *Weak-form latent space dynamics identification*, Computer Methods in Applied Mechanics and Engineering, 427 (2024), p. 116998.
- [52] S. VOLKWEIN, *Model reduction using proper orthogonal decomposition*, Lecture Notes, Institute of Mathematics and Scientific Computing, University of Graz, 1025 (2011).
- [53] Q. WANG, J. S. HESTHAVEN, AND D. RAY, *Non-intrusive reduced order modeling of unsteady flows using artificial neural networks with application to a combustion problem*, Journal of computational physics, 384 (2019), pp. 289–307.
- [54] Z. WANG, D. XIAO, F. FANG, R. GOVINDAN, C. C. PAIN, AND Y. GUO, *Model identification of reduced order fluid dynamics systems using deep learning*, International Journal for Numerical Methods in Fluids, 86 (2018), pp. 255–268.
- [55] F. X.-F. YE, S. YANG, AND M. MAGGIONI, *Nonlinear model reduction for slow-fast stochastic systems near unknown invariant manifolds*, Journal of Nonlinear Science, 34 (2024), p. 22.
- [56] K. ZENG, C. E. P. DE JESÚS, A. J. FOX, AND M. D. GRAHAM, *Autoencoders for discovering manifold dimension and coordinates in data from complex dynamical systems*, Machine Learning: Science and Technology, 5 (2024), p. 025053.

888     **Appendix A. Helper Approximation Results.** In this section, we introduce  
 889 some approximation results for ReLU networks. They are later used to prove our main  
 890 results.

891     First, we establish that ReLU networks (with a given size) can implement the  
 892 identity function.

893     PROPOSITION A.1. *Let  $d, L, W \in \mathbb{N}$ , with  $W \geq 2d$ . Then there is a  $f_{\text{NN}} \in$   
 894  $\mathcal{F}_{\text{NN}}(d, d, L, W)$  such that  $f_{\text{NN}}(\mathbf{x}) = \mathbf{x}$  for all  $\mathbf{x} \in \mathbb{R}^d$ .*

895     *Proof of Proposition A.1.* We denote by  $I_d$  the identity function on  $\mathbb{R}^d$ . First,  
 896 suppose  $L = 1$  and  $W = 2d$ . Consider the weights and biases  $\mathfrak{W}_1 = \begin{bmatrix} I_d \\ -I_d \end{bmatrix}$ ,  $\mathfrak{b}_1 = \vec{0}_{2d}$ ,  
 897  $\mathfrak{W}_2 = [I_d \quad -I_d]$ , and  $\mathfrak{b}_2 = \vec{0}_d$ . Then the one-hidden-layer ReLU network with the  
 898 above weights and biases is in  $\mathcal{F}_{\text{NN}}(d, d, 1, 2d)$  and exactly implements the identity  
 899 (since  $\sigma(x) - \sigma(-x) = x$  for  $\sigma(x) = \max\{x, 0\}$ ).

900     Next, if  $L > 1$  and  $W = 2d$ , define  $\mathfrak{W}_1 = \begin{bmatrix} I_d \\ -I_d \end{bmatrix}$  and  $\mathfrak{W}_{L+1} = [I_d \quad -I_d]$ . Also for  
 901 any  $1 < \ell < L+1$ , let  $\mathfrak{W}_\ell = I_{2d}$ . Let all bias vectors be zero. Then the  $L$ -hidden-layer  
 902 ReLU network with the above weights and biases is in  $\mathcal{F}_{\text{NN}}(d, d, L, 2d)$  and exactly  
 903 implements the identity (as the output of the  $L$ th hidden layer is the same as the  
 904 output of the first hidden layer).

905     Finally, suppose  $W > 2d$ . It is easy to see that augmenting weight matrices  
 906  $\mathfrak{W}_1, \dots, \mathfrak{W}_L$  defined above with zeros in all new rows and columns (with the original  
 907 weight matrix being in the top left corner) does not affect output of the neural network.  
 908 Thus we have constructed an  $L$ -hidden-layer ReLU network with width exactly  $W$   
 909 that implements the identity in  $\mathbb{R}^d$ .  $\square$

910     PROPOSITION A.2. *Let  $f_{\text{NN}} \in \mathcal{F}_{\text{NN}}(d_{\text{in}}, d_{\text{out}}, L, W)$  be a feedforward ReLU net-  
 911 work. Suppose  $W > d_{\text{out}}$ .*

- 912       1. *For any  $L' > L$ , there exists a FNN  $g_{\text{NN}} \in \mathcal{F}_{\text{NN}}(d_{\text{in}}, d_{\text{out}}, L', 2W)$  with depth  
 913 exactly  $L'$  such that  $g_{\text{NN}}(\mathbf{x}) = f_{\text{NN}}(\mathbf{x})$  for all  $\mathbf{x} \in \mathbb{R}^{d_{\text{in}}}$ .*
- 914       2. *Let  $A = \prod_{i=1}^{d_{\text{out}}} [a_i, b_i]$  be a cube. Then there is a  $h_{\text{NN}} \in \mathcal{F}_{\text{NN}}(d_{\text{in}}, d_{\text{out}}, L +  
 915 2, \max\{d_{\text{out}}, W\})$  such that for all  $\mathbf{x} \in \mathbb{R}^{d_{\text{in}}}$ ,  $h_{\text{NN}}(\mathbf{x})$  is vector valued and  
 916  $= (h_{\text{NN}}^{(1)}(\mathbf{x}), \dots, h_{\text{NN}}^{(d_{\text{out}})}(\mathbf{x})) \in \mathbb{R}^{d_{\text{out}}}$  with  $h_{\text{NN}}^{(i)}(\mathbf{x}) = \min\{\max\{f_{\text{NN}}^{(i)}(\mathbf{x}), a_i\}, b_i\}$ .*
- 917       3. *Let  $f_{\text{NN}}^1, \dots, f_{\text{NN}}^m$  be FNNs such that  $f_{\text{NN}}^j \in \mathcal{F}_{\text{NN}}(d_{\text{in}}, d_{\text{out}}, L_j, W_j)$  for all  $j \in [m]$ . Then there is a  $g_{\text{NN}} \in \mathcal{F}_{\text{NN}}(d_{\text{in}}, d_{\text{out}}, \sum_{j=1}^m L_j, 2d_{\text{in}} + 2d_{\text{out}} + \max_j W_j)$   
 918 such that for all  $\mathbf{x} \in \mathbb{R}^{d_{\text{in}}}$ , we have  $g_{\text{NN}}(\mathbf{x}) = \sum_{j=1}^m f_{\text{NN}}^j(\mathbf{x})$ .*

920     Proposition A.2 gives rise to basic properties of feedforward ReLU networks. The  
 921 first property is that, any neural network can be realized by a deeper one; The second  
 922 property says that, feedforward ReLU networks allows one to truncate or restrict  
 923 its output to be in the set of  $A = \prod_{i=1}^{d_2} [a_i, b_i]$ ; The third property allows one to  
 924 concatenate several feedforward ReLU networks and then output its summation.

925     *Proof of Proposition A.2. Part 1:* Denote the number of neurons in the last  
 926 hidden layer of  $f_{\text{NN}}$  by  $r$  (note  $r \leq W$ ). Let  $L_f$  be the number of layers of  $f_{\text{NN}}$ , and  
 927 by the hypothesis we have  $L_f < L'$ .

928     First, consider the function  $\phi_1 : \mathbb{R}^r \rightarrow \mathbb{R}^{2r}$  defined by  $\phi_1(x_1, \dots, x_r) = [\sigma(x_1),  
 929 \sigma(-x_1), \dots, \sigma(x_r), \sigma(-x_r)]$ . Clearly,  $\phi_1$  can be implemented by a single ReLU layer.  
 930 Next, consider the function  $\phi_2$  on  $\mathbb{R}^{2r}$  defined by  $\phi_2(x_1, \dots, x_{2r}) = [\sigma(x_1), \dots, \sigma(x_{2r})]$ .  
 931 Clearly,  $\phi_2$  can also be implemented by a single ReLU layer. Note that  $\bigcirc_{i=1}^n \phi_2 =  
 932 \phi_2$  for any  $n \in \mathbb{N}$ . We will construct  $g_{\text{NN}}$  from  $f_{\text{NN}}$  by appending a hidden layer

933 implementing  $\phi_1$ , appending  $L' - L_f - 1$  hidden layers implementing  $\phi_2$ , and then  
 934 adjusting the weights in the final layer of  $f_{\text{NN}}$ .

935 Denote the final weight matrix of  $f_{\text{NN}}$  by  $\mathfrak{W} \in \mathbb{R}^{r \times d_{\text{out}}}$ . We will adjust this  
 936 matrix to handle inputs of size  $2r$ . Define the matrix  $\hat{\mathfrak{W}} \in \mathbb{R}^{2r \times d_{\text{out}}}$  by  $\hat{\mathfrak{W}}_{i,j} =$   
 937  $(-1)^{j+1} \mathfrak{W}_{i,\lceil j/2 \rceil}$ , for all  $i \in [d_{\text{out}}]$ ,  $j \in [2r]$ . This means  $\hat{\mathfrak{W}}(x_1, x_2, \dots, x_{2r-1}, x_{2r}) =$   
 938  $\mathfrak{W}(x_1 - x_2, \dots, x_{2r-1} - x_{2r})$ .

939 Note that for all  $i \in [r]$ , we have that  $\sigma(x_i) - \sigma(-x_i) = x_i$ . This means (where  
 940 we use denote multiplication vector multiplication by parenthesis)

$$\begin{aligned} 941 \quad & \hat{\mathfrak{W}} \left( (\bigcirc_{i=1}^{L'-L_f-1} \phi_2) \circ \phi_1(x_1, \dots, x_r) \right) \\ 942 \quad & = \hat{\mathfrak{W}}(\sigma(x_1), \sigma(-x_1), \dots, \sigma(x_r), \sigma(-x_r)) = \mathfrak{W}(x_1, \dots, x_r). \end{aligned}$$

943 Finally, we can describe  $g_{\text{NN}}$ . Before the final weight matrix of  $f_{\text{NN}}$ , we will  
 944 append one layer implementing  $\phi_1$  and  $L' - L_f - 1$  layers implementing  $\phi_2$  to the net-  
 945 work. Then, we change the final weight matrix to  $\hat{\mathfrak{W}}$ . Thus,  $g_{\text{NN}}$  and  $f_{\text{NN}}$  implement  
 946 the same function, and  $g_{\text{NN}} \in \mathcal{F}(d_{\text{in}}, d_{\text{out}}, L', 2W)$ .

947 **Part 2:** Note that the function  $b - \sigma(b - a - \sigma(x - a)) = \min(\max(x, a), b)$ .

948 Let  $\mathbf{b} \in \mathbb{R}^{d_{\text{out}}}$  be the bias vector for the final layer. Define a new bias vector  
 949  $\tilde{\mathbf{b}} \in \mathbb{R}^{d_{\text{out}}}$  by  $\tilde{\mathbf{b}}_i = \mathbf{b}_i - a_i$ , for all  $i \in [d_{\text{out}}]$ .

950 If we denote  $\mathbf{a} = [a_1, \dots, a_{d_{\text{out}}}]$  and  $\mathbf{b} = [b_1, \dots, b_{d_{\text{out}}}]$ , then note that for all  
 951  $x \in \mathbb{R}^{d_{\text{in}}}$ ,

$$952 \quad (\text{A.1}) \quad -\sigma(-(\sigma(f_{\text{NN}}(x) - \mathbf{a}) + \mathbf{b} - \mathbf{a}) + \mathbf{b} = \min(\max(f_{\text{NN}}(x), \mathbf{a}), \mathbf{b}),$$

953 where  $\sigma$  and  $\max, \min$  are applied componentwise. We can implement the LHS of  
 954 Equation (A.1) starting from  $f_{\text{NN}}$  using the following steps (where  $I_{d_{\text{out}}}$  is the identity  
 955 matrix for  $\mathbb{R}^{d_{\text{out}}}$ ):

- 956 1. Modify the last bias vector of  $f_{\text{NN}}$  to  $\tilde{\mathbf{b}}$ . This implements  $f_{\text{NN}}(x) - \mathbf{a}$ .
- 957 2. Add a new layer with weight matrix equal to  $-I_{d_{\text{out}}}$  and bias vector equal to  
 958  $\mathbf{b} - \mathbf{a}$ . This implements  $-\sigma(f_{\text{NN}}(x) - \mathbf{a}) + \mathbf{b} - \mathbf{a}$ .
- 959 3. Add a new layer with weight matrix equal to  $-I_{d_{\text{out}}}$  (identity) and bias equal  
 960 to  $\mathbf{b}$ . This implements the LHS of (A.1).  $\square$

961 Thus we have constructed a neural network  $h_{\text{NN}}$  by modifying the bias vector in  
 962 the final layer of  $f_{\text{NN}}$  and adding two more layers.

963 **Part 3:** We add extra layers that store the running sum. For every network  $f_{\text{NN}}^j$ ,  
 964 we add  $2d_{\text{in}}$  extra neurons that store the positive and negative values of the inputs  
 965 every layer. For the final layer of each network, we store the positive and negative  
 966 values of the output with  $2d_{\text{out}}$  neurons.

967 We then concatenate each of these networks. We collect the running sum in the  
 968  $2d_{\text{out}}$  neurons, and we preserve the input in the  $2d_{\text{in}}$  extra intermediate neurons. The  
 969 depth of this concatenated network is the sum of the depths of each network, and the  
 970 width is the maximum width of the networks plus  $2d_{\text{out}} + 2d_{\text{in}}$  neurons.

971 Each layer will need  $2d_{\text{out}}$  neurons to capture the positive and negative values of  
 972 the  $d_{\text{out}}$  outputs in the sum. Specifically, we modify the final layer of  $f_{\text{NN}}$  to have  
 973 double the number of output neurons as  $d_{\text{out}}$ .

Now, let's consider to approximate Hölder functions using ReLU networks. For

any function  $f : \Omega \rightarrow \mathbb{R}$ , we define it 1-Hölder norm:

$$\|f\|_{\mathcal{H}^1(\Omega; \mathbb{R})} = \max \left\{ \sup_{x \in \Omega} |f(x)|, \sup_{x, y \in \Omega} \frac{|f(x) - f(y)|}{\|x - y\|_2} \right\}$$

974 and

$$\mathcal{H}^1(\Omega; \mathbb{R}) = \{f : \|f\|_{\mathcal{H}^1(\Omega; \mathbb{R})} < \infty\}.$$

976 PROPOSITION A.3. Let  $f^* \in \mathcal{H}^1([0, 1]^d; \mathbb{R})$ . Then,  $\forall M \in \mathbb{N}$ , there is a  $f_{\text{NN}} \in$   
977  $\mathcal{F}_{\text{NN}}(d, 1, L, W)$  such that

$$978 \|f^* - f_{\text{NN}}\|_{L^\infty([0, 1]^d; \mathbb{R})} \leq CM^{-\frac{1}{d}}.$$

979 Here  $L = O(\log(M))$  and  $W = O(M)$ , where the big  $O$  and  $C$  hides constants only  
980 depending on  $d$  and  $\|f^*\|_{\mathcal{H}^1([0, 1]^d; \mathbb{R})}$ .

981 *Proof.* This is a restatement of [35, Lemma 7].  $\square$

982 We extend Proposition A.3 to functions with vector-valued outputs.

983 PROPOSITION A.4. Let  $f^* \in \mathcal{H}^1([0, 1]^{d_{in}}; A)$ , where  $A = \mathbb{R}^{d_{out}}$  is the Euclidean  
984 space, or  $A = \prod_{i=1}^{d_{out}} [a_i, b_i]$  is a cube. Then,  $\forall M \in \mathbb{N}, \exists f_{\text{NN}} \in \mathcal{F}_{\text{NN}}(d_{in}, d_{out}, L, W)$   
985 such that

$$986 \|f^* - f_{\text{NN}}\|_{L^\infty([0, 1]^{d_{in}}; \mathbb{R}^{d_{out}})} \leq CM^{-\frac{1}{d_{in}}}.$$

987 In addition, the range of  $f_{\text{NN}}$  is contained in  $A$ . Here  $L = O(\log(M))$  and  $W = O(M)$ ,  
988 where big  $O$  and  $C$  hide constants depending on  $d_{in}$ ,  $d_{out}$ , and  $\|f^*\|_{\mathcal{H}^1([0, 1]^{d_{in}}; \mathbb{R}^{d_{out}})}$ .

989 *Proof of Proposition A.4.* For all  $i \in [d_{out}]$ , we denote the  $i$ th component function  
990 of  $f^*$  by  $f^i$ . Then since  $f^i \in \mathcal{H}^1([0, 1]^{d_{in}}, \mathbb{R})$ , by Proposition A.3 there exists an FNN  
991  $f_{\text{NN}}^i \in \mathcal{F}_{\text{NN}}(d_{in}, 1, L_i, W_i)$  such that  $\|f^i - f_{\text{NN}}^i\|_{L^\infty([0, 1]^{d_{in}}; \mathbb{R})} \leq C_i M^{-\frac{1}{d_{in}}}$  for some  $C_i$   
992 depending on  $d$  and  $\|f^i\|_{\mathcal{H}^1([0, 1]^{d_{in}}; \mathbb{R})}$ . Let  $L = \max_i L_i$ , and  $W = \sum_{i=1}^{d_{out}} W_i$ . We  
993 will add extra layers implementing the identity so that each  $f_{\text{NN}}^i$  has exactly  $L$  layers.  
994 This is done by applying Proposition A.2 Part 1, which will double the width of each  
995  $f_{\text{NN}}^i$ , leaving the term inside the big  $O$  notation unchanged.

996 Now, define  $f_{\text{NN}}$  as the network implementing each of the  $f_{\text{NN}}^i$  as component  
997 functions. Then  $f_{\text{NN}} \in \mathcal{F}_{\text{NN}}(d_{in}, d_{out}, L, W)$  with  $L = O(\log(M))$  and  $W = O(M)$ .  
998 For all  $x \in [0, 1]^{d_{in}}$  we estimate

$$\begin{aligned} 999 \|f^*(x) - f_{\text{NN}}(x)\|_{\mathbb{R}^{d_{out}}} &= \left( \sum_{i=1}^{d_{out}} (f^i(x) - f_{\text{NN}}^i(x))^2 \right)^{\frac{1}{2}} \\ 1000 &\leq \left( \sum_{i=1}^{d_{out}} C_i^2 M^{-\frac{2}{d_{in}}} \right)^{\frac{1}{2}} \\ 1001 &\leq \left( d_{out} \left( \max_i C_i^2 \right) M^{-\frac{2}{d_{in}}} \right)^{\frac{1}{2}} \\ 1002 &= \sqrt{d_{out} \left( \max_i C_i^2 \right)} M^{-\frac{1}{d_{in}}}. \end{aligned}$$

1003 Finally, if  $A = \prod_{i=1}^{d_{out}} [a_i, b_i]$ , we replace each  $f_{\text{NN}}^i$  by  $\sigma(-\sigma(-f_{\text{NN}}^i + b_i - a_i) + a_i)$  so  
1004 that for any  $x \in [0, 1]^{d_{in}}$ ,  $f_{\text{NN}}^i(x) \in [a_i, b_i]$  and thus  $f_{\text{NN}}(x) \in A$ .  $\square$

We use Proposition A.4 to approximate target functions up to  $\epsilon$  accuracy. Setting the approximation error as  $\epsilon$  in Proposition A.4 gives rise to the following proposition.

**PROPOSITION A.5.** *Let  $\epsilon > 0$  and  $f^* \in \mathcal{H}^1([0, 1]^{d_{in}}; A)$ , where  $A = \mathbb{R}^{d_{out}}$  is Euclidean space or  $A = \prod_{i=1}^{d_{out}} [a_i, b_i]$  is a cube. Then  $\exists f_{\text{NN}} \in \mathcal{F}_{\text{NN}}(d_{in}, d_{out}, L, W)$  such that*

$$\|f^* - f_{\text{NN}}\|_{L^\infty([0, 1]^{d_{in}}; \mathbb{R}^{d_{out}})} < \epsilon.$$

In addition, the range of  $f_{\text{NN}}$  is contained in  $A$ . Here  $L = O(\log(\epsilon^{-1}))$  and  $W = O(\epsilon^{-d_{in}})$ , where the big  $O$  hides constants depending only on  $d_{in}$ ,  $d_{out}$ , and  $\|f^*\|_{\mathcal{H}^1([0, 1]^{d_{in}}; \mathbb{R}^{d_{out}})}$ .

*Proof of Proposition A.5.* This follows directly from Proposition A.4.  $\square$

WeldNet uses a composition of neural networks (i.e. propagators and transcoders) that are sequentially trained. In order to control the approximation error of such a composition, we use the following proposition.

**PROPOSITION A.6 (Composition).** *Let  $f_1, \dots, f_n$  be Lipschitz functions, such that  $f_j : A_{j-1} \rightarrow A_j$  for all  $j \in [n]$ . Suppose for each  $j \in [n]$  and  $\epsilon_j > 0$ , there is an approximating function  $\hat{f}_j : A_{j-1} \rightarrow A_j$  such that*

$$\sup_{x \in A_{j-1}} \|f_j(x) - \hat{f}_j(x)\| < \epsilon_j.$$

Then

$$(A.2) \quad \sup_{x \in A_0} \|(\bigcirc_{j=1}^n f_j)(x) - (\bigcirc_{j=1}^n \hat{f}_j)(x)\| < \sum_{j=1}^n \text{Lip}(\bigcirc_{\ell=j+1}^n f_\ell) \epsilon_i.$$

*Proof of Proposition A.6.* We have

$$\begin{aligned} & \sup_{x \in A_0} \|(\bigcirc_{j=1}^n f_j)(x) - (\bigcirc_{j=1}^n \hat{f}_j)(x)\| \\ & \leq \sup_{x \in A_0} \sum_{j=1}^n \|(\bigcirc_{\ell=j+1}^n f_\ell \circ f_j \circ \bigcirc_{\ell=1}^{j-1} \hat{f}_\ell)(x) - (\bigcirc_{\ell=j+1}^n f_\ell \circ \hat{f}_j \circ \bigcirc_{\ell=1}^{j-1} \hat{f}_\ell)(x)\| \\ & \leq \sum_{j=1}^n \sup_{x \in A_0} \|(\bigcirc_{\ell=j+1}^n f_\ell \circ f_j \circ \bigcirc_{\ell=1}^{j-1} \hat{f}_\ell)(x) - (\bigcirc_{\ell=j+1}^n f_\ell \circ \hat{f}_j \circ \bigcirc_{\ell=1}^{j-1} \hat{f}_\ell)(x)\| \\ & \leq \sum_{j=1}^n \text{Lip}(\bigcirc_{\ell=j+1}^n f_\ell) \sup_{x \in A_0} \|(f_j \circ \bigcirc_{\ell=1}^{j-1} \hat{f}_\ell)(x) - (\hat{f}_j \circ \bigcirc_{\ell=1}^{j-1} \hat{f}_\ell)(x)\| \\ & \leq \sum_{j=1}^n \text{Lip}(\bigcirc_{\ell=j+1}^n f_\ell) \sup_{y \in A_{j-1}} \|f_j(y) - \hat{f}_j(y)\| \\ & \leq \sum_{j=1}^n \text{Lip}(\bigcirc_{\ell=j+1}^n f_\ell) \epsilon_j. \end{aligned} \quad \square$$

## Appendix B. Proof of Main Results.

**B.1. Proof of Lemma 3.7.** In this section, we prove Lemma 3.7. In the following, the windows are the same as the segments. First, we prove lemmas that establish approximation results for the autoencoders, propagators, and transcoders.

1031 LEMMA B.1 (Autoencoder Approximation). *Suppose Assumption 2 holds, and  
 1032 let  $0 < \epsilon_1, \epsilon_2 < \min\{1, \tau(\mathcal{M}([0, T]))/2\}$ . Then for any  $i \in [S]$ , there exists  $\mathcal{E}_{\text{NN}}^i \in$   
 1033  $\mathcal{F}_{\text{NN}}(D, d+1, L_{\mathcal{E}^i}, W_{\mathcal{E}^i})$  and  $\mathcal{D}_{\text{NN}}^i \in \mathcal{F}_{\text{NN}}(d+1, D, L_{\mathcal{D}^i}, W_{\mathcal{D}^i})$  with parameters*

$$\begin{aligned} 1034 \quad L_{\mathcal{E}^i} &= O(\log^2(\epsilon_1^{-1})), \quad W_{\mathcal{E}^i} = O(D\epsilon_1^{-(d+1)}) \\ 1035 \quad L_{\mathcal{D}^i} &= O(\log(\epsilon_2^{-1})), \quad W_{\mathcal{D}^i} = O(D\epsilon_2^{-(d+1)}) \end{aligned}$$

1036 such that

$$1037 \quad \sup_{\mathbf{x} \in \mathcal{M}([s_i, s_{i+1}])} \|\mathcal{E}_{\text{NN}}(\mathbf{x}) - \mathcal{E}_*^i(\mathbf{x})\|_{\mathbb{R}^{d+1}} \leq \epsilon_1 \quad \text{and} \quad \sup_{\mathbf{z} \in \mathcal{Z}([s_i, s_{i+1}])} \|\mathcal{D}_{\text{NN}}(\mathbf{z}) - \mathcal{D}_*(\mathbf{z})\|_{\mathbb{R}^D} \leq \epsilon_2.$$

1038 Moreover, the range of  $\mathcal{E}_{\text{NN}}$  is contained in  $\mathcal{Z}([0, T])$ . The constants hidden in  $O$   
 1039 depend on (for each  $i \in [S]$ )  $\log D$ ,  $d$ ,  $\tau(\mathcal{M}([s_i, s_{i+1}]))$ ,  $s_{i+1} - s_i$ ,  $\text{Lip}(\mathcal{E}_*^i)$ ,  $\text{Lip}(\mathcal{D}_*^i)$ ,  
 1040 the volume of  $\mathcal{M}([s_i, s_{i+1}])$ , and  $\sup_{\mathbf{x} \in \mathcal{M}([s_i, s_{i+1}])} \|\mathbf{x}\|_{\mathbb{R}^D}$ .

1041 Proof of Lemma B.1. This follows from [30, Lemma 6 and 8], with one additional step to restrict the range of  $\mathcal{E}_{\text{NN}}$ . This is done by applying Proposition A.2(Part 2), which increases the number of layers by 2.  $\square$   
 1042

1043 LEMMA B.2 (Propagator Approximation with Latent Dynamics). *Suppose there  
 1044 is a Lipschitz function  $g : \mathcal{Z}([a, b]) \rightarrow \mathbb{R}^{d+1}$  such that for all  $\mathbf{z}(a) \in \mathcal{Z}(a)$ , if we denote  
 1045  $\mathbf{z}(t) = \mathcal{P}_*(\mathbf{z}(a), t - a)$  for all  $t \in [a, b]$ , then these latent codes satisfy:*

$$1047 \quad (\text{B.1}) \quad \frac{\partial \mathbf{z}}{\partial t}(t) = g(\mathbf{z}(t)).$$

1048 Denote  $\text{Lip}(g) = \sup_{t \in [a, b]} \text{Lip}_{\mathcal{Z}(t)}(g(\mathbf{z}(t)))$  and  $T = b - a$ . Suppose  $a = t_1 < \dots <$   
 1049  $t_{T-1} < t_T = b$  is an equally spaced time grid, i.e.  $t_{k-1} - t_k = \Delta t$  is the same for  
 1050 all  $k \in [\mathsf{T} - 1]$ . For any  $\epsilon > 0$  and  $\mathsf{T} > 1 + \text{Lip}(g)T^2 e^{\text{Lip}(g)T} \epsilon^{-1}$ , then there is a  
 1051  $\mathcal{P}_{\text{NN}} \in \mathcal{F}_{\text{NN}}(d+1, d+1, L, W)$ , with range contained in  $\mathcal{Z}([a, b])$ , such that for any  
 1052  $k \in [\mathsf{T} - 1]$  and  $i \in [\mathsf{T} - k]$ ,

$$1053 \quad \sup_{\mathbf{z}(t_k) \in \mathcal{Z}(t_k)} \|(\bigcirc_{j=1}^i \mathcal{P}_{\text{NN}})(\mathbf{z}(t_k)) - \mathcal{P}_*(\mathbf{z}(t_k), t_{k+i} - t_k)\|_{\mathbb{R}^{d+1}} < \epsilon.$$

1054 Here  $L = O(\log(\frac{1}{\epsilon}))$ ,  $W = O(\epsilon^{-(d+1)})$ , where the constants hidden in  $O$  depend on  
 1055  $d$ ,  $\text{Lip}(g)$ ,  $\|g\|_{L^\infty(\mathcal{Z}([0, T]))}$ , and  $T$ .

1056 Proof of Lemma B.2. Since  $g$  is Lipschitz, by Proposition A.5 there is a  $g_{\text{NN}} \in$   
 1057  $g_{\text{NN}}(d+1, d+1, L_g, W_g)$  such that

$$1058 \quad \sup_{\mathbf{z} \in \mathcal{Z}([a, b])} \|g_{\text{NN}}(\mathbf{z}) - g(\mathbf{z})\|_{\mathbb{R}^{d+1}} < \delta := \frac{\epsilon}{6Te^{\text{Lip}(g)T}}.$$

1059 Here  $L_g = cO(\log(\frac{1}{\epsilon}))$  and  $W_g = O(\epsilon^{-(d+1)})$ . We define  $\mathcal{P}_{\text{NN}} : \mathcal{Z}([a, b]) \rightarrow \mathcal{Z}([a, b])$   
 1060 component-wise. For all  $i \in [d]$ , define the  $i$ th component of  $\mathcal{P}_{\text{NN}}$  as

$$1061 \quad \mathcal{P}_{\text{NN}}^i(\mathbf{z}) = \max(0, \min(\mathbf{z}_i + \Delta t(g_{\text{NN}}(\mathbf{z})))^i)$$

1062 and the  $(d+1)$ th component as  $\mathcal{P}_{\text{NN}}^{d+1}(\mathbf{z}_1, \dots, \mathbf{z}_{d+1}) = \mathbf{z}_{d+1} + \Delta t$ . Then note we  
 1063 can implement  $\mathcal{P}_{\text{NN}}$  as a neural network by modifying the last layer of  $\mathcal{F}_{\text{NN}}$  and  
 1064 then adding two layers using Proposition A.2(Part 2). Then  $\mathcal{P}_{\text{NN}} \in g_{\text{NN}}(d+1, d+1,  
 1065 1, L_{\mathcal{P}}, W_{\mathcal{P}})$  with  $L_{\mathcal{P}} = O(\log(\frac{1}{\epsilon}))$  and  $W_{\mathcal{P}} = O(\epsilon^{-(d+1)})$ .

We next use classical argument for the convergence of the Euler method for ODEs to finish the proof (c.f. [27] Chapter 6). Let  $\mathbf{z}(a) \in \mathcal{Z}(a)$ . For all  $k \in [\mathsf{T}]$ , define  $\mathbf{z}(t_k) = \mathcal{P}_*(\mathbf{z}(a), t_k - a)$ , and define  $\hat{\mathbf{z}}_k$  using the iterative formula

$$\hat{\mathbf{z}}_k = \mathcal{P}_{\text{NN}}(\hat{\mathbf{z}}_{k-1}), \quad \hat{\mathbf{z}}_1 = \mathbf{z}(a_i).$$

Note that

$$\begin{aligned} \mathbf{z}(t_{k+1}) &= \mathbf{z}(t_k) + \Delta t \frac{\partial \mathbf{z}}{\partial t}(t_k) + \int_{t_k}^{t_{k+1}} \left( \frac{\partial \mathbf{z}}{\partial t}(s) - \frac{\partial \mathbf{z}}{\partial t}(t_k) \right) ds \\ &= \mathbf{z}(t_k) + \Delta t g(\mathbf{z}(t_k)) + \int_{t_k}^{t_{k+1}} (g(\mathbf{z}(s)) - g(\mathbf{z}(t_k))) ds. \end{aligned}$$

This means

$$\begin{aligned} \left\| \frac{1}{\Delta t} (\mathbf{z}(t_{k+1}) - \mathbf{z}(t_k)) - g(\mathbf{z}(t_k)) \right\|_{\mathbb{R}^d} &= \frac{1}{\Delta t} \left\| \int_{t_k}^{t_{k+1}} (g(\mathbf{z}(s)) - g(\mathbf{z}(t_k))) ds \right\|_{\mathbb{R}^{d+1}} \\ &\leq \frac{1}{\Delta t} \int_{t_k}^{t_{k+1}} \|g(\mathbf{z}(s)) - g(\mathbf{z}(t_k))\|_{\mathbb{R}^{d+1}} ds \\ &\leq \frac{1}{\Delta t} \int_{t_k}^{t_{k+1}} \text{Lip}(g) \|\mathbf{z}(s) - \mathbf{z}(t_k)\|_{\mathbb{R}^{d+1}} ds \\ &\leq \frac{\text{Lip}(g)}{\Delta t} \int_{t_k}^{t_{k+1}} \|g\|_{L^\infty(\mathcal{Z}([a,b]))} (s - t_k) ds \\ &\leq \frac{\text{Lip}(g)}{\Delta t} \|g\|_{L^\infty(\mathcal{Z}([a,b]))} \frac{(\Delta t)^2}{2} \leq C \Delta t, \end{aligned}$$

where  $C = \text{Lip}(g) \|g\|_{L^\infty(\mathcal{Z}([a,b]))} / 2$ .

Now, let  $\mathbf{e}_k = \frac{1}{\Delta t} (\mathbf{z}(t_{k+1}) - \mathbf{z}(t_k)) - g_{\text{NN}}(\mathbf{z}(t_k))$  be the step  $k$  truncation error for Euler method of solving the ODE (B.1) while the governing equation is  $g_{\text{NN}}$  instead of  $g$ . We have

$$\mathbf{z}(t_{k+1}) = \mathbf{z}(t_k) + \Delta t g_{\text{NN}}(\mathbf{z}(t_k)) + \Delta t \mathbf{e}_k.$$

This local truncation error satisfies the error bound:

$$\begin{aligned} \|\mathbf{e}_k\|_{\mathbb{R}^{d+1}} &= \left\| \frac{1}{\Delta t} (\mathbf{z}(t_{k+1}) - \mathbf{z}(t_k)) - g_{\text{NN}}(\mathbf{z}(t_k)) \right\|_{\mathbb{R}^{d+1}} \\ &\leq \left\| \frac{1}{\Delta t} (\mathbf{z}(t_{k+1}) - \mathbf{z}(t_k)) - g(\mathbf{z}(t_k)) \right\|_{\mathbb{R}^{d+1}} + \|g(\mathbf{z}(t_k)) - g_{\text{NN}}(\mathbf{z}(t_k))\|_{\mathbb{R}^{d+1}} \\ &\leq C \Delta t + \delta. \end{aligned}$$

Then we decompose

$$\begin{aligned} \mathbf{z}(t_{k+1}) - \hat{\mathbf{z}}_{k+1} &= \mathbf{z}(t_k) + \Delta t g_{\text{NN}}(\mathbf{z}(t_k)) + \Delta t \mathbf{e}_k - \hat{\mathbf{z}}_k - \Delta t g_{\text{NN}}(\hat{\mathbf{z}}_k) \\ &= (\mathbf{z}(t_k) - \hat{\mathbf{z}}_k) + \Delta t (g_{\text{NN}}(\mathbf{z}(t_k)) - g_{\text{NN}}(\hat{\mathbf{z}}_k)) + \Delta t \mathbf{e}_k. \end{aligned}$$

This means

$$\|\mathbf{z}(t_{k+1}) - \hat{\mathbf{z}}_{k+1}\|_{\mathbb{R}^{d+1}} \leq \|\mathbf{z}(t_k) - \hat{\mathbf{z}}_k\|_{\mathbb{R}^{d+1}} + \Delta t \|g_{\text{NN}}(\mathbf{z}(t_k)) - g_{\text{NN}}(\hat{\mathbf{z}}_k)\|_{\mathbb{R}^{d+1}} + \Delta t \|\mathbf{e}_k\|_{\mathbb{R}^{d+1}}.$$

Now  $\|g_{\text{NN}}(\mathbf{z}(t_k)) - g_{\text{NN}}(\hat{\mathbf{z}}_k)\|_{\mathbb{R}^{d+1}} \leq \|g_{\text{NN}}(\mathbf{z}(t_k)) - g(\mathbf{z}(t_k))\|_{\mathbb{R}^{d+1}} + \|g(\mathbf{z}(t_k)) - g(\hat{\mathbf{z}}_k)\|_{\mathbb{R}^{d+1}} + \|g(\hat{\mathbf{z}}_k) - g_{\text{NN}}(\hat{\mathbf{z}}_k)\|_{\mathbb{R}^{d+1}} \leq 2\delta + \text{Lip}(g) \|\mathbf{z}(t_k) - \hat{\mathbf{z}}_k\|_{\mathbb{R}^{d+1}}$ . So we have the inequality

1096       $\|\mathbf{z}(t_{k+1}) - \hat{\mathbf{z}}_{k+1}\|_{\mathbb{R}^{d+1}} \leq (1 + \Delta t \text{Lip}(g)) \|\mathbf{z}(t_k) - \hat{\mathbf{z}}_k\|_{\mathbb{R}^{d+1}} + \Delta t(2\delta + \|\mathbf{e}_k\|_{\mathbb{R}^{d+1}}).$

1097      Since  $\|\mathbf{z}(t_1) - \hat{\mathbf{z}}_1\|_{\mathbb{R}^d} = 0$ , this implies

1098       $\|\mathbf{z}(t_{k+1}) - \hat{\mathbf{z}}_{k+1}\|_{\mathbb{R}^{d+1}}$   
1099       $\leq \Delta t \sum_{i=1}^k (1 + \Delta t \text{Lip}(g))^{k-i} (2\delta + \|\mathbf{e}_i\|_{\mathbb{R}^{d+1}}) \leq e^{\text{Lip}(g)T} \Delta t \sum_{i=1}^k (2\delta + \|\mathbf{e}_i\|_{\mathbb{R}^{d+1}}).$

1100      Therefore (since  $k\Delta t < T$  and  $\|\mathbf{e}_k\|_{\mathbb{R}^{d+1}} \leq C\Delta t + \delta$ )

1101       $\|\mathbf{z}(t_{k+1}) - \hat{\mathbf{z}}_{k+1}\|_{\mathbb{R}^{d+1}} \leq e^{\text{Lip}(g)T} \Delta t \sum_{i=1}^k (3\delta + C\Delta t)$   
1102       $= e^{\text{Lip}(g)T} (k\Delta t)(3\delta + C\Delta t) \leq 3Te^{\text{Lip}(g)T} \delta + CT e^{\text{Lip}(g)T} \Delta t.$

1103      Now since  $\Delta t = \frac{T}{\tau-1}$ , if we have  $\tau-1 > \frac{2CT^2 e^{\text{Lip}(g)T}}{\epsilon}$ , then

1104       $3Te^{\text{Lip}(g)T} \delta + CT e^{\text{Lip}(g)T} \Delta t$   
1105       $\leq \frac{3Te^{\text{Lip}(g)T} \epsilon}{6Te^{\text{Lip}(g)T}} + \frac{CT^2 e^{\text{Lip}(g)T}}{\tau-1} \leq \frac{\epsilon}{2} + (CT^2 e^{\text{Lip}(g)T}) \frac{\epsilon}{2CT^2 e^{\text{Lip}(g)T}} = \epsilon. \square$

1106      LEMMA B.3 (Transcoder Approximation). *Suppose Assumption 2 holds. Fix an*  
1107 *i ∈ [S], and for any ε > 0, there exists a*  $\mathcal{T}_{\text{NN}} \in \mathcal{F}_{\text{NN}}(d+1, d+1, L_{\mathcal{T}}, W_{\mathcal{T}})$  *with*  
1108 *parameters*

1109       $L_{\mathcal{T}} = O(\log(\epsilon^{-1})), \quad W_{\mathcal{T}} = O(\epsilon^{-d})$

1110      such that

1111       $\sup_{\mathbf{z}(s_{i+1}) \in \mathcal{Z}(s_{i+1})} \|\mathcal{T}_{\text{NN}}(\mathbf{z}(s_{i+1})) - \mathcal{E}_*^{i+1}(\mathcal{D}_*^i(\mathbf{z}(s_{i+1})))\|_{\mathbb{R}^{d+1}} < \epsilon$

1112      Moreover, the range of  $\mathcal{T}_{\text{NN}}$  is contained in  $\mathcal{Z}(s_{i+1})$ . Here,  $L_{\mathcal{T}} = O(\log(\frac{1}{\epsilon}))$  and  
1113  $W_{\mathcal{T}} = O(\epsilon^{-\frac{1}{d}})$ . The constants hidden in  $O$  depend on  $d, \tau_{\mathcal{M}}, \text{Lip}_{\mathcal{M}(s_{i+1})}(\mathcal{E}_*^{i+1}),$   
1114  $\text{Lip}_{\mathcal{Z}(s_{i+1})}(\mathcal{D}_*^i)$ , the volume of  $\mathcal{M}(s_{i+1})$ , and  $\sup_{x \in \mathcal{M}(s_{i+1})} \|x\|_{\mathbb{R}^D}$ .

1115      *Proof of Lemma B.3.* Consider the oracle transcoder  $\mathcal{T}_* = \mathcal{E}_*^{i+1} \circ \mathcal{D}_*^i$ , which satisfies the Lipschitz condition:

1117       $\text{Lip}_{\mathcal{Z}(s_{i+1})}(\mathcal{T}_*) = \text{Lip}_{\mathcal{Z}(s_{i+1})}(\mathcal{E}_*^{i+1} \circ \mathcal{D}_*^i) \leq \text{Lip}_{\mathcal{Z}(s_{i+1})}(\mathcal{E}_*^{i+1}) \text{Lip}_{\mathcal{M}(s_{i+1})}(\mathcal{D}_*^i) < \infty.$

1118      The proof then follows by applying Proposition A.5 to  $\mathcal{T}_*$ . The width of the  
1119      neural network size is exponential in  $d$  instead of  $d+1$ , because the function  $\mathcal{T}_*$  leaves  
1120      the  $(d+1)$ th component unchanged (as that component is the time), so it can be  
1121      exactly represented by a neural network with no need for approximation.  $\square$

1122      Now we prove Lemma 3.7. We refer to Table 1 for some important notations to  
1123      be used in the proof.

1124     *Proof of Lemma 3.7.* For any  $i \in [S]$ , denote  $T_i = |\mathbb{T} \cap (s_i, s_{i+1}]|$ . Suppose  $t_k$   
1125 is the  $j$ th element of  $\mathbb{T} \cap (s_i, s_{i+1}]$ , i.e. it is in window/segment number  $i$ . We will  
1126 construct encoder, decoder, propagator, and transcoder networks such that for any  
1127  $t_k$ ,

$$(B.2) \quad \left\| \mathcal{D}_{NN}^i \circ \bigcirc_{\ell=1}^j \mathcal{P}_{NN}^i \circ \bigcirc_{w=1}^{i-1} \left( \mathcal{T}_{NN}^w \circ \bigcirc_{\ell=1}^{T_w} \mathcal{P}_{NN}^w \right) \circ \mathcal{E}_{NN}^1 - \mathcal{F}(\cdot, t_k) \right\|_{L^\infty(\mathcal{M}(0))} < \epsilon.$$

1129     We denote  $t_\ell^w$  as the  $\ell$ th element of  $\mathbb{T} \cap [s_w, s_{w+1}]$ . Note that the evolution operator  
1130 to  $t_k$  is

$$\begin{aligned} & \mathcal{F}(\cdot, t_k) \\ &= \bigcirc_{\ell=1}^j \mathcal{F}(\cdot, t_{\ell+1}^i - t_\ell^i) \circ \bigcirc_{w=1}^{i-1} \bigcirc_{\ell=1}^{T_w} \mathcal{F}(\cdot, t_{\ell+1}^w - t_\ell^w) \\ &= \mathcal{D}_*^i \circ \bigcirc_{\ell=1}^j (\mathcal{E}_*^i \circ \mathcal{F}(\cdot, t_{\ell+1}^i - t_\ell^i) \circ \mathcal{D}_*^i) \circ \\ & \quad \bigcirc_{w=1}^{i-1} \left( \mathcal{E}_*^{w+1} \circ \mathcal{D}_*^w \circ \bigcirc_{\ell=1}^{T_w} (\mathcal{E}_*^w \circ \mathcal{F}(\cdot, t_{\ell+1}^w - t_\ell^w) \circ \mathcal{D}_*^w) \right) \circ \mathcal{E}_*^1. \end{aligned}$$

1135     We will approximate each term of the composition in this expression by neural  
1136 networks, and then apply Proposition A.6.

1137     For all  $w \in [S]$ , consider the encoders and decoders  $\mathcal{E}_*^w$  and  $\mathcal{D}_*^w$ . By Lemma B.1,  
1138 there is a neural network  $\mathcal{E}_{NN}^w$  with  $O(\log^2(\frac{S}{\epsilon}))$  layers and width  $O(D(\frac{S}{\epsilon})^{d+1})$  and  
1139 a neural network  $\mathcal{D}_{NN}^w$  with  $O(\log(\frac{S}{\epsilon}))$  layers and width  $O(D(\frac{S}{\epsilon})^{d+1})$  such that

$$(B.3) \quad \begin{aligned} \|\mathcal{E}_{NN}^w - \mathcal{E}_*^w\|_{L^\infty(\mathcal{M}([s_w, s_{w+1}]))} &< \frac{\epsilon}{(2S+1)\text{Lip}(\mathcal{D}_*)\text{Lip}(\mathcal{F})}, \\ \|\mathcal{D}_*^w - \mathcal{D}_{NN}^w\|_{L^\infty(\mathcal{Z}([s_w, s_{w+1}]))} &< \frac{\epsilon}{2S+1}, \end{aligned}$$

1141 where  $\text{Lip}(\mathcal{D}_*) = \max_w \text{Lip}(\mathcal{D}_*^w)$ . Note that the range of  $\mathcal{E}_{NN}^w$  is contained in  $\mathcal{Z}([s_w, s_{w+1}])$ . ■

1142 Next, for all  $w \in [W]$  consider  $\mathcal{P}_*^w(\mathbf{z}, t) = \mathcal{E}_*^w(\mathcal{F}(\mathcal{D}_*^w(\mathbf{z}), t))$  which is called the  
1143 oracle propagator. Since  $\mathcal{P}_*^w$  is Lipschitz as a composition of Lipschitz functions,  
1144 by Lemma B.2, there is a neural network  $\mathcal{P}_{NN}^w$  with  $O(\log(\frac{S}{\epsilon}))$  layers and width  
1145  $O((\frac{S}{\epsilon})^{d+1})$  such that

$$(B.4) \quad \|\bigcirc_{\ell=1}^{T_w} \mathcal{P}_{NN}^w - \bigcirc_{\ell=1}^{T_w} \mathcal{P}_*^w(\cdot, t_{\ell+1}^w - t_\ell^w)\|_{L^\infty(\mathcal{Z}(s_w))} < \frac{\epsilon}{(2S+1)\text{Lip}(\mathcal{D})\text{Lip}(\mathcal{F})}.$$

1147 The composition above makes sense, since the range of  $\mathcal{P}_{NN}^w$  is contained in  $\mathcal{Z}([s_w, s_{w+1}])$ . ■  
1148 by construction.

1149 Finally, for all  $w \in [S-1]$ , consider  $\mathcal{E}_*^{w+1} \circ \mathcal{D}_*^w$ . By Lemma B.3, there is a neural  
1150 network  $\mathcal{T}_{NN}^w$  with  $O(\log(\frac{S}{\epsilon}))$  layers and width  $O((\frac{S}{\epsilon})^d)$  such that

$$(B.5) \quad \|\mathcal{T}_{NN}^w - \mathcal{E}_*^{w+1} \circ \mathcal{D}_*^w\|_{L^\infty(\mathcal{Z}(s_{w+1}))} < \frac{\epsilon}{(2S+1)\text{Lip}(\mathcal{D})\text{Lip}(\mathcal{F})},$$

1152 with the range of  $\mathcal{T}_{NN}^w$  contained in  $\mathcal{Z}(s_{w+1}) = [0, 1]^d \times \{s_{w+1}\}$ . Theorem 3.9 then  
1153 follows by applying Proposition A.6 to the expression (where we indicate matching  
1154 terms with the same letter and subscript):

$$1155 \quad \underbrace{\mathcal{D}_{NN}^i}_{a} \circ \underbrace{\bigcirc_{\ell=1}^j \mathcal{P}_{NN}^i}_{b} \circ \bigcirc_{w=1}^{i-1} \underbrace{(\mathcal{T}_{NN}^w \circ \bigcirc_{\ell=1}^{T_w} \mathcal{P}_{NN}^w)}_{c_w} \circ \underbrace{\mathcal{E}_{NN}^1}_{d_w} - \underbrace{\mathcal{E}_{NN}^1}_{e}$$

$$1156 \quad 1157 \quad \underbrace{\mathcal{D}_*^i \circ \bigcirc_{\ell=1}^j \mathcal{P}_*^w(\cdot, t_{\ell+1}^i - t_\ell)}_a \circ \bigcirc_{w=1}^{i-1} \underbrace{(\mathcal{E}_*^{w+1} \circ \mathcal{D}_*^w)}_{c_w} \circ \bigcirc_{\ell=1}^{T_w} \mathcal{P}_*^w(\cdot, t_{\ell+1}^w - t_{\ell+1}^w) \circ \underbrace{\mathcal{E}_*^1}_{d_w} \parallel, \quad e$$

1158 where the norm is taken over the space  $L^\infty(\mathcal{M}(0))$ . Note that we can do this because  
 1159 the range of each neural network in the composition is contained in the domain of  
 1160 the subsequent neural network (by construction). To compute the bound given in  
 1161 Proposition A.6, we compute Lipschitz constants of compositions of functions. Both  
 1162 terms in the difference involve  $2+2(i-1)+1 = 2i+1$  functions. Each term in Equation  
 1163 (A.2) is the product of the approximation error of a network (encoder, propagator,  
 1164 transcoder, or decoder), and the Lipschitz constant of the oracle maps that come later  
 1165 in the approximation. We will show that each term in the sum is less than  $\frac{\epsilon}{2i+1}$ , so  
 1166 that the total error is less than  $\epsilon$ . Recall  $t_{\ell+1}^i = t_k$  and  $a_1 = t_1$ .

- 1167 (a) **Decoder.** The decoder  $\mathcal{D}_*^i$  corresponds to the  $(2i+1)$ th and last term of the  
 1168 sum in Equation (A.2). The Lipschitz factor is over an empty composition,  
 1169 which is defined to be 1 by convention. Thus the decoder term contributes  
 1170 an error of  $\frac{\epsilon}{2S+1}$  by Equation (B.3), which is less than  $\frac{\epsilon}{2i+1}$  since  $i \leq S$ .
- 1171 (b) **Final Propagator.** The final propagator is  $\mathcal{P}_*^i$ , and it corresponds to the  
 1172 second to last term of the sum in Equation (A.2). The Lipschitz factor is the  
 1173 Lipschitz constant of the decoder  $\mathcal{D}_*^i$ . Thus this term contributes an error of  
 1174 (by the approximation in (B.4))

$$1175 \quad \text{Lip}_{\mathcal{Z}^i(t_k)}(\mathcal{D}_*^i) \cdot \frac{\epsilon}{(2S+1)\text{Lip}(\mathcal{D})\text{Lip}(\mathcal{F})} \leq \frac{\epsilon}{(2S+1)\text{Lip}(\mathcal{F})} \leq \frac{\epsilon}{2i+1}.$$

- 1176 (c) **Transcoder.** For any window  $m \in [i-1]$ , the  $m$ th transcoder corresponds  
 1177 to a middle term in the sum in Equation (A.2). For each of these terms,  
 1178 the error is again given by the approximation error of the corresponding  
 1179 transcoder network, scaled by the Lipschitz constant of the neural networks  
 1180 that come after. The Lipschitz constant of the composition of all functions  
 1181 after the  $m$ th transcoder are

$$\begin{aligned} 1182 \quad & \text{Lip}_{\mathcal{Z}^m(s_{m+1})} \left( \mathcal{D}_*^i \circ \bigcirc_{\ell=1}^j (\mathcal{P}_*^i(\cdot, t_{\ell+1}^i - t_\ell^i)) \circ \bigcirc_{w=m+1}^{i-1} \left( \mathcal{E}_*^{w+1} \circ \mathcal{D}_*^w \circ \bigcirc_{\ell=1}^{T_w} (\mathcal{P}_*^w(\cdot, t_{\ell+1}^w - t_\ell^w)) \right) \right) \\ 1183 \quad & = \text{Lip}_{\mathcal{Z}^m(s_{m+1})} \left( \mathcal{D}_*^i \circ \bigcirc_{\ell=1}^j (\mathcal{E}_*^i \circ \mathcal{F}(\cdot, t_{\ell+1}^i - t_\ell^i) \circ \mathcal{D}_*^i) \circ \bigcirc_{w=m+1}^{i-1} \left( \mathcal{E}_*^{w+1} \circ \mathcal{D}_*^w \circ \bigcirc_{\ell=1}^{T_w} (\mathcal{E}_*^w \circ \mathcal{F}(\cdot, t_{\ell+1}^w - t_\ell^w) \circ \mathcal{D}_*^w) \right) \right) \\ 1184 \quad & = \text{Lip}_{\mathcal{Z}^m(s_{m+1})} \left( \mathcal{D}_*^i \circ \mathcal{E}_*^i \circ \mathcal{F}(\cdot, t_{j+1}^i - s_i) \circ \mathcal{D}_*^i \circ \bigcirc_{w=m+1}^{i-1} \left( \mathcal{E}_*^{w+1} \circ \mathcal{D}_*^w \circ \mathcal{E}_*^w \circ \mathcal{F}(\cdot, s_{w+1} - s_w) \circ \mathcal{D}_*^w \right) \right) \\ 1185 \quad & = \text{Lip}_{\mathcal{Z}^m(s_{m+1})} \left( \mathcal{F}(\cdot, t_k - s_i) \circ \mathcal{D}_*^i \circ \bigcirc_{w=m+1}^{i-1} \left( \mathcal{E}_*^{w+1} \circ \mathcal{F}(\cdot, s_{w+1} - s_w) \circ \mathcal{D}_*^w \right) \right) \\ 1186 \quad & = \text{Lip}_{\mathcal{Z}^m(s_{m+1})} \left( \mathcal{F}(\cdot, t_k - s_i) \circ \mathcal{D}_*^i \circ \mathcal{E}_*^i \circ \mathcal{F}(\cdot, s_i - s_{m+1}) \circ \mathcal{D}_*^{m+1} \right) \\ 1187 \quad & = \text{Lip}_{\mathcal{Z}^m(s_{m+1})} \left( \mathcal{F}(\cdot, t_k - s_i) \circ \mathcal{F}(\cdot, s_i - s_{m+1}) \circ \mathcal{D}_*^{m+1} \right) \\ 1188 \quad & = \text{Lip}_{\mathcal{Z}^m(s_{m+1})} \left( \mathcal{F}(\cdot, t_k - s_{m+1}) \circ \mathcal{D}_*^{m+1} \right) \leq \text{Lip}_{\mathcal{M}(s_{m+1})}(\mathcal{F}(\cdot, t_k - s_{m+1})) \text{Lip}_{\mathcal{Z}^{m+1}(s_{m+1})}(\mathcal{D}_*^{m+1}). \quad \blacksquare \end{aligned}$$

1189 Thus, by scaling the approximation guarantee in Equation (B.5), we see that  
 1190 each transcoder term contributes error at most

$$\begin{aligned} 1191 \quad & \text{Lip}_{\mathcal{M}(s_{m+1})}(\mathcal{F}(\cdot, t_k - s_{m+1})) \text{Lip}_{\mathcal{Z}^{m+1}(s_{m+1})}(\mathcal{D}_*^{m+1}) \frac{\epsilon}{(2S+1)\text{Lip}(\mathcal{D})\text{Lip}(\mathcal{F})} \\ 1192 \quad & \leq \frac{\epsilon}{2S+1} \leq \frac{\epsilon}{2i+1}. \end{aligned}$$

- 1193 (d) **Intermediate Propagators.** For any window  $m \in [i-1]$ , the  $m$ th propagator  
 1194 corresponds to a middle term in the sum in Equation (A.2), similar to the  
 1195 transcoder case. The Lipschitz constant of the composition of all functions  
 1196 after the  $m$ th oracle propagator are

$$\begin{aligned}
 & \text{Lip}_{\mathcal{Z}^m(s_{m+1})}(\mathcal{D}_*^i \circ \bigcirc_{\ell=1}^j (\mathcal{P}_*^i(\cdot, t_{\ell+1}^i - t_\ell^i) \circ \mathcal{D}_*^i) \circ \bigcirc_{w=m+1}^{i-1} (\mathcal{E}_*^{w+1} \circ \mathcal{D}_*^w \circ \bigcirc_{\ell=1}^{T_w} (\mathcal{P}_*^w(\cdot, t_{\ell+1}^w - t_\ell^w)) \circ \mathcal{E}_*^{m+1} \circ \mathcal{D}_*^m) \\
 & = \text{Lip}_{\mathcal{Z}^m(s_{m+1})}(\mathcal{D}_*^i \circ \bigcirc_{\ell=1}^j (\mathcal{E}_*^i \circ \mathcal{F}(\cdot, t_{\ell+1}^i - t_\ell^i) \circ \mathcal{D}_*^i) \circ \bigcirc_{w=m+1}^{i-1} (\mathcal{E}_*^{w+1} \circ \mathcal{D}_*^w \circ \bigcirc_{\ell=1}^{T_w} (\mathcal{E}_*^w \circ \mathcal{F}(\cdot, t_{\ell+1}^w - t_\ell^w) \circ \mathcal{D}_*^w)) \circ \mathcal{E}_*^{m+1} \circ \mathcal{D}_*^m) \\
 & = \text{Lip}_{\mathcal{Z}^m(s_{m+1})}(\mathcal{D}_*^i \circ \mathcal{E}_*^i \circ \mathcal{F}(\cdot, t_{j+1}^i - s_i) \circ \mathcal{D}_*^i \circ \bigcirc_{w=m+1}^{i-1} (\mathcal{E}_*^{w+1} \circ \mathcal{D}_*^w \circ \mathcal{E}_*^w \circ \mathcal{F}(\cdot, s_{w+1} - s_w) \circ \mathcal{D}_*^w) \circ \mathcal{E}_*^{m+1} \circ \mathcal{D}_*^m) \\
 & = \text{Lip}_{\mathcal{Z}^m(s_{m+1})}(\mathcal{F}(\cdot, t_k - s_i) \circ \mathcal{D}_*^i \circ \bigcirc_{w=m+1}^{i-1} (\mathcal{E}_*^{w+1} \circ \mathcal{F}(\cdot, s_{w+1} - s_w) \circ \mathcal{D}_*^w) \circ \mathcal{E}_*^{m+1} \circ \mathcal{D}_*^m) \\
 & = \text{Lip}_{\mathcal{Z}^m(s_{m+1})}(\mathcal{F}(\cdot, t_k - s_i) \circ \mathcal{D}_*^i \circ \mathcal{E}_*^i \circ \mathcal{F}(\cdot, s_i - s_{m+1}) \circ \mathcal{D}_*^{m+1} \circ \mathcal{E}_*^{m+1} \circ \mathcal{D}_*^m) \\
 & = \text{Lip}_{\mathcal{Z}^m(s_{m+1})}(\mathcal{F}(\cdot, t_k - s_i) \circ \mathcal{F}(\cdot, s_i - s_{m+1}) \circ \mathcal{D}_*^m) \\
 & = \text{Lip}_{\mathcal{Z}^m(s_{m+1})}(\mathcal{F}(\cdot, t_k - s_{m+1}) \circ \mathcal{D}_*^m) \leq \text{Lip}_{\mathcal{M}(s_{m+1})}(\mathcal{F}(\cdot, t_k - s_{m+1})) \text{Lip}_{\mathcal{Z}^m(s_{m+1})}(\mathcal{D}_*^m). \quad \blacksquare
 \end{aligned}$$

1204 Thus, by scaling the approximation guarantee in Equation (B.4), we see that  
 1205 each intermediate propagator term contributes error at most

$$\begin{aligned}
 & \text{Lip}_{\mathcal{M}(s_{m+1})}(\mathcal{F}(\cdot, t_k - s_{m+1})) \text{Lip}_{\mathcal{Z}^m(s_{m+1})}(\mathcal{D}_*^m) \frac{\epsilon}{(2W+1)\text{Lip}(\mathcal{D})\text{Lip}(\mathcal{F})} \\
 & \leq \frac{\epsilon}{2W+1} \leq \frac{\epsilon}{2i+1}.
 \end{aligned}$$

- 1208 (e) **Encoder.** The encoder term  $\mathcal{E}^1$  corresponds to the first term of the sum in  
 1209 Equation (A.2). The Lipschitz factor can be computed as:

$$\begin{aligned}
 & \text{Lip}_{\mathcal{Z}^1(t_1)}(\mathcal{D}_*^i \circ \bigcirc_{\ell=1}^j (\mathcal{P}_*^i(\cdot, t_{\ell+1}^i - t_\ell^i) \circ \mathcal{D}_*^i) \circ \bigcirc_{w=1}^{i-1} (\mathcal{E}_*^{w+1} \circ \mathcal{D}_*^w \circ \bigcirc_{\ell=1}^{T_w} (\mathcal{P}_*^w(\cdot, t_{\ell+1}^w - t_\ell^w))) \\
 & = \text{Lip}_{\mathcal{Z}^1(t_1)}(\mathcal{D}_*^i \circ \bigcirc_{\ell=1}^j (\mathcal{E}_*^i \circ \mathcal{F}(\cdot, t_{\ell+1}^i - t_\ell^i) \circ \mathcal{D}_*^i) \circ \bigcirc_{w=1}^{i-1} (\mathcal{E}_*^{w+1} \circ \mathcal{D}_*^w \circ \bigcirc_{\ell=1}^{T_w} (\mathcal{E}_*^w \circ \mathcal{F}(\cdot, t_{\ell+1}^w - t_\ell^w) \circ \mathcal{D}_*^w)) \\
 & = \text{Lip}_{\mathcal{Z}^1(t_1)}(\mathcal{D}_*^i \circ \mathcal{E}_*^i \circ \mathcal{F}(\cdot, t_{j+1}^i - s_i) \circ \mathcal{D}_*^i \circ \bigcirc_{w=1}^{i-1} (\mathcal{E}_*^{w+1} \circ \mathcal{D}_*^w \circ \mathcal{E}_*^w \circ \mathcal{F}(\cdot, s_{w+1} - s_w) \circ \mathcal{D}_*^w)) \\
 & = \text{Lip}_{\mathcal{Z}^1(t_1)}(\mathcal{F}(\cdot, t_k - s_i) \circ \mathcal{D}_*^i \circ \bigcirc_{w=1}^{i-1} (\mathcal{E}_*^{w+1} \circ \mathcal{F}(\cdot, s_{w+1} - s_w) \circ \mathcal{D}_*^w)) \\
 & = \text{Lip}_{\mathcal{Z}^1(t_1)}(\mathcal{F}(\cdot, t_k - s_i) \circ \mathcal{D}_*^i \circ \mathcal{E}_*^i \circ \mathcal{F}(\cdot, s_i - s_1) \circ \mathcal{D}_*^1) \\
 & = \text{Lip}_{\mathcal{Z}^1(t_1)}(\mathcal{F}(\cdot, t_k - t_1) \circ \mathcal{F}(\cdot, s_i - s_1) \circ \mathcal{D}_*^1) \\
 & = \text{Lip}_{\mathcal{Z}^1(t_1)}(\mathcal{F}(\cdot, t_k - t_1) \circ \mathcal{D}_*^1) \leq \text{Lip}_{\mathcal{M}(t_1)}(\mathcal{F}(\cdot, t_k - t_1)) \text{Lip}_{\mathcal{Z}^1(t_1)}(\mathcal{D}_*^1). \quad \blacksquare
 \end{aligned}$$

1217 Thus, by scaling the approximation guarantee in Equation (B.3), we see that  
 1218 the encoder term contributes error at most

$$\begin{aligned}
 & \text{Lip}_{\mathcal{M}(t_1)}(\mathcal{F}(\cdot, t_k - t_1)) \text{Lip}_{\mathcal{Z}^1(t_1)}(\mathcal{D}_*^1) \frac{\epsilon}{(2S+1)\text{Lip}(\mathcal{D})\text{Lip}(\mathcal{F})} \\
 & \leq \frac{\epsilon}{2S+1} \leq \frac{\epsilon}{2i+1}. \quad \square
 \end{aligned}$$

- 1221 **B.2. Proof of Theorem 3.9.** In this section, we prove Theorem 3.9, which  
 1222 follows quickly from Lemma 3.7. The key is to construct a WeldNet model with  
 1223 more windows that implements exactly the same function as the WeldNet constructed  
 1224 before.

1225 *Proof of Theorem 3.9.* Let  $\pi : [W] \rightarrow [S]$  be the function that indicates (with the  
 1226 index) which segment each window falls inside. Consider the S-window WeldNet from

1227 Lemma 3.7 denoted  $\overline{\mathcal{W}_{\text{NN}}}$  with components (for all  $s \in [S]$ )  $\overline{\mathcal{E}_{\text{NN}}^s}$ ,  $\overline{\mathcal{D}_{\text{NN}}^s}$ ,  $\overline{\mathcal{P}_{\text{NN}}^s}$ , and (for  
 1228  $s < S$ )  $\overline{\mathcal{T}_{\text{NN}}^s}$  such that for any  $k \in [T]$  with  $k$  being the  $j$ th element of  $T \cap [s_i, s_{i+1}]$ ),

$$1229 \quad \sup_{\mathbf{x}(0) \in \mathcal{M}(0)} \|\overline{\mathcal{W}_{\text{NN}}}(\mathbf{x}(0), t_k) - \mathcal{F}(\mathbf{x}(0), t_k)\|_{\mathbb{R}^D} < \epsilon.$$

1230 We now construct a  $W$ -window WeldNet that implements the same function as  
 1231 the  $S$ -window  $\overline{\mathcal{W}_{\text{NN}}}$ . For each  $i \in [W]$ , consider the neural networks

$$1232 \quad \mathcal{E}_{\text{NN}}^i = \overline{\mathcal{E}_{\text{NN}}^{\pi(i)}}, \mathcal{D}_{\text{NN}}^i = \overline{\mathcal{D}_{\text{NN}}^{\pi(i)}}, \mathcal{P}_{\text{NN}}^i = \overline{\mathcal{P}_{\text{NN}}^{\pi(i)}},$$

1233 and for each  $i \in [W]$  such that  $w_{i+1} = \underline{s_{j+1}}$  for some  $j \in [S]$  (i.e windows that end at  
 1234 the end of a segment), we define  $\mathcal{T}_{\text{NN}}^i = \mathcal{T}_{\text{NN}}^j$ , and we construct every other transcoder  
 1235 to be the identity function in  $\mathbb{R}^{d+1}$  given by Proposition A.1.

1236 Since for all  $\mathbf{x}(0) \in \mathcal{M}(0)$ , we have  $\|\overline{\mathcal{W}_{\text{NN}}}(\mathbf{x}(0), t_k) - \mathcal{F}(\mathbf{x}(0), t_k)\|_{\mathbb{R}^D} < \epsilon$ , we  
 1237 can complete the proof by showing that  $\mathcal{W}_{\text{NN}}(\mathbf{x}(0), t_k) = \overline{\mathcal{W}_{\text{NN}}}(\mathbf{x}(0), t_k)$ . We denote  
 1238  $T_i = |\mathbb{T} \cap (w_i, w_{i+1})|$  and  $\overline{T}_i = |\mathbb{T} \cap (s_i, s_{i+1})|$ . Suppose that  $t_k$  is the  $j$ th element of  
 1239  $\mathbb{T} \cap (w_i, w_{i+1})$  and also that  $t_k$  is the  $j'$ -th element of  $\mathbb{T} \cap (s_{i'}, s_{i'+1})$  (which means  
 1240  $\pi(i) = i'$ ). Then we have

$$1241 \quad \mathcal{W}_{\text{NN}}(\mathbf{x}(0), t_k) = \mathcal{D}_{\text{NN}}^i \circ \bigcirc_{\ell=1}^j \mathcal{P}_{\text{NN}}^i \circ \bigcirc_{w=1}^{i-1} \left( \mathcal{T}_{\text{NN}}^w \circ \bigcirc_{\ell=1}^{T_w} \mathcal{P}_{\text{NN}}^w \right) \circ \mathcal{E}_{\text{NN}}^1(\mathbf{x}(0)) \text{ and}$$

$$1242 \quad \overline{\mathcal{W}_{\text{NN}}}(\mathbf{x}(0), t_k) = \overline{\mathcal{D}_{\text{NN}}^{i'}} \circ \bigcirc_{\ell=1}^j \overline{\mathcal{P}_{\text{NN}}^{i'}} \circ \bigcirc_{s=1}^{i'-1} \left( \overline{\mathcal{T}_{\text{NN}}^s} \circ \bigcirc_{\ell=1}^{\overline{T}_s} \overline{\mathcal{P}_{\text{NN}}^s} \right) \circ \overline{\mathcal{E}_{\text{NN}}^1}(\mathbf{x}(0)).$$

1243 Note that  $\mathcal{E}_{\text{NN}}^1 = \overline{\mathcal{E}_{\text{NN}}^{\pi(1)}}$  since the first window is inside of the first segment. Next,  
 1244 note that  $\mathcal{D}_{\text{NN}}^i = \overline{\mathcal{D}_{\text{NN}}^{\pi(i)}} = \overline{\mathcal{D}_{\text{NN}}^{i'}}$  and  $\mathcal{P}_{\text{NN}}^i = \overline{\mathcal{P}_{\text{NN}}^{\pi(i)}} = \overline{\mathcal{P}_{\text{NN}}^{i'}}$ . Finally, we will decompose  
 1245 the composition of propagators and encoders to be over each segment before being  
 1246 over each window.

1247 Let  $n_s^m$  denote the index of the  $m$ th window in the  $s$ th segment, and let  $\mathbf{n}_s$   
 1248 denote the index of the final window in the  $s$ th segment. Note that all encoders  
 1249 except for the last window in each segment implements the identity, i.e.  $\mathcal{T}_{\text{NN}}^{n_s^m} = I_{d+1}$   
 1250 if  $m < |\pi^{-1}(\{s\})|$ , where  $\pi^{-1}(\{s\})$  is the pre-image of  $s$  by  $\pi$ , which is the set of  
 1251 all window indices inside of Segment  $s$ , and the last encoder within segment  $s$  of  
 1252  $\mathcal{W}_{\text{NN}}$  is equal to the  $s$ th encoder of  $\overline{\mathcal{W}_{\text{NN}}}$ , i.e.  $\mathcal{T}_{\text{NN}}^{\mathbf{n}_s} = \overline{\mathcal{T}_{\text{NN}}^s}$ . Also by construction  
 1253  $\mathcal{P}_{\text{NN}}^{n_s^m} = \overline{\mathcal{P}_{\text{NN}}^s}$ . Then we decompose the non encoder/decoder terms of the above  
 1254 composition

$$\begin{aligned}
& \bigcirc_{\ell=1}^j \mathcal{P}_{\text{NN}}^i \circ \bigcirc_{w=1}^{i-1} \left( \mathcal{T}_{\text{NN}}^w \circ \bigcirc_{\ell=1}^{\mathbf{T}_w} \mathcal{P}_{\text{NN}}^w \right) \\
&= \bigcirc_{\ell=1}^j \overline{\mathcal{P}_{\text{NN}}^i} \circ \bigcirc_{w=1}^{i-1} \left( \mathcal{T}_{\text{NN}}^w \circ \bigcirc_{\ell=1}^{\mathbf{T}_w} \mathcal{P}_{\text{NN}}^w \right) \\
&= \bigcirc_{\ell=1}^j \overline{\mathcal{P}_{\text{NN}}^{i'}} \circ \bigcirc_{s=1}^{i'-1} \left[ \bigcirc_{m=1}^{|\pi^{-1}(s)|} \left( \mathcal{T}_{\text{NN}}^{n_s^m} \circ \bigcirc_{\ell=1}^{\mathbf{T}_{n_s^m}} \mathcal{P}_{\text{NN}}^{n_s^m} \right) \right] \\
&= \bigcirc_{\ell=1}^j \overline{\mathcal{P}_{\text{NN}}^{i'}} \circ \bigcirc_{s=1}^{i'-1} \left[ \bigcirc_{m=1}^{|\pi^{-1}(s)|} \left( \mathcal{T}_{\text{NN}}^{n_s^m} \circ \bigcirc_{\ell=1}^{\mathbf{T}_{n_s^m}} \overline{\mathcal{P}_{\text{NN}}^s} \right) \right] \\
&= \bigcirc_{\ell=1}^j \overline{\mathcal{P}_{\text{NN}}^{i'}} \circ \bigcirc_{s=1}^{i'-1} \left[ \left( \mathcal{T}_{\text{NN}}^{n_s} \circ \bigcirc_{\ell=1}^{\mathbf{T}_{n_s}} \overline{\mathcal{P}_{\text{NN}}^s} \right) \circ \bigcirc_{m=1}^{|\pi^{-1}(s)|-1} \left( \mathcal{T}_{\text{NN}}^{n_s^m} \circ \bigcirc_{\ell=1}^{\mathbf{T}_{n_s^m}} \overline{\mathcal{P}_{\text{NN}}^s} \right) \right] \\
&= \bigcirc_{\ell=1}^j \overline{\mathcal{P}_{\text{NN}}^{i'}} \circ \bigcirc_{s=1}^{i'-1} \left[ \overline{\mathcal{T}_{\text{NN}}^s} \circ \bigcirc_{\ell=1}^{\mathbf{T}_{n_s}} \overline{\mathcal{P}_{\text{NN}}^s} \circ \bigcirc_{m=1}^{|\pi^{-1}(s)|-1} \left( I_{d+1} \circ \bigcirc_{\ell=1}^{\mathbf{T}_{n_s^m}} \overline{\mathcal{P}_{\text{NN}}^s} \right) \right] \\
&= \bigcirc_{\ell=1}^j \overline{\mathcal{P}_{\text{NN}}^{i'}} \circ \bigcirc_{s=1}^{i'-1} \left[ \overline{\mathcal{T}_{\text{NN}}^s} \circ \bigcirc_{\ell=1}^{\mathbf{T}_{n_s}} \overline{\mathcal{P}_{\text{NN}}^s} \circ \bigcirc_{m=1}^{|\pi^{-1}(s)|-1} \bigcirc_{\ell=1}^{\mathbf{T}_{n_s^m}} \overline{\mathcal{P}_{\text{NN}}^s} \right] \\
&= \bigcirc_{\ell=1}^j \overline{\mathcal{P}_{\text{NN}}^{i'}} \circ \bigcirc_{s=1}^{i'-1} \left[ \overline{\mathcal{T}_{\text{NN}}^s} \circ \bigcirc_{\ell=1}^{\mathbf{T}_s} \overline{\mathcal{P}_{\text{NN}}^s} \right].
\end{aligned}$$

1263 Thus,

$$\begin{aligned}
1264 \quad \mathcal{W}_{\text{NN}}(\mathbf{x}(0), t_k) &= \mathcal{D}_{\text{NN}}^i \circ \bigcirc_{\ell=1}^j \mathcal{P}_{\text{NN}}^i \circ \bigcirc_{w=1}^{i-1} \left( \mathcal{T}_{\text{NN}}^w \circ \bigcirc_{\ell=1}^{\mathbf{T}_w} \mathcal{P}_{\text{NN}}^w \right) \circ \mathcal{E}_{\text{NN}}^1 \\
1265 &= \mathcal{D}_{\text{NN}}^i \circ \bigcirc_{\ell=1}^j \overline{\mathcal{P}_{\text{NN}}^{i'}} \circ \bigcirc_{s=1}^{i'-1} \left[ \overline{\mathcal{T}_{\text{NN}}^s} \circ \bigcirc_{\ell=1}^{\mathbf{T}_s} \overline{\mathcal{P}_{\text{NN}}^s} \right] \circ \mathcal{E}_{\text{NN}}^1 \\
1266 &= \overline{\mathcal{D}_{\text{NN}}^{i'}} \circ \bigcirc_{\ell=1}^j \overline{\mathcal{P}_{\text{NN}}^{i'}} \circ \bigcirc_{s=1}^{i'-1} \left[ \overline{\mathcal{T}_{\text{NN}}^s} \circ \bigcirc_{\ell=1}^{\mathbf{T}_s} \overline{\mathcal{P}_{\text{NN}}^s} \right] \circ \overline{\mathcal{E}_{\text{NN}}^1} \\
1267 &= \overline{\mathcal{W}_{\text{NN}}}(\mathbf{x}(0), t_k).
\end{aligned}$$

1268 Now note that the size of all autoencoders, all propagators, and  $S - 1$  transcoders  
1269 of  $\mathcal{W}_{\text{NN}}$  are exactly equal to the size of (corresponding) components in  $\overline{\mathcal{W}_{\text{NN}}}$ . For  
1270 the  $W - S$  identity transcoders of  $\mathcal{W}_{\text{NN}}$ , we can use any network size as a result of  
1271 Proposition A.1. This completes the proof.  $\square$

1272 **B.3. Proof of Theorem 3.10.** In this section, we prove Theorem 3.10. First,  
1273 we prove a version of Lemma B.2 but without the dynamics assumption.

1274 **LEMMA B.4** (Propagator Approximation in General Case). *Let  $\mathcal{P}_* : \mathcal{Z}([0, T]) \times$   
1275  $[0, T] \rightarrow \mathcal{Z}([0, T])$  be a Lipschitz function such that*

- *For all  $x \in \mathcal{Z}([0, T])$ ,  $\mathcal{P}_*(x, 0) = x$ ,*
- *For all  $x \in \mathcal{Z}([0, T])$ ,  $t \in [0, T]$ ,  $s \in [0, T-t]$ ,  $\mathcal{P}_*(\mathcal{P}_*(x, t), s) = \mathcal{P}_*(x, t+s)$ .*

1276 Let  $0 = t_1 < t_2 < \dots < t_T = T$  be a grid for  $[0, T]$  of  $T$  points, and  $\epsilon > 0$ . Then there  
1277 is a function  $\mathcal{P}_{\text{NN}} \in \mathcal{F}_{\text{NN}}(d+1, d+1, L, W)$  with parameters

$$1280 \quad L = O\left(T \log\left(\frac{T}{\epsilon}\right)\right), \quad W = O\left(\left(\frac{T}{\epsilon}\right)^d\right),$$

1281 such that for all  $k \in [\mathsf{T} - 1]$  and  $i \in [\mathsf{T} - k]$ , we have

$$1282 \quad \sup_{\mathbf{z}(t_k) \in \mathcal{Z}(t_k)} \| (\bigcirc_{j=1}^i \mathcal{P}_{\text{NN}}) (\mathbf{z}(t_k)) - \mathcal{P}_*(\mathbf{z}(t_k), t_{k+i} - t_k) \|_{\mathbb{R}^{d+1}} \leq \varepsilon.$$

1283 Moreover, the range of  $\mathcal{P}_{\text{NN}}$  is contained in  $\mathcal{Z}([0, 1])$ . The constants hidden in  $O$   
1284 depend on  $d$  and  $\text{Lip}_*(\mathcal{P}) = \sup_{j \in [\mathsf{T}-1], i \leq \mathsf{T}-j} \text{Lip}_{\mathcal{Z}(t_j)}(\mathcal{P}_*(\cdot, t_{j+i} - t_j))$ .

1285     Proof of Lemma B.4. For all  $j \in [\mathsf{T} - 1]$ , define  $p_*^j : [0, 1]^d \rightarrow [0, 1]^d$  by  $p_*^j(\mathbf{v}) =$   
1286      $[\mathcal{P}_*((\mathbf{v}, t_j), t_{j+1} - t_j)]_{1, \dots, d}$  for all  $\mathbf{v} \in [0, 1]^d$ . In other words,  $p_*^j$  returns the first  
1287      $d$ -components (i.e. without time) of the code propagated from time  $t_j$  to  $t_{j+1}$ , and  
1288     note that this is Lipschitz.

1289     By Proposition A.5, there is a neural network  $p_{\text{NN}}^j \in \mathcal{F}_{\text{NN}}(d, d, L_j, W_j)$  such that  
1290      $\|p_{\text{NN}}^j - p_*^j\|_{L^\infty([0, 1]^d)} < \frac{\epsilon}{(\mathsf{T}-1)\text{Lip}(\mathcal{P}_*)}$ , with  $O(\log(\frac{\mathsf{T}}{\epsilon}))$  layers and width  $O((\frac{\mathsf{T}}{\epsilon})^d)$ .  
1291     Moreover, the range of  $p_{\text{NN}}^j$  can be restricted to  $[0, 1]^d$  by Proposition A.2(Part 2).  
1292     For all  $k \in [\mathsf{T} - 1]$  and  $i \in [\mathsf{T} - k]$ , note that  $\bigcirc_{j=k}^{k+i-1} p_*^j = \bigcirc_{j=k}^{k+i-1} (\mathcal{P}_*(\cdot, t_{j+1} - t_j)) = \mathcal{P}_*(\cdot, t_{k+i} - t_k)$ . By Proposition A.6, we have

$$\begin{aligned} 1294 \quad & \sup_{\mathbf{v} \in [0, 1]^d} \left\| \left( \bigcirc_{j=k}^{k+i-1} p_{\text{NN}}^j \right) (\mathbf{v}) - \mathcal{P}_*((\mathbf{v}, t_k), t_{k+i} - t_k) \right\| \\ 1295 \quad & \leq \sum_{j=k}^{k+i-1} \text{Lip}_{\mathcal{Z}(t_{j+1})} \left( \bigcirc_{\ell=j+1}^{k+i-1} \mathcal{P}_*(\cdot, t_{\ell+1} - t_\ell) \right) \|p_{\text{NN}}^j - p_*^j\|_{L^\infty([0, 1]^d)} \\ 1296 \quad & = \sum_{j=k}^{k+i-1} \text{Lip}_{\mathcal{Z}(t_{j+1})} (\mathcal{P}_*(\cdot, t_{k+i} - t_{j+1})) \|p_{\text{NN}}^j - p_*^j\|_{L^\infty([0, 1]^d)} \\ 1297 \quad & \leq \sum_{j=k}^{k+i-1} \frac{\text{Lip}_{\mathcal{Z}(t_{j+1})} (\mathcal{P}_*(\cdot, t_{k+i} - t_{j+1}))}{(\mathsf{T}-1)\text{Lip}(\mathcal{P}_*)} \epsilon \leq \sum_{j=k}^{k+i-1} \frac{\epsilon}{\mathsf{T}} = \frac{i}{\mathsf{T}} \epsilon < \epsilon, \end{aligned}$$

1298     We finish the proof by constructing a single neural network, denoted  $\mathcal{P}_{\text{NN}}$ , which  
1299     will exactly represent each of the  $p_{\text{NN}}^j$ , selecting the correct network based on a time  
1300     parameter. For all  $j \in [\mathsf{T} - 1]$ , let  $\delta_j = \min\{|t_{j+1} - t_j|, |t_j - t_{j-1}|\}$  (and  $\delta_1 = |t_2 - t_1|$ )  
1301     represent the distance between the time point  $t_j$  and the closest other time point.

1302     We will use the notation  $\mathbf{v}$  to represent the first  $d$ -dimensions of latent code in  
1303      $\mathcal{Z}([a, b])$ , and we define the function  $\tilde{p}_{\text{NN}}^j : [0, 1]^d \times [0, T] \rightarrow [0, 1]^d$  (where we recall  
1304     that  $\sigma$  is the ReLU function)

$$1305 \quad \tilde{p}_{\text{NN}}^j((\mathbf{v}, t)) := \sigma \left( \left( p_{\text{NN}}^j((\mathbf{v}, t)), t_j \right) - \frac{\sqrt{d}}{\delta_j} (t - t_j) \mathbf{1}_{d+1} \right).$$

1306     Here,  $\mathbf{1}_{d+1} \in \mathbb{R}^{d+1}$  is the vector of all 1s. Note that  $\tilde{p}_{\text{NN}}^j \in \mathcal{F}_{\text{NN}}(d+1, d+1, 1 +$   
1307      $L_j, W_j)$ , and for all  $\mathbf{v} \in [0, 1]^d$  and  $k \in [\mathsf{T} - 1]$ .

$$1308 \quad \tilde{p}_{\text{NN}}^j((\mathbf{v}, t_k)) = \begin{cases} (p_{\text{NN}}^k((\mathbf{v}, t_k)), t_{k+1}), & j = k, \\ \mathbf{0}_{d+1}, & j \neq k. \end{cases}$$

1309     Finally, we define the propagator network for all  $\mathbf{z} \in \mathbb{R}^{d+1}$  as

$$1310 \quad \mathcal{P}_{\text{NN}}(\mathbf{z}) = \sum_{j=1}^{\mathsf{T}} \tilde{p}_{\text{NN}}^j(\mathbf{z}).$$

1311     Then for all  $k \in [\mathsf{T} - 1]$ , we have that  $\mathcal{P}_{\text{NN}}((\mathbf{v}, t_k)) = p_{\text{NN}}^k(\mathbf{z})$ . In addition, we can

1312 see that for all  $j \in [\mathsf{T} - 1]$ , we have

$$\begin{aligned} 1313 \quad (\bigcirc_{j=1}^k \mathcal{P}_{\text{NN}}) (\mathbf{v}, t_1) &= (\bigcirc_{j=1}^{k-1} \mathcal{P}_{\text{NN}}) (\tilde{p}_{\text{NN}}^1(\mathbf{v}, t_1)) \\ 1314 \quad &= (\bigcirc_{j=1}^{k-1} \mathcal{P}_{\text{NN}}) (p_{\text{NN}}^1(\mathbf{v}), t_2) \\ 1315 \quad &= (\bigcirc_{j=1}^{k-2} \mathcal{P}_{\text{NN}}) (\tilde{p}_{\text{NN}}^2(p_{\text{NN}}^1(\mathbf{v}), t_2)) \\ 1316 \quad &= (\bigcirc_{j=1}^{k-2} \mathcal{P}_{\text{NN}}) (p_{\text{NN}}^2(p_{\text{NN}}^1(\mathbf{v})), t_3) \\ 1317 \quad &= \dots \\ 1318 \quad &= p_{\text{NN}}^k(p_{\text{NN}}^{k-1}(\dots(p_{\text{NN}}^1(\mathbf{v}))), t_{k+1}). \end{aligned}$$

1319 We finish by computing the size of  $\mathcal{P}_{\text{NN}}$ . According to Proposition A.2 Part  
 1320 3, the number of layers of  $\mathcal{P}_{\text{NN}}$  is the sum of the number of layers of each  $\tilde{p}_{\text{NN}}$ ,  
 1321 so it is  $\sum_{t=j}^{\mathsf{T}} (1 + O(\log(\frac{1}{\epsilon}))) = O(\mathsf{T} \log(\frac{1}{\epsilon}))$ , and the width is given by  $d + 1 +$   
 1322  $2 \max_j O\left((\frac{1}{\epsilon})^d\right) = O\left((\frac{1}{\epsilon})^d\right)$ .  $\square$

1323 Now we prove Theorem 3.10. In the proof, we first derive a result analogous to  
 1324 Lemma 3.7, and then apply the argument in Theorem 3.9.

1325 *Proof of Theorem 3.10.* We first assume that  $\mathsf{W} = \mathsf{S}$ , i.e. the windows are equal  
 1326 to the segments. This is done by following proof of Lemma 3.7, but using Lemma  
 1327 B.4 instead of Lemma B.2 to construct  $\mathcal{P}_{\text{NN}}$ . We use the exact same construction as  
 1328 Lemma 3.7 for the autoencoders and transcoders.

1329 For all  $w \in [\mathsf{W}]$ , consider the oracle propagator  $\mathcal{P}_*^w(x, t) = \mathcal{E}_*^w(\mathcal{F}(\mathcal{D}_*^w(x), t))$ .  
 1330 By Lemma B.4, there is a neural network  $\mathcal{P}_{\text{NN}}^w$  with  $O(\mathbf{T}_w \log(\frac{\mathsf{WT}_w}{\epsilon}))$  layers and  
 1331 width  $O\left((\frac{\mathsf{WT}_w}{\epsilon})^d\right)$  such that

$$1332 \quad \|\bigcirc_{\ell=1}^{\mathbf{T}_w} \mathcal{P}_{\text{NN}}^w - \bigcirc_{\ell=1}^{\mathbf{T}_w} \mathcal{P}_*^w(\cdot, t_{\ell+1}^w - t_\ell^w)\|_{L^\infty(\mathcal{Z}(a_w))} < \frac{\epsilon}{(2\mathsf{W} + 1)\text{Lip}(\mathcal{D})\text{Lip}(\mathcal{F})}.$$

1333 The composition above makes sense, since the range of  $\mathcal{P}_{\text{NN}}^w$  is contained in  
 1334  $\mathcal{Z}([s_w, s_{w+1}])$  by construction. This establishes the result for the case that  $\mathsf{W} = \mathsf{S}$ .  
 1335 To handle the case that  $\mathsf{W} > \mathsf{S}$ , we can use the same argument as in Theorem 3.9 to  
 1336 complete the proof.  $\square$

### 1337 Appendix C. Comparison Model Details.

1338 **C.1. LDNet.** We implemented LDNet which was proposed in [43], but we mod-  
 1339 ified the architecture to be grid-dependent. Specifically, since WeldNet has inputs  
 1340 and outputs on a fixed grid, for a more direct comparison we implemented a grid-  
 1341 dependent LDNet such that the reconstruction network has to predict the values on  
 1342 the grid. The original LDNet implementation uses a reconstruction network that in-  
 1343 puts the query location and outputs the value of the output field at that location. In  
 1344 other words, we would have to call the reconstruction network  $D$  times to output val-  
 1345 ues on a size  $D$  grid. Training the grid based LDNet implementation is subsequently  
 1346 much faster and requires much lower computational resources (such as memory); the  
 1347 original LDNet implementation has significant memory requirements for our data. We  
 1348 found that the original model sizes used in [43] were too small to perform well on our  
 1349 test problems, so we increased the network size to be comparable to the models in  
 1350 LDNet (e.g. width 500 networks), but this lead to a significant memory requirement.

1351 We present a comparison between the original LDNet implementation and the  
 1352 grid based LDNet model. The original model was trained on an NVIDIA A100 GPU

	Original-LDNet	Grid-LDNet
bscale	3.81% / 5.12%	0.78% / 1.36%
bshift	6.13% / 6.87%	10.4% / 6.38%
tscale	4.12% / 4.90%	4.72% / 8.29%
tshift	0.75% / 0.89%	5.98% / 10.3%
kscale	>100% / >100%	11.1% / 22.7%
kshift	0.95% / 1.89%	0.43% / 0.72%

TABLE 3

*Comparison of middle and final time relative test errors for original LDNet (with 101 time steps) and grid LDNet (with 301 time steps). Each cell is formatted as MiddleError / FinalError.*

1353 with 80 GBs of memory, while the grid based model was trained on an NVIDIA  
 1354 RTX6000 (i.e. we train with the same resources as WeldNet and other models). Due  
 1355 to memory limitations, we trained the original LDNet using data on a time grid that  
 1356 was **three** times as coarse as the time used for grid LDNet. With all of these changes,  
 1357 the total training time on the Burgers' scale dataset (with initial conditions from  
 1358 (4.6)) for the original LDNet model is 1335.3s, while the grid LDNet took 1343.6s to  
 1359 train (on a dataset with 3x as many time steps).

1360 Table 3 shows a comparison of the midway (i.e. middle of the time grid) and  
 1361 final time test errors between the original LDNet and the grid LDNet models. We  
 1362 used the tanh activation function for the reconstruction network for the Burgers' and  
 1363 KdV examples, but we used a ReLU activation for the transport examples. Superior  
 1364 performance of the original LDNet model is observed on the transport examples, but  
 1365 it is matched or outperformed with the grid LDNet. Recall that the original LDNet  
 1366 is trained to predict the evolution of 101 time steps but the grid LDNet is trained  
 1367 to predict 301 time steps. Clearly, a grid based LDNet performs roughly similarly or  
 1368 better to the original LDNet implementation for most of our problems using much less  
 1369 computational resources. We note that the original LDNet implementation has the  
 1370 advantage of being grid independent, so the output predictions can be queried at all  
 1371 locations regardless of the grid it was trained on. This is also advantageous for some  
 1372 datasets such as for the two-hat-shaped initial conditions we used for the transport  
 1373 equation. However, for our purpose of grid based surrogate modeling, a grid based  
 1374 LDNet is faster to train and performs better on lower resources (and on a more fine  
 1375 time grid), so we use that for comparison in the rest of the paper.

1376 **C.2. DeepONet.** Latent-DON consists of an autoencoder (which is identical to  
 1377 the autoencoder used in a one-window WeldNet model) and a DeepONet that predicts  
 1378 the evolution of a latent code from its initial time to a given future time. DeepONet  
 1379 is an operator learning architecture [32] that can be used in this case as time can be  
 1380 considered as the “input” to the operator. Specifically, a DeepONet is a function of  
 1381 the form (for some  $p \in \mathbb{N}$ ):

1382 (C.1) 
$$\phi(\vec{z})(t) = \tau_{\text{NN}}(t) \cdot B_{\text{NN}}(\vec{z})$$

1383 Here,  $B_{\text{NN}}$  is a neural network with input dimension  $k$  (i.e. the latent space dimension)  
 1384 and output dimension  $pk$  which is then reshaped to be dimension  $p \times k$ . On the other  
 1385 hand,  $\tau_{\text{NN}}$  is a neural network with input dimension 1 and output dimension  $p$ , so the  
 1386 product in (C.1) outputs a vector of dimension  $k$ .

1387 This is the formulation used in [24]. We implement it in Pytorch such that all  
 1388 components of latent DeepONet (i.e. encoder, decoder, branch net, trunk net) are

	Weld-1	Weld-2	Weld-2*	Weld-4	Weld-4*
Time	668.5s	547.7s	373.2s	528.1s	267.6s

TABLE 4

Total training time for each model on  $\mathcal{M}_{bscale}$  data for WeldNet models. \* indicates training with multiple GPUs.

	LDON	LDNet	Time Input	HDP	WLaSDI
Time	65.2s	1355.3s	52.0s	56.1s	2285s

TABLE 5

Total training time for each model on  $\mathcal{M}_{bscale}$  data for comparison models.

1389 width 400 and depth 3 feedforward ReLU networks. We use  $p = 10$  for the Latent-  
 1390 DON comparison in this work.

1391 **Appendix D. Error Tables.** For each dataset, we compute the relative test  
 1392 errors at different times and show them in Table 6-12. The best model error in the  
 1393 final time is bolded. A dash “–” indicates a relative test over of 10 or higher.

1394 We also display the total training time for each model on the Burgers’ scale  
 1395 dataset in Table 4 for WeldNet models and Table 5 for comparison models. The total  
 1396 training times for other models are similar. We use NVIDIA RTX 6000 GPUs. For  
 1397 WeldNet-2 and WeldNet-4, we use one GPU per window since the training of each  
 1398 window’s models is completely independent from the other windows, but we use only  
 1399 one GPU for other models (there is no easy way to distribute the training among  
 1400 multiple GPUs for those models).

Model	30	60	90	120	150	180	210	240	270	300
FF-Weld-1	0.0091	0.0095	0.0111	0.0133	0.0145	0.0178	0.0220	0.0239	0.0241	0.0343
FF-Weld-2	0.0068	0.0045	0.0050	0.0053	0.0058	0.0122	0.0088	0.0089	0.0100	0.0101
FF-Weld-4	0.0063	0.0071	0.0059	0.0071	0.0175	0.0073	0.0081	0.0082	0.0084	0.0137
PCA-WeldNet-1	0.1875	0.1177	0.1078	0.1176	0.1307	0.1561	0.1798	0.1990	0.2174	0.2404
PCA-WeldNet-2	0.1077	0.0420	0.0656	0.0578	0.0909	0.1675	0.1408	0.1204	0.1220	0.1551
PCA-WeldNet-4	0.0445	0.0254	0.0292	0.0524	0.0440	0.1145	0.0885	0.1223	0.1213	0.1133
Conv-Weld-1	0.0099	0.0070	0.0064	0.0065	0.0074	0.0072	0.0076	0.0076	0.0080	0.0091
Conv-Weld-2	0.0090	0.0056	0.0055	0.0064	0.0066	0.0121	0.0082	0.0084	0.0090	0.0094
Conv-Weld-4	0.0159	0.0119	0.0101	0.0067	0.0069	0.0089	0.0063	0.0074	0.0057	<b>0.0060</b>
Latent-DON	0.0059	0.0065	0.0069	0.0072	0.0078	0.0083	0.0092	0.0091	0.0092	0.0153
LDNet	0.0440	0.0373	0.0370	0.0375	0.0381	0.0374	0.0369	0.0360	0.0388	0.0512
Grid-LDNet	0.0063	0.0070	0.0072	0.0071	0.0078	0.0083	0.0097	0.0097	0.0095	0.0136
TimeInput	0.0258	0.0135	0.0116	0.0126	0.0123	0.0126	0.0145	0.0159	0.0178	0.0176
HDP	0.0209	0.0313	0.0403	0.0553	0.0825	0.0967	0.1116	0.1313	0.1323	0.1433
WLaSDI	1.045	1.042	1.055	1.072	1.097	1.137	1.192	1.263	1.359	1.486

TABLE 6  
*Test Errors for Burgers' Scale*

Model	30	60	90	120	150	180	210	240	270	300
FF-Weld-2	0.0232	0.0181	0.0246	0.0427	0.0455	0.0484	0.0363	0.0287	0.0247	0.0242
FF-Weld-4	0.0151	0.0364	0.0570	0.0669	0.0623	0.0557	0.0492	0.0536	0.0468	0.0472
FF-AENet	0.0218	0.0310	0.0468	0.0542	0.0444	0.0364	0.0312	0.0284	0.0289	0.0340
Conv-Weld-2	0.0245	0.0147	0.0170	0.0202	0.0205	0.0252	0.0178	0.0155	0.0171	<b>0.0201</b>
Conv-Weld-4	0.0199	0.0116	0.0153	0.0237	0.0278	0.0295	0.0229	0.0253	0.0227	0.0205
Conv-AENet	0.0328	0.0204	0.0291	0.0442	0.0554	0.0627	0.0738	0.0867	0.0998	0.1119
PCA-WeldNet-2	0.1500	0.1628	0.3493	0.4804	0.4964	0.4950	0.4884	0.4806	0.4727	0.4651
PCA-WeldNet-4	0.1484	0.1568	0.3466	0.4830	0.4980	0.4949	0.4882	0.4805	0.4736	0.4719
PCA-AENet	0.1521	0.2058	0.3729	0.4967	0.5127	0.5107	0.5060	0.5009	0.4983	0.5021
Latent-DON	0.0209	0.0291	0.0400	0.0460	0.0386	0.0320	0.0272	0.0248	0.0261	0.0318
LDNet	0.0254	0.0399	0.1043	0.1158	0.1041	0.0919	0.0811	0.0722	0.0657	0.0638
TimeInput	0.0291	0.0222	0.0323	0.0548	0.0550	0.0451	0.0366	0.0303	0.0267	0.0265
HDP	0.2155	0.4972	0.8060	1.2015	1.7128	2.3725	3.2132	4.2975	5.6602	7.3479
WLaSDI	0.782	0.839	0.854	0.835	0.844	0.910	0.978	1.045	1.199	1.380

TABLE 7  
*Test Errors for Burgers' Shift*

Model	30	60	90	120	150	180	210	240	270	300
FF-Weld-2	0.0095	0.0125	0.0103	0.0139	0.0120	0.0085	0.0140	0.0118	0.0113	0.0142
FF-Weld-4	0.0136	0.0112	0.0116	0.0141	0.0125	0.0160	0.0135	0.0094	0.0141	<b>0.0102</b>
FF-AENet	0.0215	0.0204	0.0223	0.0196	0.0195	0.0224	0.0245	0.0217	0.0249	0.0362
Conv-Weld-2	0.0221	0.0234	0.0235	0.0248	0.0181	0.0198	0.0248	0.0250	0.0225	0.0268
Conv-Weld-4	0.0168	0.0163	0.0155	0.0212	0.0202	0.0248	0.0311	0.0292	0.0190	0.0203
Conv-AENet	0.0298	0.0427	0.0346	0.0516	0.0445	0.0439	0.0271	0.0324	0.0300	0.0234
PCA-WeldNet-2	0.5815	0.4339	0.4880	0.5128	0.5090	0.5813	0.4336	0.4875	0.5101	0.5010
PCA-WeldNet-4	0.2867	0.1783	0.2379	0.2582	0.1697	0.2860	0.1790	0.2388	0.2571	0.1705
PCA-AENet	0.8214	0.5968	0.5885	0.6678	0.6158	0.6106	0.6205	0.6837	0.6580	0.5872
Latent-DON	0.0156	0.0155	0.0218	0.0217	0.0165	0.0193	0.0174	0.0204	0.0207	0.0318
LDNet	0.0359	0.0416	0.0492	0.0459	0.0472	0.0476	0.0453	0.0504	0.0463	0.0829
TimeInput	0.0580	0.0668	0.0591	0.0669	0.0679	0.0610	0.0640	0.0703	0.0661	0.0921
HDP	0.0189	0.0348	0.0767	0.1062	0.1366	0.1480	0.1629	0.1946	0.1881	0.0697
WLaSDI	1.210	1.671	1.521	3.540	9.388	5.028	7.136	—	—	—

TABLE 8  
*Test Errors for Transport Scale*

Model	30	60	90	120	150	180	210	240	270	300
FF-Weld-2	0.0283	0.0295	0.0276	0.0308	0.0305	0.0405	0.1727	0.1273	0.0631	0.0599
FF-Weld-4	0.0310	0.0297	0.0353	0.0372	0.0432	0.0357	0.0374	0.0457	0.0559	0.0587
FF-AENet	0.7640	0.8940	0.9187	0.9386	0.9492	0.9243	0.8009	0.8958	0.8003	1.1303
Conv-Weld-2	0.0207	0.0186	0.0168	0.0189	0.0187	0.0208	0.0213	0.0221	0.0188	0.0210
Conv-Weld-4	0.0173	0.0151	0.0154	0.0186	0.0183	0.0180	0.0213	0.0231	0.0200	0.0227
Conv-AENet	0.0312	0.0279	0.0330	0.0293	0.0276	0.0275	0.0254	0.0320	0.0321	0.0319
PCA-WeldNet-2	0.7819	0.7177	0.6850	0.6988	0.7210	0.7889	0.7401	0.6998	0.7150	0.7326
PCA-WeldNet-4	0.6773	0.6055	0.5905	0.6087	0.6159	0.6777	0.6055	0.5906	0.6095	0.6176
PCA-AENet	0.8650	0.8727	0.8775	0.9071	0.8685	0.8806	0.8619	0.8546	0.8281	0.8190
Latent-DON	0.1120	0.1088	0.1519	0.1396	0.1161	0.1438	0.1252	0.1412	0.1473	0.3016
LDNet	0.0503	0.0521	0.0565	0.0589	0.0598	0.0616	0.0609	0.0618	0.0664	0.1029
TimeInput	0.1435	0.1065	0.0966	0.1035	0.1085	0.1059	0.1188	0.1048	0.0955	0.1031
HDP	0.3034	2.3470	5.1519	8.4939	—	—	—	—	—	—
WLaSDI	6.216	—	—	—	—	—	—	—	—	—

TABLE 9  
Test Errors for Transport Shift

Model	30	60	90	120	150	180	210	240	270	300
FF-Weld-2	0.0160	0.0249	0.0329	0.0273	0.0281	0.0277	0.0366	0.0332	0.0338	0.0371
FF-Weld-4	0.0096	0.0101	0.0127	0.0128	0.0240	0.0166	0.0187	0.0216	0.0272	0.0395
FF-AENet	0.0795	0.1304	0.1255	0.1307	0.1696	0.1798	0.2025	0.2515	0.4120	0.6350
Conv-Weld-2	0.0261	0.0231	0.0225	0.0235	0.0240	0.0232	0.0288	0.0283	0.0308	0.0316
Conv-Weld-4	0.0106	0.0112	0.0128	0.0173	0.0186	0.0227	0.0206	0.0234	0.0283	<b>0.0317</b>
Conv-AENet	0.0187	0.0514	0.0864	0.1098	0.1423	0.1818	0.2392	0.2971	0.3628	0.3851
PCA-WeldNet-2	0.5672	0.3454	0.3709	0.3762	0.4345	0.6621	0.6137	0.5912	0.5849	0.6024
PCA-WeldNet-4	0.3723	0.2473	0.2717	0.5046	0.4217	0.5999	0.5589	0.5352	0.6763	0.6782
PCA-AENet	0.6709	0.4938	0.4819	0.4796	0.4769	0.4968	0.5399	0.5860	0.6234	0.6611
Latent-DON	0.2379	0.1528	0.1858	0.1949	0.1685	0.1378	0.1225	0.1438	0.1914	0.3889
LDNet	0.0735	0.0793	0.0890	0.1011	0.1105	0.1262	0.1393	0.1550	0.1734	0.2272
TimeInput	0.1673	0.1505	0.1362	0.1506	0.1570	0.1700	0.1767	0.1878	0.2060	0.2472
HDP	0.0770	0.1302	0.1688	0.2621	0.3905	0.6109	0.8120	0.9843	1.0881	1.1590
WLaSDI	0.993	1.417	2.495	5.950	—	—	—	—	—	—

TABLE 10  
Test Errors for KdV Scale

Model	30	60	90	120	150	180	210	240	270	300
FF-Weld-2	0.0028	0.0027	0.0028	0.0025	0.0027	0.0031	0.0030	0.0030	0.0030	<b>0.0028</b>
FF-Weld-4	0.0026	0.0028	0.0028	0.0027	0.0041	0.0030	0.0032	0.0033	0.0032	0.0046
FF-AENet	0.0071	0.0069	0.0065	0.0069	0.0063	0.0061	0.0058	0.0051	0.0053	0.0087
Conv-Weld-2	0.0044	0.0034	0.0036	0.0032	0.0036	0.0042	0.0034	0.0031	0.0031	0.0032
Conv-Weld-4	0.0040	0.0030	0.0031	0.0034	0.0031	0.0038	0.0030	0.0030	0.0028	<b>0.0028</b>
Conv-AENet	0.0050	0.0053	0.0046	0.0042	0.0046	0.0042	0.0039	0.0046	0.0047	0.0053
PCA-WeldNet-2	0.0444	0.0233	0.0192	0.0215	0.0230	0.0482	0.0210	0.0177	0.0233	0.0186
PCA-WeldNet-4	0.0199	0.0119	0.0055	0.0051	0.0098	0.0168	0.0105	0.0057	0.0059	0.0090
PCA-AENet	0.1276	0.0832	0.0672	0.0551	0.0476	0.0542	0.0629	0.0587	0.0473	0.0628
Latent-DON	0.0047	0.0039	0.0044	0.0044	0.0035	0.0039	0.0038	0.0032	0.0032	0.0051
LDNet	0.0049	0.0041	0.0044	0.0045	0.0043	0.0043	0.0042	0.0038	0.0037	0.0072
TimeInput	0.0130	0.0071	0.0081	0.0099	0.0067	0.0074	0.0074	0.0064	0.0073	0.0071
HDP	0.0115	0.0226	0.0212	0.0253	0.0273	0.0321	0.0319	0.0332	0.0360	0.0426
WLaSDI	0.464	0.463	0.476	0.499	0.521	0.520	0.523	0.520	0.505	0.497

TABLE 11  
Test Errors for KdV Shift

Model	10	20	30	40	50	60	70	80	90	100
FF-Weld-2	0.0051	0.0048	0.0043	0.0042	0.0041	0.0030	0.0030	0.0030	0.0032	0.0034
FF-Weld-4	0.0051	0.0045	0.0037	0.0034	0.0031	0.0027	0.0026	0.0027	0.0026	0.0028
FF-AENet	0.0064	0.0056	0.0051	0.0051	0.0045	0.0047	0.0045	0.0042	0.0036	0.0043
Latent-DON	0.0080	0.0063	0.0052	0.0046	0.0044	0.0047	0.0047	0.0049	0.0045	0.0061
LDNet	0.0086	0.0074	0.0072	0.0072	0.0066	0.0065	0.0063	0.0063	0.0073	0.0095
TimeInput	0.0155	0.0129	0.0121	0.0120	0.0118	0.0119	0.0121	0.0115	0.0140	0.0194
HDP	0.0389	0.0312	0.0298	0.0294	0.0274	0.0262	0.0267	0.0269	0.0250	0.0316

TABLE 12  
Test Errors for Shallow Water