
Physics-Informed Regularization for Domain-Agnostic Dynamical System Modeling

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Learning complex physical dynamics purely from data is challenging due to the
2 intrinsic properties of systems to be satisfied. Incorporating physics-informed
3 priors, such as in Hamiltonian Neural Networks (HNNs), achieves high-precision
4 modeling for energy-conservative systems. However, real-world systems often
5 deviate from strict energy conservation and follow different physical priors. To ad-
6 dress this, we present a framework that achieves high-precision modeling for a wide
7 range of dynamical systems from the numerical aspect, by enforcing *Time-Reversal*
8 *Symmetry* (TRS) via a novel regularization term. It helps preserve energies for
9 conservative systems while serving as a strong inductive bias for non-conservative,
10 reversible systems. While TRS is a domain-specific physical prior, we present the
11 *first* theoretical proof that TRS loss can universally improve modeling accuracy by
12 minimizing higher-order Taylor terms in ODE integration, which is numerically
13 beneficial to various systems regardless of their properties, even for irreversible
14 systems. By integrating the TRS loss within neural ordinary differential equation
15 models, the proposed model TREAT demonstrates superior performance on diverse
16 physical systems. It achieves a significant 11.5% MSE improvement in a challeng-
17 ing chaotic triple-pendulum scenario, underscoring TREAT’s broad applicability
18 and effectiveness. Code and further details are available at [here](#).

19 1 Introduction

20 Dynamical systems, spanning applications from physical simulations (Kipf et al., 2018; Wang et al.,
21 2020; Lu et al., 2022) to robotic control (Li et al., 2022; Ni and Qureshi, 2022), are challenging
22 to model due to intricate dynamic patterns and potential interactions under multi-agent settings.
23 Traditional numerical simulators require extensive domain knowledge for design, which is sometimes
24 unknown (Sanchez-Gonzalez et al., 2020), and can consume significant computational resources.
25 Therefore, directly learning dynamics from the observational data becomes an attractive alternative.

26 Existing deep learning approaches (Sanchez-Gonzalez et al., 2020; Pfaff et al., 2021) usually learn a
27 fixed-step transition function to predict system dynamics from timestamp t to timestamp $t + 1$ and
28 rollout trajectories recursively. The transition function can have different inductive biases, such as
29 Graph Neural Networks (GNNs) (Lam et al., 2023) for capturing pair-wise interactions among agents
30 through message passing. Most recently, neural ordinary differential equations (ODEs) (Chen et al.,
31 2018; Rubanova et al., 2019) have emerged as a potent solution for modeling system dynamics in a
32 continuous manner, which offer superior prediction accuracy over discrete models in the long-range,
33 and can handle systems with partial observations. In particular, GraphODEs (Huang et al., 2020;
34 Luo et al., 2023; Zang and Wang, 2020; Jiang et al., 2023) extend NeuralODEs to model interacting
35 (multi-agent) dynamical systems, where agents co-evolve and form trajectories jointly.

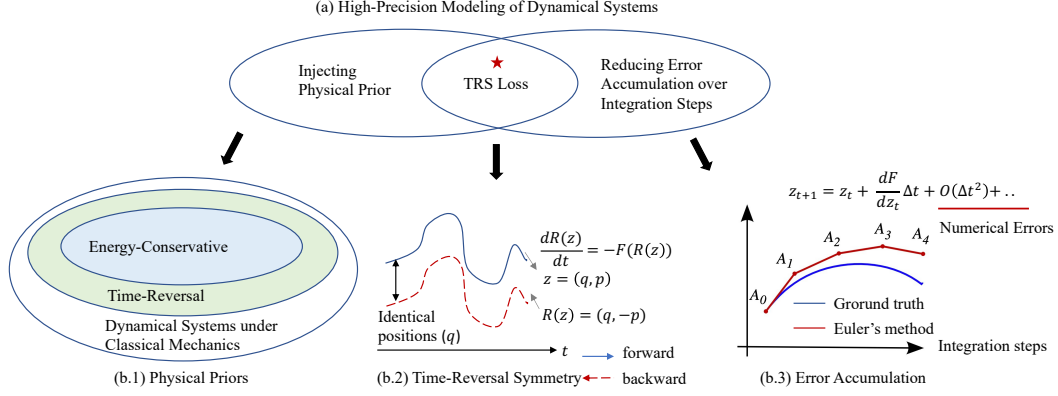


Figure 1: (a) High-precision modeling for dynamical systems; (b.1) Classification of classical mechanical systems based on (Tolman, 1938; Lamb and Roberts, 1998); (b.2) Time-Reversal Symmetry illustration; (b.3) Error accumulation in numerical solvers.

However, the complexity of dynamical systems necessitates large amounts of data. Models trained on limited data risk violating fundamental physical principles such as energy conservation. A promising strategy to improve modeling accuracy involves incorporating physical inductive biases (Raissi et al., 2019; Cranmer et al., 2020). Existing models like Hamiltonian Neural Networks (HNNs) (Greydanus et al., 2019; Sanchez-Gonzalez et al., 2019) strictly enforce energy conservation, yielding more accurate predictions for energy-conservative systems. However, not all real-world systems strictly adhere to energy conservation, and they may adhere to various physical priors. Such system diversity largely limits the usage of existing models which are designed for individual physical prior.

To address this, we present a framework that achieves high-precision modeling for a wide range of dynamical systems from the numerical aspect, by enforcing Time-Reversal Symmetry (TRS) via a novel regularization term. Specifically, TRS posits that a system’s dynamics should remain invariant when time is reversed (Lamb and Roberts, 1998). To incorporate TRS, we propose a simple-yet-effective self-supervised regularization term that acts as a soft constraint. This term aligns *forward and backward trajectories* predicted by a neural network and we use GraphODE as the backbone. We theoretically prove that the TRS loss effectively minimizes higher-order Taylor expansion terms during ODE integration, offering a general numerical advantage for improving modeling accuracy across a wide array of systems, regardless of their physical properties. It forces the model to capture fine-grained physical properties such as jerk (the derivatives of accelerations) and provides more regularization for long-term prediction. We also justify our TRS design choice, showing case its superior performance both analytically and empirically. We name the model as TREAT (**T**ime-**R**eversal Symmetry ODE).

Note that TRS itself is a physical prior, that is broader than energy conservation as depicted in Figure 1(b.1). It covers classical energy-conservative systems such as Newtonian mechanics, and also non-conservative, reversible systems like Stokes flow (Pozrikidis, 2001), commonly encountered in microfluidics (Kim and Karrila, 2013; Cao and Li, 2018; Cao et al., 2019). Therefore, TRS loss achieves high-precision modeling from both the physical aspect, and the numerical aspect as shown in Figure 1(a), making it domain-agnostic and widely applicable to various dynamical systems. We systematically conduct experiments across 9 diverse datasets spanning across 1.) single-agent, multi-agent systems; 2.) simulated and real-world systems; and 3.) systems with different physical priors. TREAT consistently outperforms state-of-the-art baselines, affirming its effectiveness and versatility across various dynamic scenarios.

Our primary contributions can be summarized as follows:

- We introduce TREAT, a powerful framework that achieves high-precision modeling for a wide range of systems from the numerical aspect, by enforcing Time-Reversal Symmetry (TRS) via a regularization term.
- We establish the *first* theoretical proof that the time-reversal symmetry loss could in general help learn more fine-grained and long-context system dynamics from the numerical aspect,

regardless of systems’ physical properties (even irreversible systems). This bridges the specific physical implication and the general numerical benefits of the physical prior -TRS.

- We present empirical evidence of TREAT’s state-of-the-art performance in a variety of systems over 9 datasets, including real-world & simulated systems, etc. It yields a significant MSE improvement of 11.5% on the challenging chaotic triple-pendulum system.

2 Preliminaries and Related Work

We represent a dynamical system as a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} denotes the node set of N agents¹ and \mathcal{E} denotes the set of edges representing their physical interactions. For simplicity, we assumed \mathcal{G} to be static over time. Single-agent dynamical system is a special case where the graph only has one node. In the following, we use the multi-agent setting by default to illustrate our model. We denote $\mathbf{X}(t) \in \mathbb{R}^{N \times d}$ as the feature matrix at timestamp t for all agents, with d as the feature dimension. Model input consists of trajectories of feature matrices over M historical timestamps $X(t_{-M:-1}) = \{\mathbf{X}(t_{-M}), \dots, \mathbf{X}(t_{-1})\}$ and \mathcal{G} . The timestamps $t_{-1}, \dots, t_{-M} < 0$ can have non-uniform intervals and take any continuous values. Our goal is to learn a neural simulator $f_\theta(\cdot) : [X(t_{-M:-1}), \mathcal{G}] \rightarrow Y(t_{0:K})$, which predicts node dynamics $\mathbf{Y}(t)$ in the future on timestamps $0 = t_0 < \dots < t_K = T$ sampled within $[0, T]$. We use $\mathbf{y}_i(t)$ to denote the targeted dynamic vector of agent i at time t . In some cases when we are only predicting system feature trajectories, $\mathbf{Y}(\cdot) \equiv \mathbf{X}(\cdot)$.

2.1 NeuralODE for Dynamical Systems

NeuralODEs (Chen et al., 2018; Rubanova et al., 2019) are a family of continuous models that define the evolution of dynamical systems by ordinary differential equations (ODEs). The state evolution can be described as: $\dot{\mathbf{z}}_i(t) := \frac{d\mathbf{z}_i(t)}{dt} = g(\mathbf{z}_1(t), \mathbf{z}_2(t) \dots \mathbf{z}_N(t))$, where $\mathbf{z}_i(t) \in \mathbb{R}^d$ denotes the latent state variable for agent i at timestamp t . The ODE function g is parameterized by a neural network such as Multi-Layer Perception (MLP), which is automatically learned from data. GraphODEs (Poli et al., 2019; Huang et al., 2020; Luo et al., 2023; Wen et al., 2022) are special cases of NeuralODEs, where g is a Graph Neural Network (GNN) to capture the continuous interaction among agents.

GraphODEs have been shown to achieve superior performance, especially in long-range predictions and can handle data irregularity issues. They usually follow the encoder-processor-decoder architecture, where an encoder first computes the latent initial states $\mathbf{z}_1(t_0), \dots, \mathbf{z}_N(t_0)$ for all agents simultaneously based on their historical observations as in Eqn 1.

$$\mathbf{z}_1(t_0), \mathbf{z}_2(t_0), \dots, \mathbf{z}_N(t_0) = f_{\text{ENC}}(X(t_{-M:-1}), \mathcal{G}) \quad (1)$$

Then the GNN-based ODE predicts the latent trajectories starting from the learned initial states. The latent state $\mathbf{z}_i(t)$ can be computed at any desired time using a numerical solver such as Runge-Kutta (Schober et al., 2019) as:

$$\mathbf{z}_i(t) = \text{ODE-Solver}(g, [\mathbf{z}_1(t_0), \dots, \mathbf{z}_N(t_0)], t) = \mathbf{z}_i(t_0) + \int_{t_0}^t g(\mathbf{z}_1(t), \mathbf{z}_2(t) \dots \mathbf{z}_N(t)) dt. \quad (2)$$

Finally, a decoder extracts the predicted dynamics $\hat{\mathbf{y}}_i(t)$ based on the latent states $\mathbf{z}_i(t)$ for any timestamp t :

$$\hat{\mathbf{y}}_i(t) = f_{\text{DEC}}(\mathbf{z}_i(t)). \quad (3)$$

However, vanilla GraphODEs can violate physical properties of a system, resulting in unrealistic predictions. We therefore propose to inject physics-informed regularization term to make more accurate predictions.

2.2 Time-Reversal Symmetry (TRS)

¹Following (Kipf et al., 2018), we use “agents” to denote “objects” in dynamical systems, which is different from “intelligent agent” in AI.

111 Consider a dynamical system de-
 112 scribed in the form of $\frac{d\mathbf{x}(t)}{dt} =$
 113 $F(\mathbf{x}(t))$, where $\mathbf{x}(t) \in \Omega$ is the ob-
 114 served states such as positions. The
 115 system is said to follow the *Time-*
 116 *Reversal Symmetry* if there exists a
 117 reversing operator $R : \Omega \mapsto \Omega$ such
 118 that (Lamb and Roberts, 1998):

$$\frac{d(R \circ \mathbf{x}(t))}{dt} = -F(R \circ \mathbf{x}(t)), \quad (4)$$

119 where \circ denote the action of func-
 120 tional R on the function \mathbf{x} .

121 Intuitively, we can assume $\mathbf{x}(t)$ is the position of a flying ball and the conventional reversing operator
 122 is defined as $R : \mathbf{x} \mapsto R \circ \mathbf{x}$, $R \circ \mathbf{x}(t) = \mathbf{x}(-t)$. This implies when $\mathbf{x}(t)$ is a forward trajectory
 123 position with initial position $\mathbf{x}(0)$, $\mathbf{x}(-t)$ is then a position in the time-reversal trajectory, where
 124 $\mathbf{x}(-t)$ is calculated using the same function F , but with the integration time reversed, i.e. $dt \mapsto d(-t)$.
 125 Eqn 4 shows how to create the reverse trajectory of a flying ball: at each position, the velocity (i.e.,
 126 the derivative of position with respect to time) should be the opposite. In neural networks, we usually
 127 model trajectories in the latent space via \mathbf{z} (Sanchez-Gonzalez et al., 2020), which can be decoded
 128 back to real observation state i.e. positions. Therefore, we apply the reversal operator for \mathbf{z} .

129 Now we introduce a time evolution operator ϕ_τ such that $\phi_\tau \circ \mathbf{z}(t) = \mathbf{z}(t + \tau)$ for arbitrary $t, \tau \in \mathbb{R}$.
 130 It satisfies $\phi_{\tau_1} \circ \phi_{\tau_2} = \phi_{\tau_1 + \tau_2}$, where \circ denotes composition. The time evolution operator helps us
 131 to move forward (when $\tau > 0$) or backward (when $\tau < 0$) through time, thus forming a trajectory.
 132 Based on (Lamb and Roberts, 1998), in terms of the evolution operator, Eqn 4 implies:

$$R \circ \phi_t = \phi_{-t} \circ R = \phi_t^{-1} \circ R, \quad (5)$$

133 which means that moving forward t steps and then turning to the opposite direction is equivalent
 134 to firstly turning to the opposite direction and then moving backwards t steps². Eqn 5 has been
 135 widely used to describe time-reversal symmetry in existing literature (Huh et al., 2020; Valperga
 136 et al., 2022). Nevertheless, we propose the following lemma, which is more intuitive to understand
 137 and straightforward to guide the design of our time-reversal regularizer.

138 **Lemma 2.1.** *Eqn 5 is equivalent to $R \circ \phi_t \circ R \circ \phi_t = I$, where I denotes identity mapping.*

139 Lemma 2.1 means if we move t steps forward, then turn to the opposite direction, and then move
 140 forward for t more steps, it shall restore back to the same state. This is illustrated in Figure 2 where
 141 the reverse trajectory should be the same as the forward trajectory.³ It can be understood as rewinding
 142 a video to the very beginning. The proof of Lemma 2.1 is in Appendix A.2.

143 3 Method: TREAT

144 We present a novel framework TREAT that achieves high-precision modeling for a wide range of
 145 systems from the numerical aspect, by enforcing Time-Reversal Symmetry (TRS) via a regularization
 146 term. It improves modeling accuracy regardless of systems' physical properties. We first introduce
 147 our architecture design, followed by theoretical analysis to explain its numerical benefits.

148 TREAT uses GraphODE (Huang et al., 2020) as the backbone and flexibly incorporates TRS as a
 149 regularization term based on Lemma 2.1. This term aligns model forward and reverse trajectories. In
 150 practice, our model predicts the forward trajectories at a series of timestamps $\{t_k\}_{k=0}^K$ as ground truth
 151 observations are discrete, where $0 = t_0 < t_1 < \dots < t_K = T$. The reverse trajectories are also at the
 152 same series of K timestamps so as to be aligned with the forward one, which we denote as $\{t'_k\}_{k=0}^K$
 153 satisfying $0 = t'_0 < t'_1 < \dots < t'_K = T$. It's important to note that the values of the time variable

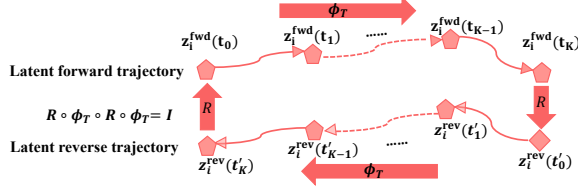


Figure 2: Illustration of time-reversal symmetry based on Lemma 2.1. The total length of the trajectory is $t_K - t_0 = T$. t'_k is the time index in the reverse trajectory, which points to the same time as t_{K-k} in the forward trajectory.

²Time-reversal symmetry is a property of physical systems, which requires the forward and reverse trajectories to be generated by the same mechanism $F(\cdot)$. It differs from reversibility of neural networks (Chang et al., 2018; Liu et al., 2019), which is a property of machine learning models and ensures the recovery of input from output via a reversed operator $f^{-1}(\cdot)$. We highlight the detailed discussions in Appendix F.

³We explain Figure 2 with implementation in Appendix A.1.

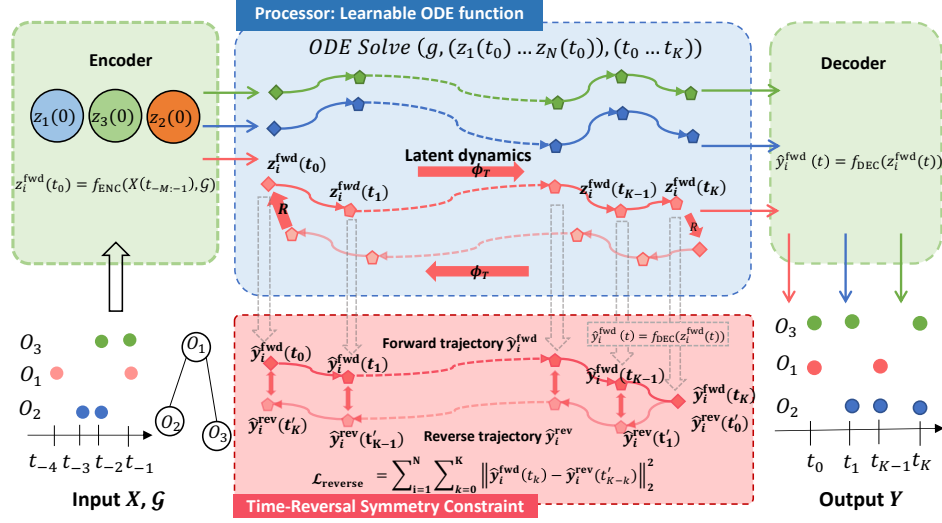


Figure 3: Overall framework of TREAT. O_1, O_2, O_3 are connected agents. It follows the encoder-processor-decoder architecture introduced in Sec 2.1. A novel TRS loss is incorporated to improve modeling accuracy across systems from the numerical aspect, regardless of their physical properties.

154 t'_k in the reverse trajectories do not represent real time, but serve as indexes of reverse trajectories.
 155 This leads to the relation $t'_{K-k} = T - t_k$, which means the reverse trajectories at timestamp t'_{K-k}
 156 correspond to the forward trajectories at time t_k . For example, $t'_0 = T - t_K = 0$. It indicates t'_0 and
 157 t_K are both pointing to the same real time T , which is the ending point of the forward trajectory as
 158 shown in Figure 3. Based on Lemma 2.1, the difference of the two trajectories at any observed time
 159 should be small, i.e. $z_i^{fwd}(t_k) \approx z_i^{rev}(t'_{K-k})$. This serves as the guideline for our regularizer design.
 160 The weight of the regularizer is also adjustable to adapt different systems. The overall framework is
 161 depicted in Figure 3.

162 3.1 Time-Reversal Symmetry Loss and Training

163 **Forward Trajectory Prediction and Reconstruction Loss.** For multi-agent systems, we utilize
 164 the GNN operator described in (Kipf et al., 2018) as our ODE function $g(\cdot)$, which drives the system
 165 to move forward and output the forward trajectories for latent states $z_i^{fwd}(t)$ at each continuous time
 166 $t \in [0, T]$ and each agent i . We then employ a Multilayer Perceptron (MLP) as a decoder to predict
 167 output trajectories $\hat{y}_i^{fwd}(t)$ based on the latent states. We summarize the whole procedure as:

$$\begin{aligned} \dot{z}_i^{fwd}(t) &:= \frac{dz_i^{fwd}(t)}{dt} = g(z_1^{fwd}(t), z_2^{fwd}(t), \dots, z_N^{fwd}(t)), \\ z_i^{fwd}(t_0) &= f_{ENC}(X(t_{-M:-1}), \mathcal{G}), \quad \hat{y}_i^{fwd}(t) = f_{DEC}(z_i^{fwd}(t)). \end{aligned} \quad (6)$$

168 To train the model, we use the reconstruction loss that minimizes the L2 distance between predicted
 169 forward trajectories $\{\hat{y}_i^{fwd}(t_k)\}_{k=0}^K$ and the ground truth trajectories $\{y_i(t_k)\}_{k=0}^K$ as :

$$\mathcal{L}_{pred} = \sum_{i=1}^N \sum_{k=0}^K \left\| y_i(t_k) - \hat{y}_i^{fwd}(t_k) \right\|_2^2. \quad (7)$$

170 **Reverse Trajectory Prediction and Regularization Loss.** We design a novel time-reversal symme-
 171 try loss as a soft constraint to flexibly regulate systems' behavior based on Lemma 2.1. Specifically,
 172 we first compute the latent reverse trajectories $z_i^{rev}(t)$ by starting from the ending state of the forward
 173 one, traversed back over time. We then employ the decoder to output dynamic trajectories $y^{rev}(t)$.

$$\begin{aligned} \dot{z}_i^{rev}(t) &:= \frac{dz_i^{rev}(t)}{dt} = -g(z_1^{rev}(t), z_2^{rev}(t), \dots, z_N^{rev}(t)), \\ z_i^{rev}(t'_0) &= z_i^{fwd}(t_K), \quad \hat{y}_i^{rev}(t) = f_{DEC}(z_i^{rev}(t)). \end{aligned} \quad (8)$$

Next, based on Lemma 2.1, if the system follows *Time-Reversal Symmetry*, the forward and backward trajectories shall be exactly overlap. We thus design the reversal loss by minimizing the L2 distances between model forward and backward trajectories decoded from the latent trajectories:

$$\mathcal{L}_{reverse} = \sum_{i=1}^N \sum_{k=0}^K \left\| \hat{\mathbf{y}}_i^{\text{fwd}}(t_k) - \hat{\mathbf{y}}_i^{\text{rev}}(t'_{K-k}) \right\|_2^2. \quad (9)$$

Finally, we jointly train TREAT as a weighted combination of the two losses:

$$\mathcal{L} = \mathcal{L}_{pred} + \alpha \mathcal{L}_{reverse} = \sum_{i=1}^N \sum_{k=0}^K \left\| \mathbf{y}_i(t_k) - \hat{\mathbf{y}}_i^{\text{fwd}}(t_k) \right\|_2^2 + \alpha \sum_{i=1}^N \sum_{k=0}^K \left\| \hat{\mathbf{y}}_i^{\text{fwd}}(t_k) - \hat{\mathbf{y}}_i^{\text{rev}}(t'_{K-k}) \right\|_2^2, \quad (10)$$

where α is a positive coefficient to balance the two losses based on different targeted systems.

Remark. The computational time of $\mathcal{L}_{reverse}$ is of the same scale as the reconstruction loss \mathcal{L}_{pred} . As the computation process of the reversal loss is to first use the ODE solver to generate the reverse trajectories, which has the same computational overhead as computing the forward trajectories, and then compute the L2 distances.

3.2 Theoretical Analysis of Time-Reversal Symmetry Loss

We next theoretically show that the time-reversal symmetry loss numerically helps to improve prediction accuracy in general, regardless of systems' physical properties. Specifically, we show that it minimizes higher-order Taylor expansion terms during the ODE integration steps.

Theorem 3.1. *Let Δt denote the integration step size in an ODE solver and T be the prediction length. The reconstruction loss \mathcal{L}_{pred} defined in Eqn 7 is $\mathcal{O}(T^3 \Delta t^2)$. The time-reversal loss $\mathcal{L}_{reverse}$ defined in Eqn 9 is $\mathcal{O}(T^5 \Delta t^4)$.*

We prove Theorem 3.1 in Appendix A.3. From Theorem 3.1, we can see two nice properties of our proposed time-reversal loss: 1) Regarding the relationship to Δt , $\mathcal{L}_{reverse}$ is optimizing a high-order term Δt^4 , which forces the model to predict fine-grained physical properties such as jerk (the derivatives of accelerations). In comparison, the reconstruction loss optimizes Δt^2 , which mainly guides the model to predict the locations/velocities accurately. Therefore, the combined loss enables our model to be more noise-tolerable; 2) Regarding the relationship to T , $\mathcal{L}_{reverse}$ is more sensitive to total sequence length (T^5), thus it provides more regularization for long-context prediction, a key challenge for dynamic modeling.

TRS Loss Design Choice. We define $\mathcal{L}_{reverse}$ as the distance between model forward trajectories and backward trajectories. Based on the definition of TRS in Sec. 2.2, there are other implementation choices. One prior work TRS-ODE (Huh et al., 2020) designed a TRS loss based on Eqn 5, where a reverse trajectory shares the same starting point as the forward one. However, we show that our implementation based on Lemma 2.1 to approximate time-reversal symmetry has a lower maximum error compared to their implementation below, supported by empirical experiments in Sec. 4.2.

Lemma 3.2. *Let $\mathcal{L}_{reverse}$ be the TRS implementation of TREAT based on Lemma 2.1, $\mathcal{L}_{reverse2}$ be the one in (Huh et al., 2020) based on Eqn 5. When the reconstruction loss defined in Eqn 7 of both methods are equal, and the two TRS losses are equal, i.e. $\mathcal{L}_{reverse} = \mathcal{L}_{reverse2}$, the maximum error between the reversal and ground truth trajectory for each agent, i.e. $\text{MaxError}_{gt_rev} = \max_{k \in [K]} \left\| \mathbf{y}_i(t_k) - \hat{\mathbf{y}}_i^{\text{rev}}(t'_{K-k}) \right\|_2$ for $i = 1, 2 \dots N$, made by TREAT is smaller.*

We prove Lemma 3.2 in Appendix A.4. Another implementation is to minimize the distances between model backward trajectories and ground truth trajectories. When both forward and backward trajectories are close to ground-truth, they are implicitly symmetric. The major drawback is that at the early stage of learning when the forward is far away from ground truth (\mathcal{L}_{pred}), such implicit regularization does not force time-reversal symmetry, but introduces more noise.

4 Experiments

Datasets. We conduct systematic evaluations over five multi-agent systems including three 5-body spring systems (Kipf et al., 2018), a complex chaotic pendulum system and a real-world motion

capture dataset (CMU, 2003); and four single-agent systems including three spring systems (with only one node) and a chaotic strange attractors system (Huh et al., 2020).

The settings of spring systems include: 1) conservative, i.e. no interactions with the environments, we call it *Simple Spring*; 2) non-conservative with frictions, we call it *Damped Spring*; 3) non-conservative with periodic external forces, we call it *Forced Spring*. The *Pendulum* system contains three connected sticks in a 2D plane. It is highly sensitive to initial states, with minor disturbances leading to significantly different trajectories (Shinbrot et al., 1992; Awrejcewicz et al., 2008). The real-world motion capture dataset (CMU, 2003) describes the walking trajectories of a person, each tracking a single joint. We call it *Human Motion*. The strange attractor consists of symmetric attractor/repellor force pairs and is chaotic (Sprott and Clinton, 2015). It is also highly sensitive to the initial states (Koppe et al., 2019). We call it *Attractor*.

Towards physical properties, *Simple Spring* and *Pendulum* are conservative and reversible; *Force Spring* and *Attractor* are reversible but non-conservative; *Damped Spring* are irreversible and non-conservative. For *Human Motion*, it does not adhere to specific physical laws since it is a real-world dataset. Details of the datasets and generation pipelines can be found in Appendix C.

Task Setup. We conduct evaluation by splitting trajectories into two halves: $[t_1, t_M]$, $[t_{M+1}, t_K]$ where timestamps can be irregular. We condition the first half of observations to make predictions for the second half as in (Rubanova et al., 2019). For spring datasets and *Pendulum*, we generate irregular-sampled trajectories and set the training samples to be 20,000 and testing samples to be 5,000 respectively. For *Attractor*, We generate 1,000 and 50 trajectories for training and testing respectively following Huh et al. (2020). 10% of training samples are used as validation sets and the maximum trajectory prediction length is 60. Details can be found in Appendix C.

Baselines. We compare TREAT against three baseline types: 1) pure data-driven approaches including LG-ODE (Huang et al., 2020) and LatentODE (Rubanova et al., 2019), where the first one is a multi-agent approach considering pair-wise interactions, and the second one is a single-agent approach that predicts each trajectory independently; 2) energy-preserving HODEN (Greydanus et al., 2019); and 3) time-reversal TRS-ODEN (Huh et al., 2020).

The latter two are single-agent approaches and require initial states as given input. To handle missing initial states in our dataset, we approximate the initial states for the two methods via linear spline interpolation (Endre Süli, 2003). In addition, we substitute the ODE network in TRS-ODEN with a GNN (Kipf et al., 2018) as TRS-ODEN_{GNN}, which serves as a new multi-agent approach for fair comparison. HODEN cannot be easily extended to the multi-agent setting as replacing the ODE function with a GNN can violate energy conservation of the original HODEN. For running LGODE and TREAT on single-agent datasets, we only include self-loop edges in the graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, which makes the ODE function g a simple MLP. Implementation details can be found in Appendix D.2.

Table 1: Evaluation results on MSE (10^{-2}). Best results are in **bold** numbers and second-best results are in underline numbers. *Human Motion* is a real-world dataset and all others are simulated datasets.

Dataset	Multi-Agent Systems					Single-Agent Systems			
	<i>Simple Spring</i>	<i>Forced Spring</i>	<i>Damped Spring</i>	<i>Pendulum</i>	<i>Human Motion</i>	<i>Simple Spring</i>	<i>Forced Spring</i>	<i>Damped Spring</i>	<i>Attractor</i>
LatentODE	5.2622	5.0277	3.3419	2.6894	2.9061	5.7957	0.4563	1.3012	0.58394
HODEN	3.0039	4.0668	8.7950	741.2296	1.9855	3.2119	4.004	1.5675	54.2912
TRS-ODEN	3.6785	4.4465	1.7595	741.4988	0.5400	3.0271	0.4056	1.5667	2.2683
TRS-ODEN _{GNN}	<u>1.4115</u>	<u>2.1102</u>	<u>0.5951</u>	596.0319	<u>0.2609</u>	/	/	/	/
LG-ODE	1.7429	1.8929	0.9718	1.4156	0.7610	1.6156	0.1465	1.1223	0.6942
TREAT	1.1178	1.4525	0.5944	1.2527	0.2192	1.6026	0.0960	1.0750	0.5581
(—Ablation of our method with different implementation of $L_{reverse}$ —)									
TREAT _{$\mathcal{L}_{rev}=gt\text{-}rev$}	1.1313	1.5254	0.6171	1.6158	0.2495	1.6190	0.1104	1.1205	0.6364
TREAT _{$\mathcal{L}_{rev}=rev2$}	1.6786	1.9786	0.9692	1.5631	0.8785	1.6901	0.0983	1.0952	0.7286

4.1 Main Results

Table 1 shows the prediction performance on both multi-agent systems and single-agent systems measured by mean squared error (MSE). We can see that TREAT consistently surpasses other models, highlighting its generalizability and the efficacy of the proposed TRS loss.

For multi-agent systems, approaches that consider interactions among agents (LG-ODE, TRS-ODEN_{GNN}, TREAT) consistently outperform single-agent baselines (LatentODE, HODEN, TRS-ODEN), and TREAT achieves the best performance across datasets.

The chaotic nature of the *Pendulum* system and the *Attractor* system, with their sensitivity to initial states⁴, poses extreme challenges for dynamic modeling. This leads to highly unstable predictions for models like HODEN and TRS-ODEN, as they estimate initial states via inaccurate linear spline interpolation (Endre Süli, 2003). In contrast, LatentODE, LG-ODE, and TREAT employ advanced encoders that infer latent states from observed data and demonstrate superior accuracy. Among them, TREAT achieves the most accurate predictions, further showing its robust generalization capabilities.

We observe that misapplied inductive biases can degrade results, which limits the usage of physics-informed methods that are designed for individual physical prior such as HODEN. HODEN only excels on energy-conservative systems, such as *Simple Spring* compared with LatentODE and TRS-ODEN in the multi-agent setting. Its performance drop dramatically on *Force Spring*, *Damped Spring*, and *Attractor*. Note that HODEN naively forces each agent to be energy-conservative, instead of the whole system. Therefore, it performs poorly than LG-ODE, TREAT in the multi-agent settings.

For the *Human Motion* dataset, characterized by its dynamic ambiguity as it does not adhere to specific physical laws, we cannot directly determine whether it is conservative or time-reversal. For such a system with an unknown nature, TREAT outperforms other purely data-driven methods significantly, showcasing its strong numerical benefits in improving prediction accuracy across diverse system types. This is also shown by its superior performance on *Damped Spring*, which is irreversible.

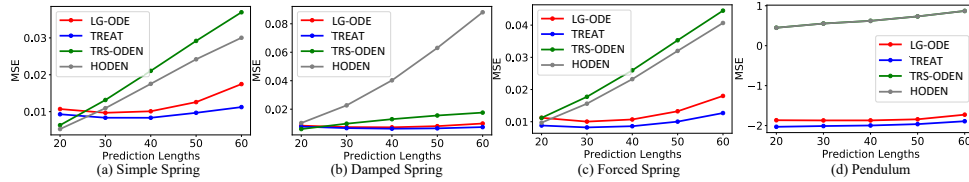


Figure 4: Varying prediction lengths across multi-agent datasets (Pendulum MSE is in log values).

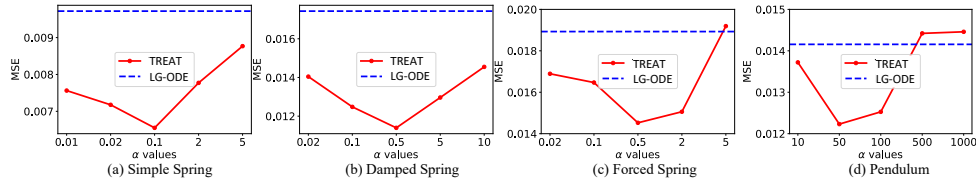


Figure 5: Varying α values across multi-agent datasets.

4.2 Ablation and Sensitivity Analysis

Ablation on implementation of $\mathcal{L}_{reverse}$. We conduct two ablation by changing the implementation of $\mathcal{L}_{reverse}$ discussed in Sec. 3.2: 1) $\text{TREAT}_{\mathcal{L}_{rev}=\text{gt-rev}}$, which computes the reversal loss as the L2 distance between ground truth trajectories to model backward trajectories; 2) $\text{TREAT}_{\mathcal{L}_{rev}=\text{rev2}}$, which implements the TRS loss based on Eqn 5 as in TRS-ODEN (Huh et al., 2020). From the last block of Table 1, we can clearly see that our implementation achieves the best performance against the two.

Evaluation across prediction lengths. We vary the maximum prediction lengths from 20 to 60 and report model performance as shown in Figure 4. As the prediction step increases, TREAT consistently maintains optimal prediction performance, while other baselines exhibit significant error accumulations. The performance gap between TREAT and baselines widens when making long-range predictions, highlighting the superior predictive capability of TREAT.

Evaluation across different α . We vary the values of the coefficient α defined in Eqn 10, which balances the reconstruction loss and the TRS loss. Figure 5 demonstrates that the optimal α values being neither too high nor too low. This is because when α is too small, the model tends to neglect the TRS physical bias, resulting in error accumulations. Conversely, when α becomes too large, the

⁴Video to show *Pendulum* is highly sensitive to initial states.

model can emphasize TRS at the cost of accuracy. Nonetheless, across different α values, TREAT consistently surpasses the purely data-driven LG-ODE, showcasing its superiority and flexibility in modeling diverse dynamical systems.

Finally, we study its sensitivity towards solver choice and observation ratios in Appendix E.

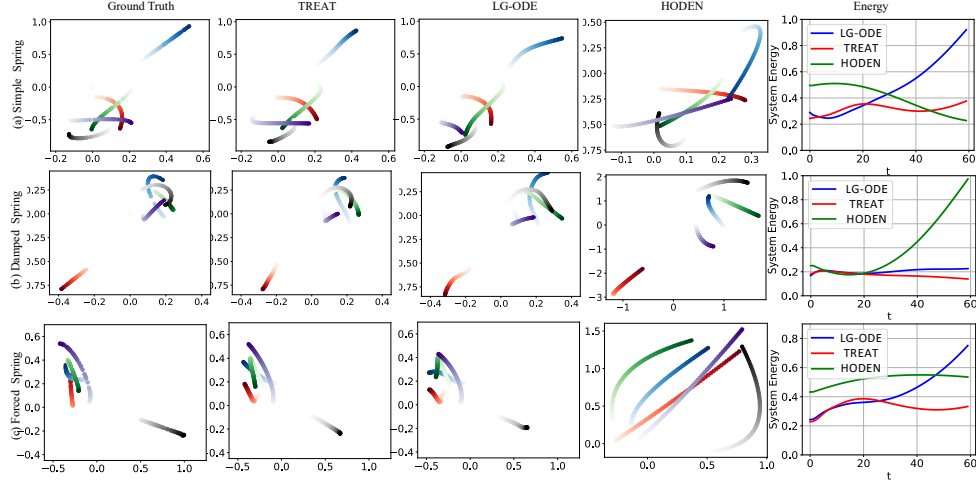


Figure 6: Visualization for 5-body spring systems (trajectory starts from light to dark colors).

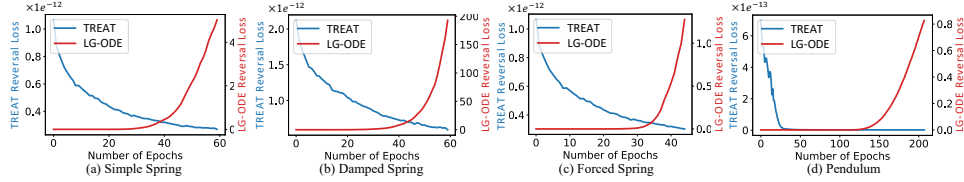


Figure 7: TRS loss visualization across multi-agent datasets (scales of two y-axes are different).

4.3 Visualizations

Trajectory Visualizations. Model predictions and ground truth are visualized in Figure 6. As HODEN is a single-agent baseline that individually forces every agent’s energy to be constant over time which is not valid, the predicted trajectories is having the largest errors and systems’ total energy is not conserved for all datasets. The purely data-driven LG-ODE exhibits unrealistic energy patterns, as seen in the energy spikes in *Simple Spring* and *Force Spring*. In contrast, TREAT, incorporating reversal loss, generates realistic energy trends, and consistently produces trajectories closest to the ground truth, showing its superior performance.

Reversal Loss Visualizations To illustrate the issue of energy explosion from the purely data-driven LG-ODE, we visualize the TRS loss over training epochs from LG-ODE⁵ and TREAT in Figure 7. As results suggest, LG-ODE has increased TRS loss over training epochs, meaning it is violating the time-reversal symmetry sharply, in contrast to TREAT which has decreased reversal loss over epochs.

5 Conclusions

We propose TREAT, a deep learning framework that achieves high-precision modeling for a wide range of dynamical systems by injecting time-reversal symmetry as an inductive bias. TREAT features a novel regularization term to softly enforce time-reversal symmetry by aligning predicted forward and reverse trajectories from a GraphODE model. Notably, we theoretically prove that the regularization term effectively minimizes higher-order Taylor expansion terms during the ODE integration, which serves as a general numerical benefit widely applicable to various systems (even irreversible systems) regardless of their physical properties. Empirical evaluations on different kinds of datasets illustrate TREAT’s superior efficacy in accurately capturing real-world system dynamics.

⁵There is no reversal loss backpropagation in LG-ODE, we just compute its value along training.

References

- J. Awrejcewicz, G. Kudra, and G. Wasilewski. 2008. Chaotic zones in triple pendulum dynamics observed experimentally and numerically. *Applied Mechanics and Materials* (2008), 1–17.
- Y. Cao, X. Gao, and R. Li. 2019. A liquid plug moving in an annular pipe—Heat transfer analysis. *International Journal of Heat and Mass Transfer* 139 (2019), 1065–1076.
- Y. Cao and R. Li. 2018. A liquid plug moving in an annular pipe—Flow analysis. *Physics of Fluids* 30, 9 (2018).
- B. Chang, L. Meng, E. Haber, L. Ruthotto, D. Begert, and E. Holtham. 2018. Reversible architectures for arbitrarily deep residual neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 32.
- R. T. Q. Chen, B. Amos, and M. Nickel. 2021. Learning Neural Event Functions for Ordinary Differential Equations. *International Conference on Learning Representations* (2021).
- T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud. 2018. Neural Ordinary Differential Equations. In *Advances in Neural Information Processing Systems*.
- CMU. 2003. Carnegie-Mellon Motion Capture Database. <http://mocap.cs.cmu.edu>
- M. Cranmer, S. Greydanus, S. Hoyer, P. Battaglia, D. Spergel, and S. Ho. 2020. Lagrangian neural networks. *arXiv preprint arXiv:2003.04630* (2020).
- D. F. M. Endre Süli. 2003. *An Introduction to Numerical Analysis*. Cambridge University Press. 293 pages.
- S. Greydanus, M. Dzamba, and J. Yosinski. 2019. Hamiltonian neural networks. *Advances in Neural Information Processing Systems* (2019).
- Z. Hu, Y. Dong, K. Wang, and Y. Sun. 2020. Heterogeneous Graph Transformer. In *Proceedings of the 2020 World Wide Web Conference*.
- Z. Huang, Y. Sun, and W. Wang. 2020. Learning Continuous System Dynamics from Irregularly-Sampled Partial Observations. In *Advances in Neural Information Processing Systems*.
- Z. Huang, Y. Sun, and W. Wang. 2021. Coupled Graph ODE for Learning Interacting System Dynamics. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*.
- I. Huh, E. Yang, S. J. Hwang, and J. Shin. 2020. Time-Reversal Symmetric ODE Network. In *Advances in Neural Information Processing Systems*.
- S. Jiang, Z. Huang, X. Luo, and Y. Sun. 2023. CF-GODE: Continuous-Time Causal Inference for Multi-Agent Dynamical Systems. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*.
- S. Kim and S. J. Karrila. 2013. *Microhydrodynamics: principles and selected applications*. Courier Corporation.
- T. Kipf, E. Fetaya, K. Wang, M. Welling, and R. Zemel. 2018. Neural Relational Inference for Interacting Systems. *arXiv preprint arXiv:1802.04687* (2018).
- Georgia Koppe, Hazem Toutounji, Peter Kirsch, Stefanie Lis, and Daniel Durstewitz. 2019. Identifying nonlinear dynamical systems via generative recurrent neural networks with applications to fMRI. *PLoS computational biology* 15, 8 (2019), e1007263.
- Remi Lam, Alvaro Sanchez-Gonzalez, Matthew Willson, Peter Wirsberger, Meire Fortunato, Ferran Alet, Suman Ravuri, Timo Ewalds, Zach Eaton-Rosen, Weihua Hu, Alexander Meroze, Stephan Hoyer, George Holland, Oriol Vinyals, Jacklynn Stott, Alexander Pritzel, Shakir Mohamed, and Peter Battaglia. 2023. Learning skillful medium-range global weather forecasting. *Science* 382, 6677 (2023), 1416–1421.

361 J. S. Lamb and J. A. Roberts. 1998. Time-reversal symmetry in dynamical systems: a survey. *Physica*
362 *D: Nonlinear Phenomena* (1998), 1–39.

363 C. Li, F. Xia, R. Martín-Martín, M. Lingelbach, S. Srivastava, B. Shen, K. E. Vainio, C. Gokmen, G.
364 Dharan, T. Jain, A. Kurenkov, K. Liu, H. Gweon, J. Wu, L. Fei-Fei, and S. Savarese. 2022. iGibson
365 2.0: Object-Centric Simulation for Robot Learning of Everyday Household Tasks. In *Proceedings*
366 *of the 5th Conference on Robot Learning*.

367 J. Liu, A. Kumar, J. Ba, J. Kiros, and K. Swersky. 2019. Graph normalizing flows. *Advances in*
368 *Neural Information Processing Systems* 32 (2019).

369 I. Loshchilov and F. Hutter. 2019. Decoupled weight decay regularization. In *The International*
370 *Conference on Learning Representations*.

371 Yupu Lu, Shijie Lin, Guanqi Chen, and Jia Pan. 2022. ModLaNets: Learning Generalisable Dynamics
372 via Modularity and Physical Inductive Bias. In *Proceedings of the 39th International Conference on*
373 *Machine Learning (Proceedings of Machine Learning Research, Vol. 162)*, Kamalika Chaudhuri,
374 Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (Eds.). PMLR, 14384–
375 14397.

376 X. Luo, J. Yuan, Z. Huang, H. Jiang, Y. Qin, W. Ju, M. Zhang, and Y. Sun. 2023. HOPE: High-
377 order Graph ODE For Modeling Interacting Dynamics. In *Proceedings of the 40th International*
378 *Conference on Machine Learning*.

379 Ruiqi Ni and Ahmed H Qureshi. 2022. Ntfields: Neural time fields for physics-informed robot motion
380 planning. *arXiv preprint arXiv:2210.00120* (2022).

381 J. North. 2021. Formulations of classical mechanics. *Forthcoming in A companion to the philosophy*
382 *of physics*. Routledge (2021).

383 T. Pfaff, M. Fortunato, A. Sanchez-Gonzalez, and P. Battaglia. 2021. Learning Mesh-Based Simula-
384 tion with Graph Networks. In *International Conference on Learning Representations*.

385 M. Poli, S. Massaroli, J. Park, A. Yamashita, H. Asama, and J. Park. 2019. Graph neural ordinary
386 differential equations. *arXiv preprint arXiv:1911.07532* (2019).

387 C. Pozrikidis. 2001. Interfacial dynamics for Stokes flow. *J. Comput. Phys.* 169, 2 (2001), 250–301.

388 M. Raissi, P. Perdikaris, and G. E. Karniadakis. 2019. Physics-informed neural networks: A
389 deep learning framework for solving forward and inverse problems involving nonlinear partial
390 differential equations. *Journal of Computational physics* 378 (2019), 686–707.

391 Y. Rubanova, R. T. Chen, and D. K. Duvenaud. 2019. Latent ordinary differential equations for
392 irregularly-sampled time series. In *Advances in Neural Information Processing Systems*.

393 A. Sanchez-Gonzalez, V. Bapst, K. Cranmer, and P. Battaglia. 2019. Hamiltonian Graph Networks
394 with ODE Integrators. In *Advances in Neural Information Processing Systems*.

395 A. Sanchez-Gonzalez, J. Godwin, T. Pfaff, R. Ying, J. Leskovec, and P. W. Battaglia. 2020. Learning
396 to Simulate Complex Physics with Graph Networks. In *Proceedings of the 37th International*
397 *Conference on Machine Learning*.

398 V. G. Satorras, E. Hoogeboom, and M. Welling. 2021. E (n) equivariant graph neural networks. In
399 *International conference on machine learning*. PMLR, 9323–9332.

400 M. Schober, S. Särkkä, and P. Hennig. 2019. A probabilistic model for the numerical solution of
401 initial value problems. In *Statistics and Computing*. 99–122.

402 H. Sepp and S. Jürgen. 1997. Long Short-term Memory. *Neural computation* (1997).

403 T. Shinbrot, C. Grebogi, J. Wisdom, and J. A. Yorke. 1992. Chaos in a double pendulum. *American*
404 *Journal of Physics* 6 (1992), 491–499.

405 Sprott and Julien Clinton. 2015. Symmetric time-reversible flows with a strange attractor. *Interna-*
406 *tional Journal of Bifurcation and Chaos* 25, 05 (2015), 1550078.

- 407 T. Stachowiak and T. Okada. 2006. A numerical analysis of chaos in the double pendulum. *Chaos,*
408 *Solitons & Fractals* 2 (2006), 417–422.
- 409 E. C. Tolman. 1938. The Determiners of Behavior at a Choice Point. *Psychological Review* 45, 1
410 (1938), 1–41.
- 411 R. Valperga, K. Webster, D. Turaev, V. Klein, and J. Lamb. 2022. Learning Reversible Symplectic
412 Dynamics. In *Proceedings of The 4th Annual Learning for Dynamics and Control Conference*.
- 413 A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, ' U. Kaiser, and I. Polosukhin.
414 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*.
- 415 R. Wang, K. Kashinath, M. Mustafa, A. Albert, and R. Yu. 2020. Towards physics-informed deep
416 learning for turbulent flow prediction. In *Proceedings of the 26th ACM SIGKDD International*
417 *Conference on Knowledge Discovery and Data Mining*.
- 418 S. Wen, H. Wang, and D. Metaxas. 2022. Social ODE: Multi-agent Trajectory Forecasting with
419 Neural Ordinary Differential Equations. In *European Conference on Computer Vision*.
- 420 C. Zang and F. Wang. 2020. Neural dynamics on complex networks. In *Proceedings of the 26th ACM*
421 *SIGKDD International Conference on Knowledge Discovery and Data Mining*.

A Theoretical Analysis

A.1 Implementation of the Time-Reversal Symmetry Loss

Algorithm 1 The implementation of $\mathcal{L}_{reverse}$

Require: latent initial states $\mathbf{z}_i^{\text{fwd}}(t_0)$; the ODE function $g(\cdot)$; number of agents N :

- 1: **for** each $i \in N$ **do**
 - 2: Compute the latent forward trajectory at timestamps $\{t_k\}_{k=0}^K$:
 $\mathbf{z}_i^{\text{fwd}}(t_k) = \text{ODE-Solver}(g, [\mathbf{z}_1^{\text{fwd}}(t_0), \mathbf{z}_2^{\text{fwd}}(t_0) \dots \mathbf{z}_N^{\text{fwd}}(t_0)], t_k)$. Reach the final state $\mathbf{z}_i^{\text{fwd}}(t_K)$.
 - 3: The initial state of the reverse trajectory is defined as $\mathbf{z}_i^{\text{rev}}(t'_0) = \mathbf{z}_i^{\text{fwd}}(t_K)$, and the dynamics of the system which is the ODE function $g(\cdot)$ is also reversed as $-g(\cdot)$.
 - 4: Compute the latent reverse trajectory at timestamps $\{t'_k\}_{k=0}^K$,
 $\mathbf{z}_i^{\text{rev}}(t'_k) = \text{ODE-Solver}(g, [\mathbf{z}_1^{\text{rev}}(t'_0), \mathbf{z}_2^{\text{rev}}(t'_0) \dots \mathbf{z}_N^{\text{rev}}(t'_0)], t'_k)$.
 - 5: $\hat{\mathbf{y}}_i^{\text{fwd}}(t_k) = f_{\text{DEC}}(\mathbf{z}_i^{\text{fwd}}(t_k))$, $\hat{\mathbf{y}}_i^{\text{rev}}(t'_k) = f_{\text{DEC}}(\mathbf{z}_i^{\text{rev}}(t'_k))$
 - 6: **end for**
 - 7: $\mathcal{L}_{reverse} = \sum_{i=1}^N \sum_{k=0}^K \left\| \hat{\mathbf{y}}_i^{\text{fwd}}(t_k) - \hat{\mathbf{y}}_i^{\text{rev}}(t'_{K-k}) \right\|_2^2$
-

A.2 Proof of Lemma 1

Proof. The definition of time-reversal symmetry is given by:

$$R \circ \phi_t = \phi_{-t} \circ R = \phi_t^{-1} \circ R \quad (11)$$

Here, R is an involution operator, which means $R \circ R = \text{I}$.

First, we apply the time evolution operator ϕ_t to both sides of Eqn 11:

$$\phi_t \circ R \circ \phi_t = \phi_t \circ \phi_t^{-1} \circ R \quad (12)$$

Simplifying, we obtain:

$$\phi_t \circ R \circ \phi_t = R \quad (13)$$

Next, we apply the involution operator R to both sides of the equation:

$$R \circ \phi_t \circ R \circ \phi_t = R \circ R \quad (14)$$

Since $R \circ R = \text{I}$, we finally arrive at:

$$R \circ \phi_t \circ R \circ \phi_t = \text{I} \quad (15)$$

which means the trajectories can overlap when evolving backward from the final state. \square

A.3 Proof of Theorem 3.1

Let Δt denote the integration step size in an ODE solver and T be the prediction length. The time stamps of the ODE solver are $\{t_j\}_{j=0}^T$, where $t_{j+1} - t_j = \Delta t$ for $j = 0, \dots, T(T > 1)$. Next suppose during the forward evolution, the updates go through states $\mathbf{z}^{\text{fwd}}(t_j) = (\mathbf{q}^{\text{fwd}}(t_j), \mathbf{p}^{\text{fwd}}(t_j))$ for $j = 0, \dots, T$, where $\mathbf{q}^{\text{fwd}}(t_j)$ is position, $\mathbf{p}^{\text{fwd}}(t_j)$ is momentum, while during the reverse evolution they go through states $\mathbf{z}^{\text{rev}}(t_j) = (\mathbf{q}^{\text{rev}}(t_j), \mathbf{p}^{\text{rev}}(t_j))$ for $j = 0, \dots, T$, in reverse order. The ground truth trajectory is $\mathbf{z}^{\text{gt}}(t_j) = (\mathbf{q}^{\text{gt}}(t_j), \mathbf{p}^{\text{gt}}(t_j))$ for $j = 0, \dots, T$.

For the sake of brevity in the ensuing proof, we denote $\mathbf{z}^{\text{gt}}(t_j)$ by \mathbf{z}_j^{gt} , $\mathbf{z}^{\text{fwd}}(t_j)$ by $\mathbf{z}_j^{\text{fwd}}$ and $\mathbf{z}^{\text{rev}}(t_j)$ by $\mathbf{z}_j^{\text{rev}}$, and we will use Mathematical Induction to prove the theorem.

A.3.1 Reconstruction Loss (\mathcal{L}_{pred}) Analysis.

First, we bound the forward loss $\sum_{j=0}^T \|\mathbf{z}_j^{\text{fwd}} - \mathbf{z}_j^{\text{gt}}\|_2^2$. Since our method models the momentum and position of the system, we can write the following Taylor expansion of the forward process, where

for any $0 \leq j < T$:

$$\begin{cases} \mathbf{q}_{j+1}^{\text{fwd}} = \mathbf{q}_j^{\text{fwd}} + (\mathbf{p}_j^{\text{fwd}}/m)\Delta t + (\dot{\mathbf{p}}_j^{\text{fwd}}/2m)\Delta t^2 + \mathcal{O}(\Delta t^3), \\ \mathbf{p}_{j+1}^{\text{fwd}} = \mathbf{p}_j^{\text{fwd}} + \dot{\mathbf{p}}_j^{\text{fwd}}\Delta t + \mathcal{O}(\Delta t^2), \\ \dot{\mathbf{p}}_{j+1}^{\text{fwd}} = \dot{\mathbf{p}}_j^{\text{fwd}} + \mathcal{O}(\Delta t), \end{cases} \quad \begin{aligned} (16a) \\ (16b) \\ (16c) \end{aligned}$$

and for the ground truth process, we also have from Taylor expansion that

$$\begin{cases} \mathbf{q}_{j+1}^{\text{gt}} = \mathbf{q}_j^{\text{gt}} + (\mathbf{p}_j^{\text{gt}}/m)\Delta t + (\dot{\mathbf{p}}_j^{\text{gt}}/2m)\Delta t^2 + \mathcal{O}(\Delta t^3), \\ \mathbf{p}_{j+1}^{\text{gt}} = \mathbf{p}_j^{\text{gt}} + \dot{\mathbf{p}}_j^{\text{gt}}\Delta t + \mathcal{O}(\Delta t^2), \\ \dot{\mathbf{p}}_{j+1}^{\text{gt}} = \dot{\mathbf{p}}_j^{\text{gt}} + \mathcal{O}(\Delta t). \end{cases} \quad \begin{aligned} (17a) \\ (17b) \\ (17c) \end{aligned}$$

With these, we aim to prove that for any $k = 0, 1, \dots, T$, the following hold :

$$\begin{cases} \|\mathbf{q}_k^{\text{fwd}} - \mathbf{q}_k^{\text{gt}}\|_2 \leq C_2^{\text{fwd}} k^2 \Delta t^2, \\ \|\mathbf{p}_k^{\text{fwd}} - \mathbf{p}_k^{\text{gt}}\|_2 \leq C_1^{\text{fwd}} k \Delta t, \end{cases} \quad \begin{aligned} (18a) \\ (18b) \end{aligned}$$

442 where C_1^{fwd} and C_2^{fwd} are constants.

443 **Base Case $k = 0$:** Based on the initialization rules, it is obvious that $\|\mathbf{q}_0^{\text{fwd}} - \mathbf{q}_0^{\text{gt}}\|_2 = 0$ and
444 $\|\mathbf{p}_0^{\text{fwd}} - \mathbf{p}_0^{\text{gt}}\|_2 = 0$, thus (18a) and (18b) both hold for $k = 0$.

Inductive Hypothesis: Assume (18a) and (18b) hold for $k = j$, which means:

$$\begin{cases} \|\mathbf{q}_j^{\text{fwd}} - \mathbf{q}_j^{\text{gt}}\|_2 \leq C_2^{\text{fwd}} j^2 \Delta t^2, \\ \|\mathbf{p}_j^{\text{fwd}} - \mathbf{p}_j^{\text{gt}}\|_2 \leq C_1^{\text{fwd}} j \Delta t, \end{cases} \quad \begin{aligned} (19a) \\ (19b) \end{aligned}$$

445 **Inductive Proof:** We need to prove (18a) and (18b) hold for $k = j + 1$.

446 First, using (16c) and (17c), we have

$$\|\dot{\mathbf{p}}_{j+1}^{\text{fwd}} - \dot{\mathbf{p}}_{j+1}^{\text{gt}}\|_2 = \|\dot{\mathbf{p}}_j^{\text{fwd}} - \dot{\mathbf{p}}_j^{\text{gt}}\|_2 + \mathcal{O}(\Delta t) = \|\dot{\mathbf{p}}_0^{\text{fwd}} - \dot{\mathbf{p}}_0^{\text{gt}}\|_2 + \mathcal{O}((j+1)\Delta t) = \mathcal{O}(1), \quad (20)$$

447 where we iterate through $j, j-1, \dots, 0$ in the second equality. Then using (17b) and (16b), we get
448 for $j+1$ that

$$\begin{aligned} \|\mathbf{p}_{j+1}^{\text{fwd}} - \mathbf{p}_{j+1}^{\text{gt}}\|_2 &= \|(\mathbf{p}_j^{\text{fwd}} + \dot{\mathbf{p}}_j^{\text{fwd}}\Delta t) - (\mathbf{p}_j^{\text{gt}} + \dot{\mathbf{p}}_j^{\text{gt}}\Delta t) + \mathcal{O}(\Delta t^2)\|_2 \\ &\leq \|\mathbf{p}_j^{\text{fwd}} - \mathbf{p}_j^{\text{gt}}\|_2 + \|\dot{\mathbf{p}}_j^{\text{fwd}} - \dot{\mathbf{p}}_j^{\text{gt}}\|_2 \Delta t + \mathcal{O}(\Delta t^2) \\ &\leq [C_1^{\text{fwd}} j + \mathcal{O}(1)] \Delta t, \end{aligned}$$

449 where the first inequality uses the triangle inequality, and in the second inequality we use (19b) as
450 well as (20). We can see there exists C_1^{fwd} such that the final expression above is upper bounded by
451 $C_1^{\text{fwd}}(j+1)\Delta t$, with which the claim holds for $j+1$.

452 Next for (18a), using (17a) and (16a), we get for any j that

$$\begin{aligned} \|\mathbf{q}_{j+1}^{\text{fwd}} - \mathbf{q}_{j+1}^{\text{gt}}\|_2 &= \|(\mathbf{q}_j^{\text{fwd}} + (\mathbf{p}_j^{\text{fwd}}/m)\Delta t + (\dot{\mathbf{p}}_j^{\text{fwd}}/2m)\Delta t^2) - (\mathbf{q}_j^{\text{gt}} + (\mathbf{p}_j^{\text{gt}}/m)\Delta t + (\dot{\mathbf{p}}_j^{\text{gt}}/2m)\Delta t^2) + \mathcal{O}(\Delta t^3)\|_2 \\ &\leq \|\mathbf{q}_j^{\text{fwd}} - \mathbf{q}_j^{\text{gt}}\|_2 + \frac{1}{m} \|\mathbf{p}_j^{\text{fwd}} - \mathbf{p}_j^{\text{gt}}\|_2 \Delta t + \frac{1}{2m} \|\dot{\mathbf{p}}_j^{\text{fwd}} - \dot{\mathbf{p}}_j^{\text{gt}}\|_2 \Delta t^2 + \mathcal{O}(\Delta t^3) \\ &\leq \left[C_2^{\text{fwd}} j^2 + \frac{C_1^{\text{fwd}}}{m} j + \mathcal{O}(1) \right] \Delta t^2, \end{aligned}$$

453 where the first inequality uses the triangle inequality, and in the second inequality we use (19a) and
454 (19b) as well as (20). Thus with an appropriate C_2^{fwd} , we have the final expression above is upper
455 bounded by $C_2^{\text{fwd}}(j+1)^2 \Delta t^2$, and so the claim holds for $j+1$.

456 Since both the base case and the inductive step have been proven, by the principle of mathematical
457 induction, (18a) and (18b) holds for all $k = 0, 1, \dots, T$.

458 With this, we finish the forward proof by plugging (18a) and (18b) into the loss function:

$$\begin{aligned}
\sum_{j=0}^T \|\mathbf{z}_j^{\text{fwd}} - \mathbf{z}_j^{\text{gt}}\|_2^2 &= \sum_{j=0}^T \|\mathbf{p}_j^{\text{fwd}} - \mathbf{p}_j^{\text{gt}}\|_2^2 + \sum_{j=0}^T \|\mathbf{q}_j^{\text{fwd}} - \mathbf{q}_j^{\text{gt}}\|_2^2 \\
&\leq (C_1^{\text{fwd}})^2 \sum_{j=0}^T j^2 \Delta t^2 + (C_2^{\text{fwd}})^2 \sum_{j=0}^T j^4 \Delta t^4 \\
&= \mathcal{O}(T^3 \Delta t^2).
\end{aligned}$$

459 A.3.2 Reversal Loss ($\mathcal{L}_{\text{reverse}}$) Analysis.

460 Next we analyze the reversal loss $\sum_{j=0}^T \|R(\mathbf{z}_j^{\text{rev}}) - \mathbf{z}_j^{\text{fwd}}\|_2^2$. For this, we need to refine the Taylor
 461 expansion residual terms for a more in-depth analysis.

First reconsider the forward process. Since the process is generated from the learned network, we may assume that for some constants c_1 , c_2 , and c_3 , the states satisfy the following for any $0 \leq j < T$:

$$\begin{cases} \mathbf{q}_j^{\text{fwd}} = \mathbf{q}_{j+1}^{\text{fwd}} - (\mathbf{p}_{j+1}^{\text{fwd}}/m)\Delta t + (\dot{\mathbf{p}}_{j+1}^{\text{fwd}}/2m)\Delta t^2 + \mathbf{rem}_j^{\text{fwd},3}, \end{cases} \quad (21a)$$

$$\begin{cases} \mathbf{p}_j^{\text{fwd}} = \mathbf{p}_{j+1}^{\text{fwd}} - \dot{\mathbf{p}}_{j+1}^{\text{fwd}}\Delta t + \mathbf{rem}_j^{\text{fwd},2}, \end{cases} \quad (21b)$$

$$\begin{cases} \dot{\mathbf{p}}_j^{\text{fwd}} = \dot{\mathbf{p}}_{j+1}^{\text{fwd}} + \mathbf{rem}_j^{\text{fwd},1}, \end{cases} \quad (21c)$$

where the remaining terms $\|\mathbf{rem}_j^{\text{fwd},i}\|_2 \leq c_i \Delta t^i$ for $i = 1, 2, 3$. Similarly, we have approximate Taylor expansions for the reverse process:

$$\begin{cases} \mathbf{q}_j^{\text{rev}} = \mathbf{q}_{j+1}^{\text{rev}} + (\mathbf{p}_{j+1}^{\text{rev}}/m)\Delta t + (\dot{\mathbf{p}}_{j+1}^{\text{rev}}/2m)\Delta t^2 + \mathbf{rem}_j^{\text{rev},3}, \end{cases} \quad (22a)$$

$$\begin{cases} \mathbf{p}_j^{\text{rev}} = \mathbf{p}_{j+1}^{\text{rev}} + \dot{\mathbf{p}}_{j+1}^{\text{rev}}\Delta t + \mathbf{rem}_j^{\text{rev},2}, \end{cases} \quad (22b)$$

$$\begin{cases} \dot{\mathbf{p}}_j^{\text{rev}} = \dot{\mathbf{p}}_{j+1}^{\text{rev}} + \mathbf{rem}_j^{\text{rev},1}, \end{cases} \quad (22c)$$

462 where $\|\mathbf{rem}_j^{\text{rev},i}\|_2 \leq c_i \Delta t^i$ for $i = 1, 2, 3$.

We will prove via induction that for $k = T, T-1, \dots, 0$,

$$\begin{cases} \|R(\mathbf{q}_k^{\text{rev}}) - \mathbf{q}_k^{\text{fwd}}\|_2 \leq C_3^{\text{rev}}(T-k)^3 \Delta t^3, \end{cases} \quad (23a)$$

$$\begin{cases} \|R(\mathbf{p}_k^{\text{rev}}) - \mathbf{p}_k^{\text{fwd}}\|_2 \leq C_2^{\text{rev}}(T-k)^2 \Delta t^2, \end{cases} \quad (23b)$$

$$\begin{cases} \|R(\dot{\mathbf{p}}_k^{\text{rev}}) - \dot{\mathbf{p}}_k^{\text{fwd}}\|_2 \leq C_1^{\text{rev}}(T-k)\Delta t, \end{cases} \quad (23c)$$

463 where C_1^{rev} , C_2^{rev} and C_3^{rev} are constants.

464 The entire proof process is analogous to the previous analysis of Reconstruction Loss.

465 **Base Case $k = T$:** Since the reverse process is initialized by the forward process variables at $k = T$,
 466 it is obvious that $\|\mathbf{q}_T^{\text{fwd}} - \mathbf{q}_T^{\text{ev}}\|_2 = \|\mathbf{p}_T^{\text{fwd}} - \mathbf{p}_T^{\text{rev}}\|_2 = \|\dot{\mathbf{p}}_T^{\text{fwd}} - \dot{\mathbf{p}}_T^{\text{rev}}\|_2 = 0$. Thus (23a), (23b) and
 467 (23c) all hold for $k = 0$.

Inductive Hypothesis: Assume the inequalities (23b), (23a) and (23c) hold for $k = j+1$, which means:

$$\begin{cases} \|R(\mathbf{q}_{j+1}^{\text{rev}}) - \mathbf{q}_{j+1}^{\text{fwd}}\|_2 \leq C_3^{\text{rev}}(T-(j+1))^3 \Delta t^3, \end{cases} \quad (24a)$$

$$\begin{cases} \|R(\mathbf{p}_{j+1}^{\text{rev}}) - \mathbf{p}_{j+1}^{\text{fwd}}\|_2 \leq C_2^{\text{rev}}(T-(j+1))^2 \Delta t^2, \end{cases} \quad (24b)$$

$$\begin{cases} \|R(\dot{\mathbf{p}}_{j+1}^{\text{rev}}) - \dot{\mathbf{p}}_{j+1}^{\text{fwd}}\|_2 \leq C_1^{\text{rev}}(T-(j+1))\Delta t, \end{cases} \quad (24c)$$

468 **Inductive Proof:** We need to prove (23b) (23a) and (23c) holds for $k = j$.

469 First, for (23c), using (21c) and (22c), we get for any j that

$$\begin{aligned}
&\|R(\dot{\mathbf{p}}_j^{\text{rev}}) - \dot{\mathbf{p}}_j^{\text{fwd}}\|_2 \\
&= \|(\dot{\mathbf{p}}_{j+1}^{\text{rev}} + \mathbf{rem}_j^{\text{rev},1}) - (\dot{\mathbf{p}}_{j+1}^{\text{fwd}} + \mathbf{rem}_j^{\text{fwd},1})\|_2 \\
&\leq \|R(\dot{\mathbf{p}}_{j+1}^{\text{rev}}) - \dot{\mathbf{p}}_{j+1}^{\text{fwd}}\|_2 + \|\mathbf{rem}_j^{\text{rev},1}\|_2 + \|\mathbf{rem}_j^{\text{fwd},1}\|_2 \\
&\leq C_1^{\text{rev}}(T-j-1)\Delta t + 2c_1\Delta t,
\end{aligned}$$

where the first inequality uses the triangle inequality, and the second inequality plugs in (24c). Thus taking $C_1^{\text{rev}} = 2c_1$, the above is upper bounded by $C_1^{\text{rev}}(T-j)\Delta t$, and (23b) holds for j .

Second, for (24b), using (21b) and (22b), we get

$$\begin{aligned} \|R(\mathbf{p}_j^{\text{rev}}) - \mathbf{p}_j^{\text{fwd}}\|_2 &= \|(\mathbf{p}_{j+1}^{\text{rev}} + \dot{\mathbf{p}}_{j+1}^{\text{rev}}\Delta t + \mathbf{rem}_j^{\text{rev},2}) - (\mathbf{p}_{j+1}^{\text{fwd}} - \dot{\mathbf{p}}_{j+1}^{\text{fwd}}\Delta t + \mathbf{rem}_j^{\text{fwd},2})\|_2 \\ &\leq \|R(\mathbf{p}_{j+1}^{\text{rev}}) - \mathbf{p}_{j+1}^{\text{fwd}}\|_2 + \|\dot{\mathbf{p}}_{j+1}^{\text{rev}} - \dot{\mathbf{p}}_{j+1}^{\text{fwd}}\|_2\Delta t + \|\mathbf{rem}_j^{\text{rev},2}\|_2 + \|\mathbf{rem}_j^{\text{fwd},2}\|_2 \\ &\leq [C_2^{\text{rev}}(T-j-1)^2 + C_1^{\text{rev}}(T-j-1) + 2c_2]\Delta t^2, \end{aligned}$$

where the first inequality uses the triangle inequality, and in the second inequality we use (24a) and (24b). Thus taking $C_2^{\text{rev}} = \max\{C_1^{\text{rev}}/2, 2c_2\}$, we have the final expression above is upper bounded by $C_2^{\text{rev}}(T-j)^2\Delta t^2$, and so the claim holds for j .

Finally, for (24a), we use (21a) and (22a) to get

$$\begin{aligned} &\|R(\mathbf{q}_j^{\text{rev}}) - \mathbf{q}_j^{\text{fwd}}\|_2 \\ &= \|(\mathbf{q}_{j+1}^{\text{rev}} + (\dot{\mathbf{p}}_{j+1}^{\text{rev}}/m)\Delta t + (\dot{\mathbf{p}}_{j+1}^{\text{rev}}/2m)\Delta t^2 + \mathbf{rem}_j^{\text{rev},3}) - (\mathbf{q}_{j+1}^{\text{fwd}} - (\dot{\mathbf{p}}_{j+1}^{\text{fwd}}/m)\Delta t + (\dot{\mathbf{p}}_{j+1}^{\text{fwd}}/2m)\Delta t^2 + \mathbf{rem}_j^{\text{fwd},3})\|_2 \\ &\leq \|R(\mathbf{q}_{j+1}^{\text{rev}}) - \mathbf{q}_{j+1}^{\text{fwd}}\|_2 + \frac{1}{m}\|R(\dot{\mathbf{p}}_{j+1}^{\text{rev}}) - \dot{\mathbf{p}}_{j+1}^{\text{fwd}}\|_2\Delta t + \frac{1}{2m}\|R(\dot{\mathbf{p}}_{j+1}^{\text{rev}}) - \dot{\mathbf{p}}_{j+1}^{\text{fwd}}\|_2\Delta t^2 + \|\mathbf{rem}_j^{\text{rev},3}\|_2 + \|\mathbf{rem}_j^{\text{fwd},3}\|_2 \\ &\leq \left[C_3^{\text{rev}}(T-j-1)^3 + \frac{C_2^{\text{rev}}}{m}(T-j-1)^2 + \frac{C_1^{\text{rev}}}{2m}(T-j-1) + 2c_3 \right] \Delta t^3, \end{aligned}$$

where the first inequality uses the triangle inequality, and in the second inequality we use (24a), (24b) and (24c). Thus taking $C_3^{\text{rev}} = \max\{C_2^{\text{rev}}/3m, C_1^{\text{rev}}/6m, 2c_3\}$, we have the final expression above is upper bounded by $C_3^{\text{rev}}(T-j)^3\Delta t^3$, and so the claim holds for j .

Since both the base case and the inductive step have been proven, by the principle of mathematical induction, (23b), (23a) and (23c) hold for all $k = T, T-1, \dots, 0$.

With this we finish the proof by plugging (23b) and (23a) into the loss function:

$$\begin{aligned} \sum_{j=0}^T \|R(\mathbf{z}_j^{\text{rev}}) - \mathbf{z}_j^{\text{fwd}}\|_2^2 &= \sum_{j=0}^T \|R(\mathbf{p}_j^{\text{rev}}) - \mathbf{p}_j^{\text{fwd}}\|_2^2 + \sum_{j=0}^T \|R(\mathbf{q}_j^{\text{rev}}) - \mathbf{q}_j^{\text{fwd}}\|_2^2 \\ &\leq (C_2^{\text{rev}})^2 \sum_{j=0}^T (T-j)^4 \Delta t^4 + (C_3^{\text{rev}})^2 \sum_{j=0}^T (T-j)^6 \Delta t^6 \\ &= \mathcal{O}(T^5 \Delta t^4). \end{aligned} \tag{25}$$

A.4 Proof of Lemma 3.2

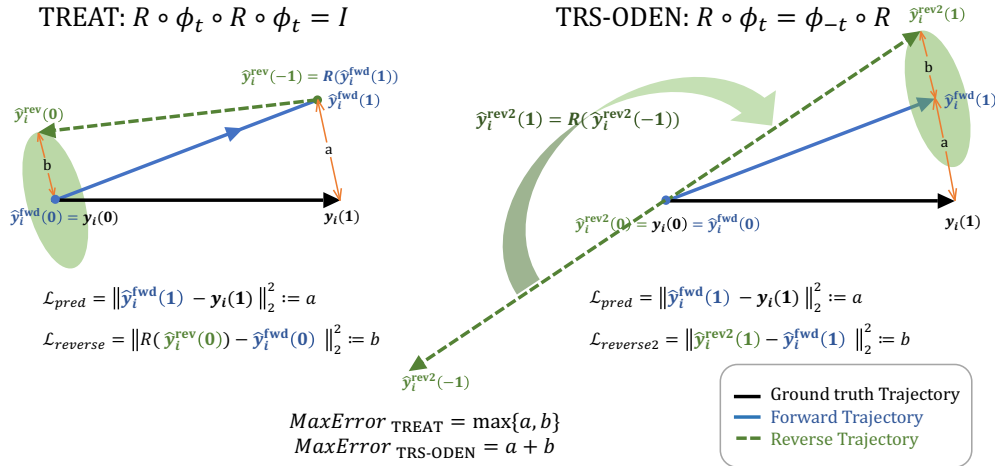


Figure 8: Comparison between two reversal loss implementation

484 We expect an ideal model to align both the predicted forward and reverse trajectories with the ground
 485 truth. As shown in Figure 8, we integrate one step from the initial state $\hat{\mathbf{y}}_i^{\text{fwd}}(0)$ (which is the same as
 486 $\mathbf{y}_i(0)$) and reach the state $\hat{\mathbf{y}}_i^{\text{fwd}}(1)$.

487 The first reverse loss implementation (ours) follows Lemma 2.1 as $R \circ \Phi_t \circ R \circ \Phi_t = \mathbf{I}$, which means
 488 when we evolve forward and reach the state $\hat{\mathbf{y}}_i^{\text{fwd}}(1)$ we reverse it into $\hat{\mathbf{y}}_i^{\text{rev}}(-1) = R(\hat{\mathbf{y}}_i^{\text{fwd}}(1))$ and go
 489 back to reach $\hat{\mathbf{y}}_i^{\text{rev}}(0)$, then reverse it to get $R(\hat{\mathbf{y}}_i^{\text{rev}}(0))$, which ideally should be the same as $\hat{\mathbf{y}}_i^{\text{fwd}}(0)$.

490 The second reverse loss implementation follows Eqn 5as $R \circ \Phi_t = \Phi_{-t} \circ R$, which means we first
 491 reverse the initial state as $\hat{\mathbf{y}}_i^{\text{rev2}}(0) = R(\mathbf{y}_i(0))$, then evolve the reverse trajectory in the opposite
 492 direction to reach $\hat{\mathbf{y}}_i^{\text{rev2}}(-1)$, and then perform a symmetric operation to reach $\hat{\mathbf{y}}_i^{\text{rev2}}(1)$, aligning it
 493 with the forward trajectory.

494 We assume the two reconstruction losses $\mathcal{L}_{\text{pred}} = \|\hat{\mathbf{y}}_i^{\text{fwd}}(1) - \mathbf{y}_i(1)\|_2^2 := a$ are the same. For the
 495 time-reversal losses, we also assume they have reached the same value b :

$$\begin{aligned}\mathcal{L}_{\text{reverse}} &= \|R(\hat{\mathbf{y}}_i^{\text{rev}}(0)) - \hat{\mathbf{y}}_i^{\text{fwd}}(0)\|_2^2 + \|R(\hat{\mathbf{y}}_i^{\text{rev}}(-1)) - \hat{\mathbf{y}}_i^{\text{fwd}}(1)\|_2^2 = \|R(\hat{\mathbf{y}}_i^{\text{rev}}(0)) - \hat{\mathbf{y}}_i^{\text{fwd}}(0)\|_2^2 := b, \\ \mathcal{L}_{\text{reverse2}} &= \|\hat{\mathbf{y}}_i^{\text{rev2}}(0) - \hat{\mathbf{y}}_i^{\text{fwd}}(0)\|_2^2 + \|\hat{\mathbf{y}}_i^{\text{rev2}}(1) - \hat{\mathbf{y}}_i^{\text{fwd}}(1)\|_2^2 = \|\hat{\mathbf{y}}_i^{\text{rev2}}(1) - \hat{\mathbf{y}}_i^{\text{fwd}}(1)\|_2^2 := b,\end{aligned}$$

496 As shown in Figure 8 where we illustrate the worst case scenario $\text{MaxError}_{\text{gt_rev}} =$
 497 $\max_{k \in [K]} \|\mathbf{y}_i(t_k) - \hat{\mathbf{y}}_i^{\text{rev}}(t'_{K-k})\|_2$ of TREAT and TRS-ODEN, we can see that in our implementation
 498 the worst error is the maximum of two loss, while the TRS-ODEN's implementation has the risk of
 499 accumulating the error together, making the worst error being the sum of both:

$$\begin{aligned}\text{MaxError}_{\text{TREAT}} &= \max \{ \|R(\hat{\mathbf{y}}_i^{\text{rev}}(0)) - \mathbf{y}_i(0)\|_2, \|R(\hat{\mathbf{y}}_i^{\text{rev}}(-1)) - \mathbf{y}_i(1)\|_2 \} = \max\{a, b\}, \\ \text{MaxError}_{\text{TRS-ODEN}} &= \max \{ \|\hat{\mathbf{y}}_i^{\text{rev2}}(0) - \mathbf{y}_i(0)\|_2, \|\hat{\mathbf{y}}_i^{\text{rev2}}(1) - \mathbf{y}_i(1)\|_2 \} \\ &= \max \{ 0, \|R(\hat{\mathbf{y}}_i^{\text{rev}}(-1)) - \mathbf{y}_i(1)\|_2 \} \\ &= \|\hat{\mathbf{y}}_i^{\text{rev2}}(1) - \hat{\mathbf{y}}_i^{\text{fwd}}(1)\|_2 + \|\hat{\mathbf{y}}_i^{\text{fwd}}(1) - \mathbf{y}_i(1)\|_2 = a + b,\end{aligned}\tag{26}$$

500 So it is obvious that $\text{MaxError}_{\text{TREAT}}$ made by TREAT is smaller., which means our model achieves
 501 a smaller error of the maximum distance between the reversal and ground truth trajectory.

502 B Example of varying dynamical systems

503 We illustrate the energy conservation and time reversal of the three n-body spring systems used in our
 504 experiments. We use the Hamiltonian formalism of systems under classical mechanics to describe
 505 their dynamics and verify their energy conservation and time-reversibility characteristics.

506 The scalar function that describes a system's motion is called the Hamiltonian, \mathcal{H} , and is typically
 507 equal to the total energy of the system, that is, the potential energy plus the kinetic energy (North,
 508 2021). It describes the phase space equations of motion by following two first-order ODEs called
 509 Hamilton's equations:

$$\frac{d\mathbf{q}}{dt} = \frac{\partial \mathcal{H}(\mathbf{q}, \mathbf{p})}{\partial \mathbf{p}}, \quad \frac{d\mathbf{p}}{dt} = -\frac{\partial \mathcal{H}(\mathbf{q}, \mathbf{p})}{\partial \mathbf{q}},\tag{27}$$

510 where $\mathbf{q} \in \mathbb{R}^n$, $\mathbf{p} \in \mathbb{R}^n$, and $\mathcal{H} : \mathbb{R}^{2n} \mapsto \mathbb{R}$ are positions, momenta, and Hamiltonian of the system.

511 Under this formalism, energy conservative is defined by $d\mathcal{H}/dt = 0$, and the time-reversal symmetry
 512 is defined by $\mathcal{H}(q, p, t) = \mathcal{H}(q, -p, -t)$ (Lamb and Roberts, 1998).

513 B.1 Conservative and reversible systems.

514 A simple example is the isolated n-body spring system, which can be described by :

$$\begin{aligned}\frac{d\mathbf{q}_i}{dt} &= \frac{\mathbf{p}_i}{m} \\ \frac{d\mathbf{p}_i}{dt} &= \sum_{j \in N_i} -k(\mathbf{q}_i - \mathbf{q}_j),\end{aligned}\tag{28}$$

515 where $\mathbf{q} = (\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_N)$ is a set of positions of each object, $\mathbf{p} = (\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_N)$ is a set of
 516 momenta of each object, m_i is mass of each object, k is spring constant.

517 The Hamilton's equations are:

$$\begin{aligned}\frac{\partial \mathcal{H}(\mathbf{q}, \mathbf{p})}{\partial \mathbf{p}_i} &= \frac{d\mathbf{q}_i}{dt} = \frac{\mathbf{p}_i}{m} \\ \frac{\partial \mathcal{H}(\mathbf{q}, \mathbf{p})}{\partial \mathbf{q}_i} &= -\frac{d\mathbf{p}_i}{dt} = \sum_{j \in N_i} k(\mathbf{q}_i - \mathbf{q}_j),\end{aligned}\quad (29)$$

518 Hence, we can obtain the Hamiltonian through the integration of the above equation.

$$\mathcal{H}(\mathbf{q}, \mathbf{p}) = \sum_{i=1}^N \frac{\mathbf{p}_i^2}{2m_i} + \frac{1}{2} \sum_{i=1}^N \sum_{j \in N_i} \frac{1}{2} k(\mathbf{q}_i - \mathbf{q}_j)^2, \quad (30)$$

519 **Verify the systems' energy conservation**

$$\frac{d\mathcal{H}(\mathbf{q}, \mathbf{p})}{dt} = \frac{1}{dt} \left(\sum_{i=1}^N \frac{\mathbf{p}_i^2}{2m_i} \right) + \frac{1}{dt} \left(\frac{1}{2} \sum_{i=1}^N \sum_{j \in N_i} \frac{1}{2} k(\mathbf{q}_i - \mathbf{q}_j)^2 \right) = 0, \quad (31)$$

520 So it is conservative.

521 **Verify the systems' time-reversal symmetry** We do the transformation $R : (\mathbf{q}, \mathbf{p}, t) \mapsto (\mathbf{q}, -\mathbf{p}, -t)$.

$$\begin{aligned}\mathcal{H}(\mathbf{q}, \mathbf{p}) &= \sum_{i=1}^N \frac{\mathbf{p}_i^2}{2m_i} + \frac{1}{2} \sum_{i=1}^N \sum_{j \in N_i} \frac{1}{2} k(\mathbf{q}_i - \mathbf{q}_j)^2, \\ \mathcal{H}(\mathbf{q}, -\mathbf{p}) &= \sum_{i=1}^N \frac{(-\mathbf{p}_i)^2}{2m_i} + \frac{1}{2} \sum_{i=1}^N \sum_{j \in N_i} \frac{1}{2} k(\mathbf{q}_i - \mathbf{q}_j)^2,\end{aligned}\quad (32)$$

522 It is obvious $\mathcal{H}(\mathbf{q}, \mathbf{p}) = \mathcal{H}(\mathbf{q}, -\mathbf{p})$, so it is reversible

523 **B.2 Non-conservative and reversible systems.**

524 A simple example is a n-body spring system with periodical external force, which can be described
525 by:

$$\begin{aligned}\frac{d\mathbf{q}_i}{dt} &= \frac{\mathbf{p}_i}{m} \\ \frac{d\mathbf{p}_i}{dt} &= \sum_{j \in N_i} -k(\mathbf{q}_i - \mathbf{q}_j) - k_1 \cos \omega t,\end{aligned}\quad (33)$$

526 The Hamilton's equations are:

$$\begin{aligned}\frac{\partial \mathcal{H}(\mathbf{q}, \mathbf{p})}{\partial \mathbf{p}_i} &= \frac{d\mathbf{q}_i}{dt} = \frac{\mathbf{p}_i}{m} \\ \frac{\partial \mathcal{H}(\mathbf{q}, \mathbf{p})}{\partial \mathbf{q}_i} &= -\frac{d\mathbf{p}_i}{dt} = \sum_{j \in N_i} k(\mathbf{q}_i - \mathbf{q}_j) + k_1 \cos \omega t,\end{aligned}\quad (34)$$

527 Hence, we can obtain the Hamiltonian through the integration of the above equation:

$$\mathcal{H}(\mathbf{q}, \mathbf{p}) = \sum_{i=1}^N \frac{\mathbf{p}_i^2}{2m_i} + \frac{1}{2} \sum_{i=1}^N \sum_{j \in N_i} \frac{1}{2} k(\mathbf{q}_i - \mathbf{q}_j)^2 + \sum_{i=1}^N q_i * k_1 \cos \omega t, \quad (35)$$

528 **Verify the systems' energy conservation**

$$\begin{aligned}\frac{d\mathcal{H}(\mathbf{q}, \mathbf{p})}{dt} &= \frac{1}{dt} \left(\sum_{i=1}^N \frac{\mathbf{p}_i^2}{2m_i} \right) + \frac{1}{dt} \left(\frac{1}{2} \sum_{i=1}^N \sum_{j \in N_i} \frac{1}{2} k(\mathbf{q}_i - \mathbf{q}_j)^2 \right) + \frac{1}{dt} \left(\sum_{i=1}^N q_i * k_1 \cos \omega t \right) \\ &= 0 + \frac{1}{dt} \left(\sum_{i=1}^N q_i k_1 \cos \omega t \right) \\ &= \left(\sum_{i=1}^N -\omega q_i k_1 \sin \omega t \right) \neq 0\end{aligned}\quad (36)$$

529 So it is non-conservative.

530 **Verify the systems' time-reversal symmetry** We do the transformation $R : (\mathbf{q}, \mathbf{p}, t) \mapsto (\mathbf{q}, -\mathbf{p}, -t)$.

$$\begin{aligned}\mathcal{H}(\mathbf{q}, \mathbf{p}) &= \sum_{i=1}^N \frac{\mathbf{p}_i^2}{2m_i} + \frac{1}{2} \sum_{i=1}^N \sum_{j \in N_i} \frac{1}{2} k(\mathbf{q}_i - \mathbf{q}_j)^2 + \sum_{i=1}^N q_i * k_1 \cos \omega t, \\ \mathcal{H}(\mathbf{q}, -\mathbf{p}) &= \sum_{i=1}^N \frac{(-\mathbf{p}_i)^2}{2m_i} + \frac{1}{2} \sum_{i=1}^N \sum_{j \in N_i} \frac{1}{2} k(\mathbf{q}_i - \mathbf{q}_j)^2 + \sum_{i=1}^N q_i * k_1 \cos \omega(-t),\end{aligned}\tag{37}$$

531 It is obvious $\mathcal{H}(\mathbf{q}, \mathbf{p}, t) = \mathcal{H}(\mathbf{q}, -\mathbf{p}, t)$, so it is reversible

532 **B.3 Non-conservative and irreversible systems.**

533 A simple example is an n-body spring system with frictions proportional to its velocity, γ is the
534 coefficient of friction, which can be described by:

$$\begin{aligned}\frac{d\mathbf{q}_i}{dt} &= \frac{\mathbf{p}_i}{m} \\ \frac{d\mathbf{p}_i}{dt} &= -k_0 \mathbf{q}_i - \gamma \frac{\mathbf{p}_i}{m}\end{aligned}\tag{38}$$

535 The Hamilton's equations are:

$$\begin{aligned}\frac{\partial \mathcal{H}(\mathbf{q}, \mathbf{p})}{\partial \mathbf{p}_i} &= \frac{d\mathbf{q}_i}{dt} = \frac{\mathbf{p}_i}{m} \\ \frac{\partial \mathcal{H}(\mathbf{q}, \mathbf{p})}{\partial \mathbf{q}_i} &= -\frac{d\mathbf{p}_i}{dt} = \sum_{j \in N_i} k(\mathbf{q}_i - \mathbf{q}_j) + \gamma \frac{\mathbf{p}_i}{m}\end{aligned}\tag{39}$$

536 Hence, we can obtain the Hamiltonian through the integration of the above equation:

$$\mathcal{H}(\mathbf{q}, \mathbf{p}) = \sum_{i=1}^N \frac{\mathbf{p}_i^2}{2m_i} + \frac{1}{2} \sum_{i=1}^N \sum_{j \in N_i} \frac{1}{2} k(\mathbf{q}_i - \mathbf{q}_j)^2 + \sum_{i=1}^N \frac{\gamma}{m} \int_0^t \frac{\mathbf{p}_i^2}{m} dt,\tag{40}$$

537 **Verify the systems' energy conservation**

$$\begin{aligned}\frac{d\mathcal{H}(\mathbf{q}, \mathbf{p})}{dt} &= \frac{1}{dt} \left(\sum_{i=1}^N \frac{\mathbf{p}_i^2}{2m_i} \right) + \frac{1}{dt} \left(\frac{1}{2} \sum_{i=1}^N \sum_{j \in N_i} \frac{1}{2} k(\mathbf{q}_i - \mathbf{q}_j)^2 \right) + \frac{1}{dt} \left(\sum_{i=1}^N \frac{\gamma}{m} \int_0^t \frac{\mathbf{p}_i^2}{m} dt \right) \\ &= 0 + \frac{1}{dt} \left(\sum_{i=1}^N \frac{\gamma}{m} \int_0^t \frac{\mathbf{p}_i^2}{m} dt \right) \\ &= \left(\sum_{i=1}^N \frac{\gamma}{m} \frac{\mathbf{p}_i^2}{m} \right) \neq 0\end{aligned}\tag{41}$$

538 So it is non-conservative.

539 **Verify the systems' time-reversal symmetry** We do the transformation $R : (\mathbf{q}, \mathbf{p}, t) \mapsto (\mathbf{q}, -\mathbf{p}, -t)$.

$$\begin{aligned}\mathcal{H}(\mathbf{q}, \mathbf{p}) &= \sum_{i=1}^N \frac{\mathbf{p}_i^2}{2m_i} + \frac{1}{2} \sum_{i=1}^N \sum_{j \in N_i} \frac{1}{2} k(\mathbf{q}_i - \mathbf{q}_j)^2 + \sum_{i=1}^N \frac{\gamma}{m} \int_0^t \frac{\mathbf{p}_i^2}{m} dt, \\ \mathcal{H}(\mathbf{q}, -\mathbf{p}) &= \sum_{i=1}^N \frac{(-\mathbf{p}_i)^2}{2m_i} + \frac{1}{2} \sum_{i=1}^N \sum_{j \in N_i} \frac{1}{2} k(\mathbf{q}_i - \mathbf{q}_j)^2 + \sum_{i=1}^N \frac{\gamma}{m} \int_0^{(-t)} \frac{\mathbf{p}_i^2}{m} d(-t),\end{aligned}\tag{42}$$

540 It is obvious $\mathcal{H}(\mathbf{q}, \mathbf{p}, t) \neq \mathcal{H}(\mathbf{q}, -\mathbf{p}, t)$, so it is irreversible

541 C Dataset

542 In our experiments, all datasets are synthesized from ground-truth physical law via simulation. We
 543 generate five simulated datasets: three n -body spring systems under damping, periodic, or no external
 544 force, one chaotic tripe pendulum dataset with three sequentially connected stiff sticks that form and
 545 a chaotic strange attractor. We name the first three as *Sipmle Spring*, *Forced Spring*, and *Damped*
 546 *Spring* respectively. For multi-agent systems, all n -body spring systems contain 5 interacting balls,
 547 with varying connectivities. Each *Pendulum* system contains 3 connected stiff sticks. For single-agent
 548 systems, all spring systems contain only one ball. For the chaotic single *Attractor*, we follow the
 549 setting of (Huh et al., 2020).

550 For the n -body spring system, we randomly sample whether a pair of objects are connected, and
 551 model their interaction via forces defined by Hooke’s law. In the *Damped spring*, the objects have an
 552 additional friction force that is opposite to their moving direction and whose magnitude is proportional
 553 to their speed. In the *Forced spring*, all objects have the same external force that changes direction
 554 periodically. We show in Figure 1(a), the energy variation in both of the *Damped spring* and *Forced*
 555 *spring* is significant. For the chaotic tripe *Pendulum*, the equations governing the motion are
 556 inherently nonlinear. Although this system is deterministic, it is also highly sensitive to the initial
 557 condition and numerical errors (Shinbrot et al., 1992; Awrejcewicz et al., 2008; Stachowiak and
 558 Okada, 2006). This property is often referred to as the "butterfly effect", as depicted in Figure 9.
 559 Unlike for n -body spring systems, where the forces and equations of motion can be easily articulated,
 560 for the *Pendulum*, the explicit forces cannot be directly defined, and the motion of objects can only
 561 be described through Lagrangian formulations (North, 2021), making the modeling highly complex
 and raising challenges for accurate learning. We simulate the trajectories by using Euler’s method for

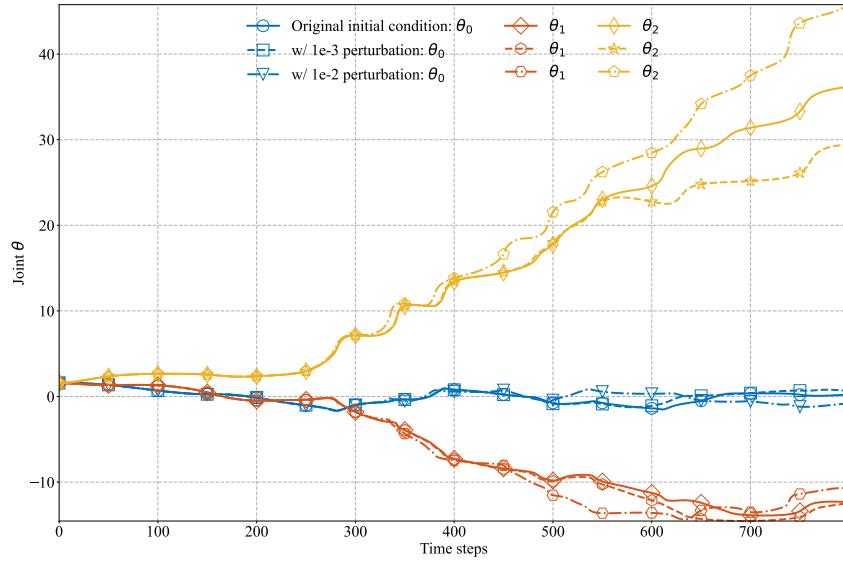


Figure 9: Illustration to show the pendulum is highly-sensitive to initial states

562 n -body spring systems and using the 4th order Runge-Kutta (RK4) method for the *Pendulum* and
 563 *Attractor*. For all spring systems and *Pendulum*, We integrate with a fixed step size and subsample
 564 every 100 steps. For training, we use a total of 6000 forward steps. To generate irregularly sampled
 565 partial observations, we follow (Huang et al., 2020) and sample the number of observations n from a
 566 uniform distribution $U(40, 52)$ and draw the n observations uniformly for each object. For testing, we
 567 additionally sample 40 observations following the same procedure from PDE steps [6000, 12000],
 568 besides generating observations from steps [1, 6000]. The above sampling procedure is conducted
 569 independently for each object. We generate 20k training samples and 5k testing samples for each
 570 dataset. For *Attractor*, we integrate a total of 600 forward steps for training and subsample every
 571 10 steps. For testing, we additionally sample 40 observations from step [600, 1200]. The irregularly
 572 sampled partial observations generation is the same as above. We generate 1000 training samples
 573 and 50 testing samples following (Huh et al., 2020). Therefore, for all datasets, condition length is

60 steps and prediction length is 40s steps. The features (position/velocity) are normalized to the maximum absolute value of 1 across training and testing datasets.

We also compute the Maximum Lyapunov Exponent (MLE) to assess the chaos level of the systems, using the formula:

$$\lambda = \max_{t \rightarrow \inf} \left(\frac{1}{t} \ln \frac{\|\delta(t)\|}{\|\delta(0)\|} \right).$$

We set fixed initial values for each dataset and generate 10 trajectories by perturbing the initial values with random noise (0, 0.0001). We calculate the Maximum Lyapunov Exponent (MLE) between any two trajectories. Finally, we compute the average and std of MLE from all pairs to gauge the chaotic behavior of each dataset. The data is presented in the table below:

Table 2: MLE of different Multi-agent Systems

Dataset	<i>Simple Spring</i>	<i>Forced Spring</i>	<i>Damped Spring</i>	<i>Pendulum</i>
MLE(in 60 steps)	0.4031 ± 0.3944	1.0087 ± 1.0577	0.6307 ± 0.7065	34.1832 ± 30.1846

From the table, it's evident that the order of MLE values is: *Pendulum* » three *Spring* datasets. This observation is consistent with the evaluation results based on MSE presented in our previous responses in Table 3 which indicates that as the prediction length(steps*step size) increases, there is a more significant performance degradation of all models on *Pendulum* dataset.

In the following subsections, we show the dynamical equations of each dataset in detail.

C.1 Spring Systems

C.1.1 Simple Spring

The dynamical equations of *simple spring* are as follows:

$$\begin{aligned} \frac{d\mathbf{q}_i}{dt} &= \frac{\mathbf{p}_i}{m} \\ \frac{d\mathbf{p}_i}{dt} &= \sum_{j \in N_i}^N -k(\mathbf{q}_i - \mathbf{q}_j) \end{aligned} \quad (43)$$

where where $\mathbf{q} = (\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_N)$ is a set of positions of each object, $\mathbf{p} = (\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_N)$ is a set of momenta of each object. We set the mass of each object $m = 1$, the spring constant $k = 0.1$.

C.1.2 Damped Spring

The dynamical equations of *damped spring* are as follows:

$$\begin{aligned} \frac{d\mathbf{q}_i}{dt} &= \frac{\mathbf{p}_i}{m} \\ \frac{d\mathbf{p}_i}{dt} &= \sum_{j \in N_i}^N -k(\mathbf{q}_i - \mathbf{q}_j) - \gamma \frac{\mathbf{p}_i}{m} \end{aligned} \quad (44)$$

where where $\mathbf{q} = (\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_N)$ is a set of positions of each object, $\mathbf{p} = (\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_N)$ is a set of momenta of each object, We set the mass of each object $m = 1$, the spring constant $k = 0.1$, the coefficient of friction $\gamma = 10$.

C.1.3 Forced Spring

The dynamical equations of *forced spring* system are as follows:

$$\begin{aligned} \frac{d\mathbf{q}_i}{dt} &= \frac{\mathbf{p}_i}{m} \\ \frac{d\mathbf{p}_i}{dt} &= \sum_{j \in N_i}^N -k(\mathbf{q}_i - \mathbf{q}_j) - k_1 \cos \omega t, \end{aligned} \quad (45)$$

600 where $\mathbf{q} = (\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_N)$ is a set of positions of each object, $\mathbf{p} = (\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_N)$ is a
 601 set of momenta of each object. We set the mass of each object $m = 1$, the spring constant $k = 0.1$,
 602 the external strength $k_1 = 10$ and the frequency of variation $\omega = 1$

603 We simulate the positions and momentums of three spring systems by using Euler methods as follows:

$$\begin{aligned}\mathbf{q}_i(t+1) &= \mathbf{q}_i(t) + \frac{d\mathbf{q}_i}{dt} \Delta t \\ \mathbf{p}_i(t+1) &= \mathbf{p}_i(t) + \frac{d\mathbf{p}_i}{dt} \Delta t\end{aligned}\tag{46}$$

604 where $\frac{d\mathbf{q}_i}{dt}$ and $\frac{d\mathbf{p}_i}{dt}$ were defined as above for each datasets, and $\Delta t = 0.001$ is the integration steps.

605 C.2 Chaotic Pendulum

606 In this section, we demonstrate how to derive the dynamics equations for a chaotic triple pendulum
 607 using the Lagrangian formalism.

608 The moment of inertia of each stick about the centroid is

$$I = \frac{1}{12}ml^2\tag{47}$$

609 The position of the center of gravity of each stick is as follows:

$$\begin{aligned}x_1 &= \frac{l}{2} \sin \theta_1, \quad y_1 = -\frac{l}{2} \cos \theta_1 \\ x_2 &= l(\sin \theta_1 + \frac{1}{2} \sin \theta_2), \quad y_2 = -l(\cos \theta_1 + \frac{1}{2} \cos \theta_2) \\ x_3 &= l(\sin \theta_1 + \sin \theta_2 + \frac{1}{2} \sin \theta_3), \quad y_3 = -l(\cos \theta_1 + \cos \theta_2 + \frac{1}{2} \cos \theta_3)\end{aligned}\tag{48}$$

610 The change in the center of gravity of each stick is:

$$\begin{aligned}\dot{x}_1 &= \frac{l}{2} \cos \theta_1 \cdot \dot{\theta}_1, \quad \dot{y}_1 = \frac{l}{2} \sin \theta_1 \cdot \dot{\theta}_1 \\ \dot{x}_2 &= l(\cos \theta_1 \cdot \dot{\theta}_1 + \frac{1}{2} \cos \theta_2 \cdot \dot{\theta}_2), \quad \dot{y}_2 = l(\sin \theta_1 \cdot \dot{\theta}_1 + \frac{1}{2} \sin \theta_2 \cdot \dot{\theta}_2) \\ \dot{x}_3 &= l(\cos \theta_1 \cdot \dot{\theta}_1 + \cos \theta_2 \cdot \dot{\theta}_2 + \frac{1}{2} \cos \theta_3 \cdot \dot{\theta}_3), \quad \dot{y}_3 = l(\sin \theta_1 \cdot \dot{\theta}_1 + \sin \theta_2 \cdot \dot{\theta}_2 + \frac{1}{2} \sin \theta_3 \cdot \dot{\theta}_3)\end{aligned}\tag{49}$$

611 The Lagrangian \mathcal{L} of this triple pendulum system is:

$$\begin{aligned}\mathcal{L} &= T - V \\ &= \frac{1}{2}m(\dot{x}_1^2 + \dot{x}_2^2 + \dot{x}_3^2 + \dot{y}_1^2 + \dot{y}_2^2 + \dot{y}_3^2) + \frac{1}{2}I(\dot{\theta}_1^2 + \dot{\theta}_2^2 + \dot{\theta}_3^2) - mg(y_1 + y_2 + y_3) \\ &= \frac{1}{6}ml(9\dot{\theta}_2\dot{\theta}_1l \cos(\theta_1 - \theta_2) + 3\dot{\theta}_3\dot{\theta}_1l \cos(\theta_1 - \theta_3) + 3\dot{\theta}_2\dot{\theta}_3l \cos(\theta_2 - \theta_3) + 7\dot{\theta}_1^2l + 4\dot{\theta}_2^2l + \dot{\theta}_3^2l \\ &\quad + 15g \cos(\theta_1) + 9g \cos(\theta_2) + 3g \cos(\theta_3))\end{aligned}\tag{50}$$

612 The Lagrangian equation is defined as follows:

$$\frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{\theta}} - \frac{\partial \mathcal{L}}{\partial \theta} = \mathbf{0}\tag{51}$$

613 and we also have:

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \dot{\theta}} &= \frac{\partial T}{\partial \dot{\theta}} = p \\ \dot{p} &= \frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{\theta}} = \frac{\partial \mathcal{L}}{\partial \theta}\end{aligned}\tag{52}$$

614 where p is the Angular Momentum.

615 We can list the equations for each of the three sticks separately:

$$\begin{aligned}p_1 &= \frac{\partial \mathcal{L}}{\partial \dot{\theta}_1} \quad \dot{p}_1 = \frac{\partial \mathcal{L}}{\partial \theta_1} \\ p_2 &= \frac{\partial \mathcal{L}}{\partial \dot{\theta}_2} \quad \dot{p}_2 = \frac{\partial \mathcal{L}}{\partial \theta_2} \\ p_3 &= \frac{\partial \mathcal{L}}{\partial \dot{\theta}_3} \quad \dot{p}_3 = \frac{\partial \mathcal{L}}{\partial \theta_3}\end{aligned}\tag{53}$$

616 Finally, we have :

$$\left\{ \begin{array}{l} \dot{\theta}_1 = \frac{6(9p_1 \cos(2(\theta_2 - \theta_3)) + 27p_2 \cos(\theta_1 - \theta_2) - 9p_2 \cos(\theta_1 + \theta_2 - 2\theta_3) + 21p_3 \cos(\theta_1 - \theta_3) - 27p_3 \cos(\theta_1 - 2\theta_2 + \theta_3) - 23p_1)}{ml^2(81 \cos(2(\theta_1 - \theta_2)) - 9 \cos(2(\theta_1 - \theta_3)) + 45 \cos(2(\theta_2 - \theta_3)) - 169)} \\ \dot{\theta}_2 = \frac{6(27p_1 \cos(\theta_1 - \theta_2) - 9p_1 \cos(\theta_1 + \theta_2 - 2\theta_3) + 9p_2 \cos(2(\theta_1 - \theta_3)) - 27p_3 \cos(2\theta_1 - \theta_2 - \theta_3) + 57p_3 \cos(\theta_2 - \theta_3) - 47p_2)}{ml^2(81 \cos(2(\theta_1 - \theta_2)) - 9 \cos(2(\theta_1 - \theta_3)) + 45 \cos(2(\theta_2 - \theta_3)) - 169)} \\ \dot{\theta}_3 = \frac{6(21p_1 \cos(\theta_1 - \theta_3) - 27p_1 \cos(\theta_1 - 2\theta_2 + \theta_3) - 27p_2 \cos(2\theta_1 - \theta_2 - \theta_3) + 57p_2 \cos(\theta_2 - \theta_3) + 81p_3 \cos(2(\theta_1 - \theta_2)) - 143p_3)}{ml^2(81 \cos(2(\theta_1 - \theta_2)) - 9 \cos(2(\theta_1 - \theta_3)) + 45 \cos(2(\theta_2 - \theta_3)) - 169)} \\ p_1 = -\frac{1}{2}ml \left(3\dot{\theta}_2\dot{\theta}_1 l \sin(\theta_1 - \theta_2) + \dot{\theta}_1\dot{\theta}_3 l \sin(\theta_1 - \theta_3) + 5g \sin(\theta_1) \right) \\ p_1 = -\frac{1}{2}ml \left(-3\dot{\theta}_1\dot{\theta}_2 l \sin(\theta_1 - \theta_2) + \dot{\theta}_2\dot{\theta}_3 l \sin(\theta_2 - \theta_3) + 3g \sin(\theta_2) \right) \\ p_1 = -\frac{1}{2}ml \left(\dot{\theta}_1\dot{\theta}_3 l \sin(\theta_1 - \theta_3) + \dot{\theta}_2\dot{\theta}_3 l \sin(\theta_2 - \theta_3) - g \sin(\theta_3) \right) \end{array} \right. \quad (54)$$

617 We simulate the angular of the three sticks by using the Runge-Kutta 4th Order Method as follows:

$$\begin{aligned} \Delta\theta^1(t) &= \dot{\theta}(t, \theta(t)) \cdot \Delta t \\ \Delta\theta^2(t) &= \dot{\theta}\left(t + \frac{\Delta t}{2}, \theta(t) + \frac{\Delta\theta^1(t)}{2}\right) \cdot \Delta t \\ \Delta\theta^3(t) &= \dot{\theta}\left(t + \frac{\Delta t}{2}, \theta(t) + \frac{\Delta\theta^2(t)}{2}\right) \cdot \Delta t \\ \Delta\theta^4(t) &= \dot{\theta}(t + \Delta t, \theta(t) + \Delta\theta^3(t)) \cdot \Delta t \\ \Delta\theta(t) &= \frac{1}{6}(\Delta\theta^1(t) + \Delta\theta^2(t) + \Delta\theta^3(t) + \Delta\theta^4(t)) \\ \theta(t + 1) &= \theta(t) + \Delta\theta(t) \end{aligned} \quad (55)$$

618 where $\dot{\theta}$ was defined as above , and $\Delta t = 0.0001$ is the integration steps.

619 C.3 Chaotic Strange Attractor

620 The dynamical equations of this reversible strange attractor are as follows:

$$\begin{aligned} \frac{dx}{dt} &= 1 + yz, \\ \frac{dy}{dt} &= -xz, \\ \frac{dz}{dt} &= y^2 + 2yz, \\ x, y, z &\in \mathbb{R} \end{aligned} \quad (56)$$

621 The above equations can be presented as $(\dot{x}(t), \dot{y}(t), \dot{z}(t)) = \text{Dynamic}(x(t), y(t), z(t))$.

622 We simulate $K(t) = (x(t), y(t), z(t))$ by using the Runge-Kutta 4th Order Method as follows:

$$\begin{aligned} \Delta K_1(t) &= \text{Dynamic}(K(t)) * \Delta t \\ \Delta K_2(t) &= \text{Dynamic}\left(K(t) + \frac{\Delta K_1(t)}{2}\right) * \Delta t \\ \Delta K_3(t) &= \text{Dynamic}\left(K(t) + \frac{\Delta K_2(t)}{2}\right) * \Delta t \\ \Delta K_4(t) &= \text{Dynamic}(K(t) + \Delta K_3(t)) * \Delta t \\ \Delta K(t) &= \frac{1}{6}(\Delta K_1(t) + \Delta K_2(t) + \Delta K_3(t) + \Delta K_4(t)) \\ K(t + 1) &= K(t) + \Delta K(t) \end{aligned} \quad (57)$$

623 We sampling $z(t_0)$ randomly from uniform distribution $[1, 3]$ while fixing $x(t_0) = y(t_0) = 0$. We set
624 the trajectory lengths of both training and test dataset to 600, with regular time-step size $\Delta t = 0.03$
625 and the sample frequency of 10. We add Gaussian noise $0.05n$, $n \sim \mathcal{N}(0, 1)$ to training trajectories.

626 C.4 Human Motion

627 For the real-world motion capture dataset(CMU, 2003), we focus on the walking sequences of subject
 628 35. Each sample in this dataset is represented by 31 trajectories, each corresponding to the movement
 629 of a single joint. For each joint, we first randomly sample the number of observations from a uniform
 630 distribution $\mathcal{U}(30, 42)$ and then sample uniformly from the first 50 frames for training and validation
 631 trajectories. For testing, we additionally sampled 40 observations from frames [51, 99]. We split
 632 different walking sequences into training (15 trials) and test sets (7 trials). For each walking sequence,
 633 we further split it into several non-overlapping small sequences with maximum length 50 for training,
 634 and maximum length 100 for testing. In this way, we generate total 120 training samples and 27
 635 testing samples. We normalize all features (position/velocity) to maximum absolute value of 1 across
 636 training and testing datasets.

637 D Model Details

638 In the following we introduce in details how we implement our model and each baseline.

639 D.1 Initial State Encoder

640 For multi-agent systems, the initial state encoder computes the latent node initial states $z_i(t)$ for all
 641 agents simultaneously considering their mutual interaction. Specifically, it first fuses all observations
 642 into a temporal graph and conducts dynamic node representation through a spatial-temporal GNN as
 643 in (Huang et al., 2020):

$$\begin{aligned} \mathbf{h}_{j(t)}^{l+1} &= \mathbf{h}_{j(t)}^l + \sigma \left(\sum_{i(t') \in \mathcal{N}_{j(t)}} \alpha_{i(t') \rightarrow j(t)}^l \times \mathbf{W}_v \hat{\mathbf{h}}_{i(t')}^{l-1} \right) \\ \alpha_{i(t') \rightarrow j(t)}^l &= \left(\mathbf{W}_k \hat{\mathbf{h}}_{i(t')}^{l-1} \right)^T \left(\mathbf{W}_q \mathbf{h}_{j(t)}^{l-1} \right) \cdot \frac{1}{\sqrt{d}}, \quad \hat{\mathbf{h}}_{i(t')}^{l-1} = \mathbf{h}_{i(t')}^{l-1} + \text{TE}(t' - t) \\ \text{TE}(\Delta t)_{2i} &= \sin \left(\frac{\Delta t}{10000^{2i/d}} \right), \quad \text{TE}(\Delta t)_{2i+1} = \cos \left(\frac{\Delta t}{10000^{2i/d}} \right), \end{aligned} \quad (58)$$

644 where \parallel denotes concatenation; $\sigma(\cdot)$ is a non-linear activation function; d is the dimension of node
 645 embeddings. The node representation is computed as a weighted summation over its neighbors
 646 plus residual connection where the attention score is a transformer-based (Vaswani et al., 2017)
 647 dot-product of node representations by the use of value, key, query projection matrices $\mathbf{W}_v, \mathbf{W}_k, \mathbf{W}_q$.
 648 Here $\mathbf{h}_{j(t)}^l$ is the representation of agent j at time t in the l -th layer. $i(t')$ is the general index for
 649 neighbors connected by temporal edges (where $t' \neq t$) and spatial edges (where $t = t'$ and $i \neq j$).
 650 The temporal encoding (Hu et al., 2020) is added to a neighborhood node representation in order
 651 to distinguish its message delivered via spatial and temporal edges. Then, we stack L layers to get
 652 the final representation for each observation node: $\mathbf{h}_i^t = \mathbf{h}_{i(t)}^L$. Finally, we employ a self-attention
 653 mechanism to generate the sequence representation \mathbf{u}_i for each agent as their latent initial states:

$$\mathbf{u}_i = \frac{1}{K} \sum_t \sigma \left(\mathbf{a}_i^T \hat{\mathbf{h}}_i^t \hat{\mathbf{h}}_i^t \right), \quad \mathbf{a}_i = \tanh \left(\left(\frac{1}{K} \sum_t \hat{\mathbf{h}}_i^t \right) \mathbf{W}_a \right), \quad (59)$$

654 where \mathbf{a}_i is the average of observation representations with a nonlinear transformation \mathbf{W}_a and
 655 $\hat{\mathbf{h}}_i^t = \mathbf{h}_i^t + \text{TE}(t)$. K is the number of observations for each trajectory. Compared with recurrent
 656 models such as RNN, LSTM (Sepp and Jürgen, 1997), it offers better parallelization for accelerating
 657 training speed and in the meanwhile alleviates the vanishing/exploding gradient problem brought by
 658 long sequences. For single-agent Systems, there only left the self-attention mechanism component.

659 Given the latent initial states, the dynamics of the whole system are determined by the ODE function
 660 g which we parametrize as a GNN as in (Huang et al., 2020) for Multi-Agent Systems to capture the
 661 continuous interaction among agents. For single-agent systems, we only include self-loop edges in
 662 the graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, which makes the ODE function g a simple MLP.

663 We then employ Multilayer Perceptron (MLP) as a decoder to predict the trajectories $\hat{\mathbf{y}}_i(t)$ from the
 664 latent states $z_i(t)$.

$$\begin{aligned} z_1(t), z_2(t), z_3(t) \cdots z_N(t) &= \text{ODEsolver}(g, [z_1(t_0), z_2(t_0) \cdots z_N(t_0)], (t_0, t_1 \cdots t_K)) \\ \hat{y}_i(t) &= f_{dec}(z_i(t)) \end{aligned} \quad (60)$$

D.2 Implementation Details

TREAT

For multi-agent systems, our implementation of TREAT follows GraphODE pipeline. We implement the initial state encoder using a 2-layer GNN with a hidden dimension of 64 across all datasets. We use ReLU for nonlinear activation. For the sequence self-attention module, we set the output dimension to 128. The encoder’s output dimension is set to 16, and we add 64 additional dimensions initialized with all zeros to the latent states $z_i(t)$ to stabilize the training processes as in (Huang et al., 2021). The GNN ODE function is implemented with a single-layer GNN from (Kipf et al., 2018) with hidden dimension 128. For single-agent systems, we only include self-loop edges in the graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, which makes the ODE function g a simple MLP. To compute trajectories, we use the Runge-Kutta method from torchdiffeq python package s(Chen et al., 2021) as the ODE solver and a one-layer MLP as the decoder.

We implement our model in pytorch. Encoder, generative model, and the decoder parameters are jointly optimized with AdamW optimizer (Loshchilov and Hutter, 2019) using a learning rate of 0.0001 for spring datasets and 0.00001 for *Pendulum*. The batch size for all datasets is set to 512.

TREAT $_{\mathcal{L}_{rev}=gt-rev}$ and TREAT $_{\mathcal{L}_{rev}=rev2}$ share the same architecture and hyparameters as TREAT, with different implementations of the loss function. In TREAT $_{\mathcal{L}_{rev}=gt-rev}$, instead of comparing forward and reverse trajectories, we look at the L2 distance between the ground truth and reverse trajectories when computing the reversal loss.

For TREAT $_{\mathcal{L}_{rev}=rev2}$, we implement the reversal loss following (Huh et al., 2020) with one difference: we do not apply the reverse operation to the momentum portion of the initial state to the ODE function. This is because the initial hidden state is an output of the encoder that mixes position and momentum information. Note that we also remove the additional dimensions to the latent state that TREAT has. To reproduce our model’s results, we provide our code implementation link here.

LatentODE

We implement the Latent ODE sequence to sequence model as specified in (Rubanova et al., 2019). We use a 4-layer ODE function in the recognition ODE, and a 2-layer ODE function in the generative ODE. The recognition and generative ODEs use Euler and Dopri5 as solvers (Chen et al., 2021), respectively. The number of units per layer is 1000 in the ODE functions and 50 in GRU update networks. The dimension of the recognition model is set to 100. The model is trained with a learning rate of 0.001 with an exponential decay rate of 0.999 across different experiments. Note that since latentODE is a single-agent model, we compute the trajectory of each object independently when applying it to multi-agent systems.

HODEN

To adapt HODEN, which requires full initial states of all objects, to systems with partial observations, we compute each object’s initial state via linear spline interpolation if it is missing. Following the setup in (Huh et al., 2020), we have two 2-layer linear networks with Tanh activation in between as ODE functions, in order to model both positions and momenta. Each network has a 1000-unit layer followed by a single-unit layer. The model is trained with a learning rate of 0.00001 using a cosine scheduler. HODEN is a single-agent model, we compute the trajectory of each object independently when applying it to multi-agent systems.

TRS-ODEN

Similar to HODEN, we compute each object’s initial state via linear spline interpolation if it is missing. As in (Huh et al., 2020), we use a 2-layer linear network with Tanh activation in between as the ODE functions, and the Leapfrog method for solving ODEs. The network has 1000 hidden units and is trained with a learning rate of 0.00001 using a cosine scheduler. TRS-ODEN is a single-agent model, we compute the trajectory of each object independently when applying it to multi-agent systems.

713 TRS-ODEN_{GNN}

714 For TRSODEN_{GNN}, we substitute the ODE function in TRS-ODEN with a GraphODE network. The
 715 GraphODE generative model is implemented with a single-layer GNN with hidden dimension 128.
 716 As in HODEN and TRS-ODEN, we compute each object’s missing initial state via linear spline
 717 interpolation and the Leapfrog method for solving ODE. For all datasets, we use 0.5 as the coefficient
 718 for the reversal loss in (Huh et al., 2020), and 0.0002 as the learning rate under cosine scheduling.

719 LGODE

720 Our implementation follows (Huang et al., 2020) except we remove the Variational Autoencoder
 721 (VAE) from the initial state encoder. Instead of using the output from the encoder GNN as the
 722 mean and std of the VAE, we directly use it as the latent initial state. That is, the initial states are
 723 deterministic instead of being sampled from a distribution. We use the same architecture as in TREAT
 724 and train the model using an AdamW optimizer with a learning rate of 0.0001 across all datasets.

725 E Additional Experiments

726 E.1 Comparison of different solvers

727 We next show our model’s sensitivity regarding solvers with different precisions. Specifically, we
 728 compare against Euler and Runge-Kutta (RK4) where the latter is a higher-precision solver. We show
 729 the comparison against LGODE and TREAT in Table 3.

730 We can firstly observe that TREAT consistently outperforms LGODE, which is our strongest baseline
 731 across different solvers and datasets, indicating the effectiveness of the proposed time-reversal
 732 symmetry loss. Secondly, we compute the improvement ratio as $\frac{LGODE - TREAT}{LGODE}$. We can see that
 733 the improvement ratios get larger when using RK4 over Euler. This can be understood as our reversal
 734 loss is minimizing higher-order Tayler expansion terms (Theorem 3.1) thus compensating numerical
 735 errors brought by ODE solvers.

Table 3: Evaluation results on MSE (10^{-2}) over different solvers for multi-agent systems.

Dataset Solvers	<i>Simple Spring</i>		<i>Forced Spring</i>		<i>Damped Spring</i>		<i>Pendulum</i>	
	Euler	RK4	Euler	RK4	Euler	RK4	Euler	RK4
LGODE	1.8443	1.7429	2.0462	1.8929	1.1686	0.9718	1.4634	1.4156
TREAT	1.4864	1.1178	1.6058	1.4525	0.8070	0.5944	1.3093	1.2527
% Improvement	19.4057	35.8655	21.5228	23.2659	30.9430	38.8352	10.5303	11.5075

736 E.2 Evaluation across observation ratios.

737 For LG-ODE and TREAT, the encoder computes the initial states from observed trajectories. To show
 738 models’ sensitivity towards data sparsity, we randomly mask out 40% and 80% historical observations
 739 and compare model performance. As shown in Table 4, when changing the ratios from 80% to 40%,
 740 we observe that TREAT has a smaller performance drop compared with LG-ODE, especially on the
 741 more complex Pendulum dataset (LG-ODE decreases 22.04% while TREAT decreases 1.62%). This
 742 indicates that TREAT is less sensitive toward data sparsity.

Table 4: Results of varying observation ratios on MSE (10^{-2}) of multi-agent datasets.

Dataset Observation Ratios	<i>Simple Spring</i>		<i>Forced Spring</i>		<i>Damped Spring</i>		<i>Pendulum</i>	
	0.8	0.4	0.8	0.4	0.8	0.4	0.8	0.4
LG-ODE	1.7054	1.6889	1.7554	2.0370	0.9305	1.0217	1.4314	1.7469
TREAT	1.1176	1.1429	1.3611	1.5109	0.6920	0.6964	1.2309	1.2110

743 E.3 Evaluation for ablation across 5 run

744 The second model variant TREAT _{$\mathcal{L}_{rev}=gt\text{-}rev$} is by computing the reversal loss $\mathcal{L}_{reverse}$ as between
 745 model backward predictions to ground truth, in contrast with our proposed loss between model

backward and forward predictions used in TREAT. In Table 1, we can see that $\text{TREAT}_{\mathcal{L}_{rev}=\text{gt-rev}}$ decreases the performance by 1.20%, 5.02%, 3.82%, and 28.98% for the four systems respectively. We further repeated our experiments multiple rounds and provided std in the following. We observed that TREAT consistently outperforms $\text{TREAT}_{\mathcal{L}_{rev}=\text{gt-rev}}$ and in general has smaller stds.

Table 5: Evaluation Results on MSE (10^{-2}) across 5 runs

Model	<i>Simple Spring</i>	<i>Forced Spring</i>	<i>Damped Spring</i>	<i>Pendulum</i>
TREAT	1.1101 \pm 0.0159	1.4565 \pm 0.0176	0.6023 \pm 0.0112	1.2561 \pm 0.0021
$\text{TREAT}_{\mathcal{L}_{rev}=\text{gt-rev}}$	1.1113 \pm 0.0162	1.5865 \pm 0.0451	0.6209 \pm 0.0160	1.6254 \pm 0.0150

E.4 Comparison of different solver step sizes.

F Discussion about Reversible Neural Networks

In literature, there is another line of research about building reversible neural networks (NNs). For example, (Chang et al., 2018) formulates three architectures for reversible neural networks to address the stability issue and achieve arbitrary deep lengths, motivated by dynamical system modeling. (Liu et al., 2019) employs normalizing flow to create a generative model of graph structures. They all propose novel architectures to construct reversible NN where intermediate states across layer depths do not need to be stored, thus improving memory efficiency.

However, we’d like to clarify that reversible NNs (RevNet) do not resolve the time-reversal symmetry problem that we’re studying. The core of RevNet is that input can be recovered from output via a reversible operation (which is another operator), similar as any linear operator $W(\cdot)$ have a reversed projector $W^{-1}(\cdot)$. In the contrary, what we want to study is that the **same operator** can be used for both forward and backward prediction over time, and keep the trajectory the same. That being said, to generate the forward and backward trajectories, we are using the same $g(\cdot)$, instead of $g(\cdot), g^{-1}(\cdot)$ respectively.

In summary, though both reversible NN and time-reversal symmetry share similar insights and intuition, they’re talking about different things: reversible NNs make every operator $g(\cdot)$ having a $g^{-1}(\cdot)$, while time-reversible assume the trajectory get from $\hat{z}^{fwd} = g(z)$ and $\hat{z}^{bwd} = -g(z)$ to be closer. Making g to be reversible cannot make the system to be time-reversible.

G Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. TREAT is trained upon physical simulation data (e.g., , spring and pendulum) and implemented by public libraries in PyTorch. During the modeling, we neither introduces any social/ethical bias nor amplify any bias in the data. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

H Limitations

Currently, TREAT only incorporates inductive bias from the temporal aspect, while there are many important properties in the spatial aspect such as translation and rotation equivariance (Satorras et al., 2021). Future endeavors that combine biases from both temporal and spatial dimensions could unveil a new frontier in dynamical systems modeling.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: **[TODO]**

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: Limitations are discussed in Appendix H

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: Proofs are in Appendix A.3, A.2 and A.4.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: The Datasets, Task Setup, Baselines description are in Sec. 4. Pseudo code for the implementation of the Time-Reversal Symmetry Loss is in Appendix A.1. More Model Details and Implementation Details are in Appendix D.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The Datasets, Task Setup, Baselines description are in Sec. 4. Pseudo code for the implementation of the Time-Reversal Symmetry Loss is in Appendix A.1. More Dataset descriptions are in Appendix C. More Model Details and Implementation Details are in Appendix D.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The Datasets, Task Setup, Baselines description are in Sec. 4. More Implementation Details are in Appendix D.2.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: In Appendix E.3

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: [TODO]

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: [TODO]

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: In Appendix G.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: [TODO]

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: [TODO]

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: [TODO]

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: [TODO]

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: [TODO]

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.