DISCOVERING PHYSICS LAWS OF DYNAMICAL SYS TEMS VIA INVARIANT FUNCTION LEARNING

Anonymous authors

004

006

008

009 010

011

012

013

014

015

016

017

018

019

021

024

025

026

027

Paper under double-blind review

ABSTRACT

We consider learning underlying laws of dynamical systems governed by ordinary differential equations (ODE). A key challenge is how to discover intrinsic dynamics across multiple environments while circumventing environment-specific mechanisms. Unlike prior work, we tackle more complex environments where changes extend beyond function coefficients to entirely different function forms. For example, we demonstrate the discovery of ideal pendulum's natural motion $\alpha^2 \sin \theta_t$ by observing pendulum dynamics in different environments, such as the damped environment $\alpha^2 \sin(\theta_t) - \rho \omega_t$ and powered environment $\alpha^2 \sin(\theta_t) + \rho \frac{\omega_t}{|\omega_t|}$ Here, we formulate this problem as an invariant function learning task and propose a new method, known as Disentanglement of Invariant Functions (DIF), that is grounded in causal analysis. We propose a causal graph and design an encoder-decoder hypernetwork that explicitly disentangles invariant functions from environment-specific dynamics. The discovery of invariant functions is guaranteed by our information-based principle that enforces the independence between extracted invariant functions and environments. Quantitative comparisons with meta-learning and invariant learning baselines on three ODE systems demonstrate the effectiveness and efficiency of our method. Furthermore, symbolic regression explanation results highlight the ability of our framework to uncover intrinsic laws.

028 1 INTRODUCTION

Deep neural networks (Goodfellow et al., 2016) have been widely used for predicting dynamical systems (Aussem, 1999; Singh et al., 2012; Wang et al., 2016; Lusch et al., 2018; Yeo & Melnyk, 2019; Giannakis, 2019). Numerous efforts (Kirchmeyer et al., 2022; Wang et al., 2022; Mouli et al., 2024) have focused on modeling dynamical systems by forecasting future states from past observations, often emphasizing rapid adaptation to new systems or improved model architectures. However, an important scenario in scientific discovery has been overlooked, *i.e.*, learning invariant mechanisms, which aims to identify shared motion patterns across dynamical systems observed in multiple environments. This task not only facilitates scientific equation discovery but also holds potential for advancing the understanding and extraction of physical laws from observational data—such as images and videos—where physical laws are highly entangled. As the first step in invariant function learning, this paper focuses on ordinary differential equation (ODE) systems.

040 The need for invariant function learning arises because data collected is often observed under 041 varying environments and entangled with multiple factors. For instance, the oscillation of a simple 042 pendulum (Yin et al., 2021b) is commonly influenced by air friction; a prey-predator system (Ahmad, 043 1993) can be affected by limited resources. These factors significantly hinder deep models from 044 learning the true and invariant dynamics. Instead of capturing invariant dynamical patterns, deep models tend to be sensitive to trivial information and spurious correlations, leading to failures in identifying the true and isolated mechanisms. In light of this challenge, we explore an innovative 046 setting called *invariant function learning*, which aims to extract intrinsic mechanisms from data 047 observed in multiple environments. Unlike prior work, we aim to tackle broader and more complex 048 environments where changes extend beyond function coefficients to entirely different function forms. For example, we target the discovery of ideal pendulum's natural motion $\alpha^2 \sin \theta_t$ by observing pendulum dynamics in different environments such as the damped environment $\alpha^2 \sin(\theta_t) - \rho \omega_t$ and 051 powered environment $\alpha^2 \sin(\theta_t) + \rho \frac{\omega_t}{|\omega_t|}$, as shown in our motivation example 1.

052

Invariant function learning presents two key challenges. Firstly, invariant mechanisms are not isolated entities, and being intertwined with varying initial conditions, system parameters, and time makes



Figure 1: Multi-environment Pendulum ODE systems. In this example, ODEs with different coefficients and function forms are used to extract their corresponding invariant functions (green).

them difficult to define or disentangle. Secondly, existing invariant learning techniques (Arjovsky et al., 2019; Lu et al., 2021a; Rosenfeld et al., 2020; Peters et al., 2016), which are primarily designed 071 for categorical tasks, do not extend effortlessly to dynamical systems, requiring the design of new 072 invariant principles. To overcome these challenges, we formulate our invariant function learning 073 framework as a causal graph where functions are parameterized and isolated to be learned and disen-074 tangled. Furthermore, we propose an invariant function learning principle and the implementation 075 of an encoder-decoder hypernetwork (Ha et al., 2016) to identify the true invariant mechanism. 076 Specifically, our contributions are listed as follows. (1) We introduce a new task, invariant function 077 learning, aimed at scientific discovery. (2) We formulate its framework with causal foundations. (3) 078 We propose an invariant function learning principle to identify true invariant mechanisms and design 079 a method, Disentanglement of Invariant Function (DIF), with hypernetwork implementation. (4) 080 To facilitate comprehensive benchmarking, we design several new baselines by adapting existing meta-learning and invariant learning techniques to our function learning framework and conducting 081 extensive experimental studies. 082

2 INVARIANT FUNCTION LEARNING FOR DYNAMICAL SYSTEM

In this section, we first provide the background on ODEs, followed by the introduction and formulation of our invariant function learning task, along with the causal analyses that underpin our proposal.

2.1 ORDINARY DIFFERENTIAL EQUATION DYNAMICAL SYSTEM

We describe a dynamical system using an ordinary differential equation (ODE) as:

$$\frac{d\boldsymbol{x}_t}{dt} = f(\boldsymbol{x}_t),\tag{1}$$

where $x_t \in \mathcal{X} \subseteq \mathbb{R}^d$ includes *d* hidden states of the system at time *t*. $f \in \mathcal{F} : \mathcal{X} \mapsto T\mathcal{X}$ is the derivative function of the dynamical system mapping the hidden states to their tangent space, where \mathcal{F} is the function space containing functions that describe all dynamical systems with *d* hidden states.

Given proper time discretization, we consider T time steps denoted as $t = t_0, t_1, \ldots, t_{T-1}$. The corresponding T-length trajectory can be written as a matrix $X = [\mathbf{x}_{t_0}, \mathbf{x}_{t_1}, \ldots, \mathbf{x}_{t_{T-1}}] \in \mathbb{R}^{d \times T}$. Given the system hidden states before a certain time step $T_c \in \mathbb{N}$, denoted as $X_p = X_{:,0:T_c} = [\mathbf{x}_{t_0}, \ldots, \mathbf{x}_{t_{T-1}}] \in \mathbb{R}^{d \times T_c}$, the forecasting task aims to predict the future trajectory $X_{:,T_c:T} = [\mathbf{x}_{t_c}, \ldots, \mathbf{x}_{t_{T-1}}]$. For theoretical analysis, we represent random variables in boldface, *e.g.*, the matrix-valued random variable corresponding to X is denoted as \mathbf{X} ; the function f is a realization of the function variable f from a given function space. Full notations are detailed in Tab. 3.

104 105

106

054

055

056

058

060

061

062 063

064 065 066

083 084

085

087 088

089 090

- 2.2 INVARIANT FUNCTION LEARNING
- In this paper, we introduce a new task, **invariant function learning (IFL)**. Specifically, given the prior distribution of trajectories, denoted as $p(\mathbf{X})$, we consider a scenario where trajectories are

Table 1: Comparison of environments with previous works in the pendulum ODE system. We
 list examples from 2 environments, the pendulum states and coefficients distribution within one
 environment, and inference time targets. The blue factors are changing across different environments.
 Full function environments are provided in Appx. C.

Туре	Environment $e = 1$	Environment $e = 2$	Distribution	Inference
Coefficient environment Function environment	$\begin{vmatrix} f_1 = -\alpha_1^2 \sin(\theta_t) - \rho_1 \omega_t \\ f_1 = -\alpha^2 \sin(\theta_t) - \rho \omega_t \end{vmatrix}$	$f_2 = -\alpha_2^2 \sin(\theta_t) - \rho_2 \omega_t$ $f_2 = -\alpha^2 \sin(\theta_t) + \rho_{\frac{\omega_t}{ \omega_t }}$	$p(oldsymbol{ heta}_0,oldsymbol{\omega}_0) \ p(oldsymbol{ heta}_0,oldsymbol{\omega}_0,oldsymbol{lpha},oldsymbol{ ho})$	$f_3 = -\alpha_3^2 \sin(\theta_t) - \rho_3 \omega_t$ $f_c = -\alpha^2 \sin(\theta_t)$

117 observed under multiple environments. The trajectories observed in an environment $e \in \mathcal{E}$ are 118 sampled from the conditional distribution $p(\mathbf{X}|\mathbf{e}=e)$. Given a trajectory X from environment e, our 119 goal is to discover its invariant function f_c , which generates the invariant trajectory X^c . f_c and X^c 120 only include the shared mechanisms across all environments, thus capturing the underlying natural 121 laws unaffected by environmental factors. For instance, as illustrated in Fig. 1, in the case of a 122 pendulum system with varying environmental effects such as frictions or power, the goal is to extract 123 the natural motion of an ideal pendulum when excluding these external influences. A set of specific 124 examples for the task is provided below, more examples are available in Appx. C.

125 **Function environments.** The environments in this paper are different from those defined in 126 CoDA (Kirchmeyer et al., 2022), LEAD (Yin et al., 2021a), and MetaphysiCa (Mouli et al., 2024). 127 As shown in Tab. 1, the environments/tasks used in previous works, namely, coefficient environments, 128 are defined by the changes on the function coefficients α and ρ , *i.e.*, each environment contains only 129 one function while different environments include functions with different coefficients. In contrast, 130 we consider more complex cases and define environments as the interventions on function forms, *i.e.*, 131 each environment can contain functions with the same function form and different coefficients, while 132 different environments differ in function forms. Specifically, in Tab. 1, while coefficient environment 1 consists of a single function f_1 , our function environment 1 includes all the functions in form of 133 $-\alpha^2 \sin(\theta_t) - \rho \omega_t$ where $\alpha \sim p(\alpha), \rho \sim p(\rho)$. Here, we model these functions as a function random 134 variable f_1 . 135

136 **Challenges.** However, extracting such invariant dynamics presents two significant challenges. Firstly, 137 invariant mechanisms are intertwined with varying initial conditions, system parameters, and time, 138 making them particularly difficult to isolate or define. Secondly, in dynamical systems, state values 139 and their derivatives evolve over time, meaning that there is no single *invariant representation* fixed across time steps. This requires defining invariant factors in a function space, where functions can 140 cover dynamical states. Conventional invariant learning techniques are not directly applicable in this 141 context, as they are typically not designed to capture invariance in function spaces. These challenges 142 necessitate the development of new function representations and the development of a novel invariant 143 learning principle tailored to dynamical systems. 144

145 146

147

2.2.1 CAUSALITY-BASED DEFINITIONS

148 First challenge: causal formulation. In light of the first challenge, we aim to formulate the invariant 149 function learning and the dynamical system forecasting problem from a causal perspective. As 150 shown in Fig. 2, we formulate the trajectory data generation process as a Structural Causal Model 151 (SCM) (Pearl, 2009), where c, e, and X_0 are exogenous variables. All endogenous variables, except 152 the function composition step $f_c, f_e \to f$, are generated with extra random noises to model complex 153 real-world scenarios. Please refer to Appx. B.1 for more details. The optimization goal of the 154 forecasting task is to estimate the true distribution $p(\mathbf{X})$. As can be observed from the causal graph, our function learning framework can be considered as two phases, namely, function prediction and 155 *forecasting*. For the function prediction phase, similar to inverse problems (Lu et al., 2021b), given 156 the observed trajectory \mathbf{X}_{p} , the target is to reversely infer the invariant derivative function f that can 157 represent the dynamics of the system, *i.e.*, fitting $p(f|\mathbf{X}_p)$. Intuitively, taking Fig. 1 as an example, this 158 phase aims at the reasoning of the function basis $\sin(\theta_t)$, $-\omega_t$, $\frac{\omega_t}{|\omega_t|}$, and the coefficients α, ρ . After 159 obtaining the derivative function \hat{f} , the forecasting phase feeds it into a numerical integrator with the 160 initial condition X_0 for the X forecasting, which can be demonstrated as $p(X|f, X_0)$. Note that the 161 bold font variables are random variables instead of individual realizations.

162 Second challenge: function disentanglement. To han-163 dle the second challenge of extracting invariant mecha-164 nisms in dynamical systems, we aim to define invariant 165 representation in the function space, which requires the 166 access to intermediate functions and the disentanglement formulations. With the two-phase function prediction to 167 forecasting process, we explicitly expose the derivative 168 function, allowing us to isolate the function prediction process and disentangle the learning of invariant functions. 170 In the invariant function learning problem, we decompose 171 the exposed function variable f into an invariant function 172 variable f_c and an environment-specific function variable 173 f_e , as in Fig. 2. These two functions are caused by the 174 exogenous factors c and e, respectively. Intuitively, the c 175 variable includes the common and invariant mechanism, 176 while e is the environment variable determined by the observational environment of the trajectory, e.g., the pen-177 dulum dynamics can be observed in different mediums 178 (environments), such as air and water. 179



Figure 2: **Structural causal model.** The causal data generation process includes two phases: function generation and trajectory generation, which correspond to our two learning phases in parentheses, namely, function prediction and forecasting. The gray nodes in the causal graph indicate observable variables.

Formally, the invariant function learning target is the dis-

181 covery of f_c , which eliminates the effect of environments and obtains the invariant mechanism. 182 However, from the d-separation perspective (Pearl, 2009), since X_p is the descendent of the collider f 183 of f_c and f_e , given X_p , f_c and f_e are correlated/biased and not distinguishable, preventing the direct fit-184 ting of f_c given X_p . Therefore, we aim to identify this intermediate hidden variable by characterizing 185 the unique properties of the prior distribution of $p(f_c)$. Theoretically, since f is the collider between f_c 186 and e, f_c is expected to be independent of e. In addition, among all the functions that are independent 187 of e, f_c should be the most informative with respect to the observed trajectories, which will be proved 188 in Thm. 3.1. This forms the foundation of invariant function learning.

- 189
- 190 191

3 A METHOD: DISENTANGLEMENT OF INVARIANT FUNCTION

Following the two-phase function learning framework, we now propose the first method for IFL, Disentanglement of Invariant Function (DIF), with hypernetwork-based implementations of the two corresponding networks. For the forecasting network, similar to traditional representation learning tasks, we aim to learn a function $f \in \mathcal{F} : \mathbb{R}^d \mapsto \mathbb{R}^d$. For the function prediction, however, it requires learning a function that returns a function, $h \in \mathcal{H} : \mathbb{R}^{d \times T_c} \mapsto \mathcal{F}$, *i.e.*, learning a hyper-function, which is enabled using a hypernetwork.

198 199

200

3.1 Hypernetwork Design

201 **Function prediction.** To quantify the objective of the hyper-function, we approximate its output 202 function as a neural network with m parameters. The function space \mathcal{F} consists of all possible 203 neural networks with m parameters, and a function $f \in \mathcal{F}$ can be represented as a vector in \mathbb{R}^m . 204 Thus this parameterization process introduces a hypernetwork structure (Ha et al., 2016) into the 205 implementation, as shown in Fig. 3. Note that since our parameterization transfers functions into the 206 real number space, it is now possible to apply invariant learning techniques such as IRM (Arjovsky et al., 2019) and VREx (Krueger et al., 2021), where invariant (function) representations need to be 207 extracted. In the following sections, we use \mathcal{F} and omit \mathbb{R}^m for simplicity. 208

Practically, since the number of parameters in a network is generally large, we propose to compress the invariant function representations into hidden representations, thus forming an encoder-decoder framework. Specifically, as shown in the Fig. 3, the trajectory encoder is a transformer-based network with positional embedding design, denoted as $h_{\theta_{enc}} : \mathbb{R}^{d \times T_c} \mapsto \mathbb{Z}$. Given the hidden representation from the encoder, we further encode an invariant function embedding $\hat{z}_c \in \mathbb{Z}$ and an environment function embedding $\hat{z}_e \in \mathbb{Z}$ using two multilayer perceptrons (MLPs), denoted as $h_{\theta_{inv}}$ and $h_{\theta_{env}}$, respectively. Then, aligning with our causal graph as Fig. 2, we combine \hat{z}_c and \hat{z}_e by summing them as the function representation $\hat{z} \in \mathbb{Z}$, which can be used for full dynamics prediction. Finally, we



Forecasting. Given the produced neural network function f, we apply a numerical integrator as our forecasting. Given the produced neural network function f, we apply a numerical integrator as our forecasting, a function g_{int} that takes a derivative function \hat{f} and initial states X_0 as inputs, to obtain $\hat{X} = g_{int}(\hat{f}, X_0) + \epsilon$ where ϵ is sampled from a Gaussian noise $\mathcal{N}(\mathbf{X}; 0, \sigma^2 I)$ introduced by calculation deviations. This forecasting formulation enables the following probability modeling, where we obtain the forecasting given realizations X_0 and \hat{f} as a Gaussian distribution $\mathcal{N}(\mathbf{X}; g_{int}(\hat{f}, X_0), \sigma^2 I)$ denoted as $\mathbf{x}(\mathbf{X}|\hat{f}, X_0)$. Therefore, in probability modeling, \hat{X} is sampled from $\mathbf{n}(\mathbf{X}|\hat{f}, X_0)$. It is worth

denoted as $p(\mathbf{X}|\hat{f}, X_0)$. Therefore, in probability modeling, \hat{X} is sampled from $p(\mathbf{X}|\hat{f}, X_0)$. It is worth noting that unlike in inference time and analyses, it is time-consuming to use numerical integrators during training; therefore, we follow Mouli et al. (2024) and train the model by fitting the derivative $\frac{d\hat{X}_t}{dt} = \hat{f}(X_t)$ with numerical derivatives from the ground-truth X instead. For simplicity, we denote $\hat{f}(\cdot)$ as a neural network based derivative function with parameters $\hat{f} \in \mathbb{R}^m$.

244 245 3.2 DISCOVERY OF INVARIANT FUNCTION

253

254

255

256 257

258

259

With the above hypernetwork design, we can now propose the discovery of the invariant function f_c . Following the independence and information properties of f_c discussed in Sec. 2.2.1, we achieve invariant function learning with the following theorem, which is our main theoretical result.

Theorem 3.1 (Invariant function learning principle). Given the causal graph in Fig. 2, and the predicted function random variable $\hat{\mathbf{f}}_c = h_{\theta_c}(\mathbf{X}_p)$, it follows that the true invariant function random variable $\mathbf{f}_c = h_{\theta_c^*}(\mathbf{X}_p)$, where θ_c^* is the solution of the following optimization process, described as

$$\theta_c^* = \operatorname*{arg\,max}_{\theta_c} I(h_{\theta_c}(\mathbf{X}_p); \mathbf{f} | \mathbf{X}_0) \quad s.t. \quad h_{\theta_c}(\mathbf{X}_p) \perp \mathbf{I} \mathbf{e}, \tag{2}$$

where $I(\cdot; \cdot)$ is mutual information that measures the information overlap between the predicted invariant function random variable $h_{\theta_r}(\mathbf{X}_p)$ and the true full-dynamics function random variable f.

The proof is available in Appx. B.2. Thm. 3.1 establish guarantees and conditions for the function output \hat{f}_c of the hypernetwork to be the invariant function f_c , fulfilling the goal of the IFL task.

260 3.2.1 IMPLEMENTATION OF INVARIANT FUNCTION LEARNING PRINCIPLE

Following Thm. 3.1, next we introduce the implementation and optimization process of our proposed networks. We first train the encoder and decoder of our hypernetwork by approximating the trajectory distribution $p(\mathbf{X})$, parameterized as $p(\mathbf{X}|h_{\theta}(\mathbf{X}_p), \mathbf{X}_0)$, where we apply the cross-entropy minimization. Given that our supervision signals only come from the ground-truth trajectories, we introduce a simple lemma for our optimization processes. The proof is provided in Appx. B.3.

267 Lemma 3.2 (ODE cross-entropy minimization). *Given forecasting model* $p(\mathbf{X}|h_{\theta}(\mathbf{X}_{p}), \mathbf{X}_{0})$, *it follows that the cross-entropy minimization between the data distribution* $p(\mathbf{X})$ *and* $p(\mathbf{X}|h_{\theta}(\mathbf{X}_{p}), \mathbf{X}_{0})$ **269** *is equivalent to minimizing mean square error* $\min_{\theta} \mathbb{E}_{X \sim p} ||X - \hat{X}||_{2}^{2}$, where \hat{X} is sampled from $p(\mathbf{X}|h_{\theta}(X_{p}), X_{0})$. In order to discover invariant functions, we apply the invariant function learning principle, which requires maximizing the conditional mutual information between the predicted function random variable $\hat{\mathbf{f}}_c = h_{\theta_c}(\mathbf{X}_p)$ and f. Based on the derivation of Lemma 3.2, we have the following proposition.

Proposition 3.3 (ODE conditional mutual information maximization). Given forecasting model $p(\mathbf{X}|h_{\theta_c}(X_p), X_0)$, it follows that the conditional mutual information maximization max_{$\theta_c} I(h_{\theta_c}(\mathbf{X}_p); \mathbf{f}|\mathbf{X}_0)$ is equivalent to minimizing mean square error min_{$\theta_c} <math>\mathbb{E}_{\mathbf{X}\sim P} ||X - \hat{X}^c||_2^2$, where \hat{X}^c is the predicted trajectory sampled from $p(\mathbf{X}|h_{\theta_c}(X_p), X_0)$.</sub></sub>

277

278 The proof is provided in Appx. B.4. Lemma 3.2 and Prop. 3.3 essentially transfers the mutual information maximization of Thm. 3.1 into an implementable optimization of mean square error 279 (MSE) loss, enabling the practical use of the invariant function learning principle. In addition, 280 the independence constraint in Thm. 3.1 requires that the extracted functions should not be over-281 informative or contain biased information from environments e. This independence constraint can be 282 implemented in an adversarial way, where we require the environment prediction $P(\mathbf{e}|\mathbf{\hat{f}}_c)$ to be as 283 less informative as possible, *i.e.*, minimizing the mutual information of $I(e; f_c)$. Thus we introduce 284 an environment discriminator g_{ϕ} , a.k.a., $P_{\phi}(e|f)$, which aims to distinguish the environment of any 285 function from \mathcal{F} . The hypernetwork is trained adversarially to enforce \hat{f}_c as indistinguishable as 286 possible. The theoretically analysis of this independence training is provided in Appx. B.5. 287

Objective. The overall optimization objective can be obtained with three training strategies. First, the training of the discriminator is conducted on both \hat{f}_c and \hat{f}_e to fully capture environment patterns. Second, we use the corresponding hidden representations of \hat{f}_c and \hat{f}_e , \hat{z}_c and \hat{z}_e , as the input of the discriminator. Third, as we mentioned in Sec.3.1, during training, we fit derivatives instead of using an integrator.

293

295

296

$$\min_{\theta} \mathbb{E}_{\mathbf{X} \sim p} \sum_{t} \left\| \frac{dX_{t}}{dt} - \hat{f}(X_{t}) \right\|_{2}^{2} + \lambda_{c} \cdot \min_{\theta_{c}} \mathbb{E}_{\mathbf{X} \sim p} \sum_{t} \left\| \frac{dX_{t}}{dt} - \hat{f}_{c}(X_{t}) \right\|_{2}^{2} + \lambda_{dis} \cdot \left[\min_{\phi} -\mathbb{E}_{\mathbf{X} \sim p} \log g_{\phi}(\hat{z}_{c}) + \min_{\phi, \hat{\theta}_{c}} -\mathbb{E}_{\mathbf{X} \sim p} \log g_{\phi}(\hat{z}_{c}) \right] + \lambda_{adv} \cdot \max_{\hat{\theta}_{c}} -\mathbb{E}_{\mathbf{X} \sim p} \log g_{\phi}(\hat{z}_{c})$$
(3)

297 298

300

301

302

303 304

305 306 307

308

299 where we denote $\bar{\theta}_e = \{\theta_{enc}, \theta_{env}\}; \bar{\theta}_c = \{\theta_{enc}, \theta_{inv}\}$. Please refer to Appx. D.2.2 for more details.

Efficient hypernetwork implementation. Last but not least, one of the major challenges that limits the usage of hypernetworks is the implementation complexity. In this work, we propose a reference-based hypernetwork implementation to accelerate the running speed using only PyTorch Paszke et al. (2019) without re-implementing basic neural networks. The speedup compared to the naïve implementation and the vectorized functional implementation are 16.8x and 2x, respectively. Please refer to Appx. F for implementation and experimental details.

4 RELATED WORK

This work is inspired by the ideas and limitations of previous research in dynamical system forecasting,meta-learning, and invariant learning.

311 Deep learning models are widely applied in many physical applications (Lusch et al., 2018; Yeo & 312 Melnyk, 2019; Kochkov et al., 2021; Chen et al., 2018) including partial differential equations (PDEs) 313 with the focus on the multi-scale (Li et al., 2020; Stachenfeld et al., 2021), multi-resolution (Kochkov 314 et al., 2021; Wu et al., 2022), and long-term stability (Li et al., 2021; Lippe et al., 2023) issues. 315 Operator learning and neural operators (Gupta et al., 2021; Kovachki et al., 2023) are popular for PDE estimations. Although the ODE dynamical system does not contain the multi-scale problem that 316 Fourier neural operator (Kovachki et al., 2023) tried to solve, our framework can be considered as a 317 kind of operator learning. 318

Meta-learning methods (Finn et al., 2017; Rusu et al., 2018; Li et al., 2017; Zintgraf et al., 2019;
Perez et al., 2018) aim to learn meta-parameters that can be used across multiple tasks, where the
meta parameters are generally learned to make rapid adaptations. In previous meta-learning studies
on dynamical systems (Kirchmeyer et al., 2022; Wang et al., 2022; Yin et al., 2021a), the objective
was to find a meta-function that could quickly adapt to multiple new systems, where hypernetworks
are only employed as low-rank adaptors for new dynamical system trajectories, similar to the idea

of LoRA (Hu et al., 2021). Our work differs from such meta-learning approaches in two key ways. First, we focus on discovering invariant functions rather than quickly adaptable ones. Second, while meta-learning methods seek to learn a single meta-function, our framework learns multiple functions, represented by an invariant function random variable f_c . This distinction stems from our more complex environment definition, detailed in Sec. 2.2. From another aspect, learning an invariant function distribution instead of a single function can be considered as generalized meta-learning with an invariant function learning goal.

331 Current invariant learning methods (Arjovsky et al., 2019; Lu et al., 2021a; Rosenfeld et al., 2020; 332 Krueger et al., 2021; Sagawa et al., 2019) follow the framework of invariant risk minimization 333 (IRM) (Arjovsky et al., 2019), which was inspired by invariant causal predictor (Peters et al., 2016). 334 This invariant learning framework aims to learn a hidden invariant representation that generalizes across multiple environments, ensuring out-of-distribution performance. However, this approach 335 cannot work on dynamical forecast tasks due to the lack of invariant function definition and the 336 violation of the categorical data assumption. To be more specific, first, invariant functions cannot 337 be naturally defined in the real number vector space. Second, invariant learning commonly assumes 338 the prediction results are categorical, where a single invariant representation can fully determine the 339 corresponding label. However, this assumption is violated in dynamical system forecasting, where 340 the invariant mechanism is only partially responsible for the output. In this case, the IRM principle 341 can not hold even when the invariant function ground truth is provided. To address the issues, we 342 introduce the causal assumption (see Fig. 2) that defines the invariant function space, and propose the 343 corresponding invariant function learning principle and implementation. 344

A related line of research involves symbolic regression for ordinary differential equations (ODEs),
 where transformer models have shown significant success (Becker et al., 2023; d'Ascoli et al., 2023;
 Seifner et al., 2024). While these approaches primarily focus on deriving symbolic expressions for individual trajectories, rather than identifying invariant functions across groups of trajectories, exploring the interplay between these two directions presents an exciting avenue for future research.

5 EXPERIMENTS

We conduct experiments to address the following research questions. **RQ1**: Are existing metalearning and invariant learning techniques effective for extracting invariant functions? **RQ2**: Can the proposed invariant function learning principle outperform baseline techniques? **RQ3**: How do the full functions f and the invariant functions f_c differ in performance? **RQ4**: Are the extracted invariant functions explainable and aligned with the true invariant mechanisms? **RQ5**: How will performance change given different lengths of inputs and types of environments? (See Appx. E) **RQ6**: Is the proposed hypernetwork implementation more efficient than previous implementations? (See Appx. F)

359 360

361

350

351

5.1 DATASETS

362 In our experiments, we introduce three multi-environment (ME) datasets, ME-Pendulum, ME-Lotka-363 Volterra, and ME-SIREpidemic. These three datasets are generated by simulators modified from the DampedPendulum (Yin et al., 2021b), Lotka-Volterra (Ahmad, 1993), and SIREpidemic (Wang et al., 364 2021). Specifically, each of the dataset's training sets includes four environments with 200 samples 365 for each environment. Specifically, each environment corresponds to one specific environmental 366 effect. ME-Pendulum contains three types of friction and one effect with external energy. ME-Lotka-367 Volterra modified the common predatory relationship into four modified relationships, e.g., adding 368 resource limits. ME-SIREpidemic produces four conceptual epidemiology models with the same 369 susceptible population to infected population relationship. In addition to the training set, we generate 370 200 samples with 50 samples for each environment as an in-distribution test set. Please refer to 371 Appx. C for more details. 372

373 374

5.2 EXPERIMENTAL SETUP

To quantitatively evaluate the invariant function extraction performance, we need to remove the environment-related effects to generate invariant trajectories X^c as the invariant function groundtruth, *e.g.*, we simulate new data by eliminating $-\rho\omega_t$ from $-\alpha^2 \sin(\theta_t) - \rho\omega_t$ in the ME-Pendulum dataset (Fig. 1). To be more specific, a generated invariant trajectory X^c , aligning the causal graph, has the same system parameters as the corresponding biased trajectory X for the invariant part controlled by c, *i.e.*, they have the same α in the pendulum example. This invariant trajectory generation is being done on the in-distribution test set so that each trajectory X in this test set has its special corresponding invariant trajectory ground truth X^c .

This test set design enables us to mimic the situation of scientific discoveries, where we only observe environment-biased data but are required to find and evaluate invariant function candidates. Specifically, given the biased X_p , the hypernetwork is supposed to predict the corresponding invariant derivative function \hat{f}_c . Then, with a numeral integrator, the output of this invariant forecaster \hat{X}^c will be evaluated by comparing with the corresponding invariant trajectory ground truth X^c using normalized root mean square error (NRMSE).

388 389

390

5.3 PROPOSED META-LEARNING AND INVARIANT LEARNING BASELINES

General dynamical system forecasting is different from invariant function learning significantly,
 where they focus on how to adapt to new trajectories, which commonly requires further optimization,
 e.g., test-time adaptation (Mouli et al., 2024) or adaptations with meta information (Wang et al.,
 2022; Kirchmeyer et al., 2022). Unfortunately, under our scientific discovery setting, there is no extra
 information provided at test time, making them inapplicable to this setting. Therefore, we construct 4
 new baseline settings by transplanting the techniques of previous meta-learning and invariant learning
 to our proposed framework detailed in Appx. D.1.

We first adopt the meta-learning baseline MAML (Finn et al., 2017), where we use its learned metaparameters for invariant learning to evaluate whether the fastest adapted parameter is the invariant function parameter. Our second meta-learning baseline is CoDA (Kirchmeyer et al., 2022), where we replace its hypernetwork decoder with the full encoder-decoder hypernetwork in our framework to fit in our task. Aligning with the original CoDA paper, we set the dimension of the hidden representation to be 2. Similar to MAML, we eliminate the adaptation part and use only the learned meta-parameter for invariant state prediction.

For invariant learning baselines, we adopted the two most typical techniques, IRM (Ahuja et al., 2021) and VREx (Krueger et al., 2021). These two techniques are applied to the proposed framework, where IRM stands for the most typical definition of invariant learning, while VREx stands for the distributionally robust optimization baseline, which can be considered as the generalization of GroupDRO (Sagawa et al., 2019).

409 410 411

5.4 QUANTITATIVE RESULTS

Similar to other scientific discovery tasks, such as drug discovery, constructing a proper validation set is challenging. Instead, with only observational data available, we generate invariant function candidates that can be further validated in real experimental settings, *e.g.*, through the introduction of interventions (Pearl, 2009).

416 To quantitatively compare these methods, we provide the corresponding hyper-parameter search 417 spaces for each technique in Appx. D and plot the results of random hyper-parameter sampling as 418 distributions using Boxen plots. As shown in Fig. 4, we compare the quality of the invariant function 419 candidates based on their median, best result, and quantiles. Specifically, the median performance 420 of our proposed method surpasses the middle candidates of all other approaches. The performance 421 gaps are particularly notable on the ME-Pendulum and ME-SIR-Epidemic datasets. For example, 422 on ME-Pendulum, over 75% of our method's candidates outperform the best results of MAML and 423 CoDA, and more than 93.75% candidates of IRM and VREx. On ME-Lotka-Volterra, the median of our candidates still outperforms nearly all candidates from other methods. In addition, as shown in 424 the visualizations on ME-Pendulum 5, our learned invariant function f_c eliminates environmental 425 resistances from the original trajectory (Fig. 5a) and obtain a simple pendulum motion without 426 attenuation (Fig. 5b). Both quantitative and visualization results demonstrate the superior capability 427 of our method in extracting invariant functions (RQ2). 428

To address the first research question (RQ1), we observe that the invariant learning techniques, IRM and VREx, are generally more stable than the meta-learning baselines. Although IRM and VREx do not surpass MAML on ME-Lotka-Volterra, they outperform MAML on 2 out of 3 datasets and are consistently better than CoDA. However, when compared to our proposed method, the best function



Figure 4: Invariant trajectory prediction errors on 5 methods under 3 multi-environment ODE systems. For each method, we provide model candidates with 80+ random hyper-parameter selections in their searching spaces, *i.e.*, more than 1200 models in the figure.



Figure 5: Visualization of trajectory predictions on ME-Pendulum Table 2: Invariant function validation and symbolic regression. NAN denotes that the result is not applicable or not of interest.

Target	Function		ME-Pendulum		ME-Lotka-Volterra		ME-SIREpidemic	
Iuger		NRMSE	SR Explanation	NRMSE	SR Explanation	NRMSE	SR Explanation	
W ^C	$\hat{\mathbf{f}}_c$	0.3561	$\frac{d\theta_t}{dt} = 0.99\omega_t$ $\frac{d\omega_t}{dt} = -0.97\alpha^2\sin\left(\theta_t\right)$	0.6194	$\frac{dp_t}{dt} = 1.254p_t - 0.38q_t p_t$ $\frac{dq_t}{dt} = 4.1p_t - 0.30q_t - \gamma$	0.0652	$\frac{dS_t}{dt} = -1.7S_t I_t$ $\frac{dI_t}{dt} = 0.42S_t I_t$ $\frac{dR_t}{dt} = -0.0088$	
X	f	0.7884	$\frac{\frac{d\theta_t}{dt} = \omega_t \cos\left(\frac{\omega_t}{e^{\alpha}}\right)}{\frac{d\omega_t}{dt} = \theta_t \alpha \left(-\alpha + \rho\right)}$	0.7919	$\frac{dp_t}{dt} = -0.76p_t$ $\frac{dq_t}{dt} = \frac{p_t}{0.36} - \gamma$	0.9867	$\left \begin{array}{l} \displaystyle \frac{dS_t}{dt} = -0.24\beta I_t S_t - \\ \displaystyle \frac{dI_t}{dt} = 0.40S_t \\ \displaystyle \frac{dR_t}{dt} = 0.66\gamma \end{array} \right $	
x	$\hat{\mathbf{f}}_c$	0.7994	NAN	0.6912	NAN	0.7641	NAN	
A		0.1700	NAN	0.3881	NAN	0.0212	NAN	
f	g GT	NAN	$\frac{d\theta_t}{dt} = \omega_t$ $\frac{d\omega_t}{dt} = -\alpha^2 \sin\left(\theta_t\right)$	NAN	$\frac{dp_t}{dt} = \alpha p_t - \beta p_t q_t$ $\frac{dq_t}{dt} = \delta p_t q_t - \gamma q_t$	NAN	$\begin{vmatrix} \frac{dS_t}{dt} = -\beta \frac{S_t I_t}{S_t + I_t + t} \\ \frac{dI_t}{dt} = \beta \frac{S_t I_t}{S_t + I_t + R} \\ \frac{dR_t}{dt} = 0 \end{vmatrix}$	

candidates from these invariant learning techniques are suboptimal. This confirms that the general invariant learning principles fall short in the context of invariant function extraction, aligning with the discussions in Sec. 4.

5.5 FULL FUNCTION V.S. INVARIANT FUNCTION

To analyze **RO3**, we benchmark our best invariant function learning models on the invariant state ground truth \mathbf{X}^{c} and the multi-environment state ground truth \mathbf{X} , comparing their results using the predicted invariant function f_c and the full function f. As shown in Tab. 2, the performance on \mathbf{X}^c using \hat{f}_c represents the core results of our invariant function learning approach. In contrast, the predictions



on X^c using \hat{f} serve as a baseline for the *ablation study*, where no invariant function learning principle is applied. Additionally, the prediction error on X using \hat{f}_c implies the environmental information eliminated by the invariant function learning principle, while the NRMSEs on X using \hat{f} reflect standard in-distribution (ID) test errors. We observe that the ME-Lotka-Volterra dataset is the most challenging, with an NRMSE of 0.3881 in the ID test. This result is consistent with general deep learning outcomes in Mouli et al. (2024), given that ME-Lotka-Volterra is more complex than its original version.

As expected, \hat{f} performs well on X, while \hat{f}_c excels in predicting X^c. In our ablation study, we compare the performance of invariant state predictions X^c across all datasets. The predicted invariant functions \hat{f}_c significantly outperform the full predicted functions \hat{f} in terms of NRMSE, validating the effectiveness of the proposed invariant function learning principle.Furthermore, we conduct more strict ablation study with independent \hat{f} and \hat{f}_c training in Appx. E.1.

498 499 500

5.6 SYMBOLIC REGRESSION EXPLANATION

501 Furthermore, to address **RQ4**, we analyze the extracted invariant functions f_c by applying symbolic 502 regression using PySR (Cranmer, 2023). As shown in Tab. 2, we compare the symbolic regression (SR) explanations of the extracted invariant functions f_c with the true invariant functions f_c . On 504 the ME-Pendulum dataset, the frictionless pendulum function is nearly perfectly extracted, with 505 $0.99\omega_t$ matching ω_t and $-0.97\alpha^2\sin(\theta_t)$ closely approximating the true $-\alpha^2\sin(\theta_t)$. Given the 506 complexity of the ME-Lotka-Volterra dataset, the extracted invariant functions \hat{f}_c are non-trivial 507 and significantly outperform the full function f. On the ME-SIREpidemic dataset, the near-perfect 508 NRMSE for the invariant state indicates that the invariant function must have been correctly extracted. However, although the expression for $\frac{dR_t}{dt}$ is correct, the extracted expressions for $\frac{dS_t}{dt}$ and $\frac{dI_t}{dt}$ do not precisely match the expected f_c . Specifically, the coefficient from $\frac{dS_t}{dt}$ does not equal the inverse 509 510 of the coefficient from $\frac{dI_t}{dt}$, though this discrepancy should be constant. These mismatches attribute 511 to the limitations of PySR given the large number of variables and samples.¹ Future work could 512 explore incorporating stronger inductive biases, similar to physics-informed machine learning (PIML) 513 methods (Mouli et al., 2024; Yin et al., 2021b; Cranmer et al., 2020), to address these challenges. 514 Please refer to Appx. E.3 for symbolic comparisons with all baseline. 515

516 517

6 LIMITATIONS

While this work provides a foundation for invariant function learning in dynamical systems, several limitations and opportunities for future exploration remain. These include extending the framework to more complex entanglements, exploring applications in PDE systems, strengthening theoretical guarantees with advanced error bounds, and developing comprehensive benchmarks. Additionally, broader applications such as generalizable physics learning and foundational model development represent exciting directions for further research. Please refer to Appx. I for more discussions.

7 CONCLUSION

525 526 527

In this work, we target addressing the challenge of the invariant mechanism discovery in ODE 528 dynamical systems by extending invariant learning into function spaces. We introduce a new task, invariant function learning, which aims to extract the invariant dynamics across all environments with 529 different environment-specific function forms. We design a causal analysis based disentanglement 530 framework DIF to expose the underlying invariant functions. Additionally, we propose an invari-531 ant function learning principle with theoretical guarantees to optimize the framework and ensure 532 effective invariant function discovery. Our experiments, including invariant trajectory validations, 533 visualizations, ablation studies, and symbolic regression analyses, demonstrate the effectiveness of 534 our method. Finally, as discussed in Sec. 6, the introduced invariant function learning task has wide application scenarios and many challenges remain to be addressed. We expect that our work will 536 shed light on numerous future explorations in this field.

¹Superior PySR explanations indicate great invariant function learning results, but effective invariant function learning results might not lead to good PySR explanations.

540 REPRODUCIBILITY STATEMENT 541

To make sure the reproducibility of this work, for our theoretical results, all assumptions and proofs
are included in Appx. B. For the datasets, the system parameters, simulator expressions, and datasets
visualizations are provided in Appx. C. For the model, we provide full details including all the training
setup, architecture, objectives, and hyper-parameter searching spaces in Appx. D. The finalized code
will be released upon acceptance.

548 **REFERENCES**

547

549

550

551

- Shair Ahmad. On the nonautonomous volterra-lotka competition equations. *Proceedings of the American Mathematical Society*, 117(1):199–204, 1993.
- Kartik Ahuja, Ethan Caballero, Dinghuai Zhang, Jean-Christophe Gagnon-Audet, Yoshua Bengio,
 Ioannis Mitliagkas, and Irina Rish. Invariance principle meets information bottleneck for out-of distribution generalization. *Advances in Neural Information Processing Systems*, 34:3438–3450,
 2021.
 - Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization.
 arXiv preprint arXiv:1907.02893, 2019.
- Alex Aussem. Dynamical recurrent neural networks towards prediction and modeling of dynamical systems. *Neurocomputing*, 28(1-3):207–232, 1999.
- Sören Becker, Michal Klein, Alexander Neitz, Giambattista Parascandolo, and Niki Kilbertus.
 Predicting ordinary differential equations with transformers. In *International Conference on Machine Learning*, pp. 1978–2002. PMLR, 2023.
- Olivier Catoni. Pac-bayesian supervised classification: the thermodynamics of statistical learning.
 arXiv preprint arXiv:0712.0248, 2007.
- Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary
 differential equations. *Advances in neural information processing systems*, 31, 2018.
- Kyunghyun Cho. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- 572 Miles Cranmer. Interpretable machine learning for science with pysr and symbolic regression. jl.
 573 arXiv preprint arXiv:2305.01582, 2023.
- 574
 575
 576
 576
 576
 577
 576
 576
 577
 577
 578
 579
 579
 570
 570
 570
 571
 572
 573
 574
 574
 574
 574
 574
 575
 576
 577
 577
 576
 577
 576
 577
 577
 578
 578
 579
 579
 579
 570
 570
 571
 572
 574
 575
 574
 575
 576
 577
 577
 578
 578
 578
 579
 579
 579
 570
 570
 571
 572
 572
 574
 574
 574
 575
 575
 576
 577
 577
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 579
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
- Stéphane d'Ascoli, Sören Becker, Alexander Mathis, Philippe Schwaller, and Niki Kilbertus.
 Odeformer: Symbolic regression of dynamical systems with transformers. *arXiv preprint arXiv:2310.05573*, 2023.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of
 deep networks. In *International conference on machine learning*, pp. 1126–1135. PMLR, 2017.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François
 Laviolette, Mario March, and Victor Lempitsky. Domain-adversarial training of neural networks. *Journal of machine learning research*, 17(59):1–35, 2016.
- Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. Convolutional
 sequence to sequence learning. In *International conference on machine learning*, pp. 1243–1252.
 PMLR, 2017.
- Dimitrios Giannakis. Data-driven spectral decomposition and forecasting of ergodic dynamical systems. *Applied and Computational Harmonic Analysis*, 47(2):338–396, 2019.
- 593 Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.

594 595 596	Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. <i>Communications of the ACM</i> , 63(11):139–144, 2020.
597 598 599	Gaurav Gupta, Xiongye Xiao, and Paul Bogdan. Multiwavelet-based operator learning for differential equations. <i>Advances in neural information processing systems</i> , 34:24048–24062, 2021.
600	David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. arXiv preprint arXiv:1609.09106, 2016.
601 602 603 604	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> , 2021.
605 606 607	Matthieu Kirchmeyer, Yuan Yin, Jérémie Donà, Nicolas Baskiotis, Alain Rakotomamonjy, and Patrick Gallinari. Generalizing to new physical systems via context-informed dynamics model. In <i>International Conference on Machine Learning</i> , pp. 11283–11301. PMLR, 2022.
608 609 610 611	Dmitrii Kochkov, Jamie A Smith, Ayya Alieva, Qing Wang, Michael P Brenner, and Stephan Hoyer. Machine learning–accelerated computational fluid dynamics. <i>Proceedings of the National Academy</i> of Sciences, 118(21):e2101784118, 2021.
612 613 614	Nikola Kovachki, Zongyi Li, Burigede Liu, Kamyar Azizzadenesheli, Kaushik Bhattacharya, Andrew Stuart, and Anima Anandkumar. Neural operator: Learning maps between function spaces with applications to pdes. <i>Journal of Machine Learning Research</i> , 24(89):1–97, 2023.
615 616 617	David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghuai Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (REx). In <i>International Conference on Machine Learning</i> , pp. 5815–5826. PMLR, 2021.
619 620	Zhenguo Li, Fengwei Zhou, Fei Chen, and Hang Li. Meta-sgd: Learning to learn quickly for few shot learning. <i>CoRR</i> , abs/1707.09835, 2017. URL http://arxiv.org/abs/1707.09835.
621 622 623	Zongyi Li, Nikola Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Kaushik Bhattacharya, Andrew Stuart, and Anima Anandkumar. Fourier neural operator for parametric partial differential equations. <i>arXiv preprint arXiv:2010.08895</i> , 2020.
625 626 627	Zongyi Li, Miguel Liu-Schiaffini, Nikola Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Kaushik Bhattacharya, Andrew Stuart, and Anima Anandkumar. Learning dissipative dynamics in chaotic systems. <i>arXiv preprint arXiv:2106.06898</i> , 2021.
628 629 630	Phillip Lippe, Bastiaan S Veeling, Paris Perdikaris, Richard E Turner, and Johannes Brandstetter. Pde- refiner: Achieving accurate long rollouts with neural pde solvers. <i>arXiv preprint arXiv:2308.05732</i> , 2023.
631 632 633 634	Chaochao Lu, Yuhuai Wu, José Miguel Hernández-Lobato, and Bernhard Schölkopf. Invariant causal representation learning for out-of-distribution generalization. In <i>International Conference on Learning Representations</i> , 2021a.
635 636 637	Lu Lu, Raphael Pestourie, Wenjie Yao, Zhicheng Wang, Francesc Verdugo, and Steven G Johnson. Physics-informed neural networks with hard constraints for inverse design. <i>SIAM Journal on Scientific Computing</i> , 43(6):B1105–B1132, 2021b.
638 639 640	Bethany Lusch, J Nathan Kutz, and Steven L Brunton. Deep learning for universal linear embeddings of nonlinear dynamics. <i>Nature communications</i> , 9(1):4950, 2018.
641 642	David A McAllester. Some pac-bayesian theorems. In <i>Proceedings of the eleventh annual conference</i> on Computational learning theory, pp. 230–234, 1998.
643 644 645	S Chandra Mouli, Muhammad Alam, and Bruno Ribeiro. Metaphysica: Improving ood robustness in physics-informed machine learning. In <i>The Twelfth International Conference on Learning Representations</i> , 2024.
647	Jose Javier Gonzalez Ortiz, John Guttag, and Adrian Dalca. Non-proportional parametrizations for stable hypernetwork learning. <i>arXiv:2304.07645</i> , 2023.

- 648 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor 649 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, 650 high-performance deep learning library. Advances in neural information processing systems, 32, 651 2019. 652 Judea Pearl. Causality. Cambridge university press, 2009. 653 654 Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual 655 reasoning with a general conditioning layer. In Proceedings of the AAAI conference on artificial intelligence, volume 32, 2018. 656 657 Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant 658 prediction: identification and confidence intervals. Journal of the Royal Statistical Society: Series 659 B (Statistical Methodology), 78(5):947–1012, 2016. 660 Elan Rosenfeld, Pradeep Ravikumar, and Andrej Risteski. The risks of invariant risk minimization. 661 *arXiv preprint arXiv:2010.05761*, 2020. 662 663 Andrei A Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osin-664 dero, and Raia Hadsell. Meta-learning with latent embedding optimization. arXiv preprint 665 arXiv:1807.05960, 2018. 666 Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust 667 neural networks for group shifts: On the importance of regularization for worst-case generalization. 668 arXiv preprint arXiv:1911.08731, 2019. 669 670 Patrick Seifner, Kostadin Cvejoski, Antonia Körner, and Ramsés J. Sánchez. Foundational inference 671 models for dynamical systems, 2024. URL https://arxiv.org/abs/2402.07594. 672 Satinder Singh, Michael James, and Matthew Rudary. Predictive state representations: A new theory 673 for modeling dynamical systems. arXiv preprint arXiv:1207.4167, 2012. 674 675 Kimberly Stachenfeld, Drummond B Fielding, Dmitrii Kochkov, Miles Cranmer, Tobias Pfaff, 676 Jonathan Godwin, Can Cui, Shirley Ho, Peter Battaglia, and Alvaro Sanchez-Gonzalez. Learned coarse models for efficient turbulence simulation. In International Conference on Learning 677 Representations, 2021. 678 679 Shyam Sudhakaran Sudhakaran. hyper-nn. https://github.com/shyamsn97/hyper-nn, 2022. 680 Vladimir N Vapnik and A Ya Chervonenkis. On the uniform convergence of relative frequencies of 681 events to their probabilities. In Measures of complexity: festschrift for alexey chervonenkis, pp. 682 11-30. Springer, 2015. 683 684 A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017. 685 Johannes von Oswald, Christian Henning, Benjamin F. Grewe, and João Sacramento. Continual 686 learning with hypernetworks. In International Conference on Learning Representations, 2020. 687 URL https://arxiv.org/abs/1906.00695. 688 689 Rui Wang, Danielle Maddix, Christos Faloutsos, Yuyang Wang, and Rose Yu. Bridging physics-based 690 and data-driven modeling for learning dynamical systems. In Learning for dynamics and control, 691 pp. 385–398. PMLR, 2021. 692 Rui Wang, Robin Walters, and Rose Yu. Meta-learning dynamics forecasting using task inference. 693 Advances in Neural Information Processing Systems, 35:21640–21653, 2022. 694 695 Wen-Xu Wang, Ying-Cheng Lai, and Celso Grebogi. Data based identification and prediction of nonlinear and complex dynamical systems. *Physics Reports*, 644:1–76, 2016. 696 697 Tailin Wu, Takashi Maruyama, Qingqing Zhao, Gordon Wetzstein, and Jure Leskovec. Learning controllable adaptive simulation for multi-scale physics. In NeurIPS 2022 AI for Science: Progress 699 and Promises, 2022. URL https://openreview.net/forum?id=PhktEpJHU3. 700
- 701 Kyongmin Yeo and Igor Melnyk. Deep learning algorithm for data-driven simulation of noisy dynamical system. *Journal of Computational Physics*, 376:1212–1231, 2019.

Yuan Yin, Ibrahim Ayed, Emmanuel de Bézenac, Nicolas Baskiotis, and Patrick Gallinari. Leads: Learning dynamical systems that generalize across environments. Advances in Neural Information Processing Systems, 34:7561-7573, 2021a. Yuan Yin, Vincent Le Guen, Jérémie Dona, Emmanuel de Bézenac, Ibrahim Ayed, Nicolas Thome, and Patrick Gallinari. Augmenting physical models with deep networks for complex dynamics forecasting. Journal of Statistical Mechanics: Theory and Experiment, 2021(12):124012, 2021b. Luisa Zintgraf, Kyriacos Shiarli, Vitaly Kurin, Katja Hofmann, and Shimon Whiteson. Fast context adaptation via meta-learning. In International Conference on Machine Learning, pp. 7693–7702. PMLR, 2019.

Appendix of Discovering Physics Laws of Dynamical Systems via Invariant Function Learning

CONTENTS

763	A	A Notations						
764 765	B	Inva	riant function learning foundation	16				
766		B .1	Structural causal model	16				
767		B.2	Proof of Invariant function learning principle	17				
769		B.3	Proof of ODE cross-entropy minimization	18				
770 771		B.4	Proof of ODE conditional mutual information maximization	19				
772		B.5	Theoretical Justification for Adversarial Training	20				
774	С	Data	asets	21				
775 776		C.1	Basic setup	21				
777		C.2	ME-Pendulum	21				
778		C.3	ME-Lotka-Volterra	21				
779		C.4	ME-SIREpidemic	22				
781		0.1						
782	D	Exp	erimental details	23				
784		D.1	Baselines	24				
785		D.2	Disentanglement of Invariant Function setup	24				
786 787			D.2.1 Architecture	24				
788			D.2.2 Training objectives	25				
789			D.2.3 Metric	25				
790 791		D.3	Software and Hardware	26				
792 793	Е	Sup	plementary experiments	26				
794 795		E.1	Ablation study	26				
796		E.2	Input length and environment analysis	26				
797 798		E.3	Symbolic regression explanation comparisons	28				
799 800	F	Effic	cient Hypernetwork Implementation	28				
801		F.1	Efficiency comparisons	28				
803	G	Sup	plimentary explanations	29				
804 805	•••	T 74		•••				
806	Н	Visu	alizations on ME-Lotka-Volterra and ME-SIREpidemic	29				
807 808	Ι	Lim	itations and Future Work	30				
500								

810 NOTATIONS А 811

812

813 814 For the ease of reading, we providing a table including major notations below for reference.

814	Table 3: Notation table					
815						
816	Notation	Explanation				
817						
818	x_t	Dynamical system states at time t				
819	\mathcal{X}	Dynamical system state space				
820	TX	langent state space				
821	R	Real number space				
822		The number of state				
823	ι	The first future time step				
824		A trajectory: A state matrix with T step states				
825		A indjectory, A state matrix with 1-step states Dest states before time step T				
225	X_p Y^c	An invariant trajectory				
020	\hat{X} \hat{Y} \hat{Y} \hat{Y}^c	The predicted trajectory of $V_{-}V_{-}$ and V^{c} (by a model)				
027	$egin{array}{c} {oldsymbol{\Lambda}},{oldsymbol{\Lambda}}_p,{oldsymbol{\Lambda}}\\ {oldsymbol{X}}{oldsymbol{X}}{oldsymbol{X}}{oldsymbol{X}}^c \end{array}$	The matrix-valued random variable of X and X^c				
020	$\hat{\mathbf{X}}, \hat{\mathbf{X}}_p, \hat{\mathbf{X}}$ $\hat{\mathbf{X}} \ \hat{\mathbf{X}} \ \hat{\mathbf{X}}^c$	The predicted matrix valued random variable of X , X_p , and X				
829	$\mathbf{M}, \mathbf{M}_p, \mathbf{M}_p$	X^c				
830	f	A derivative function (of an underlying dynamical system)				
831	f_{c}	An invariant derivative function				
832	f_{e}	An environment derivative function				
833	\hat{f}_{i} , \hat{f}_{a} , \hat{f}_{a}	The predicted functions of f_1 , f_2 , and f_3 (by a model)				
834	f. fc. fe	The derivative function random variable of f , f_c , and f_e				
835	\hat{f}_{-} \hat{f}_{-}	The predicted derivative function random variable of f_{1}				
836	1,10,10	and f_{-}				
837	F	Function space/Functional vector space				
838	$\hat{z}/\hat{z}_c/\hat{z}_e$	A predicted full/invariant/environment hidden function rep-				
839	, 0, 0	resentation				
840	\mathcal{Z}	Hidden function space/Hidden functional vector space				
841	h	A hypernetwork				
842	\mathcal{H}	Hypernetwork function space				
843	$p(\cdot)$	A probability distribution over a random variable				
844	$I(\cdot; \cdot)$	Mutual information between random variables				
845	$H(\mathbf{X})$	Shannon entropy of the matrix-valued random variable X				
846	$\mathbb{E}_{\mathbf{X}\sim P}[f(X)] \text{ or } \mathbb{E}f(X)$	X) Expectation of $f(X)$ with respect to $p(\mathbf{X})$				
847	$\mathbf{X} \sim P$	Random variable \mathbf{X} has distribution P				
848	$X \sim p(\mathbf{X})$	X is sampled from distribution $p(\mathbf{X})$				
849	$\mathbb{E}_{\mathbf{X}\sim P}[f(\mathbf{X})] \text{ or } \mathbb{E}f(.$	X) Expectation of $f(\mathbf{X})$ with respect to $p(\mathbf{X})$				
850	{·}	An assignment/A substitution rule				
050	$\{\alpha \rightarrow\}$	A symbol with an assigned value				
1 60	$\{\alpha \to a\}$	A symbol with an assigned value a				

85 852

853

854

855 856

857 858

859

В **INVARIANT FUNCTION LEARNING FOUNDATION**

B.1 STRUCTURAL CAUSAL MODEL

In this section, we discuss the trajectory generation process under the Structural Causal Model (SCM) assumption in Fig. 6. To begin with, this SCM is a directed acyclic graph (DAG) with the following components:

- 860 861
- 862 863
- Exogenous variables $U = \{c, e, X_0, \epsilon_c, \epsilon_e, \epsilon_p, \epsilon\}$ are not caused by any variables within the model and are from their own independent distributions. Here $\epsilon_c, \epsilon_e, \epsilon_p, \epsilon$ are noise terms introduced during the function generation and trajectory integral.



which contradicts the independence condition $f'_c \perp \mu$ e, as this would require $H(f_e) = H(f_e|f'_c)$. Therefore, a solution θ_c^* exists, satisfying $\mathbf{f}_c = h_{\theta_c^*}(\mathbf{X}_p)$.

Uniqueness: We now prove that for any solution θ_c^* of the optimization process, it holds that $\mathbf{f}_c = h_{\theta_{\perp}^*}(\mathbf{X}_p).$

We use a proof by contradiction. Assume that there exists another solution $f'_c \neq f_c$ that satisfies the independence constraint and achieves the maximum mutual information. By assumption, we have $H(\mathbf{e}) = H(\mathbf{e}|\mathbf{f}_c)$ and $I(\mathbf{f}_c;\mathbf{f}|\mathbf{X}_0) = I(\mathbf{f}_c;\mathbf{f}|\mathbf{X}_0)$, which implies $H(\mathbf{f}|\mathbf{f}_c) = H(\mathbf{f}|\mathbf{f}_c)$.

Since $f = g_f(f_c, f_e)$, we expand the entropy terms:

$$H(\mathbf{f}|\mathbf{f}_c) = H(\mathbf{f}_c, \mathbf{f}_e|\mathbf{f}_c) = H(\mathbf{f}_e),\tag{9}$$

and

$$H(\mathbf{f}|\mathbf{f}_{c}') = H(\mathbf{f}_{c}, \mathbf{f}_{e}|\mathbf{f}_{c}') = H(\mathbf{f}_{c}|\mathbf{f}_{c}') + H(\mathbf{f}_{e}|\mathbf{f}_{c}') = H(\mathbf{f}_{c}|\mathbf{f}_{c}') + H(\mathbf{f}_{e}).$$
(10)

Substituting $H(\mathbf{f}|\mathbf{f}_c) = H(\mathbf{f}|\mathbf{f}_c)$ into these equations, we find:

$$H(\mathbf{f}_c|\mathbf{f}_c') = 0.$$
 (11)

Based on this, we now aim to prove that $f_c = f'_c$. First, given $H(f'_c|f_c) \ge 0$, we need to show that $H(\mathbf{f}_c'|\mathbf{f}_c) = 0$. Assume that $H(\mathbf{f}_c'|\mathbf{f}_c) > 0$. In this case, \mathbf{f}_c' can determine \mathbf{f}_c while containing more information than f_c , all while remaining independent of f_e . This would imply that H(f'|f) > 0, where $f' = g_f(f', f_e)$, which would in turn affect the corresponding prediction, leading to $H(\mathbf{X}'|\mathbf{X}) > 0$. Such a situation would violate the MSE minimization condition.

Therefore, we must have both $H(f'_c|f_c) = 0$ and $H(f_c|f'_c) = 0$. This implies that f_c and f'_c are isomorphic functions, i.e., there exists a bijective function g_b such that $f_c = g_b(f'_c)$.

Furthermore, since minimizing the MSE of X is equivalent to minimizing the MSE of $\frac{dX}{dt}$ given X_0 , this minimization ensures that, for the same f_e , $f_c(X) = f'_c(X)$ for all X. As a result, the bijective function g_b must be the identity mapping, and we conclude that $f'_c = f_c$ in the support of $p(\mathbf{X})$.

Thus, for any solution θ_c^* of the optimization process, it follows that $f_c = h_{\theta_c^*}(\mathbf{X}_p)$.

B.3 PROOF OF ODE CROSS-ENTROPY MINIMIZATION

Lemma B.2 (ODE cross-entropy minimization 3.2). *Given our forecasting model* $p(\mathbf{X}|h_{\theta}(\mathbf{X}_{p}), \mathbf{X}_{0})$, *it* follows that the cross-entropy minimization between the data distribution $p(\mathbf{X})$ and $p(\mathbf{X}|h_{\theta}(\mathbf{X}_{p}),\mathbf{X}_{0})$ is equivalent to minimizing mean square error $\min_{\theta} \mathbb{E}_{X \sim p} ||X - \hat{X}||_2^2$, where \hat{X} is sampled from $p(\mathbf{X}|h_{\theta}(X_{p}), X_{0}).$

Proof. The forecasting optimization goal is to use our framework to approximate the data distribution $p(\mathbf{X})$, parameterized as $p(\mathbf{X}|h_{\theta}(\mathbf{X}_{p}), \mathbf{X}_{0})$, where we apply the cross-entropy minimization, *i.e.*, $H(p(\mathbf{X}), p(\mathbf{X}|h_{\theta}(\mathbf{X}_{p}), \mathbf{X}_{0})) = -\mathbb{E}_{\mathbf{X} \sim p} \left[\log p(X|h_{\theta}(X_{p}), X_{0})\right]$. Furthermore, this negative log-likelihood optimization can be further reduced to the common mean squared error (MSE).

$$\min_{\theta} \mathbb{E}_{\mathbf{X} \sim p} \| X - \hat{X} \|_2^2, \tag{12}$$

where $\hat{X} \sim p(\mathbf{X}|h_{\theta}(X_p), X_0)$. Since the distribution $p(\mathbf{X}|h_{\theta}(X_p), X_0)$ is modeled as a Gaussian $\mathcal{N}(\mathbf{X}; g_{int}(h_{\theta}(X_p), X_0), \sigma^2 I)$ (Sec. 3.1), we have $\hat{X} = \mu + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$.

$$\min_{\theta} -\mathbb{E}_{\mathbf{X} \sim p}[\log p(X|h_{\theta}(X_p), X_0)]$$

969
970
$$= \min_{\theta} \mathbb{E}_{\mathbf{X} \sim p} \left[\frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|X - \mu\|_2^2 \right]$$
(13)

971
$$= \min_{\theta} \mathbb{E}_{\mathbf{X} \sim p} \left[\frac{1}{2\sigma^2} \| X - \mu \|_2^2 \right] + \frac{n}{2} \log(2\pi\sigma^2)$$

Since $\frac{n}{2}\log(2\pi\sigma^2)$ is a constant, we ignore it in the minimization process.

$$\min_{\theta} \mathbb{E}_{\mathbf{X} \sim p} \left[\frac{1}{2\sigma^{2}} \|X - \mu\|_{2}^{2} \right] + \frac{n}{2} \log(2\pi\sigma^{2})$$

$$= \min_{\theta} \mathbb{E}_{\mathbf{X} \sim p} \left[\frac{1}{2\sigma^{2}} \|X - (\hat{X} - \epsilon)\|_{2}^{2} \right]$$

$$= \min_{\theta} \mathbb{E}_{\mathbf{X} \sim p} \left[\frac{1}{2\sigma^{2}} \|X - \hat{X} + \epsilon\|_{2}^{2} \right]$$

$$= \min_{\theta} \mathbb{E}_{\mathbf{X} \sim p} \left[\frac{1}{2\sigma^{2}} (\|X - \hat{X}\|_{2}^{2} + \|\epsilon\|_{2}^{2} + 2(X - \hat{X})^{T} \epsilon) \right]$$

$$= \min_{\theta} \mathbb{E}_{\mathbf{X} \sim p} \left[\frac{1}{2\sigma^{2}} \|X - \hat{X}\|_{2}^{2} \right] + \frac{1}{2\sigma^{2}} \mathbb{E}_{\epsilon} [\|\epsilon\|_{2}^{2}] + \frac{1}{\sigma^{2}} \mathbb{E}_{\mathbf{X} \sim p} (X - \hat{X})^{T} \mathbb{E}_{\epsilon} [\epsilon]$$
(14)

Here, $\frac{1}{2\sigma^2} \mathbb{E}_{\mathbf{X} \sim p}[\|\epsilon\|_2^2]$ is a constant; $\mathbb{E}_{\mathbf{X} \sim p}\left[\frac{1}{\sigma^2}(X - \hat{X})^T\epsilon\right] = \frac{1}{\sigma^2} \mathbb{E}_{\mathbf{X} \sim p}(X - \hat{X})^T \mathbb{E}_{\epsilon}[\epsilon] = 0$ since ϵ is independently sampled with a zero mean. Therefore,

$$\min_{\theta} \mathbb{E}_{\mathbf{X} \sim p} \left[\frac{1}{2\sigma^{2}} \| X - \hat{X} \|_{2}^{2} \right] + \frac{1}{2\sigma^{2}} \mathbb{E}_{\epsilon} [\|\epsilon\|_{2}^{2}] + \frac{1}{\sigma^{2}} \mathbb{E}_{\mathbf{X} \sim p} (X - \hat{X})^{T} \mathbb{E}_{\epsilon} [\epsilon]$$

$$= \min_{\theta} \mathbb{E}_{\mathbf{X} \sim p} \left[\frac{1}{2\sigma^{2}} \| X - \hat{X} \|_{2}^{2} \right]$$

$$= \min_{\theta} \mathbb{E}_{\mathbf{X} \sim p} \left[\| X - \hat{X} \|_{2}^{2} \right]$$
(15)

 This reduction connects the information theory and the practical MSE optimization, which further helps us to transform the mutual information maximization into a similar MSE optimization below.

PROOF OF ODE CONDITIONAL MUTUAL INFORMATION MAXIMIZATION **B**.4

Proposition B.3 (Proof of ODE conditional mutual information maximization 3.3). Given fore-casting model $p(\mathbf{X}|h_{\theta_c}(X_p), X_0)$, it follows that the conditional mutual information maximization $\max_{\theta_{\alpha}} I(h_{\theta_{\alpha}}(\mathbf{X}_{p}); \mathbf{f}|\mathbf{X}_{0})$ is equivalent to minimizing mean square error $\min_{\theta_{\alpha}} \mathbb{E}_{\mathbf{X} \sim P} ||X - X^{c}||_{2}^{2}$, where \hat{X}^c is the predicted trajectory sampled from $p(\mathbf{X}|h_{\theta_c}(X_p), X_0)$.

Proof. According to Lemma 3.2, we can reduce a negative log-likelihood minimization $\min_{\theta_c} -\mathbb{E}_{\mathbf{X}\sim p} \log p(X|h_{\theta_c}(X_p), X_0) \text{ to } \min_{\theta_c} \mathbb{E}_{\mathbf{X}\sim P} ||X - \hat{X}^c||_2^2.$

Therefore, we only need to prove the equivalence between the $-\mathbb{E}_{\mathbf{X}\sim p} \log p(X|h_{\theta_c}(X_p), X_0)$ mini-mization and the $I(h_{\theta_c}(\mathbf{X}_p); \mathbf{f}|\mathbf{X}_0)$ maximization. It follows that

1014
1015
1016
$$\begin{aligned}
\max_{\theta_c} I(h_{\theta_c}(\mathbf{X}_p); \mathbf{f} | \mathbf{X}_0) \\
= \max_{\theta_c} H(\mathbf{f} | \mathbf{X}_0) - H(\mathbf{f} | h_{\theta_c}(\mathbf{X}_p), \mathbf{X}_0).
\end{aligned}$$
(16)

Since $H(f|\mathbf{X}_0)$ is a constant, we ignore it in the maximization process and obtain

- $\max_{\theta_c} H(\mathbf{f}|\mathbf{X}_0) - H(\mathbf{f}|h_{\theta_c}(\mathbf{X}_p), \mathbf{X}_0)$
- $= \max_{\theta_c} H(\mathbf{f}|h_{\theta_c}(\mathbf{X}_p), \mathbf{X}_0)$

$$=\min_{o} H(\mathbf{f}|h_{\theta_c}(\mathbf{X}_p), \mathbf{X}_0)$$

Since $\mathbf{X} = g_{int}(\mathbf{f}, \mathbf{X}_0) + \epsilon$ where ϵ is an independent random noise, $H(\mathbf{X}|h_{\theta_c}(\mathbf{X}_p), \mathbf{X}_0, \mathbf{f}) =$ $H(\epsilon | h_{\theta_c}(\mathbf{X}_p)) = H(\epsilon)$. Since given specific θ_c , $H(h_{\theta_c}(\mathbf{X}_p) | \mathbf{X}) = 0$, with the conditional entropy

(17)

1026 chain rule, it follows that 1027 $\min_{\theta_c} H(\mathbf{f}|h_{\theta_c}(\mathbf{X}_p), \mathbf{X}_0)$ 1028 $= \min_{\theta_c} H(\mathbf{f}, \mathbf{X} | h_{\theta_c}(\mathbf{X}_p), \mathbf{X}_0) - H(\mathbf{X} | h_{\theta_c}(\mathbf{X}_p), \mathbf{X}_0, \mathbf{f})$ 1029 1030 $= \min_{\theta_{\tau}} H(\mathbf{f}, \mathbf{X} | h_{\theta_c}(\mathbf{X}_p), \mathbf{X}_0) - H(\epsilon)$ 1031 1032 $=\min_{a} H(\mathbf{f}, \mathbf{X} | h_{\theta_c}(\mathbf{X}_p), \mathbf{X}_0)$ 1033 $= \min_{\theta_c} H(\mathbf{f}|h_{\theta_c}(\mathbf{X}_p), \mathbf{X}_0, \mathbf{X}) + H(\mathbf{X}|h_{\theta_c}(\mathbf{X}_p), \mathbf{X}_0)$ 1034 (18)1035 $= \min_{\theta_c} H(\mathbf{f}|\mathbf{X}_0, \mathbf{X}) + H(\mathbf{X}|h_{\theta_c}(\mathbf{X}_p), \mathbf{X}_0)$ 1036 $= \min_{\theta_c} H(\mathbf{X}|h_{\theta_c}(\mathbf{X}_p), \mathbf{X}_0)$ $=\min_{\theta_{\alpha}} -\mathbb{E}_{\mathbf{X},\mathbf{X}_{p},\mathbf{X}_{0}\sim p}\log p(X|h_{\theta_{c}}(X_{p}),X_{0})$ 1039 1040 $= \min_{\theta} - \mathbb{E}_{\mathbf{X} \sim p} \log p(X|h_{\theta_c}(X_p), X_0)$ 1041 where $H(f|\mathbf{X}_0, \mathbf{X})$ is a constant. Here, we finish building the equivalence between the function

where $H(f|X_0, X)$ is a constant. Here, we finish building the equivalence between the function conditional mutual information maximization and the trajectory negative log-likelihood minimization. Using the proof in Lemma 3.2, our final optimization goal can be reduced to

$$\min_{\theta_c} \mathbb{E}_{\mathbf{X} \sim P} \| X - \hat{X}^c \|_2^2.$$
(19)

1048 1049 B.5 THEORETICAL JUSTIFICATION FOR ADVERSARIAL TRAINING

To incorporate the independence constraint, we enforce the condition $\hat{\mathbf{f}}_c \perp \mathbf{I}$ e, where $\hat{\mathbf{f}}_c = h_{\theta_c}(\mathbf{X}_p)$ is the predicted function random variable, not a realization. Since $\hat{\mathbf{f}}_c \perp \mathbf{I}$ e is equivalent to $I(\mathbf{e}; \hat{\mathbf{f}}_c) = 0$, and $I(\mathbf{e}; \hat{\mathbf{f}}_c) \ge 0$, the objective $I(\mathbf{e}; \hat{\mathbf{f}}_c)$ reaches its minimum when and only when $\hat{\mathbf{f}}_c \perp \mathbf{I}$ e. This leads us to define minimizing $I(\mathbf{e}; \hat{\mathbf{f}}_c)$ as our training criterion.

Definition 1. The environment independence training criterion is

1045

1046

1047

1056 1057

1066

1071 1072

1074

$$\theta_c^* = \operatorname*{arg\,min}_{\theta_c} I(\mathbf{e}; \hat{\mathbf{f}}_c). \tag{20}$$

1058 This training criterion can be used directly, since the mutual information $I(\mathbf{e}; \hat{\mathbf{f}}_c) = \mathbb{E}\left[\log \frac{P(\mathbf{e}|\hat{\mathbf{f}}_c)}{P(\mathbf{e})}\right]$ 1059 while $P(\mathbf{e}|\hat{\mathbf{f}}_c)$ is unknown. Following Proposition 1 in GAN Goodfellow et al. (2020), we introduce 1061 an optimal discriminator $g_{\phi} : \mathcal{F} \mapsto \mathcal{E}$ with parameters ϕ to approximate the unknown $P(\mathbf{e}|\hat{\mathbf{f}}_c)$ as 1062 $P_{\phi}(\mathbf{e}|\hat{\mathbf{f}}_c)$, minimizing the negative log-likelihood $-\mathbb{E}\left[\log P_{\phi}(\mathbf{e}|\hat{\mathbf{f}}_c)\right]$. We then have the following two 1063 propositions:

Proposition B.4. For θ_c fixed, the optimal discriminator ϕ is

$$= \arg\min_{\phi} -\mathbb{E}\left[\log P_{\phi}(\mathbf{e}|\hat{\mathbf{f}}_{c})\right].$$
(21)

This proposition can be proved straightforwardly by applying the cross-entropy training criterion. **Proposition B.5.** Denoting KL-divergence as $KL[\cdot || \cdot]$, for θ_c fixed, the optimal discriminator ϕ is ϕ^* , such that

$$\operatorname{KL} \left| P(\mathbf{e}|\mathbf{f}_c) \| P_{\phi^*}(\mathbf{e}|\mathbf{f}_c) \right| = 0.$$
(22)

⁷³ *Proof.* Given a fixed θ_c , both $I(e; f_c)$ and H(e) are constants. Therefore, we have:

 ϕ

1075 1076 1076 1077 1078 1079 $\phi^* = \arg\min_{\phi} -\mathbb{E}\left[\log P_{\phi}(\mathbf{e}|\hat{\mathbf{f}}_c)\right]$ $= \arg\min_{\phi} I(\mathbf{e}; \hat{\mathbf{f}}_c) - \mathbb{E}\left[\log P_{\phi}(\mathbf{e}|\hat{\mathbf{f}}_c)\right] - H(\mathbf{e}) \quad (23)$ $= \arg\min_{\phi} \mathrm{KL}\left[P(\mathbf{e}|\hat{\mathbf{f}}_c) \| P_{\phi}(\mathbf{e}|\hat{\mathbf{f}}_c)\right].$ Thus, minimizing the negative log-likelihood of $P_{\phi}(\mathbf{e}|\hat{\mathbf{f}}_c)$ is equivalent to minimizing the KL divergence between $P(\mathbf{e}|\hat{\mathbf{f}}_c)$ and its approximation $P_{\phi}(\mathbf{e}|\hat{\mathbf{f}}_c)$. Since KL divergence is bounded by 0, we have KL $\left[P(\mathbf{e}|\hat{\mathbf{f}}_c) \| P_{\phi^*}(\mathbf{e}|\hat{\mathbf{f}}_c)\right] = 0$. This concludes the proof.

1084 1085 With these propositions, the mutual information can be computed with the help of the optimal 1086 discriminator ϕ^* . According to Proposition B.5, we have:

$$I(\mathbf{e}; \hat{\mathbf{f}}_c) = \mathbb{E}\left[\log P_{\phi^*}(\mathbf{e}|\hat{\mathbf{f}}_c)\right] + H(\mathbf{e}) + \mathrm{KL}\left[P(\mathbf{e}|\hat{\mathbf{f}}_c) \| P_{\phi^*}(\mathbf{e}|\hat{\mathbf{f}}_c)\right]$$

$$= \mathbb{E}\left[\log P_{\phi^*}(\mathbf{e}|\hat{\mathbf{f}}_c)\right] + H(\mathbf{e}) + 0.$$
(24)

Thus, by disregarding the constant H(e), the training criterion becomes:

$$\theta_{c}^{*} = \arg\min_{\theta_{c}} I(\mathbf{e}; \hat{\mathbf{f}}_{c})$$

$$= \arg\min_{\theta_{c}} \mathbb{E} \left[\log P_{\phi^{*}}(\mathbf{e} | \hat{\mathbf{f}}_{c}) \right]$$

$$= \arg\min_{\theta_{c}} \left\{ \max_{\phi} \mathbb{E} \left[\log P_{\phi}(\mathbf{e} | \hat{\mathbf{f}}_{c}) \right] \right\},$$
(25)

where $P_{\phi}(\mathbf{e}|\mathbf{f}_c)$ is the probability modeling of g_{ϕ} . Therefore, the log-likelihood adversarial training can enforce the independence $h_{\theta_c}(\mathbf{X}_p) \perp \mathbf{e}$.

1102 1103 C DATASETS

1104 1105 C.1 BASIC SETUP

1106 We conduct experiments on the proposed three multi-environment datasets ME-Pendulum, ME-Lotka-1107 Volterra, and ME-SIREpidemic. Each of these datasets includes 1000 samples, where 800 and 200 1108 samples are assigned to training set and test set, respectively. Each training set has 4 environments 1109 where 200 samples are generated in each environment. Each sample is observed over 10 units of time, 1110 and each time is discretized by regularly-spaced discrete time steps from t_0 to t_T , where T = 99, *i.e.*, there are 100 time intervals of 0.1 unit of time each. A sample generation process is controlled by a 1111 set of ODEs with common parameters $\mathbf{W}^c \sim \mathcal{U}(W_{low}^c, W_{high}^c)$ and environment-specific parameters 1112 $\mathbf{W}^e \sim \mathcal{U}(W^e_{low}, W^e_{high})$ sampled from their uniform distributions, *e.g.*, in the pendulum system, 1113 $\mathbf{W}^c = \{ \boldsymbol{\alpha} \}$ and $\mathbf{W}^e = \{ \boldsymbol{\rho} \}$, where $\boldsymbol{\alpha} \sim \mathcal{U}(\alpha_{low}, \alpha_{high})$ and $\boldsymbol{\rho} \sim \mathcal{U}(\rho_{low}, \rho_{high})$. Note that each 1114 environment has its specific function form with its environment-specific parameters \mathbf{W}^{e} . 1115

1116 The environment split is the same in the test set. The only difference is that in the test set each 1117 sample X has one additional prediction target X^c with only the invariant dynamics. For example, the 1118 invariant trajectory X^c of the one generated by $-\alpha^2 \sin \theta_t - \rho \omega_t$ will be created by $-\alpha^2 \sin \theta_t$.

1119

1093 1094 1095

1099

1100

1101

1120 C.2 ME-PENDULUM

1121 1122 ME-Pendulum is motivated by the DampedPendulum system (Yin et al., 2021b). The state $X_t = [\theta_t, \omega_t] \in \mathbb{R}^2$ are the angle and angular velocity of the pendulum at time t, where we have $\mathbf{W}^c = \{\alpha\}$, 1124 $\mathbf{W}^e = \{\rho\}$, and $T_C = \frac{T}{3}$. The underlying invariant ODE is $\frac{d\theta_t}{dt} = \omega_t, \frac{d\omega_t}{dt} = -\alpha^2 \sin(\theta_t)$. As shown in Tab. 4, this invariant ODE is entangled with different environmental factors, forming four environments, namely, *damped*, *powered*, *spring*, *air*.

- 1126 1127
- C.3 ME-LOTKA-VOLTERRA

1128

1129 Motivated by the Lotka-Volterra system (Ahmad, 1993), the state $X_t = [p_t, q_t] \in \mathbb{R}^2$ are the population 1130 of preys and predators at time t, where we have $\mathbf{W}^c = \{\alpha, \beta, \gamma, \delta\}$, $\mathbf{W}^e = \{\alpha', \beta', \gamma', \delta'\}$, and 1131 $T_C = \frac{T}{2}$. The underlying invariant ODEs are $\frac{dp}{dt} = \alpha p - \beta pq$, $\frac{dq}{dt} = \delta pq - \gamma q$. As shown in Tab. 5, 1132 these invariant ODEs are entangled in 4 environments, *i.e.*, save, fight, resource, omnivore. The 1133 save environment ODEs simulate the decrease of food wastage along with the increase of predators. The fight environment ODEs simulate the decrease of hunting efficiency along with the increase of

Environment	ODE for θ_t	ODE for ω_t	Distribution of Parameter
Damped	$\frac{d\theta_t}{dt} = \omega_t$	$\frac{d\omega_t}{dt} = -\boldsymbol{\alpha}^2 \sin\left(\theta_t\right) - \boldsymbol{\rho}\omega_t$	
Powered	$\frac{d\tilde{\theta}_t}{dt} = \omega_t$	$\frac{d\omega_t}{dt} = -\boldsymbol{\alpha}^2 \sin\left(\theta_t\right) + \boldsymbol{\rho} \frac{\omega_t}{ \omega_t }$	$2 + \frac{1}{2} \left(1 + \frac{1}{2} + \frac{1}{2} \right)$
Spring	$\frac{d\theta_t}{dt} = \omega_t$	$\frac{d\omega_t}{dt} = -\boldsymbol{\alpha}^2 \sin\left(\theta_t\right) - \boldsymbol{\rho}\theta_t$	$\boldsymbol{\alpha} \sim \mathcal{U}(1.0, 2.0)$
Air	$rac{d heta_t}{dt} = \omega_t$	$\frac{d\omega_t}{dt} = -\boldsymbol{\alpha}^2 \sin\left(\theta_t\right) - \boldsymbol{\rho} \omega_t \omega_t$	$\rho \sim u(0.2, 0.4)$
Invariant	$\frac{d\theta_t}{dt} = \omega_t$	$\frac{d\omega_t}{dt} = -\boldsymbol{lpha}^2 \sin\left(\theta_t\right)$	-

Table 4: **ME-Pendulum ODEs** with $\theta_0 \sim \mathcal{U}(0, \frac{\pi}{2})$ and $\omega_0 \sim \mathcal{U}(-1, 0)$.

predators. The resource environment ODEs limit the increase rate of the prey population. In the omnivore environment ODEs, the predators are omnivores that can build the population without preys under certain resource limits.

We plot the trajectories X in the training set, and the invariant trajectories X^c in the test set in Fig. 7.

Invariant Test: ME-Pendulum Training Set: ME-Pendulum 0.75 0.50 0.25 0.0 -0.4 -0.6 -0.8

Figure 7: ME-Pendulum trajectories.

Table 5: ME-Lotka-Volterra ODEs with $p_0 \sim \mathcal{U}(1000, 2000)$ and $q_0 \sim \mathcal{U}(10, 20)$.

Environment	ODE for p_t	ODE for q_t	Distribution of Parameters
Save Fight Resource Omnivore	$ \begin{array}{l} \frac{dp}{dt} = \boldsymbol{\alpha}p - \boldsymbol{\beta}pq - \boldsymbol{\beta}'pq \cdot 10\exp\left(-\frac{q}{10}\right) \\ \frac{dp}{dt} = \boldsymbol{\alpha}p - \boldsymbol{\beta}pq \\ \frac{dp}{dt} = \boldsymbol{\alpha}p - \boldsymbol{\alpha}'\frac{p^2}{2000} - \boldsymbol{\beta}pq \\ \frac{dp}{dt} = \boldsymbol{\alpha}p - \boldsymbol{\alpha}' - \boldsymbol{\beta}pq \end{array} $	$ \begin{array}{c} \frac{dq}{dt} = \pmb{\delta}pq - \pmb{\gamma}q \\ \frac{dq}{dt} = \pmb{\delta}pq + \pmb{\delta}'pq \cdot 10 \exp\left(-\frac{q}{10}\right) - \pmb{\gamma}q \\ \frac{dq}{dt} = \pmb{\delta}pq - \pmb{\gamma}q \\ \frac{dq}{dt} = \pmb{\delta}pq + 20\pmb{\gamma}' \left(1 - \frac{q}{100}\right) - \pmb{\gamma}q \end{array} $	$ \begin{aligned} \boldsymbol{\alpha}, \boldsymbol{\alpha}' &\sim \mathcal{U}(1.2, 2.4) \\ \boldsymbol{\beta}, \boldsymbol{\beta}' &\sim \mathcal{U}(6e-2, 1.2e-1) \\ \boldsymbol{\gamma}, \boldsymbol{\gamma}' &\sim \mathcal{U}(0.48, 0.96) \\ \boldsymbol{\delta}, \boldsymbol{\delta}' &\sim \mathcal{U}(4.8e-4, 9.6e-4) \end{aligned} $
Invariant	$rac{dp}{dt} = oldsymbol{lpha} p - oldsymbol{eta} p q$	$rac{dq}{dt} = oldsymbol{\delta} pq - oldsymbol{\gamma} q$	

We plot the trajectories X in the training set, and the invariant trajectories X^c in the test set in Fig. 8.

C.4 ME-SIREPIDEMIC

In the SIREpidemic (Wang et al., 2021) system, the states $X_t = [S_t, I_t, R_t] \in \mathbb{R}^3$ are the suscepti-ble, infected, and recovered individuals at time t, respectively. In this adapted ME-SIREpidemic system, we have $\mathbf{W}^c = \{\beta\}$, $\mathbf{W}^e = \{\gamma\}$, and $T_c = \frac{T}{2}$. The underlying invariant ODEs are $\frac{dS}{dt} = -\beta \frac{SI}{S+I+R}, \frac{dI}{dt} = \beta \frac{SI}{S+I+R}, \frac{dR}{dt} = 0$, where we only care about the S to I transformation relationship. As shown in Tab. 6, we introduce 4 environments, *origin, enlarge, loop, negative*. These four environments describe four different models, where some of them are only for math modeling. The origin environment is the same as the original SIREpidemic model. The enlarge environment ODEs expand the epidemic range. The loop environment ODEs include deaths and second-time infections. The negative environment ODEs is a pure math model allowing negative numbers.

We plot the trajectories X in the training set, and the invariant trajectories X^c in the test set in Fig. 9.



1239 We conduct experiments on 800-sample training sets with a training batch size of 32, which leads 1240 to 25 iterations per epoch. For each run, we optimize the neural network with 2,000 epochs, which 1241 is equivalent to 50,000 iterations. Given fixed learning iterations, the learning rate is selected from U(1e - 4, 1e - 3).

1242 D.1 BASELINES

In our experiments, we design four adapted baselines, since this new task has never been explored before. The selection of baselines is based on the following two questions.

- As invariant learning has been successfully applied in many representation learning tasks, can general invariant learning principle still work for invariant function learning?
- Meta-learning techniques has been well designed to solved problems in dynamical systems due to their quick adaptation characteristics. However, there is no evidence that the quick adaptable functions are the invariant function that shares across environments. Are they?

For the first question, since our method is based on enforcing independence which shares a similar philosophy as domain adversarial neural network Ganin et al. (2016), we consider other invariant methods going different ways. There are many methods well-known in the field of invariant learning, but considering our invariant function learning formulation, we found the discovery of invariant function requires independence that is different from the most typical "well performed across all environments" invariant learning requirements. Therefore, as a validation of our guess, we adapt two widely used and known invariant learning methods, IRM Arjovsky et al. (2019) and VREx Krueger et al. (2021). These two techniques are directly applied to our framework with the same architecture and their typical hyper-parameters searching spaces:

1261 1262 1263

1264

1246

1247

1248 1249

1250

1251

1252

•
$$\lambda_{irm} \sim \mathcal{U}(1e-2, 1e2)$$

• $\lambda_{vrex} \sim \mathcal{U}(1e - 1, 1e3)$

As shown in our result plot 4, the results are not surprising. That means that the function can perform well across multiple environments is not the invariant function.

For the second question, our initial guess is that meta-learning is very sensitive to the distribution of 1268 the training set. If certain pattern exists in multiple environments (not all), the meta-learning methods 1269 are prone to capture it as a part of the meta-function, which satisfies their quick adaptation goals. 1270 In our experiments, we choose MAML Finn et al. (2017) and CoDA Kirchmeyer et al. (2022) as 1271 our adapted baselines. CoDA is adapted since it uses hypernetwork as a full network adaptor which 1272 is similar to our framework. However, CoDA is applied on coefficient-environments and requires 1273 test-time adaptation without a trajectory encoder, leading to significant architecture differences. 1274 Therefore, we apply CoDA as meta-learning techniques focusing on its low-dimension (2-dimension) 1275 environment representation and regularization. MAML is selected as the most typical meta-learning 1276 baseline, which does not require the use of hypernetwork, and only learns a meta function used to 1277 predict invariant trajectories. Their typical hyper-parameters searching spaces are shown as follows.

1278 1279

1280

1281

- $\lambda_{coda} \sim \mathcal{U}(1e-5, 1e-3)$
 - $\lambda_{maml} \sim \mathcal{U}(1e-3,1)$ (Meta learning rate)

1282 1283 D.2 DISENTANGLEMENT OF INVARIANT FUNCTION SETUP

1284 D.2.1 ARCHITECTURE

Transformer. For the trajectory encoder in our hypernetwork, we apply a 6-layer 8-head 256-dimension FFN transformer Vaswani (2017) with frequency positional encoding Gehring et al. (2017).
We tried different architectures like GRUs Cho (2014), but the transformer encoder can provide the best in-domain test performance easily. We also sweeped to the depth, width, and number of heads, and found that 6-layer 8-head 256-dimension FFN transformer is strong enough for our ODE systems without making training difficult.

Function embedding. In the hidden function embedding space, we select the function embedding dimension to be 32 or 64, while these two selections perform quite similar. For the MLPs used to disentangle and decode hidden function embedding, we use 3-layer MLPs with ReLU as the activations. While for the decoder, the last layer projects the hidden function embedding to parameterized function space \mathbb{R}^m where *m* is the number of parameters in the derivative neural network. **Derivative function network.** The derivative neural network is a 4-layer or 5-layer MLP with width 16 or 32, which takes $X_t \in \mathbb{R}^d$ as input and output $\frac{dX_t}{dt} \in \mathbb{R}^d$. This neural network is transformed to be a functional in our implementation.

Discriminator. Our discriminator is a 3 to 6 layer MLP with width 64 or 128. The size of discriminator can be easily chosen since the main goal of it is to discriminate the environment information in the hidden function space. Therefore, the simplest way to filter non-qualified discriminators is using the prediction entropy of $P_{\phi}(e|\hat{f}_e)$, since \hat{f}_e should contain rich environment information, different to the prediction from \hat{f}_c . Our experiments also validate that the failure to distinguish \hat{f}_e will always cause the failure of invariant function discovery, which is natural, since the adversarial training is based on the optimal discriminators Goodfellow et al. (2020).

1307 Note that for all our MLPs, we apply one LayerNorm before each activation.²

1308 1309

D.2.2 TRAINING OBJECTIVES

1310

1311

1312 We restate our three additional strategies here. Firstly, the adversarial training of f_c will cause the 1313 loss of environment information in \hat{f}_c , leading to the training difficulty of the discriminator; therefore, 1314 this discriminator is not only trained on \hat{f}_c but also on \hat{f}_e with the same hyper-parameter λ_{dis} , *i.e.*, 1315 $\min_{\phi,\theta} - \mathbb{E}_{\mathbf{X},e\sim P} \left| \log P_{\phi}(e|\hat{f}_e) \right|$. Secondly, instead of using the large \hat{f}_c and \hat{f}_e as the input of the 1316 discriminator, we input the the corresponding embeddings z_c and z_c . Thirdly, to avoid the use of a 1317 numerical or neural integrator which causes long training time, we follow Mouli et al. (2024) to fit 1318 derivatives only. That is, instead of using the inference forecaster $p(X|h_{\theta_c}(X_p), X_0)$, we calculate the 1319 derivatives of X using \hat{f} and \hat{f}_c , and replace the MSE over trajectory matrices with the MSE over 1320 derivative matrices. Note that this modification only eliminates the use of integrator for stability 1321 during training and thus does not affect our analysis and optimization goal. 1322

We introduce our hyper-parameter searching space as follows.

1324 1325

1326

1327

- $\lambda_c \sim \mathcal{U}(1e-7, 1e-4)$
- $\lambda_{dis} \sim \mathcal{U}(1e-1,1)$
 - λ'_{adv} ~ U(1e2, 1e6)
 λ_{adv} = λ_c · λ'_{adv}
- 1328 1329
- 1330

1339

1340

1344 1345

1347 1348 1349

1331 The most critical hyper-parameters are λ_c and λ_{adv} which control the information overlap between 1332 f_c and f. Conceptually, λ_c controls the conditional mutual information (MI) maximization in our 1333 invariant function learning principle, while λ_{adv} enforces the independence constraint. The intensity 1334 of the independence enforcing λ_{adv} is dependent on the intensity of MI maximization λ_c ; thus, we set 1335 λ_{adv} according to λ_c , leading to $\lambda_{adv} = \lambda_c \cdot \lambda'_{adv}$.

1336 λ_{dis} is only discriminator training, which is relatively trivial according to our discriminator descriptions in Appx. D.2.1.

D.2.3 METRIC

Root mean square deviation (RMSE) is a commonly used metric, but it suffers difficulties when comparing datasets with different value scales. Therefore, we normalize it using its standard deviation.

$$NRMSE = \frac{\sqrt{\mathbb{E}_{\mathbf{X} \sim p} \|X - \hat{X}\|_2^2}}{Std(\mathbf{X})}$$
(26)

²The code will be released upon acceptance.



Figure 10: Ablation study on 3 DIF variants under 3 multi-environment ODE systems. For each pipeline, we provide model candidates with 50+ random hyper-parameter selections in their searching spaces, *i.e.*, more than 450 models in the figure.



Figure 11: **Trajectory input length study** on models trained with different training input length factors under 3 multi-environment ODE systems.

1377 D.3 SOFTWARE AND HARDWARE

Our implementation is under the architecture of PyTorch Paszke et al. (2019). The deployment
environments are Ubuntu 20.04 with 48 Intel(R) Xeon(R) Silver, 4214R CPU @ 2.40GHz, 755GB
RAM, and graphics cards NVIDIA RTX 2080Ti.

1382 1383

1384

1374

1375 1376

1360

1361

1362

E SUPPLEMENTARY EXPERIMENTS

1385 E.1 ABLATION STUDY

1387 As a complementary ablation study of Sec. 5.5, we train two models by eliminating two important components from the original model, namely, f pipeline and f_c pipeline. The f pipeline removes 1388 the discriminator and only output \hat{f} to the forecasting, which neglects the disentanglement process. 1389 The f_c pipeline prunes the \hat{f} output while maintaining the adversarial training process. As shown 1390 in Fig. 10, both f and f_c pipelines fail to perform valid invariant function learning aligning with 1391 our theoretical results. Specifically, the f pipeline faces difficulties in extracting invariant functions 1392 without environment information. The unsatisfactory performance of the f_c pipeline is attributed 1393 to the discriminator's training failure. This is because the training of the discriminator requires the 1394 capture of environment information, but the elimination of the \hat{f} part also removes the training of \hat{z}_e , 1395 the critical environment information captor. Therefore, the discriminator loses the most important 1396 environment information input, leading to training failure.

1397

1398 E.2 INPUT LENGTH AND ENVIRONMENT ANALYSIS

1400 In order to ensure fairness, we fix the input length and the number of environments in our experiments. 1401 However, it is also interesting to figure out the effects of the input trajectory length T_c and the 1402 number of environments on the model performance. Fig. 11 shows the performance of DIF given 1403 different input length factor l_t , where $T_c = \frac{T}{l_t}$. The results indicate the input length does not affect the performance significantly, where the only variances are attributed to the training difficulty of the



Figure 12: Environment analysis on models trained with different numbers of training environments
 under 3 multi-environment ODE systems.

Method	ME-Pendulum			ME-Lotka-Volterra	ME-SIREpidemic	
memou	NRMSE	SR Explanation	NRMSE	SR Explanation	NRMSE	SR Explanation
MAML	0.9704	$\begin{aligned} \frac{d\theta_t}{dt} &= \omega_t - \frac{0.036\left(\theta_t + \omega_t\right)}{e^{\omega_t}} \\ \frac{d\omega_t}{dt} &= \frac{\sin\left(\theta_t\right)}{-0.67} - 0.48\sin\left(\left(\theta_t + 0.68\right)\sin\left(\omega_t\right)\right) \end{aligned}$	0.6774	$\frac{dp_t}{dt} = q_t \frac{1}{p_t + 0.63} (-0.022)$ $\frac{dq_t}{dt} = \frac{p_t}{\frac{p_t}{0.17} - q_t + q_t q_t} + \delta$	0.2673	$ \begin{vmatrix} \frac{dS_t}{dt} = \frac{-S_t I_t - I_t + \cos\left(t\right)}{0.69} \\ \frac{dI_t}{dt} = 0.49 S_t + e^{\sin\left(0.45 S_t\right)} \\ \frac{dR_t}{dt} = \cos\left(I_t \sin\left(S_t\right)\right) \left(-0.45 S_t\right) \\ \frac{dR_t}{dt} = \cos\left(I_t \sin\left(S_t\right)\right) $
CoDA	0.9695	$\frac{d\theta_t}{dt} = \omega_t \rho \sin\left(\theta_t + \omega_t + 0.94\right) + \omega_t + 0.074$ $\frac{d\omega_t}{dt} = \frac{\sin\left(-\theta_t + \omega_t\left(-0.20\right) + 0.37\right)}{0.49} - 0.57$	0.7097	$\frac{dp_t}{dt} = (1.1 - p_t p_t) \ 0.54$ $\frac{dq_t}{dt} = (p_t + p_t) (p_t + q_t (-0.26)) - 0.66$	0.3184	$\begin{vmatrix} \frac{dS_t}{dt} = \frac{-S_t + \sin(S_t \cdot 0.4)}{0.44} \\ \frac{dI_t}{dt} = \frac{3.5S_t}{S_t + \frac{I_t}{S_t - 0.54}} + 0.0 \\ \frac{dR_t}{dt} = \frac{S_t}{I_t + I_t + \frac{eS_t}{I_t}} + 0.0 \end{aligned}$
IRM	0.7042	$\frac{d\theta_t}{dt} = \omega_t \cdot 0.93$ $\frac{d\omega_t}{dt} = -\theta_t \alpha + \rho$	0.6989	$\frac{dp_t}{dt} = -0.012$ $\frac{dq_t}{dt} = 0.083$	0.9768	$\left \begin{array}{c} \displaystyle \frac{dS_t}{dt} = e^{-S_t + \sin{(S_t)}} - 1.7 \\ \displaystyle \frac{dI_t}{dt} = \sin{\left(\frac{S_t}{S_t + S_t + \frac{\beta}{S_t}}\right)} \\ \displaystyle \frac{dR_t}{dt} = 0.33 - 0.021S_t \end{array} \right $
VREx	0.7274	$\frac{d\theta_t}{dt} = \omega_t \cdot 0.92$ $\frac{d\omega_t}{dt} = \alpha (-1.1) \theta_t$	0.6877	$\frac{dp_t}{dt} = -0.032$ $\frac{dq_t}{dt} = 0.12$	0.4652	$ \begin{vmatrix} \frac{dS_t}{dt} = S_t \left(-I_t - 0.10\beta \right) \\ \frac{dI_t}{dt} = \frac{S_t\beta}{S_t + S_t + \frac{\beta}{S_t}} + 0.0 \\ \frac{dR_t}{dt} = 0.070 + \frac{\sin(\beta)}{e^{I_t}} \end{vmatrix} $
Ours	0.3561	$\frac{d\theta_t}{dt} = 0.99\omega_t$ $\frac{d\omega_t}{dt} = -0.97\alpha^2\sin\left(\theta_t\right)$	0.6194	$egin{aligned} rac{dp_t}{dt} &= 1.254p_t - 0.38q_tp_t \ rac{dq_t}{dt} &= 4.1p_t - 0.30q_t - \gamma \end{aligned}$	0.0652	$\frac{\frac{dS_t}{dt} = -1.7S_t I_t}{\frac{dI_t}{dt} = 0.42S_t I_t}$ $\frac{\frac{dR_t}{dt} = -0.0088$
GT	NAN	$\begin{aligned} \frac{d\theta_t}{dt} &= \omega_t \\ \frac{d\omega_t}{dt} &= -\alpha^2 \sin\left(\theta_t\right) \end{aligned}$	NAN	$rac{dp_t}{dt} = lpha p_t - eta p_t q_t \ rac{dq_t}{dt} = \delta p_t q_t - \gamma q_t$	NAN	$\frac{dS_t}{dt} = -\beta \frac{S_t I_t}{S_t + I_t + R_t}$ $\frac{dI_t}{dt} = \beta \frac{S_t I_t}{S_t + I_t + R_t}$ $\frac{dR_t}{dR_t} = 0$

Table 7: Symbolic regression explanation comparisons.

1441

1416 1417

transformer given different input lengths. Therefore, a shorter input length can perform slightly bettergiven the same training steps.

1444 For the environment analysis, in addition to evaluations on the full set of environments, we benchmark 1445 model performance on datasets with three and two training environments. Specifically, we select 1446 [Powered, Air, Spring] and [Powered, Air] for ME-Pendulum; [Save, Fight, Resource] and [Save, 1447 Fight] for ME-Lotka-Volterra; and [Negative, Origin, Enlarge] and [Negative, Origin] for ME-1448 SIREpidemic. While not all possible environment combinations are evaluated, these selections provide intriguing insights. As illustrated in Fig. 12, changes in the set of environments weakly 1449 affect model performance on ME-Lotka-Volterra. For ME-Pendulum, however, the inclusion of 1450 each additional environment consistently improves model performance. On ME-SIREpidemic, the 1451 performance boost observed with "3 envs" underscores the critical role of the environment Enlarge. 1452

Two key observations regarding the ME-SIREpidemic are worth noting. First, the average performance degradation on "4 envs" suggests a reduced focus on the important environment Enlarge due to the addition of the final environment Loop. Second, the improvement in the best performance candidate demonstrates the additional benefits of the environment Loop. These findings illustrate that while adding environments can enhance the best possible discovery of invariant functions, it also increases the average training complexity that may cause average performance degradations.

1458 E.3 SYMBOLIC REGRESSION EXPLANATION COMPARISONS

1460 To further evaluate the performance differences between the proposed method and the baselines, we apply PySR to the four baseline methods and obtain analytical explanations, as summarized 1461 in Tab. 7. The symbolic regression results provide an intuitive understanding of performance in 1462 relation to different NRMSE values. Specifically, when the NRMSE approaches 1, the resulting 1463 explanations are largely meaningless. As the NRMSE decreases to around 0.7, the explanations 1464 become more interpretable but may sometimes converge to oversimplified expressions, such as the 1465 IRM and VREx results on ME-Lotka-Volterra. When the NRMSE approaches zero, the expressions 1466 become more reasonable but are not always ideal due to the inherent limitations of PySR. This 1467 suggests that strong model performance does not necessarily guarantee high-quality explanations, 1468 highlighting the performance constraints of the explainer (PySR).

1469 1470 1471

1472

F EFFICIENT HYPERNETWORK IMPLEMENTATION

One of the major challenges that limits the usage of hypernetworks 1473 is the implementation complexity. Most current implementations 1474 requires either re-implementing basic neural networks (von Oswald 1475 et al., 2020) or assigning predicted weights to the main function 1476 (forecaster) one by one for each forward pass (Ortiz et al., 2023; 1477 Sudhakaran, 2022; Kirchmeyer et al., 2022). To overcome these 1478 issues, we propose a Reference-based hypernetwork implementation 1479 technique that uses pure PyTorch without introducing any new mod-1480 ules or CUDA kernels. Our proposed technique does not require 1481 reassigning weights for each sample in one forward pass, *i.e.*, for any continuous N training iterations with batch size of B and a fore-1482 caster with M parameter variables, our computation complexity is 1483 O(NM + BM), instead of O(BMN) as previous implementations. 1484

1485 Specifically, for every forward pass, we create a function parameter 1486 vector buffer $\in \mathbb{R}^m$ with fixed storage space, instead of reshaping and assigning the predicted function parameters with complexity 1487 O(BM). As shown in Fig. 13, we consider the derivative neural 1488 network parameter variables as storage space pointers, *i.e.*, the net-1489 work stores references instead of matrices. The fractions of function 1490 parameter vector buffer are pointed by these pointers; thus, once the 1491 buffer's values change by the predicted function parameters, e.g., \hat{f} , 1492 the derivative network's parameters will be changed automatically 1493 without any assignment operators. To maintain the buffer's fixed 1494 storage space, several in-place operations are applied to maintains 1495 computational graphs and gradients. 1496



Figure 13: **Reference-based** hypernetwork implementation.

F.1 EFFICIENCY COMPARISONS

	• • • • •	ee • •
Table X. Hypernetwork	implementation et	fficiency comparisons
indic of injection	mprementation	includy comparisons

501	Implementation	Vectorization	Сору	Reference	First Step Time (s)	Avg Time \pm Std (s)	Speedup
02	Non-vectorized	× ×	1	×	0.2466	0.1818 ± 0.0601	1x
03	Module-based	1	1	×	0.1768	0.1513 ± 0.0737	1.2x
	Functional-based	1	X	×	0.2013	0.0198 ± 0.0007	9.2x
5	Ours	1	×	1	0.1805	0.0108 ± 0.0006	16.8x

1506

1497

1498 1499

1500

To evaluate the efficiency of our hypernetwork implementation, we compare it against several common implementation approaches. Specifically, we measure the forward pass time of our model over 200 continuous iterations in training mode, recording both the time for the first iteration and the average time for the subsequent iterations. While the first iteration typically takes a similar amount of time across all implementations, their performance diverges significantly in the subsequent iterations. As shown in Tab. 8, the *Non-vectorized* implementation represents methods that do not vectorize the



Figure 14: A function learning example.

derivative function and therefore must run different derivative functions sequentially. Approaches
like CoDA (Kirchmeyer et al., 2022) attempt to vectorize the model by employing group-based
convolution networks. However, these module-based implementations rely on stateful PyTorch
modules, requiring the derivative function module to be replicated during each forward pass, which
slows down the process. While these *Module-based* implementations offer a slight improvement over
non-vectorized methods due to the vectorization benefits, the performance gain is limited.

In contrast, our vectorized *Functional-based* implementation leverages PyTorch's functional methods, achieving 9.2x speedup by avoiding the overhead associated with stateful modules. Note that vector-izing hypernetworks using libraries such as hypnettorch (von Oswald et al., 2020) can deliver similar speedups. Finally, our *Reference-based* implementation, which eliminates parameter assignment after the first iteration, nearly doubles the forward pass speed (16.8x) compared to implementations that require such assignments. Notably, this optimization remains applicable for potential future CUDA-based parallel hypernetwork implementations.

1533

1518

1519

1534 G SUPPLIMENTARY EXPLANATIONS

1536 The diagram 14 illustrates the conceptual function learning process associated with the DIF frame-1537 work, exemplified by the prediction of $\frac{d\omega}{dt}$. Each (hidden) function representation, such as \hat{z}_c , \hat{z}_e , \hat{f}_c , 1538 and f_e , *implicitly* encodes both a symbolic function form and a value-assignment/substitution rule. 1539 For instance, f_c represents two types of information: (1) the function form, e.g., $-\alpha^2 \sin(\theta_t)$, and (2) 1540 a partial assignment of values, such as $\{\alpha \to 1.5, \theta_t \to \}$, where α is assigned a value of 1.5 while θ_t 1541 remains a symbolic variable. This implies that \hat{f}_c represents the function $-1.5^2 \sin(\theta_t)$, which can 1542 accept input values in the form $\{\theta_t \rightarrow ?\}$, effectively evaluated as $f_c(?)$. Similarly, f is a function defined by two symbolic variables and accepts inputs in the form $\{\theta_t \rightarrow ?, \omega_t \rightarrow ?\}$. In this specific 1543 example, the input to \hat{f} is $\{\theta_t \to 1.2, \omega_t \to 0.4\}$, resulting in the evaluation f(1.2, 0.4) = -2.23. 1544

1545 1546

1547

H VISUALIZATIONS ON ME-LOTKA-VOLTERRA AND ME-SIREPIDEMIC



155



In this section, we present visualization comparisons for ME-Lotka-Volterra (Fig. 15) and ME-SIREpidemic (Fig. 16). The results for ME-SIREpidemic closely align with its quantitative findings. For the more challenging task of ME-Lotka-Volterra, our method's predicted trajectories remain closer to the ground truth. In the X^c predictions, where most methods fail, our predicted trajectory has turning points closest to the ground truth in terms of timing, although there are deviations in magnitude (Fig. 15b). The complexity of the ME-Lotka-Volterra task arises from several factors, including the introduction of exponential functions within environments, the distribution of environments, and

