

---

# Physics-Integrated Variational Autoencoders for Robust and Interpretable Generative Modeling

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Integrating physics models within machine learning models holds considerable  
2 promise toward learning robust models with improved interpretability and abilities  
3 to extrapolate. In this work, we focus on the integration of incomplete physics  
4 models into deep generative models. In particular, we introduce an architecture of  
5 variational autoencoders (VAEs) in which a part of the latent space is grounded by  
6 physics. A key technical challenge is to strike a balance between the incomplete  
7 physics and trainable components such as neural networks for ensuring that the  
8 physics part is used in a meaningful manner. To this end, we propose a regularized  
9 learning method that controls the effect of the trainable components and preserves  
10 the semantics of the physics-based latent variables as intended. We not only  
11 demonstrate generative performance improvements over a set of synthetic and real-  
12 world datasets, but we also show that we learn robust models that can consistently  
13 extrapolate beyond the training distribution in a meaningful manner. Moreover, we  
14 show that we can control the generative process in an interpretable manner.

## 15 1 Introduction

16 Data-driven modeling is often opposed to theory-driven modeling, yet their integration has also been  
17 recognized as an important approach known as *gray-box* or *hybrid* modeling. In statistical machine  
18 learning, incorporation of physics (in a broad sense; including knowledge of biology, economics, etc.)  
19 has also been attracting attention. Gray-box / hybrid modeling in machine learning holds considerable  
20 promise toward learning robust models with improved abilities to extrapolate beyond the distributions  
21 that they have been exposed to during training. Moreover, it can bring significant benefits in terms of  
22 model interpretability since parts of a model get semantically grounded to concrete inductive bias.

23 A technical challenge in gray-box deep modeling is to ensure an appropriate use of physics models.  
24 A careless design of models and learning can lead to an erratic behavior of the components meant to  
25 represent physics (e.g., with erroneous estimation of physics parameters), and eventually, the overall  
26 gray-box model just learns to ignore them. This is particularly the case when we bring together  
27 simplified or imperfect physics models with very expressive data-driven machine learning models  
28 such as deep neural networks. Such cases call for principled methods for striking an appropriate  
29 balance between physics and data-driven models to prevent the detrimental effects during learning.

30 Integration of physics models into machine learning has been considered in various contexts (see,  
31 e.g., [41, 38] and our Section 4), but most existing studies focus on prediction or forecasting tasks  
32 and are not directly applicable to other tasks. More importantly, hardly any have addressed the  
33 careful orchestration of physics-based and data-driven components to avoid the detrimental effects.  
34 A notable exception is Yin et al. [44], in which they proposed a method to harness the action of  
35 trainable components of a hybrid model of differential equations. Their method has been developed  
36 for dynamics forecasting and is limited to additive combinations of physics and trainable models.

37 In this work, we aim at the integration of (incomplete) physics models into deep generative models,  
 38 variational autoencoders (VAEs, [15, 28]) in particular, while the basic idea is applicable to other  
 39 models. In our VAE, the decoder comprises physics-based models and trainable neural networks, and  
 40 some of the latent variables are semantically grounded to the parameters of the physics models. Such  
 41 a VAE, if appropriately trained, is by construction partly interpretable. Moreover, since it can by  
 42 construction capture the underlying physics, it will be robust in out-of-distribution regime and exhibit  
 43 meaningful extrapolation properties. We propose a regularized learning framework for ensuring the  
 44 meaningful use of the physics models and the preservation of the semantics of the latent variables in  
 45 the physics-integrated VAEs. We empirically demonstrate that our method can learn a model that  
 46 exhibits better generalization, and more importantly, can extrapolate robustly in out-of-distribution  
 47 regime. In addition, we show how the direct access to the physics-grounded latent variables allows us  
 48 to alter properties of generation meaningfully and explore counterfactual scenarios.

## 49 2 Physics-integrated VAEs

50 We first describe the structure of VAEs we consider, which comprise physics models and machine  
 51 learning models such as neural nets. We suppose that the physics models can be solved analytically  
 52 or numerically with a reasonable cost, and the (approximate) solution is differentiable with regard to  
 53 the quantities on which the solution depends. This assumption holds in most physics models known  
 54 in practice, which come in different forms such as algebraic and differential equations. If there is no  
 55 closed-form solution of algebraic equations, we can utilize differentiable optimizers [3] as a layer of  
 56 the model. For differential equations, differentiable integrators [see, e.g., 7] will constitute a layer.  
 57 Handling non-differentiable and/or overly-complex simulators remains an important open challenge.

### 58 2.1 Example

59 We start with an example to demonstrate the main concepts. Let us suppose that data comprise  
 60 time-series of the angle of pendulums following an ordinary differential equation (ODE):

$$\underbrace{d^2\theta(t)/dt^2 + \omega^2 \sin \theta(t)}_{\text{given as prior knowledge, } f_P} + \underbrace{\zeta d\theta(t)/dt - u(t)}_{\text{to be learned by NN, } f_A} = 0, \quad (1)$$

61 where  $\theta$  is the pendulum’s angle, and  $\omega$ ,  $\zeta$ , and  $u$  are the pendulum’s angular velocity, damping  
 62 coefficient, and external force, respectively. We suppose that a data point  $x$  is a sequence of  $\theta(t)$ , i.e.,  
 63  $x = [\theta(0) \theta(\Delta t) \cdots \theta((\tau - 1)\Delta t)]^T \in \mathbb{R}^\tau$  for some  $\Delta t \in \mathbb{R}$  and  $\tau \in \mathbb{N}$ , where  $\theta(t)$  is the solution  
 64 of (1) with a particular configuration of  $\omega$ ,  $\zeta$ , and  $u$ . In this example, we learn a VAE on a dataset  
 65 comprising such  $x$  with different configurations of  $\omega$ ,  $\zeta$ , and  $u$ .

66 Suppose that the first two terms of (1) are given as prior knowledge, i.e., we know that the governing  
 67 equation should contain  $f_P(\theta, \omega) := \ddot{\theta} + \omega^2 \sin \theta$ . We will use such prior knowledge,  $f_P$ , by  
 68 incorporating it in the decoder of the VAE. Since  $f_P$  misses some effects of the true pendulum system  
 69 (1), we complete it by augmenting the decoder with an auxiliary function  $f_A(\theta, z_A)$ , which we model  
 70 with a neural network. The VAE’s latent variable will have two parts,  $z_P$  and  $z_A$ , respectively linked  
 71 to  $f_P$  and  $f_A$ . One one hand,  $z_A$  works as an ordinary VAE’s latent variable since  $f_A$  is a neural  
 72 net, and we suppose  $z_A \in \mathbb{R}^d$ ,  $p(z_A) := \mathcal{N}(\mathbf{0}, \mathbf{I})$ . On the other hand, we semantically ground  $z_P$   
 73 to the parameter of  $f_P$ , that is,  $z_P := \omega \in \mathbb{R}$  in this example. In summary, the augmented decoder  
 74 here is  $\mathbb{E}[x] = \text{ODEsolve}_\theta [f_P(\theta(t), z_P) + f_A(\theta(t), z_A) = 0]$ , where  $\text{ODEsolve}_\theta$  denotes some  
 75 differentiable solver of an ODE with regard to  $\theta$ . The encoder will have corresponding recognition  
 76 networks for  $z_P$  and  $z_A$ . The situation in this example will be numerically examined in Section 5.1.

### 77 2.2 General formulation

78 We now present the concept of our physics-integrated VAEs in a general form. Note that our interest is  
 79 not limited to additive cases nor ODEs. In fact, the general formulation below subsumes non-additive  
 80 augmentation of various physics models (i.e., not only ODEs). The notation introduced in this section  
 81 will be used to explain the proposed regularized learning method later in Section 3.

82 For clarity, we suppose that a VAE decoder comprises two parts: a physics-based model  $f_P$  and a  
 83 trainable auxiliary function  $f_A$ . More general cases, for example with multiple trainable functions  
 84  $f_{A,1}, f_{A,2}, \dots$  used in different ways, are handled in Appendix A.

### 85 2.2.1 Latent variables and priors

86 We consider two types of latent variables,  $\mathbf{z}_P \in \mathcal{Z}_P$  and  $\mathbf{z}_A \in \mathcal{Z}_A$ , which respectively will be used in  
 87  $f_P$  and  $f_A$ . The latent variables can be in any space, but for simplicity of discussion, we suppose  $\mathcal{Z}_P$   
 88 and  $\mathcal{Z}_A$  are (subsets of) the Euclidean space and set their prior distribution as multivariate normal:

$$p(\mathbf{z}_P) := \mathcal{N}(\mathbf{z}_P \mid \mathbf{m}_P, v_P^2 \mathbf{I}) \quad \text{and} \quad p(\mathbf{z}_A) := \mathcal{N}(\mathbf{z}_A \mid \mathbf{0}, \mathbf{I}), \quad (2)$$

89 where  $\mathbf{m}_P$  and  $v_P^2$  are defined in accordance with prior knowledge of  $f_P$ 's parameters. Note that  $\mathbf{z}_P$   
 90 will be directly interpretable as they will be semantically grounded to the parameters of the physics  
 91 model  $f_P$ ; for example in Section 2.1,  $\mathbf{z}_P := \omega$  was the angular velocity of a pendulum.

### 92 2.2.2 Decoder

93 The decoder of a physics-integrated VAE comprises two types of functions<sup>1</sup>,  $f_P: \mathcal{Z}_P \rightarrow \mathcal{Y}_P$  and  
 94  $f_A: \mathcal{Y}_P \times \mathcal{Z}_A \rightarrow \mathcal{Y}_A$ . For notational convenience, we consider a functional  $\mathcal{F}$  that evaluates  $f_P$  and  
 95  $f_A$ , solve an equation if any, and finally gives observation  $\mathbf{x} \in \mathcal{X}$ .  $\mathcal{X}$  may be the space of sequences,  
 96 images, and so on. Assuming Gaussian observation noise, we write the observation model as

$$p_\theta(\mathbf{x} \mid \mathbf{z}_P, \mathbf{z}_A) := \mathcal{N}(\mathbf{x} \mid \mathcal{F}[f_A(f_P(\mathbf{z}_P), \mathbf{z}_A)], \Sigma_x). \quad (3)$$

97 Note that  $f$  may have other arguments besides  $\mathbf{z}$ , but they are omitted for simplicity. We denote the  
 98 set of trainable parameters of  $f_A$  and  $f_P$  (and  $\Sigma_x$ ) by  $\theta$ , while  $f_P$  may have no trainable parameters.

99 Let us see the semantics of  $\mathcal{F}$  first in the light of the example of Section 2.1. Recall that there we  
 100 considered the additive augmentation of ODE (as in [44] and other studies). It is subsumed by the  
 101 expression (3) by setting  $f_A(f_P(\mathbf{z}_P), \mathbf{z}_A) := f_P(\mathbf{z}_P) + f_{A'}(\mathbf{z}_A)$  and  $\mathcal{F}[f] := \text{ODEsolve}[f = 0]$ ,  
 102 where  $f_{A'}$  is a neural network. Let us generalize the idea. Our definition of the decoder in (3) allows  
 103 not only additive augmentation of ODE but also broader range of architectures. The composition  
 104 of  $f_P$  and  $f_A$  is *not* limited to be additive because we consider general function composition  
 105  $f_A(f_P(\mathbf{z}_P), \mathbf{z}_A)$ . Moreover, the form of the physics model is *not* limited to ODEs:

- 106 • If equation  $f_P = 0$  has a closed-form solution  $S_{f_P}$ , then  $\mathcal{F}$  is simply, e.g.,  $\mathcal{F}[f_P, f_A] := f_A(S_{f_P})$ .
- 107 • If an algebraic equation  $f_P = 0$  or  $f_A \circ f_P = 0$  has no closed-form solution, then  $\mathcal{F}$  will have a  
 108 differentiable optimizer, e.g.,  $\mathcal{F}[f_P, f_A] := f_A(\arg \min \|f_P\|^2)$  or  $\mathcal{F} := \arg \min \|f_A \circ f_P\|^2$ .
- 109 •  $f_P = 0$  or  $f_A \circ f_P = 0$  can be a stochastic differential equation (and  $\mathcal{F}$  contains its solver), for  
 110 which  $\mathbf{z}_P$  and/or  $\mathbf{z}_A$  would become a sequence encoding the realization of the process noise.

111 The role of  $f_A$  can also be diverse; it can work not only as a complement of physics models inside  
 112 equations, but also as correction of numerical errors of solvers or optimizers, downsampling or  
 113 upsampling, and observables (e.g., from angle sequence to video of a pendulum).

### 114 2.2.3 Encoder

115 The encoder of a physics-integrated VAE accordingly comprises two parts: posterior inference of  $\mathbf{z}_P$   
 116 and that of  $\mathbf{z}_A$ . We consider the following decomposition of the approximated posterior:

$$q_\psi(\mathbf{z}_P, \mathbf{z}_A, \mid \mathbf{x}) := q_\psi(\mathbf{z}_A \mid \mathbf{x}) q_\psi(\mathbf{z}_P \mid \mathbf{x}, \mathbf{z}_A), \quad (4)$$

where  $q_\psi(\mathbf{z}_A \mid \mathbf{x}) := \mathcal{N}(\mathbf{z}_A \mid g_A(\mathbf{x}), \Sigma_A)$ ,  $q_\psi(\mathbf{z}_P \mid \mathbf{x}, \mathbf{z}_A) := \mathcal{N}(\mathbf{z}_P \mid g_P(\mathbf{x}, \mathbf{z}_A), \Sigma_P)$ .

117  $g_A: \mathcal{X} \rightarrow \mathcal{Z}_A$  and  $g_P: \mathcal{X} \times \mathcal{Z}_A \rightarrow \mathcal{Z}_P$  are recognition networks. We denote the trainable parameters  
 118 of  $g_A$  and  $g_P$  (and  $\Sigma_A$  and  $\Sigma_P$ ) as  $\psi$ . This particular dependency is for our regularization method in  
 119 Section 3.2, where  $g_P$  should first remove the information of  $\mathbf{z}_A$  from  $\mathbf{x}$  and then infer  $\mathbf{z}_P$ .

### 120 2.3 Evidence lower bound

121 The VAE is to be learned as usual by maximizing the lower bound of the marginal log likelihood  
 122 known as evidence lower bound (ELBO). In our case, it is straightforward to derive:

$$\begin{aligned} \text{ELBO}(\theta, \psi; \mathbf{x}) &= \mathbb{E}_{q_\psi(\mathbf{z}_P, \mathbf{z}_A \mid \mathbf{x})} \log p_\theta(\mathbf{x} \mid \mathbf{z}_P, \mathbf{z}_A) \\ &\quad + D_{\text{KL}}[q_\psi(\mathbf{z}_A \mid \mathbf{x}) \parallel p(\mathbf{z}_A)] + \mathbb{E}_{q_\psi(\mathbf{z}_A \mid \mathbf{x})} D_{\text{KL}}[q_\psi(\mathbf{z}_P \mid \mathbf{x}, \mathbf{z}_A) \parallel p(\mathbf{z}_P)]. \end{aligned} \quad (5)$$

<sup>1</sup>The distinction between  $f_P$  and  $f_A$  depends on the origin of the functional forms (and not if trainable or not). The form of  $f_P$  depends on physics' insight and thus fixed. On the other hand, the form of  $f_A$  is determined only from utility as a function approximator, and we can use whatever useful (e.g., feed-forward NNs, RNNs, etc.).

### 123 3 Regularizing physics-integrated VAEs

124 We propose a regularized learning objective for physics-integrated VAEs. It comprises two types of  
 125 regularizers. The first one is for harnessing unnecessary flexibility of function approximators like  
 126 neural networks and presented in Section 3.1. The second ones are for grounding encoder’s output to  
 127 physics parameters and presented in Section 3.2. The overall objective is summarized in Section 3.3.

#### 128 3.1 Harnessing trainable functions by PPC-like procedure

129 If the trainable component of the physics-integrated VAE (i.e.,  $f_A$ ) has rich expression capability,  
 130 as is often the case with deep neural networks, merely maximizing the ELBO in (5) provides no  
 131 guarantee that the physics-based component (i.e.,  $f_P$ ) will be used in a meaningful manner; e.g.,  $f_P$   
 132 may just be ignored. We want to ensure that  $f_A$  does not unnecessarily dominate the behavior of the  
 133 entire model and that  $f_P$  is not ignored. To this end, we borrow an idea from the *posterior predictive*  
 134 *check* (PPC), a procedure to check the validity of a statistical model [see, e.g., 9]. Whereas the  
 135 standard PPCs examine the discrepancy between model’s and data’s posterior predictive distributions,  
 136 we compute the discrepancy between those of the physics-integrated model and its “physics-only”  
 137 reduced model for monitoring and balancing the contributions of parts of the model.

138 For the sake of argument, suppose that a given physics model  $f_P$  is completely correct for given data.  
 139 Then, the discrepancy between the original model and its “physics-only” reduced model (where  $f_A$  is  
 140 somehow invalidated) should be close to zero because the decoder of both the original model (with  
 141  $f_P$  and  $f_A$  working) and the reduced model (with only  $f_P$  working) should coincide in an ideal limit  
 142 with the true data-generating process. Even if  $f_P$  captures only a part of the truth, the discrepancy  
 143 should be kept small, if not zero, to ensure meaningful use of the physics models in the overall model.

144 The “physics-only” reduced model is created as follows. Recall that the original VAE is defined by  
 145 Eqs. (3) and (4). We define the decoder of the reduced model by replacing  $f_A: \mathcal{Y}_P \times \mathcal{Z}_A \rightarrow \mathcal{Y}_A$  of  
 146 (3) with a *baseline function*  $h_A: \mathcal{Y}_P \rightarrow \mathcal{Y}_A$ . That is, the reduced observation model is

$$147 p_{\theta, \theta^r}^r(\mathbf{x} \mid \mathbf{z}_P, \mathbf{z}_A) := \mathcal{N}(\mathbf{x} \mid \mathcal{F}[h_A(f_P(\mathbf{z}_P))], \Sigma_{\mathbf{x}}). \quad (3r)$$

148 We denote the set of the trainable parameters of  $h_A$  as  $\theta^r$ , while it may often be empty. The  
 149 corresponding encoder is defined as follows. Recall that in the original model, posterior distributions  
 150 of both  $\mathbf{z}_P$  and  $\mathbf{z}_A$  are inferred in (4) and then used for reconstructing each input  $\mathbf{x}$  in (3). On the  
 151 other hand, in the “physics-only” reduced model,  $\mathbf{z}_A$  is not referred to by (3r), which makes it less  
 152 meaningful to place a particular posterior of  $\mathbf{z}_A$  for each  $\mathbf{x}$ . Hence, we define the “physics-only”  
 153 encoder by marginalizing out  $\mathbf{z}_A$  and using prior<sup>2</sup>  $p(\mathbf{z}_A)$  instead. That is, the reduced posterior is

$$154 q_{\psi}^r(\mathbf{z}_A, \mathbf{z}_P \mid \mathbf{x}) := p(\mathbf{z}_A) \int q_{\psi}(\mathbf{z}_P, \mathbf{z}_A, \mid \mathbf{x}) d\mathbf{z}_A. \quad (4r)$$

155 Below we give a guideline for the choice of the baseline function,  $h_A$ :

- 156 • If the ranges of  $f_P$  and  $f_A$  are the same (i.e.,  $\mathcal{Y}_P = \mathcal{Y}_A$ ), then  $h_A$  can be an identity function  
 157  $h_A = \text{Id}$ . Note that in the additive case  $f_A \circ f_P = f_P + f_{A'}$ , where  $f_{A'}$  is a trainable function,  
 158 replacing  $f_A$  with  $h_A = \text{Id}$  is equivalent to replacing  $f_{A'}$  with  $h_{A'} = 0$ .
- 159 • If  $\mathcal{Y}_P \neq \mathcal{Y}_A$ , then  $h_A$  can be a linear or affine map from  $\mathcal{Y}_P$  to  $\mathcal{Y}_A$ . For example, if  $\mathcal{Y}_P = \mathbb{R}^{d_P}$  and  
 160  $\mathcal{Y}_A = \mathbb{R}^{d_A}$  ( $d_P \neq d_A$ ), then we can set  $h_A(f_P(\mathbf{z}_P)) = \mathbf{W} f_P(\mathbf{z}_P)$  where  $\mathbf{W} \in \mathbb{R}^{d_A \times d_P}$ .

161 The idea is to minimize the discrepancy between the full model and the “physics-only” reduced  
 162 model. In particular, we minimize the discrepancy between the posterior predictive distributions

$$163 D_{\text{KL}}[p_{\theta, \psi}(\tilde{\mathbf{x}} \mid X) \parallel p_{\theta, \theta^r, \psi}^r(\tilde{\mathbf{x}} \mid X)], \quad \text{where}$$

$$164 p_{\theta, \psi}(\tilde{\mathbf{x}} \mid X) = \int p_{\theta}(\tilde{\mathbf{x}} \mid \mathbf{z}_P, \mathbf{z}_A) q_{\psi}(\mathbf{z}_P, \mathbf{z}_A \mid \mathbf{x}) p_d(\mathbf{x} \mid X) d\mathbf{z}_P d\mathbf{z}_A d\mathbf{x}, \quad (6)$$

$$165 p_{\theta, \theta^r, \psi}^r(\tilde{\mathbf{x}} \mid X) = \int p_{\theta, \theta^r}^r(\tilde{\mathbf{x}} \mid \mathbf{z}_P, \mathbf{z}_A) q_{\psi}^r(\mathbf{z}_P, \mathbf{z}_A \mid \mathbf{x}) p_d(\mathbf{x} \mid X) d\mathbf{z}_P d\mathbf{z}_A d\mathbf{x}.$$

166  $p_d(\mathbf{x} \mid X)$  is the empirical distribution with the support on data  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ . We use  $\tilde{\mathbf{x}}$   
 167 (instead of  $\mathbf{x}$ ) just for avoiding notational confusion by clarifying the target of integral  $\int d\mathbf{x}$ .

<sup>2</sup>It is just for defining  $q_{\psi}^r$  on the common support with  $q_{\psi}$ . Any non-informative distributions of  $\mathbf{z}_A$  are fine.

163 Unfortunately, analytically computing (6) is usually intractable. Hence, we take the following upper  
 164 bound of (6) (a proof is in Appendix B, and further remarks are in Appendix C):

165 **Proposition 1.** *Let  $p_\theta$  and  $p_\theta^r$  be the shorthand of  $p_\theta(\tilde{\mathbf{x}} \mid \mathbf{z}_P, \mathbf{z}_A)$  in (3) and  $p_{\theta^r, \psi}^r(\tilde{\mathbf{x}} \mid \mathbf{z}_P, \mathbf{z}_A)$  in  
 166 (3r), respectively. Let  $p_P$  and  $p_A$  be some distributions of  $\mathbf{z}_P$  and  $\mathbf{z}_A$ , e.g.,  $p(\mathbf{z}_P)$  and  $p(\mathbf{z}_A)$  using  
 167 the priors in (2), respectively. The KL divergence in (6) can be upper bounded as follows:*

$$D_{\text{KL}}[p_{\theta, \psi}(\tilde{\mathbf{x}} \mid X) \parallel p_{\theta^r, \psi}^r(\tilde{\mathbf{x}} \mid X)] \leq \mathbb{E}_{p_d(\mathbf{x} \mid X)} \left[ \mathbb{E}_{q_\psi(\mathbf{z}_P, \mathbf{z}_A \mid \mathbf{x})} D_{\text{KL}}[p_\theta \parallel p_\theta^r] \right. \\ \left. + D_{\text{KL}}[q_\psi(\mathbf{z}_A \mid \mathbf{x}) \parallel p_A] + \mathbb{E}_{q_\psi(\mathbf{z}_A \mid \mathbf{x})} D_{\text{KL}}[q_\psi(\mathbf{z}_P \mid \mathbf{z}_A, \mathbf{x}) \parallel p_P] \right]. \quad (7)$$

168 **Definition 1.** Let us denote the upper bound (7) by  $\mathbb{E}_{p_d(\mathbf{x} \mid X)} \hat{D}(\theta, \theta^r, \psi; \mathbf{x})$ . The regularization for  
 169 harnessing unnecessary flexibility of trainable functions is defined as minimization of

$$R_{\text{PPC}}(\theta, \theta^r, \psi) := \mathbb{E}_{p_d(\mathbf{x} \mid X)} \hat{D}(\theta, \theta^r, \psi; \mathbf{x}). \quad (8)$$

170 *Remark 1.* When multiple trainable functions are differently used in a model (e.g., inside *and* outside  
 171 an equation solver), which is often the case in practice, the definition of  $R_{\text{PPC}}$  should be generalized  
 172 to consider marginal contribution of every trainable function. See Appendix A.

### 173 3.2 Grounding physics encoder by physics-based data augmentation

174 Toward properly learning physics-integrated VAEs, minimizing  $R_{\text{PPC}}$  solely may not be enough  
 175 because inferred  $\mathbf{z}_P$  may be still meaningless but makes  $R_{\text{PPC}}$  not that large (e.g., with solution of  
 176  $f_P$  fluctuating around the mean pattern of data). Though it is difficult to avoid such a local solution  
 177 perfectly, we can alleviate the situation by considering additional objectives to encourage a proper use  
 178 of the physics. The idea is to use the physics model as a source of information for data augmentation,  
 179 which helps us to ground the output of the recognition network,  $g_P$  in (4), to the parameters of  $f_P$ .

180 Let  $\mathbf{z}_P^*$  be a sample drawn from some distribution of  $\mathbf{z}_P$  (e.g., prior  $p(\mathbf{z}_P)$ ). We artificially generate sig-  
 181 nals  $\mathbf{x}^*$  by feeding  $\mathbf{z}_P^*$  to the “physics-only” decoding process in (3r), that is,  $\mathbf{x}_{\mathbf{z}_P^*}^* := \mathcal{F}[h_A(f_P(\mathbf{z}_P^*))]$ .  
 182 We want the physics-part recognition network,  $g_P$ , to successfully estimate  $\mathbf{z}_P^*$  given the correspond-  
 183 ing  $\mathbf{x}_{\mathbf{z}_P^*}^*$ , which is necessary to say that the result of the inference by  $g_P$  is grounded to the parameters  
 184 of  $f_P$ . However, in general, real data  $\mathbf{x}$  and the augmented data  $\mathbf{x}^*$  have different natures because  $f_P$   
 185 may miss some aspects of the true data-generating process. We handle this issue by considering a  
 186 specific design of the physics-part recognition network,  $g_P$ .

187 Let us decompose  $g_P$  as  $g_P(\mathbf{x}, \mathbf{z}_A) = g_{P,2}(g_{P,1}(\mathbf{x}, \mathbf{z}_A))$  without loss of generality. On one hand,  
 188  $g_{P,1}$  should transform real data  $\mathbf{x}$  into  $\mathbf{x}'$  such that  $\mathbf{x}'$  *resembles the physics-based augmented signal*  
 189  $\mathbf{x}^*$ . In other words,  $g_{P,1}$  should “cleanse” real data into a virtual “physics-only” counterpart. On  
 190 the other hand,  $g_{P,2}$  should receive such “cleansed” data  $\mathbf{x}'$  and return the (sufficient statistics of)  
 191 posterior of  $\mathbf{z}_P$ . As  $g_{P,2}$  works on  $\mathbf{x}'$ , which should resemble  $\mathbf{x}^*$ , we can directly self-supervise  $g_{P,2}$   
 192 with  $\mathbf{x}^*$ . We define a couple of regularizers for setting such functionality of  $g_{P,1}$  and  $g_{P,2}$  as follows:

193 **Definition 2.** Let  $\text{sg}(\cdot)$  be the stop-gradient operator. Let  $\mathbf{x}' := \mathcal{F}[h_A(f_P(g_P(\mathbf{x}, \mathbf{z}_A)))]$ . The  
 194 regularization for the physics-based data augmentation is defined as minimization of

$$R_{\text{DA},1}(\psi) := \mathbb{E}_{p_d(\mathbf{x} \mid X)q(\mathbf{z}_A \mid \mathbf{x})} \|g_{P,1}(\mathbf{x}, \mathbf{z}_A) - \text{sg } \mathbf{x}'\|_2^2 \quad \text{and} \quad (9)$$

$$R_{\text{DA},2}(\psi) := \mathbb{E}_{\mathbf{z}_P^*} \|g_{P,2}(\text{sg } \mathbf{x}_{\mathbf{z}_P^*}^*) - \mathbf{z}_P^*\|_2^2. \quad (10)$$

195 *Remark 2.* If both  $g_{P,1}$  and  $g_{P,2}$  work as intended (i.e., both  $R_{\text{DA},1}$  and  $R_{\text{DA},2}$  are small enough),  
 196  $\mathbf{x}'$  is the virtual “physics-only” counterpart of  $\mathbf{x}$ .  $R_{\text{DA},1}$  is for ensuring the functionality of  $g_{P,1}$  to  
 197 “cleanse”  $\mathbf{x}$  to  $\mathbf{x}'$ .  $R_{\text{DA},2}$  is for giving the supervision to  $g_{P,2}$  with the augmented data ( $\mathbf{z}_P^*, \mathbf{x}^*$ ).

### 198 3.3 Overall regularized learning objective

199 The overall regularized learning problem of the proposed physics-integrated VAEs is as follows:

$$\underset{\theta, \theta^r, \psi}{\text{minimize}} \quad -\mathbb{E}_{p_d(\mathbf{x} \mid X)} \text{ELBO}(\theta, \psi; \mathbf{x}) + \alpha R_{\text{PPC}}(\theta, \theta^r, \psi) + \beta R_{\text{DA},1}(\psi) + \gamma R_{\text{DA},2}(\psi),$$

200 where each term appears in (5), (8), (9), and (10), respectively. Recall that  $\theta$ ,  $\psi$ , and  $\theta^r$  are the sets of  
 201 the parameters of the full model’s decoder (3), encoder (4), and the reduced model’s decoder (3r),  
 202 respectively, while  $\theta^r$  may be empty. If we cannot specify a reasonable sampling distribution of  $\mathbf{z}_P^*$   
 203 needed in (10), we do not compute  $R_{\text{DA},1}$  and  $R_{\text{DA},2}$  and set  $\beta = \gamma = 0$ ; it may happen when the  
 204 semantics of  $\mathbf{z}_P$  are not inherently grounded, e.g., when  $f_P$  is a *neural* Hamilton’s equation [37].

## 205 4 Related work

206 The integration of theory-driven and data-driven methodologies has been sought in various ways.  
207 Ones in model design, which we followed, are one of the key approaches. Other approaches have also  
208 been studied; e.g., physics-informed neural nets (PINNs) [27] incorporate physics knowledge in the  
209 definition of loss function. We overview these perspectives in this section and more in Appendix D.

210 **Physics+ML in model design** Integration in model design, often called gray-box or hybrid mod-  
211 eling, has been a subject of study for decades [e.g., 24, 29, 36] and is still active, with deep neural  
212 networks employed in various applications [e.g., 45, 26, 21, 39, 23, 1, 2, 8, 46, 40, 32, 16, 22, 5,  
213 33, 25, 19]. Most recent studies focus on prediction, and the generative modeling has been less  
214 investigated. Moreover, mechanisms to harness trainable components have hardly been addressed.

215 The work of Yin et al. [44] is notable here because they consider a mechanism to harness a trainable  
216 component to preserve the utility of physics in the model, even though it is only focused on dynamics  
217 learning for forecasting. They learn an additive hybrid ODE model  $\dot{x} = f_P(x) + f_A(x)$ , where  $f_P$  is  
218 a prescribed physics model, and  $f_A$  is a neural network. Such a model is subsumed in our architecture  
219 as exemplified in Section 2. Moreover, Yin et al. [44] propose to harness  $f_A$  by minimizing  $\|f_A\|_2^2$ .  
220 Such a term also appears in one of our regularizers,  $R_{PPC}$ ; when the observation noise is Gaussian,  
221 the first term of the rhs of (7) becomes  $\mathbb{E}\|(f_A \circ f_P) - f_P\|_2^2 = \mathbb{E}\|f_P + f_{A'} - f_P\|_2^2 = \mathbb{E}\|f_{A'}\|_2^2$ .  
222 Therefore, we get a “VAE variant” of Yin et al. [44] by switching off a part of  $R_{PPC}$  and the other  
223 regularizers,  $R_{DA,1}$  and  $R_{DA,2}$ . We examine cases similar to it in our experiment for comparison.

224 Yıldız et al. [43] and Linial et al. [20] developed VAEs whose latent variable follows ODEs. Linial  
225 et al. [20] also suggest grounding the semantics of the latent variable by providing sparse supervision  
226 on it. It is feasible only when we have a chance to observe the latent variable (e.g., with an increased  
227 cost) and may often be inherently infeasible in some problem settings including ours. In our method,  
228 we never assume availability of observation of latent variables and instead use the physics models in  
229 a self-supervised manner. While direct comparison is not meaningful due to the difference of settings,  
230 we examine a baseline close to the base model of Linial et al. [20] in our experiment for comparison.

231 Toth et al. [37] propose a model where the latent variable sequence is governed by the Hamiltonian  
232 mechanics with a neural Hamiltonian. While it does not suppose very specific physics models  
233 but considers general mechanics, they can also be included in our framework; that is,  $f_P$  can be a  
234 Hamilton’s equation with a neural Hamiltonian. We try such a model in one of our experiments.

235 **Physics+ML in objective design** Another prevailing strategy is to define objective functions based  
236 on physics knowledge [e.g., 34, 14, 27, 12, 42, 13, 47, 30, 6]. In generative modeling, for example,  
237 Stinis et al. [35] use residuals from physics models as a feature of GAN’s discriminator. Golany et al.  
238 [10] regularize the generation from GANs by forcing it close to a prescribed physics relation. These  
239 approaches are often easy to deploy, but an inherent limitation is that given physics knowledge should  
240 be complete to some extent, otherwise a physics-based loss is not well-defined.

## 241 5 Experiments

242 We performed experiments on two synthetic datasets and two real-world datasets, for which we  
243 prepared instances of physics-integrated VAEs. We show each particular architecture of physics-  
244 integrated VAEs and the corresponding results; some details are deferred to Appendix E. While direct  
245 comparison is impossible due to the differences of the problem settings, the baseline methods we  
246 examined (listed below) are similar to some existing methods [4, 43, 37, 20, 44].

247 NN-only	Ordinary VAE [15, 28]; the decoder is $\mathbb{E}x = f_A(z_A)$ , where $f_A$ is a neural net.
248 Phys-only	Physics VAE; the decoder is $\mathbb{E}x = \mathcal{F}(f_P(z_P))$ , while the encoder is with neural 249 nets as usual. Almost equivalent to Aragon-Calvo and Carvajal [4] in Section 5.3.
250 NN+solver	VAE with physics solvers; the decoder is $\mathbb{E}x = \mathcal{F}(f_A(z_A))$ , where $f_A$ is a neural 251 net, and $\mathcal{F}$ includes some equation-solving process (e.g., ODE/PDE solver). It is 252 similar to the methods of, for example, Yıldız et al. [43] and Toth et al. [37].
253 NN+phys	Physics-integrated VAE learned without the regularizers (i.e., $\alpha = \beta = \gamma = 0$ ); 254 almost equivalent to the base model of Linial et al. [20]. Finer ablations are also 255 studied, among which the cases with $\beta = 0$ or $\gamma = 0$ are similar to Yin et al. [44].

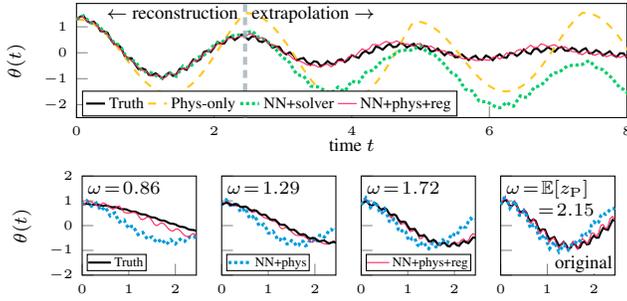


Figure 1: Reconstruction and extrapolation of a test sample of the pendulum data. Range  $0 \leq t < 2.5$  is reconstruction, whereas  $t \geq 2.5$  is extrapolation.

Figure 2: Counterfactual generation for the pendulum data. Horizontal axis is time  $t$ . The center panel shows the original data, and the rest is the generation with  $z_P$  (i.e.,  $\omega$ ) altered while  $z_A$  fixed.

256 **NN+phys+reg** Our proposal; physics-integrated VAE learned with the proposed regularizers.

257 We aligned the total dimensionality of the latent variables of each method (except **phys-only**);  
 258 when  $\dim z_A = d_A$  and  $\dim z_P = d_P$  in **NN+phys+reg**, we set  $\dim z_A = d_A + d_P$  in **NN-only** and  
 259 **NN+solver**. The hyperparameters,  $\alpha$ ,  $\beta$ , and  $\gamma$ , were chosen with validation set performance. We  
 260 investigated the performance sensitivity to them. No large degradation of performance was observed  
 261 even if we changed the values by  $\times 10$  or  $\times \frac{1}{10}$  from the chosen values; details are in Appendix F.

## 262 5.1 Forced damped pendulum

263 **Dataset** We generated data from (1) with  $u(t) = A\omega^2 \cos(2\pi\phi t)$ . Each data-point  $\mathbf{x}$  is a sequence  
 264  $\mathbf{x} := [\theta_1 \cdots \theta_\tau] \in \mathbb{R}^\tau$ , where  $\theta_j$  is the value of a solution  $\theta(t_j)$  at  $t_j := (j-1)\Delta t$ . We randomly  
 265 drew a sample of the initial condition  $\theta_1$  (with  $\dot{\theta}_1 = 0$  fixed) and the values of  $\omega$ ,  $\zeta$ ,  $A$ , and  $\phi$  for each  
 266 sequence. We generated 2,500 sequences of length  $\tau = 50$  with  $\Delta t = 0.05$  and separated them into  
 267 a training, validation, and test sets with 1,000, 500, and 1,000 sequences, respectively.

268 **Setting** We set  $f_P$  as in Section 2.1, i.e.,  $f_P(\theta, z_P) := \ddot{\theta} + z_P^2 \sin(\theta)$ , where  $z_P \in \mathbb{R}$  should  
 269 work as angular velocity  $\omega$ . We augmented it by  $f_{A,1}(\theta, z_{A,1})$  additively, where  $f_{A,1}$  was a multi-  
 270 layer perceptron (MLP) and  $z_{A,1} \in \mathbb{R}$ . The ODE  $f_P + f_{A,1} = 0$  is solved with the Euler update  
 271 scheme in the model. The model has another MLP<sup>3</sup>  $f_{A,2}$  with another latent variable  $z_{A,2} \in$   
 272  $\mathbb{R}^2$  for further modifying the solution of the ODE. In summary, the decoding process is  $\mathcal{F} :=$   
 273  $f_{A,2}(\text{solve}_\theta[f_P(\theta, z_P) + f_{A,1}(\theta, z_{A,1}) = 0], z_{A,2})$ . The construction of the proposed regularizer for  
 274 such multiple  $f_A$ 's is elaborated in Appendix A. We used  $h_{A,1} = 0$  and  $h_{A,2} = \text{Id}$  as the baseline  
 275 functions. The recognition networks  $g$  were modeled with MLPs. We used the initial element of each  
 276  $\mathbf{x}$  as an estimation of the initial condition  $\theta_1$ .

277 **Results** Figure 1 demonstrates a unique benefit of the hybrid modeling. We show an example  
 278 of reconstruction with extrapolation. Recall that the training data comprise sequences of range  
 279  $0 \leq t < 2.5$  only; so the results in  $t \geq 2.5$  are extrapolation (in time) rather than mere reconstruction.  
 280 We can observe that while **NN+solver** cannot extrapolate even if it is equipped with a neural ODE,  
 281 **NN+phys+reg** can reconstruct and extrapolate correctly.

282 Figure 2 illustrates well the advantage of the proposed regularizers. We show an example of generation  
 283 from learned models with  $z_P$  manipulated. Recall that  $z_P$  is expected to work as pendulum's angular  
 284 velocity  $\omega$ . We took a test sample with  $\omega \approx \mathbb{E}[z_P] \approx 2.15$  and generated signals with the original  
 285 and different values of  $z_P$ , keeping the values of  $z_A$  to be the original posterior mean. We can see  
 286 that the generation from **NN+phys+reg** matches better with the signals from the true process.

287 Table 1 (left half) summarizes the performance in terms of the reconstruction error and the inference  
 288 error of physics parameter  $\omega$  on the test set. The errors are reported in mean absolute errors (MAEs).  
 289 The inference error of  $\omega$  is evaluated by  $|\mathbb{E}[z_P] - \omega_{\text{true}}|$ . **NN+phys+reg** achieves small values in *both*  
 290 reconstruction error and inference error. The MAE of  $\omega$  inferred by **NN+phys** is significantly worse  
 291 than the others, which indicates the importance of the proposed regularizers.

<sup>3</sup>We used MLP as the data are fixed length. The same holds hereafter. Extension to other networks is easy.

Table 1: Reconstruction errors and inference errors on test sets of the pendulum data and the advection-diffusion data. Averages (and SDs) over 20 random trials are reported.

		Pendulum				Advection-diffusion			
		MAE of reconstr.		MAE of inferred $\omega$		MAE of reconstr.		MAE of inferred $a$	
NN-only		0.438	( $2.9 \times 10^{-2}$ )	–	–	0.0396	( $2.2 \times 10^{-4}$ )	–	–
Phys-only		1.55	( $7.1 \times 10^{-4}$ )	0.232	( $5.9 \times 10^{-3}$ )	0.393	( $9.5 \times 10^{-4}$ )	0.0103	( $1.5 \times 10^{-3}$ )
NN+solver		0.439	( $2.3 \times 10^{-2}$ )	–	–	0.0388	( $1.7 \times 10^{-4}$ )	–	–
NN+phys		0.370	( $4.3 \times 10^{-2}$ )	1.04	( $2.2 \times 10^{-1}$ )	0.0404	( $1.2 \times 10^{-2}$ )	0.258	( $3.2 \times 10^{-1}$ )
NN+phys+reg		0.363	( $4.8 \times 10^{-2}$ )	0.229	( $3.8 \times 10^{-2}$ )	0.0437	( $1.5 \times 10^{-3}$ )	0.00951	( $6.2 \times 10^{-3}$ )
Ablations	$\alpha = 0$	0.396	( $4.3 \times 10^{-2}$ )	0.889	( $1.9 \times 10^{-1}$ )	0.0461	( $1.3 \times 10^{-2}$ )	0.0444	( $1.4 \times 10^{-2}$ )
	$\beta = 0$	0.372	( $4.1 \times 10^{-2}$ )	0.223	( $3.6 \times 10^{-2}$ )	0.0747	( $2.4 \times 10^{-2}$ )	0.199	( $2.3 \times 10^{-1}$ )
	$\gamma = 0$	0.381	( $4.1 \times 10^{-2}$ )	0.276	( $4.2 \times 10^{-2}$ )	0.0588	( $9.1 \times 10^{-4}$ )	0.0548	( $9.4 \times 10^{-7}$ )

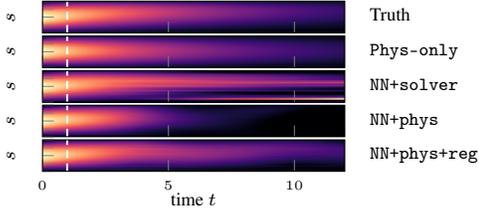


Figure 3: Reconstruction and extrapolation of a test sample of the advection-diffusion data. Range  $0 \leq t < 1$  is reconstruction, whereas  $t \geq 1$  is extrapolation; dashed line is the border.

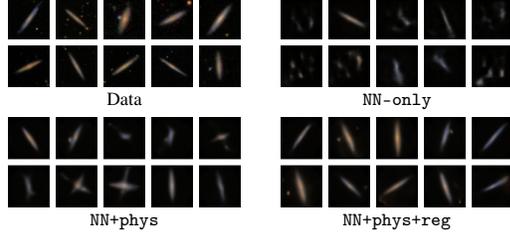


Figure 4: (left) Subset of the galaxy image data. (right three) Random generation from the NN-only model and the NN+phys+reg models.

## 292 5.2 Advection-diffusion system

293 **Dataset** We generated data from advection-diffusion PDE  $\partial T / \partial t - a \cdot \partial^2 T / \partial s^2 + b \cdot \partial T / \partial s = 0$ ,  
 294 where  $s$  is the 1-D spatial dimension. We approximated the solution  $T(s, t)$  on the 12-point even grid  
 295 from  $s = 0$  to  $s = s_{\max}$ , so each data-point  $\mathbf{x}$  is a sequence of 12-dim vectors, i.e.,  $\mathbf{x} := [T_1 \cdots T_\tau] \in$   
 296  $\mathbb{R}^{12 \times \tau}$ , where  $T_j := [T(0, t_j) \cdots T(s_{\max}, t_j)]^T$  at  $t_j := (j - 1)\Delta t$ . We set the boundary condition  
 297 as  $T(0, t) = T(s_{\max}, t) = 0$  and the initial condition as  $T(s, 0) = c \sin(\pi s / s_{\max})$ . We randomly  
 298 drew  $a, b$ , and  $c$  for each  $\mathbf{x}$ . We generated 2,500 sequences with  $\tau = 50$  and  $\Delta t = 0.02$  and separated  
 299 them into a training, validation, and test sets with 1,000, 500, and 1,000 sequences, respectively.

300 **Setting** We set  $f_P$  as the diffusion PDE, i.e.,  $f_P(T, z_P) := \partial T / \partial t - z_P \partial^2 T / \partial s^2$ , where  $z_P \in \mathbb{R}$   
 301 should work as diffusion coefficient  $a$ . We augmented it by  $f_A(T, z_A)$  additively, where  $f_A$  was an  
 302 MLP and  $z_A \in \mathbb{R}^4$ . Hence, the decoding process is  $\mathcal{F} := \text{solve}_T[f_P(T, z_P) + f_A(T, z_A) = 0]$ . We  
 303 used  $h_A = 0$  as the baseline function. The recognition networks  $g$  were modeled with MLPs. We  
 304 used the initial snapshot of each sequence  $\mathbf{x}$  as an estimation of the initial condition  $T_1$ .

305 **Results** Figure 3 shows an example of reconstruction with extrapolation. As the training data  
 306 only comprise sequences of range  $0 \leq t < 1$ , the remaining range  $t \geq 1$  is extrapolation. Only  
 307 NN+phys+reg (the bottom panel) achieves adequate extrapolation; phys-only lacks advection,  
 308 NN+solver has unnatural artifacts, and NN+phys infers  $z_P$  (i.e., diffusion coefficient  $a$ ) wrongly.

309 Table 1 (right half) summarizes the reconstruction and inference errors, which are consistent with  
 310 the results in the pendulum example. We also show the performance of ablations of NN+phys+reg,  
 311 where either of the regularizers was turned off (i.e.,  $\alpha = 0$ ,  $\beta = 0$ , or  $\gamma = 0$ ). Not surprisingly their  
 312 performance is worse than the full regularization, especially in terms of the inference error.

## 313 5.3 Galaxy images

314 **Dataset** We used images of galaxy of the Galaxy10 dataset [18]. We selected the 589 images of the  
 315 “Disk, Edge-on, No Bulge” class and separated them into training, validation, and test sets with 400,  
 316 100, and 89 images, respectively. Each image is of size  $69 \times 69$  with three channels. We performed  
 317 data augmentation with random rotation and increased the size of the training set by 20 times.

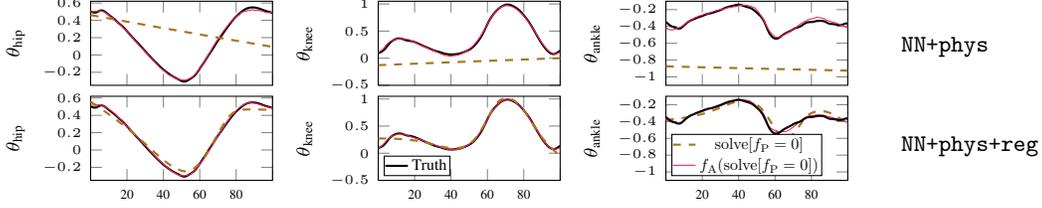


Figure 5: Reconstruction of a test sample of the gait data. Horizontal axis is normalized time.

318 **Setting** We set  $f_P: \mathbb{R}_{>0}^4 \rightarrow \mathbb{R}^{69 \times 69}$  as an exponential profile of the light distribution of galaxies  
319 [see 4, and references therein] whose input is  $z_P := [I_0 \ A \ B \ \theta]^T \in \mathbb{R}_{>0}^4$ . Let  $[f_P(z_P)]_{i,j}$  denote  
320 the  $(i, j)$ -element of the output of  $f_P$ . Then, for  $1 \leq i, j \leq 69$ ,  $[f_P(z_P)]_{i,j} := I_0 \exp(-r_{i,j})$ , where  
321  $r_{i,j}^2 := (X_j \cos \theta - Y_i \sin \theta)^2 / A^2 + (X_j \sin \theta + Y_i \cos \theta)^2 / B^2$ , and  $(X_j, Y_i)$  is the coordinate on  
322 the  $69 \times 69$  even grid on  $[-1, 1] \times [-1, 1]$ . We modify the output of  $f_P$  using a U-Net-like neural  
323 network  $f_A: \mathbb{R}^{69 \times 69} \times \mathbb{R}^{\dim z_A} \rightarrow \mathbb{R}^{69 \times 69 \times 3}$ . Thus, the decoding process is  $\mathcal{F} := f_A(f_P(z_P), z_A)$ .  
324 We set  $\dim z_A = 2$  for NN+phys+reg. We set  $h_A: \mathbb{R}^{69 \times 69} \rightarrow \mathbb{R}^{69 \times 69 \times 3}$  to be the repeat operator  
325 along the channel axis. The encoding process is as follows: first, features are extracted from an image  
326  $x$  by a convolutional net like [4]. The extracted features are flattened and fed to MLPs  $g_P$  and  $g_A$ .

327 **Results** Figure 4 shows an example of original data and random generation from the learned models.  
328 NN-only tends to generate non-realistic images, and NN+phys generates slightly better but still spuriously,  
329 whereas NN+phys+reg consistently generates galaxy-like images. More results (reconstruction,  
330 counterfactual generation, and inspection of latent variable) are deferred to Appendix F.

## 331 5.4 Human gait

332 **Dataset** We used a part of the dataset provided by [17], which contains measurements of locomotion  
333 at different speeds of 50 subjects. We extracted the angles of hip, knee, and ankle in the sagittal plane.  
334 Data originally comprise sequences of each stride normalized to be 100 steps, so each data-point  $x$  is  
335 a sequence  $x := [\theta_1 \ \dots \ \theta_{100}] \in \mathbb{R}^{3 \times 100}$ , where  $\theta_j := [\theta_{\text{hip},j} \ \theta_{\text{knee},j} \ \theta_{\text{ankle},j}]^T$ . We used different 400,  
336 100, and 344 sequences as training, validation, and test sets, respectively.

337 **Setting** Biomechanical modeling of gait is a long-standing problem [see, e.g., 31]. We did not  
338 choose a specific model but let  $f_P$  be a trainable Hamilton’s equation as in [37, 11].  $z_P \in \mathbb{R}^{2d_H}$   
339 works as the initial conditions of it, where  $d_H$  is the dimensionality of the generalized position.  
340 We let  $d_H = 3$  and modeled the neural Hamiltonian with an MLP. The solution of  $f_P = 0$  is  
341 transformed by  $f_A$  that also takes  $z_A \in \mathbb{R}^{15}$  as an argument. In summary, the decoding process  
342 is  $\mathcal{F} = f_A(\text{solve}[f_P = 0], z_A)$ . We set  $h_A$  to be an affine transform at each timestep, which has a  
343 weight matrix and a bias as  $\theta^r$ . The recognition networks  $g$  were modeled with MLPs.

344 **Results** Figure 5 is for visually comparing the difference of the learned models’ behavior due to  
345 the proposed regularizers. We compare the reconstructions by NN+phys and NN+phys+reg. The  
346 dashed lines show an intermediate of the decoding process, i.e.,  $\text{solve}[f_P = 0]$ , and the red solid  
347 lines show the final reconstruction, i.e.,  $f_A(\text{solve}[f_P = 0])$ . Without the regularization (upper row),  
348  $\text{solve}[f_P = 0]$  returns almost meaningless signals, and  $f_A$  bears the most effort of reconstruction. On  
349 the other hand, with the regularization (lower row),  $\text{solve}[f_P = 0]$  already matches well the data, and  
350  $f_A$  modifies it only slightly. Superiority of the regularized model was also confirmed quantitatively;  
351 the average test reconstruction errors were 0.273 with NN+phys and 0.259 with NN+phys+reg.

## 352 6 Conclusion

353 Physics-integrated VAEs by construction attain partial interpretability as some of the latent variables  
354 are semantically grounded to the physics models, and thus we can generate signals in a controlled  
355 manner. Moreover, they have extrapolation capability due to the physics models. In this work, we  
356 proposed a regularized learning objective for ensuring a proper functionality of the integrated physics  
357 models. We empirically validated the aforementioned unique capability of physics-integrated VAEs  
358 and the importance of the proposed regularization method.

## References

- 359
- 360 [1] A. Ajay, J. Wu, N. Fazeli, M. Bauza, L. P. Kaelbling, J. B. Tenenbaum, and A. Rodriguez.  
361 Augmenting physical simulators with stochastic neural networks: Case study of planar pushing  
362 and bouncing. In *Proceedings of the 2018 IEEE/RSJ International Conference on Intelligent*  
363 *Robots and Systems*, pages 3066–3073, 2018.
- 364 [2] A. Ajay, M. Bauza, J. Wu, N. Fazeli, J. B. Tenenbaum, A. Rodriguez, and L. P. Kaelbling.  
365 Combining physical simulators and object-based networks for control. In *Proceedings of the*  
366 *2019 IEEE International Conference on Robotics and Automation*, pages 3217–3223, 2019.
- 367 [3] B. Amos and J. Z. Kolter. OptNet: Differentiable optimization as a layer in neural networks. In  
368 *Proceedings of the 34th International Conference on Machine Learning*, pages 136–145, 2017.
- 369 [4] M. A. Aragon-Calvo and J. C. Carvajal. Self-supervised learning with physics-aware neural  
370 networks – I. Galaxy model fitting. *Monthly Notices of the Royal Astronomical Society*, 498(3):  
371 3713–3719, 2020.
- 372 [5] F. d. A. Belbute-Peres, T. D. Economon, and J. Z. Kolter. Combining differentiable PDE solvers  
373 and graph neural networks for fluid flow prediction. In *Proceedings of the 37th International*  
374 *Conference on Machine Learning*, pages 2402–2411, 2020.
- 375 [6] C. Chen, G. Zheng, H. Wei, and Z. Li. Physics-informed generative adversarial networks for  
376 sequence generation with limited data. NeurIPS Workshop on Interpretable Inductive Biases  
377 and Physically Structured Learning, 2020.
- 378 [7] T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud. Neural ordinary differential  
379 equations. In *Advances in Neural Information Processing Systems 31*, pages 6572–6583, 2018.
- 380 [8] E. de Bézenac, A. Pajot, and P. Gallinari. Deep learning for physical processes: Incorporating  
381 prior scientific knowledge. *Journal of Statistical Mechanics: Theory and Experiment*, 2019(12):  
382 124009, 2019.
- 383 [9] A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtar, and D. B. Rubin. *Bayesian Data*  
384 *Analysis*. Chapman and Hall/CRC, 3rd edition, 2013.
- 385 [10] T. Golany, D. Freedman, and K. Radinsky. SimGANs: Simulator-based generative adversarial  
386 networks for ECG synthesis to improve deep ECG classification. In *Proceedings of the 37th*  
387 *International Conference on Machine Learning*, pages 3597–3606, 2020.
- 388 [11] S. Greydanus, M. Dzamba, and J. Yosinski. Hamiltonian neural networks. In *Advances in*  
389 *Neural Information Processing Systems 32*, pages 15379–15389, 2019.
- 390 [12] X. Jia, J. Willard, A. Karpatne, J. Read, J. Zwart, M. Steinbach, and V. Kumar. Physics guided  
391 RNNs for modeling dynamical systems: A case study in simulating lake temperature profiles.  
392 In *Proceedings of the 2019 SIAM International Conference on Data Mining*, pages 558–566,  
393 2019.
- 394 [13] S. Kaltenbach and P.-S. Koutsourelakis. Incorporating physical constraints in a deep probabilistic  
395 machine learning framework for coarse-graining dynamical systems. *Journal of Computational*  
396 *Physics*, 419:109673, 2020.
- 397 [14] A. Karpatne, W. Watkins, J. Read, and V. Kumar. Physics-guided neural networks (PGNN): An  
398 application in lake temperature modeling. arXiv:1710.11431, 2017.
- 399 [15] D. P. Kingma and M. Welling. Auto-encoding variational Bayes. In *Proceedings of the 2nd*  
400 *International Conference on Learning Representations*, 2014.
- 401 [16] V. Le Guen and N. Thome. Disentangling physical dynamics from unknown factors for  
402 unsupervised video prediction. In *Proceedings of the 2020 IEEE/CVF Conference on Computer*  
403 *Vision and Pattern Recognition*, pages 11471–11481, 2020.
- 404 [17] T. Lencioni, I. Carpinella, M. Rabuffetti, A. Marzegan, and M. Ferrarin. Human kinematic,  
405 kinetic and EMG data during different walking and stair ascending and descending tasks.  
406 *Scientific Data*, 6(1):309, 2019.

- 407 [18] H. W. Leung and J. Bovy. Deep learning of multi-element abundances from high-resolution  
408 spectroscopic data. *Monthly Notices of the Royal Astronomical Society*, 483(3):3255–3277,  
409 2018.
- 410 [19] L. Li, S. Hoyer, R. Pederson, R. Sun, E. D. Cubuk, P. Riley, and K. Burke. Kohn-Sham equations  
411 as regularizer: Building prior knowledge into machine-learned physics. *Physical Review Letters*,  
412 126(3):036401, 2020.
- 413 [20] O. Linial, D. Eytan, and U. Shalit. Generative ODE modeling with known unknowns.  
414 arXiv:2003.10775, 2020.
- 415 [21] Y. Long and X. She. HybridNet: Integrating model-based and data-driven learning to predict  
416 evolution of dynamical systems. In *Proceedings of the 2nd Conference on Robot Learning*,  
417 pages 551–560, 2018.
- 418 [22] N. Muralidhar, J. Bu, Z. Cao, L. He, N. Ramakrishnan, D. Tafti, and A. Karpatne. PhyNet:  
419 Physics guided neural networks for particle drag force prediction in assembly. In *Proceedings*  
420 *of the 2020 SIAM International Conference on Data Mining*, pages 559–567, 2020.
- 421 [23] A. Nutkiewicz, Z. Yang, and R. K. Jain. Data-driven Urban Energy Simulation (DUE-S): A  
422 framework for integrating engineering simulation and machine learning methods in a multi-scale  
423 urban energy modeling workflow. *Applied Energy*, 225:1176–1189, 2018.
- 424 [24] D. C. Psychogios and L. H. Ungar. A hybrid neural network-first principles approach to process  
425 modeling. *AIChE Journal*, 38(10):1499–1511, 1992.
- 426 [25] C. Rackauckas, Y. Ma, J. Martensen, C. Warner, K. Zubov, R. Supekar, D. Skinner, A. Ra-  
427 madhan, and A. Edelman. Universal differential equations for scientific machine learning.  
428 arXiv:2001.04385, 2020.
- 429 [26] M. Raissi. Deep hidden physics models: Deep learning of nonlinear partial differential equations.  
430 *Journal of Machine Learning Research*, 19(25):1–24, 2018.
- 431 [27] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics-informed neural networks: A deep  
432 learning framework for solving forward and inverse problems involving nonlinear partial  
433 differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- 434 [28] D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate  
435 inference in deep generative models. In *Proceedings of the 31st International Conference on*  
436 *Machine Learning*, pages 1278–1286, 2014.
- 437 [29] R. Rico-Martínez, J. S. Anderson, and I. G. Kevrekidis. Continuous-time nonlinear signal  
438 processing: A neural network based approach for gray box identification. In *Proceedings of the*  
439 *IEEE Workshop on Neural Networks for Signal Processing*, pages 596–605, 1994.
- 440 [30] M. Rixner and P.-S. Koutsourelakis. A probabilistic generative model for semi-supervised train-  
441 ing of coarse-grained surrogates and enforcing physical constraints through virtual observables.  
442 arXiv:2006.01789, 2020.
- 443 [31] D. G. E. Robertson, G. E. Caldwell, J. Hamill, G. Kamen, and S. N. Whittlesey. *Research*  
444 *Methods in Biomechanics*. Human Kinetics, 2nd edition, 2014.
- 445 [32] M. A. Roehrl, T. A. Runkler, V. Brandtstetter, M. Tokic, and S. Obermayer. Modeling  
446 system dynamics with physics-informed neural networks based on Lagrangian mechanics.  
447 arXiv:2005.14617, 2020.
- 448 [33] U. Sengupta, M. Amos, J. S. Hosking, C. E. Rasmussen, M. Juniper, and P. J. Young. Ensembling  
449 geophysical models with Bayesian neural networks. In *Advances in Neural Information*  
450 *Processing Systems 33*, 2020.
- 451 [34] R. Stewart and S. Ermon. Label-free supervision of neural networks with physics and domain  
452 knowledge. In *Proceedings of the 31st AAAI Conference on Artificial Intelligence*, pages  
453 2576–2582, 2017.

- 454 [35] P. Stinis, T. Hagge, A. M. Tartakovsky, and E. Yeung. Enforcing constraints for interpolation  
455 and extrapolation in generative adversarial networks. *Journal of Computational Physics*, 397:  
456 108844, 2019.
- 457 [36] M. L. Thompson and M. A. Kramer. Modeling chemical processes using prior knowledge and  
458 neural networks. *AIChE Journal*, 40(8):1328–1340, 1994.
- 459 [37] P. Toth, D. J. Rezende, A. Jaegle, S. Racanière, A. Botev, and I. Higgins. Hamiltonian generative  
460 networks. In *Proceedings of the 8th International Conference on Learning Representations*,  
461 2020.
- 462 [38] L. von Rueden, S. Mayer, K. Beckh, B. Georgiev, S. Giesselbach, R. Heese, B. Kirsch, J. Pfrom-  
463 mer, A. Pick, R. Ramamurthy, M. Walczak, J. Garcke, C. Bauckhage, and J. Schuecker. Informed  
464 machine learning – A taxonomy and survey of integrating knowledge into learning systems.  
465 arXiv:1903.12394v2, 2020.
- 466 [39] Z. Y. Wan, P. Vlachas, P. Koumoutsakos, and T. Sapsis. Data-assisted reduced-order modeling  
467 of extreme events in complex dynamical systems. *PLOS ONE*, 13(5):e0197704, 2018.
- 468 [40] Q. Wang, F. Li, Y. Tang, and Y. Xu. Integrating model-driven and data-driven methods for  
469 power system frequency stability assessment and control. *IEEE Transactions on Power Systems*,  
470 34(6):4557–4568, 2019.
- 471 [41] J. Willard, X. Jia, S. Xu, M. Steinbach, and V. Kumar. Integrating physics-based modeling with  
472 machine learning: A survey. arXiv:2003.04919, 2020.
- 473 [42] Z. Yang, J.-L. Wu, and H. Xiao. Enforcing deterministic constraints on generative adversarial  
474 networks for emulating physical systems. arXiv:1911.06671, 2019.
- 475 [43] Ç. Yıldız, M. Heinonen, and H. Lähdesmäki. ODE2VAE: Deep generative second order ODEs  
476 with Bayesian neural networks. In *Advances in Neural Information Processing Systems 32*,  
477 pages 13412–13421, 2019.
- 478 [44] Y. Yin, V. Le Guen, J. Dona, I. Ayed, E. de Bézenac, N. Thome, and P. Gallinari. Augmenting  
479 physical models with deep networks for complex dynamics forecasting. In *Proceedings of the*  
480 *9th International Conference on Learning Representations*, 2021.
- 481 [45] C.-C. Young, W.-C. Liu, and M.-C. Wu. A physically based and machine learning hybrid  
482 approach for accurate rainfall-runoff modeling during extreme typhoon events. *Applied Soft*  
483 *Computing*, 53:205–216, 2017.
- 484 [46] A. Zeng, S. Song, J. Lee, A. Rodriguez, and T. Funkhouser. TossingBot: Learning to throw  
485 arbitrary objects with residual physics. In *Proceedings of Robotics: Science and Systems*, 2019.
- 486 [47] J. Zhang, C. Wei, and C. Wu. Thermodynamic consistent neural networks for learning material  
487 interfacial mechanics. arXiv:2011.14172, 2020.

## 488 Checklist

- 489 1. For all authors...
- 490 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
491 contributions and scope? [Yes]
- 492 (b) Did you describe the limitations of your work? [Yes] See Section 2; e.g., integration of  
493 overly-complex physics models is an open challenge.
- 494 (c) Did you discuss any potential negative societal impacts of your work? [N/A] This  
495 paper is on general methodology, and we do not think we can discuss concrete social  
496 impacts at this layer of research.
- 497 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
498 them? [Yes]
- 499 2. If you are including theoretical results...

- 500 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Section 3  
501 and Appendix B.
- 502 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix B.
- 503 3. If you ran experiments...
- 504 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
505 mental results (either in the supplemental material or as a URL)? [Yes]
- 506 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
507 were chosen)? [Yes] See Section 5 and Appendix E.
- 508 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
509 ments multiple times)? [Yes] See Section 5 and Appendix F.
- 510 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
511 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 5 and Appendix F.
- 512 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 513 (a) If your work uses existing assets, did you cite the creators? [Yes] See Section 5 and  
514 Appendix E.
- 515 (b) Did you mention the license of the assets? [No] License of the two existing assets we  
516 used is already manifested at the original sources.
- 517 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]  
518
- 519 (d) Did you discuss whether and how consent was obtained from people whose data you're  
520 using/curating? [N/A]
- 521 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
522 information or offensive content? [N/A]
- 523 5. If you used crowdsourcing or conducted research with human subjects...
- 524 (a) Did you include the full text of instructions given to participants and screenshots, if  
525 applicable? [N/A]
- 526 (b) Did you describe any potential participant risks, with links to Institutional Review  
527 Board (IRB) approvals, if applicable? [N/A]
- 528 (c) Did you include the estimated hourly wage paid to participants and the total amount  
529 spent on participant compensation? [N/A]