

---

# Heavy Ball Neural Ordinary Differential Equations

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We propose heavy ball neural ordinary differential equations (HBNODEs), lever-  
2 aging the continuous limit of the classical momentum accelerated gradient descent,  
3 to improve neural ODEs (NODEs) training and inference. HBNODEs have two  
4 properties that imply practical advantages over NODEs: (i) The adjoint state of an  
5 HBNODE also satisfies an HBNODE, accelerating both forward and backward  
6 ODE solvers, thus significantly reducing the number of function evaluations (NFEs)  
7 and improving the utility of the trained models. (ii) The spectrum of HBNODEs  
8 is well structured, enabling effective learning of long-term dependencies from  
9 complex sequential data. We verify the advantages of HBNODEs over NODEs on  
10 benchmark tasks, including image classification, learning complex dynamics, and  
11 sequential modeling. Our method requires remarkably fewer forward and backward  
12 NFEs, is more accurate, and learns long-term dependencies more effectively than  
13 the other ODE-based neural network models.

## 1 Introduction

15 Neural ordinary differential equations (NODEs) are a family of continuous-depth machine learn-  
16 ing (ML) models whose forward and backward propagations rely on solving an ODE and its ad-  
17 joint equation [4]. NODEs model the dynamics of hidden features  $\mathbf{h}(t) \in \mathbb{R}^N$  using an ODE,  
18 which is parametrized by a neural network  $f(\mathbf{h}(t), t, \theta) \in \mathbb{R}^N$  with learnable parameters  $\theta$ , i.e.,  
19

$$\frac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}(t), t, \theta). \quad (1)$$

20 Starting from the input  $\mathbf{h}(t_0)$ , NODEs obtain the out-  
21 put  $\mathbf{h}(T)$  by solving (1) for  $t_0 \leq t \leq T$  with the ini-  
22 tial value  $\mathbf{h}(t_0)$ , using a black-box numerical ODE  
23 solver. The number of function evaluations (NFEs)  
24 that the black-box ODE solver requires in a single  
25 forward pass is an analogue for the continuous-depth  
26 models [4] to the depth of networks in ResNets [17].  
27 The loss between NODE prediction  $\mathbf{h}(T)$  and the  
28 ground truth is denoted by  $\mathcal{L}(\mathbf{h}(T))$ ; we update pa-  
29 rameters  $\theta$  using the following gradient [41]

$$\frac{d\mathcal{L}(\mathbf{h}(T))}{d\theta} = \int_{t_0}^T \mathbf{a}(t) \frac{\partial f(\mathbf{h}(t), t, \theta)}{\partial \theta} dt, \quad (2)$$

30 where  $\mathbf{a}(t) := \partial \mathcal{L} / \partial \mathbf{h}(t)$  is the adjoint state, which  
31 satisfies the adjoint equation

$$\frac{d\mathbf{a}(t)}{dt} = -\mathbf{a}(t) \frac{\partial f(\mathbf{h}(t), t, \theta)}{\partial \mathbf{h}}. \quad (3)$$

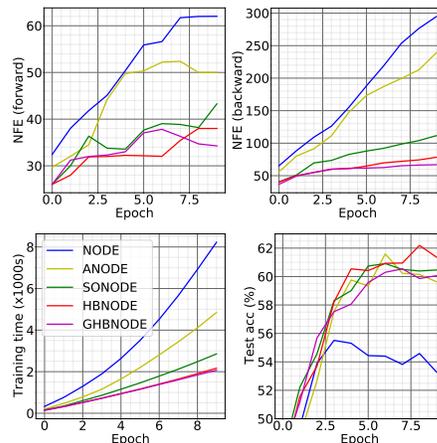


Figure 1: Contrasting NODE, ANODE, SONODE, HBNODE, and GHBNODE for CIFAR10 classification in NFEs, training time, and test accuracy. (Tolerance:  $10^{-5}$ , see Sec. 5.2 for experimental details.)

32 NODEs are flexible in learning from irregularly-sampled sequential data and particularly suitable  
 33 for learning complex dynamical systems [4, 44, 58, 37, 9, 24], which can be trained by efficient  
 34 algorithms [42, 7, 59]. NODE-based continuous generative models have computational advantages  
 35 over the classical normalizing flows [4, 15, 57, 13]. NODEs have also been generalized to neural  
 36 stochastic differential equations, stochastic processes, and graph NODEs [22, 29, 39, 52, 21, 36].  
 37 The drawback of NODEs is also prominent. In many ML tasks, NODEs require very high NFEs in  
 38 both training and inference, especially in high accuracy settings where a lower tolerance is needed.  
 39 The NFEs increase rapidly with training; high NFEs reduce computational speed and accuracy of  
 40 NODEs and can lead to blow-ups in the worst-case scenario [15, 10, 31, 37]. As an illustration, we  
 41 train NODEs for CIFAR10 classification using the same model and experimental settings as in [10],  
 42 except using a tolerance of  $10^{-5}$ ; Fig. 1 shows both forward and backward NFEs and the training  
 43 time of different ODE-based models; we see that NFEs and computational times increase very rapidly  
 44 for NODE, ANODE [10], and SONODE [37]. More results on the large NFE and degrading utility  
 45 issues for different benchmark experiments are available in Sec. 5. Another issue is that NODEs  
 46 often fail to effectively learn long-term dependencies in sequential data [27], discussed in Sec. 4.

## 47 1.1 Contribution

48 We propose heavy ball neural ODEs (HBNODEs), leveraging the continuous limit of the classical  
 49 momentum accelerated gradient descent, to improve NODE training and inference. At the core of  
 50 HBNODE is replacing the first-order ODE (1) with a heavy ball ODE (HBODE), i.e., a second-order  
 51 ODE with an appropriate damping term. HBNODEs have two theoretical properties that imply  
 52 practical advantages over NODEs:

- 53 • The adjoint equation used for training a HBNODE is also a HBNODE (see Prop. 1 and Prop. 2),  
 54 accelerating both forward and backward propagation, thus significantly reducing both forward and  
 55 backward NFEs. The reduction in NFE using HBNODE over existing benchmark ODE-based  
 56 models becomes more aggressive as the error tolerance of the ODE solvers decreases.
- 57 • The spectrum of the HBODE is well-structured (see Prop. 4), alleviating the vanishing gradient  
 58 issue in back-propagation and enabling the model to effectively learn long-term dependencies from  
 59 sequential data.

60 To mitigate the potential blow-up problem in training HBNODEs, we further propose generalized  
 61 HBNODEs (GHBNODEs) by integrating skip connections [18] and gating mechanisms [20] into the  
 62 HBNODE. See Sec. 3 for details.

## 63 1.2 Organization

64 We organize the paper as follows: In Secs 2 and 3, we present our motivation, algorithm, and analysis  
 65 of HBNODEs and GHBNODEs, respectively. We analyze the spectrum structure of the adjoint  
 66 equation of HBNODEs/GHBNODEs in Sec. 4, which indicates that HBNODEs/GHBNODEs can  
 67 learn long-term dependency effectively. We test the performance of HBNODEs and GHBNODEs on  
 68 benchmark point cloud separation, image classification, learning dynamics, and sequential modeling  
 69 in Sec. 5. We discuss more related work in Sec. 6, followed by concluding remarks. Technical proofs  
 70 and more experimental details are provided in the appendix.

# 71 2 Heavy Ball Neural Ordinary Differential Equations

## 72 2.1 Heavy ball ordinary differential equation

73 Classical momentum method, a.k.a., the heavy ball method, has achieved remarkable success in  
 74 accelerating gradient descent [40] and has significantly improved the training of deep neural networks  
 75 [49]. As the continuous limit of the classical momentum method, heavy ball ODE (HBODE) has  
 76 been studied in various settings and has been used to analyze the acceleration phenomenon of the  
 77 momentum methods [26, 46]. For the ease of reading and completeness, we derive the HBODE  
 78 from the classical momentum method. Starting from initial points  $\mathbf{x}^0$  and  $\mathbf{x}^1$ , gradient descent with  
 79 classical momentum searches a minimum of the function  $F(\mathbf{x})$  through the following iteration

$$\mathbf{x}^{k+1} = \mathbf{x}^k - s\nabla F(\mathbf{x}^k) + \beta(\mathbf{x}^k - \mathbf{x}^{k-1}), \quad (4)$$

80 where  $s > 0$  is the step size and  $0 \leq \beta < 1$  is the momentum hyperparameter. For any fixed step  
 81 size  $s$ , let  $\mathbf{m}^k := (\mathbf{x}^{k+1} - \mathbf{x}^k)/\sqrt{s}$ , and let  $\beta := 1 - \gamma\sqrt{s}$ , where  $\gamma \geq 0$  is another hyperparameter.  
 82 Then we can rewrite (4) as

$$\mathbf{m}^{k+1} = (1 - \gamma\sqrt{s})\mathbf{m}^k - \sqrt{s}\nabla F(\mathbf{x}^k); \quad \mathbf{x}^{k+1} = \mathbf{x}^k + \sqrt{s}\mathbf{m}^{k+1}. \quad (5)$$

83 Let  $s \rightarrow 0$  in (5); we obtain the following system of first-order ODEs,

$$\frac{d\mathbf{x}(t)}{dt} = \mathbf{m}(t); \quad \frac{d\mathbf{m}(t)}{dt} = -\gamma\mathbf{m}(t) - \nabla F(\mathbf{x}(t)). \quad (6)$$

84 This can be further rewritten as a second-order heavy ball ODE (HBODE), which also models a  
85 damped oscillator,

$$\frac{d^2\mathbf{x}(t)}{dt^2} + \gamma\frac{d\mathbf{x}(t)}{dt} = -\nabla F(\mathbf{x}(t)). \quad (7)$$

86 We compare the dynamics of HBODE (7) and the  
87 following ODE limit of the gradient descent (GD)

$$\frac{d\mathbf{x}}{dt} = -\nabla F(\mathbf{x}). \quad (8)$$

88 In particular, we solve the ODEs (7) and (8) with  
89  $F(\mathbf{x})$  defined as a Rosenbrock [43] or Beale [14]  
90 function (see Appendix E.6 for experimental de-  
91 tails). Fig. 2 shows that with the same numerical  
92 ODE solver, HBODE converges to the stationary  
93 point (marked by stars) faster than (8). The fact that  
94 HBODE can accelerate the dynamics of the ODE

95 for a gradient system motivates us to propose HBNODE to accelerate forward propagation of NODE.

## 96 2.2 Heavy ball neural ordinary differential equations

97 Similar to NODE, we parameterize  $-\nabla F$  in (7) using a neural network  $f(\mathbf{h}(t), t, \theta)$ , resulting in the  
98 following HBNODE with initial position  $\mathbf{h}(t_0)$  and momentum  $\mathbf{m}(t_0) := d\mathbf{h}/dt(t_0)$ ,

$$\frac{d^2\mathbf{h}(t)}{dt^2} + \gamma\frac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}(t), t, \theta), \quad (9)$$

99 where  $\gamma \geq 0$  is the damping parameter, which can be set as a tunable or a learnable hyperparameter  
100 with positivity constraint. In the trainable case, we use  $\gamma = \epsilon \cdot \text{sigmoid}(\omega)$  for a trainable  $\omega \in \mathbb{R}$  and  
101 a fixed tunable upper bound  $\epsilon$  (we set  $\epsilon = 1$  below). According to (6), HBNODE (9) is equivalent to

$$\frac{d\mathbf{h}(t)}{dt} = \mathbf{m}(t); \quad \frac{d\mathbf{m}(t)}{dt} = -\gamma\mathbf{m}(t) + f(\mathbf{h}(t), t, \theta). \quad (10)$$

102 Equation (9) (or equivalently, the system (10)) defines the forward ODE for the HBNODE, and we  
103 can use either the first-order (Prop. 2) or the second-order (Prop. 1) adjoint sensitivity method to  
104 update the parameter  $\theta$  [37].

105 **Proposition 1** (Adjoint equation for HBNODE). *The adjoint state  $\mathbf{a}(t) := \partial\mathcal{L}/\partial\mathbf{h}(t)$  for the  
106 HBNODE (9) satisfies the following HBODE with the same damping parameter  $\gamma$  as that in (9),*

$$\frac{d^2\mathbf{a}(t)}{dt^2} - \gamma\frac{d\mathbf{a}(t)}{dt} = \mathbf{a}(t)\frac{\partial f}{\partial\mathbf{h}}(\mathbf{h}(t), t, \theta). \quad (11)$$

107 **Remark 1.** *Note that we solve the adjoint equation (11) from time  $t = T$  to  $t = t_0$  in the backward  
108 propagation. By letting  $\tau = T - t$  and  $\mathbf{b}(\tau) = \mathbf{a}(T - \tau)$ , we can rewrite (11) as follows,*

$$\frac{d^2\mathbf{b}(\tau)}{d\tau^2} + \gamma\frac{d\mathbf{b}(\tau)}{d\tau} = \mathbf{b}(\tau)\frac{\partial f}{\partial\mathbf{h}}(\mathbf{h}(T - \tau), T - \tau, \theta). \quad (12)$$

109 *Therefore, the adjoint of the HBNODE is also a HBNODE and they have the same damping parameter.*

110 We can also employ (10) and its adjoint for the forward and backward propagations, respectively.

111 **Proposition 2** (Adjoint equations for the first-order HBNODE system). *The adjoint states  $\mathbf{a}_h(t)$   
112  $:= \partial\mathcal{L}/\partial\mathbf{h}(t)$  and  $\mathbf{a}_m(t) := \partial\mathcal{L}/\partial\mathbf{m}(t)$  for the first-order HBNODE system (10) satisfy*

$$\frac{d\mathbf{a}_h(t)}{dt} = -\mathbf{a}_m(t)\frac{\partial f}{\partial\mathbf{h}}(\mathbf{h}(t), t, \theta); \quad \frac{d\mathbf{a}_m(t)}{dt} = -\mathbf{a}_h(t) + \gamma\mathbf{a}_m(t). \quad (13)$$

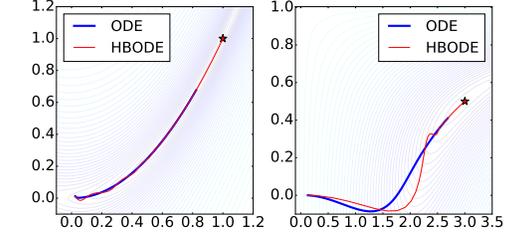


Figure 2: Comparing the trajectory of ODE and HBODE when  $F(\mathbf{x})$  is the Rosenbrock (left) and Beale (right) functions.

113 **Remark 2.** Let  $\tilde{\mathbf{a}}_m(t) = d\mathbf{a}_m(t)/dt$ , then  $\mathbf{a}_m(t)$  and  $\tilde{\mathbf{a}}_m(t)$  satisfies the following first-order  
 114 heavy ball ODE system

$$\frac{d\mathbf{a}_m(t)}{dt} = \tilde{\mathbf{a}}_m(t); \quad \frac{d\tilde{\mathbf{a}}_m(t)}{dt} = \mathbf{a}_m(t) \frac{\partial f}{\partial \mathbf{h}}(\mathbf{h}(t), t, \theta) + \gamma \tilde{\mathbf{a}}_m(t). \quad (14)$$

115 Note that we solve this system backward in time in back-propagation. Moreover, we have  $\mathbf{a}_h(t) =$   
 116  $\gamma \mathbf{a}_m(t) - \tilde{\mathbf{a}}_m(t)$ .

117 Similar to [37], we use the coupled first-order HBNODE system (10) and its adjoint first-order  
 118 HBNODE system (13) for practical implementation, since the entangled representation permits faster  
 119 computation [37] of the gradients of the coupled ODE systems.

### 120 3 Generalized Heavy Ball Neural Ordinary Differential Equations

121 In this section, we propose a generalized version of HBNODE (GHBNODE), see (15), to mitigate  
 122 the potential blow-up issue in training ODE-based models. In our experiments, we observe that  
 123  $\mathbf{h}(t)$  of ANODEs [10], SONODEs [37], and HBNODEs (10) usually grows much faster than that of  
 124 NODEs. The fast growth of  $\mathbf{h}(t)$  can lead to finite-time blow up. As an illustration, we compare the  
 125 performance of NODE, ANODE, SONODE, HBNODE, and GHBNODE on the Silverbox task as  
 126 in [37]. The goal of the task is to learn the voltage of an electronic circuit that resembles a Duffing  
 127 oscillator, where the input voltage  $V_1(t)$  is used to predict the output  $V_2(t)$ . Similar to the setting  
 128 in [37], we first augment ANODE by 1 dimension with 0-augmentation and augment SONODE,  
 129 HBNODE, and GHBNODE with a dense network. We use a simple dense layer to parameterize  $f$   
 130 for all five models, with an extra input term for  $V_1(t)$ <sup>1</sup>. For both HBNODE and GHBNODE, we  
 131 set the damping parameter  $\gamma$  to be sigmoid(-3). For GHBNODE (15) below, we set  $\sigma(\cdot)$  to be the  
 132 hardtanh function with bound  $[-5, 5]$  and  $\xi = \ln(2)$ . The detailed architecture can be found in  
 133 Appendix E. As shown in Fig. 3, compared to the vanilla NODE, the  $\ell_2$  norm of  $\mathbf{h}(t)$  grows much  
 134 faster when a higher order NODE is used, which leads to blow-up during training. Similar issues arise  
 135 in the time series experiments (see Sec. 5.4), where SONODE blows up during long term integration  
 136 in time, and HBNODE suffers from the same issue with some initialization.

137 To alleviate the problem above, we propose the fol-  
 138 lowing generalized HBNODE

$$\begin{aligned} \frac{d\mathbf{h}(t)}{dt} &= \sigma(\mathbf{m}(t)), \\ \frac{d\mathbf{m}(t)}{dt} &= -\gamma \mathbf{m}(t) + f(\mathbf{h}(t), t, \theta) - \xi \mathbf{h}(t), \end{aligned} \quad (15)$$

139 where  $\sigma(\cdot)$  is a nonlinear activation, which is set as  
 140 tanh in our experiments. The positive hyperparam-  
 141 eters  $\gamma, \xi > 0$  are tunable or learnable. In the trainable  
 142 case, we let  $\gamma = \epsilon \cdot \text{sigmoid}(\omega)$  as in HBNODE, and  
 143  $\xi = \text{softplus}(\chi)$  to ensure that  $\gamma, \xi \geq 0$ . Here, we  
 144 integrate two main ideas into the design of GHBN-  
 145 ODE: (i) We incorporate the gating mechanism used  
 146 in LSTM [20] and GRU [6], which can suppress the aggregation of  $\mathbf{m}(t)$ ; (ii) Following the idea  
 147 of skip connection [18], we add the term  $\xi \mathbf{h}(t)$  into the governing equation of  $\mathbf{m}(t)$ , which benefits  
 148 training and generalization of GHBNODEs. Fig. 3 shows that GHBNODE can indeed control the  
 149 growth of  $\mathbf{h}(t)$  effectively.

150 **Proposition 3** (Adjoint equations for GHBNODEs). *The adjoint states  $\mathbf{a}_h(t) := \partial \mathcal{L} / \partial \mathbf{h}(t)$ ,*  
 151  *$\mathbf{a}_m(t) := \partial \mathcal{L} / \partial \mathbf{m}(t)$  for the GHBNODE (15) satisfy the following first-order ODE system*

$$\frac{\partial \mathbf{a}_h(t)}{\partial t} = -\mathbf{a}_m(t) \left( \frac{\partial f}{\partial \mathbf{h}}(\mathbf{h}(t), t, \theta) - \xi \mathbf{I} \right), \quad \frac{\partial \mathbf{a}_m(t)}{\partial t} = -\mathbf{a}_h(t) \sigma'(\mathbf{m}(t)) + \gamma \mathbf{a}_m(t). \quad (16)$$

152 Though the adjoint state of the GHBNODE (16) does not satisfy the exact heavy ball ODE, based on  
 153 our empirical study, it also significantly reduces the backward NFES.

<sup>1</sup>Here, we exclude an  $\mathbf{h}^3$  term that appeared in the original Duffing oscillator model because including it would result in finite-time explosion.

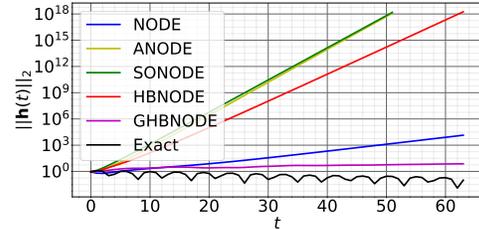


Figure 3: Contrasting  $\mathbf{h}(t)$  for different models.  $\mathbf{h}(t)$  in ANODE, SONODE, and HBNODE grows much faster than that in NODE. GHBNODE controls the growth of  $\mathbf{h}(t)$  effectively when  $t$  is large.

#### 154 4 Learning long-term dependencies – Vanishing gradient

155 It is known that the vanishing and exploding gradients are two bottlenecks for training recurrent  
 156 neural networks (RNNs) with long-term dependencies [2, 38] (see Appendix C for a brief review on  
 157 the exploding and vanishing gradient issues in training RNNs). The exploding gradients issue can be  
 158 effectively resolved via gradient clipping, training loss regularization, etc [38]. Thus in practice the  
 159 vanishing gradient is the major issue for learning long-term dependencies [38]. As the continuous  
 160 analogue of RNN, NODEs as well as their hybrid ODE-RNN models, may also suffer from vanishing  
 161 in the adjoint state  $\mathbf{a}(t) := \partial\mathcal{L}/\partial\mathbf{h}(t)$  [27]. When the vanishing gradient issue happens,  $\mathbf{a}(t)$  goes  
 162 to  $\mathbf{0}$  quickly as  $T - t$  increases, then  $d\mathcal{L}/d\theta$  in (2) will be independent of these  $\mathbf{a}(t)$ . We have the  
 163 following expressions for the adjoint states of the NODE and HBNODE (see Appendix C for detailed  
 164 derivation):

165 • For NODE, we have

$$\frac{\partial\mathcal{L}}{\partial\mathbf{h}_t} = \frac{\partial\mathcal{L}}{\partial\mathbf{h}_T} \frac{\partial\mathbf{h}_T}{\partial\mathbf{h}_t} = \frac{\partial\mathcal{L}}{\partial\mathbf{h}_T} \exp\left\{-\int_T^t \frac{\partial f}{\partial\mathbf{h}}(\mathbf{h}(s), s, \theta) ds\right\}. \quad (17)$$

166 • For GHBNODE<sup>2</sup>, from (13) we can derive

$$\begin{bmatrix} \frac{\partial\mathcal{L}}{\partial\mathbf{h}_t} & \frac{\partial\mathcal{L}}{\partial\mathbf{m}_t} \end{bmatrix} = \begin{bmatrix} \frac{\partial\mathcal{L}}{\partial\mathbf{h}_T} & \frac{\partial\mathcal{L}}{\partial\mathbf{m}_T} \end{bmatrix} \begin{bmatrix} \frac{\partial\mathbf{h}_T}{\partial\mathbf{h}_t} & \frac{\partial\mathbf{h}_T}{\partial\mathbf{m}_t} \\ \frac{\partial\mathbf{m}_T}{\partial\mathbf{h}_t} & \frac{\partial\mathbf{m}_T}{\partial\mathbf{m}_t} \end{bmatrix} = \begin{bmatrix} \frac{\partial\mathcal{L}}{\partial\mathbf{h}_T} & \frac{\partial\mathcal{L}}{\partial\mathbf{m}_T} \end{bmatrix} \exp\left\{-\int_T^t \underbrace{\begin{bmatrix} \mathbf{0} & \frac{\partial\sigma}{\partial\mathbf{m}} \\ \frac{\partial f}{\partial\mathbf{h}} - \xi\mathbf{I} & -\gamma\mathbf{I} \end{bmatrix}}_{:=M} ds\right\}. \quad (18)$$

167 Note that the matrix exponential is directly related to its eigenvalues. By Schur decomposition, there  
 168 exists an orthogonal matrix  $\mathbf{Q}$  and an upper triangular matrix  $\mathbf{U}$ , where the diagonal entries of  $\mathbf{U}$  are  
 169 eigenvalues of  $\mathbf{Q}$  ordered by their real parts, such that

$$-M = \mathbf{Q}\mathbf{U}\mathbf{Q}^\top \implies \exp\{-M\} = \mathbf{Q}\exp\{\mathbf{U}\}\mathbf{Q}^\top. \quad (19)$$

170 Let  $\mathbf{v}^\top := \begin{bmatrix} \frac{\partial\mathcal{L}}{\partial\mathbf{h}_T} & \frac{\partial\mathcal{L}}{\partial\mathbf{m}_T} \end{bmatrix} \mathbf{Q}$ , then (18) can be rewritten as

$$\begin{bmatrix} \frac{\partial\mathcal{L}}{\partial\mathbf{h}_t} & \frac{\partial\mathcal{L}}{\partial\mathbf{m}_t} \end{bmatrix} = \begin{bmatrix} \frac{\partial\mathcal{L}}{\partial\mathbf{h}_T} & \frac{\partial\mathcal{L}}{\partial\mathbf{m}_T} \end{bmatrix} \exp\{-M\} = \begin{bmatrix} \frac{\partial\mathcal{L}}{\partial\mathbf{h}_T} & \frac{\partial\mathcal{L}}{\partial\mathbf{m}_T} \end{bmatrix} \mathbf{Q}\exp\{\mathbf{U}\}\mathbf{Q}^\top = \mathbf{v}^\top \exp\{\mathbf{U}\}\mathbf{Q}^\top. \quad (20)$$

171 By taking the  $\ell_2$  norm in (20) and dividing both sides by  $\|\begin{bmatrix} \frac{\partial\mathcal{L}}{\partial\mathbf{h}_T} & \frac{\partial\mathcal{L}}{\partial\mathbf{m}_T} \end{bmatrix}\|_2$ , we arrive at

$$\frac{\|\begin{bmatrix} \frac{\partial\mathcal{L}}{\partial\mathbf{h}_t} & \frac{\partial\mathcal{L}}{\partial\mathbf{m}_t} \end{bmatrix}\|_2}{\|\begin{bmatrix} \frac{\partial\mathcal{L}}{\partial\mathbf{h}_T} & \frac{\partial\mathcal{L}}{\partial\mathbf{m}_T} \end{bmatrix}\|_2} = \frac{\|\mathbf{v}^\top \exp\{\mathbf{U}\}\mathbf{Q}^\top\|_2}{\|\mathbf{v}^\top \mathbf{Q}^\top\|_2} = \frac{\|\mathbf{v}^\top \exp\{\mathbf{U}\}\|_2}{\|\mathbf{v}\|_2} = \|\mathbf{e}^\top \exp\{\mathbf{U}\}\|_2, \quad (21)$$

172 i.e.,  $\|\begin{bmatrix} \frac{\partial\mathcal{L}}{\partial\mathbf{h}_t} & \frac{\partial\mathcal{L}}{\partial\mathbf{m}_t} \end{bmatrix}\|_2 = \|\mathbf{e}^\top \exp\{\mathbf{U}\}\|_2 \|\begin{bmatrix} \frac{\partial\mathcal{L}}{\partial\mathbf{h}_T} & \frac{\partial\mathcal{L}}{\partial\mathbf{m}_T} \end{bmatrix}\|_2$  where  $\mathbf{e} = \mathbf{v}/\|\mathbf{v}\|_2$ .

173 **Proposition 4.** *The eigenvalues of  $-M$  can be paired so that the sum of each pair equals  $(t - T)\gamma$ .*

174 For a given constant  $a > 0$ , we can group the upper triangular matrix  $\exp\{\mathbf{U}\}$  as follows

$$\exp\{\mathbf{U}\} := \begin{bmatrix} \exp\{\mathbf{U}_L\} & \mathbf{P} \\ \mathbf{0} & \exp\{\mathbf{U}_V\} \end{bmatrix}, \quad (22)$$

175 where the diagonal of  $\mathbf{U}_L$  ( $\mathbf{U}_V$ ) contains eigenvalues of  $-M$  that are no less (greater) than  $(t - T)a$ .  
 176 Then, we have  $\|\mathbf{e}^\top \exp\{\mathbf{U}\}\|_2 \geq \|\mathbf{e}_L^\top \exp\{\mathbf{U}_L\}\|_2$  where the vector  $\mathbf{e}_L$  denotes the first  $m$  columns  
 177 of  $\mathbf{e}$  with  $m$  be the number of columns of  $\mathbf{U}_L$ . By choosing  $0 \leq \gamma \leq 2a$ , for every pair of eigenvalues  
 178 of  $-M$  there is at least one eigenvalue whose real part is no less than  $(t - T)a$ . Therefore,  $\exp\{\mathbf{U}_L\}$   
 179 decays at a rate at most  $(t - T)a$ , and the dimension of  $\mathbf{U}_L$  is at least  $N \times N$ . We avoid exploding  
 180 gradients by clipping the  $\ell_2$  norm of the adjoint states similar to that used for training RNNs.

181 In contrast, all eigenvalues of the matrix  $\int_T^t \partial f / \partial \mathbf{h} ds$  in (17) for NODE can be very positive  
 182 or negative, resulting in exploding or vanishing gradients. As an illustration, we consider the  
 183 benchmark Walker2D kinematic simulation task that requires learning long-term dependencies  
 184 effectively [27, 3]. We train ODE-RNN [44] and (G)HBNODE-RNN on this benchmark dataset, and  
 185 the detailed experimental settings are provided in Sec. 5.4. Figure 4 plots  $\|\partial\mathcal{L}/\partial\mathbf{h}_t\|_2$  for ODE-RNN  
 186 and  $\|\partial\mathcal{L}/\partial\mathbf{h}_t \partial\mathcal{L}/\partial\mathbf{m}_t\|_2$  for (G)HBNODE-RNN, showing that the adjoint state of ODE-RNN  
 187 vanishes quickly, while that of (G)HBNODE-RNN does not vanish even when the gap between  $T$   
 188 and  $t$  is very large.

<sup>2</sup>HBNODE can be seen as a special GHBNODE with  $\xi = 0$  and  $\sigma$  be the identity map.

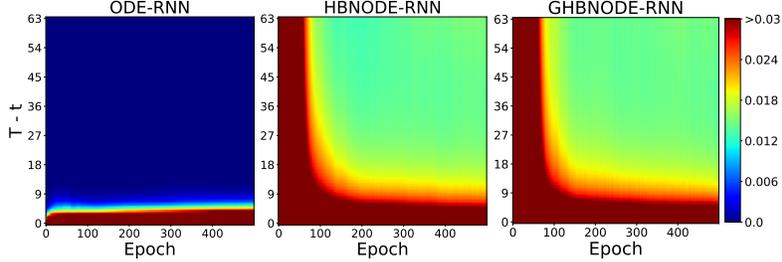


Figure 4: Plot of the  $\ell_2$ -norm of the adjoint states for ODE-RNN and (G)HBNODE-RNN back-propagated from the last time stamp. The adjoint state of ODE-RNN vanishes quickly when the gap between the final time  $T$  and intermediate time  $t$  becomes larger, while the adjoint states of (G)HBNODE-RNN decays much more slowly. This implies that (G)HBNODE-RNN is more effective in learning long-term dependency than ODE-RNN.

## 189 5 Experimental Results

190 In this section, we compare the performance of the proposed HBNODE and GHBNODE with existing  
 191 ODE-based models, including NODE [4], ANODE [10], and SONODE [37] on the benchmark point  
 192 cloud separation, image classification, learning dynamical systems, and kinematic simulation. For all  
 193 the experiments, we use Adam [25] as the benchmark optimization solver (the learning rate and batch  
 194 size for each experiment are listed in Table 1) and Dormand–Prince-45 as the numerical ODE solver.  
 195 For HBNODE and GHBNODE, we set  $\gamma = \text{sigmoid}(\theta)$ , where  $\theta$  is a trainable weight initialized as  
 196  $\theta = -3$ . The network architecture used to parameterize  $f(\mathbf{h}(t), t, \theta)$  for each experiment below are  
 197 described in Appendix E. All experiments are conducted on a server with 2 NVIDIA Titan Xp GPUs.

Table 1: The batch size and learning rate for different datasets.

Dataset	Point Cloud	MNIST	CIFAR10	Plane Vibration	Walker2D
Batch Size	50	64	64	64	256
Learning Rate	0.01	0.001	0.001	0.0001	0.003

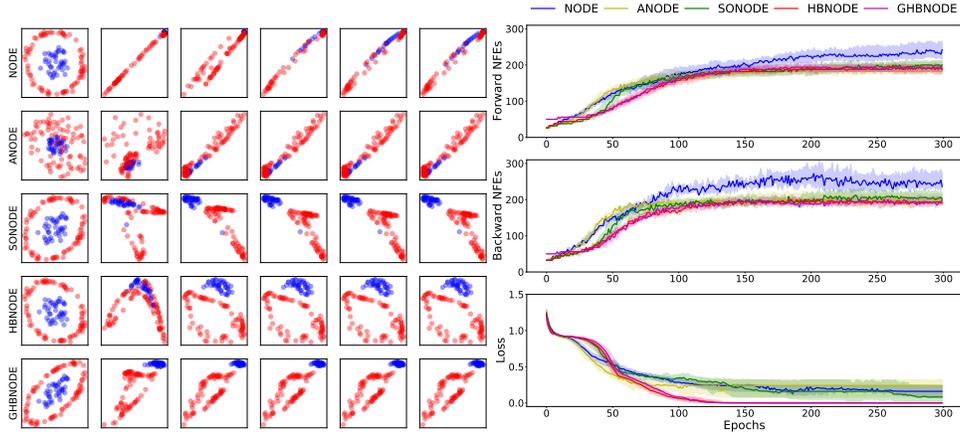


Figure 5: Comparison between NODE, ANODE, SONODE, HBNODE, and GHBNODE for two-dimensional point cloud separation. HBNODE and GHBNODE converge better and require less NFEs in both forward and backward propagation than the other benchmark models.

### 198 5.1 Point cloud separation

199 In this subsection, we consider the two-dimensional point cloud separation benchmark. A total of 120  
 200 points are sampled, in which 40 points are drawn uniformly from the circle  $\|\mathbf{r}\| < 0.5$ , and 80 points  
 201 are drawn uniformly from the annulus  $0.85 < \|\mathbf{r}\| < 1.0$ . This experiment aims to learn effective  
 202 features to classify these two point clouds. Following [10], we use a three-layer neural network to  
 203 parameterize the right-hand side of each ODE-based model, integrate the ODE-based model from  
 204  $t_0 = 0$  to  $T = 1$ , and pass the integration results to a dense layer to generate the classification  
 205 results. We set the size of hidden layers so that the models have similar sizes, and the number of

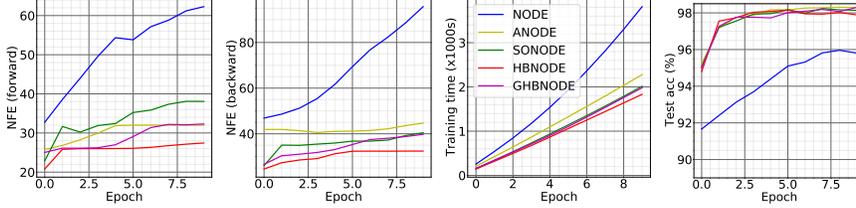


Figure 6: Contrasting NODE [4], ANODE [10], SONODE [37], HBNODE, and GHBNODE for MNIST classification in NFE, training time, and test accuracy. (Tolerance:  $10^{-5}$ ).

206 parameters of NODE, ANODE, SONODE, HBNODE, and GHBNODE are 525, 567, 528, 568,  
 207 and 568, respectively. To avoid the effects of numerical error of the black-box ODE solver we set  
 208 tolerance of ODE solver to be  $10^{-7}$ . Figure 5 plots a randomly selected evolution of the point cloud  
 209 separation for each model; we also compare the forward and backward NFEs and the training loss of  
 210 these models (100 independent runs). HBNODE and GHBNODE improve training as the training  
 211 loss consistently goes to zero over different runs, while ANODE and SONODE often get stuck at  
 212 local minima, and NODE cannot separate the point cloud since it preserves the topology [10].

## 213 5.2 Image classification

214 We compare the performance of HBNODE and GHBNODE with the existing ODE-based models  
 215 on MNIST and CIFAR10 classification tasks using the same setting as in [10]. We parameterize  
 216  $f(\mathbf{h}(t), t, \theta)$  using a 3-layer convolutional network for each ODE-based model, and the total number  
 217 of parameters for each model is listed in Table 2. For a given input image of the size  $c \times h \times w$ , we  
 218 first augment the number of channel from  $c$  to  $c + p$  with the augmentation dimension  $p$  dependent  
 219 on each method<sup>3</sup>. Moreover, for SONODE, HBNODE and GHBNODE, we further include velocity  
 220 or momentum with the same shape as the augmented state.

Table 2: The number of parameters for each models for image classification.

Model	NODE	ANODE	SONODE	HBNODE	GHBNODE
#Params (MNIST)	85,315	85,462	86,179	85,931	85,235
#Params (CIFAR10)	173,611	172,452	171,635	172,916	172,916

221 **NFEs.** As shown in Figs. 1 and 6, the NFEs grow rapidly with training of the NODE, resulting  
 222 in an increasingly complex model with reduced performance and the possibility of blow up. Input  
 223 augmentation has been verified to effectively reduce the NFEs, as both ANODE and SONODE  
 224 require fewer forward NFEs than NODE for the MNIST and CIFAR10 classification. However,  
 225 input augmentation is less effective in controlling their backward NFEs. HBNODE and GHBNODE  
 226 require much fewer NFEs than the existing benchmarks, especially for backward NFEs. In practice,  
 227 reducing NFEs implies reducing both training and inference time, as shown in Figs. 1 and 6.

228 **Accuracy.** We also compare the accuracy of different ODE-based models for MNIST and CIFAR10  
 229 classification. As shown in Figs. 1 and 6, HBNODE and GHBNODE have slightly better classification  
 230 accuracy than the other three models; this resonates with the fact that less NFEs lead to simpler  
 231 models which generalize better [10, 37].

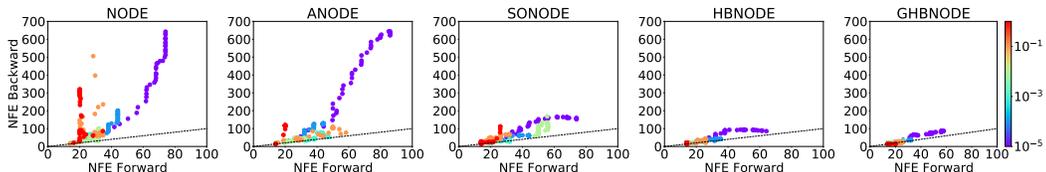


Figure 7: NFE vs. tolerance (shown in the colorbar) for training ODE-based models for CIFAR10 classification. Both forward and backward NFEs of HBNODE and GHBNODE grow much more slowly than that of NODE, ANODE, and SONODE; especially the backward NFEs. As the tolerance decreases, the advantage of HBNODE and GHBNODE in reducing NFEs becomes more significant.

<sup>3</sup>We set  $p = 0, 5, 4, 4, 5/0, 10, 9, 9, 9$  on MNIST/CIFAR10 for NODE, ANODE, SONODE, HBNODE, and GHBNODE, respectively.

232 **NFEs vs. tolerance.** We further study the NFEs for different ODE-based models under different  
 233 tolerances of the ODE solver using the same approach as in [4]. Figure 7 depicts the forward  
 234 and backward NFEs for different models under different tolerances. We see that (i) both forward  
 235 and backward NFEs grow quickly when tolerance is decreased, and HBNODE and GHBNODE  
 236 require much fewer NFEs than other models; (ii) under different tolerances, the backward NFEs of  
 237 NODE, ANODE, and SONODE are much larger than the forward NFEs, and the difference becomes  
 238 larger when the tolerance decreases. In contrast, the forward and backward NFEs of HBNODE  
 239 and GHBNODE scale almost linearly with each other. This reflects that the advantage in NFEs of  
 240 (G)HBNODE over the benchmarks become more significant when a smaller tolerance is used.

### 241 5.3 Learning dynamical systems from irregularly-sampled time series

242 In this subsection, we learn dynamical systems from experimental measurements. In particular, we  
 243 use the ODE-RNN framework [4, 44], with the recognition model being set to different ODE-based  
 244 models, to study the vibration of an airplane dataset [35]. The dataset was acquired, from time 0 to  
 245 73627, by attaching a shaker underneath the right wing to provide input signals, and 5 attributes are  
 246 recorded per time stamp; these attributes include voltage of input signal, force applied to aircraft,  
 247 and acceleration at 3 different spots of the airplane. We randomly take out 10% of the data to  
 248 make the time series irregularly-sampled. We use the first 50% of data as our train set, the next  
 249 25% as validation set, and the rest as test set. We divide each set into non-overlapping segments of  
 250 consecutive 65 time stamps of the irregularly-sampled time series, with each input instance consisting  
 251 of 64 time stamps of the irregularly-sampled time series, and we aim to forecast 8 consecutive time  
 252 stamps starting from the last time stamp of the segment. The input is fed though the the hybrid  
 253 methods in a recurrent fashion; by changing the time duration of the last step of the ODE integration,  
 254 we can forecast the output in the different time stamps. The output of the hybrid method is passed  
 255 to a single dense layer to generate the output time series. In our experiments, we compare different  
 256 ODE-based models hybrid with RNNs. The ODE of each model is parametrized by a 3-layer network  
 257 whereas the RNN is parametrized by a simple dense network; the total number of parameters for  
 258 ODE-RNN, ANODE-RNN, SONODE-RNN, HBNODE-RNN, and GHBNODE-RNN with 16, 22,  
 259 14, 15, 15 augmented dimensions are 15,986, 16,730, 16,649, 16,127, and 16,127, respectively. To  
 260 avoid potential error due to the ODE solver, we use a tolerance of  $10^{-7}$ .

261 In training those hybrid models, we regularize the models by penalizing the L2 distance between the  
 262 RNN output and the values of the next time stamp. Due to the second-order natural of the underlying  
 263 dynamics [37], ODE-RNN and ANODE-RNN learn the dynamics very poorly with much larger  
 264 training and test losses than the other models even they take smaller NFEs. HBNODE-RNN and  
 265 GHBNODE-RNN give better prediction than SONODE-RNN using less backward NFEs.

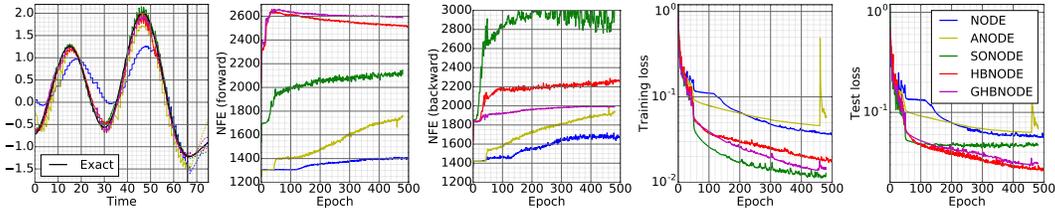


Figure 8: Contrasting ODE-RNN, ANODE-RNN, SONODE-RNN, HBNODE-RNN, and GHBNODE-RNN for learning a vibrational dynamical system. Left most: The learned curves of each model vs. the ground truth (Time: <66 for training, 66-75 for testing).

### 266 5.4 Walker2D kinematic simulation

267 In this subsection, we evaluate the performance of HBNODE-RNN and GHBNODE-RNN on the  
 268 Walker2D kinematic simulation task, which requires learning long-term dependency effectively [27].  
 269 The dataset [3] consists of a dynamical system from kinematic simulation of a person walking from  
 270 a pre-trained policy, aiming to learn the kinematic simulation of the MuJoCo physics engine [51].  
 271 The dataset is irregularly-sampled where 10% of data are removed from the simulation. Each input  
 272 is consisted of 64 time stamps and fed though the the hybrid methods in a recurrent fashion, and  
 273 the outputs of hybrid methods is passed to a single dense layer to generate the output time series.  
 274 The target is to provide auto-regressive forecast so that the output time series is as close as the  
 275 input sequence shifted 1 time stamp to the right. We compare ODE-RNN (with 7 augmentation),  
 276 ANODE-RNN (with 7 ANODE style augmentation), HBNODE-RNN (with 7 augmentation), and

277 GHBNODE-RNN (with 7 augmentation)<sup>4</sup>. The RNN is parametrized by a 3-layer network whereas  
 278 the ODE is parametrized by a simple dense network. The number of parameters of the above four  
 279 models are 8,729, 8,815, 8,899, and 8,899, respectively. In Fig. 9, we compare the performance of  
 280 the above four models on the Walker2D benchmark; HBNODE-RNN and GHBNODE-RNN not  
 281 only require significantly less NFEs in both training (forward and backward) and in testing than  
 282 ODE-RNN and ANODE-RNN, but also have much smaller training and test losses.

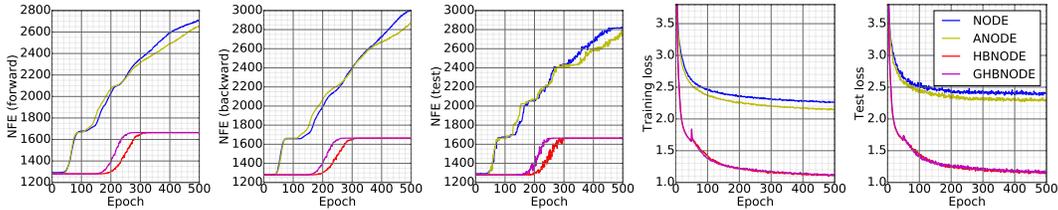


Figure 9: Contrasting ODE-RNN, ANODE-RNN, SONODE-RNN, HBNODE-RNN, and GHBNODE-RNN for the Walker-2D kinematic simulation.

## 283 6 Related Work

284 **Reducing NFEs in training NODEs.** Several techniques have been developed to reduce the  
 285 NFEs for the forward solvers in NODEs, including weight decay [15], input augmentation [10],  
 286 regularizing the learning dynamics [13], high-order ODE [37], data control [31], and depth-variance  
 287 [31]. HBNODEs can reduce both forward and backward NFEs at the same time.

288 **Second-order ODE accelerated dynamics.** It has been noticed in both optimization and sampling  
 289 communities that second-order ODEs with an appropriate damping term, e.g., the classical momentum  
 290 and Nesterov’s acceleration in discrete regime, can significantly accelerate the first-order gradient  
 291 dynamics (gradient descent), e.g., [40, 34, 5, 48, 55]. Also, these second-order ODEs have been  
 292 discretized via some interesting numerical schemes to design fast optimization schemes, e.g., [47].

293 **Learning long-term dependencies.** Learning long-term dependency is one of the most important  
 294 goals for learning from sequential data. Most of the existing works focus on mitigating exploding  
 295 or vanishing gradient issues in training RNNs, e.g., [1, 56, 23, 54, 32, 19, 50]. Attention-based  
 296 models are proposed for learning on sequential data concurrently with the effective accommodation  
 297 of learning long-term dependency [53, 8]. Recently, NODEs have been integrated with long-short  
 298 term memory model [20] to learn long-term dependency for irregularly-sampled time series [27].  
 299 HBNODEs directly enhance learning long-term dependency from sequential data.

300 **Momentum in neural network design.** As a line of orthogonal work, the momentum has also been  
 301 studied in designing neural network architecture, e.g., [33, 50, 45], which can also help accelerate  
 302 training and learn long-term dependencies. These techniques can be considered as changing the  
 303 neural network  $f$  in (1). We leave the synergistic integration of adding momentum to  $f$  with our  
 304 work on changing the left-hand side of (1) as a future work.

305 **ResNet-style models.** Interpreting ResNet as an ODE model has been an interesting research  
 306 direction [11, 16], which has lead to interesting neural network architectures and analysis from the  
 307 numerical ODE solvers and differential equation theory viewpoints, e.g., [16, 30, 28].

## 308 7 Concluding Remarks

309 We proposed HBNODEs to reduce the NFEs in solving both forward and backward ODEs, which  
 310 also improve generalization performance over the existing benchmark models. Moreover, HBNODEs  
 311 alleviate vanishing gradients in training NODEs, making HBNODEs able to learn long-term depen-  
 312 dency effectively from sequential data. In the optimization community, Nesterov acceleration [34]  
 313 is also a famous algorithm for accelerating gradient descent, that achieves an optimal convergence  
 314 rate for general convex optimization problems. The ODE counterpart of the Nesterov’s acceleration  
 315 corresponds to (9) with  $\gamma$  being replaced by a time-dependent damping parameter, e.g.,  $t/3$  [48]. The  
 316 adjoint equation of the Nesterov’s ODE [48] is no longer a Nesterov’s ODE. We notice that directly  
 317 using Nesterov’s ODE cannot improve the performance of the vanilla neural ODE. How to integrate  
 318 Nesterov’s ODE with neural ODE is an interesting future direction.

<sup>4</sup>Here, we do not compare with SONODE-RNN since SONODE has some initialization problem on this dataset, and the ODE solver encounters failure due to exponential growth over time. This issue is originally tackled by re-initialization [37]. We re-initialized SONODE 100 times; all failed due to initialisation problems.

319 **References**

- 320 [1] Martin Arjovsky, Amar Shah, and Yoshua Bengio. Unitary evolution recurrent neural networks.  
321 In *International Conference on Machine Learning*, pages 1120–1128, 2016.
- 322 [2] Yoshua Bengio, Patrice Simard, and Paolo Frasconi. Learning long-term dependencies with  
323 gradient descent is difficult. *IEEE Transactions on Neural Networks*, 5(2):157–166, 1994.
- 324 [3] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang,  
325 and Wojciech Zaremba. OpenAI Gym, 2016. cite arxiv:1606.01540.
- 326 [4] Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary dif-  
327 ferential equations. In *Proceedings of the 32nd International Conference on Neural Information*  
328 *Processing Systems*, pages 6572–6583, 2018.
- 329 [5] Tianqi Chen, Emily Fox, and Carlos Guestrin. Stochastic gradient Hamiltonian Monte Carlo.  
330 In *International conference on machine learning*, pages 1683–1691, 2014.
- 331 [6] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares,  
332 Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-  
333 decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- 334 [7] Talgat Daulbaev, Alexandr Katrutsa, Larisa Markeeva, Julia Gusak, Andrzej Cichocki, and  
335 Ivan Oseledets. Interpolation technique to speed up gradients propagation in neural odes. In  
336 H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural*  
337 *Information Processing Systems*, volume 33, pages 16689–16700. Curran Associates, Inc.,  
338 2020.
- 339 [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of  
340 deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*,  
341 2018.
- 342 [9] Jianzhun Du, Joseph Futoma, and Finale Doshi-Velez. Model-based reinforcement learning  
343 for semi-markov decision processes with neural odes. In H. Larochelle, M. Ranzato, R. Had-  
344 sell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*,  
345 volume 33, pages 19805–19816. Curran Associates, Inc., 2020.
- 346 [10] Emilien Dupont, Arnaud Doucet, and Yee Whye Teh. Augmented neural odes. In H. Wallach,  
347 H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in*  
348 *Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- 349 [11] Weinan E. A proposal on machine learning via dynamical systems. *Communications in*  
350 *Mathematics and Statistics*, 5:1–11, 2017.
- 351 [12] Jeffrey L Elman. Finding structure in time. *Cognitive Science*, 14(2):179–211, 1990.
- 352 [13] Chris Finlay, Joern-Henrik Jacobsen, Levon Nurbekyan, and Adam Oberman. How to train  
353 your neural ODE: the world of Jacobian and kinetic regularization. In Hal Daumé III and Aarti  
354 Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume  
355 119 of *Proceedings of Machine Learning Research*, pages 3154–3164. PMLR, 13–18 Jul 2020.
- 356 [14] Amilcar dos Santos Gonçalves. A Version of Beale’s Method Avoiding the Free-Variables. In  
357 *Proceedings of the 1971 26th Annual Conference*, ACM ’71, page 433–441, New York, NY,  
358 USA, 1971. Association for Computing Machinery.
- 359 [15] Will Grathwohl, Ricky T. Q. Chen, Jesse Bettencourt, and David Duvenaud. Scalable re-  
360 versible generative models with free-form continuous dynamics. In *International Conference*  
361 *on Learning Representations*, 2019.
- 362 [16] Eldad Haber and Lars Ruthotto. Stable architectures for deep neural networks. *Inverse Problems*,  
363 34(1):014004, 2017.
- 364 [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
365 recognition. *arXiv preprint arXiv:1512.03385*, 2015.

- 366 [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual  
367 networks. In *European conference on computer vision*, pages 630–645. Springer, 2016.
- 368 [19] Kyle Helfrich, Devin Willmott, and Qiang Ye. Orthogonal recurrent neural networks with  
369 scaled Cayley transform. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th  
370 International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning  
371 Research*, pages 1969–1978, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR.
- 372 [20] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*,  
373 9(8):1735–1780, 1997.
- 374 [21] Zijie Huang, Yizhou Sun, and Wei Wang. Learning continuous system dynamics from  
375 irregularly-sampled partial observations. In *Advances in Neural Information Processing Systems*,  
376 2020.
- 377 [22] Junteng Jia and Austin R Benson. Neural jump stochastic differential equations. In H. Wallach,  
378 H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in  
379 Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- 380 [23] Li Jing, Yichen Shen, Tena Dubcek, John Peurifoy, Scott Skirlo, Yann LeCun, Max Tegmark,  
381 and Marin Soljačić. Tunable efficient unitary neural networks (eunn) and their application to  
382 rnns. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*,  
383 pages 1733–1741. JMLR. org, 2017.
- 384 [24] Patrick Kidger, James Morrill, James Foster, and Terry J. Lyons. Neural controlled differential  
385 equations for irregular time series. In *NeurIPS*, 2020.
- 386 [25] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint  
387 arXiv:1412.6980*, 2014.
- 388 [26] Nikola B. Kovachki and Andrew M. Stuart. Continuous time analysis of momentum methods.  
389 *Journal of Machine Learning Research*, 22(17):1–40, 2021.
- 390 [27] Mathias Lechner and Ramin Hasani. Learning long-term dependencies in irregularly-sampled  
391 time series. *arXiv preprint arXiv:2006.04418*, 2020.
- 392 [28] Qianxiao Li, Ting Lin, and Zuwei Shen. Deep learning via dynamical systems: An approxima-  
393 tion perspective. *arXiv preprint arXiv:1912.10382*, 2019.
- 394 [29] Xuechen Li, Ting-Kam Leonard Wong, Ricky T. Q. Chen, and David Duvenaud. Scalable  
395 gradients for stochastic differential equations. In Silvia Chiappa and Roberto Calandra, editors,  
396 *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statis-  
397 tics*, volume 108 of *Proceedings of Machine Learning Research*, pages 3870–3882. PMLR,  
398 26–28 Aug 2020.
- 399 [30] Yiping Lu, Aoxiao Zhong, Quanzheng Li, and Bin Dong. Beyond finite layer neural networks:  
400 Bridging deep architectures and numerical differential equations. In *International Conference  
401 on Machine Learning*, pages 3276–3285. PMLR, 2018.
- 402 [31] Stefano Massaroli, Michael Poli, Jinkyoo Park, Atsushi Yamashita, and Hajime Asma. Dissect-  
403 ing neural odes. In *34th Conference on Neural Information Processing Systems, NeurIPS 2020*.  
404 The Neural Information Processing Systems, 2020.
- 405 [32] Zakaria Mhammedi, Andrew Hellicar, Ashfaque Rahman, and James Bailey. Efficient orthogonal  
406 parametrisation of recurrent neural networks using householder reflections. In *Proceedings of  
407 the 34th International Conference on Machine Learning-Volume 70*, pages 2401–2409. JMLR.  
408 org, 2017.
- 409 [33] Thomas Moreau and Joan Bruna. Understanding the learned iterative soft thresholding algorithm  
410 with matrix factorization. *arXiv preprint arXiv:1706.01338*, 2017.
- 411 [34] Yurii E Nesterov. A method for solving the convex programming problem with convergence  
412 rate  $o(1/k^2)$ . In *Dokl. Akad. Nauk Sssr*, volume 269, pages 543–547, 1983.

- 413 [35] Jean-Philippe Noël and M Schoukens. F-16 aircraft benchmark based on ground vibration test  
414 data. In *2017 Workshop on Nonlinear System Identification Benchmarks*, pages 19–23, 2017.
- 415 [36] Alexander Norcliffe, Cristian Bodnar, Ben Day, Jacob Moss, and Pietro Liò. Neural {ode}  
416 processes. In *International Conference on Learning Representations*, 2021.
- 417 [37] Alexander Norcliffe, Cristian Bodnar, Ben Day, Nikola Simidjievski, and Pietro Liò. On second  
418 order behaviour in augmented neural odes. In *Advances in Neural Information Processing*  
419 *Systems*, 2020.
- 420 [38] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent  
421 neural networks. In *International Conference on Machine Learning*, pages 1310–1318, 2013.
- 422 [39] Michael Poli, Stefano Massaroli, Junyoung Park, Atsushi Yamashita, Hajime Asama, and  
423 Jinkyoo Park. Graph neural ordinary differential equations. *arXiv preprint arXiv:1911.07532*,  
424 2019.
- 425 [40] Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *USSR*  
426 *Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964.
- 427 [41] Lev Semenovich Pontryagin. *Mathematical theory of optimal processes*. Routledge, 2018.
- 428 [42] Alessio Quaglino, Marco Gallieri, Jonathan Masci, and Jan Koutník. Snode: Spectral dis-  
429 cretization of neural odes for system identification. In *International Conference on Learning*  
430 *Representations*, 2020.
- 431 [43] H. H. Rosenbrock. An Automatic Method for Finding the Greatest or Least Value of a Function.  
432 *The Computer Journal*, 3(3):175–184, 01 1960.
- 433 [44] Yulia Rubanova, Ricky T. Q. Chen, and David K Duvenaud. Latent ordinary differential  
434 equations for irregularly-sampled time series. In H. Wallach, H. Larochelle, A. Beygelzimer,  
435 F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing*  
436 *Systems*, volume 32. Curran Associates, Inc., 2019.
- 437 [45] Michael E. Sander, Pierre Ablin, Mathieu Blondel, and Gabriel Peyré. Momentum residual  
438 neural networks. *arXiv preprint arXiv:2102.07870*, 2021.
- 439 [46] Bin Shi, Simon S. Du, Michael I. Jordan, and Weijie J. Su. Understanding the acceleration  
440 phenomenon via high-resolution differential equations. *arXiv preprint arXiv:1810.08907*, 2018.
- 441 [47] Bin Shi, Simon S Du, Weijie Su, and Michael I Jordan. Acceleration via symplectic discretiza-  
442 tion of high-resolution differential equations. In H. Wallach, H. Larochelle, A. Beygelzimer,  
443 F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing*  
444 *Systems*, volume 32. Curran Associates, Inc., 2019.
- 445 [48] Weijie Su, Stephen Boyd, and Emmanuel Candes. A differential equation for modeling nes-  
446 terov’s accelerated gradient method: Theory and insights. In *Advances in Neural Information*  
447 *Processing Systems*, pages 2510–2518, 2014.
- 448 [49] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initial-  
449 ization and momentum in deep learning. In *International Conference on Machine Learning*,  
450 pages 1139–1147, 2013.
- 451 [50] Tan M. Nguyen and Richard G. Baraniuk and Andrea L. Bertozzi and Stanley J. Osher and Bao  
452 Wang. MomentumRNN: Integrating momentum into recurrent neural networks. In *Advances in*  
453 *Neural Information Processing Systems (NeurIPS)*, pages 9154–9164, 2020.
- 454 [51] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based  
455 control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages  
456 5026–5033, 2012.
- 457 [52] Belinda Tzen and Maxim Raginsky. Neural stochastic differential equations: Deep latent  
458 gaussian models in the diffusion limit. *arXiv preprint arXiv:1905.09883*, 2019.

- 459 [53] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
460 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg,  
461 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural  
462 Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 463 [54] Eugene Vorontsov, Chiheb Trabelsi, Samuel Kadoury, and Chris Pal. On orthogonality and learn-  
464 ing recurrent networks with long term dependencies. In *Proceedings of the 34th International  
465 Conference on Machine Learning-Volume 70*, pages 3570–3578. JMLR. org, 2017.
- 466 [55] Ashia C. Wilson, Benjamin Recht, and Michael I. Jordan. A Lyapunov Analysis of Momentum  
467 Methods in Optimization. *arXiv preprint arXiv:1611.02635*, 2018.
- 468 [56] Scott Wisdom, Thomas Powers, John Hershey, Jonathan Le Roux, and Les Atlas. Full-capacity  
469 unitary recurrent neural networks. In *Advances in Neural Information Processing Systems*,  
470 pages 4880–4888, 2016.
- 471 [57] Cagatay Yildiz, Markus Heinonen, and Harri Lahdesmaki. ODE2VAE: Deep generative second  
472 order ODEs with Bayesian neural networks. In H. Wallach, H. Larochelle, A. Beygelzimer,  
473 F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing  
474 Systems*, volume 32. Curran Associates, Inc., 2019.
- 475 [58] Tianjun Zhang, Zhewei Yao, Amir Gholami, Joseph E Gonzalez, Kurt Keutzer, Michael W  
476 Mahoney, and George Biros. ANODEV2: A Coupled Neural ODE Framework. In H. Wallach,  
477 H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in  
478 Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- 479 [59] Juntang Zhuang, Nicha C Dvornek, sekhar tatikonda, and James s Duncan. {MALI}: A memory  
480 efficient and reverse accurate integrator for neural {ode}s. In *International Conference on  
481 Learning Representations*, 2021.

482 **Checklist**

- 483 1. For all authors...
- 484 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
485 contributions and scope? [Yes]
- 486 (b) Did you describe the limitations of your work? [Yes] See Section 4.1.
- 487 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 488 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
489 them? [Yes]
- 490 2. If you are including theoretical results...
- 491 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Section  
492 4.1.
- 493 (b) Did you include complete proofs of all theoretical results? [Yes] See Supplementary  
494 Materials
- 495 3. If you ran experiments...
- 496 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
497 mental results (either in the supplemental material or as a URL)? [Yes]
- 498 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
499 were chosen)? [Yes]
- 500 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
501 ments multiple times)? [Yes]
- 502 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
503 of GPUs, internal cluster, or cloud provider)? [Yes]
- 504 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 505 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 506 (b) Did you mention the license of the assets? [Yes]
- 507 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 508
- 509 (d) Did you discuss whether and how consent was obtained from people whose data you're  
510 using/curating? [N/A]
- 511 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
512 information or offensive content? [N/A]
- 513 5. If you used crowdsourcing or conducted research with human subjects...
- 514 (a) Did you include the full text of instructions given to participants and screenshots, if  
515 applicable? [N/A]
- 516 (b) Did you describe any potential participant risks, with links to Institutional Review  
517 Board (IRB) approvals, if applicable? [N/A]
- 518 (c) Did you include the estimated hourly wage paid to participants and the total amount  
519 spent on participant compensation? [N/A]