
Heavy Ball Neural Ordinary Differential Equations

Hedi Xia*

Department of Mathematics
University of California, Los Angeles

Vai Suliafu *

Scientific Computing and Imaging (SCI) Institute
University of Utah, Salt Lake City, UT, USA

Hangjie Ji

Department of Mathematics
University of California, Los Angeles

Tan M. Nguyen

Department of Mathematics
University of California, Los Angeles

Andrea L. Bertozzi

Department of Mathematics
University of California, Los Angeles

Stanley J. Osher

Department of Mathematics
University of California, Los Angeles

Bao Wang [†]

Department of Mathematics
Scientific Computing and Imaging (SCI) Institute
University of Utah, Salt Lake City, UT, USA

Abstract

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

1 Introduction

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

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

Starting from the input $\mathbf{h}(t_0)$, NODEs obtain the output $\mathbf{h}(T)$ by solving (1) for $t_0 \leq t \leq T$ with the initial value $\mathbf{h}(t_0)$, using a black-box numerical ODE solver. The

*Co-first author

[†]Please correspond to: wangbaonj@gmail.com

number of function evaluations (NFEs) that the black-box ODE solver requires in a single forward pass is an analogue for the continuous-depth models [4] to the depth of networks in ResNets [17]. The loss between NODE prediction $\mathbf{h}(T)$ and the ground truth is denoted by $\mathcal{L}(\mathbf{h}(T))$; we update parameters θ using the following gradient [44]

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

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

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

NODEs are flexible in learning from irregularly-sampled sequential data and particularly suitable for learning complex dynamical systems [4, 47, 62, 39, 9, 25], which can be trained by efficient algorithms [45, 7, 64]. NODE-based continuous generative models have computational advantages over the classical normalizing flows [4, 15, 61, 12]. NODEs have also been generalized to neural stochastic differential equations, stochastic processes, and graph NODEs [22, 31, 42, 55, 21, 38]. The drawback of NODEs is also prominent. In many ML tasks, NODEs require very high NFEs in both training and inference, especially in high accuracy settings where a lower tolerance is needed. The NFEs increase rapidly with training; high NFEs reduce computational speed and accuracy of NODEs and can lead to blow-ups in the worst-case scenario [15, 10, 33, 39]. As an illustration, we train NODEs for CIFAR10 classification using the same model and experimental settings as in [10], except using a tolerance of 10^{-5} ; Fig. 1 shows both forward and backward NFEs and the training time of different ODE-based models; we see that NFEs and computational times increase very rapidly for NODE, ANODE [10], and SONODE [39]. More results on the large NFE and degrading utility issues for different benchmark experiments are available in Sec. 5. Another issue is that NODEs often fail to effectively learn long-term dependencies in sequential data [28], discussed in Sec. 4.

1.1 Contribution

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

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

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

1.2 Organization

We organize the paper as follows: In Secs 2 and 3, we present our motivation, algorithm, and analysis of HBNODEs and GHBNODEs, respectively. We analyze the spectrum structure of the adjoint equation of HBNODEs/GHBNODEs in Sec. 4, which indicates that HBNODEs/GHBNODEs can

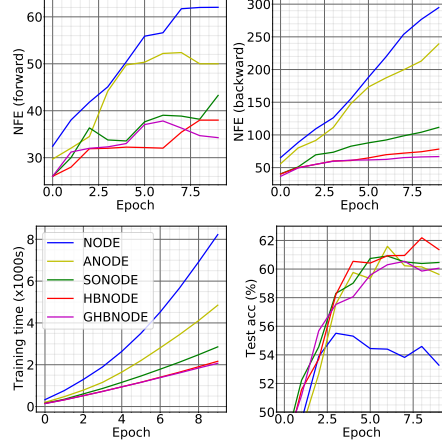


Figure 1: Contrasting NODE, ANODE, SONODE, HBNODE, and GHBNODE for CIFAR10 classification in NFEs, training time, and test accuracy. (Tolerance: 10^{-5} , see Sec. 5.2 for experimental details.)

learn long-term dependency effectively. We test the performance of HBNODEs and GHBNODEs on benchmark point cloud separation, image classification, learning dynamics, and sequential modeling in Sec. 5. We discuss more related work in Sec. 6, followed by concluding remarks. Technical proofs and more experimental details are provided in the appendix.

2 Heavy Ball Neural Ordinary Differential Equations

2.1 Heavy ball ordinary differential equation

Classical momentum method, a.k.a., the heavy ball method, has achieved remarkable success in accelerating gradient descent [43] and has significantly improved the training of deep neural networks [52]. As the continuous limit of the classical momentum method, heavy ball ODE (HBODE) has been studied in various settings and has been used to analyze the acceleration phenomenon of the momentum methods. For the ease of reading and completeness, we derive the HBODE from the classical momentum method. Starting from initial points \mathbf{x}^0 and \mathbf{x}^1 , gradient descent with classical momentum searches a minimum of the function $F(\mathbf{x})$ through the following iteration

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

where $s > 0$ is the step size and $0 \leq \beta < 1$ is the momentum hyperparameter. For any fixed step 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. Then we can rewrite (4) as

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

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

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

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

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

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

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

In particular, we solve the ODEs (7) and (8) with $F(\mathbf{x})$ defined as a Rosenbrock [46] or Beale [14] function (see Appendix E.6 for experimental details). Fig. 2 shows that with the same numerical ODE solver, HBODE converges to the stationary point (marked by stars) faster than (8). The fact that HBODE can accelerate the dynamics of the ODE for a gradient system motivates us to propose HBNODE to accelerate forward propagation of NODE.

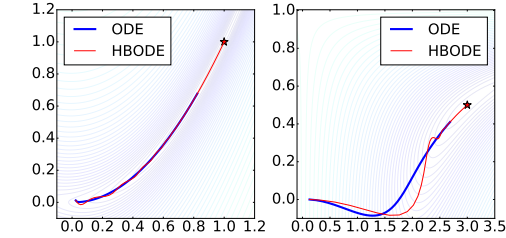


Figure 2: Comparing the trajectory of ODE and HBODE when $F(\mathbf{x})$ is the Rosenbrock (left) and Beale (right) functions.

2.2 Heavy ball neural ordinary differential equations

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

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

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

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

Equation (9) (or equivalently, the system (10)) defines the forward ODE for the HBNODE, and we can use either the first-order (Prop. 2) or the second-order (Prop. 1) adjoint sensitivity method to update the parameter θ [39].

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

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

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

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

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

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

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

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

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 heavy ball ODE system*

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

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

Similar to [39], we use the coupled first-order HBNODE system (10) and its adjoint first-order HBNODE system (13) for practical implementation, since the entangled representation permits faster computation [39] of the gradients of the coupled ODE systems.

3 Generalized Heavy Ball Neural Ordinary Differential Equations

In this section, we propose a generalized version of HBNODE (GHBNODE), see (15), to mitigate the potential blow-up issue in training ODE-based models. In our experiments, we observe that $\mathbf{h}(t)$ of ANODEs [10], SONODEs [39], and HBNODEs (10) usually grows much faster than that of NODEs. The fast growth of $\mathbf{h}(t)$ can lead to finite-time blow up. As an illustration, we compare the performance of NODE, ANODE, SONODE, HBNODE, and GHBNODE on the Silverbox task as in [39]. The goal of the task is to learn the voltage of an electronic circuit that resembles a Duffing oscillator, where the input voltage $V_1(t)$ is used to predict the output $V_2(t)$. Similar to the setting in [39], we first augment ANODE by 1 dimension with 0-augmentation and augment SONODE, HBNODE, and GHBNODE with a dense network. We use a simple dense layer to parameterize f for all five models, with an extra input term for $V_1(t)^3$. For both HBNODE and GHBNODE, we set the damping parameter γ to be $\text{sigmoid}(-3)$. For GHBNODE (15) below, we set $\sigma(\cdot)$ to be the hardtanh function with bound $[-5, 5]$ and $\xi = \ln(2)$. The detailed architecture can be found in Appendix E. As shown in Fig. 3, compared to the vanilla NODE, the ℓ_2 norm of $\mathbf{h}(t)$ grows much faster when a higher order NODE is used, which leads to blow-up during training. Similar issues arise in the time series experiments (see Sec. 5.4), where SONODE blows up during long term integration in time, and HBNODE suffers from the same issue with some initialization.

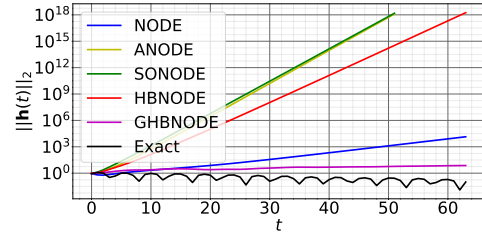


Figure 3: Contrasting $\mathbf{h}(t)$ for different models. $\mathbf{h}(t)$ in ANODE, SONODE, and HBNODE grows much faster than that in NODE. GHBNODE controls the growth of $\mathbf{h}(t)$ effectively when t is large.

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

To alleviate the problem above, we propose the following generalized HBNODE

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

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

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

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

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

4 Learning long-term dependencies – Vanishing gradient

It is known that the vanishing and exploding gradients are two bottlenecks for training recurrent neural networks (RNNs) with long-term dependencies [2, 41] (see Appendix C for a brief review on the exploding and vanishing gradient issues in training RNNs). The exploding gradients issue can be effectively resolved via gradient clipping, training loss regularization, etc [41, 11]. Thus in practice the vanishing gradient is the major issue for learning long-term dependencies [41]. As the continuous analogue of RNN, NODEs as well as their hybrid ODE-RNN models, may also suffer from vanishing in the adjoint state $\mathbf{a}(t) := \partial\mathcal{L}/\partial\mathbf{h}(t)$ [28]. When the vanishing gradient issue happens, $\mathbf{a}(t)$ goes 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 following expressions for the adjoint states of the NODE and HBNODE (see Appendix C for detailed derivation):

- For NODE, we have

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

- For GHBNODE⁴, from (13) we can derive

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

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

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

Let $\mathbf{v}^\top := \left[\frac{\partial\mathcal{L}}{\partial\mathbf{h}_T} \quad \frac{\partial\mathcal{L}}{\partial\mathbf{m}_T} \right] \mathbf{Q}$, then (18) can be rewritten as

$$\left[\frac{\partial\mathcal{L}}{\partial\mathbf{h}_t} \quad \frac{\partial\mathcal{L}}{\partial\mathbf{m}_t} \right] = \left[\frac{\partial\mathcal{L}}{\partial\mathbf{h}_T} \quad \frac{\partial\mathcal{L}}{\partial\mathbf{m}_T} \right] \exp\{-\mathbf{M}\} = \left[\frac{\partial\mathcal{L}}{\partial\mathbf{h}_T} \quad \frac{\partial\mathcal{L}}{\partial\mathbf{m}_T} \right] \mathbf{Q} \exp\{\mathbf{U}\} \mathbf{Q}^\top = \mathbf{v}^\top \exp\{\mathbf{U}\} \mathbf{Q}^\top.\tag{20}$$

By taking the ℓ_2 norm in (20) and dividing both sides by $\left\| \left[\frac{\partial\mathcal{L}}{\partial\mathbf{h}_T} \quad \frac{\partial\mathcal{L}}{\partial\mathbf{m}_T} \right] \right\|_2$, we arrive at

$$\frac{\left\| \left[\frac{\partial\mathcal{L}}{\partial\mathbf{h}_t} \quad \frac{\partial\mathcal{L}}{\partial\mathbf{m}_t} \right] \right\|_2}{\left\| \left[\frac{\partial\mathcal{L}}{\partial\mathbf{h}_T} \quad \frac{\partial\mathcal{L}}{\partial\mathbf{m}_T} \right] \right\|_2} = \frac{\left\| \mathbf{v}^\top \exp\{\mathbf{U}\} \mathbf{Q}^\top \right\|_2}{\left\| \mathbf{v}^\top \mathbf{Q}^\top \right\|_2} = \frac{\left\| \mathbf{v}^\top \exp\{\mathbf{U}\} \right\|_2}{\left\| \mathbf{v} \right\|_2} = \left\| \mathbf{e}^\top \exp\{\mathbf{U}\} \right\|_2,\tag{21}$$

i.e., $\left\| \left[\frac{\partial\mathcal{L}}{\partial\mathbf{h}_t} \quad \frac{\partial\mathcal{L}}{\partial\mathbf{m}_t} \right] \right\|_2 = \left\| \mathbf{e}^\top \exp\{\mathbf{U}\} \right\|_2 \left\| \left[\frac{\partial\mathcal{L}}{\partial\mathbf{h}_T} \quad \frac{\partial\mathcal{L}}{\partial\mathbf{m}_T} \right] \right\|_2$ where $\mathbf{e} = \mathbf{v}/\|\mathbf{v}\|_2$.

⁴HBNODE can be seen as a special GHBNODE with $\xi = 0$ and σ be the identity map.

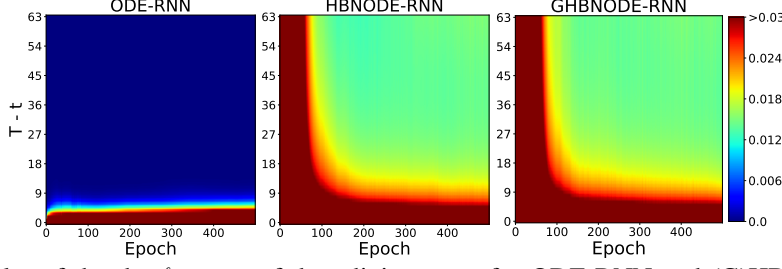


Figure 4: Plot of the ℓ_2 -norm of the adjoint states for ODE-RNN and (G)HBNODE-RNN back-propagated from the last time stamp. The adjoint state of ODE-RNN vanishes quickly when the gap between the final time T and intermediate time t becomes larger, while the adjoint states of (G)HBNODE-RNN decays much more slowly. This implies that (G)HBNODE-RNN is more effective in learning long-term dependency than ODE-RNN.

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

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

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

where the diagonal of \mathbf{U}_L (\mathbf{U}_V) contains eigenvalues of $-\mathbf{M}$ that are no less (greater) than $(t - T)a$. 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 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 of $-\mathbf{M}$ there is at least one eigenvalue whose real part is no less than $(t - T)a$. Therefore, $\exp\{\mathbf{U}_L\}$ 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 gradients by clipping the ℓ_2 norm of the adjoint states similar to that used for training RNNs.

In contrast, all eigenvalues of the matrix $\int_T^t \partial f / \partial \mathbf{h} ds$ in (17) for NODE can be very positive or negative, resulting in exploding or vanishing gradients. As an illustration, we consider the benchmark Walker2D kinematic simulation task that requires learning long-term dependencies effectively [28, 3]. We train ODE-RNN [47] and (G)HBNODE-RNN on this benchmark dataset, and the detailed experimental settings are provided in Sec. 5.4. Figure 4 plots $\|\partial \mathcal{L} / \partial \mathbf{h}_t\|_2$ for ODE-RNN 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 vanishes quickly, while that of (G)HBNODE-RNN does not vanish even when the gap between T and t is very large.

5 Experimental Results

In this section, we compare the performance of the proposed HBNODE and GHBNODE with existing ODE-based models, including NODE [4], ANODE [10], and SONODE [39] on the benchmark point cloud separation, image classification, learning dynamical systems, and kinematic simulation. For all the experiments, we use Adam [26] as the benchmark optimization solver (the learning rate and batch size for each experiment are listed in Table 1) and Dormand–Prince-45 as the numerical ODE solver. For HBNODE and GHBNODE, we set $\gamma = \text{sigmoid}(\theta)$, where θ is a trainable weight initialized as $\theta = -3$. The network architecture used to parameterize $f(\mathbf{h}(t), t, \theta)$ for each experiment below are 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

5.1 Point cloud separation

In this subsection, we consider the two-dimensional point cloud separation benchmark. A total of 120 points are sampled, in which 40 points are drawn uniformly from the circle $\|\mathbf{r}\| < 0.5$, and 80 points are drawn uniformly from the annulus $0.85 < \|\mathbf{r}\| < 1.0$. This experiment aims to learn effective

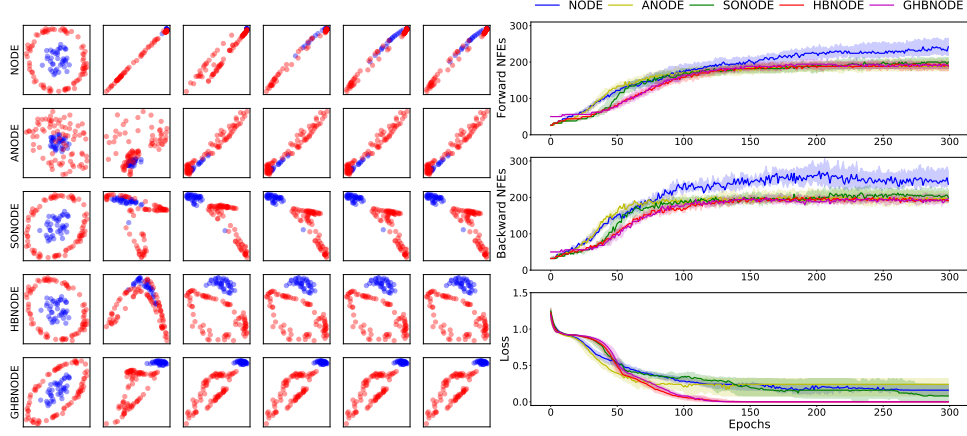


Figure 5: Comparison between NODE, ANODE, SONODE, HBNODE, and GHBNODE for two-dimensional point cloud separation. HBNODE and GHBNODE converge better and require less NFEs in both forward and backward propagation than the other benchmark models.

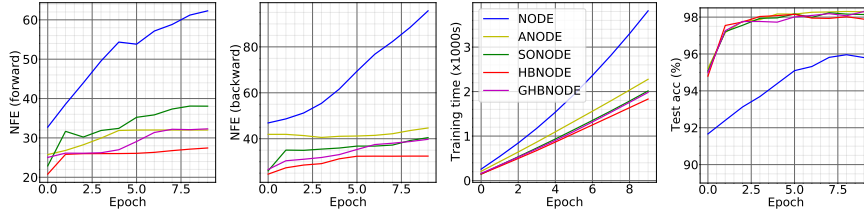


Figure 6: Contrasting NODE [4], ANODE [10], SONODE [39], HBNODE, and GHBNODE for MNIST classification in NFE, training time, and test accuracy. (Tolerance: 10^{-5}).

features to classify these two point clouds. Following [10], we use a three-layer neural network to parameterize the right-hand side of each ODE-based model, integrate the ODE-based model from $t_0 = 0$ to $T = 1$, and pass the integration results to a dense layer to generate the classification results. We set the size of hidden layers so that the models have similar sizes, and the number of parameters of NODE, ANODE, SONODE, HBNODE, and GHBNODE are 525, 567, 528, 568, and 568, respectively. To avoid the effects of numerical error of the black-box ODE solver we set tolerance of ODE solver to be 10^{-7} . Figure 5 plots a randomly selected evolution of the point cloud separation for each model; we also compare the forward and backward NFEs and the training loss of these models (100 independent runs). HBNODE and GHBNODE improve training as the training loss consistently goes to zero over different runs, while ANODE and SONODE often get stuck at local minima, and NODE cannot separate the point cloud since it preserves the topology [10].

5.2 Image classification

We compare the performance of HBNODE and GHBNODE with the existing ODE-based models on MNIST and CIFAR10 classification tasks using the same setting as in [10]. We parameterize $f(\mathbf{h}(t), t, \theta)$ using a 3-layer convolutional network for each ODE-based model, and the total number of parameters for each model is listed in Table 2. For a given input image of the size $c \times h \times w$, we first augment the number of channel from c to $c + p$ with the augmentation dimension p dependent on each method⁵. Moreover, for SONODE, HBNODE and GHBNODE, we further include velocity 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

NFEs. As shown in Figs. 1 and 6, the NFEs grow rapidly with training of the NODE, resulting in an increasingly complex model with reduced performance and the possibility of blow up. Input

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

augmentation has been verified to effectively reduce the NFEs, as both ANODE and SONODE require fewer forward NFEs than NODE for the MNIST and CIFAR10 classification. However, input augmentation is less effective in controlling their backward NFEs. HBNODE and GHBNODE require much fewer NFEs than the existing benchmarks, especially for backward NFEs. In practice, reducing NFEs implies reducing both training and inference time, as shown in Figs. 1 and 6.

Accuracy. We also compare the accuracy of different ODE-based models for MNIST and CIFAR10 classification. As shown in Figs. 1 and 6, HBNODE and GHBNODE have slightly better classification accuracy than the other three models; this resonates with the fact that less NFEs lead to simpler models which generalize better [10, 39].

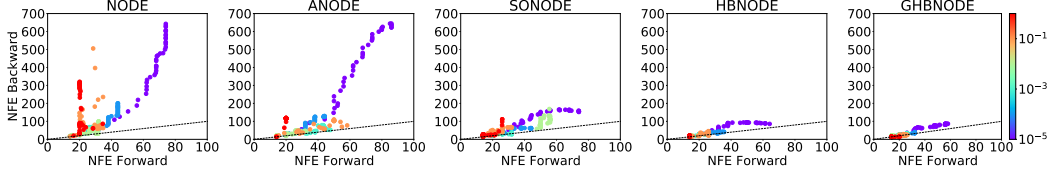


Figure 7: NFE vs. tolerance (shown in the colorbar) for training ODE-based models for CIFAR10 classification. Both forward and backward NFEs of HBNODE and GHBNODE grow much more slowly than that of NODE, ANODE, and SONODE; especially the backward NFEs. As the tolerance decreases, the advantage of HBNODE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

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

5.3 Learning dynamical systems from irregularly-sampled time series

In this subsection, we learn dynamical systems from experimental measurements. In particular, we use the ODE-RNN framework [4, 47], with the recognition model being set to different ODE-based models, to study the vibration of an airplane dataset [37]. The dataset was acquired, from time 0 to 73627, by attaching a shaker underneath the right wing to provide input signals, and 5 attributes are recorded per time stamp; these attributes include voltage of input signal, force applied to aircraft, and acceleration at 3 different spots of the airplane. We randomly take out 10% of the data to make the time series irregularly-sampled. We use the first 50% of data as our train set, the next 25% as validation set, and the rest as test set. We divide each set into non-overlapping segments of consecutive 65 time stamps of the irregularly-sampled time series, with each input instance consisting of 64 time stamps of the irregularly-sampled time series, and we aim to forecast 8 consecutive time stamps starting from the last time stamp of the segment. The input is fed through the the hybrid methods in a recurrent fashion; by changing the time duration of the last step of the ODE integration, we can forecast the output in the different time stamps. The output of the hybrid method is passed to a single dense layer to generate the output time series. In our experiments, we compare different ODE-based models hybrid with RNNs. The ODE of each model is parametrized by a 3-layer network whereas the RNN is parametrized by a simple dense network; the total number of parameters for ODE-RNN, ANODE-RNN, SONODE-RNN, HBNODE-RNN, and GHBNODE-RNN with 16, 22, 14, 15, 15 augmented dimensions are 15,986, 16,730, 16,649, 16,127, and 16,127, respectively. To avoid potential error due to the ODE solver, we use a tolerance of 10^{-7} .

In training those hybrid models, we regularize the models by penalizing the L2 distance between the RNN output and the values of the next time stamp. Due to the second-order natural of the underlying dynamics [39], ODE-RNN and ANODE-RNN learn the dynamics very poorly with much larger training and test losses than the other models even they take smaller NFEs. HBNODE-RNN and GHBNODE-RNN give better prediction than SONODE-RNN using less backward NFEs.

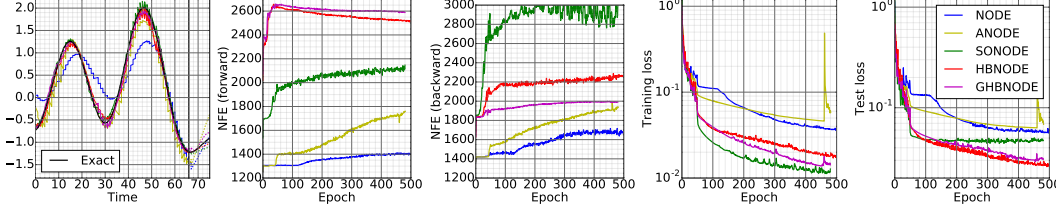


Figure 8: Contrasting ODE-RNN, ANODE-RNN, SONODE-RNN, HBNODE-RNN, and GHBNODE-RNN for learning a vibrational dynamical system. Left most: The learned curves of each model vs. the ground truth (Time: <66 for training, 66-75 for testing).

5.4 Walker2D kinematic simulation

In this subsection, we evaluate the performance of HBNODE-RNN and GHBNODE-RNN on the Walker2D kinematic simulation task, which requires learning long-term dependency effectively [28]. The dataset [3] consists of a dynamical system from kinematic simulation of a person walking from a pre-trained policy, aiming to learn the kinematic simulation of the MuJoCo physics engine [54]. The dataset is irregularly-sampled where 10% of data are removed from the simulation. Each input is consisted of 64 time stamps and fed though the the hybrid methods in a recurrent fashion, and the outputs of hybrid methods is passed to a single dense layer to generate the output time series. The target is to provide auto-regressive forecast so that the output time series is as close as the input sequence shifted 1 time stamp to the right. We compare ODE-RNN (with 7 augmentation), ANODE-RNN (with 7 ANODE style augmentation), HBNODE-RNN (with 7 augmentation), and GHBNODE-RNN (with 7 augmentation)⁶. The RNN is parametrized by a 3-layer network whereas the ODE is parametrized by a simple dense network. The number of parameters of the above four models are 8,729, 8,815, 8,899, and 8,899, respectively. In Fig. 9, we compare the performance of the above four models on the Walker2D benchmark; HBNODE-RNN and GHBNODE-RNN not only require significantly less NFEs in both training (forward and backward) and in testing than ODE-RNN and ANODE-RNN, but also have much smaller training and test losses.

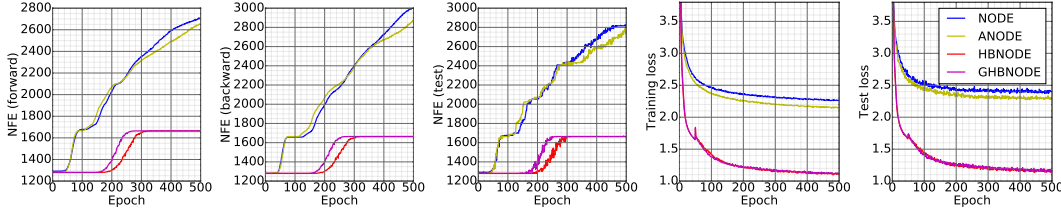


Figure 9: Contrasting ODE-RNN, ANODE-RNN, SONODE-RNN, HBNODE-RNN, and GHBNODE-RNN for the Walker-2D kinematic simulation.

6 Related Work

Reducing NFEs in training NODEs. Several techniques have been developed to reduce the NFEs for the forward solvers in NODEs, including weight decay [15], input augmentation [10], regularizing solvers and learning dynamics [12, 24, 13, 40], high-order ODE [39], data control [33], and depth-variance [33]. HBNODEs can reduce both forward and backward NFEs at the same time.

Second-order ODE accelerated dynamics. It has been noticed in both optimization and sampling communities that second-order ODEs with an appropriate damping term, e.g., the classical momentum and Nesterov’s acceleration in discrete regime, can significantly accelerate the first-order gradient dynamics (gradient descent), e.g., [43, 36, 5, 51, 59]. Also, these second-order ODEs have been discretized via some interesting numerical schemes to design fast optimization schemes, e.g., [50].

Learning long-term dependencies. Learning long-term dependency is one of the most important goals for learning from sequential data. Most of the existing works focus on mitigating exploding or vanishing gradient issues in training RNNs, e.g., [1, 60, 23, 57, 34, 19, 53]. Attention-based models are proposed for learning on sequential data concurrently with the effective accommodation

⁶Here, 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 [39]. We re-initialized SONODE 100 times; all failed due to initialisation problems.

of learning long-term dependency [56, 8]. Recently, NODEs have been integrated with long-short term memory model [20] to learn long-term dependency for irregularly-sampled time series [28]. HBNODEs directly enhance learning long-term dependency from sequential data.

Momentum in neural network design. As a line of orthogonal work, the momentum has also been studied in designing neural network architecture, e.g., [35, 53, 29, 48], which can also help accelerate training and learn long-term dependencies. These techniques can be considered as changing the neural network f in (1). We leave the synergistic integration of adding momentum to f with our work on changing the left-hand side of (1) as a future work.

ResNet-style models. Interpreting ResNet as an ODE model has been an interesting research direction which has lead to interesting neural network architectures and analysis from the numerical ODE solvers and differential equation theory viewpoints, e.g., [16, 32, 30, 58].

7 Concluding Remarks

We proposed HBNODEs to reduce the NFEs in solving both forward and backward ODEs, which also improve generalization performance over the existing benchmark models. Moreover, HBNODEs alleviate vanishing gradients in training NODEs, making HBNODEs able to learn long-term dependency effectively from sequential data. In the optimization community, Nesterov acceleration [36] is also a famous algorithm for accelerating gradient descent, that achieves an optimal convergence rate for general convex optimization problems. The ODE counterpart of the Nesterov’s acceleration corresponds to (9) with γ being replaced by a time-dependent damping parameter, e.g., $t/3$ [51]. The adjoint equation of the Nesterov’s ODE [51] is no longer a Nesterov’s ODE. We notice that directly using Nesterov’s ODE cannot improve the performance of the vanilla neural ODE. How to integrate Nesterov’s ODE with neural ODE is an interesting future direction. Another interesting direction is connecting HBNODE with symplectic ODE-net [63] through an appropriate change of variables.

8 Acknowledgement

This material is based on research sponsored by the NSF grant DMS-1924935 and DMS-1952339, the DOE grant DE-SC0021142, and the ONR grant N00014-18-1-2527 and ONR MURI grant N00014-20-1-2787.

References

- [1] Martin Arjovsky, Amar Shah, and Yoshua Bengio. Unitary evolution recurrent neural networks. In *International Conference on Machine Learning*, pages 1120–1128, 2016.
- [2] Yoshua Bengio, Patrice Simard, and Paolo Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, 5(2):157–166, 1994.
- [3] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. OpenAI Gym, 2016. cite arxiv:1606.01540.
- [4] Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary differential equations. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 6572–6583, 2018.
- [5] Tianqi Chen, Emily Fox, and Carlos Guestrin. Stochastic gradient Hamiltonian Monte Carlo. In *International conference on machine learning*, pages 1683–1691, 2014.
- [6] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- [7] Talgat Daulbaev, Alexandr Katrutsa, Larisa Markeeva, Julia Gusak, Andrzej Cichocki, and Ivan Oseledets. Interpolation technique to speed up gradients propagation in neural odes. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 16689–16700. Curran Associates, Inc., 2020.
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [9] Jianzhun Du, Joseph Futoma, and Finale Doshi-Velez. Model-based reinforcement learning for semi-markov decision processes with neural odes. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 19805–19816. Curran Associates, Inc., 2020.
- [10] Emilien Dupont, Arnaud Doucet, and Yee Whye Teh. Augmented neural odes. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [11] N. Benjamin Erichson, Omri Azencot, Alejandro Queiruga, Liam Hodgkinson, and Michael W. Mahoney. Lipschitz recurrent neural networks. In *International Conference on Learning Representations*, 2021.
- [12] Chris Finlay, Joern-Henrik Jacobsen, Levon Nurbekyan, and Adam Oberman. How to train your neural ODE: the world of Jacobian and kinetic regularization. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 3154–3164. PMLR, 13–18 Jul 2020.
- [13] Arnab Ghosh, Harkirat Behl, Emilien Dupont, Philip Torr, and Vinay Namboodiri. Steer : Simple temporal regularization for neural ode. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 14831–14843. Curran Associates, Inc., 2020.
- [14] Amílcar dos Santos Gonçalves. A Version of Beale’s Method Avoiding the Free-Variables. In *Proceedings of the 1971 26th Annual Conference*, ACM ’71, page 433–441, New York, NY, USA, 1971. Association for Computing Machinery.
- [15] Will Grathwohl, Ricky T. Q. Chen, Jesse Bettencourt, and David Duvenaud. Scalable reversible generative models with free-form continuous dynamics. In *International Conference on Learning Representations*, 2019.

- [16] Eldad Haber and Lars Ruthotto. Stable architectures for deep neural networks. *Inverse Problems*, 34(1):014004, 2017.
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015.
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *European conference on computer vision*, pages 630–645. Springer, 2016.
- [19] Kyle Helfrich, Devin Willmott, and Qiang Ye. Orthogonal recurrent neural networks with scaled Cayley transform. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1969–1978, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR.
- [20] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
- [21] Zijie Huang, Yizhou Sun, and Wei Wang. Learning continuous system dynamics from irregularly-sampled partial observations. In *Advances in Neural Information Processing Systems*, 2020.
- [22] Junteng Jia and Austin R Benson. Neural jump stochastic differential equations. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [23] Li Jing, Yichen Shen, Tena Dubcek, John Peurifoy, Scott Skirlo, Yann LeCun, Max Tegmark, and Marin Soljačić. Tunable efficient unitary neural networks (eunn) and their application to rnns. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1733–1741. JMLR. org, 2017.
- [24] Jacob Kelly, Jesse Bettencourt, Matthew J Johnson, and David K Duvenaud. Learning differential equations that are easy to solve. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 4370–4380. Curran Associates, Inc., 2020.
- [25] Patrick Kidger, James Morrill, James Foster, and Terry J. Lyons. Neural controlled differential equations for irregular time series. In *NeurIPS*, 2020.
- [26] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [27] Nikola B. Kovachki and Andrew M. Stuart. Continuous time analysis of momentum methods. *Journal of Machine Learning Research*, 22(17):1–40, 2021.
- [28] Mathias Lechner and Ramin Hasani. Learning long-term dependencies in irregularly-sampled time series. *arXiv preprint arXiv:2006.04418*, 2020.
- [29] Huan Li, Yibo Yang, Dongmin Chen, and Zhouchen Lin. Optimization algorithm inspired deep neural network structure design. *arXiv preprint arXiv:1810.01638*, 2018.
- [30] Qianxiao Li, Ting Lin, and Zuowei Shen. Deep learning via dynamical systems: An approximation perspective. *arXiv preprint arXiv:1912.10382*, 2019.
- [31] Xuechen Li, Ting-Kam Leonard Wong, Ricky T. Q. Chen, and David Duvenaud. Scalable gradients for stochastic differential equations. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 3870–3882. PMLR, 26–28 Aug 2020.
- [32] Yiping Lu, Aoxiao Zhong, Quanzheng Li, and Bin Dong. Beyond finite layer neural networks: Bridging deep architectures and numerical differential equations. In *International Conference on Machine Learning*, pages 3276–3285. PMLR, 2018.

- [33] Stefano Massaroli, Michael Poli, Jinkyoo Park, Atsushi Yamashita, and Hajime Asma. Dissecting neural odes. In *34th Conference on Neural Information Processing Systems, NeurIPS 2020*. The Neural Information Processing Systems, 2020.
- [34] Zakaria Mhammedi, Andrew Hellicar, Ashfaqur Rahman, and James Bailey. Efficient orthogonal parametrisation of recurrent neural networks using householder reflections. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2401–2409. JMLR.org, 2017.
- [35] Thomas Moreau and Joan Bruna. Understanding the learned iterative soft thresholding algorithm with matrix factorization. *arXiv preprint arXiv:1706.01338*, 2017.
- [36] Yurii E Nesterov. A method for solving the convex programming problem with convergence rate $O(1/k^2)$. In *Dokl. Akad. Nauk SSSR*, volume 269, pages 543–547, 1983.
- [37] Jean-Philippe Noël and M Schoukens. F-16 aircraft benchmark based on ground vibration test data. In *2017 Workshop on Nonlinear System Identification Benchmarks*, pages 19–23, 2017.
- [38] Alexander Norcliffe, Cristian Bodnar, Ben Day, Jacob Moss, and Pietro Liò. Neural {ode} processes. In *International Conference on Learning Representations*, 2021.
- [39] Alexander Norcliffe, Cristian Bodnar, Ben Day, Nikola Simidjievski, and Pietro Liò. On second order behaviour in augmented neural odes. In *Advances in Neural Information Processing Systems*, 2020.
- [40] Avik Pal, Yingbo Ma, Viral Shah, and Christopher V Rackauckas. Opening the blackbox: Accelerating neural differential equations by regularizing internal solver heuristics. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8325–8335. PMLR, 18–24 Jul 2021.
- [41] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In *International Conference on Machine Learning*, pages 1310–1318, 2013.
- [42] Michael Poli, Stefano Massaroli, Junyoung Park, Atsushi Yamashita, Hajime Asama, and Jinkyoo Park. Graph neural ordinary differential equations. *arXiv preprint arXiv:1911.07532*, 2019.
- [43] Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964.
- [44] Lev Semenovich Pontryagin. *Mathematical theory of optimal processes*. Routledge, 2018.
- [45] Alessio Quaglino, Marco Gallieri, Jonathan Masci, and Jan Koutník. Snode: Spectral discretization of neural odes for system identification. In *International Conference on Learning Representations*, 2020.
- [46] H. H. Rosenbrock. An Automatic Method for Finding the Greatest or Least Value of a Function. *The Computer Journal*, 3(3):175–184, 01 1960.
- [47] Yulia Rubanova, Ricky T. Q. Chen, and David K Duvenaud. Latent ordinary differential equations for irregularly-sampled time series. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [48] Michael E. Sander, Pierre Ablin, Mathieu Blondel, and Gabriel Peyré. Momentum residual neural networks. *arXiv preprint arXiv:2102.07870*, 2021.
- [49] Bin Shi, Simon S. Du, Michael I. Jordan, and Weijie J. Su. Understanding the acceleration phenomenon via high-resolution differential equations. *arXiv preprint arXiv:1810.08907*, 2018.
- [50] Bin Shi, Simon S Du, Weijie Su, and Michael I Jordan. Acceleration via symplectic discretization of high-resolution differential equations. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.

- [51] Weijie Su, Stephen Boyd, and Emmanuel Candes. A differential equation for modeling nesterov’s accelerated gradient method: Theory and insights. In *Advances in Neural Information Processing Systems*, pages 2510–2518, 2014.
- [52] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *International Conference on Machine Learning*, pages 1139–1147, 2013.
- [53] Tan M. Nguyen and Richard G. Baraniuk and Andrea L. Bertozzi and Stanley J. Osher and Bao Wang. MomentumRNN: Integrating momentum into recurrent neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 9154–9164, 2020.
- [54] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033, 2012.
- [55] Belinda Tzen and Maxim Raginsky. Neural stochastic differential equations: Deep latent gaussian models in the diffusion limit. *arXiv preprint arXiv:1905.09883*, 2019.
- [56] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [57] Eugene Vorontsov, Chiheb Trabelsi, Samuel Kadoury, and Chris Pal. On orthogonality and learning recurrent networks with long term dependencies. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3570–3578. JMLR. org, 2017.
- [58] Bao Wang, Binjie Yuan, Zuoqiang Shi, and Stanley J. Osher. EnResNet: ResNets ensemble via the Feynman–Kac formalism for adversarial defense and beyond. *SIAM Journal on Mathematics of Data Science*, 2(3):559–582, 2020.
- [59] Ashia C. Wilson, Benjamin Recht, and Michael I. Jordan. A Lyapunov Analysis of Momentum Methods in Optimization. *arXiv preprint arXiv:1611.02635*, 2018.
- [60] Scott Wisdom, Thomas Powers, John Hershey, Jonathan Le Roux, and Les Atlas. Full-capacity unitary recurrent neural networks. In *Advances in Neural Information Processing Systems*, pages 4880–4888, 2016.
- [61] Cagatay Yildiz, Markus Heinonen, and Harri Lahdesmaki. ODE2VAE: Deep generative second order ODEs with Bayesian neural networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [62] Tianjun Zhang, Zhewei Yao, Amir Gholami, Joseph E Gonzalez, Kurt Keutzer, Michael W Mahoney, and George Biros. ANODEV2: A Coupled Neural ODE Framework. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [63] Yaofeng Desmond Zhong, Biswadip Dey, and Amit Chakraborty. Symplectic ode-net: Learning hamiltonian dynamics with control. In *International Conference on Learning Representations*, 2020.
- [64] Juntang Zhuang, Nicha C Dvornek, sekhar tatikonda, and James s Duncan. {MALI}: A memory efficient and reverse accurate integrator for neural {ode}s. In *International Conference on Learning Representations*, 2021.

Checklist

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

Supplementary Material for Heavy Ball Neural Ordinary Differential Equations

A Review of the Adjoint Equation for the First- and Second-order ODEs

The adjoint sensitivity method is the key to assuring constant memory usage in training neural ODEs [4]. In this section, we present two different proofs for the first-order adjoint sensitivity equations. The differentiation proof in Appendix A.1.1 is adapted from the proof by Norcliffe et al [39]. We provide a new integral proof in Appendix A.1.2 to extend theoretical support for the Lipschitz continuous functions. We also revisit the proof of the second-order adjoint sensitivity equations by Norcliffe et al [39].

A.1 First-order Adjoint Sensitivity Equation

A Neural ODE for hidden features $\mathbf{h}(t) \in \mathbb{R}^N$ takes the form

$$\frac{\partial \mathbf{h}}{\partial t} = f(\mathbf{h}(t), t, \theta), \quad \mathbf{h}(t_0) = \mathbf{h}_{t_0}, \quad \mathbf{h}(T) = \mathbf{h}_T, \quad (23)$$

where $f(\mathbf{h}(t), t, \theta) \in \mathbb{R}^N$ is a neural network with learnable parameters θ . The corresponding adjoint equation, with \mathcal{L} being a scalar loss function, is defined by the following ODE,

$$\frac{\partial \mathbf{A}(t)}{\partial t} = -\mathbf{A}(t) \frac{\partial f}{\partial \mathbf{h}}, \quad \mathbf{A}(T) = -\mathbf{I}, \quad \mathbf{a}(t) = -\frac{d\mathcal{L}}{d\mathbf{h}_T} \mathbf{A}(t). \quad (24)$$

For gradient-based optimization, we need to compute the following derivatives

$$\frac{d\mathcal{L}}{d\theta} = \frac{d\mathcal{L}}{d\mathbf{h}_T} \frac{d\mathbf{h}_T}{d\theta}, \quad \frac{d\mathcal{L}}{d\mathbf{h}_{t_0}} = \frac{d\mathcal{L}}{d\mathbf{h}_T} \frac{d\mathbf{h}_T}{d\mathbf{h}_{t_0}}. \quad (25)$$

In the following sections, we show that

$$\frac{d\mathbf{h}_T}{d\theta} = \int_T^{t_0} \mathbf{A} \frac{\partial f}{\partial \theta} dt, \quad \frac{d\mathbf{h}_T}{d\mathbf{h}_{t_0}} = -\mathbf{A}(t_0). \quad (26)$$

By linearity, we immediately arrive at the following adjoint sensitivity equations

$$\frac{d\mathcal{L}}{d\theta} = \int_{t_0}^T \mathbf{a} \frac{\partial f}{\partial \theta} dt, \quad \frac{d\mathcal{L}}{d\mathbf{h}_{t_0}} = \mathbf{a}(t_0). \quad (27)$$

A.1.1 Proof of the First-Order Adjoint Sensitivity Equation: Differentiation Approach

We adapt the proof of the adjoint sensitivity equations from Norcliffe et al [39]. Assume that $f \in C^1$, ϕ is either θ or \mathbf{h}_{t_0} , then the following equations hold

$$\frac{\partial \mathbf{A}(t)}{\partial t} = -\mathbf{A}(t) \frac{\partial f}{\partial \mathbf{h}}, \quad \frac{\partial^2 \mathbf{h}}{\partial \phi \partial t} = \frac{\partial f}{\partial \theta} \frac{d\theta}{d\phi} + \frac{\partial f}{\partial \mathbf{h}} \frac{d\mathbf{h}}{d\phi}, \quad \frac{\partial(\mathbf{A} \frac{\partial \mathbf{h}}{\partial \phi})}{\partial t} = \frac{\partial \mathbf{A}}{\partial t} \frac{\partial \mathbf{h}}{\partial \phi} + \mathbf{A} \frac{\partial^2 \mathbf{h}}{\partial \phi \partial t}. \quad (28)$$

Combining the three equations in (28) yields the differential equation

$$\frac{\partial(\mathbf{A} \frac{\partial \mathbf{h}}{\partial \phi})}{\partial t} = \frac{\partial \mathbf{A}}{\partial t} \frac{\partial \mathbf{h}}{\partial \phi} + \mathbf{A} \frac{\partial^2 \mathbf{h}}{\partial \phi \partial t} = -\mathbf{A}(t) \frac{\partial f}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \phi} + \mathbf{A} \left(\frac{\partial f}{\partial \theta} \frac{d\theta}{d\phi} + \frac{\partial f}{\partial \mathbf{h}} \frac{d\mathbf{h}}{d\phi} \right) = \mathbf{A} \frac{\partial f}{\partial \theta} \frac{d\theta}{d\phi}. \quad (29)$$

Integrating both sides of (29) in t from T to t_0 , we arrive at the integral equation

$$\left(\mathbf{A} \frac{\partial \mathbf{h}}{\partial \phi} \right) \Big|_T^{t_0} = \int_T^{t_0} \mathbf{A} \frac{\partial f}{\partial \theta} \frac{d\theta}{d\phi} dt. \quad (30)$$

Using the conditions $\mathbf{A}(T) = -\mathbf{I}$, $\mathbf{h}(t_0) = \mathbf{h}_{t_0}$, $\mathbf{h}(T) = \mathbf{h}_T$, we rewrite the equation (30) as

$$\frac{d\mathbf{h}_T}{d\phi} = -\mathbf{A}(t_0) \frac{d\mathbf{h}_{t_0}}{d\phi} + \int_T^{t_0} \mathbf{A} \frac{\partial f}{\partial \theta} \frac{d\theta}{d\phi} dt. \quad (31)$$

Substituting $\phi = \mathbf{h}_{t_0}$ and $\phi = \theta$ respectively in (31) leads to

$$\frac{d\mathbf{h}_T}{d\mathbf{h}_{t_0}} = -\mathbf{A}(t_0), \quad \frac{d\mathbf{h}_T}{d\theta} = \int_T^{t_0} \mathbf{A} \frac{\partial f}{\partial \theta} dt. \quad (32)$$

This proof is adapted from the proof provided by Norcliffe et al [39] for general second-order neural ODEs by differentiation and this proof only holds for $f \in C^1$.

A.1.2 Proof of the First-Order Adjoint Sensitivity Equations: Integration Approach

The proof in Appendix A.1.1 requires that $f \in C^1$. However, with activation functions like ReLU, f may not be smooth enough to satisfy this requirement. Meanwhile, the adjoint equation (24) that \mathbf{A} satisfies may not have a continuous right hand side, which can fail the Picard-Lindelöf theorem that guarantees the existence and uniqueness of solutions to the adjoint equation.

To circumvent these deficiencies, we propose a new proof based on integration. Assume that $f(\mathbf{h}, t, \theta)$ is continuous in t and Lipschitz continuous in \mathbf{h}, θ , and there exists some open ball around $\mathbf{h}_{t_0} = \mathbf{s}_0$, $\theta = \theta_0$ such that for every pair of initial condition and parameters in the open ball, there exists a unique solution for $t \in [t_0, T]$. We denote the solution starting from $\mathbf{h}_{t_0} = \mathbf{s}_0$, $\theta = \theta_0$ as \mathbf{h}_0 . In order to avoid difficulties in proving the existence and uniqueness of the solution, we explicitly define the adjoint equation through the following matrix exponential

$$\mathbf{A}(t) = -\exp \left\{ -\int_T^t \frac{\partial f}{\partial \mathbf{h}}(\mathbf{h}_0(\tau), \tau, \theta_0) d\tau \right\}. \quad (33)$$

By definition, \mathbf{A} is Lipschitz continuous and satisfies the differential equation almost everywhere

$$\frac{d\mathbf{A}(t)}{dt} = -\mathbf{A}(t) \frac{\partial f}{\partial \mathbf{h}}(\mathbf{h}_0(t), t, \theta_0). \quad (34)$$

Since $\mathbf{h} \in C^1(t)$ and $\frac{d\mathbf{A}}{dt} \in L^1(t)$, we obtain the following using integration by parts,

$$\mathbf{A}\mathbf{h}|_{t_0}^T = \int_{t_0}^T \mathbf{A} \frac{\partial \mathbf{h}}{\partial t} dt + \int_{t_0}^T \frac{d\mathbf{A}}{dt} \mathbf{h} dt. \quad (35)$$

Taking partial derivatives with respect to ϕ on both sides of (35), as $\mathbf{A}(t)$ is only a function of t , we have

$$\mathbf{A} \frac{\partial \mathbf{h}}{\partial \phi} \Big|_{t_0}^T = \frac{\partial}{\partial \phi} \int_{t_0}^T \mathbf{A} \frac{\partial \mathbf{h}}{\partial t} dt + \frac{\partial}{\partial \phi} \int_{t_0}^T \frac{d\mathbf{A}}{dt} \mathbf{h} dt. \quad (36)$$

In order to exchange integral and derivatives, we use the dominated convergence theorem. Because f is Lipschitz continuous on \mathbf{h} , \mathbf{h} is Lipschitz continuous on ϕ , and thus $\frac{\partial \mathbf{h}}{\partial \phi}$ is Lebesgue integrable. Therefore, by chain rule, the following equation holds almost everywhere,

$$\frac{\partial^2 \mathbf{h}}{\partial t \partial \phi} = \frac{\partial^2 \mathbf{h}}{\partial \phi \partial t} = \frac{df}{d\phi} = \frac{\partial f}{\partial \theta} \frac{d\theta}{d\phi} + \frac{\partial f}{\partial \mathbf{h}} \frac{d\mathbf{h}}{d\phi}. \quad (37)$$

Because t is bounded, the right hand side of equation (37) is Lebesgue integrable, and so is the left hand side. Because both $\frac{\partial \mathbf{h}}{\partial \phi}$ and $\frac{\partial^2 \mathbf{h}}{\partial t \partial \phi}$ are Lebesgue integrable, by dominated convergence theorem, we have the following exchange of integrals and derivatives

$$\frac{\partial}{\partial \phi} \int_{t_0}^T \mathbf{A} \frac{\partial \mathbf{h}}{\partial t} dt = \int_{t_0}^T \mathbf{A} \frac{\partial^2 \mathbf{h}}{\partial t \partial \phi} dt, \quad \frac{\partial}{\partial \phi} \int_{t_0}^T \frac{d\mathbf{A}}{dt} \mathbf{h} dt = \int_{t_0}^T \frac{d\mathbf{A}}{dt} \frac{\partial \mathbf{h}}{\partial \phi} dt. \quad (38)$$

Combining equation (36) with (38) gives us

$$\mathbf{A} \frac{\partial \mathbf{h}}{\partial \phi} \Big|_{t_0}^T = \int_{t_0}^T \mathbf{A} \frac{\partial^2 \mathbf{h}}{\partial t \partial \phi} dt + \int_{t_0}^T \frac{d\mathbf{A}}{dt} \frac{\partial \mathbf{h}}{\partial \phi} dt. \quad (39)$$

By taking Lebesgue integral of equation (37), we have the equation

$$\int_{t_0}^T \mathbf{A} \frac{\partial^2 \mathbf{h}}{\partial t \partial \phi} dt = \int_{t_0}^T \mathbf{A} \left(\frac{\partial f}{\partial \theta} \frac{d\theta}{d\phi} + \frac{\partial f}{\partial \mathbf{h}} \frac{d\mathbf{h}}{d\phi} \right) dt. \quad (40)$$

Meanwhile, at \mathbf{h}_0 , we can integrate equation (34) to a similar form as

$$\int_{t_0}^T \frac{d\mathbf{A}}{dt} \frac{\partial \mathbf{h}}{\partial \phi} dt = - \int_{t_0}^T \mathbf{A} \frac{\partial f}{\partial \mathbf{h}} \frac{d\mathbf{h}}{d\phi} dt. \quad (41)$$

Consequently, at \mathbf{h}_0 , we can sum up equations (39), (40), and (41) and arrive at

$$\mathbf{A} \frac{\partial \mathbf{h}}{\partial \phi} \Big|_{t_0}^T = \int_{t_0}^T \mathbf{A} \frac{\partial f}{\partial \theta} \frac{d\theta}{d\phi} dt, \quad (42)$$

which is the same integral equation as equation (30) in the differentiation proof in Appendix A.1.1. Thus, plugging in the initial conditions provides us with the same result.

A.1.3 Corollary of the First-order Gradient Propagation

An immediate corollary of the above proof is that combining equations (26) and (33) results in

$$\frac{d\mathbf{h}_T}{d\mathbf{h}_{t_0}} = -\mathbf{A}(t_0) = \exp \left\{ - \int_{t_0}^T \frac{\partial f}{\partial \mathbf{h}}(\mathbf{h}_0(\tau), \tau, \theta_0) d\tau \right\}. \quad (43)$$

As (43) is true for every choice of t_0 , we can also generalize it to

$$\frac{d\mathbf{h}_T}{d\mathbf{h}_t} = \exp \left\{ - \int_t^T \frac{\partial f}{\partial \mathbf{h}}(\mathbf{h}_0(\tau), \tau, \theta_0) d\tau \right\}, \quad (44)$$

which shows the relative gradient between different times in integral.

A.2 Second-order Adjoint Sensitivity Equation

A SONODE satisfies the following equations

$$\frac{\partial \mathbf{h}}{\partial t} = \mathbf{v}, \quad \frac{\partial \mathbf{v}}{\partial t} = f(\mathbf{h}(t), \mathbf{v}(t), t, \theta), \quad \mathbf{h}(t_0) = \mathbf{h}_{t_0}, \quad \mathbf{v}(t_0) = \mathbf{v}_{t_0}, \quad (45)$$

which can be viewed as a coupled first-order ODE system of the form

$$\frac{\partial}{\partial t} \begin{bmatrix} \mathbf{h} \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} \mathbf{v} \\ f(\mathbf{h}(t), \mathbf{v}(t), t, \theta) \end{bmatrix}, \quad \begin{bmatrix} \mathbf{h} \\ \mathbf{v} \end{bmatrix} (t_0) = \begin{bmatrix} \mathbf{h}_{t_0} \\ \mathbf{v}_{t_0} \end{bmatrix}. \quad (46)$$

Denote $\mathbf{z} = \begin{bmatrix} \mathbf{h} \\ \mathbf{v} \end{bmatrix}$ and final state as

$$\begin{bmatrix} \mathbf{h}(T) \\ \mathbf{v}(T) \end{bmatrix} = \begin{bmatrix} \mathbf{h}_T \\ \mathbf{v}_T \end{bmatrix} = \mathbf{z}_T. \quad (47)$$

Using the conclusions from Appendix A.1, then the adjoint equation is given by

$$\frac{\partial \mathbf{A}(t)}{\partial t} = -\mathbf{A}(t) \begin{bmatrix} \mathbf{0} & \mathbf{I} \\ \frac{\partial f}{\partial \mathbf{h}} & \frac{\partial f}{\partial \mathbf{v}} \end{bmatrix}, \quad \mathbf{A}(T) = -\mathbf{I}, \quad \mathbf{a}(t) = -\frac{d\mathcal{L}}{d\mathbf{z}_T} \mathbf{A}(t). \quad (48)$$

By rewriting $\mathbf{A} = [\mathbf{A}_h \quad \mathbf{A}_v]$, we have the following differential equations

$$\frac{\partial \mathbf{A}_h(t)}{\partial t} = -\mathbf{A}_v(t) \frac{\partial f}{\partial \mathbf{h}}, \quad \frac{\partial \mathbf{A}_v(t)}{\partial t} = -\mathbf{A}_h(t) - \mathbf{A}_v(t) \frac{\partial f}{\partial \mathbf{v}}, \quad (49)$$

with initial conditions

$$\mathbf{A}_h(T) = -\begin{bmatrix} \mathbf{I} \\ \mathbf{0} \end{bmatrix}, \quad \mathbf{A}_v(T) = -\begin{bmatrix} \mathbf{0} \\ \mathbf{I} \end{bmatrix}, \quad (50)$$

and adjoint states

$$\mathbf{a}_h(t) = \frac{d\mathcal{L}}{d\mathbf{z}_T} \mathbf{A}_h(t), \quad \mathbf{a}_v(t) = \frac{d\mathcal{L}}{d\mathbf{z}_T} \mathbf{A}_v(t). \quad (51)$$

The gradient equations becomes

$$\frac{d\mathcal{L}}{d\theta} = \int_{t_0}^T \mathbf{a} \begin{bmatrix} \mathbf{0} \\ \frac{\partial f}{\partial \theta} \end{bmatrix} dt = \int_{t_0}^T \mathbf{a}_v \frac{\partial f}{\partial \theta} dt, \quad \frac{d\mathcal{L}}{d\mathbf{h}_{t_0}} = \mathbf{a}_h(t_0), \quad \frac{d\mathcal{L}}{d\mathbf{v}_{t_0}} = \mathbf{a}_v(t_0). \quad (52)$$

In SONODE, \mathbf{h}_{t_0} is fixed, and thus \mathbf{a}_h disappears in gradient computation. Therefore, we are only interested in \mathbf{a}_v . Thus the adjoint \mathbf{A}_v satisfies the following second-order ODE

$$\frac{\partial^2 \mathbf{A}_v(t)}{\partial t^2} = \mathbf{A}_v(t) \frac{\partial f}{\partial \mathbf{h}} - \frac{\partial(\mathbf{A}_v(t) \frac{\partial f}{\partial \mathbf{v}})}{\partial t}, \quad (53)$$

and thus

$$\frac{\partial^2 \mathbf{a}_v(t)}{\partial t^2} = \mathbf{a}_v(t) \frac{\partial f}{\partial \mathbf{h}} - \frac{\partial(\mathbf{a}_v(t) \frac{\partial f}{\partial \mathbf{v}})}{\partial t}, \quad (54)$$

with initial conditions

$$\mathbf{a}_v(T) = -\frac{d\mathcal{L}}{dz} \mathbf{A}_v(T) = \frac{d\mathcal{L}}{dv_T}, \quad \frac{\partial \mathbf{a}_v(T)}{\partial t} = -\frac{d\mathcal{L}}{d\mathbf{h}_T} - \mathbf{a}_v(T) \frac{\partial f}{\partial \mathbf{v}}(T). \quad (55)$$

This proves the second order adjoint equations for \mathbf{a}_v .

B Proof of Propositions 1, 2, and 3

B.1 Proof of Adjoint Equation for HBNODE (Propositions 1)

As HBNODE takes the form

$$\frac{d^2 \mathbf{h}(t)}{dt^2} + \gamma \frac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}(t), t, \theta), \quad (56)$$

which can also be viewed as a SONODE. By applying the adjoint equation of SONODE (54), we arrive at

$$\frac{\partial^2 \mathbf{a}(t)}{\partial t^2} = \mathbf{a}(t) \frac{\partial f}{\partial \mathbf{h}} + \gamma \frac{\partial \mathbf{a}(t)}{\partial t}. \quad (57)$$

As HBNODE only carries its state \mathbf{h} to the loss \mathcal{L} , we have $\frac{d\mathcal{L}}{dv_T} = 0$, and thus the initial conditions in equation (55) become

$$\mathbf{a}(T) = \mathbf{0}, \quad \frac{\partial \mathbf{a}(T)}{\partial t} = -\frac{d\mathcal{L}}{d\mathbf{h}_T}. \quad (58)$$

B.2 Proof of Adjoint Equation for First-order HBNODE (Proposition 2)

The coupled form of HBNODE is a coupled first-order ODE system of the form

$$\frac{\partial}{\partial t} \begin{bmatrix} \mathbf{h} \\ \mathbf{m} \end{bmatrix} = \begin{bmatrix} \mathbf{m} \\ -\gamma \mathbf{m} + f(\mathbf{h}(t), t, \theta) \end{bmatrix}, \quad \begin{bmatrix} \mathbf{h} \\ \mathbf{m} \end{bmatrix} (t_0) = \begin{bmatrix} \mathbf{h}_{t_0} \\ \mathbf{m}_{t_0} \end{bmatrix}. \quad (59)$$

Denote the final state as

$$\begin{bmatrix} \mathbf{h}(T) \\ \mathbf{m}(T) \end{bmatrix} = \begin{bmatrix} \mathbf{h}_T \\ \mathbf{m}_T \end{bmatrix} = \mathbf{z}. \quad (60)$$

Using the conclusions from Appendix A.1, we have the adjoint equation

$$\frac{\partial \mathbf{A}(t)}{\partial t} = -\mathbf{A}(t) \begin{bmatrix} \mathbf{0} & \mathbf{I} \\ \frac{\partial f}{\partial \mathbf{h}} & -\gamma \mathbf{I} \end{bmatrix}, \quad \mathbf{A}(T) = -\mathbf{I}, \quad \mathbf{a}(t) = -\frac{d\mathcal{L}}{dz} \mathbf{A}(t). \quad (61)$$

Let $[\mathbf{a}_h \quad \mathbf{a}_m] = \mathbf{a}$, by linearity we have

$$\frac{\partial [\mathbf{a}_h \quad \mathbf{a}_m]}{\partial t} = -[\mathbf{a}_h \quad \mathbf{a}_m] \begin{bmatrix} \mathbf{0} & \mathbf{I} \\ \frac{\partial f}{\partial \mathbf{h}} & -\gamma \mathbf{I} \end{bmatrix}, \quad [\mathbf{a}_h(T) \quad \mathbf{a}_m(T)] = \begin{bmatrix} \frac{d\mathcal{L}}{d\mathbf{h}_T} & \frac{d\mathcal{L}}{d\mathbf{m}_T} \end{bmatrix}, \quad (62)$$

which gives us the initial conditions at $t = T$, and the simplified first-order ODE system

$$\frac{\partial \mathbf{a}_h}{\partial t} = -\mathbf{a}_m \frac{\partial f}{\partial \mathbf{h}}, \quad \frac{\partial \mathbf{a}_m}{\partial t} = -\mathbf{a}_h + \gamma \mathbf{a}_m. \quad (63)$$

B.3 Proof of Adjoint Equation for GHBNODE (Proposition 3)

The coupled form of GHBNODE is a first-order ODE system of the form

$$\frac{\partial}{\partial t} \begin{bmatrix} \mathbf{h} \\ \mathbf{m} \end{bmatrix} = \begin{bmatrix} \sigma(\mathbf{m}) \\ -\gamma\mathbf{m} + f(\mathbf{h}(t), t, \theta) - \xi\mathbf{h}(t) \end{bmatrix}, \quad \begin{bmatrix} \mathbf{h} \\ \mathbf{m} \end{bmatrix}(t_0) = \begin{bmatrix} \mathbf{h}_{t_0} \\ \mathbf{m}_{t_0} \end{bmatrix}. \quad (64)$$

Denote the final state as

$$\begin{bmatrix} \mathbf{h}(T) \\ \mathbf{m}(T) \end{bmatrix} = \begin{bmatrix} \mathbf{h}_T \\ \mathbf{m}_T \end{bmatrix} = \mathbf{z}_T. \quad (65)$$

Using the conclusions from Appendix A.1, we have the adjoint equation

$$\frac{\partial \mathbf{A}(t)}{\partial t} = -\mathbf{A}(t) \begin{bmatrix} \mathbf{0} & \sigma'(\mathbf{m}) \\ \frac{\partial f}{\partial \mathbf{h}} - \xi \mathbf{I} & -\gamma \mathbf{I} \end{bmatrix}, \quad \mathbf{A}(T) = -\mathbf{I}, \quad \mathbf{a}(t) = -\frac{d\mathcal{L}}{d\mathbf{z}_T} \mathbf{A}(t). \quad (66)$$

Let $[\mathbf{a}_h \quad \mathbf{a}_m] = \mathbf{a}$, by linearity we have

$$\frac{\partial [\mathbf{a}_h \quad \mathbf{a}_m]}{\partial t} = -[\mathbf{a}_h \quad \mathbf{a}_m] \begin{bmatrix} \mathbf{0} & \sigma'(\mathbf{m}) \\ \frac{\partial f}{\partial \mathbf{h}} - \xi \mathbf{I} & -\gamma \mathbf{I} \end{bmatrix}, \quad [\mathbf{a}_h(T) \quad \mathbf{a}_m(T)] = \left[\frac{d\mathcal{L}}{d\mathbf{h}_T} \quad \frac{d\mathcal{L}}{d\mathbf{m}_T} \right], \quad (67)$$

which gives us the initial conditions at $t = T$, and the simplified first-order ODE system

$$\frac{\partial \mathbf{a}_h}{\partial t} = -\mathbf{a}_m \left(\frac{\partial f}{\partial \mathbf{h}} - \xi \mathbf{I} \right), \quad \frac{\partial \mathbf{a}_m}{\partial t} = -\mathbf{a}_h \sigma'(\mathbf{m}) + \gamma \mathbf{a}_m. \quad (68)$$

C Vanishing and Exploding Gradients in Training RNNs

Recurrent cells are the building blocks of RNNs. A recurrent cell can be mathematically written as

$$\mathbf{h}_t = \sigma(\mathbf{U}\mathbf{h}_{t-1} + \mathbf{W}\mathbf{x}_t + \mathbf{b}), \quad \mathbf{x}_t \in \mathbb{R}^d, \text{ for } t = 1, 2, \dots, T, \quad (69)$$

where $\mathbf{h}_t \in \mathbb{R}^h$ is the hidden state, $\mathbf{U} \in \mathbb{R}^{h \times h}$, $\mathbf{W} \in \mathbb{R}^{h \times d}$, and $\mathbf{b} \in \mathbb{R}^h$ are trainable parameters; $\sigma(\cdot)$ is a nonlinear activation function, e.g., sigmoid. Backpropagation through time is a popular algorithm for training RNNs, which usually results in exploding or vanishing gradients [2]. Thus RNNs may fail to learn long term dependencies. As an illustration, let \mathbf{h}_T and \mathbf{h}_t be the state vectors at the timestamps T and t ($T \gg t$), respectively. Assume \mathcal{L} is the loss to minimize, then

$$\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}_T} \cdot \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_t} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}_T} \cdot \prod_{k=t}^{T-1} \frac{\partial \mathbf{h}_{k+1}}{\partial \mathbf{h}_k} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}_T} \cdot \prod_{k=t}^{T-1} (\mathbf{D}_k \mathbf{U}^\top), \quad (70)$$

where $\mathbf{D}_k = \text{diag}(\sigma'(\mathbf{U}\mathbf{h}_k + \mathbf{W}\mathbf{x}_{k+1}))$ is a diagonal matrix with $\sigma'(\mathbf{U}\mathbf{h}_k + \mathbf{W}\mathbf{x}_{k+1})$ being its diagonal entries. $\|\prod_{k=t}^{T-1} (\mathbf{D}_k \mathbf{U}^\top)\|_2$ tends to either vanish or explode [2].

When applying RNNs to sequence applications with $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_T)$ be an input sequence of length T and $\mathbf{y} = (y_1, \dots, y_T)$ be the sequence of labels, we let \mathcal{L}_t be the loss at the timestamp t and the total loss on the whole sequence be

$$\mathcal{L} = \sum_{t=1}^T \mathcal{L}_t, \quad (71)$$

the vanishing or exploding issue can be shown following (70).

For neural ODEs, note that the adjoint state $\mathbf{a}(t)$ is defined as $\partial \mathcal{L} / \partial \mathbf{h}(t)$, which also tends to explode or vanish during training.

C.1 Derivation of equation (18)

GHBNODE can be viewed as a system of higher dimensional NODE as in equation (64). With $z = \begin{bmatrix} \mathbf{h} \\ \mathbf{m} \end{bmatrix}$, z satisfies NODE equation, and therefore it also satisfies the equation for relative gradient information as in (44),

$$\frac{dz_T}{dz_t} = \exp \left\{ - \int_T^t \frac{\partial f}{\partial z}(z_0(\tau), \tau, \theta_0) d\tau \right\}. \quad (72)$$

By definition of multivariate derivatives, we have

$$\frac{dz_T}{dz_t} = \begin{bmatrix} \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_t} & \frac{\partial \mathbf{h}_T}{\partial \mathbf{m}_t} \\ \frac{\partial \mathbf{m}_T}{\partial \mathbf{h}_t} & \frac{\partial \mathbf{m}_T}{\partial \mathbf{m}_t} \end{bmatrix}, \quad (73)$$

and

$$\frac{\partial f}{\partial z} = \begin{bmatrix} \mathbf{0} & \frac{\partial \sigma}{\partial \mathbf{m}} \\ \frac{\partial f}{\partial \mathbf{h}} - \xi \mathbf{I} & -\gamma \mathbf{I} \end{bmatrix}. \quad (74)$$

With equations (73) and (74), we can rewrite equation (72) in terms of \mathbf{h} and \mathbf{m} as

$$\begin{bmatrix} \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_t} & \frac{\partial \mathbf{h}_T}{\partial \mathbf{m}_t} \\ \frac{\partial \mathbf{m}_T}{\partial \mathbf{h}_t} & \frac{\partial \mathbf{m}_T}{\partial \mathbf{m}_t} \end{bmatrix} = \exp \left\{ - \int_T^t \begin{bmatrix} \mathbf{0} & \frac{\partial \sigma}{\partial \mathbf{m}} \\ \frac{\partial f}{\partial \mathbf{h}} - \xi \mathbf{I} & -\gamma \mathbf{I} \end{bmatrix} ds \right\}. \quad (75)$$

In particular, since HBNODEs are GHBNODEs with $\xi = 0$ and σ being the identity map, the gradient equation of HBNODEs takes the form

$$\begin{bmatrix} \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_t} & \frac{\partial \mathbf{h}_T}{\partial \mathbf{m}_t} \\ \frac{\partial \mathbf{m}_T}{\partial \mathbf{h}_t} & \frac{\partial \mathbf{m}_T}{\partial \mathbf{m}_t} \end{bmatrix} = \exp \left\{ - \int_T^t \begin{bmatrix} \mathbf{0} & \mathbf{I} \\ \frac{\partial f}{\partial \mathbf{h}} & -\gamma \mathbf{I} \end{bmatrix} ds \right\}. \quad (76)$$

This concludes the derivation of equation (18).

D Proof of Proposition 4

Proof. Let $\mathbf{F} = \frac{1}{t-T} \int_T^t \frac{\partial f}{\partial \mathbf{h}}(\mathbf{h}(s), s, \theta) ds - \xi \mathbf{I}$, $\mathbf{J} = \frac{1}{t-T} \int_T^t \frac{\partial \sigma}{\partial \mathbf{m}}(\mathbf{m}(s)) ds$, and $\mathbf{H} = \frac{1}{t-T} \mathbf{M}$, then we have the following equation

$$\mathbf{H} = \frac{1}{t-T} \mathbf{M} = \begin{bmatrix} \mathbf{0} & \mathbf{J} \\ \mathbf{F} & -\gamma \mathbf{I} \end{bmatrix}. \quad (77)$$

As $(\lambda + \gamma)\mathbf{I}$ commutes with any matrix \mathbf{F} , the characteristics polynomials of \mathbf{H} and \mathbf{JF} satisfy the relation

$$ch_{\mathbf{H}}(\lambda) = \det(\lambda \mathbf{I} - \mathbf{H}) = \det \begin{bmatrix} \lambda \mathbf{I} & -\mathbf{J} \\ -\mathbf{F} & (\lambda + \gamma) \mathbf{I} \end{bmatrix} = \det(\lambda(\lambda + \gamma) \mathbf{I} - \mathbf{JF}) = -ch_{\mathbf{JF}}(\lambda(\lambda + \gamma)). \quad (78)$$

Since the characteristics polynomial of \mathbf{JF} splits in the field \mathbb{C} of complex numbers, i.e. $ch_{\mathbf{JF}}(x) = \prod_{i=1}^n (x - \lambda_{\mathbf{JF},i})$, we have

$$ch_{\mathbf{H}}(\lambda) = -ch_{\mathbf{JF}}(\lambda(\lambda + \gamma)) = - \prod_{i=1}^n (\lambda(\lambda + \gamma) - \lambda_{\mathbf{JF},i}). \quad (79)$$

Therefore, the eigenvalues of \mathbf{H} appear in n pairs with each pair satisfying the quadratic equation

$$\lambda(\lambda + \gamma) - \lambda_{\mathbf{JF},i} = 0. \quad (80)$$

By Vieta's formulas, the sum of these pairs are all $-\gamma$. Therefore, the eigenvalues of \mathbf{M} comes in n pairs and the sum of each pair is $-(t - T)\gamma$. \square

E Experimental details

We first list some common settings below:

- NODE and ANODE do not have initial layers.
- For SONODE $n^* = 2n$, and for other ones $n^* = n$.
- Hyper parameters are listed in Table 3
- HTanh: HardTanh(-5, 5)
- LReLU: LeakyReLU(0.3)
- tpad: Padding with time t within ODE. i.e., transform the shape $c \times x \times y$ to $(c+1) \times x \times y$ by concatenating with a tensor of shape $1 \times x \times y$ filled with all t .
- For all tasks, we use learnable γ with $\epsilon = 1$ for both HBNODE and GHBNODE, and learnable ξ .
- fc_n : a fully connected layer with output dimension to be n .

Table 3: The hyper-parameters for each models.

Model	NODE	ANODE	SONODE	HBNODE	GHBNODE
n (Initialization)	1	2	1	1	1
h (Initialization)	22	22	22	22	22
n (Point Cloud)	2	3	2	2	2
h (Point Cloud)	20	20	13	14	14
n (MNIST)	1	6	5	5	6
h (MNIST)	92	64	50	50	45
n (CIFAR)	3	13	12	12	12
h (CIFAR)	125	64	50	51	51

Table 4: The hyper-parameters for ODE-RNN integration models.

Model	ODE-RNN	ANODE-RNN	SONODE-RNN	HBNODE-RNN	GHBNODE-RNN
d	1	1	2	2	2
n (Plane Vibration)	21	27	19	20	20
h_1 (Plane Vibration)	63	83	19	20	20
h_2 (Plane Vibration)	84	108	19	20	20
n (Walker 2D)	24	24	23	24	24
h_1 (Walker 2D)	72	72	46	48	48
h_2 (Walker 2D)	48	48	46	48	48

E.1 Network architecture used in Section 3 Initialization Test

- ODE : $\text{input}_{n^*+1} \rightarrow \text{fc}_n$

E.2 Experimental details for 5.1

- Initial Velocity : $\text{input}_2 \rightarrow \text{fc}_h \rightarrow \text{HTanh} \rightarrow \text{fc}_h \rightarrow \text{HTanh} \rightarrow \text{fc}_n$
- ODE : $\text{input}_{n^*} \rightarrow \text{fc}_h \rightarrow \text{ELU} \rightarrow \text{fc}_h \rightarrow \text{ELU} \rightarrow \text{fc}_n$
- Output : $\text{input}_n \rightarrow \text{fc}_1 \rightarrow \text{Tanh}$

E.3 Experimental details for 5.2

E.3.1 MNIST

- Initial Velocity : $\text{input}_{1 \times 28 \times 28} \rightarrow \text{conv}_{h,1} \rightarrow \text{LReLU} \rightarrow \text{conv}_{h,3} \rightarrow \text{LReLU} \rightarrow \text{conv}_{2n-1,1}$
- ODE : $\text{input}_{n^* \times 28 \times 28} \rightarrow \text{tpad} \rightarrow \text{conv}_{h,1} \rightarrow \text{ReLU} \rightarrow \text{tpad} \rightarrow \text{conv}_{h,3} \rightarrow \text{ReLU} \rightarrow \text{tpad} \rightarrow \text{conv}_{n,1}$
- Output : $\text{input}_{n \times 28 \times 28} \rightarrow \text{fc}_{10}$

E.3.2 CIFAR

- Initial Velocity : $\text{input}_{3 \times 28 \times 28} \rightarrow \text{conv}_{h,1} \rightarrow \text{LReLU} \rightarrow \text{conv}_{h,3} \rightarrow \text{LReLU} \rightarrow \text{conv}_{2n-3,1}$
- ODE : $\text{input}_{n^* \times 32 \times 32} \rightarrow \text{tpad} \rightarrow \text{conv}_{h,1} \rightarrow \text{ReLU} \rightarrow \text{tpad} \rightarrow \text{conv}_{h,3} \rightarrow \text{ReLU} \rightarrow \text{tpad} \rightarrow \text{conv}_{n,1}$
- Output : $\text{input}_{n \times 32 \times 32} \rightarrow \text{fc}_{10}$

E.4 Experimental details for 5.3

- ODE : $\text{input}_{n^*} \rightarrow \text{fc}_{h_1} \rightarrow \text{ReLU} \rightarrow \text{fc}_{h_2} \rightarrow \text{ReLU} \rightarrow \text{fc}_n$
- RNN : $\text{input}_{dn+k} \rightarrow \text{fc}_{dn}$
- Output : $\text{input}_n \rightarrow \text{fc}_5$

E.5 Experimental details for 5.4

- ODE : $\text{input}_{n^*} \rightarrow \text{fc}_n$
- RNN : $\text{input}_{dn+k} \rightarrow \text{fc}_{h_1} \rightarrow \text{Tanh} \rightarrow \text{fc}_{h_2} \rightarrow \text{Tanh} \rightarrow \text{fc}_{dn}$
- Output : $\text{input}_n \rightarrow \text{fc}_{17}$

E.6 Experimental details for ODE vs. HBNODE on benchmarks

To numerically show that the HBNODE (7) converges faster to the stationary point than the ODE limit of gradient descent (8), we apply the Dormand–Prince-45 ODE solver, which is the default solver for NODEs, to solve both ODEs. We set $F(\mathbf{x})$ to be two celebrated functions used in the optimization community, namely, the Rosenbrock and the Beale functions.

Rosenbrock function. The Rosenbrock function is given by

$$F(\mathbf{x}) := F(x, y) = 100(y - x^2)^2 + (1 - x)^2,$$

which has the minimum $(x, y) = (1, 1)$. Starting from $(0, 0)$, we apply Dormand–Prince-45 solver using a step size $\Delta t = 0.001$ to solve both ODEs (7) and (8) for t from 0 to 1. For the HBODE, we set $\gamma = 0.9$ and set the initial value of $d\mathbf{x}/dt = (0, 0)$.

Beale function. The Beale function is given by

$$F(x, y) = (1.5 - x + xy)^2 + (2.25 - x + xy^2)^2 + (2.625 - x + xy^3)^2$$

which has the minimum $(x, y) = (3, 0.5)$. Starting from $(0, 0)$, we apply Dormand–Prince-45 solver using a step size 0.01 to solve both ODEs in (7) and (8) for t from 0 to 2. For the HBODE, we set $\gamma = 0.7$ and set the initial value of $d\mathbf{x}/dt = (0, 0)$.