
Improving Neural Ordinary Differential Equations with Nesterov’s Accelerated Gradient Method

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We propose the Nesterov neural ordinary differential equations (NesterovNODEs),
2 whose layers solve the second-order ordinary differential equations (ODEs) limit
3 of Nesterov’s accelerated gradient (NAG) method, and a generalization called
4 GNesterovNODEs. Taking the advantage of the convergence rate $\mathcal{O}(1/k^2)$ of the
5 NAG scheme, GNesterovNODEs speed up training and inference by reducing the
6 number of function evaluations (NFEs) needed to solve the ODEs. We also prove
7 that the adjoint state of a GNesterovNODEs also satisfies a GNesterovNODEs,
8 thus accelerating both forward and backward ODE solvers and allowing the model
9 to be scaled up for large-scale tasks. We empirically corroborate the advantage of
10 GNesterovNODEs on a wide range of practical applications, including point cloud
11 separation, image classification, and sequence modeling. Compared to NODEs,
12 GNesterovNODEs require a significantly smaller number of NFEs while achieving
13 better accuracy across our experiments.

14 1 Introduction

15 Dynamical systems have been recently integrated into deep neural networks for modeling high-
16 dimensional data. The advantage of this approach is that well-developed mathematical modeling
17 techniques from dynamical systems can be employed to improve neural networks. Along this research
18 direction, the correspondence between residual networks, a popular class of neural networks with skip
19 connections, and the numerical solution of ordinary differential equations (ODEs) have been vastly
20 studied in [Haber & Ruthotto, 2017; Weinan, 2017; Ruthotto & Haber, 2019; Chen et al., 2018]. The
21 resulting Neural ODEs (NODEs) model when taking the the discretization step to zero have shown
22 great promises in a wide range of applications including scientific discovery [Köhler et al., 2020;
23 Zhong et al., 2020], irregular time series modeling [Rubanova et al., 2019; De Brouwer et al., 2019],
24 mean-field games [Ruthotto et al., 2020], and generative modeling [Grathwohl et al., 2019b; Yildiz
25 et al., 2019]. NODEs model the dynamics of hidden state $\mathbf{h}(t) \in \mathbb{R}^N$ in a neural network by an ODE

$$\frac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}(t), t, \theta), \quad \mathbf{h}(0) = \mathbf{h}(t_0), \quad (1)$$

26 where function f captures the dynamics and is chosen to be a neural network with parameters θ that
27 are learned from the data. Starting from the input $\mathbf{h}(t_0)$ at the initial time t_0 , NODEs compute the
28 output $\mathbf{h}(T)$ at time T by solving the Initial Value Problem in Eq. (1) for some time $T \geq t_0$. NODEs
29 are trained by optimizing the loss $L(\mathbf{h}(T))$ between the prediction $\mathbf{h}(T)$ and the ground truth where
30 the parameters θ are updated using the following gradient [Pontryagin, 1987]

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

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

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

33 NODEs solve both the ODE Eq. (1) in its forward pass and the ODEs (2) and (3) in its backward
 34 pass using black-box numerical ODE solvers. The literature refers to this approach as either the
 35 continuous adjoint method or optimise-then-discretise [Kidger, 2022]. The number of function
 36 evaluations (NFEs) that these solvers need in a single forward and backward pass is among the main
 37 factors that decide the computational efficiency of the model, i.e. how fast the model is. Unfortunately,
 38 in many applications, NODEs require high NFEs in both training and inference, especially when the
 39 error tolerances of the solvers are set to small values for obtaining high accuracy. Furthermore, the
 40 NFEs increase rapidly when training progresses. High NFEs deteriorate the efficiency of NODEs,
 41 reduce the accuracy of the trained model, and results in instability during training, making it difficult
 42 to scale up the models to large-scale tasks [Grathwohl et al., 2019a; Dupont et al., 2019b; Massaroli
 43 et al., 2020; Norcliffe et al., 2020; Finlay et al., 2020].

44 1.1 Contribution

45 We propose the Nesterov Neural ODEs (NesterovNODEs) that leverage the continuous limit of
 46 the Nesterov’s accelerated gradient (NAG) descent and the convergence rate $\mathcal{O}(1/k^2)$ of the NAG
 47 scheme to enhance NODE training and inference. Our contributions are four-fold:

- 48 1. We formulate the NesterovNODE that solves Nesterov ODEs, i.e. second-order ODEs with
 49 a time-dependent damping term, instead of first-order ODEs (1). To improve computational
 50 efficiency of the model, we convert these second-order ODEs into equivalent systems of
 51 first-order differential-algebraic equations that are solved in both forward and backward
 52 propagations of the NesterovNODE.
- 53 2. To eliminate the potential blow-up problem in training NesterovNODEs, we further develop
 54 the Generalized NesterovNODEs (GNesterovNODEs) by introducing skip connections [He
 55 et al., 2016] and gating mechanisms [Hochreiter & Schmidhuber, 1997] into NesterovN-
 56 ODEs. In general, GNesterovNODEs form a wide class of neural differential equations
 57 which are represented by differential-algebraic systems and contain NesterovNODEs as a
 58 subclass.
- 59 3. We prove that the continuous adjoint equation used to compute the gradients for updating
 60 the parameters θ in a GNesterovNODE also follows a generalized Nesterov ODE. Thus, the
 61 NFEs in both forward and backward passes of GNesterovNODEs are significantly reduced,
 62 especially when the solvers are used with small error tolerances.
- 63 4. We prove that the spectrum of the GNesterovNODE is well-structured. This property of
 64 GNesterovNODEs helps alleviate the vanishing gradient issue during training and allows
 65 the model to capture long-term dependencies in the data.

66 We empirically demonstrate the advantages of the NesterovNODEs/GNesterovNODEs over the
 67 baseline NODE and the state-of-the-art neural ODE models including the Heavy Ball NODEs
 68 (HBNODEs), which solve the continuous limit of the heavy ball momentum accelerated gradient
 69 descent [Xia et al., 2021] on a wide range of applications including point cloud separation, image
 70 classification, and kinetic simulation. In all experiments, our proposed models achieve better accuracy
 71 and smaller NFEs than the baselines.

72 1.2 Organization

73 We structure this paper as follows: In Section 2, we review HBNODEs and NesterovODEs. In Section
 74 3, we present the algorithm and analysis of the NesterovNODEs and GNesterovNODEs. We study
 75 the spectrum structure of the adjoint equations of NesterovNODEs/GNesterovNODEs to show that
 76 NesterovNODEs/GNesterovNODEs can learn long-term dependencies effectively in Section 4. In
 77 Section 5 and 6, we empirically validate the advantages of NesterovNODEs/GNesterovNODEs and
 78 analyze our models with ablation studies. We discuss related works in Section 7. The paper ends
 79 with concluding remarks. Proofs and additional experimental details are provided in the Appendix.

80 2 An Integration of Nesterov ODEs into NODEs

81 We first establish a connection between NODEs and gradient descent (GD), then review the Heavy
 82 Ball Neural ODEs (HBNODEs), and motivate the integration of Nesterov ODEs into NODEs.

83 **ODE limit of gradient descent and connections to NODEs** Gradient descent (GD) has been among
 84 the methods of choice in optimization and machine learning for training complex systems. Starting
 85 from initial point $\mathbf{x}_0 \in \mathbb{R}^d$, GD iterates as $\mathbf{x}_k = \mathbf{x}_{k-1} - s \nabla F(\mathbf{x}_k)$ with $s > 0$ being the step size in order
 86 to find a minimum of the function $F(\mathbf{x})$. Let $s \rightarrow 0$, we obtain the following ODE limit of the GD

$$\frac{d\mathbf{x}}{dt} = -\nabla F(\mathbf{x}_t). \quad (4)$$

88 Comparing Eq. (1) and (4), we observe that a NODE
89 solves the ODE limit of the GD where the gradi-
90 ent $-\nabla F(\mathbf{x}_t)$ is parameterized by a neural network
91 $f(\mathbf{x}(t), t, \theta)$.

92 **Heavy ball neural ordinary differential equations**
93 HBNODEs are proposed in [Xia et al., 2021]. This model
94 takes advantage of the acceleration of heavy ball (HB)
95 momentum [Polyak, 1964] to reduce the NFEs needed
96 for solving the ODEs and speed up NODEs. In particu-
97 lar, HBNODEs replace the first-order ODE limit of GD
98 by the following second-order ODE limit of heavy ball
99 momentum method:

$$\frac{d^2\mathbf{x}(t)}{dt^2} + \gamma \frac{d\mathbf{x}(t)}{dt} = -\nabla F(\mathbf{x}(t)). \quad (5)$$

100 Similar to NODEs, [Xia et al., 2021] parameterize $-\nabla F(\mathbf{x}(t))$ by a neural network $f(\mathbf{x}(t), t, \theta)$ and
101 formulate HBNODEs as follows

$$\frac{d^2\mathbf{x}(t)}{dt^2} + \gamma \frac{d\mathbf{x}(t)}{dt} = f(\mathbf{x}(t), t, \theta), \quad (6)$$

102 where $\gamma > 0$ is the damping parameter, which can be a hyperparameter or a learnable parameter.

103 **Nesterov’s accelerated gradient (NAG) momentum** Even though HB improves the convergence
104 and accelerates GD, both GD and HB share the same convergence rate of $O(1/k)$ for convex smooth
105 optimization. A breakthrough due to Nesterov [Nesterov, 1983] replaces the constant momentum γ
106 with $(k-1)/(k+2)$, a.k.a. NAG momentum, improves the convergence rate to $O(1/k^2)$, which is
107 proved to be optimal for convex and smooth objective functions [Nesterov, 1983; Su et al., 2014].
108 We demonstrate the faster convergence of NAG in comparison with GD and HB on a quadratic
109 optimization problem in Figure 1. The much faster convergence rate of NAG than that of GD and
110 HB motivates us to incorporate the second-order ODE limit of NAG into a NODE and propose the
111 NesterovNODE. Nesterov acceleration gradient method [Nesterov, 1983] takes the following form:
112 given initial points $x_0 \in \mathbb{R}^N$ and $y_0 = x_0$, the sequence $\{(x_k, y_k)\}_k$ is defined inductively as:

$$\begin{cases} x_k = y_{k-1} - s \nabla F(y_{k-1}), \\ y_k = x_k + \frac{k-1}{k+2} (x_k - x_{k-1}). \end{cases} \quad (7)$$

113 The continuous limit of the Nesterov scheme is obtained by setting $x_k = \mathbf{h}(k\sqrt{s}) = \mathbf{h}(t)$ with
114 $t = k\sqrt{s}$ and some smooth function \mathbf{h} from \mathbb{R} to \mathbb{R}^N . According to [Su et al., 2014], the function \mathbf{h}
115 satisfies the Nesterov ODE

$$\mathbf{h}''(t) + \frac{3}{t} \mathbf{h}'(t) + \nabla F(\mathbf{h}(t)) = 0, \quad (8)$$

116 with the initial conditions $\mathbf{h}(0) = \mathbf{h}_0$, $\mathbf{h}'(0) = 0$.

117 **Remark 1** (Nesterov factor). *The constant 3 in the coefficient of $\mathbf{h}'(t)$ in Eq. (8) is originally from the*
118 *approximation $(k-1)/(k+2) = 1 - 3/k + O(1/k^2)$. This constant will be replaced by a constant*
119 *r if the factor $(k-1)/(k+2)$ in Eq. (7) is replaced by $(k-1)/(k+r-1)$. It is proved in [Su et al.,*
120 *2014] that the Nesterov ODE still holds the quadratic convergence rate when 3 is replaced by any*
121 *number $r > 3$.*

122 **Remark 2** (Numerical stability). *If the Euler method is used, then to keep the numerical solution*
123 *close to the exact solution, the step size chosen in the Euler method must be small enough. However,*
124 *the smaller the step size, the more expensive the computation. The maximum stable step size, which is*
125 *the maximum value for which the step size can be chosen so that the numerical solution remains close*
126 *to the exact solution, reflects the numerical stability of the ODEs. It is proved in [Su et al., 2014]*
127 *that the maximum stable step size of the Nesterov ODE is much larger than that of the ODE (4), thus*
128 *showing the numerical stability advantage of the Nesterov ODE.*

129 3 Generalize Nesterov ODEs to Differential-Algebraic Systems

130 One can parameterize $\nabla F(\mathbf{h}(t))$ in Eq. (8) by a neural network $f(\mathbf{h}(t), t, \theta)$ with learnable parameters
131 θ in a similar way as NODE. This results in the following Nesterov Neural ODE (NesterovNODE)

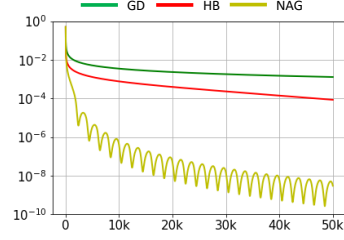


Figure 1: Comparing the convergence of GD, HB and NAG for solving the optimization problem $\min_{\mathbf{x}} F(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{L} \mathbf{x} - \mathbf{x}^\top \mathbf{b}$ where $\mathbf{L} \in \mathbb{R}^{d \times d}$ is the Laplacian of a cycle graph and \mathbf{b} is a d -dimensional vector whose first entry is 1 and all the other entries are 0.

132

$$\mathbf{h}''(t) + \frac{3}{t}\mathbf{h}'(t) + f(\mathbf{h}(t), t, \theta) = 0, \quad (9)$$

133 This second-order NesterovNODE can be written in term
134 of the first-order NesterovNODE as

$$\begin{cases} \mathbf{h}'(t) = \mathbf{m}(t), \\ \mathbf{m}'(t) = -\frac{3}{t}\mathbf{m}(t) - f(\mathbf{h}(t), t, \theta). \end{cases} \quad (10)$$

135 However, because of the singularity created by the coef-
136 ficient $\frac{3}{t}$, the training process based directly on Eq. (9)
137 and (10) will be unstable. To avoid the instability issue,
138 we set $\mathbf{h}(t) = k(t)\mathbf{x}(t)$ with $k(t) = t^{-\frac{3}{2}}e^{\frac{t}{2}}$. Then Eq. (8)
139 becomes

$$k(t)\mathbf{x}''(t) + \left(2k'(t) + \frac{3}{t}k(t)\right)\mathbf{x}'(t) + \left(k''(t) + \frac{3}{t}k'(t)\right)\mathbf{x}(t) + \nabla F(\mathbf{h}(t)) = 0. \quad (11)$$

140 We observe that $2k'(t) + \frac{3}{t}k(t) = k(t)$. By dividing both sides of Eq. (11) by $k(t)$, we obtain:

$$\mathbf{x}''(t) + \mathbf{x}'(t) + f(\mathbf{h}(t), t) = 0, \quad (12)$$

141 where

$$f(\mathbf{h}(t), t) = \frac{1}{4}(t^2 - 3)t^{-\frac{1}{2}}e^{-\frac{t}{2}}\mathbf{h}(t) + t^{\frac{3}{2}}e^{\frac{-t}{2}}\nabla F(\mathbf{h}(t)).$$

142 Let $\mathbf{m}(t) = \mathbf{x}'(t)$. Then Eq. (8) is equivalent to the following differential-algebraic system

$$\begin{cases} \mathbf{h}(t) = t^{-\frac{3}{2}}e^{\frac{t}{2}}\mathbf{x}(t), \\ \mathbf{x}'(t) = \mathbf{m}(t), \\ \mathbf{m}'(t) = -\mathbf{m}(t) - f(\mathbf{h}(t), t). \end{cases} \quad (13)$$

143 We parameterize $f(\mathbf{h}(t), t)$ as a neural network, recalled as $f(\mathbf{h}(t), t, \theta)$ with learnable parameter θ ,
144 and obtain the differential-algebraic version of NesterovNODE

$$\begin{cases} \mathbf{h}(t) = t^{-\frac{3}{2}}e^{\frac{t}{2}}\mathbf{x}(t), \\ \mathbf{x}'(t) = \mathbf{m}(t), \\ \mathbf{m}'(t) = -\mathbf{m}(t) - f(\mathbf{h}(t), t, \theta). \end{cases} \quad (14)$$

145 The singularity at $t = 0$ of the NesterovNODE in Eqs. (9) and (10) is now moved to the algebraic
146 equation of the above differential-algebraic system.

147 It is often the case when training ODE-based models that some functions diverge or explode at a
148 finite time. This phenomenon is called blow-up, which we demonstrate in Fig. 2. In order to alleviate
149 the blow-up problem, we introduce a generalized version of NesterovNODE, termed Generalized
150 NesterovNODE (GNesterovNODE). For the NesterovNODE, the potential blow-up is due to the
151 oscillation inherited from NAG scheme as can be seen in Fig. 1. The blow-up can also occur because
152 of the singularity caused by the factor $t^{-\frac{3}{2}}e^{\frac{t}{2}}$ in the algebraic equation. Following the techniques
153 presented in [Xia et al., 2021], we address these potential blow-up by applying an activation function σ
154 to the function $f(\mathbf{h}(t), t, \theta)$, the factor $t^{-\frac{3}{2}}e^{\frac{t}{2}}$, and the momentum state $\mathbf{m}(t)$ of the NesterovNODE.
155 In addition, the residual term $\xi\mathbf{h}(t)$, which stands for a skip connection [He et al., 2016], is also added
156 into the governing equation of $\mathbf{m}(t)$, which benefits training and generalization of GNesterovNODEs.
157 The GNesterovNODE differential-algebraic system is then given by

$$\begin{cases} \mathbf{h}(t) = \sigma(t^{-\frac{3}{2}}e^{\frac{t}{2}})\mathbf{x}(t), \\ \mathbf{x}'(t) = \sigma(\mathbf{m}(t)), \\ \mathbf{m}'(t) = -\mathbf{m}(t) - \sigma(f(\mathbf{h}(t), t, \theta)) - \xi\mathbf{h}(t), \end{cases} \quad (15)$$

158 with the initial conditions $\mathbf{h}(0) = \mathbf{h}_0$, $\mathbf{x}(0) = \mathbf{x}_0$, $\mathbf{m}(0) = \mathbf{m}_0$. It is noted that, from $\mathbf{h}(0) = \mathbf{h}_0$
159 and $\mathbf{h}'(0) = 0$, we must have $\mathbf{x}_0 = \mathbf{h}_0 \lim_{t \rightarrow 0} \sigma(t^{-\frac{3}{2}}e^{\frac{t}{2}})^{-1}$ and $\mathbf{m}_0 = \mathbf{h}_0 \lim_{t \rightarrow 0} \frac{d}{dt} \sigma(t^{-\frac{3}{2}}e^{\frac{t}{2}})^{-1}$. Fig. 2
160 shows that GNesterovNODE can indeed control the growth of $\mathbf{h}(t)$ effectively.

161 In general, GNesterovNODEs form a wide class of neural differential equations which are represented
162 by differential-algebraic systems and contain NesterovNODEs as a subclass. Eqs. (9) and (15)
163 define the forward ODE for the (G)NesterovNODE. To efficiently update the parameters during the
164 training process based on (G)NesterovNODEs, we use the continuous adjoint sensitivity given in
165 Propositions 1 and 2 below.

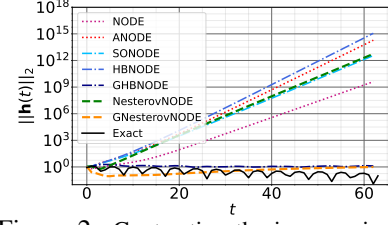


Figure 2: Contrasting the increase in the l_2 -norm of $\mathbf{h}(t)$ for NODE, ANODE, SONODE, HBNODE, GHBNODE, NesterovNODE, and GNesterovNODE over long integration time on the Silverbox Initialization task (More details in Appendix D.1).

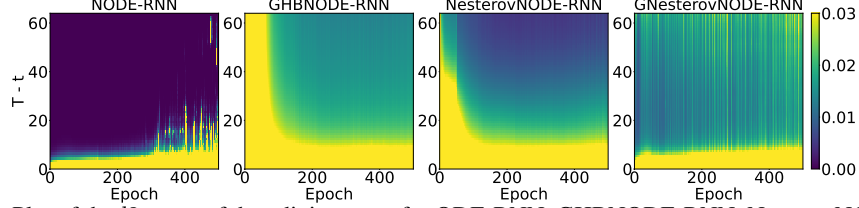


Figure 3: Plot of the l_2 -norm of the adjoint states for ODE-RNN, GHBNODE-RNN, NesterovNODE-RNN and GNesterovNODE-RNN back-propagated from the last time stamp. The term $T - t$ demonstrates the gap between the final time T and intermediate time t . When the gap $T - t$ becomes larger, NesterovNODE-RNN and GNesterovNODE-RNN can address the vanishing gradient problem due to the adjoint states of these methods' decay slowly.

Proposition 1 (Continuous adjoint equation for the second-order NesterovNODE). *The continuous adjoint state $\mathbf{a}(t) = \frac{\partial \mathcal{L}}{\partial \mathbf{h}(t)}$ of the NesterovNODE given in Eq. (9) satisfies the following NesterovN-ODE*

$$\mathbf{a}''(t) - \frac{3}{t}\mathbf{a}'(t) + \mathbf{a}(t) \frac{\partial f(\mathbf{h}(t), t, \theta)}{\partial \mathbf{h}} + \frac{3}{t^2}\mathbf{a}(t) = 0. \quad (16)$$

Proposition 2 (Continuous adjoint equation for GNesterovNODE). *The continuous adjoint state functions $\mathbf{a}_h(t) := \frac{\partial \mathcal{L}}{\partial \mathbf{h}(t)}$, $\mathbf{a}_x(t) := \frac{\partial \mathcal{L}}{\partial \mathbf{x}(t)}$ and $\mathbf{a}_m(t) := \frac{\partial \mathcal{L}}{\partial \mathbf{m}(t)}$ of the GNesterovNODE system given in Eq. (15) satisfy the following differential-algebraic adjoint system*

$$\begin{cases} \mathbf{a}_h(t) = \sigma(t^{-\frac{3}{2}} e^{\frac{t}{2}})^{-1} \mathbf{a}_x(t), \\ \mathbf{a}_x'(t) = t^{-\frac{3}{2}} e^{\frac{t}{2}} \mathbf{a}_m(t) \left[\frac{\partial \sigma(f(\mathbf{h}(t), t))}{\partial \mathbf{h}} + \xi \mathbf{I} \right], \\ \mathbf{a}_m'(t) = \mathbf{a}_m(t) - \mathbf{a}_x(t) \sigma'(\mathbf{m}(t)), \end{cases} \quad (17)$$

with the final value conditions $\mathbf{a}_h(T) = \frac{\partial \mathcal{L}}{\partial \mathbf{h}_T}$, $\mathbf{a}_x(T) = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_T}$ and $\mathbf{a}_m(T) = \frac{\partial \mathcal{L}}{\partial \mathbf{m}_T}$.

4 The Effectiveness of GNesterovNODE in Alleviating Vanishing Gradients

In neural networks with many layers, the vanishing gradient is one of the major issues [Pascanu et al., 2013]. In the cases of NODEs and their hybrid ODE-RNN variants, the vanishing gradient issue may occur when the adjoint state $\mathbf{a}(t) := \frac{\partial \mathcal{L}}{\partial \mathbf{h}(t)}$ goes to 0 quickly as $T - t$ increases. In this section, we will prove that this vanishing gradient issue can be avoided in GNesterovNODEs.

For the GNesterovNODE given in Eq. (15), the gradient $\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t}$ can be determined from $\frac{\partial \mathcal{L}}{\partial \mathbf{x}_t}$ via the algebraic relation:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_t} \frac{\partial \mathbf{x}_t}{\partial \mathbf{h}_t} = \sigma(t^{-\frac{3}{2}} e^{\frac{t}{2}})^{-1} \frac{\partial \mathcal{L}}{\partial \mathbf{x}_t}.$$

While the gradients $\frac{\partial \mathcal{L}}{\partial \mathbf{x}_t}$ and $\frac{\partial \mathcal{L}}{\partial \mathbf{m}_t}$ satisfy the following proposition.

Proposition 3. *For every $t \in (0, T)$, there exist a unit length row vector $v \in \mathbb{C}^N$ and an upper triangular matrix $U \in \mathbb{C}^{N \times N}$ such that*

$$\left\| \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{x}_t} & \frac{\partial \mathcal{L}}{\partial \mathbf{m}_t} \end{bmatrix} \right\|_2 = \|v \exp(U)\|_2 \left\| \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{x}_T} & \frac{\partial \mathcal{L}}{\partial \mathbf{m}_T} \end{bmatrix} \right\|_2,$$

and at least $\frac{N}{2}$ complex values in the diagonal of U have the real parts greater than or equal to $\frac{t-T}{2}$.

The term $v \exp(U)$ in the above proposition plays the essential role in the nonvanishing gradient issue.

Without loss of generality, we can assume that $U = \begin{bmatrix} U_{\text{large}} & P \\ \mathbf{0} & U_{\text{small}} \end{bmatrix}$ where all complex numbers in the diagonal of U_{large} (resp. U_{small}) have the real parts greater than or equal to (resp. smaller than) $\frac{1}{2}(t - T)$. According to Proposition 3, the size of U_{large} is at least $N/2$. Then we have,

$$\exp(U) = \begin{bmatrix} \exp(U_{\text{large}}) & \tilde{P} \\ \mathbf{0} & \exp(U_{\text{small}}) \end{bmatrix} \quad \text{and} \quad \|v \exp(U)\|_2 \geq \|v_{\text{large}} \exp(U_{\text{large}})\|_2.$$

Here, the vector v_{large} is the first m columns of v , and m is the size of U_{large} . Since the real parts of elements in the diagonal of U_{large} is no less than $\frac{1}{2}(t - T)$, $\exp(U_{\text{large}})$ decays at a rate at most

Table 1: The parameters count for the models in point cloud separation, image classifications, and the Walker2D kinematic simulation tasks and the test accuracy on CIFAR10/MNIST. Our methods are able to reach similar or better test accuracy than the baseline methods on CIFAR10 while retaining a similar test accuracy on MNIST.

<i>Model</i>	Number of parameters				Test Accuracy	
	<i>PC</i>	<i>MNIST</i>	<i>CIFAR10</i>	<i>Walker2D</i>	<i>CIFAR10</i>	<i>MNIST</i>
NODE	545	85316	173611	9929	0.5466 ± 0.0051	0.9531 ± 0.0042
ANODE	587	85462	172452	10019	0.6025 ± 0.0032	0.9816 ± 0.0024
SONODE	541	86179	171635	11471	0.6132 ± 0.0073	0.9824 ± 0.0013
HBNODE	582	85931	172916	10099	0.5989 ± 0.0035	0.9814 ± 0.0011
GHBNODE	582	85931	172916	10099	0.6085 ± 0.0050	0.9817 ± 0.0005
NesterovNODE	581	85930	172915	10098	0.5996 ± 0.0033	0.9824 ± 0.0015
GNesterovNODE	581	85930	172915	10098	0.6172 ± 0.0064	0.9807 ± 0.0013

$\frac{1}{2}(t - T)$. This results in the nonvanishing gradient of $\left[\frac{\partial \mathcal{L}}{\partial \mathbf{x}_t} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{m}_t} \right]$, hence so is $\left[\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{x}_t} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{m}_t} \right]$.

To illustrate, we take the Walker2D kinematic simulation task [Brockman et al., 2016a] in consideration, which requires learning long-term dependency effectively [Lechner & Hasani, 2020b]. We train ODE-RNN [Rubanova et al., 2019], GHBNODE-RNN [Xia et al., 2021], NesterovNODE-RNN and GNesterovNODE-RNN on this benchmark dataset (More details in Appendix D.4). Fig. 3 plots $\left\| \frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} \right\|_2$ for ODE-RNN, $\left\| \left[\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{m}_t} \right] \right\|_2$ for GHBNODE-RNN and $\left\| \left[\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{x}_t} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{m}_t} \right] \right\|_2$ for NesterovNODE-RNN and GNesterovNODE-RNN, showing that when the gap between the final time T and intermediate time t becomes larger, the adjoint states of NesterovNODE-RNN, GNesterovNODE-RNN and GHBNODE-RNN decay much more slowly than NODE-RNN. Thus, NesterovNODE-RNN and GNesterovNODE-RNN have the ability to tackle the vanishing gradient issue.

Remark 3. The gradient exploding problem can be effectively resolved via gradient clipping, training loss regularization, etc. [Pascanu et al., 2013; Erichson et al., 2020]. Therefore, in practice the vanishing gradient problem is the major issue for training deep neural networks [Pascanu et al., 2013].

5 Experimental Results

In this section, we empirically study the advantages of our proposed NesterovNODE/GNesterovNODE over the baseline NODEs and other popular NODE-based architectures, including the augmented NODE (ANODE) [Dupont et al., 2019a], the Second Order NODE (SONODE) [Norcliffe et al., 2020], HBNODE/GHBNODE [Xia et al., 2021] on a variety of benchmarks including point cloud separation, image classification, and kinetic simulation which involve different data modalities ranging from point cloud to images and time series. ANODEs augments the space on which the ODE is solved while SONODEs and (G)HBNODEs solve a second-order ODE. We aim to show that: (i) NesterovNODEs/GNesterovNODEs require significantly fewer NFEs while attaining similar or even better accuracy as the baselines; (ii) GNesterovNODEs avoid the blow-up of $\mathbf{h}(t)$ and thus improve over NesterovNODEs; (iii) NesterovNODEs/GNesterovNODEs capture better long-term dependencies than the baselines and achieve better results in long-sequence modeling tasks.

For all the experiments, we use Adam [Kingma & Ba, 2015] as the optimizer and Dormand-Prince 5(4) [Dormand & Prince, 1980] as the numerical ODE solver. We choose the network architecture used to parameterize $f(\mathbf{h}(t), t, \theta)$ so that our proposed models and the baselines have similar numbers of parameters in our experiments as shown in Table 1. Other training/model/dataset details are provided in Appendix D. All results are averaged over 5 runs with different seeds. We conduct the experiments on a server with 6 NVIDIA 2080Ti GPUs with 11GB of GPU memory.

5.1 Image classification

We validate the accuracy and efficiency advantage of NesterovNODE/GNesterovNODE for image classification on MNIST [Deng, 2012] and CIFAR10 [Krizhevsky et al., 2009] in comparison with other ODE-based baselines. We follow the same training and model settings as in [Xia et al., 2021].

NFEs. As shown in Fig. 4, our NesterovNODE and GNesterovNODE reduce the NFEs in both the forward and the backward propagations compared to the baseline models. Although the augmented

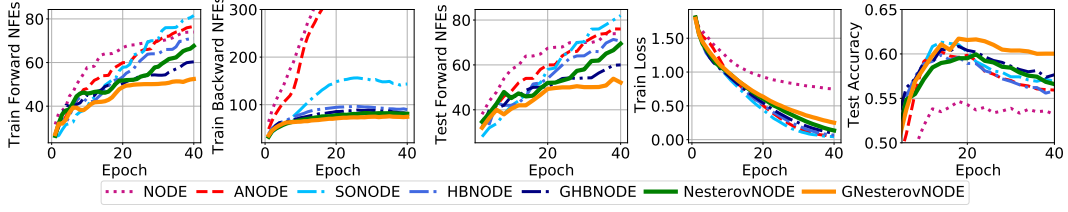


Figure 4: Contrasting the NFEs and accuracy of NODE-based baselines and our methods NesterovNODE/GNesterovNODE on the CIFAR10 dataset (Tolerance: 10^{-5}).

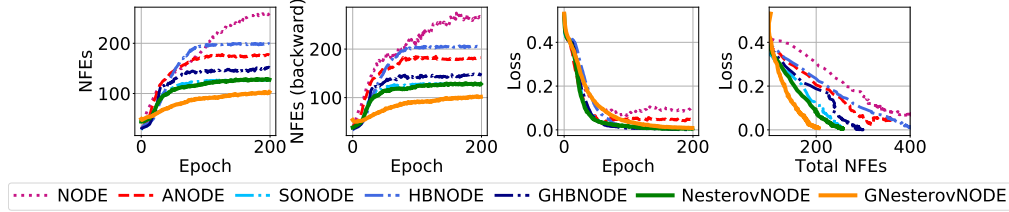


Figure 5: Contrasting the NFEs and training loss of NODE-based baselines and our NesterovNODE/GNesterovNODE on the point cloud benchmark. Results are averaged over 50 runs (Tolerance: 10^{-7}).

input dimensions in ANODE help reduce the NFEs compared to NODE, second-order methods reduce the NFEs more significantly. Compared to the other second-order methods i.e. SONODE, HBNODE, and GHBNODE, GNesterovNODE achieve much better NFE reductions on both MNIST (see Fig 11 in the Appendix) and CIFAR10, indicating the improvement in efficiency and stability of our methods over the other second-order baseline models. Such advancement is an essential step to scale GNesterovNODE to larger and more complex practical tasks.

Accuracy. Table 1 shows that our GNesterovNODE achieves the highest test accuracy on CIFAR10. On MNIST, NesterovNODE attains the second-highest test accuracy and very close to the best result from SONODE while being much more efficient than SONODE. This advantage in terms of accuracy can be associated with the small number of NFEs needed by NesterovNODE and GNesterovNODE above, which reduces the model complexity and leads to better generalization.

Note that GNesterovNODE improves over NesterovNODE in efficiency and accuracy (on larger and more challenging benchmark like CIFAR10). This justifies the effectiveness of our solution that introduces an additional bounded activation function σ and the residual term $\xi h(t)$ to prevent the blow-up of $h(t)$ as explained in Section 3.

5.2 Point cloud separation

We perform experiments on a point cloud separation task in order to verify that NesterovNODEs/GNesterovNODEs can learn effective features to separate two sets of point clouds. The first set consists of 40 points drawn from a circle with the radius $\|r\| < 0.5$, while the second set comprises 80 points drawn from an annulus with the inner and outer radius of 0.85 and 1, respectively, i.e. $0.85 < \|r\| < 1.0$. Fig. 5 shows that our NesterovNODE and GNesterovNODE models are able to converge to 0 loss consistently while other methods have difficulty to reach 0 loss. In addition, NesterovNODE and especially the GNesterovNODE require significantly fewer NFEs in forward and backward passes compared to the baselines. Thus, NesterovNODE and GNesterovNODE help improve both the training and the efficiency of the model. We plot the evolution of the point cloud separation through 100 epochs for a random run of each model in Appendix D.2. Like SONODE, HBNODE/GHBNODE, NesterovNODE/GNesterovNODE learn effective features that allow good separation between the two classes of point clouds in these experiments while NODE and ANODE fail for this task.

5.3 Walker2D kinematic simulation

In this section, we investigate NesterovNODE/GNesterovNODE when applied on time-series data. In particular, we use the ODE-RNN framework [Rubanova et al., 2019], with the recognition model being set to different ODE-based models, to study Walker2D kinematic simulation task, which requires learning long-term dependency effectively [Lechner & Hasani, 2020a]. As shown in Fig. 6, our NesterovNODE-RNN and GNesterovNODE-RNN not only reduce the NFEs both in the forward and the backward stages, but also achieve smaller test loss compared to the baseline models that we compare with, including (G)HBNODE-RNN, ANODE-RNN, and NODE-RNN. Although

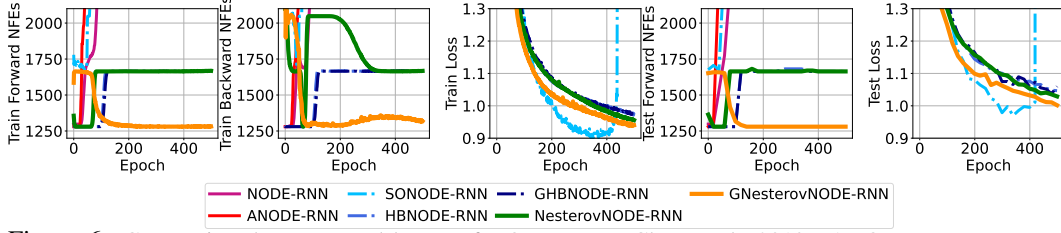


Figure 6: Contrasting the NFEs and losses of NODE-RNN [Chen et al., 2018], ANODE-RNN [Dupont et al., 2019a], HBNODE-RNN/GHBNODE-RNN [Xia et al., 2021], and our methods NesterovNODE-RNN/GNesterovNODE-RNN on the Walker2D dataset (Tolerance: 10^{-7}).

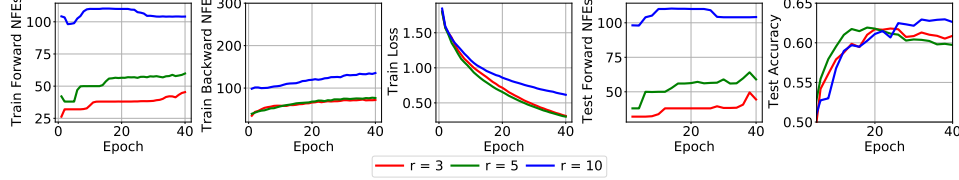


Figure 7: Contrasting different values of the factor r on CIFAR10 with GNesterovNODE (Tolerance: 10^{-5}).

SONODE-RNN achieves much smaller losses, its NFEs are too high, thus training SONODE-RNN for this task is much more time-consuming compared to our Nesterov-based methods.

6 Observed Properties of NesterovNODE in GNesterovNODE

In this section, we verify empirically that with the addition of the activation function σ and the skip connection, GNesterovNODEs still preserve important properties of NesterovNODEs as stated in Remark 1 and Remark 2. This hints that GNesterovNODEs have the same behaviors as NesterovNODEs while enjoying more stable training.

Effect of the Nesterov factor on NesterovNODE (Remark 1)

As stated in [Su et al., 2014], the “magic constant” 3 can be replaced by a constant $r > 3$ while maintaining a convergence rate of $O(1/k^2)$. In this work’s experiments, a larger r leads to a better test loss eventually despite a smaller r outperforming in the beginning, but their experiments are based on Nesterov ODE. This section investigates whether this behavior extends to our generalized method GNesterovNODE. As shown in Fig. 7 and Table 2, a larger r also leads to a higher test accuracy, and for $r = 3$ and $r = 5$, the model converges to its best test accuracy sooner than $r = 10$. Although the training loss of $r = 10$ is larger than $r = 3$ and $r = 5$, we observe that all three values of r reach their best testing accuracy when the training loss is approximately in the range of $[0.65, 0.75]$, which hints that decreasing the loss further leads to overfitting. One more interesting observation is that a higher value of r increases the forward and backward NFEs. Intuitively speaking, the term $-\frac{r}{t}\mathbf{h}'(t)$ in the NesterovNODE, which is opposite to the product of the damping parameter r/t and the velocity $\mathbf{h}'(t)$, represents the friction force of the model (see [Su et al., 2014, Section 4]). When r is larger, the friction force resists the movement of $\mathbf{h}(t)$ along the trajectory stronger, which slows down the convergence of training loss.

Stability of NesterovNODE (Remark 2)

The NesterovNODE is more numerically stable than NODE in the sense that the step size in the Euler method for solving the NesterovNODE can be chosen larger while the stability of the numerical solution is still guaranteed (see Remark 2). We hypothesize that this fact still holds for the GNesterovNODE in comparison with both NODE and GHBNODE. To illustrate this fact, we perform experiments on CIFAR10 using Euler solvers with large step sizes (0.1, 0.2, 0.5). Due to the instability of large step size, GHBNODE moves fast to the maximum accuracy points and then goes down, as shown in Fig. 8, while NODE achieves high train loss and low test accuracy. Even when the step size is big, GNesterovNODE’s training is stable as the curve evolves gradually. More interestingly, as shown in Fig. 9, GNesterovNODE solved with rk4 (Fourth-order Runge-Kutta with 3/8 rule) solver using a large step size outperforms the same model solved with the adaptive step size solvers like dopri5 (Dormand-Prince of order 5) solver. This shows the promise of using large step sizes in GNesterovNODE to obtain models with better accuracy and efficiency.

Table 2: Test accuracy for GNesterovNODE on CIFAR10 with varying Nesterov factors r .

Value of r	Test accuracy
3	0.6180
5	0.6192
10	0.6296

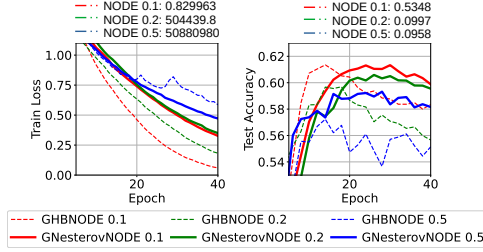


Figure 8: Train loss and test accuracy for large step sizes of Euler solver on CIFAR10 training with NODE, GHBNODE and GNesterovNODE. In the legend, we place the step size next to the model’s name.

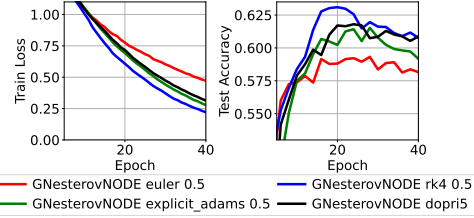


Figure 9: Train loss and test accuracy for various solvers with large step size 0.5 on CIFAR10 training with GNesterovNODE. In the legend, we place the step size next to the model’s name.

7 Related Work

Increasing efficiency of training NODEs. Several methods have been proposed to reduce the NFEs in NODEs and increase the model efficiency. Among them are works that use weight decay [Grathwohl et al., 2019a] and other regularizers applied on the solver and the learned dynamics [Finlay et al., 2020; Kelly et al., 2020; Ghosh et al., 2020; Pal et al., 2021]. Other works employ input augmentation [Dupont et al., 2019b], data-control [Massaroli et al., 2020] and depth-variance [Massaroli et al., 2020] to reduce the NFEs needed to compute the ODE solution. Another method replaces norms with seminorms in the backward continuous adjoint stage to reduce backward NFEs [Kidger et al., 2021]. NesterovNODE solves second-order Nesterov ODE to reduce both forward and backward NFEs.

Second-order dynamical systems. Second-order ODEs have also been employed to speed up NODEs. SONODE [Norcliffe et al., 2020] replaces the first-order ODE in Eq. (1) by a second ODE which can be solved as a system of first-order ODEs. HBNODE [Xia et al., 2021] also solves a second-order ODE but with an additional constant damping parameter, which corresponds to the ODE limit of the HB momentum method. NesterovNODE solves the second-order Nesterov ODE with a time-dependent damping parameter, which corresponds to the ODE limit of the NAG momentum method. These momentum-based systems have also been employed in designing deep neural networks as in [Moreau & Bruna, 2017; Li et al., 2018; Sander et al., 2021]

Learning long-term dependencies. The ability of a model to learn long-term dependencies is highly needed to scale up the model to large-scale tasks that involve very long sequences. Existing works try to alleviate exploding or vanishing gradient issues happened during the training of recurrent neural networks, including [Arjovsky et al., 2016; Wisdom et al., 2016; Jing et al., 2017; Vorontsov et al., 2017; Mhammedi et al., 2017; Helfrich et al., 2018]. Recently, learning long-term dependencies with NODEs has been explored. For example, [Lechner & Hasani, 2020a] integrate a long-short term memory cell into NODEs. Among the hallmarks of NesterovNODEs is that our proposed models can directly capture long-term dependencies in long sequences.

8 Concluding Remarks

In this paper, we propose the NesterovNODE and its generalized version GNesterovNODE that solve the second-order ODE limit of NAG. These models take advantage of the convergence rate $\mathcal{O}(1/k^2)$ of the NAG scheme to gain acceleration over NODEs and the existing NODE-based models such as HBNODE by reducing the NFEs in solving both forward and backward ODEs. Our Nesterov-based NODEs also achieve better accuracy than NODEs and outperform or at least are on par with other NODE-based models in our experiments while requiring much fewer NFEs. We also prove that NesterovNODEs and GNesterovNODEs can avoid the vanishing gradient issue and can capture long-term dependencies in long sequences effectively. It is worth mentioning that NesterovNODEs/GNesterovNODEs do not introduce any inherently negative societal impact. Also, the high-resolution ODE in [Shi et al., 2019, 2021] is an alternative way to take a continuous-time limit of the NAG scheme and heavy-ball momentum method. This gradient correction allows the NAG scheme to achieve an inverse cubic rate for minimizing the squared gradient norm, which is much better than the inverse square rate in the low-resolution ODE our method use. This is a limitation of our method and it is interesting to explore how to incorporate this construction into NesterovNODEs.

References

- Arjovsky, M., Shah, A., and Bengio, Y. Unitary evolution recurrent neural networks. In *International Conference on Machine Learning*, pp. 1120–1128, 2016.
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. Openai gym. *CoRR*, abs/1606.01540, 2016a. URL <http://arxiv.org/abs/1606.01540>.
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016b.
- Chen, R. T. Q., Rubanova, Y., Bettencourt, J., and Duvenaud, D. K. Neural ordinary differential equations. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL <https://proceedings.neurips.cc/paper/2018/file/69386f6bb1dfed68692a24c8686939b9-Paper.pdf>.
- De Brouwer, E., Simm, J., Arany, A., and Moreau, Y. Gru-ode-bayes: Continuous modeling of sporadically-observed time series. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/455cb2657aaa59e32fad80cb0b65b9dc-Paper.pdf>.
- Deng, L. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- Dormand, J. R. and Prince, P. J. A family of embedded runge-kutta formulae. *Journal of computational and applied mathematics*, 6(1):19–26, 1980.
- Dupont, E., Doucet, A., and Teh, Y. W. Augmented neural odes. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019a. URL <https://proceedings.neurips.cc/paper/2019/file/21be9a4bd4f81549a9d1d241981cec3c-Paper.pdf>.
- Dupont, E., Doucet, A., and Teh, Y. W. Augmented neural odes. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019b. URL <https://proceedings.neurips.cc/paper/2019/file/21be9a4bd4f81549a9d1d241981cec3c-Paper.pdf>.
- Erichson, N. B., Azencot, O., Queiruga, A., Hodgkinson, L., and Mahoney, M. W. Lipschitz recurrent neural networks. *arXiv preprint arXiv:2006.12070*, 2020.
- Finlay, C., Jacobsen, J.-H., Nurbekyan, L., and Oberman, A. M. How to train your neural ode: the world of jacobian and kinetic regularization. In *ICML*, pp. 3154–3164, 2020. URL <http://proceedings.mlr.press/v119/finlay20a.html>.
- Ghosh, A., Behl, H., Dupont, E., Torr, P., and Namboodiri, V. Steer : Simple temporal regularization for neural ode. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 14831–14843. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/a9e18cb5dd9d3ab420946fa19ebbbf52-Paper.pdf>.
- Grathwohl, W., Chen, R. T., Bettencourt, J., and Duvenaud, D. Scalable reversible generative models with free-form continuous dynamics. In *International Conference on Learning Representations*, 2019a.
- Grathwohl, W., Chen, R. T. Q., Bettencourt, J., and Duvenaud, D. Scalable reversible generative models with free-form continuous dynamics. In *International Conference on Learning Representations*, 2019b. URL <https://openreview.net/forum?id=rJxgkNCcK7>.
- Haber, E. and Ruthotto, L. Stable architectures for deep neural networks. *Inverse problems*, 34(1):014004, 2017.

388 He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings*
389 *of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.

390 Helfrich, K., Willmott, D., and Ye, Q. Orthogonal recurrent neural networks with scaled Cayley
391 transform. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on*
392 *Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1969–1978,
393 Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR. URL [http://proceedings.mlr.](http://proceedings.mlr.press/v80/helfrich18a.html)
394 [press/v80/helfrich18a.html](http://proceedings.mlr.press/v80/helfrich18a.html).

395 Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural computation*, 9(8):1735–1780,
396 1997.

397 Jing, L., Shen, Y., Dubcek, T., Peurifoy, J., Skirlo, S., LeCun, Y., Tegmark, M., and Soljačić, M.
398 Tunable efficient unitary neural networks (eunn) and their application to rnns. In *Proceedings of*
399 *the 34th International Conference on Machine Learning-Volume 70*, pp. 1733–1741. JMLR. org,
400 2017.

401 Kelly, J., Bettencourt, J., Johnson, M. J., and Duvenaud, D. K. Learning differential equa-
402 tions that are easy to solve. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F.,
403 and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp.
404 4370–4380. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper/2020/file/](https://proceedings.neurips.cc/paper/2020/file/2e255d2d6bf9bb33030246d31f1a79ca-Paper.pdf)
405 [2e255d2d6bf9bb33030246d31f1a79ca-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/2e255d2d6bf9bb33030246d31f1a79ca-Paper.pdf).

406 Kidger, P. On neural differential equations. *arXiv preprint arXiv:2202.02435*, 2022.

407 Kidger, P., Chen, R. T., and Lyons, T. J. ” hey, that’s not an ode”: Faster ode adjoints via seminorms.
408 In *ICML*, pp. 5443–5452, 2021.

409 Kingma, D. and Ba, J. Adam: A method for stochastic optimization. In *International Conference on*
410 *Learning Representations*, 2015.

411 Köhler, J., Klein, L., and Noé, F. Equivariant flows: exact likelihood generative learning for symmetric
412 densities. In *International Conference on Machine Learning*, pp. 5361–5370. PMLR, 2020.

413 Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.

414 Langley, P. Crafting papers on machine learning. In Langley, P. (ed.), *Proceedings of the 17th*
415 *International Conference on Machine Learning (ICML 2000)*, pp. 1207–1216, Stanford, CA, 2000.
416 Morgan Kaufmann.

417 Lechner, M. and Hasani, R. Learning long-term dependencies in irregularly-sampled time series,
418 2020a.

419 Lechner, M. and Hasani, R. M. Learning long-term dependencies in irregularly-sampled time series.
420 *CoRR*, abs/2006.04418, 2020b. URL <https://arxiv.org/abs/2006.04418>.

421 Li, H., Yang, Y., Chen, D., and Lin, Z. Optimization algorithm inspired deep neural network structure
422 design. In *Asian Conference on Machine Learning*, pp. 614–629. PMLR, 2018.

423 Luštrek, M., Kaluža, B., Piltaver, R., Krivec, J., and Vidulin, V. Localization data for person activity
424 data set, 2010.

425 Massaroli, S., Poli, M., Park, J., Yamashita, A., and Asama, H. Dissecting neural odes. In Larochelle,
426 H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin, H. (eds.), *Advances in Neural Information*
427 *Processing Systems*, volume 33, pp. 3952–3963. Curran Associates, Inc., 2020. URL [https:](https://proceedings.neurips.cc/paper/2020/file/293835c2cc75b585649498ee74b395f5-Paper.pdf)
428 [/proceedings.neurips.cc/paper/2020/file/293835c2cc75b585649498ee74b395f5-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/293835c2cc75b585649498ee74b395f5-Paper.pdf).

429 Mhammedi, Z., Hellicar, A., Rahman, A., and Bailey, J. Efficient orthogonal parametrisation of
430 recurrent neural networks using householder reflections. In *Proceedings of the 34th International*
431 *Conference on Machine Learning-Volume 70*, pp. 2401–2409. JMLR. org, 2017.

432 Moreau, T. and Bruna, J. Understanding the learned iterative soft thresholding algorithm with matrix
433 factorization. *arXiv preprint arXiv:1706.01338*, 2017.

- 434 Nesterov, Y. E. A method for solving the convex programming problem with convergence rate $o(1/k^2)$. In *Dokl. Akad. Nauk Sssr*, volume 269, pp. 543–547, 1983.
- 435
- 436 Norcliffe, A., Bodnar, C., Day, B., Simidjievski, N., and Lió, P. On second order behaviour
437 in augmented neural odes. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F.,
438 and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp.
439 5911–5921. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper/2020/file/](https://proceedings.neurips.cc/paper/2020/file/418db2ea5d227a9ea8db8e5357ca2084-Paper.pdf)
440 [418db2ea5d227a9ea8db8e5357ca2084-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/418db2ea5d227a9ea8db8e5357ca2084-Paper.pdf).
- 441 Pal, A., Ma, Y., Shah, V., and Rackauckas, C. V. Opening the blackbox: Accelerating neural
442 differential equations by regularizing internal solver heuristics. In Meila, M. and Zhang, T.
443 (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of
444 *Proceedings of Machine Learning Research*, pp. 8325–8335. PMLR, 18–24 Jul 2021. URL
445 <https://proceedings.mlr.press/v139/pal21a.html>.
- 446 Pascanu, R., Mikolov, T., and Bengio, Y. On the difficulty of training recurrent neural networks. In
447 *International Conference on Machine Learning*, pp. 1310–1318, 2013.
- 448 Polyak, B. T. Some methods of speeding up the convergence of iteration methods. *USSR Computa-*
449 *tional Mathematics and Mathematical Physics*, 4(5):1–17, 1964.
- 450 Pontryagin, L. S. *Mathematical theory of optimal processes*. CRC press, 1987.
- 451 Rackauckas, C. and Nie, Q. Differentialequations.jl – a performant and feature-rich ecosystem
452 for solving differential equations in julia. *The Journal of Open Research Software*, 5(1), 2017.
453 doi: 10.5334/jors.151. URL [https://app.dimensions.ai/details/publication/pub.1085583166andhttp://](https://app.dimensions.ai/details/publication/pub.1085583166andhttp://openresearchsoftware.metajnl.com/articles/10.5334/jors.151/galley/245/download/)
454 openresearchsoftware.metajnl.com/articles/10.5334/jors.151/galley/245/download/. Exported
455 from <https://app.dimensions.ai> on 2019/05/05.
- 456 Rackauckas, C., Innes, M., Ma, Y., Bettencourt, J., White, L., and Dixit, V. Diffeqflux.jl - A julia
457 library for neural differential equations. *CoRR*, abs/1902.02376, 2019. URL [http://arxiv.org/abs/](http://arxiv.org/abs/1902.02376)
458 [1902.02376](http://arxiv.org/abs/1902.02376).
- 459 Rubanova, Y., Chen, R. T. Q., and Duvenaud, D. K. Latent ordinary differential equations for
460 irregularly-sampled time series. In Wallach, H., Larochelle, H., Beygelzimer, A., d’Alché-Buc,
461 F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, vol-
462 ume 32. Curran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper/2019/file/](https://proceedings.neurips.cc/paper/2019/file/42a6845a557bef704ad8ac9cb4461d43-Paper.pdf)
463 [42a6845a557bef704ad8ac9cb4461d43-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/42a6845a557bef704ad8ac9cb4461d43-Paper.pdf).
- 464 Ruthotto, L. and Haber, E. Deep neural networks motivated by partial differential equations. *Journal*
465 *of Mathematical Imaging and Vision*, pp. 1–13, 2019.
- 466 Ruthotto, L., Osher, S. J., Li, W., Nurbekyan, L., and Fung, S. W. A machine learning framework for
467 solving high-dimensional mean field game and mean field control problems. *Proceedings of the*
468 *National Academy of Sciences*, 117(17):9183–9193, 2020. ISSN 0027-8424. doi: 10.1073/pnas.
469 1922204117. URL <https://www.pnas.org/content/117/17/9183>.
- 470 Sander, M. E., Ablin, P., Blondel, M., and Peyré, G. Momentum residual neural networks. In Meila,
471 M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*,
472 volume 139 of *Proceedings of Machine Learning Research*, pp. 9276–9287. PMLR, 18–24 Jul
473 2021. URL <https://proceedings.mlr.press/v139/sander21a.html>.
- 474 Shi, B., Du, S. S., Su, W., and Jordan, M. I. Acceleration via symplectic discretization of high-
475 resolution differential equations. *Advances in Neural Information Processing Systems*, 32, 2019.
- 476 Shi, B., Du, S. S., Jordan, M. I., and Su, W. J. Understanding the acceleration phenomenon via
477 high-resolution differential equations. *Mathematical Programming*, pp. 1–70, 2021.
- 478 Su, W., Boyd, S., and Candes, E. A differential equation for modeling nesterov’s accelerated
479 gradient method: Theory and insights. In *Advances in Neural Information Processing Systems*, pp.
480 2510–2518, 2014.

- 481 Todorov, E., Erez, T., and Tassa, Y. Mujoco: A physics engine for model-based control. In *Intelligent*
482 *Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, pp. 5026–5033. IEEE,
483 2012. URL <https://ieeexplore.ieee.org/abstract/document/6386109/>.
- 484 Tsitouras, C. Runge–kutta pairs of order 5 (4) satisfying only the first column simplifying assumption.
485 *Computers & Mathematics with Applications*, 62(2):770–775, 2011.
- 486 Vorontsov, E., Trabelsi, C., Kadoury, S., and Pal, C. On orthogonality and learning recurrent networks
487 with long term dependencies. In *Proceedings of the 34th International Conference on Machine*
488 *Learning-Volume 70*, pp. 3570–3578. JMLR. org, 2017.
- 489 Weinan, E. A proposal on machine learning via dynamical systems. *Communications in Mathematics*
490 *and Statistics*, 5(1):1–11, 2017.
- 491 Wigren, T. and Schoukens, J. Three free data sets for development and benchmarking in nonlinear
492 system identification. In *2013 European control conference (ECC)*, pp. 2933–2938. IEEE, 2013.
- 493 Wisdom, S., Powers, T., Hershey, J., Le Roux, J., and Atlas, L. Full-capacity unitary recurrent neural
494 networks. In *Advances in Neural Information Processing Systems*, pp. 4880–4888, 2016.
- 495 Xia, H., Suliafu, V., Ji, H., Nguyen, T. M., Bertozzi, A., Osher, S., and Wang, B. Heavy ball neural
496 ordinary differential equations. In Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W.
497 (eds.), *Advances in Neural Information Processing Systems*, 2021. URL [https://openreview.net/](https://openreview.net/forum?id=fYLfs9yrtMQ)
498 [forum?id=fYLfs9yrtMQ](https://openreview.net/forum?id=fYLfs9yrtMQ).
- 499 Yildiz, C., Heinonen, M., and Lahdesmaki, H. Ode2vae: Deep generative second order odes
500 with bayesian neural networks. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc,
501 F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, vol-
502 *ume 32*. Curran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper/2019/file/](https://proceedings.neurips.cc/paper/2019/file/99a401435dcb65c4008d3ad22c8cdad0-Paper.pdf)
503 [99a401435dcb65c4008d3ad22c8cdad0-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/99a401435dcb65c4008d3ad22c8cdad0-Paper.pdf).
- 504 Zhong, Y. D., Dey, B., and Chakraborty, A. Symplectic ode-net: Learning hamiltonian dynamics
505 with control. In *International Conference on Learning Representations*, 2020. URL [https://](https://openreview.net/forum?id=ryxmb1rKDS)
506 openreview.net/forum?id=ryxmb1rKDS.

507 Checklist

- 508 1. For all authors...
- 509 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
510 contributions and scope? [Yes]
- 511 (b) Did you describe the limitations of your work? [Yes] See Section 8
- 512 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See
513 Section 8
- 514 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
515 them? [Yes]
- 516 2. If you are including theoretical results...
- 517 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 518 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix
- 519 3. If you ran experiments...
- 520 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
521 mental results (either in the supplemental material or as a URL)? [Yes]
- 522 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
523 were chosen)? [Yes]
- 524 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
525 ments multiple times)? [Yes]
- 526 (d) Did you include the total amount of compute and the type of resources used (e.g., type
527 of GPUs, internal cluster, or cloud provider)? [Yes]
- 528 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 529 (a) If your work uses existing assets, did you cite the creators? [N/A]
530 (b) Did you mention the license of the assets? [N/A]
531 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
532
533 (d) Did you discuss whether and how consent was obtained from people whose data you're
534 using/curating? [N/A]
535 (e) Did you discuss whether the data you are using/curating contains personally identifiable
536 information or offensive content? [N/A]
537 5. If you used crowdsourcing or conducted research with human subjects...
538 (a) Did you include the full text of instructions given to participants and screenshots, if
539 applicable? [N/A]
540 (b) Did you describe any potential participant risks, with links to Institutional Review
541 Board (IRB) approvals, if applicable? [N/A]
542 (c) Did you include the estimated hourly wage paid to participants and the total amount
543 spent on participant compensation? [N/A]

Supplementary materials for Improving Neural Ordinary Differential Equations with Nesterov's Accelerated Gradient Method

A Review of the adjoint equation and the gradient for the first-order NODEs

A NODE for a hidden feature $\mathbf{z} : \mathbb{R} \rightarrow \mathbb{R}^N$ takes of the form

$$\mathbf{z}'(t) = g(\mathbf{z}(t), t, \theta), \quad \mathbf{z}(0) = \mathbf{z}_0, \quad (18)$$

where $g(\mathbf{z}(t), t, \theta) \in \mathbb{R}^N$ is a neural network with learnable parameters θ . For a scalar loss function \mathcal{L} , the adjoint state $\mathbf{a}(t) := \frac{\partial \mathcal{L}}{\partial \mathbf{z}(t)}$ satisfies the following differential equation

$$\mathbf{a}'(t) = -\mathbf{a}(t) \frac{\partial g(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}}, \quad \mathbf{a}(T) = \frac{\partial \mathcal{L}}{\partial \mathbf{z}(T)}. \quad (19)$$

B The adjoint equations for the NesterovNODEs and GNesterovNODEs

The adjoint equation for the NesterovNODEs (Proposition 1) is an implication of [Norcliffe et al., 2020, Proposition 3.1] for Nesterov differential equations. In this section, we give the proofs for Proposition 2. The proofs will be intrinsically based on Eq. (19).

Proof of Proposition 2 - the adjoint equation for GNesterovNODE

The GNesterovNODE is formulated as the following differential-algebraic system:

$$\begin{cases} \mathbf{h}(t) = \sigma(k(t))\mathbf{x}(t), \\ \mathbf{x}'(t) = \sigma(\mathbf{m}(t)), \\ \mathbf{m}'(t) = -\mathbf{m}(t) - \sigma(f(\mathbf{h}(t), t)) - \xi \mathbf{h}(t), \end{cases} \quad (20)$$

where $k(t) = t^{-\frac{3}{2}} e^{\frac{t}{2}}$. The adjoints of this system are given by $\mathbf{a}_h(t) := \frac{\partial \mathcal{L}}{\partial \mathbf{h}(t)}$, $\mathbf{a}_x(t) := \frac{\partial \mathcal{L}}{\partial \mathbf{x}(t)}$ and $\mathbf{a}_m(t) := \frac{\partial \mathcal{L}}{\partial \mathbf{m}(t)}$. From the first equation in the system, we have

$$\mathbf{a}_h(t) = \frac{\partial \mathcal{L}}{\partial \mathbf{x}(t)} \frac{\partial \mathbf{x}(t)}{\partial \mathbf{h}(t)} = \frac{1}{\sigma(k(t))} \mathbf{a}_x(t). \quad (21)$$

To determine the dynamics of $\mathbf{a}_x(t)$ and $\mathbf{a}_m(t)$, we rewrite the last two equations in system (20) as the first-order system

$$\begin{cases} \mathbf{x}'(t) = \sigma(\mathbf{m}(t)), \\ \mathbf{m}'(t) = -\mathbf{m}(t) - \sigma(f(k(t)\mathbf{x}(t), t)) - \xi k(t)\mathbf{x}(t). \end{cases}$$

Set $\mathbf{z}(t) = [\mathbf{x}(t) \quad \mathbf{m}(t)]^T$ and

$$g(\mathbf{z}(t), t, \theta) = \begin{bmatrix} g_1 \\ g_2 \end{bmatrix} = \begin{bmatrix} \sigma(\mathbf{m}(t)) \\ -\mathbf{m}(t) - \sigma(f(k(t)\mathbf{x}(t), t)) - \xi k(t)\mathbf{x}(t) \end{bmatrix},$$

then the first-order NesterovNODE can be rewritten as

$$\mathbf{z}'(t) = g(\mathbf{z}(t), t, \theta).$$

In this case, we have

$$\begin{aligned} \frac{\partial g(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}} &= \begin{bmatrix} \frac{\partial g_1}{\partial \mathbf{x}} & \frac{\partial g_1}{\partial \mathbf{m}} \\ \frac{\partial g_2}{\partial \mathbf{x}} & \frac{\partial g_2}{\partial \mathbf{m}} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{0} & \sigma'(\mathbf{m}(t)) \\ -k(t) \frac{\partial \sigma[f(\mathbf{h}(t), t, \theta)]}{\partial \mathbf{h}} - \xi k(t) \mathbf{I} & -\mathbf{I} \end{bmatrix}. \end{aligned}$$

Thus, by using Eq. (19), we obtain the first-order differential system for the adjoints \mathbf{a}_x and \mathbf{a}_m as

$$\begin{cases} \mathbf{a}_x'(t) = \mathbf{a}_m(t) \left[k(t) \frac{\partial \sigma(f(\mathbf{h}(t), t, \theta))}{\partial \mathbf{h}} + \xi k(t) \mathbf{I} \right], \\ \mathbf{a}_m'(t) = -\mathbf{a}_x(t) \sigma'(\mathbf{m}(t)) + \mathbf{a}_m(t). \end{cases} \quad (22)$$

Together with Eq. (21), we obtain the differential-algebraic system for the adjoints of the GNesterovNODE in Proposition 2.

568 **C Proof of Proposition 3 - the nonvanishing gradient for GNesterovNODEs**

569 Following the lines in [Xia et al., 2021], for a NODE given by $\mathbf{z}'(t) = g(\mathbf{z}(t), t, \theta)$, we have

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}_t} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}_T} \exp \left\{ - \int_T^t \frac{\partial g(\mathbf{z}(s), s, \theta)}{\partial \mathbf{z}} ds \right\}. \quad (23)$$

570 For the GNesterovNODE given in Eq. (15), the gradient $\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t}$ can be determined from $\frac{\partial \mathcal{L}}{\partial \mathbf{x}_t}$ via the
571 algebraic relation:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_t} \frac{\partial \mathbf{x}_t}{\partial \mathbf{h}_t} = \sigma(t^{-\frac{3}{2}} e^{\frac{t}{2}})^{-1} \frac{\partial \mathcal{L}}{\partial \mathbf{x}_t}.$$

572 While the gradients $\frac{\partial \mathcal{L}}{\partial \mathbf{x}_t}$ and $\frac{\partial \mathcal{L}}{\partial \mathbf{m}_t}$ can be determined by using Eq. (23) as

$$\begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{x}_t} & \frac{\partial \mathcal{L}}{\partial \mathbf{m}_t} \end{bmatrix} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{x}_T} & \frac{\partial \mathcal{L}}{\partial \mathbf{m}_T} \end{bmatrix} \cdot \exp(M), \quad (24)$$

573 where

$$M = - \int_T^t \begin{bmatrix} \mathbf{0} & \sigma'(\mathbf{m}(s)) \\ \frac{\partial \sigma(f)}{\partial \mathbf{x}} - \xi t^{-\frac{3}{2}} e^{\frac{t}{2}} \mathbf{I} & -\mathbf{I} \end{bmatrix} ds.$$

574 Let $M = QUQ^\top$ be a Schur decomposition of M where Q is an orthogonal matrix and U is an
575 upper triangular matrix with eigenvalues of M in the diagonal. Then from Eq. (24), we have

$$\left\| \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{x}_t} & \frac{\partial \mathcal{L}}{\partial \mathbf{m}_t} \end{bmatrix} \right\|_2 = \left\| \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{x}_T} & \frac{\partial \mathcal{L}}{\partial \mathbf{m}_T} \end{bmatrix} Q \exp(U) Q^\top \right\|_2.$$

576 Set $v = \frac{1}{\left\| \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{x}_T} & \frac{\partial \mathcal{L}}{\partial \mathbf{m}_T} \end{bmatrix} \right\|_2} \cdot \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{x}_T} & \frac{\partial \mathcal{L}}{\partial \mathbf{m}_T} \end{bmatrix} Q$, then v is a unit length vector and

$$\left\| \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{x}_t} & \frac{\partial \mathcal{L}}{\partial \mathbf{m}_t} \end{bmatrix} \right\|_2 = \|v \exp(U)\|_2 \left\| \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{x}_T} & \frac{\partial \mathcal{L}}{\partial \mathbf{m}_T} \end{bmatrix} \right\|_2.$$

577 To finish the proof of Proposition 3, we need to prove that at least half of complex numbers in the
578 diagonal of U have the real parts greater than or equal to $\frac{t-T}{2}$.

We first claim that the eigenvalues of M can be paired in such a way that each pair has the sum $t - T$.

To prove the claim, we write $M = \begin{bmatrix} \mathbf{0} & A \\ B & (t-T)\mathbf{I} \end{bmatrix}$, where

$$A = - \int_T^t \sigma'(\mathbf{m}(s)) ds, \quad \text{and} \quad B = - \int_T^t \left[\frac{\partial \sigma(f)}{\partial \mathbf{x}} - \xi t^{-\frac{3}{2}} e^{\frac{t}{2}} \mathbf{I} \right] ds.$$

579 Since the matrices $(t-T)\mathbf{I}$ and A commute, the characteristic polynomial of M can be determined
580 as

$$\det(\lambda \mathbf{I} - M) = \det(\lambda(\lambda - t + T)\mathbf{I} - BA)$$

581 which is the value of the characteristic polynomial of the matrix BA at $\lambda(\lambda - t + T)$. Over the field
582 of complex numbers, the characteristic polynomial of BA is splitted completely and it has the form

$$\det(\lambda \mathbf{I} - BA) = \prod_{i=1}^N (\lambda - \lambda_{BA,i}),$$

583 where $\lambda_{BA,i}$ are the eigenvalues of BA . Then the characteristic polynomial of M has the form

$$\det(\lambda \mathbf{I} - M) = \prod_{i=1}^N (\lambda(\lambda - t + T) - \lambda_{BA,i}).$$

584 Thus the eigenvalues of M can be paired in such a way that each pair is the roots of the quadratic
585 equation

$$\lambda(\lambda - t + T) - \lambda_{BA,i}.$$

586 Such pairs always have the sum $t - T$. The claim is proved.

587 Since the diagonal of U are exactly the eigenvalues of M , the claim implies that at least half of
588 complex numbers in the diagonal of U are greater than or equal to $\frac{t-T}{2}$. The proposition is then
589 proved.

Table 3: The hyperparameters for the models.

Model	NODE	ANODE	SONODE	(G)HBNODE	(G)NesterovNODE
n (Initialization)	1	2	1	1	1
n (Point Cloud)	2	3	2	2	2
h (Point Cloud)	20	20	13	14	14
p (MNIST)	0	5	4	4	4
n (MNIST)	1	6	5	5	5
h (MNIST)	92	64	50	50	50
p (CIFAR10)	0	10	9	9	9
n (CIFAR10)	3	13	12	12	12
h (CIFAR10)	125	64	50	51	51

Table 4: The hyperparameters for ODE-RNN integration models.

Model	NODE	ANODE	SONODE	(G)HBNODE	(G)NesterovNODE
d	1	1	2	2	2
n (Walker2D)	24	24	23	24	24
h_1 (Walker2D)	72	72	46	48	48
h_2 (Walker2D)	48	48	46	48	48

Table 5: The ξ hyperparameters for GHBNODE and GNesterovNODE.

Model	ξ
Initialization	learnable
MNIST	learnable
CIFAR	1.5
Point Cloud	2
Walker 2D	learnable

D Experimental details

We list some experimental details that are common to all experiments before presenting more details for each experiment. For ODE solvers in the experiments (`dopri5`, `euler`, `rk4`, `explicit_adams`), we use the implementation provided by `torchdiffeq`¹.

- fc_n is a fully-connected layer with the output dimension of size n .
- HTanh: $\text{HardTanh}(-5, 5)$.
- LReLU: $\text{LeakyReLU}(0.3)$.
- tpad: Padding a set of features with the value of the time t . Specifically, this is equivalent to augmenting the feature tensor with shape $c \times x \times y$ to with a time tensor with shape $1 \times x \times y$ filled with the value of t , creating a time-augmented feature tensor with shape $(c + 1) \times x \times y$.
- We use a learnable γ with for HBNODE/GHBNODE for all tasks.
- The values of ξ used for GHBNODE and GNesterovNODE used in all tasks are in Table 5.
- Except for the experiments in Section 6 where we investigate the effect of the variation of the Nesterov factor r , we choose $r = 3$ for the remaining experiments.
- $n^* = n$ for all models except for SONODE where $n^* = 2n$. The bounded activation function σ is chosen to be either `tanh` or `hardtanh` in PyTorch.
- The tolerance (for adaptive solvers) and step sizes (for fixed step-size solvers) used in the experiments are in Table 6.
- Other hyperparameters are in Table 3 and Table 4.

¹<https://github.com/rtqichen/torchdiffeq>

D.1 Experimental details used in Section 3 Silverbox Initialization Test

The Silverbox dataset [Wigren & Schoukens, 2013] is as follows: Given the input voltage $V1(t)$, the models must predict the output voltage $V2(t)$ and the experiment is evaluated over 64 time steps. Similar to [Xia et al., 2021], we parameterize the dynamic f of all methods with a dense layer. The network architecture is:

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

D.2 Experimental details used in the Point Cloud benchmark in Section 5.2

Following the setting in [Xia et al., 2021; Dupont et al., 2019a], we use a 3-layer neural network to parameterize the function $f(\mathbf{h}(t), t, \theta)$ on the right hand side of the ODE-based models in our study. We integrate the ODE from $t_0 = 1$ to $T = 2$, and pass the output $\mathbf{h}(T)$ at time T to a dense classifier. We also set the tolerance of the ODE solvers to be 10^{-7} so that the effect of numerical error is minimized. Visualization of the point cloud evolution for a random run of each model is in Fig. 10.

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

D.3 Experimental details for MNIST/CIFAR10 in Section 5.1

In our experiment, we parameterize $f(\mathbf{h}(t), t, \theta)$ in the NODE-based layer using a 3-layer neural network. The output of this NODE-based layer is then passed through a dense layer to perform the classification task. Following the augmentation approach in ANODE [Dupont et al., 2019b], we add p additional channels to input image, augmenting the number of image channels from c to $c + p$ where p is differently chosen for each method. The values of p chosen for each model are specified in Table 3. For SONODE, HBNODE, GHBNODE, NesterovNODE and GNesterovNODE, we also incorporate velocity or momentum of the same shape as the augment state. We also present extra experiment results for MNIST in Fig. 11.

The architecture for MNIST:

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

The architecture for CIFAR10:

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

D.4 Experimental details for Walker2D in Section 5.3

The dataset [Brockman et al., 2016b] 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 [Todorov et al., 2012]. Following the procedure in HBNODE [Xia et al., 2021], we randomly take out 10% of the data to make the time series irregularly-sampled. Each input sequence consists of 64 timestamps, which are recurrently fed through a hybrid technique, with the output of the hybrid method being transferred to a single dense layer to form the output time series. The goal is to generate an auto-regressive forecast with an output time series that is as close as the input sequence when shifted one time stamp to the right. The RNN and ODE are parameterized by a 3-layer network. The network architecture is:

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

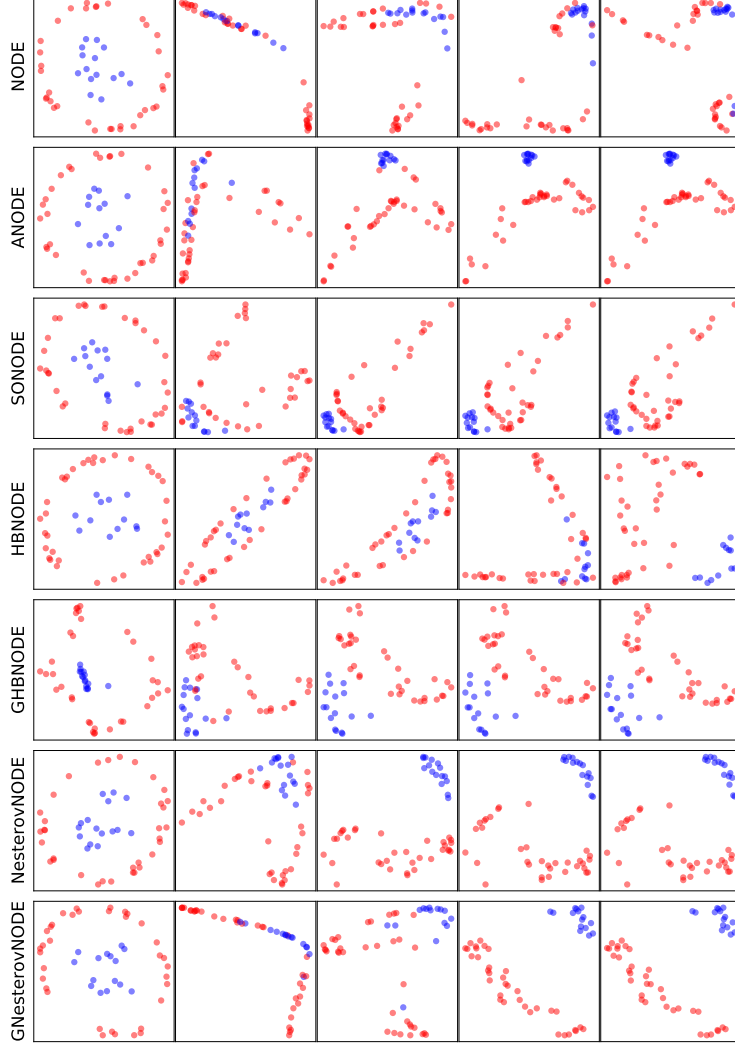


Figure 10: Visualization of the features transform by NODE, ANODE, SONODE, HBNODE, GHBNODE, NesterovNODE and GNesterovNODE after the first 100 epochs of a random run.

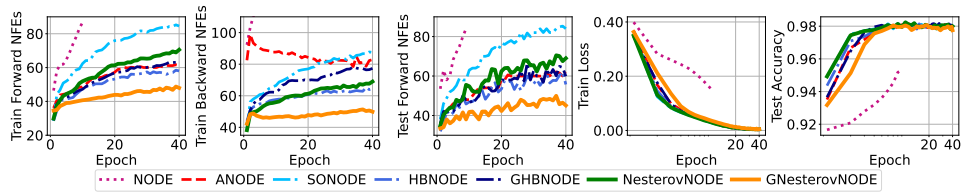


Figure 11: Contrasting the NFEs, training time, training loss, and test accuracy of NODE [Chen et al., 2018], ANODE [Dupont et al., 2019a], HBNODE/GHBNODE [Xia et al., 2021], and our methods NesterovNODE/GNesterovNODE on the MNIST dataset (Tolerance: 10^{-5}). The run for NODE is stopped due to long running time. The x-axes of the plots for training loss and testing accuracy are scaled logarithmically for visibility.

659 E Additional experiments

660 E.1 Integration time

661 In the the differential-algebraic version of NesterovNODE, we reparametrize from \mathbf{h} to \mathbf{x} in order
 662 to eliminate the time-dependent damping coefficient $\frac{3}{t}$ in the Nesterov ODE given by Eq. 8. In
 663 particular, we set $\mathbf{h}(t) = k(t)\mathbf{x}(t)$ and then find $k(t)$ such that Eq. 8 written in terms of $\mathbf{x}(t)$ matches

Table 6: Solvers, tolerance, and step sizes used for the experiments.

Experiments	Solvers	Tolerance	Step sizes
Silverbox Initialization (Section 3)	dopri5	10^{-7}	N/A
Point Cloud separation (Section 5.2)	dopri5	10^{-7}	N/A
MNIST/CIFAR10 (Section 5.1)	dopri5	10^{-5}	N/A
Walker2D (Section 5.3)	dopri5	10^{-7}	N/A
Stability of NesterovNODE (Section 6)	Euler, RK4, Explicit Adams, dopri5	10^{-5} for dopri5	0.1, 0.2, 0.5 (for details, refer to Fig. 8 and Fig. 9)

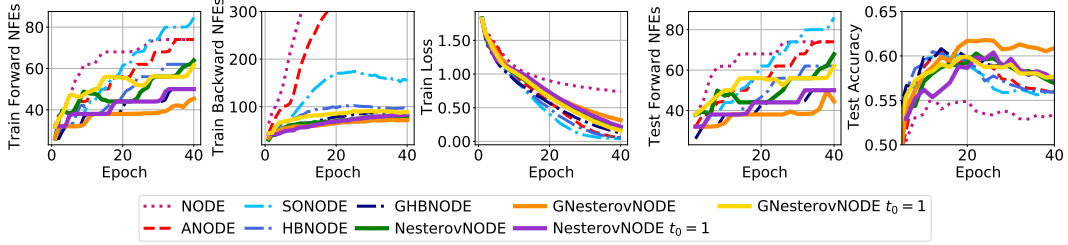


Figure 12: Contrasting the NFEs and accuracy of NODE [Chen et al., 2018], ANODE [Dupont et al., 2019a], HBNODE/GHBNODE [Xia et al., 2021], our methods NesterovNODE/GNesterovNODE and the methods NesterovNODE/GNesterovNODE with starting integration time $t_0 = 1$ on the CIFAR10 dataset (Tolerance: 10^{-5}).

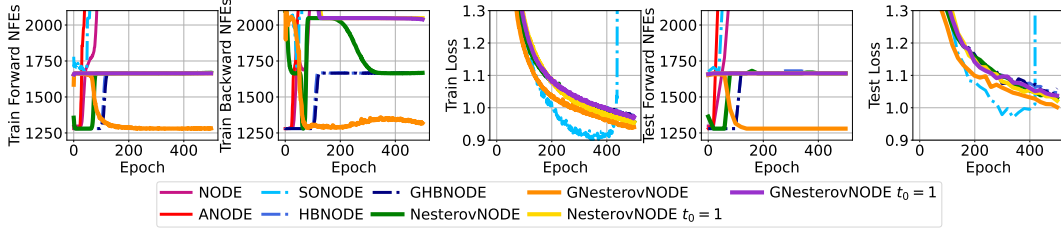


Figure 13: Contrasting the NFEs and losses of NODE [Chen et al., 2018], ANODE [Dupont et al., 2019a], HBNODE/GHBNODE [Xia et al., 2021], our methods NesterovNODE/GNesterovNODE and the methods NesterovNODE/GNesterovNODE with starting integration time $t_0 = 1$ on the Walker2D dataset (The tolerance is 10^{-7}).

the Heavy-Ball ODE with a constant damping coefficient given by Eq. 5. We find that $k(t) = t^{-\frac{3}{2}}e^{\frac{t}{2}}$ satisfies this requirement. Solving this Heavy-Ball ODE with respect to $\mathbf{x}(t)$ is easier and more stable than solving the Nesterov ODE with respect to $\mathbf{h}(t)$. To strengthen our proposed trick, we conduct two experiments (CIFAR 10 and Walker2d) using the original version of Nesterov in Eq. 10 with the starting integration time $t_0 = 1$. We have to change the starting integration to $t_0 = 1$ due to the singularity in $t_0 = 0$. The empirical results on CIFAR10 (Fig. 12) and Walker2d (Fig. 13) benchmarks show that the approach of changing starting integration time shows much worse accuracy and efficiency than using our reparametrization trick.

E.2 The Effect of Nesterov factor regarding the solver’s tolerance

We run another CIFAR10 experiment using smaller tolerance (10^{-7}) in order to study the effect of Nesterov factor regarding the solver’s tolerance. As shown in Fig. 14, when the tolerance are shrunk small enough, the forward NFEs of the smaller factor (3, 5 compared to 10) are still smaller, although the opposite is true for backward NFEs. Moreover, the larger factor converges to a better solution (higher accuracy) despite converging to its best solution later than the smaller factor.

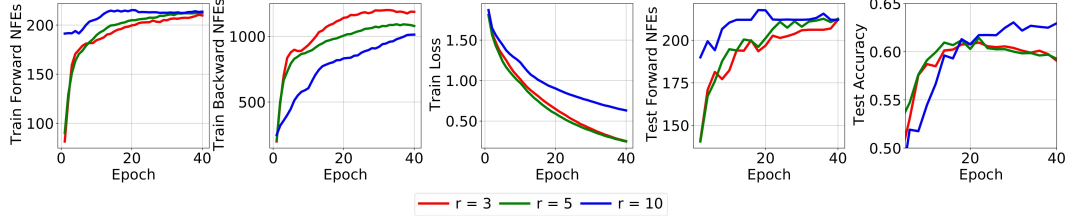


Figure 14: Contrasting the NFEs and accuracy of NODE-based baselines and our methods NesterovNODE/GNesterovNODE on the CIFAR10 dataset (Tolerance: 10^{-7}).

Table 7: GPU memory consumption of the ODE-based method on the CIFAR10 dataset.

Method	Maximum CUDA memory consumption (GiB)
NODE	3.2046
ANODE	2.3292
SONODE	2.2034
HBNODE	2.2346
GHBNODE	2.2456
NesterovNODE	2.2346
GNesterovNODE	2.2580

E.3 GPU memory consumption

We present the maximum CUDA memory consumption of the models used in the CIFAR10 experiments in Table 7. We extract the numbers using the function `max_memory_allocated`² in PyTorch.

E.4 Wall-clock time advantage of NesterovNODEs/GNesterovNODEs

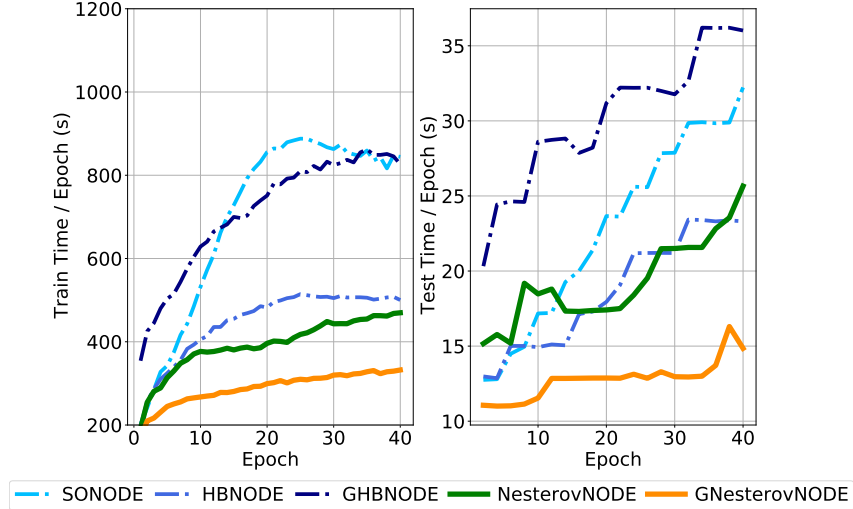


Figure 15: The total wall-clock training and testing time of SONODE, HBNODE, GHBNODE, NesterovNODE, and GNesterovNODE.

In this Figs. 15 and 16, we show the advantage of NesterovNODEs in wall-clock time. NesterovNODEs/GNesterovNODEs uses less time to achieve a comparable accuracy compared to other methods.

E.5 The Effect of using Seminorm for reducing backward NFEs on GNesterovNODE

We tested the effect of the seminorm method [Kidger et al., 2021] on GNesterovNODE. It is surprising that GNesterovNODE with seminorm is using much more NFEs while not having better test accuracy.

²https://pytorch.org/docs/stable/generated/torch.cuda.max_memory_allocated.html

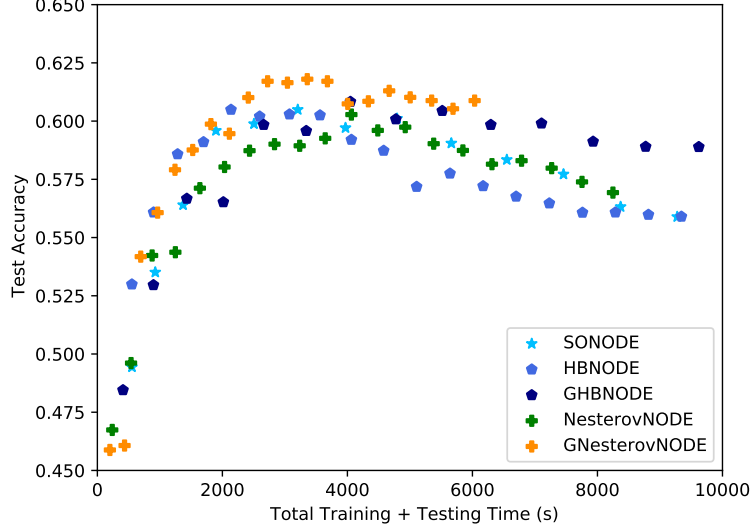


Figure 16: Test accuracy vs the corresponding wall-clock time for SONODE, HBNODE, NesterovNODE, and GNesterovNODE.

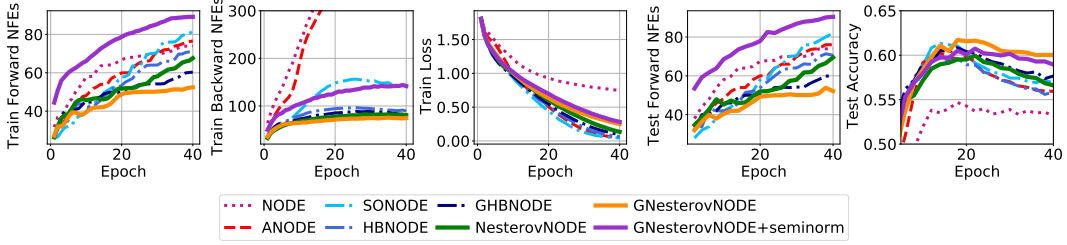


Figure 17: CIFAR10 with the seminorm method.

It is an open question why this is the case, as the seminorm method [Kidger et al., 2021] has shown to be consistently useful in reducing backward NFEs. The result is in Fig. 17.

E.6 Human Activity Dataset

We verify the advantage of the GNesterovNODE over the baseline NODE and GHBNODE for time series data on the Human Activity dataset [Luštrek et al., 2010]. This dataset consists of time series from five individuals doing various activities: walking, sitting, lying, etc. The data contains 3D positions of tags attached to those individuals’ belt, chest and ankles (12 features in total). After preprocessing, the dataset has 6554 sequences of 211 time points. In this task, the model classifies each time point into one of seven types of activities (walking, sitting, etc.). We use the ODE-RNN architecture described in Section 4.5 of [Rubanova et al., 2019] as the baseline model, with dopri5 adaptive solver and tolerance 10^{-7} . For the GNesterovNODE, we use $\sigma = \tanh$, and ξ is a learnable scalar. Figure 18 and Table 8 show that GNesterovNODE-RNN achieves the best accuracy and the smallest NFEs in both forward and backward pass.

Table 8: Accuracy of NODE-RNN, GHBNODE-RNN and our method GNesterovNODE-RNN on the Human Activity benchmark (Per-time-point classification) [Luštrek et al., 2010] (Tolerance: 10^{-7}).

Method	Accuracy
NODE-RNN	0.829 ± 0.016
GHBNODE-RNN	0.838 ± 0.017
GNesterovNODE-RNN	0.840 ± 0.016

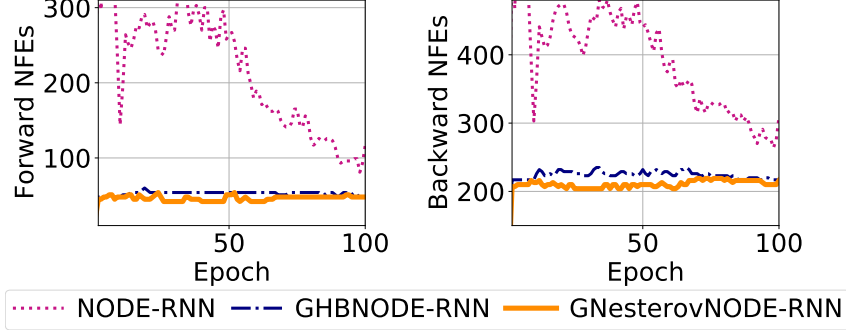


Figure 18: Contrasting the NFEs of NODE-RNN, GHBNODE-RNN and our GNesterovNODE-RNN on the Human Activity benchmark (Per-time-point classification) [Luštrek et al., 2010] (Tolerance: 10^{-7}).

702 E.7 Continuous Normalizing Flows for MNIST

703 We compare the GNesterovNODE with the baseline NODE and GHBNODE for use in variational
 704 inference with the continuous normalizing flow model trained on the MNIST dataset. The continuous
 705 normalizing flow model we use is the FFJORD in [Grathwohl et al., 2019a]. We summarize our
 706 results in Table 9 and Figure 19. Compared to the FFJORD-NODE and FFJORD-GHBNODE, the
 707 FFJORD-GNesterovNODE significantly reduces the NFEs in both forward and backward passes
 708 while improving the negative ELBO on the test set. Our training and the baseline FFJORD-NODE
 709 model follow the setting in Section 4.3 in [Grathwohl et al., 2019a]. We use the dopri5 solver with a
 710 tolerance of 10^{-5} . We use the identity function as σ , and $\xi = 2$ for GHBNODE and GNesterovNODE.

Table 9: Negative ELBO on the test data (lower is better) of FFJORD-NODE, FFJORD-HBNODE and our method FFJORD-GNesterovNODE for use in variational inference with a continuous normalizing flow model, i.e. FFJORD [Grathwohl et al., 2019a], trained on the binarized MNIST dataset. We also include the reported results from [Grathwohl et al., 2019a] (in parentheses) in addition to our reproduced results. (Tolerance: 10^{-5}).

Method	Negative ELBO
FFJORD-NODE	88.30 (82.82)
FFJORD-GHBNODE	75.87
FFJORD-GNesterovNODE	72.16

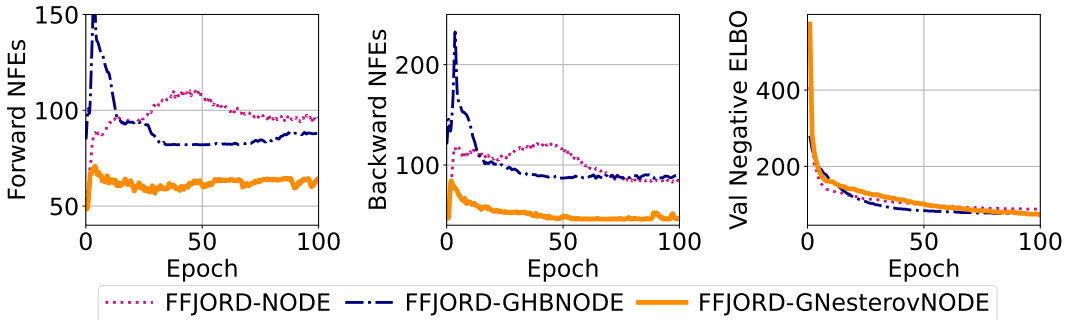


Figure 19: Contrasting the NFEs and the validation negative ELBO of the FFJORD-NODE, the FFJORD-HBNODE, and our FFJORD-GNesterovNODE for the variational inference task with a continuous normalizing flow model, i.e. FFJORD [Grathwohl et al., 2019a], on the binarized MNIST dataset (Tolerance: 10^{-5}).

711 **F Solving Neural Differential-Algebraic Equations**

712 Our implementation uses an ODE to calculate the DAE. An alternative approach is using DAE
713 solvers to solve the DAE, which are implemented by the `DifferentialEquations.jl` library
714 [Rackauckas & Nie, 2017; Rackauckas et al., 2019]. Although Dormand-Prince 5(4) [Dormand &
715 Prince, 1980] is a common choice for adaptive ODE solver, Tsitouras 5(4) [Tsitouras, 2011] is a more
716 efficient method, which the libraries `DifferentialEquations.jl` ³ [Rackauckas & Nie, 2017]
717 and `DiffraX` ⁴ [Kidger, 2022] have implemented.

³<https://github.com/SciML/DifferentialEquations.jl>

⁴<https://github.com/patrick-kidger/diffrax>