

SUPPLEMENTARY MATERIAL

Anonymous authors

Paper under double-blind review

CONTENTS

A Code Availability	2
B Proofs	2
B.1 Proof of Proposition 1	2
C Gradient Calculation Using Discrete Adjoint Methods	2
C.1 Deriving the discrete adjoint for a time-stepping algorithm	2
C.2 Deriving the gradient to parameters and incorporating integrals	3
C.3 Discrete adjoints of exemplar timestepping methods	3
C.3.1 Backward Euler	4
C.3.2 Theta methods	4
C.3.3 Explicit Runge–Kutta methods	4
C.4 Checkpointing	4
D Details For Numerical Examples	5
D.1 Image classification	5
D.2 Continuous normalizing flow	5
D.3 Time series regression	6

A CODE AVAILABILITY

Our code is publicly available. For reviewing purposes, it is temporarily hosted at <https://github.com/pnode-dev/pnode>.

B PROOFS

B.1 PROOF OF PROPOSITION 1

Proposition 1. Assuming $\tilde{\lambda}_{n+1} = \lambda_{n+1}$, the local discrepancy between the continuous adjoint sensitivity and the discrete adjoint sensitivity for forward Euler with the same step size h is

$$\|\tilde{\lambda}_n - \lambda_n\| = h^2 \|\mathcal{H}(\mathbf{u}_n + \epsilon(\mathbf{u}_{n+1} - \mathbf{u}_n), \boldsymbol{\theta}, t_n) f(\mathbf{u}_n, \boldsymbol{\theta}, t_n)\| \|\lambda_{n+1}\|, \quad (1)$$

where \mathcal{H} is the Hessian of f and ϵ is a constant in $(0, 1)$.

Proof. Discretizing the continuous adjoint equation

$$\frac{d\tilde{\lambda}}{dt} = - \left(\frac{\partial f}{\partial \mathbf{u}} \right)^T \tilde{\lambda} \quad (2)$$

with forward Euler gives

$$\tilde{\lambda}_n = \tilde{\lambda}_{n+1} + h \left(\frac{\partial f(\mathbf{u}_{n+1}, \boldsymbol{\theta}, t_{n+1})}{\partial \mathbf{u}} \right)^T \tilde{\lambda}_{n+1}. \quad (3)$$

The discrete adjoint of the forward Euler is

$$\lambda_n = \lambda_{n+1} + h \left(\frac{\partial f(\mathbf{u}_n, \boldsymbol{\theta}, t_n)}{\partial \mathbf{u}} \right)^T \lambda_{n+1}. \quad (4)$$

With (3) and (4) and the assumption $\tilde{\lambda}_{n+1} = \lambda_{n+1}$, we have

$$\|\tilde{\lambda}_n - \lambda_n\| = h \left\| \left(\frac{\partial f(\mathbf{u}_{n+1}, \boldsymbol{\theta}, t_{n+1})}{\partial \mathbf{u}} \right)^T - \left(\frac{\partial f(\mathbf{u}_n, \boldsymbol{\theta}, t_n)}{\partial \mathbf{u}} \right)^T \right\| \|\lambda_{n+1}\|. \quad (5)$$

Using the Taylor's series expansion, we obtain

$$\begin{aligned} \left(\frac{\partial f(\mathbf{u}_{n+1}, \boldsymbol{\theta}, t_{n+1})}{\partial \mathbf{u}} \right)^T &= \left(\frac{\partial f(\mathbf{u}_n, \boldsymbol{\theta}, t_n)}{\partial \mathbf{u}} \right)^T + \mathcal{H}(\mathbf{u}_n + \epsilon(\mathbf{u}_{n+1} - \mathbf{u}_n)) (\mathbf{u}_{n+1} - \mathbf{u}_n) \\ &= \left(\frac{\partial f(\mathbf{u}_n, \boldsymbol{\theta}, t_n)}{\partial \mathbf{u}} \right)^T + \mathcal{H}(\mathbf{u}_n + \epsilon(\mathbf{u}_{n+1} - \mathbf{u}_n)) (hf(\mathbf{u}_n)), \end{aligned} \quad (6)$$

where \mathcal{H} is the Hessian of f and ϵ is a constant in $(0, 1)$.

Plugging (6) into (5), we obtain (1) and complete the proof. \square

C GRADIENT CALCULATION USING DISCRETE ADJOINT METHODS

C.1 DERIVING THE DISCRETE ADJOINT FOR A TIME-STEPPING ALGORITHM

Let us define the Lagrange multipliers $\lambda_n \in \mathbb{R}^{N_d}$, $n = 0, \dots, N$ and define the Lagrangian

$$\mathcal{L}(\boldsymbol{\eta}) = \phi(\mathbf{u}_N) - \lambda_0^T (\mathbf{u}_0 - \boldsymbol{\eta}) - \sum_{n=0}^{N-1} \lambda_{n+1}^T (\mathbf{u}_{n+1} - \mathcal{N}(\mathbf{u}_n)), \quad (7)$$

which depends on the initial state $\boldsymbol{\eta}$ (input data) of an ODE system. Taking the total derivative of Equation (7) with respect to $\boldsymbol{\eta}$ leads to

$$\frac{d\mathcal{L}}{d\boldsymbol{\eta}} = \lambda_0^T - \left(\frac{d\phi}{d\mathbf{u}}(\mathbf{u}_N) - \lambda_N^T \right) \frac{d\mathbf{u}_N}{d\boldsymbol{\eta}} - \sum_{n=0}^{N-1} \left(\lambda_n^T - \lambda_{n+1}^T \frac{d\mathcal{N}}{d\mathbf{u}}(\mathbf{u}_n) \right) \frac{d\mathbf{u}_n}{d\boldsymbol{\eta}}. \quad (8)$$

The discrete adjoint equation is defined as

$$\begin{aligned}\lambda_n &= \left(\frac{d\mathcal{N}}{d\mathbf{u}}(\mathbf{u}_n) \right)^T \lambda_{n+1}, \quad n = N-1, \dots, 0, \\ \lambda_N &= \left(\frac{d\phi}{d\mathbf{u}}(\mathbf{u}_N) \right)^T,\end{aligned}\tag{9}$$

in order to make the last two terms in (8) vanish so that the total derivative can be obtained without computing the forward sensitivities $\frac{d\mathbf{u}_n}{d\boldsymbol{\eta}}$.

C.2 DERIVING THE GRADIENT TO PARAMETERS AND INCORPORATING INTEGRALS

We consider a system of ODEs in the general form

$$\mathcal{M}\dot{\mathbf{u}} = f(t, \boldsymbol{\theta}, \mathbf{u}), \quad t \in [t_0, t_F],\tag{10}$$

where \mathcal{M} is a mass matrix. For explicit ODEs, \mathcal{M} is an identity matrix \mathcal{I} of size $N_d \times N_d$.

The loss function with an integral term can be written as

$$\mathcal{L} = \phi(\mathbf{u}(t_F)) + \int_{t_0}^{t_F} q(\mathbf{u}(t), t) dt.\tag{11}$$

To obtain the gradient of \mathcal{L} to parameters (weights of the neural network), we can extend the ODE system (10) by augmenting the state vector with the parameters and the integrand q , and we obtain a larger system,

$$\underline{\mathcal{M}}\dot{\underline{\mathbf{u}}} = \underline{F}(t, \underline{\boldsymbol{\theta}}, \underline{\mathbf{u}}), \quad t \in [t_0, t_F],\tag{12}$$

where

$$\underline{\mathcal{M}} = \begin{bmatrix} \mathcal{M} & & \\ & \mathbf{I}_{N_p \times N_p} & \\ & & 1 \end{bmatrix}, \quad \underline{\mathbf{u}} = \begin{bmatrix} \mathbf{u} \\ \boldsymbol{\theta} \\ q \end{bmatrix}, \quad \underline{F} = \begin{bmatrix} f \\ \mathbf{0}_{N_p \times 1} \\ r \end{bmatrix}.$$

The second equation indicates that the parameters are constant, and the last equation results from the following transformation of the integral:

$$r = \int_{t_0}^{t_F} q(\mathbf{u}(t), t) dt.\tag{13}$$

The initial condition for the augmented ODE system is $\underline{\boldsymbol{\eta}}_0 = [\boldsymbol{\eta} \ \boldsymbol{\theta} \ 0]^T$.

For notional brevity, we denote the Jacobian with respect to \mathbf{u} and $\boldsymbol{\theta}$ as $\mathbf{f}_{\mathbf{u}}$ and $\mathbf{f}_{\boldsymbol{\theta}}$, respectively. Then the extended Jacobian can be written as

$$\underline{F}_{\underline{\mathbf{u}}} = \begin{bmatrix} \mathbf{f}_{\mathbf{u}} & \mathbf{f}_{\boldsymbol{\theta}} & \mathbf{0}_{N_d \times 1} \\ \mathbf{0}_{N_p \times N_d} & \mathbf{0}_{N_p \times N_p} & \mathbf{0}_{N_p \times 1} \\ r_{\mathbf{u}} & r_{\boldsymbol{\theta}} & 0 \end{bmatrix}.$$

Now the adjoint variable expands to the combination of three variables, corresponding to the partial derivative of the loss function with respect to the initial system state, the parameters, and the initial value of q , respectively. In the later part of this document, we will use $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$ to denote the first and second adjoint variables. The terminal conditions for these two variables are

$$\boldsymbol{\lambda}_N = \left(\frac{\partial \phi}{\partial \mathbf{u}}(\mathbf{u}_N) \right)^T, \quad \boldsymbol{\mu}_N = \left(\frac{\partial \phi}{\partial \boldsymbol{\theta}}(\mathbf{u}_N) \right)^T.\tag{14}$$

The third adjoint variable has a constant value 1 because of the zeros in the last column of $\underline{F}_{\underline{\mathbf{u}}}$ and thus does not need to be computed.

C.3 DISCRETE ADJOINTS OF EXEMPLAR TIMESTEPING METHODS

This section lists the discrete adjoints for three families of timestepping methods that are widely used. Note that these discrete adjoints are derived and implemented for only selective timestepping methods in PETSC, which is also a limitation in this work. Nevertheless, the discrete adjoint methods can be easily extended to other timestepping methods in PETSC.

C.3.1 BACKWARD EULER

Forward propagation

$$\mathcal{M}u_{n+1} = \mathcal{M}u_n + h_n \theta f(u_{n+1}) \quad (15)$$

Adjoint propagation

$$\begin{aligned} \mathcal{M}^T \lambda_s &= \lambda_{n+1} + h_n f_u^T(u_{n+1}) \lambda_s + h_n r_u^T(u_{n+1}), \\ \lambda_n &= \lambda_{n+1} + h_n f_u^T(u_{n+1}) \lambda_s + h_n r_u^T(u_{n+1}), \\ \mu_n &= \mu_{n+1} + h_n f_\theta^T(u_{n+1}) \lambda_s + h_n r_\theta^T(u_{n+1}) \end{aligned}$$

C.3.2 THETA METHODS

Forward propagation

$$\mathcal{M}u_{n+1} = \mathcal{M}u_n + h_n(1 - \theta)f(u_n) + h_n\theta f(u_{n+1}) \quad (16)$$

Adjoint propagation

$$\begin{aligned} \mathcal{M}^T \lambda_s &= \lambda_{n+1} + h_n \theta f_u^T(u_{n+1}) \lambda_s + h_n \theta r_u^T(t_{n+1}, u_{n+1}), \\ \lambda_n &= \mathcal{M}^T \lambda_s + h_n(1 - \theta) f_u^T(u_n) \lambda_s + h_n(1 - \theta) r_u^T(t_n, u_n), \\ \mu_n &= \mu_{n+1} + h_n \theta (f_p^T(u_{n+1}) \lambda_s + r_p^T(u_{n+1})) + h_n(1 - \theta) (f_p^T(u_n) \lambda_s + r_p^T(u_n)) \end{aligned} \quad (17)$$

C.3.3 EXPLICIT RUNGE-KUTTA METHODS

Forward propagation

$$\begin{aligned} U_i &= u_n + \sum_{j=1}^{i-1} h_n a_{ij} F(U_j), \quad i = 1, \dots, s, \\ u_{n+1} &= u_n + \sum_{i=1}^s h_n b_i F(U_i) \end{aligned} \quad (18)$$

Adjoint propagation

$$\begin{aligned} \lambda_{s,i} &= h_n f_u^T(U_i) \left(b_i \lambda_{n+1} + \sum_{j=i+1}^s a_{ji} \lambda_{s,j} \right) + h_n b_i r_u^T(U_i), \quad i = s, \dots, 1, \\ \mu_{s,i} &= h_n f_p^T(U_i) \left(b_i \lambda_{n+1} + \sum_{j=i+1}^s a_{ji} \lambda_{s,j} \right) + h_n b_i r_p^T(U_i), \quad i = s, \dots, 1, \\ \lambda_n &= \lambda_{n+1} + \sum_{j=1}^s \lambda_{s,j}, \\ \mu_N &= \mu_{n+1} + \sum_{j=1}^s \mu_{s,j} \end{aligned} \quad (19)$$

C.4 CHECKPOINTING

PETSc offers transparent and optimal checkpointing strategies on high-performance computing platforms. All of the strategies are accessible through our PNODE framework. In addition to the classic Revolve algorithm Griewank & Walther (2000), PETSc provides extended algorithms Zhang & Constantinescu (2021) that are tailored for multistage time integration methods. Figure 1 illustrates an optimal schedule for the adjoint calculation given a memory budget for storing 3 checkpoints, with each checkpoint consisting of a solution vector (numbered circles) and the stage values (black dots). This schedule requires a minimal number (6) of extra recomputations in order to reverse 10 time steps. PETSc also supports many other checkpointing algorithms, such as online algorithms Stumm & Walther (2010); Wang et al. (2009), multistage algorithms Stumm & Walther (2009), and multilevel algorithms Schanen et al. (2016) for heterogeneous platforms.

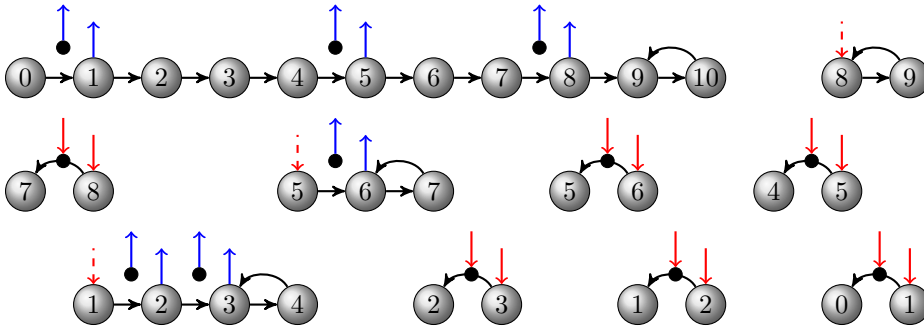


Figure 1: Modified Revolve algorithm for multistage time integration methods. Reprinted from Zhang & Constantinescu (2021). A numbered circle stands for the solution at the corresponding time index. A block dot stands for the stage values associated with the time step. The up arrow and down arrow stand for “store” operation and “restore” operation, respectively. When a stack is used for holding the checkpoints, the arrows with solid lines correspond to push and pop operations. The down arrow with dashed line indicates to read the top element on the stack without removing it.

D DETAILS FOR NUMERICAL EXAMPLES

D.1 IMAGE CLASSIFICATION

We modify a SqueezeNext network into a neural ODE by replacing all nonactivation blocks with ODE blocks. Our experiment is based on model files from the ANODE Gholaminejad et al. (2019) code repository¹. All models are trained with the SGD optimizer for 200 epochs, with an initial learning rate 0.1 decayed by a factor of 10 at the 150th epoch.

To further demonstrate the stability and accuracy of PNODE, we train SqueezeNext models using one-step forward Euler and RK4 for 3 independent runs, and we report the mean and standard deviation (STD) across the runs. As shown in Table 1, training with the discrete adjoint reaches lower training loss and higher test accuracy when compared with the neural ODE implementation using the continuous adjoint method (NODE cont).

Training loss				
	NODE disc	ANODE	NODE cont	PNODE
Euler	0.0006 ± 0.0000	0.0008 ± 0.0002	0.0030 ± 0.0002	0.0006 ± 0.0001
RK4	0.0005 ± 0.0000	0.0006 ± 0.0001	0.0006 ± 0.0000	0.0005 ± 0.0001
Test accuracy				
	NODE disc	ANODE	NODE cont	PNODE
Euler	0.8830 ± 0.0085	0.8740 ± 0.0157	0.5303 ± 0.0702	0.8780 ± 0.0036
RK4	0.8803 ± 0.0032	0.8850 ± 0.0050	0.8780 ± 0.0010	0.8803 ± 0.0068

Table 1: Mean and STD of training loss and test accuracy on classification on CIFAR-10 dataset after 200 epochs of training across 3 runs with different random seeds.

D.2 CONTINUOUS NORMALIZING FLOW

We apply PNODE on FFJORD, which is a free-form continuous generative model, and we follow the architecture and hyperparameter settings chosen by the authors in Grathwohl et al. (2019). The details are given in Table 2.

¹<https://github.com/amirgholami/anode>

Dataset	Dimension	Activation	Number of layers in NN	Hidden dim. multiplier	Number of ODE blocks	Batchsize
POWER	6	tanh	3	10	5	10000
BSDS300	63	softplus	3	20	2	10000 (1000*)
MINIBOONE	43	softplus	2	10	1	1000

Table 2: Information of hyperparameters chosen by Grathwohl et al. (2019). *The batch size is reduced to 1000 when using the discrete adjoint, so the memory can fit in the GPU.

D.3 TIME SERIES REGRESSION

We use the same model as in Rackauckas et al. (2020) and Onken & Ruthotto (2020):

$$\frac{d}{dt} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = A \begin{bmatrix} x_1^3 \\ x_2^3 \end{bmatrix}, \quad A = \begin{bmatrix} -0.10 & 2.00 \\ -2.00 & -0.10 \end{bmatrix},$$

with an initial condition $\mathbf{x} = [2, 0]$. The data are generated on 160 equidistant time grid points in $[0, 16]$ by integrating the ODE using an adaptive dopri5 solver. For the training model, we use a single linear layer with 2D input and output applied to the elementwise cubic of x ; therefore, the true underlying model can be represented exactly. We train the model on $t \in [0, 1]$ with a step size of 0.1. In each iteration, we use an equal batch size of 20. The loss function is defined as the average of the absolute values of the residuals. The parameters in the model are updated through an RMSprop minimizer with a learning rate of 0.01. The test loss is then evaluated on all the samples.

Here we provide the comparison in terms of loss functions using two second-order schemes—explicit midpoint and Crank–Nicolson (implicit)—using three different random seeds, as shown in Figure 2. The solution using explicit solver blows up during tests for all three runs, while training with implicit solvers always reaches a lower training loss and the solutions always remain finite during testing.

REFERENCES

- Amir Gholaminejad, Kurt Keutzer, and George Biros. ANODE: Unconditionally accurate memory-efficient gradients for neural ODEs. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pp. 730–736. International Joint Conferences on Artificial Intelligence Organization, 7 2019. doi: 10.24963/ijcai.2019/103.
- 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.
- Andreas Griewank and Andrea Walther. Algorithm 799: Revolve: An implementation of checkpointing for the reverse or adjoint mode of computational differentiation. *ACM Trans. Math. Softw.*, 26(1):19–45, March 2000. ISSN 0098-3500. doi: 10.1145/347837.347846.
- Derek Onken and Lars Ruthotto. Discretize-optimize vs. optimize-discretize for time-series regression and continuous normalizing flows. *arXiv e-prints*, 2020.
- Christopher Rackauckas, Yingbo Ma, Julius Martensen, Collin Warner, Kirill Zubov, Rohit Supekar, Dominic Skinner, and Ali Ramadhan. Universal differential equations for scientific machine learning. *arXiv preprint arXiv:2001.04385*, 2020.
- Michel Schanen, Oana Marin, Hong Zhang, and Mihai Anitescu. Asynchronous two-level checkpointing scheme for large-scale adjoints in the spectral-element solver Nek5000. *Procedia Computer Science*, 80: 1147–1158, 2016. ISSN 18770509. doi: 10.1016/j.procs.2016.05.444.
- Philipp Stumm and Andrea Walther. MultiStage Approaches for Optimal Offline Checkpointing. *SIAM Journal on Scientific Computing*, 31(3):1946–1967, 2009. ISSN 1064-8275. doi: 10.1137/080718036.
- Philipp Stumm and Andrea Walther. New algorithms for optimal online checkpointing. *SIAM Journal on Scientific Computing*, 32(2):836–854, 2010. ISSN 1064-8275. doi: 10.1137/080742439.

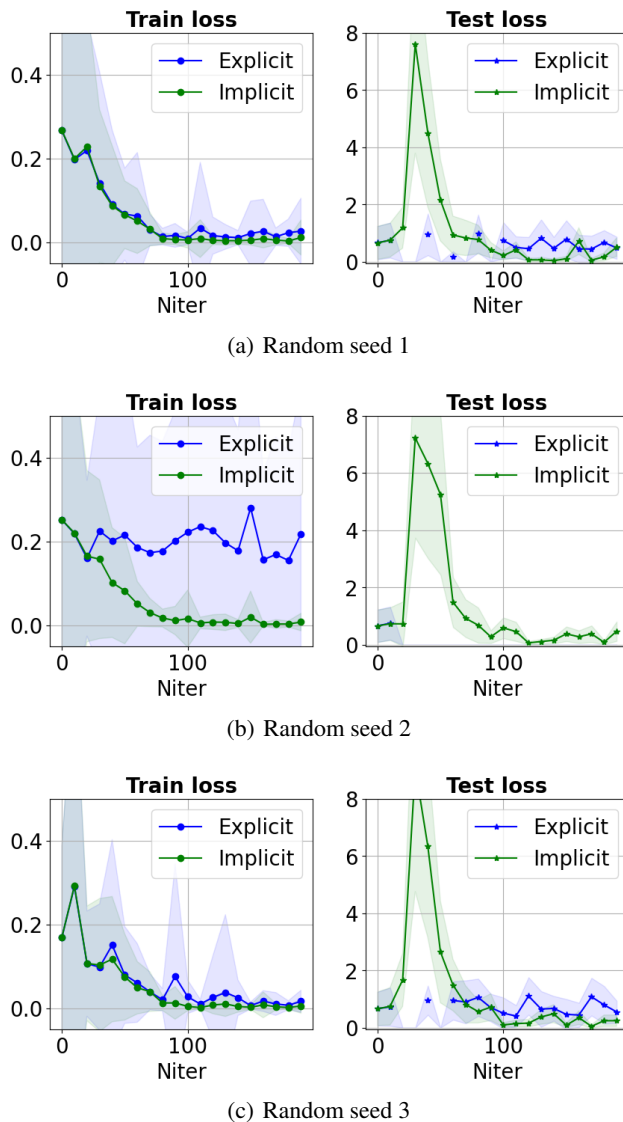


Figure 2: Train and test loss with explicit midpoint and Crank–Nicolson (implicit) solvers as functions of number of iterations using different random seeds.

Qiqi Wang, Parviz Moin, and Gianluca Iaccarino. Minimal repetition dynamic checkpointing algorithm for unsteady adjoint calculation. *SIAM Journal on Scientific Computing*, 31(4):2549–2567, 2009. ISSN 1064-8275. doi: 10.1137/080727890.

Hong Zhang and Emil M. Constantinescu. Revolve-based adjoint checkpointing for multistage time integration. In *International Conference on Computational Science*, (online), in main track, 2021.