

# “Hey, that’s not an ODE”: Faster ODE Adjoint with 12 Lines of Code

Patrick Kidger, Ricky T. Q. Chen, Terry Lyons



## SUMMARY

We train neural ODEs, using the adjoint method, much faster: with a median of 40% fewer function evaluations.

This is done by observing that the backward (adjoint) equations exhibit a particular structure, which we may adjust the ODE solver to exploit. This change is simple to implement: at least using `torchdiffEq`, it requires only 12 lines of code.

## BACKGROUND: NEURAL ODES

- Differential equations with neural network vector fields.
- Use ODEs as components of a differentiable computation graph.

For example: supervised learning of the map  $x \mapsto y$  via

$$z(0) = \ell_{\theta}^1(x) \quad z(t) = z(0) + \int_0^t f_{\theta}(z(s)) ds \quad y \approx \ell_{\theta}^2(z(T))$$

with  $\ell_{\theta}^1, \ell_{\theta}^2$  linear, and  $f_{\theta}$  a neural network.

## BACKGROUND: ADJOINT EQUATIONS

Neural ODEs may be backpropagated through via the adjoint equations. For a loss function  $L$ , then  $\frac{dL}{d\theta} = a_{\theta}(0)$  and  $\frac{dL}{dz(0)} = a_z(0)$ , where

$$a_z(T) = \frac{dL}{dz(T)}, \quad a_z(t) = a_z(T) - \int_T^t a_z(s) \cdot \frac{\partial f}{\partial z}(s, z(s), \theta) ds,$$

$$a_{\theta}(T) = 0, \quad a_{\theta}(t) = a_{\theta}(T) - \int_T^t a_z(s) \cdot \frac{\partial f}{\partial \theta}(s, z(s), \theta) ds.$$

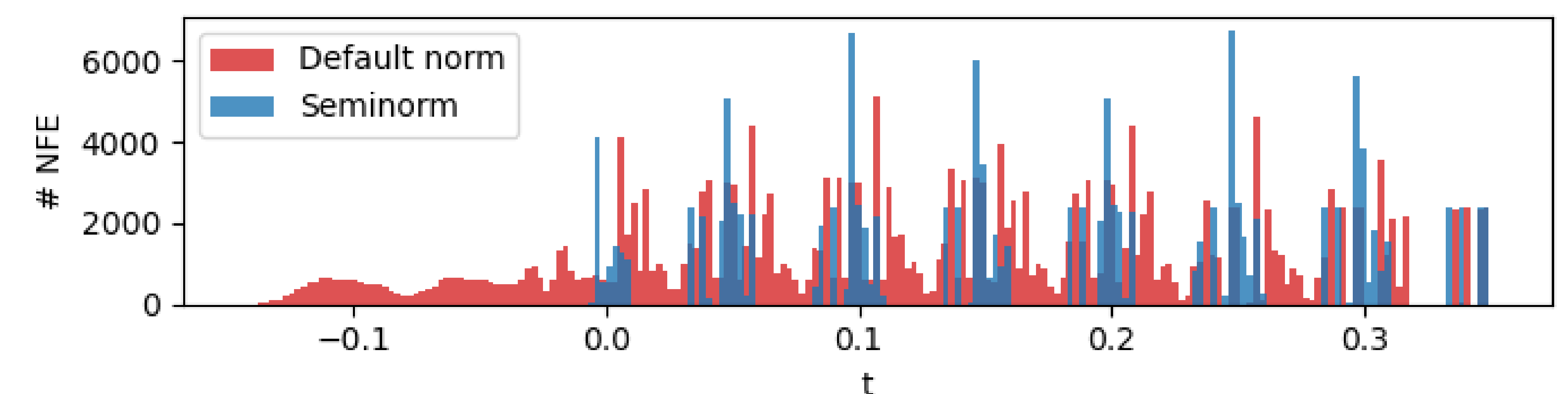
## METHOD

$a_{\theta}$  does not appear on the right hand side of the adjoint equations. Thus it is just an integral—not an ODE. Using an ODE solver, it assumes that errors now propagate to create errors later, and takes steps small enough to compensate—in this case taking too many steps.

By switching the norm (used to measure error) for a seminorm that does not see  $a_{\theta}$ , then many fewer steps are needed.

## RESULTS

Hamiltonian neural network: histogram of step locations.



Neural controlled differential equation: NFE for varying tolerances

RTOL, ATOL	Default norm		Seminorm	
	Accuracy (%)	Bwd. NFE ( $10^6$ )	Accuracy (%)	Bwd. NFE ( $10^6$ )
$10^{-3}, 10^{-6}$	92.6±0.4	14.36±1.09	92.5±0.5	<b>8.67±1.60</b>
$10^{-4}, 10^{-7}$	92.8±0.4	30.67±2.48	92.5±0.5	<b>12.75±2.00</b>
$10^{-5}, 10^{-8}$	92.4±0.7	77.95±4.47	92.9±0.4	<b>29.39±0.80</b>

## FIND OUT MORE:

Full paper: <https://arxiv.org/abs/2009.09457>

Code: <https://github.com/patrick-kidger/FasterNeuralDiffEq>

Email: [kidger@maths.ox.ac.uk](mailto:kidger@maths.ox.ac.uk)