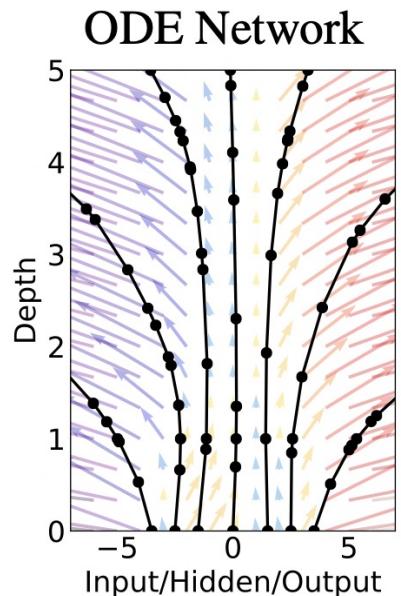
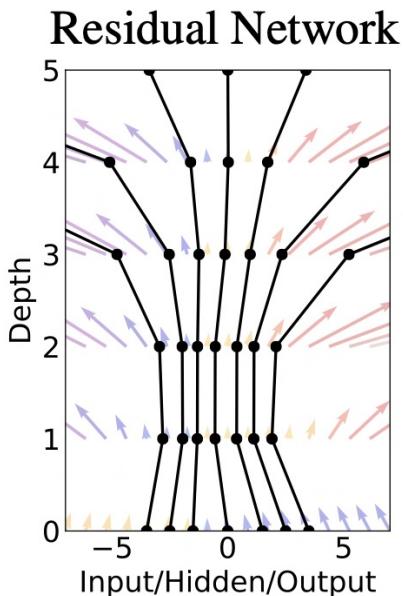


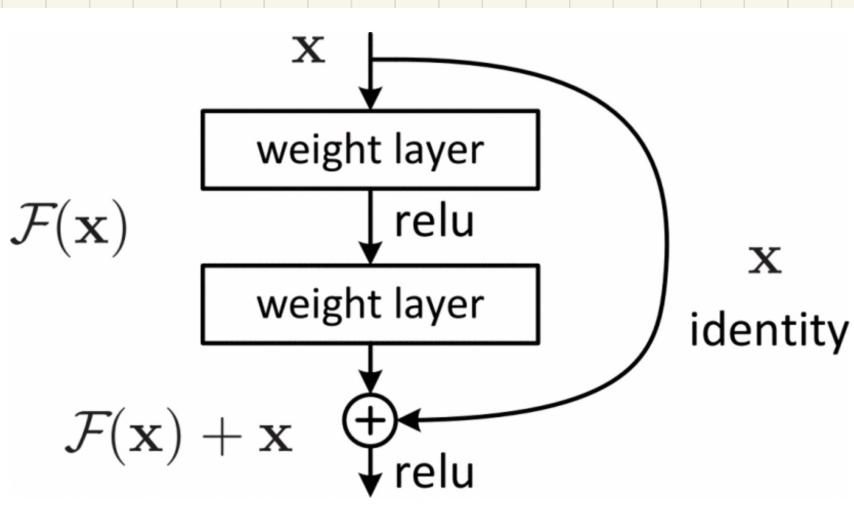
[Submitted on 19 Jun 2018 (v1), last revised 14 Dec 2019 (this version, v5)]

Neural Ordinary Differential Equations

Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, David Duvenaud



Resnets



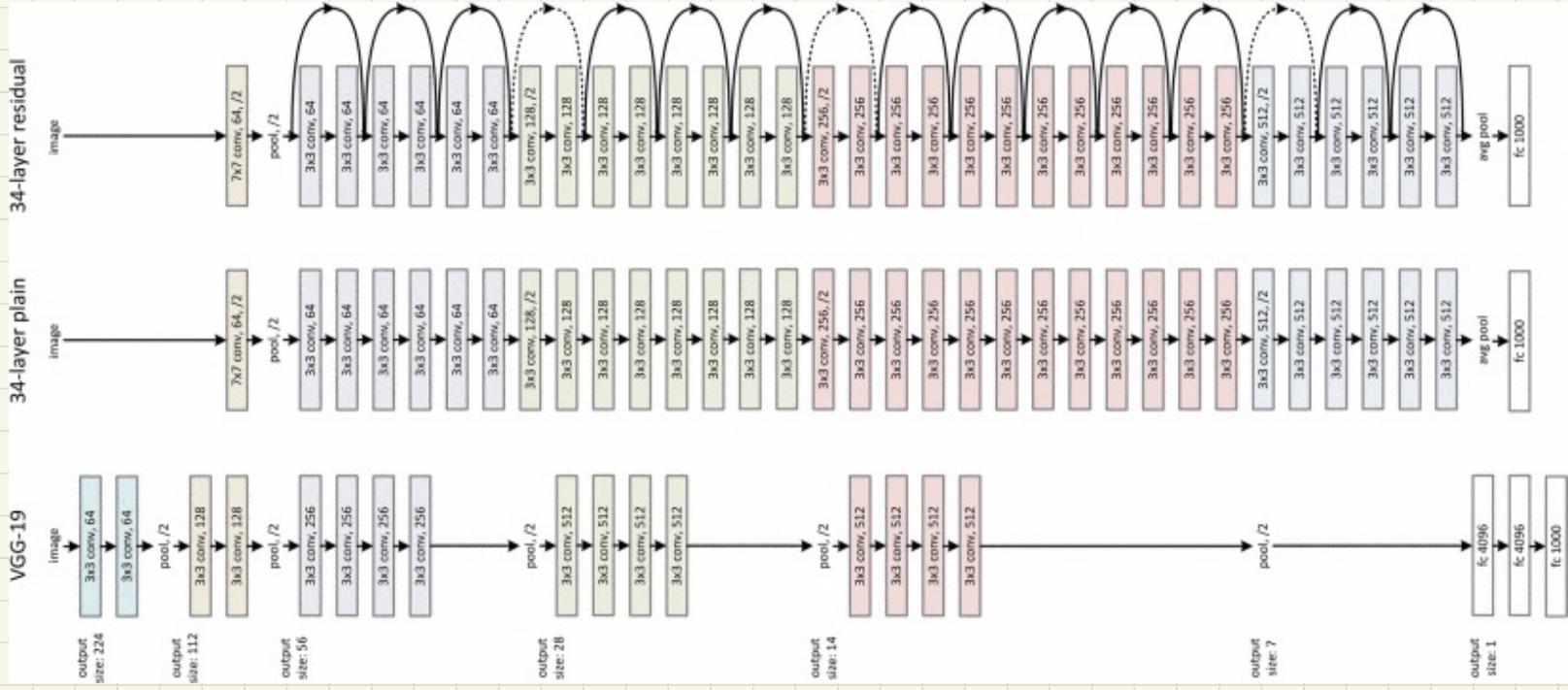
$$\text{Plain NN} = x_{n+1} = F(x_n)$$

$$\text{ResNet} = x_{n+1} = x_n + F(x_n)$$

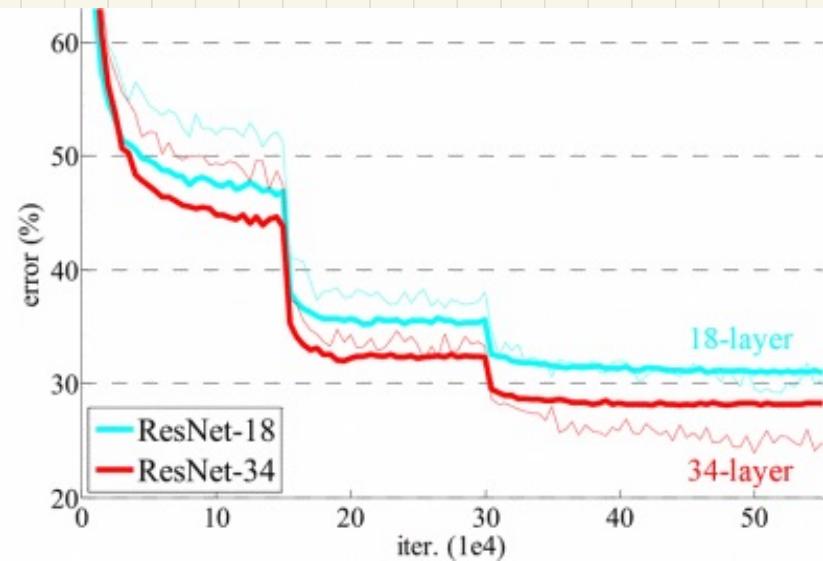
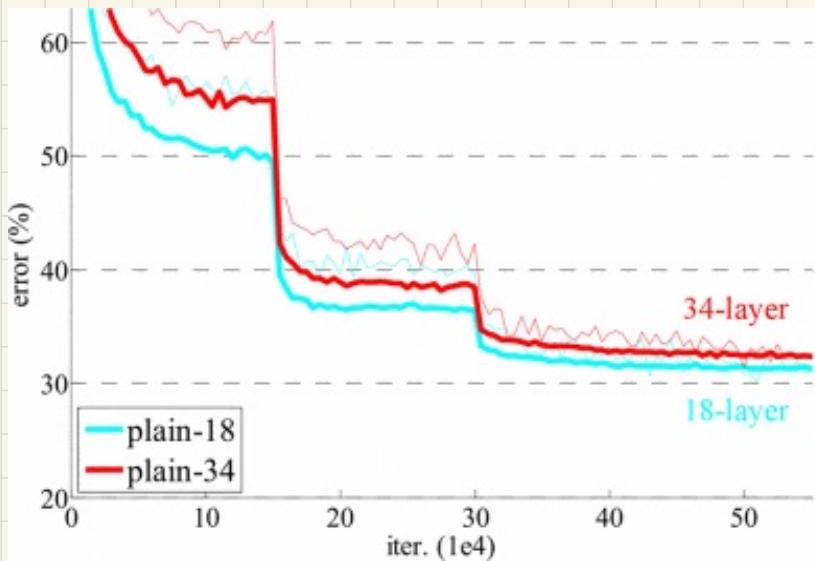
↑

residual

K. He, X. Zhang, S. Ren and J. Sun, "Deep Residual Learning for Image Recognition," 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 770-778, doi: 10.1109/CVPR.2016.90.



ImageNet results:



$$x_1 = f_1(x_0)$$

$$x_2 = f_2(x_1) = f_2(f_1(x_0))$$

$$\nabla_{\theta} x_2 = \nabla_{\theta} f_2(f_1(x_0)) + \underline{\nabla_x f_2} \cdot \underline{\nabla_{\theta} f_1(x_0)}$$

$$x_2 = x_1 + f_2(x_0 + f_1(x_0))$$

$$\nabla_{\theta} x_2 = \nabla_{\theta} x_1 + \nabla_x f_2 \cdot (\nabla_{\theta} x_0 + \nabla_{\theta} f_1)$$

$$\text{Resnets: } z_{n+k} = z_n + F_\theta(z_n)$$

$$\text{Neural ODE: } z_{n+1} = z_n + \Delta t f_\theta(z_n) \quad \leftarrow \text{Euler discretization of } \dot{z} = f_\theta(z)$$

$$\boxed{\dot{z} = f_\theta(z)}$$

↑
NN with params θ

$$\text{Solve: } z(t_1) = z(t_0) + \int_{t_0}^{t_1} f(z(t)) dt$$

$$L = L(z(t_0), z(t_1), z(t_2), \dots)$$

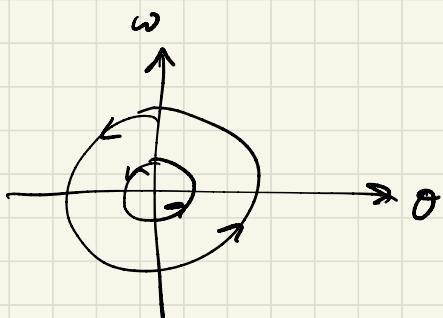
Ex pendulum (linearized)

$$\dot{\theta} = \omega$$

$$\ddot{\omega} = -\sin \theta$$

$$\dot{\theta} = \omega$$

$$\ddot{\omega} = -\theta$$



$$z = \begin{bmatrix} 0 \\ \omega \end{bmatrix}$$

$$\dot{z} = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} z$$

$$NODE: \dot{z} = \begin{bmatrix} \dot{\theta} \\ \dot{\omega} \end{bmatrix} = f_\theta(z) = f_\theta \left(\begin{bmatrix} 0 \\ \omega \end{bmatrix} \right)$$

$$= f_\theta \left(\begin{bmatrix} 0 \\ \omega \end{bmatrix} \right)$$

Adjoints - {
 x state $\in \mathbb{R}^d$
 $Ax = b$ $A = A(\theta)$ $b = b(\theta)$
 loss $L(x)$
 find θ to minimize $L(x)$

$$\frac{d}{d\theta} L = \nabla_x L \cdot \nabla_\theta x$$

$$\begin{matrix} \uparrow \\ \mathbb{R} \end{matrix} = L_x x_\theta$$

$$\frac{d}{d\theta} L = L_x A^{-1} (b_\theta - A_\theta x)$$

$1 \times p$ $1 \times d$ $d \times d$ $d \times p$

$O(d^2 p)$
 $\underbrace{\lambda^T}_{1 \times d}$

$$Ax = b$$

$$A_\theta x + A_\theta x_\theta = b_\theta$$

$$x_\theta = A^{-1} (b_\theta - A_\theta x)$$

$$\lambda^+ = L_x A^{-1}$$

$$\lambda = A^{-T} L_x^T$$

$$A^T \lambda = L_x^+$$

adjoint system

A B C

$d \times 1$ $|x_n$ $n \times 1$

$\underbrace{}$
 $O(dn)$

$\underbrace{}$

$O(dn)$

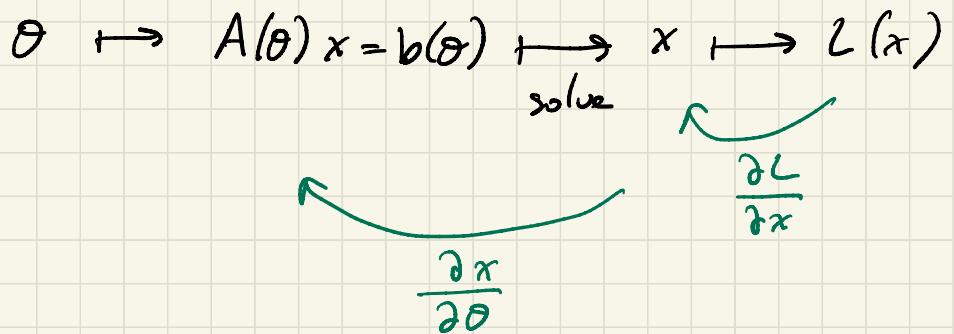
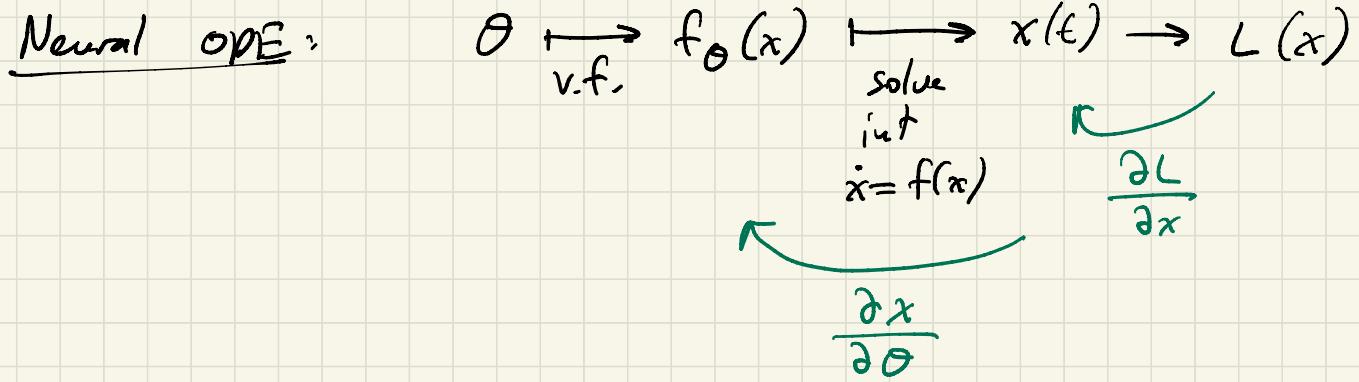
A B C

$d \times 1$ $|x_n$ $n \times 1$

$\underbrace{}$
 $O(n)$

$\underbrace{}$

$O(d+n)$



Linear

$$Ax = b \quad A = A(\theta) \quad b = b(\theta)$$

$$A_\theta x + A_\theta x_\theta = b_\theta$$

$$A_\theta x_\theta = (b_\theta - A_\theta x)$$

$$x_\theta = A^{-1}(b_\theta - A_\theta x)$$

want $L_\theta = L_x x_\theta$

$$= L_x \underbrace{A^{-1}}_{\lambda^T} (b_\theta - A_\theta x)$$

$A^T \lambda = L_x^T$

adjoint eqn

Nonlinear

$$f(x, \theta) = 0$$

$$f_x x_\theta + f_\theta = 0$$

$\curvearrowleft A$

$$x_\theta = -f_x^{-1} f_\theta$$

want $L_\theta = L_x x_\theta$

$$= -L_x \underbrace{f_x^{-1}}_{\lambda^T} f_\theta$$

$f_x^T \lambda = L_x^T$

Linear ODEs

$$\dot{x} = Bx \quad x(0) = b \quad t \in [0, T]$$

$$\bar{x} = x(T) = e^{BT}b$$

$$e^{-B^T} \bar{x} = b$$

↑
A

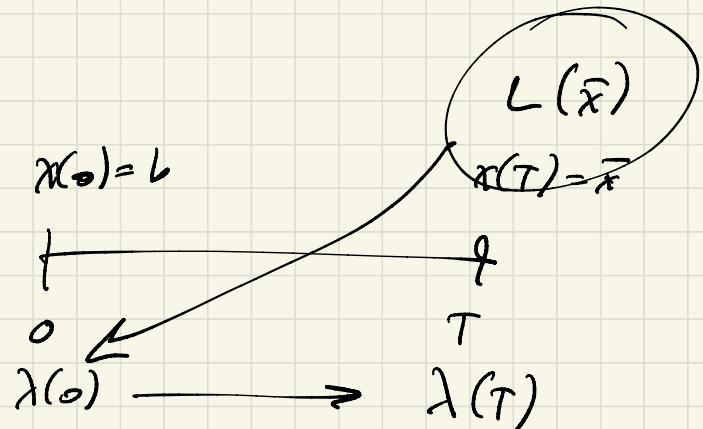
$$\text{Assume } L = L(\bar{x})$$

$$A^T \lambda = L_{\bar{x}}^T$$

$$e^{-B^T T} \lambda = L_{\bar{x}}^T$$

$$\lambda(T) = \lambda = e^{B^T T} L_{\bar{x}}^T$$

$$\dot{\lambda} = B^T \lambda \quad \lambda(0) = L_{\bar{x}}^T$$



$$\text{define } \varepsilon = -t \quad \frac{\partial}{\partial t} = \frac{\partial \varepsilon}{\partial t} \frac{\partial}{\partial \varepsilon} = -\frac{\partial}{\partial \varepsilon}$$

$$\frac{d}{d\varepsilon} \lambda = -B^T \lambda$$

$$\lambda(T) = L_{\bar{x}}^T$$

want $\lambda(0)$

Nonlinear ODEs

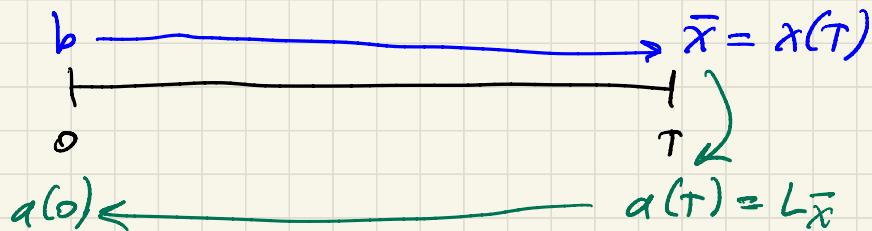
$$\dot{x} = f(x) \quad x(0) = b$$

linearize = $\delta \dot{x} = \underbrace{\frac{\partial f}{\partial x}(x(t))}_{B} \delta x$

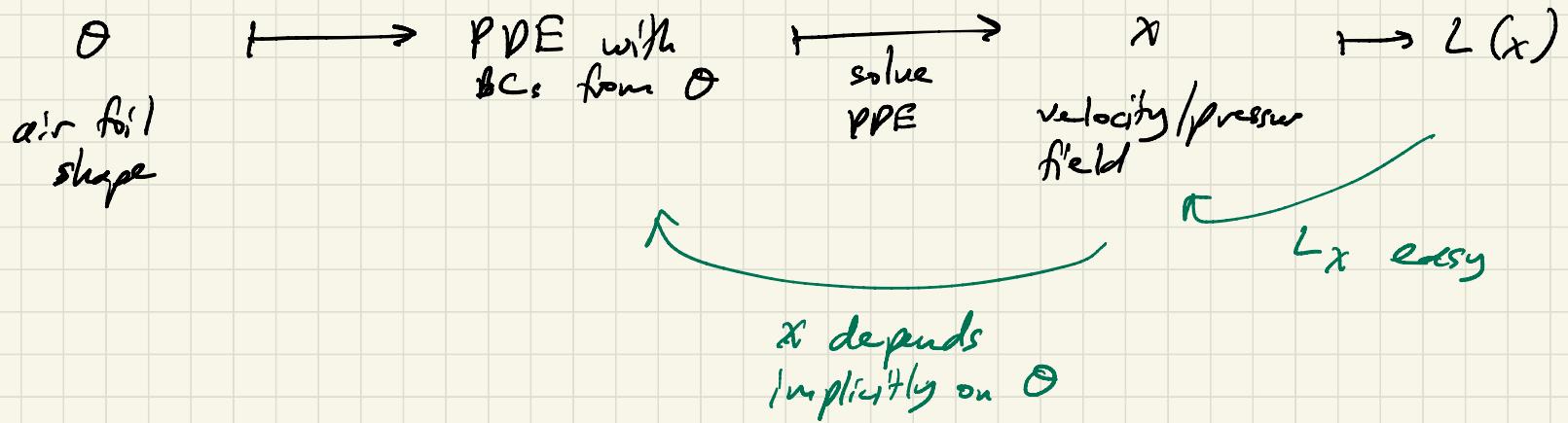
$$\Rightarrow \text{adjoint eqn: } \frac{d}{dt} \alpha = - \left(\frac{\partial f}{\partial x} \right)^T \alpha$$

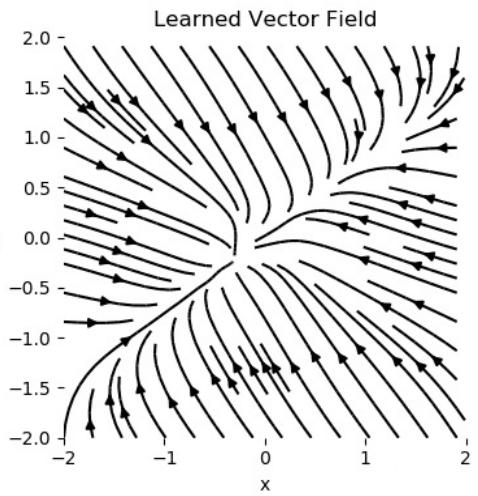
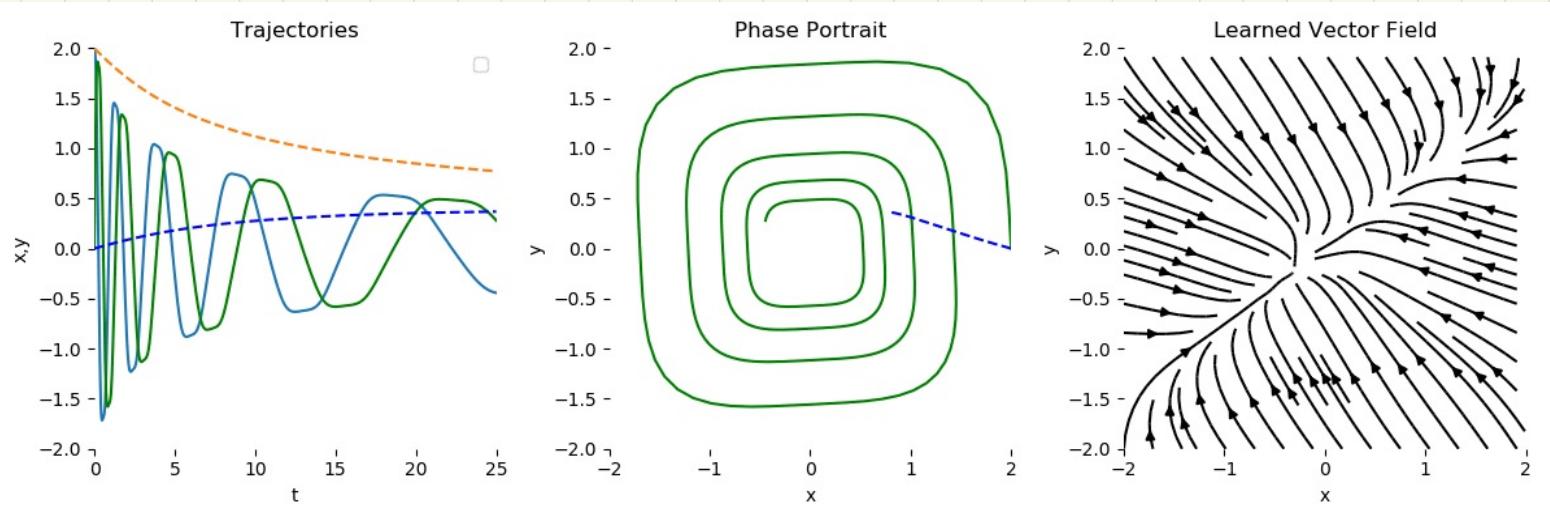
or $\alpha = \lambda^T \Rightarrow \frac{d}{dt} \alpha = - \alpha \frac{\partial f}{\partial x}(x(t))$

$$\alpha(T) = L_{\bar{x}}$$



PDE optimization (e.g. airfoils)





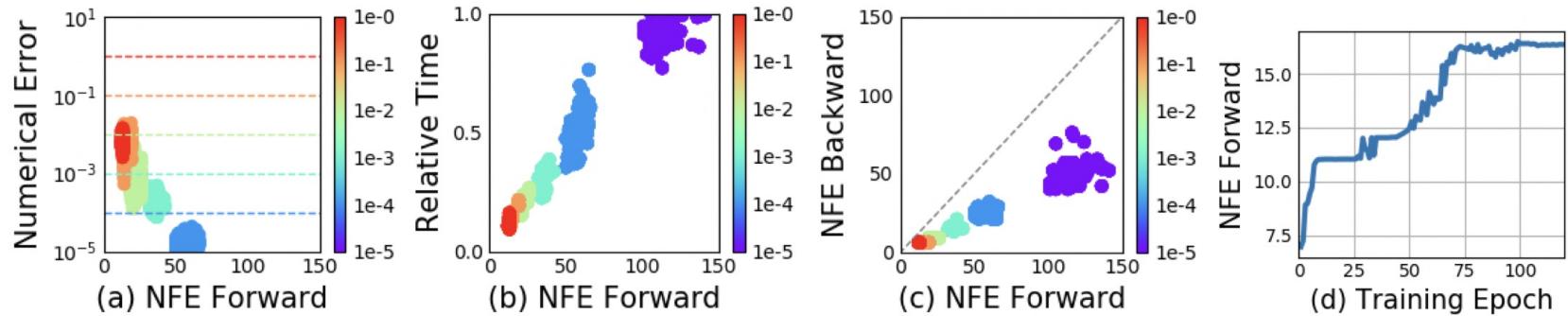
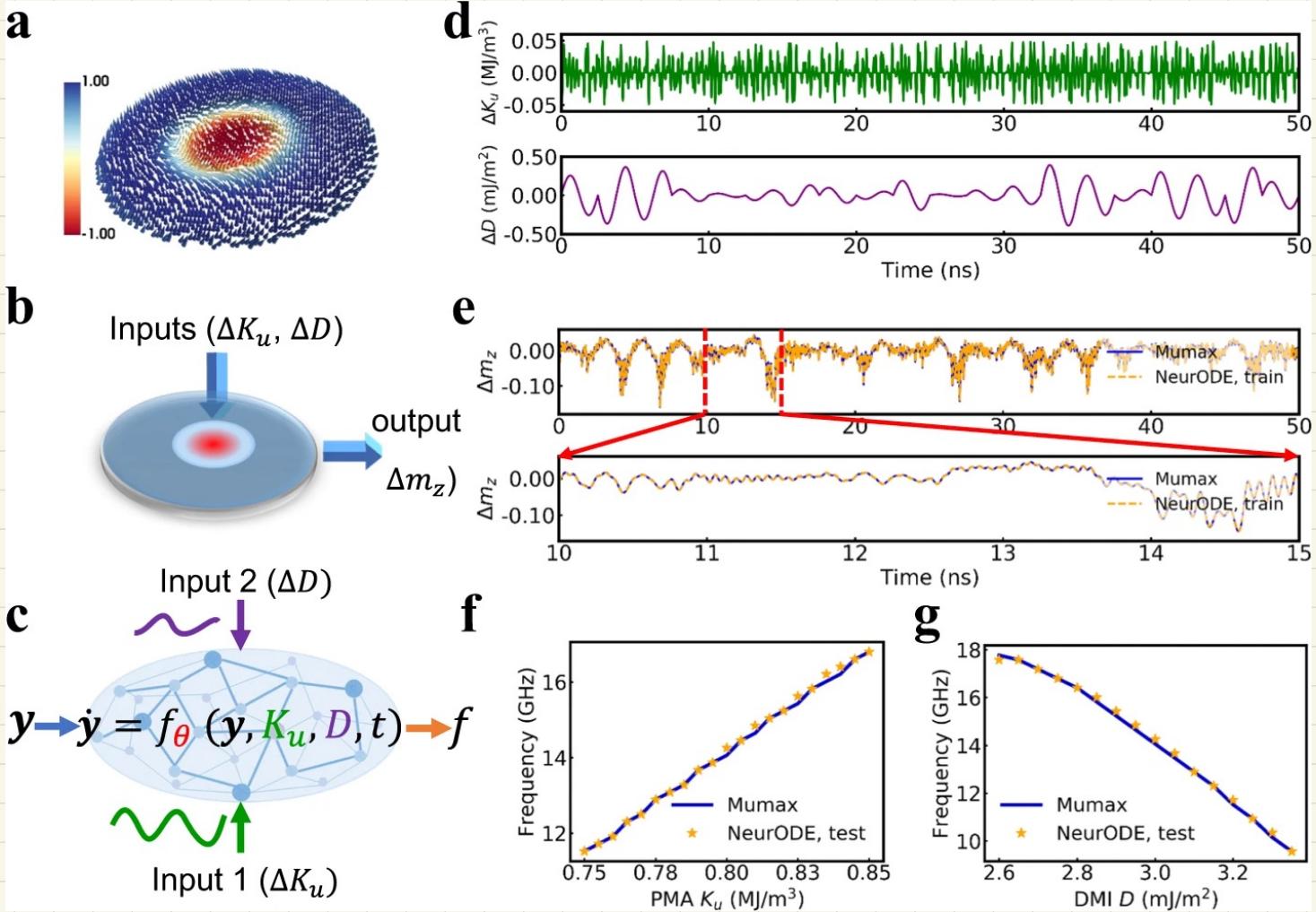


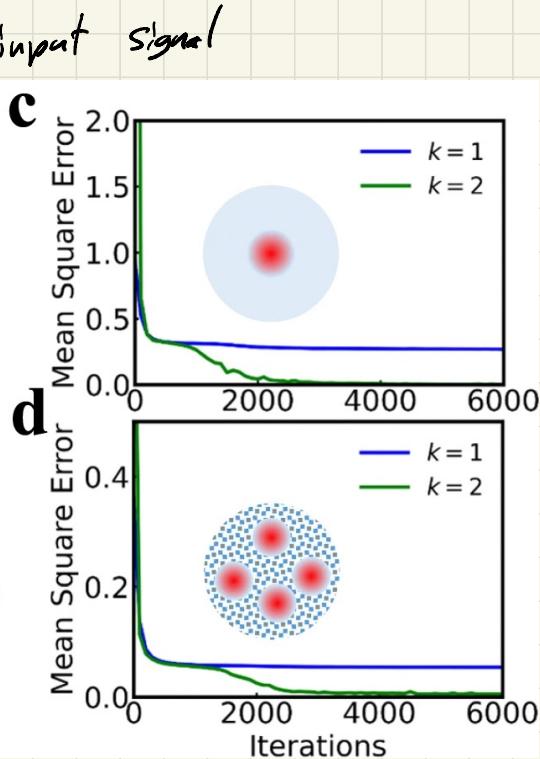
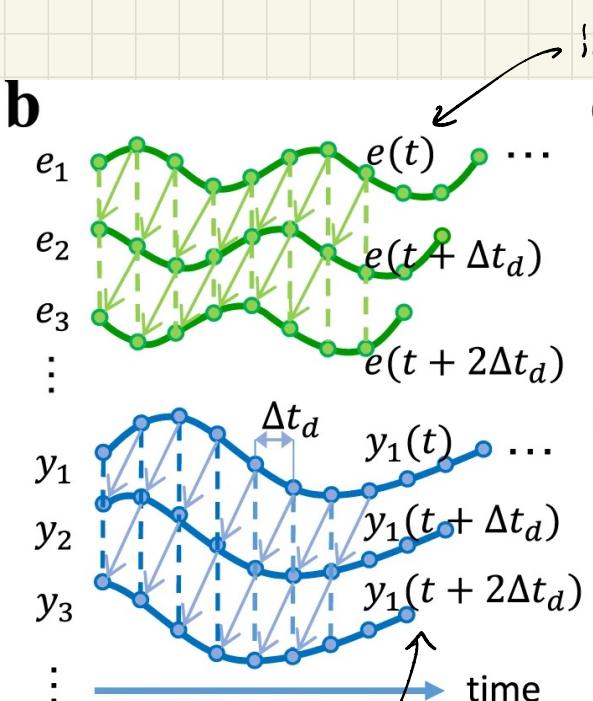
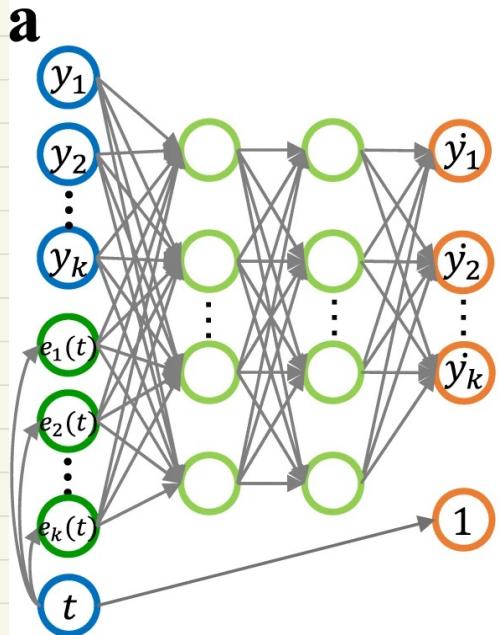
Figure 3: Statistics of a trained ODE-Net. (NFE = number of function evaluations.)

Forecasting the outcome of spintronic experiments with Neural Ordinary Differential Equations

[Xing Chen](#), [Flavio Abreu Araujo](#), [Mathieu Riou](#), [Jacob Torrejon](#), [Dafiné Ravelosona](#), [Wang Kang](#),
[Weisheng Zhao](#), [Julie Grollier](#) & [Damien Querlioz](#) 

[Nature Communications](#) **13**, Article number: 1016 (2022) | [Cite this article](#)





system state

$k = \# \text{ delayed states}$

$k=2 \Rightarrow (y_1(t), y_1(t+\Delta t))$

Time-delay embeddings

2D ODE

$$\begin{pmatrix} \dot{x} \\ \dot{y} \end{pmatrix} = f \begin{pmatrix} x \\ y \end{pmatrix}$$

observe only x_0, x_1, x_2, \dots

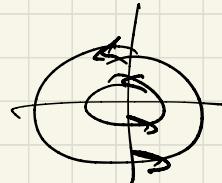
do NOT observe y_0, y_1, y_2, \dots

homeomorphic

Thus (Takens embedding) there exists an equivalent,

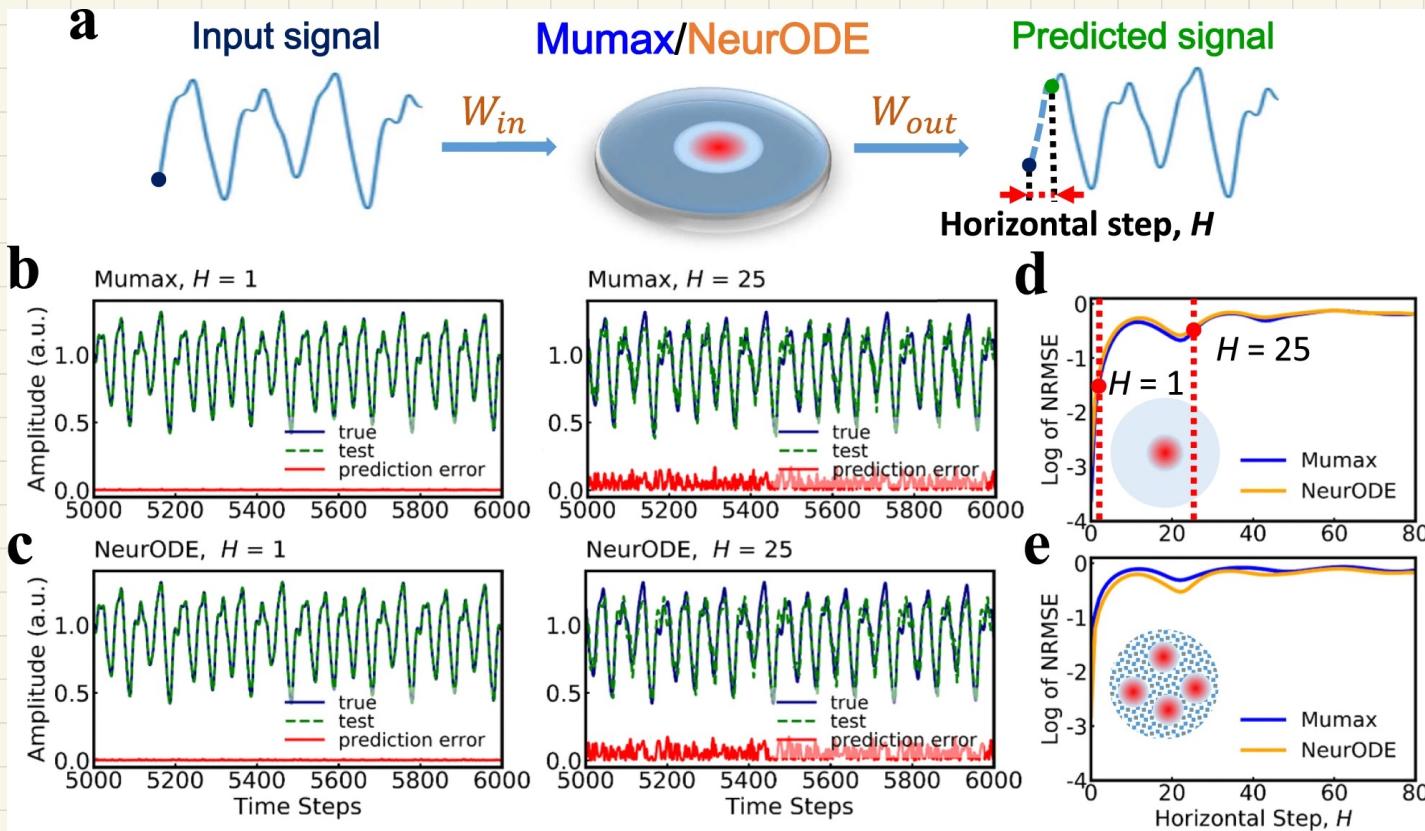
system $\begin{pmatrix} x_t \\ x_{t+1} \end{pmatrix} = \tilde{f} \begin{pmatrix} x_{t-1} \\ x_t \end{pmatrix}$

Ex $\begin{pmatrix} \dot{\theta} \\ \dot{\omega} \end{pmatrix} = \begin{pmatrix} \omega \\ -\sin \theta \end{pmatrix}$



$$\begin{pmatrix} \theta_t \\ \theta_{t+1} \end{pmatrix} = \tilde{f} \begin{pmatrix} \theta_{t-1} \\ \theta_t \end{pmatrix}$$

$$\omega_t \approx \frac{\theta_t - \theta_{t-1}}{st}$$



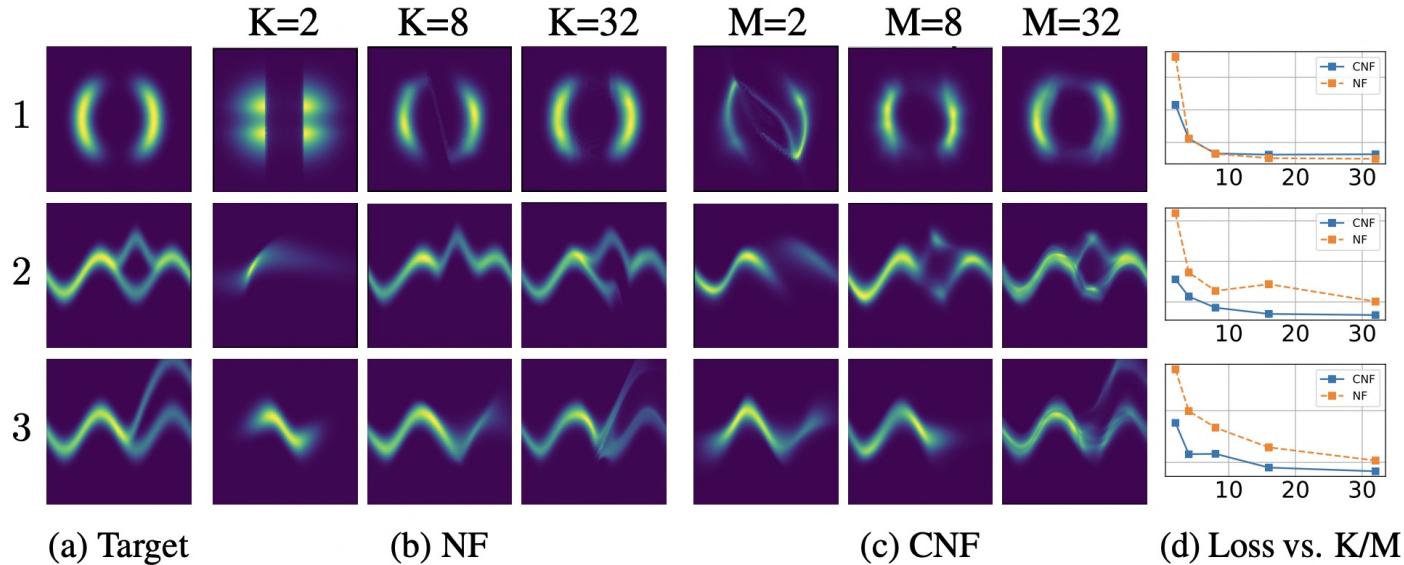
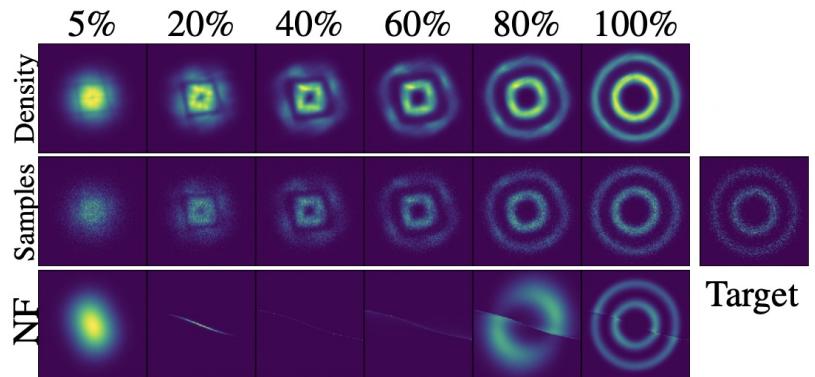
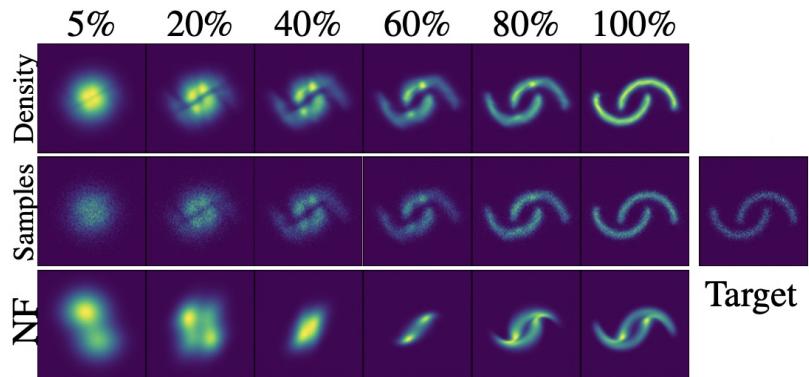


Figure 4: Comparison of normalizing flows versus continuous normalizing flows. The model capacity of normalizing flows is determined by their depth (K), while continuous normalizing flows can also increase capacity by increasing width (M), making them easier to train.



(a) Two Circles



(b) Two Moons

Figure 5: Visualizing the transformation from noise to data. Continuous-time normalizing flows are reversible, so we can train on a density estimation task and still be able to sample from the learned density efficiently.

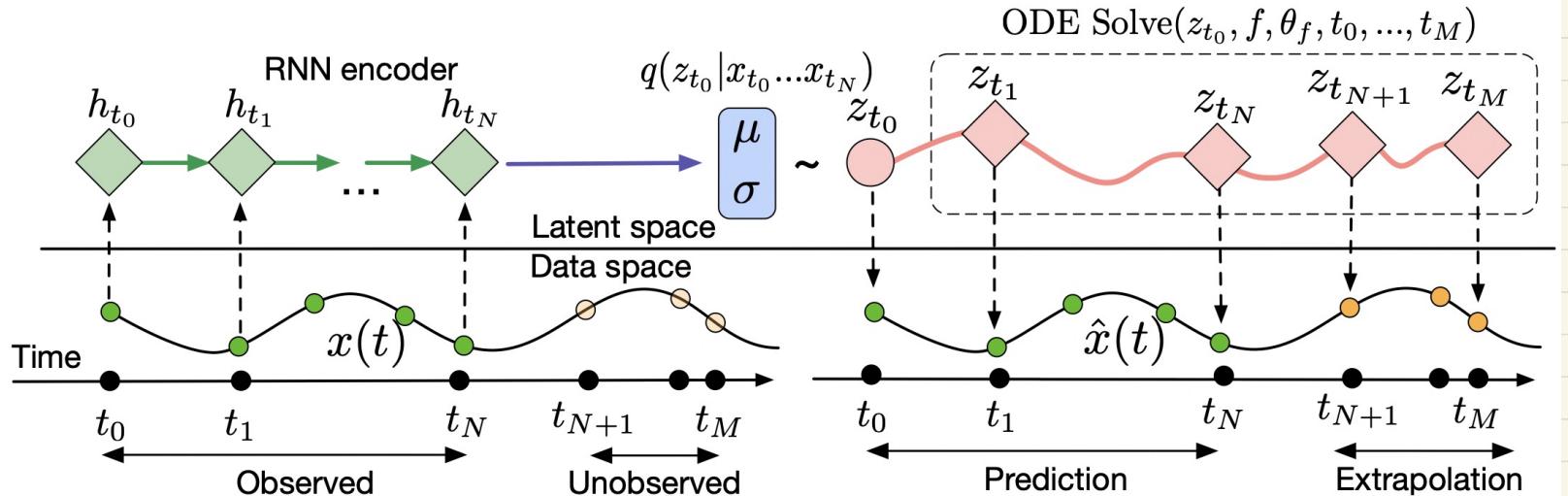
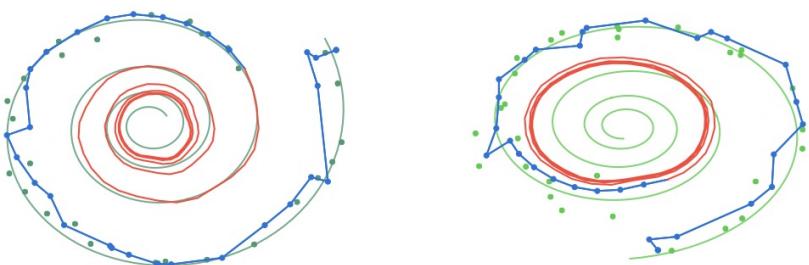
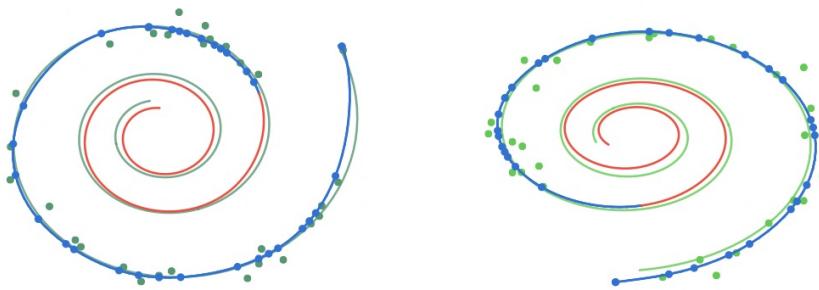


Figure 6: Computation graph of the latent ODE model.



(a) Recurrent Neural Network



(b) Latent Neural Ordinary Differential Equation

Algorithm 2 Complete reverse-mode derivative of an ODE initial value problem

Input: dynamics parameters θ , start time t_0 , stop time t_1 , final state $\mathbf{z}(t_1)$, loss gradient $\partial L / \partial \mathbf{z}(t_1)$

$$\frac{\partial L}{\partial t_1} = \left(\frac{\partial L}{\partial \mathbf{z}(t_1)} \right)^\top f(\mathbf{z}(t_1), t_1, \theta) \quad \triangleright \text{Compute gradient w.r.t. } t_1$$

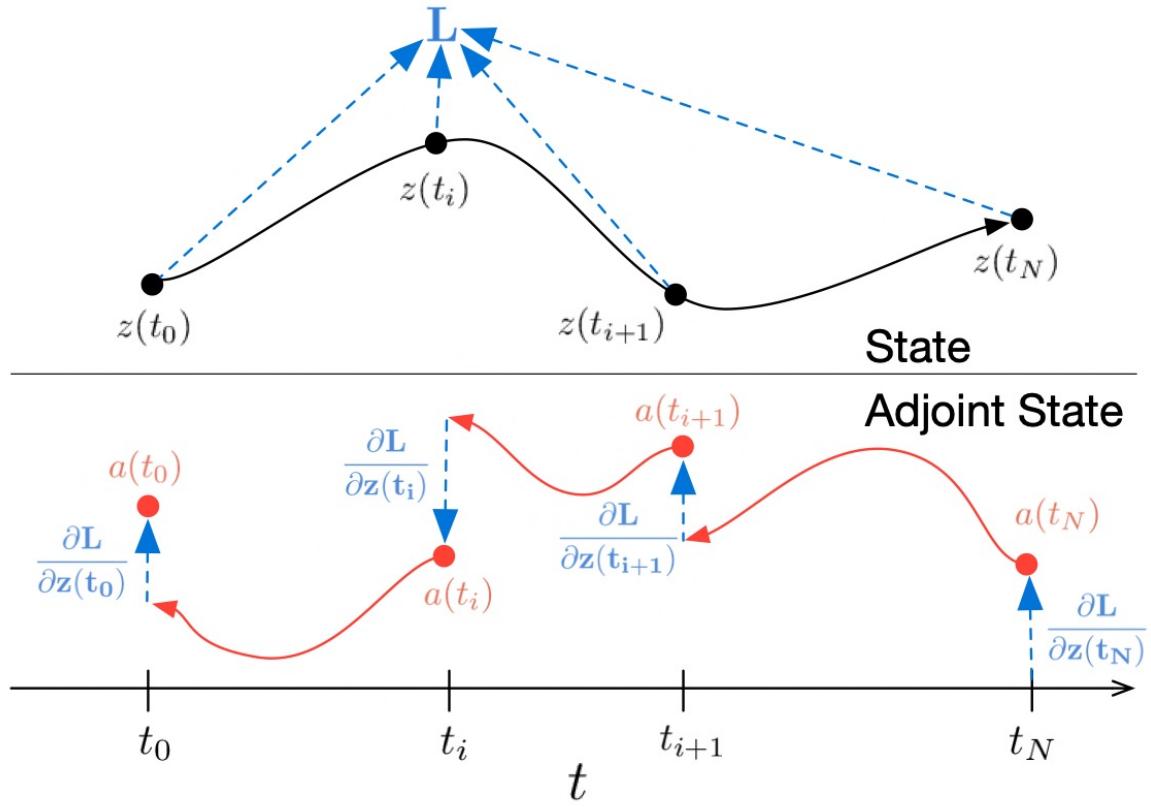
$$s_0 = [\mathbf{z}(t_1), \frac{\partial L}{\partial \mathbf{z}(t_1)}, \mathbf{0}_{|\theta|}, -\frac{\partial L}{\partial t_1}] \quad \triangleright \text{Define initial augmented state}$$

def aug_dynamics($[\mathbf{z}(t), \mathbf{a}(t), \cdot, \cdot], t, \theta$): ▷ Define dynamics on augmented state

return $[f(\mathbf{z}(t), t, \theta), -\mathbf{a}(t)^\top \frac{\partial f}{\partial \mathbf{z}}, -\mathbf{a}(t)^\top \frac{\partial f}{\partial \theta}, -\mathbf{a}(t)^\top \frac{\partial f}{\partial t}]$ ▷ Compute vector-Jacobian products

$[\mathbf{z}(t_0), \frac{\partial L}{\partial \mathbf{z}(t_0)}, \frac{\partial L}{\partial \theta}, \frac{\partial L}{\partial t_0}] = \text{ODESolve}(s_0, \text{aug_dynamics}, t_1, t_0, \theta)$ ▷ Solve reverse-time ODE

return $\frac{\partial L}{\partial \mathbf{z}(t_0)}, \frac{\partial L}{\partial \theta}, \frac{\partial L}{\partial t_0}, \frac{\partial L}{\partial t_1}$ ▷ Return all gradients



Resources

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



<https://github.com/diffeqml/torchdyn>