

## Overview

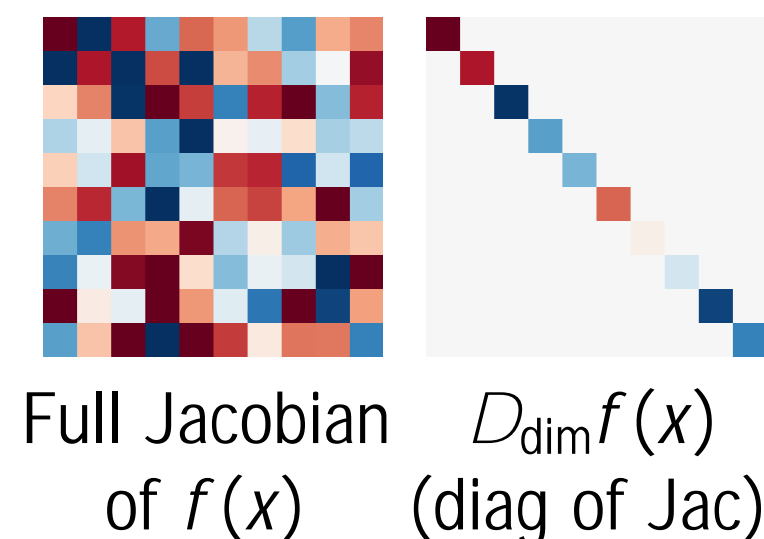
Given a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ , we seek to obtain a vector containing its *dimension-wise*  $k$ -th order derivatives,

$$D_{\dim}^k f(x) := \begin{bmatrix} \frac{\partial^k f_1(x)}{\partial x_1^k} & \dots & \frac{\partial^k f_d(x)}{\partial x_d^k} \end{bmatrix}^T \in \mathbb{R}^d \quad (1)$$

using only  $k$  evaluations of automatic differentiation regardless of the dimension  $d$ .

This has applications in:

- I. Solving differential equations.
- II. Continuous Normalizing Flows.
- III. Learning stochastic differential equations.



## Reverse-mode Automatic Differentiation

Deep learning software rely on automatic differentiation (AD) to compute gradients. More generally, AD computes vector-Jacobian products.

$$v^T \frac{\partial f(x)}{\partial x} = \sum_{i=1}^d v_i \frac{\partial f_i(x)}{\partial x} \quad (2)$$

However, **computing the Jacobian trace (ie.  $\sum D_{\dim}$ ) is as expensive as the full Jacobian.** One evaluation of AD can only compute one row of the Jacobian.

## Forward: Network Structure

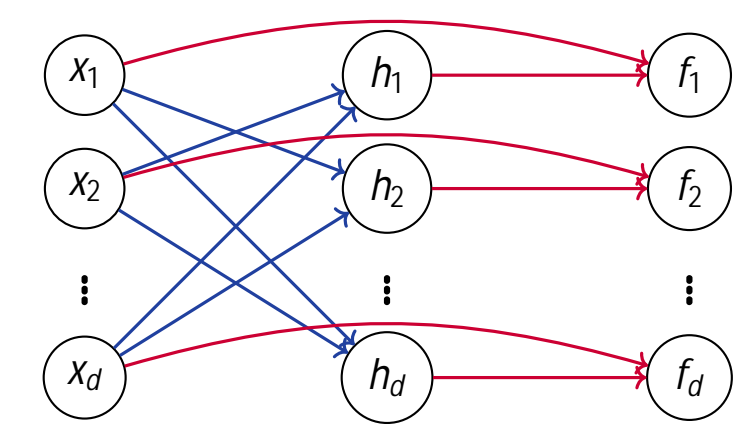
Our approach:

Dimension-wise derivatives can be efficiently computed for *restricted* architectures with simple *modifications to the AD procedure* in the backward pass.

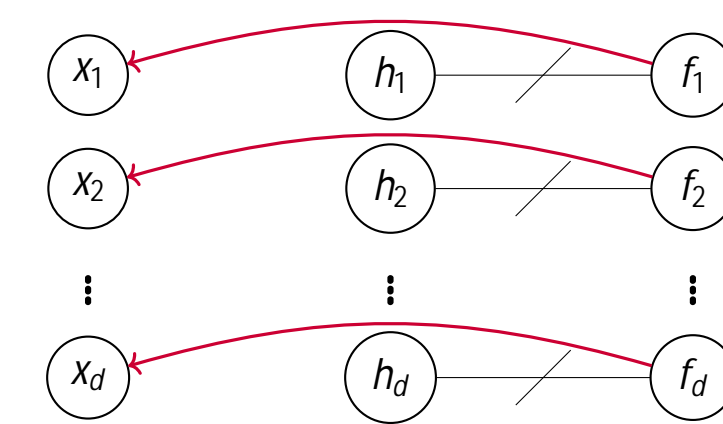
**Conditioner**  $h_i = c_i(x_{-i})$ . The  $i$ -th hidden dimension depends on all inputs except the  $i$ -th input dimension. Can be computed in parallel using masked neural networks (e.g. MADE, PixelCNN).

**Transformer**  $f_i(x) = \sigma_i(x_i; h_i)$ .  $\sigma_i : \mathbb{R}^d \rightarrow \mathbb{R}$  outputs the  $i$ -th dimension given the concatenated vector. Can be computed in parallel if composed of matrix multiplication and element-wise operations.

## Backward: Modified Computation Graph



Forward computation graph



Backward computation graph

Let  $\hat{h}$  be stop\_gradient( $h$ ) and  $\hat{f} = (x; \hat{h})$ , so

$$\frac{\partial \hat{f}_i(x)}{\partial x_j} = \frac{\partial \sigma_i(x_i; \hat{h}_i)}{\partial x_j} = \begin{cases} \frac{\partial f_i(x)}{\partial x_i} & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases} \quad (3)$$

**Dimension-wise derivatives.**

$$\mathbf{1}^T \frac{\partial \hat{f}(x)}{\partial x} = D_{\dim} \hat{f} = \begin{bmatrix} \frac{\partial f_1(x)}{\partial x_1} \\ \dots \\ \frac{\partial f_d(x)}{\partial x_d} \end{bmatrix}^T = D_{\dim} f \quad (4)$$

**Higher orders.**

$$\mathbf{1}^T \frac{\partial D_{\dim}^k \hat{f}(x)}{\partial x} = D_{\dim}^k \hat{f}(x) = D_{\dim}^k f(x) \quad (5)$$

**Backpropagating through dim-wise derivatives.**

$$\frac{\partial D_{\dim}^k \hat{f}}{\partial w} + \frac{\partial D_{\dim}^k \hat{f}}{\partial \hat{h}} \frac{\partial \hat{h}}{\partial w} = \frac{\partial D_{\dim}^k f}{\partial w} \quad (6)$$

## App I: Linear Multistep ODE Solvers

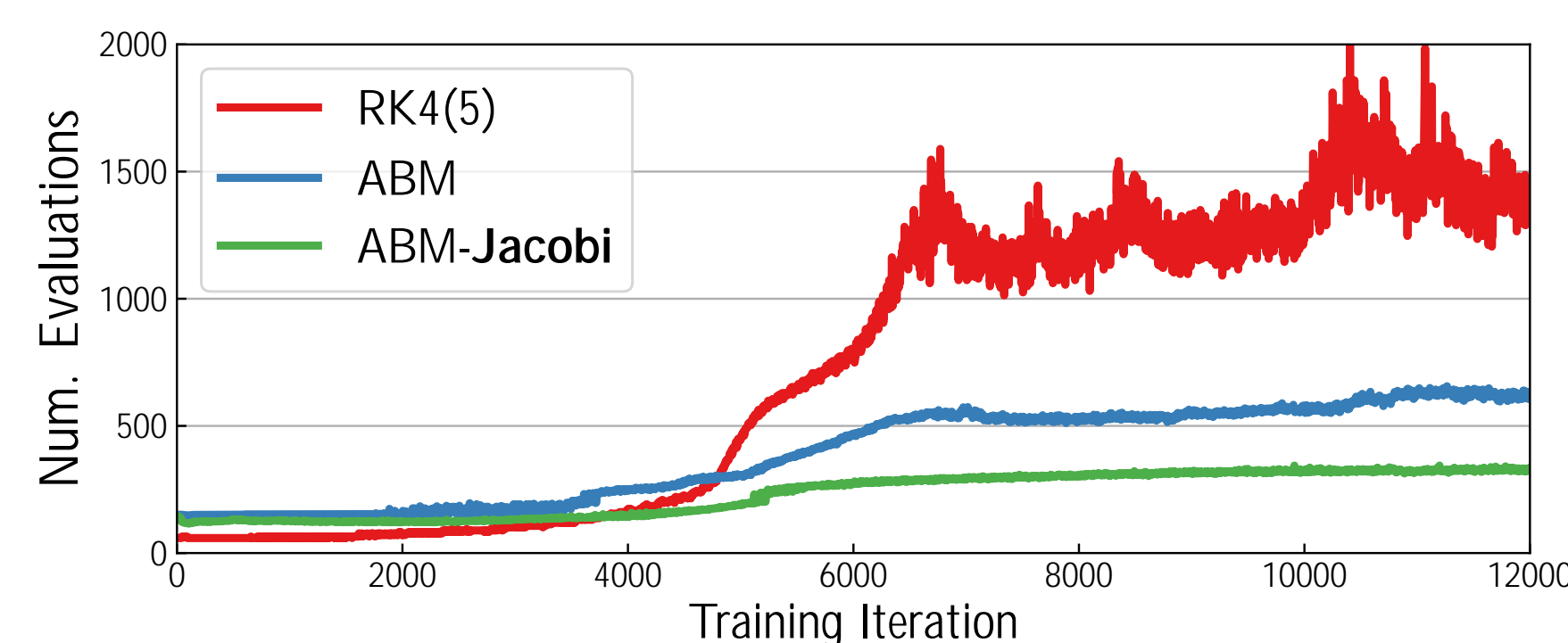
Implicit ODE solvers need to solve an optimization sub-problem in the inner loop.

General idea: replace Newton-Raphson

$$y^{(k+1)} = y^{(k)} - \frac{\partial F(y^{(k)})}{\partial y^{(k)}}^{-1} F(y^{(k)}) \quad (7)$$

with Jacobi-Newton

$$y^{(k+1)} = y^{(k)} - [D_{\dim} F(y)]^{-1} F(y^{(k)}) \quad (8)$$



## App II: Continuous Normalizing Flows

If  $\frac{dx}{dt} = f(t; x)$ , then  $\frac{\partial \log p(x)}{\partial t} = -\text{tr} \frac{\partial f}{\partial x}$ . (See "Neural ODEs")

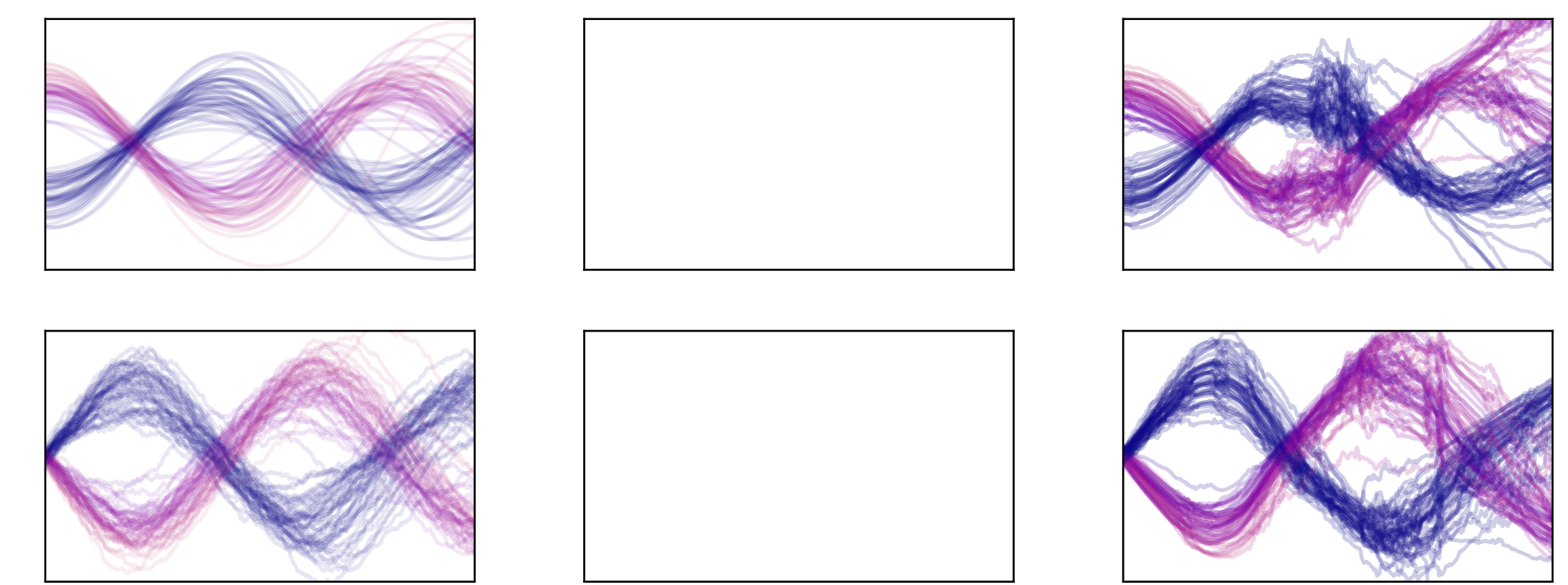
Model	MNIST		Omniglot	
	ELBO	NLL #	ELBO	NLL #
VAE	-86.55	82.14	-104.28	97.25
Planar	-86.06	81.91	-102.65	96.04
IAF	-84.20	80.79	-102.41	96.08
Sylvester	-83.32	<b>80.22</b>	-99.00	<b>93.77</b>
FFJORD	-82.82		-98.33	
Di Op-CNF	<b>-82.37</b>	<b>80.22</b>	<b>-97.42</b>	<b>93.90</b>

Table: Evidence lower bound and negative log-likelihood for static MNIST and Omniglot.

**Exact trace vs stochastic trace.** Converges faster and results in easier to solve dynamical systems.

## App III: Stochastic Differential Eqs

$$dx(t) = f(x(t); t)dt + g(x(t); t)dW \quad (9)$$



Idea: match left- and right-hand-side of the Fokker-Planck equation, which describes the change in density.

$$\frac{\partial p(t; x)}{\partial t} = \sum_{i=1}^d \frac{\partial}{\partial x_i} [f_i(t; x)p(t; x)] + \frac{1}{2} \sum_{i=1}^d \frac{\partial^2}{\partial x_i^2} g_{ii}^2(t; x)p(t; x) \quad (10)$$

## References

- Germain *et al.* "MADE: Masked autoencoder for distribution estimation." (2015)
- Huang *et al.* "Neural Autoregressive Flows." (2018)
- Chen *et al.* "Neural ODEs." (2018)