

# CS330 Automatic differentiation

Matt Johnson, [mattjj@google.com](mailto:mattjj@google.com), Sept 28 2020



# Why study **AD?**

# Why study AD?

1. You use it **every day!**

# Why study **AD**?

1. You use it **every day!**
2. **DL revolution** = data + compute + architectures + **AD**

# Why study **AD**?

1. You use it **every day!**
2. **DL revolution** = data + compute + architectures + **AD**
3. **Research frontiers** use advanced AD
  - iMAML, Neural ODEs, DEQs, OptNet, RevNets, Reformer, new optimizers...

# Why study **AD**?

1. You use it **every day!**
2. **DL revolution** = data + compute + architectures + **AD**
3. **Research frontiers** use advanced AD
  - iMAML, Neural ODEs, DEQs, OptNet, RevNets, Reformer, new optimizers...
4. AD is **interesting in its own right**

# What is autodiff?

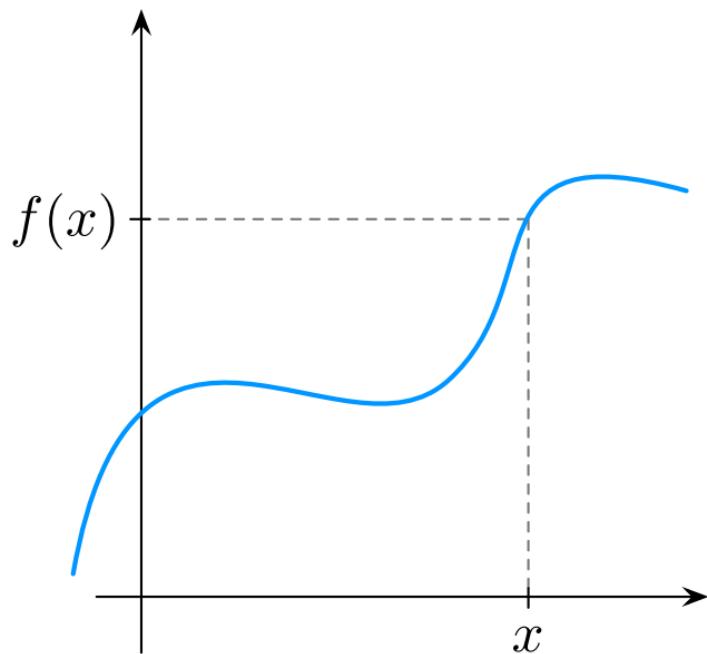
## Demo!

# This lecture

- **Fundamentals**
  - Forward- and reverse-modes
  - Some sense for how they're implemented
- **Advanced techniques**
  - Efficient AD of equation solvers, fixed-points, optimizers
  - Efficient AD of long-running iterative processes
  - Very high-order AD via jet

# Math and notation

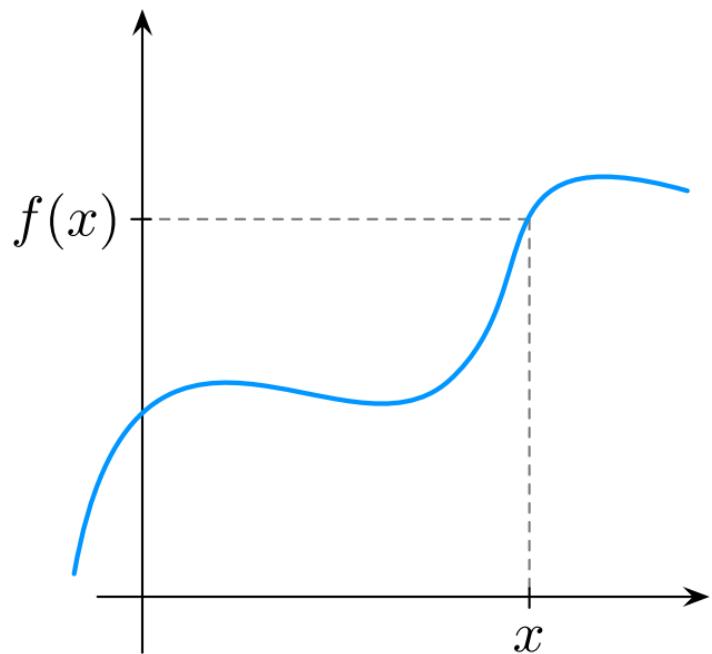
$$\textcolor{blue}{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$$



# Math and notation

$$\textcolor{blue}{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$\partial f : \mathbb{R}^n \rightarrow (\mathbb{R}^n \rightarrow \mathbb{R}^m)$$

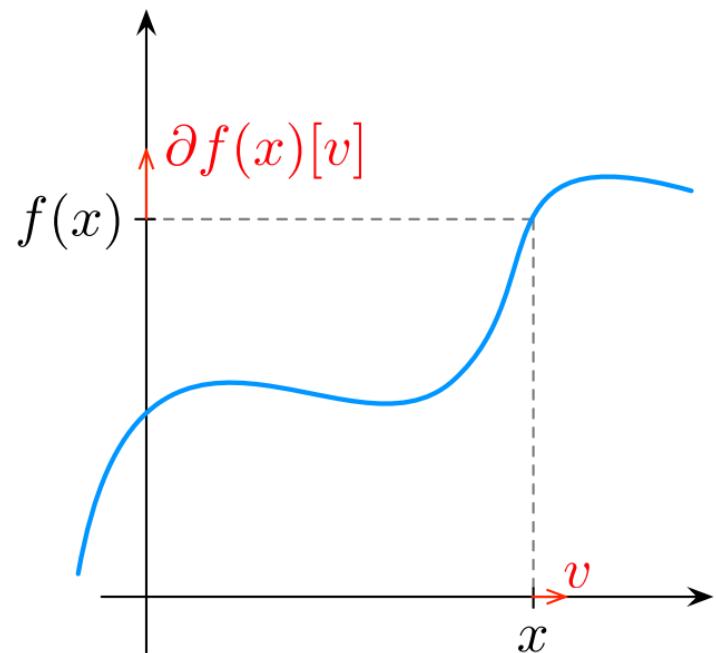


# Math and notation

$$\textcolor{blue}{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$\partial f : \mathbb{R}^n \rightarrow (\mathbb{R}^n \rightarrow \mathbb{R}^m)$$

$$\partial f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$$



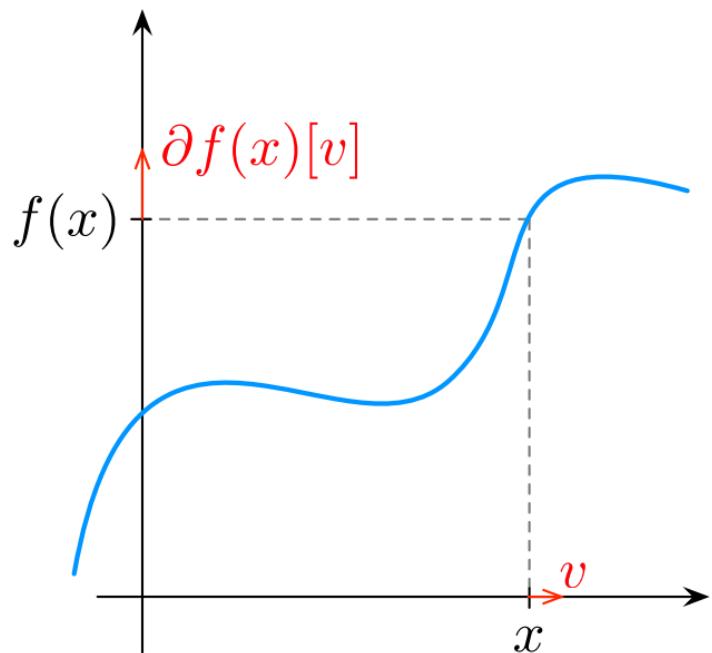
# Math and notation

$$\textcolor{blue}{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$\partial f : \mathbb{R}^n \rightarrow (\mathbb{R}^n \rightarrow \mathbb{R}^m)$$

$$\partial f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$\partial f(x) \in \mathbb{R}^{m \times n}$$

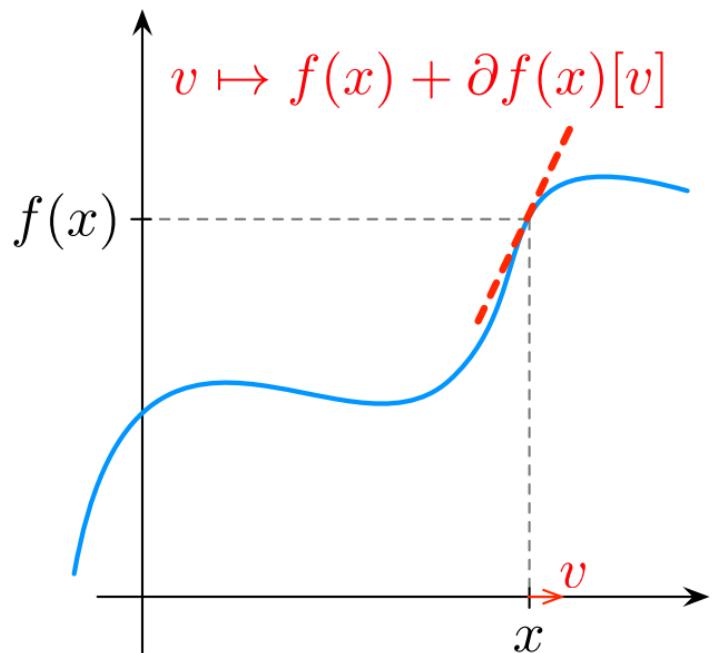


# Math and notation

$$\textcolor{blue}{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$f(x + v) = f(x) + \partial f(x)[v]$$

$$+ \mathcal{O}(\|v\|^2)$$



# Math and notation: comparison

$$\frac{\partial y}{\partial x}$$

# Math and notation: comparison

$$y = f(x)$$

$$\frac{\partial y}{\partial x}$$

# Math and notation: comparison

$$y = f(x)$$

$$\frac{\partial y}{\partial x} = \partial f$$

# Math and notation: comparison

$$y = f(x)$$

$$\frac{\partial y}{\partial x} = \partial f$$

$$\left. \frac{\partial y}{\partial x} \right|_{x=a} = \partial f(a)$$

# Math and notation: gradients

$\nabla f(x)$  is the vector such that

$$\langle \nabla f(x), v \rangle = \partial f(x)[v]$$

# Math and notation: chain rule

$$f = g \circ h \quad f : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

# Math and notation: chain rule

$$f = g \circ h \quad f : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$h : \mathbb{R}^n \rightarrow \mathbb{R}^p \quad g : \mathbb{R}^p \rightarrow \mathbb{R}^m$$

# Math and notation: chain rule

$$f = g \circ h \quad f : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$h : \mathbb{R}^n \rightarrow \mathbb{R}^p \quad g : \mathbb{R}^p \rightarrow \mathbb{R}^m$$

$$\partial f(x) = \partial g(h(x)) \circ \partial h(x)$$

# Math and notation: chain rule

$$f = g \circ h \quad f : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$h : \mathbb{R}^n \rightarrow \mathbb{R}^p \quad g : \mathbb{R}^p \rightarrow \mathbb{R}^m$$

$$\partial f(x) = \partial g(h(x)) \circ \partial h(x)$$

$$\partial f(x)[v] = \partial g(h(x))[\partial h(x)[v]]$$

# Math and notation: chain rule

$$f = g \circ h \quad f : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$h : \mathbb{R}^n \rightarrow \mathbb{R}^p \quad g : \mathbb{R}^p \rightarrow \mathbb{R}^m$$

$$\left. \frac{\partial y}{\partial x} \right|_{x=a} = \left. \frac{\partial y}{\partial z} \right|_{z=h(a)} \left. \frac{\partial z}{\partial x} \right|_{x=a}$$

$z = h(x)$   
 $y = g(z)$

# **Do we need a DAG rule?**

# Do we need a DAG rule?

$$f(x) = (x, x)$$

# Do we need a DAG rule?

$$f(x) = (x, x)$$

$$f(x) = \begin{bmatrix} I \\ I \end{bmatrix} x$$

# The two linear maps of AD

$$1. \ v \mapsto \partial f(x)v \quad \text{or} \quad v \mapsto \partial f(x)[v]$$

JVP / push-forward / forward-mode

build Jacobian one column at a time

# The two linear maps of AD

1.  $v \mapsto \partial f(x)v$  or  $v \mapsto \partial f(x)[v]$

JVP / push-forward / forward-mode

build Jacobian one column at a time

2.  $w^T \mapsto w^T \partial f(x)$  or  $w^T \mapsto \partial f(x)^T[w^T]$

VJP / pull-back / reverse-mode

build Jacobian one row at a time

# The two linear maps of AD

1.  $v \mapsto \partial f(x)v$     or     $v \mapsto \partial f(x)[v]$

# The two linear maps of AD

$$1. \ v \mapsto \partial f(x)v \quad \text{or} \quad v \mapsto \partial f(x)[v]$$

Let's say you have

$$v = \frac{\partial x}{\partial \theta}$$

and you want to compute

$$w = \frac{\partial y}{\partial \theta} \quad \text{where} \quad y = f(x)$$

# The two linear maps of AD

$$1. \ v \mapsto \partial f(x)v \quad \text{or} \quad v \mapsto \partial f(x)[v]$$

Let's say you have

$$v = \frac{\partial x}{\partial \theta}$$


scalar parameter

and you want to compute

$$w = \frac{\partial y}{\partial \theta} \quad \text{where} \quad y = f(x)$$

# The two linear maps of AD

1.  $v \mapsto \partial f(x)v$  or  $v \mapsto \partial f(x)[v]$

Let's say you have

$$v = \frac{\partial x}{\partial \theta} \quad \text{scalar parameter}$$

and you want to compute

$$w = \frac{\partial y}{\partial \theta} \quad \text{where} \quad y = f(x)$$

$$w = \partial f(x)[v] = \frac{\partial y}{\partial x} \frac{\partial x}{\partial \theta}.$$

So  $\partial f(x)$  “pushes forward” perturbation information.

# The two linear maps of AD

2.  $w^T \mapsto w^T \partial f(x)$  or  $w^T \mapsto \partial f(x)^T [w^T]$

# The two linear maps of AD

2.  $w^T \mapsto w^T \partial f(x)$  or  $w^T \mapsto \partial f(x)^T [w^T]$

Let's say you have

$$w^T = \frac{\partial \ell}{\partial y}$$

and you want to compute

$$u^T = \frac{\partial \ell}{\partial x} \quad \text{where} \quad y = f(x)$$

# The two linear maps of AD

2.  $w^T \mapsto w^T \partial f(x)$  or  $w^T \mapsto \partial f(x)^T [w^T]$

Let's say you have

$$w^T = \frac{\partial \ell}{\partial y} \quad \text{scalar loss value}$$

and you want to compute

$$u^T = \frac{\partial \ell}{\partial x} \quad \text{where} \quad y = f(x)$$

# The two linear maps of AD

2.  $w^T \mapsto w^T \partial f(x)$  or  $w^T \mapsto \partial f(x)^T [w^T]$

Let's say you have

$$w^T = \frac{\partial \ell}{\partial y} \quad \text{scalar loss value}$$

and you want to compute

$$u^T = \frac{\partial \ell}{\partial x} \quad \text{where} \quad y = f(x)$$

$$u^T = \partial f(x)^T [w^T] = \frac{\partial \ell}{\partial y} \frac{\partial y}{\partial x}.$$

So  $\partial f(x)^T$  “pulls back” sensitivity information.

# From math to code

(math)

$$f(x) = \sin^2(x)$$

impl

(code)

```
def f(x):  
    return sin(x) ** 2
```

↓ differentiate

$$\nabla f(x)$$

impl

↓ AD

grad(f)

# Two kinds of function: primitive and composite

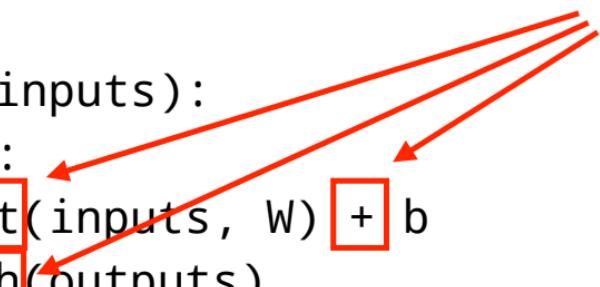
```
import jax.numpy as jnp

def predict(params, inputs):
    for W, b in params:
        outputs = jnp.dot(inputs, W) + b
        inputs = jnp.tanh(outputs)
    return outputs
```

# Two kinds of function: primitive and composite

```
import jax.numpy as jnp  
  
def predict(params, inputs):  
    for W, b in params:  
        outputs = jnp.dot(inputs, W) + b  
        inputs = jnp.tanh(outputs)  
    return outputs
```

primitives

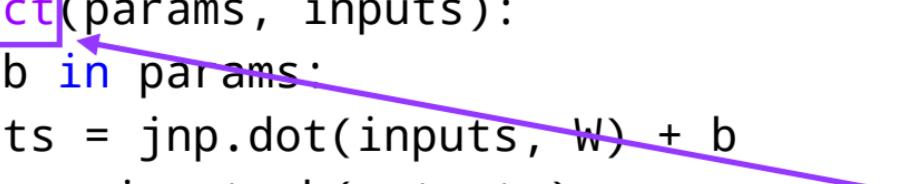


The diagram consists of three red arrows originating from the word "primitives" in red text at the top right. Each arrow points to one of the three highlighted code snippets in the Python code below. The first arrow points to the expression `jnp.dot(inputs, W)`. The second arrow points to the addition operator `+`. The third arrow points to the expression `jnp.tanh(outputs)`.

# Two kinds of function: primitive and composite

```
import jax.numpy as jnp

def predict(params, inputs):
    for W, b in params:
        outputs = jnp.dot(inputs, W) + b
        inputs = jnp.tanh(outputs)
    return outputs
```



composite

# Two kinds of function: primitive and composite

```
import jax.numpy as jnp

def predict(params, inputs):
    for W, b in params:
        outputs = jnp.dot(inputs, W) + b
        inputs = jnp.tanh(outputs)
    return outputs
```

For **primitive** functions we have a lookup table

For **composite** functions we use the chain rule

# JVPs in code: types

$$(x, v) \mapsto (f(x), \partial f(x)[v])$$

# JVPs in code: types

$$(x, v) \mapsto (f(x), \partial f(x)[v])$$

```
jvp :: (a -> b) -> (a, T a) -> (b, T b)
```

# JVPs in code: types

$$(x, v) \mapsto (f(x), \partial f(x)[v])$$

jvp :: (a -> b) -> (a, T a) -> (b, T b)

given a function

# JVPs in code: types

$$(x, v) \mapsto (f(x), \partial f(x)[v])$$

jvp :: ( $a \rightarrow b$ )  $\rightarrow$  ( $a, T a$ )  $\rightarrow$  ( $b, T b$ )

given a function

a “primal” input, and  
“tangent” input perturbation

# JVPs in code: types

$$(x, v) \mapsto (f(x), \partial f(x)[v])$$

jvp :: ( $a \rightarrow b$ )  $\rightarrow (a, T a)$   $\rightarrow (b, T b)$

given a function

a “primal” input, and  
“tangent” input perturbation

return a primal output, and  
tangent output perturbation

# JVPs in code: types

$$(x, v) \mapsto (f(x), \partial f(x)[v])$$

jvp :: ( $a \rightarrow b$ )  $\rightarrow (a, T a)$   $\rightarrow (b, T b)$

given a function

a “primal” input, and  
“tangent” input perturbation

return a primal output, and  
tangent output perturbation

$$T a \equiv a$$

# JVPs in code: abstract version

```
jvp :: (a -> b) -> (a, T a) -> (b, T b)
```

# JVPs in code: abstract version

```
jvp :: (a -> b) -> (a, T a) -> (b, T b)
```

```
jvp sin (x, x_dot) =  
  let y = sin x  
      y_dot = (cos x) * x_dot  
  in (y, y_dot)
```

# JVPs in code: abstract version

```
jvp :: (a -> b) -> (a, T a) -> (b, T b)
```

```
jvp (f . g) (x, x_dot) =  
  let (y, y_dot) = jvp g (x, x_dot)  
      (z, z_dot) = jvp f (y, y_dot)  
  in (z, z_dot)
```

# JVPs in code: Python version (sketch)

# JVPs in code: Python version (sketch)

```
DualNumber = namedtuple('DualNumber',
                       ['primal', 'tangent'])
```

```
def add(x, y):
    t = (type(x), type(y))
    if t == (DualNumber, DualNumber):
        return DualNumber(add(x.primal, y.primal),
                           add(x.tangent, y.tangent))
    elif ...:
```

# JVPs in code: Python version (sketch)

```
DualNumber = namedtuple('DualNumber',  
                       ['primal', 'tangent'])
```

```
def add(x, y):  
    t = (type(x), type(y))  
    if t == (DualNumber, DualNumber):  
        return DualNumber(add(x.primal, y.primal),  
                           add(x.tangent, y.tangent))  
    elif ...
```

what if these represent  
different perturbations?

# Demo!

# Linearization in code: types

$$x \mapsto (f(x), \partial f(x))$$

# Linearization in code: types

$$x \mapsto (f(x), \partial f(x))$$

```
lin :: (a -> b) -> a -> (b, T a --o T b)
```

# Linearization in code: types

$$x \mapsto (f(x), \partial f(x))$$

lin :: ( $a \rightarrow b$ )  $\rightarrow a \rightarrow (b, T a \dashv\vdash T b)$

given a function

# Linearization in code: types

$$x \mapsto (f(x), \partial f(x))$$

lin :: ( $a \rightarrow b$ )  $\rightarrow$   $a$   $\rightarrow$  ( $b$ , T  $a$  --o T  $b$ )

given a function

and a primal input

# Linearization in code: types

$$x \mapsto (f(x), \partial f(x))$$

lin :: ( $a \rightarrow b$ )  $\rightarrow$   $a$   $\rightarrow$   $(b, T a \dashv\vdash T b)$

given a function

and a primal input

return a primal output, and  
a linearized function

# Linearization in code: types

$$x \mapsto (f(x), \partial f(x))$$

lin :: ( $a \rightarrow b$ )  $\rightarrow$   $a$   $\rightarrow$   $(b, T a \dashv\! \dashv T b)$

given a function

and a primal input

return a primal output, and  
a linearized function

$\dashv\! \dashv$  means “*linear* function”

# Linearization in code: abstract version

```
lin ::  (a -> b) -> a -> (b, T a --o T b)
```

# Linearization in code: abstract version

```
lin :: (a -> b) -> a -> (b, T a --o T b)
```

```
lin sin x =
  (sin x, lambda x_dot: (cos x) * x_dot)
```

# Linearization in code: abstract version

```
lin :: (a -> b) -> a -> (b, T a --o T b)
```

```
lin (f . g) x =
  let y, g_lin = lin g x
      z, f_lin = lin f y
  in (z, f_lin . g_lin)
```

# VJPs in code: types

$$x \mapsto (f(x), \partial f(x)^\top)$$

# VJPs in code: types

$$x \mapsto (f(x), \partial f(x)^\top)$$

vjp :: ( $a \rightarrow b$ )  $\rightarrow a \rightarrow (b, \text{CT } b \dashv\circ \text{CT } a)$

# VJPs in code: types

$$x \mapsto (f(x), \partial f(x)^\top)$$

vjp :: ( $a \rightarrow b$ )  $\rightarrow a \rightarrow (b, \text{CT } b \dashv\! \dashv \text{CT } a)$

given a function

# VJPs in code: types

$$x \mapsto (f(x), \partial f(x)^\top)$$

vjp :: ( $a \rightarrow b$ )  $\rightarrow a$   $\rightarrow (b, \text{CT } b \text{ --o } \text{CT } a)$

given a function

and a primal input

# VJPs in code: types

$$x \mapsto (f(x), \partial f(x)^\top)$$

vjp :: ( $a \rightarrow b$ )  $\rightarrow a$   $\rightarrow (b, \text{CT } b \text{ --o } \text{CT } a)$

given a function

and a primal input

return a primal output, and  
the transposed linearized function

# VJPs in code: types

$$x \mapsto (f(x), \partial f(x)^\top)$$

vjp :: ( $a \rightarrow b$ )  $\rightarrow a$   $\rightarrow (b, CT\ b \dashv\circ CT\ a)$

given a function

and a primal input

return a primal output, and  
the transposed linearized function

$$CT\ a \equiv a$$

# VJPs in code: abstract version

```
vjp :: (a -> b) -> a -> (b, CT b --o CT a)
```

# VJPs in code: abstract version

```
vjp :: (a -> b) -> a -> (b, CT b --o CT a)
```

```
vjp sin x =  
  let y = sin x  
    sin_vjp = lambda y_bar: (cos x) * y_bar  
  in (y, sin_vjp)
```

# VJPs in code: abstract version

```
vjp :: (a -> b) -> a -> (b, CT b --o CT a)
```

```
vjp (f . g) x =
  let (y, g_vjp) = vjp g x
      (z, f_vjp) = vjp f y
  in (z, g_vjp . f_vjp)
```

# VJPs in code: Python version

<https://github.com/mattjj/autodidact>

[http://videolectures.net/deeplearning2017\\_johnson\\_automatic\\_differentiation/](http://videolectures.net/deeplearning2017_johnson_automatic_differentiation/)

# The two linear maps of AD

`jvp :: (a -> b) -> (a, T a) -> (b, T b)`

build Jacobian one **column** at a time

costs  $\mathcal{O}(1)$  times the FLOPs, and  $\sim 2x$  memory

# The two linear maps of AD

`jvp :: (a -> b) -> (a, T a) -> (b, T b)`

build Jacobian one **column** at a time

costs  $\mathcal{O}(1)$  times the FLOPs, and  $\sim 2x$  memory

`vjp :: (a -> b) -> a -> (b, CT b --o CT a)`

build Jacobian one **row** at a time

costs  $\mathcal{O}(1)$  times the FLOPs, and  $\mathcal{O}(\text{depth})$  memory

**What about **grad**?**

# What about **grad**?

$$\ell : \mathbb{R}^{10^9} \rightarrow \mathbb{R}$$

# What about **grad**?

$$\ell : \mathbb{R}^{10^9} \rightarrow \mathbb{R} \quad \partial \ell(\theta) \in \mathbb{R}^{1 \times 10^9}$$

# What about **grad**?

$$\ell : \mathbb{R}^{10^9} \rightarrow \mathbb{R} \quad \partial \ell(\theta) \in \mathbb{R}^{1 \times 10^9}$$

```
def grad(x):
    def gradfun(*args):
        _, f_vjp = vjp(f, *args)
        return f_vjp(1.)
    return gradfun
```

# Demo!

# Differentiating solvers and optimizers

$$x^*(\textcolor{red}{a}) \triangleq \arg \min_x f(x, \textcolor{red}{a})$$

$$x^*(\textcolor{red}{a}) \triangleq \text{solve } g(x, \textcolor{red}{a}) = 0$$

# Differentiating solvers and optimizers

$$x^*(\textcolor{red}{a}) \triangleq \arg \min_x f(x, \textcolor{red}{a})$$



optimality conditions

$$x^*(\textcolor{red}{a}) \triangleq \text{solve } g(x, \textcolor{red}{a}) = 0$$

# Differentiating solvers and optimizers

$$x^*(\textcolor{red}{a}) \triangleq \text{solve } g(x, \textcolor{red}{a}) = 0$$

$$g(x^*, a_0) = 0$$

# Differentiating solvers and optimizers

$$x^*(\textcolor{red}{a}) \triangleq \text{solve } g(x, \textcolor{red}{a}) = 0$$

$$g(x^*(\textcolor{red}{a}), \textcolor{red}{a}) = 0 \quad \forall \textcolor{red}{a} \text{ near } a_0$$

# Differentiating solvers and optimizers

$$x^*(\textcolor{red}{a}) \triangleq \text{solve } g(x, \textcolor{red}{a}) = 0$$

$$g(x^*(\textcolor{red}{a}), \textcolor{red}{a}) = 0 \quad \forall \textcolor{red}{a} \text{ near } a_0$$

$$\partial_0 g(x^*(\textcolor{red}{a}), \textcolor{red}{a}) \partial x^*(\textcolor{red}{a}) + \partial_1 g(x^*(\textcolor{red}{a}), \textcolor{red}{a}) = 0$$

# Differentiating solvers and optimizers

$$x^*(\textcolor{red}{a}) \triangleq \text{solve } g(x, \textcolor{red}{a}) = 0$$

$$g(x^*(\textcolor{red}{a}), \textcolor{red}{a}) = 0 \quad \forall \textcolor{red}{a} \text{ near } a_0$$

$$\partial_0 g(x^*(\textcolor{red}{a}), \textcolor{red}{a}) \partial x^*(\textcolor{red}{a}) + \partial_1 g(x^*(\textcolor{red}{a}), \textcolor{red}{a}) = 0$$

$$\partial x^*(\textcolor{red}{a}) = -\partial_0 g(x^*, \textcolor{red}{a})^{-1} \partial_1 g(x^*, \textcolor{red}{a})$$

# Differentiating solvers and optimizers

$$x^*(\textcolor{red}{a}) \triangleq \text{solve } g(x, \textcolor{red}{a}) = 0$$

$$g(x^*(\textcolor{red}{a}), \textcolor{red}{a}) = 0 \quad \forall \textcolor{red}{a} \text{ near } a_0$$

$$\partial_0 g(x^*(\textcolor{red}{a}), \textcolor{red}{a}) \partial x^*(\textcolor{red}{a}) + \partial_1 g(x^*(\textcolor{red}{a}), \textcolor{red}{a}) = 0$$

$$\partial x^*(\textcolor{red}{a}) = -\partial_0 g(x^*, \textcolor{red}{a})^{-1} \partial_1 g(x^*, \textcolor{red}{a})$$

$$\textcolor{violet}{w}^\top \partial x^*(\textcolor{red}{a}) = - (\textcolor{violet}{w}^\top \partial_0 g(x^*, \textcolor{red}{a})^{-1}) \partial_1 g(x^*, \textcolor{red}{a})$$

# Differentiating iterative fixed points

$$x^*(\textcolor{red}{a}) \triangleq \text{solve } x = f(x, \textcolor{red}{a})$$

# Differentiating iterative fixed points

$$x^*(\textcolor{red}{a}) \triangleq \text{solve } x = f(x, \textcolor{red}{a})$$

$$x^{k+1} = f(x^k, \textcolor{red}{a})$$

# Differentiating iterative fixed points

$$x^*(\textcolor{red}{a}) \triangleq \text{solve } x = f(x, \textcolor{red}{a})$$

$$x^{k+1} = f(x^k, \textcolor{red}{a})$$

$$\partial x^*(\textcolor{red}{a}) = \partial_0 f(x^*(\textcolor{red}{a}), \textcolor{red}{a}) \partial x^*(\textcolor{red}{a}) + \partial_1 f(x^*(\textcolor{red}{a}), \textcolor{red}{a})$$

# Differentiating iterative fixed points

$$x^*(\textcolor{red}{a}) \triangleq \text{solve } x = f(x, \textcolor{red}{a})$$

$$x^{k+1} = f(x^k, \textcolor{red}{a})$$

$$\partial x^*(\textcolor{red}{a}) = \partial_0 f(x^*(\textcolor{red}{a}), \textcolor{red}{a}) \partial x^*(\textcolor{red}{a}) + \partial_1 f(x^*(\textcolor{red}{a}), \textcolor{red}{a})$$

$$A \triangleq \partial_0 f(x^*(\textcolor{red}{a}), \textcolor{red}{a})$$

$$B \triangleq \partial_1 f(x^*(\textcolor{red}{a}), \textcolor{red}{a})$$

# Differentiating iterative fixed points

$$x^*(\textcolor{red}{a}) \triangleq \text{solve } x = f(x, \textcolor{red}{a})$$

$$\partial x^*(\textcolor{red}{a}) = A \partial x^*(\textcolor{red}{a}) + B$$

# Differentiating iterative fixed points

$$x^*(\textcolor{red}{a}) \triangleq \text{solve } x = f(x, \textcolor{red}{a})$$

$$\partial x^*(\textcolor{red}{a}) = A \partial x^*(\textcolor{red}{a}) + B$$

$$\partial x^*(\textcolor{red}{a}) = (I - A)^{-1} B$$

# Differentiating iterative fixed points

$$x^*(\textcolor{red}{a}) \triangleq \text{solve } x = f(x, \textcolor{red}{a})$$

$$\partial x^*(\textcolor{red}{a}) = A \partial x^*(\textcolor{red}{a}) + B$$

$$\textcolor{violet}{w}^\top \partial x^*(\textcolor{red}{a}) = \textcolor{violet}{w}^\top (I - A)^{-1} B$$

# Differentiating iterative fixed points

$$x^*(\textcolor{red}{a}) \triangleq \text{solve } x = f(x, \textcolor{red}{a})$$

$$\partial x^*(\textcolor{red}{a}) = A \partial x^*(\textcolor{red}{a}) + B$$

$$\textcolor{violet}{w}^\top \partial x^*(\textcolor{red}{a}) = \textcolor{violet}{w}^\top (I - A)^{-1} B$$

$$= \textcolor{violet}{u}^\top B$$

$$\text{where } \textcolor{violet}{u}^\top = \textcolor{violet}{w}^\top (I - A)^{-1}$$

# Differentiating iterative fixed points

$$x^*(\textcolor{red}{a}) \triangleq \text{solve } x = f(x, \textcolor{red}{a})$$

$$\partial x^*(\textcolor{red}{a}) = A \partial x^*(\textcolor{red}{a}) + B$$

$$\textcolor{violet}{w}^\top \partial x^*(\textcolor{red}{a}) = \textcolor{violet}{w}^\top (I - A)^{-1} B$$

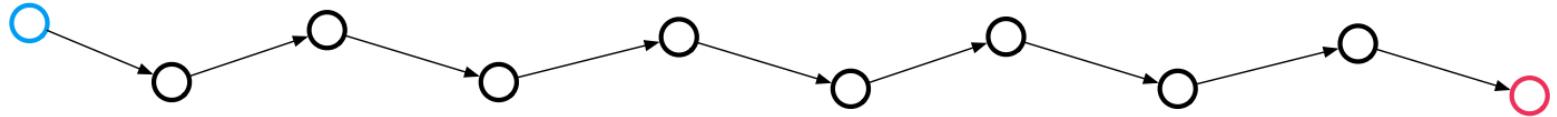
$$= \textcolor{violet}{u}^\top B$$

where  $\textcolor{violet}{u}^\top = \textcolor{violet}{w}^\top + \textcolor{violet}{u}^\top A$     backward pass is  
a (linear) fixed point!

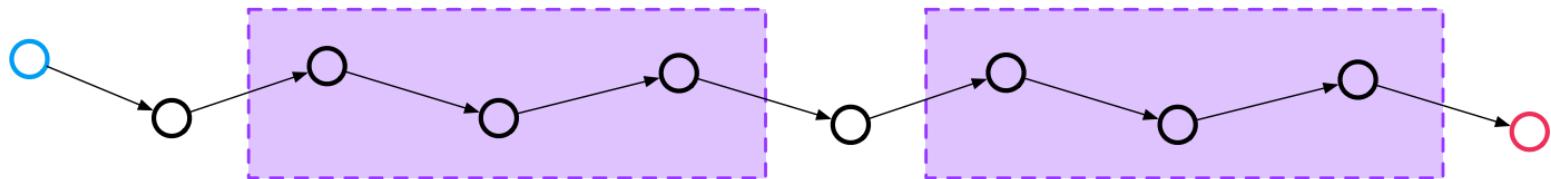
# Demo!

# **Differentiating long iterations: checkpoints**

# Differentiating long iterations: checkpoints



# Differentiating long iterations: checkpoints



# Differentiating long iterations: checkpoints



# Differentiating long iterations: checkpoints



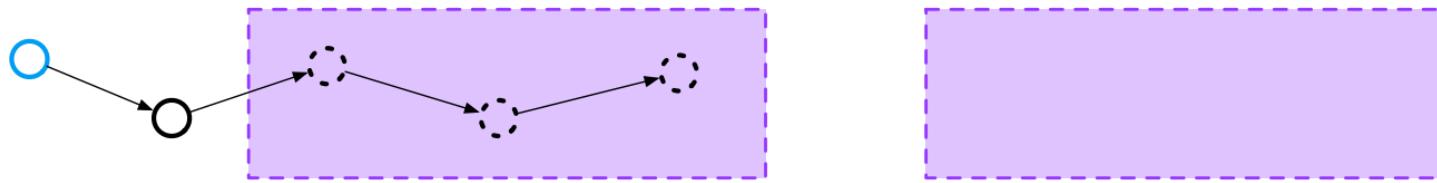
# Differentiating long iterations: checkpoints



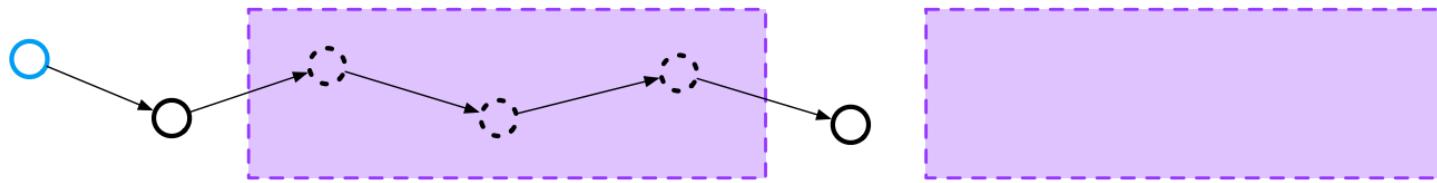
# Differentiating long iterations: checkpoints



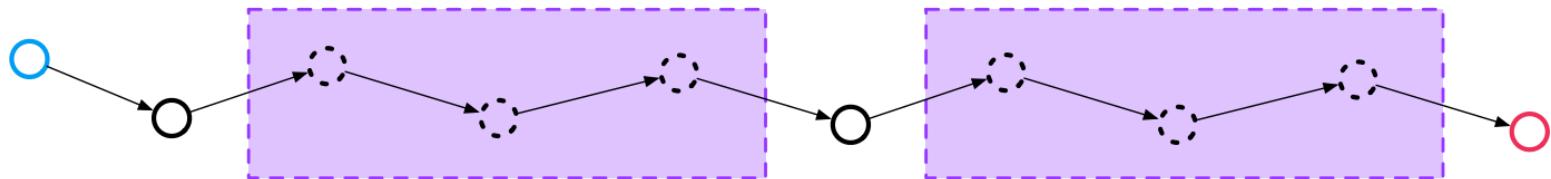
# Differentiating long iterations: checkpoints



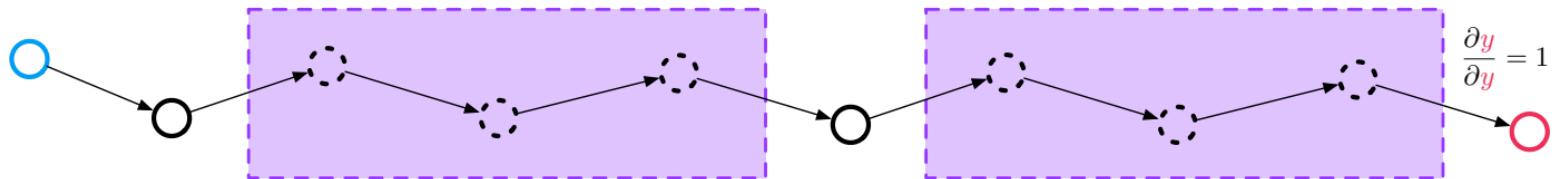
# Differentiating long iterations: checkpoints



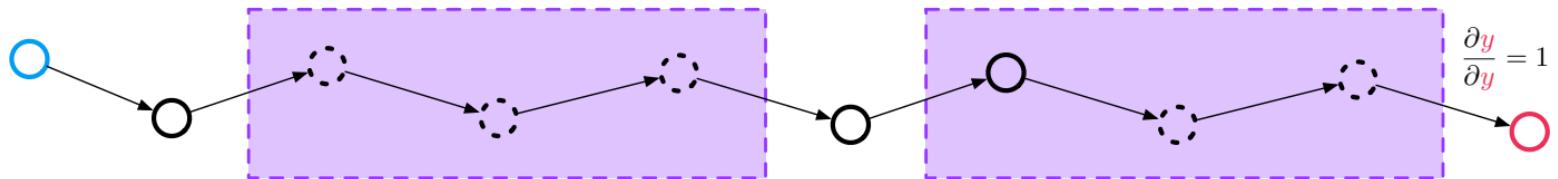
# Differentiating long iterations: checkpoints



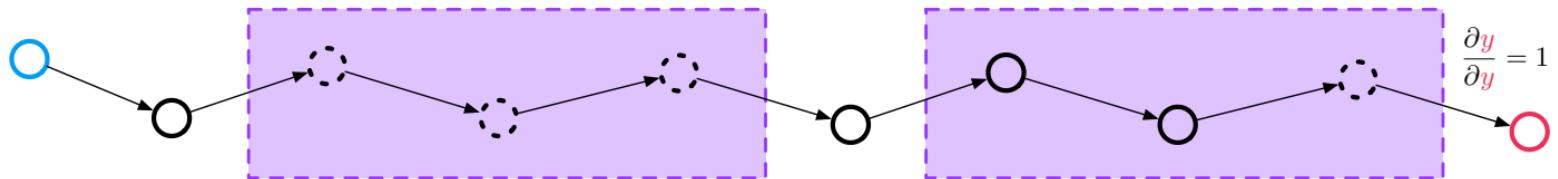
# Differentiating long iterations: checkpoints



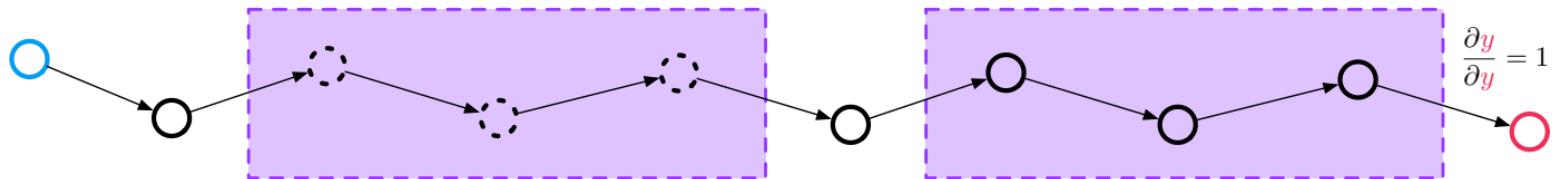
# Differentiating long iterations: checkpoints



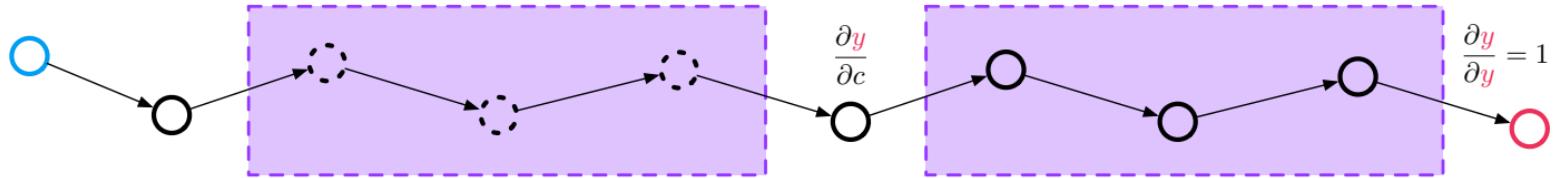
# Differentiating long iterations: checkpoints



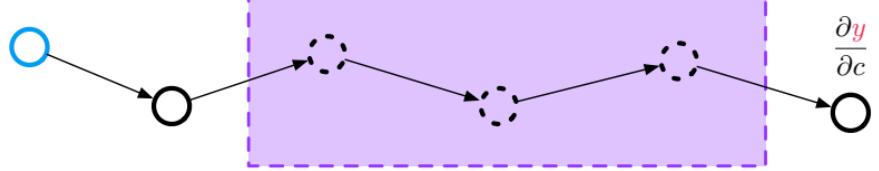
# Differentiating long iterations: checkpoints



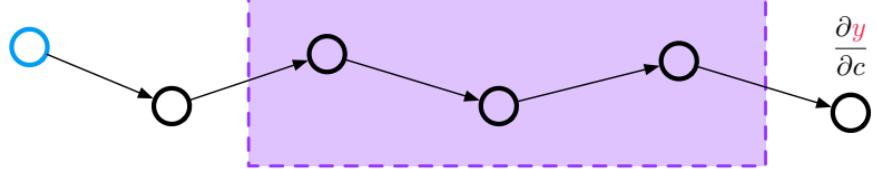
# Differentiating long iterations: checkpoints



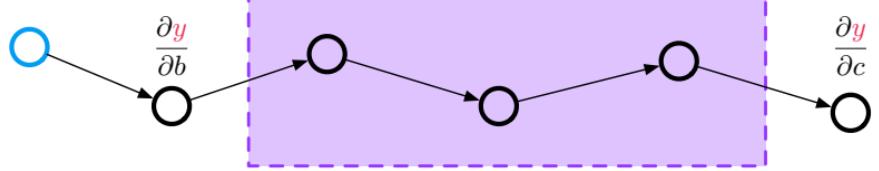
# Differentiating long iterations: checkpoints



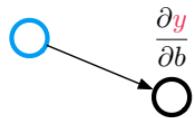
# Differentiating long iterations: checkpoints



# Differentiating long iterations: checkpoints



# Differentiating long iterations: checkpoints



# Demo!

# Differentiating long iterations: checkpoints

Computing VJPs uses  $\mathcal{O}(D)$  memory for depth  $D$ .

# Differentiating long iterations: checkpoints

Computing VJPs uses  $\mathcal{O}(D)$  memory for depth  $D$ .

With recursive checkpointing, only need  $\mathcal{O}(\log D)$  memory,  
at overhead of  $\mathcal{O}(\log D)$  times as much computation.

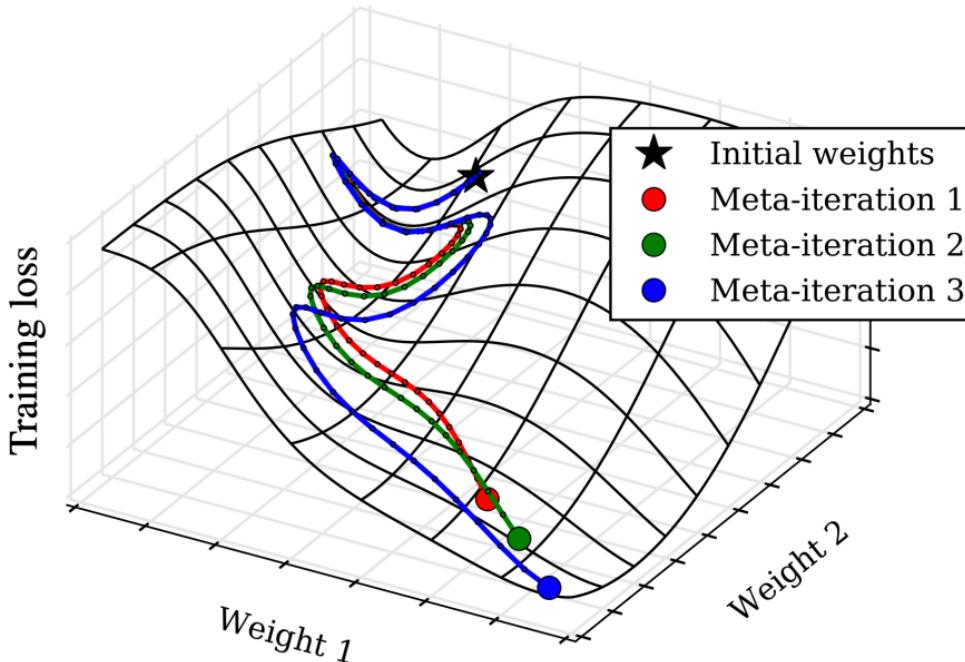
# Differentiating long iterations: checkpoints

Computing VJPs uses  $\mathcal{O}(D)$  memory for depth  $D$ .

With recursive checkpointing, only need  $\mathcal{O}(\log D)$  memory,  
at overhead of  $\mathcal{O}(\log D)$  times as much computation.

... can we ever get  $\mathcal{O}(1)$  memory and  $\mathcal{O}(1)$  overhead?

# Differentiating long iterations: reversibility



MacLaurin\*, Duvenaud\*, and Adams. *Gradient-based hyperparameter optimization through reversible learning*. <https://arxiv.org/pdf/1502.03492.pdf> ICML 2015.

# Differentiating long iterations: reversibility

---

**Algorithm 1** Stochastic gradient descent with momentum

---

```
1: input: initial  $\mathbf{w}_1$ , decays  $\gamma$ , learning rates  $\alpha$ , loss function  $L(\mathbf{w}, \theta, t)$ 
2: initialize  $\mathbf{v}_1 = \mathbf{0}$ 
3: for  $t = 1$  to  $T$  do
4:    $\mathbf{g}_t = \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$             $\triangleright$  evaluate gradient
5:    $\mathbf{v}_{t+1} = \gamma_t \mathbf{v}_t - (1 - \gamma_t) \mathbf{g}_t$      $\triangleright$  update velocity
6:    $\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t \mathbf{v}_t$                   $\triangleright$  update position
7: end for
8: output trained parameters  $\mathbf{w}_T$ 
```

---

# Differentiating long iterations: reversibility

---

## Algorithm 1 Stochastic gradient descent with momentum

---

- 1: **input:** initial  $\mathbf{w}_1$ , decays  $\gamma$ , learning rates  $\alpha$ , loss function  $L(\mathbf{w}, \theta, t)$
- 2: initialize  $\mathbf{v}_1 = \mathbf{0}$
- 3: **for**  $t = 1$  to  $T$  **do**
- 4:    $\mathbf{g}_t = \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$  ▷ evaluate gradient
- 5:    $\mathbf{v}_{t+1} = \gamma_t \mathbf{v}_t - (1 - \gamma_t) \mathbf{g}_t$  ▷ update velocity
- 6:    $\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t \mathbf{v}_t$  ▷ update position
- 7: **end for**
- 8: **output** trained parameters  $\mathbf{w}_T$

---

---

## Algorithm 2 Reverse-mode differentiation of SGD

---

- 1: **input:**  $\mathbf{w}_T, \mathbf{v}_T, \gamma, \alpha$ , train loss  $L(\mathbf{w}, \theta, t)$ , loss  $f(\mathbf{w})$
- 2: initialize  $d\mathbf{v} = \mathbf{0}, d\theta = \mathbf{0}, d\alpha_t = \mathbf{0}, d\gamma = \mathbf{0}$
- 3: initialize  $d\mathbf{w} = \nabla_{\mathbf{w}} f(\mathbf{w}_T)$
- 4: **for**  $t = T$  counting down to 1 **do**
- 5:    $d\alpha_t = d\mathbf{w}^T \mathbf{v}_t$
- 6:    $\mathbf{w}_{t-1} = \mathbf{w}_t - \alpha_t \mathbf{v}_t$
- 7:    $\mathbf{g}_t = \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$
- 8:    $\mathbf{v}_{t-1} = [\mathbf{v}_t + (1 - \gamma_t) \mathbf{g}_t] / \gamma_t$
- 9:    $d\mathbf{v} = d\mathbf{v} + \alpha_t d\mathbf{w}$
- 10:    $d\gamma_t = d\mathbf{v}^T (\mathbf{v}_t + \mathbf{g}_t)$
- 11:    $d\mathbf{w} = d\mathbf{w} - (1 - \gamma_t) d\mathbf{v} \nabla_{\mathbf{w}} \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$
- 12:    $d\theta = d\theta - (1 - \gamma_t) d\mathbf{v} \nabla_{\theta} \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$
- 13:    $d\mathbf{v} = \gamma_t d\mathbf{v}$
- 14: **end for**
- 15: **output** gradient of  $f(\mathbf{w}_T)$  w.r.t  $\mathbf{w}_1, \mathbf{v}_1, \gamma, \alpha$  and  $\theta$

---

# Differentiating long iterations: reversibility

---

## Algorithm 1 Stochastic gradient descent with momentum

---

- 1: **input:** initial  $\mathbf{w}_1$ , decays  $\gamma$ , learning rates  $\alpha$ , loss function  $L(\mathbf{w}, \theta, t)$
- 2: initialize  $\mathbf{v}_1 = \mathbf{0}$
- 3: **for**  $t = 1$  to  $T$  **do**
- 4:    $\mathbf{g}_t = \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$  ▷ evaluate gradient
- 5:    $\mathbf{v}_{t+1} = \gamma_t \mathbf{v}_t - (1 - \gamma_t) \mathbf{g}_t$  ▷ update velocity
- 6:    $\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t \mathbf{v}_t$  ▷ update position
- 7: **end for**
- 8: **output** trained parameters  $\mathbf{w}_T$

---

---

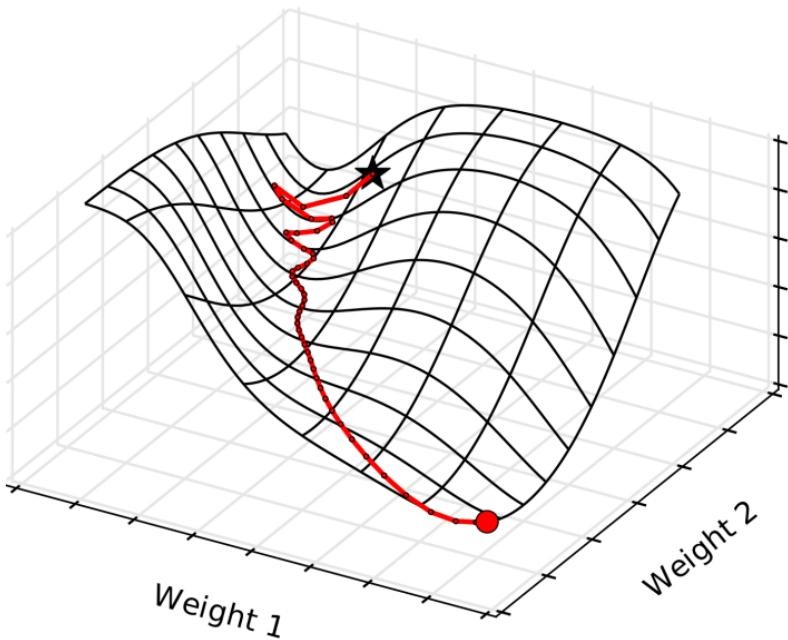
## Algorithm 2 Reverse-mode differentiation of SGD

---

- 1: **input:**  $\mathbf{w}_T, \mathbf{v}_T, \gamma, \alpha$ , train loss  $L(\mathbf{w}, \theta, t)$ , loss  $f(\mathbf{w})$
- 2: initialize  $d\mathbf{v} = \mathbf{0}, d\theta = \mathbf{0}, d\alpha_t = \mathbf{0}, d\gamma = \mathbf{0}$
- 3: initialize  $d\mathbf{w} = \nabla_{\mathbf{w}} f(\mathbf{w}_T)$
- 4: **for**  $t = T$  counting down to 1 **do**
- 5:    $d\alpha_t = d\mathbf{w}^T \mathbf{v}_t$
- 6:    $\mathbf{w}_{t-1} = \mathbf{w}_t - \alpha_t \mathbf{v}_t$
- 7:    $\mathbf{g}_t = \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$
- 8:    $\mathbf{v}_{t-1} = [\mathbf{v}_t + (1 - \gamma_t) \mathbf{g}_t] / \gamma_t$
- 9:    $d\mathbf{v} = d\mathbf{v} + \alpha_t d\mathbf{w}$
- 10:    $d\gamma_t = d\mathbf{v}^T (\mathbf{v}_t + \mathbf{g}_t)$
- 11:    $d\mathbf{w} = d\mathbf{w} - (1 - \gamma_t) d\mathbf{v} \nabla_{\mathbf{w}} \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$
- 12:    $d\theta = d\theta - (1 - \gamma_t) d\mathbf{v} \nabla_{\theta} \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$
- 13:    $d\mathbf{v} = \gamma_t d\mathbf{v}$
- 14: **end for**
- 15: **output** gradient of  $f(\mathbf{w}_T)$  w.r.t  $\mathbf{w}_1, \mathbf{v}_1, \gamma, \alpha$  and  $\theta$

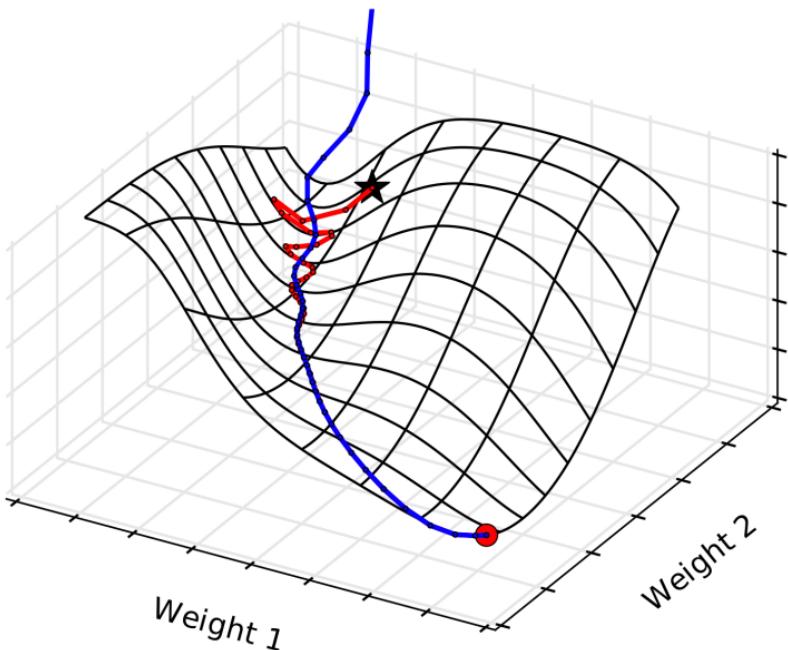
---

# Differentiating long iterations: reversibility



MacLaurin\*, Duvenaud\*, and Adams. *Gradient-based hyperparameter optimization through reversible learning.* <https://arxiv.org/pdf/1502.03492.pdf> ICML 2015.

# Differentiating long iterations: reversibility



MacLaurin\*, Duvenaud\*, and Adams. *Gradient-based hyperparameter optimization through reversible learning.* <https://arxiv.org/pdf/1502.03492.pdf> ICML 2015.

# Differentiating long iterations: reversibility

---

## Algorithm 1 Stochastic gradient descent with momentum

---

- 1: **input:** initial  $\mathbf{w}_1$ , decays  $\gamma$ , learning rates  $\alpha$ , loss function  $L(\mathbf{w}, \theta, t)$
- 2: initialize  $\mathbf{v}_1 = \mathbf{0}$
- 3: **for**  $t = 1$  to  $T$  **do**
- 4:    $\mathbf{g}_t = \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$  ▷ evaluate gradient
- 5:    $\mathbf{v}_{t+1} = \gamma_t \mathbf{v}_t - (1 - \gamma_t) \mathbf{g}_t$  ▷ update velocity
- 6:    $\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t \mathbf{v}_t$  ▷ update position
- 7: **end for**
- 8: **output** trained parameters  $\mathbf{w}_T$

---

---

## Algorithm 2 Reverse-mode differentiation of SGD

---

- 1: **input:**  $\mathbf{w}_T, \mathbf{v}_T, \gamma, \alpha$ , train loss  $L(\mathbf{w}, \theta, t)$ , loss  $f(\mathbf{w})$
- 2: initialize  $d\mathbf{v} = \mathbf{0}, d\theta = \mathbf{0}, d\alpha_t = \mathbf{0}, d\gamma = \mathbf{0}$
- 3: initialize  $d\mathbf{w} = \nabla_{\mathbf{w}} f(\mathbf{w}_T)$
- 4: **for**  $t = T$  counting down to 1 **do**
- 5:    $d\alpha_t = d\mathbf{w}^T \mathbf{v}_t$
- 6:    $\mathbf{w}_{t-1} = \mathbf{w}_t - \alpha_t \mathbf{v}_t$
- 7:    $\mathbf{g}_t = \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$
- 8:    $\mathbf{v}_{t-1} = [\mathbf{v}_t + (1 - \gamma_t) \mathbf{g}_t] / \gamma_t$
- 9:    $d\mathbf{v} = d\mathbf{v} + \alpha_t d\mathbf{w}$
- 10:    $d\gamma_t = d\mathbf{v}^T (\mathbf{v}_t + \mathbf{g}_t)$
- 11:    $d\mathbf{w} = d\mathbf{w} - (1 - \gamma_t) d\mathbf{v} \nabla_{\mathbf{w}} \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$
- 12:    $d\theta = d\theta - (1 - \gamma_t) d\mathbf{v} \nabla_{\theta} \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$
- 13:    $d\mathbf{v} = \gamma_t d\mathbf{v}$
- 14: **end for**
- 15: **output** gradient of  $f(\mathbf{w}_T)$  w.r.t  $\mathbf{w}_1, \mathbf{v}_1, \gamma, \alpha$  and  $\theta$

---

# Differentiating long iterations: reversibility

---

## Algorithm 1 Stochastic gradient descent with momentum

```
1: input: initial  $\mathbf{w}_1$ , decays  $\gamma$ , learning rates  $\alpha$ , loss function  $L(\mathbf{w}, \theta, t)$ 
2: initialize  $\mathbf{v}_1 = \mathbf{0}$ 
3: for  $t = 1$  to  $T$  do
4:    $\mathbf{g}_t = \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$             $\triangleright$  evaluate gradient
5:    $\mathbf{v}_{t+1} = \gamma_t \mathbf{v}_t - (1 - \gamma_t) \mathbf{g}_t$      $\triangleright$  update velocity
6:    $\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t \mathbf{v}_t$                  $\triangleright$  update position
7: end for
8: output trained parameters  $\mathbf{w}_T$ 
```

---

finite precision means information loss!

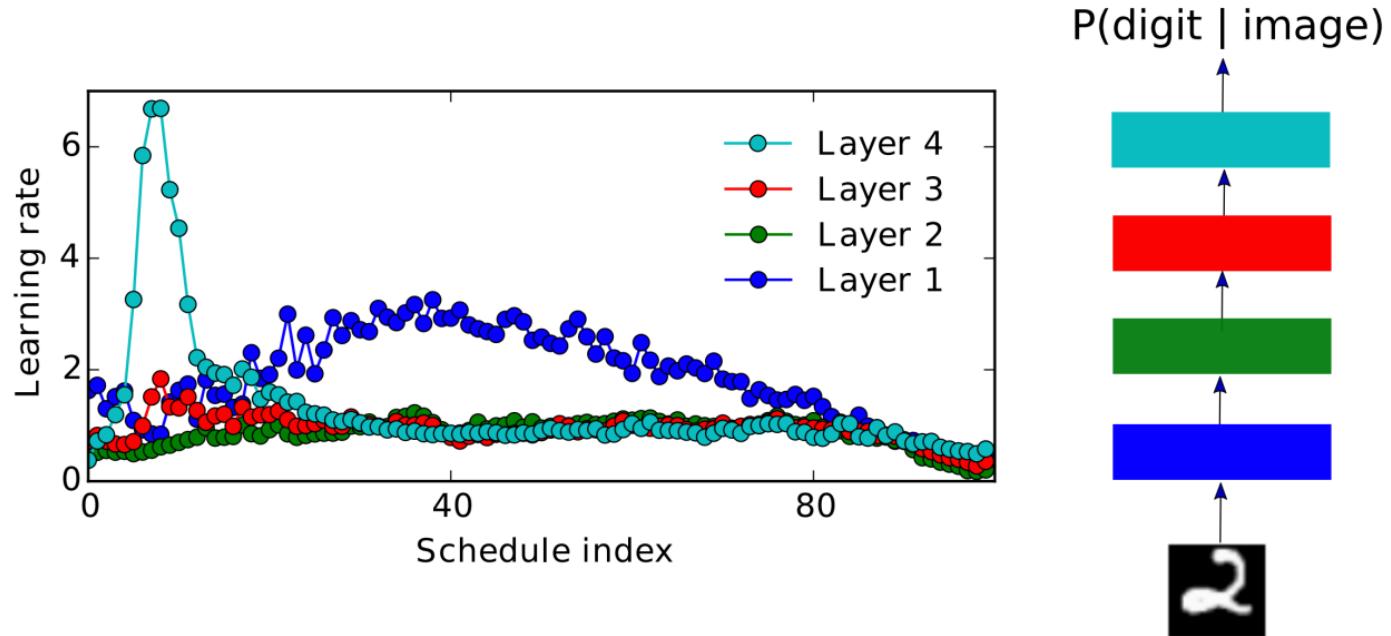
---

## Algorithm 2 Reverse-mode differentiation of SGD

```
1: input:  $\mathbf{w}_T, \mathbf{v}_T, \gamma, \alpha$ , train loss  $L(\mathbf{w}, \theta, t)$ , loss  $f(\mathbf{w})$ 
2: initialize  $d\mathbf{v} = \mathbf{0}, d\theta = \mathbf{0}, d\alpha_t = \mathbf{0}, d\gamma = \mathbf{0}$ 
3: initialize  $d\mathbf{w} = \nabla_{\mathbf{w}} f(\mathbf{w}_T)$ 
4: for  $t = T$  counting down to 1 do
5:    $d\alpha_t = d\mathbf{w}^T \mathbf{v}_t$ 
6:    $\mathbf{w}_{t-1} = \mathbf{w}_t - \alpha_t \mathbf{v}_t$             $\left. \begin{array}{l} \text{exactly reverse} \\ \text{gradient descent} \end{array} \right\} \text{operations}$ 
7:    $\mathbf{g}_t = \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$ 
8:    $\mathbf{v}_{t-1} = [\mathbf{v}_t + (1 - \gamma_t) \mathbf{g}_t] / \gamma_t$ 
9:    $d\mathbf{v} = d\mathbf{v} + \alpha_t d\mathbf{w}$ 
10:   $d\gamma_t = d\mathbf{v}^T (\mathbf{v}_t + \mathbf{g}_t)$ 
11:   $d\mathbf{w} = d\mathbf{w} - (1 - \gamma_t) d\mathbf{v} \nabla_{\mathbf{w}} \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$ 
12:   $d\theta = d\theta - (1 - \gamma_t) d\mathbf{v} \nabla_{\theta} \nabla_{\mathbf{w}} L(\mathbf{w}_t, \theta, t)$ 
13:   $d\mathbf{v} = \gamma_t d\mathbf{v}$ 
14: end for
15: output gradient of  $f(\mathbf{w}_T)$  w.r.t  $\mathbf{w}_1, \mathbf{v}_1, \gamma, \alpha$  and  $\theta$ 
```

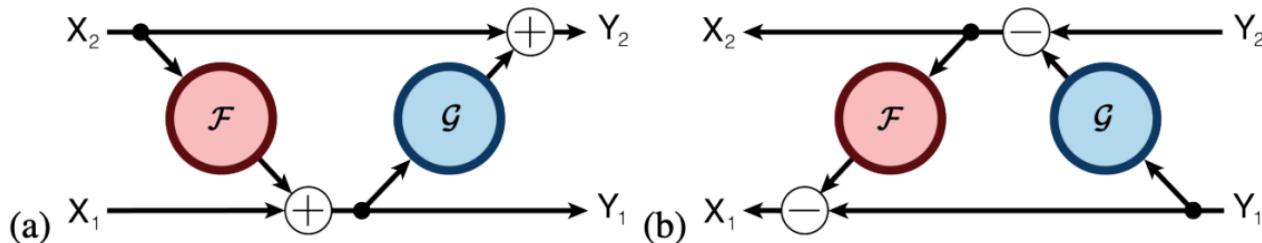
---

# Differentiating long iterations: reversibility



MacLaurin\*, Duvenaud\*, and Adams. *Gradient-based hyperparameter optimization through reversible learning*. <https://arxiv.org/pdf/1502.03492.pdf> ICML 2015.

# Differentiating long iterations: reversibility



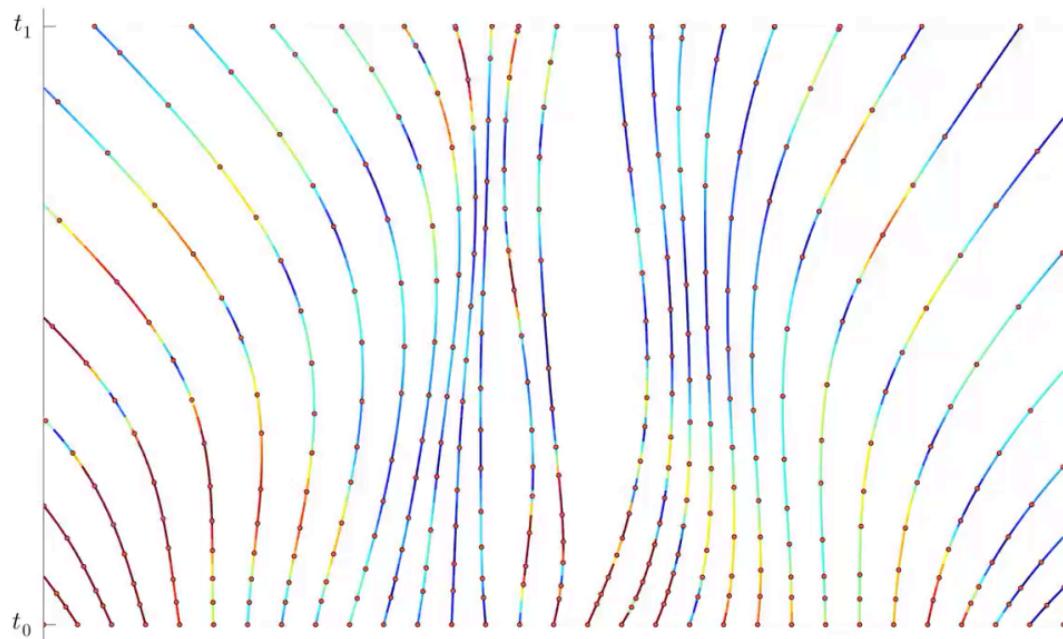
```
from jax import reversible
```

Gomez, Ren, Urtasun, and Grosse. *The reversible residual network: backpropagation without storing activations.* <https://arxiv.org/pdf/1707.04585.pdf> NeurIPS 2017.

Kitaev, Kaiser, and Levskaya. *Reformer: the efficient Transformer.* <https://arxiv.org/pdf/2001.04451.pdf> ICLR 2019.

# Very high-order AD with jet

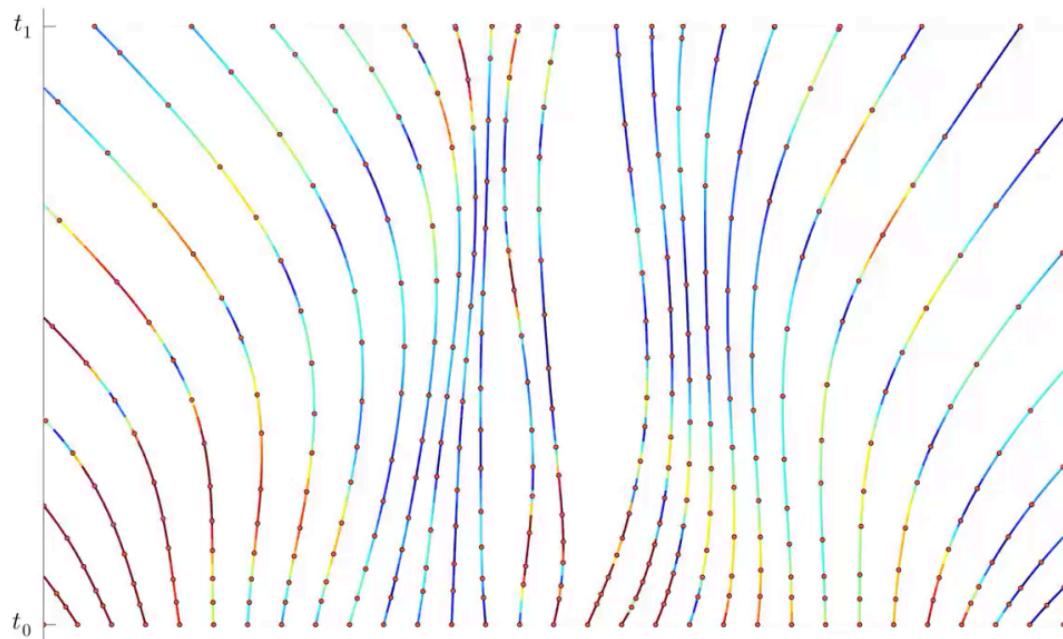
Neural ODEs can be **slow** at test time...



Idea: speed them up by  $\overset{z}{\text{regularizing}}$  higher derivatives,  
learning ODEs that are easy to solve!

# Very high-order AD with jet

Neural ODEs can be **slow** at test time...



Idea: speed them up by  $\overset{\sim}{\text{regularizing}}$  higher derivatives,  
learning ODEs that are easy to solve!

# Very high-order AD with `jet`

$$f = g \circ h$$

given  $(h(x), \partial h(x)[v])$   
compute  $(f(x), \partial f(x)[v])$

using  $g$  and  $\partial g$

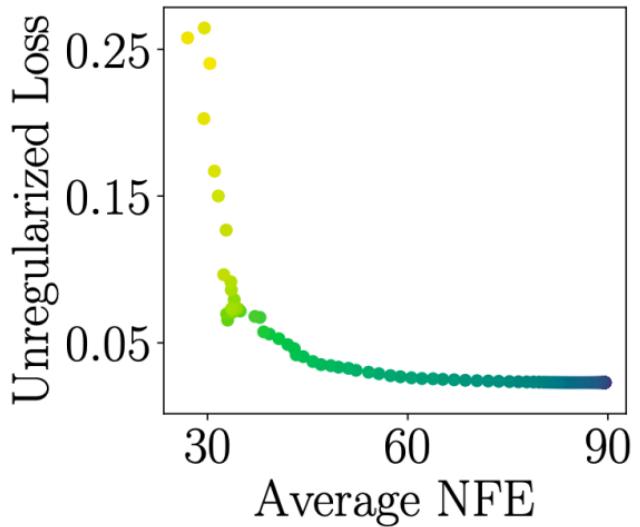
# Very high-order AD with `jet`

$$f = g \circ h$$

given  $(h(x), \partial h(x)[v], \frac{1}{2}\partial^2 h(x)[v, v], \dots, \frac{1}{K!}\partial^K h(x)[v, \dots, v])$   
compute  $(f(x), \partial f(x)[v], \frac{1}{2}\partial^2 f(x)[v, v], \dots, \frac{1}{K!}\partial^K f(x)[v, \dots, v])$

using  $g, \partial g, \partial^2 g, \dots$ , and  $\partial^K g$

# Very high-order AD with jet



(a) MNIST Classification

# Very high-order AD with jet



<https://twitter.com/DavidDuvenaud/status/1284181673496776706>

# Thanks!

