



# Automatic differentiation with JAX



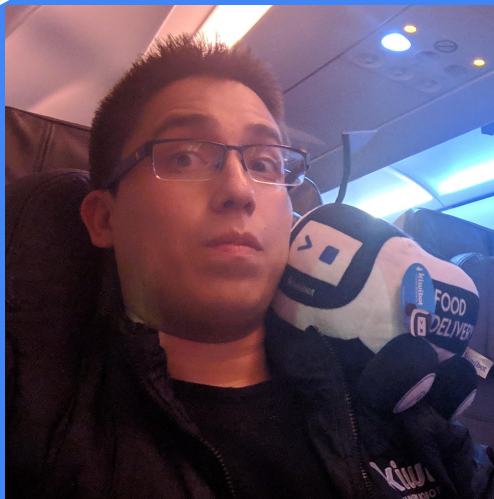
**David Cardozo**

ML Engineer



@\_davidcardozo

Please  
reach  
out with  
questions



# Our Awesome new world

- TensorFlow, Stan, Theano, Edward, PyTorch, MinPy
- Only need to specify forward model
- Autodiff + optimization / inference done for you
- loops? branching? recursion? closures? data structures?
- debugger?
- a second compiler/interpreter to satisfy



## JAX Overview

- What is JAX?
- Why JAX?
- Overview of JAX

## Jacobians and the chain rule

Forward and reverse accumulation

## Autograd implementation

Fully closed tracing autodiff in python

## References

### Models

Image Transformer  
DCGAN



**JAX**



# JAX

- differentiates native Python code
- handles most of Numpy + Scipy
- loops, branching, recursion, closures
- arrays, tuples, lists, dicts, classes, ...
- derivatives of derivatives
- a one-function API
- small and easy to extend



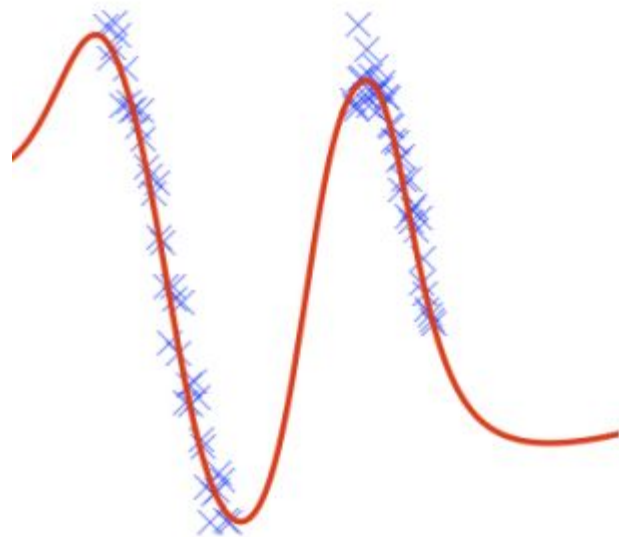
**XLA**

**The Linear  
algebra  
compiler**



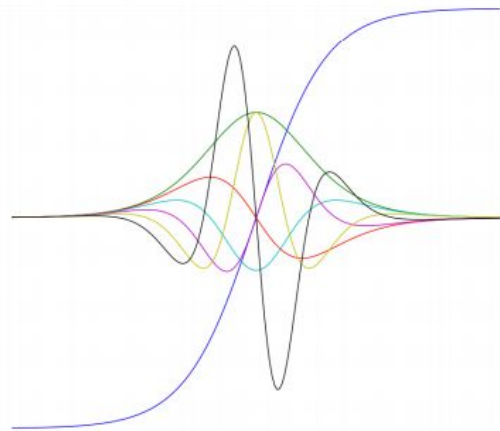
■ ACCELERATED LINEAR ALGEBRA ■

```
1 import jax.numpy as np
2 import jax.random as npr
3 from jax import grad
4
5 def predict(weights, inputs):
6     for W, b in weights:
7         outputs = np.dot(inputs, W) + b
8         inputs = np.tanh(outputs)
9     return outputs
10
11 def init_params(scale, sizes):
12     return [(npr.randn(m, n) * scale,
13             npr.randn(n) * scale)
14             for m, n in zip(sizes[:-1], sizes[1:])]
15
16 def logprob_fun(params, inputs, targets):
17     preds = predict(weights, inputs)
18     return np.sum((preds - targets)**2)
19
20 gradient_fun = grad(logprob_fun)
```



# JAX Examples

```
1 import jax.numpy as np
2 from jax import grad, vmap
3 import matplotlib.pyplot as plt
4
5 x = np.linspace(-7, 7, 200)
6 plt.plot(x, np.tanh(x),
7          x, vmap(grad(np.tanh))(x),
8          x, vmap(grad(grad(np.tanh)))(x),
9          x, vmap(grad(grad(grad(np.tanh))))(x),
10         x, vmap(grad(grad(grad(grad(np.tanh)))))(x), )
```





# Hessians and HVP

```
1  from jax import grad, jacobian
2
3  def hessian(fun, argnum=0):
4      ... return jacobian(jacobian(fun, argnum), argnum)
5  def hvp(fun):
6      ... def grad_dot_vector(arg, vector):
7          ...     ... return np.dot(grad(fun)(arg), vector)
8      ... return grad(grad_dot_vector)
```

$$\nabla^2 f(x) \cdot v = \nabla_x (\nabla_x f(x) \cdot v)$$



# Black-box inference in a tweet



**Ryan Adams** @ryan\_p\_adams · 7 Nov 2015

@DavidDuvenaud

```
def elbo(p, lp, D, N):  
    v=exp(p[D:])  
    s=randn(N,D)*sqrt(v)+p[:D]  
    return mvn.entropy(0, diag(v))+mean(lp(s))  
gf = grad(elbo)
```



1



7



22

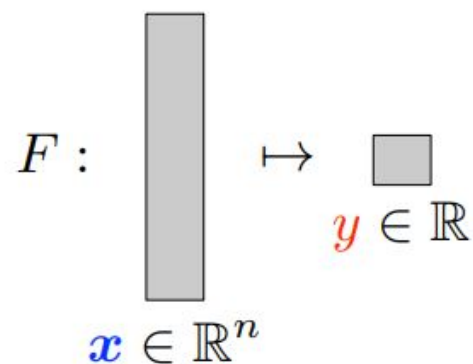


# Jacobians and the chain rule

Forward and  
reverse  
accumulation



$$F : \mathbb{R}^n \rightarrow \mathbb{R}$$



$$F = D \circ C \circ B \circ A$$

$$y = F(x) = D(C(B(A(x))))$$

$$y = D(c), \quad c = C(b), \quad b = B(a), \quad a = A(x)$$



$$\textcolor{red}{y} = D(\textcolor{red}{c}), \quad \textcolor{red}{c} = C(\textcolor{red}{b}), \quad \textcolor{red}{b} = B(\textcolor{red}{a}), \quad \textcolor{red}{a} = A(\textcolor{blue}{x})$$

$$F'(\textcolor{blue}{x}) = \frac{\partial \textcolor{red}{y}}{\partial \textcolor{blue}{x}} = \left[ \frac{\partial \textcolor{red}{y}}{\partial x_1} \quad \cdots \quad \frac{\partial \textcolor{red}{y}}{\partial x_n} \right]$$

$$y = D(c), \quad c = C(b), \quad b = B(a), \quad a = A(x)$$

$$F'(x) = \frac{\partial y}{\partial x} = \left[ \frac{\partial y}{\partial x_1} \quad \cdots \quad \frac{\partial y}{\partial x_n} \right]$$

$$F'(x) = \begin{matrix} \frac{\partial y}{\partial c} & \frac{\partial c}{\partial b} & \frac{\partial b}{\partial a} & \frac{\partial a}{\partial x} \end{matrix}$$

$$y = D(c), \quad c = C(b), \quad b = B(a), \quad a = A(x)$$

$$F'(x) = \frac{\partial y}{\partial x} = \left[ \frac{\partial y}{\partial x_1} \quad \cdots \quad \frac{\partial y}{\partial x_n} \right]$$

$$F'(x) = \begin{matrix} \frac{\partial y}{\partial c} & \frac{\partial c}{\partial b} & \frac{\partial b}{\partial a} & \frac{\partial a}{\partial x} \end{matrix}$$

$$\frac{\partial y}{\partial c} = D'(c)$$



$$y = D(c), \quad c = C(b), \quad b = B(a), \quad a = A(x)$$

$$F'(x) = \frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \cdots & \frac{\partial y}{\partial x_n} \end{bmatrix}$$

$$F'(x) = \begin{matrix} \frac{\partial y}{\partial c} & \frac{\partial c}{\partial b} & \frac{\partial b}{\partial a} & \frac{\partial a}{\partial x} \end{matrix}$$

$$\frac{\partial y}{\partial c} = D'(c)$$

$$\frac{\partial c}{\partial b} = C'(b)$$



**c** **b**



$$y = D(c), \quad c = C(b), \quad b = B(a), \quad a = A(x)$$

$$F'(x) = \frac{\partial y}{\partial x} = \left[ \frac{\partial y}{\partial x_1} \quad \cdots \quad \frac{\partial y}{\partial x_n} \right]$$

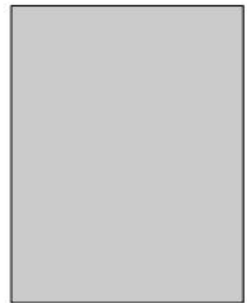
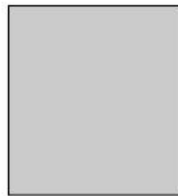
$$F'(x) = \frac{\partial y}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial x}$$

$$\frac{\partial y}{\partial c} = D'(c)$$

$$\frac{\partial c}{\partial b} = C'(b)$$

$$\frac{\partial b}{\partial a} = B'(a)$$

$$\frac{\partial a}{\partial x} = A'(x)$$



$$F'(\mathbf{x}) = \frac{\partial y}{\partial \mathbf{c}} \left( \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \left( \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \quad \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \right) \right)$$



$$F'(\mathbf{x}) = \frac{\partial y}{\partial \mathbf{c}} \underbrace{\left( \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \left( \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \quad \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \right) \right)}$$

$$\frac{\partial \mathbf{b}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial b_1}{\partial x_1} & \dots & \frac{\partial b_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial b_m}{\partial x_1} & \dots & \frac{\partial b_m}{\partial x_n} \end{bmatrix}$$

$$F'(\mathbf{x}) = \frac{\partial y}{\partial \mathbf{c}} \underbrace{\left( \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \left( \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \quad \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \right) \right)}$$

$$\frac{\partial \mathbf{b}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial b_1}{\partial x_1} & \dots & \frac{\partial b_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial b_m}{\partial x_1} & \dots & \frac{\partial b_m}{\partial x_n} \end{bmatrix}$$

Forward  
accumulation



$$F'(\mathbf{x}) = \frac{\partial \mathbf{y}}{\partial \mathbf{c}} \left( \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \underbrace{\left( \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \quad \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \right)} \right)$$

$$\frac{\partial \mathbf{b}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial b_1}{\partial x_1} & \dots & \frac{\partial b_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial b_m}{\partial x_1} & \dots & \frac{\partial b_m}{\partial x_n} \end{bmatrix}$$

Forward  
accumulation

$$F'(\mathbf{x}) = \left( \left( \frac{\partial \mathbf{y}}{\partial \mathbf{c}} \quad \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \right) \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \right) \frac{\partial \mathbf{a}}{\partial \mathbf{x}}$$



$$F'(\mathbf{x}) = \frac{\partial y}{\partial \mathbf{c}} \underbrace{\left( \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \left( \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \right) \right)}$$

$$\frac{\partial \mathbf{b}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial b_1}{\partial x_1} & \cdots & \frac{\partial b_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial b_m}{\partial x_1} & \cdots & \frac{\partial b_m}{\partial x_n} \end{bmatrix}$$

Forward  
accumulation

$$F'(\mathbf{x}) = \underbrace{\left( \left( \frac{\partial y}{\partial \mathbf{c}} \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \right) \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \right)} \frac{\partial \mathbf{a}}{\partial \mathbf{x}}$$



$$\frac{\partial y}{\partial \mathbf{b}} = \begin{bmatrix} \frac{\partial y}{\partial b_1} & \cdots & \frac{\partial y}{\partial b_m} \end{bmatrix}$$

$$F'(\mathbf{x}) = \frac{\partial \mathbf{y}}{\partial \mathbf{c}} \left( \underbrace{\frac{\partial \mathbf{c}}{\partial \mathbf{b}} \left( \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \quad \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \right)}_{\text{Forward accumulation}} \right)$$

$$\frac{\partial \mathbf{b}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial b_1}{\partial x_1} & \cdots & \frac{\partial b_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial b_m}{\partial x_1} & \cdots & \frac{\partial b_m}{\partial x_n} \end{bmatrix}$$

Forward  
accumulation

$$F'(\mathbf{x}) = \left( \underbrace{\left( \frac{\partial \mathbf{y}}{\partial \mathbf{c}} \quad \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \right) \frac{\partial \mathbf{b}}{\partial \mathbf{a}}}_{\text{Reverse accumulation}} \right) \frac{\partial \mathbf{a}}{\partial \mathbf{x}}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{b}} = \begin{bmatrix} \frac{\partial y}{\partial b_1} & \cdots & \frac{\partial y}{\partial b_m} \end{bmatrix}$$

Reverse  
accumulation



$$F'(\mathbf{x}) \mathbf{v} = \frac{\partial \mathbf{y}}{\partial \mathbf{c}} \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \mathbf{v}$$

$$F'(\mathbf{x}) \mathbf{v} = \frac{\partial \mathbf{y}}{\partial \mathbf{c}} \left( \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \left( \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \left( \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \mathbf{v} \right) \right) \right)$$

Forward accumulation  $\leftrightarrow$  Jacobian-vector products

Build Jacobian one column at a time

$$F'(\mathbf{x}) = \frac{\partial \mathbf{y}}{\partial \mathbf{c}} \left( \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \left( \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \left( \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial \mathbf{x}} \right) \right) \right)$$





$$\mathbf{v}^\top F'(\mathbf{x}) = \mathbf{v}^\top \frac{\partial \mathbf{y}}{\partial \mathbf{c}} \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial \mathbf{x}}$$

$$\mathbf{v}^\top F'(\mathbf{x}) = \left( \left( \left( \mathbf{v}^\top \frac{\partial \mathbf{y}}{\partial \mathbf{c}} \right) \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \right) \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \right) \frac{\partial \mathbf{a}}{\partial \mathbf{x}}$$

Reverse accumulation  $\leftrightarrow$  vector-Jacobian products

Build Jacobian one row at a time

$$F'(\mathbf{x}) = \left( \left( \left( \frac{\partial \mathbf{y}}{\partial \mathbf{y}} \quad \frac{\partial \mathbf{y}}{\partial \mathbf{c}} \right) \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \right) \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \right) \frac{\partial \mathbf{a}}{\partial \mathbf{x}}$$



# Forward and reverse accumulation

- Forward accumulation
  - Jacobian-vector products
  - “push-forward”
  - build Jacobian matrix one column at a time
- Reverse accumulation
  - vector-Jacobian products
  - “pull-back”
  - build Jacobian matrix one row at a time



# Non-chain composition

Fan-in

$$y = F(\mathbf{x}_1, \mathbf{x}_2)$$

$$\frac{\partial y}{\partial \mathbf{x}_1} = F'_1(\mathbf{x}_1, \mathbf{x}_2)$$

$$\frac{\partial y}{\partial \mathbf{x}_2} = F'_2(\mathbf{x}_1, \mathbf{x}_2)$$

Fan-out

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

$$G'(\mathbf{x}) = \begin{bmatrix} I \\ I \end{bmatrix}$$

$$\mathbf{v}^\top G'(\mathbf{x}) = \begin{bmatrix} \mathbf{v}_1^\top & \mathbf{v}_2^\top \end{bmatrix} \begin{bmatrix} I \\ I \end{bmatrix} = \mathbf{v}_1^\top + \mathbf{v}_2^\top$$



# Autodiff differentiation

## Read and generate source code ahead-of-time

---

- Source and target language could be python
- Monitor function execution at runtime

## Monitor function execution at runtime

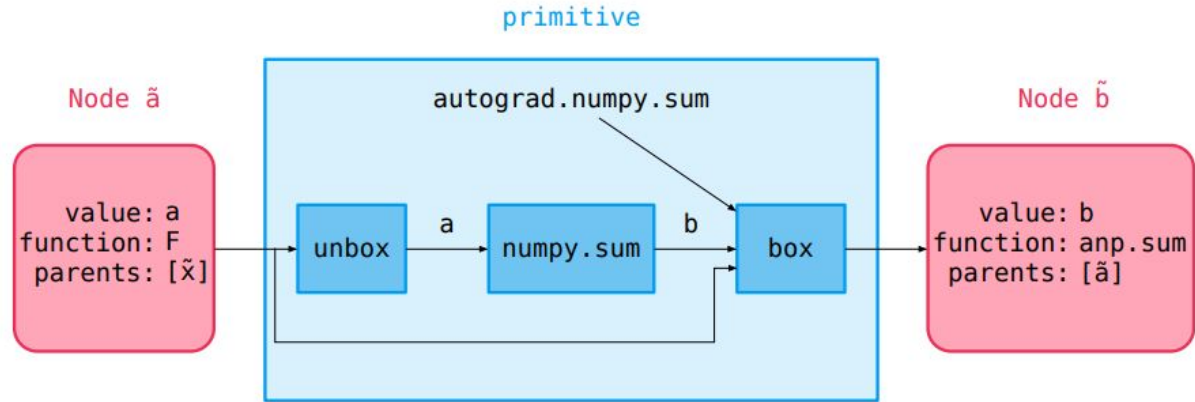
---

Strategy at Higher level languages

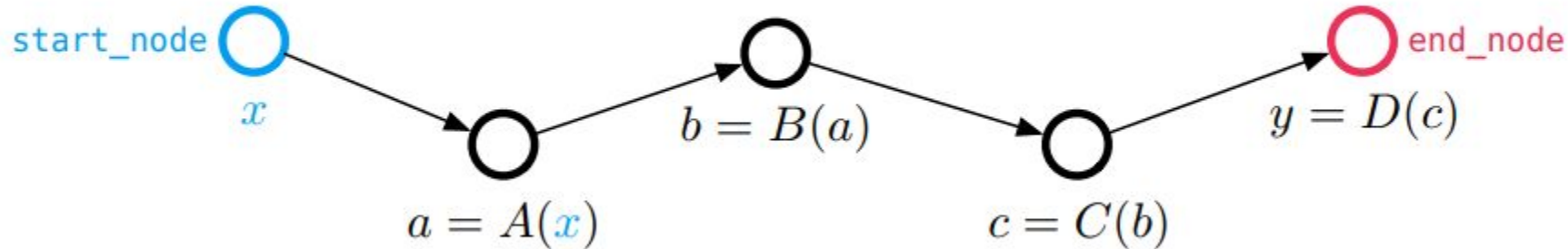
Pytorch mynpy



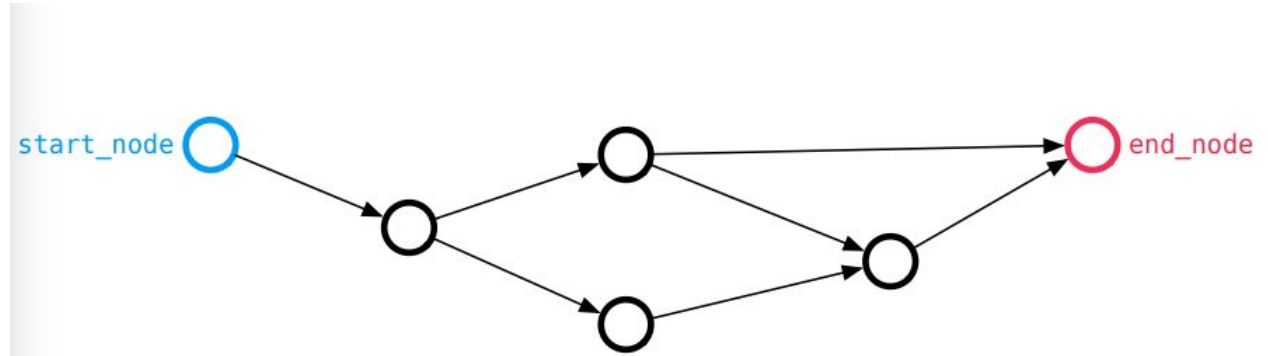
# Tracing the composition of primitive functions



# Tracing on chains



# Graphs on DAGs



No control flow!



# Defining a vector-jacobian product for each primitive

vector-Jacobian product



$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial a} \cdot A'(x)$$



$x$

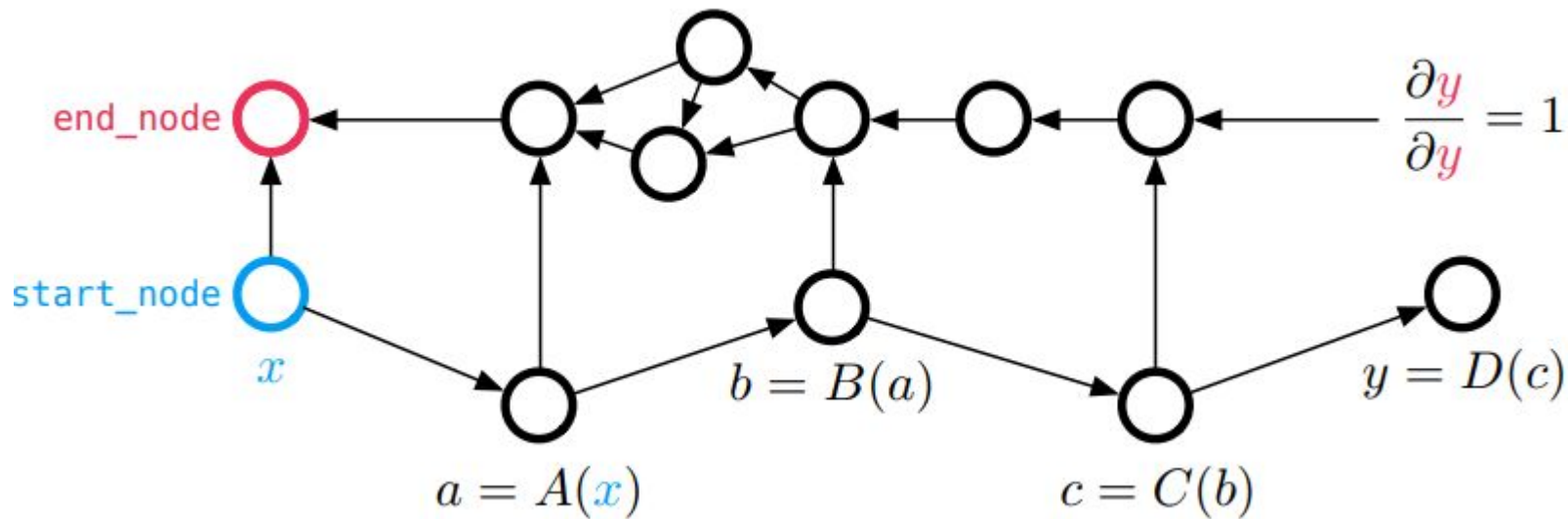


$a = A(x)$





# Composing VJPs backwards



# Autograd's ingredients

1. Tracing the composition of primitive functions

`Node`, `primitive`, `forward_pass`

2. Defining a vector-Jacobian product (VJP) operator for each primitive

`defvjp`

3. Composing VJPs backward

`backward_pass`, `make_vjp`, `grad`



# JAX examples



# JIT in JAX

```
1  import jax.numpy as np
2  from jax import jit
3
4  def slow_f(x):
5      ·· # Element-wise ops see a large benefit from fusion
6      ·· return x * x + x * 2.0
7
8  x = np.ones((5000, 5000))
9  fast_f = jit(slow_f)
10 %timeit -n10 -r3 fast_f(x) ·· # ~ 4.5 ms / loop on Titan X
11 %timeit -n10 -r3 slow_f(x) ·· # ~ 14.5 ms / loop (also on GPU via JAX)
```



# Stax in JAX

```
1 import jax.numpy as np
2 from jax import random
3 from jax.experimental import stax
4 from jax.experimental.stax import Conv, Dense, MaxPool, Relu, Flatten, LogSoftmax
5
6 # Use stax to set up network initialization and evaluation functions
7 net_init, net_apply = stax.serial(
8     .... Conv(32, (3, 3), padding='SAME'), Relu,
9     .... Conv(64, (3, 3), padding='SAME'), Relu,
10    .... MaxPool((2, 2)), Flatten,
11    .... Dense(128), Relu,
12    .... Dense(10), LogSoftmax,
13 )
14
15 # Initialize parameters, not committing to a batch shape
16 rng = random.PRNGKey(0)
17 in_shape = (-1, 28, 28, 1)
18 out_shape, net_params = net_init(rng, in_shape)
19
20 # Apply network to dummy inputs
21 inputs = np.zeros((128, 28, 28, 1))
22 predictions = net_apply(net_params, inputs)
```



# Conclusions

`trace` :: Traceable a b -> Abstract a -> Jaxpr a b

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

Traceable a b

```
{ lambda inputs ; params.  
  let (b, c) = params  
      (d, e) = b  
      g = dot inputs d  
      h = add g e  
      i = tanh h  
      (j, k) = c  
      l = dot i j  
      m = add l k  
  in m }
```

Jaxpr a b

fwd\_deriv,  
rev\_deriv,  
xla\_compile,  
batch

`lift` :: Jaxpr a b -> Traceable a b



# Thank you !

**David Cardozo**

ML Engineer

 @\_davidcardozo

*Please  
reach  
out with  
questions*

