# CS 4803 / 7643: Deep Learning

Topics:

– Automatic Differentiation
  – Patterns in backprop
  – Jacobians in FC+ReLU NNs

Dhruv Batra

Georgia Tech

# Administrativia

- HW1 Reminder
  - Due: 09/26, 11:55pm

- Project Teams Google Doc
  - https://docs.google.com/spreadsheets/d/1ouD6ctaemV_3nb2MQHs7rUOAaW9DFLu8I5Zd3yOFs7E/edit?usp=sharing
  - Project Title
  - 1-3 sentence project summary TL;DR
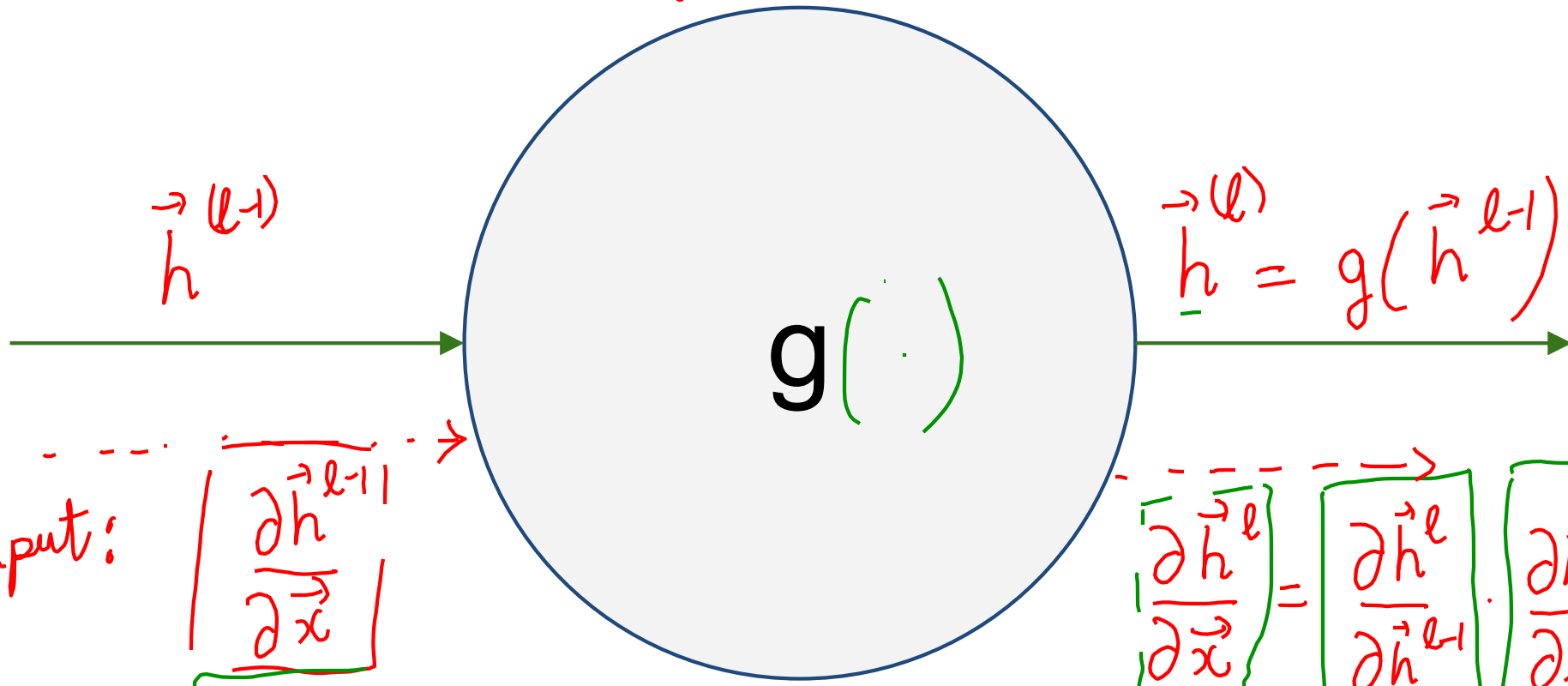  - Team member names

# Recap from last time

# Deep Learning = Differentiable Programming

- ## Computation = Graph
  - Input = Data + Parameters
  - Output = Loss
  - Scheduling = Topological ordering

- ## Auto-Diff
  - A family of algorithms for
    implementing chain-rule on computation graphs

# Forward mode AD

layer $l$

$\vec{h}^{(l-1)}$

$\vec{h}^{(l)} = g(\vec{h}^{l-1})$

$g(\cdot)$

Input: $\left| \frac{\partial \vec{h}^{l-1}}{\partial \vec{x}} \right|$

$\frac{\partial \vec{h}^{l}}{\partial \vec{x}} = \frac{\partial \vec{h}^{l}}{\partial \vec{h}^{l-1}} \cdot \frac{\partial \vec{h}^{l-1}}{\partial \vec{x}}$

Jacobian   Input
of $g$

# Reverse mode AD

$\vec{y} = A\vec{x}$

$\dfrac{\partial \vec{y}}{\partial \vec{x}} = A$

$\vec{h}^{\ell} = W \vec{h}^{\ell-1}$

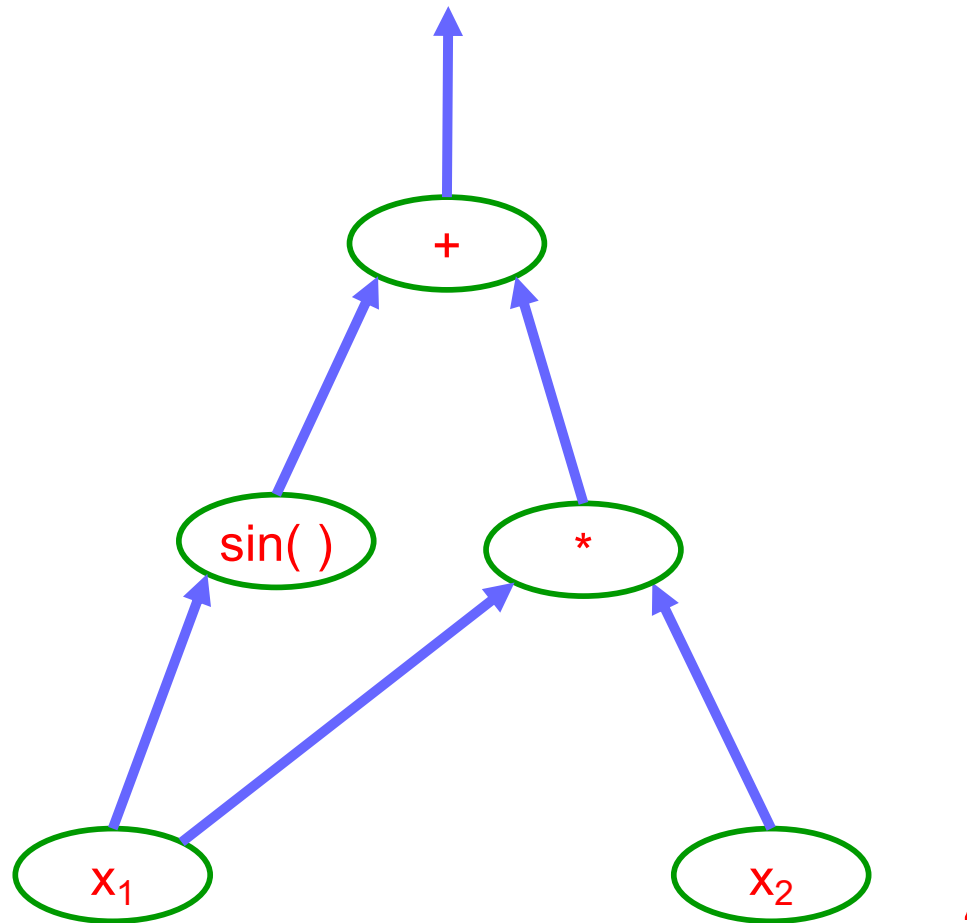$\dfrac{\partial \vec{h}^{\ell}}{\partial \vec{h}^{\ell-1}} = W$

Goal: $\boxed{\dfrac{\partial L}{\partial \vec{x}}}$

$\vec{h}^{\ell-1}$

$g$

$\vec{h}^{\ell} = g(\vec{h}^{\ell-1})$

$\boxed{\dfrac{\partial L}{\partial \vec{h}^{\ell}}}$ Input

Output

$\dfrac{\partial L}{\partial \vec{h}^{\ell-1}} = \boxed{\dfrac{\partial L}{\partial \vec{h}^{\ell}}} \boxed{\dfrac{\partial \vec{h}^{\ell}}{\partial \vec{h}^{\ell-1}}}$
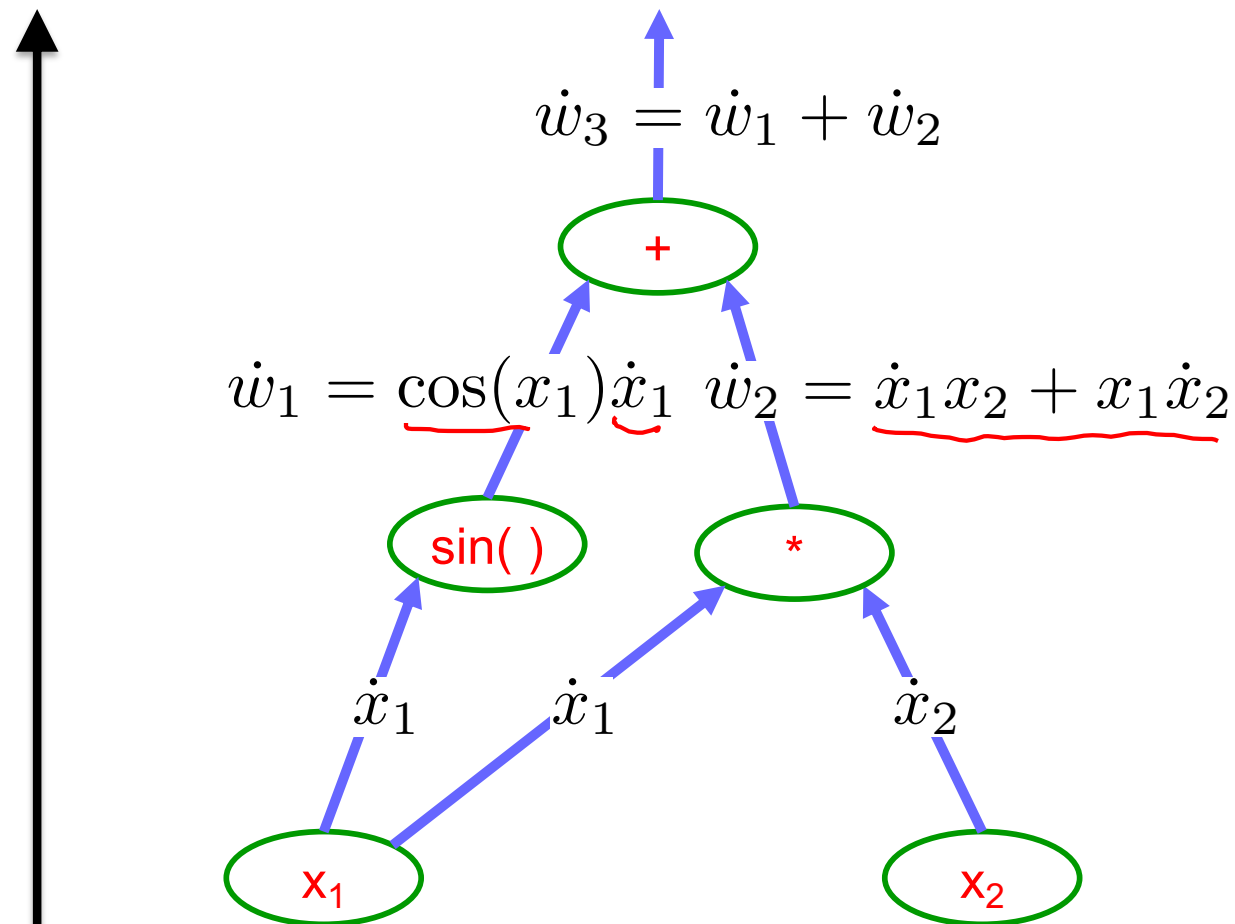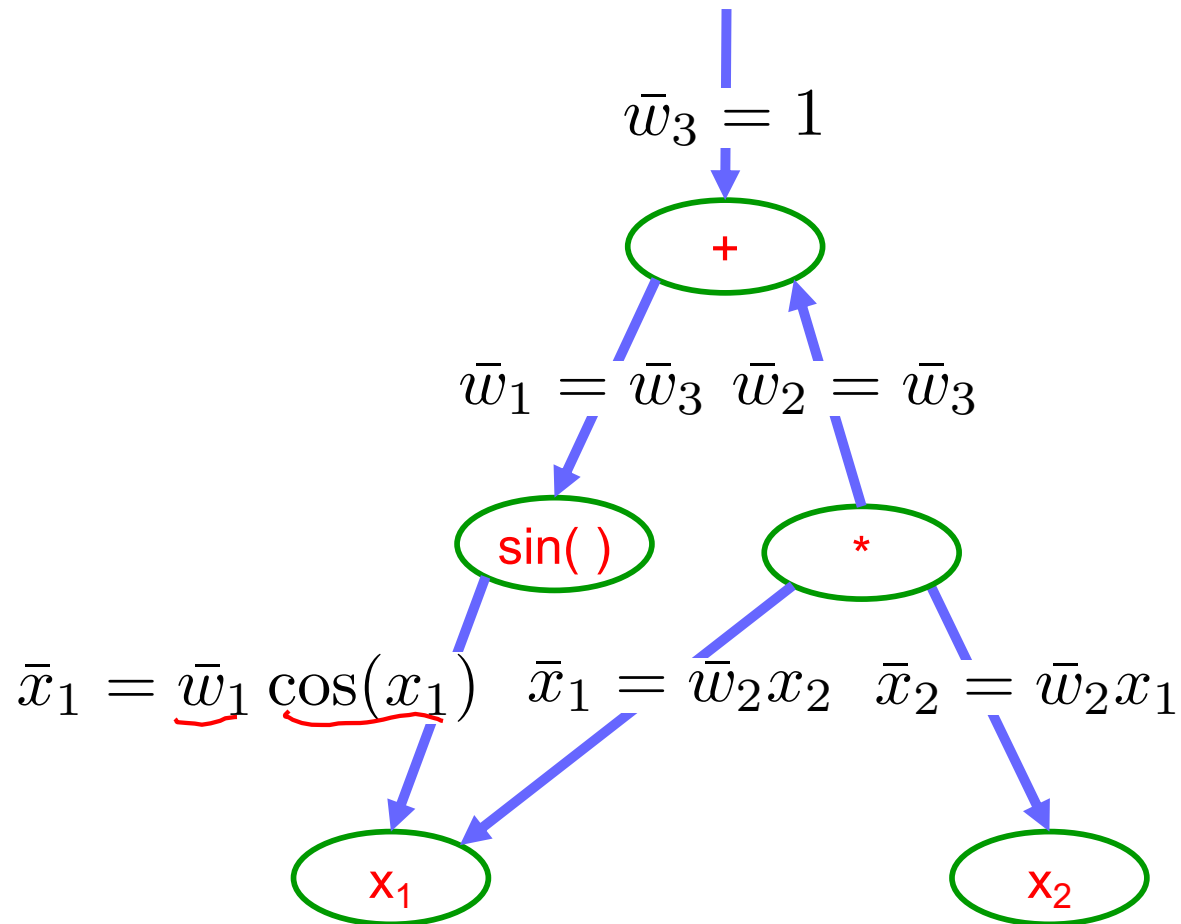
Input    Jacobian of $g$

# Example: Forward mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

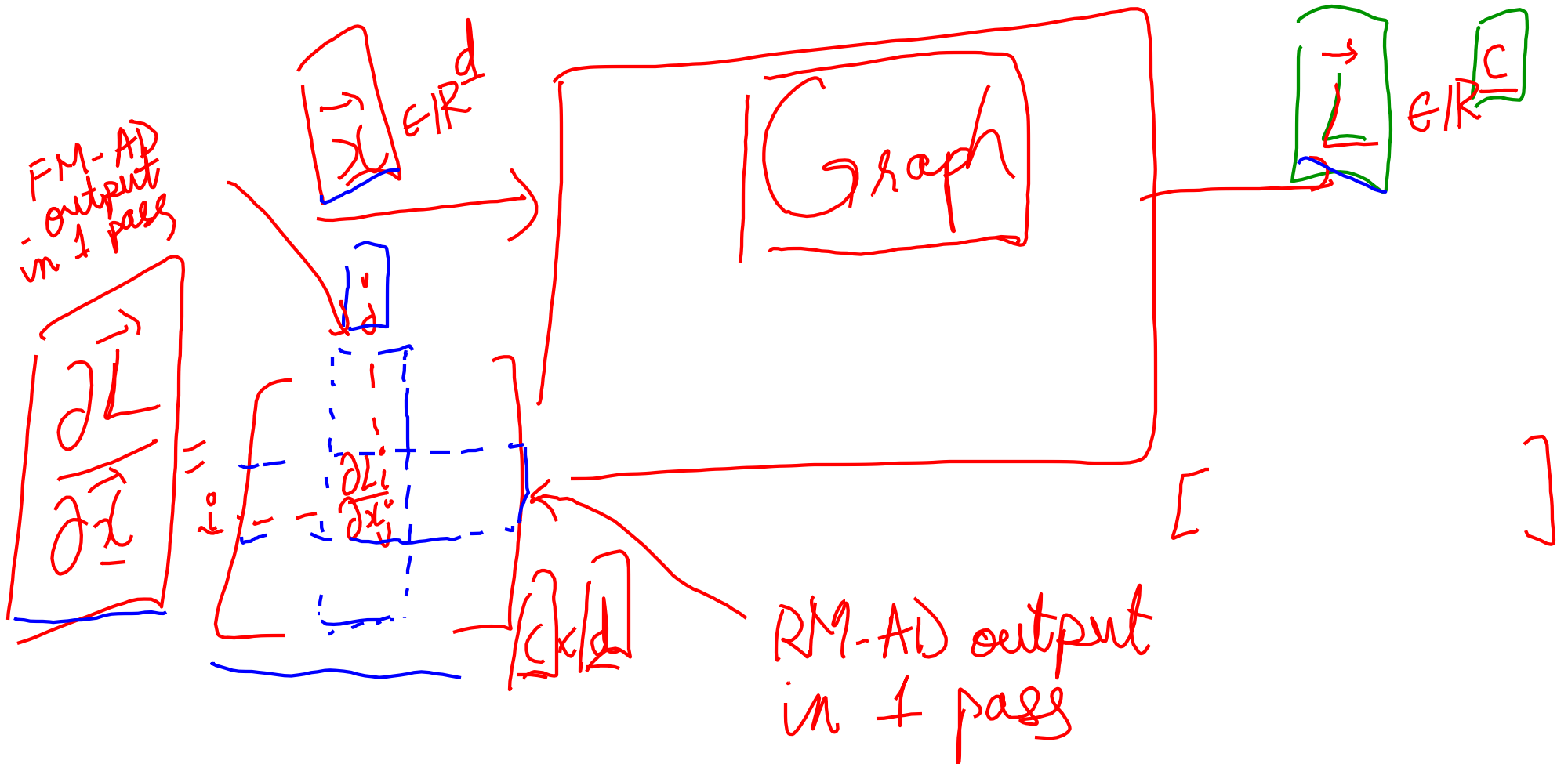# Example: Forward mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$



$$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$$

$$+$$

$$\dot{w}_1 = \cos(x_1)\dot{x}_1 \quad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$$

$$\sin(\ )$$

$$*$$

$$\dot{x}_1 \qquad \dot{x}_1 \qquad \dot{x}_2$$

$$x_1 \qquad \qquad x_2$$

# Example: Reverse mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

# Forward mode vs Reverse Mode

- x → Graph → L
- Intuition of Jacobian
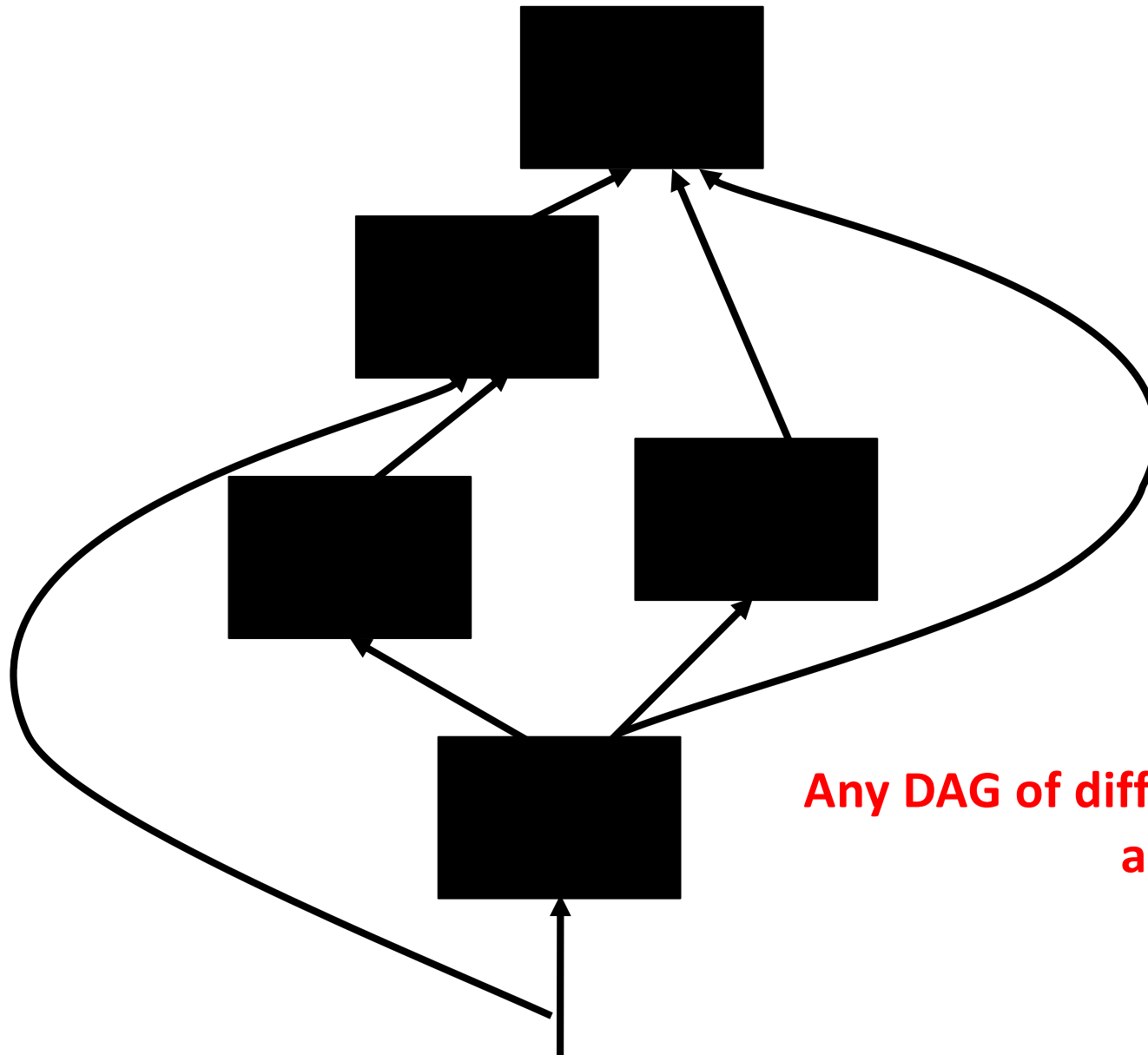
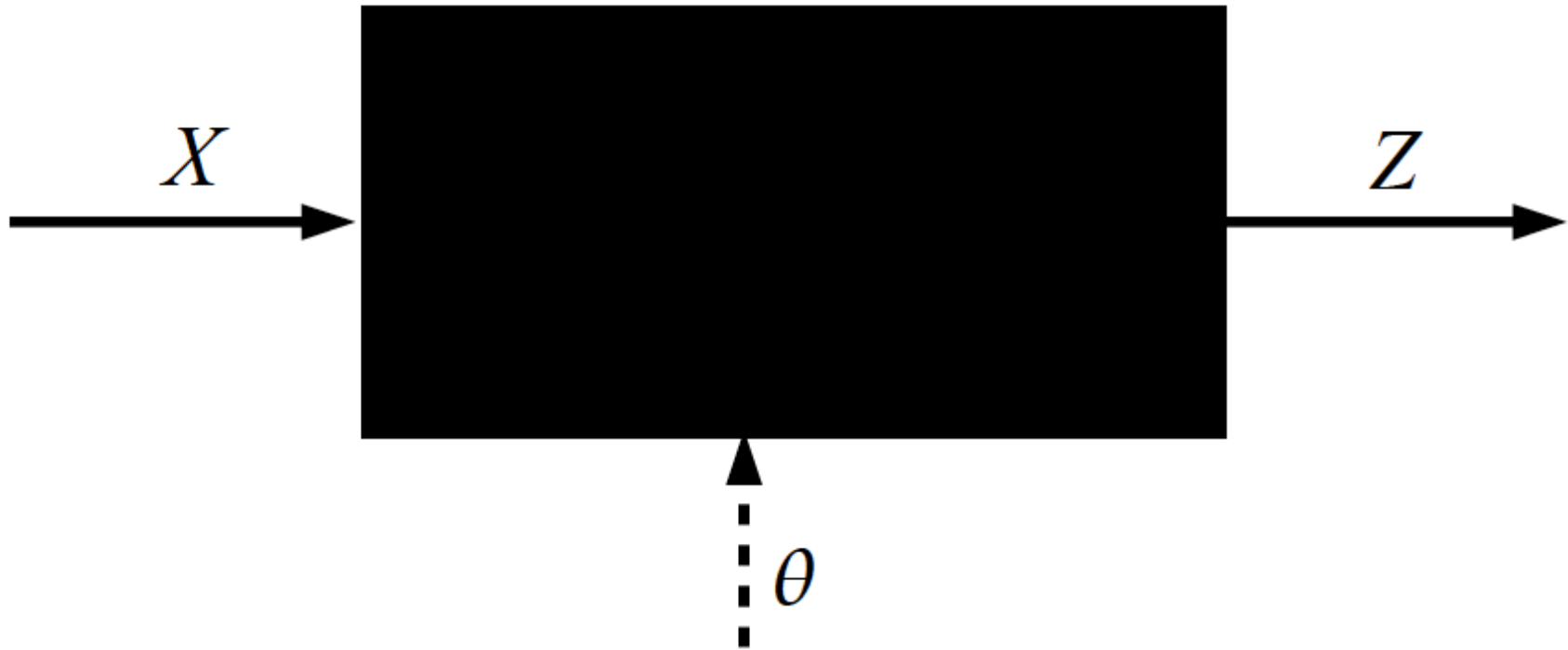# Forward mode vs Reverse Mode

- What are the differences?

- Which one is faster to compute?
  - Forward or backward?

- Which one is more memory efficient (less storage)?
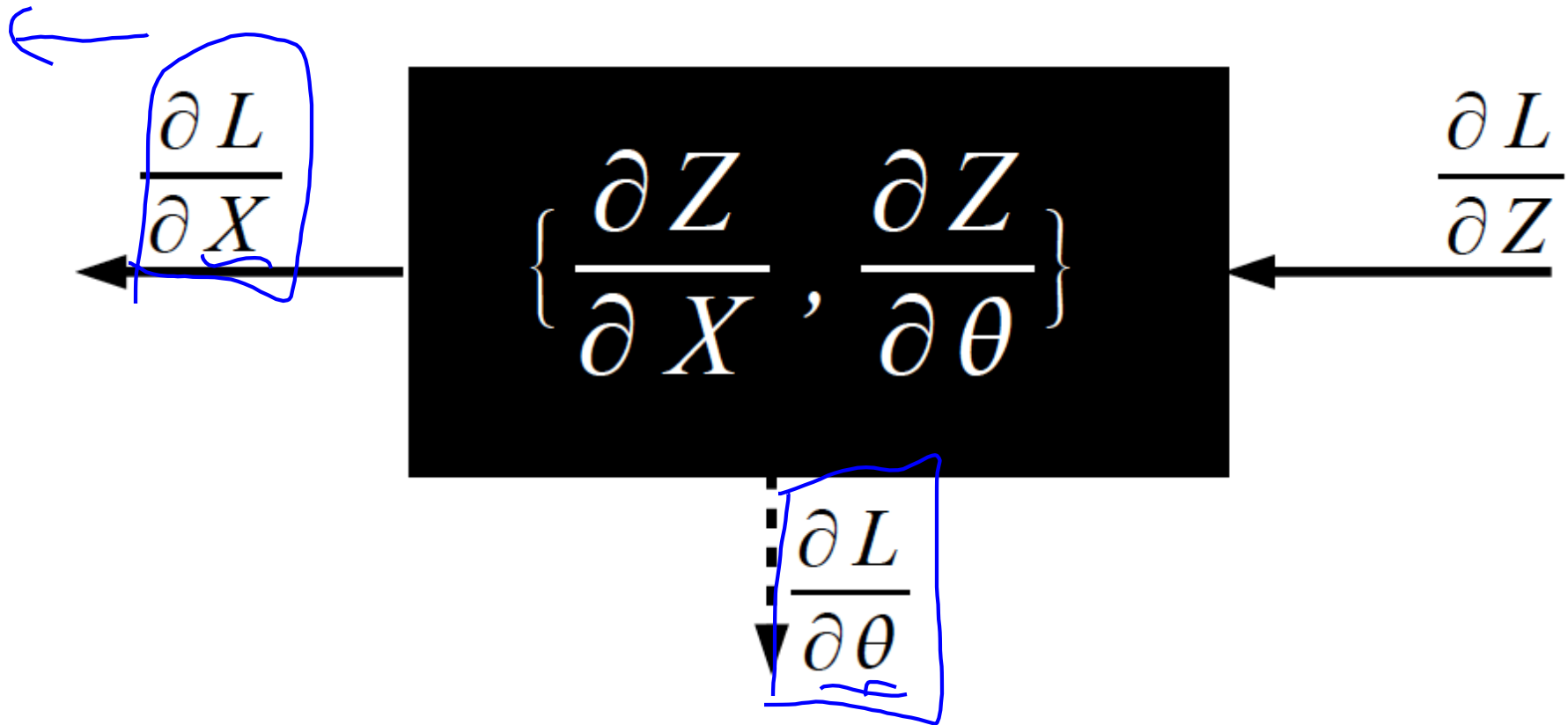  - Forward or backward?

# Computational Graph



**Any DAG of differentiable modules is allowed!**
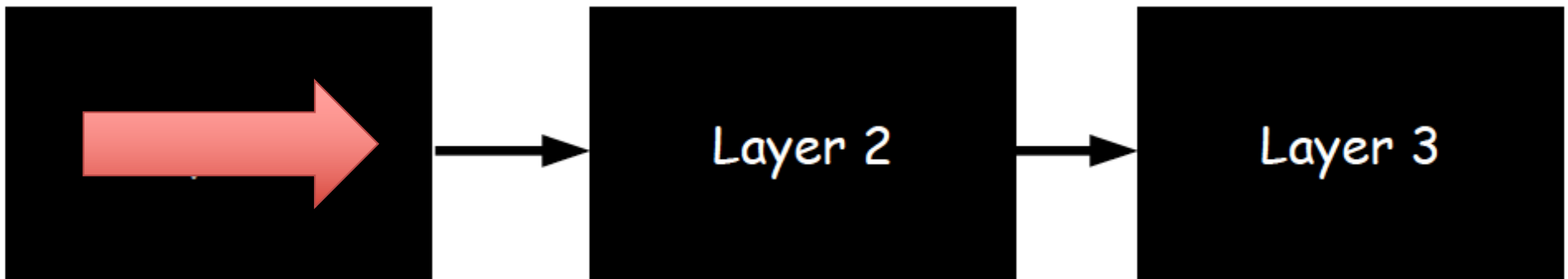
# Key Computation: Forward-Prop

# Key Computation: Back-Prop



$$\frac{\partial L}{\partial X}$$

$$\left\{ \frac{\partial Z}{\partial X}, \frac{\partial Z}{\partial \theta} \right\}$$

$$\frac{\partial L}{\partial Z}$$

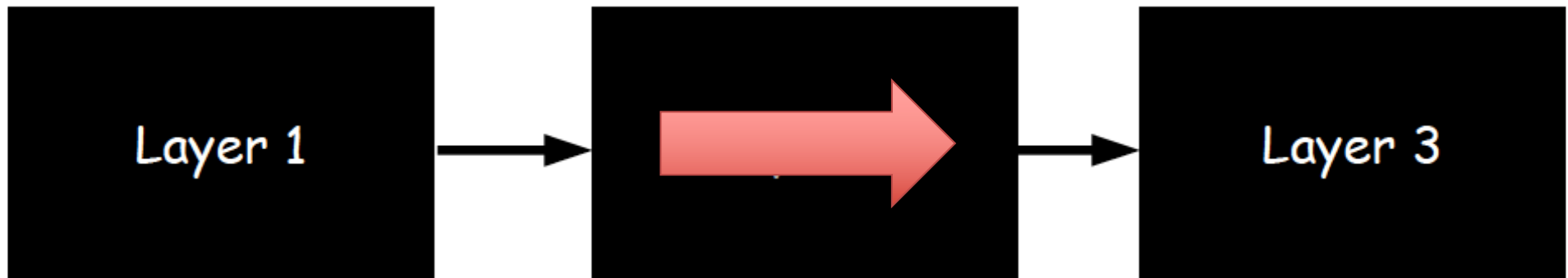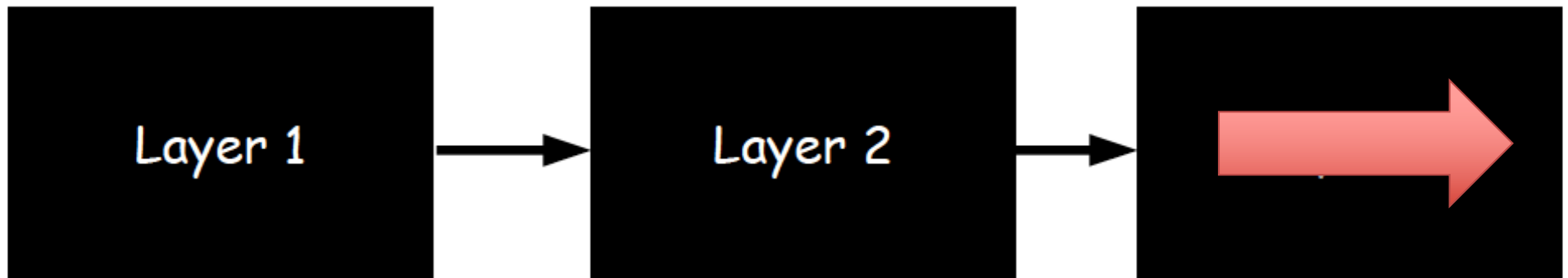$$\frac{\partial L}{\partial \theta}$$

# Neural Network Training

- Step 1: Compute Loss on mini-batch  [F-Pass]

# Neural Network Training

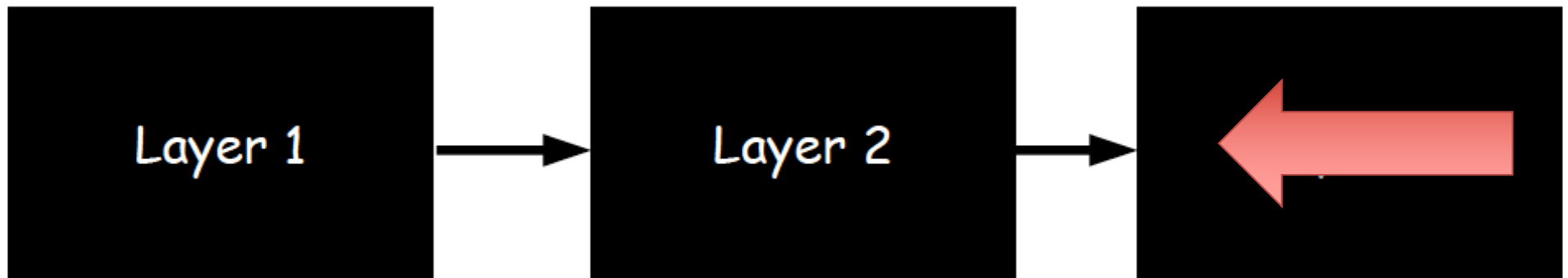- Step 1: Compute Loss on mini-batch        [F-Pass]

# Neural Network Training

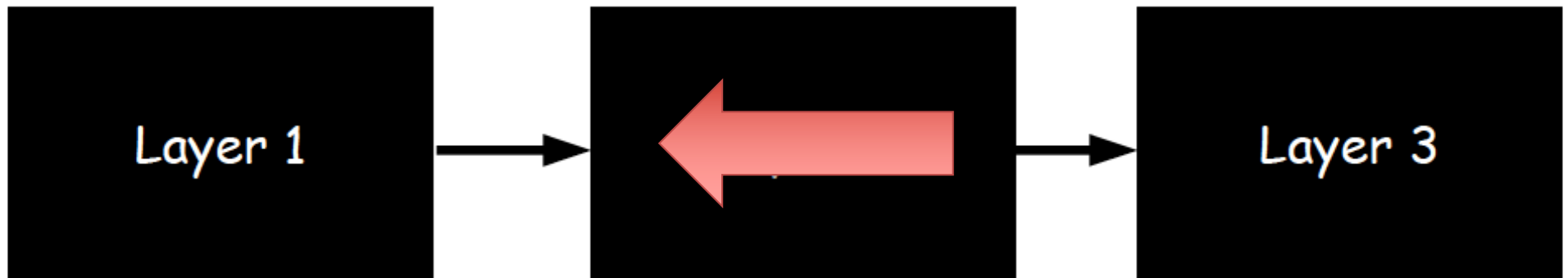- Step 1: Compute Loss on mini-batch        [F-Pass]

# Neural Network Training

- Step 1: Compute Loss on mini-batch     [F-Pass]
- Step 2: Compute gradients wrt parameters   [B-Pass]

# Neural Network Training

- Step 1: Compute Loss on mini-batch        [F-Pass]
- Step 2: Compute gradients wrt parameters   [B-Pass]
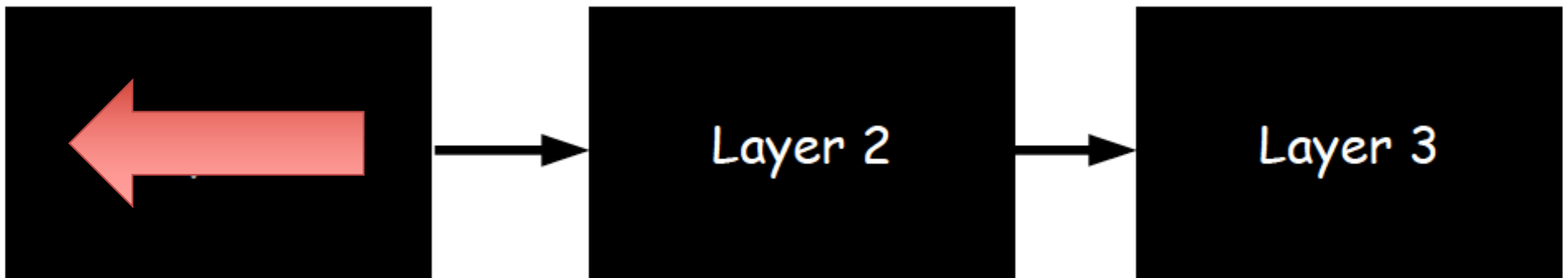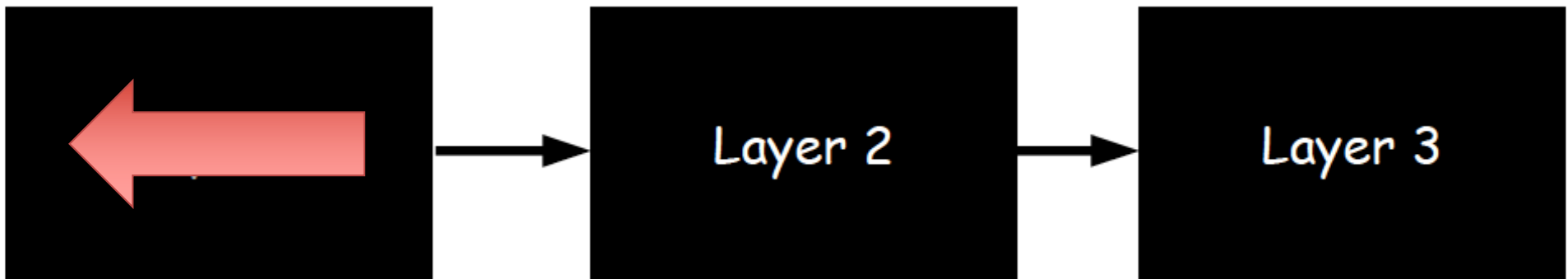
# Neural Network Training

- Step 1: Compute Loss on mini-batch    [F-Pass]
- Step 2: Compute gradients wrt parameters   [B-Pass]

# Neural Network Training

- Step 1: Compute Loss on mini-batch    [F-Pass]
- Step 2: Compute gradients wrt parameters   [B-Pass]
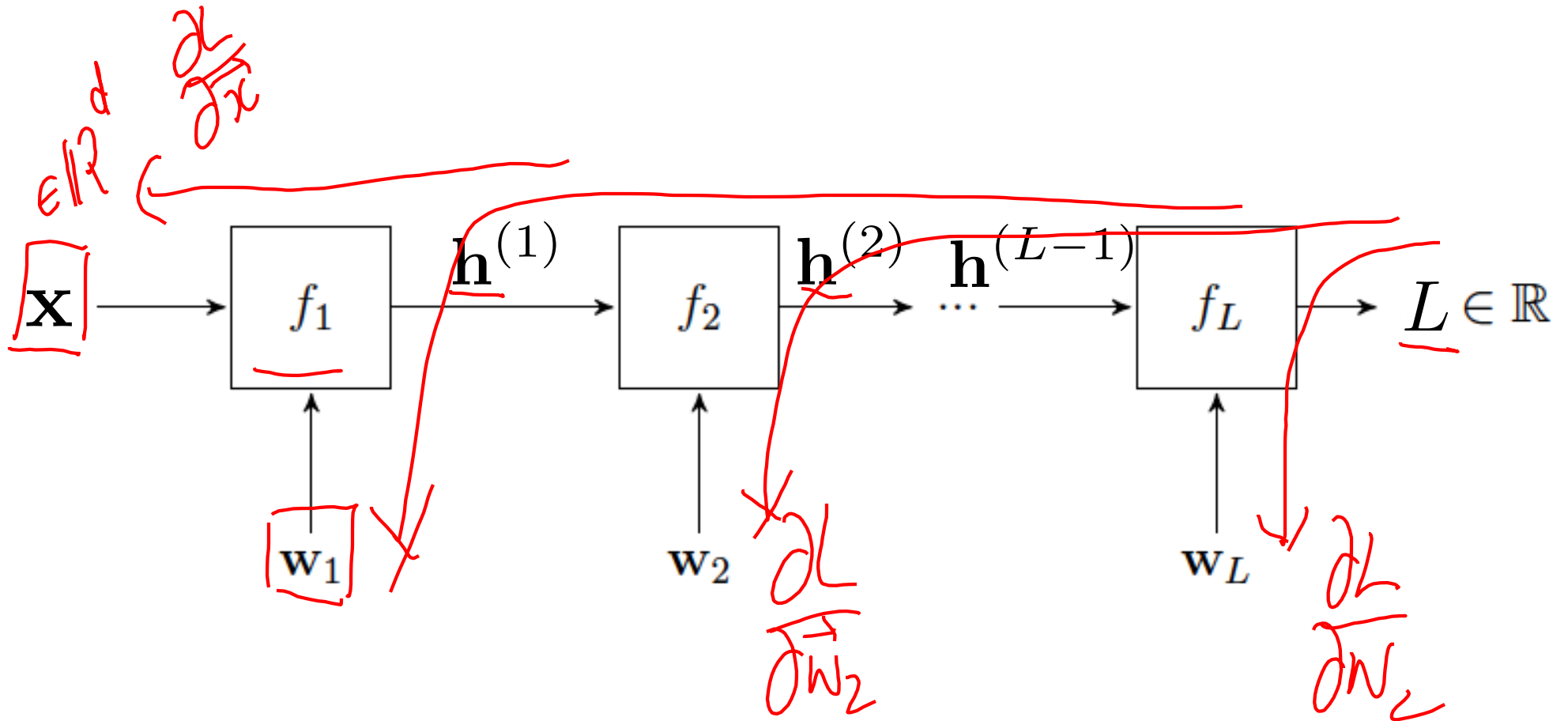- Step 3: Use gradient to update parameters

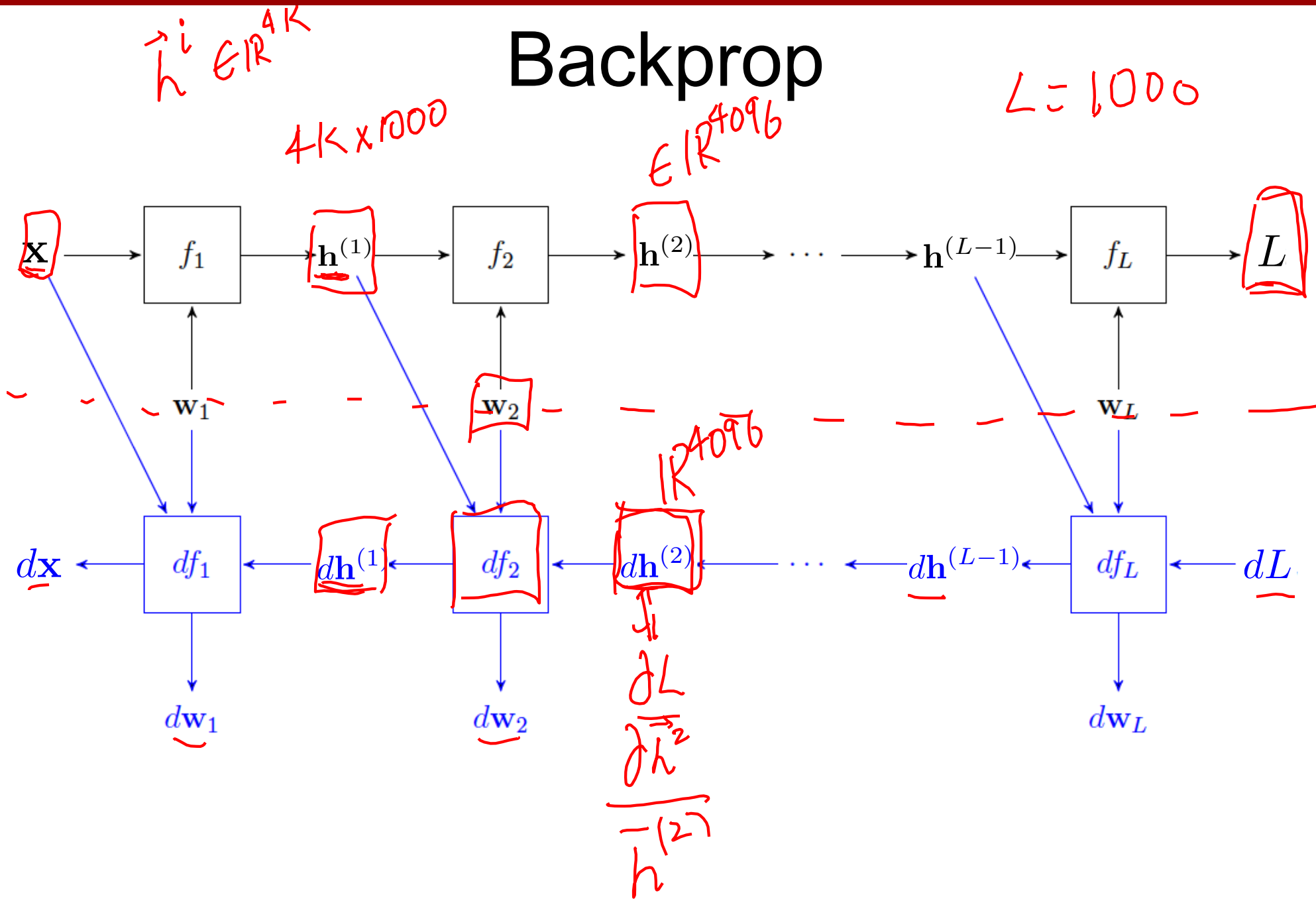

$$\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$$

# Plan for Today

- Automatic Differentiation
  - Patterns in backprop
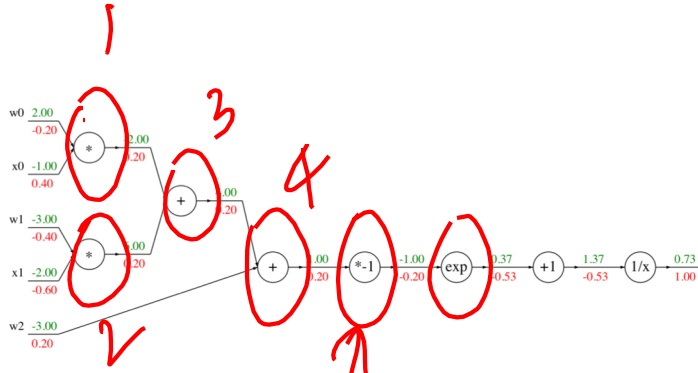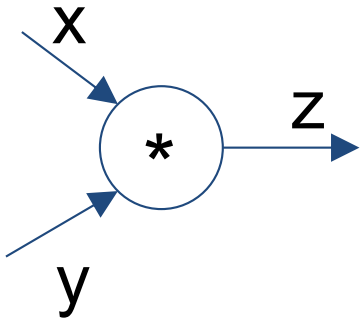  - Jacobians in FC+ReLU NNs

# Neural Network Computation Graph



$$\mathbf{x} \to \boxed{f_1} \xrightarrow{\mathbf{h}^{(1)}} \boxed{f_2} \xrightarrow{\mathbf{h}^{(2)}} \cdots \xrightarrow{\mathbf{h}^{(L-1)}} \boxed{f_L} \to L \in \mathbb{R}$$

$$\mathbf{x} \in \mathbb{R}^d \qquad \frac{\partial L}{\partial x}$$

$$\mathbf{w}_1 \qquad \mathbf{w}_2 \qquad \frac{\partial L}{\partial \vec{w}_2} \qquad \mathbf{w}_L \qquad \frac{\partial L}{\partial w_L}$$

Figure Credit: Andrea Vedaldi

# Backprop



$\vec{h}^i \in \mathbb{R}^{4K}$

$4K \times 1000$

$\in \mathbb{R}^{4096}$

$L = 1000$

$\in \mathbb{R}^{4096}$

$\dfrac{\partial L}{\partial \vec{h}^2} \quad \dfrac{\partial \vec{h}^2}{\vec{h}^{(2)}}$

# Modularized implementation: forward / backward API



Graph (or Net) object  *(rough psuedo code)*

```python
class ComputationalGraph(object):
    #...
    def forward(inputs):
        # 1. [pass inputs to input gates...]
        # 2. forward the computational graph:
        for gate in self.graph.nodes_topologically_sorted():
            gate.forward()
        return loss # the final gate in the graph outputs the loss
    def backward():
        for gate in reversed(self.graph.nodes_topologically_sorted()):
            gate.backward() # little piece of backprop (chain rule applied)
        return inputs_gradients
```

# Modularized implementation: forward / backward API



x

z

*

y

(x,y,z are scalars)

```python
class MultiplyGate(object):
    def forward(x,y):
        z = x*y
        return z
    def backward(dz):
        # dx = ... #todo
        # dy = ... #todo
        return [dx, dy]
```
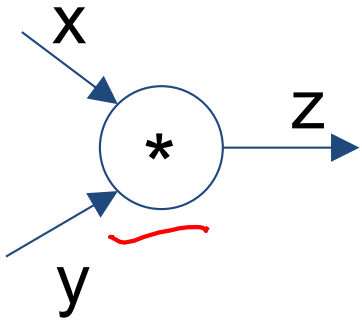
$$\frac{\partial L}{\partial z}$$

$$\frac{\partial L}{\partial x}$$

# Modularized implementation: forward / backward API

x

z

*

y

(x,y,z are scalars)

```python
class MultiplyGate(object):
    def forward(x,y):
        z = x*y
        self.x = x # must keep these around!
        self.y = y
        return z

    def backward(dz):
        dx = self.y * dz # [dz/dx * dL/dz]
        dy = self.x * dz # [dz/dy * dL/dz]
        return [dx, dy]
```

# Example: Caffe layers



Caffe is licensed under BSD 2-Clause

# Caffe Sigmoid Layer

```cpp
1   #include <cmath>
2   #include <vector>
3
4   #include "caffe/layers/sigmoid_layer.hpp"
5
6   namespace caffe {
7
8   template <typename Dtype>
9   inline Dtype sigmoid(Dtype x) {
10    return 1. / (1. + exp(-x));
11  }
12
13  template <typename Dtype>
14  void SigmoidLayer<Dtype>::Forward_cpu(const vector<Blob<Dtype>*>& bottom,
15      const vector<Blob<Dtype>*>& top) {
16    const Dtype* bottom_data = bottom[0]->cpu_data();
17    Dtype* top_data = top[0]->mutable_cpu_data();
18    const int count = bottom[0]->count();
19    for (int i = 0; i < count; ++i) {
20      top_data[i] = sigmoid(bottom_data[i]);
21    }
22  }
23
24  template <typename Dtype>
25  void SigmoidLayer<Dtype>::Backward_cpu(const vector<Blob<Dtype>*>& top,
26      const vector<bool>& propagate_down,
27      const vector<Blob<Dtype>*>& bottom) {
28    if (propagate_down[0]) {
29      const Dtype* top_data = top[0]->cpu_data();
30      const Dtype* top_diff = top[0]->cpu_diff();
31      Dtype* bottom_diff = bottom[0]->mutable_cpu_diff();
32      const int count = bottom[0]->count();
33      for (int i = 0; i < count; ++i) {
34        const Dtype sigmoid_x = top_data[i];
35        bottom_diff[i] = top_diff[i] * sigmoid_x * (1. - sigmoid_x);
36      }
37    }
38  }
39
40  #ifdef CPU_ONLY
41  STUB_GPU(SigmoidLayer);
42  #endif
43
44  INSTANTIATE_CLASS(SigmoidLayer);
45
46
47  }  // namespace caffe
```

Caffe is licensed under BSD 2-Clause

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$(1 - \sigma(x))\,\sigma(x)$$ * top_diff  (chain rule)

# Plan for Today

- Automatic Differentiation
  - Patterns in backprop
  - Jacobians in FC+ReLU NNs

# Backpropagation: a simple example

$$f(x, y, z, w) = 2 \left[ \underbrace{xy}_{w_1} + \underbrace{\text{max}\{z, w\}}_{w_2} \right]$$



x 3.00

y -4.00

$w_1$

z 2.00

$w_2$

w -1.00

$w_3$

*2

$$\bar{w}_3 = \frac{\partial f}{\partial w_3}$$

$= 2$

$$\bar{f} = \frac{\partial f}{\partial f} = 1$$

# Backpropagation: a simple example

# Patterns in backward flow

# Patterns in backward flow



$w_3 = w_1 + w_2$

Q: What is an **add** gate?

$\overline{w_1} = \dfrac{\partial f}{\partial w_1} = \dfrac{\partial f}{\partial w_3} \cdot \dfrac{\partial w_3}{\partial w_1}$

$= \overline{w_3} \left[ 1 \right]$

$= \overline{w_3}$

$\overline{w_2} = \overline{w_3}$

x   3.00

y   -4.00

*

-12.00

$w_1$

z   2.00

w   -1.00

max

2.00

$w_2$

+

-10.00
2.00

$w_3$

*2

-20.00
1.00

$\overline{f}$

# Patterns in backward flow

**add** gate: gradient distributor

# Patterns in backward flow

**add** gate: gradient distributor

Q: What is a **max** gate?

$$W_2 = \begin{cases} |z| & \text{if } z > w \\ w & \text{else} \end{cases}$$

$$\bar{z} = \frac{\partial f}{\partial z} = \left|\frac{\partial f}{\partial w_2}\right| \cdot \left|\frac{\partial w_2}{\partial z}\right|$$

$$\overline{W_2}$$

$$\begin{cases} +1 & \text{if } z > w \\ 0 & \text{else} \end{cases}$$

x  3.00

-12.00
2.00

y  -4.00

$\ast$

-10.00
2.00

$\ast 2$

-20.00
1.00

+

z  2.00   $\bar{z} = 1$

2.00
2.00

max

w  -1.00   $\bar{w} = 0$

$W_2 = \max\{z, w\}$

$\max\{0, x\}$

Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

# Patterns in backward flow

**add** gate: gradient distributor

**max** gate: gradient router

# Patterns in backward flow

$$\bar{x} = \frac{\partial f}{\partial x} = \left|\frac{\partial f}{\partial w_1}\right| \left|\frac{\partial w_1}{\partial x}\right| = y$$

$$x = \bar{w}_1 y$$

$$\bar{y} = \bar{w}_1 x$$

$$w_1 = xy$$

**add** gate: gradient distributor

**max** gate: gradient router

Q: What is a **mul** gate?



x  3.00

y  -4.00

*  -12.00
   2.00

$w_1$

z  2.00
   2.00

w  -1.00
   0.00

max  2.00
     2.00

$\bar{w}_2$

+

-10.00
2.00

*2

-20.00
1.00

# Patterns in backward flow

**add** gate: gradient distributor

**max** gate: gradient router

**mul** gate: gradient switcher

3,000

x 3.00
-8.00

y -4.00
6.00

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

z 2.00
2.00

w -1.00
0.00

* → -12.00
2.00

max → 2.00
2.00

+ → -10.00
2.00

*2 → -20.00
1.00

$h_i^{\ell} = \vec{w_i}^T \vec{x}$

$= \sum_{j=1} w_{ij} x_j$

$\vec{x} \rightarrow \boxed{\phantom{}} \rightarrow h^{\ell} \cdots \boxed{\phantom{}} \boxed{\phantom{}} \xrightarrow{h} \boxed{\phantom{}} \rightarrow$

$\frac{\partial L}{\partial h} \ll c$

# Gradients add at branches

# Duality in Fprop and Bprop

# Plan for Today

- Automatic Differentiation
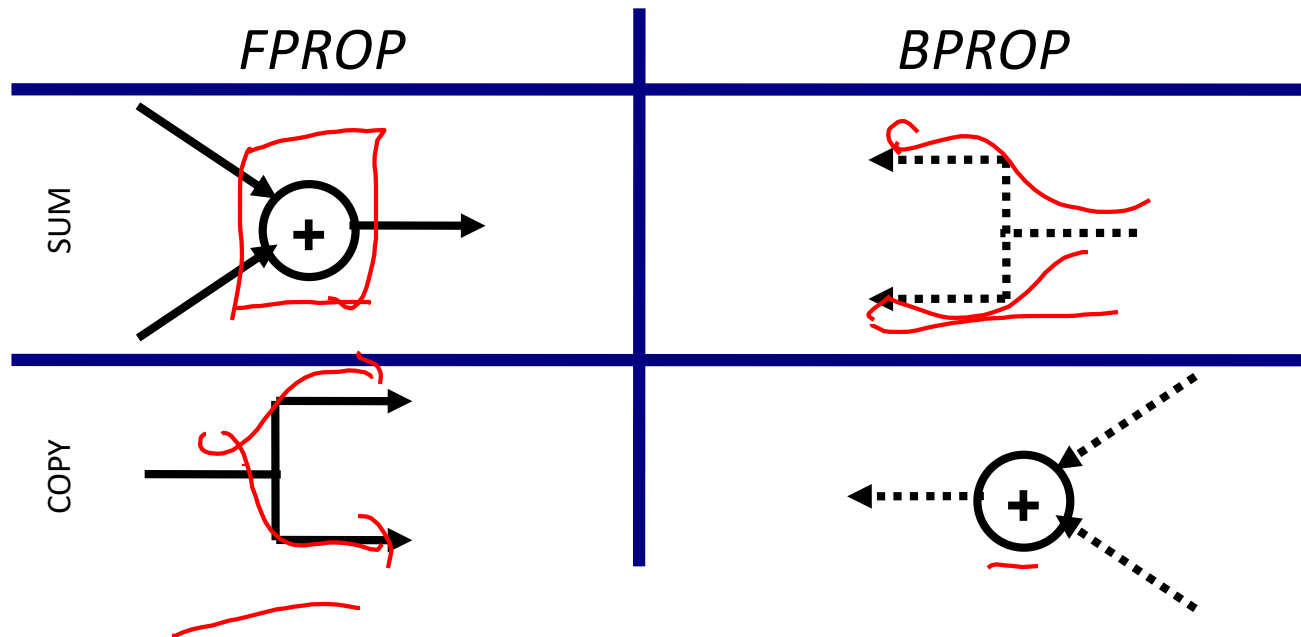    - Patterns in backprop
    - Jacobians in FC+ReLU NNs

# Backprop

# Jacobian of ReLU $\max\{0, x\}$

*layer $\ell$*

$\vec{h}^{\ell-1} \in \mathbb{R}^{4096}$

$\vec{h}^{\ell} \in \mathbb{R}^{4096}$

$\vec{h}^{\ell} = \max(0, \vec{h}^{\ell-1})$

4096-d
input vector

g(x) = max(0,x)
*(elementwise)*

4096-d
output vector

$$\left[ \frac{\partial \vec{h}^{\ell}}{\partial \vec{h}^{\ell-1}} \right] \quad 4096 \times 4096$$

# Jacobian of ReLU

4096-d
input vector

$g(x) = \max(0,x)$
*(elementwise)*

4096-d
output vector

Q: what is the
size of the
Jacobian matrix?

# Jacobian of ReLU

4096-d
input vector
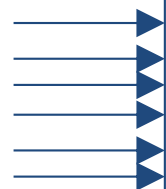
$g(x) = \max(0,x)$
*(elementwise)*

4096-d
output vector

Q: what is the size of the Jacobian matrix? [4096 x 4096!]

# Jacobian of ReLU



4096-d input vector $\rightarrow$ g(x) = max(0,x) *(elementwise)* $\rightarrow$ 4096-d output vector

Q: what is the size of the Jacobian matrix? [4096 x 4096!]

Q2: what does it look like?

$\dfrac{\partial \vec{h}^{\ell}}{\partial \vec{h}^{\ell-1}}$

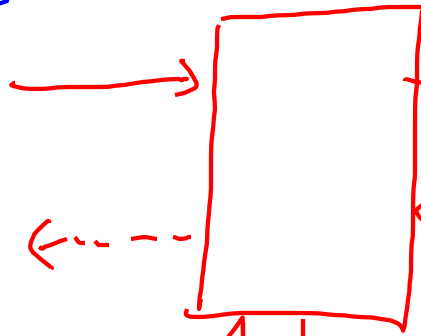$\dfrac{\partial L}{\partial h^{\ell}}$

$\dfrac{\partial h_i^{\ell}}{\partial h_j^{\ell-1}}$

0

$J = \begin{bmatrix} & 0 \\ 0 & \end{bmatrix}$

# Jacobians of FC-Layer

$$\vec{W}_i^{(t+1)} = \vec{w}_i^{(t)} - \eta \frac{\partial L}{\partial \vec{w}_i}$$

$$\vec{h}^{\ell-1} \in \mathbb{R}^{C_1}$$

$$\vec{h}^{\ell} \in \mathbb{R}^{C_2}$$

$$\frac{\partial L}{\partial \vec{h}^{\ell-1}} \Big|_{\times C_1}$$

$$\frac{\partial L}{\partial \vec{h}^{\ell}} \Big|_{1 \times C_2}$$

$$\underline{W} \quad \frac{\partial L}{\partial W} \Big|_{C_2 \times C_1}$$

$$\text{If } \vec{W}_i^{(0)} = \vec{W}_j^{(0)}$$

$$\frac{\partial L}{\partial \vec{h}^{\ell}} \qquad h_i^{\ell} = \vec{w}_i^{T} \vec{h}^{\ell-1}$$
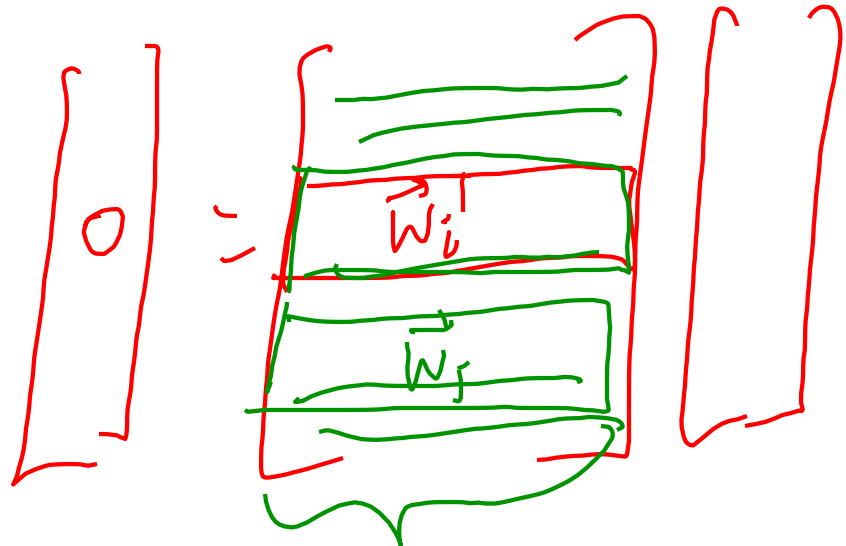
$$h_j = \vec{w}_j^{T} \vec{h}^{\ell-1}$$
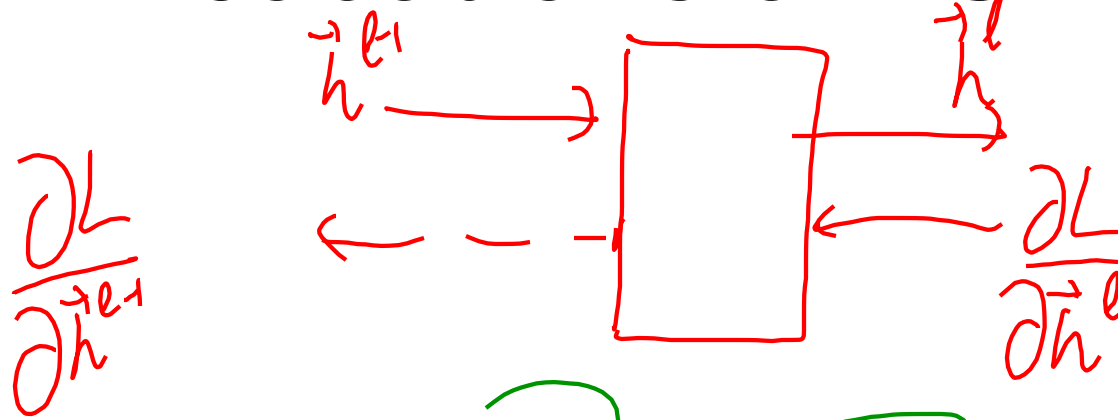
$$\vec{h}^{\ell} = W \vec{h}^{\ell-1}$$

$$\frac{\partial L}{\partial \vec{w}_i} = \frac{\partial L}{\partial h_i} \frac{\partial h_i}{\partial \vec{w}_i}$$

Scalar input $\vec{h}^{\ell-1}$

$$h_i^{\ell} = \vec{w}_i^{T} \vec{h}^{\ell-1}$$

$$0 = \frac{\vec{w}_i}{\vec{w}_j}$$

# Jacobians of FC-Layer



$$\vec{h}^{\ell-1} \longrightarrow \boxed{\phantom{xxx}} \longrightarrow \vec{h}^{\ell}$$

$$\frac{\partial L}{\partial \vec{h}^{\ell-1}} \longleftarrow \qquad \longleftarrow \frac{\partial L}{\partial \vec{h}^{\ell}}$$
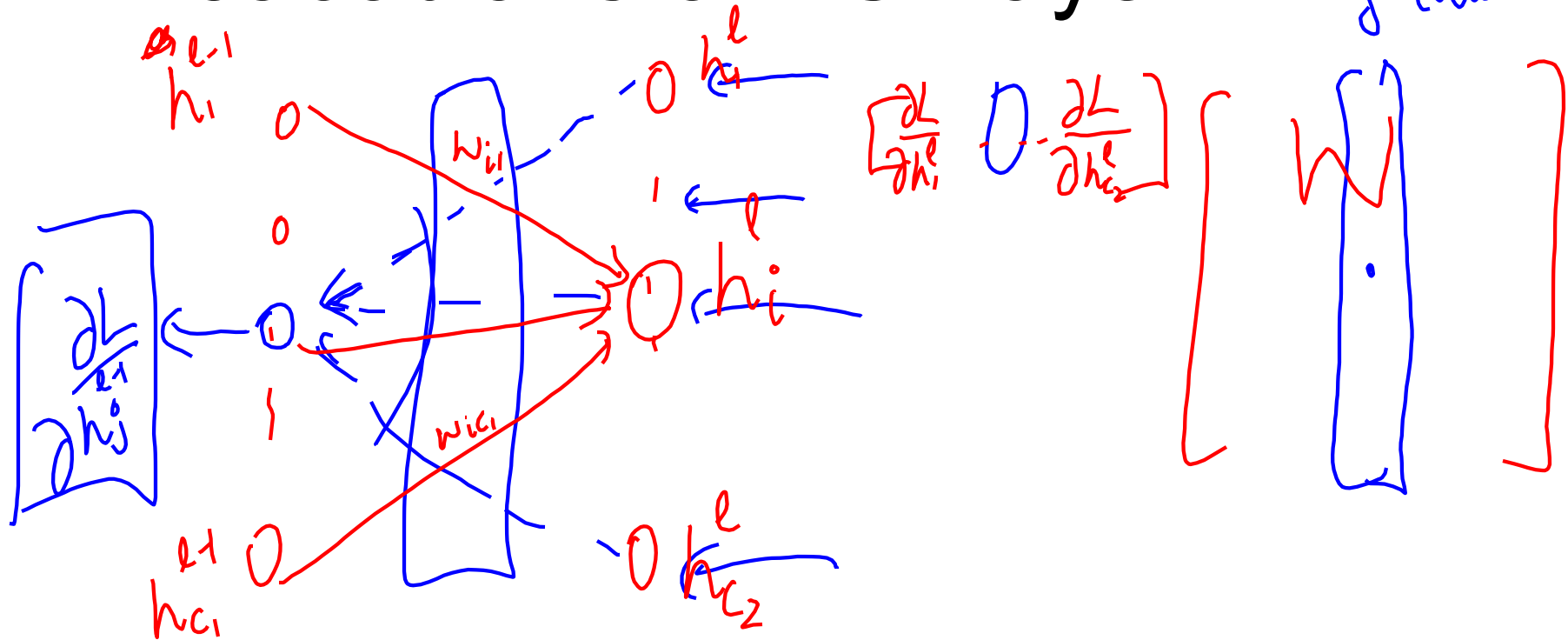
$$\vec{h}^{(\ell)} = W \vec{h}^{\ell-1}$$

$$\frac{\partial \vec{h}^{\ell}}{\partial \vec{h}^{\ell-1}} = W$$

$$\frac{\partial L}{\partial \vec{h}^{\ell-1}} = \left[ \frac{\partial L}{\partial \vec{h}^{\ell}} \right] \left[ \frac{\partial \vec{h}^{\ell}}{\partial \vec{h}^{\ell-1}} \right]$$

$$\text{input}$$

$$= [\text{input}] \cdot W$$

# Jacobians of FC-Layer

# Jacobians of FC-Layer