

# Lecture 6: Backpropagation

## Waitlist update

Over the past few days everyone on the waitlist got an override!

**Reminder: A2**

**Due Monday, 9/30, 11:59pm**

Justin Johnson

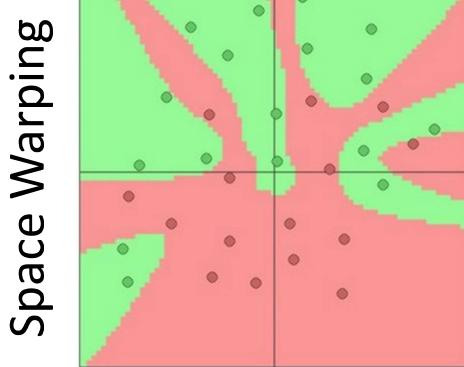
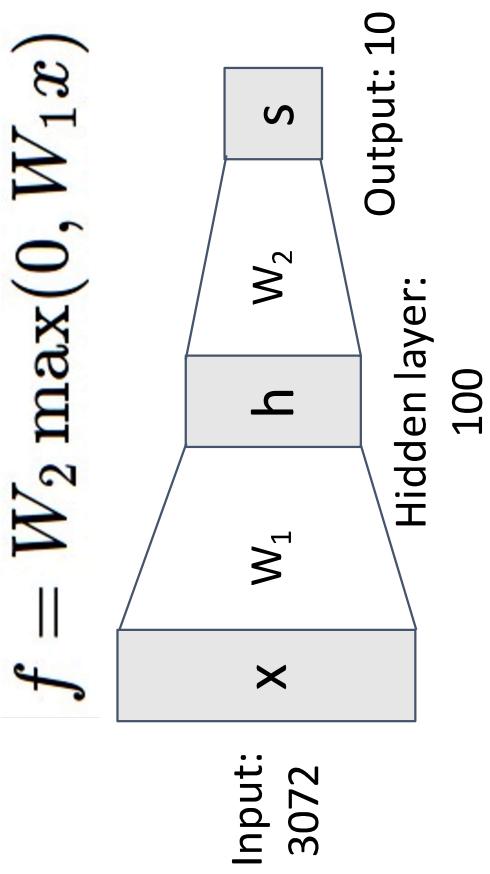
Lecture 6 - 3

September 23, 2019

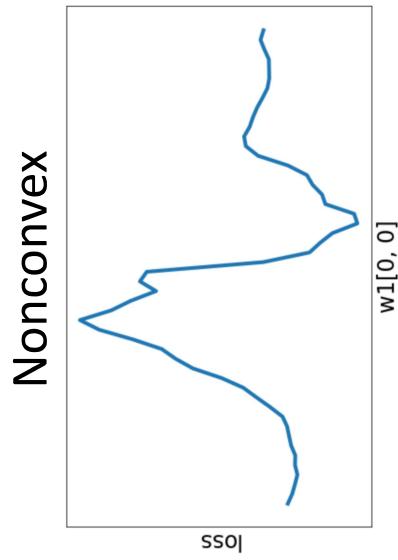
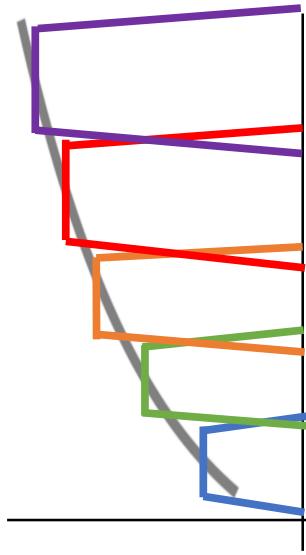
# Last time: Neural Networks

From linear classifiers to  
fully-connected networks

$$f = W_2 \max(0, W_1 x)$$



Universal Approximation



Nonconvex

## Problem: How to compute gradients?

$$s = f(x; W_1, W_2) = W_2 \max(0, W_1 x) \quad \text{Nonlinear score function}$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \quad \text{SVM Loss on predictions}$$

$$R(W) = \sum_k W_k^2 \quad \text{Regularization}$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i + \lambda R(W_1) + \lambda R(W_2) \quad \text{Total loss: data loss + regularization}$$

If we can compute  $\frac{\partial L}{\partial W_1}, \frac{\partial L}{\partial W_2}$  n we can learn  $W_1$  and  $W_2$

## (Bad) Idea: Derive $\nabla_W L$ on paper

$$s = f(x; W) = Wx$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$= \sum_{j \neq y_i} \max(0, W_{j,:} \cdot x + W_{y_i,:} \cdot x + 1)$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i + \lambda \sum_k W_k^2$$

$$= \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, W_{j,:} \cdot x + W_{y_i,:} \cdot x + 1) + \lambda \sum_k W_k^2$$

$$\nabla_W L = \nabla_W \left( \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, W_{j,:} \cdot x + W_{y_i,:} \cdot x + 1) + \lambda \sum_k W_k^2 \right)$$

**Problem:** Very tedious: Lots of matrix calculus, need lots of paper

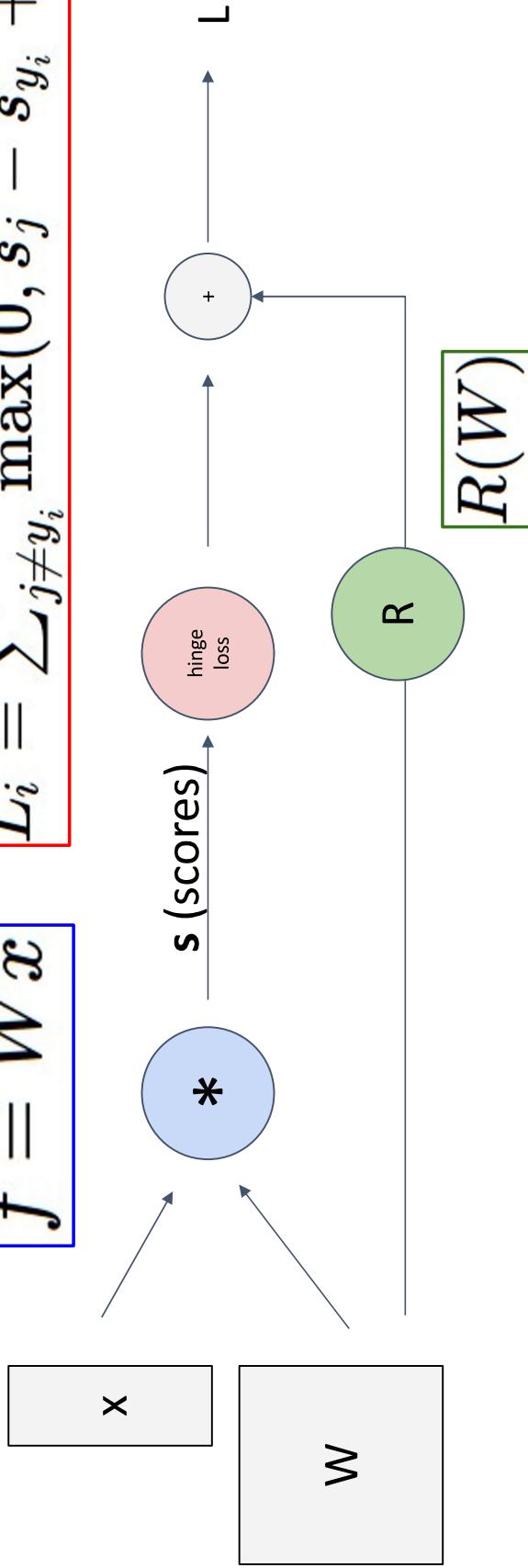
**Problem:** What if we want to change loss? E.g. use softmax instead of SVM? Need to re-derive from scratch. Not modular!

**Problem:** Not feasible for very complex models!

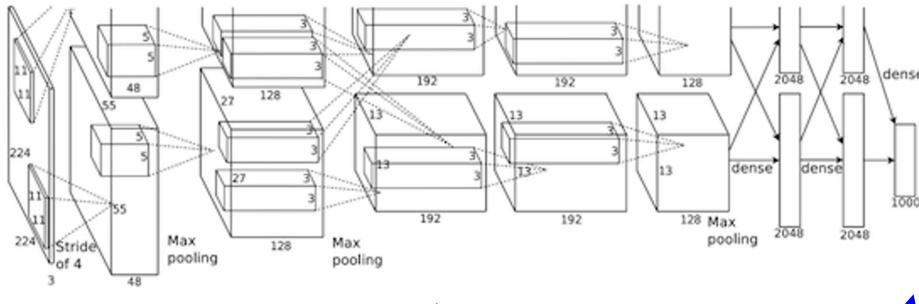
# Better Idea: Computational Graphs

$$f = Wx$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$



# Deep Network (AlexNet)



input image

weights

loss

# Neural Turing Machine

input image

loss

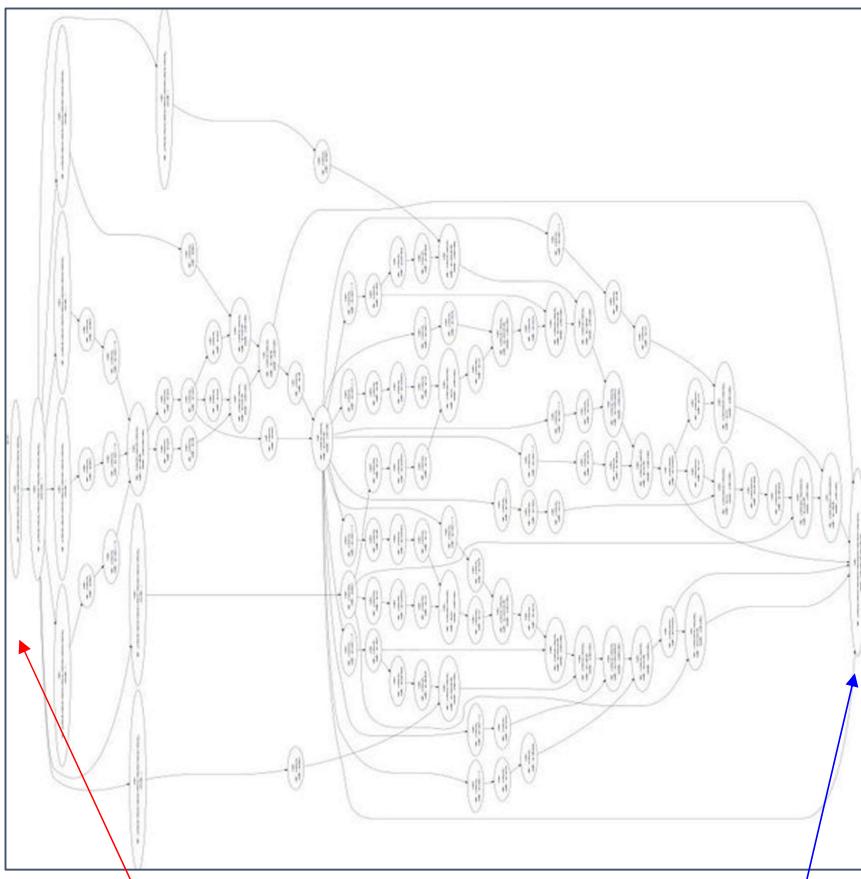


Figure reproduced with permission from a [Twitter post](#) by Andrej Karpathy.

# Neural Turing Machine



Graves et al, arXiv 2014

Justin Johnson

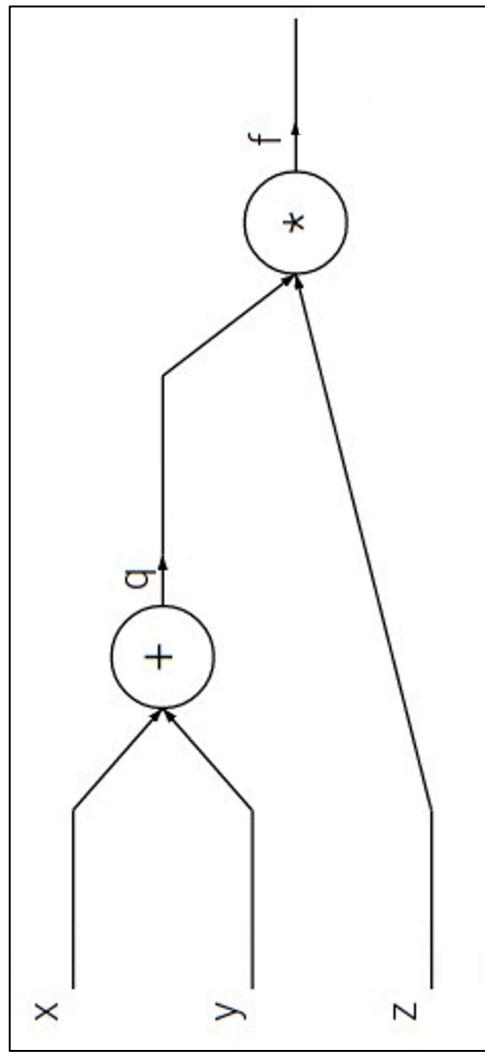
Lecture 6 - 10

September 23, 2019

Figure reproduced with permission from a [Twitter post](#) by Andrej Karpathy.

## Backpropagation: Simple Example

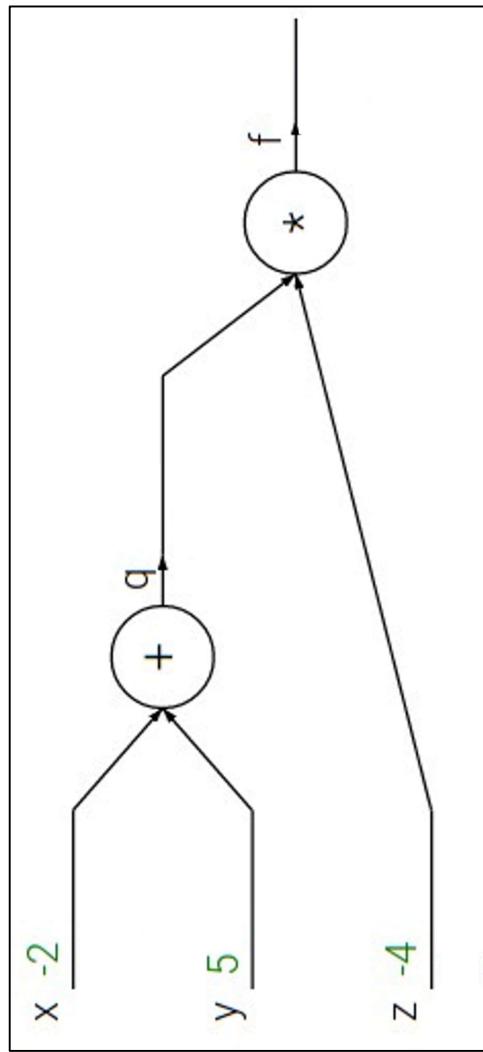
$$f(x, y, z) = (x + y)z$$



## Backpropagation: Simple Example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$



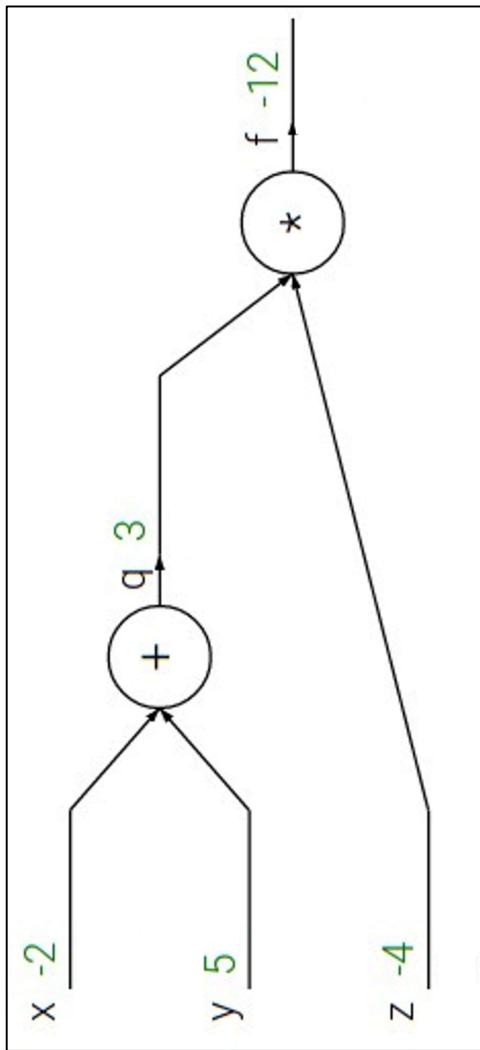
## Backpropagation: Simple Example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

**1. Forward pass:** Compute outputs

$$q = x + y \quad f = qz$$



## Backpropagation: Simple Example

$$f(x, y, z) = (x + y)z$$

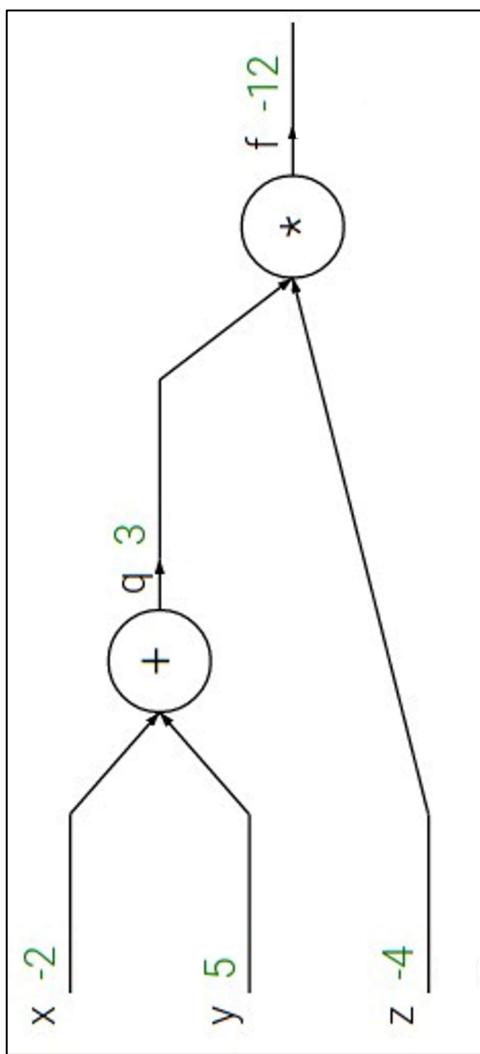
e.g.  $x = -2, y = 5, z = -4$

**1. Forward pass:** Compute outputs

$$q = x + y \quad f = qz$$

**2. Backward pass:** Compute derivatives

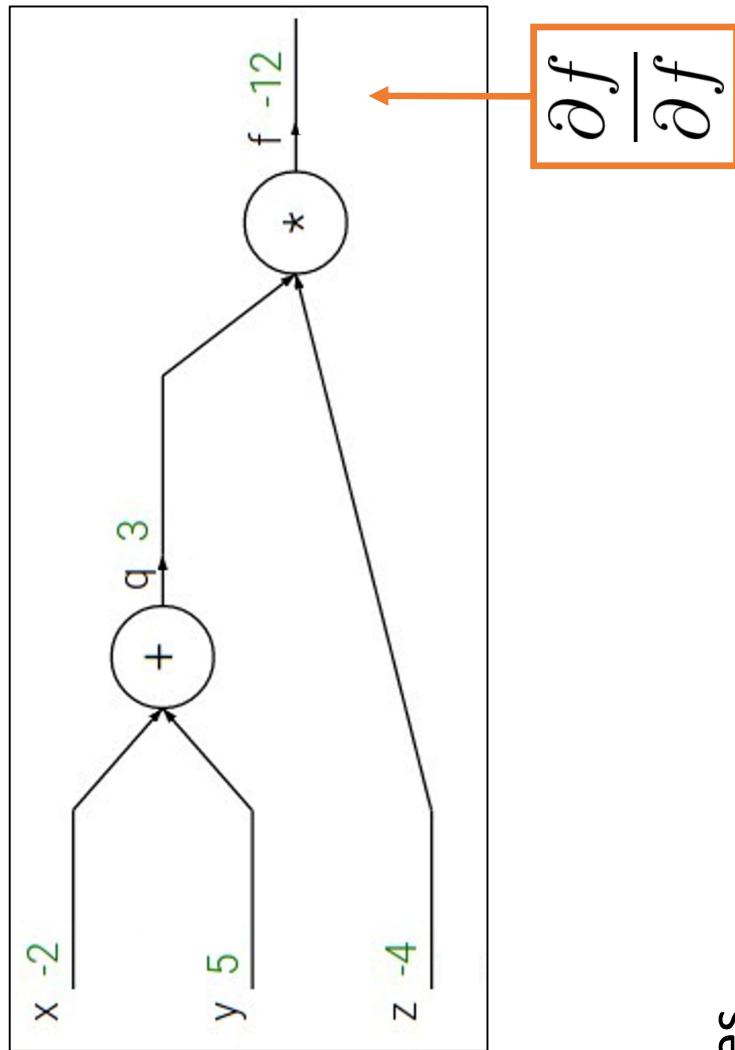
Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



## Backpropagation: Simple Example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$



**1. Forward pass:** Compute outputs

$$q = x + y \quad f = qz$$

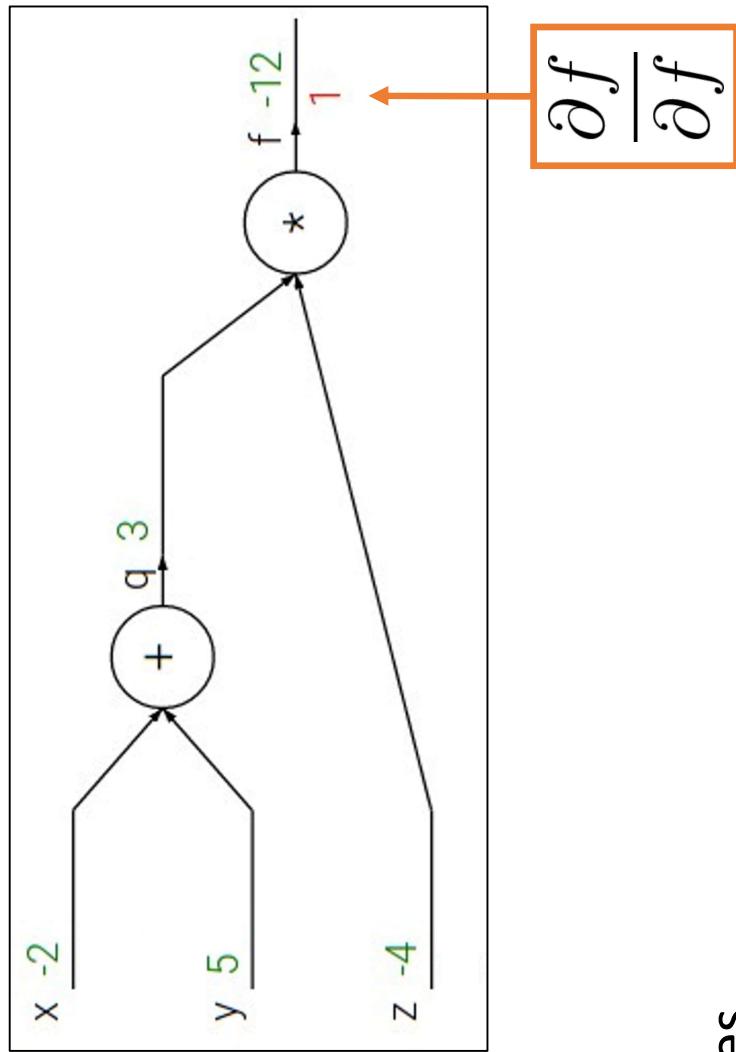
**2. Backward pass:** Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

## Backpropagation: Simple Example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$



**1. Forward pass:** Compute outputs

$$q = x + y \quad f = qz$$

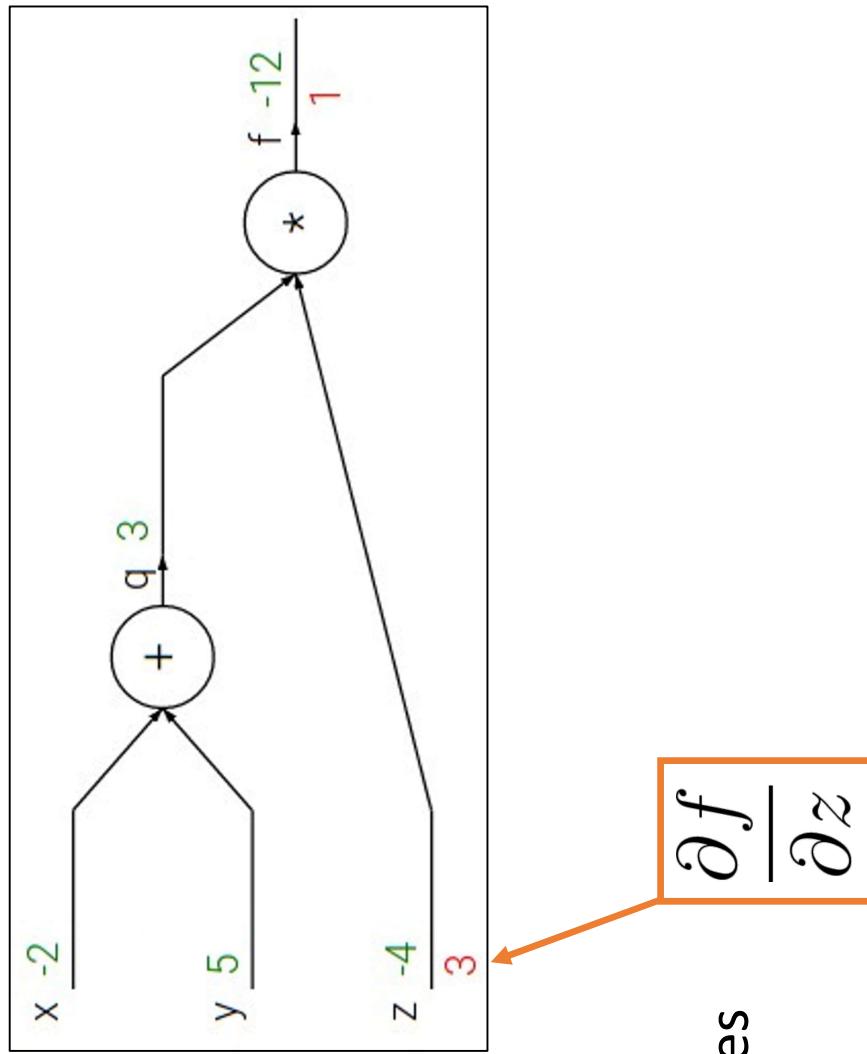
**2. Backward pass:** Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

## Backpropagation: Simple Example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$



**1. Forward pass:** Compute outputs

$$f = qz$$

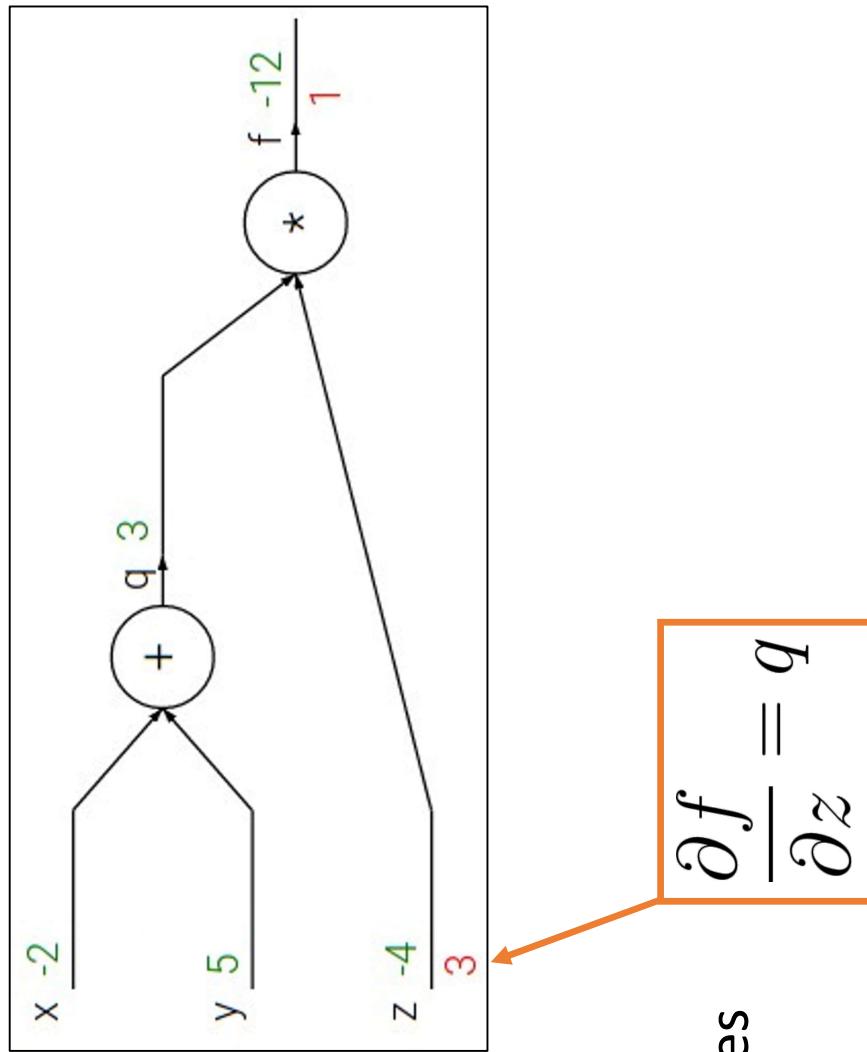
**2. Backward pass:** Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

## Backpropagation: Simple Example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$



**1. Forward pass:** Compute outputs

$$f = qz$$

**2. Backward pass:** Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

## Backpropagation: Simple Example

$$f(x, y, z) = (x + y)z$$

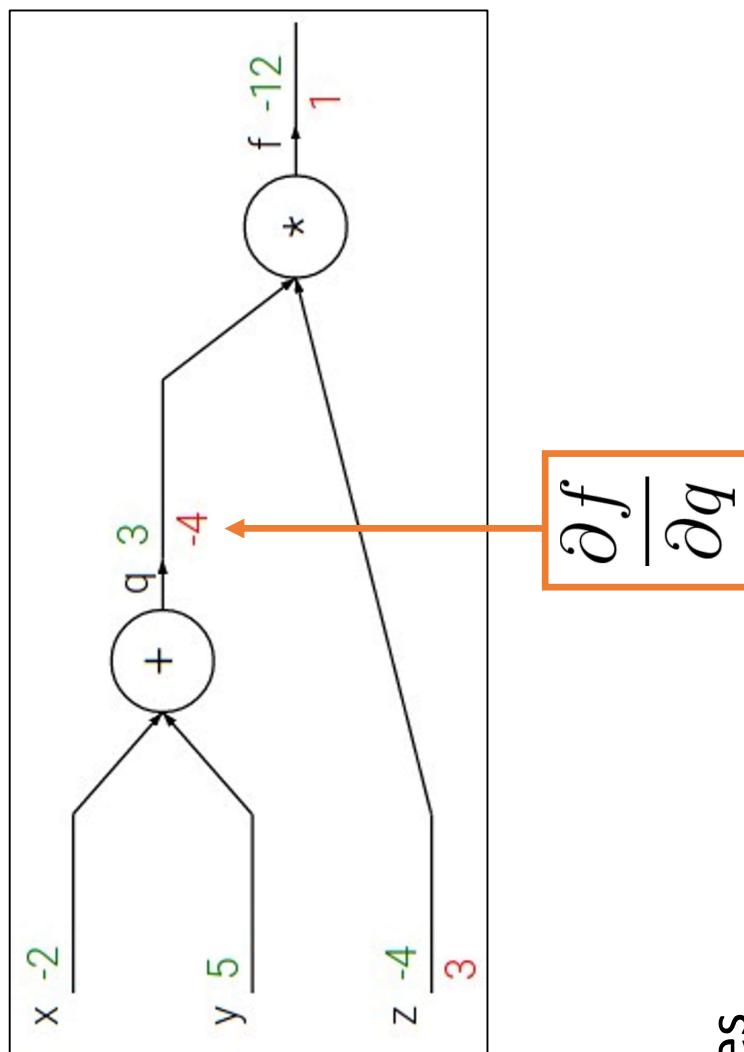
e.g.  $x = -2, y = 5, z = -4$

**1. Forward pass:** Compute outputs

$$q = x + y \quad f = qz$$

**2. Backward pass:** Compute derivatives

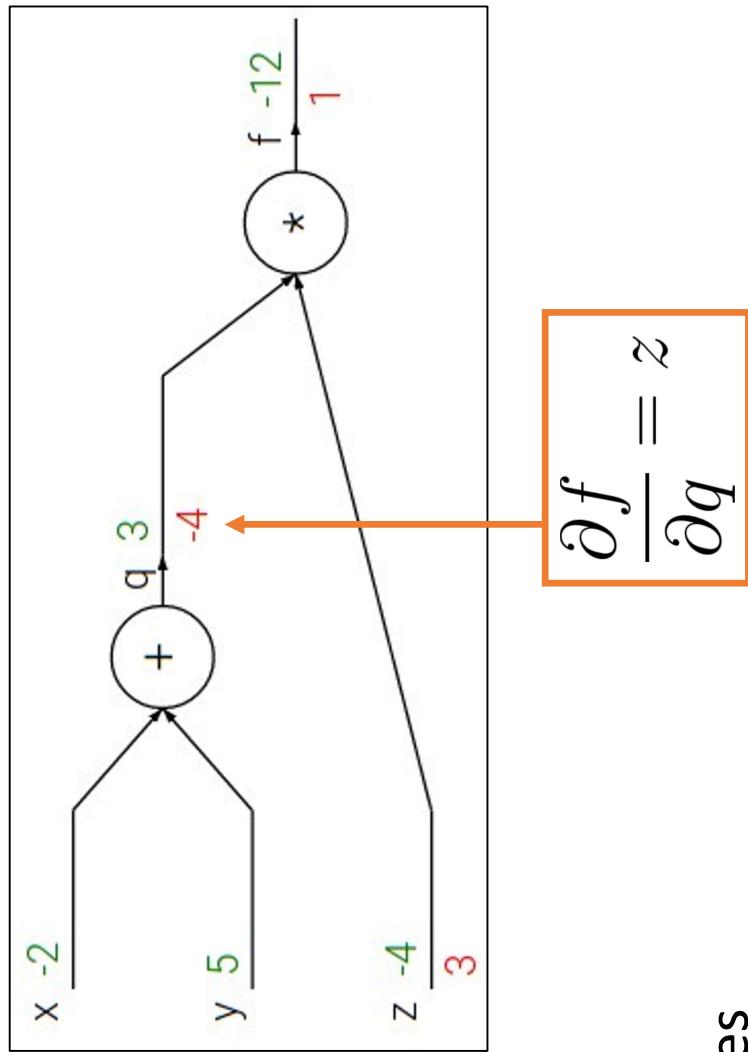
Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



## Backpropagation: Simple Example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$



**1. Forward pass:** Compute outputs

$$f = qz$$

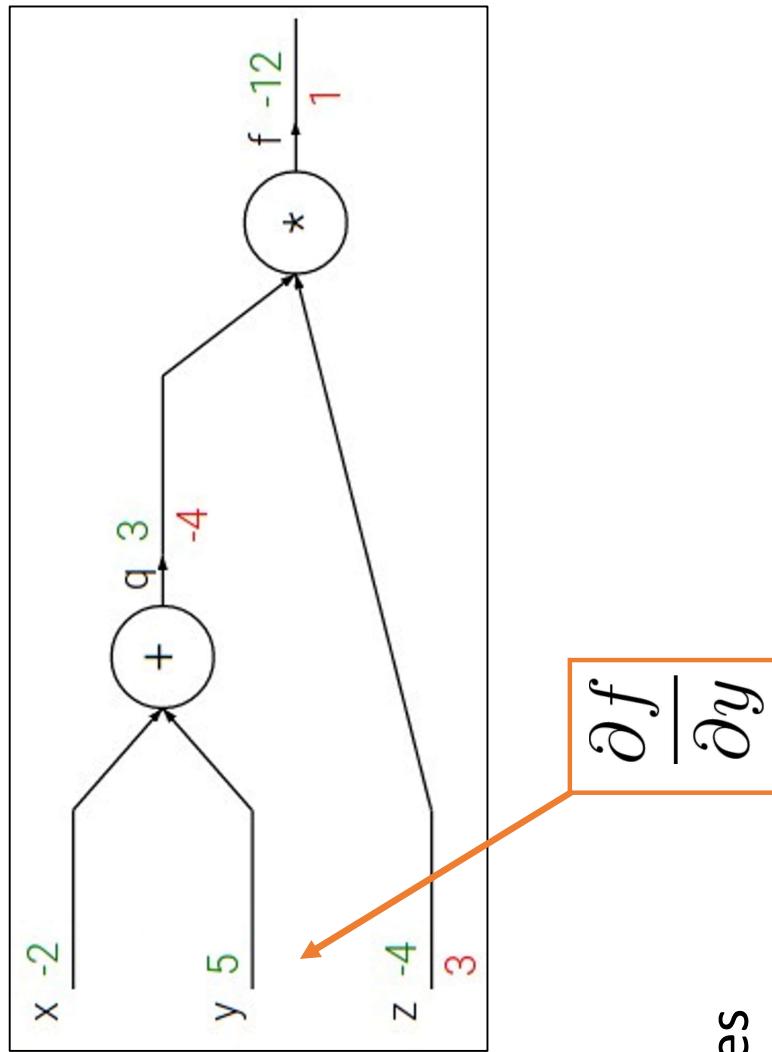
**2. Backward pass:** Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

## Backpropagation: Simple Example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$



**1. Forward pass:** Compute outputs

$$q = x + y \quad f = qz$$

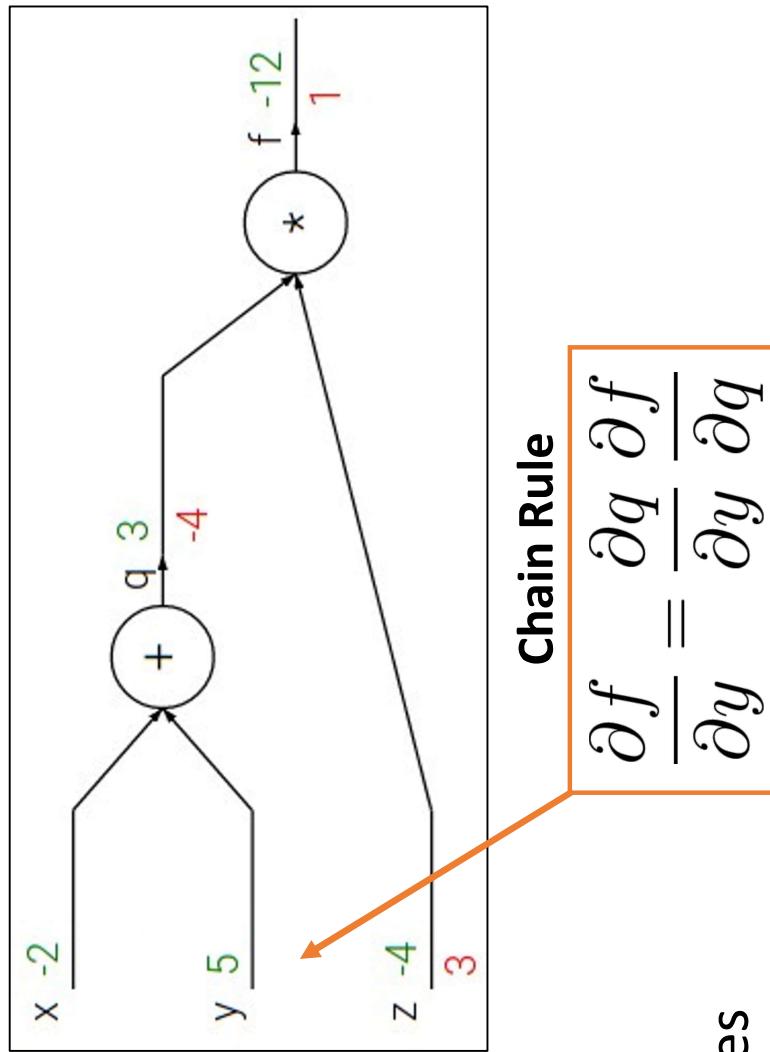
**2. Backward pass:** Compute derivatives

$$\text{Want: } \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$

## Backpropagation: Simple Example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$



**1. Forward pass:** Compute outputs

$$q = x + y \quad f = qz$$

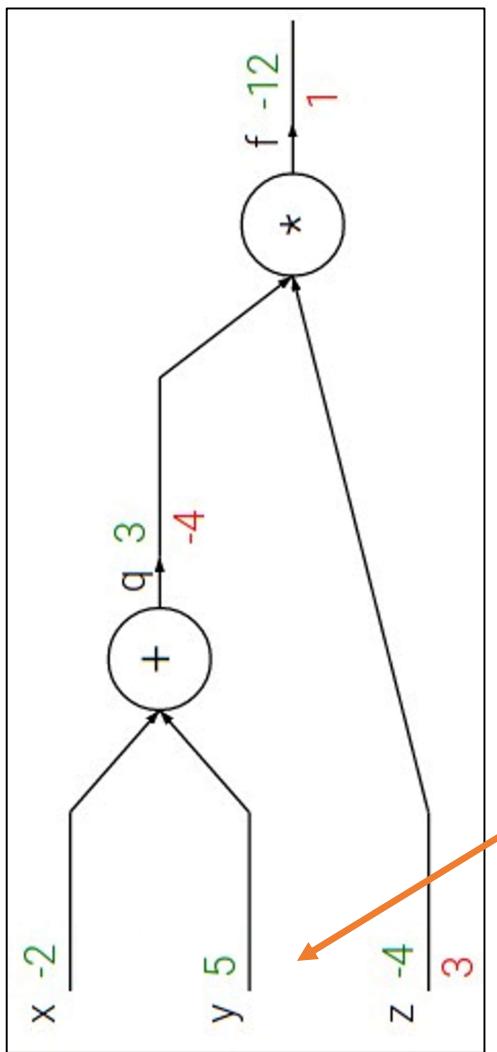
**2. Backward pass:** Compute derivatives

$$\text{Want: } \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$

## Backpropagation: Simple Example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$



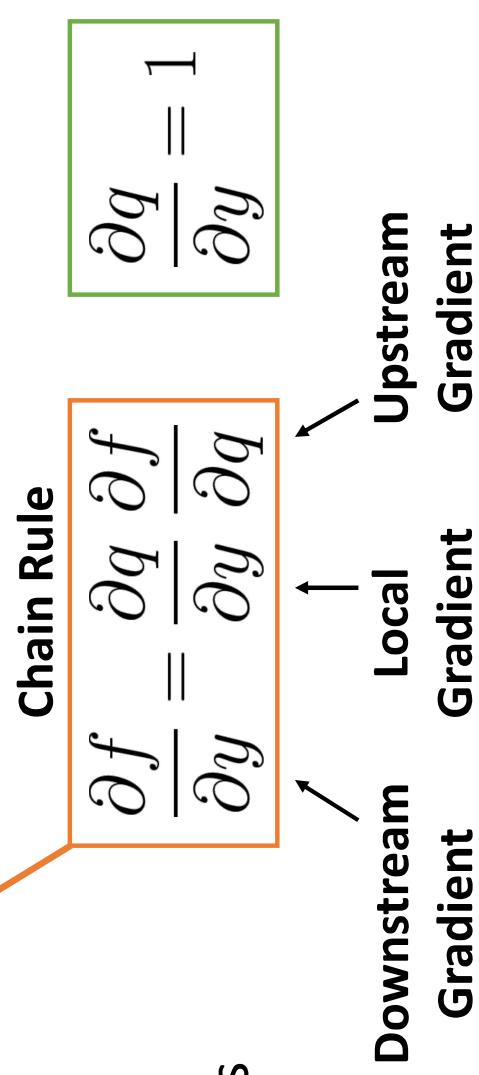
**1. Forward pass:** Compute outputs

$$q = x + y$$

$$f = qz$$

**2. Backward pass:** Compute derivatives

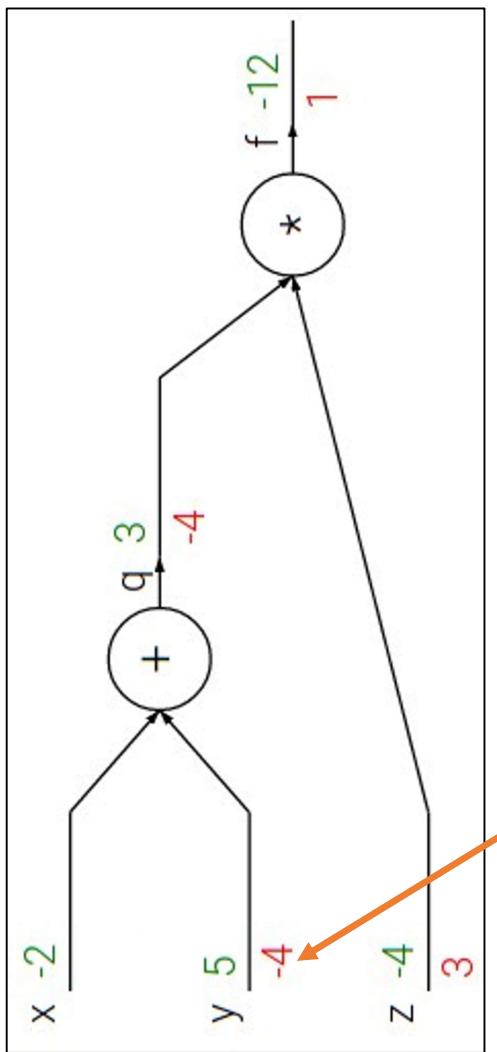
$$\text{Want: } \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$



## Backpropagation: Simple Example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$



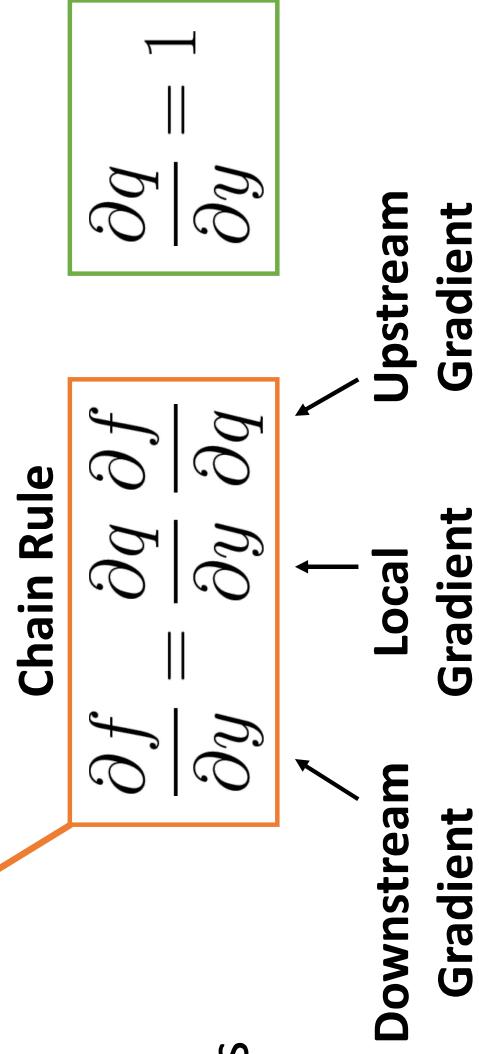
**1. Forward pass:** Compute outputs

$$q = x + y$$

$$f = qz$$

**2. Backward pass:** Compute derivatives

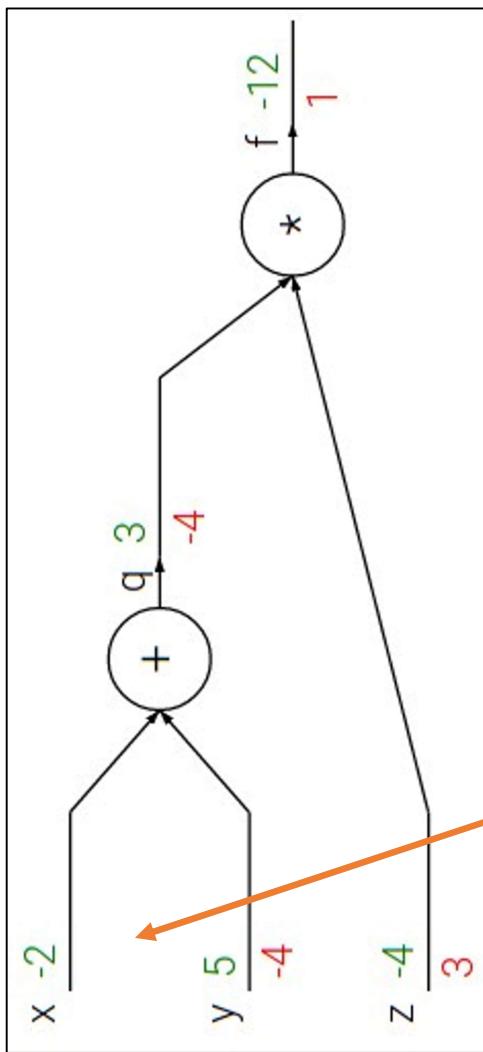
$$\text{Want: } \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$



## Backpropagation: Simple Example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$



**1. Forward pass:** Compute outputs

$$q = x + y$$

$$f = qz$$

**2. Backward pass:** Compute derivatives

$$\text{Want: } \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$

$$\frac{\partial q}{\partial x} = 1$$

$$\frac{\partial q}{\partial y} = 1$$

$$\frac{\partial q}{\partial z} = 1$$

Downstream  
Gradient

Local  
Gradient

Upstream  
Gradient

# Backpropagation: Simple Example

$$f(x, y, z) = (x + y)z$$

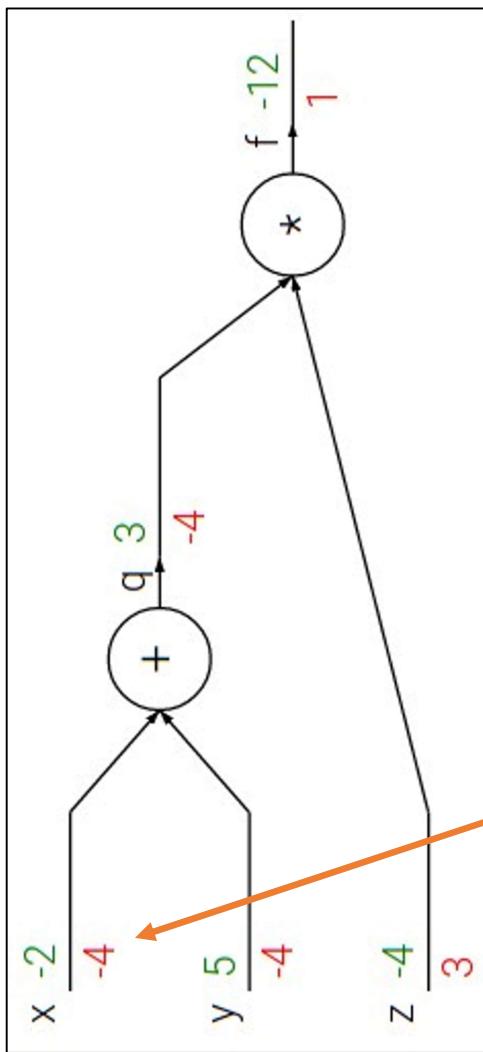
e.g.  $x = -2, y = 5, z = -4$

## 1. Forward pass: Compute outputs

$$q = x + y$$

## 2. Backward pass: Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

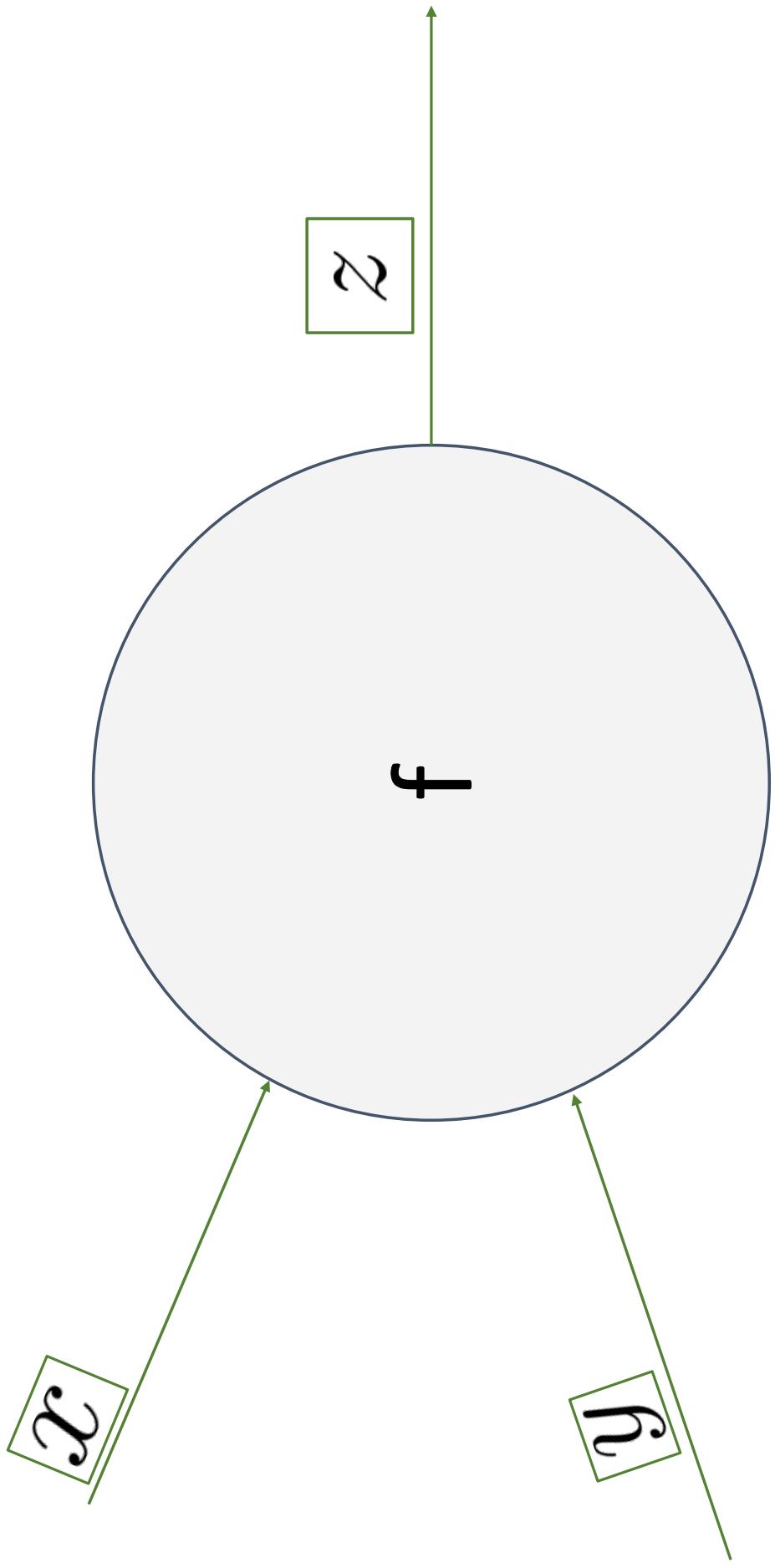


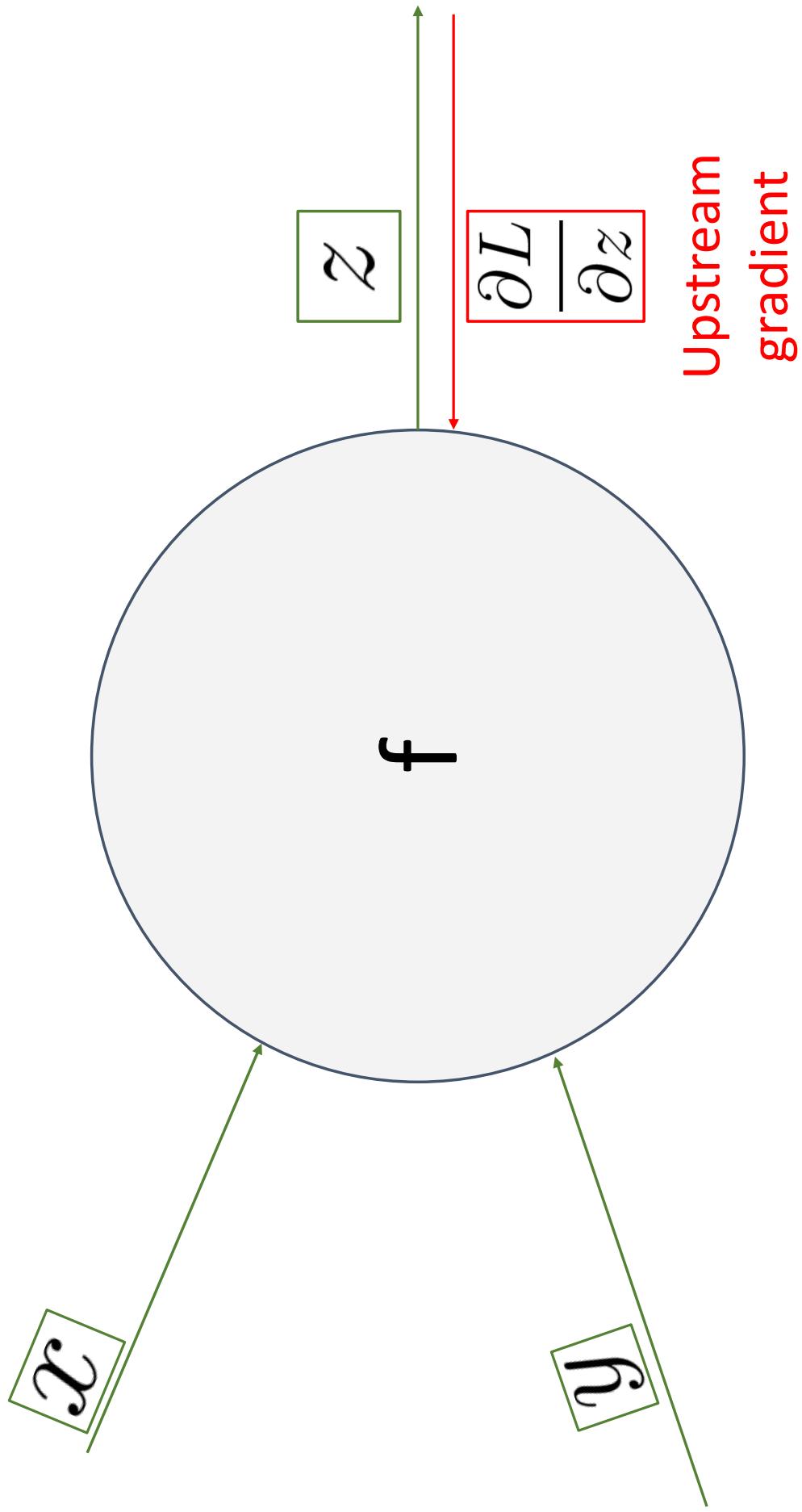
## Chain Rule

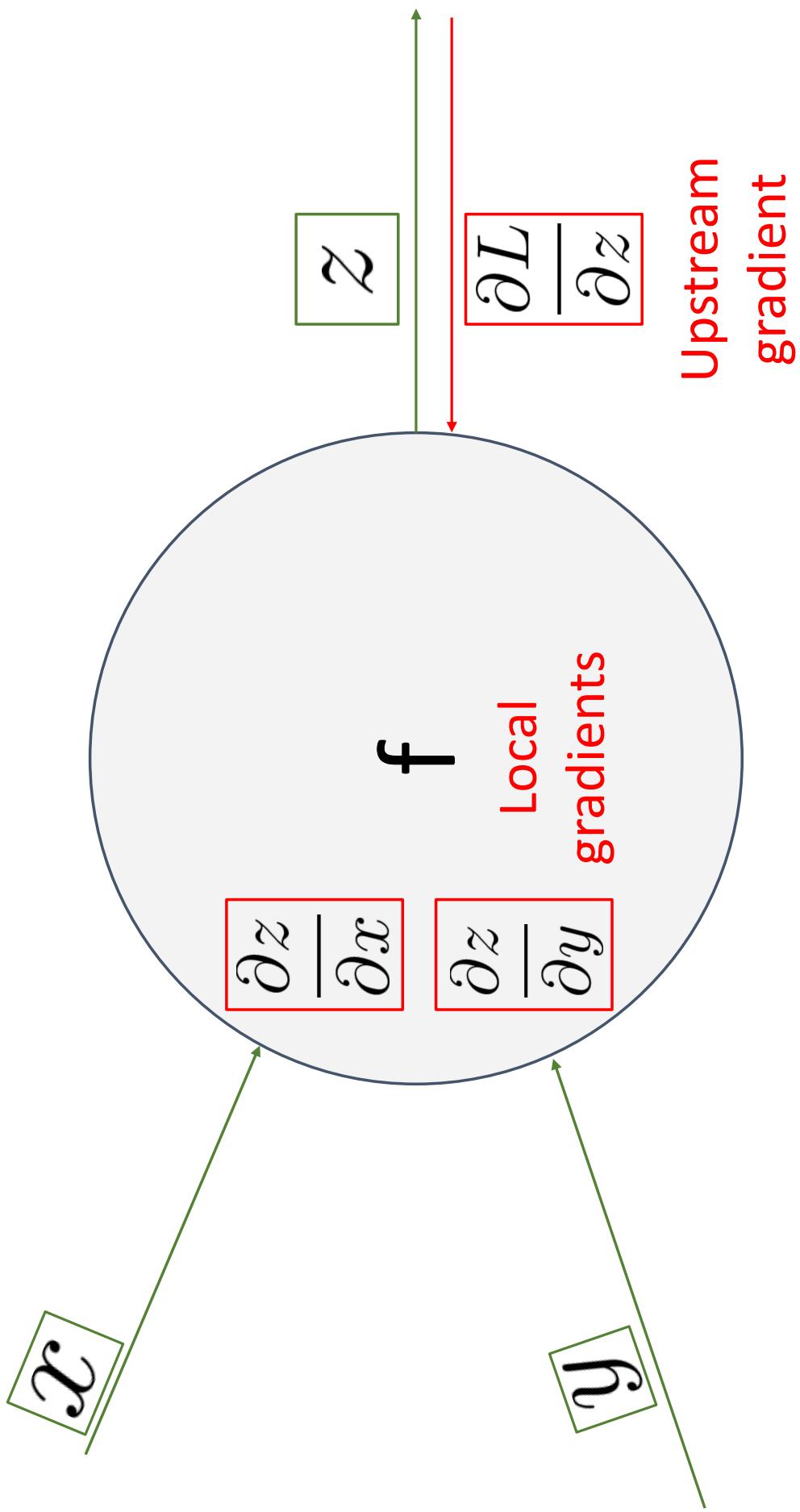
$$\frac{\partial q}{\partial x} = 1$$

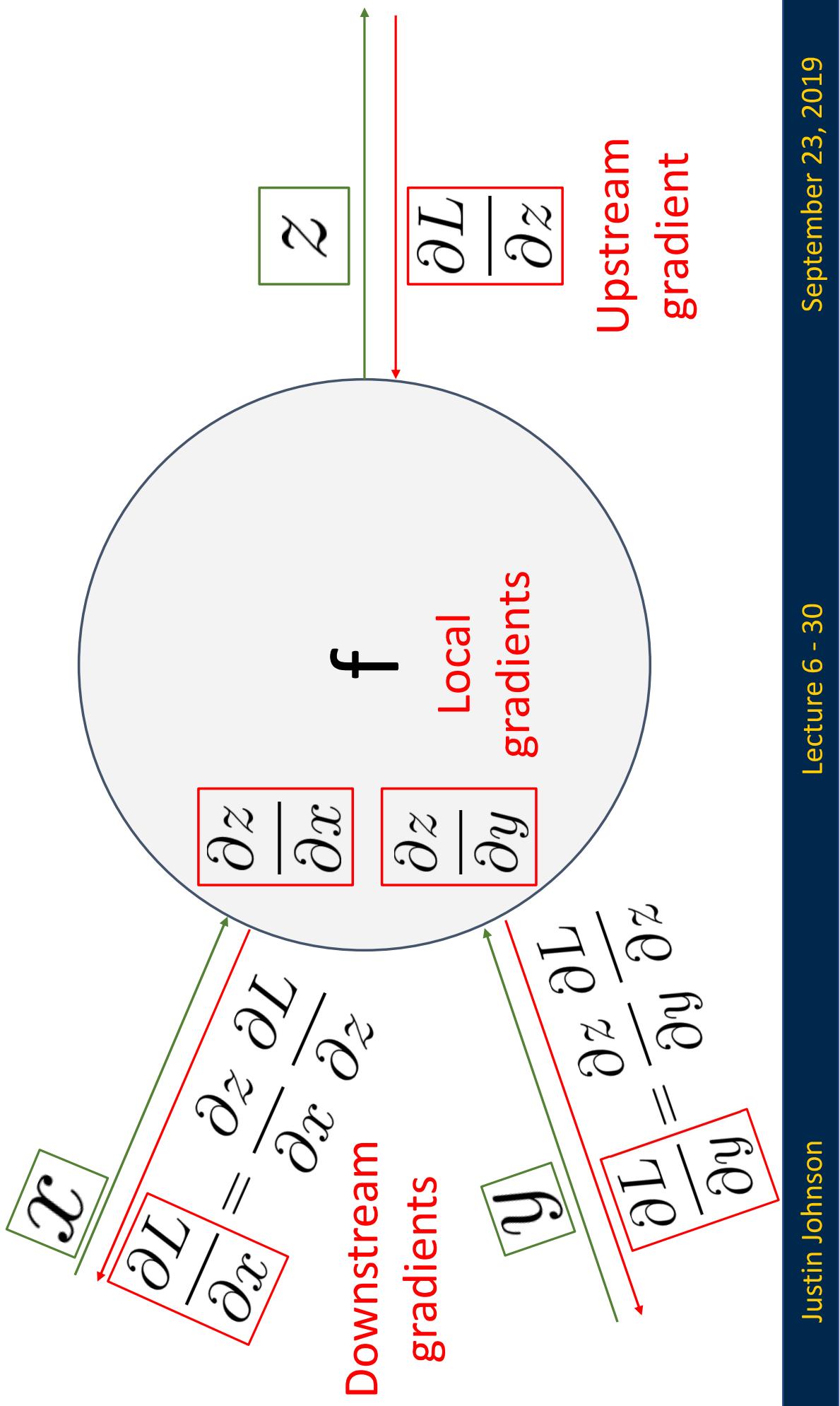
```

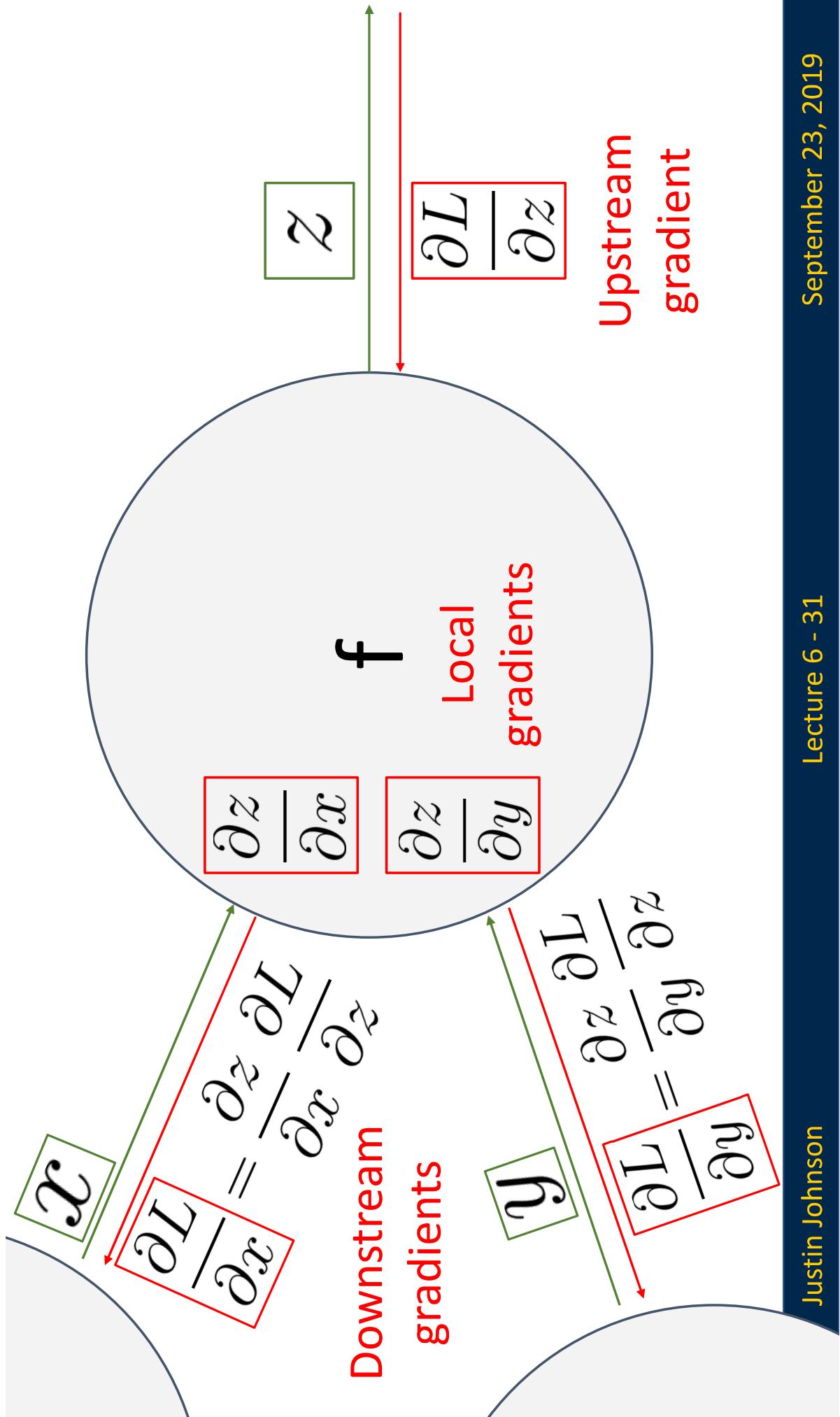
graph TD
    DG[Downstream Gradient] --> LG[Local Gradient]
    LG --> UG[Upstream Gradient]
  
```



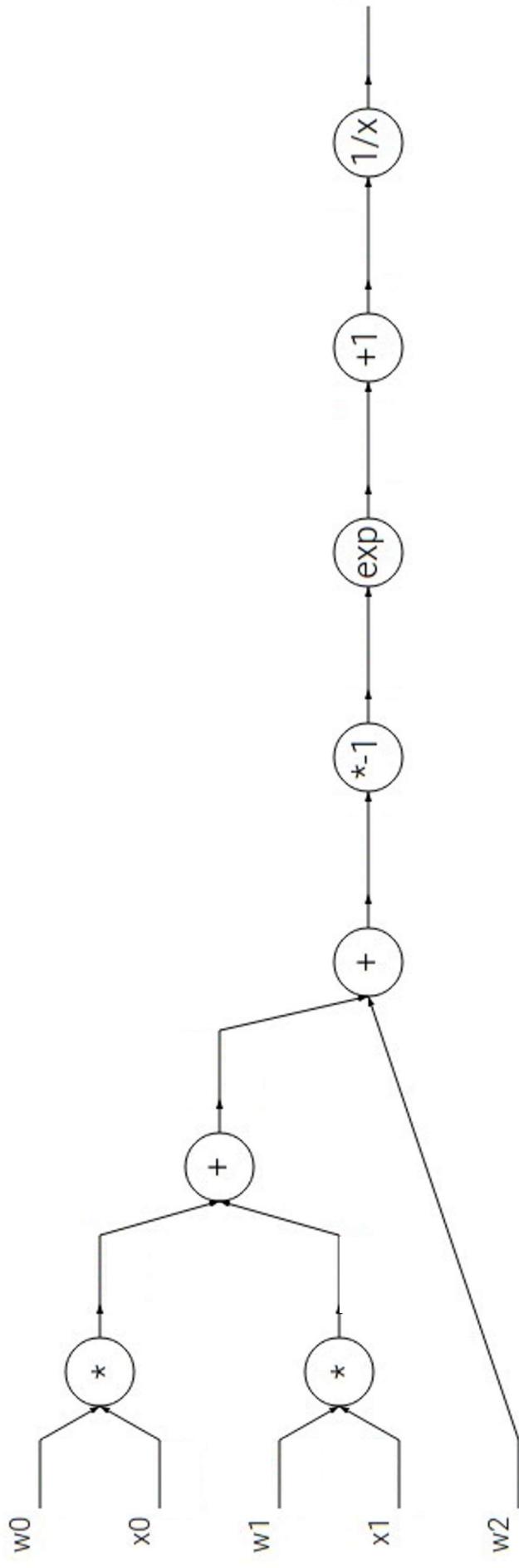






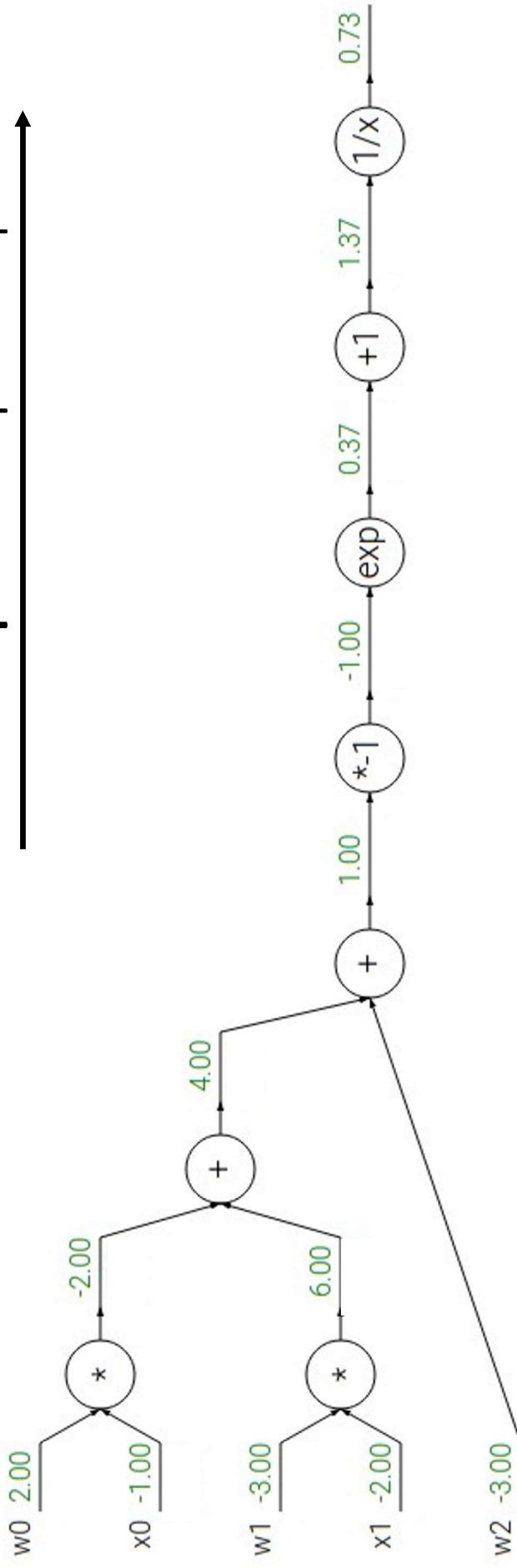


$$\text{Another Example } f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$



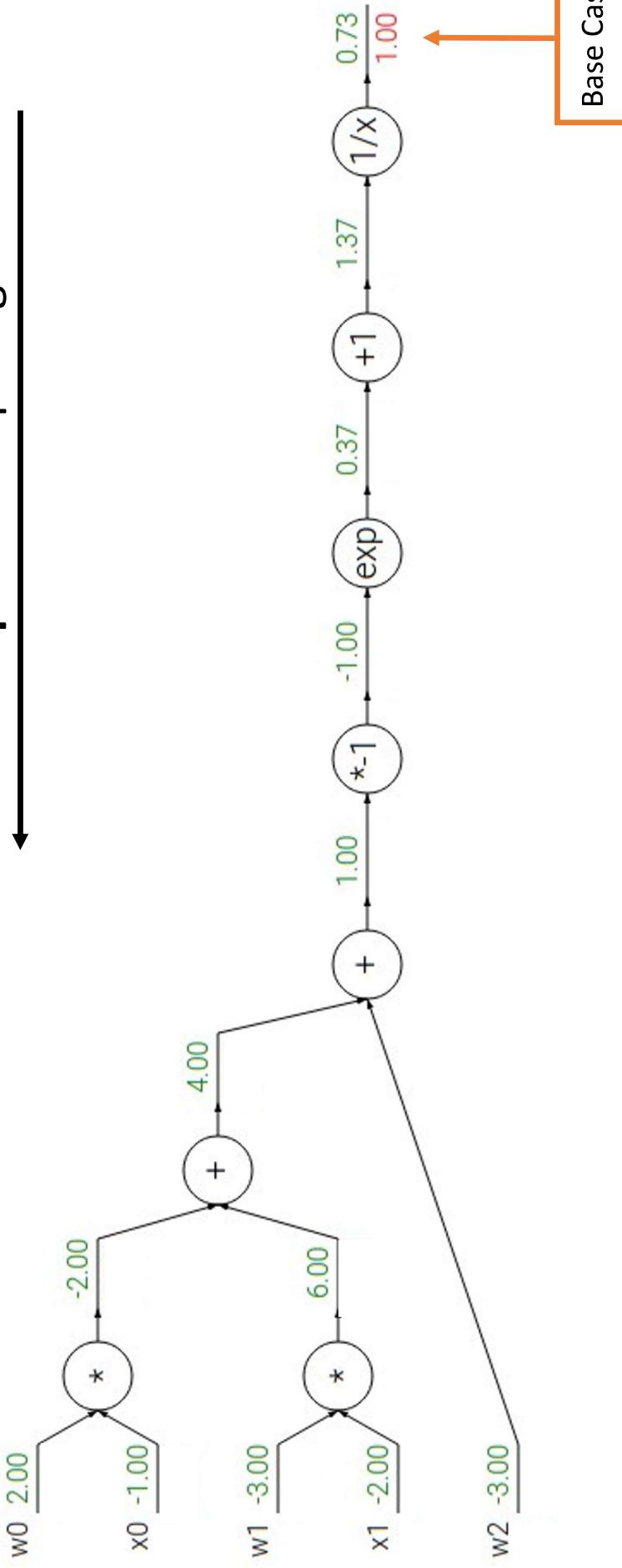
$$\text{Another Example } f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

**Forward pass: Compute outputs**



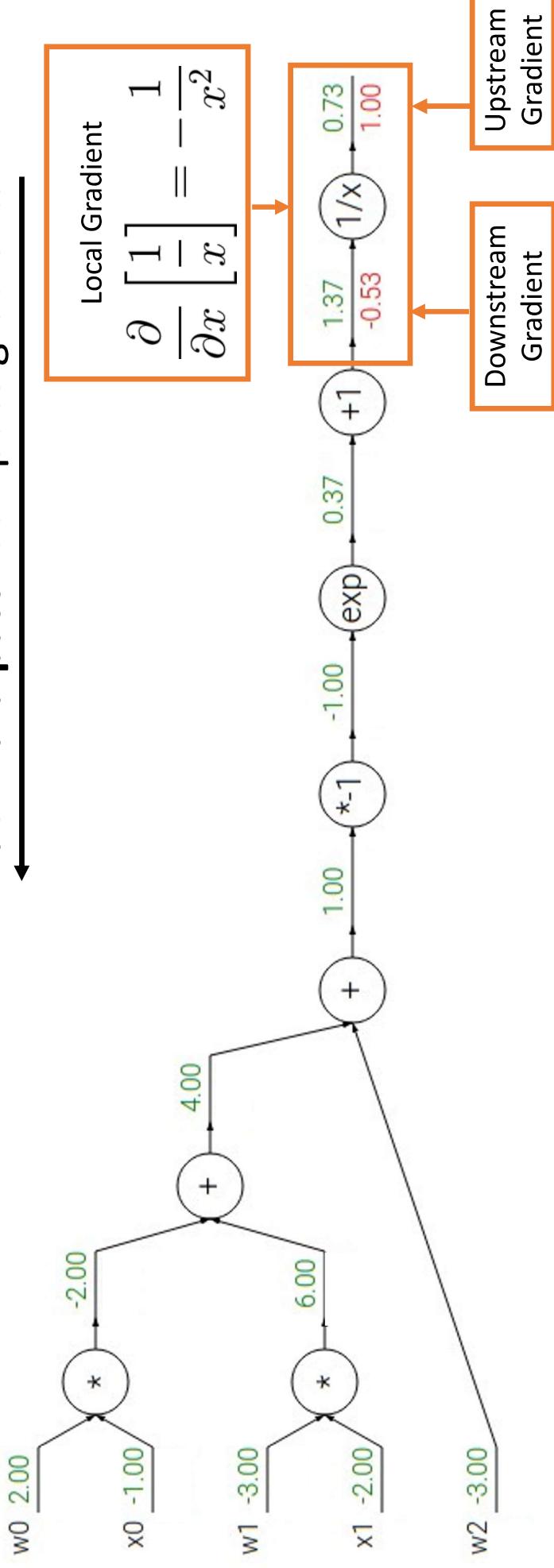
$$\text{Another Example } f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

**Backward pass: Compute gradients**



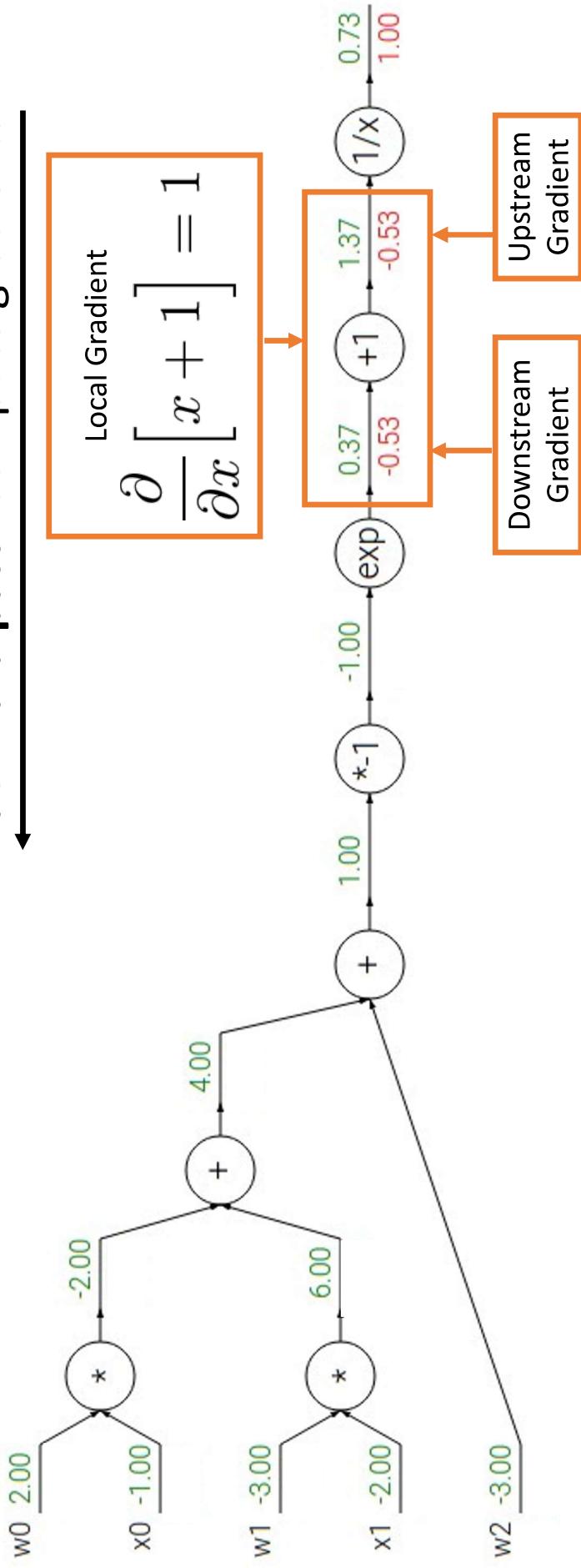
$$f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

**Backward pass: Compute gradients**



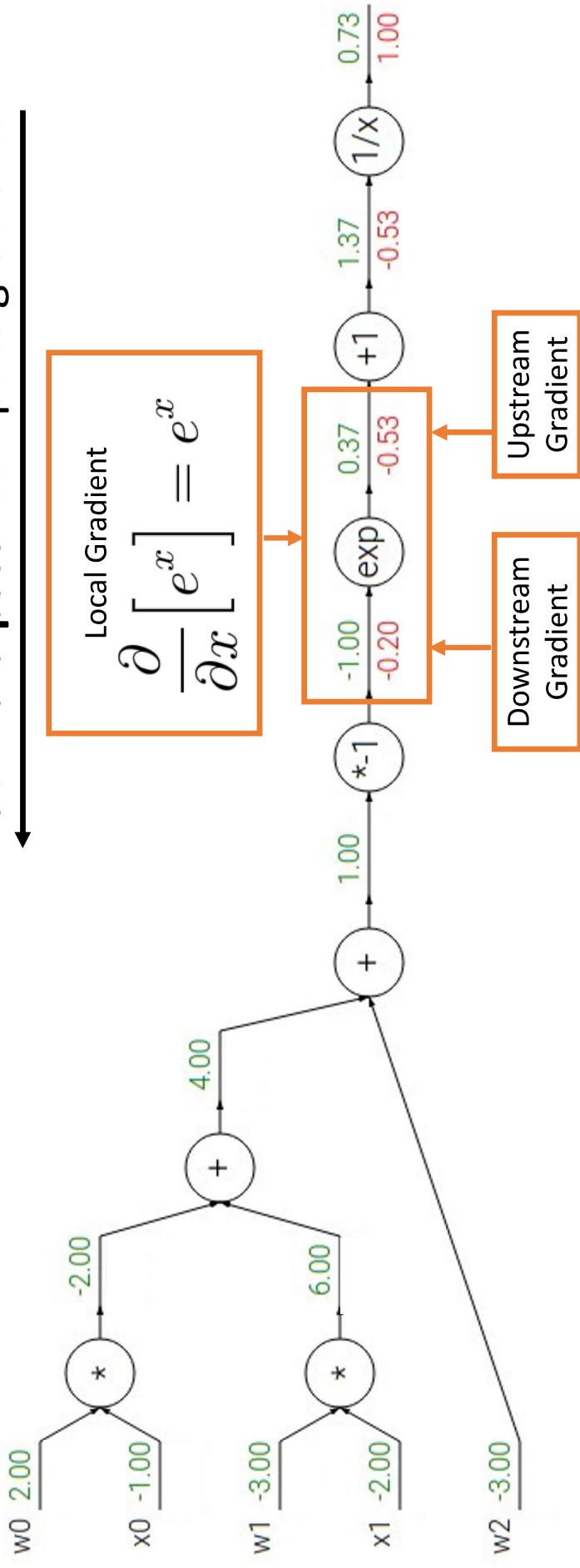
Another Example  $f(x, w) = \frac{1}{1 + e^{-(w_0x_0+w_1x_1+w_2)}}$

**Backward pass:** Compute gradients



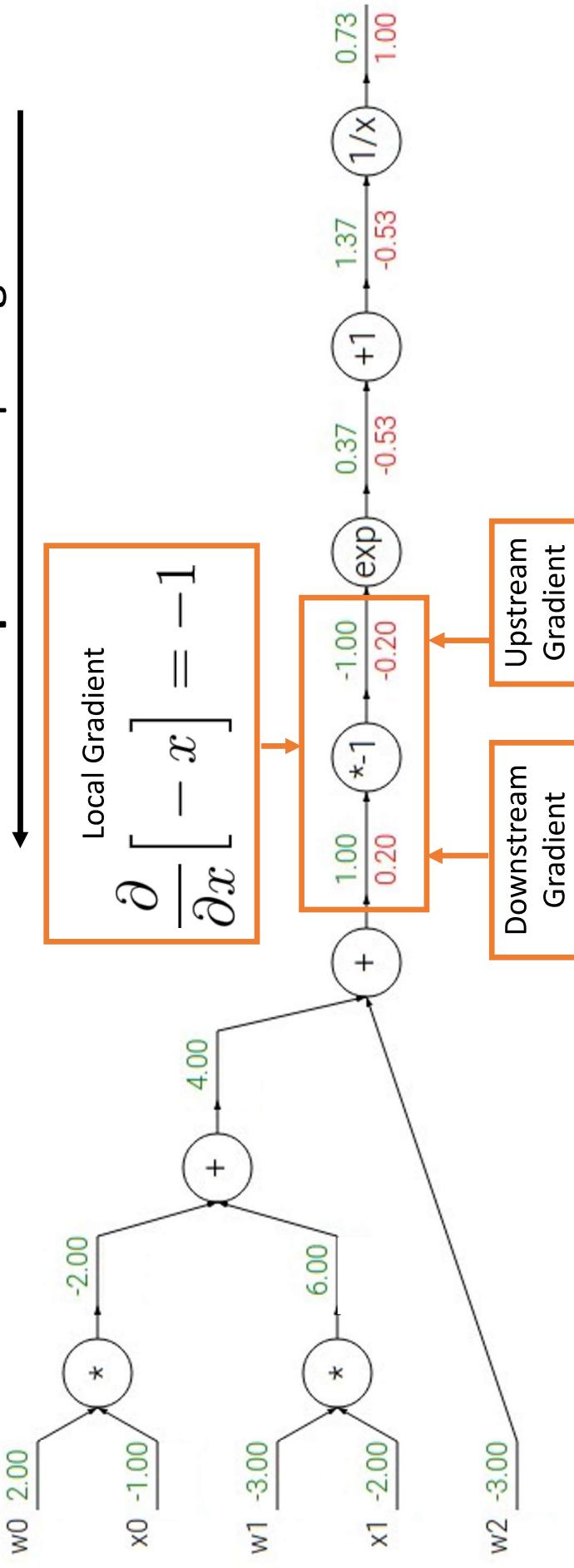
$$\text{Another Example } f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

**Backward pass: Compute gradients**



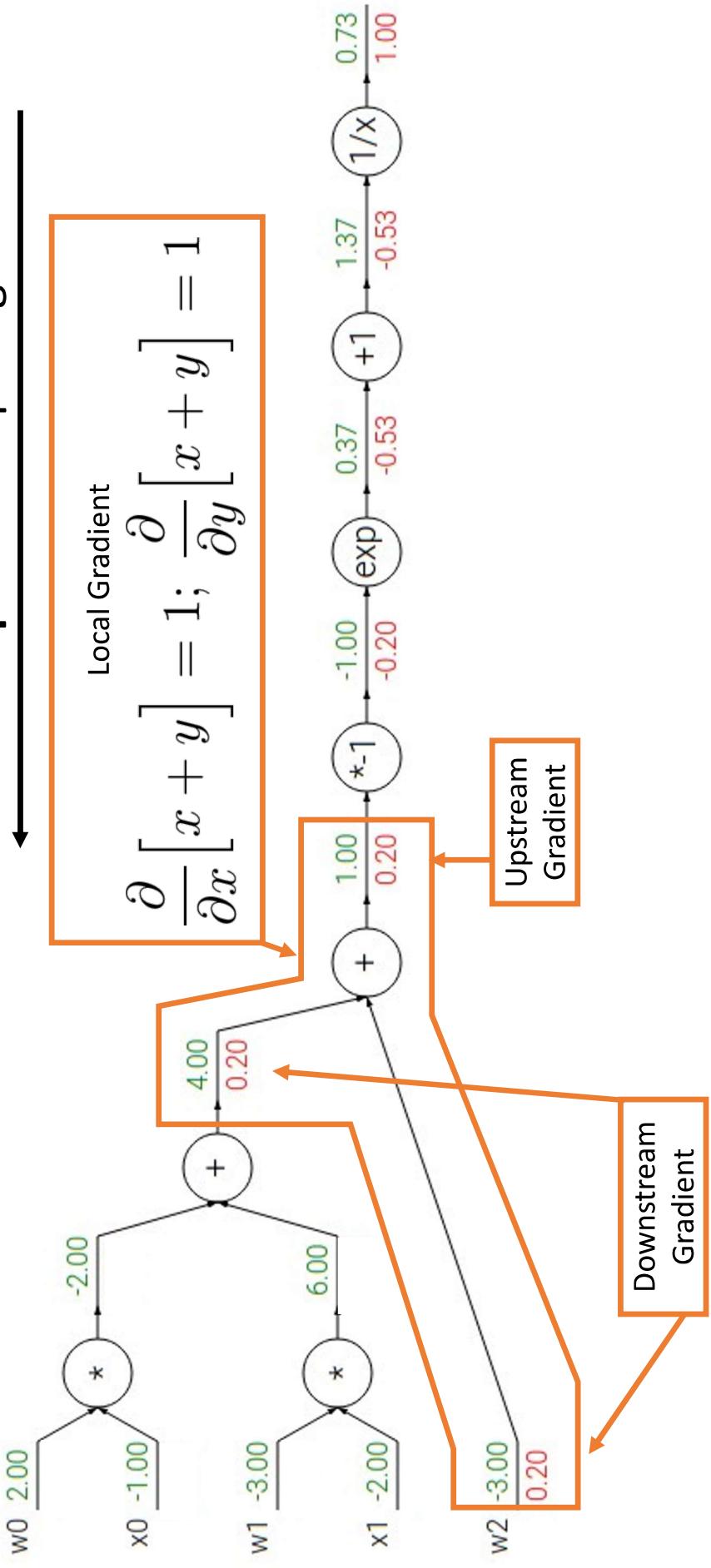
$$\text{Another Example } f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

**Backward pass: Compute gradients**



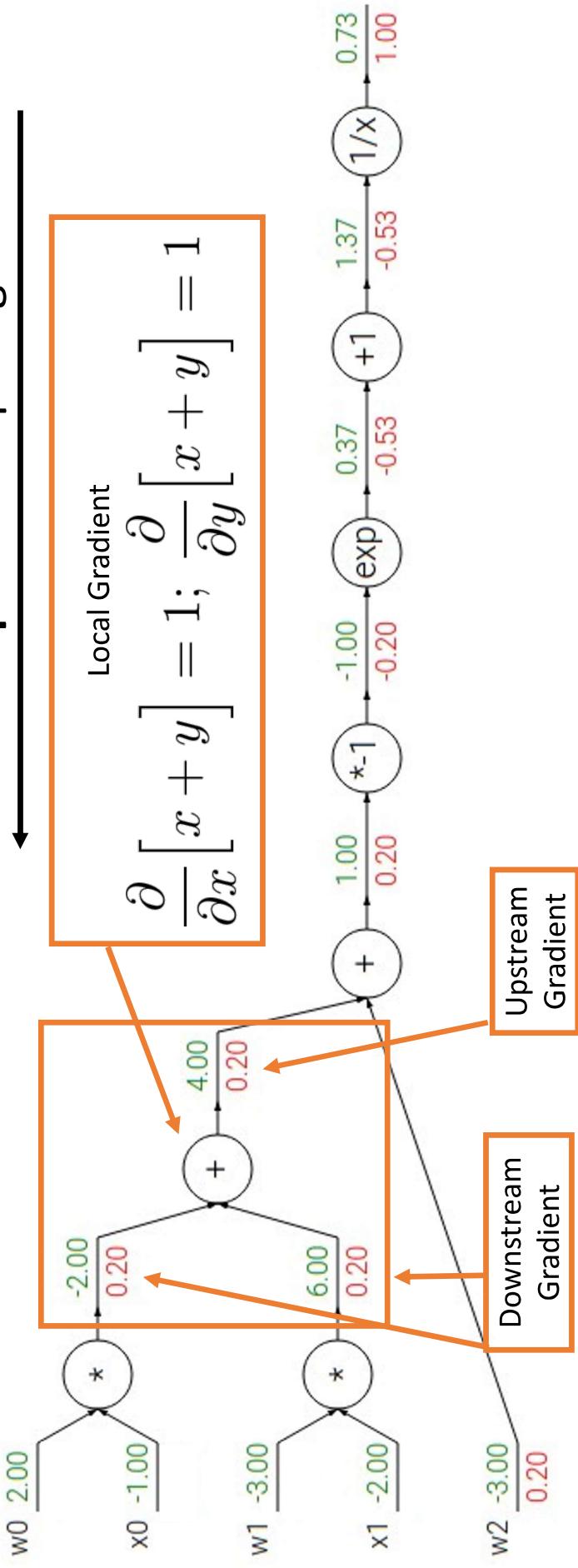
$$\text{Another Example } f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

**Backward pass: Compute gradients**



$$\text{Another Example } f(x, w) = \frac{1}{1 + e^{-(w_0x_0+w_1x_1+w_2)}}$$

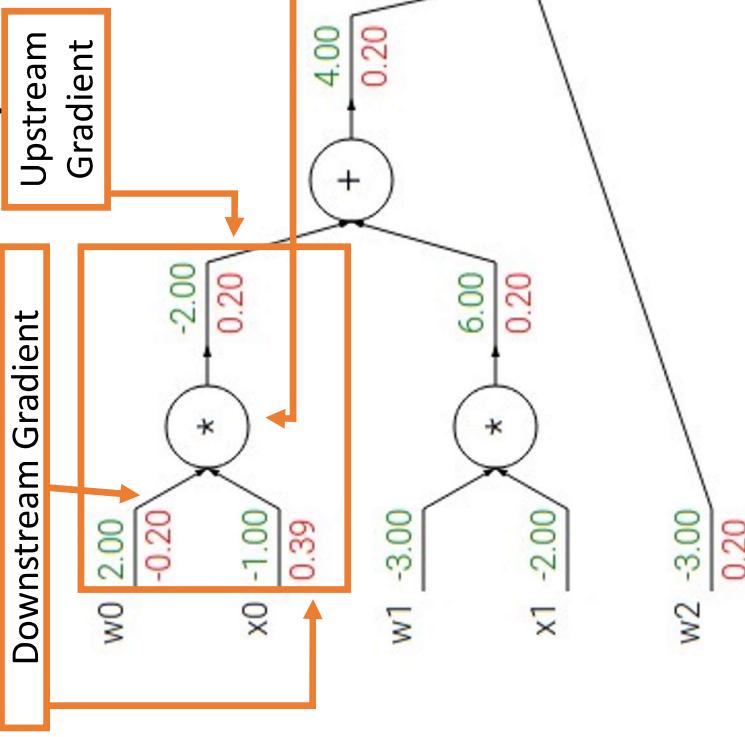
**Backward pass: Compute gradients**



## Another Example

$$f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

**Backward pass: Compute gradients**



Local Gradient

$$\frac{\partial}{\partial x} [xy] = y; \quad \frac{\partial}{\partial y} [xy] = x$$

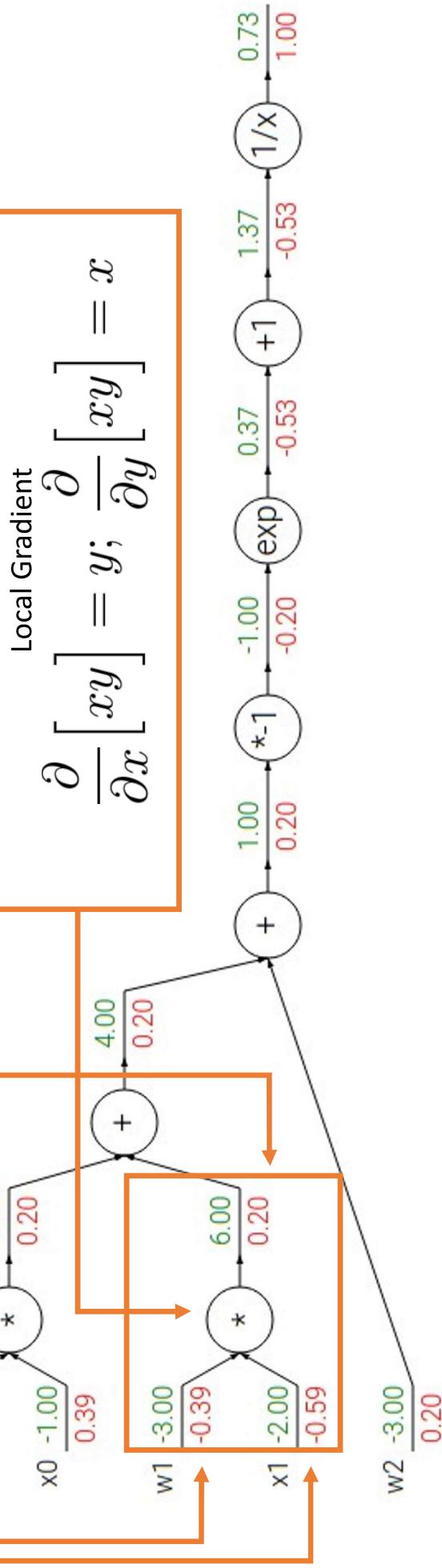
## Another Example

$$f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Downstream Gradient

Upstream Gradient

**Backward pass: Compute gradients**

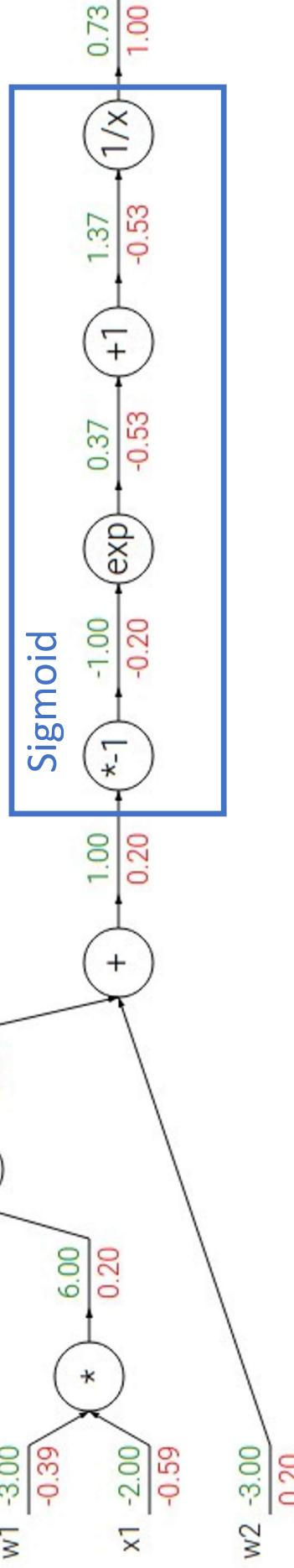


$$\text{Another Example } f(x, w) = \frac{1}{1 + e^{-(w_0x_0+w_1x_1+w_2)}} = \sigma(w_0x_0 + w_1x_1 + w_2)$$

**Backward pass: Compute gradients**

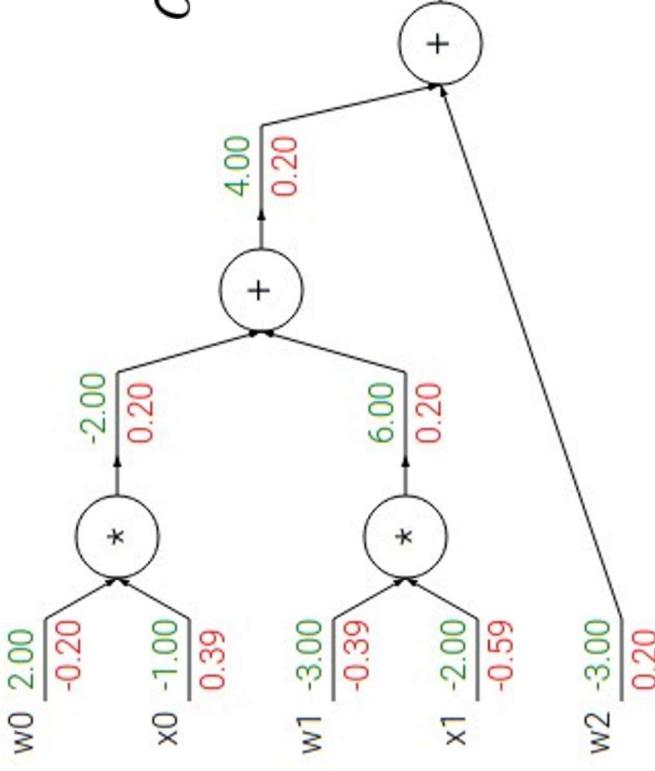
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Computational graph is not unique: we can use primitives that have simple local gradients

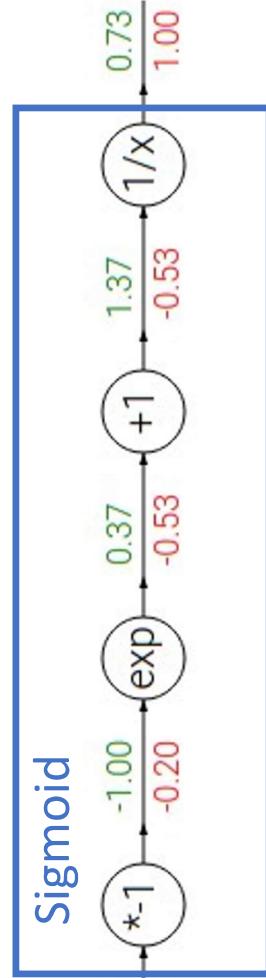


$$\text{Another Example } f(x, w) = \frac{1}{1 + e^{-(w_0x_0+w_1x_1+w_2)}} = \sigma(w_0x_0 + w_1x_1 + w_2)$$

**Backward pass: Compute gradients**



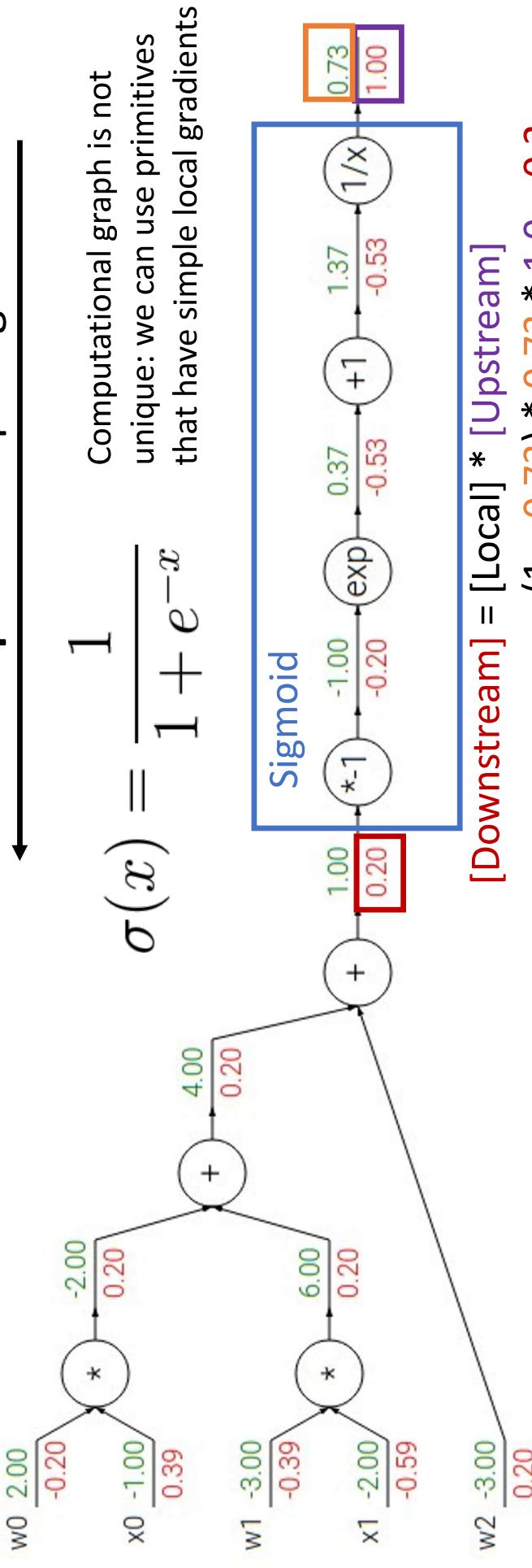
Computational graph is not unique: we can use primitives that have simple local gradients



Sigmoid local gradient:  $\frac{\partial}{\partial x} [\sigma(x)] = \frac{e^{-x}}{(1 + e^{-x})^2} = \left( \frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left( \frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x))\sigma(x)$

$$\text{Another Example } f(x, w) = \frac{1}{1 + e^{-(w_0x_0+w_1x_1+w_2)}} = \sigma(w_0x_0 + w_1x_1 + w_2)$$

**Backward pass: Compute gradients**

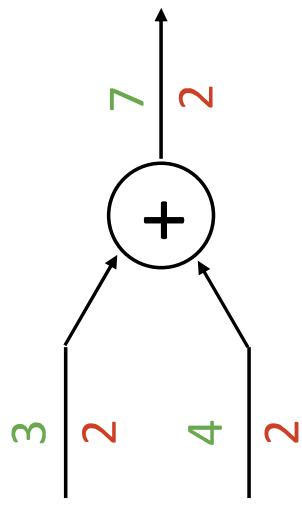


$$\begin{aligned} [\text{Downstream}] &= [\text{Local}] * [\text{Upstream}] \\ &= (1 - 0.73) * 0.73 * 1.0 = 0.2 \end{aligned}$$

$$\text{Sigmoid local gradient: } \frac{\partial}{\partial x} [\sigma(x)] = \frac{e^{-x}}{(1 + e^{-x})^2} = \left( \frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left( \frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x))\sigma(x)$$

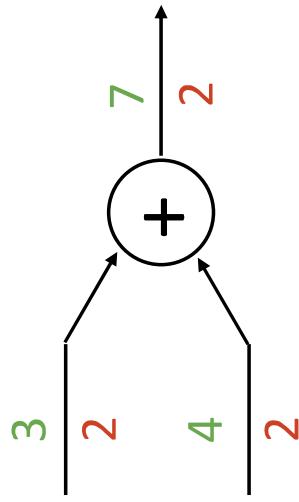
# Patterns in Gradient Flow

**add gate:** gradient distributor

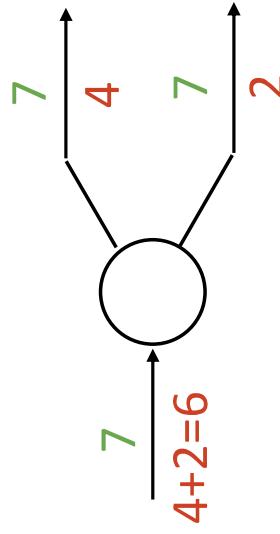


# Patterns in Gradient Flow

**add gate:** gradient distributor

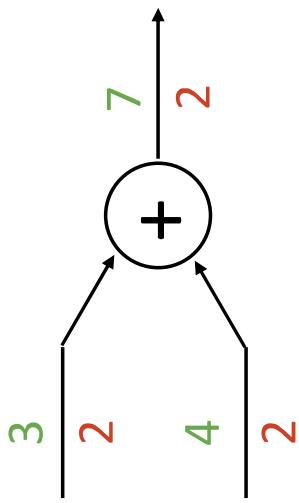


**copy gate:** gradient adder

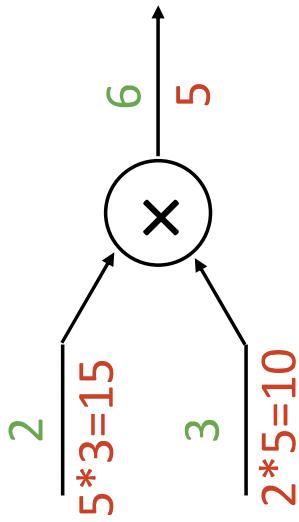


# Patterns in Gradient Flow

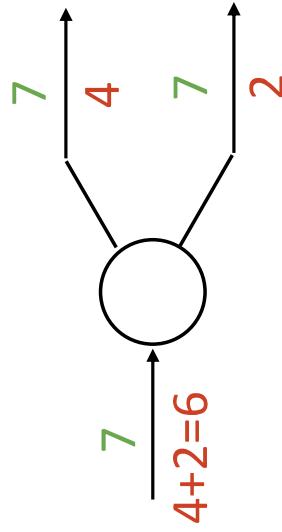
**add gate:** gradient distributor



**mul gate:** “swap multiplier”

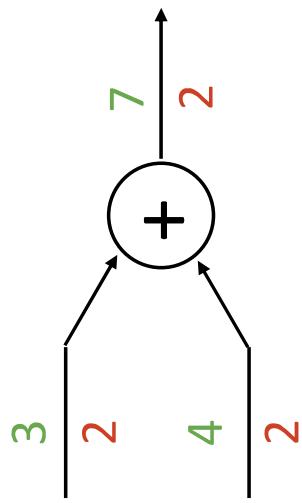


**copy gate:** gradient adder

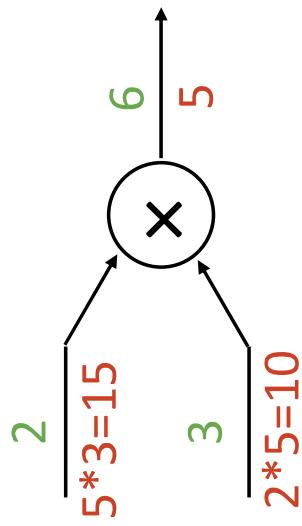


# Patterns in Gradient Flow

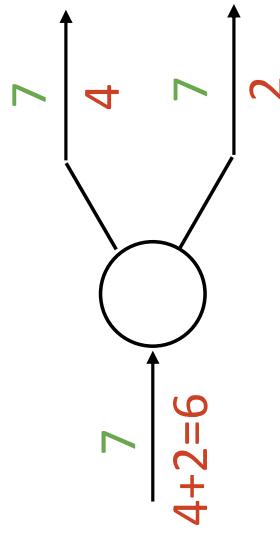
**add gate:** gradient distributor



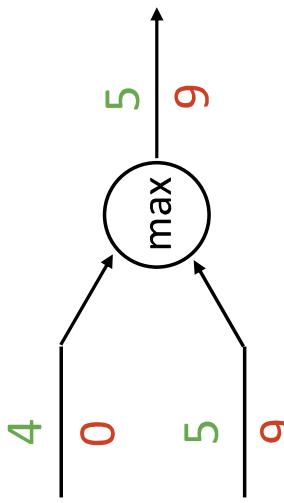
**mul gate:** “swap multiplier”



**copy gate:** gradient adder



**max gate:** gradient router



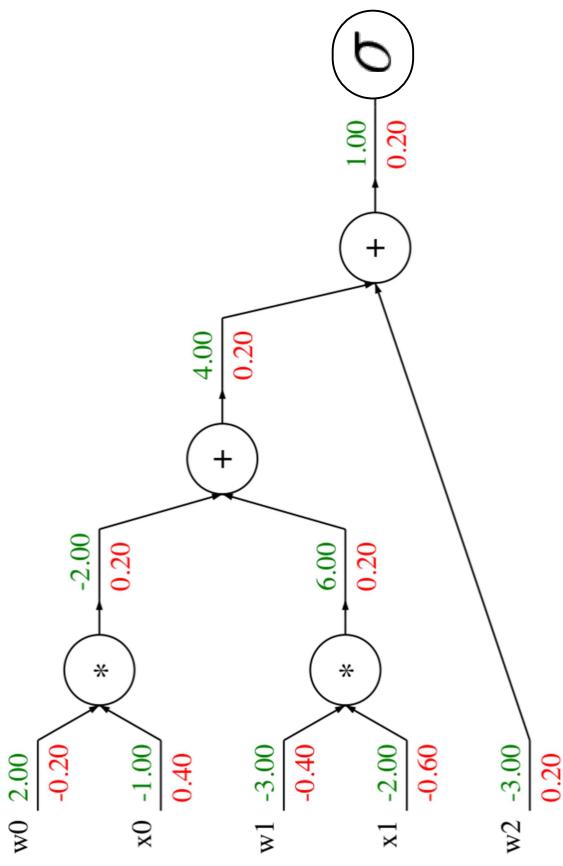
Backprop Implementation:

“Flat” gradient code:

Forward pass:  
Compute output

```
def f(w0, x0, w1, x1, w2):
```

```
s0 = w0 * x0  
s1 = w1 * x1  
s2 = s0 + s1  
s3 = s2 + w2  
L = sigmoid(s3)
```

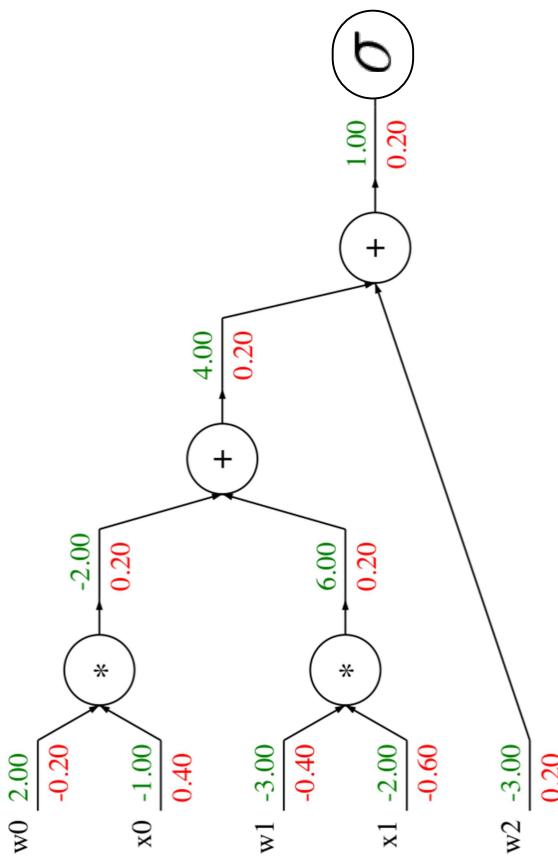


Backprop Implementation:

“Flat” gradient code:

Forward pass:  
Compute output

```
def f(w0, x0, w1, x1, w2):  
    s0 = w0 * x0  
    s1 = w1 * x1  
    s2 = s0 + s1  
    s3 = s2 + w2  
    L = sigmoid(s3)
```



Backward pass:  
Compute grads

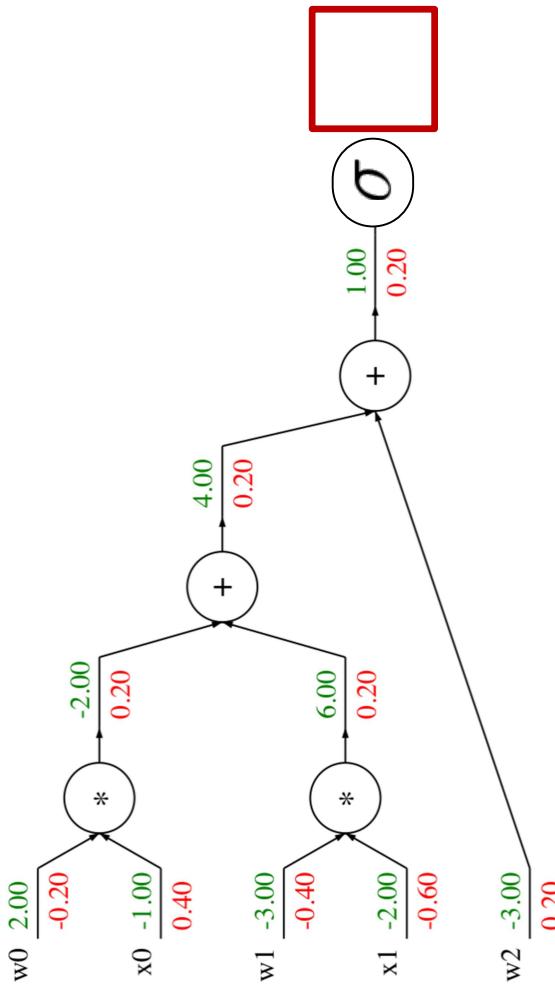
```
grad_L = 1.0  
grad_s3 = grad_L * (1 - L) * L  
grad_w2 = grad_s3  
grad_s2 = grad_s3  
grad_s0 = grad_s2  
grad_s1 = grad_s2  
grad_w1 = grad_s1 * x1  
grad_x1 = grad_s1 * w1  
grad_w0 = grad_s0 * x0  
grad_x0 = grad_s0 * w0
```

Backprop Implementation:

”Flat” gradient code:

Forward pass:  
Compute output

```
def f(w0, x0, w1, x1, w2):  
    s0 = w0 * x0  
    s1 = w1 * x1  
    s2 = s0 + s1  
    s3 = s2 + w2  
    L = sigmoid(s3)
```



Base Case

```
grad_L = 1.0  
grad_s3 = grad_L * (1 - L) * L  
grad_w2 = grad_s3  
grad_s2 = grad_s3  
grad_s0 = grad_s2  
grad_s1 = grad_s2  
grad_w1 = grad_s1 * x1  
grad_x1 = grad_s1 * w1  
grad_w0 = grad_s0 * x0  
grad_x0 = grad_s0 * w0
```

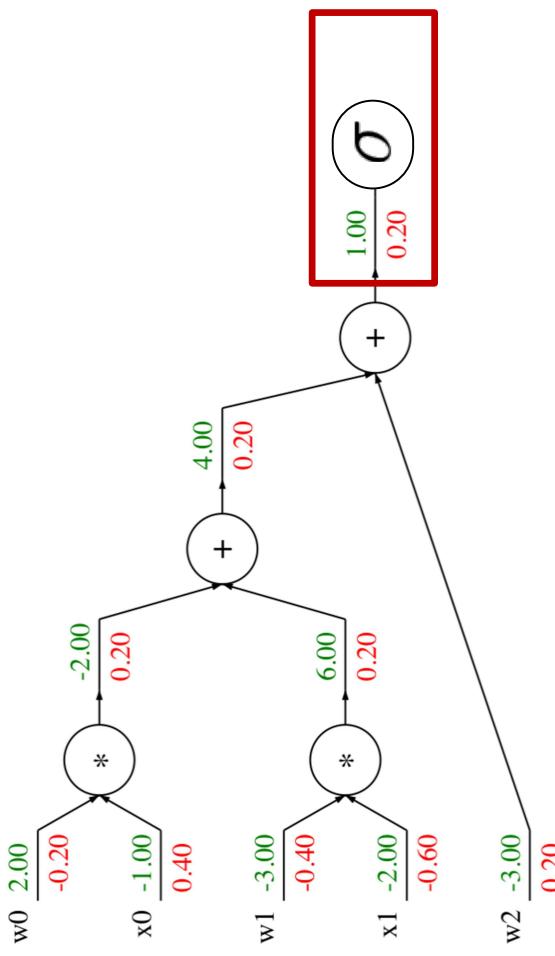
Backprop Implementation:

”Flat” gradient code:

Forward pass:  
Compute output

```
def f(w0, x0, w1, x1, w2):
```

```
    s0 = w0 * x0
    s1 = w1 * x1
    s2 = s0 + s1
    s3 = s2 + w2
    L = sigmoid(s3)
```



Sigmoid

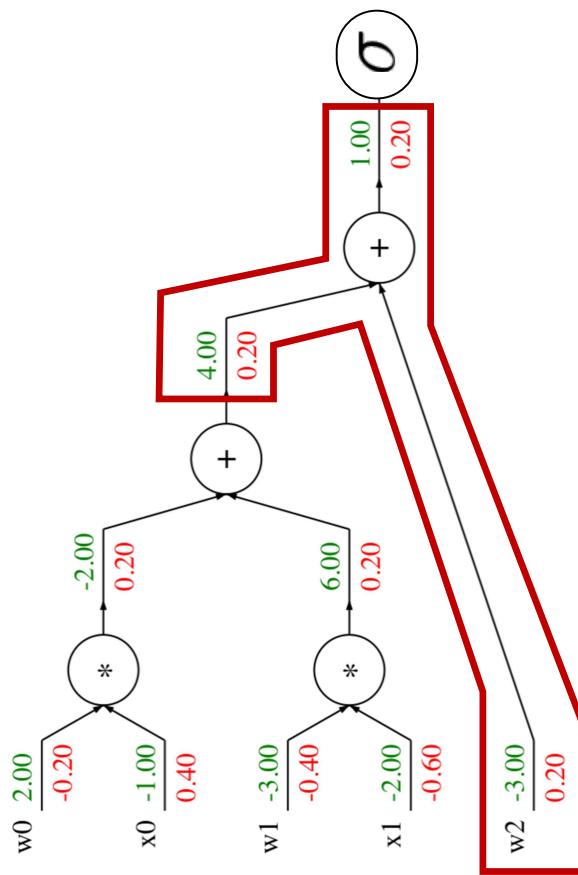
```
grad_L = 1.0
grad_s3 = grad_L * (1 - L) * L
grad_w2 = grad_s3
grad_s2 = grad_s3
grad_s0 = grad_s2
grad_s1 = grad_s2
grad_w1 = grad_s1 * x1
grad_x1 = grad_s1 * w1
grad_w0 = grad_s0 * x0
grad_x0 = grad_s0 * w0
```

Backprop Implementation:

“Flat” gradient code:

Forward pass:  
Compute output

```
def f(w0, x0, w1, x1, w2):  
    s0 = w0 * x0  
    s1 = w1 * x1  
    s2 = s0 + s1  
    s3 = s2 + w2  
    L = sigmoid(s3)
```



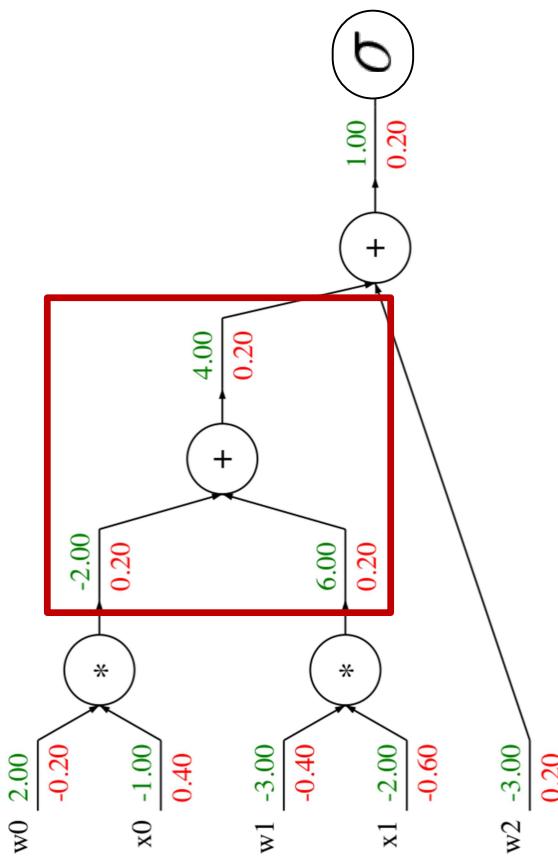
Backprop Implementation:

“Flat” gradient code:

Forward pass:  
Compute output

```
def f(w0, x0, w1, x1, w2):
```

```
s0 = w0 * x0  
s1 = w1 * x1  
s2 = s0 + s1  
s3 = s2 + w2  
L = sigmoid(s3)
```



```
grad_L = 1.0  
grad_s3 = grad_L * (1 - L) * L  
grad_w2 = grad_s3  
grad_s2 = grad_s3  
grad_s0 = grad_s2  
grad_s1 = grad_s2  
grad_w1 = grad_s1 * x1  
grad_x1 = grad_s1 * w1  
grad_w0 = grad_s0 * x0  
grad_x0 = grad_s0 * w0
```

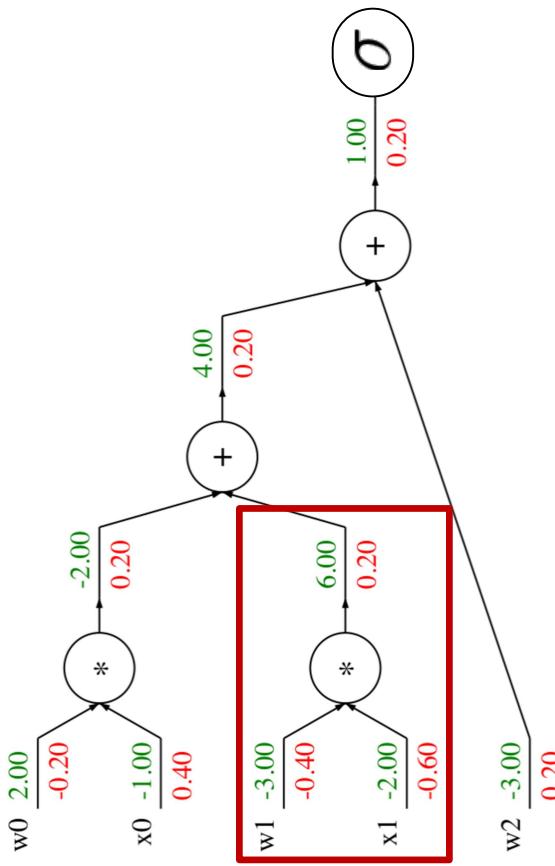
Backprop Implementation:

“Flat” gradient code:

Forward pass:  
Compute output

```
def f(w0, x0, w1, x1, w2):
```

```
s0 = w0 * x0  
s1 = w1 * x1  
s2 = s0 + s1  
s3 = s2 + w2  
L = sigmoid(s3)
```

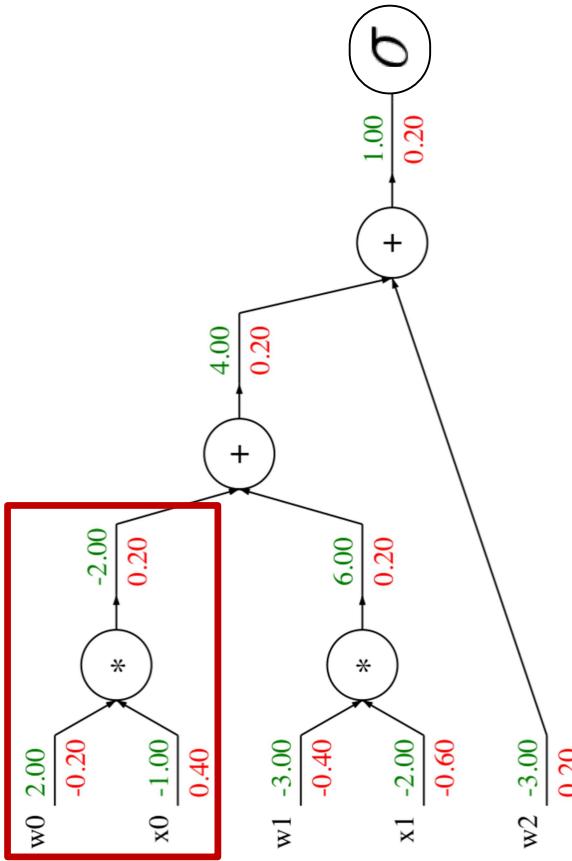


Multiply

```
grad_L = 1.0  
grad_s3 = grad_L * (1 - L) * L  
grad_w2 = grad_s3  
grad_s2 = grad_s3  
grad_s0 = grad_s2  
grad_s1 = grad_s2  
grad_w1 = grad_s1 * x1  
grad_x1 = grad_s1 * w1  
grad_w0 = grad_s0 * x0  
grad_x0 = grad_s0 * w0
```

# Backprop Implementation:

”Flat” gradient code:



```
def f(w0, x0, w1, x1, w2):
```

$$\begin{aligned}
 s_0 &= w_0 * x_0 \\
 s_1 &= w_1 * x_1 \\
 s_2 &= s_0 + s_1 \\
 s_3 &= s_2 + w_2 \\
 L &= \text{sigmoid}(s_3)
 \end{aligned}$$

**Forward pass:**

## Compute output

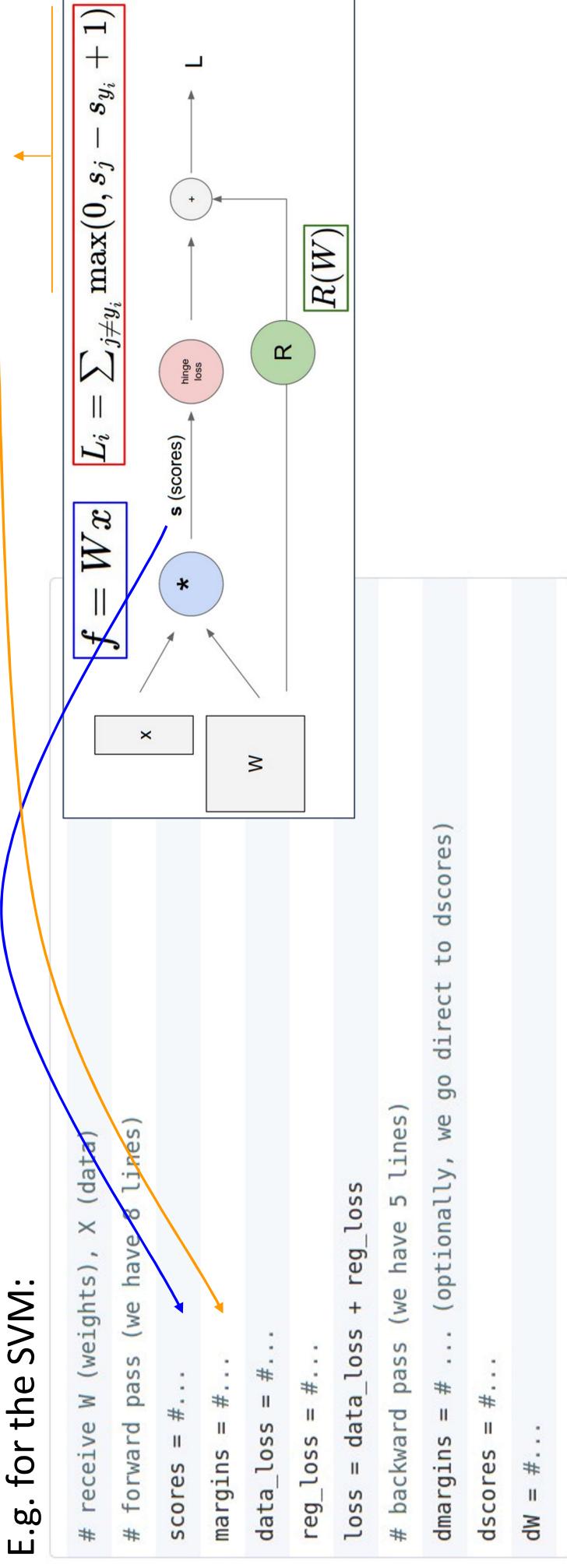
```
grad_L = 1.0  
grad_s3 = grad_L * (1 - L) * L  
grad_w2 = grad_s3  
grad_s2 = grad_s3  
grad_s0 = grad_s2  
grad_s1 = grad_s2  
grad_w1 = grad_s1 * x1  
grad_x1 = grad_s1 * w1  
grad_w0 = grad_s0 * x0  
grad_x0 = grad_s0 * w0
```

## Multiply

# “Flat” Backprop: Do this for Assignment 2!

Your gradient code should look like a “reversed version” of your forward pass!

E.g. for the SVM:



## “Flat” Backprop: Do this for Assignment 2!

Your gradient code should look like a “reversed version” of your forward pass!

E.g. for two-layer neural net:

```
# receive w1,w2,b1,b2 (weights/biases), X (data)

# forward pass:
h1 = #... function of X,W1,b1
scores = #... function of h1,W2,b2
loss = #... (several lines of code to evaluate Softmax loss)

# backward pass:
dscores = #...
dh1,dW2,db2 = #...
dW1,db1 = #...
```

# Backprop Implementation: Modular API

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

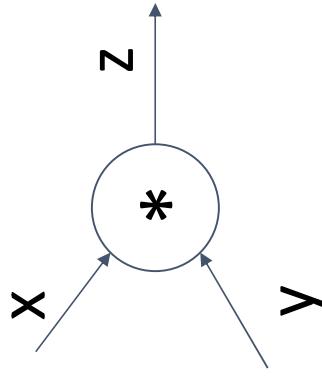
```
class ComputationalGraph(object):
    ...
    def forward(inputs):
        # 1. [pass inputs to input gates...]
        # 2. forward the computational graph:
        for gate in self.graph.nodes_topologically_sorted():
            gate.forward()
        return loss # the final gate in the graph outputs 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
```

The diagram illustrates a computational graph with the following structure:

- Inputs:  $w_0$  (2.00),  $x_0$  (-1.00, 0.40),  $w_1$  (-3.00),  $x_1$  (-2.00, -0.60), and  $w_2$  (-3.00, 0.20).
- Operations:
  - $w_0$  is multiplied by  $x_0$  with weight -0.20.
  - The result of the multiplication is multiplied by  $x_0$  again with weight 0.20.
  - The result of the second multiplication is added to the result of  $w_1 \times x_1$  with weight 4.00.
  - The result of the addition is multiplied by  $x_1$  with weight 0.20.
  - The result of the multiplication is added to the result of  $w_2 \times x_2$  with weight 1.00.
  - The final result is multiplied by  $x_2$  with weight 0.20.
- Output: The final output is labeled  $\sigma$ .

# Example: PyTorch Autograd Functions



```
class Multiply(torch.autograd.Function):
    @staticmethod
    def forward(ctx, x, y):
        ctx.save_for_backward(x, y)
        z = x * y
        return z

    @staticmethod
    def backward(ctx, grad_z):
        x, y = ctx.saved_tensors
        grad_x = y * grad_z # dz/dx * dL/dz
        grad_y = x * grad_z # dz/dy * dL/dz
        return grad_x, grad_y
```

Annotations for the code:

- Annotations for the forward pass:
  - An annotation "Need to stash some values for use in backward" points to the line `ctx.save_for_backward(x, y)`.
  - An annotation "Upstream gradient" points to the line `grad_x = y * grad_z`.
  - An annotation "Multiply upstream and local gradients" points to the line `grad_y = x * grad_z`.
- Annotations for the backward pass:
  - An annotation "Upstream gradient" points to the line `grad_x = y * grad_z`.
  - An annotation "Multiply upstream and local gradients" points to the line `grad_y = x * grad_z`.

( $x, y, z$  are scalars)

# Example: PyTorch operators

The screenshot shows a GitHub search results page for the PyTorch repository. The search term is 'canonicalize'. There are 6,340 results. The results are listed in chronological order, with each commit being 4 months ago. The commits are mostly from the 'PyTorch / aten / src / THNN / generic /' directory. The commits are titled with various PyTorch operator names like 'AbsCriterion.c', 'BCECriterion.c', 'ClassNLLCriterion.c', 'Co2m.c', 'ELU.c', 'FeatureFPooling.c', 'GatedLinearUnit.c', 'Hardtanh.c', 'Im2Col.c', 'IndexLinear.c', 'LeakyReLU.c', 'LogSigmoid.c', 'MSECriterion.c', 'MultiLabelMarginCriterion.c', 'MultiMarginCriterion.c', 'RReLU.c', 'Sigmoid.c', 'SmoothL1Criterion.c', 'SoftMarginCriterion.c', 'Sorpus.c', 'SoftShrink.c', 'SparseLinear.c', 'SpatialAdaptiveAveragePooling.c', 'SpatialAdaptiveMaxPooling.c', and 'SpatialAveragePooling.c'. Each commit has a detailed description of what it does, such as 'Canonicalize all includes in PyTorch. (#r4849)'.

Commit	Description
SpatialClassNLLCriterion.c	Canonicalize all includes in PyTorch. (#r4849)
SpatialConvolutionNLLCriterion.c	Canonicalize all includes in PyTorch. (#r4849)
SpatialDilatedConvolution.c	Canonicalize all includes in PyTorch. (#r4849)
SpatialDilatedMaxPooling.c	Canonicalize all includes in PyTorch. (#r4849)
SpatialFractionalMaxPooling.c	Canonicalize all includes in PyTorch. (#r4849)
SpatialFullDilatedConvolution.c	Canonicalize all includes in PyTorch. (#r4849)
SpatialMaxUpSampling.c	Canonicalize all includes in PyTorch. (#r4849)
SpatialReflectionPadding.c	Canonicalize all includes in PyTorch. (#r4849)
SpatialReplicationPadding.c	Canonicalize all includes in PyTorch. (#r4849)
SpatialUpSamplingLinear.c	Canonicalize all includes in PyTorch. (#r4849)
SpatialUpSamplingNearest.c	Canonicalize all includes in PyTorch. (#r4849)
THNN.h	Canonicalize all includes in PyTorch. (#r4849)
Tanh.c	Canonicalize all includes in PyTorch. (#r4849)
TemporalReflectionPadding.c	Canonicalize all includes in PyTorch. (#r4849)
TemporalReplicationPadding.c	Canonicalize all includes in PyTorch. (#r4849)
TemporalRowConvolution.c	Canonicalize all includes in PyTorch. (#r4849)
TemporalUpSamplingLinear.c	Canonicalize all includes in PyTorch. (#r4849)
TemporalUpSamplingNearest.c	Canonicalize all includes in PyTorch. (#r4849)
VolumetricAdaptiveAvgPool... VolumetricAdaptiveMaxPooling.c	Canonicalize all includes in PyTorch. (#r4849)
VolumetricRowConvolution.c	Canonicalize all includes in PyTorch. (#r4849)
VolumetricUpSamplingLinear.c	Canonicalize all includes in PyTorch. (#r4849)
VolumetricCorrelationMM.c	Canonicalize all includes in PyTorch. (#r4849)
VolumetricDilatedConvolution.c	Canonicalize all includes in PyTorch. (#r4849)
VolumetricDilatedMaxPooling.c	Canonicalize all includes in PyTorch. (#r4849)
VolumetricFracMaxPooling.c	Canonicalize all includes in PyTorch. (#r4849)
VolumetricInterpolatePooling.c	Canonicalize all includes in PyTorch. (#r4849)
VolumetricConvolution.c	Canonicalize all includes in PyTorch. (#r4849)
VolumetricUpSampling.c	Canonicalize all includes in PyTorch. (#r4849)
VolumetricReplicationPadding.c	Canonicalize all includes in PyTorch. (#r4849)
VolumetricSamplingNearest.c	Canonicalize all includes in PyTorch. (#r4849)
VolumetricSamplingTrilinear.c	Canonicalize all includes in PyTorch. (#r4849)
linear_upsampling_h pooling_shape_h	Implement nn.functional.interpolate based on upsample. (#r8591) Use integer math to compute output size of pooling operations (#r4405)
unfold.c	Canonicalize all includes in PyTorch. (#r4849)

```
1 #ifndef TH_GENERIC_FILE
2 #define TH_GENERIC_FILE "THNN/generic/Sigmoid.c"
3 #else
4
5 void THNN_(Sigmoid_updateOutput) (
6     THNNState *state,
7     THTensor *input,
8     THTensor *output)
9 {
10     THTensor_(sigmoid)(output, input);
11 }
12
13 void THNN_(Sigmoid_updateGradInput) (
14     THNNState *state,
15     THTensor *gradOutput,
16     THTensor *gradInput,
17     THTensor *output)
18 {
19     THNN_CHECK_NELEMENT(output, gradOutput);
20     THTensor_(resizeAs)(gradInput, output);
21     TH_TENSOR_APPLY3(scalar_t, gradInput, scalar_t, gradOutput, scalar_t, output,
22     scalar_t z = *output_data;
23     *gradInput_data = *gradOutput_data * (1. - z) * z;
24 );
25 }
26
27 #endif
```

## PyTorch sigmoid layer

```
1 #ifndef TH_GENERIC_FILE
2 #define TH_GENERIC_FILE "THNN/generic/Sigmoid.c"
3 #else
```

## PyTorch sigmoid layer

```
4
5 void THNN_(Sigmoid_updateOutput)(
6     THNNState *state,
7     THTensor *input,
8     THTensor *output)
9 {
10     THTensor_(sigmoid)(output, input);
11 }
```

Forward

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

```
12
13 void THNN_(Sigmoid_updateGradInput)(
14     THNNState *state,
15     THTensor *gradOutput,
16     THTensor *gradInput,
17     THTensor *output)
18 {
19     THNN_CHECK_NELEMENT(output, gradOutput);
20     THTensor_(resizeAs)(gradInput, output);
21     TH_TENSOR_APPLY3(scalar_t, gradInput, scalar_t, gradOutput, scalar_t, output,
22     scalar_t z = *output_data;
23     *gradInput_data = *gradOutput_data * (1. - z) * z;
24 );
25 }
26
27 #endif
```

Justin Johnson

Lecture 6 - 64

September 23, 2019

[Source](#)

# PyTorch sigmoid layer

```
1 #ifndef TH_GENERIC_FILE
2 #define TH_GENERIC_FILE "THNN/generic/Sigmoid.c"
3 #else
4
5 void THNN_(Sigmoid_updateOutput)(
6     THNNState *state,
7     THTensor *input,
8     THTensor *output)
9 {
10     THTensor_(sigmoid)(output, input);
11 }
12
13 void THNN_(Sigmoid_updateGradInput)(
14     THNNState *state,
15     THTensor *gradOutput,
16     THTensor *gradInput,
17     THTensor *output)
18 {
19     THNN_CHECK_NELEMENT(output, gradOutput);
20     THTensor_(resizeAs)(gradInput, output);
21     TH_TENSOR_APPLY3(scalar_t, gradInput, scalar_t, gradOutput, scalar_t, output,
22         scalar_t z = *output_data;
23         *gradInput_data = *gradOutput_data * (1. - z) * z;
24     );
25 }
26
27 #endif
```

**Forward**

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

**static void sigmoid\_kernel(TensorIterator& iter) {**

**AT\_DISPATCH\_FLOATING\_TYPES(iter.dtype(), "sigmoid\_cpu", [&]() {**

**unary\_kernel\_vec(**

**iter,**

**[=](scalar\_t a) -> scalar\_t { **return (1 / (1 + std::exp((-a))));****

**[=](Vec256<scalar\_t> a) {**

**a = Vec256<scalar\_t>((scalar\_t)(0)) - a;**

**a = a.exp();**

**a = Vec256<scalar\_t>((scalar\_t)(1)) + a;**

**a = a.reciprocal();**

****return a;****

**};**

**);**

**};** **Forward actually defined elsewhere...**

****return (1 / (1 + std::exp((-a))));****

**}**

# PyTorch sigmoid layer

```
1 #ifndef TH_GENERIC_FILE
2 #define TH_GENERIC_FILE "THNN/generic/Sigmoid.c"
3 #else
```

Forward

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

```
4
5 void THNN_(Sigmoid_updateOutput)(
6     THNNState *state,
7     THTensor *input,
8     THTensor *output)
9 {
10     THTensor_(sigmoid)(output, input);
11 }
12
```

```
13 void THNN_(Sigmoid_updateGradInput)(
14     THNNState *state,
15     THTensor *gradOutput,
16     THTensor *gradInput,
17     THTensor *output)
18 {
19     THNN_CHECK_NELEMENT(output, gradOutput);
20     THTensor_(resizeAs)(gradInput, output);
21     TH_TENSOR_APPLY3(scalar_t, gradInput, scalar_t, gradOutput, scalar_t, output,
22     scalar_t z = *output_data;
23     *gradInput_data = *gradOutput_data * (1. - z) * z;
24 );
25 }
```

```
26 #endif
```

Backward

$$(1 - \sigma(x)) \sigma(x)$$

Justin Johnson

Lecture 6 - 66

Source

September 23, 2019

What about vector-valued functions?

So far: backprop with scalars

## Recap: Vector Derivatives

$$x \in \mathbb{R}, y \in \mathbb{R}$$

Regular derivative:

$$\frac{\partial y}{\partial x} \in \mathbb{R}$$

If  $x$  changes by a small amount, how much will  $y$  change?

## Recap: Vector Derivatives

$$x \in \mathbb{R}, y \in \mathbb{R}$$

$$x \in \mathbb{R}^N, y \in \mathbb{R}$$

Regular derivative:

$$\frac{\partial y}{\partial x} \in \mathbb{R} \quad \frac{\partial y}{\partial x} \in \mathbb{R}^N \quad \left( \frac{\partial y}{\partial x} \right)_n = \frac{\partial y}{\partial x_n}$$

If  $x$  changes by a small amount, how much will  $y$  change?

For each element of  $x$ , if it changes by a small amount then how much will  $y$  change?

## Recap: Vector Derivatives

$$x \in \mathbb{R}, y \in \mathbb{R}$$

$$x \in \mathbb{R}^N, y \in \mathbb{R}$$

Regular derivative:

Derivative is **Gradient**:

Derivative is **Jacobian**:

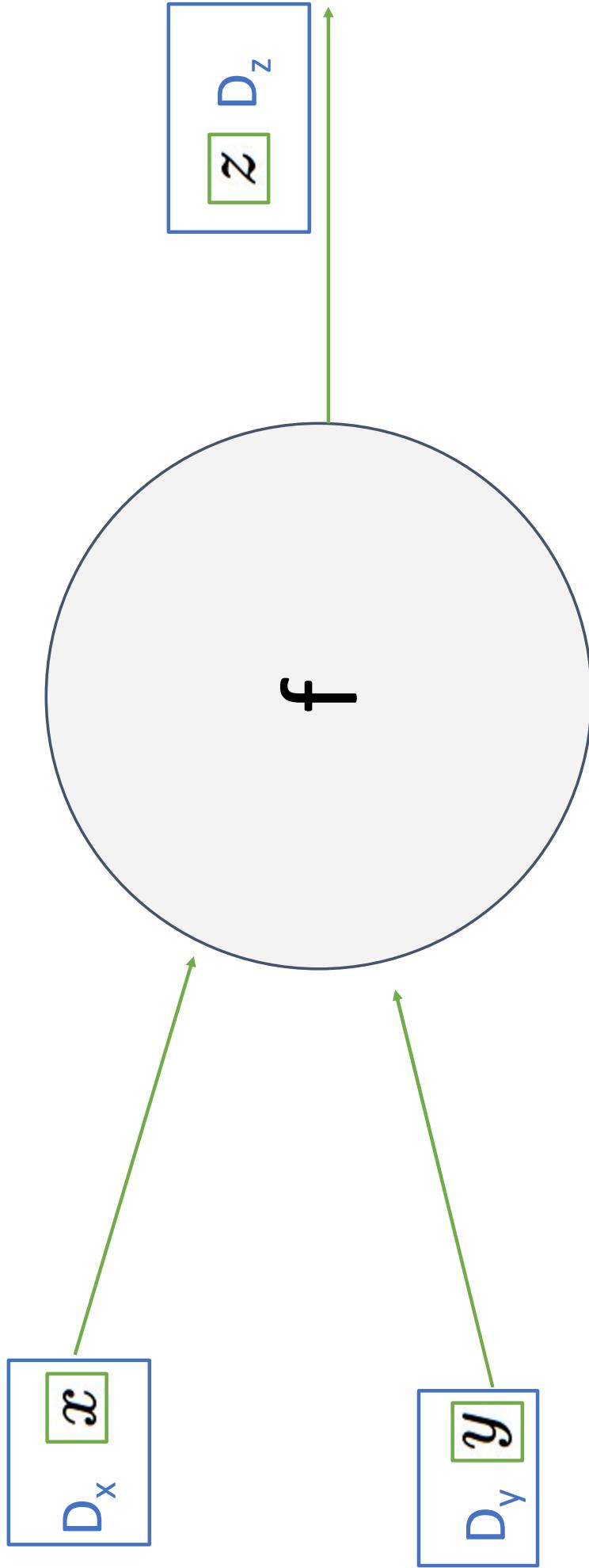
$$\frac{\partial y}{\partial x} \in \mathbb{R} \quad \frac{\partial y}{\partial x} \in \mathbb{R}^N \quad \left( \frac{\partial y}{\partial x} \right)_n = \frac{\partial y}{\partial x_n} \quad \frac{\partial y}{\partial x} \in \mathbb{R}^{N \times M} \quad \left( \frac{\partial y}{\partial x} \right)_{n,m} = \frac{\partial y_m}{\partial x_n}$$

If  $x$  changes by a small amount, how much will  $y$  change?

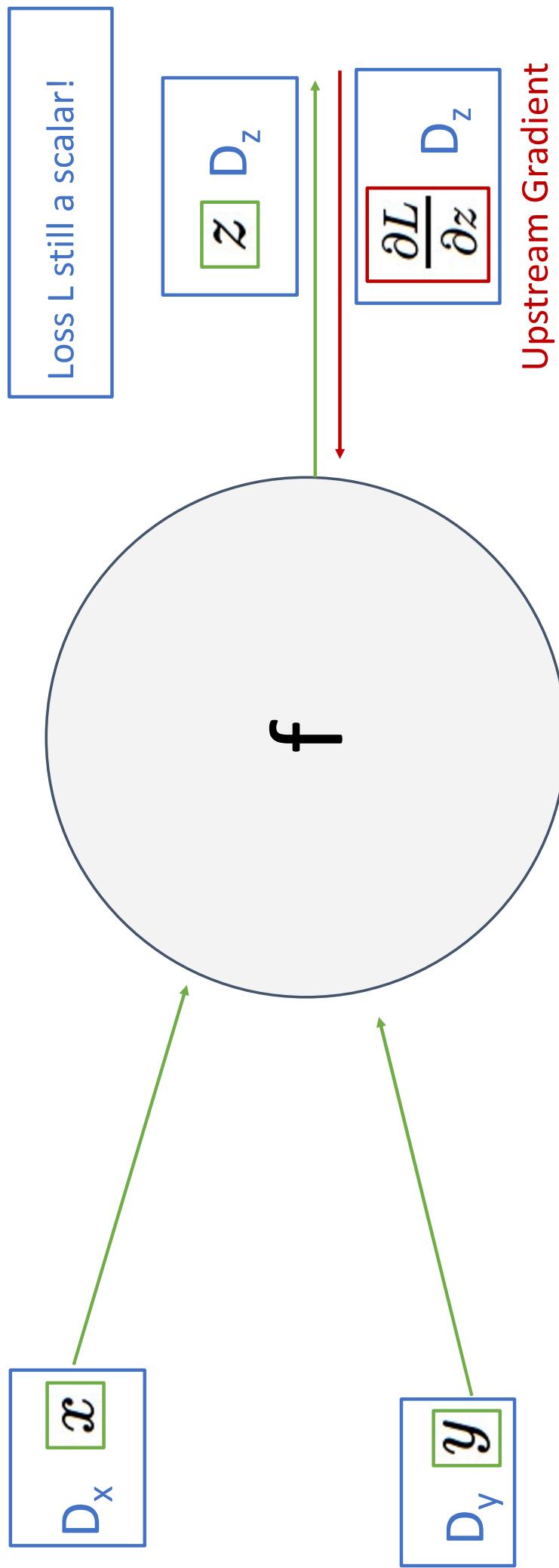
For each element of  $x$ , if it changes by a small amount then how much will  $y$  change?

For each element of  $x$ , if it changes by a small amount then how much will each element of  $y$  change?

# Backprop with Vectors

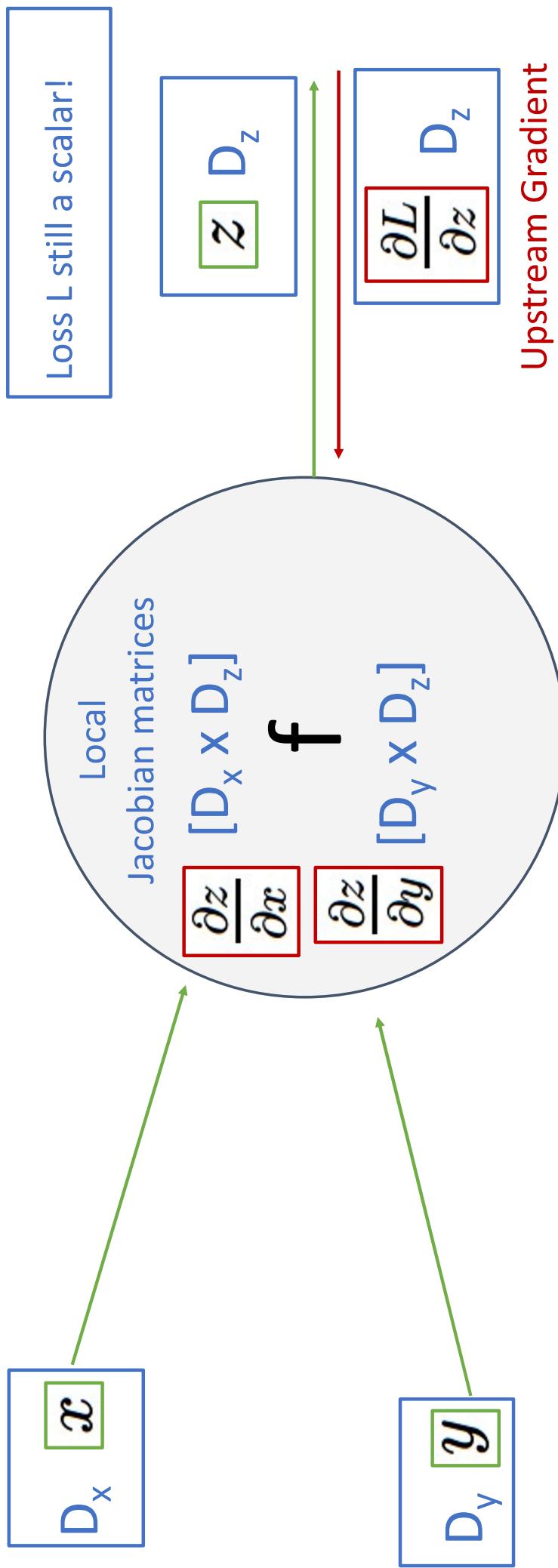


## Backprop with Vectors

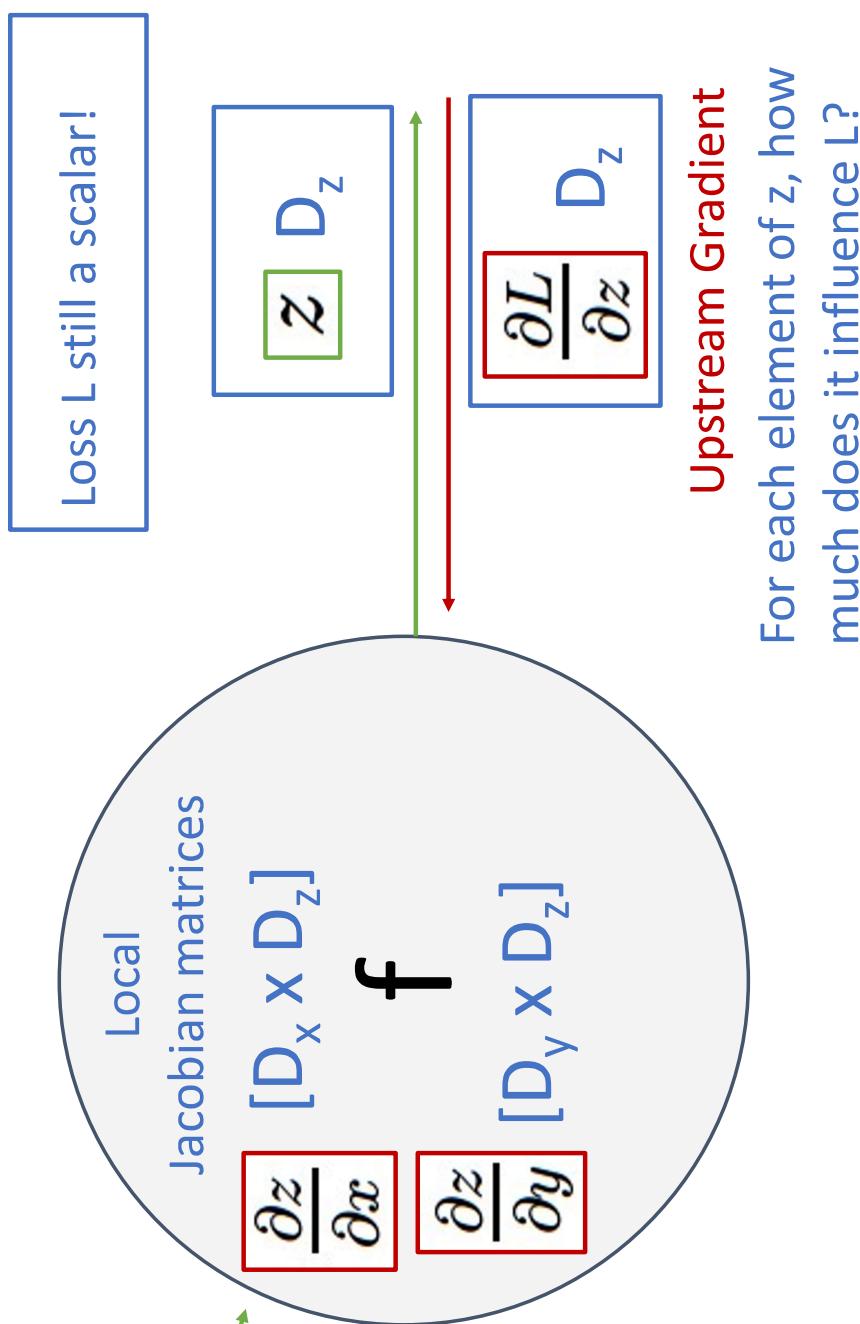
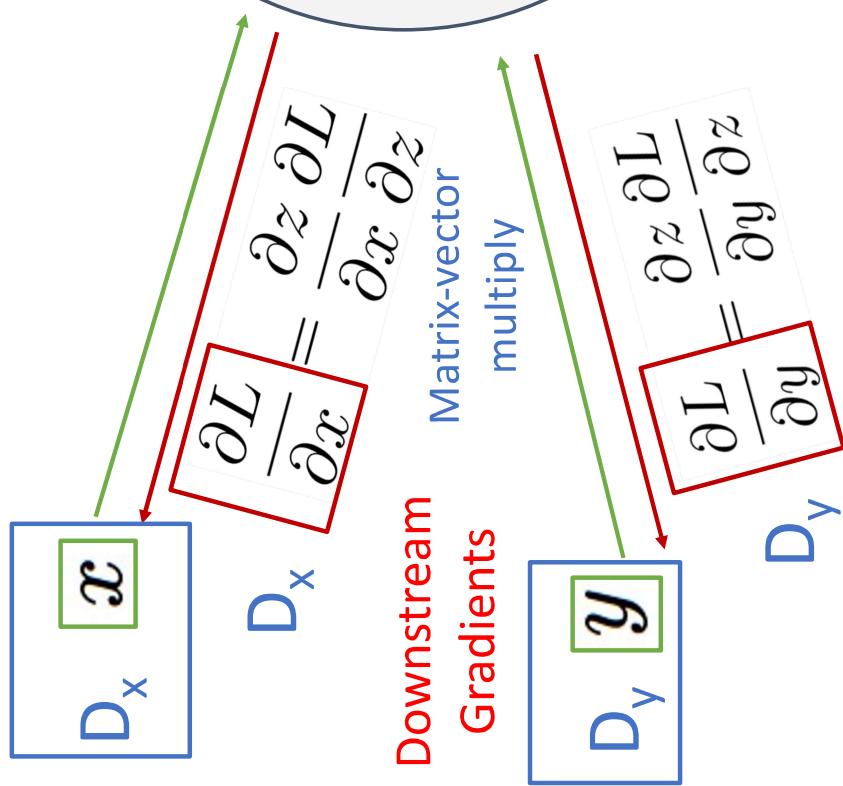


**Upstream Gradient**  
For each element of  $z$ , how much does it influence  $L$ ?

# Backprop with Vectors



# Backprop with Vectors



# Backprop with Vectors

4D input x:

$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix}$$

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

4D output y:

$$\begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$$

# Backprop with Vectors

4D input x:

$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix}$$

4D output y:

$$\begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$$

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

4D  $dL/dy$ :

$$\begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$$

Upstream  
gradient

# Backprop with Vectors

4D input x:

$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix}$$

4D output y:

$$\begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$$

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

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

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

4D  $\frac{\partial L}{\partial y}$ :

$$\begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$$

Upstream  
gradient

# Backprop with Vectors

4D input x:

$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix}$$

4D output y:

$$\begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$$

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

[dy/dx] [dL/dy]

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$$

4D dL/dy:

$$\begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$$

Upstream  
gradient

# Backprop with Vectors

4D input x:

$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix}$$

4D output y:

$$\begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$$

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

4D  $dL/dx$ :

$$\begin{bmatrix} 4 \\ 0 \\ 5 \\ 0 \end{bmatrix} \quad \begin{bmatrix} dL/dx \end{bmatrix} \quad \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad \begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$$

4D  $dL/dy$ :

$$\begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix} \quad \begin{bmatrix} dL/dy \end{bmatrix} \quad \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad \begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$$

Upstream  
gradient

# Backprop with Vectors

Jacobian is **sparse**: off-diagonal entries all zero! Never **explicitly** form Jacobian; instead use **implicit** multiplication

4D input x:

$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix}$$

4D output y:

$$\begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$$

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

(elementwise)

4D  $dL/dx$ :

$$\begin{bmatrix} 4 \\ 0 \\ 5 \\ 0 \end{bmatrix}$$

$[dy/dx]$  [ $dL/dy$ ]

4D  $dL/dy$ :

$$\begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$$

Upstream

gradient

# Backprop with Vectors

Jacobian is **sparse**: off-diagonal entries all zero! Never **explicitly** form Jacobian; instead use **implicit** multiplication

4D input x:

$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix}$$

4D output y:

$$\begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$$

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

(elementwise)

4D  $dL/dx$ :

$$\begin{bmatrix} 4 \\ 0 \\ 5 \\ 0 \end{bmatrix}$$

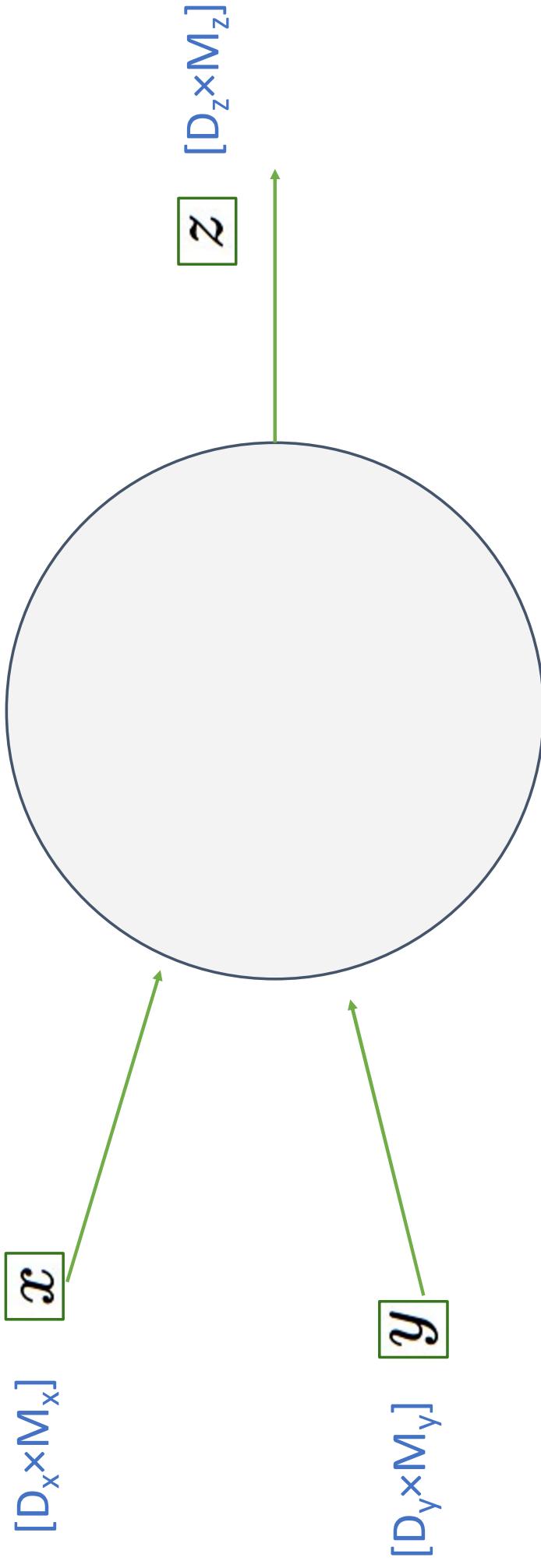
[ $dy/dx$ ] [ $dL/dy$ ]

4D  $dL/dy$ :

$$\begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$$

Upstream  
gradient

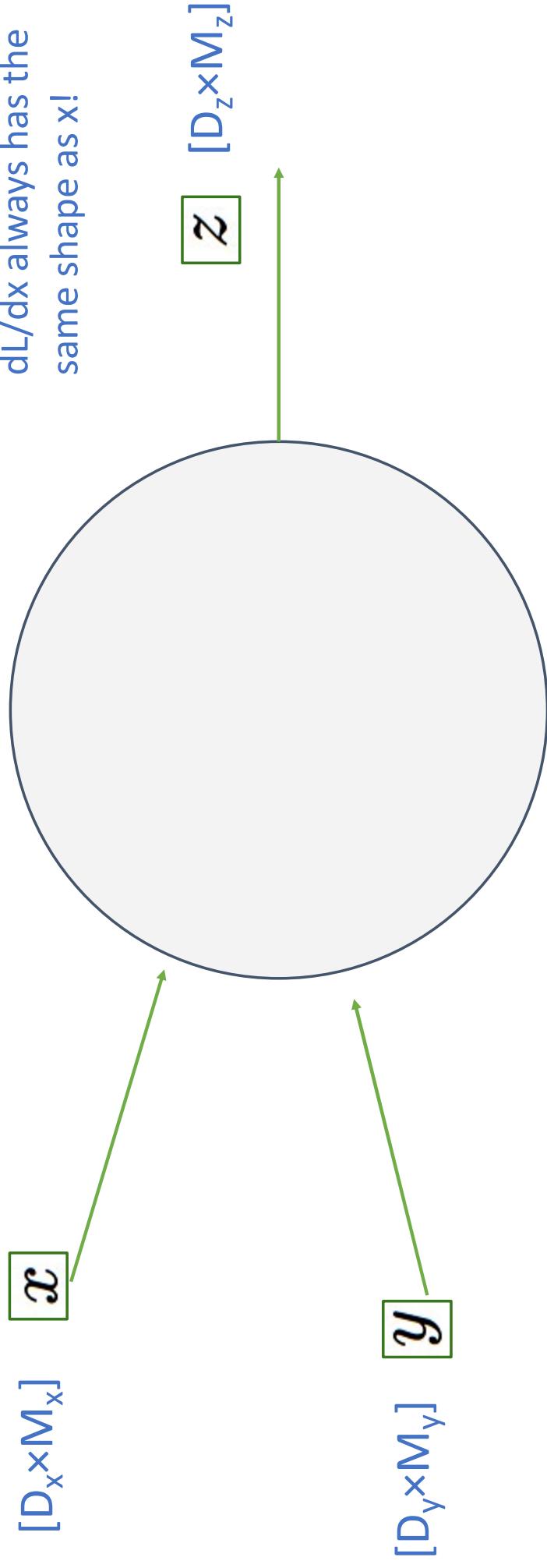
Backprop with Matrices (or Tensors):



## Backprop with Matrices (or Tensors):

Loss L still a scalar!

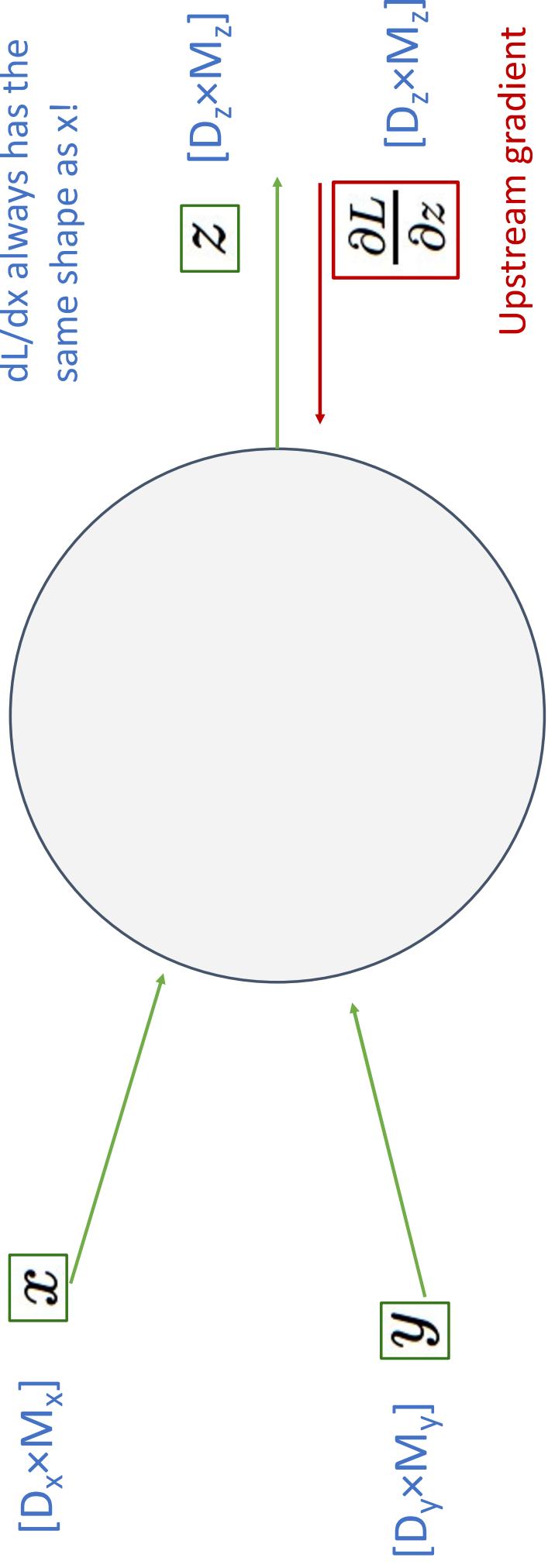
$dL/dx$  always has the same shape as  $x$ !



## Backprop with Matrices (or Tensors):

Loss L still a scalar!

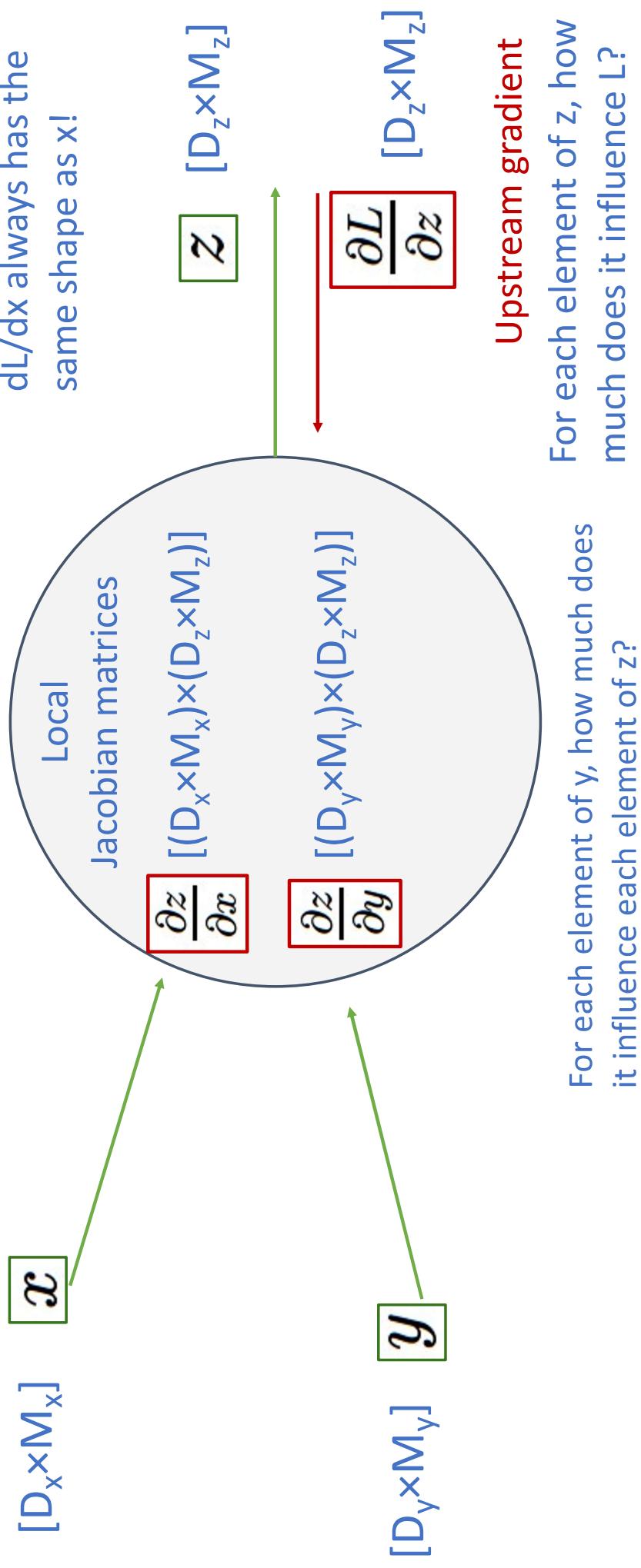
$dL/dx$  always has the same shape as x!



Upstream gradient  
For each element of z, how much does it influence L?

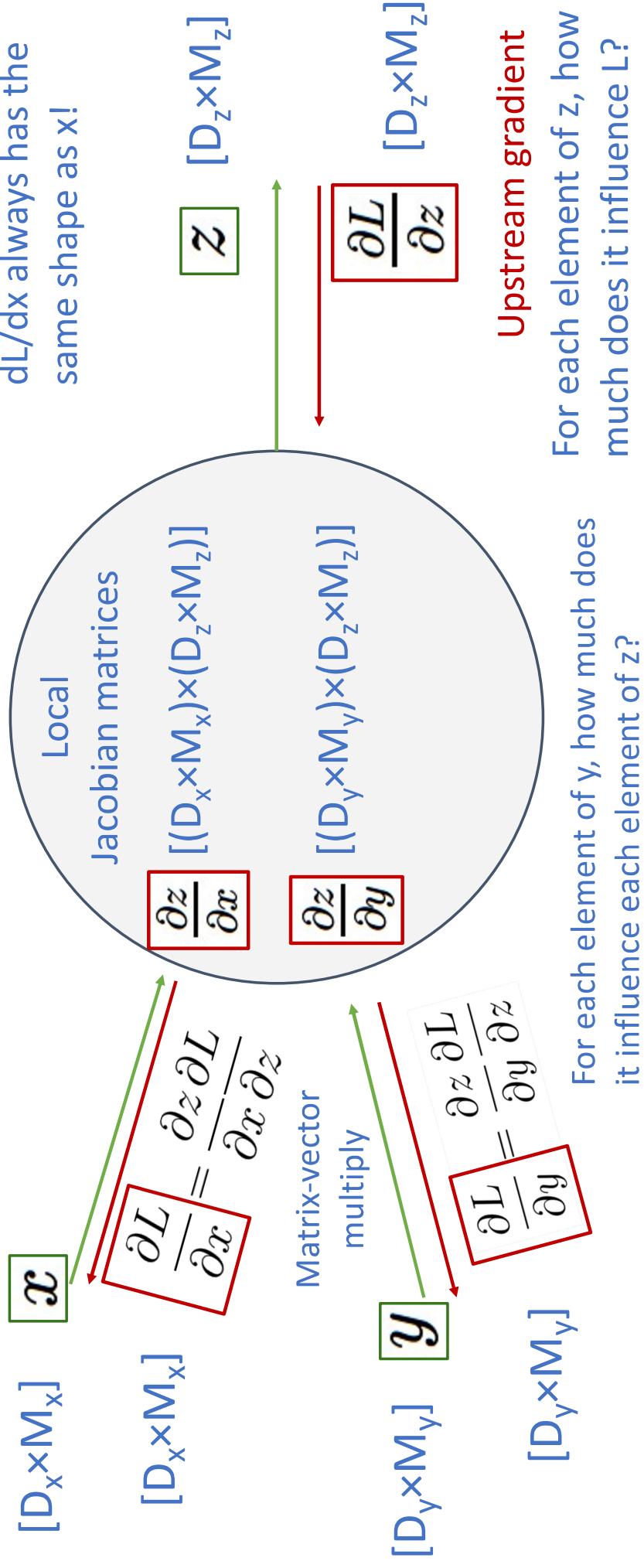
## Backprop with Matrices (or Tensors):

Loss L still a scalar!



## Backprop with Matrices (or Tensors):

Loss L still a scalar!



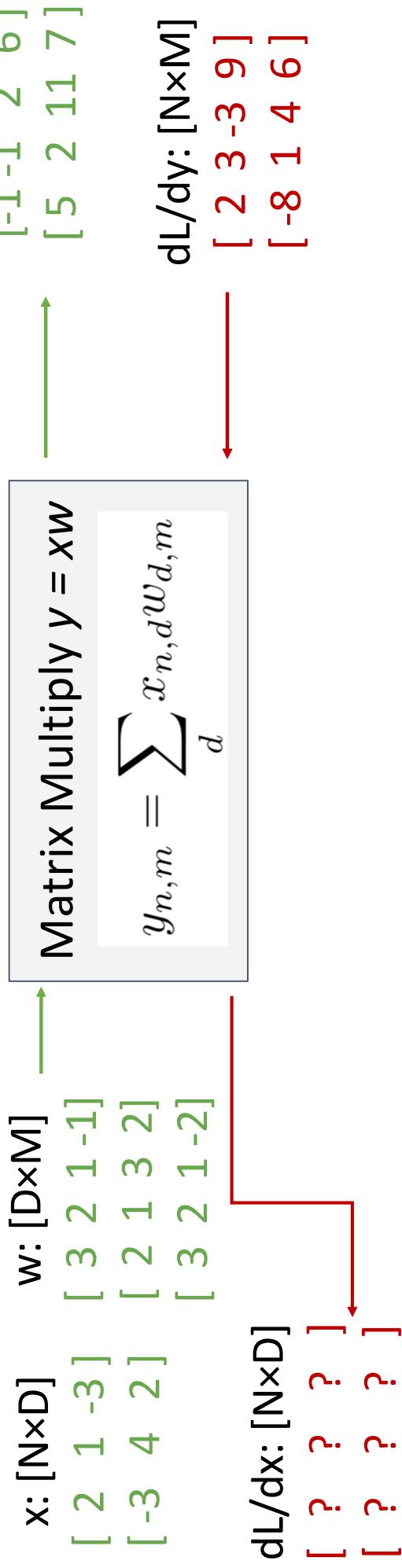
## Example: Matrix Multiplication

$$x: [N \times D] \quad w: [D \times M]$$
$$\begin{bmatrix} 2 & 1 & -3 \\ -3 & 4 & 2 \end{bmatrix} \quad \begin{bmatrix} 3 & 2 & 1 & -1 \\ 2 & 1 & 3 & 2 \\ 3 & 2 & 1 & -2 \end{bmatrix}$$

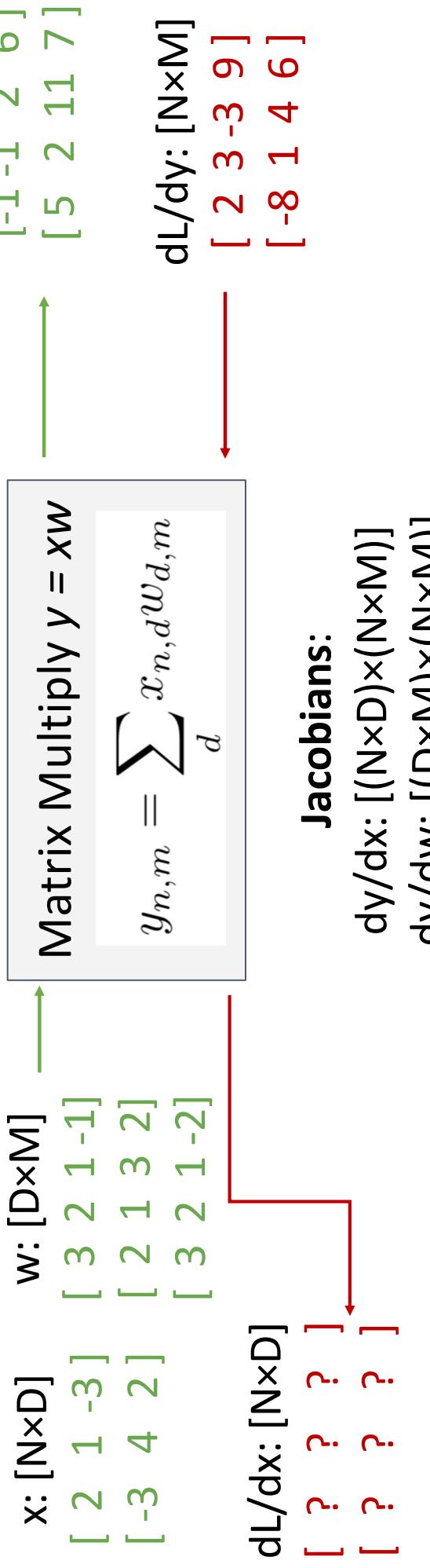
$$\boxed{\text{Matrix Multiply } y = xw}$$
$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

$$y: [N \times M]$$
$$\begin{bmatrix} -1 & -1 & 2 & 6 \\ 5 & 2 & 11 & 7 \end{bmatrix}$$

## Example: Matrix Multiplication

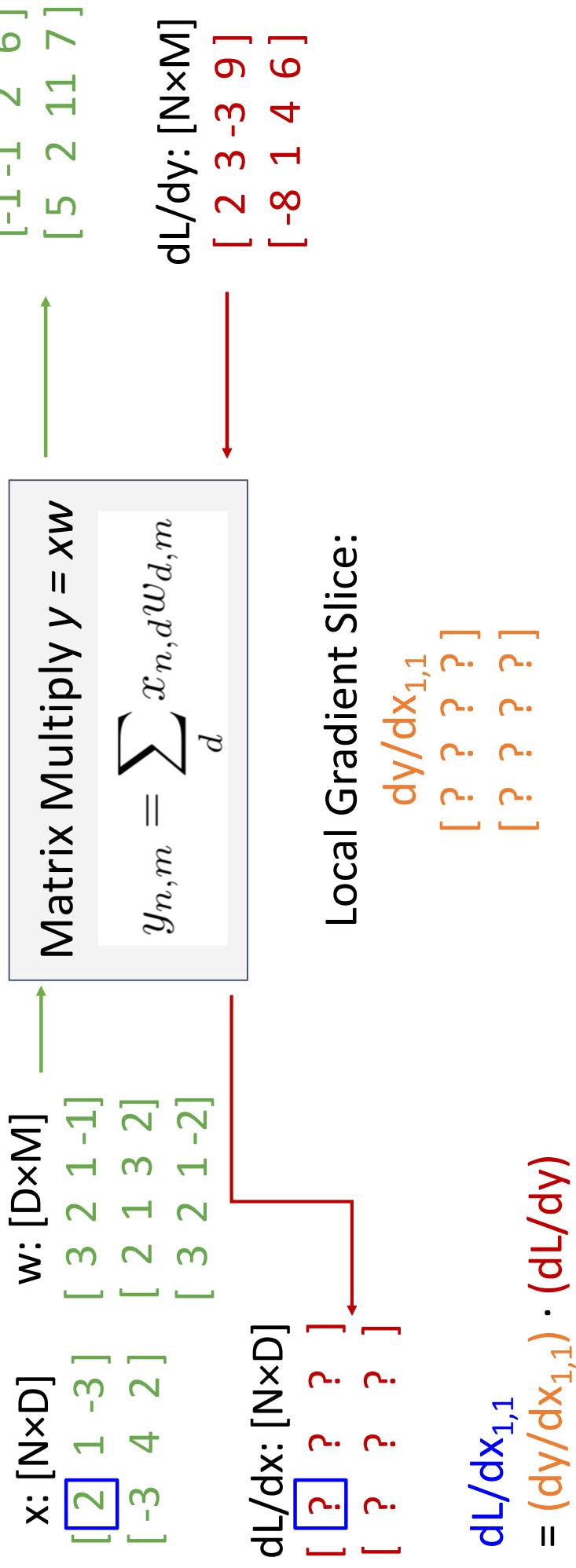


## Example: Matrix Multiplication

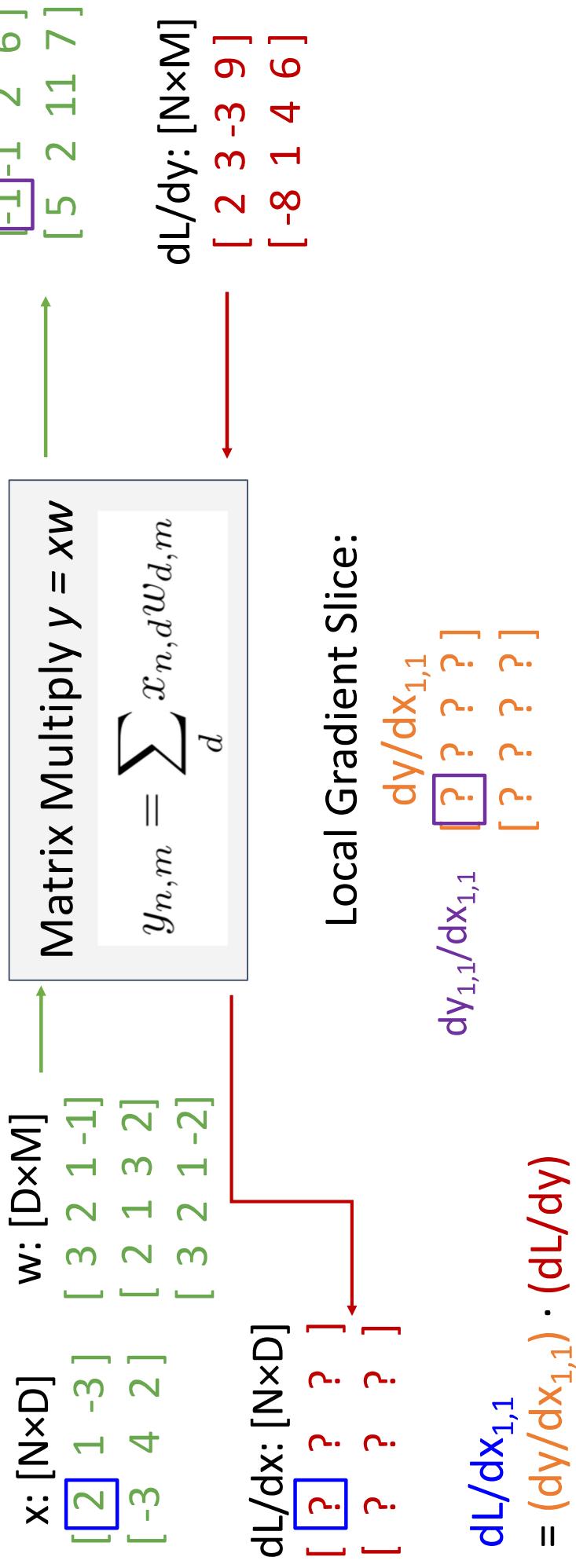


For a neural net we may have  
 $N=64, D=M=4096$   
 Each Jacobian takes 256 GB of memory! Must work with them implicitly!

## Example: Matrix Multiplication



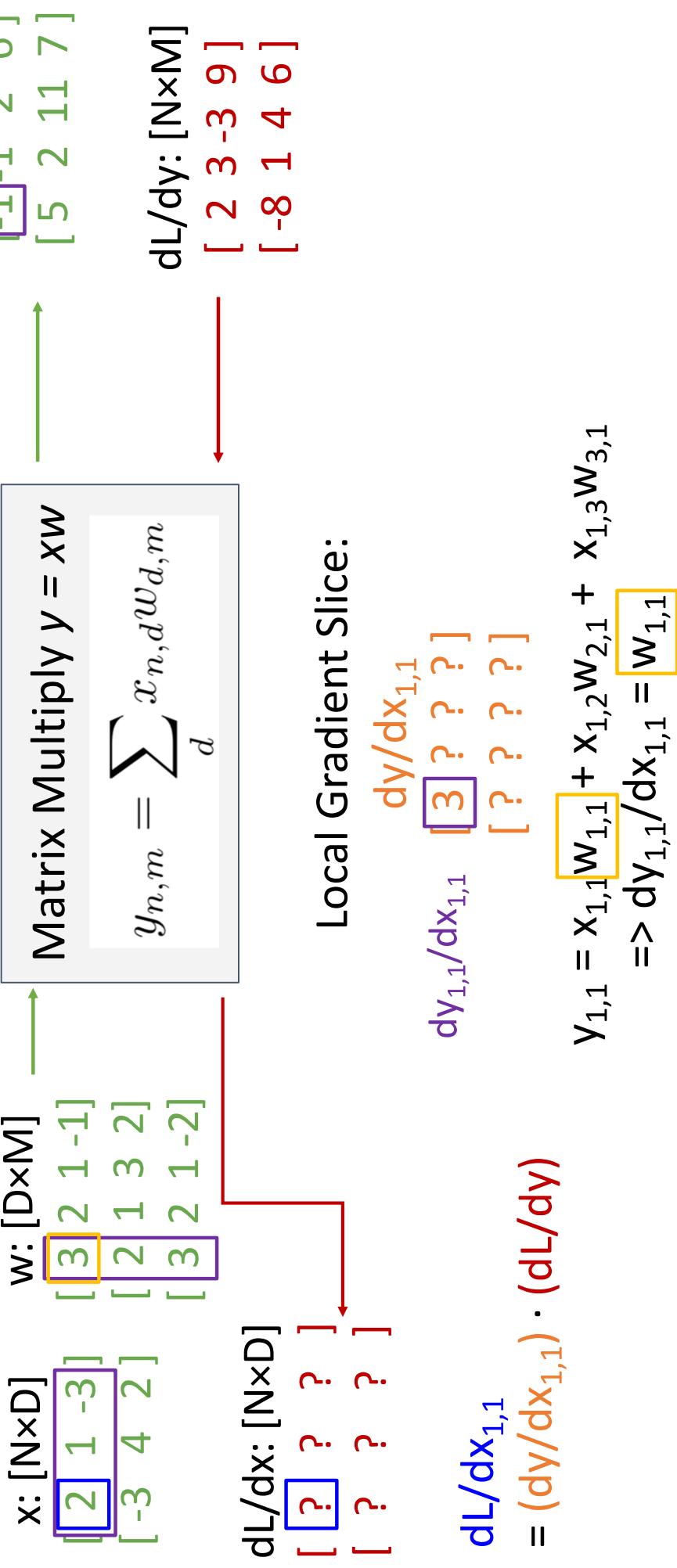
## Example: Matrix Multiplication



## Example: Matrix Multiplication



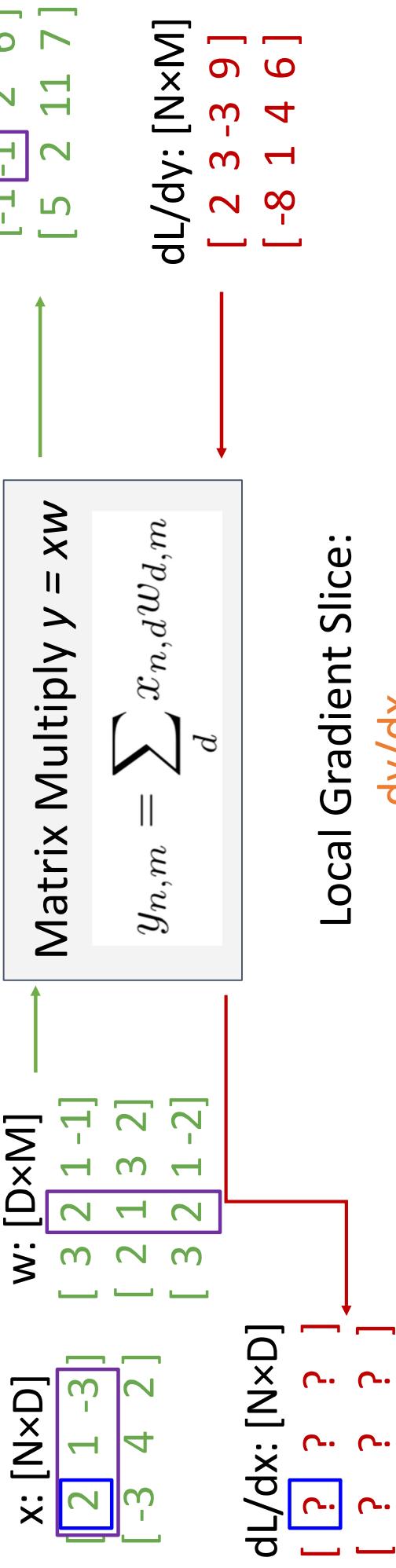
## Example: Matrix Multiplication



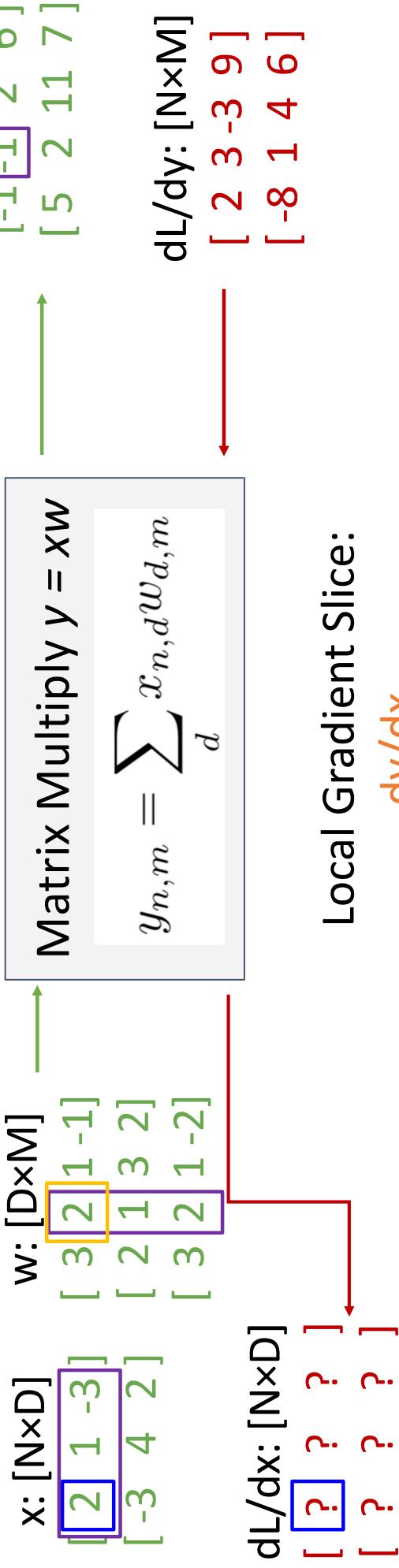
## Example: Matrix Multiplication



## Example: Matrix Multiplication

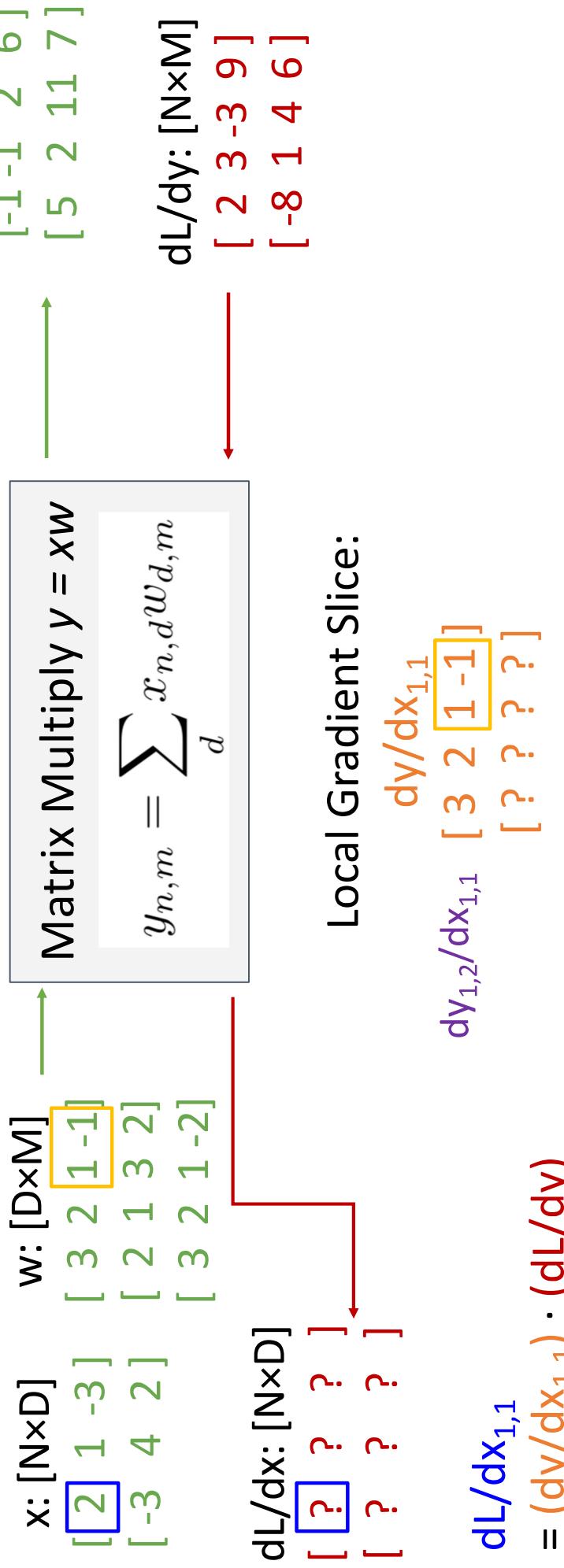


## Example: Matrix Multiplication

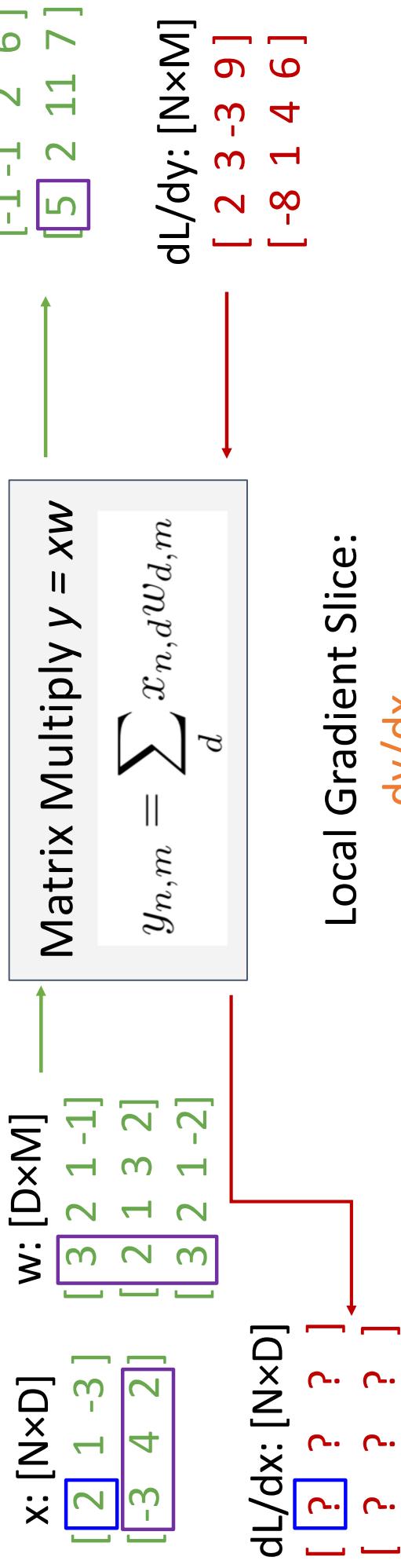


$$\begin{aligned}
 & dy/dx_{1,1} = \frac{\partial y_{1,2}}{\partial x_{1,1}} = \frac{\partial}{\partial x_{1,1}} \left[ \sum_d x_{n,d} w_{d,m} \right] \\
 & = \left[ \begin{array}{c} 3 \\ 2 \\ ? \end{array} \right] \cdot \left[ \begin{array}{c} 2 \\ 3 \\ -3 \end{array} \right] = \left[ \begin{array}{c} -3 \\ 6 \\ -8 \end{array} \right]
 \end{aligned}$$

## Example: Matrix Multiplication



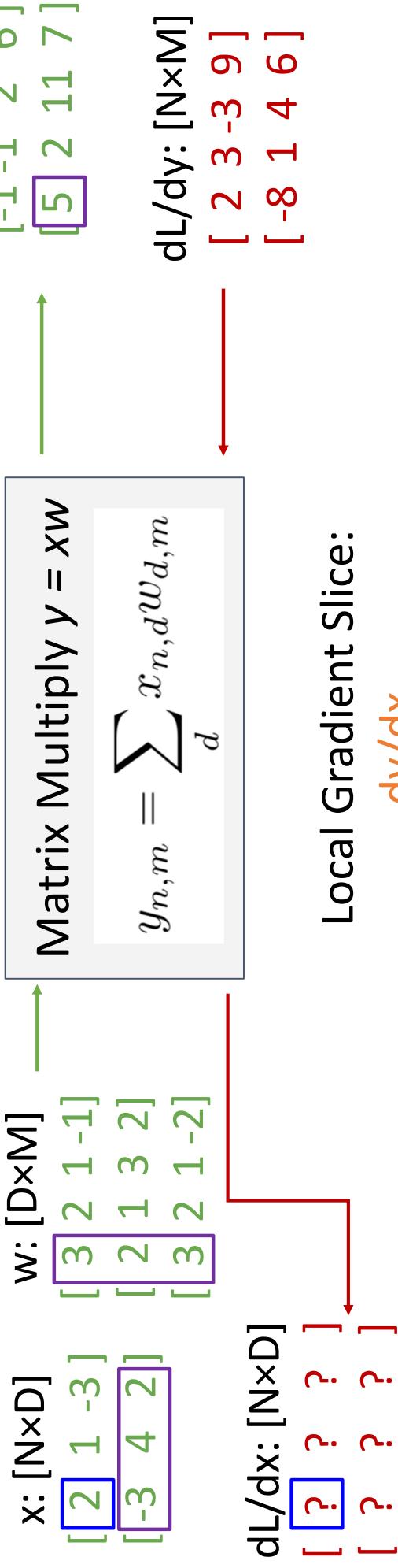
## Example: Matrix Multiplication



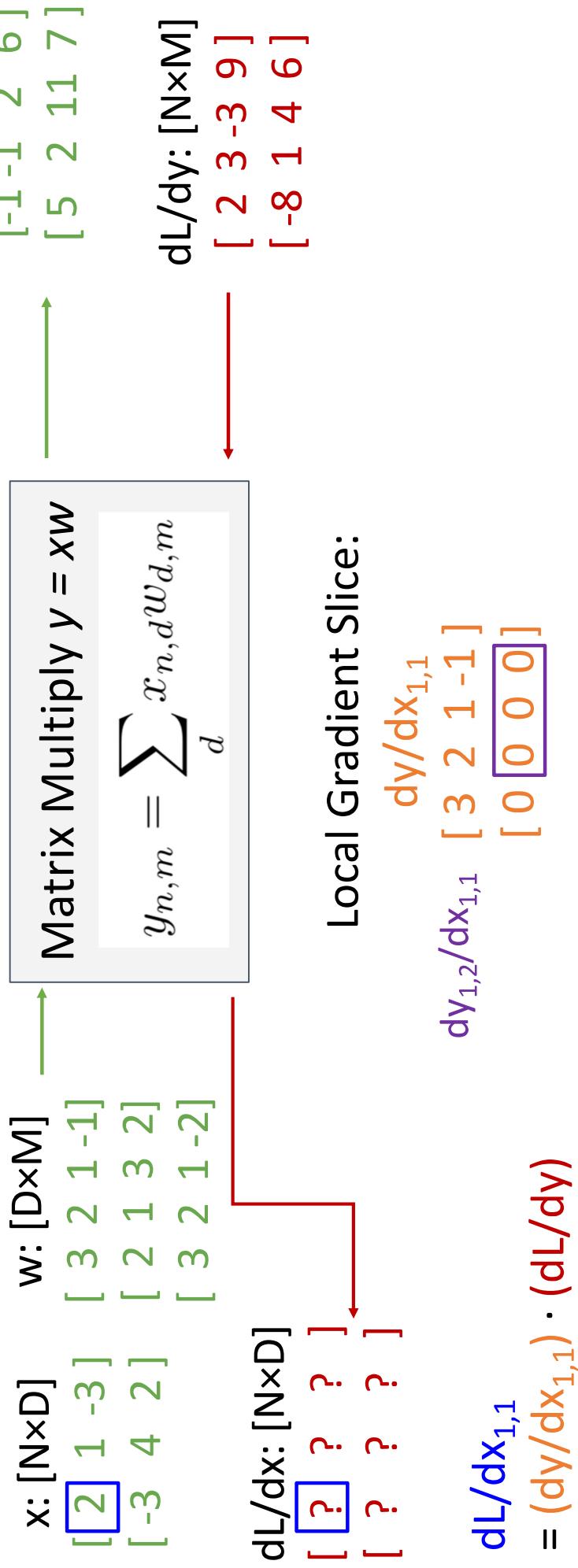
$$dL/dx_{1,1} = (dy/dx_{1,1}) \cdot (dL/dy)$$

$$y_{2,1} = x_{2,1}w_{1,1} + x_{2,2}w_{2,1} + x_{2,3}w_{3,1}$$

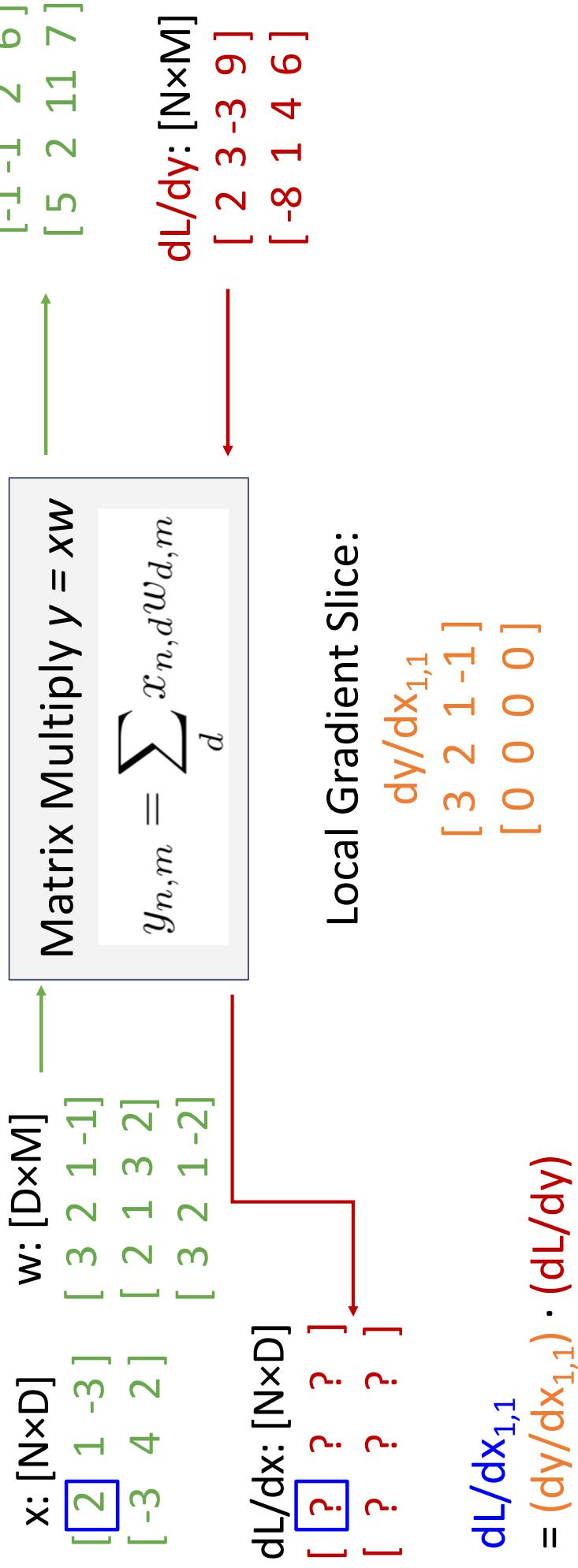
## Example: Matrix Multiplication



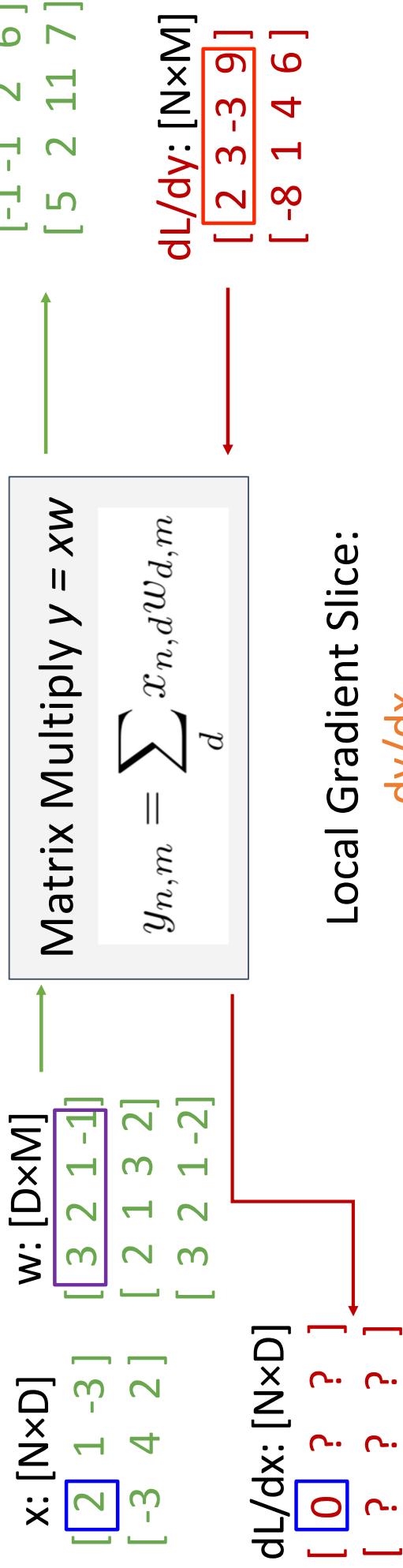
## Example: Matrix Multiplication



## Example: Matrix Multiplication



## Example: Matrix Multiplication



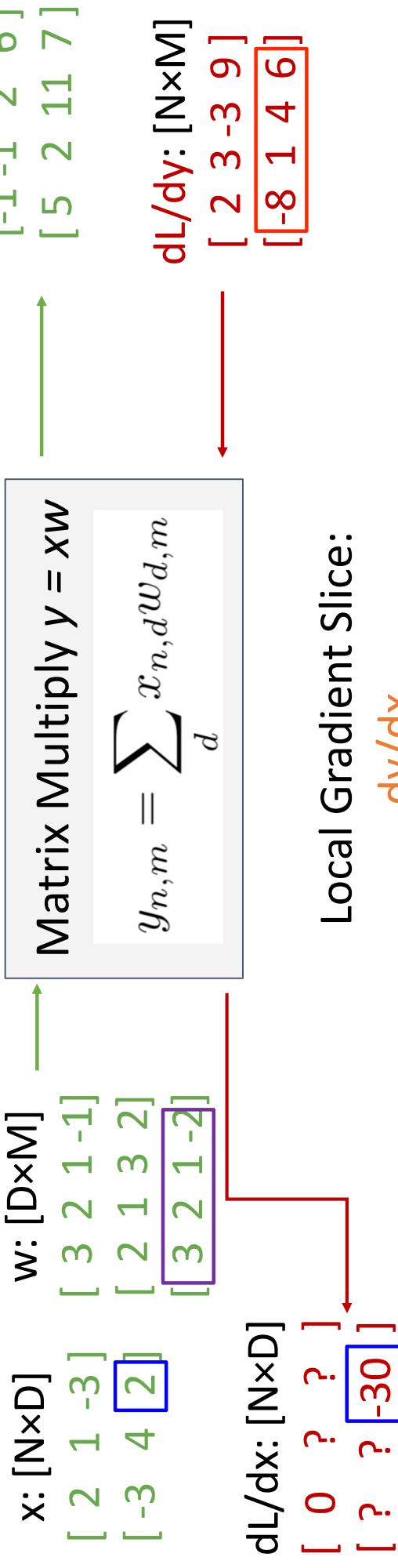
$$\begin{aligned}
 dL/dx_{1,1} &= (\frac{dy}{dx_{1,1}}) \cdot (\frac{dL}{dy}) \\
 &= (w_{1,:}) \cdot (\frac{dL}{dy_{1,:}}) \\
 &= 3*2 + 2*3 + 1*(-3) + (-1)*9 = 0
 \end{aligned}$$

## Example: Matrix Multiplication



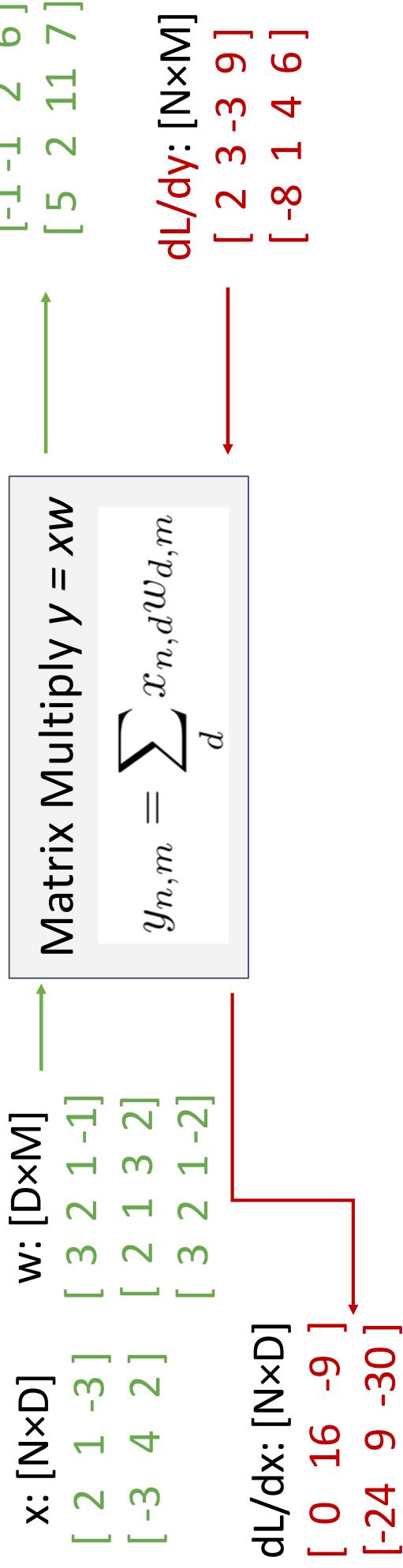
$$dL/dx_{2,3} = (dy/dx_{2,3}) \cdot (dL/dy)$$

## Example: Matrix Multiplication



$$\begin{aligned}
 dL/dx_{2,3} &= (\mathbf{dy}/d\mathbf{x}_{2,3}) \cdot (\mathbf{dL}/dy) \\
 &= (\mathbf{w}_{3,:}) \cdot (\mathbf{dL}/dy_{2,:}) \\
 &= 3 * (-8) + 2 * 1 + 1 * 4 + (-2) * 6 = -30
 \end{aligned}$$

## Example: Matrix Multiplication



$$\begin{aligned} dL/dx_{ij} &= (dy/dx_{ij}) \cdot (dL/dy) \\ &= (w_{j,:}) \cdot (dL/dy_{ij}) \end{aligned}$$

## Example: Matrix Multiplication

$$\begin{array}{l}
 \text{x: [NxD]} \quad \text{w: [DxM]} \quad \text{y: [NxM]} \\
 \left[ \begin{matrix} 2 & 1 & -3 \\ -3 & 4 & 2 \end{matrix} \right] \quad \left[ \begin{matrix} 3 & 2 & 1 & -1 \\ 2 & 1 & 3 & 2 \\ 3 & 2 & 1 & -2 \end{matrix} \right] \quad \left[ \begin{matrix} -1 & -1 & 2 & 6 \\ 5 & 2 & 11 & 7 \end{matrix} \right]
 \end{array}$$

Matrix Multiply  $y = xw$

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

$$\begin{array}{l}
 \text{dL/dx: [NxD]} \quad \text{dL/dy: [NxM]} \\
 \left[ \begin{matrix} 0 & 16 & -9 \\ -24 & 9 & -30 \end{matrix} \right] \quad \left[ \begin{matrix} 2 & 3 & -3 & 9 \\ -8 & 1 & 4 & 6 \end{matrix} \right]
 \end{array}$$

$$\text{dL/dx} = (\text{dL/dy}) w^T$$

$$\begin{array}{l}
 [\text{NxD}] \quad [\text{N x M}] \quad [\text{M x D}] \\
 \text{Easy way to remember:} \\
 \text{It's the only way the shapes work out!} \\
 \text{dL/dx}_{ij} \\
 = (\text{dy/dx}_{ij}) \cdot (\text{dL/dy}) \\
 = (w_{j,:}) \cdot (\text{dL/dy}_{i,:})
 \end{array}$$

## Example: Matrix Multiplication

$$\begin{array}{l}
 \text{x: [NxD]} \quad \text{w: [DxM]} \quad \text{y: [NxM]} \\
 \left[ \begin{matrix} 2 & 1 & -3 \\ -3 & 4 & 2 \end{matrix} \right] \quad \left[ \begin{matrix} 3 & 2 & 1 & -1 \\ 2 & 1 & 3 & 2 \\ 3 & 2 & 1 & -2 \end{matrix} \right] \quad \left[ \begin{matrix} -1 & -1 & 2 & 6 \\ 5 & 2 & 11 & 7 \end{matrix} \right]
 \end{array}$$

Matrix Multiply  $y = xw$

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

$$\begin{array}{l}
 \text{dL/dx: [NxD]} \\
 \left[ \begin{matrix} 0 & 16 & -9 \\ -24 & 9 & -30 \end{matrix} \right]
 \end{array}
 \quad
 \begin{array}{l}
 \text{dL/dy: [NxM]} \\
 \left[ \begin{matrix} 2 & 3 & -3 & 9 \\ -8 & 1 & 4 & 6 \end{matrix} \right]
 \end{array}$$

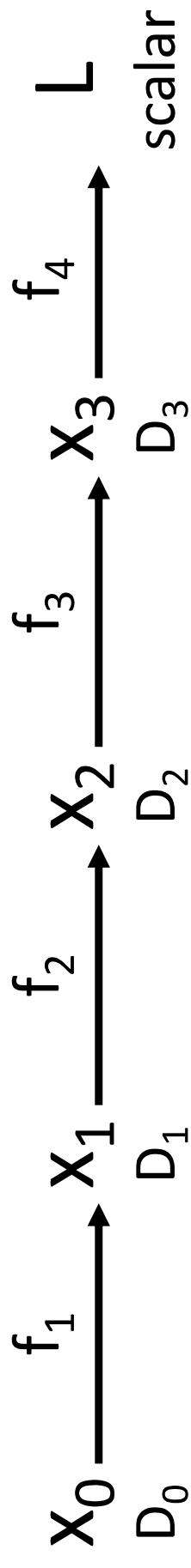
$$\text{dL/dw} = x^T (\text{dL/dy})$$

$$[\text{DxD}] \quad [\text{NxM}] \quad [\text{MxD}]$$

Easy way to remember:  
It's the only way the shapes work out!

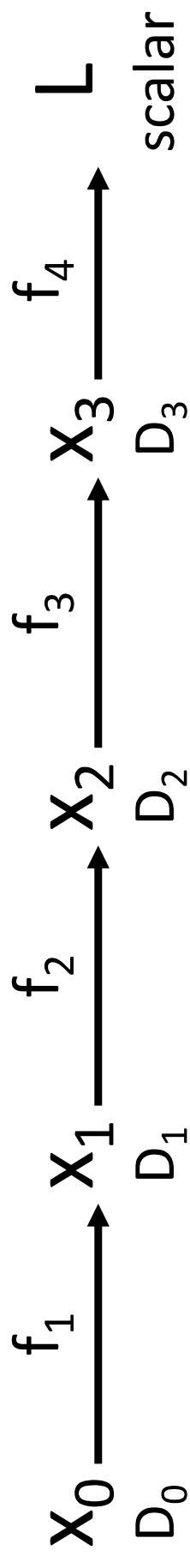
$$\begin{array}{l}
 \text{dL/dx} = (\text{dL/dy}) w^T \\
 [\text{NxD}] \quad [\text{NxM}] \quad [\text{MxD}]
 \end{array}$$

## Backpropagation: Another View



$$\text{Chain rule} \quad \frac{\partial L}{\partial x_0} = \left( \frac{\partial x_1}{\partial x_0} \right) \left( \frac{\partial x_2}{\partial x_1} \right) \left( \frac{\partial x_3}{\partial x_2} \right) \left( \frac{\partial L}{\partial x_3} \right)$$

## Backpropagation: Another View

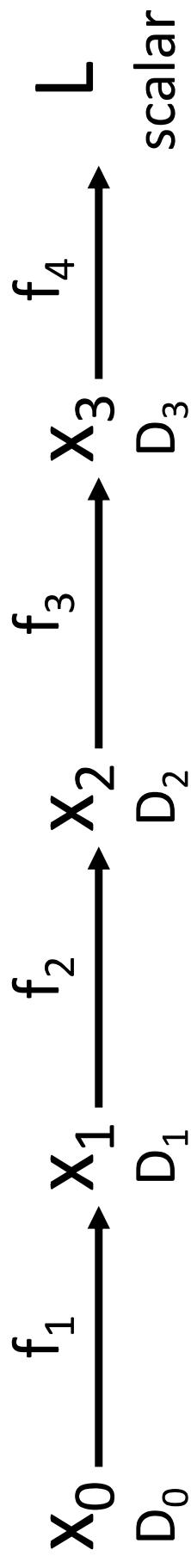


Matrix multiplication is **associative**: we can compute products in any order

Chain rule

$$\frac{\partial L}{\partial x_0} = \left( \frac{\partial x_1}{\partial x_0} \right) \left( \frac{\partial x_2}{\partial x_1} \right) \left( \frac{\partial x_3}{\partial x_2} \right) \left( \frac{\partial L}{\partial x_3} \right)$$
$$D_0 \times D_1 \quad D_1 \times D_2 \quad D_2 \times D_3 \quad D_3$$

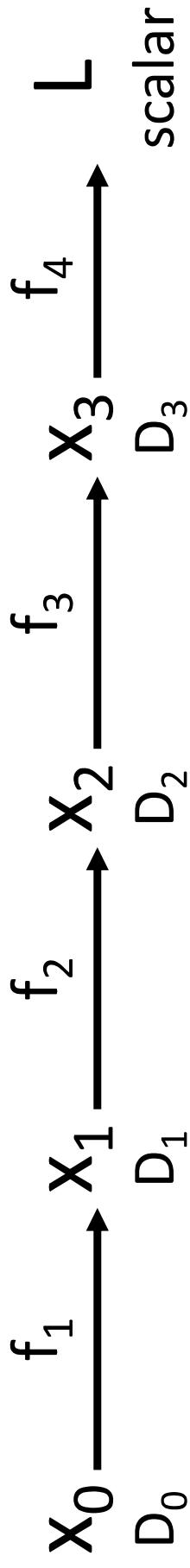
## Reverse-Mode Automatic Differentiation



Matrix multiplication is **associative**: we can compute products in any order  
Computing products right-to-left avoids matrix-matrix products; only needs matrix-vector

$$\text{Chain rule} \quad \frac{\partial L}{\partial x_0} = \left( \frac{\partial x_1}{\partial x_0} \right) \left( \frac{\partial x_2}{\partial x_1} \right) \left( \frac{\partial x_3}{\partial x_2} \right) \left( \frac{\partial L}{\partial x_3} \right)$$
$$D_0 \times D_1 \quad D_1 \times D_2 \quad D_2 \times D_3 \quad D_3$$

# Reverse-Mode Automatic Differentiation



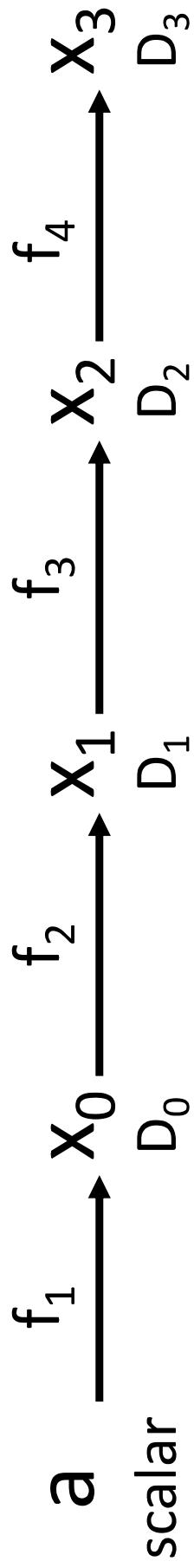
Matrix multiplication is **associative**: we can compute products in any order  
Computing products right-to-left avoids matrix-matrix products; only needs matrix-vector



$$\text{Chain rule} \quad \frac{\partial L}{\partial x_0} = \left( \frac{\partial x_1}{\partial x_0} \right) \left( \frac{\partial x_2}{\partial x_1} \right) \left( \frac{\partial x_3}{\partial x_2} \right) \left( \frac{\partial L}{\partial x_3} \right)$$

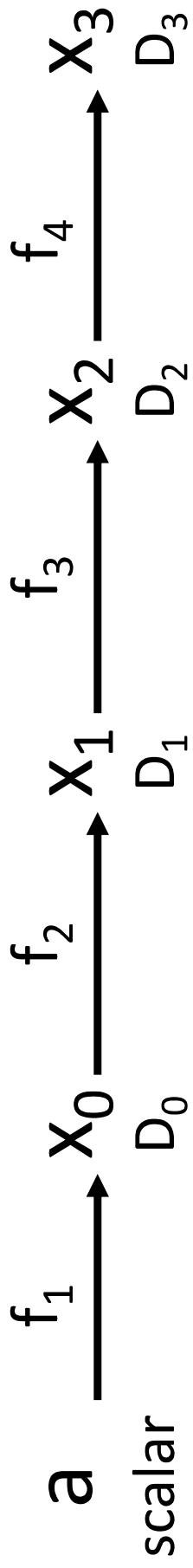
Compute grad of scalar output w/respect to all vector inputs  
 $D_0 \times D_1 \quad D_1 \times D_2 \quad D_2 \times D_3 \quad D_3$

## Forward-Mode Automatic Differentiation



$$\text{Chain rule} \quad \frac{\partial x_3}{\partial a} = \left( \frac{\partial x_0}{\partial a} \right) \left( \frac{\partial x_1}{\partial x_0} \right) \left( \frac{\partial x_2}{\partial x_1} \right) \left( \frac{\partial x_3}{\partial x_2} \right)$$
$$D_0 \quad D_0 \times D_1 \quad D_1 \times D_2 \quad D_2 \times D_3$$

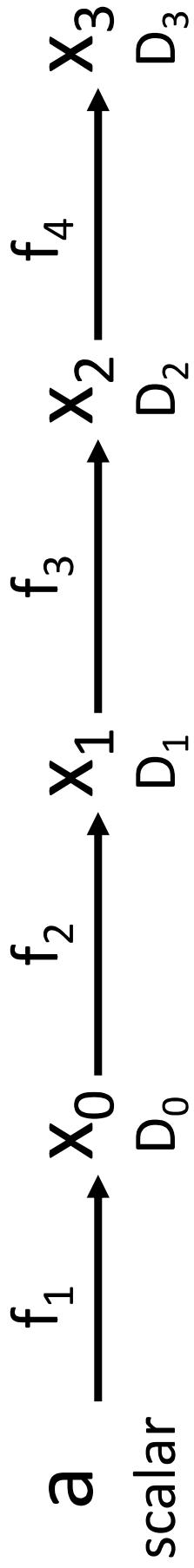
## Forward-Mode Automatic Differentiation



Computing products left-to-right avoids matrix-matrix products; only needs matrix-vector

$$\text{Chain rule} \quad \frac{\partial x_3}{\partial a} = \left( \frac{\partial x_0}{\partial a} \right) \left( \frac{\partial x_1}{\partial x_0} \right) \left( \frac{\partial x_2}{\partial x_1} \right) \left( \frac{\partial x_3}{\partial x_2} \right)$$
$$D_0 \quad D_0 \times D_1 \quad D_1 \times D_2 \quad D_2 \times D_3$$

## Forward-Mode Automatic Differentiation



Computing products left-to-right avoids matrix-matrix products; only needs matrix-vector

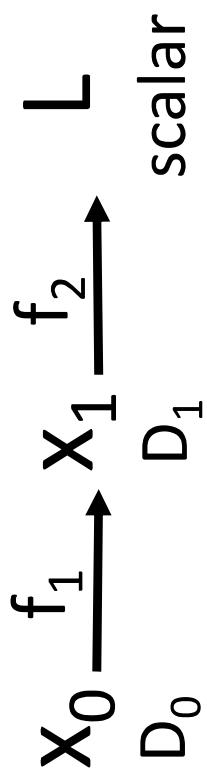
Not implemented in PyTorch / TensorFlow = (

$$\text{Chain rule} \quad \frac{\partial x_3}{\partial a} = \left( \frac{\partial x_0}{\partial a} \right) \left( \frac{\partial x_1}{\partial x_0} \right) \left( \frac{\partial x_2}{\partial x_1} \right) \left( \frac{\partial x_3}{\partial x_2} \right)$$

D<sub>0</sub>      D<sub>0</sub> × D<sub>1</sub>      D<sub>1</sub> × D<sub>2</sub>      D<sub>2</sub> × D<sub>3</sub>

But you can implement forward-mode AD using two calls to reverse-mode AD!  
(Inefficient but elegant)

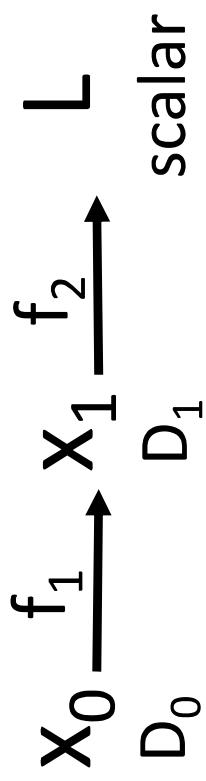
## Backprop: Higher-Order Derivatives



$\frac{\partial^2 L}{\partial x_0^2}$  Hessian matrix  
H of second  
derivatives.

$$D_0 \times D_0$$

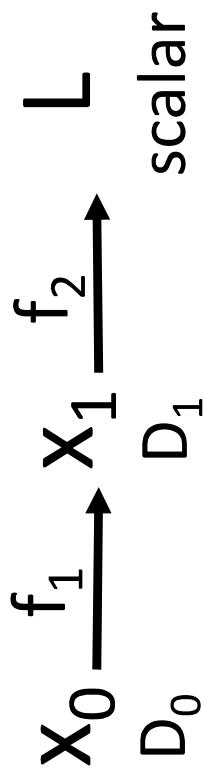
## Backprop: Higher-Order Derivatives



Hessian / vector multiply

$$\frac{\partial^2 L}{\partial x_0^2} \quad \begin{matrix} \text{Hessian matrix} \\ \text{H of second} \\ \text{derivatives.} \end{matrix} \quad \frac{\partial^2 L}{\partial x_0^2} v$$
$$D_0 \times D_0 \quad D_0$$

## Backprop: Higher-Order Derivatives



Hessian / vector multiply

$$\frac{\partial^2 L}{\partial x_0^2} v = \frac{\partial}{\partial x_0} \left[ \frac{\partial L}{\partial x_0} \cdot v \right]$$

(if  $v$  doesn't depend on  $x_0$ )

$D_0 \times D_0 \quad D_0$

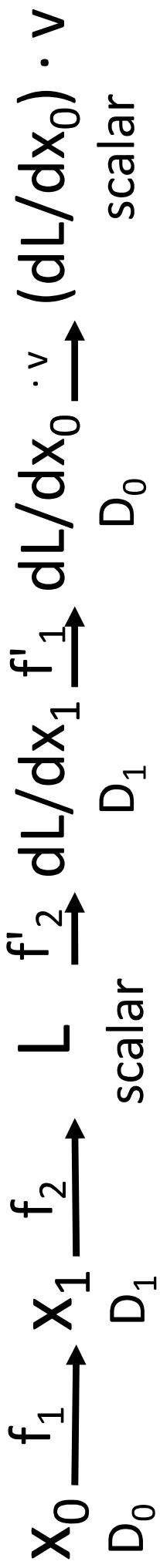
## Backprop: Higher-Order Derivatives

$$\begin{matrix} x_0 & \xrightarrow{f_1} & x_1 & \xrightarrow{f_2} & L & \xrightarrow{f'_2} dL/dx_1 & \xrightarrow{f'_1} dL/dx_0 \cdot v \\ & & D_1 & & \text{scalar} & & D_1 \\ & & & & & & D_0 \\ & & & & & & \text{scalar} \end{matrix}$$

Hessian / vector multiply

$$\frac{\partial^2 L}{\partial x_0^2} v = \frac{\partial}{\partial x_0} \left[ \frac{\partial L}{\partial x_0} \cdot v \right] \quad (\text{if } v \text{ doesn't depend on } x_0)$$
$$D_0 \times D_0 \quad D_0$$

## Backprop: Higher-Order Derivatives



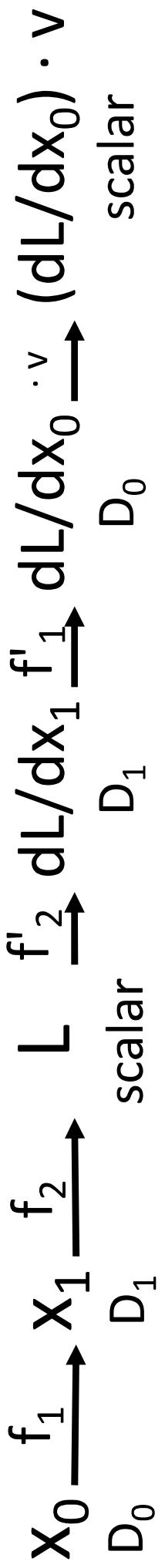
Backprop!

### Hessian / vector multiply

$$\frac{\partial^2 L}{\partial x_0^2} v = \frac{\partial}{\partial x_0} \left[ \frac{\partial L}{\partial x_0} \cdot v \right] \quad (\text{if } v \text{ doesn't depend on } x_0)$$

$D_0 \times D_0 \quad D_0$

## Backprop: Higher-Order Derivatives



Backprop!

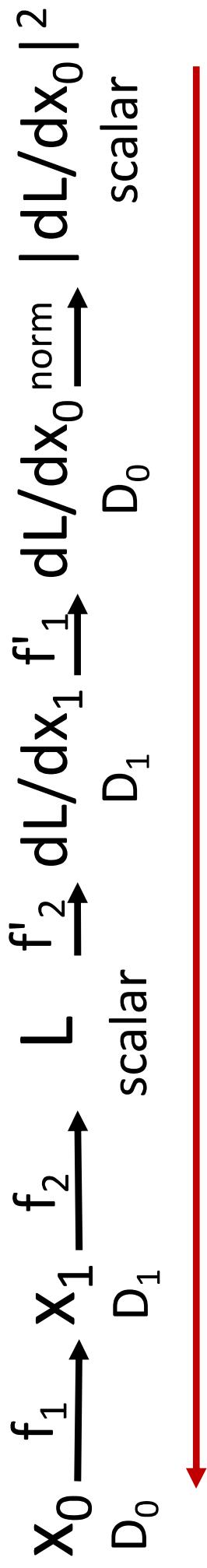
This is implemented in PyTorch / Tensorflow!

## Hessian / vector multiply

$$\frac{\partial^2 L}{\partial x_0^2} v = \frac{\partial}{\partial x_0} \left[ \frac{\partial L}{\partial x_0} \cdot v \right] \quad (\text{if } v \text{ doesn't depend on } x_0)$$

$D_0 \times D_0 \quad D_0$

## Backprop: Higher-Order Derivatives



Backprop!  
This is implemented in PyTorch / Tensorflow!

**Example:** Regularization to penalize the norm of the gradient

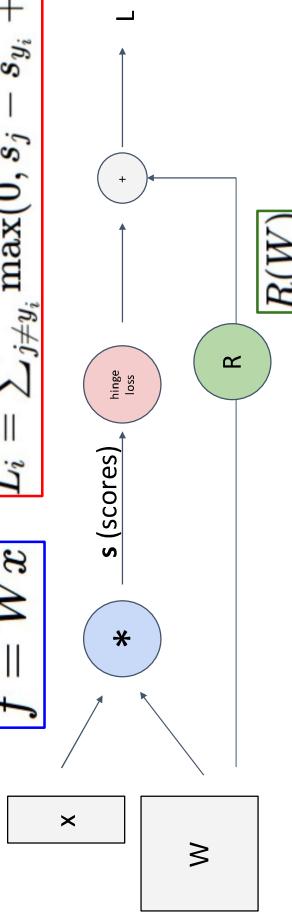
$$R(W) = \left\| \frac{\partial L}{\partial W} \right\|_2^2 = \left( \frac{\partial L}{\partial W} \right) \cdot \left( \frac{\partial L}{\partial W} \right) \quad \frac{\partial}{\partial x_0} [R(W)] = 2 \left( \frac{\partial^2 L}{\partial x_0^2} \right) \left( \frac{\partial L}{\partial x_0} \right)$$

Gulrajani et al, "Improved Training of Wasserstein GANS", NeurIPS 2017

# Summary

Represent complex expressions  
as computational graphs

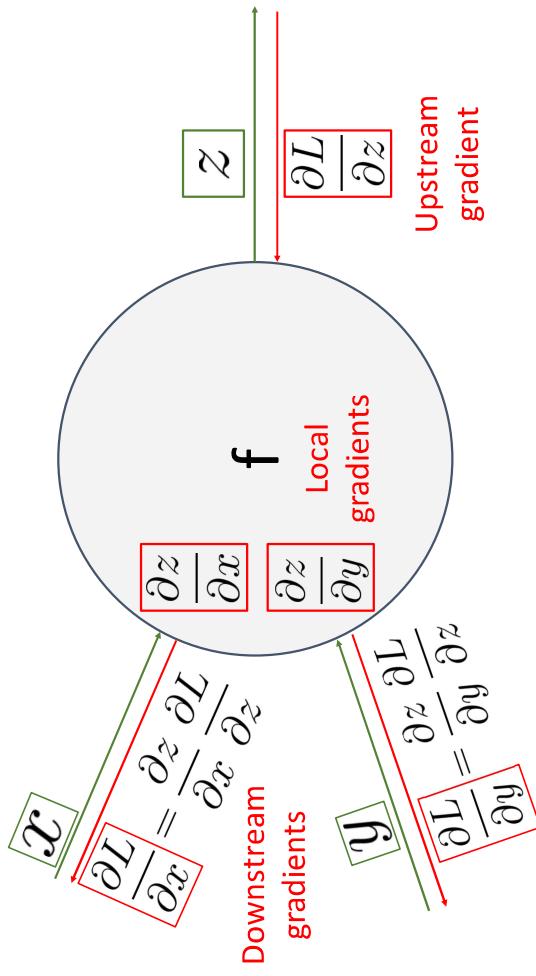
$$f = Wx \quad L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$



Forward pass computes outputs

Backward pass computes gradients

During the backward pass, each node in the graph receives **upstream gradients** and multiplies them by local gradients to compute **downstream gradients**



# Summary

Backprop can be implemented with “flat” code where the backward pass looks like forward pass reversed (Use this for A2!)

```
def f(w0, x0, w1, x1, w2):
    s0 = w0 * x0
    s1 = w1 * x1
    s2 = s0 + s1
    s3 = s2 + w2
    L = sigmoid(s3)

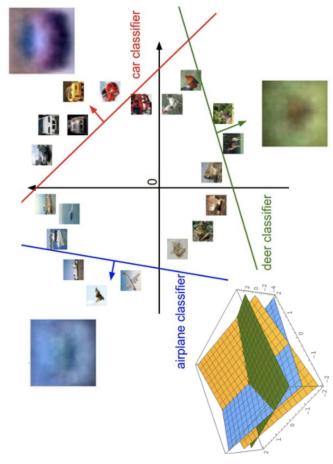
    grad_L = 1.0
    grad_s3 = grad_L * (1 - L) * L
    grad_w2 = grad_s3
    grad_s2 = grad_s3
    grad_s0 = grad_s2
    grad_s1 = grad_s2
    grad_w1 = grad_s1 * x1
    grad_x1 = grad_s1 * w1
    grad_w0 = grad_s0 * x0
    grad_x0 = grad_s0 * w0
```

Backprop can be implemented with a modular API, as a set of paired forward/backward functions (We will do this on A3!)

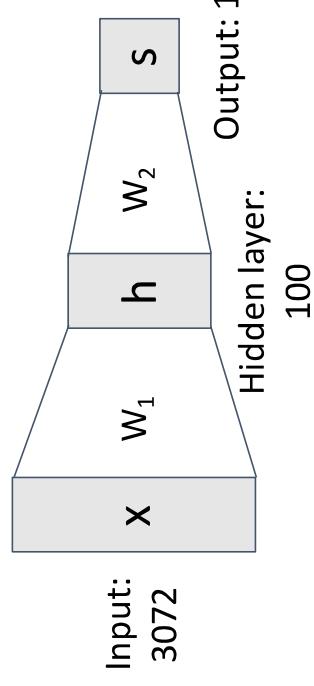
```
class Multiply(torch.autograd.Function):
    @staticmethod
    def forward(ctx, x, y):
        ctx.save_for_backward(x, y)
        z = x * y
        return z

    @staticmethod
    def backward(ctx, grad_z):
        x, y = ctx.saved_tensors
        grad_x = y * grad_z # dz/dx * dL/dz
        grad_y = x * grad_z # dz/dy * dL/dz
        return grad_x, grad_y
```

$$f(x, W) = Wx$$

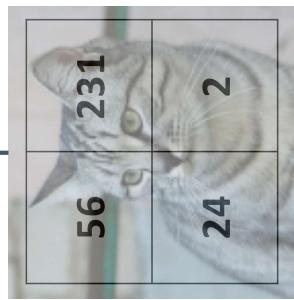


$$f = W_2 \max(0, W_1 x)$$



Stretch pixels into column

**Problem:** So far our classifiers don't respect the spatial structure of images!



Input image  
(2, 2)  
(4, )

56  
231  
24  
2

Next time:  
Convolutional Neural Networks