# DeepRob

**Lecture 6**
**Backpropagation**
**University of Minnesota**

$$\frac{\partial L}{\partial W_{\ell_1}} \longleftarrow \frac{\partial L}{\partial W_{\ell_2}} \longleftarrow \frac{\partial L}{\partial W_{\ell_3}} \longleftarrow \frac{\partial L}{\partial W_{\ell_4}} \longleftarrow \frac{\partial L}{\partial W_{\ell_5}} \longleftarrow \frac{\partial L}{\partial \text{Out}}$$

# Project 1—Reminder

- Instructions and code available on the website

  - Here: https://rpm-lab.github.io/CSCI5980-F24-DeepRob/projects/project1/

- Uses Python, PyTorch and Google Colab

- Implement KNN, linear SVM, and linear softmax classifiers

- **Autograder is available!**

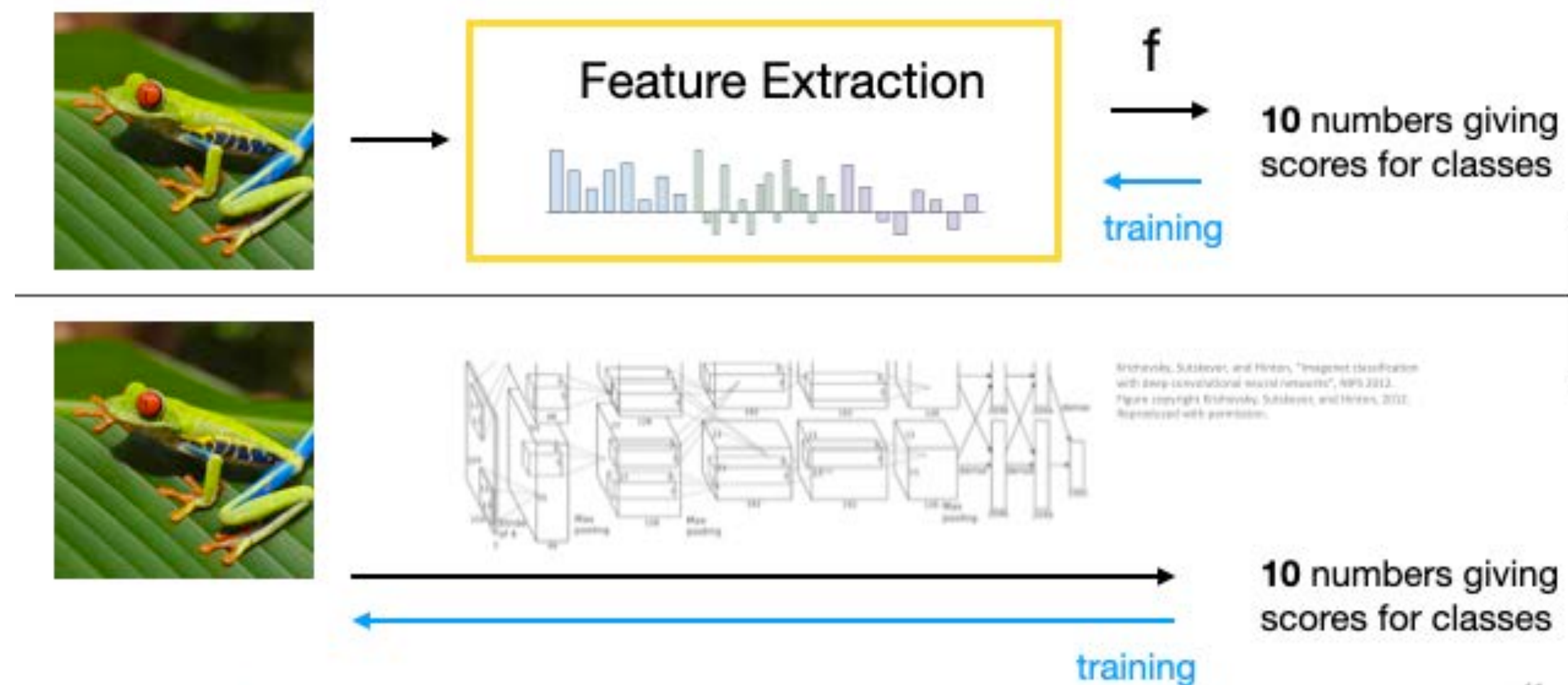- **Due Monday, Sept 30th 11:59 PM CT**

# Recap from Previous Lecture

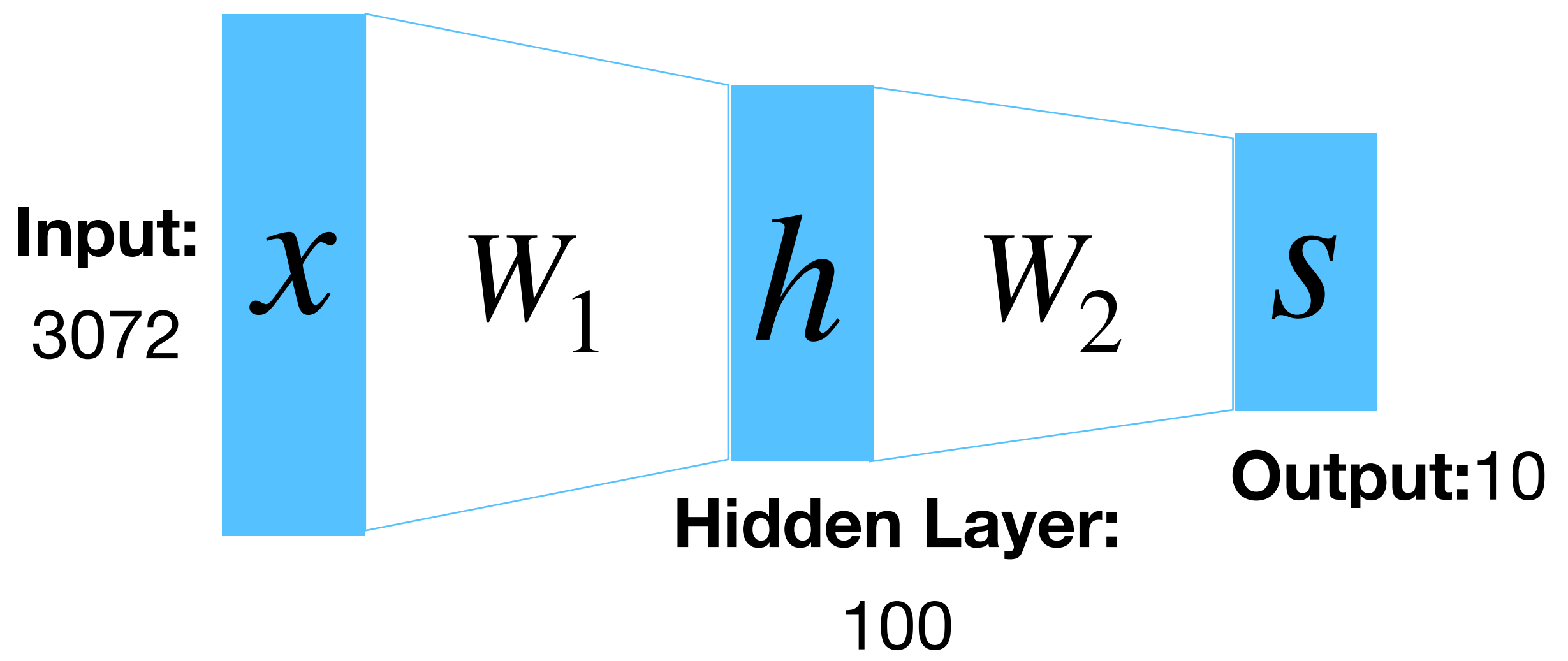Feature transform + Linear classifier allows nonlinear decision boundaries
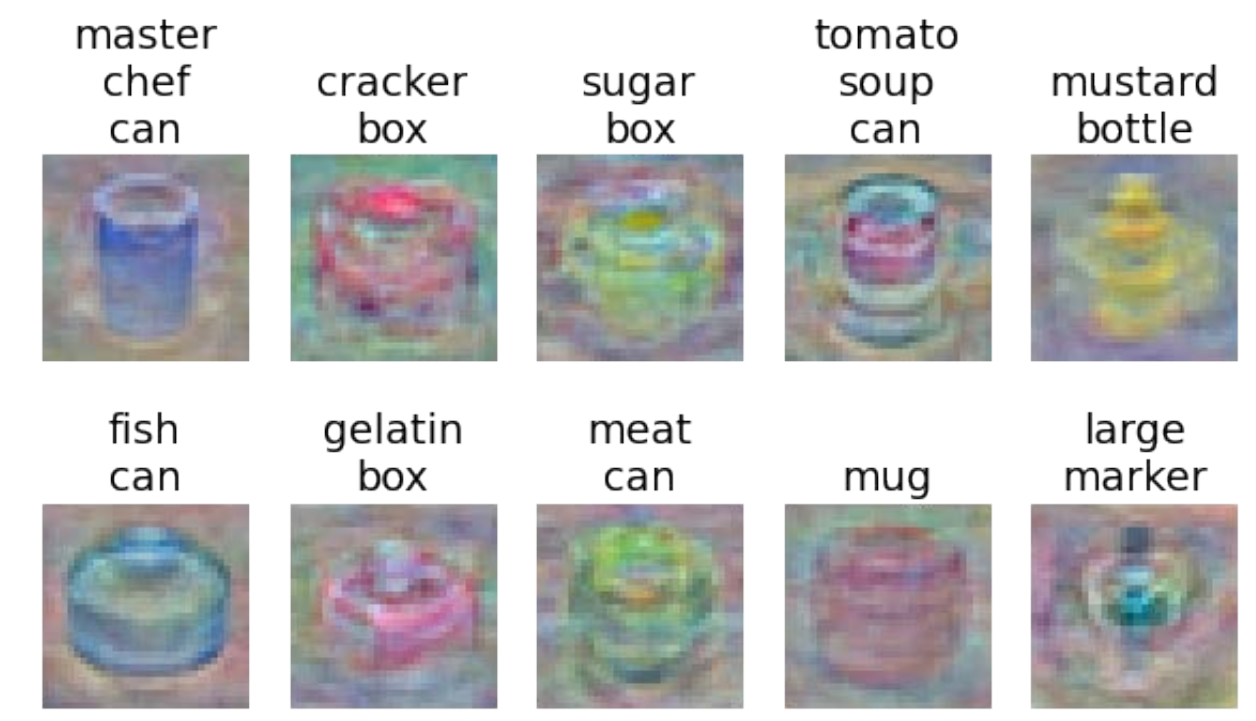
Neural Networks as learnable feature transforms

# Recap from Previous Lecture

From linear classifiers to
fully-connected networks

$$f(x) = W_2 \max(0, W_1 x + b_1) + b_2$$

**Input:**
3072

$x$ $W_1$ $h$ $W_2$ $s$

**Hidden Layer:**
100

**Output:** 10

Linear classifier: One template per class
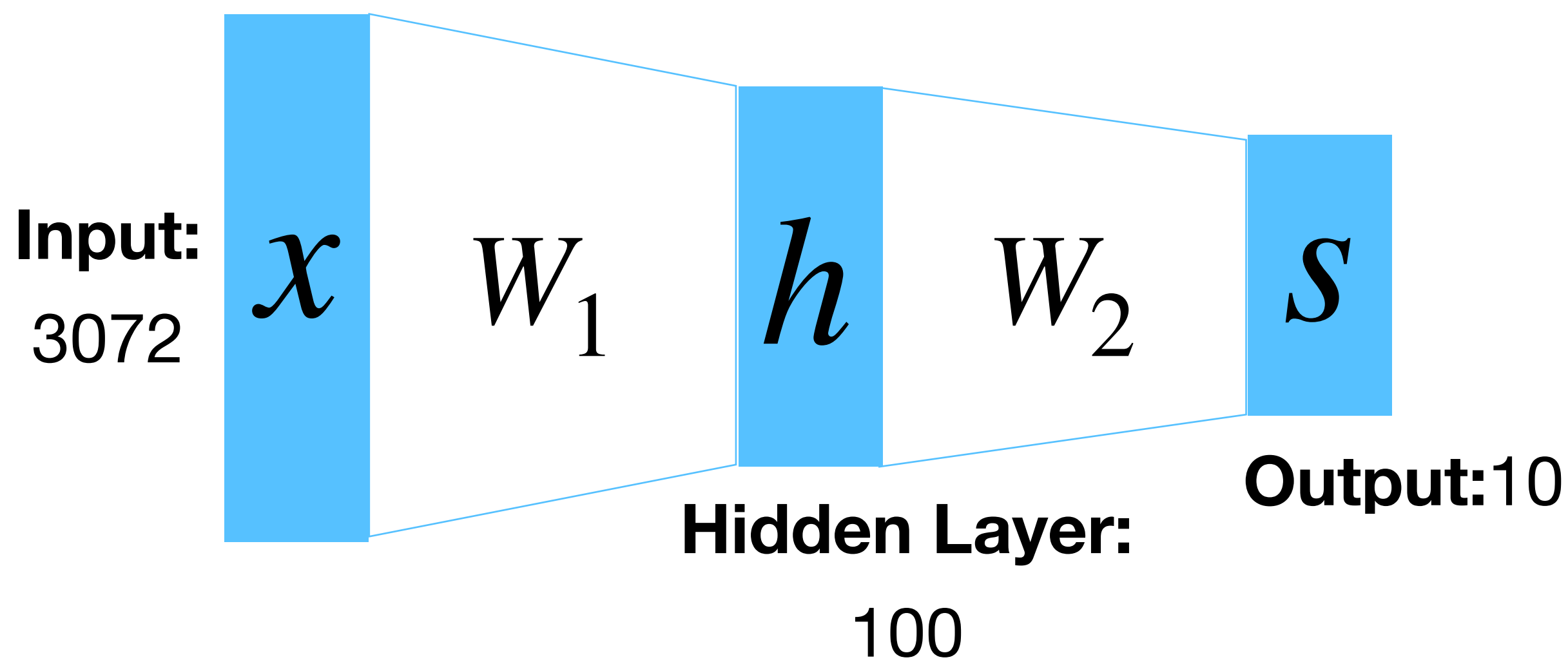


Neural networks: Many reusable templates

# Recap from Previous Lecture

From linear classifiers to fully-connected networks

$$f : X \subseteq \mathbb{R}^N \to \mathbb{R}$$

$$f(x) = W_2 \max(0, W_1 x + b_1) + b_2 \leq t f(tx_1 + (1-t)x_2) \leq t f$$

**Input:** $x$ $W_1$ $h$ $W_2$ $s$

3072

**Hidden Layer:**

100

**Output:** 10

### Space Warping



### Universal approximation

$$x_1, x_2 \in X, t \in [0, 1]$$

$$- t) f(x_2)$$



### Nonconvex



loss

w1[0, 0]

# Problem: How to compute gradients?

**ReLU activation**

$$s = W_2 \boxed{\max(0, W_1 x + b_1)} + b_2 \qquad \text{Nonlinear score function}$$

$$L_i = \boxed{\sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)} \qquad \text{Per-element data loss}$$

**Hinge loss**

$$R(W) = \sum_k W_k^2 \qquad \text{L2 regularization}$$

**Data loss**

$$L(W_1, W_2, b_1, b_2) = \boxed{\frac{1}{N} \sum_{i=1}^{N} L_i} + \boxed{\lambda R(W_1) + \lambda R(W_2)} \quad \text{Total loss}$$

**Regularization term**

If we can compute $\dfrac{\delta L}{\delta W_1}, \dfrac{\delta L}{\delta W_2}, \dfrac{\delta L}{\delta b_1}, \dfrac{\delta L}{\delta b_2}$ then we can optimize with SGD

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

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

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

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

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

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

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

**Problem**: Very tedious with lots of matrix calculus

**Problem**: What if we want to change the 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

$s = Wx$

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

$x$

$W$

$*$

**Hinge loss**

$+$

$R$

**L**

$R(W)$

# Deep Network (AlexNet)



**Input image**

**Weights**

**Loss**

# Backpropagation: Simple Example

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

# Backpropagation: Simple Example

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

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

# Backpropagation: Simple Example

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

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

**1. Forward pass**: Compute outputs

$$q = x + y \qquad f = q \cdot z$$

# Backpropagation: Simple Example

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

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

**1. Forward pass**: Compute outputs

$$q = x + y \qquad f = q \cdot z$$

**2. Backward pass**: Compute derivatives

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

# Backpropagation: Simple Example

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

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



**1. Forward pass**: Compute outputs

$$q = x + y \qquad \boxed{f = q \cdot z}$$

**2. Backward pass**: Compute derivatives

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

$$\boxed{\frac{\partial f}{\partial f}}$$

# Backpropagation: Simple Example

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

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

**1. Forward pass**: Compute outputs

$$q = x + y \qquad \boxed{f = q \cdot z}$$

**2. Backward pass**: Compute derivatives

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
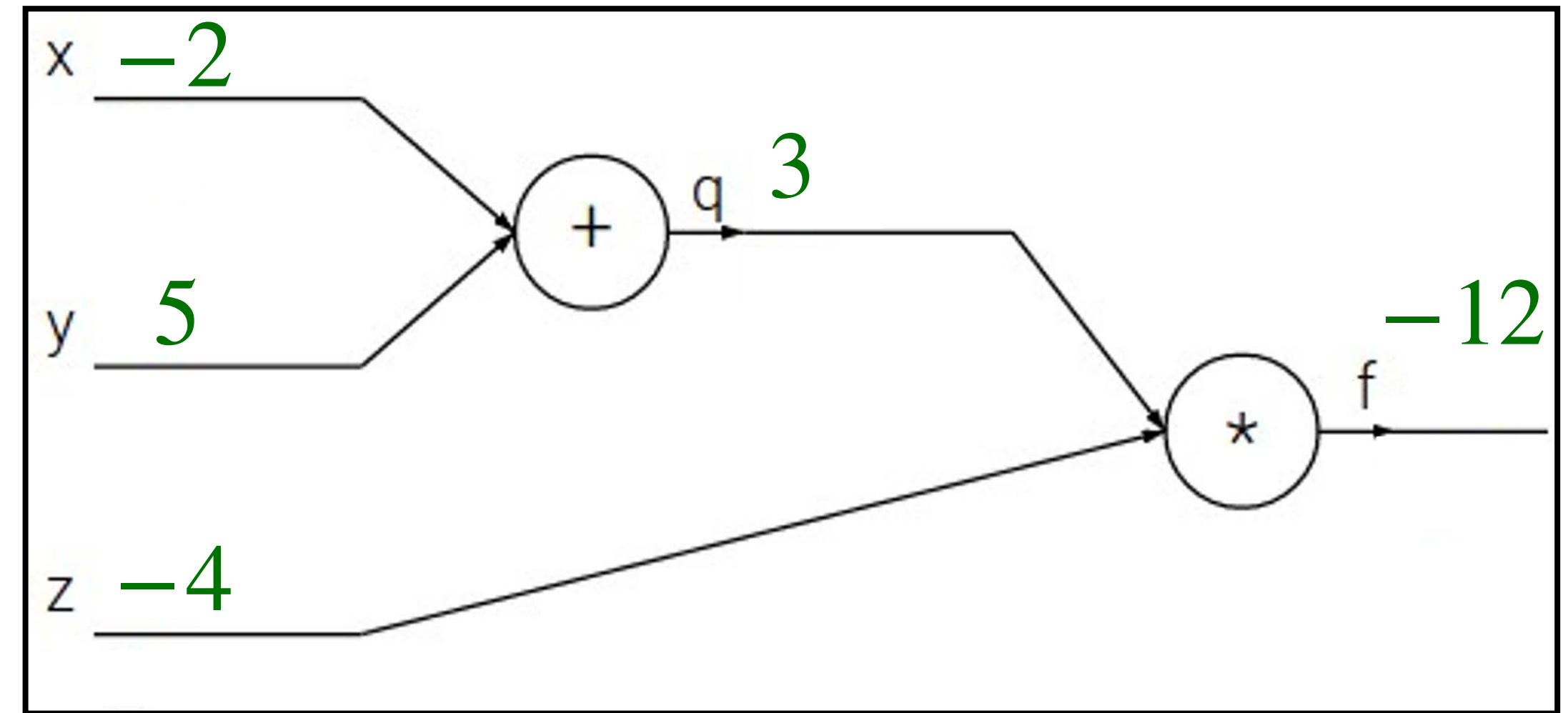
# Backpropagation: Simple Example
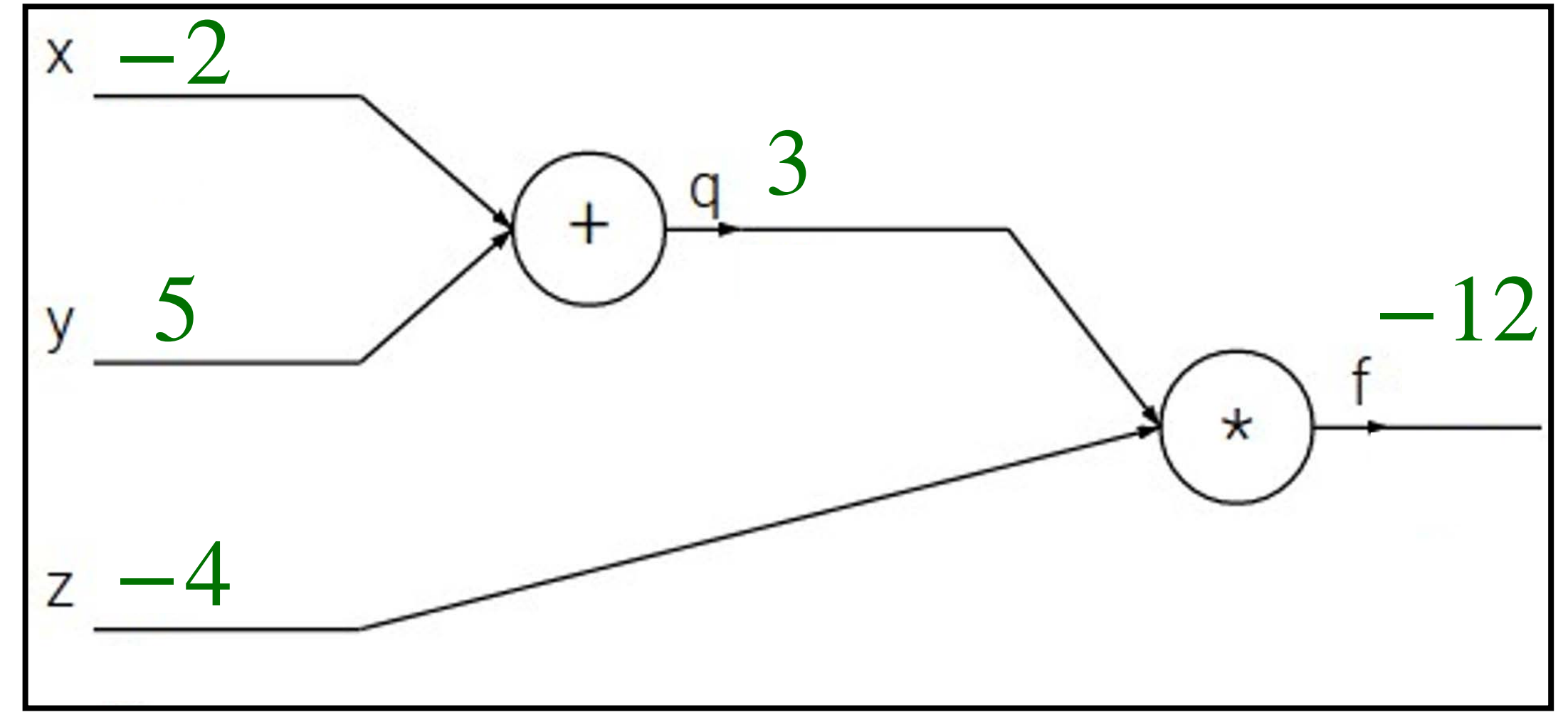
$$f(x, y, z) = (x + y) \cdot z$$

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

**1. Forward pass**: Compute outputs

$$q = x + y \qquad \boxed{f = q \cdot z}$$

**2. Backward pass**: Compute derivatives

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

# Backpropagation: Simple Example

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

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

**1. Forward pass**: Compute outputs

$$q = x + y \qquad \boxed{f = q \cdot z}$$

**2. Backward pass**: Compute derivatives

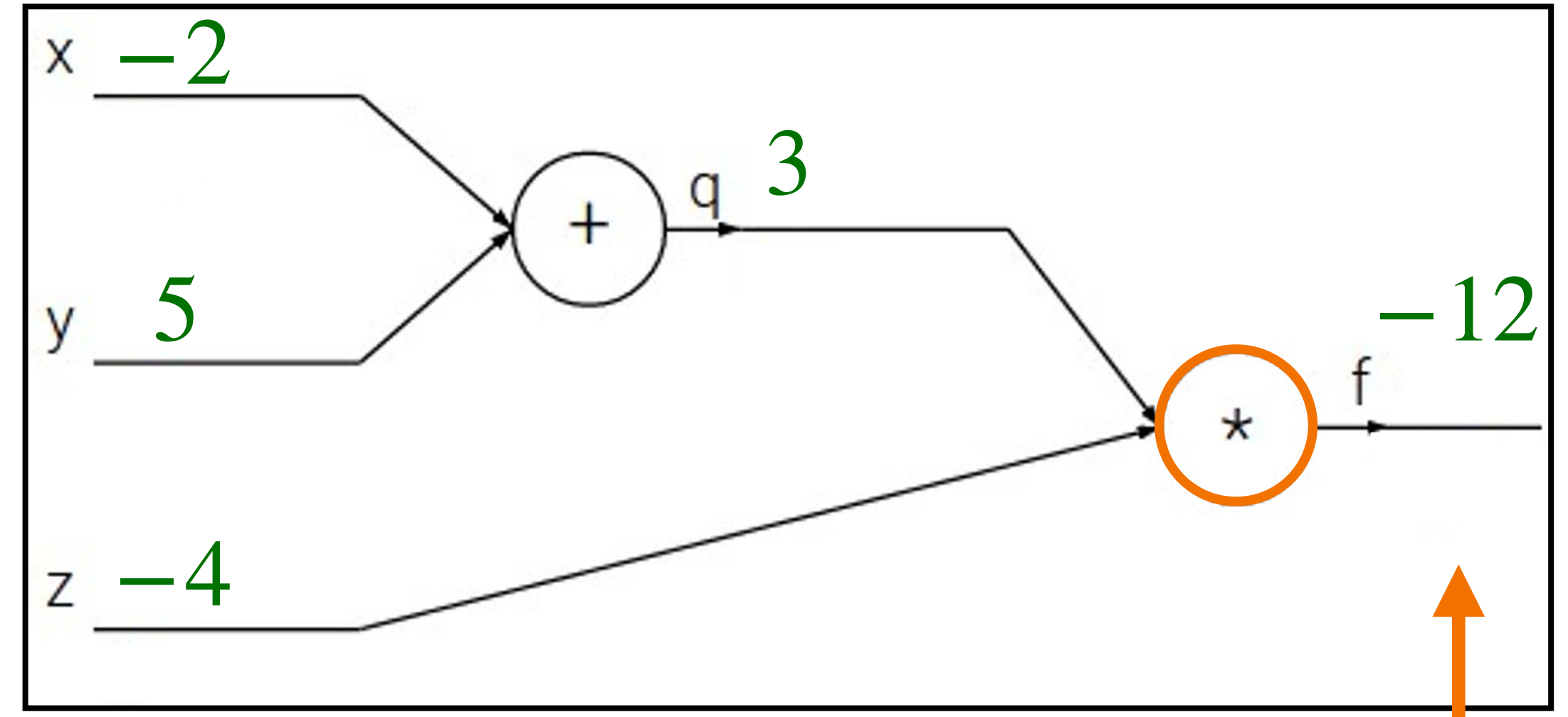Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



$$\boxed{\dfrac{\partial f}{\partial z}}$$

# Backpropagation: Simple Example

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

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

**1. Forward pass**: Compute outputs

$$q = x + y \qquad \boxed{f = q \cdot z}$$

**2. Backward pass**: Compute derivatives

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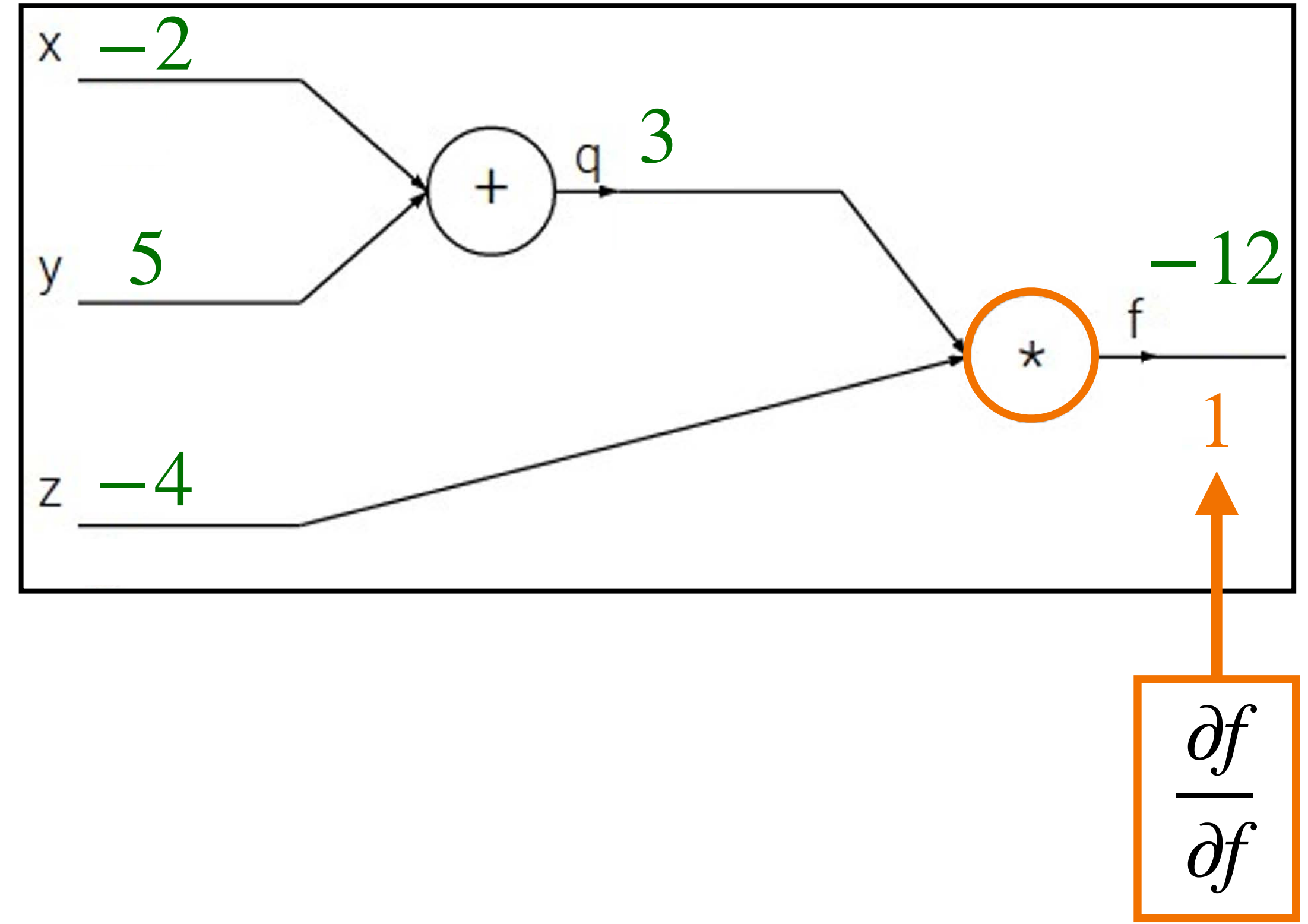


$$\dfrac{\partial f}{\partial z}$$

# Backpropagation: Simple Example

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

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

**1. Forward pass**: Compute outputs

$$q = x + y \qquad \boxed{f = q \cdot z}$$

**2. Backward pass**: Compute derivatives

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



$$\boxed{\dfrac{\partial f}{\partial z} = q}$$

# Backpropagation: Simple Example

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

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

**1. Forward pass**: Compute outputs

$$q = x + y \qquad \boxed{f = q \cdot z}$$

**2. Backward pass**: Compute derivatives

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
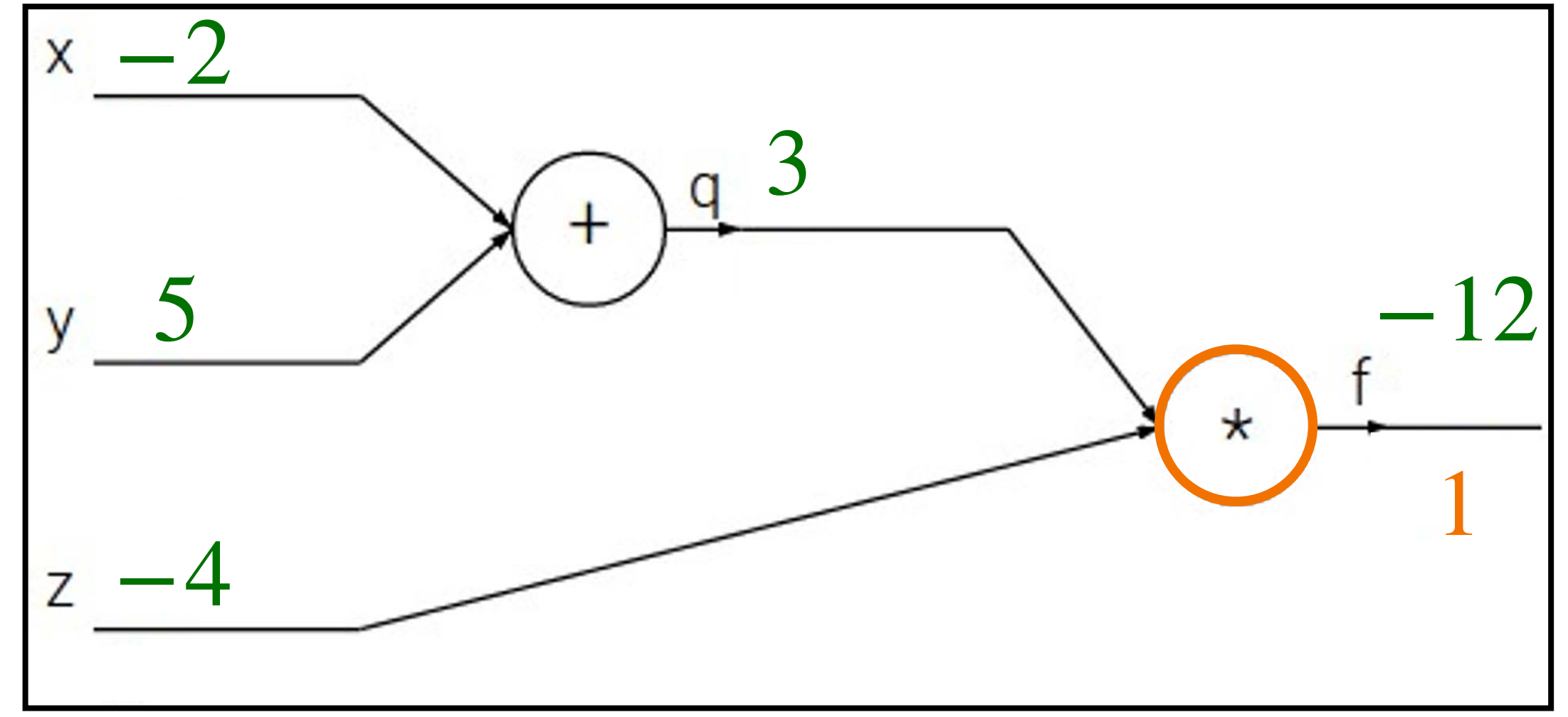
# Backpropagation: Simple Example
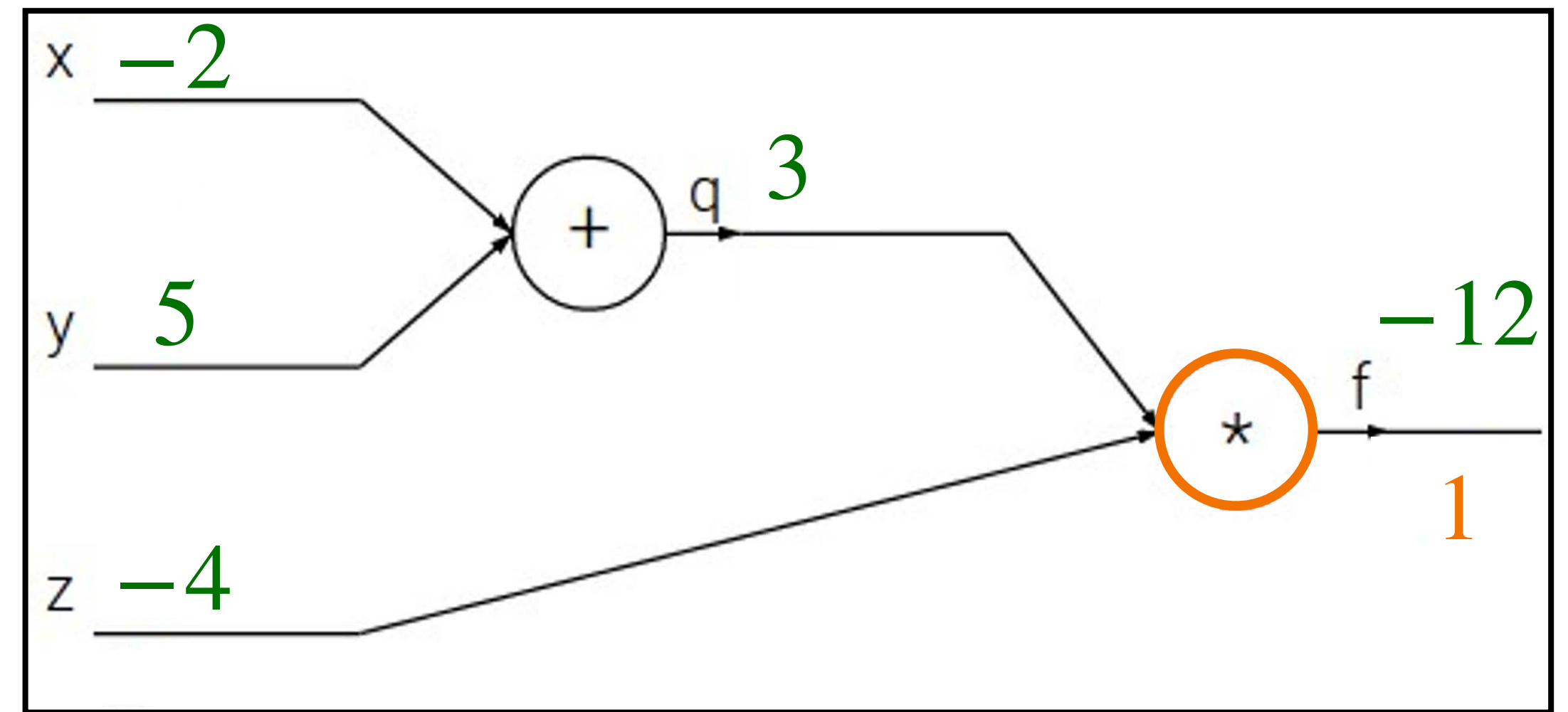
$$f(x, y, z) = (x + y) \cdot z$$

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

**1. Forward pass**: Compute outputs

$$q = x + y \qquad \boxed{f = q \cdot z}$$

**2. Backward pass**: Compute derivatives

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



$$\frac{\partial f}{\partial q}$$

# Backpropagation: Simple Example

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

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

**1. Forward pass**: Compute outputs

$$q = x + y \qquad \boxed{f = q \cdot z}$$

**2. Backward pass**: Compute derivatives

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



$$\boxed{\frac{\partial f}{\partial q} = z}$$

# Backpropagation: Simple Example

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

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

**1. Forward pass**: Compute outputs

$$\boxed{q = x + y} \quad f = q \cdot z$$

**2. Backward pass**: Compute derivatives

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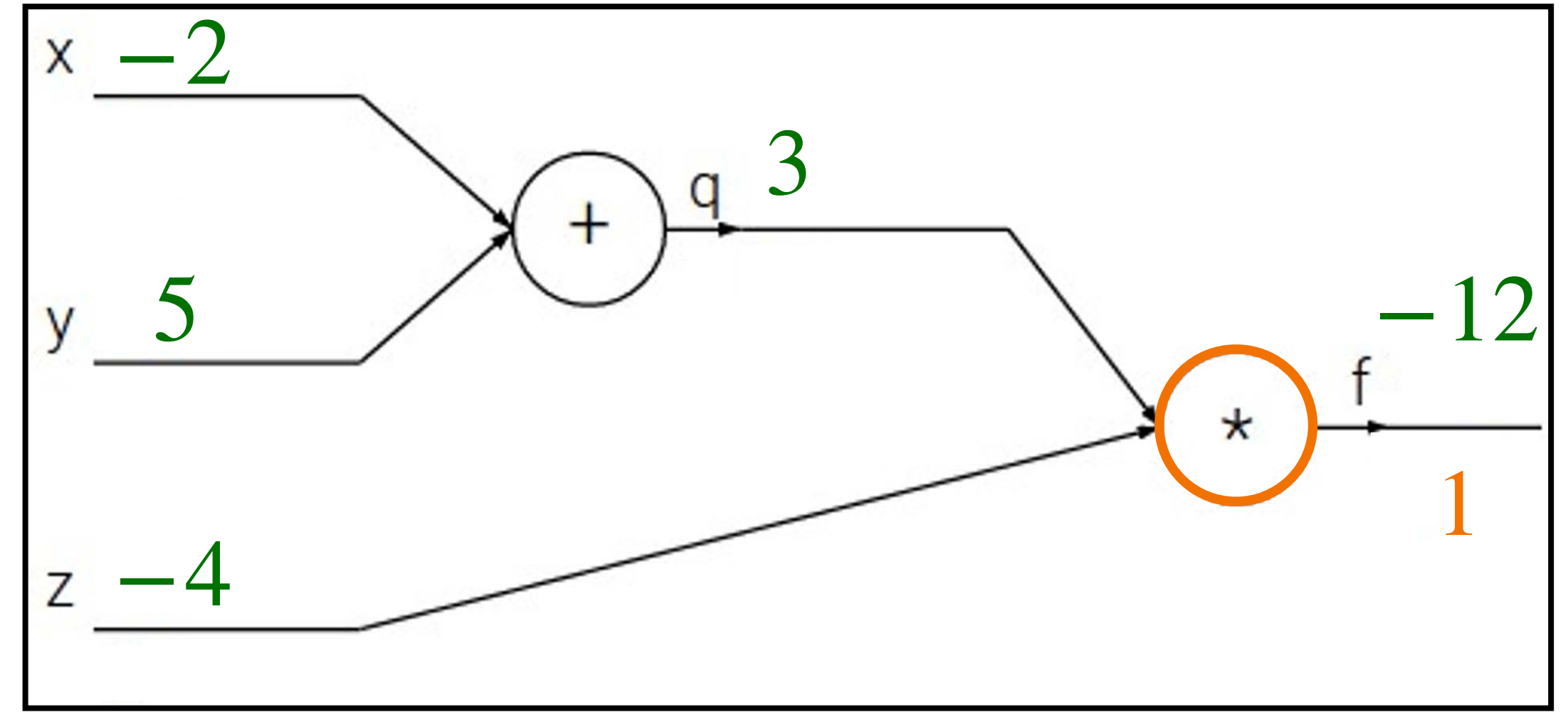


$$\frac{\partial f}{\partial y}$$

# Backpropagation: Simple Example

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

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

**1. Forward pass**: Compute outputs

$$\boxed{q = x + y} \qquad f = q \cdot z$$

**2. Backward pass**: Compute derivatives

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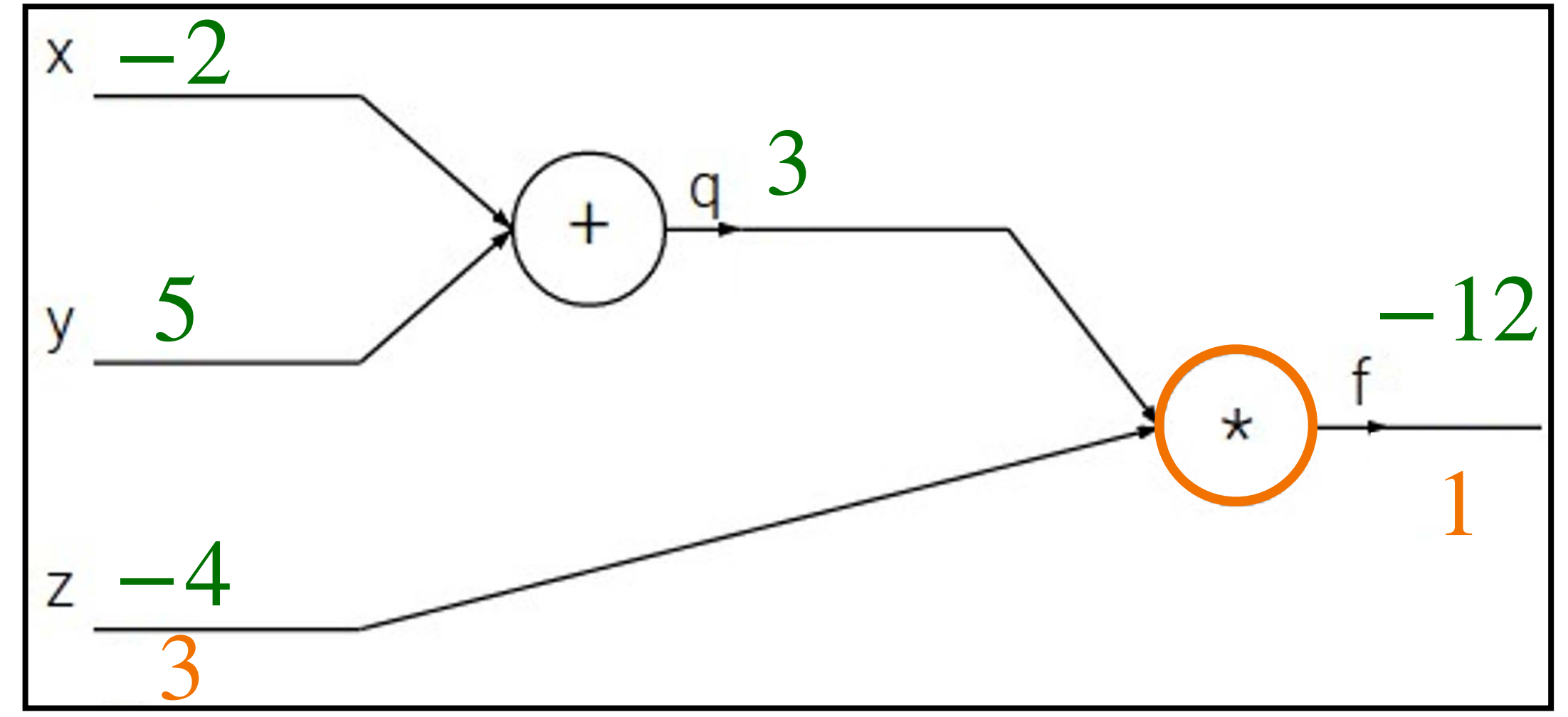


$$\frac{\partial f}{\partial y}$$

# Backpropagation: Simple Example

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

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

**1. Forward pass**: Compute outputs

$$\boxed{q = x + y} \qquad f = q \cdot z$$

**2. Backward pass**: Compute derivatives

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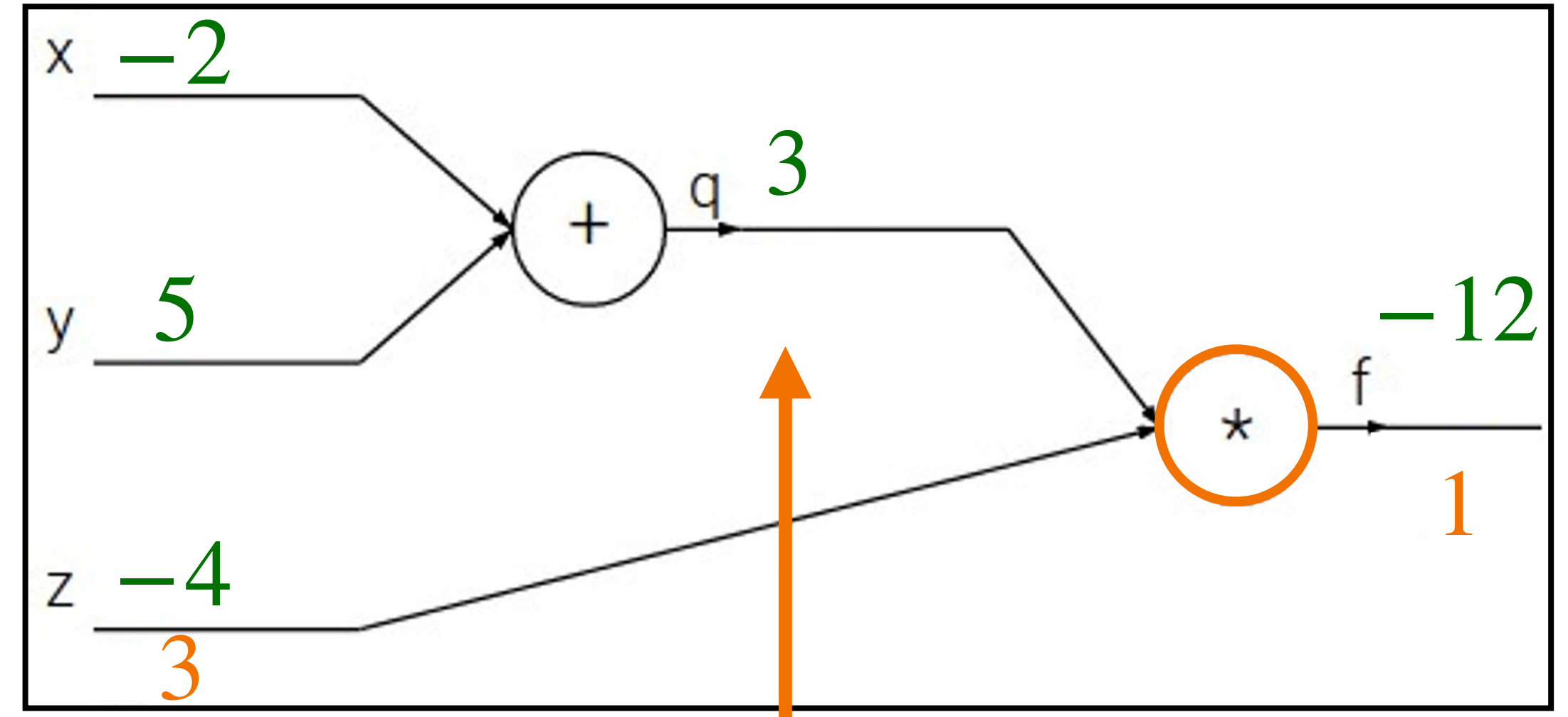


$$\frac{\partial f}{\partial y} = \frac{\partial q}{\partial y} \frac{\partial f}{\partial q}$$

Downstream Gradient    Local Gradient    Upstream Gradient

# Backpropagation: Simple Example

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

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

**1. Forward pass**: Compute outputs

$$\boxed{q = x + y} \qquad f = q \cdot z$$

**2. Backward pass**: Compute derivatives

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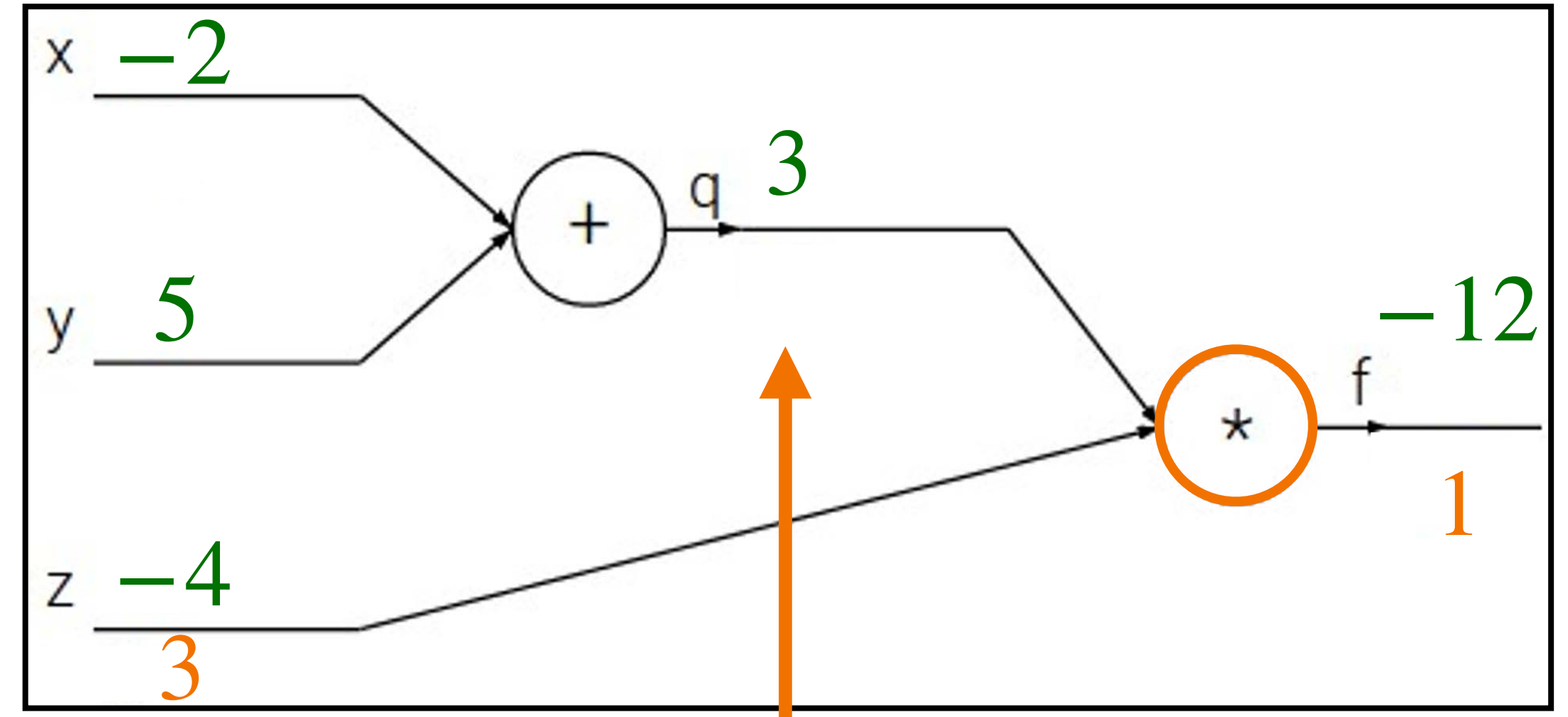


$$\frac{\partial f}{\partial y} = \frac{\partial q}{\partial y} \frac{\partial f}{\partial q} \qquad \frac{\partial q}{\partial y} = 1$$

Downstream Gradient    Local Gradient    Upstream Gradient

# Backpropagation: Simple Example

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

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

**1. Forward pass**: Compute outputs

$$\boxed{q = x + y} \qquad f = q \cdot z$$

**2. Backward pass**: Compute derivatives

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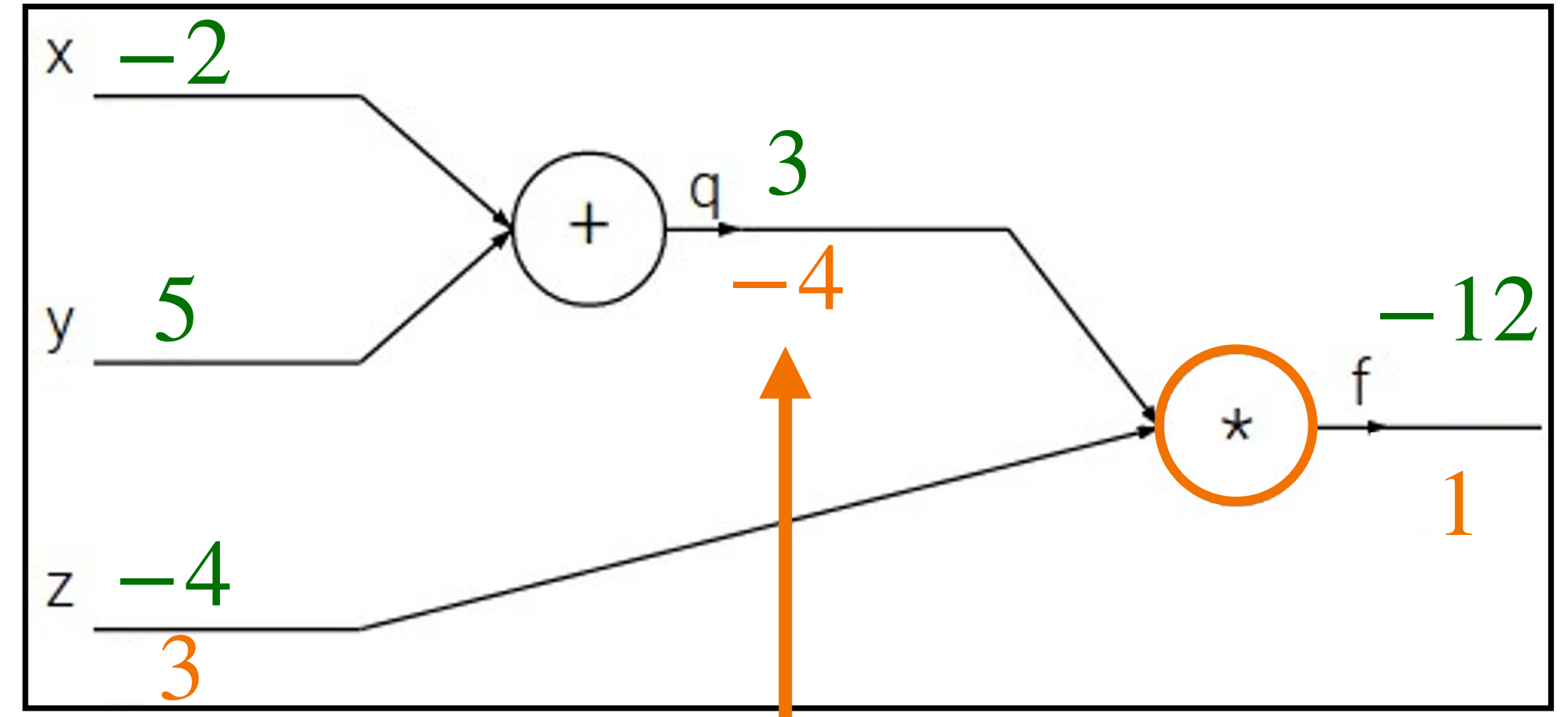


$$\frac{\partial f}{\partial y} = \frac{\partial q}{\partial y} \frac{\partial f}{\partial q}$$

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

Downstream Gradient    Local Gradient    Upstream Gradient

# Backpropagation: Simple Example

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

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

**1. Forward pass**: Compute outputs

$$q = x + y \qquad f = q \cdot z$$

**2. Backward pass**: Compute derivatives

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



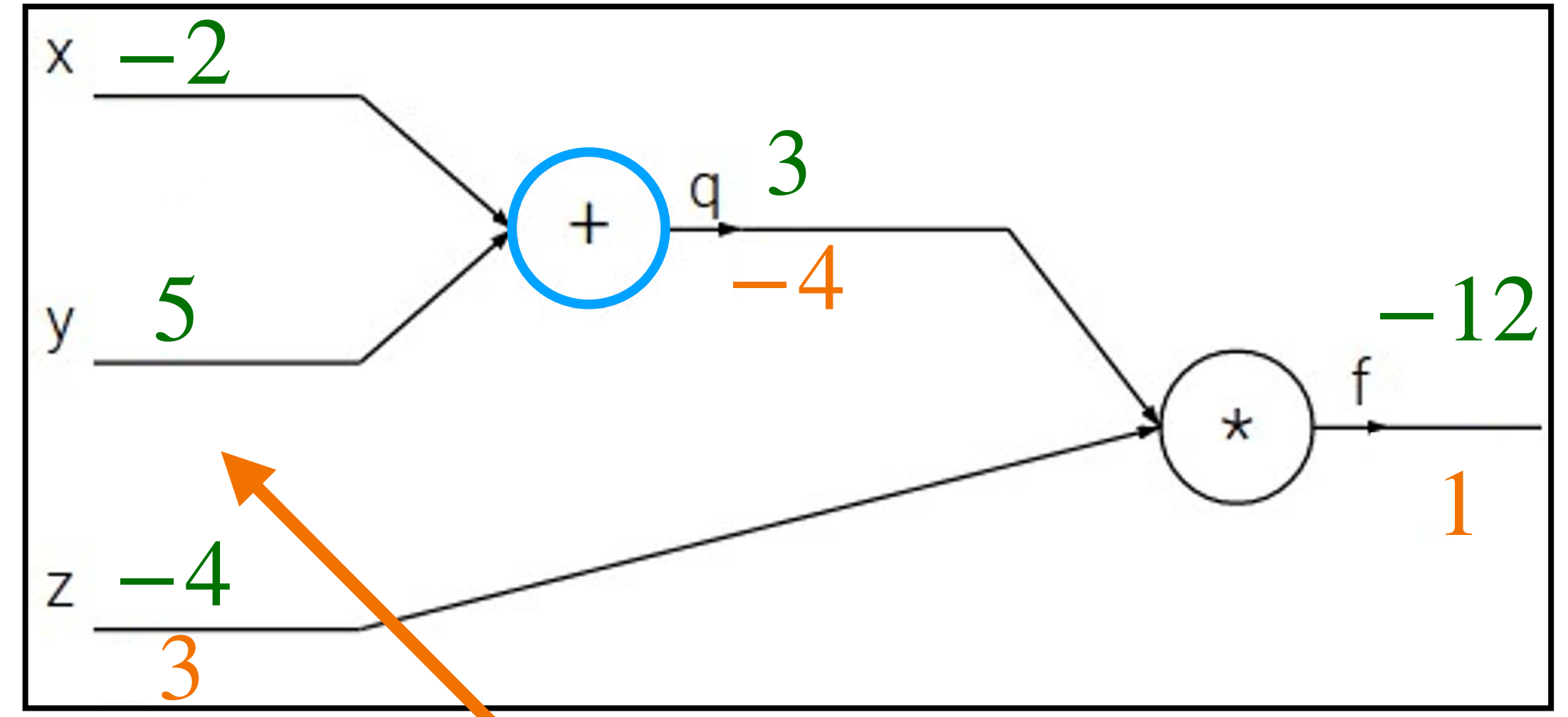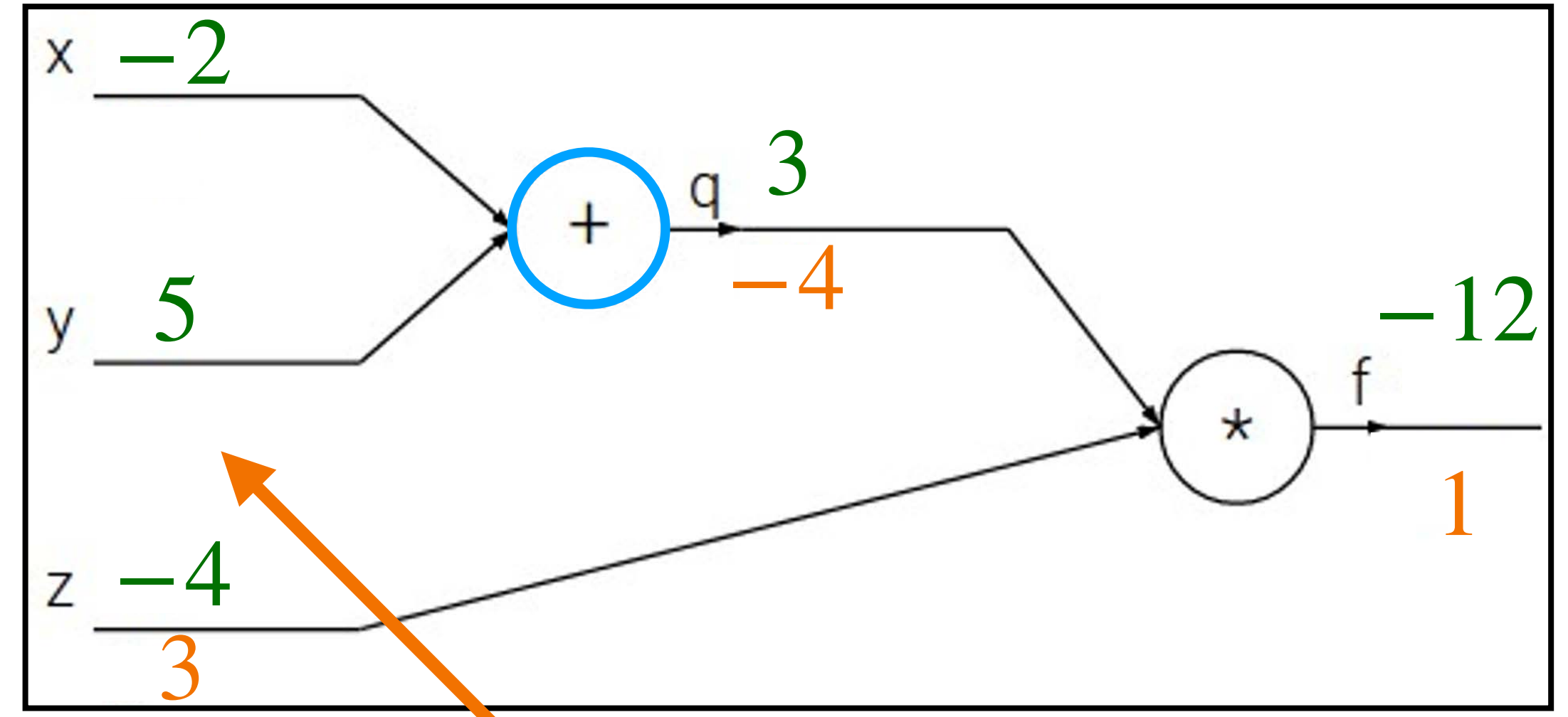$$\frac{\partial f}{\partial x} = \frac{\partial q}{\partial x} \frac{\partial f}{\partial q} \qquad \frac{\partial q}{\partial x} = 1$$

Downstream Gradient      Local Gradient      Upstream Gradient
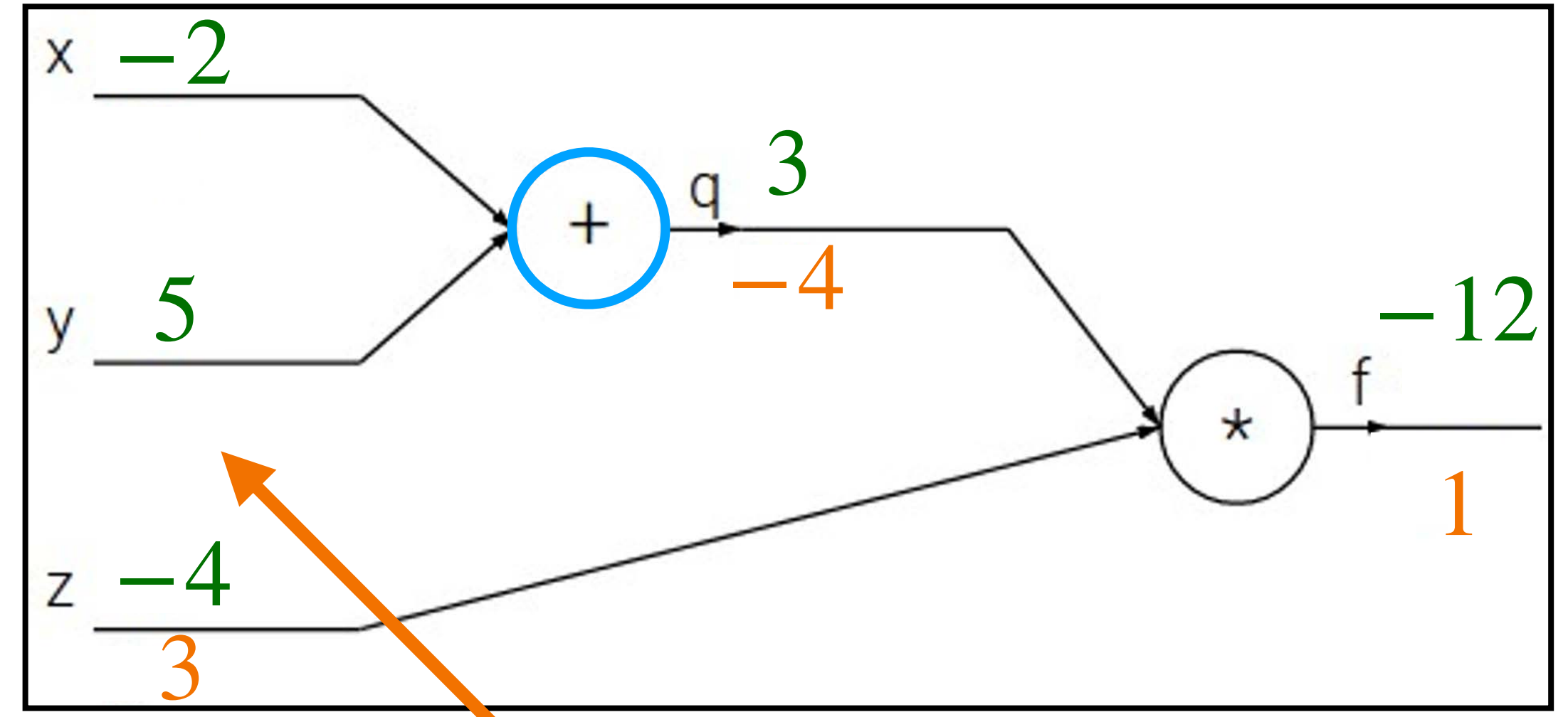
# Backpropagation: Simple Example

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

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

**1. Forward pass**: Compute outputs

$$\boxed{q = x + y} \quad f = q \cdot z$$

**2. Backward pass**: Compute derivatives

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

x $-2$
$-4$

y $5$
$-4$

q $3$
$-4$

z $-4$
$3$

$-12$

f

$*$

$1$

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

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

Downstream
Gradient

Local
Gradient

Upstream
Gradient

# Local Properties of Backpropagation



$x$

$y$

$f$

$z$

# Local Properties of Backpropagation

$x$

$y$

$f$

$z$

# Local Properties of Backpropagation



$x$

$y$

$f$

$z$

$\dfrac{\partial L}{\partial z}$   Upstream Gradient

$$x$$

$$y$$

$$\frac{\partial z}{\partial x}$$

$$\frac{\partial z}{\partial y}$$

$$f$$

Local Gradients

$$z$$

$$\frac{\partial L}{\partial z}$$

Upstream Gradient

# Local Properties of Backpropagation



$x$

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

Downstream
Gradients

$y$

$$\frac{\partial L}{\partial y} = \frac{\partial z}{\partial y}\frac{\partial L}{\partial z}$$

$$\frac{\partial z}{\partial x}$$

$$\frac{\partial z}{\partial y}$$

$f$

Local
Gradients

$z$

$$\frac{\partial L}{\partial z}$$

Upstream
Gradient

34

# Local Properties of Backpropagation

$$\boxed{x}$$

$$\boxed{\frac{\partial L}{\partial x} = \frac{\partial z}{\partial x}\frac{\partial L}{\partial z}}$$

Downstream Gradients

$$\boxed{y}$$

$$\boxed{\frac{\partial L}{\partial y} = \frac{\partial z}{\partial y}\frac{\partial L}{\partial z}}$$

$$\boxed{\frac{\partial z}{\partial x}}$$

$$\boxed{\frac{\partial z}{\partial y}}$$

$f$

Local Gradients

$$\boxed{z}$$

$$\boxed{\frac{\partial L}{\partial z}}$$

Upstream Gradient

# Another example

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

# Another example

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

**1. Forward pass**: Compute outputs

# Another example

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

**1. Forward pass**: Compute outputs



w0  2.00
x0  −1.00
  *  −2.00
w1  −3.00
x1  −2.00
  *  6.00
  +  4.00
  +  1.00
  *-1  −1.00
  exp  0.37
  +1  1.37
  1/x  0.73
w2  −3.00

# Another example

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

**1. Forward pass**: Compute outputs

**2. Backward pass**: Compute gradients

# Another example

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

**1. Forward pass**: Compute outputs

**2. Backward pass**: Compute gradients



w0  2.00

x0  −1.00

*  −2.00

w1  −3.00

x1  −2.00

*  6.00

+  4.00

w2  −3.00

+  1.00

*-1  −1.00

exp  0.37

+1  1.37

1/x  0.73  1.00

Base Case

# Another example

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

**1. Forward pass**: Compute outputs

**2. Backward pass**: Compute gradients

Local Gradient

$$\frac{\partial}{\partial x}\left[\frac{1}{x}\right] = -\frac{1}{x^2}$$

w0  2.00

x0  −1.00

−2.00

w1  −3.00

x1  −2.00

6.00

4.00

w2  −3.00

1.00     −1.00     0.37     1.37     0.73

*-1     exp     +1     1/x

−0.53     1.00

Downstream Gradient     Upstream Gradient

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

**1. Forward pass**: Compute outputs

**2. Backward pass**: Compute gradients



Local Gradient

$$\frac{\partial}{\partial x}\left[x + 1\right] = 1$$

Downstream Gradient

Upstream Gradient

43

# Another example

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

**1. Forward pass**: Compute outputs

**2. Backward pass**: Compute gradients

w0  2.00

x0  −1.00

\*  −2.00

w1  −3.00

x1  −2.00

\*  6.00

+  4.00

w2  −3.00

+  1.00

\*-1

−1.00
−0.20

exp

0.37
−0.53

+1

1.37
−0.53

1/x

0.73
1.00

Local Gradient

$$\frac{\partial}{\partial x}\left[ e^x \right] = e^x$$

Downstream Gradient

Upstream Gradient

# Another example

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

**1. Forward pass**: Compute outputs

**2. Backward pass**: Compute gradients

w0  2.00

−2.00

x0 −1.00

4.00

Local Gradient

$$\frac{\partial}{\partial x}\left[ -1 \cdot x \right] = -1$$

w1 −3.00

6.00

x1 −2.00

1.00    *-1    −1.00    exp    0.37    +1    1.37    1/x    0.73

0.20         −0.20       −0.53       −0.53       1.00

w2 −3.00

Downstream Gradient    Upstream Gradient

# Another example

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

**1. Forward pass**: Compute outputs

**2. Backward pass**: Compute gradients



Local Gradient

$$\frac{\partial}{\partial x}\left[x + y\right] = 1 \qquad \frac{\partial}{\partial y}\left[x + y\right] = 1$$

w0  2.00

x0  −1.00

w1  −3.00

x1  −2.00

w2  −3.00
0.20

−2.00

6.00

4.00
0.20

1.00
0.20

−1.00
−0.20

0.37
−0.53

1.37
−0.53

0.73
1.00

Downstream Gradient

Upstream Gradient

# Another example

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

**1. Forward pass**: Compute outputs

**2. Backward pass**: Compute gradients



w0  2.00

x0  −1.00

−2.00
0.20

w1  −3.00

x1  −2.00

6.00
0.20

4.00
0.20

w2  −3.00
0.20

1.00
0.20

−1.00
−0.20

0.37
−0.53

1.37
−0.53

0.73
1.00

Local Gradient

$$\frac{\partial}{\partial x}\left[x + y\right] = 1 \qquad \frac{\partial}{\partial y}\left[x + y\right] = 1$$

Downstream Gradient

Upstream Gradient

# Another example

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

**1. Forward pass**: Compute outputs

**2. Backward pass**: Compute gradients



Local Gradient

$$\frac{\partial}{\partial x}\left[x \cdot y\right] = y \qquad \frac{\partial}{\partial y}\left[x \cdot y\right] = x$$

w0  2.00
−0.20
x0  −1.00   * → −2.00  0.20
0.39

Downstream Gradient

w1  −3.00
x1  −2.00   * → 6.00  0.20

+ → 4.00  0.20

w2  −3.00
0.20

Upstream Gradient

+ → 1.00  0.20 → *-1 → −1.00  −0.20 → exp → 0.37  −0.53 → +1 → 1.37  −0.53 → 1/x → 0.73  1.00
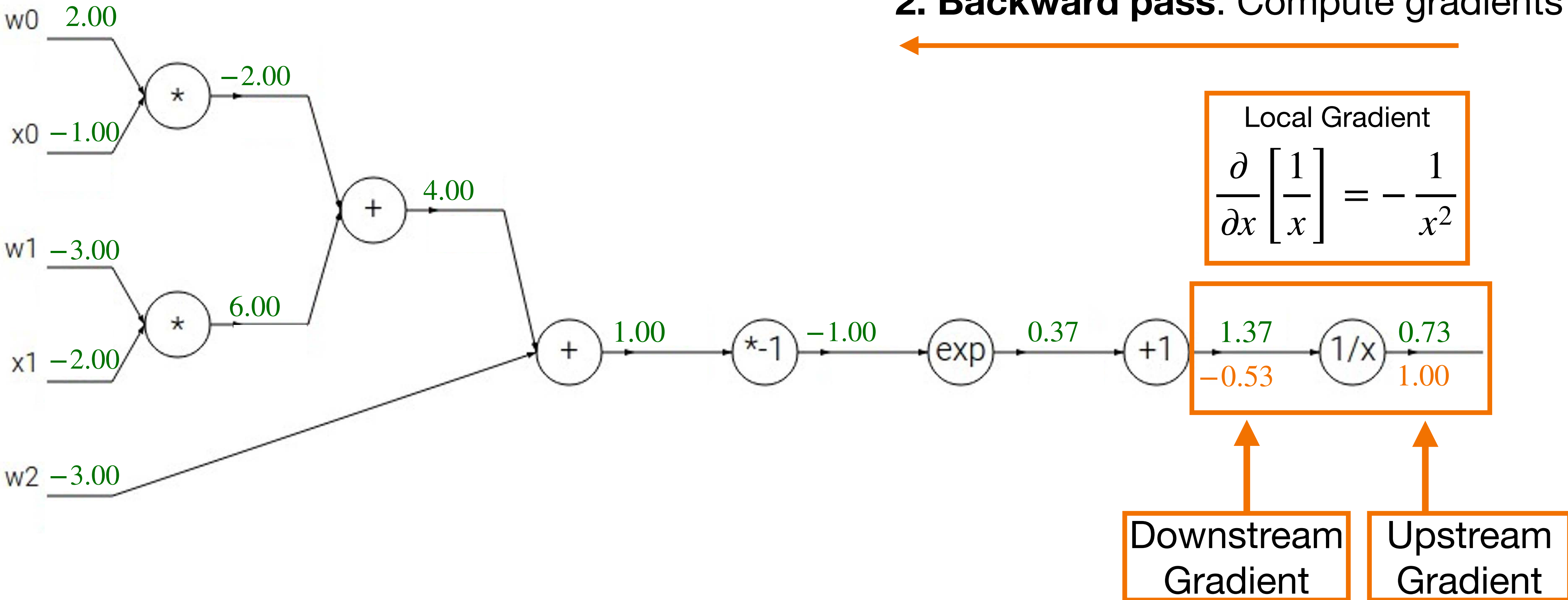
48

# Another example

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

**1. Forward pass**: Compute outputs

**2. Backward pass**: Compute gradients



w0   2.00
−0.20

x0 −1.00
0.39

*   −2.00
0.20

Local Gradient

$$\frac{\partial}{\partial x}\left[x \cdot y\right] = y \qquad \frac{\partial}{\partial y}\left[x \cdot y\right] = x$$

+   4.00
0.20

w1 −3.00
−0.39

x1 −2.00
−0.59

*   6.00
0.20

Downstream
Gradient

Upstream
Gradient

w2 −3.00
0.20

+   1.00
0.20

*-1   −1.00
−0.20

exp   0.37
−0.53

+1   1.37
−0.53

1/x   0.73
1.00

# Another example

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$$= \sigma(w_0 x_0 + w_1 x_1 + w_2)$$

**1. Forward pass**: Compute outputs

**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



w0  2.00
   −0.20

x0  −1.00
   0.39

* → −2.00 / 0.20

w1  −3.00
   −0.39

x1  −2.00
   −0.59

* → 6.00 / 0.20

+ → 4.00 / 0.20

w2  −3.00
   0.20

+ → 1.00 / 0.20

**Sigmoid**

*-1 → −1.00 / −0.20

exp → 0.37 / −0.53

+1 → 1.37 / −0.53

1/x → 0.73 / 1.00

# Another example

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$$= \sigma(w_0 x_0 + w_1 x_1 + w_2)$$

**1. Forward pass**: Compute outputs

**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

Sigmoid

w0  2.00
    −0.20

x0 −1.00
    0.39

* −2.00
  0.20

w1 −3.00
   −0.39

x1 −2.00
   −0.59

* 6.00
  0.20

+ 4.00
  0.20

w2 −3.00
   0.20

+ 1.00
  0.20

*-1  −1.00
     −0.20

exp  0.37
     −0.53

+1  1.37
    −0.53

1/x  0.73
     1.00

Sigmoid local gradient:

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

51

# Another example

$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$$= \sigma(w_0 x_0 + w_1 x_1 + w_2)$$

**1. Forward pass**: Compute outputs

**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

w0  2.00
     −0.20

x0  −1.00
     0.39

* → −2.00 / 0.20

w1  −3.00
     −0.39

x1  −2.00
     −0.59

* → 6.00 / 0.20

+ → 4.00 / 0.20

w2  −3.00
     0.20

+ → 1.00 / 0.20

**Sigmoid**

*-1 → −1.00 / −0.20

exp → 0.37 / −0.53

+1 → 1.37 / −0.53

1/x → 0.73 / 1.00

[Downstream]=[Local]·[Upstream]

$$=(1-0.73) \cdot 0.73 \cdot 1.00 = 0.20$$

Sigmoid local gradient:

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

# Patterns in Gradient Flow

**add** gate: gradient distributor

**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 outputs

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



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

# Backprop Implementation: "Flat" gradient code

**Forward pass**:
Compute outputs

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

w0  2.00
    -0.20

x0  -1.00
    0.40

*  -2.00
   0.20

w1  -3.00
    -0.40

x1  -2.00
    -0.60

*  6.00
   0.20

+  4.00
   0.20

w2  -3.00
    0.20

+  1.00
   0.20

$\sigma$

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

```
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
```

**Backward pass**:
Compute gradients

# Backprop Implementation: "Flat" gradient code



**Forward pass**:
Compute outputs

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

Base case

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

```
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
```

**Backward pass**:
Compute gradients

**Forward pass**:
Compute outputs

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

Sigmoid

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

```
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
```

**Backward pass**:
Compute gradients

# Backprop Implementation: "Flat" gradient code



**Forward pass**:
Compute outputs

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

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

Add

```
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
```

**Backward pass**:
Compute gradients

# Backprop Implementation: "Flat" gradient code



**Forward pass**: Compute outputs

```
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
```

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

Add

**Backward pass**: Compute gradients

# Backprop Implementation: "Flat" gradient code

**Forward pass**:
Compute outputs

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

w0  2.00
    -0.20

x0  -1.00
    0.40

* → -2.00
    0.20

w1  -3.00
    -0.40

x1  -2.00
    -0.60

* → 6.00
    0.20

+ → 4.00
    0.20

w2  -3.00
    0.20

+ → 1.00
    0.20

σ

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

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
```

**Backward pass**:
Compute gradients

# Backprop Implementation: "Flat" gradient code



**Forward pass**:
Compute outputs

```
def f(w0, x0, w1, x1, w2):
    s0 = w0 * x0
    s1 = w1 * x1
    s2 = s0 + s1
    s3 = s2 + w2
    L = sigmoid(s3)
```
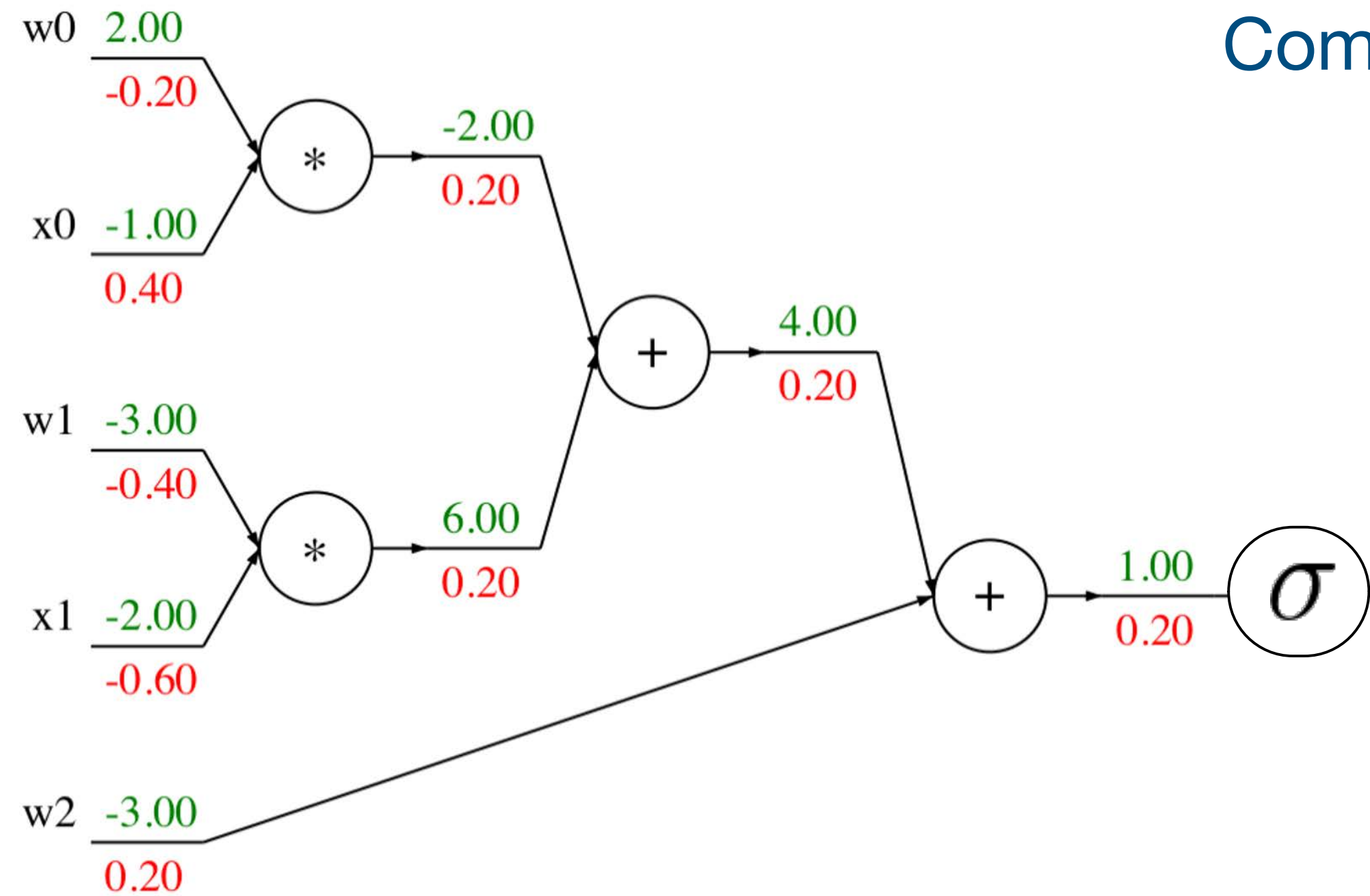
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

```
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

**Backward pass**:
Compute gradients

64

# "Flat" Backprop: Do this for Project 1 & 2

**Forward pass**:
Compute outputs

**Backward pass**:
Compute gradients

```
###############################################################
# TODO:                                                       #
# Implement a vectorized version of the structured SVM loss, storing the   #
# result in loss.                                             #
###############################################################
# Replace "pass" statement with your code
num_classes = W.shape[1]
num_train = X.shape[0]
score = # ...
correct_class_score = # ...
margin = # ...
data_loss = # ...
reg_loss = # ...
loss += data_loss + reg_loss

###############################################################
#                        END OF YOUR CODE                     #
###############################################################
```
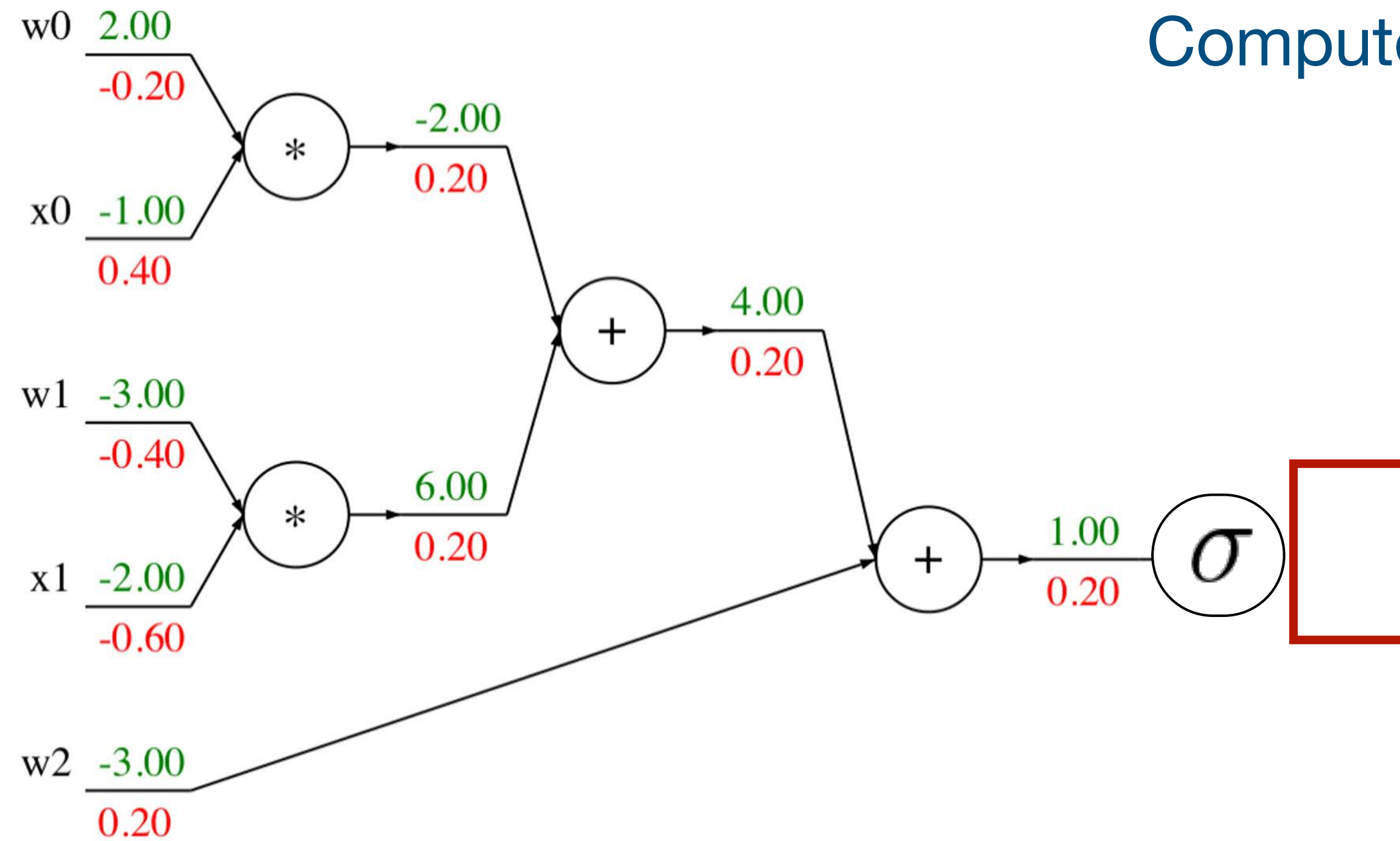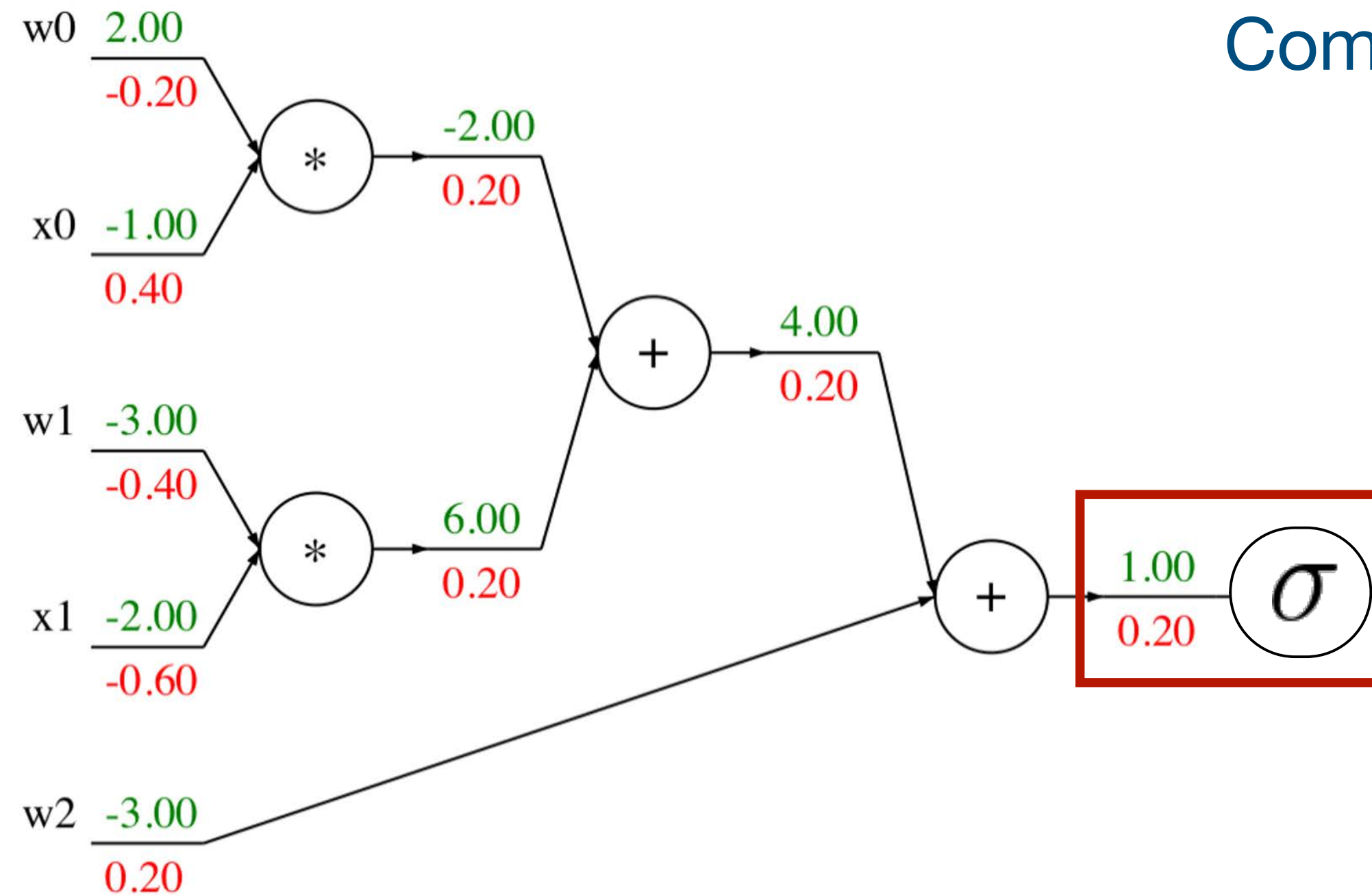
```
###############################################################
# TODO:                                                       #
# Implement a vectorized version of the gradient for the structured SVM    #
# loss, storing the result in dW.                             #
#                                                             #
# Hint: Instead of computing the gradient from scratch, it may be easier    #
# to reuse some of the intermediate values that you used to compute the    #
# loss.                                                       #
###############################################################
# Replace "pass" statement with your code
dmargins = # ...
dscores = # ...
dW = # ...
###############################################################
#                        END OF YOUR CODE                     #
###############################################################
```

$s = Wx$

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

$x$

$W$

$*$ → Hinge loss → $+$ → L

$R$

$R(W)$

# Backprop Implementation: Modular API

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

# Example: PyTorch Autograd Functions

x

*

z

y

(x,y,z are scalars)

```python
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
```

Need to stash some values for use in backward

Upstream gradient

Multiply upstream and local gradients

So far:  backprop with scalars

What about vector-valued functions?

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

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

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

Regular derivative:

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

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

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

Derivative is **Gradient**:

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

$$\left(\frac{\partial y}{\partial x}\right)_i = \frac{\partial y}{\partial x_i}$$

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

70

# 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?

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

Derivative is **Gradient**:

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

$$\left(\frac{\partial y}{\partial x}\right)_i = \frac{\partial y}{\partial x_i}$$

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

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

Derivative is **Jacobian**:

$$\frac{\partial y}{\partial x} \in \mathbb{R}^{N \times M}$$

$$\left(\frac{\partial y}{\partial x}\right)_{i,j} = \frac{\partial y_j}{\partial x_i}$$

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

$\boxed{x}\,D_x$

Loss L still a scalar!

$\boxed{z}\,D_z$

$f$

$\boxed{y}\,D_y$

# Backprop with Vectors

$x$ $D_x$

$y$ $D_y$

$f$

$z$ $D_z$

Loss L still a scalar!

$$\frac{\partial L}{\partial z}$$ $D_z$

Upstream
Gradient

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

74

# Backprop with Vectors

$x$ $D_x$

Local Jacobians (matrices)

Loss L still a scalar!

$$\frac{\partial z}{\partial x}$$ $D_x \times D_z$

$f$

$z$ $D_z$

$$\frac{\partial z}{\partial y}$$ $D_y \times D_z$

$y$ $D_y$

$$\frac{\partial L}{\partial z}$$ $D_z$

Upstream Gradient

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

75

# Backprop with Vectors

$$x \boxed{} D_x$$

Local Jacobians
(matrices)

Loss L still a scalar!

Matrix-vector
multiply

$$D_x \quad \frac{\partial L}{\partial x} = \frac{\partial z}{\partial x} \frac{\partial L}{\partial z}$$

$$\frac{\partial z}{\partial x} \quad D_x \times D_z$$

$$f$$

$$z \boxed{} D_z$$

$$\frac{\partial z}{\partial y} \quad D_y \times D_z$$

Downstream
Gradients

$$y \boxed{} D_y$$

$$D_y \quad \frac{\partial L}{\partial y} = \frac{\partial z}{\partial y} \frac{\partial L}{\partial z}$$

$$\frac{\partial L}{\partial z} \quad D_z$$

Upstream
Gradient

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

76

# Backprop with Vectors

4D input x:

[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

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

4D output y:

[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

# Backprop with Vectors

4D input x:

[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

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

4D output y:

[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

4D dL/dy:

[ 4 ]
[ -1 ]
[ 5 ]
[ 9 ]

Upstream gradient

# Backprop with Vectors

4D input x:
[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

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

4D output y:
[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

[dy/dx] [dL/dy]
[ 1 0 0 0 ] [ 4 ]
[ 0 0 0 0 ] [ -1 ]
[ 0 0 1 0 ] [ 5 ]
[ 0 0 0 0 ] [ 9 ]

4D dL/dy:
[ 4 ]
[ -1 ]
[ 5 ]
[ 9 ]

Upstream gradient

# Backprop with Vectors

4D input x:

[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

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

4D output y:

[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

4D dL/dx:

[ 4 ]
[ 0 ]
[ 5 ]
[ 0 ]

[dy/dx] [dL/dy]

[ 1 0 0 0 ] [ 4 ]
[ 0 0 0 0 ] [ -1 ]
[ 0 0 1 0 ] [ 5 ]
[ 0 0 0 0 ] [ 9 ]

4D dL/dy:

[ 4 ]
[ -1 ]
[ 5 ]
[ 9 ]

Upstream
gradient

# Backprop with Vectors

4D input x:

[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

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

4D output y:

[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

Jacobian is **sparse**: off-diagonal entries all zero!

4D dL/dx:        [dy/dx] [dL/dy]        4D dL/dy:

[ 4 ]    ⟵    [ 1 0 0 0 ] [ 4 ]    ⟵    [ 4 ]  ⟵
[ 0 ]    ⟵    [ 0 0 0 0 ] [ -1 ]   ⟵    [ -1 ] ⟵        Upstream
[ 5 ]    ⟵    [ 0 0 1 0 ] [ 5 ]    ⟵    [ 5 ]  ⟵        gradient
[ 0 ]    ⟵    [ 0 0 0 0 ] [ 9 ]    ⟵    [ 9 ]  ⟵

# Backprop with Vectors

4D input x:

[ 1 ] ⟶
[ -2 ] ⟶
[ 3 ] ⟶
[ -1 ] ⟶

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

4D output y:

⟶ [ 1 ]
⟶ [ 0 ]
⟶ [ 3 ]
⟶ [ 0 ]

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

4D dL/dx:

[ 4 ] ⟵
[ 0 ] ⟵
[ 5 ] ⟵
[ 0 ] ⟵

[dy/dx] [dL/dy]

$$\left(\frac{\partial L}{\partial x}\right)_i = \begin{cases} \left(\frac{\partial L}{\partial y}\right)_i, & if \ x_i > 0 \\ 0, & otherwise \end{cases}$$

4D dL/dy:

⟵ [ 4 ] ⟵
⟵ [ -1 ] ⟵
⟵ [ 5 ] ⟵
⟵ [ 9 ] ⟵

Upstream gradient

# Backprop with Matrices (or Tensors)

$x$ $D_x \times M_x$

$y$ $D_y \times M_y$

$f$

$z$ $D_z \times M_z$

Loss L still a scalar!

dL/dx always has the same shape as x!

# Backprop with Matrices (or Tensors)



$x$ $D_x \times M_x$

$y$ $D_y \times M_y$

$f$

$z$ $D_z \times M_z$

$\dfrac{\partial L}{\partial z}$ $D_z \times M_z$

Upstream Gradient

Loss L still a scalar!

dL/dx always has the same shape as x!

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

# Backprop with Matrices (or Tensors)

$x$   $D_x \times M_x$

Local Jacobians (matrices)

Loss L still a scalar!

dL/dx always has the same shape as x!

$$\frac{\partial z}{\partial x} \quad (D_x \times M_x) \times (D_z \times M_z)$$

$f$

$z$   $D_z \times M_z$

$$\frac{\partial z}{\partial y} \quad (D_y \times M_y) \times (D_z \times M_z)$$

$y$   $D_y \times M_y$

$$\frac{\partial L}{\partial z} \quad D_z \times M_z$$

Upstream Gradient

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

# Backprop with Matrices (or Tensors)



**DR**

$x$   $D_x \times M_x$

Loss L still a scalar!

dL/dx always has the same shape as x!

Local Jacobians (matrices)

Matrix-vector multiply

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

$D_x \times M_x$

Downstream Gradients

$y$   $D_y \times M_y$

$\frac{\partial z}{\partial x}$   $(D_x \times M_x) \times (D_z \times M_z)$

$\frac{\partial z}{\partial y}$   $(D_y \times M_y) \times (D_z \times M_z)$

$f$

$z$   $D_z \times M_z$

$\frac{\partial L}{\partial z}$   $D_z \times M_z$

$D_y \times M_y$

$$\frac{\partial L}{\partial y} = \frac{\partial z}{\partial y} \frac{\partial L}{\partial z}$$

Upstream Gradient

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

86

# Example: Matrix Multiplication

x: [N×D]

[ 2  1 -3 ]
[-3  4  2 ]

w: [D×M]

[ 3 2 1 -1]
[ 2 1 3  2]
[ 3 2 1 -2]

Matrix Multiply *y = xw*

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

y: [N×M]

[-1 -1  2  6 ]
[ 5  2 11  7 ]

# Example: Matrix Multiplication

x: [N×D]

[ 2  1  -3 ]
[ -3  4  2 ]

w: [D×M]

[ 3  2  1  -1]
[ 2  1  3  2]
[ 3  2  1  -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

y: [N×M]

[-1 -1  2  6 ]
[ 5  2  11  7 ]

dL/dy: [N×M]

[ 2  3  -3  9 ]
[ -8  1  4  6 ]

dL/dx: [N×D]

[ ?  ?  ?  ]
[ ?  ?  ?  ]

# Example: Matrix Multiplication

y: [N×M]
[-1 -1  2  6 ]
[ 5  2 11  7 ]

x: [N×D]
[ 2  1 -3 ]
[-3  4  2 ]

w: [D×M]
[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

dL/dy: [N×M]
[ 2  3 -3  9 ]
[-8  1  4  6 ]

dL/dx: [N×D]
[ ?  ?  ?  ]
[ ?  ?  ?  ]

**Jacobians**:
dy/dx: [(N×D)×(N×M)]
dy/dw: [(D×M)×(N×M)]

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

y: [N×M]

[-1 -1  2  6 ]
[ 5  2 11  7 ]

x: [N×D]

[ 2  1 -3 ]
[-3  4  2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

dL/dy: [N×M]

[ 2  3 -3  9 ]
[-8  1  4  6 ]

dL/dx: [N×D]

[ ?  ?  ? ]
[ ?  ?  ? ]

Local Gradient Slice:

dy/dx$_{1,1}$
[ ? ? ? ? ]
[ ? ? ? ? ]

dL/dx$_{1,1}$
= (dy/dx$_{1,1}$) · (dL/dy)

# Example: Matrix Multiplication

x: [N×D]

[ 2  1 -3 ]
[ -3  4  2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_{k} x_{i,k} w_{k,j}$$

[-1 -1  2  6 ]
[ 5  2 11  7 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[ -8  1  4  6 ]

dL/dx: [N×D]

[ ?  ?  ? ]
[ ?  ?  ? ]

**Local Gradient Slice:**

$dy/dx_{1,1}$

$dy_{1,1}/dx_{1,1}$   [ ?  ?  ?  ? ]
[ ?  ?  ?  ? ]

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

# Example: Matrix Multiplication

x: [N×D]

[ 2  1 -3]
[-3  4  2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

[-1 -1  2  6 ]
[ 5  2 11  7 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[-8  1  4  6 ]

dL/dx: [N×D]

[  ?   ?   ?  ]
[  ?   ?   ?  ]

Local Gradient Slice:

dy/dx$_{1,1}$

dy$_{1,1}$/dx$_{1,1}$    [  ?  ?  ?  ? ]
[ ?  ?  ?  ? ]

dL/dx$_{1,1}$
= (dy/dx$_{1,1}$) · (dL/dy)

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

# Example: Matrix Multiplication

x: [N×D]

[ 2  1 -3]
[ -3  4  2]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

[ -1 -1  2  6 ]
[ 5  2  11  7 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[ -8  1  4  6 ]

dL/dx: [N×D]

[ ?  ?  ? ]
[ ?  ?  ? ]

dL/dx$_{1,1}$
= (dy/dx$_{1,1}$) · (dL/dy)

Local Gradient Slice:

dy/dx$_{1,1}$

dy$_{1,1}$/dx$_{1,1}$  [ 3  ?  ?  ? ]
[ ?  ?  ?  ? ]

$y_{1,1} = x_{1,1}w_{1,1} + x_{1,2}w_{2,1} + x_{1,3}w_{3,1}$
$\Rightarrow$ dy$_{1,1}$/dx$_{1,1}$ = $w_{1,1}$

# Example: Matrix Multiplication

x: [N×D]

[ 2  1 -3 ]
[ -3  4  2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

[-1 -1  2  6 ]
[ 5  2  11  7 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[ -8  1  4  6 ]

dL/dx: [N×D]

[ ?  ?  ? ]
[ ?  ?  ? ]

Local Gradient Slice:

dy/dx$_{1,1}$

dy$_{1,2}$/dx$_{1,1}$   [ 3  ?  ?  ? ]
[ ?  ?  ?  ? ]

dL/dx$_{1,1}$
= (dy/dx$_{1,1}$) · (dL/dy)

# Example: Matrix Multiplication

x: [N×D]

[ 2  1 -3]
[ -3  4  2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

[-1 -1  2  6 ]
[ 5  2 11  7 ]

dL/dy: [N×M]

[ 2 3 -3 9 ]
[ -8 1 4 6 ]

dL/dx: [N×D]

[ ?  ?  ? ]
[ ?  ?  ? ]

dL/dx$_{1,1}$
= (dy/dx$_{1,1}$) · (dL/dy)

Local Gradient Slice:

dy/dx$_{1,1}$

dy$_{1,2}$/dx$_{1,1}$  [ 3  ?  ?  ? ]
[ ?  ?  ?  ? ]

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

# Example: Matrix Multiplication

x: [N×D]
[ 2  1 -3]
[ -3  4  2 ]

w: [D×M]
[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

[-1 -1  2  6 ]
[ 5  2  11  7 ]

dL/dy: [N×M]
[ 2  3 -3  9 ]
[ -8  1  4  6 ]

dL/dx: [N×D]
[ ?  ?  ?  ]
[ ?  ?  ?  ]

Local Gradient Slice:

$dy/dx_{1,1}$
$dy_{1,2}/dx_{1,1}$ [ 3  2  ?  ? ]
[ ?  ?  ?  ? ]

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

$y_{1,2} = x_{1,1}w_{1,2} + x_{1,2}w_{2,2} + x_{1,3}w_{3,2}$
$=> dy_{1,2}/dx_{1,1} = w_{1,2}$

# Example: Matrix Multiplication

x: [N×D]

[ 2  1 -3 ]
[ -3  4  2 ]

w: [D×M]

[ 3 2 1 -1]
[ 2 1 3 2]
[ 3 2 1 -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

[-1 -1  2  6 ]
[ 5  2 11  7 ]

dL/dy: [N×M]

[ 2 3 -3 9 ]
[ -8 1 4 6 ]

dL/dx: [N×D]

[ ?  ?  ? ]
[ ?  ?  ? ]

Local Gradient Slice:

dy/dx$_{1,1}$

dy$_{1,2}$/dx$_{1,1}$  [ 3  2  1 -1 ]
[ ? ? ? ? ]

dL/dx$_{1,1}$
= (dy/dx$_{1,1}$) · (dL/dy)

# Example: Matrix Multiplication

x: [N×D]

[ 2  1 -3 ]
[ -3  4  2]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

[ -1 -1  2  6 ]
[ 5  2 11  7 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[ -8  1  4  6 ]

dL/dx: [N×D]

[ ?  ?  ? ]
[ ?  ?  ? ]

**Local Gradient Slice:**

$dy/dx_{1,1}$

$dy_{1,2}/dx_{1,1}$  [ 3  2  1 -1 ]

[ ? ? ? ? ]

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

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

# Example: Matrix Multiplication

x: [N×D]

[ 2  1 -3 ]
[-3  4  2]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

$$\text{Matrix Multiply } y = xw$$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

[-1 -1  2  6 ]
[ 5  2  11  7 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[-8  1  4  6 ]

dL/dx: [N×D]

[  ?   ?   ? ]
[  ?   ?   ? ]

Local Gradient Slice:

$dy/dx_{1,1}$

$dy_{1,2}/dx_{1,1}$  [ 3  2  1 -1 ]

[ 0  ?  ?  ? ]

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

$y_{2,1} = x_{2,1}w_{1,1} + x_{2,2}w_{2,1} + x_{2,3}w_{3,1}$
$\Rightarrow dy_{2,1}/dx_{1,1} = 0$

# Example: Matrix Multiplication

x: [N×D]

[ 2  1 -3 ]
[ -3  4  2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

[ -1 -1  2  6 ]
[ 5  2  11  7 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[ -8  1  4  6 ]

dL/dx: [N×D]

[ ?  ?  ? ]
[ ?  ?  ? ]

Local Gradient Slice:

dy/dx$_{1,1}$

dy$_{1,2}$/dx$_{1,1}$  [ 3  2  1 -1 ]
[ 0  0  0  0 ]

dL/dx$_{1,1}$
= (dy/dx$_{1,1}$) · (dL/dy)

# Example: Matrix Multiplication

x: [N×D]

[ 2  1 -3 ]
[ -3  4  2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

[-1 -1  2  6 ]
[ 5  2  11  7 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[ -8  1  4  6 ]

dL/dx: [N×D]

[ ?  ?  ? ]
[ ?  ?  ? ]

Local Gradient Slice:

dy/dx_{1,1}
[ 3  2  1 -1 ]
[ 0  0  0  0 ]

dL/dx_{1,1}
= (dy/dx_{1,1}) · (dL/dy)

# Example: Matrix Multiplication

x: [N×D]

[ 2  1 -3 ]
[ -3  4  2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

[-1 -1  2  6 ]
[ 5  2  11  7 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[ -8  1  4  6 ]

dL/dx: [N×D]

[ 0  ?  ? ]
[ ?  ?  ? ]

Local Gradient Slice:

dy/dx$_{1,1}$
[ 3  2  1 -1 ]
[ 0  0  0  0 ]

dL/dx$_{1,1}$
= (dy/dx$_{1,1}$) · (dL/dy)
= (w$_{1,:}$) · (dL/dy$_{1,:}$)
= 3*2 + 2*3 + 1*(-3) + (-1)*9 = 0

# Example: Matrix Multiplication

x: [N×D]

[ 2  1 -3 ]
[-3  4  2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply *y = xw*

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

[-1 -1  2  6 ]
[ 5  2 11  7 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[-8  1  4  6 ]

dL/dx: [N×D]

[ 0   ?   ?  ]
[ ?   ?  -30 ]

Local Gradient Slice:

dy/dx$_{2,3}$

[ 0  0  0  0 ]
[ 3  2  1 -2 ]

dL/dx$_{2,3}$
= (dy/dx$_{2,3}$) · (dL/dy)

# Example: Matrix Multiplication

x: [N×D]

[ 2   1  -3 ]
[-3   4   2 ]

w: [D×M]

[ 3   2   1  -1]
[ 2   1   3   2]
[ 3   2   1  -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

[-1 -1   2   6 ]
[ 5   2  11   7 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[-8  1  4  6 ]

dL/dx: [N×D]

[ 0   ?   ?  ]
[ ?   ?  -30 ]

Local Gradient Slice:

dy/dx$_{2,3}$

[ 0  0  0  0 ]
[ 3  2  1 -2 ]

dL/dx$_{2,3}$
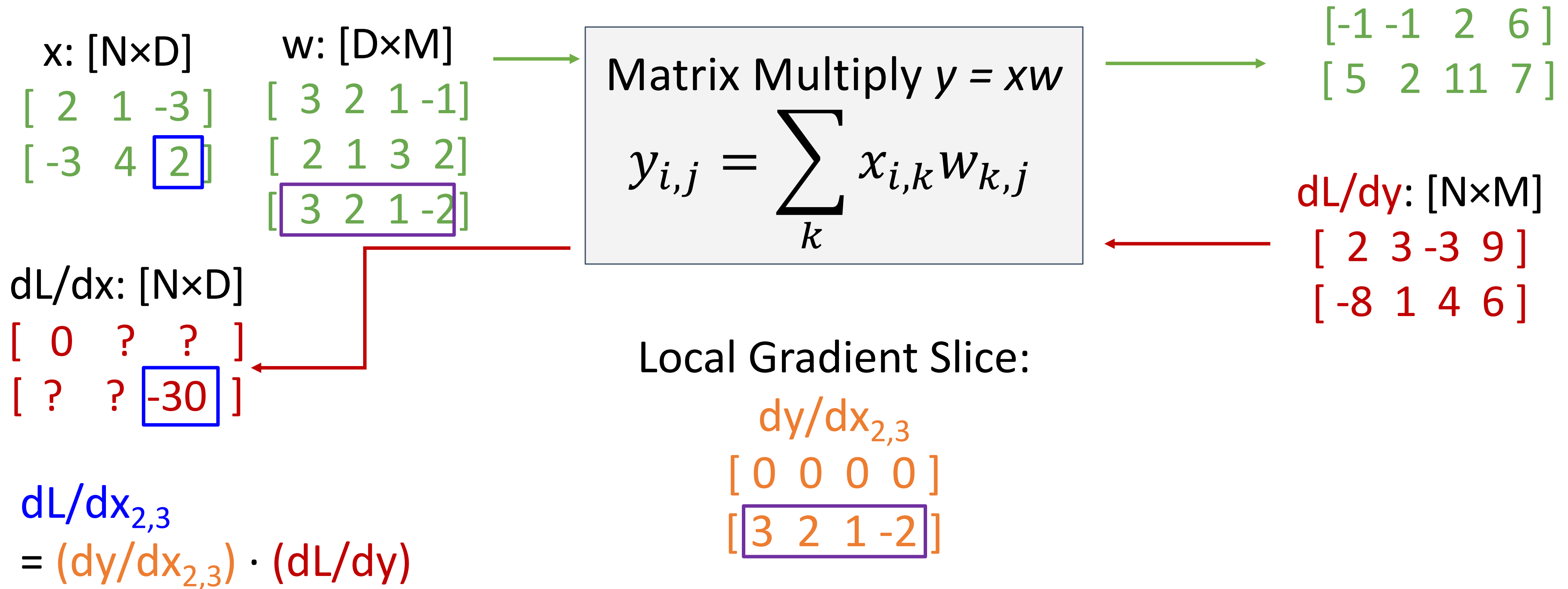= (dy/dx$_{2,3}$) · (dL/dy)
= (w$_{3,:}$) · (dL/dy$_{2,:}$)
= 3*(-8) + 2*1 + 1*4 + (-2)*6 = -30

# Example: Matrix Multiplication

x: [N×D]

[ 2  1 -3 ]

[-3  4  2 ]

w: [D×M]

[ 3  2  1 -1]

[ 2  1  3  2]

[ 3  2  1 -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

[-1 -1  2  6 ]

[ 5  2 11  7 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]

[-8  1  4  6 ]

dL/dx: [N×D]

[ 0  16  -9 ]

[-24  9 -30 ]

$dL/dx_{i,j}$

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

$= (w_{j,:}) \cdot (dL/dy_{i,:})$

# Example: Matrix Multiplication

x: [N×D]

[ 2  1 -3 ]

[-3  4  2 ]

w: [D×M]

[ 3  2  1 -1]

[ 2  1  3  2]

[ 3  2  1 -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

[-1 -1  2  6 ]

[ 5  2 11  7 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]

[-8  1  4  6 ]

dL/dx: [N×D]

[ 0  16  -9 ]

[-24  9 -30 ]

dL/dx = (dL/dy) w$^T$

[N x D]      [N x M]  [M x D]

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

dL/dx$_{i,j}$

= (dy/dx$_{i,j}$) · (dL/dy)

= (w$_{j,:}$) · (dL/dy$_{i,:}$)

# Example: Matrix Multiplication

x: [N×D]

[ 2  1 -3 ]

[-3  4  2 ]

w: [D×M]

[ 3 2 1 -1]

[ 2 1 3 2]

[ 3 2 1 -2]

Matrix Multiply $y = xw$

$$y_{i,j} = \sum_k x_{i,k} w_{k,j}$$

[-1 -1  2  6 ]

[ 5  2 11  7 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]

[-8  1  4  6 ]

dL/dx: [N×D]

[ 0  16  -9 ]
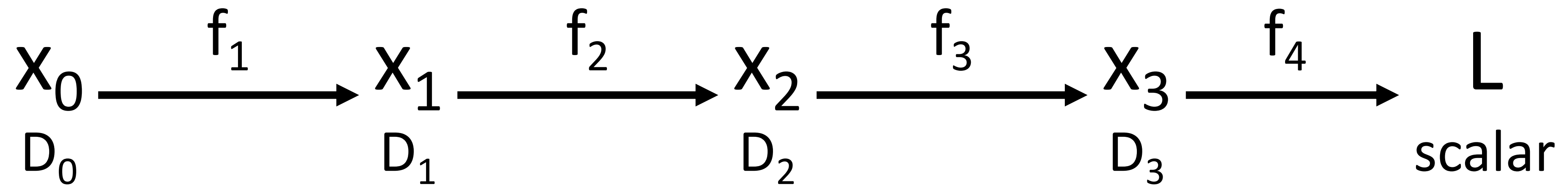
[-24  9 -30 ]

dL/dx = (dL/dy) w$^T$

[N x D]      [N x M]  [M x D]

dL/dw = x$^T$ (dL/dy)

[D x M]   [D x N] [N x M]

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

# Backpropagation: Another View

$$x_0 \xrightarrow{f_1} x_1 \xrightarrow{f_2} x_2 \xrightarrow{f_3} x_3 \xrightarrow{f_4} L$$

$D_0$      $D_1$      $D_2$      $D_3$      scalar

Chain rule

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

# Backpropagation: Another View

$$x_0 \xrightarrow{\quad f_1 \quad} x_1 \xrightarrow{\quad f_2 \quad} x_2 \xrightarrow{\quad f_3 \quad} x_3 \xrightarrow{\quad f_4 \quad} L$$

$D_0 \qquad\qquad D_1 \qquad\qquad D_2 \qquad\qquad D_3 \qquad\qquad$ scalar

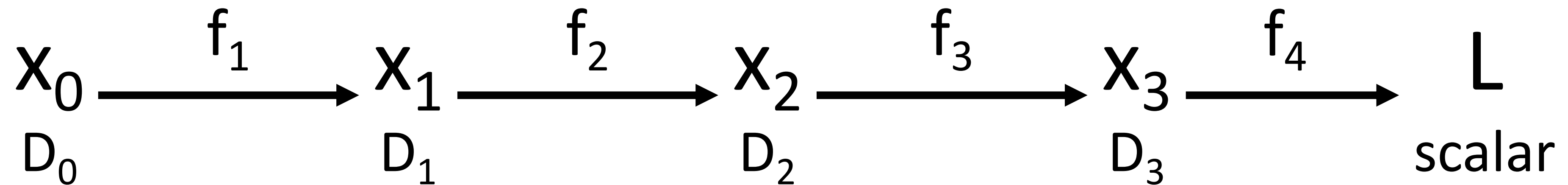Matrix multiplication is **associative**: we can compute products in any order

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

$$x_0 \xrightarrow{f_1} x_1 \xrightarrow{f_2} x_2 \xrightarrow{f_3} x_3 \xrightarrow{f_4} L$$

$D_0 \qquad\qquad D_1 \qquad\qquad D_2 \qquad\qquad D_3 \qquad\qquad$ scalar

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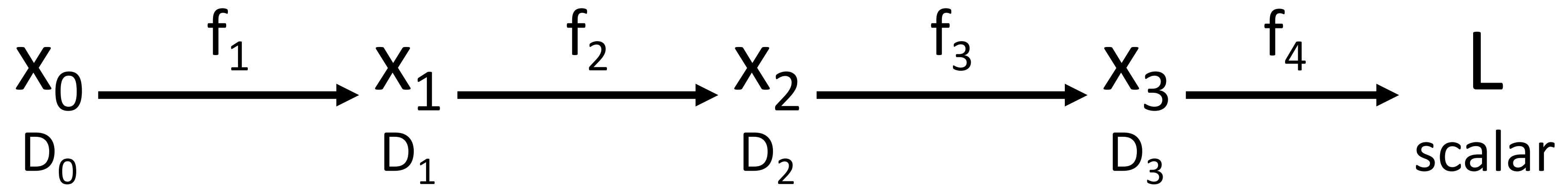Computing products right-to-left avoids matrix-matrix products; only needs matrix-vector

Chain
rule

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

$[D_0 \times D_1] \quad [D_1 \times D_2] \quad [D_2 \times D_3] \quad [D_3]$

$$x_0 \xrightarrow{f_1} x_1 \xrightarrow{f_2} x_2 \xrightarrow{f_3} x_3 \xrightarrow{f_4} L$$

$D_0 \qquad D_1 \qquad D_2 \qquad D_3 \qquad$ scalar

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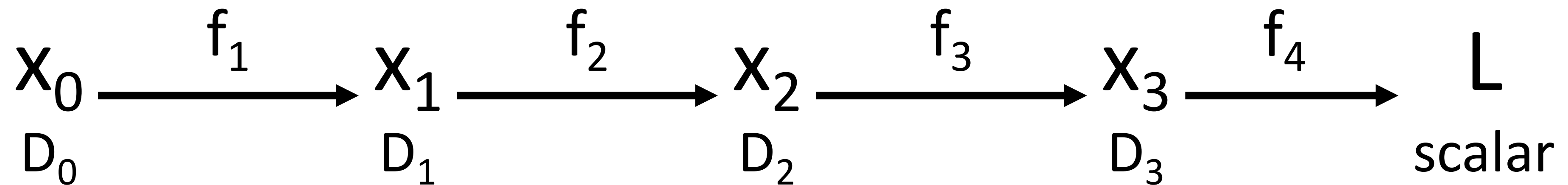
Chain rule

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

Compute grad of scalar <u>output</u> w/respect to all vector <u>inputs</u>

$[D_0 \times D_1] \quad [D_1 \times D_2] \quad [D_2 \times D_3] \quad [D_3]$

# Reverse-Mode Automatic Differentiation

$$x_0 \xrightarrow{f_1} x_1 \xrightarrow{f_2} x_2 \xrightarrow{f_3} x_3 \xrightarrow{f_4} L$$

$D_0 \qquad\qquad D_1 \qquad\qquad D_2 \qquad\qquad D_3 \qquad\qquad$ scalar

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

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

Chain rule

$$\frac{\partial L}{\partial x_0} = \left(\frac{\partial x_1}{\partial x_0}\right)\left(\frac{\partial x_2}{\partial x_1}\right)\left(\frac{\partial x_3}{\partial x_2}\right)\left(\frac{\partial L}{\partial x_3}\right)$$
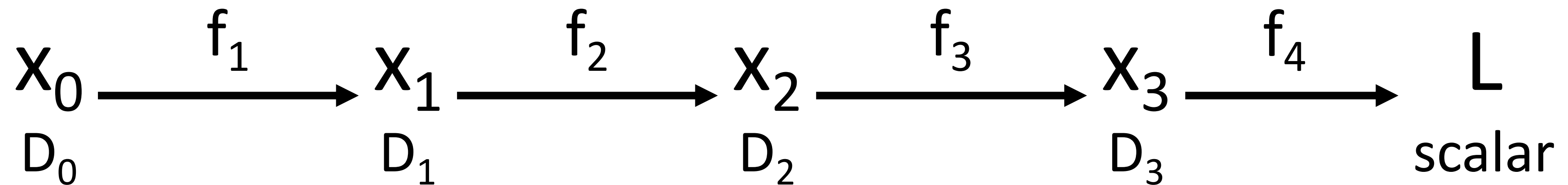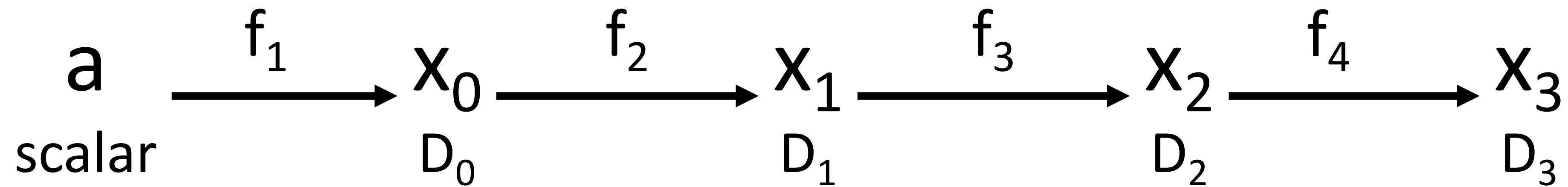
What if we want grads of scalar <u>input</u> w/respect to vector <u>outputs</u>?

Compute grad of scalar <u>output</u> w/respect to all vector <u>inputs</u>

$[D_0 \times D_1] \quad [D_1 \times D_2] \quad [D_2 \times D_3] \quad [D_3]$

# Forward-Mode Automatic Differentiation

$$a \xrightarrow{f_1} x_0 \xrightarrow{f_2} x_1 \xrightarrow{f_3} x_2 \xrightarrow{f_4} x_3$$

a
scalar

$x_0$
$D_0$

$x_1$
$D_1$

$x_2$
$D_2$

$x_3$
$D_3$

Chain
rule

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

$[D_0]$   $[D_0 \times D_1]$   $[D_1 \times D_2]$   $[D_2 \times D_3]$

# Forward-Mode Automatic Differentiation

$$a \xrightarrow{f_1} x_0 \xrightarrow{f_2} x_1 \xrightarrow{f_3} x_2 \xrightarrow{f_4} x_3$$

scalar          $D_0$          $D_1$          $D_2$          $D_3$

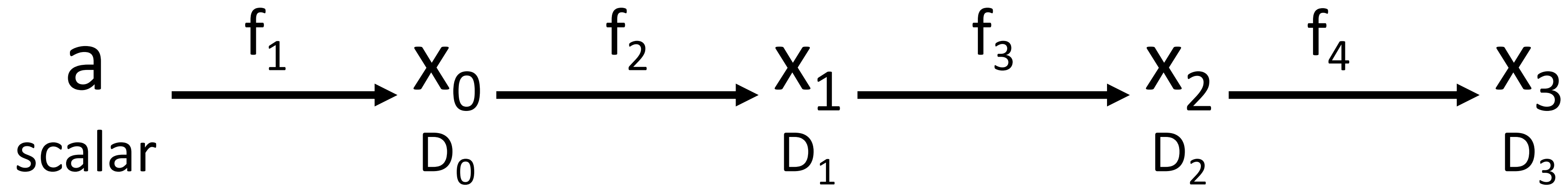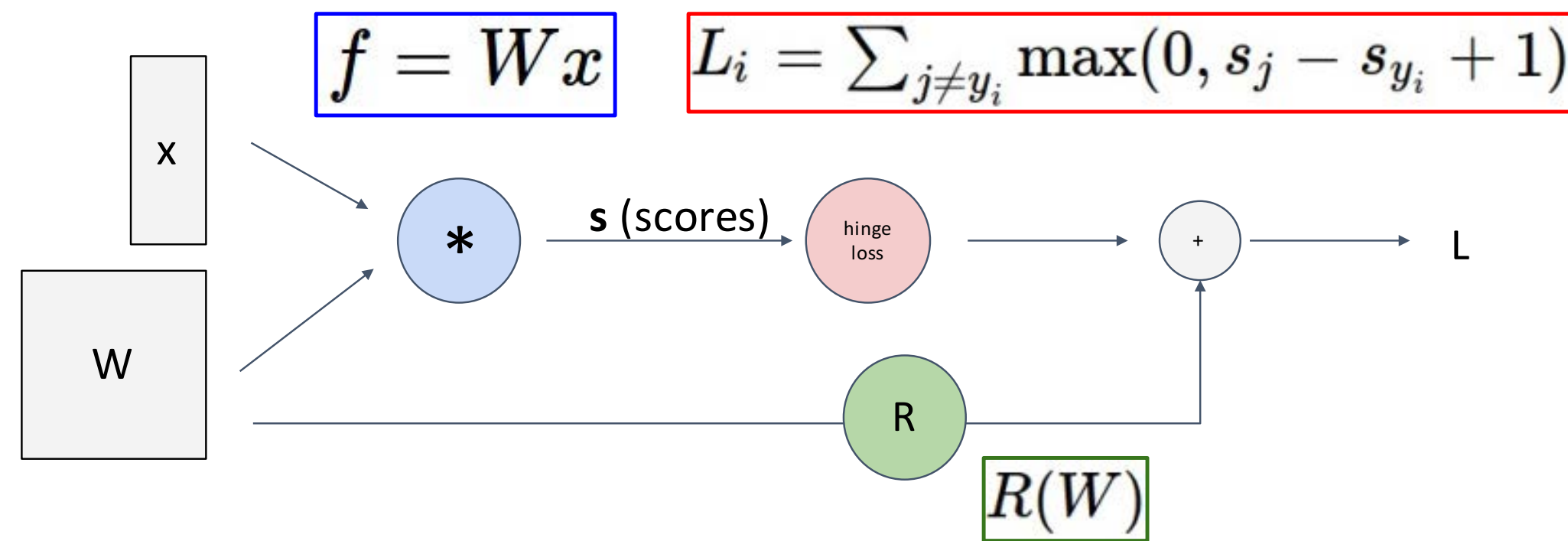Computing products <u>left-to-right</u> avoids matrix-matrix products; only needs matrix-vector

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

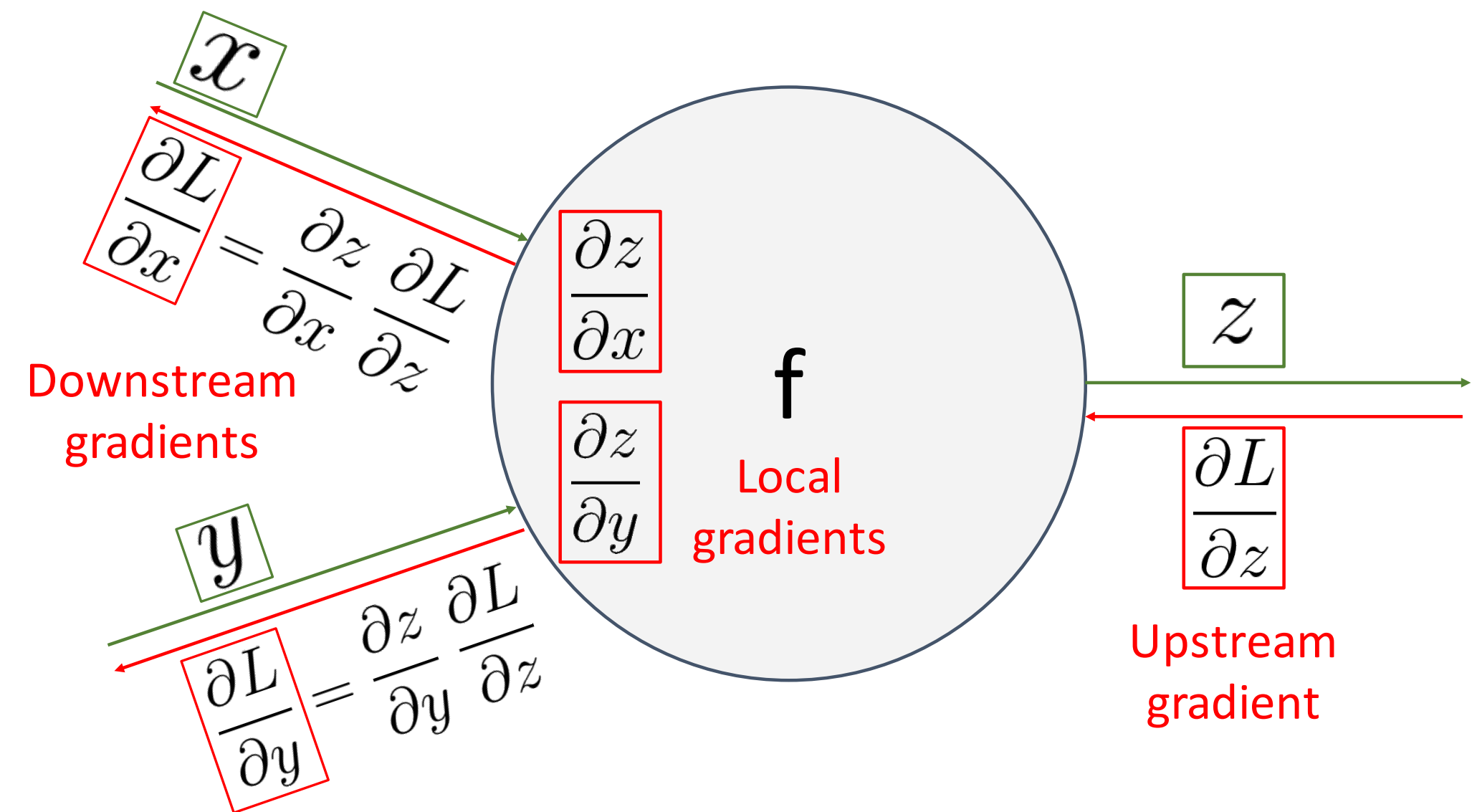$[D_0] \quad [D_0 \times D_1] \quad [D_1 \times D_2] \quad [D_2 \times D_3]$

# Summary

Represent complex expressions as **computational graphs**

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

$$f = Wx$$

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

x

W

*

**s** (scores)

hinge loss

+

L

R

$$R(W)$$

Forward pass computes outputs

Backward pass computes gradients

$$x$$

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

Downstream gradients

$$\frac{\partial z}{\partial x}$$

$$\frac{\partial z}{\partial y}$$

**f**

Local gradients

$$z$$

$$\frac{\partial L}{\partial z}$$

Upstream gradient

$$y$$

$$\frac{\partial L}{\partial y} = \frac{\partial z}{\partial y} \frac{\partial L}{\partial z}$$

# Summary

Backprop can be implemented with "flat" code where the backward pass looks like forward pass reversed
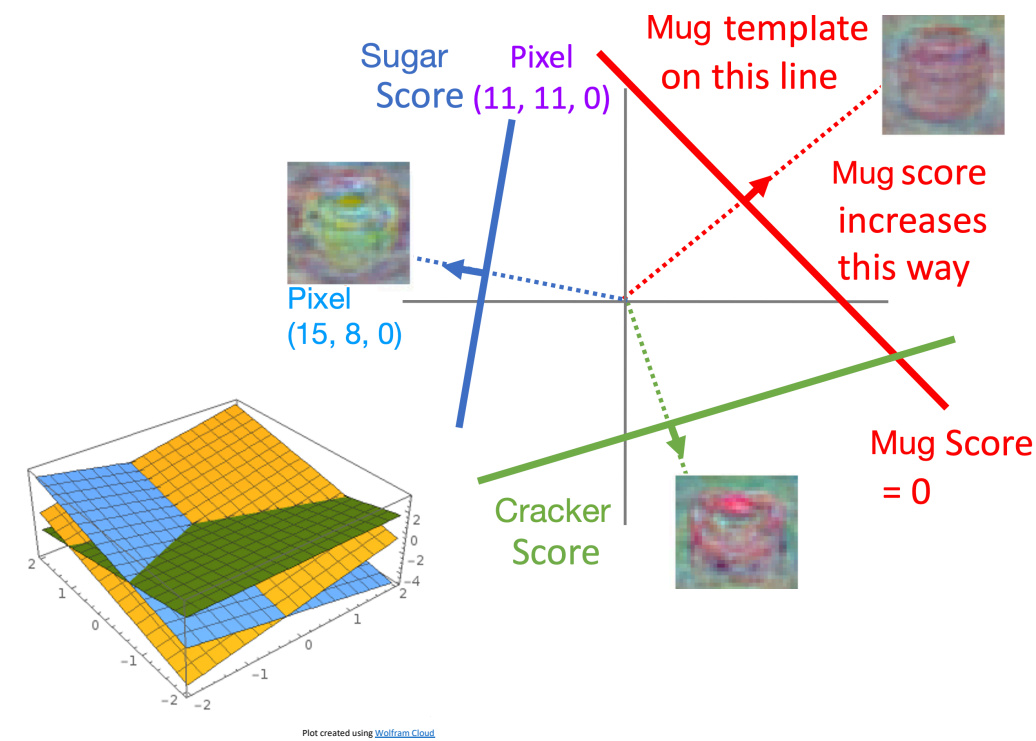
```
def f(w0, x0, w1, x1, w2):
  s0 = w0 * x0
  s1 = w1 * x1
  s2 = s0 + s1
  s3 = s2 + w2
  L = sigmoid(s3)

  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

```
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
```
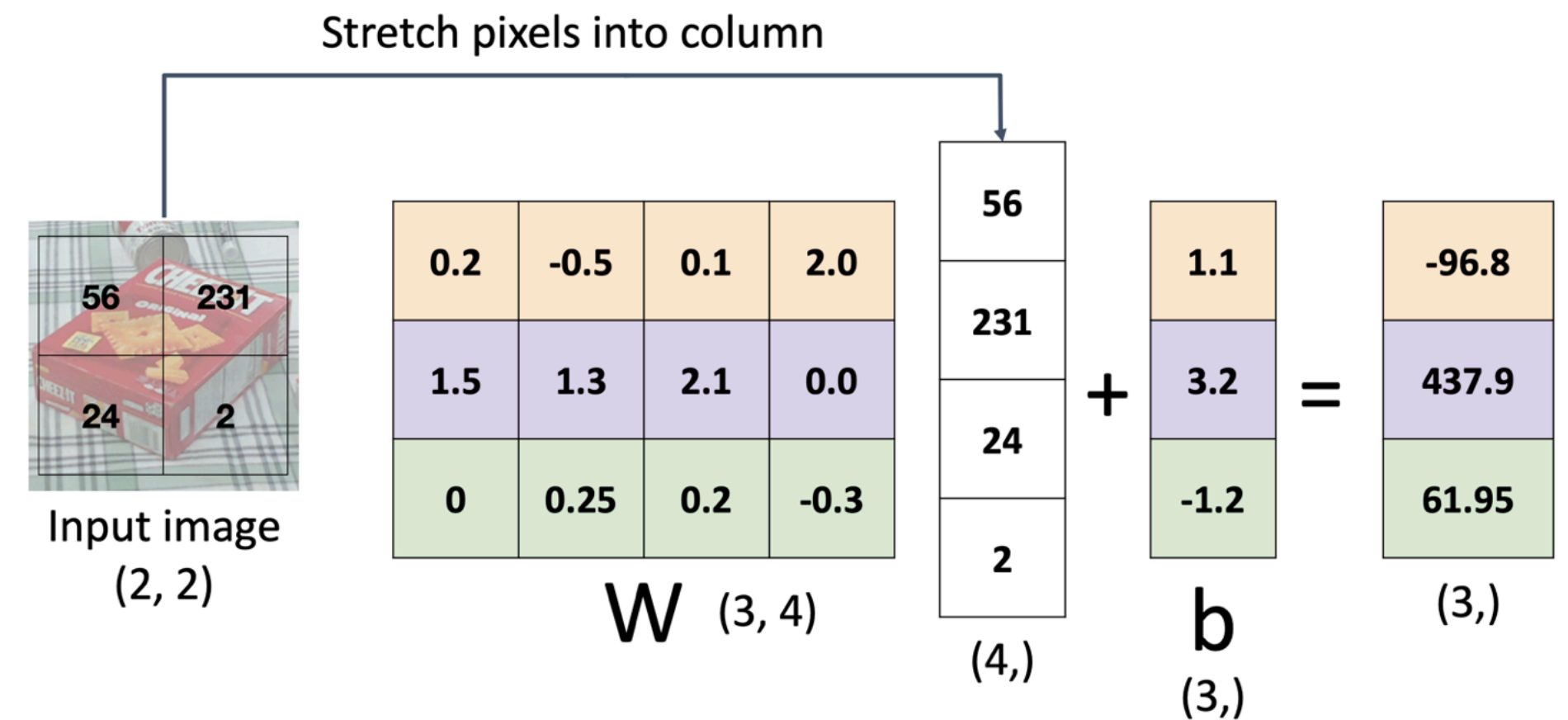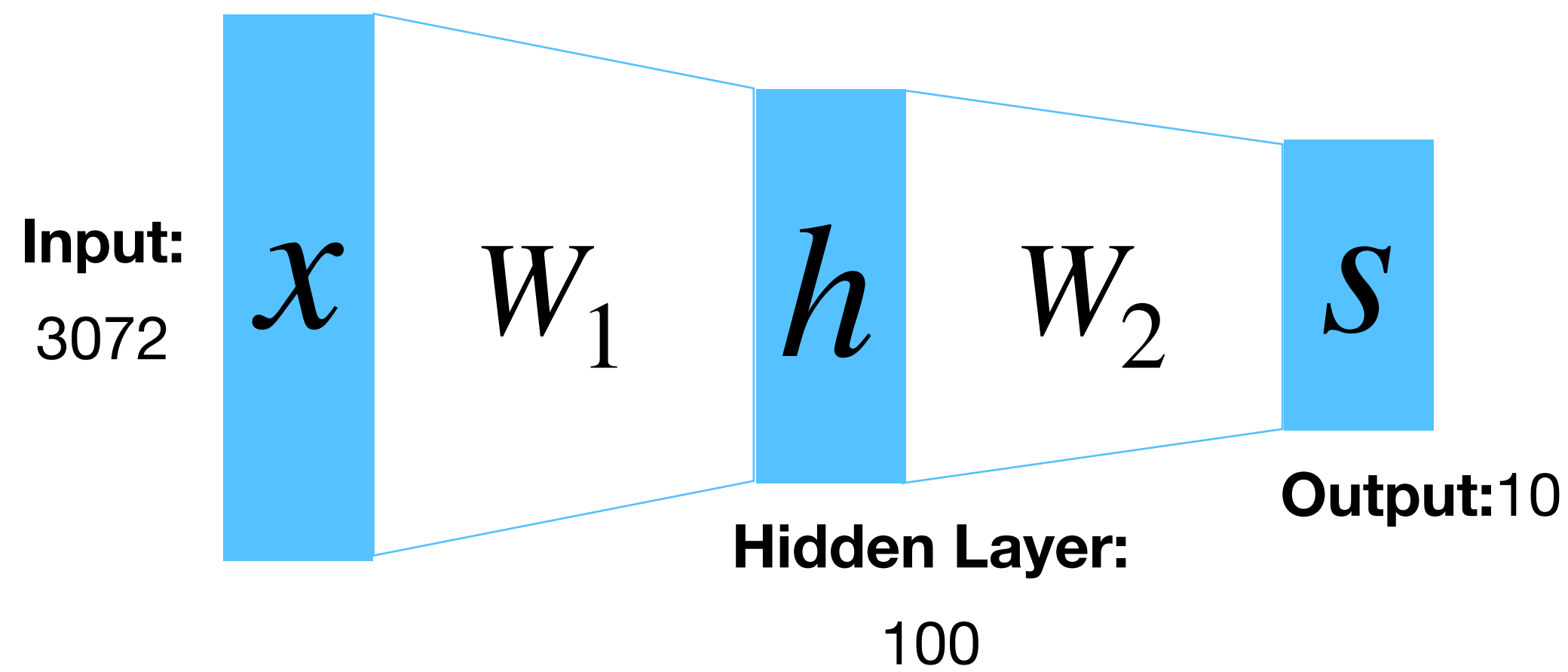
# Summary



$$f(x) = W_2 \max(0, W_1 x + b_1) + b_2$$

**Input:** 3072

$x$  $W_1$  $h$  $W_2$  $s$

**Hidden Layer:** 100

**Output:** 10

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

Stretch pixels into column



Input image (2, 2)

| 0.2 | -0.5 | 0.1 | 2.0 |
| 1.5 | 1.3 | 2.1 | 0.0 |
| 0 | 0.25 | 0.2 | -0.3 |

W (3, 4)

| 56 |
| 231 |
| 24 |
| 2 |

(4,)

+

| 1.1 |
| 3.2 |
| -1.2 |

b (3,)

=

| -96.8 |
| 437.9 |
| 61.95 |

(3,)

**DR**

# Next time: Convolutional Neural Networks

# Next: Individual brainstorming task

- Pick one of the 3 papers and write a) one 1-page summary, b) one thing from this paper that you can use for your ideated task.
- Data:
  - Where will the data for your work come from?
    - Real-robot vs. Simulation environment?
    - Existing dataset or new data collection?
    - Benchmarking task?
- Network:
  - What is the network architecture you found to be suitable from reading the papers?
    - ResNet, PointNet, Transformers
  - What training strategy would be applied?
    - Supervised, Self-supervised, Semi-supervised
  - Is there an existing code that you can build on?
- Compute:
  - What compute requirements do you have? Memory, GPU etc.
  - What compute resources do you have? Compute heavy laptop/desktop, MSI, etc.
- Evaluation:
  - How do you know if a method could solve your problem?
    - Baselines from existing literature or your own baselines
  - How do you measure how well your method works?
    - Evaluation metrics (existing ones vs. new ones)
  - How will you choose hyperparameters in your project?
    - Ablation study

# DeepRob

**Lecture 6**
**Backpropagation**
**University of Minnesota**

$$\frac{\partial L}{\partial W_{\ell_1}} \longleftarrow \frac{\partial L}{\partial W_{\ell_2}} \longleftarrow \frac{\partial L}{\partial W_{\ell_3}} \longleftarrow \frac{\partial L}{\partial W_{\ell_4}} \longleftarrow \frac{\partial L}{\partial W_{\ell_5}} \longleftarrow \frac{\partial L}{\partial \text{Out}}$$