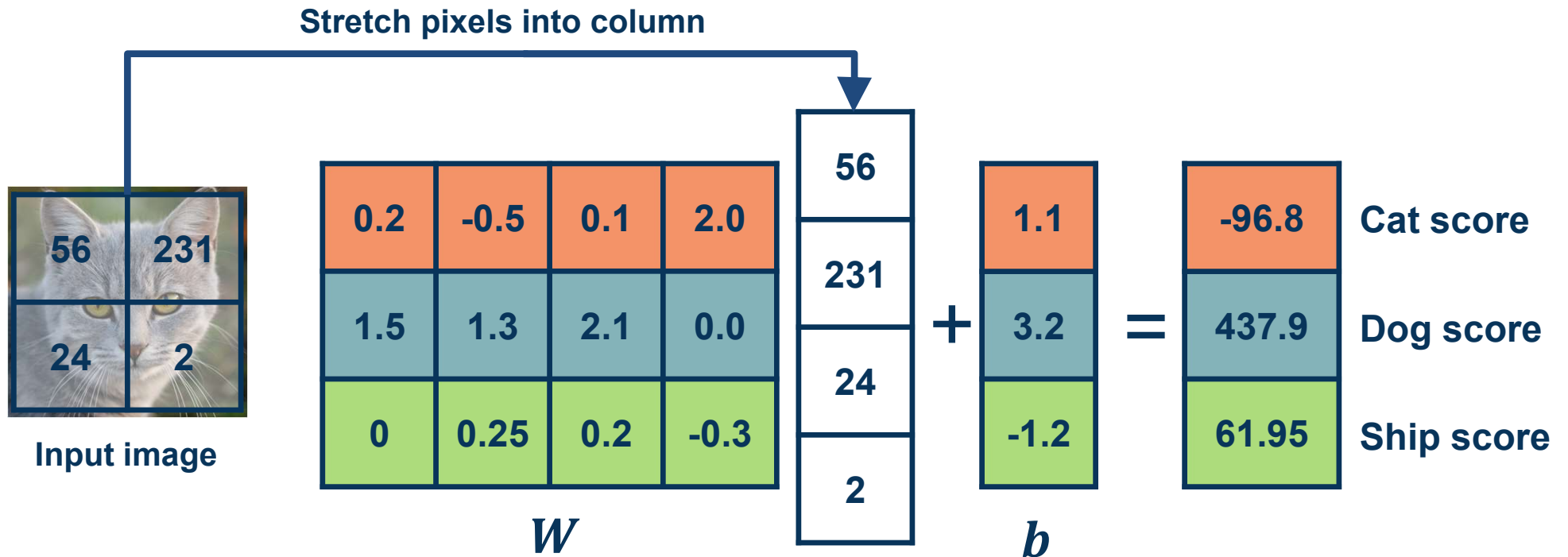Topics:

- Backpropagation
- Matrix/Linear Algebra view

**CS 4644-DL / 7643-A**
**ZSOLT KIRA**

- **Assignment 1 out!**
  - **Due Feb 3rd (with grace period 5$^{th}$)**
  - Start now, start now, start now!
  - Start now, start now, start now!
  - Start now, start now, start now!

- Resources:
  - These lectures
  - Matrix calculus for deep learning
  - Gradients notes and MLP/ReLU Jacobian notes.
  - Assignment 1 (@67) and matrix calculus (@86), convex optimization (@89)

- Piazza: Project teaming thread
  - Project proposal overview during my OH (Thursday 3pm ET)

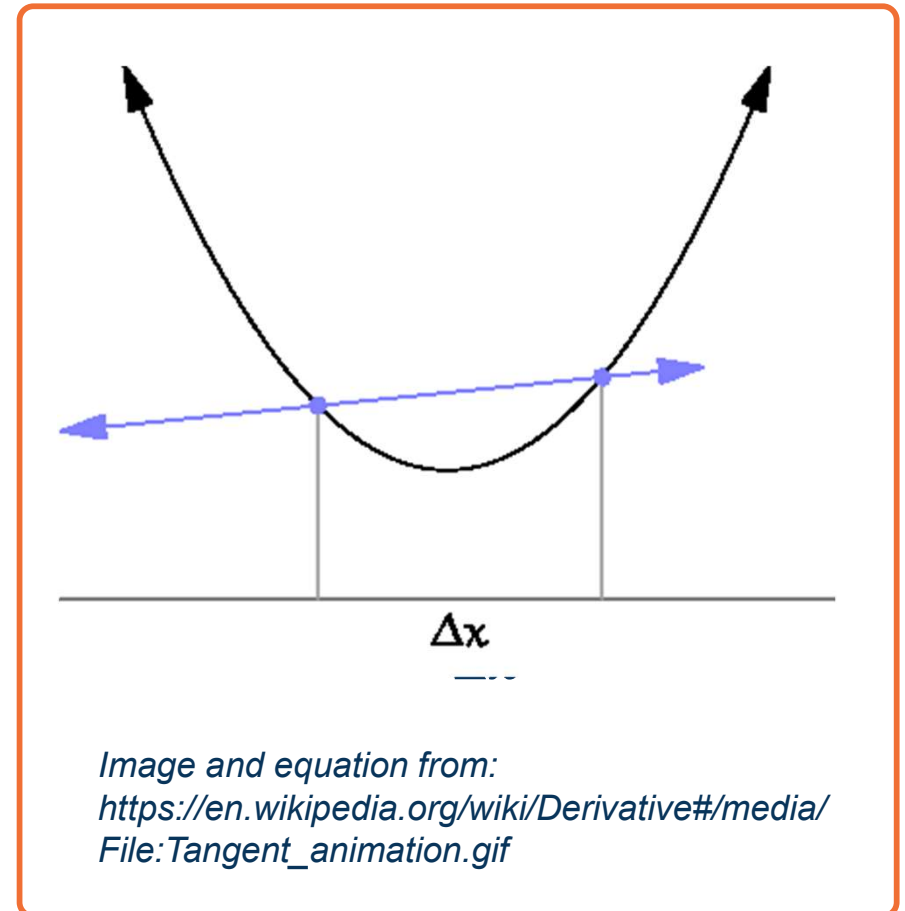# Example with an image with **4 pixels**, and **3 classes (cat/dog/ship)**

Stretch pixels into column



Input image

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

$W$

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

$+$

| 1.1 |
|-----|
| 3.2 |
| -1.2 |

$b$

$=$

| -96.8 | Cat score |
|-------|-----------|
| 437.9 | Dog score |
| 61.95 | Ship score |

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

**Example**

Georgia Tech

- We can find the steepest descent direction by computing the **derivative (gradient):**

$$f'(a) = \lim_{h \to 0} \frac{f(a+h) - f(a)}{h}$$

- Steepest descent direction is the **negative gradient**

- **Intuitively:** Measures how the function changes as the argument a changes by a small step size

  - As step size goes to zero

- **In Machine Learning:** Want to know how the **loss function** changes **as weights** are varied

  - Can consider each parameter separately by taking **partial derivative** of loss function with respect to that parameter



*Image and equation from: https://en.wikipedia.org/wiki/Derivative#/media/File:Tangent_animation.gif*

**Derivatives**

The same two-layered neural network **corresponds to adding another weight matrix**

- We will prefer the linear algebra view, but use some terminology from neural networks (& biology)
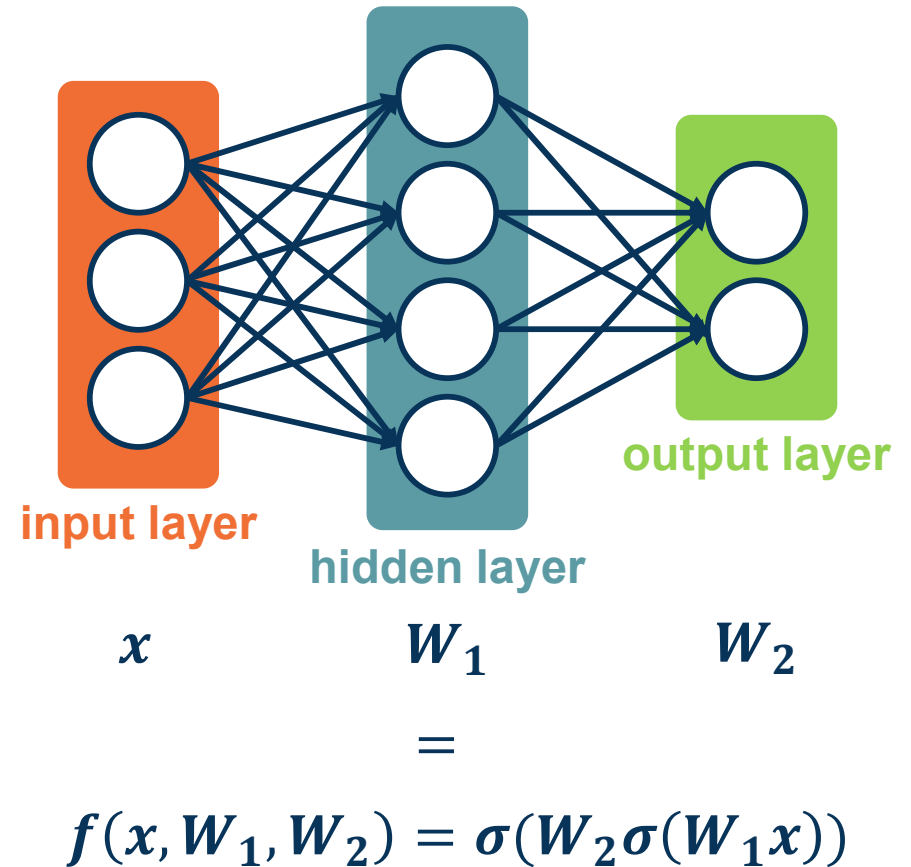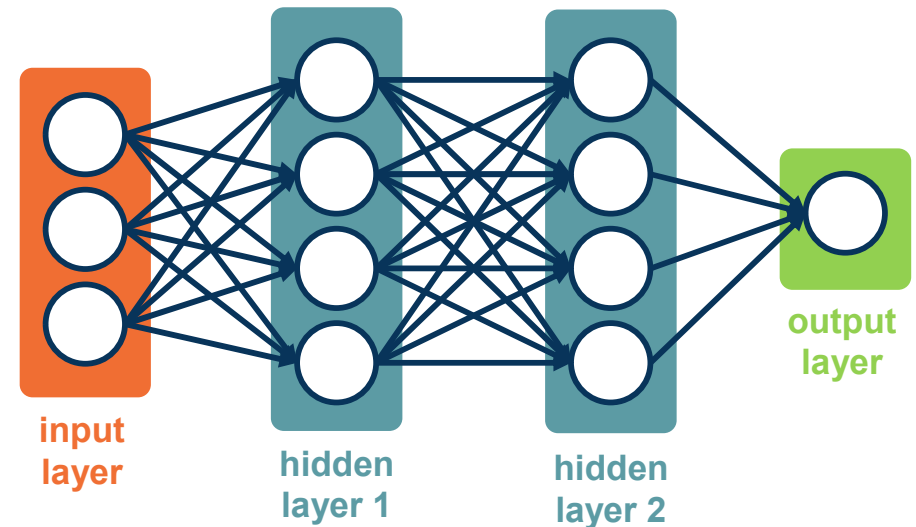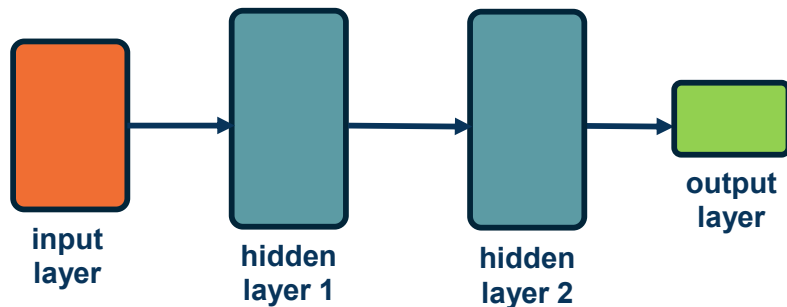
**input layer**

**hidden layer**

**output layer**

$$x \qquad W_1 \qquad W_2$$

$$=$$

$$f(x, W_1, W_2) = \sigma(W_2 \sigma(W_1 x))$$

*Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n*

**The Linear Algebra View**

Georgia Tech

**Large (deep) networks** can be built by adding more and more layers

Three-layered neural networks can represent **any function**

- The number of nodes could grow unreasonably (exponential or worse) with respect to the complexity of the function

We will show them **without edges**:



$$f(x, W_1, W_2, W_3) = \sigma(W_2 \sigma(W_1 x))$$

*Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n*
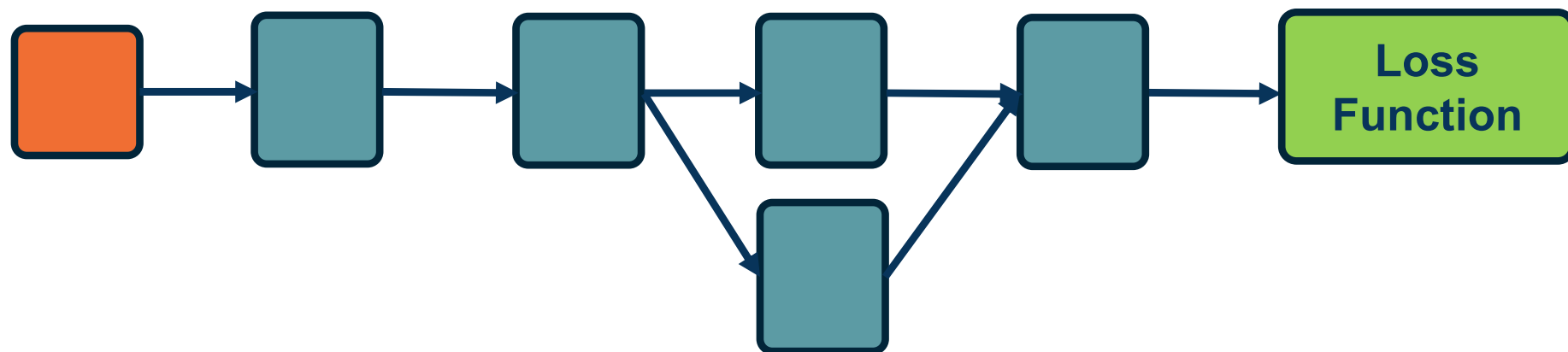
**Adding More Layers!**

Georgia Tech

Functions can be made **arbitrarily complex** (subject to memory and computational limits), e.g.:

$$f(x, W) = \sigma(W_5\sigma(W_4\sigma(W_3\sigma(W_2\sigma(W_1 x))))$$

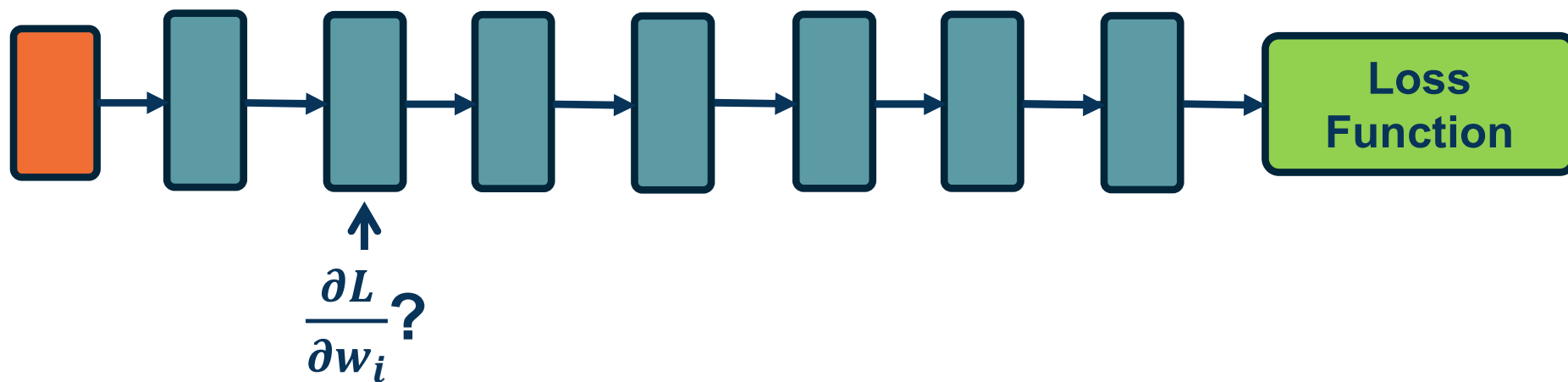We can use **any type of differentiable function (layer)** we want!

◆ At the end, **add the loss function**

Composition can have **some structure**

Georgia Tech

- We are learning **complex models** with significant amount of parameters (millions or billions)

- How do we compute the gradients of the **loss** (at the end) with respect to **internal** parameters?

- Intuitively, want to understand how **small changes** in weight deep inside **are propagated** to affect the **loss function** at the end
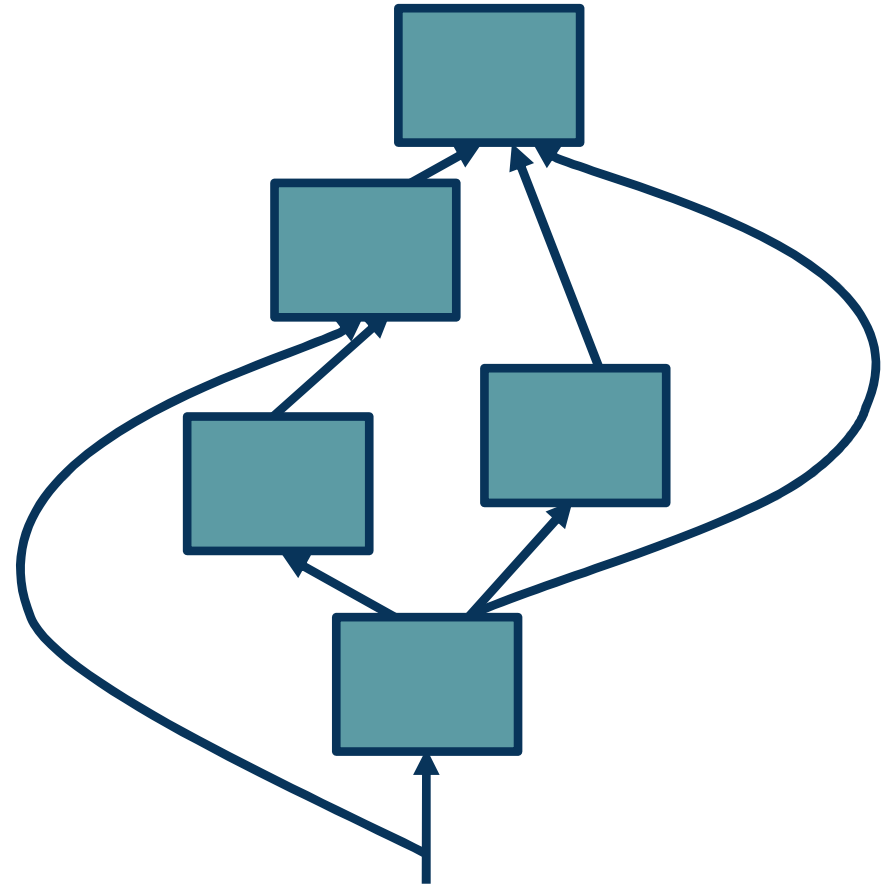


$$\frac{\partial L}{\partial w_i}?$$

To develop a general algorithm for this, we will view the function as a **computation graph**

Graph can be any **directed acyclic graph (DAG)**

◆ Modules must be differentiable to support gradient computations for gradient descent

A **training algorithm** will then process this graph, **one module at a time**



*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

**A General Framework**

Georgia Tech

**Step 1:** Compute Loss on Mini-Batch: **Forward Pass**

Layer 2

Layer 3

*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

Georgia Tech

**Step 1:** Compute Loss on Mini-Batch: **Forward Pass**

Layer 1 → Layer 3

*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

**Neural Network Training**

Georgia Tech

Note that we must store the **intermediate outputs of all layers**!

◆ This is because we will need them to **compute the gradients** (the gradient equations will have terms with the output values in them)

*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

**Neural Network Training**

Georgia Tech

**Step 1:** Compute Loss on Mini-Batch: **Forward Pass**

**Step 2:** Compute Gradients wrt parameters: **Backward Pass**



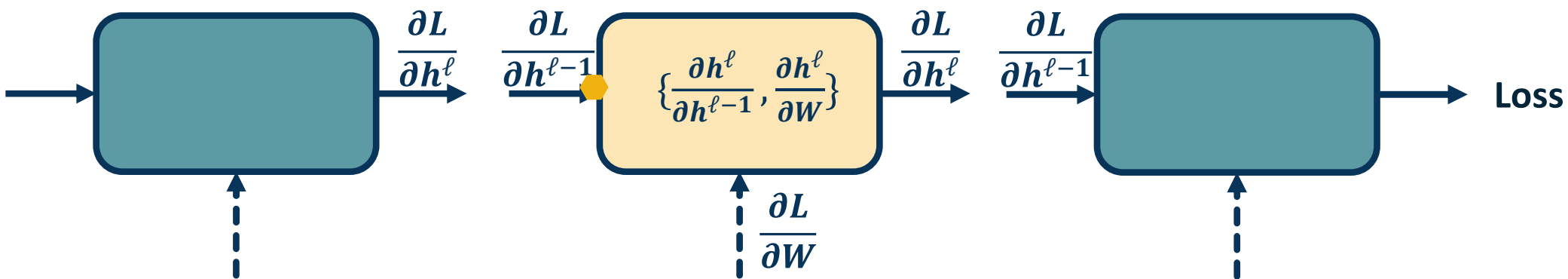*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

Georgia Tech

**Step 1:** Compute Loss on Mini-Batch: **Forward Pass**

**Step 2:** Compute Gradients wrt parameters: **Backward Pass**



*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

Georgia Tech

**Step 1:** Compute Loss on Mini-Batch: **Forward Pass**

**Step 2:** Compute Gradients wrt parameters: **Backward Pass**

Layer 2

Layer 3

*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

**Neural Network Training**

- We want to compute: $\left\{\dfrac{\partial L}{\partial h^{\ell-1}}, \dfrac{\partial L}{\partial W}\right\}$

$$\frac{\partial L}{\partial h^{\ell}} \quad \frac{\partial L}{\partial h^{\ell-1}} \quad \left\{\frac{\partial h^{\ell}}{\partial h^{\ell-1}}, \frac{\partial h^{\ell}}{\partial W}\right\} \quad \frac{\partial L}{\partial h^{\ell}} \quad \frac{\partial L}{\partial h^{\ell-1}} \quad \text{Loss}$$

$$\frac{\partial L}{\partial W}$$

- We will use the *chain rule* to do this:

**Chain Rule:** $\dfrac{\partial z}{\partial x} = \dfrac{\partial z}{\partial y} \cdot \dfrac{\partial y}{\partial x}$

**Computing the Gradients of Loss**

Georgia Tech

**Step 1:** Compute Loss on Mini-Batch: **Forward Pass**

**Step 2:** Compute Gradients wrt parameters: **Backward Pass**

**Step 3:** Use **gradient** to update **all parameters** at the end

Layer 2        Layer 3

$$w_i = w_i - \alpha \frac{\partial L}{\partial w_i}$$

**Backpropagation is the application of gradient descent to a computation graph via the chain rule!**

*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

Georgia Tech

- We can compute **local gradients**: $\{\frac{\partial h^\ell}{\partial h^{\ell-1}}, \frac{\partial h^\ell}{\partial W}\}$

- This is just the **derivative of our function** with respect to its parameters and inputs!

**Example:**    If $h^\ell = Wh^{\ell-1}$

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

and $\frac{\partial h^\ell}{\partial w_i} = h^{\ell-1,T}$    (a sparse matrix with $h^{\ell-1,T}$ in the *i*-th row

Georgia Tech

◆ We will use the **chain rule** to compute: $\{\frac{\partial L}{\partial h^{\ell-1}}, \frac{\partial L}{\partial W}\}$

◆ **Gradient of loss w.r.t. inputs:** $\frac{\partial L}{\partial h^{\ell-1}} = \frac{\partial L}{\partial h^\ell} \frac{\partial h^\ell}{\partial h^{\ell-1}}$

Given by upstream module **(upstream gradient)**

◆ **Gradient of loss w.r.t. weights:** $\frac{\partial L}{\partial W} = \frac{\partial L}{\partial h^\ell} \frac{\partial h^\ell}{\partial W}$

$\frac{\partial L}{\partial h^{\ell-1}}$        $\frac{\partial L}{\partial h^\ell}$

$\frac{\partial L}{\partial W}$

*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

# Backpropagation: a simple example

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

# Backpropagation: a simple example

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

Georgia Tech

# Backpropagation: a simple example

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

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

Georgia Tech

# Backpropagation: a simple example

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

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



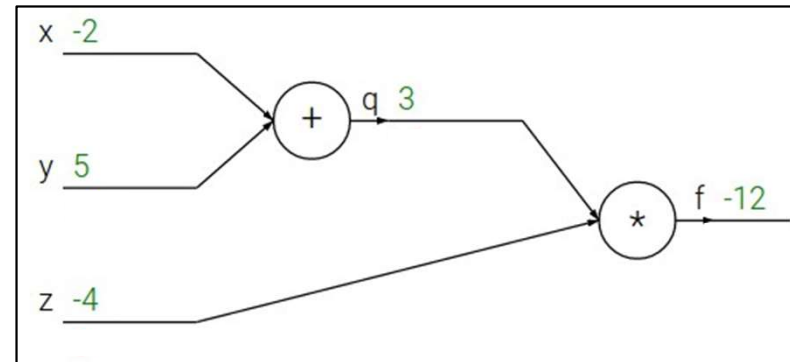Want:  $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$

Georgia
Tech

# Backpropagation: a simple example

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

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

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$



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

*Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n*

Georgia
Tech

# Backpropagation: a simple example

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

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

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

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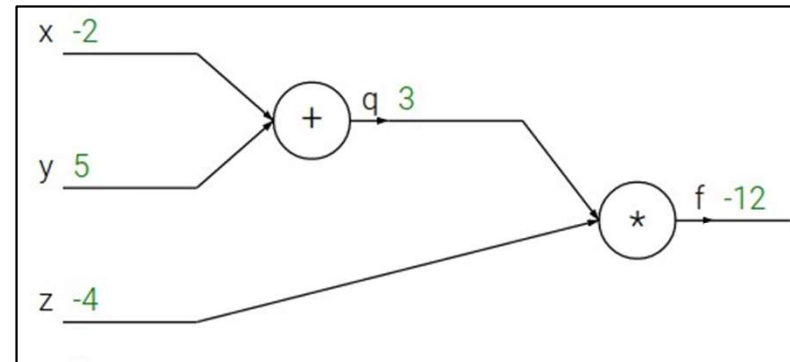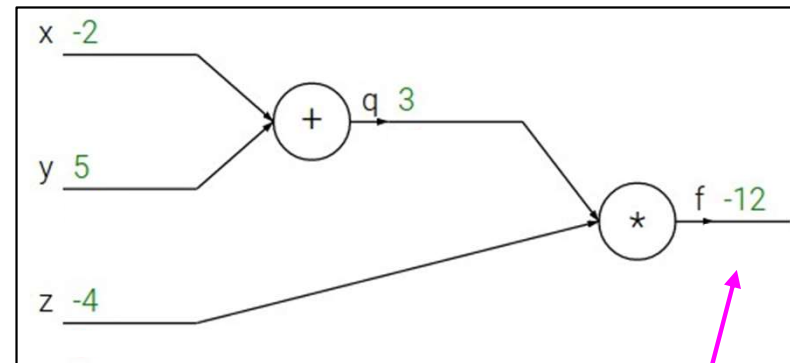


*Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n*

Georgia Tech

# Backpropagation: a simple example

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

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

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

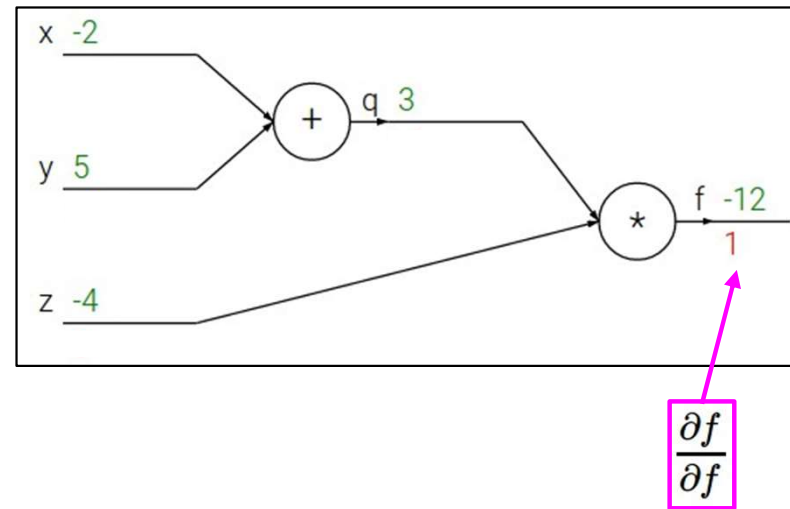Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



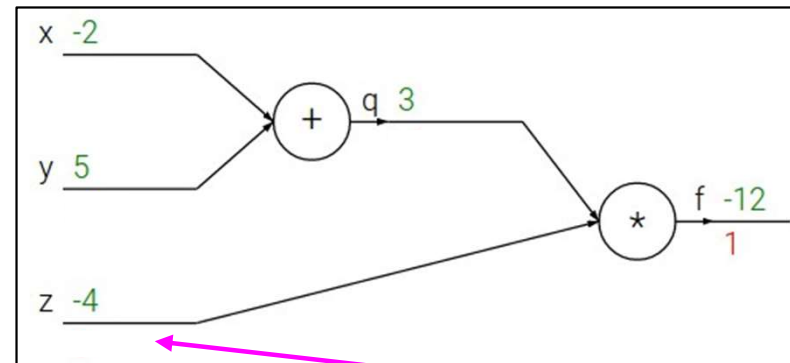$$\frac{\partial f}{\partial f}$$

# Backpropagation: a simple example

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

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

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

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



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

*Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n*
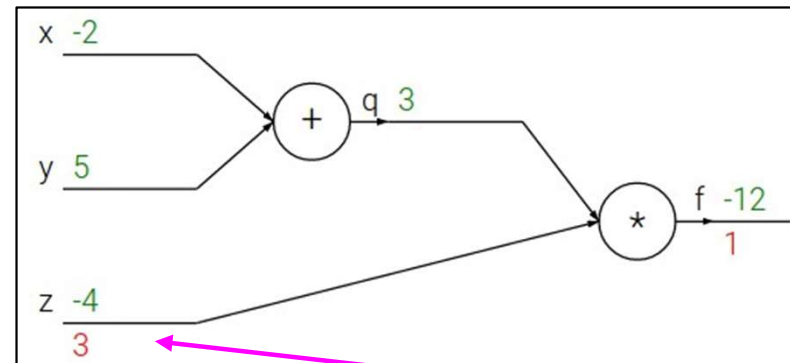
Georgia Tech

# Backpropagation: a simple example

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

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

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

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



$$\frac{\partial f}{\partial z}$$

Georgia Tech

# Backpropagation: a simple example

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

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

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

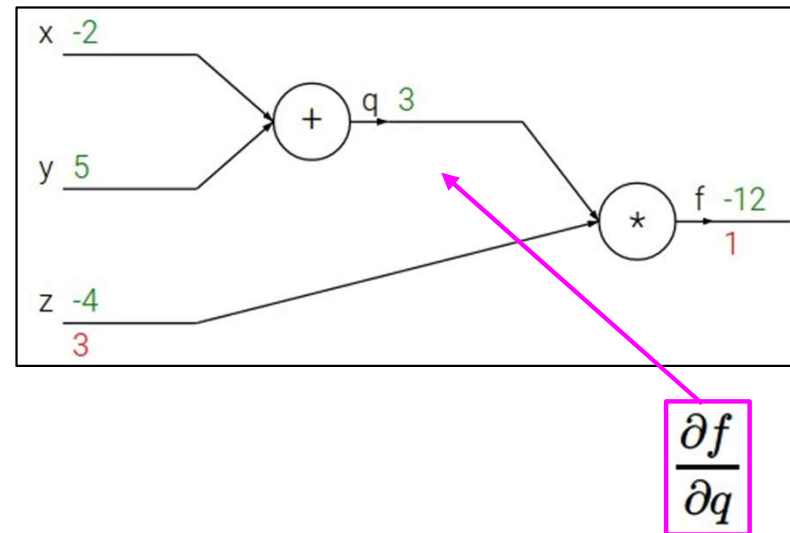Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial z}$$

Georgia Tech

# Backpropagation: a simple example

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

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

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

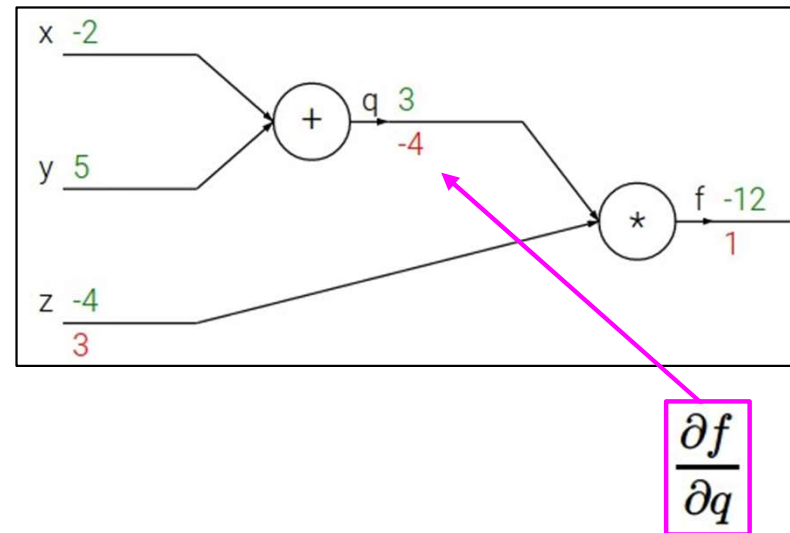Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



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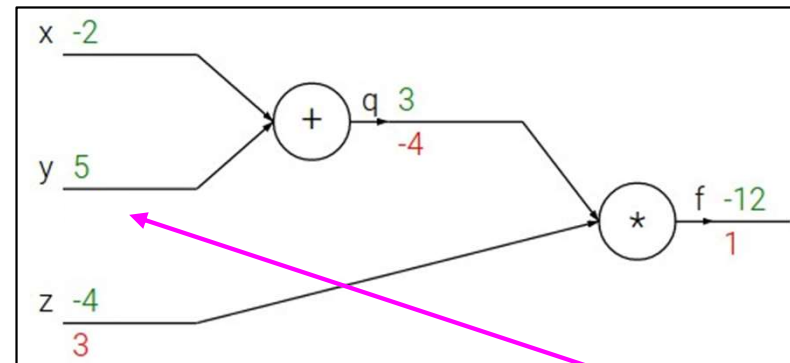
Georgia Tech

# Backpropagation: a simple example

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

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

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

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



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

*Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n*
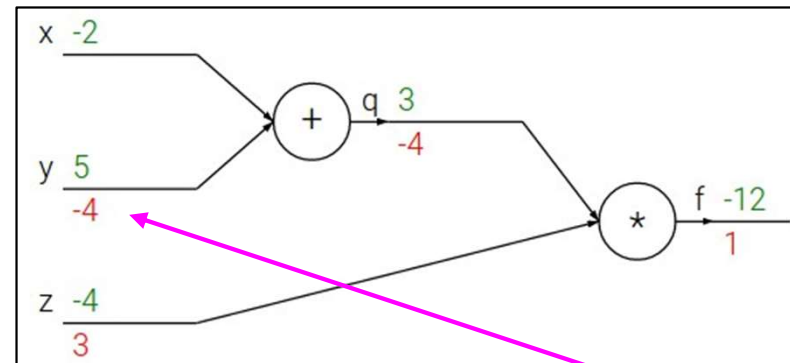
Georgia Tech

# Backpropagation: a simple example

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

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

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

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



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

Chain rule:

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

Upstream gradient    Local gradient

*Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n*
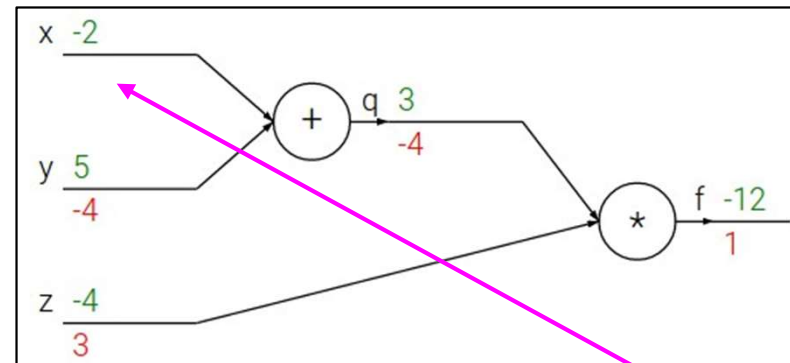
Georgia Tech

# Backpropagation: a simple example

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

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

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

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



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

Chain rule:

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

Upstream gradient     Local gradient

*Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n*
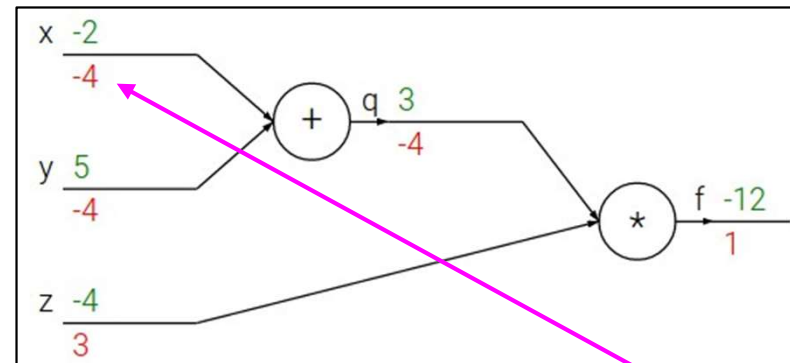
Georgia Tech

# Backpropagation: a simple example

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

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

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

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



$\dfrac{\partial f}{\partial x}$

Chain rule:

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

Upstream gradient  Local gradient

Georgia Tech

# Backpropagation: a simple example

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

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

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

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



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

Chain rule:

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

Upstream gradient    Local gradient

# Backpropagation: a simple example

Georgia Tech

# Backpropagation: a simple example

Georgia Tech

# Patterns in backward flow

# Patterns in backward flow

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

Georgia Tech

# Patterns in backward flow

**add** gate: gradient distributor

Georgia Tech

# Patterns in backward flow

**add** gate: gradient distributor

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

# Patterns in backward flow

**add** gate: gradient distributor

**max** gate: gradient router

Georgia
Tech

# Patterns in backward flow

**add** gate: gradient distributor

**max** gate: gradient router
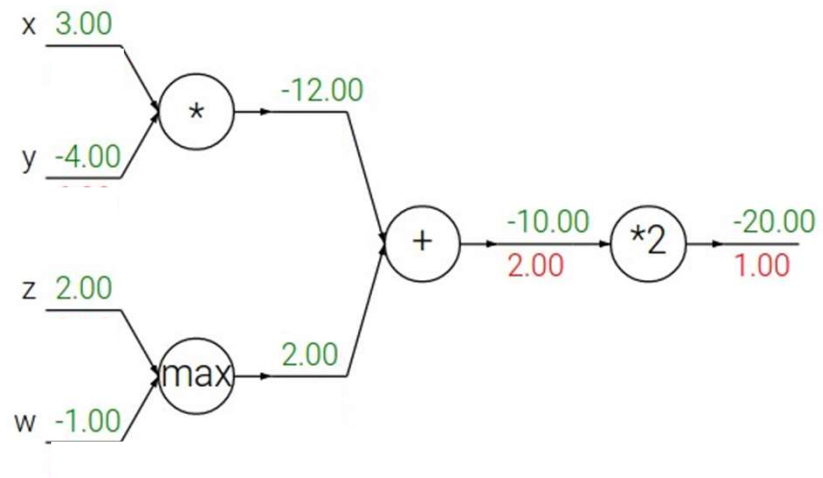
Q: What is a **mul** gate?



*Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n*

Georgia
Tech

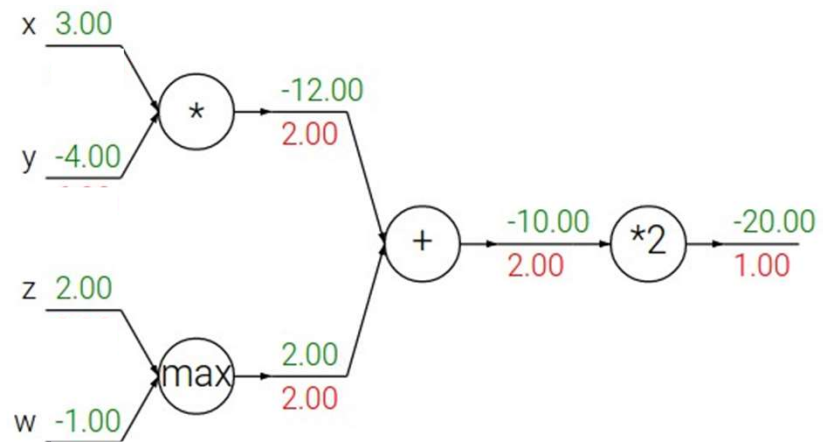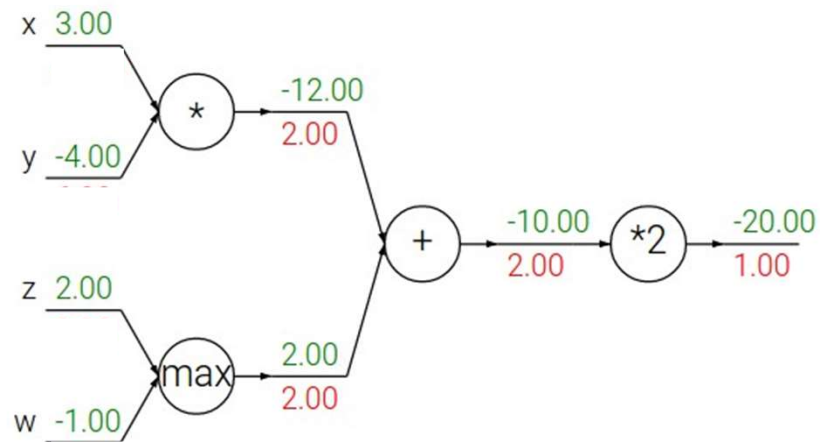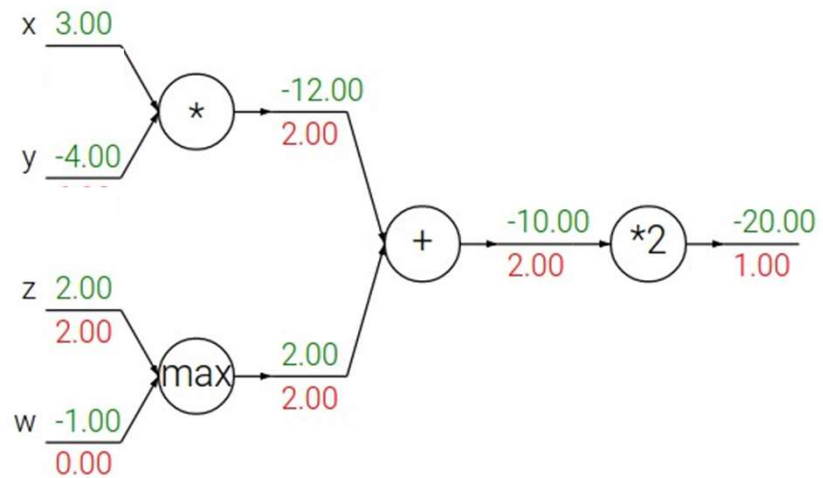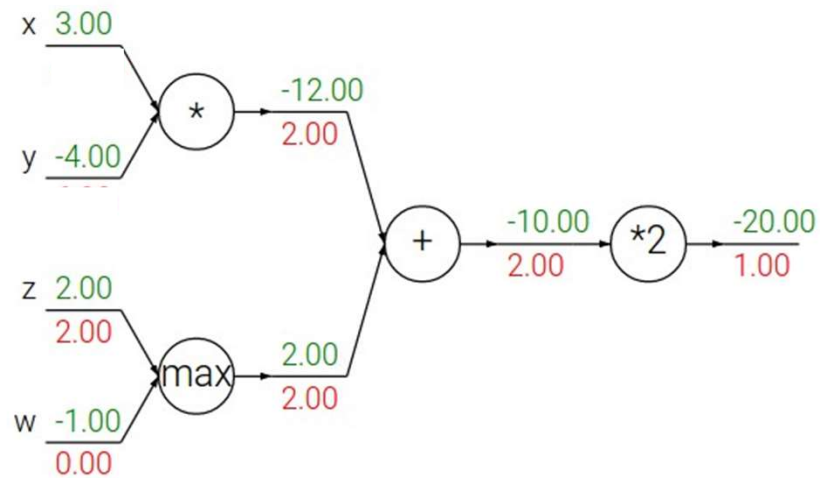# Patterns in backward flow

**add** gate: gradient distributor

**max** gate: gradient router

**mul** gate: gradient switcher

Georgia
Tech

# Gradients add at branches

Georgia Tech

# Duality in Fprop and Bprop

- Neural networks involves composing simple functions into a **computation graph**

- Optimization (updating weights) of this graph is through backpropagation
  - Recursive algorithm: Gradient descent (partial derivatives) plus chain rule

- Remaining questions:
  - How does this work with vectors, matrices, tensors?
    - Across a composed function?
  - How can we implement this algorithmically to make these calculations automatic? **Automatic Differentiation**

Georgia
Tech

# Linear Algebra View:
# Vector and Matrix Sizes

$$
\begin{bmatrix}
w_{11} & w_{12} & \cdots & w_{1m} & b1 \\
w_{21} & w_{22} & \cdots & w_{2m} & b2 \\
w_{31} & w_{32} & \cdots & w_{3m} & b3
\end{bmatrix}
\begin{bmatrix}
x_1 \\
x_2 \\
\vdots \\
x_m \\
1
\end{bmatrix}
$$

$$W \qquad\qquad x$$

**Sizes:** $[c \times (d + 1)]$  $[(d + 1) \times 1]$

Where $c$ is number of classes

$d$ is dimensionality of input

**Conventions:**

- Size of derivatives for scalars, vectors, and matrices:
Assume we have scalar $s \in \mathbb{R}^1$, vector $v \in \mathbb{R}^m$, i.e. $v = [v_1, v_2, \dots, v_m]^T$ and matrix $M \in \mathbb{R}^{k \times \ell}$

|  | $S$ $[\ ]$ | $V$ $\begin{bmatrix} \\ \end{bmatrix}$ | $M$ $\begin{bmatrix} & \\ & \end{bmatrix}$ |
|---|---|---|---|
| $S$ | $\dfrac{\partial s_1}{\partial s_2}$ $[\ ]$ | $\dfrac{\partial s}{\partial v}$ $[\qquad]$ | $\dfrac{\partial s}{\partial M}$ $\begin{bmatrix} & \\ & \end{bmatrix}$ |
| $V$ | $\dfrac{\partial v}{\partial s}$ $\begin{bmatrix} \\ \end{bmatrix}$ | $\dfrac{\partial v_1}{\partial v_2}$ $\begin{bmatrix} \\ \end{bmatrix}$ | |
| $M$ | $\dfrac{\partial M}{\partial s}$ $\begin{bmatrix} & \\ & \end{bmatrix}$ | | |

**Tensors**

Georgia Tech

**Conventions:**

- Size of derivatives for scalars, vectors, and matrices:
  Assume we have scalar $s \in \mathbb{R}^1$, vector $v \in \mathbb{R}^m$, i.e. $v = [v_1, v_2, \ldots, v_m]^T$
  and matrix $M \in \mathbb{R}^{k \times \ell}$

- What is the size of $\frac{\partial v}{\partial s}$ ? $\mathbb{R}^{m \times 1}$ (column vector of size $m$)

- What is the size of $\frac{\partial s}{\partial v}$ ? $\mathbb{R}^{1 \times m}$ (row vector of size $m$)

$$\begin{bmatrix} \dfrac{\partial v_1}{\partial s} \\ \dfrac{\partial v_2}{\partial s} \\ \vdots \\ \dfrac{\partial v_m}{\partial s} \end{bmatrix}$$

$$\begin{bmatrix} \dfrac{\partial s}{\partial v_1} & \dfrac{\partial s}{\partial v_1} & \cdots & \dfrac{\partial s}{\partial v_m} \end{bmatrix}$$

Georgia Tech

# Conventions:

⬡ What is the size of $\frac{\partial v^1}{\partial v^2}$ ? A matrix:

$$
\text{Row } i \quad
\begin{bmatrix}
\dfrac{\partial v_1^1}{\partial v_1^2} & \dots & \dots & \dots & \dots & \\
\dots & & \dots & \dots & \dots & \dots \\
\dfrac{\partial v_i^1}{\partial v_1^2} & \dots & \dfrac{\partial v_i^1}{\partial v_j^2} & \dots & \dfrac{\partial v_i^1}{\partial v_{m_2}^2} & \\
\dots & \dots & \dots & \dots & \dots & \\
\dots & \dots & \dots & \dots & \dots &
\end{bmatrix}
\begin{array}{c} \textbf{Col } j \end{array}
$$

$$m_1 \times m_2$$

⬡ This matrix of partial derivatives is called a **Jacobian**

(Note this is slightly different convention than on Wikipedia). Also, computationally other conventions are used.

**Dimensionality of Derivatives**

# Conventions:

◆ What is the size of $\frac{\partial s}{\partial M}$ ? A matrix:

$$\begin{bmatrix} \dfrac{\partial s}{\partial m_{[1,1]}} & \cdots & \cdots & \cdots & \cdots & \\ \cdots & & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \dfrac{\partial s}{\partial m_{[i,j]}} & \cdots & \cdots & \\ \cdots & \cdots & \cdots & \cdots & \cdots & \\ \cdots & \cdots & \cdots & \cdots & \cdots & \end{bmatrix}$$

(Note this is slightly different convention than on Wikipedia). Also, computationally other conventions are used.

**Example 1:**

$$y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} x \\ x^2 \end{bmatrix} \qquad \frac{\partial y}{\partial x} = \begin{bmatrix} 1 \\ 2x \end{bmatrix}$$

**Example 2:**

$$y = w^T x = \sum_k w_k x_k$$

$$\frac{\partial y}{\partial x} = \left[ \frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_m} \right]$$

$$= [w_1, \dots, w_m] \qquad \text{because} \qquad \frac{\partial \left( \sum_k w_k x_k \right)}{\partial x_i} = w_i$$

$$= w^T$$

Georgia
Tech

## Example 3:

$$y = Wx \qquad \frac{\partial y}{\partial x} = W$$

**Col** $j$

$$\text{Row } i \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \cdots & \cdots & \cdots \\ \cdots & & \cdots & \cdots & \cdots \\ \cdots & \cdots & \frac{\partial y_i}{\partial x_j} & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots \end{bmatrix} = \begin{bmatrix} \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & w_{ij} & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots \end{bmatrix} \qquad y_i = \sum_j w_{ij} x_j$$

## Example 4:

$$\frac{\partial(wAw)}{\partial w} = 2w^T A \text{ (assuming A is symmetric)}$$

- What is the size of $\frac{\partial L}{\partial W}$ ?

  - Remember that loss is a **scalar** and $W$ is a matrix:

$$\begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1m} & b1 \\ w_{21} & w_{22} & \cdots & w_{2m} & b2 \\ w_{31} & w_{32} & \cdots & w_{3m} & b3 \end{bmatrix}$$

Jacobian is also a matrix:

$$W$$

$$\begin{bmatrix} \frac{\partial L}{\partial w_{11}} & \frac{\partial L}{\partial w_{12}} & \cdots & \frac{\partial L}{\partial w_{1m}} & \frac{\partial L}{\partial b_1} \\ \frac{\partial L}{\partial w_{21}} & \cdots & \cdots & \frac{\partial L}{\partial w_{2m}} & \frac{\partial L}{\partial b_2} \\ \cdots & \cdots & \cdots & \frac{\partial L}{\partial w_{3m}} & \frac{\partial L}{\partial b_3} \end{bmatrix}$$

Georgia
Tech

Batches of data are **matrices** or **tensors** (multi-dimensional matrices)

**Examples:**

- Each instance is a vector of size $m$, our batch is of size $[B \times m]$

- Each instance is a matrix (e.g. grayscale image) of size $W \times H$, our batch is $[B \times W \times H]$

- Each instance is a multi-channel matrix (e.g. color image with R,B,G channels) of size $C \times W \times H$, our batch is $[B \times C \times W \times H]$
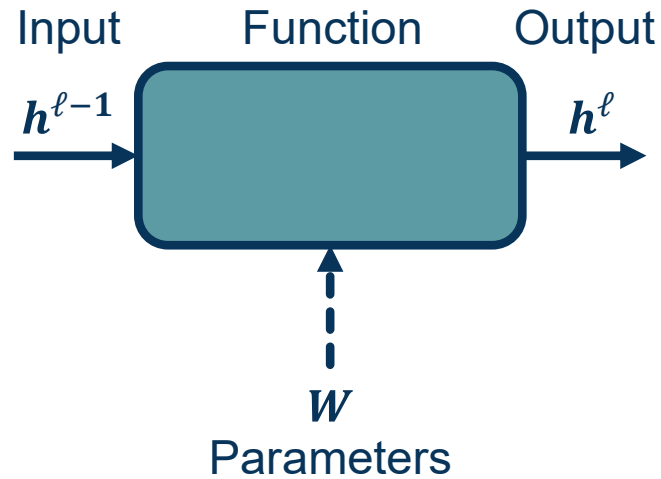
**Jacobians become tensors which is complicated**

- Instead, flatten input to a vector and get a vector of derivatives!

- This can also be done for partial derivatives between two vectors, two matrices, or two tensors

$$\begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nn} \end{bmatrix}$$

**Flatten** ⬇

$$\begin{bmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{21} \\ x_{22} \\ \vdots \\ x_{n1} \\ \vdots \\ x_{nn} \end{bmatrix}$$

**Jacobians of Batches**

Georgia Tech

Input      Function     Output

$h^{\ell-1}$ → [Function] → $h^{\ell}$

$W$

Parameters

**Define:**

$h_i^{\ell} = w_i^T h^{\ell-1}$

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

$\leftarrow w_i^T \rightarrow$

$|h^{\ell}| \times 1$    $|h^{\ell}| \times |h^{\ell-1}|$    $|h^{\ell-1}| \times 1$

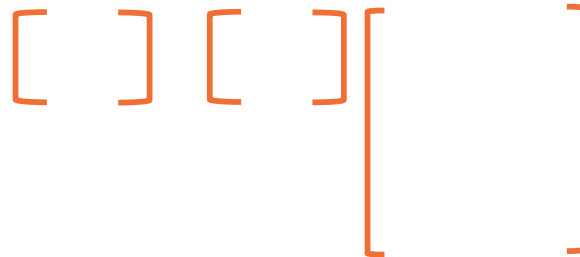**Fully Connected (FC) Layer: Forward Function**

Georgia Tech

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

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

**Define:**

$$h_i^\ell = w_i^T h^{\ell-1}$$

$$\frac{\partial L}{\partial h^{\ell-1}}$$

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

$$\frac{\partial L}{\partial W}$$

$$\boxed{\frac{\partial L}{\partial h^{\ell-1}}} = \frac{\partial L}{\partial h^\ell} \frac{\partial h^\ell}{\partial h^{\ell-1}}$$

$$\begin{bmatrix} \end{bmatrix} \begin{bmatrix} \end{bmatrix} \begin{bmatrix} \end{bmatrix}$$

$$1 \times |h^{\ell-1}| \qquad 1 \times |h^\ell| \qquad |h^\ell| \times |h^{\ell-1}|$$

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

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

**Define:**

$$h_i^\ell = w_i^T h^{\ell-1}$$

$$\frac{\partial h_i^\ell}{\partial w_i^T} = h^{(\ell-1),T}$$

$$\frac{\partial L}{\partial h^{\ell-1}} \qquad \frac{\partial L}{\partial h^\ell}$$

$$\frac{\partial L}{\partial W}$$

$$\frac{\partial L}{\partial w_i^T} = \frac{\partial L}{\partial h^\ell} \frac{\partial h^\ell}{\partial w_i^T}$$

$$\begin{bmatrix} \; \end{bmatrix} \begin{bmatrix} \; \end{bmatrix} \begin{bmatrix} \leftarrow & 0 & \rightarrow \\ \leftarrow & \frac{\partial h_i^\ell}{\partial w_i^T} & \rightarrow \\ \leftarrow & 0 & \rightarrow \end{bmatrix}$$

$$1 \times |h^{\ell-1}| \quad 1 \times |h^\ell| \quad |h^\ell| \times |h^{\ell-1}|$$
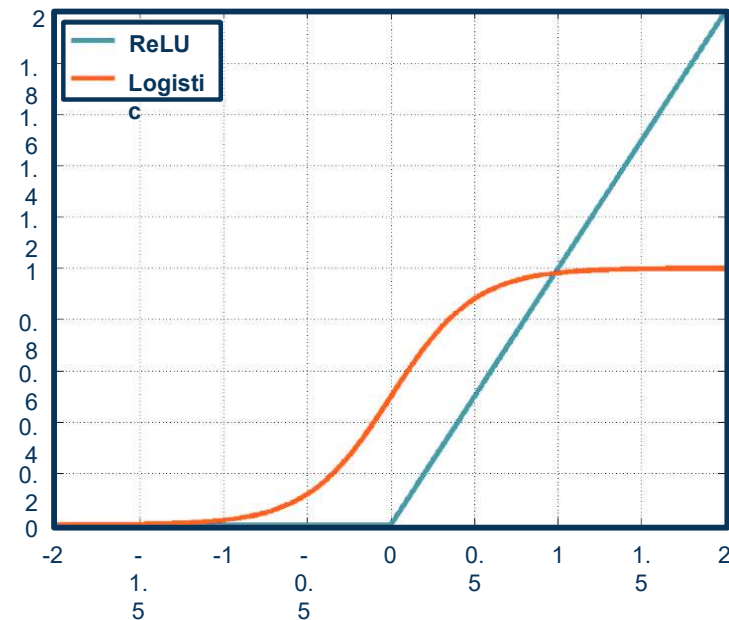
**Fully Connected (FC) Layer**

Georgia Tech

We can employ **any differentiable (or piecewise differentiable) function**
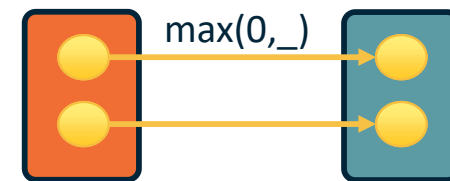
A common choice is the **Rectified Linear Unit**

- Provides non-linearity but better gradient flow than sigmoid
- Performed **element-wise**

**How many** parameters for this layer?

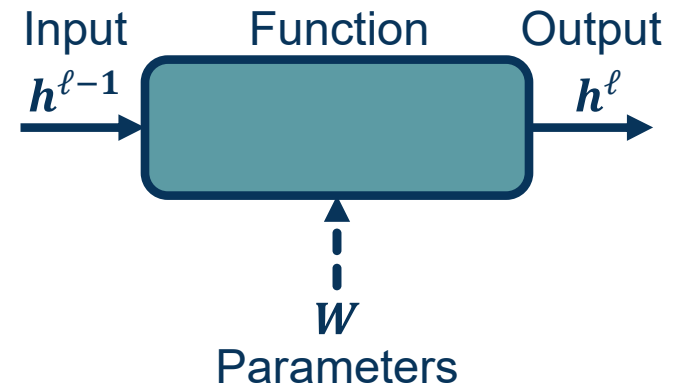

$$h^\ell = \max(0, h^{\ell-1})$$

Full Jacobian of ReLU layer is **large** (output dim x input dim)

- But again it is **sparse**

- Only **diagonal values non-zero** because it is element-wise

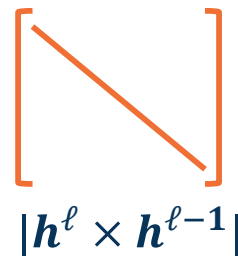- An output value affected only by **corresponding input value**

Max function **funnels gradients through selected max**

- Gradient will be **zero** if input **<= 0**

Input $\quad$ Function $\quad$ Output

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

$W$

Parameters

**Forward:** $h^{\ell} = \max(0, h^{\ell-1})$

**Backward:** $\frac{\partial L}{\partial h^{\ell-1}} = \frac{\partial L}{\partial h^{\ell}} \frac{\partial h^{\ell}}{\partial h^{\ell-1}}$
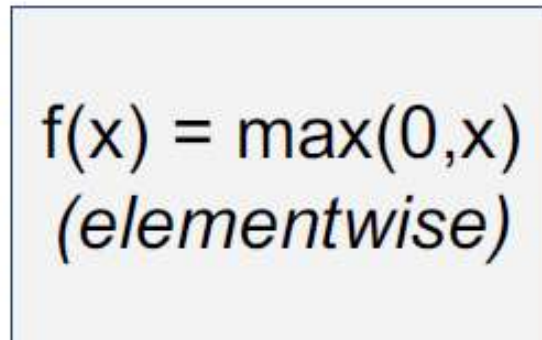
For diagonal

$|h^{\ell} \times h^{\ell-1}|$

$$\frac{\partial h^{\ell}}{\partial h^{\ell-1}} = \begin{cases} 1 & if\ h^{\ell-1} > 0 \\ 0 & otherwise \end{cases}$$

**Jacobian of ReLU**

4D input x:

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

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

4D output z:

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

What does $\dfrac{\partial z}{\partial x}$ look like?

4D dL/dz:

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

Upstream gradient

Georgia Tech

For element-wise ops, jacobian is **sparse**: off-diagonal entries always zero!
Never **explicitly** form Jacobian -- instead use elementwise multiplication

- Neural networks involves composing simple functions into a **computation graph**

- Optimization (updating weights) of this graph is through backpropagation
  - Recursive algorithm: Gradient descent (partial derivatives) plus chain rule

- Remaining questions:
  - How does this work with vectors, matrices, tensors?
    - Across a composed function? **Next Time!**
  - How can we implement this algorithmically to make these calculations automatic? **Automatic Differentiation**

**Summary**

Georgia Tech