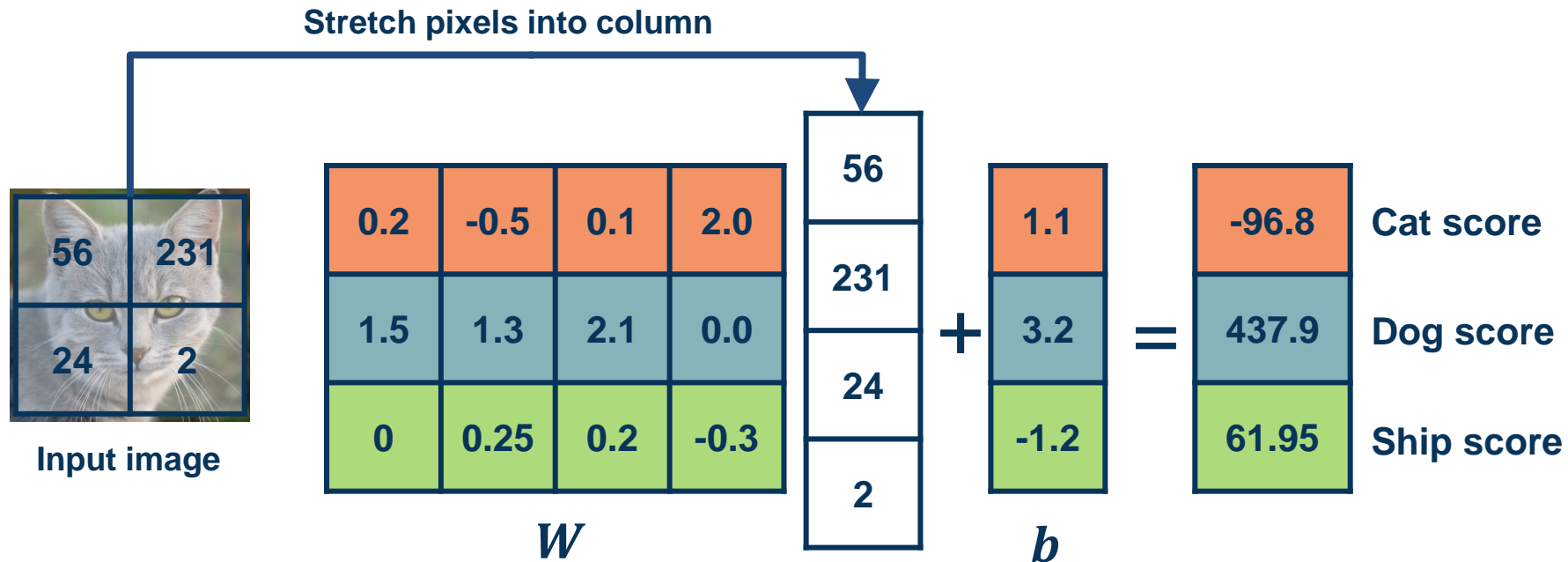Topics:

- Backpropagation
- Matrix/Linear Algebra view

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

- **Assignment 1 out!**
  - **Due Feb 2$^{nd}$ (with grace period 4$^{th}$)**
  - 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 (@57) and matrix calculus/computation graph (TBD)

- Piazza: Project teaming thread
  - Project proposal overview during my OH (Thursday 4pm ET)

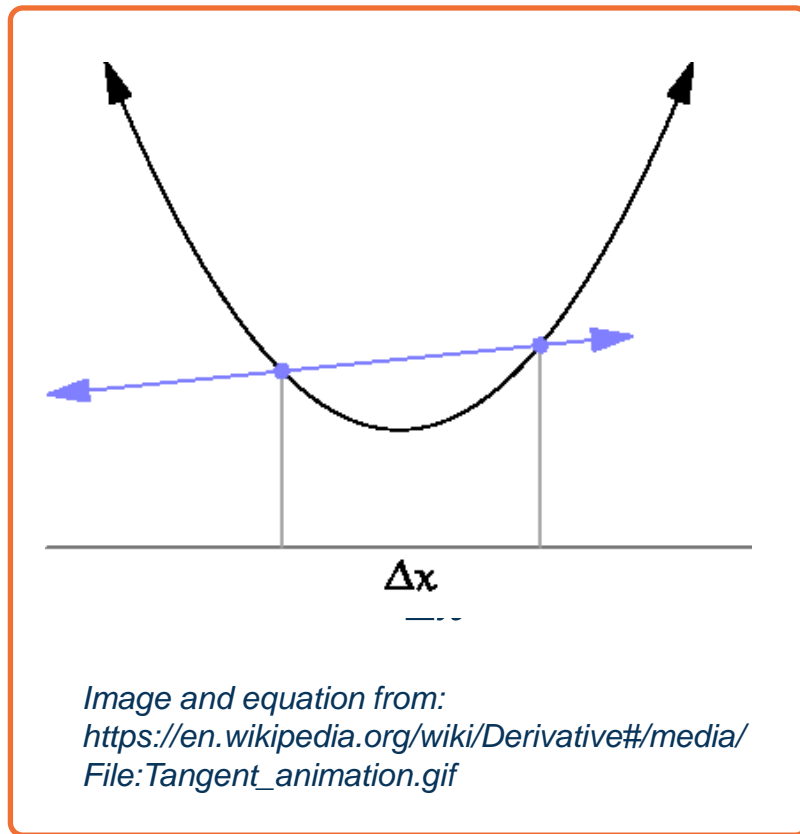# Example with an image with **4 pixels**, and **3 classes (cat/dog/ship)**

**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



$\Delta x$

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

**Derivatives**

Georgia Tech

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

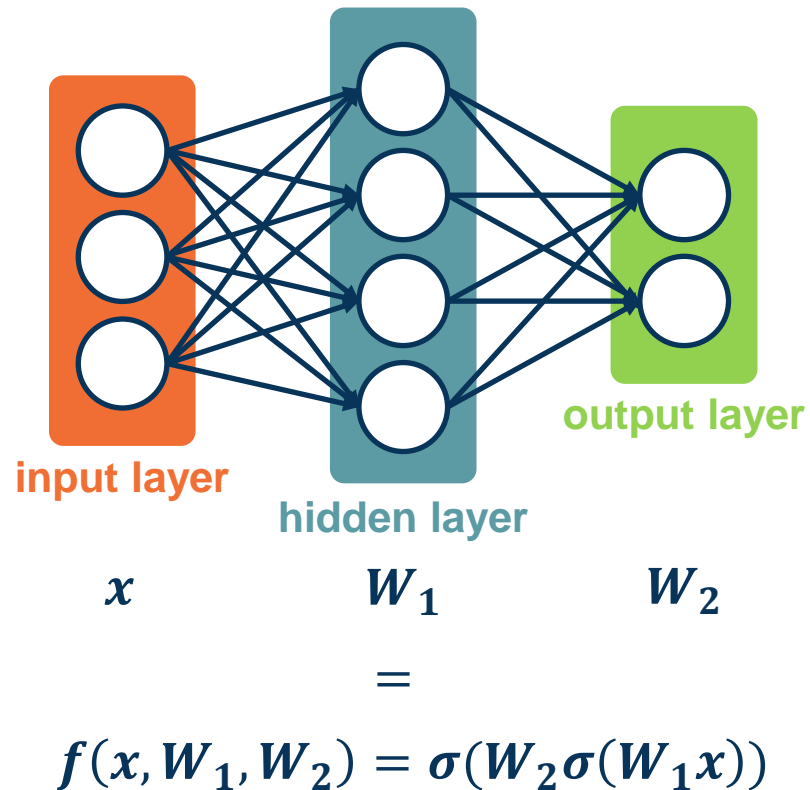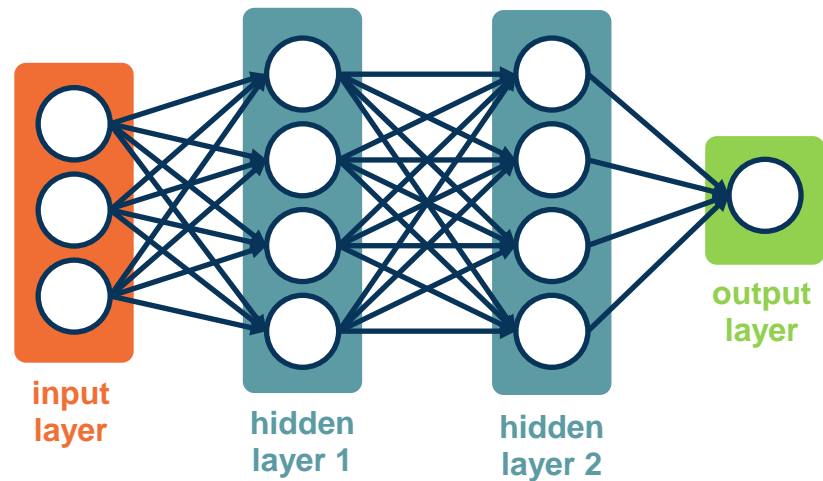◆ We will prefer the linear algebra view, but use some terminology from neural networks (& biology)
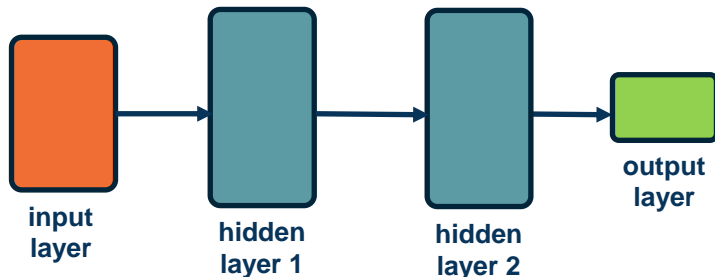


**input layer**

**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
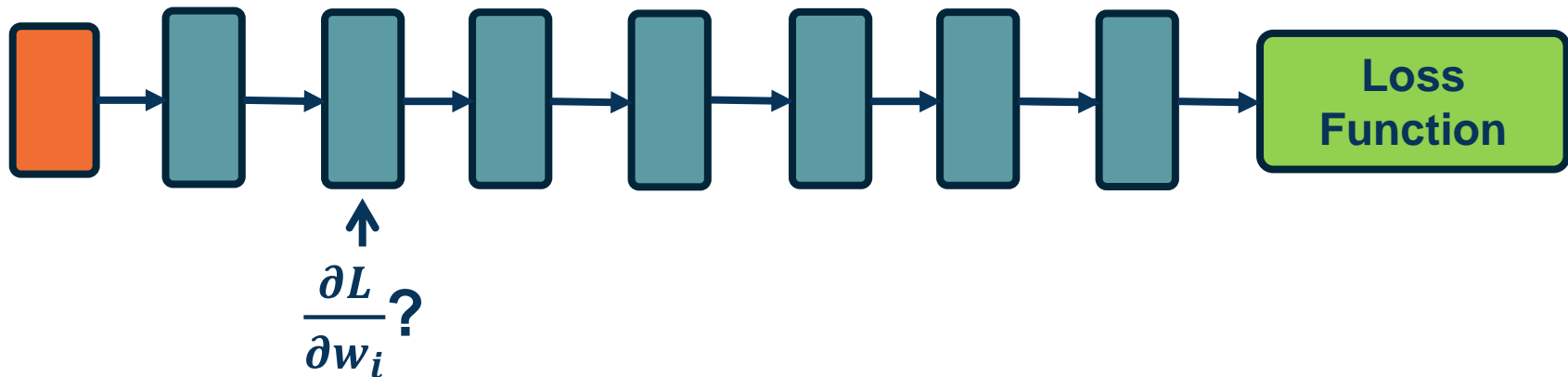
We will show them **without edges**:





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

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

**Adding More Layers!**

- 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
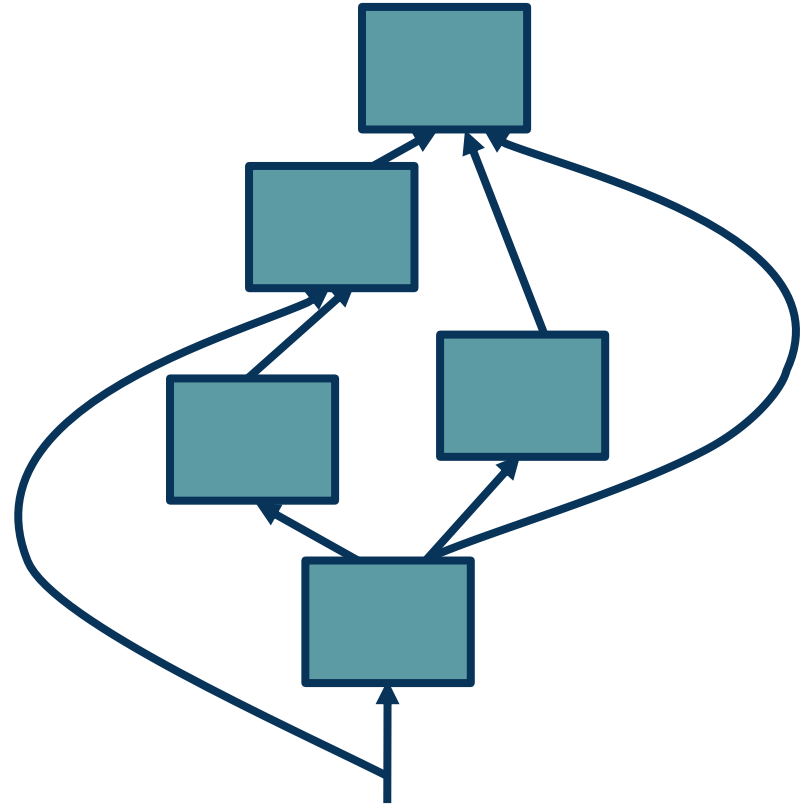


$$\frac{\partial L}{\partial w_i}?$$

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

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

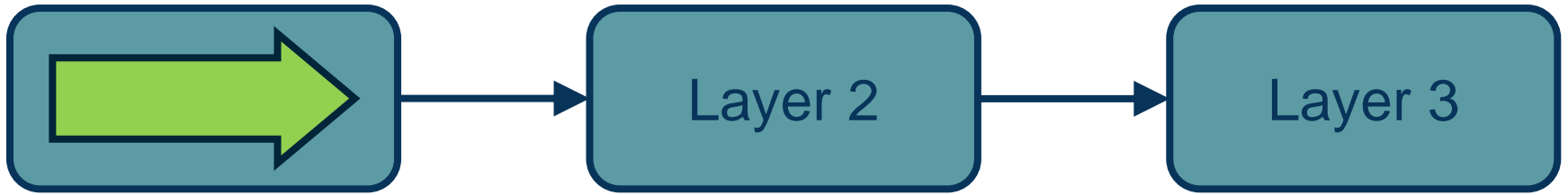◆ Modules must be differentiable to support gradient computations for gradient descent

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



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

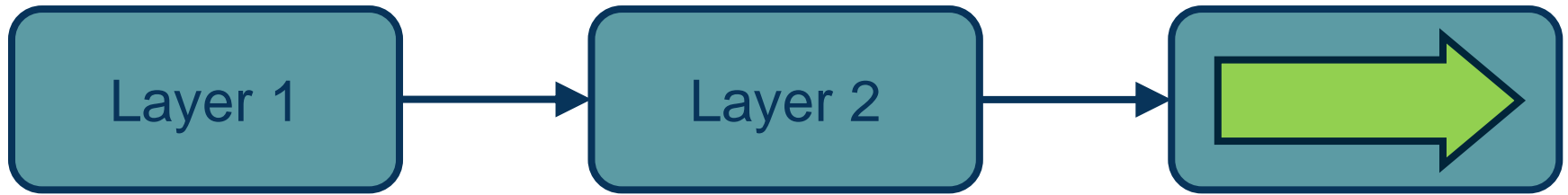Georgia Tech

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

Layer 2

Layer 3

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

**Neural Network Training**

Georgia Tech

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

Layer 1

Layer 3

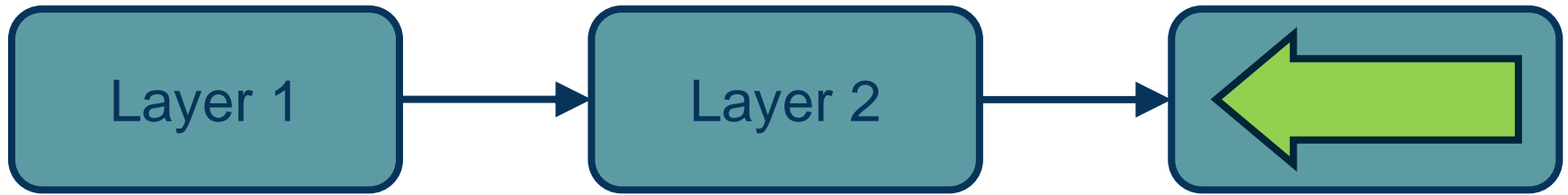*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

**Neural Network Training**

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*

**Neural Network Training**

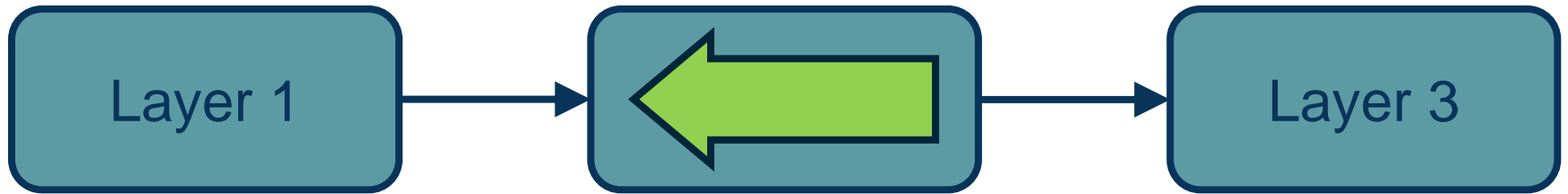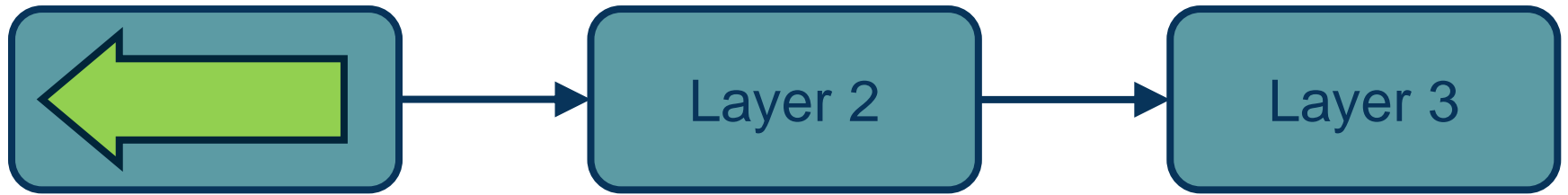**Step 1:** Compute Loss on Mini-Batch: **Forward Pass**

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



*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*
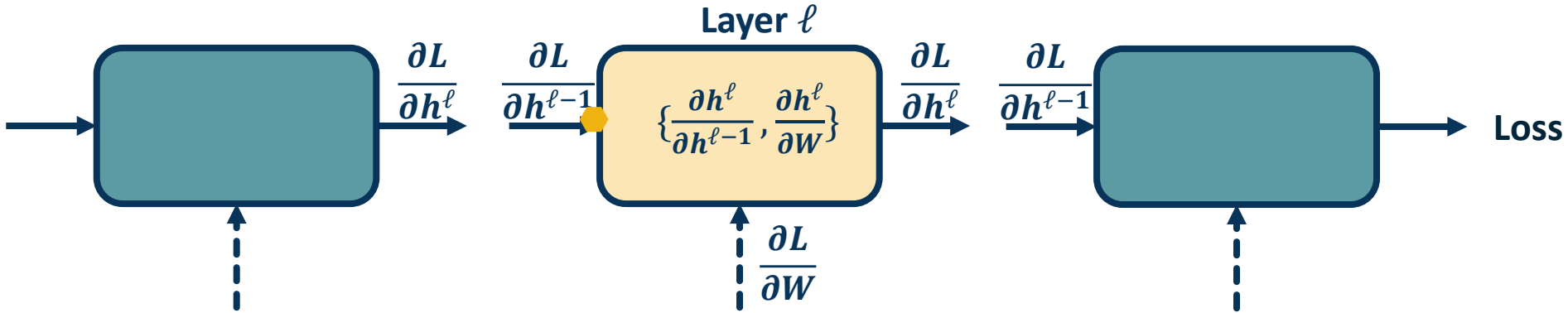
**Step 1:** Compute Loss on Mini-Batch: **Forward Pass**

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



Layer 2

Layer 3

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

**Neural Network Training**

Georgia Tech

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

**Layer** $\ell$

$\dfrac{\partial L}{\partial h^{\ell}}$  $\dfrac{\partial L}{\partial h^{\ell-1}}$  $\left\{ \dfrac{\partial h^{\ell}}{\partial h^{\ell-1}}, \dfrac{\partial h^{\ell}}{\partial W} \right\}$  $\dfrac{\partial L}{\partial h^{\ell}}$  $\dfrac{\partial L}{\partial h^{\ell-1}}$  **Loss**

$\dfrac{\partial L}{\partial W}$

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

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

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

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

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



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

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

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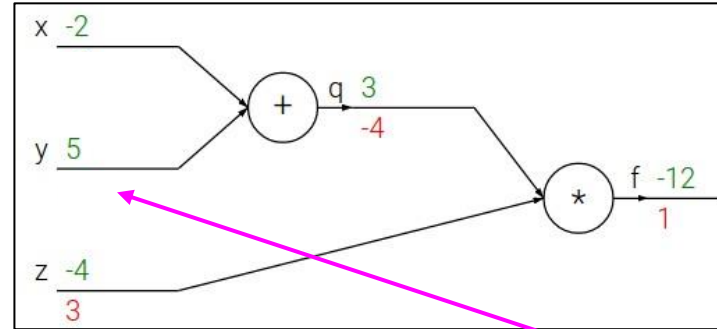
Georgia Tech

# 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

# Patterns in backward flow

**add** gate: gradient distributor

**max** gate: gradient router
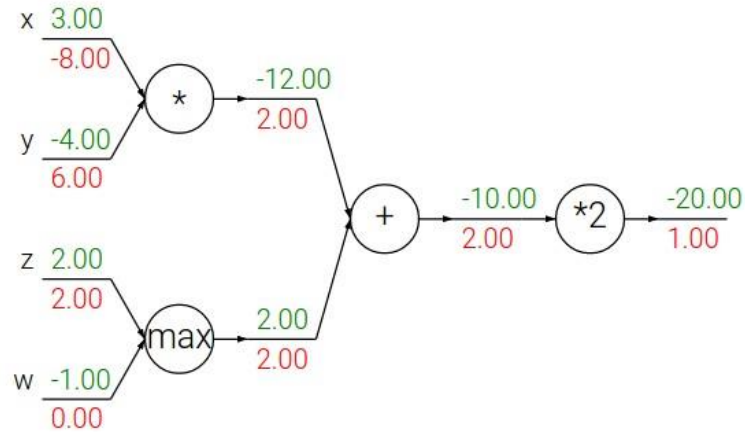
**mul** gate: gradient switcher



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

- 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

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

$$W \qquad\qquad\qquad x$$

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

Where $c$ is number of classes

$m$ 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. $v = [v_1, v_2, \ldots, v_m]^T$
and matrix $M \in \mathbb{R}^{m_1 \times m_2}$

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

**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 $M \in \mathbb{R}^{m_1 \times m_2}$

◆ 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:

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

Row $i$

$m_1 \times m_2$

- 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**

Georgia Tech

# 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.

Georgia Tech

**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:**

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

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

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

$$= w^T$$

Georgia Tech

# Example 3:

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

**Col $j$**

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

# Example 4:

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

Georgia Tech

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:

$$W$$

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

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 $[B \times C \times W \times H]$
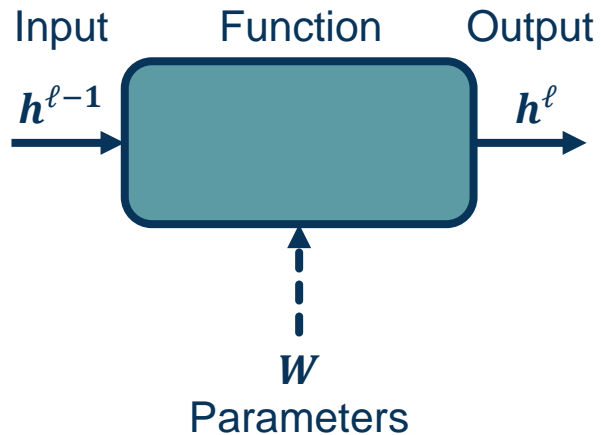
**Jacobians become tensors which is complicated**

- 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

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

$W$
Parameters

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

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

$$\overleftrightarrow{w_i^T}$$

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

Georgia Tech

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

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

**Define:**

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

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

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

Georgia Tech

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

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

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

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

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

Note doing this on full $W$ matrix would result in Jacobian tensor!

But it is *sparse* – each output only affected by corresponding weight row

Georgia Tech

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

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

**Define:**

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

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



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

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

Note doing this on full $W$ matrix would result in Jacobian tensor!

But it is *sparse* – each output only affected by corresponding weight row

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

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

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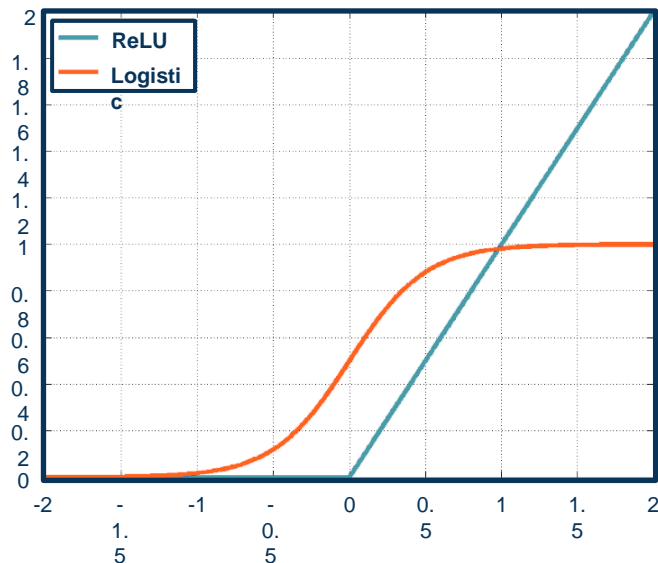
**Fully Connected (FC) Layer**

Georgia Tech

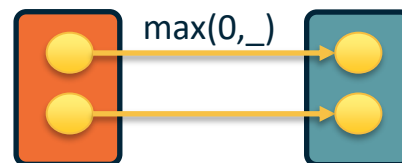We can employ **any differentiable (or piecewise differentiable) function**

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?
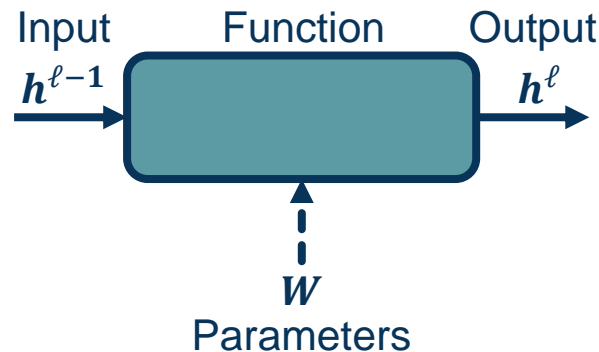


$$h^{\ell} = \max(0, h^{\ell-1})$$

max(0,_)

Georgia Tech

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

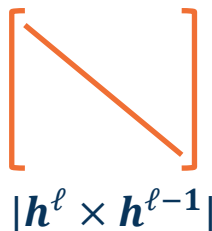- An output value affected only by **corresponding input value**

Max function **funnels gradients through selected max**

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

Input     Function     Output

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

$W$

Parameters

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

**Backward:** $\dfrac{\partial L}{\partial h^{\ell-1}} = \dfrac{\partial L}{\partial h^{\ell}} \quad \dfrac{\partial h^{\ell}}{\partial h^{\ell-1}}$

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

For diagonal

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

**Jacobian of ReLU**

Georgia Tech

4D input x:

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

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

4D output z:

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

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

4D dL/dz:

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

Upstream
gradient

Georgia
Tech

4D input x:

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

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

4D output z:

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

4D dL/dx:

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

[dz/dx] [dL/dz]

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

4D dL/dz:

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

Upstream gradient

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

Georgia
Tech

- 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!**
  - How can we implement this algorithmically to make these calculations automatic? **Automatic Differentiation**

Georgia
Tech

**Composition of Functions:** $f\big(g(x)\big) = (f \circ g)(x)$

**A complex function (e.g. defined by a neural network):**

$$f(x) = g_\ell\big(g_{\ell-1}(\ldots g_1(x))\big)$$

$$f(x) = g_\ell \circ g_{\ell-1} \ldots \circ g_1(x)$$

(Many of these will be parameterized)
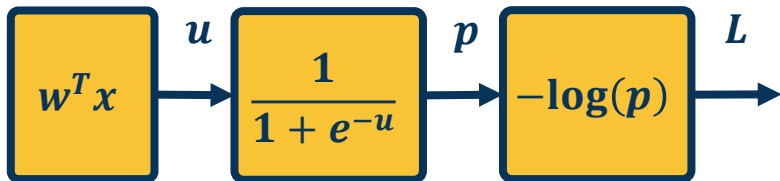
(Note you might find the opposite notation as well!)

**Composition of Functions & Chain Rule**

Georgia Tech

$$\bar{L} = 1$$

$$\bar{p} = \frac{\partial L}{\partial p} = -\frac{1}{p}$$

where $p = \sigma(w^T x)$ and $\sigma(x) = \frac{1}{1+e^{-x}}$

$$\bar{u} = \frac{\partial L}{\partial u} = \frac{\partial L}{\partial p}\frac{\partial p}{\partial u} = \bar{p}\,\sigma(1-\sigma)$$
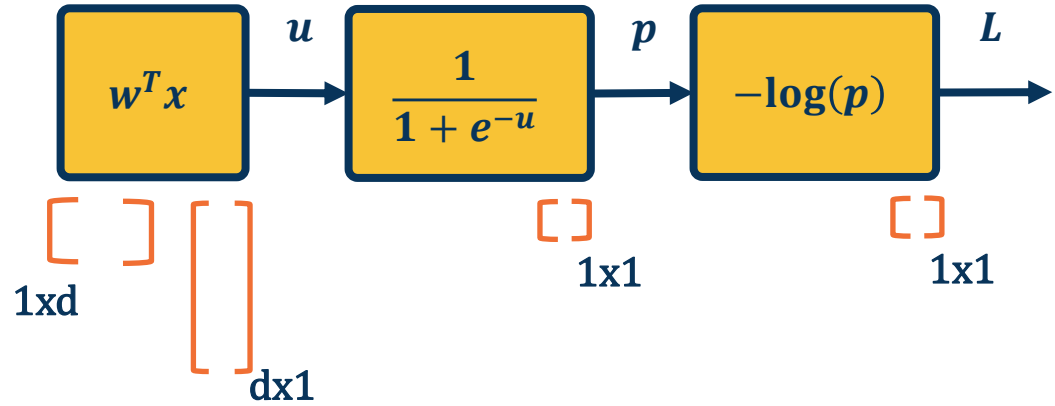
$$\bar{w} = \frac{\partial L}{\partial w} = \frac{\partial L}{\partial u}\frac{\partial u}{\partial w} = \bar{u}x^T$$

**We can do this in a combined way to see all terms together:**

$$\bar{w} = \frac{\partial L}{\partial p}\frac{\partial p}{\partial u}\frac{\partial u}{\partial w} = -\frac{1}{\sigma(w^T x)}\sigma(w^T x)(1-\sigma(w^T x))x^T$$

$$= -\left(1-\sigma(w^T x)\right)x^T$$

**This effectively shows gradient flow along path from $L$ to $w$**

Georgia Tech

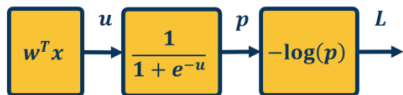The chain rule can be computed as a **series of scalar, vector, and matrix linear algebra operations**

$$w^T x$$  $u$  $\dfrac{1}{1 + e^{-u}}$  $p$  $-\log(p)$  $L$

[ ]  1xd

[ ]  dx1

[ ]  1x1

[ ]  1x1

**Extremely efficient** in graphics processing units (GPUs)

$$\overline{w} = -\frac{1}{\sigma(w^T x)} \sigma(w^T x)(1 - \sigma(w^T x))x^T$$

[ ]  1x1

[ ]  1x1

[ ]  1x1

[ ]  1xd

Georgia Tech

$$L = 1$$
$$\bar{p} = \frac{\partial L}{\partial p} = -\frac{1}{p}$$

where $p = \sigma(w^T x)$ and $\sigma(x) = \frac{1}{1+e^{-x}}$

$$\bar{u} = \frac{\partial L}{\partial u} = \frac{\partial L}{\partial p}\frac{\partial p}{\partial u} = \bar{p}\,\sigma(1-\sigma)$$

$$\bar{w} = \frac{\partial L}{\partial w} = \frac{\partial L}{\partial u}\frac{\partial u}{\partial w} = \bar{u}x^T$$
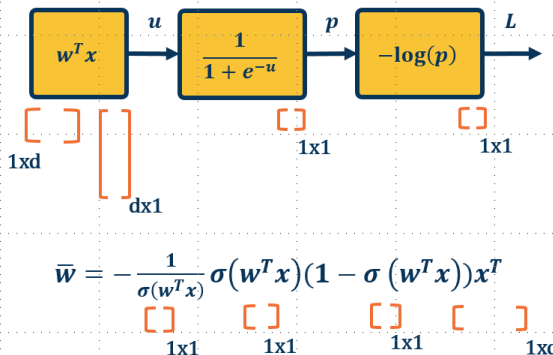
**We can do this in a combined way to see all terms together:**

$$\bar{w} = \frac{\partial L}{\partial p}\frac{\partial p}{\partial u}\frac{\partial u}{\partial w} = -\frac{1}{\sigma(w^T x)}\sigma(w^T x)(1-\sigma(w^T x))x^T$$
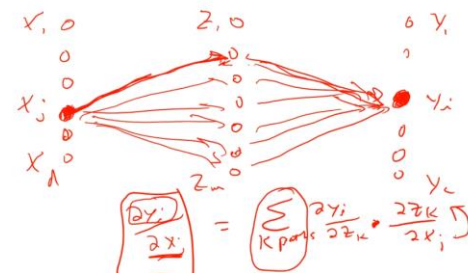$$= -\big(1-\sigma(w^T x)\big)x^T$$

**This effectively shows gradient flow along path from $L$ to $w$**
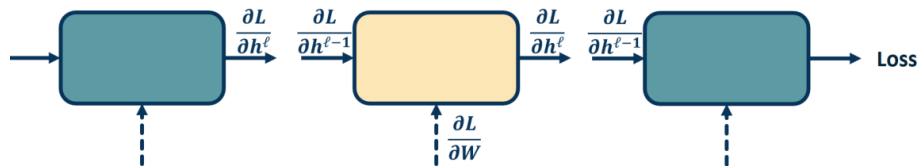
**Computation Graph / Global View of Chain Rule**

**Computational / Tensor View**

1xd    dx1    1x1    1x1

$$\bar{w} = -\frac{1}{\sigma(w^T x)}\sigma(w^T x)(1-\sigma(w^T x))x^T$$

1x1    1x1    1x1    1xd

**Graph View**

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

$\dfrac{\partial L}{\partial h^{\ell}}$    $\dfrac{\partial L}{\partial h^{\ell-1}}$    $\dfrac{\partial L}{\partial h^{\ell}}$    $\dfrac{\partial L}{\partial h^{\ell-1}}$    Loss
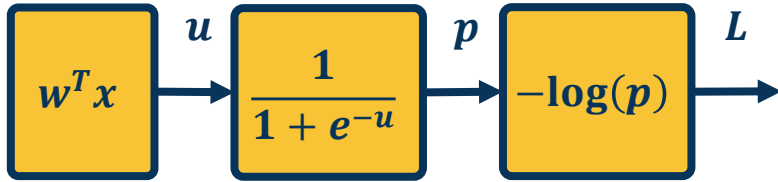
$\dfrac{\partial L}{\partial W}$

**Backpropagation View (Recursive Algorithm)**

**Different Views of Equivalent Ideas**

Georgia Tech

- **Backpropagation:** Recursive, modular algorithm for chain rule + gradient descent

- **When we move to vectors and matrices:**
  - Composition of functions (scalar)
  - Composition of functions (vectors/matrices)
  - Jacobian view of chain rule
  - Can view entire set of calculations as linear algebra operations (matrix-vector or matrix-matrix multiplication)

- **Automatic differentiation**:
  - Reduction of modules to simple operations we know (simple multiplication, etc.)
  - Automatically build computation graph in background as write code
  - Automatically compute gradients via backward pass

**Summary**

Georgia Tech

**Automatic differentiation:**

- Carries out this procedure for us on arbitrary graphs

- Knows derivatives of primitive functions

- As a result, we just define these (forward) functions **and don't even need to specify the gradient (backward) functions!**

Georgia Tech