

Topics:

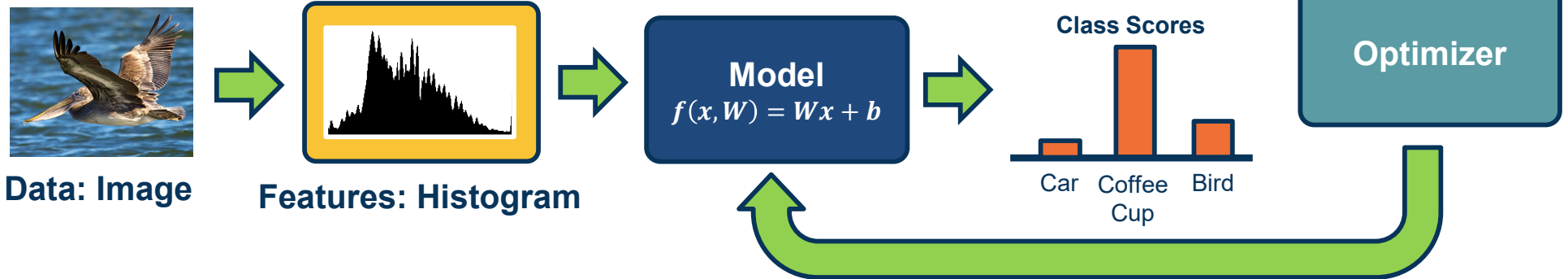
- Gradient Descent
- Neural Networks

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

**ZSOLT KIRA**

- **Assignment 1 out!**
  - **Due Feb 3<sup>th</sup> (with grace period Feb 5<sup>th</sup>)**
  - Start now, start now, start now!
  - Start now, start now, start now!
  - Start now, start now, start now!
- **Piazza**
  - Be active!!!
- **Office hours**
  - Lots of special topics (e.g. Assignment 1, Matrix Calculus, etc. )
- Note: Course will start to get math heavy!

- Input (and representation)
- Functional form of the model
  - Including parameters
- Performance measure to improve
  - Loss or objective function
- Algorithm for finding best parameters
  - Optimization algorithm

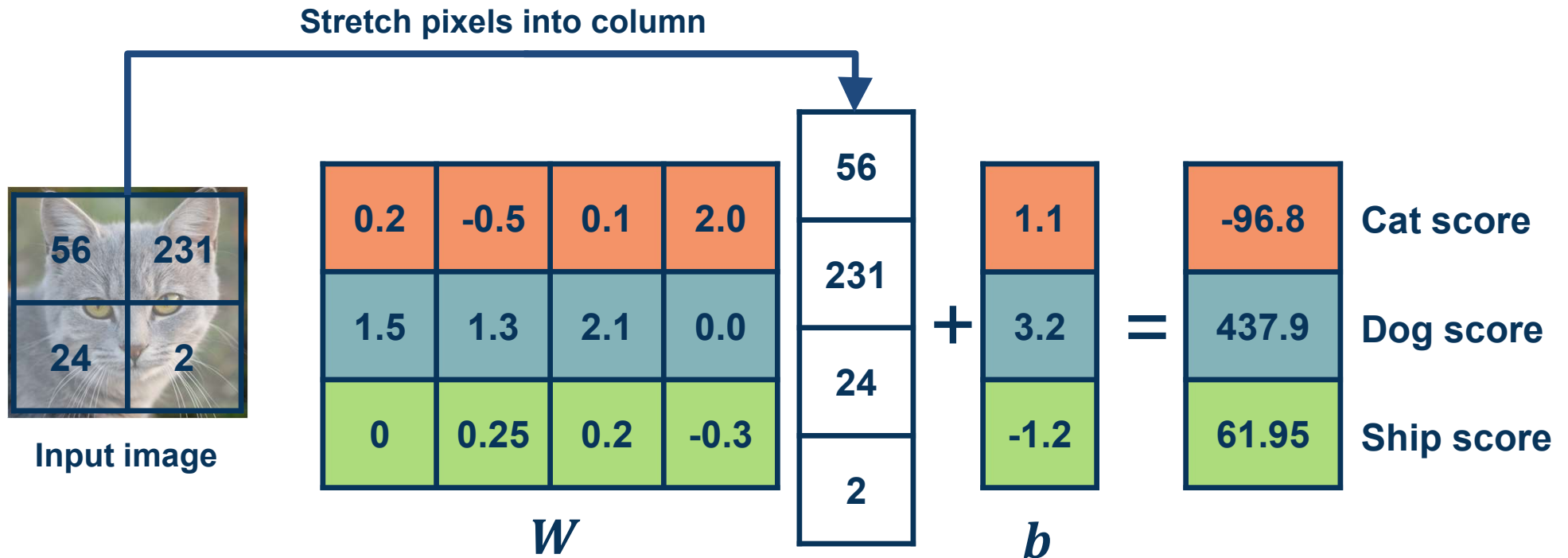


## Components of a Parametric Model

- ◆ Input: **Vector**
- ◆ Functional form of the model: **Softmax(Wx)**
- ◆ Performance measure to improve: **Cross-Entropy**
- ◆ Algorithm for finding best parameters: **Gradient Descent**
  - ◆ Compute  $\frac{\partial L}{\partial w_i}$
  - ◆ Update Weights  $w_i = w_i - \alpha \frac{\partial L}{\partial w_i}$

So far

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



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

Example

# Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

Probabilities must be  $\geq 0$

$$P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax Function

Probabilities must sum to 1

$$L_i = -\log P(Y = y_i|X = x_i)$$

cat  
car  
frog

3.2  
5.1  
-1.7

Unnormalized log-probabilities / logits

exp

24.5  
164.0  
0.18

Unnormalized probabilities

normalize

0.13  
0.87  
0.00

Probabilities



$$L_i = -\log(0.13)$$

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

Cross-Entropy Loss Example

Often, we add a **regularization term** to the loss function

### L1 Regularization

$$L_i = |y - Wx_i|^2 + |W|$$

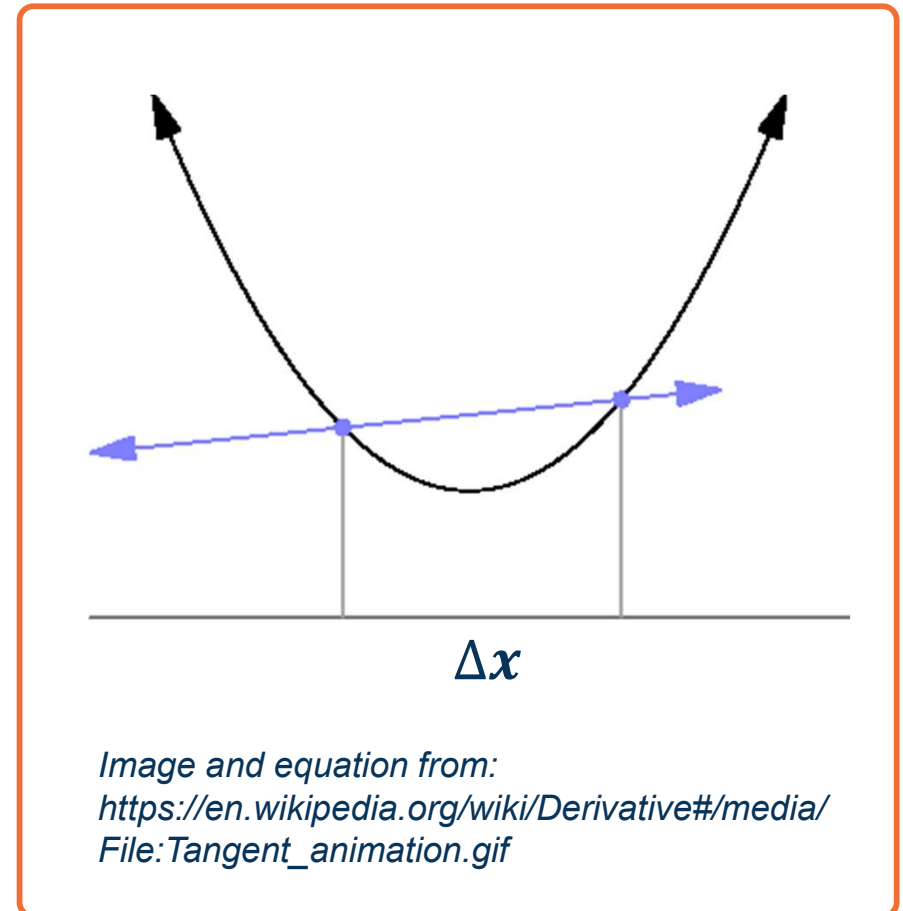
### Example regularizations:

- ◆ L1/L2 on weights (encourage small values)

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

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

- 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





This idea can be turned into an **algorithm (gradient descent)**

1. Choose a model:  $f(x, W) = Wx$
2. Choose loss function:  $L_i = (y - Wx_i)^2$
3. Calculate partial derivative for each parameter:  $\frac{\partial L}{\partial w_i}$
4. Update the parameters:  $w_i = w_i - \frac{\partial L}{\partial w_i}$

Add learning rate to prevent too big of a step:  $w_i = w_i - \alpha \frac{\partial L}{\partial w_i}$

5. Repeat (from Step 3)

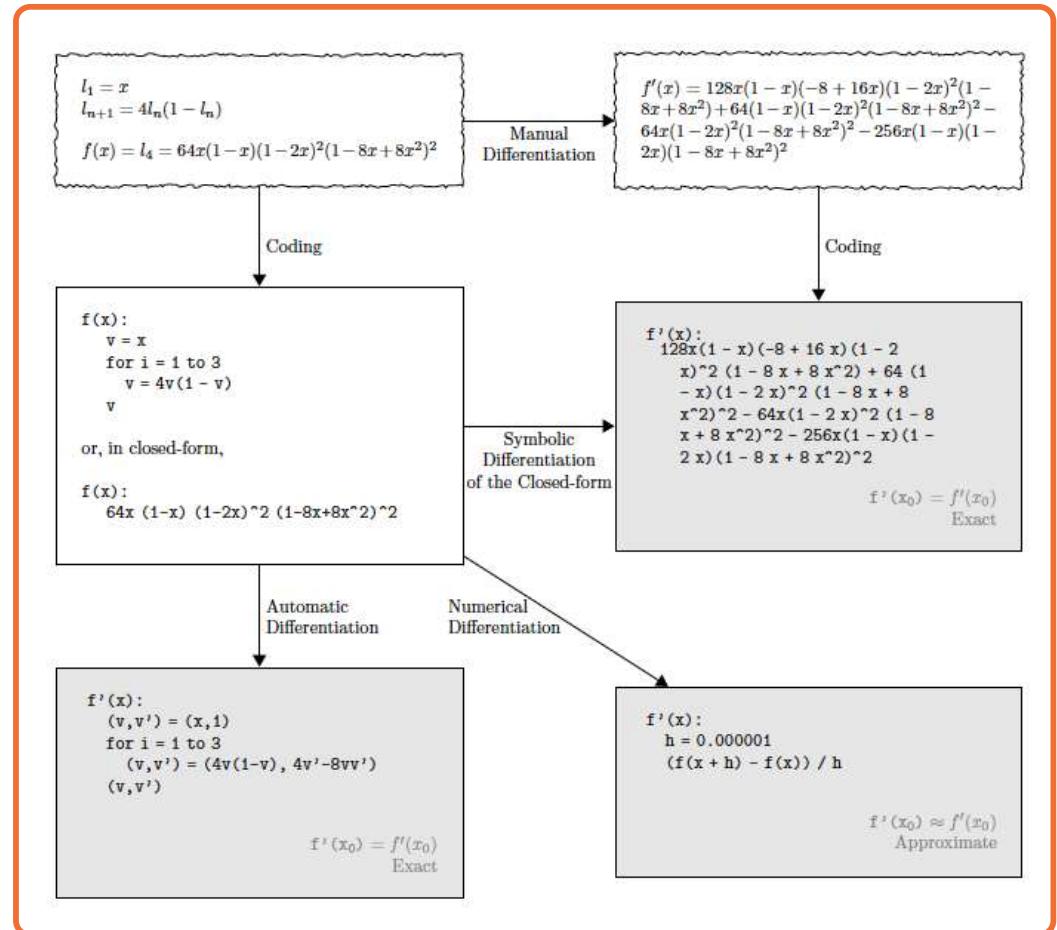
Gradient descent is guaranteed to converge under some conditions

- ◆ For example, learning rate has to be appropriately reduced throughout training
- ◆ It will converge to a *local* minima
  - ◆ Small changes in weights would not decrease the loss
- ◆ It turns out that some of the local minima that it finds in practice (if trained well) are still pretty good!

We know how to compute the **model output and loss function**

Several ways to compute  $\frac{\partial L}{\partial w_i}$

- ◆ Manual differentiation
- ◆ Symbolic differentiation
- ◆ Numerical differentiation
- ◆ Automatic differentiation



For some functions, we can analytically derive the partial derivative

## Example:

### Function

$$f(\mathbf{w}, \mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i$$

(Assume  $\mathbf{w}$  and  $\mathbf{x}_i$  are column vectors, so same as  $\mathbf{w} \cdot \mathbf{x}_i$ )

### Loss

$$\sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

**Dataset:** N examples (indexed by  $k$ )

### Update Rule

$$\mathbf{w}_j \leftarrow \mathbf{w}_j + 2\alpha \sum_{i=1}^N \delta_i \mathbf{x}_{ij}$$

## Derivation of Update Rule

$$L = \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

Gradient descent tells us we should update  $\mathbf{w}$  as follows to minimize  $L$ :

$$\mathbf{w}_j \leftarrow \mathbf{w}_j - \alpha \frac{\partial L}{\partial \mathbf{w}_j}$$

So what's  $\frac{\partial L}{\partial \mathbf{w}_j}$ ?

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{w}_j} &= \sum_{i=1}^N \frac{\partial}{\partial \mathbf{w}_j} (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \\ &= \sum_{k=1}^N 2(y_i - \mathbf{w}^T \mathbf{x}_i) \frac{\partial}{\partial \mathbf{w}_j} (y_i - \mathbf{w}^T \mathbf{x}_i) \\ &= -2 \sum_{i=1}^N \delta_i \frac{\partial}{\partial \mathbf{w}_j} \mathbf{w}^T \mathbf{x}_i \\ &= -2 \sum_{i=1}^N \delta_i \frac{\partial}{\partial \mathbf{w}_j} \sum_{k=1}^N \mathbf{w}_k \mathbf{x}_{ik} \\ &= -2 \sum_{i=1}^N \delta_i \mathbf{x}_{ij} \end{aligned}$$

...where...  
 $\delta_i = y_i - \mathbf{w}^T \mathbf{x}_i$

If we add a **non-linearity (sigmoid)**, derivation is more complex

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

First, one can derive that:  $\sigma'(x) = \sigma(x)(1 - \sigma(x))$

$$f(\mathbf{x}) = \sigma\left(\sum_k w_k x_k\right)$$

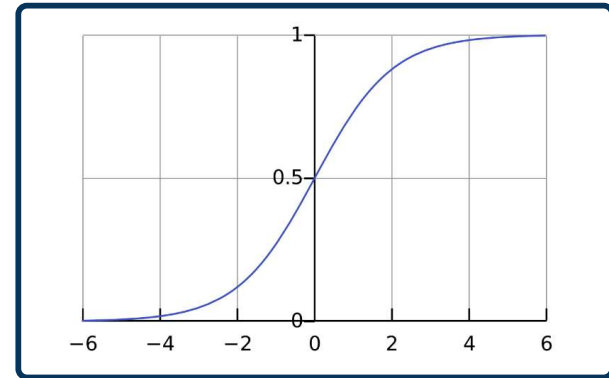
$$L = \sum_i \left( y_i - \sigma\left(\sum_k w_k x_{ik}\right) \right)^2$$

$$\frac{\partial L}{\partial w_j} = \sum_i 2 \left( y_i - \sigma\left(\sum_k w_k x_{ik}\right) \right) \left( -\frac{\partial}{\partial w_j} \sigma\left(\sum_k w_k x_{ik}\right) \right)$$

$$= \sum_i -2 \left( y_i - \sigma\left(\sum_k w_k x_{ik}\right) \right) \sigma'\left(\sum_k w_k x_{ik}\right) \frac{\partial}{\partial w_j} \sum_k w_k x_{ik}$$

$$= \sum_i -2 \delta_i \sigma(\mathbf{d}_i) (1 - \sigma(\mathbf{d}_i)) x_{ij}$$

where  $\delta_i = y_i - f(x_i)$        $\mathbf{d}_i = \sum_k w_k x_{ik}$



The sigmoid perception update rule:

$$w_j \leftarrow w_j + 2\alpha \sum_{k=1}^N \delta_i \sigma_i (1 - \sigma_i) x_{ij}$$

where  $\sigma_i = \sigma\left(\sum_{j=1}^d w_j x_{ij}\right)$

$$\delta_i = y_i - \sigma_i$$

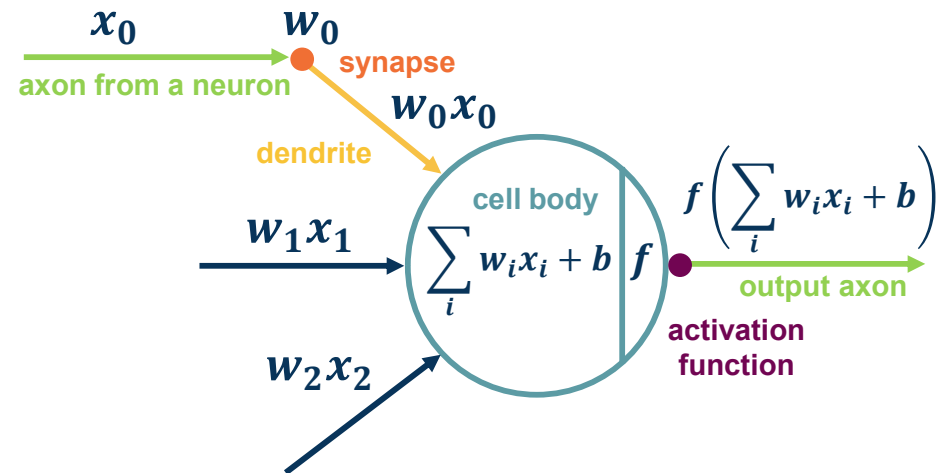
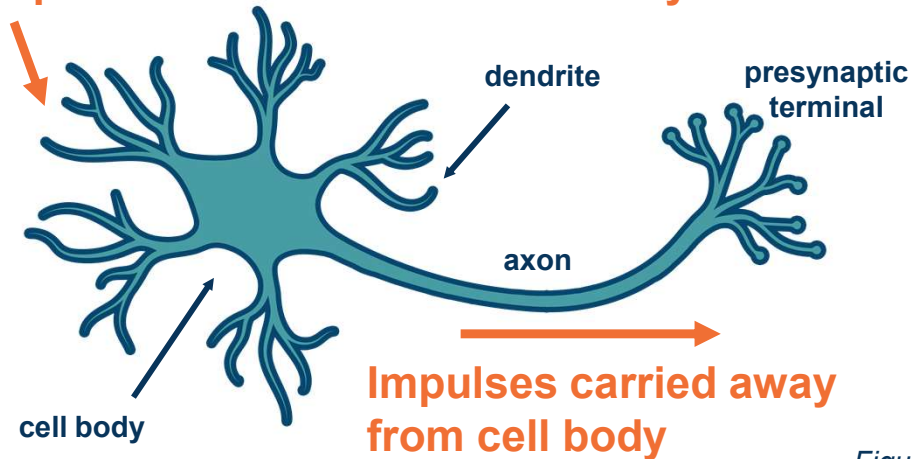
Adding a Non-Linear Function

**Neural  
Network  
View of a  
Linear  
Classifier**

A simple **neural network** has similar structure as our linear classifier:

- A neuron takes input (firings) from other neurons (-> **input to linear classifier**)
- The inputs are summed in a weighted manner (-> **weighted sum**)
  - Learning is through a modification of the weights
- If it receives enough input, it “fires” (threshold or if weighted sum plus bias is high enough)

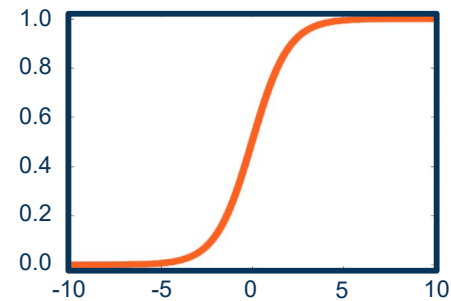
Impulses carried toward cell body



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

## Origins of the Term Neural Network

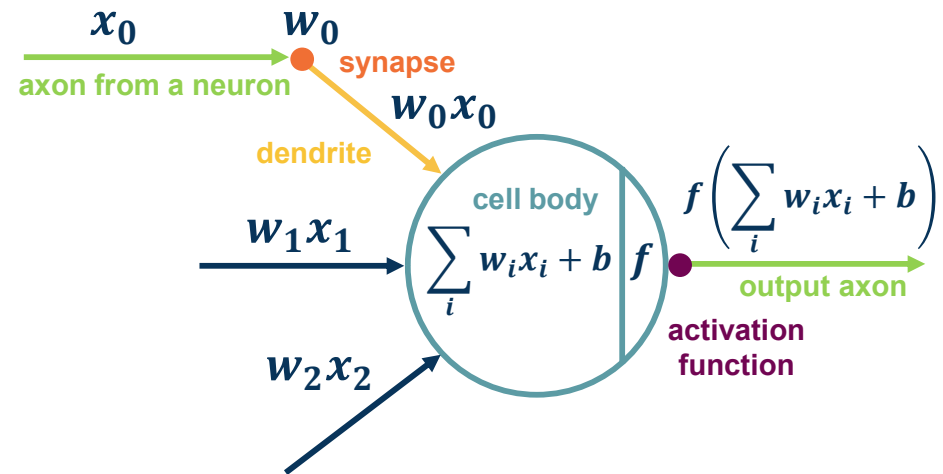
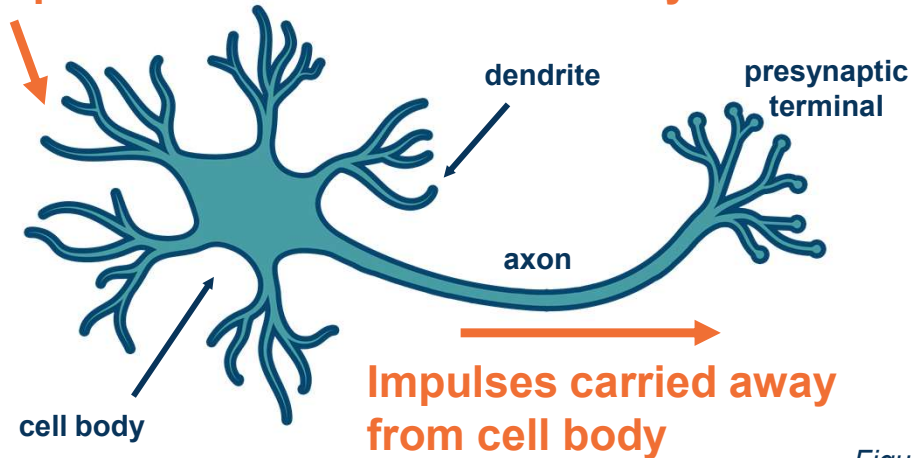
As we did before, the output of a neuron can be modulated by a non-linear function (e.g. sigmoid)



**Sigmoid  
Activation  
Function**

$$\frac{1}{1 + e^{-x}}$$

Impulses carried toward cell body



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

## Adding Non-Linearities



We can have **multiple** neurons connected to the same input

Corresponds to a multi-class classifier

- Each output node outputs the score for a class

$$f(x, W) = \sigma(Wx + b) \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1m} & b_1 \\ w_{21} & w_{22} & \cdots & w_{2m} & b_2 \\ w_{31} & w_{32} & \cdots & w_{3m} & b_3 \end{bmatrix}$$

- Often called fully connected layers
  - Also called a linear *projection layer*

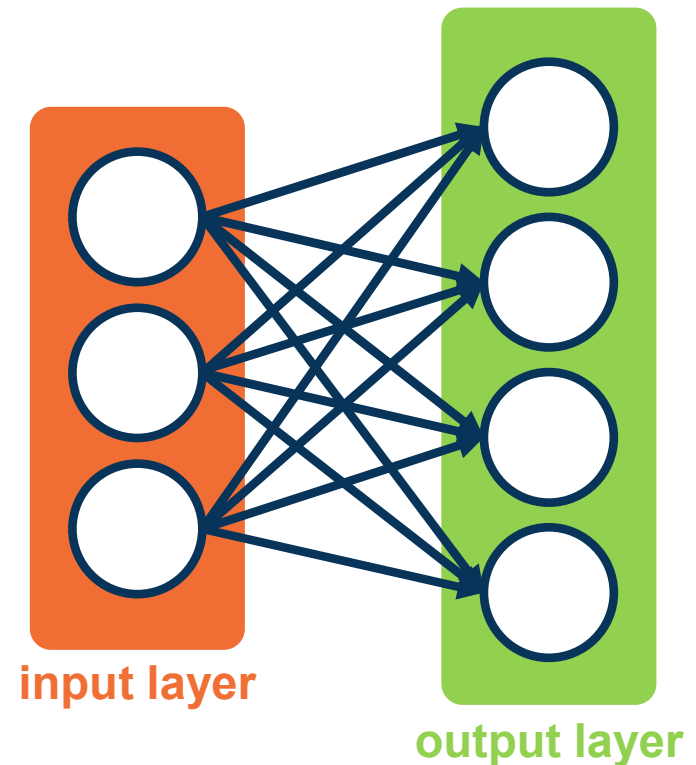


Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

- Each input/output is a **neuron (node)**
- A linear classifier (+ optional non-linearity) is called a **fully connected layer**
- Connections are represented as **edges**
- Output of a particular neuron is referred to as **activation**
- This will be expanded as we view computation in a neural network as a **graph**

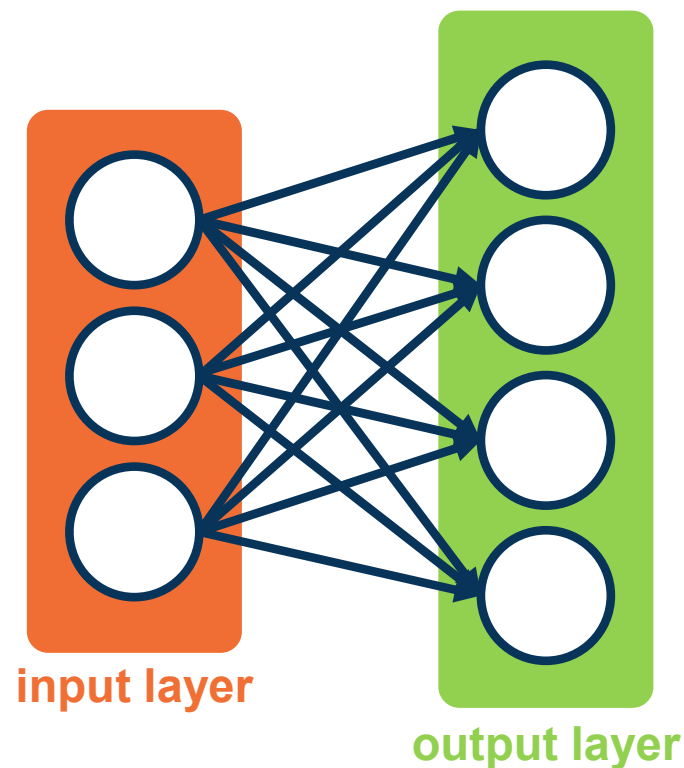


Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

We can **stack** multiple layers together

- ◆ Input to second layer is output of first layer

Called a **2-layered neural network** (input is not counted)

Because the middle layer is neither input or output, and we don't know what their values represent, we call them **hidden** layers

- ◆ We will see that they end up learning effective features

This **increases** the representational power of the function!

- ◆ Two layered networks can represent any continuous function

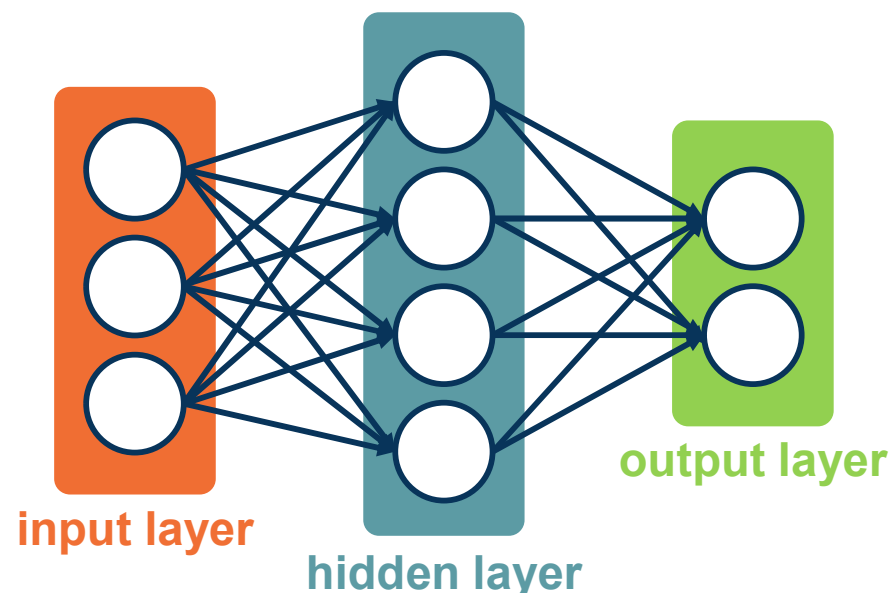


Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

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)

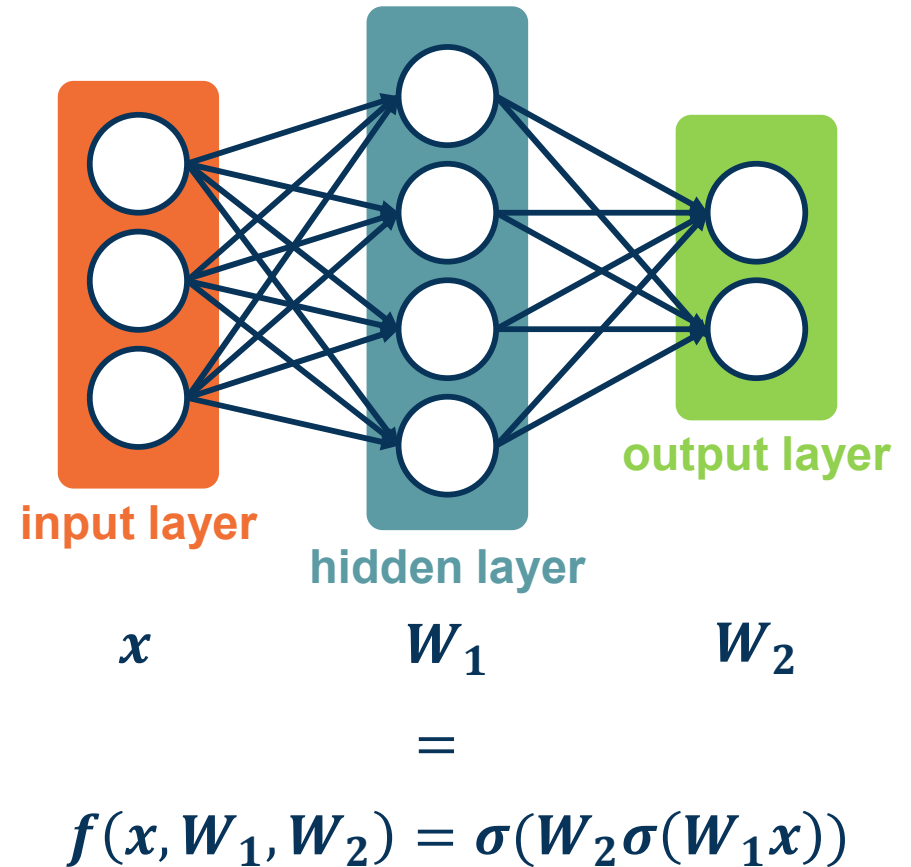


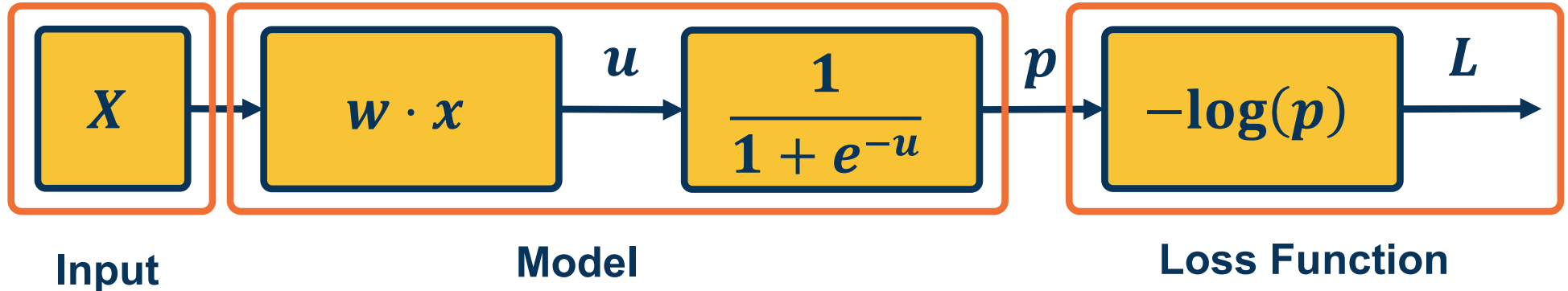
Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

## The Linear Algebra View

A **linear classifier** can be broken down into:

- ◆ Input
- ◆ A function of the input
- ◆ A loss function

It's all just one function that can be **decomposed** into building blocks



What Does a Linear Classifier Consist of?

**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**:

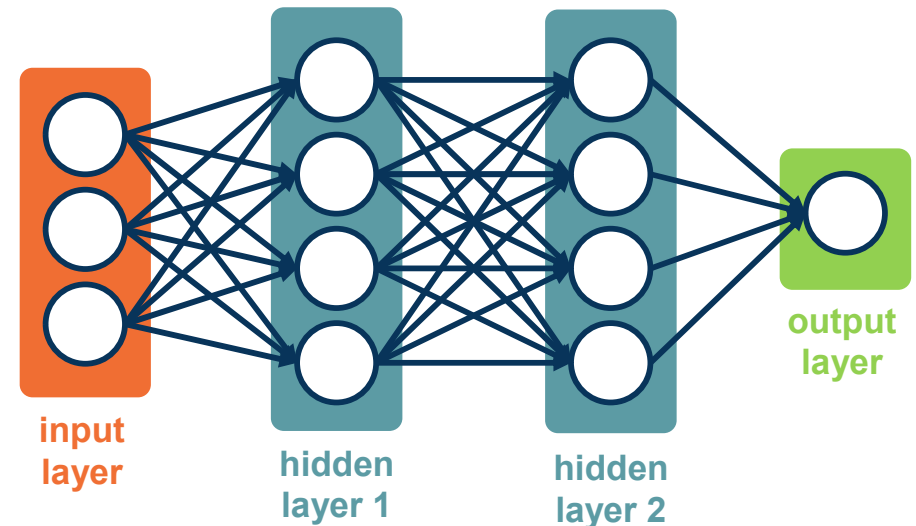
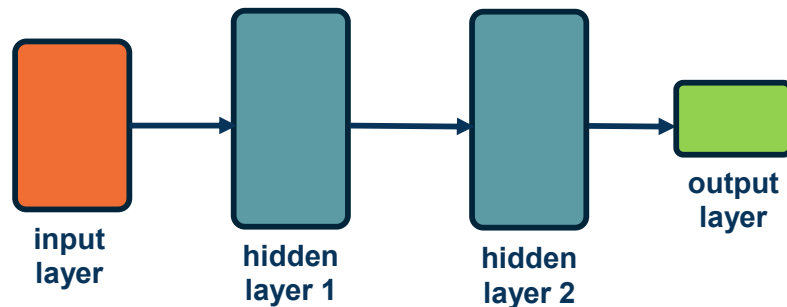


Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

**Adding More Layers!**

# Computation Graphs

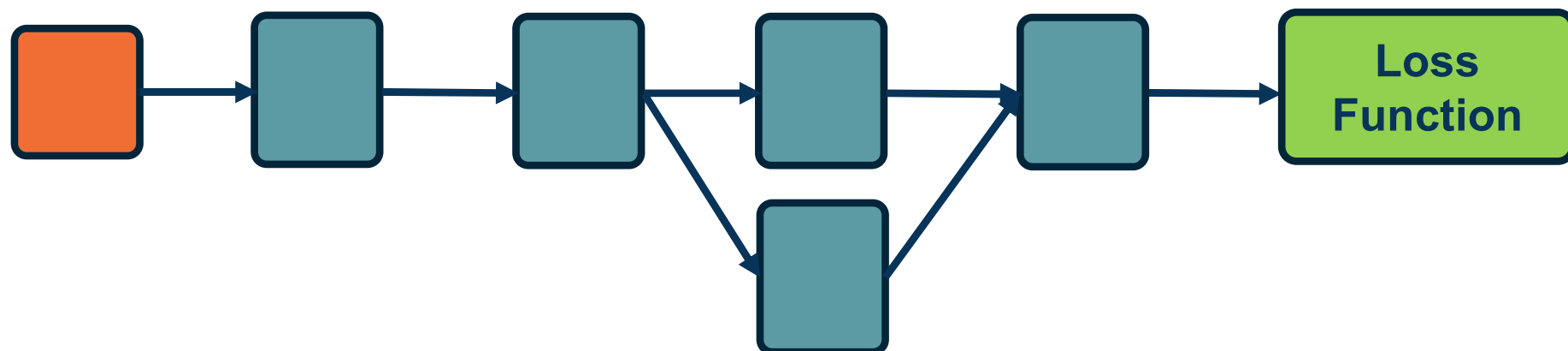
Functions can be made **arbitrarily complex** (subject to memory and computational limits), e.g.:

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

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

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

Composition can have **some structure**



**Adding Even More Layers**

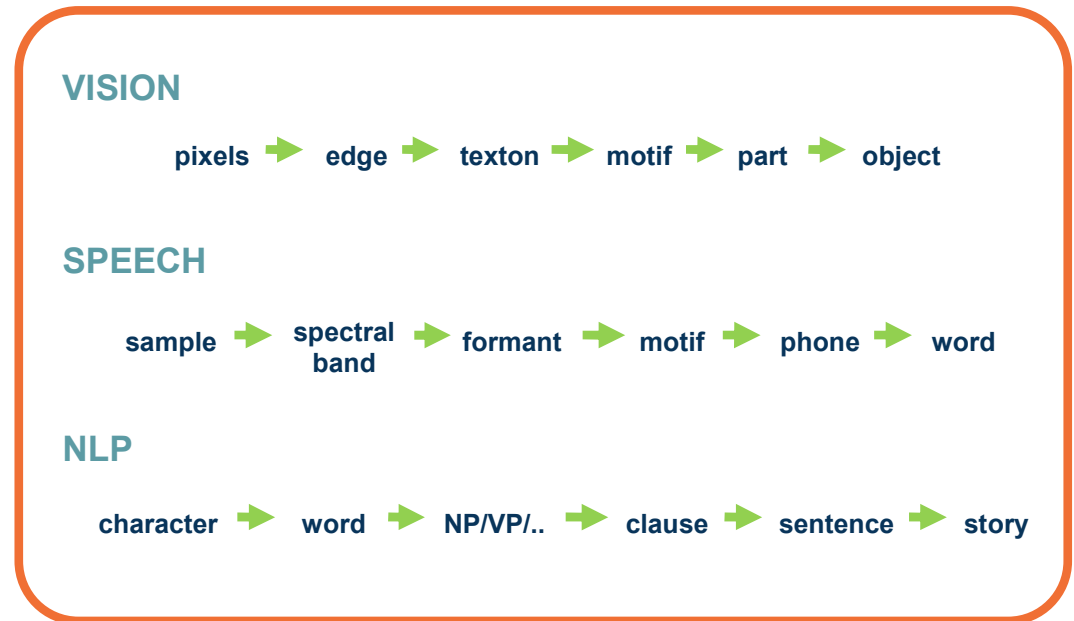


The world is **compositional!**

We want our **model** to reflect this

Empirical and theoretical evidence that it makes **learning complex functions easier**

Note that **prior state of art engineered features** often had this compositionality as well

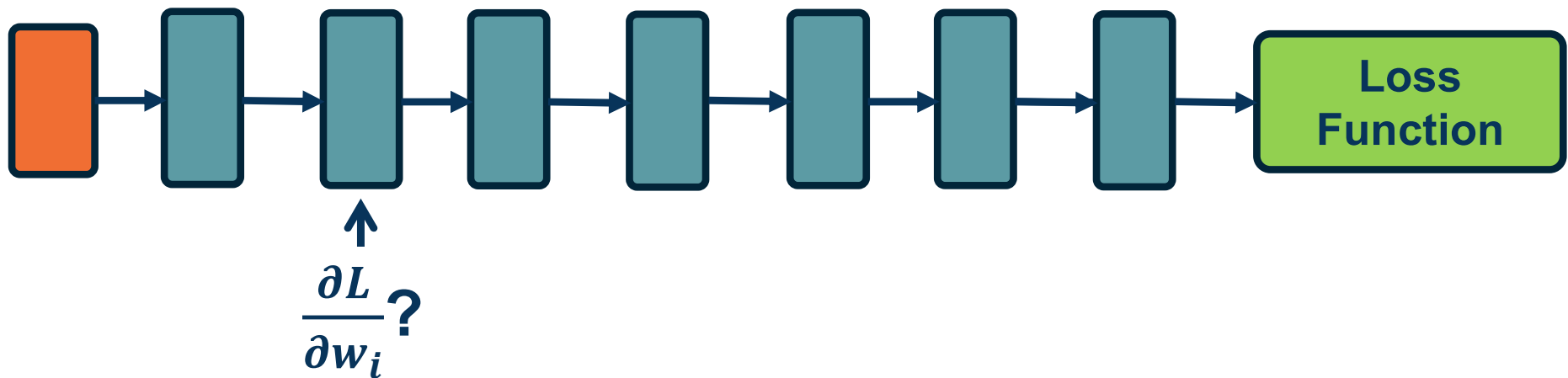


*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

◆ **Pixels -> edges -> object parts -> objects**

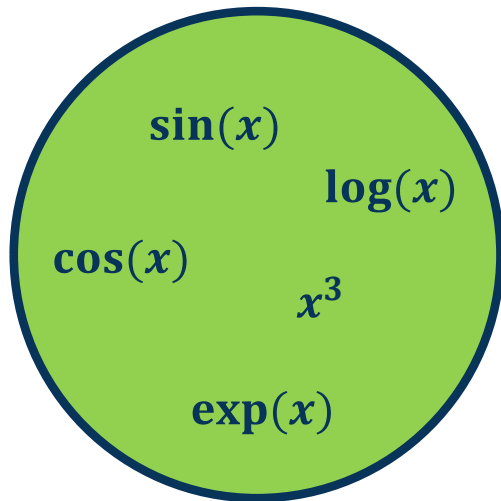
**Compositionality**

- ◆ 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



## Computing Gradients in Complex Function

Given a library of simple functions



Compose into a  
→  
complicate function

$$-\log\left(\frac{1}{1 + e^{-w \cdot x}}\right)$$



*Adapted from slides by: Marc'Aurelio Ranzato, Yann LeCun*

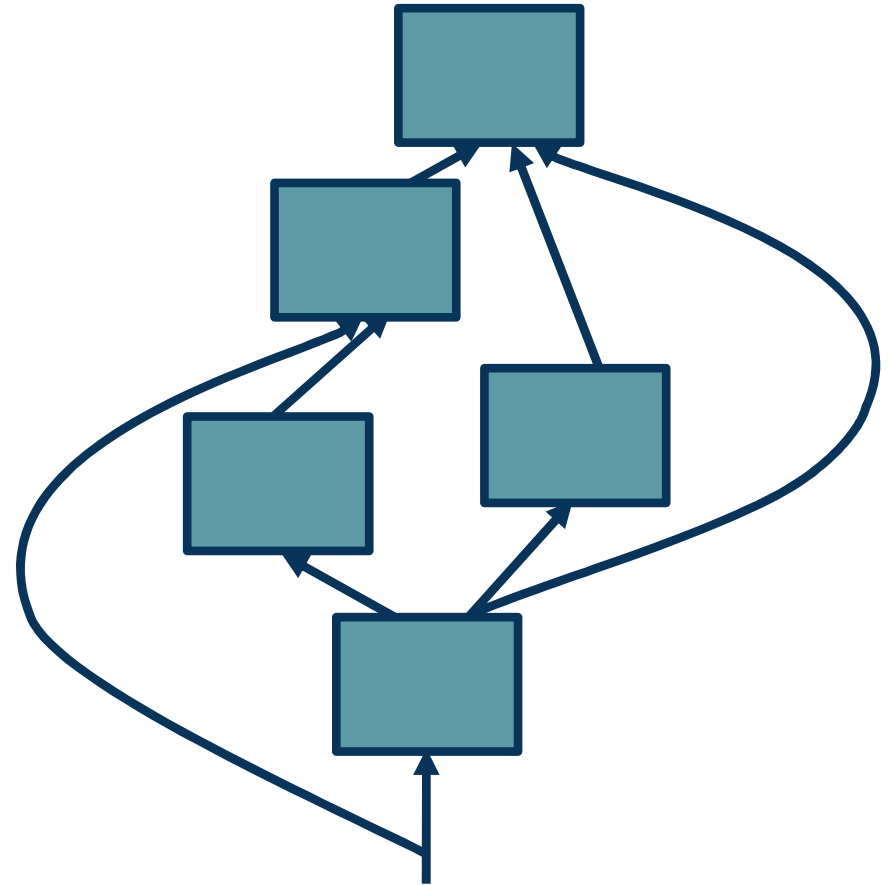
## Decomposing a Function

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

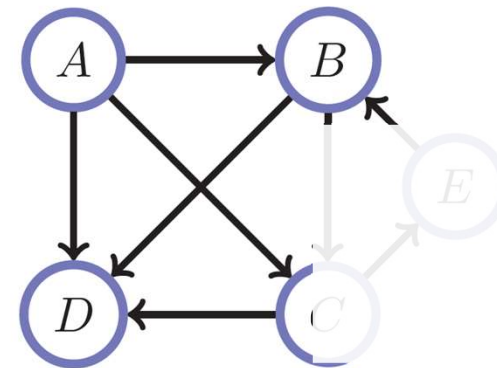
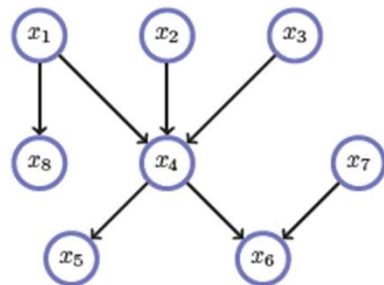
A **training algorithm** will then process this graph, **one module at a time**



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

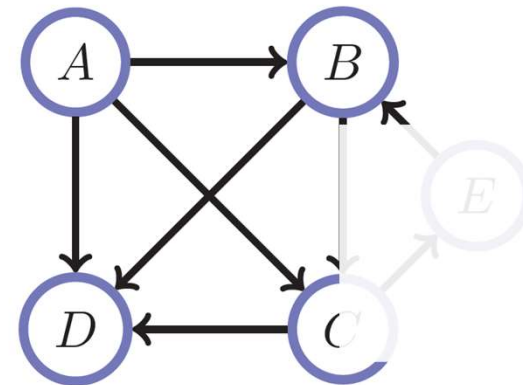
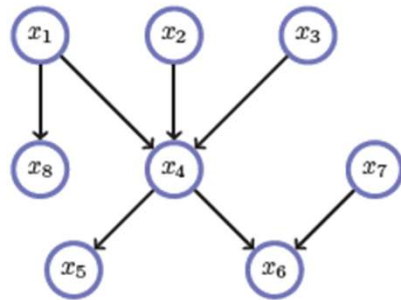
# Directed Acyclic Graphs (DAGs)

- Exactly what the name suggests
  - Directed edges
  - No (directed) cycles
  - Underlying undirected cycles okay

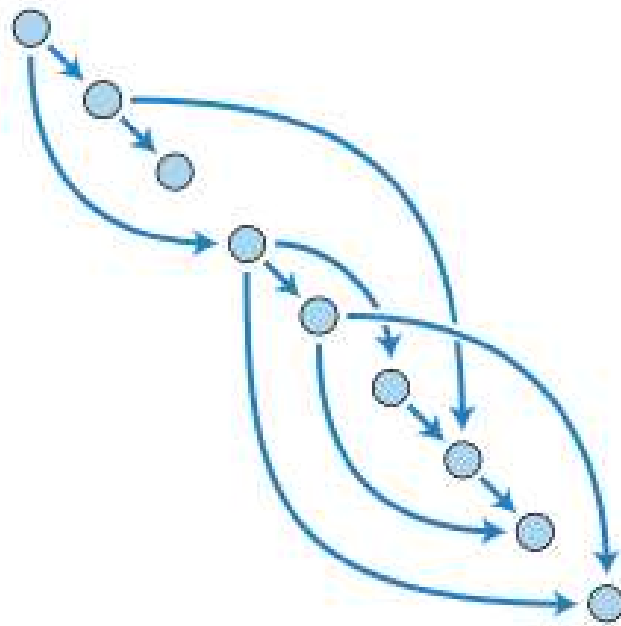


# Directed Acyclic Graphs (DAGs)

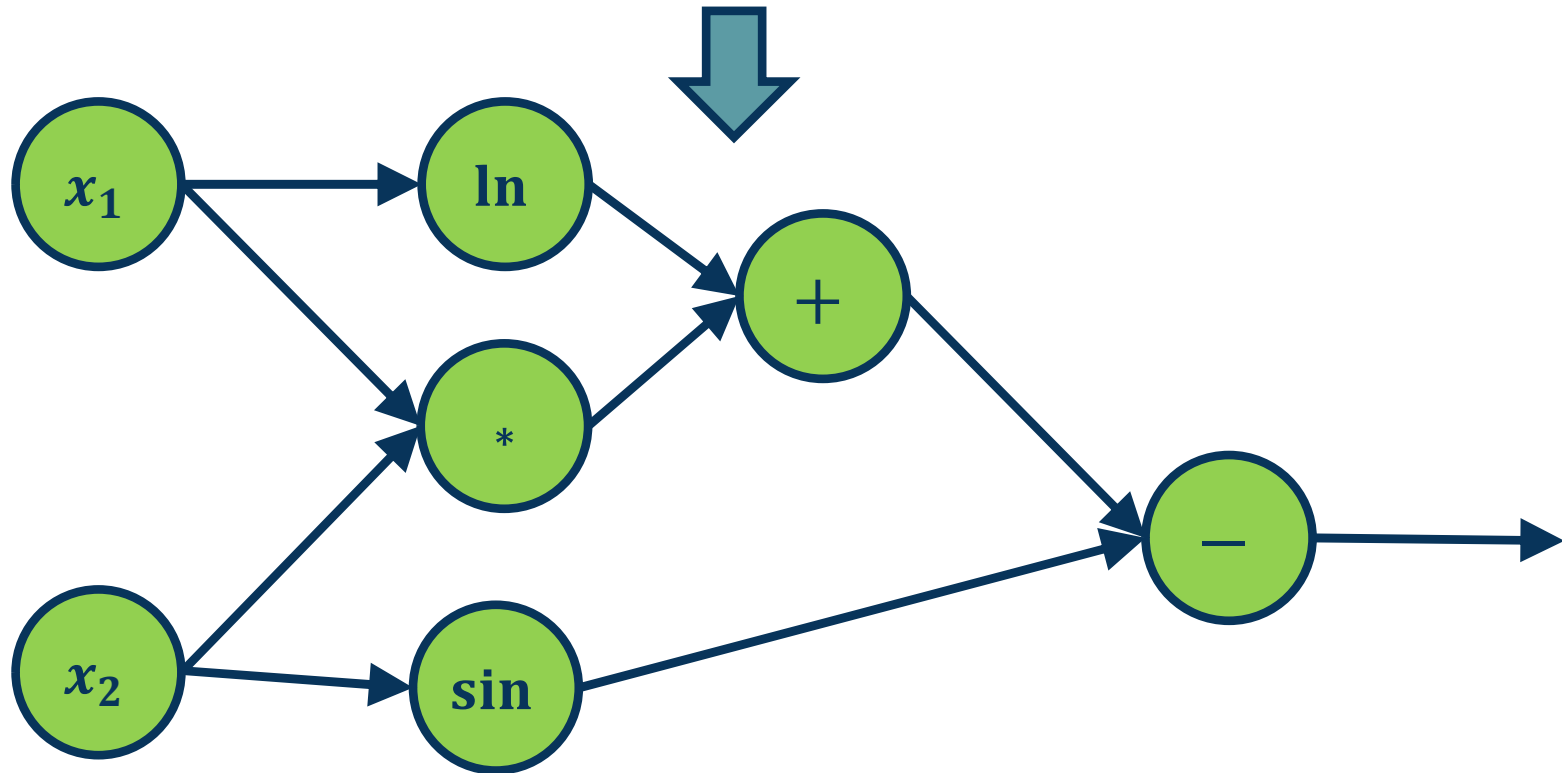
- Concept
  - Topological Ordering



# Directed Acyclic Graphs (DAGs)



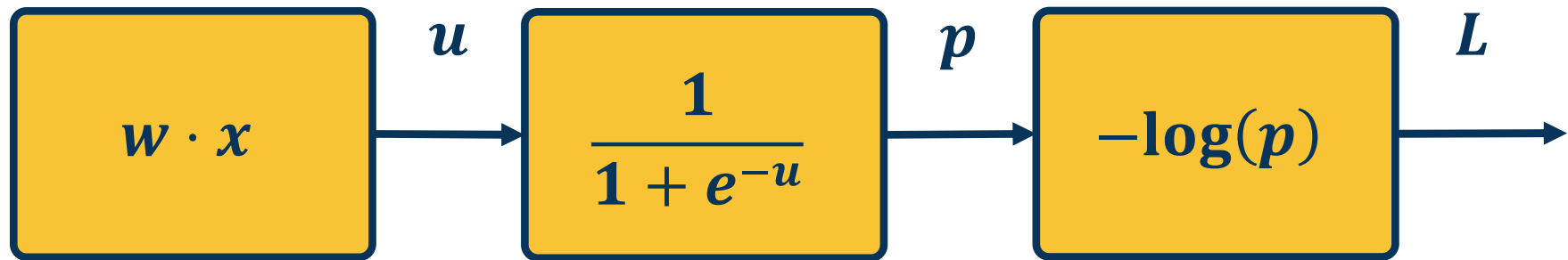
$$f(x_1, x_2) = \ln(x_1) + x_1x_2 - \sin(x_2)$$



Example



$$-\log\left(\frac{1}{1 + e^{-w \cdot x}}\right)$$



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

# Backpropagation

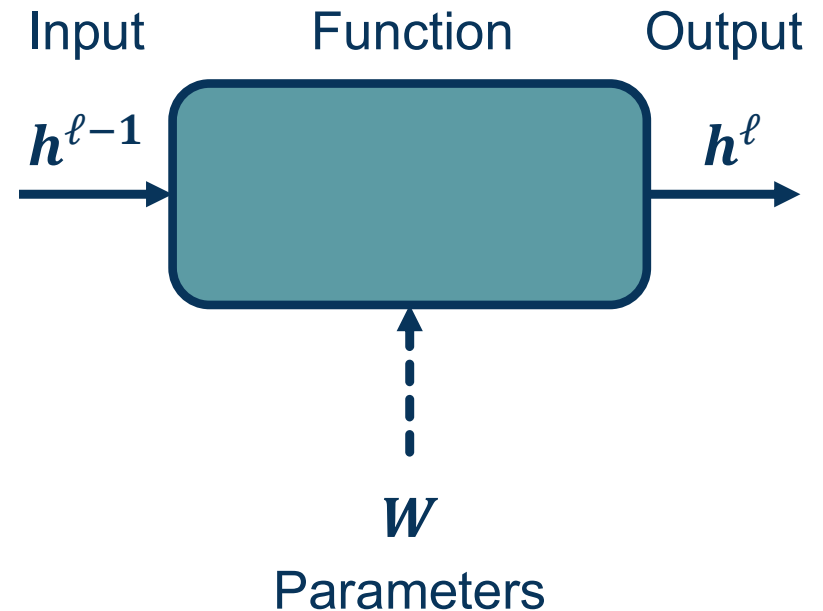
Given this computation graph, the training algorithm will:

- Calculate the current model's outputs (called the **forward pass**)
- Calculate the gradients for each module (called the **backward pass**)

Backward pass is a recursive algorithm that:

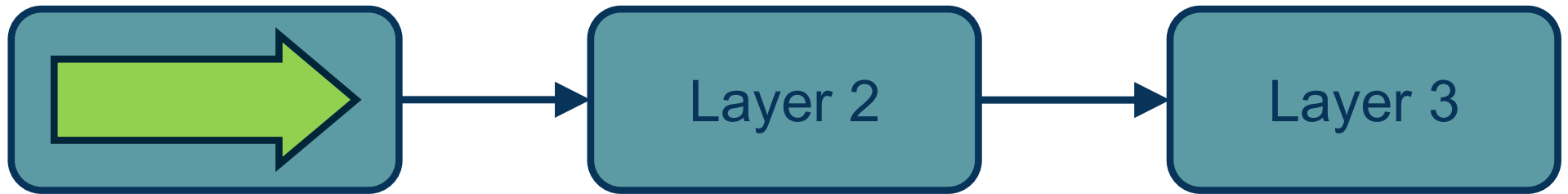
- Starts at **loss function** where we know how to calculate the gradients
- Progresses back through the modules
- Ends in the **input layer** where we do not need gradients (no parameters)

This algorithm is called **backpropagation**



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

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



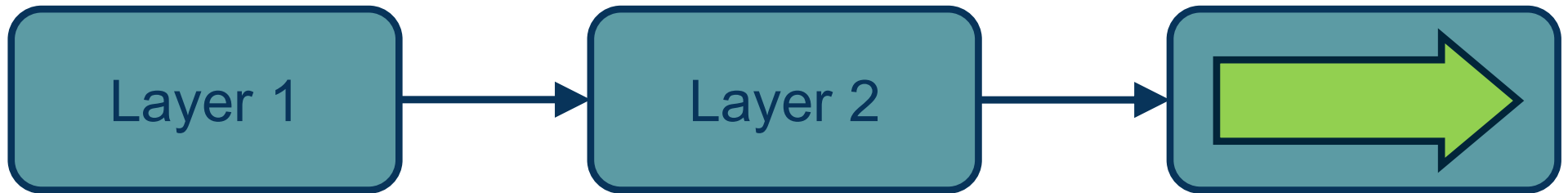
*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

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



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

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



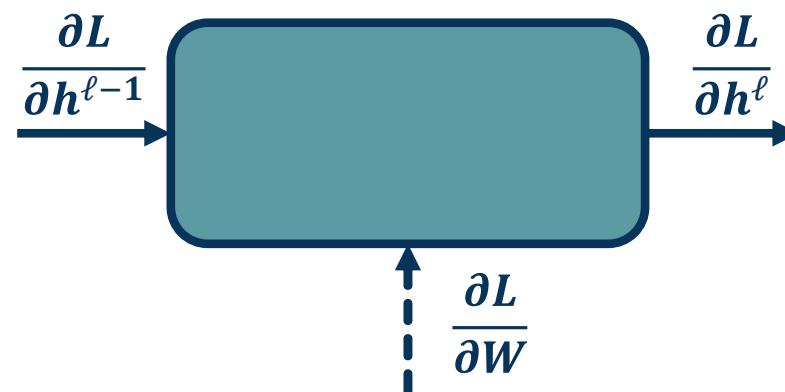
Note that we must store the **intermediate outputs of all layers!**

- ◆ 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*

In the **backward pass**, we seek to calculate the gradients of the loss with respect to the module's parameters

- Assume that we have the gradient of the loss with respect to the **module's outputs** (given to us by upstream module)
- We will also pass the gradient of the loss with respect to the **module's inputs**
  - This is not required for update the module's weights, but passes the gradients back to the previous module

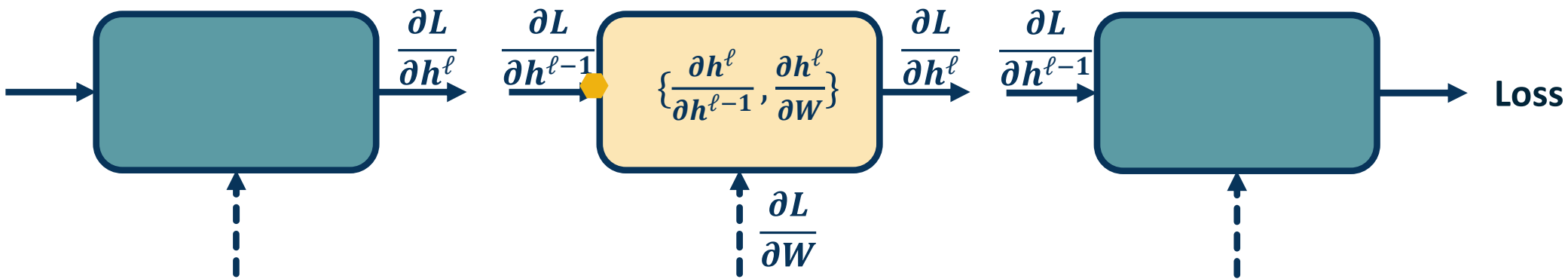


### Problem:

- We can compute local gradients:  $\left\{ \frac{\partial h^{\ell}}{\partial h^{\ell-1}}, \frac{\partial h^{\ell}}{\partial W} \right\}$
- We are given:  $\frac{\partial L}{\partial h^{\ell}}$
- Compute:  $\left\{ \frac{\partial L}{\partial h^{\ell-1}}, \frac{\partial L}{\partial W} \right\}$

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

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



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

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



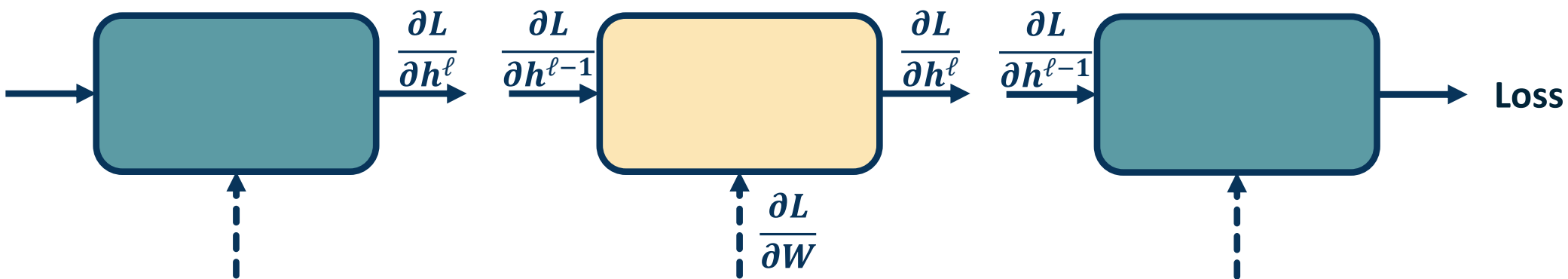
- ◆ We can compute **local gradients**:  $\left\{ \frac{\partial h^\ell}{\partial h^{\ell-1}}, \frac{\partial h^\ell}{\partial W} \right\}$
- ◆ This is just the **derivative of our function** with respect to its parameters and inputs!

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

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

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

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



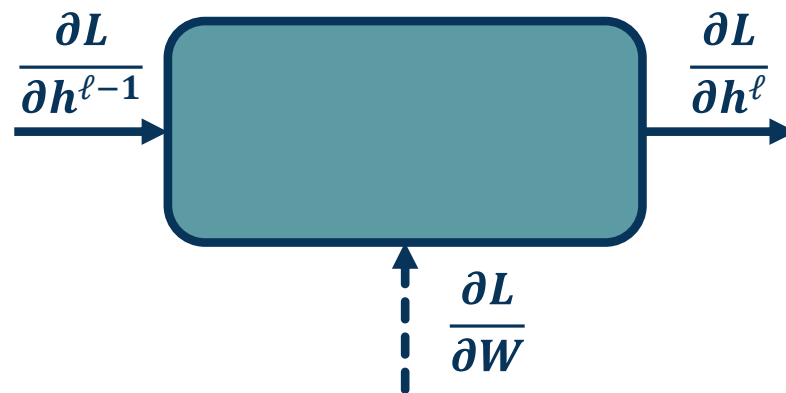
- ◆ We will use the *chain rule* to do this:

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

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

- **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**)

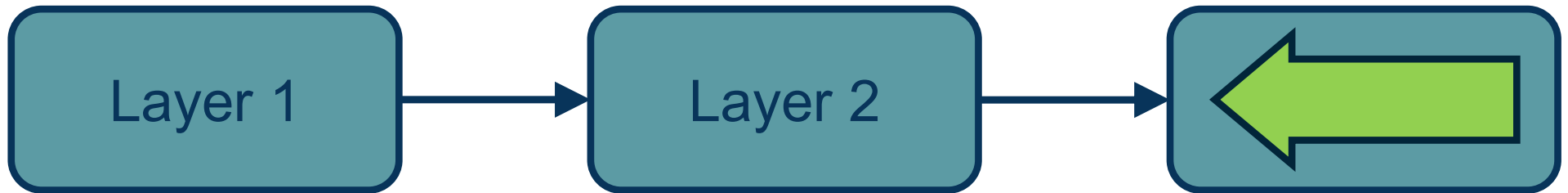
- **Gradient of loss w.r.t. weights:**  $\frac{\partial L}{\partial W} = \frac{\partial L}{\partial h^{\ell}} \frac{\partial h^{\ell}}{\partial W}$



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

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

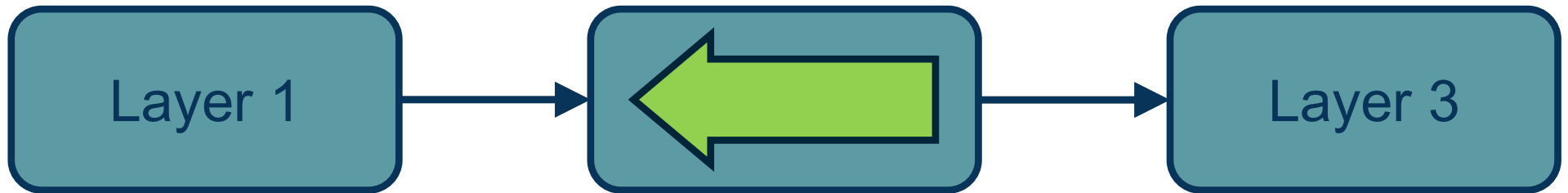
**Step 2: Compute Gradients wrt parameters: Backward Pass**



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

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

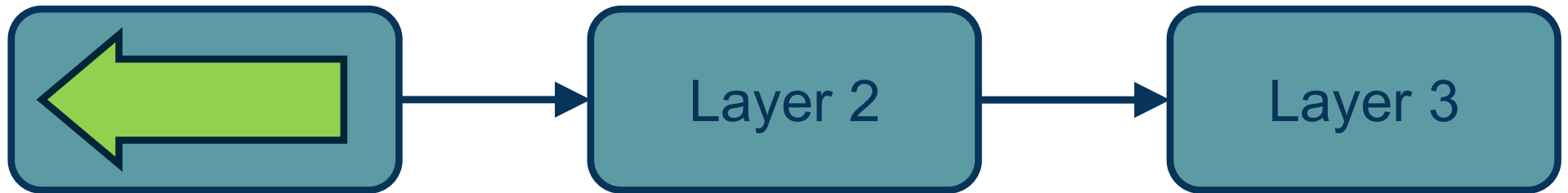
**Step 2: Compute Gradients wrt parameters: Backward Pass**



*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**

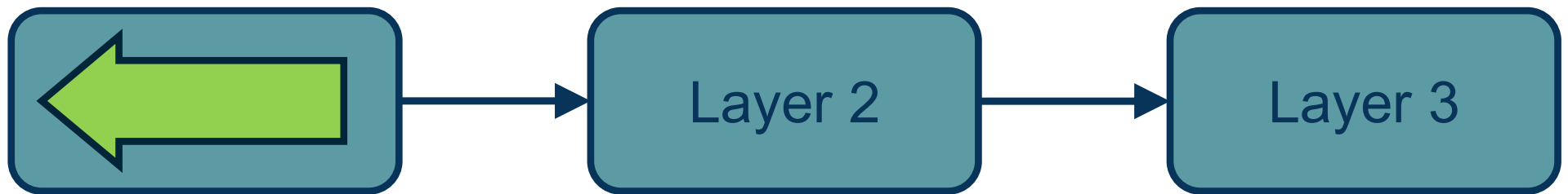


*Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun*

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

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

**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*

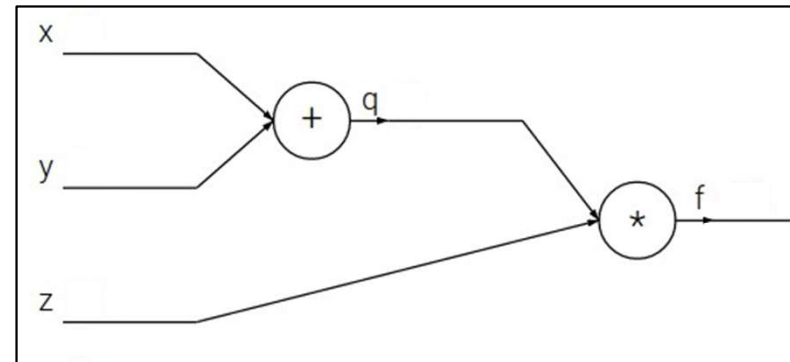
# Backpropagation: a simple example

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



# Backpropagation: a simple example

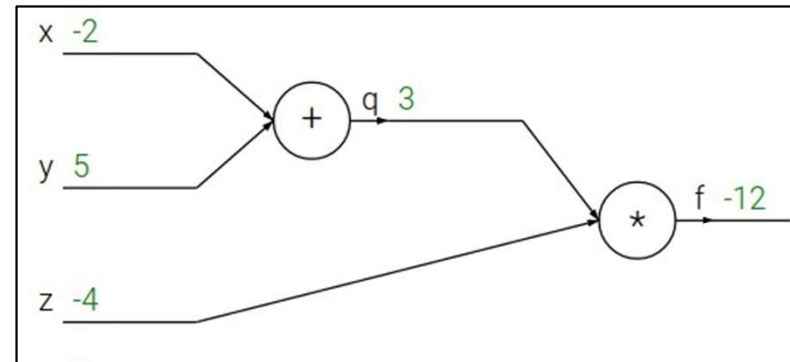
$$f(x, y, z) = (x + y)z$$



# Backpropagation: a simple example

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

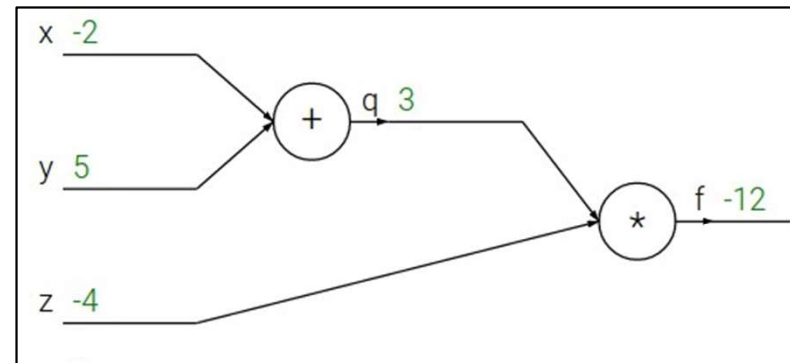
e.g.  $x = -2, y = 5, z = -4$



# Backpropagation: a simple example

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

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



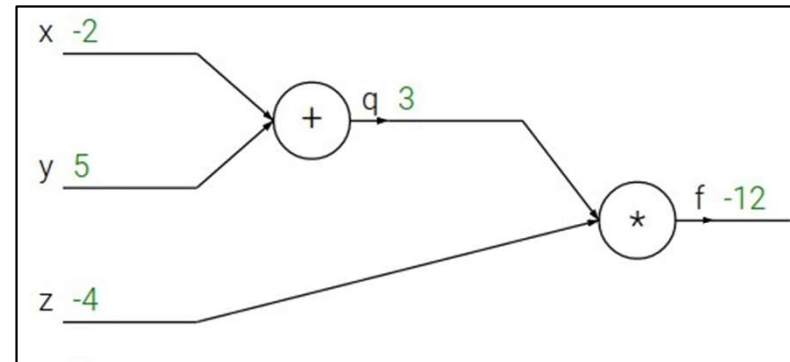
Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

# Backpropagation: a simple example

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

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

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



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

# Backpropagation: a simple example

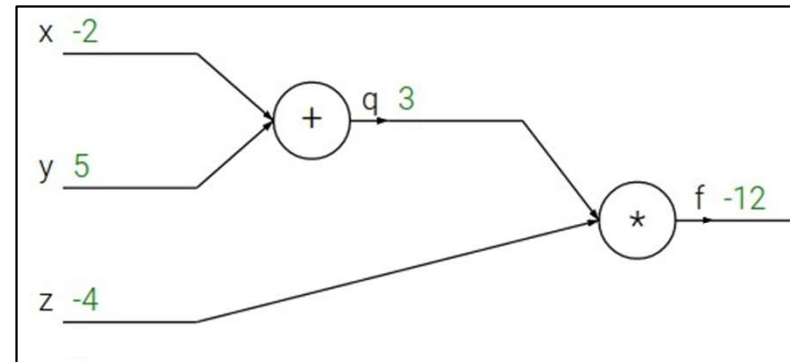
$$f(x, y, z) = (x + y)z$$

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

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

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

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



# Backpropagation: a simple example

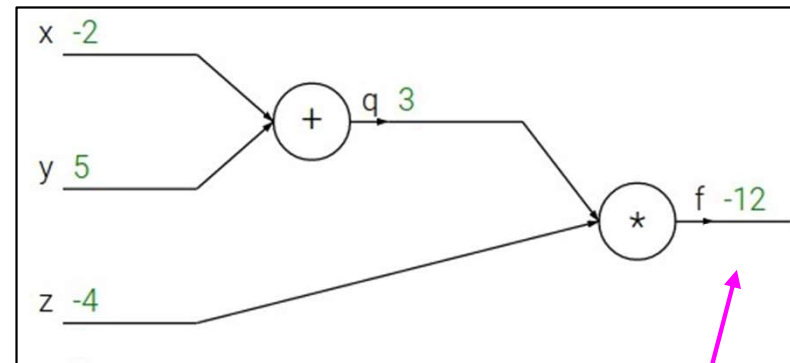
$$f(x, y, z) = (x + y)z$$

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

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

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

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



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

# Backpropagation: a simple example

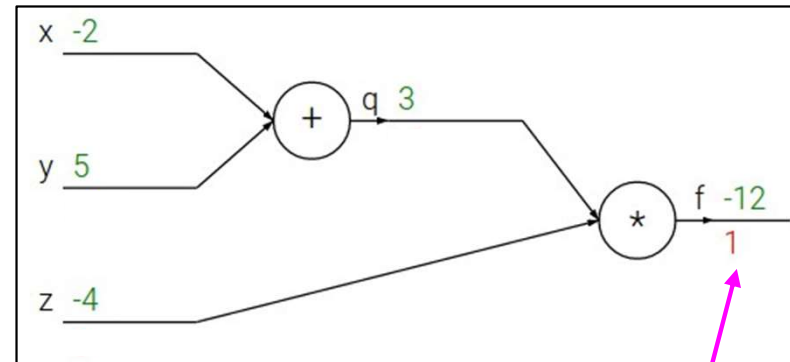
$$f(x, y, z) = (x + y)z$$

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

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

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

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



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

# Backpropagation: a simple example

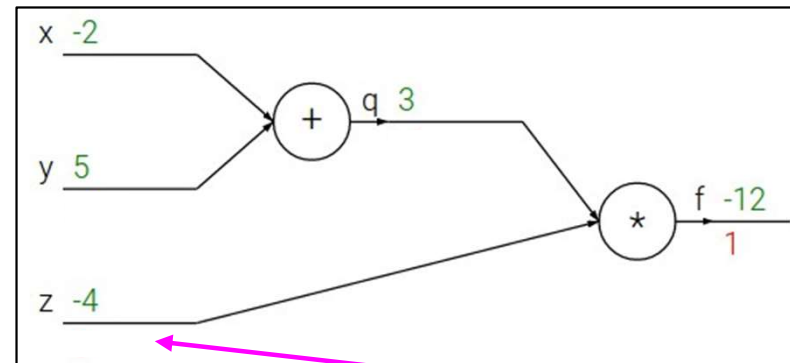
$$f(x, y, z) = (x + y)z$$

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

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

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

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



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



# Backpropagation: a simple example

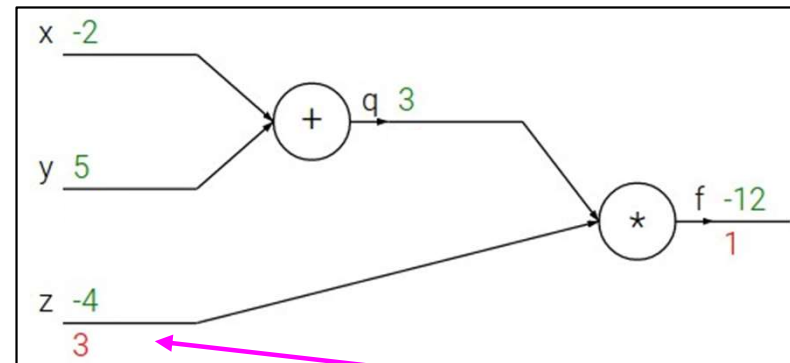
$$f(x, y, z) = (x + y)z$$

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

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

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

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



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

# Backpropagation: a simple example

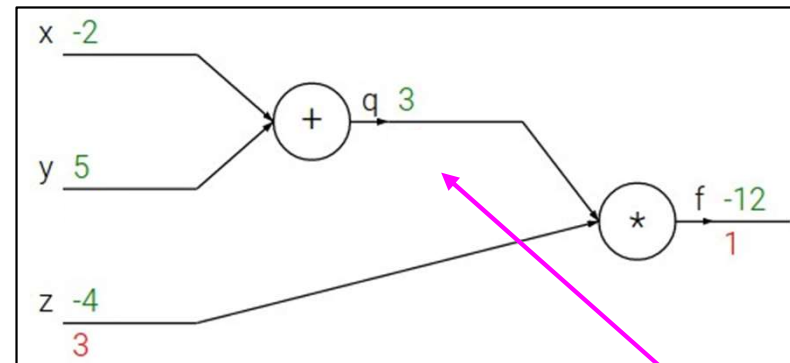
$$f(x, y, z) = (x + y)z$$

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

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

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

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



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

# Backpropagation: a simple example

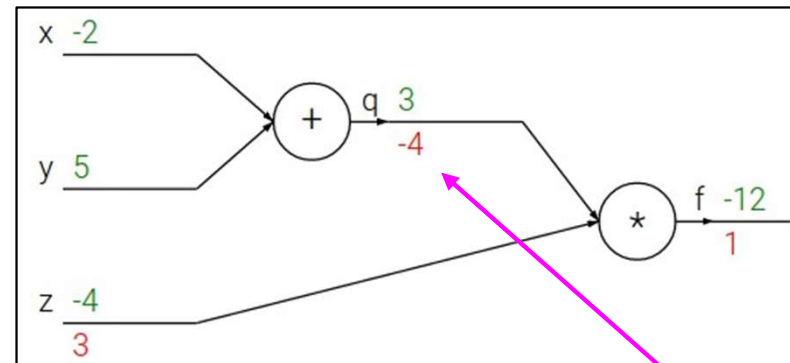
$$f(x, y, z) = (x + y)z$$

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

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

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

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



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

# Backpropagation: a simple example

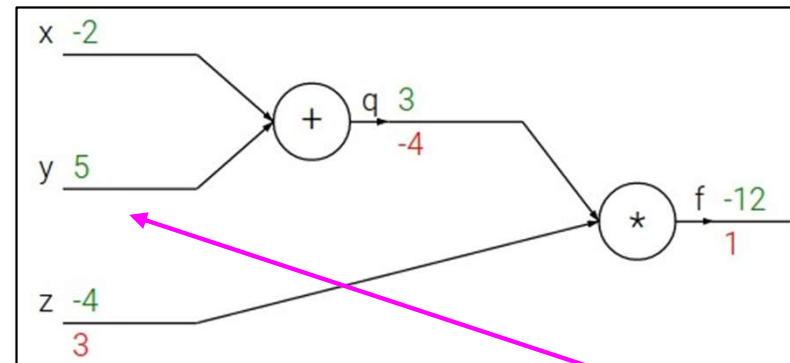
$$f(x, y, z) = (x + y)z$$

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

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

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

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

# Backpropagation: a simple example

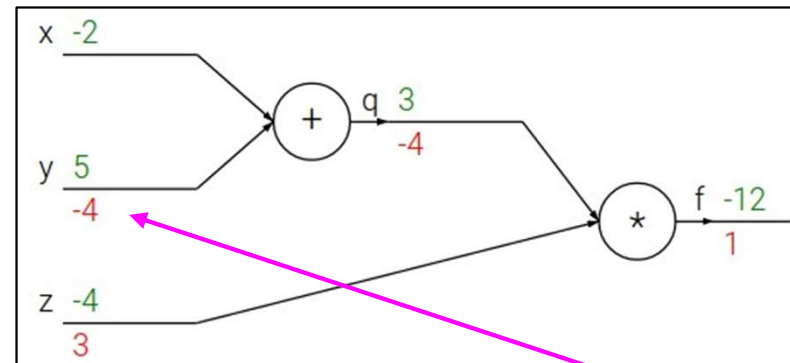
$$f(x, y, z) = (x + y)z$$

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

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

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

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

# Backpropagation: a simple example

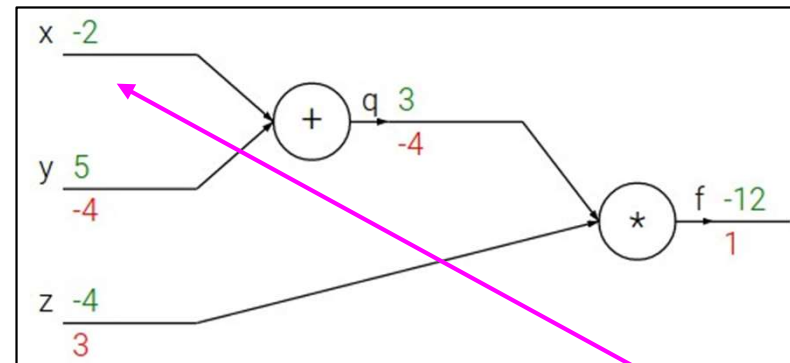
$$f(x, y, z) = (x + y)z$$

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

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

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

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



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

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

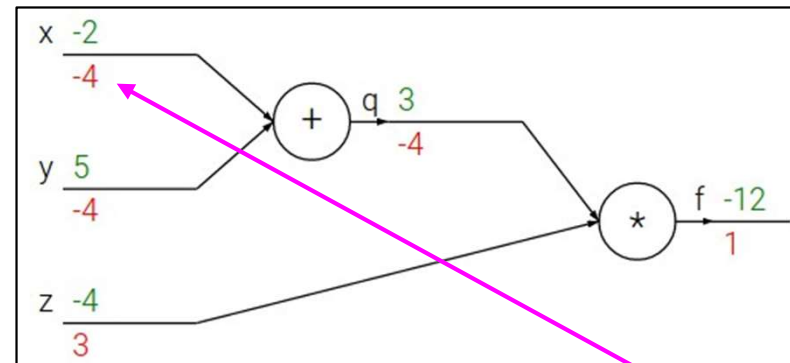
$$f(x, y, z) = (x + y)z$$

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

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

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

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



Chain rule:

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

Upstream  
gradient

Local  
gradient

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