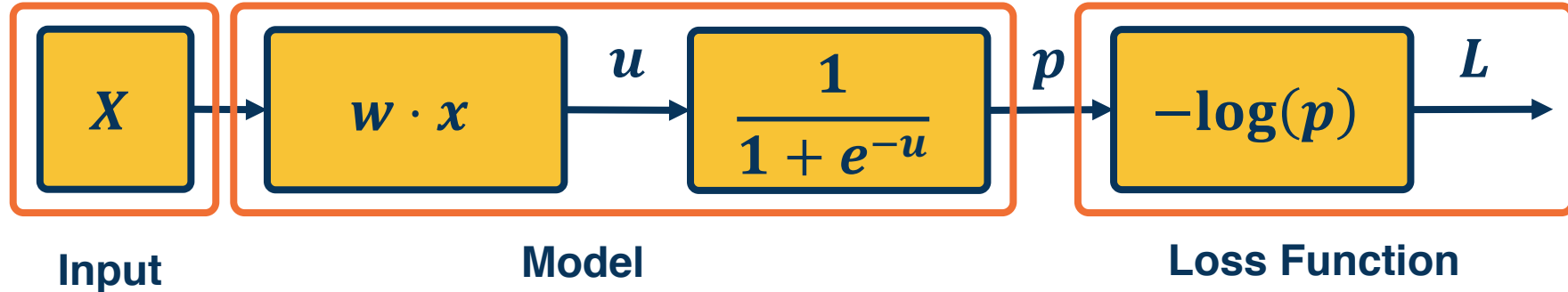


Neural Network View of a Linear Classifier

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

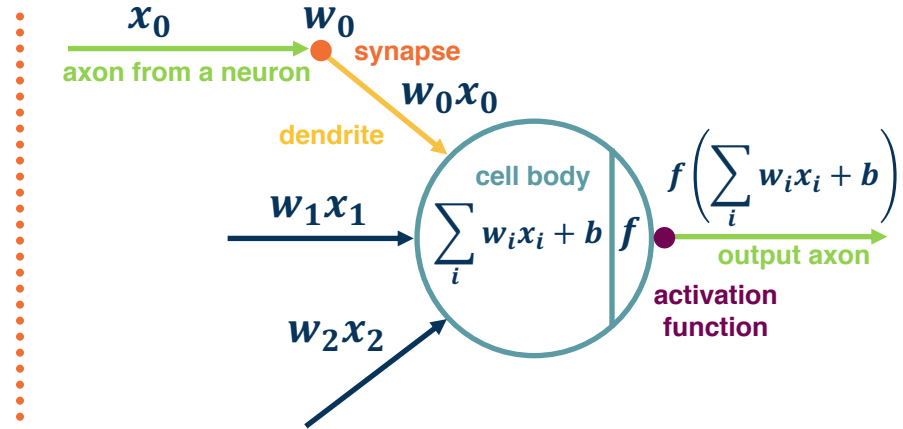
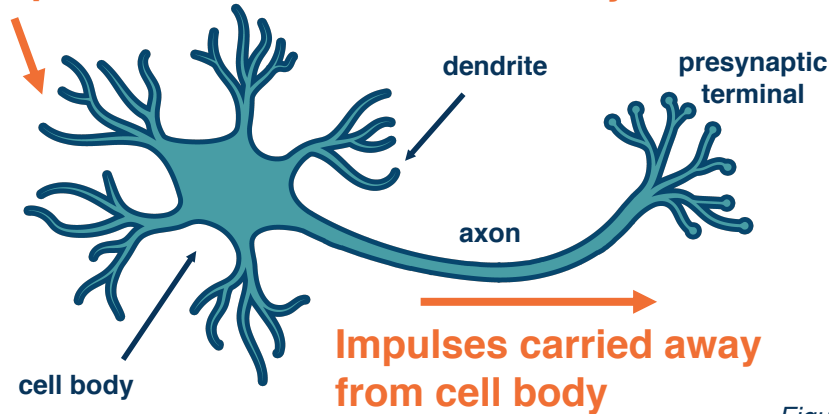


What Does a Linear Classifier Consist of?

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

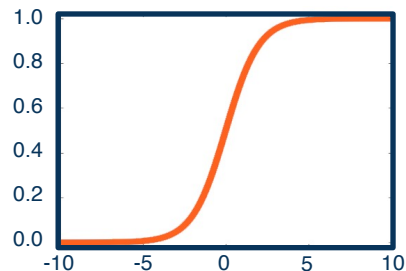
- A neuron takes input (firings) from other neurons (-> **input to linear classifier**)
- The inputs are summed in a weighted manner (-> **weighted sum**)
 - Learning is through a modification of the weights
- If it receives enough input, it “fires” (threshold or if weighted sum plus bias is high enough)

Impulses carried toward cell body



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

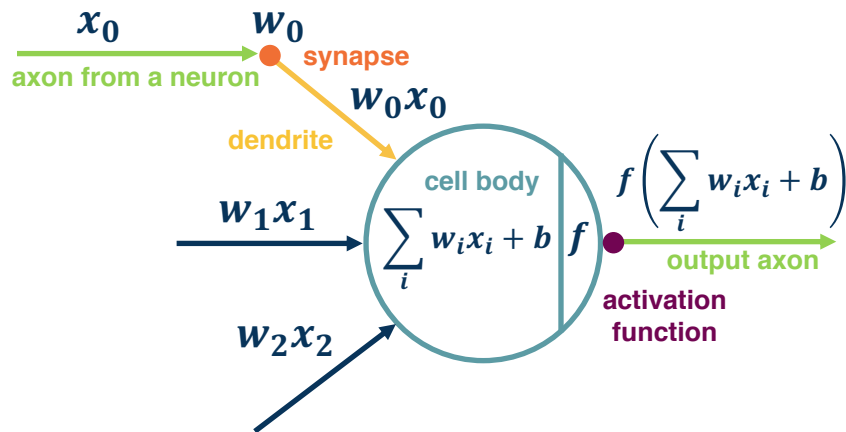
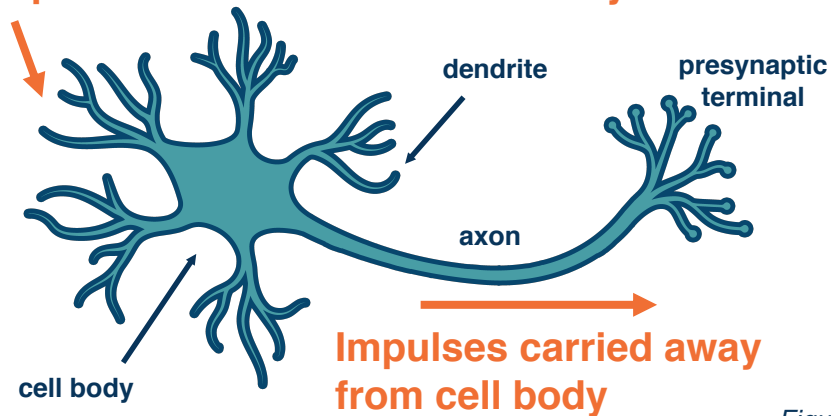
As we did before, the output of a neuron can be modulated by a non-linear function (e.g. sigmoid)



**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

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

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

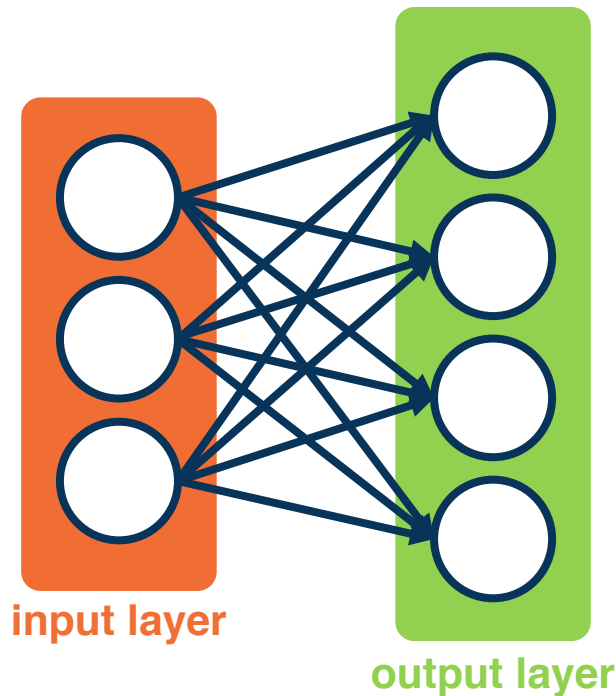


Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

- Each input/output is a **neuron (node)**
- A linear classifier 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**

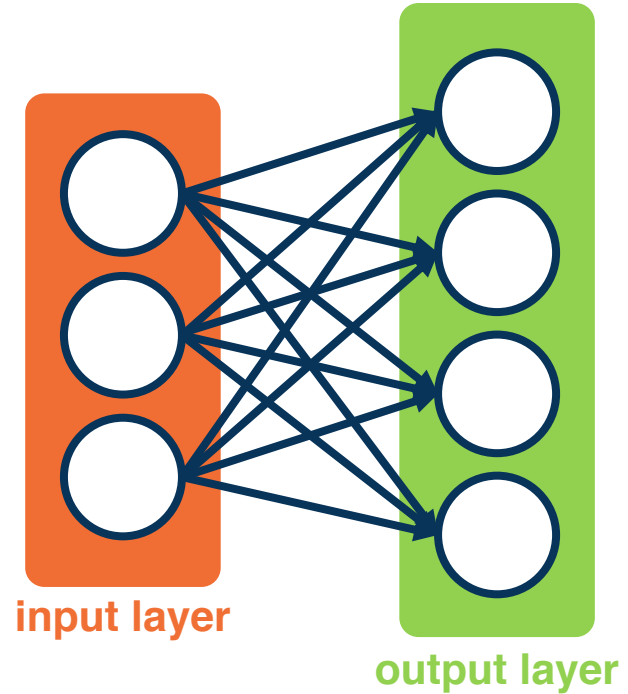


Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

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

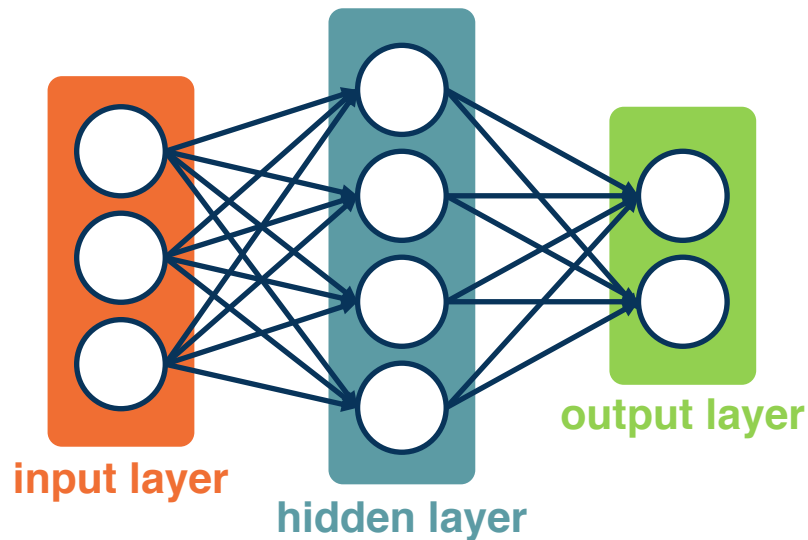


Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

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)

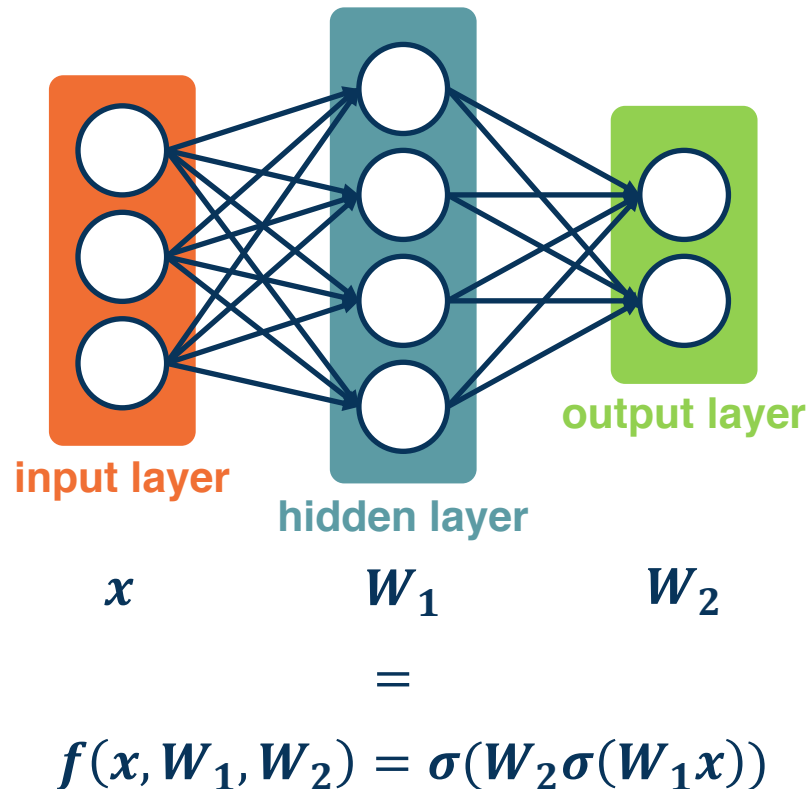


Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

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**:

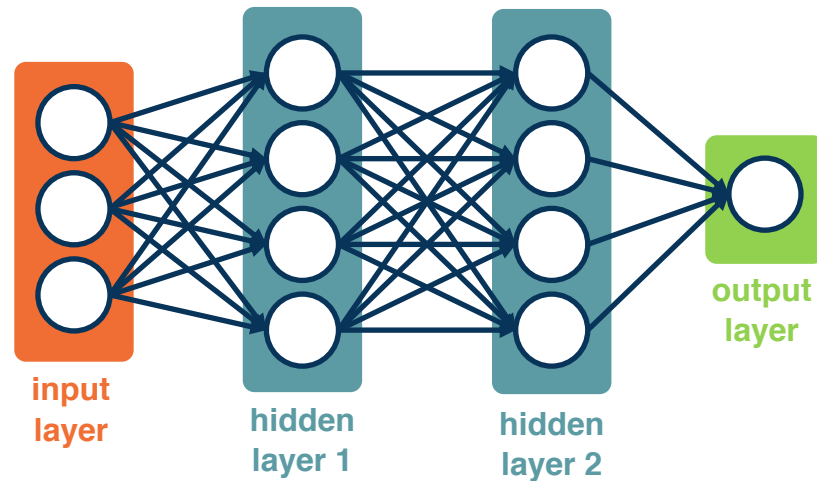
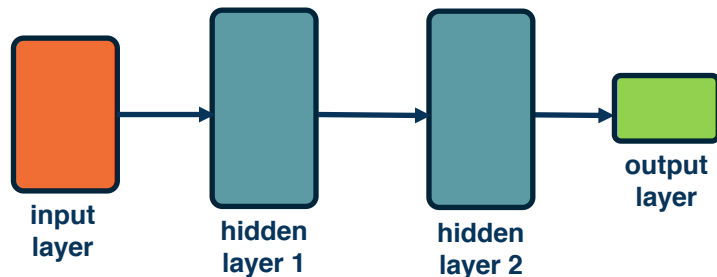


Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Adding More Layers!

Computation Graphs

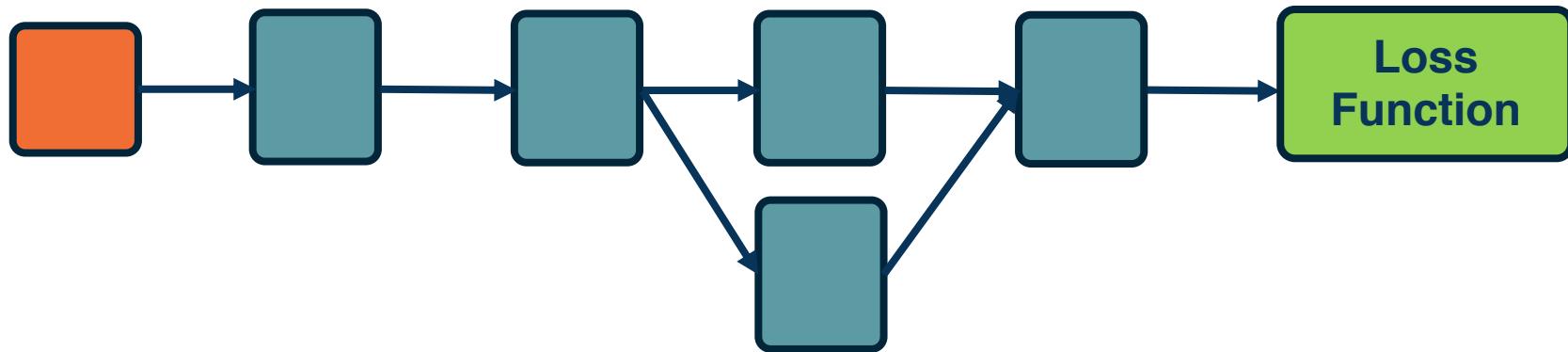
Functions can be made **arbitrarily complex** (subject to memory and computational limits), e.g.:

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

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

VISION

pixels → edge → texton → motif → part → object

SPEECH

sample → spectral band → formant → motif → phone → word

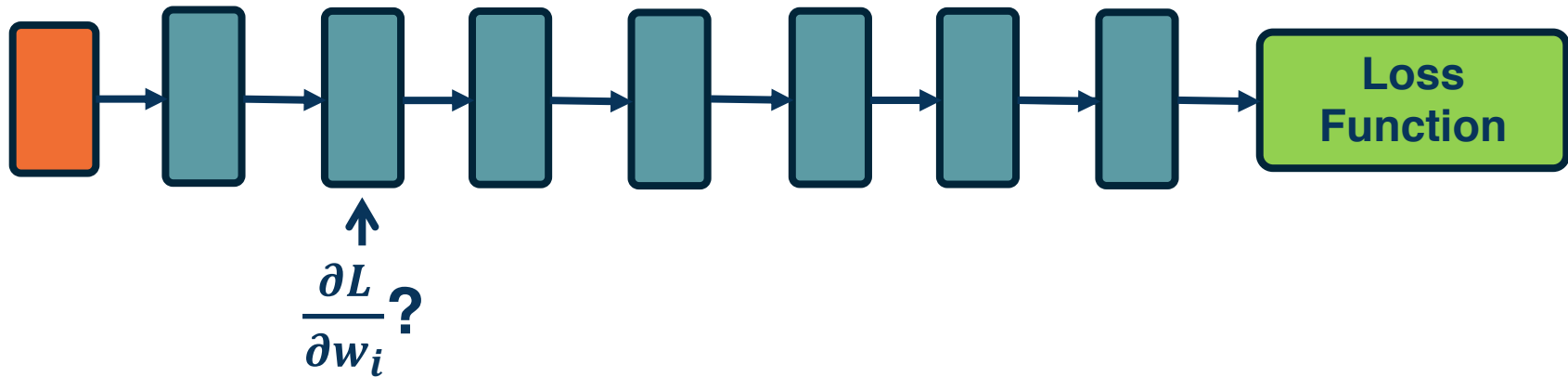
NLP

character → word → NP/VP/.. → clause → sentence → story

Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

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

- 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

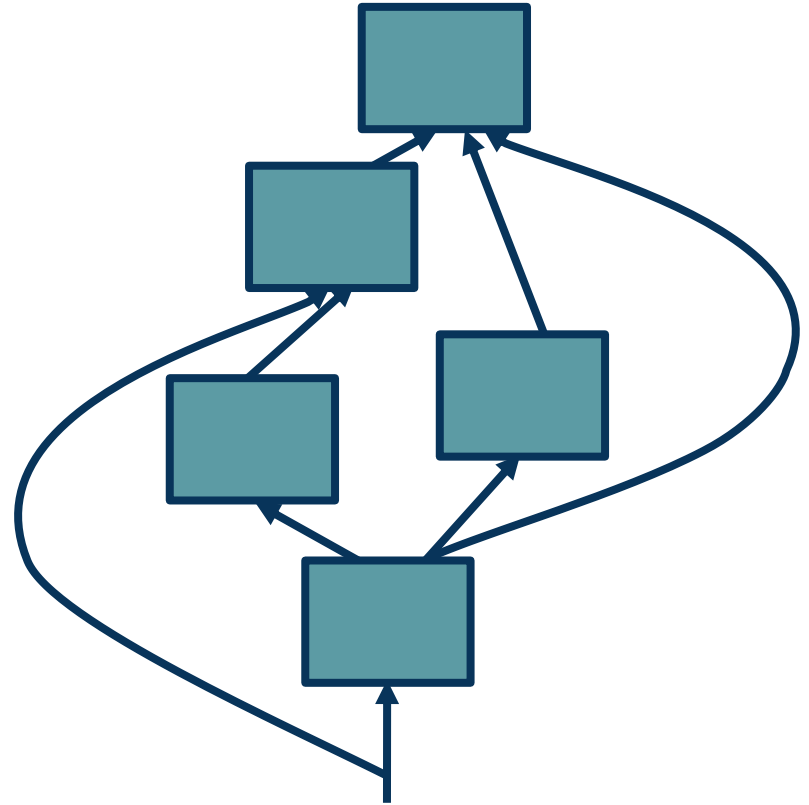


To develop a general algorithm for this, we will view the function as a **computation graph**

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

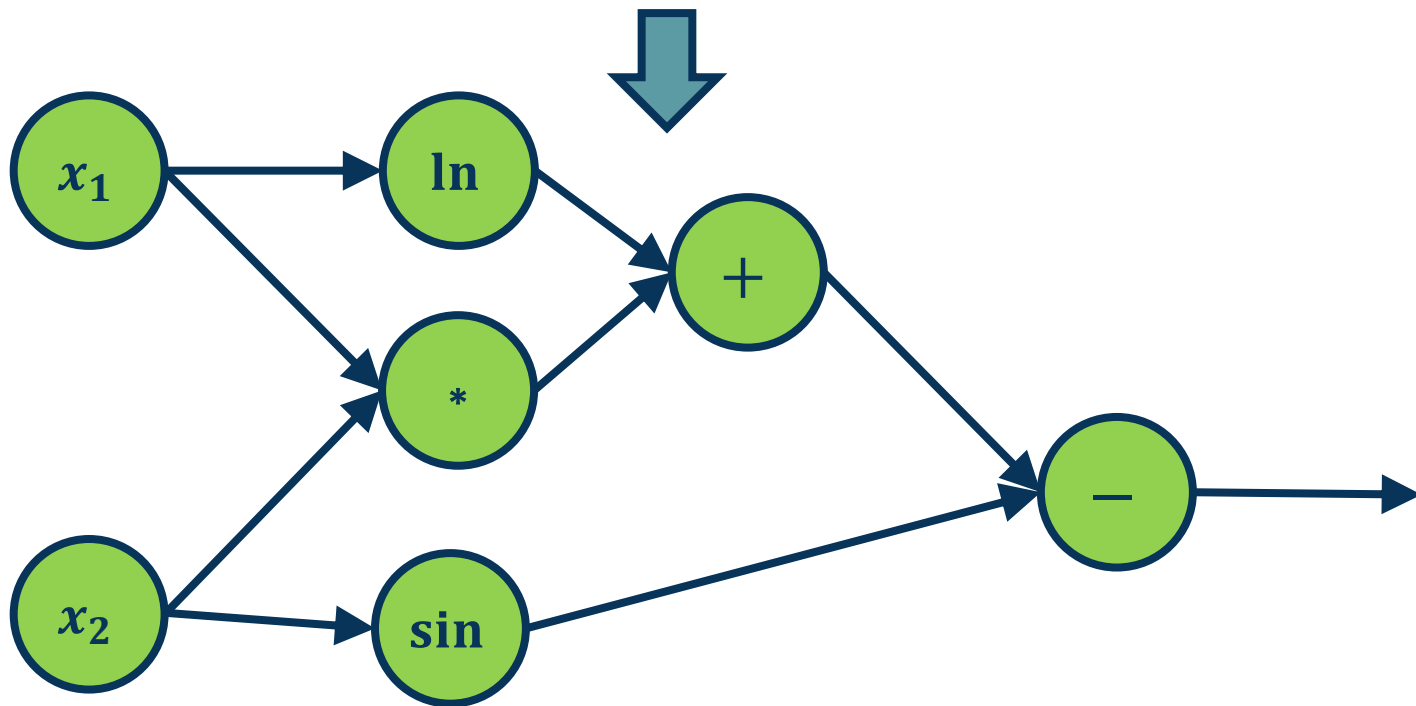
- Modules must be differentiable to support gradient computations for gradient descent

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



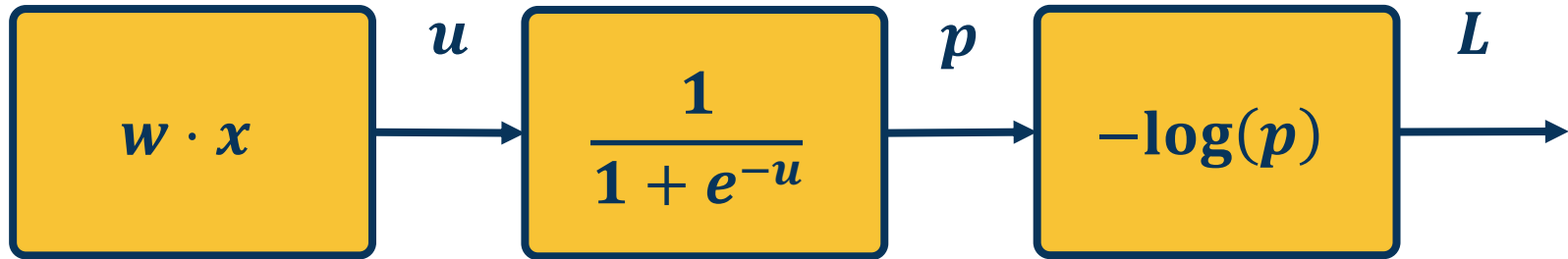
Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

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



Example

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



Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

Backpropagation

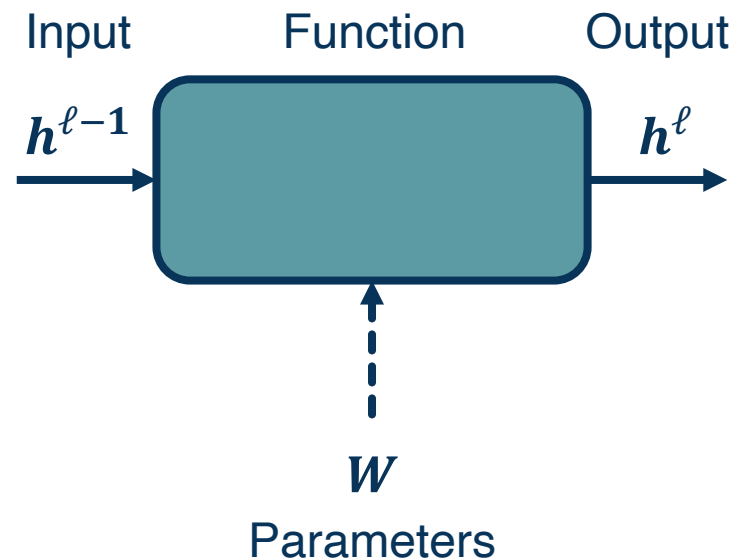
Given this computation graph, the training algorithm will:

- 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**



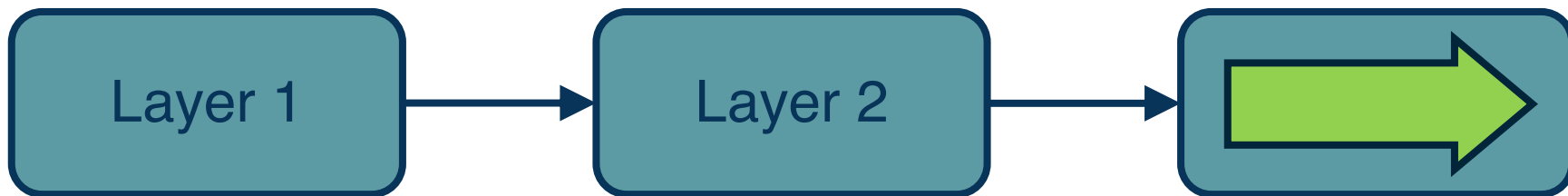
Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

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



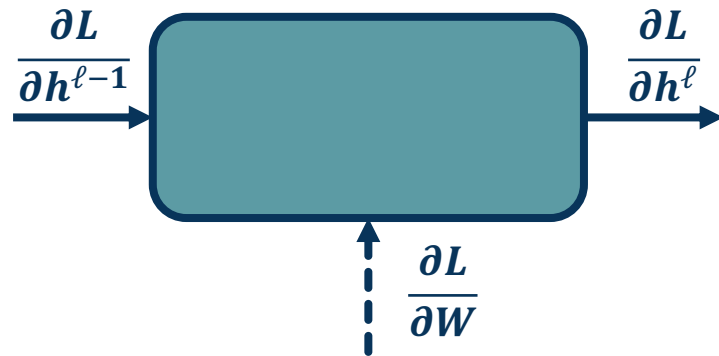
Note that we must store the **intermediate outputs of all layers**!

- 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

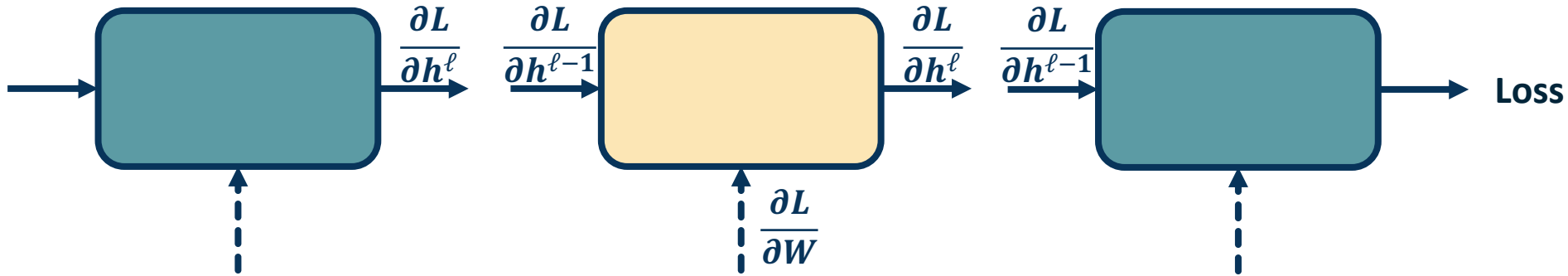
- We can compute **local gradients**: $\{\frac{\partial h^\ell}{\partial h^{\ell-1}}, \frac{\partial h^\ell}{\partial W}\}$
- This is just the **derivative of our function** with respect to its parameters and inputs!

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

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

$$\text{and } \frac{\partial h^\ell}{\partial W} = h^{\ell-1,T}$$

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



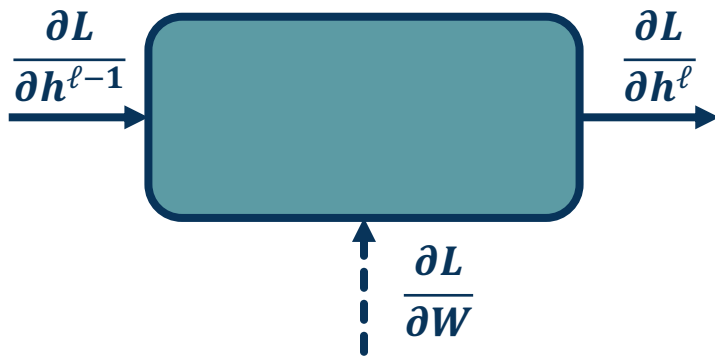
- We will use the *chain rule* to do this:

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**)

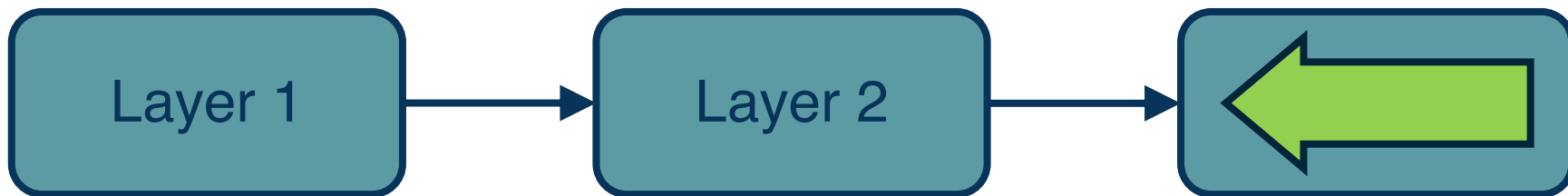
- **Gradient of loss w.r.t. weights:** $\frac{\partial L}{\partial W} = \frac{\partial L}{\partial h^{\ell}} \frac{\partial h^{\ell}}{\partial W}$



Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

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

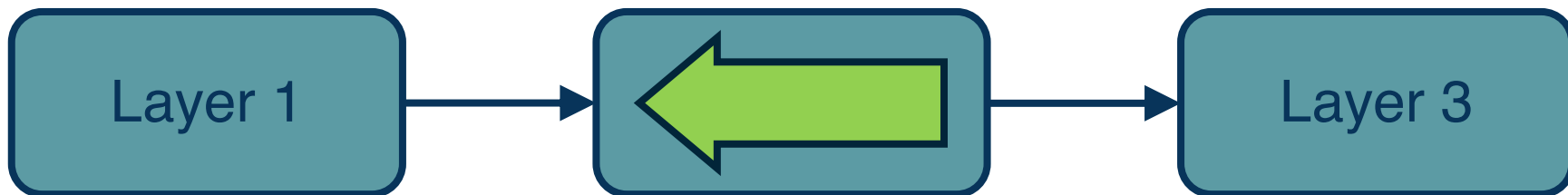
Step 2: Compute Gradients wrt parameters: **Backward Pass**



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**

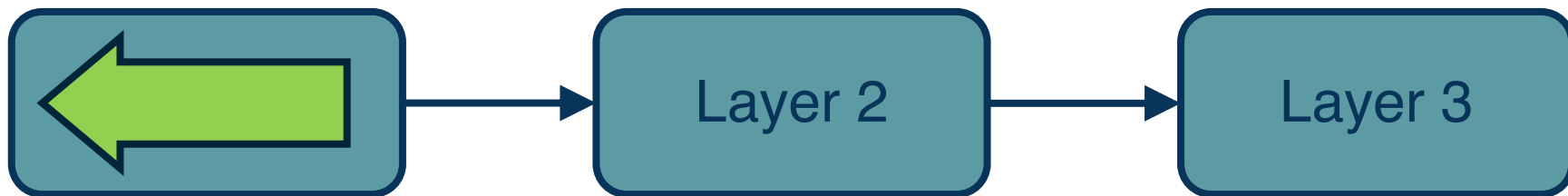


Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun

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 and Automatic Differentiation

Backpropagation does not really spell out how to **efficiently** carry out the necessary computations

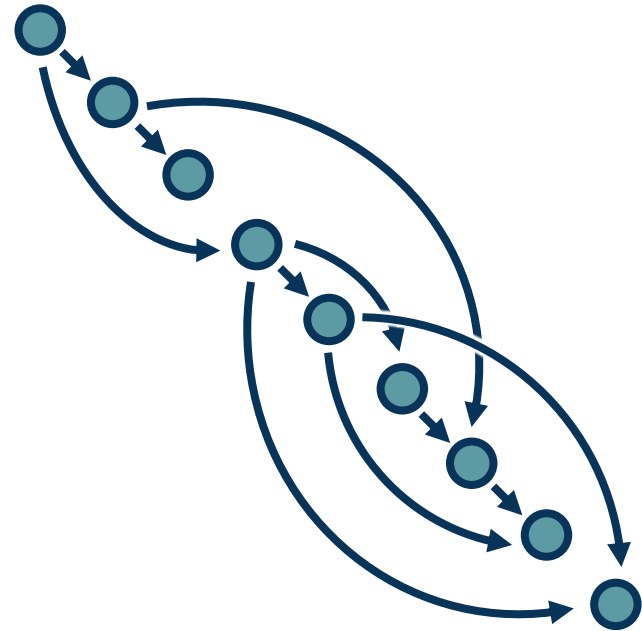
But the idea can be applied to **any directed acyclic graph (DAG)**

- Graph represents an **ordering constraining** which paths must be calculated first

Given an ordering, we can then iterate from the last module backwards, **applying the chain rule**

- We will store, for each node, its **gradient outputs for efficient computation**

This is called reverse-mode **automatic differentiation**



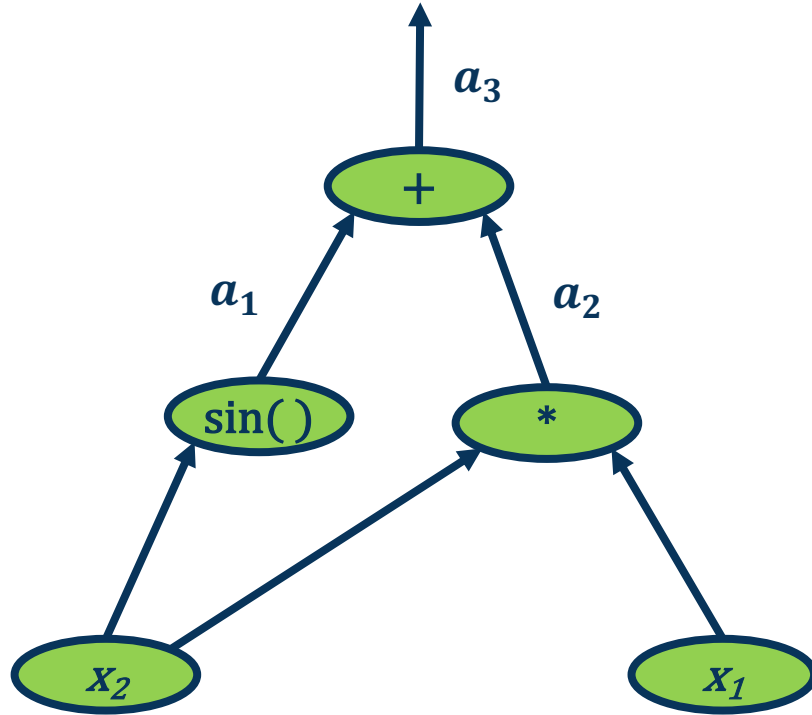
Computation = Graph

- Input = Data + Parameters
- Output = Loss
- Scheduling = Topological ordering

Auto-Diff

- A family of algorithms for implementing chain-rule on computation graphs

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



We want to find the **partial derivative of output f** (output) with respect to **all intermediate variables**

- Assign intermediate variables

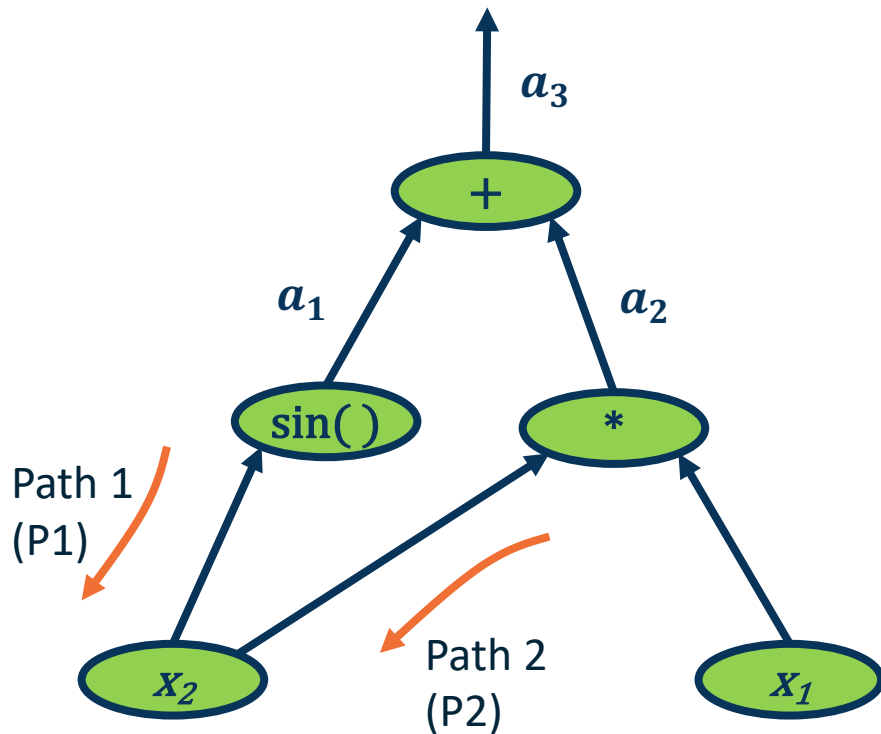
Simplify notation:

Denote bar as: $\bar{a}_3 = \frac{\partial f}{\partial a_3}$

- Start at **end** and move **backward**

Example

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



$$\overline{a_3} = \frac{\partial f}{\partial a_3} = 1$$

$$\overline{a_1} = \frac{\partial f}{\partial a_1} = \frac{\partial f}{\partial a_3} \frac{\partial a_3}{\partial a_1} = \frac{\partial f}{\partial a_3} \frac{\partial(a_1 + a_2)}{\partial a_1} = \frac{\partial f}{\partial a_3} 1 = \overline{a_3}$$

$$\overline{a_2} = \frac{\partial f}{\partial a_2} = \frac{\partial f}{\partial a_3} \frac{\partial a_3}{\partial a_2} = \overline{a_3}$$

$$\overline{x_2^{P1}} = \frac{\partial f}{\partial a_1} \frac{\partial a_1}{\partial x_2} = \overline{a_1} \cos(x_2) +$$

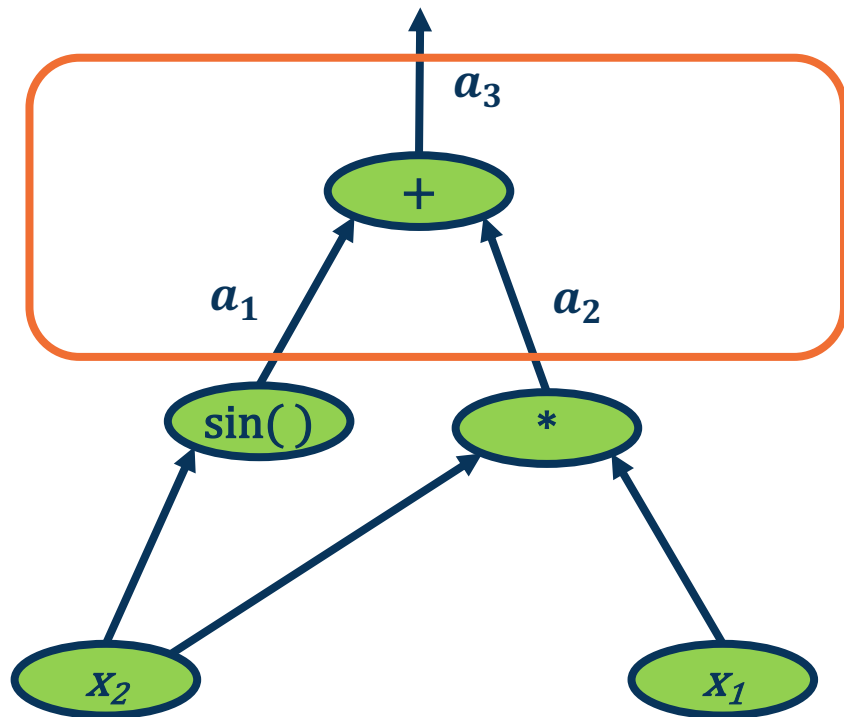
$$\overline{x_2^{P2}} = \frac{\partial f}{\partial a_2} \frac{\partial a_2}{\partial x_2} = \frac{\partial f}{\partial a_2} \frac{\partial(x_1x_2)}{\partial x_2} = \overline{a_2}x_1$$

$$\overline{x_1} = \frac{\partial f}{\partial a_2} \frac{\partial a_2}{\partial x_1} = \overline{a_2}x_2$$

Gradients
from multiple
paths
summed

Example

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

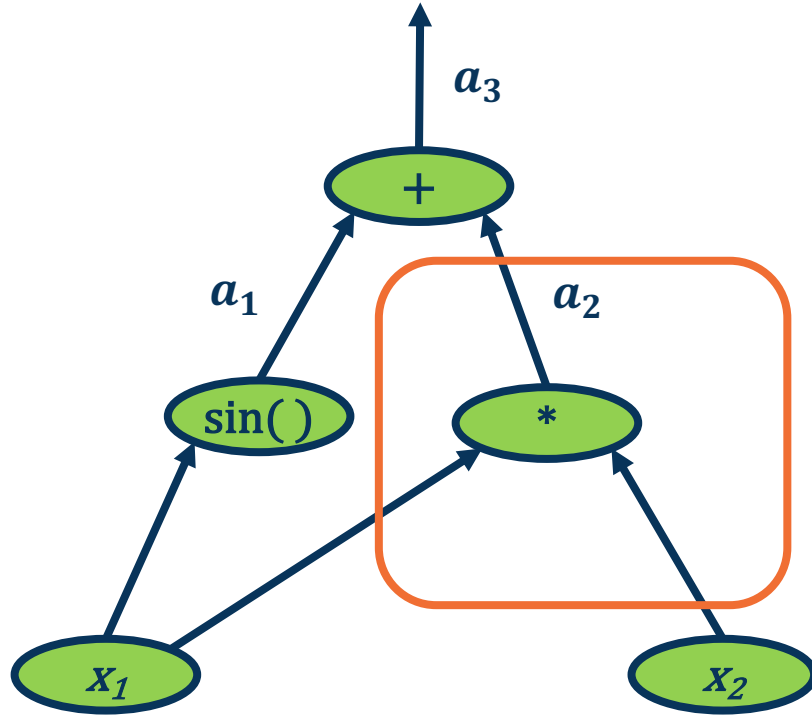


$$\overline{a_1} = \frac{\partial f}{\partial a_1} = \frac{\partial f}{\partial a_3} \frac{\partial a_3}{\partial a_1} = \frac{\partial f}{\partial a_3} \frac{\partial(a_1 + a_2)}{\partial a_1} = \frac{\partial f}{\partial a_3} 1 = \overline{a_3}$$

$$\overline{a_2} = \frac{\partial f}{\partial a_2} = \frac{\partial f}{\partial a_3} \frac{\partial a_3}{\partial a_2} = \overline{a_3}$$

Addition operation distributes gradients along all paths!

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



Multiplication operation is a gradient switcher (multiplies it by the values of the other term)

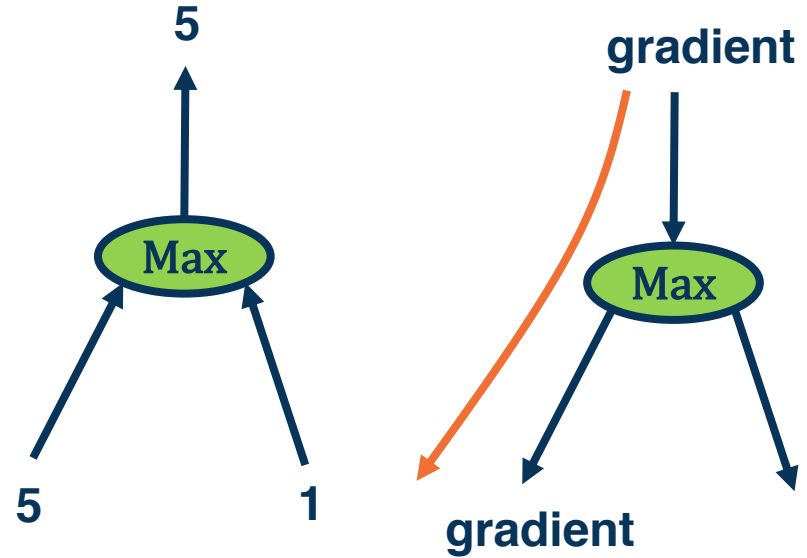
$$\overline{x_2} = \frac{\partial f}{\partial a_2} \frac{\partial a_2}{\partial x_2} = \frac{\partial f}{\partial a_2} \frac{\partial (x_1x_2)}{\partial x_2} = \overline{a_2}x_1$$

$$\overline{x_1} = \frac{\partial f}{\partial a_2} \frac{\partial a_2}{\partial x_1} = \overline{a_2}x_2$$

Several other patterns as well, e.g.:

Max operation **selects** which path to push the gradients through

- Gradient flows along the path that was “selected” to be max
- This information must be recorded in the forward pass

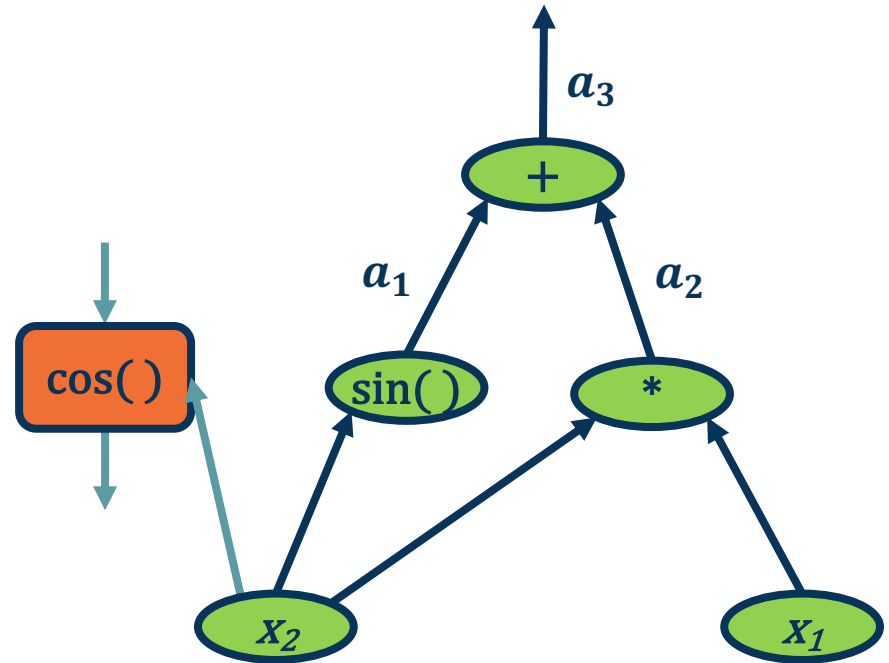


The flow of gradients is one of the **most important aspects** in deep neural networks

- If gradients **do not flow backwards properly**, learning slows or stops!

- Key idea is to **explicitly store computation graph** in memory and **corresponding gradient functions**
- Nodes** broken down to **basic primitive computations** (addition, multiplication, log, etc.) for which **corresponding derivative is known**

$$\overline{x_2} = \frac{\partial f}{\partial a_1} \frac{\partial a_1}{\partial x_2} = \overline{a_1} \cos(x_2)$$

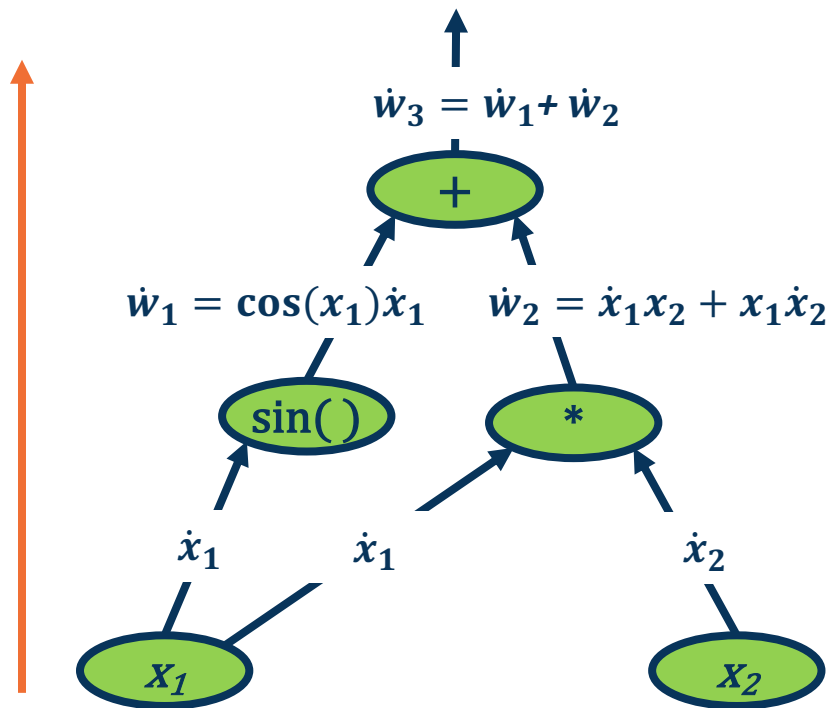


Note that we can also do **forward mode** automatic differentiation

Start from **inputs** and propagate gradients forward

Complexity is proportional to input size

- However, in most cases our **inputs** (images) are large and **outputs** (loss) are small



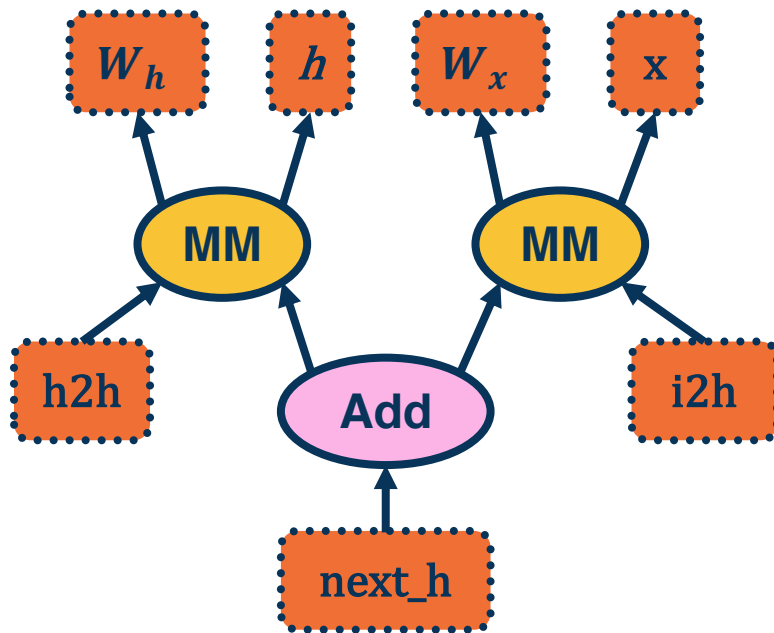
A graph is created on the fly

```
from torch.autograd import Variable
```

```
x = Variable(torch.randn(1, 20))  
prev_h = Variable(torch.randn(1, 20))  
W_h = Variable(torch.randn(20, 20))  
W_x = Variable(torch.randn(20, 20))
```

```
i2h = torch.mm(W_x, x.t())  
h2h = torch.mm(W_h, prev_h.t())  
next_h = i2h + h2h
```

(Note above)



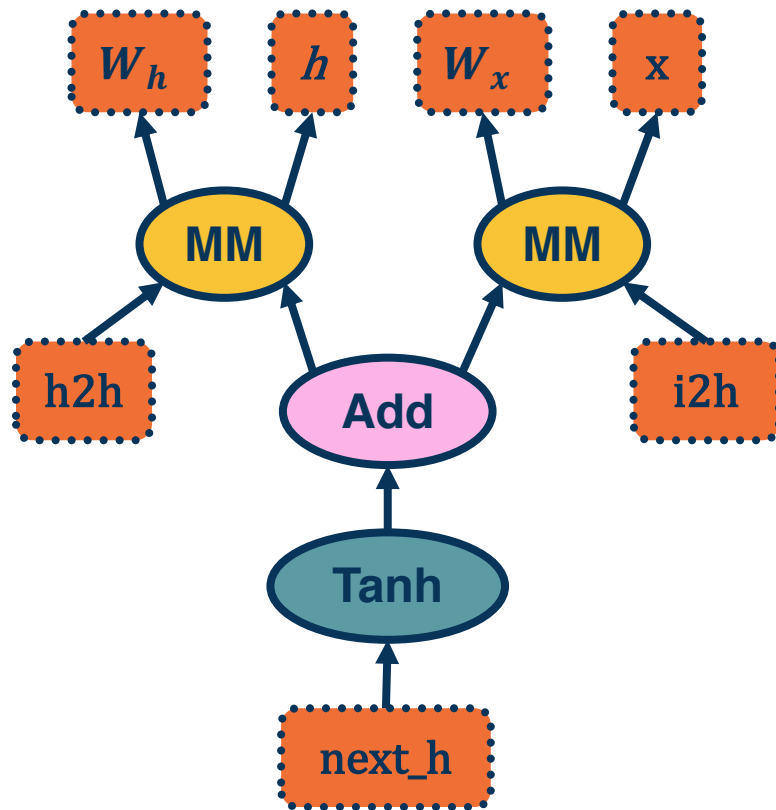
Back-propagation uses the dynamically built graph

```
from torch.autograd import Variable
```

```
x = Variable(torch.randn(1, 20))  
prev_h = Variable(torch.randn(1, 20))  
W_h = Variable(torch.randn(20, 20))  
W_x = Variable(torch.randn(20, 20))
```

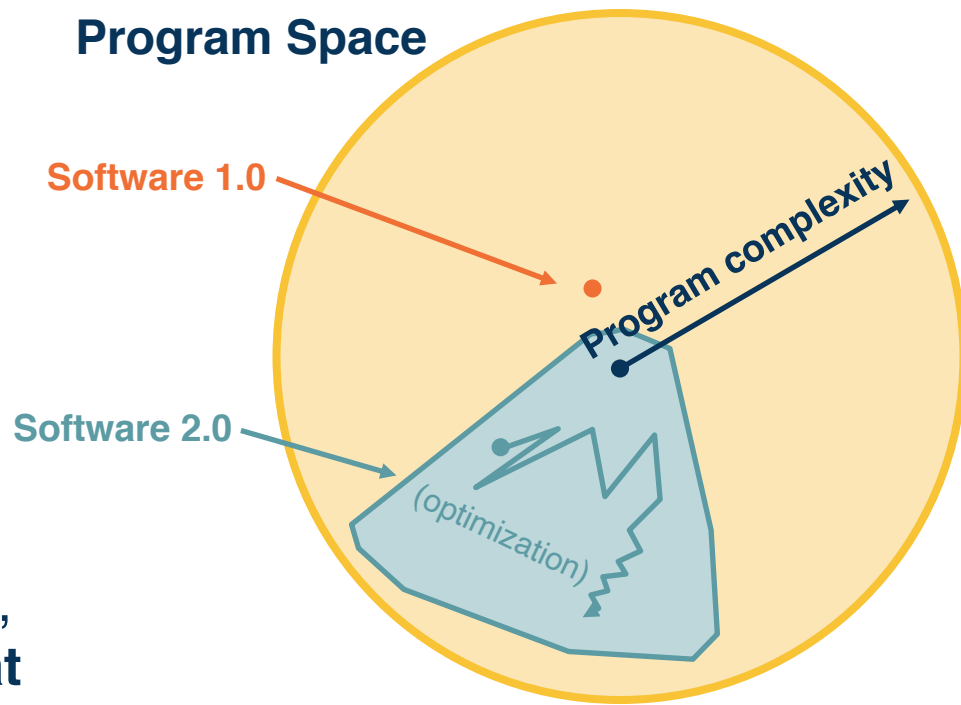
```
i2h = torch.mm(W_x, x.t())  
h2h = torch.mm(W_h, prev_h.t())  
next_h = i2h + h2h  
next_h = next_h.tanh()
```

```
next_h.backward(torch.ones(1, 20))
```



From pytorch.org

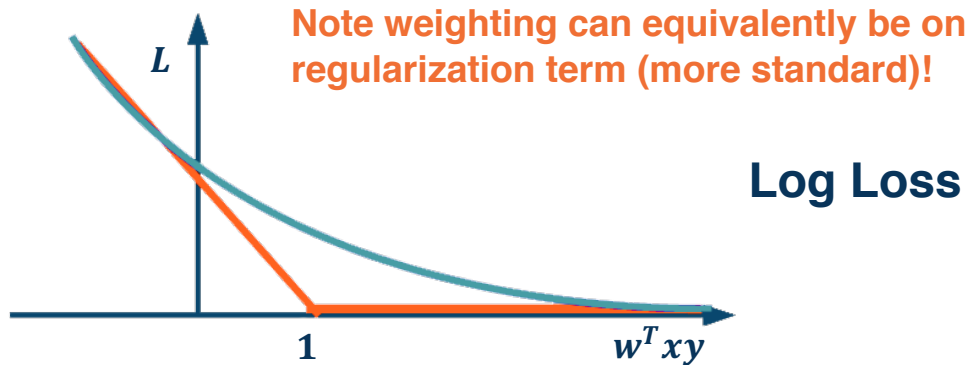
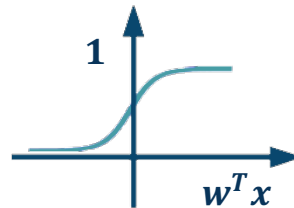
- Computation graphs are **not limited to mathematical functions!**
- Can have **control flows** (if statements, loops) and **backpropagate** through **algorithms!**
- Can be done **dynamically** so that **gradients are computed**, then **nodes are added**, repeat
- Differentiable programming**



Adapted from figure by Andrej Karpathy

Computation Graph Example for Logistic Regression

- Input: $x \in R^D$
- Binary label: $y \in \{-1, +1\}$
- Parameters: $w \in R^D$
- Output prediction: $p(y = 1|x) = \frac{1}{1+e^{-w^T x}}$
- Loss: $L = \frac{1}{2} \|w\|^2 - \lambda \log(p(y|x))$

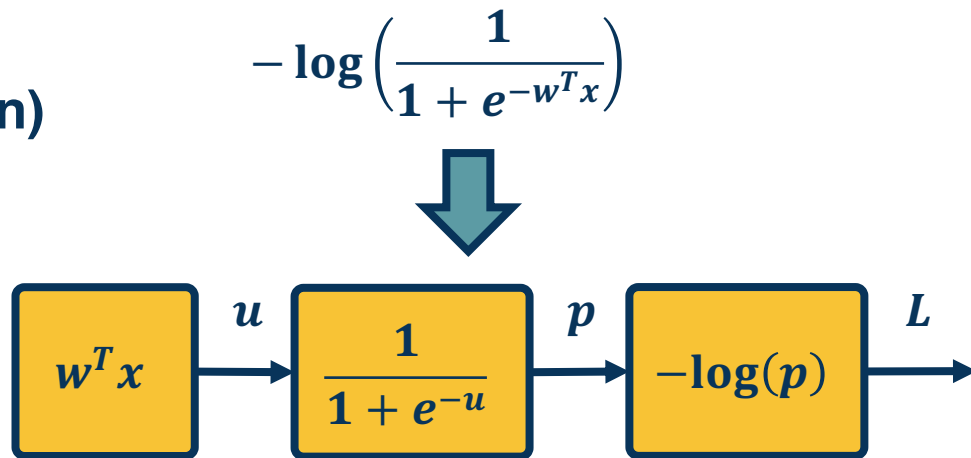


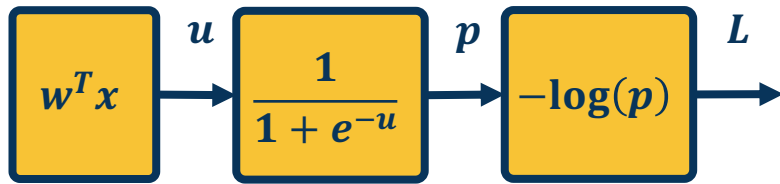
Adapted from slide by Marc'Aurelio Ranzato

We have discussed **computation graphs for generic functions**

Machine Learning functions
(input -> model -> loss function)
is also a computation graph

We can use the **computed gradients from backprop/automatic differentiation** to update the weights!





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!**

$$\bar{L} = 1$$

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

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

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

$$\bar{w} = \frac{\partial L}{\partial w} = \frac{\partial L}{\partial u} \frac{\partial u}{\partial w} = \bar{u} x^T$$

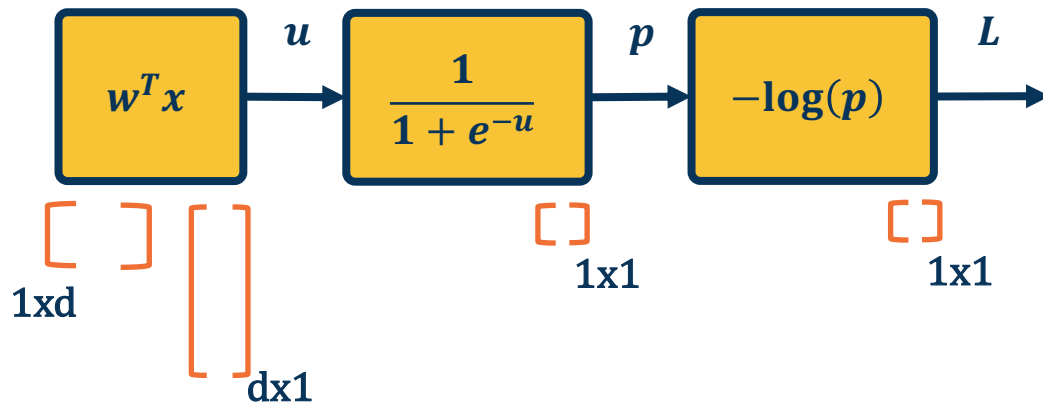
We can do this in a combined way to see all terms together:

$$\begin{aligned} \bar{w} &= \frac{\partial L}{\partial p} \frac{\partial p}{\partial u} \frac{\partial u}{\partial w} = -\frac{1}{\sigma(w^T x)} \sigma(w^T x) (1 - \sigma(w^T x)) x^T \\ &= -\left(1 - \sigma(w^T x)\right) x^T \end{aligned}$$

This effectively shows gradient flow along path from L to w

Vectorization and Jacobians of Simple Layers

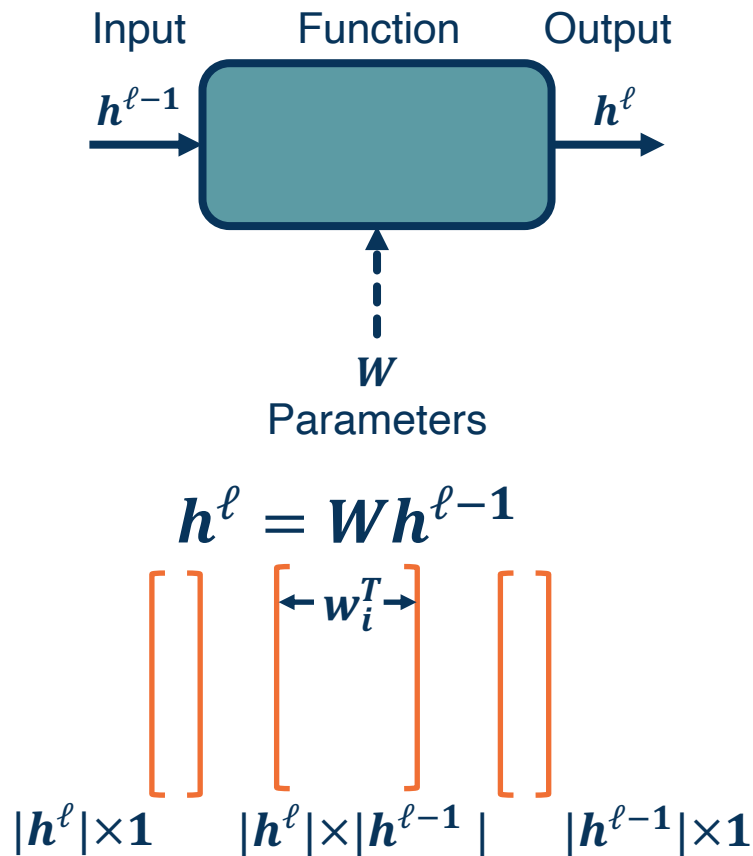
The chain rule can be computed as a **series of scalar, vector, and matrix linear algebra operations**



Extremely efficient in graphics processing units (GPUs)

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

$\begin{matrix} \text{[]} & \text{[]} & \text{[]} & \text{[]} & \text{[]} \\ 1 \times 1 & 1 \times 1 & 1 \times 1 & 1 \times 1 & 1 \times d \end{matrix}$

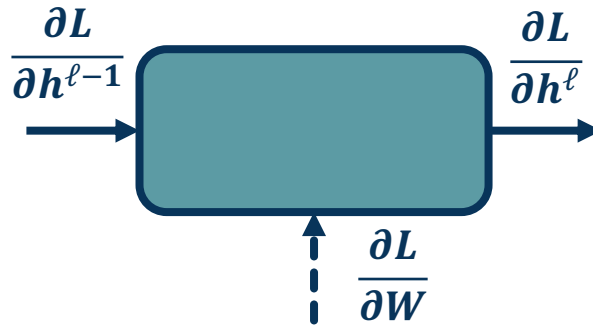


Fully Connected (FC) Layer: Forward Function

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

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

(other elements zeros)



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

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

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

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

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

$$1 \times |h^\ell|$$

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

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

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

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

$$1 \times |h^\ell|$$

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

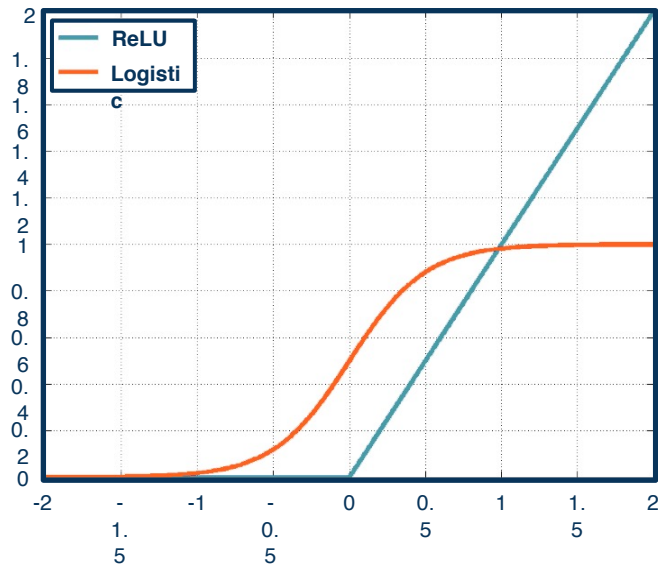
Fully Connected (FC) Layer

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

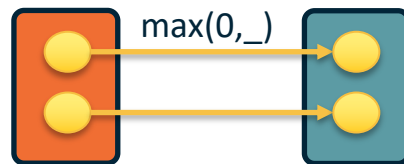
A common choice is the **Rectified Linear Unit**

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

How many parameters for this layer?



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



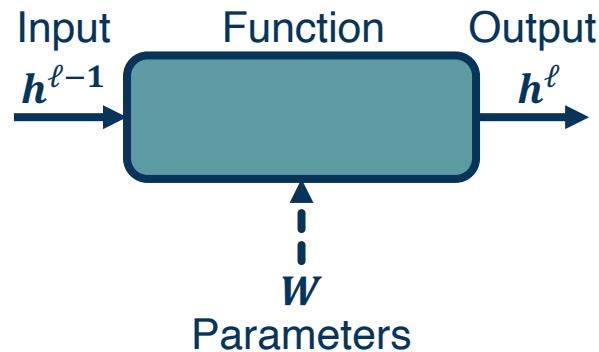
Rectified Linear Unit (ReLU)

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



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

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

Diagram illustrating the Jacobian matrix structure for ReLU. The matrix is $|h^{\ell} \times h^{\ell-1}|$ in size, showing a sparse structure with a single non-zero diagonal element.

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