

Topics:

- Jacobians/Matrix Calculus continued
- Backpropagation / Automatic Differentiation

**CS 4644 / 7643-A**

**ZSOLT KIRA**

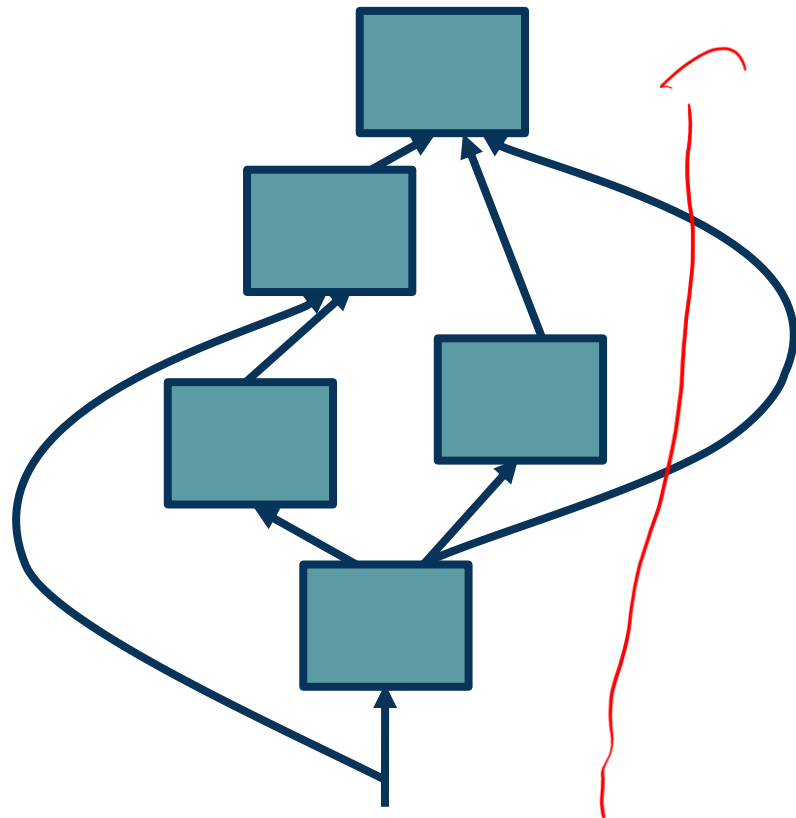
- **Assignment 1 out!**
  - **Due Jan 31<sup>st</sup> (with grace period Feb 2<sup>nd</sup>)**
  - Hopefully you have started!
- Resources:
  - These lectures
  - [Matrix calculus for deep learning](#)
  - [Gradients notes](#) and [MLP/ReLU Jacobian notes](#).
  - Assignment 1 and matrix calculus (@93, @109)
- Project:
  - Project proposal overview out, due by **Feb 14<sup>th</sup>**
  - See recorded office hours for discussion about requirements, answers to some questions, etc.
  - Groups of 5 requires permission: Start private piazza thread

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

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

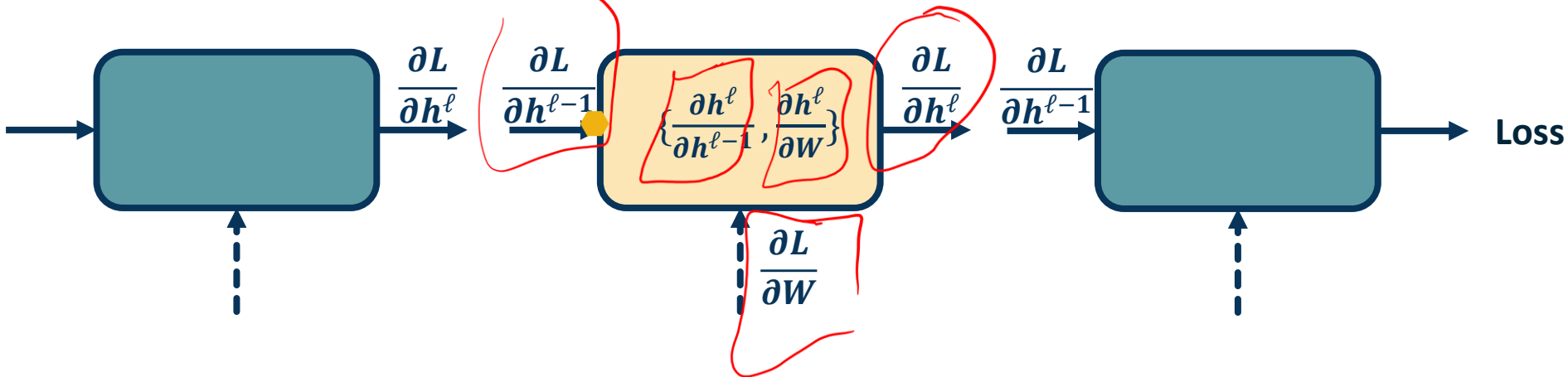
- ◆ Modules must be differentiable to support gradient computations for gradient descent

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



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

- 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:** 
$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial x}$$

## 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  $M \in \mathbb{R}^{k \times \ell}$

	$S$ [ ]	$V$ [ ]	$M$ [ ]
$S$	$\frac{\partial s_1}{\partial s_2}$ [ ]	$\frac{\partial s}{\partial v}$ [ ]	$\frac{\partial s}{\partial M}$ [ ]
$V$	$\frac{\partial v}{\partial s}$ [ ]	$\frac{\partial v_1}{\partial v_2}$ [ ]	
$M$	$\frac{\partial M}{\partial s}$ [ ]		

**Tensors**

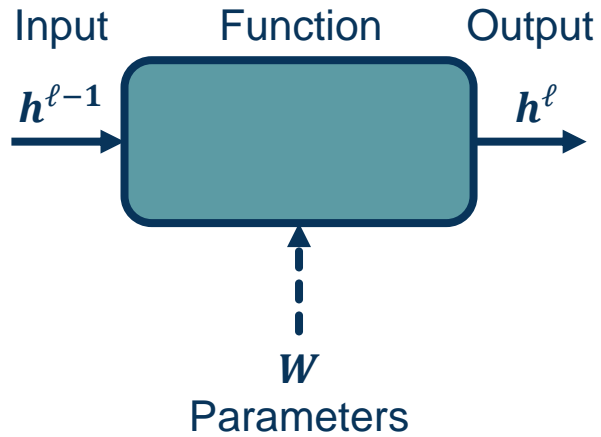
What is the size of  $\frac{\partial L}{\partial W}$ ?

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

$$\begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1m} & b_1 \\ w_{21} & w_{22} & \cdots & w_{2m} & b_2 \\ w_{31} & w_{32} & \cdots & w_{3m} & b_3 \end{bmatrix}$$

Jacobian is also a matrix:

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



Define:

$$h_i^l = w_i^T h^{l-1}$$

$$h^l = W h^{l-1}$$

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

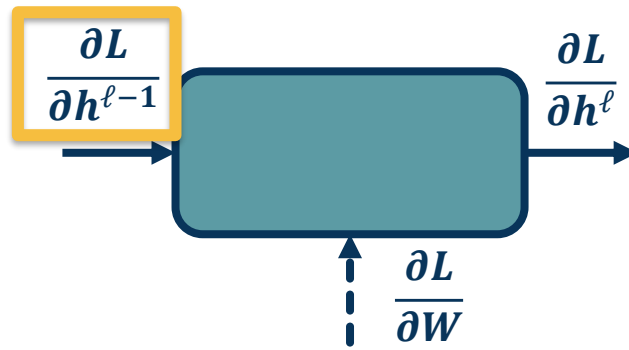
## Fully Connected (FC) Layer: Forward Function

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

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

Define:

$$\mathbf{h}_i^\ell = \mathbf{w}_i^T \mathbf{h}^{\ell-1}$$



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

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

Fully Connected (FC) Layer



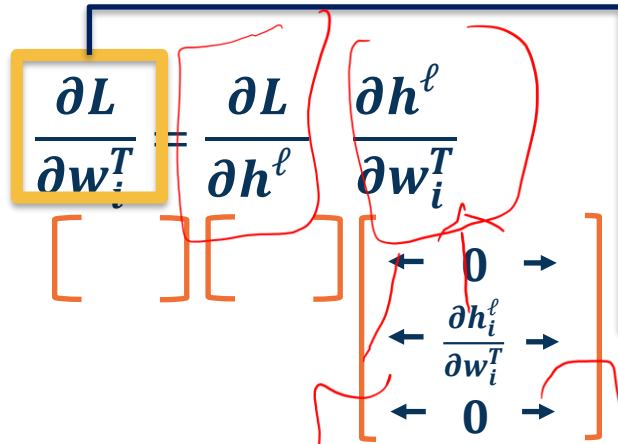
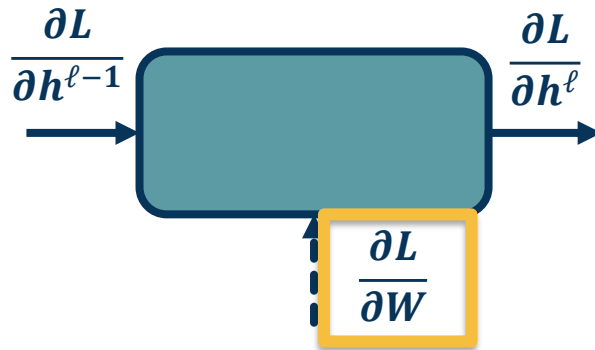
$$\mathbf{h}^\ell = \mathbf{W}\mathbf{h}^{\ell-1}$$

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

Define:

$$\mathbf{h}_i^\ell = \mathbf{w}_i^T \mathbf{h}^{\ell-1}$$

$$\frac{\partial \mathbf{h}_i^\ell}{\partial \mathbf{w}_i^T} = \mathbf{h}^{(\ell-1),T}$$



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

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

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

$$\frac{\partial \mathbf{h}^\ell}{\partial \mathbf{w}} = \begin{bmatrix} \frac{\partial h_1^\ell}{\partial \mathbf{w}} \\ \frac{\partial h_2^\ell}{\partial \mathbf{w}} \\ \vdots \end{bmatrix}$$

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

Iterate and populate  
Note can simplify/vectorize!

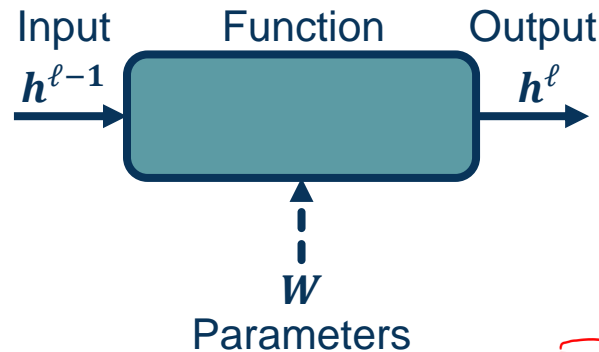
## Fully Connected (FC) Layer

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**

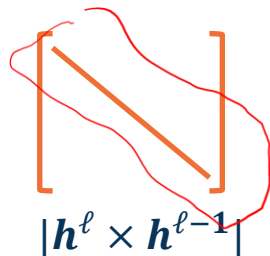
Max function funnels gradients through selected max

- Gradient will be **zero** if input  $\leq 0$



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

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



For diagonal

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

4D input x:

$$\begin{bmatrix} 1 \\ -2 \\ 3 \\ -1 \end{bmatrix}$$

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

(elementwise)

4D output z:

$$\begin{bmatrix} 1 \\ 0 \\ 3 \\ 0 \end{bmatrix}$$

4D dL/dx:

$$\begin{bmatrix} 4 \\ 0 \\ 5 \\ 0 \end{bmatrix}$$


[dz/dx] [dL/dz]

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$$

4D dL/dz:

$$\begin{bmatrix} 4 \\ -1 \\ 5 \\ 9 \end{bmatrix}$$
Upstream  
gradient

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

# Vectorization in Function Compositions

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

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

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

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

(Many of these will be parameterized)

$\times \mathbb{R} \rightarrow \mathbb{O} \rightarrow \mathbb{O} \rightarrow \dots$

(Note you might find the opposite notation as well!)

$$x \in \mathbb{R}^1 \xrightarrow{g_1} z \in \mathbb{R}^1 \xrightarrow{g_2} y \in \mathbb{R}^1$$

$$y = g_2(g_1(x))$$

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial x}$$

↑  
scalar  
mult

$$\vec{x} \in \mathbb{R}^d \xrightarrow{g_1} \vec{z} \in \mathbb{R}^m \xrightarrow{g_2} \vec{y} \in \mathbb{R}^c$$

$$\mathbb{R}^d \rightarrow \mathbb{R}^m \qquad \mathbb{R}^m \rightarrow \mathbb{R}^c$$

$$\left[ \frac{\vec{y}}{\vec{x}} \right] = \begin{bmatrix} \vec{y} \\ \vec{z} \end{bmatrix} = \begin{bmatrix} \vec{y} \\ \vec{z} \end{bmatrix} \begin{bmatrix} \vec{z} \\ \vec{x} \end{bmatrix}$$

matrix multiplication



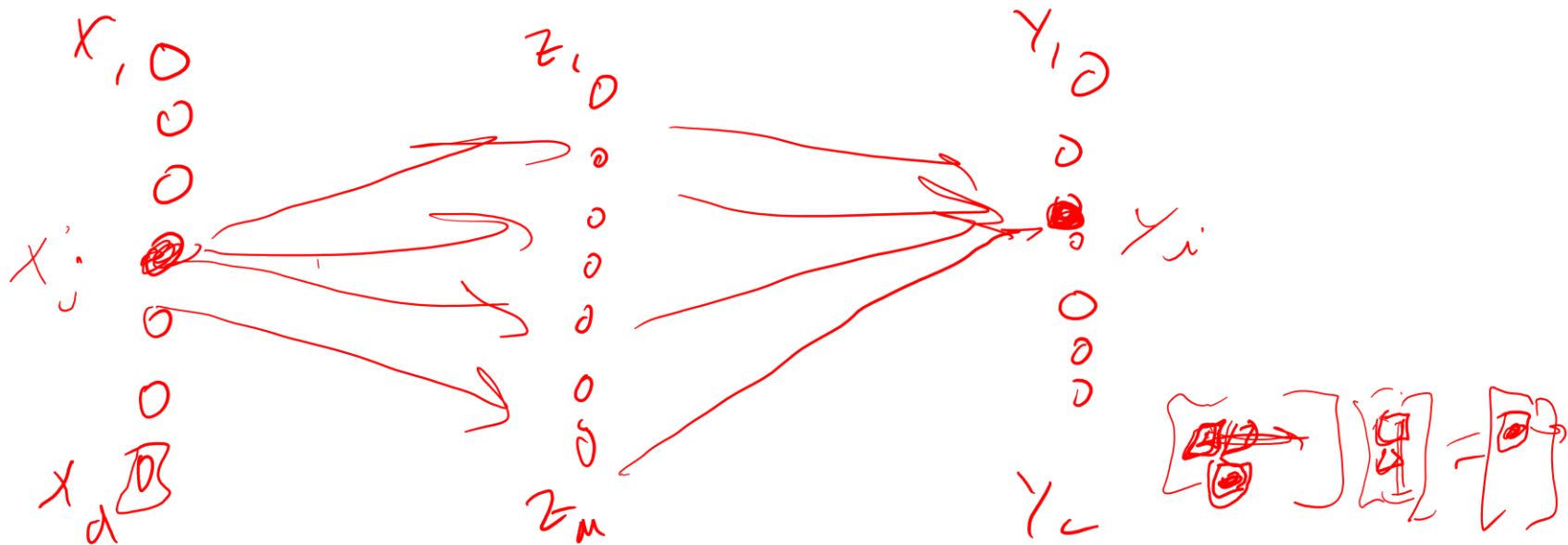
$$\frac{\partial \vec{y}}{\partial \vec{x}} =$$

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

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

$$\left[ \frac{\partial y_i}{\partial x_j} \right] = \left[ \frac{\partial y_i}{\partial z_k} \right] \left[ \frac{\partial z_k}{\partial x_j} \right]$$

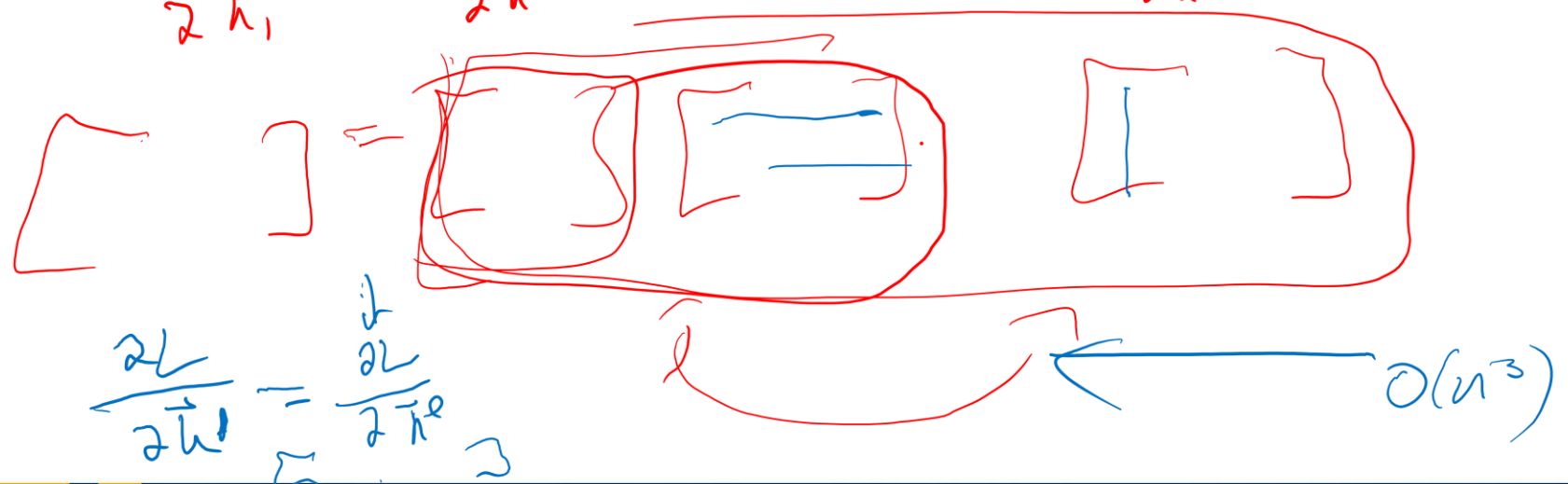
$$\frac{\partial y_i}{\partial x_j} = \sum_k \frac{\partial y_i}{\partial z_k} \cdot \frac{\partial z_k}{\partial x_j}$$



$$\frac{\partial y_i}{\partial x_j} = \sum_{k \text{ path } k_s} \frac{\partial y_i}{\partial z_k} \cdot \frac{\partial z_k}{\partial x_j}$$

$$\vec{h}^1 \in \mathbb{R}^d \rightarrow \vec{h}^2 \rightarrow \dots \rightarrow \vec{h}^l \rightarrow L$$

$$\frac{\partial \vec{h}^l}{\partial \vec{h}^1} = \frac{\partial \vec{h}^l}{\partial \vec{h}^{l-1}} \cdot \frac{\partial \vec{h}^{l-1}}{\partial \vec{h}^{l-2}} \cdot \dots \cdot \frac{\partial \vec{h}^2}{\partial \vec{h}^1}$$

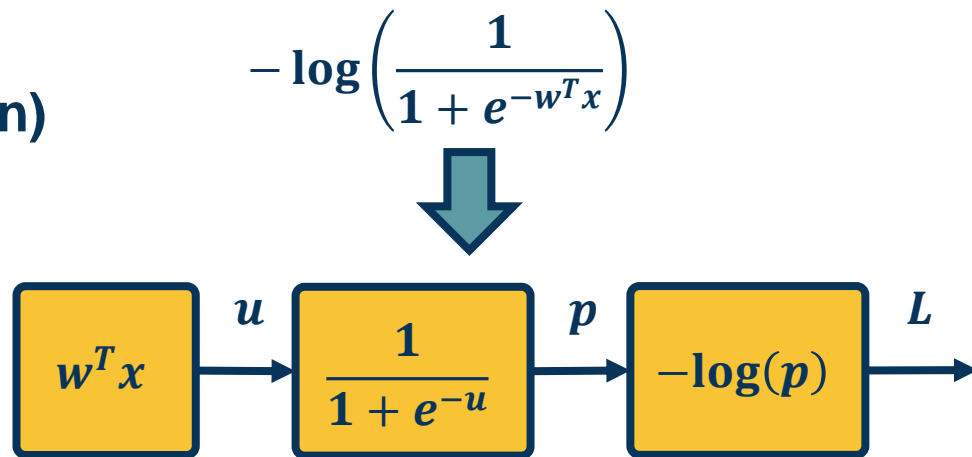


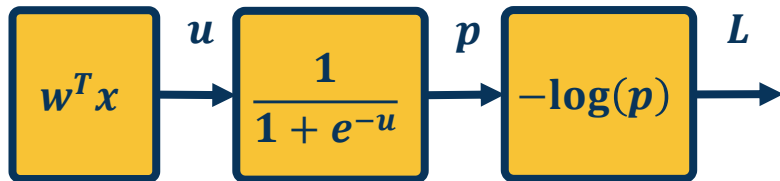
## Chain Rule: Cascaded

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!





$$\bar{L} = 1$$

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

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

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

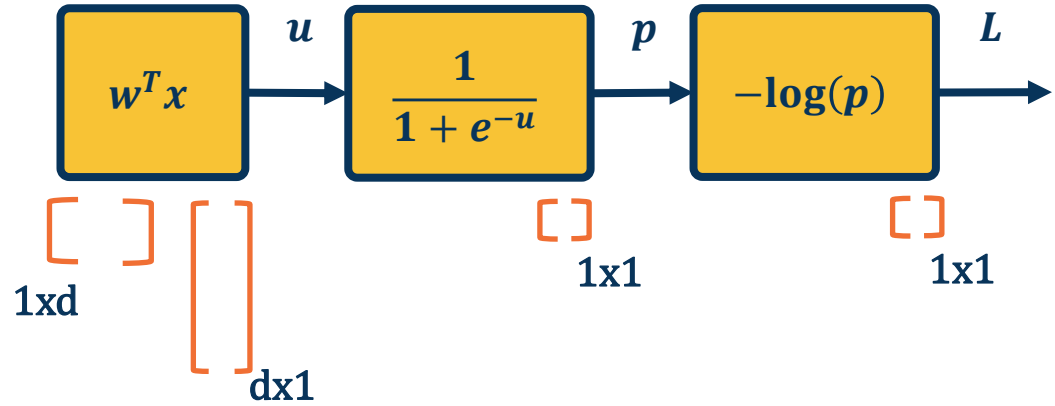
$$\bar{w} = \frac{\partial L}{\partial w} = \frac{\partial L}{\partial u} \frac{\partial u}{\partial w} = \bar{u} x^T$$

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

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

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

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



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

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

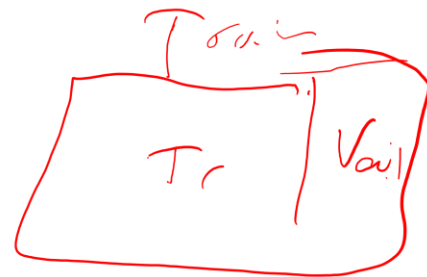
$\begin{bmatrix} \phantom{0} \end{bmatrix}$   
1x1
 $\begin{bmatrix} \phantom{0} \end{bmatrix}$   
1x1
 $\begin{bmatrix} \phantom{0} \end{bmatrix}$   
1x1
 $\begin{bmatrix} \phantom{0} \end{bmatrix}$   
1xd

Many **standard regularization methods** still apply!

## L1 Regularization

$$L = |y - Wx_i|^2 + \lambda|W|$$

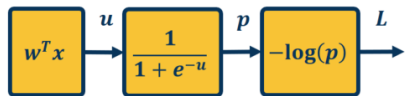
where  $|W|$  is element-wise



**Example regularizations:**

- ◆ L1/L2 on weights (encourage small values)
- ◆ L2:  $L = |y - Wx_i|^2 + \lambda|W|^2$  (weight decay)
- ◆ Elastic L1/L2:  $|y - Wx_i|^2 + \alpha|W|^2 + \beta|W|$





$$L = \frac{1}{1 + e^{-u}}$$

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

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:

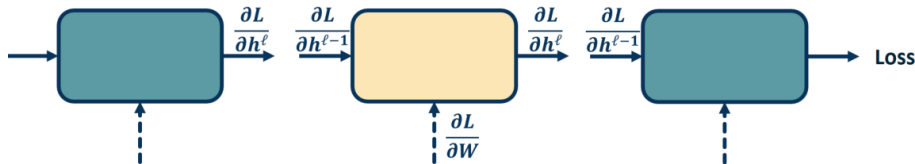
$$\bar{w} = \frac{\partial L}{\partial p} \frac{\partial p}{\partial u} \frac{\partial u}{\partial w} = -\frac{1}{\sigma(w^T x)} \sigma(w^T x) (1 - \sigma(w^T x)) x^T$$

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

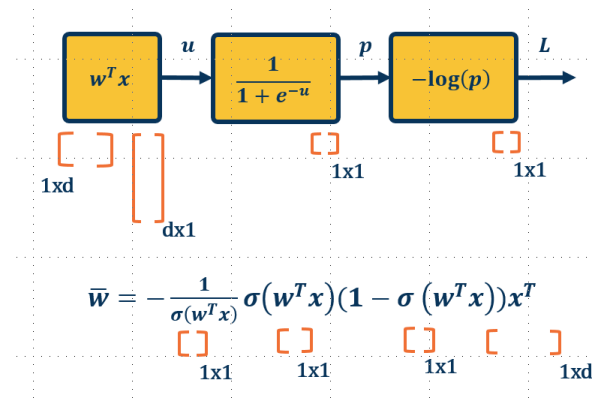
This effectively shows gradient flow along path from  $L$  to  $w$

## Computation Graph of primitives (automatic differentiation)

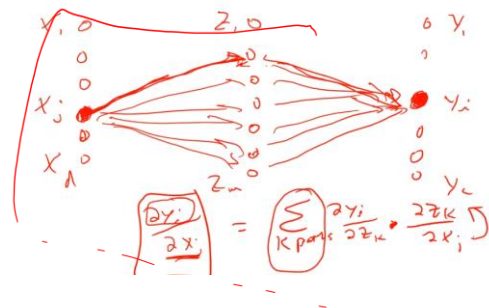
- We want to compute:  $\left\{ \frac{\partial L}{\partial h^{\ell-1}}, \frac{\partial L}{\partial W} \right\}$



## Backpropagation View (Recursive Algorithm)



## Computational / Tensor View




## Graph View

# Different Views of Equivalent Ideas

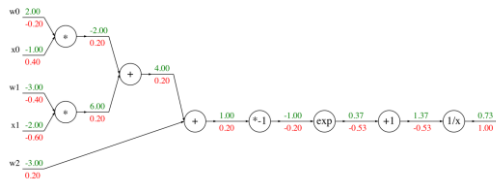


# **Backpropagation and Automatic Differentiation**

# Deep Learning = Differentiable Programming

- Computation = Graph
  - Input = Data + Parameters
  - Output = Loss
  - Scheduling = Topological ordering
- What do we need to do?
  - Generic code for representing the graph of modules 
  - Specify modules (both forward and backward function)

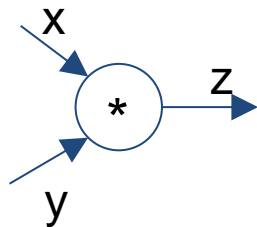
# Modularized implementation: forward / backward API



Graph (or Net) object (*rough psuedo code*)

```
class ComputationalGraph(object):  
    #...  
    def forward(inputs):  
        # 1. [pass inputs to input gates...]  
        # 2. forward the computational graph:  
        for gate in self.graph.nodes_topologically_sorted():  
            gate.forward()  
        return loss # the final gate in the graph outputs the loss  
    def backward():  
        for gate in reversed(self.graph.nodes_topologically_sorted()):  
            gate.backward() # little piece of backprop (chain rule applied)  
        return inputs_gradients
```

# Modularized implementation: forward / backward API



(x,y,z are scalars)

```
class MultiplyGate(object):
```

```
    def forward(x,y):
```

```
        z = x*y
```

```
        return z
```

```
    def backward(dz):
```

```
        # dx = ... #todo
```

```
        # dy = ... #todo
```

```
        return [dx, dy]
```

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

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

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

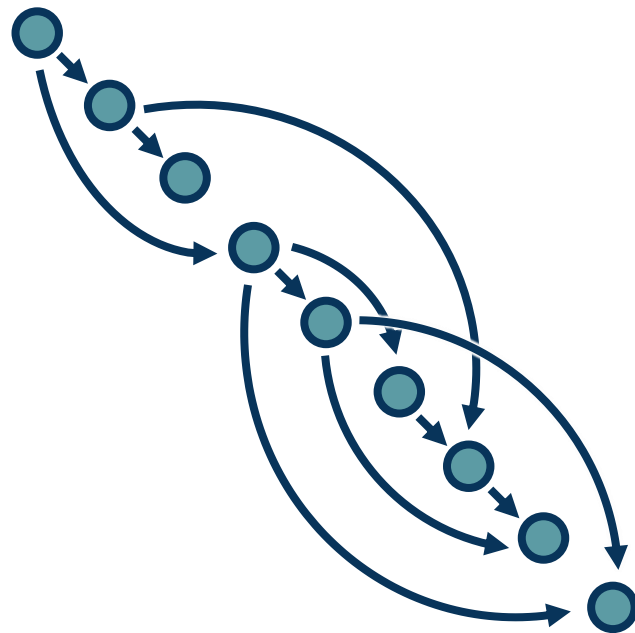
But the idea can be applied to **any directed acyclic graph (DAG)**

- 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 **local gradient function/computation for efficiency**
- We will do this **automatically** by computing backwards function for primitives and as you write code, express the function with them

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

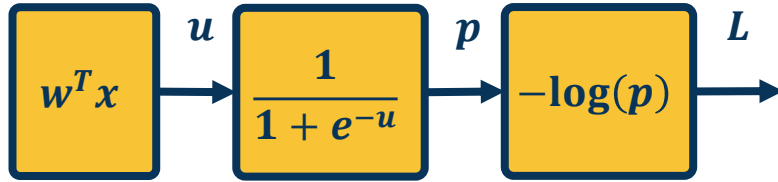


## Computation = Graph

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

## Auto-Diff

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



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

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

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

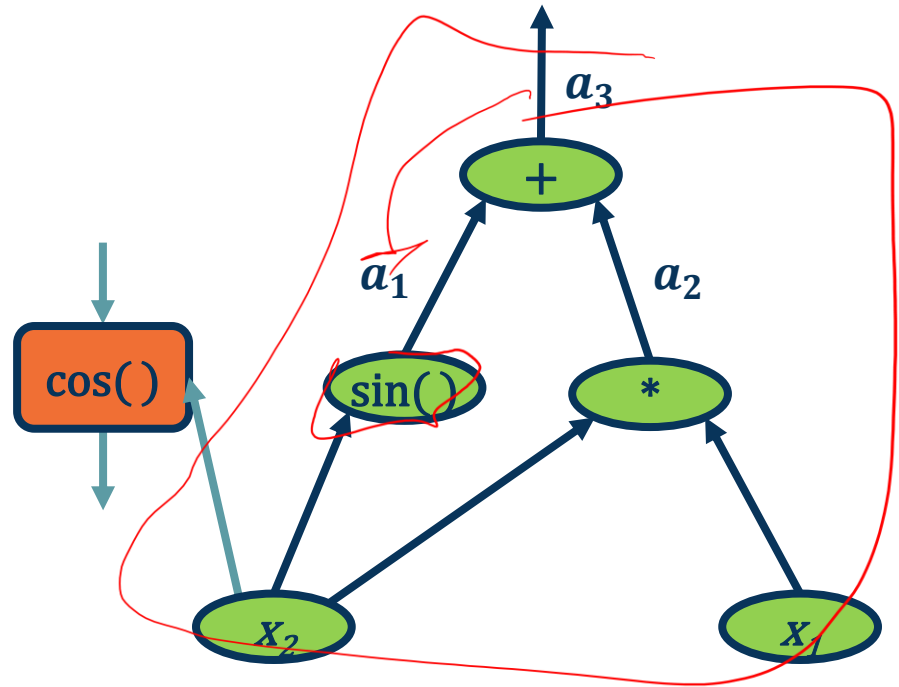
We can do this in a combined way to see all terms together:

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

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

- 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**

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





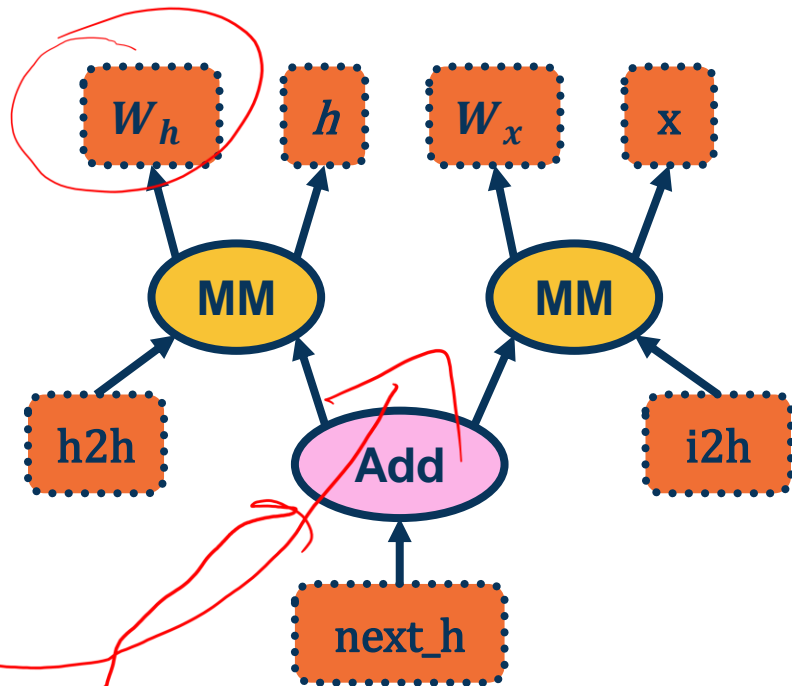
## A graph is created on the fly

```
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)



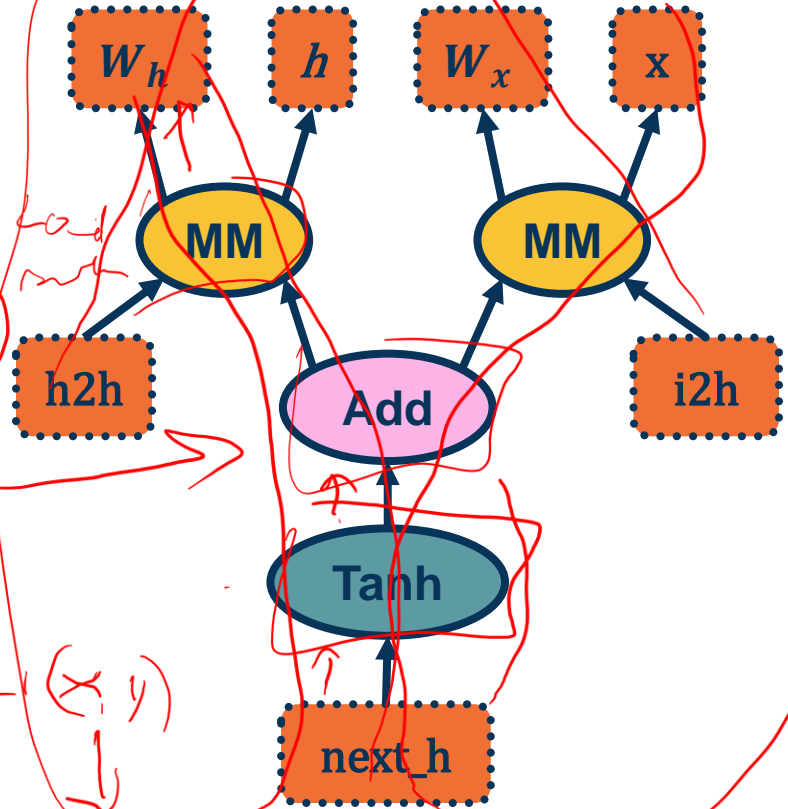
# Back-propagation uses the dynamically built graph

```
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))
```

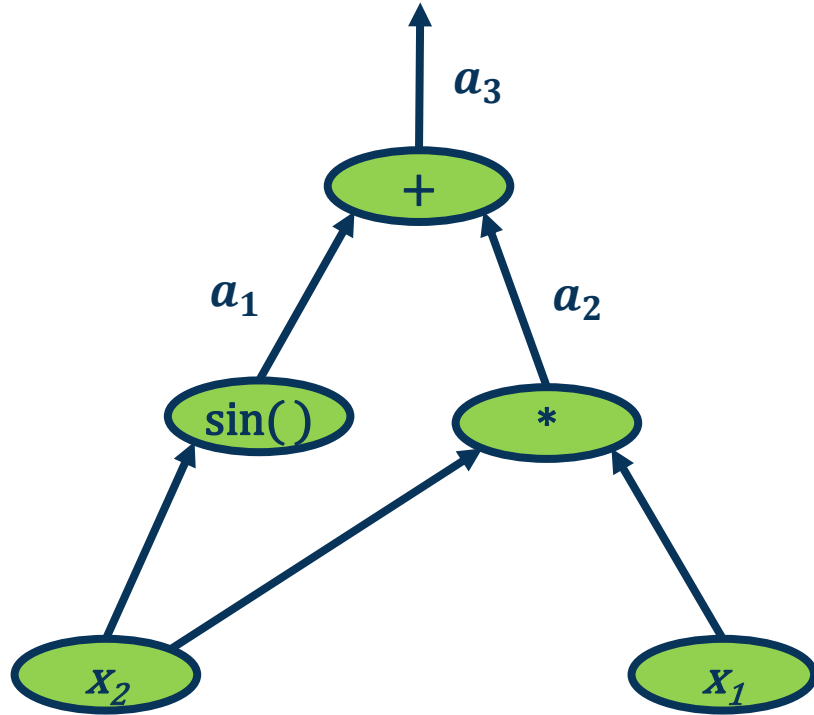
```
i2h = torch.mm(W_x, x.t())  
h2h = torch.mm(W_h, prev_h.t())  
next_h = i2h + h2h  
next_h = next_h.tanh()
```

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



From pytorch.org

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



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

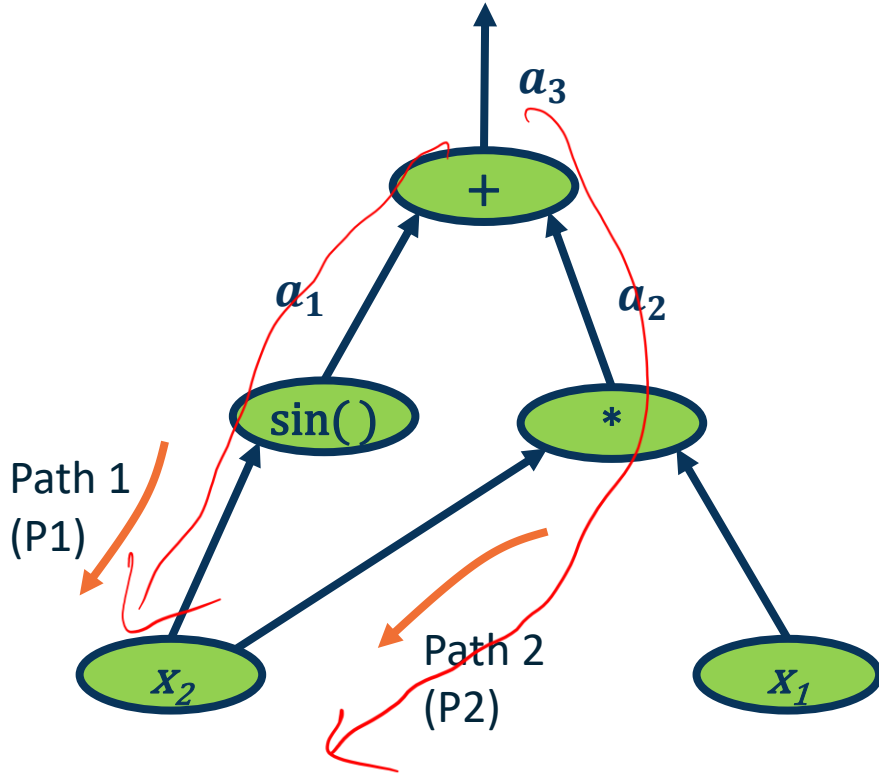
- Assign intermediate variables

**Simplify notation:**

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

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

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



$$\bar{a}_3 = \frac{\partial f}{\partial a_3} = 1$$

$$\bar{a}_1 = \frac{\partial f}{\partial a_1} = \frac{\partial f}{\partial a_3} \frac{\partial a_3}{\partial a_1} = \frac{\partial f}{\partial a_3} \frac{\partial(a_1+a_2)}{\partial a_1} = \frac{\partial f}{\partial a_3} 1 = \bar{a}_3$$

$$\bar{a}_2 = \frac{\partial f}{\partial a_2} = \frac{\partial f}{\partial a_3} \frac{\partial a_3}{\partial a_2} = \bar{a}_3$$

$$\bar{x}_2^{P1} = \frac{\partial f}{\partial a_1} \frac{\partial a_1}{\partial x_2} = \bar{a}_1 \cos(x_2)$$

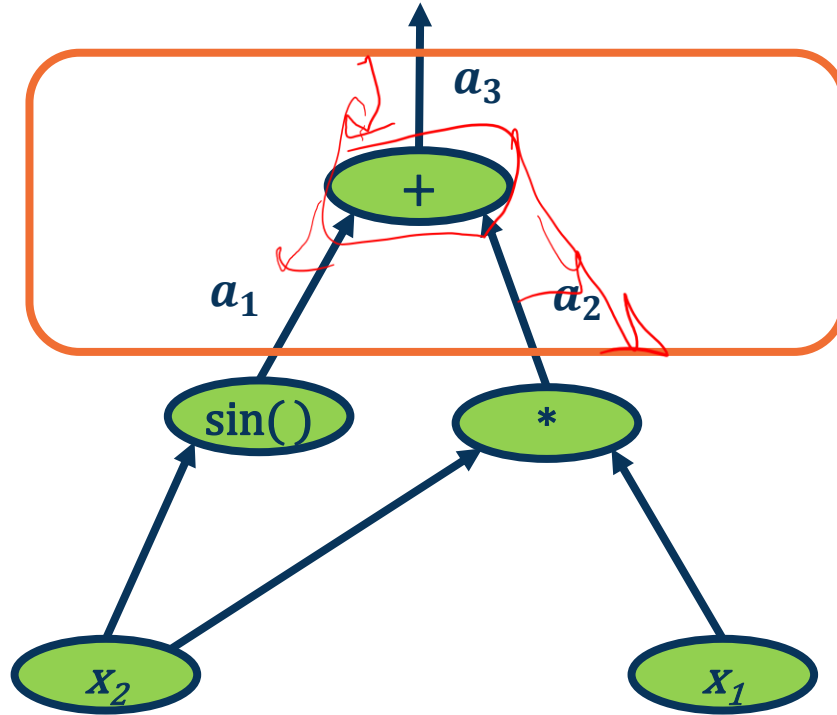
$$\bar{x}_2^{P2} = \frac{\partial f}{\partial a_2} \frac{\partial a_2}{\partial x_2} = \frac{\partial f}{\partial a_2} \frac{\partial(x_1x_2)}{\partial x_2} = \bar{a}_2 x_1$$

$$\bar{x}_1 = \frac{\partial f}{\partial a_2} \frac{\partial a_2}{\partial x_1} = \bar{a}_2 x_2$$

Gradients from multiple paths summed

Example

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

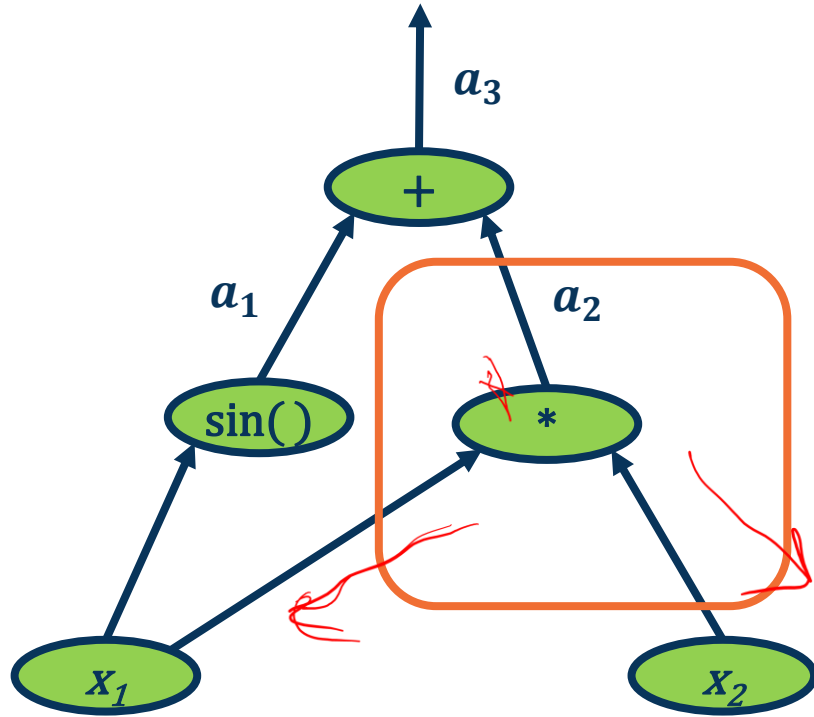


$$\bar{a}_1 = \frac{\partial f}{\partial a_1} = \frac{\partial f}{\partial a_3} \frac{\partial a_3}{\partial a_1} = \frac{\partial f}{\partial a_3} \frac{\partial(a_1+a_2)}{\partial a_1} = \frac{\partial f}{\partial a_3} \cdot 1 = \bar{a}_3$$

$$\bar{a}_2 = \frac{\partial f}{\partial a_2} = \frac{\partial f}{\partial a_3} \frac{\partial a_3}{\partial a_2} = \bar{a}_3$$

**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)

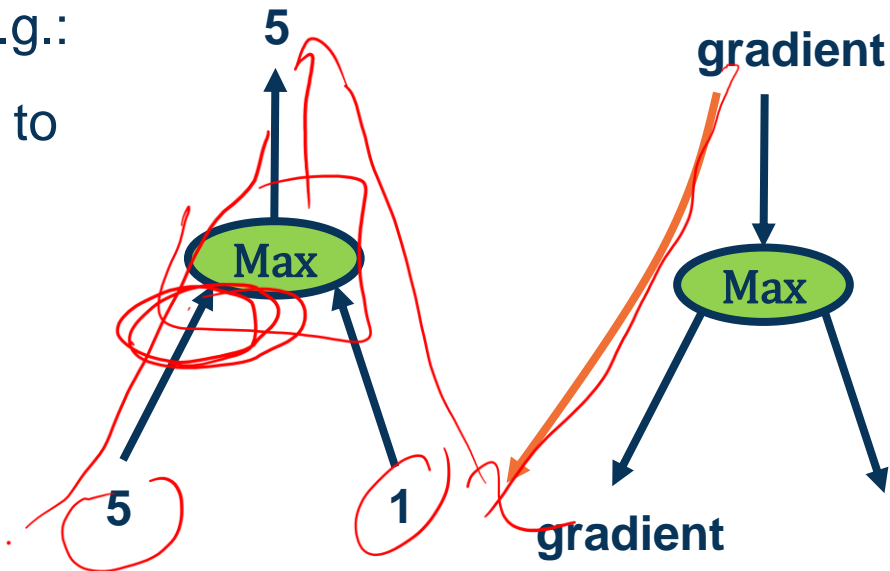
$$\overline{x_2} = \frac{\partial f}{\partial a_2} \frac{\partial a_2}{\partial x_2} = \frac{\partial f}{\partial a_2} \frac{\partial (x_1x_2)}{\partial x_2} = \overline{a_2}x_1$$

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

**Several other patterns** as well, e.g.:

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

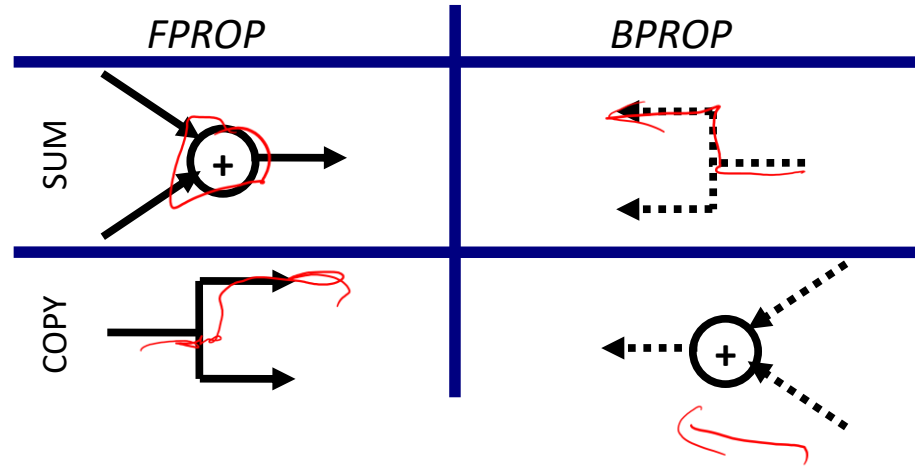
- Gradient flows along the path that was “selected” to be max
- This information must be recorded in the forward pass



**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!

# Duality in Fprop and Bprop





# Convolutional network (AlexNet)

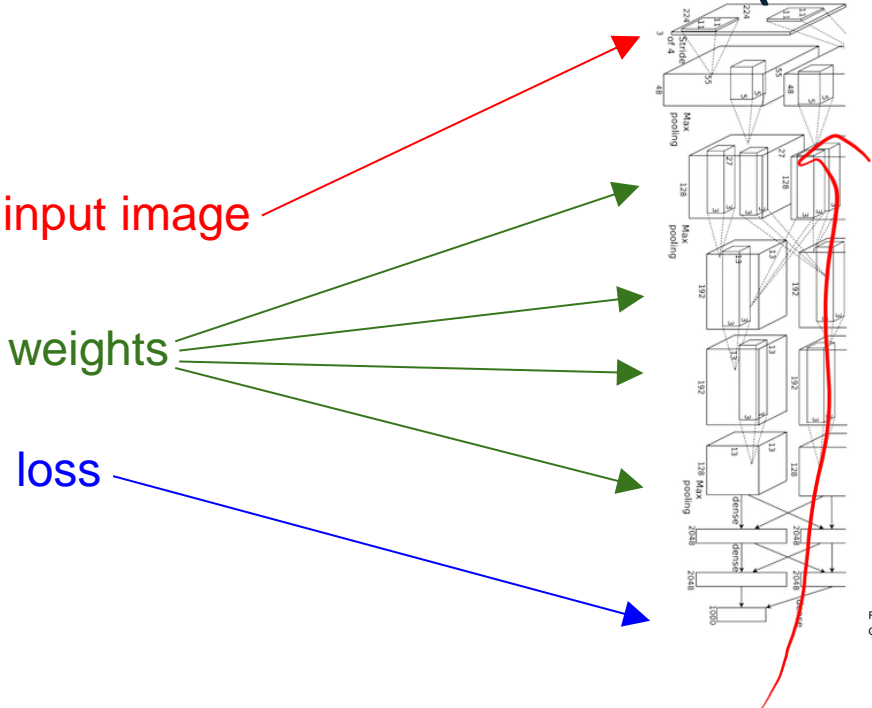


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

# Neural Turing Machine

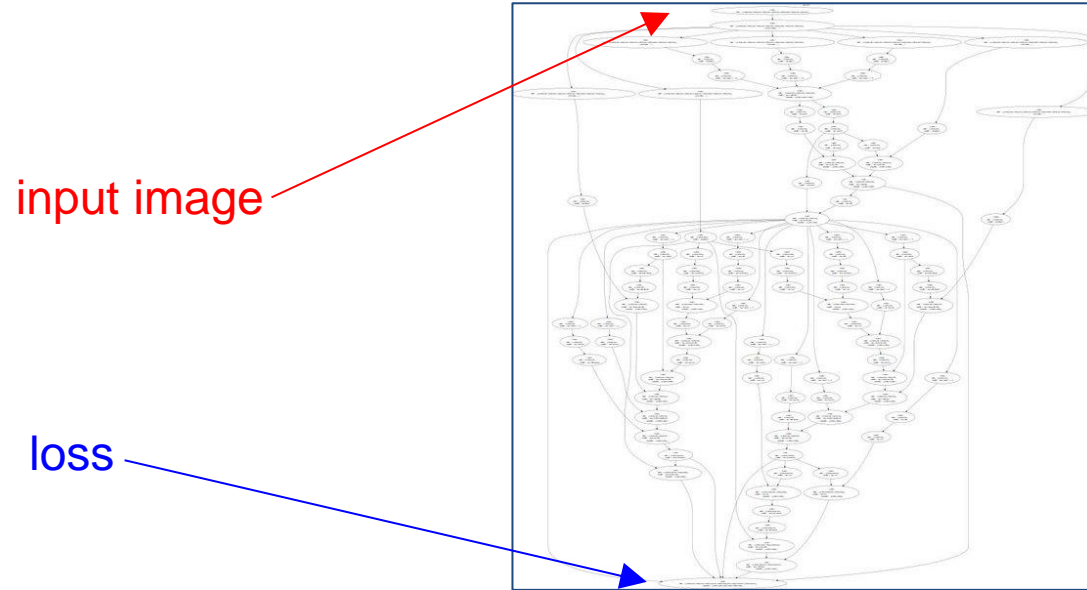
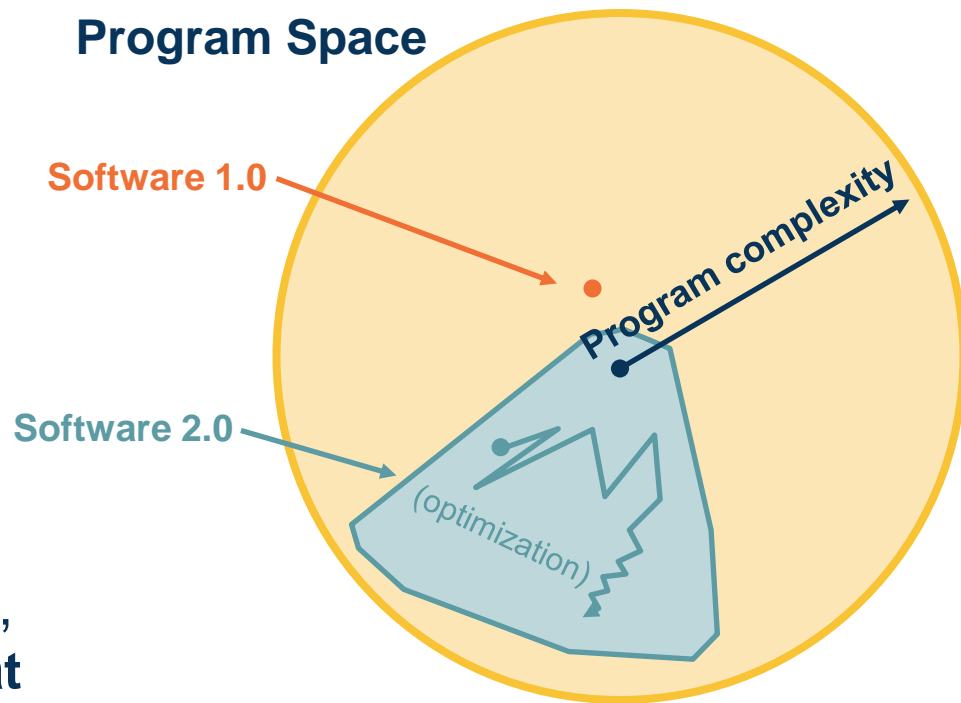


Figure reproduced with permission from a [Twitter post](#) by Andrej Karpathy.

- 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*