Topics:

- Jacobians/Matrix Calculus continued
- Backpropagation / Automatic Differentiation
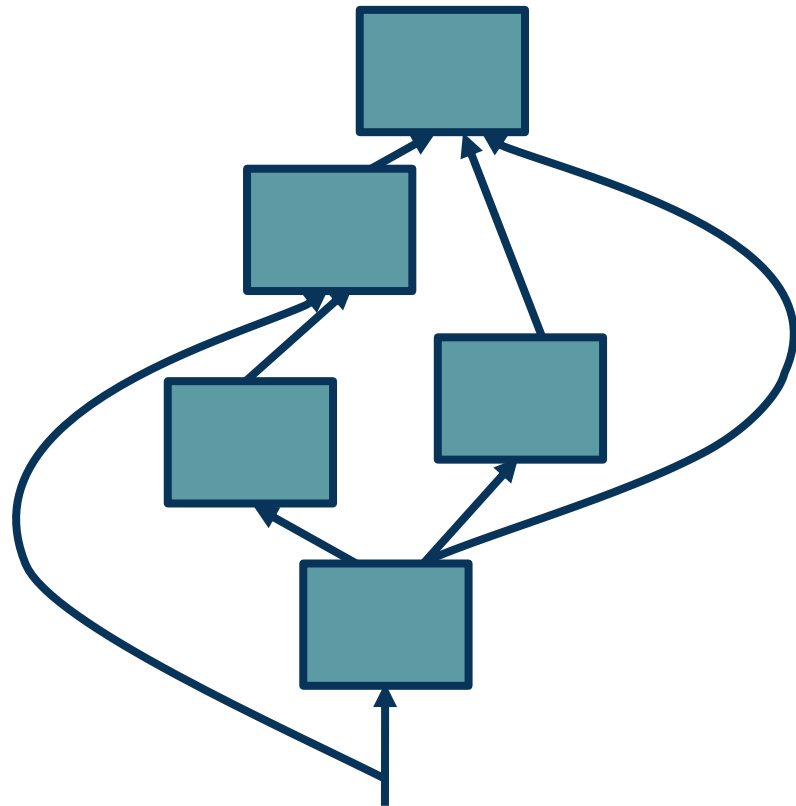
# CS 4644 / 7643-A
# ZSOLT KIRA

- **Assignment 1 out!**
  - **Due Feb 4th**
  - 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.
  - **Topic OH:** Assignment 1 and Matrix Calculus

- **In-class Quiz (30 mins) – Feb 11**

- Piazza: Project teaming thread
  - **Project Proposal: Feb. 14th**, **Project Check-in: Mar. 14th.**
  - Project proposal overview during my OH (Thursday 2pm ET, recorded)

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

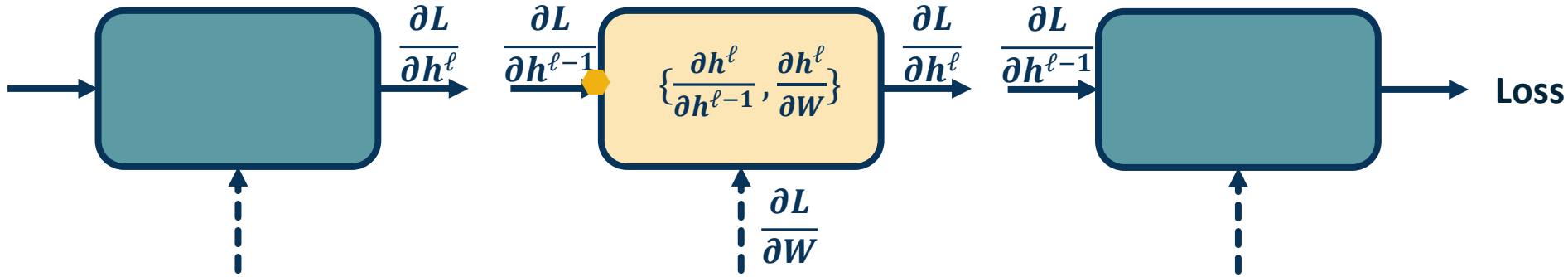Graph can be any **directed acyclic graph (DAG)**

◆ 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

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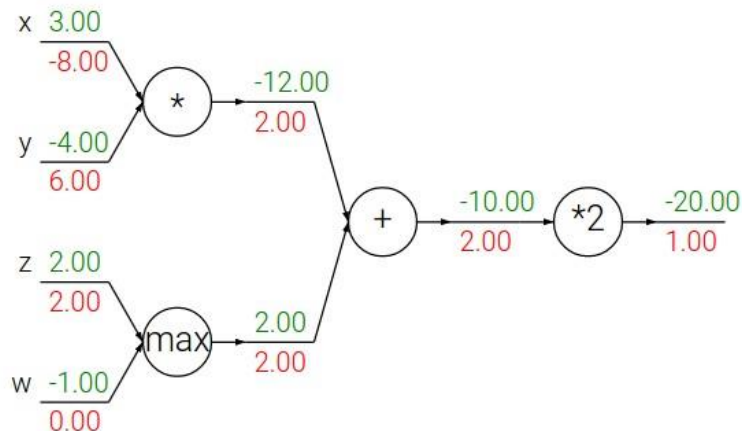


We will use the *chain rule* to do this:

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

**Computing the Gradients of Loss**

Georgia Tech

# Patterns in backward flow

**add** gate: gradient distributor

**max** gate: gradient router

**mul** gate: gradient switcher

Georgia Tech

# Conventions:

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

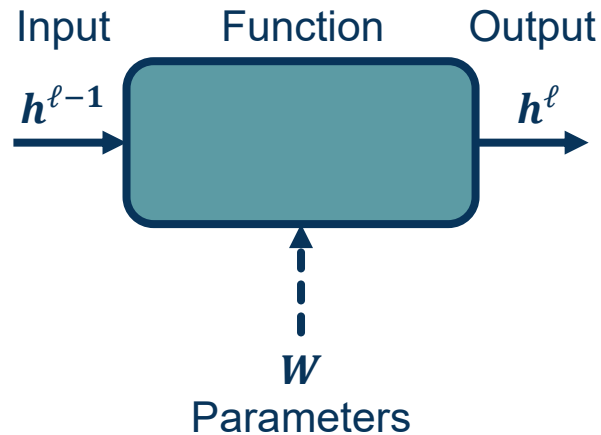|  | $S$ $[\quad]$ | $V$ $\begin{bmatrix} \\ \end{bmatrix}$ | $M$ $\begin{bmatrix} & \\ & \end{bmatrix}$ |
|---|---|---|---|
| $S$ | $\dfrac{\partial s_1}{\partial s_2}$ $[\quad]$ | $\dfrac{\partial s}{\partial v}$ $[\quad\quad]$ | $\dfrac{\partial s}{\partial M}$ $\begin{bmatrix} & \\ & \end{bmatrix}$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

**Tensors**

Georgia Tech

- 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:

$$W$$

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

**Dimensionality of Derivatives in ML**

Input        Function        Output

$h^{\ell-1}$ → [        ] → $h^\ell$

↑
$W$
Parameters

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

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

$$\leftarrow w_i^T \rightarrow$$

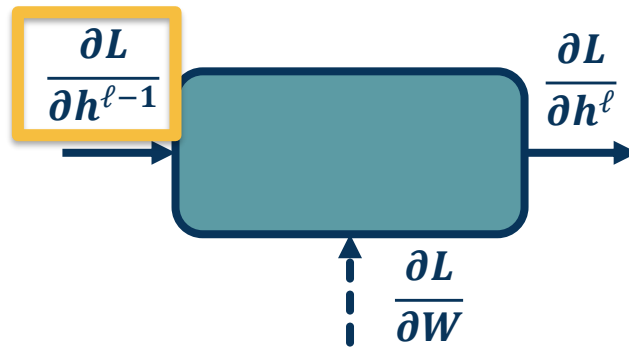$|h^\ell| \times 1$        $|h^\ell| \times |h^{\ell-1}|$        $|h^{\ell-1}| \times 1$

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

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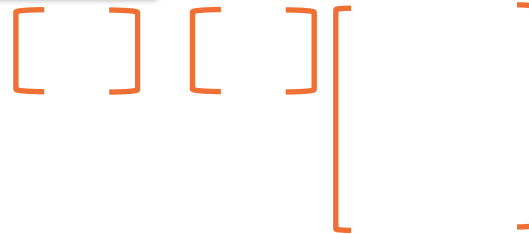
$$h^\ell = W h^{\ell-1}$$

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

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

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

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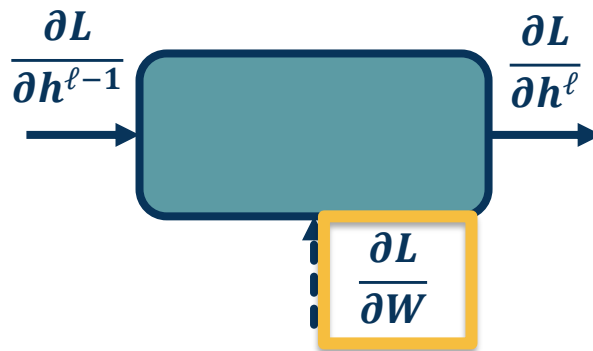
Georgia Tech

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

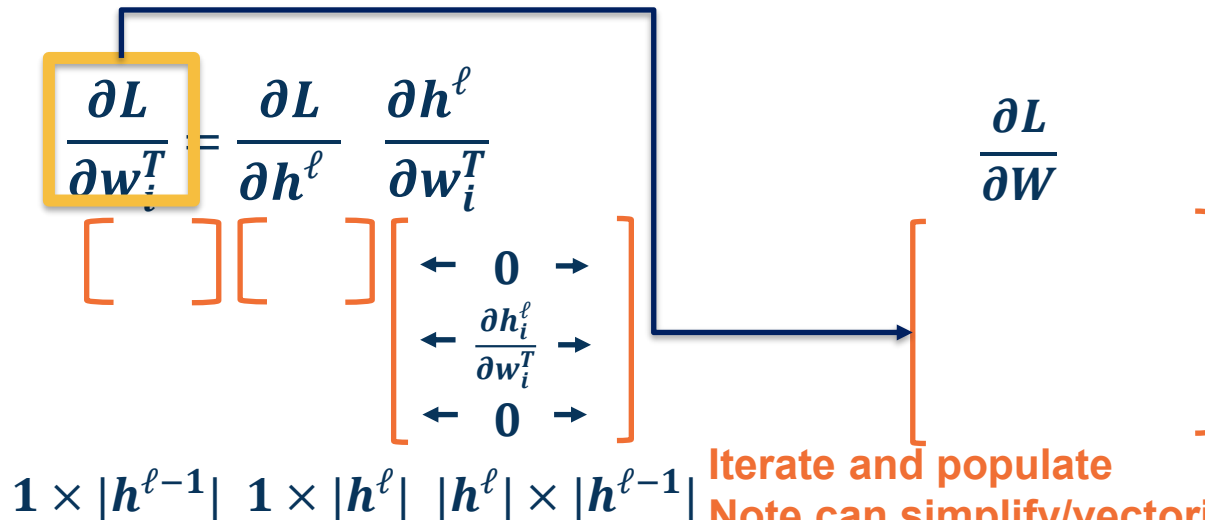$$\frac{\partial h^\ell}{\partial h^{\ell-1}} = W$$

**Define:**

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

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

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

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

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

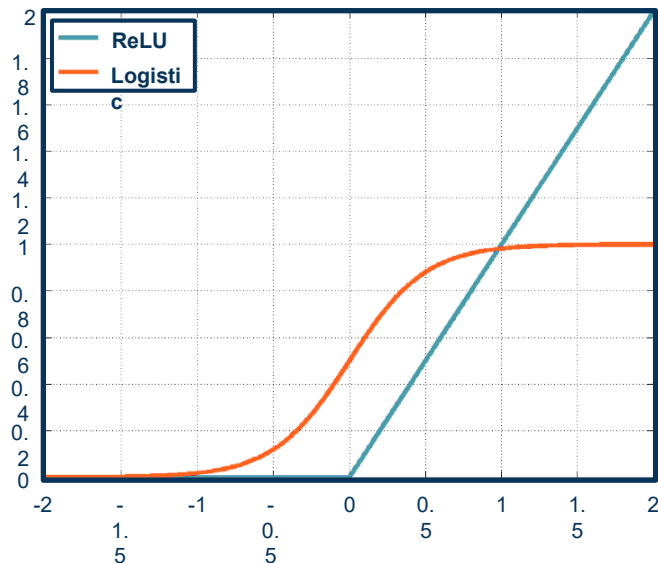But it is *sparse* – each output only affected by corresponding weight row

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

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

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

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

Iterate and populate
Note can simplify/vectorize!

**Fully Connected (FC) Layer**

Georgia Tech

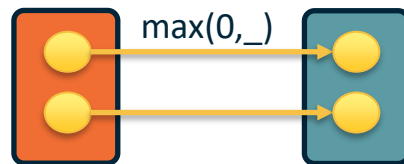We can employ **any differentiable (or piecewise differentiable) function**

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?

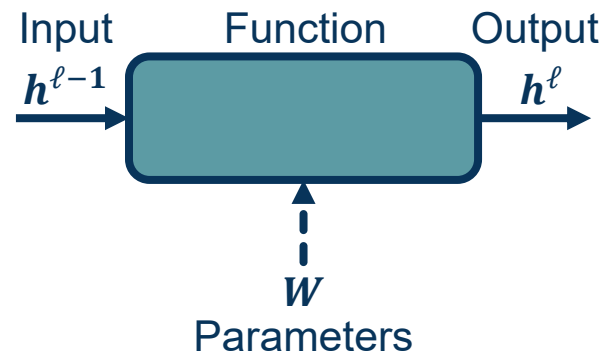

$$h^\ell = \max(0, h^{\ell-1})$$



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

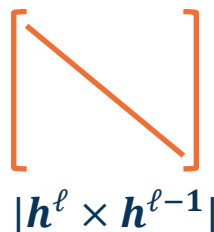- An output value affected only by **corresponding input value**

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

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

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 ]

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

**Summary**

Georgia
Tech

# Vectorization in Function Compositions

**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

$$\mathbf{x} \in \mathbb{R}^{1} \longrightarrow z \in \mathbb{R}^{1} \longrightarrow y \in \mathbb{R}^{1}$$

$$\boldsymbol{g_1} \qquad\qquad \boldsymbol{g_2}$$

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

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

**Scalar Multiplication**

Georgia Tech

$$\vec{x}\{\in \mathbb{R}^d \longrightarrow \vec{z}\{\in \mathbb{R}^m \longrightarrow \vec{y}\{\in \mathbb{R}^c$$

$$g_1 \qquad\qquad g_2$$

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

$$\left[\frac{\partial \vec{y}}{\partial \vec{x}}\right] = \left[\frac{\partial \vec{y}}{\partial \vec{z}}\right]\left[\frac{\partial \vec{z}}{\partial \vec{x}}\right]$$

$$J_{g_1 \circ g_2} \qquad\qquad J_{g_1} \qquad\qquad J_{g_2}$$

**Matrix Multiplication**

Georgia Tech

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

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

**Jacobian View of Chain Rule**

Georgia Tech

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

$k$ paths

**Graphical View of Chain Rule**

Georgia Tech

$$h^0 \in \mathbb{R}^d \longrightarrow h^1 \in \mathbb{R}^d \longrightarrow \quad \dots \quad \longrightarrow h^l \in \mathbb{R}^d$$

$$\frac{\partial h^l}{\partial h^1} = \frac{\partial h^l}{\partial h^{l-1}} \; \frac{\partial h^{l-1}}{\partial h^{l-2}} \; \dots \; \frac{\partial h^2}{\partial h^1}$$
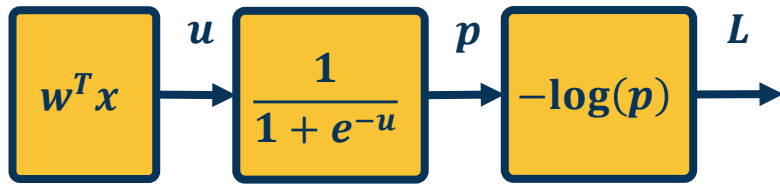
$$\begin{bmatrix} \phantom{xx} \end{bmatrix} = \begin{bmatrix} \phantom{xx} \end{bmatrix} \begin{bmatrix} \phantom{xx} \end{bmatrix} \begin{bmatrix} \phantom{xx} \end{bmatrix}$$

Georgia Tech

$$h^0 \in \mathbb{R}^d \longrightarrow h^1 \in \mathbb{R}^d \longrightarrow \ldots \longrightarrow h^l \in \mathbb{R}^d \longrightarrow L \in \mathbb{R}^1$$

$$\frac{\partial L}{\partial h^1} = \frac{\partial L}{\partial h^l} \frac{\partial h^l}{\partial h^{l-1}} \frac{\partial h^{l-1}}{\partial h^{l-2}} \ldots \frac{\partial h^2}{\partial h^1}$$

$$[\quad] = [\quad][\quad][\quad][\quad]$$

**Which directions is more efficient to multiply?**

Georgia Tech

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!

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

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

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:**

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

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

**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**

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

[ ] 1xd  [ ] dx1  [ ] 1x1  [ ] 1x1

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

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

[ ] 1x1  [ ] 1x1  [ ] 1x1  [ ] 1xd

Georgia Tech
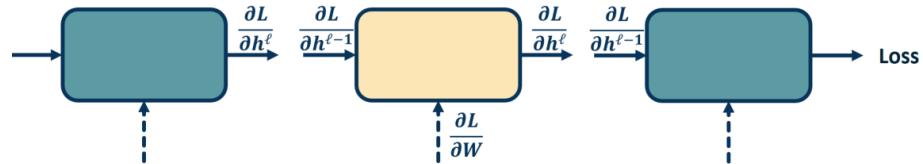
Many **standard regularization methods** still apply!

---

### 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: $|y - Wx_i|^2 + \alpha|W|^2 + \beta|W|$
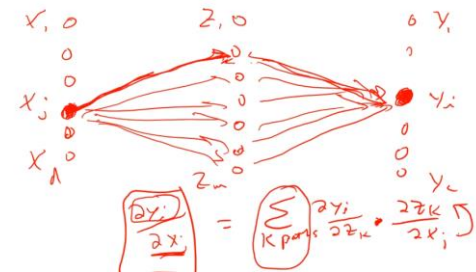
**Regularization**

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

**Backpropagation View (Recursive Algorithm)**

$L = 1$

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

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:

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

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

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

**Computation Graph of primitives (automatic differentiation)**

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

**Computational / Tensor View**

**Graph View**

**Different Views of Equivalent Ideas**

Georgia Tech

# 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
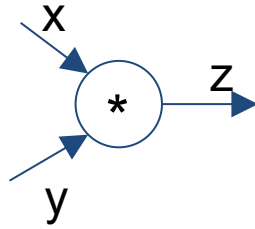  - Specify modules (both forward and backward function)

# Modularized implementation: forward / backward API



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

```python
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

*

z

y

(x,y,z are scalars)

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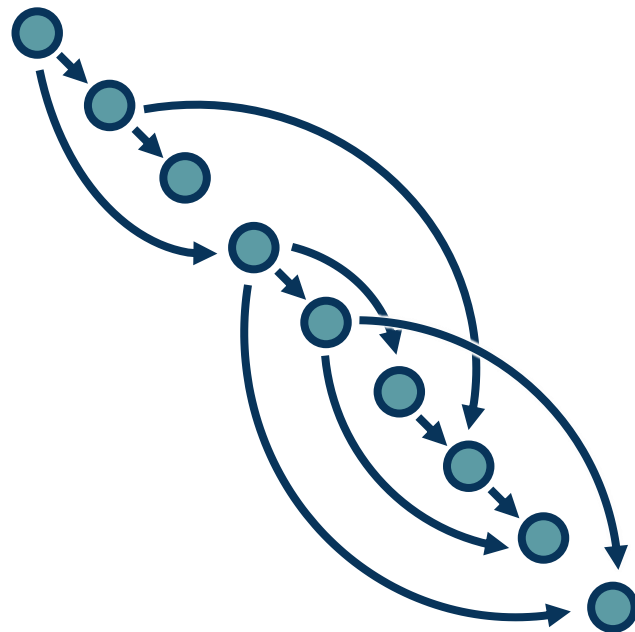
But the idea can be applied to **any directed acyclic graph (DAG)**

- 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

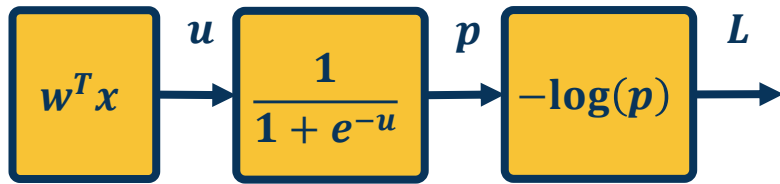This is called reverse-mode **automatic differentiation**

**Computation = Graph**

- Input = Data + Parameters

- Output = Loss

- Scheduling = Topological ordering

**Auto-Diff**

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

Georgia
Tech

$$\bar{L} = 1$$

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

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

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

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

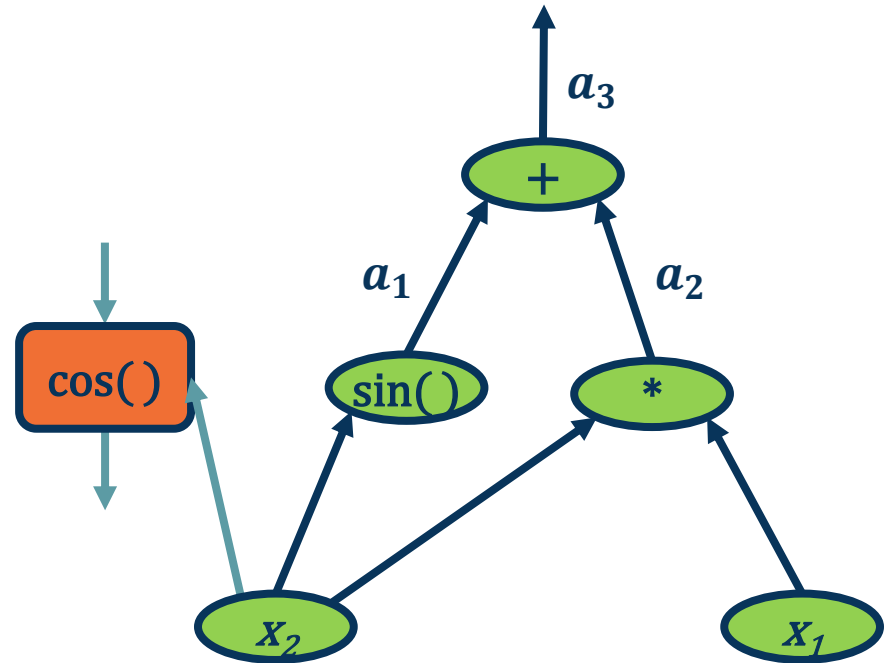**We can do this in a combined way to see all terms together:**

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

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

**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!**

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**

$$\overline{x_2} = \frac{\partial f}{\partial a_1} \ \frac{\partial a_1}{\partial x_2} = \overline{a_1} \ \cos(x_2)$$
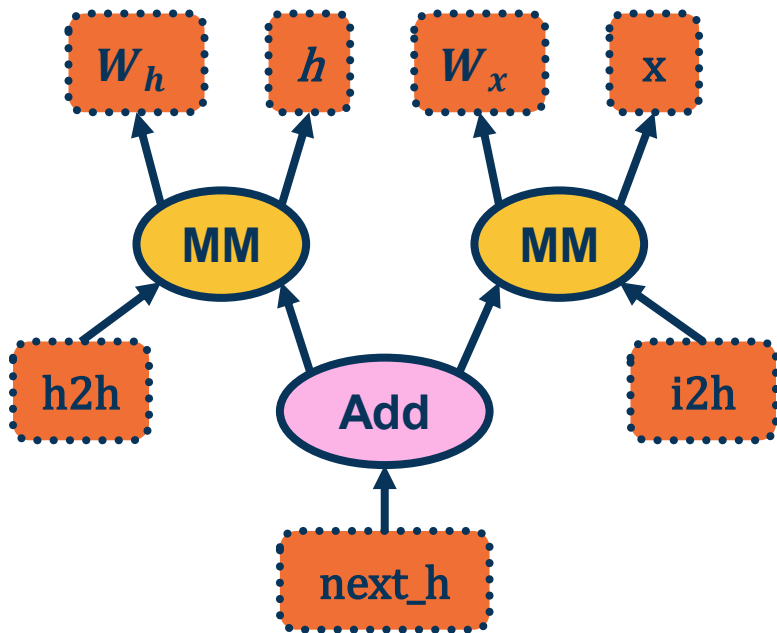
# A graph is created on the fly

```python
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
```
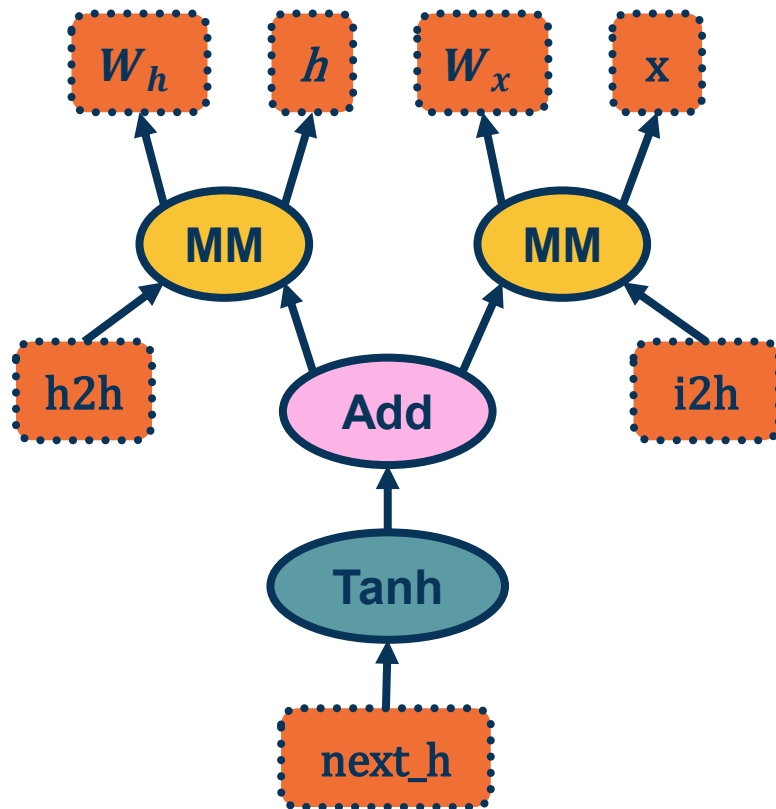
(Note above)

# Back-propagation uses the dynamically built graph

```python
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))
```
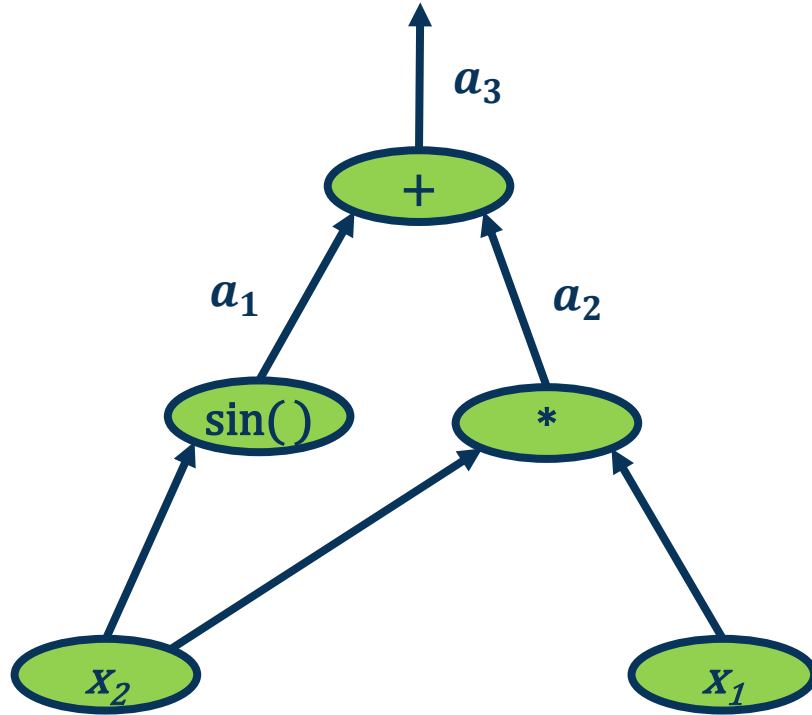


*From pytorch.org*

Georgia Tech

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

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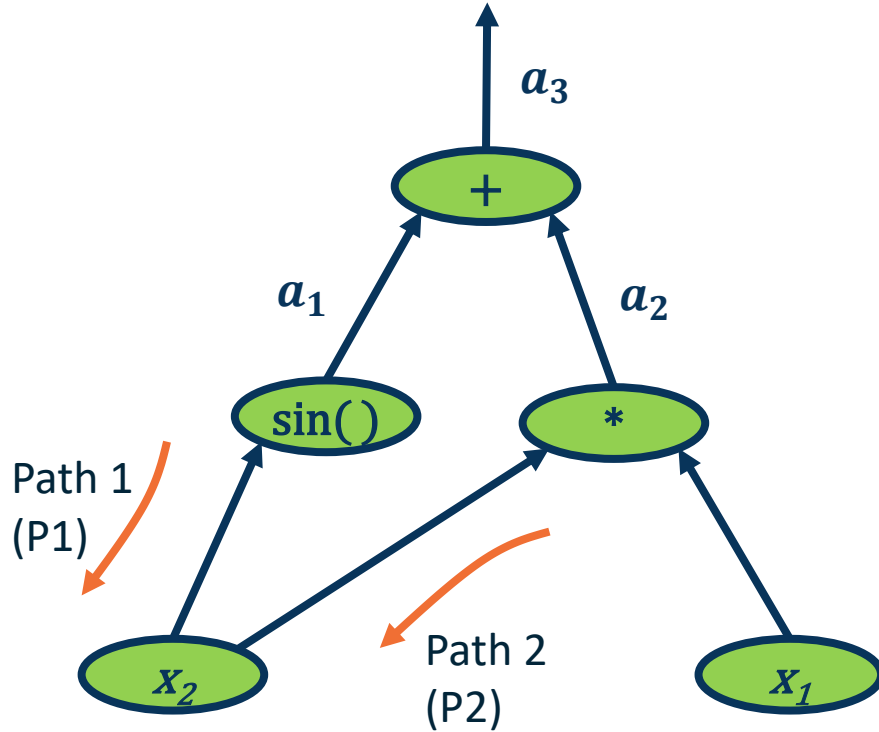
- Assign intermediate variables

  **Simplify notation:**
  **Denote bar as:** $\overline{a_3} = \dfrac{\partial f}{\partial a_3}$

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

**Example**

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

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

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

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

$$\overline{x_2^{P1}} = \frac{\partial f}{\partial a_1} \frac{\partial a_1}{\partial x_2} = \overline{a_1} \; \cos(x_2)$$
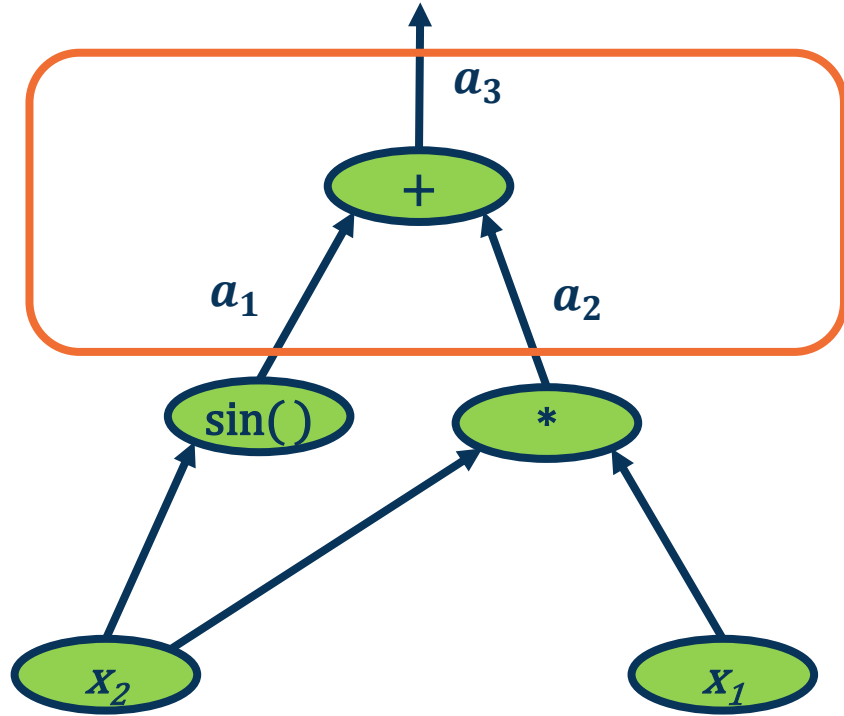
$$\overline{x_2^{P2}} = \frac{\partial f}{\partial a_2} \frac{\partial a_2}{\partial x_2} = \frac{\partial f}{\partial a_2} \frac{\partial (x_1 x_2)}{\partial x_2} = \overline{a_2} x_1$$

Gradients from multiple paths **summed**

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

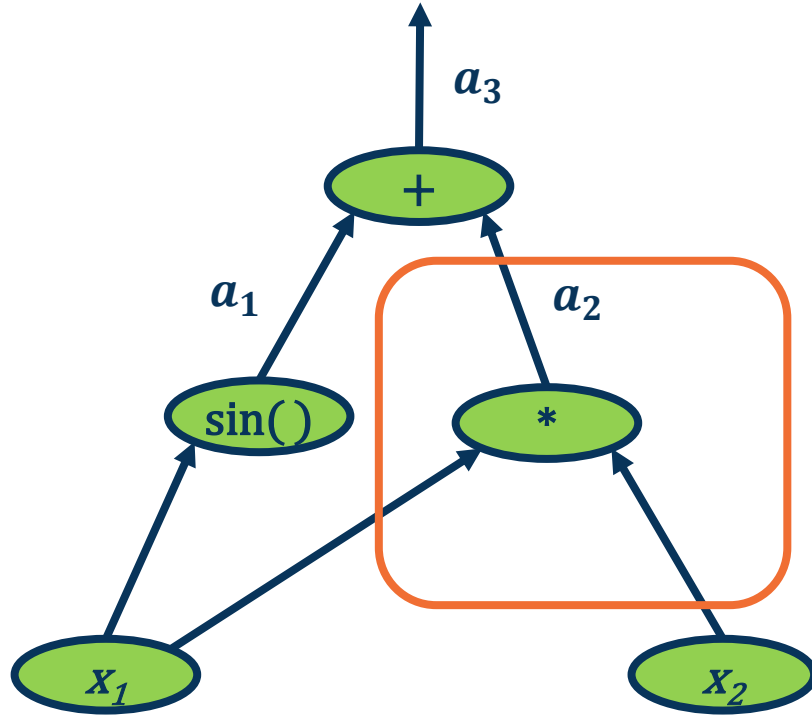**Example**

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



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

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

**Addition operation distributes gradients along all paths!**

Georgia Tech

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

$a_3$

$a_1$    $a_2$

sin()    *

$X_1$    $X_2$

**Multiplication operation is a gradient switcher (multiplies it by the values of the other term)**
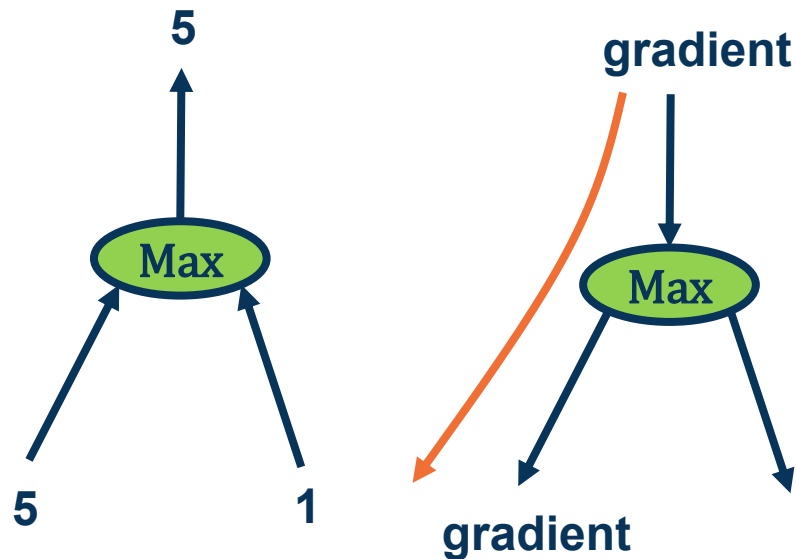
$$\overline{x_2} = \frac{\partial f}{\partial a_2}\ \frac{\partial a_2}{\partial x_2} = \frac{\partial f}{\partial a_2}\ \frac{\partial (x1x2)}{\partial x_2} = \overline{a_2} x_1$$

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

**Patterns of Gradient Flow: Multiplication**

Georgia Tech

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

Max operation **selects** which path to push the gradients through
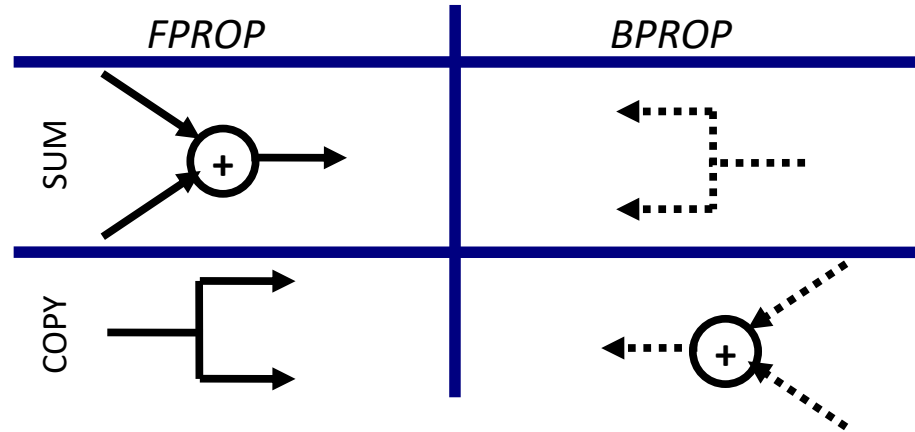
- ◆ Gradient flows along the path that was "selected" to be max

- ◆ 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!

5

Max

5          1

gradient

Max

gradient

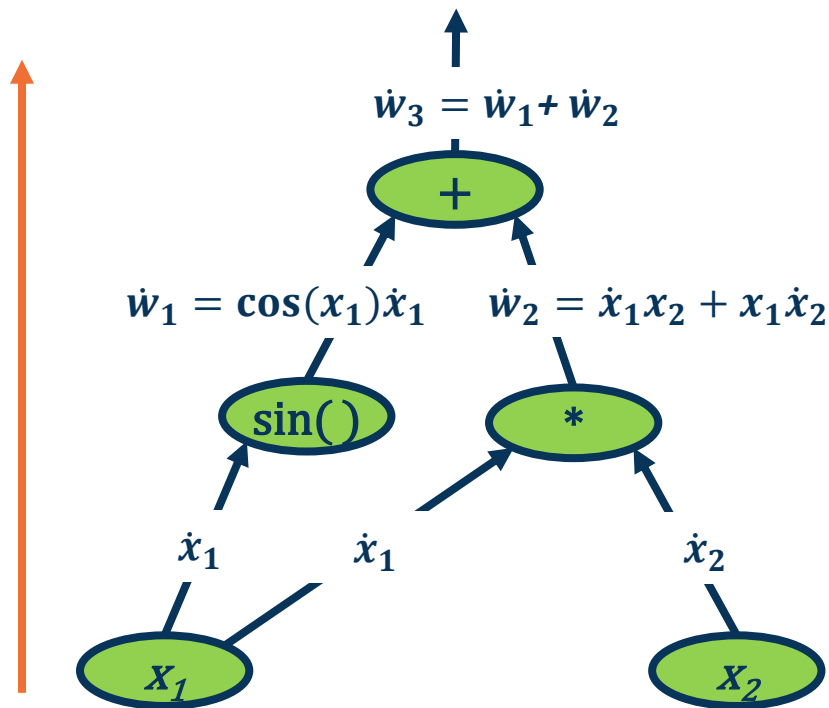Georgia Tech

# Duality in Fprop and Bprop

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

Start from **inputs** and propagate gradients forward

Complexity is proportional to input size

- Memory savings (all forward pass, no need to store activations)

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



$$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$$

$+$

$$\dot{w}_1 = \cos(x_1)\dot{x}_1 \qquad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$$

$\sin()$ $*$

$\dot{x}_1$ $\dot{x}_1$ $\dot{x}_2$

$X_1$ $X_2$

Georgia Tech
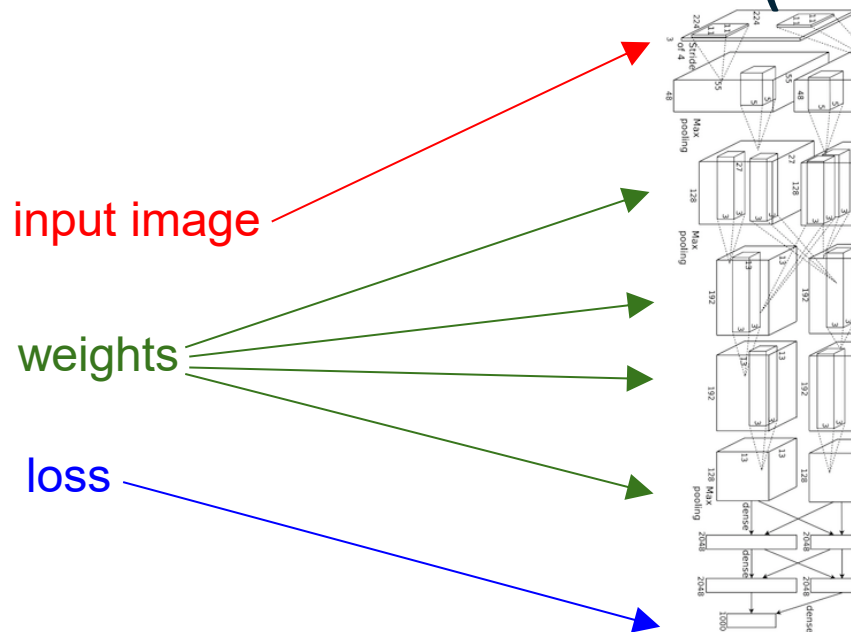
# Convolutional network (AlexNet)



input image

weights

loss

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.
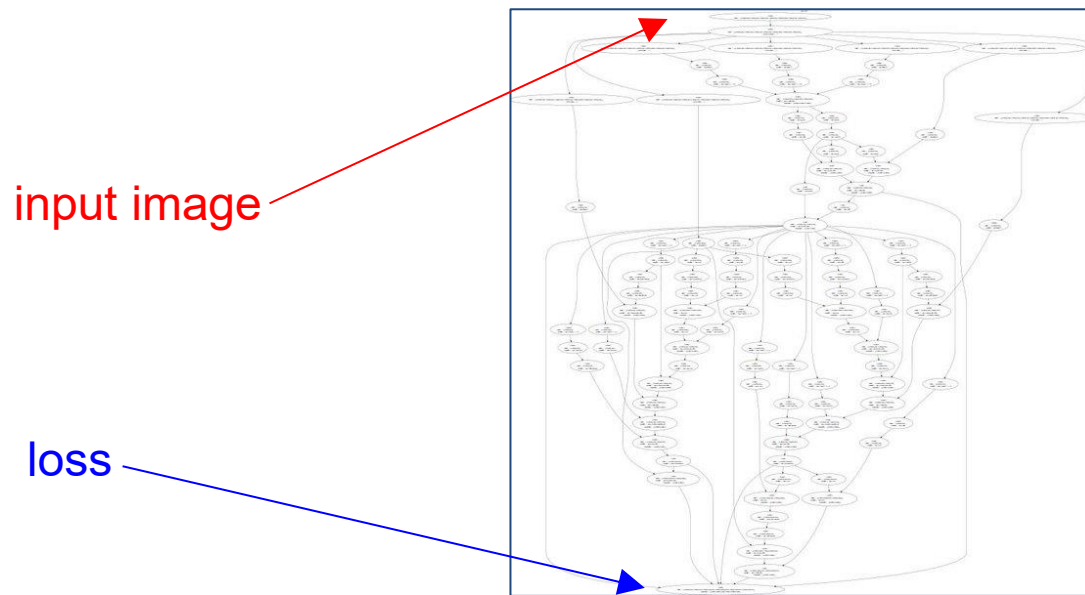
# Neural Turing Machine



input image

loss
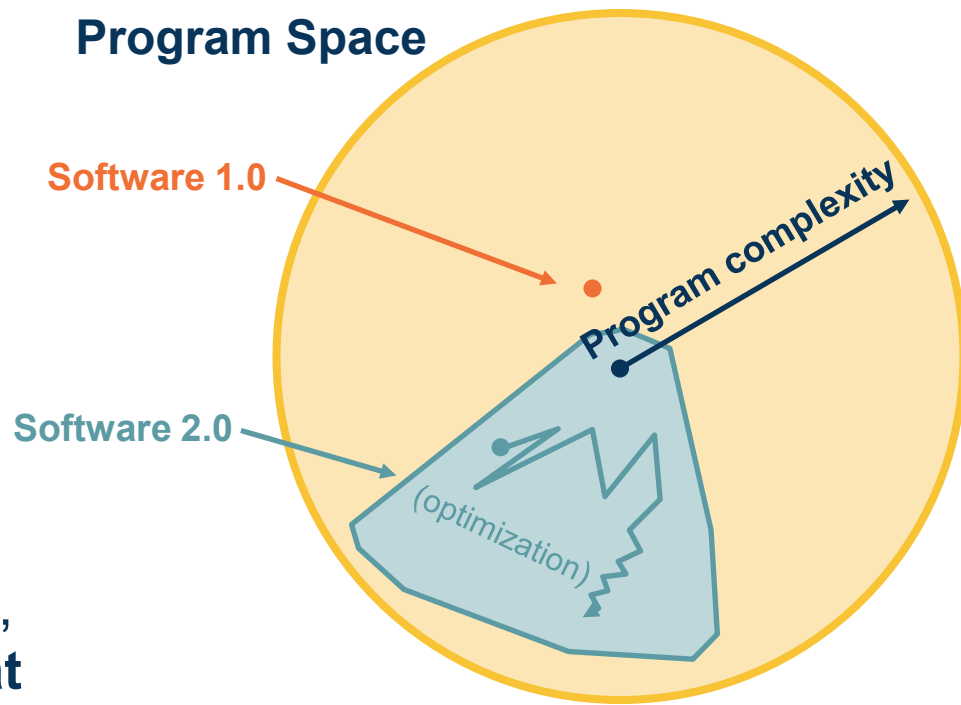
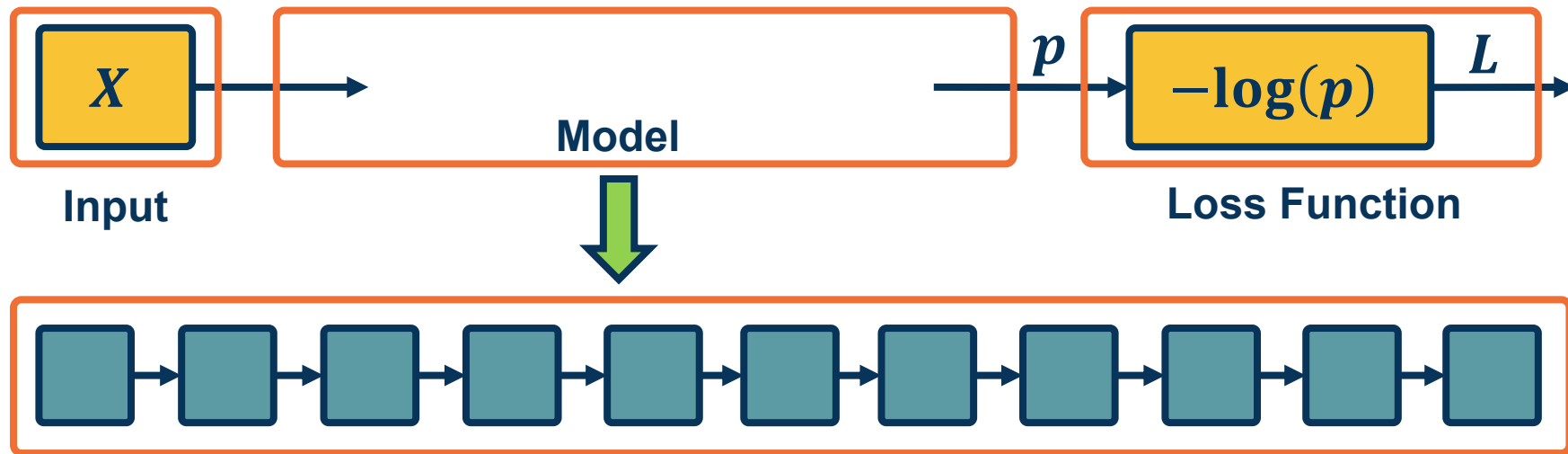Figure reproduced with permission from a Twitter post by Andrej Karpathy.

- 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**

**Program Space**

Software 1.0

Program complexity

Software 2.0

(optimization)

*Adapted from figure by Andrej Karpathy*

**Power of Automatic Differentiation**

Georgia Tech

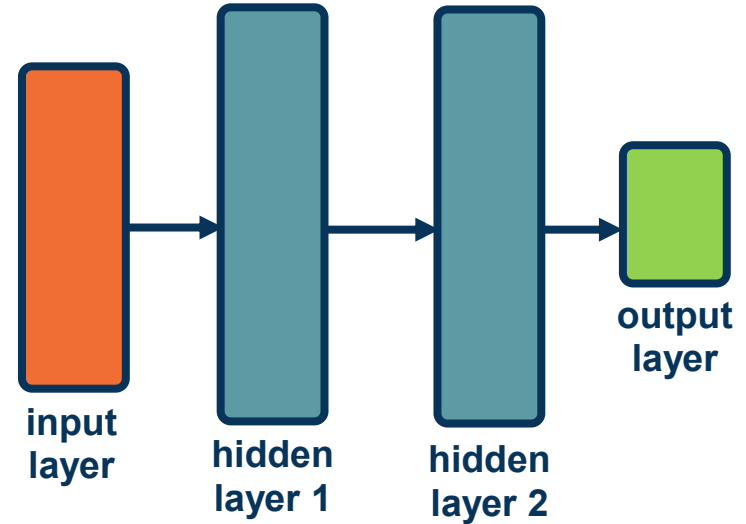Backpropagation, and automatic differentiation, allows us to optimize **any** function composed of differentiable blocks

⬡ **No need to modify** the learning algorithm!

⬡ The complexity of the function is only limited by **computation and memory**

A network with two or more hidden layers is often considered a **deep** model

**Depth is important:**

◆ Structure the model to represent an inherently compositional world

◆ Theoretical evidence that it leads to parameter efficiency

◆ Gentle dimensionality reduction (if done right)



input layer

hidden layer 1

hidden layer 2

output layer

Georgia Tech

There are still many design decisions that must be made:

◆ **Architecture**

◆ **Data Considerations**

◆ **Training and Optimization**

◆ **Machine Learning Considerations**

**?**

IM GENET

Local Minima

Georgia Tech