Topics:

- Jacobians/Matrix Calculus continued

- Backpropagation / Automatic Differentiation

# CS 4644 / 7643-A
# ZSOLT KIRA

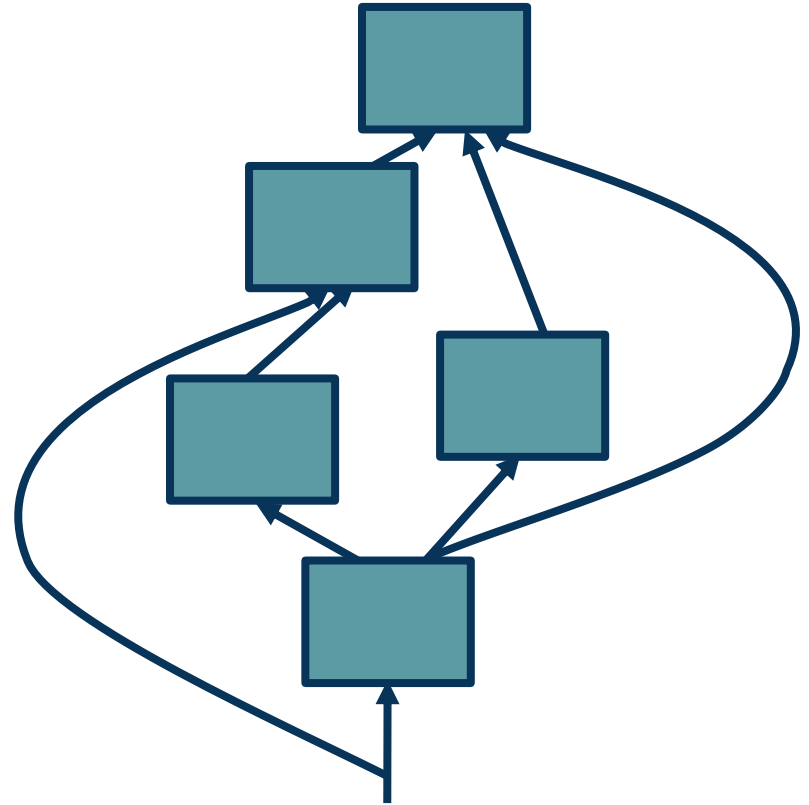- **Assignment 1 out!**
  - **Due Feb 2$^{nd}$ (with grace period 4$^{th}$ )**
  - Start now, start now, start now!
  - Start now, start now, start now!
  - Start now, start now, start now!

- Resources:
  - These lectures
  - Matrix calculus for deep learning
  - Gradients notes and MLP/ReLU Jacobian notes.
  - Assignment 1 (@57) and matrix calculus (@80), convex optimization (@82)

- Piazza: Project teaming thread
  - Will post video of project overview

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

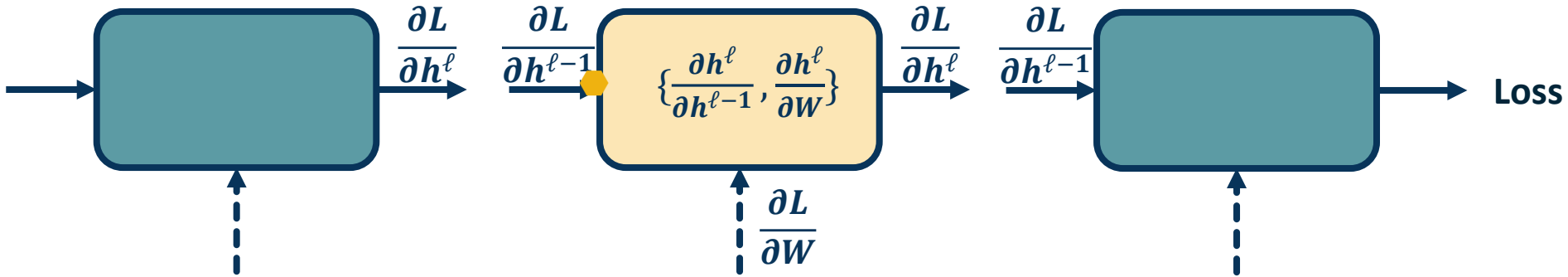Graph can be any **directed acyclic graph (DAG)**

◆ 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*

**A General Framework**

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

$$\frac{\partial L}{\partial h^{\ell}} \qquad \frac{\partial L}{\partial h^{\ell-1}} \quad \left\{\frac{\partial h^{\ell}}{\partial h^{\ell-1}}, \frac{\partial h^{\ell}}{\partial W}\right\} \qquad \frac{\partial L}{\partial h^{\ell}} \qquad \frac{\partial L}{\partial h^{\ell-1}} \qquad \text{Loss}$$

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

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

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

**Computing the Gradients of Loss**
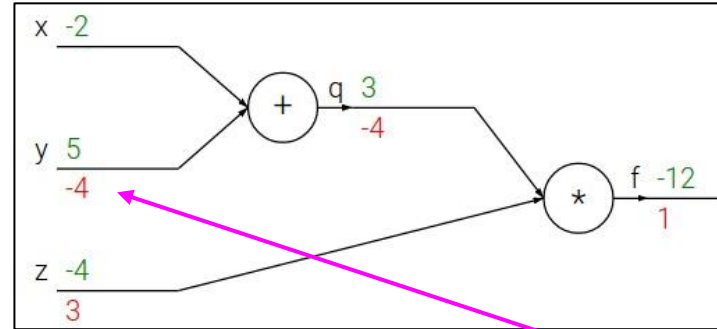
Georgia Tech

# Backpropagation: a simple example

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

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

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

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

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

Chain rule:

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

Upstream gradient    Local gradient

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

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

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

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:

$$W$$

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

Input    Function    Output

$h^{\ell-1}$ → [ Function ] → $h^{\ell}$

$W$
Parameters

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

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

$\leftarrow w_i^T \rightarrow$

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

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

Georgia Tech

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

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

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

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

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

$$\begin{bmatrix} \end{bmatrix} \qquad \begin{bmatrix} \end{bmatrix} \qquad \begin{bmatrix} \end{bmatrix}$$

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

**Fully Connected (FC) Layer**

Georgia Tech

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

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

**Define:**

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

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

$$\frac{\partial L}{\partial h^{\ell-1}} \quad \boxed{\phantom{xxxx}} \quad \frac{\partial L}{\partial h^{\ell}}$$

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

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

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

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

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

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

$$\begin{bmatrix} \phantom{xxxxxx} \end{bmatrix}$$

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

Iterate and populate
Note can simplify/vectorize!

**Fully Connected (FC) Layer**

Georgia Tech

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

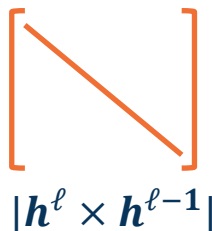- An output value affected only by **corresponding input value**

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

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

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**

- 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**

Georgia
Tech

# Vectorizaiton 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):**

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

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

(Many of these will be parameterized)

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

Georgia Tech

$$X \in \mathbb{R}^1 \xrightarrow[g_1()]{} Z \in \mathbb{R}^1 \xrightarrow[g_2()]{} Y \in \mathbb{R}^1 \quad \swarrow^{LOSS}$$

$$Y = g_2\left(g_1(x)\right)$$

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

Scalar
mult

Georgia
Tech

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

$$g_1() \qquad\qquad g_2()$$

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

$$\overset{\text{matrix}}{\left[\dfrac{\partial \vec{y}}{\partial \vec{x}}\right]} = \overset{\text{matrix}}{\left[\dfrac{\partial \vec{y}}{\partial \vec{z}}\right]} \cdot \left[\dfrac{\partial \vec{z}}{\partial \vec{x}}\right]$$

$$J_{g_2 \circ g_1} = J_{g_2} \cdot J_{g_1}$$

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!

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



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

**Neural Network Computation Graph**

Georgia Tech

$$\bar{L} = 1$$

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

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

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

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

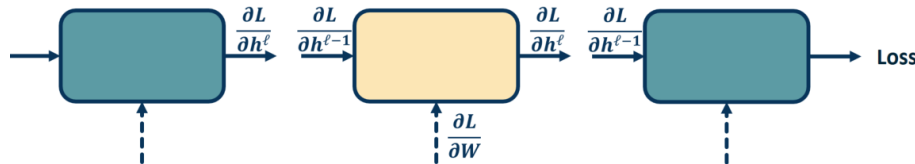**We can do this in a combined way to see all terms together:**

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

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

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

**Example Gradient Computations**

Georgia Tech

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

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

$$[ \quad ]$$
1xd

$$\begin{bmatrix} \\ \\ \end{bmatrix}$$
dx1

$$[ \ ]$$
1x1

$$[ \ ]$$
1x1

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

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

$$[ \ ]$$
1x1

$$[ \ ]$$
1x1

$$[ \ ]$$
1x1

$$[ \quad ]$$
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: $|y - Wx_i|^2 + \alpha|W|^2 + \beta|W|$

Georgia Tech

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

**Backpropagation View (Recursive Algorithm)**

$L = 1$

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

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

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

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

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

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

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

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

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

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

**Computational / Tensor View**

**Graph View**

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)*

```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]
```

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

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

# Modularized implementation: forward / backward API

x

*

z

y

(x,y,z are scalars)

```
class MultiplyGate(object):
    def forward(x,y):
        z = x*y
        self.x = x # must keep these around!
        self.y = y
        return z
    def backward(dz):
        dx = self.y * dz # [dz/dx * dL/dz]
        dy = self.x * dz # [dz/dy * dL/dz]
        return [dx, dy]
```

# Example: Caffe layers

# Caffe Sigmoid Layer

```cpp
1   #include <cmath>
2   #include <vector>
3
4   #include "caffe/layers/sigmoid_layer.hpp"
5
6   namespace caffe {
7
8   template <typename Dtype>
9   inline Dtype sigmoid(Dtype x) {
10    return 1. / (1. + exp(-x));
11  }
12
13  template <typename Dtype>
14  void SigmoidLayer<Dtype>::Forward_cpu(const vector<Blob<Dtype>*>& bottom,
15      const vector<Blob<Dtype>*>& top) {
16    const Dtype* bottom_data = bottom[0]->cpu_data();
17    Dtype* top_data = top[0]->mutable_cpu_data();
18    const int count = bottom[0]->count();
19    for (int i = 0; i < count; ++i) {
20      top_data[i] = sigmoid(bottom_data[i]);
21    }
22  }
23
24  template <typename Dtype>
25  void SigmoidLayer<Dtype>::Backward_cpu(const vector<Blob<Dtype>*>& top,
26      const vector<bool>& propagate_down,
27      const vector<Blob<Dtype>*>& bottom) {
28    if (propagate_down[0]) {
29      const Dtype* top_data = top[0]->cpu_data();
30      const Dtype* top_diff = top[0]->cpu_diff();
31      Dtype* bottom_diff = bottom[0]->mutable_cpu_diff();
32      const int count = bottom[0]->count();
33      for (int i = 0; i < count; ++i) {
34        const Dtype sigmoid_x = top_data[i];
35        bottom_diff[i] = top_diff[i] * sigmoid_x * (1. - sigmoid_x);
36      }
37    }
38  }
39
40  #ifdef CPU_ONLY
41  STUB_GPU(SigmoidLayer);
42  #endif
43
44  INSTANTIATE_CLASS(SigmoidLayer);
45
46
47  } // namespace caffe
```

Caffe is licensed under BSD 2-Clause

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

$(1 - \sigma(x))\sigma(x)$ * top_diff  (chain rule)

Georgia Tech

Backpropagation does not really spell out how to **efficiently** carry out the necessary computations
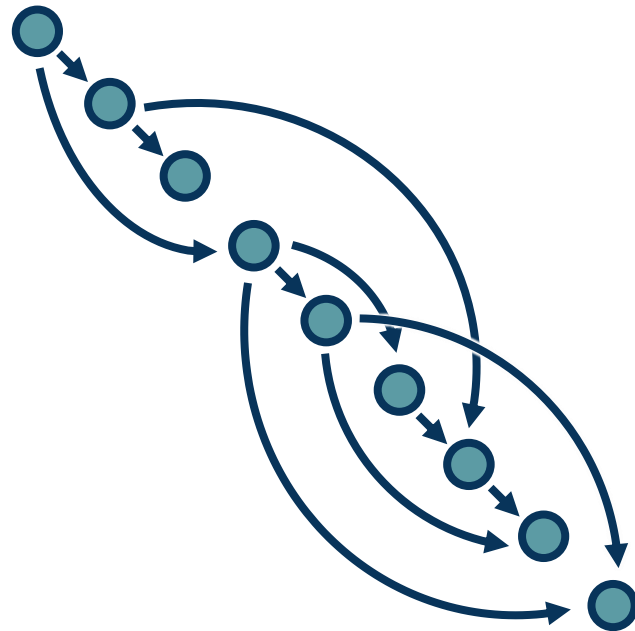
But the idea can be applied to **any directed acyclic graph (DAG)**

- 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**

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

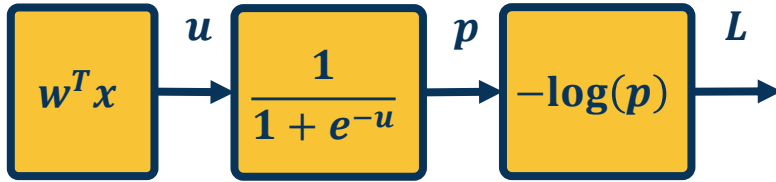This is called reverse-mode **automatic differentiation**

# 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

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

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

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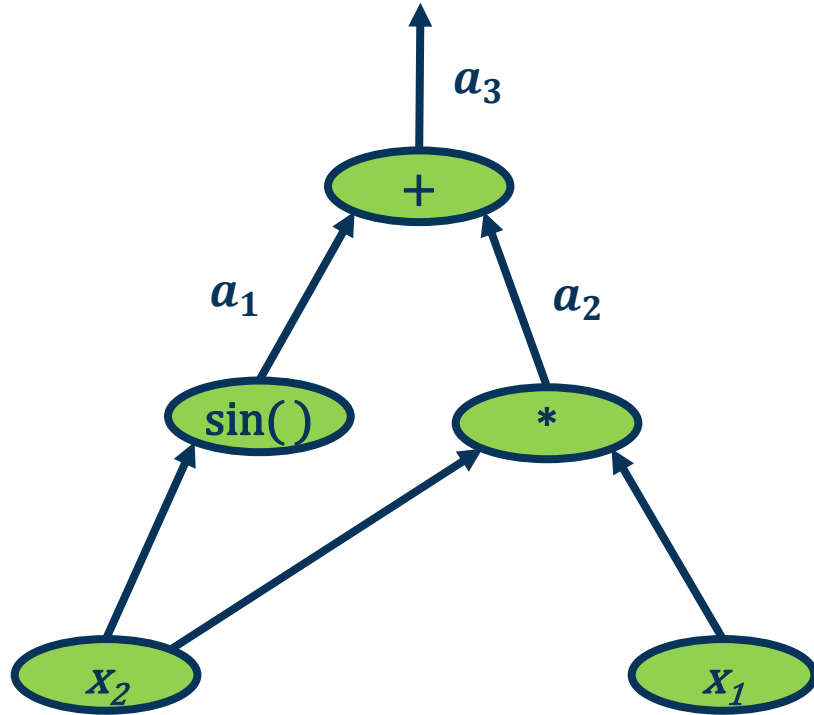
$$\bar{w} = \frac{\partial L}{\partial w} = \frac{\partial L}{\partial u}\frac{\partial u}{\partial w} = \bar{u}x^T$$

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

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

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

# Example Gradient Computations

Georgia Tech

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

$a_3$

$+$

$a_1$  $a_2$

$\sin()$  $*$

$x_2$  $x_1$

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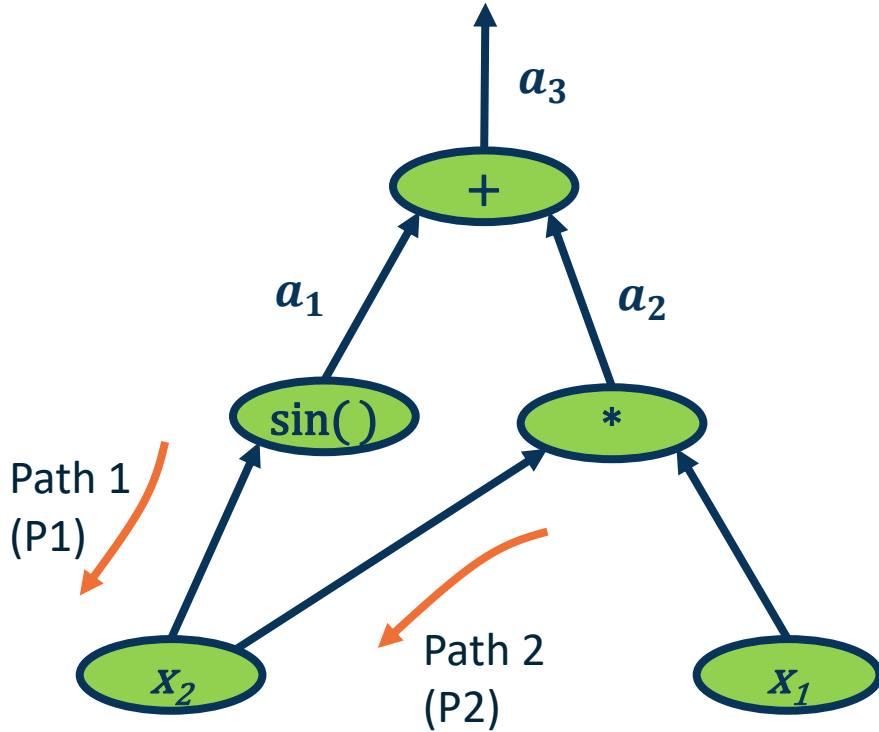
◆ Assign intermediate variables

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

◆ Start at **end** and move **backward**

Georgia Tech

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

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

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

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

$$\overline{x_2^{P1}} = \frac{\partial f}{\partial a_1} \frac{\partial a_1}{\partial x_2} = \overline{a_1} \; \cos(x_2)$$
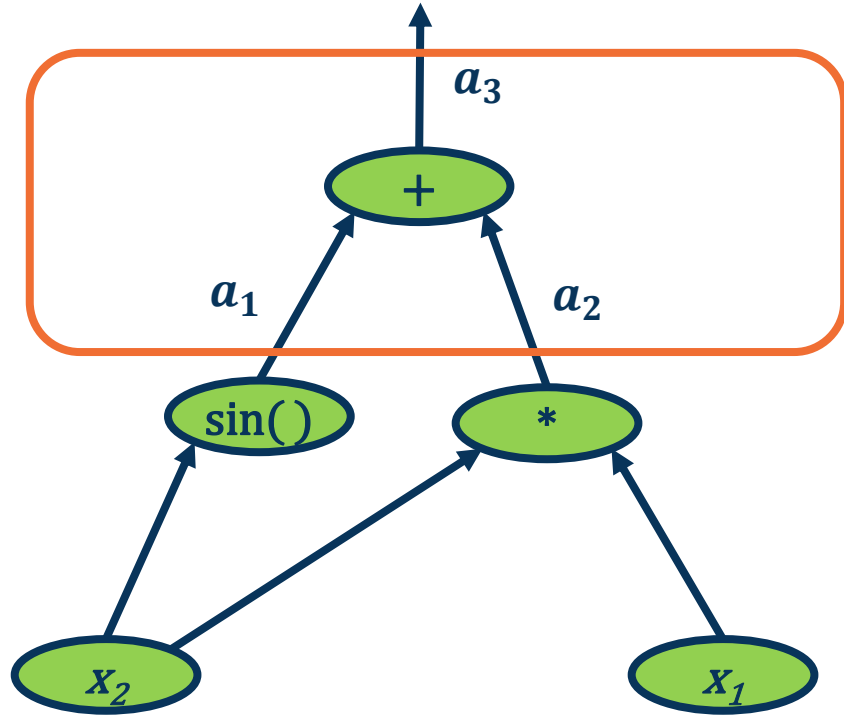
$$\overline{x_2^{P2}} = \frac{\partial f}{\partial a_2} \frac{\partial a_2}{\partial x_2} = \frac{\partial f}{\partial a_2} \frac{\partial(x1x2)}{\partial x_2} = \overline{a_2} x_1$$

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**

Georgia Tech

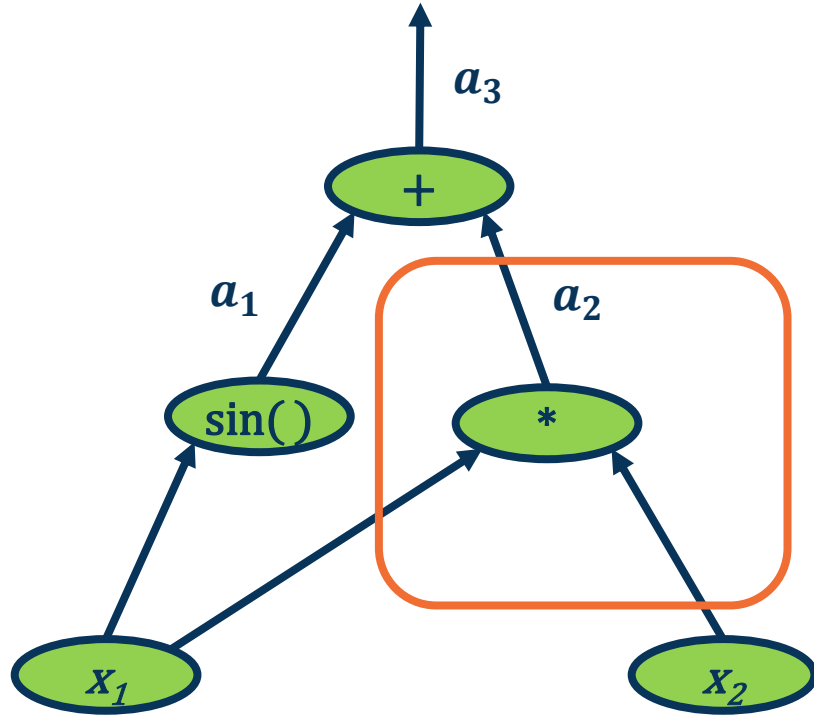$$f(x_1, x_2) = x_1 x_2 + \sin(x_2)$$



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

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

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

**Patterns of Gradient Flow: Addition**

Georgia Tech

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



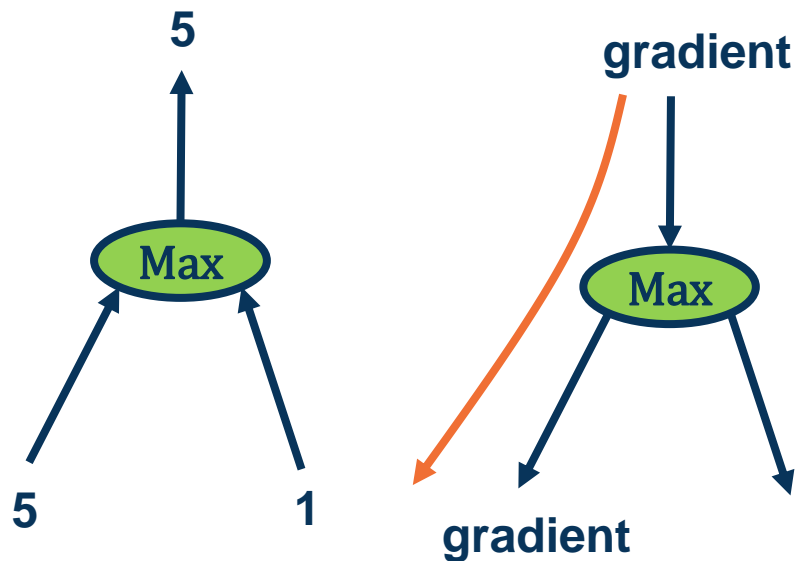Multiplication operation is a gradient switcher (multiplies it by the values of the other term)

$$\overline{x_2} = \frac{\partial f}{\partial a_2} \frac{\partial a_2}{\partial x_2} = \frac{\partial f}{\partial a_2} \frac{\partial (x1x2)}{\partial x_2} = \overline{a_2} x_1$$

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

**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

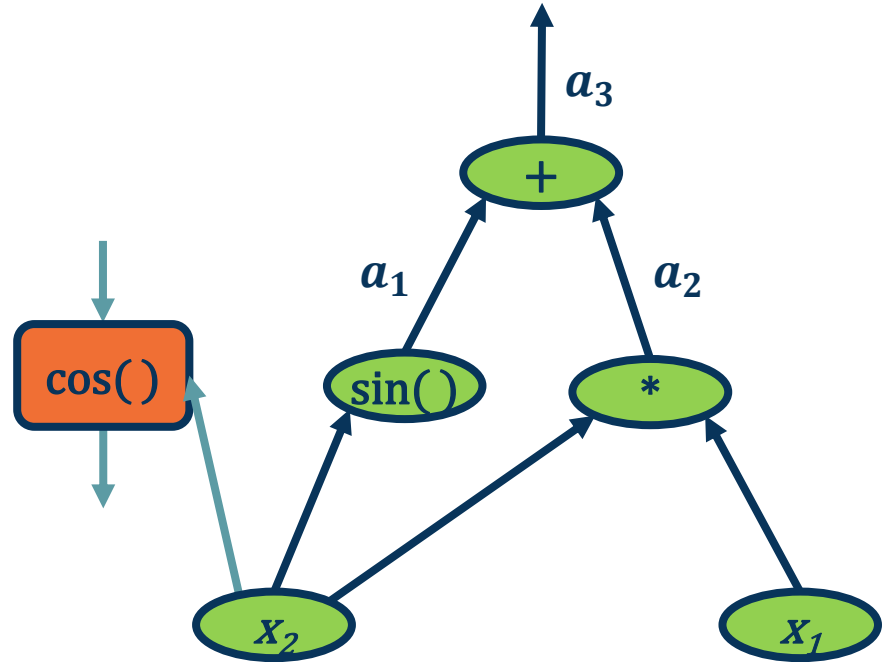◆ This information must be recorded in the forward pass



**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!

Georgia Tech

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**

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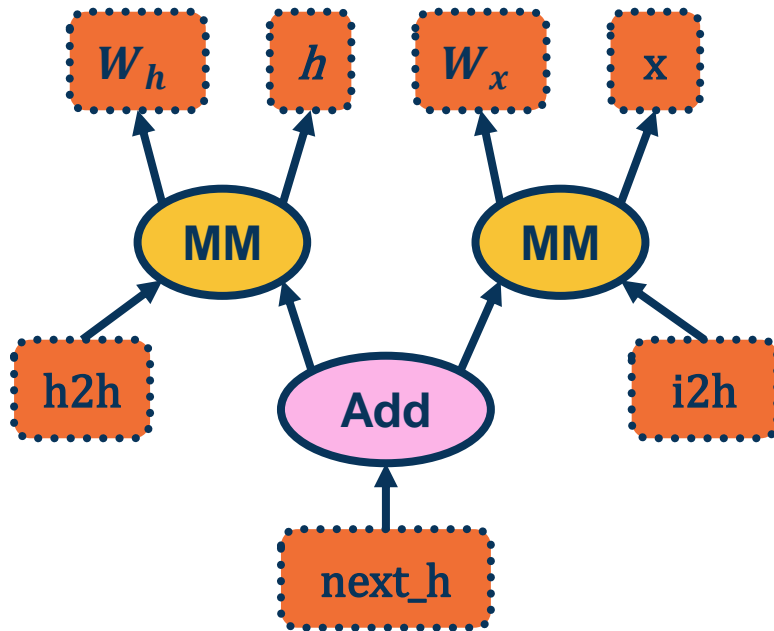
# A graph is created on the fly

```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
```
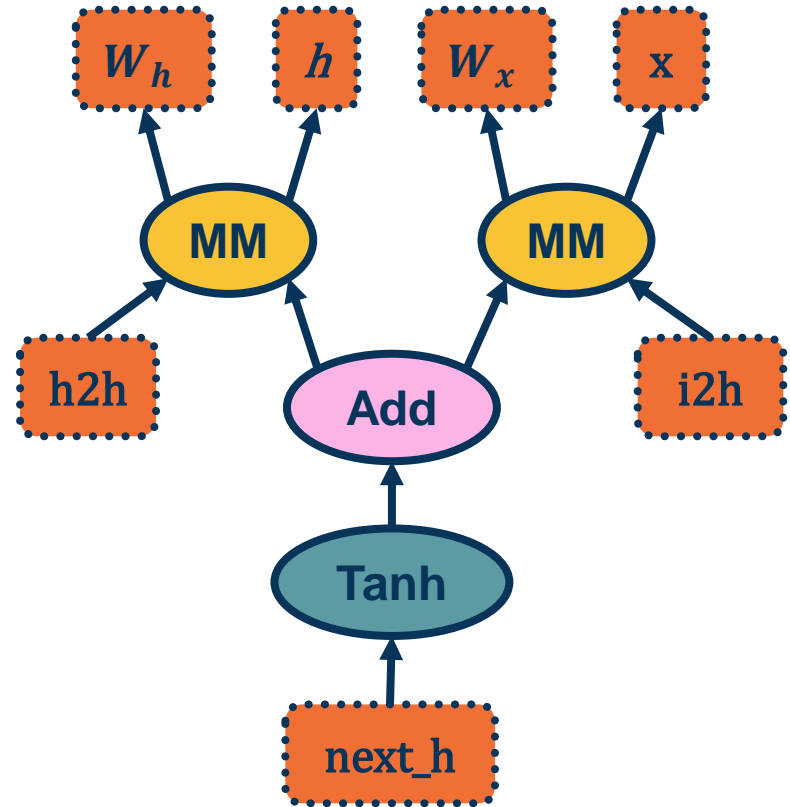
(Note above)

# 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))
```
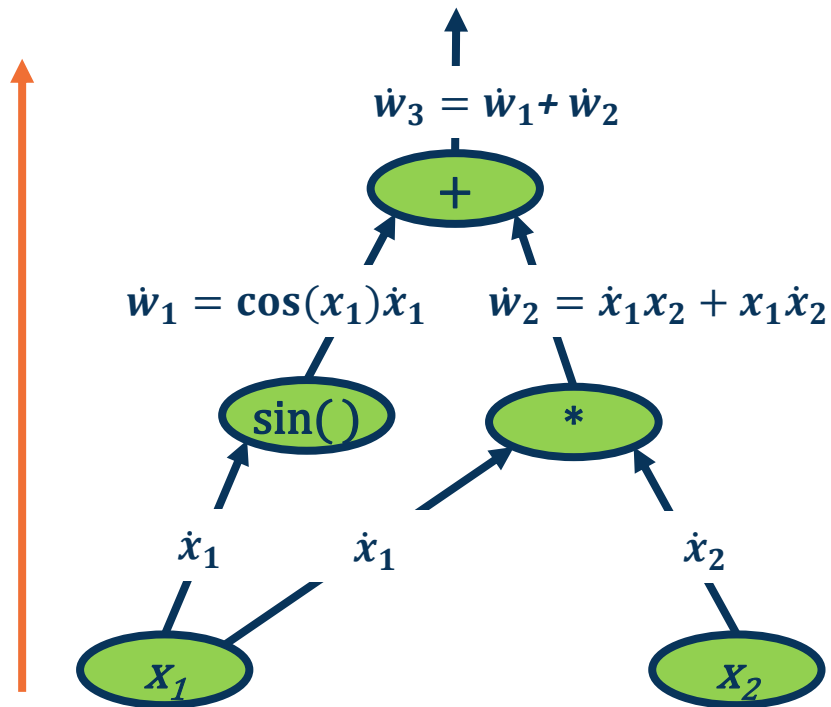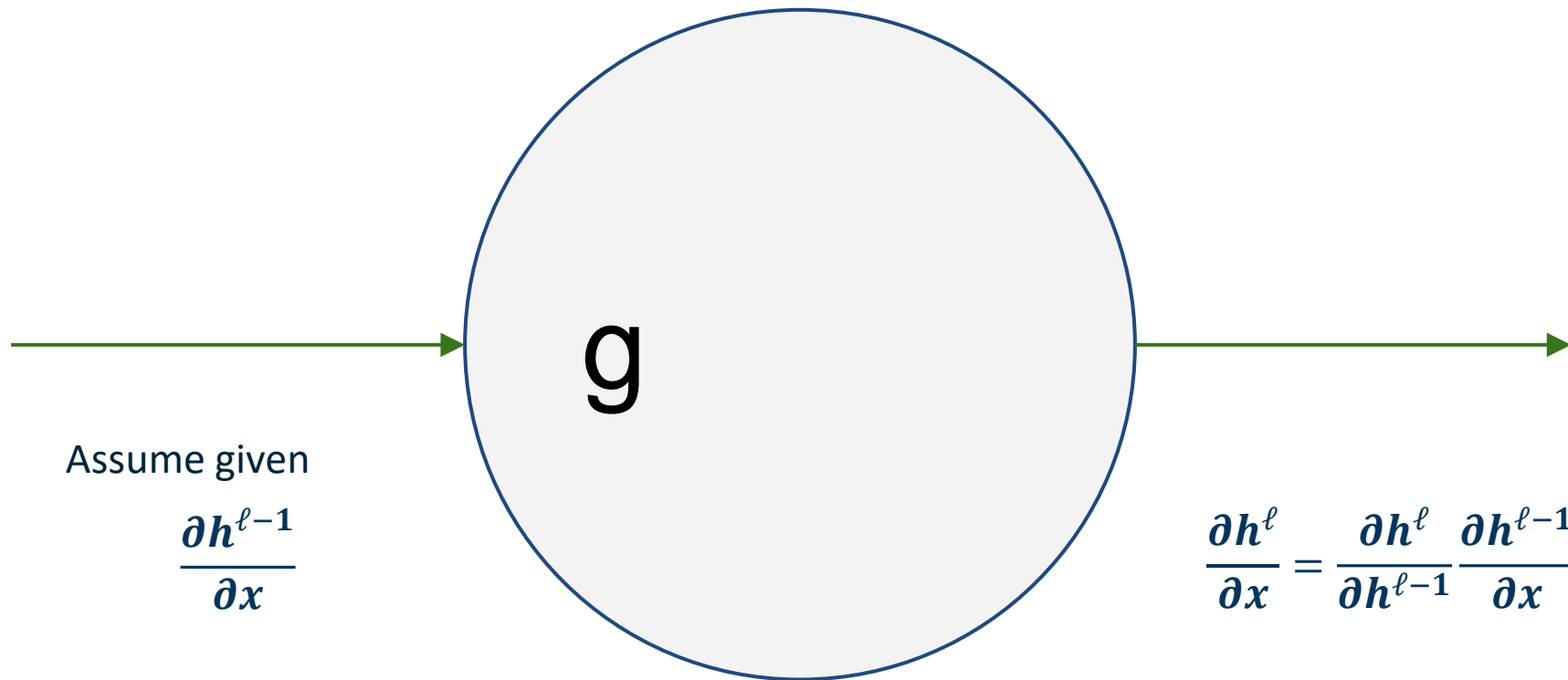


*From pytorch.org*

Georgia Tech

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

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

$+$

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

$\sin()$      $*$

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

$X_1$      $X_2$

g

Assume given

$$\frac{\partial h^{\ell-1}}{\partial x}$$

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

See https://www.cc.gatech.edu/classes/AY2020/cs7643_spring/slides/autodiff_forward_reverse.pdf

**Forward Mode Autodifferentiation**

Georgia Tech

# Convolutional network (AlexNet)

input image

weights

loss

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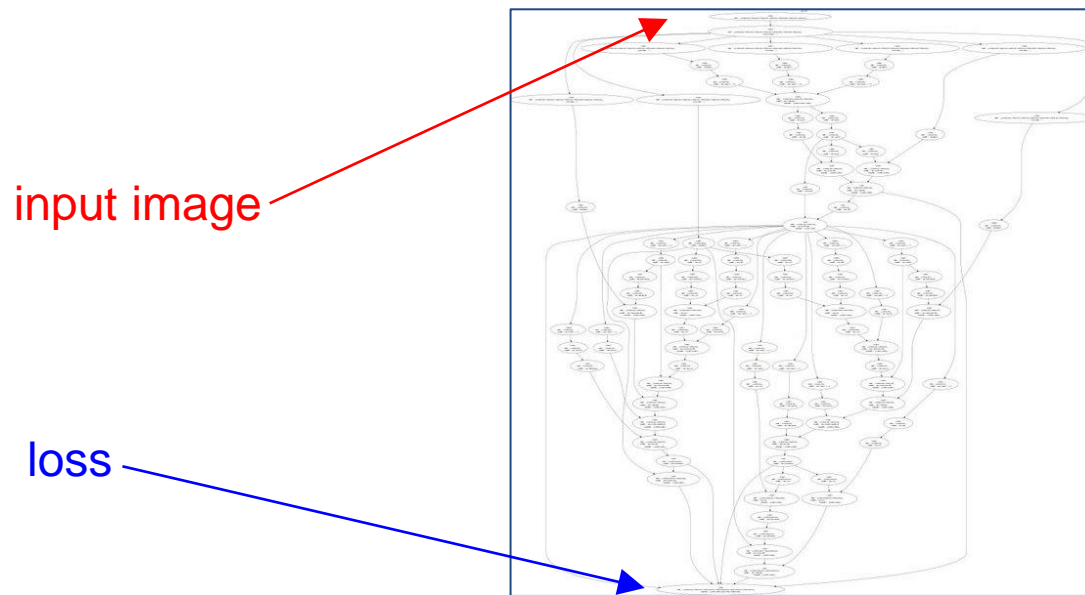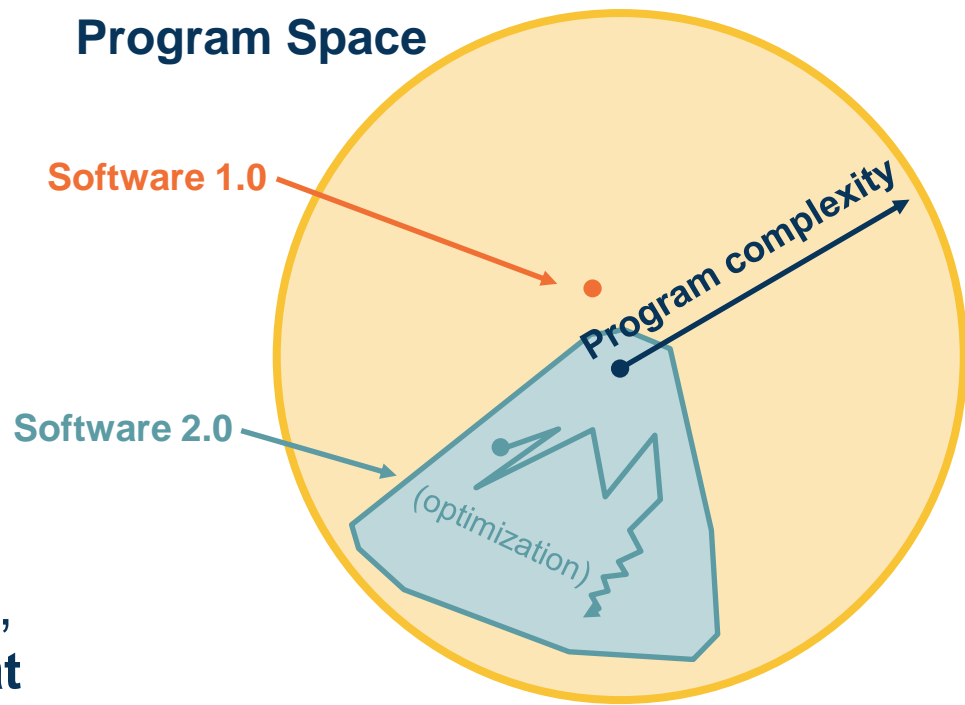
Georgia Tech

# 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**
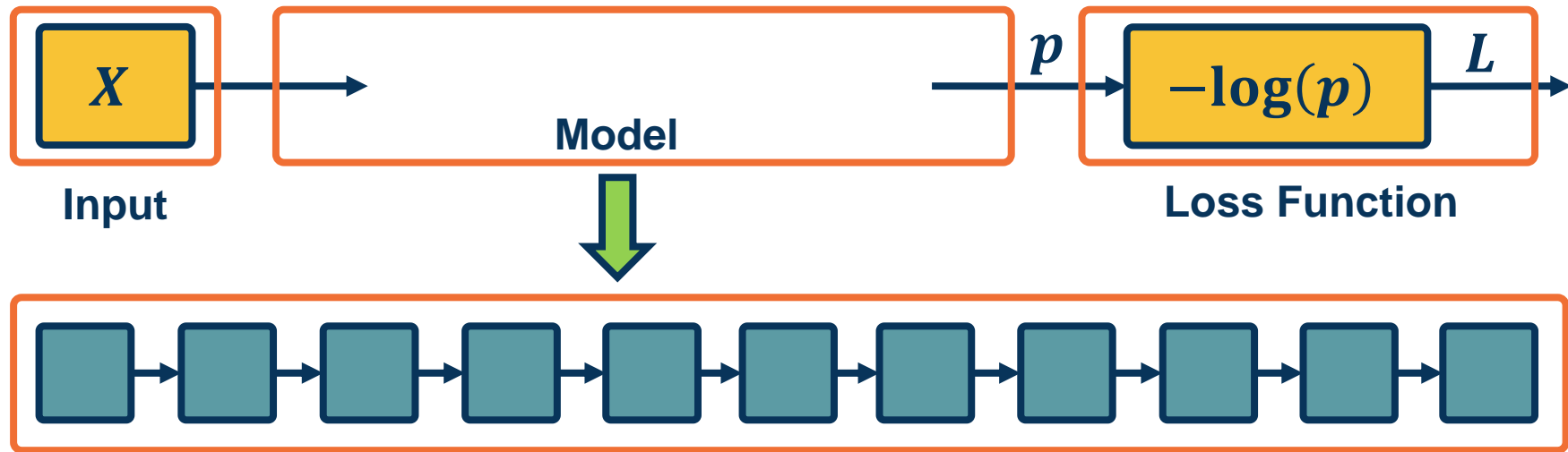
- **Differentiable programming**

**Program Space**

Software 1.0

Program complexity

Software 2.0

(optimization)

*Adapted from figure by Andrej Karpathy*

**Power of Automatic Differentiation**

Georgia Tech

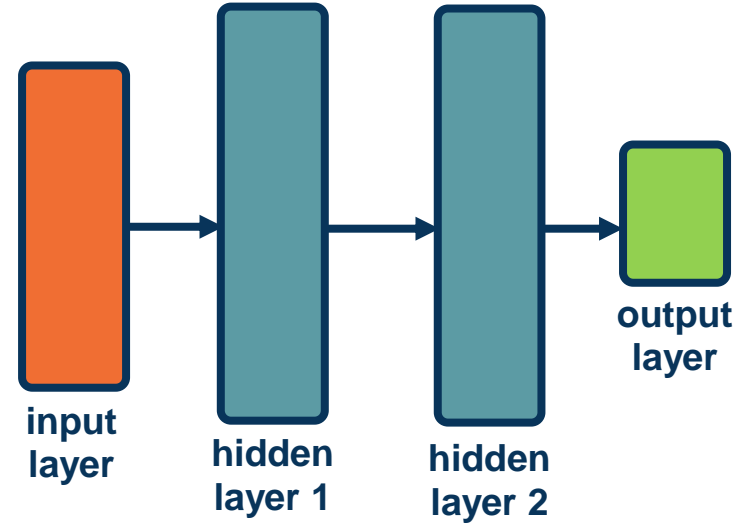Backpropagation, and automatic differentiation, allows us to optimize **any** function composed of differentiable blocks

- **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
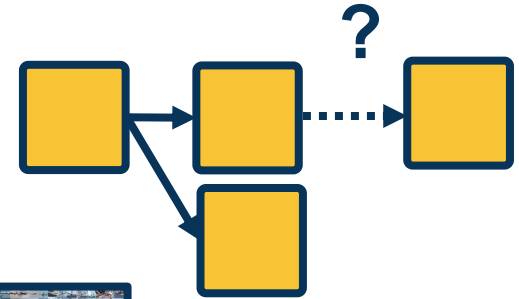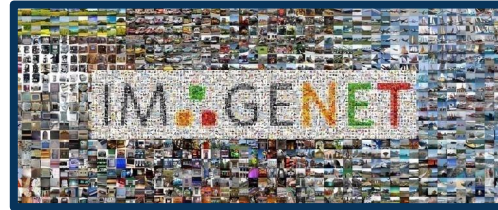
**Depth is important:**

⬡ Structure the model to represent an inherently compositional world

⬡ Theoretical evidence that it leads to parameter efficiency
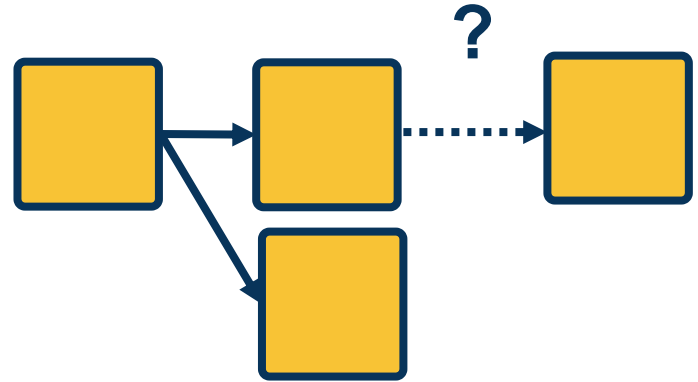
⬡ Gentle dimensionality reduction (if done right)



input layer

hidden layer 1

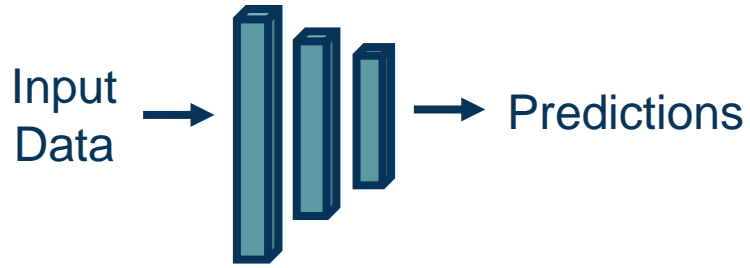hidden layer 2

output layer

There are still many design decisions that must be made:

◆ **Architecture**

◆ **Data Considerations**

◆ **Training and Optimization**

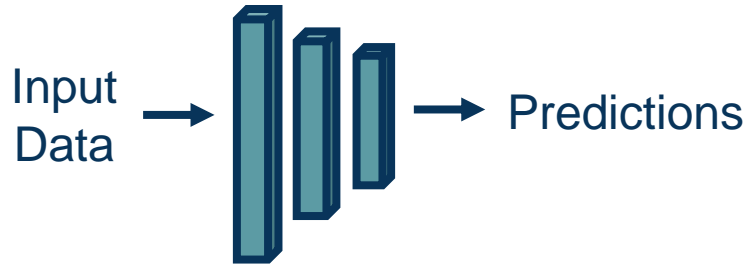◆ **Machine Learning Considerations**

**Local Minima**

Georgia Tech

We must design the **neural network architecture:**

- What **modules (layers)** should we use?

- How should they **be connected together**?

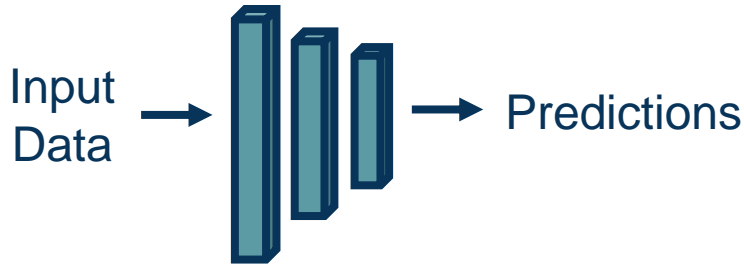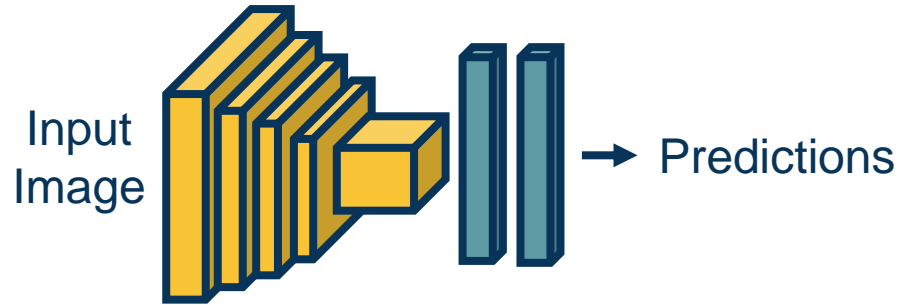- Can we use our **domain knowledge** to add architectural biases?

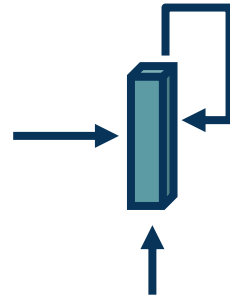Input Data → Predictions

**Fully Connected Neural Network**

Input Data → Predictions

**Fully Connected Neural Network**

Input Image → Predictions

**Convolutional Neural Networks**

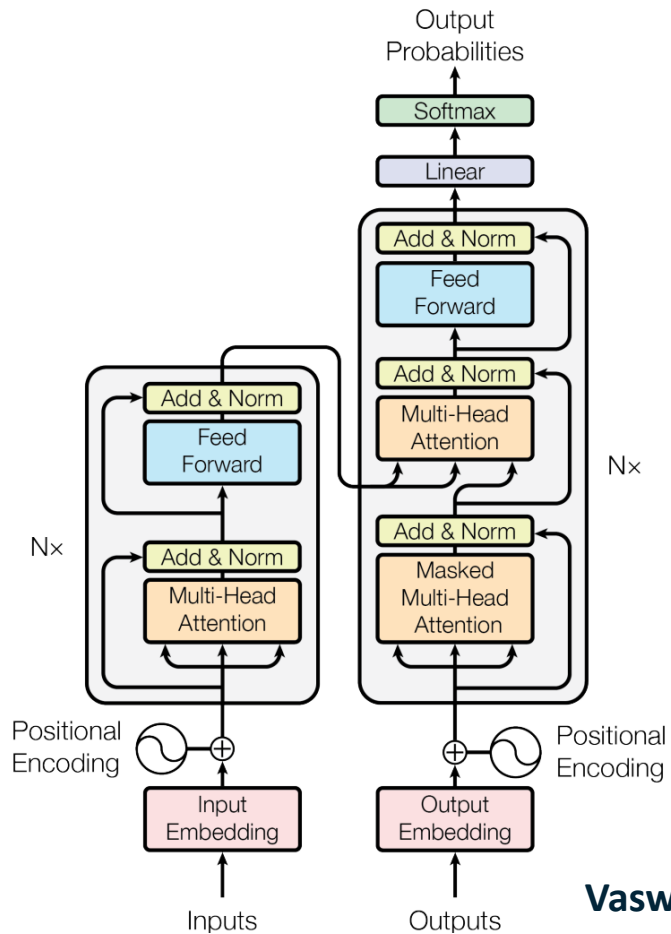Input Data → **Fully Connected Neural Network** → Predictions

Input Image → **Convolutional Neural Networks** → Predictions

**Recurrent Neural Network**

Before: Different architectures are suitable for different applications or types of input

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Add & Norm

Masked Multi-Head Attention

N×

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs

**Now: Transformers for all input modalities!**

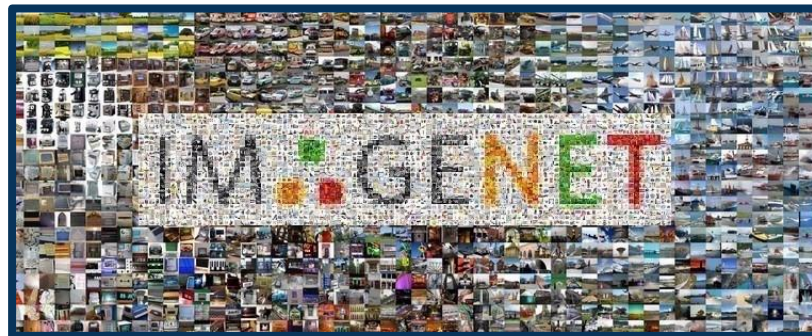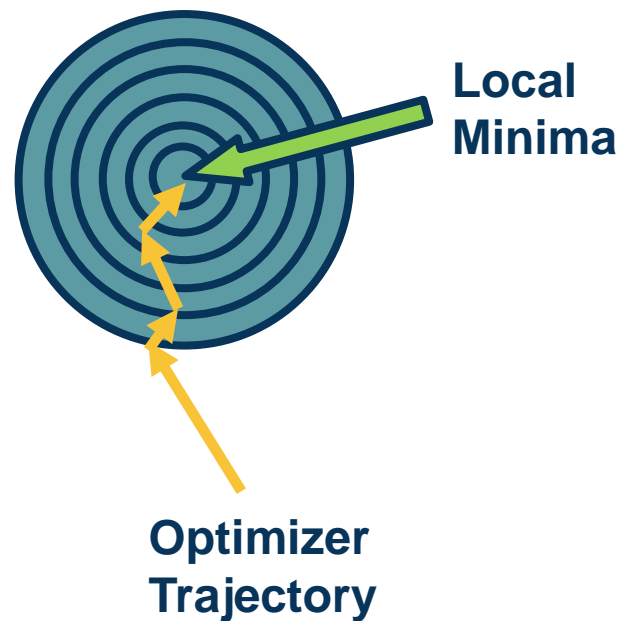**Vaswani et al., Attention Is All You Need**

As in traditional machine learning, **data** is key:

⬡ Should we **pre-process** the data?

⬡ Should we **normalize** it?

⬡ Can we **augment** our data by adding noise or other perturbations?

Even given a good neural network architecture, we need a **good optimization algorithm to find good weights**

- What **optimizer** should we use?

  - Different optimizers make **different weight updates** depending on the gradients

- How should we **initialize** the weights?

- What **regularizers** should we use?

- What **loss function** is appropriate?



**Local Minima**

**Optimizer Trajectory**

Georgia Tech

## Machine Learning Considerations

**The practice of machine learning is complex:** For your particular application you have to **trade off** all of the considerations together

- Trade-off between **model capacity** (e.g. measured by # of parameters) and **amount of data**

- Adding **appropriate biases** based on knowledge of the domain