

# CS 4803 / 7643: Deep Learning

## Topics:

- (Finish) Automatic Differentiation
  - Patterns in backprop
  - Jacobians in FC+ReLU NNs

Dhruv Batra  
Georgia Tech

# Administrativa

- HW2 out
  - Due: 09/23 11:59pm
  - Theory: Gradient descent, Hessians, Auto-diff, Convolutions
    - <https://www.overleaf.com/project/5f4c143e06061d00013dd4f0>
  - Implementation: ConvNets in Python and PyTorch
  - Bonus: Challenge on EvalAI
    - <https://evalai.cloudcv.org/web/challenges/challenge-page/684/overview>

# Project

- Goal
  - Chance to take on something open-ended
  - Encouraged to apply to your research (computer vision, NLP, robotics,...)
- Main categories
  - **Reproducibility**
    - Pick a paper from a recent conference. Attempt to reproduce the method and validate claims.
  - **Application/Survey**
    - Compare a collection of existing algorithms on a new application domain of your interest.
  - **Formulation/Development**
    - Formulate a new model or algorithm for a new or old problem
  - **Theory**
    - Theoretically analyze an existing algorithm

# Project

- Rules

- **Combine with other classes / research / credits / anything**
- You have our blanket permission
- Get permission from other instructors; delineate different parts
- Must be done this semester.
- Groups of 3-4

- Expectations

- 20% of final grade = individual effort equivalent to 1 HW
- Expectation scales with team size
- Most work will be done in Nov but please plan early.

# Project Ideas

- ML Reproducibility Challenge 2020
  - <https://paperswithcode.com/rc2020>

Browse State-of-the-Art   Methods   Trends   About   RC2020

▸ RC2020

[Registration](#)

[Task](#)

[Resources](#)

## ML Reproducibility Challenge 2020

Welcome to the ML Reproducibility Challenge 2020! This is already the fourth edition of this event (see [V1](#), [V2](#), [V3](#)), and we are excited this year to announce that we are broadening our coverage of conferences and papers to cover several new top venues, including: [NeurIPS](#), [ICML](#), [ICLR](#), [ACL](#), [EMNLP](#), [CVPR](#) and [ECCV](#).

The primary goal of this event is to encourage the publishing and sharing of scientific results that are reliable and reproducible. In support of this, the objective of this challenge is to investigate reproducibility of papers accepted for publication at top conferences by inviting members of the community at large to select a paper, and verify the empirical results and claims in the paper by reproducing the computational experiments, either via a new implementation or using code/data or other information provided by the authors.

All submitted reports will be peer reviewed and shown next to the original papers on [Papers with Code](#). Reports will be peer-reviewed via [OpenReview](#). Every year, a small number of these reports, selected for their clarity, thoroughness, correctness and insights, are selected for publication in a special edition of the journal [ReScience](#). (see [J1](#), [J2](#)).

# TAs



Sameer Dharur



Joanne Truong



Yihao Chen



Michael Pisen



Hrishikesh Kale



Tianyu Zhan



Prabhav Chawla



Guillermo Nicolas Grande

# Computing

- Major bottleneck
  - GPUs
- Options
  - Your own / group / advisor's resources
  - Google Cloud Credits
    - \$50 credits to every registered student courtesy Google
  - Google Colab
    - jupyter-notebook + free GPU instance

# Administrativa

- Project Teams

- [https://gtvault-my.sharepoint.com/:x:/g/personal/dbatra8\\_gatech\\_edu/EY4\\_65XOzWtOkXSSz2WgpoUBY8ux2gY9PsRzR6KnglIFEQ?e=4tnKWI](https://gtvault-my.sharepoint.com/:x:/g/personal/dbatra8_gatech_edu/EY4_65XOzWtOkXSSz2WgpoUBY8ux2gY9PsRzR6KnglIFEQ?e=4tnKWI)
- Project Title
- 1-3 sentence project summary TL;DR
- Team member names

*] change ok.*



# Recap from last time

# Deep Learning = Differentiable Programming

- Computation = Graph
  - Input = Data + Parameters
  - Output = Loss
  - Scheduling = Topological ordering
- Auto-Diff
  - A family of algorithms for implementing chain-rule on computation graphs

# Forward mode AD

Goal:  $\frac{\partial L}{\partial \theta} \left[ \frac{\partial L}{\partial x} \right]$

layer  $l$

$\vec{h}^{(l-1)}$

$\vec{h}^{(l)} = g(\vec{h}^{(l-1)})$

$g(\cdot)$

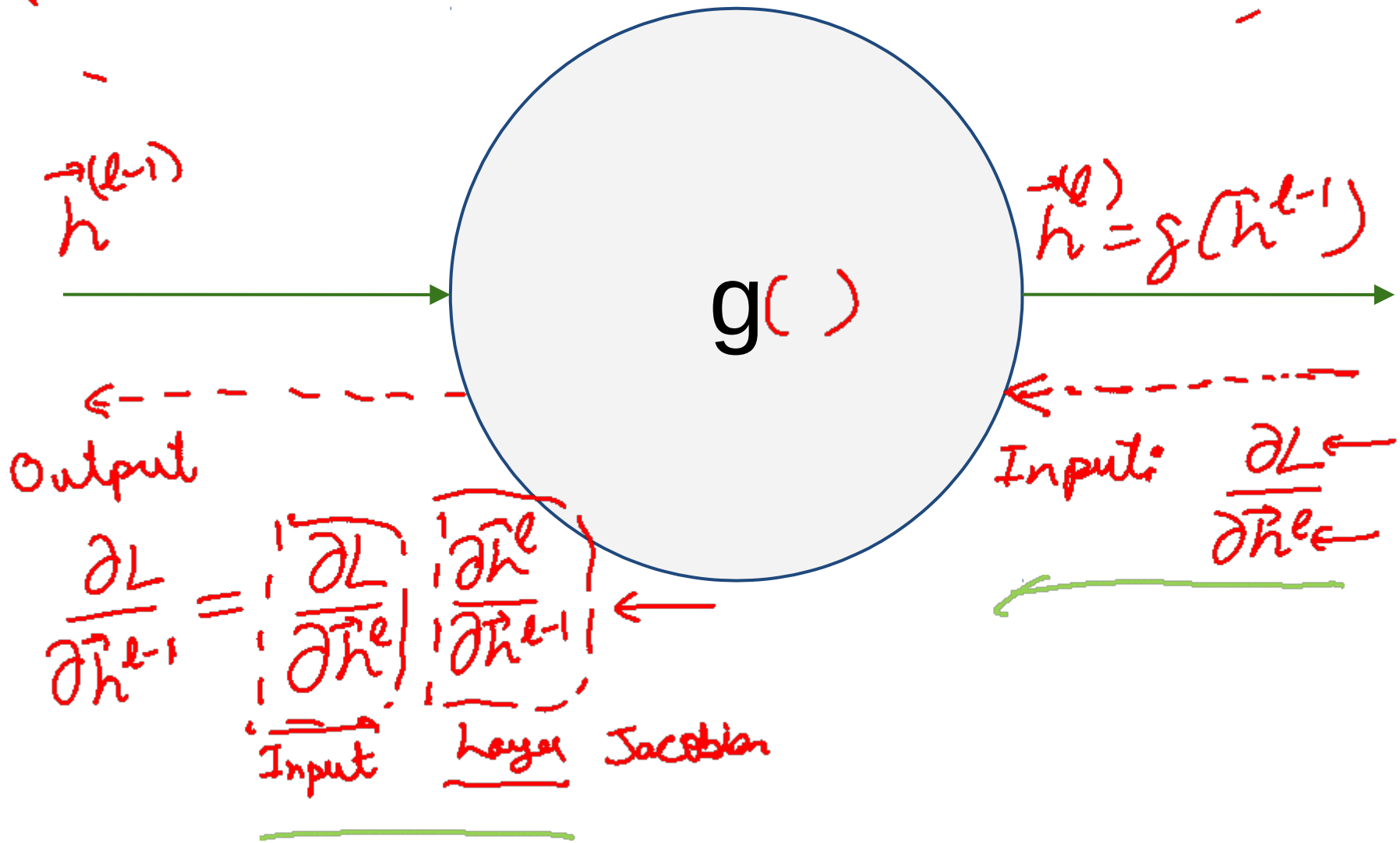
Input:

$$\frac{\partial \vec{h}^{(l-1)}}{\partial \vec{x}}$$

$$\frac{\partial \vec{h}^{(l)}}{\partial \vec{x}} = \frac{\partial \vec{h}^{(l)}}{\partial \vec{h}^{(l-1)}} \frac{\partial \vec{h}^{(l-1)}}{\partial \vec{x}}$$

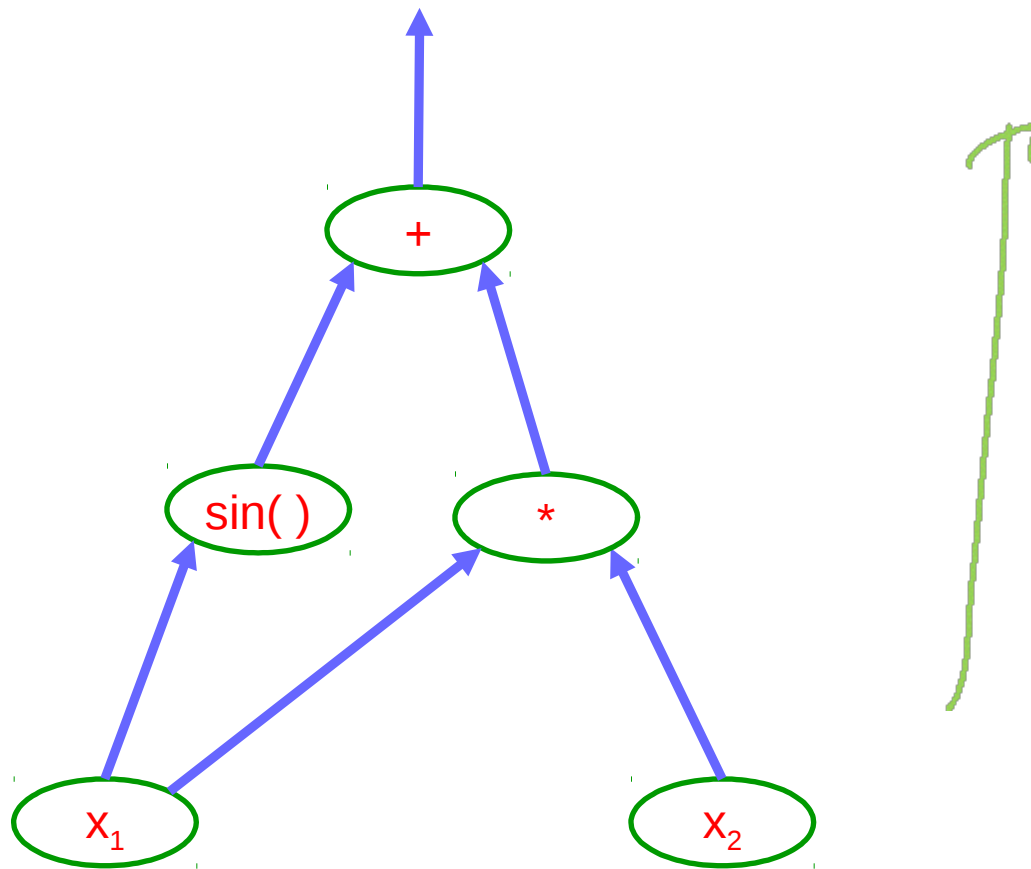
Layer Jacobian      Input FM-A

# Reverse mode AD Goal: $\frac{\partial L}{\partial x}$



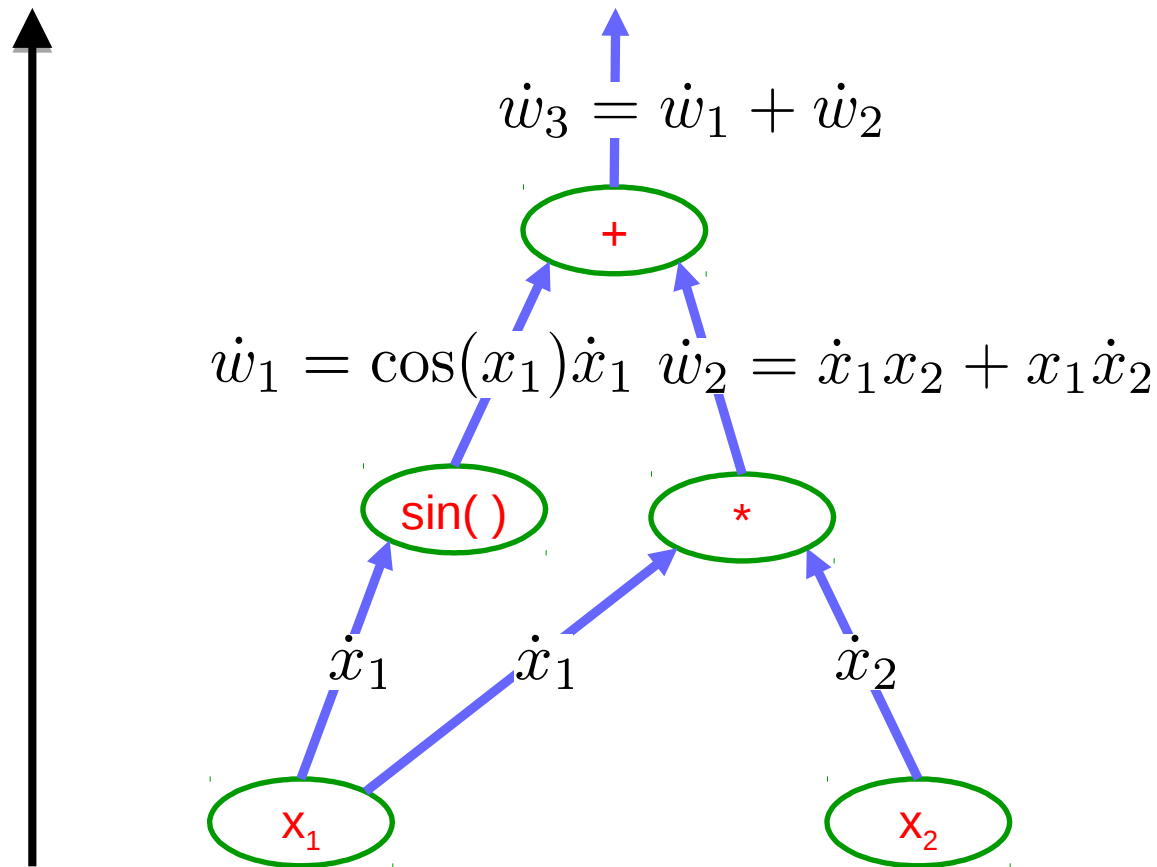
# Example: Forward mode AD

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



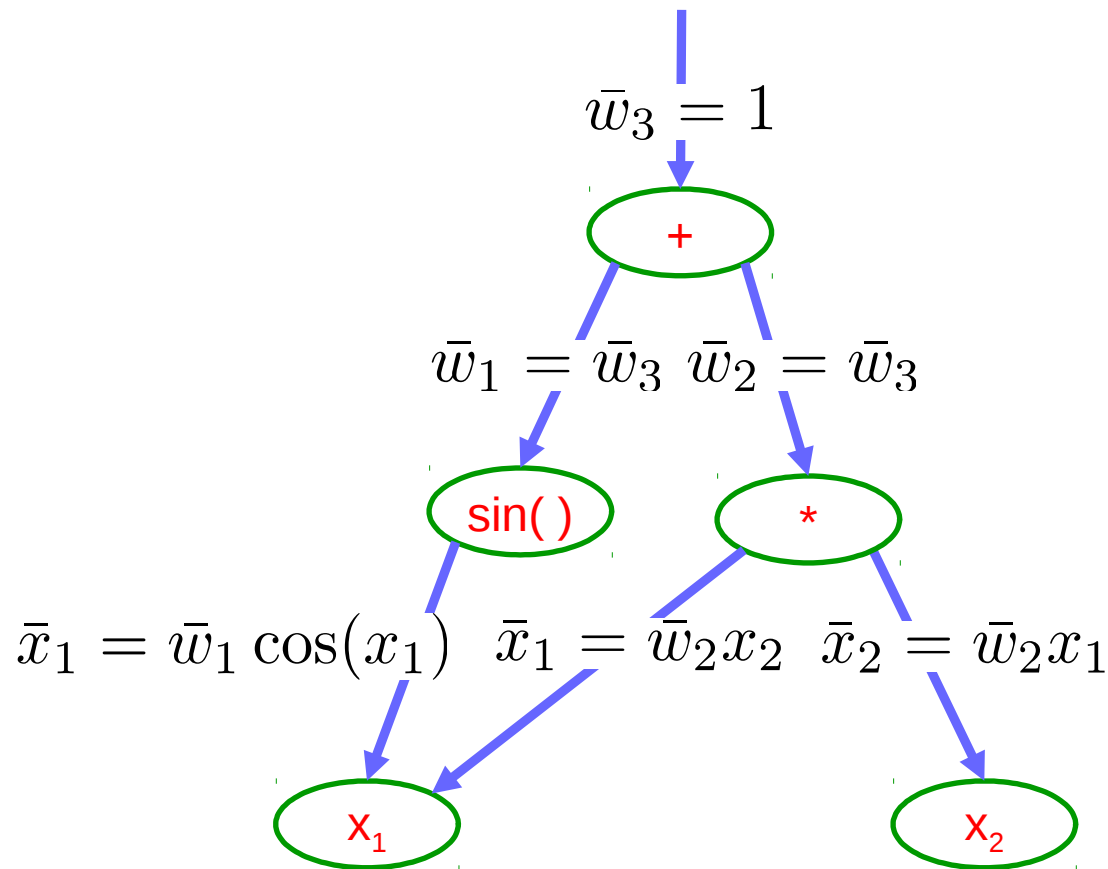
# Example: Forward mode AD

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



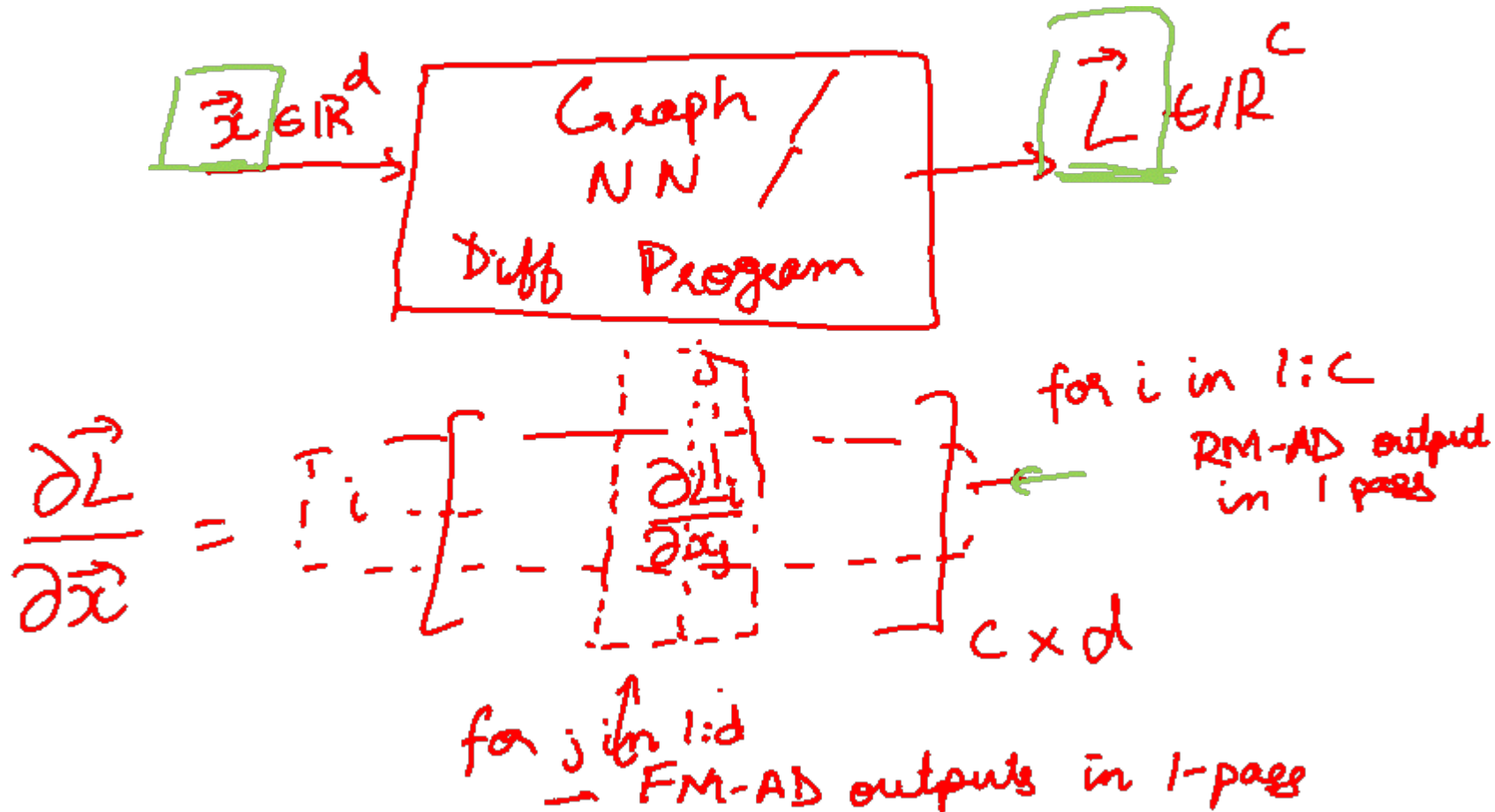
# Example: Reverse mode AD

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



# Forward mode vs Reverse Mode

- x  $\rightarrow$  Graph  $\rightarrow$  L
- Intuition of Jacobian

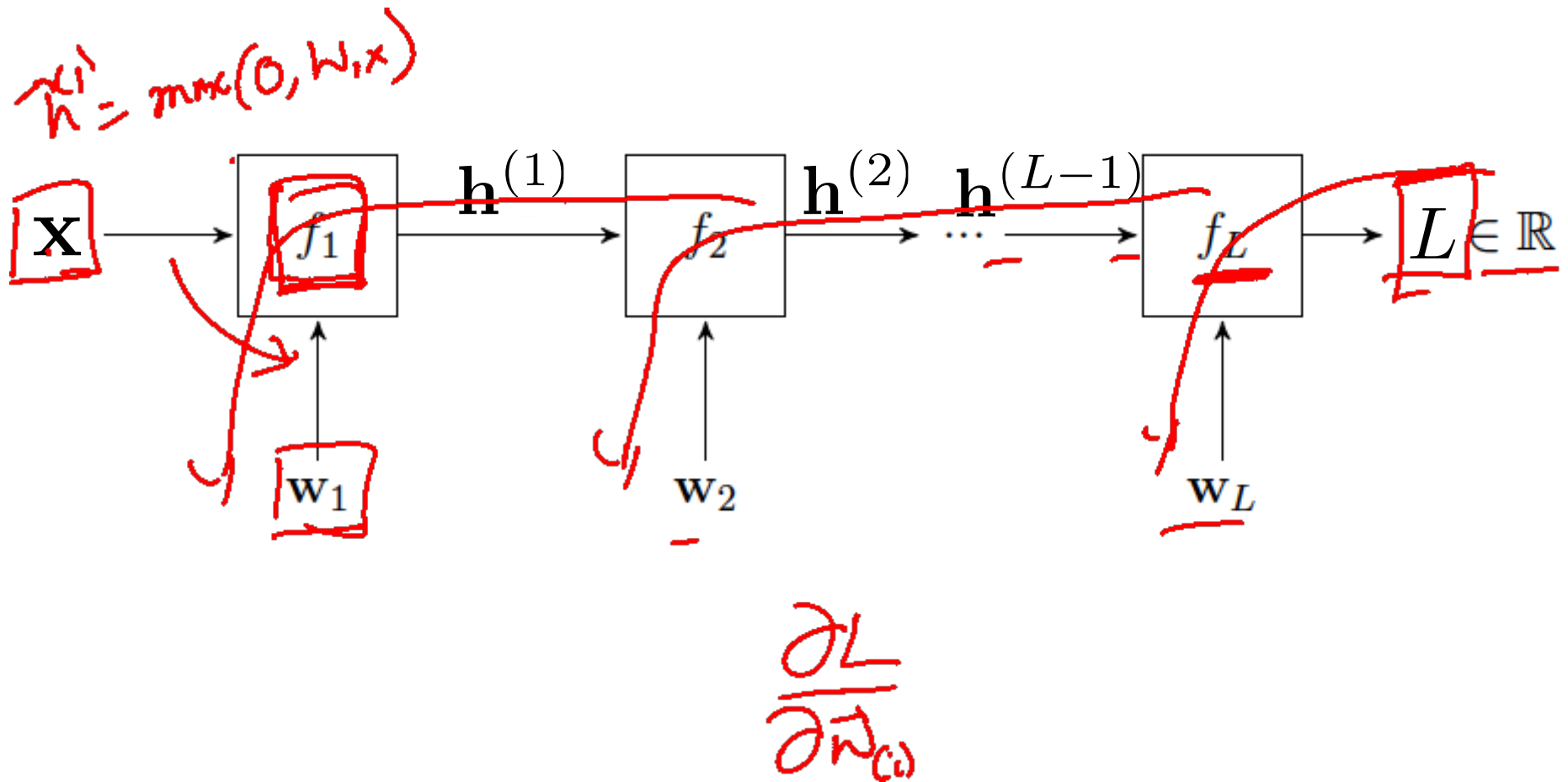




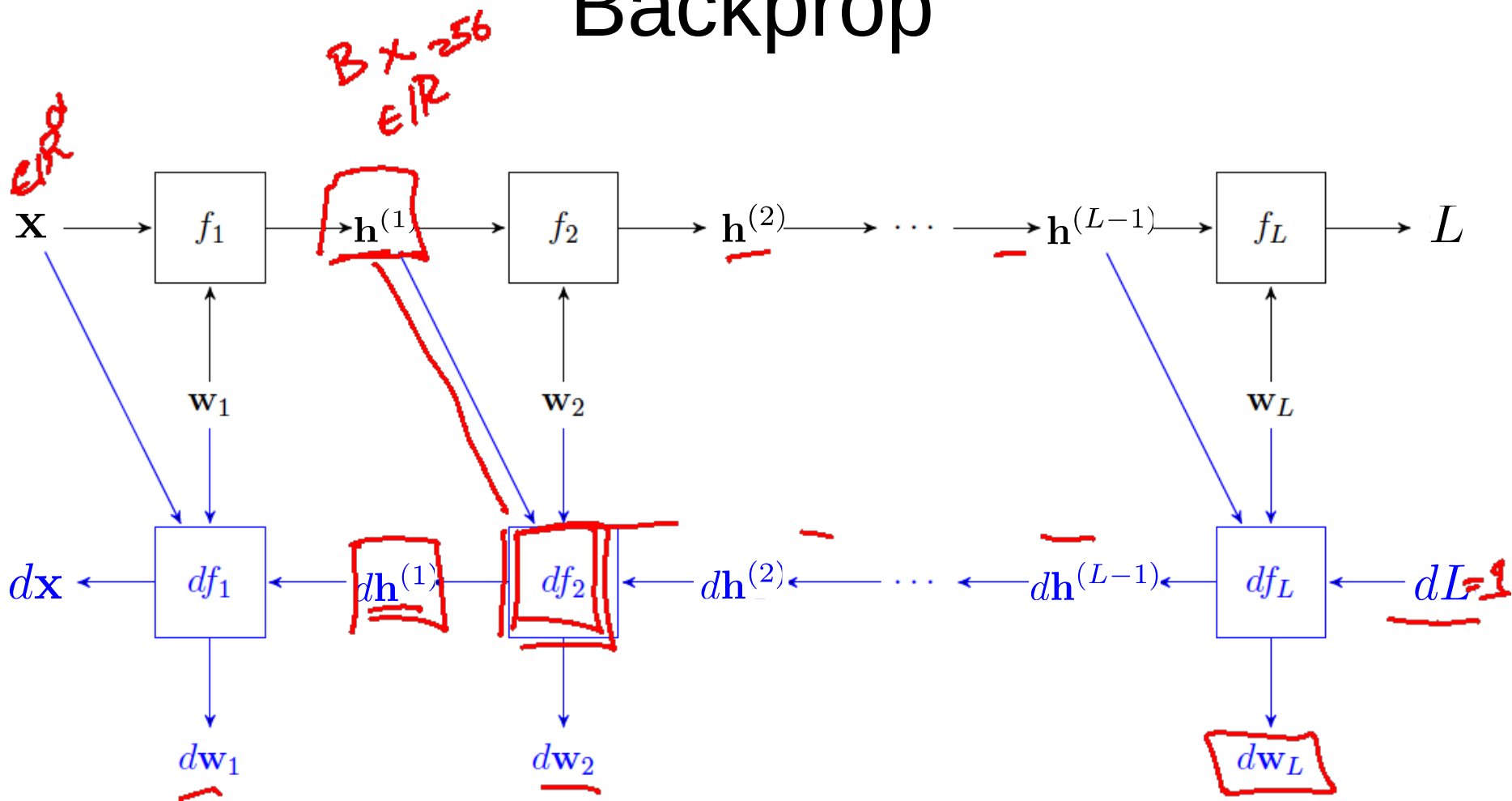
# Forward mode vs Reverse Mode

- What are the differences?
- Which one is faster to compute?
  - Forward or backward?
- Which one is more memory efficient (less storage)?
  - Forward or backward?

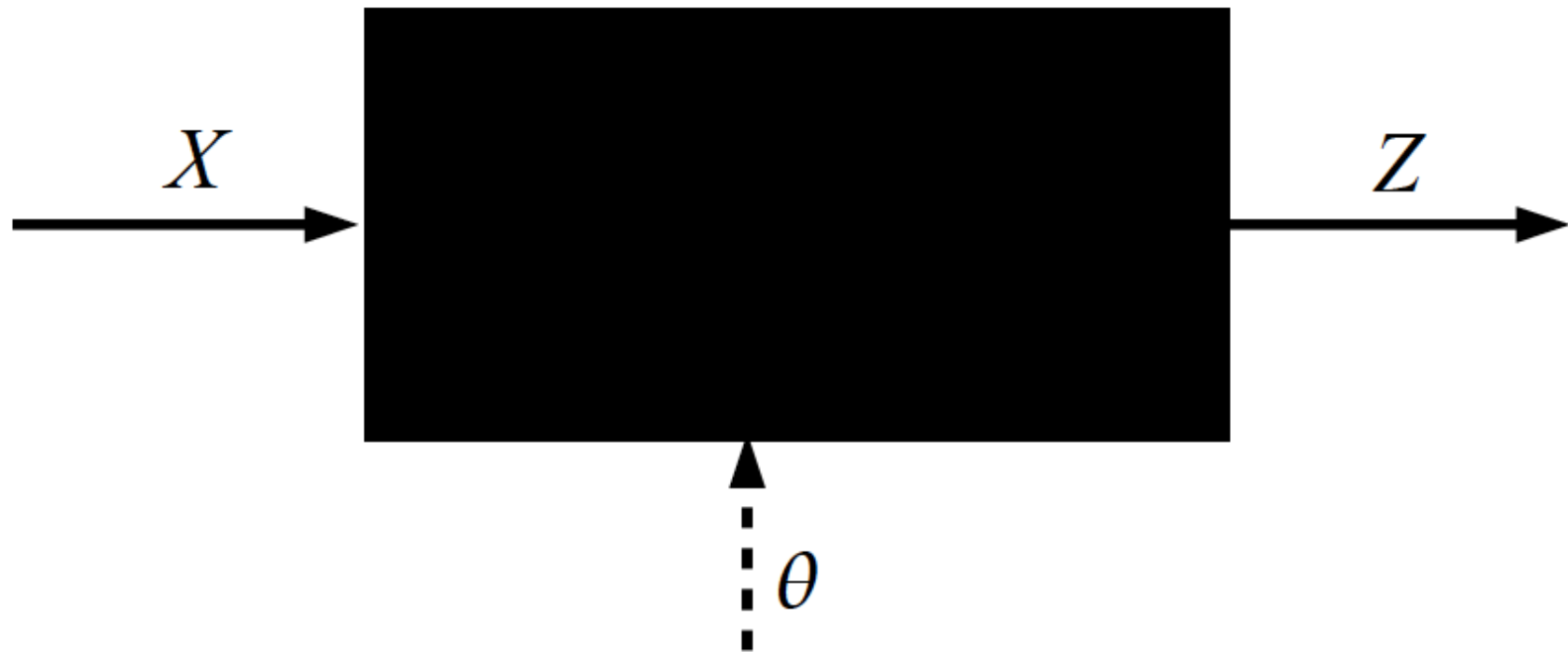
# Neural Network Computation Graph



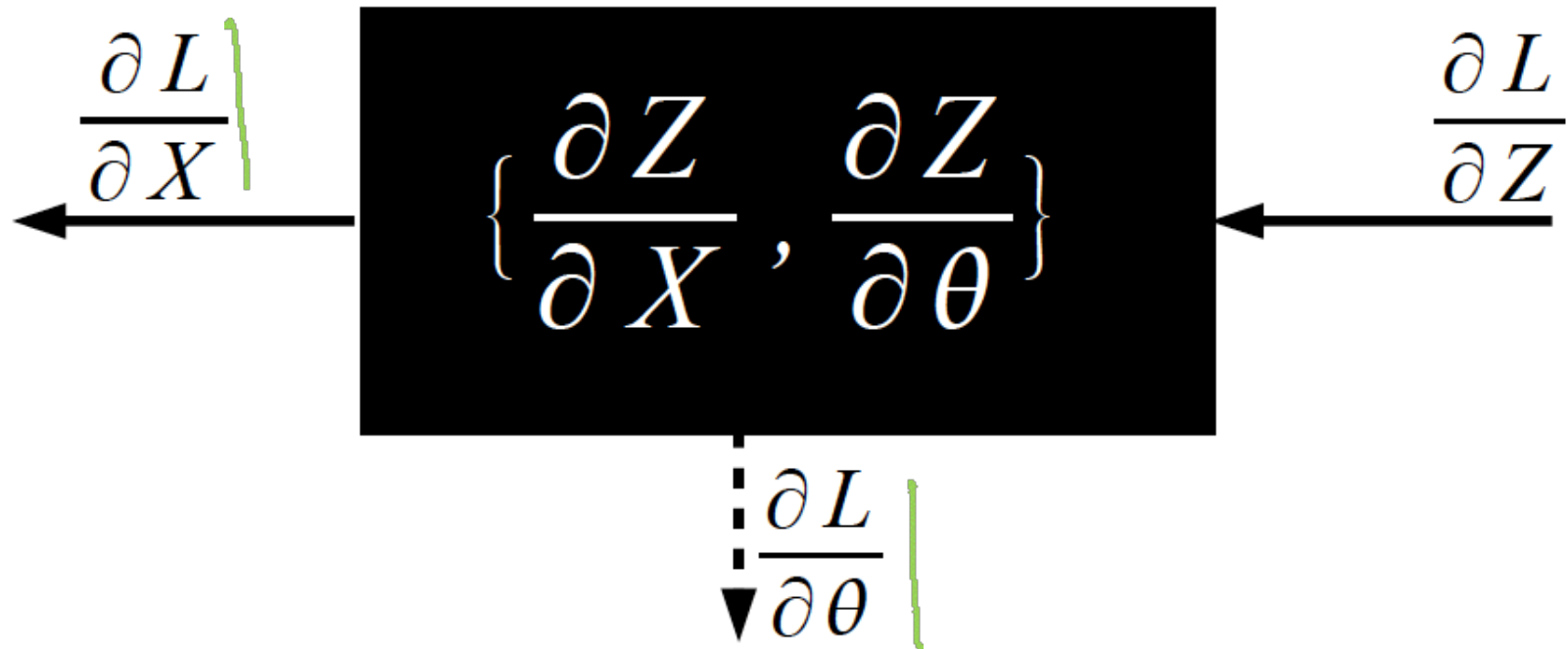
# Backprop



# Key Computation: Forward-Prop

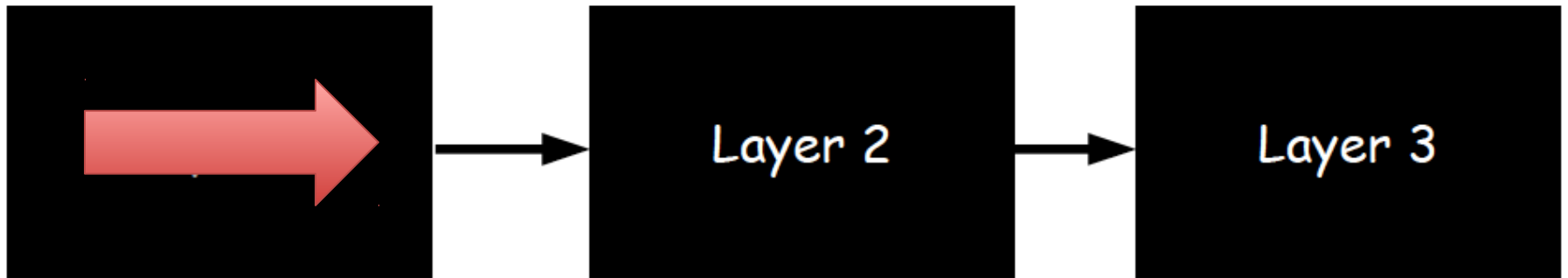


# Key Computation: Back-Prop



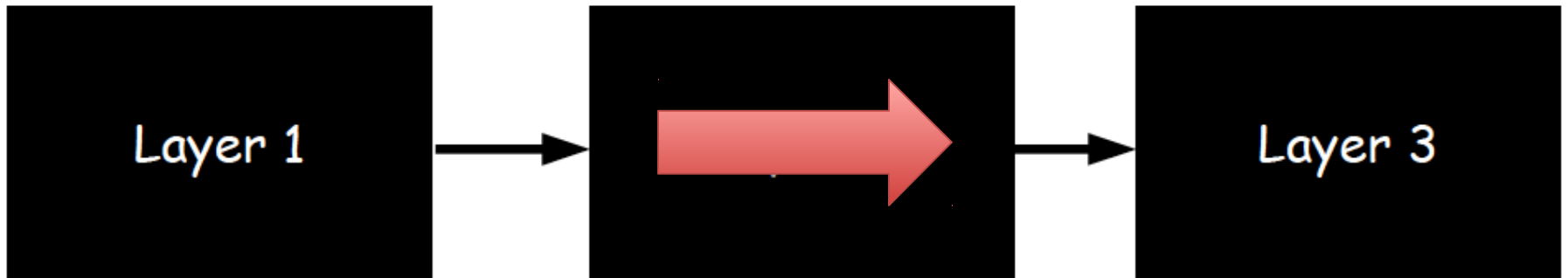
# Neural Network Training

- Step 1: Compute Loss on mini-batch [F-Pass]



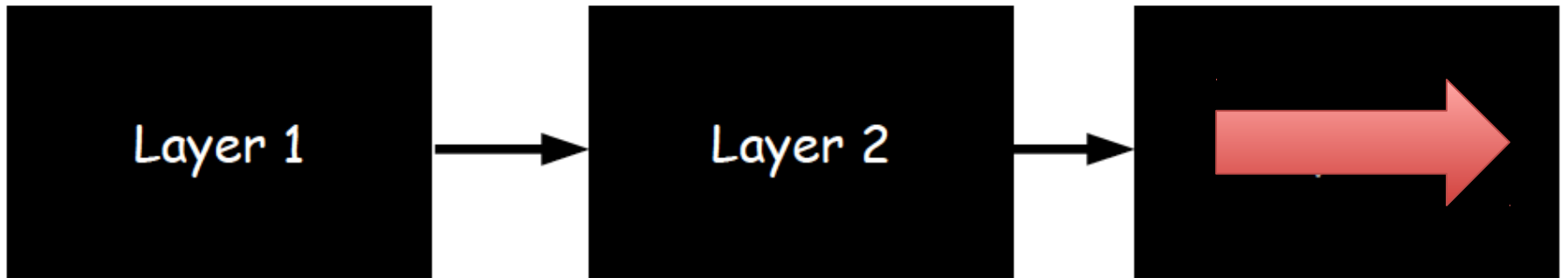
# Neural Network Training

- Step 1: Compute Loss on mini-batch [F-Pass]



# Neural Network Training

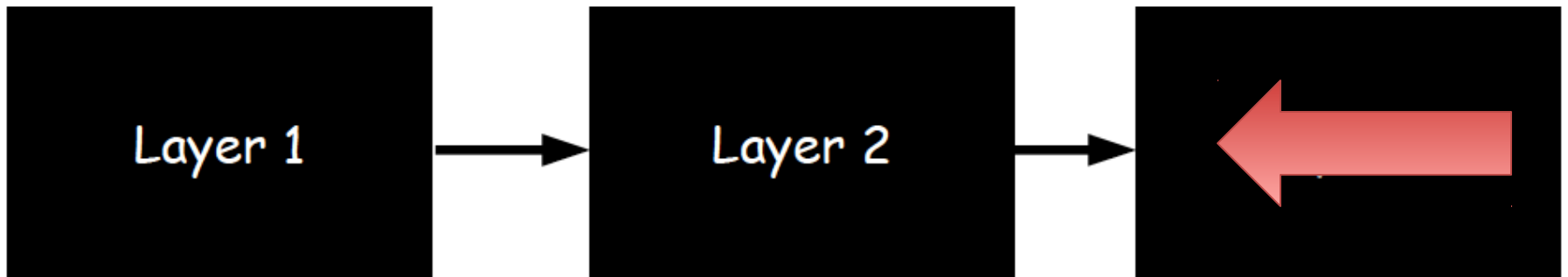
- Step 1: Compute Loss on mini-batch [F-Pass]





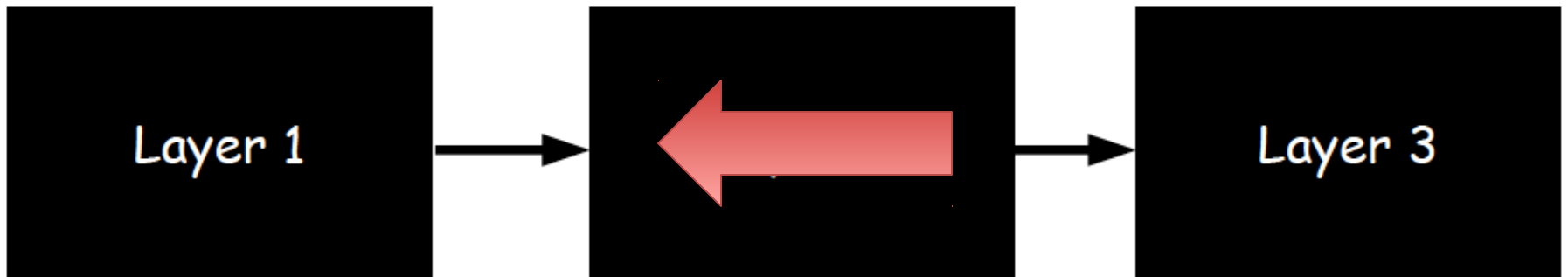
# Neural Network Training

- Step 1: Compute Loss on mini-batch [F-Pass]
- Step 2: Compute gradients wrt parameters [B-Pass]



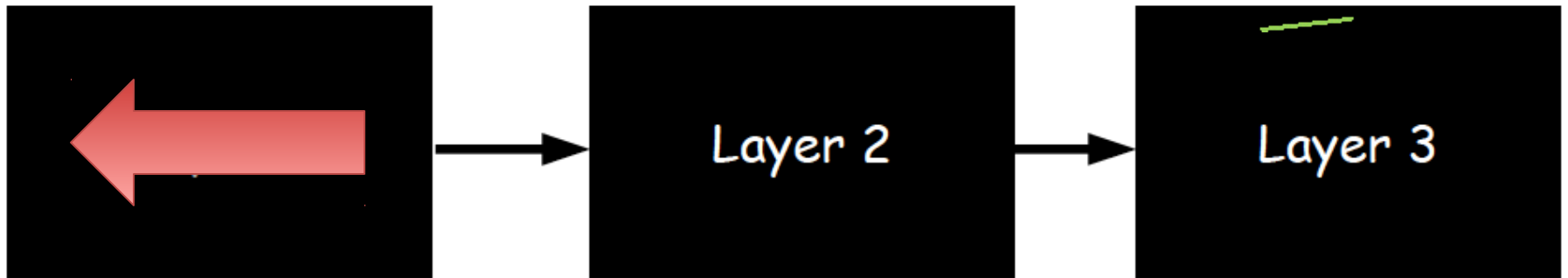
# Neural Network Training

- Step 1: Compute Loss on mini-batch [F-Pass]
- Step 2: Compute gradients wrt parameters [B-Pass]



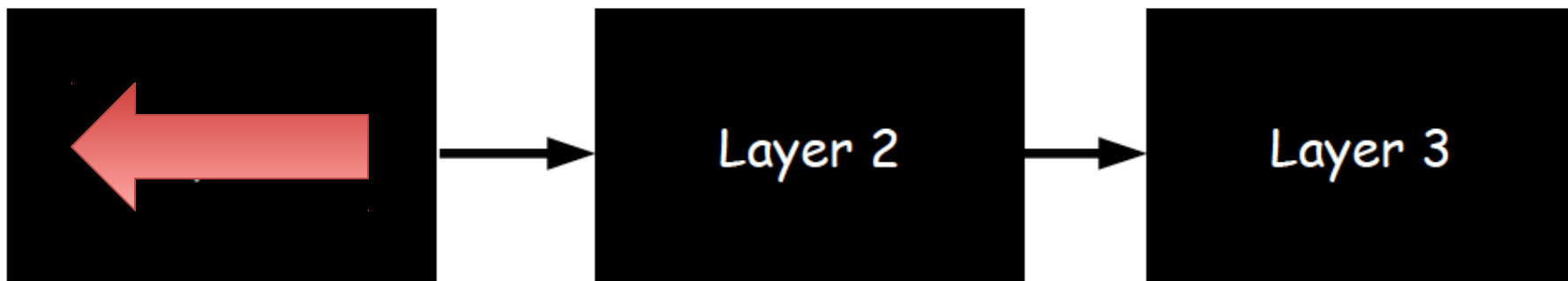
# Neural Network Training

- Step 1: Compute Loss on mini-batch [F-Pass]
- Step 2: Compute gradients wrt parameters [B-Pass]



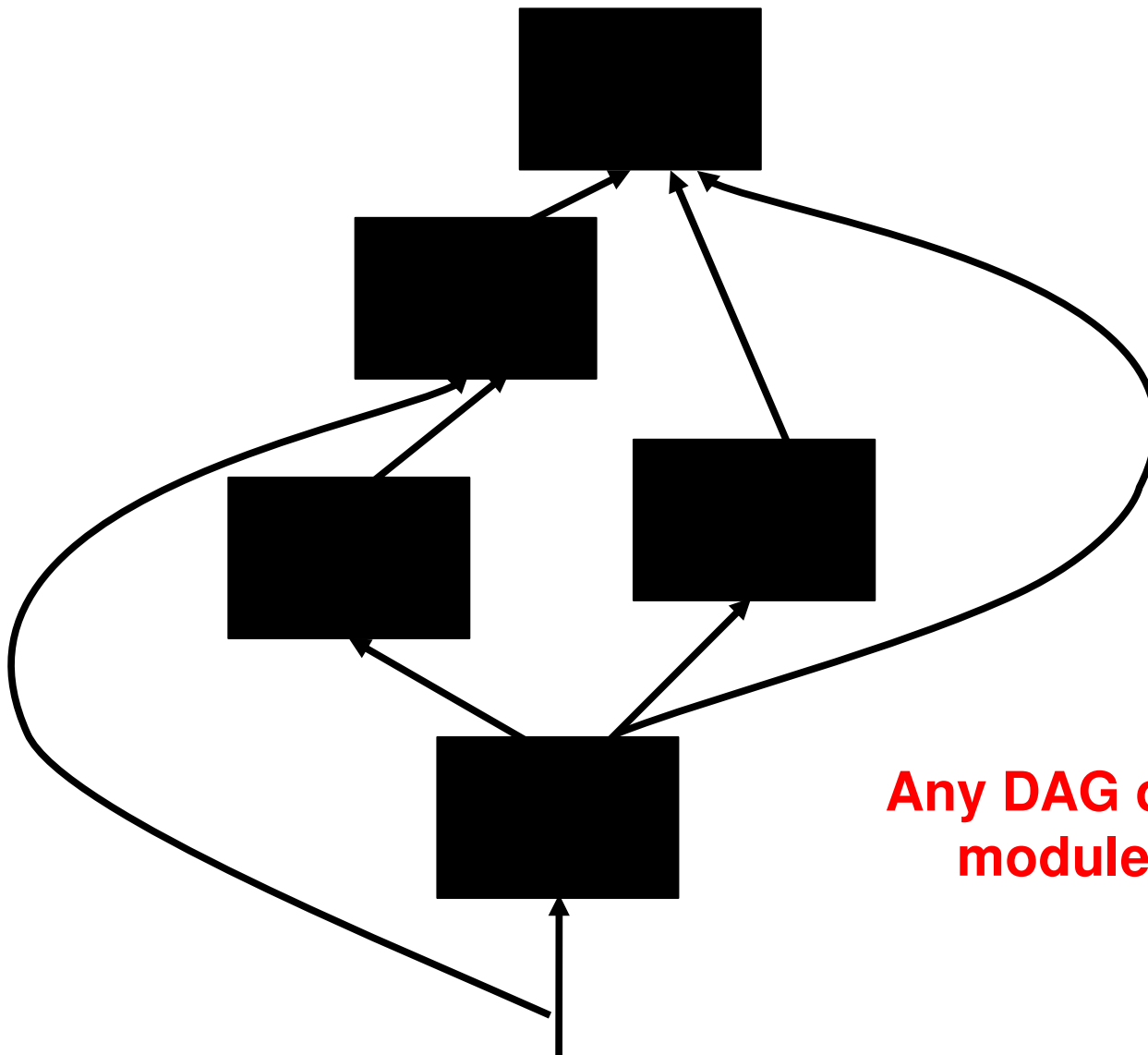
# [Neural Network Training]

- Step 1: Compute Loss on mini-batch [F-Pass] ←
- Step 2: Compute gradients wrt parameters [B-Pass] ←
- Step 3: Use gradient to update parameters



$$\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$$

# Computational Graph



**Any DAG of differentiable modules is allowed!**

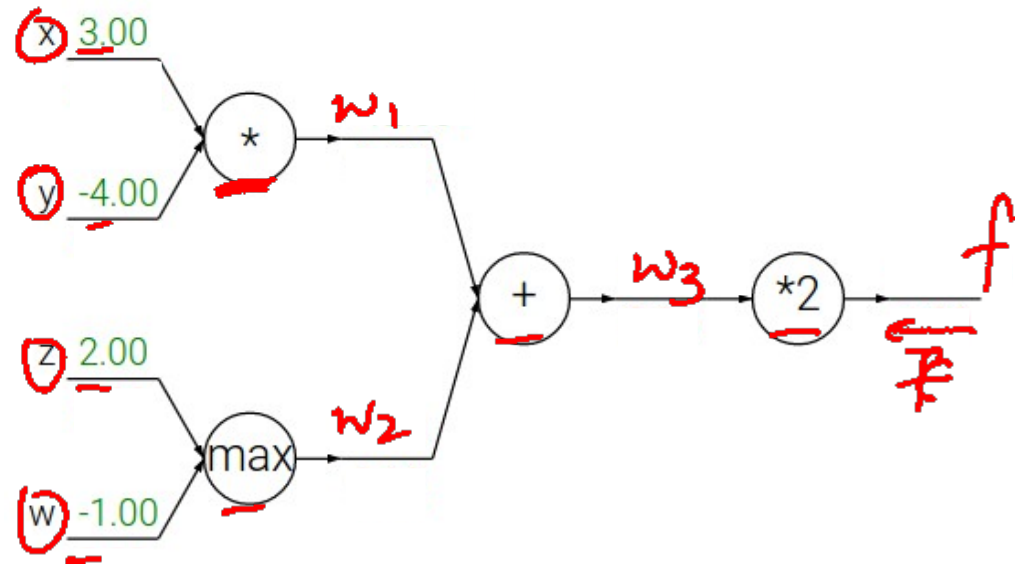
# Plan for Today

- Automatic Differentiation

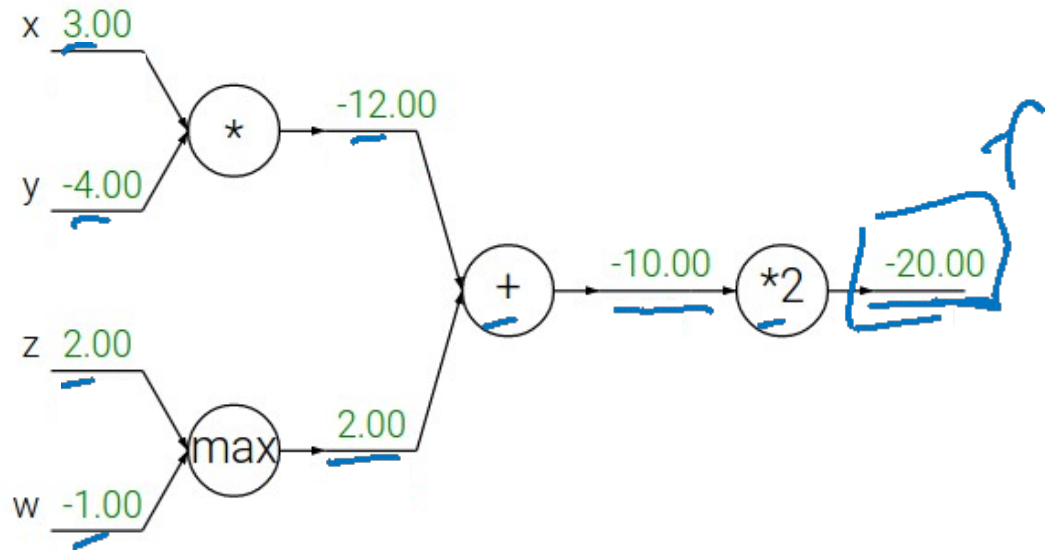
- Patterns in backprop
- Jacobians in FC+ReLU NNs

# Backpropagation: a simple example

$$f(x, y, z, w) = 2 \left( \underbrace{xy}_{w_1} + \underbrace{\max\{z, w\}}_{w_2} \right)$$

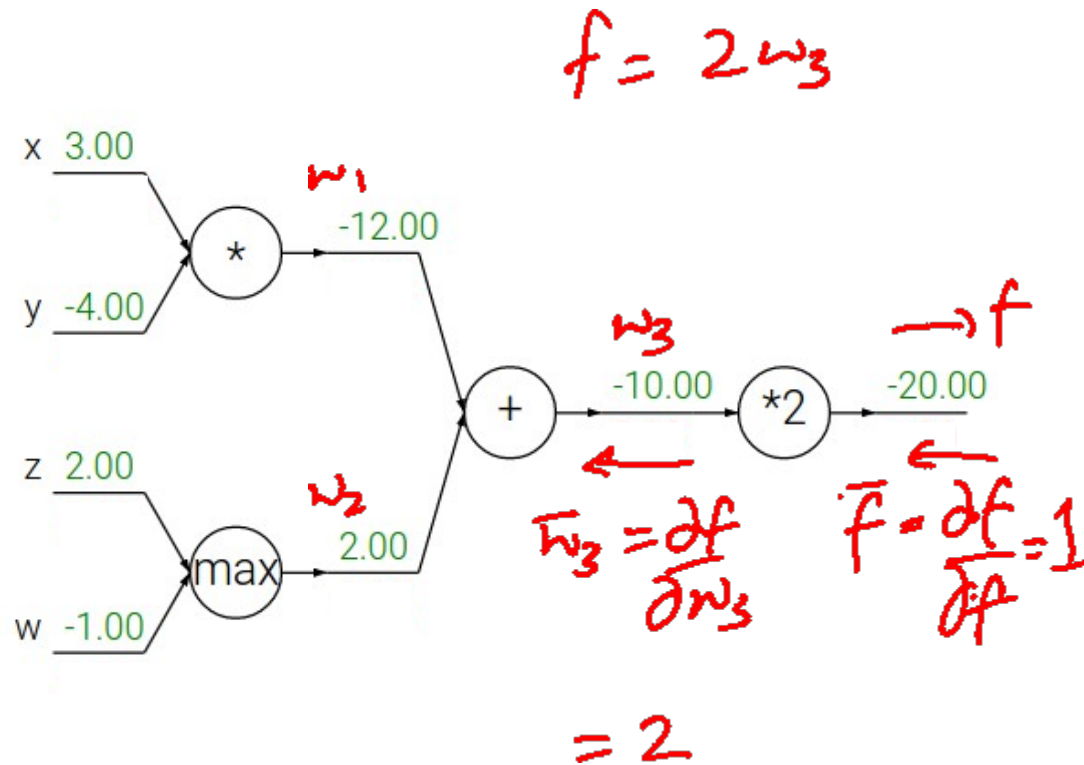


# Forward pass





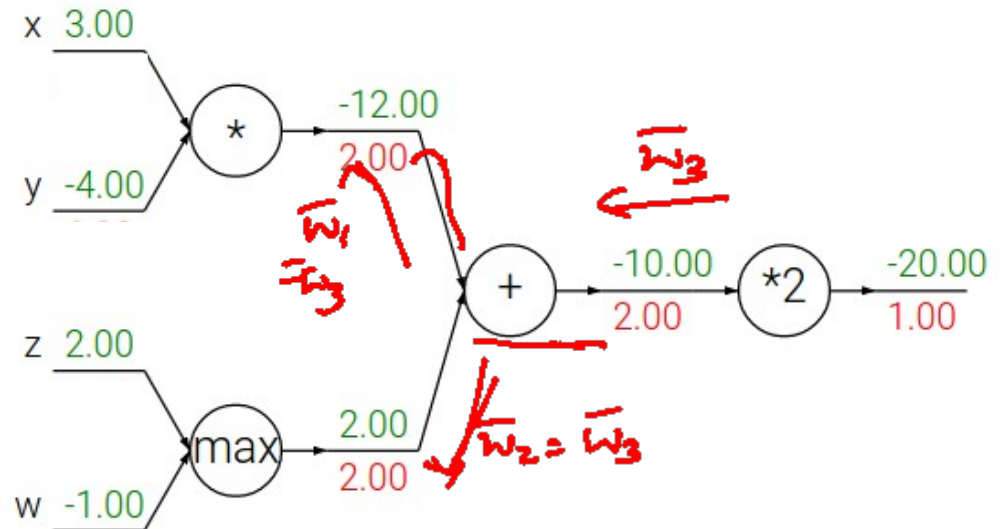
# Patterns in backprop





# Patterns in backprop

**add gate:** gradient distributor



# Patterns in backprop



add gate: gradient distributor

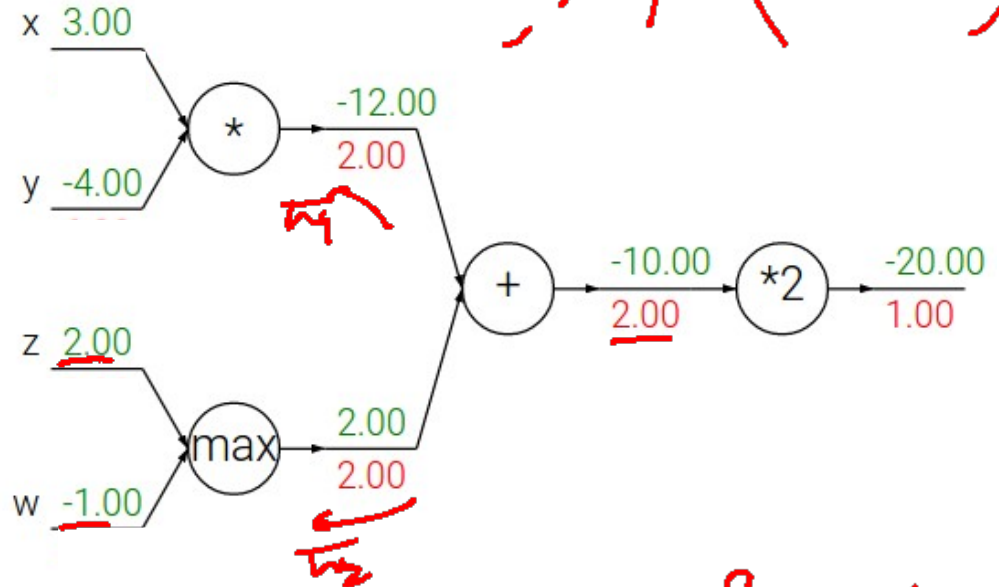
Q: What is a **max** gate? *function*

$$\bar{z} = \frac{\partial f}{\partial z} = \begin{bmatrix} \frac{\partial f}{\partial z} \\ \frac{\partial f}{\partial w_2} \end{bmatrix} = \begin{bmatrix} \frac{\partial f}{\partial z} \\ \frac{\partial f}{\partial w_2} \end{bmatrix}$$

$F_2$

$$\frac{\partial w_2}{\partial z} = \begin{cases} +1 & \text{if } z > w \\ 0 & \text{if } z < w \\ \text{Not defined} & \text{if } z = w \end{cases}$$

$$\frac{\partial w_2}{\partial w}$$



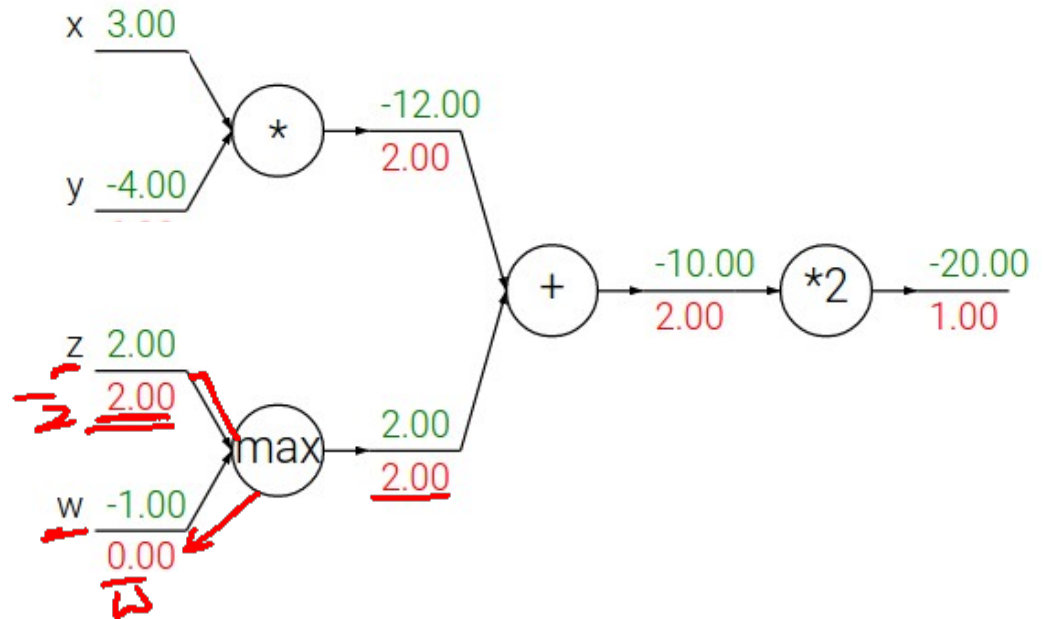
$$w_2 = \max\{z, w\} = \begin{cases} z & \text{if } z > w \\ w & \text{else} \end{cases}$$

# Patterns in backprop



**add** gate: gradient distributor

**max** gate: gradient router

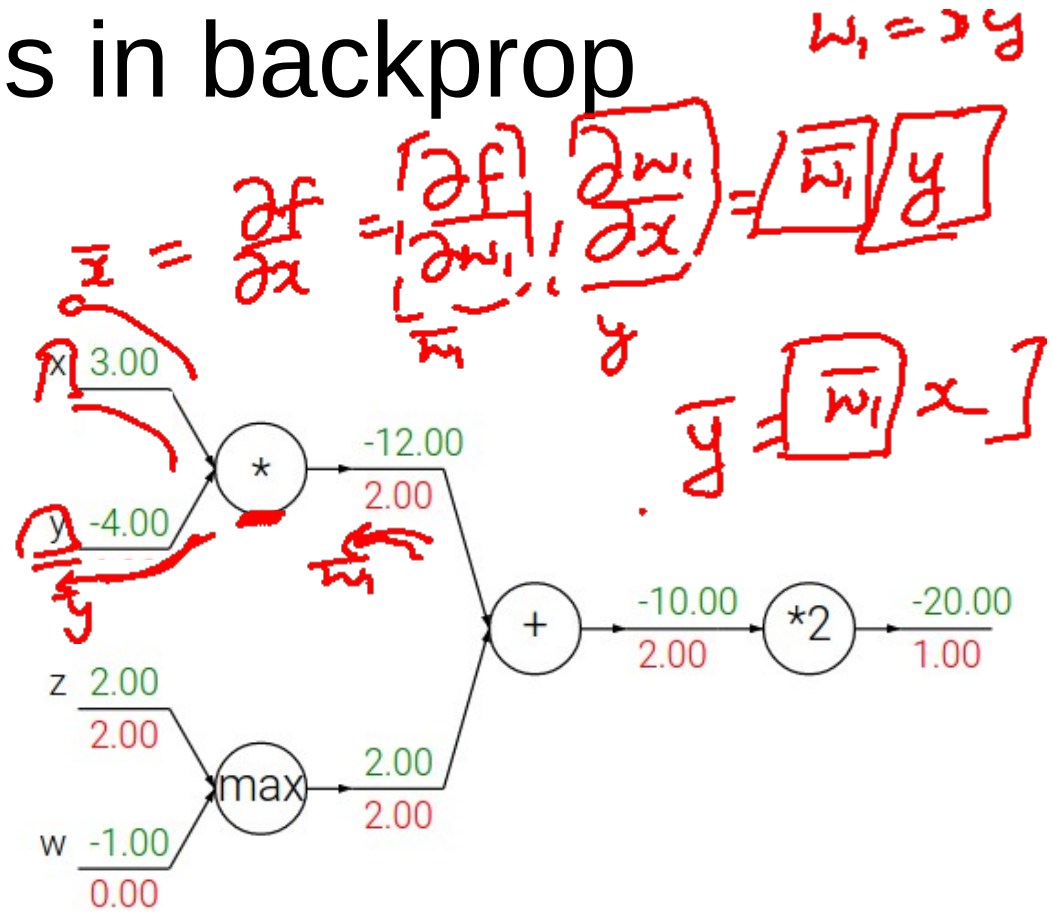


# Patterns in backprop

**add** gate: gradient distributor

**max** gate: gradient router

Q: What is a mul gate?



# Patterns in backprop

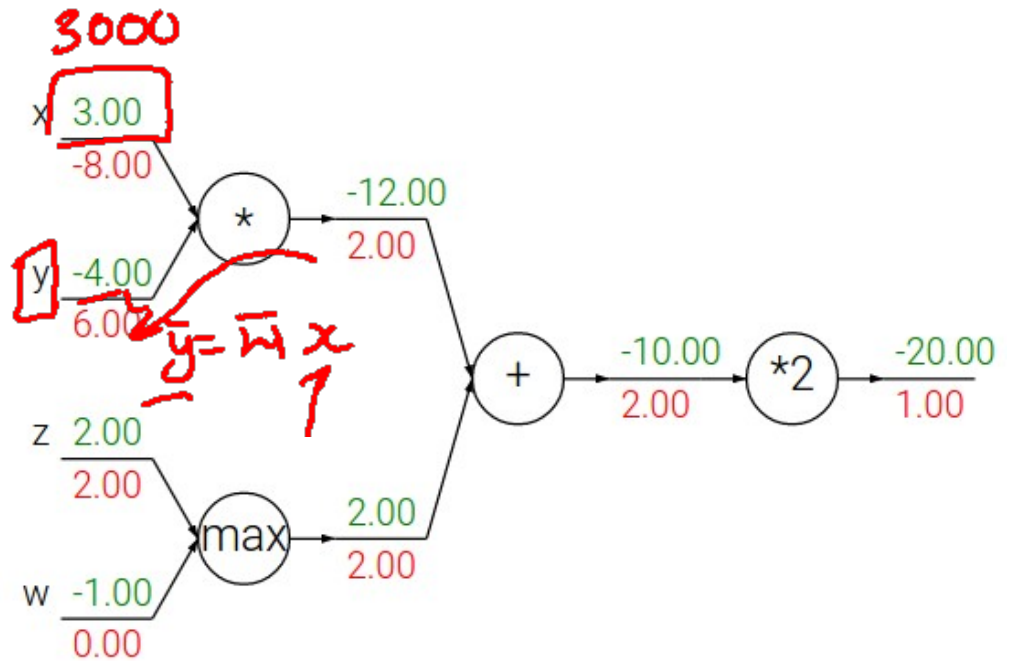


$$h_i^l = \vec{w}_i^T \vec{x}$$

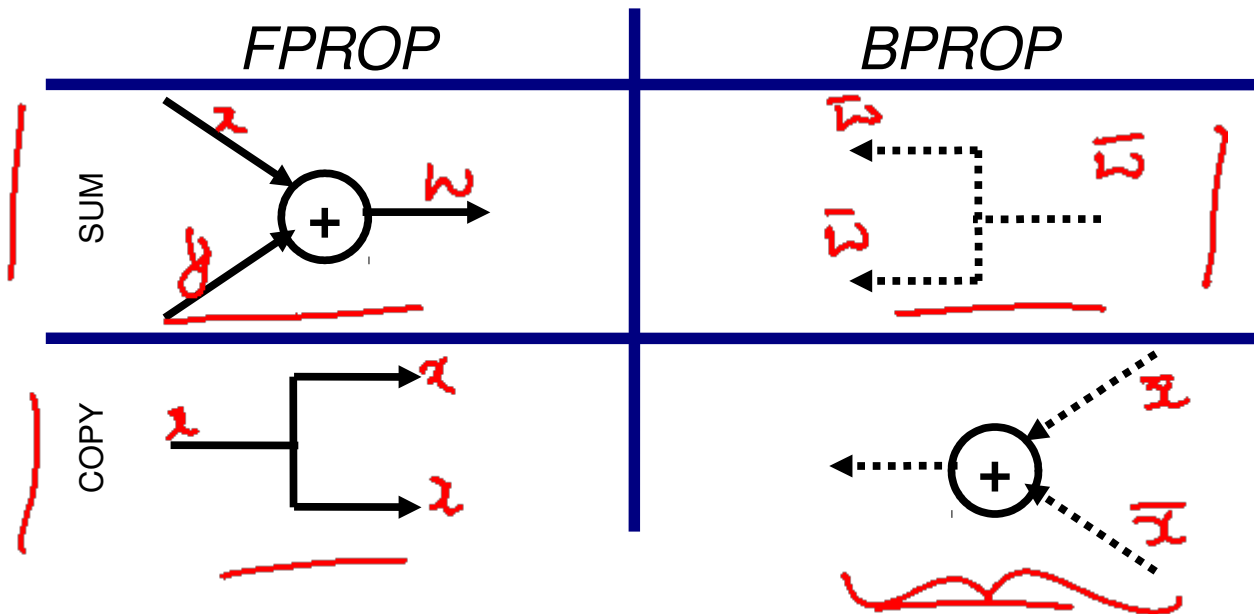
**add gate:** gradient distributor

**max gate:** gradient router

**mul gate:** gradient switcher

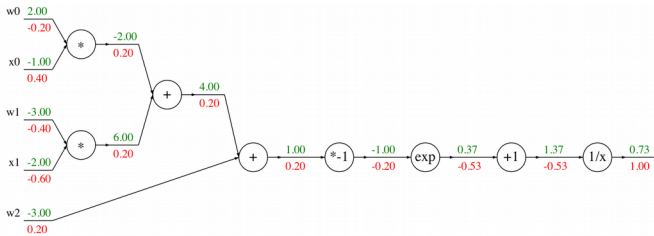


# Duality in Fprop and Bprop





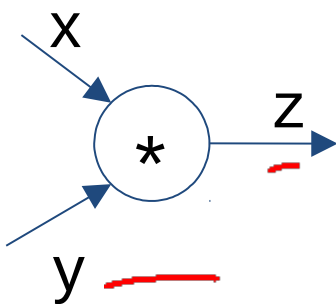
# Modularized implementation: forward / backward API



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



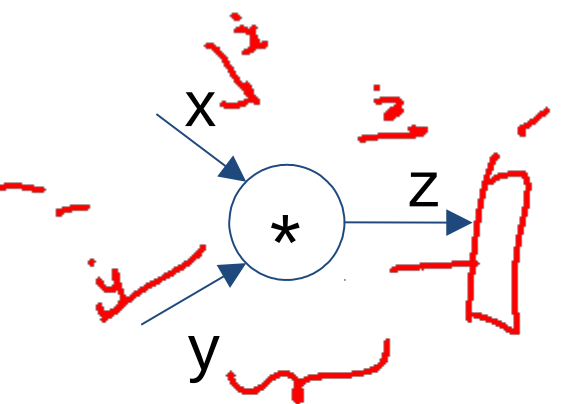
(x,y,z are scalars)

```
class MultiplyGate(object):  
    def forward(x,y):  
        z = x*y  
        return z  
    def backward(dz):  
        # dx = ... #todo  
        # dy = ... #todo  
        return [dx, dy]
```

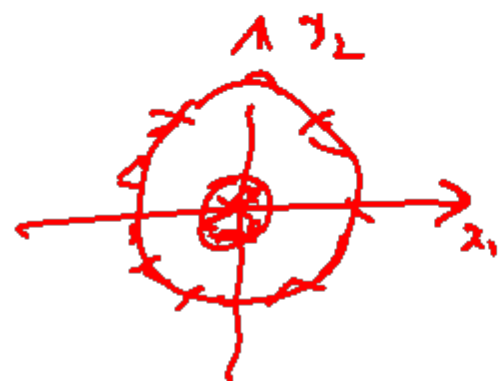
$$\frac{\partial L}{\partial z}$$

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

# Modularized implementation: forward / backward API



(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]
```

memory overhead  
of reverse  
mode  
AD.

# Example: Caffe layers

Branch: master | [caffe / src / caffe / layers /](#) | [Create new file](#) | [Upload files](#) | [Find file](#) | [History](#)

shelhamer committed on GitHub Merge pull request #4630 from BiGene/load\_hdf5\_fix Latest commit e687a71 21 days ago

<a href="#">absval_layer.cpp</a>	dismantle layer headers	a year ago
<a href="#">absval_layer.cu</a>	dismantle layer headers	a year ago
<a href="#">accuracy_layer.cpp</a>	dismantle layer headers	a year ago
<a href="#">argmax_layer.cpp</a>	dismantle layer headers	a year ago
<a href="#">base_conv_layer.cpp</a>	enable dilated deconvolution	a year ago
<a href="#">base_data_layer.cpp</a>	Using default from proto for prefetch	3 months ago
<a href="#">base_data_layer.cu</a>	Switched multi-GPU to NCCL	3 months ago
<a href="#">batch_norm_layer.cpp</a>	Add missing spaces besides equal signs in batch_norm_layer.cpp	4 months ago
<a href="#">batch_norm_layer.cu</a>	dismantle layer headers	a year ago
<a href="#">batch_reindex_layer.cpp</a>	dismantle layer headers	a year ago
<a href="#">batch_reindex_layer.cu</a>	dismantle layer headers	a year ago
<a href="#">bias_layer.cpp</a>	Remove incorrect cast of gemm int arg to Dtype in BiasLayer	a year ago
<a href="#">bias_layer.cu</a>	Separation and generalization of ChannelwiseAffineLayer into BiasLayer	a year ago
<a href="#">bnll_layer.cpp</a>	dismantle layer headers	a year ago
<a href="#">bnll_layer.cu</a>	dismantle layer headers	a year ago
<a href="#">concat_layer.cpp</a>	dismantle layer headers	a year ago
<a href="#">concat_layer.cu</a>	dismantle layer headers	a year ago
<a href="#">contrastive_loss_layer.cpp</a>	dismantle layer headers	a year ago
<a href="#">contrastive_loss_layer.cu</a>	dismantle layer headers	a year ago
<a href="#">conv_layer.cpp</a>	add support for 2D dilated convolution	a year ago
<a href="#">conv_layer.cu</a>	dismantle layer headers	a year ago
<a href="#">crop_layer.cpp</a>	remove redundant operations in Crop layer (#5138)	2 months ago
<a href="#">crop_layer.cu</a>	remove redundant operations in Crop layer (#5138)	2 months ago
<a href="#">cudnn_conv_layer.cpp</a>	dismantle layer headers	a year ago
<a href="#">cudnn_conv_layer.cu</a>	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago

<a href="#">cudnn_lcn_layer.cpp</a>	dismantle layer headers	a year ago
<a href="#">cudnn_lcn_layer.cu</a>	dismantle layer headers	a year ago
<a href="#">cudnn_lrn_layer.cpp</a>	dismantle layer headers	a year ago
<a href="#">cudnn_lrn_layer.cu</a>	dismantle layer headers	a year ago
<a href="#">cudnn_pooling_layer.cpp</a>	dismantle layer headers	a year ago
<a href="#">cudnn_pooling_layer.cu</a>	dismantle layer headers	a year ago
<a href="#">cudnn_relu_layer.cpp</a>	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
<a href="#">cudnn_relu_layer.cu</a>	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
<a href="#">cudnn_sigmoid_layer.cpp</a>	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
<a href="#">cudnn_sigmoid_layer.cu</a>	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
<a href="#">cudnn_softmax_layer.cpp</a>	dismantle layer headers	a year ago
<a href="#">cudnn_softmax_layer.cu</a>	dismantle layer headers	a year ago
<a href="#">cudnn_tanh_layer.cpp</a>	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
<a href="#">cudnn_tanh_layer.cu</a>	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
<a href="#">data_layer.cpp</a>	Switched multi-GPU to NCCL	3 months ago
<a href="#">deconv_layer.cpp</a>	enable dilated deconvolution	a year ago
<a href="#">deconv_layer.cu</a>	dismantle layer headers	a year ago
<a href="#">dropout_layer.cpp</a>	supporting N-D Blobs in Dropout layer Reshape	a year ago
<a href="#">dropout_layer.cu</a>	dismantle layer headers	a year ago
<a href="#">dummy_data_layer.cpp</a>	dismantle layer headers	a year ago
<a href="#">eltwise_layer.cpp</a>	dismantle layer headers	a year ago
<a href="#">eltwise_layer.cu</a>	dismantle layer headers	a year ago
<a href="#">elu_layer.cpp</a>	ELU layer with basic tests	a year ago
<a href="#">elu_layer.cu</a>	ELU layer with basic tests	a year ago
<a href="#">embed_layer.cpp</a>	dismantle layer headers	a year ago
<a href="#">embed_layer.cu</a>	dismantle layer headers	a year ago
<a href="#">euclidean_loss_layer.cpp</a>	dismantle layer headers	a year ago
<a href="#">euclidean_loss_layer.cu</a>	dismantle layer headers	a year ago
<a href="#">exp_layer.cpp</a>	Solving issue with exp layer with base e	a year ago
<a href="#">exp_layer.cu</a>	dismantle layer headers	a year ago

Caffe is licensed under [BSD 2-Clause](#)

# Caffe Sigmoid Layer

```
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 } // namespace caffe
```

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

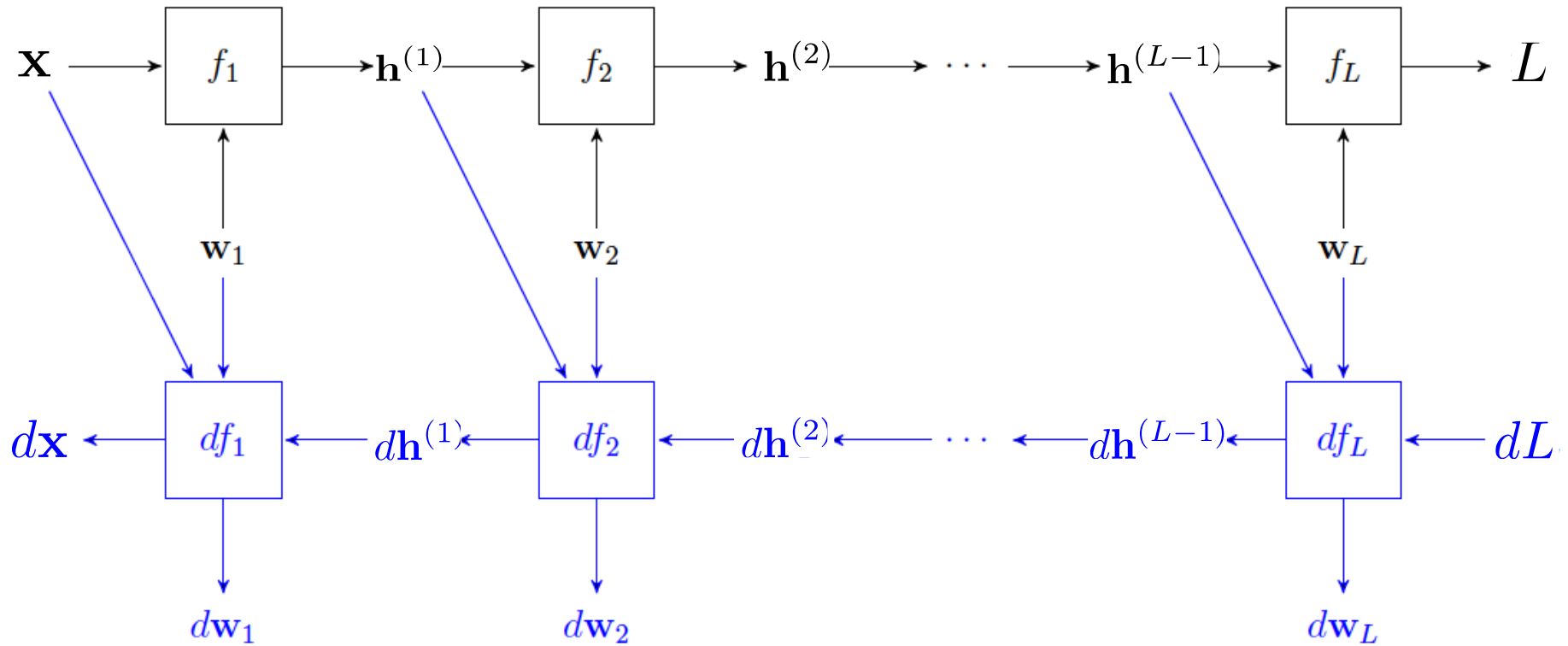
$$(1 - \sigma(x)) \sigma(x) * \text{top\_diff} \text{ (chain rule)}$$

Caffe is licensed under [BSD 2-Clause](#)

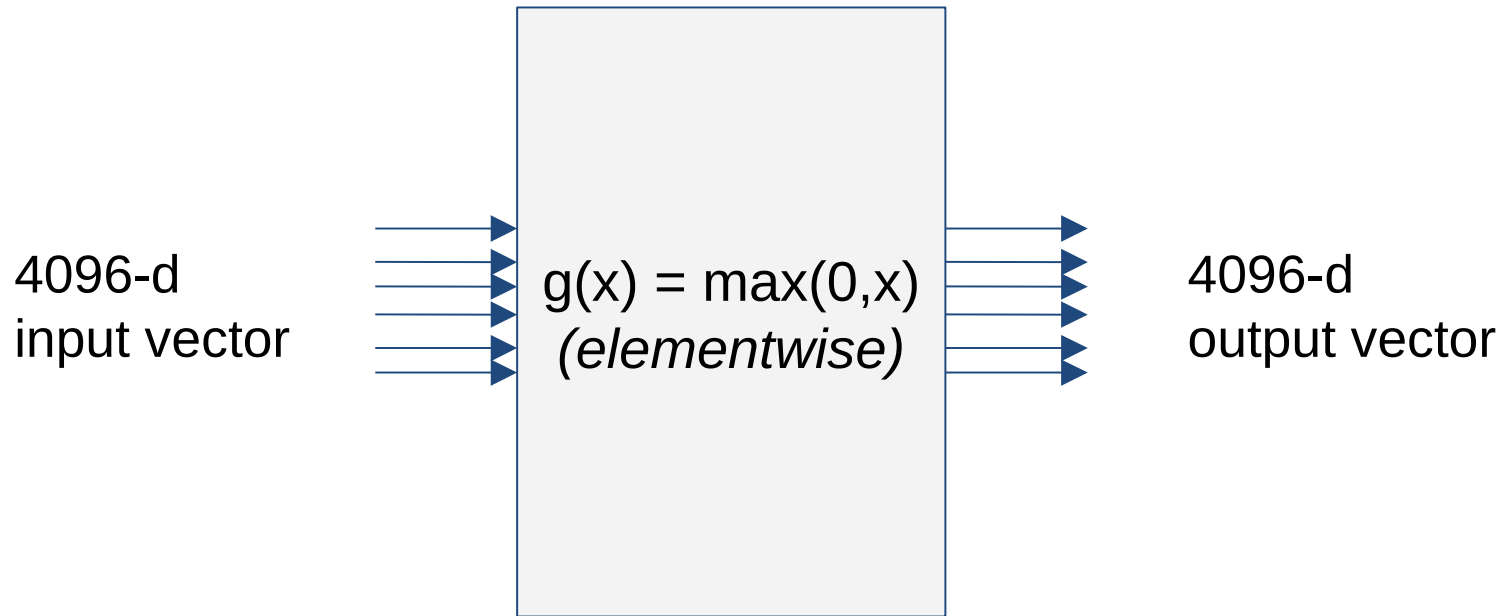
# Plan for Today

- Automatic Differentiation
  - Patterns in backprop
  - Jacobians in FC+ReLU NNs

# Backprop

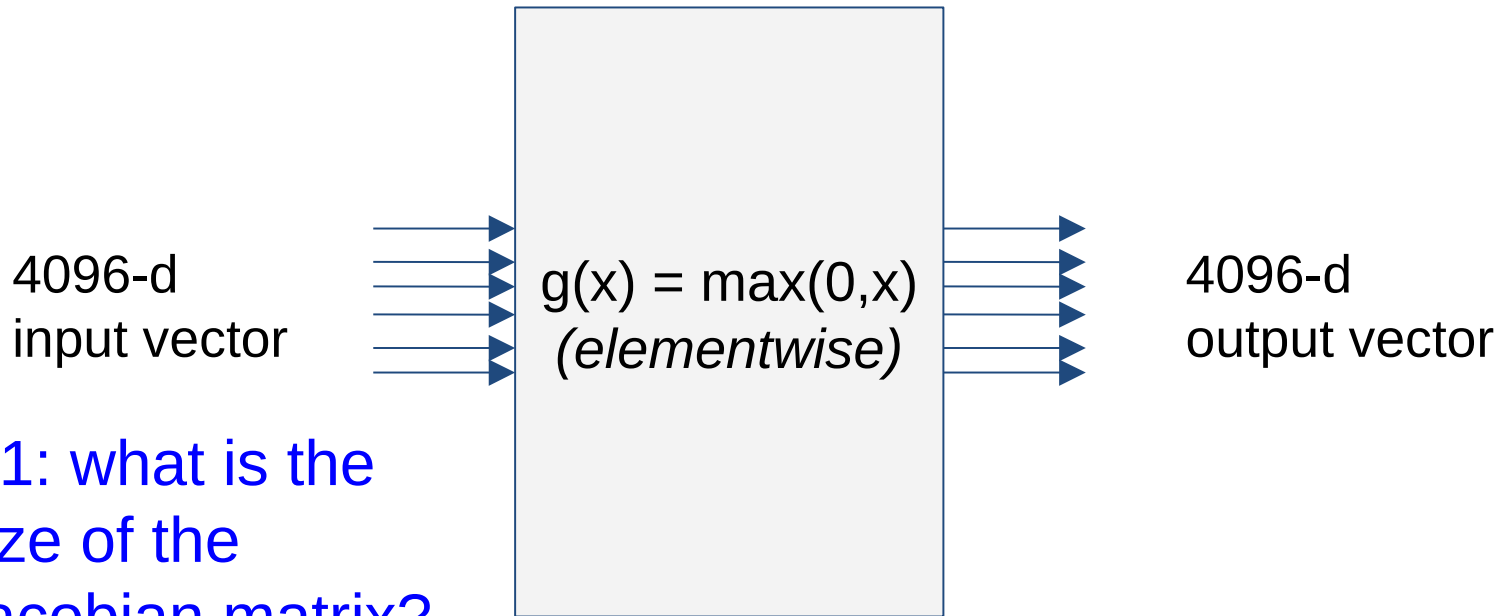


# Jacobian of ReLU



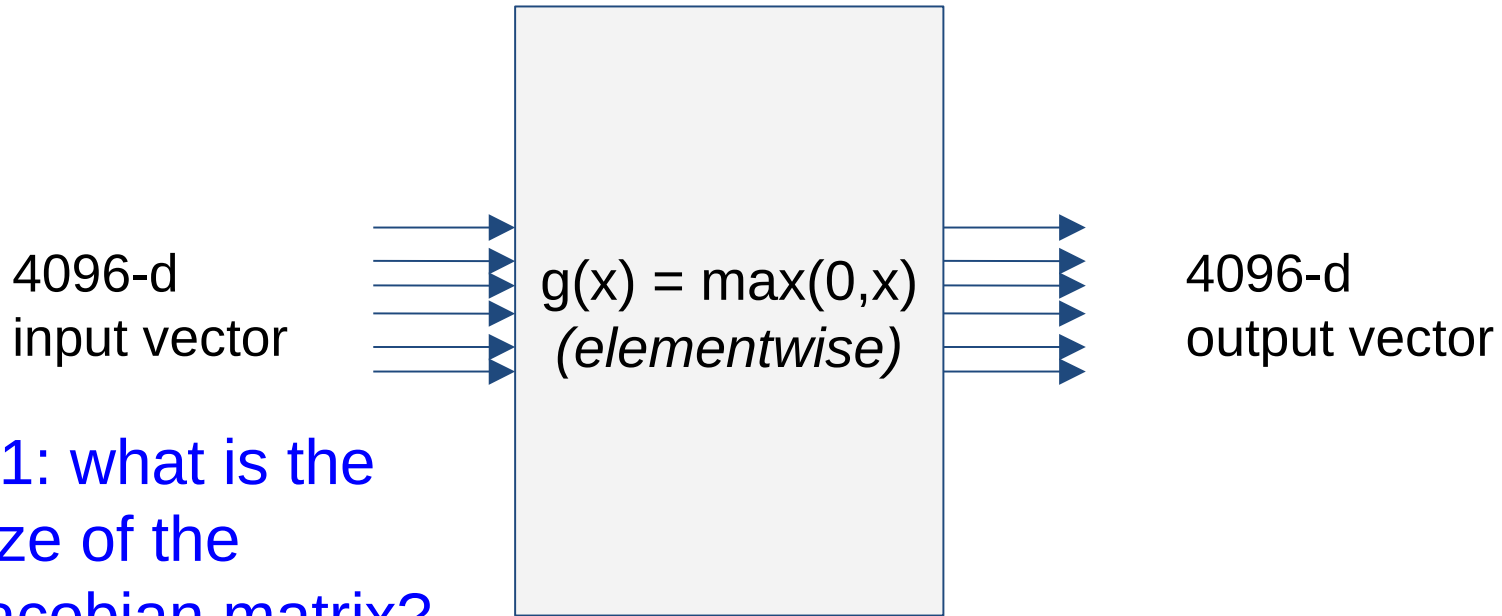


# Jacobian of ReLU



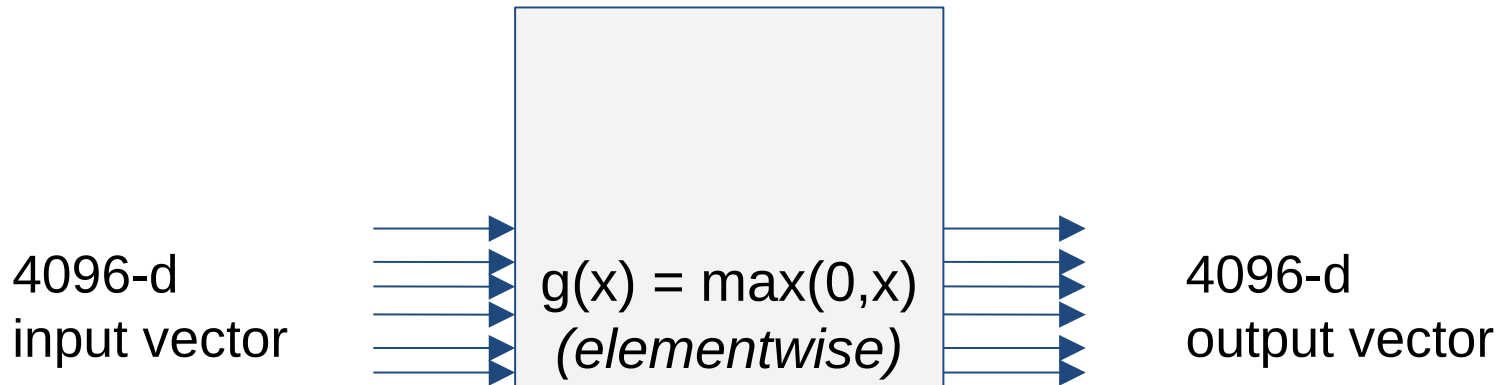
Q1: what is the  
size of the  
Jacobian matrix?

# Jacobian of ReLU



Q1: what is the  
size of the  
Jacobian matrix?  
[4096 x 4096!]

# Jacobian of ReLU

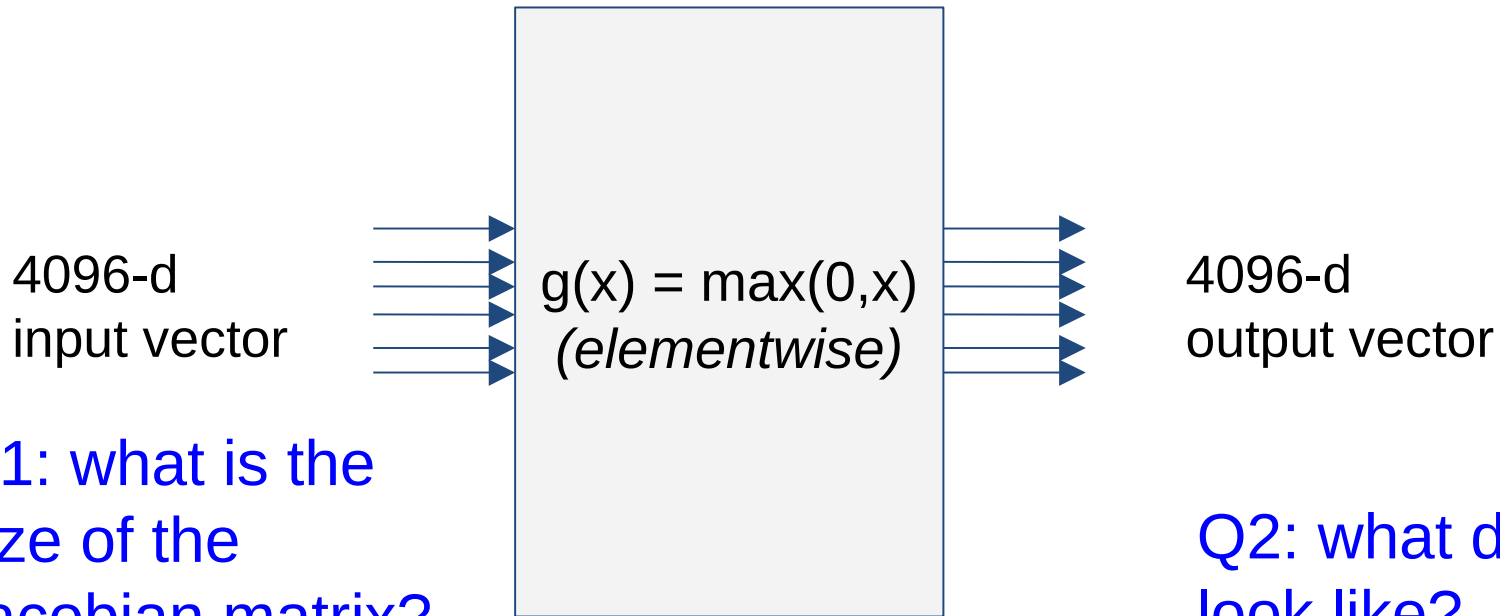


Q1: what is the  
size of the  
Jacobian matrix?  
[4096 x 4096!]

in practice we process an  
entire minibatch (e.g. 100)  
of examples at one time:

i.e. Jacobian would technically be a  
[409,600 x 409,600] matrix :)

# Jacobian of ReLU



Q1: what is the size of the Jacobian matrix?  
[4096 x 4096!]

Q2: what does it look like?

# Jacobians of FC-Layer

# Jacobians of FC-Layer

# Jacobians of FC-Layer

# Jacobians of FC-Layer