

# CS6423: Scalable Computing

---

**Gregory Provan**

Spring 2020

Lecture 8: TensorFlow and Computation Graphs

Based on notes from Hung-Yi Lee,  
Andrej Karpathy, Fei-Fei Li, Justin Johnson

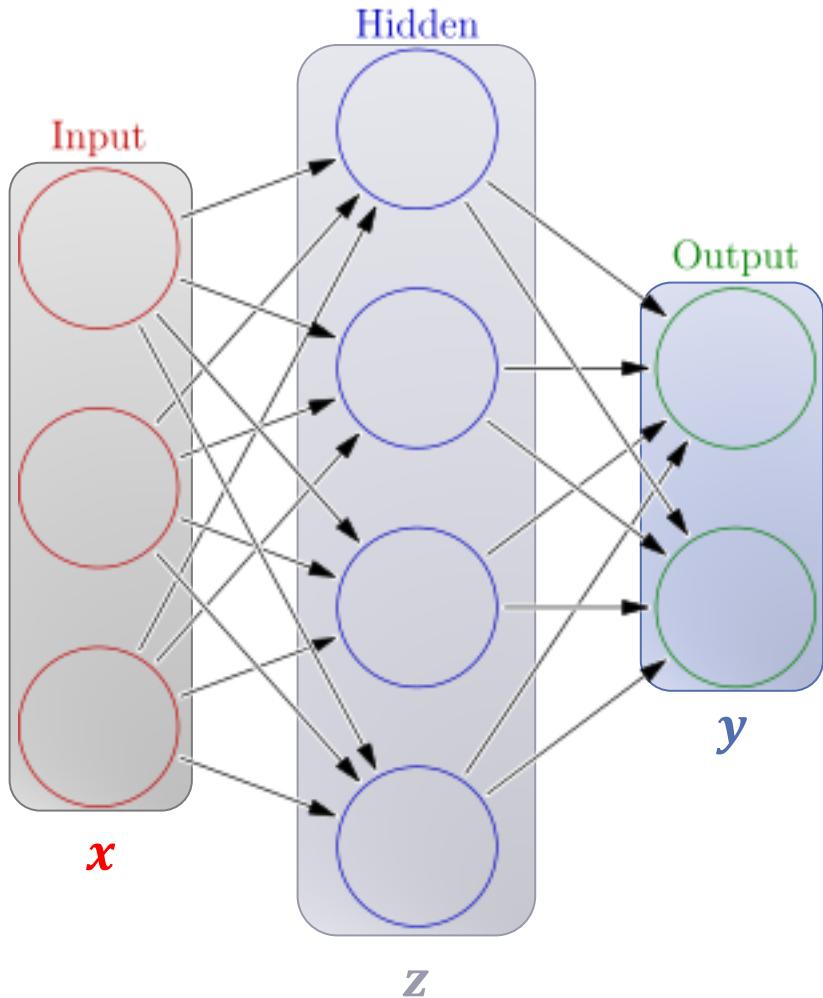
# Motivation

- Deep network
  - an inter-connected set of modules
- Network inference uses this modular decomposition
  - All DL platforms: TensorFlow, Theano, etc.
- Computation graph
  - Mathematical foundations for inference

# Overview

- Modularity and Computation Graphs
- Inference in Computation Graphs
- Backpropagation
  - Computation Graphs viewpoint
- DL Platforms
  - TensorFlow, Theano, etc.

# Neural Network: Definition



Weights

$$z = f(W_1x + b_1)$$
$$y = g(W_2z + b_2)$$

Activation functions

$4 + 2 = 6$  neurons (not counting inputs)

$[3 \times 4] + [4 \times 2] = 20$  weights

$4 + 2 = 6$  biases

---

26 learnable parameters

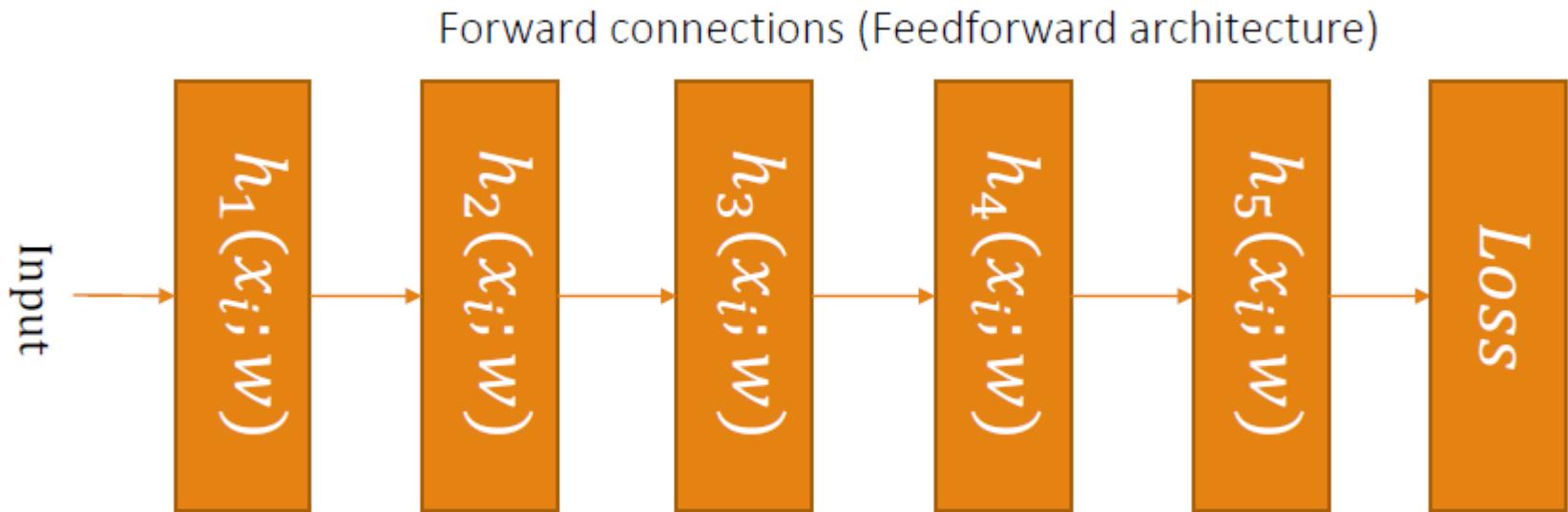
Demo

# Neural Network: Definition

- A family of parametric, non-linear and hierarchical representation learning functions
  - massively optimized with stochastic gradient descent to encode domain knowledge, i.e. domain invariances, stationarity.
- $a^L(x; w^1, \dots, w^L) = h^L(h^{L-1} \dots h^1(x, w^1), w^{L-1}), w^L)$ 
  - $x$ : input,
  - $w^l$ : parameters for layer  $l$ ,
  - $a^l = h^l(x, w^l)$ : (non-) linear function
- Given training corpus  $\{X, Y\}$  find optimal parameters
  - $w^* \leftarrow \operatorname{argmin}_w \sum_{(x,y) \subseteq (X,Y)} L(y, a^L(x))$

# Architectural View of Deep Networks

- A neural network model is a series of hierarchically connected functions
- These hierarchies can be very complex



# What is TensorFlow?

 tensorflow / tensorflow

Watch ▾ 7,777   Star 96,717   Fork 61,507

 Code   Issues 1,313   Pull requests 196   Projects 0   Insights

Computation using data flow graphs for scalable machine learning <https://tensorflow.org>

tensorflow   machine-learning   python   deep-learning   deep-neural-networks   neural-network   ml   distributed

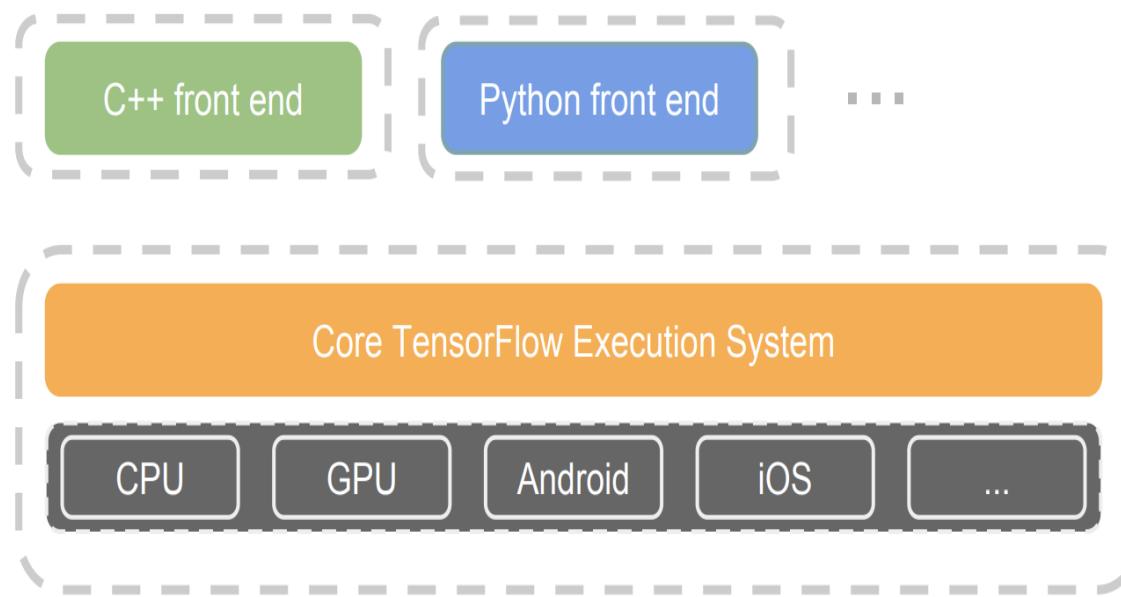
 31,895 commits    31 branches    54 releases    1,435 contributors    Apache-2.0



- Open source library for numerical computation using **computation graphs**
- Developed by Google Brain Team to conduct machine learning research
  - Based on DisBelief used internally at Google since 2011
- “TensorFlow is an interface for expressing machine learning algorithms, and an implementation for executing such algorithms”

# TensorFlow Architecture

- Core in C++
  - Very low overhead
- Different front ends for specifying/driving the computation
  - Python and C++ today, easy to add more



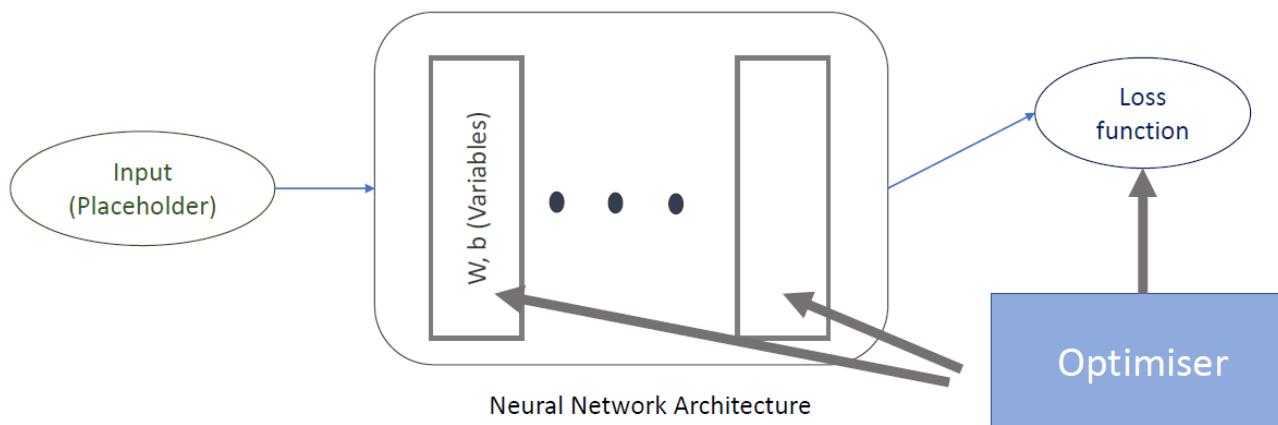
# TF Learn

- TensorFlow with an SkLearn API
  - `model.fit()`
  - `model.predict()`
  - `model.evaluate()`
  - `model.save()` and `model.load()`
- GPU accelerated deep learning in 8 lines of Python

```
import tflearn  
  
tflearn.init_graph()  
  
net = tflearn.input_data(shape=[None, 3])  
net = tflearn.fully_connected(net, 4,  
                             activation='relu')  
net = tflearn.fully_connected(net, 2,  
                             activation='relu')  
net = tflearn.regression(net, optimizer='adam')  
  
model = tflearn.DNN(net)  
model.fit(X, Y)
```

# How Does TensorFlow Work?

- TensorFlow input elements
  - Input data (placeholder)
  - Deep-network architecture
  - Loss function
- Applies optimising compiler
  - Uses structure of network
  - Distributes inference to available processors (CPU, GPU, etc.)
  - Generates *computation graph* for optimising inference

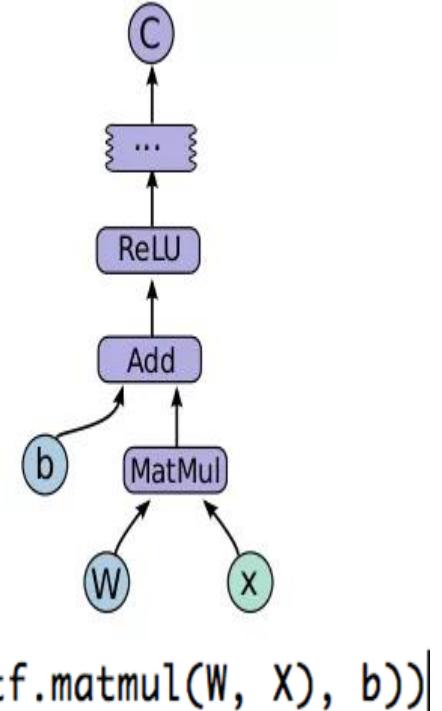


# TensorFlow

- Large-Scale Machine Learning Across Heterogeneous Distributed Systems
- Normally, we run programs (largely) sequentially
- TensorFlow:
  - builds a *computation graph*
  - assigns operations to optimal hardware
  - **then** runs the computations

```
import tensorflow  
import numpy
```

```
C = tf.relu(tf.add(tf.matmul(W, X), b))
```

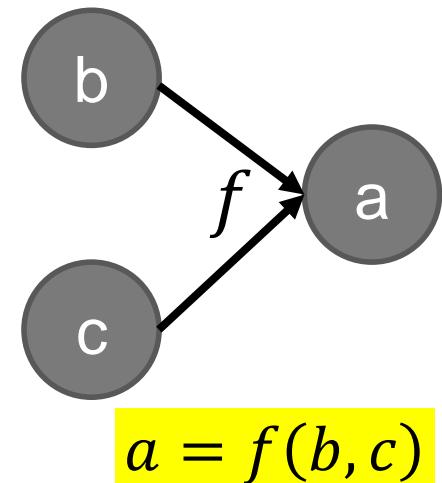


# What is TensorFlow

- **Key idea:** express a numeric computation as a **graph**
- Graph nodes are **operations** with any number of inputs and outputs
- Graph edges are **tensors** which flow between nodes

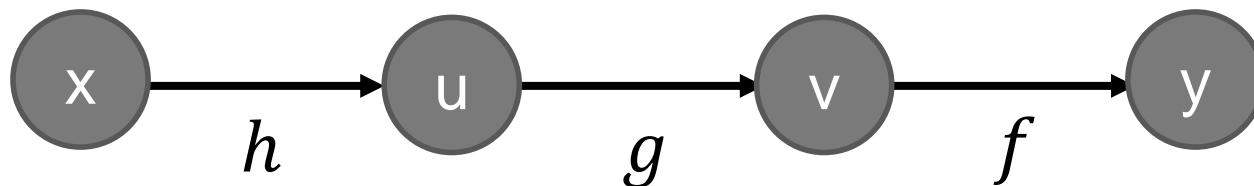
# Computation Graph

- A “language” describing a function
  - **Node**: variable (scalar, vector, tensor .....)
  - **Edge**: operation (simple function)

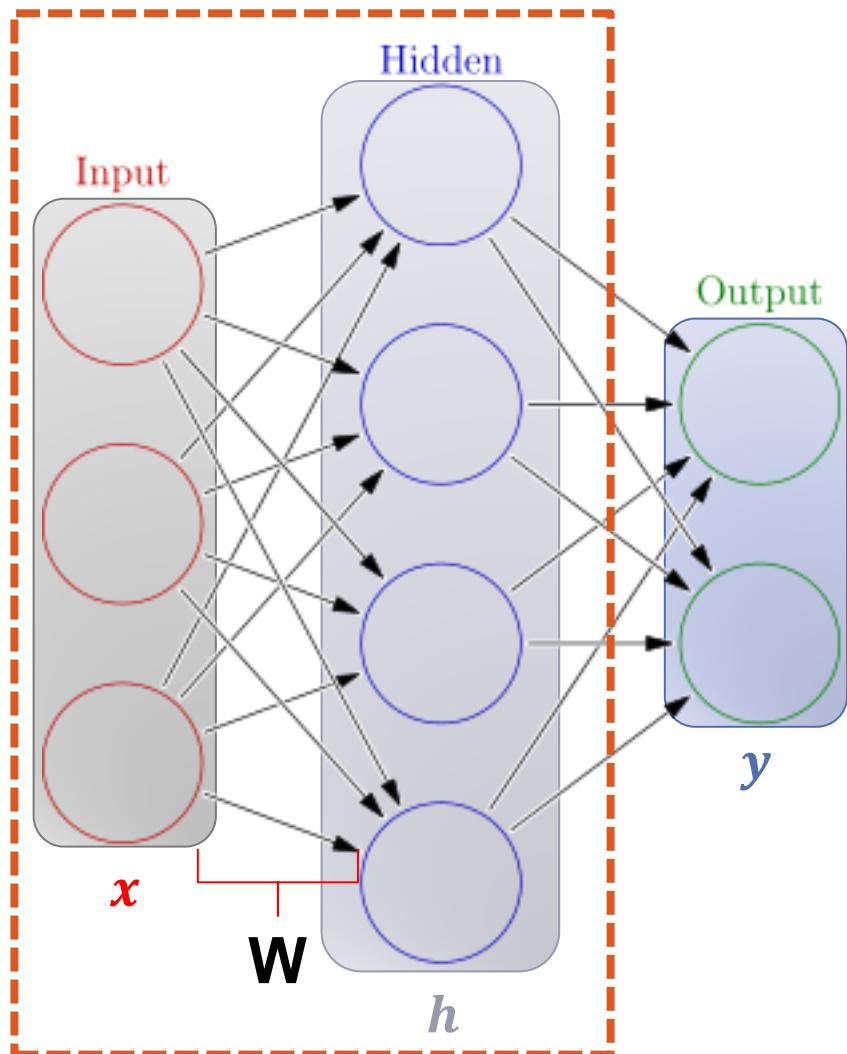


**Example**  $y = f(g(h(x)))$

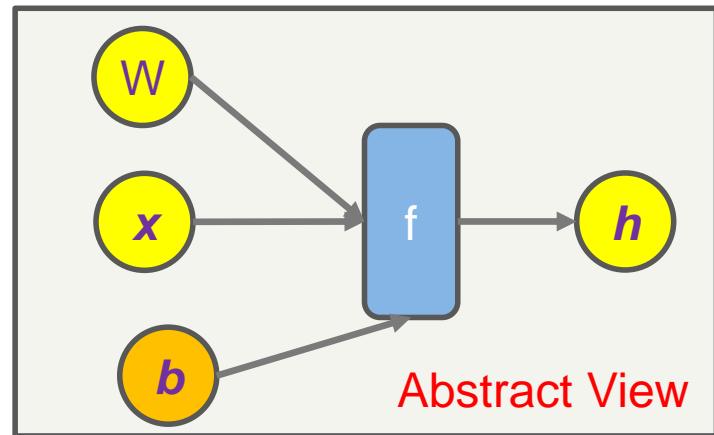
$$u = h(x) \quad v = g(u) \quad y = f(v)$$



# Modular Operations: Computation Graph Definition

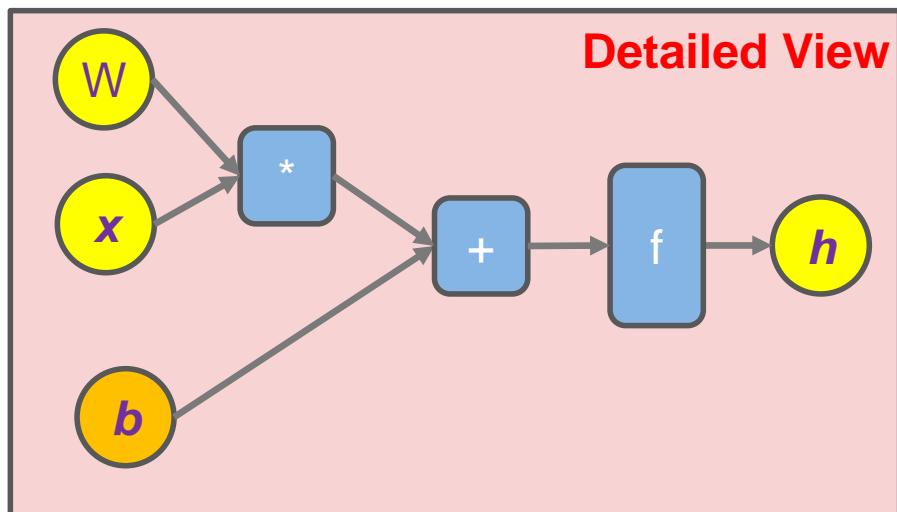


$$h = f(Wx + b)$$



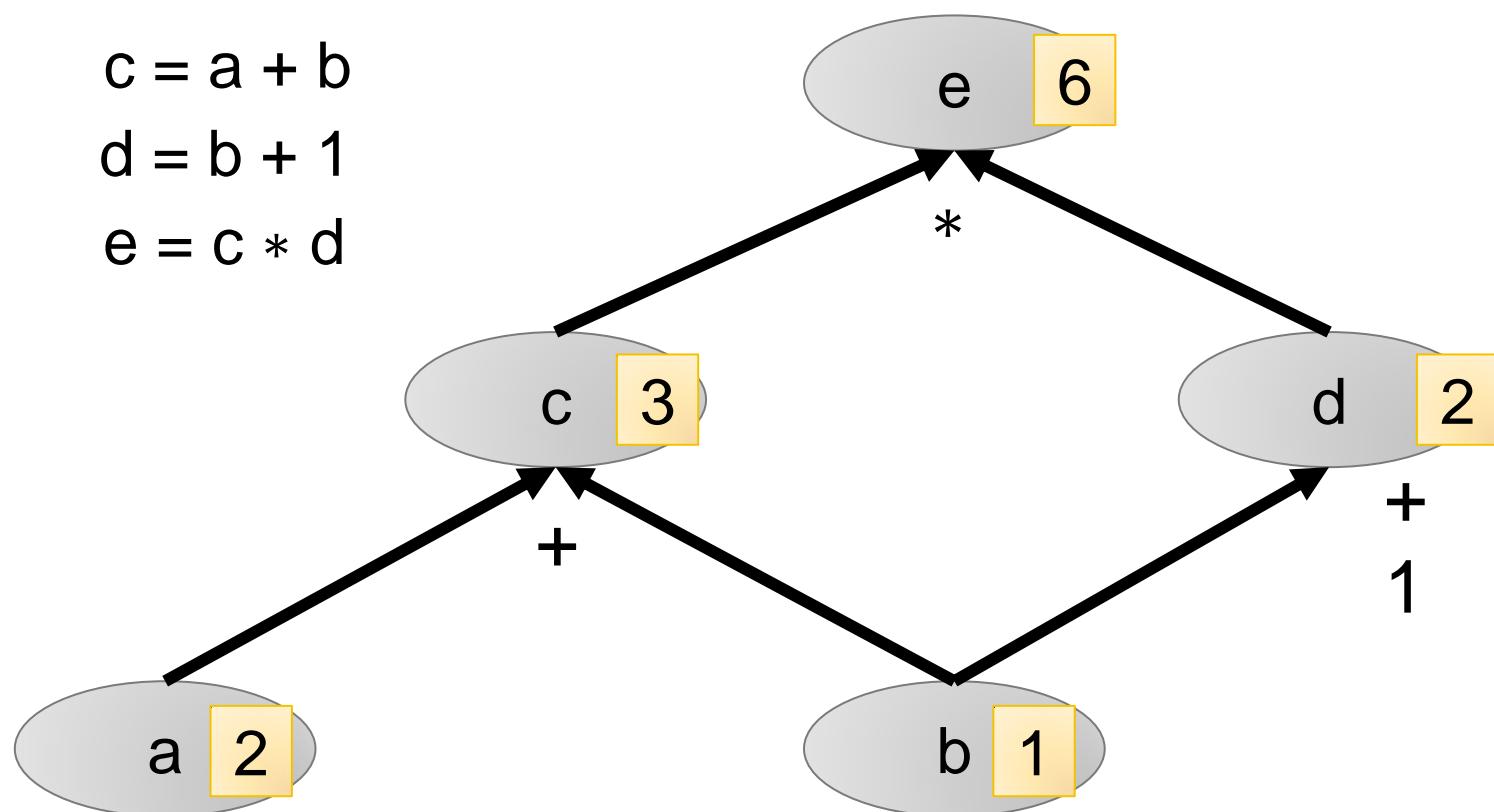
Abstract View

Detailed View



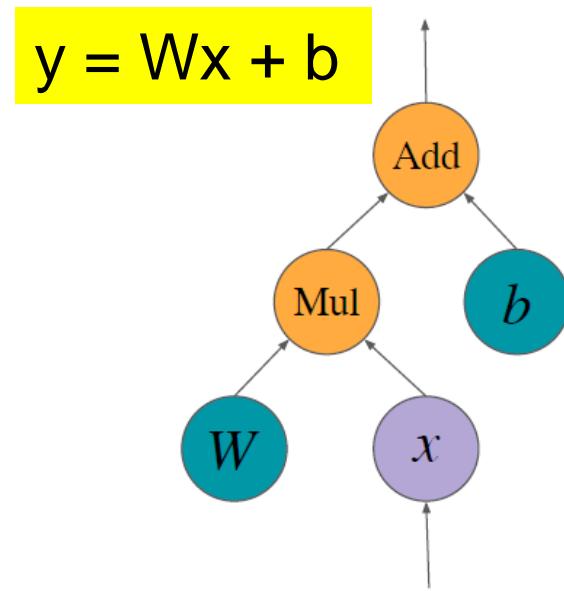
# Computational Graph

- Example:  $e = (a+b) * (b+1)$



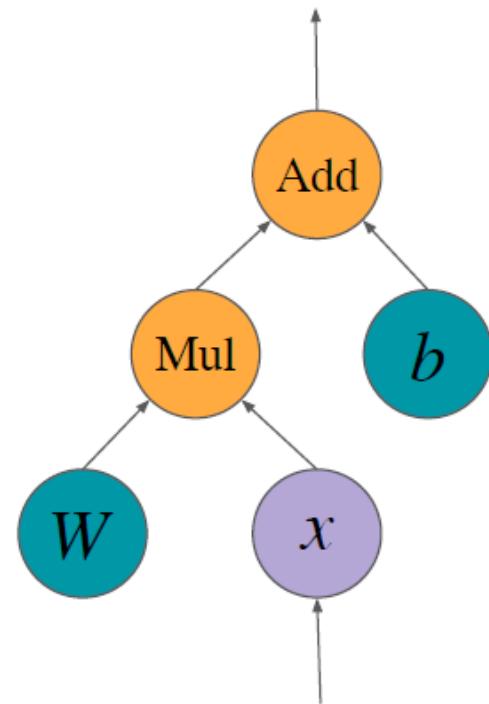
# Basic Code Structure - Graphs

- Nodes are operators (ops), variables, and constants
- Edges are tensors
  - 0-d is a scalar
  - 1-d is a vector
  - 2-d is a matrix
  - Etc.
- TensorFlow = Tensor + Flow = Data + Flow



# Computation Graph

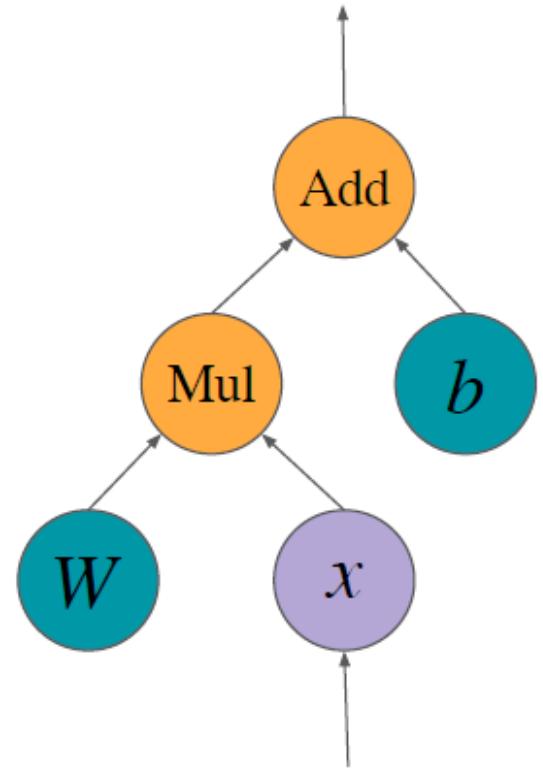
- **Constants**
  - fixed value tensors - not trainable
- **Variables**
  - tensors initialized in a session - trainable
- **Placeholders**
  - tensors of values that are unknown during the graph construction
  - passed as input during a session
- **Ops**
  - functions on tensors



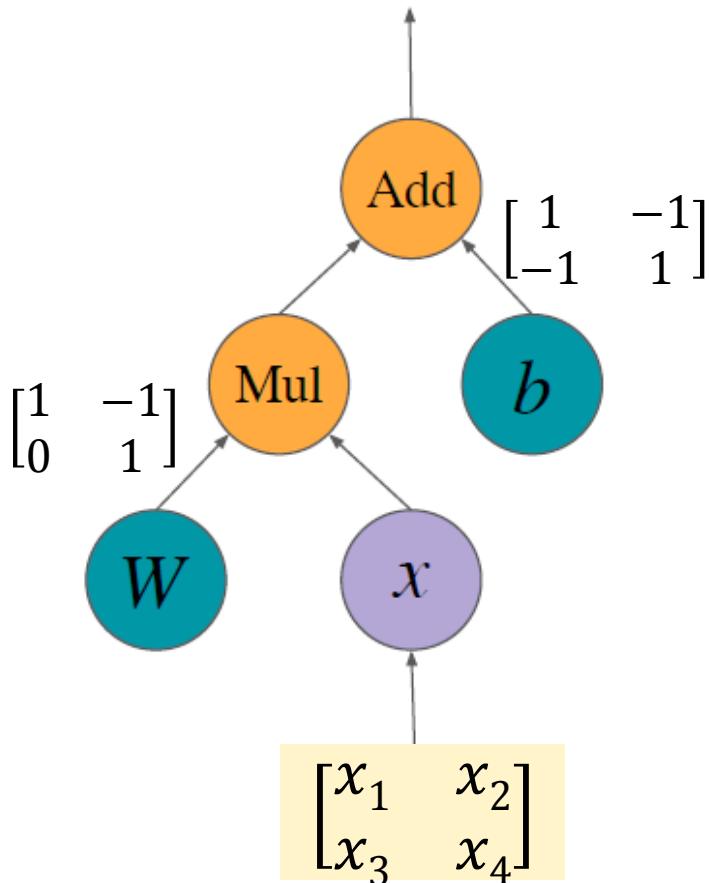
# Computation Graph Structure

$$\hat{y} = Wx + b$$

variable  
ops  
placeholder



# Example



$$y = Wx + b$$

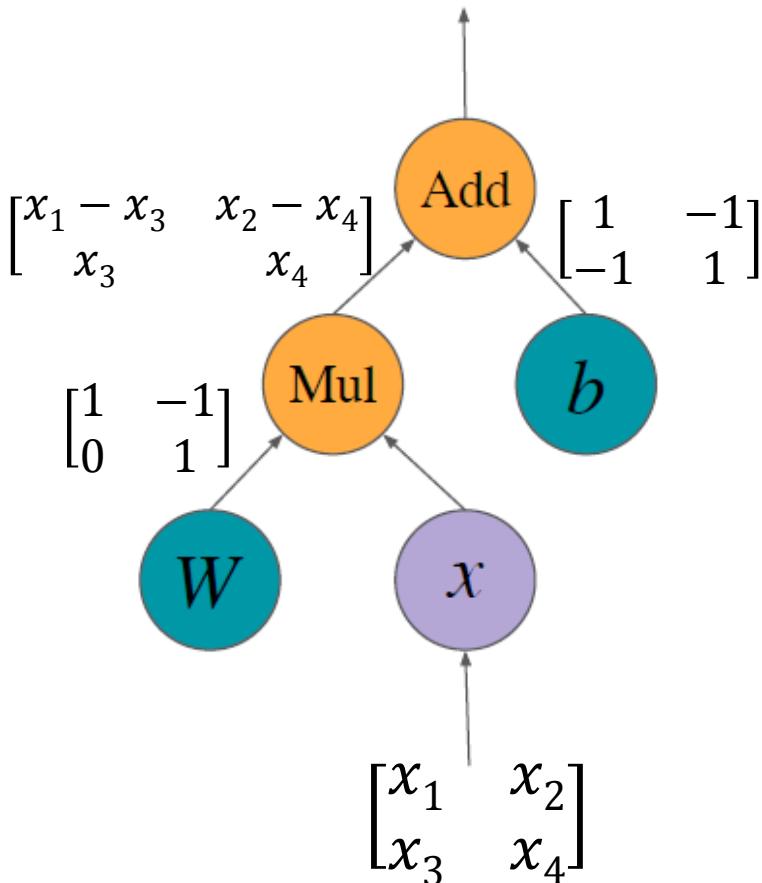
$$\begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 & x_2 \\ x_3 & x_4 \end{bmatrix} + \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$$

W

x

b

# Evaluating Mult



$$y = Wx + b$$

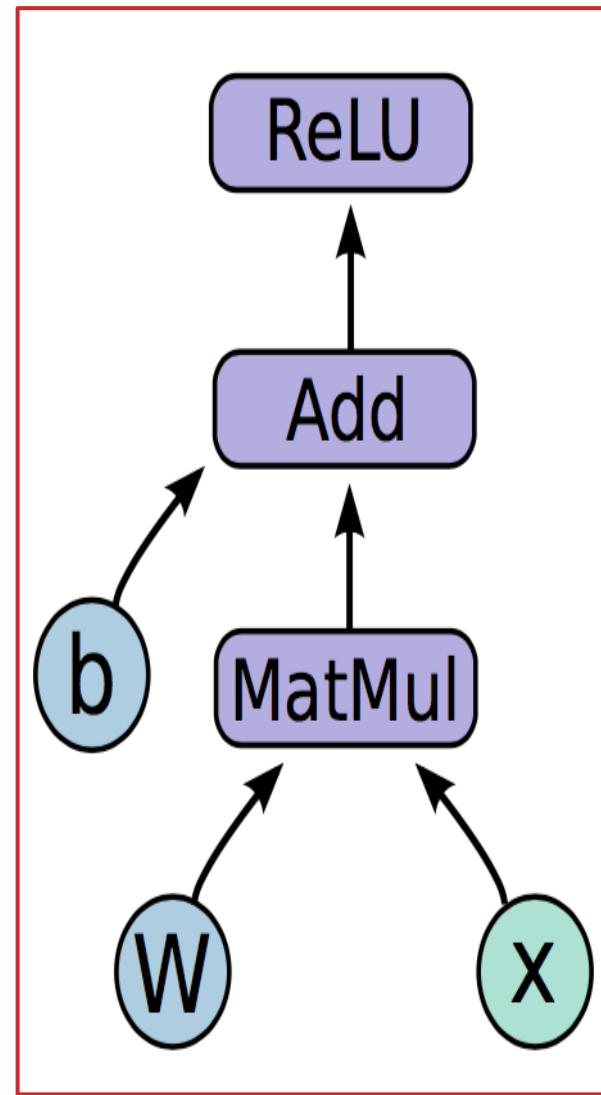
$$\begin{bmatrix} x_1 - x_3 & x_2 - x_4 \\ x_3 & x_4 \end{bmatrix} + \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$$

$$Wx$$

$$b$$

# Programming Model: Computation Graph

$$h = \text{ReLU}(Wx + b)$$



# Programming Model : Steps for DL Implementation

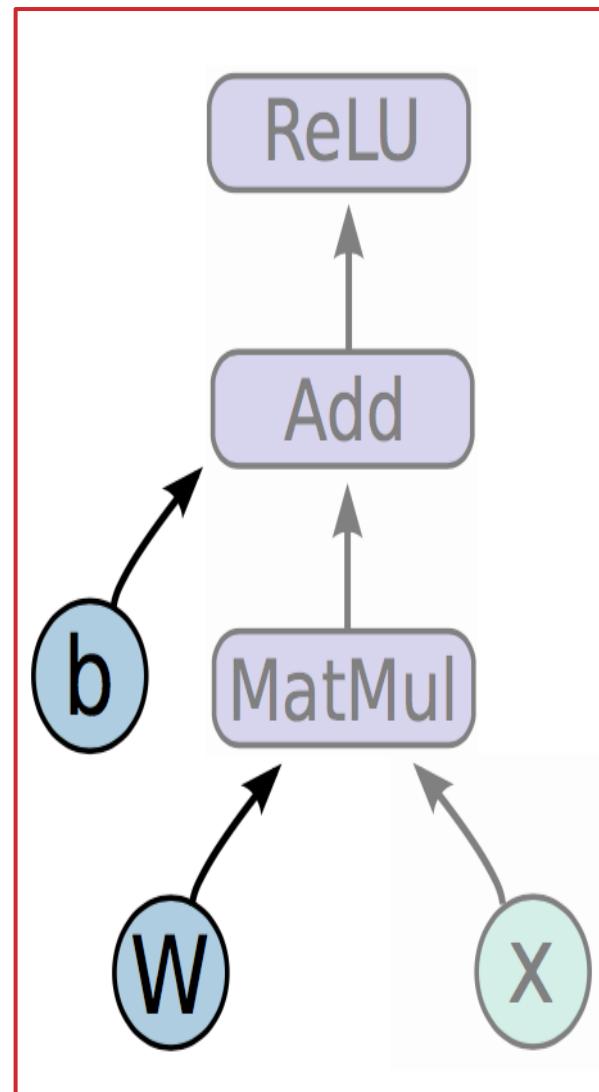
1. Define the Computation Graph
2. Initialize your data
3. Assign to hardware (CPU, GPU, etc.)

# Step 1: Computation Graph

$$h = \text{ReLU}(Wx + b)$$

**Variables** are stateful nodes which output their current value. State is retained across multiple executions of a graph

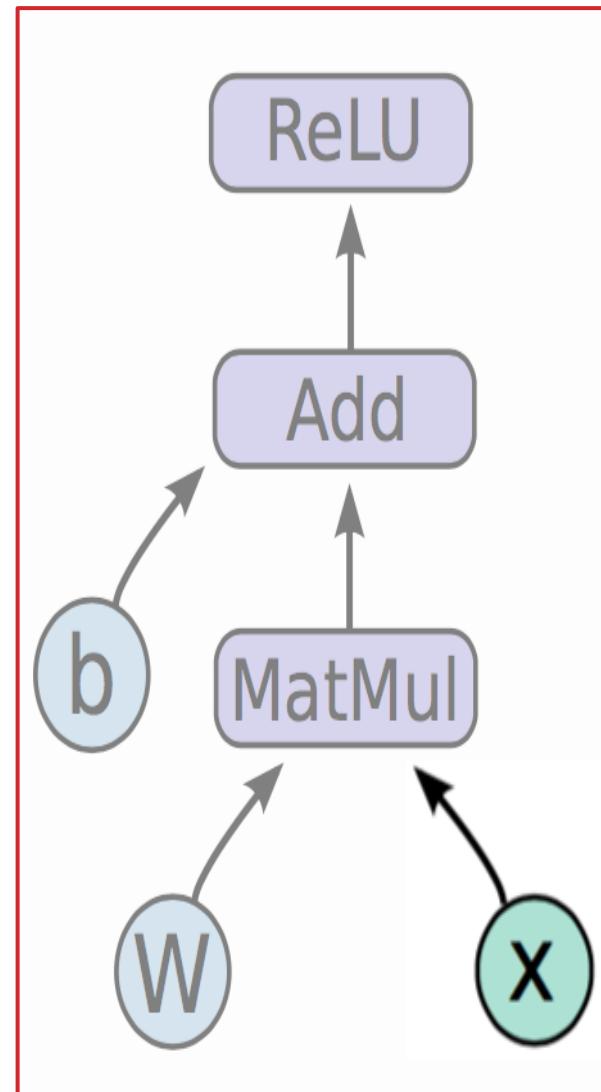
(mostly parameters)



# Step 1: Computation Graph

$$h = \text{ReLU}(Wx + b)$$

**Placeholders** are nodes  
whose value is fed in at  
execution time  
(inputs, labels, ...)



# Step 1: Computation Graph

$$h = \text{ReLU}(Wx + b)$$

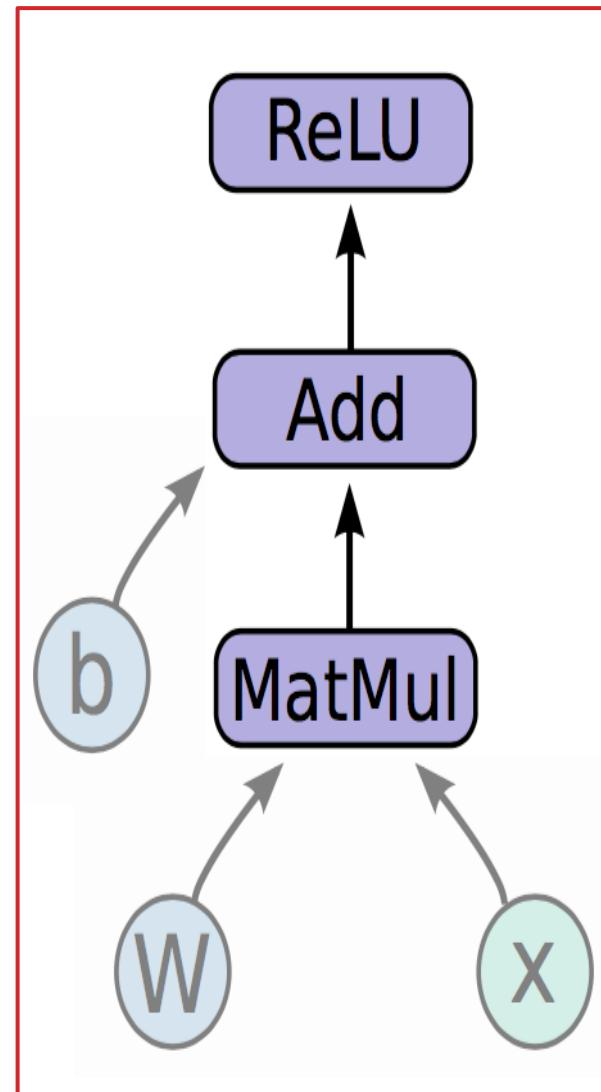
## Mathematical operations:

**MatMul**: Multiply two matrices

**Add**: Add elementwise

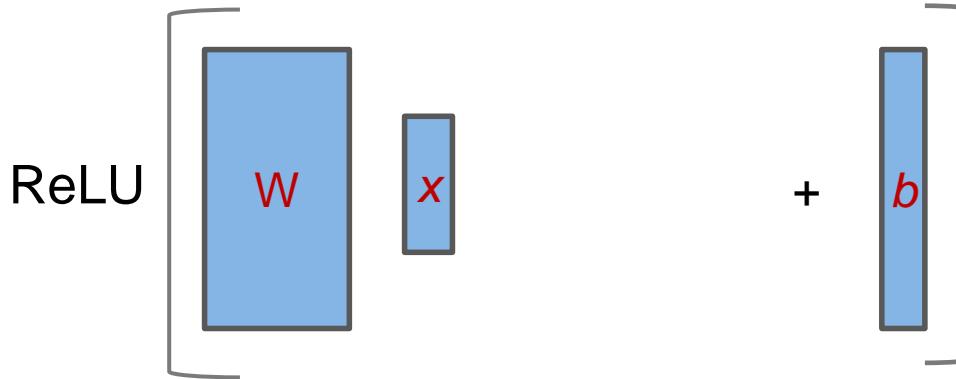
**ReLU**: Activate with elementwise  
rectified linear function

$$\text{ReLU}(x) = \begin{cases} 0, & x \leq 0 \\ x, & x > 0 \end{cases}$$



# Step 2: Code and Variable Assignment

$$h = \text{ReLU}(Wx + b)$$

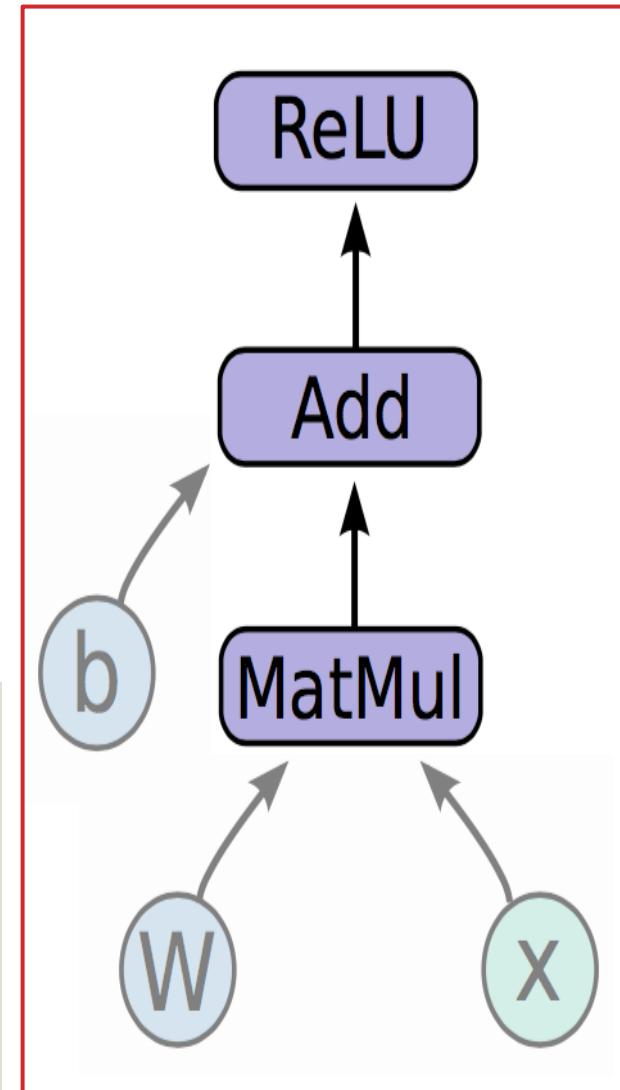


```
import tensorflow as tf

b = tf.Variable(tf.zeros((100,)))
W = tf.Variable(tf.random_uniform((784, 100), -1, 1))

x = tf.placeholder(tf.float32, (1, 784))

h = tf.nn.relu(tf.matmul(x, W) + b)
```



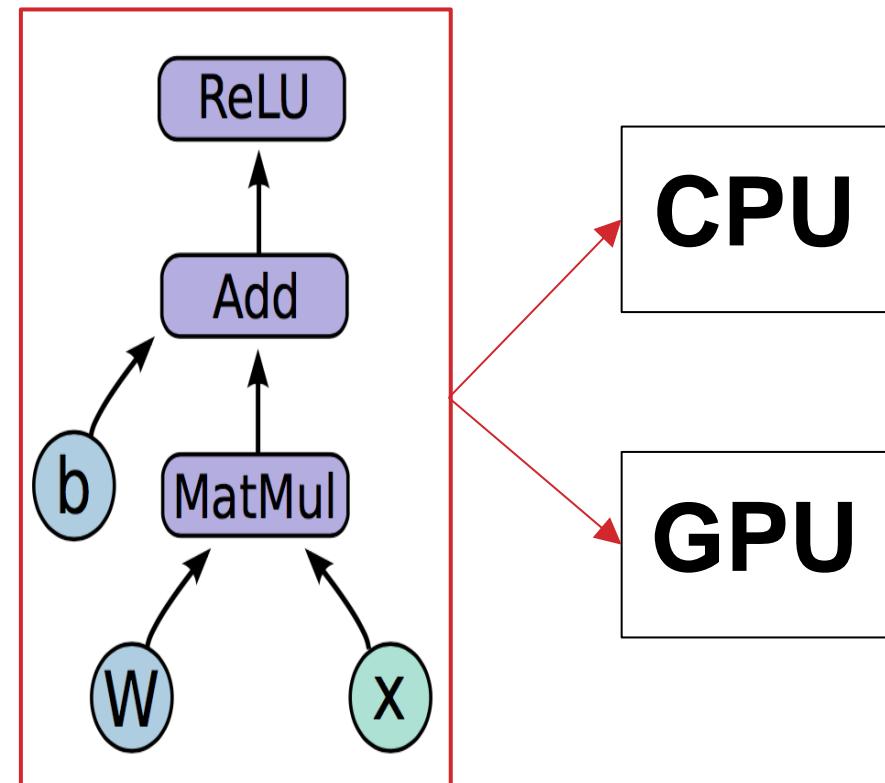
# Step 3: Running the graph

Deploy graph with a **session**:

a binding to a particular

execution context

- e.g. CPU, GPU



# End-to-end

- So far:
  1. Built a **graph** using **variables** and **placeholders**
  2. **Assign variable bindings**
  3. Deploy the graph onto a **session**, i.e., **execution environment**
- Next: train model
  - Define loss function
  - Compute gradients

# Defining loss

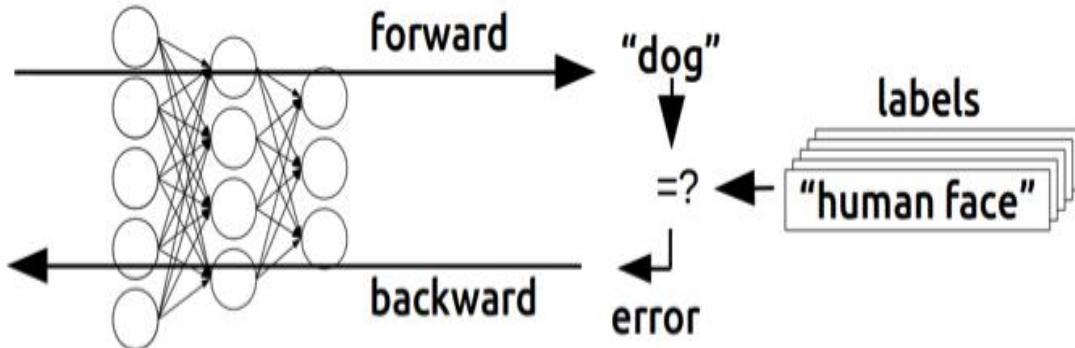
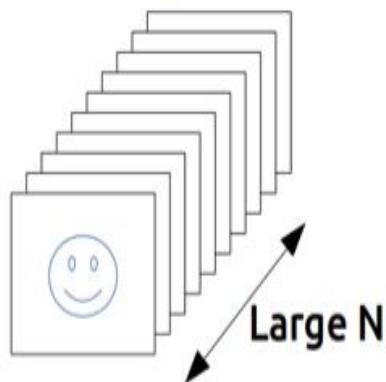
- Use **placeholder** for **labels**
- Build loss node using **labels** and **prediction**

```
prediction = tf.nn.softmax(...) #Output of neural network
label = tf.placeholder(tf.float32, [100, 10])

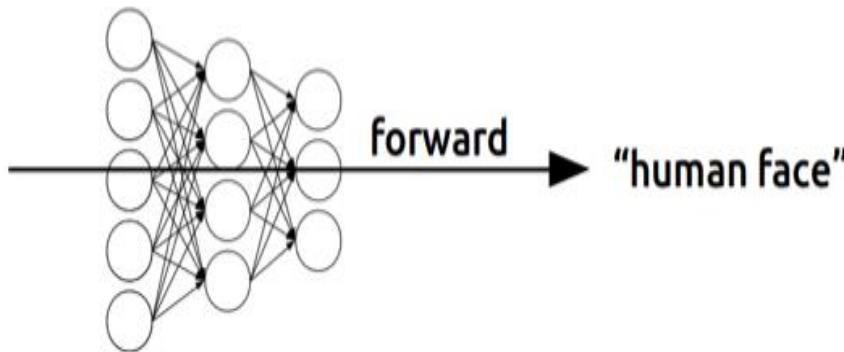
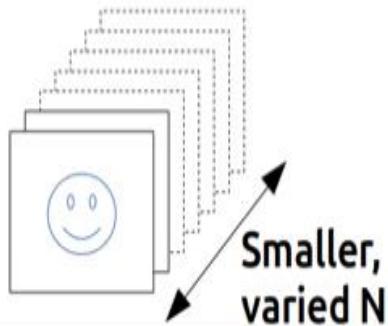
cross_entropy = -tf.reduce_sum(label * tf.log(prediction), axis=1)
```

# Learning: Backpropagation

## Training



## Inference



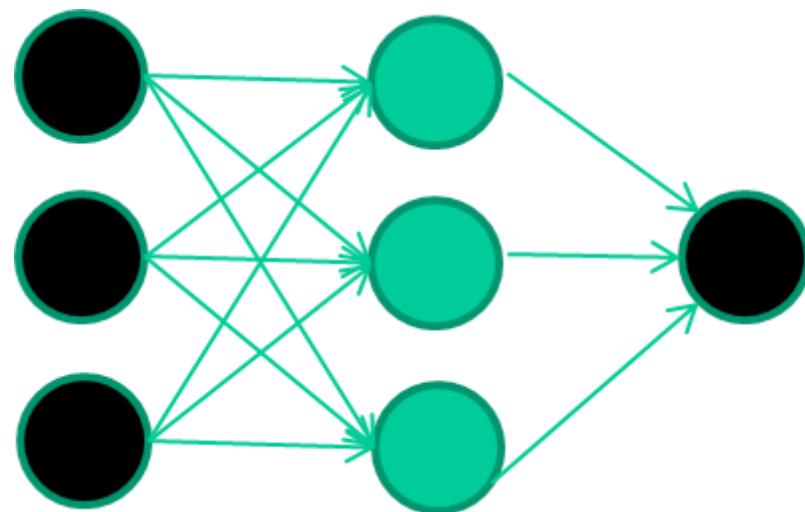
# Example: Backpropagation

*A dataset*

*Fields*            *class*

1.4	2.7	1.9	0
3.8	3.4	3.2	0
6.4	2.8	1.7	1
4.1	0.1	0.2	0

etc ...

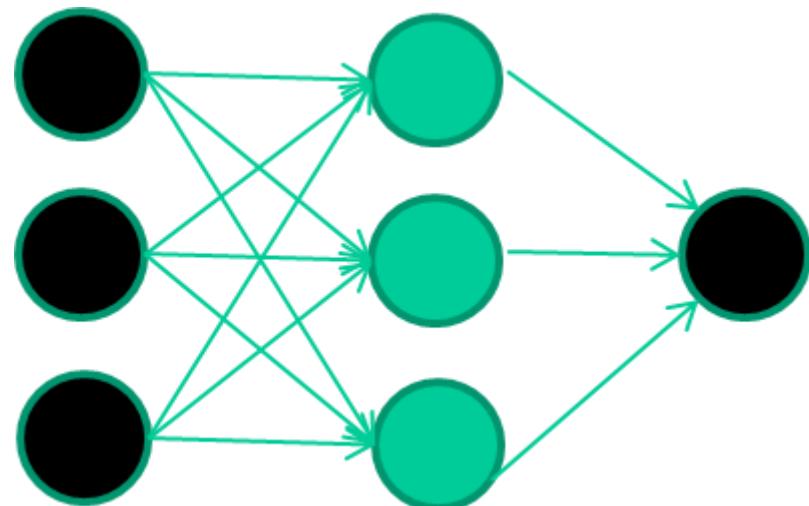


# Example: Backpropagation

*Training the neural network*

*Fields*              *class*

1.4	2.7	1.9	0
3.8	3.4	3.2	0
6.4	2.8	1.7	1
4.1	0.1	0.2	0
etc ...			

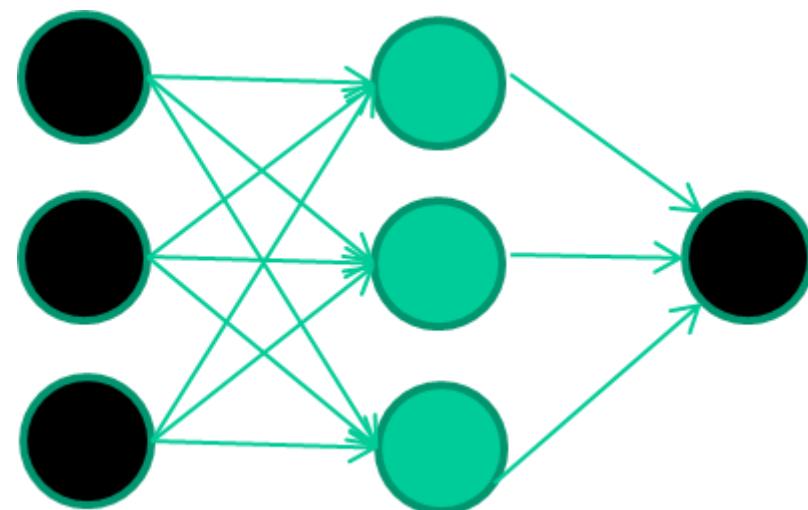


# Example: Backpropagation

*Training data*

<i>Fields</i>	<i>class</i>		
1.4	2.7	1.9	0
3.8	3.4	3.2	0
6.4	2.8	1.7	1
4.1	0.1	0.2	0
etc ...			

Initialise with random weights

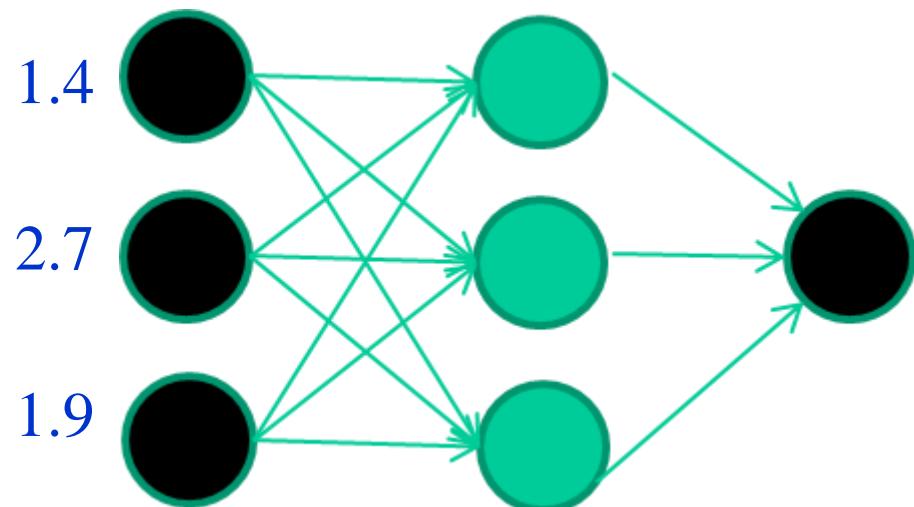


# Example: Backpropagation

*Training data*

<i>Fields</i>	<i>class</i>		
1.4	2.7	1.9	0
3.8	3.4	3.2	0
6.4	2.8	1.7	1
4.1	0.1	0.2	0
etc ...			

Present a training pattern

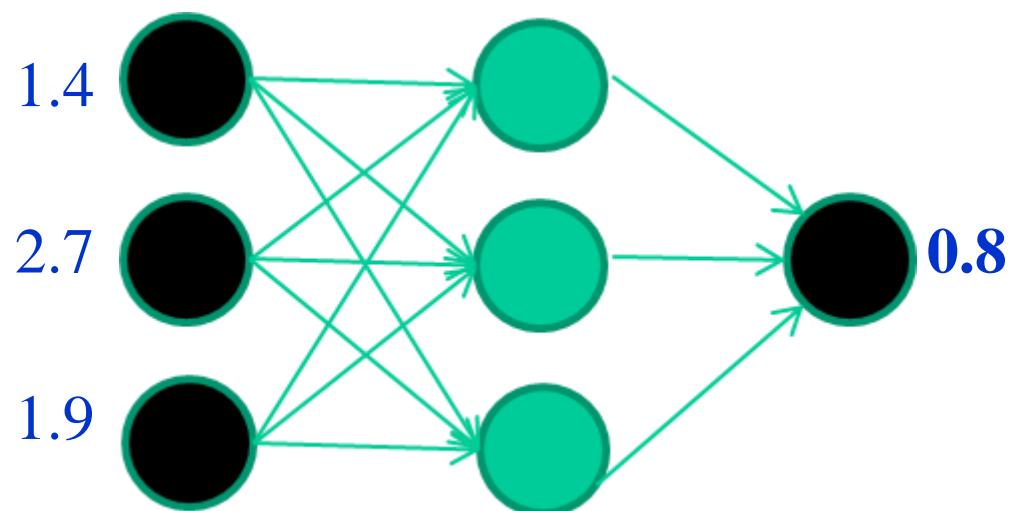


# Example: Backpropagation

*Training data*

<i>Fields</i>	<i>class</i>		
1.4	2.7	1.9	0
3.8	3.4	3.2	0
6.4	2.8	1.7	1
4.1	0.1	0.2	0
etc ...			

Feed it through to get output

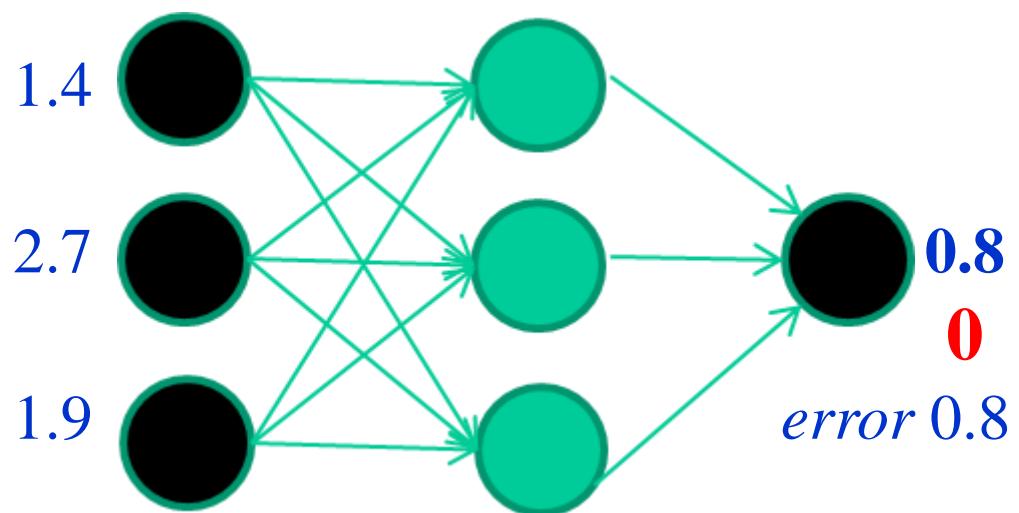


# Example: Backpropagation

Training data

Fields	class		
1.4	2.7	1.9	0
3.8	3.4	3.2	0
6.4	2.8	1.7	1
4.1	0.1	0.2	0
etc ...			

Compare with target output

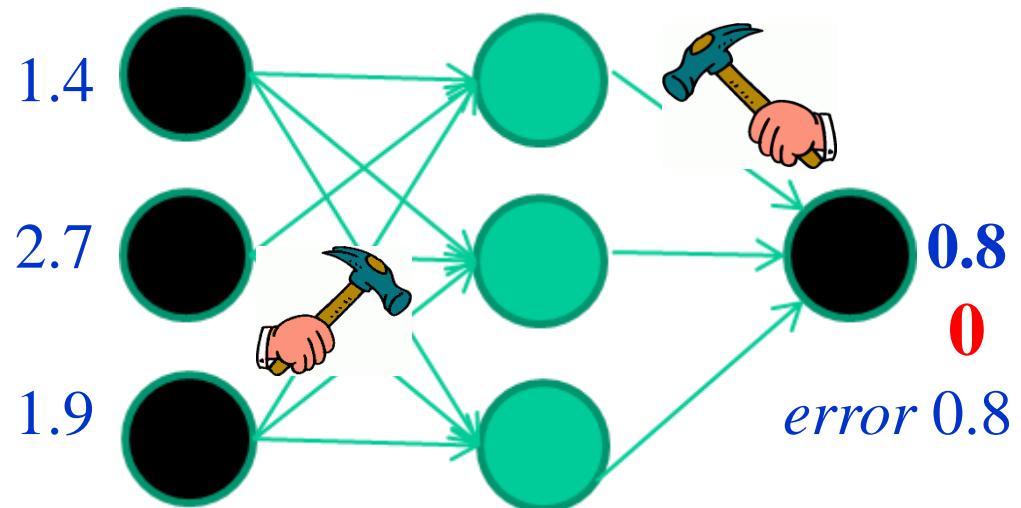


# Example: Backpropagation

*Training data*

<i>Fields</i>	<i>class</i>		
1.4	2.7	1.9	0
3.8	3.4	3.2	0
6.4	2.8	1.7	1
4.1	0.1	0.2	0
etc ...			

Adjust weights based on error



# Example: Backpropagation

*Training data*

*Fields*                    *class*

1.4 2.7 1.9 0

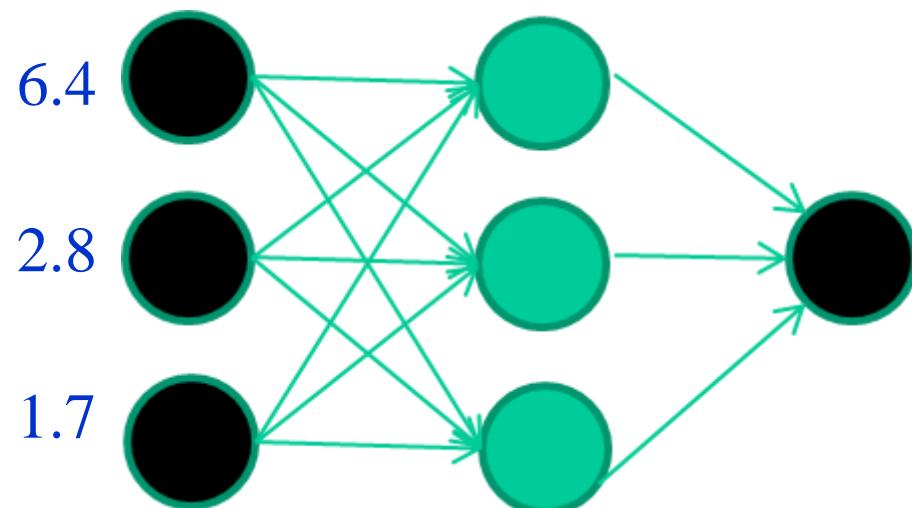
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

Present a training pattern



# Example: Backpropagation

*Training data*

*Fields*                    *class*

1.4 2.7 1.9 0

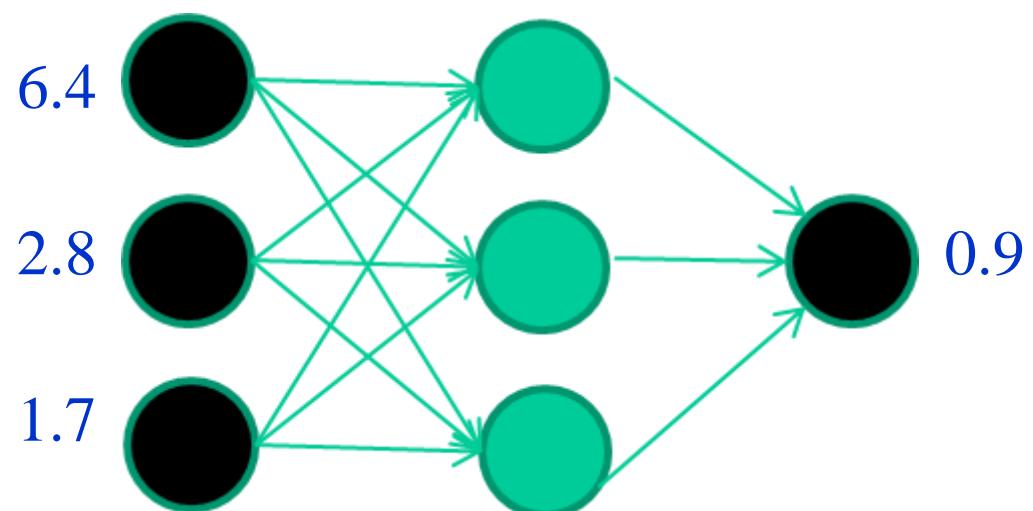
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

Feed it through to get output



# Example: Backpropagation

*Training data*

*Fields*              *class*

1.4 2.7 1.9 0

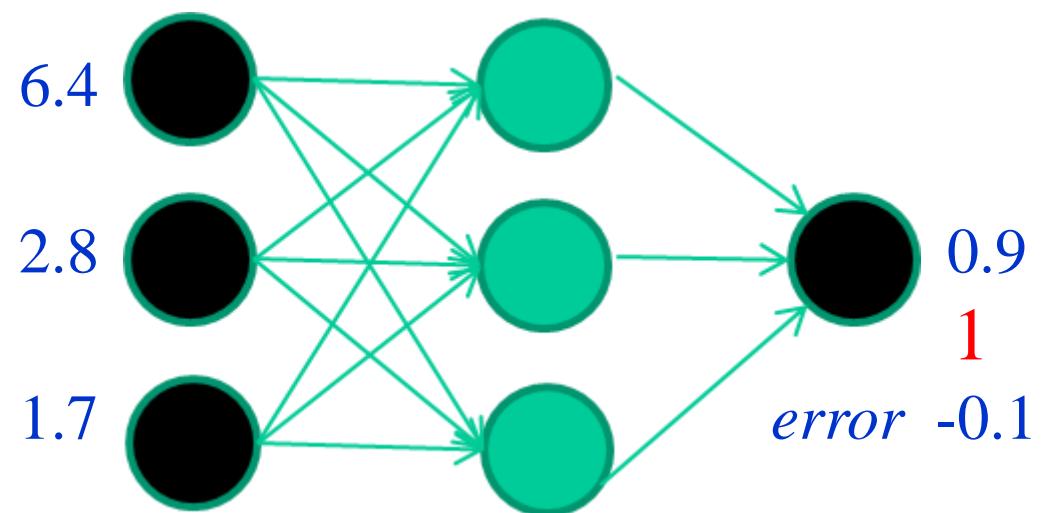
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

Compare with target output



# Example: Backpropagation

*Training data*

*Fields*                    *class*

1.4 2.7 1.9 0

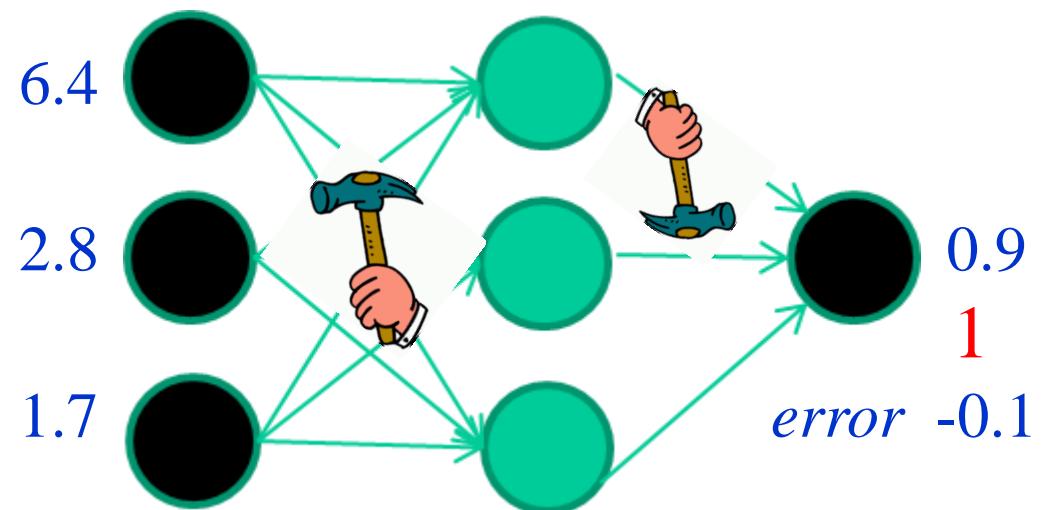
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

Adjust weights based on error



# Example: Backpropagation

*Training data*

*Fields*                    *class*

1.4 2.7 1.9 0

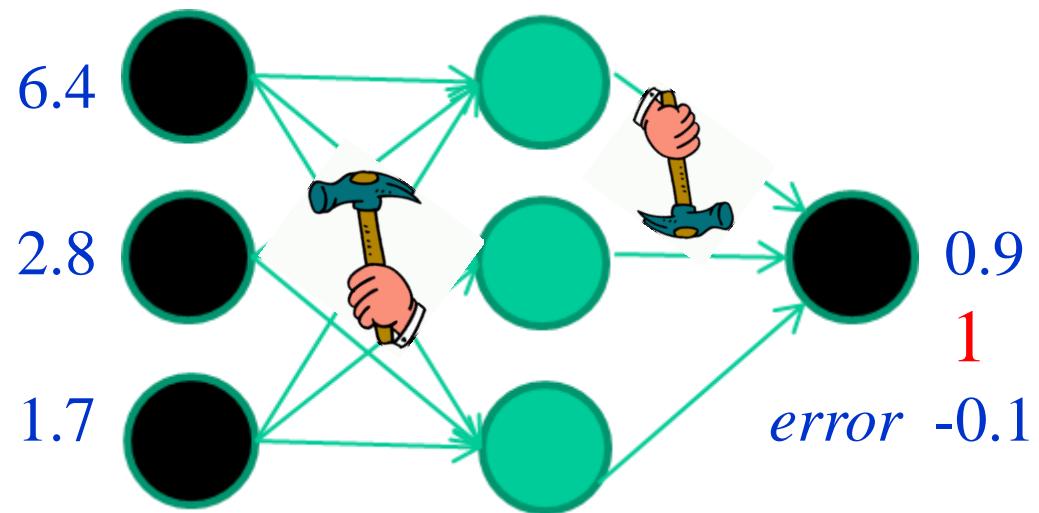
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

And so on ....



Repeat this thousands, maybe millions of times – each time taking a random training instance, and making slight weight adjustments

*Algorithms for weight adjustment are designed to make changes that will reduce the error*



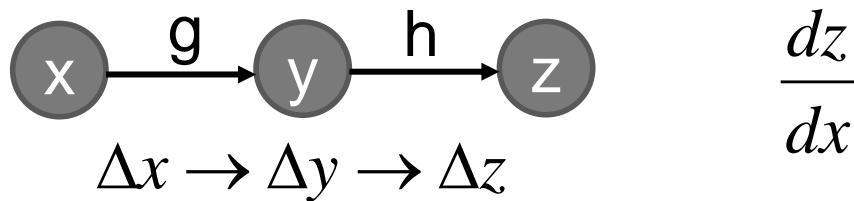
# Backpropagation and Gradients

- Backpropagation: an efficient way to compute the gradient
- Prerequisite
  - Backpropagation for feedforward net:
    - [http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS\\_2015\\_2/Lecture/DNN%20backprop.ecm.mp4/index.html](http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS_2015_2/Lecture/DNN%20backprop.ecm.mp4/index.html)
    - Simple version: <https://www.youtube.com/watch?v=ibJpTrp5mcE>
  - Backpropagation through time for RNN:  
[http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS\\_2015\\_2/Lecture/RNN%20training%20\(v6\).ecm.mp4/index.html](http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS_2015_2/Lecture/RNN%20training%20(v6).ecm.mp4/index.html)
- Understanding backpropagation by computational graph
  - Tensorflow, Theano, CNTK, etc.

# Review: Chain Rule

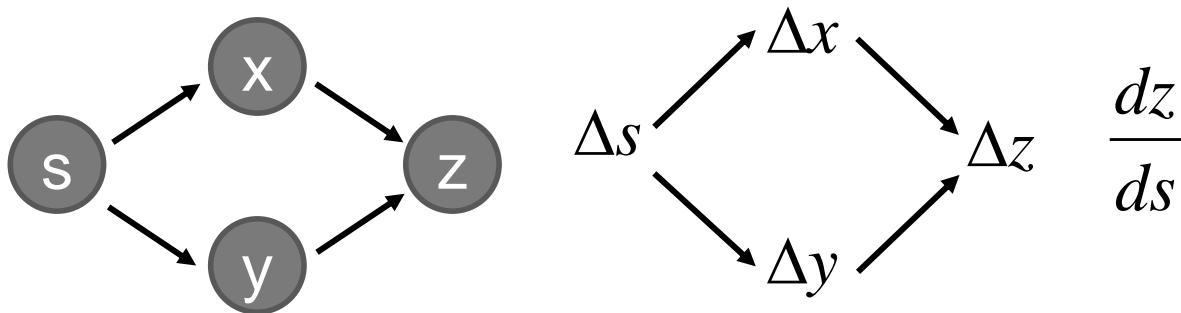
Case  
1

$$z = f(x) \rightarrow y = g(x) \quad z = h(y)$$



Case  
2

$$z = f(s) \rightarrow x = g(s) \quad y = h(s) \quad z = k(x, y)$$



# Computational Graph

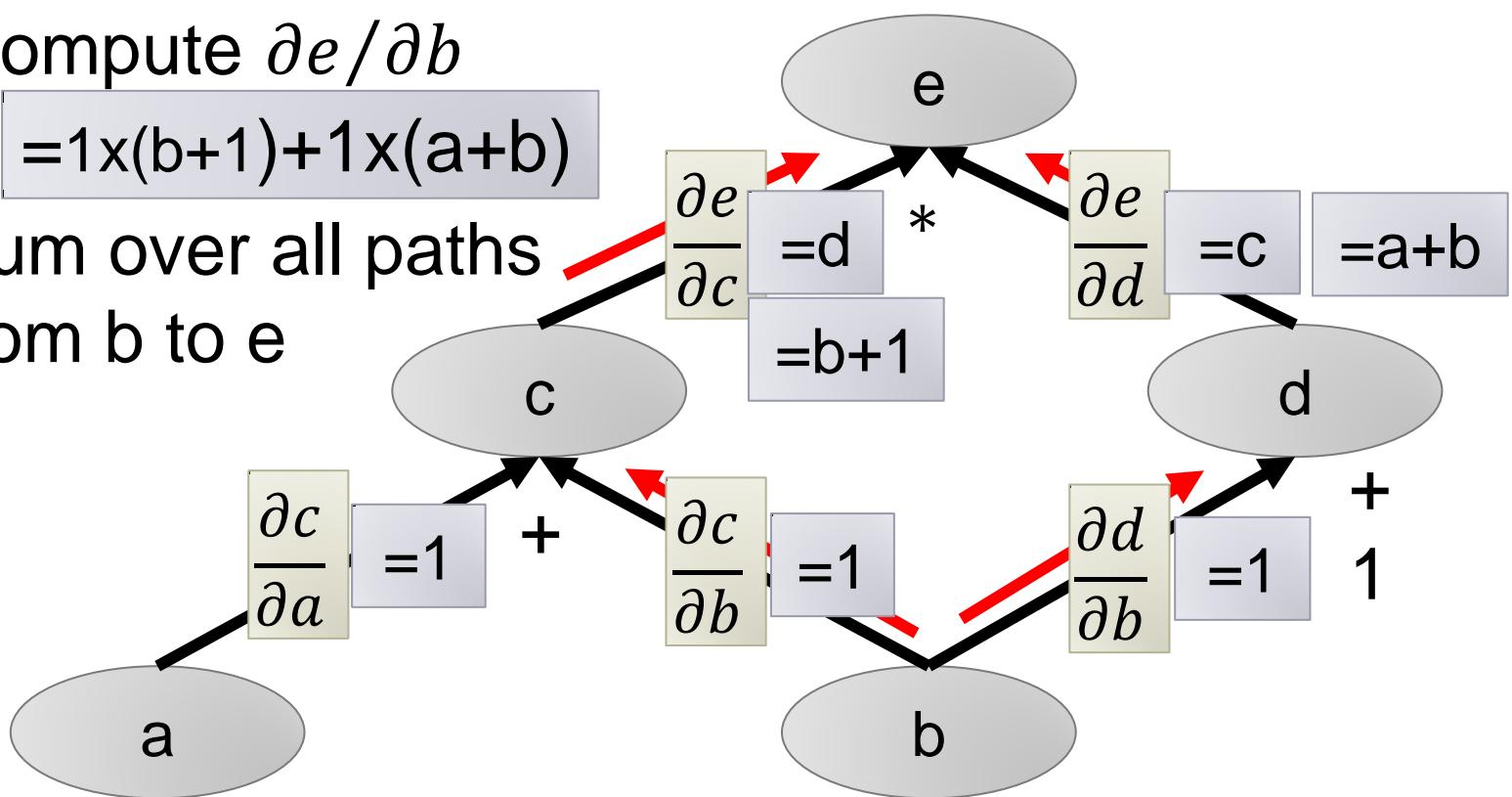
- Example:  $e = (a+b) * (b+1)$

$$\frac{\partial e}{\partial b} = \frac{\partial e}{\partial c} \frac{\partial c}{\partial b} + \frac{\partial e}{\partial d} \frac{\partial d}{\partial b}$$

Compute  $\frac{\partial e}{\partial b}$

$$=1 \times (b+1) + 1 \times (a+b)$$

Sum over all paths  
from b to e

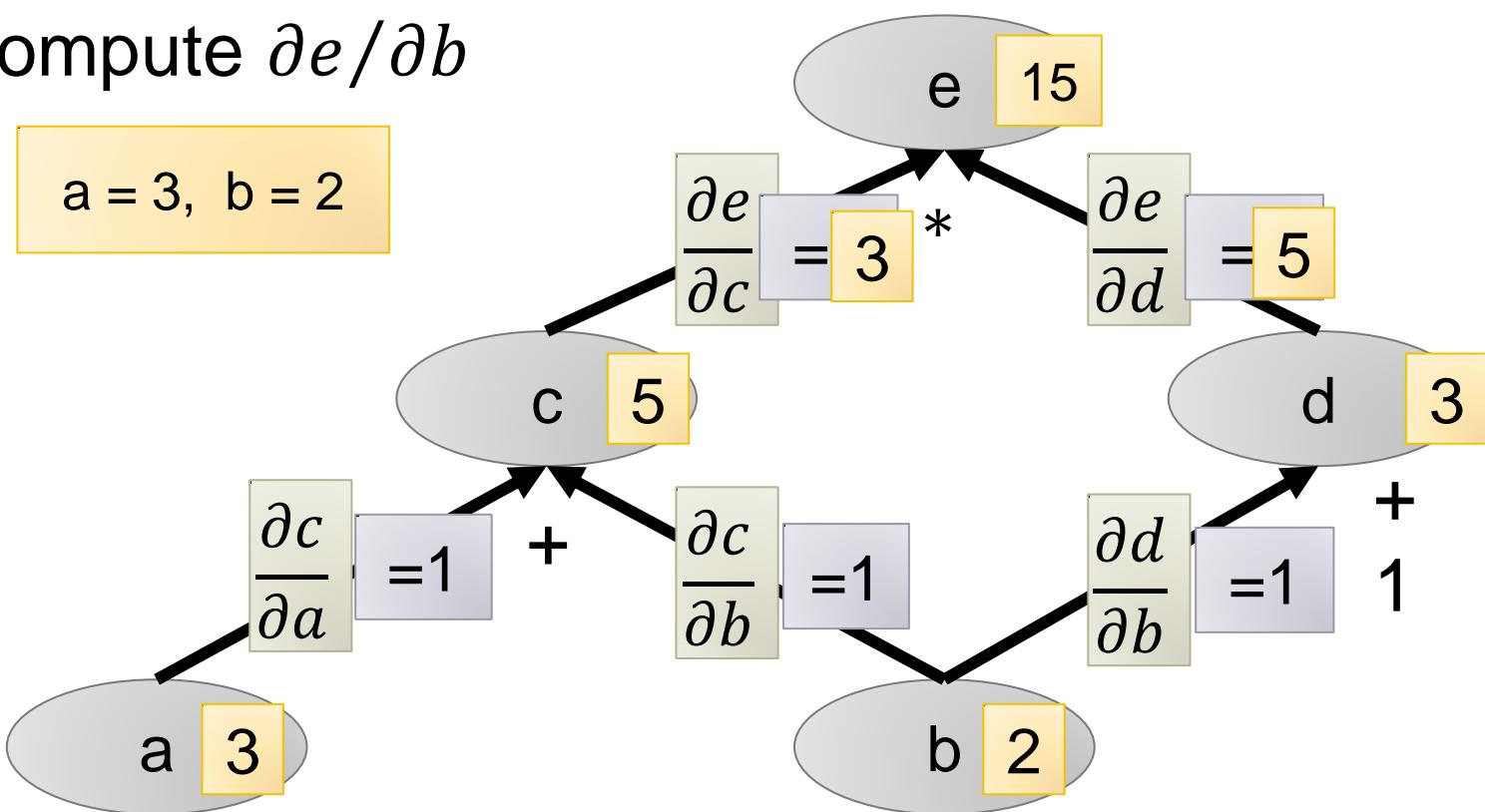


# Computational Graph

- Example:  $e = (a+b) * (b+1)$

$$\frac{\partial e}{\partial b} = \frac{\partial e}{\partial c} \frac{\partial c}{\partial b} + \frac{\partial e}{\partial d} \frac{\partial d}{\partial b}$$

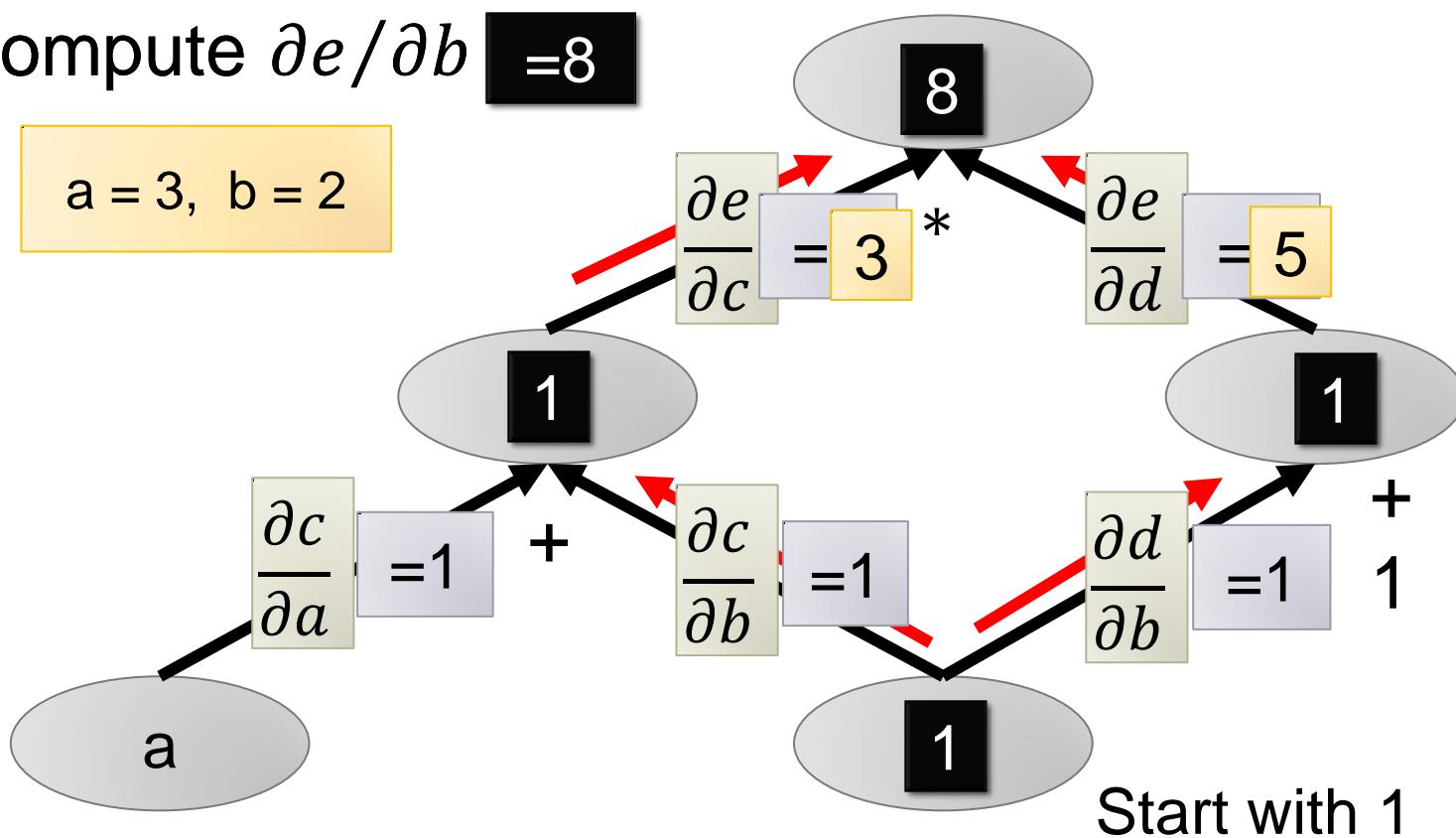
Compute  $\frac{\partial e}{\partial b}$



# Computational Graph

- Example:  $e = (a+b) * (b+1)$

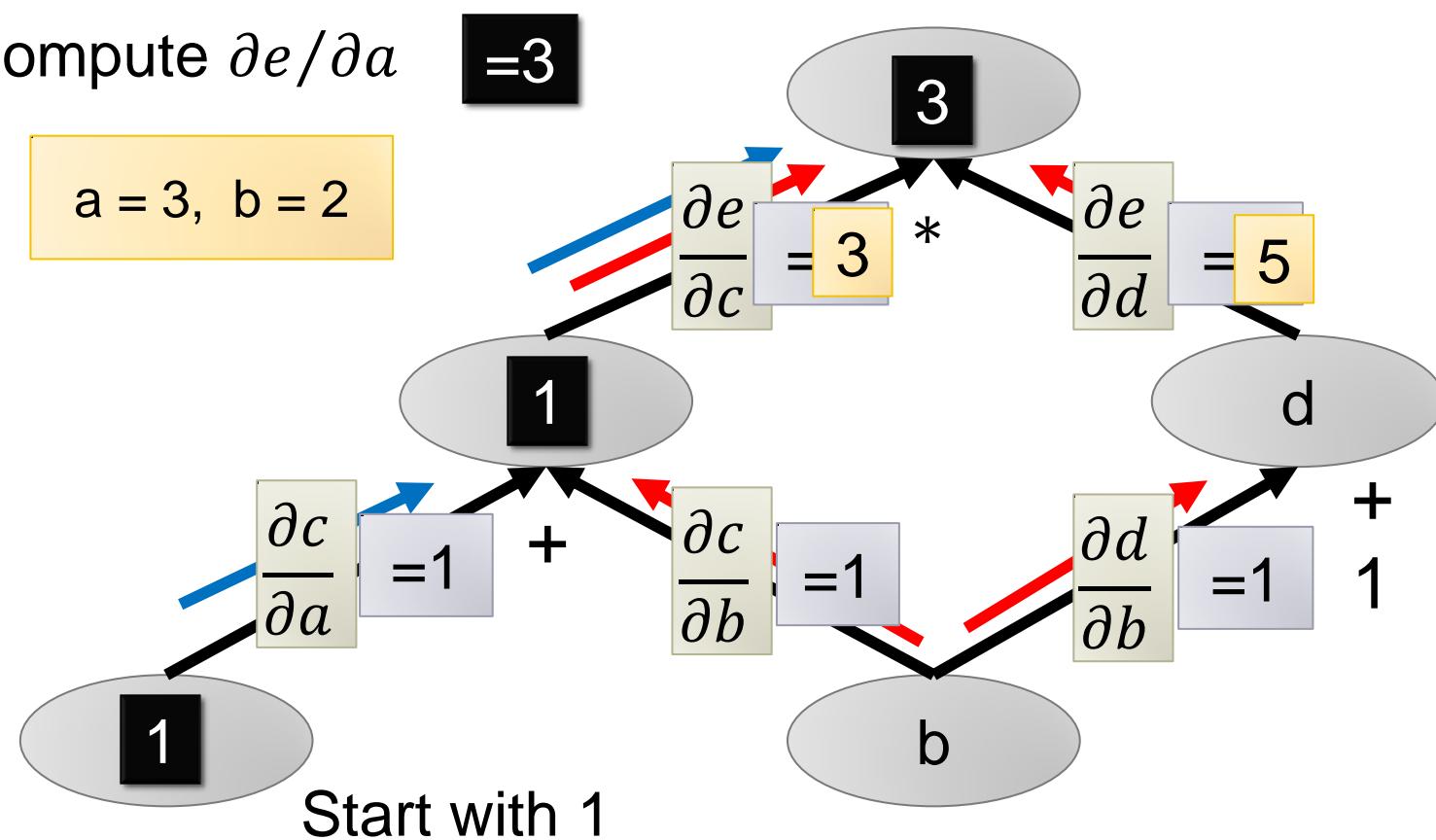
Compute  $\partial e / \partial b = 8$



# Computational Graph

- Example:  $e = (a+b) * (b+1)$

Compute  $\partial e / \partial a$

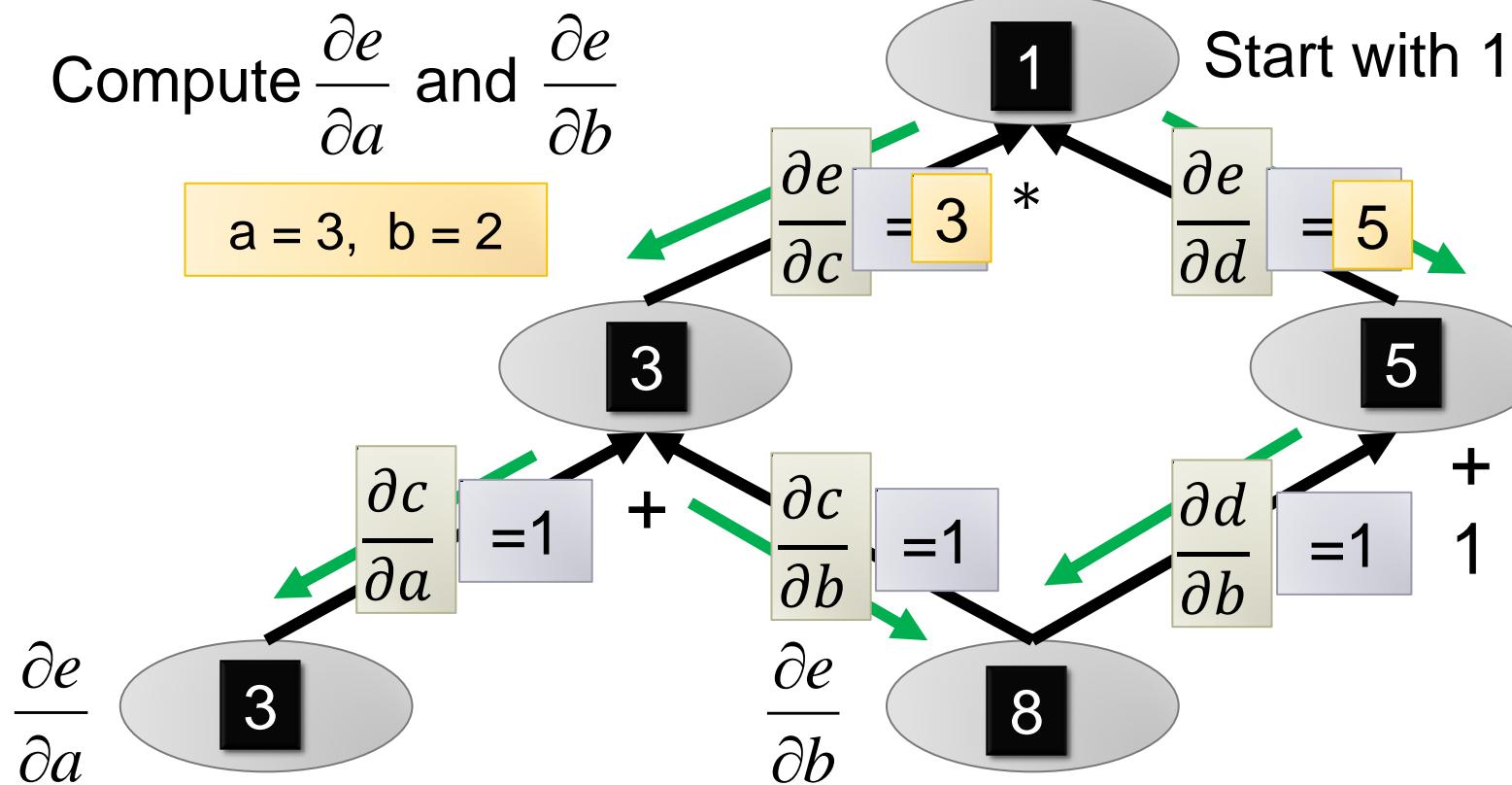


# Computational Graph

- Example:  $e = (a+b) * (b+1)$

Reverse mode

What is the benefit?



# Computational Graph

- **Parameter sharing:** the same parameters appearing in different nodes

$$y = xe^{x^2} \quad \frac{\partial y}{\partial x} = ? \quad e^{x^2} + x \cdot e^{x^2} \cdot 2x$$

