

Deep Learning 工具介绍

Introduction to Tensorflow, PyTorch and Caffe

Johnson
Sep. 2nd, 2018

Outline

- Overview of Caffe
- Introduction to Tensorflow/Keras
- Introduction to PyTorch
- Comparisons of Tensorflow and Pytorch
- MNIST Classification with logistic regression using TensorFlow
- Sentiment Analysis with RNN using PyTorch

Caffe

- What is Caffe?
- Core written in C++
- Has Python and Matlab bindings
- Good for training and fine-tuning feedforward classification models
- Often no need to write code
- Not used as much in research anymore

Caffe

- Training/Fine-tuning

No need to write code!

- Convert data (run a script)
- Define network (edit prototxt)
- Define solver (edit prototxt)
- Train (with pretrained weights, run a script)

Caffe

- Training/Fine-tuning

Convert data:

- DataLayer reading from LMDB is the easiest
- Create LMDB using convert_imageset
- Need text file where each line is: [path/to/image.jpeg]
[label]
- Create HDF5 file using h5py

Caffe

- Training/Fine-tuning

Define network (prototxt)

```
name: "LogisticRegressionNet"
layers {
    top: "data"
    top: "label"
    name: "data"
    type: HDF5_DATA
    hdf5_data_param {
        source: "examples/hdf5_classification/data/train.txt"
        batch_size: 10
    }
    include {
        phase: TRAIN
    }
}
layers {
    bottom: "data"
    top: "fc1"
    name: "fc1"
    type: INNER_PRODUCT
    blobs_lr: 1
    blobs_lr: 2
    weight_decay: 1
    weight_decay: 0
}
inner_product_param {
    num_output: 2
    weight_filler {
        type: "gaussian"
        std: 0.01
    }
    bias_filler {
        type: "constant"
        value: 0
    }
}
layers {
    bottom: "fc1"
    bottom: "label"
    top: "loss"
    name: "loss"
    type: SOFTMAX_LOSS
}
```

Caffe

- Training/Fine-tuning

Define network (prototxt)

resnet-152: 6775 lines!

```
1 name: "ResNet-152"
2 input: "data"
3 input_dim: 1
4 input_dim: 3
5 input_dim: 224
6 input_dim: 224
7
8 layer {
9     bottom: "data"
10    top: "conv1"
11    name: "conv1"
12    type: "Convolution"
13    convolution_param {
14        num_output: 64
15        kernel_size: 7
16        pad: 3
17        stride: 2
18        bias_term: false
19    }
20}
21
22 layer {
23    bottom: "conv1"
24    top: "conv1"
25    name: "bn_conv1"
26    type: "BatchNorm"
27    batch_norm_param {
28        use_global_stats: true
29    }
30}
```

```
6747 layer {
6748     bottom: "res5c"
6749     top: "pools"
6750     name: "pools"
6751     type: "Pooling"
6752     pooling_param {
6753         kernel_size: 7
6754         stride: 1
6755         pool: AVE
6756     }
6757 }
6758
6759 layer {
6760     bottom: "pools"
6761     top: "fc1000"
6762     name: "fc1000"
6763     type: "InnerProduct"
6764     inner_product_param {
6765         num_output: 1000
6766     }
6767 }
6768
6769 layer {
6770     bottom: "fc1000"
6771     top: "prob"
6772     name: "prob"
6773     type: "Softmax"
6774 }
```

Caffe

- Training/Fine-tuning

Define solver (prototxt)

- Write a prototxt defining SolverParameter
- If fining-tuning, copy existing solver.prototxt file
 - Change net to be your net
 - Reduce base learning rate
 - Change max_iter

```
1 net: "models/bvlc_alexnet/train_val.prototxt"
2 test_iter: 1000
3 test_interval: 1000
4 base_lr: 0.01
5 lr_policy: "step"
6 gamma: 0.1
7 stepsize: 100000
8 display: 20
9 max_iter: 450000
10 momentum: 0.9
11 weight_decay: 0.0005
12 snapshot: 10000
13 snapshot_prefix: "models/bvlc_alexnet/caffe_alexnet_train"
14 solver_mode: GPU
```

Caffe

- Training/Fine-tuning

Training!

```
./build/tools/caffe train \
-gpu 0 \
-model path/to/trainval.prototxt \
-solver path/to/solver.prototxt \
-weights path/to/pretrained_weights.caffemodel
```

Caffe

- Training/Fine-tuning

Training!

```
./build/tools/caffe train \
  -gpu 0 \
  -model path/to/trainval.prototxt \
  -solver path/to/solver.prototxt \
  -weights path/to/pretrained_weights.caffemodel
```

|
-gpu -1 for CPU-only
-gpu all for multi-gpu

TensorFlow

- What's TensorFlow?
 - Open source library for numeric and symbolic computation with GPU support
 - Developed by the Google Brain Team for machine learning related research

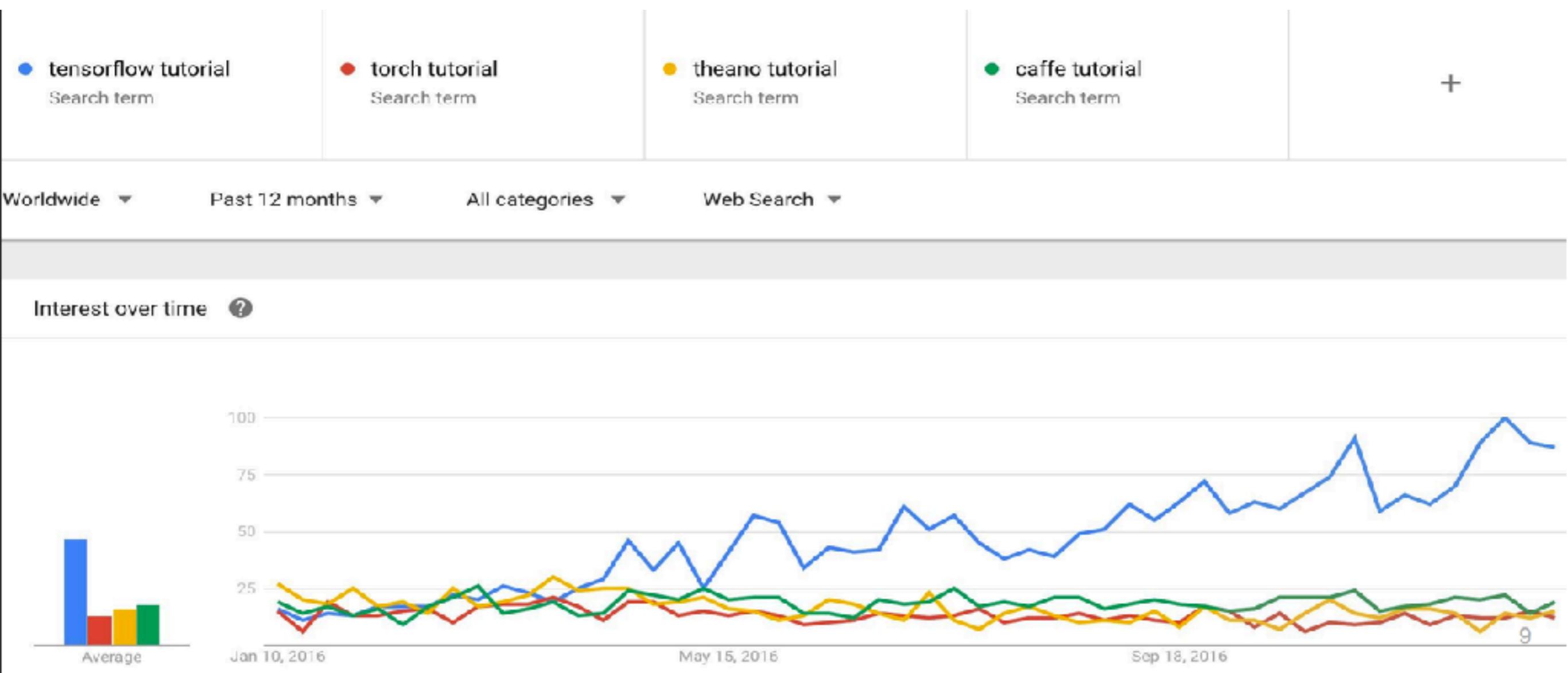
We will compare TensorFlow with Numpy to highlight its key features

TensorFlow

- Besides TensorFlow
- Caffe/Caffe 2 (Berkeley/Facebook)
- Torch (NYU)
- Theano (University of Montreal)
- CNTK (Microsoft)
- Paddle (Baidu)
- MXNet (Amazon)
- PyTorch (Facebook)

TensorFlow

- Why TensorFlow?



TensorFlow

- Why TensorFlow?
- Python API
- Portability: easy to deploy computations over one or more CPUs/GPUs, with the same API
- Flexibility: easy to extend to mobile devices, including Android, iOS, etc.
- Visualization: TensorBoard is great!
- **Auto-differentiation: autodiff, no need to compute the gradient manually**
- Large community: > 10,000 commits and > 3,000 TF-related repos in 1 year

TensorFlow

- Companies that use TensorFlow
 - Google
 - DeepMind
 - Dropbox
 - Snapchat
 - Uber
 - eBay
 - OpenAI
- ...

TensorFlow

- Goals of this lecture
- Understand TF's computation graph approach
- Explore TF's built-in functions
- Be familiar with the pipeline of a typical machine learning project
(MNIST image classification using TF)

TensorFlow

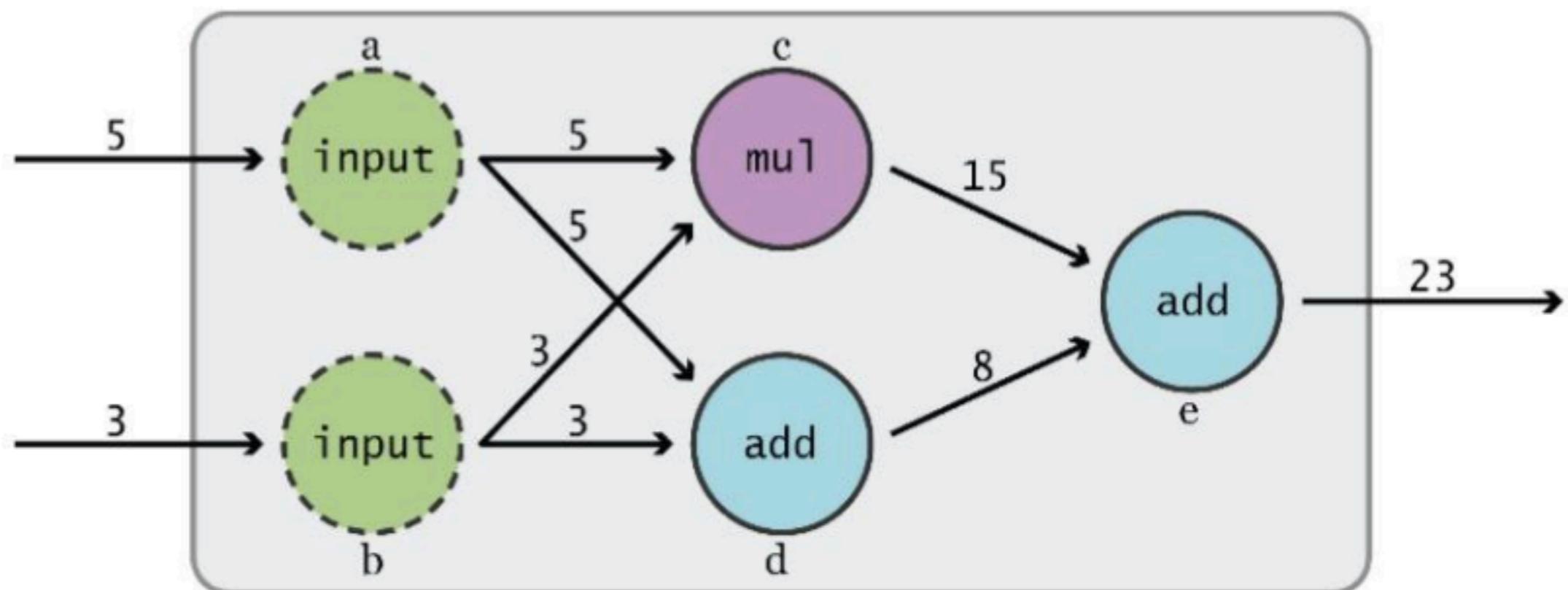
```
import tensorflow as tf
```

TensorFlow

- Even higher level abstraction of TF:
- TF Learn (`tf.contrib.learn`): simplified interface of TensorFlow, similar to scikit-learn
- TF Slim (`tf.contrib.slim`): lightweight library for defining and running complex models in TensorFlow
- Even more: **Keras**

TensorFlow

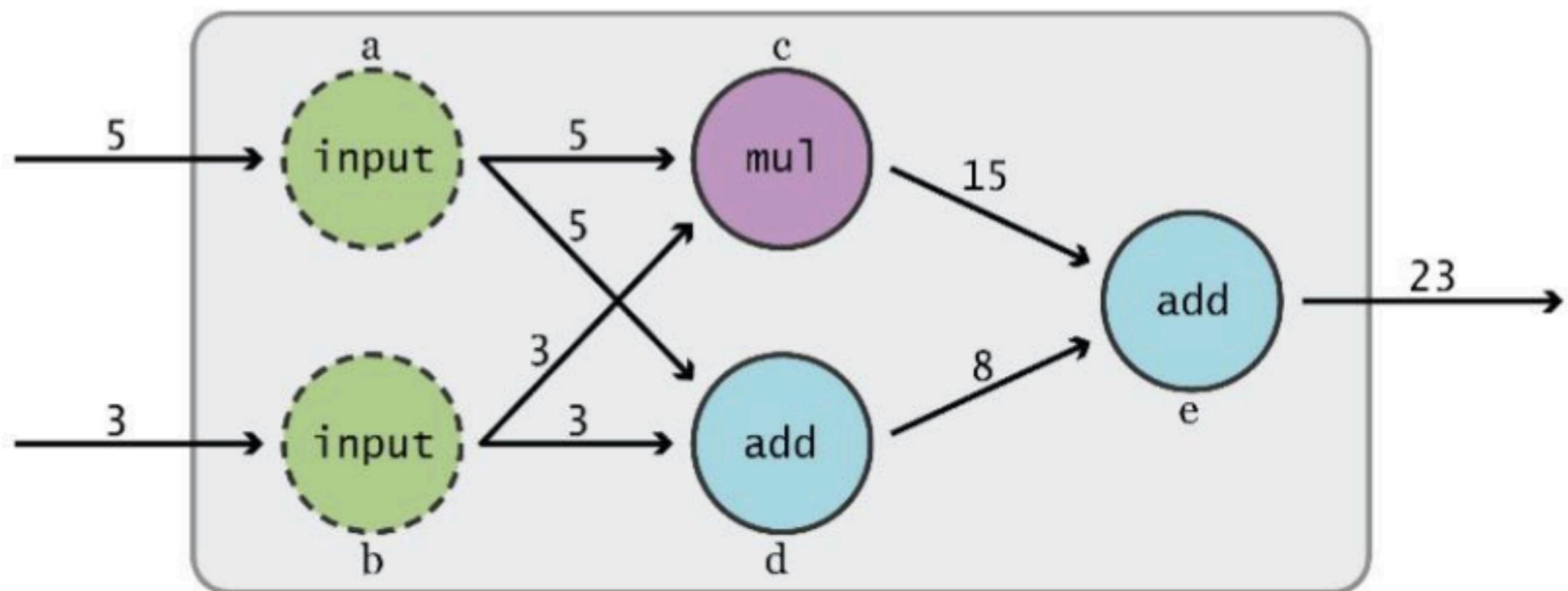
- Data flow graphs



Defining computation \neq Execution of Computation

TensorFlow

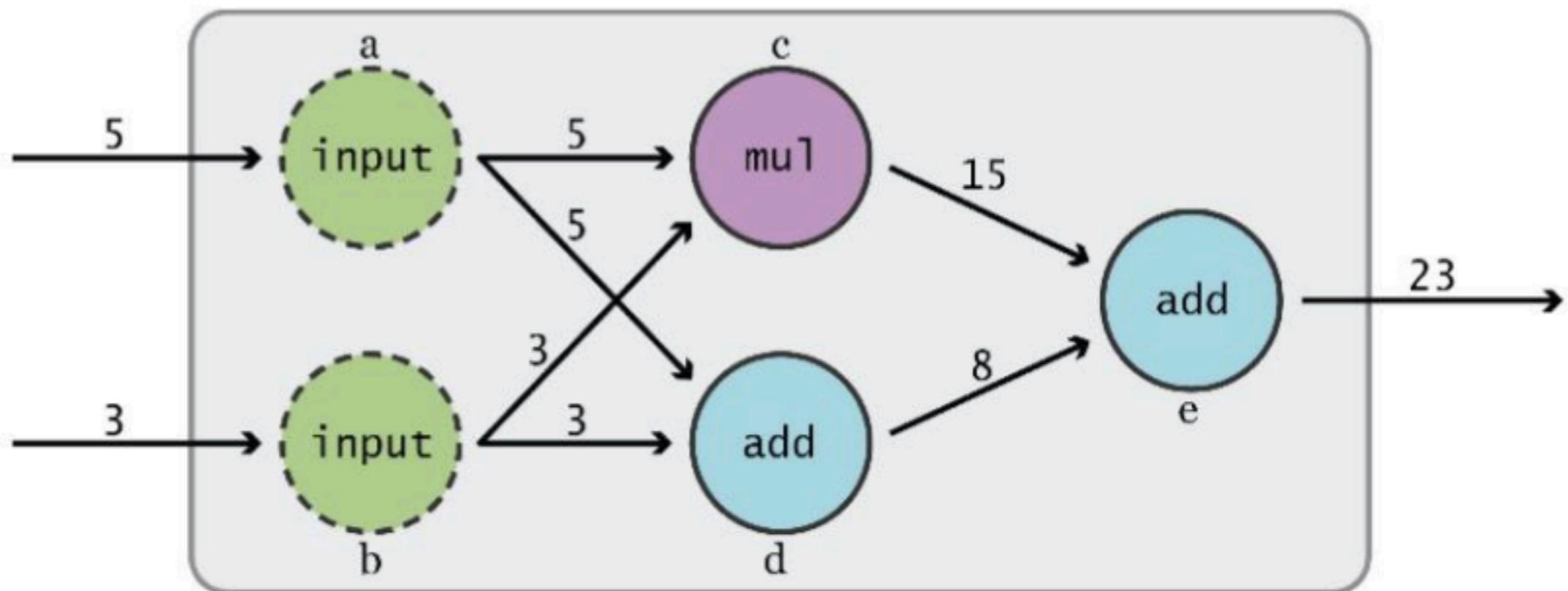
- Data flow graphs



- Typical pipeline in TF:
 - Define the computation graph
 - Run session to execute the computational graph

TensorFlow

- Data flow graphs



- TensorFlow = Tensor + Flow:
 - Tensor = data
 - Flow = operators

Data are transmitted and transformed by operators (op) in the computational graph

TensorFlow

- What is a tensor?
- An n-dimensional array
 - 0-d tensor: scalar (number)
 - 1-d tensor: vector
 - 2-d tensor: matrix
 - ...

TensorFlow

- Numpy vs TensorFlow

Numpy	TensorFlow
<code>a = np.zeros((2,2)); b = np.ones((2,2))</code>	<code>a = tf.zeros((2,2)); b = tf.ones((2,2))</code>
<code>np.sum(b, axis=1)</code>	<code>tf.reduce_sum(b, reduction_indices=[1])</code>
<code>a.shape</code>	<code>a.get_shape()</code>
<code>np.reshape(a, (1, 4))</code>	<code>tf.reshape(a, (1, 4))</code>
<code>5*b+1</code>	<code>5*b+1</code>
<code>np.dot(a, b)</code>	<code>tf.matmul(a, b)</code>
<code>a[1, 1], a[:, 1], a[1, :]</code>	<code>a[1, 1], a[:, 1], a[1, :]</code>

TensorFlow

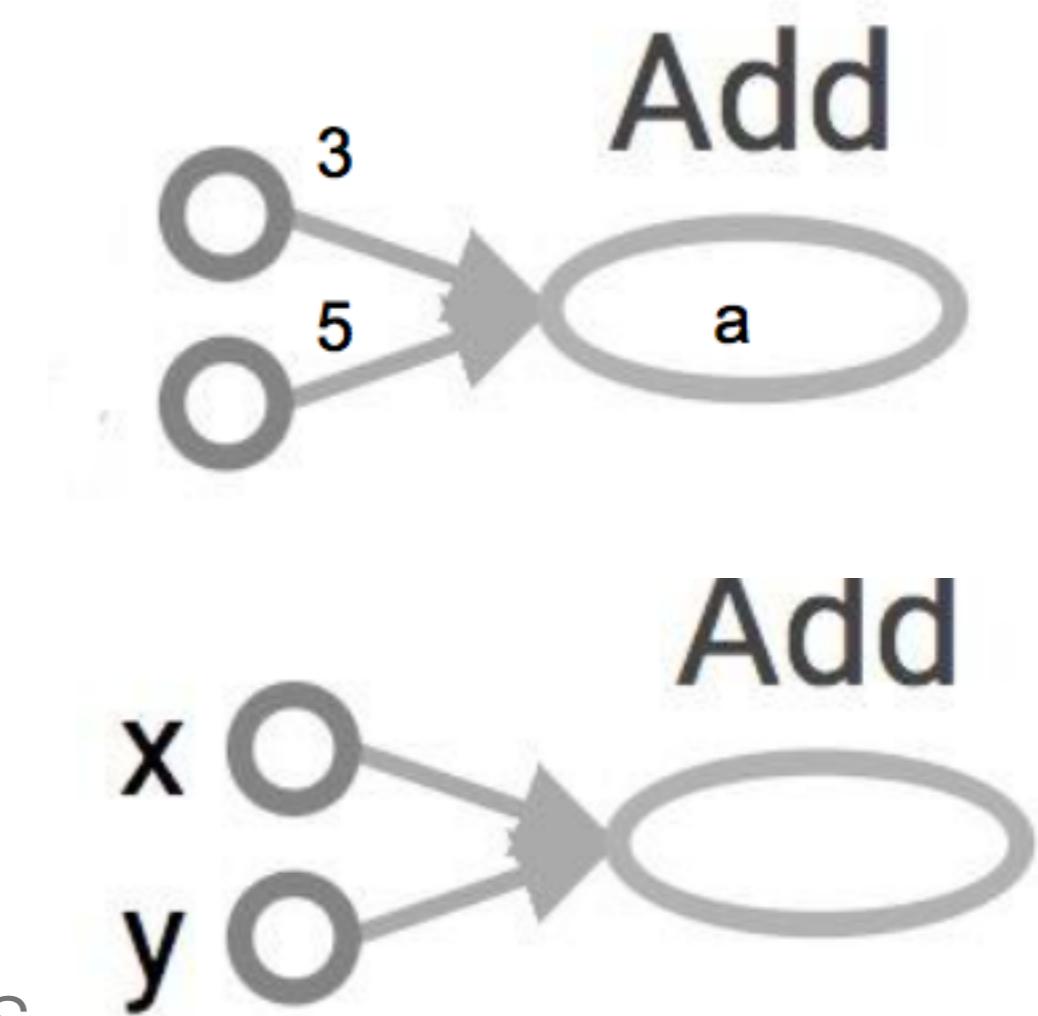
- Data flow graphs

```
import tensorflow as tf  
a = tf.add(3, 5)
```

What is x, y?

$x = 3, y = 5;$

TF will automatically name variables



Visualization from TensorBoard

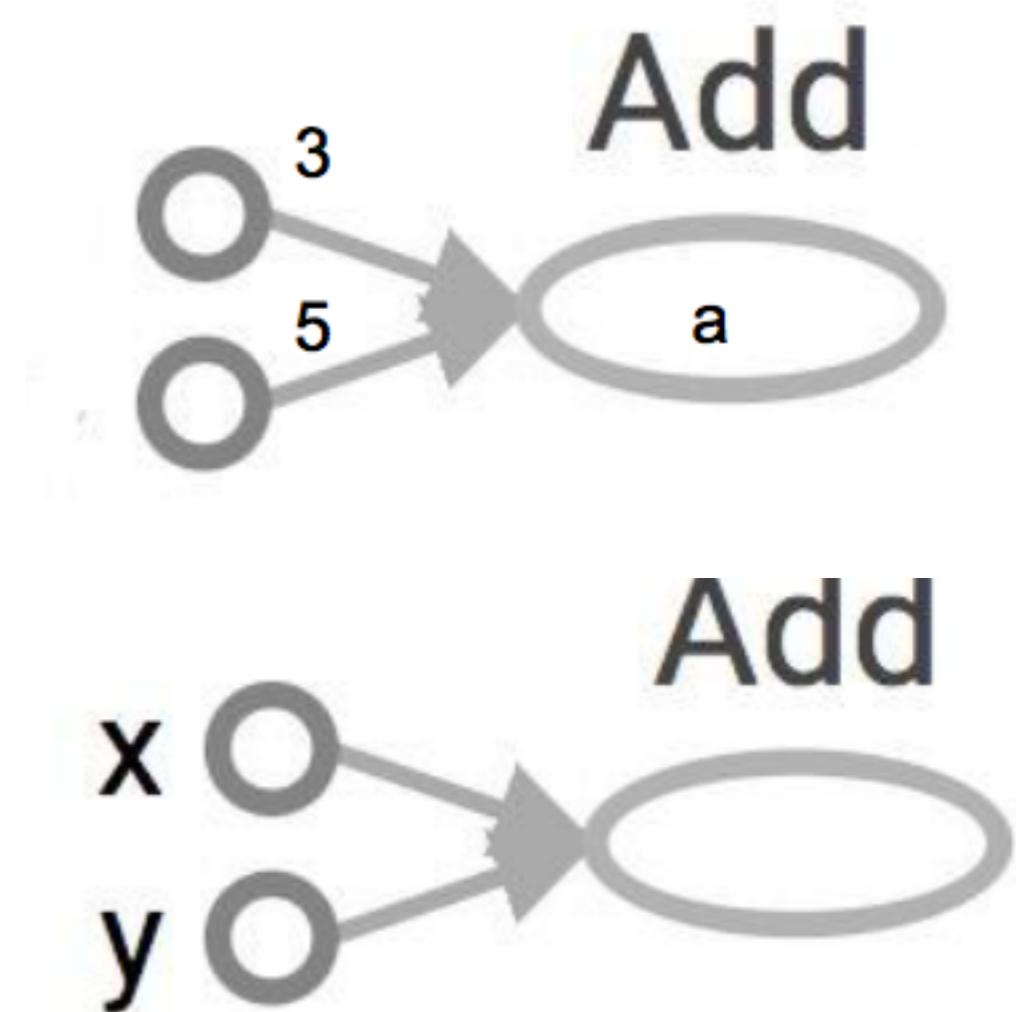
Nodes: operators, variables, or constants

Edges: tensors

TensorFlow

- Data flow graphs

```
import tensorflow as tf  
a = tf.add(3, 5)  
print a
```



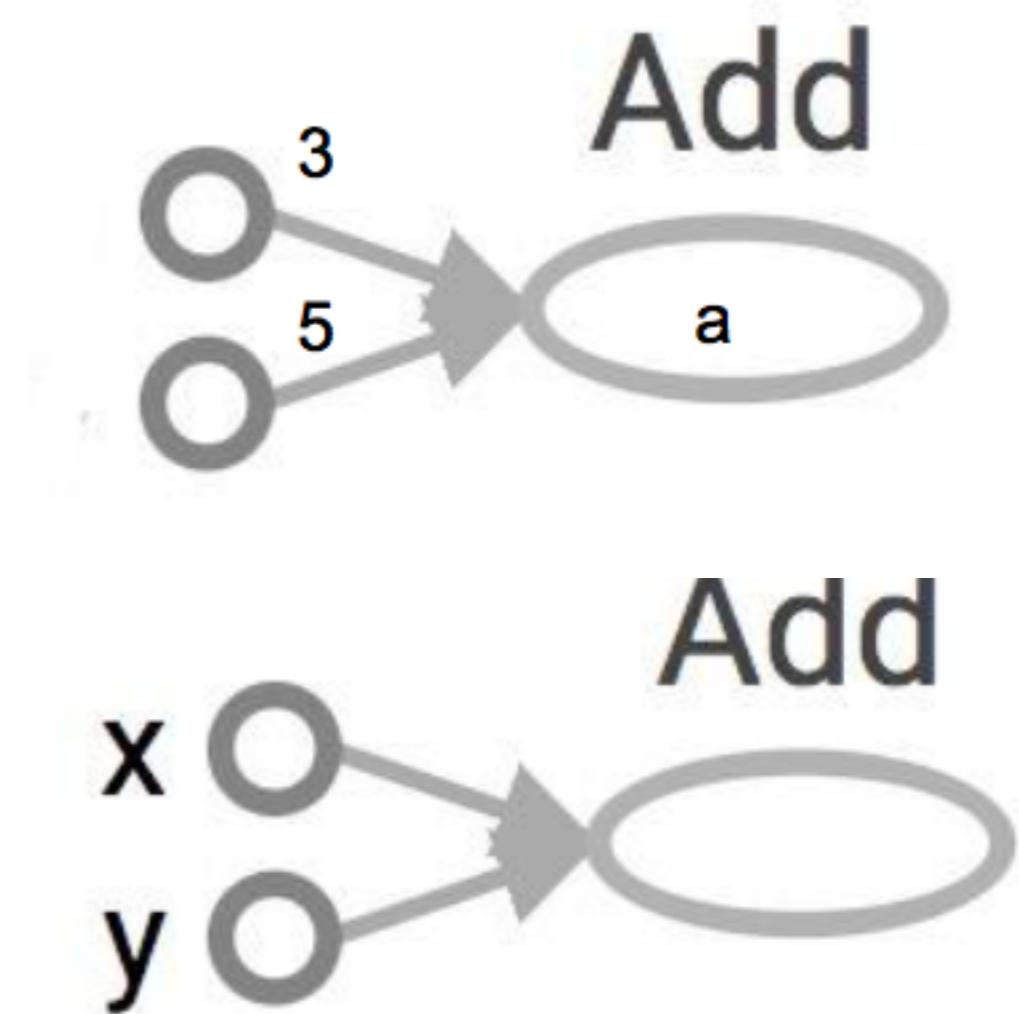
Visualization from TensorBoard

```
>> Tensor("Add:0", shape=(), dtype=int32)  
(Not 8)  
Why?
```

TensorFlow

- Data flow graphs

```
import tensorflow as tf  
a = tf.add(3, 5)  
print a
```



Visualization from TensorBoard

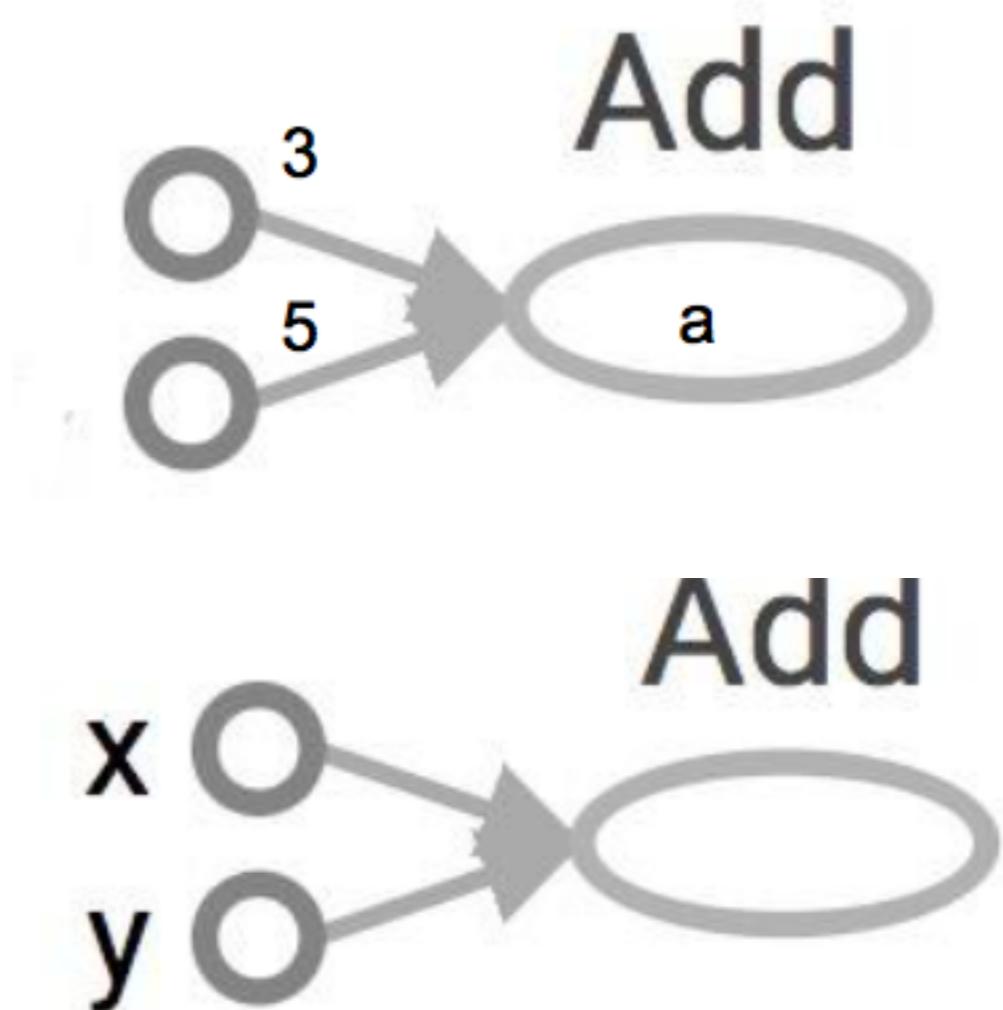
```
>> Tensor("Add:0", shape=(), dtype=int32)  
(Not 8)
```

Symbolic variable! How to get the value of a?

TensorFlow

- Data flow graphs

```
import tensorflow as tf  
a = tf.add(3, 5)  
print a
```



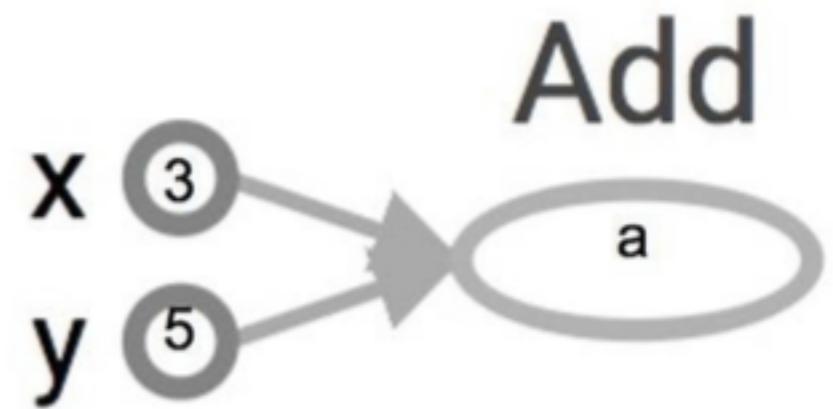
Visualization from TensorBoard

We need to create a **session** in order to get the value of a

TensorFlow

- Create a session

```
import tensorflow as tf  
a = tf.add(3, 5)  
sess = tf.Session()  
print sess.run(a)      >> output 8  
sess.close()
```

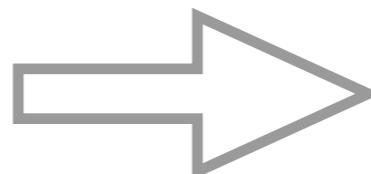


Session will find the dependency of `a`, and computes all the nodes that lead to `a`

TensorFlow

- Create a session (Recommended practice)

```
import tensorflow as tf  
a = tf.add(3, 5)  
sess = tf.Session()  
print sess.run(a)  
sess.close()
```



```
import tensorflow as tf  
a = tf.add(3, 5)  
with tf.Session() as sess:  
    print sess.run(a)
```

Session will find the dependency of a, and computes all the nodes that lead to a

TensorFlow

- Summary

A Session object encapsulates the running environment such that operators are executed and tensors are evaluated

TensorFlow

- More examples

```
import tensorflow as tf
```

```
x = 2
```

```
y = 3
```

```
op1 = tf.add(x, y)
```

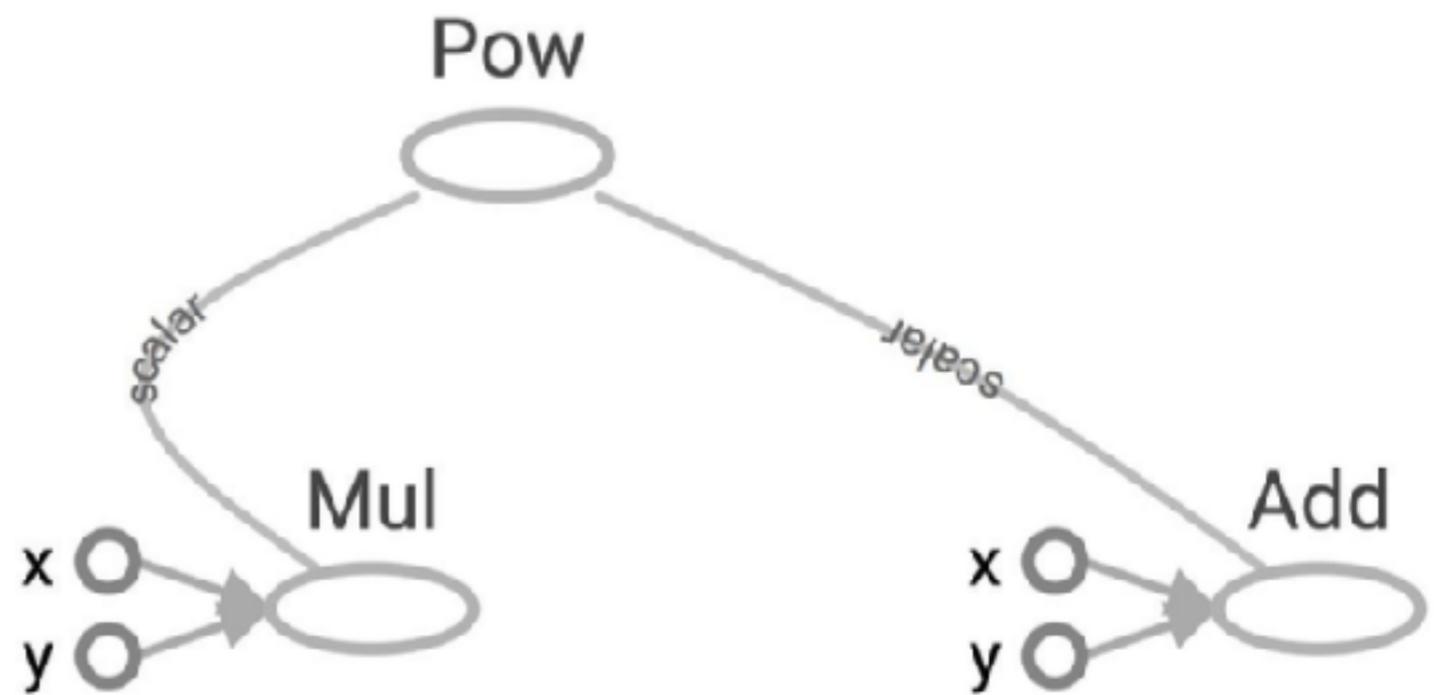
```
op2 = tf.multiply(x, y)
```

```
op3 = tf.pow(op2, op1)
```

```
with tf.Session() as sess:
```

```
    op3 = sess.run(op3)
```

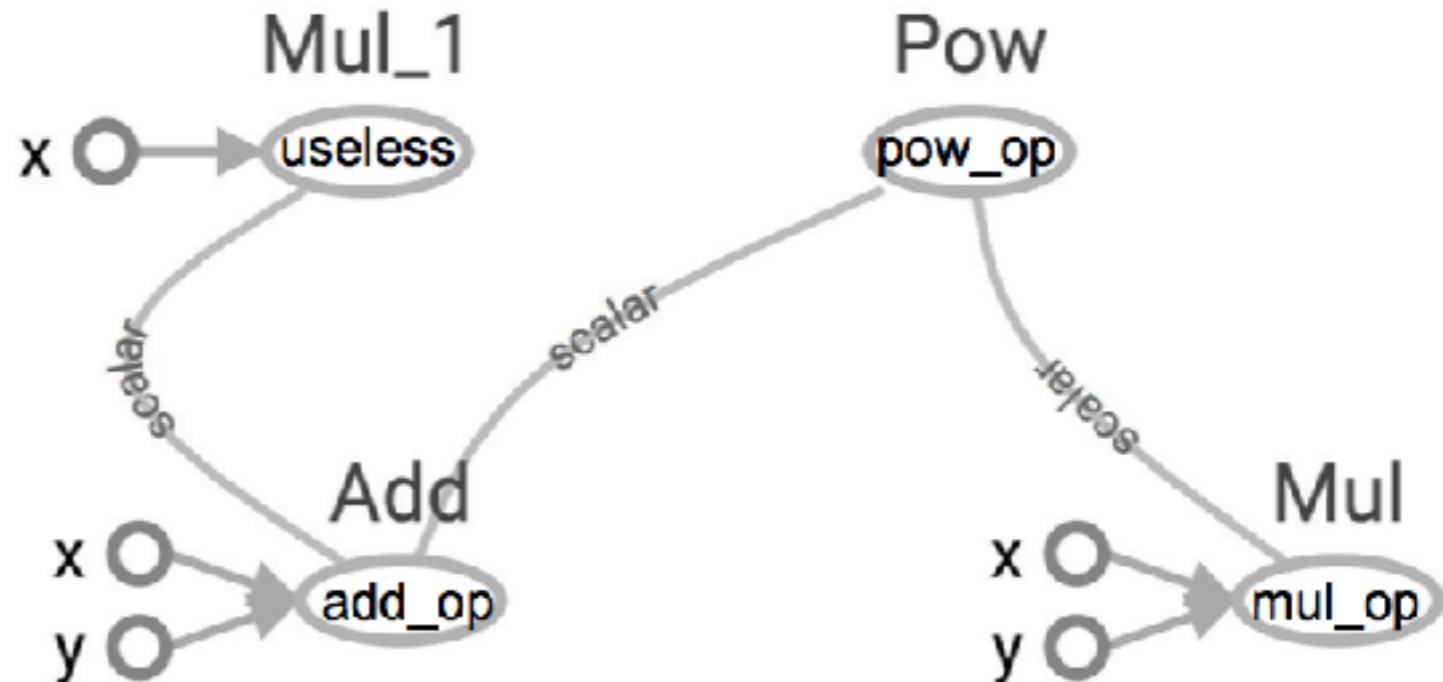
```
print op3  >> output 7776
```



TensorFlow

- More examples

```
import tensorflow as tf  
  
x = 2  
y = 3  
  
add_op = tf.add(x, y)  
  
mul_op = tf.multiply(x, y)  
  
useless = tf.multiply(x, add_op)  
  
pow_op = tf.pow(add_op, mul_op)  
  
with tf.Session() as sess:  
    z = sess.run(pow_op)  
  
print z    >> output 15625
```

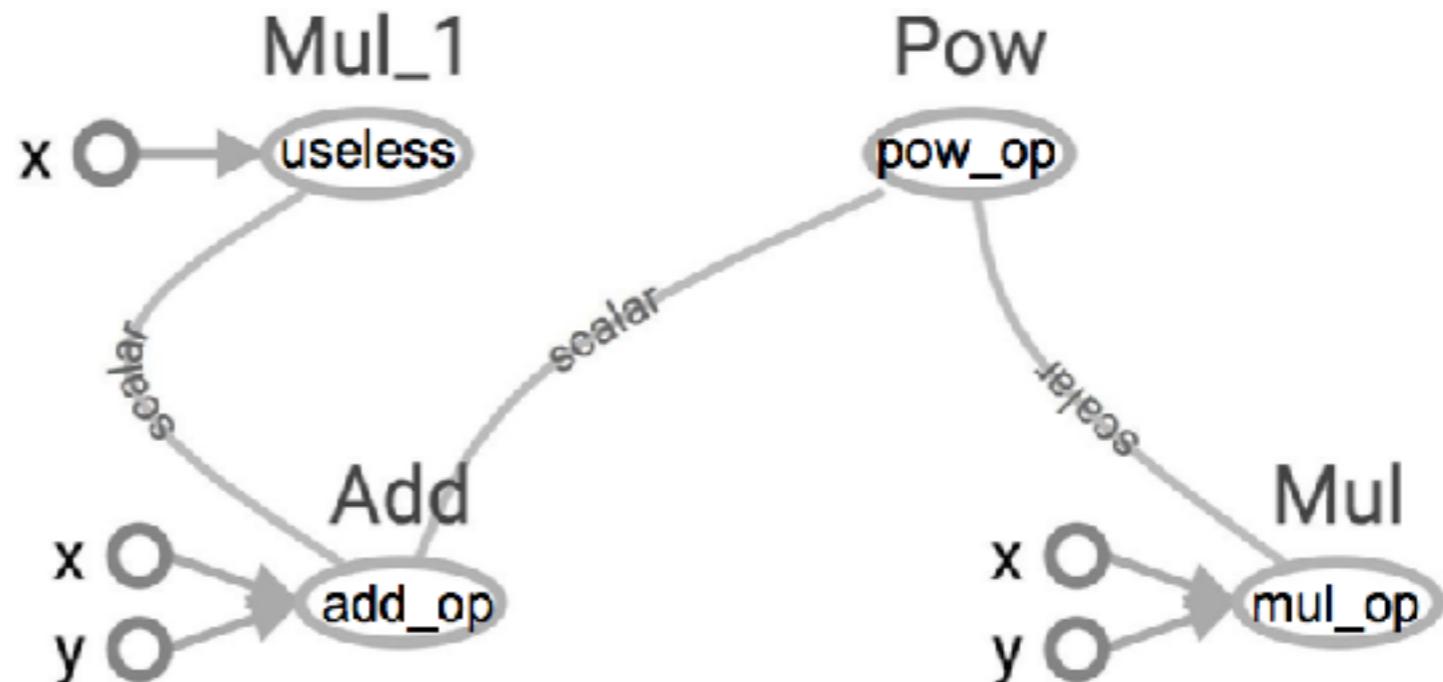


Think: will session also compute the value of **useless**?

TensorFlow

- More examples

```
import tensorflow as tf  
  
x = 2  
y = 3  
  
add_op = tf.add(x, y)  
  
mul_op = tf.multiply(x, y)  
  
useless = tf.multiply(x, add_op)  
  
pow_op = tf.pow(add_op, mul_op)  
  
with tf.Session() as sess:  
    z = sess.run(pow_op)  
  
print z    >> output 15625
```

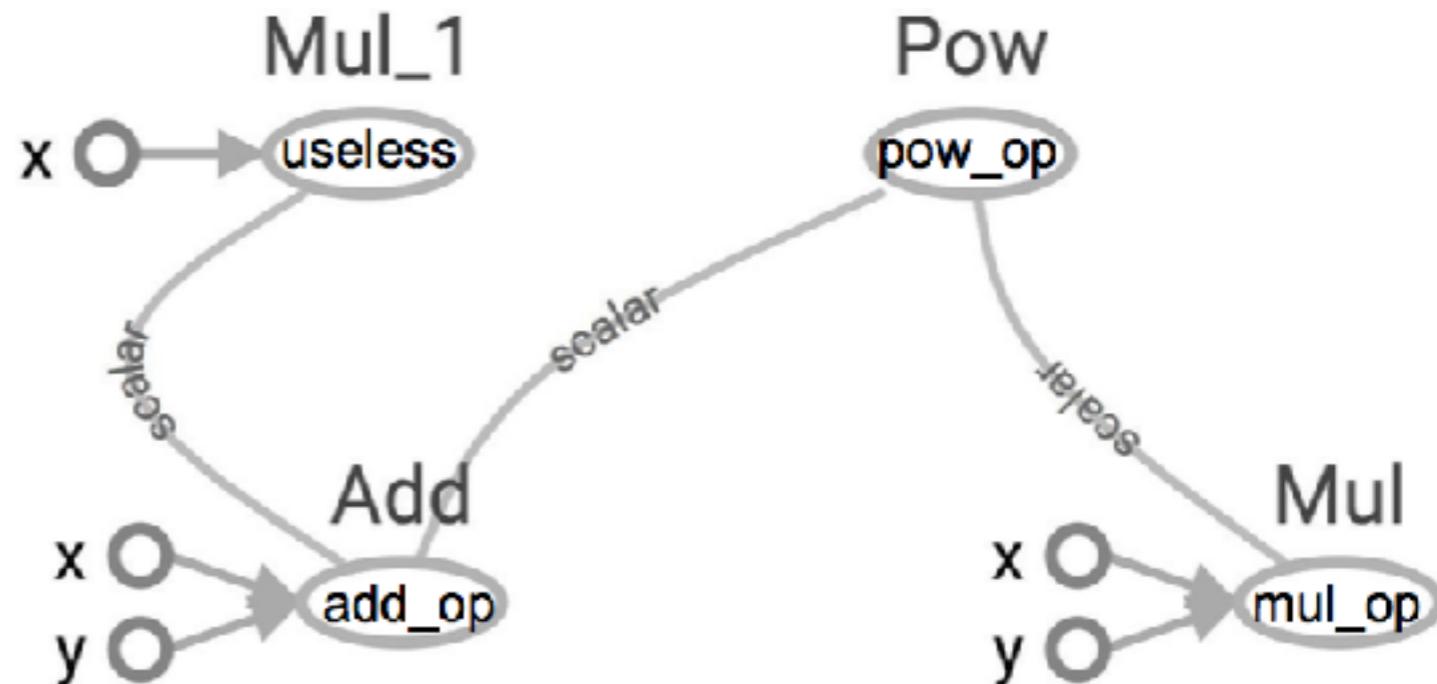


No: print useless >> Tensor("Mul_3:0", shape=(), dtype=int32)

TensorFlow

- More examples

```
import tensorflow as tf  
  
x = 2  
y = 3  
  
add_op = tf.add(x, y)  
  
mul_op = tf.multiply(x, y)  
  
useless = tf.multiply(x, add_op)  
  
pow_op = tf.pow(add_op, mul_op)  
  
with tf.Session() as sess:  
    z, w = sess.run([pow_op, useless])  
  
print z, w    >> output 15625, 10
```



API: `tf.Session.run(fetches, feed_dict=None, options=None, run_metadata=None)`

Pass all the variables you want to evaluate to a list in fetches

TensorFlow

- Run session with specific device:

```
import tensorflow as tf  
  
# build computational graph.  
  
with tf.device("/gpu:0"):  
    a = tf.constant([1.0, 2.0, 3.0, 4.0], shape=(2, 2), name="a")  
    b = tf.constant([2.0, 4.0, 6.0, 8.0], shape=(2, 2), name="b")  
    c = tf.matmul(a, b)  
  
  
# build session, set log_device_placement=True  
with tf.Session(config=tf.ConfigProto(allow_soft_placement=True)):  
    print sess.run(c)
```

TensorFlow

- Why computational graph?
- Save computations (only evaluates subgraphs which lead to the values you're interested in)
- Facilitate distributed computation (model parallelization)
- Directed acyclic graph is required in order to implement auto-differentiation

TensorFlow

- TensorBoard for visualization

```
import tensorflow as tf  
a = tf.constant(2)  
b = tf.constant(3)  
x = tf.add(a, b)  
with tf.Session() as sess:  
    # add this line to use TensorBoard  
    writer = tf.summary.FileWriter("./graphs", sess.graph)  
    print sess.run(x)  
writer.close()
```

Create the summary writer after graph definition but before running the session

TensorFlow

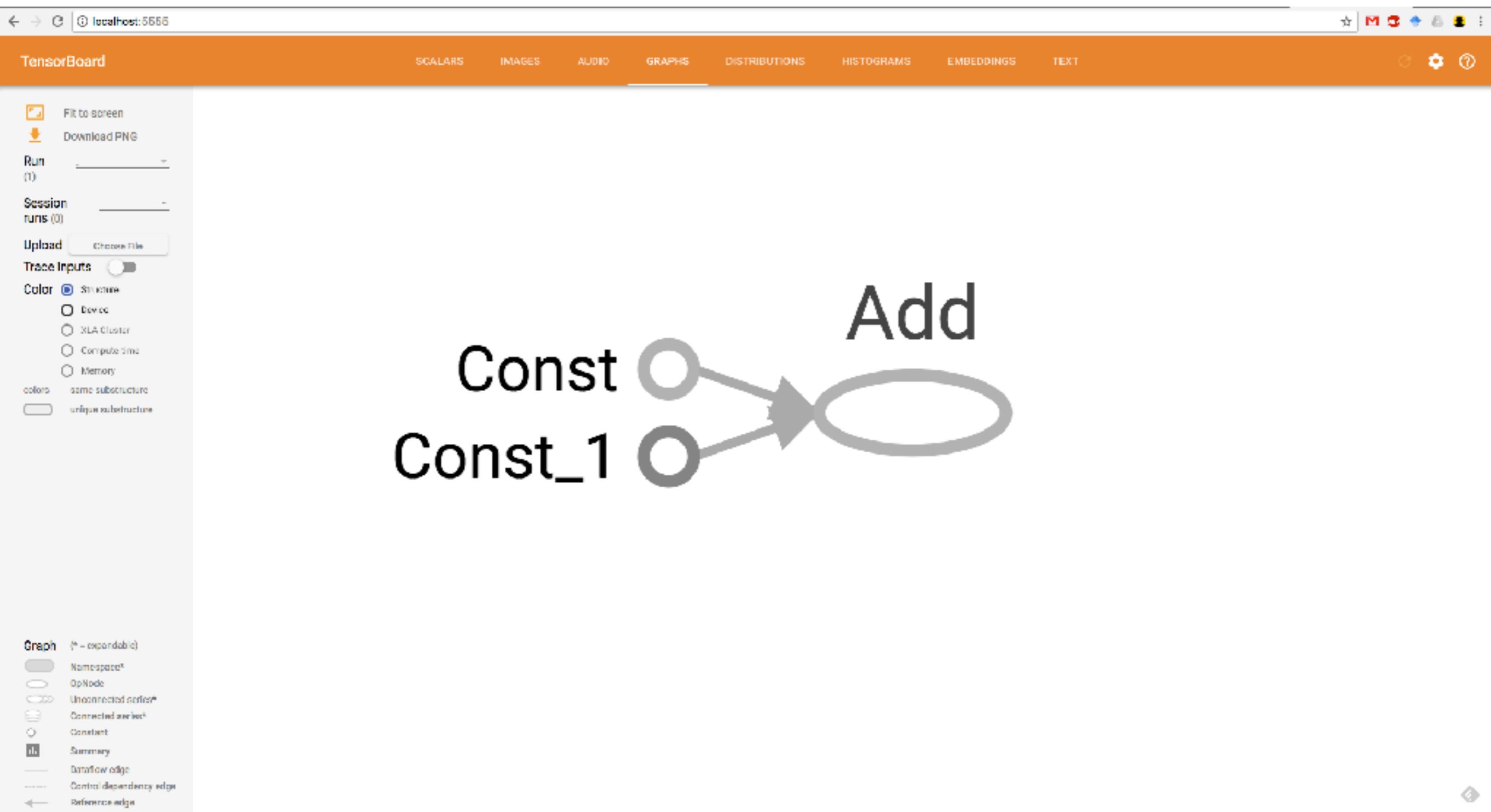
- TensorBoard for visualization

Open terminal, run:

```
$ python [thisprogram].py
```

```
$ tensor board --logdir="./graphs" --port 5555
```

TensorFlow



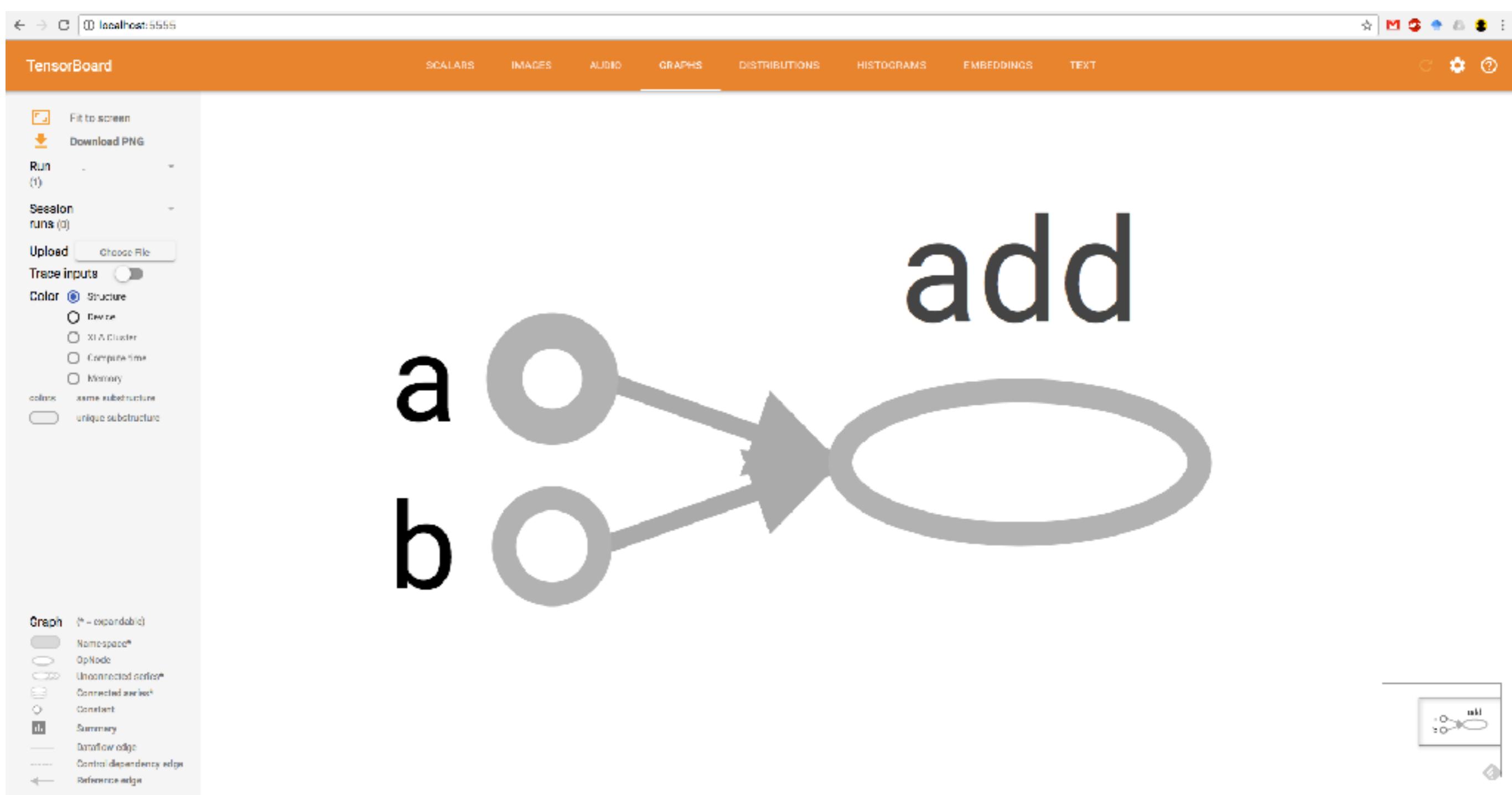
TensorFlow

- TensorBoard for visualization

We can name the variables, as well as the operators:

```
import tensorflow as tf  
  
a = tf.constant(2, name="a")  
  
b = tf.constant(3, name="b")  
  
x = tf.add(a, b, name="add")  
  
with tf.Session() as sess:  
  
    # add this line to use TensorBoard  
    writer = tf.summary.FileWriter("./graphs", sess.graph)  
  
    print sess.run(x)  
  
writer.close()
```

TensorFlow



TensorFlow

- Constants

tf.constant(value, dtype=None, shape=None, name=“Const”, verify_shape=False)

- Very similar to that of Numpy

- tf.zeros
- tf.zeros_like
- tf.ones
- tf.ones_like
- tf.fill
- tf.constant
- tf.linspace
- tf.range

TensorFlow

- Random variables
- `tf.random_normal(shape, mean=0.0, stddev=1.0, dtype=tf.float32, seed=None, name=None)`
- `tf.truncated_normal(shape, mean=0.0, stddev=1.0, dtype=tf.float32, seed=None, name=None)`
- `tf.random_uniform(shape, minval=0, maxval=None, dtype=tf.float32, seed=None, name=None)`
- `tf.random_shuffle(value, seed=None, name=None)`
- `tf.random_crop(value, size, seed=None, name=None)`
- `tf.multinomial(logits, num_samples, seed=None, name=None)`
- `tf.random_gamma(shape, alpha, beta=None, dtype=tf.float32, seed=None, name=None)`
- Set random seed: `tf.set_random_seed(seed)`

https://www.tensorflow.org/api_guides/python/constant_op#Random_Tensors

TensorFlow

- In TensorFlow we can perform all the usual matrix operations in Numpy

Category	Examples
Element-wise mathematical operations	Add, Sub, Mul, Div, Exp, Log, Greater, Less, Equal, ...
Array operations	Concat, Slice, Split, Constant, Rank, Shape, Shuffle, ...
Matrix operations	MatMul, MatrixInverse, MatrixDeterminant, ...
Stateful operations	Variable, Assign, AssignAdd, ...
Neural network building blocks	SoftMax, Sigmoid, ReLU, Convolution2D, MaxPool, ...
Checkpointing operations	Save, Restore
Queue and synchronization operations	Enqueue, Dequeue, MutexAcquire, MutexRelease, ...
Control flow operations	Merge, Switch, Enter, Leave, NextIteration

TensorFlow

- Operations

```
import tensorflow as tf  
  
a = tf.constant([3, 6])  
  
b = tf.constant([2, 2])  
  
tf.add(a, b)    >> [5, 8]  
  
tf.add_n([a, b, b]) >> [7, 10] = a + b + b  
  
tf.multiply(a, b)  >> [6, 12], elementwise multiplication  
  
tf.matmul(a, b)   >> Error, shape inconsistency for matrix multiplication  
  
tf.matmul(tf.reshape(a, [1, 2]), tf.reshape(b, [2, 1]))  >> 18  
  
tf.div(a, b)     >> [1, 3], elementwise division  
  
tf.mod(a, b)    >> [1, 0], elementwise modulus
```

- More math operations at:

https://www.tensorflow.org/api_guides/python/math_ops

TensorFlow

- TensorFlow data types

Data type	Python type	Description
<code>DT_FLOAT</code>	<code>tf.float32</code>	32 bits floating point.
<code>DT_DOUBLE</code>	<code>tf.float64</code>	64 bits floating point.
<code>DT_INT8</code>	<code>tf.int8</code>	8 bits signed integer.
<code>DT_INT16</code>	<code>tf.int16</code>	16 bits signed integer.
<code>DT_INT32</code>	<code>tf.int32</code>	32 bits signed integer.
<code>DT_INT64</code>	<code>tf.int64</code>	64 bits signed integer.
<code>DT_UINT8</code>	<code>tf.uint8</code>	8 bits unsigned integer.
<code>DT_UINT16</code>	<code>tf.uint16</code>	16 bits unsigned integer.
<code>DT_STRING</code>	<code>tf.string</code>	Variable length byte arrays. Each element of a Tensor is a byte array.
<code>DT_BOOL</code>	<code>tf.bool</code>	Boolean.
<code>DT_COMPLEX64</code>	<code>tf.complex64</code>	Complex number made of two 32 bits floating points: real and imaginary parts.
<code>DT_COMPLEX128</code>	<code>tf.complex128</code>	Complex number made of two 64 bits floating points: real and imaginary parts.
<code>DT_QINT8</code>	<code>tf.qint8</code>	8 bits signed integer used in quantized Ops.
<code>DT_QINT32</code>	<code>tf.qint32</code>	32 bits signed integer used in quantized Ops.
<code>DT_QUINT8</code>	<code>tf.quint8</code>	8 bits unsigned integer used in quantized Ops.

https://www.tensorflow.org/programmers_guide/dims_types#data_types

TensorFlow

- Variables (tf.Variable is a class, tf. constant is an op)
- Constants are stored in graph definition, but variables are not!

```
# create variable with a scalar value  
a = tf.Variable(2, name="scalar")  
  
# create variable with a vector value  
b = tf.Variable([2, 3], name="vector")  
  
# create variable with a matrix value  
c = tf.Variable([[1, 2], [3, 4]], name="matrix")  
  
# create variable with zeros  
W = tf.Variable(tf.zeros([784, 10]))
```

https://www.tensorflow.org/programmers_guide/variables

TensorFlow

- Variables contain operations:

```
x = tf.Variable(...)
```

```
x.initializer # init op
```

```
x.value() # read op
```

```
x.assign(...) # write op
```

```
x.assign_add(...) # and more
```

https://www.tensorflow.org/programmers_guide/variables

TensorFlow

- Variables should be initialized before running

```
# The easiest way to initialize all the variables
```

```
init = tf.global_variables_initializer()
```

```
with tf.Session() as sess:
```

```
    sess.run(init)
```

```
# Initialize a subset of variables
```

```
init_ab = tf.variables_initializer([a, b], name="init_ab")
```

```
with tf.Session() as sess:
```

```
    sess.run(init_ab)
```

```
# Initialize a single variable
```

```
W = tf.Variable(tf.zeros([784, 10]))
```

```
with tf.Session() as sess:
```

```
    sess.run(W.initializer)
```

TensorFlow

- Print the values of variables

```
W = tf.Variable(tf.random_normal([784, 10]))
```

```
with tf.Session() as sess:
```

```
    sess.run(W.initializer)
```

```
    print W
```

```
    print W.eval()
```

- >> <tf.Variable ‘Variable_1:0’, shape=(784, 10), dtype=float32>
- >> [[-0.4778471, -1.3822577, ...]]

TensorFlow

- Placeholders
- Symbolic variables that do not take actual value when defined
- Can be filled with actual values when executed
- Examples:

$$f(x, y) = 2x + y$$

- The x , y in the above function are placeholders – they don't have specific values when defined
- However we can evaluate $f(x, y)$ by giving them specific values

TensorFlow

- Placeholders
- **tf.placeholder(dtype, shape=None, name=None)**

```
# create a placeholder of float32, shape is 1x3
a = tf.placeholder(tf.float32, shape=[3])

# create a constant of type float32, shape is 1x3
b = tf.constant([5, 5, 5], tf.float32)

c = a + b
```

with tf.Session() as sess:

```
print sess.run(c)      >> Error?
```

Before running c, we need to provide a with specific values!

TensorFlow

- Placeholders
- **tf.placeholder(dtype, shape=None, name=None)**

```
# create a placeholder of float32, shape is 1x3
a = tf.placeholder(tf.float32, shape=[3])
# create a constant of type float32, shape is 1x3
b = tf.constant([5, 5, 5], tf.float32)
c = a + b
```

with tf.Session() as sess:

```
# feed [1, 2, 3] to placeholder a via a dictionary
print sess.run(c, {a: [1, 2, 3]})  >> the tensor a can be used as a key, [6, 7, 8]
```

Note: shape=None means the placeholder can have any shape

TensorFlow

- Logistic regression

The screenshot shows a Jupyter Notebook interface with the title "Logistic regression with TensorFlow". The notebook has a header with "File Edit View Insert Cell Kernel Help" and a toolbar with various icons. On the right, there are "Logout" and "Python 2" buttons. The notebook content is as follows:

Logistic Regression with TensorFlow
By Ivan

0. Imports

```
In [37]: import numpy as np
import tensorflow as tf
import time
import matplotlib.pyplot as plt
from tensorflow.examples.tutorials.mnist import input_data
```

1. Download MNIST dataset

```
In [38]: # Using tensorflow's default tools to fetch data, this is the same as what we did in the first homework assignment.
mnist = input_data.read_data_sets('./mnist', one_hot=True)

Extracting ./mnist/train-images-idx3-ubyte.gz
Extracting ./mnist/train-labels-idx1-ubyte.gz
Extracting ./mnist/t10k-images-idx3-ubyte.gz
Extracting ./mnist/t10k-labels-idx1-ubyte.gz
```

```
In [39]: # Check the dimension of training, validation and test sets.
mnist.train.images.shape
```

```
Out[39]: (55000, 784)
```

```
In [40]: mnist.train.labels.shape
```

```
Out[40]: (55000, 10)
```

```
In [41]: mnist.validation.images.shape
```

```
Out[41]: (5000, 784)
```

```
In [42]: mnist.validation.labels.shape
```

```
Out[42]: (5000, 10)
```

```
In [43]: mnist.test.images.shape
```

```
Out[43]: (10000, 784)
```

Homework

- Image classification on MNIST with a three layer MLP:
 - Image size 28x28, size of MLP: 784-500-10, with ReLU as the nonlinear activation function
 - batch size = 200
 - learning rate = 0.1
 - # iterations = 20
 - We provide an initial python script for you to work on
- What you need to implement:
 - Use TensorFlow to build the computational graph
 - Build necessary operator for SGD optimization
 - Run the graph to train the model
 - If implemented correctly, you should see a test set classification accuracy ~ 0.975
 - Running on GTX-1080, taking around ~ 80 seconds (graph compilation take time).
- You need to install numpy and tensorflow

PyTorch Introduction

Necessary points of a good deep learning framework:

- Easy and flexible to build large computational graph
- Easy to compute gradients in computational graph, automatically
- Easy and efficient to run on GPU (using cuDNN, cuBLAS, etc.)

PyTorch Introduction

Example using Numpy:

Computational Graphs

Numpy

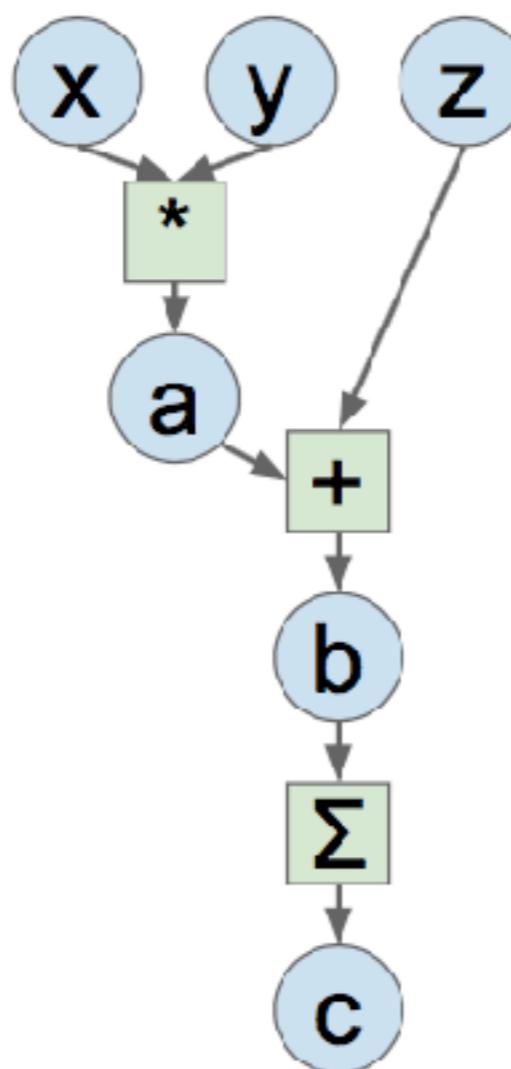
```
import numpy as np
np.random.seed(0)

N, D = 3, 4

x = np.random.randn(N, D)
y = np.random.randn(N, D)
z = np.random.randn(N, D)

a = x * y
b = a + z
c = np.sum(b)

grad_c = 1.0
grad_b = grad_c * np.ones((N, D))
grad_a = grad_b.copy()
grad_z = grad_b.copy()
grad_x = grad_a * y
grad_y = grad_a * x
```



PyTorch Introduction

Example using Numpy:

Computational Graphs

Numpy

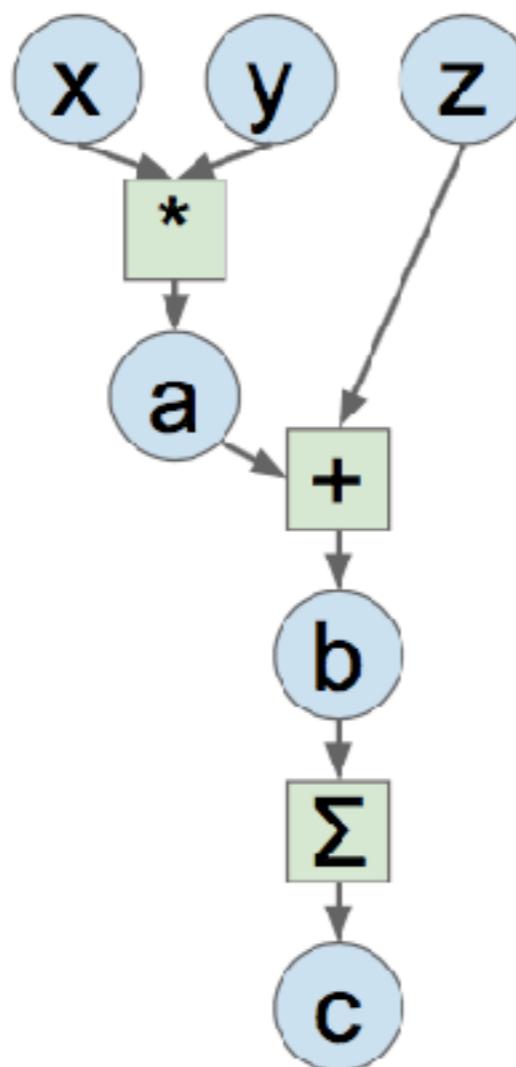
```
import numpy as np
np.random.seed(0)

N, D = 3, 4

x = np.random.randn(N, D)
y = np.random.randn(N, D)
z = np.random.randn(N, D)

a = x * y
b = a + z
c = np.sum(b)

grad_c = 1.0
grad_b = grad_c * np.ones((N, D))
grad_a = grad_b.copy()
grad_z = grad_b.copy()
grad_x = grad_a * y
grad_y = grad_a * x
```



Problems:

- Can't run on GPU
- Have to compute our own gradients

PyTorch Introduction

Computational Graphs

Numpy

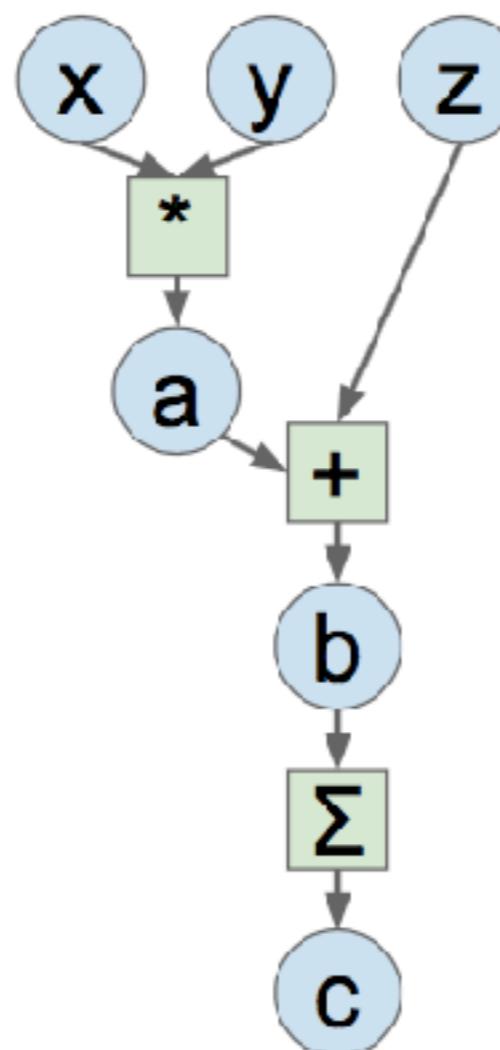
```
import numpy as np
np.random.seed(0)

N, D = 3, 4

x = np.random.randn(N, D)
y = np.random.randn(N, D)
z = np.random.randn(N, D)

a = x * y
b = a + z
c = np.sum(b)

grad_c = 1.0
grad_b = grad_c * np.ones((N, D))
grad_a = grad_b.copy()
grad_z = grad_b.copy()
grad_x = grad_a * y
grad_y = grad_a * x
```



TensorFlow

```
# Basic computational graph
import numpy as np
np.random.seed(0)
import tensorflow as tf

N, D = 3, 4

x = tf.placeholder(tf.float32)
y = tf.placeholder(tf.float32)
z = tf.placeholder(tf.float32)

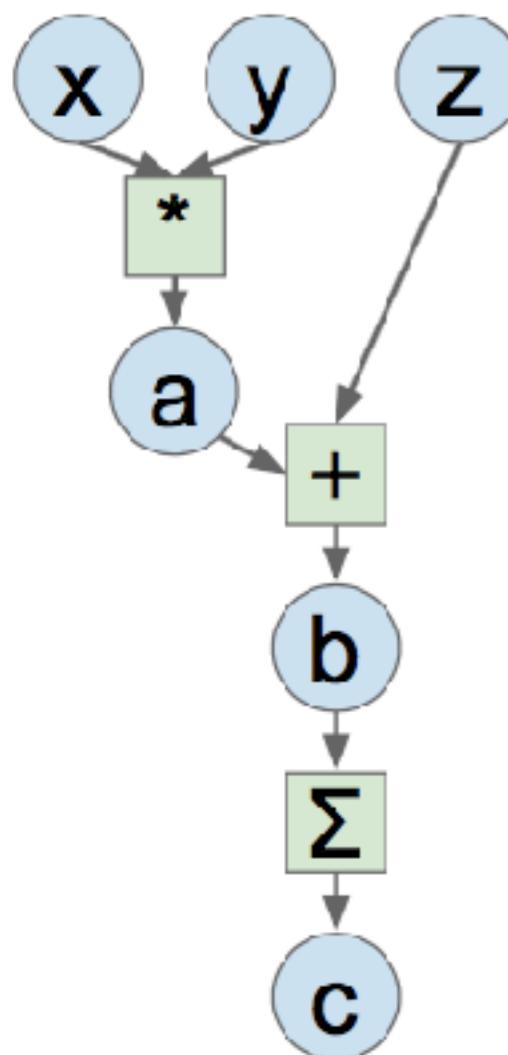
a = x * y
b = a + z
c = tf.reduce_sum(b)

grad_x, grad_y, grad_z = tf.gradients(c, [x, y, z])

with tf.Session() as sess:
    values = {
        x: np.random.randn(N, D),
        y: np.random.randn(N, D),
        z: np.random.randn(N, D),
    }
    out = sess.run([c, grad_x, grad_y, grad_z],
                  feed_dict=values)
    c_val, grad_x_val, grad_y_val, grad_z_val = out
```

PyTorch Introduction

Computational Graphs



Create forward
computational graph

TensorFlow

```
# Basic computational graph
import numpy as np
np.random.seed(0)
import tensorflow as tf

N, D = 3, 4

x = tf.placeholder(tf.float32)
y = tf.placeholder(tf.float32)
z = tf.placeholder(tf.float32)

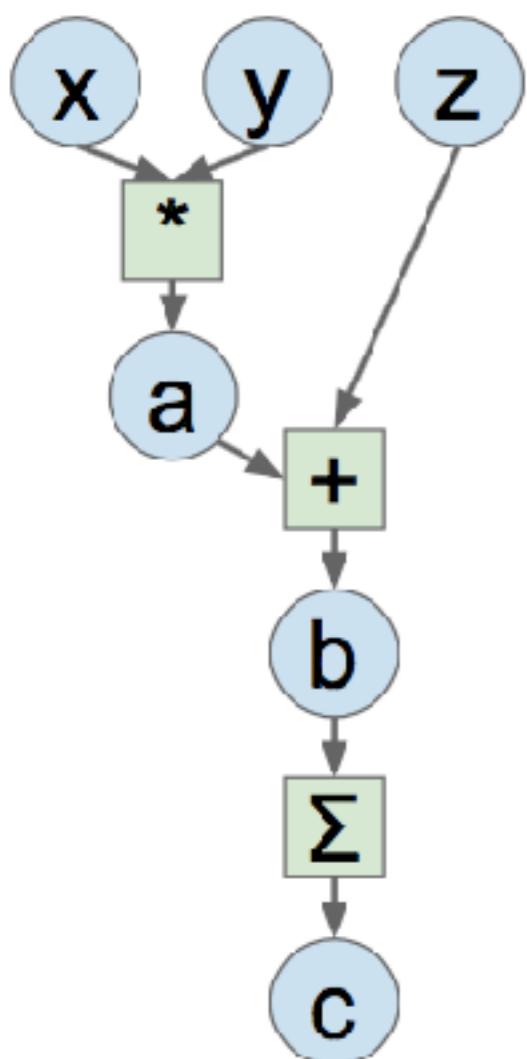
a = x * y
b = a + z
c = tf.reduce_sum(b)

grad_x, grad_y, grad_z = tf.gradients(c, [x, y, z])

with tf.Session() as sess:
    values = {
        x: np.random.randn(N, D),
        y: np.random.randn(N, D),
        z: np.random.randn(N, D),
    }
    out = sess.run([c, grad_x, grad_y, grad_z],
                  feed_dict=values)
    c_val, grad_x_val, grad_y_val, grad_z_val = out
```

PyTorch Introduction

Computational Graphs



Ask TensorFlow to compute gradients

TensorFlow

```
# Basic computational graph
import numpy as np
np.random.seed(0)
import tensorflow as tf

N, D = 3, 4

x = tf.placeholder(tf.float32)
y = tf.placeholder(tf.float32)
z = tf.placeholder(tf.float32)

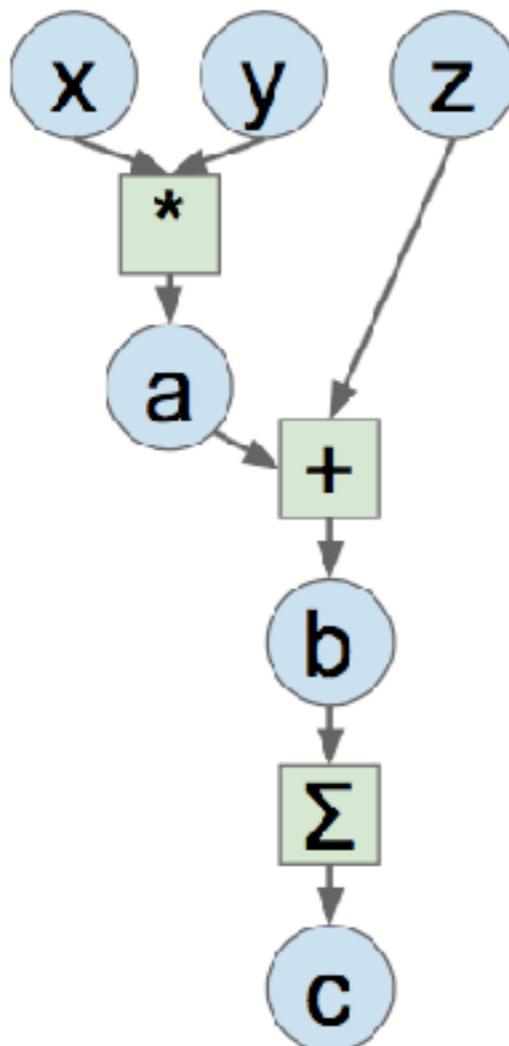
a = x * y
b = a + z
c = tf.reduce_sum(b)

grad_x, grad_y, grad_z = tf.gradients(c, [x, y, z])

with tf.Session() as sess:
    values = {
        x: np.random.randn(N, D),
        y: np.random.randn(N, D),
        z: np.random.randn(N, D),
    }
    out = sess.run([c, grad_x, grad_y, grad_z],
                  feed_dict=values)
    c_val, grad_x_val, grad_y_val, grad_z_val = out
```

PyTorch Introduction

Computational Graphs



Tell
TensorFlow
to run on **CPU**

TensorFlow

```
import numpy as np
np.random.seed(0)
import tensorflow as tf

N, D = 3000, 4000

with tf.device('/cpu:0'):
    x = tf.placeholder(tf.float32)
    y = tf.placeholder(tf.float32)
    z = tf.placeholder(tf.float32)

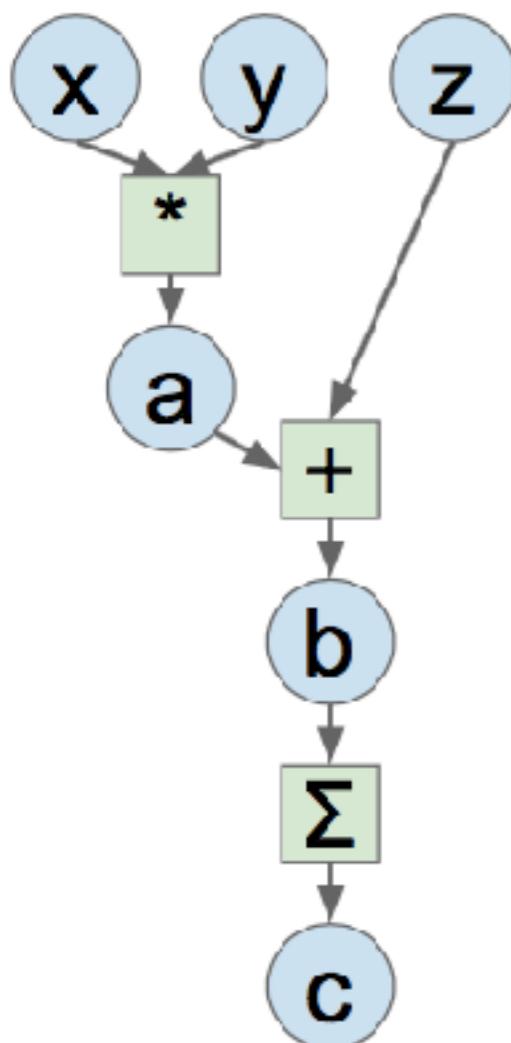
    a = x * y
    b = a + z
    c = tf.reduce_sum(b)

grad_x, grad_y, grad_z = tf.gradients(c, [x, y, z])

with tf.Session() as sess:
    values = {
        x: np.random.randn(N, D),
        y: np.random.randn(N, D),
        z: np.random.randn(N, D),
    }
    out = sess.run([c, grad_x, grad_y, grad_z],
                  feed_dict=values)
    c_val, grad_x_val, grad_y_val, grad_z_val = out
```

PyTorch Introduction

Computational Graphs



Tell
TensorFlow
to run on GPU

TensorFlow

```
import numpy as np
np.random.seed(0)
import tensorflow as tf

N, D = 3000, 4000

with tf.device('/gpu:0'):
    x = tf.placeholder(tf.float32)
    y = tf.placeholder(tf.float32)
    z = tf.placeholder(tf.float32)

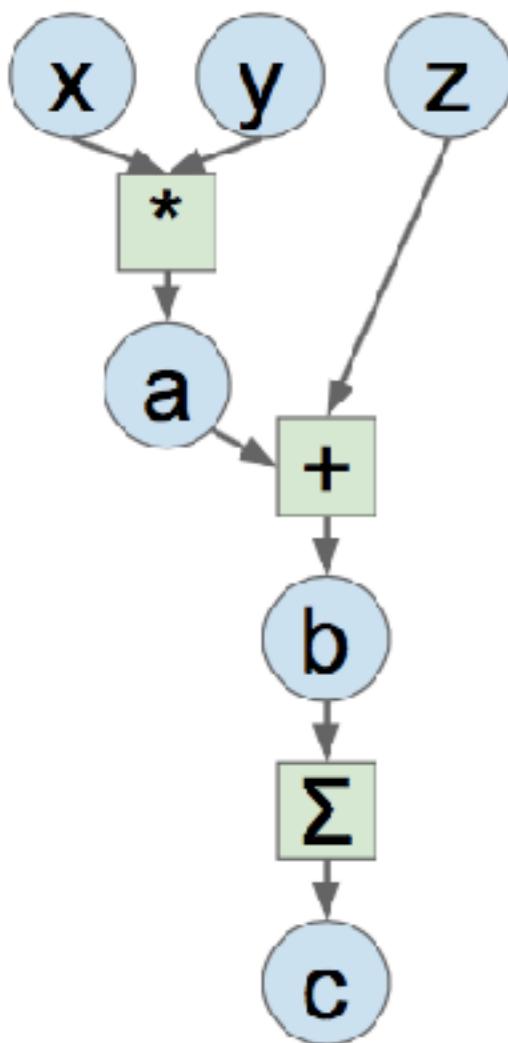
    a = x * y
    b = a + z
    c = tf.reduce_sum(b)

grad_x, grad_y, grad_z = tf.gradients(c, [x, y, z])

with tf.Session() as sess:
    values = {
        x: np.random.randn(N, D),
        y: np.random.randn(N, D),
        z: np.random.randn(N, D),
    }
    out = sess.run([c, grad_x, grad_y, grad_z],
                  feed_dict=values)
    c_val, grad_x_val, grad_y_val, grad_z_val = out
```

PyTorch Introduction

Computational Graphs



PyTorch

```
import torch
from torch.autograd import Variable

N, D = 3, 4

x = Variable(torch.randn(N, D),
             requires_grad=True)
y = Variable(torch.randn(N, D),
             requires_grad=True)
z = Variable(torch.randn(N, D),
             requires_grad=True)

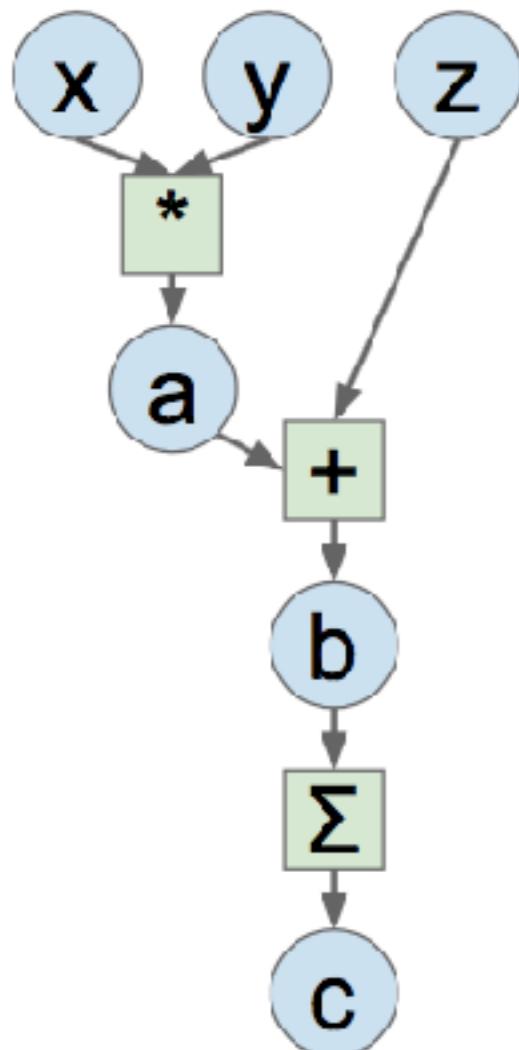
a = x * y
b = a + z
c = torch.sum(b)

c.backward()

print(x.grad.data)
print(y.grad.data)
print(z.grad.data)
```

PyTorch Introduction

Computational Graphs



Define **Variables** to start building a computational graph

PyTorch

```
import torch
from torch.autograd import Variable

N, D = 3, 4

x = Variable(torch.randn(N, D),
             requires_grad=True)
y = Variable(torch.randn(N, D),
             requires_grad=True)
z = Variable(torch.randn(N, D),
             requires_grad=True)

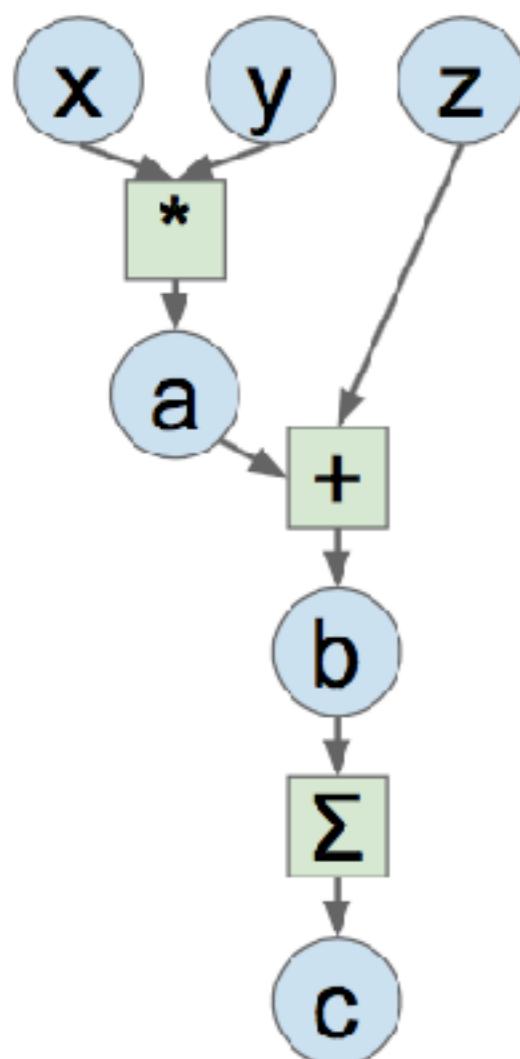
a = x * y
b = a + z
c = torch.sum(b)

c.backward()

print(x.grad.data)
print(y.grad.data)
print(z.grad.data)
```

PyTorch Introduction

Computational Graphs



Forward pass
looks just like
numpy

PyTorch

```
import torch
from torch.autograd import Variable

N, D = 3, 4

x = Variable(torch.randn(N, D),
             requires_grad=True)
y = Variable(torch.randn(N, D),
             requires_grad=True)
z = Variable(torch.randn(N, D),
             requires_grad=True)

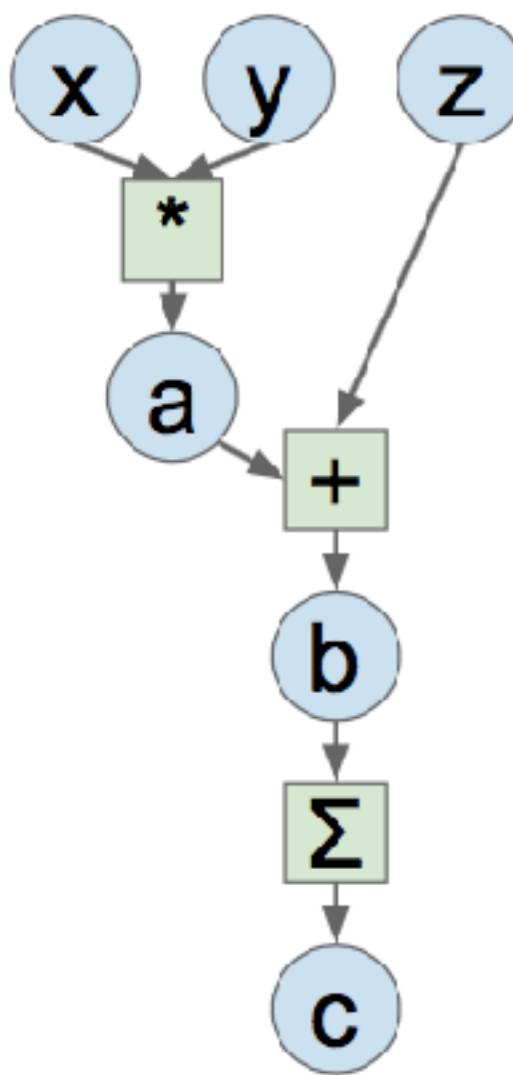
a = x * y
b = a + z
c = torch.sum(b)

c.backward()

print(x.grad.data)
print(y.grad.data)
print(z.grad.data)
```

PyTorch Introduction

Computational Graphs



Calling `c.backward()`
computes all
gradients

PyTorch

```
import torch
from torch.autograd import Variable

N, D = 3, 4

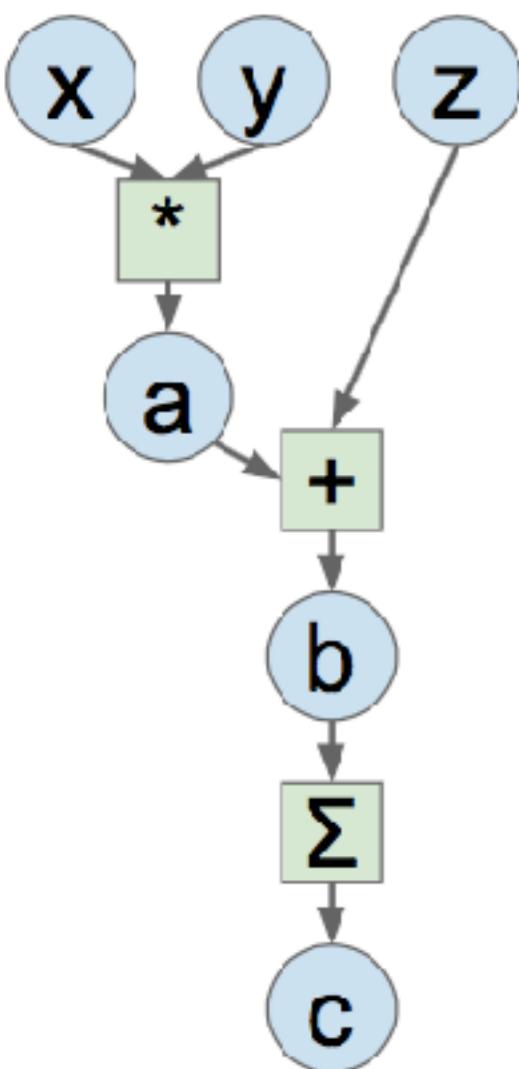
x = Variable(torch.randn(N, D),
             requires_grad=True)
y = Variable(torch.randn(N, D),
             requires_grad=True)
z = Variable(torch.randn(N, D),
             requires_grad=True)

a = x * y
b = a + z
c = torch.sum(b)

c.backward()
print(x.grad.data)
print(y.grad.data)
print(z.grad.data)
```

PyTorch Introduction

Computational Graphs



Run on GPU by
casting to .cuda()

PyTorch

```
import torch
from torch.autograd import Variable

N, D = 3, 4

x = Variable(torch.randn(N, D).cuda(),
             requires_grad=True)
y = Variable(torch.randn(N, D).cuda(),
             requires_grad=True) # Line 5
z = Variable(torch.randn(N, D).cuda(),
             requires_grad=True)

a = x * y
b = a + z
c = torch.sum(b)

c.backward()

print(x.grad.data)
print(y.grad.data)
print(z.grad.data)
```

PyTorch Introduction

Numpy

```
import numpy as np
np.random.seed(0)

N, D = 3, 4

x = np.random.randn(N, D)
y = np.random.randn(N, D)
z = np.random.randn(N, D)

a = x * y
b = a + z
c = np.sum(b)

grad_c = 1.0
grad_b = grad_c * np.ones((N, D))
grad_a = grad_b.copy()
grad_z = grad_b.copy()
grad_x = grad_a * y
grad_y = grad_a * x
```

TensorFlow

```
import numpy as np
np.random.seed(0)
import tensorflow as tf

N, D = 3, 4

with tf.device('/gpu:0'):
    x = tf.placeholder(tf.float32)
    y = tf.placeholder(tf.float32)
    z = tf.placeholder(tf.float32)

    a = x * y
    b = a + z
    c = tf.reduce_sum(b)

grad_x, grad_y, grad_z = tf.gradients(c, [x, y, z])

with tf.Session() as sess:
    values = {
        x: np.random.randn(N, D),
        y: np.random.randn(N, D),
        z: np.random.randn(N, D),
    }
    out = sess.run([c, grad_x, grad_y, grad_z],
                  feed_dict=values)
    c_val, grad_x_val, grad_y_val, grad_z_val = out
```

PyTorch

```
import torch
from torch.autograd import Variable

N, D = 3, 4

x = Variable(torch.randn(N, D).cuda(),
             requires_grad=True)
y = Variable(torch.randn(N, D).cuda(),
             requires_grad=True)
z = Variable(torch.randn(N, D).cuda(),
             requires_grad=True)

a = x * y
b = a + z
c = torch.sum(b)

c.backward()

print(x.grad.data)
print(y.grad.data)
print(z.grad.data)
```

PyTorch Introduction

Install Anaconda + PyTorch:

The screenshot shows the Anaconda download page with three main download links at the top: "Download for Windows", "Download for macOS", and "Download for Linux". The "Download for macOS" link is highlighted with a blue background.

Anaconda 4.4.0

For macOS

macOS 10.12.2 users: To prevent permissions problems, we recommend that you upgrade to macOS 10.12.3 or later before installing Anaconda.

Anaconda is BSD licensed which gives you permission to use Anaconda commercially and for redistribution.

[Changelog](#)

Graphical Installer

1. Download the graphical installer
2. Double-click the downloaded .pkg file and follow the instructions

Command Line Installer

Python 3.6 version

GRAPHICAL INSTALLER (442M)

COMMAND-LINE INSTALLER (380M)

64-Bit

Python 2.7 version

GRAPHICAL INSTALLER (438M)

COMMAND-LINE INSTALLER (375M)

Anaconda includes all the necessary packages, including numpy, scipy, sklearn, etc.

PyTorch Introduction

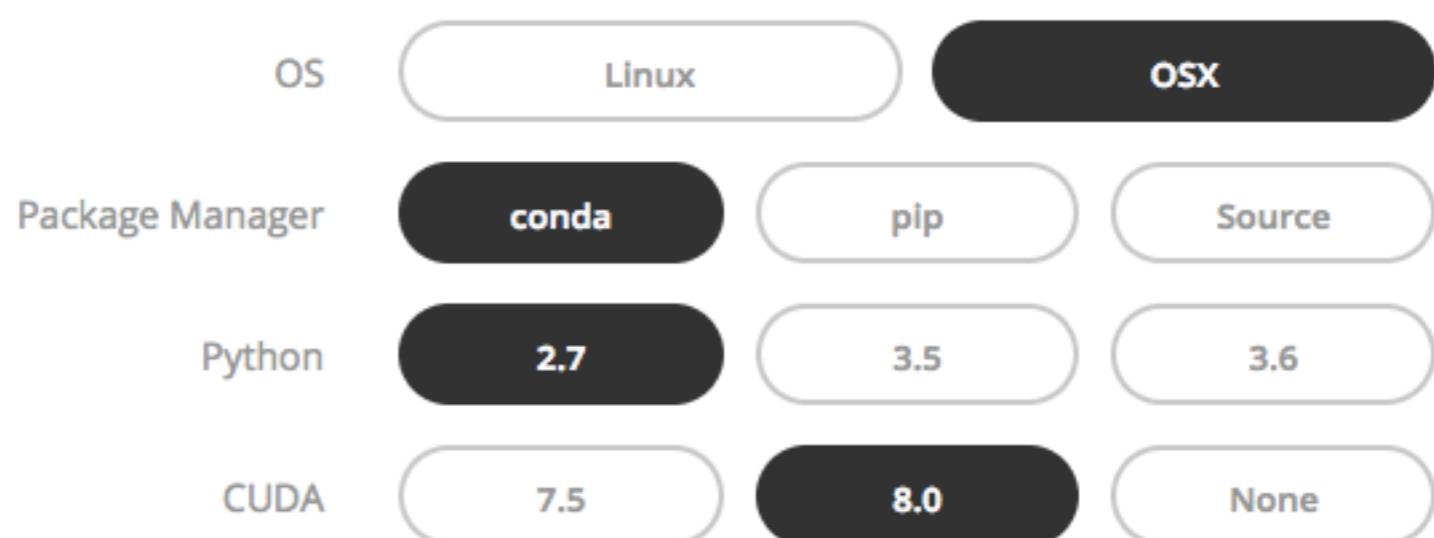
Install Anaconda + PyTorch:

Get Started.

Select your preferences, then run the PyTorch install command.

Please ensure that you are on the latest pip and numpy packages.

Anaconda is our recommended package manager



Run this command:

```
conda install pytorch torchvision -c soumith  
# OSX Binaries dont support CUDA, install from source if CUDA is needed
```

Also needs to install cuda-8 + cuDNN from Nvidia website

PyTorch Introduction

PyTorch: Three levels of abstraction

- Tensor: Similar to ndarray in numpy, but runs on GPU
- Variable: Node in a computational graph, can remember operations that lead to current state; stores data and gradient
- Module: Layers in neural networks, including dense layer, convolution layer, recurrent layer, etc.; layers include learnable weights

PyTorch Introduction

- Like bumpy arrays, but can run on GPU
- Here we build a two-layer DNN

```
import torch

dtype = torch.FloatTensor

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in).type(dtype)
y = torch.randn(N, D_out).type(dtype)
w1 = torch.randn(D_in, H).type(dtype)
w2 = torch.randn(H, D_out).type(dtype)

learning_rate = 1e-6
for t in range(500):
    h = x.mm(w1)
    h_relu = h.clamp(min=0)
    y_pred = h_relu.mm(w2)
    loss = (y_pred - y).pow(2).sum()

    grad_y_pred = 2.0 * (y_pred - y)
    grad_w2 = h_relu.t().mm(grad_y_pred)
    grad_h_relu = grad_y_pred.mm(w2.t())
    grad_h = grad_h_relu.clone()
    grad_h[h < 0] = 0
    grad_w1 = x.t().mm(grad_h)

    w1 -= learning_rate * grad_w1
    w2 -= learning_rate * grad_w2
```

PyTorch Introduction

PyTorch: Tensors

Create random tensors
for data and weights



```
import torch

dtype = torch.FloatTensor

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in).type(dtype)
y = torch.randn(N, D_out).type(dtype)
w1 = torch.randn(D_in, H).type(dtype)
w2 = torch.randn(H, D_out).type(dtype)

learning_rate = 1e-6
for t in range(500):
    h = x.mm(w1)
    h_relu = h.clamp(min=0)
    y_pred = h_relu.mm(w2)
    loss = (y_pred - y).pow(2).sum()

    grad_y_pred = 2.0 * (y_pred - y)
    grad_w2 = h_relu.t().mm(grad_y_pred)
    grad_h_relu = grad_y_pred.mm(w2.t())
    grad_h = grad_h_relu.clone()
    grad_h[h < 0] = 0
    grad_w1 = x.t().mm(grad_h)

    w1 -= learning_rate * grad_w1
    w2 -= learning_rate * grad_w2
```

PyTorch Introduction

PyTorch: Tensors

Forward pass: compute predictions and loss



```
import torch

dtype = torch.FloatTensor

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in).type(dtype)
y = torch.randn(N, D_out).type(dtype)
w1 = torch.randn(D_in, H).type(dtype)
w2 = torch.randn(H, D_out).type(dtype)

learning_rate = 1e-6
for t in range(500):
    h = x.mm(w1)
    h_relu = h.clamp(min=0)
    y_pred = h_relu.mm(w2)
    loss = (y_pred - y).pow(2).sum()

    grad_y_pred = 2.0 * (y_pred - y)
    grad_w2 = h_relu.t().mm(grad_y_pred)
    grad_h_relu = grad_y_pred.mm(w2.t())
    grad_h = grad_h_relu.clone()
    grad_h[h < 0] = 0
    grad_w1 = x.t().mm(grad_h)

    w1 -= learning_rate * grad_w1
    w2 -= learning_rate * grad_w2
```

PyTorch Introduction

PyTorch: Tensors

Backward pass:
manually compute
gradients

```
import torch

dtype = torch.FloatTensor

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in).type(dtype)
y = torch.randn(N, D_out).type(dtype)
w1 = torch.randn(D_in, H).type(dtype)
w2 = torch.randn(H, D_out).type(dtype)

learning_rate = 1e-6
for t in range(500):
    h = x.mm(w1)
    h_relu = h.clamp(min=0)
    y_pred = h_relu.mm(w2)
    loss = (y_pred - y).pow(2).sum()

    grad_y_pred = 2.0 * (y_pred - y)
    grad_w2 = h_relu.t().mm(grad_y_pred)
    grad_h_relu = grad_y_pred.mm(w2.t())
    grad_h = grad_h_relu.clone()
    grad_h[h < 0] = 0
    grad_w1 = x.t().mm(grad_h)

    w1 -= learning_rate * grad_w1
    w2 -= learning_rate * grad_w2
```

PyTorch Introduction

PyTorch: Tensors

```
import torch

dtype = torch.FloatTensor

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in).type(dtype)
y = torch.randn(N, D_out).type(dtype)
w1 = torch.randn(D_in, H).type(dtype)
w2 = torch.randn(H, D_out).type(dtype)

learning_rate = 1e-6
for t in range(500):
    h = x.mm(w1)
    h_relu = h.clamp(min=0)
    y_pred = h_relu.mm(w2)
    loss = (y_pred - y).pow(2).sum()

    grad_y_pred = 2.0 * (y_pred - y)
    grad_w2 = h_relu.t().mm(grad_y_pred)
    grad_h_relu = grad_y_pred.mm(w2.t())
    grad_h = grad_h_relu.clone()
    grad_h[h < 0] = 0
    grad_w1 = x.t().mm(grad_h)

    w1 -= learning_rate * grad_w1
    w2 -= learning_rate * grad_w2
```

Gradient descent
step on weights

PyTorch Introduction

PyTorch: Tensors

To run on GPU, just cast tensors to a cuda datatype!



```
import torch

dtype = torch.cuda.FloatTensor

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in).type(dtype)
y = torch.randn(N, D_out).type(dtype)
w1 = torch.randn(D_in, H).type(dtype)
w2 = torch.randn(H, D_out).type(dtype)

learning_rate = 1e-6
for t in range(500):
    h = x.mm(w1)
    h_relu = h.clamp(min=0)
    y_pred = h_relu.mm(w2)
    loss = (y_pred - y).pow(2).sum()

    grad_y_pred = 2.0 * (y_pred - y)
    grad_w2 = h_relu.t().mm(grad_y_pred)
    grad_h_relu = grad_y_pred.mm(w2.t())
    grad_h = grad_h_relu.clone()
    grad_h[h < 0] = 0
    grad_w1 = x.t().mm(grad_h)

    w1 -= learning_rate * grad_w1
    w2 -= learning_rate * grad_w2
```

PyTorch Introduction

PyTorch: Autograd

A PyTorch **Variable** is a node in a computational graph

`x.data` is a Tensor

`x.grad` is a Variable of gradients
(same shape as `x.data`)

`x.grad.data` is a Tensor of gradients

```
import torch
from torch.autograd import Variable

N, D_in, H, D_out = 64, 1000, 100, 10
x = Variable(torch.randn(N, D_in), requires_grad=False)
y = Variable(torch.randn(N, D_out), requires_grad=False)
w1 = Variable(torch.randn(D_in, H), requires_grad=True)
w2 = Variable(torch.randn(H, D_out), requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    if w1.grad: w1.grad.data.zero_()
    if w2.grad: w2.grad.data.zero_()
    loss.backward()

    w1.data -= learning_rate * w1.grad.data
    w2.data -= learning_rate * w2.grad.data
```

PyTorch Introduction

PyTorch: Autograd

PyTorch Tensors and Variables
have the same API!

Variables remember how they were
created (for backprop)

```
import torch
from torch.autograd import Variable

N, D_in, H, D_out = 64, 1000, 100, 10
x = Variable(torch.randn(N, D_in), requires_grad=False)
y = Variable(torch.randn(N, D_out), requires_grad=False)
w1 = Variable(torch.randn(D_in, H), requires_grad=True)
w2 = Variable(torch.randn(H, D_out), requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    if w1.grad: w1.grad.data.zero_()
    if w2.grad: w2.grad.data.zero_()
    loss.backward()

    w1.data -= learning_rate * w1.grad.data
    w2.data -= learning_rate * w2.grad.data
```

PyTorch Introduction

PyTorch: Autograd

We will not want gradients
(of loss) with respect to data

Do want gradients with
respect to weights

```
import torch
from torch.autograd import Variable

N, D_in, H, D_out = 64, 1000, 100, 10
x = Variable(torch.randn(N, D_in), requires_grad=False)
y = Variable(torch.randn(N, D_out), requires_grad=False)
w1 = Variable(torch.randn(D_in, H), requires_grad=True)
w2 = Variable(torch.randn(H, D_out), requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    if w1.grad: w1.grad.data.zero_()
    if w2.grad: w2.grad.data.zero_()
    loss.backward()

    w1.data -= learning_rate * w1.grad.data
    w2.data -= learning_rate * w2.grad.data
```

PyTorch Introduction

PyTorch: Autograd

Forward pass looks exactly the same as the Tensor version, but everything is a variable now



```
import torch
from torch.autograd import Variable

N, D_in, H, D_out = 64, 1000, 100, 10
x = Variable(torch.randn(N, D_in), requires_grad=False)
y = Variable(torch.randn(N, D_out), requires_grad=False)
w1 = Variable(torch.randn(D_in, H), requires_grad=True)
w2 = Variable(torch.randn(H, D_out), requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    if w1.grad: w1.grad.data.zero_()
    if w2.grad: w2.grad.data.zero_()
    loss.backward()

    w1.data -= learning_rate * w1.grad.data
    w2.data -= learning_rate * w2.grad.data
```

PyTorch Introduction

PyTorch: Autograd

Compute gradient of loss
with respect to w1 and w2
(zero out grads first)

```
import torch
from torch.autograd import Variable

N, D_in, H, D_out = 64, 1000, 100, 10
x = Variable(torch.randn(N, D_in), requires_grad=False)
y = Variable(torch.randn(N, D_out), requires_grad=False)
w1 = Variable(torch.randn(D_in, H), requires_grad=True)
w2 = Variable(torch.randn(H, D_out), requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    if w1.grad: w1.grad.data.zero_()
    if w2.grad: w2.grad.data.zero_()
    loss.backward()

    w1.data -= learning_rate * w1.grad.data
    w2.data -= learning_rate * w2.grad.data
```

PyTorch Introduction

PyTorch: Autograd

```
import torch
from torch.autograd import Variable

N, D_in, H, D_out = 64, 1000, 100, 10
x = Variable(torch.randn(N, D_in), requires_grad=False)
y = Variable(torch.randn(N, D_out), requires_grad=False)
w1 = Variable(torch.randn(D_in, H), requires_grad=True)
w2 = Variable(torch.randn(H, D_out), requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    if w1.grad: w1.grad.data.zero_()
    if w2.grad: w2.grad.data.zero_()
    loss.backward()

    w1.data -= learning_rate * w1.grad.data
    w2.data -= learning_rate * w2.grad.data
```

Make gradient step on weights



PyTorch Introduction

PyTorch: New Autograd Functions

Define your own autograd functions by writing forward and backward for Tensors

(similar to modular layers in A2)

```
class ReLU(torch.autograd.Function):
    def forward(self, x):
        self.save_for_backward(x)
        return x.clamp(min=0)

    def backward(self, grad_y):
        x, = self.saved_tensors
        grad_input = grad_y.clone()
        grad_input[x < 0] = 0
        return grad_input
```

PyTorch Introduction

PyTorch: New Autograd Functions

```
class ReLU(torch.autograd.Function):
    def forward(self, x):
        self.save_for_backward(x)
        return x.clamp(min=0)

    def backward(self, grad_y):
        x, = self.saved_tensors
        grad_input = grad_y.clone()
        grad_input[x < 0] = 0
        return grad_input
```

```
N, D_in, H, D_out = 64, 1000, 100, 10
x = Variable(torch.randn(N, D_in), requires_grad=False)
y = Variable(torch.randn(N, D_out), requires_grad=False)
w1 = Variable(torch.randn(D_in, H), requires_grad=True)
w2 = Variable(torch.randn(H, D_out), requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    relu = ReLU()
    y_pred = relu(x.mm(w1)).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    if w1.grad: w1.grad.data.zero_()
    if w2.grad: w2.grad.data.zero_()
    loss.backward()

    w1.data -= learning_rate * w1.grad.data
    w2.data -= learning_rate * w2.grad.data
```

Can use our new autograd
function in the forward pass



PyTorch Introduction

PyTorch: nn

Higher-level wrapper for working with neural nets

Similar to Keras and friends ...
but only one, and it's good =)

```
import torch
from torch.autograd import Variable

N, D_in, H, D_out = 64, 1000, 100, 10
x = Variable(torch.randn(N, D_in))
y = Variable(torch.randn(N, D_out), requires_grad=False)

model = torch.nn.Sequential(
    torch.nn.Linear(D_in, H),
    torch.nn.ReLU(),
    torch.nn.Linear(H, D_out))
loss_fn = torch.nn.MSELoss(size_average=False)

learning_rate = 1e-4
for t in range(500):
    y_pred = model(x)
    loss = loss_fn(y_pred, y)

    model.zero_grad()
    loss.backward()

    for param in model.parameters():
        param.data -= learning_rate * param.grad.data
```

PyTorch Introduction

PyTorch: nn

Define our model as a sequence of layers

nn also defines common loss functions

```
import torch
from torch.autograd import Variable

N, D_in, H, D_out = 64, 1000, 100, 10
x = Variable(torch.randn(N, D_in))
y = Variable(torch.randn(N, D_out), requires_grad=False)

model = torch.nn.Sequential(
    torch.nn.Linear(D_in, H),
    torch.nn.ReLU(),
    torch.nn.Linear(H, D_out))

loss_fn = torch.nn.MSELoss(size_average=False)

learning_rate = 1e-4
for t in range(500):
    y_pred = model(x)
    loss = loss_fn(y_pred, y)

    model.zero_grad()
    loss.backward()

    for param in model.parameters():
        param.data -= learning_rate * param.grad.data
```

PyTorch Introduction

PyTorch: nn

Forward pass: feed data
to model, and prediction
to loss function

```
import torch
from torch.autograd import Variable

N, D_in, H, D_out = 64, 1000, 100, 10
x = Variable(torch.randn(N, D_in))
y = Variable(torch.randn(N, D_out), requires_grad=False)

model = torch.nn.Sequential(
    torch.nn.Linear(D_in, H),
    torch.nn.ReLU(),
    torch.nn.Linear(H, D_out))
loss_fn = torch.nn.MSELoss(size_average=False)

learning_rate = 1e-4
for t in range(500):
    y_pred = model(x)
    loss = loss_fn(y_pred, y)

    model.zero_grad()
    loss.backward()

    for param in model.parameters():
        param.data -= learning_rate * param.grad.data
```

PyTorch Introduction

PyTorch: nn

```
import torch
from torch.autograd import Variable

N, D_in, H, D_out = 64, 1000, 100, 10
x = Variable(torch.randn(N, D_in))
y = Variable(torch.randn(N, D_out), requires_grad=False)

model = torch.nn.Sequential(
    torch.nn.Linear(D_in, H),
    torch.nn.ReLU(),
    torch.nn.Linear(H, D_out))
loss_fn = torch.nn.MSELoss(size_average=False)

learning_rate = 1e-4
for t in range(500):
    y_pred = model(x)
    loss = loss_fn(y_pred, y)

    model.zero_grad()
    loss.backward()

    for param in model.parameters():
        param.data -= learning_rate * param.grad.data
```

Backward pass:
compute all gradients



```
model.zero_grad()
loss.backward()
```

```
for param in model.parameters():
    param.data -= learning_rate * param.grad.data
```

PyTorch Introduction

PyTorch: nn

```
import torch
from torch.autograd import Variable

N, D_in, H, D_out = 64, 1000, 100, 10
x = Variable(torch.randn(N, D_in))
y = Variable(torch.randn(N, D_out), requires_grad=False)

model = torch.nn.Sequential(
    torch.nn.Linear(D_in, H),
    torch.nn.ReLU(),
    torch.nn.Linear(H, D_out))
loss_fn = torch.nn.MSELoss(size_average=False)

learning_rate = 1e-4
for t in range(500):
    y_pred = model(x)
    loss = loss_fn(y_pred, y)

    model.zero_grad()
    loss.backward()

    for param in model.parameters():
        param.data -= learning_rate * param.grad.data
```

Make gradient step on
each model parameter



PyTorch Introduction

PyTorch: optim

Use an **optimizer** for different update rules



```
import torch
from torch.autograd import Variable

N, D_in, H, D_out = 64, 1000, 100, 10

x = Variable(torch.randn(N, D_in))
y = Variable(torch.randn(N, D_out), requires_grad=False)

model = torch.nn.Sequential(
    torch.nn.Linear(D_in, H),
    torch.nn.ReLU(),
    torch.nn.Linear(H, D_out))
loss_fn = torch.nn.MSELoss(size_average=False)

learning_rate = 1e-4
optimizer = torch.optim.Adam(model.parameters(),
                             lr=learning_rate)

for t in range(500):
    y_pred = model(x)
    loss = loss_fn(y_pred, y)

    optimizer.zero_grad()
    loss.backward()

    optimizer.step()
```

PyTorch Introduction

PyTorch: optim

```
import torch
from torch.autograd import Variable

N, D_in, H, D_out = 64, 1000, 100, 10

x = Variable(torch.randn(N, D_in))
y = Variable(torch.randn(N, D_out), requires_grad=False)

model = torch.nn.Sequential(
    torch.nn.Linear(D_in, H),
    torch.nn.ReLU(),
    torch.nn.Linear(H, D_out))
loss_fn = torch.nn.MSELoss(size_average=False)

learning_rate = 1e-4
optimizer = torch.optim.Adam(model.parameters(),
                             lr=learning_rate)

for t in range(500):
    y_pred = model(x)
    loss = loss_fn(y_pred, y)

    optimizer.zero_grad()
    loss.backward()

    optimizer.step()
```

Update all parameters
after computing gradients

optimizer.step()

PyTorch Introduction

PyTorch: nn Define new Modules

A PyTorch **Module** is a neural net layer; it inputs and outputs **Variables**

Modules can contain weights (as Variables) or other Modules

You can define your own Modules using autograd!

```
import torch
from torch.autograd import Variable

class TwoLayerNet(torch.nn.Module):
    def __init__(self, D_in, H, D_out):
        super(TwoLayerNet, self).__init__()
        self.linear1 = torch.nn.Linear(D_in, H)
        self.linear2 = torch.nn.Linear(H, D_out)

    def forward(self, x):
        h_relu = self.linear1(x).clamp(min=0)
        y_pred = self.linear2(h_relu)
        return y_pred

N, D_in, H, D_out = 64, 1000, 100, 10

x = Variable(torch.randn(N, D_in))
y = Variable(torch.randn(N, D_out), requires_grad=False)

model = TwoLayerNet(D_in, H, D_out)

criterion = torch.nn.MSELoss(size_average=False)
optimizer = torch.optim.SGD(model.parameters(), lr=1e-4)
for t in range(500):
    y_pred = model(x)
    loss = criterion(y_pred, y)

    optimizer.zero_grad()
    loss.backward()
    optimizer.step()
```

PyTorch Introduction

PyTorch: nn Define new Modules

Define our whole model
as a single Module



```
import torch
from torch.autograd import Variable

class TwoLayerNet(torch.nn.Module):
    def __init__(self, D_in, H, D_out):
        super(TwoLayerNet, self).__init__()
        self.linear1 = torch.nn.Linear(D_in, H)
        self.linear2 = torch.nn.Linear(H, D_out)

    def forward(self, x):
        h_relu = self.linear1(x).clamp(min=0)
        y_pred = self.linear2(h_relu)
        return y_pred

N, D_in, H, D_out = 64, 1000, 100, 10

x = Variable(torch.randn(N, D_in))
y = Variable(torch.randn(N, D_out), requires_grad=False)

model = TwoLayerNet(D_in, H, D_out)

criterion = torch.nn.MSELoss(size_average=False)
optimizer = torch.optim.SGD(model.parameters(), lr=1e-4)
for t in range(500):
    y_pred = model(x)
    loss = criterion(y_pred, y)

    optimizer.zero_grad()
    loss.backward()
    optimizer.step()
```

PyTorch Introduction

PyTorch: nn Define new Modules

Initializer sets up two children (Modules can contain modules)



```
import torch
from torch.autograd import Variable

class TwoLayerNet(torch.nn.Module):
    def __init__(self, D_in, H, D_out):
        super(TwoLayerNet, self).__init__()
        self.linear1 = torch.nn.Linear(D_in, H)
        self.linear2 = torch.nn.Linear(H, D_out)

    def forward(self, x):
        h_relu = self.linear1(x).clamp(min=0)
        y_pred = self.linear2(h_relu)
        return y_pred

N, D_in, H, D_out = 64, 1000, 100, 10

x = Variable(torch.randn(N, D_in))
y = Variable(torch.randn(N, D_out), requires_grad=False)

model = TwoLayerNet(D_in, H, D_out)

criterion = torch.nn.MSELoss(size_average=False)
optimizer = torch.optim.SGD(model.parameters(), lr=1e-4)
for t in range(500):
    y_pred = model(x)
    loss = criterion(y_pred, y)

    optimizer.zero_grad()
    loss.backward()
    optimizer.step()
```

PyTorch Introduction

PyTorch: nn Define new Modules

Define forward pass using
child modules and
autograd ops on Variables

No need to define
backward - autograd will
handle it



```
import torch
from torch.autograd import Variable

class TwoLayerNet(torch.nn.Module):
    def __init__(self, D_in, H, D_out):
        super(TwoLayerNet, self).__init__()
        self.linear1 = torch.nn.Linear(D_in, H)
        self.linear2 = torch.nn.Linear(H, D_out)

    def forward(self, x):
        h_relu = self.linear1(x).clamp(min=0)
        y_pred = self.linear2(h_relu)
        return y_pred

N, D_in, H, D_out = 64, 1000, 100, 10

x = Variable(torch.randn(N, D_in))
y = Variable(torch.randn(N, D_out), requires_grad=False)

model = TwoLayerNet(D_in, H, D_out)

criterion = torch.nn.MSELoss(size_average=False)
optimizer = torch.optim.SGD(model.parameters(), lr=1e-4)
for t in range(500):
    y_pred = model(x)
    loss = criterion(y_pred, y)

    optimizer.zero_grad()
    loss.backward()
    optimizer.step()
```

PyTorch Introduction

PyTorch: nn Define new Modules

Construct and train an instance of our model



```
import torch
from torch.autograd import Variable

class TwoLayerNet(torch.nn.Module):
    def __init__(self, D_in, H, D_out):
        super(TwoLayerNet, self).__init__()
        self.linear1 = torch.nn.Linear(D_in, H)
        self.linear2 = torch.nn.Linear(H, D_out)

    def forward(self, x):
        h_relu = self.linear1(x).clamp(min=0)
        y_pred = self.linear2(h_relu)
        return y_pred

N, D_in, H, D_out = 64, 1000, 100, 10

x = Variable(torch.randn(N, D_in))
y = Variable(torch.randn(N, D_out), requires_grad=False)

model = TwoLayerNet(D_in, H, D_out)

criterion = torch.nn.MSELoss(size_average=False)
optimizer = torch.optim.SGD(model.parameters(), lr=1e-4)
for t in range(500):
    y_pred = model(x)
    loss = criterion(y_pred, y)

    optimizer.zero_grad()
    loss.backward()
    optimizer.step()
```

PyTorch Introduction

PyTorch: DataLoaders

A **DataLoader** wraps a **Dataset** and provides minibatching, shuffling, multithreading, for you

When you need to load custom data, just write your own Dataset class

```
import torch
from torch.autograd import Variable
from torch.utils.data import TensorDataset, DataLoader

N, D_in, H, D_out = 64, 1000, 100, 10

x = torch.randn(N, D_in)
y = torch.randn(N, D_out)

loader = DataLoader(TensorDataset(x, y), batch_size=8)

model = TwoLayerNet(D_in, H, D_out)

criterion = torch.nn.MSELoss(size_average=False)
optimizer = torch.optim.SGD(model.parameters(), lr=1e-4)
for epoch in range(10):
    for x_batch, y_batch in loader:
        x_var, y_var = Variable(x), Variable(y)
        y_pred = model(x_var)
        loss = criterion(y_pred, y_var)

        optimizer.zero_grad()
        loss.backward()
        optimizer.step()
```



PyTorch Introduction

PyTorch: DataLoaders

Iterate over loader to form minibatches

Loader gives Tensors so you need to wrap in Variables

```
import torch
from torch.autograd import Variable
from torch.utils.data import TensorDataset, DataLoader

N, D_in, H, D_out = 64, 1000, 100, 10

x = torch.randn(N, D_in)
y = torch.randn(N, D_out)

loader = DataLoader(TensorDataset(x, y), batch_size=8)

model = TwoLayerNet(D_in, H, D_out)

criterion = torch.nn.MSELoss(size_average=False)
optimizer = torch.optim.SGD(model.parameters(), lr=1e-4)

for epoch in range(10):
    for x_batch, y_batch in loader:
        x_var, y_var = Variable(x), Variable(y)
        y_pred = model(x_var)
        loss = criterion(y_pred, y_var)

        optimizer.zero_grad()
        loss.backward()
        optimizer.step()
```

PyTorch Introduction

PyTorch: Pretrained Models

Super easy to use pretrained models with torchvision
<https://github.com/pytorch/vision>

```
import torch
import torchvision

alexnet = torchvision.models.alexnet(pretrained=True)
vgg16 = torchvision.models.vgg16(pretrained=True)
resnet101 = torchvision.models.resnet101(pretrained=True)
```

PyTorch Introduction

What's the advantage of PyTorch over
Tensorflow?

PyTorch Introduction

Static vs Dynamic Graphs

TensorFlow: Build graph once, then run many times (**static**)

```
N, D, H = 64, 1000, 100
x = tf.placeholder(tf.float32, shape=(N, D))
y = tf.placeholder(tf.float32, shape=(N, D))
w1 = tf.Variable(tf.random_normal((D, H)))
w2 = tf.Variable(tf.random_normal((H, D)))

h = tf.maximum(tf.matmul(x, w1), 0)
y_pred = tf.matmul(h, w2)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_sum(diff ** 2, axis=1))
grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

learning_rate = 1e-5
new_w1 = w1.assign(w1 - learning_rate * grad_w1)
new_w2 = w2.assign(w2 - learning_rate * grad_w2)
updates = tf.group(new_w1, new_w2)

with tf.Session() as sess:
    sess.run(tf.global_variables_initializer())
    values = {x: np.random.randn(N, D),
              y: np.random.randn(N, D)}
    losses = []
    for t in range(50):
        loss_val, _ = sess.run([loss, updates],
                              feed_dict=values)
```

Build graph

Run each iteration

PyTorch: Each forward pass defines a new graph (**dynamic**)

```
import torch
from torch.autograd import Variable

N, D_in, H, D_out = 64, 1000, 100, 10
x = Variable(torch.randn(N, D_in), requires_grad=False)
y = Variable(torch.randn(N, D_out), requires_grad=False)
w1 = Variable(torch.randn(D_in, H), requires_grad=True)
w2 = Variable(torch.randn(H, D_out), requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    if w1.grad: w1.grad.data.zero_()
    if w2.grad: w2.grad.data.zero_()
    loss.backward()

    w1.data -= learning_rate * w1.grad.data
    w2.data -= learning_rate * w2.grad.data
```

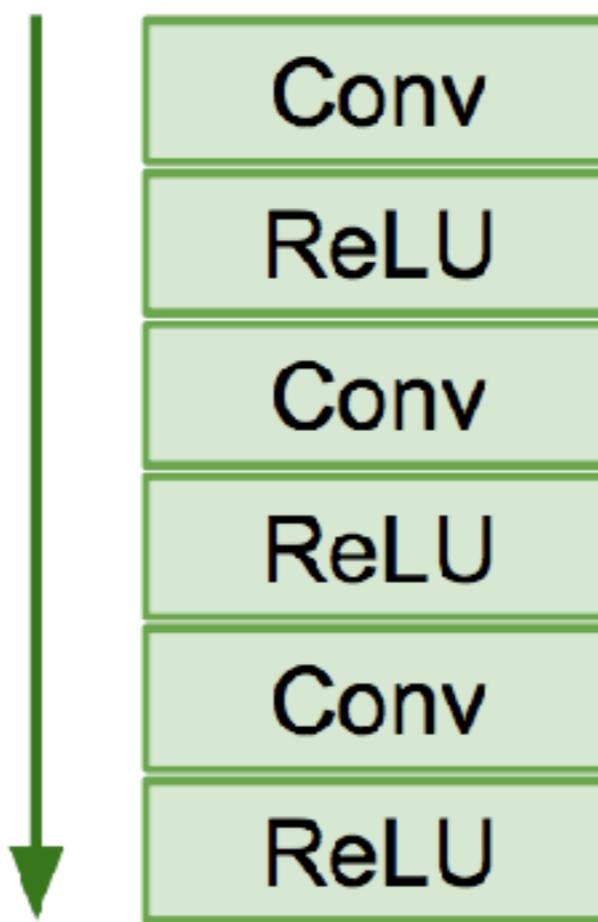
New graph each iteration

PyTorch Introduction

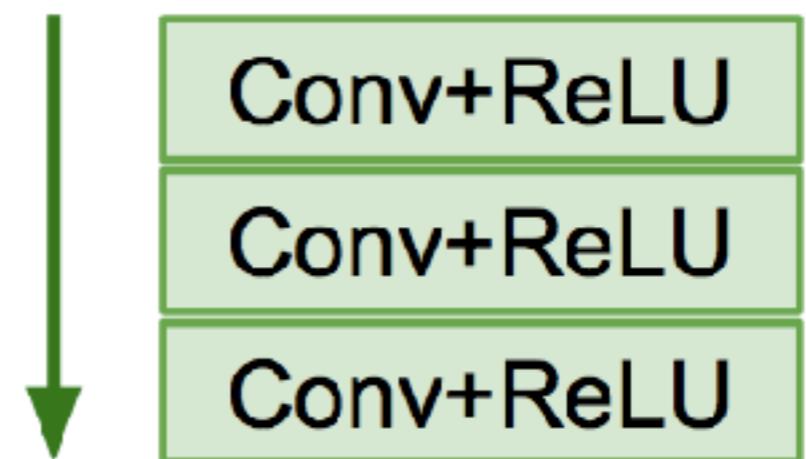
Static vs Dynamic: Optimization

With static graphs, framework can **optimize** the graph for you before it runs!

The graph you wrote



Equivalent graph with **fused operations**



PyTorch Introduction

Static vs Dynamic: Serialization

Static

Once graph is built, can **serialize** it and run it without the code that built the graph!

Dynamic

Graph building and execution are intertwined, so always need to keep code around

PyTorch Introduction

Static vs Dynamic: Conditional

$$y = \begin{cases} w_1 * x & \text{if } z > 0 \\ w_2 * x & \text{otherwise} \end{cases}$$

PyTorch Introduction

Static vs Dynamic: Conditional

$$y = \begin{cases} w1 * x & \text{if } z > 0 \\ w2 * x & \text{otherwise} \end{cases}$$

PyTorch: Normal Python

```
N, D, H = 3, 4, 5

x = Variable(torch.randn(N, D))
w1 = Variable(torch.randn(D, H))
w2 = Variable(torch.randn(D, H))

z = 10
if z > 0:
    y = x.mm(w1)
else:
    y = x.mm(w2)
```

PyTorch Introduction

Static vs Dynamic: Conditional

$$y = \begin{cases} w_1 * x & \text{if } z > 0 \\ w_2 * x & \text{otherwise} \end{cases}$$

PyTorch: Normal Python

```
N, D, H = 3, 4, 5

x = Variable(torch.randn(N, D))
w1 = Variable(torch.randn(D, H))
w2 = Variable(torch.randn(D, H))

z = 10
if z > 0:
    y = x.mm(w1)
else:
    y = x.mm(w2)
```

TensorFlow: Special TF control flow operator!

```
N, D, H = 3, 4, 5
x = tf.placeholder(tf.float32, shape=(N, D))
z = tf.placeholder(tf.float32, shape=None)
w1 = tf.placeholder(tf.float32, shape=(D, H))
w2 = tf.placeholder(tf.float32, shape=(D, H))

def f1(): return tf.matmul(x, w1)
def f2(): return tf.matmul(x, w2)
y = tf.cond(tf.less(z, 0), f1, f2)

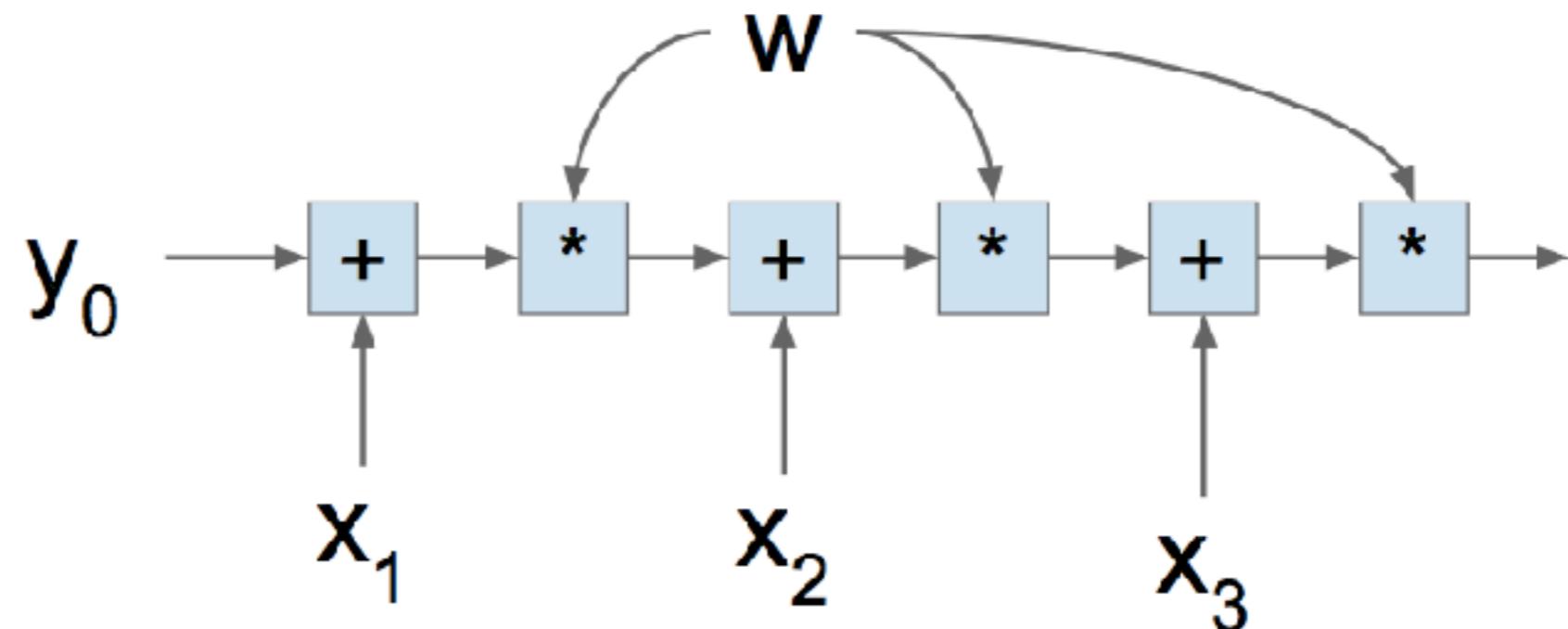
with tf.Session() as sess:
    values = {
        x: np.random.randn(N, D),
        z: 10,
        w1: np.random.randn(D, H),
        w2: np.random.randn(D, H),
    }
    y_val = sess.run(y, feed_dict=values)
```



PyTorch Introduction

Static vs Dynamic: Loops

$$y_t = (y_{t-1} + x_t) * w$$



PyTorch Introduction

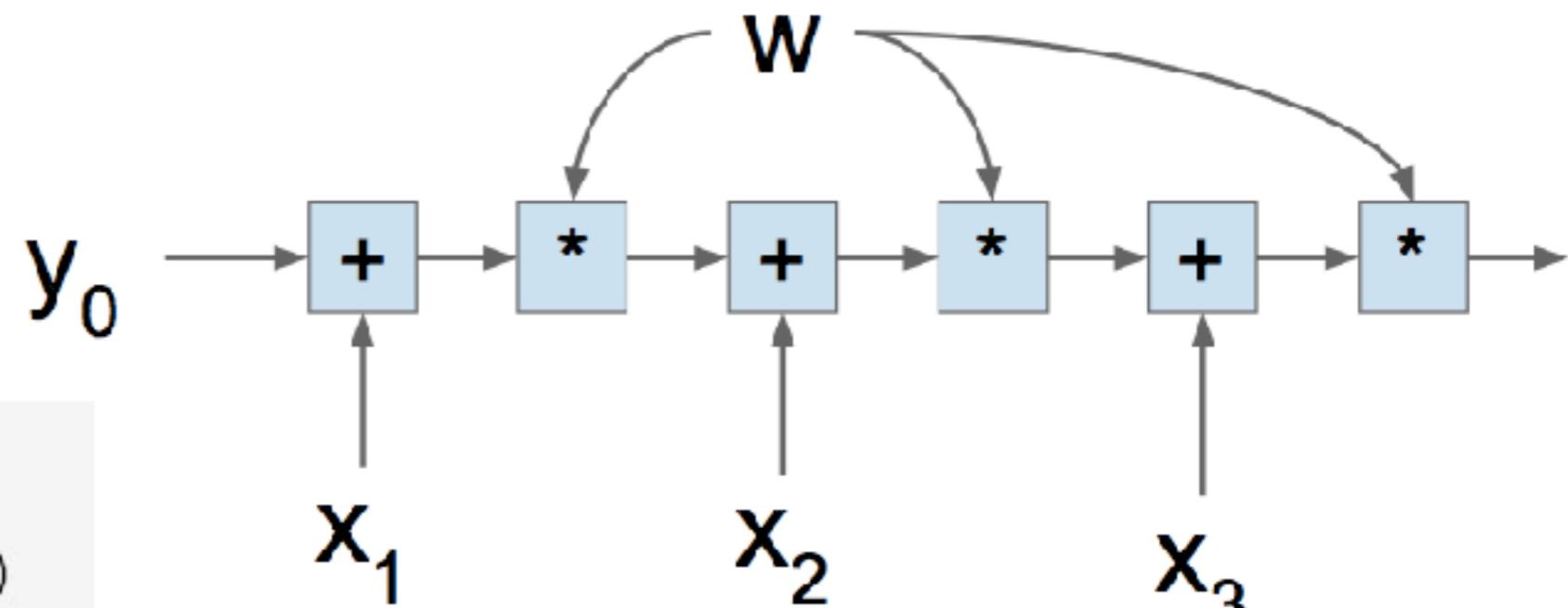
Static vs Dynamic: Loops

$$y_t = (y_{t-1} + x_t) * w$$

PyTorch: Normal Python

```
T, D = 3, 4
y0 = Variable(torch.randn(D))
x = Variable(torch.randn(T, D))
w = Variable(torch.randn(D))

y = [y0]
for t in range(T):
    prev_y = y[-1]
    next_y = (prev_y + x[t]) * w
    y.append(next_y)
```



PyTorch Introduction

Static vs Dynamic: Loops

$$y_t = (y_{t-1} + x_t) * w$$

PyTorch: Normal Python

```
T, D = 3, 4
y0 = Variable(torch.randn(D))
x = Variable(torch.randn(T, D))
w = Variable(torch.randn(D))

y = [y0]
for t in range(T):
    prev_y = y[-1]
    next_y = (prev_y + x[t]) * w
    y.append(next_y)
```

TensorFlow: Special TF control flow

```
T, N, D = 3, 4, 5
x = tf.placeholder(tf.float32, shape=(T, D))
y0 = tf.placeholder(tf.float32, shape=(D,))
w = tf.placeholder(tf.float32, shape=(D,))

def f(prev_y, cur_x):
    return (prev_y + cur_x) * w

y = tf.foldl(f, x, y0)

with tf.Session() as sess:
    values = {
        x: np.random.randn(T, D),
        y0: np.random.randn(D),
        w: np.random.randn(D),
    }
    y_val = sess.run(y, feed_dict=values)
```

PyTorch Introduction

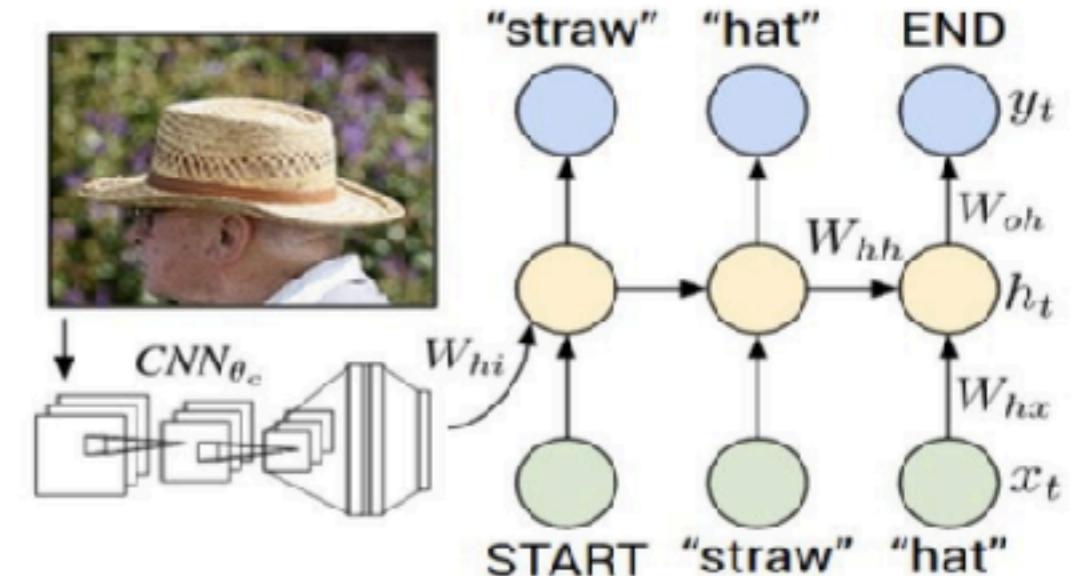
Dynamic Graphs in TensorFlow

TensorFlow Fold make dynamic graphs easier in TensorFlow through **dynamic batching**

PyTorch Introduction

Dynamic Graph Applications

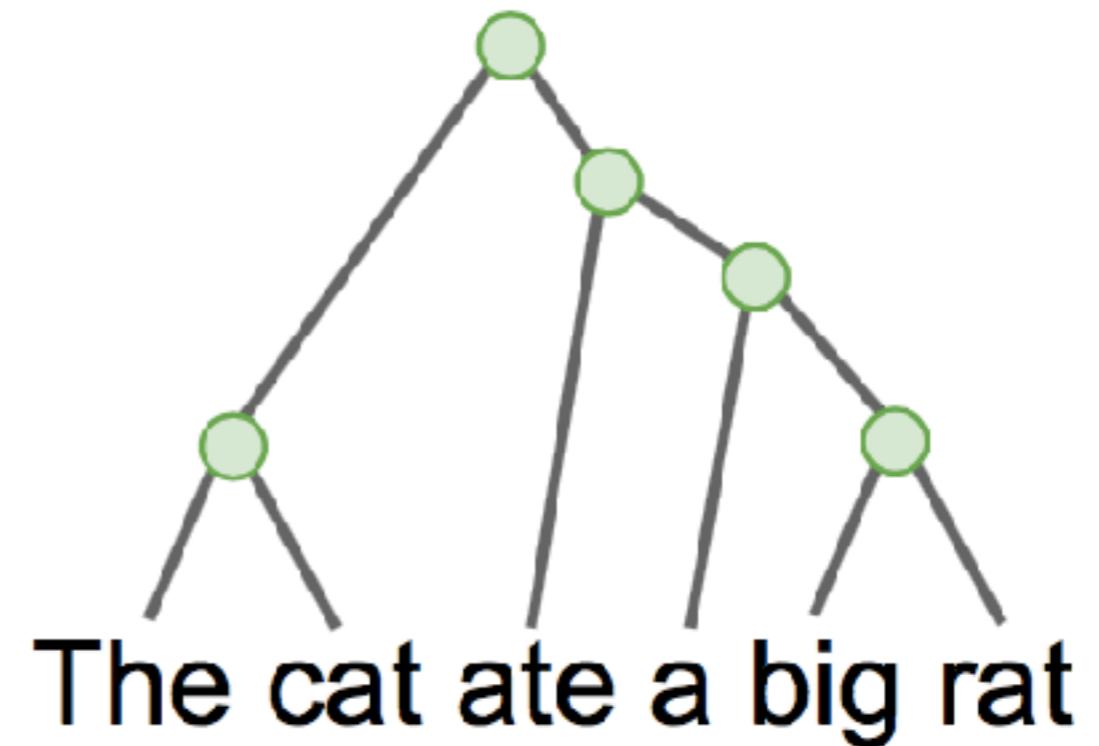
- Recurrent networks



PyTorch Introduction

Dynamic Graph Applications

- Recurrent networks
- Recursive networks

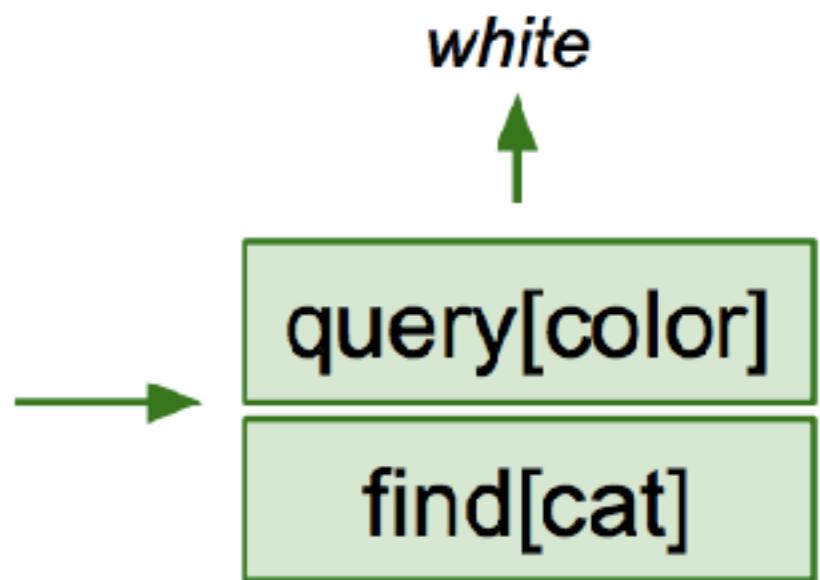


PyTorch Introduction

Dynamic Graph Applications

- Recurrent networks
- Recursive networks
- Modular Networks

*What color
is the cat?*

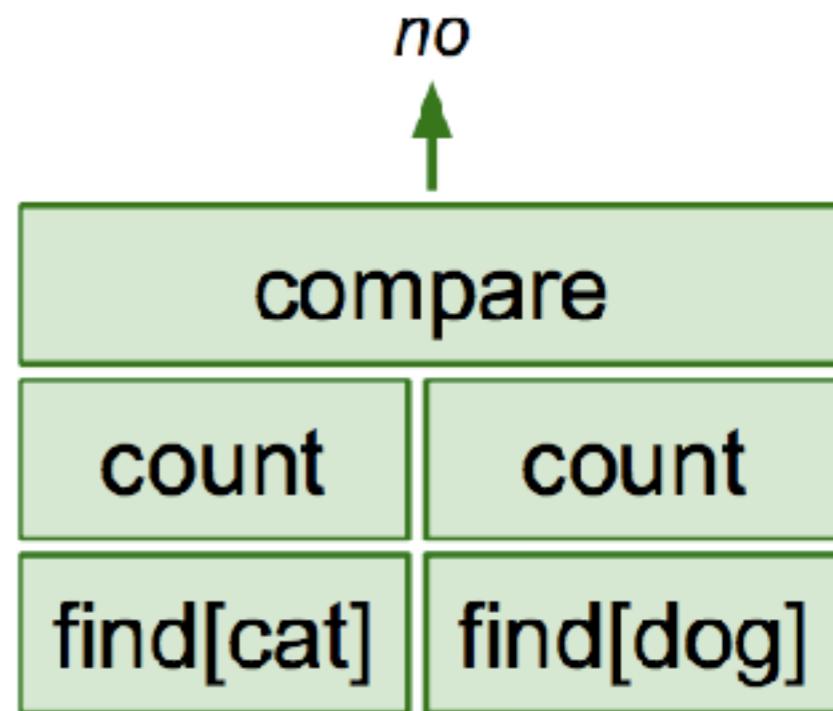


PyTorch Introduction

Dynamic Graph Applications

- Recurrent networks
- Recursive networks
- Modular Networks

*Are there
more cats
than dogs?*



PyTorch Introduction

Dynamic graph application:

- Recurrent networks
- Recursive networks
- Modular networks
- More creative ideas...