# RNN LSTM and Deep Learning Libraries

## UDRC Summer School
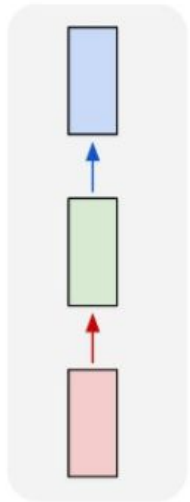
Muhammad Awais
m.a.rana@surrey.ac.uk

# Outline

➢ Recurrent Neural Network
➢ Application of RNN
➢ LSTM
➢ Caffe
➢ Torch
➢ Theano
➢ TensorFlow

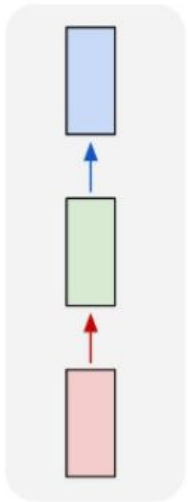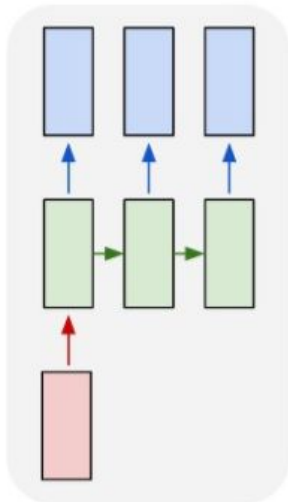# Flexibility of Recurrent Neural Networks

one to one



**Vanilla Neural Networks**

# Flexibility of Recurrent Neural Networks
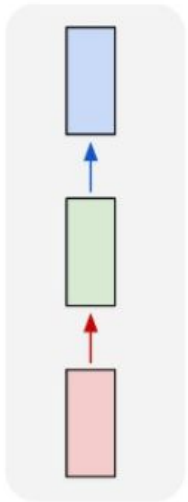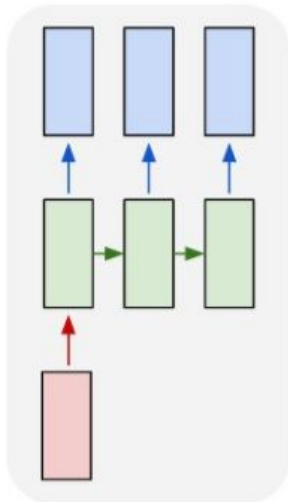
one to one

one to many

e.g. **Image Captioning**
image -> sequence of words

# Flexibility of Recurrent Neural Networks



one to one   one to many   many to one

e.g. **Sentiment Classification**
sequence of words -> sentiment

# Flexibility of Recurrent Neural Networks



e.g. **Machine Translation**
seq of words -> seq of words

# Flexibility of Recurrent Neural Networks



one to one | one to many | many to one | many to many | many to many

e.g. **Video classification on frame level**

# Recurrent Neural Networks

# Recurrent Neural Networks



y

RNN

x

usually want to predict a vector at some time steps

# Recurrent Neural Networks

We can process a sequence of vectors **x** by applying a recurrence formula at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$

new state
some function with parameters W
old state
input vector at some time step

# Recurrent Neural Networks

We can process a sequence of vectors **x** by
applying a recurrence formula at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$

new state      old state    input vector at
                                     some time step

some function
with parameters W

Notice: the same function and the same set
of parameters are used at every time step.

y

RNN

x

# Recurrent Neural Networks

The state consists of a single *"hidden"* vector **h**:



$$h_t = f_W(h_{t-1}, x_t)$$

$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t)$$

$$y_t = W_{hy} h_t$$

# Recurrent Neural Networks

**Character-level language model example**

Vocabulary:
[h,e,l,o]

Example training sequence:
**"hello"**

# Recurrent Neural Networks

**Character-level language model example**

Vocabulary:
[h,e,l,o]

Example training sequence:
**"hello"**

# Recurrent Neural Networks

**Character-level language model example**

Vocabulary: [h,e,l,o]

Example training sequence: **"hello"**

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

# Recurrent Neural Networks

**Character-level language model example**

Vocabulary: [h,e,l,o]

Example training sequence: **"hello"**

# Recurrent Neural Networks

## Image Captioning



Explain Images with Multimodal Recurrent Neural Networks, Mao et al.
Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei
Show and Tell: A Neural Image Caption Generator, Vinyals et al.
Long-term Recurrent Convolutional Networks for Visual Recognition and Description, Donahue et al.
Learning a Recurrent Visual Representation for Image Caption Generation, Chen and Zitnick

# Recurrent Neural Networks

**Recurrent Neural Network**



**Convolutional Neural Network**

# Recurrent Neural Networks



test image

# Recurrent Neural Networks

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096
FC-1000
softmax

test image

# Recurrent Neural Networks

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096
FC-1000
softmax

test image

# Recurrent Neural Networks



image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096

test image

x0
<START>

<START>

# Recurrent Neural Networks

test image



**before:**

$h = \tanh(Wxh * x + Whh * h)$

**now:**

$h = \tanh(Wxh * x + Whh * h + \mathbf{Wih * v})$

**Wih**

v

y0

h0

x0
<STA
RT>

<START>

# Recurrent Neural Networks

# Recurrent Neural Networks



test image

# Recurrent Neural Networks

# Recurrent Neural Networks



test image

# Recurrent Neural Networks



test image

sample
<END> token
=> finish.

<START>

# Recurrent Neural Networks

# Image Sentence Datasets

a man riding a bike on a dirt path through a forest.
bicyclist raises his fist as he rides on desert dirt trail.
this dirt bike rider is smiling and raising his fist in triumph.
a man riding a bicycle while pumping his fist in the air.
a mountain biker pumps his fist in celebration.



Microsoft COCO
*[Tsung-Yi Lin et al. 2014]*
mscoco.org

currently:
~120K images
~5 sentences each

# Recurrent Neural Networks



"man in black shirt is playing guitar."

"construction worker in orange safety vest is working on road."

"two young girls are playing with lego toy."

"boy is doing backflip on wakeboard."

# Recurrent Neural Networks



"man in black shirt is playing guitar."

"construction worker in orange safety vest is working on road."

"two young girls are playing with lego toy."

"boy is doing backflip on wakeboard."

"a young boy is holding a baseball bat."

"a cat is sitting on a couch with a remote control."

"a woman holding a teddy bear in front of a mirror."

"a horse is standing in the middle of a road."

# Recurrent Neural Networks

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$h \in \mathbb{R}^n.$     $W^l \; [n \times 2n]$



depth

time

# Recurrent Neural Networks

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$h \in \mathbb{R}^n. \qquad W^l \ [n \times 2n]$

LSTM:

$W^l \ [4n \times 2n]$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

depth

time

# Long Short Term Memory (LSTM)



$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

x

h

# Long Short Term Memory (LSTM)

*[Hochreiter et al., 1997]*

vector from below (**x**)

vector from before (**h**)

W

4n x 2n

sigmoid → i

sigmoid → f

sigmoid → o

tanh → g

4n

4*n

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \mathrm{sigm} \\ \mathrm{sigm} \\ \mathrm{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^{l} \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

# Long Short Term Memory (LSTM)

**cell state c**



$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$c_t^l = \boxed{f \odot c_{t-1}^l} + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

# Long Short Term Memory (LSTM)

**cell state c**

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^{l} \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

# Long Short Term Memory (LSTM)



cell state c

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^{l} \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

# Long Short Term Memory (LSTM)



**cell state c**

higher layer, or prediction

c

h

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$

# Long Short Term Memory (LSTM)

# Long Short Term Memory (LSTM)

## Summary

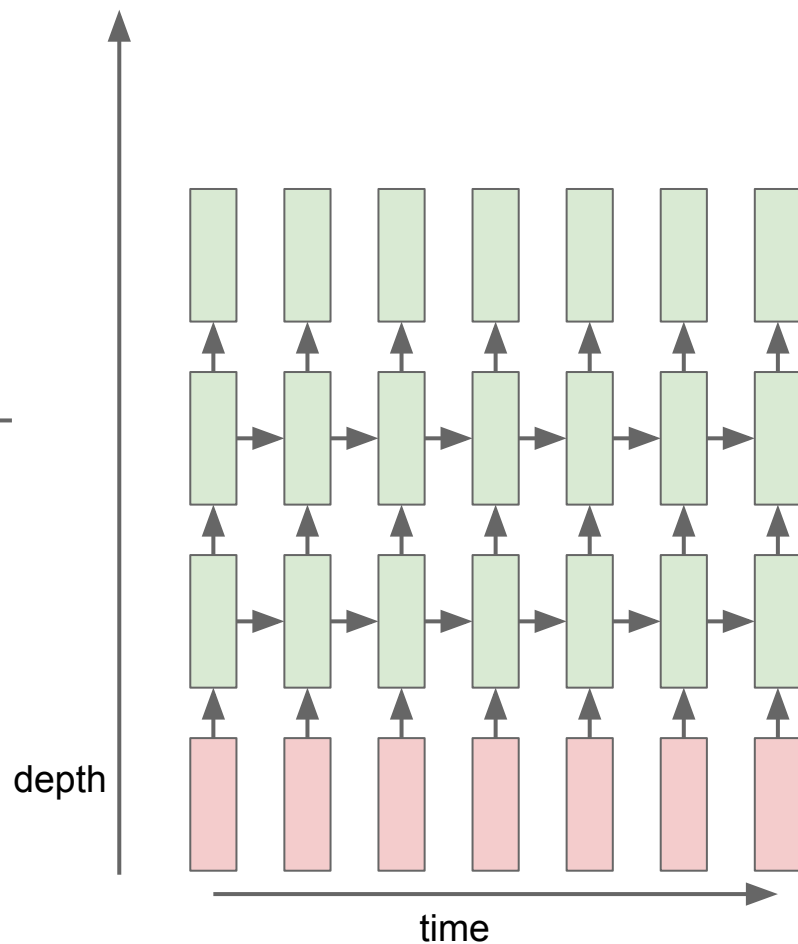- RNNs allow a lot of flexibility in architecture design
- Vanilla RNNs are simple but don't work very well
- Common to use LSTM or GRU: their additive interactions improve gradient flow
- Backward flow of gradients in RNN can explode or vanish. Exploding is controlled with gradient clipping. Vanishing is controlled with additive interactions (LSTM)
- Better/simpler architectures are a hot topic of current research
- Better understanding (both theoretical and empirical) is needed.

# Deep Learning Libraries
## Caffe, Torch, Theano, TensorFlow

# Caffe

http://caffe.berkeleyvision.org

# Caffe overview

From U.C. Berkeley
Written in C++
Has Python and MATLAB bindings
Good for training or finetuning feedforward models

# Caffe

## **Main classes**

**Blob**: Stores data and derivatives (header source)

**Layer**: Transforms bottom blobs to top blobs (header + source)

**Net**: Many layers; computes gradients via forward / backward (header source)

**Solver**: Uses gradients to update weights (header source)



45

# Protocol Buffers

"Typed JSON"
   from Google

Define "message types" in
   .proto files

**.proto file**

```
message Person {
  required string name = 1;
  required int32 id = 2;
  optional string email = 3;
}
```

https://developers.google.com/protocol-buffers/

# Protocol Buffers

"Typed JSON"
from Google

Define "message types" in
.proto files

Serialize instances to text
files (.prototxt)

**.proto file**

```
message Person {
  required string name = 1;
  required int32 id = 2;
  optional string email = 3;
}
```

**.prototxt file**

```
name: "John Doe"
id: 1234
email: "jdoe@example.com"
```

https://developers.google.com/protocol-buffers/

# Protocol Buffers

"Typed JSON"
    from Google

Define "message types" in
    .proto files

Serialize instances to text
    files (.prototxt)

Compile classes for
    different languages

**.proto file**

```
message Person {
  required string name = 1;
  required int32 id = 2;
  optional string email = 3;
}
```

**Java class**

```
Person john = Person.newBuilder()
    .setId(1234)
    .setName("John Doe")
    .setEmail("jdoe@example.com")
    .build();
output = new FileOutputStream(args[0]);
john.writeTo(output);
```

**.prototxt file**

```
name: "John Doe"
id: 1234
email: "jdoe@example.com"
```

**C++ class**

```
Person john;
fstream input(argv[1],
    ios::in | ios::binary);
john.ParseFromIstream(&input);
id = john.id();
name = john.name();
email = john.email();
```

# Caffe

# Protocol Buffers

```
64  message NetParameter {
65    optional string name = 1; // consider giving the network a name
66    // The input blobs to the network.
67    repeated string input = 3;
68    // The shape of the input blobs.
69    repeated BlobShape input_shape = 8;
70
71    // 4D input dimensions -- deprecated.  Use "shape" instead.
72    // If specified, for each input blob there should be four
73    // values specifying the num, channels, height and width of the input blob.
74    // Thus, there should be a total of (4 * #input) numbers.
75    repeated int32 input_dim = 4;
76
77    // Whether the network will force every layer to carry out backward operation.
78    // If set False, then whether to carry out backward is determined
79    // automatically according to the net structure and learning rates.
80    optional bool force_backward = 5 [default = false];
81    // The current "state" of the network, including the phase, level, and stage.
82    // Some layers may be included/excluded depending on this state and the states
83    // specified in the layers' include and exclude fields.
84    optional NetState state = 6;
85
86    // Print debugging information about results while running Net::Forward,
87    // Net::Backward, and Net::Update.
88    optional bool debug_info = 7 [default = false];
```

```
102  message SolverParameter {
103    /////////////////////////////////////////////////////////////////////////
104    // Specifying the train and test networks
105    //
106    // Exactly one train net must be specified using one of the following fields:
107    //    train_net_param, train_net, net_param, net
108    // One or more test nets may be specified using any of the following fields:
109    //    test_net_param, test_net, net_param, net
110    // If more than one test net field is specified (e.g., both net and
111    // test_net are specified), they will be evaluated in the field order given
112    // above: (1) test_net_param, (2) test_net, (3) net_param/net.
113    // A test_iter must be specified for each test_net.
114    // A test_level and/or a test_stage may also be specified for each test_net.
115    /////////////////////////////////////////////////////////////////////////
116
117    // Proto filename for the train net, possibly combined with one or more
118    // test nets.
119    optional string net = 24;
120    // Inline train net param, possibly combined with one or more test nets.
121    optional NetParameter net_param = 25;
122
123    optional string train_net = 1; // Proto filename for the train net
```

https://github.com/BVLC/caffe/blob/master/src/caffe/proto/caffe.proto
**<- All Caffe proto types defined here, good documentation!**

# Caffe

# **Training / Finetuning**

No need to write code!

1.  Convert data (run a script)
2.  Define net (edit prototxt)
3.  Define solver (edit prototxt)
4.  Train (with pretrained weights) (run a script)

# Caffe

## Step 1: 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 yourself using h5py

# Step 2: Define Net

```
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
}
```

# Step 2: Define Net

```
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
```

Layers and Blobs often have same name!

```
  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
}
```

# Step 2: Define Net

```
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
```

Layers and Blobs often have same name!

Learning rates (weight + bias)

Regularization (weight + bias)

```
  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

# Step 2: Define Net

Number of output classes

```
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
```

Layers and Blobs often have same name!

Learning rates (weight + bias)

Regularization (weight + bias)

```
  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

## Step 2: Define Net

Number of output classes

```
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
```

Layers and Blobs often have same name!

Set these to 0 to freeze a layer

Learning rates (weight + bias)

Regularization (weight + bias)

```
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

## Step 2: Define Net

- .prototxt can get ugly for big models

- ResNet-152 prototxt is 6775 lines long!

- Not "compositional"; can't easily define a residual block and reuse

```
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: "pool5"
6750          name: "pool5"
6751          type: "Pooling"
6752          pooling_param {
6753                  kernel_size: 7
6754                  stride: 1
6755                  pool: AVE
6756          }
6757  }
6758
6759  layer {
6760          bottom: "pool5"
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  }
```

https://github.com/KaimingHe/deep-residual-networks/blob/master/prototxt/ResNet-152-deploy.prototxt

# Step 2: Define Net (finetuning)

**Original prototxt:**
```
layer {
  name: "fc7"
  type: "InnerProduct"
  inner_product_param {
    num_output: 4096
  }
}
[... ReLU, Dropout]
layer {
  name: "fc8"
  type: "InnerProduct"
  inner_product_param {
    num_output: 1000
  }
}
```

Same name:
weights copied

**Pretrained weights:**
```
"fc7.weight": [values]
"fc7.bias": [values]
"fc8.weight": [values]
"fc8.bias": [values]
```

**Modified prototxt:**
```
layer {
  name: "fc7"
  type: "InnerProduct"
  inner_product_param {
    num_output: 4096
  }
}
[... ReLU, Dropout]
layer {
  name: "my-fc8"
  type: "InnerProduct"
  inner_product_param {
    num_output: 10
  }
}
```

# Step 2: Define Net (finetuning)

**Original prototxt:**

```
layer {
  name: "fc7"
  type: "InnerProduct"
  inner_product_param {
    num_output: 4096
  }
}
[... ReLU, Dropout]
layer {
  name: "fc8"
  type: "InnerProduct"
  inner_product_param {
    num_output: 1000
  }
}
```

Same name:
weights copied

**Pretrained weights:**

```
"fc7.weight": [values]
"fc7.bias": [values]
"fc8.weight": [values]
"fc8.bias": [values]
```

**Modified prototxt:**

```
layer {
  name: "fc7"
  type: "InnerProduct"
  inner_product_param {
    num_output: 4096
  }
}
[... ReLU, Dropout]
layer {
  name: "my-fc8"
  type: "InnerProduct"
  inner_product_param {
    num_output: 10
  }
}
```

# Step 2: Define Net (finetuning)

**Original prototxt:**

```
layer {
  name: "fc7"
  type: "InnerProduct"
  inner_product_param {
    num_output: 4096
  }
}
[... ReLU, Dropout]
layer {
  name: "fc8"
  type: "InnerProduct"
  inner_product_param {
    num_output: 1000
  }
}
```

Same name:
weights copied

**Pretrained weights:**

```
"fc7.weight": [values]
"fc7.bias": [values]
"fc8.weight": [values]
"fc8.bias": [values]
```

Different name:
weights reinitialized

**Modified prototxt:**

```
layer {
  name: "fc7"
  type: "InnerProduct"
  inner_product_param {
    num_output: 4096
  }
}
[... ReLU, Dropout]
layer {
  name: "my-fc8"
  type: "InnerProduct"
  inner_product_param {
    num_output: 10
  }
}
```

# Caffe

# Step 3: Define Solver

Write a prototxt file defining a
[SolverParameter](SolverParameter)

If finetuning, copy existing
solver.prototxt file

    Change net to be your net

    Change snapshot_prefix to your
      output

    Reduce base learning rate (divide
      by 100)

    Maybe change max_iter and
      snapshot

```
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

## Step 4: Train!

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

https://github.com/BVLC/caffe/blob/master/tools/caffe.cpp

# Caffe

# Step 4: Train!

```
./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 mode

https://github.com/BVLC/caffe/blob/master/tools/caffe.cpp

# Caffe

# Step 4: Train!

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

**-gpu all** for multi-GPU data parallelism

https://github.com/BVLC/caffe/blob/master/tools/caffe.cpp

# Caffe

## Pros / Cons

(+) Good for feedforward networks
(+) Good for finetuning existing networks
(+) Train models without writing any code!
(+) Python and matlab interfaces are pretty useful!
(-) Need to write C++ / CUDA for new GPU layers
(-) Not good for recurrent networks
(-) Cumbersome for big networks (GoogLeNet, ResNet)

# Torch

# Torch

From NYU + IDIAP
Written in C and Lua
Used a lot a Facebook, DeepMind

# Torch

## Lua

High level scripting language, easy to interface with C

Similar to Javascript:

    One data structure:

        table == JS object

    Prototypical inheritance

        metatable == JS prototype

    First-class functions

Some gotchas:

    1-indexed =(

    Variables global by default =(

    Small standard library

---

### Learn Lua in 15 Minutes
*more or less*

*For a more in-depth Lua tutorial, watch [this video](this video) or check out [a transcript of the video](a transcript of the video).*

```lua
-- Two dashes start a one-line comment.

--[[
    Adding two ['s and ]'s makes it a
    multi-line comment.
--]]

----------------------------------------------------
-- 1. Variables and flow control.
----------------------------------------------------

num = 42  -- All numbers are doubles.
-- Don't freak out, 64-bit doubles have 52 bits for
-- storing exact int values; machine precision is
-- not a problem for ints that need < 52 bits.

s = 'walternate'  -- Immutable strings like Python.
t = "double-quotes are also fine"
u = [[ Double brackets
       start and end
       multi-line strings.]]
t = nil  -- Undefines t; Lua has garbage collection.

-- Blocks are denoted with keywords like do/end:
while num < 50 do
  num = num + 1  -- No ++ or += type operators.
end
```

http://tylerneylon.com/a/learn-lua/

# Torch

# Tensors

Torch tensors are just like numpy arrays

# Tensors

Torch tensors are just like numpy arrays

```python
1  import numpy as np
2
3  # Simple feedforward network (no biases) in numpy
4
5  # Batch size, input dim, hidden dim, num classes
6  N, D, H, C = 100, 1000, 100, 10
7
8  # First and second layer weights
9  w1 = np.random.randn(D, H)
10 w2 = np.random.randn(H, C)
11
12 # Random input data
13 x = np.random.randn(N, D)
14
15 # Forward pass
16 a = x.dot(w1)          # First layer
17 a = np.maximum(a, 0) # In-place ReLU
18 scores = a.dot(w2)    # Second layer
19
20 print scores
```

# Tensors

Torch tensors are just like numpy arrays

```python
1  import numpy as np
2
3  # Simple feedforward network (no biases) in numpy
4
5  # Batch size, input dim, hidden dim, num classes
6  N, D, H, C = 100, 1000, 100, 10
7
8  # First and second layer weights
9  w1 = np.random.randn(D, H)
10 w2 = np.random.randn(H, C)
11
12 # Random input data
13 x = np.random.randn(N, D)
14
15 # Forward pass
16 a = x.dot(w1)        # First layer
17 a = np.maximum(a, 0) # In-place ReLU
18 scores = a.dot(w2)   # Second layer
19
20 print scores
```

```lua
1  require 'torch'
2
3  -- Simple feedforward network (no biases) in torch
4
5  -- Batch size, input dim, hidden dim, num classes
6  local N, D, H, C = 100, 1000, 100, 10
7
8  -- First and second layer weights
9  local w1 = torch.randn(D, H)
10 local w2 = torch.randn(H, C)
11
12 -- Random input data
13 local x = torch.randn(N, D)
14
15 -- Forward pass
16 local a = torch.mm(x, w1)      -- First layer
17 a:cmax(0)                      -- In-place ReLU
18 local scores = torch.mm(a, w2) -- Second layer
19
20 print(scores)
```

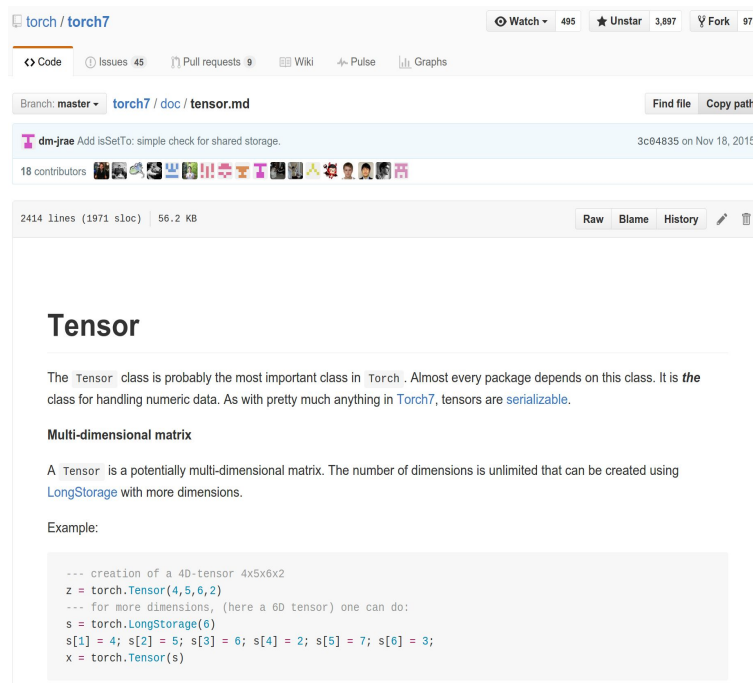# Tensors

Like numpy, can easily change data type:

```python
1 import numpy as np
2
3 |
4 # Simple feedforward network (no biases) in numpy
5
6 dtype = np.float32  # Use 32-bit floats
7
8 # Batch size, input dim, hidden dim, num classes
9 N, D, H, C = 100, 1000, 100, 10
10
11 # First and second layer weights
12 w1 = np.random.randn(D, H).astype(dtype)
13 w2 = np.random.randn(H, C).astype(dtype)
14
15 # Random input data
16 x = np.random.randn(N, D).astype(dtype)
17
18 # Forward pass
19 a = x.dot(w1)        # First layer
20 a = np.maximum(a, 0) # In-place ReLU
21 scores = a.dot(w2)   # Second layer
22
23 print scores
```

```lua
1 require 'torch'
2
3
4 -- Simple feedforward network (no biases) in torch
5
6 local dtype = 'torch.FloatTensor'  -- Use 32-bit floats
7
8 -- Batch size, input dim, hidden dim, num classes
9 local N, D, H, C = 100, 1000, 100, 10
10
11 -- First and second layer weights
12 local w1 = torch.randn(D, H):type(dtype)
13 local w2 = torch.randn(H, C):type(dtype)
14
15 -- Random input data
16 local x = torch.randn(N, D):type(dtype)
17
18 -- Forward pass
19 local a = torch.mm(x, w1)      -- First layer
20 a:cmax(0)                      -- In-place ReLU
21 local scores = torch.mm(a, w2) -- Second layer
22
23 print(scores)
```

# Tensors

Unlike numpy, GPU is just a datatype away:

```python
1  import numpy as np
2
3  |
4  # Simple feedforward network (no biases) in numpy
5
6  dtype = np.float32  # Use 32-bit floats
7
8  # Batch size, input dim, hidden dim, num classes
9  N, D, H, C = 100, 1000, 100, 10
10
11 # First and second layer weights
12 w1 = np.random.randn(D, H).astype(dtype)
13 w2 = np.random.randn(H, C).astype(dtype)
14
15 # Random input data
16 x = np.random.randn(N, D).astype(dtype)
17
18 # Forward pass
19 a = x.dot(w1)        # First layer
20 a = np.maximum(a, 0) # In-place ReLU
21 scores = a.dot(w2)   # Second layer
22
23 print scores
```

```lua
1  require 'torch'
2  require 'cutorch'
3
4  -- Simple feedforward network (no biases) in torch
5  |
6  local dtype = 'torch.CudaTensor'  -- Use CUDA
7
8  -- Batch size, input dim, hidden dim, num classes
9  local N, D, H, C = 100, 1000, 100, 10
10
11 -- First and second layer weights
12 local w1 = torch.randn(D, H):type(dtype)
13 local w2 = torch.randn(H, C):type(dtype)
14
15 -- Random input data
16 local x = torch.randn(N, D):type(dtype)
17
18 -- Forward pass
19 local a = torch.mm(x, w1)       -- First layer
20 a:cmax(0)                       -- In-place ReLU
21 local scores = torch.mm(a, w2)  -- Second layer
22
23 print(scores)
```

# Torch

# Tensors

## Documentation on GitHub:





https://github.com/torch/torch7/blob/master/doc/tensor.md

https://github.com/torch/torch7/blob/master/doc/maths.md

# Torch

## nn

nn module lets you easily build
and train neural nets

```lua
 1 require 'torch'
 2 require 'nn'
 3
 4
 5 -- Batch size, input dim, hidden dim, num classes
 6 local N, D, H, C = 100, 1000, 100, 10
 7
 8 -- Build a one-layer ReLU network
 9 local net = nn.Sequential()
10 net:add(nn.Linear(D, H))
11 net:add(nn.ReLU())
12 net:add(nn.Linear(H, C))
13
14 -- Collect all weights and gradients in a single Tensor
15 local weights, grad_weights = net:getParameters()
16
17 -- Loss functions are called "criterions"
18 local crit = nn.CrossEntropyCriterion()  -- Softmax loss
19
20 -- Generate some random input data
21 local x = torch.randn(N, D)
22 local y = torch.Tensor(N):random(C)
23
24 -- Forward pass: Compute scores and loss
25 local scores = net:forward(x)
26 local loss = crit:forward(scores, y)
27
28 -- Backward pass: compute gradients
29 grad_weights:zero()
30 local dscores = crit:backward(scores, y)
31 local dx = net:backward(x, dscores)
32
33 -- Make a gradient step
34 local learning_rate = 1e-3
35 weights:add(-learning_rate, grad_weights)
36 |
```

# Torch

## nn

nn module lets you easily build and train neural nets

Build a two-layer ReLU net

```
1  require 'torch'
2  require 'nn'
3
4
5  -- Batch size, input dim, hidden dim, num classes
6  local N, D, H, C = 100, 1000, 100, 10
7
8  -- Build a one-layer ReLU network
9  local net = nn.Sequential()
10 net:add(nn.Linear(D, H))
11 net:add(nn.ReLU())
12 net:add(nn.Linear(H, C))
13
14 -- Collect all weights and gradients in a single Tensor
15 local weights, grad_weights = net:getParameters()
16
17 -- Loss functions are called "criterions"
18 local crit = nn.CrossEntropyCriterion()  -- Softmax loss
19
20 -- Generate some random input data
21 local x = torch.randn(N, D)
22 local y = torch.Tensor(N):random(C)
23
24 -- Forward pass: Compute scores and loss
25 local scores = net:forward(x)
26 local loss = crit:forward(scores, y)
27
28 -- Backward pass: compute gradients
29 grad_weights:zero()
30 local dscores = crit:backward(scores, y)
31 local dx = net:backward(x, dscores)
32
33 -- Make a gradient step
34 local learning_rate = 1e-3
35 weights:add(-learning_rate, grad_weights)
36 |
```

# Torch

## nn

nn module lets you easily build and train neural nets

Get weights and gradient for entire network

```
1  require 'torch'
2  require 'nn'
3
4
5  -- Batch size, input dim, hidden dim, num classes
6  local N, D, H, C = 100, 1000, 100, 10
7
8  -- Build a one-layer ReLU network
9  local net = nn.Sequential()
10 net:add(nn.Linear(D, H))
11 net:add(nn.ReLU())
12 net:add(nn.Linear(H, C))
13
14 -- Collect all weights and gradients in a single Tensor
15 local weights, grad_weights = net:getParameters()
16
17 -- Loss functions are called "criterions"
18 local crit = nn.CrossEntropyCriterion()  -- Softmax loss
19
20 -- Generate some random input data
21 local x = torch.randn(N, D)
22 local y = torch.Tensor(N):random(C)
23
24 -- Forward pass: Compute scores and loss
25 local scores = net:forward(x)
26 local loss = crit:forward(scores, y)
27
28 -- Backward pass: compute gradients
29 grad_weights:zero()
30 local dscores = crit:backward(scores, y)
31 local dx = net:backward(x, dscores)
32
33 -- Make a gradient step
34 local learning_rate = 1e-3
35 weights:add(-learning_rate, grad_weights)
36 |
```

# Torch

## nn

nn module lets you easily build and train neural nets

Use a softmax loss function

```
1  require 'torch'
2  require 'nn'
3
4
5  -- Batch size, input dim, hidden dim, num classes
6  local N, D, H, C = 100, 1000, 100, 10
7
8  -- Build a one-layer ReLU network
9  local net = nn.Sequential()
10 net:add(nn.Linear(D, H))
11 net:add(nn.ReLU())
12 net:add(nn.Linear(H, C))
13
14 -- Collect all weights and gradients in a single Tensor
15 local weights, grad_weights = net:getParameters()
16
17 -- Loss functions are called "criterions"
18 local crit = nn.CrossEntropyCriterion()  -- Softmax loss
19
20 -- Generate some random input data
21 local x = torch.randn(N, D)
22 local y = torch.Tensor(N):random(C)
23
24 -- Forward pass: Compute scores and loss
25 local scores = net:forward(x)
26 local loss = crit:forward(scores, y)
27
28 -- Backward pass: compute gradients
29 grad_weights:zero()
30 local dscores = crit:backward(scores, y)
31 local dx = net:backward(x, dscores)
32
33 -- Make a gradient step
34 local learning_rate = 1e-3
35 weights:add(-learning_rate, grad_weights)
36 |
```

# Torch

## nn

nn module lets you easily build and train neural nets

Generate random data →

```
1  require 'torch'
2  require 'nn'
3
4
5  -- Batch size, input dim, hidden dim, num classes
6  local N, D, H, C = 100, 1000, 100, 10
7
8  -- Build a one-layer ReLU network
9  local net = nn.Sequential()
10 net:add(nn.Linear(D, H))
11 net:add(nn.ReLU())
12 net:add(nn.Linear(H, C))
13
14 -- Collect all weights and gradients in a single Tensor
15 local weights, grad_weights = net:getParameters()
16
17 -- Loss functions are called "criterions"
18 local crit = nn.CrossEntropyCriterion()  -- Softmax loss
19
20 -- Generate some random input data
21 local x = torch.randn(N, D)
22 local y = torch.Tensor(N):random(C)
23
24 -- Forward pass: Compute scores and loss
25 local scores = net:forward(x)
26 local loss = crit:forward(scores, y)
27
28 -- Backward pass: compute gradients
29 grad_weights:zero()
30 local dscores = crit:backward(scores, y)
31 local dx = net:backward(x, dscores)
32
33 -- Make a gradient step
34 local learning_rate = 1e-3
35 weights:add(-learning_rate, grad_weights)
36 |
```

# nn

nn module lets you easily build and train neural nets

**Forward pass**: compute scores and loss

```lua
1  require 'torch'
2  require 'nn'
3
4
5  -- Batch size, input dim, hidden dim, num classes
6  local N, D, H, C = 100, 1000, 100, 10
7
8  -- Build a one-layer ReLU network
9  local net = nn.Sequential()
10 net:add(nn.Linear(D, H))
11 net:add(nn.ReLU())
12 net:add(nn.Linear(H, C))
13
14 -- Collect all weights and gradients in a single Tensor
15 local weights, grad_weights = net:getParameters()
16
17 -- Loss functions are called "criterions"
18 local crit = nn.CrossEntropyCriterion()  -- Softmax loss
19
20 -- Generate some random input data
21 local x = torch.randn(N, D)
22 local y = torch.Tensor(N):random(C)
23
24 -- Forward pass: Compute scores and loss
25 local scores = net:forward(x)
26 local loss = crit:forward(scores, y)
27
28 -- Backward pass: compute gradients
29 grad_weights:zero()
30 local dscores = crit:backward(scores, y)
31 local dx = net:backward(x, dscores)
32
33 -- Make a gradient step
34 local learning_rate = 1e-3
35 weights:add(-learning_rate, grad_weights)
36 |
```

# Torch

## nn

nn module lets you easily build and train neural nets

**Backward pass**: Compute gradients. Remember to set weight gradients to zero!

```lua
1  require 'torch'
2  require 'nn'
3
4
5  -- Batch size, input dim, hidden dim, num classes
6  local N, D, H, C = 100, 1000, 100, 10
7
8  -- Build a one-layer ReLU network
9  local net = nn.Sequential()
10 net:add(nn.Linear(D, H))
11 net:add(nn.ReLU())
12 net:add(nn.Linear(H, C))
13
14 -- Collect all weights and gradients in a single Tensor
15 local weights, grad_weights = net:getParameters()
16
17 -- Loss functions are called "criterions"
18 local crit = nn.CrossEntropyCriterion()  -- Softmax loss
19
20 -- Generate some random input data
21 local x = torch.randn(N, D)
22 local y = torch.Tensor(N):random(C)
23
24 -- Forward pass: Compute scores and loss
25 local scores = net:forward(x)
26 local loss = crit:forward(scores, y)
27
28 -- Backward pass: compute gradients
29 grad_weights:zero()
30 local dscores = crit:backward(scores, y)
31 local dx = net:backward(x, dscores)
32
33 -- Make a gradient step
34 local learning_rate = 1e-3
35 weights:add(-learning_rate, grad_weights)
36 |
```

# Torch

## nn

nn module lets you easily build and train neural nets

**Update**: Make a gradient descent step

```
1  require 'torch'
2  require 'nn'
3
4
5  -- Batch size, input dim, hidden dim, num classes
6  local N, D, H, C = 100, 1000, 100, 10
7
8  -- Build a one-layer ReLU network
9  local net = nn.Sequential()
10 net:add(nn.Linear(D, H))
11 net:add(nn.ReLU())
12 net:add(nn.Linear(H, C))
13
14 -- Collect all weights and gradients in a single Tensor
15 local weights, grad_weights = net:getParameters()
16
17 -- Loss functions are called "criterions"
18 local crit = nn.CrossEntropyCriterion()  -- Softmax loss
19
20 -- Generate some random input data
21 local x = torch.randn(N, D)
22 local y = torch.Tensor(N):random(C)
23
24 -- Forward pass: Compute scores and loss
25 local scores = net:forward(x)
26 local loss = crit:forward(scores, y)
27
28 -- Backward pass: compute gradients
29 grad_weights:zero()
30 local dscores = crit:backward(scores, y)
31 local dx = net:backward(x, dscores)
32
33 -- Make a gradient step
34 local learning_rate = 1e-3
35 weights:add(-learning_rate, grad_weights)
36
```

# Torch

## cunn

Running on GPU is easy:

```lua
require 'torch'
require 'cutorch'
require 'nn'
require 'cunn'

-- Batch size, input dim, hidden dim, num classes
local N, D, H, C = 100, 1000, 100, 10

local dtype = 'torch.CudaTensor'

-- Build a one-layer ReLU network
local net = nn.Sequential()
net:add(nn.Linear(D, H))
net:add(nn.ReLU())
net:add(nn.Linear(H, C))
net:type(dtype)

-- Collect all weights and gradients in a single Tensor
local weights, grad_weights = net:getParameters()

-- Loss functions are called "criterions"
local crit = nn.CrossEntropyCriterion()   -- Softmax loss
crit:type(dtype)

-- Generate some random input data
local x = torch.randn(N, D):type(dtype)
local y = torch.Tensor(N):random(C):type(dtype)

-- Forward pass: Compute scores and loss
local scores = net:forward(x)
local loss = crit:forward(scores, y)

-- Backward pass: compute gradients
grad_weights:zero()
local dscores = crit:backward(scores, y)
local dx = net:backward(x, dscores)

-- Make a gradient step
local learning_rate = 1e-3
weights:add(-learning_rate, grad_weights)
```

# cunn

Running on GPU is easy:

Import a few new packages

```
 1 require 'torch'
 2 require 'cutorch'
 3 require 'nn'
 4 require 'cunn'
 5
 6 -- Batch size, input dim, hidden dim, num classes
 7 local N, D, H, C = 100, 1000, 100, 10
 8
 9 local dtype = 'torch.CudaTensor'
10
11 -- Build a one-layer ReLU network
12 local net = nn.Sequential()
13 net:add(nn.Linear(D, H))
14 net:add(nn.ReLU())
15 net:add(nn.Linear(H, C))
16 net:type(dtype)
17
18 -- Collect all weights and gradients in a single Tensor
19 local weights, grad_weights = net:getParameters()
20
21 -- Loss functions are called "criterions"
22 local crit = nn.CrossEntropyCriterion()   -- Softmax loss
23 crit:type(dtype)
24
25 -- Generate some random input data
26 local x = torch.randn(N, D):type(dtype)
27 local y = torch.Tensor(N):random(C):type(dtype)
28
29 -- Forward pass: Compute scores and loss
30 local scores = net:forward(x)
31 local loss = crit:forward(scores, y)
32
33 -- Backward pass: compute gradients
34 grad_weights:zero()
35 local dscores = crit:backward(scores, y)
36 local dx = net:backward(x, dscores)
37
38 -- Make a gradient step
39 local learning_rate = 1e-3
40 weights:add(-learning_rate, grad_weights)
```

# Torch

## cunn

Running on GPU is easy:

Import a few new packages

Cast network and criterion

```
1  require 'torch'
2  require 'cutorch'
3  require 'nn'
4  require 'cunn'
5
6  -- Batch size, input dim, hidden dim, num classes
7  local N, D, H, C = 100, 1000, 100, 10
8
9  local dtype = 'torch.CudaTensor'
10
11 -- Build a one-layer ReLU network
12 local net = nn.Sequential()
13 net:add(nn.Linear(D, H))
14 net:add(nn.ReLU())
15 net:add(nn.Linear(H, C))
   net:type(dtype)
17
18 -- Collect all weights and gradients in a single Tensor
19 local weights, grad_weights = net:getParameters()
20
21 -- Loss functions are called "criterions"
22 local crit = nn.CrossEntropyCriterion()   -- Softmax loss
   crit:type(dtype)
24
25 -- Generate some random input data
26 local x = torch.randn(N, D):type(dtype)
27 local y = torch.Tensor(N):random(C):type(dtype)
28
29 -- Forward pass: Compute scores and loss
30 local scores = net:forward(x)
31 local loss = crit:forward(scores, y)
32
33 -- Backward pass: compute gradients
34 grad_weights:zero()
35 local dscores = crit:backward(scores, y)
36 local dx = net:backward(x, dscores)
37
38 -- Make a gradient step
39 local learning_rate = 1e-3
40 weights:add(-learning_rate, grad_weights)
```

# Torch

## cunn

Running on GPU is easy:

Import a few new packages

Cast network and criterion

Cast data and labels

```
1  require 'torch'
2  require 'cutorch'
3  require 'nn'
4  require 'cunn'
5
6  -- Batch size, input dim, hidden dim, num classes
7  local N, D, H, C = 100, 1000, 100, 10
8
9  local dtype = 'torch.CudaTensor'
10
11 -- Build a one-layer ReLU network
12 local net = nn.Sequential()
13 net:add(nn.Linear(D, H))
14 net:add(nn.ReLU())
15 net:add(nn.Linear(H, C))
16 net:type(dtype)
17
18 -- Collect all weights and gradients in a single Tensor
19 local weights, grad_weights = net:getParameters()
20
21 -- Loss functions are called "criterions"
22 local crit = nn.CrossEntropyCriterion()  -- Softmax loss
23 crit:type(dtype)
24
25 -- Generate some random input data
26 local x = torch.rand(N, D):type(dtype)
27 local y = torch.Tensor(N):random(C):type(dtype)
28
29 -- Forward pass: Compute scores and loss
30 local scores = net:forward(x)
31 local loss = crit:forward(scores, y)
32
33 -- Backward pass: compute gradients
34 grad_weights:zero()
35 local dscores = crit:backward(scores, y)
36 local dx = net:backward(x, dscores)
37
38 -- Make a gradient step
39 local learning_rate = 1e-3
40 weights:add(-learning_rate, grad_weights)
```

# Torch

## optim

optim package implements different update rules: momentum, Adam, etc

```lua
1  require 'torch'
2  require 'nn'
3  require 'optim'
4
5  -- Batch size, input dim, hidden dim, num classes
6  local N, D, H, C = 100, 1000, 100, 10
7
8  -- Build a one-layer ReLU network
9  local net = nn.Sequential()
10 net:add(nn.Linear(D, H))
11 net:add(nn.ReLU())
12 net:add(nn.Linear(H, C))
13
14 -- Collect all weights and gradients in a single Tensor
15 local weights, grad_weights = net:getParameters()
16
17 -- Loss functions are called "criterions"
18 local crit = nn.CrossEntropyCriterion()  -- Softmax loss
19
20 -- Calback to interface with optim methods
21 local function f(w)
22   assert(w == weights)
23
24   -- Generate some random input data
25   local x = torch.randn(N, D)
26   local y = torch.Tensor(N):random(C)
27
28   -- Forward pass: Compute scores and loss
29   local scores = net:forward(x)
30   local loss = crit:forward(scores, y)
31
32   -- Backward pass: compute gradients
33   grad_weights:zero()
34   local dscores = crit:backward(scores, y)
35   local dx = net:backward(x, dscores)
36
37   return loss, grad_weights
38 end
39
40 -- Make a step using Adam
41 local state = {learningRate=1e-3}
42 optim.adam(f, weights, state)
```

# Torch

## optim

optim package implements different update rules: momentum, Adam, etc

Import optim package

```
1  require 'torch'
2  require 'nn'
3  require 'optim'
4
5  -- Batch size, input dim, hidden dim, num classes
6  local N, D, H, C = 100, 1000, 100, 10
7
8  -- Build a one-layer ReLU network
9  local net = nn.Sequential()
10 net:add(nn.Linear(D, H))
11 net:add(nn.ReLU())
12 net:add(nn.Linear(H, C))
13
14 -- Collect all weights and gradients in a single Tensor
15 local weights, grad_weights = net:getParameters()
16
17 -- Loss functions are called "criterions"
18 local crit = nn.CrossEntropyCriterion()  -- Softmax loss
19
20 -- Calback to interface with optim methods
21 local function f(w)
22   assert(w == weights)
23
24   -- Generate some random input data
25   local x = torch.randn(N, D)
26   local y = torch.Tensor(N):random(C)
27
28   -- Forward pass: Compute scores and loss
29   local scores = net:forward(x)
30   local loss = crit:forward(scores, y)
31
32   -- Backward pass: compute gradients
33   grad_weights:zero()
34   local dscores = crit:backward(scores, y)
35   local dx = net:backward(x, dscores)
36
37   return loss, grad_weights
38 end
39
40 -- Make a step using Adam
41 local state = {learningRate=1e-3}
42 optim.adam(f, weights, state)
```

# optim

optim package implements different update rules: momentum, Adam, etc

Import optim package

Write a callback function that returns loss and gradients

```lua
 1 require 'torch'
 2 require 'nn'
 3 require 'optim'
 4
 5 -- Batch size, input dim, hidden dim, num classes
 6 local N, D, H, C = 100, 1000, 100, 10
 7
 8 -- Build a one-layer ReLU network
 9 local net = nn.Sequential()
10 net:add(nn.Linear(D, H))
11 net:add(nn.ReLU())
12 net:add(nn.Linear(H, C))
13
14 -- Collect all weights and gradients in a single Tensor
15 local weights, grad_weights = net:getParameters()
16
17 -- Loss functions are called "criterions"
18 local crit = nn.CrossEntropyCriterion()  -- Softmax loss
19
20 -- Calback to interface with optim methods
21 local function f(w)
22   assert(w == weights)
23
24   -- Generate some random input data
25   local x = torch.randn(N, D)
26   local y = torch.Tensor(N):random(C)
27
28   -- Forward pass: Compute scores and loss
29   local scores = net:forward(x)
30   local loss = crit:forward(scores, y)
31
32   -- Backward pass: compute gradients
33   grad_weights:zero()
34   local dscores = crit:backward(scores, y)
35   local dx = net:backward(x, dscores)
36
37   return loss, grad_weights
38 end
39
40 -- Make a step using Adam
41 local state = {learningRate=1e-3}
42 optim.adam(f, weights, state)
```

# Torch

## optim

optim package implements different update rules: momentum, Adam, etc

Import optim package

Write a callback function that returns loss and gradients

state variable holds hyperparameters, cached values, etc; pass it to adam

```lua
1  require 'torch'
2  require 'nn'
3  require 'optim'
4
5  -- Batch size, input dim, hidden dim, num classes
6  local N, D, H, C = 100, 1000, 100, 10
7
8  -- Build a one-layer ReLU network
9  local net = nn.Sequential()
10 net:add(nn.Linear(D, H))
11 net:add(nn.ReLU())
12 net:add(nn.Linear(H, C))
13
14 -- Collect all weights and gradients in a single Tensor
15 local weights, grad_weights = net:getParameters()
16
17 -- Loss functions are called "criterions"
18 local crit = nn.CrossEntropyCriterion()  -- Softmax loss
19
20 -- Calback to interface with optim methods
21 local function f(w)
22   assert(w == weights)
23
24   -- Generate some random input data
25   local x = torch.randn(N, D)
26   local y = torch.Tensor(N):random(C)
27
28   -- Forward pass: Compute scores and loss
29   local scores = net:forward(x)
30   local loss = crit:forward(scores, y)
31
32   -- Backward pass: compute gradients
33   grad_weights:zero()
34   local dscores = crit:backward(scores, y)
35   local dx = net:backward(x, dscores)
36
37   return loss, grad_weights
38 end
39
4  -- Make a step using Adam
4  local state = {learningRate=1e-3}
4  optim.adam(f, weights, state)
```

# Torch

## Modules

Caffe has Nets and Layers;
Torch just has Modules

# Torch

## Modules

Caffe has Nets and Layers;
Torch just has Modules

Modules are classes written in
Lua; easy to read and write

Forward / backward written in Lua
using Tensor methods

Same code runs on CPU / GPU

```lua
1   local Linear, parent = torch.class('nn.Linear', 'nn.Module')
2
3   function Linear:__init(inputSize, outputSize, bias)
4       parent.__init(self)
5       local bias = ((bias == nil) and true) or bias
6       self.weight = torch.Tensor(outputSize, inputSize)
7       self.gradWeight = torch.Tensor(outputSize, inputSize)
8       if bias then
9           self.bias = torch.Tensor(outputSize)
10          self.gradBias = torch.Tensor(outputSize)
11      end
12      self:reset()
13  end
14
```

https://github.com/torch/nn/blob/master/Linear.lua

# Torch

## Modules

Caffe has Nets and Layers; Torch just has Modules

Modules are classes written in Lua; easy to read and write

**updateOutput**: Forward pass; compute output

```lua
local Linear, parent = torch.class('nn.Linear', 'nn.Module')

function Linear:updateOutput(input)
   if input:dim() == 1 then
      self.output:resize(self.weight:size(1))
      if self.bias then self.output:copy(self.bias) else self.output:zero() end
      self.output:addmv(1, self.weight, input)
   elseif input:dim() == 2 then
      local nframe = input:size(1)
      local nElement = self.output:nElement()
      self.output:resize(nframe, self.weight:size(1))
      if self.output:nElement() ~= nElement then
         self.output:zero()
      end
      self.addBuffer = self.addBuffer or input.new()
      if self.addBuffer:nElement() ~= nframe then
         self.addBuffer:resize(nframe):fill(1)
      end
      self.output:addmm(0, self.output, 1, input, self.weight:t())
      if self.bias then self.output:addr(1, self.addBuffer, self.bias) end
   else
      error('input must be vector or matrix')
   end

   return self.output
end
```

https://github.com/torch/nn/blob/master/Linear.lua

# Torch

## Modules

Caffe has Nets and Layers; Torch just has Modules

Modules are classes written in Lua; easy to read and write

**updateGradInput:** Backward; compute gradient of input

```lua
local Linear, parent = torch.class('nn.Linear', 'nn.Module')
```

```lua
function Linear:updateGradInput(input, gradOutput)
   if self.gradInput then

      local nElement = self.gradInput:nElement()
      self.gradInput:resizeAs(input)
      if self.gradInput:nElement() ~= nElement then
         self.gradInput:zero()
      end
      if input:dim() == 1 then
         self.gradInput:addmv(0, 1, self.weight:t(), gradOutput)
      elseif input:dim() == 2 then
         self.gradInput:addmm(0, 1, gradOutput, self.weight)
      end

      return self.gradInput
   end
end
```

https://github.com/torch/nn/blob/master/Linear.lua

# Torch

## Modules

Caffe has Nets and Layers;
Torch just has Modules

Modules are classes written in
Lua; easy to read and write

**accGradParameters:** Backward;
compute gradient of weights

```lua
1  local Linear, parent = torch.class('nn.Linear', 'nn.Module')

82  function Linear:accGradParameters(input, gradOutput, scale)
83      scale = scale or 1
84      if input:dim() == 1 then
85          self.gradWeight:addr(scale, gradOutput, input)
86          if self.bias then self.gradBias:add(scale, gradOutput) end
87      elseif input:dim() == 2 then
88          self.gradWeight:addmm(scale, gradOutput:t(), input)
89          if self.bias then
90              self.gradBias:addmv(scale, gradOutput:t(), self.addBuffer)
91          end
92      end
93  end
```

https://github.com/torch/nn/blob/master/Linear.lua

# Torch

# Modules

Tons of built-in modules and loss functions

| | |
|---|---|
| Abs.lua | |
| AbsCriterion.lua | |
| Add.lua | |
| AddConstant.lua | |
| BCECriterion.lua | |
| BatchNormalization.lua | |
| Bilinear.lua | |
| CAddTable.lua | |
| CDivTable.lua | |
| CMakeLists.txt | |
| CMul.lua | |
| CMulTable.lua | |

- TemporalConvolution.lua
- TemporalMaxPooling.lua
- TemporalSubSampling.lua
- Threshold.lua
- Transpose.lua
- View.lua
- VolumetricAveragePooling.lua
- VolumetricConvolution.lua
- VolumetricDropout.lua
- VolumetricFullConvolution.lua
- VolumetricMaxPooling.lua
- VolumetricMaxUnpooling.lua
- WeightedEuclidean.lua
- WeightedMSECriterion.lua

- MarginCriterion.lua
- MarginRankingCriterion.lua
- Max.lua
- Mean.lua
- Min.lua
- MixtureTable.lua
- Module.lua
- Mul.lua
- MulConstant.lua
- MultiCriterion.lua
- MultiLabelMarginCriterion.lua
- MultiLabelSoftMarginCriterion.lua
- MultiMarginCriterion.lua
- Narrow.lua

- SparseLinear.lua
- SpatialAdaptiveMaxPooling.lua
- SpatialAveragePooling.lua
- SpatialBatchNormalization.lua
- SpatialContrastiveNormalization.lua
- SpatialConvolution.lua
- SpatialConvolutionLocal.lua
- SpatialConvolutionMM.lua
- SpatialConvolutionMap.lua
- SpatialCrossMapLRN.lua
- SpatialDivisiveNormalization.lua
- SpatialDropout.lua
- SpatialFractionalMaxPooling.lua
- SpatialFullConvolution.lua
- SpatialFullConvolutionMap.lua
- SpatialLPPooling.lua
- SpatialMaxPooling.lua
- SpatialMaxUnpooling.lua

- ClassSimplexCriterion.lua
- Concat.lua
- ConcatTable.lua
- Container.lua
- Contiguous.lua
- Copy.lua
- Cosine.lua
- CosineDistance.lua
- CosineEmbeddingCriterion.lua
- Criterion.lua
- CriterionTable.lua
- CrossEntropyCriterion.lua
- DepthConcat.lua
- DistKLDivCriterion.lua
- DotProduct.lua
- Dropout.lua
- ELU.lua

https://github.com/torch/nn

# Torch

## Modules

Writing your own modules is easy!

```lua
TimesTwo.lua
1  require 'nn'
2
3  local times_two, parent = torch.class('nn.TimesTwo', 'nn.Module')
4
5
6  function times_two:__init()
7    parent.__init(self)
8  end
9
10
11 function times_two:updateOutput(input)
12   self.output:mul(input, 2)
13   return self.output
14 end
15
16
17 function times_two:updateGradInput(input, gradOutput)
18   self.gradInput:mul(gradOutput, 2)
19   return self.gradInput
20 end
```

```lua
times_two_example.lua
1  require 'nn'
2
3  require 'TimesTwo'
4
5  local times_two = nn.TimesTwo()
6
7  local input = torch.randn(4, 5)
8  local output = times_two:forward(input)
9
10 print('here is input:')
11 print(input)
12
13 print('here is output:')
14 print(output)
15
16 local gradOutput = torch.randn(4, 5)
17 local gradInput = times_two:backward(input, gradOutput)
18
19 print('here is gradOutput:')
20 print(gradOutput)
21
22 print('here is gradInput')
23 print(gradInput)
```

# Torch

## Modules

*Container* modules allow you to combine multiple modules

# Modules

*Container* modules allow you to combine multiple modules

```lua
local seq = nn.Sequential()
seq:add(mod1)
seq:add(mod2)
local out = seq:forward(x)
```

# Modules

*Container* modules allow you to combine multiple modules

```
local seq = nn.Sequential()
seq:add(mod1)
seq:add(mod2)
local out = seq:forward(x)
```

```
local concat = nn.ConcatTable()
concat:add(mod1)
concat:add(mod2)
local out = concat:forward(x)
```

# Modules

*Container* modules allow you to combine multiple modules

```lua
local seq = nn.Sequential()
seq:add(mod1)
seq:add(mod2)
local out = seq:forward(x)
```

```lua
local concat = nn.ConcatTable()
concat:add(mod1)
concat:add(mod2)
local out = concat:forward(x)
```

```lua
local parallel = nn.ParallelTable()
parallel:add(mod1)
parallel:add(mod2)
local out = parallel:forward({x1, x2})
```

# Torch

# nngraph

Use nngraph to build modules
that combine their inputs in
complex ways

**Inputs**: x, y, z
**Outputs**: c
a = x + y
b = a $\odot$ z
c = a + b

# Torch

## nngraph

Use nngraph to build modules that combine their inputs in complex ways

> **Inputs**: x, y, z
> **Outputs**: c
> a = x + y
> b = a $\odot$ z
> c = a + b

# Torch

## nngraph

Use nngraph to build modules that combine their inputs in complex ways

**Inputs**: x, y, z
**Outputs**: c
a = x + y
b = a ⊙ z
c = a + b



```lua
1  require 'torch'
2  require 'nn'
3  require 'nngraph'
4
5  local function build_module()
6    local x = nn.Identity()()
7    local y = nn.Identity()()
8    local z = nn.Identity()()
9
10   local a = nn.CAddTable()({x, y})
11   local b = nn.CMulTable()({a, z})
12   local c = nn.CAddTable()({a, b})
13
14   local inputs = {x, y, z}
15   local outputs = {c}
16   return nn.gModule(inputs, outputs)
17 end
18
19 local mod = build_module()
20
21 local x = torch.randn(4, 5)
22 local y = torch.randn(4, 5)
23 local z = torch.randn(4, 5)
24
25 local c = mod:forward({x, y, z})
```

# Torch

# Pretrained Models

**loadcaffe**: Load pretrained Caffe models: AlexNet, VGG, some others
https://github.com/szagoruyko/loadcaffe

**GoogLeNet v1**: https://github.com/soumith/inception.torch

**GoogLeNet v3**: https://github.com/Moodstocks/inception-v3.torch

**ResNet**: https://github.com/facebook/fb.resnet.torch

# Torch

# Package Management

After installing torch, use luarocks
to install or update Lua packages

(Similar to pip install from Python)

```
luarocks install torch
luarocks install nn
luarocks install optim
luarocks install lua-cjson
```

# Torch: Other useful packages

**torch.cudnn**: Bindings for NVIDIA cuDNN kernels
https://github.com/soumith/cudnn.torch

**torch-hdf5**: Read and write HDF5 files from Torch
https://github.com/deepmind/torch-hdf5

**lua-cjson**: Read and write JSON files from Lua
https://luarocks.org/modules/luarocks/lua-cjson

**cltorch, clnn**: OpenCL backend for Torch, and port of nn
https://github.com/hughperkins/cltorch, https://github.com/hughperkins/clnn

**torch-autograd**: Automatic differentiation; sort of like more powerful nngraph, similar to Theano or TensorFlow
https://github.com/twitter/torch-autograd

**fbcunn**: Facebook: FFT conv, multi-GPU (DataParallel, ModelParallel)
https://github.com/facebook/fbcunn

# Torch

## Pros / Cons

(-) Lua
(-) Less plug-and-play than Caffe
    You usually write your own training code
(+) Lots of modular pieces that are easy to combine
(+) Easy to write your own layer types and run on GPU
(+) Most of the library code is in Lua, easy to read
(+) Lots of pretrained models!
(-) Not great for RNNs

# Theano

http://deeplearning.net/software/theano/

# Theano

From Yoshua Bengio's group at University of Montreal

Embracing computation graphs, symbolic computation

High-level wrappers: Keras, Lasagne

# Computational Graphs

# Theano

# Computational Graphs



```python
import theano
import theano.tensor as T

# Define symbolic variables
x = T.matrix('x')
y = T.matrix('y')
z = T.matrix('z')

# Compute some other values symbolically
a = x + y
b = a * z
c = a + b

# Compile a function that computes c
f = theano.function(
    inputs=[x, y, z],
    outputs=c
    )

# Evaluate the compiled function
# on some real values
xx = np.random.randn(4, 5)
yy = np.random.randn(4, 5)
zz = np.random.randn(4, 5)
print f(xx, yy, zz)

# Repeat the same computation
# explicitly using numpy ops
aa = xx + yy
bb = aa * zz
cc = aa + bb
print cc
```

# Theano

## Computational Graphs



```python
import theano
import theano.tensor as T

# Define symbolic variables
x = T.matrix('x')
y = T.matrix('y')
z = T.matrix('z')

# Compute some other values symbolically
a = x + y
b = a * z
c = a + b

# Compile a function that computes c
f = theano.function(
    inputs=[x, y, z],
    outputs=c
    )

# Evaluate the compiled function
# on some real values
xx = np.random.randn(4, 5)
yy = np.random.randn(4, 5)
zz = np.random.randn(4, 5)
print f(xx, yy, zz)

# Repeat the same computation
# explicitly using numpy ops
aa = xx + yy
bb = aa * zz
cc = aa + bb
print cc
```

Define symbolic variables; these are inputs to the graph

# Theano

# Computational Graphs



```python
import theano
import theano.tensor as T

# Define symbolic variables
x = T.matrix('x')
y = T.matrix('y')
z = T.matrix('z')

# Compute some other values symbolically
a = x + y
b = a * z
c = a + b

# Compile a function that computes c
f = theano.function(
    inputs=[x, y, z],
    outputs=c
    )

# Evaluate the compiled function
# on some real values
xx = np.random.randn(4, 5)
yy = np.random.randn(4, 5)
zz = np.random.randn(4, 5)
print f(xx, yy, zz)

# Repeat the same computation
# explicitly using numpy ops
aa = xx + yy
bb = aa * zz
cc = aa + bb
print cc
```

Compute intermediates
and outputs symbolically

# Theano

## Computational Graphs



```python
import theano
import theano.tensor as T

# Define symbolic variables
x = T.matrix('x')
y = T.matrix('y')
z = T.matrix('z')

# Compute some other values symbolically
a = x + y
b = a * z
c = a + b

# Compile a function that computes c
f = theano.function(
    inputs=[x, y, z],
    outputs=c
    )

# Evaluate the compiled function
# on some real values
xx = np.random.randn(4, 5)
yy = np.random.randn(4, 5)
zz = np.random.randn(4, 5)
print f(xx, yy, zz)

# Repeat the same computation
# explicitly using numpy ops
aa = xx + yy
bb = aa * zz
cc = aa + bb
print cc
```

Compile a function that produces c from x, y, z (generates code)

# Theano

## Computational Graphs



```python
import theano
import theano.tensor as T

# Define symbolic variables
x = T.matrix('x')
y = T.matrix('y')
z = T.matrix('z')

# Compute some other values symbolically
a = x + y
b = a * z
c = a + b

# Compile a function that computes c
f = theano.function(
    inputs=[x, y, z],
    outputs=c
    )

# Evaluate the compiled function
# on some real values
xx = np.random.randn(4, 5)
yy = np.random.randn(4, 5)
zz = np.random.randn(4, 5)
print f(xx, yy, zz)

# Repeat the same computation
# explicitly using numpy ops
aa = xx + yy
bb = aa * zz
cc = aa + bb
print cc
```

Run the function, passing some numpy arrays (may run on GPU)
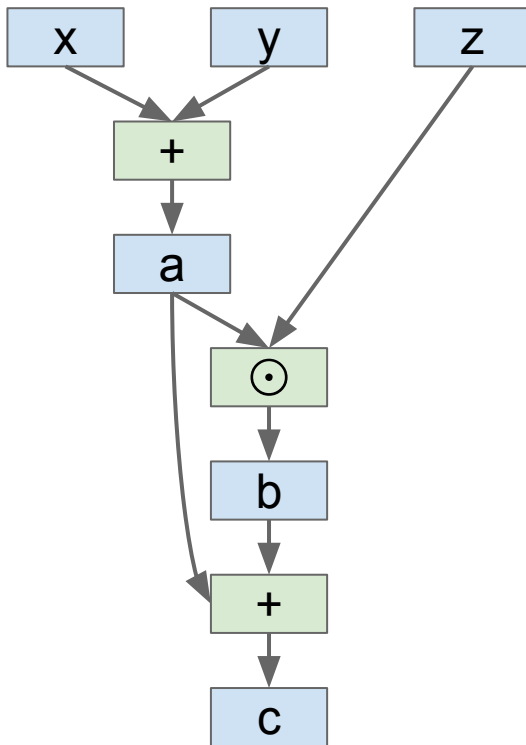
# Theano

# Computational Graphs



```python
import theano
import theano.tensor as T

# Define symbolic variables
x = T.matrix('x')
y = T.matrix('y')
z = T.matrix('z')

# Compute some other values symbolically
a = x + y
b = a * z
c = a + b

# Compile a function that computes c
f = theano.function(
        inputs=[x, y, z],
        outputs=c
    )

# Evaluate the compiled function
# on some real values
xx = np.random.randn(4, 5)
yy = np.random.randn(4, 5)
zz = np.random.randn(4, 5)
print f(xx, yy, zz)

# Repeat the same computation
# explicitly using numpy ops
aa = xx + yy
bb = aa * zz
cc = aa + bb
print cc
```
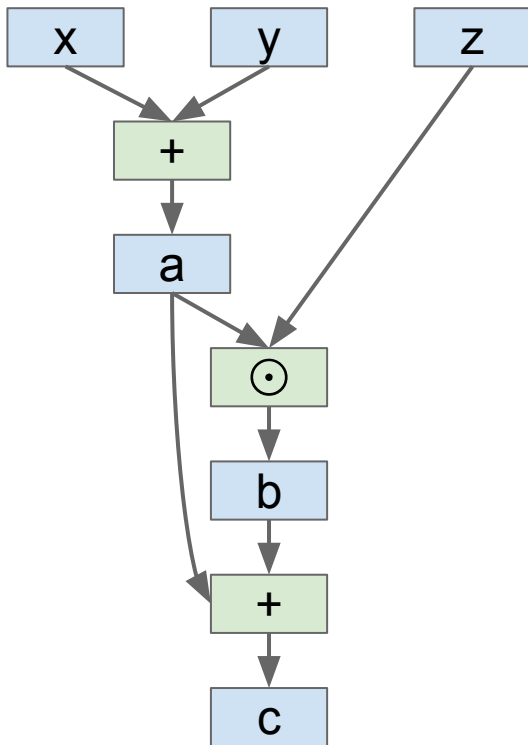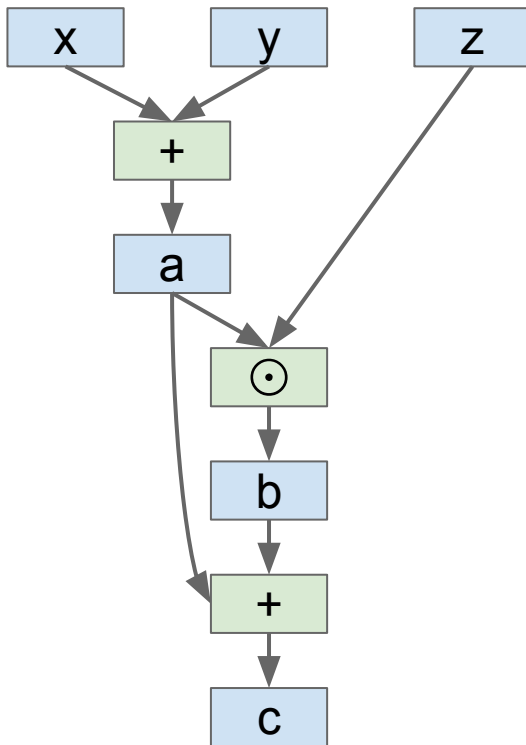
Repeat the same computation using numpy operations (runs on CPU)

# Theano

## Simple Neural Net

```python
import theano
import theano.tensor as T

# Batch size, input dim, hidden dim, num classes
N, D, H, C = 64, 1000, 100, 10

x = T.matrix('x')
y = T.vector('y', dtype='int64')
w1 = T.matrix('w1')
w2 = T.matrix('w2')

# Forward pass: Compute scores
a = x.dot(w1)
a_relu = T.nnet.relu(a)
scores = a_relu.dot(w2)

# Forward pass: compute softmax loss
probs = T.nnet.softmax(scores)
loss = T.nnet.categorical_crossentropy(probs, y).mean()

# Compile a function to compute loss, scores
f = theano.function(
    inputs=[x, y, w1, w2],
    outputs=[loss, scores],
)

# Run the function
xx = np.random.randn(N, D)
yy = np.random.randint(C, size=N)
ww1 = 1e-3 * np.random.randn(D, H)
ww2 = 1e-3 * np.random.randn(H, C)

loss, scores = f(xx, yy, ww1, ww2)
print loss
```

# Theano

# Simple Neural Net

Define symbolic variables:

    x = data

    y = labels

    w1 = first-layer weights

    w2 = second-layer weights

```python
import theano
import theano.tensor as T

# Batch size, input dim, hidden dim, num classes
N, D, H, C = 64, 1000, 100, 10

x = T.matrix('x')
y = T.vector('y', dtype='int64')
w1 = T.matrix('w1')
w2 = T.matrix('w2')

# Forward pass: Compute scores
a = x.dot(w1)
a_relu = T.nnet.relu(a)
scores = a_relu.dot(w2)

# Forward pass: compute softmax loss
probs = T.nnet.softmax(scores)
loss = T.nnet.categorical_crossentropy(probs, y).mean()

# Compile a function to compute loss, scores
f = theano.function(
        inputs=[x, y, w1, w2],
        outputs=[loss, scores],
    )

# Run the function
xx = np.random.randn(N, D)
yy = np.random.randint(C, size=N)
ww1 = 1e-3 * np.random.randn(D, H)
ww2 = 1e-3 * np.random.randn(H, C)

loss, scores = f(xx, yy, ww1, ww2)
print loss
```

# Theano

# Simple Neural Net

**Forward**: Compute scores
(symbolically)

```python
import theano
import theano.tensor as T

# Batch size, input dim, hidden dim, num classes
N, D, H, C = 64, 1000, 100, 10

x = T.matrix('x')
y = T.vector('y', dtype='int64')
w1 = T.matrix('w1')
w2 = T.matrix('w2')

# Forward pass: Compute scores
a = x.dot(w1)
a_relu = T.nnet.relu(a)
scores = a_relu.dot(w2)

# Forward pass: compute softmax loss
probs = T.nnet.softmax(scores)
loss = T.nnet.categorical_crossentropy(probs, y).mean()

# Compile a function to compute loss, scores
f = theano.function(
    inputs=[x, y, w1, w2],
    outputs=[loss, scores],
)

# Run the function
xx = np.random.randn(N, D)
yy = np.random.randint(C, size=N)
ww1 = 1e-3 * np.random.randn(D, H)
ww2 = 1e-3 * np.random.randn(H, C)

loss, scores = f(xx, yy, ww1, ww2)
print loss
```

# Theano

# Simple Neural Net

**Forward**: Compute probs, loss
(symbolically)

```python
import theano
import theano.tensor as T

# Batch size, input dim, hidden dim, num classes
N, D, H, C = 64, 1000, 100, 10

x = T.matrix('x')
y = T.vector('y', dtype='int64')
w1 = T.matrix('w1')
w2 = T.matrix('w2')

# Forward pass: Compute scores
a = x.dot(w1)
a_relu = T.nnet.relu(a)
scores = a_relu.dot(w2)

# Forward pass: compute softmax loss
probs = T.nnet.softmax(scores)
loss = T.nnet.categorical_crossentropy(probs, y).mean()

# Compile a function to compute loss, scores
f = theano.function(
    inputs=[x, y, w1, w2],
    outputs=[loss, scores],
    )

# Run the function
xx = np.random.randn(N, D)
yy = np.random.randint(C, size=N)
ww1 = 1e-3 * np.random.randn(D, H)
ww2 = 1e-3 * np.random.randn(H, C)

loss, scores = f(xx, yy, ww1, ww2)
print loss
```

# Theano

# Simple Neural Net

Compile a function that computes loss, scores

```python
import theano
import theano.tensor as T

# Batch size, input dim, hidden dim, num classes
N, D, H, C = 64, 1000, 100, 10

x = T.matrix('x')
y = T.vector('y', dtype='int64')
w1 = T.matrix('w1')
w2 = T.matrix('w2')

# Forward pass: Compute scores
a = x.dot(w1)
a_relu = T.nnet.relu(a)
scores = a_relu.dot(w2)

# Forward pass: compute softmax loss
probs = T.nnet.softmax(scores)
loss = T.nnet.categorical_crossentropy(probs, y).mean()

# Compile a function to compute loss, scores
f = theano.function(
    inputs=[x, y, w1, w2],
    outputs=[loss, scores],
    )

# Run the function
xx = np.random.randn(N, D)
yy = np.random.randint(C, size=N)
ww1 = 1e-3 * np.random.randn(D, H)
ww2 = 1e-3 * np.random.randn(H, C)

loss, scores = f(xx, yy, ww1, ww2)
print loss
```

# Theano

# Simple Neural Net

```python
import theano
import theano.tensor as T

# Batch size, input dim, hidden dim, num classes
N, D, H, C = 64, 1000, 100, 10

x = T.matrix('x')
y = T.vector('y', dtype='int64')
w1 = T.matrix('w1')
w2 = T.matrix('w2')

# Forward pass: Compute scores
a = x.dot(w1)
a_relu = T.nnet.relu(a)
scores = a_relu.dot(w2)

# Forward pass: compute softmax loss
probs = T.nnet.softmax(scores)
loss = T.nnet.categorical_crossentropy(probs, y).mean()

# Compile a function to compute loss, scores
f = theano.function(
    inputs=[x, y, w1, w2],
    outputs=[loss, scores],
)

# Run the function
xx = np.random.randn(N, D)
yy = np.random.randint(C, size=N)
ww1 = 1e-3 * np.random.randn(D, H)
ww2 = 1e-3 * np.random.randn(H, C)

loss, scores = f(xx, yy, ww1, ww2)
print loss
```

Stuff actual numpy arrays into the function

# Theano

# Computing Gradients

```python
import theano
import theano.tensor as T

# Batch size, input dim, hidden dim, num classes
N, D, H, C = 64, 1000, 100, 10

x = T.matrix('x')
y = T.vector('y', dtype='int64')
w1 = T.matrix('w1')
w2 = T.matrix('w2')

# Forward pass: Compute scores
a = x.dot(w1)
a_relu = T.nnet.relu(a)
scores = a_relu.dot(w2)

# Forward pass: compute softmax loss
probs = T.nnet.softmax(scores)
loss = T.nnet.categorical_crossentropy(probs, y).mean()

# Backward pass: compute gradients
dw1, dw2 = T.grad(loss, [w1, w2])

f = theano.function(
    inputs=[x, y, w1, w2],
    outputs=[loss, scores, dw1, dw2],
    )
```

# Computing Gradients

```python
import theano
import theano.tensor as T

# Batch size, input dim, hidden dim, num classes
N, D, H, C = 64, 1000, 100, 10

x = T.matrix('x')
y = T.vector('y', dtype='int64')
w1 = T.matrix('w1')
w2 = T.matrix('w2')

# Forward pass: Compute scores
a = x.dot(w1)
a_relu = T.nnet.relu(a)
scores = a_relu.dot(w2)

# Forward pass: compute softmax loss
probs = T.nnet.softmax(scores)
loss = T.nnet.categorical_crossentropy(probs, y).mean()

# Backward pass: compute gradients
dw1, dw2 = T.grad(loss, [w1, w2])

f = theano.function(
    inputs=[x, y, w1, w2],
    outputs=[loss, scores, dw1, dw2],
    )
```

Same as before: define variables, compute scores and loss symbolically

# Theano

# Computing Gradients

```python
import theano
import theano.tensor as T

# Batch size, input dim, hidden dim, num classes
N, D, H, C = 64, 1000, 100, 10

x = T.matrix('x')
y = T.vector('y', dtype='int64')
w1 = T.matrix('w1')
w2 = T.matrix('w2')

# Forward pass: Compute scores
a = x.dot(w1)
a_relu = T.nnet.relu(a)
scores = a_relu.dot(w2)

# Forward pass: compute softmax loss
probs = T.nnet.softmax(scores)
loss = T.nnet.categorical_crossentropy(probs, y).mean()

# Backward pass: compute gradients
dw1, dw2 = T.grad(loss, [w1, w2])

f = theano.function(
    inputs=[x, y, w1, w2],
    outputs=[loss, scores, dw1, dw2],
    )
```

Theano computes gradients for us symbolically!

# Theano

# Computing Gradients

```python
import theano
import theano.tensor as T

# Batch size, input dim, hidden dim, num classes
N, D, H, C = 64, 1000, 100, 10

x = T.matrix('x')
y = T.vector('y', dtype='int64')
w1 = T.matrix('w1')
w2 = T.matrix('w2')

# Forward pass: Compute scores
a = x.dot(w1)
a_relu = T.nnet.relu(a)
scores = a_relu.dot(w2)

# Forward pass: compute softmax loss
probs = T.nnet.softmax(scores)
loss = T.nnet.categorical_crossentropy(probs, y).mean()

# Backward pass: compute gradients
dw1, dw2 = T.grad(loss, [w1, w2])

f = theano.function(
    inputs=[x, y, w1, w2],
    outputs=[loss, scores, dw1, dw2],
    )
```

Now the function returns loss, scores, and gradients

# Theano

# Computing Gradients

```python
import theano
import theano.tensor as T

# Batch size, input dim, hidden dim, num classes
N, D, H, C = 64, 1000, 100, 10

x = T.matrix('x')
y = T.vector('y', dtype='int64')
w1 = T.matrix('w1')
w2 = T.matrix('w2')

# Forward pass: Compute scores
a = x.dot(w1)
a_relu = T.nnet.relu(a)
scores = a_relu.dot(w2)

# Forward pass: compute softmax loss
probs = T.nnet.softmax(scores)
loss = T.nnet.categorical_crossentropy(probs, y).mean()

# Backward pass: compute gradients
dw1, dw2 = T.grad(loss, [w1, w2])

f = theano.function(
    inputs=[x, y, w1, w2],
    outputs=[loss, scores, dw1, dw2],
)
```
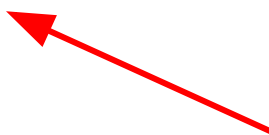
```python
# Run the function
xx = np.random.randn(N, D)
yy = np.random.randint(C, size=N)
ww1 = 1e-2 * np.random.randn(D, H)
ww2 = 1e-2 * np.random.randn(H, C)

learning_rate = 1e-1
for t in xrange(50):
    loss, scores, dww1, dww2 = f(xx, yy, ww1, ww2)
    print loss
    ww1 -= learning_rate * dww1
    ww2 -= learning_rate * dww2
```

Use the function to perform gradient descent!

# Theano

## Pros / Cons

(+) Python + numpy
(+) Computational graph is nice abstraction
(+) RNNs fit nicely in computational graph
(-) Raw Theano is somewhat low-level
(+) High level wrappers (Keras, Lasagne) ease the pain
(-) Error messages can be unhelpful
(-) Large models can have long compile times
(-) Much "fatter" than Torch; more magic
(-) Patchy support for pretrained models

# TensorFlow

https://www.tensorflow.org

# TensorFlow

From Google

Very similar to Theano - all about computation graphs

Easy visualizations (TensorBoard)

Multi-GPU and multi-node training

# TensorFlow

# TensorFlow: Two-Layer Net

```python
1  import tensorflow as tf
2  import numpy as np
3
4  N, D, H, C = 64, 1000, 100, 10
5
6  x = tf.placeholder(tf.float32, shape=[None, D])
7  y = tf.placeholder(tf.float32, shape=[None, C])
8
9  w1 = tf.Variable(1e-3 * np.random.randn(D, H).astype(np.float32))
10 w2 = tf.Variable(1e-3 * np.random.randn(H, C).astype(np.float32))
11
12 a = tf.matmul(x, w1)
13 a_relu = tf.nn.relu(a)
14 scores = tf.matmul(a_relu, w2)
15 probs = tf.nn.softmax(scores)
16 loss = -tf.reduce_sum(y * tf.log(probs))
17
18 learning_rate = 1e-2
19 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(loss)
20
21 xx = np.random.randn(N, D).astype(np.float32)
22 yy = np.zeros((N, C)).astype(np.float32)
23 yy[np.arange(N), np.random.randint(C, size=N)] = 1
24
25 with tf.Session() as sess:
26   sess.run(tf.initialize_all_variables())
27
28   for t in xrange(100):
29     _, loss_value = sess.run([train_step, loss],
30                             feed_dict={x: xx, y: yy})
31     print loss_value
```

# TensorFlow

# TensorFlow: Two-Layer Net

Create placeholders for data and labels: These will be fed to the graph

```python
1  import tensorflow as tf
2  import numpy as np
3
4  N, D, H, C = 64, 1000, 100, 10
5
6  x = tf.placeholder(tf.float32, shape=[None, D])
7  y = tf.placeholder(tf.float32, shape=[None, C])
8
9  w1 = tf.Variable(1e-3 * np.random.randn(D, H).astype(np.float32))
10 w2 = tf.Variable(1e-3 * np.random.randn(H, C).astype(np.float32))
11
12 a = tf.matmul(x, w1)
13 a_relu = tf.nn.relu(a)
14 scores = tf.matmul(a_relu, w2)
15 probs = tf.nn.softmax(scores)
16 loss = -tf.reduce_sum(y * tf.log(probs))
17
18 learning_rate = 1e-2
19 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(loss)
20
21 xx = np.random.randn(N, D).astype(np.float32)
22 yy = np.zeros((N, C)).astype(np.float32)
23 yy[np.arange(N), np.random.randint(C, size=N)] = 1
24
25 with tf.Session() as sess:
26   sess.run(tf.initialize_all_variables())
27
28   for t in xrange(100):
29     _, loss_value = sess.run([train_step, loss],
30                              feed_dict={x: xx, y: yy})
31     print loss_value
```

# TensorFlow: Two-Layer Net

Create Variables to hold weights; similar to Theano shared variables

Initialize variables with numpy arrays

```python
1  import tensorflow as tf
2  import numpy as np
3
4  N, D, H, C = 64, 1000, 100, 10
5
6  x = tf.placeholder(tf.float32, shape=[None, D])
7  y = tf.placeholder(tf.float32, shape=[None, C])
8
9  w1 = tf.Variable(1e-3 * np.random.randn(D, H).astype(np.float32))
10 w2 = tf.Variable(1e-3 * np.random.randn(H, C).astype(np.float32))
11
12 a = tf.matmul(x, w1)
13 a_relu = tf.nn.relu(a)
14 scores = tf.matmul(a_relu, w2)
15 probs = tf.nn.softmax(scores)
16 loss = -tf.reduce_sum(y * tf.log(probs))
17
18 learning_rate = 1e-2
19 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(loss)
20
21 xx = np.random.randn(N, D).astype(np.float32)
22 yy = np.zeros((N, C)).astype(np.float32)
23 yy[np.arange(N), np.random.randint(C, size=N)] = 1
24
25 with tf.Session() as sess:
26   sess.run(tf.initialize_all_variables())
27
28   for t in xrange(100):
29     _, loss_value = sess.run([train_step, loss],
30                       feed_dict={x: xx, y: yy})
31     print loss_value
32
```

# TensorFlow

# TensorFlow: Two-Layer Net

**Forward**: Compute scores, probs, loss (symbolically)

```
1  import tensorflow as tf
2  import numpy as np
3
4  N, D, H, C = 64, 1000, 100, 10
5
6  x = tf.placeholder(tf.float32, shape=[None, D])
7  y = tf.placeholder(tf.float32, shape=[None, C])
8
9  w1 = tf.Variable(1e-3 * np.random.randn(D, H).astype(np.float32))
10 w2 = tf.Variable(1e-3 * np.random.randn(H, C).astype(np.float32))
11
12 a = tf.matmul(x, w1)
13 a_relu = tf.nn.relu(a)
14 scores = tf.matmul(a_relu, w2)
15 probs = tf.nn.softmax(scores)
16 loss = -tf.reduce_sum(y * tf.log(probs))
17
18 learning_rate = 1e-2
19 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(loss)
20
21 xx = np.random.randn(N, D).astype(np.float32)
22 yy = np.zeros((N, C)).astype(np.float32)
23 yy[np.arange(N), np.random.randint(C, size=N)] = 1
24
25 with tf.Session() as sess:
26   sess.run(tf.initialize_all_variables())
27
28   for t in xrange(100):
29     _, loss_value = sess.run([train_step, loss],
30                              feed_dict={x: xx, y: yy})
31     print loss_value
```

# TensorFlow: Two-Layer Net

```python
1  import tensorflow as tf
2  import numpy as np
3
4  N, D, H, C = 64, 1000, 100, 10
5
6  x = tf.placeholder(tf.float32, shape=[None, D])
7  y = tf.placeholder(tf.float32, shape=[None, C])
8
9  w1 = tf.Variable(1e-3 * np.random.randn(D, H).astype(np.float32))
10 w2 = tf.Variable(1e-3 * np.random.randn(H, C).astype(np.float32))
11
12 a = tf.matmul(x, w1)
13 a_relu = tf.nn.relu(a)
14 scores = tf.matmul(a_relu, w2)
15 probs = tf.nn.softmax(scores)
16 loss = -tf.reduce_sum(y * tf.log(probs))
17
18 learning_rate = 1e-2
19 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(loss)
20
21 xx = np.random.randn(N, D).astype(np.float32)
22 yy = np.zeros((N, C)).astype(np.float32)
23 yy[np.arange(N), np.random.randint(C, size=N)] = 1
24
25 with tf.Session() as sess:
26   sess.run(tf.initialize_all_variables())
27
28   for t in xrange(100):
29     _, loss_value = sess.run([train_step, loss],
30                       feed_dict={x: xx, y: yy})
31     print loss_value
```

Running train_step will use SGD to minimize loss

# TensorFlow: Two-Layer Net

```python
1  import tensorflow as tf
2  import numpy as np
3
4  N, D, H, C = 64, 1000, 100, 10
5
6  x = tf.placeholder(tf.float32, shape=[None, D])
7  y = tf.placeholder(tf.float32, shape=[None, C])
8
9  w1 = tf.Variable(1e-3 * np.random.randn(D, H).astype(np.float32))
10 w2 = tf.Variable(1e-3 * np.random.randn(H, C).astype(np.float32))
11
12 a = tf.matmul(x, w1)
13 a_relu = tf.nn.relu(a)
14 scores = tf.matmul(a_relu, w2)
15 probs = tf.nn.softmax(scores)
16 loss = -tf.reduce_sum(y * tf.log(probs))
17
18 learning_rate = 1e-2
19 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(loss)
20
21 xx = np.random.randn(N, D).astype(np.float32)
22 yy = np.zeros((N, C)).astype(np.float32)
23 yy[np.arange(N), np.random.randint(C, size=N)] = 1
24
25 with tf.Session() as sess:
26   sess.run(tf.initialize_all_variables())
27
28   for t in xrange(100):
29     _, loss_value = sess.run([train_step, loss],
30                              feed_dict={x: xx, y: yy})
31     print loss_value
```

Create an artificial dataset; y is one-hot like Keras

# TensorFlow

# TensorFlow: Two-Layer Net

```python
1  import tensorflow as tf
2  import numpy as np
3
4  N, D, H, C = 64, 1000, 100, 10
5
6  x = tf.placeholder(tf.float32, shape=[None, D])
7  y = tf.placeholder(tf.float32, shape=[None, C])
8
9  w1 = tf.Variable(1e-3 * np.random.randn(D, H).astype(np.float32))
10 w2 = tf.Variable(1e-3 * np.random.randn(H, C).astype(np.float32))
11
12 a = tf.matmul(x, w1)
13 a_relu = tf.nn.relu(a)
14 scores = tf.matmul(a_relu, w2)
15 probs = tf.nn.softmax(scores)
16 loss = -tf.reduce_sum(y * tf.log(probs))
17
18 learning_rate = 1e-2
19 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(loss)
20
21 xx = np.random.randn(N, D).astype(np.float32)
22 yy = np.zeros((N, C)).astype(np.float32)
23 yy[np.arange(N), np.random.randint(C, size=N)] = 1
24
25 with tf.Session() as sess:
26   sess.run(tf.initialize_all_variables())
27
28   for t in xrange(100):
29     _, loss_value = sess.run([train_step, loss],
30                              feed_dict={x: xx, y: yy})
31     print loss_value
```
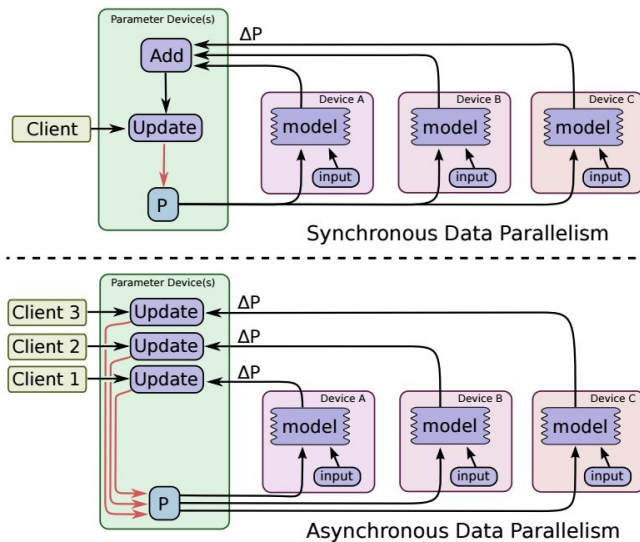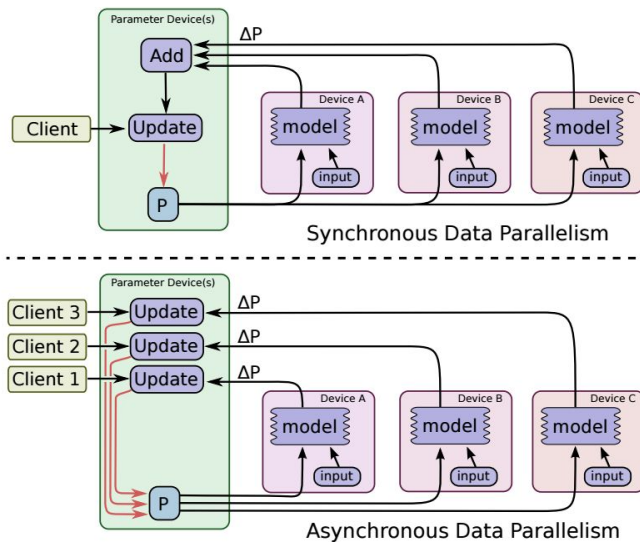
Actually train the model

# TensorFlow: Multi-GPU

**Data parallelism**:
synchronous or asynchronous
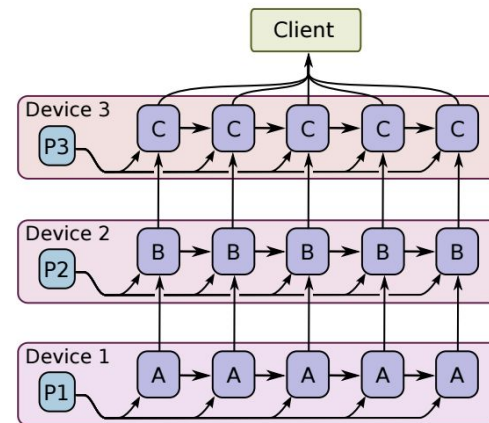
# TensorFlow

# TensorFlow: Multi-GPU

**Data parallelism**:
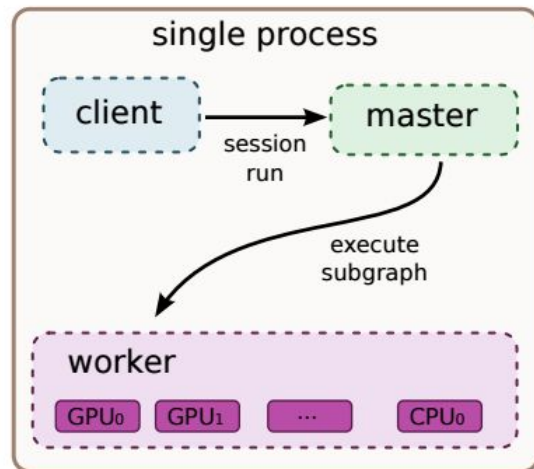synchronous or asynchronous

**Model parallelism**:
Split model across GPUs



Synchronous Data Parallelism
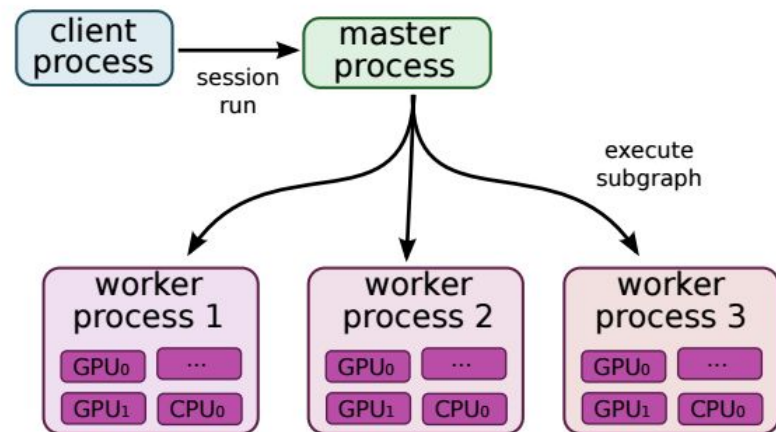
Asynchronous Data Parallelism

# TensorFlow: Distributed

**Single machine**:
Like other frameworks

**Many machines**:
Not open source (yet) =(

# TensorFlow: Pros / Cons

(+) Python + numpy
(+) Computational graph abstraction, like Theano; great for RNNs
(+) Much faster compile times than Theano
(+) Slightly more convenient than raw Theano?
(+) TensorBoard for visualization
(+) Data AND model parallelism; best of all frameworks
(+/-) Distributed models, but not open-source yet
(-) Slower than other frameworks right now
(-) Much "fatter" than Torch; more magic
(-) Not many pretrained models

# Comparison between Libraries

| | **Caffe** | **Torch** | **Theano** | **TensorFlow** |
|---|---|---|---|---|
| **Language** | C++, Python | Lua | Python | Python |
| **Pretrained** | Yes ++ | Yes ++ | Yes (Lasagne) | Inception |
| **Multi-GPU: Data parallel** | Yes | Yes cunn.DataParallelTable | Yes platoon | Yes |
| **Multi-GPU: Model parallel** | No | Yes fbcunn.ModelParallel | Experimental | Yes (best) |
| **Readable source code** | Yes (C++) | Yes (Lua) | No | No |
| **Good at RNN** | No | Mediocre | Yes | Yes (best) |

Any Question???
Thanks