

Deep Learning Libraries

Badri Narayana Patro¹

¹Department of Electrical Engineering
Indian Institute of Technology, Kanpur

August 8, 2016





- I have referred few slides of Prof. Fei-Fei Li and Prof. Anoop M. Namboodiri.
- This presentation is meant for educational purpose only.



Table of Contents

- 1 Introduction
- 2 Caffe
- 3 Torch
- 4 Theano
- 5 Tensor Flow



Table of Contents

1 Introduction

2 Caffe

3 Torch

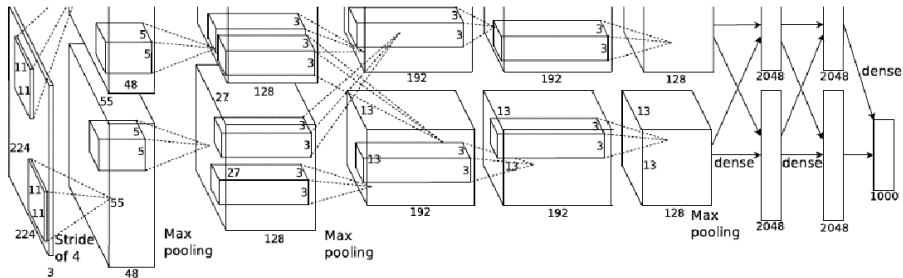
4 Theano

5 Tensor Flow

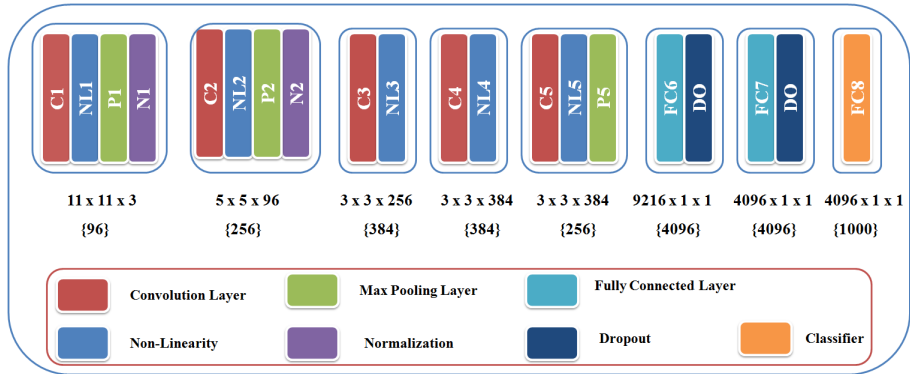


Deep Network: CNN

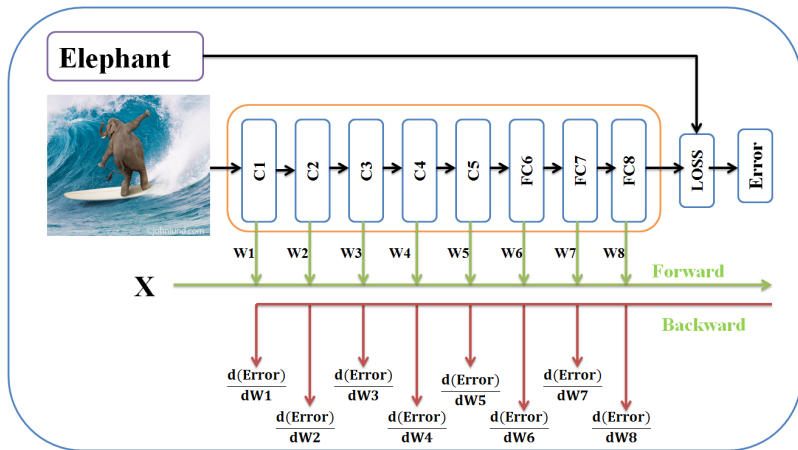
- Similar framework to LeCun'98 but:
 - Bigger model (7 hidden layers, 60,000,000 params)
 - More data (10^6 vs 10^3 images)
 - GPU implementation (50x speedup over CPU)
 - Trained on two GPUs for a week



Deep Network: CNN



Training CNN



How to code for such a network?

- How do I code each type of layer(what is a layer ?)



How to code for such a network?

- How do I code each type of layer(what is a layer ?)
- How do I connect them together



How to code for such a network?

- How do I code each type of layer(what is a layer ?)
- How do I connect them together
- Define a loss function



How to code for such a network?

- How do I code each type of layer(what is a layer ?)
- How do I connect them together
- Define a loss function
- Compute the gradient for each layer



How to code for such a network?

- How do I code each type of layer(what is a layer ?)
- How do I connect them together
- Define a loss function
- Compute the gradient for each layer
- Code for updating the weights of each node/layer



How to code for such a network?

- How do I code each type of layer(what is a layer ?)
- How do I connect them together
- Define a loss function
- Compute the gradient for each layer
- Code for updating the weights of each node/layer
- Take care of normalization



How to code for such a network?

- How do I code each type of layer(what is a layer ?)
- How do I connect them together
- Define a loss function
- Compute the gradient for each layer
- Code for updating the weights of each node/layer
- Take care of normalization
- Take care of learning rates



How to code for such a network?

- How do I code each type of layer(what is a layer ?)
- How do I connect them together
- Define a loss function
- Compute the gradient for each layer
- Code for updating the weights of each node/layer
- Take care of normalization
- Take care of learning rates
- Code the optimizer/solver (SGD?)



How to code for such a network?

- How do I code each type of layer(what is a layer ?)
- How do I connect them together
- Define a loss function
- Compute the gradient for each layer
- Code for updating the weights of each node/layer
- Take care of normalization
- Take care of learning rates
- Code the optimizer/solver (SGD?)
- Several tricks to improve learning (Batch-norm, Dropout, etc.)



How to code for such a network?

- How do I code each type of layer(what is a layer ?)
- How do I connect them together
- Define a loss function
- Compute the gradient for each layer
- Code for updating the weights of each node/layer
- Take care of normalization
- Take care of learning rates
- Code the optimizer/solver (SGD?)
- Several tricks to improve learning (Batch-norm, Dropout, etc.)
- Do this on CPU/GPU/Distributed-systems



Would be nice if we could say

- Add a convolutional layer (params)



Would be nice if we could say

- Add a convolutional layer (params)
- Add a ReLU activation layer



Would be nice if we could say

- Add a convolutional layer (params)
- Add a ReLU activation layer
- Add a max-pooling layer



Would be nice if we could say

- Add a convolutional layer (params)
- Add a ReLU activation layer
- Add a max-pooling layer
- Add a soft-max layer



Would be nice if we could say

- Add a convolutional layer (params)
- Add a ReLU activation layer
- Add a max-pooling layer
- Add a soft-max layer
- Use cross-entropy loss function



Would be nice if we could say

- Add a convolutional layer (params)
- Add a ReLU activation layer
- Add a max-pooling layer
- Add a soft-max layer
- Use cross-entropy loss function
- Train in batches of 64 images



Would be nice if we could say

- Add a convolutional layer (params)
- Add a ReLU activation layer
- Add a max-pooling layer
- Add a soft-max layer
- Use cross-entropy loss function
- Train in batches of 64 images
- Use the following learning rate schedule



- Easily configure a network(Sequential and parallel connection)



- Easily configure a network(Sequential and parallel connection)
- Easily Add and remove layers in a network



- Easily configure a network(Sequential and parallel connection)
- Easily Add and remove layers in a network
- Make use of GPUs, CPUs



- Easily configure a network(Sequential and parallel connection)
- Easily Add and remove layers in a network
- Make use of GPUs, CPUs
- Easily train a network



- Easily configure a network(Sequential and parallel connection)
- Easily Add and remove layers in a network
- Make use of GPUs, CPUs
- Easily train a network
- Change various parameter





Some of Popular Libraries

- Torch [Lua]





Some of Popular Libraries

- Torch [Lua]
- Theano [Python]





Some of Popular Libraries

- Torch [Lua]
- Theano [Python]
- Caffe





Some of Popular Libraries

- Torch [Lua]
- Theano [Python]
- Caffe
- TensorFlow





Some of Popular Libraries

- Torch [Lua]
- Theano [Python]
- Caffe
- TensorFlow
- Matconvnet



Table of Contents

1 Introduction

2 Caffe

3 Torch

4 Theano

5 Tensor Flow



- Created at UC Berkeley



- Created at UC Berkeley
- Written in C++, Python wrappers.



- Created at UC Berkeley
- Written in C++, Python wrappers.
- Creator:
Yangqing Jia(UCB, Google, Facebook)



- Created at UC Berkeley
- Written in C++, Python wrappers.
- Creator:
Yangqing Jia(UCB, Google, Facebook)
- Current Lead Developer: Evan Shelhamer: (UCB)



- Created at UC Berkeley
- Written in C++, Python wrappers.
- Creator:
Yangqing Jia(UCB, Google, Facebook)
- Current Lead Developer: Evan Shelhamer: (UCB)
- In short
 - Simple and fast
 - Very easy for beginners/users



Pros:

- Good for feedforward networks
- Good for finetuning existing networks
- Train models without writing any code!
- Python interface is pretty useful!

Cons:

- Need to write C++ / CUDA for new GPU layers
- Not good for recurrent networks
- Not Good big networks (GoogLeNet, ResNet)



Table of Contents

- 1 Introduction
- 2 Caffe
- 3 Torch**
- 4 Theano
- 5 Tensor Flow





- Developed at NYU and IDIAP
- Written in C and Lua.
- Current Version: Torch7 (1,3,5,7)
- Current Maintainers:
 - Ronan Collobert: Facebook [IDIAP]
 - Clement Farabet: Twitter [NYU]
 - Koray Kavukcuoglu: Google DeepMind
 - Soumith Chintala: Facebook
- In short
 - Flexible and fast





- Main data structure:
 - table in Lua; like object in Javascript
 - tensor in Torch: like Numpy array in Python
- Core Torch just has Tensors; others in packages
 - neural network layers in nn package
 - optimizers in optim package
 - nnggraph for complex architectures
- Torch layers are called modules; operate on Tensor



- Build a two-layer ReLU net



- Build a two-layer ReLU net
- Get weights and gradient for entire network



- Build a two-layer ReLU net
- Get weights and gradient for entire network
- Use a softmax loss function



- Build a two-layer ReLU net
- Get weights and gradient for entire network
- Use a softmax loss function
- Generate random data



- Build a two-layer ReLU net
- Get weights and gradient for entire network
- Use a softmax loss function
- Generate random data
- Forward pass: compute scores and loss





- Build a two-layer ReLU net
- Get weights and gradient for entire network
- Use a softmax loss function
- Generate random data
- Forward pass: compute scores and loss
- Backward pass: Compute gradients. Remember to set weight gradients to zero



- Build a two-layer ReLU net
- Get weights and gradient for entire network
- Use a softmax loss function
- Generate random data
- Forward pass: compute scores and loss
- Backward pass: Compute gradients. Remember to set weight gradients to zero
- Update: Make a gradient descent step



Torch: Training

```
1 --Step 0:----- Load all the necessary packages-----
2 require 'cutorch'
3 require 'cunn'
4 require 'optm'
5 require 'gnuplot'
6 -- load file which have user defined functions
7 dofile './util.lua'
8
9 --step 1:-----Intillize variables-----
10
11 --Batch size, input dimenston, hidden dimenston, number of classes
12 local batch_size, input_data_dim, No_of_neurons_hidden_layer, no_of_classes = 100, 1000, 100, 10
13
14 --step 2:-----Build Network:: two layer ReLU Network-----
15
16 --note: variable cann't start with number, 2l_network== worgn, correct one is l2_network.
17 local l2_networks = nn.Sequential()
18 l2_networks:add(nn.Linear(input_data_dim, No_of_neurons_hidden_layer))
19 l2_networks:add(nn.ReLU())
20 l2_networks:add(nn.Linear(No_of_neurons_hidden_layer, no_of_classes))
21
22 --step 3:-----Get Weight and gradient for entire Network-----
23
24 --Collect all the weights and gradients in a single tensor
25 local weights, gradient_weights = l2_networks:getParameters()
26
27 --step 4:-----Use Softmax Loss function-----
28
29 -- loss function are called criterions
30 local criterion = nn.CrossEntropyCriterion() --Softmax loss
31
32 --step 5:-----Generate Random Data-----
33
34 -- Generate some Random Input Data
35 local input_x = torch.randn(batch_size, input_data_dim)
36 local output_y = torch.Tensor(batch_size):random(no_of_classes)
37
38 --step 6:-----Forward pass -: compute scores and loss-----
39
40 local scores = l2_networks:forward(input_x)
41 local loss = criterion:forward(scores, output_y)
42
43 --step 7:-----Backward pass -: Compute gradients-----
44
45 gradient_weights:zero()-- set gradient weight to zero
46
47 local dscores = criterion:backward(scores, output_y) --ist find grad of loss function then
48 local dx = l2_networks:backward(input_x, dscores) -- find grad of network
49
50 --step 8:-----update weights -: Make a gradient descent step-----
51
52 local learning_rate = 1e-3
53 weights:add(-learning_rate, gradient_weights)
54
```



Pros:

- 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!

Cons:

- Lua
- Less plug-and-play than Caffe
- You usually write your own training code
- Not great for RNNs



Table of Contents

1 Introduction

2 Caffe

3 Torch

4 Theano

5 Tensor Flow



- Developed at Toronto and Montreal



- Developed at Toronto and Montreal
- BSD License



- Developed at Toronto and Montreal
- BSD License
- Written in Python.



- Developed at Toronto and Montreal
- BSD License
- Written in Python.
- Current Maintainers: Yoshua Bengio's group at University of Montreal

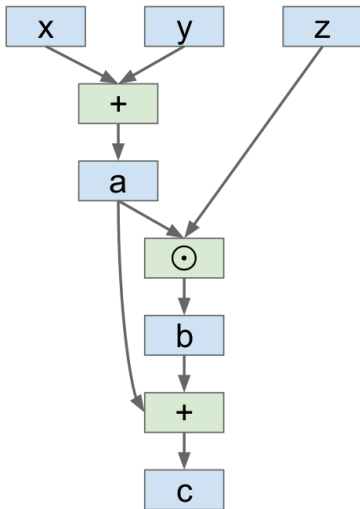




- Developed at Toronto and Montreal
- BSD License
- Written in Python.
- Current Maintainers: Yoshua Bengio's group at University of Montreal
- In short
 - Complex code to do simple things
 - High-level wrappers: Keras, Lasagne



Theano: Computational Graphs



```
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 code

```
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 "criteria"
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 |
```

Image credit: Fei-Fei Li, et al.



Theano: Pros / Cons

Pros:

- Python + numpy
- Computational graph is nice abstraction
- RNNs fit nicely in computational graph
- High level wrappers (Keras, Lasagne) ease the pain

Cons:

- Raw Theano is somewhat low-level
- Error messages can be unhelpful
- Large models can have long compile times
- Much “fatter” than Torch; more magic
- Patchy support for pretrained models



Table of Contents

- 1 Introduction
- 2 Caffe
- 3 Torch
- 4 Theano
- 5 Tensor Flow**



- Developed by Google



- Developed by Google
- Written in C++ and Python



- Developed by Google
- Written in C++ and Python
- BSD License



- Developed by Google
- Written in C++ and Python
- BSD License
- Released in 2015



- Developed by Google
- Written in C++ and Python
- BSD License
- Released in 2015
- In short
 - Good for large, distributed systems
 - Data and model parallelism



- Very similar to Theano - all about computation graphs



- Very similar to Theano - all about computation graphs
- Easy visualisation tool: TensorBoard



- Very similar to Theano - all about computation graphs
- Easy visualisation tool: TensorBoard
- Multi-GPU and multi-node training



- Very similar to Theano - all about computation graphs
- Easy visualisation tool: TensorBoard
- Multi-GPU and multi-node training
- TensorBoard makes it easy to visualise training process, in terms of loss, weights, activations, and even the whole DAG can be seen!

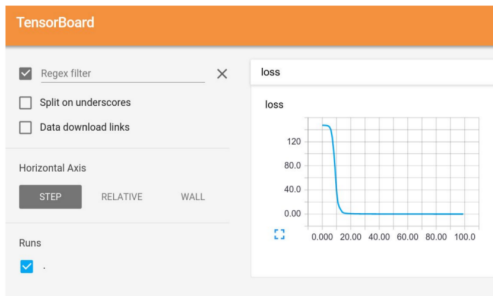


- Very similar to Theano - all about computation graphs
- Easy visualisation tool: TensorBoard
- Multi-GPU and multi-node training
- TensorBoard makes it easy to visualise training process, in terms of loss, weights, activations, and even the whole DAG can be seen!
- Matconvnet



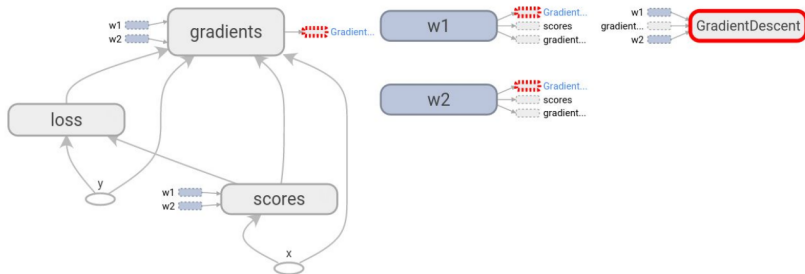
TensorFlow: Tensorboard

Start Tensorboard server, and we get graphs!



TensorFlow: TensorBoard

Tensorboard shows the graph!



TensorFlow: TensorBoard

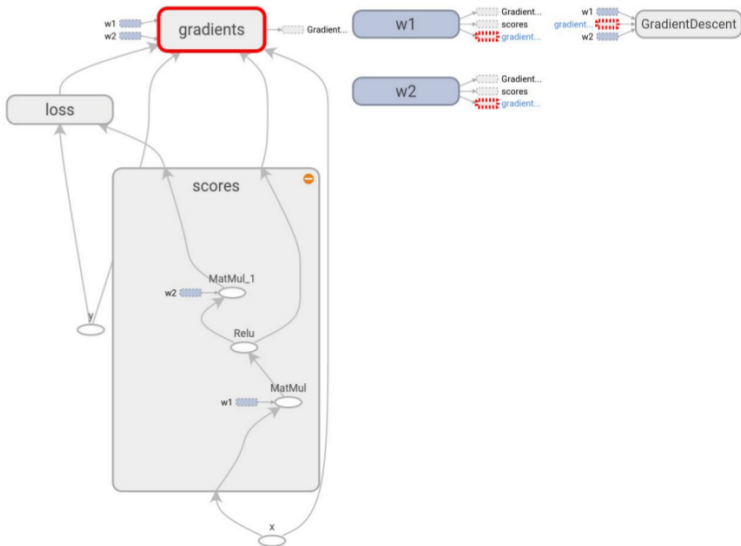
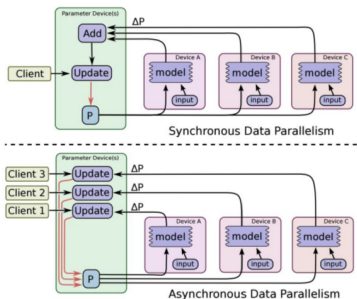


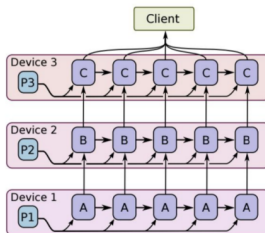
Image courtesy: Fei-Fei Li et al.



Data parallelism:
synchronous or asynchronous

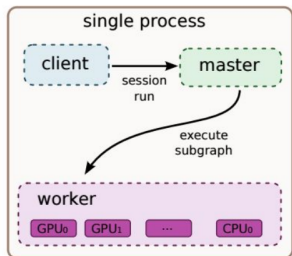


Model parallelism:
Split model across GPUs



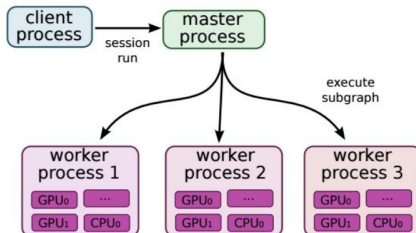
Single machine:

Like other frameworks



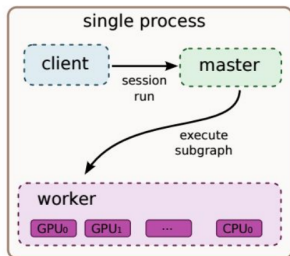
Many machines:

Not open source (yet) =(



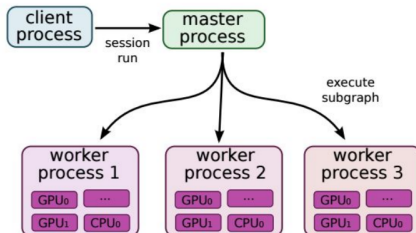
Single machine:

Like other frameworks



Many machines:

Not open source (yet) =(



Pros:

- 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

Cons:

- Distributed models, but not open-source yet
- Slower than other frameworks right now
- Much “fatter” than Torch; more magic
- Not many pretrained models



Comparison

	Caffe	Torch	Theano	TensorFlow
Language	C++,Python	Lua	Python	Python
Pretrained	Very good	Very good	OK	Poor
Multi-GPU: Data parallel	Yes	Yes	Ok	Yes
Multi-GPU: Model parallel	No	Yes	Experimental	Yes (best)
Readable	Yes(C++)	Yes(lua)	V. Poor	V. Poor
Good at RNN	No	Mediocre	Yes	Yes (best)
Visualization	No	OK	No	Yes
Automatic differentiation	No	Yes(nngraph)	Yes	Yes



- Use pretrained models for features? Caffe



- Use pretrained models for features? Caffe
- Fine-tune AlexNet for new classes? Use Caffe





- Use pretrained models for features? Caffe
- Fine-tune AlexNet for new classes? Use Caffe
- Segmentation? (Classify every pixel)
 - Need pretrained model (Caffe, Torch, Lasagna)
 - If loss function exists in Caffe: Use Caffe
 - If you want to write your own loss: Use Torch





- Use pretrained models for features? Caffe
- Fine-tune AlexNet for new classes? Use Caffe
- Segmentation? (Classify every pixel)
 - Need pretrained model (Caffe, Torch, Lasagna)
 - If loss function exists in Caffe: Use Caffe
 - If you want to write your own loss: Use Torch
- Object Detection?
 - Need pretrained model (Torch, Caffe, Lasagne)
 - Need lots of custom imperative code (NOT Lasagne)
 - Use Caffe + Python or Torch



- Use pretrained models for features? Caffe
- Fine-tune AlexNet for new classes? Use Caffe
- Segmentation? (Classify every pixel)
 - Need pretrained model (Caffe, Torch, Lasagna)
 - If loss function exists in Caffe: Use Caffe
 - If you want to write your own loss: Use Torch
- Object Detection?
 - Need pretrained model (Torch, Caffe, Lasagne)
 - Need lots of custom imperative code (NOT Lasagne)
 - Use Caffe + Python or Torch
- Implement BatchNorm?
 - Don't want to derive gradient? Theano or TensorFlow
 - Implement efficient backward pass? Use Torch



- Image Captioning with finetuning?
 - Need pretrained models (Caffe, Torch, Lasagne)
 - Need RNNs (Torch or Lasagne)
 - Use Torch or Lasagna





- Image Captioning with finetuning?
 - Need pretrained models (Caffe, Torch, Lasagne)
 - Need RNNs (Torch or Lasagne)
 - Use Torch or Lasagna
- Language modeling with new RNN structure?
 - Need easy recurrent nets (NOT Caffe, Torch)
 - No need for pretrained models
 - Use Theano or TensorFlow





- Image Captioning with finetuning?
 - Need pretrained models (Caffe, Torch, Lasagne)
 - Need RNNs (Torch or Lasagne)
 - Use Torch or Lasagna
- Language modeling with new RNN structure?
 - Need easy recurrent nets (NOT Caffe, Torch)
 - No need for pretrained models
 - Use Theano or TensorFlow
- Huge parallelism, distributed training required? Tensorflow





- Image Captioning with finetuning?
 - Need pretrained models (Caffe, Torch, Lasagne)
 - Need RNNs (Torch or Lasagne)
 - Use Torch or Lasagna
- Language modeling with new RNN structure?
 - Need easy recurrent nets (NOT Caffe, Torch)
 - No need for pretrained models
 - Use Theano or TensorFlow
- Huge parallelism, distributed training required? Tensorflow
- Simple parallelism required? Torch or Caffe





- Image Captioning with finetuning?
 - Need pretrained models (Caffe, Torch, Lasagne)
 - Need RNNs (Torch or Lasagne)
 - Use Torch or Lasagna
- Language modeling with new RNN structure?
 - Need easy recurrent nets (NOT Caffe, Torch)
 - No need for pretrained models
 - Use Theano or TensorFlow
- Huge parallelism, distributed training required? Tensorflow
- Simple parallelism required? Torch or Caffe
- Reinforcement learning? Torch or Tensorflow



- Feature extraction / finetuning existing models: Use Caffe



- Feature extraction / finetuning existing models: Use Caffe
- Complex uses of pretrained models: Use Lasagne or Torch



- Feature extraction / finetuning existing models: Use Caffe
- Complex uses of pretrained models: Use Lasagne or Torch
- Write your own layers: Use Torch



- Feature extraction / finetuning existing models: Use Caffe
- Complex uses of pretrained models: Use Lasagne or Torch
- Write your own layers: Use Torch
- Crazy RNNs: Use Theano or Tensorflow



- Feature extraction / finetuning existing models: Use Caffe
- Complex uses of pretrained models: Use Lasagne or Torch
- Write your own layers: Use Torch
- Crazy RNNs: Use Theano or Tensorflow
- Huge model, need model parallelism: Use TensorFlow



Thank you

Thank you.

