

# Assignment Lecture 3

Convolutional Neural Networks

Transfer Learning

Deep learning frameworks

Neural Networks on GPU

# Practical

- Assignment 1 grades are out!
- Additional lab hour tomorrow 14-16
- Group sign up sheet is up!

# Deep Learning Hardware

CPU, GPU, TPU

# Computer Hardware



# We have two hardware choices:

- NVIDIA GPU
- Google Tensor Processing Unit (TPU)
- AMD? Really not used.

# Nvidia Country

IDI HPC newest arrival

Tesla V100



NVIDIA DGX-2

Explore the powerful components of DGX-2.

- ① NVIDIA TESLA V100 32GB, SXM3
- ② 16 TOTAL GPUs FOR BOTH BOARDS, 5120B TOTAL HBM2 MEMORY  
Each GPU board with 8 NVIDIA Tesla V100.
- ③ 12 TOTAL NVSWITCHES  
High Speed Interconnect, 2.4 TB/sec bisection bandwidth.
- ④ 8 EDR INFINIBAND/100 GbE ETHERNET  
1600 Gb/sec Bi-directional Bandwidth and Low-Latency.
- ⑤ PCIE SWITCH COMPLEX
- ⑥ TWO INTEL XEON PLATINUM CPUS
- ⑦ 1.5 TB SYSTEM MEMORY
- ⑧ DUAL 10/25 GbE ETHERNET
- ⑨ 30 TB NVME SSDS INTERNAL STORAGE

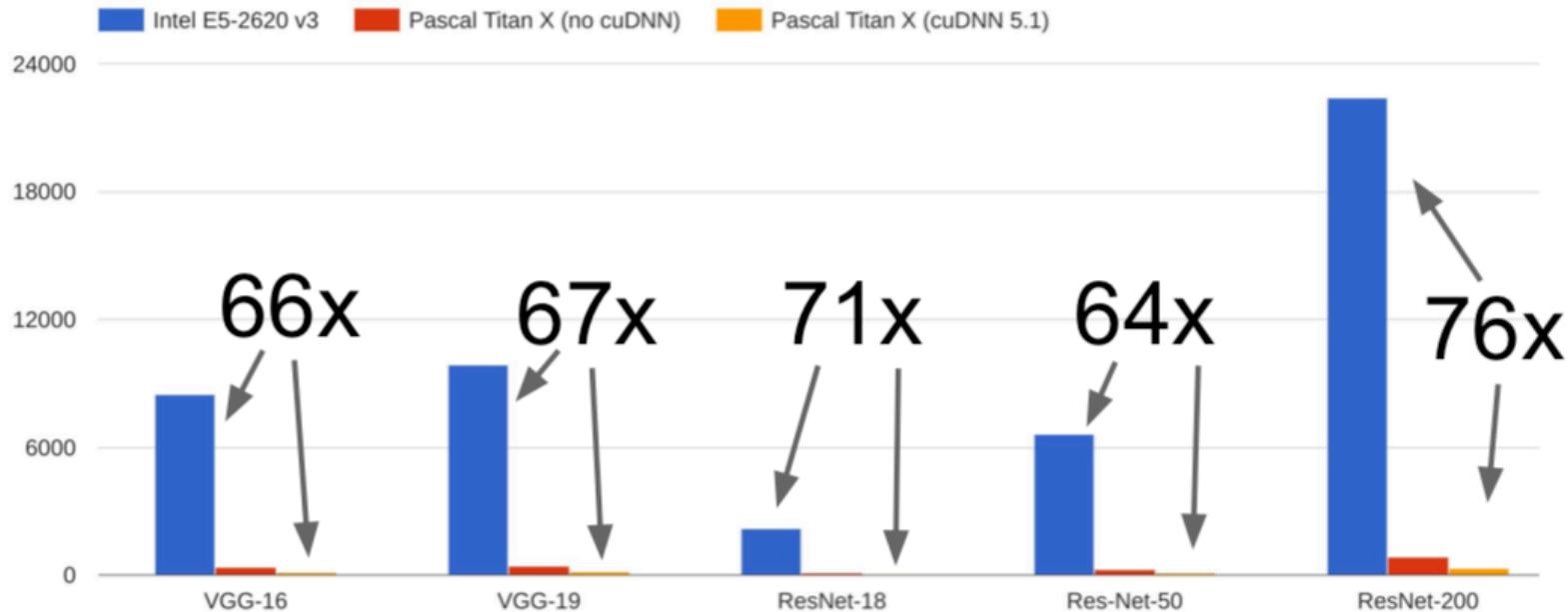
# GPU vs CPU

	Cores	Clock Speed	Memory	Price	Speed
<b>CPU</b> (Intel Core i7-7700k)	4 (8 threads with hyperthreading)	4.2 GHz	System RAM	\$339	~540 GFLOPs FP32
<b>GPU</b> (NVIDIA GTX 1080 Ti)	3584	1.6 GHz	11 GB GDDR5 X	\$699	~11.4 TFLOPs FP32

GPU: thousands of “dumb” cores: Great for parallel tasks

Neural Networks: “Only” matrix multiplication, easy in parallel

# GPU vs CPU for CNNs



# GPU vs CPU: ResNet-200

- Forward pass:
  - Pascal Titan X: 104ms
  - CPU: Dual Xeon E5-2630 v3: 8,666ms (**83x slower**)
- Backward pass:
  - Pascal Titan X: 191 ms
  - CPU: Dual Xeon E5-2630 v3: 13,758 ms (**72x slower**)

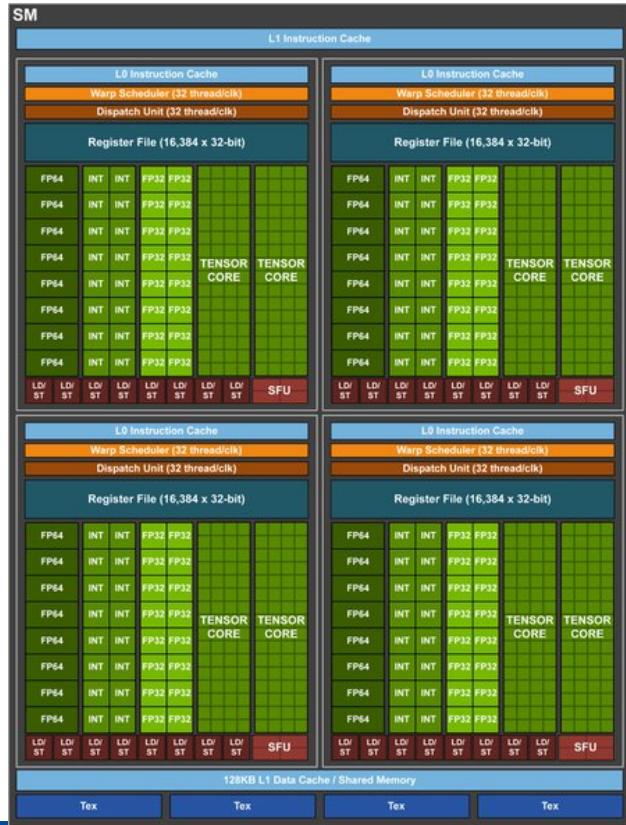
# New in GPU hardware

- CUDA cores
  - General parallelization cores in GPU
- **Tensor cores:**
  - Cores specialized for neural networks

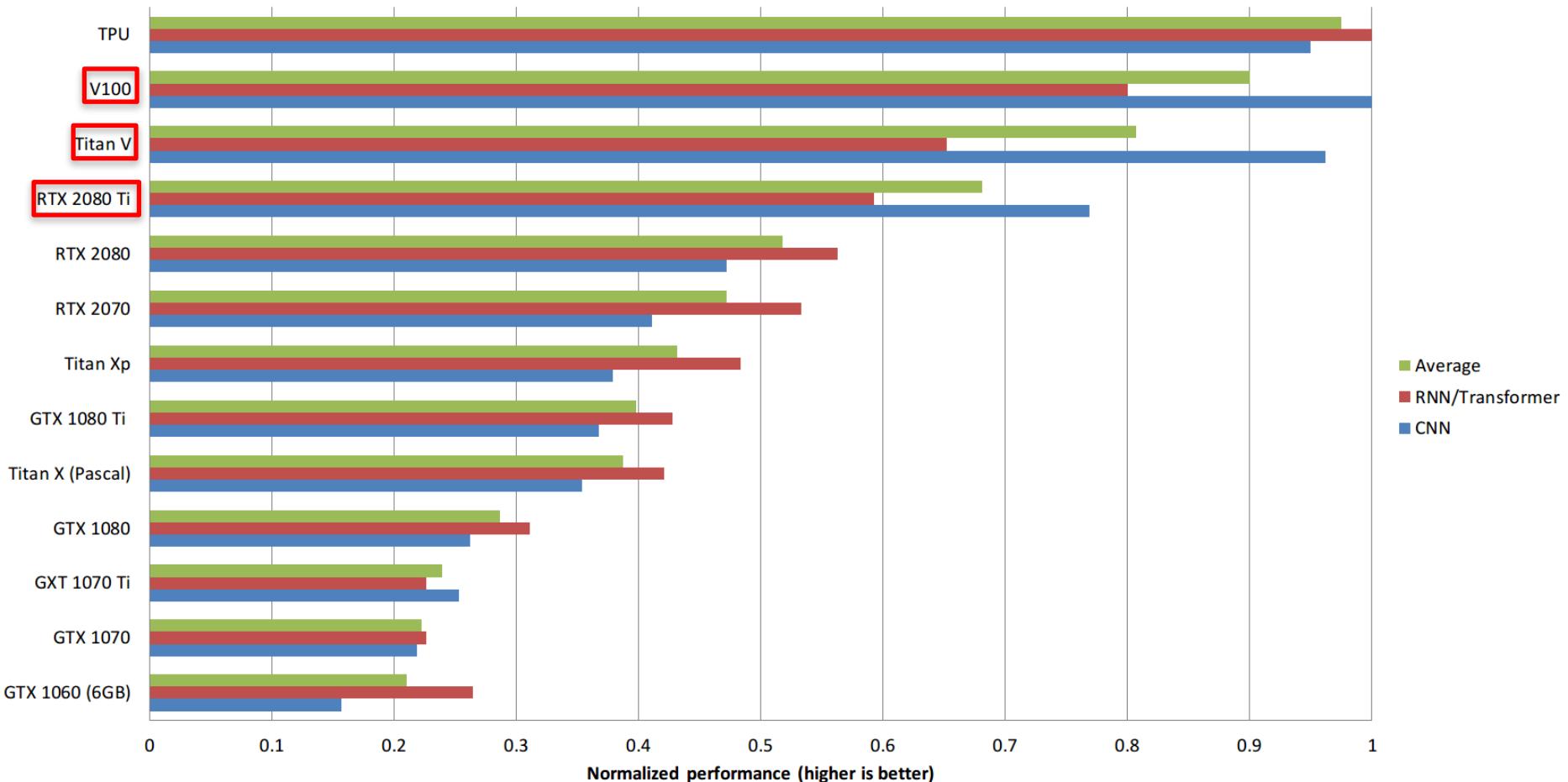
# Pascal Architecture



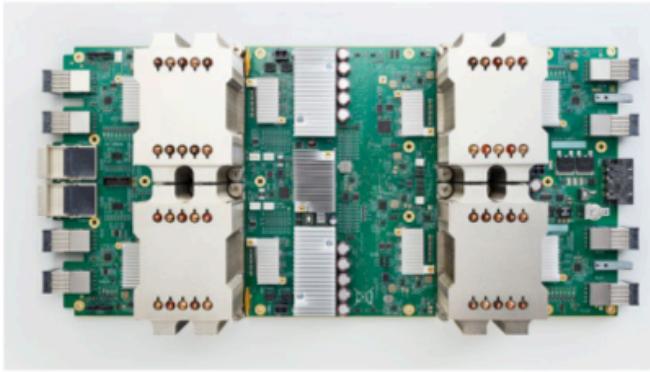
# Volta Architecture



# Performance



# Tensor Processing Units (TPU)



Google Cloud TPU  
= 180 TFLOPs of compute!

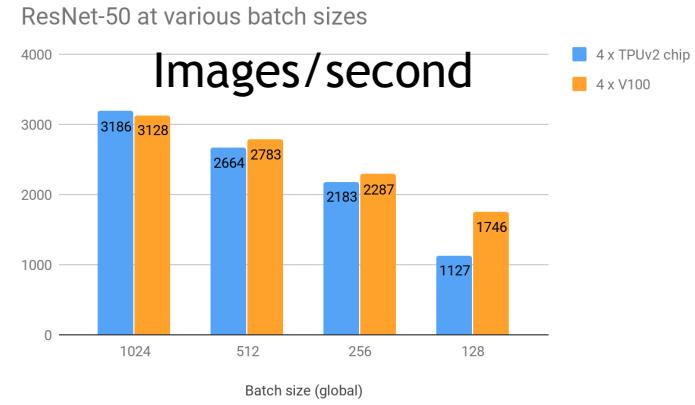


NVIDIA Tesla V100  
= 125 TFLOPs of compute

NVIDIA Tesla P100 = 11 TFLOPs of compute  
GTX 580 = 0.2 TFLOPs

# Tensor Processing Units (TPU)

- Special made for neural networks
- Produced by Google
- Not possible to buy, only rent on GCP!

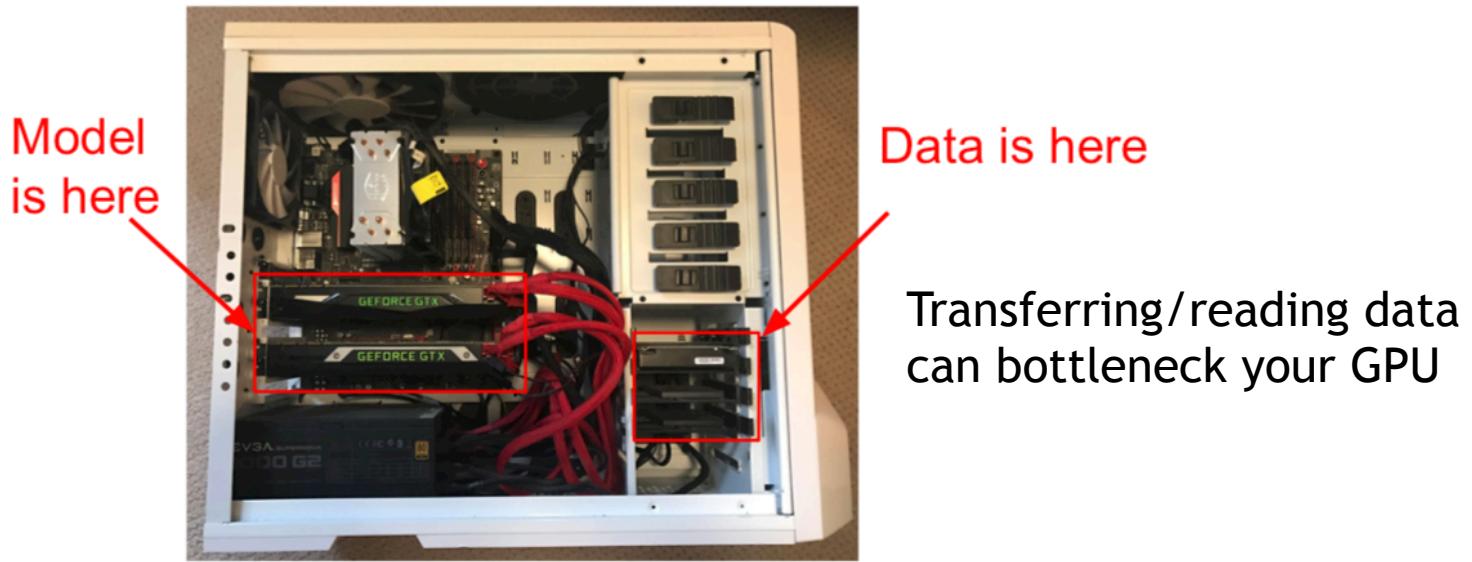


# Issue: CPU <-> GPU Communication



Justin Johnson's Computer!

# Issue: CPU <-> GPU Communication



# Issue: CPU <-> GPU Communication



Solution:

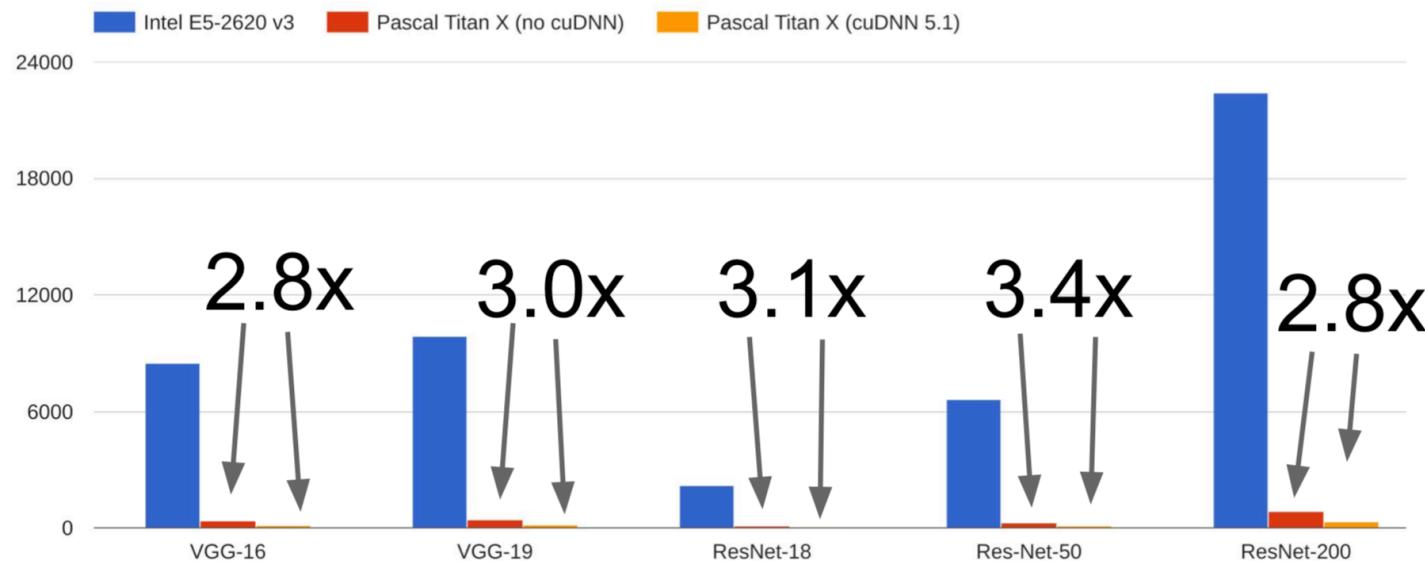
- Read data into RAM
- Use SSD
- Prefetch data with threads
- Minimize GPU<->CPU communication

# Deep Learning Software

# Deep Learning Software

- Deep learning software runs on (NVIDIA):

- CUDA
- cudNN



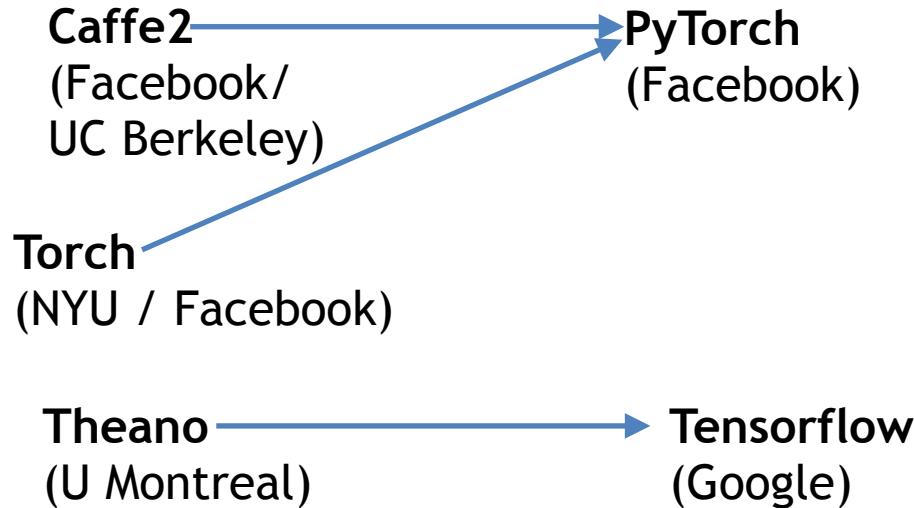
# Frameworks

**Caffe2**  
(Facebook/  
UC Berkeley)

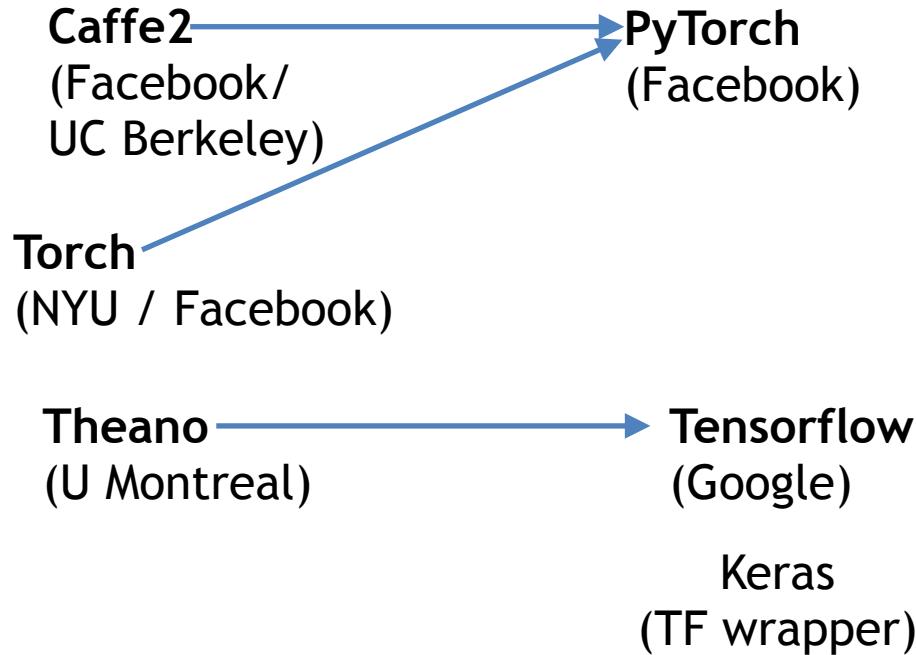
**Torch**  
(NYU / Facebook)

**Theano**  
(U Montreal)

# Frameworks



# Frameworks



Less popular stuff

PaddlePaddle  
(Baidu)

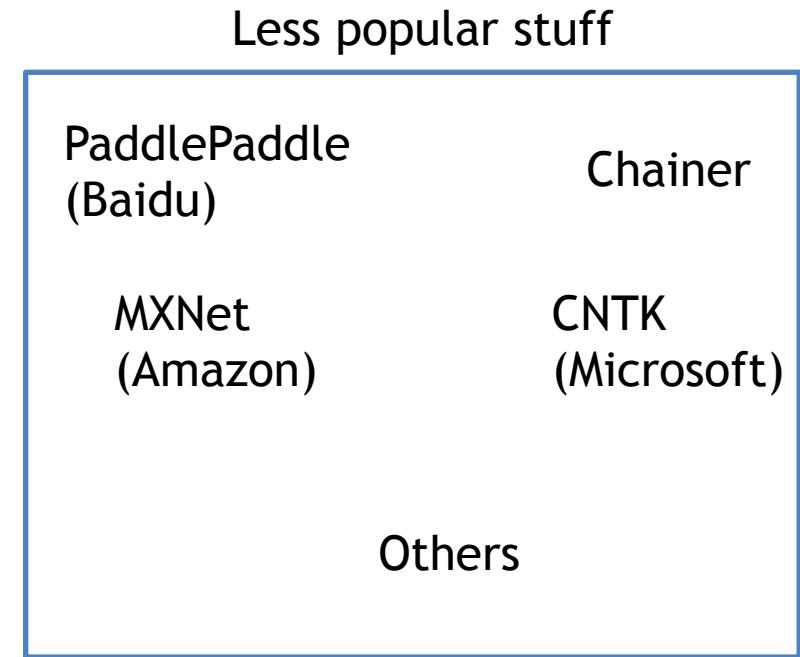
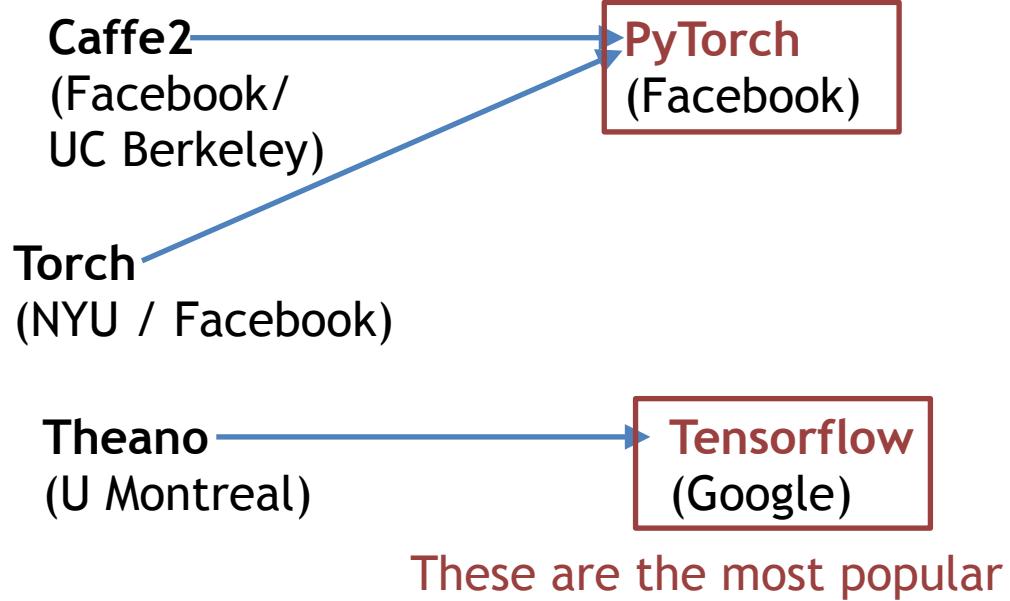
Chainer

MXNet  
(Amazon)

CNTK  
(Microsoft)

Others

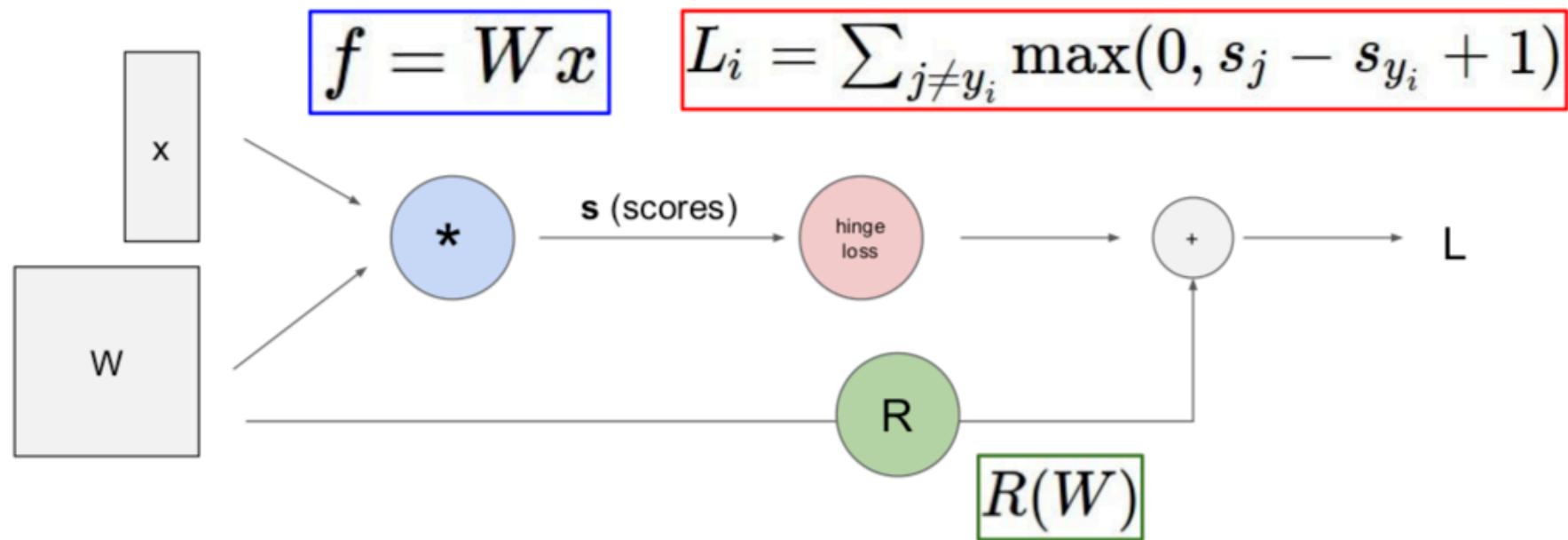
# Frameworks



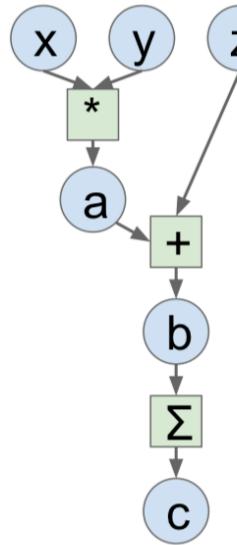
# Requirement of deep learning frameworks

- Quickly implement and test ideas
- Automatically compute gradients
- Run it efficient

# Neural Networks = Directed Acyclic Graphs



# Computational Graph



# Computational Graph

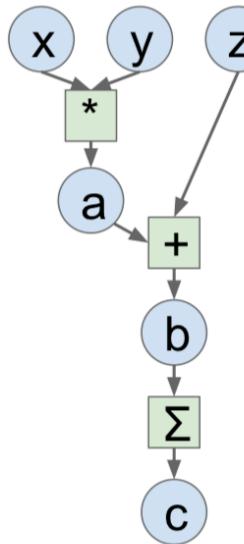
Numpy

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

N, D = 3, 4

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

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



# Computational Graph

Numpy

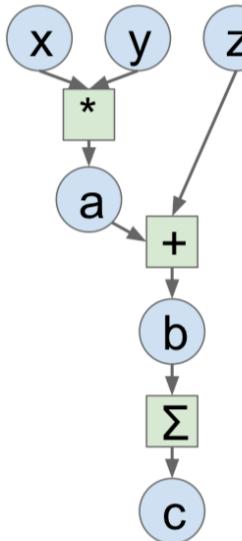
```
import numpy as np
np.random.seed(0)

N, D = 3, 4

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

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

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



# Computational Graph

## Numpy

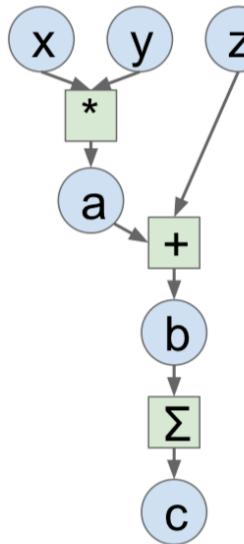
```
import numpy as np
np.random.seed(0)

N, D = 3, 4

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

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

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



Good:

- Simple, clean API

Bad:

- Have to compute gradients ourselves
- Can't run on GPU

# Computational Graph

Numpy

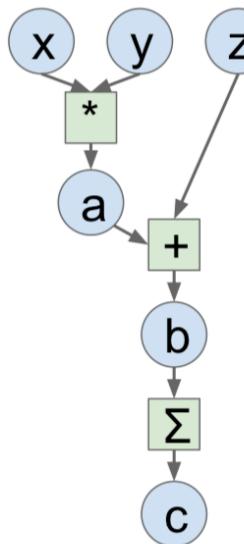
```
import numpy as np
np.random.seed(0)

N, D = 3, 4

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

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

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



PyTorch

```
import torch

N, D = 3, 4
x = torch.randn(N, D)
y = torch.randn(N, D)
z = torch.randn(N, D)

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

# Computational Graph

Numpy

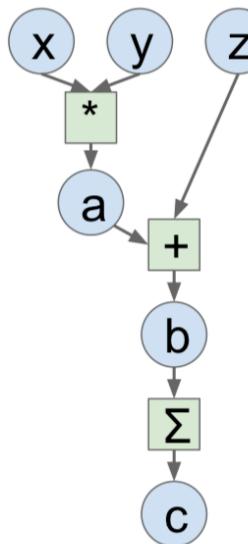
```
import numpy as np
np.random.seed(0)

N, D = 3, 4

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

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

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



PyTorch

```
import torch

N, D = 3, 4
x = torch.randn(N, D, requires_grad=True)
y = torch.randn(N, D)
z = torch.randn(N, D)

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

c.backward()
print(x.grad)
```

Torch handles all gradients!

# **Small sidetrack: Tensorflow vs Pytorch**

# Computational Graph

- Both tensorflow and pytorch define a **Directed Acyclic Graph**
- Tensorflow/Keras: **Static graph**, defined before training.
- Pytorch: **Dynamic graph**, defined during training

# Pytorch: Dynamic Graph

- Pros:
  - Intuitive code
  - “Numpy on GPU”
  - Easy debugging
  - “Pythonic code”
- Cons:
  - Dynamic graph is slower  
(really though?)
  - Smaller community, but its growing!

A graph is created on the fly



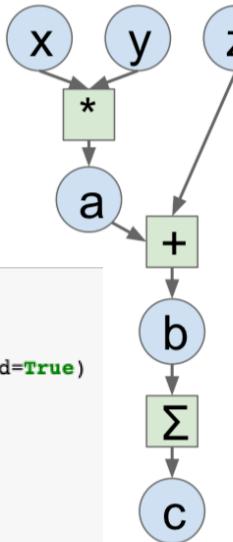
```
W_h = torch.randn(20, 20, requires_grad=True)
W_x = torch.randn(20, 10, requires_grad=True)
x = torch.randn(1, 10)
prev_h = torch.randn(1, 20)
```



# Computational Graph

PyTorch

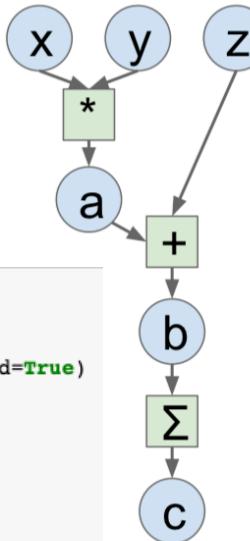
```
1 import torch
2 N, D = 3, 4
3
4 x = torch.randn(N, D, requires_grad=True)
5 y = torch.randn(N, D)
6 z = torch.randn(N, D)
7
8 a = x * y
9 b = a + z
10 c = b.sum()
11
12 c.backward()
13 print(x.grad)
```



# Computational Graph

PyTorch

```
1 import torch
2 N, D = 3, 4
3
4 x = torch.randn(N, D, requires_grad=True)
5 y = torch.randn(N, D)
6 z = torch.randn(N, D)
7
8 a = x * y
9 b = a + z
10 c = b.sum()
11
12 c.backward()
13 print(x.grad)
```



Tensorflow

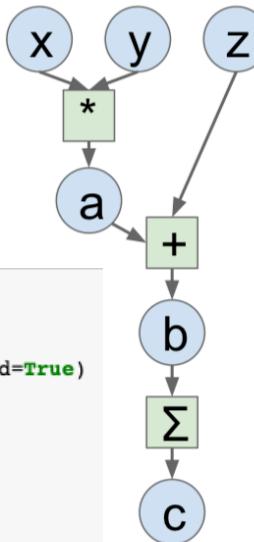
```
1 N, D = 3, 4
2
3 x = tf.placeholder(tf.float32, shape=(N,D))
4 y = tf.placeholder(tf.float32, shape=(N,D))
5 z = tf.placeholder(tf.float32, shape=(N,D))
6
7 a = tf.multiply(x,y)
8 b = tf.add(a,z)
9 c = tf.reduce_sum(b)
10 grad_x = tf.gradients(c, [x])
```

First, define the graph

# Computational Graph

PyTorch

```
1 import torch
2 N, D = 3, 4
3
4 x = torch.randn(N, D, requires_grad=True)
5 y = torch.randn(N, D)
6 z = torch.randn(N, D)
7
8 a = x * y
9 b = a + z
10 c = b.sum()
11
12 c.backward()
13 print(x.grad)
```



Tensorflow

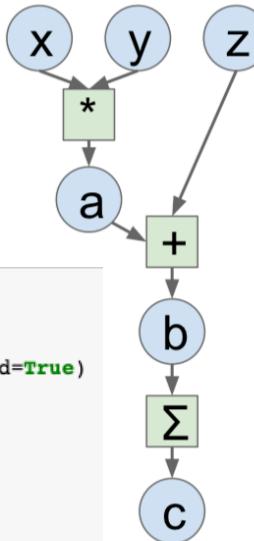
```
1 N, D = 3, 4
2
3 x = tf.placeholder(tf.float32, shape=(N,D))
4 y = tf.placeholder(tf.float32, shape=(N,D))
5 z = tf.placeholder(tf.float32, shape=(N,D))
6
7 a = tf.multiply(x,y)
8 b = tf.add(a,z)
9 c = tf.reduce_sum(b)
10 grad_x = tf.gradients(c, [x])
11
12 with tf.Session() as sess:
13     sess.run(tf.global_variables_initializer())
```

Then, initialize a `tf.Session`

# Computational Graph

PyTorch

```
1 import torch
2 N, D = 3, 4
3
4 x = torch.randn(N, D, requires_grad=True)
5 y = torch.randn(N, D)
6 z = torch.randn(N, D)
7
8 a = x * y
9 b = a + z
10 c = b.sum()
11
12 c.backward()
13 print(x.grad)
```



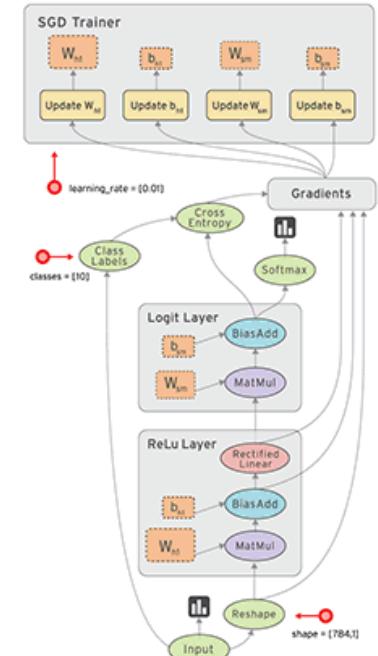
Tensorflow

```
1 N, D = 3, 4
2
3 x = tf.placeholder(tf.float32, shape=(N,D))
4 y = tf.placeholder(tf.float32, shape=(N,D))
5 z = tf.placeholder(tf.float32, shape=(N,D))
6
7 a = tf.multiply(x,y)
8 b = tf.add(a,z)
9 c = tf.reduce_sum(b)
10 grad_x = tf.gradients(c, [x])
11
12 with tf.Session() as sess:
13     sess.run(tf.global_variables_initializer())
14     values = {
15         x: np.random.randn(N,D),
16         y: np.random.randn(N,D),
17         z: np.random.randn(N,D)
18     }
19     gradient_x = sess.run(grad_x, feed_dict=values)
```

Then, perform the forward pass

# Tensorflow/Keras: Static Graph

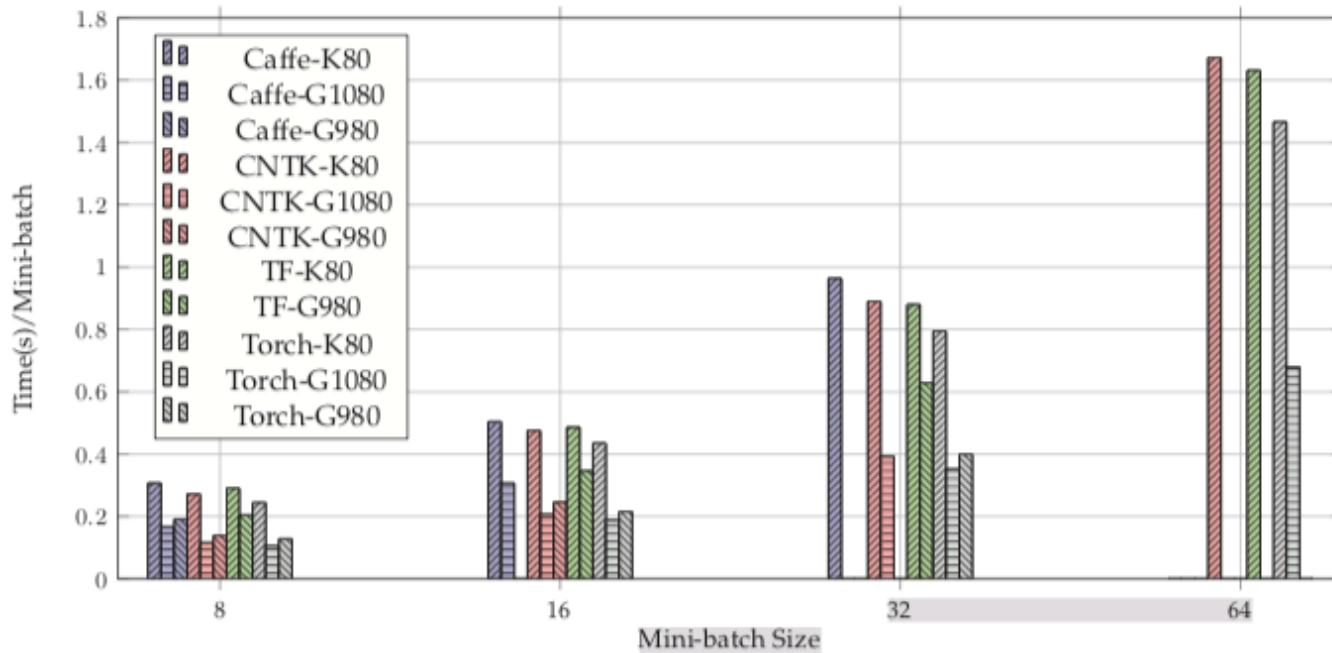
- Pros:
  - Faster?
  - Large community
  - Keras is a simpler API
- Cons:
  - Not intuitive
  - Hard to debug



# Dynamic vs Static Graphs

- Tensorflow is adding dynamic graphs in TF 2.0
- Pytorch 1.0 includes static graph
- Lines are blurring....

# Framework speed comparison

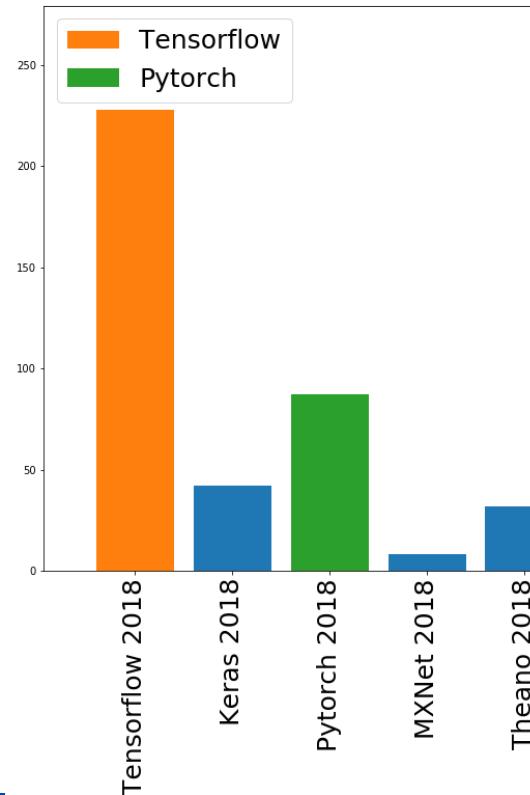


(d) ResNet-50 training speed on GPUs.

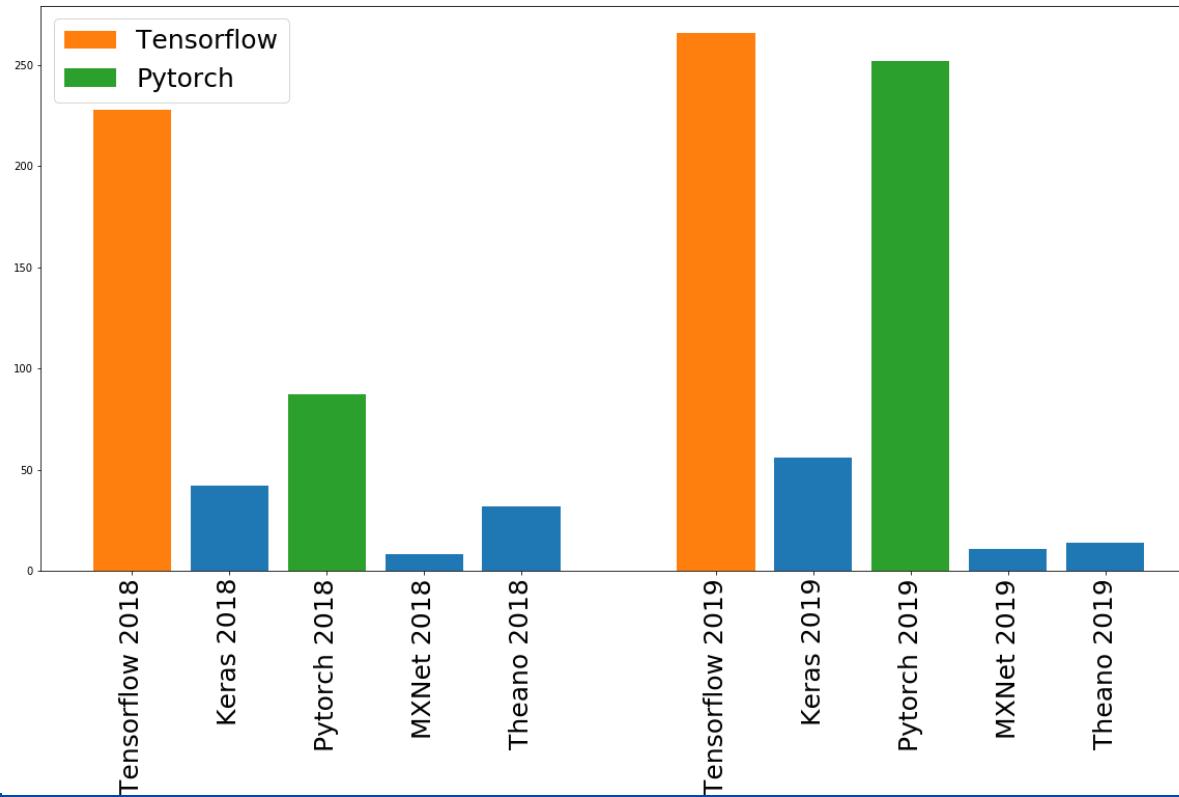
# Framework comparison

“There is no single software tool that can consistently outperform others”

# What is actually used? (ICLR 2018/19)



# What is actually used? (ICLR 2018/19)



# My advice:

- **Pytorch(My favorite):** Has an intuitive API that works for both high level and low level neural networks. Dynamic graphs make it easy to develop and debug.
- **Tensorflow/Keras:** Has a larger community and a wide usage in companies. Is hard to debug and graph setup is initially hard to understand. Probably have to use Keras.

# **PyTorch** **(in-depth)**

# Pytorch: Fundamental Concepts

- Tensor: Like a numpy array, but can run on GPUs
- Module: A neural network layer; stores states and learnable weights

# Pytorch version

- We recommend Pytorch 1.0, which was released November 2018
- Pytorch 0.4.0/0.4.1 is fine too, not much syntax change
- **NB:** Be careful when looking at older Pytorch code!

# Pytorch: Tensors

Example: The 2-layer neural network from Assignment 2

```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 w1 = torch.randn(I, J, requires_grad=True)
6 w2 = torch.randn(J, C, requires_grad=True)
7
8 loss_function = torch.nn.NLLLoss()
9 losses = []
10 for i in range(10):
11     for (X,Y) in dataloader:
12         X = pre_process_image(X)
13         # forward pass
14         z_j = X.mm(w1)
15         a_j = torch.sigmoid(z_j)
16         z_k = a_j.mm(w2)
17         y_k = torch.softmax(z_k, dim=1)
18         # Compute loss
19         loss = loss_function(y_k, Y)
20         losses.append(loss)
21         # Backpropagation
22         loss.backward()
23         with torch.no_grad():
24             w1 -= learning_rate * w1.grad
25             w2 -= learning_rate * w2.grad
26             w1.grad.zero_()
27             w2.grad.zero_()
```

# Pytorch: Tensors

Initialize the weights randomly

```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 w1 = torch.randn(I, J, requires_grad=True)
6 w2 = torch.randn(J, C, requires_grad=True)
7
8 loss_function = torch.nn.NLLLoss()
9 losses = []
10 for i in range(10):
11     for (X,Y) in dataloader:
12         X = pre_process_image(X)
13         # forward pass
14         z_j = X.mm(w1)
15         a_j = torch.sigmoid(z_j)
16         z_k = a_j.mm(w2)
17         y_k = torch.softmax(z_k, dim=1)
18         # Compute loss
19         loss = loss_function(y_k, Y)
20         losses.append(loss)
21         # Backpropagation
22         loss.backward()
23         with torch.no_grad():
24             w1 -= learning_rate * w1.grad
25             w2 -= learning_rate * w2.grad
26             w1.grad.zero_()
27             w2.grad.zero_()
```

# Pytorch: Tensors

Do want to save gradients w.r.t weights

```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 w1 = torch.randn(I, J, requires_grad=True)
6 w2 = torch.randn(J, C, requires_grad=True)
7
8 loss_function = torch.nn.NLLLoss()
9 losses = []
10 for i in range(10):
11     for (X,Y) in dataloader:
12         X = pre_process_image(X)
13         # forward pass
14         z_j = X.mm(w1)
15         a_j = torch.sigmoid(z_j)
16         z_k = a_j.mm(w2)
17         y_k = torch.softmax(z_k, dim=1)
18         # Compute loss
19         loss = loss_function(y_k, Y)
20         losses.append(loss)
21         # Backpropagation
22         loss.backward()
23         with torch.no_grad():
24             w1 -= learning_rate * w1.grad
25             w2 -= learning_rate * w2.grad
26             w1.grad.zero_()
27             w2.grad.zero_()
```

# Pytorch: Tensors

Define our loss function:

Cross Entropy Loss

```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 w1 = torch.randn(I, J, requires_grad=True)
6 w2 = torch.randn(J, C, requires_grad=True)
7
8 loss_function = torch.nn.NLLLoss()
9 losses = []
10 for i in range(10):
11     for (X,Y) in dataloader:
12         X = pre_process_image(X)
13         # forward pass
14         z_j = X.mm(w1)
15         a_j = torch.sigmoid(z_j)
16         z_k = a_j.mm(w2)
17         y_k = torch.softmax(z_k, dim=1)
18         # Compute loss
19         loss = loss_function(y_k, Y)
20         losses.append(loss)
21         # Backpropagation
22         loss.backward()
23         with torch.no_grad():
24             w1 -= learning_rate * w1.grad
25             w2 -= learning_rate * w2.grad
26             w1.grad.zero_()
27             w2.grad.zero_()
28
```

# Pytorch: Tensors

Pytorch uses “dataloaders” to efficiently load datasets

```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 w1 = torch.randn(I, J, requires_grad=True)
6 w2 = torch.randn(J, C, requires_grad=True)
7
8 loss_function = torch.nn.NLLLoss()
9 losses = []
10 for i in range(10):
11     for (X,Y) in dataloader:
12         X = pre_process_image(X)
13         # forward pass
14         z_j = X.mm(w1)
15         a_j = torch.sigmoid(z_j)
16         z_k = a_j.mm(w2)
17         y_k = torch.softmax(z_k, dim=1)
18         # Compute loss
19         loss = loss_function(y_k, Y)
20         losses.append(loss)
21         # Backpropagation
22         loss.backward()
23         with torch.no_grad():
24             w1 -= learning_rate * w1.grad
25             w2 -= learning_rate * w2.grad
26             w1.grad.zero_()
27             w2.grad.zero_()
28
```

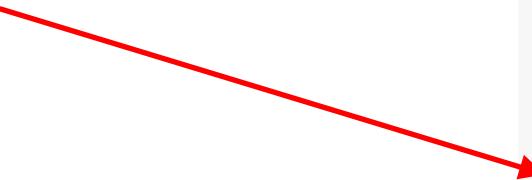
# Pytorch: Tensors

Bias trick and normalization

```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 w1 = torch.randn(I, J, requires_grad=True)
6 w2 = torch.randn(J, C, requires_grad=True)
7
8 loss_function = torch.nn.NLLLoss()
9 losses = []
10 for i in range(10):
11     for (X,Y) in dataloader:
12         X = pre process image(X) # Line 12
13         # forward pass
14         z_j = X.mm(w1)
15         a_j = torch.sigmoid(z_j)
16         z_k = a_j.mm(w2)
17         y_k = torch.softmax(z_k, dim=1)
18         # Compute loss
19         loss = loss_function(y_k, Y)
20         losses.append(loss)
21         # Backpropagation
22         loss.backward()
23         with torch.no_grad():
24             w1 -= learning_rate * w1.grad
25             w2 -= learning_rate * w2.grad
26             w1.grad.zero_()
27             w2.grad.zero_()
```

# Pytorch: Tensors

Define the forward pass



```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 w1 = torch.randn(I, J, requires_grad=True)
6 w2 = torch.randn(J, C, requires_grad=True)
7
8 loss_function = torch.nn.NLLLoss()
9 losses = []
10 for i in range(10):
11     for (X,Y) in dataloader:
12         X = pre_process_image(X)
13         # forward pass
14         z_j = X.mm(w1)
15         a_j = torch.sigmoid(z_j)
16         z_k = a_j.mm(w2)
17         y_k = torch.softmax(z_k, dim=1)
18         # Compute loss
19         loss = loss_function(y_k, Y)
20         losses.append(loss)
21         # Backpropagation
22         loss.backward()
23         with torch.no_grad():
24             w1 -= learning_rate * w1.grad
25             w2 -= learning_rate * w2.grad
26             w1.grad.zero_()
27             w2.grad.zero_()
28
```

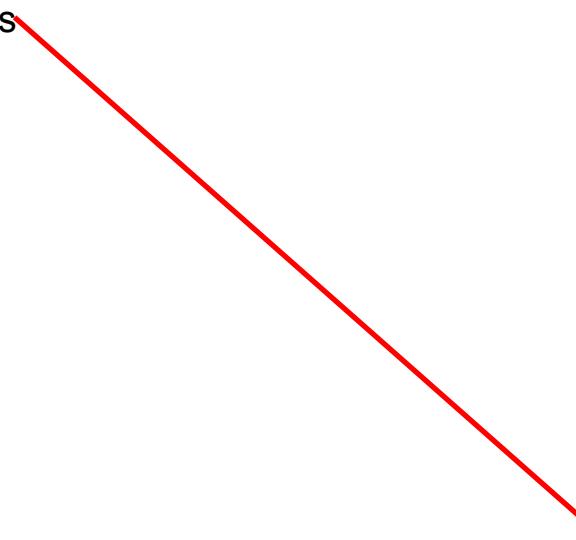
# Pytorch: Tensors

Compute the loss

```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 w1 = torch.randn(I, J, requires_grad=True)
6 w2 = torch.randn(J, C, requires_grad=True)
7
8 loss_function = torch.nn.NLLLoss()
9 losses = []
10 for i in range(10):
11     for (X,Y) in dataloader:
12         X = pre_process_image(X)
13         # forward pass
14         z_j = X.mm(w1)
15         a_j = torch.sigmoid(z_j)
16         z_k = a_j.mm(w2)
17         y_k = torch.softmax(z_k, dim=1)
18         # Compute loss
19         loss = loss_function(y_k, Y)
20         losses.append(loss)
21         # Backpropagation
22         loss.backward()
23         with torch.no_grad():
24             w1 -= learning_rate * w1.grad
25             w2 -= learning_rate * w2.grad
26             w1.grad.zero_()
27             w2.grad.zero_()
28
```

# Pytorch: Tensors

Backpropagate the loss



```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 w1 = torch.randn(I, J, requires_grad=True)
6 w2 = torch.randn(J, C, requires_grad=True)
7
8 loss_function = torch.nn.NLLLoss()
9 losses = []
10 for i in range(10):
11     for (X,Y) in dataloader:
12         X = pre_process_image(X)
13         # forward pass
14         z_j = X.mm(w1)
15         a_j = torch.sigmoid(z_j)
16         z_k = a_j.mm(w2)
17         y_k = torch.softmax(z_k, dim=1)
18         # Compute loss
19         loss = loss_function(y_k, Y)
20         losses.append(loss)
21         # Backpropagation
22         loss.backward()
23         with torch.no_grad():
24             w1 -= learning_rate * w1.grad
25             w2 -= learning_rate * w2.grad
26             w1.grad.zero_()
27             w2.grad.zero_()
28
```

# Pytorch: Tensors

Make gradient step on weights

`torch.no_grad()` means “don’t build a computational graph here

```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 w1 = torch.randn(I, J, requires_grad=True)
6 w2 = torch.randn(J, C, requires_grad=True)
7
8 loss_function = torch.nn.NLLLoss()
9 losses = []
10 for i in range(10):
11     for (X,Y) in dataloader:
12         X = pre_process_image(X)
13         # forward pass
14         z_j = X.mm(w1)
15         a_j = torch.sigmoid(z_j)
16         z_k = a_j.mm(w2)
17         y_k = torch.softmax(z_k, dim=1)
18         # Compute loss
19         loss = loss_function(y_k, Y)
20         losses.append(loss)
21         # Backpropagation
22         loss.backward()
23         with torch.no_grad():
24             w1 -= learning_rate * w1.grad
25             w2 -= learning_rate * w2.grad
26             w1.grad.zero_()
27             w2.grad.zero_()
```

# Pytorch: torch.nn

Higher lever wrapper for defining neural networks

```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 model = nn.Sequential(
6     nn.Linear(I, J),
7     nn.Sigmoid(),
8     nn.Linear(J, C)
9     # No need for softmax, since its included in
10    # nn.CrossEntropyLoss()
11 )
12
13 loss_function = torch.nn.CrossEntropyLoss()
14 losses = []
15 for epoch in range(10):
16     for (X,Y) in dataloader:
17         X = pre_process_image(X)
18         # forward pass
19         y_k = model(X)
20         # Compute loss
21         loss = loss_function(y_k, Y)
22         losses.append(loss)
23         # Backpropagation
24         loss.backward()
25         with torch.no_grad():
26             for param in model.parameters():
27                 param -= learning_rate * param.grad
```

# Pytorch: torch.nn

Higher lever wrapper for defining neural networks

Define each layer in model.

Each layer is a nn.Module() object, containing learnable weights.

```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 model = nn.Sequential(
6     nn.Linear(I, J),
7     nn.Sigmoid(),
8     nn.Linear(J, C)
9     # No need for softmax, since its included in
10    # nn.CrossEntropyLoss()
11 )
12
13 loss_function = torch.nn.CrossEntropyLoss()
14 losses = []
15 for epoch in range(10):
16     for (X,Y) in dataloader:
17         X = pre_process_image(X)
18         # forward pass
19         y_k = model(X)
20         # Compute loss
21         loss = loss_function(y_k, Y)
22         losses.append(loss)
23         # Backpropagation
24         loss.backward()
25         with torch.no_grad():
26             for param in model.parameters():
27                 param -= learning_rate * param.grad
```

# Pytorch: torch.nn

Changed loss function to  
nn.CrossEntropyLoss.  
This includes the softmax!

```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 model = nn.Sequential(
6     nn.Linear(I, J),
7     nn.Sigmoid(),
8     nn.Linear(J, C)
9     # No need for softmax, since its included in
10    # nn.CrossEntropyLoss()
11 )
12
13 loss_function = torch.nn.CrossEntropyLoss()
14 losses = []
15 for epoch in range(10):
16     for (X,Y) in dataloader:
17         X = pre_process_image(X)
18         # forward pass
19         y_k = model(X)
20         # Compute loss
21         loss = loss_function(y_k, Y)
22         losses.append(loss)
23         # Backpropagation
24         loss.backward()
25         with torch.no_grad():
26             for param in model.parameters():
27                 param -= learning_rate * param.grad
```

# Pytorch: torch.nn

Simplifies our forward pass!

```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 model = nn.Sequential(
6     nn.Linear(I, J),
7     nn.Sigmoid(),
8     nn.Linear(J, C)
9     # No need for softmax, since its included in
10    # nn.CrossEntropyLoss()
11 )
12
13 loss_function = torch.nn.CrossEntropyLoss()
14 losses = []
15 for epoch in range(10):
16     for (X,Y) in dataloader:
17         X = pre_process_image(X)
18         # forward pass
19         y_k = model(X)
20         # Compute loss
21         loss = loss_function(y_k, Y)
22         losses.append(loss)
23         # Backpropagation
24         loss.backward()
25         with torch.no_grad():
26             for param in model.parameters():
27                 param -= learning_rate * param.grad
```

# Pytorch: torch.nn

Compute loss and perform backward pass

Each weight in model has  
requires\_grad=True by default

```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 model = nn.Sequential(
6     nn.Linear(I, J),
7     nn.Sigmoid(),
8     nn.Linear(J, C)
9     # No need for softmax, since its included in
10    # nn.CrossEntropyLoss()
11 )
12
13 loss_function = torch.nn.CrossEntropyLoss()
14 losses = []
15 for epoch in range(10):
16     for (X,Y) in dataloader:
17         X = pre_process_image(X)
18         # forward pass
19         y_k = model(X)
20         # Compute loss
21         loss = loss_function(y_k, Y)
22         losses.append(loss)
23         # Backpropagation
24         loss.backward()
25         with torch.no_grad():
26             for param in model.parameters():
27                 param -= learning_rate * param.grad
```

# Pytorch: torch.nn

Perform our gradient step  
(and disable gradients)

```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 model = nn.Sequential(
6     nn.Linear(I, J),
7     nn.Sigmoid(),
8     nn.Linear(J, C)
9     # No need for softmax, since its included in
10    # nn.CrossEntropyLoss()
11 )
12
13 loss_function = torch.nn.CrossEntropyLoss()
14 losses = []
15 for epoch in range(10):
16     for (X,Y) in dataloader:
17         X = pre_process_image(X)
18         # forward pass
19         y_k = model(X)
20         # Compute loss
21         loss = loss_function(y_k, Y)
22         losses.append(loss)
23         # Backpropagation
24         loss.backward()
25         with torch.no_grad():
26             for param in model.parameters():
27                 param -= learning_rate * param.grad
```

# Pytorch: torch.optim

Final piece you need to know

Implements **Stochastic Gradient Descent**

Input: our learnable parameters (weights + biases)  
+ learning rate

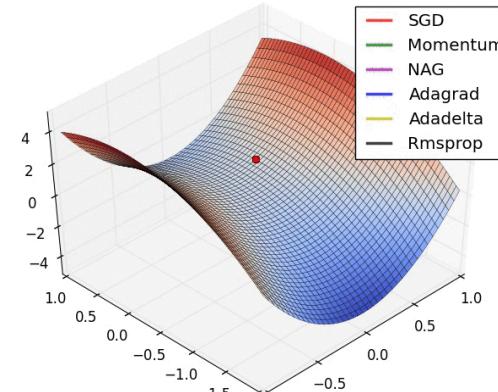
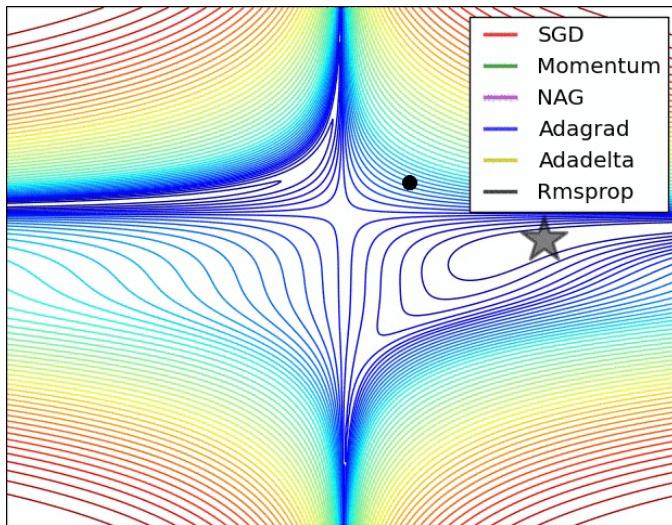
```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 model = nn.Sequential(
6     nn.Linear(I, J),
7     nn.Sigmoid(),
8     nn.Linear(J, C)
9     # No need for softmax, since its included in
10    # nn.CrossEntropyLoss()
11 )
12
13 loss_function = torch.nn.CrossEntropyLoss()
14 optimizer = torch.optim.SGD(model.parameters(),
15                             lr=learning_rate)
16 losses = []
17 for epoch in range(10):
18     for (X,Y) in dataloader:
19         X = pre_process_image(X)
20         # forward pass
21         y_k = model(X)
22         # Compute loss
23         loss = loss_function(y_k, Y)
24         losses.append(loss)
25         # Backpropagation
26         loss.backward()
27         optimizer.step()
28         optimizer.zero_grad()
```

# Pytorch: torch.optim

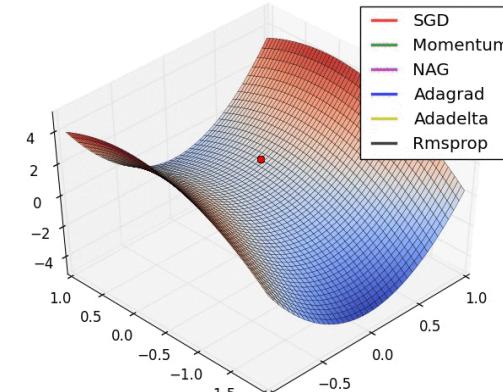
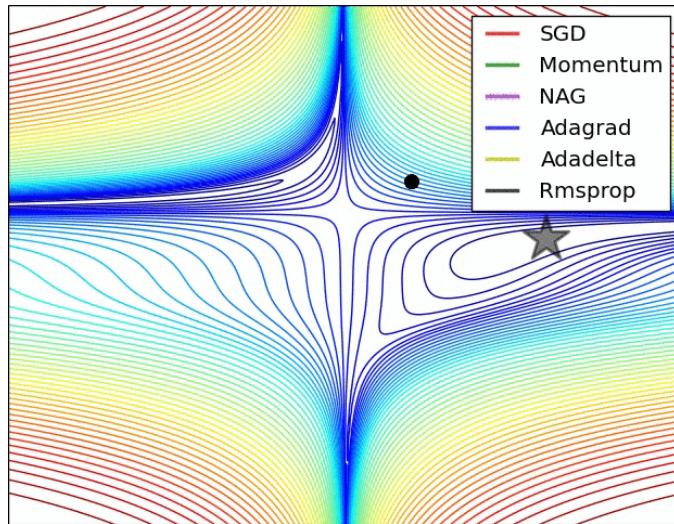
Perform gradient step and reset the gradients

```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 model = nn.Sequential(
6     nn.Linear(I, J),
7     nn.Sigmoid(),
8     nn.Linear(J, C)
9     # No need for softmax, since its included in
10    # nn.CrossEntropyLoss()
11 )
12
13 loss_function = torch.nn.CrossEntropyLoss()
14 optimizer = torch.optim.SGD(model.parameters(),
15                             lr=learning_rate)
16 losses = []
17 for epoch in range(10):
18     for (X,Y) in dataloader:
19         X = pre_process_image(X)
20         # forward pass
21         y_k = model(X)
22         # Compute loss
23         loss = loss_function(y_k, Y)
24         losses.append(loss)
25         # Backpropagation
26         loss.backward()
27         optimizer.step()
28         optimizer.zero_grad()
```

# The choice of optimizer is important!



# The choice of optimizer is important!



Try out `optim.Adam` in your assignment!

# Pytorch: nn.Module

A PyTorch **Module** is a neural network layer; it inputs and outputs tensors

Can contain weights or other modules

Required for more complex layers

Easily customizable layers

```
1 class TwoLayerNet(nn.Module):
2     def __init__(self):
3         super().__init__()
4         I, J, C = 785, 64, 10
5         self.layer1 = nn.Sequential(
6             nn.Linear(I,J),
7             nn.Sigmoid()
8         )
9         self.layer2 = nn.Linear(J, C)
10    def forward(self, x):
11        x = self.layer1(x)
12        x = self.layer2(x)
13        return x
14
15 learning_rate = 1e-3
16 model = TwoLayerNet()
17
18 loss_function = torch.nn.CrossEntropyLoss()
19 optimizer = torch.optim.SGD(model.parameters(),
20                             lr=learning_rate)
21 losses = []
22 for epoch in range(2):
23     for (X,Y) in dataloader:
24         X = pre_process_image(X)
25         # forward pass
26         y_k = model(X)
27         # Compute loss
28         loss = loss_function(y_k, Y)
29         losses.append(loss)
30         # Backpropagation
31         loss.backward()
32         optimizer.step()
33         optimizer.zero_grad()
```

# Pytorch: nn.Module

Start with defining the model

```
1 class TwoLayerNet(nn.Module):
2     def __init__(self):
3         super().__init__()
4         I, J, C = 785, 64, 10
5         self.layer1 = nn.Sequential(
6             nn.Linear(I,J),
7             nn.Sigmoid()
8         )
9         self.layer2 = nn.Linear(J, C)
10    def forward(self, x):
11        x = self.layer1(x)
12        x = self.layer2(x)
13        return x
14
15 learning_rate = 1e-3
16 model = TwoLayerNet()
17
18 loss_function = torch.nn.CrossEntropyLoss()
19 optimizer = torch.optim.SGD(model.parameters(),
20                             lr=learning_rate)
21 losses = []
22 for epoch in range(2):
23     for (X,Y) in dataloader:
24         X = pre_process_image(X)
25         # forward pass
26         y_k = model(X)
27         # Compute loss
28         loss = loss_function(y_k, Y)
29         losses.append(loss)
30         # Backpropagation
31         loss.backward()
32         optimizer.step()
33         optimizer.zero_grad()
```

# Pytorch: nn.Module

Called when we initialize our model

```
1 class TwoLayerNet(nn.Module):
2     def __init__(self):
3         super().__init__()
4         I, J, C = 785, 64, 10
5         self.layer1 = nn.Sequential(
6             nn.Linear(I,J),
7             nn.Sigmoid()
8         )
9         self.layer2 = nn.Linear(J, C)
10    def forward(self, x):
11        x = self.layer1(x)
12        x = self.layer2(x)
13        return x
14
15    learning_rate = 1e-3
16    model = TwoLayerNet()
17
18    loss_function = torch.nn.CrossEntropyLoss()
19    optimizer = torch.optim.SGD(model.parameters(),
20                                lr=learning_rate)
21    losses = []
22    for epoch in range(2):
23        for (X,Y) in dataloader:
24            X = pre_process_image(X)
25            # forward pass
26            y_k = model(X)
27            # Compute loss
28            loss = loss_function(y_k, Y)
29            losses.append(loss)
30            # Backpropagation
31            loss.backward()
32            optimizer.step()
33            optimizer.zero_grad()
```

# Pytorch: nn.Module

Called when we perform forward pass

```
1 class TwoLayerNet(nn.Module):
2     def __init__(self):
3         super().__init__()
4         I, J, C = 785, 64, 10
5         self.layer1 = nn.Sequential(
6             nn.Linear(I,J),
7             nn.Sigmoid()
8         )
9         self.layer2 = nn.Linear(J, C)
10    def forward(self, x):
11        x = self.layer1(x)
12        x = self.layer2(x)
13        return x
14
15 learning_rate = 1e-3
16 model = TwoLayerNet()
17
18 loss_function = torch.nn.CrossEntropyLoss()
19 optimizer = torch.optim.SGD(model.parameters(),
20                             lr=learning_rate)
21 losses = []
22 for epoch in range(2):
23     for (X,Y) in dataloader:
24         X = pre_process_image(X)
25         # forward pass
26         y_k = model(X)
27         # Compute loss
28         loss = loss_function(y_k, Y)
29         losses.append(loss)
30         # Backpropagation
31         loss.backward()
32         optimizer.step()
33         optimizer.zero_grad()
34
```

# Pytorch: DataLoaders

A **DataLoader** wraps a dataset and provides features such as:

- Data augmentation
- Data pre-processing
- mini-batch shuffling and splitting

```
1 class TwoLayerNet(nn.Module):
2     def __init__(self):
3         super().__init__()
4         I, J, C = 785, 64, 10
5         self.layer1 = nn.Sequential(
6             nn.Linear(I,J),
7             nn.Sigmoid()
8         )
9         self.layer2 = nn.Linear(J, C)
10    def forward(self, x):
11        x = self.layer1(x)
12        x = self.layer2(x)
13        return x
14
15 learning_rate = 1e-3
16 batch_size=32
17 model = TwoLayerNet()
18
19 dataloader_train, dataloader_test = load_mnist(batch_size)
20
21 loss_function = torch.nn.CrossEntropyLoss()
22 optimizer = torch.optim.SGD(model.parameters(),
23                             lr=learning_rate)
24 losses = []
25 for epoch in range(2):
26     for (X_batch,Y_batch) in dataloader_train:
27         X_batch = pre_process_image(X_batch)
28         # forward pass
29         y_k = model(X_batch)
30         # Compute loss
31         loss = loss_function(y_k, Y_batch)
32         losses.append(loss)
33         # Backpropagation
34         loss.backward()
35         optimizer.step()
36         optimizer.zero_grad()
```

# Pytorch: DataLoaders

Iterates over each batch in a epoch

```
1 class TwoLayerNet(nn.Module):
2     def __init__(self):
3         super().__init__()
4         I, J, C = 785, 64, 10
5         self.layer1 = nn.Sequential(
6             nn.Linear(I,J),
7             nn.Sigmoid()
8         )
9         self.layer2 = nn.Linear(J, C)
10    def forward(self, x):
11        x = self.layer1(x)
12        x = self.layer2(x)
13        return x
14
15 learning_rate = 1e-3
16 batch_size=32
17 model = TwoLayerNet()
18
19 dataloader_train, dataloader_test = load_mnist(batch_size)
20
21 loss_function = torch.nn.CrossEntropyLoss()
22 optimizer = torch.optim.SGD(model.parameters(),
23                             lr=learning_rate)
24 losses = []
25 for epoch in range(2):
26     for (X_batch,Y_batch) in dataloader_train:
27         X_batch = pre_process_image(X_batch)
28         # forward pass
29         y_k = model(X_batch)
30         # Compute loss
31         loss = loss_function(y_k, Y_batch)
32         losses.append(loss)
33         # Backpropagation
34         loss.backward()
35         optimizer.step()
36         optimizer.zero_grad()
```

# Pytorch: On GPU

Utilizing GPU resources is simple!

```
1 class TwoLayerNet(nn.Module):
2     def __init__(self):
3         super().__init__()
4         I, J, C = 785, 64, 10
5         self.layer1 = nn.Sequential(
6             nn.Linear(I,J),
7             nn.Sigmoid()
8         )
9         self.layer2 = nn.Linear(J, C)
10    def forward(self, x):
11        x = self.layer1(x)
12        x = self.layer2(x)
13        return x
14
15 learning_rate = 1e-3
16 batch_size=32
17 model = TwoLayerNet().cuda()
18
19 dataloader_train, dataloader_test = load_mnist(batch_size)
20
21 loss_function = torch.nn.CrossEntropyLoss()
22 optimizer = torch.optim.SGD(model.parameters(),
23                             lr=learning_rate)
24 losses = []
25 for epoch in range(2):
26     for (X_batch,Y_batch) in dataloader_train:
27         X_batch = pre_process_image(X_batch)
28         X_batch, Y_batch = X_batch.cuda(), Y_batch.cuda()
29         # forward pass
30         y_k = model(X_batch)
31         # Compute loss
32         loss = loss_function(y_k, Y_batch)
33         losses.append(loss)
34         # Backpropagation
35         loss.backward()
36         optimizer.step()
37         optimizer.zero_grad()
```

# Pytorch: On GPU

Utilizing GPU resources is simple!

.cuda() transfers weights/tensors to GPU VRAM

```
1 class TwoLayerNet(nn.Module):
2     def __init__(self):
3         super().__init__()
4         I, J, C = 785, 64, 10
5         self.layer1 = nn.Sequential(
6             nn.Linear(I,J),
7             nn.Sigmoid()
8         )
9         self.layer2 = nn.Linear(J, C)
10    def forward(self, x):
11        x = self.layer1(x)
12        x = self.layer2(x)
13        return x
14
15    learning_rate = 1e-3
16    batch_size=32
17    model = TwoLayerNet().cuda()  
18
19    dataloader_train, dataloader_test = load_mnist(batch_size)
20
21    loss_function = torch.nn.CrossEntropyLoss()
22    optimizer = torch.optim.SGD(model.parameters(),
23                                lr=learning_rate)
24    losses = []
25    for epoch in range(2):
26        for (X_batch,Y_batch) in dataloader_train:
27            X_batch = pre_process_image(X_batch)
28            X_batch, Y_batch = X_batch.cuda(), Y_batch.cuda()
29            # forward pass
30            y_k = model(X_batch)
31            # Compute loss
32            loss = loss_function(y_k, Y_batch)
33            losses.append(loss)
34            # Backpropagation
35            loss.backward()
36            optimizer.step()
37            optimizer.zero_grad()
```

# Pytorch: On GPU

Utilizing GPU resources is simple!

.cuda() transfers weights/tensors to GPU VRAM

```
1 class TwoLayerNet(nn.Module):
2     def __init__(self):
3         super().__init__()
4         I, J, C = 785, 64, 10
5         self.layer1 = nn.Sequential(
6             nn.Linear(I,J),
7             nn.Sigmoid()
8         )
9         self.layer2 = nn.Linear(J, C)
10    def forward(self, x):
11        x = self.layer1(x)
12        x = self.layer2(x)
13        return x
14
15    learning_rate = 1e-3
16    batch_size=32
17    model = TwoLayerNet().cuda()
18
19    dataloader_train, dataloader_test = load_mnist(batch_size)
20
21    loss_function = torch.nn.CrossEntropyLoss()
22    optimizer = torch.optim.SGD(model.parameters(),
23                                lr=learning_rate)
24    losses = []
25    for epoch in range(2):
26        for (X_batch,Y_batch) in dataloader_train:
27            X_batch = pre_process_image(X_batch)
28            X_batch, Y_batch = X_batch.cuda(), Y_batch.cuda()
29            # forward pass
30            y_k = model(X_batch)
31            # Compute loss
32            loss = loss_function(y_k, Y_batch)
33            losses.append(loss)
34            # Backpropagation
35            loss.backward()
36            optimizer.step()
37            optimizer.zero_grad()
```

# Pytorch: On GPU

Utilizing GPU resources is simple!

.cuda() transfers weights/tensors to GPU VRAM

**CAREFUL:** Calling .cuda() without a NVIDIA GPU available will cause error!

```
1 class TwoLayerNet(nn.Module):
2     def __init__(self):
3         super().__init__()
4         I, J, C = 785, 64, 10
5         self.layer1 = nn.Sequential(
6             nn.Linear(I,J),
7             nn.Sigmoid()
8         )
9         self.layer2 = nn.Linear(J, C)
10    def forward(self, x):
11        x = self.layer1(x)
12        x = self.layer2(x)
13        return x
14
15    learning_rate = 1e-3
16    batch_size=32
17    model = TwoLayerNet().cuda()
18
19    dataloader_train, dataloader_test = load_mnist(batch_size)
20
21    loss_function = torch.nn.CrossEntropyLoss()
22    optimizer = torch.optim.SGD(model.parameters(),
23                                lr=learning_rate)
24    losses = []
25    for epoch in range(2):
26        for (X_batch,Y_batch) in dataloader_train:
27            X_batch = pre_process_image(X_batch)
28            X_batch, Y_batch = X_batch.cuda(), Y_batch.cuda()
29            # forward pass
30            y_k = model(X_batch)
31            # Compute loss
32            loss = loss_function(y_k, Y_batch)
33            losses.append(loss)
34            # Backpropagation
35            loss.backward()
36            optimizer.step()
37            optimizer.zero_grad()
```

# Pytorch: On GPU

Instead: Implement a `to_cuda()` function

```
1 def to_cuda(elements):
2     if torch.cuda.is_available():
3         if type(elements) == tuple or type(elements) == list:
4             return [x.cuda() for x in elements]
5         return elements.cuda()
6     return elements
```

```
1 class TwoLayerNet(nn.Module):
2     def __init__(self):
3         super().__init__()
4         I, J, C = 785, 64, 10
5         self.layer1 = nn.Sequential(
6             nn.Linear(I,J),
7             nn.Sigmoid()
8         )
9         self.layer2 = nn.Linear(J, C)
10    def forward(self, x):
11        x = self.layer1(x)
12        x = self.layer2(x)
13        return x
14
15 learning_rate = 1e-3
16 batch_size=32
17 model = to_cuda(TwoLayerNet())
18
19 dataloader_train, dataloader_test = load_mnist(batch_size)
20
21 loss_function = torch.nn.CrossEntropyLoss()
22 optimizer = torch.optim.SGD(model.parameters(),
23                             lr=learning_rate)
24 losses = []
25 for epoch in range(2):
26     for (X_batch,Y_batch) in dataloader_train:
27         X_batch = pre_process_image(X_batch)
28         X_batch, Y_batch = to_cuda([X_batch, Y_batch])
29         # forward pass
30         y_k = model(X_batch)
31         # Compute loss
32         loss = loss_function(y_k, Y_batch)
33         losses.append(loss)
34         # Backpropagation
35         loss.backward()
36         optimizer.step()
37         optimizer.zero_grad()
```

# Convolutional Neural Networks

Short Recap

# LeNet (LeCun et al. 1998)

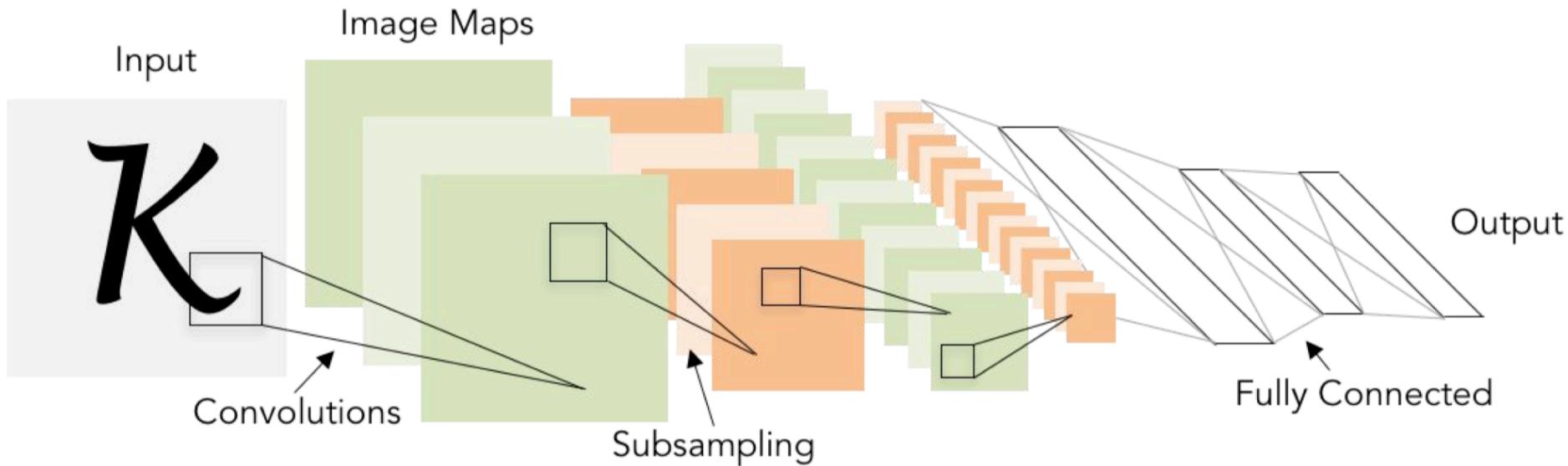
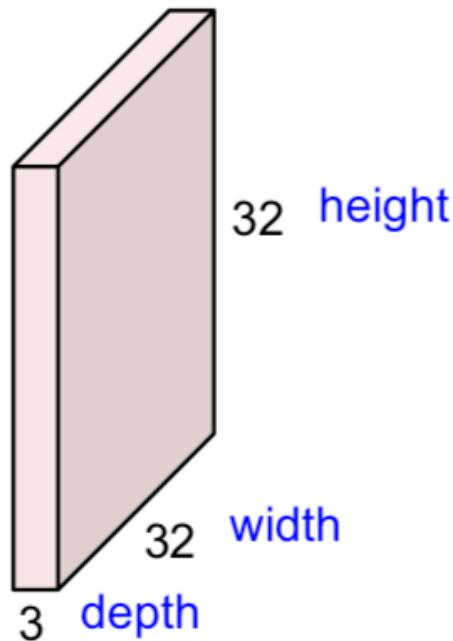


Illustration of LeCun et al. 1998 from CS231n 2017 Lecture 1

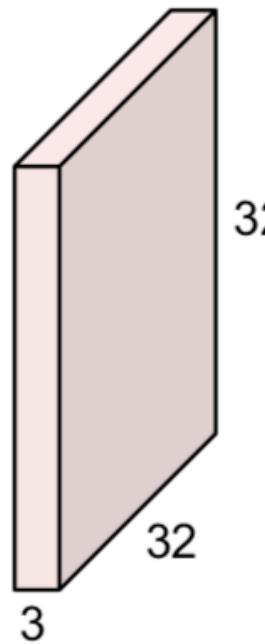
# Convolution Layer

32x32x3 image -> preserve spatial structure

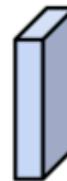


# Convolution Layer

32x32x3 image



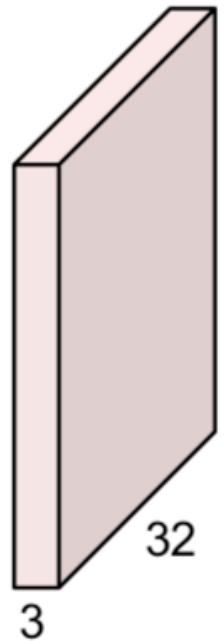
5x5x3 filter



**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

# Convolution Layer

32x32x3 image



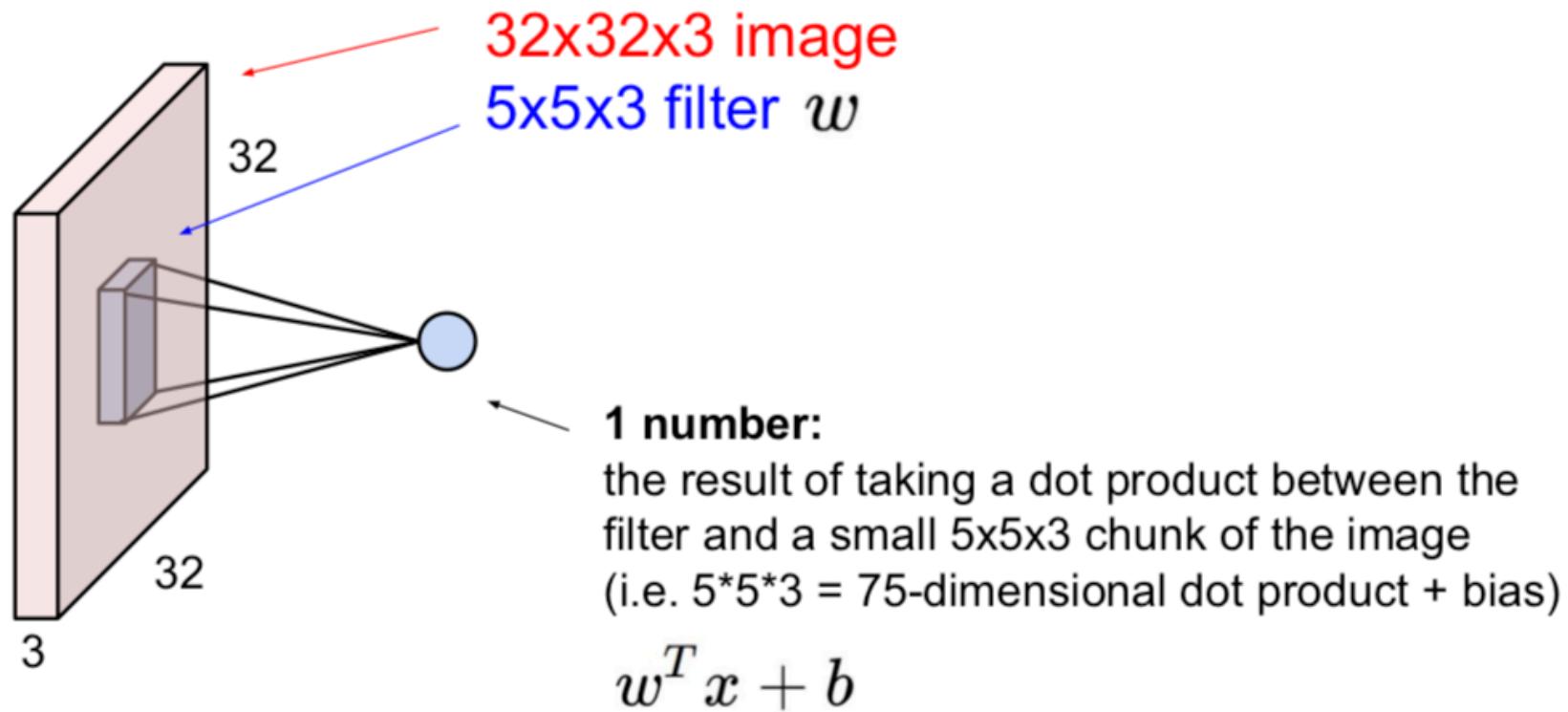
5x5x3 filter



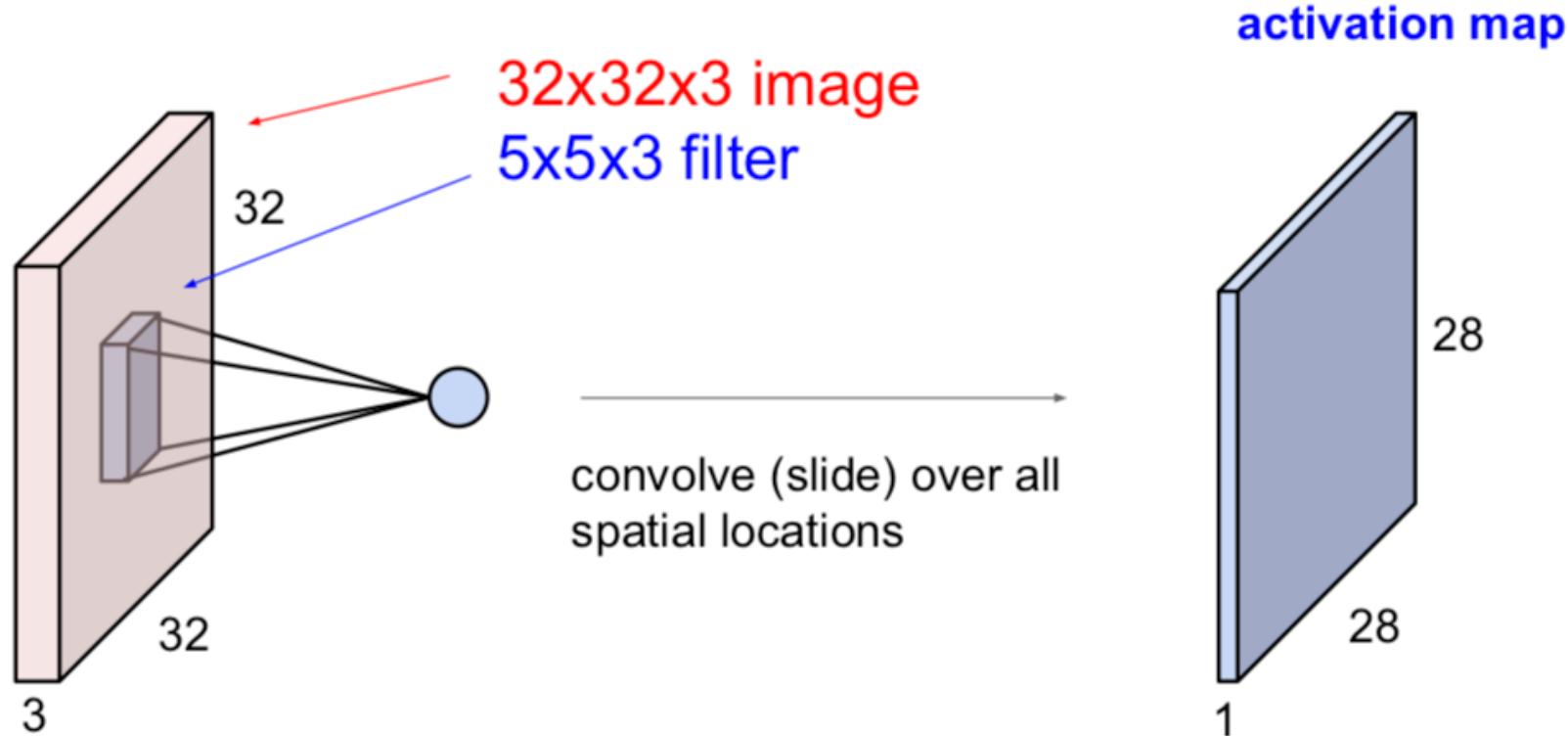
Filters always extend the full depth of the input volume

**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

# Convolution Layer

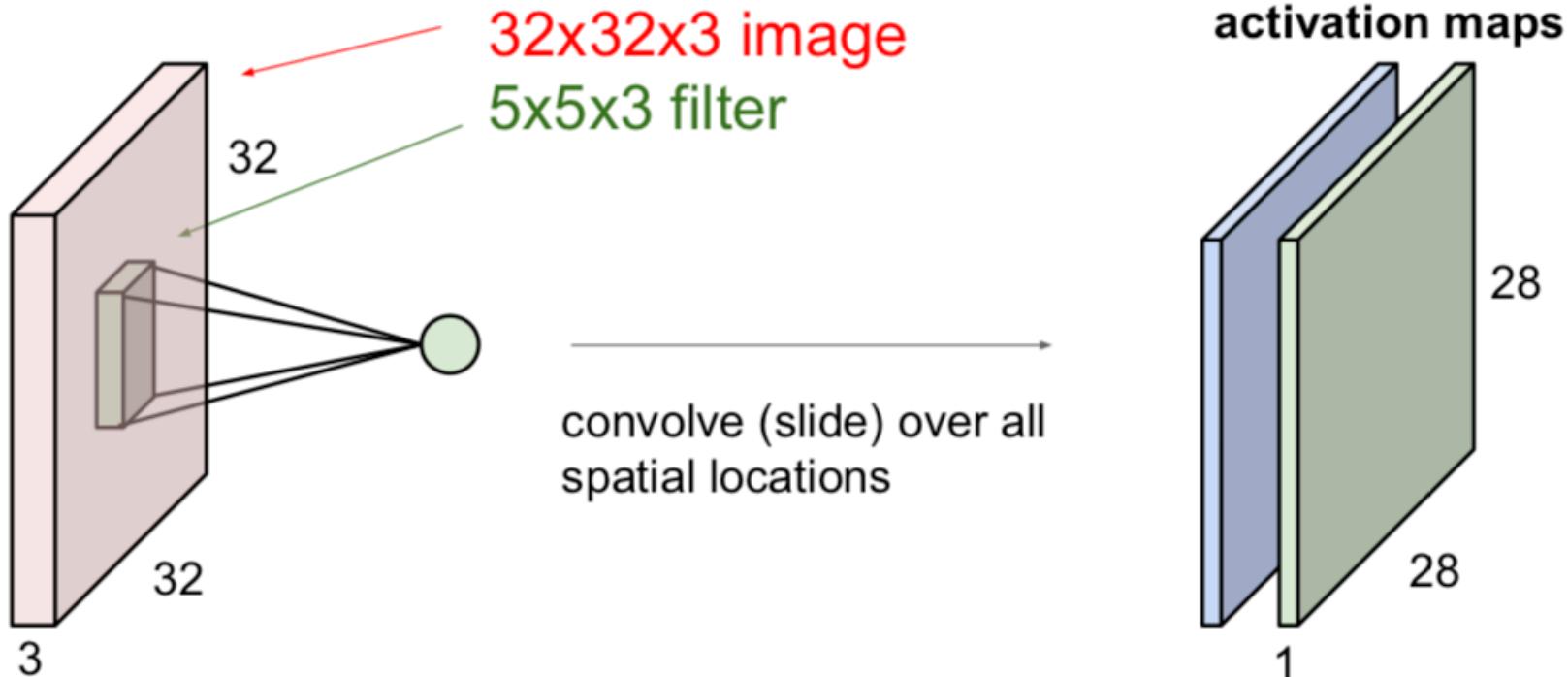


# Convolution Layer

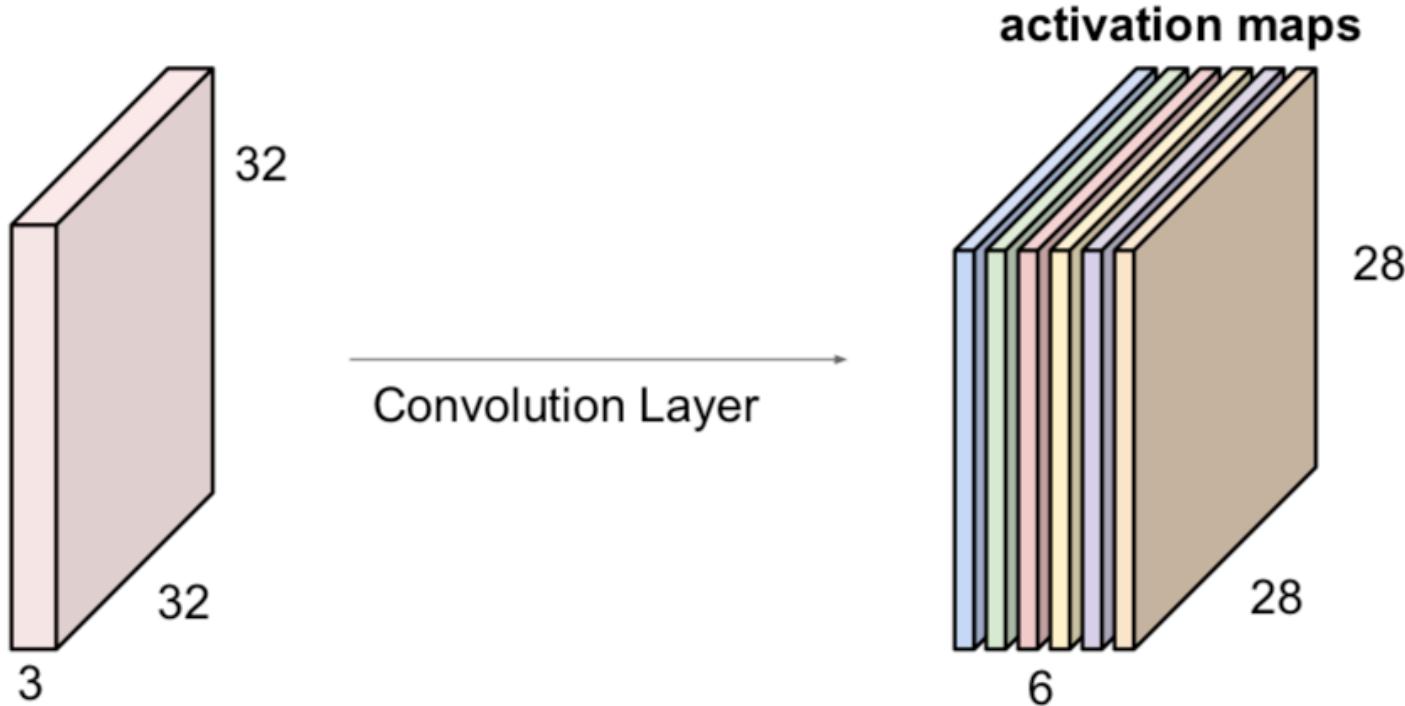


# Convolution Layer

consider a second, green filter



For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a “new image” of size  $28 \times 28 \times 6$ !

## Common settings:

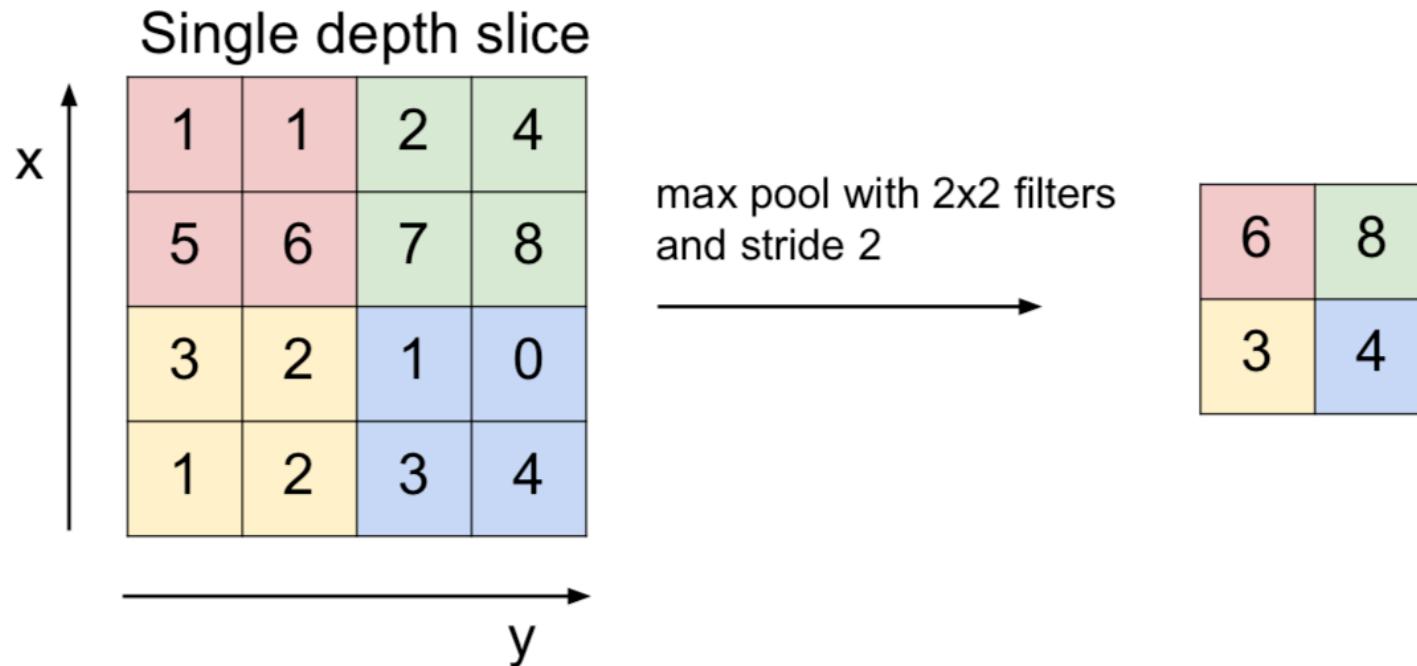
**Summary.** To summarize, the Conv Layer:

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
  - Number of filters  $K$ ,
  - their spatial extent  $F$ ,
  - the stride  $S$ ,
  - the amount of zero padding  $P$ .
- Produces a volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F + 2P)/S + 1$
  - $H_2 = (H_1 - F + 2P)/S + 1$  (i.e. width and height are computed equally by symmetry)
  - $D_2 = K$
- With parameter sharing, it introduces  $F \cdot F \cdot D_1$  weights per filter, for a total of  $(F \cdot F \cdot D_1) \cdot K$  weights and  $K$  biases.
- In the output volume, the  $d$ -th depth slice (of size  $W_2 \times H_2$ ) is the result of performing a valid convolution of the  $d$ -th filter over the input volume with a stride of  $S$ , and then offset by  $d$ -th bias.

$K = (\text{powers of 2, e.g. } 32, 64, 128, 512)$

- $F = 3, S = 1, P = 1$
- $F = 5, S = 1, P = 2$
- $F = 5, S = 2, P = ?$  (whatever fits)
- $F = 1, S = 1, P = 0$

# Max Pooling



# Max Pooling

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
  - their spatial extent  $F$ ,
  - the stride  $S$ ,
- Produces a volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F)/S + 1$
  - $H_2 = (H_1 - F)/S + 1$
  - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

# Pytorch: CNNs

Define a set of convolutional layers to extract **features**

Requires:

- in\_channels=1, since MNIST is grayscale
- out\_channels= number of filters in layer
- kernel\_size=filter width and height
- padding = number of 0's to pad on the side of image

```
1 class TwoLayerNet(nn.Module):  
2     def __init__(self):  
3         super().__init__()  
4         I, J, C = 785, 64, 10  
5         self.feature_extractor = nn.Sequential(  
6             nn.Conv2d(in_channels=1,  
7                     out_channels=32,  
8                     kernel_size=3,  
9                     padding=1),  
10            nn.MaxPool2d(kernel_size=2, stride=2),  
11            nn.Conv2d(in_channels=32,  
12                     out_channels=64,  
13                     kernel_size=3,  
14                     padding=1),  
15            nn.MaxPool2d(kernel_size=2, stride=2)  
16        )  
17        self.classifier = nn.Sequential(  
18            nn.Linear(64*7*7, 64),  
19            nn.ReLU(),  
20            nn.Linear(64, 10)  
21        )  
22    def forward(self, x):  
23        x = self.feature_extractor(x)  
24        x = x.view(-1, 64*7*7) # Flatten image  
25        x = self.classifier(x)  
26        return x
```

## Pytorch: CNNs

Need to flatten image before passing it to our fully-connected layer (classifier)

```
1 class TwoLayerNet(nn.Module):
2     def __init__(self):
3         super().__init__()
4         I, J, C = 785, 64, 10
5         self.feature_extractor = nn.Sequential(
6             nn.Conv2d(in_channels=1,
7                     out_channels=32,
8                     kernel_size=3,
9                     padding=1),
10            nn.MaxPool2d(kernel_size=2, stride=2),
11            nn.Conv2d(in_channels=32,
12                     out_channels=64,
13                     kernel_size=3,
14                     padding=1),
15            nn.MaxPool2d(kernel_size=2, stride=2)
16        )
17        self.classifier = nn.Sequential(
18            nn.Linear(64*7*7, 64),
19            nn.ReLU(),
20            nn.Linear(64, 10)
21        )
22    def forward(self, x):
23        x = self.feature_extractor(x)
24        x = x.view(-1, 64*7*7) # Flatten image
25        x = self.classifier(x)
26        return x
```

# Pytorch: CNNs

```
28 learning_rate = 1e-3
29 batch_size=32
30 model = to_cuda(TwoLayerNet())
31
32 dataloader_train, dataloader_test = load_mnist(batch_size)
33
34 loss_function = torch.nn.CrossEntropyLoss()
35 optimizer = torch.optim.Adam(model.parameters(),
36                             lr=learning_rate)
37 losses = []
38 for epoch in range(2):
39     for (X_batch,Y_batch) in dataloader_train:
40         X_batch, Y_batch = to_cuda([X_batch, Y_batch])
41         # forward pass
42         y_k = model(X_batch)
43         # Compute loss
44         loss = loss_function(y_k, Y_batch)
45         losses.append(loss)
46         # Backpropagation
47         loss.backward()
48         optimizer.step()
49         optimizer.zero_grad()
50
```

```
1 class TwoLayerNet(nn.Module):
2     def __init__(self):
3         super().__init__()
4         I, J, C = 785, 64, 10
5         self.feature_extractor = nn.Sequential(
6             nn.Conv2d(in_channels=1,
7                     out_channels=32,
8                     kernel_size=3,
9                     padding=1),
10            nn.MaxPool2d(kernel_size=2, stride=2),
11            nn.Conv2d(in_channels=32,
12                     out_channels=64,
13                     kernel_size=3,
14                     padding=1),
15            nn.MaxPool2d(kernel_size=2, stride=2)
16        )
17        self.classifier = nn.Sequential(
18            nn.Linear(64*7*7, 64),
19            nn.ReLU(),
20            nn.Linear(64, 10)
21        )
22    def forward(self, x):
23        x = self.feature_extractor(x)
24        x = x.view(-1, 64*7*7) # Flatten image
25        x = self.classifier(x)
26        return x
```

# Pytorch: CNNs

Lets compute the accuracy!

```
1 num_correct = 0
2 total_samples = 0
3 for X_test, Y_test in dataloader_test:
4     output = model(X_test)
5     prediction = output.argmax(dim=1)
6     num_correct += (prediction == Y_test).sum().item()
7     total_samples += X_test.shape[0]
8 print("Accuracy:", num_correct/total_samples)
```

## Pytorch: CNNs

Lets compute the accuracy!

```
1 num_correct = 0
2 total_samples = 0
3 for X_test, Y_test in dataloader_test:
4     output = model(X_test)
5     prediction = output.argmax(dim=1)
6     num_correct += (prediction == Y_test).sum().item()
7     total_samples += X_test.shape[0]
8 print("Accuracy:", num_correct/total_samples)
```

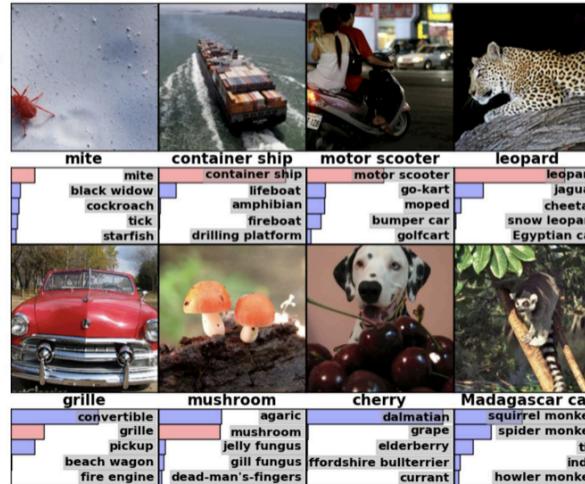
Accuracy: 0.9859

# Transfer Learning

- Learn features on ImageNet

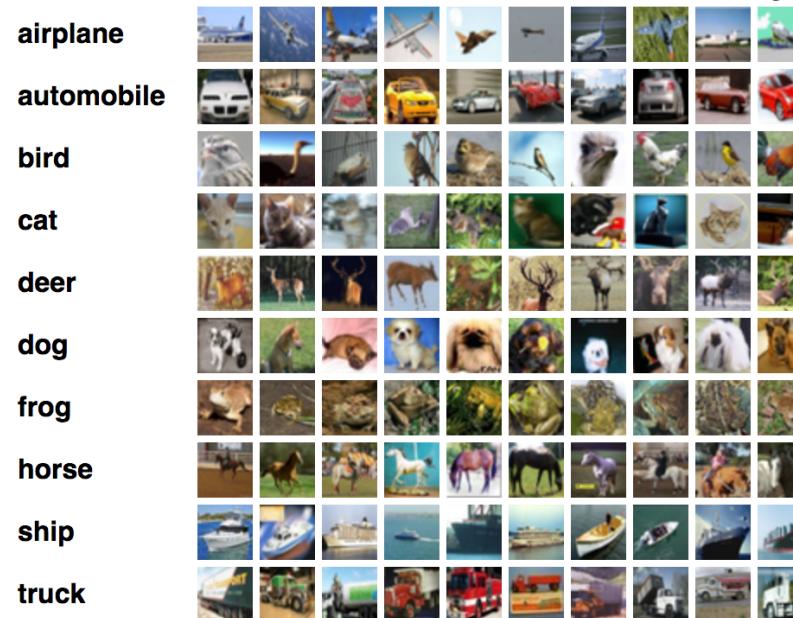
IMAGENET

- 1,000 object classes (categories).
- Images:
  - 1.2 M train
  - 100k test.

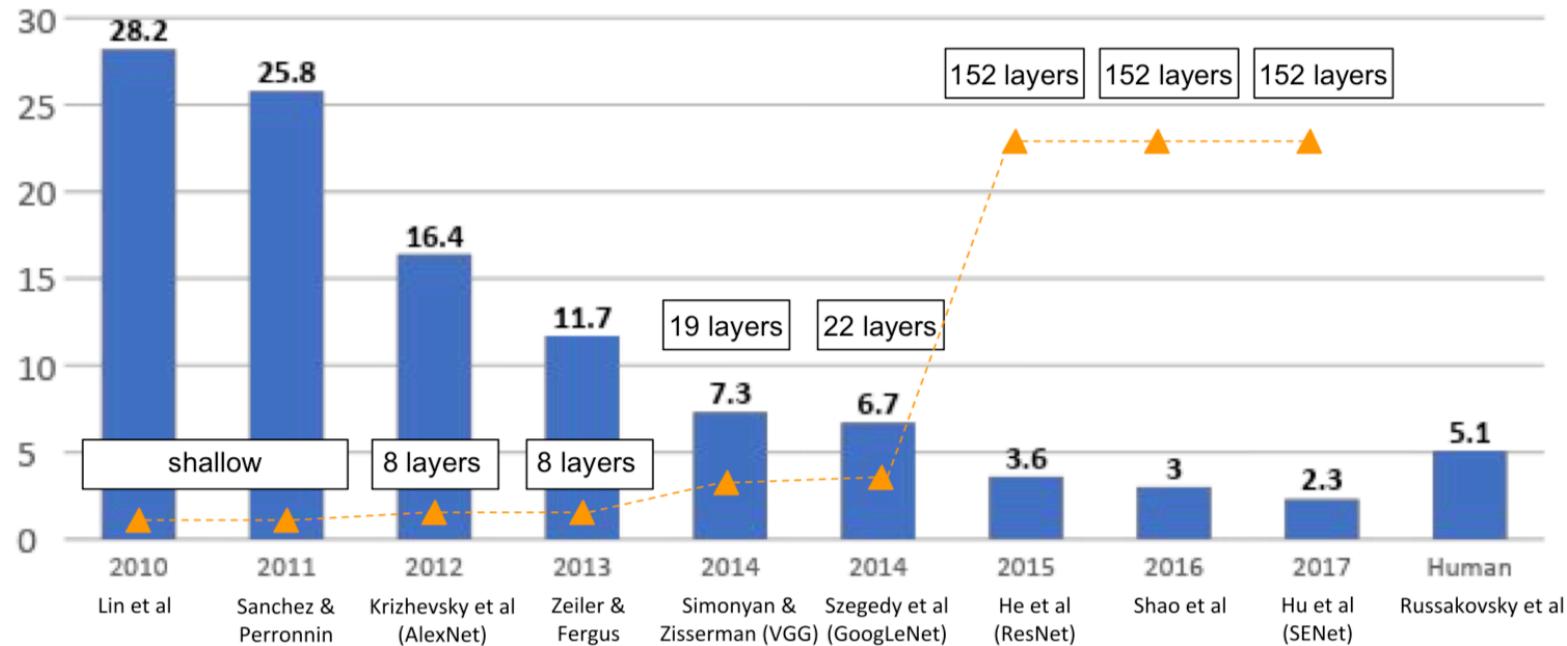


# Transfer Learning

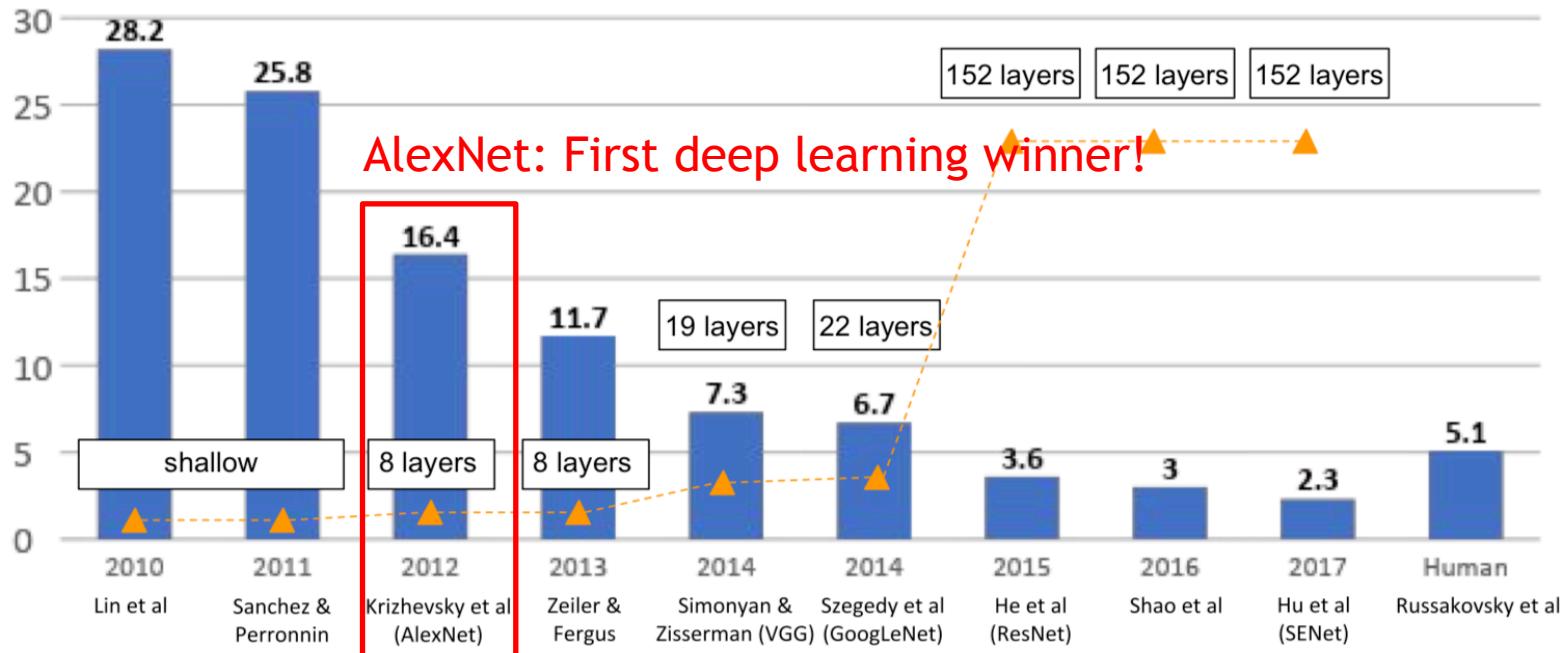
- Fine-tune on CIFAR10



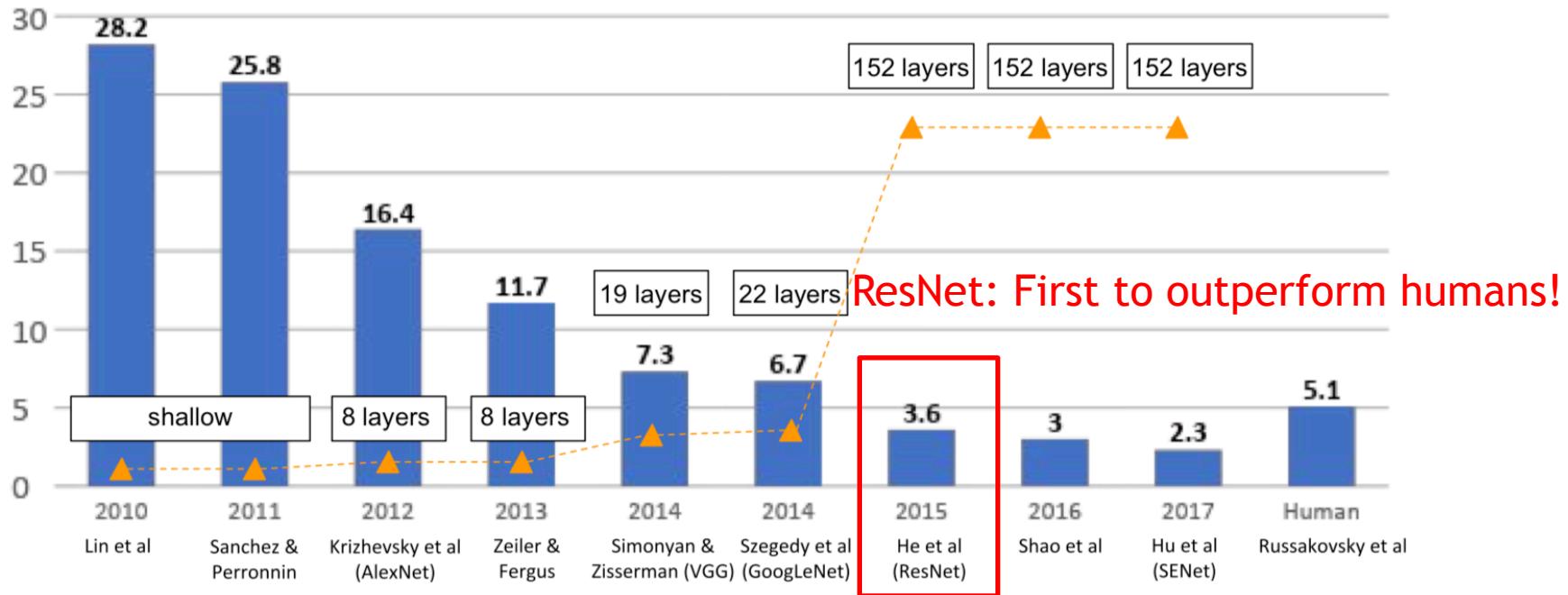
# ImageNet Large Scale Visual Recognition Challenge



# ImageNet Large Scale Visual Recognition Challenge

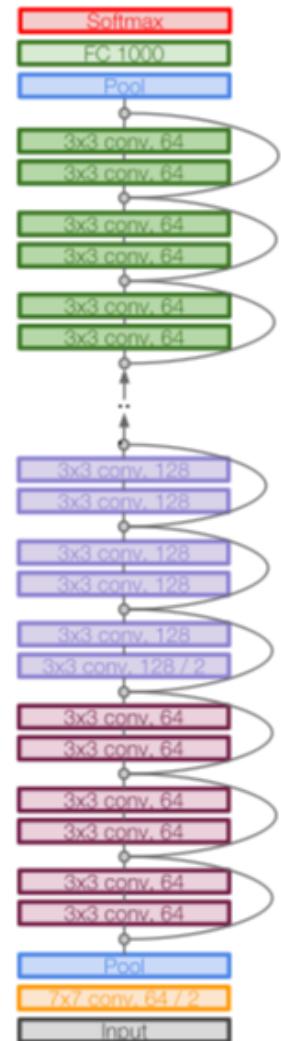
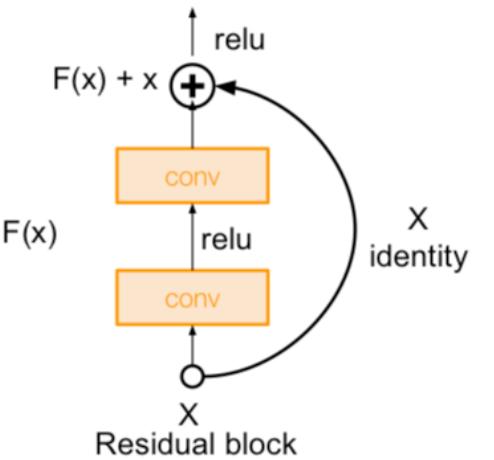


# ImageNet Large Scale Visual Recognition Challenge

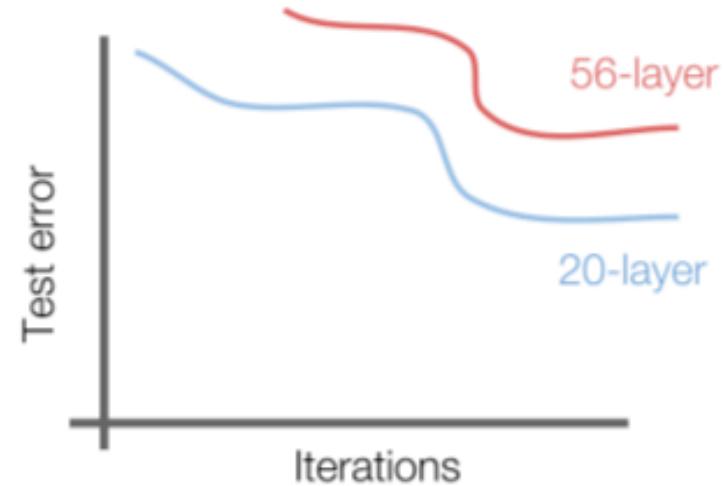
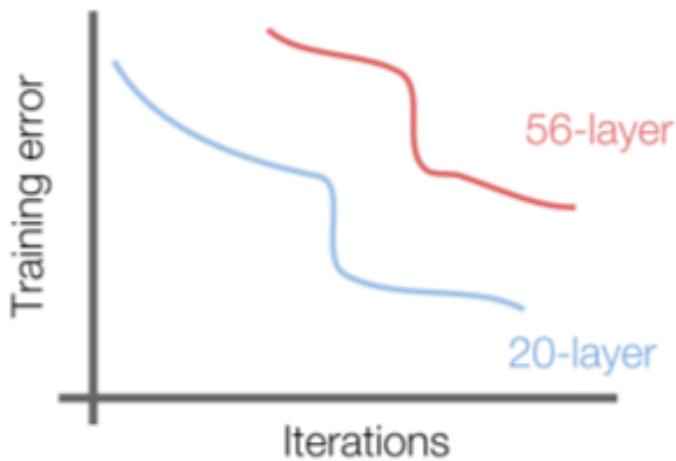


# ResNet

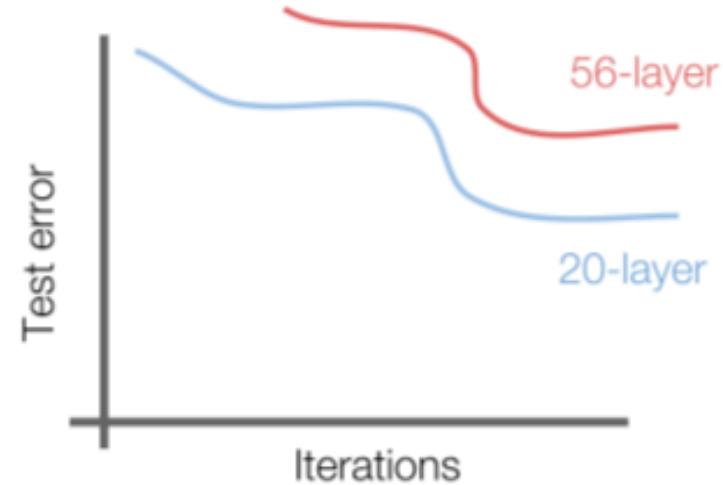
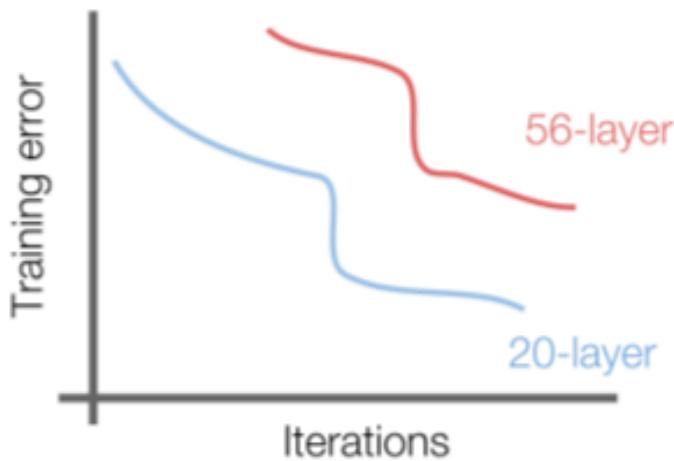
- Very deep Residual Network
  - Uses residual connections



# What happens without Residual Connections?



# What happens without Residual Connections?



Deeper -> Worse accuracy. Not caused by overfitting!

# What happens without Residual Connections?

# Why?

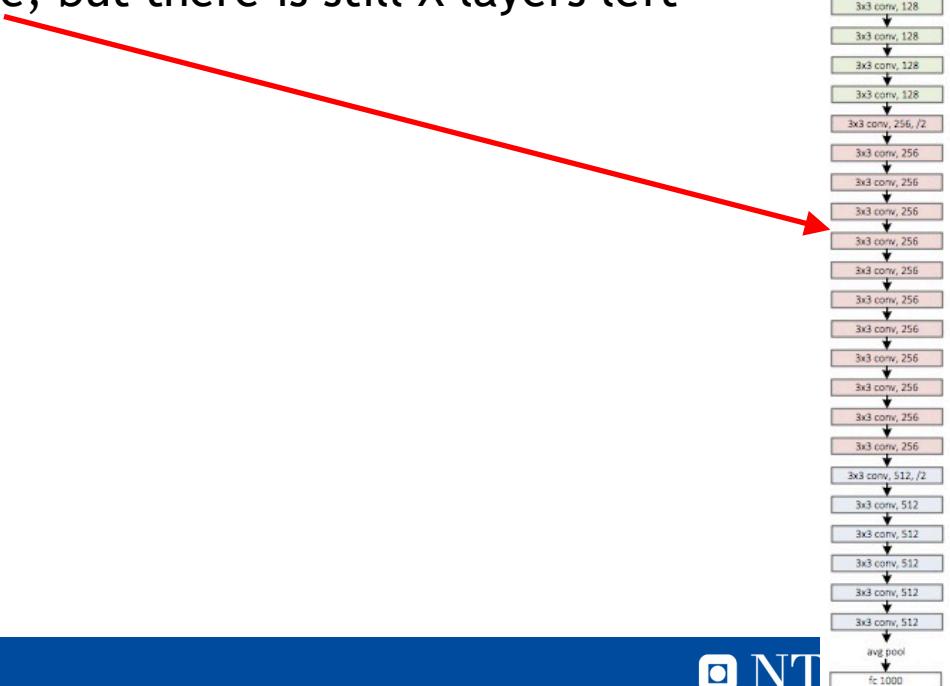
- The degradation problem



# What happens without Residual Connections?

Why?

- The degradation problem
  - We have the correct prediction here, but there is still X layers left



# What happens without Residual Connections?

Why?

- The degradation problem
  - We have the correct prediction here, but there is still X layers left
- Vanishing gradients
  - Backpropagation makes deltas smaller and smaller



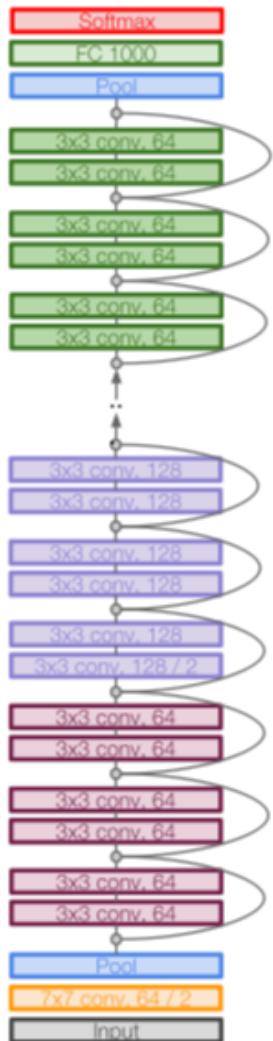
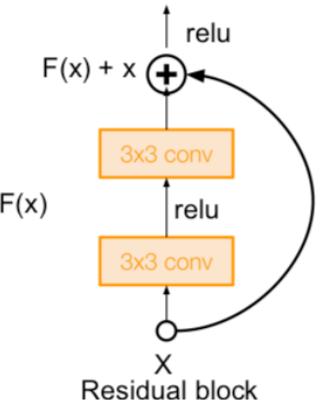
# What happens without Residual Connections?

Why?

- The degradation problem
  - We have the correct prediction here, but there is still X layers left
- Vanishing gradients
  - Backpropagation makes deltas smaller and smaller

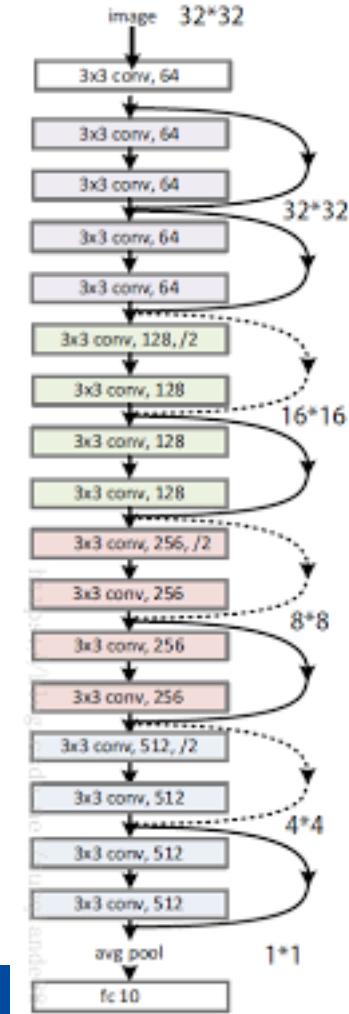
Solution:

- Add skip connections



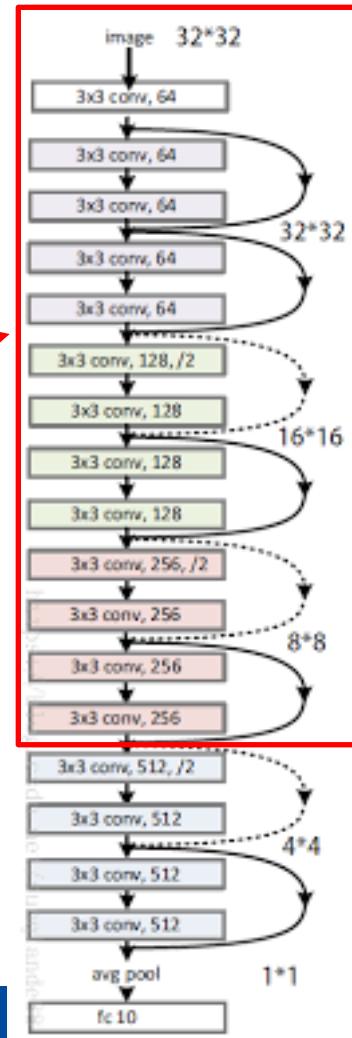
# What you will use:

- Resnet-18
- Pre-trained on ImageNet



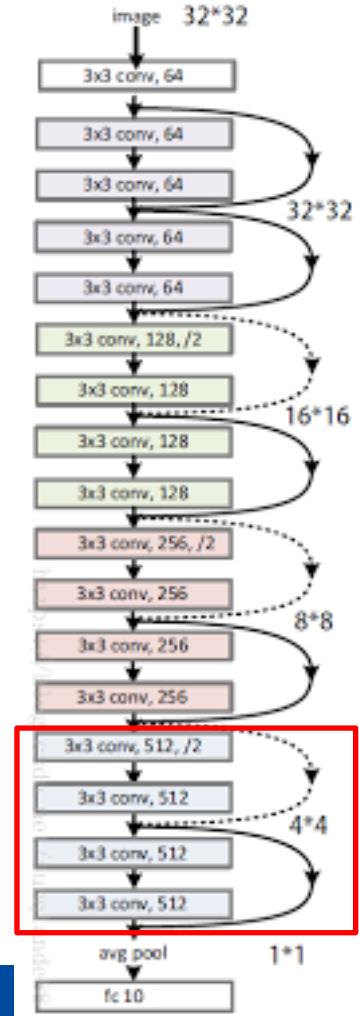
# What you will use:

- Resnet-18
- Pre-trained on ImageNet
- Freeze these layers
  - = Do not update weights



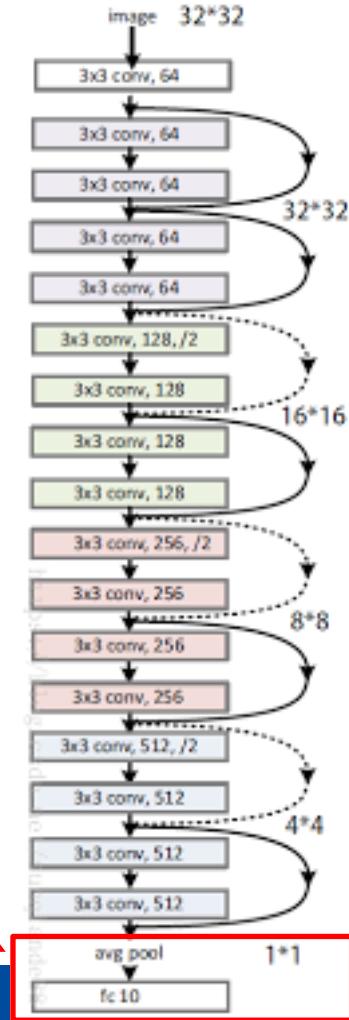
# What you will use:

- Resnet-18
- Pre-trained on ImageNet
- Unfreeze these
  - = only update these weights



# What you will use:

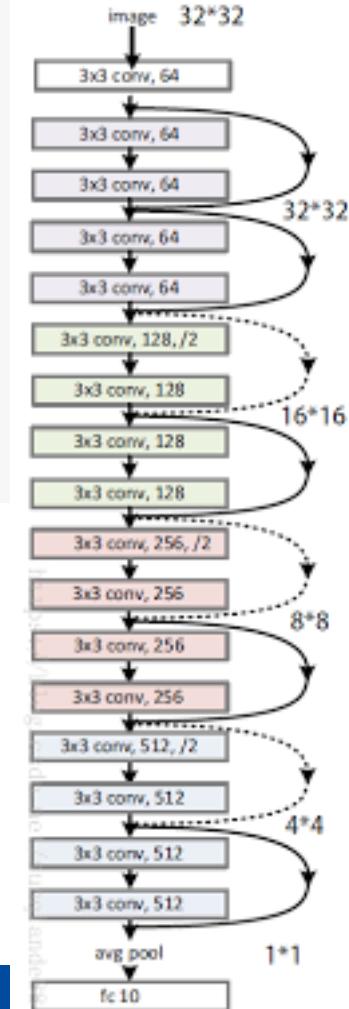
- Resnet-18
- Pre-trained on ImageNet
- Replace 1000-classes softmax  
with 10-classes softmax



# In PyTorch:

Load the pretrained model

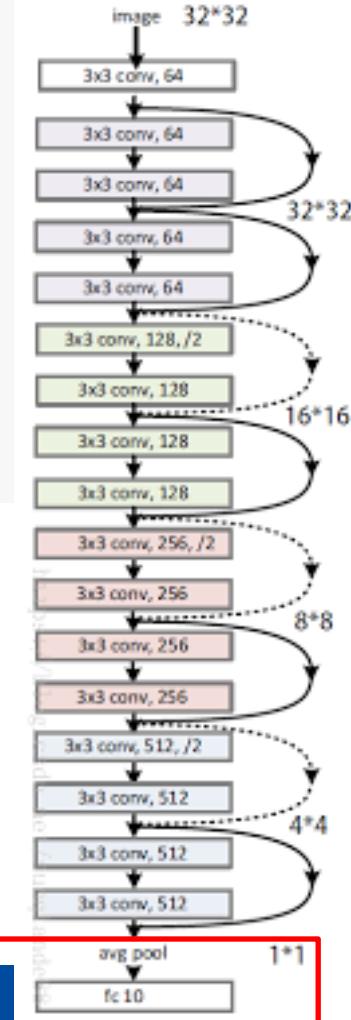
```
1 class Model(nn.Module):
2     def __init__(self):
3         super().__init__()
4         self.model = torchvision.models.resnet18(pretrained=True)
5         self.model.fc = nn.Sequential(
6             nn.Linear(2048, 10),
7         )
8         for param in self.model.parameters():
9             param.requires_grad = False
10        for param in self.model.fc.parameters():
11            param.requires_grad = True
12        for param in self.model.layer4.parameters():
13            param.requires_grad = True
14
15    def forward(self, x):
16        x = nn.functional.interpolate(x, scale_factor=8)
17        x = self.model(x)
18        return x
```



# In PyTorch:

Replace output layer  
with 10 classes

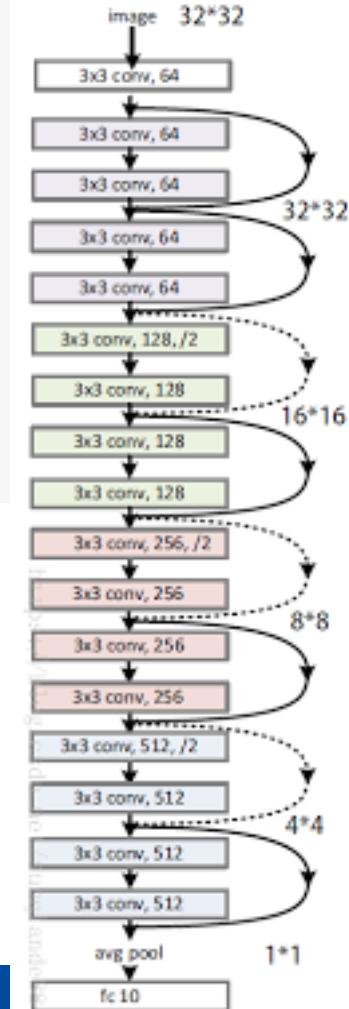
```
1 class Model(nn.Module):
2     def __init__(self):
3         super().__init__()
4         self.model = torchvision.models.resnet18(pretrained=True)
5         self.model.fc = nn.Sequential(
6             nn.Linear(2048, 10),
7         )
8         for param in self.model.parameters():
9             param.requires_grad = False
10        for param in self.model.fc.parameters():
11            param.requires_grad = True
12        for param in self.model.layer4.parameters():
13            param.requires_grad = True
14
15    def forward(self, x):
16        x = nn.functional.interpolate(x, scale_factor=8)
17        x = self.model(x)
18        return x
```



# In PyTorch:

Freeze all layers in model

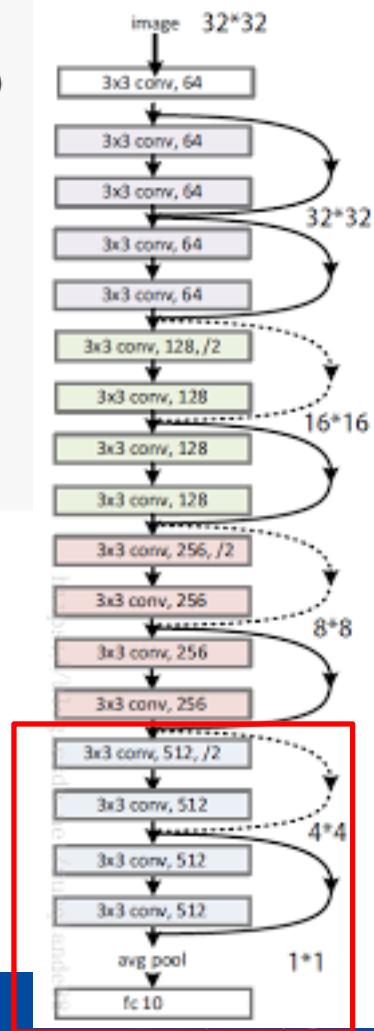
```
1 class Model(nn.Module):
2     def __init__(self):
3         super().__init__()
4         self.model = torchvision.models.resnet18(pretrained=True)
5         self.model.fc = nn.Sequential(
6             nn.Linear(2048, 10),
7         )
8         for param in self.model.parameters():
9             param.requires_grad = False
10        for param in self.model.fc.parameters():
11            param.requires_grad = True
12        for param in self.model.layer4.parameters():
13            param.requires_grad = True
14
15    def forward(self, x):
16        x = nn.functional.interpolate(x, scale_factor=8)
17        x = self.model(x)
18        return x
```



# In PyTorch:

## Unfreeze last layers

```
1 class Model(nn.Module):
2     def __init__(self):
3         super().__init__()
4         self.model = torchvision.models.resnet18(pretrained=True)
5         self.model.fc = nn.Sequential(
6             nn.Linear(2048, 10),
7         )
8         for param in self.model.parameters():
9             param.requires_grad = False
10        for param in self.model.fc.parameters():
11            param.requires_grad = True
12        for param in self.model.layer4.parameters():
13            param.requires_grad = True
14
15    def forward(self, x):
16        x = nn.functional.interpolate(x, scale_factor=8)
17        x = self.model(x)
18        return x
```



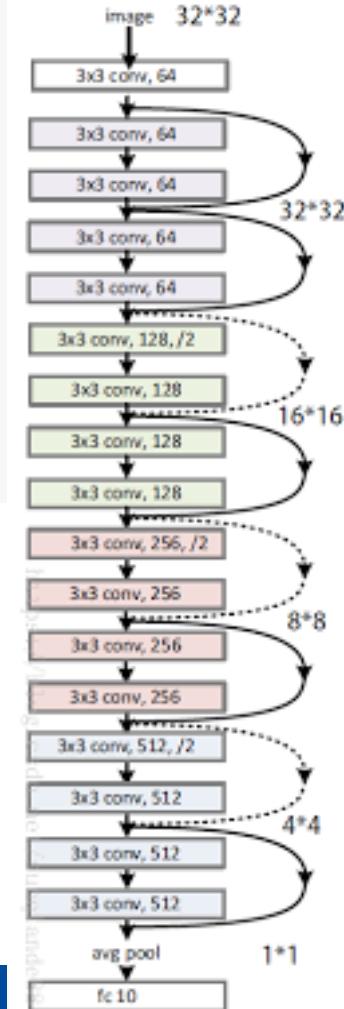
# In PyTorch:

ImageNet is trained on  
256x256

CIFAR-10 is 32x32.

We upsample the image to  
256x256

```
1 class Model(nn.Module):
2     def __init__(self):
3         super().__init__()
4         self.model = torchvision.models.resnet18(pretrained=True)
5         self.model.fc = nn.Sequential(
6             nn.Linear(2048, 10),
7         )
8         for param in self.model.parameters():
9             param.requires_grad = False
10        for param in self.model.fc.parameters():
11            param.requires_grad = True
12        for param in self.model.layer4.parameters():
13            param.requires_grad = True
14
15    def forward(self, x):
16        x = nn.functional.interpolate(x, scale factor=8)
17        x = self.model(x)
18        return x
```

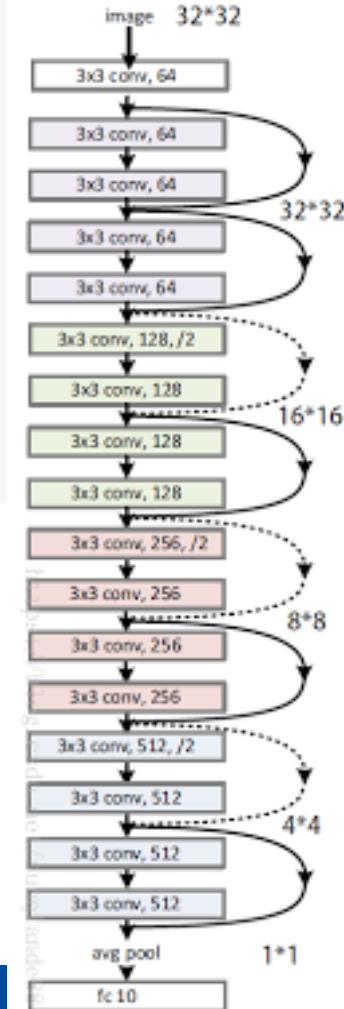


# In PyTorch:

Then:

Train of CIFAR10 until  
early stopping

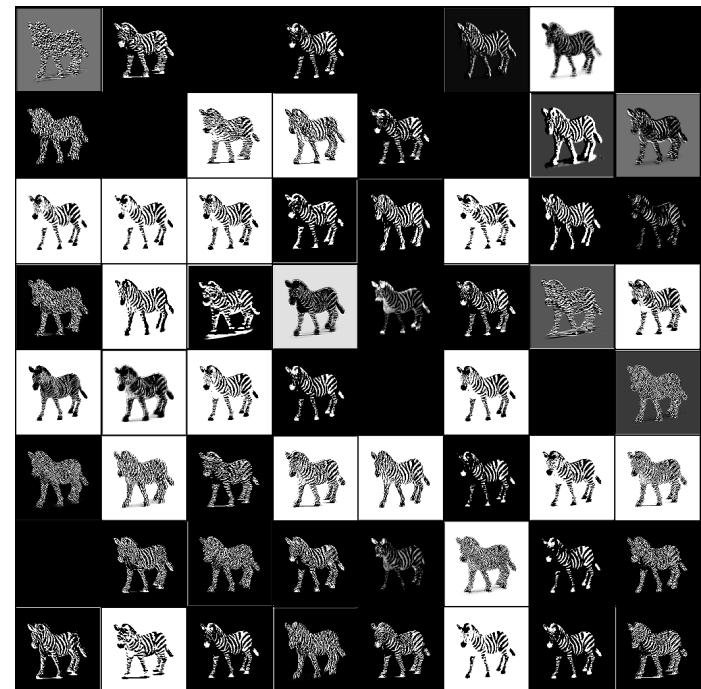
```
1 class Model(nn.Module):
2     def __init__(self):
3         super().__init__()
4         self.model = torchvision.models.resnet18(pretrained=True)
5         self.model.fc = nn.Sequential(
6             nn.Linear(2048, 10),
7         )
8         for param in self.model.parameters():
9             param.requires_grad = False
10        for param in self.model.fc.parameters():
11            param.requires_grad = True
12        for param in self.model.layer4.parameters():
13            param.requires_grad = True
14
15    def forward(self, x):
16        x = nn.functional.interpolate(x, scale_factor=8)
17        x = self.model(x)
18        return x
```



# Interpreting CNNs by visualizations



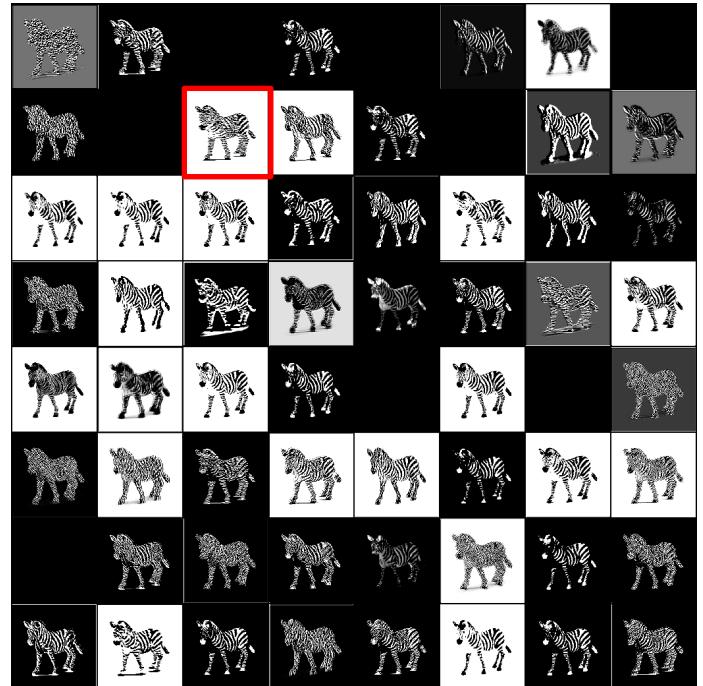
# Visualization: First Layer



# Visualization: First Layer



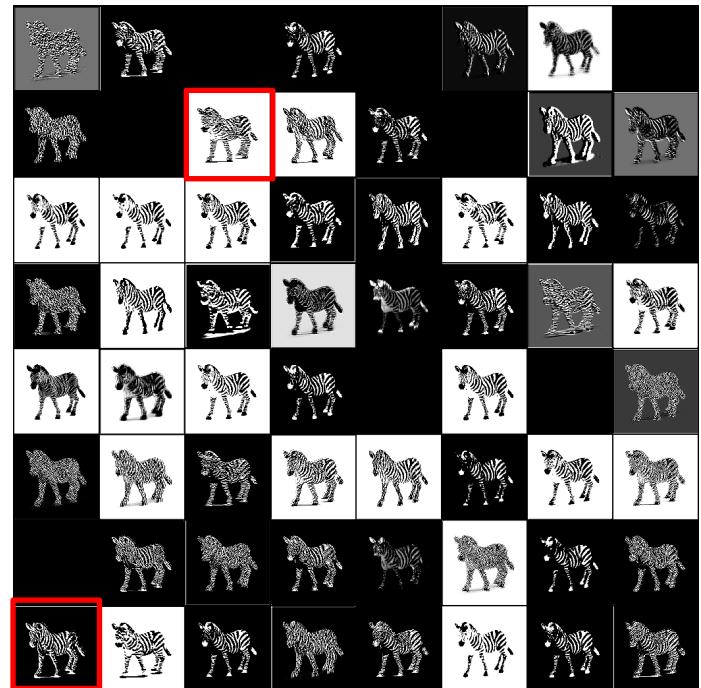
Horizontal filter?



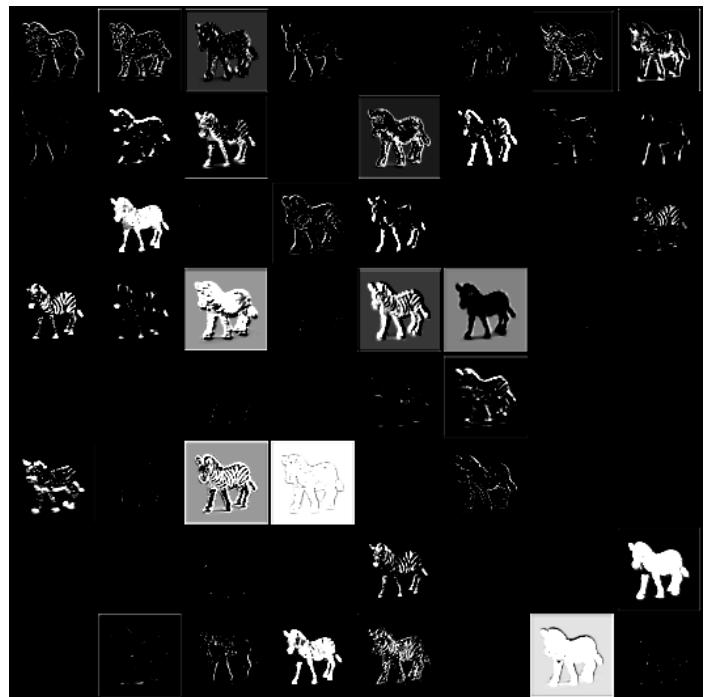
# Visualization: First Layer



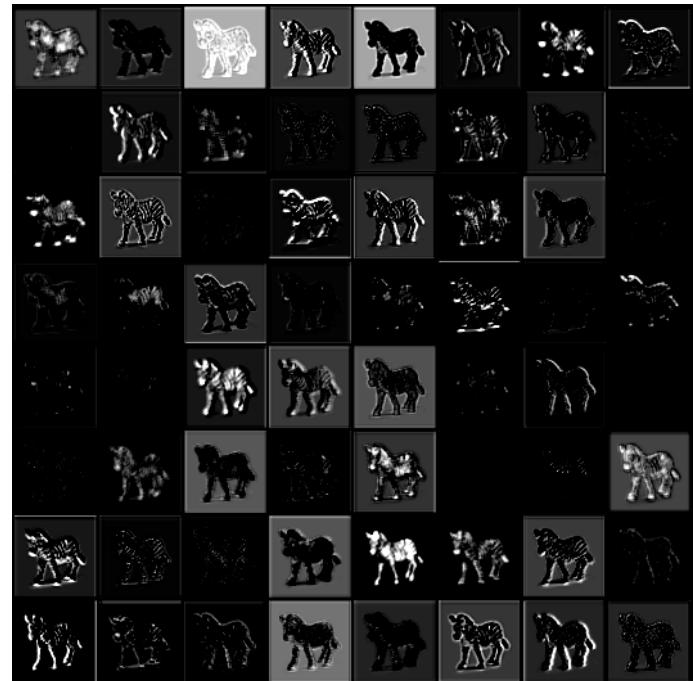
Vertical filter?



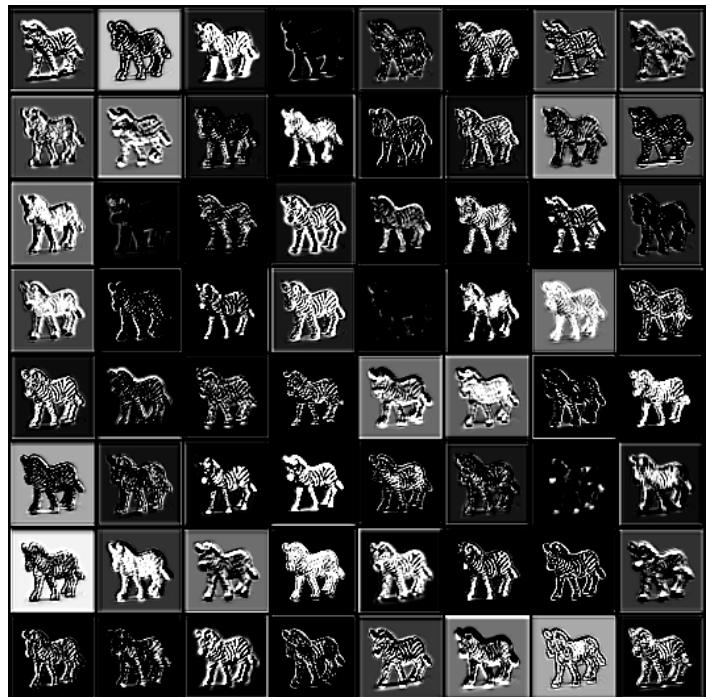
# Visualization: Second Layer



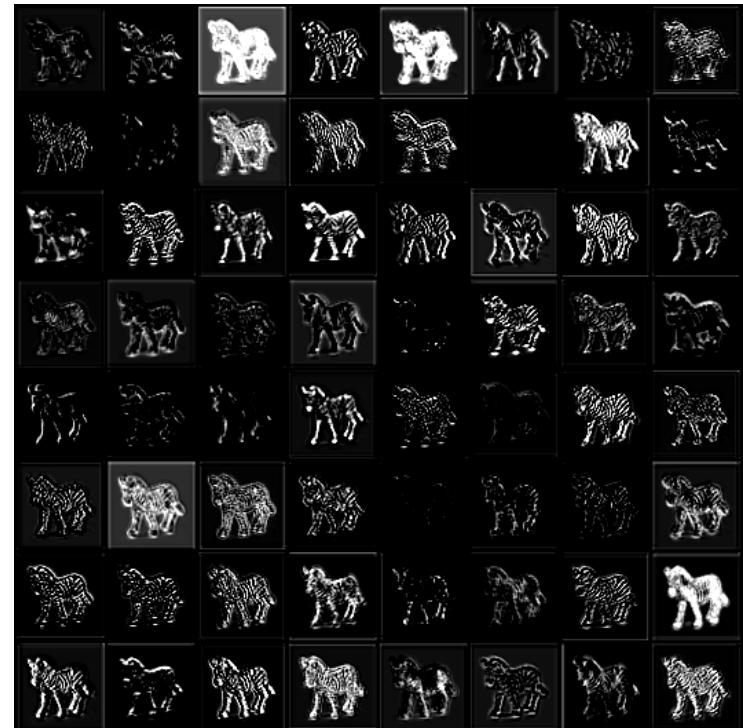
# Visualization: Third Layer



# Visualization: Fourth Layer



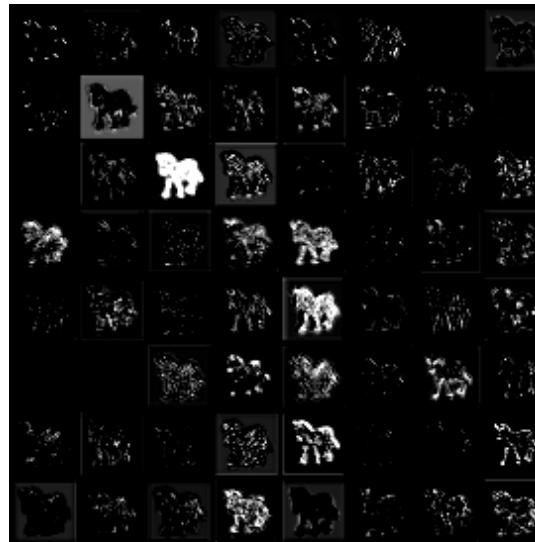
# Visualization: Fifth Layer



# Visualization: Sixth Layer



# Visualization: Seventh Layer



# Visualization: Last Layer



# In PyTorch:

```
1 from dataloaders import mean, std
2 import torchvision
3 from torchvision.transforms.functional \
4     import to_pil_image, to_tensor, normalize
5
6 image = plt.imread("test_img5.jpg")
7 image = to_tensor(image)
8 image = normalize(image.data, mean, std)
9 image = image.view(1, *image.shape)
10 image = nn.functional.interpolate(image, size=(256, 256))
11
12 model = torchvision.models.resnet18(pretrained=True)
13 first_layer_out = model.conv1(image)
14
15 to_visualize = first_layer_out.view(first_layer_out.shape[1],
16                                     1, *first_layer_out.shape[2:])
17 torchvision.utils.save_image(to_visualize,
18                               "filters_first_layer.png")
19
```

## In PyTorch:

Read the data and normalize similar as we do in dataloader

```
1 from dataloaders import mean, std
2 import torchvision
3 from torchvision.transforms.functional \
4     import to_pil_image, to_tensor, normalize
5
6 image = plt.imread("test_img5.jpg")
7 image = to_tensor(image)
8 image = normalize(image.data, mean, std)
9 image = image.view(1, *image.shape)
10 image = nn.functional.interpolate(image, size=(256, 256))
11
12 model = torchvision.models.resnet18(pretrained=True)
13 first_layer_out = model.conv1(image)
14
15 to_visualize = first_layer_out.view(first_layer_out.shape[1],
16                                     1, *first_layer_out.shape[2:])
17 torchvision.utils.save_image(to_visualize,
18                               "filters_first_layer.png")
19
```

# In PyTorch:

Resize to the expected shape

```
1 from dataloaders import mean, std
2 import torchvision
3 from torchvision.transforms.functional \
4     import to_pil_image, to_tensor, normalize
5
6 image = plt.imread("test_img5.jpg")
7 image = to_tensor(image)
8 image = normalize(image.data, mean, std)
9 image = image.view(1, *image.shape)
10 image = nn.functional.interpolate(image, size=(256, 256))
11
12 model = torchvision.models.resnet18(pretrained=True)
13 first_layer_out = model.conv1(image)
14
15 to_visualize = first_layer_out.view(first_layer_out.shape[1],
16                                     1, *first_layer_out.shape[2:])
17 torchvision.utils.save_image(to_visualize,
18                               "filters_first_layer.png")
19
```

# In PyTorch:

Load pretrained model

Forward pass through **ONLY** the first layer

```
1 from dataloaders import mean, std
2 import torchvision
3 from torchvision.transforms.functional \
4     import to_pil_image, to_tensor, normalize
5
6 image = plt.imread("test_img5.jpg")
7 image = to_tensor(image)
8 image = normalize(image.data, mean, std)
9 image = image.view(1, *image.shape)
10 image = nn.functional.interpolate(image, size=(256, 256))
11
12 model = torchvision.models.resnet18(pretrained=True)
13 first_layer_out = model.conv1(image)
14
15 to_visualize = first_layer_out.view(first_layer_out.shape[1],
16                                     1, *first_layer_out.shape[2:])
17 torchvision.utils.save_image(to_visualize,
18                               "filters_first_layer.png")
19
```

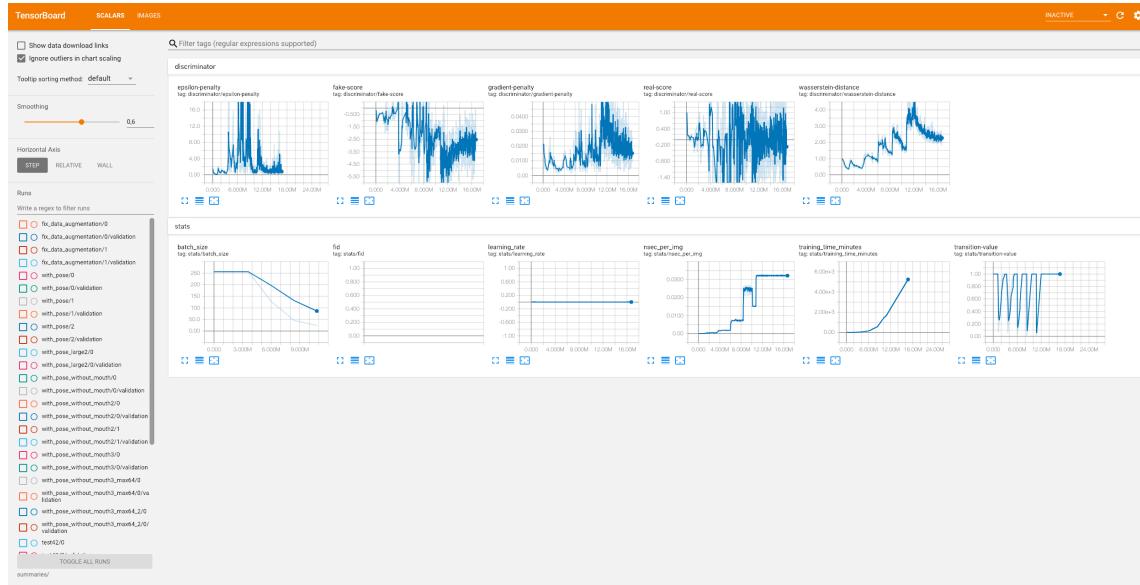
## In PyTorch:

Reshape image to expected such that num\_filters swap place with batch size and save image.

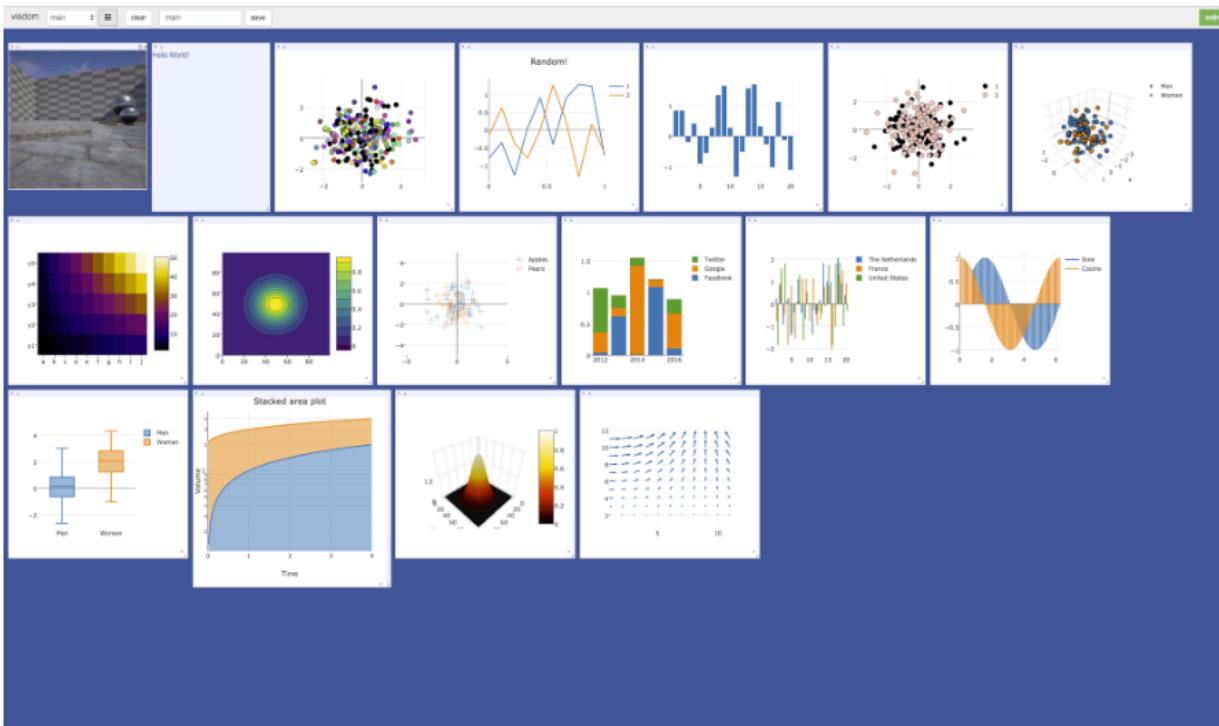
```
1 from dataloaders import mean, std
2 import torchvision
3 from torchvision.transforms.functional \
4     import to_pil_image, to_tensor, normalize
5
6 image = plt.imread("test_img5.jpg")
7 image = to_tensor(image)
8 image = normalize(image.data, mean, std)
9 image = image.view(1, *image.shape)
10 image = nn.functional.interpolate(image, size=(256, 256))
11
12 model = torchvision.models.resnet18(pretrained=True)
13 first_layer_out = model.conv1(image)
14
15 to_visualize = first_layer_out.view(first_layer_out.shape[1],
16                                     1, *first_layer_out.shape[2:])
17 torchvision.utils.save_image(to_visualize,
18                               "filters_first_layer.png")
19
```

# Tensorboard + TensorboardX(Pytorch)

My current Tensorboard:



# Pytorch Visdom



# Google Cloud Console

- [console.cloud.google.com](https://console.cloud.google.com)
-

# Recommended resources

- Stanford CS231n:
  - [http://cs231n.stanford.edu/slides/2018/cs231n\\_2018\\_lecture09.pdf](http://cs231n.stanford.edu/slides/2018/cs231n_2018_lecture09.pdf)
  - [http://cs231n.stanford.edu/slides/2018/cs231n\\_2018\\_lecture08.pdf](http://cs231n.stanford.edu/slides/2018/cs231n_2018_lecture08.pdf)
- <https://aifiddle.io>
- [CNN Cheat Sheet](#)