

# Assignment Lecture 3

Convolutional Neural Networks

Transfer Learning

Deep learning frameworks

Neural Networks on GPU

# Computer Hardware



# We have two hardware choices:

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

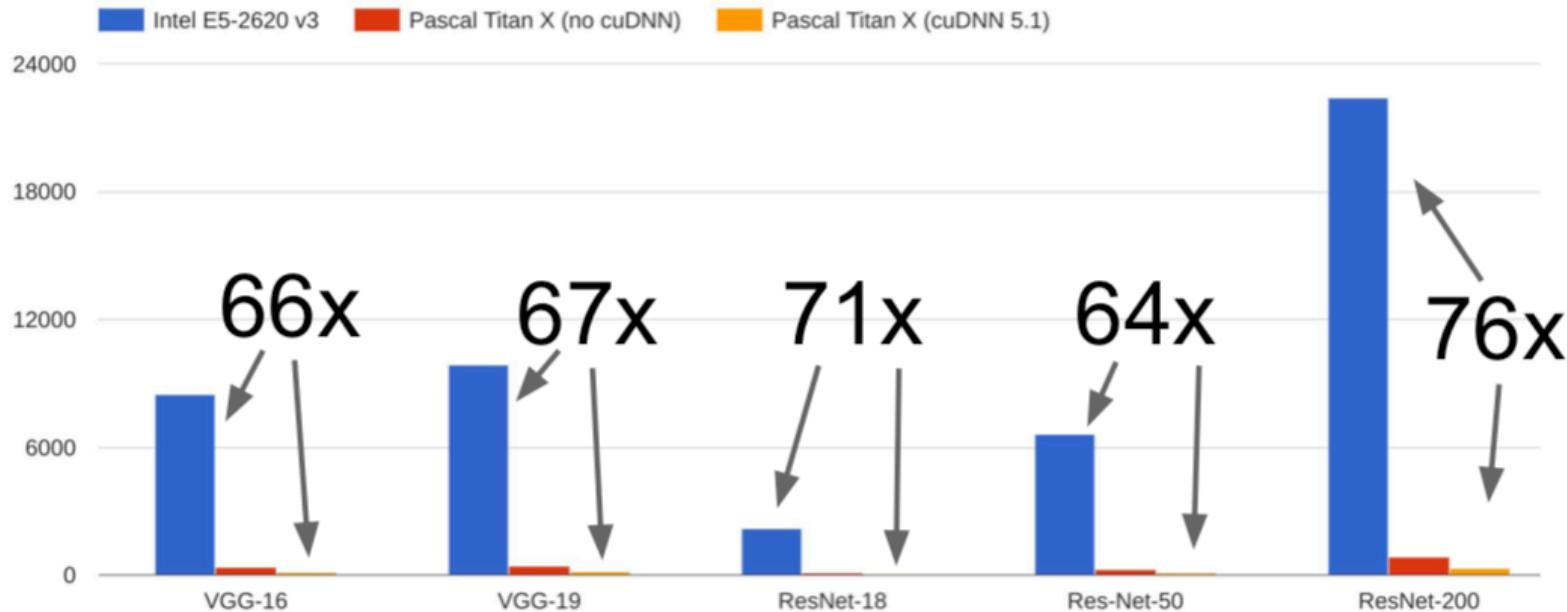
# GPU vs CPU

	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

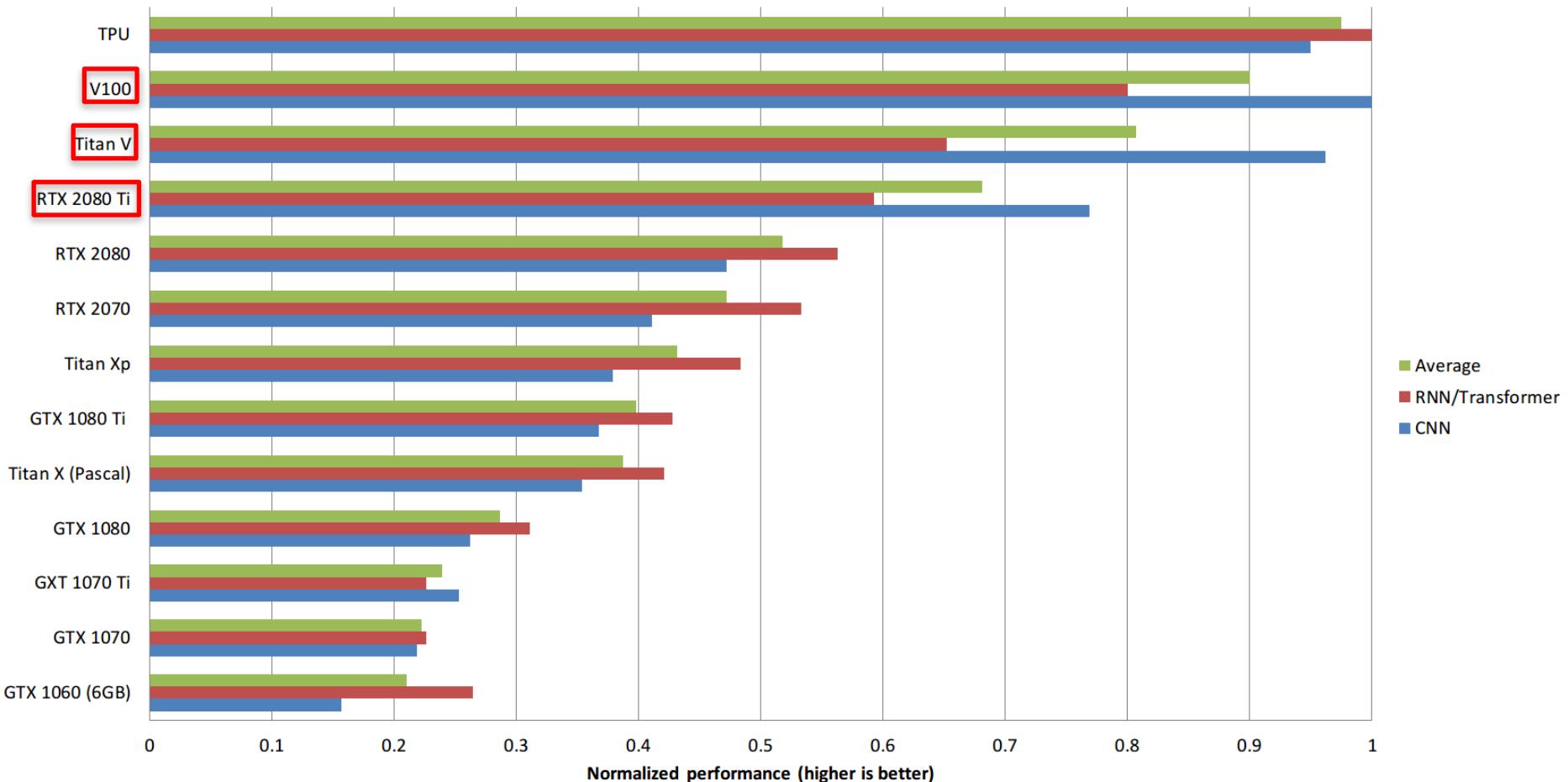
# GPU vs CPU for CNNs



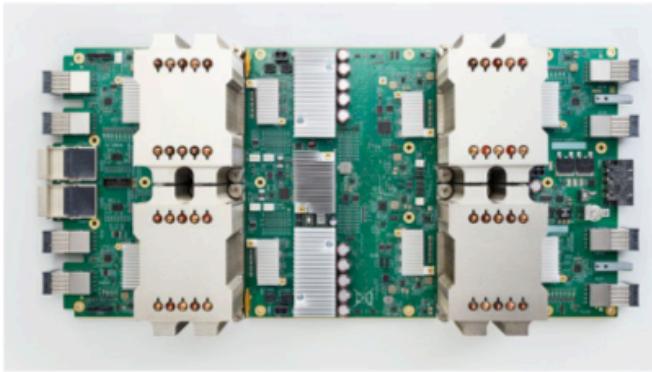
# GPU vs CPU: ResNet-200

- 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**)

# Performance



# Tensor Processing Units (TPU)



Google Cloud TPU  
= 180 TFLOPs of compute!



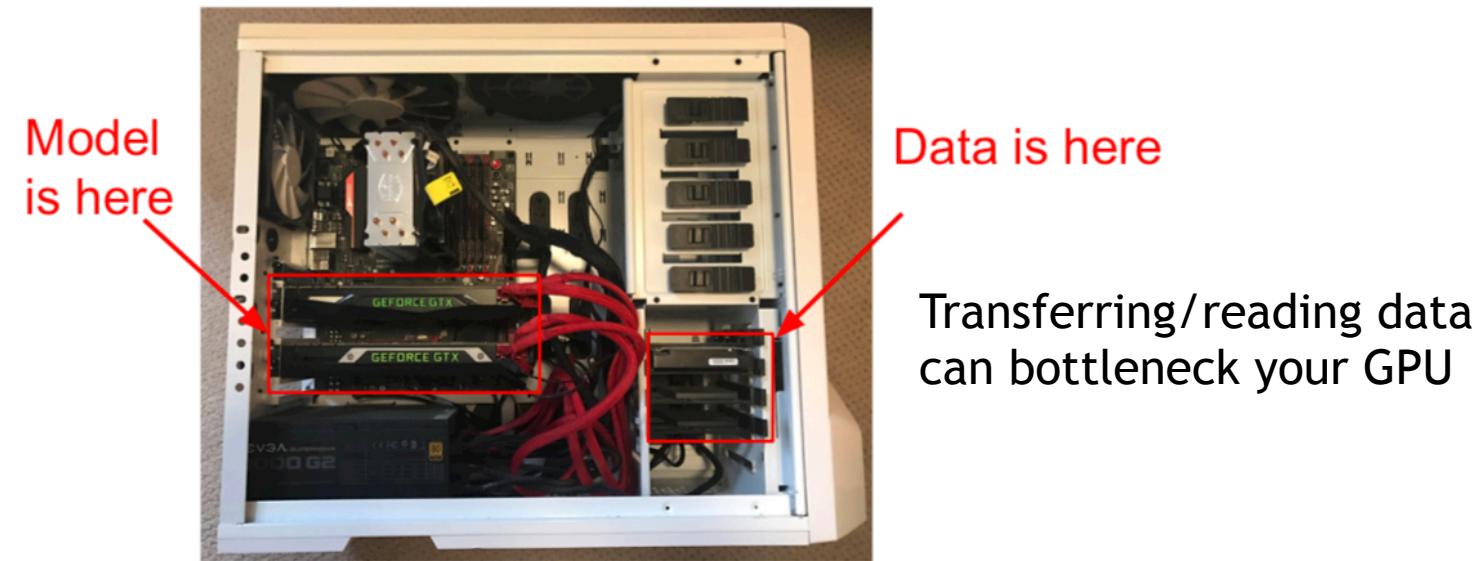
NVIDIA Tesla V100  
= 125 TFLOPs of compute

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

# Issue: CPU <-> GPU Communication



# Issue: CPU <-> GPU Communication



# Issue: CPU <-> GPU Communication



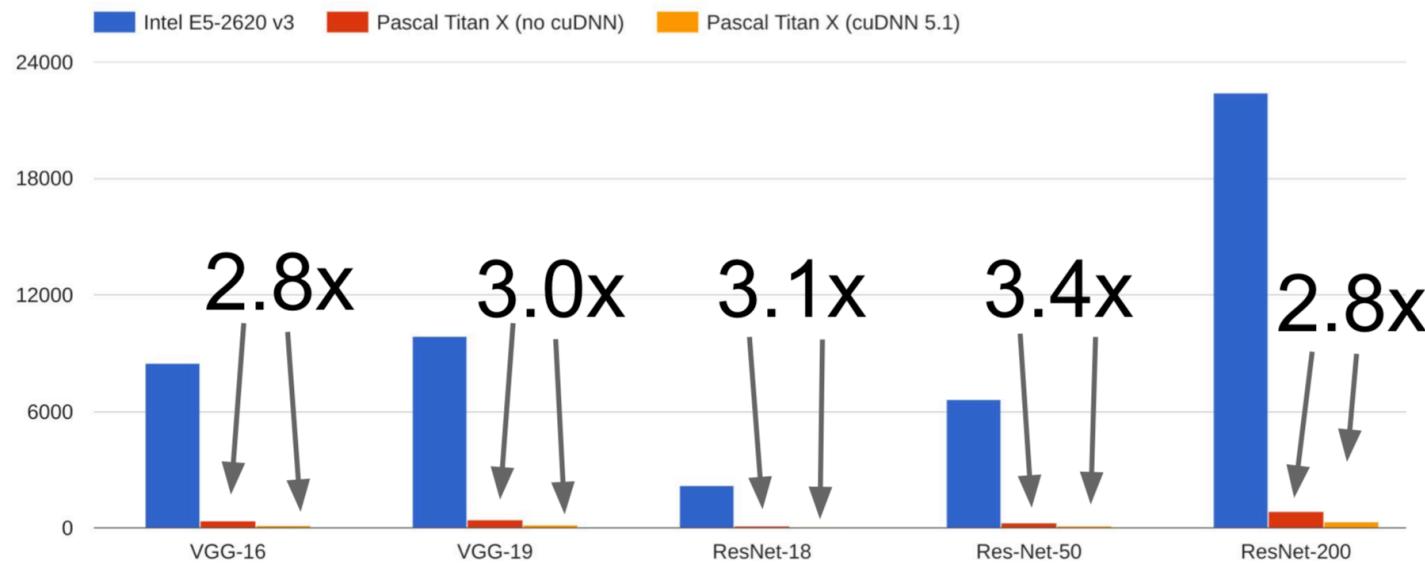
Solution:

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

# Deep Learning Software

- Deep learning software runs on (NVIDIA):

- CUDA
- cudNN



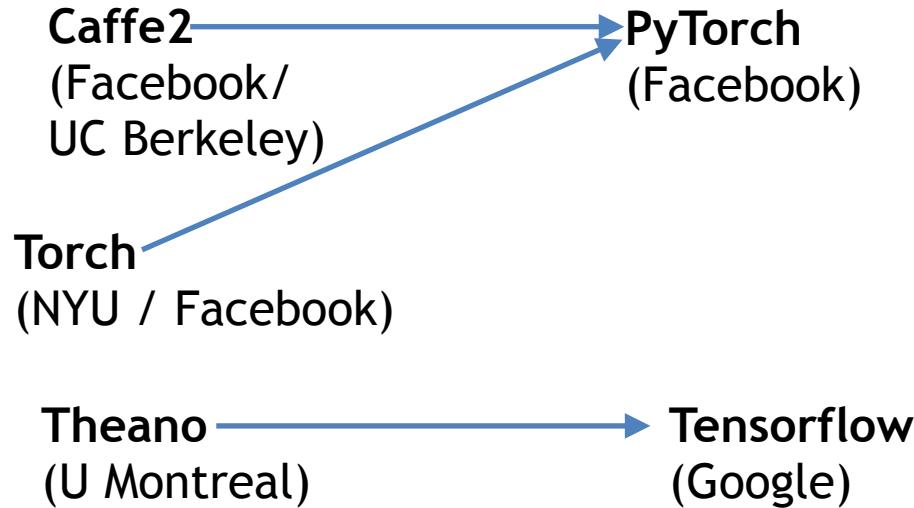
# Frameworks

**Caffe2**  
(Facebook/  
UC Berkeley)

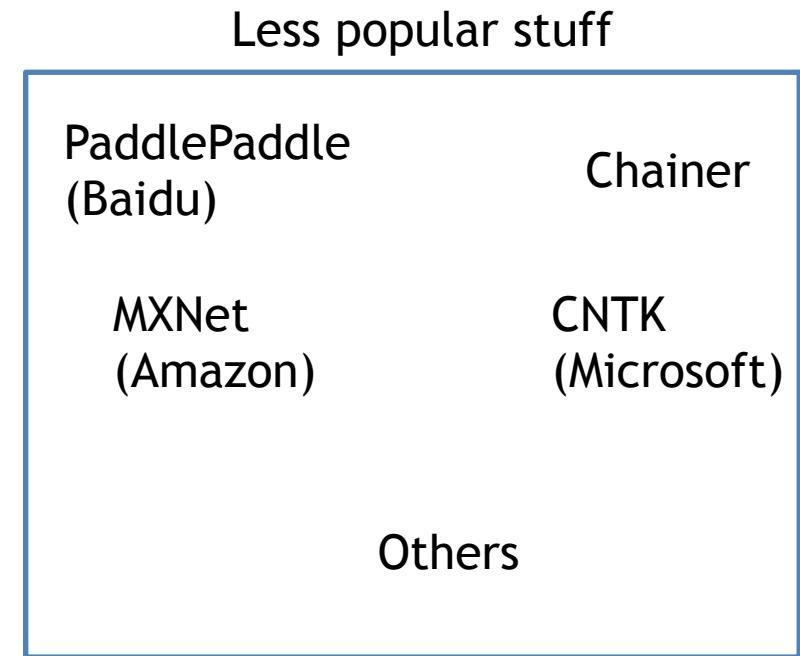
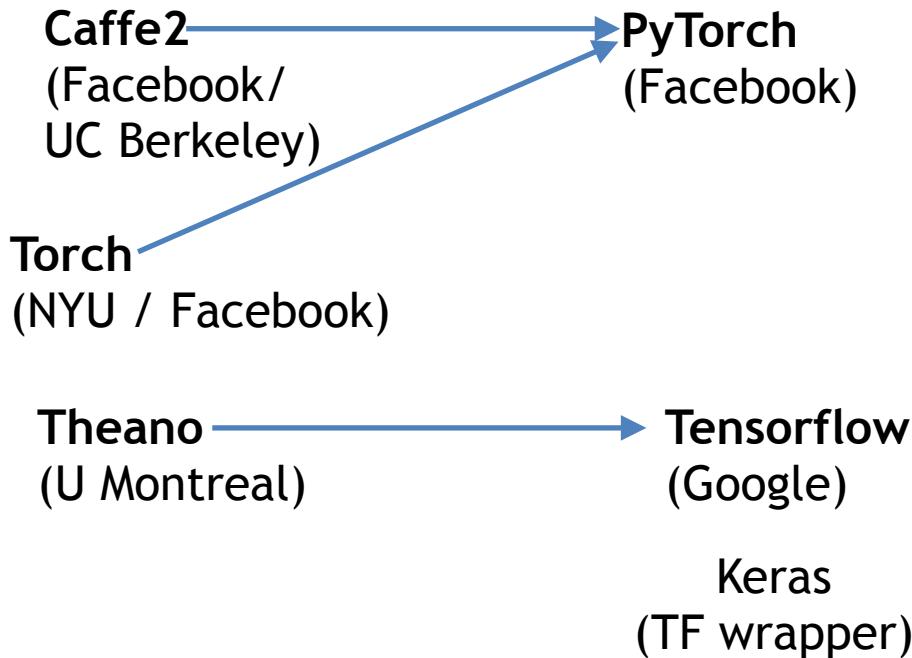
**Torch**  
(NYU / Facebook)

**Theano**  
(U Montreal)

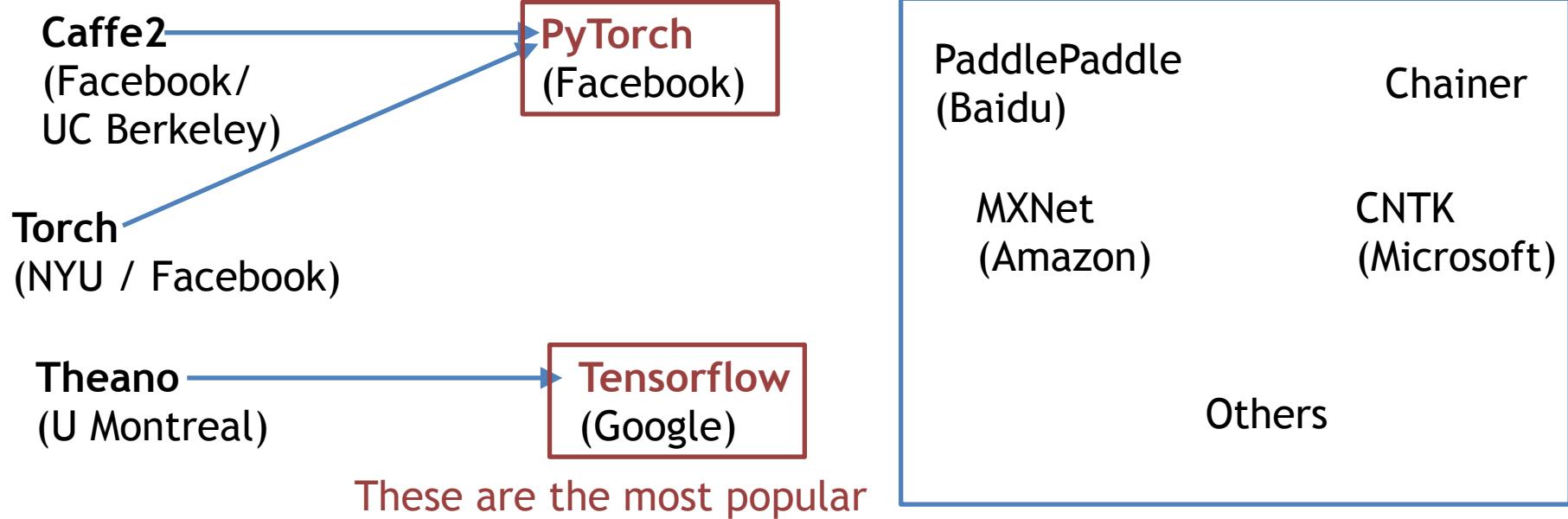
# Frameworks



# Frameworks



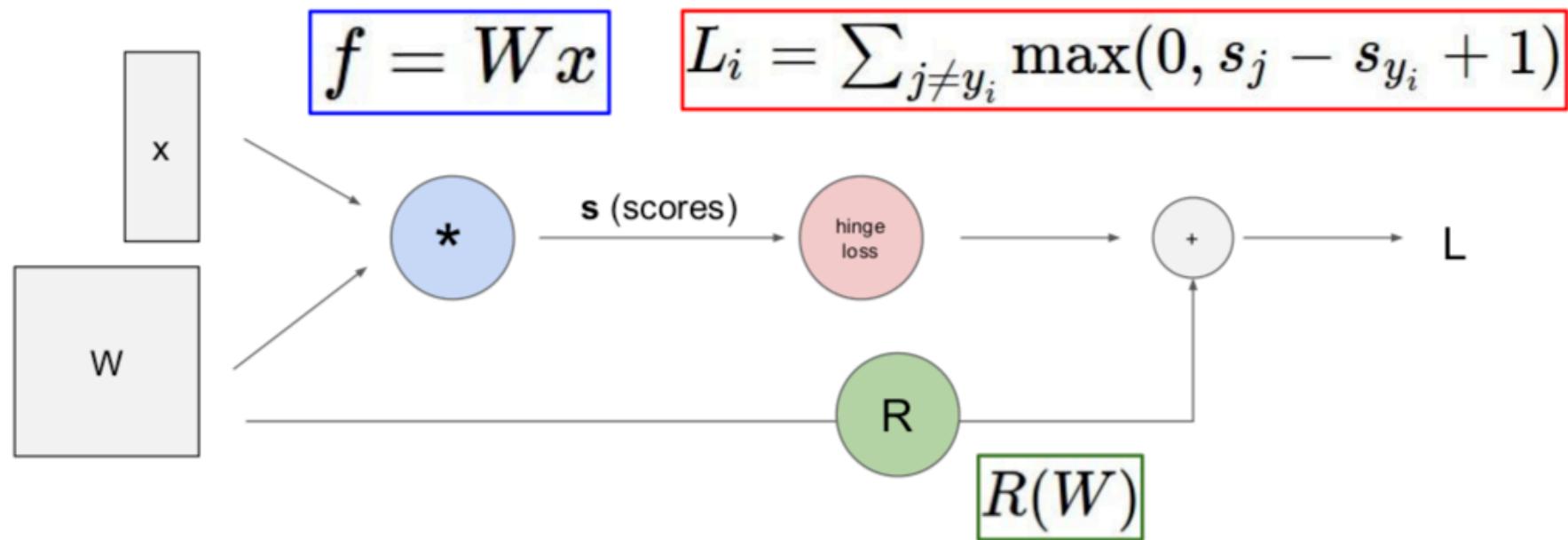
# Frameworks



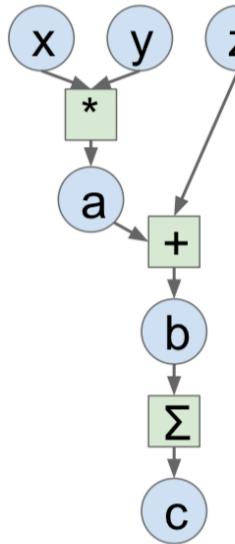
# Requirement of deep learning frameworks

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

# Neural Networks = Directed Acyclic Graphs



# Computational Graph



# Computational Graph

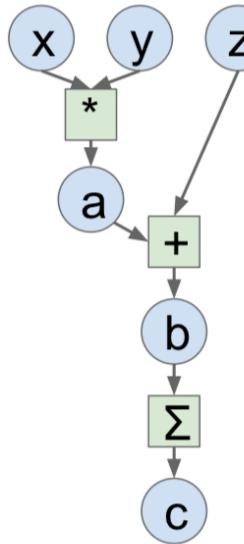
Numpy

```
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

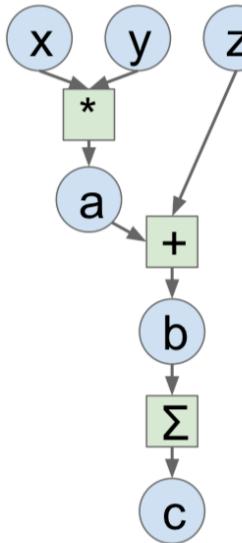
```
import numpy as np
np.random.seed(0)

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

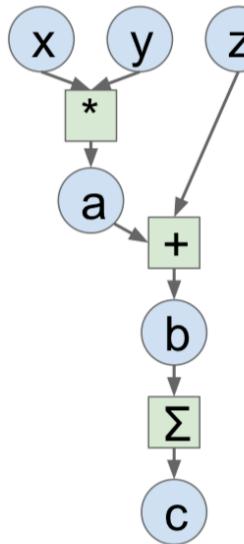
```
import numpy as np
np.random.seed(0)

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

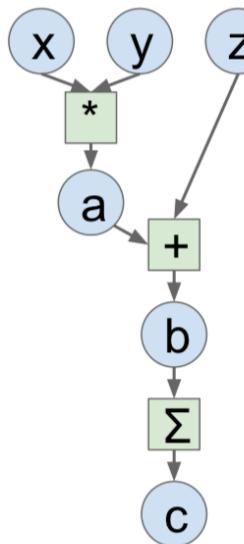
```
import numpy as np
np.random.seed(0)

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

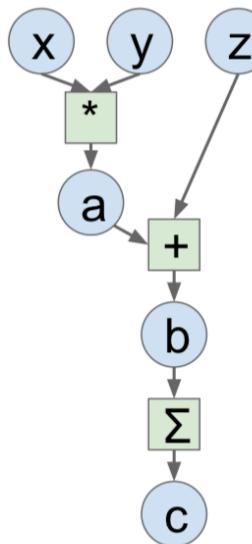
```
import numpy as np
np.random.seed(0)

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



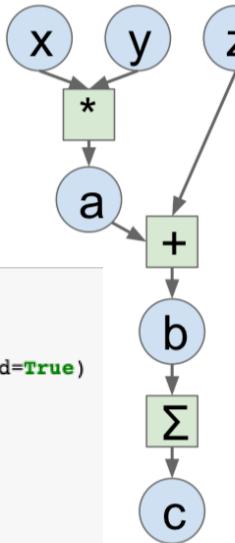
```
W_h = torch.randn(20, 20, requires_grad=True)
W_x = torch.randn(20, 10, requires_grad=True)
x = torch.randn(1, 10)
prev_h = torch.randn(1, 20)
```



# 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)
```

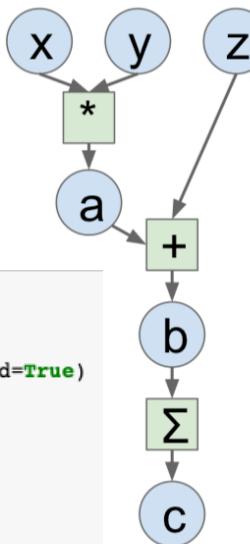


# Computational Graph

Tensorflow 1.2

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)
```



```
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])
```

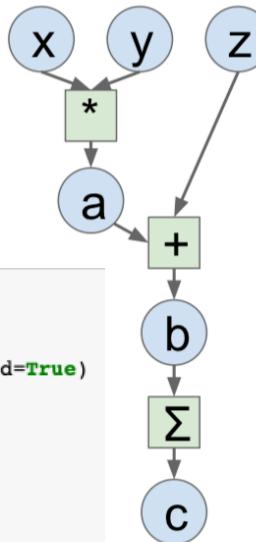
First, define the graph

# Computational Graph

Tensorflow 1.2

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)
```



```
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())
```

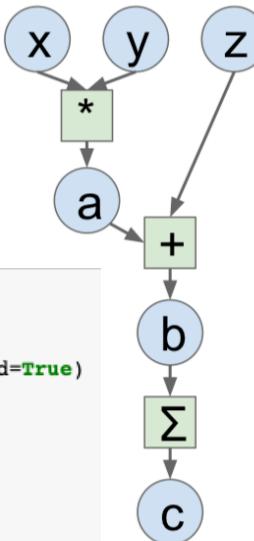
Then, initialize a `tf.Session`

# Computational Graph

Tensorflow 1.2

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)
```



```
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)
```

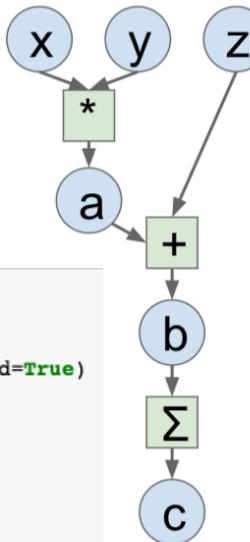
Then, perform the forward pass

# Computational Graph

Tensorflow 2.0

PyTorch

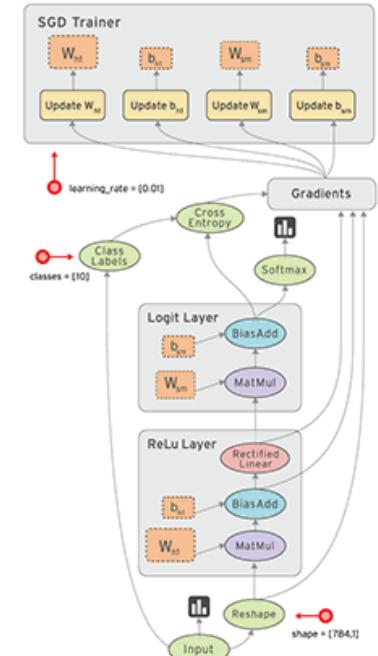
```
1 import torch
2 N, D = 3, 4
3
4 x = torch.randn(N, D, requires_grad=True)
5 y = torch.randn(N, D)
6 z = torch.randn(N, D)
7
8 a = x * y
9 b = a + z
10 c = b.sum()
11
12 c.backward()
13 print(x.grad)
```



```
1 N, D = 3, 4
2 x = tf.Variable(tf.ones(shape=(N, D)), name="x")
3 y = tf.Variable(tf.ones(shape=(N, D)), name="y")
4 z = tf.Variable(tf.ones(shape=(N, D)), name="z")
5
6
7 with tf.GradientTape() as tape:
8     a = x * y
9     b = a + z
10    c = tf.reduce_sum(b)
11    grad = tape.gradient(c, [x])
12    print(grad)
```

# Tensorflow/Keras: Static Graph

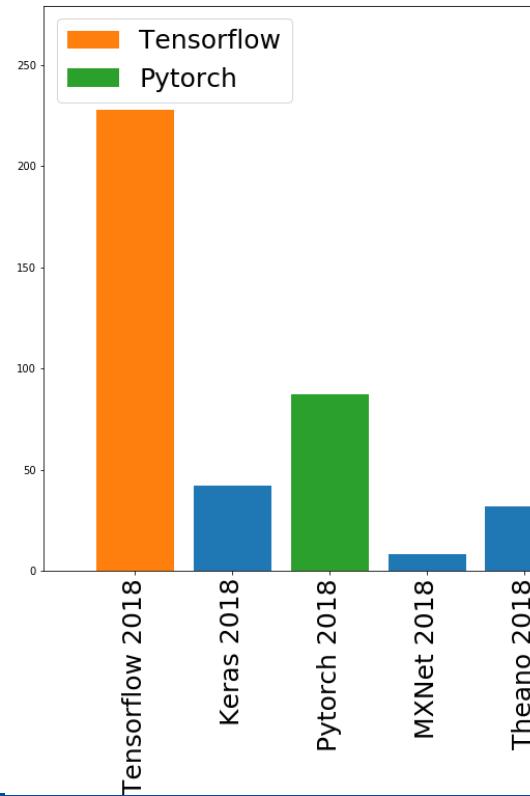
- Pros:
  - Faster?
  - Large community
  - Keras is a simpler API
- Cons:
  - Not intuitive
  - Hard to debug



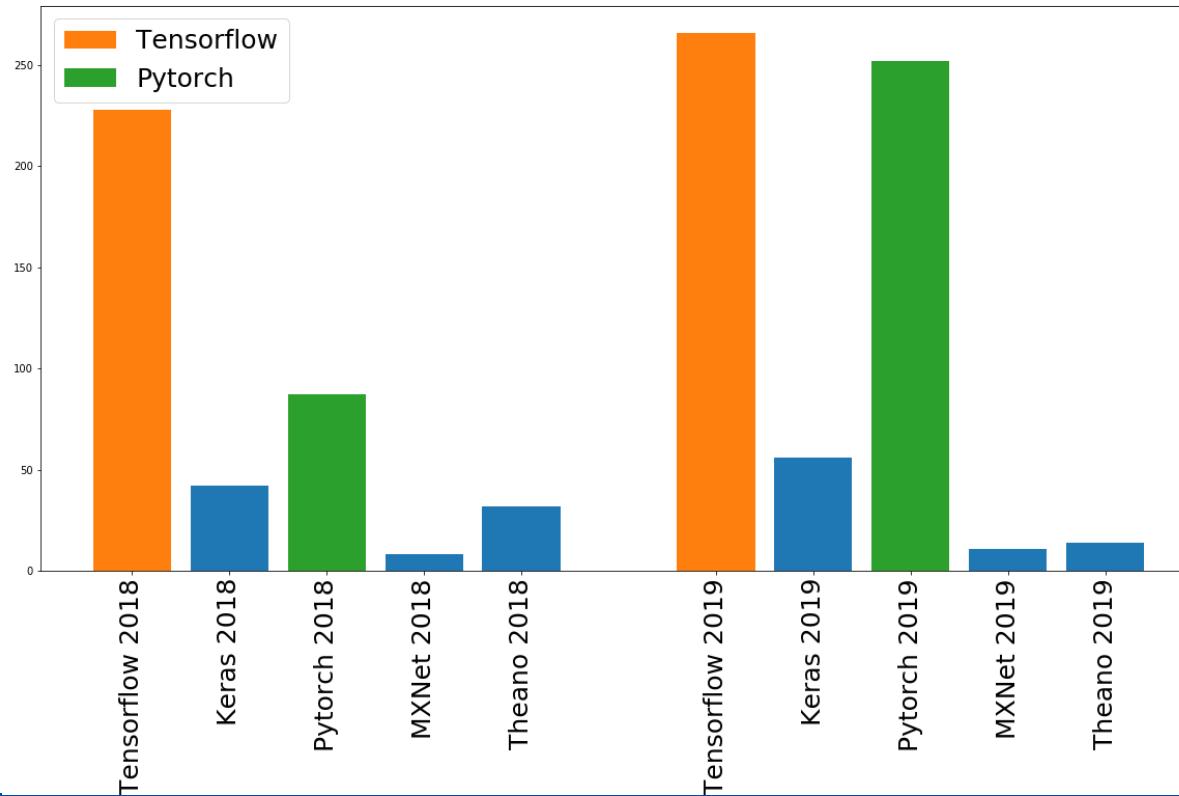
# Dynamic vs Static Graphs

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

# 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.

# The Pytorch Ecosystem

# The Pytorch Ecosystem

CUDnn

CUDA

# The Pytorch Ecosystem

Pytorch (torch)

Torchvision

CUDnn

CUDA

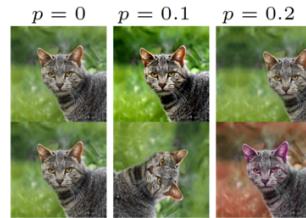
# The Pytorch Ecosystem

[Detectron2](#)



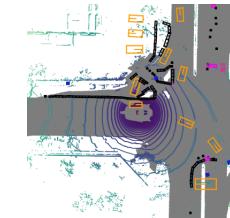
Pytorch (torch)

[Stylegan2](#)



Torchvision

[MMdetection3d](#)



CUDnn

CUDA

# 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_()
28
```

# 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_()
```

# 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_()
```

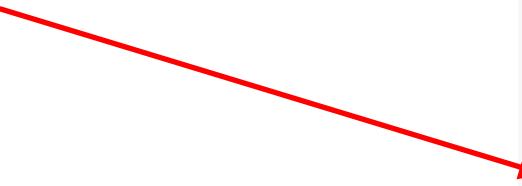
# Pytorch: Tensors

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

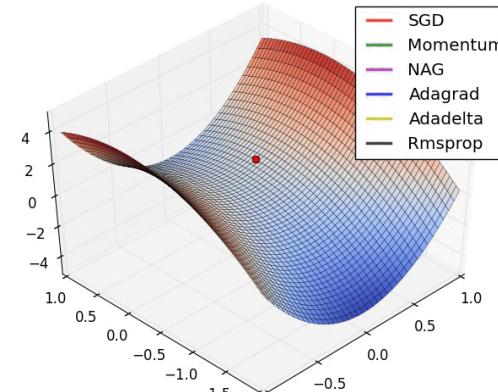
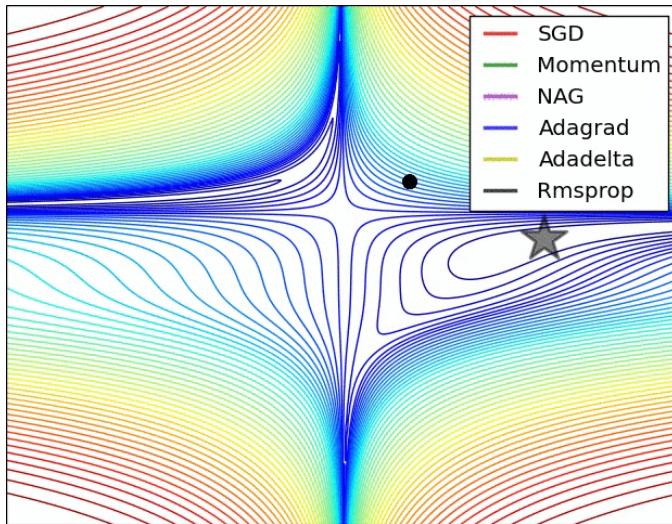
```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 model = nn.Sequential(
6     nn.Linear(I, J),
7     nn.Sigmoid(),
8     nn.Linear(J, C)
9     # No need for softmax, since its included in
10    # nn.CrossEntropyLoss()
11 )
12
13 loss_function = torch.nn.CrossEntropyLoss()
14 optimizer = torch.optim.SGD(model.parameters(),
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

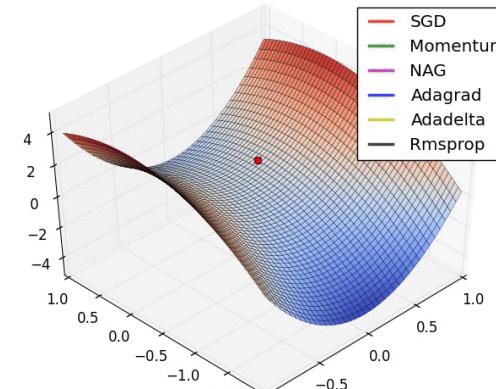
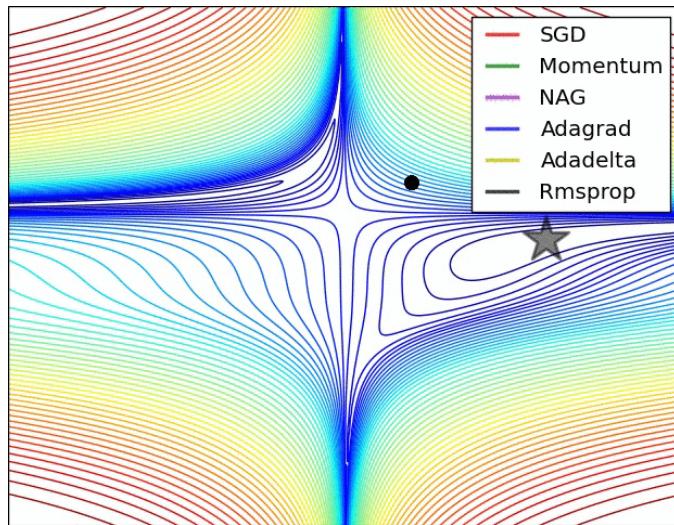
```
1 I = 785
2 J = 64
3 C = 10
4 learning_rate = 1e-5
5 model = nn.Sequential(
6     nn.Linear(I, J),
7     nn.Sigmoid(),
8     nn.Linear(J, C)
9     # No need for softmax, since its included in
10    # nn.CrossEntropyLoss()
11 )
12
13 loss_function = torch.nn.CrossEntropyLoss()
14 optimizer = torch.optim.SGD(model.parameters(),
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()
```

Perform gradient step and reset the gradients

# The choice of optimizer is important!



# The choice of optimizer is important!



Try out `optim.Adam` in your assignment!

What about more complex models?

```
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: 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()
34
```

# 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()
```

# 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
```

# Computing Resources

Three options:

- Visual Computing Cluster ([tdt4265.idi.ntnu.no/  
ssh](http://tdt4265.idi.ntnu.no/ssh))
- Cybele Computers (at campus and remote)
- Google Colab?

Tutorials in assignment text!