

ECE-GY 9143

# Introduction to High Performance Machine Learning

**Lecture 4 09/18/24**

Parijat Dube

Zehra Sura

# Agenda

- PyTorch - Deep Learning Concepts
- PyTorch Optimizer
- PyTorch Multiprocessing
- PyTorch Dataloader

# PyTorch - Deep Learning Concepts

# PyTorch

*A Python-based scientific computing package targeted at two sets of audiences:*

- A replacement for numpy to use the power of GPUs*
- A deep-learning research platform that provides maximum flexibility and speed*

*[This lesson uses material from <http://pytorch.org/tutorials/> throughout.]*

- To install: <https://github.com/pytorch/pytorch#installation>

# Tensors

- Tensors are matrix-like data structures which are essential components in deep learning libraries and efficient computation.
- GPUs are especially effective at calculating operations between tensors

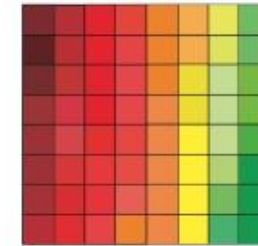
tensor = multidimensional array

vector



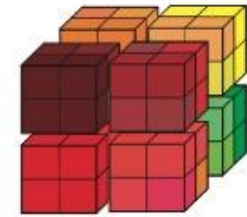
$$\mathbf{v} \in \mathbb{R}^{64}$$

matrix



$$\mathbf{X} \in \mathbb{R}^{8 \times 8}$$

tensor



$$\mathbf{X} \in \mathbb{R}^{4 \times 4 \times 4}$$

From: <https://www.slideshare.net/BertonEarnshaw/a-brief-survey-of-tensors>

- Tensor operations:
  - ones, zeros, add, dot, etc.
- PyTorch tensors can live on
  - CPU
  - GPU (speedup!)
  - Or other accelerators

# PyTorch Tensors

- Import Torch:

```
from __future__ import print_function
import torch
```

- Construct a 2x3 matrix, uninitialized:

```
x = torch.Tensor(2, 3)
print(x)
0.00e+00  0.00e+00  1.15e-24
1.58e-29  1.67e-37  2.97e-41
[torch.FloatTensor of size 2x3]
```

- Use tensor in CUDA

```
device = torch.device("cuda")
y = torch.ones_like(x, device=device) # directly create a
tensor on GPU
x = x.to(device) # or just use .to("cuda")
```

- Initialize zeros or ones tensors

```
x = torch.zeros(2,3)
x = torch.ones(2,3)
```

- Convert a numpy array

```
b = torch.from_numpy(a)
```

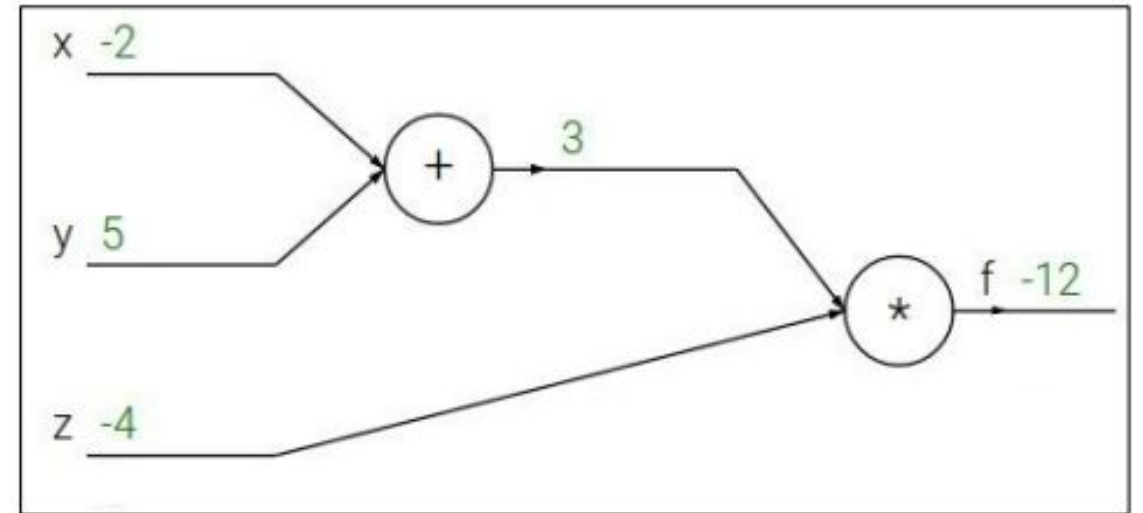
- the slice functionality is available like in numpy

```
x = torch.rand(2,3) # Initialize a tensor randomly
print x[:,1] #second column
0.6297 0.1196
[torch.FloatTensor of size 2]
print x[0,:] #first row
0.9749 0.6297 0.3045
[torch.FloatTensor of size 3]
```

# Backward Propagation Example

- Initial Function  $f(x, y, z) = (x + y) * z$
- Computation Graph Functions:
  - $q(x, y) = x + y$
  - $f(q, z) = qz$
- Inputs:  $x = -2, y = 5, z = -4$
- Want to obtain:
  - $\frac{df}{dx}, \frac{df}{dy}, \frac{df}{dz}$

Computation Graph for  $(x+y) * z$



From [http://cs231n.stanford.edu/slides/2017/cs231n\\_2017\\_lecture4.pdf](http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture4.pdf)

(Don't confuse **Computation Graph** with actual **Neural Network**!)

# Backward Propagation Example

- Computation Graph Functions:

- $q(x, y) = x + y$

- $f(q, z) = qz$

- Basic gradients :

- $\frac{dq(x,y)}{dx} = 1, \frac{dq(x,y)}{dy} = 1$

- $\frac{df(q,z)}{dq} = z, \frac{df(q,z)}{dz} = q$

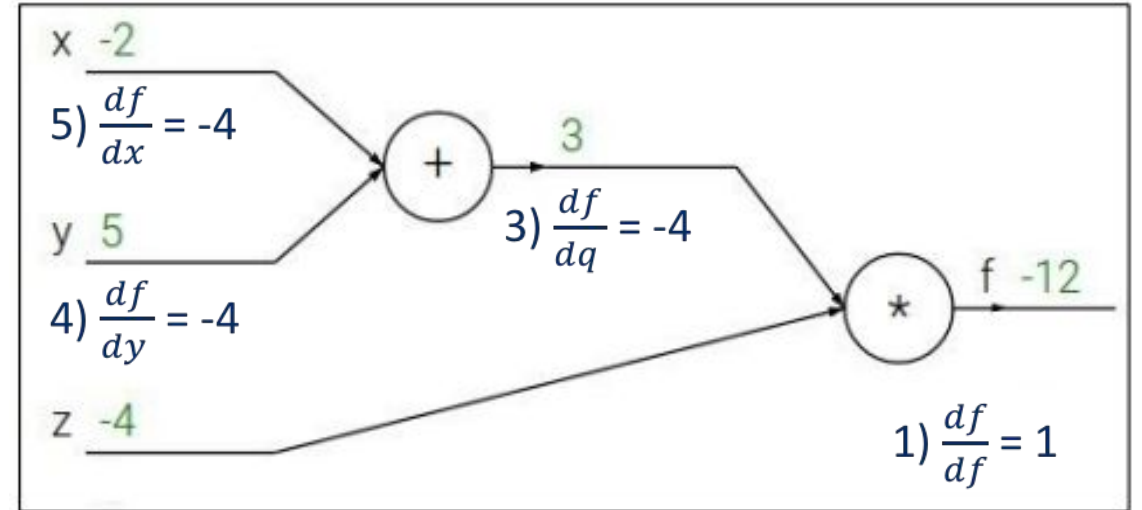
- Compute gradients with chain rule:

- $\frac{df}{dz} = q = 3$

- $\frac{df}{dx} = \frac{df}{dq} * \frac{dq}{dx} = z * 1 = -4$

- $\frac{df}{dy} = \frac{df}{dq} * \frac{dq}{dy} = z * 1 = -4$

Computation Graph for  $(x+y)*z$



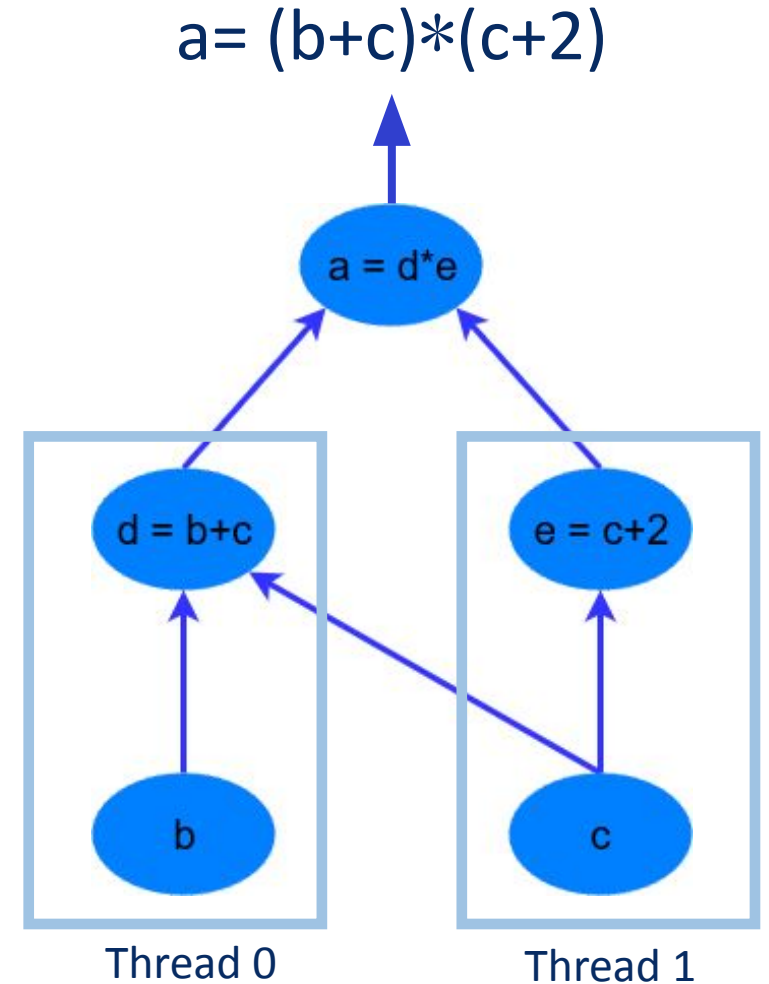
2)  $\frac{df}{dz} = 3$

From [http://cs231n.stanford.edu/slides/2017/cs231n\\_2017\\_lecture4.pdf](http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture4.pdf)



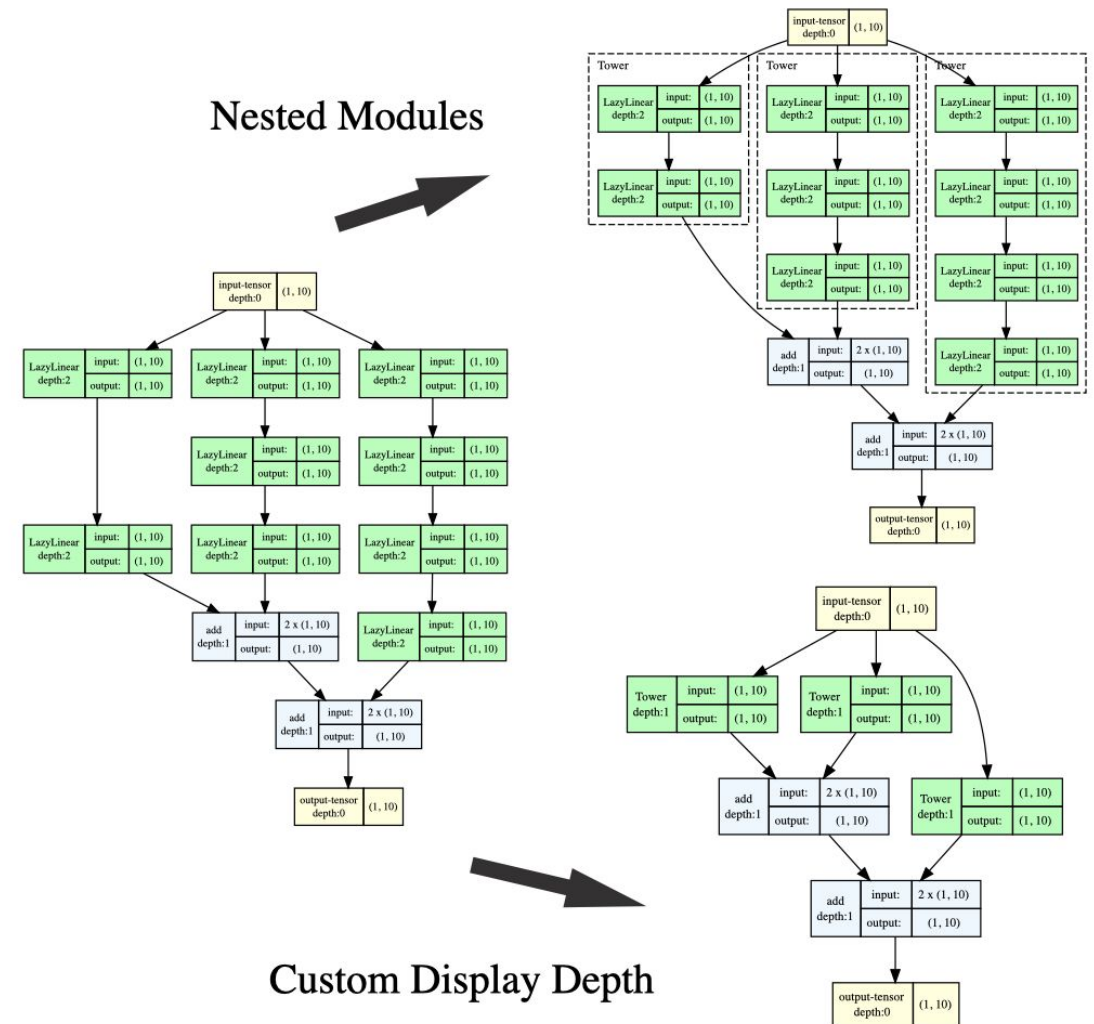
# Computation Graphs

- A computational graph represents a function in a directed acyclic graph of its component functions
- **Performance optimizations:** Computation Graph exposes parallelism!
- In PyTorch the graph construction is **dynamic**: the graph is built at run-time
  - Easier debugging
  - Better for some algorithms (RNNs)
- In TensorFlow graph construction is **static**: meaning the graph is “compiled” and then run
  - Compiler adds latency but can also apply optimizations



# Visualizing PyTorch Computational Graphs

- Example: Torchview provides visualization of PyTorch models in the form of visual graphs. Visualization includes tensors, modules, torch.functions, and info such as input/output shapes.
- Link: <https://github.com/mert-kurtutun/torchview>
- Example Notebooks
  - Introduction Notebook
  - Computer Vision Models Notebook
  - NLP Models Notebook



# Autograd in PyTorch

- **Autograd** builds the Computation Graph **Dynamically**
- The **Tensor** class is the main component of this autograd system in PyTorch (from PyTorch 0.4 version, the *Variable* class is deprecated)
- If you set its attribute *.requires\_grad* as *True*, it starts to track all operations on it
- The gradient for this tensor will be accumulated into *.grad* attribute
- Tensors allow automatic gradient computation when the *.backward()* function is called
- Based on the graph, *<variable>.backward()* computes the **gradient** and writes it in **grad**
  - Example: **b.backward()** computes  $\frac{d(y)}{dx}$

```
x = torch.randn(5, 5) # requires_grad=False by default
y = torch.randn(5, 5) # requires_grad=False by default
z = torch.randn((5, 5), requires_grad=True)
a = x + y
a.requires_grad
    False
b = a + z
b.requires_grad
    True
b.backward
```

# PyTorch Tensors, Functions and Gradients

- Create a tensor

```
x = torch.tensor(torch.ones(2, 2) * 2,  
requires_grad=True)
```

- Do a simple math equation:

```
z = 2 * (x * x) + 5 * x
```

- To get the gradient of this operation with respect to  $x$  i.e.  $dz/dx$  we can analytically calculate this.

- If all elements of  $x$  are 2, then we should expect the gradient  $dz/dx$  to be a (2, 2) shaped tensor with 13-values.
- However, first we have to run the `.backwards()` operation to compute these gradients.
- To compute gradients, we need to compute them with respect to something.
- In this case, we can supply a (2,2) tensor of 1-values to be what we compute the gradients against – so the calculation simply becomes  $d/dx$ :

```
z.backward(torch.ones(2, 2)*2)  
print(x.grad)  
      tensor([[ 13.,  13.],  
              [ 13.,  13.]])
```

# PyTorch Neural Network

- *torch.nn.module* is used to define a neural network
- Example: NN with 3 fully connected layers
  - Using RELU activation for 1<sup>st</sup> and 2<sup>nd</sup> layer
  - Input to 1<sup>st</sup> FC layer: 256 features
  - Input to 2<sup>nd</sup> FC layer: 120 values
  - Input to 3<sup>rd</sup> FC layer: 84 values
- Define only forward prop. : backward prop is automatically derived from it

```
import torch
import torch.nn as nn
import torch.nn.functional as F
```

```
#Inherit from class nn.Module
```

```
class Net(nn.Module):
```

```
    def __init__(self):
```

```
        super(Net, self).__init__()
```

```
        #  $y = Wx + b$ 
```

```
        self.fc1 = nn.Linear(256, 120)
```

```
        self.fc2 = nn.Linear(120, 84)
```

```
        self.fc3 = nn.Linear(84, 10)
```

```
    def forward(self, x):
```

```
        x = F.relu(self.fc1(x))
```

```
        x = F.relu(self.fc2(x))
```

```
        x = self.fc3(x)
```

```
        return x
```

# PyTorch Loss Function

- Loss function:
  - How “far” from the **target** is the **output** of forward propagation (prediction)
- NN provides various loss functions with syntax:
  - *loss = <loss-function>(output, target)*
  - *loss, output and target are Tensors*

```
#net is the network previously defined
#input is the input data of the network
net = Net()
input = torch.randn(256) # a dummy input
output = net(input)
#criterion is a Mean-Squared Error loss function
criterion = nn.MSELoss()

target = torch.arange(1, 11) # a dummy target
#target comes from the labelled dataset
loss = criterion(output, target)

print(loss)
```

# PyTorch Backpropagation and weights update

- First reset gradients of the network
- Compute backward prop.
  - It uses *autograd* inside
- Update the weights with SGD:  
 $weight = weight - learning\_rate * gradient$
- Different optimization algorithms are in *torch.optim*

```
# Zeroes the gradient buffer of all parameters  
net.zero_grad()
```

```
#Backpropagation step  
loss.backward()
```

```
#Stochastic Gradient Descent weights update  
learning_rate = 0.01  
for f in net.parameters():  
    f.data.sub_(f.grad.data * learning_rate)
```

# Gradient Descent Optimization Algorithms and PyTorch



# ML Performance Factors

## Algorithms Performance

- **PyTorch Optimizer (training): Momentum, Nesterov, Adagrad, AdaDelta**

## Hyperparameters Performance

- **Learning rate, Momentum, Batch size, Others**

## Implementation Performance

- **Pytorch Multiprocessing, PyTorch DataLoader, PyTorch CUDA**

## Framework Performance

- **ML Frameworks: PyTorch, TensorFlow, Caffe, MXNET**

## Libraries Performance

- **Math libraries (cuDNN), Communication Libraries (MPI, GLOO)**

## Hardware Performance

- **CPU, DRAM, GPU, HBM, Tensor Units, Disk/Filesystem, Network**

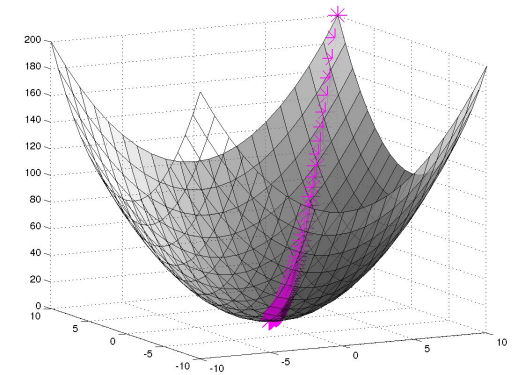
# PyTorch Optimizer

# Gradient Descent - Recap

- Simplest and very popular
- Main Idea: take a step proportional to the negative of the gradient (i.e. minimize the loss function):

$$\theta = \theta - \alpha \nabla J(\theta)$$

- Where  $\theta$  is the parameters vector,  $\alpha$  is the learning rate, and  $\nabla J(\theta)$  is the gradient of the cost
- Easy to implement
- Each iteration is relatively cheap
- **Can be slow to converge**
- Gradient descent variants:
  - **Batch gradient descent:** Update after computing all the training samples
  - **Stochastic gradient descent:** Update for each training sample
  - **Mini-batch gradient descent:** Update after a subset (mini-batch) of training samples



# Learning Rate Challenges

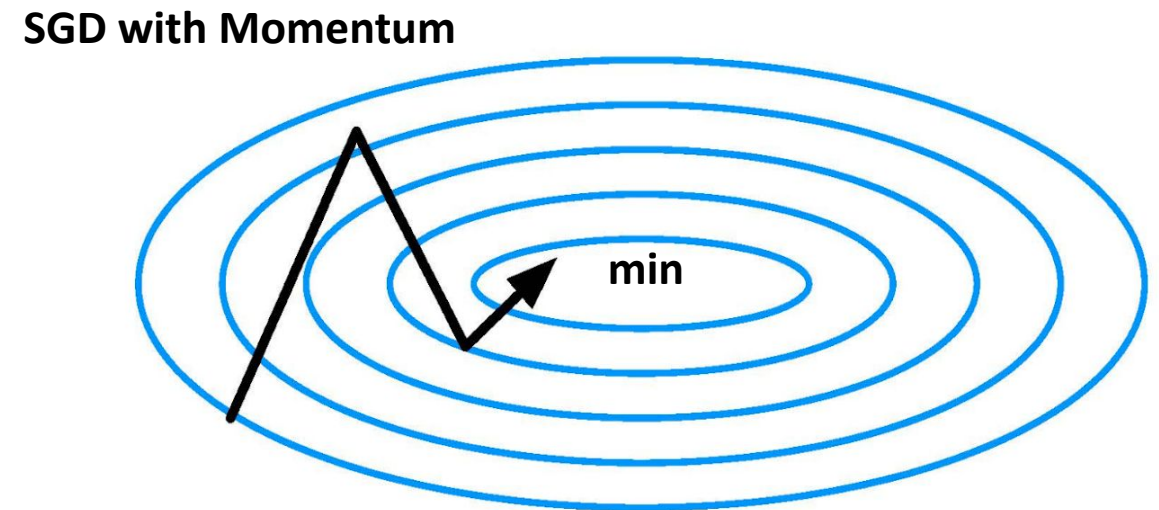
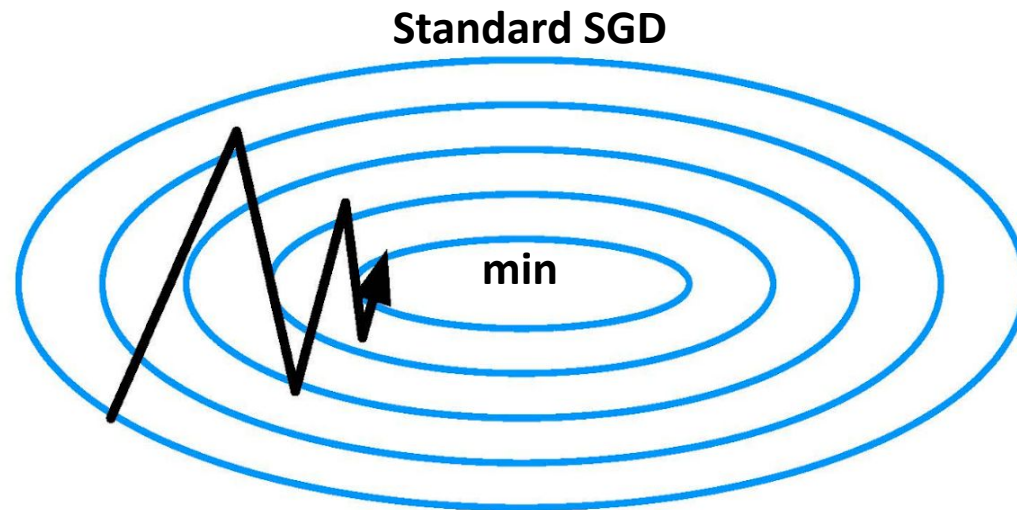
- The Learning rate has a large impact on convergence
  - Too small  $\square$  too slow
  - Too large  $\square$  oscillatory and may even diverge
- Choosing a proper (initial) learning rate
- Should learning rate be fixed or adaptive?
  - **Decaying learning rate:** drop by 10 after N iterations and then again after other M iterations, and so on...
  - How to do define an **adaptive learning rate**?
- Avoid to get trapped in local minima
- Changing learning rate for each parameter

# Learning Rate Challenges (II)

- Traversing efficiently through error differentiable functions' surfaces is an important research area today
- Some of the recently popular techniques **take the gradient history into account** such that each “move” on the error function surface relies on previous ones a bit
- Many algorithms already implemented in PyTorch
  - However, these algorithms often used as black-box tools: need to understand their strength and weakness

# Optimizer Algorithms – Momentum

- SGD with Momentum:
  - Descent with momentum keeps going in the same direction longer
  - Descent is faster because it takes less steps ( $W$  updates)



From: <http://www.del2z.com/2016/06/param-optimiz-3/>

# Optimizer Algorithms - Momentum

- Momentum is a simple method that helps accelerate SGD by adding a fraction  $\gamma$  of the update vector of the past time step to the current update vector
- In simple words momentum adds a **velocity** component to the parameter update routine

$$v_t = \gamma v_{t-1} + \alpha \nabla J(\theta)$$

$$\theta = \theta - v_t$$

- In PyTorch, momentum is implemented in the default SGD method

```
optimizer = torch.optim.SGD(model.parameters(), lr=0.1, momentum=0.9)
optimizer.zero_grad()
loss_fn(model(input), target).backward()
optimizer.step()
```

# Optimizer Algorithms - Nesterov Momentum

- Instead of computing the gradient of the current position, it computes the gradient at the approximated new position
- Use the **next approximated position's gradient** with the hope that it will give us better information when we're taking the next step:

$$v_t = \gamma v_{t-1} + \alpha \nabla J(\theta - \gamma v_{t-1})$$

$$\theta = \theta - v_t$$

- Momentum is usually set to 0.9
- In PyTorch, Nesterov momentum is implemented in the default SGD method

```
optimizer = torch.optim.SGD(model.parameters(), lr=0.1, momentum=0.9, nesterov=True)
```



# Classical vs Nesterov Momentum

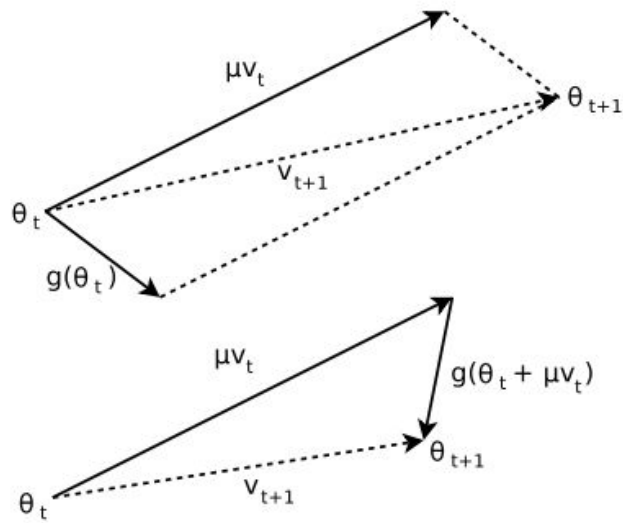
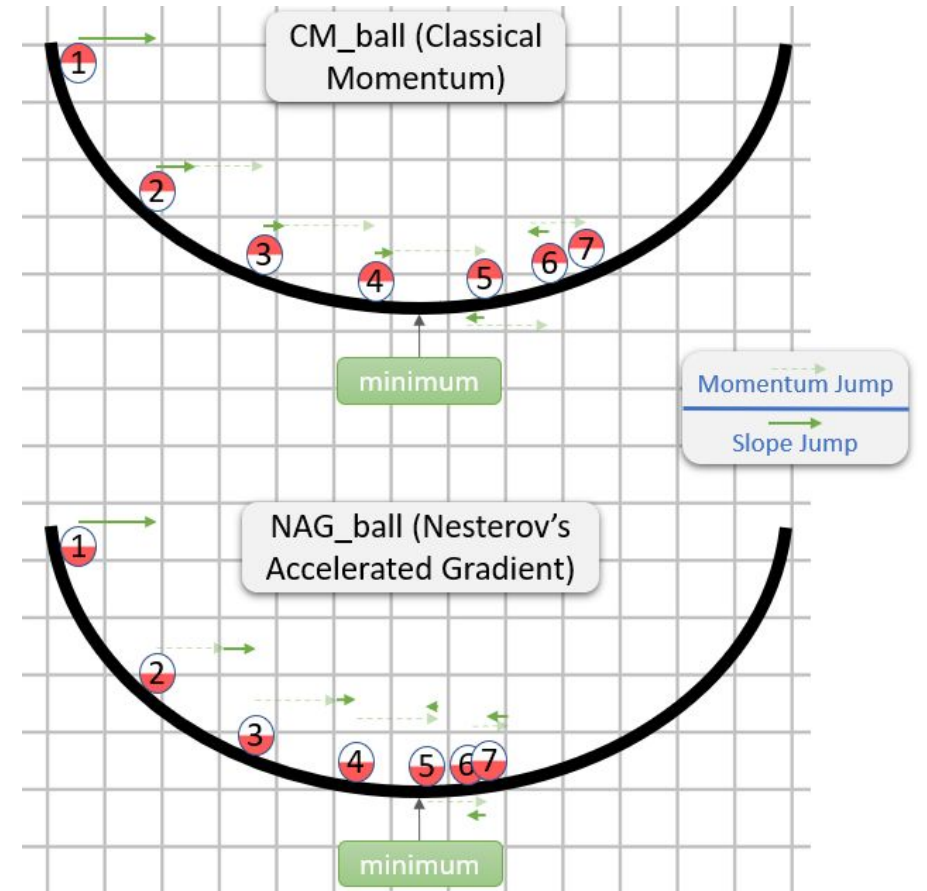


Figure 1. **(Top)** Classical Momentum **(Bottom)** Nesterov Accelerated Gradient



# Optimizer Algorithms - Adagrad

- Adapts learning rate to parameters
- AdaGrad update rule is given by the following formula:

$$\theta_{t+1} = \theta_t - \frac{\alpha}{\sqrt{\text{diag}(G_t) + \epsilon}} \odot \nabla J(\theta_t)$$

- $\theta_t$  vector of parameters at step  $t$  (size  $d$ )
- $\nabla J(\theta_t)$  vector of gradients at step  $t$
- $\odot$  is the element-wise product
- $G_t$  is the historical gradient information until step  $t$ :
  - diagonal matrix  $d \times d$  ( $d = \text{\#parameters}$ ) with sum of squares of gradients until time  $t$
- $\text{diag}()$ : Transform  $G_t$  diagonal into a vector
- $\epsilon$  is the smoothing term to avoid division by 0 (usually  $1\text{e-}8$ )
- In PyTorch:

```
optimizer = torch.optim.Adagrad(params, lr=0.01)
```

# Optimizer Algorithms - Adagrad II

- Pros:
  - It is well-suited for dealing with sparse data
  - It greatly improved the robustness of SGD
  - It eliminates the need to manually tune the learning rate
- Cons:
  - Main weakness is its accumulation of the squared gradients in the denominator
  - This causes the learning rate to shrink and become infinitesimally small
  - The algorithm can no longer acquire additional knowledge
- In PyTorch:

```
optimizer = torch.optim.Adagrad(params, lr=0.01)
```

# Optimizer Algorithms - Adadelta

- One of the inspiration for AdaDelta was to improve AdaGrad weakness of learning rate converging to zero with increase of time
- Adadelta mixes two ideas:
  1. to scale learning rate based on historical gradient while taking into account only recent time window – not the whole history, like AdaGrad
  2. to use component that serves an acceleration term, that accumulates historical updates (similar to momentum)
- Adadelta step is composed of the following phases:
  1. Compute gradient  $g_t$  at current step  $t$
  2. Accumulate gradients (AdaGrad-like step)
  3. Compute update
  4. Accumulate updates (momentum-like step)
  5. Apply the update
- In PyTorch:

---

**Algorithm 1** Computing ADADELTA update at time  $t$ 

---

**Require:** Decay rate  $\rho$ , Constant  $\epsilon$

**Require:** Initial parameter  $x_1$

- 1: Initialize accumulation variables  $E[g^2]_0 = 0, E[\Delta x^2]_0 = 0$
  - 2: **for**  $t = 1 : T$  **do** %% Loop over # of updates
  - 3:   Compute Gradient:  $g_t$
  - 4:   Accumulate Gradient:  $E[g^2]_t = \rho E[g^2]_{t-1} + (1 - \rho)g_t^2$
  - 5:   Compute Update:  $\Delta x_t = -\frac{\text{RMS}[\Delta x]_{t-1}}{\text{RMS}[g]_t} g_t$
  - 6:   Accumulate Updates:  $E[\Delta x^2]_t = \rho E[\Delta x^2]_{t-1} + (1 - \rho)\Delta x_t^2$
  - 7:   Apply Update:  $x_{t+1} = x_t + \Delta x_t$
  - 8: **end for**
- 

```
optimizer = torch.optim.Adadelta(params, lr=1.0, rho=0.9, eps=1e-06)
```

# Optimizer Algorithms - Adam

- Adam might be seen as a generalization of AdaGrad
  - AdaGrad is Adam with certain parameters choice
- The idea is to mix benefits of:
  - **AdaGrad** that maintains a per-parameter learning rate that improves performance on problems with sparse gradients (e.g. NLP and computer vision problems).
  - **RMSProp** that also maintains per-parameter learning rates that are adapted based on the average of recent magnitudes of the gradients for the weight (e.g. how quickly it is changing).
    - This means the algorithm does well on online and non-stationary problems (e.g. noisy)
- Specifically, the algorithm calculates an exponential moving average of the gradient and the squared gradient, and the parameters  $\beta_1$  and  $\beta_2$  control the decay rates of these moving averages

1.  $m_k = \beta_1 m_{k-1} + (1 - \beta_1) g_k$  The first moment

2.  $v_k = \beta_2 v_{k-1} + (1 - \beta_2) g_k^2$  The second moment

3.  $\hat{m}_k = \frac{m_k}{1 - \beta_1^k}$  Compute bias-corrected first moment estimate

4.  $\hat{v}_k = \frac{v_k}{1 - \beta_2^k}$  Compute bias-corrected second raw moment estimate

5.  $\theta_k = \theta_{k-1} - \alpha \frac{\hat{m}_k}{\sqrt{\hat{v}_k} + \epsilon}$  Update parameters

- In PyTorch:

```
torch.optim.Adam(params, lr=0.001, betas=(0.9, 0.999), eps=1e-08)
```

# And many others...

- AdaMax: variant of Adam
- Nadam
  - Adam with Nesterov momentum instead of classical momentum
- RMSprop
  - Divide the gradient by a running average of its recent magnitude
- ....
- Here some interesting readings:
  - <https://arxiv.org/abs/1609.04747>
  - <http://runder.io/optimizing-gradient-descent/>

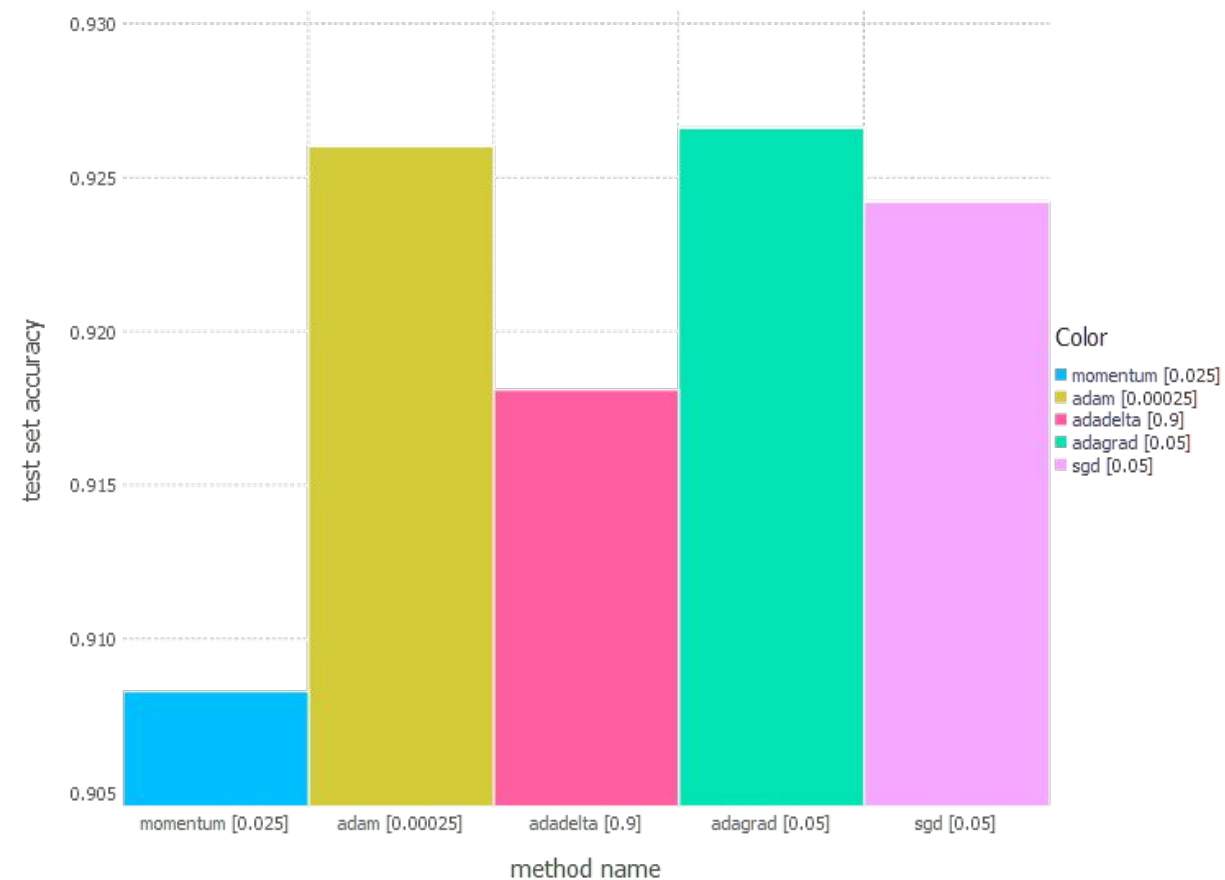
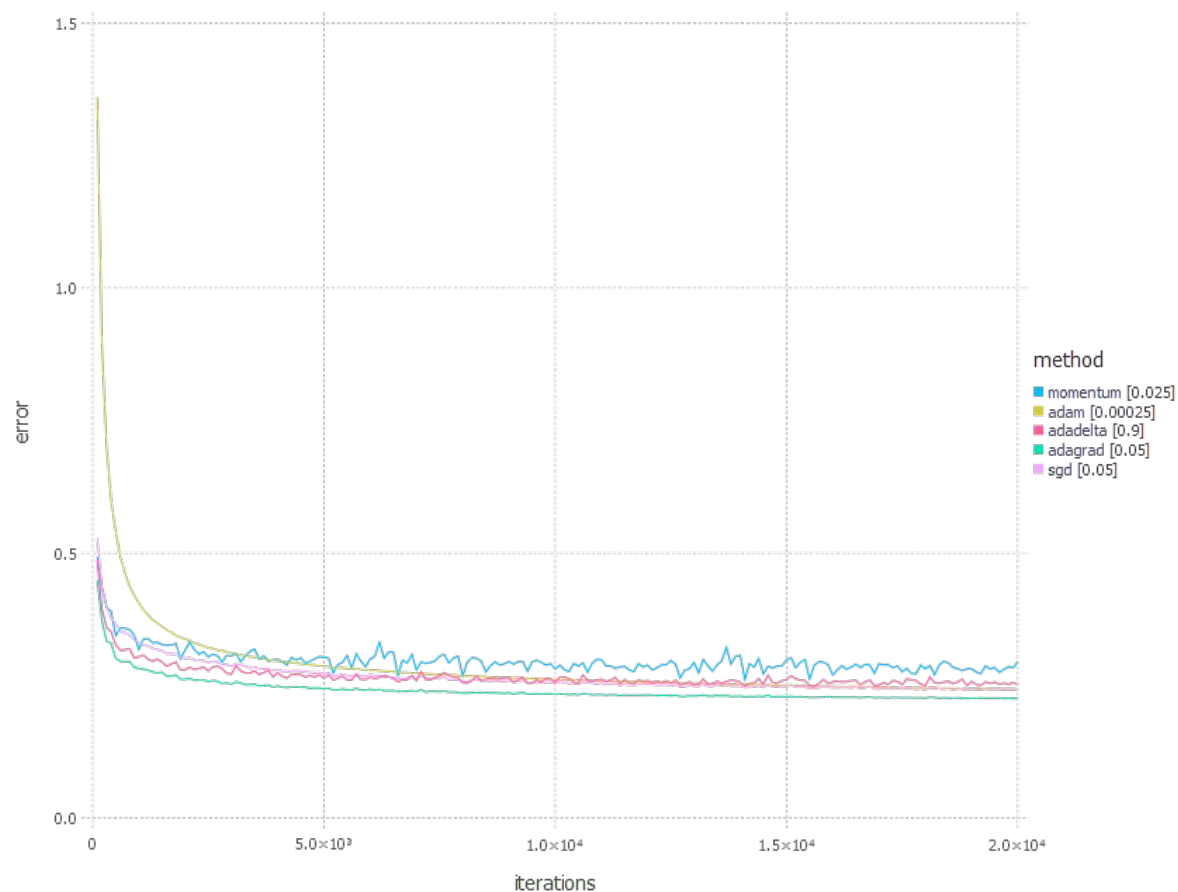
# Optimization Techniques Comparison (I)

- Toy example on MNIST Classification
- Three different network architecture tested:
  1. network with linear layer and softmax output (softmax classification)
  2. network with sigmoid layer (100 neurons), linear layer and softmax output
  3. network with sigmoid layer (300 neurons), ReLU layer (100 neurons), sigmoid layer (50 neurons) again, linear layer and softmax output
- Mini-batch size of 128
- Run the algorithm for approx. 42 epochs (20000 iterations)
- <http://int8.io/comparison-of-optimization-techniques-stochastic-gradient-descent-momentum-adagrad-and-adadelta/>



# Results on net 1

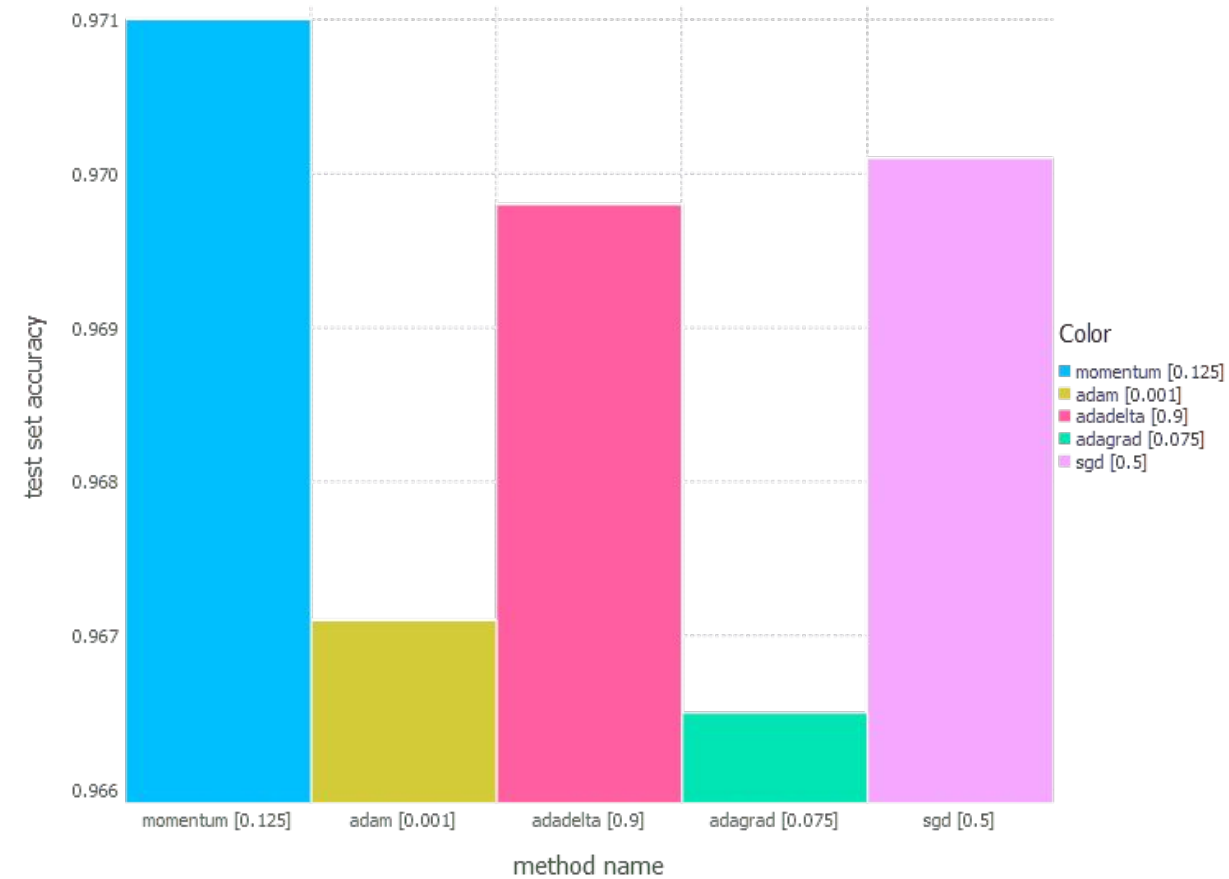
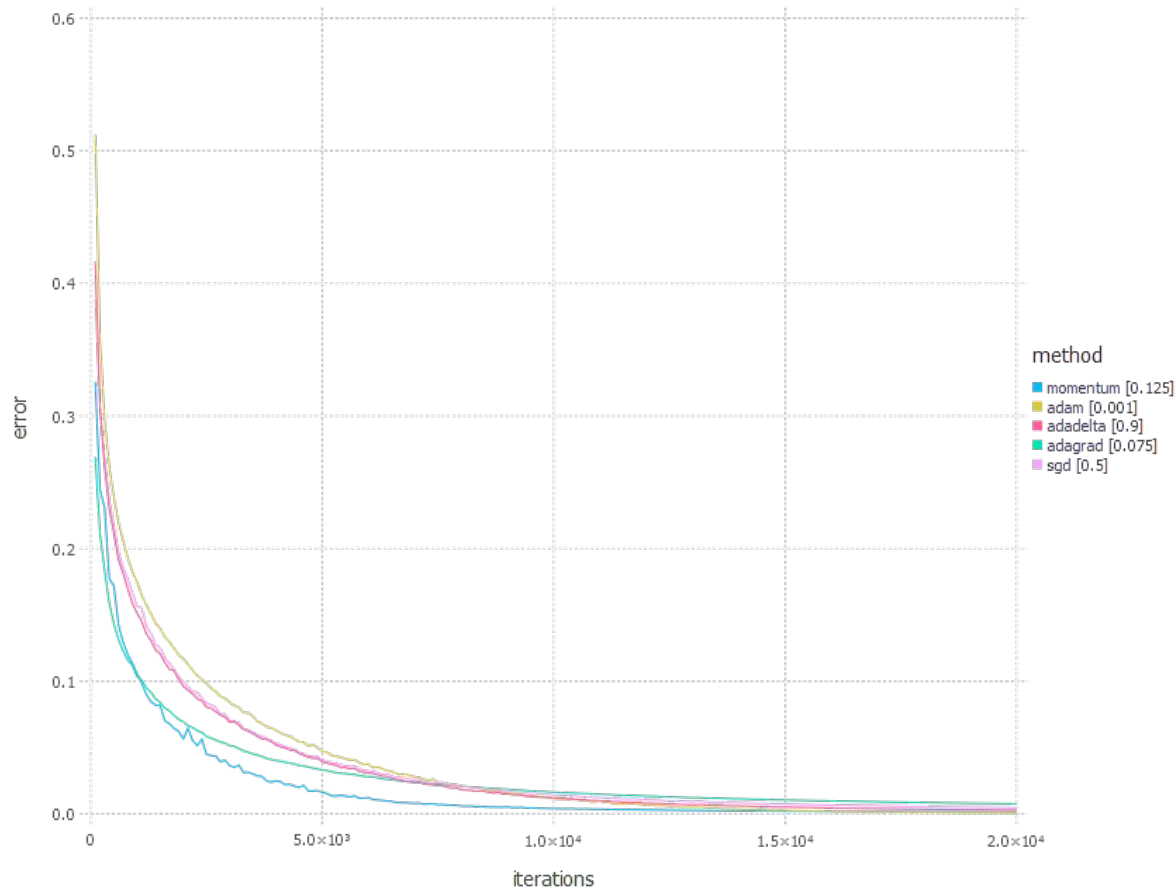
network with linear layer and softmax output (softmax classification)





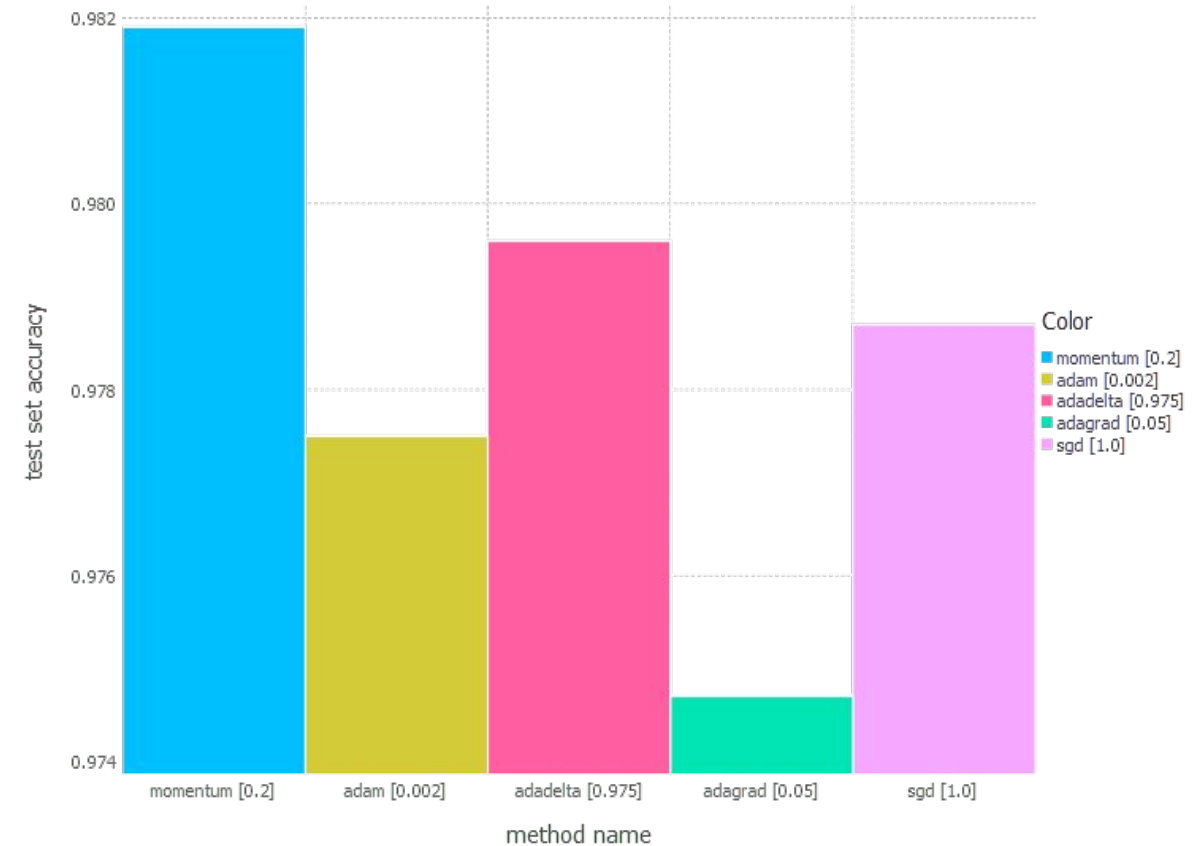
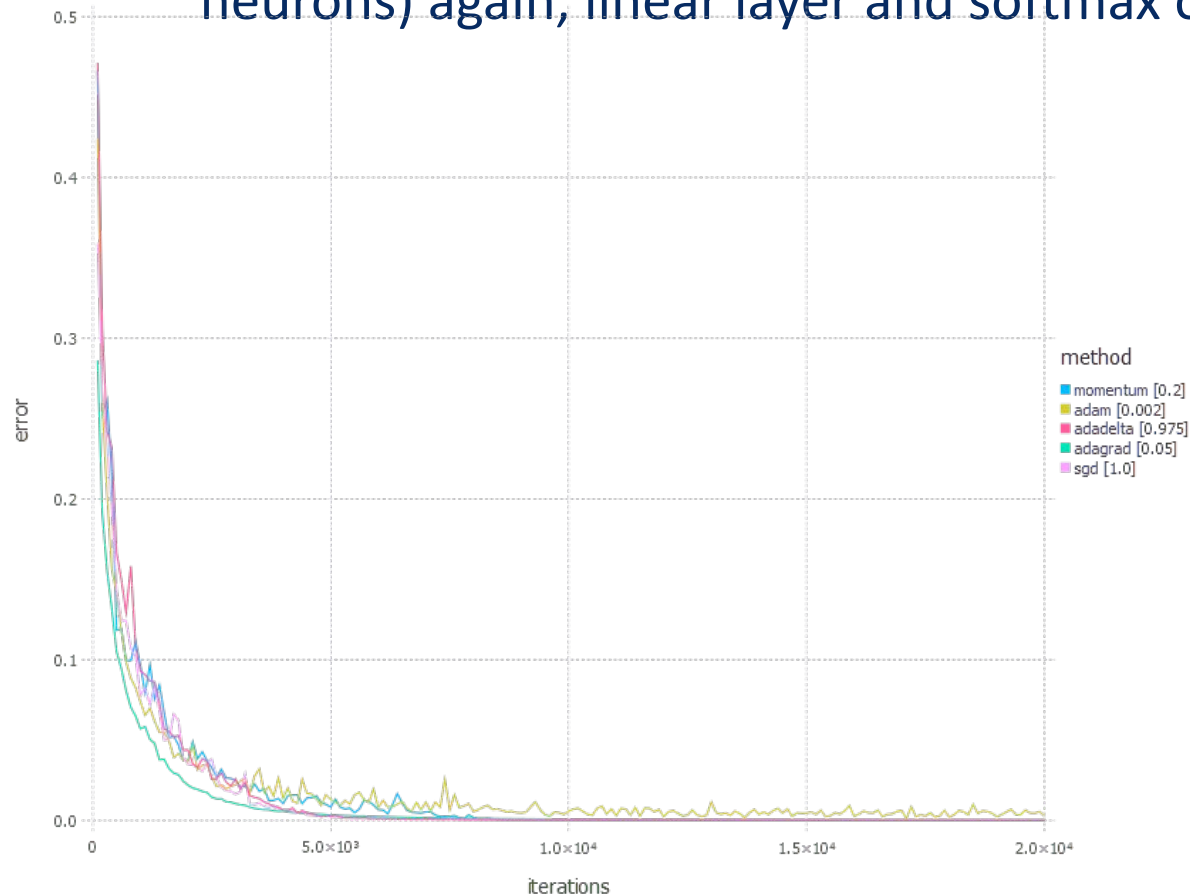
# Results on net 2

network with sigmoid layer (100 neurons), linear layer and softmax output



# Results on net 3

network with sigmoid layer (300 neurons), ReLU layer (100 neurons), sigmoid layer (50 neurons) again, linear layer and softmax output



# So...Which optimizer to use?

- Unfortunately there isn't a clear answer
- As in the toy example before, different networks have completely different behaviors
- If your input data is **sparse**...?
  - You likely achieve the best results using one of the **adaptive learning-rate** methods
  - An additional benefit is that you will not need to tune the learning rate
- Which one is the best?
  - Adagrad, Adadelata, and Adam are very similar algorithms that do well in similar circumstances
  - Insofar, **Adam** might be the best overall choice as trade off between time and accuracy
- Interestingly, many recent papers use SGD without momentum and a simple learning rate annealing schedule
  - SGD usually achieves to find a minimum but it takes much longer time than others, is much more reliant on a robust initialization and annealing schedule
  - ☐ If you care about **fast convergence** and train a deep or complex NN, you should choose one of the **adaptive learning** rate methods

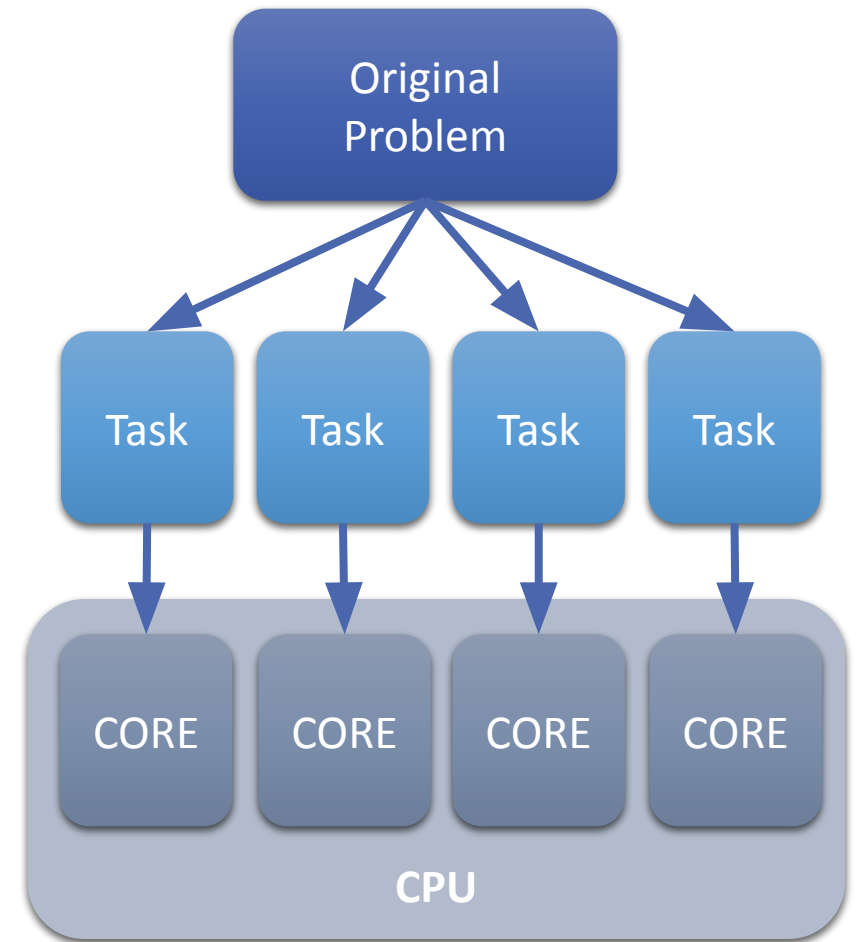
# Comments on optimizer

- All the optimizers take in the parameters and gradients and manipulate over them.
- None of the optimizers are theoretically proven to have better convergence rate than SGD, except Nesterov momentum.
- Even Nesterov momentum is only proved to work on strongly convex problems.
- In practice, SGD with momentum and a per-iteration learning rate schedule often produce better final accuracy with some tuning than automatic update algorithms.

# PyTorch Multiprocessing

# Why Multiprocessing

- **CPU** has many cores and hardware threads
  - Multi-core processors, hardware threads (hyperthreading)
  - Superscalar CPU cores with hardware multi-threading
- **When a problem can be parallelized (i.e. divided in parallel tasks) use parallel tasks to complete it in less time**
- Examples of parallelizable ML tasks:
  - Load multiple files from disk
  - Applying transformations to each dataset sample
  - Send/Receive data from multiple GPUs



# Python Threading

```
class threading.Thread(group=None, target=None, name=None, args=(), kwargs={}, *, daemon=None)
```

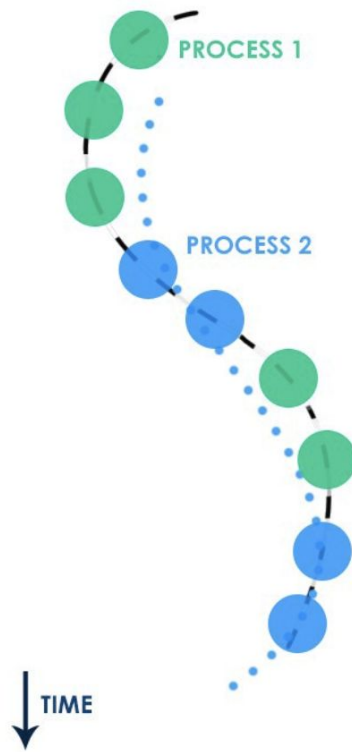
- Allows creation and handling of *Thread* objects: independent execution context
- Threads share data unless its thread local:

```
mydata = threading.local()
mydata.x = 1
```
- Threads in Python:
  - **CONCURRENCY** but **NOT PARALLELISM** (global interpreter lock--GIL)
  - Within a single process only one thread may be processing Python byte-code at any one time.
  - Memory is shared unless is specifically thread local
  - Cannot take advantage of multiprocessing
- Daemon threads:
  - They do not die with the parent (wait for interpreter shutdown) or needs to be explicitly killed

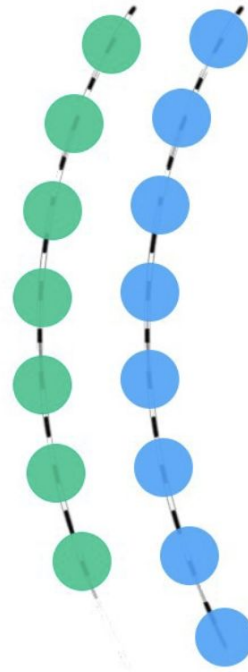
## Not useful for Parallelism Performance!

See <https://docs.python.org/3.6/library/threading.html#module-threading>  
<https://superfastpython.com/thread-local-data/>

## CONCURRENCY



## PARALLELISM





# Python Multiprocessing

- *multiprocessing.Process* class: create another python interpreter process
- Forking and spawning are two different start methods for new processes.
  - 1) SPAWN method: start new fresh process with minimal resource inherited
    - Slower, Default on Windows and MacOS
  - 2) FORK method: `fork()` a new process and inherit all resources (files, pipes, etc.)
    - Faster, default on Unix
    - Can be unstable or incompatible (need to try); If the parent process has threads owning locks that may cause problems in the child process
  - 3) FORK-SERVER method: a new server-process is created that forks new child processes
    - keeps the original process safer
    - Minimal set of resources inherited
  - In general SPAWN and FORK-SERVER methods are safer but slower
- **Creating a process is much slower and heavy than creating a thread**
- <https://docs.python.org/3.6/library/multiprocessing.html#multiprocessing-programming>

# Python Multiprocessing and Pool examples

- Example: creation of a *process*
- Example: creation of a *pool*

```
from multiprocessing import Process, Pool

def f(name):
    print('hello', name)

def f(x):
    return x*x

if __name__ == '__main__':
    p = Process(target=f, args=('bob',))
    p.start()
    p.join()

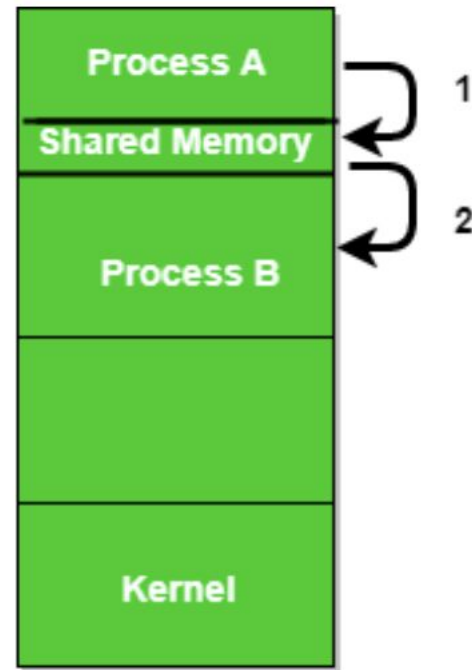
    with Pool(5) as p:
        print(p.map(f, [1, 2, 3]))
```

# PyTorch *torch.multiprocessing*

- Replaces standard Python *multiprocessing*
  - Provides ability to send tensors *efficiently*
- Why is it difficult or not efficient to send *tensors* among Python processes?
  - It is inefficient to send tensors because they need to be serialized before being sent to another process (pickle class), since the address space is different (cpython pointers cannot be passed)
  - Serializing and sending Tensors also implies a memory copy
    - This is called deep copy
- High Performance Solution: *multiprocessing.Queue*
  - Copy and Serialize tensors in a *shared memory area* and only pass handles
  - Additional performance improvement: *multiprocessing.Queue* multiple threads to serialize and send objects
- High Performance tips for torch.multiprocessing
- See <https://pytorch.org/docs/master/notes/multiprocessing.html#reuse-buffers-passed-through-a-queue>

# Shared Memory for Inter Process Communication

Producer-Consumer pattern



<https://www.geeksforgeeks.org/inter-process-communication-ipc/>

# Shared Memory

- **Shared memory** is a memory area that the OS (eg Linux) maps on the address space of the processes, allowing in this way to be simultaneously accessed by multiple programs with an intent to provide communication among them or avoid redundant copies.
  - It is an efficient means of passing data between programs
- Do not confuse this shared memory used as communication between OS processes with the concept of the shared memory in the GPU that is a HW component to reduce GPU memory access latency
  - see also: <https://devblogs.nvidia.com/using-shared-memory-cuda-cc/>

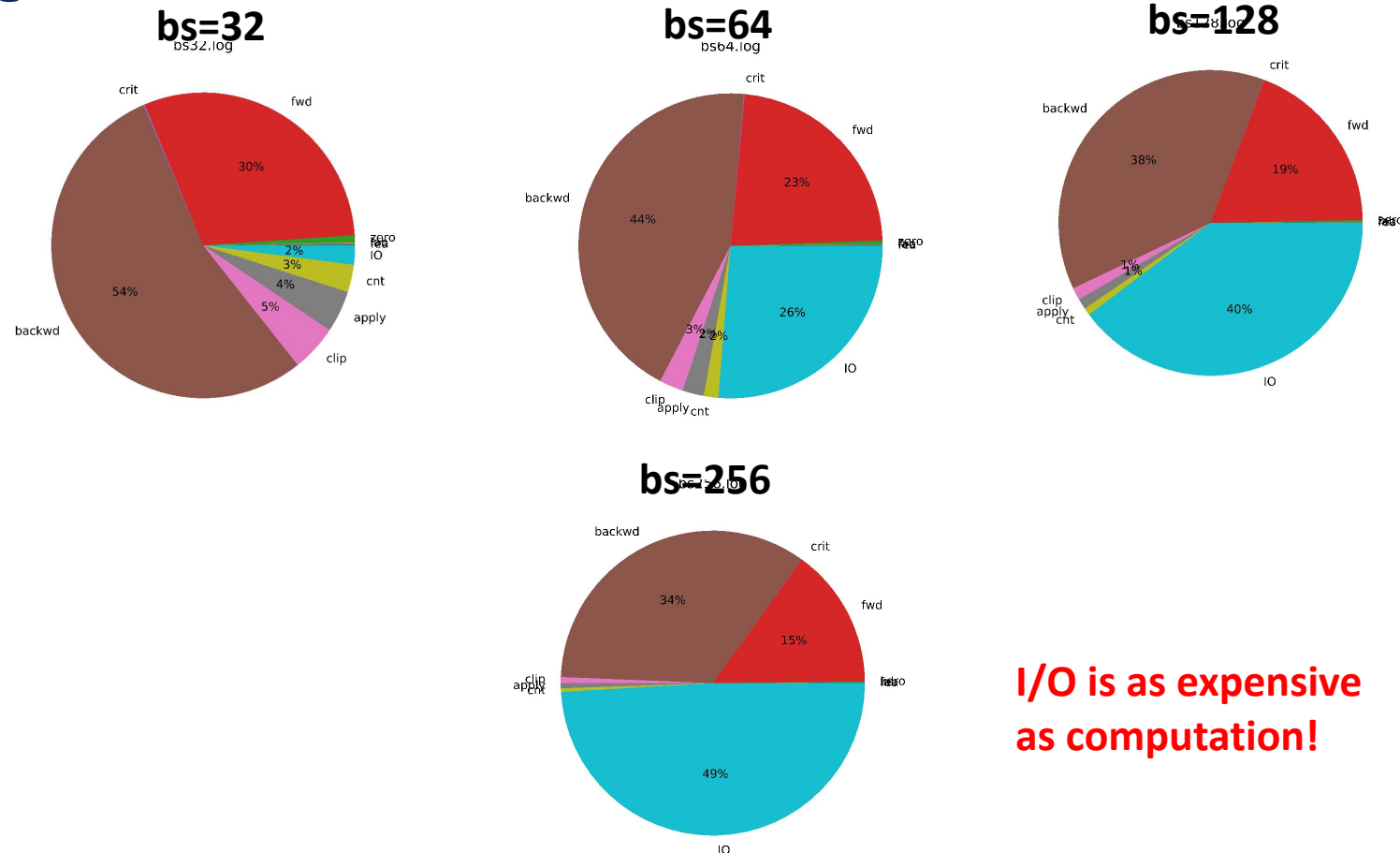
# Shared Memory

- **Shared memory** is a memory area that the OS (eg Linux) maps on the address space of the processes, allowing in this way to be simultaneously accessed by multiple programs with an intent to provide communication among them or avoid redundant copies.
  - It is an efficient means of passing data between programs
- **PyTorch HP tip:** Remember that each time you put a Tensor into a multiprocessing.Queue, it has to be moved into shared memory. If it's already shared, it is a no-op, otherwise it will incur an additional memory copy that can slow down the whole process. Even if you have a pool of processes sending data to a single one, make it send the buffers back - this is nearly free and will let you avoid a copy when sending next batch.
  - from: <https://pytorch.org/docs/stable/notes/multiprocessing.html>
- Do not confuse this shared memory used as communication between OS processes with the concept of the shared memory in the GPU that is a HW component to reduce GPU memory access latency
  - see also: <https://devblogs.nvidia.com/using-shared-memory-cuda-cc/>

# PyTorch Data loading and Preparation

# I/O could be a big problem

- As GPU gets faster, I/O becomes a bottleneck



**I/O is as expensive  
as computation!**



# Loading data - *torch.utils.data.Dataset*

- *torch.utils.data.Dataset* represents a dataset (abstract class)
- Create your own class and override the following:
  - `__len__` so that `len(dataset)` returns the size of the dataset.
  - `__getitem__` to support the indexing such that `dataset[i]` can be used to get i-th sample
  - `__getitem__` will be called to load 1 item at a time

# Custom Dataset Class Example

- Define a set of transformation primitives for each sample
- Pass the transformation primitives through the *transform* argument

```
class FaceLandmarksDataset(Dataset):
    """Face Landmarks dataset."""

    def __init__(self, csv_file, root_dir, transform=None):
        """
        Args:
            csv_file (string):
                Path to the csv file with annotations.
            root_dir (string):
                Directory with all the images.
            transform (callable, optional):
                Optional transform to be applied on a sample.
        """
        self.landmarks_frame = pd.read_csv(csv_file)
        self.root_dir = root_dir
        self.transform = transform
```

[from: http://pytorch.org/tutorials/beginner/data\\_loading\\_tutorial.html#dataset-class](http://pytorch.org/tutorials/beginner/data_loading_tutorial.html#dataset-class)

# Custom Dataset Class Example

- Transform primitives are applied in `__getitem__` to each sample after loading it with `io.imread()`

```
def __len__(self):
    return len(self.landmarks_frame)

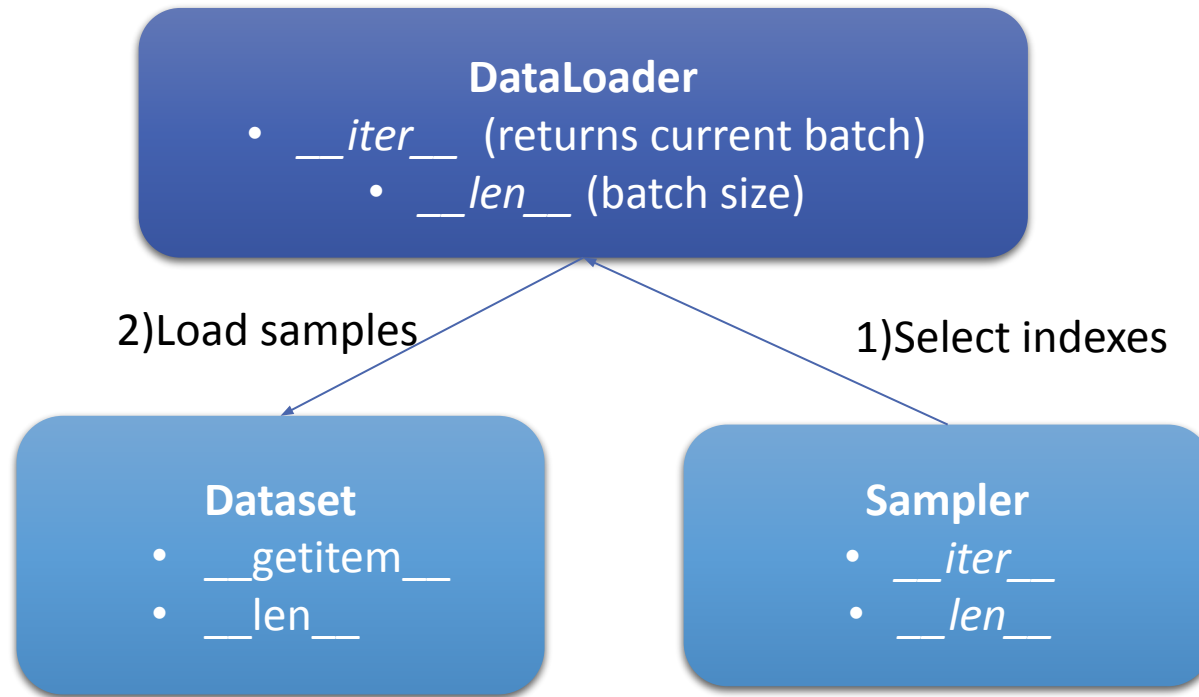
def __getitem__(self, idx):
    img_name = os.path.join(self.root_dir,
                             self.landmarks_frame.iloc[idx, 0])
    image = io.imread(img_name)
    landmarks = self.landmarks_frame.iloc[idx, 1:].as_matrix()
    landmarks = landmarks.astype('float').reshape(-1, 2)
    sample = {'image': image, 'landmarks': landmarks}
    if self.transform:
        sample = self.transform(sample)

    return sample
```

from: [http://pytorch.org/tutorials/beginner/data\\_loading\\_tutorial.html#dataset-class](http://pytorch.org/tutorials/beginner/data_loading_tutorial.html#dataset-class)

# Loading data - *torch.utils.data.DataLoader*

- Selects and Loads data from the Dataset
- Composed of a **Dataset** and a **Sampler**



# Loading data - *torch.utils.data.Dataloader*

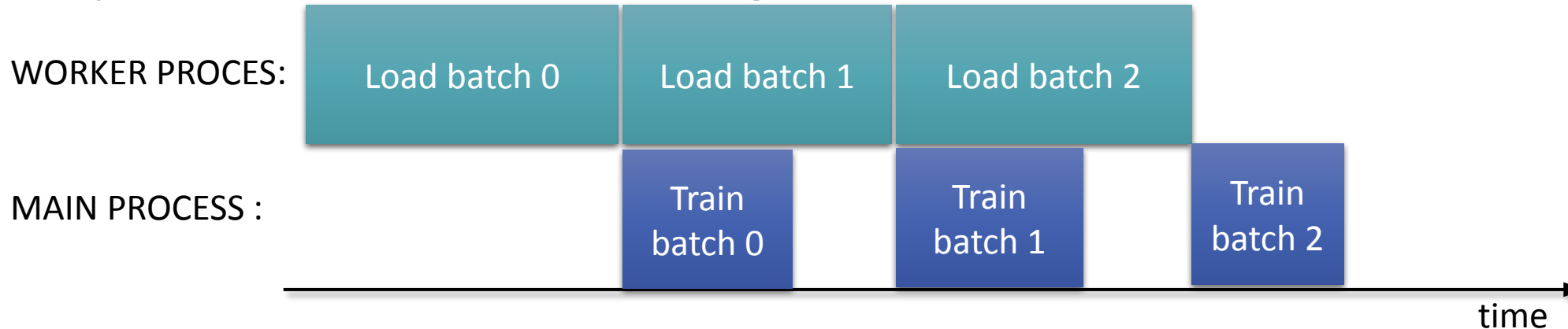
- Useful Parameters:
  - *batch\_size*: batch size
  - *sampler*: define which sampler to use
  - *pin\_memory*: copy into CUDA pinned memory before returning the data
  - *shuffle*: reshuffle data at each epoch
- Available *samplers*:
  - *SequentialSampler*
  - *RandomSampler*
  - *SubsetRandomSampler*
  - *WeightedRandomSampler*
  - Create your *CustomSampler*

# Data Prefetching

- Normal pipeline:



- Pipeline with Load Prefetching:



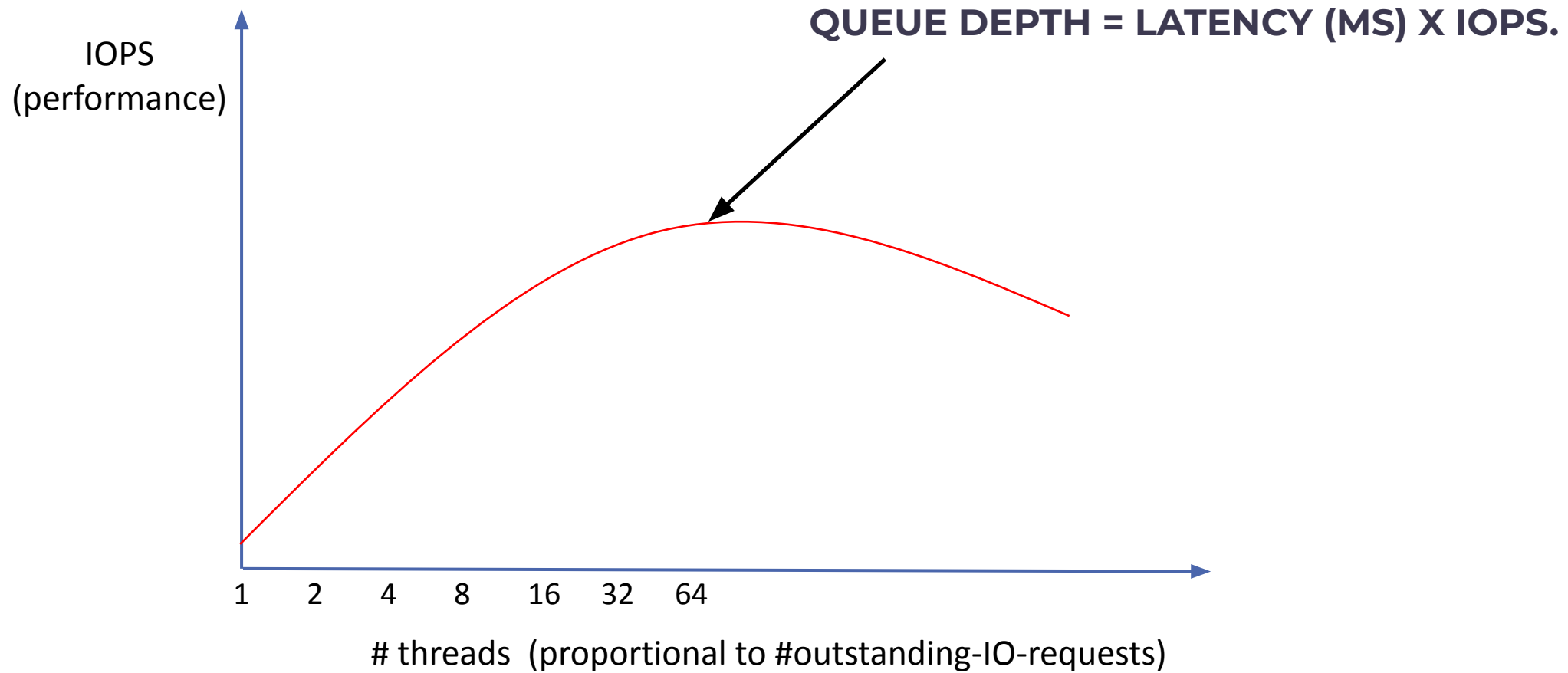
Load time still limiting factor. Can we do better?

# Disk Performance Characteristics

- IOPS: IO-ops/sec
  - read/write
  - random/sequential
- Bandwidth = IOPS \* BlockSize [bytes/sec]
  - Block size: 4KB+
- Queue Depth: maximum number of Outstanding IO requests in a device
- **Important fact: IOPS can be higher with multiple I/O outstanding requests**

<https://en.wikipedia.org/wiki/IOPS>

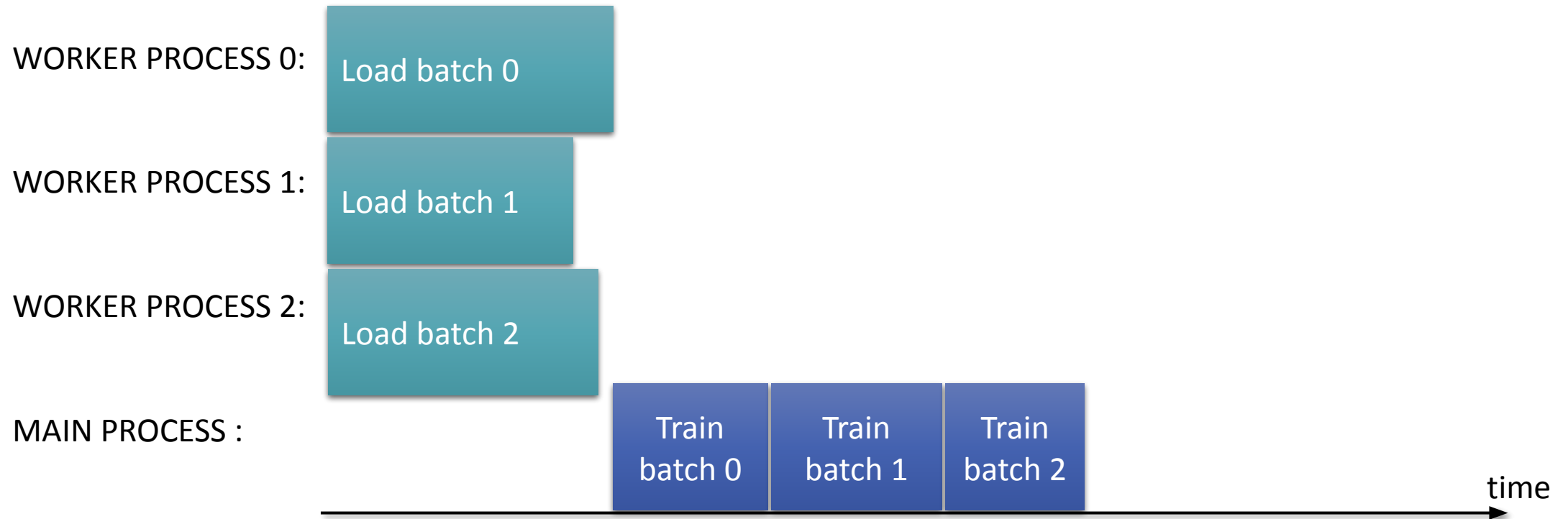
# Disk Performance Characteristics





# Multi-threaded Data Prefetching

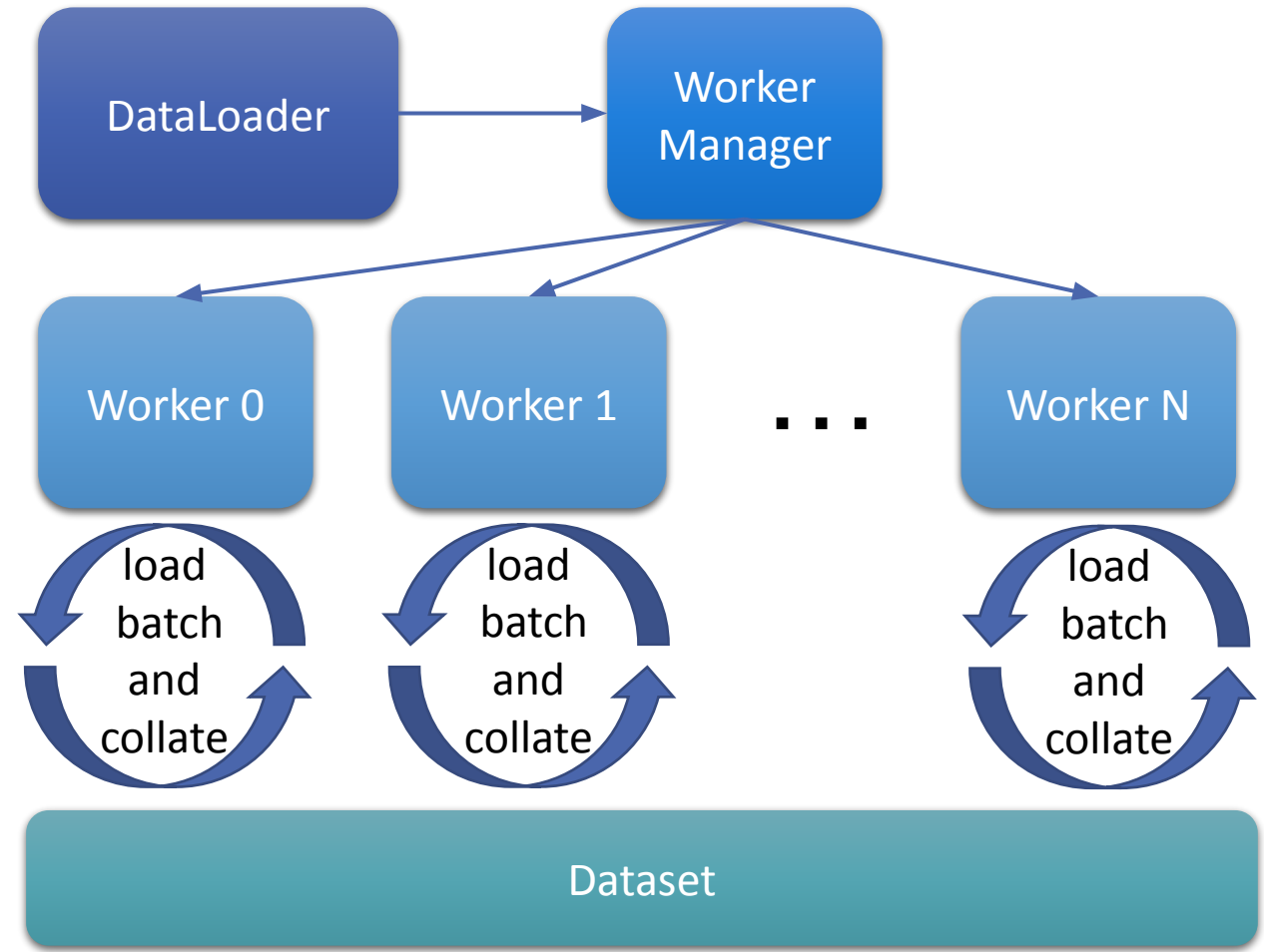
- Issue many outstanding I/O requests using multiple threads:



- Prefetching limitations?

# torch.utils.data.DataLoader Architecture

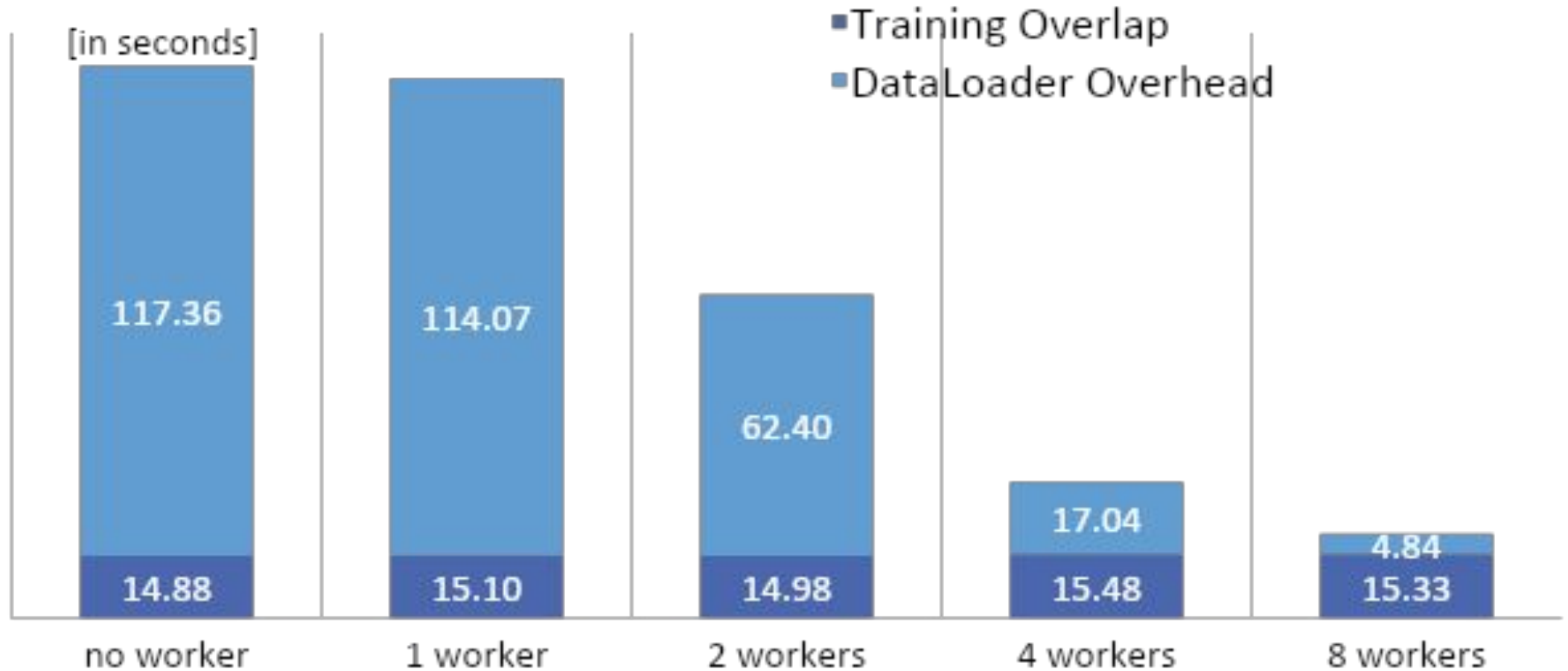
- DataLoader (main process):
  - Requests a set of batches based on *sampler*
- WorkerManager (thread):
  - Dispatch batches to workers
- Workers (processes):
  - Load all samples in each batch
  - Collate batches:
    - Build a batch Tensor with samples
- Synchronization:
  - Producer/Consumer using shared queues and Signals/Exceptions



# Example: Images loading and CNN training

- dataset = ucf101 (video dataset)
- Load images with DataLoader + CNN Network training
- batch\_size = 32, num\_workers = 0 and 4
- I/O: NFS over ethernet 100 Gbit
  - Several seconds per minibatch
  - Most time spent in `__getitem__` of data loader iterator
    - ReadSegmentFlow: 18 ~ 600ms
    - video\_transform: 13 ~ 27 ms
- How many workers is optimal?

# Example: Number of workers effect



Ideal speed-up vs actual speed-up:

- The 8x speedup should be  $(114.07 + 15.10) / 8 = 16.15$ .
- The ideal overhead with 8 workers would be  $16.15 - 15.33 = 0.82$ , but the actual measured is 4.84

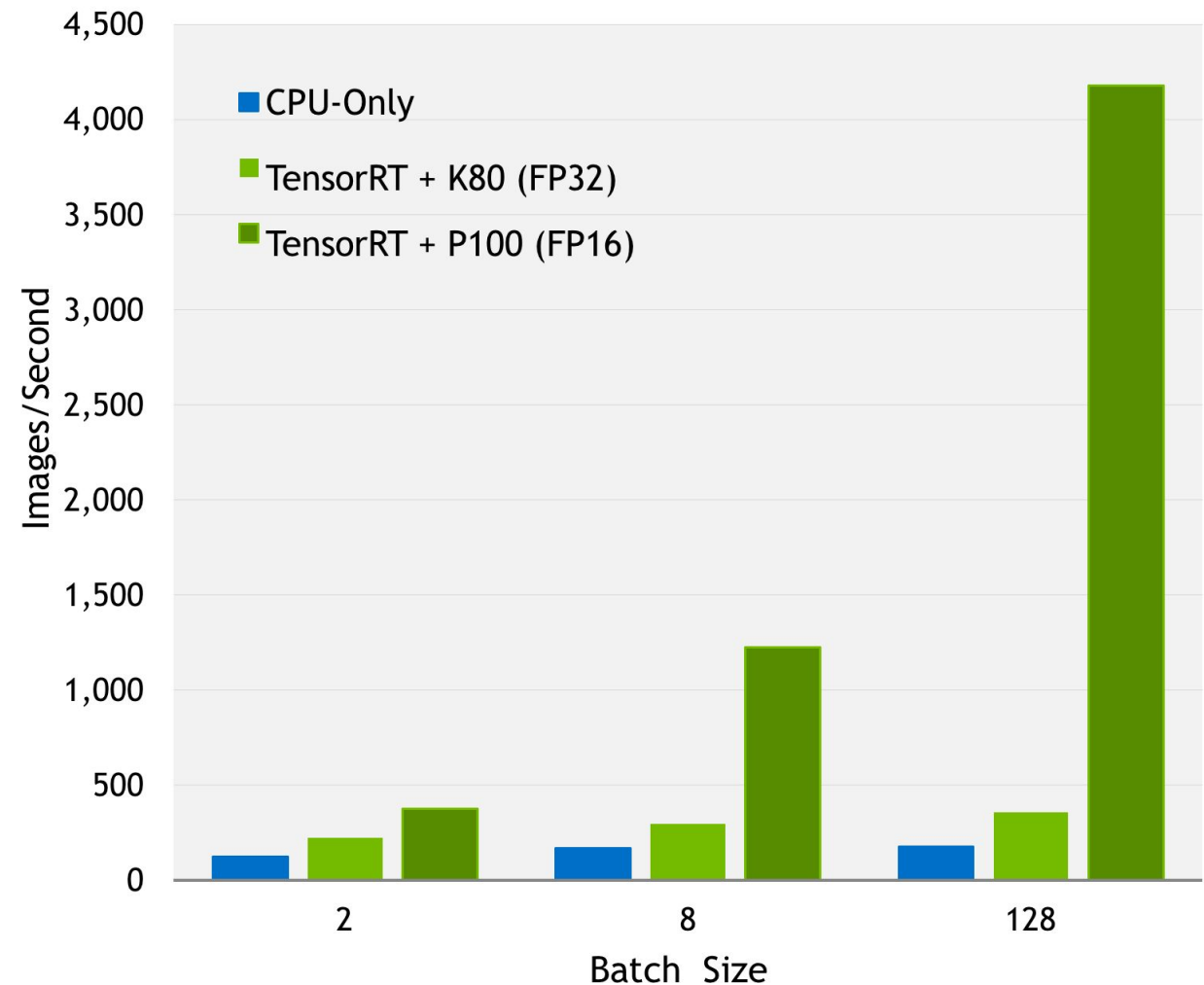
# Using GPUs

CUDA is a parallel computing platform and programming model that makes using a GPU for general purpose computing simple and elegant.

# Why GPUs - Inference

- Intel Xeon E5-2690v4 CPU
- Tesla P100 GPU
- TensorRT is a library developed by NVIDIA for faster inference on NVIDIA graphics processing units

Up To 23x More Images/sec vs. CPU-Only Inference



**GoogLeNet**, Tesla P100 + TensorRT (FP16), Tesla K80 + TensorRT (FP32), CPU-Only + Caffe (FP32)

**CPU:** 1 Socket Broadwell E5-2690 v4@2.6GHz with HT off

From:

<https://devblogs.nvidia.com/deploying-deep-learning-nvidia-tensorrt/>

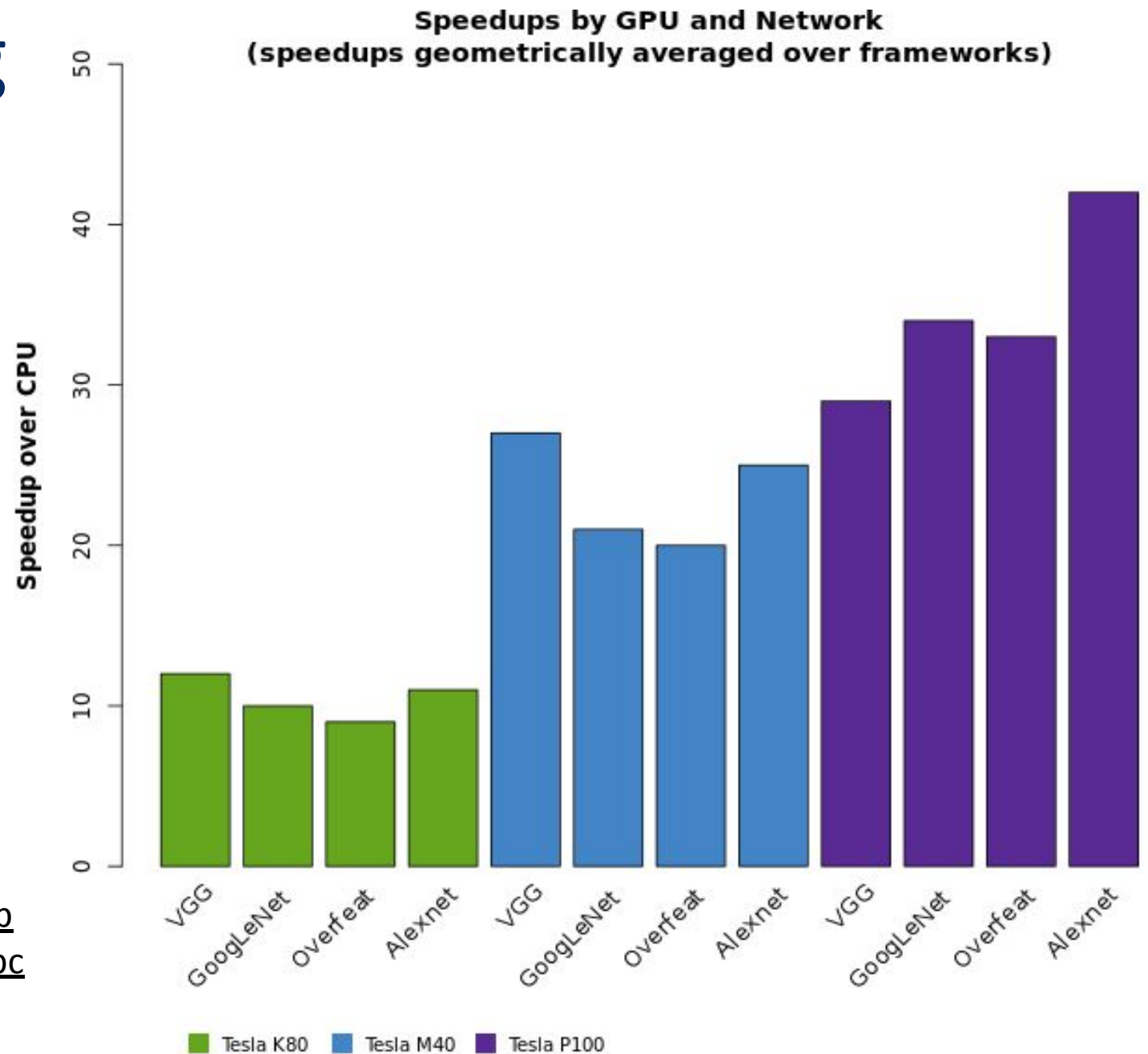
from: NVIDIA

# Why GPUs - Training

- Intel Xeon E5-2690v4 CPU
  - 256GB of DRAM
- Tesla K80 GPU
- Tesla M40 GPU
- Tesla P100 GPU

From:

<https://www.microway.com/hpc-tech-tips/deep-learning-benchmarks-nvidia-tesla-p100-16gb-pcie-tesla-k80-tesla-m40-gpus/>



# Don't give up on CPUs too quickly!

- **Caffe con Troll: Shallow Ideas to Speed Up Deep Learning**
  - <https://arxiv.org/abs/1504.04343>
  - Take-away: training speed is the same as long as FLOPs are the same, independent from hardware
- **Scaling deep learning on GPU and knights landing clusters**
  - <https://dl.acm.org/citation.cfm?id=3126912>
  - Take-away: when batch size is sufficiently large, a large computer cluster is much more economical than a gpu cluster because GPUs have high-margins.



# Lesson Key Points

- PyTorch - Deep Learning Concepts
- PyTorch Optimizer Algorithms
- PyTorch DataLoader and Disk performance
- PyTorch Multiprocessing

# Acknowledgements

- The lecture material is prepared by Giacomo Domeniconi, Parijat Dube, Kaoutar El-Maghraoui, Ulrich Finkler, Alessandro Morari, and Zehra Sura