Hi!

# CS 744: PYTORCH

Shivaram Venkataraman

Fall 2020

# ADMINISTRIVIA

Assignment 2 out! → Due 10/5 (Monday next week)

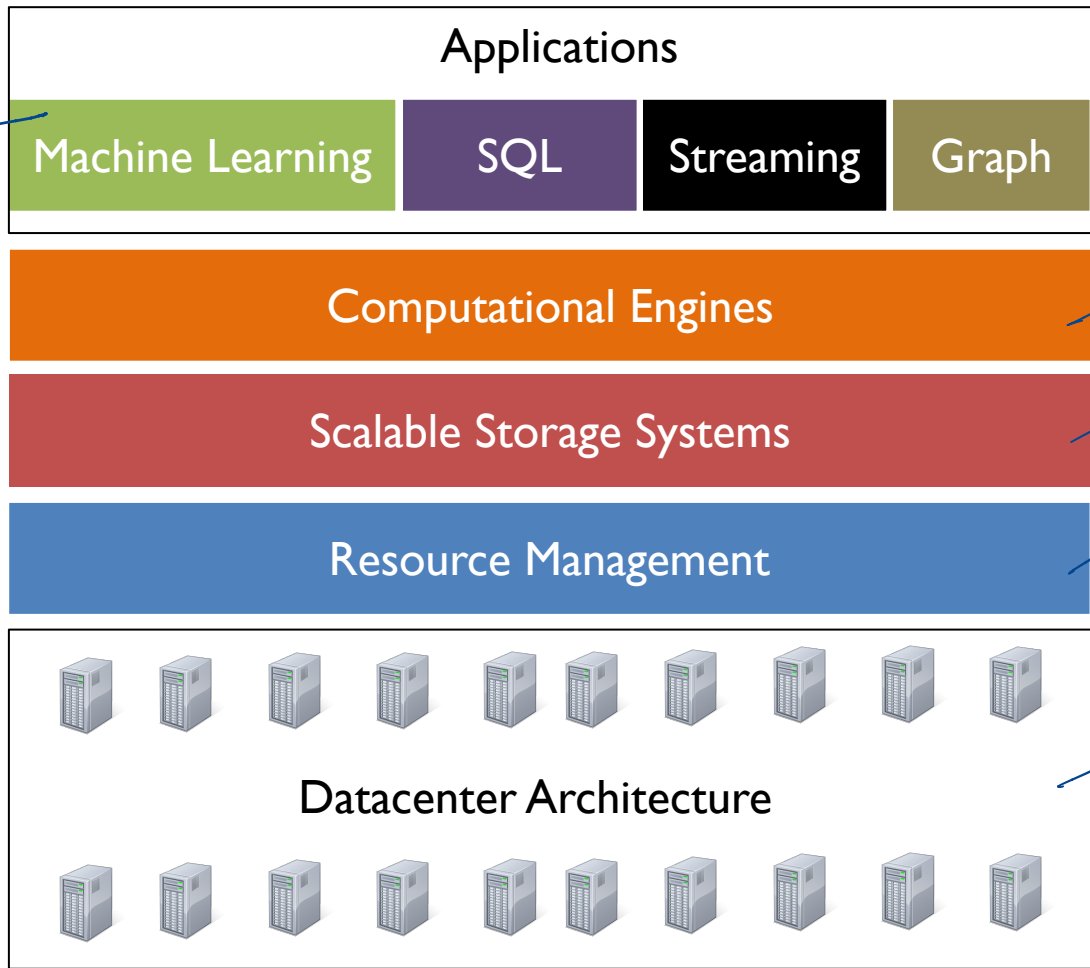Bid on topics, submit group (1 sentences) – ~~Oct 5~~ → Oct 1

Project Proposal (2 pages) – Oct 16

Piazza

Introduction

Related Work

Timeline (with eval plan)

# Applications

| Machine Learning | SQL | Streaming | Graph |
|:---:|:---:|:---:|:---:|

~ 2 - 3 weeks

## Computational Engines

→ Spark, MapReduce

## Scalable Storage Systems

→ GFS

## Resource Management

→ Mesos DRF

## Datacenter Architecture

# EMPIRICAL RISK MINIMIZATION

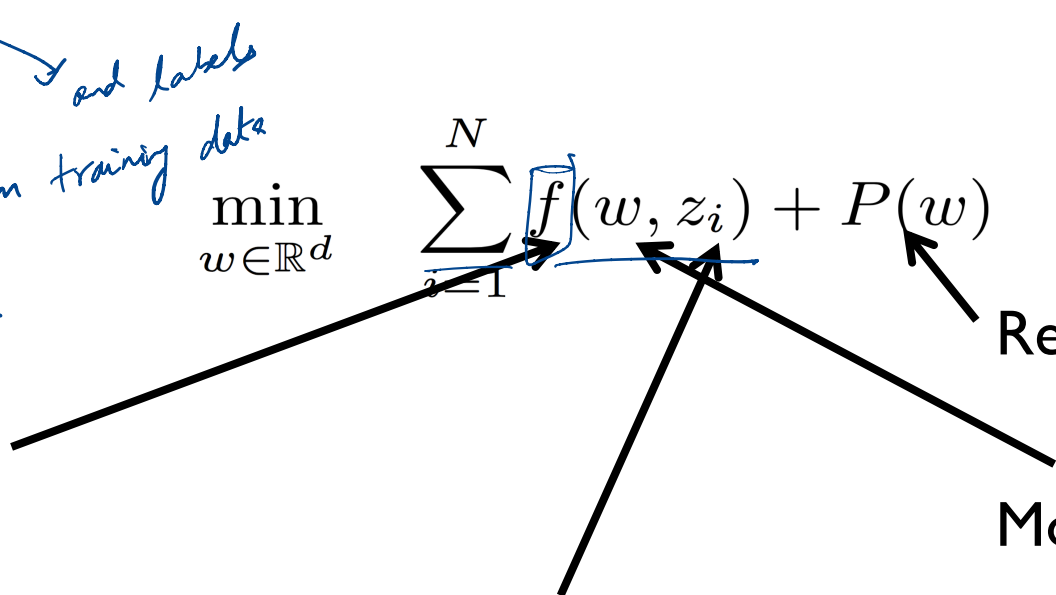Supervised learning → and labels

↳ Given training data

Fit a model

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^{N} f(w, z_i) + P(w)$$

Function

Regularization

Model

Data (Examples)

# DEEP LEARNING

→ Pytorch

weights of this layer

FC. 84 = [ 84 dim ]

[ input ]   * = [ ]

**ResNet18**

**Convolution**
**ReLU**
**MaxPool**
**Fully Connected**
**SoftMax**

con' pooling

layers of the model

composition of all layers



3d matrix width x height x 3 RGB image

layer = function ( (w_i); input)

weights or parameters of this layer

# STOCHASTIC GRADIENT DESCENT

Good fit for Spark !?

many iterations !!
iterative computation

update step

$$w^{(k+1)} = w^{(k)} - \alpha_k \nabla f(w^{(k)})$$

learning rate

gradient

Initialize w

For many iterations:

did it in parallel

Loss = Forward pass (model) → $\| f(w, input) - b \|_2$  square loss
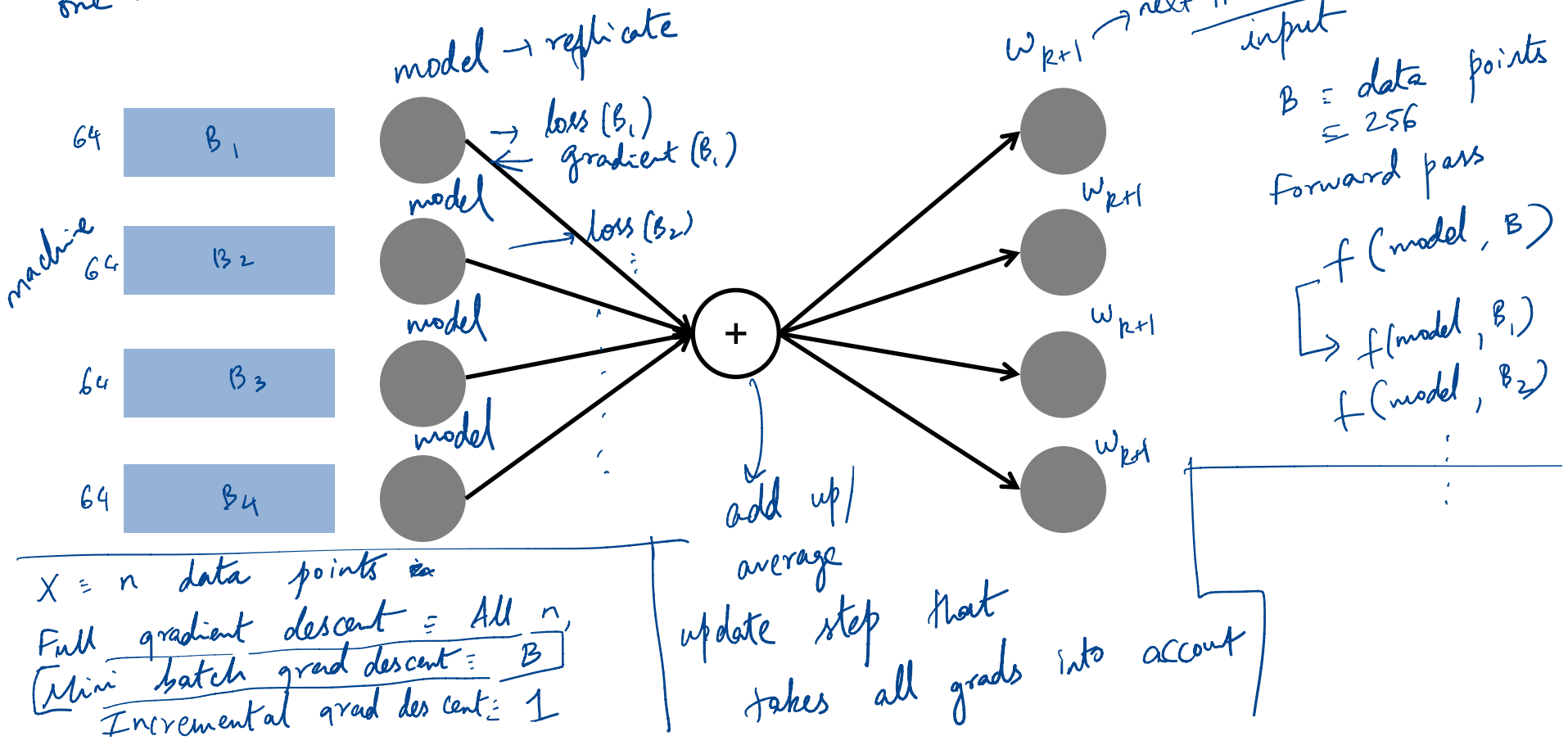
Gradient = backward (model) chain rule

Update model

End

model is shared → how do we parallelize

every iteration depends on previous

# DATA PARALLEL MODEL TRAINING

Parallelize one iteration

model → replicate

$w_{k+1}$ → next iteration input

64 | $B_1$

→ loss $(B_1)$
← gradient $(B_1)$

model

→ loss $(B_2)$

machine 64 | $B_2$

model

64 | $B_3$

model

64 | $B_4$

(+)

add up / average

update step that takes all grads into account

$w_{k+1}$

$w_{k+1}$

$w_{k+1}$

$w_{k+1}$

B = data points ≤ 256

forward pass

$f(model, B)$

→ $f(model, B_1)$

$f(model, B_2)$

X = n data points

Full gradient descent = All n,
(Mini) batch grad descent = B
Incremental grad descent = 1

# COLLECTIVE COMMUNICATION

$\bigcirc$ $g_1$
$\bigcirc$ $g_2$ $\Big\}$ reduce: $g_1 + g_2 + g_3 + g_4$
$\bigcirc$ $g_3$ $\to$ $\oplus$
$\bigcirc$ $g_4$

broadcast agg grads

Either broadcast agg grad · or · broadcast the model $\omega_{k+1}$

$\to$ MPI

send, recv $\equiv$ RPC
P2P primitives

## Broadcast, Scatter

## Gather, Reduce

MPI_Bcast

$\to$ data item

MPI_Gather (data_item, root)

$\to$ concate

opposite

vector [finely]

MPI_Scatter

MPI_Reduce

$[5, 2, 7, 4]$

MPI_SUM

$= 5 + 2 + 7 + 4$

From https://mpitutorial.com/tutorials/

# ALL REDUCE

MPI_Allreduce



From https://mpitutorial.com/tutorials/

# DISTRIBUTED DATA PARALLEL API

```
 9    # setup model and optimizer
10    net = nn.Linear(10, 10)
11    net = par.DistributedDataParallel(net)
12    opt = optim.SGD(net.parameters(), lr=0.01)
13
14    # run forward pass
15    inp = torch.randn(20, 10)
16    exp = torch.randn(20, 10)
17    out = net(inp)
18
19    # run backward pass
20    nn.MSELoss()(out, exp).backward()
21
22    # update parameters
23    opt.step()
```

*→ only line of code change*

*→ local model*

*— Non - intrusive*

*— Hooks to do optimizations in background*

# GRADIENT BUCKETING
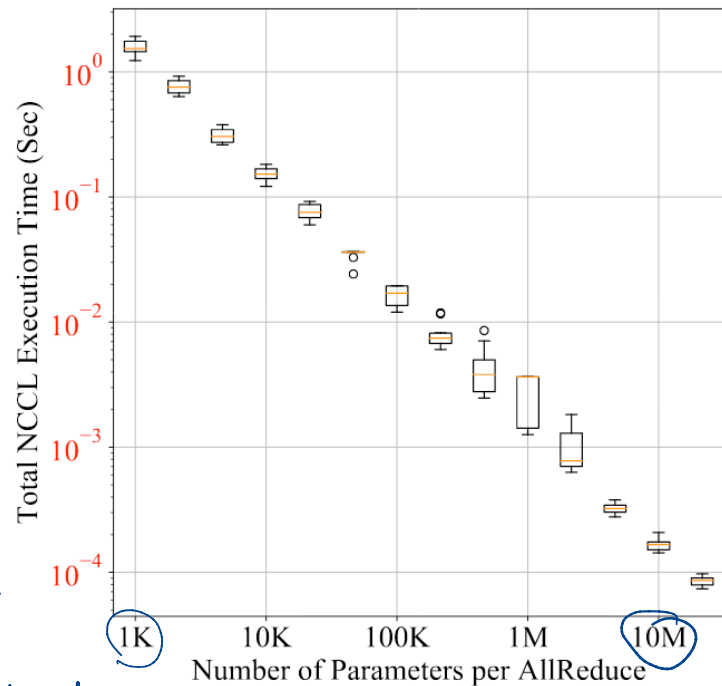
60 M parameter

## Why do we need gradient bucketing?

↳ Small tensor sizes
lead to greater time for
All Reduce

Every AllReduce
= (latency cost - fixed overhead) + bw cost how much data

→ Why not one big bucket
≡ Wait for all gradients to be ready
≡ Cannot overlap backward, AllReduce

# GRADIENT BUCKETING + ALL REDUCE



parameter ≡ layers

local model 1    DDP1    DDP2    local model 2

w₁  $g_{w1}$    $g_{b1}$  b₁    bucket2  $g_{w1}$  allreduce2  $g_{w1}$  bucket2  b₁  $g_{b1}$    $g_{w1}$  w₁

addmm1    $g_{b1}$    $g_{b1}$    addmm1

w₂  $g_{w2}$    $g_{b2}$  b₂    bucket1  $g_{w2}$  allreduce1  $g_{w2}$  bucket1  b₂  $g_{b2}$    $g_{w2}$  w₂

addmm2    $g_{b2}$    $g_{b2}$    addmm2

mse_loss    mse_loss

loss    loss

Process 1    Process 2

automatic gradient layers

As buckets become ready, we start All Reduce on them → async

In background, the gradient comp continues

Buckets are defined by size ≡ 25 MB

| Parameter | Gradient | → Autograd Edge | ⋯▶ Copy | ●—● Communication |

# GRADIENT ACCUMULATION

```
1  ddp = DistributedDataParallel(net)
2  with ddp.no_sync():
3      for inp, exp in zip(inputs, expected_outputs):
4          # no synchronization, accumulate grads
5          loss_fn(ddp(inp), exp).backward()
6  # synchronize grads
7  loss_fn(ddp(another_inp), another_exp).backward()
8  opt.step()
```

→ extra parameter

device can only hold smaller subset in memory

$B_7, B_4, B_1$   ☐ $B_1$        AllReduce

$B_8, B_5, B_2$   ☐ $B_2$     ⊕

$B_9, B_6, B_3$   ☐ $B_3$

no-sync

☐ [$B_1$  $B_4$  $B_7$]      AllReduce

☐ [$B_2$   $B_5$  $B_8$]   ⊕

☐ [$B_3$   $B_6$  $B_9$]

# IMPLEMENTATION

AllRedue 2  Port 1040

AllRedue1  Port 1234
→ kicked off

Port 1234  ① → ②
              ↓
③ ← ④

## Bucket_cap_mb ≡ Parameter that is tunable
small → overhead
large → no overlap

~ middle ≡ 25 MB

## Parameter-to-bucket mapping

1040 ↑  ① → ②
         ↓
    ③ ← ④

[layer 2] 20MB → bucket 1
[layer 3] 5MB

## Round-robin ProcessGroups

gradient
↳ math function
GPUs ≡ on a batch
CPUs ≡ data

[layer 4] → bucket 2 → 25MB
[layer] → bucket 3 → filled up 25MB

backward pass

# BREAKDOWN



Figure 6: Per Iteration Latency Breakdown

# SUMMARY

Pytorch: Framework for deep learning

DistributedDataParallel API

Gradient bucketing, AllReduce

Overlap computation and communication

# DISCUSSION

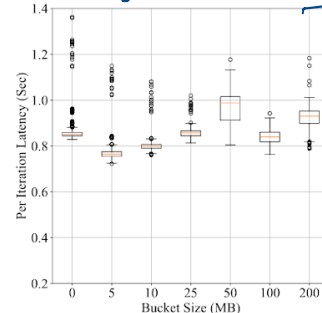https://forms.gle/6xhVBNBhdzsJ6gBE6
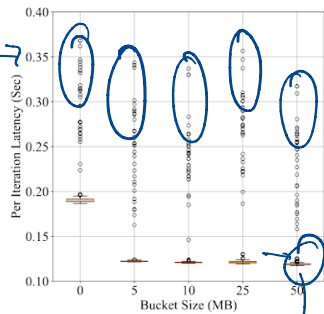
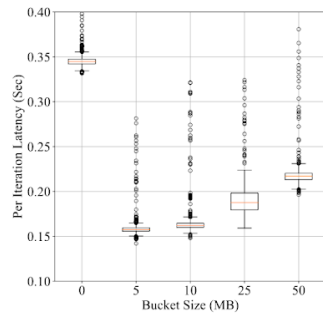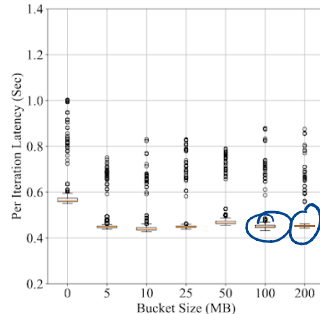(a) ResNet50 on NCCL  (b) ResNet50 on Gloo  (c) BERT on NCCL  (d) BERT on Gloo

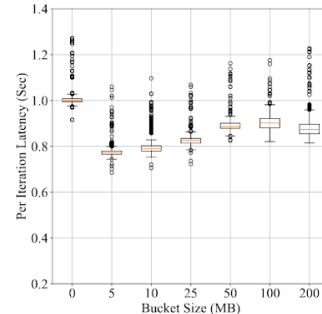**Figure 7: Per Iteration Latency vs Bucket Size on 16 GPUs**



(a) ResNet50 on NCCL  (b) ResNet50 on Gloo  (c) BERT on NCCL  (d) BERT on Gloo

**Figure 8: Per Iteration Latency vs Bucket Size on 32 GPUs**

Scales well !?

Strong scaling

$B = 256$
increase num GPUs
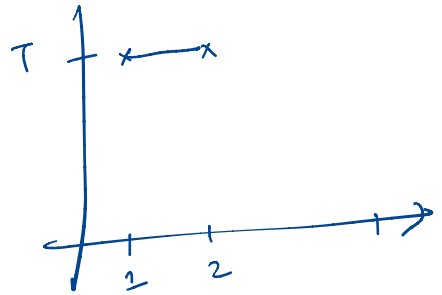
T
T/2

1   2

Weak scaling
$B_i = 64$, increase num GPUs

T

1   2

# What could be some challenges in implementing similar optimizations for AllReduce in Apache Spark?

spark = "larger workloads" ?

Each worker node on spark had dataset

    ↳ Spark needs to shuffle operation

||
more expensive than NCCL Ring reduce

○ $\nabla g_1$
○ $\nabla g_2$
○ $\nabla g_3$ — Tree Reduce using shuffles
○ $\nabla g_4$

Overlap compute / communications | Not all tasks active at same time
Task completes → shuffle

# NEXT STEPS

Next class: PipeDream

Assignment 2 is due soon!

Project Proposal

Groups by Oct 1

2 pager by Oct 16

Alloc bucket layer0 layer0 copy [- -]

2.5MB [- - -]

User program scatter

Process Group API

scatter

NCCL          Gloo

network monitoring

MB/s

time

extra signal    bitmaps & other info !?

?? P1 : 0 link
P2 : 0 — ⊕ rank0
0
0

$[\nabla_{g_1}\ \nabla_{g_2}\ \nabla_{g_3}]$

which link

input ≡ batches

layer 1 $\nabla_{g_1}$
$[\nabla_{g_2}]$ comm
layer 2
$[\nabla_{g_3}]$ ≡ comm
layer 3
loss

Reduce → brings to single machine
5
2  → + 14 ← bcast
7