

CS 744: PYTORCH

Shivaram Venkataraman

Spring 2024

ADMINISTRIVIA

Assignment 2 out! Due **Feb 23rd 10PM!**

Course Project

Topics list posted – **Feb 21st**

Propose / Bid on topics, submit group (1 sentence) – **Feb 26th**

Title confirmed – **March 1**

Project Proposal (2 pages) – **March 8**

Introduction

Related Work

Timeline (with eval plan)

Applications

Machine Learning

SQL

Streaming

Graph

Computational Engines

Scalable Storage Systems

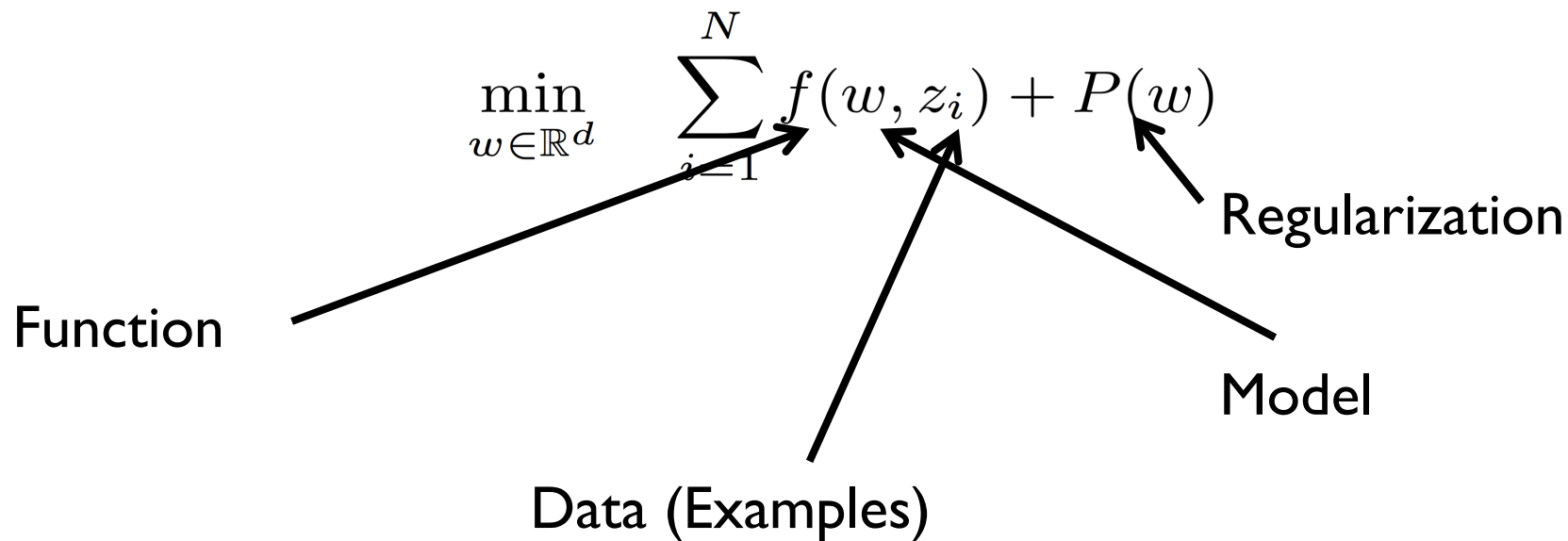
Resource Management



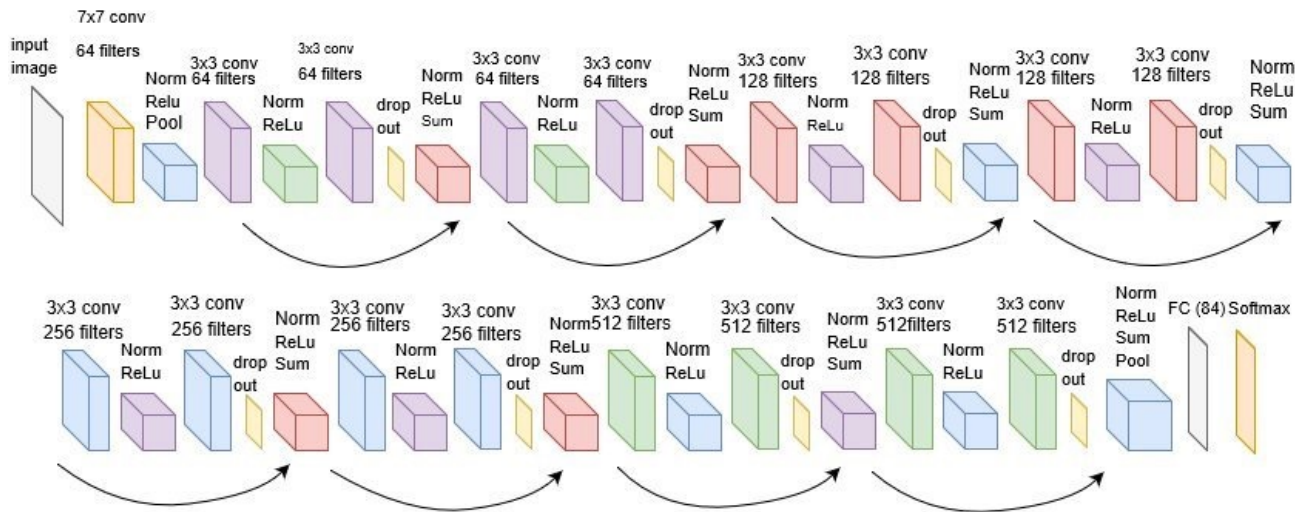
Datacenter Architecture



EMPIRICAL RISK MINIMIZATION



DEEP LEARNING



ResNet18

Convolution
ReLU
MaxPool
Fully Connected
SoftMax

STOCHASTIC GRADIENT DESCENT

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

Initialize w

For many iterations:

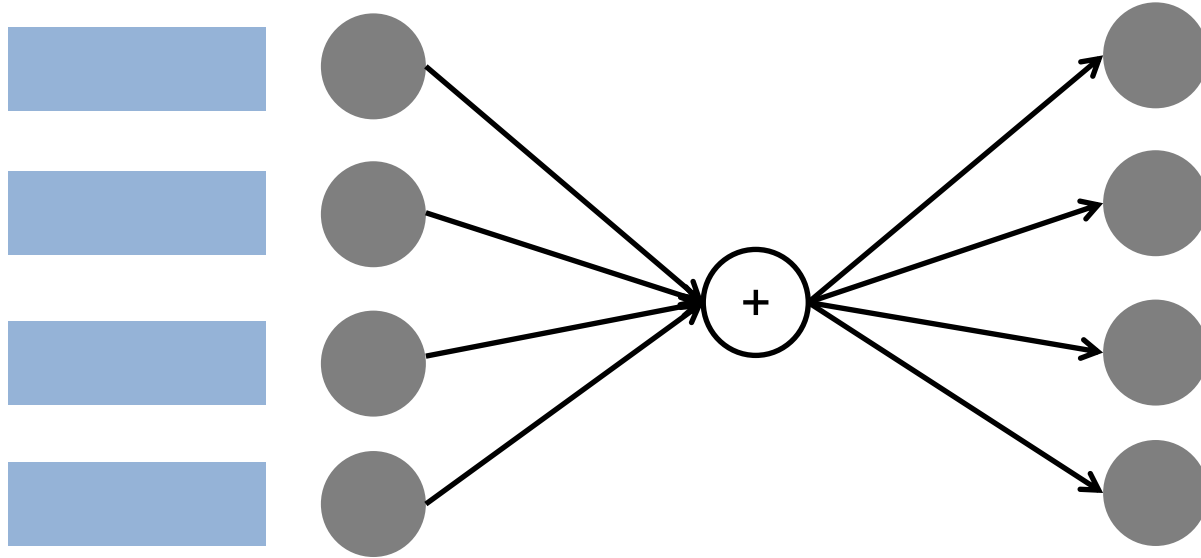
Loss = Forward pass

Gradient = backward

Update model

End

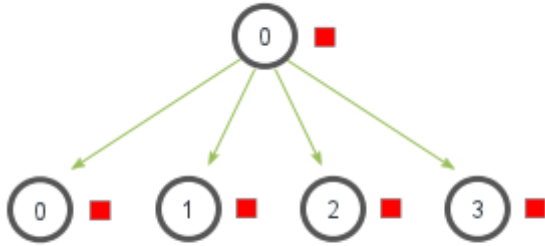
DATA PARALLEL MODEL TRAINING



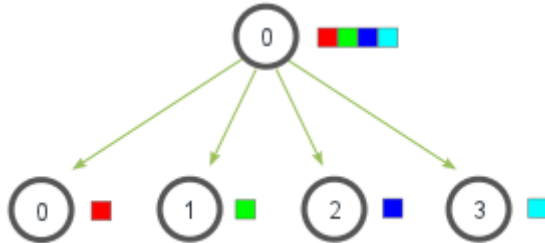
COLLECTIVE COMMUNICATION

Broadcast, Scatter

MPI_Bcast

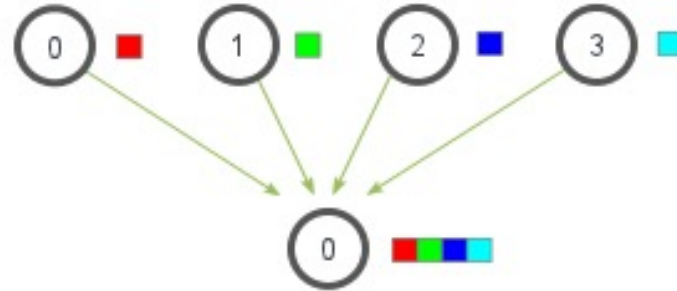


MPI_Scatter

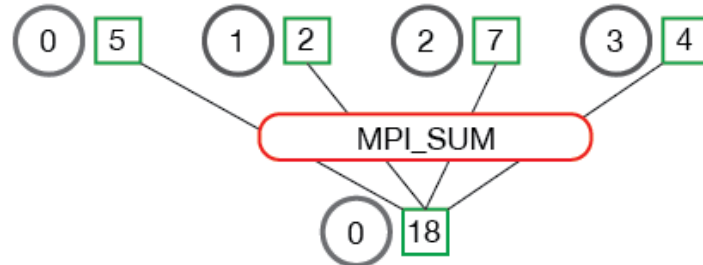


Gather, Reduce

MPI_Gather

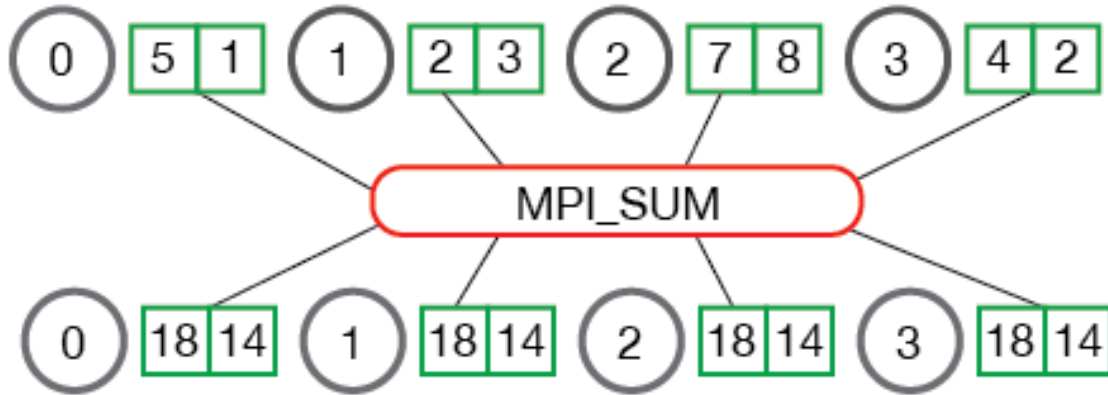


MPI_Reduce



ALL REDUCE USING A RING

MPI_Allreduce

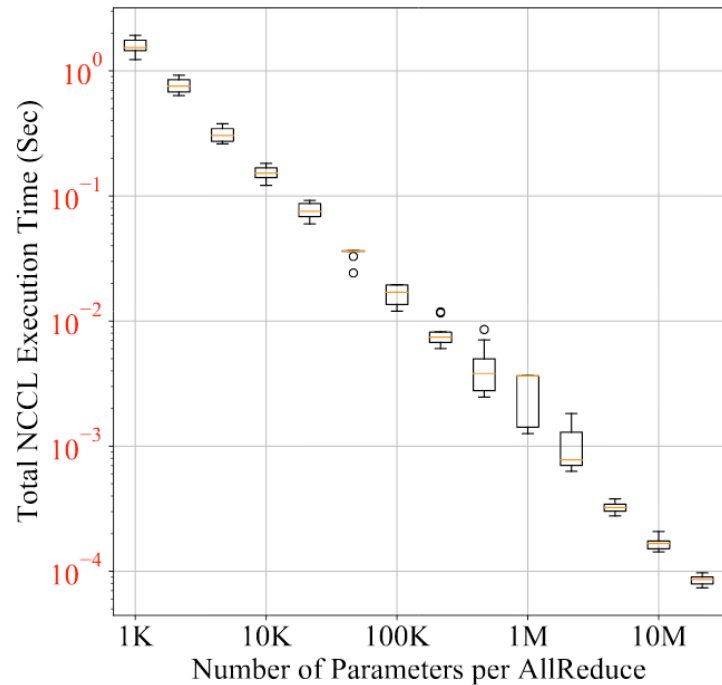


DISTRIBUTED DATA PARALLEL API

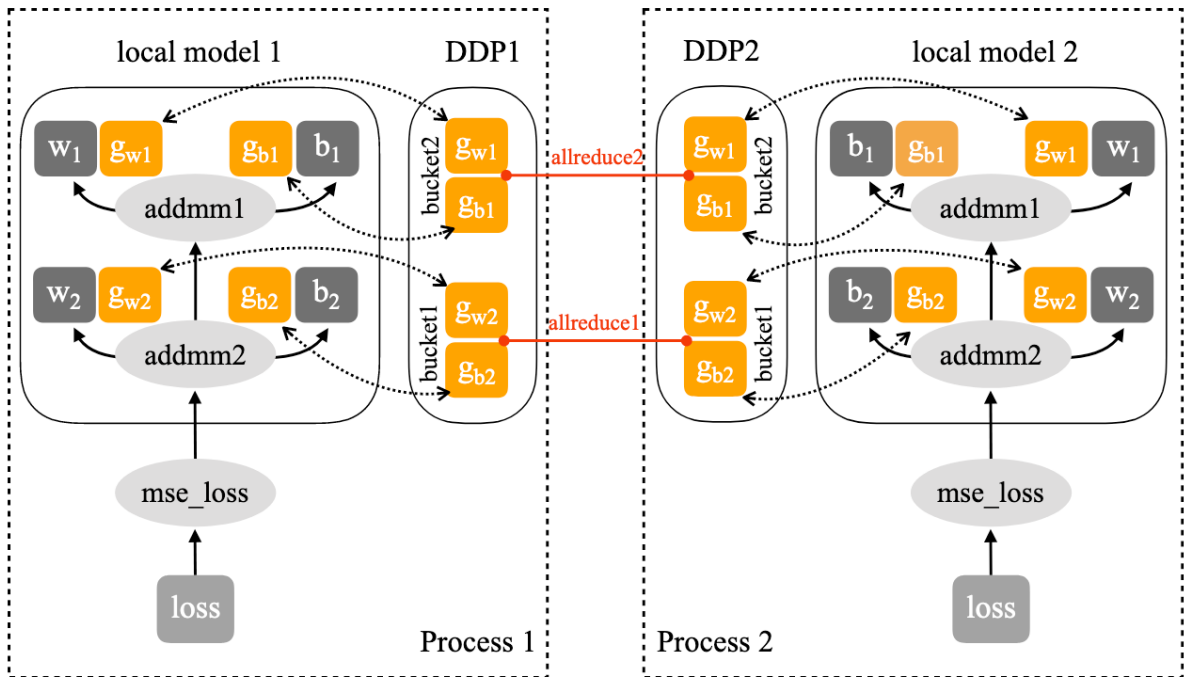
```
9  # setup model and optimizer
10 net = nn.Linear(10, 10)
11 net = par.DistributedDataParallel(net)
12 opt = optim.SGD(net.parameters(), lr=0.01)
13
14 # run forward pass
15 inp = torch.randn(20, 10)
16 exp = torch.randn(20, 10)
17 out = net(inp)
18
19 # run backward pass
20 nn.MSELoss()(out, exp).backward()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

GRADIENT BUCKETING

Why do we need gradient bucketing?



GRADIENT BUCKETING + ALL REDUCE



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()
```

IMPLEMENTATION

Bucket_cap_mb

Parameter-to-bucket mapping

Round-robin ProcessGroups

SUMMARY

Pytorch: Framework for deep learning

DistributedDataParallel API

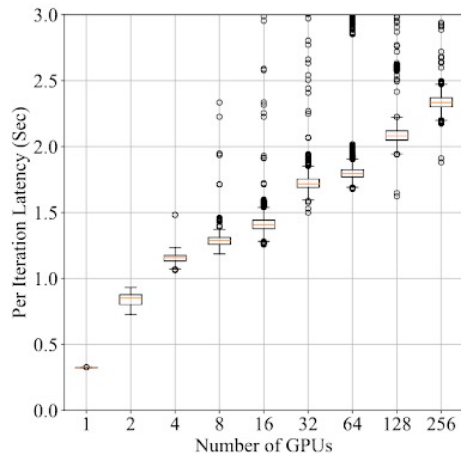
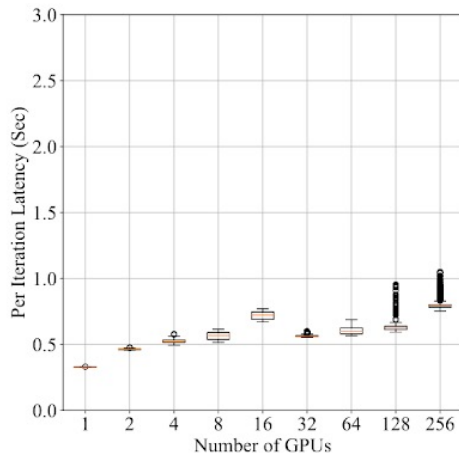
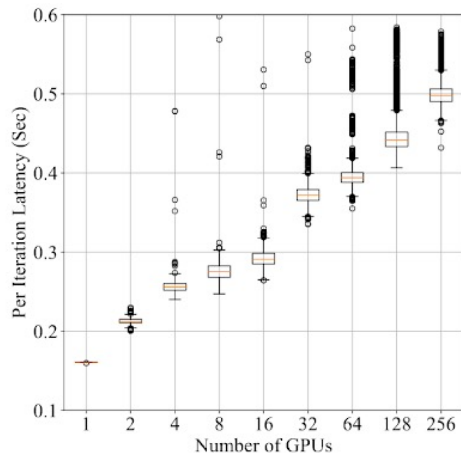
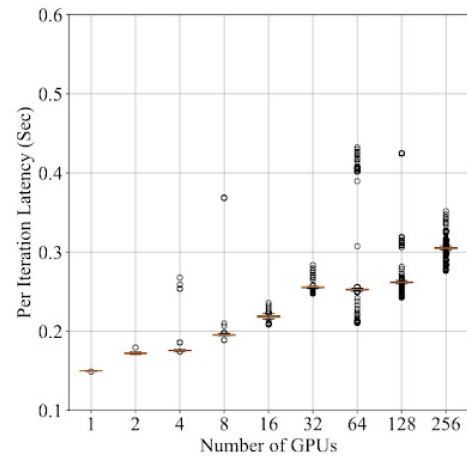
Gradient bucketing, AllReduce

Overlap computation and communication



DISCUSSION

<https://forms.gle/aUFy5fsN8KMS4Li6>



(a) ResNet50 on NCCL

(b) ResNet50 on Gloo

(c) BERT on NCCL

(d) BERT on Gloo

Figure 9: Scalability

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

NEXT STEPS

Next class: PipeDream

Assignment 2 is out!

BREAKDOWN

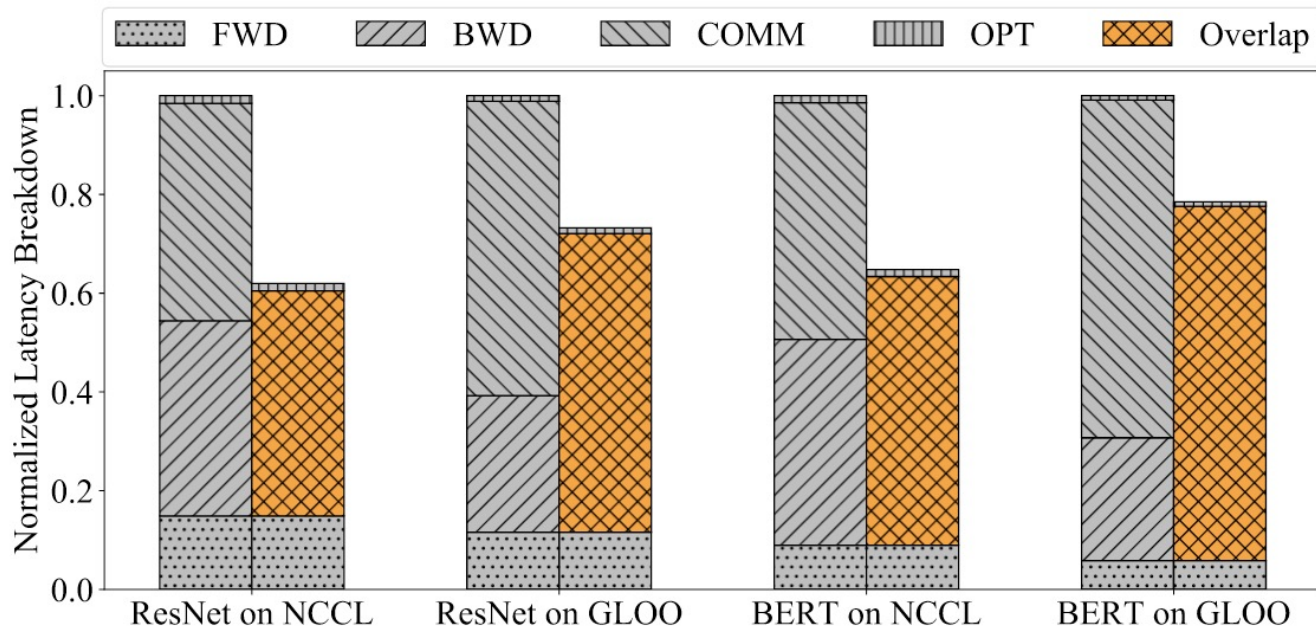


Figure 6: Per Iteration Latency Breakdown