

Good morning!!

CS 744: PYTORCH

Shivaram Venkataraman

Spring 2024

ADMINISTRIVIA

Assignment 2 out! Due **Feb 23rd 10PM!**
↳ Friday next week

Course Project

Topics list posted – **Feb 21st**

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

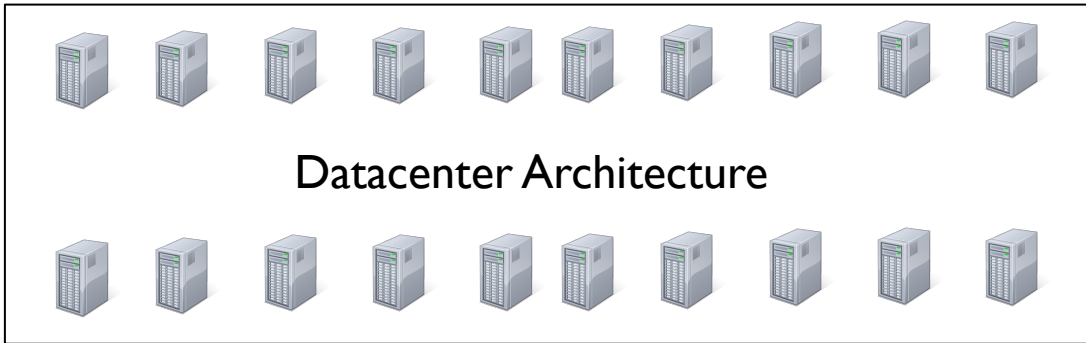
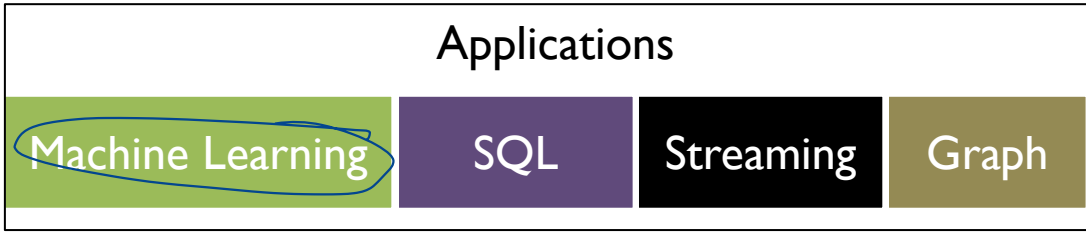
Title confirmed – **March 1**

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

Introduction

Related Work

Timeline (with eval plan)



Spark
MapReduce →

→ GFS

EMPIRICAL RISK MINIMIZATION

Supervised learning

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

loss function

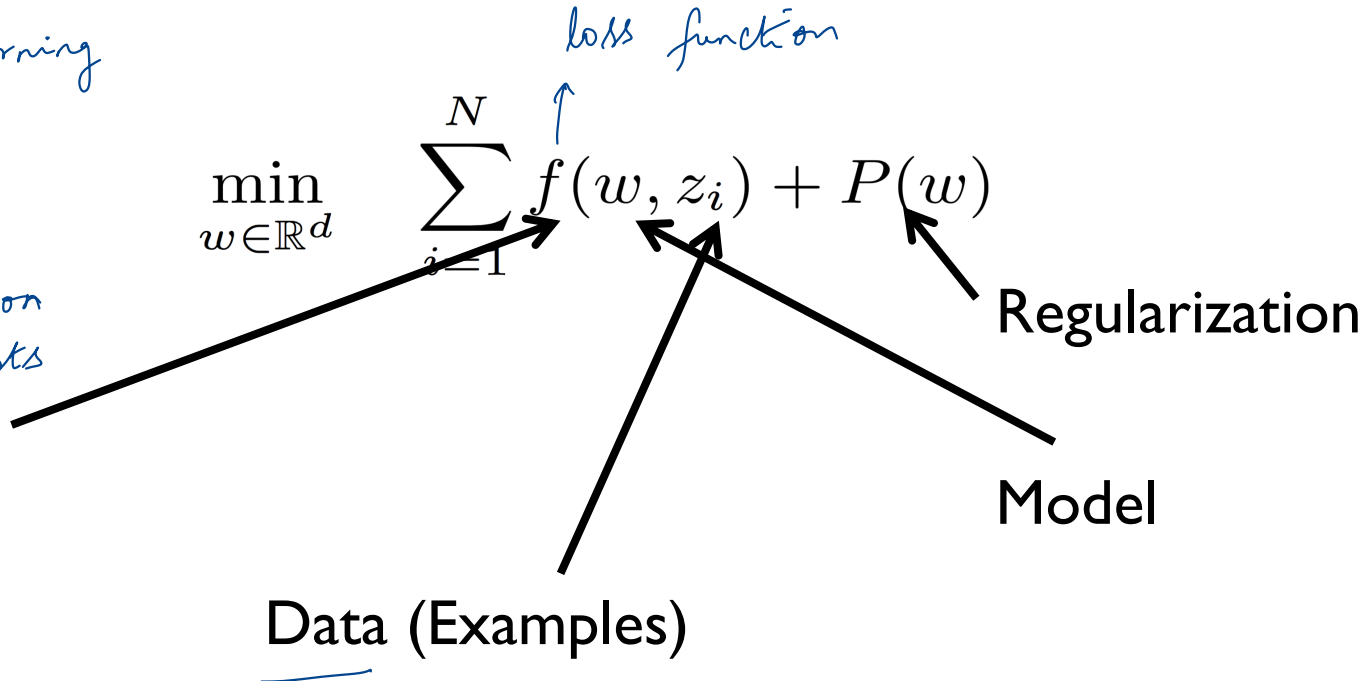
*Logistic Regression
Random Forests*

Function

Regularization

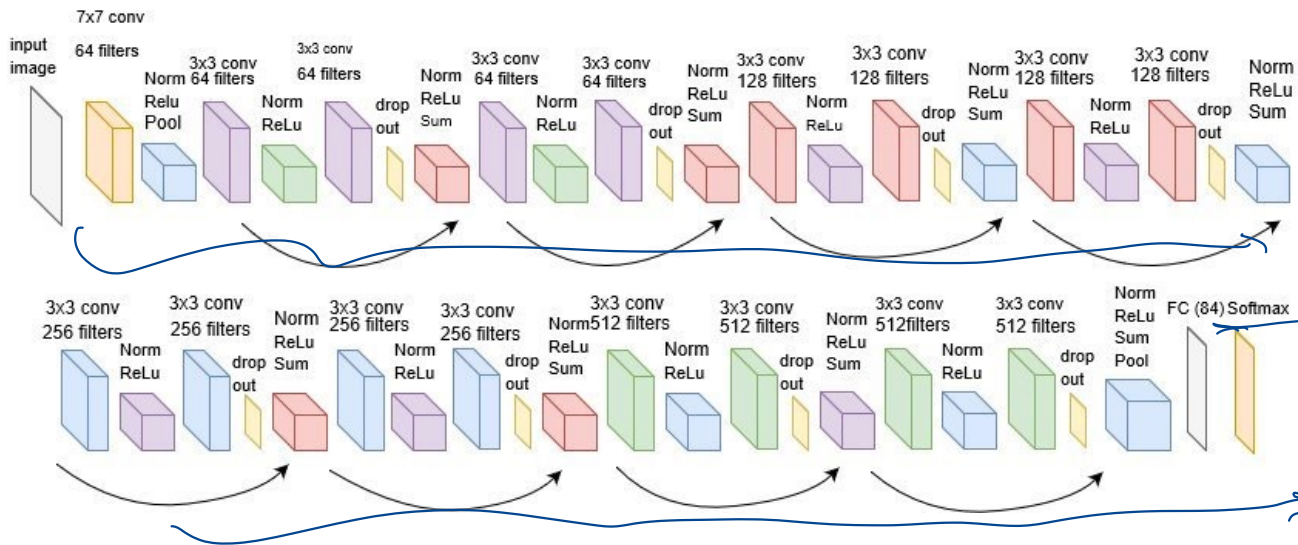
Model

Data (Examples)



DEEP LEARNING

data items



prediction

ResNet18

- Convolution
- ReLU
- MaxPool
- Fully Connected
- SoftMax

STOCHASTIC GRADIENT DESCENT

Adem

Tensors
matrices

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

↳ step size

Initialize w → randomly

For many iterations:

Loss = Forward pass

Gradient = backward

Update model

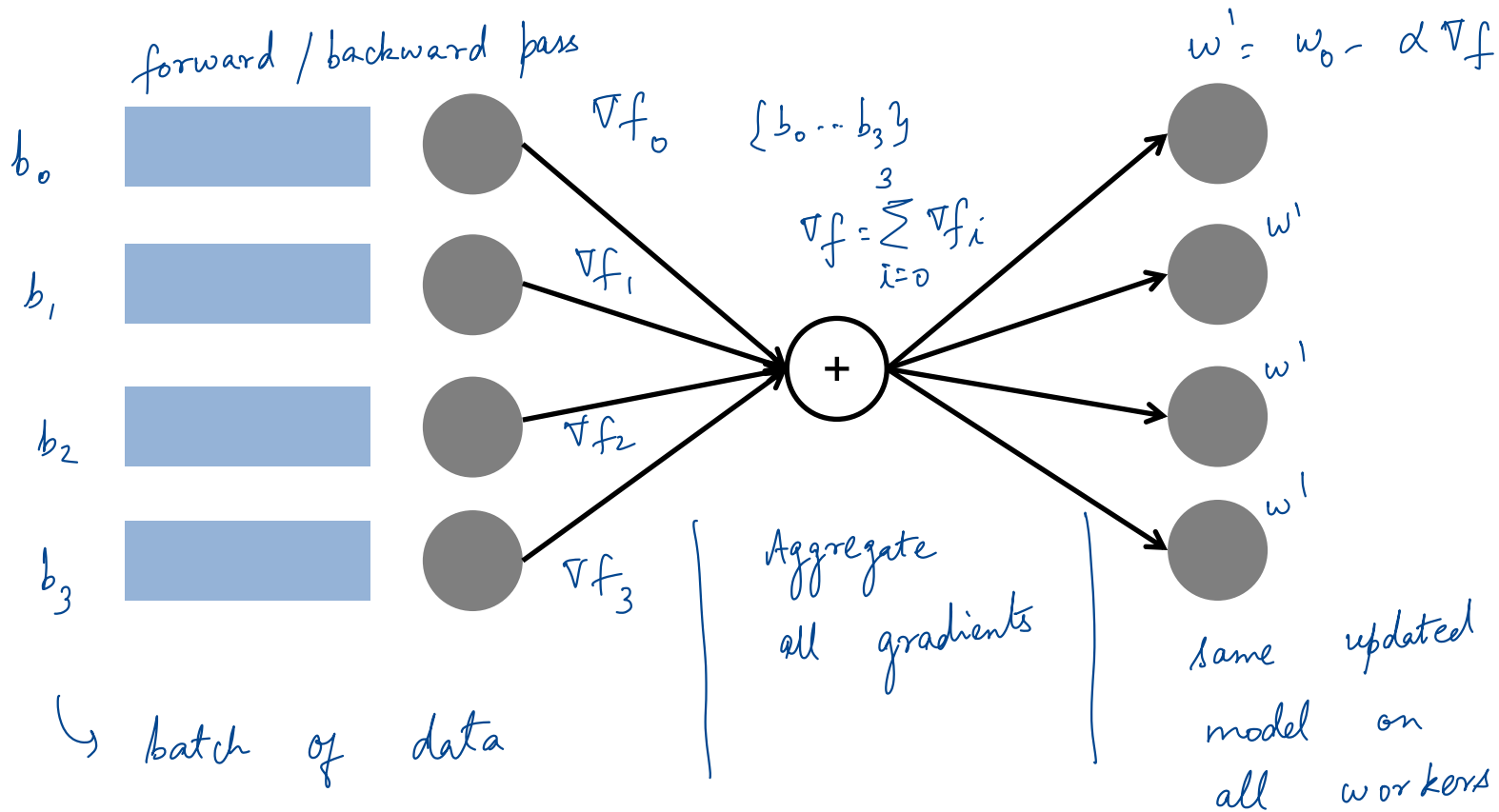
End

automatic
differentiation

Iterative
repeating
computation
as model
changes

iteration batch of data
on 128 1M

DATA PARALLEL MODEL TRAINING



COLLECTIVE COMMUNICATION

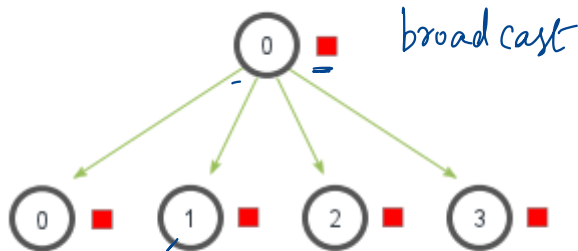


MPI_Bcast

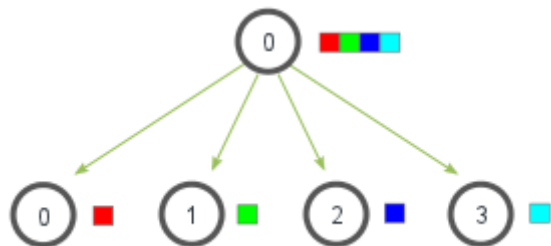
optimized

Broadcast, Scatter

MPI_Bcast

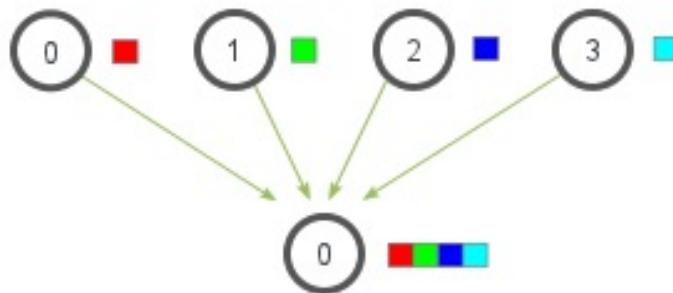


MPI_Scatter

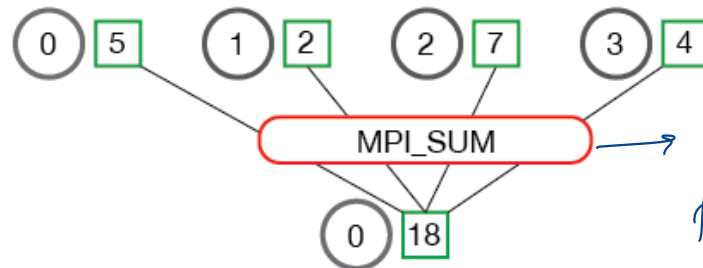


Gather, Reduce

MPI_Gather

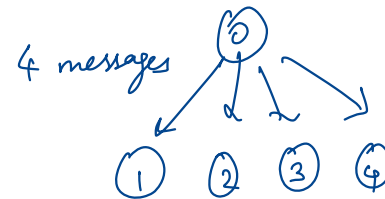
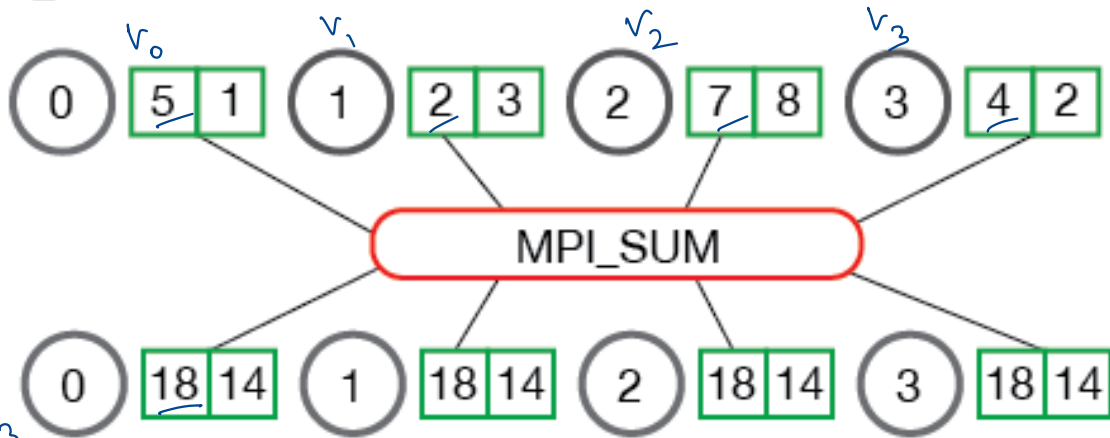


MPI_Reduce



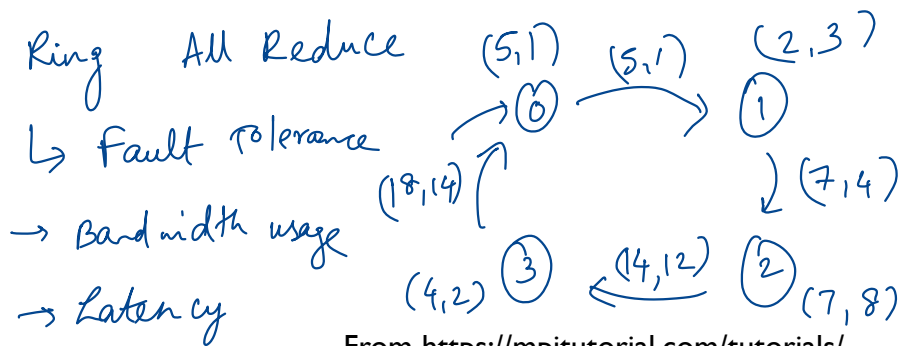
ALL REDUCE USING A RING

MPI_Allreduce



reduction
+
broadcast

$$\sum_{i=0}^3 v_i$$



1 round of the ring, you get the final value at Process 0.

2nd round broadcast

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()
```

magic

*Minimize
code changes
go from
single machine
to
Distributed
training*

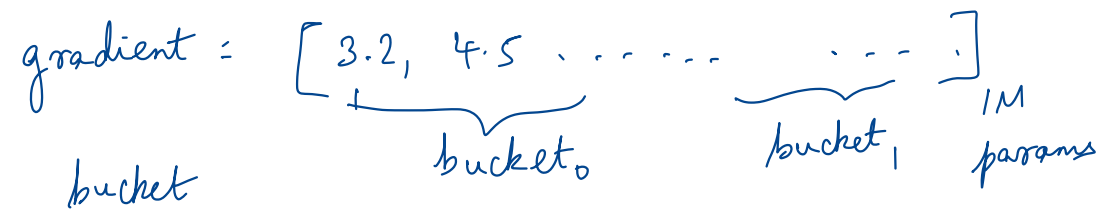
$$T = \alpha T_{\text{latency}} + \beta T_{\text{bandwidth}}$$

GRADIENT BUCKETING

1 bucket \rightarrow miss out on overlap

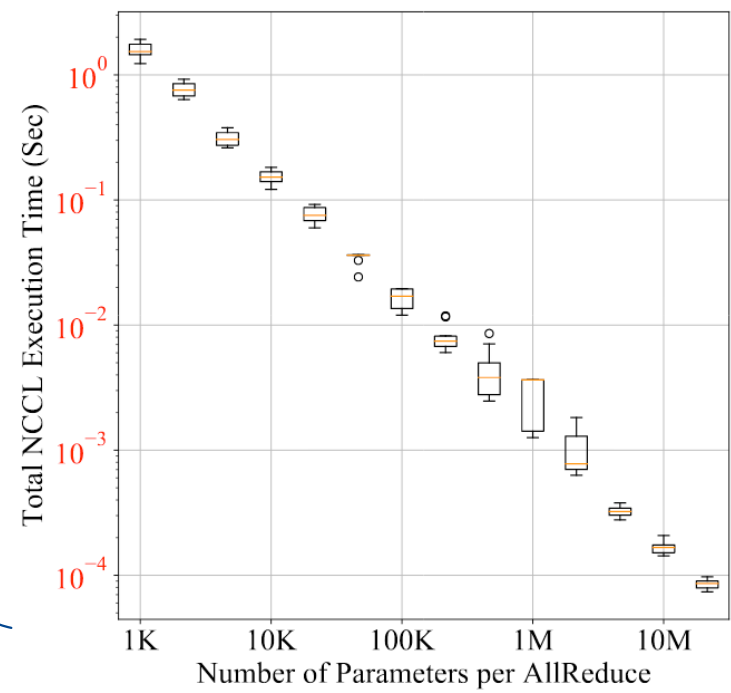
1M buckets \rightarrow overhead!!

Why do we need gradient bucketing?



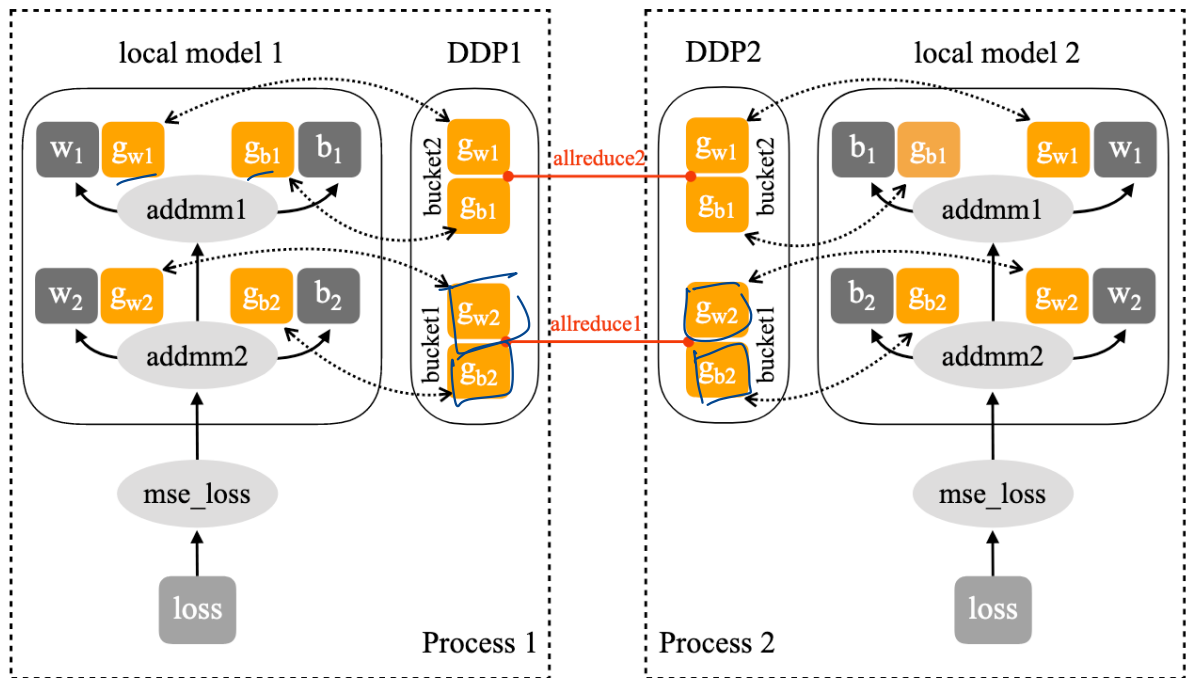
all reduce on gradient
 \hookrightarrow all reduce (bucket₀) all reduce (bucket₁)

every bucket is independent
 \hookrightarrow if bucket is ready \rightarrow trigger all reduce



GRADIENT BUCKETING + ALL REDUCE

which gradients in which bucket



■ Parameter ■ Gradient → Autograd Edge Copy ●-● Communication

fwd / backward
- gradients for final layers are available earliest

overlap ∇ calculation for first bucket with allreduce of second bucket

GRADIENT ACCUMULATION

large batch size
where

```
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()
```

no all-reduce
operations
in background

IMPLEMENTATION

Bucket_cap_mb → 25 MB is a good default

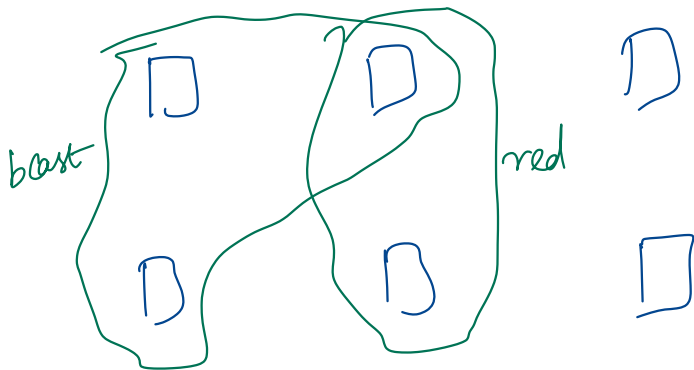
Parameter-to-bucket mapping → walk backwards using
model.parameters()

Round-robin ProcessGroups

↳ Multiple link types

→ NVlink

→ PCIe



SUMMARY

Pytorch: Framework for deep learning

DistributedDataParallel API

Gradient bucketing, AllReduce

Overlap computation and communication

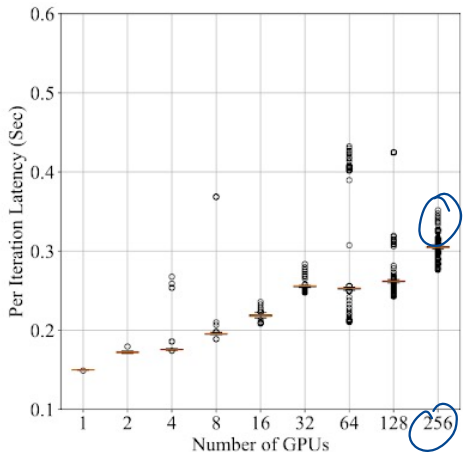


DISCUSSION

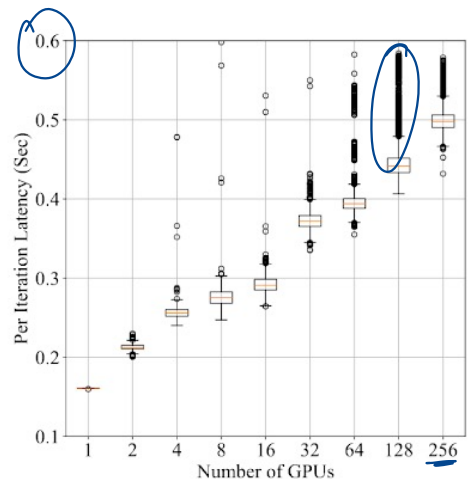
<https://forms.gle/aUFy5fsN8KMS4Li6>

More GPUs \rightarrow more time.
 \rightarrow linear vs. exponential?

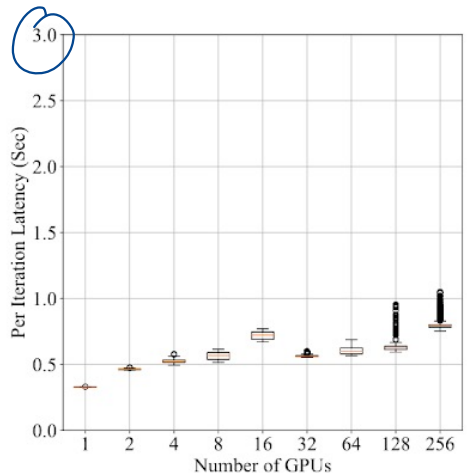
Variance for Gloo is much higher
 at 128/256 GPUs



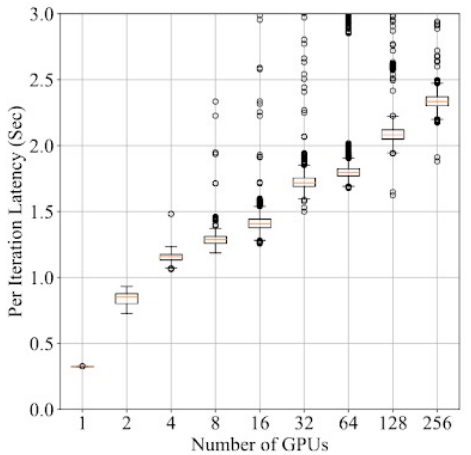
(a) ResNet50 on NCCL



(b) ResNet50 on Gloo



(c) BERT on NCCL



(d) BERT on Gloo

Figure 9: Scalability

BERT has higher latency
 \rightarrow larger model than Resnet-50

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

↳ overlap all Reduce with grad compute

- Fault Tolerance difference

↳ MPI-style vs. Spark

AllReduce Spark

Reduction tree + broadcast

- Barrier between map and reduce stages prevents overlap

map



- Slice the ~~data~~ gradient ~~is~~ doesn't fit very well with spark API

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

Spark

↳ Tasks → Not all tasks are active at the same time

0

1

2

3

assuming all are active at some time!

NEXT STEPS

Next class: PipeDream

Assignment 2 is out!

BREAKDOWN

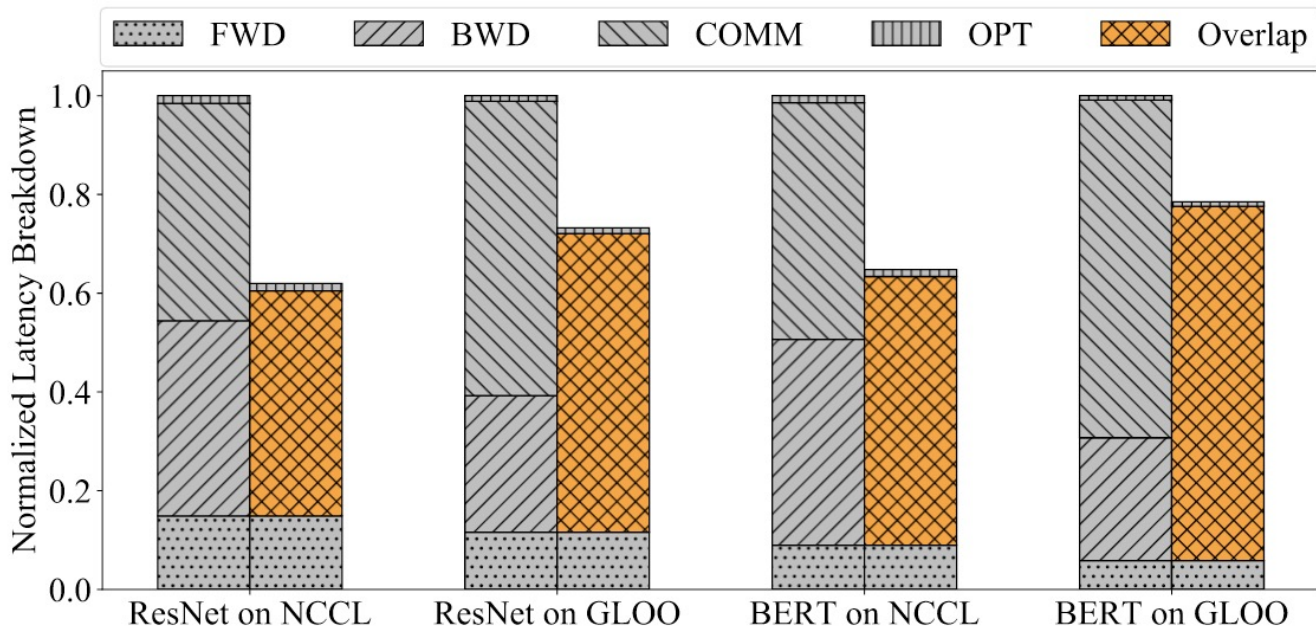


Figure 6: Per Iteration Latency Breakdown