

Hello!

CS 744: PYTORCH

Shivaram Venkataraman

Fall 2022

ADMINISTRIVIA

- Assignment 2 out! Due **Oct 12th 11 AM!** → *Maybe start soon?*
- Paper review deadline – 12:59pm on Tue/Thu → *Survey in one or two weeks*
- Bid on topics, submit group (1 sentences) – Oct 11
- Title confirmed – Oct 15
- Project Proposal (2 pages) – **Oct 25**
- Introduction
 - Related Work
 - Timeline (with eval plan)

Apache Spark

WRITING AN INTRODUCTION

→ 1 or 2 page

preface to your
paper

→ interactive data analytics

I-2 paras: what is the problem you are solving

why is it important (need citations)

→ Many industries
need this

I-2 paras: How other people solve and why they fall short

↳ MapReduce jobs and these are slow

I-2 paras: How do you plan on solving it and why your approach is better

I para: Anticipated results or what experiments you will use

→ Project proposal: plan | Final report: approach you took

RELATED WORK, EVAL PLAN

Group related work into 2 or 3 buckets (1-2 para per bucket)

Explain what the papers / projects do
Why are they different / insufficient] 8-10 related papers or projects

Eval Plan → what resources will help you succeed?

Describe what datasets, hardware you will use

Available: Cloudlab, Google Cloud (~\$150), Jetson TX2 etc.

Applications

Machine Learning

SQL

Streaming

Graph

Computational Engines

Scalable Storage Systems

Resource Management



Datacenter Architecture



*MapReduce
Spark*

GFS

Mesos / DRF

Hardware

EMPIRICAL RISK MINIMIZATION

Supervised
learning

Training data, labels

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

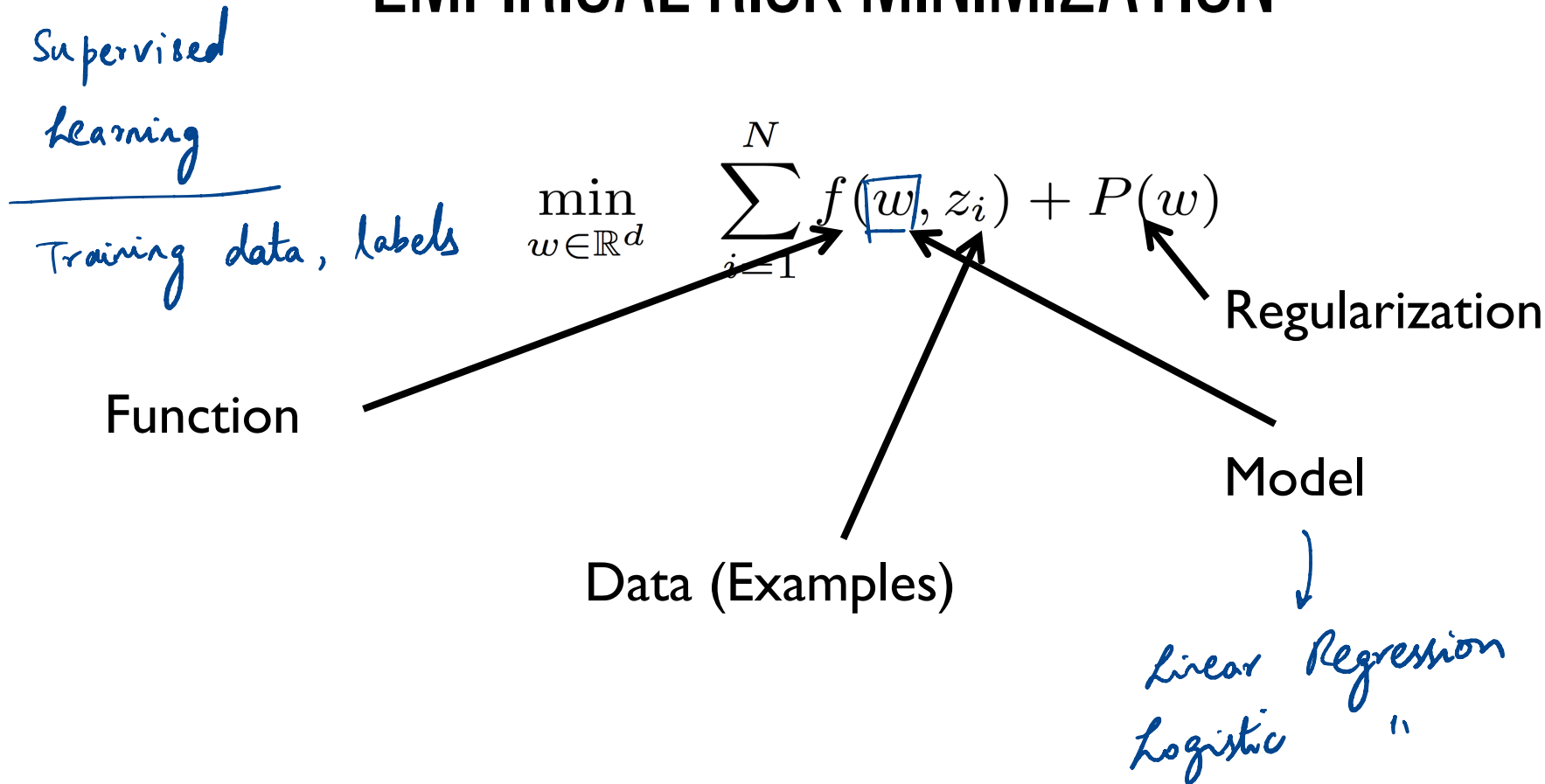
Function

Data (Examples)

Regularization

Model

Linear Regression
Logistic " "



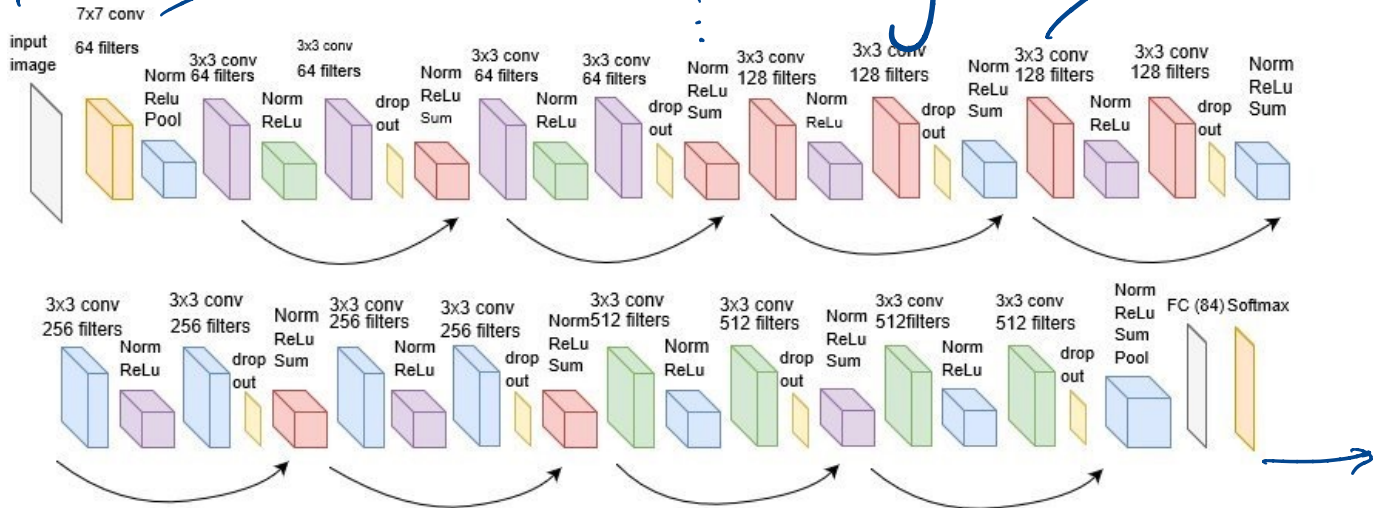
DEEP LEARNING

layers

$t_1 = l_1(\text{input})$
 $t_2 = l_2(t_1)$
 \vdots

forward pass

weights parameters



ResNet18

- Convolution
- ReLU
- MaxPool
- Fully Connected
- SoftMax

classify images or other tasks

STOCHASTIC GRADIENT DESCENT

model at iter $k+1$

model at iter k

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

SAD

Initialize w ← model initialization

For many iterations: → Batch of data

Loss = Forward pass

Gradient = backward

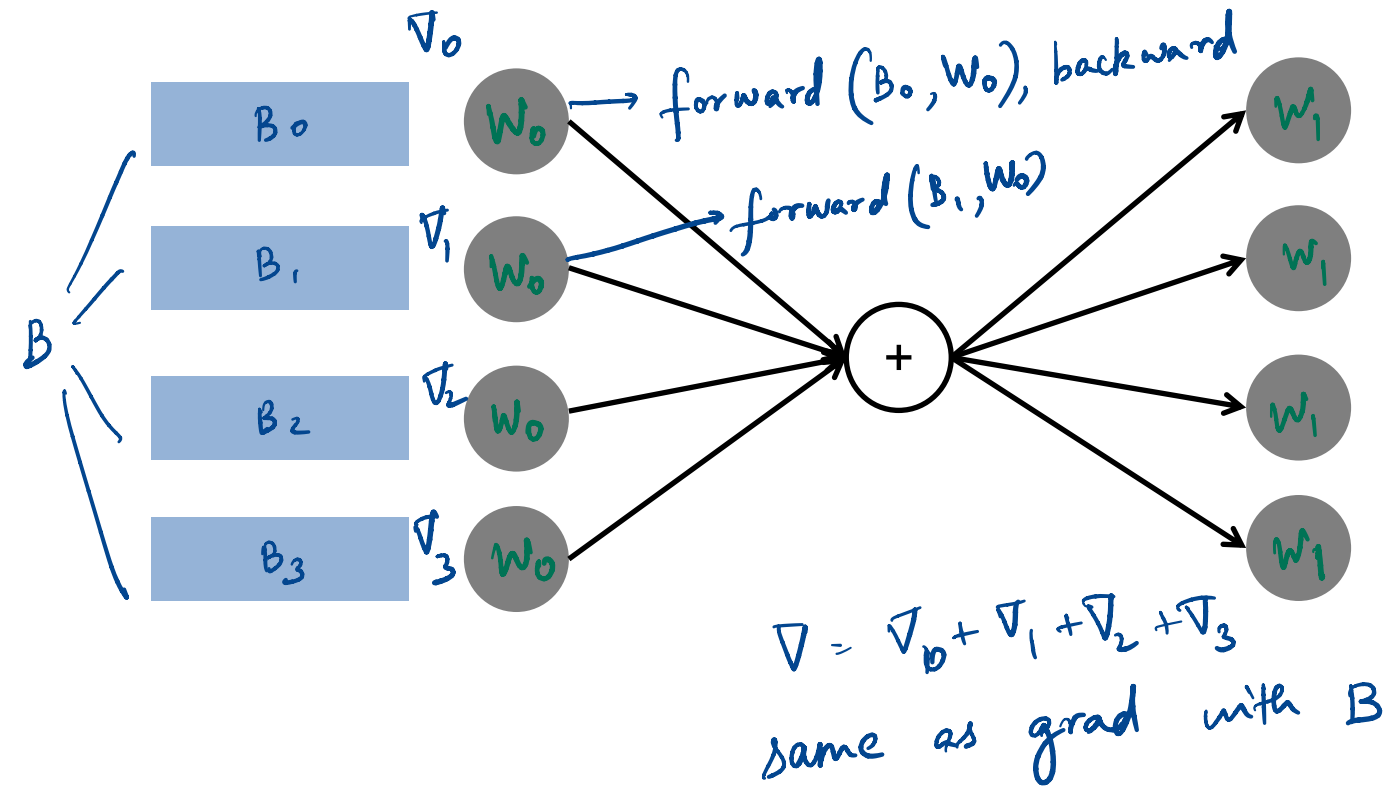
Update model

End

input and return "loss"
↳ how good are your predictions

Chain rule
differentiate
sequence of functions

DATA PARALLEL MODEL TRAINING

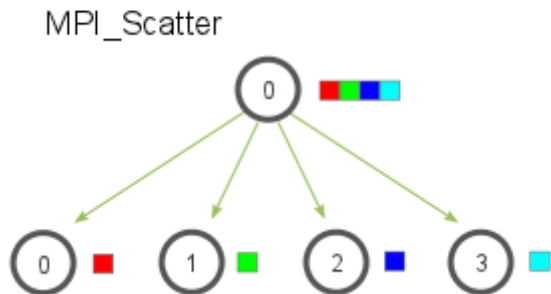
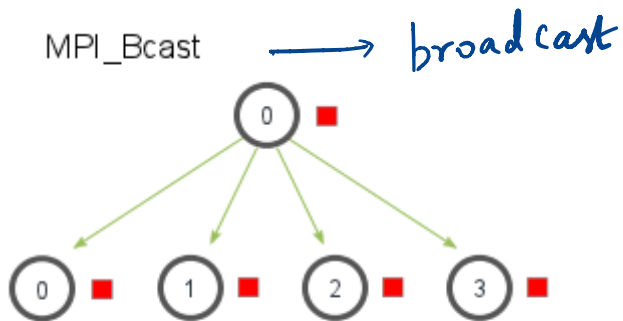


Mimic or reproduce behavior of single machine training!

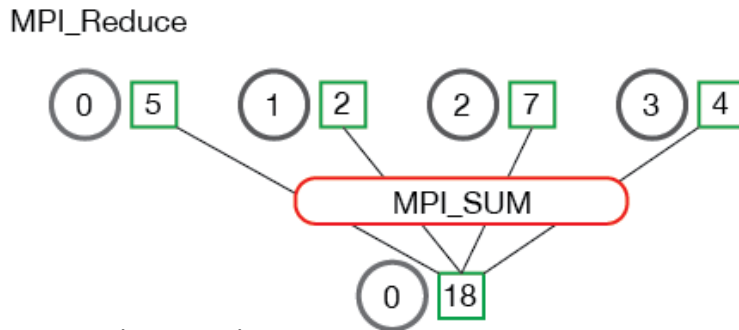
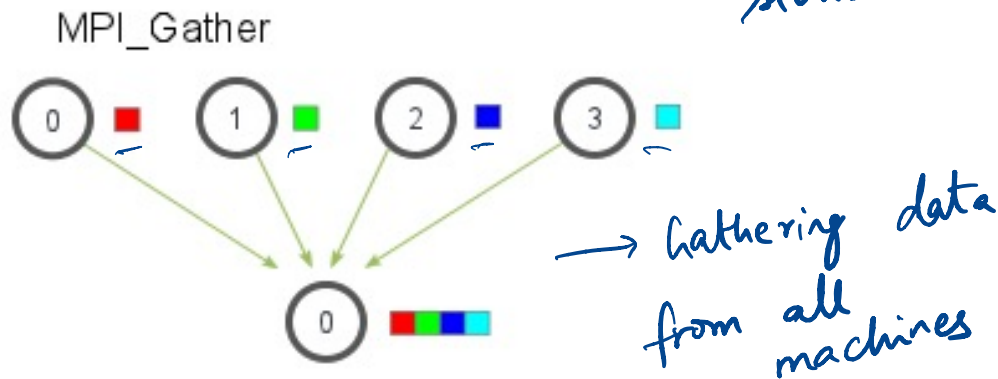
COLLECTIVE COMMUNICATION

→ 1991
first
MPI
standard

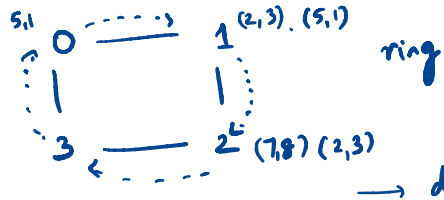
Broadcast, Scatter



Gather, Reduce

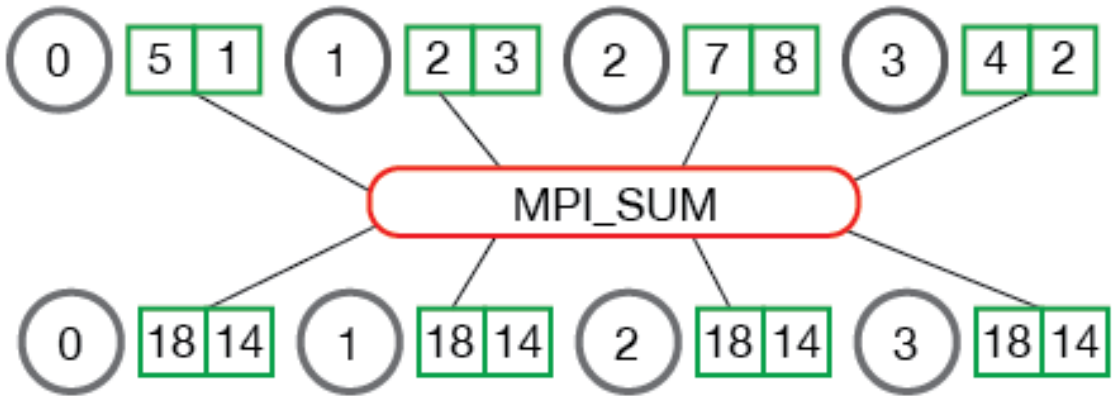


ALL REDUCE USING A RING



MPI_Allreduce

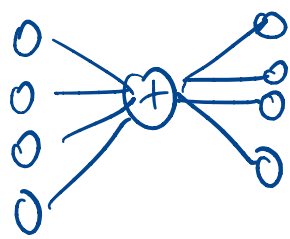
→ data transmitted per worker
→ number of connections



Reduce ≡ Aggregate data from all workers

Agg data + Broadcast agg value to all workers

reduce + broadcast



bottleneck from one machine

- All nodes broadcast to other nodes. Aggregate locally

DISTRIBUTED DATA PARALLEL API

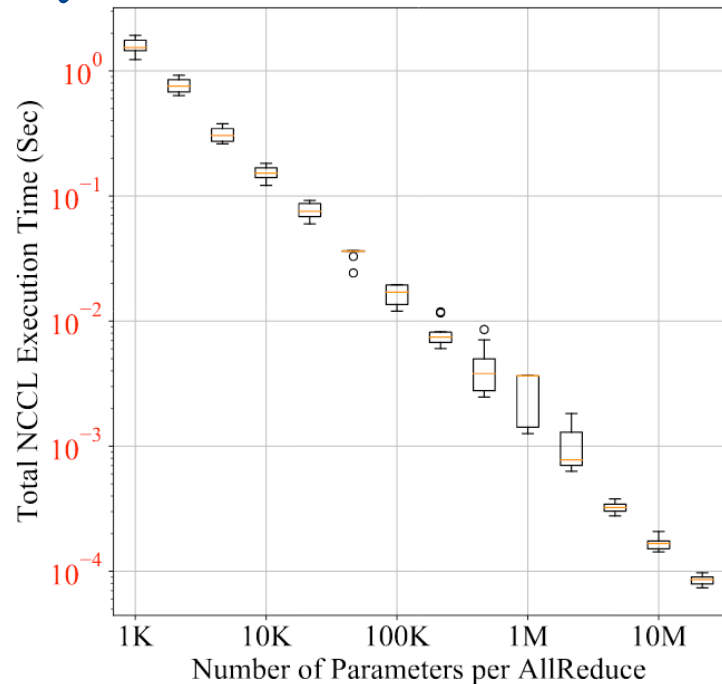
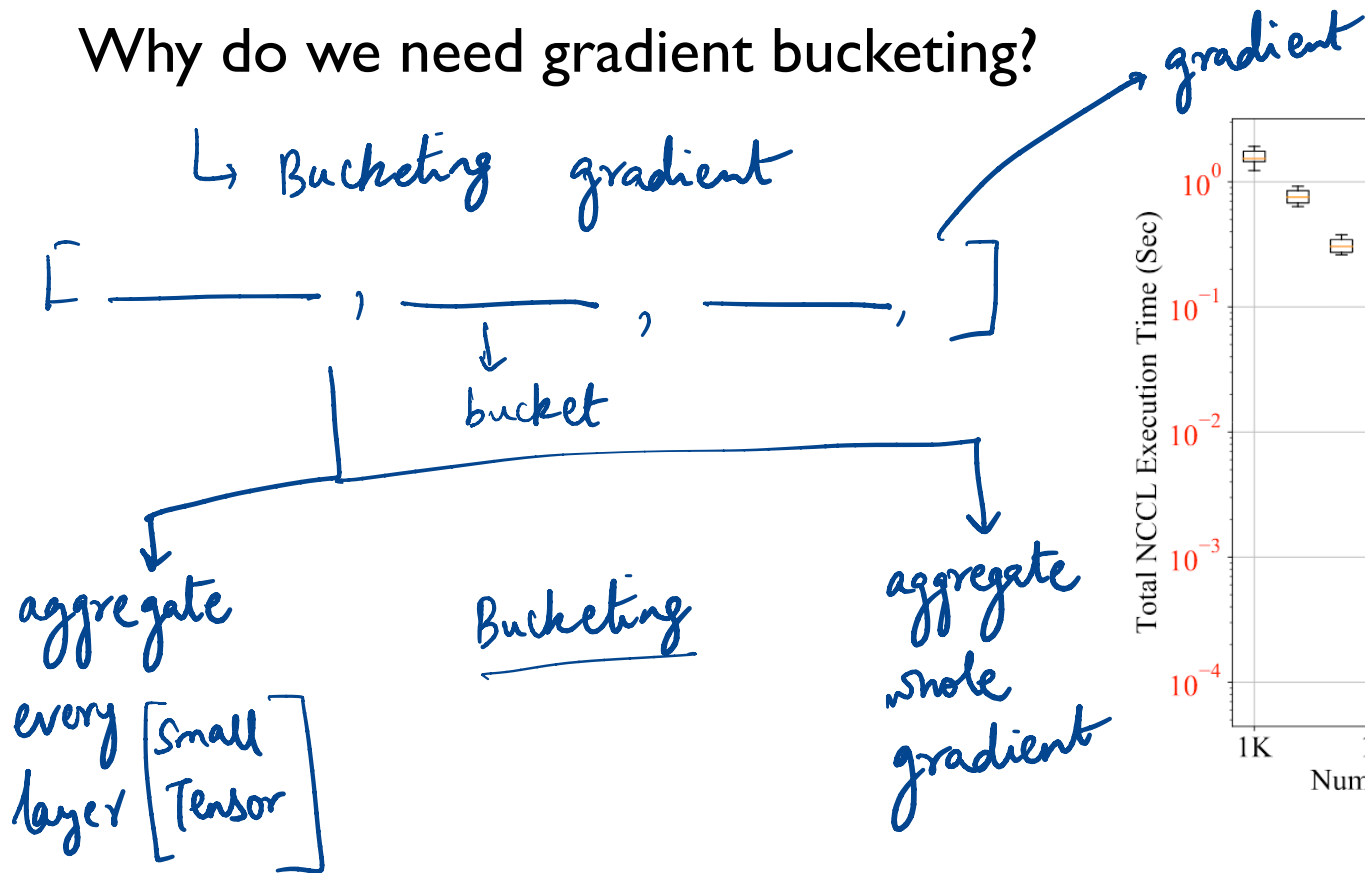
```
9 # setup model and optimizer
10 net = nn.Linear(10, 10)
11 net = par.DistributedDataParallel(net)
12 opt = optim.SGD(net.parameters(), lr=0.01)
13
14 # run forward pass
15 inp = torch.randn(20, 10)
16 exp = torch.randn(20, 10)
17 out = net(inp)
18
19 # run backward pass
20 nn.MSELoss()(out, exp).backward()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

gradient happens egg under the hood!

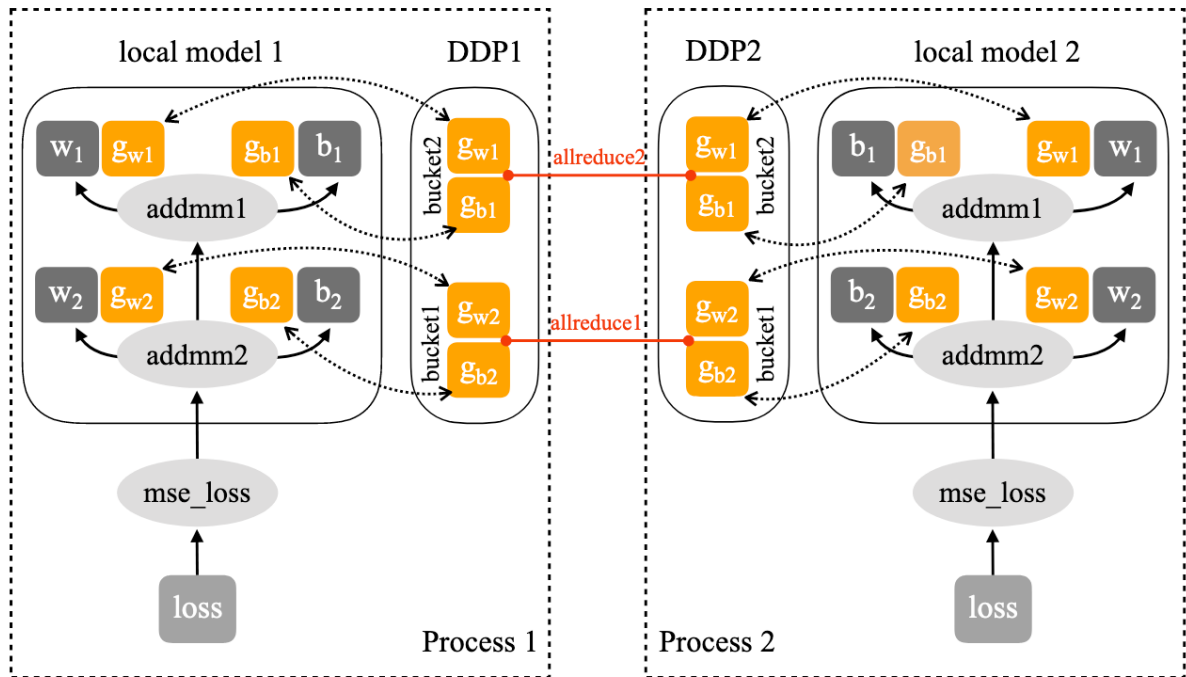
fwd, backward happens in \parallel

GRADIENT BUCKETING

Why do we need gradient bucketing?



GRADIENT BUCKETING + ALL REDUCE



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

forward
 backward → last layer
 gradient
 is available!

first

You can pipeline
 backward & gradient
 aggregation

$$T_{\text{com}} = \alpha T_{\text{lat}} + \beta T_{\text{BW}}$$

GRADIENT ACCUMULATION

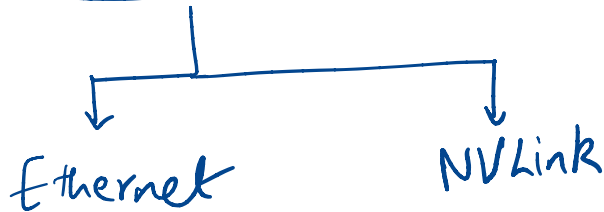
```
1 ddp = DistributedDataParallel(net)
2 with ddp.no_sync():
3     for inp, exp in zip(inputs, expected_outputs):
4         # no synchronization, accumulate grads
5         loss_fn(ddp(inp), exp).backward()
6 # synchronize grads
7 loss_fn(ddp(another_inp), another_exp).backward()
8 opt.step()
```

IMPLEMENTATION

Bucket_cap_mb → 25 MB is a good default

Parameter-to-bucket mapping → model walk backwards and assign params to buckets

Round-robin ProcessGroups → Multiple network types



SUMMARY

Pytorch: Framework for deep learning

DistributedDataParallel API

Gradient bucketing, AllReduce

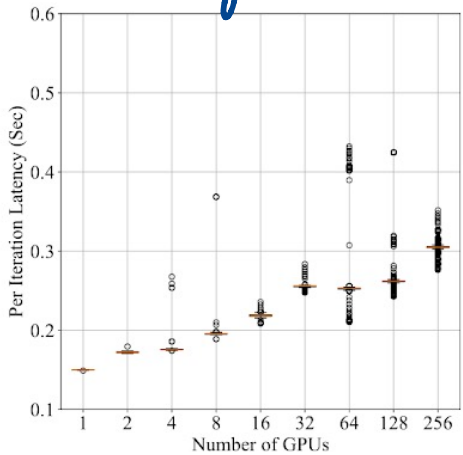
Overlap computation and communication

DISCUSSION

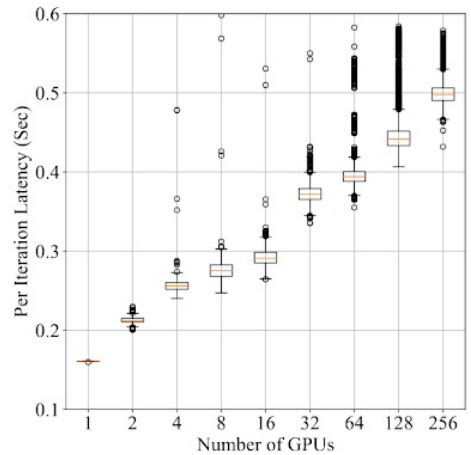
<https://forms.gle/jivzEEo5oz8tugYH9>

- NCCL is faster than Gloo \rightarrow NCCL is optimized for GPUs

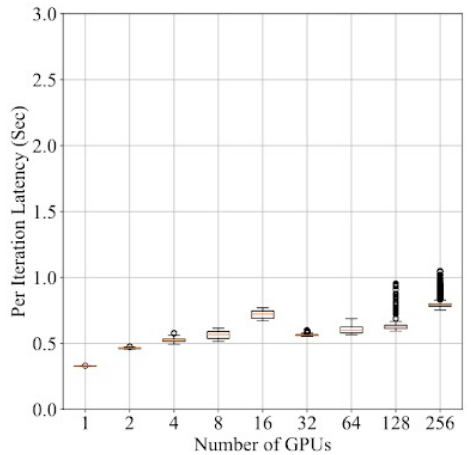
- Latency increases with num GPUs



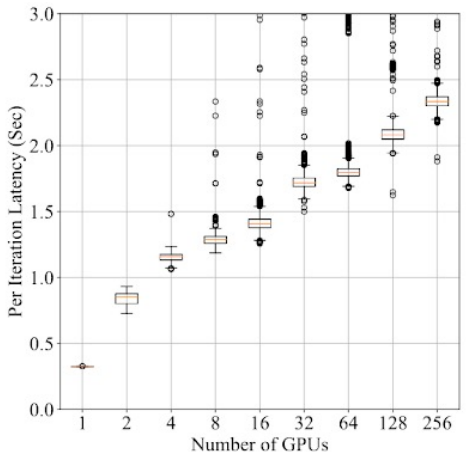
(a) ResNet50 on NCCL



(b) ResNet50 on Gloo



(c) BERT on NCCL



(d) BERT on Gloo

Figure 9: Scalability

NCCL is linear / sub-linear while Gloo has worse scaling

tail latency is higher for larger num GPUs

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

NEXT STEPS

Next class: PipeDream

Assignment 2 is out!

Project Proposal – Check Piazza!

BREAKDOWN

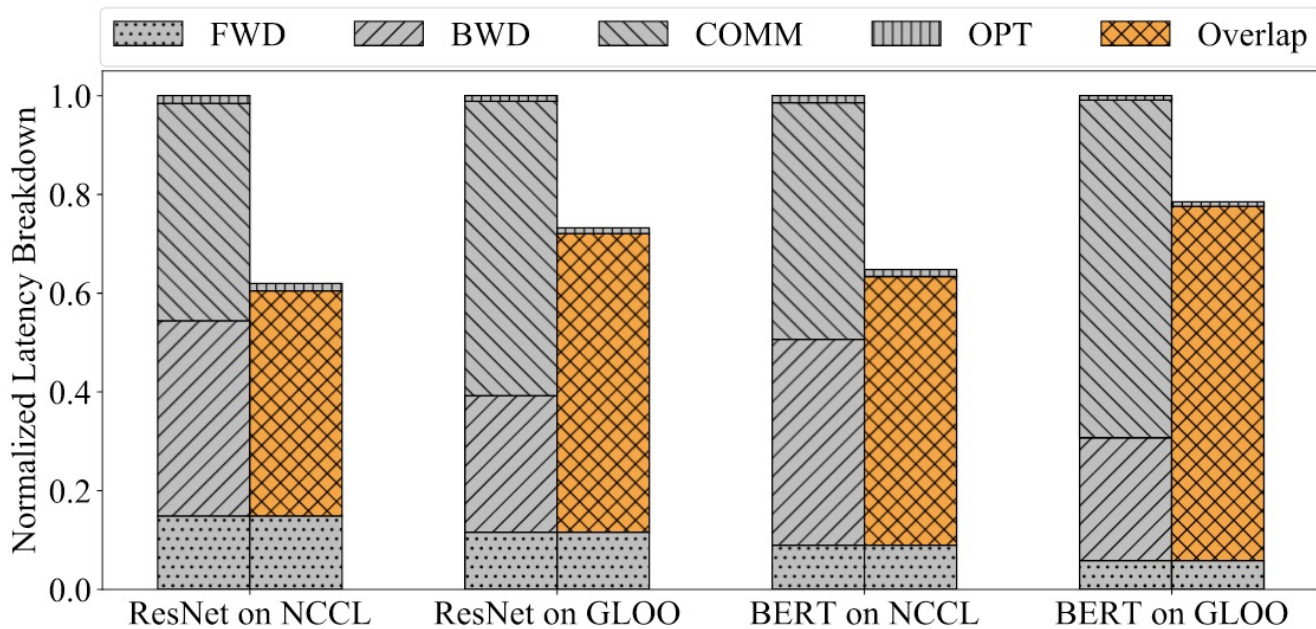


Figure 6: Per Iteration Latency Breakdown