

CS 744: PYTORCH

Shivaram Venkataraman Fall 2021

ADMINISTRIVIA

Assignment 2 out! Due Oct 13th early AM! 12th even g

Bid on topics, submit group (I sentences) – Oct II - One week! Title confirmed – Oct I4

Project Proposal (2 pages) – Oct 25

Introduction Related Work Timeline (with eval plan)

Spark WRITING AN INTRODUCTION

I-2 paras: what is the problem you are solving interactive data why is it important (need citations)

I-2 paras: How other people solve and why they fall short Is Related work and what is the downide

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

I para: Anticipated results or what experiments you will use

RELATED WORK, EVAL PLAN

Group related work into 2 or 3 buckets (I-2 para per bucket) Explain what the papers / projects do Why are they different / insufficient

Eval Plan

0-5 page

Describe what datasets, hardware you will use Available: Cloudlab, Google Cloud (~\$150), Jetson TX2 etc.





DEEP LEARNING



ResNet18

 → Convolution
 → ReLU
 → MaxPool
 → Fully Connected SoftMax



DATA PARALLEL MODEL TRAINING



B

COLLECTIVE COMMUNICATION

Broadcast, Scatter

Gather, Reduce

Ly late

1980,



ALL REDUCE USING A RING



DISTRIBUTED DATA PARALLEL API

```
-> one line of
Code change
    # setup model and optimizer
9
    net = nn.Linear(10, 10)
10
    net = par.DistributedDataParallel(net)
11
    opt = optim.SGD(net.parameters(), lr=0.01)
12
13
                                           - rest of code
se resembles
sigle machine
Gode
    # run forward pass
14
    inp = torch.randn(20, 10)
15
    exp = torch.randn(20, 10)
16
    out = net(inp)
17
18
    # run backward pass
19
    nn.MSELoss()(out, exp).backward()
20
21
    # update parameters
22
    opt.step()
23
```

GRADIENT BUCKETING



Why do we need gradient bucketing?





GRADIENT ACCUMULATION



IMPI FMFNTATION

= Profiling network topology + libraries Bucket_cap_mb

Parameter-to-bucket mapping - walk back from the end and add layers to buckets



D-V NV hink

SUMMARY

- Pytorch: Framework for deep learning
- DistributedDataParallel API
- Gradient bucketing, AllReduce
- Overlap computation and communication

DISCUSSION

https://forms.gle/YnZC8PKQyICDFJRf9





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

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

- spark job could have many stages - Reduction tree + Broadcast La plower than this paper? - How mutable models can be handled Assuming that all tasks are running at some time Sparkt Many morie tasks Sparkt than wrea

NEXT STEPS

Next class: PipeDream

Assignment 2 is out!

Project Proposal

Preferences, Groups by Oct 11 2 pager by Oct 25

BREAKDOWN



Figure 6: Per Iteration Latency Breakdown