Hil.

CS 744: NEXUS

Shivaram Venkataraman Fall 2021

ADMINISTRIVIA

Course Project Proposals

- Due Oct 25! --- Monday
- See Piazza for template

Midterm details

- Oct 28th: Includes papers from Datacenter as a Computer to Nexus
- Open book, open notes
- Held in class time 9.30-10.45am Central Time



EXAMPLE APPLICATION



Scheduler to
determine placement GOAL: HIGH GPU UTILIZATION
and hot ding
Placement

$$rightarrow on accederator like
 $rightarrow on accederator like$$

SCHEDULING BATCHED EXECUTION 15mz



A for 75 ms, B for 50 ms, A for 75 ms... (A & C cannot be aboated haterey (A) + haterey (B) will both meet SLOS (because that will lead to SLO violation Consider A, B being colocated

BATCH-AWARF SCHEDULING tput read every request latency 1 SLO/if not drop [latency at various batch sizes Inputs: Request rate, SLO for each model, Profiles at batch size Approach: Allocate "full" GPUs based on load. Handle residuals 63=6 If req. rate = 1100 reg/s. 1 GPU at bs = 16, $125 \text{ reg/s} = 3 \begin{bmatrix} 1100\\125 \end{bmatrix} = 8$ = 3 100 reg/sduty cycle d₁ -> Node * Greedy Approximation "revolust" -> Node 2 Decreaning & start by giving every residual workload its own GPU occupancy & Merge workloads (Figether => 2 GPUs > 2 GPUs > 2 GPUs > 2 GPUs batch b₂ smaller Node 1 Node 2 -> 2 GPU Merged Node

HANDLING COMPLEX QUERIES

Challenge:

How do we set latency SLOs for complex queries?

Additional challenge we don't know what fraction cells "Car" and what fraction calls "face"



SCHEDULING COMPLEX QUERIES



ADAPTIVE BATCHING

Clipper: Adapt the batch size based on the oldest request





BATCH-AWARE DISPATCH



OTHER FEATURES



GPU Multiplexing La round rubin fashion

Overlapping CPU and GPU computation Re - processing & post - processing

NEXUS ARCHITECTURE



SUMMARY

- ML Inference goals: latency SLO, GPU utilization
- Nexus: Handle multiple tenants, multiple DNNs
- Schedule using squishy bin packing
- Breakdown SLO for complex queries, adaptive batching



Pytorch Distributed

DataParallel Training API Overlap compute, communication

PipeDream

Generalize parallelism: Pipeline parallel Reduce network, maintain consistency

Ray

Reinforcement learning applications Actors and tasks, Local and global scheduler

Pollux

Scheduler ML training jobs in a cluster Co-adaptive scheduling to set batch size, LR

Nexus

System for ML Inference, scheduling

Meet latency SLOs while ensuring high utilization

DISCUSSION https://forms.gle/XQ4CfNzTTFsSrVv7A

Consider a scenario where you have a model that takes variable amount of time depending on the input. For example if a frame contains 100 cars it takes 250ms to process but if the frame has 1 car then it finishes in 10ms. What could be one shortcoming in using Nexus to schedule this model? SL0 = 250 ms



Next class: SQL