hood morning !.

## CS 744: CLIPPER

Shivaram Venkataraman Fall 2020

## ADMINISTRIVIA

Course Project Proposals

- Due on Friday! →
- See Piazza for template
- Submission instructions soon -

Midterm details

-> Oct 22 \_\_\_\_\_ Veto ML / banple motherms section / on liazza

- Open book, open notes
- Held in class time 9.30-10.45am Central Time -
- Type / Upload photos (extra 15 mins)













BATCHING, QUEUING

Model Container

Spar

Goals, Insight

- Increase latency (within SLO) batches use for improved throughput parellism Reduce RPC overhead
- GPU / BLAS acceleration

Approach

- Per container queues.

- Why?





**Fensor**Flo

GPU

Model Container

CPU

learn

## **ADAPTIVE BATCHING**





# Nordels Nordels Lip which one should I use? SINGLE MODEL SELECTION

The Usual Place

associate of the option

Multi-Arm Bandit formulation

- Explore vs Exploit
- Regret: Loss by not
- picking optimal action
- Goal: Minimize regret

#### Clipper

- Exp3 algorithm
- Single evaluation
- Scales to more models

- update weighte based on feedback

Midel Selection

GRAND OPENING!

## MULTIMODELS -> ensembles

Predict movie

blerd?

Combine & propheld 70:5 -> cat

Ensemble

- Combine output from models (weighted average)
- How do we get the weights ? linear combination TF = X Y + B Zi

**Robust Prediction** 

- React to model changes
- Binary C1 0:25 0:75 Exp4 -> update X & B C2 0:4 0:6 - Output confidence score

## STRAGGLER MITIGATION

Lo If we wait for N model containers to reply, some of them night be Maw? Lo more replices / scaling? Why do stragglers occur? gg pe lokery -) .... (2 replies Approach Ly Approx result based on whatever las finished -> Better approx non late! -> ML specific

## SUMMARY

- Clipper: ML inference Workloads + Requirements
- Layered architecture provides generality
- Caching, Batching, Replication to improve latency, throughput
- Multi-Arm bandits to improve accuracy

# DISCUSSION

https://forms.gle/FCVhPURqz7HSbDtg6

Consider a scenario where you run a model serving service that hosts a number of different applications. The traffic for some applications is sporadic (e.g. only a few hours where they are used). What are some advantages / disadvantages of using Clipper for such a service?

Diss advantages -> Rache night be contented Advantages -> Adaptive batching delaged -> tune > how to do this online a fashion? -> multiple replicas elasticity 10-20ms / request dipese RPC ppc user provided Lode e Oms

8-10 reems to a good rize - Labouring bunch of different things? entendet la terres 1 Me prore resources mets goig 300 % Ensemble Missing 100 1.00 Straggler Mitigation P99 Latency (ms) 200 150 100 50 -+ Stragglers P99 - P99 250 80 0.95 -+- Stragglers Mean -I- Mean Straggler Mitigation Mean Accuracy 60 0.90 0.85 40 20 0.80 50 0.75 0 2 6 12 14 16 10 12 14 16 6 2 Ω 6 8 12 4 Size of ensemble Size of ensemble Size of ensemble (a) Latency (b) Missing Predictions (c) Accuracy ) 16 entembles lateny inflation is very low reasonable accuracy