Hil Good morning!

#### CS 744: RAY

Shivaram Venkataraman Fall 2021

# ADMINISTRIVIA

- Assignment Two: Due tonight!
- Project proposal aka Introduction (10/25)
  Introduction

Related Work

Timeline (with eval plan)

Confirm project title in couple of days



### **REINFORCEMENT LEARNING: APPLICATIONS**





### **RL REQUIREMENTS**

low latency for serving your model / Simulation poli cy resources con also vary exe cution GRUS, CPUS etc. Training time can be variable Ly millise conds to hours Compute hours prev. iteration model Stateful computation Serving stateless -, serving simulation



ready = ray.wait(futures, k,timeout)







#### FAULT TOLERANCE

lineage, replay or re-execute tasks Tasks side - effect free the ck point the actor state periodically Actors chair replication - handle replicer failures GCS -> statelers => verpown a new scheduler Global Scheduler

# SUMMARY

Ray: Unified system for ML training, serving, simulation

Flexible API with support for

Stateless tasks

**Stateful Actors** 

Distributed scheduling, Global control store

# DISCUSSION

https://forms.gle/MnsCJA87CVhMmShs8

Consider you are implementing two apps: a deep learning model training and a sorting application. When will use tasks vs actors and why ?

Deep learning model training - Use actors, have model as state. Iteration is 4 method call - Mix actors & tasks All duce ( Data parallel - each boal gradient reduce ( Actor - optimize, stepl)

Sorting - Merge sort -> can we have data sorted so far as state? -> Partition to many actors? - Data, partition it, object store state less tasks intermediate data on object store!



#### NEXT STEPS

Next class: Pollux