# CS 744: RAY

Shivaram Venkataraman

Fall 2019

# ADMINISTRIVIA

- Assignment 1 Grades

- Assignment 2 due on Fri

- Course Project emails

# Bismarck

Supervised learning, Unified Interface
Shared memory, Model fits in memory

# Parameter Server

Large datasets, large models (PB scale)
Consistency model, Fault tolerance

Machine Learning

# Tensorflow

Need for flexible programming model
Dataflow graph
Heterogeneous accelerators

CNNs
language model
PS tasks

## WORKLOADS

### Bismarck

→ Convex optimization

→ Small datasets

simple models — SVMs, Logistic reg.

### Parameter Server

→ Larger datasets → Sparse data

→ High dim parameters → Ad Click model

### Tensorflow

↳ Dense features
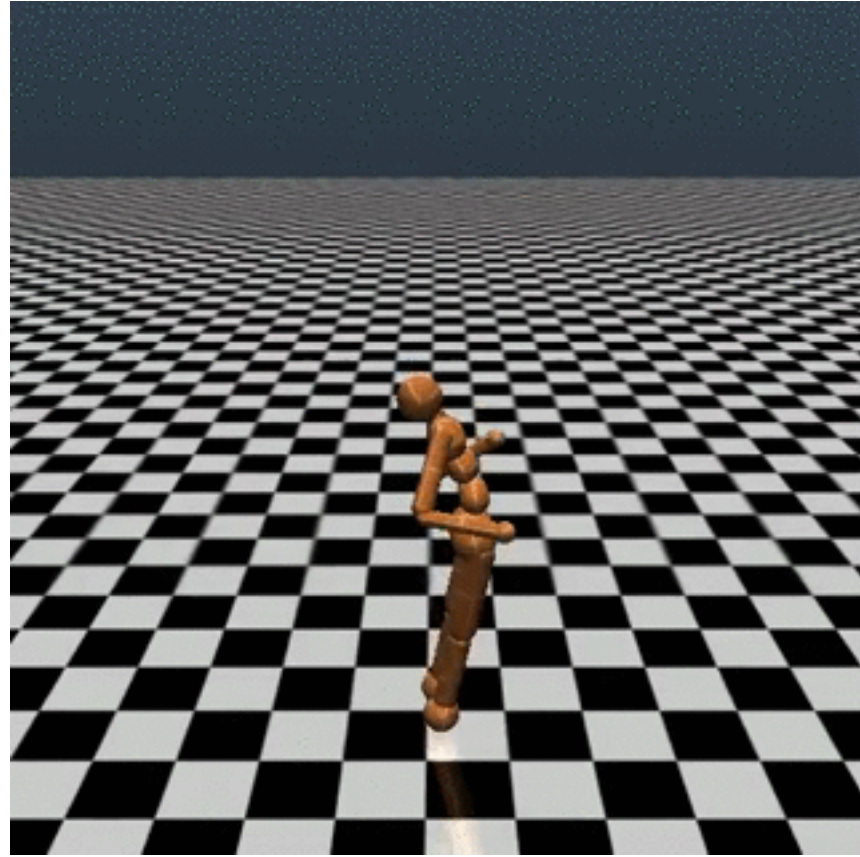
↳ Advanced models    Supervised learning

Deep learning

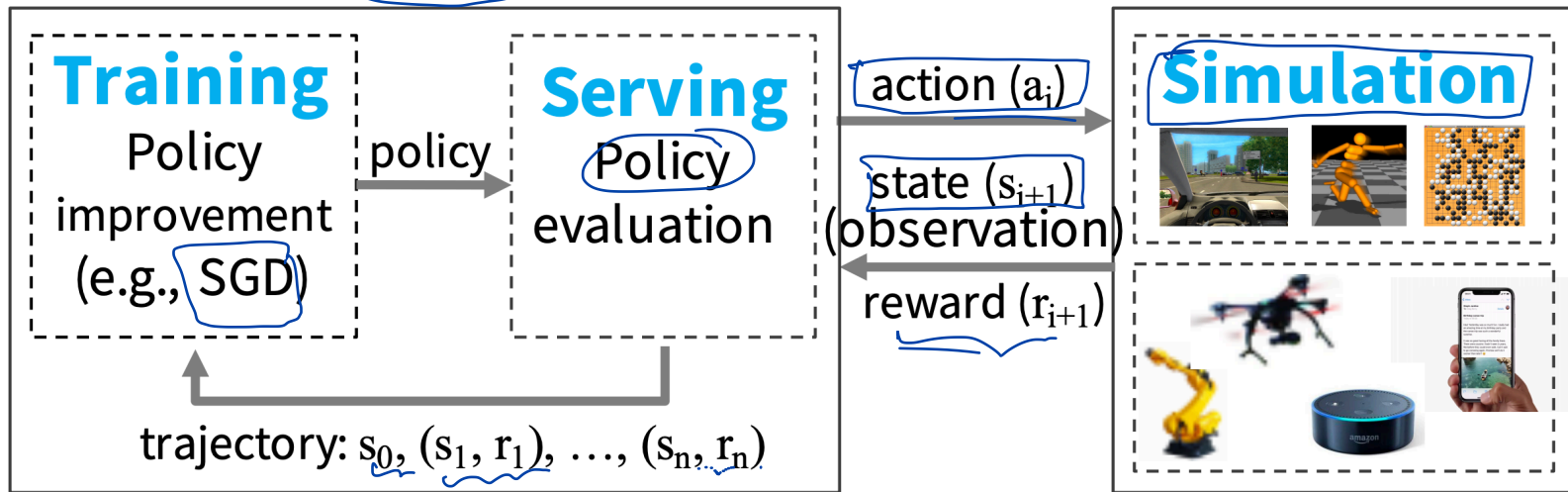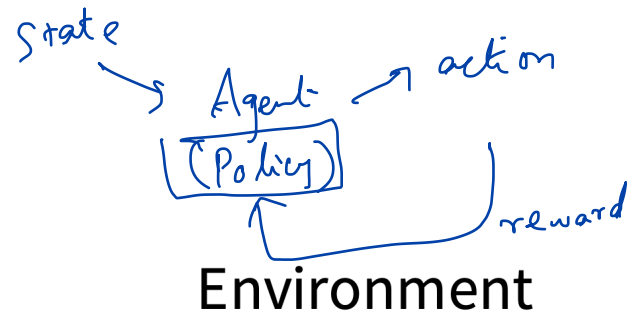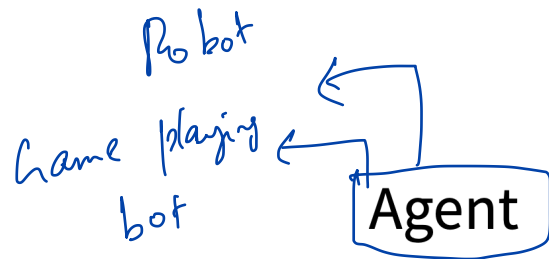↳ Image Classification   Inception v3

Supervised Learning

Training Data                    Labels

        Model

        $\longrightarrow$

        fit

# REINFORCEMENT LEARNING

# RL SETUP

state → Agent → action
(Policy)
reward

Robot

Game playing bot

Agent

Environment

**Training**
Policy improvement
(e.g., SGD)

policy

**Serving**
Policy evaluation

action ($a_i$)

state ($s_{i+1}$)
(observation)

reward ($r_{i+1}$)

**Simulation**



trajectory: $s_0$, $(s_1, r_1)$, ..., $(s_n, r_n)$

improve Policy    given sequence    of state,    reward

# RL REQUIREMENTS

Stateful $\equiv$ Game engine
+
Stateless

**Simulation**
     rollout
     $\hookrightarrow$ Tasks could be of varying length   <u>1ms</u> $\rightarrow$ seconds/mins

**Training**
     Tasks are deterministic   (state, action) $\rightarrow$ reward

     Dynamic execution $\equiv$ iteration structure
                               depends on current iteration

**Serving**
     $\hookrightarrow$ Very low latency $\equiv$ Parallelism

          Rollouts happen in parallel

# RAY API

f1 : f.remote (args 1)

## Tasks
normal python function

futures = f.remote(args)

Task which with run
f with args

→ Handle to result of task

## Actors

Erlang

Queue

Class : Actor, Shopping Cart

add

remote

variable

Output

add(item2)   add(item1)

rCmove(item1)

actor = Class.remote(args)
futures = actor.method.remote(args)

= actor. method. remote (args 1)

f (args) before    f (args1)

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

Stateless tasks

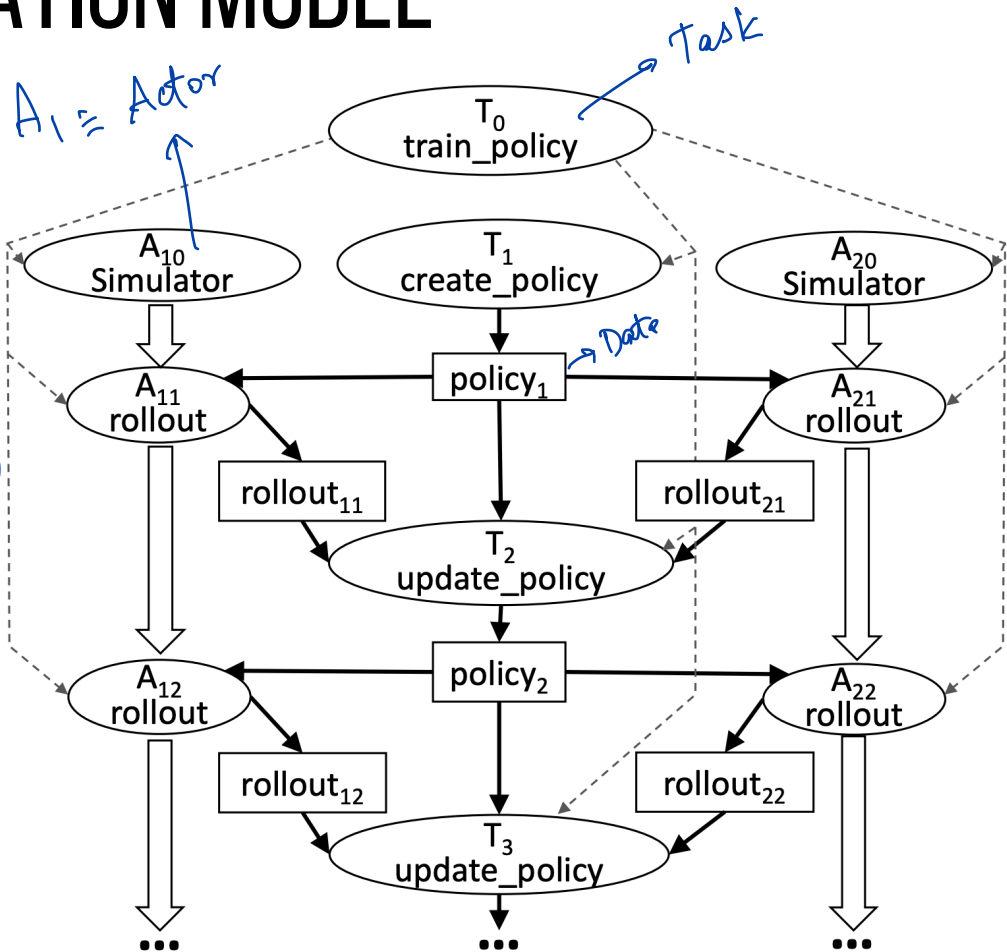# COMPUTATION MODEL



Vertices:

Data.

Tasks / Actors

$A_1 =$ Actor

def task()

= task1. remote (a)

$\rightarrow$ Task

$\rightarrow$ Data

Edges

$--\rightarrow$ Control edge

$\longrightarrow$ Data edge    Data vertex    Args
Task / Actor

$\Longrightarrow$ Stateful edge
Actor methods    to ensure
ordering

$T_0$
train_policy

$A_{10}$
Simulator

$T_1$
create_policy

$A_{20}$
Simulator

$A_{11}$
rollout

policy$_1$

$A_{21}$
rollout

rollout$_{11}$

rollout$_{21}$

$T_2$
update_policy

$A_{12}$
rollout

policy$_2$

$A_{22}$
rollout

rollout$_{12}$

rollout$_{22}$

$T_3$
update_policy

...    ...    ...

# ARCHITECTURE

# GLOBAL CONTROL STORE

Object table
↳ list of object
their locations | Namenode
metadata

Task table
↳ Lineage, tasks created
edge in Comp. graph

Function table
↳ Code blocks that
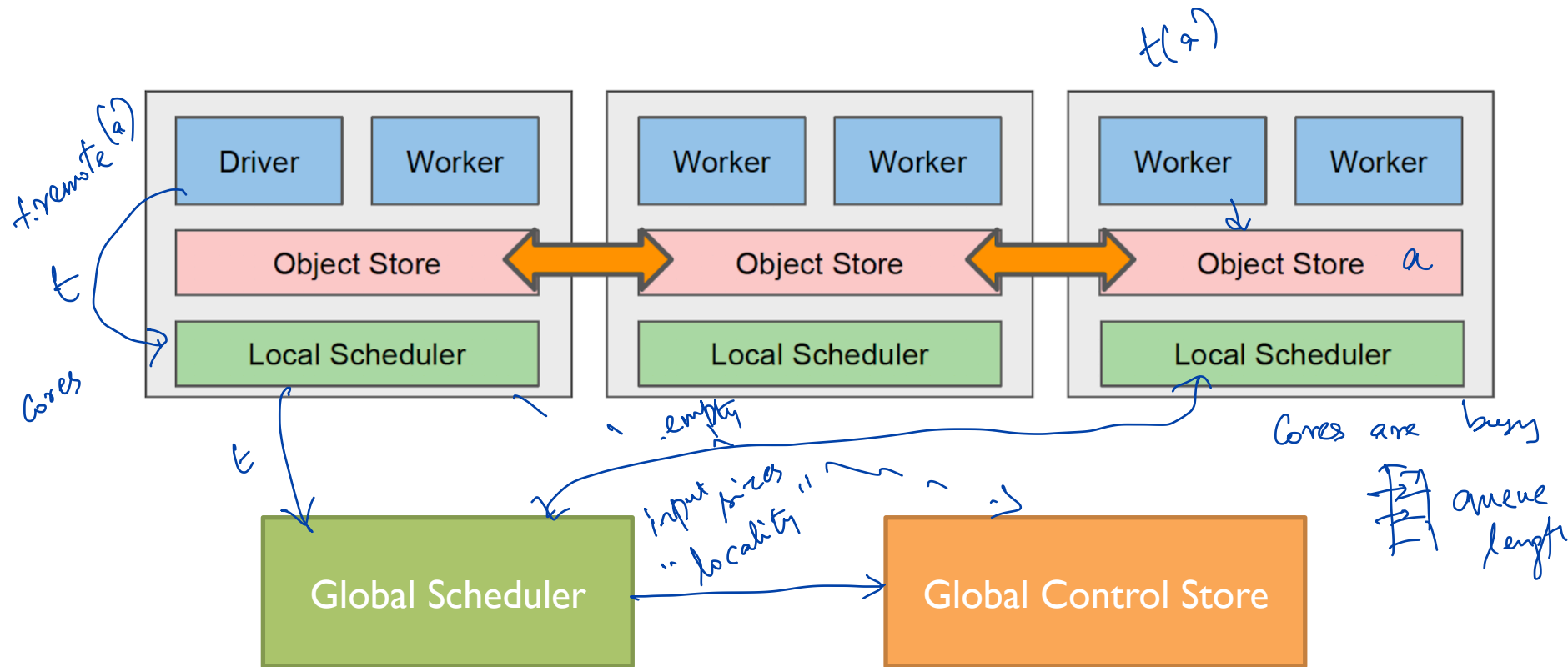are running

Externalized State
↳ Bottlenecks not in
global system
↳ Sharding

Fault tolerance
↳ chain replication

# RAY SCHEDULER

# FAULT TOLERANCE

Tasks
    ↳ lineage from GCS

Actors
    ↳ "checkpointing"

GCS
    ↳ sharded, replicated

Scheduler ⟶ Stateless,

Driver
    ↳ Computation fails

GCS → extra resources
    ↳ scalability

Object = Shard

multiple schedulers
      or
backup sched.

# DISCUSSION

https://forms.gle/QQyLbwjAufJNXWnr6

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

**Actors**

Sorting :

wait for other tasks ?

**Tasks**

You can parallelize and stateless operations

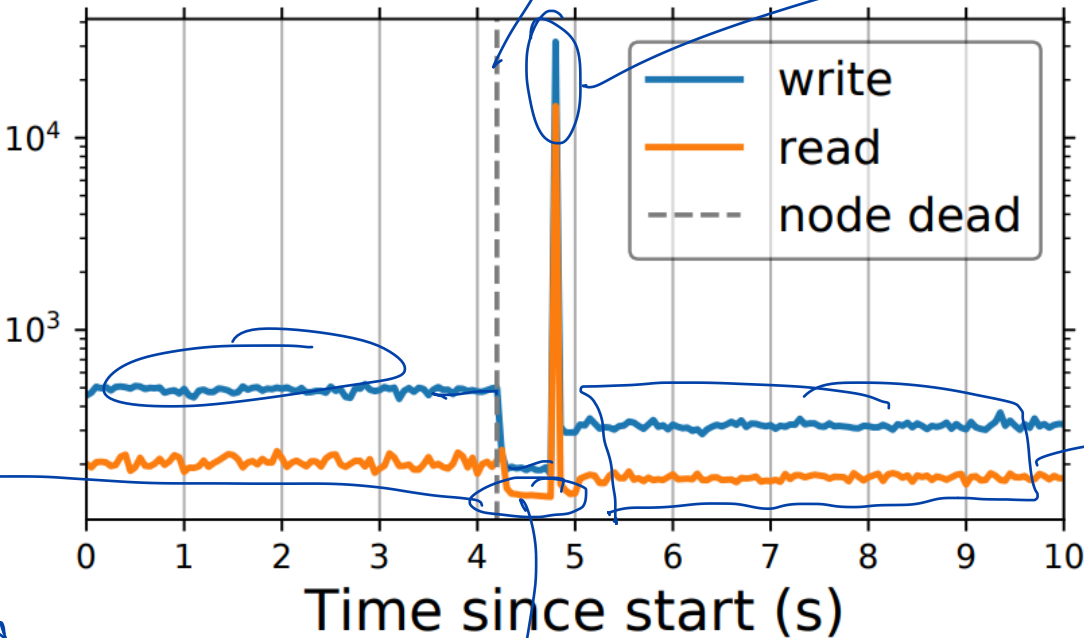Model training

Parameters are "state" don't need to broadcast

Flexible sync. fine grained recovery ?

Why is read latency going down?

① 

⓪

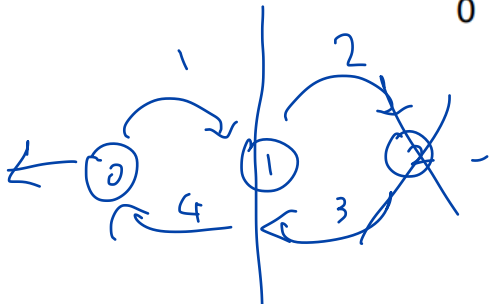one of replicas GCS is down

warming up of new replica

Latency (μs)

10⁴

10³

write
read
node dead

10⁴

10³

lower than before

0 1 2 3 4 5 6 7 8 9 10

Time since start (s)

latency goes down!

→ New replica somehow is closer? or has more resources?

Considering AllReduce using MPI as the baseline parallel programming task. Discuss the improvements made by MapReduce, Spark over MPI and discuss if/how Ray further contributes to the comparison.

# NEXT STEPS

Next class: Clipper

Assignment 2 due this week!

Course project