

Hello!

# CS 744: PARAMETER SERVERS

Shivaram Venkataraman

Fall 2022

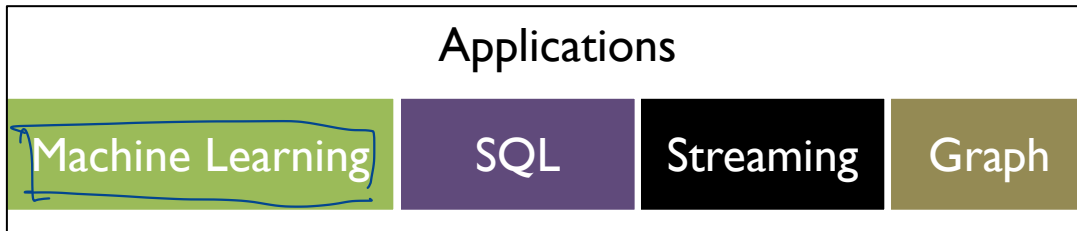
# ADMINISTRIVIA

- Assignment 2 is due on Oct 12 (Wed) at 10am! → piazza
- Course project groups are also due same time?!

↳ Your own project

↳ Seed ideas, paper pointers

PyTorch  
Distributed



Pipe Dream



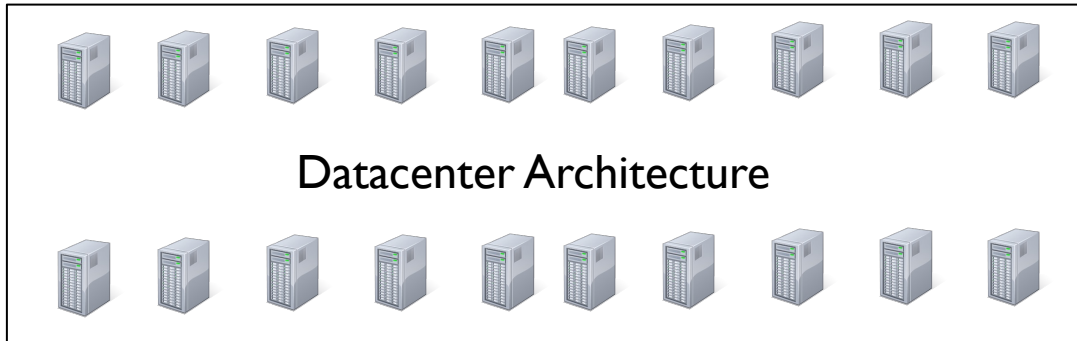
diff schemes



for  
Distributed



Training



## PyTorch

Distributed DataParallel

Easy-to-use Interface

Model replicated on every worker

1M parameters

≈ 8MB storage

## PipeDream

Model and pipeline-parallelism

Split model across workers

model (if dense)  
[0.25, 0.5, ----- ]

gradient  
[0.1, 0.02, ----- ]

Commonalities?

Training Deep Learning models

→ Dense models and

dense updates to models

# PARAMETER SERVER: MOTIVATION

- Large training data 1TB to 1PB
- Models with  $10^9$  to  $10^{12}$  parameters
- Goals
  - Efficient communication
  - Flexible synchronization
  - Elastic Scalability
  - Fault Tolerance and Durability

*sparse updates*

→ 1 training example  
compute gradient  
↳ mostly zeros  
few non-zeros

→ Both computation and  
comm. only for non-zero  
entries

# EXAMPLE WORKLOAD

## Ad Click Prediction

- Trillions of clicks per day
- Very sparse feature vectors

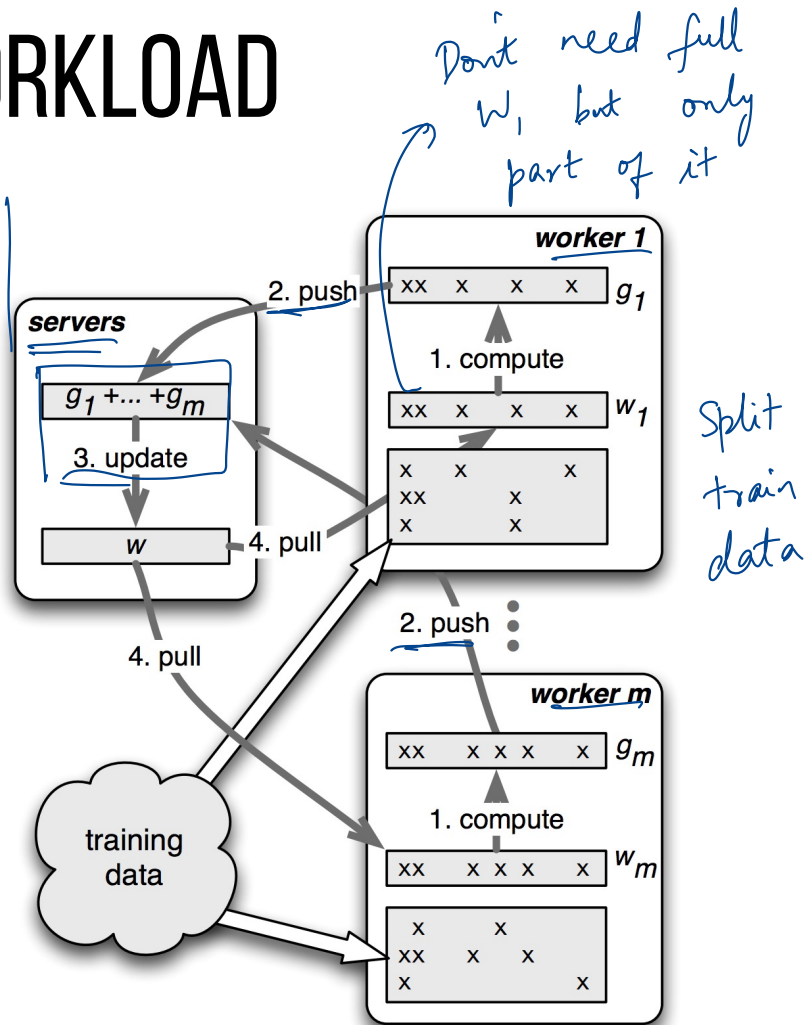
$[0 \ 0 \ 0 \ \dots \ 0.25 \ \dots \ 0]$

## Computation flow

event, "example"

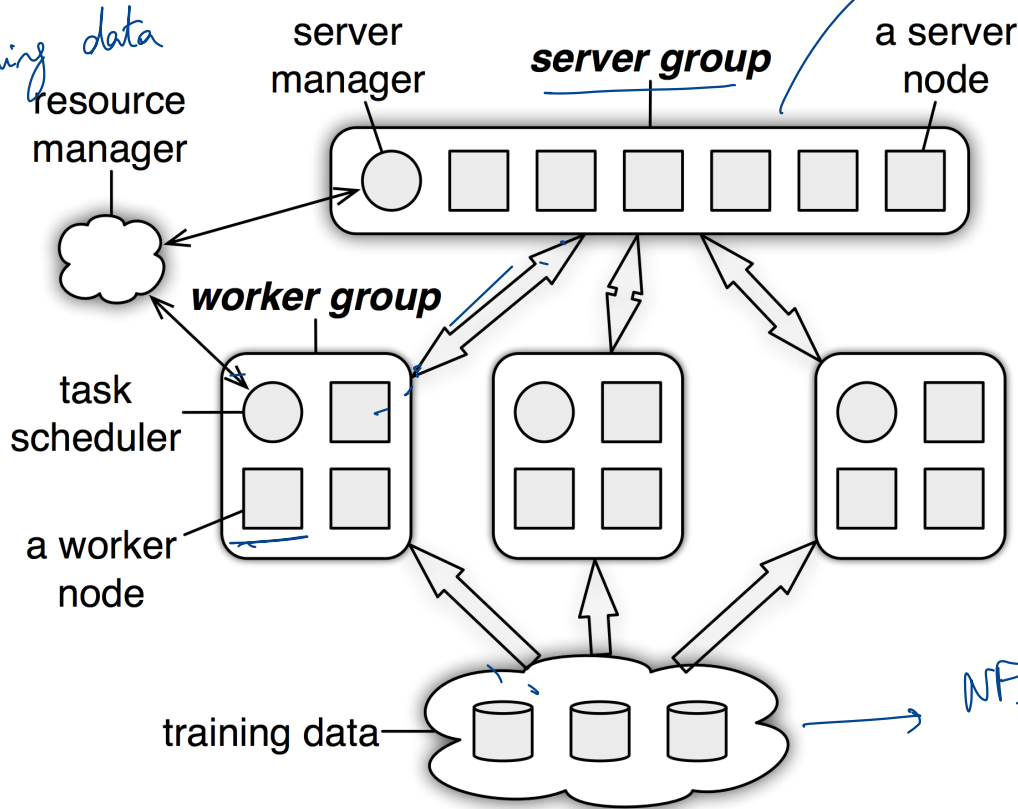
$d = \text{millions or more}$

1. Compute gradient wrt current model  
Density of  $w_i$  depends on data properties
2. Push gradients to servers
3. [Server] updated value of  $w$



# ARCHITECTURE

Based on size of training data  
Num of workers  
= Num of servers  
Based on number of parameters



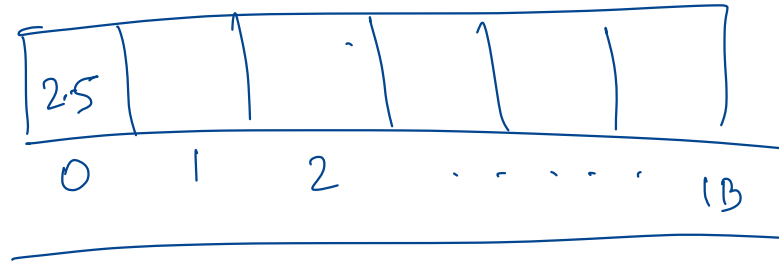
Model is sharded or partitioned

different computation  
could use diff worker group

# REPRESENTATION

push(0, 2.5)

- Key value pairs e.g., (featureID, weight)
- Assume keys are **ordered**.  
Easier to apply linear algebra operations
- Interface supports **range** push and pull  
w.push(R, dest) [1:4]
- Support for **user-defined functions** on server-side





# TASK DEPENDENCY

Model might not converge

iter 10: gradient — push & pull → ○

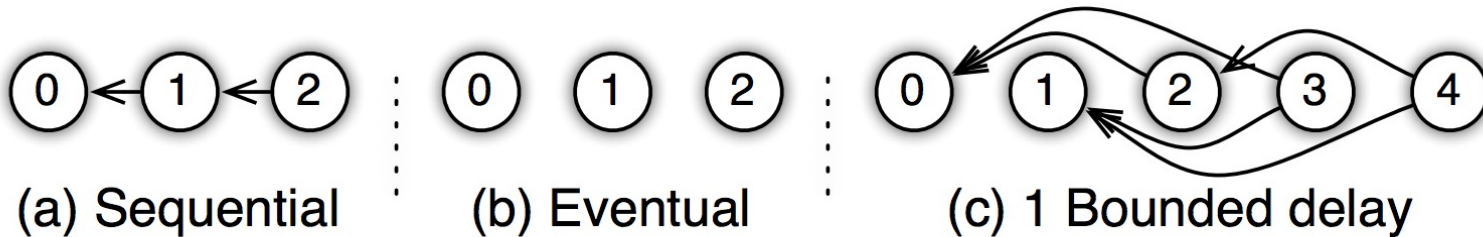
iter 11: gradient — push & pull → ○  sync across iters

If you want to mimic single thread  
- then iter  $i+1$  can start only after  
iter  $i$  completes

---

Can overlap iterations to accelerate training  
at the cost of stale model contents

# CONSISTENCY MODELS



User defined filters

Significantly modified filter

↳  $[0, 0, \dots, 0.1, 1e-5, 1e-2, \dots]$

KKT filter

Ignore any values  $< 1e-4$

↳ worker side

gradient

T- bounded delay

↳ upper bound on staleness

→ Creates dependencies

T iterations before

# IMPLEMENTATION: VECTOR CLOCKS

→ logical timestamps

Distributed Systems

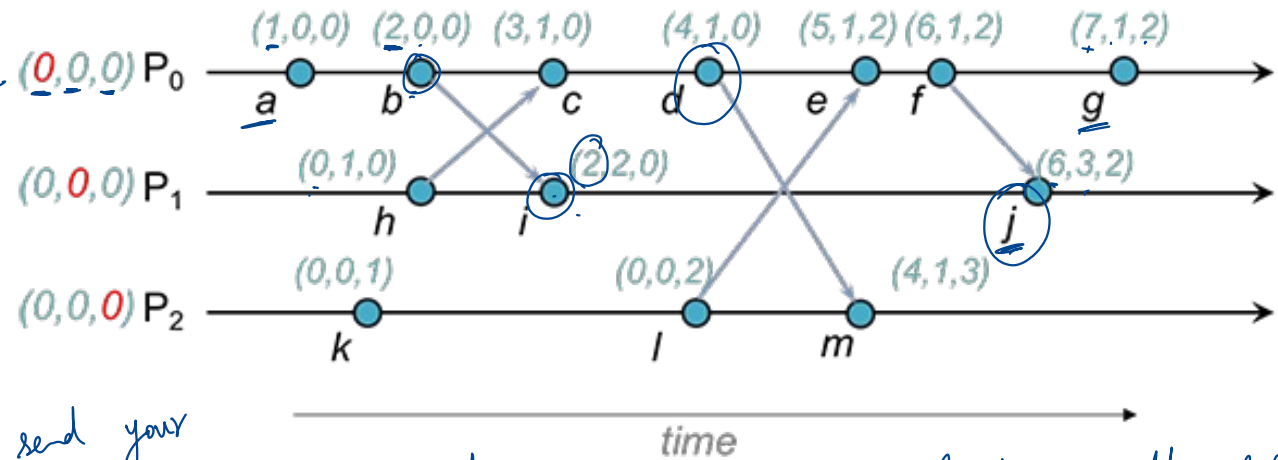
event j happen before event d

3 Processes  
3 tuple

compute - inc 1

send - inc 1, send your clock

recv - inc local, copy other process clock in the message



d before j  
g and j concurrent

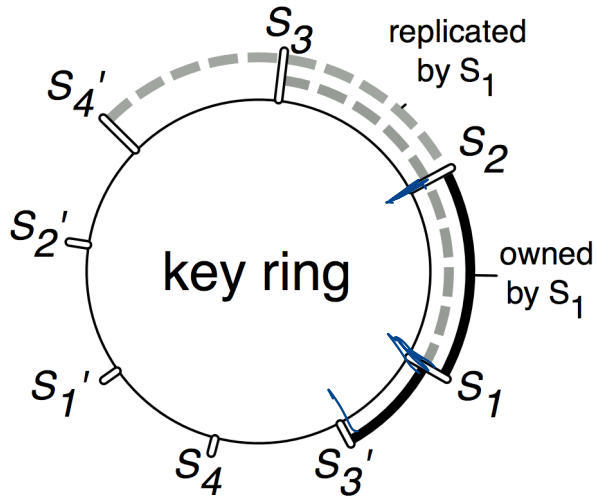
happens before: Compare all coordinates for i, j if  $TS_i \leq TS_j$  and less than in at least one of them

Each  $k, v$  pair has a vector clock associated with it

- Every push you increment vector clock.
- Every pull returns clock value.

Optimization

# IMPLEMENTATION: REPLICATION

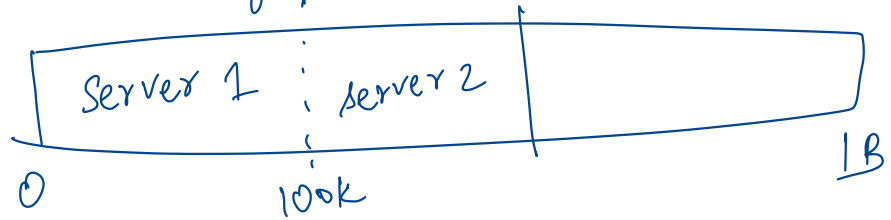


Server 2 owns  
100k - 200k  
replica 0 - 100k

Distributed Hash Tables [early 2000s]

Replication after aggregation

- Partition the key space using a ring



- Replicas next two nodes on the ring

# FAULT TOLERANCE

→ Ideas from

server addition (similar failure)

DHT

1. Server manager assigns the new node a key range to serve as master.



2. The node fetches the range of data to maintain as master and  $k$  additional ranges to keep as slave.

3. [The server manager broadcasts the node changes.]  
The recipients of the message may shrink their own data

→ From this point  
the new  
membership is in  
effect.

---

**Algorithm 3** Delayed Block Proximal Gradient [31]

---

**Scheduler:**

- 1: Partition features into  $b$  ranges  $\mathcal{R}_1, \dots, \mathcal{R}_b$
- 2: **for**  $t = 0$  **to**  $T$  **do**
- 3:     Pick random range  $\mathcal{R}_{i_t}$  and issue task to workers
- 4: **end for**

**Worker  $r$  at iteration  $t$** 

- 1: Wait until all iterations before  $t - \tau$  are finished  $\rightarrow$  bounded staleness
- 2: Compute first-order gradient  $g_r^{(t)}$  and diagonal second-order gradient  $u_r^{(t)}$  on range  $\mathcal{R}_{i_t}$
- 3: Push  $g_r^{(t)}$  and  $u_r^{(t)}$  to servers with the KKT filter
- 4: Pull  $w_r^{(t+1)}$  from servers

**Servers at iteration  $t$** 

- 1: Aggregate gradients to obtain  $g^{(t)}$  and  $u^{(t)}$
- 2: Solve the proximal operator

$$w^{(t+1)} \leftarrow \underset{u}{\operatorname{argmin}} \Omega(u) + \frac{1}{2\eta} \|w^{(t)} - \eta g^{(t)} + u\|_H^2,$$

where  $H = \operatorname{diag}(h^{(t)})$  and  $\|x\|_H^2 = x^T H x$

---

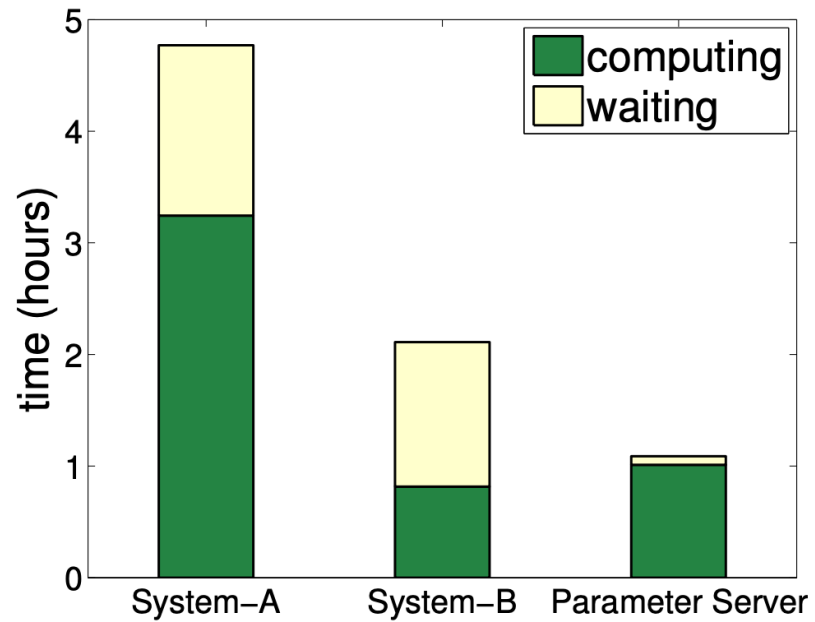
# SPARSE LR

User defined function

# DISCUSSION

<https://forms.gle/qPX1bBCAsd2fhL2i6>





## What are some of the downsides of using PS compared to implementing Gradient Descent in PyTorch / Spark?

- PS assume sparse updates. What if updates are dense?
  - Comm but also memory overheads
  - Repd and consistency format
  - Key, vector clocks etc. storage overhead
  - Gather + Scatter / Broadcast
- Async exec could affect convergence with dense updates

How would you integrate PS with a resource manager like Mesos? What would be some of the challenges?

- "Tasks" in Mesos is a worker iteration

"push" ends the task

→ wait for task allocations esp. if you want all workers to finish

→ lose cache between iterations

\*

# NEXT STEPS

Next class: Gavel

Assignment 2 is due soon!