

Hello!

CS 744: DRF

Shivaram Venkataraman

Spring 2024

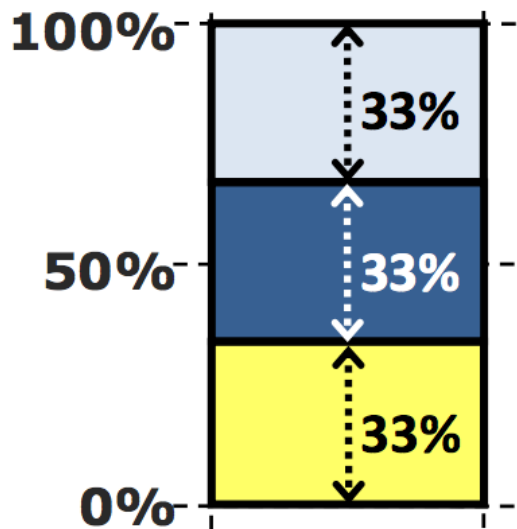
ADMINISTRIVIA

- Assignment 2 done! → Assignment 1 grading
↳ wed / Thur
 - Course Project
 - Form groups and submit project bids by tonight!
 - Assigned projects by **March 1** → preferences
 - Introductions due **March 8**
- Come up at the end of the class

SETTING: FAIR-SHARING

3 clients or users
1 resource = CPU

Equal Share

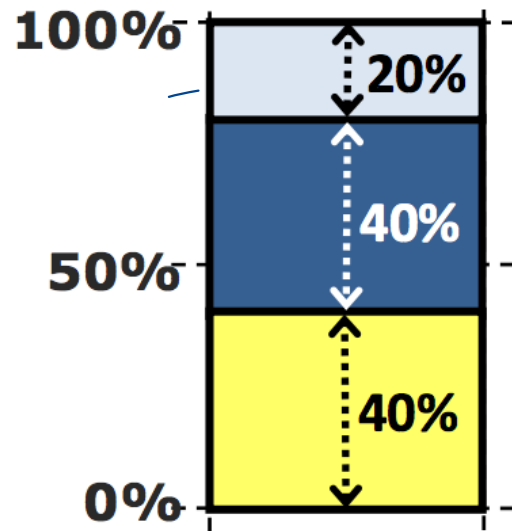


Max-Min Share

Maximize the allocation for most poorly treated users

Maximize the minimum

demands of some users is greater than others



late 1980 mid network 1990s

SLOT-BASED MODEL

Slot: Fixed quantity of CPU and memory

→ state of the art in 2009 / 2010

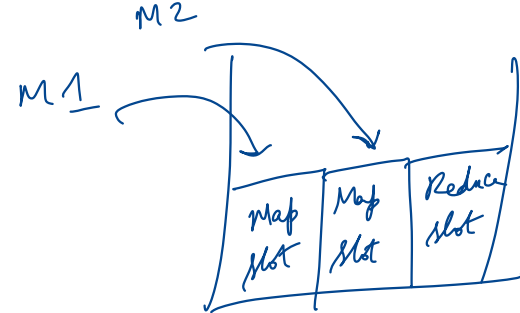
Example: Hadoop MapReduce

Mapper: 2 CPU and 1 GB

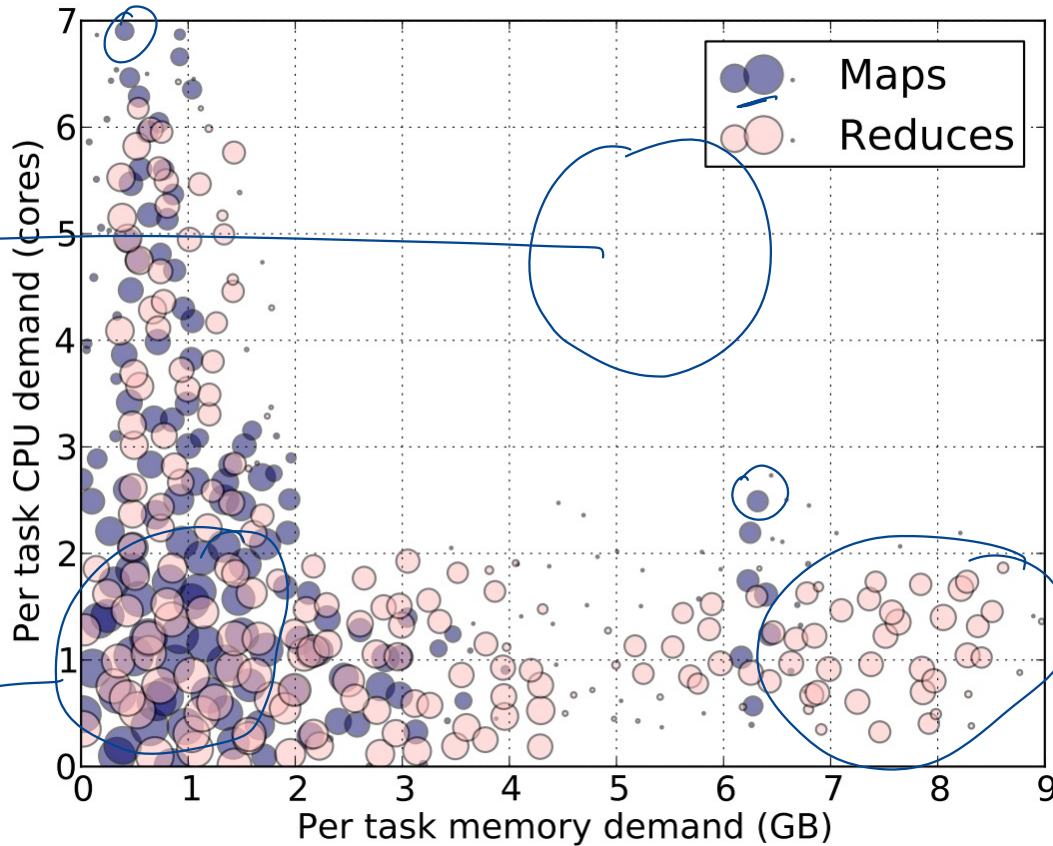
Reducer: 1 CPU and 2 GB

→ use up 1 map slot
→ 1 reduce slot

Allocate in units of slots



MOTIVATION: MULTI RESOURCES



empty

most tasks are here

diff map tasks can have very varied needs

Reduce tasks more memory intensive

DRF: MODEL

Users have a **demand vector** → generalization for multiple resources

<2, 3, 1> means user's task needs 2 R1, 3 R2, 1 R3
↓ ↓ ↘
CPU Mem Disk

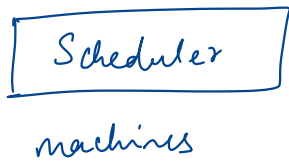
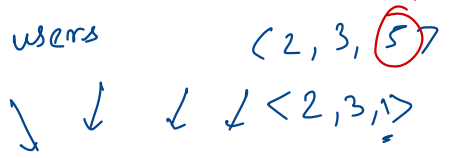
Resources given in multiples of demand vector

i.e., users might get <4,6,2>

10 map tasks in my job
<20, 30, 10> as my overall resource requirement

PROPERTIES

lies



Sharing Incentive

At least as well off as having a partitioned cluster

Strategy Proof

Users cannot benefit from lying about resources required

Pareto Efficiency

Not possible to allocate resources to one user without removing from another user.

Envy free

No user should desire another user's share

Utilization ←

PROPERTIES

Sharing Incentive

User is no worse off than a cluster with
 $1/n$ resources

Strategy Proof

User should not benefit by
lying about demands

Pareto Efficiency

Not possible to increase
one user without
decreasing another

Envy free

User should not desire the
allocation of another user

DRF: APPROACH

Dominant Resource

Resource user has the biggest share of

Total: <10 CPU, 4 GB> *Cluster*

User 1: <1 CPU, 1 GB>

Dominant resource is **memory**

$\frac{1}{10}$ th CPU
 $\frac{1}{4}$ th memory → Dominant

Dominant Share

Fraction of the dominant resource user is allocated

E.g., for User 1 this is **25% or 1/4**

Dominant share

DRF: APPROACH

Objective

→ Equalize the dominant share of users

Total: <9 CPU, 18 GB>

User1: <1 CPU, 4 GB>

dom res: mem

User2: <3 CPU, 1 GB>

dom res: CPU

| User | Allocation | Dominant Share |
|-------|---|--|
| User1 | <0 CPU, 0 GB> <1 CPU, 4 GB> <2 CPU, 8 GB> <3 CPU, 12 GB> | 0 $4/18 = 2/9$ $8/18 = 4/9$ $12/18 = 2/3$ |
| User2 | <0 CPU, 0 GB> <3 CPU, 1 GB> <6 CPU, 2 GB> | 0 $3/9 = 1/3$ $6/9 = 2/3$ |

<9 CPU, 14 GB>

<18 GB

cluster is full

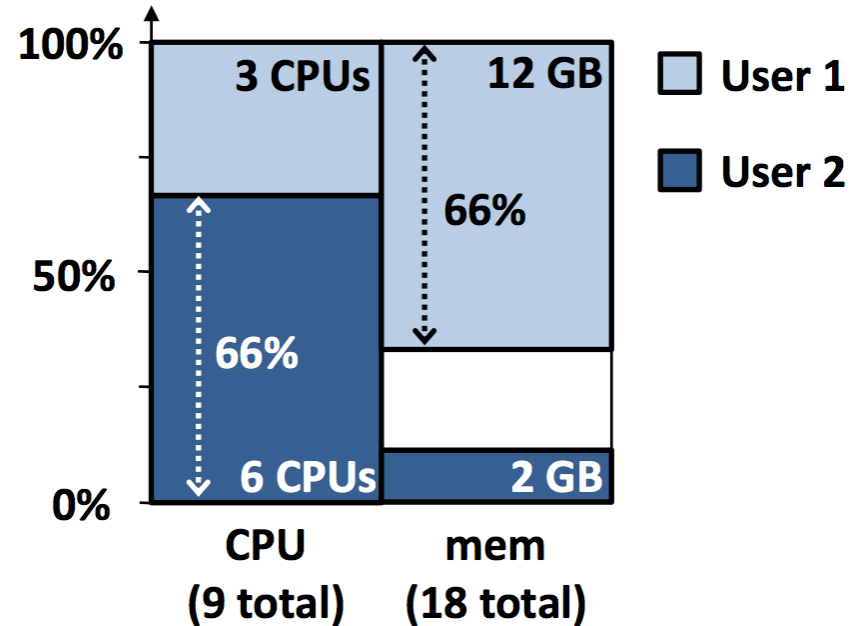


DRF: APPROACH

Total: <9 CPU, 18 GB>

User1: <1 CPU, 4 GB> per task
<3 CPU, 12 GB> for 3 tasks
dom res: mem
dom share: $12/18 = 2/3$

User2: <3 CPU, 1 GB>
<6 CPU, 2 GB> for 2 tasks
dom res: CPU
dom share: $6/9 = 2/3$



DRF ALGORITHM

Whenever there are available resources:

Schedule a task to the user with **smallest dominant share**

↓
new users who arrive
are given resources
when running tasks complete

DRF ALGORITHM

Algorithm 1 DRF pseudo-code

$R = \langle r_1, \dots, r_m \rangle$ \triangleright total resource capacities

$C = \langle c_1, \dots, c_m \rangle$ \triangleright consumed resources, initially 0

s_i ($i = 1..n$) \triangleright user i 's dominant shares, initially 0 \longrightarrow

$U_i = \langle u_{i,1}, \dots, u_{i,m} \rangle$ ($i = 1..n$) \triangleright resources given to user i , initially 0

pick user i with lowest dominant share s_i \longrightarrow

$D_i \leftarrow$ demand of user i 's next task

if $C + D_i \leq R$ **then**

$C = C + D_i$ \triangleright update consumed vector

$U_i = U_i + D_i$ \triangleright update i 's allocation vector

$s_i = \max_{j=1}^m \{u_{i,j}/r_j\}$ \longrightarrow

else

return

\triangleright the cluster is full

end if

picks user with lowest dominant share from S

is the cluster full

recompute dominant share

COMPARISON: ASSET FAIRNESS

Violating sharing incentive

Asset Fairness: Equalize each user's sum of resource shares

Consider total of 70 CPUs, 70 GB RAM

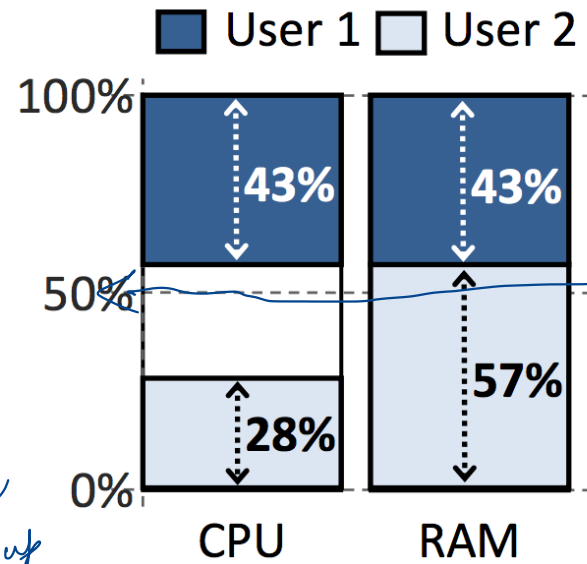
U1 needs <2 CPU, 2 GB RAM> per task

U2 needs <1 CPU, 2 GB RAM> per task

Asset Fair Allocation:

U1: 15 tasks: 30 CPU, 30 GB (Sum = 60)

U2: 20 tasks: 20 CPU, 40 GB (Sum = 60)



add up CPU & memory used by this user

COMPARISON: ASSET FAIRNESS

Asset Fairness: Equalize each user's sum of resource shares

Violates Sharing Incentive

Consider total of 70 CPUs, 70 GB RAM

U1 needs <2 CPU, 2 GB RAM> per task

U2 needs <1 CPU, 2 GB RAM> per task

Sharing incentive?

Half of the cluster is 35 CPU, 35 GB RAM

U1: 17 tasks to get 34 CPU, 34 GB

U2:

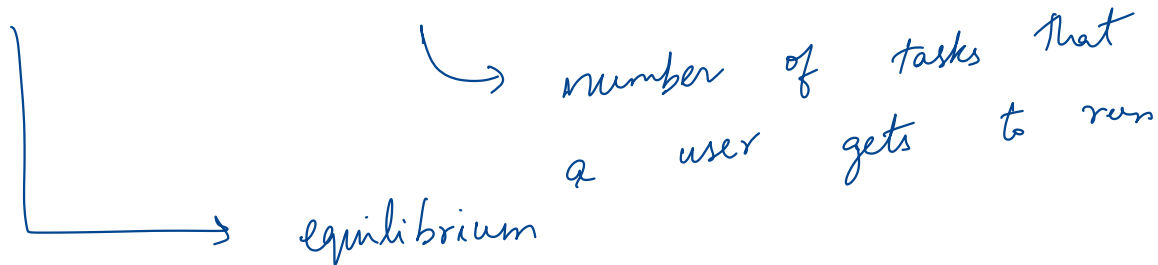
RAM ↗
17 tasks >
15 tasks with
Asset fairness

COMPARISON: CEEI

CEEI: Competitive Equilibrium from Equal Incomes

- Each user receives initially $1/n$ of every resource,
- Subsequently, each user can trade resources with other users in a perfectly competitive market
- Nash solution: Maximize **product of utilities** across users

exchange 1GB
for 1CPU



COMPARISON: CEEI

Total: <9 CPU, 18 GB>

User1: <1 CPU, 4 GB>

User2: <3 CPU, 1 GB>

x tasks

y tasks

max ($x \cdot y$) \rightarrow product utilities
subject to

CPU $\rightarrow x + 3y \leq 9$
memory $\rightarrow 4x + y \leq 18$

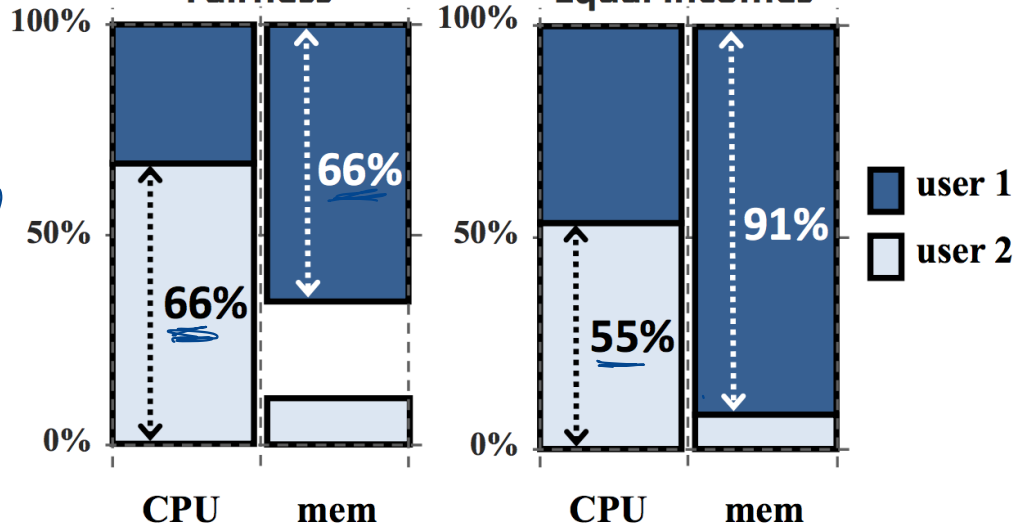
$x = 4.05$

$y = 1.62$

rounding?

Dominant Resource Fairness

Competitive Equilibrium from Equal Incomes



CEEI: STRATEGY PROOFNESS

By lying about non dominant resource user 2 gets more

Total: <9 CPU, 18 GB>

User2 Before:

CEEI: 55% CPU, 9% mem

Total: <9 CPU, 18 GB>

User1: <1 CPU, 4 GB>

User2: <3 CPU, 1 GB>

User2: <3 CPU, 2 GB>

lies about memory req

max x, y

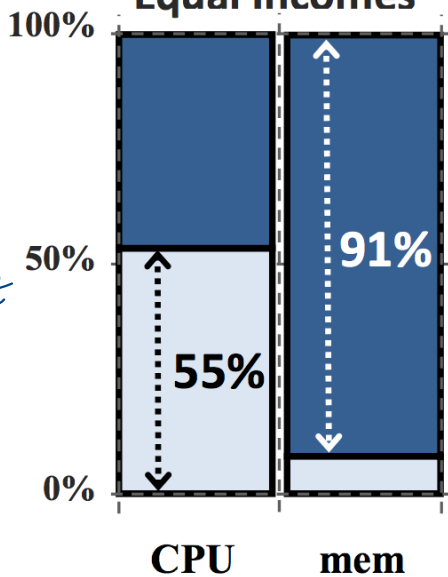
$$x + 3y \leq 9$$

$$4x + 2y \leq 18$$

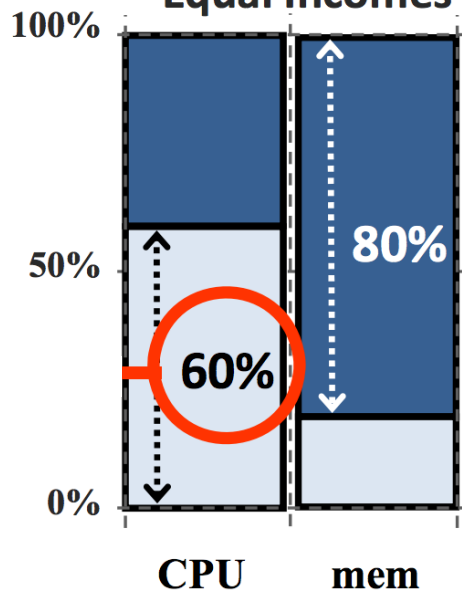
x = 3.6

y = 1.8

Competitive Equilibrium from Equal Incomes



Competitive Equilibrium from Equal Incomes



dominant resource

■ user 1
■ user 2

COMPARISON

| Property | Allocation Policy | | |
|--------------------------|-------------------|------|-----|
| | Asset | CEEI | DRF |
| Sharing Incentive | | ✓ | ✓ |
| Strategy-proofness | ✓ | | ✓ |
| Envy-freeness | ✓ | ✓ | ✓ |
| Pareto efficiency | ✓ | ✓ | ✓ |
| Single Resource Fairness | ✓ | ✓ | ✓ |
| Bottleneck Fairness | | ✓ | ✓ |
| Population Monotonicity | ✓ | | ✓ |
| Resource Monotonicity | | | ✗ |

as you
add more
resources →
to the

cluster

Table 2: Properties of Asset Fairness, CEEI and DRF.

SUMMARY

DRF: Dominant Resource Fairness

Allocation policy for scheduling

Provides multi-resource fairness

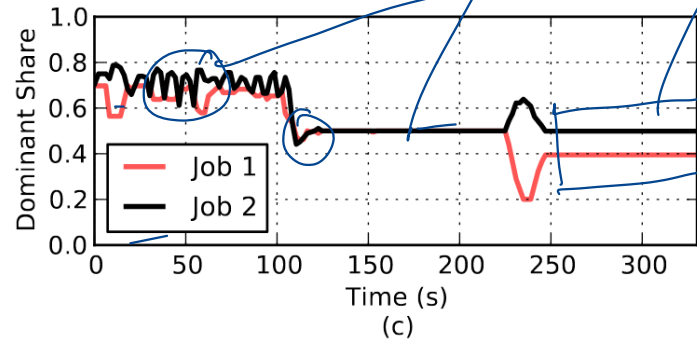
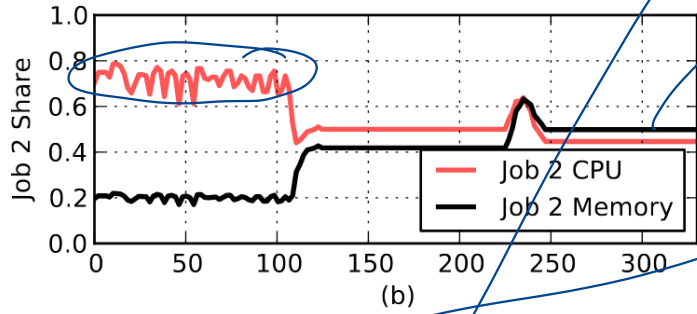
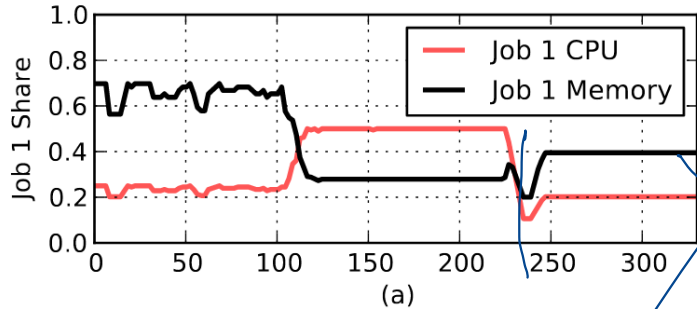
Ensures sharing incentive, strategy proofness

↳ DRF was integrated Mesos, YARN



DISCUSSION

<https://forms.gle/75faGZ4quQgVSYRQ8>



CPU is dominant share for both jobs

Mem dominant resource fragmentation / discrete allocations

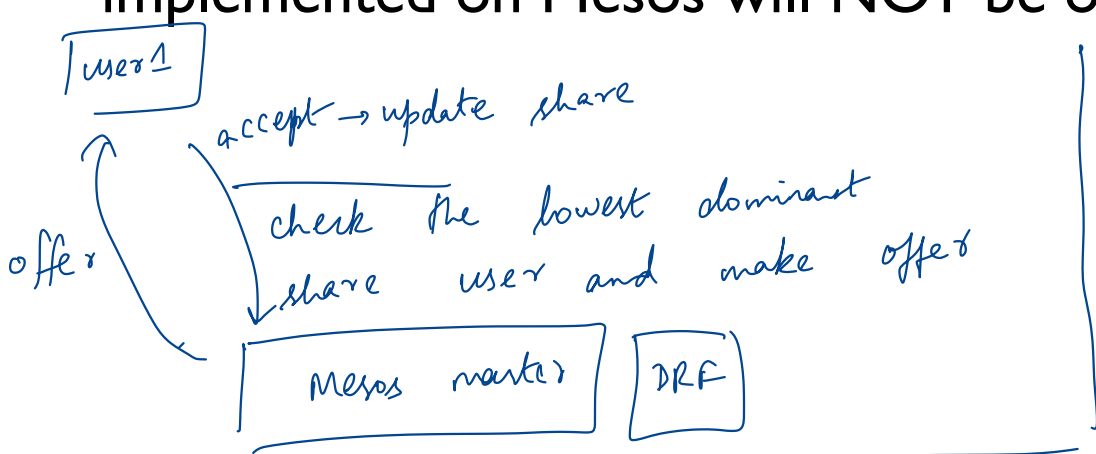
Happens when jobs are very diff.

Both dominant shares are ~0.7

DRF quickly handles change in job requirements

→ task length → pre-emption vs. waiting for tasks

What could be one workload / cluster scenario where DRF implemented on Mesos will NOT be optimal?



- DRF can only adjust shares after tasks finish
- Resources that fail back and come back frequent changes to running jobs

- Constant stream of "small" jobs starve the larger jobs

- $\langle 8 \text{ CPU}, 1 \text{ GB} \rangle$
 $\langle 5 \text{ CPU}, 4 \text{ GB} \rangle$

- Resource demands that flip very often \rightarrow might not converge
- Locality

NEXT STEPS

Next class: Machine Learning Schedulers