

# Preemptive Scheduling in Mesos Framework

Li Jin

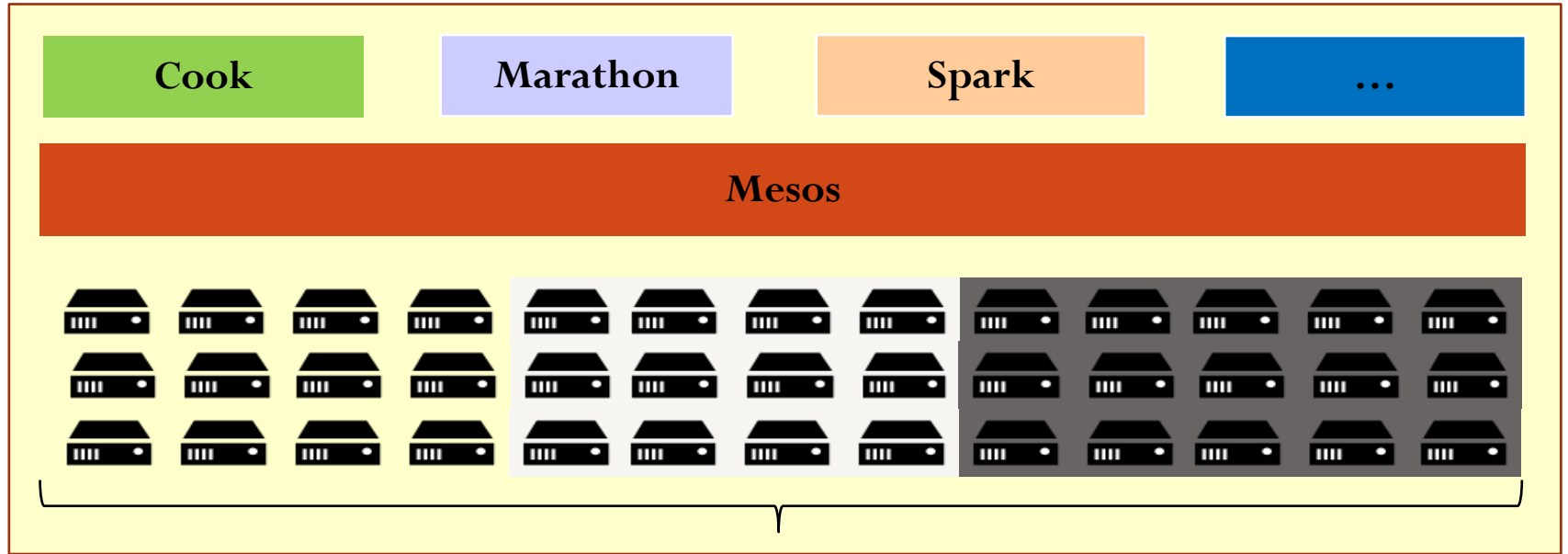
# About Two Sigma



# About Me

- Software Engineer @ Two Sigma

# Mesos @ TS



# Outline

- Cook: A Mesos Framework
- Problem: Utilization and Fairness
- Preemption: Intuition and Formalization
- Preemption in Cook: Implementation

# What is Cook

- Two Sigma's *Simulation* Platform
  - Manage tens of thousands of simulations
  - Share compute resource among users

# What is Simulation

- Idempotent, parallel, distributed, resource intensive computations
- One simulation = Multiple Mesos tasks

# What is Simulation

- Simulation task footprint
  - 10 ~ 100 GB RAM
  - 1 ~ 20 CPUs
  - 15 minutes ~ a few hours
- Simulation use cases
  - Interactive Research
  - Batch computation



# Outline

- Cook: A Mesos Framework
- Problem: Utilization and Fairness
- Preemption: Intuition and Formalization
- Preemption in Cook: Implementation

# Problem

- High computation demand
  - 5 x capacity during peak hours
- Optimize
  - Utilization
    - Use as much compute resources as possible
  - Fairness
    - Allocate resources fairly for some definition of 'fair'

# What is Fairness

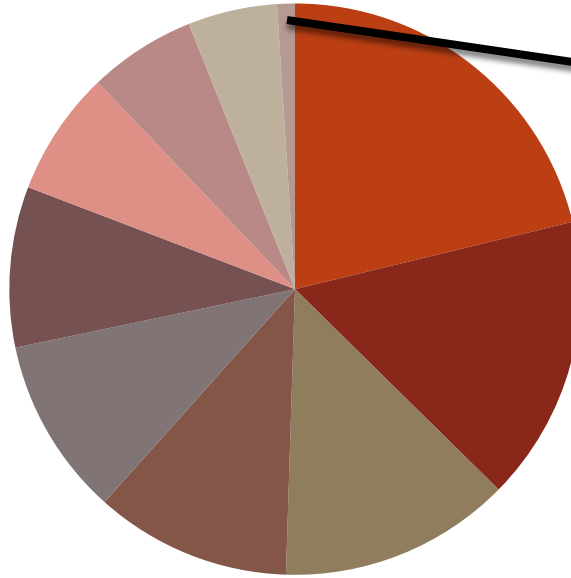
- FIFO
- Time sharing
- Throw a dice
- ...

# What is Fairness

- A story...

# What is Fairness

## Resource Allocation



# What is Fairness, Really

- Fairness is not about 'fair'
- Fairness is about **user experience**
  - User should get their share of the cluster **whenever** they need it

# Quota

- Quota = Maximum percentage of the cluster allowed
- Static
  - $100\% / \# \text{ Max concurrent users}$
- Pros:
  - Guarantee Fairness
- Cons:
  - Low Utilization

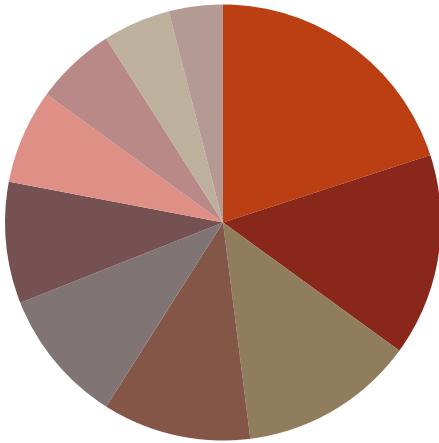
# Dynamic Quota

- Dynamic
  - $\text{Quota} * \text{Utilization Adjustment}$
- Pros:
  - Higher Utilization
- Cons:
  - Slow reaction to change of demand



# Dynamic Quota

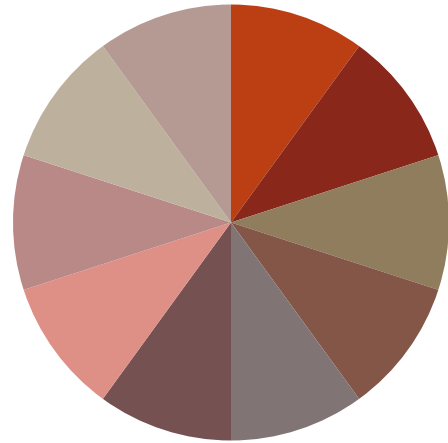
**Unfair Resource  
Allocation**



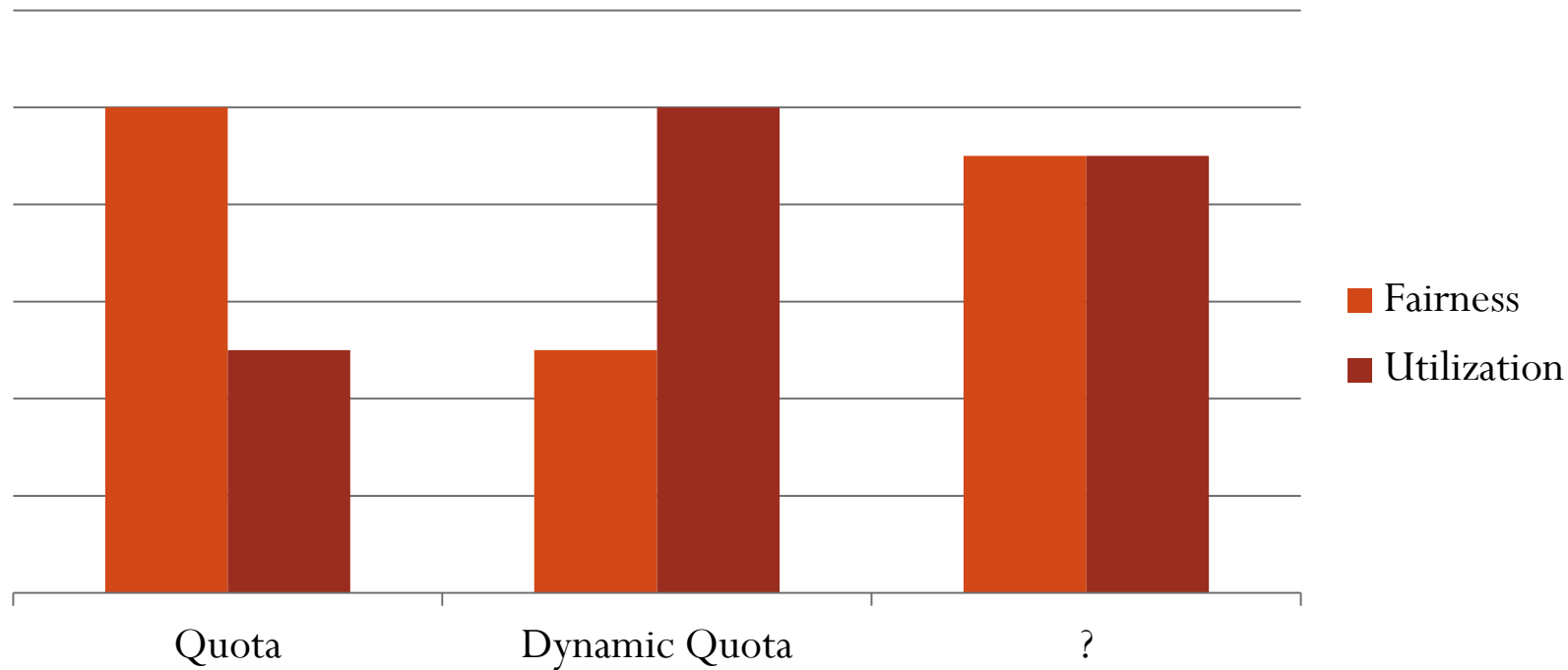
Hours...



**Fair Resource  
Allocation**



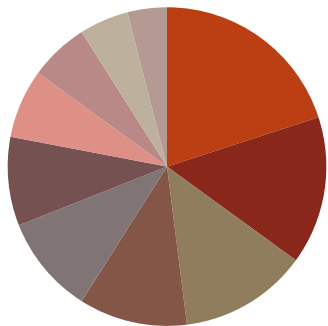
# Can we do better?



# Preemption

- Kill a Mesos task and reschedule later
- Reclaim resource faster!

**Unfair Resource Allocation**



Minutes!



**Fair Resource Allocation**



# Outline

- Cook: A Mesos Framework
- Problem: Utilization and Fairness
- Preemption: **Intuition** and Formalization
- Preemption in Cook: Implementation

# Preemption: Intuition

Waiting

Running



**JERRY**



**KEVIN**



# Preemption: Intuition

Waiting



Running



**JERRY**



**KEVIN**



**DAVE**



# Preemption: Intuition

Waiting



---

Running



**JERRY**



**KEVIN**



**DAVE**

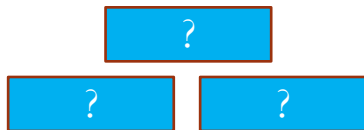


# Preemption: Intuition

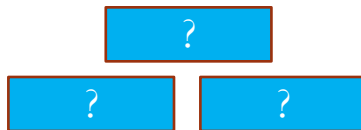
Waiting



Running



**JERRY**



**KEVIN**



**DAVE**





# Preemption: Intuition

Waiting



Running



**JERRY**



**KEVIN**

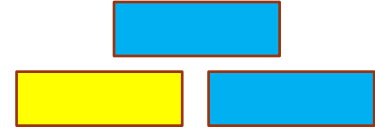


**DAVE**

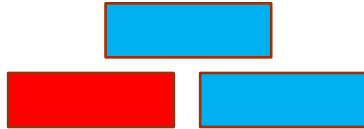


# Preemption: Intuition

Waiting



Running



**JERRY**



**KEVIN**



**DAVE**



# Preemption: Intuition

Waiting



Running



**JERRY**



**KEVIN**



**DAVE**



# Preemption: Intuition

Waiting



---

Running



**JERRY**



**KEVIN**



**DAVE**



# Preemption: Intuition

Waiting



---

Running



**JERRY**



**KEVIN**



**DAVE**



# Preemption: Intuition

Waiting



---

Running



**JERRY**



**KEVIN**



**DAVE**



# Problem

- Not all jobs are equal
  - We just preempted some important jobs!



Bad User  
Experience



# Score Function

- Reflect task's value
  - Fairness
  - Importance
- Preempt low value task for high value task



# Preemption: Intuition

Waiting

¥¥
¥¥¥
¥¥¥¥

Running

€
€€
€€€

JERRY



£££
££££
£££££

KEVIN



DAVE



# Preemption: Intuition

Waiting

\$
\$ \$
\$ \$ \$

Running

\$
\$ \$
\$ \$ \$

**JERRY**



\$
\$ \$
\$ \$ \$

**KEVIN**



**DAVE**



# Preemption: Intuition

Waiting

\$
\$\$
\$\$\$

Running

\$
\$\$
\$\$\$

**JERRY**



\$
\$\$
\$\$\$

**KEVIN**



**DAVE**



# Preemption: Intuition

Waiting

\$
\$ \$
\$ \$ \$

Running

\$
\$ \$
\$ \$ \$

**JERRY**



\$
\$ \$
\$ \$ \$

**KEVIN**



**DAVE**



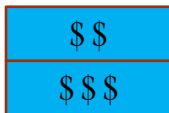
# Preemption: Intuition

Waiting

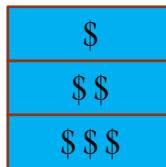


---

Running



**JERRY**



**KEVIN**



**DAVE**



# Preemption: Intuition

Waiting



---

Running



**JERRY**



**KEVIN**



**DAVE**



# Preemption: Intuition

Waiting



---

Running



**JERRY**



**KEVIN**



**DAVE**



# Outline

- Cook: A Mesos Framework
- Problem: Utilization and Fairness
- Preemption: Intuition and **Formalization**
- Preemption in Cook: Implementation



# Cumulative Resource Share (CRS)

- Assuming there is an total order of jobs for *each user*, where  $>$  means ‘has higher value than’.
  - CRS of job  $j$  is sum of all jobs of the same user that are greater than or equal to  $j$ , divided by total resource.
- $CRS(j) = \frac{1}{R_{Total}} \sum_{j' \geq j} R_{j'}$

# Preemption: Formalization

Waiting

3/6
2/6
1/6

Running

3/6
2/6
1/6

**JERRY**



3/6
2/6
1/6

**KEVIN**



**DAVE**



# Preemption: Formalization

Waiting

3/6
2/6
1/6

Running

3/6
2/6
1/6

**JERRY**



3/6
2/6
1/6

**KEVIN**



**DAVE**



# Preemption: Formalization

Waiting

3/6

3/6

2/6

Running

2/6

1/6

JERRY



3/6

2/6

1/6

KEVIN



1/6

DAVE



# Preemption: Formalization

Waiting

3/6

3/6

3/6

---

Running

2/6

1/6

**JERRY**



2/6

1/6

**KEVIN**



2/6

1/6

**DAVE**



# Cumulative Resource Share

# Multiple Resources?

- Dominant Resource Fairness: Fair Allocation of Multiple Resource Types
  - UC Berkeley in 2011
  - Used in Mesos
- Dominant Resource Share:
  - $DRS(u) = \max_R \frac{R_u}{R_{Total}}$

# Dominant Cumulative Resource Share

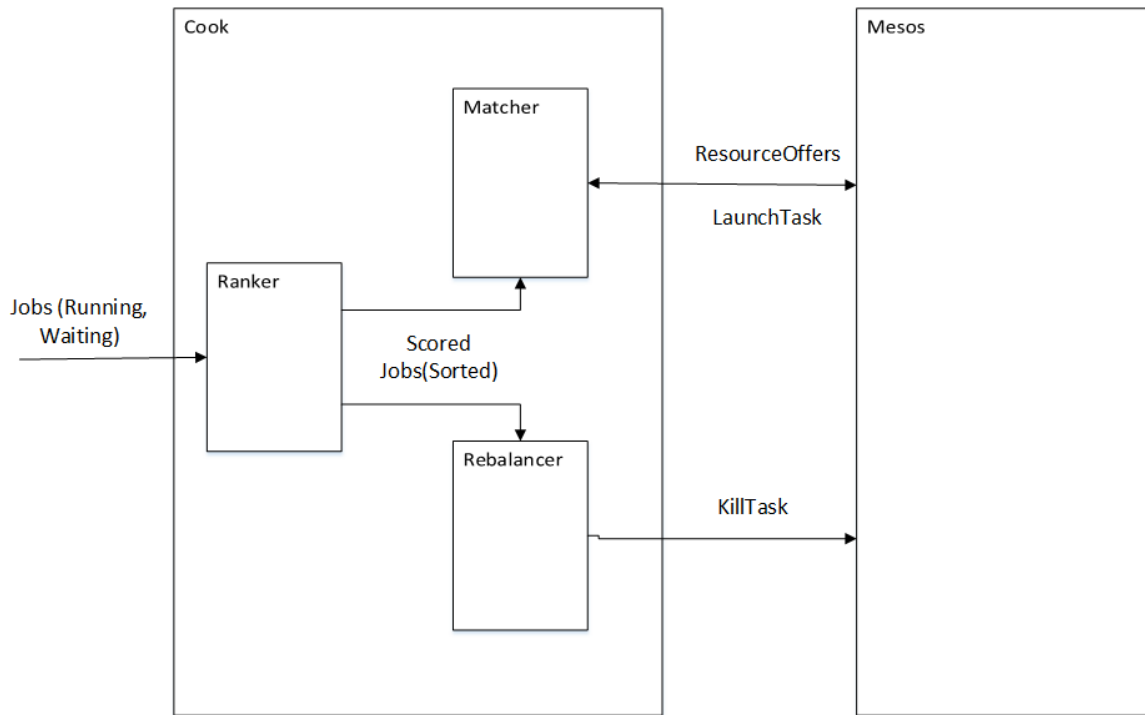
- $CRS(j) = \frac{1}{R_{Total}} \sum_{j' \geq j} R_{j'}$
- $DCRS(j) = \max_R \frac{1}{R_{Total}} \sum_{j' \geq j} R_{j'}$
- $Score(j) = -DCRS(j)$



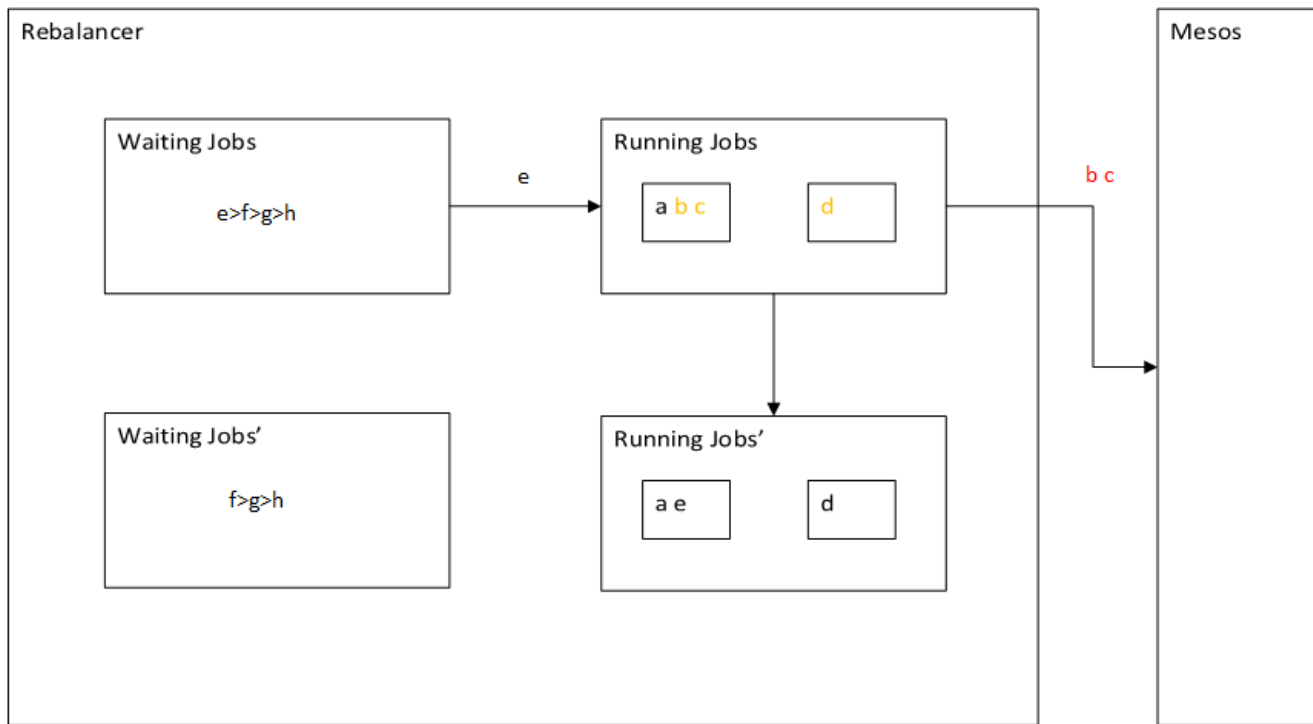
# Outline

- Cook: A Mesos Framework
- Problem: Utilization and Fairness
- Preemption: Intuition and Formalization
- Preemption in Cook: Implementation

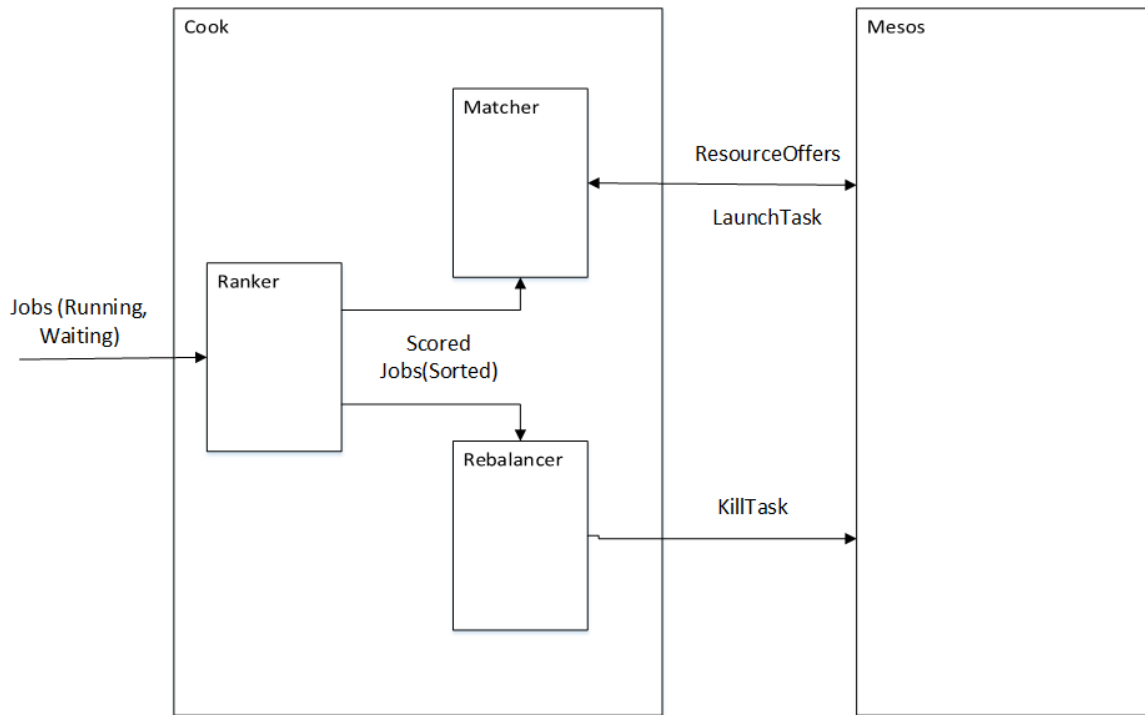
# Cook: Architecture



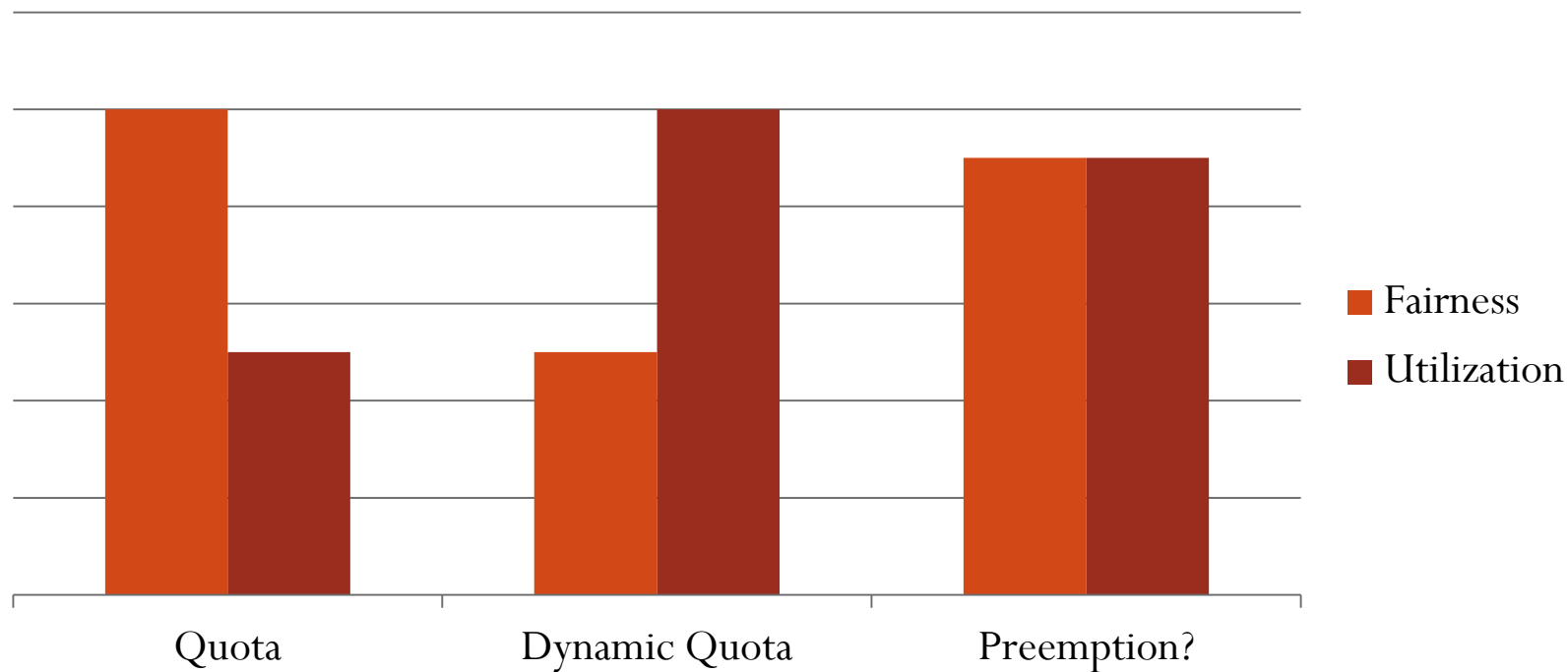
# Rebalancer



# Cook: Architecture



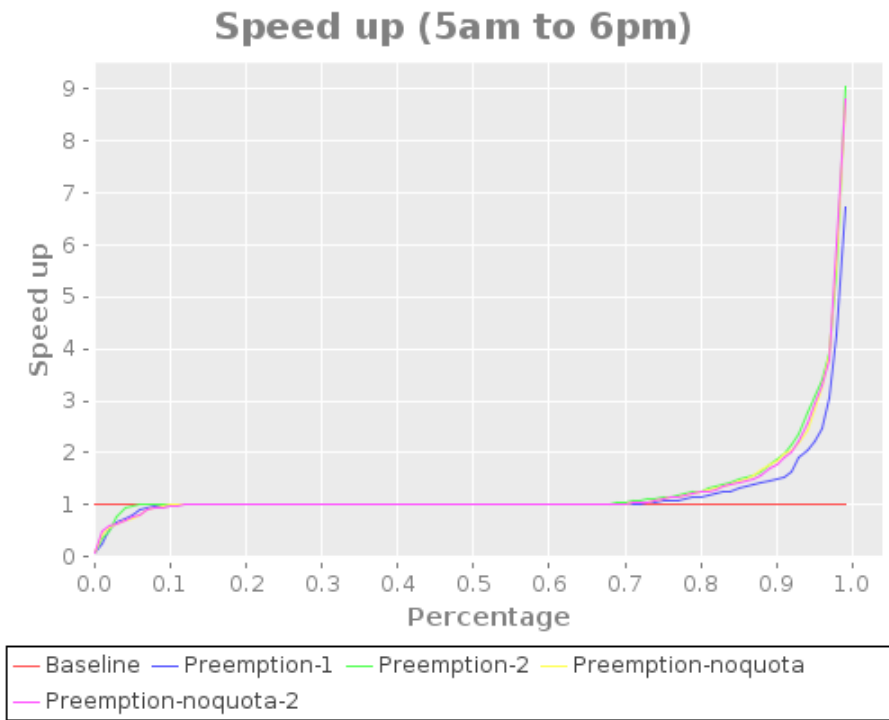
# Are we doing better?



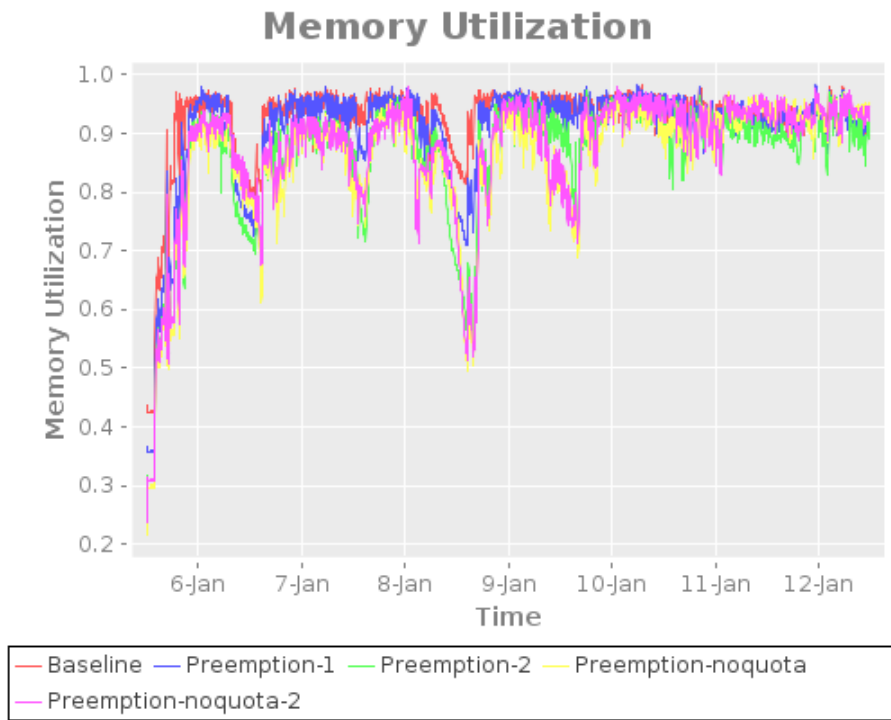
# Yes!



# It works!



# It works!





# Think Beyond Cook

- DCRS in Mesos?

# Open Source

- Not quite yet... But we are getting close
- [Github.com/twosigma](https://github.com/twosigma)