

CS162  
Operating Systems and  
Systems Programming  
Lecture 24

Chord and Berkeley Data Analytics Stack (part I)

November 20<sup>th</sup>, 2017

Prof. Ion Stoica

<http://cs162.eecs.Berkeley.edu>

Outline

- Quorum consensus

- Chord

- Mesos

11/19/17

CS162 © UCB Fall 2017

Lec 24.2

Quorum Consensus

- Improve put() and get() operation performance
- Define a replica set of size N
  - put() waits for acknowledgements from at least W replicas
  - get() waits for responses from at least R replicas
  - $W+R > N$
- Why does it work?
  - There is at least one node that contains the update
- Why might you use  $W+R > N+1$ ?

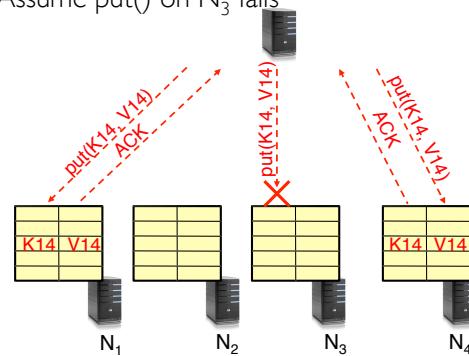
11/19/17

CS162 © UCB Fall 2017

Lec 24.3

Quorum Consensus Example

- $N=3, W=2, R=2$
- Replica set for K14:  $\{N_1, N_3, N_4\}$
- Assume put() on  $N_3$  fails



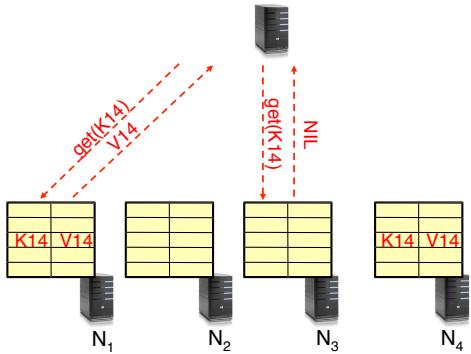
11/19/17

CS162 © UCB Fall 2017

Lec 24.4

## Quorum Consensus Example

- Now, issuing get() to any two nodes out of three will return the answer



11/19/17

CS162 © UCB Fall 2017

Lec 24.5

## Outline

- Quorum consensus

- Chord

- Mesos

11/19/17

CS162 © UCB Fall 2017

Lec 24.6

## Scaling Up Directory – Challenges

- Directory contains a number of entries equal to number of (key, value) tuples in the system
- Can be tens or hundreds of billions of entries in the system!

11/19/17

CS162 © UCB Fall 2017

Lec 24.7

## Scaling Up Directory – Strawman Solution

- Assume a system with  $m$  nodes
- Store  $(key, value)$  on node  $i = key \bmod M + 1$
- No need to keep any metadata!
- Challenge: what happened if you take away a node or add a new node?



11/19/17

CS162 © UCB Fall 2017

Lec 24.8

## Scaling Up Directory – Consistent Hashing

- Associate to each node a unique *id* in an *uni*-dimensional space  $0..2^m - 1$ 
  - Partition this space across  $m$  machines
  - Assume keys are in same uni-dimensional space
  - Each (Key, Value) is stored at the node with the smallest ID larger than Key

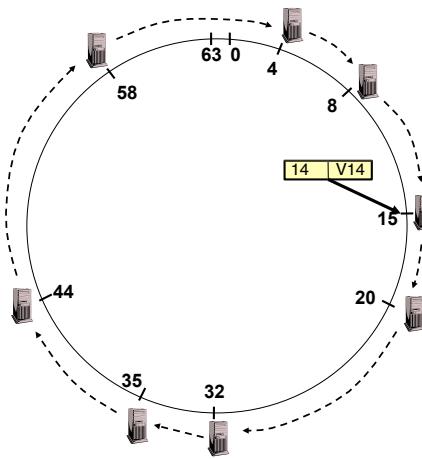
11/19/17

CS162 © UCB Fall 2017

Lec 24.9

## Key to Node Mapping Example

- $m = 6 \rightarrow$  ID space: 0..63
- Node 8 maps keys [5, 8]
- Node 15 maps keys [9, 15]
- Node 20 maps keys [16, 20]
- ...
- Node 4 maps keys [59, 4]



11/19/17

CS162 © UCB Fall 2017

Lec 24.10

## Scaling Up Directory

- With consistent hashing, directory contains only a number of entries equal to number of nodes
  - Much smaller than number of tuples
- Next challenge: every query still needs to contact the directory
- Solution: distributed directory (a.k.a. lookup) service:
  - Given a **key**, find the **node** storing value associated to the key
- Key idea: route request from node to node until reaching the node storing the request's key
- Key advantage: totally distributed
  - No point of failure; no hot spot

11/19/17

CS162 © UCB Fall 2017

Lec 24.11

## Chord: Distributed Lookup (Directory) Service

- Key design decision
  - Decouple correctness from efficiency
- Properties
  - Each node needs to know about  $O(\log(M))$ , where  $M$  is the total number of nodes
  - Guarantees that a tuple is found in  $O(\log(M))$  steps
- Many other lookup services: CAN, Tapestry, Pastry, Kademlia, ...

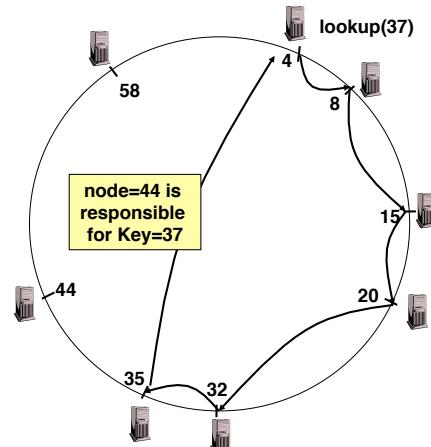
11/19/17

CS162 © UCB Fall 2017

Lec 24.12

## Lookup

- Each node maintains pointer to its successor
- Route packet (Key, Value) to the node responsible for ID using successor pointers
- E.g., node=4 looks up for node responsible for Key=37



11/19/17

CS162 © UCB Fall 2017

Lec 24.13

## Stabilization Procedure

- Periodic operation performed by each node  $n$  to maintain its successor when new nodes join the system

```

n.stabilize()
  x = succ.pred;
  if (x ∈ (n, succ))
    succ = x; // if x better successor, update
    succ.notify(n); // n tells successor about itself

n.notify(n')
  if (pred = nil or n' ∈ (pred, n))
    pred = n'; // if n' is better predecessor, update
  
```

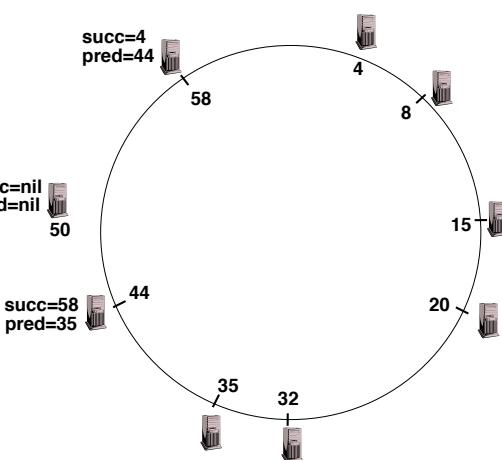
11/19/17

CS162 © UCB Fall 2017

Lec 24.14

## Joining Operation

- Node with id=50 joins the ring
- Node 50 needs to know at least one node already in the system
  - Assume known node is 15



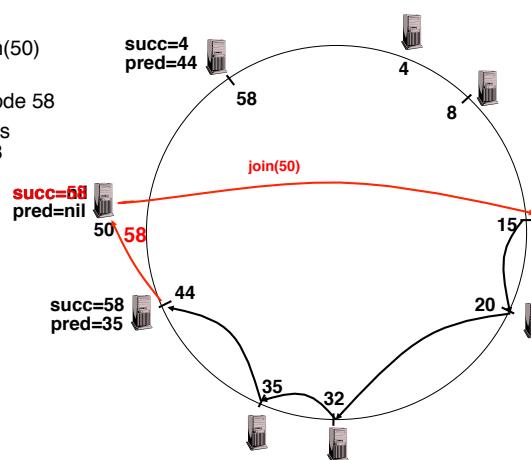
11/19/17

CS162 © UCB Fall 2017

Lec 24.15

## Joining Operation

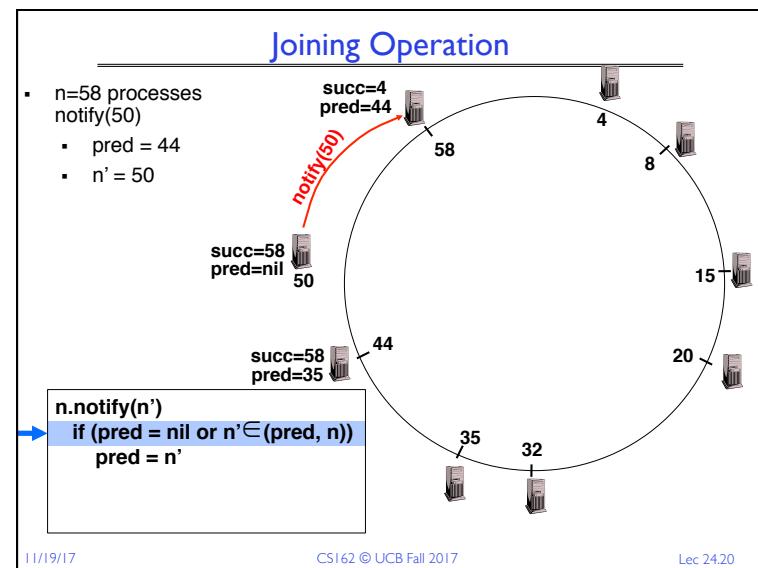
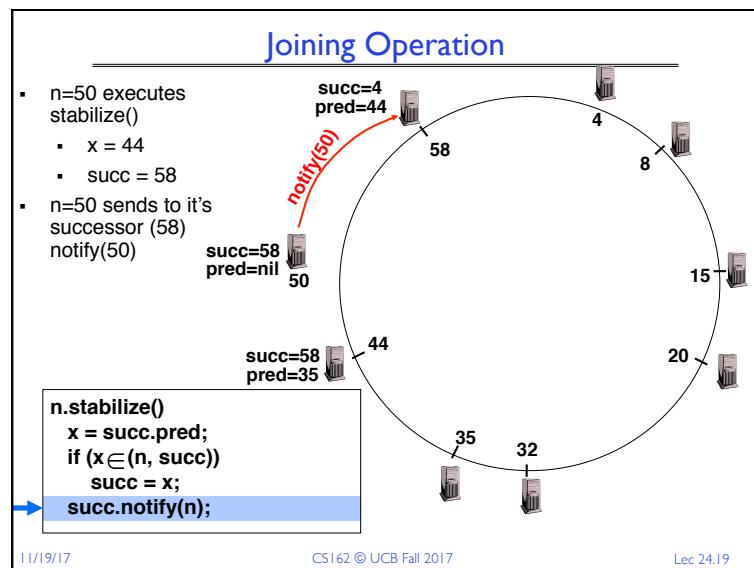
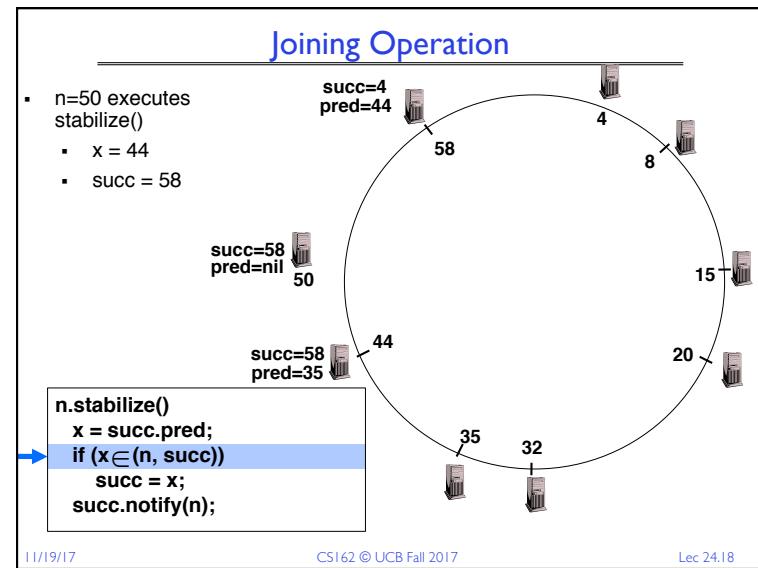
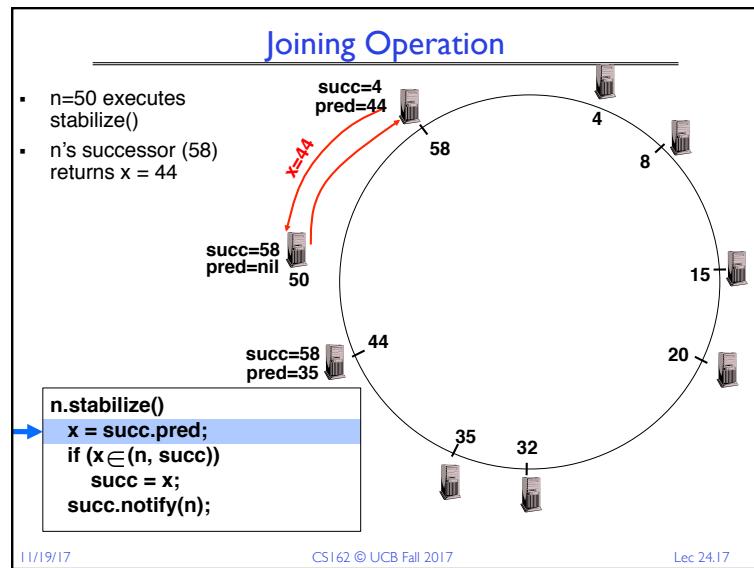
- $n=50$  sends  $\text{join}(50)$  to node 15
- $n=44$  returns node 58
- $n=50$  updates its successor to 58

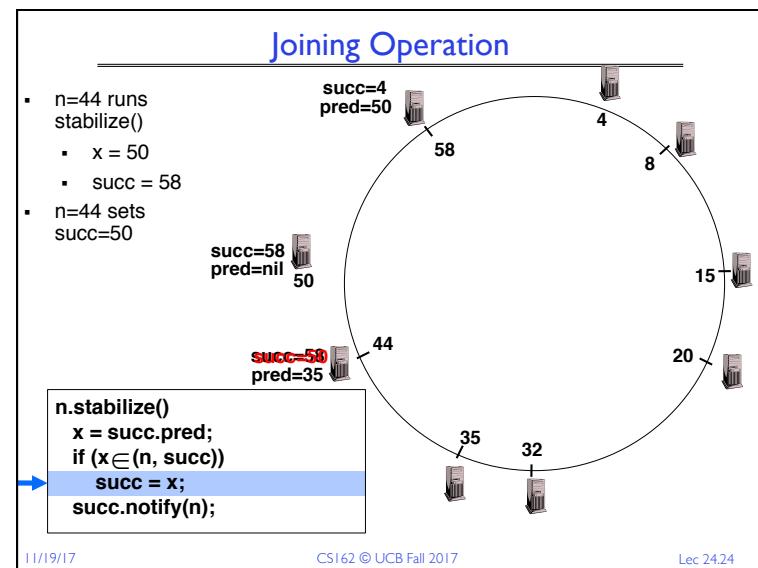
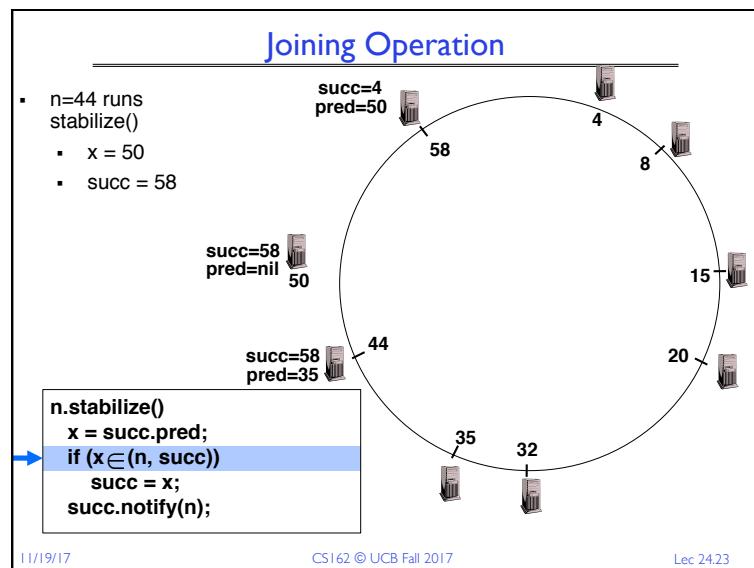
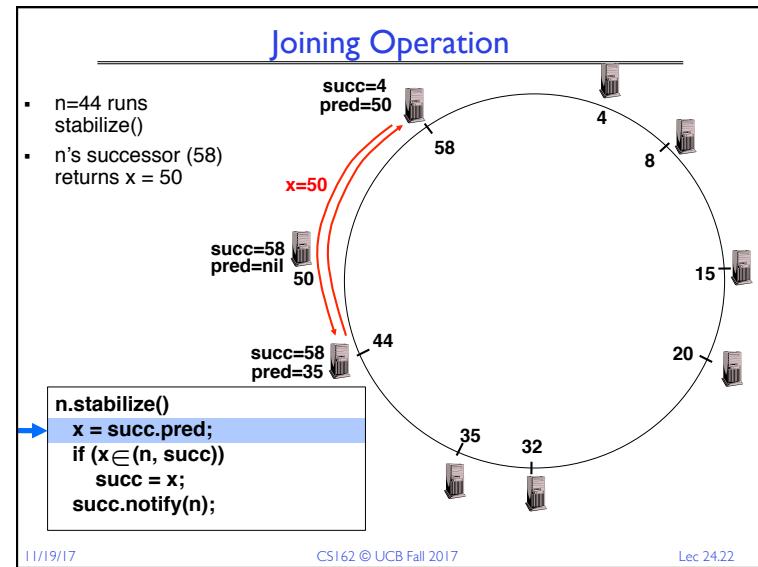
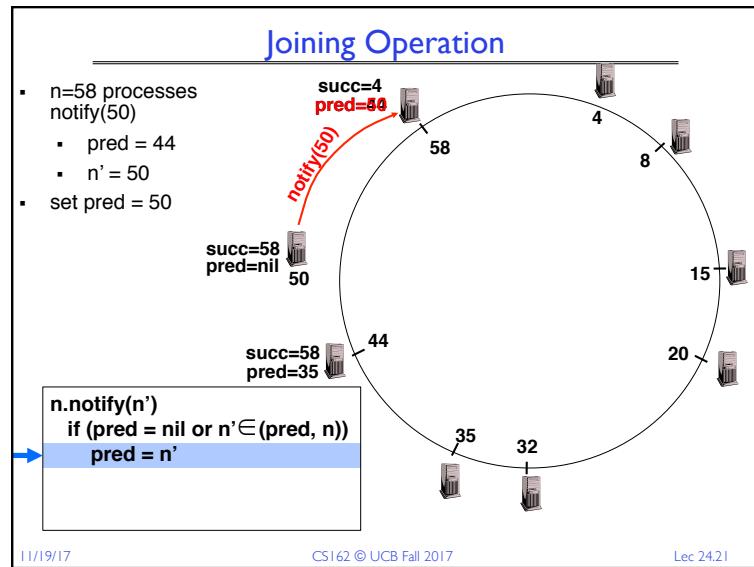


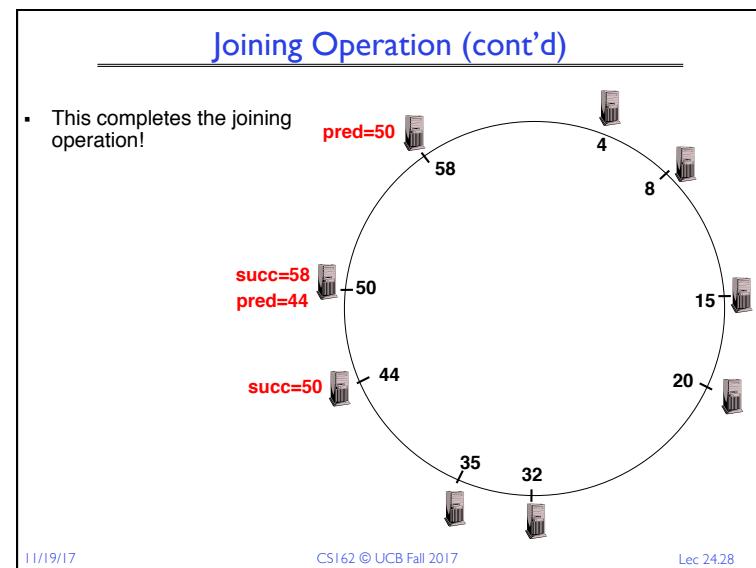
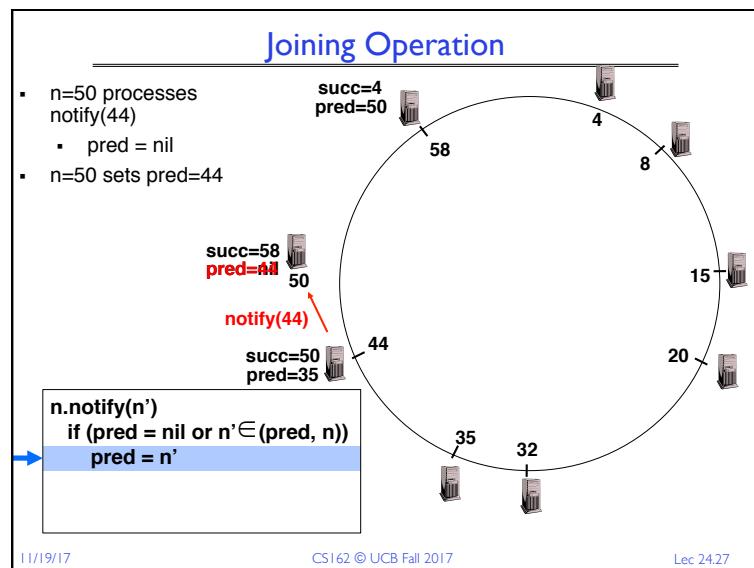
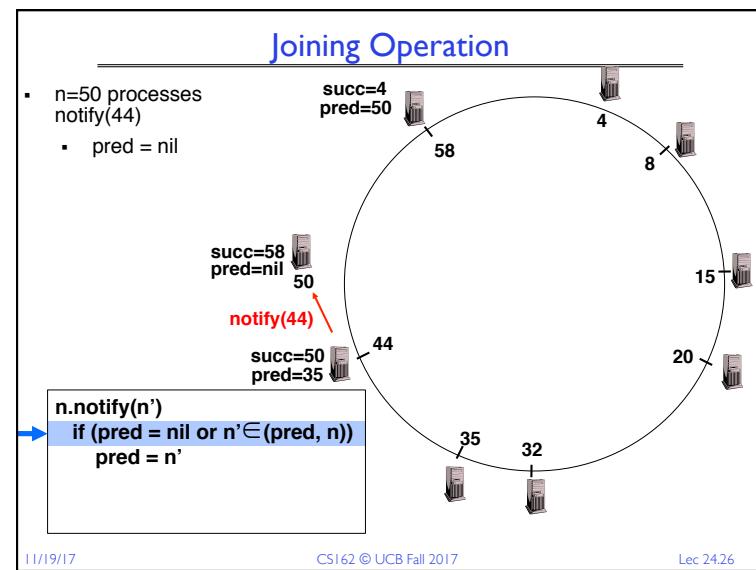
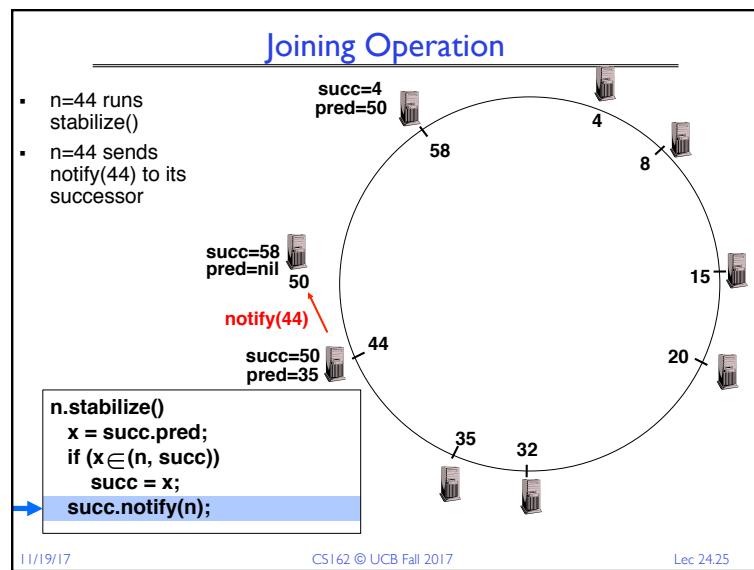
11/19/17

CS162 © UCB Fall 2017

Lec 24.16





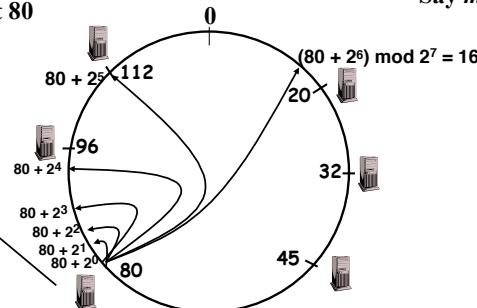


## Achieving Efficiency: finger tables

Finger Table at 80

$i$	$fi[i]$
0	96
1	96
2	96
3	96
4	96
5	112
6	20

Say  $m=7$



11/19/17

CS162 © UCB Fall 2017

Lec 24.29

## Achieving Fault Tolerance for Lookup Service

- To improve robustness each node maintains the  $k (> 1)$  immediate successors instead of only one successor
- In the `pred()` reply message, node A can send its  $k-1$  successors to its predecessor B
- Upon receiving `pred()` message, B can update its successor list by concatenating the successor list received from A with its own list
- If  $k = \log(M)$ , lookup operation works with high probability even if half of nodes fail, where M is number of nodes in the system

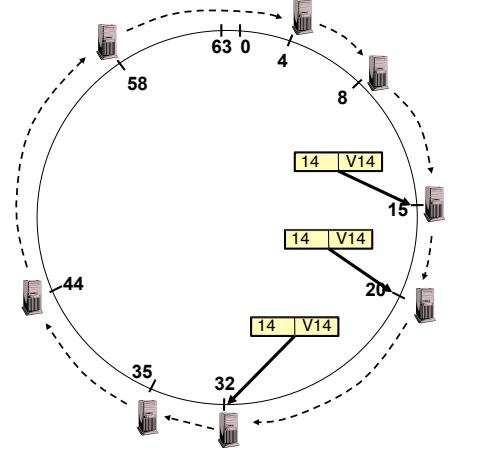
11/19/17

CS162 © UCB Fall 2017

Lec 24.30

## Storage Fault Tolerance

- Replicate tuples on successor nodes
- Example: replicate (K14, V14) on nodes 20 and 32



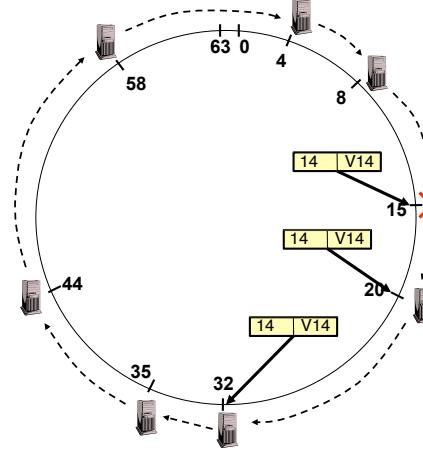
11/19/17

CS162 © UCB Fall 2017

Lec 24.31

## Storage Fault Tolerance

- If node 15 fails, no reconfiguration needed
  - Still have two replicas
  - All lookups will be correctly routed
- Will need to add a new replica on node 35



11/19/17

CS162 © UCB Fall 2017

Lec 24.32

## Administrivia

- Midterm 3 coming up on **Wen 11/29 6:30-8PM**
  - All topics up to and including Lecture 24
    - » Focus will be on Lectures 17 – 25 and associated readings, and Projects 3
    - » But expect 20-30% questions from materials from Lectures 1-16
  - Closed book
  - 2 sides hand-written notes both sides

11/19/17

CS162 © UCB Fall 2017

Lec 24.33

## BREAK

11/19/17

CS162 © UCB Fall 2017

Lec 24.34

## Outline

- Quorum consensus
- Chord
- Mesos

11/19/17

CS162 © UCB Fall 2017

Lec 24.35

## Data Deluge

- Billions of users connected through the net
  - WWW, FB, twitter, cell phones, ...
  - 80% of the data on FB was produced last year
  - FB building Exabyte ( $2^{60} \approx 10^{18}$ ) data centers
- It's all happening online – could record every:
  - Click, ad impression, billing event, server request, transaction, network msg, fault, fast forward, pause, skip, ...
- User Generated Content (Web & Mobile)
  - Facebook, Instagram, Yelp, TripAdvisor, Twitter, YouTube, ...

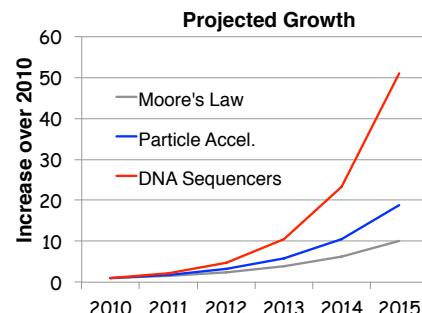


11/19/17

CS162 © UCB Fall 2017

Lec 24.36

## Data Grows Faster than Moore's Law



And Moore's law is ending!!

11/19/17

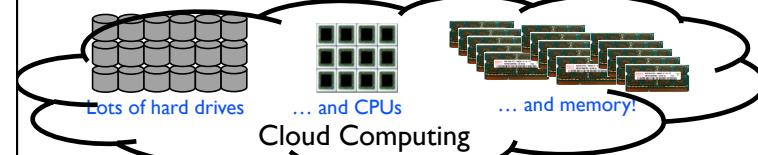
CS162 @ UCB Fall 2017

Lec 24.37

## The Big Data Solution: Cloud Computing

- One machine can not process or even store all the data!

- Solution: distribute data over cluster of cheap machines



- Cloud Computing provides:

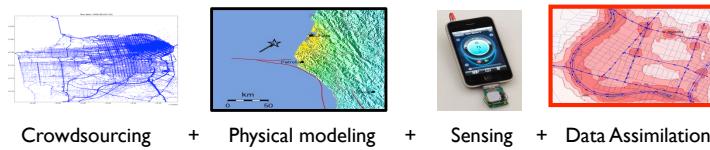
- Illusion of infinite resources
- Short-term, on-demand resource allocation
- Can be much less expensive than owning computers
- Access to latest technologies (SSDs, GPUs, ...)

11/19/17

CS162 @ UCB Fall 2017

Lec 24.38

## What Can You do with Big Data?



<http://traffic.berkeley.edu>

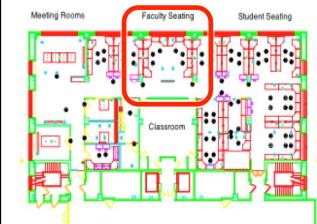
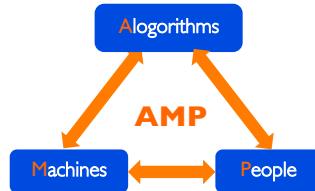
11/19/17

CS162 @ UCB Fall 2017

Lec 24.39

## The Berkeley AMPLab

- January 2011 – 2016
  - 8 faculty
  - > 50 students
  - 3 software engineer team
- Organized for collaboration



AMPCamp (since 2012)

400+ campers  
(100s companies)

Lec 24.40

## The Berkeley AMPLab

- Governmental and industrial funding:



Goal: Next generation of open source data analytics stack for industry & academia:  
Berkeley Data Analytics Stack (BDAS)

11/19/17

CS162 © UCB Fall 2017

Lec 24.41

## Generic Big Data Stack

Processing Layer

Resource Management Layer

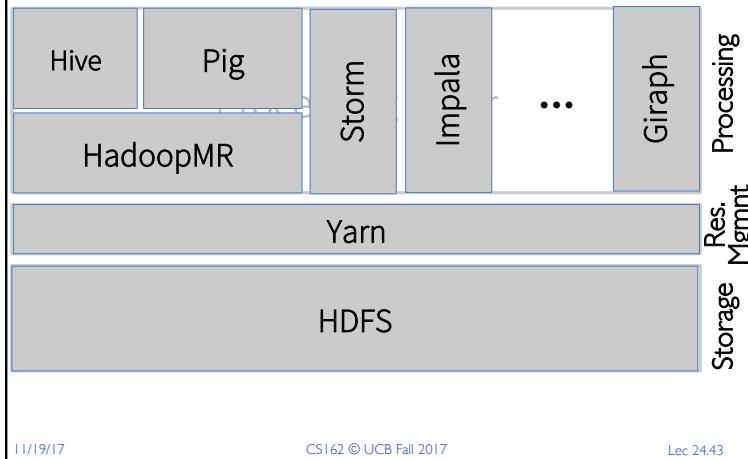
Storage Layer

11/19/17

CS162 © UCB Fall 2017

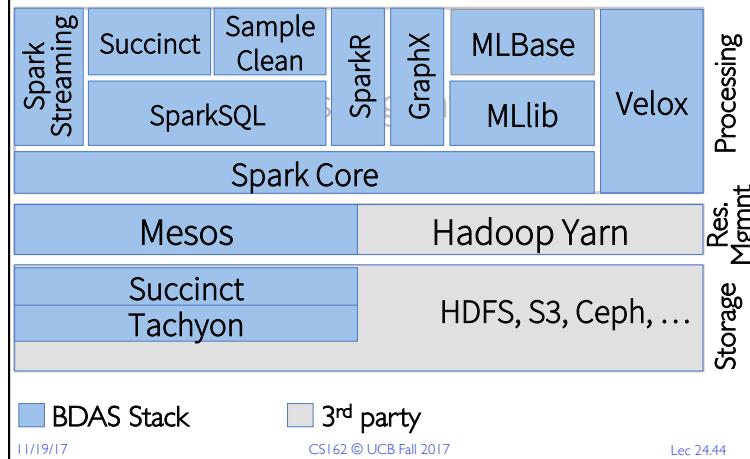
Lec 24.42

## Hadoop Stack



Lec 24.43

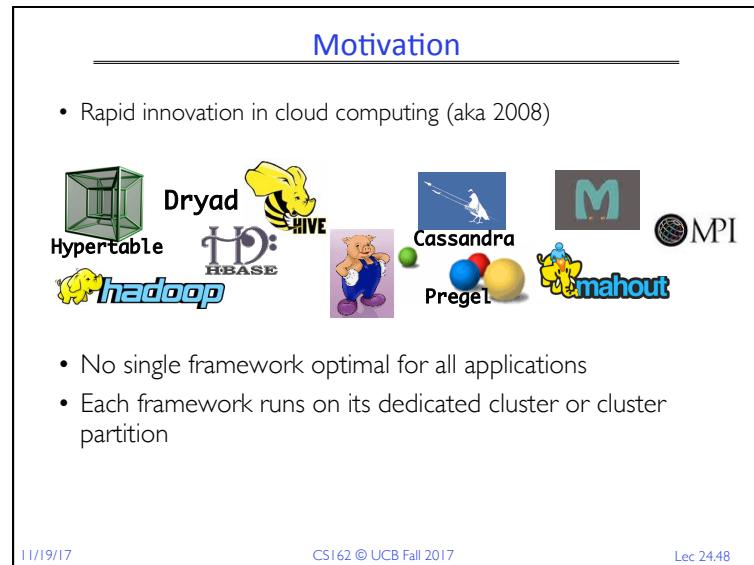
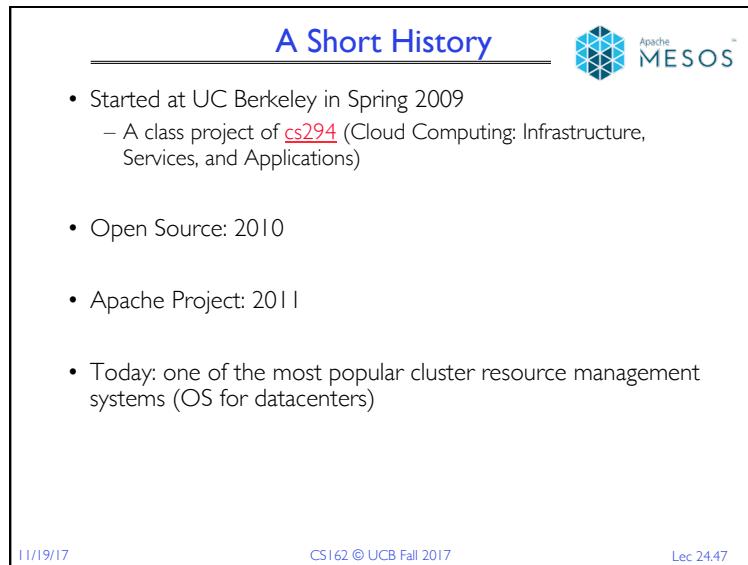
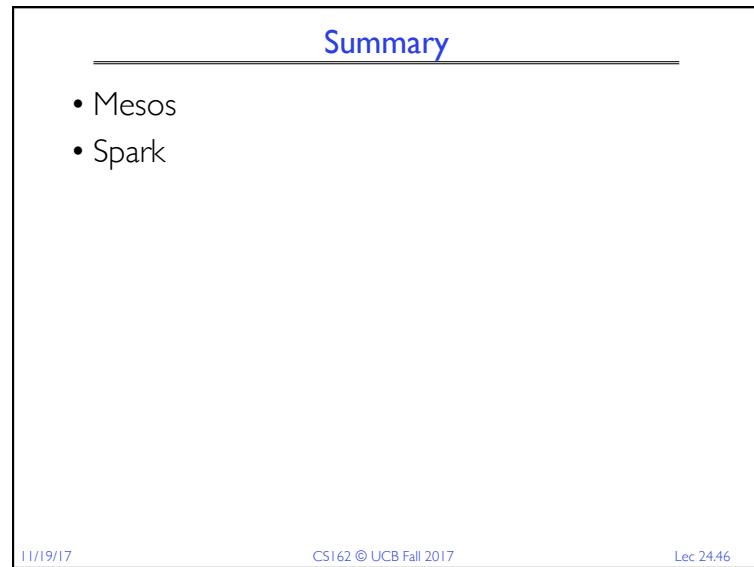
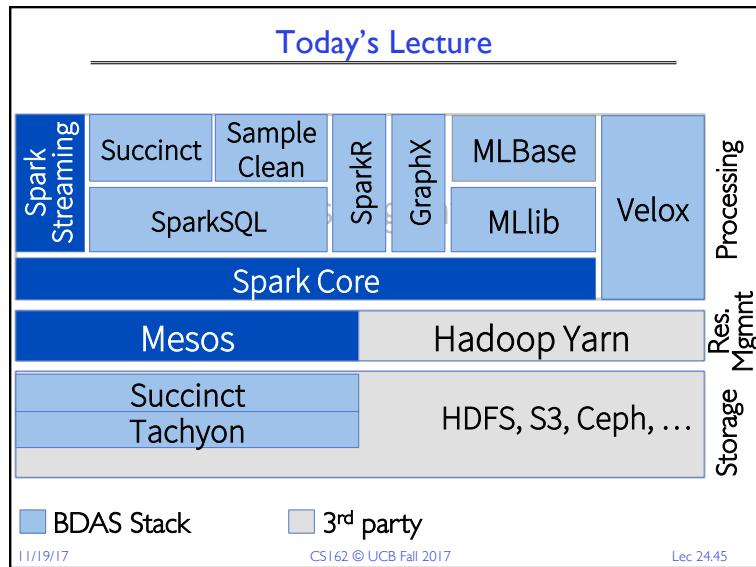
## BDAS Stack



3rd party

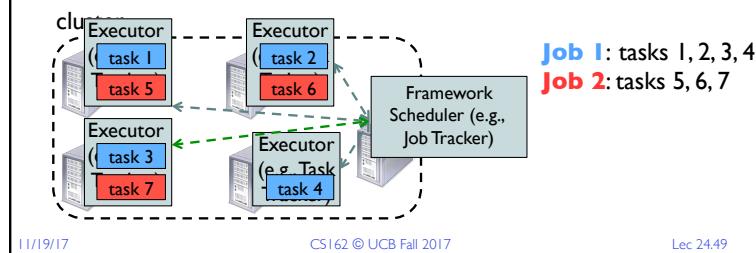
CS162 © UCB Fall 2017

Lec 24.44



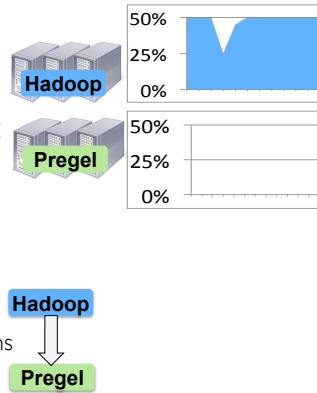
## Computation Model: Frameworks

- A **framework** (e.g., Hadoop, MPI) manages one or more **jobs** in a computer cluster
- A **job** consists of one or more **tasks**
- A **task** (e.g., map, reduce) is implemented by one or more processes running on a single machine



## One Framework Per Cluster Challenges

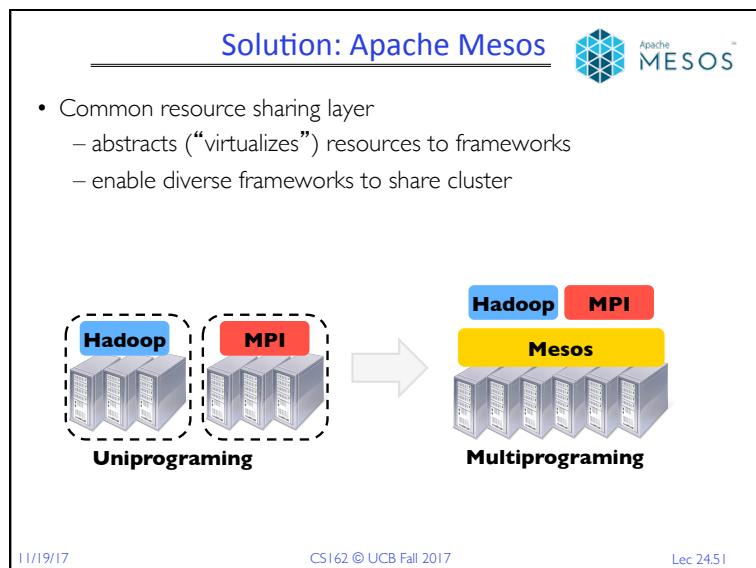
- Inefficient resource usage
  - E.g., Hadoop cannot use available resources from Pregel's cluster
  - No opportunity for stat. multiplexing
- Hard to share data
  - Copy or access remotely, expensive
- Hard to cooperate
  - E.g., Not easy for Pregel to use graphs generated by Hadoop



## Solution: Apache Mesos



- Common resource sharing layer
  - abstracts (“virtualizes”) resources to frameworks
  - enable diverse frameworks to share cluster



## Fine Grained Resource Sharing

- Task granularity both in **time & space**
  - Multiplex node/time between tasks belonging to different jobs/frameworks
- Tasks typically short; median  $\sim= 10$  sec, minutes
- Why fine grained?
  - Improve data locality
  - Easier to handle node failures

11/19/17 CS162 © UCB Fall 2017 Lec 24.52

## Mesos Goals

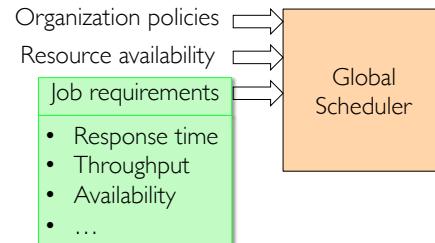
- High utilization of resources
- Support diverse frameworks (existing & future)
- Scalability to 10,000's of nodes
- Reliability in face of node failures
- Focus of this talk: *resource management & scheduling*

11/19/17

CS162 © UCB Fall 2017

Lec 24.53

## Approach: Global Scheduler

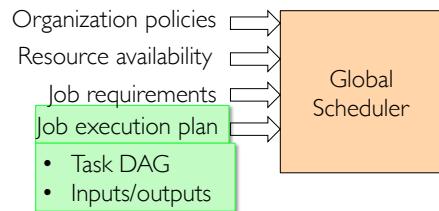


11/19/17

CS162 © UCB Fall 2017

Lec 24.54

## Approach: Global Scheduler

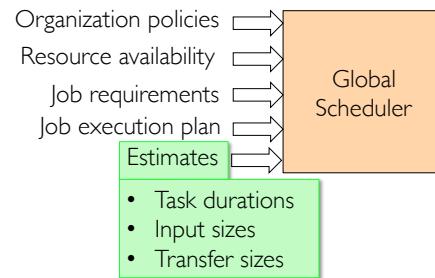


11/19/17

CS162 © UCB Fall 2017

Lec 24.55

## Approach: Global Scheduler

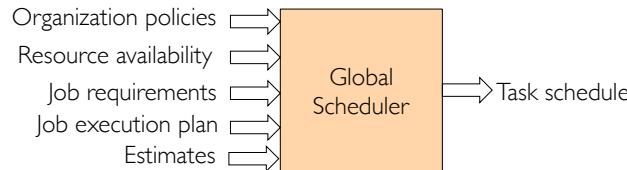


11/19/17

CS162 © UCB Fall 2017

Lec 24.56

## Approach: Global Scheduler



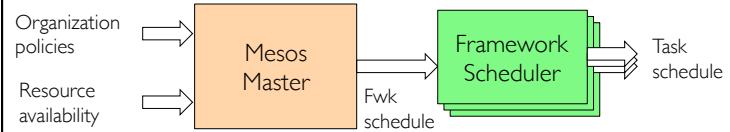
- Advantages: can achieve optimal schedule
- Disadvantages:
  - Complexity → hard to scale and ensure resilience
  - Hard to anticipate future frameworks' requirements
  - Need to refactor existing frameworks

11/19/17

CS162 © UCB Fall 2017

Lec 24.57

## Our Approach: Distributed Scheduler



- Advantages:
  - Simple → easier to scale and make resilient
  - Easy to port existing frameworks, support new ones
- Disadvantages:
  - Distributed scheduling decision → not optimal

11/19/17

CS162 © UCB Fall 2017

Lec 24.58

## Resource Offers

- Unit of allocation: **resource offer**
  - Vector of available resources on a node
  - E.g., node1: <1CPU, 1GB>, node2: <4CPU, 16GB>
- Master sends resource offers to frameworks
- Frameworks select which offers to accept and which tasks to run

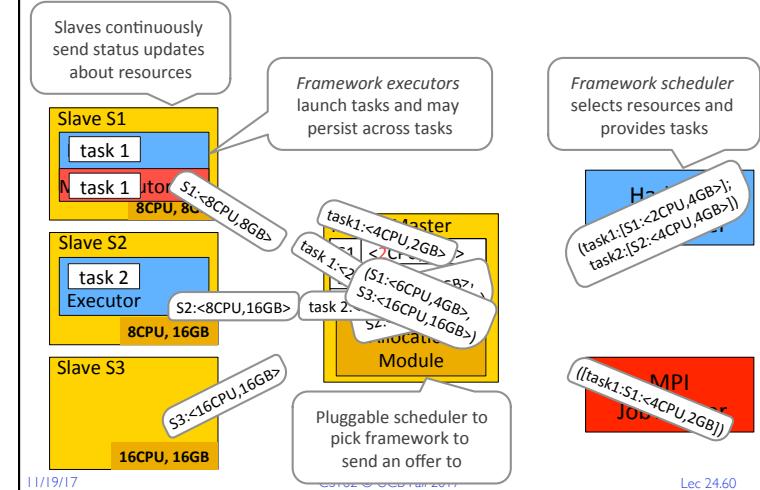
**Push task scheduling to frameworks**

11/19/17

CS162 © UCB Fall 2017

Lec 24.59

## Mesos Architecture: Example



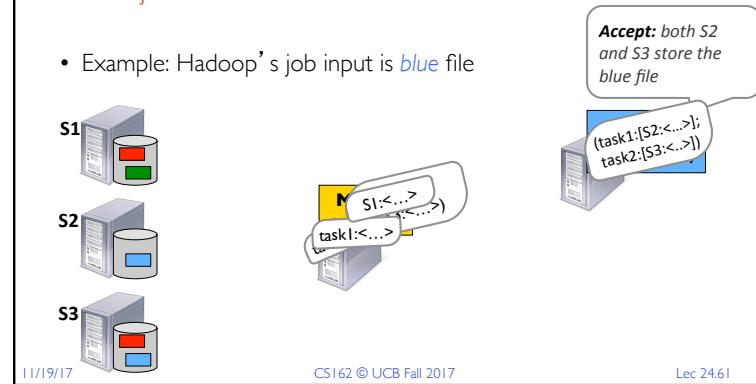
11/19/17

Lec 24.60

## Why does it Work?

- A framework can just wait for an offer that matches its constraints or preferences!
  - Reject offers it does not like

- Example: Hadoop's job input is *blue* file



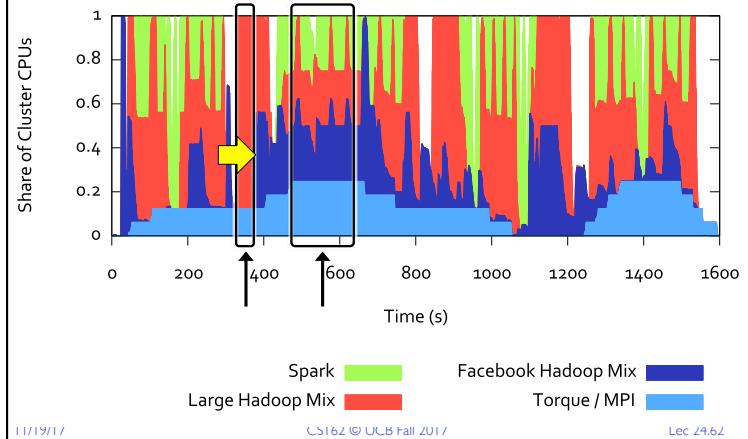
11/19/17

CS162 © UCB Fall 2017

Lec 24.61

## Dynamic Resource Sharing

- 100 node cluster



## Apache Mesos Today

- Hundreds of contributors
- Hundreds of deployments in production
  - E.g., Twitter, GE, Apple
  - Managing 10K node datacenters!
- Mesosphere, startup to commercialize Apache Spark



11/19/17

CS162 © UCB Fall 2017

Lec 24.63

## Summary (1/2)

- Quorum consensus:
  - N replicas
  - Write to W replicas, read from R, where  $W + R > N$
  - Tolerate one failure
- Chord:
  - Highly scalable distributed lookup protocol
  - Each node needs to know about  $O(\log(M))$ , where  $M$  is the total number of nodes
  - Guarantees that a tuple is found in  $O(\log(M))$  steps
  - Highly resilient: works with high probability even if half of nodes fail

11/19/17

CS162 © UCB Fall 2017

Lec 24.64

## Summary (2/2)

- Mesos: large scale resource management systems:
  - Two-level schedulers
  - First level: performance isolation across multiple frameworks, models
  - Second level: frameworks decide which tasks to schedule, when, and where

11/19/17

CS162 © UCB Fall 2017

Lec 24.65