

# Mesos

sharing the cluster.

Benjamin Hindman, **Andy Konwinski**, Matei Zaharia,  
Ali Ghodsi, Anthony Joseph, Randy Katz, Scott Shenker, Ion Stoica

University of California, Berkeley

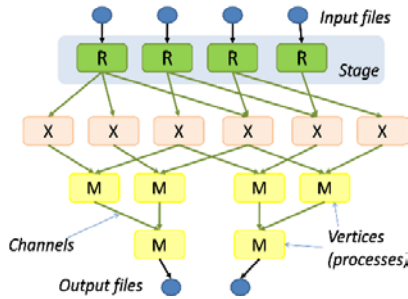




Google™  
Pregel



Pig



Dryad



CIEL

**S4** distributed stream  
computing platform

Google™  
Percolator



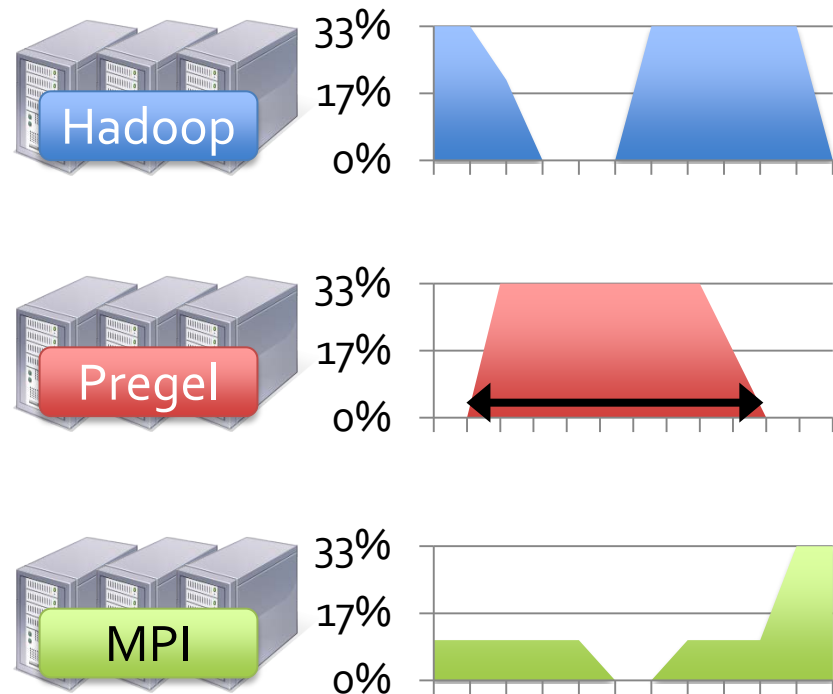
**observation.** Many distributed apps; no single one optimal for all use cases.

Want to run multiple applications in a single cluster

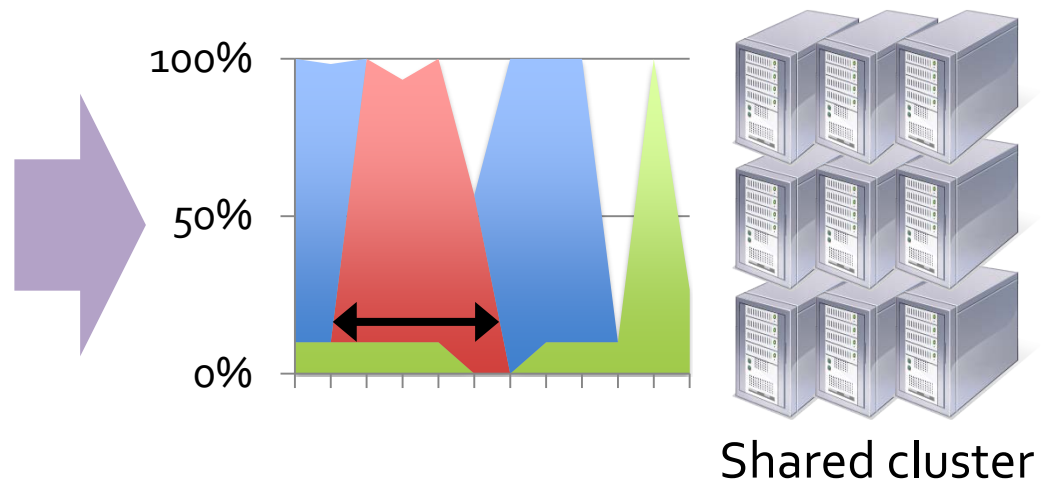
...to *maximize utilization*

...to *share data*

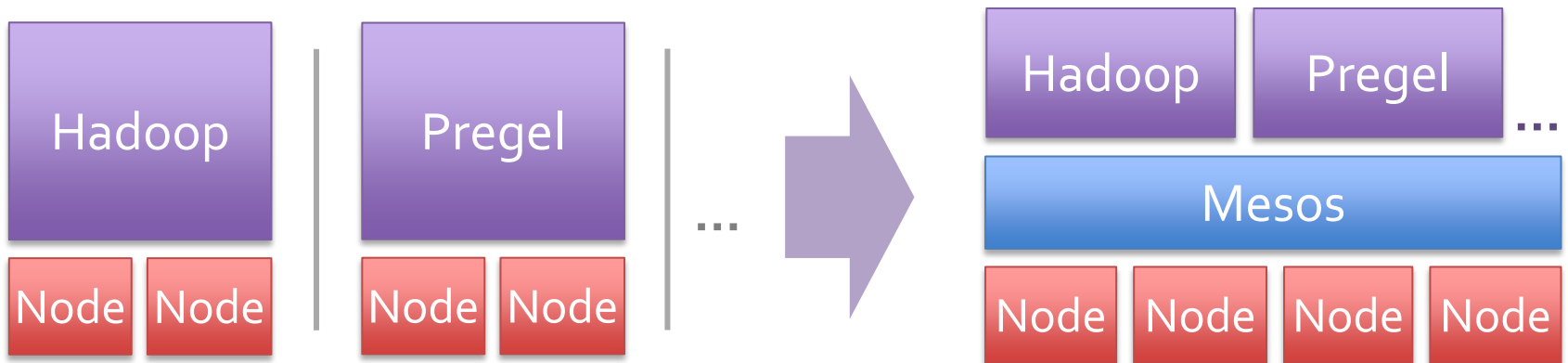
## static partitioning



## dynamic sharing



**solution.** Mesos, a common layer over which diverse applications can run.



Multiple instances of the *same* application

Build *specialized applications* targeting particular problem domains

Build *specialized applications*



# goals.

**High utilization of resources**

**Diverse applications**

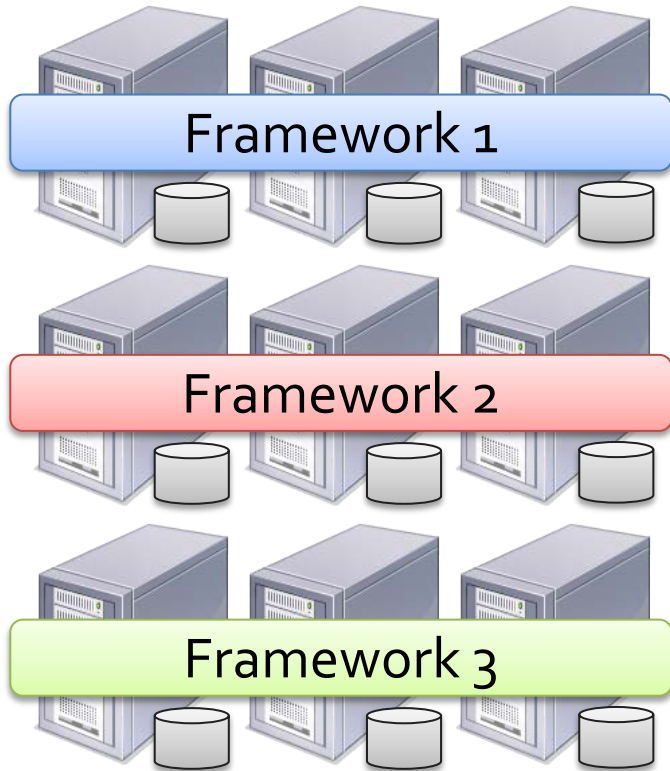
**Scalability to 50-100K nodes**

**Fault tolerance**

**resulting design.** Small microkernel-like core

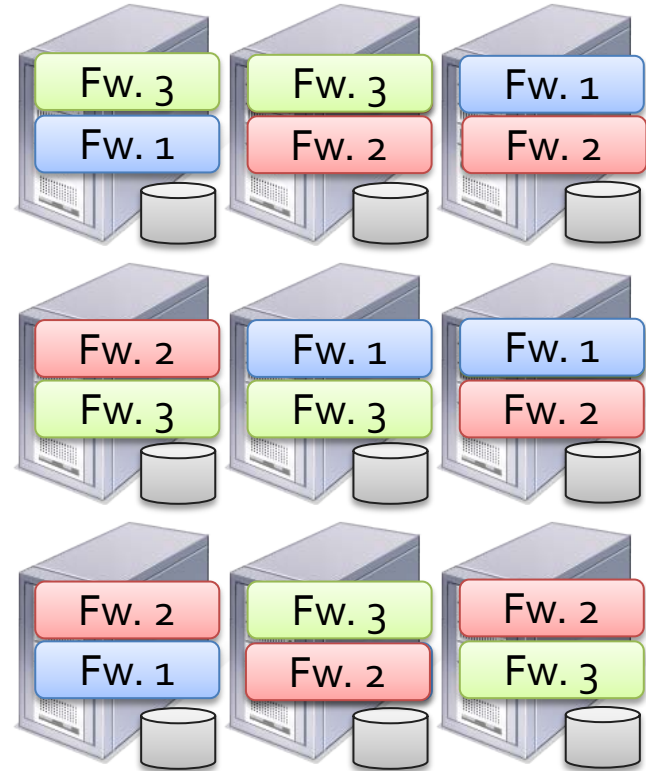
**design element 1. Fine-grained sharing**

## Coarse-Grained Sharing (HPC):



Storage System (e.g. HDFS)

## Fine-Grained Sharing (Mesos):



Storage System (e.g. HDFS)

+ Improved utilization, responsiveness, data locality

**design element 2.** Resource offers (vs. global scheduler)

## global scheduler.

Frameworks express needs in a specification language, global scheduler matches them to resources

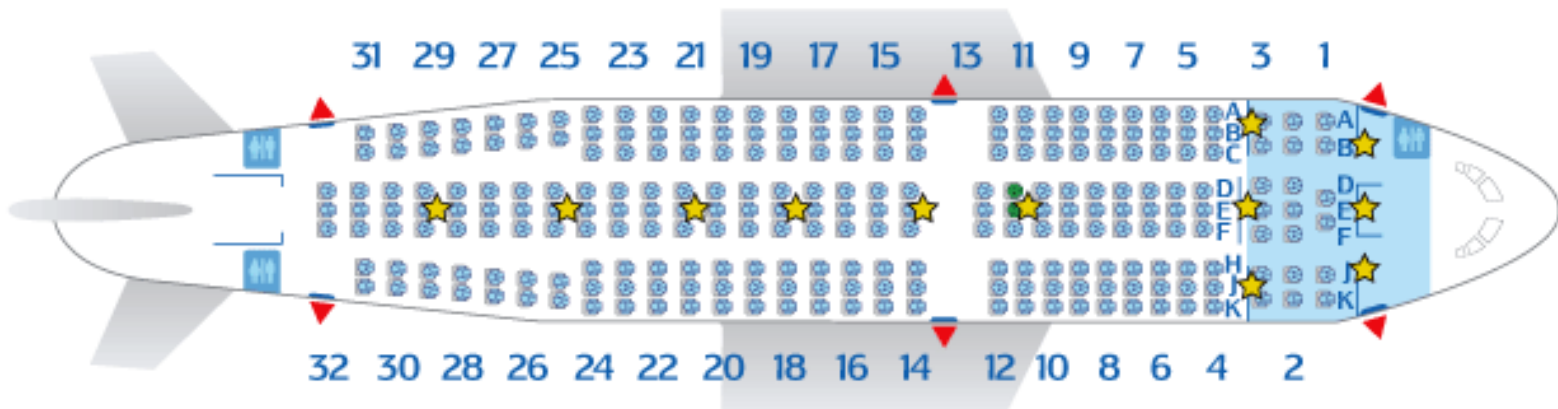
+ Can make optimal decisions

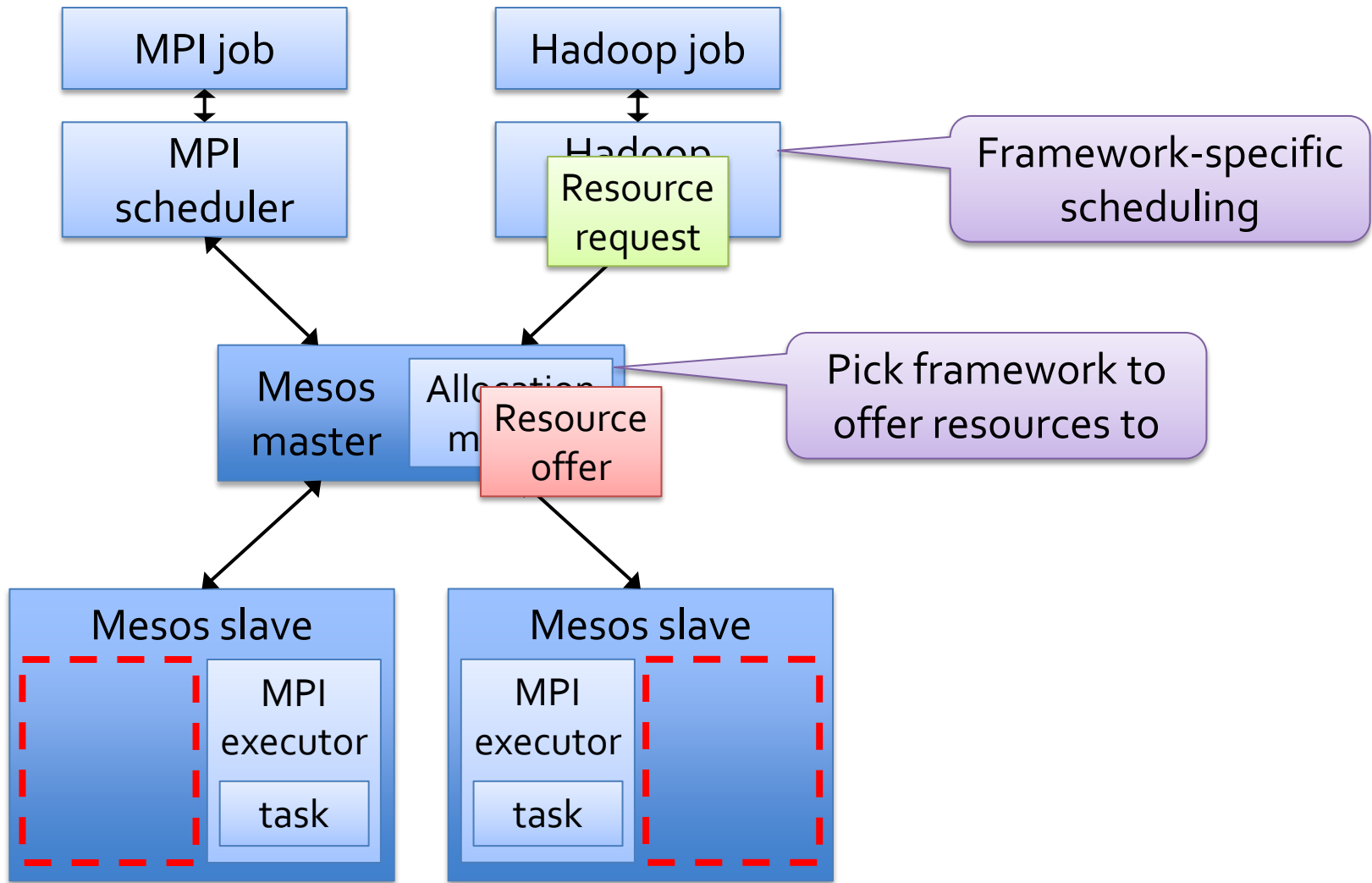
- Complex: language must support all framework needs
- Difficult to scale and to make robust
- Future frameworks may have unanticipated needs

## mesos: resource offers.

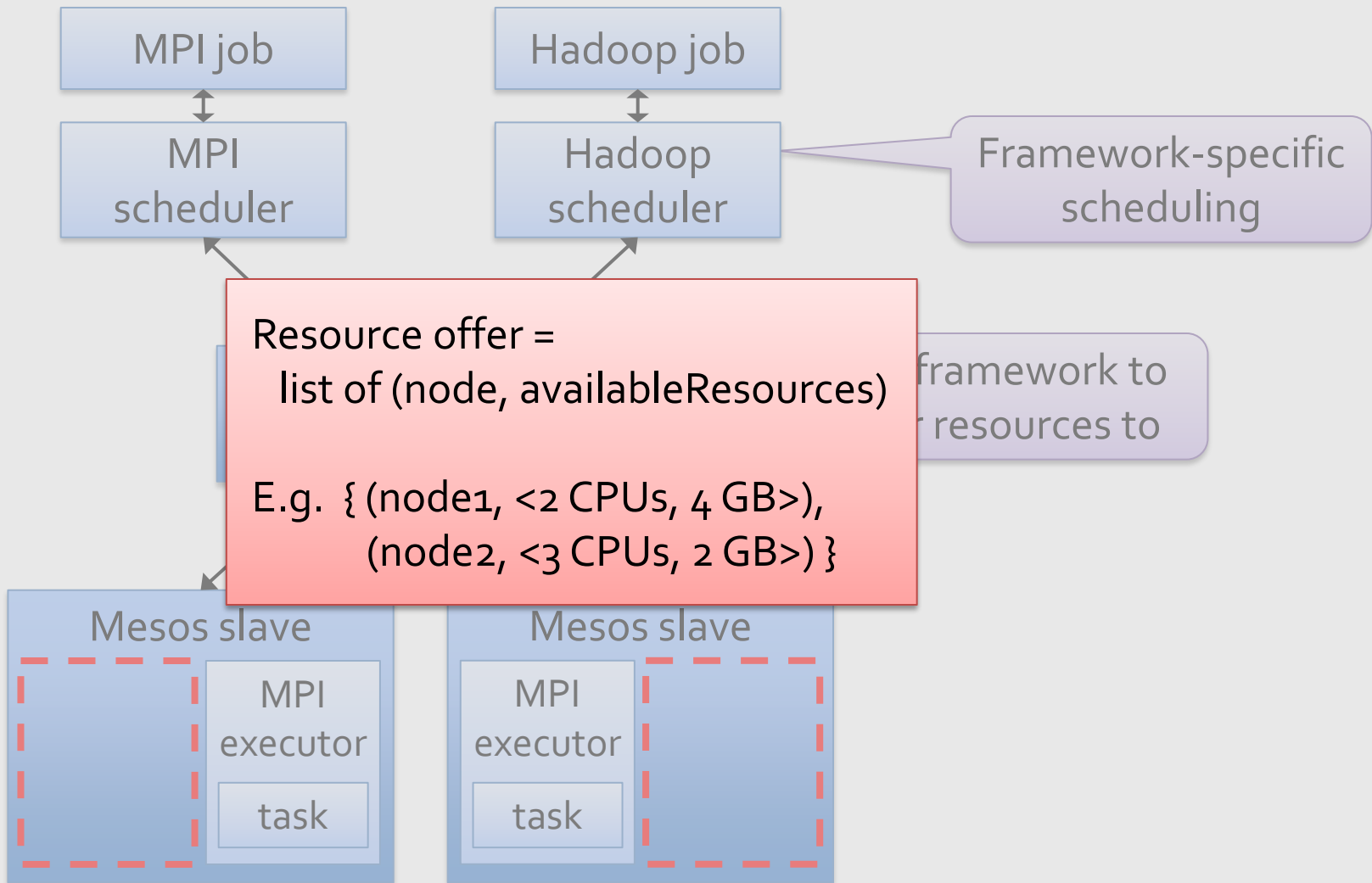
Offer available resources to frameworks, let them pick which resources to use and which tasks to launch

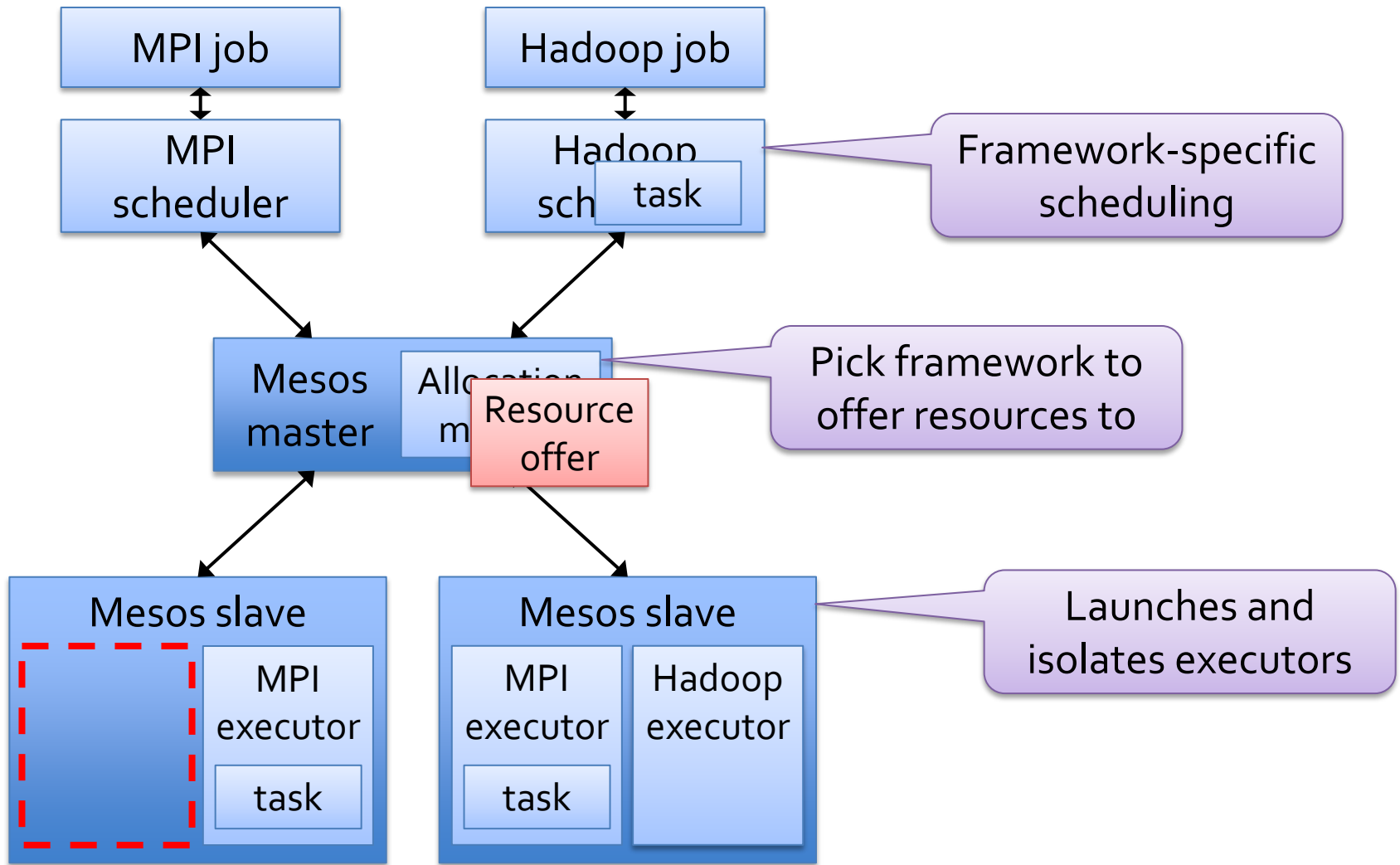
- + Keeps Mesos simple, lets it support future frameworks
- Decentralized decisions might not be optimal













200+ nodes running ~12 production services (stream processing)



Genomics researchers using Hadoop and Spark on Mesos



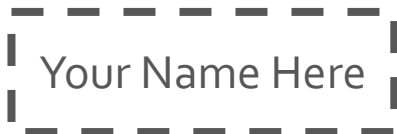
Spark in use by Y! research



Spark for Analytics

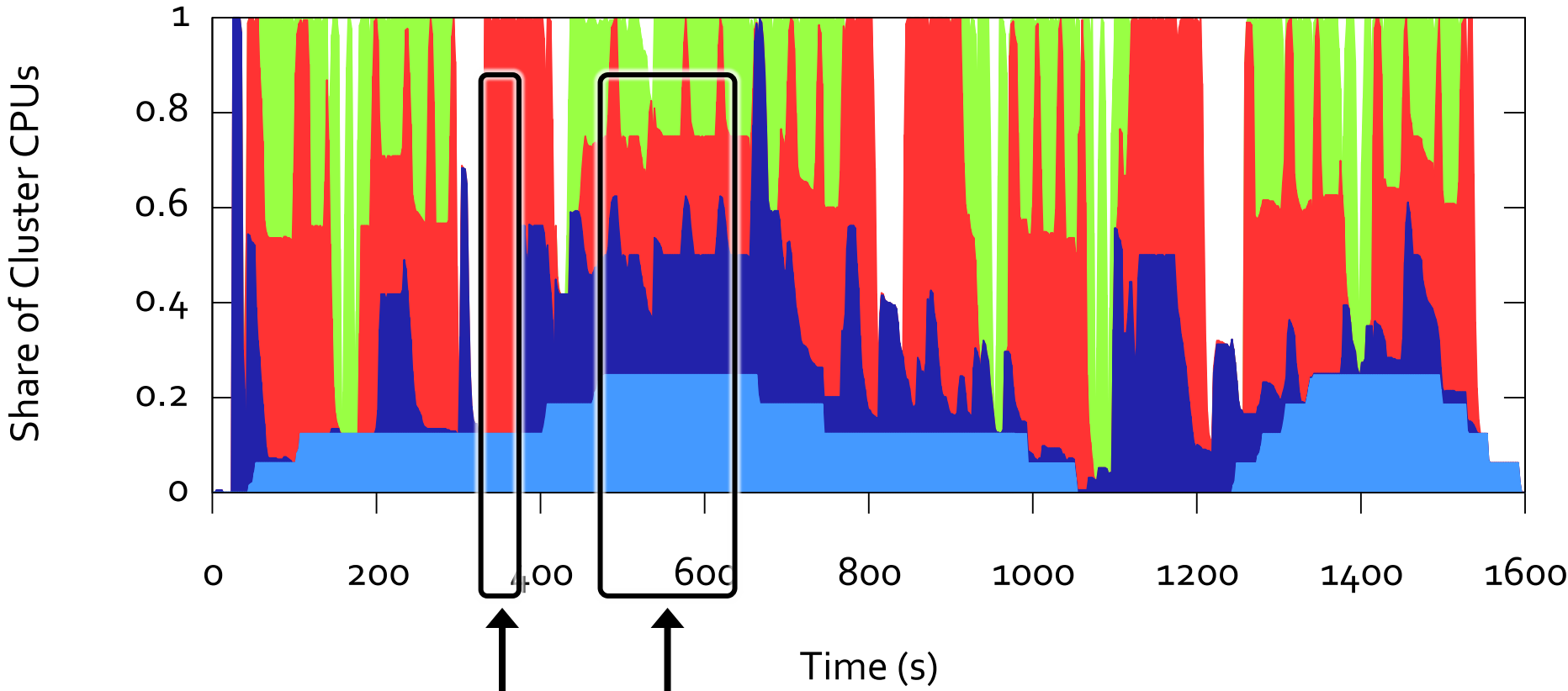


Hadoop, running Experiments, Spark for ML Grad students for research



*Cutting Edge Research!*

**evidence that sharing helps.**



Spark	<span style="display: inline-block; width: 20px; height: 15px; background-color: #90EE90; border: 1px solid black;"></span>	Facebook Hadoop Mix	<span style="display: inline-block; width: 20px; height: 15px; background-color: #00008B; border: 1px solid black;"></span>
Large Hadoop Mix	<span style="display: inline-block; width: 20px; height: 15px; background-color: #FF0000; border: 1px solid black;"></span>	Torque / MPI	<span style="display: inline-block; width: 20px; height: 15px; background-color: #6495ED; border: 1px solid black;"></span>

vs. static allocation (each app gets 25% of nodes)

Framework	Speedup on Mesos
Facebook Hadoop Mix	1.14x
Large Hadoop Mix	2.10x
Spark	1.26x
Torque / MPI	0.96x

# let's talk!

(actually we have been)

about...

1. Lessons learned: ours, Tashi's, others
2. Long running services
3. Socialized services (HDFS, HBase)
4. Priority scheduling

[mesos-project.org](http://mesos-project.org)



# Backup Slides

# Deployments



RAD Lab

# What's running?

Stream processing (tweets from the “firehose”)

» Detecting spam and malicious users

- Duplicate detection
- Repeat offender detection
- Spam account creation banning
- ... more

» Trend identification

# Conviva

35% of all analytics jobs in Spark on Mesos

Spark jobs running on data in Hive

Porting Hive on top of Spark (Shark)

# UCSF

Goal: run Hadoop and Spark to analyze genomic data for UCSF Medical Center

Scale out from 10 dedicated nodes onto ~1000 node HPC cluster running Sun Grid Engine

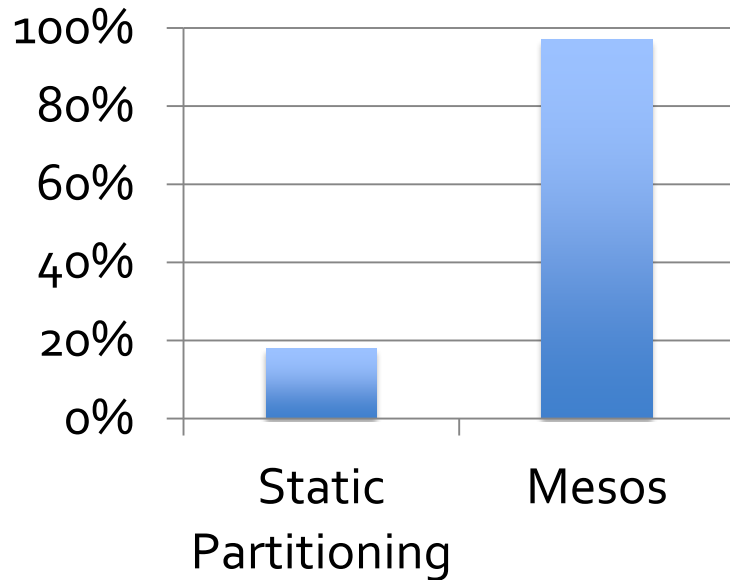
» Nodes allocated through SGE join Mesos cluster

# Data Locality with Resource Offers

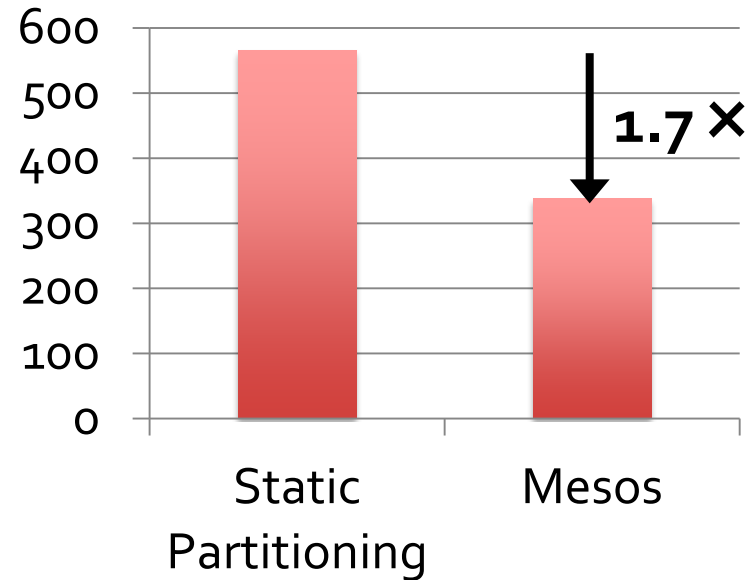
Ran 16 instances of Hadoop on a shared HDFS cluster

Used delay scheduling [EuroSys '10] in Hadoop to get locality (wait a short time to acquire data-local nodes)

## Local Map Tasks (%)



## Job Duration (s)

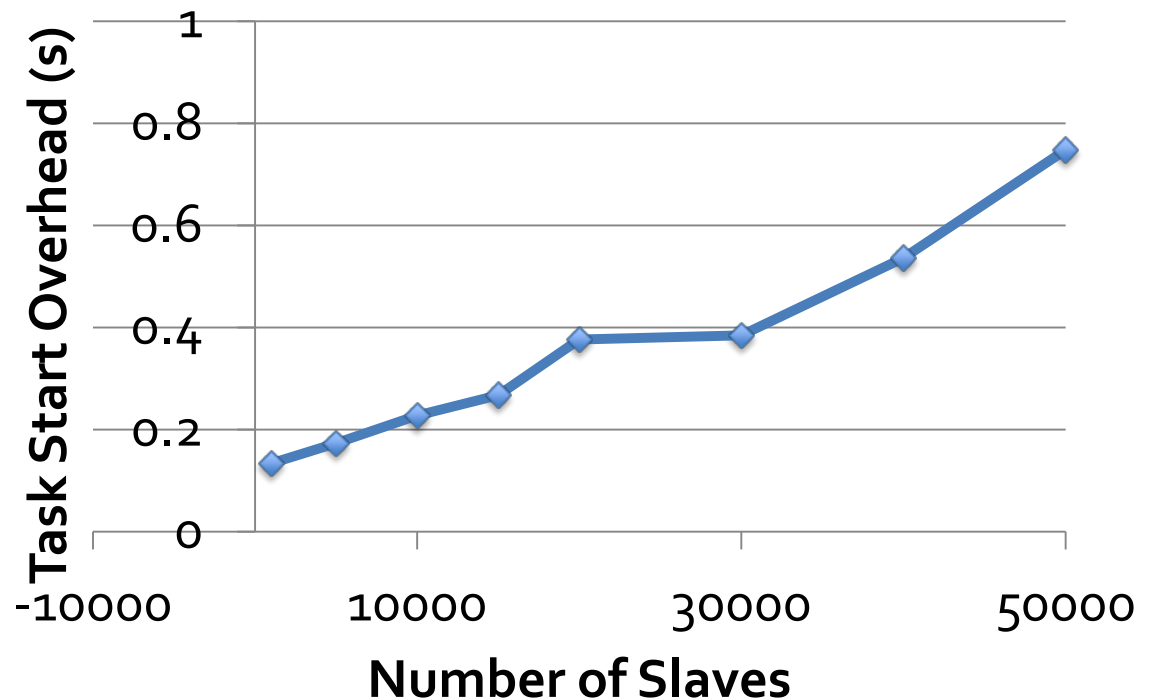


# Scalability

Mesos only performs *inter-framework* scheduling (e.g. fair sharing), which is easier than *intra-framework* scheduling

## Result:

Scaled to 50,000 emulated slaves, 200 frameworks, 100K tasks (30s len)





# Fault Tolerance

Mesos master has only *soft state*: list of currently running frameworks and tasks

Rebuild when frameworks and slaves re-register with new master after a failure

**Result:** fault detection and recovery in ~10 sec

# Related Work

HPC schedulers (e.g. Torque, LSF, Sun Grid Engine)

- » Coarse-grained sharing for inelastic jobs (e.g. MPI)

Virtual machine clouds

- » Coarse-grained sharing similar to HPC

Condor

- » Centralized scheduler based on matchmaking

Parallel work: Next-Generation Hadoop

- » Redesign of Hadoop to have per-application masters
- » Also aims to support non-MapReduce jobs
- » Based on resource request language with locality prefs

# Conclusion

Mesos shares clusters efficiently among diverse frameworks thanks to two design elements:

- » **Fine-grained sharing** at the level of tasks
- » **Resource offers**, a scalable mechanism for application-controlled scheduling

Enables co-existence of current frameworks and development of new specialized ones

In use at Twitter, UC Berkeley, Conviva and UCSF

# Framework Isolation

Mesos uses OS isolation mechanisms, such as Linux containers and Solaris projects

Containers currently support CPU, memory, IO and network bandwidth isolation

Not perfect, but much better than no isolation

# Analysis

Resource offers work well when:

- » Frameworks can scale up and down elastically
- » Task durations are homogeneous
- » Frameworks have many preferred nodes

These conditions hold in current data analytics frameworks (MapReduce, Dryad, ...)

- » Work divided into short tasks to facilitate load balancing and fault recovery
- » Data replicated across multiple nodes

# Revocation

Mesos allocation modules can revoke (kill) tasks to meet organizational SLOs

Framework given a grace period to clean up

“Guaranteed share” API lets frameworks avoid revocation by staying below a certain share

# Mesos API

## Scheduler Callbacks

resourceOffer(offerId, offers)  
offerRescinded(offerId)  
statusUpdate(taskId, status)  
slaveLost(slaveId)

## Scheduler Actions

replyToOffer(offerId, tasks)  
setNeedsOffers(bool)  
setFilters(filters)  
getGuaranteedShare()  
killTask(taskId)

## Executor Callbacks

launchTask(taskDescriptor)  
killTask(taskId)

## Executor Actions

sendStatus(taskId, status)