

Mesos

A platform for Fine-Grained
Resource Sharing in the Data Center

Big Data and Cloud Systems

Paper Presentation

Yang Young Seok

September 3, 2014

Table of Contents

1. Mesos Overview

4. Implementation

2. Mesos Architecture

5. Evaluation

3. Mesos Behavior

6. Comments

1. Mesos Overview

Mesos Website

Apache Software Foundation ▾ / Apache Mesos



[Getting Started](#)

[Documentation](#)

[Downloads](#)

[Community](#)

Learn how to [get started](#), or
[download Mesos 0.20.0](#)

Develop and run resource-efficient distributed systems

Apache Mesos is a cluster manager that simplifies the complexity of running applications on a shared pool of servers.

Mesos Adopters



Chris Fry, SVP of Engineering at Twitter

"Mesos is the cornerstone of our elastic compute infrastructure -- it's how we build all our new services and is critical for Twitter's continued success at scale. It's one of the primary keys to our data center efficiency."



Brenden Matthews, AirBnB

"At Airbnb, we're using Mesos to manage cluster resources for most of our data infrastructure. We run [Chronos](#), [Storm](#), and Hadoop on top of Mesos in order to process petabytes of data."

News

- *August 21, 2014* - Mesos 0.20.0 is released! See the [0.20.0 release notes](#)
- *July 24, 2014* - [#MesosCon Hackathon Details Announced](#)
- *July 17th, 2014* - Mesos 0.19.1 is released! See the [0.19.1 release notes](#) and [blog post announcement](#) for more details.
- *June 24th, 2014* - [Results from the 2014 Mesos community survey](#) have been published on the Mesos blog
- *June 11th, 2014* - [#MesosCon program has been announced, and registration is open!](#)
- *June 9th, 2014* - Mesos 0.19.0 is released! See the

Mesos Master Dashboard

Mesos

Frameworks

Slaves

Offers

Master 20140421-191144-2688574400-5050-21887

Cluster: (Unnamed)

Server: 192.99.64.160:5050

Built: an hour ago by a...

Started: 4 minutes ago

LOG

Slaves

| | |
|-------------|---|
| Activated | 3 |
| Deactivated | 1 |

Tasks

| | |
|----------|---|
| Staged | 0 |
| Started | 0 |
| Finished | 0 |
| Killed | 0 |
| Failed | 0 |
| Lost | 0 |

Resources

| | | |
|---------|-----|---------|
| | CPU | Mem |
| Total | 18 | 44.0 GB |
| Used | 0 | 0 B |
| Offered | 0 | 0 B |
| Idle | 18 | 44.0 GB |

Active Tasks

Find...

| ID | Name | State | Started | Host |
|------------------|------|-------|---------|------|
| No active tasks. | | | | |

Completed Tasks

Find...

| ID | Name | State | Started | Host |
|---------------------|------|-------|---------|------|
| No completed tasks. | | | | |

Mesos Master & Slave
= daemon processes

Tasks to execute in the cluster

Resources in the cluster

*screenshot from google images

(1) Fine-grained resource sharing

Goals:

- Dynamic allocation
- High data locality

Framework A

Task1

Task2

Task3

Task4

Framework B

Task1

Task2

Task3

Task4

Mesos

Cluster

Resource
(memory
/cpu)

Computer

Data

Resource
(memory
/cpu)

Computer

Data

Resource
(memory
/cpu)

Computer

Data

Resource
(memory
/cpu)

Computer

Data

Resource
(memory
/cpu)

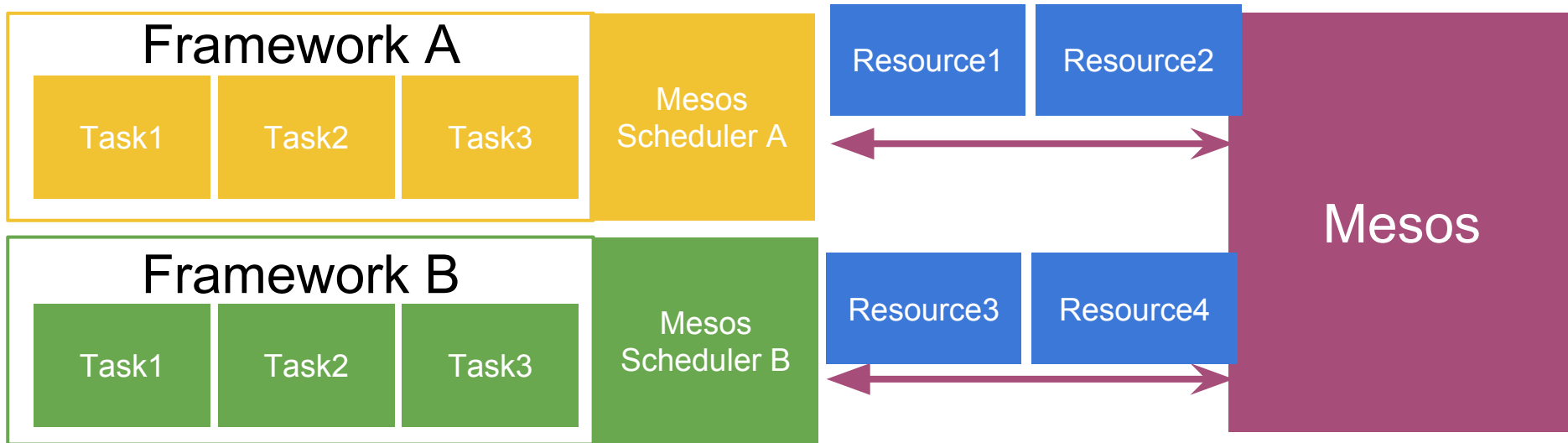
Computer

Data

(2) Distributed Scheduling

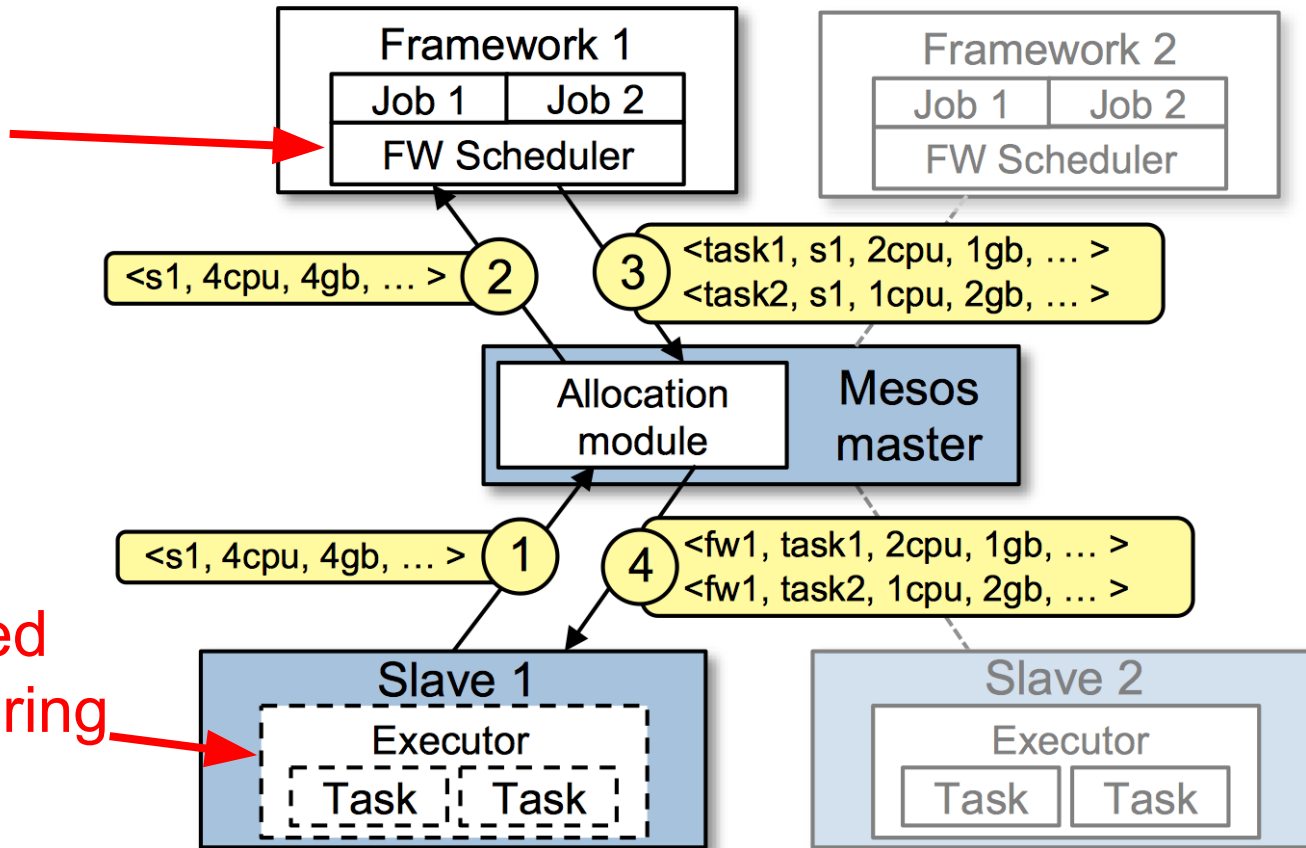
Goal:

- Meet Different scheduling needs from different frameworks



Result: Resource Offer

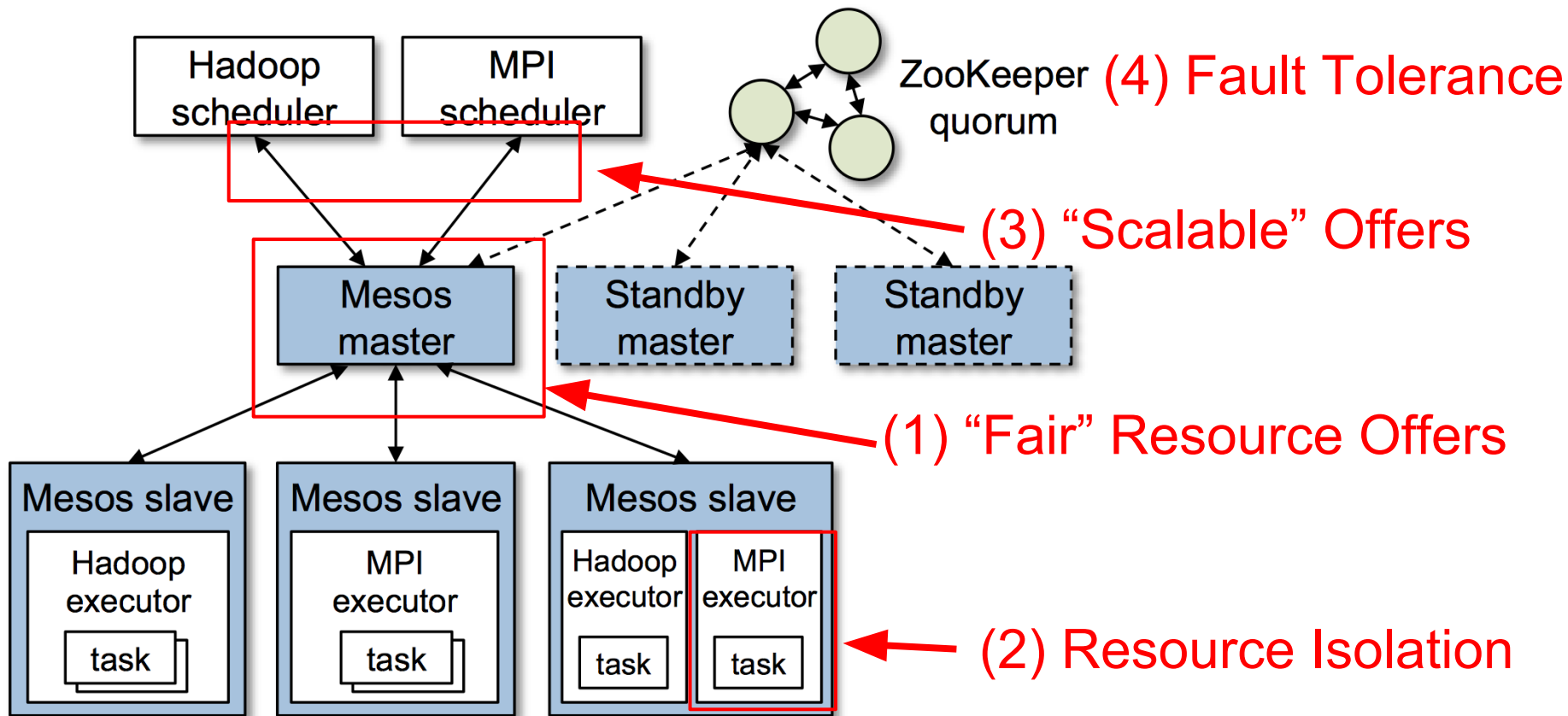
(2) Distributed Scheduling



(1) Fine-Grained Resource Sharing

2. Mesos Architecture

Mesos Architecture



(1) “Fair” Resource Offers

Dominant Resource Fairness algorithm:

Select user with least dominant share

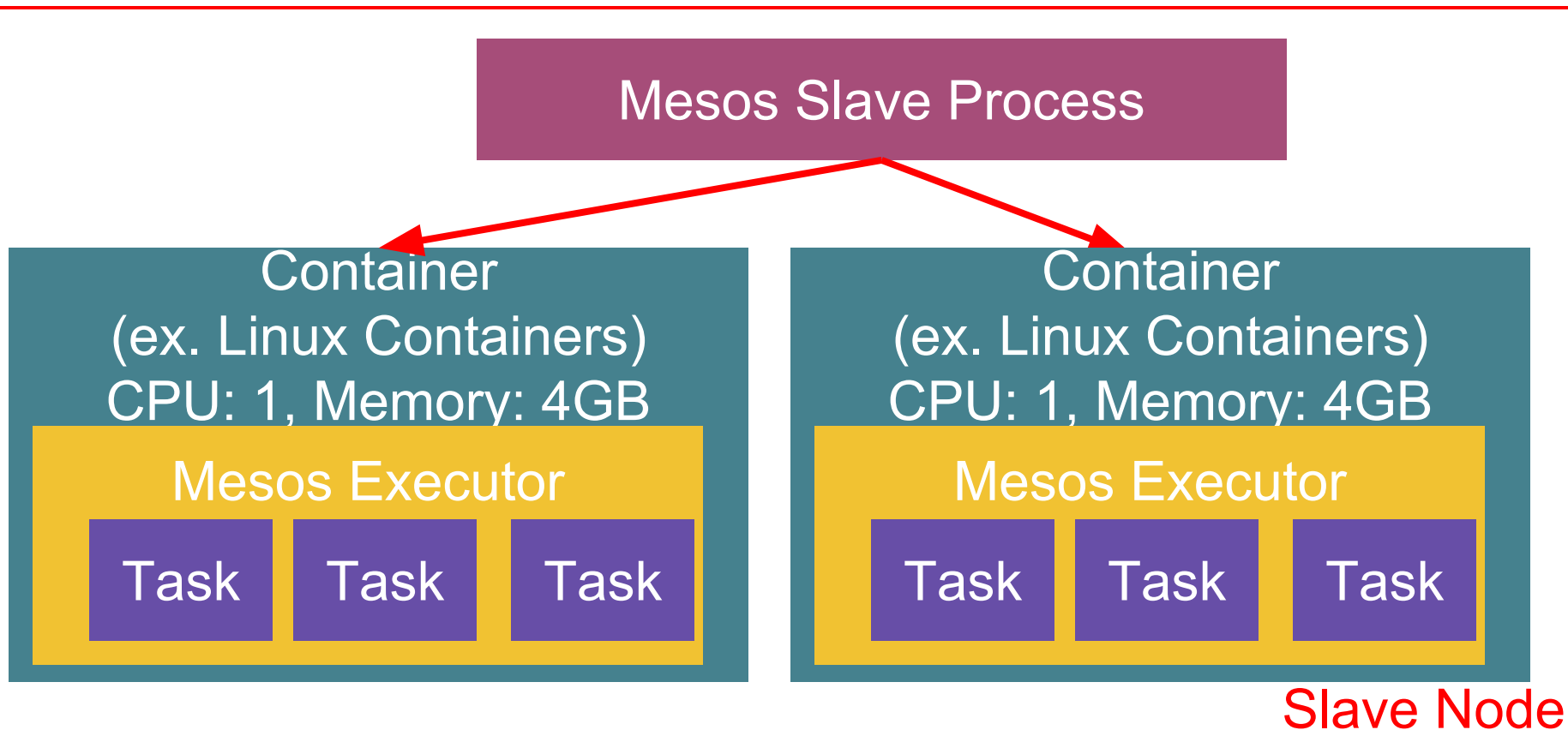
*dominant share = $\max(\text{resource shares})$

| Schedule | User <i>A</i> | | User <i>B</i> | | CPU total alloc. | RAM total alloc. |
|---------------|------------------------------|------------|-----------------------------|------------|---------------------|---------------------|
| | res. shares | dom. share | res. shares | dom. share | | |
| User <i>B</i> | $\langle 0, 0 \rangle$ | 0 | $\langle 3/9, 1/18 \rangle$ | 1/3 | 3/9 | 1/18 |
| User <i>A</i> | $\langle 1/9, 4/18 \rangle$ | 2/9 | $\langle 3/9, 1/18 \rangle$ | 1/3 | 4/9 | 5/18 |
| User <i>A</i> | $\langle 2/9, 8/18 \rangle$ | 4/9 | $\langle 3/9, 1/18 \rangle$ | 1/3 | 5/9 | 9/18 |
| User <i>B</i> | $\langle 2/9, 8/18 \rangle$ | 4/9 | $\langle 6/9, 2/18 \rangle$ | 2/3 | 8/9 | 10/18 |
| User <i>A</i> | $\langle 3/9, 12/18 \rangle$ | 2/3 | $\langle 6/9, 2/18 \rangle$ | 2/3 | 1 | 14/18 |

(1) “Fair” Resource Offers

The definition of "least share" is based on work we did at Berkeley called dominant resource fairness (DRF, see [here](https://www.usenix.org/legacy/event/nsdi11/tech/full_papers/Ghodsi.pdf)https://www.usenix.org/legacy/event/nsdi11/tech/full_papers/Ghodsi.pdf)). In a nutshell, DRF looks at the resources a framework cares the most about (it's dominant resource) and makes sure it has a "fair" share of that resource. Here's the intuition: imagine if we statically partitioned the cluster into $1/N$ sized smaller clusters for N frameworks. Then the amount of work that each framework could run in their cluster will be upper bounded by the resource they use up first, i.e., their dominant resource! Thus, at the very least, we should probably make sure we let a framework in a Mesos cluster be able to execute that much work at any point in time ... and of course the benefit to running in Mesos is the framework can also elastically scale to more than $1/N$ when those resources are idle.

(2) Resource Isolation



(3) Scalable Offers

- Filters at launchTask(ex. time to consider the unused resources refused)
→ avoids over-communication
- Counts resource offered as allocated
→ prompts quick response to the offers
- No response for a long time
→ rescinds offers

(4) Fault Tolerance(FT)

- HA Master → Soft state (only lists of active slaves, frameworks, schedulers) enables quick recovery via ZooKeeper
- Reports Node/Executor faults to Scheduler
→ enables error-handling by Scheduler
- HA Scheduler → A framework can register multiple schedulers (coordination is up to the user)

3. Mesos Behavior

Homogeneous Tasks

| | Elastic Framework | | Rigid Framework | |
|-----------------|-------------------|-------------------|-----------------------|-----------------------------|
| | Constant dist. | Exponential dist. | Constant dist. | Exponential dist. |
| Ramp-up time | T | $T \ln k$ | T | $T \ln k$ |
| Completion time | $(1/2 + \beta)T$ | $(1 + \beta)T$ | $(1 + \beta)T$ | $(\ln k + \beta)T$ |
| Utilization | 1 | 1 | $\beta/(1/2 + \beta)$ | $\beta/(\ln k - 1 + \beta)$ |

*derivations in mesos_tech_report

Assuming infinite demand,
Elastic Framework achieves full utilization

Heterogeneous Tasks

Long tasks hurting short tasks?

- if $\frac{1}{2}$ random tasks and 5 slots per node, the probability for a node filled with long tasks is $(\frac{1}{2})^5 = 0.4\%$
- Also, Mesos can “offer” maximum task duration

Placement Preferences

When Demand for node N exceeds the Supply
p slots preferred by m frameworks(intend total: s)

A central scheduler

- Weighted fair allocation policy: $p \cdot s_i / (\sum_{i=1}^m s_i)$

(Distributed) Mesos schedulers don't know s...

- Mesos Master make offers: $s_i / (\sum_{i=1}^n s_i)$

Framework Incentives

1. Short Tasks:

can use short slots, less loss at faults

2. Scale elastically:

start/complete job earlier, grab opportunities

3. Decline resources that won't be used ASAP:

these will also be counted when offers are made

Limitations

1. Fragmentation

- Short tasks waiting for long tasks
- Long tasks waiting for many short tasks

*Solution: an allocation module enforcing a maximum offer size

2. Interdependent framework constraints

- For example, tasks from 2 frameworks cannot be colocated

*Solution: None (but these cases are rare)

3. Framework complexity due to resource offers

- 1st Argument: decision(distributed) vs. preference(central)
- 2nd Argument: Online-algorithms can be easily implemented

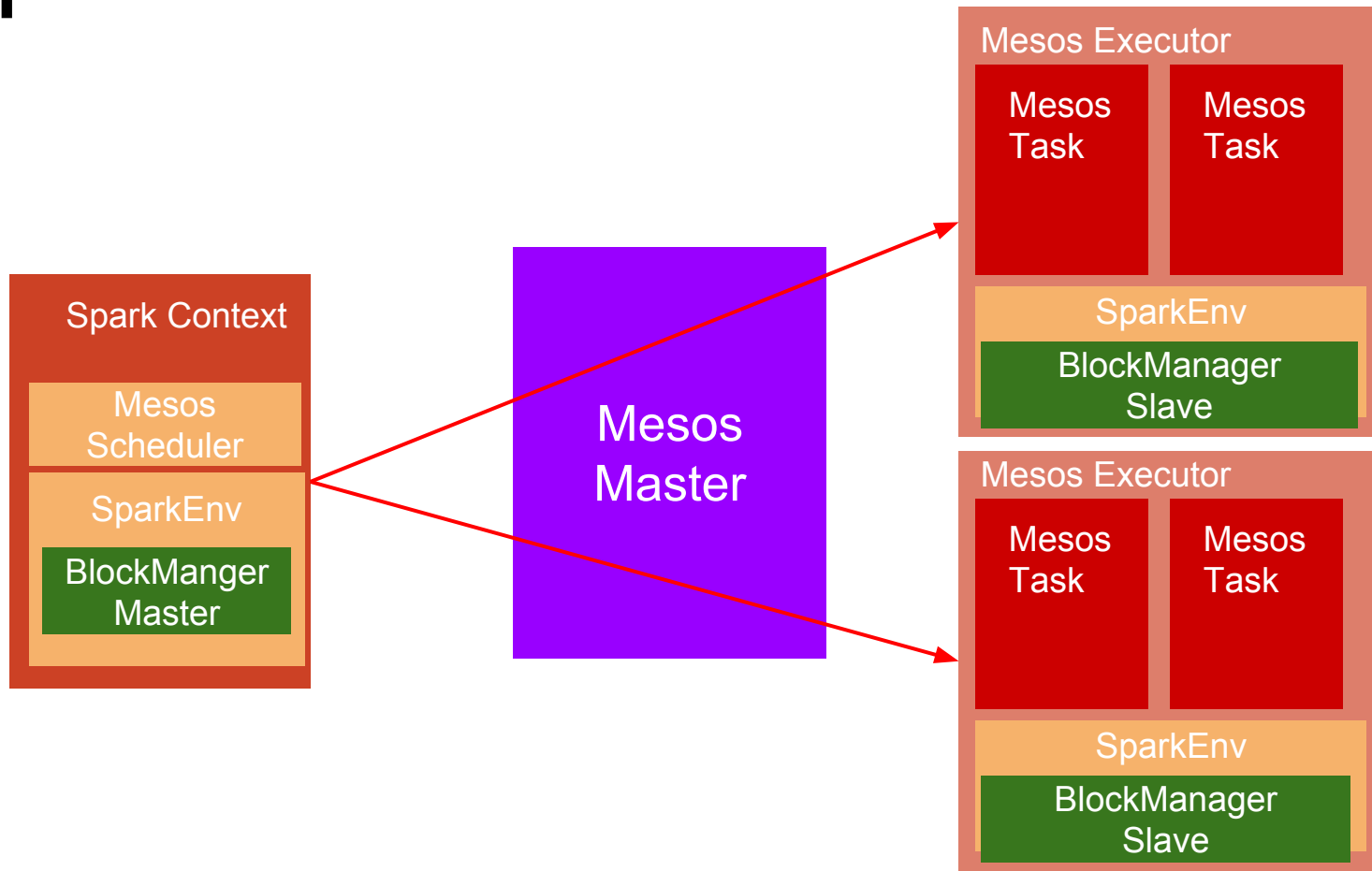
4. Implementation

Mesos Implementation

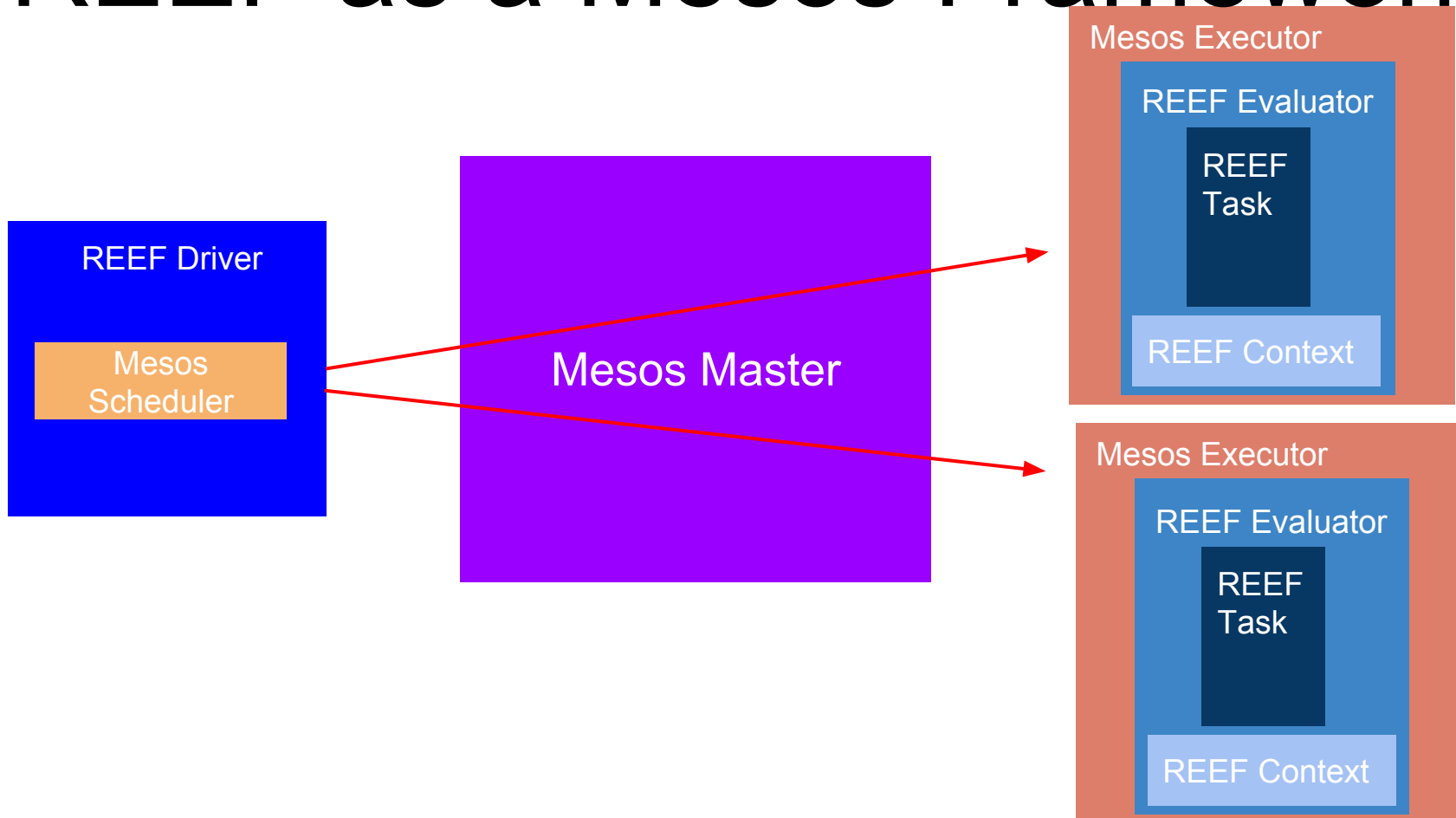
```
[3] authorizer
[4] cli
[5] common
[6] credentials
[7] deploy
[8] examples
[9] exec
[10] files
[11] hdfs
[12] java
[13] jvm
[14] launcher
[15] linux
[16] local
[17] log
[18] logging
[19] master
[20] messages
[21] python
[22] sasl
[23] sched
[24] slave
[25] state
[26] tests
[27] usage
[28] webui
[29] zookeeper
```

- 10,000 lines of C++
- libprocess for actor-based async I/O
- Zookeeper for HA Mesos Master
- Linux Containers/Solaris projects for task Isolation(CPU Cores/Memory)

Spark as a Mesos Framework



REEF as a Mesos Framework



5. Evaluation

Macrobenchmark

macrobenchmark consisting of a mix of four workloads:

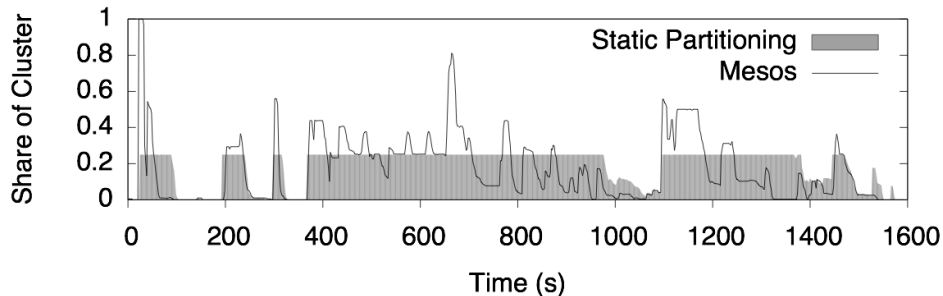
- A Hadoop instance running a mix of small and large jobs based on the workload at Facebook.
- A Hadoop instance running a set of large batch jobs.
- Spark running a series of machine learning jobs.
- Torque running a series of MPI jobs.

We compared a scenario where the workloads ran as four frameworks on a 96-node Mesos cluster using fair sharing to a scenario where they were each given a static partition of the cluster (24 nodes), and measured job response times and resource utilization in both cases. We used EC2 nodes with 4 CPU cores and 15 GB of RAM.

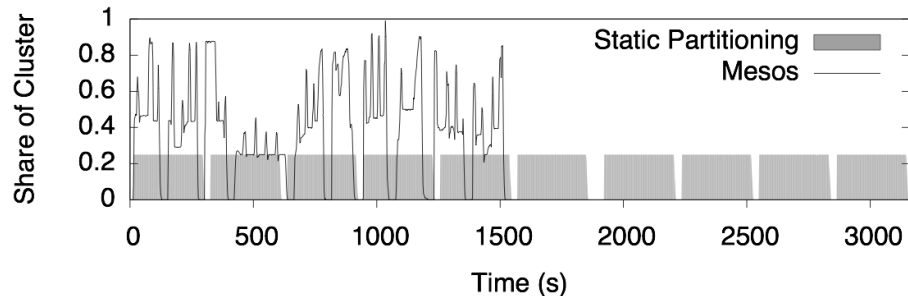
Macrobenchmark

1. higher total utilization 2. faster job execution

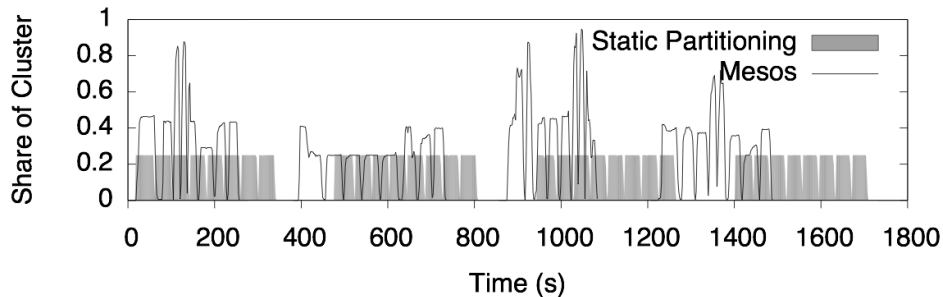
(a) Facebook Hadoop Mix



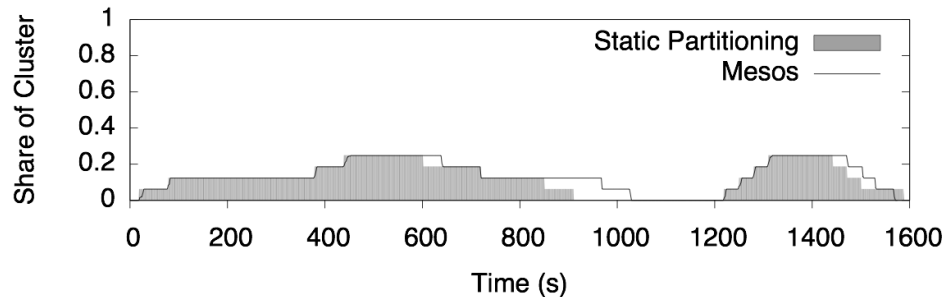
(b) Large Hadoop Mix



(c) Spark



(d) Torque / MPI



Mesos Scalability

- 200 frameworks running continuously
- Launch a test framework(10-sec task) and calculate overhead

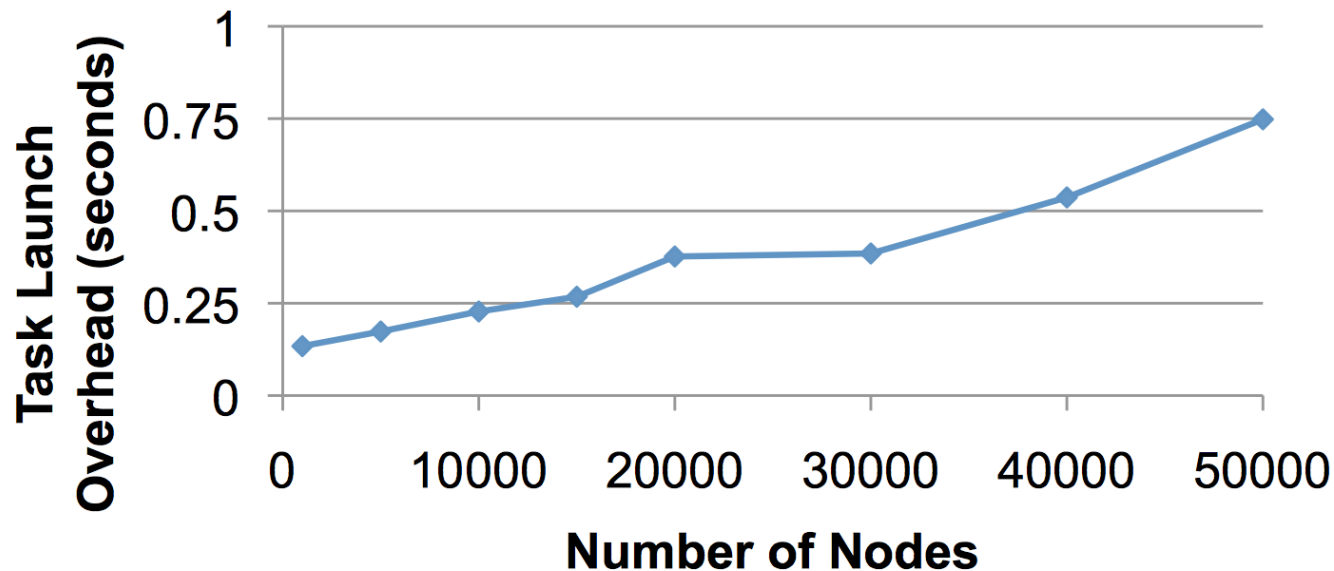


Figure 10: Mesos master's scalability versus number of slaves.

6. Comments

Comments

Pros:

- More Flexible interface than YARN
(YARN ApplicationMaster vs. Mesos Scheduler)
(YARN Container vs. Mesos Executor)
- Solid definition of “Fair Sharing”

Cons:

- No Benchmark against similar systems
(ex. Hadoop 1.0 MapReduce jobs)
- Why does Mesos Scheduler run locally?
(It could run on a slave node like YARN ApplicationMaster)