Philip Norman & Sunil Shah

# SCALING LIKE TWITTER WITH APACHE MESOS

MESOSPHERE

# MODERN INFRASTRUCTURE

**Dan the Datacenter Operator**

- Doesn't sleep very well
- Loves automation
- Wants to control what runs in his datacenter

**Alice the Application Developer**

- Finds setting up infrastructure tedious
- Wants her application to be deployed as quickly as possible

# 3 TENETS

Clean separation of responsibilities

No more 3am wake ups

Easy programmatic deployment

# CLEAN SEPARATION

**Before**

-   Dan cares about his hardware and Alice's software that runs on it

-   Alice cares about her software and what hardware Dan provides

**Now**

-   With Mesos, all the nodes are provisioned exactly the same (but may have heterogenous hardware).

-   Dan doesn't care what software is deployed since applications are well encapsulated.

-   Alice doesn't care where her software is deployed because it's easy enough to scale up and down.

# NO MORE 3AM WAKE UPS

**Before**

- Dan had to react every time an application or machine went down.

**Now**

- Mesos and Marathon monitor running tasks.

- If a task fails or is lost (due to a machine going offline), Mesos communicates that to Marathon.

- Marathon restarts the application.

- Dan gets to sleep peacefully!
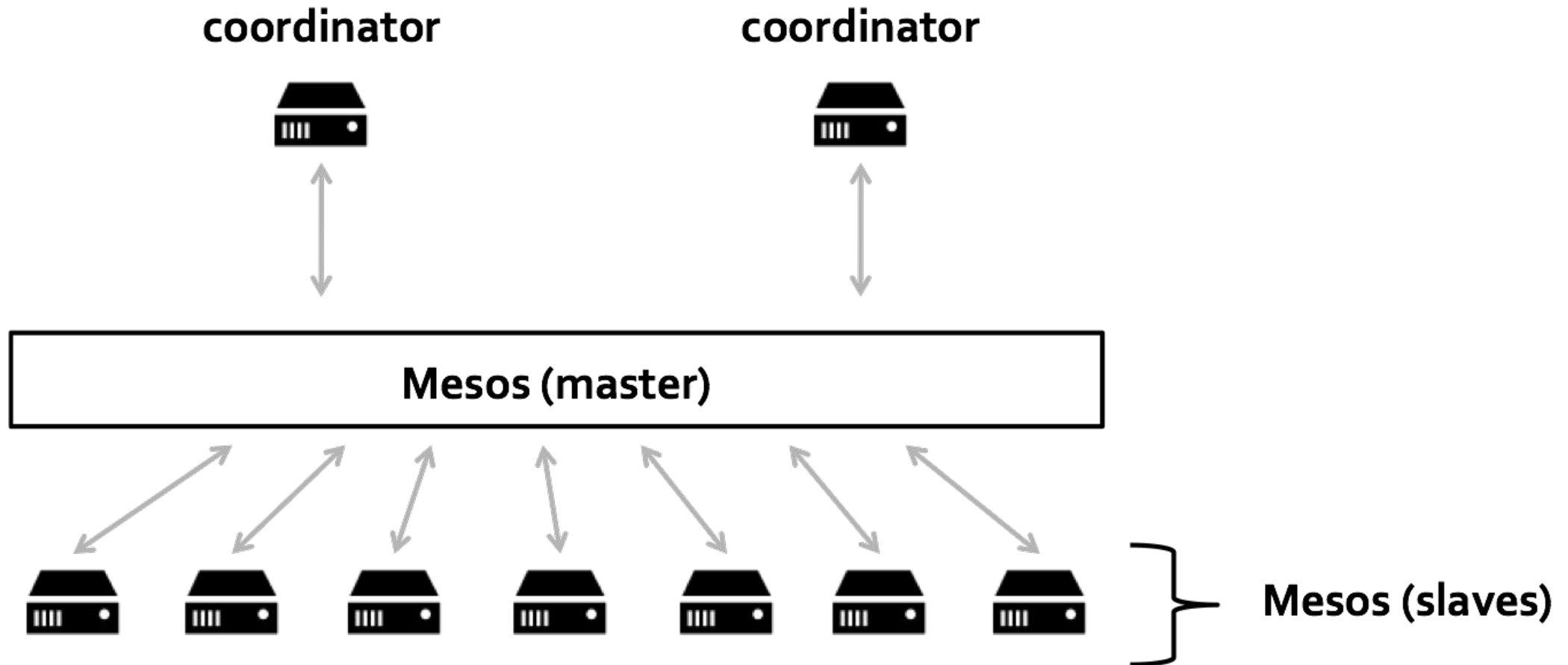
# EASY PROGRAMMATIC DEPLOYMENT

**Before**

- Servers were handcrafted.

- Deploying new or updated software would require oversight and involvement from both Alice and Dan.

**Now**

- Dan provides Alice with her own instance of Marathon that makes it hard for her to take down someone else's application.

- Running applications are isolated from each other by Mesos.

- Marathon offers a nice API that allows Alice to easily deploy new versions safely.

# LAYER OF ABSTRACTION

# INTRODUCTION

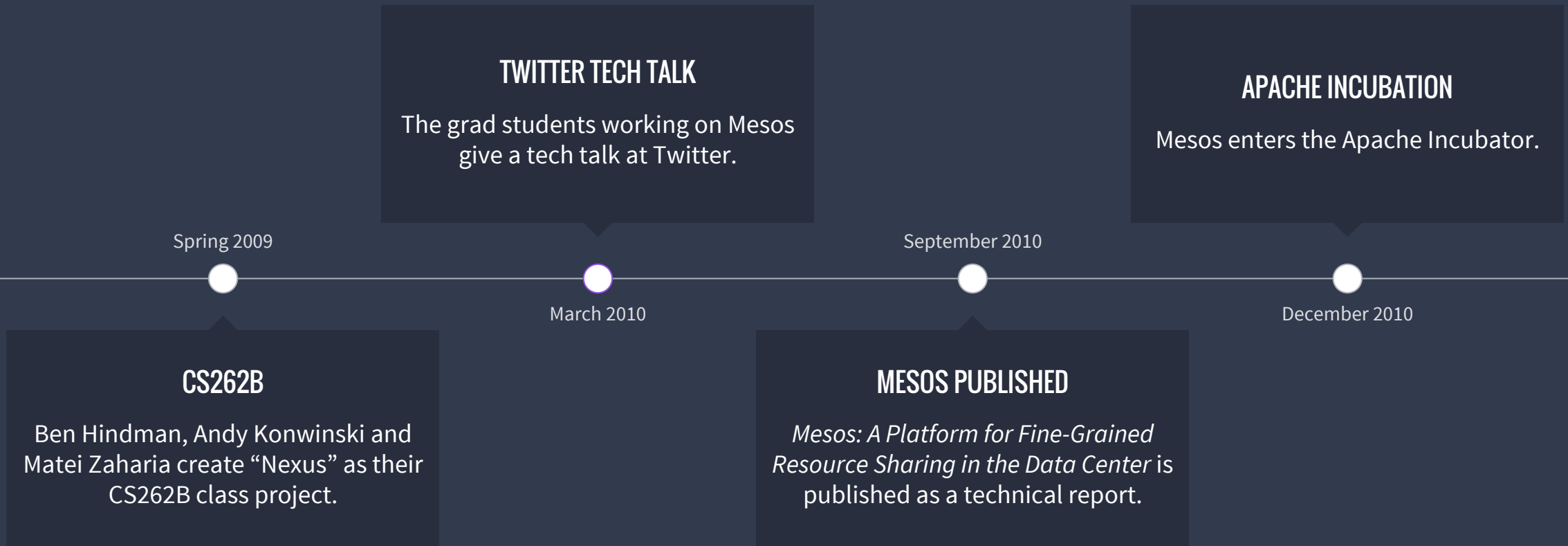Apache Mesos is a **cluster resource manager**.

It handles:

- **Aggregating resources** and **offering them to schedulers**
- **Launching tasks** (i.e. processes) on those resources
- **Communicating the state of those tasks** back to schedulers

# PRODUCTION CUSTOMERS AND MESOS USERS



Bloomberg

Twitter

Apple

ebay

airbnb

HubSpot

verizon

PayPal

Government Agencies

NETFLIX

2σ TWO SIGMA

yelp

# MESOS: ORIGINS

# THE BIRTH OF MESOS

**TWITTER TECH TALK**

The grad students working on Mesos give a tech talk at Twitter.

**APACHE INCUBATION**

Mesos enters the Apache Incubator.

Spring 2009

September 2010

March 2010

December 2010

**CS262B**

Ben Hindman, Andy Konwinski and Matei Zaharia create "Nexus" as their CS262B class project.

**MESOS PUBLISHED**

*Mesos: A Platform for Fine-Grained Resource Sharing in the Data Center* is published as a technical report.

# TECHNOLOGY

# VISION

**Mesos: A Platform for Fine-Grained Resource Sharing in the Data Center**

Benjamin Hindman,    Andy Konwinski,    Matei Zaharia,
Ali Ghodsi, Anthony D. Joseph, Randy Katz, Scott Shenker, Ion Stoica
*University of California, Berkeley*

**The Datacenter Needs an Operating System**

Matei Zaharia,   Benjamin Hindman,   Andy Konwinski,   Ali Ghodsi,
Anthony D. Joseph,   Randy Katz,   Scott Shenker,   Ion Stoica
*University of California, Berkeley*

Sharing resources between batch processing frameworks

- Hadoop
- MPI
- Spark

What does an operating system provide?

- Resource management
- Programming abstractions
- Security
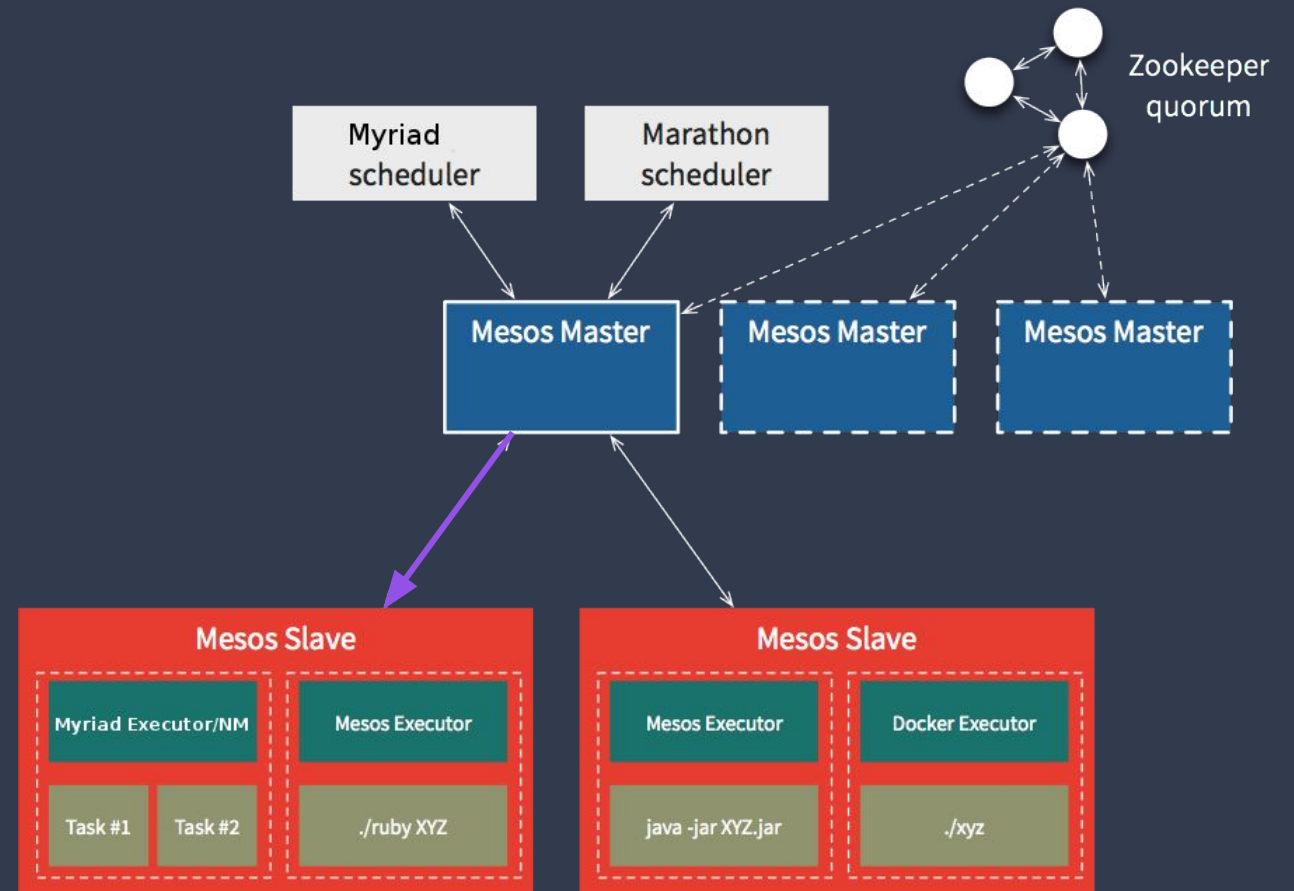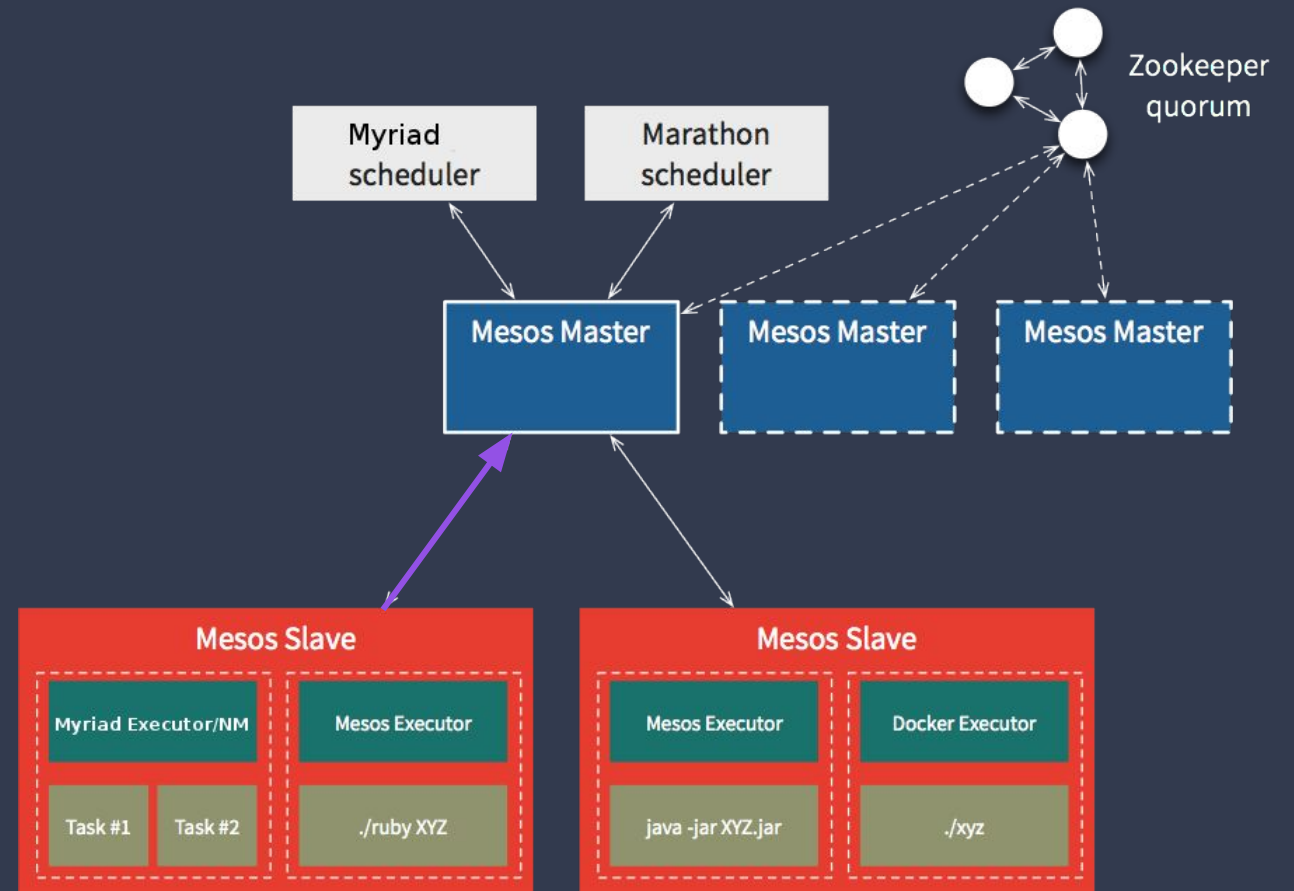- Monitoring, debugging, logging
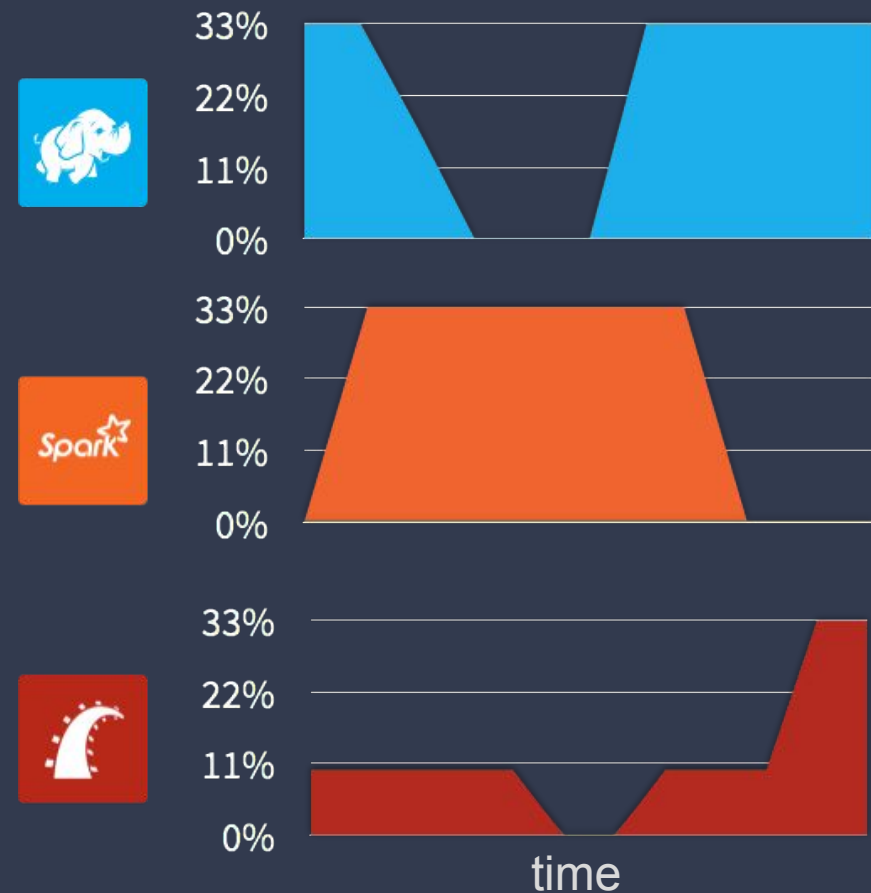
# ARCHITECTURE

# ARCHITECTURE

- Agents advertise resources to Master
- Master offers resources to Scheduler
- Scheduler rejects/uses resources
- Agents report task status to Master

# ARCHITECTURE

- Agents advertise resources to Master
- Master offers resources to Scheduler
- Scheduler rejects/uses resources
- Agents report task status to Master

# ARCHITECTURE
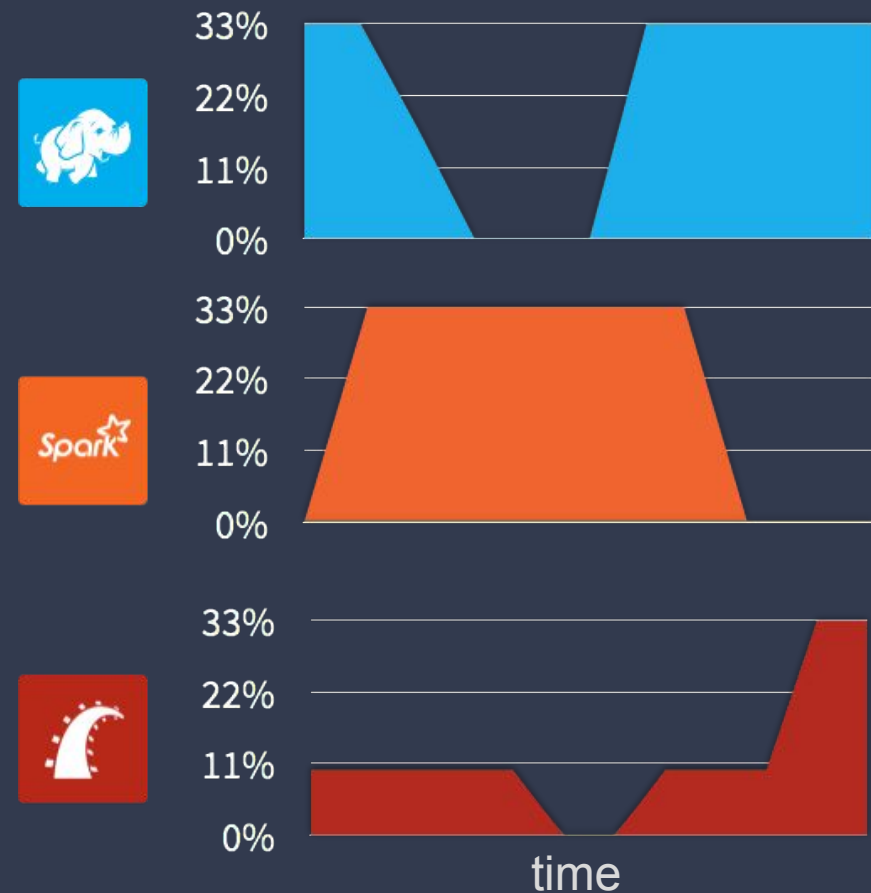
- Agents advertise resources to Master
- Master offers resources to Scheduler
- Scheduler rejects/uses resources
- Agents report task status to Master

# ARCHITECTURE

- Agents advertise resources to Master
- Master offers resources to Scheduler
- Scheduler rejects/uses resources
- Agents report task status to Master

# ARCHITECTURE

- Agents advertise resources to Master
- Master offers resources to Scheduler
- Scheduler rejects/uses resources
- Agents report task status to Master

# KEEP IT STATIC

A naive approach to handling varied app requirements: **static partitioning**.
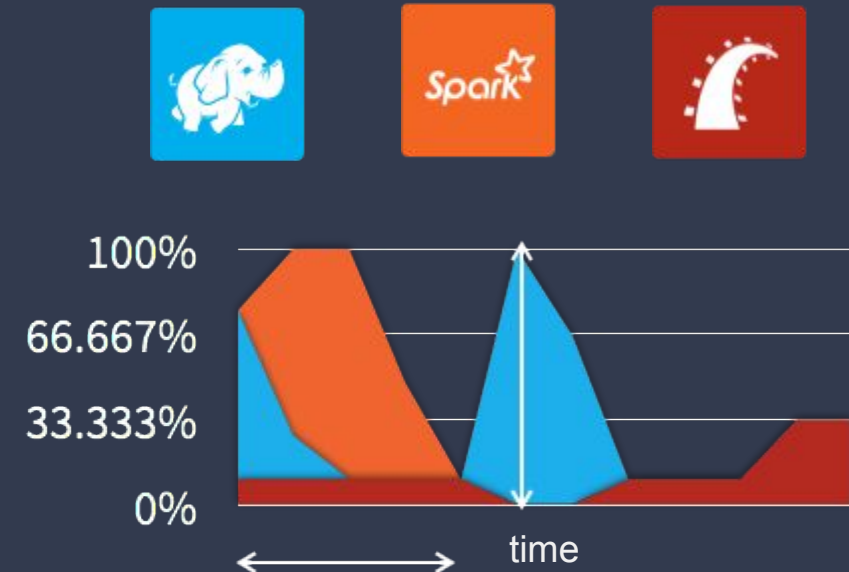
This can cope with heterogeneity, but is very expensive.

# KEEP IT STATIC

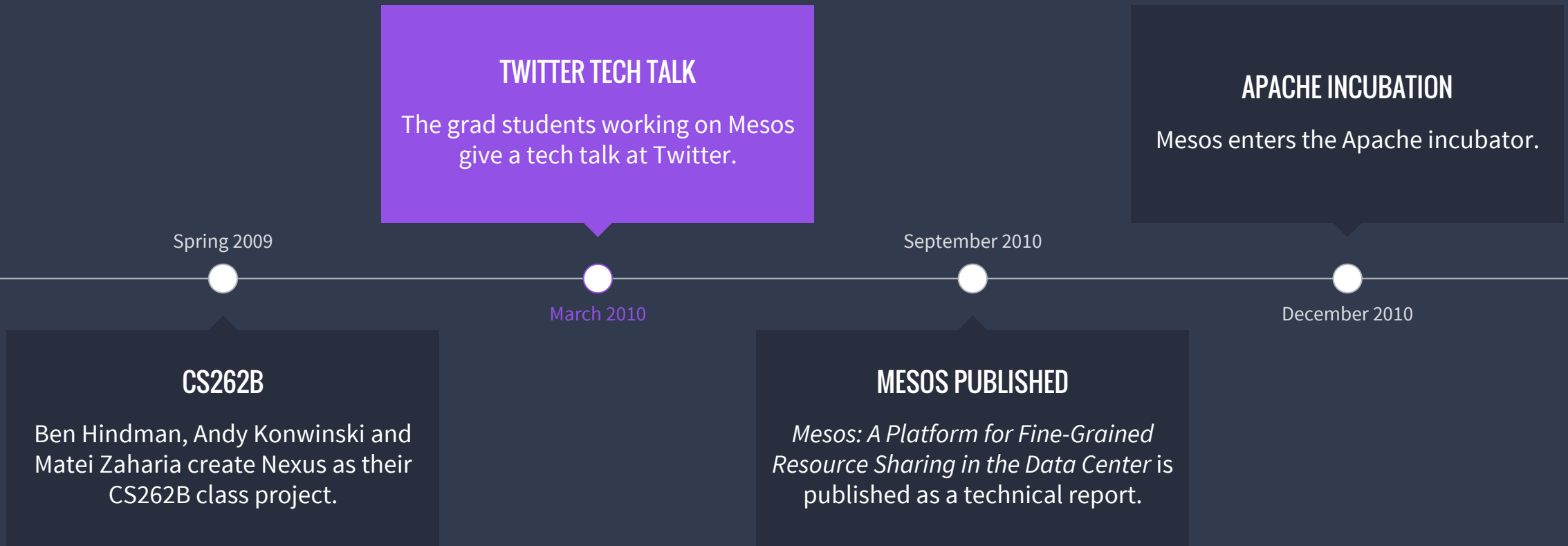Maintaining sufficient headroom to handle peak workloads on all partitions leads to **poor utilisation** overall.

# SHARED RESOURCES

Multiple frameworks can use the same cluster resources, with their share adjusting dynamically.
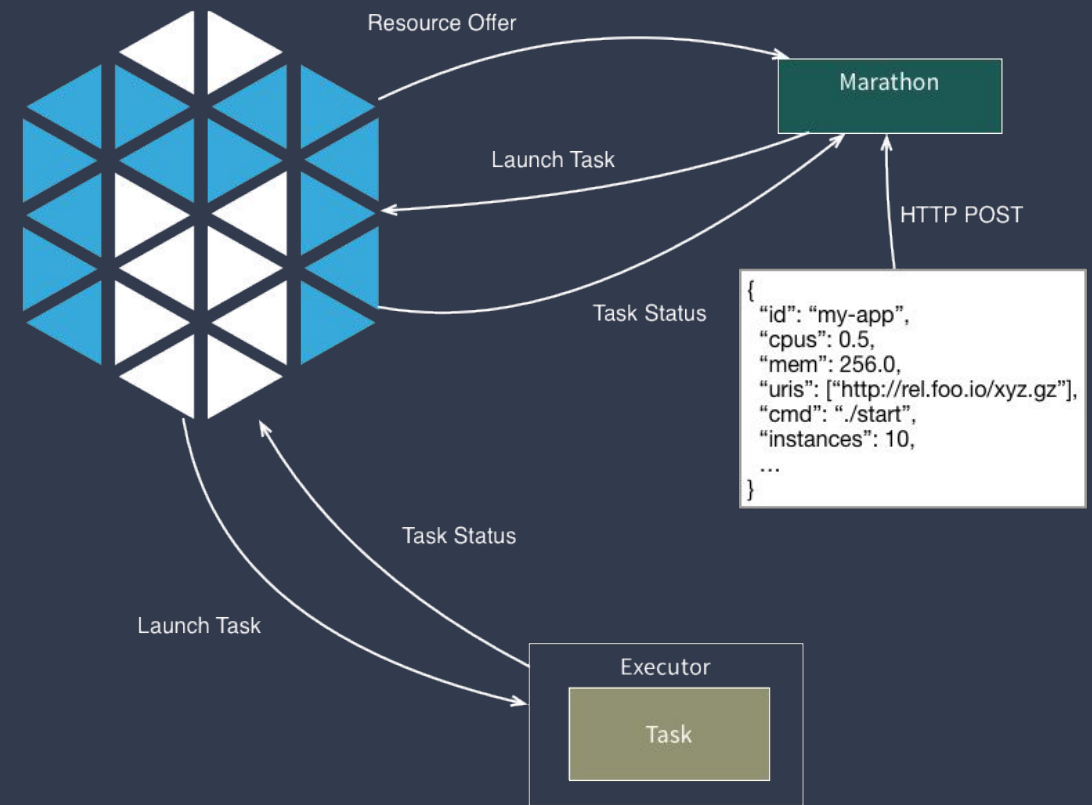
# TWITTER & MESOS

# THE BIRTH OF MESOS

**TWITTER TECH TALK**

The grad students working on Mesos give a tech talk at Twitter.

**APACHE INCUBATION**

Mesos enters the Apache incubator.

Spring 2009

September 2010

March 2010

December 2010

**CS262B**

Ben Hindman, Andy Konwinski and Matei Zaharia create Nexus as their CS262B class project.

**MESOS PUBLISHED**

*Mesos: A Platform for Fine-Grained Resource Sharing in the Data Center* is published as a technical report.

# MESOS REALLY HELPS

- Former Google engineers at Twitter thought Mesos could provide the same functionality as Borg.

- Mesos actually works pretty well for long running services.

Resource Offer

Marathon

Launch Task

HTTP POST

Task Status

{
"id": "my-app",
"cpus": 0.5,
"mem": 256.0,
"uris": ["http://rel.foo.io/xyz.gz"],
"cmd": "./start",
"instances": 10,
…
}

Task Status

Launch Task

Executor

Task

# MESOS WITH MARATHON IN PRODUCTION

# WHAT IS MARATHON?

- Service scheduler for Mesos
- `init.d` for long-running apps
- Your own private PaaS

M A R A T H O N

# WHAT IS MARATHON?

# USEFUL MARATHON FEATURES

- Start, stop, scale, update apps
- Highly available, no SPoF
- Native Docker support
- Powerful Web UI
- Fully featured REST API
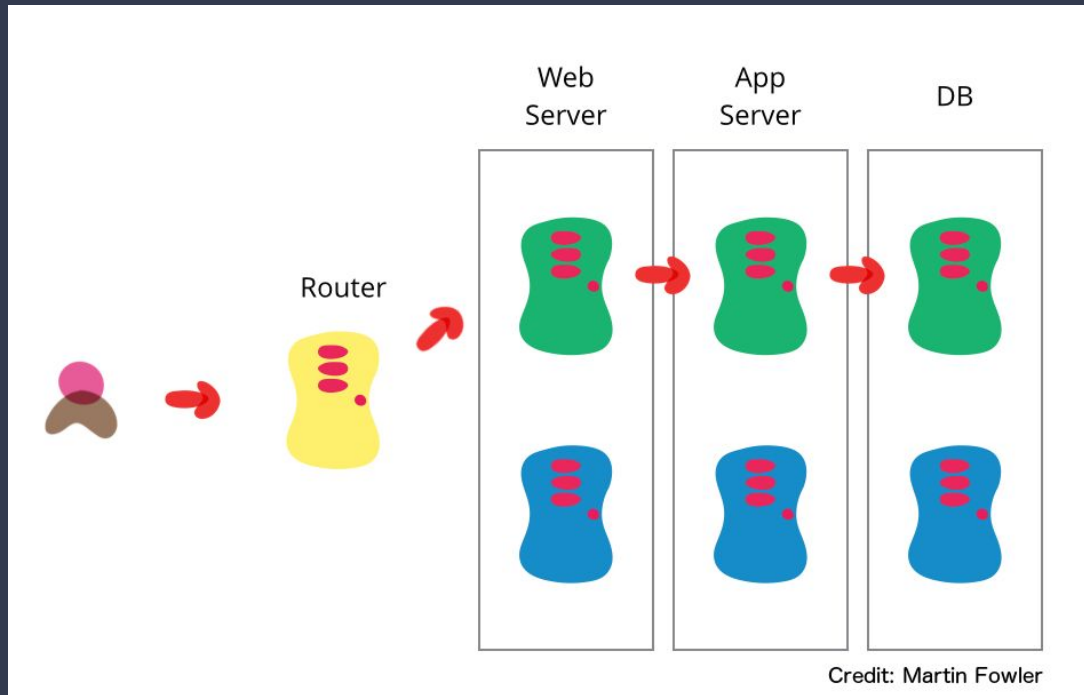- Pluggable event bus
- Artifact staging

MARATHON

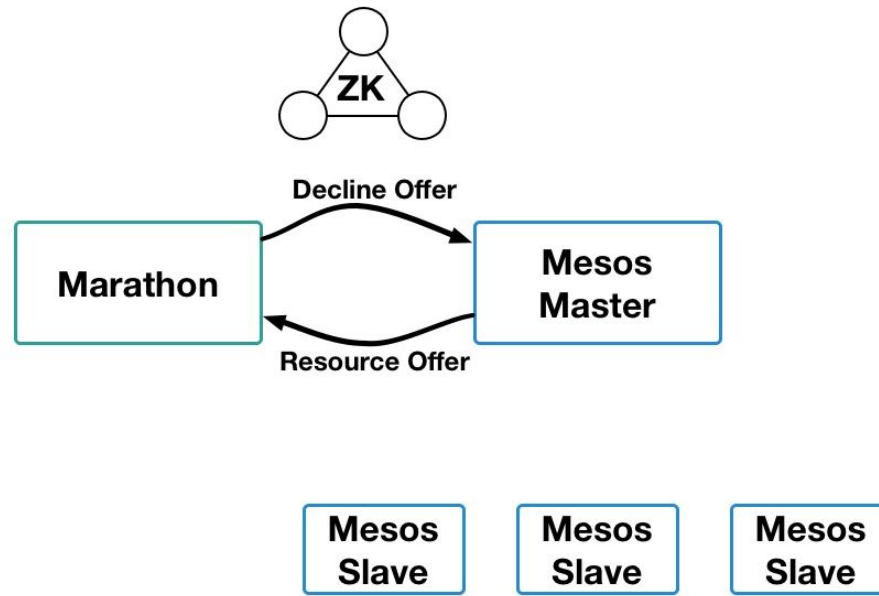# USEFUL MARATHON FEATURES: DEPLOY LIKE FACEBOOK



- Application versioning
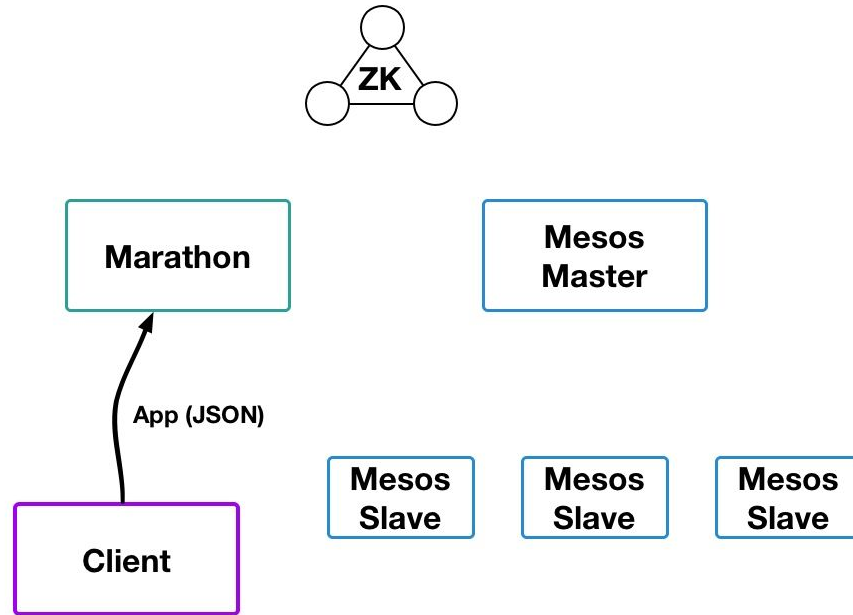- Rolling deploy / restart
- Deployment strategies

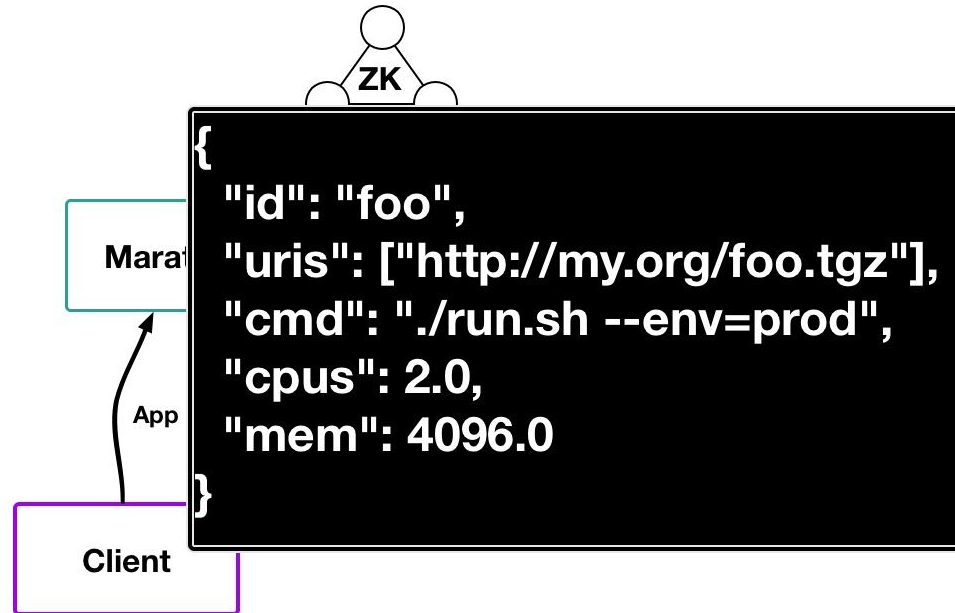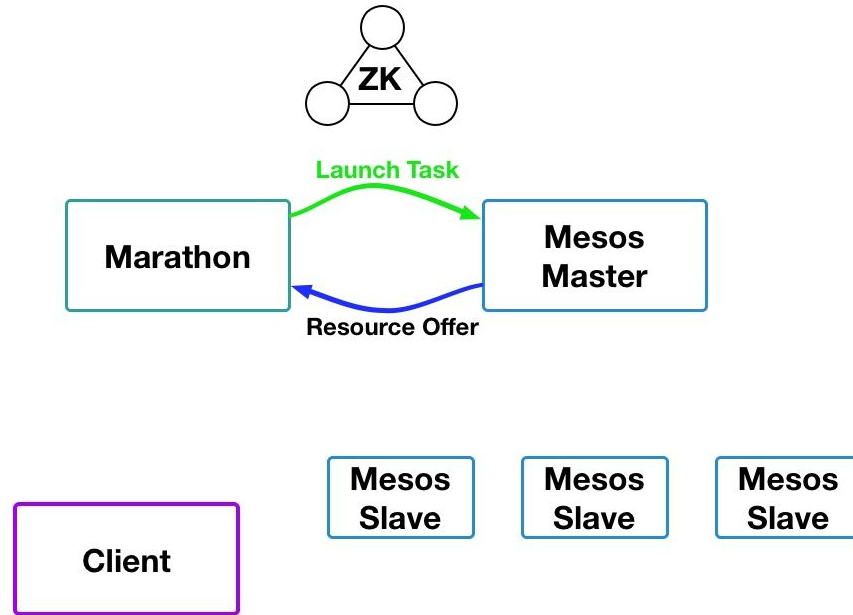# USEFUL MARATHON FEATURES: DEPLOY LIKE A TELCO



Credit: Martin Fowler

- Application versioning
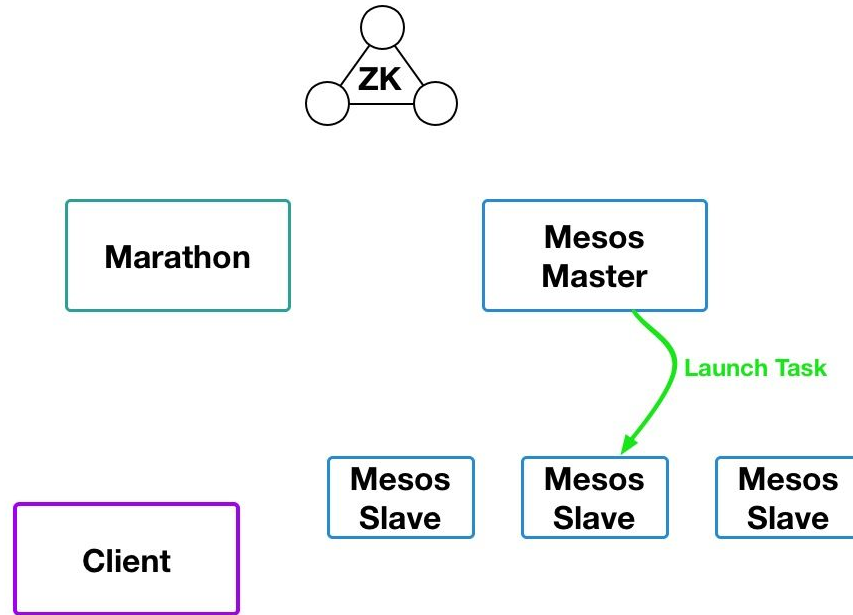- Hot/hot new/old clusters
- Authentic scale testing
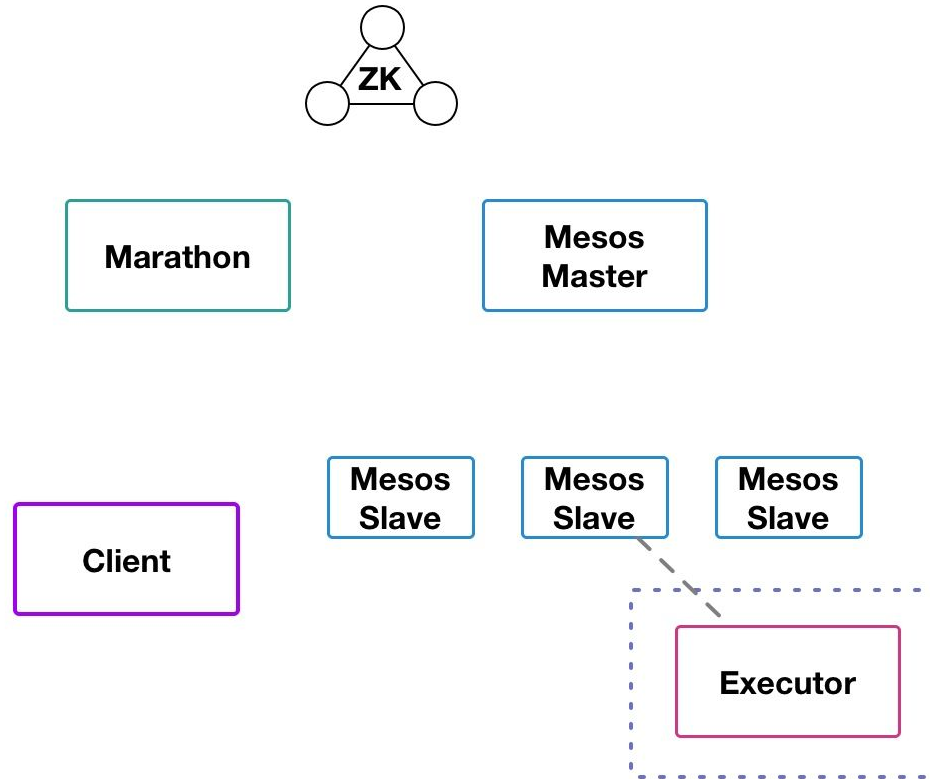- Manual *and* automated testing

# MESOS WITH MARATHON IN ACTION

ZK
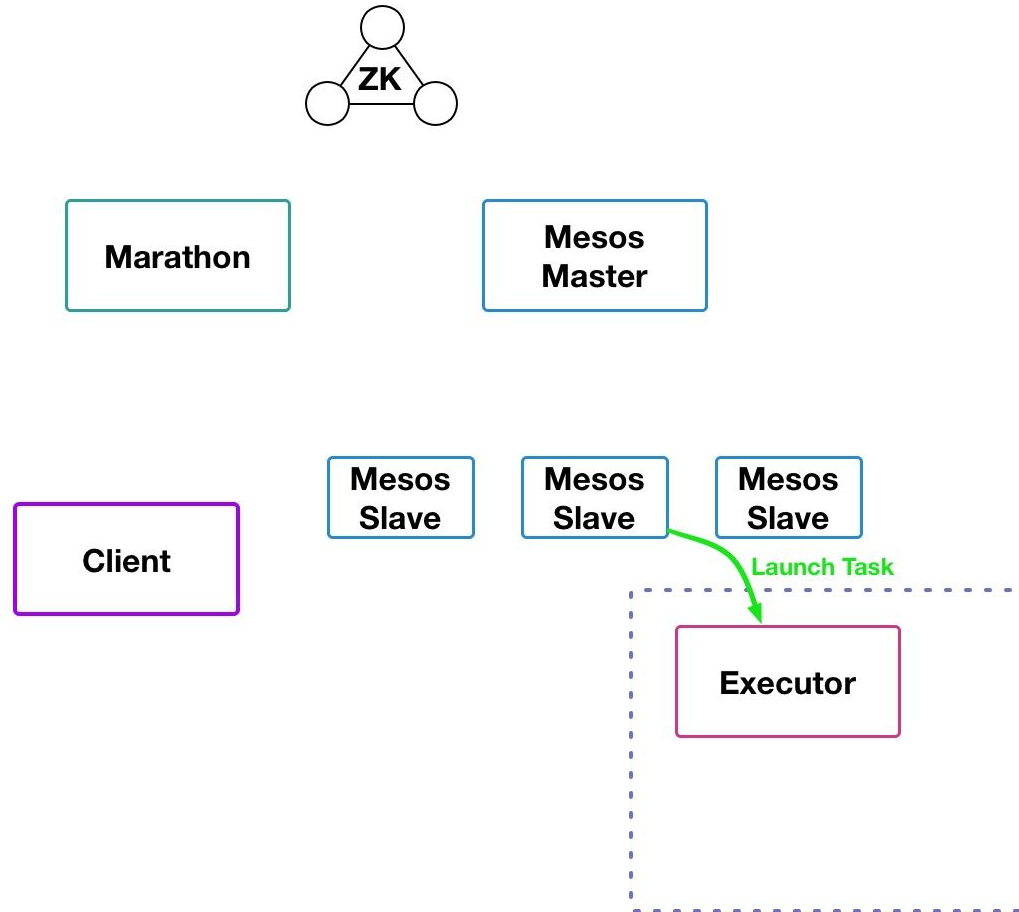
Marathon

Mesos
Master

App (JSON)

Client

Mesos
Slave

Mesos
Slave

Mesos
Slave

ZK

Marat[hon]

App

Client

```
{

  "id": "foo",
  "uris": ["http://my.org/foo.tgz"],
  "cmd": "./run.sh --env=prod",
  "cpus": 2.0,
  "mem": 4096.0

}
```

ZK

Marathon

Mesos
Master

Launch Task

Mesos
Slave

Mesos
Slave

Mesos
Slave

Client

ZK

Marathon

Mesos
Master

Mesos
Slave

Mesos
Slave

Mesos
Slave

Client

Executor

ZK

Marathon

Mesos Master

Mesos Slave

Mesos Slave

Mesos Slave

Client

Launch Task

Executor

ZK

Marathon

Mesos Master

Client

Mesos Slave

Mesos Slave

Mesos Slave

Status Update

Executor

Task

ZK

**Marathon**

**Mesos Master**

**Status Update**

**Mesos Slave**　　**Mesos Slave**　　**Mesos Slave**

**Client**

**Executor**

**Task**

ZK

Marathon

Mesos
Master

Status Update

State (JSON)

Client

Mesos
Slave

Mesos
Slave

Mesos
Slave

Executor

Task

ZK

Marathon

Mesos
Master

Mesos
Slave

Mesos
Slave

Mesos
Slave

Client

Executor

Task

Mesos with Marathon in Action

# TASK FAILURE

ZK

Marathon

Mesos
Master

**Status Update**

Status Update

Mesos
Slave

Mesos
Slave

Mesos
Slave

Client

Executor

SEGFAULT :(

Task

ZK

Marathon

Mesos
Master

Status Update

Status Update

Client

Mesos
Slave

Mesos
Slave

Mesos
Slave

Status Update

Executor

Task

ZK

Marathon

Mesos
Master

Mesos
Slave

Mesos
Slave
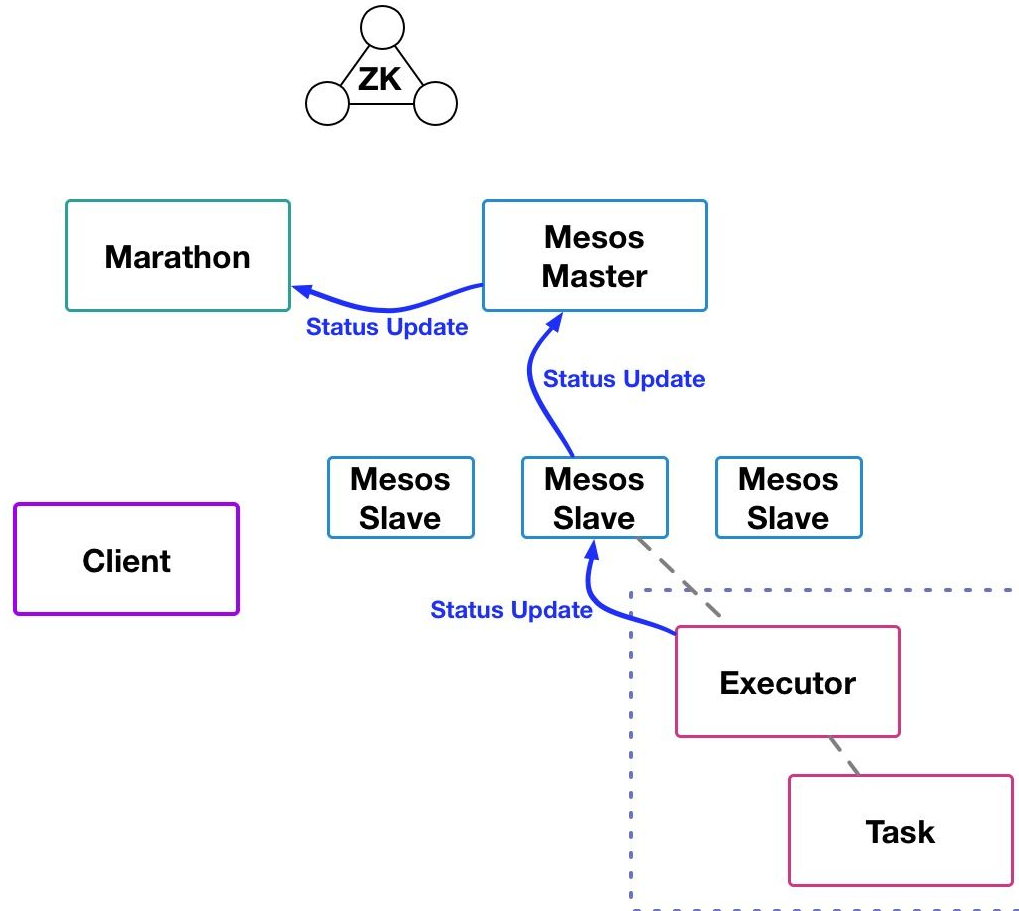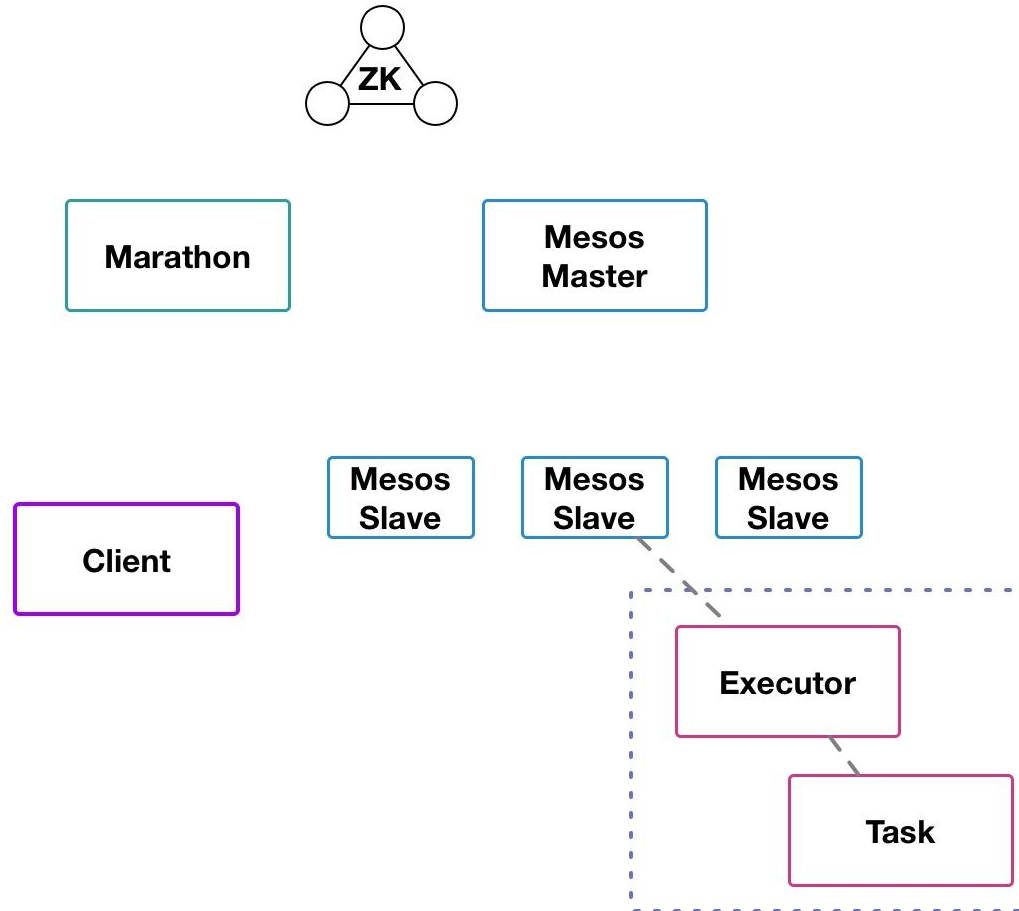
Mesos
Slave

Client

Executor
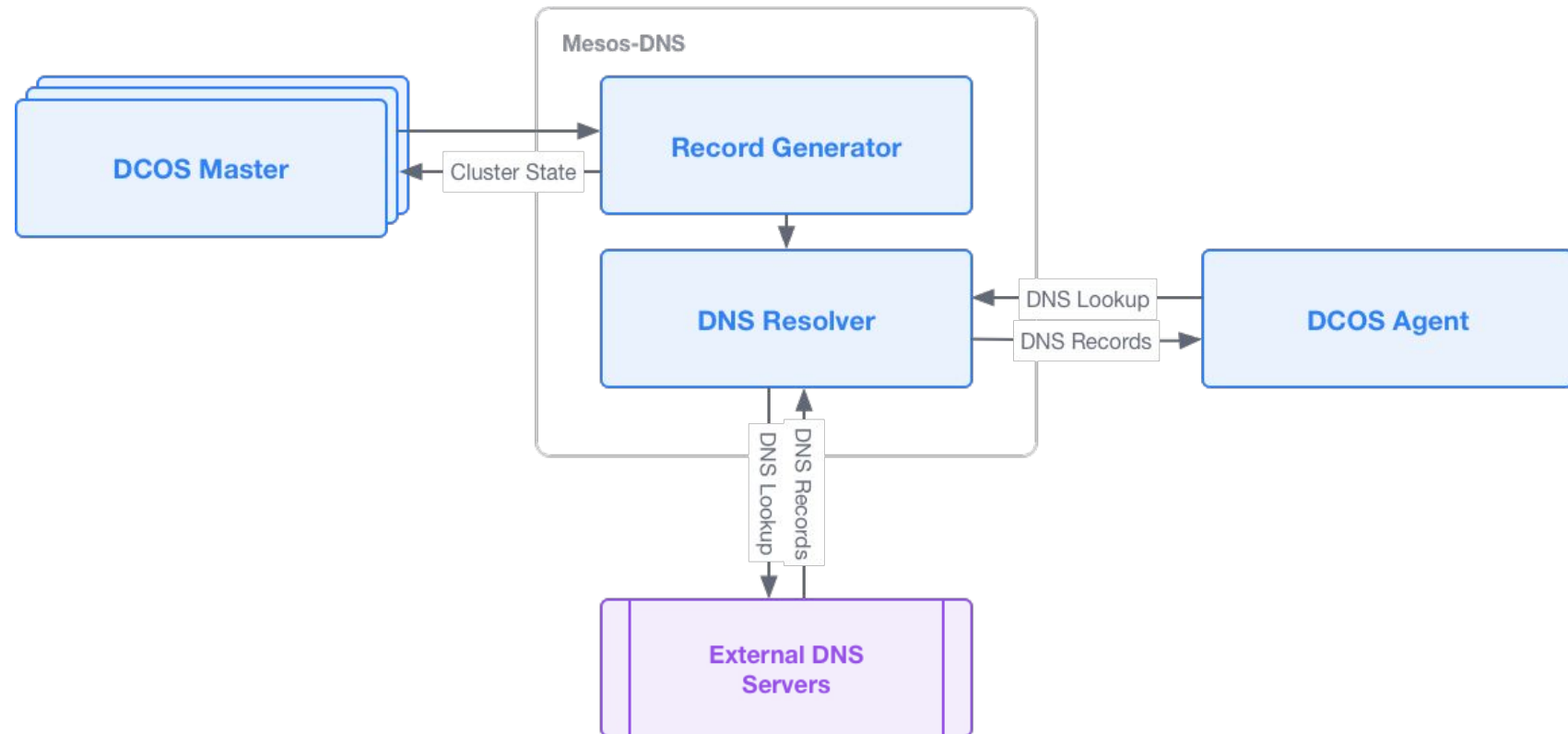
Task

# SERVICE DISCOVERY

How do my applications discover each other?

Two main service discovery mechanisms:

1. DNS based (Mesos-DNS)
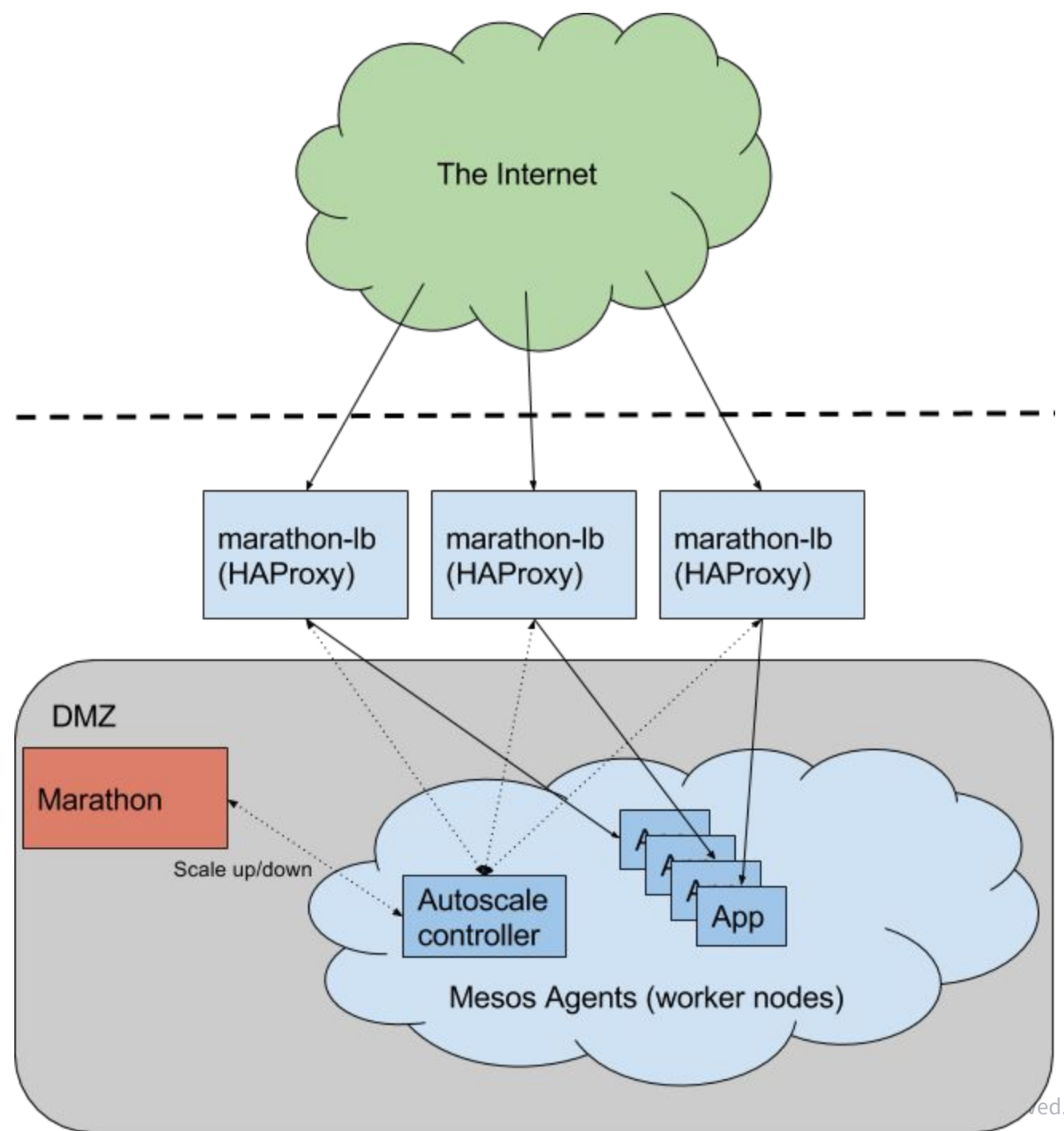2. HAProxy based (Marathon-lb)

# MESOS-DNS

- Ingests cluster state periodically.

- Uses cluster state to generate DNS records for all running Mesos tasks.

- Services query DNS server to discover IP address and port of other services.

- Primarily used for internal service discovery.

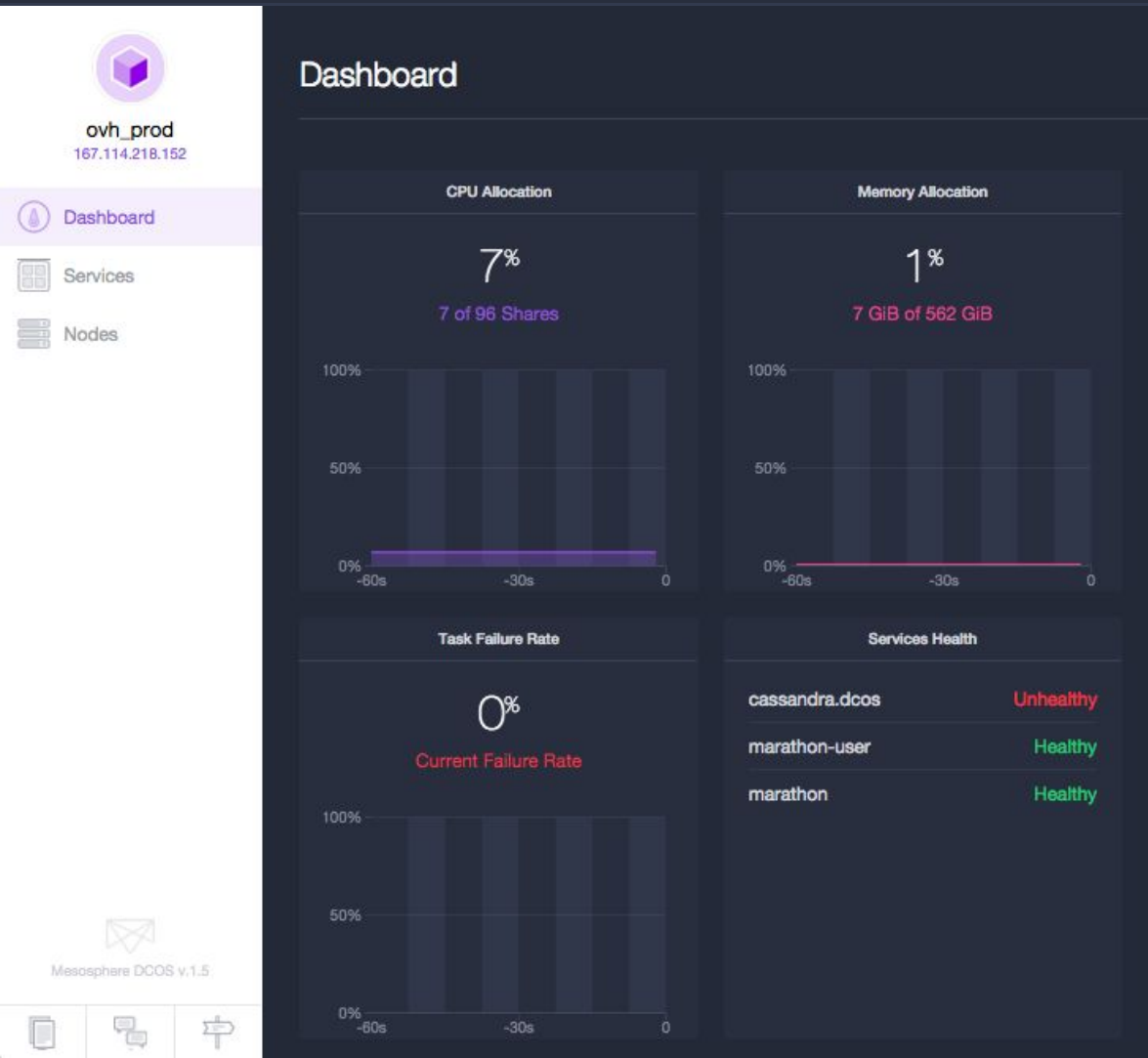- No extra configuration required!



**49**

# MARATHON-LB

- Ingests state of running Marathon applications.

- Regenerates HAProxy configuration.

- Supports virtual hosts!

- Can be used for both internal and external service discovery.

- Must add HAPROXY_GROUP and HAPROXY_0_VHOST variables to your marathon.json.
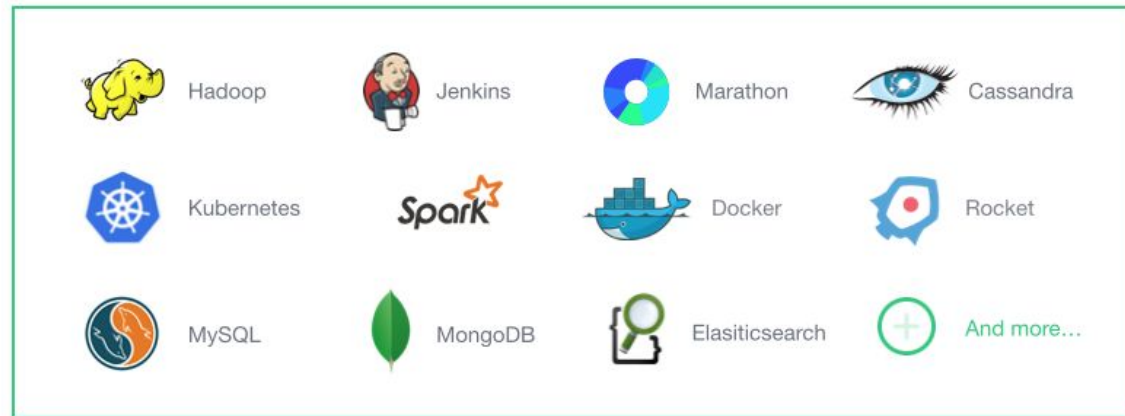
# HOW TO DEPLOY A MESOS CLUSTER (THE HARD WAY)

- Using chef/puppet/ansible  (or a reliable intern)
- Install ZooKeeper and Mesos
- Install your scheduler (Marathon)
- Deploy some long-running services.
- See https://open.mesosphere.com/getting-started/tools/ for more docs

# HOW TO DEPLOY A MESOS CLUSTER (OUR WAY)



- Visit http://mesosphere.com
- Hit the 'Get Started' button

# MESOS AS THE DATACENTER KERNEL



Services & Containers

Hadoop · Jenkins · Marathon · Cassandra

Kubernetes · Spark · Docker · Rocket

MySQL · MongoDB · Elasiticsearch · And more…

Mesosphere DCOS

Container Orchestration · Security & Governance · Monitoring & Operations · User Interface

MESOS

Existing Infrastructure

Physical · VMs · Private Cloud · Google · Amazon · Azure
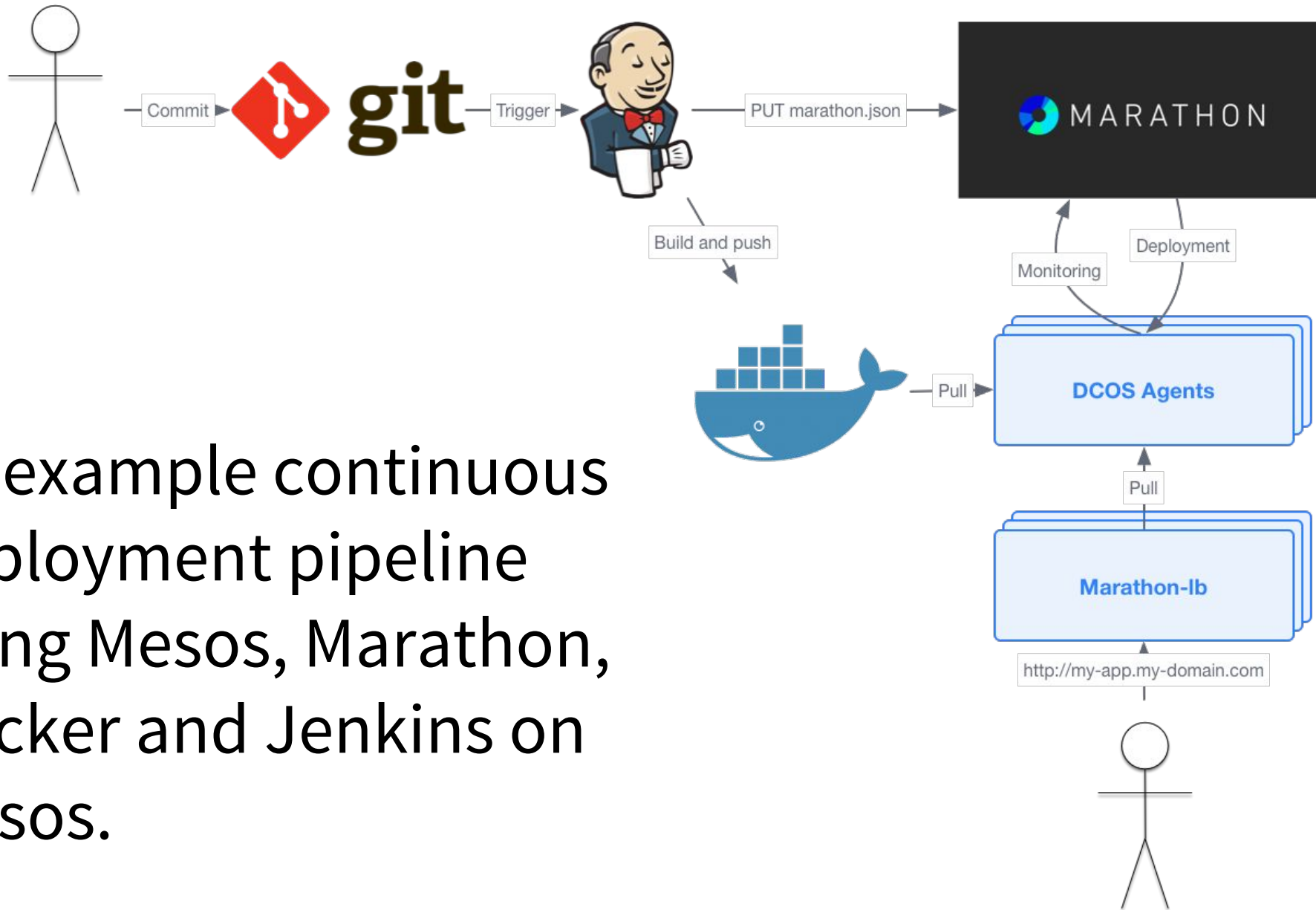
# JENKINS: BUILD RESOURCE POOLING

**Jenkins on Mesos** allows you to share build resources between multiple Jenkins masters.
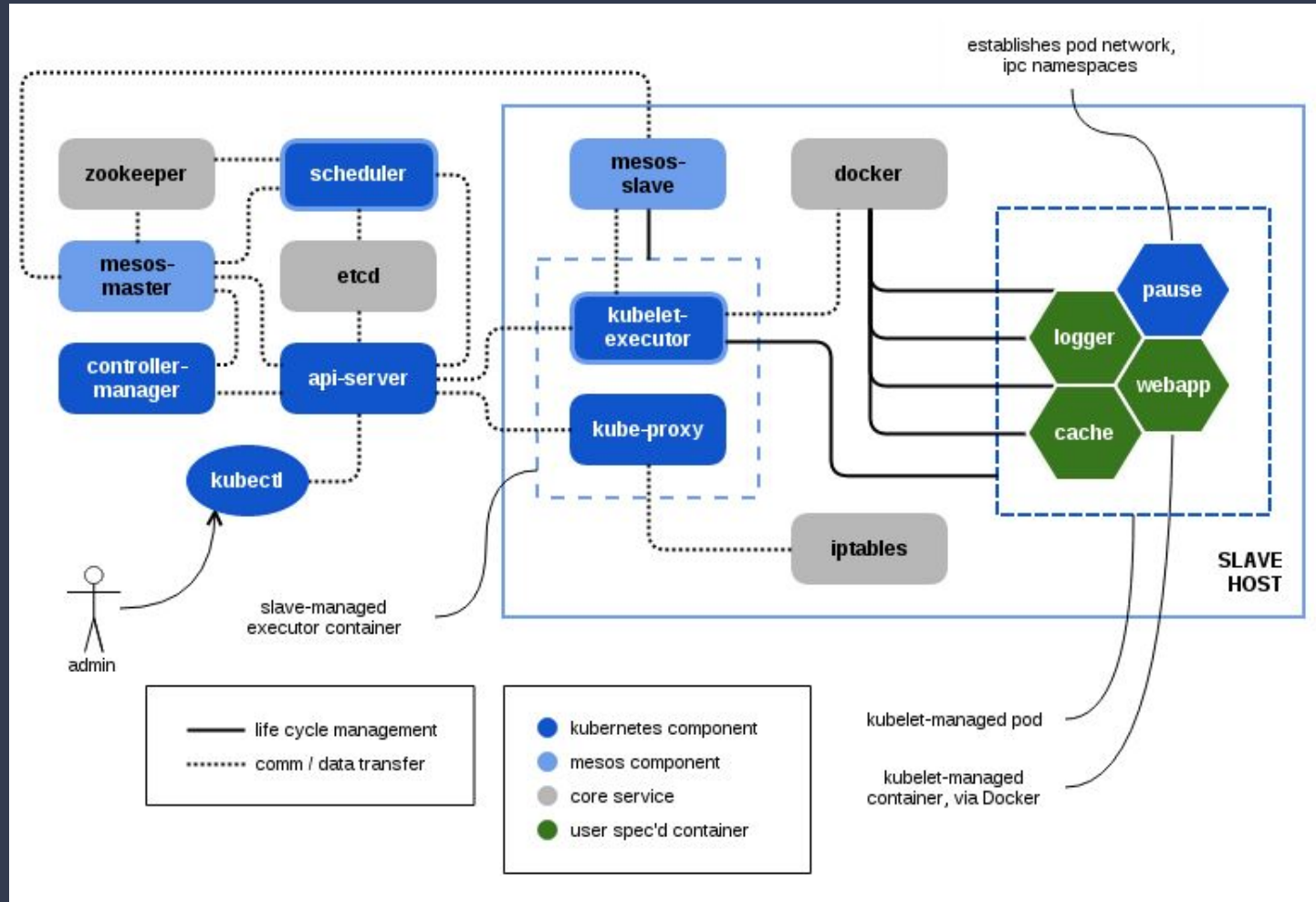
- PayPal does this with hundreds of Jenkins masters
- Between them, they use less than a hundred build slaves to service several thousand developers.
- Combining Jenkins with a PaaS like Marathon or Kubernetes allows you to practice easy continuous deployment.
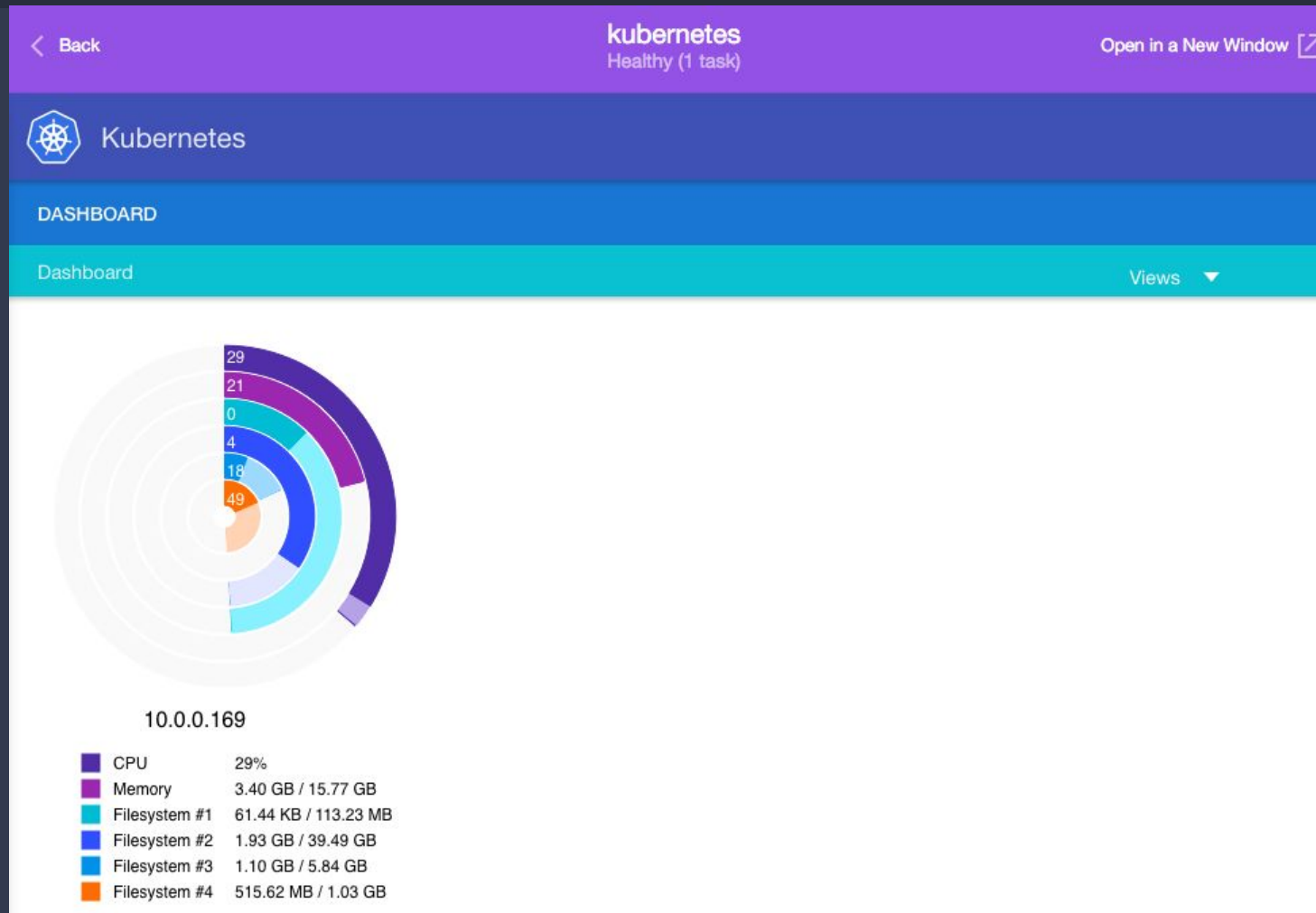
An example continuous deployment pipeline using Mesos, Marathon, Docker and Jenkins on Mesos.

# CONTINUOUS DELIVERY DEMO

# KUBERNETES ON MESOS

# KUBERNETES ON MESOS

# BIG DATA ON MESOS

Mesos was built for and is great for running big data workloads:

- Chronos (time scheduled jobs)
- Spark
- Cassandra
- Kafka
- Hadoop/YARN (via Myriad)

# QUESTIONS? THANK YOU!

Come and talk to us!
- Email us at philip@mesosphere.io, sunil@mesosphere.io
- Slides will be up at http://mesosphere.github.io/presentations