Neil Conway, Niklas Nielsen, Greg Mann & Sunil Shah
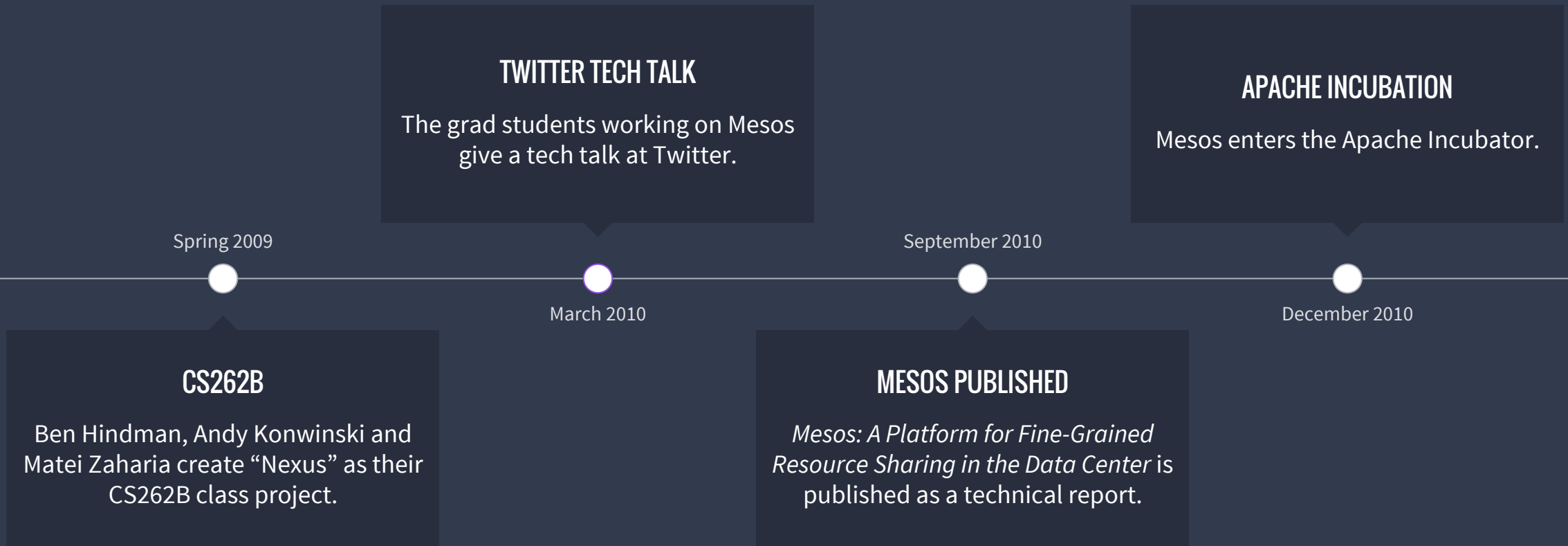
# POWERING THE INTERNET WITH APACHE MESOS

MESOSPHERE

# 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
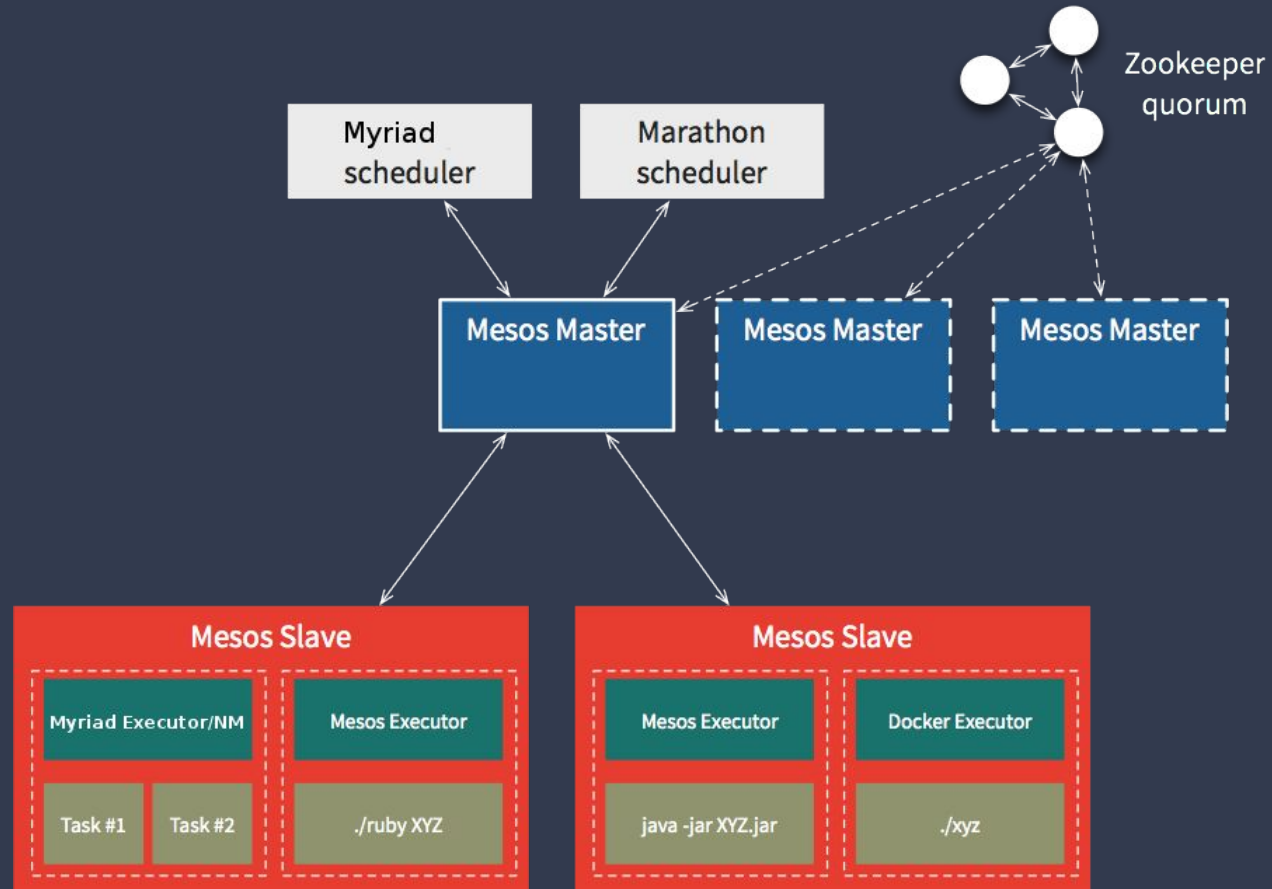*University of California, Berkeley*

**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

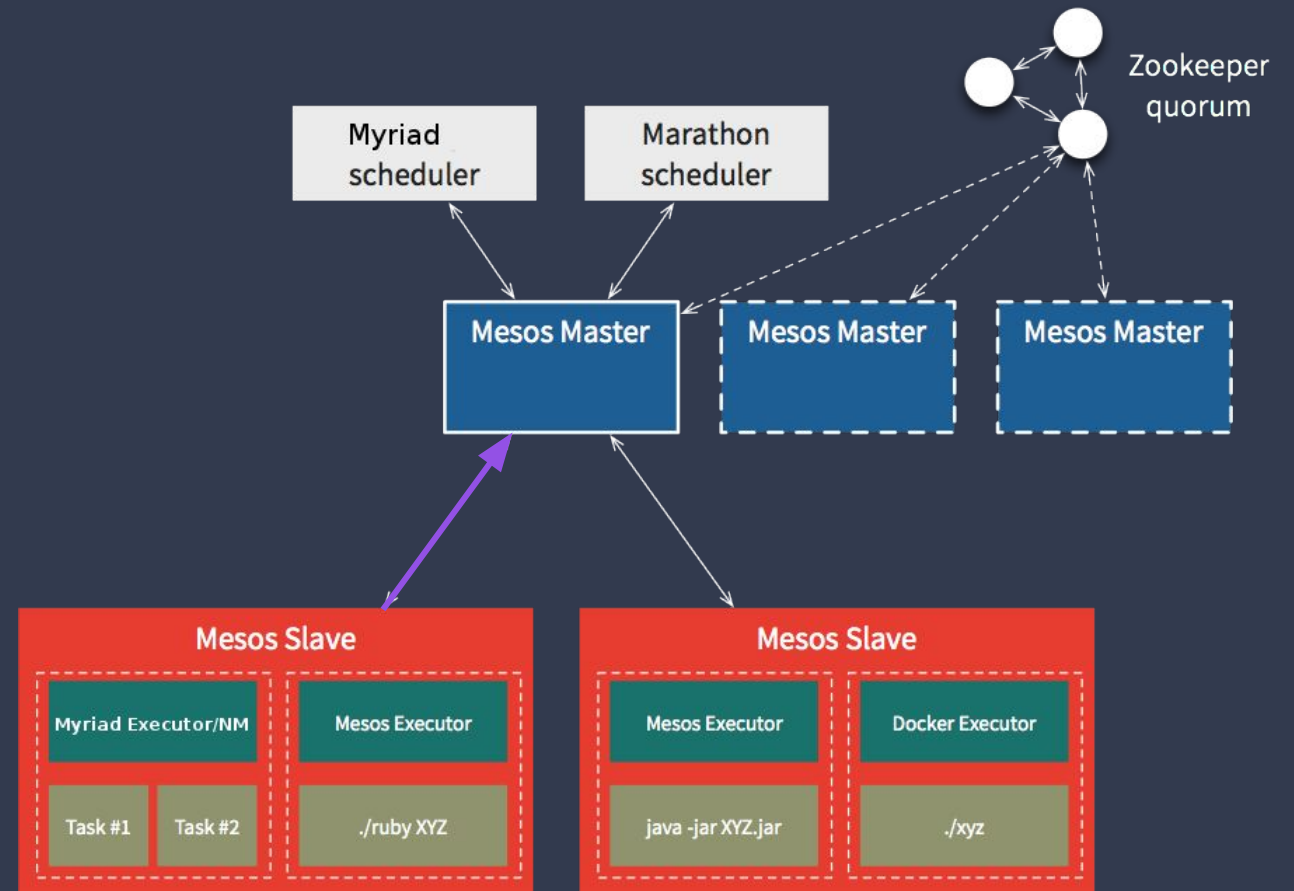What does an operating system provide?

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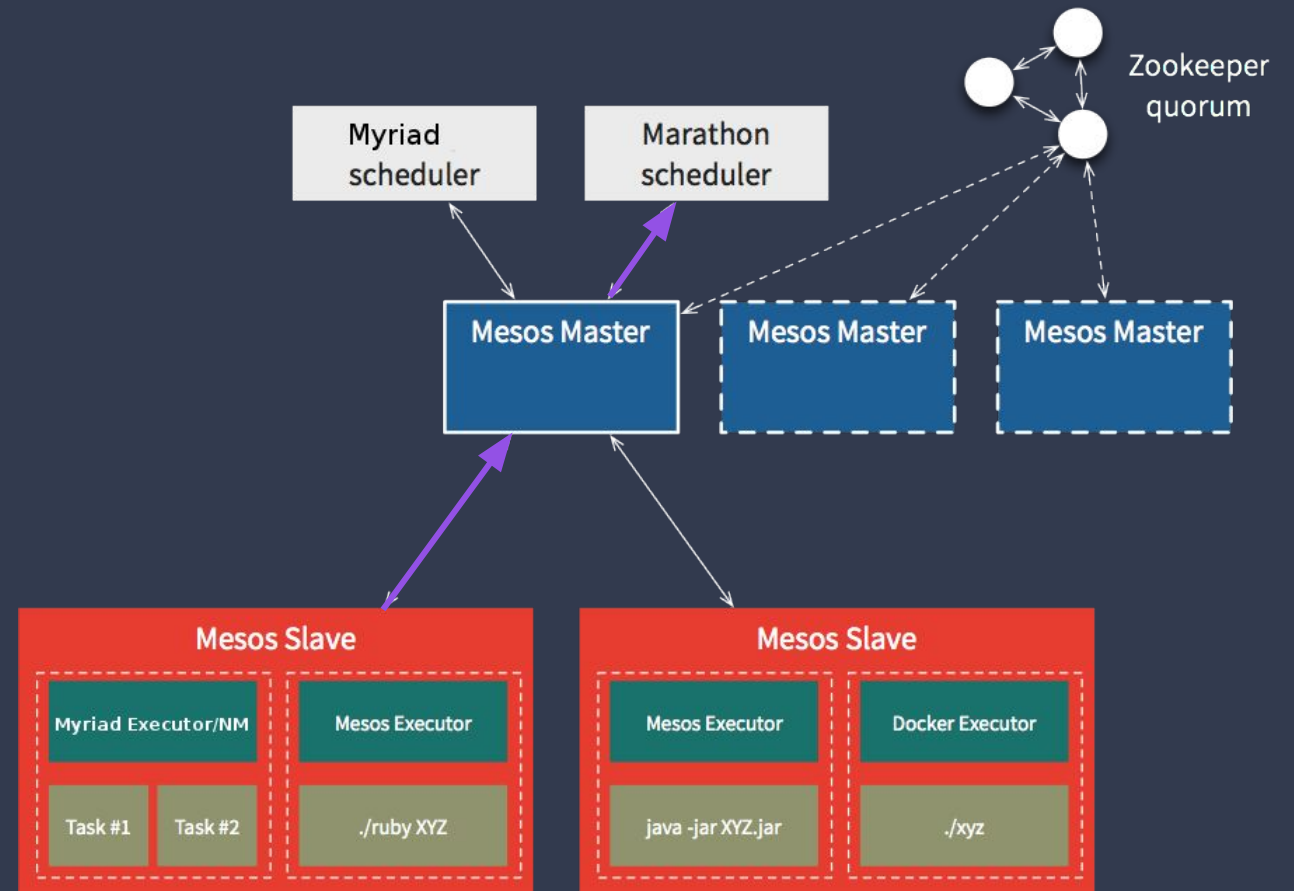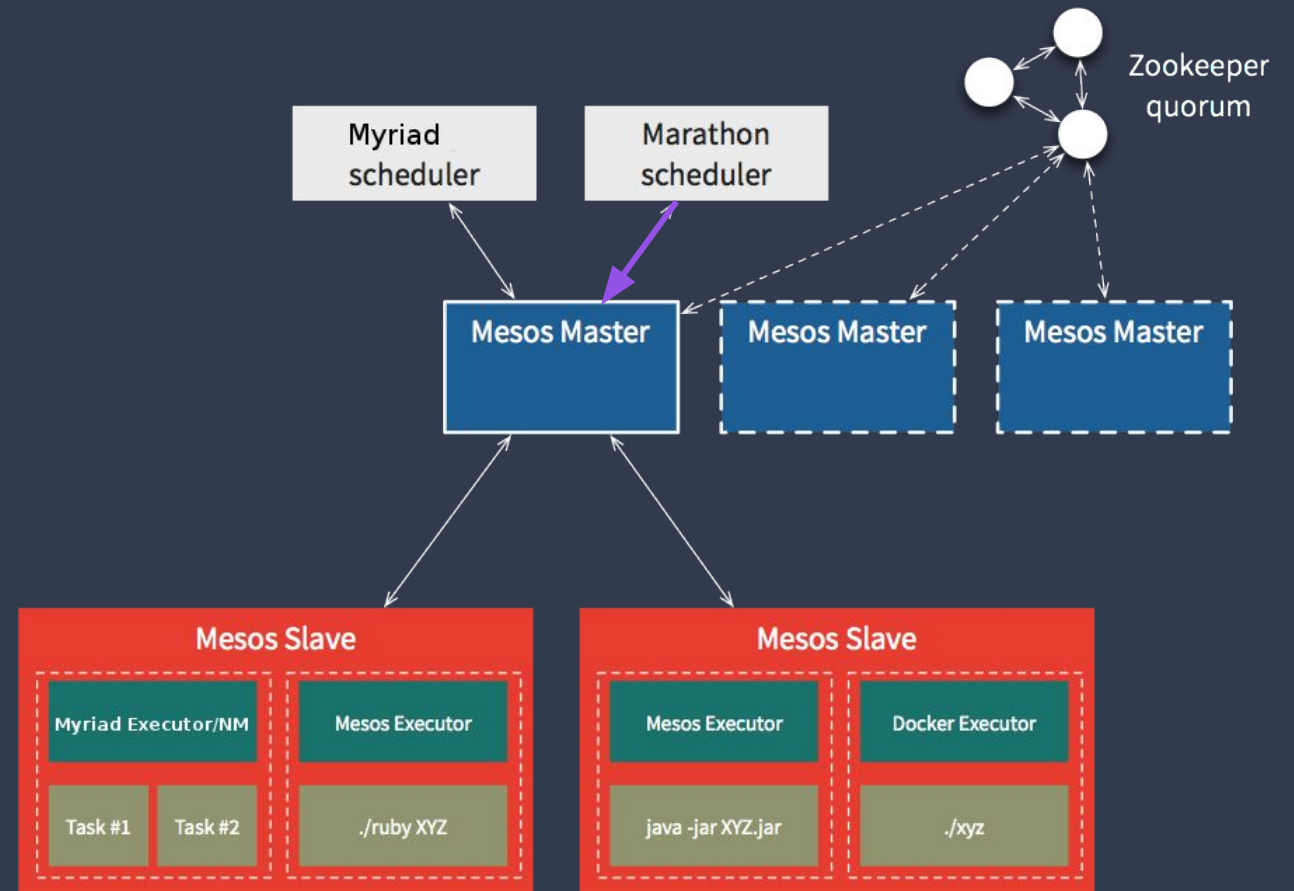- Monitoring, debugging, logging

# ARCHITECTURE

# ARCHITECTURE

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

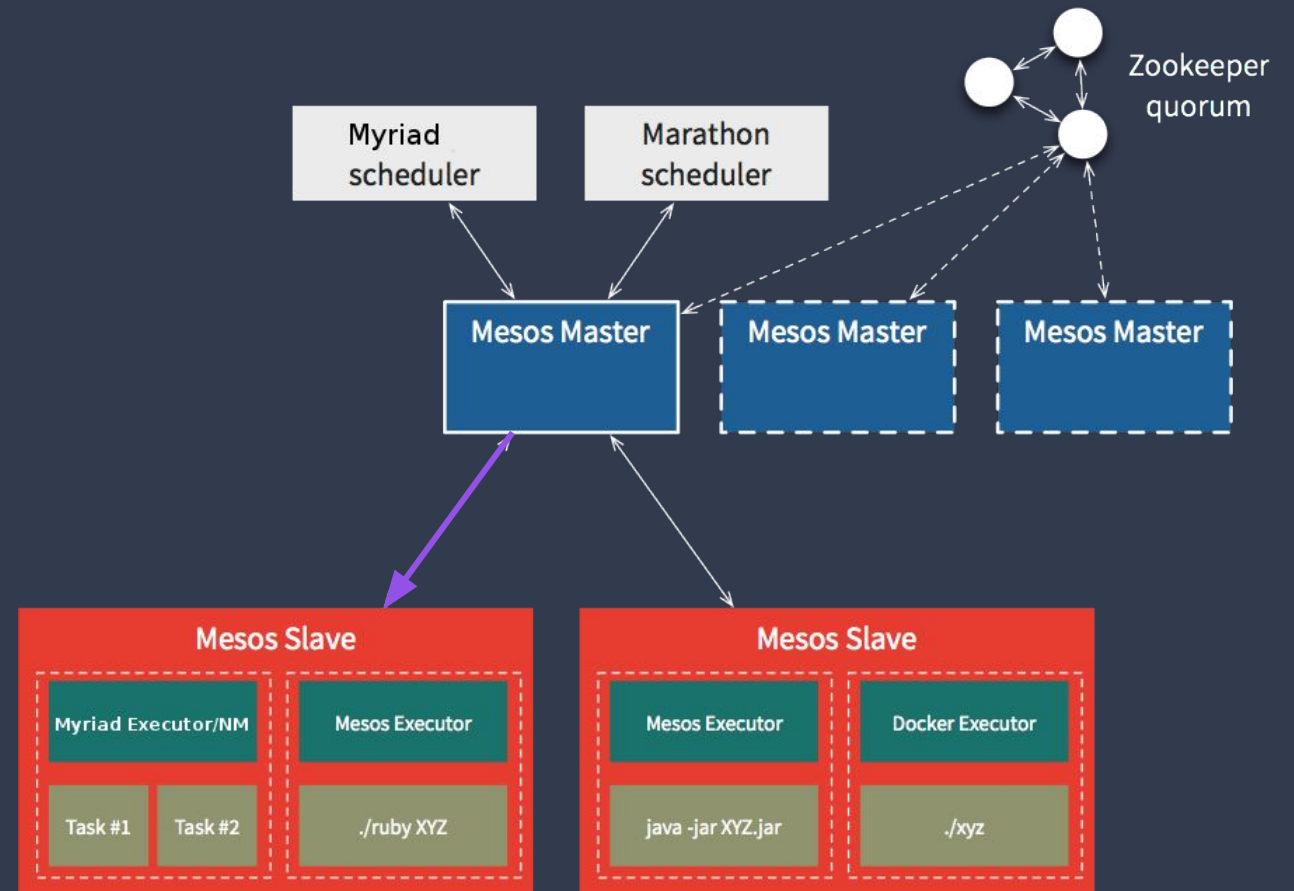# ARCHITECTURE

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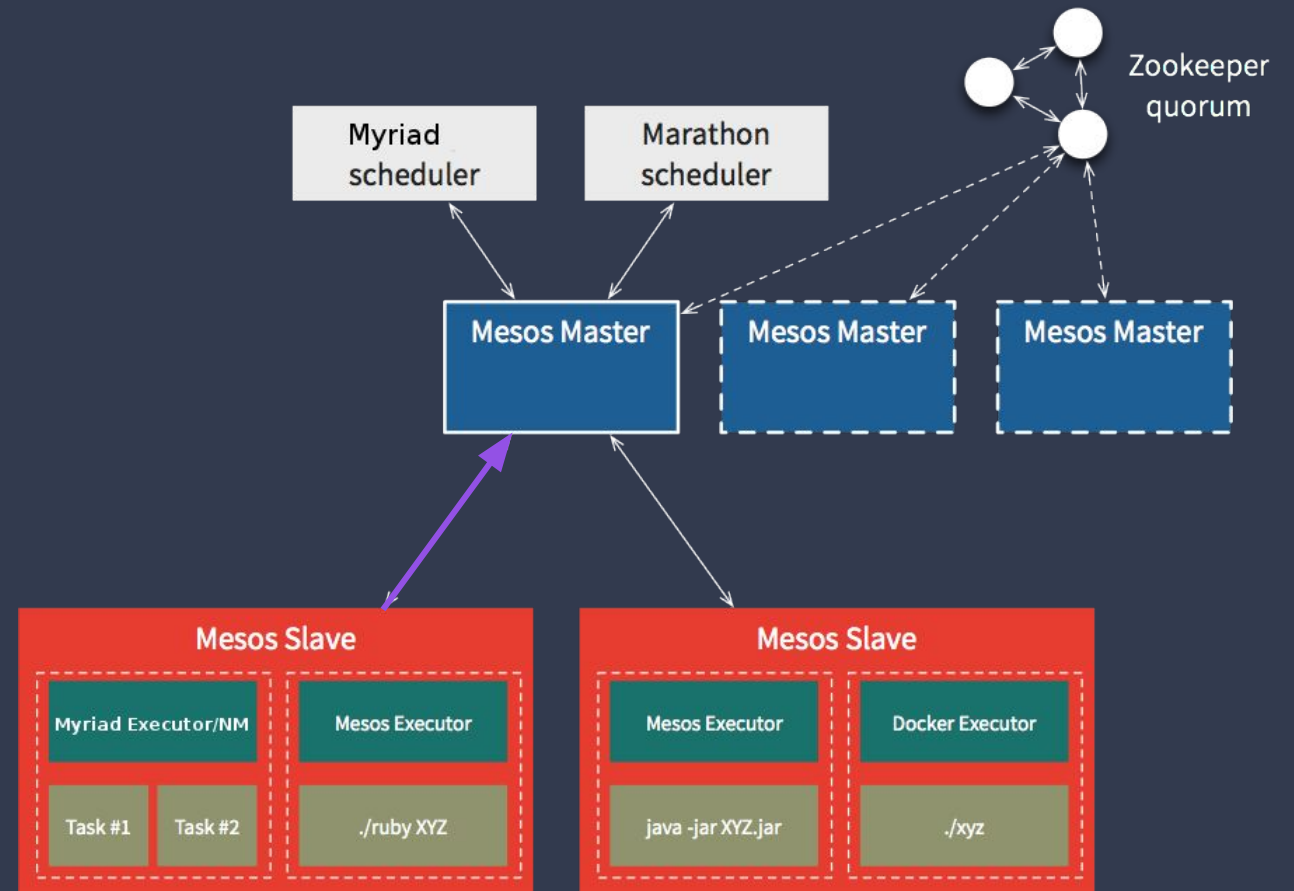- Agents report task status to Master

# ARCHITECTURE

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

# ARCHITECTURE

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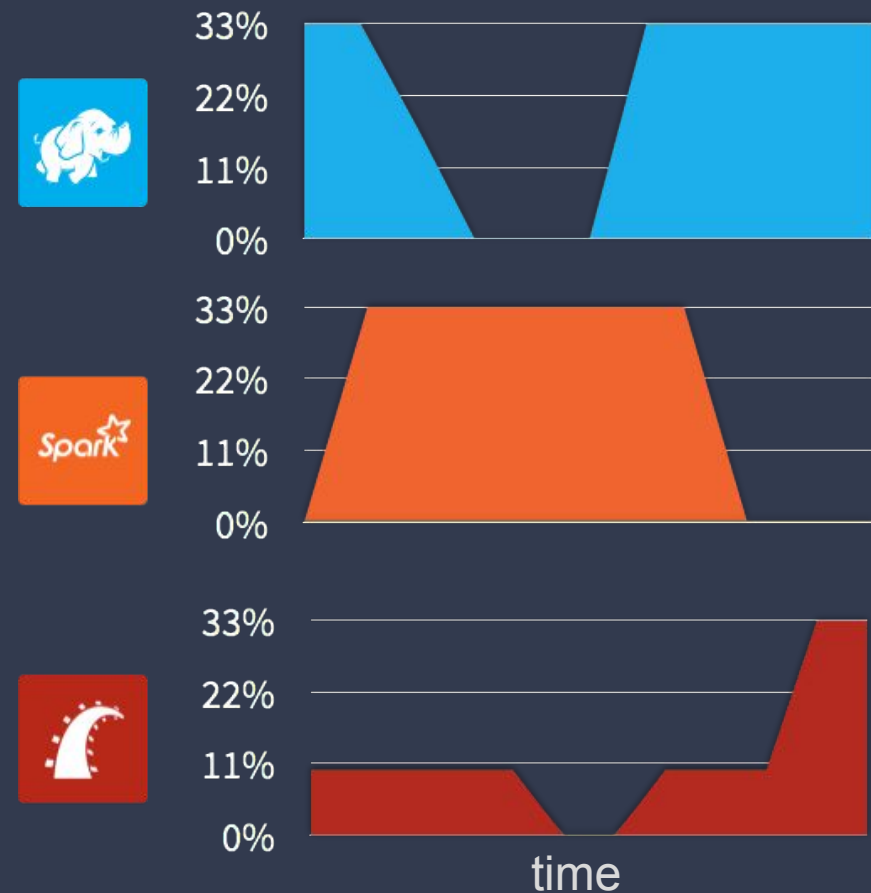- Agents report task status to Master

# ARCHITECTURE

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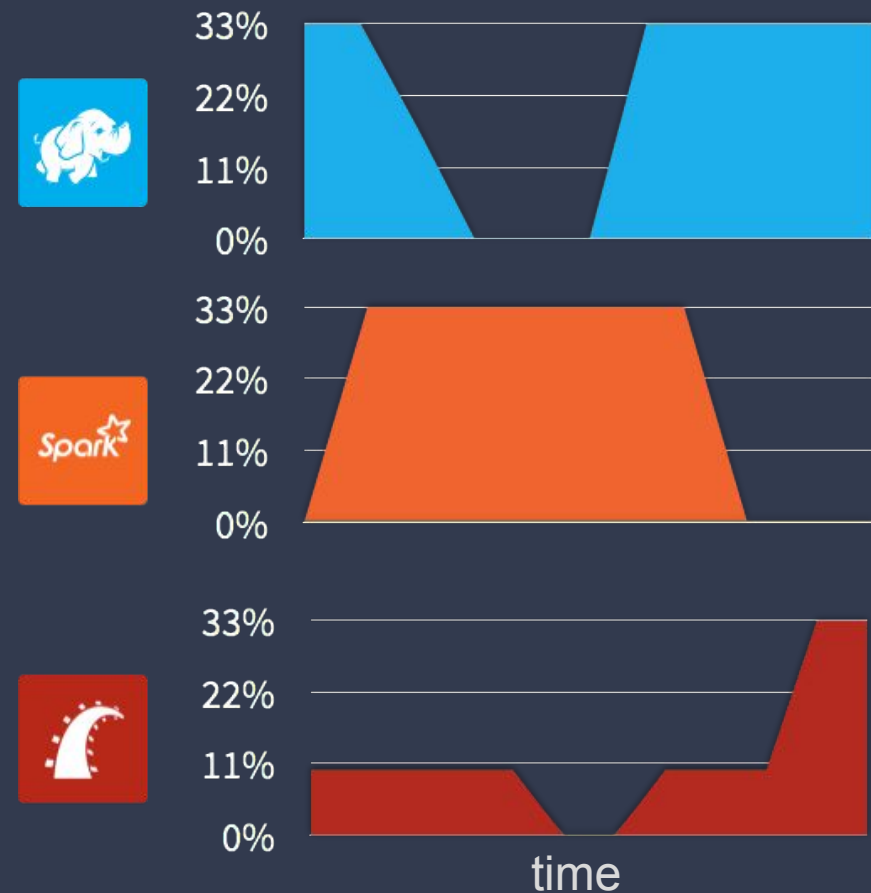- Agents report task status to Master

# KEEP IT STATIC

A naive approach to handling varied app requirements: **static partitioning**.

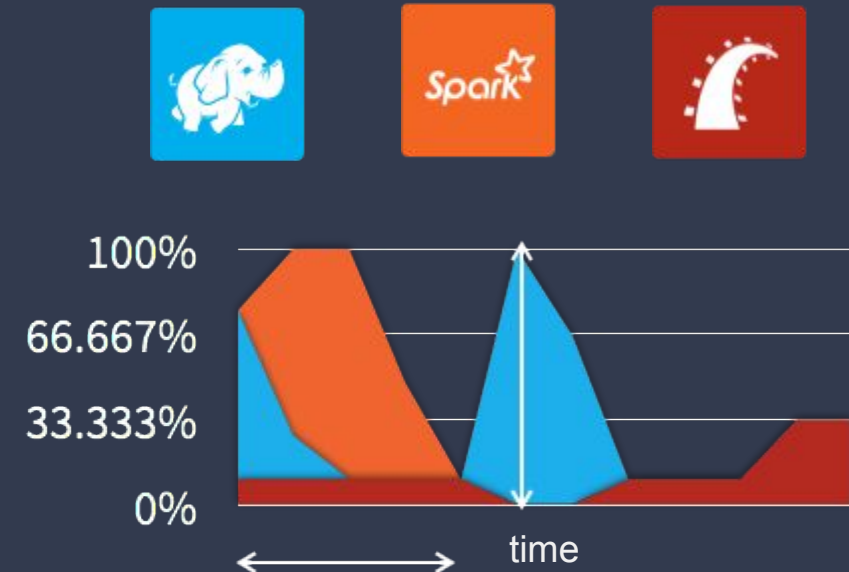This can cope with heterogeneity, but is very expensive.

# KEEP IT STATIC

Maintaining sufficient headroom to handle peak workloads on all partitions leads to **poor utilization** overall.
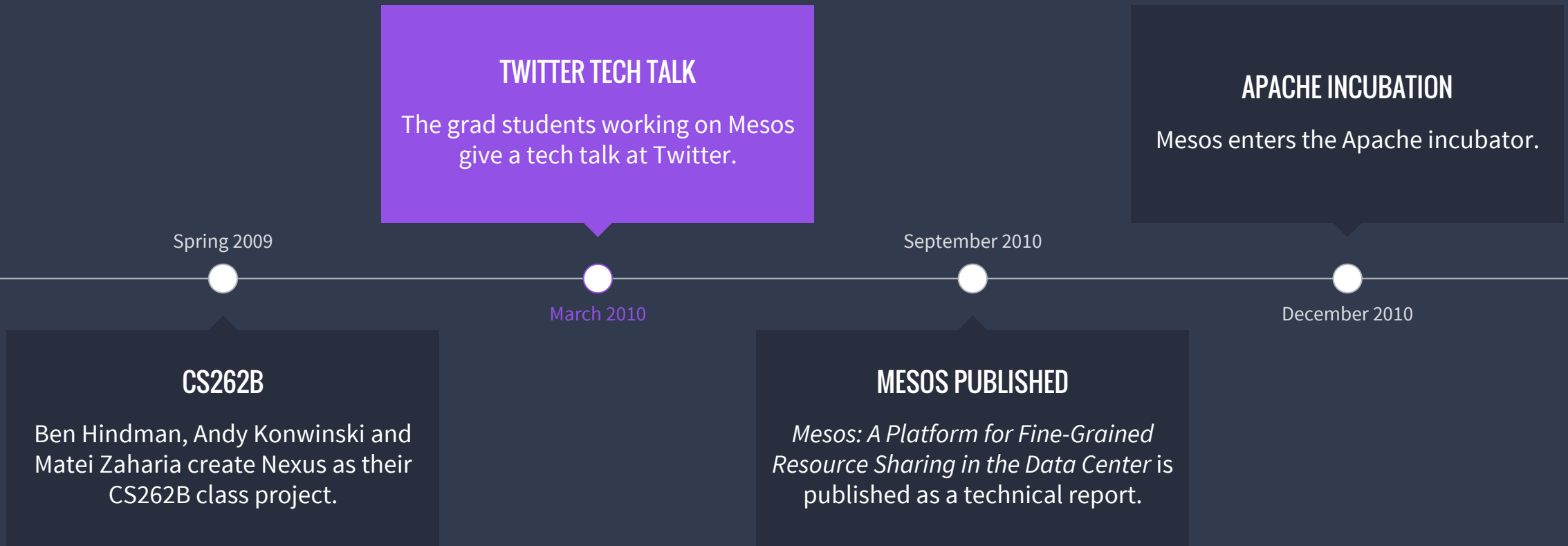
# SHARED RESOURCES

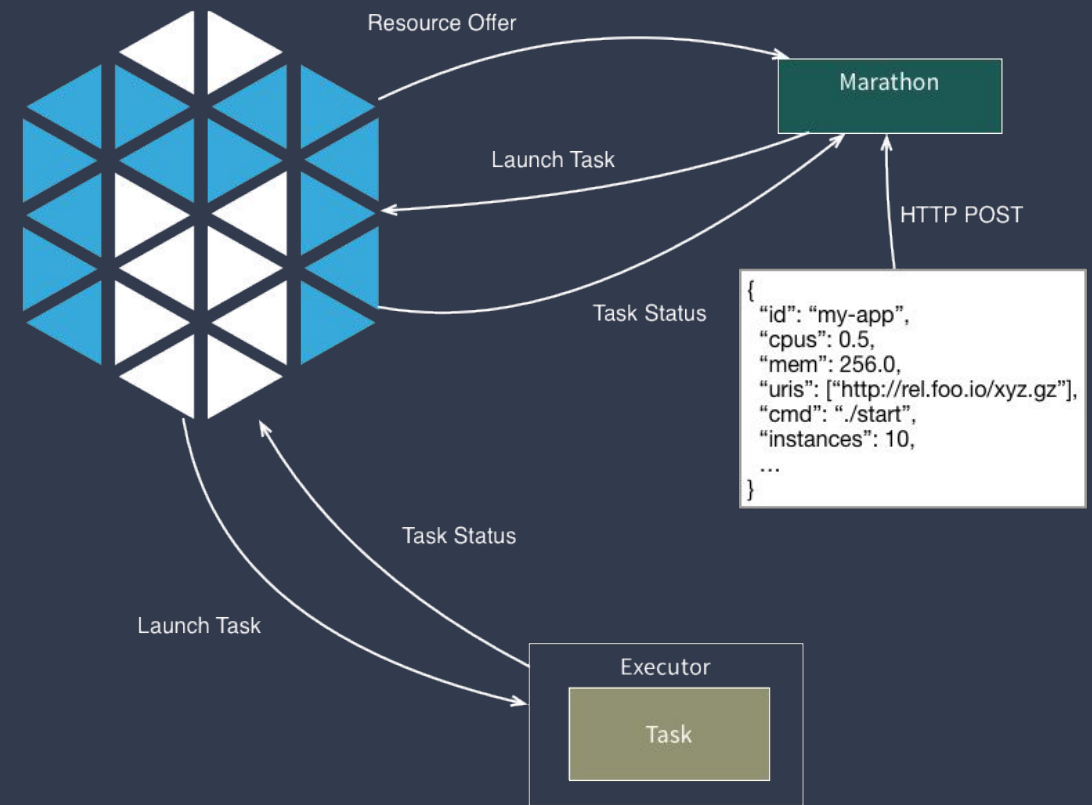Multiple frameworks can use the same cluster resources, with their share adjusting dynamically.

# 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.
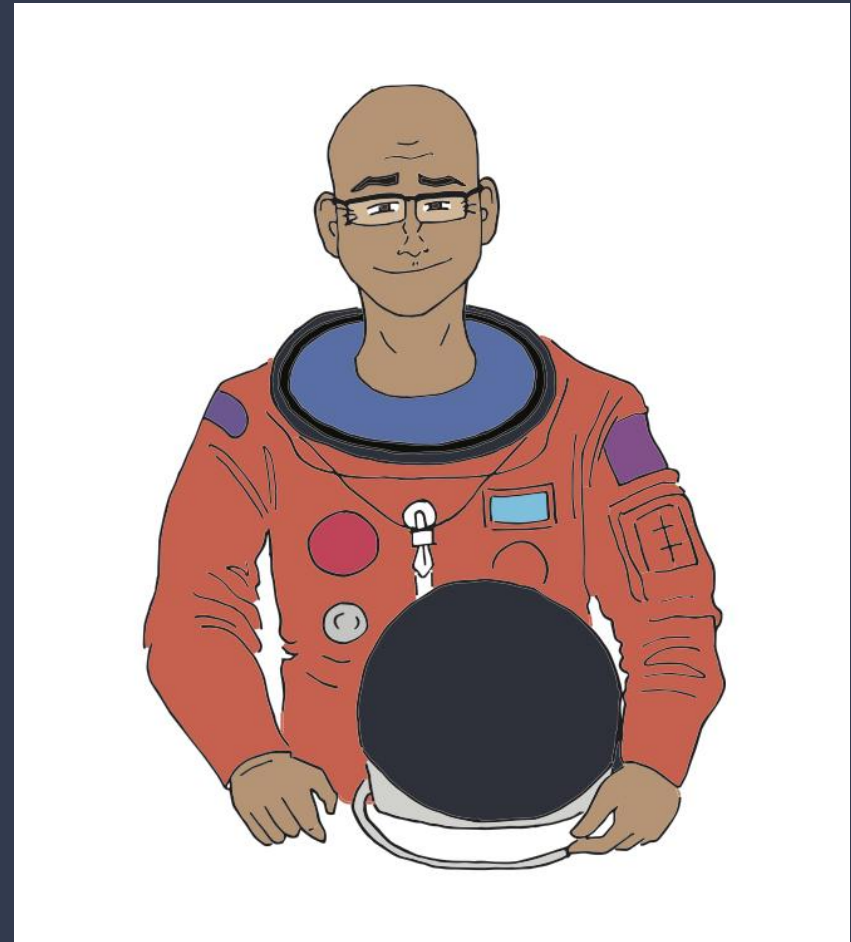
# MESOS REALLY HELPS

- 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
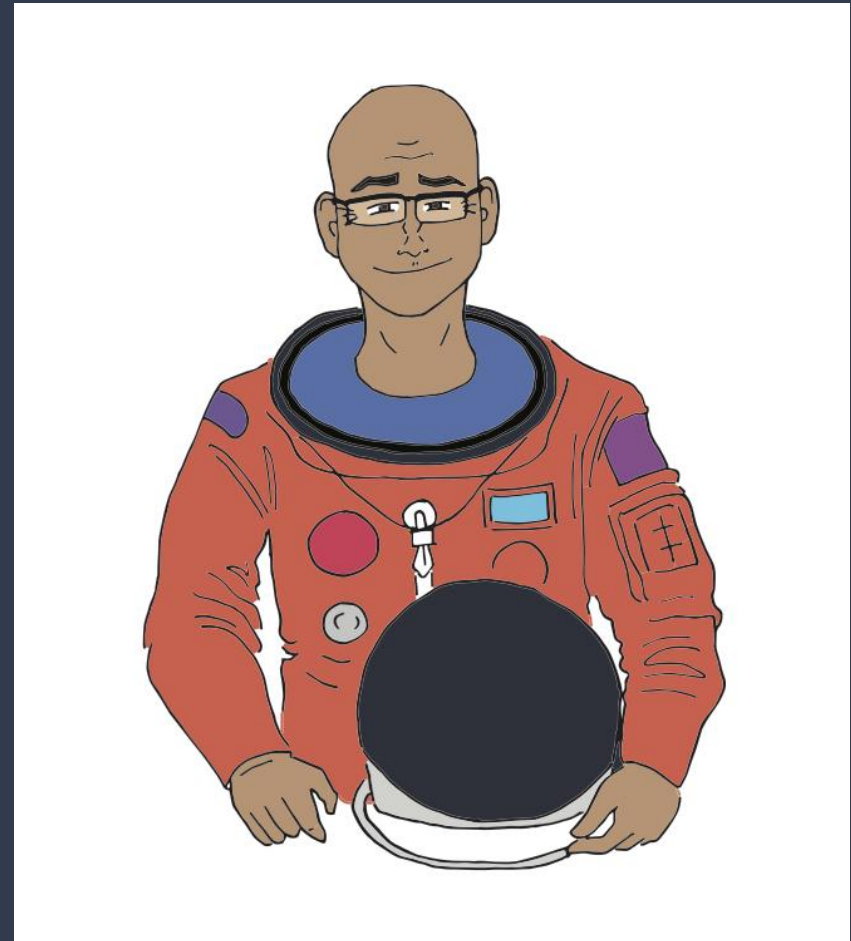
Executor

Task

# LIFE WITHOUT MESOS

# SAY HI TO DAN

- Dan is a member of operations staff in a non-Google, non-Facebook company with large and growing users and workloads.

# SO MUCH PROCESS

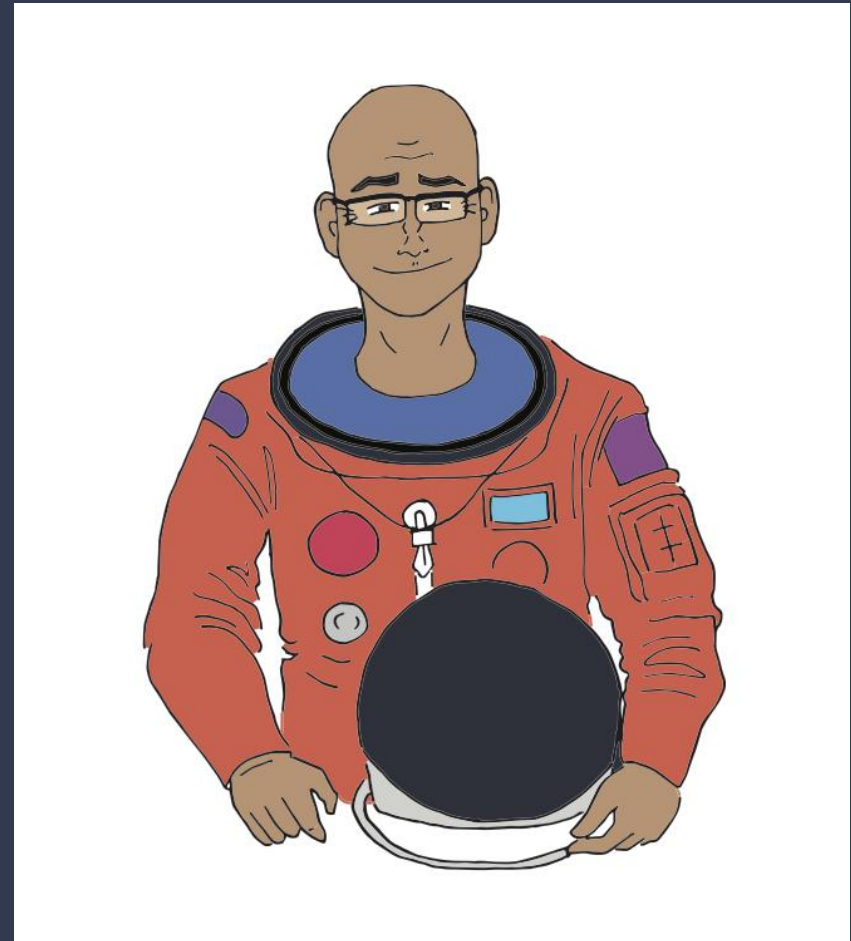You're a developer wanting to deploy a new service.

1. How many resources do you need?
   (Better overestimate, it usually takes a while to provision these.)
2. What dependencies does your application have?
3. Who monitors your applications and handles it falling over?

# CHANGE IS PAINFUL

You have more users and/or you want to upgrade your application.

1. Submit another resource request.
2. Provision new machines.
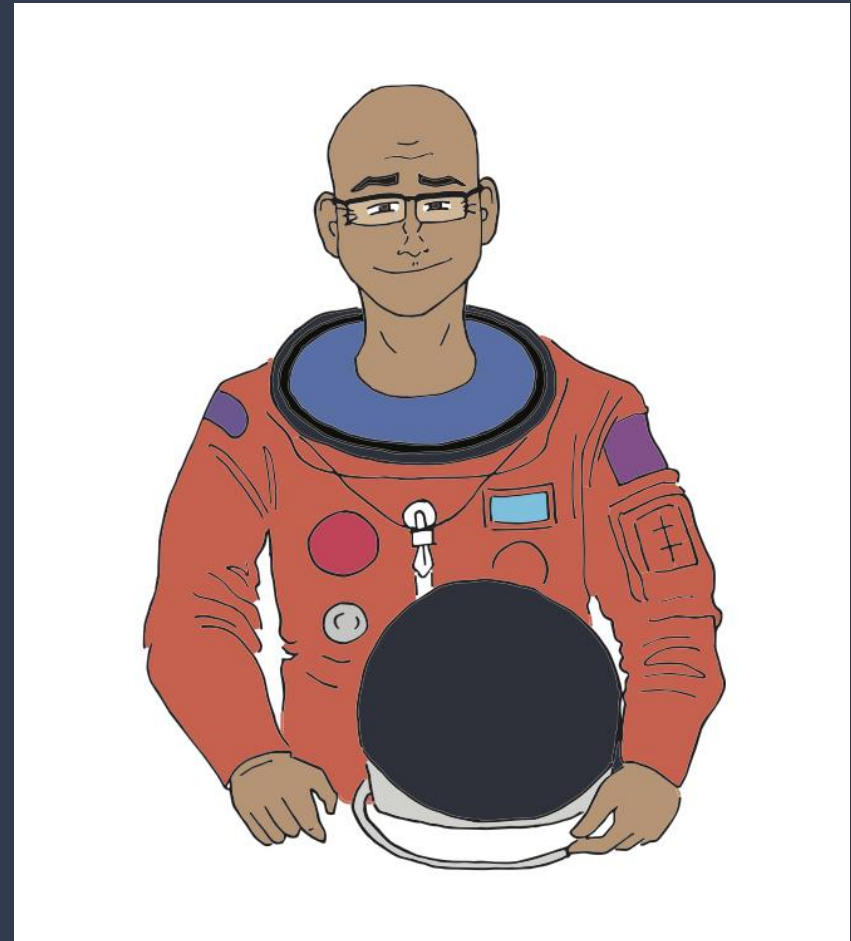3. How do we get any upgraded binaries/dependencies to existing machines?
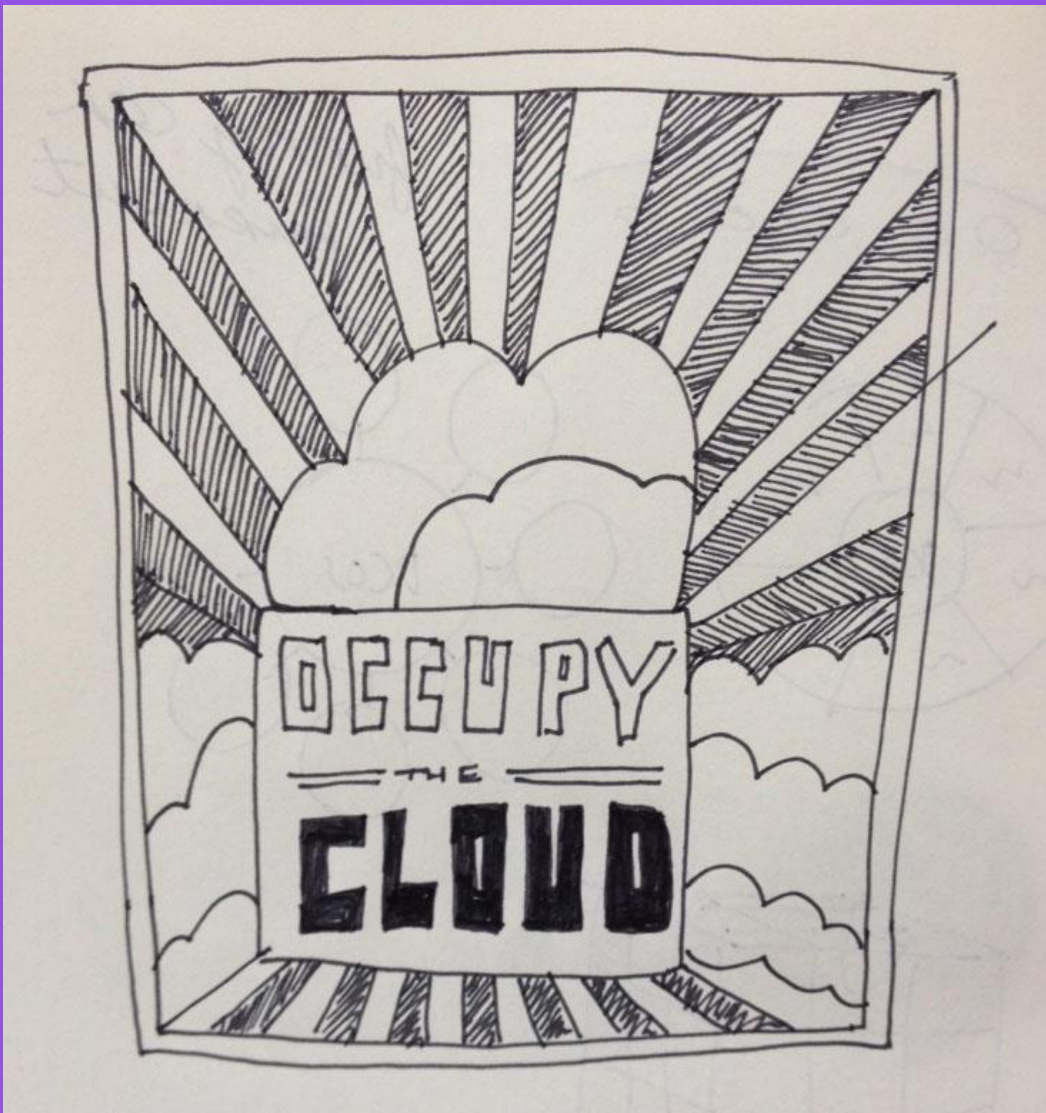
# COMPLEX WORKLOADS

Dan, our operator, is forced to partition the datacenter to accommodate these demands. **Utilisation suffers**.
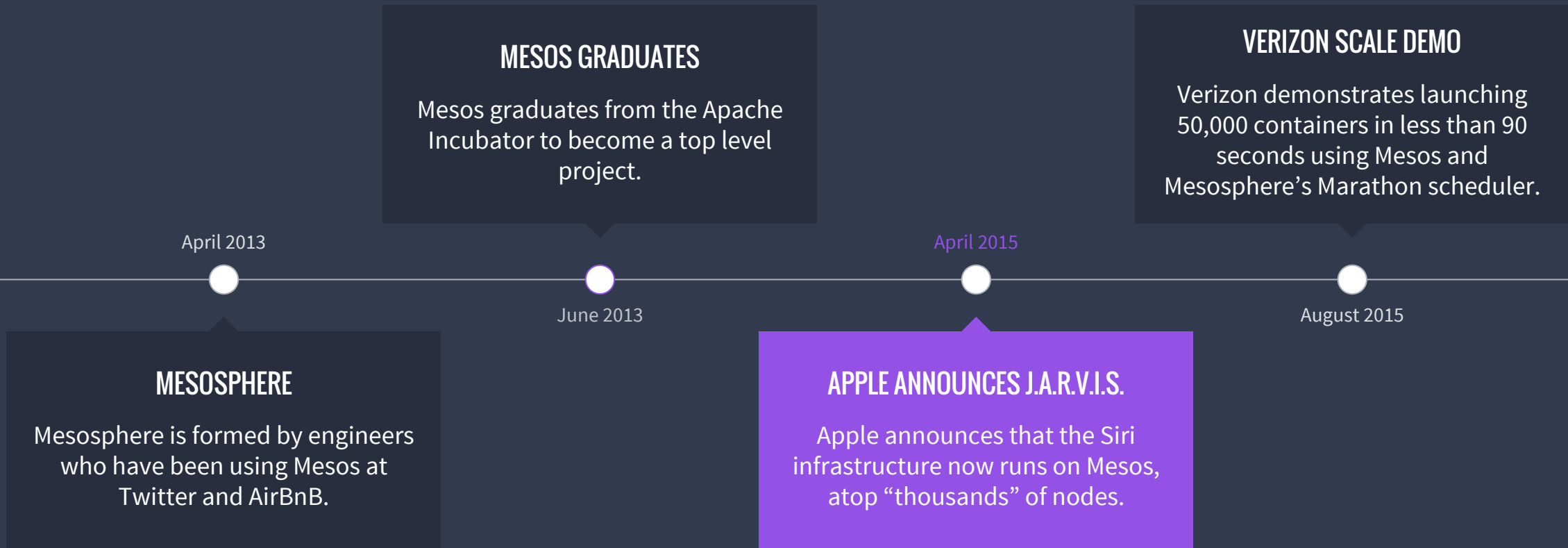
He must address errors and failures **manually**.

He has to deal with dependencies **on a one-off basis** for each of his developers' applications.
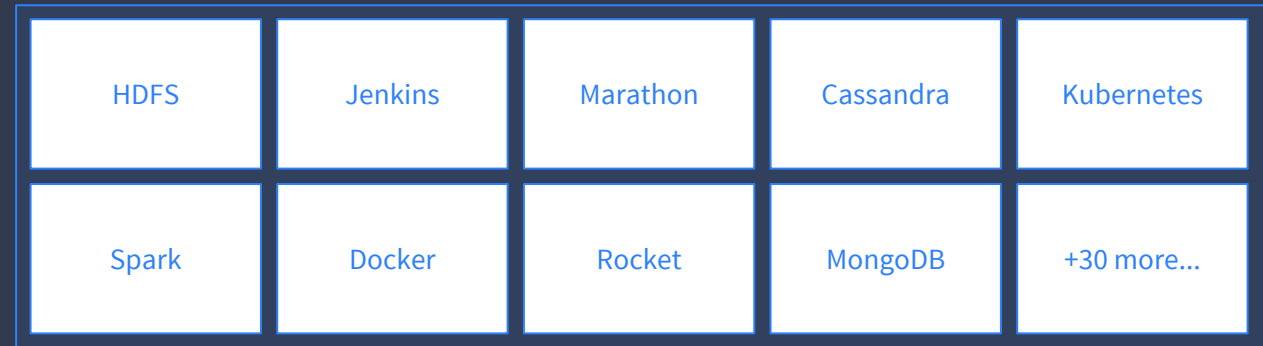
# MESOSPHERE & THE DCOS
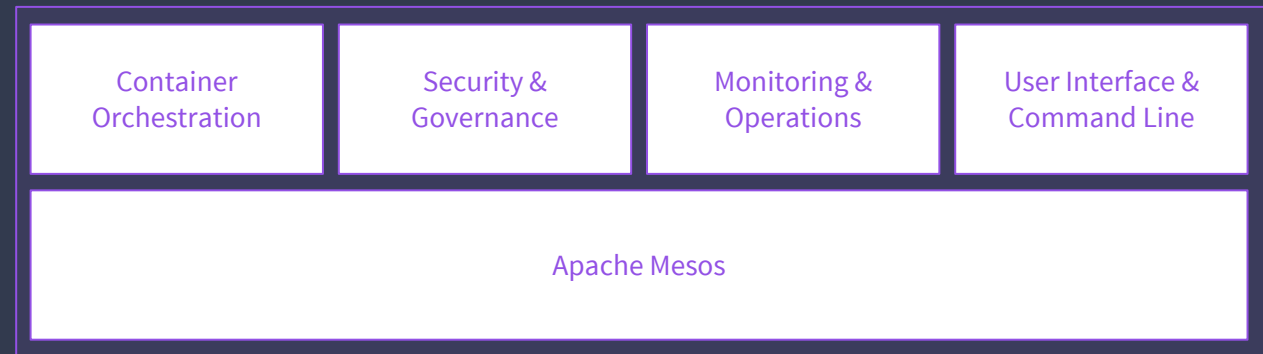
# MISSION TO THE MESOSPHERE

**MESOS GRADUATES**
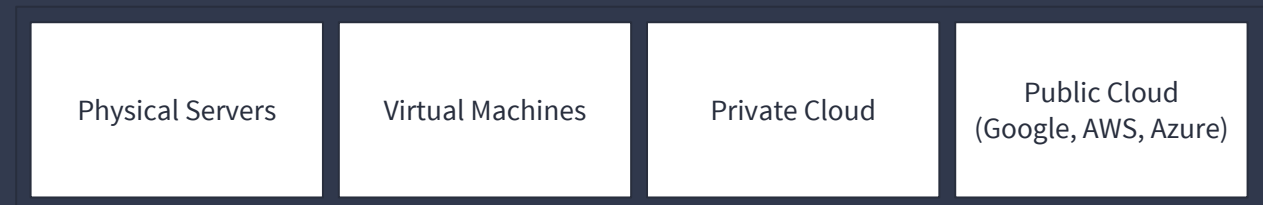
Mesos graduates from the Apache Incubator to become a top level project.

**VERIZON SCALE DEMO**

Verizon demonstrates launching 50,000 containers in less than 90 seconds using Mesos and Mesosphere's Marathon scheduler.

April 2013

April 2015

June 2013

August 2015

**MESOSPHERE**

Mesosphere is formed by engineers who have been using Mesos at Twitter and AirBnB.

**APPLE ANNOUNCES J.A.R.V.I.S.**

Apple announces that the Siri infrastructure now runs on Mesos, atop "thousands" of nodes.

# MESOSPHERE DCOS

## Services & Containers

| | | | | |
|---|---|---|---|---|
| HDFS | Jenkins | Marathon | Cassandra | Kubernetes |
| Spark | Docker | Rocket | MongoDB | +30 more... |

## Mesosphere DCOS

| | | | |
|---|---|---|---|
| Container Orchestration | Security & Governance | Monitoring & Operations | User Interface & Command Line |

Apache Mesos

## Existing Infrastructure

| | | | |
|---|---|---|---|
| Physical Servers | Virtual Machines | Private Cloud | Public Cloud (Google, AWS, Azure) |

# THE DATACENTER OPERATING SYSTEM

DCOS aims to make developing & deploying distributed apps easier.

Short-term:

- Software installation/removal
- Seamless upgrades
- Automatic failure detection, reconciliation

# A UNIFIED INTERFACE TO THE DATACENTER

# THE COMMAND LINE TO THE DATACENTER

# PRODUCTION CUSTOMERS AND MESOS USERS

# WHAT WILL IT TAKE TO MAKE DAN HAPPY?

# CONTAINERS EVERYWHERE

Many Mesos tasks run in **containers**:

- Mesos containerizer
- Docker
- Universal containerizer (in progress)

# CONTAINERS EVERYWHERE

Many Mesos tasks run in **containers**:

- Mesos containerizer
- Docker
- Universal containerizer (in progress)

Containers use standard linux features to create an isolated execution environment:

- kernel namespaces
  - process isolation
- control groups (cgroups)
  - resource isolation
- chroot
  - filesystem isolation
- seccomp
  - restricted kernel access
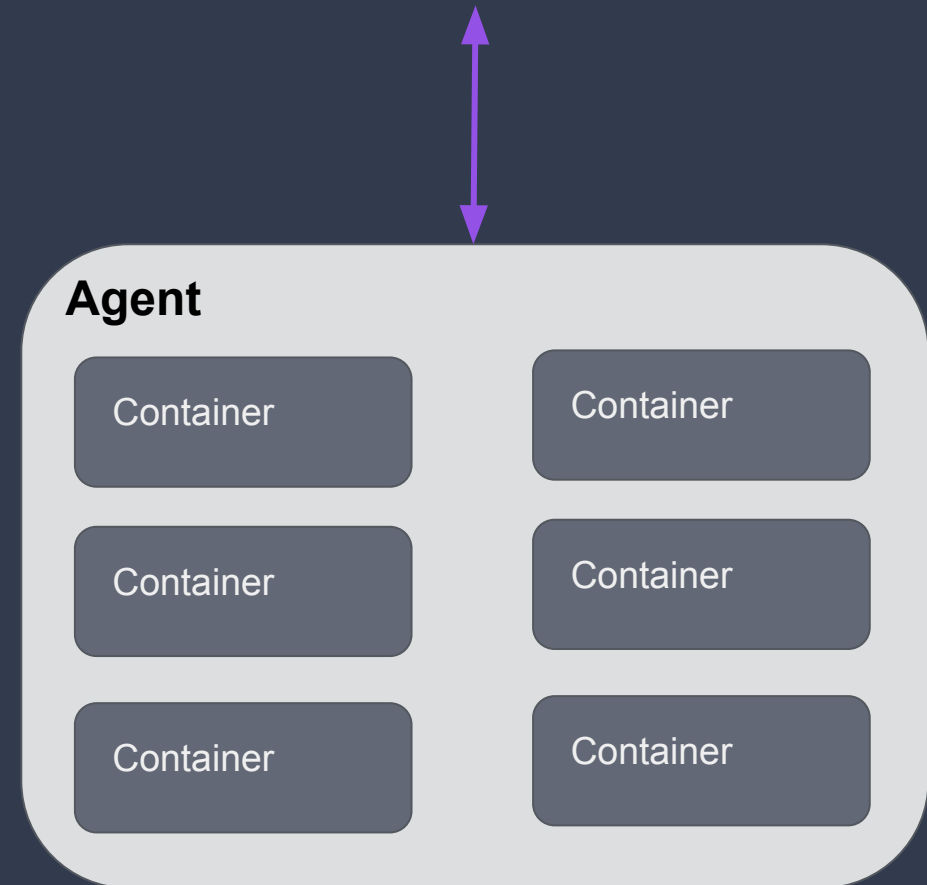
# CONTAINERS EVERYWHERE

Containers also help Dan solve his dependency
problem by giving tasks everything they need to run.

# CONTAINER NETWORKING

Containers isolate tasks on the agent, but what about their communication?
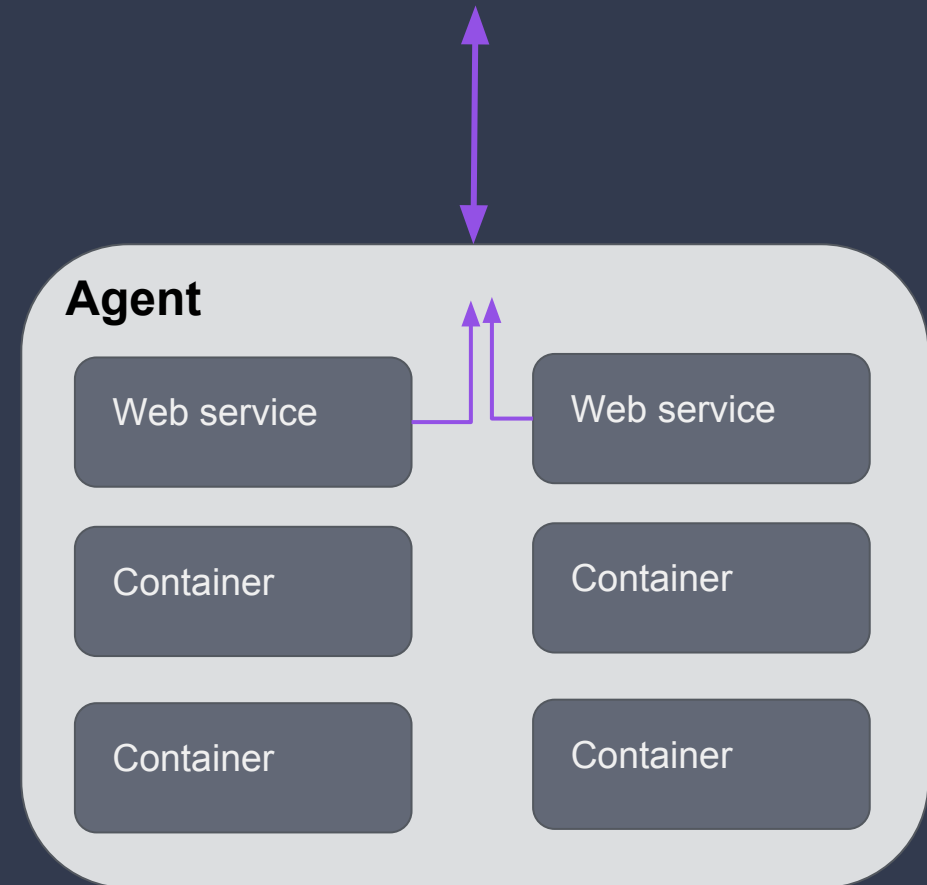
The status quo in a Mesos cluster: one IP per agent.

Many containers per agent: they must share a single IP.

**Agent**

| Container | Container |
|-----------|-----------|
| Container | Container |
| Container | Container |

# CONTAINER NETWORKING

This causes headaches:

- **Port conflicts**
- Security compromises
- Performance
- Service discovery

**Agent**

| | |
|---|---|
| Web service | Web service |
| Container | Container |
| Container | Container |

# CONTAINER NETWORKING

This causes headaches:

- Port conflicts
- Security compromises
- Performance
- Service discovery

**Agent**

| Prod. service | Test service |
|---|---|
| Container | Container |
| Container | Container |

# CONTAINER NETWORKING

This causes headaches:

- Port conflicts
- Security compromises
- Performance
- Service discovery

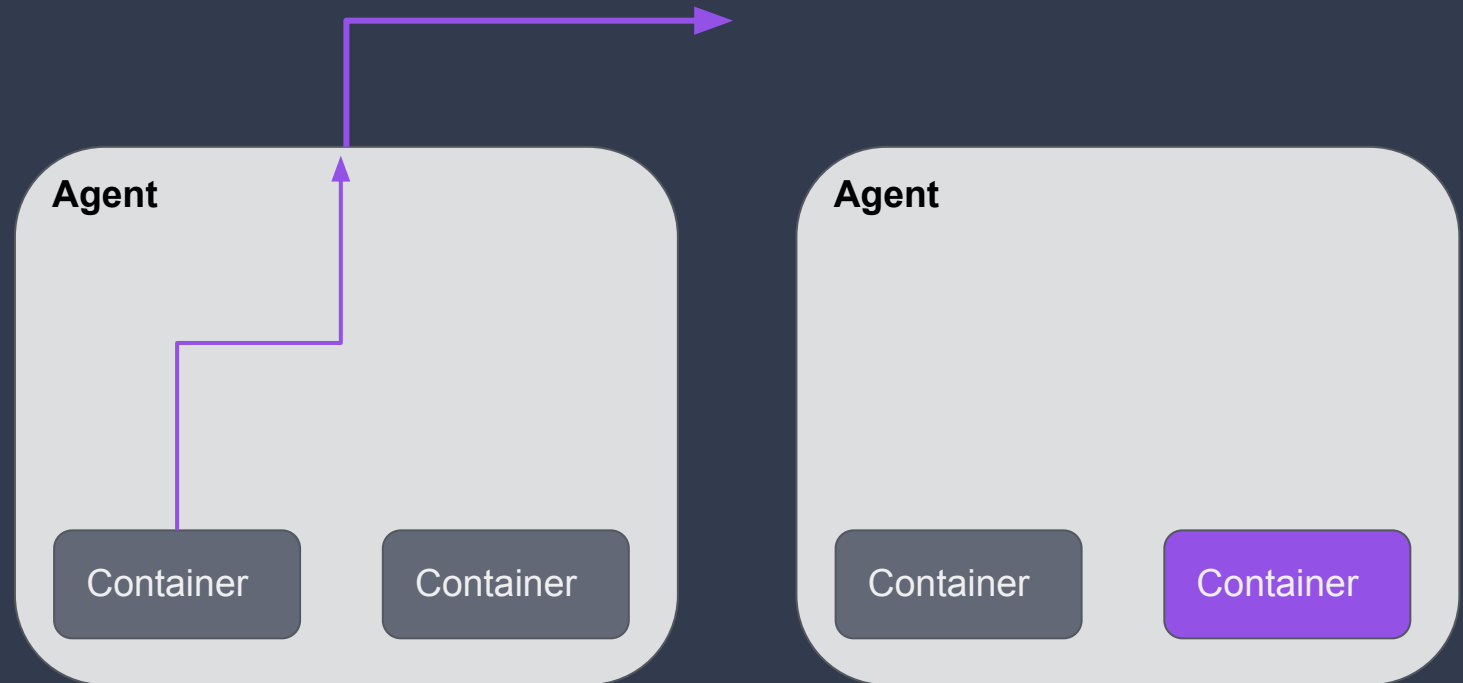# CONTAINER NETWORKING

This causes headaches:

- Port conflicts
- Security compromises
- Performance
- Service discovery

**Agent**

Container    Container

**Agent**

Container    Container

# NETWORK ISOLATION

Segregating containers' network traffic can solve these problems in an elegant, maintainable way.

# NETWORK ISOLATION

Segregating containers' network traffic can solve these problems in an elegant, maintainable way.

Implemented as Mesos modules:
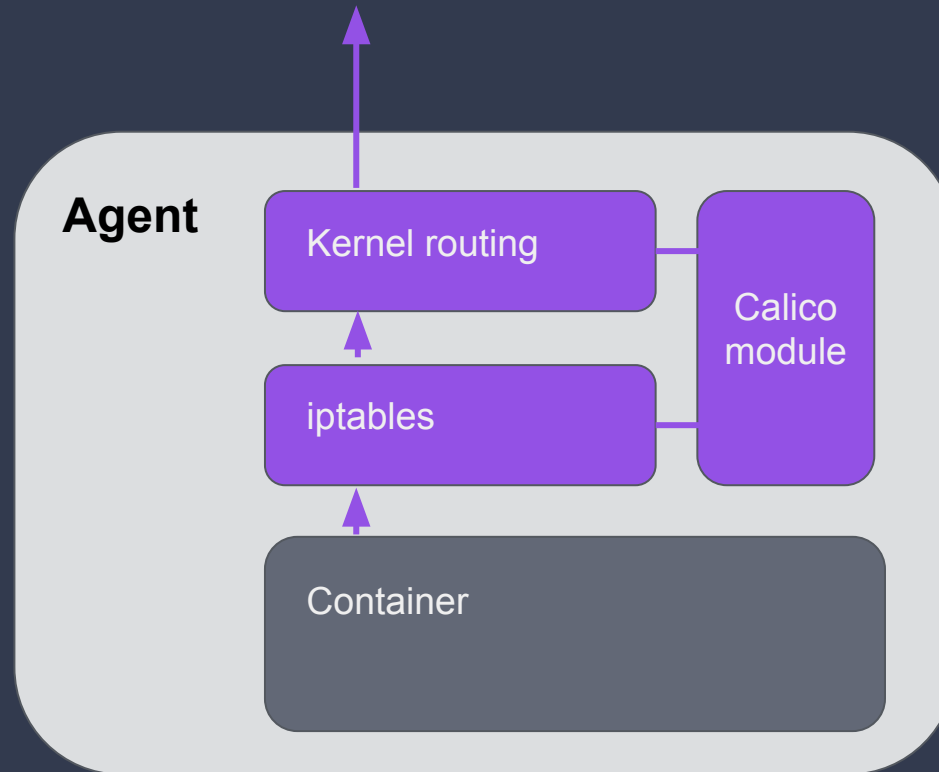
- Project Calico
- Port-mapping isolation
- …

# CALICO NETWORK ISOLATION

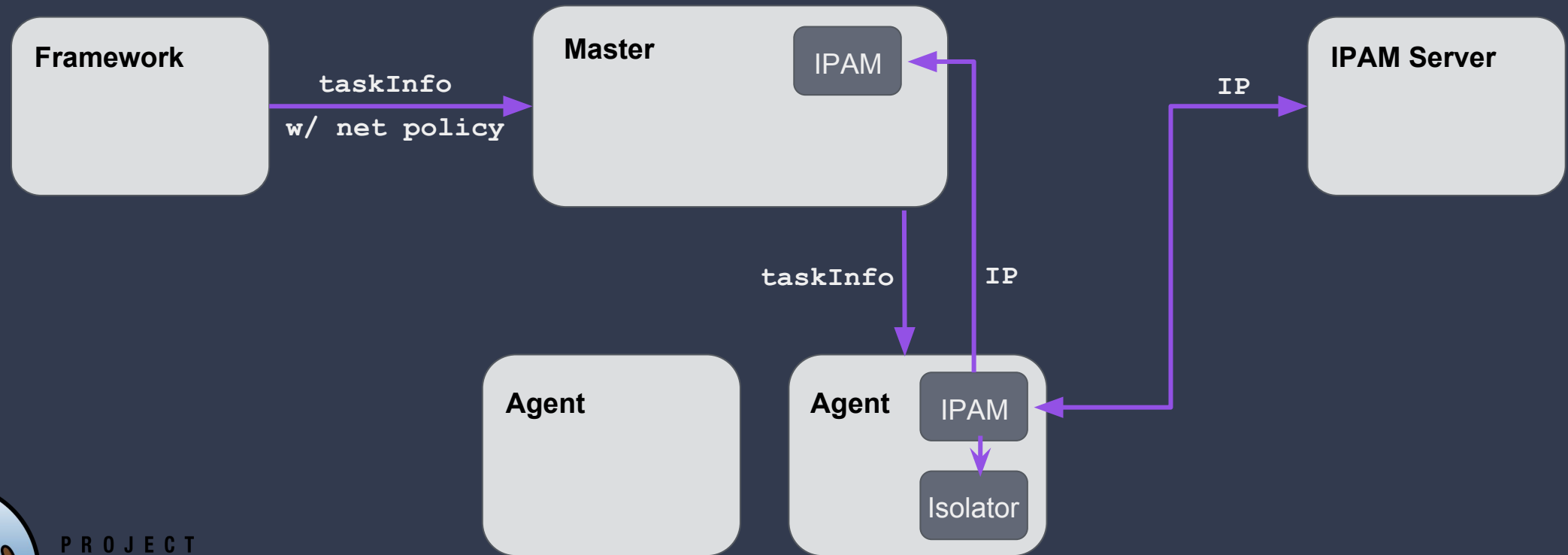Calico Network Virtualizer & IP Address Manager:

- Pure Layer-3 solution
- Uses linux features to route container traffic
- Provides security policies
- Advertises routes to local containers via BGP
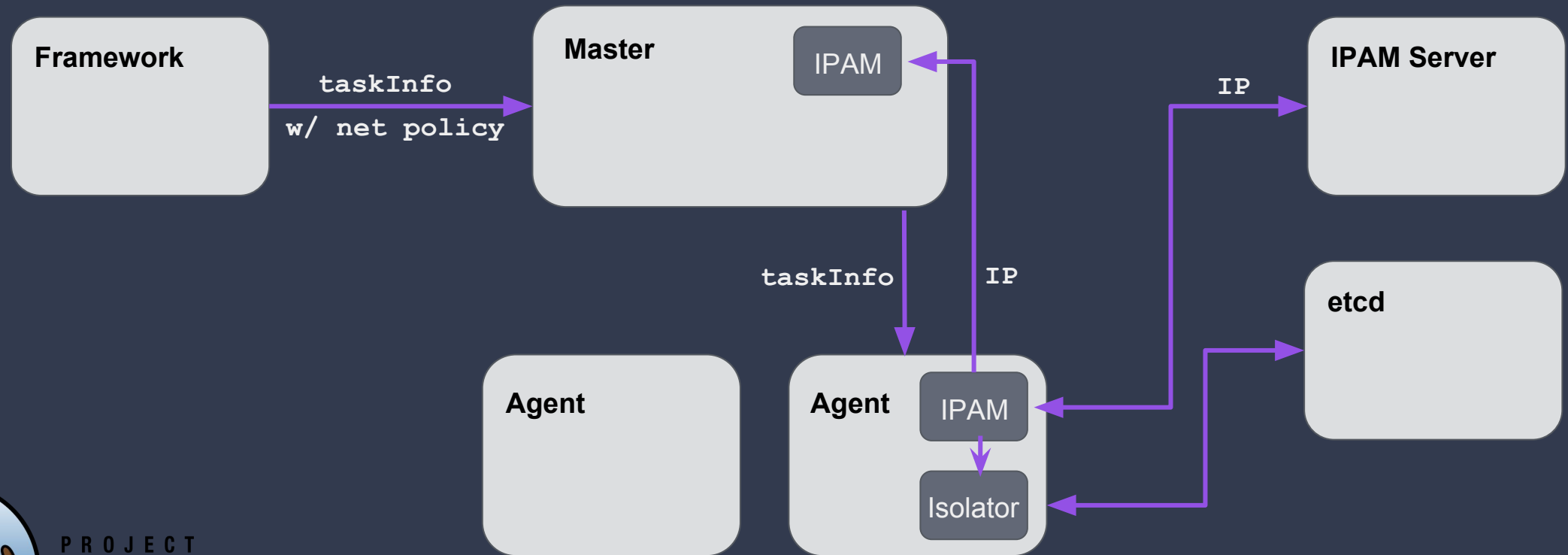- Can assign IP-per-container

# CALICO NETWORK ISOLATION

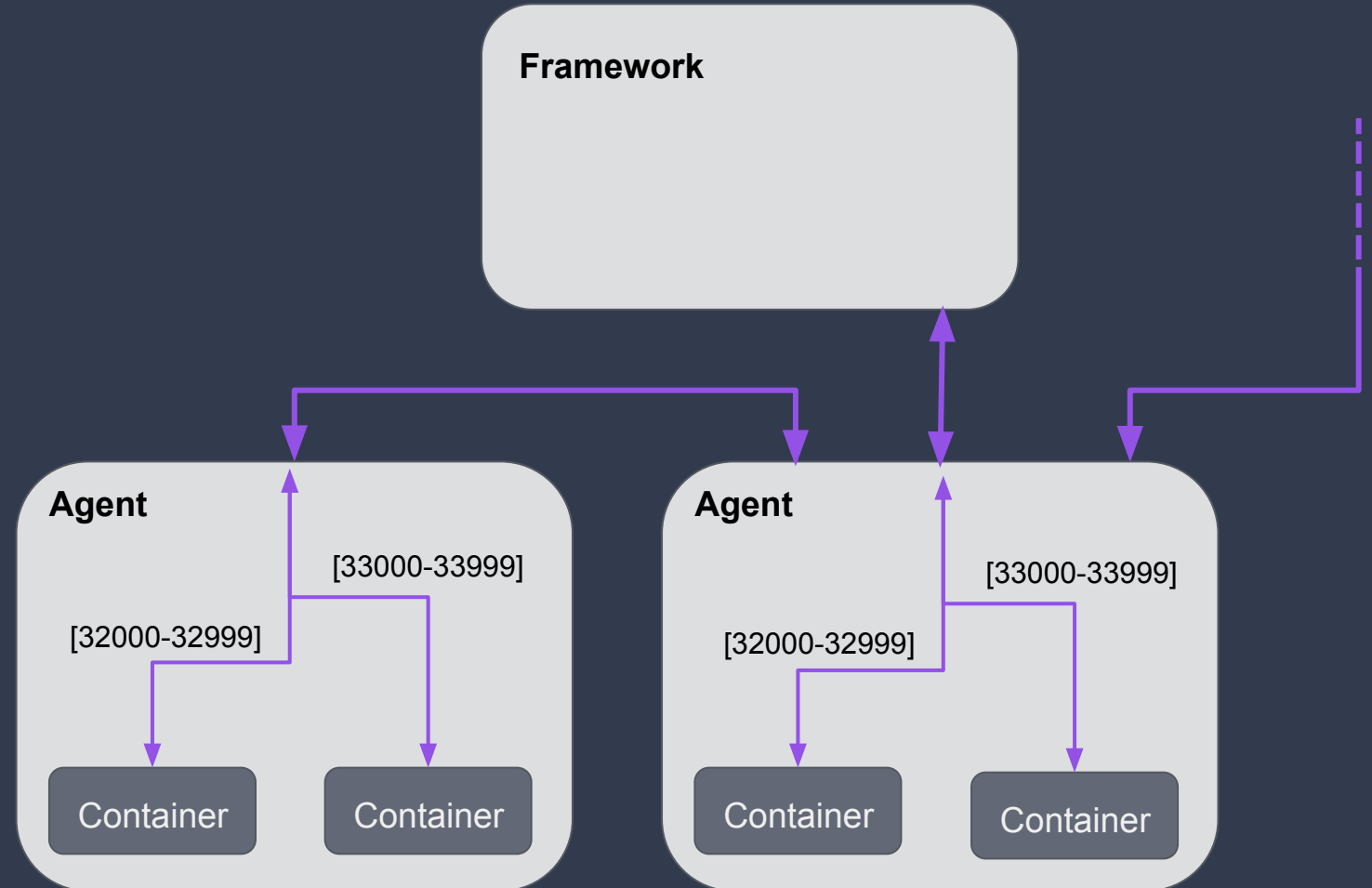# CALICO NETWORK ISOLATION

# CALICO NETWORK ISOLATION

# NETWORK ISOLATION

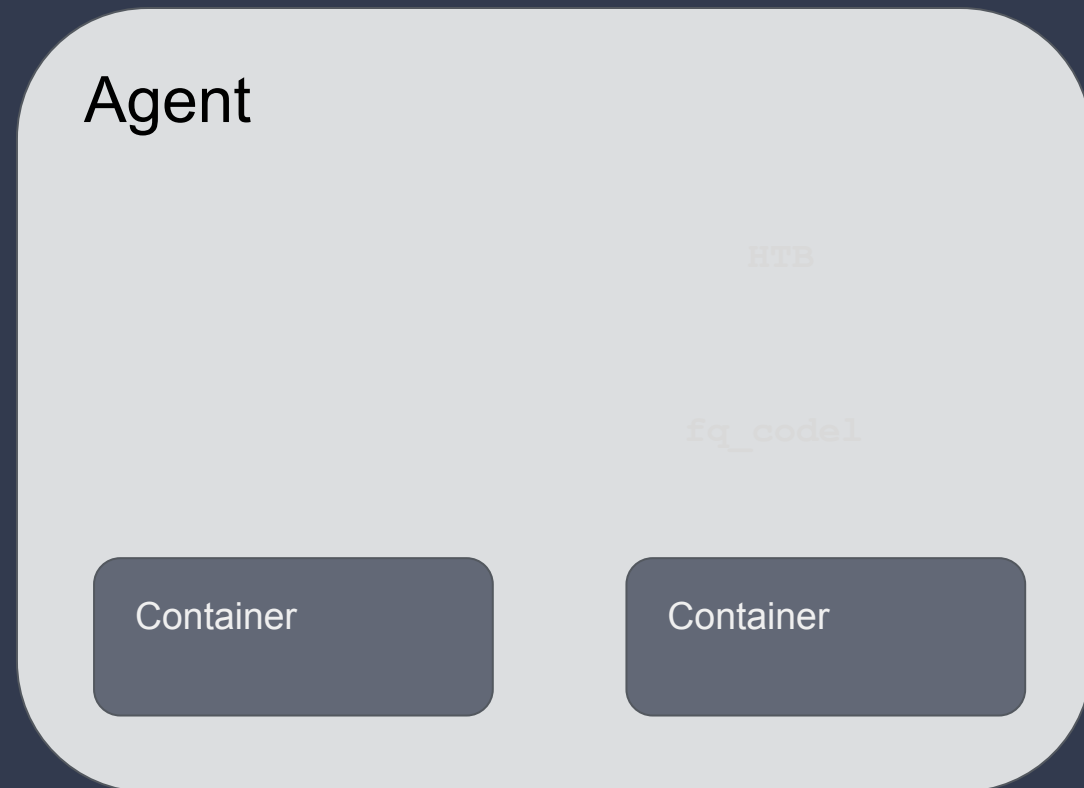What if I don't have enough IPs to go around?

# PORT-MAPPING ISOLATOR

- Ports distributed amongst containers on each agent

- Network traffic routed by port using TC rules

- Implemented with `libnl` (via netlink messages)

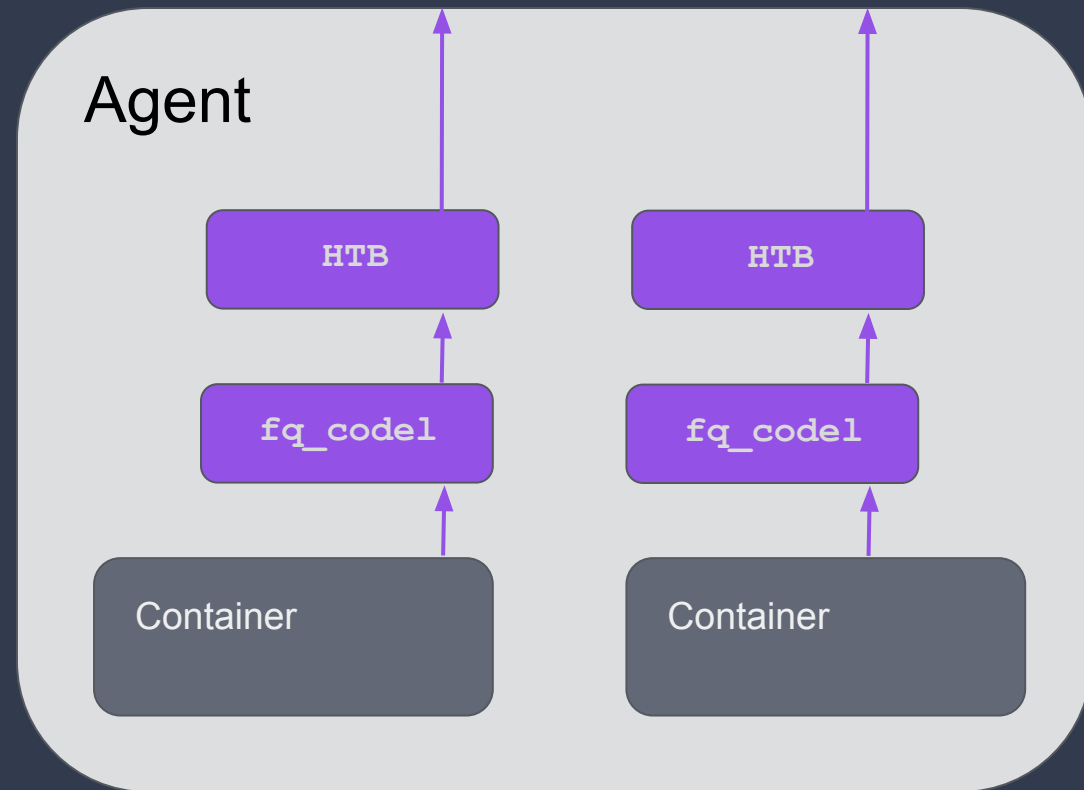- Ports assigned and tracked via scheduler (ex: Aurora)

**Framework**

**Agent**

[33000-33999]

[32000-32999]

Container

Container

**Agent**

[33000-33999]

[32000-32999]

Container

Container

# PORT-MAPPING ISOLATOR

What about performance?

Agent

HTB

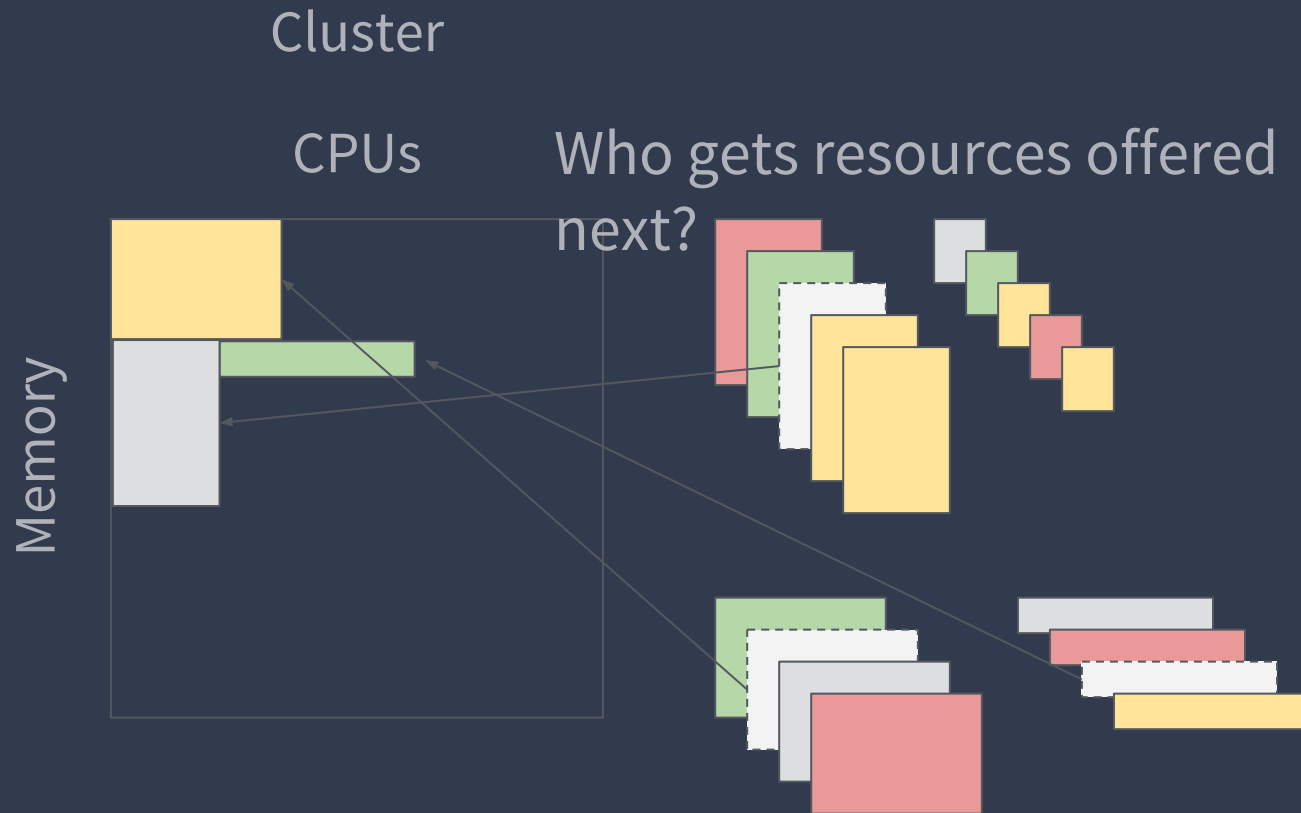fq_codel

Container

Container

# PORT-MAPPING ISOLATOR

- `fq_codel` defines discrete network flows for containers

- Separate flows prevent buffer bloat

- Hierarchical token bucket (HTB) employed to limit bandwidth

# WORKLOADS CHANGED SINCE 2009

# FAIRNESS FOR MULTI-DIMENSIONAL RESOURCES?

Cluster

CPUs

Memory

Who gets resources offered next?

While we can do fitting, how do we express fairness over different units and dimensions?

# DOMINANT RESOURCE FAIRNESS

**Algorithm 1** DRF pseudo-code

$R = \langle r_1, \cdots, r_m \rangle$ ▷ total resource capacities

$C = \langle c_1, \cdots, c_m \rangle$ ▷ consumed resources, initially 0

$s_i \ (i = 1..n)$ ▷ user $i$'s dominant shares, initially 0

$U_i = \langle u_{i,1}, \cdots, u_{i,m} \rangle \ (i = 1..n)$ ▷ resources given to user $i$, initially 0

**pick** user $i$ with lowest dominant share $s_i$

$D_i \leftarrow$ demand of user $i$'s next task

**if** $C + D_i \leq R$ **then**
    $C = C + D_i$ ▷ update consumed vector
    $U_i = U_i + D_i$ ▷ update $i$'s allocation vector
    $s_i = \max_{j=1}^m \{u_{i,j}/r_j\}$
**else**
    **return** ▷ the cluster is full
**end if**

# TIME DIMENSIONALITY HAS CHANGED!

Multitenancy now expands domains of
multiple batch schedulers with a mix of:

- Long lived services
- Storage services
- Short lived analytics tasks

Is extreme fairness what you really want?

# SEVERAL POs AT CUSTOMER SITES

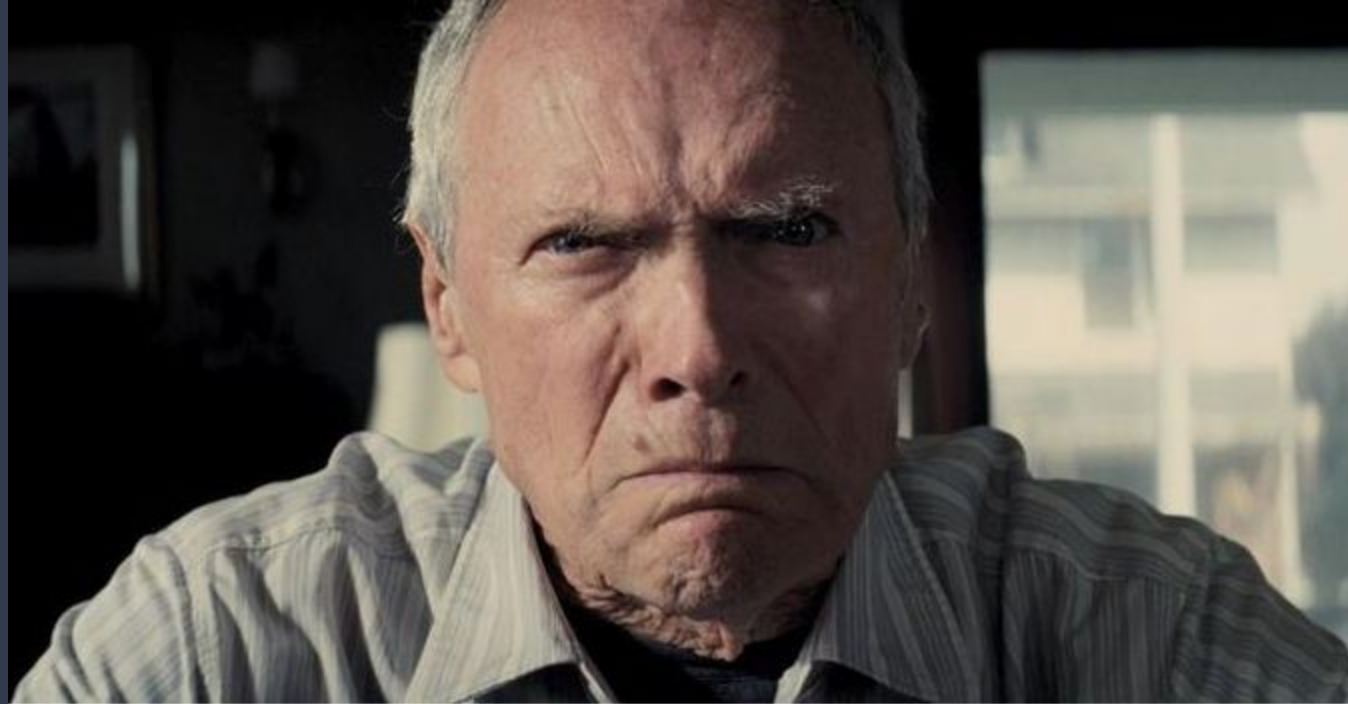"My framework is starved! Why isn't my framework receiving any resources?"

- Some frameworks has a lot of work to do, others less. All gets a fair share by default.

- Configuration is hard with weights, static reservations, etc

# EASY TO MAKE MISTAKES IN SCHEDULER IMPLEMENTATIONS

```
offer
resourceOffers(offers) {
  ...
  if (is_ok(offer)) {
    launchTasks(offer);
  }
}
```

```
resourceOffers(offers) {
  ...
  if (is_ok(offer)) {
    launchTasks(offer);
  } else {
    declineOffer(offer);
  }
}
```

# POs MAKE DAN UNHAPPY

# WHAT DID TWITTER DO?

- Uses Apache Aurora for most of its operations

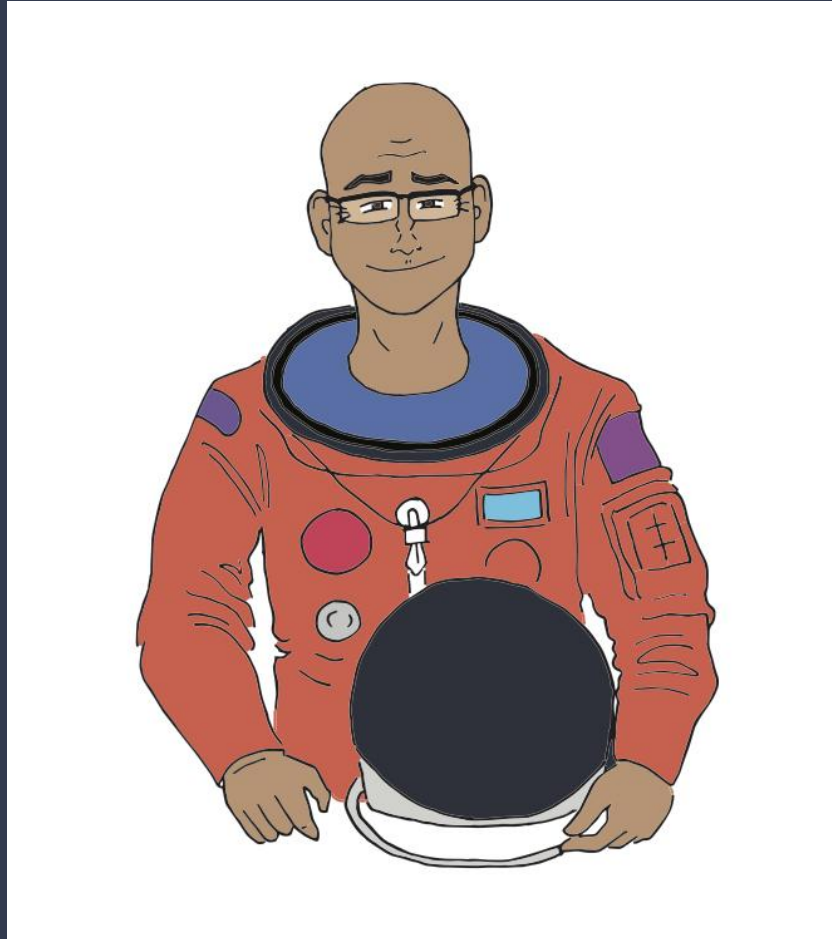- It implemented preemption assuming it was the only scheduler available

# MULTI-TENANCY BECOMES TOO RISKY FOR CRITICAL SYSTEMS

- Companies partition Mesos cluster into many smaller Mesos clusters

- Run multiple copies of the same framework on top of Mesos

- Avoid running multiple frameworks all together

- That was surely not the intent

# IN THE WORKS

- Quotas ensure minimum set of resources for frameworks

- Optimistic offers enables resource parallelism

- Cooperative preemption through Inverse offers
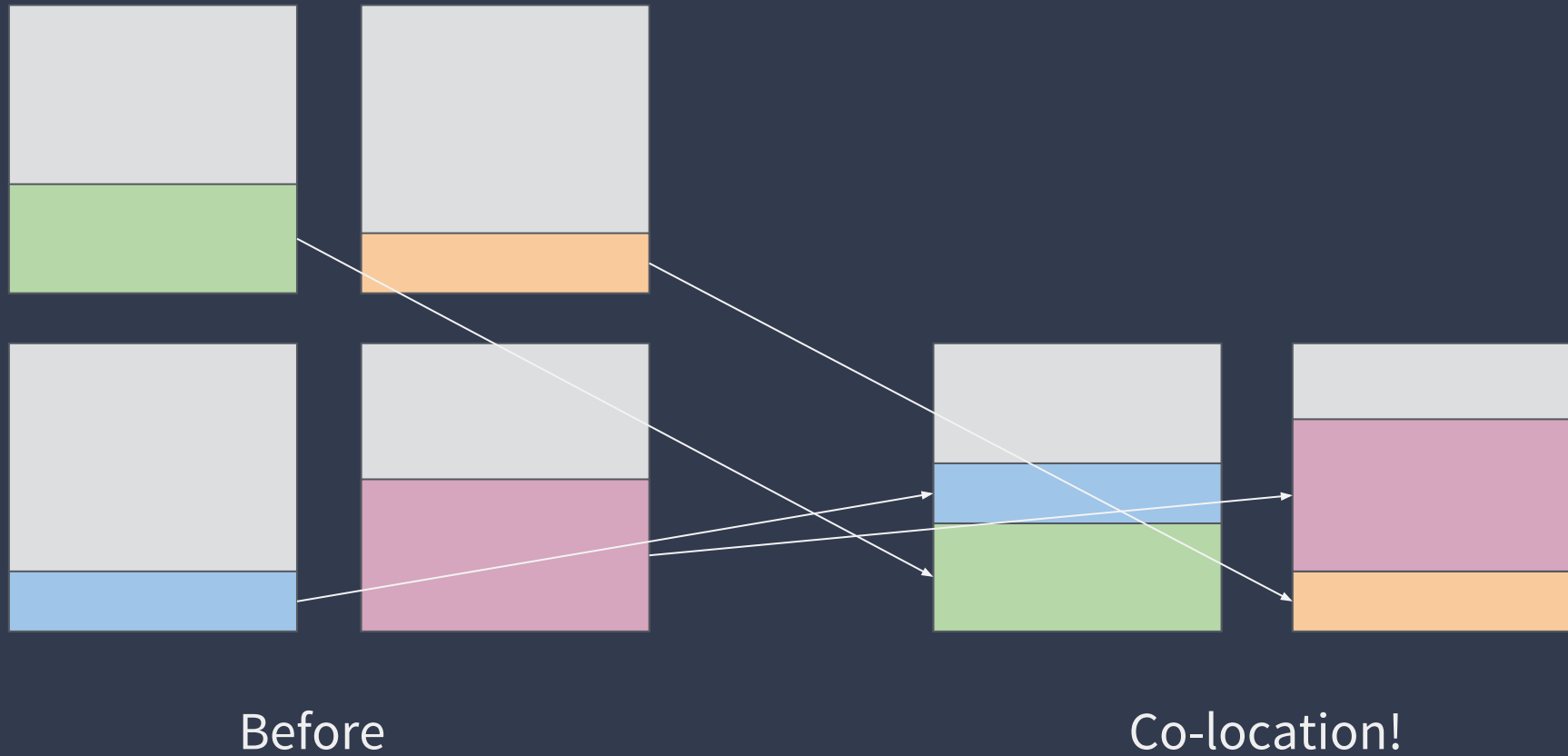
# LET'S ASSUME DAN IS HAPPY WITH HIS CLUSTER

# THAT MAKES HIS BOSS HAPPY

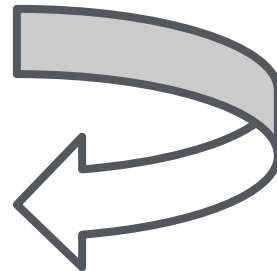# AND THEIR CFO IS HAPPY TOO

# MESOS HELPS REDUCE WASTED RESOURCES



Before

Co-location!

# ESTIMATING RESOURCES IS HARD

# MESOS ENABLES MULTIPLE SCHEDULER ALGORITHMS

**User**

**Scheduler + Allocation**

"Please run container X on Y resources"

# MESOS ENABLES MULTIPLE SCHEDULER ALGORITHMS

User
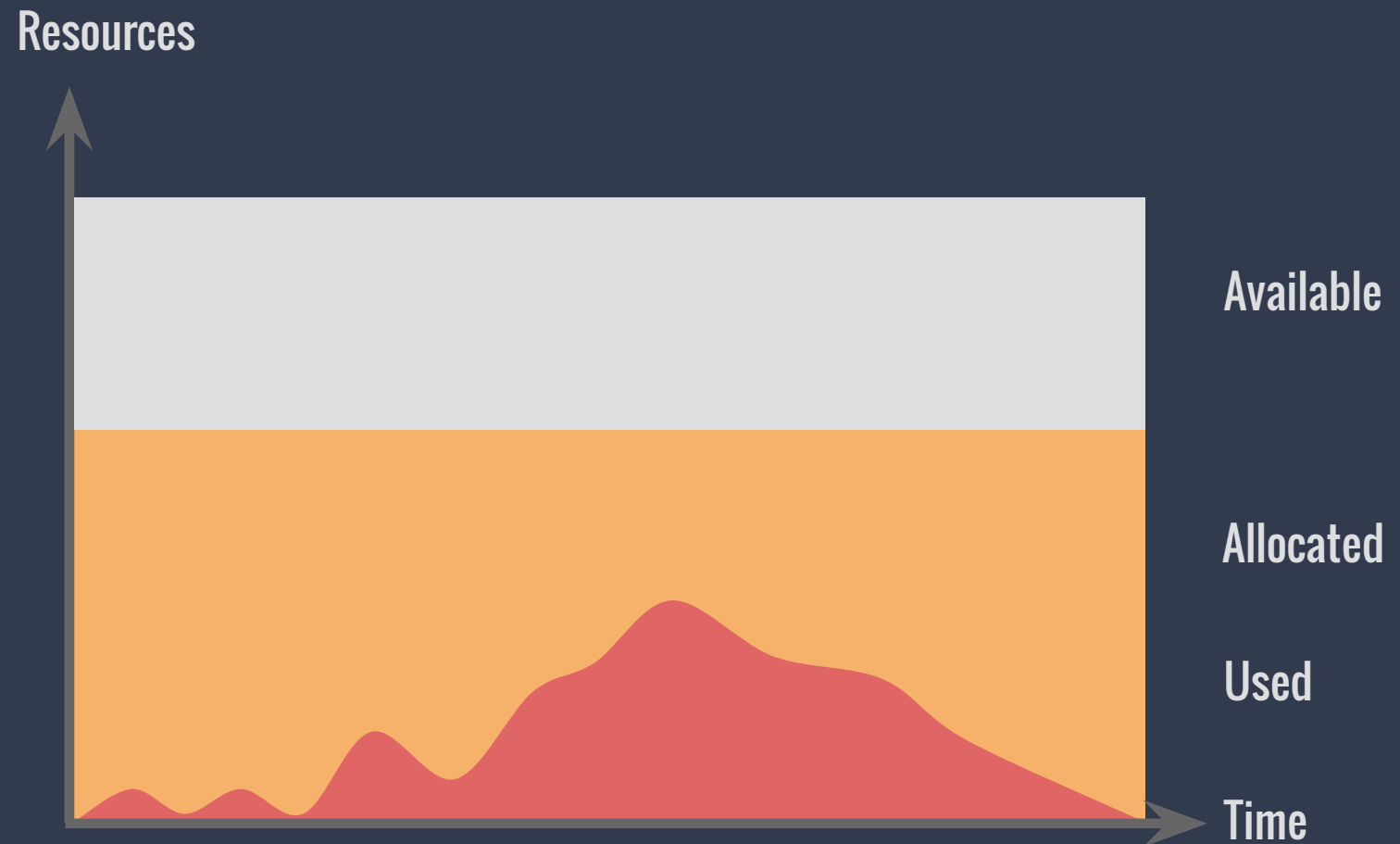
Scheduler(s)

Allocation

Work

Resource Offers

Launch Tasks

# RESOURCES REPRESENT ALLOCATION

How are users supposed to know how many resources their workload requires?
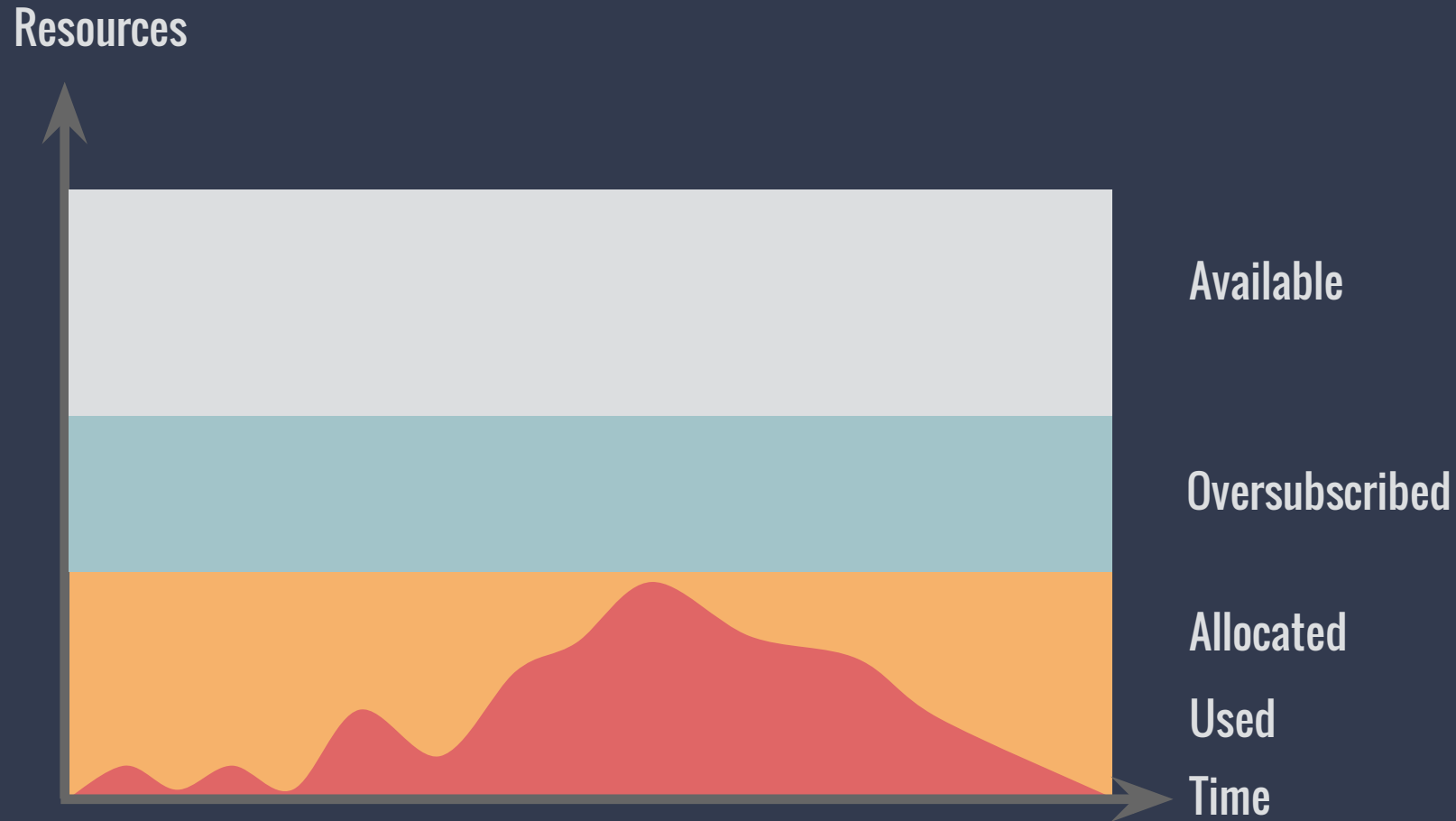
# RESOURCES REPRESENT ALLOCATION

CPUs

100%

Unallocated

90%

Allocated

Usage slack

15%

Used

Memory

# USAGE SLACK HURTS UTILISATION



Resources

Available

Allocated

Used

Time

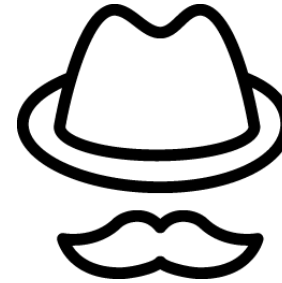# FIRST STEPS TOWARDS IMPROVED UTILISATION

# OVERSUBSCRIPTION ENABLES TASKS TO RUN ON SLACK

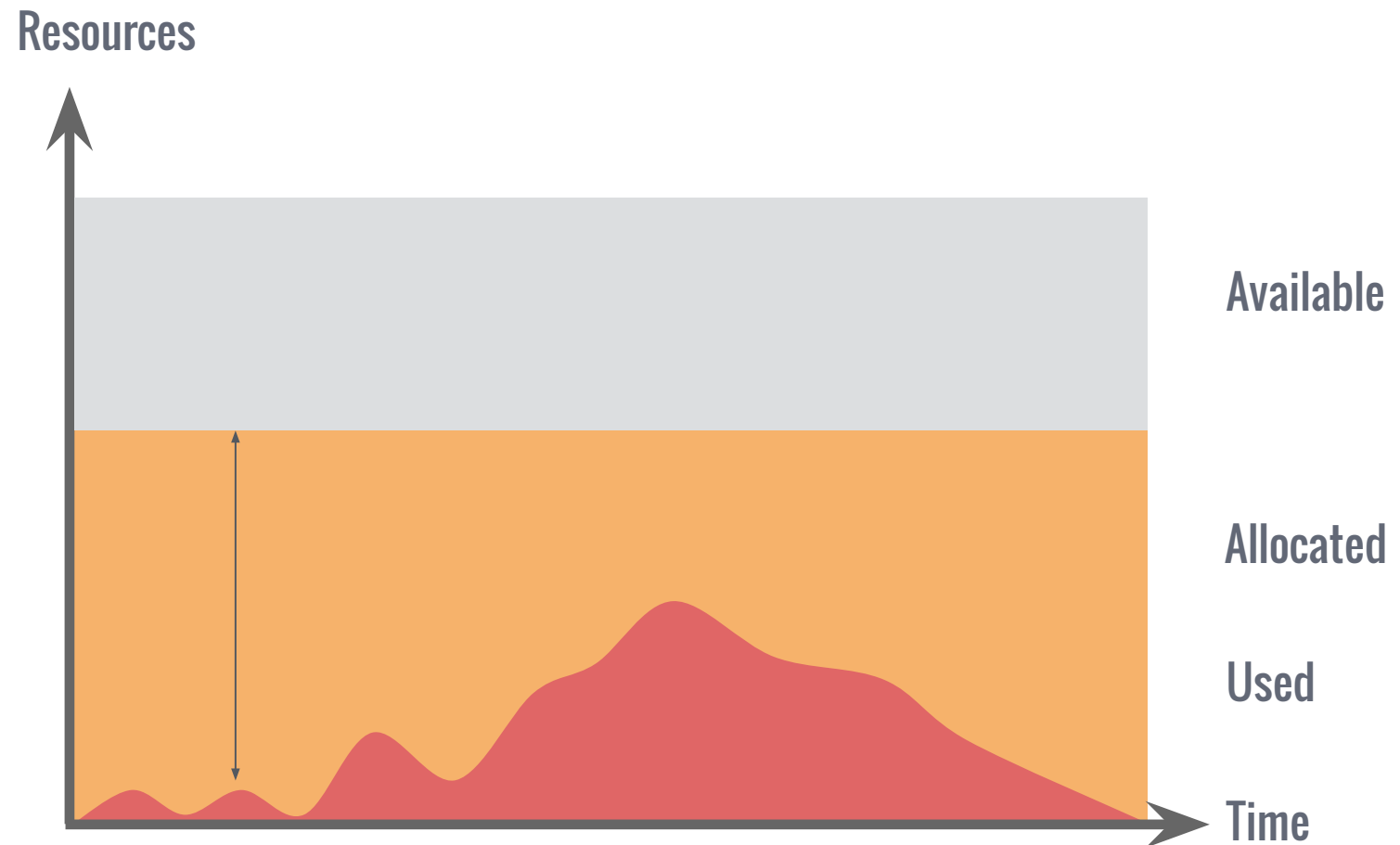# TWO COMPONENTS ENABLE OVERSUBSCRIPTION

**Resource Estimator**

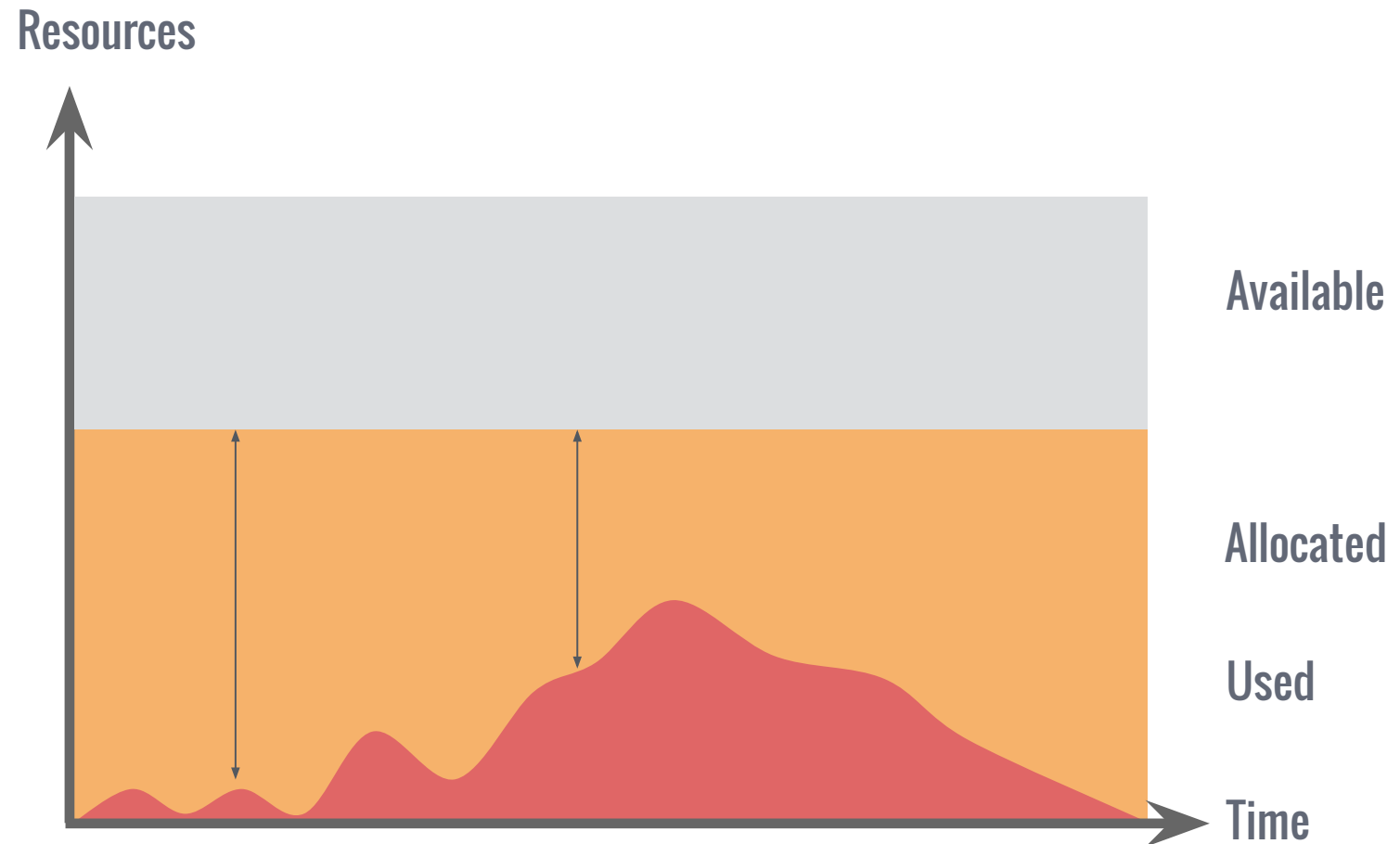**Quality of Service Controller**

# ESTIMATING OVERSUBSCRIBABLE RESOURCES

Resources

Available

Allocated

Used

Time

How many resources should be oversubscribed?

# WHAT DO WE DO ABOUT MISPREDICTIONS?

Now, what happens
when things change?

Resources

Available

Allocated

Used

Time

# THE QoS CONTROLLER

- Can shut down best effort containers
- In the future, it will be able to correct by
  - Freezing
  - Throttling
  - Resizing
  - Cooperating with the framework

# MANY RESOURCES CANNOT BE ISOLATED

- Logical units on the chip
- Last level caches
- Memory bandwidth
- I/O
- Chip power supply

# OVERSUBSCRIPTION WITH INTEL: SERENITY

## https://goo.gl/jWtu7V

# WRAPPING UP

- Mesos is being used in production at huge scale

- It forms the core of an operating system for the datacenter

- Lots of exciting work yet to do!

Slides at http://mesosphere.github.io/presentations

(P.S., we're currently hiring interns!)