# Mesos
# for Spark Users

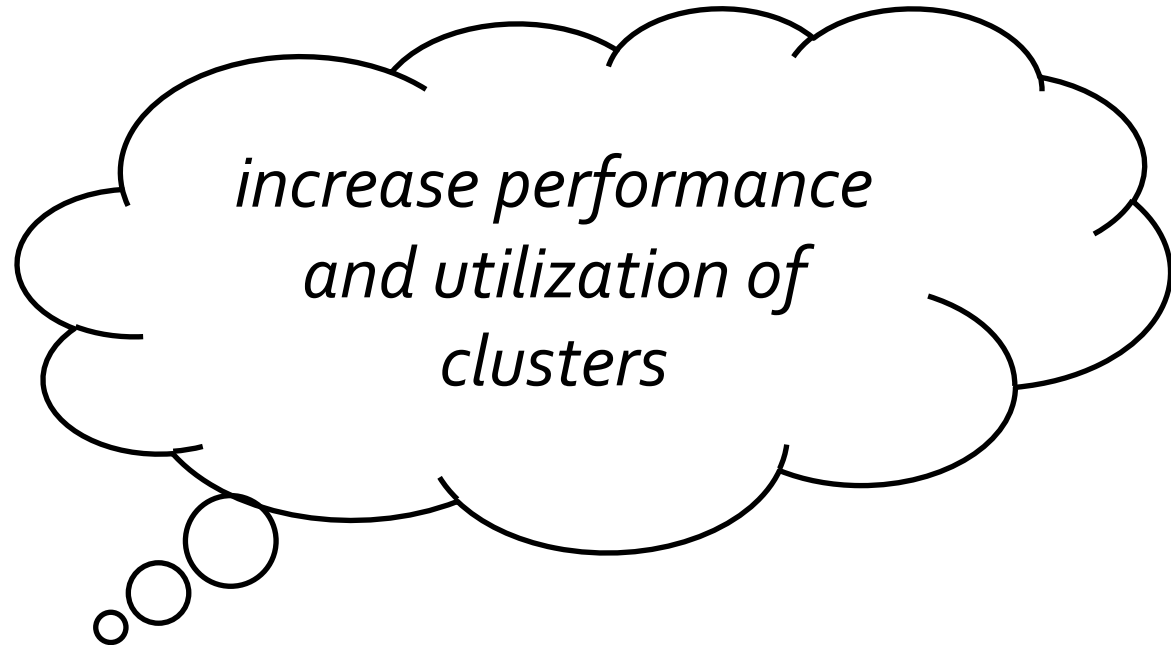mesos.apache.org

@ApacheMesos

Benjamin Hindman – @benh

# agenda

① Mesos ⟵

② Spark on Mesos

③ why Mesos?
   ① multi-tenancy
   ② fine-grained sharing
   ③ why not?

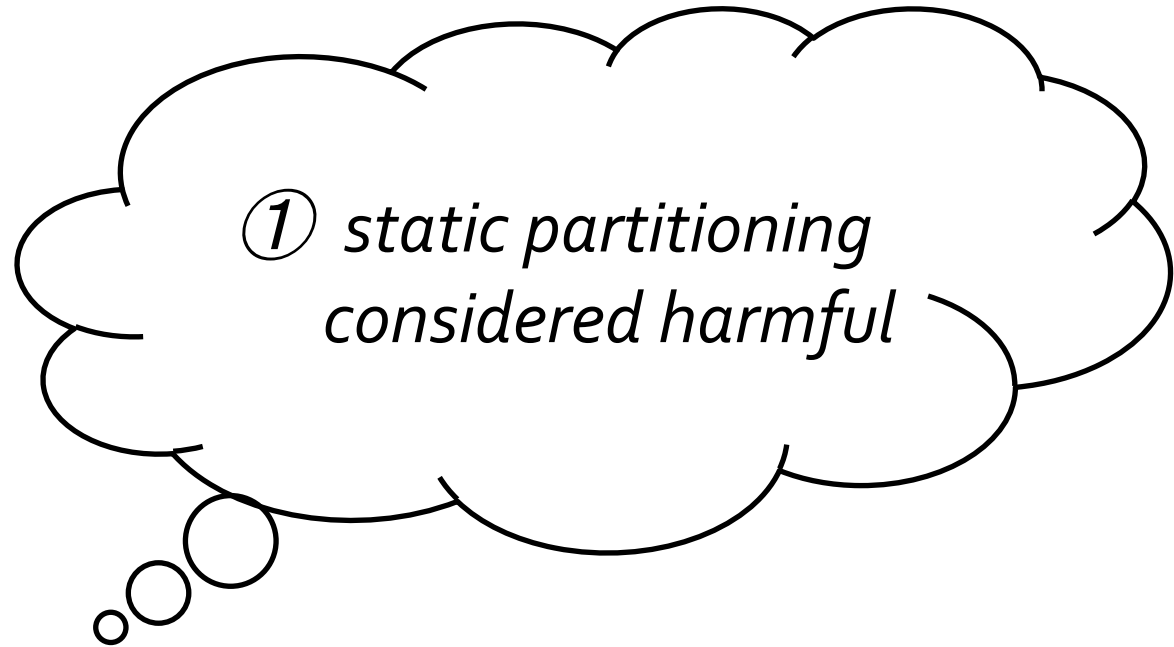④ long-lived services and other frameworks

# a little history

Mesos started as a research project at Berkeley in early 2009 by Benjamin Hindman, Andy Konwinski, Matei Zaharia, Ali Ghodsi, Anthony D. Joseph, Randy Katz, Scott Shenker, Ion Stoica
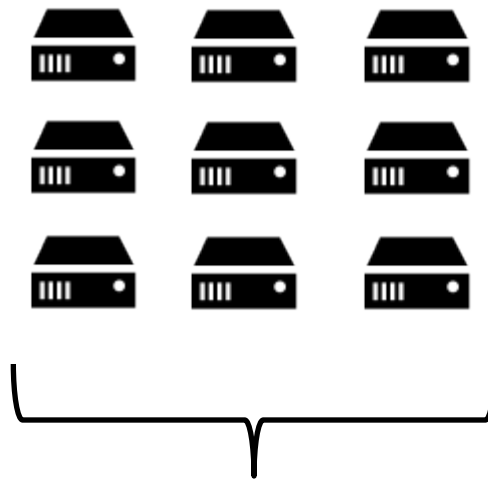
# our motivation



*increase performance and utilization of clusters*

# our intuition

① *static partitioning considered harmful*

# static partitioning considered harmful

datacenter

# static partitioning considered harmful

# static partitioning considered harmful

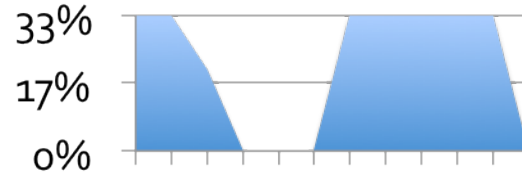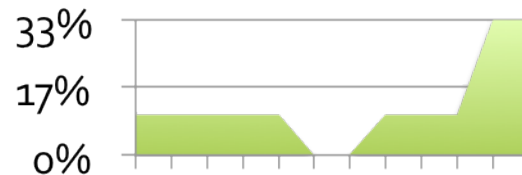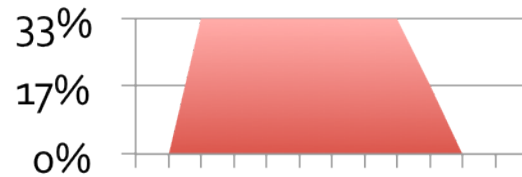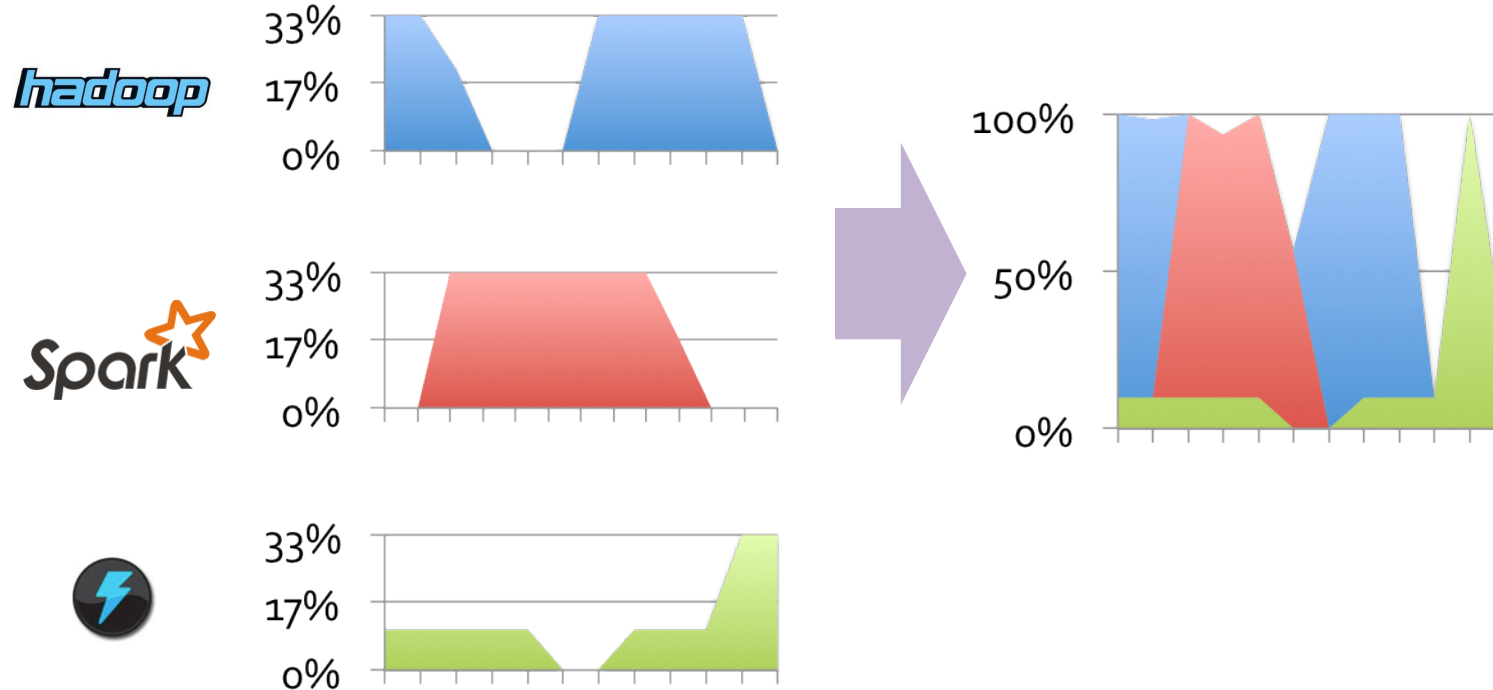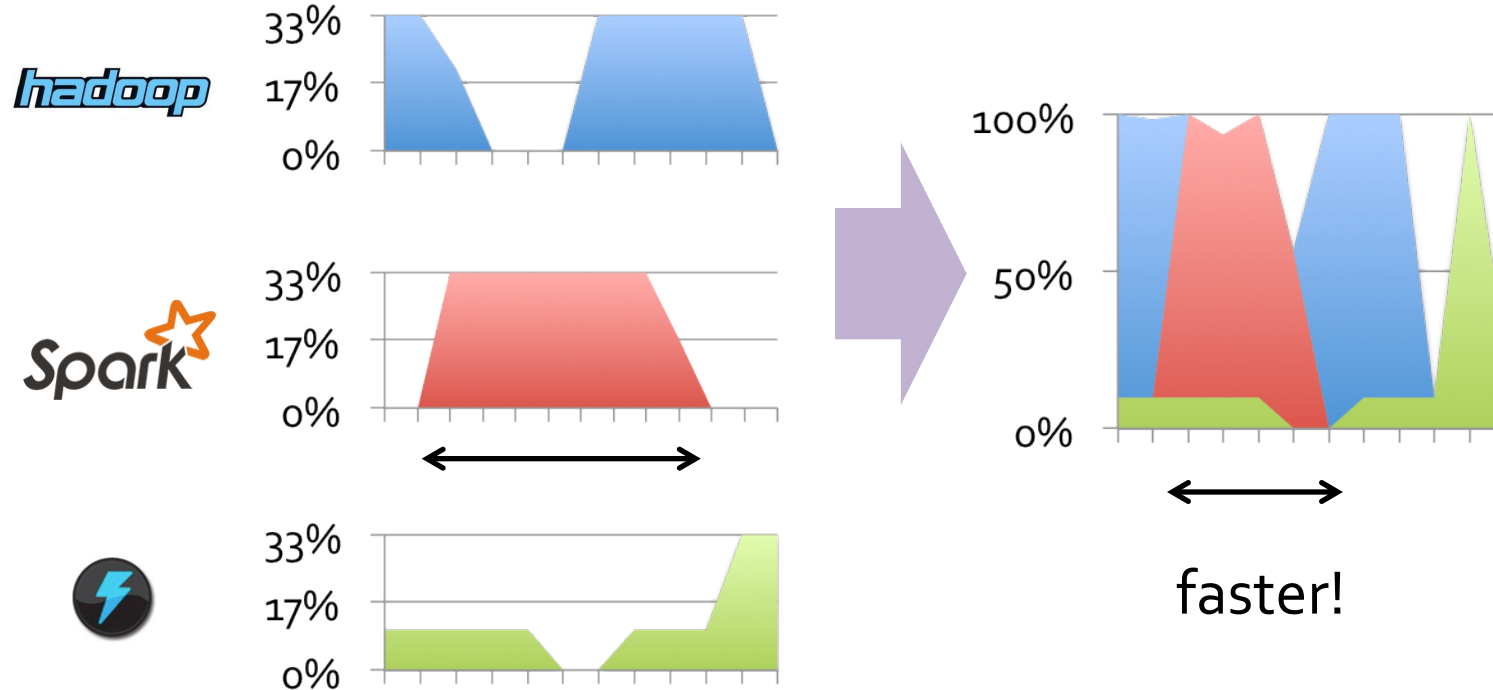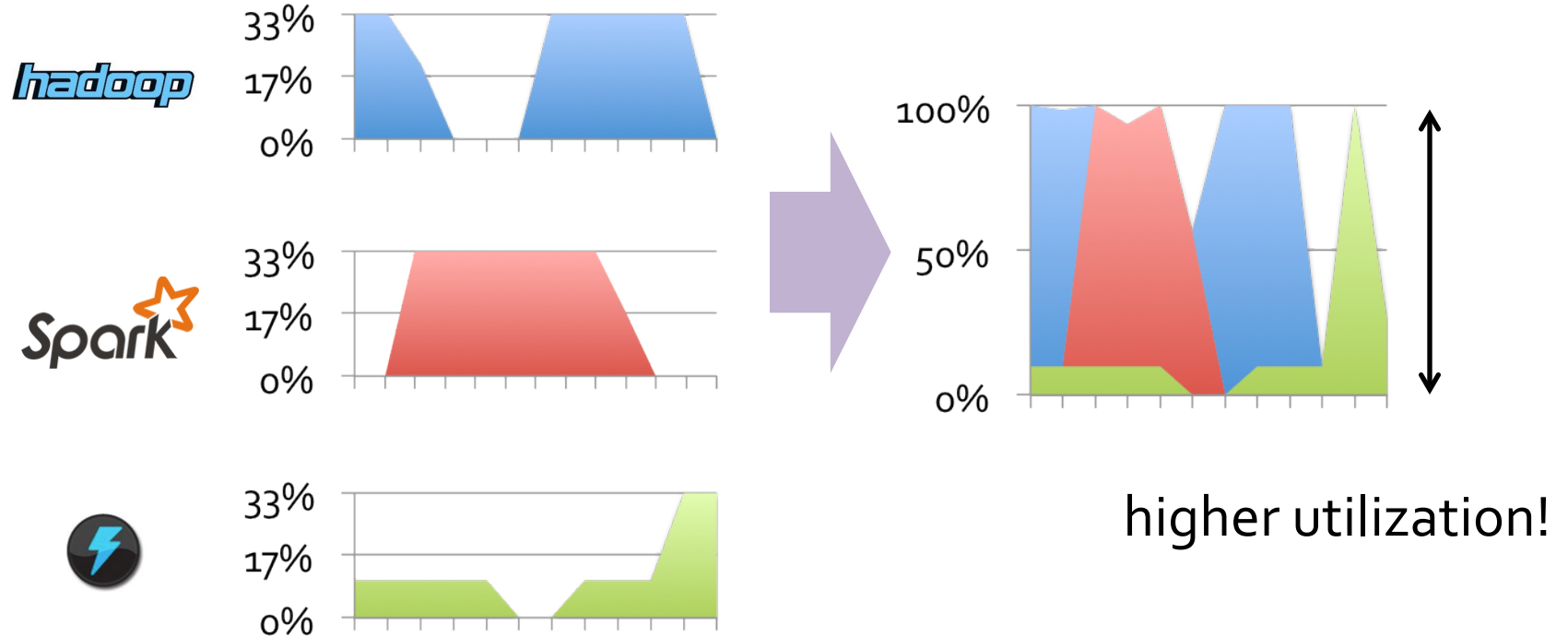# static partitioning considered harmful

# static partitioning considered harmful



faster!

# static partitioning considered harmful
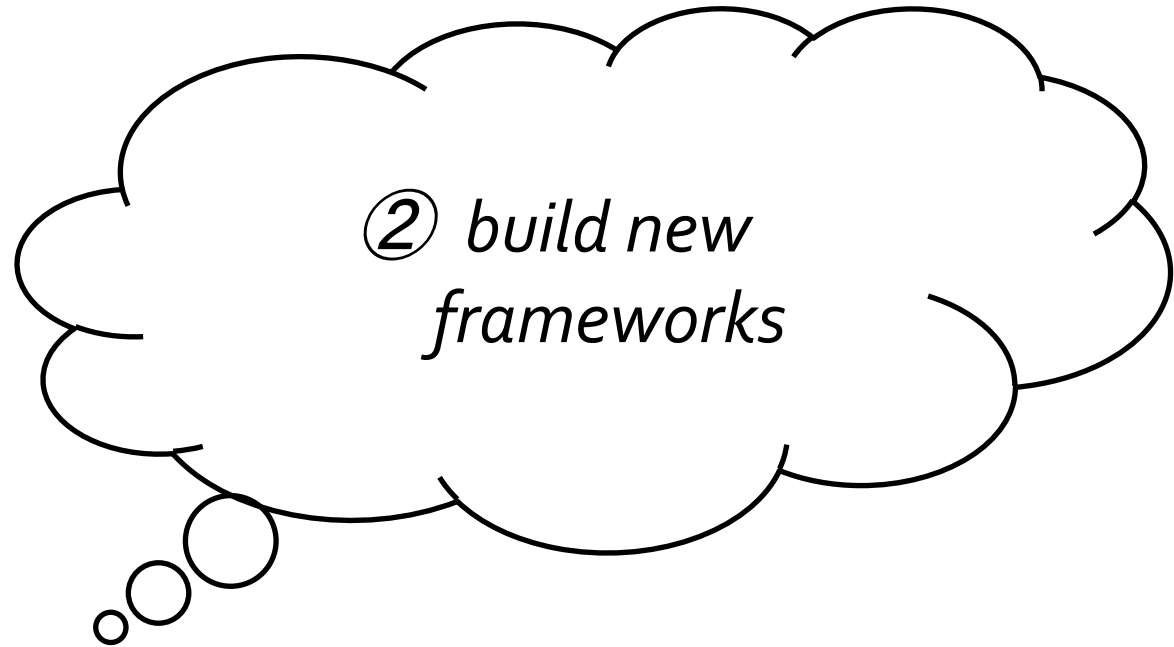


higher utilization!

# our intuition

② build new frameworks

"Map/Reduce is a big hammer, but not everything is a nail!"

↳ Spark

# anatomy of Spark

# anatomy of a framework

coordinator



workers

# framework
# ≈
# distributed system

# anatomy of a framework



coordinator

workers

# anatomy of a framework

scheduler



tasks (and/or executors)

# execution coordination
# ==
# scheduling

# Mesos: level of indirection

scheduler

tasks (and/or executors)

# Mesos: level of indirection

scheduler

Mesos (master)

Mesos (slaves)

# Mesos: a level of indirection

+ provide common functionality every new distributed system *re-implements* like failure detection, task distribution, task starting, task monitoring, task killing, task cleanup!

# Mesos: level of abstraction

scheduler

Mesos (master)

Mesos (slaves)

# Mesos: level of abstraction

**Mesos**

build and run
frameworks
using *resources*

# Mesos: level of abstraction

**Mesos**

build and run
frameworks
using *resources*

**IaaS**

provision and manage
*machines*

# Mesos: level of abstraction

**PaaS** — deploy and manage *applications/services*

**Mesos** — build and run frameworks using *resources*

**IaaS** — provision and manage *machines*

# PaaS on Mesos

| | |
|---|---|
| **PaaS** | |
| **Mesos** | |

build and run a PaaS
on top of Mesos:
Apache Aurora and Marathon

# Mesos on IaaS

Mesos

IaaS

use OpenStack or EC2
to run Mesos

# Mesos on IaaS++

| Mesos | |
|-------|---|
| hardware | IaaS |

use OpenStack or EC2
*or physical machines*
to run Mesos

# Mesos: datacenter kernel

framework

Mesos

hardware    IaaS

provide common functionality
via an API (kernel)

**Apache Mesos is a distributed system for running and building other distributed systems**
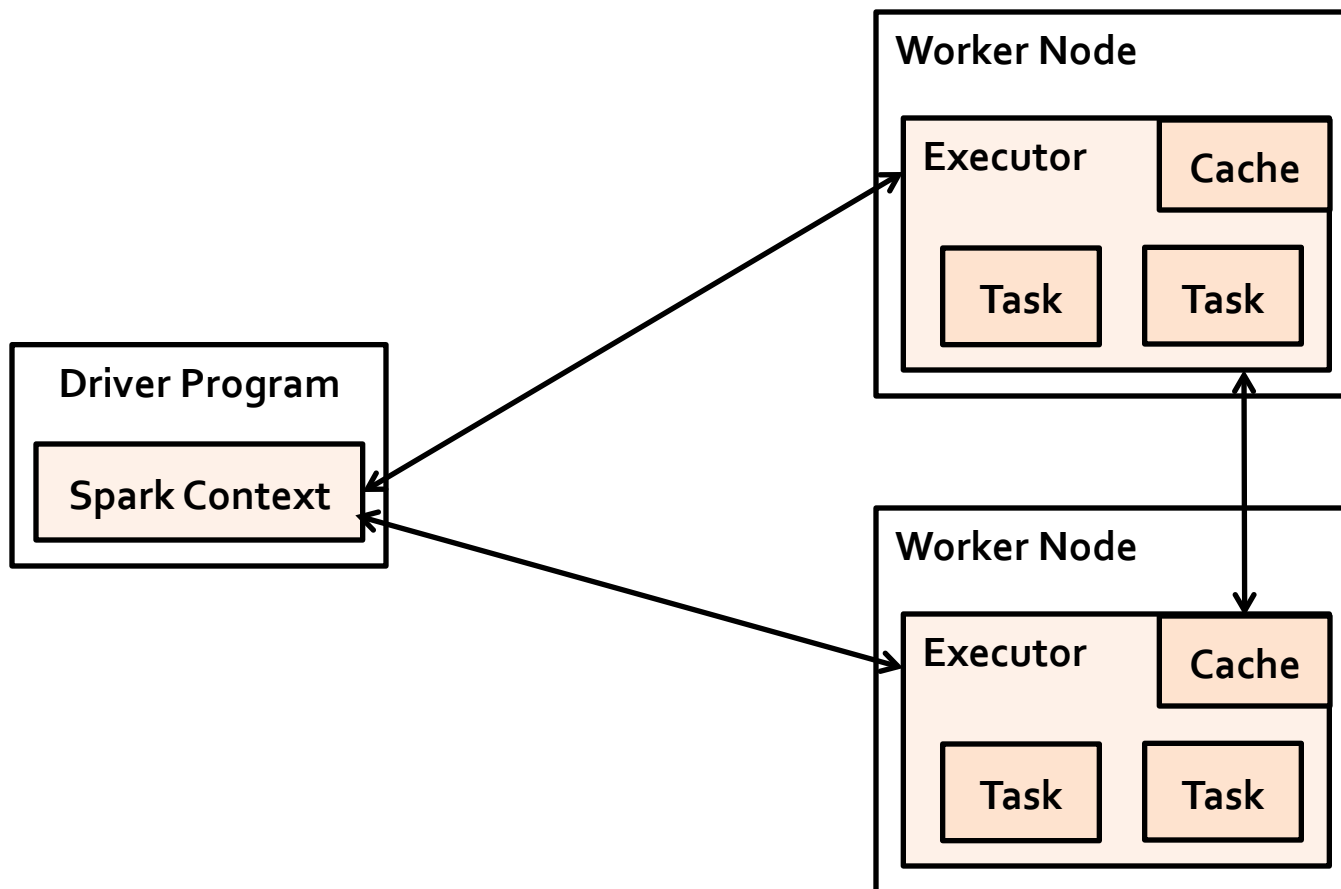
# Mesos is a cluster manager

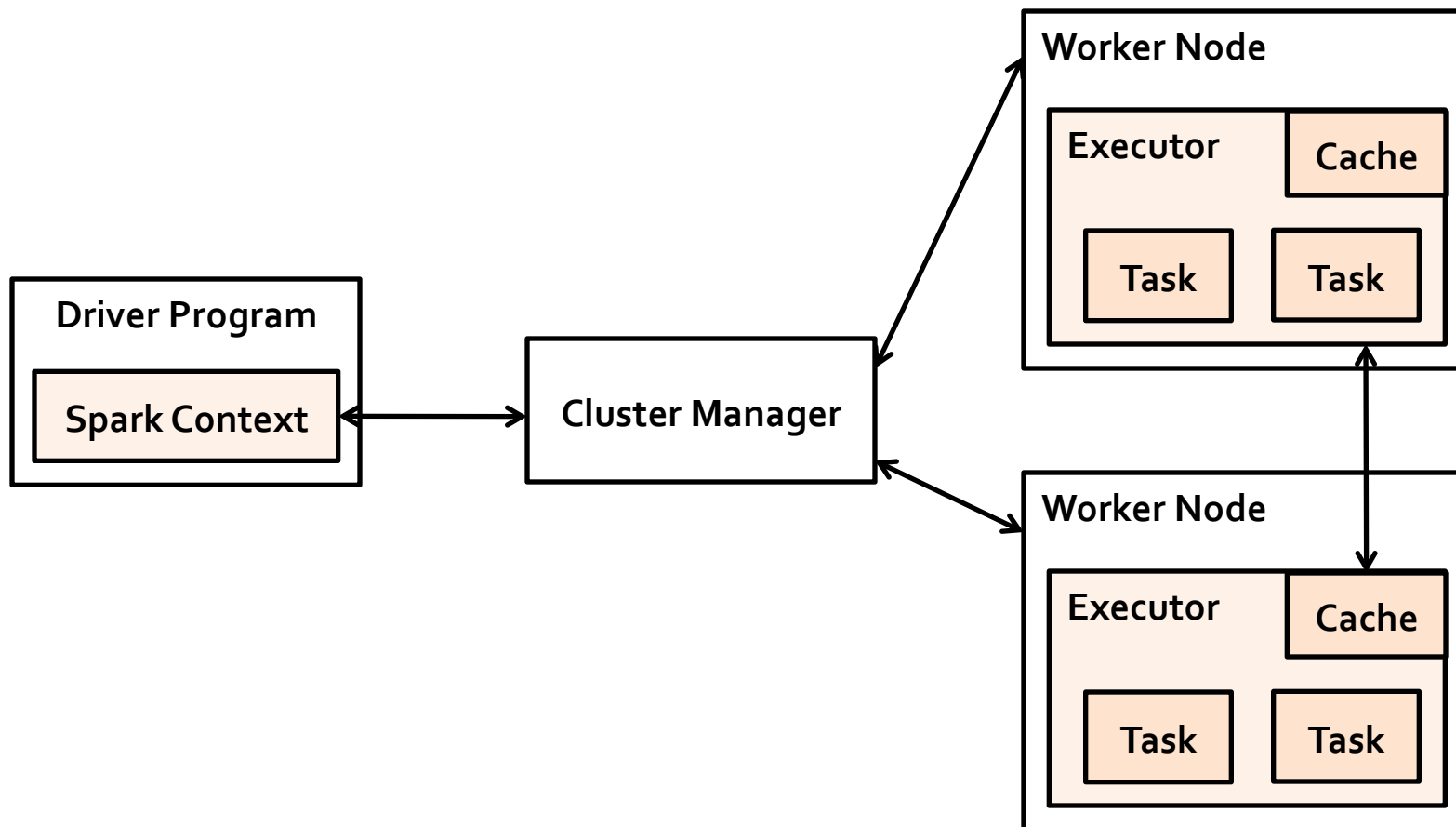# agenda

① Mesos

② Spark on Mesos     ⟵

③ why Mesos?
      ① multi-tenancy
      ② fine-grained sharing
      ③ why not?

④ long-lived services and other frameworks
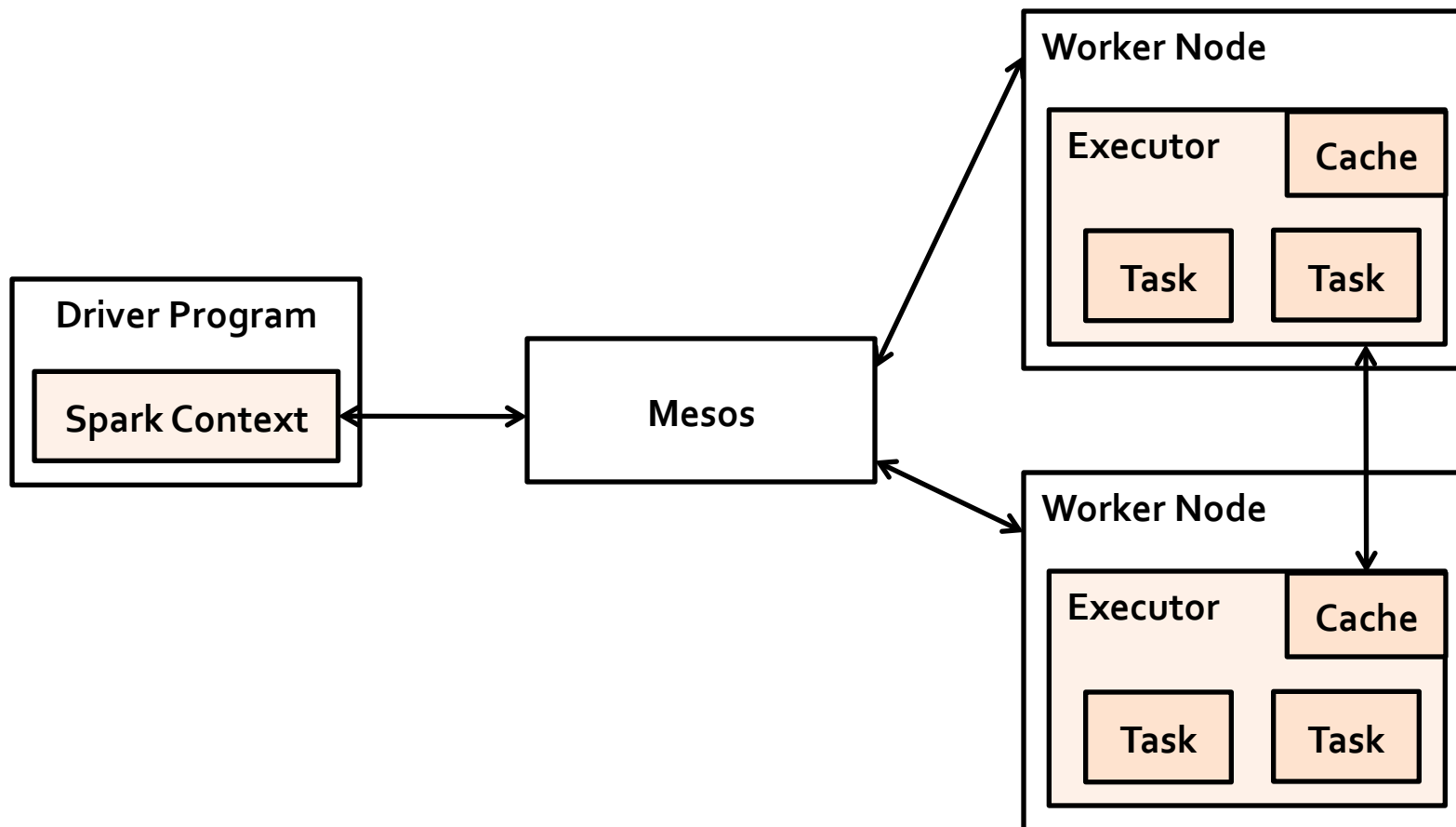
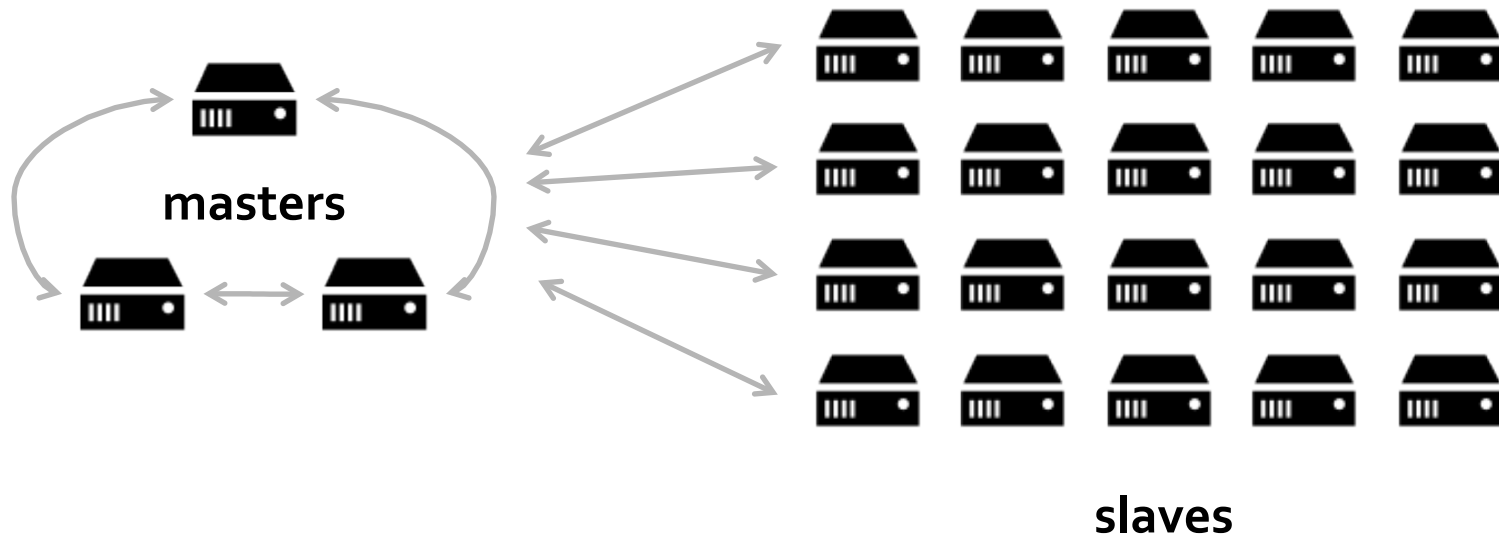# anatomy of Spark

# anatomy of Spark

# anatomy of Spark

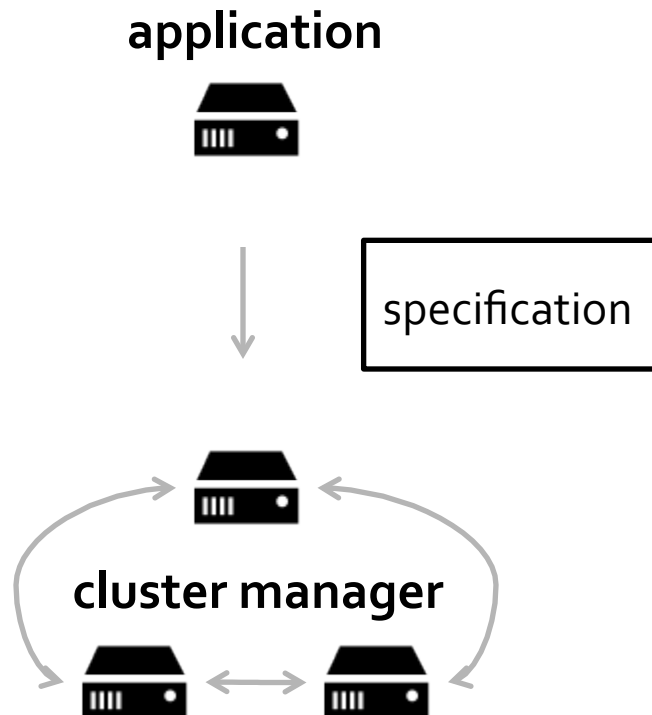# Mesos is a distributed system with a master/slave architecture



masters

slaves

**Mesos challenged
the status quo
of cluster managers**

# cluster manager status quo

**application**

**specification**
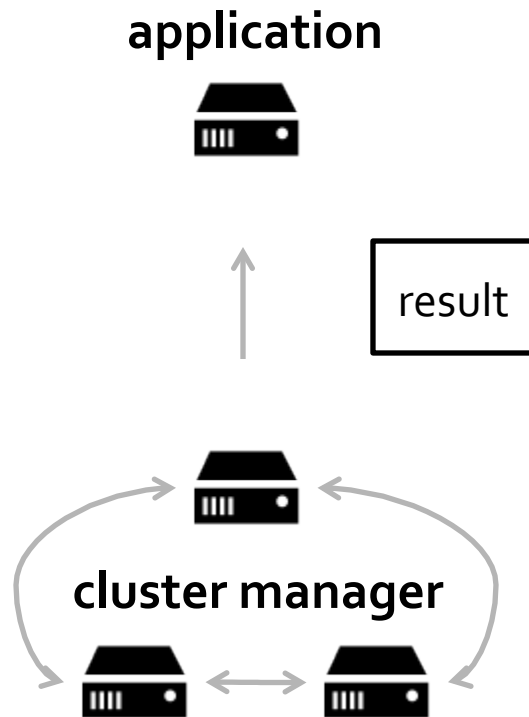
the specification includes as much information as possible to assist the cluster manager in scheduling and execution

**cluster manager**

# cluster manager status quo

application

result

cluster manager

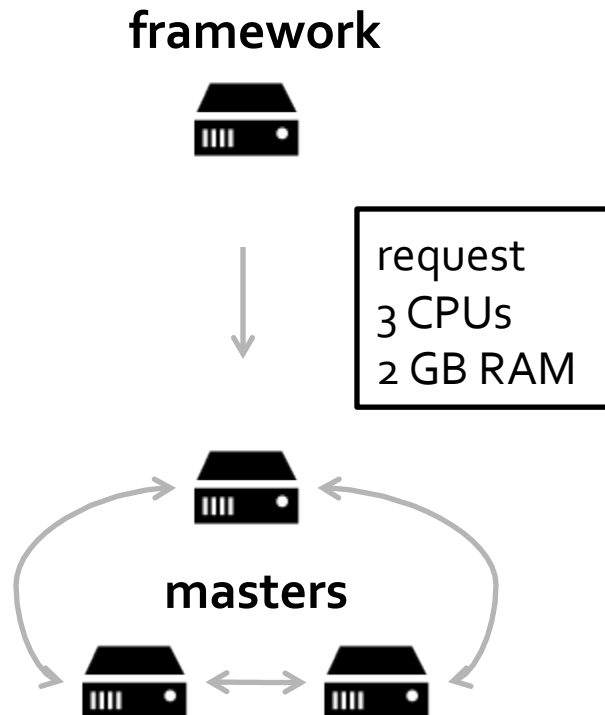# problems with specifications

① hard to specify certain desires or constraints

② hard to update specifications dynamically as tasks executed and finished/failed

# an alternative model

**framework**

request
3 CPUs
2 GB RAM

**masters**

a request is purposely simplified
subset of a specification, mainly
including the required resources

question: what should you do
if you can't satisfy a request?

**question: what should you do if you can't satisfy a request?**

① wait until you can …

**question: what should you do if you can't satisfy a request?**
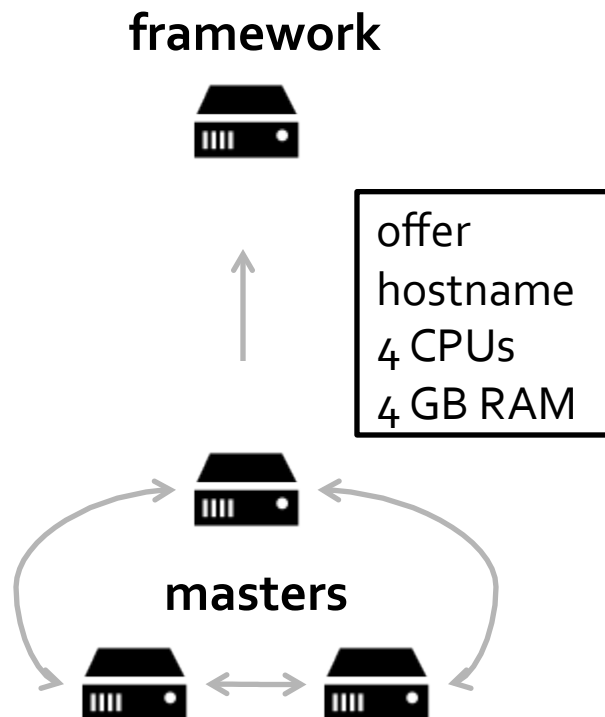
① wait until you can ...

② *offer* best you can <u>immediately</u>

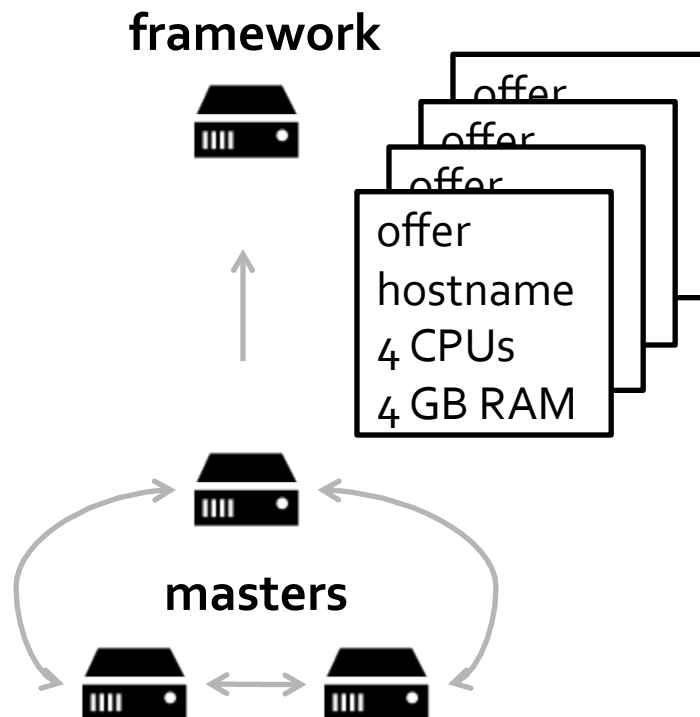**question: what should you do if you can't satisfy a request?**

① wait until you can …

② *offer* best you can <u>immediately</u>

# Mesos model

**framework**



offer
hostname
4 CPUs
4 GB RAM

**masters**

# Mesos model

framework

masters

offer
hostname
4 CPUs
4 GB RAM

# an analogue:
# non-blocking sockets

application

kernel

```
write(s, buffer, size);
```

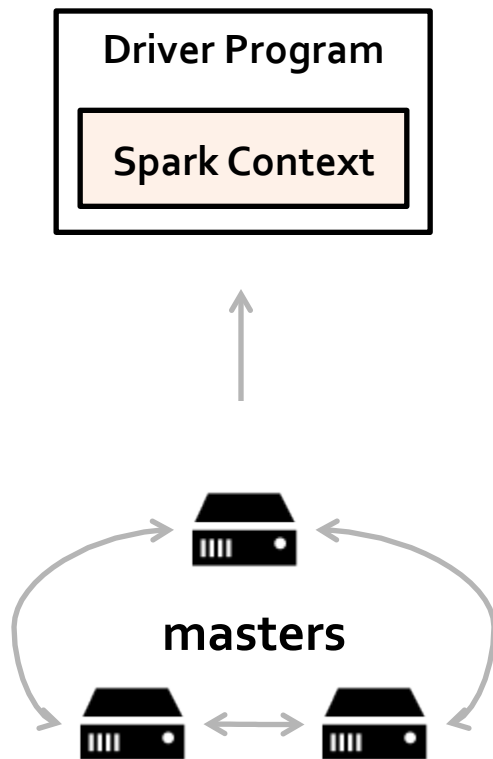# an analogue:
# non-blocking sockets
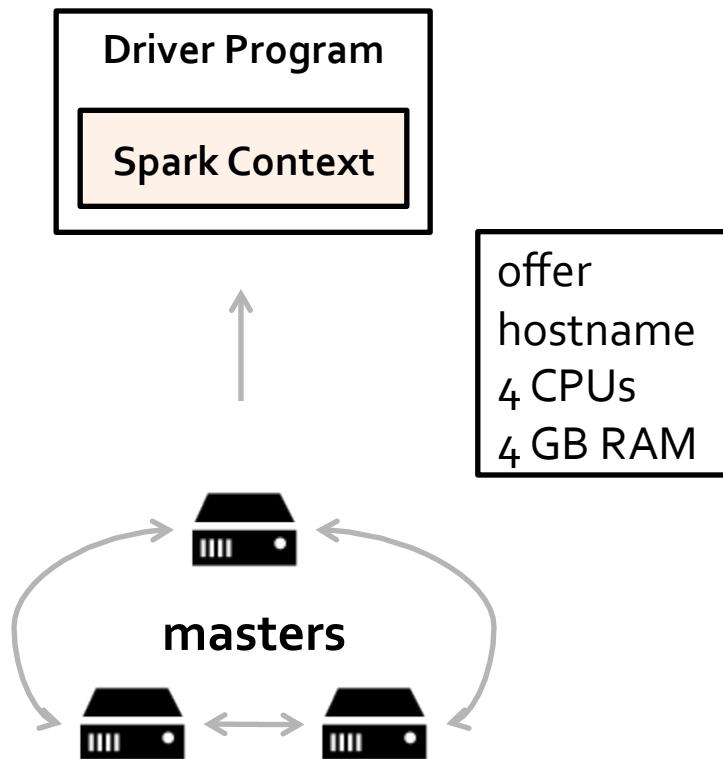
application

↑

```
42 of 100 bytes written!
```

kernel

offers represent the
current *snapshot*
of available resources
a framework can use

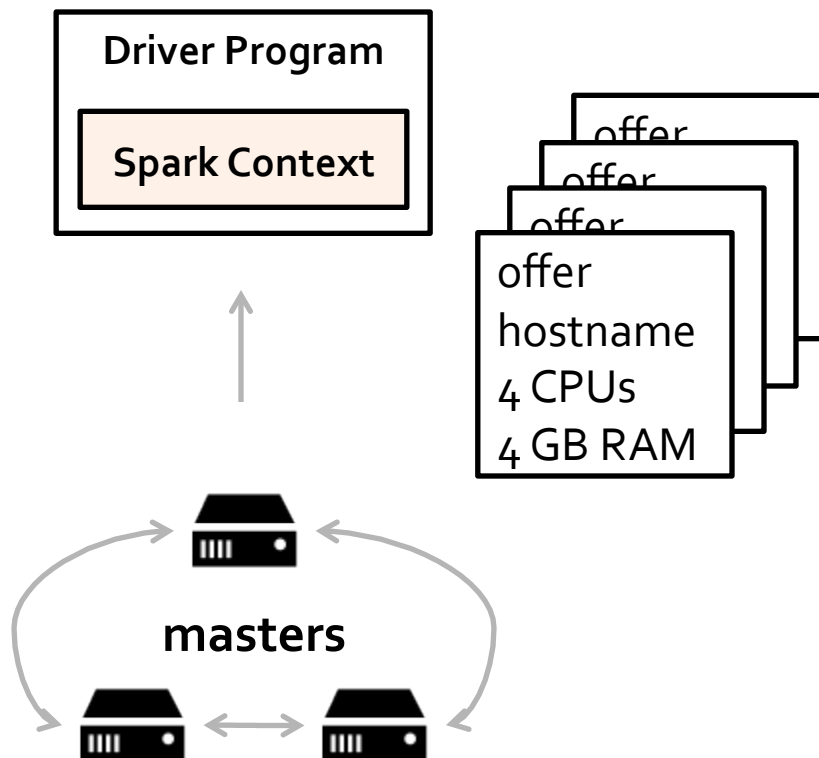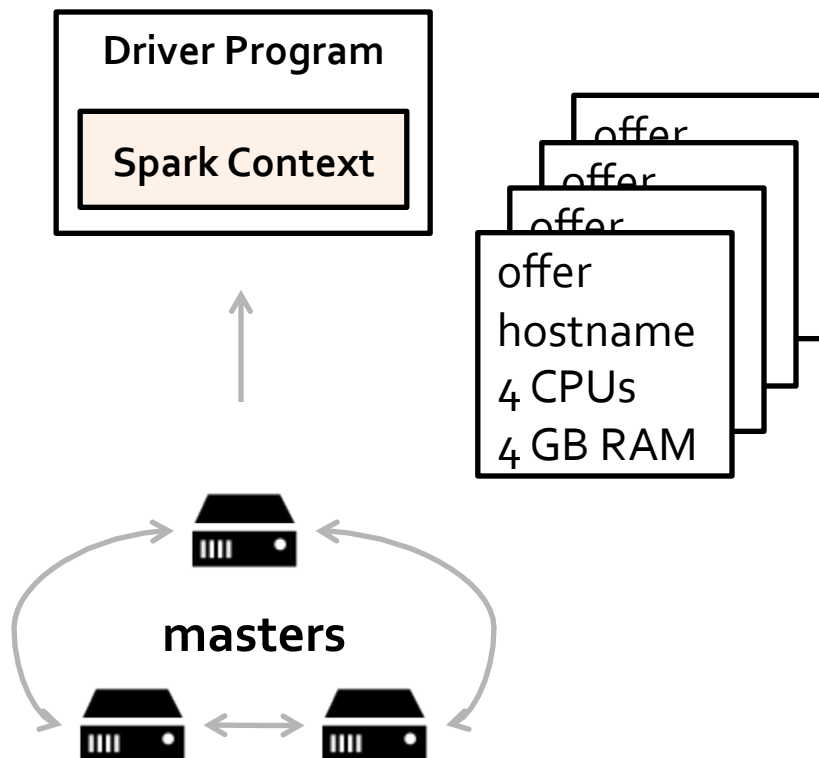(requests are complimentary, but not necessary; see Google's Omega)

# Spark on Mesos



Driver Program

Spark Context
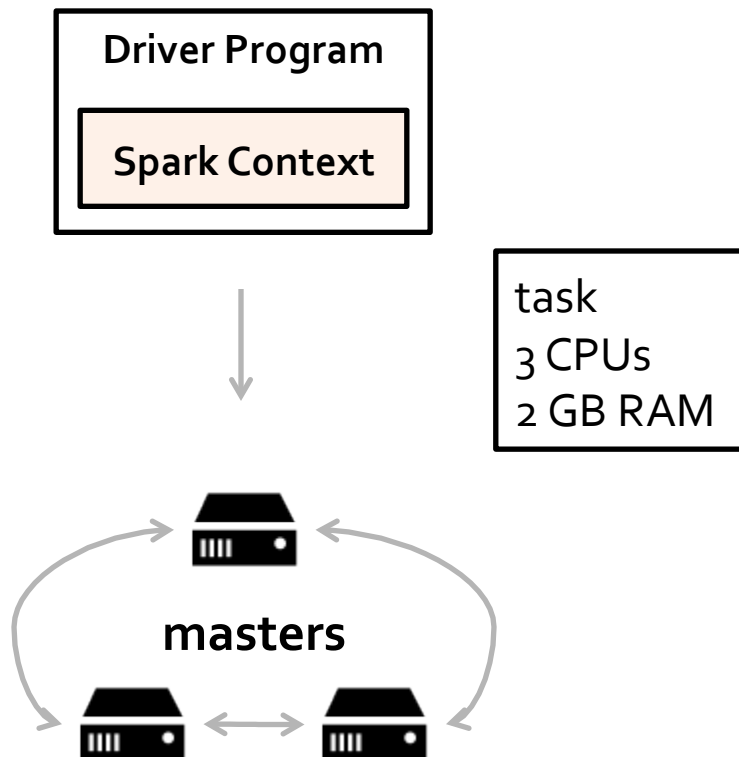
masters

# Spark on Mesos

Driver Program

Spark Context

offer
hostname
4 CPUs
4 GB RAM

masters

# Spark on Mesos

**Driver Program**

**Spark Context**

offer
offer
offer

offer
hostname
4 CPUs
4 GB RAM

**masters**

# Spark on Mesos

**Driver Program**

**Spark Context**

offer
offer
offer

offer
hostname
4 CPUs
4 GB RAM

**masters**

Spark uses the offers to perform it's own scheduling

# Spark on Mesos

Driver Program

Spark Context

task
3 CPUs
2 GB RAM

masters

Spark uses the offers to perform it's own scheduling

# Spark on Mesos

Driver Program

Spark Context

task
3 CPUs
2 GB RAM

masters

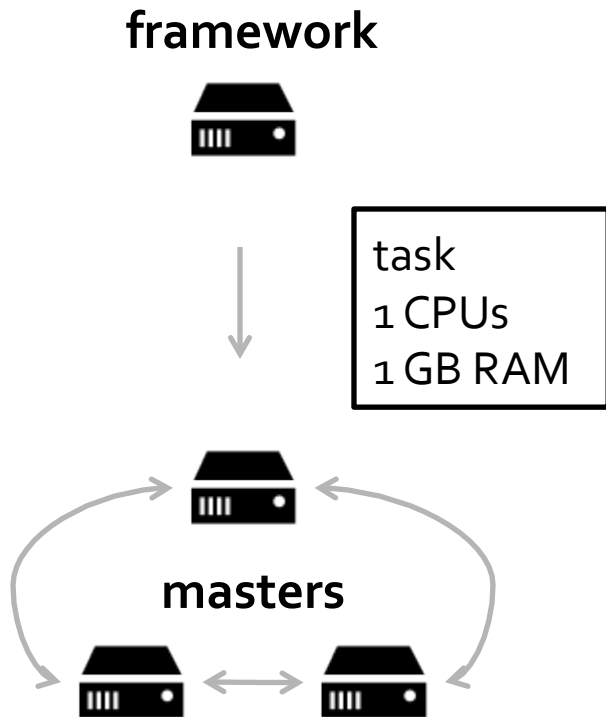Spark uses the offers to perform it's own scheduling

"two-level scheduling"

# "two-level scheduling"

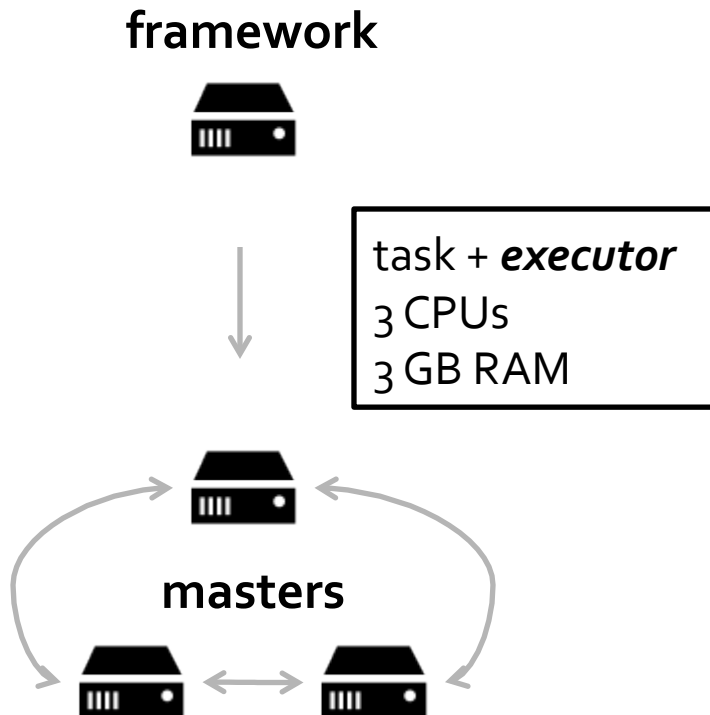Mesos: controls resource *allocations* to Spark

Spark: makes decisions about what tasks to run given available resources

# execution

**framework**



task
1 CPUs
1 GB RAM

**masters**

frameworks launch fine-grained
tasks for execution

# execution

**framework**



task + ***executor***
3 CPUs
3 GB RAM

**masters**
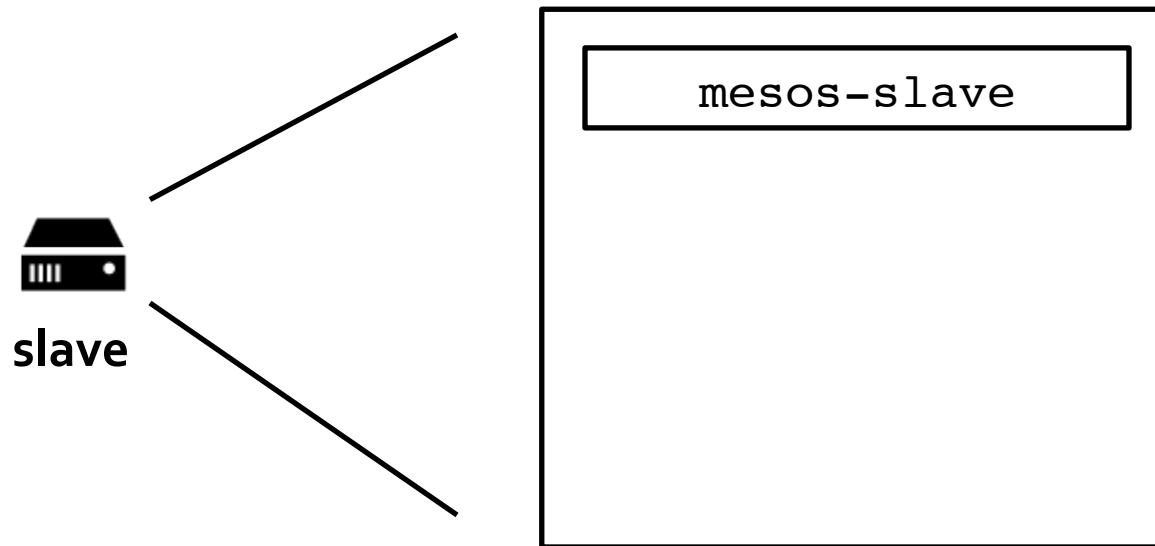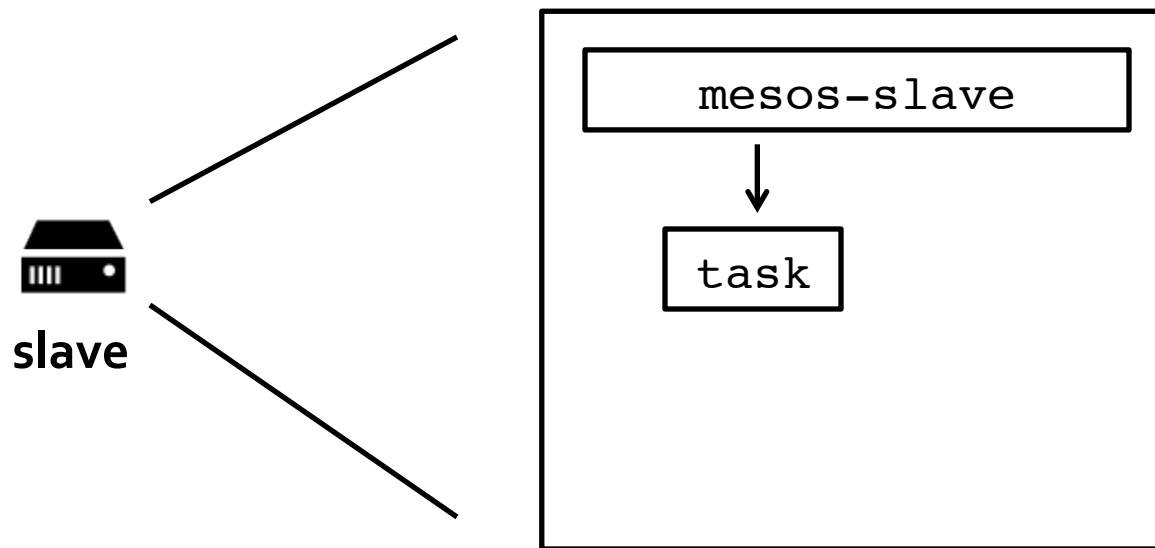
frameworks launch fine-grained
tasks for execution

if necessary, a framework can
provide an *executor* to handle the
execution of a task

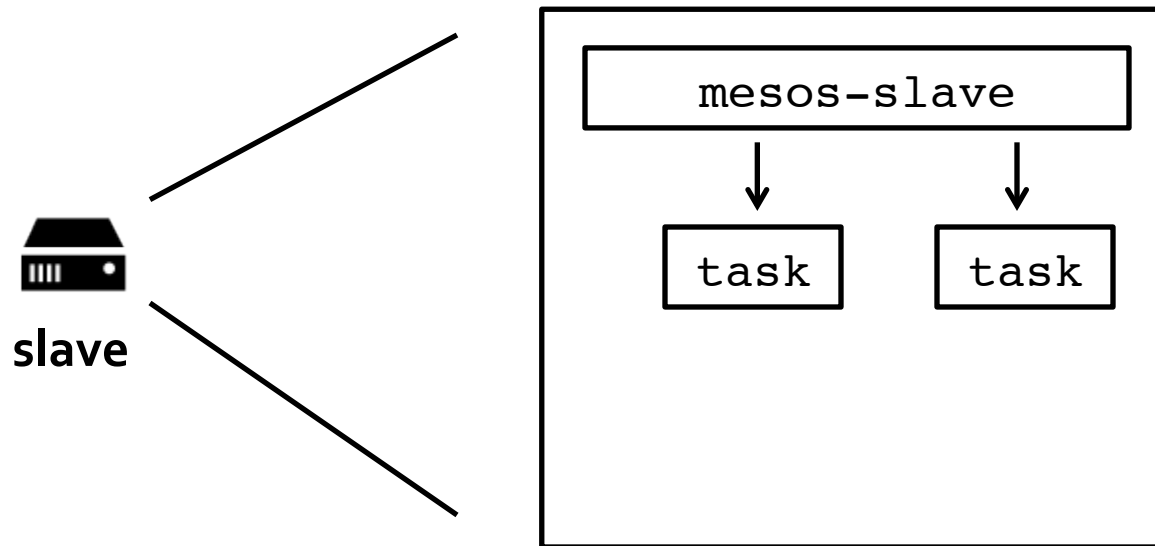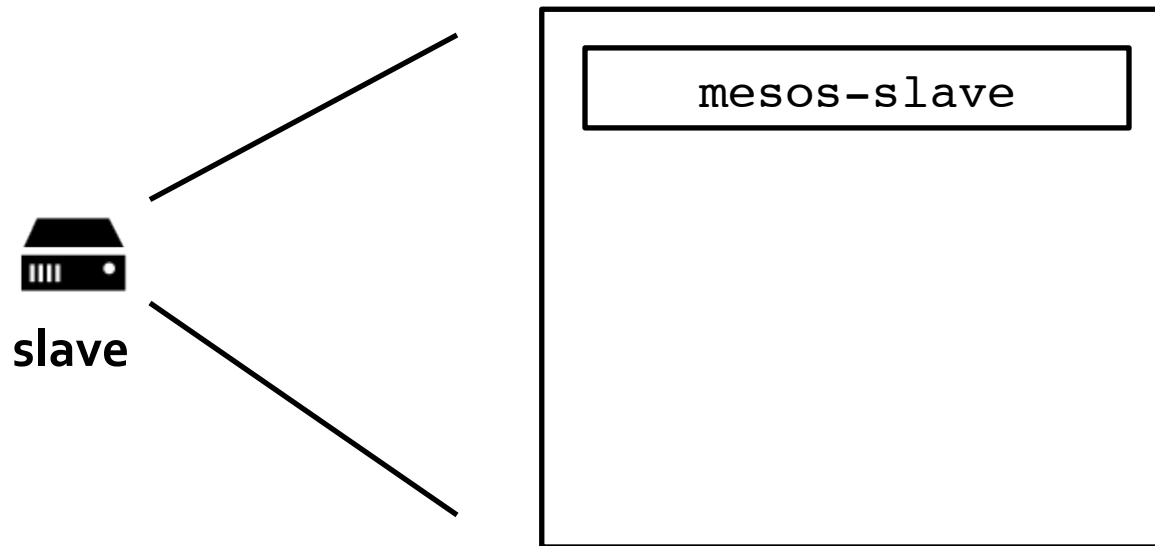# a task with a *command*



mesos-slave

slave
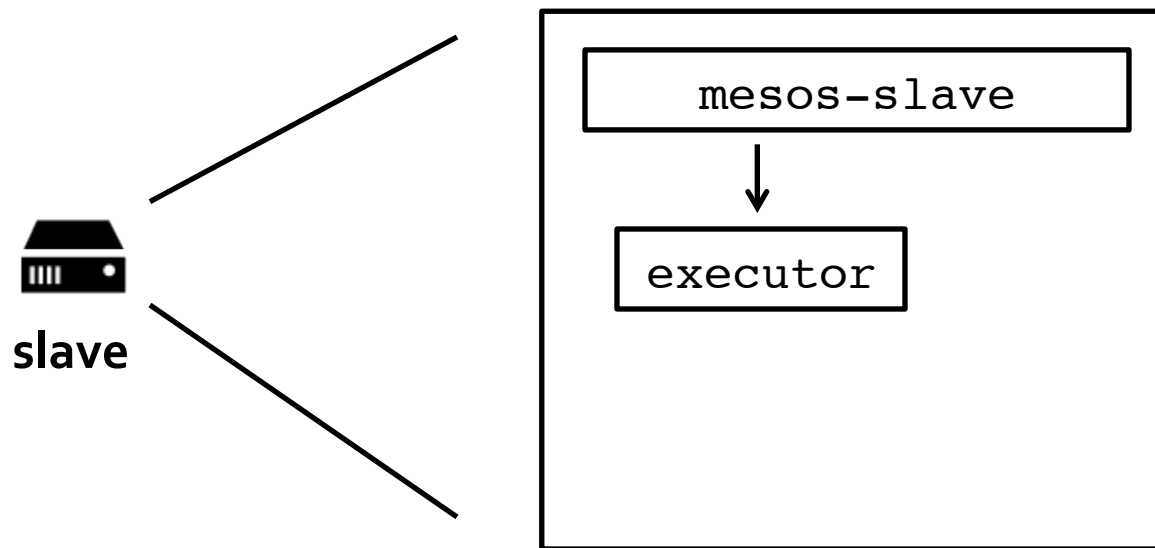
# a task with a *command*

# a task with a *command*

# a task with an *executor*
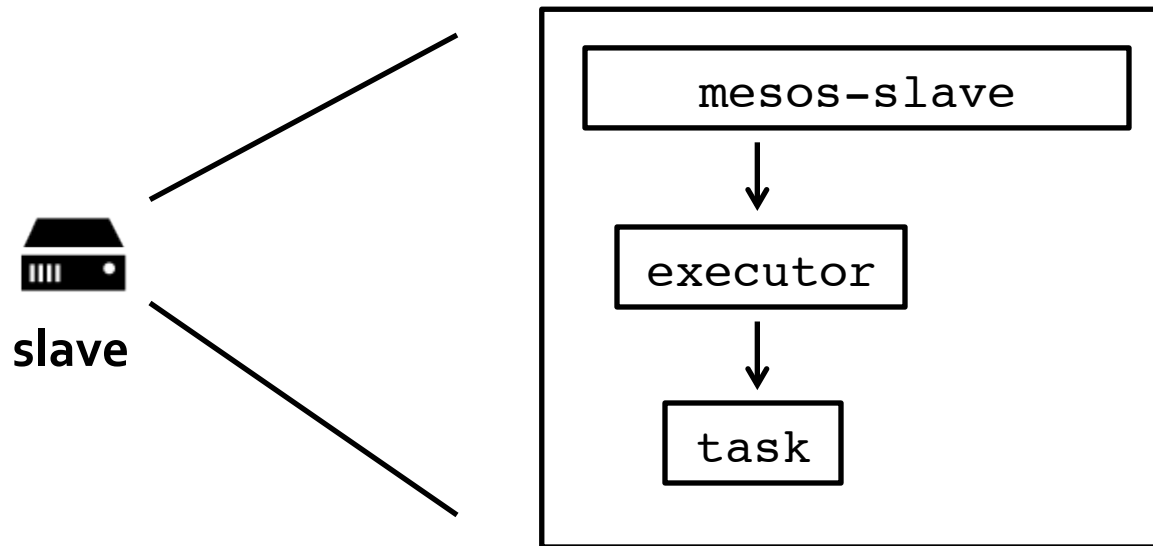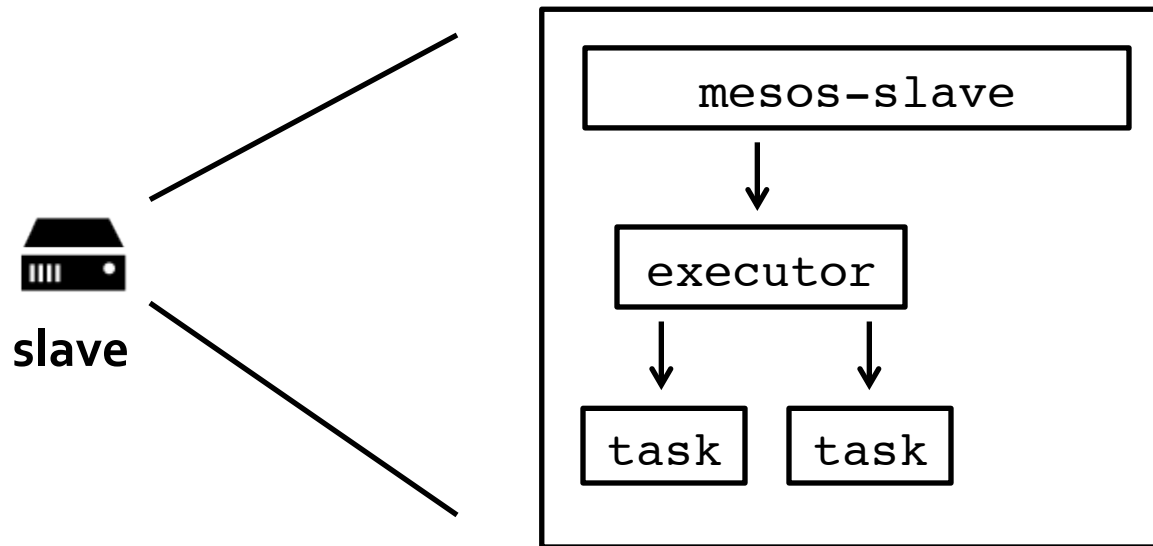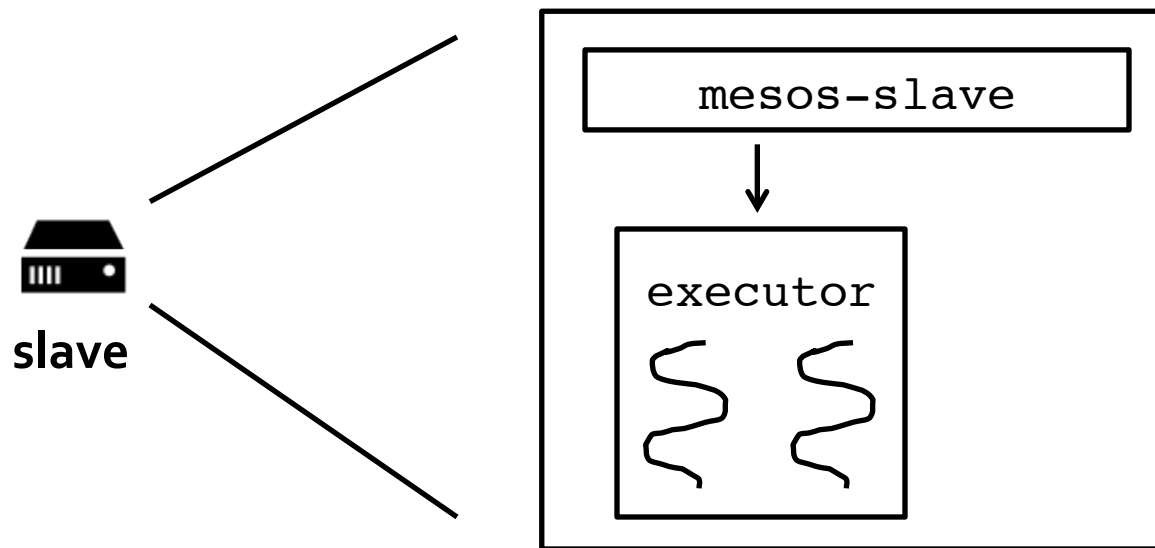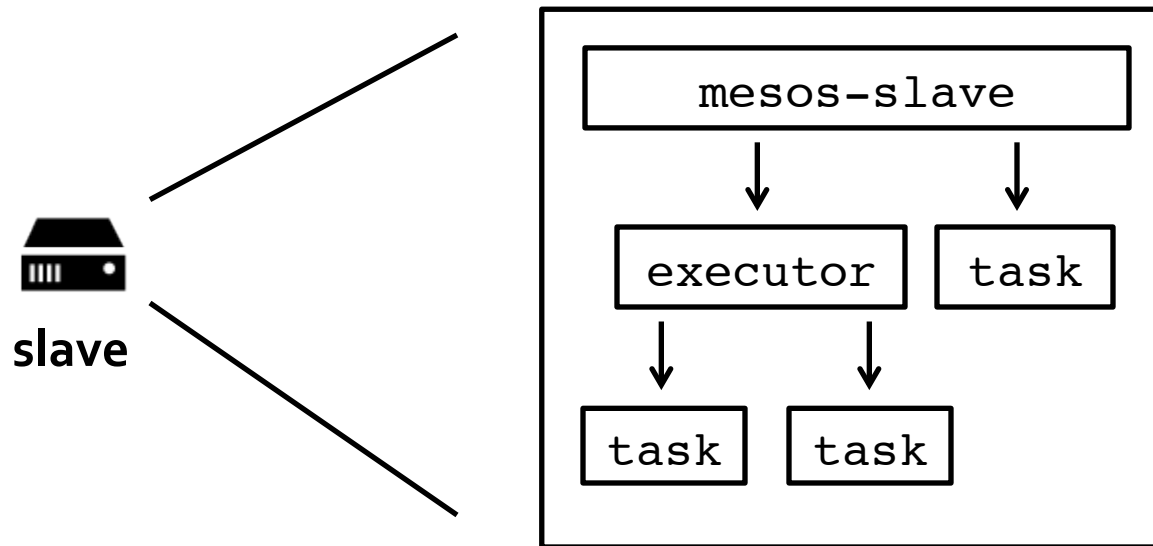


slave

mesos-slave

# a task with an *executor*

# a task with an *executor*
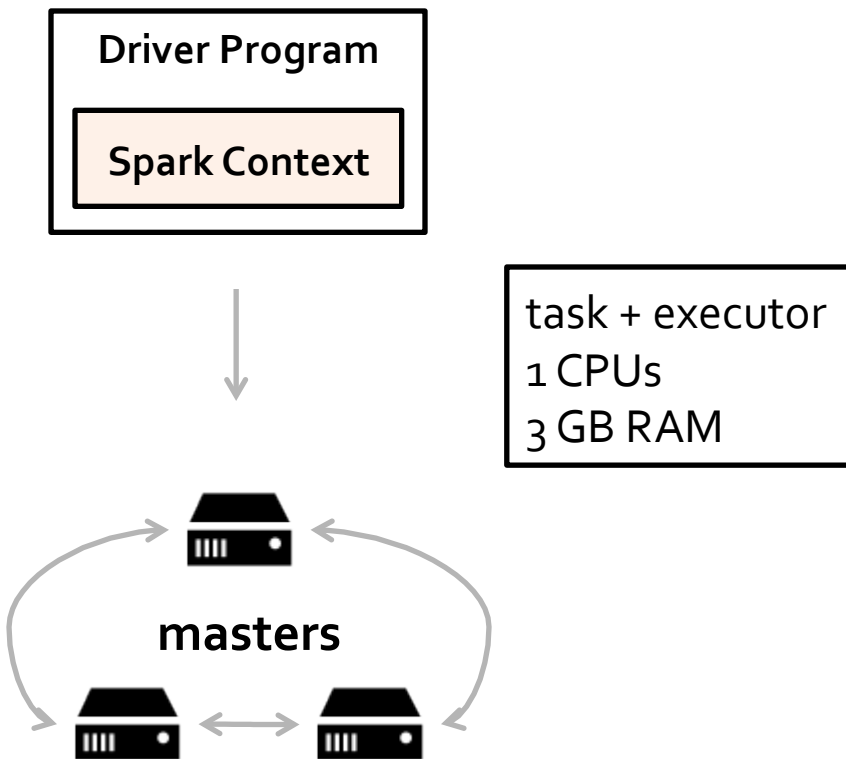
# a task with an *executor*
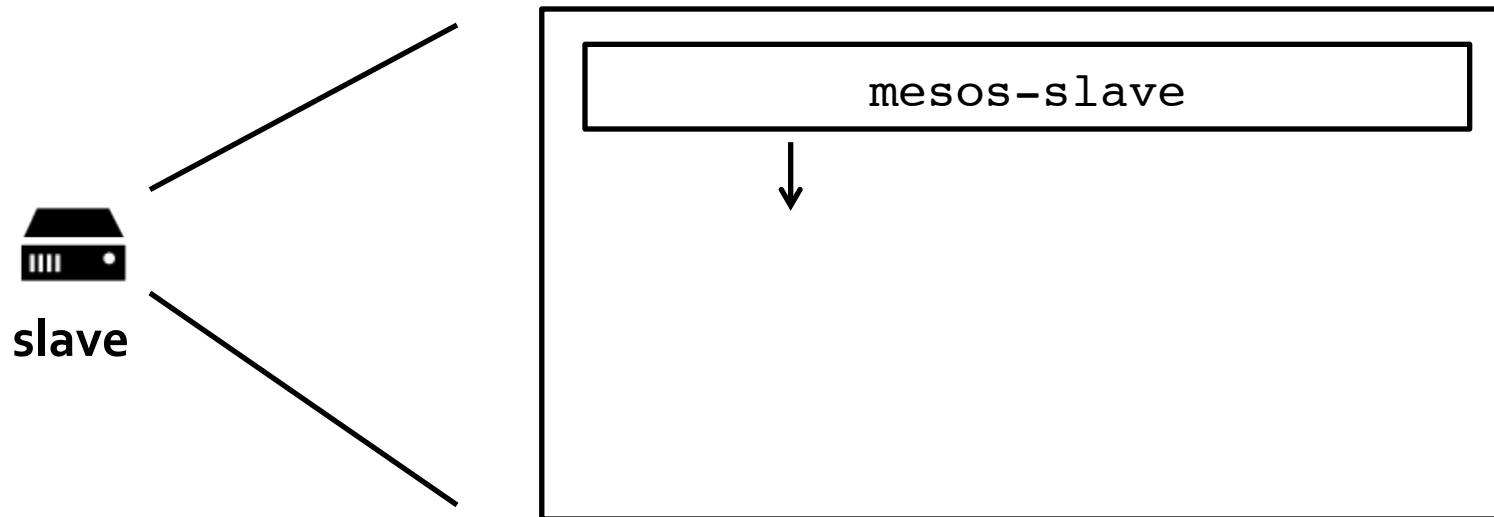
# a task with an *executor*

# a task with an *executor*

# Spark execution

Driver Program
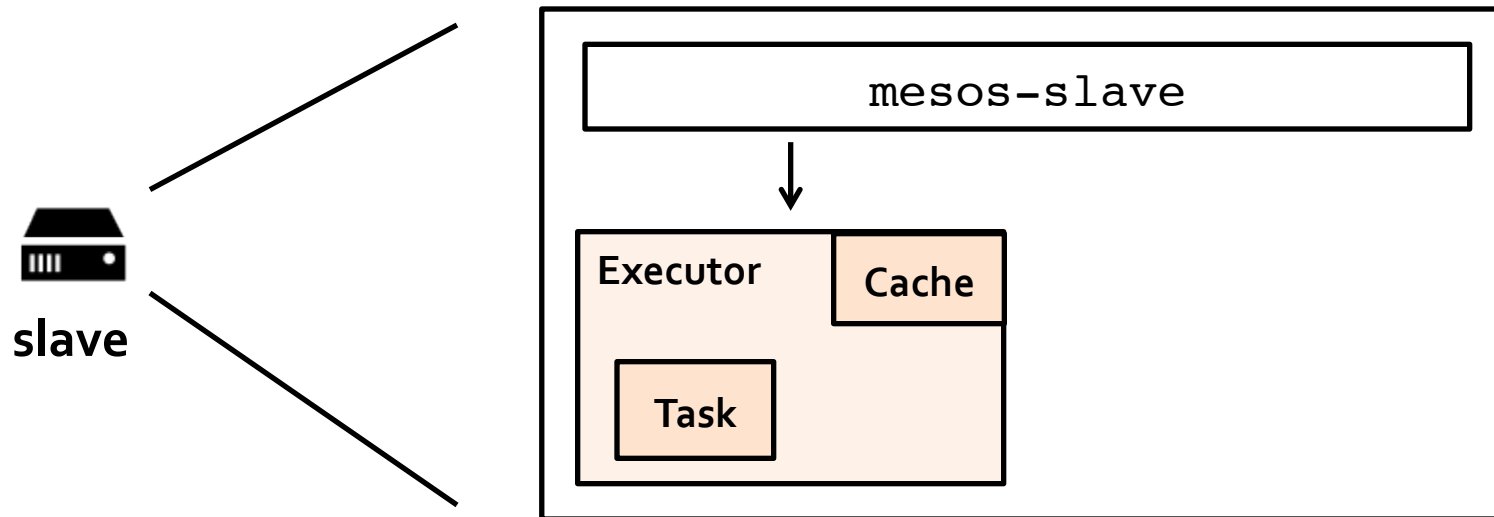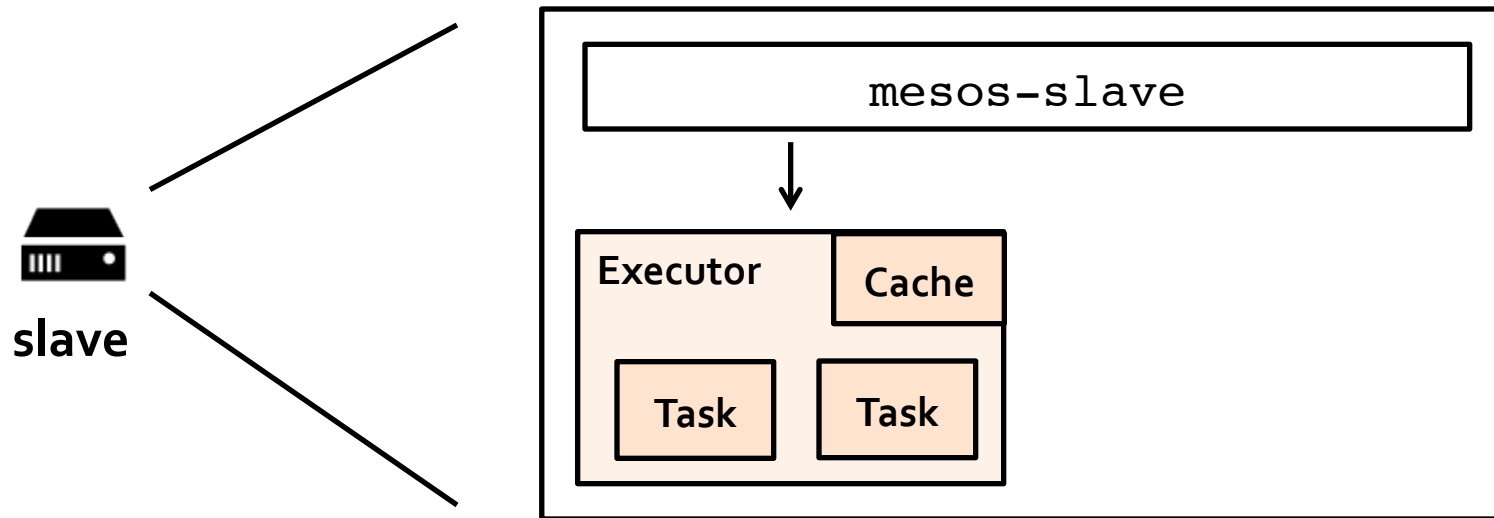
Spark Context

task + executor
1 CPUs
3 GB RAM

masters

# Spark execution

**slave**

mesos-slave

# Spark execution

**slave**

mesos-slave

Executor

Cache

# Spark execution

**slave**

mesos-slave

Executor

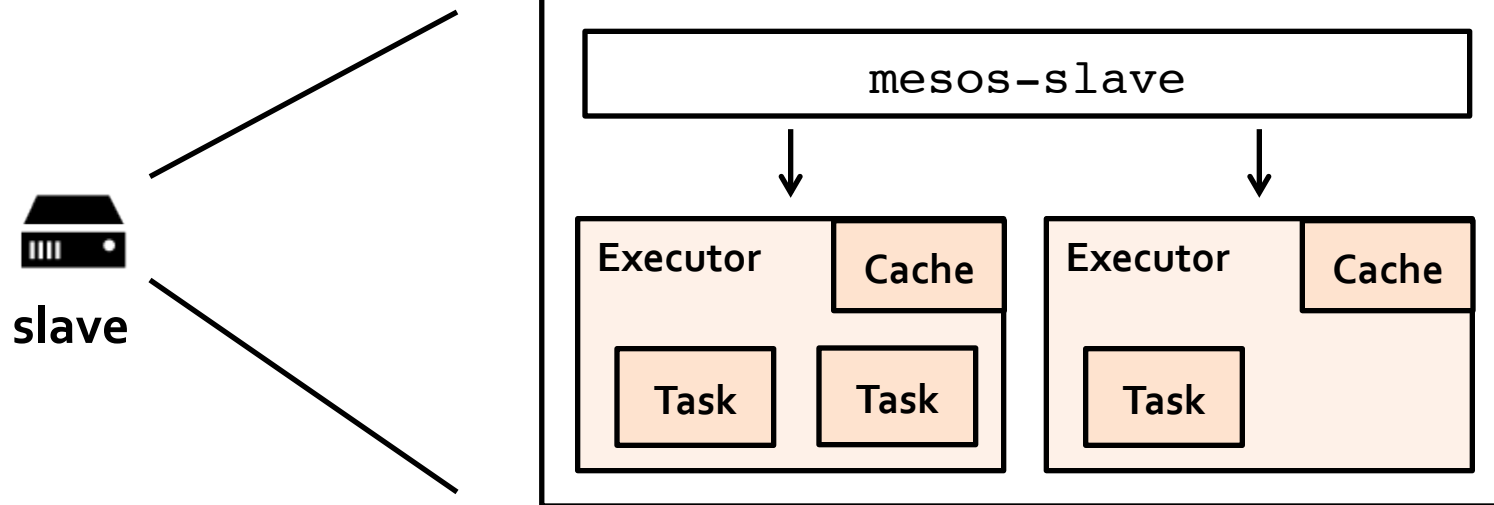Cache

Task

# Spark execution

slave

mesos-slave

Executor Cache

Task Task

# Spark execution

slave

# resource isolation

Mesos has containerization support on Linux
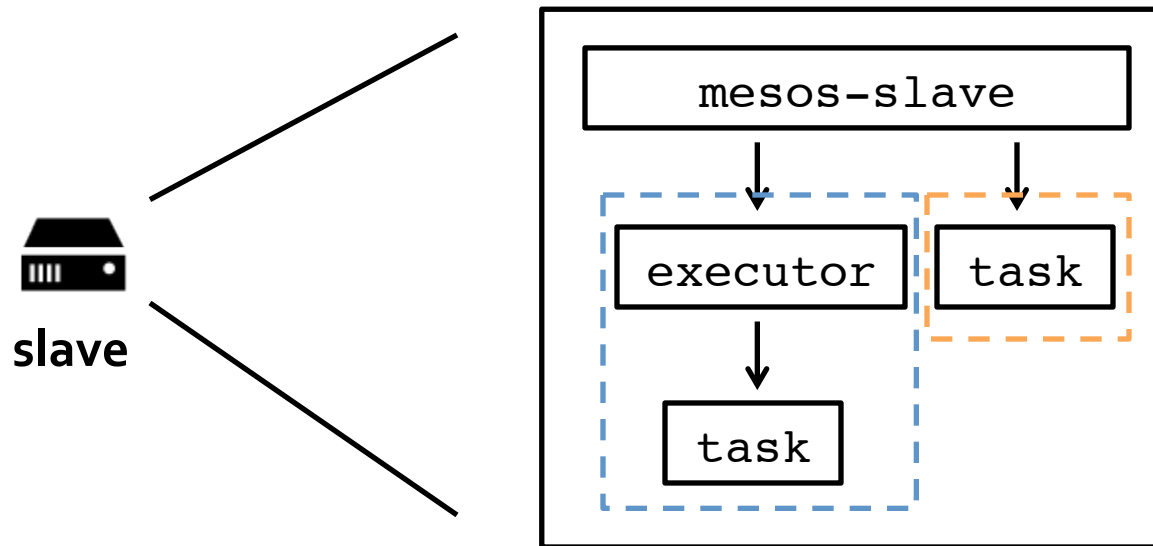(built-in usage of cgroups and namespaces)

isolator modules:

CPU (upper and lower bounds)

memory
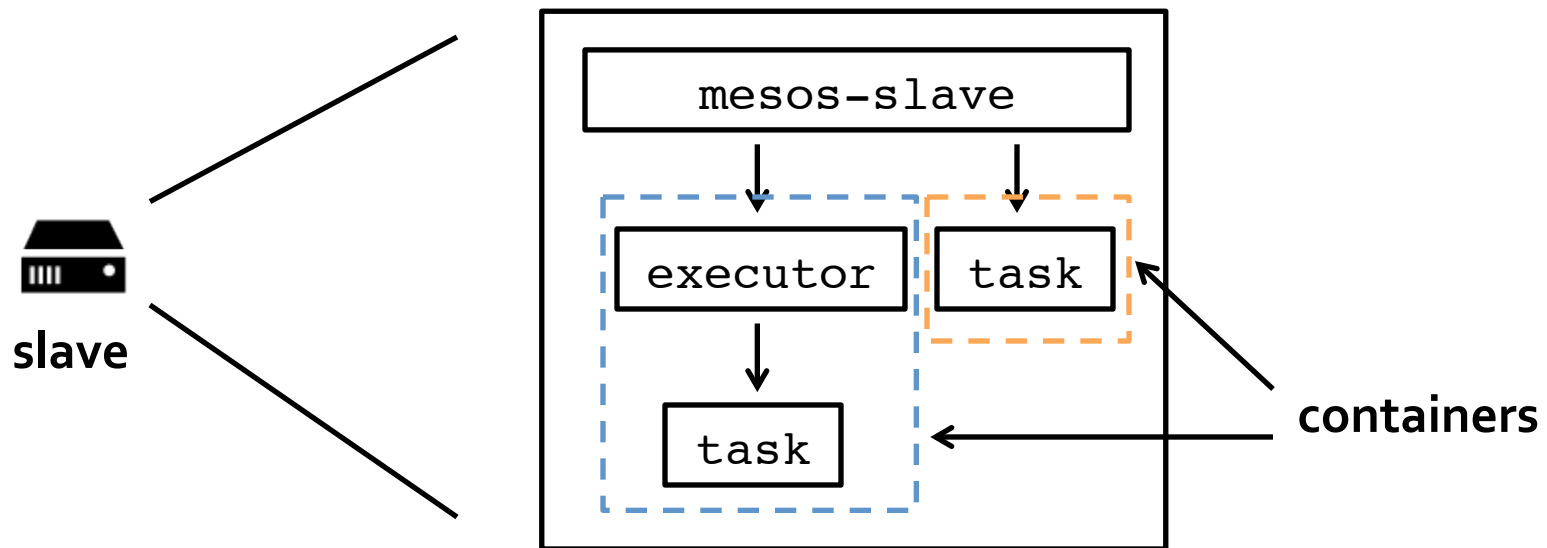
network I/O (in development)

filesystem (using LVM, planned)

# resource isolation



slave

mesos-slave

executor

task

task

# resource isolation



mesos-slave

executor

task

task

slave

containers

# resource isolation

# agenda

① Mesos

② Spark on Mesos

③ why Mesos?
  ① multi-tenancy ←
  ② fine-grained sharing
  ③ why not?

④ long-lived services and other frameworks

# multi-tenancy



slave(s)

mesos-slave

mesos-slave

# multi-tenancy (only Spark)



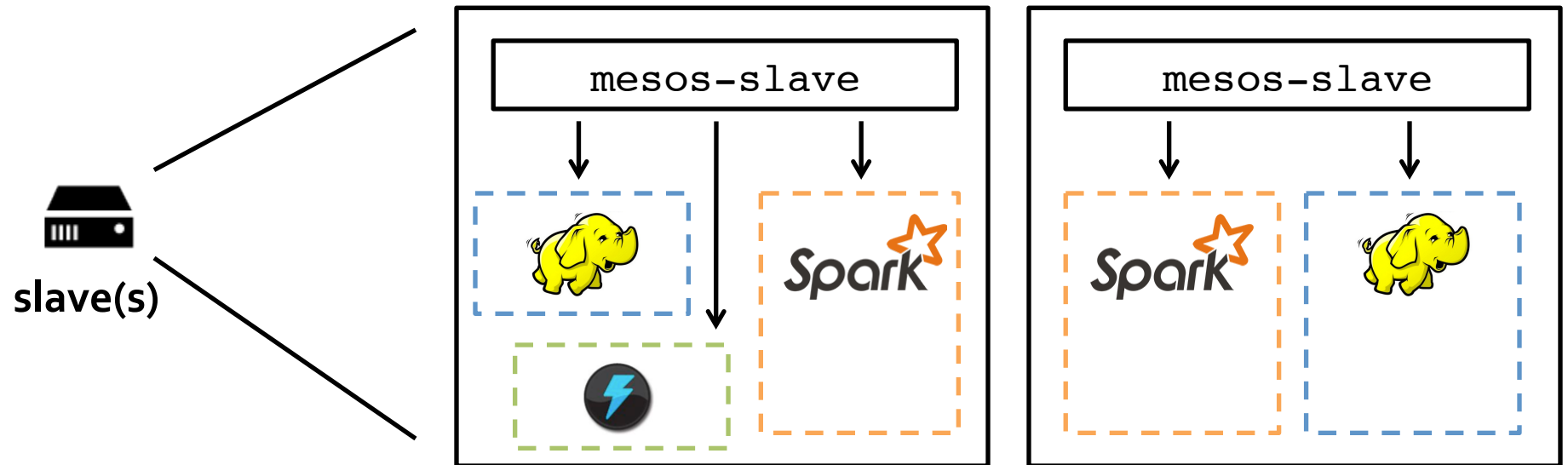(can approximate w/ standalone mode by setting max # cores per application, otherwise get FIFO execution)

# multi-tenancy
# (only Spark)



**slave(s)**

| mesos-slave | | mesos-slave | |
|---|---|---|---|
| Spark 0.9.0 | Spark 0.8.1 | Spark 0.8.1 | Spark 0.9.0 |

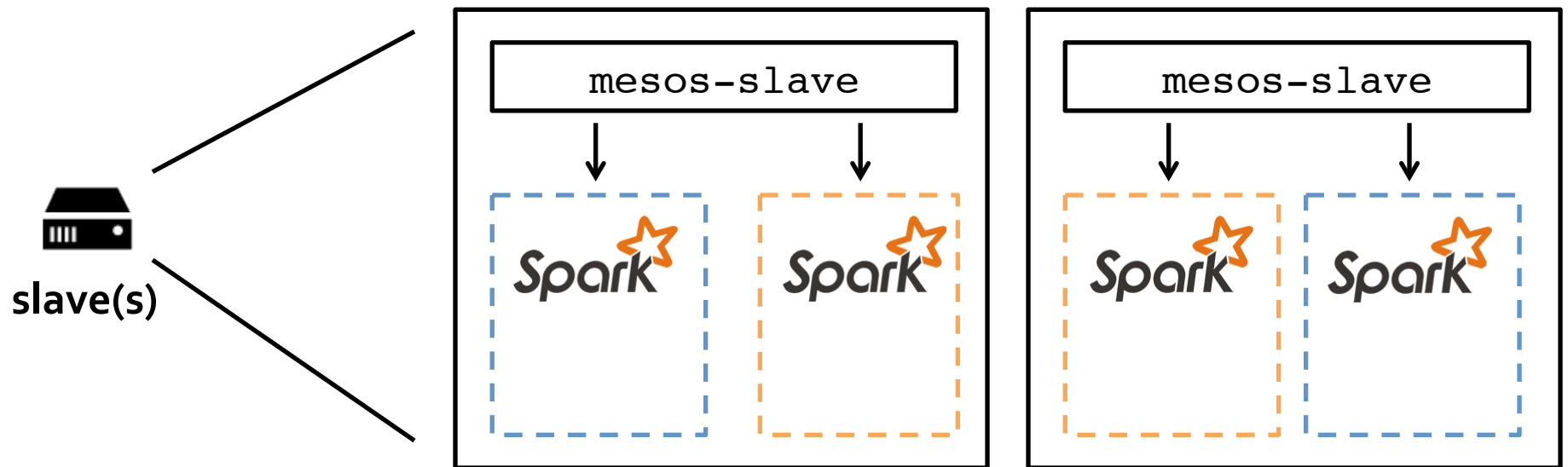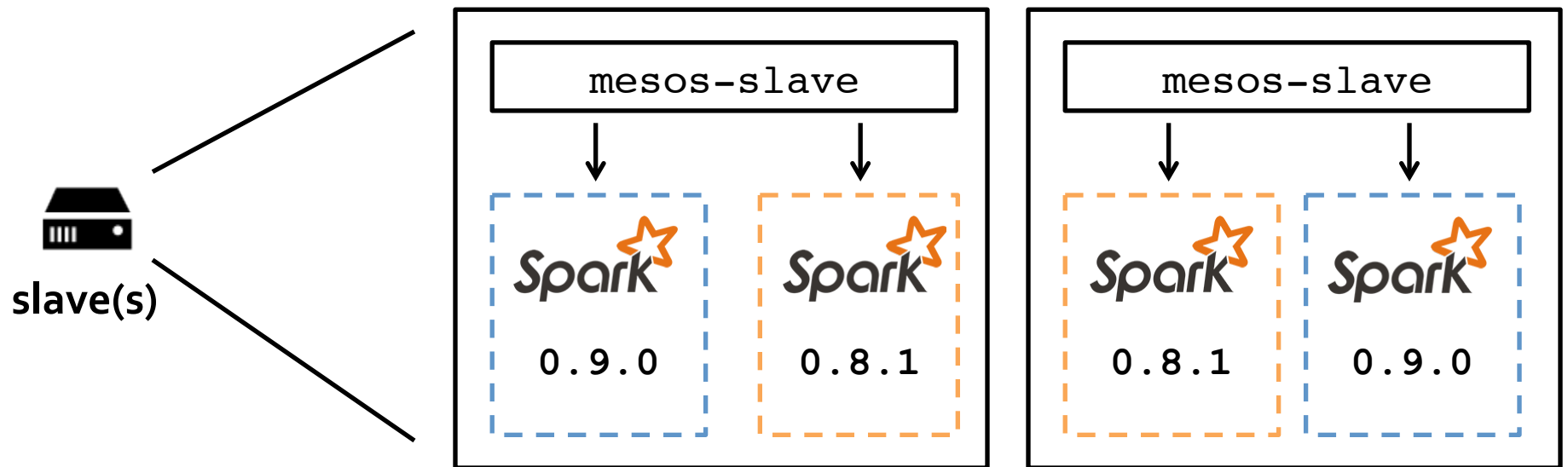(run the tried and true and test out the new at the same time!)

# agenda

① Mesos

② Spark on Mesos

③ why Mesos?
   ① multi-tenancy
   ② fine-grained sharing ⬅
   ③ why not?

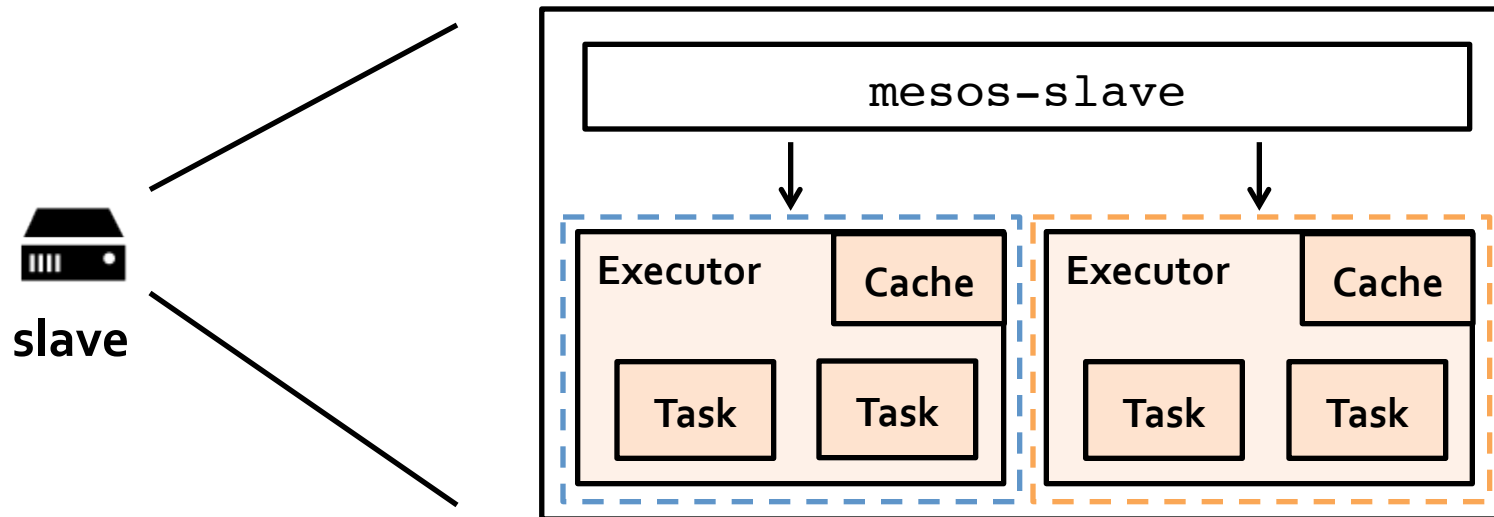④ long-lived services and other frameworks

# fine-grained sharing



**slave**

mesos-slave

Executor    Cache

Task    Task

Executor    Cache

Task    Task

Spark executors only consume memory, can share CPU between

# fine-grained sharing



slave

mesos-slave

Executor | Cache

Task | Task

Executor | Cache

Task

Spark executors only consume memory, can share CPU between

# fine-grained sharing



Spark executors only consume memory, can share CPU between

# fine-grained sharing



slave

mesos-slave

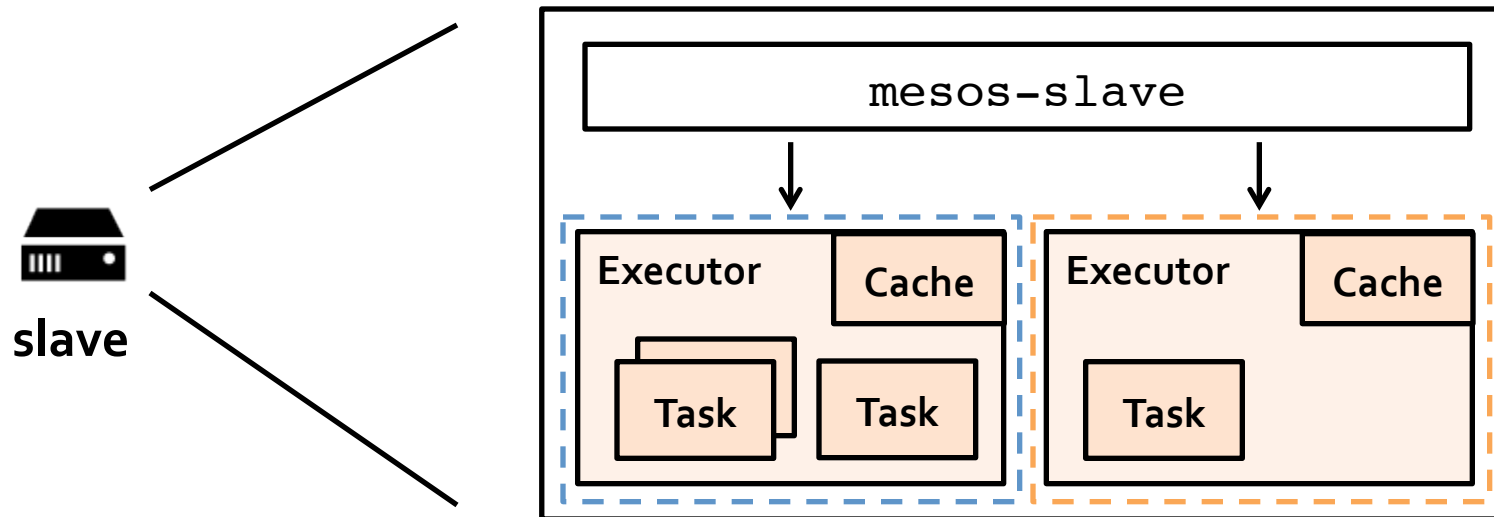Executor    Cache

Task    Task

Executor    Cache

Spark executors only consume memory, can share CPU between

# fine-grained sharing



Spark executors only consume memory, can share CPU between

# fine-grained sharing


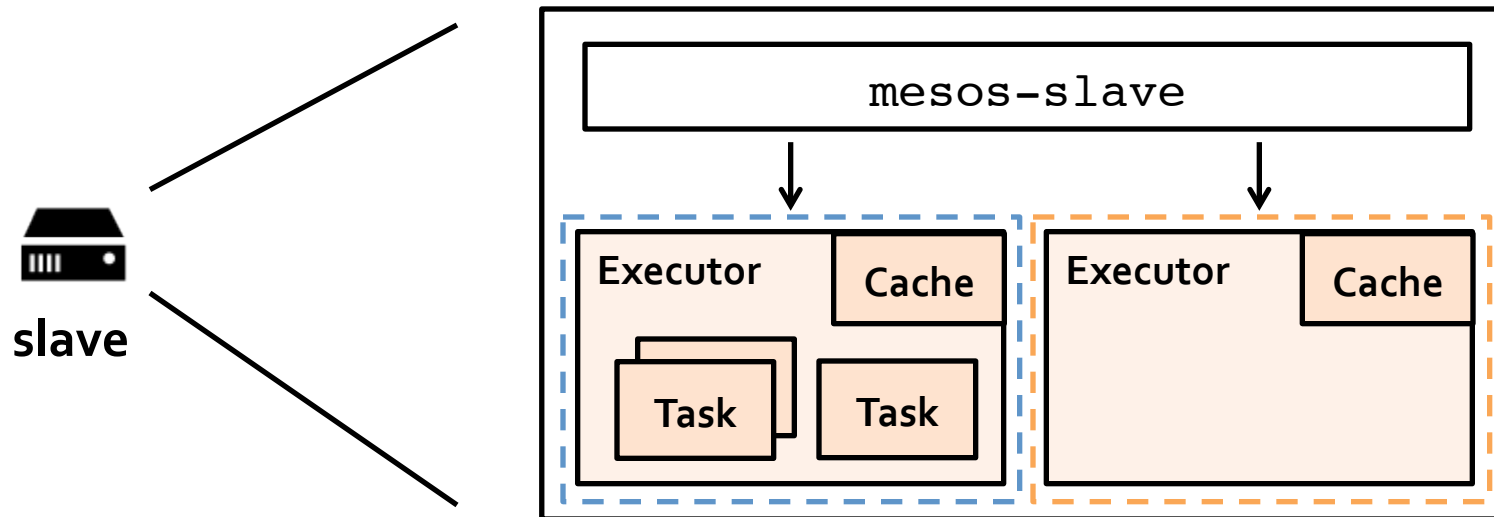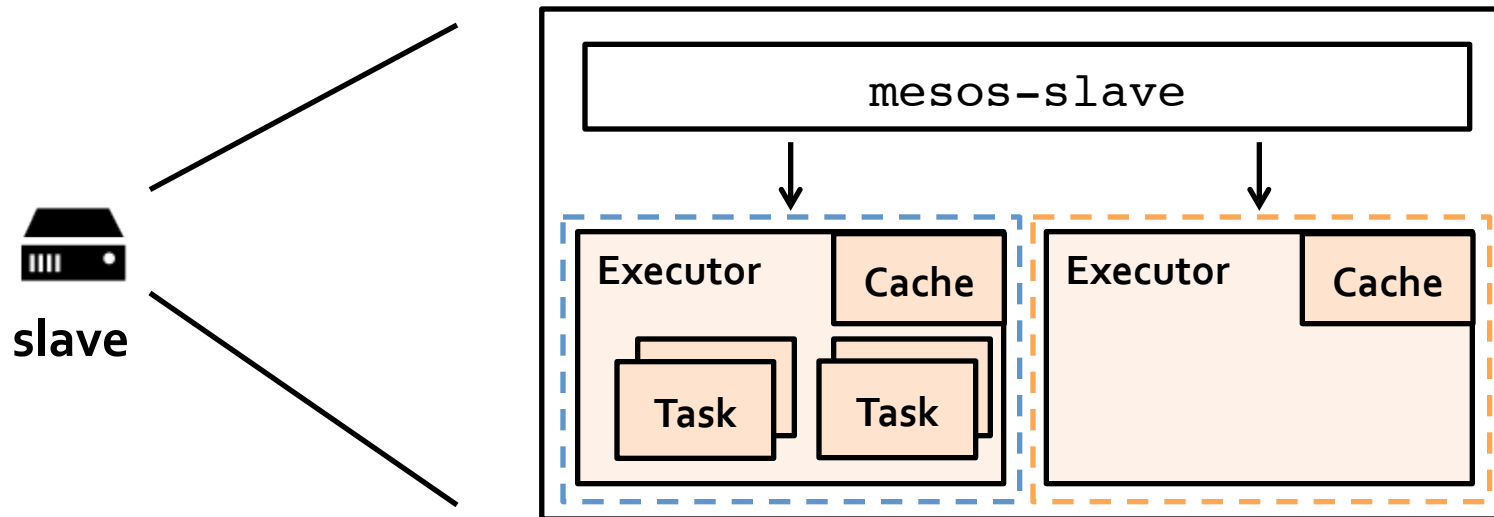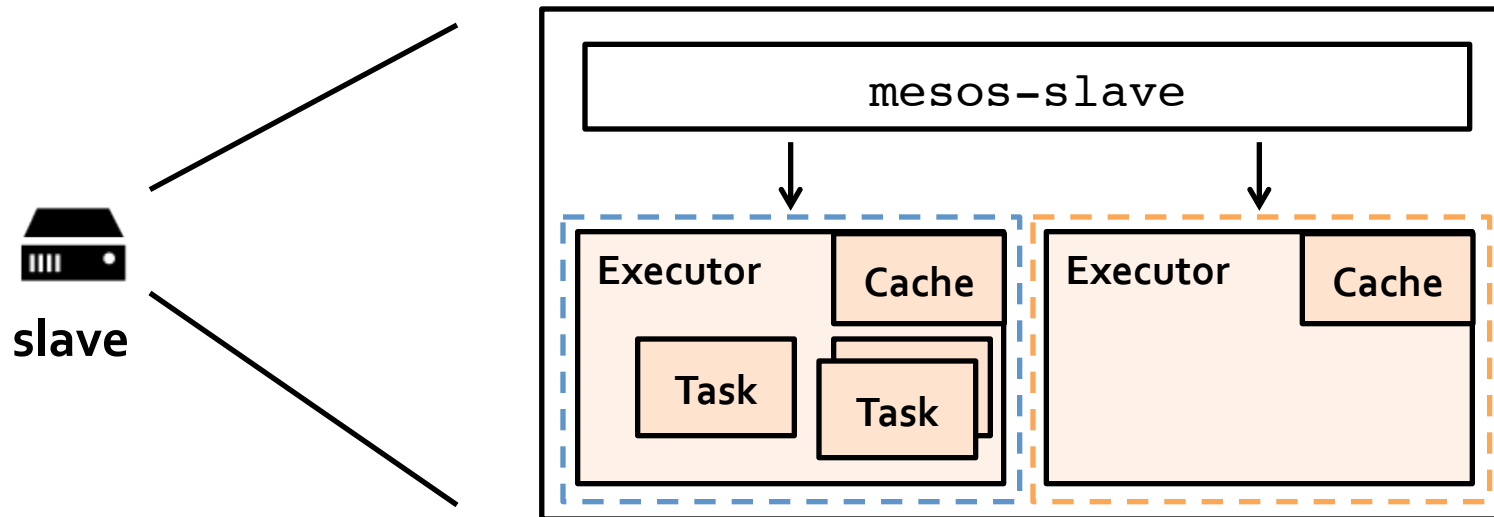
Spark executors only consume memory, can share CPU between

# fine-grained sharing



Spark executors only consume memory, can share CPU between
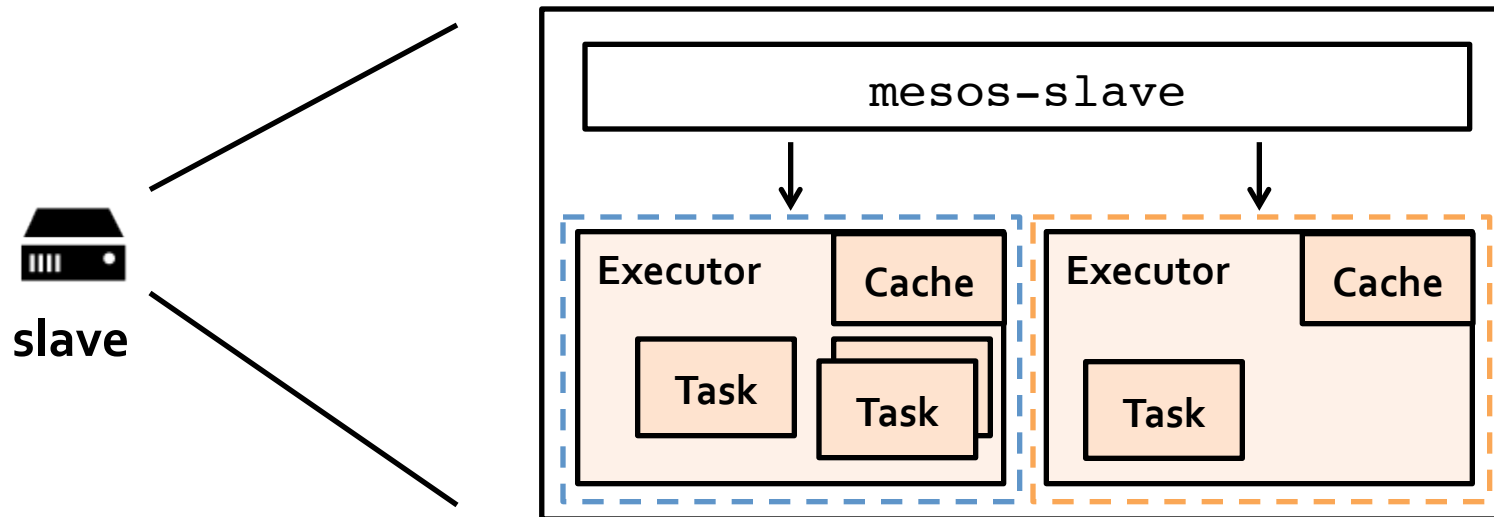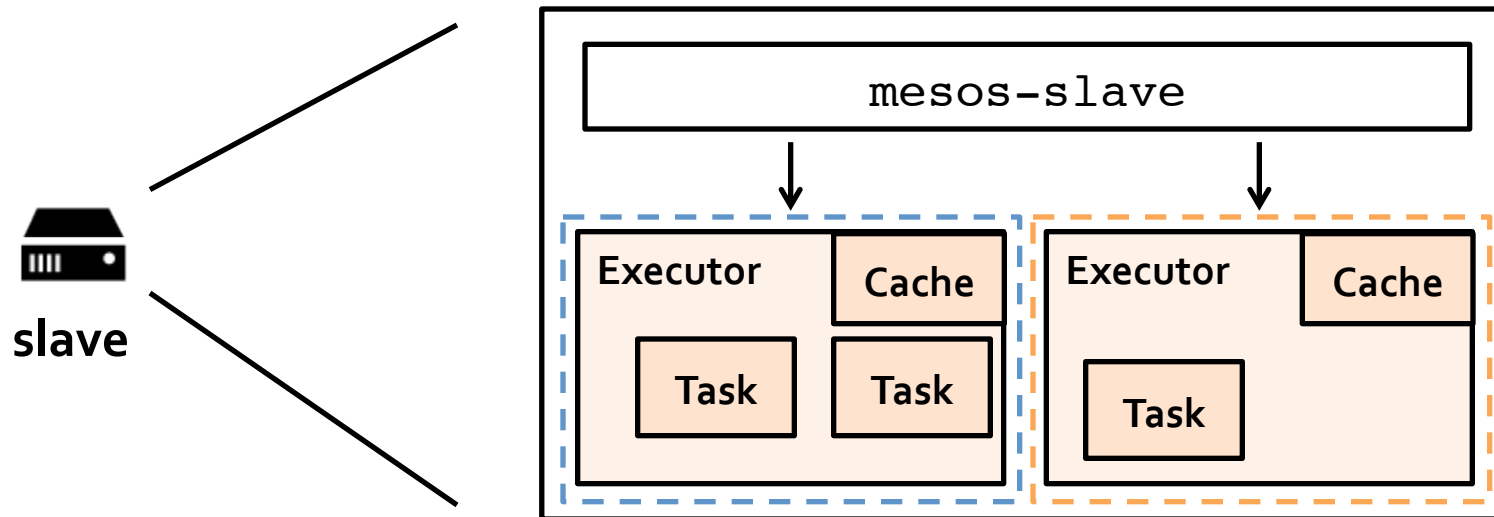
# fine-grained sharing

**slave**

```
mesos-slave
```

Executor · Cache · Task · Task

Executor · Cache · Task

Spark executors only consume memory, can share CPU between

# fine-grained sharing

slave



| mesos-slave |
|---|

Executor — Cache — Task — Task
Executor — Cache — Task — Task

Spark executors only consume memory, can share CPU between

# fine-grained sharing



**slave**

mesos-slave

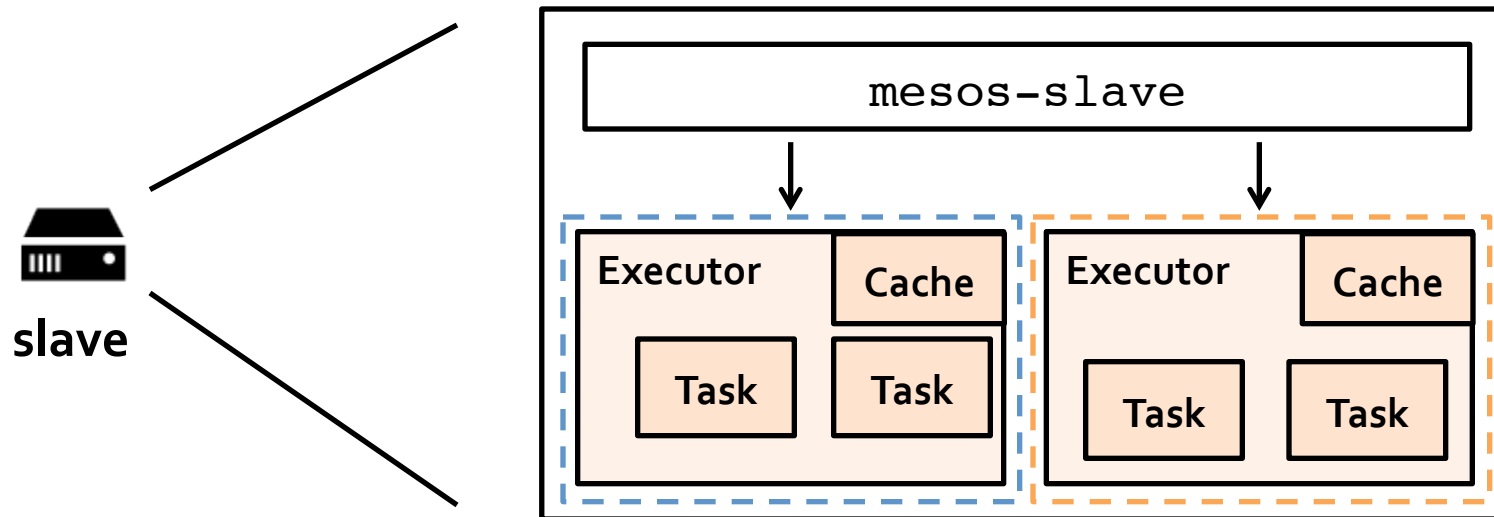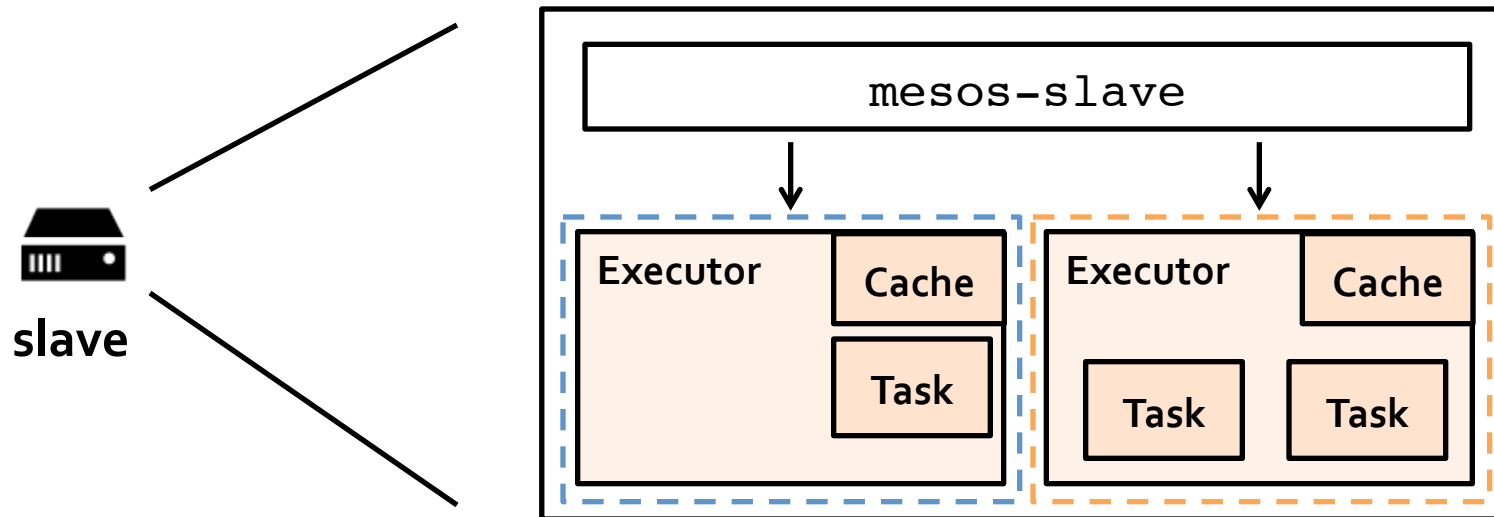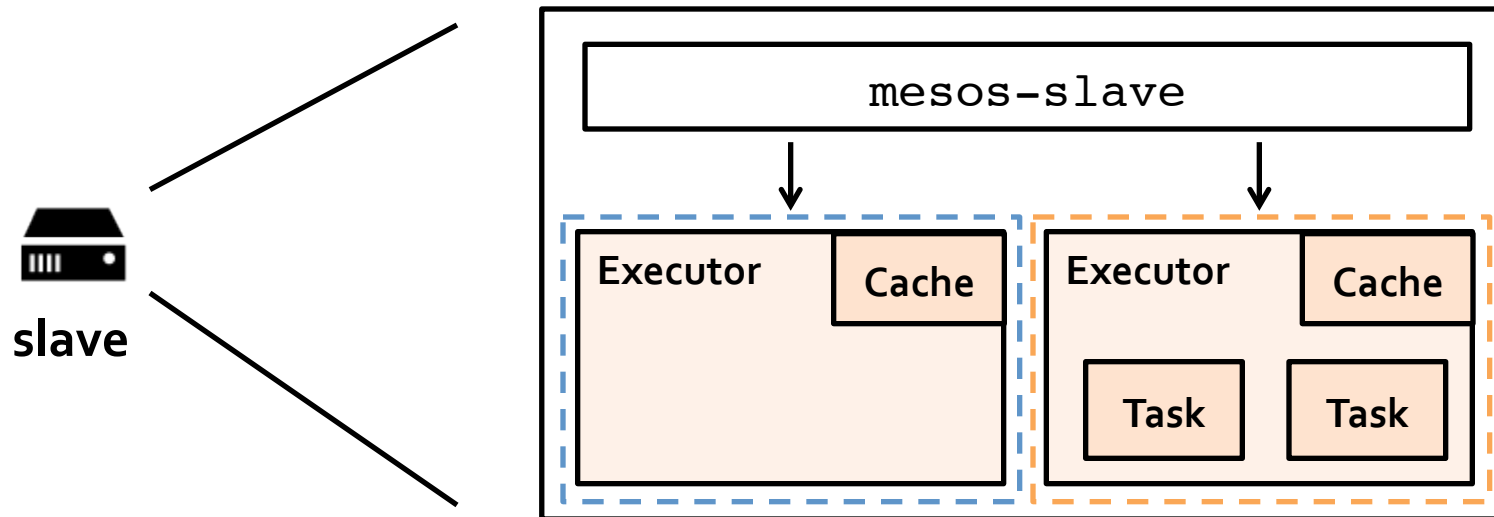| Executor | Cache |
| Task |

| Executor | Cache |
| Task | Task |

Spark executors only consume memory, can share CPU between

# fine-grained sharing



Spark executors only consume memory, can share CPU between

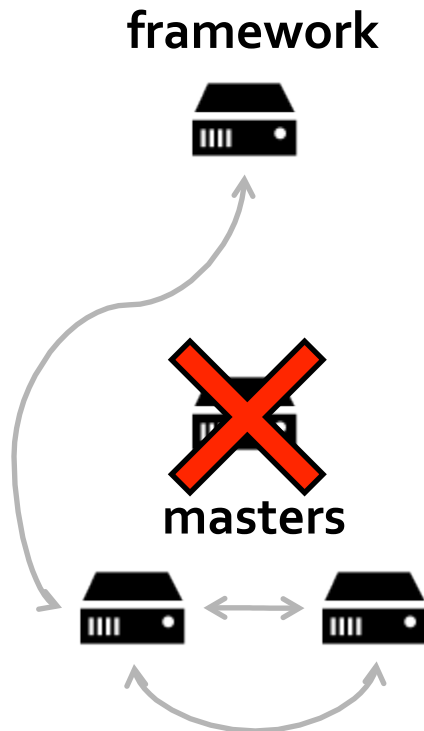# agenda

① Mesos

② Spark on Mesos

③ why Mesos?
   ① multi-tenancy
   ② fine-grained sharing
   ③ why not?  ⬅

④ long-lived services and other frameworks

# why not?

more moving pieces means more things to learn and more things that can fail …
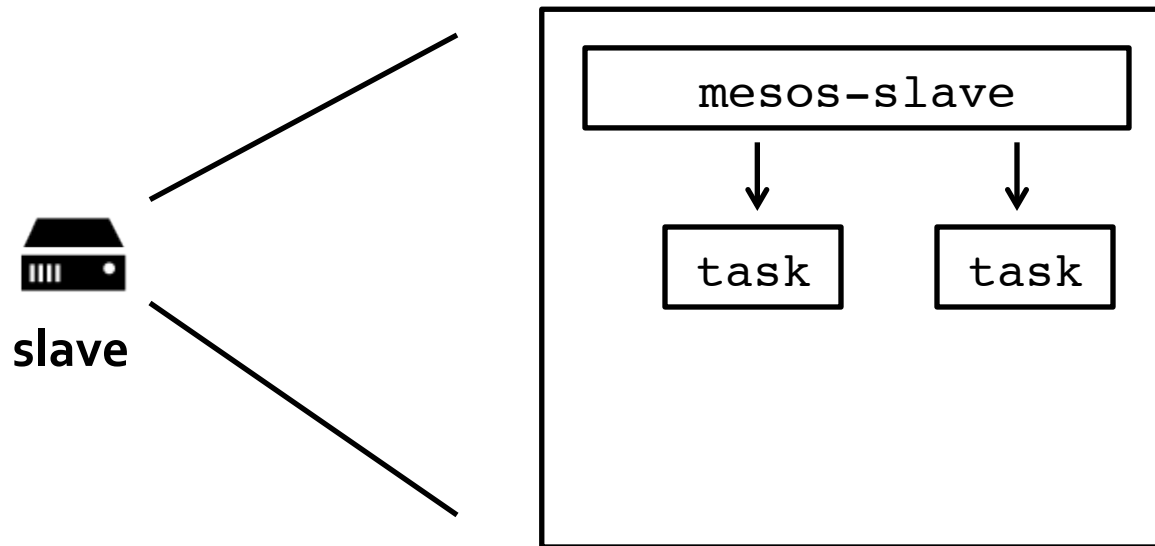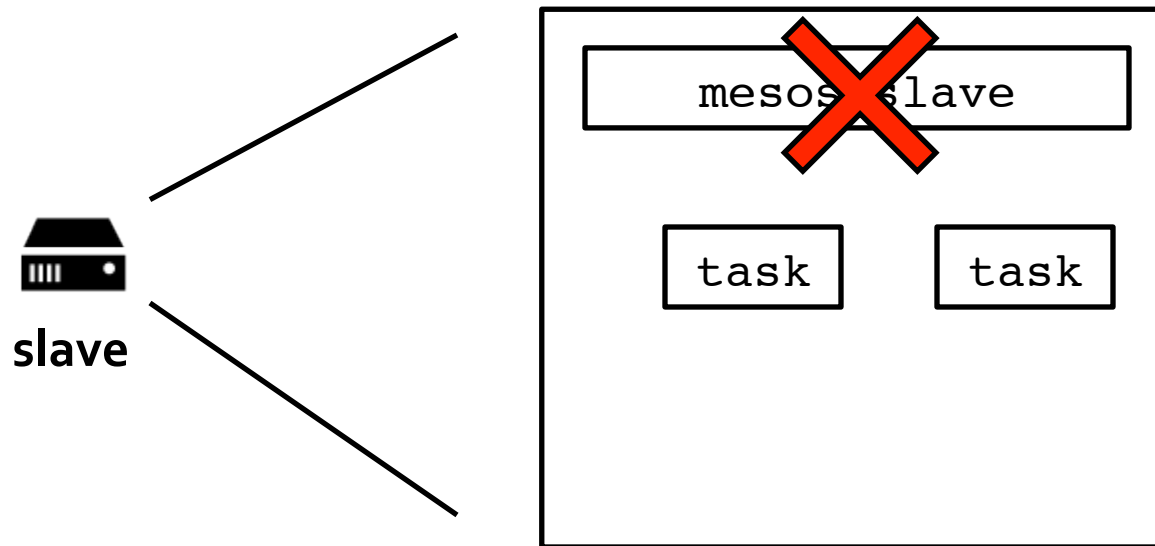
# master failover

**framework**

**masters**

after a new master is elected all frameworks and slaves connect to the new master

*all tasks keep running across master failover!*

# slave failover
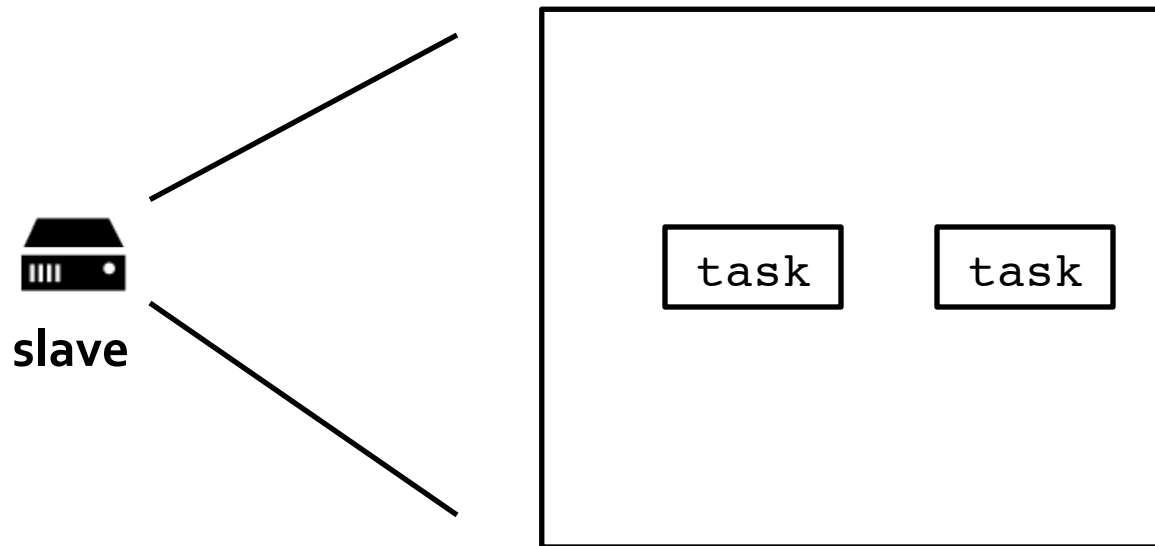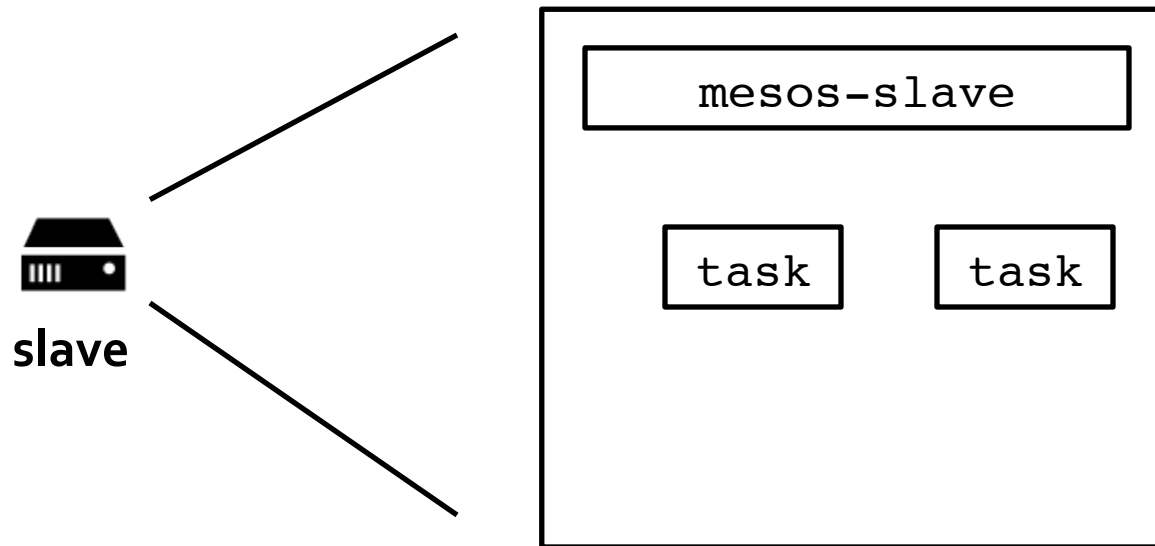
# slave failover

# slave failover



**slave**

task    task

# slave failover

# slave failover



slave

mesos-slave

task    task

# slave failover @twitter



**slave**

```
mesos-slave
```

(large in-memory services, expensive to restart)

# framework failover
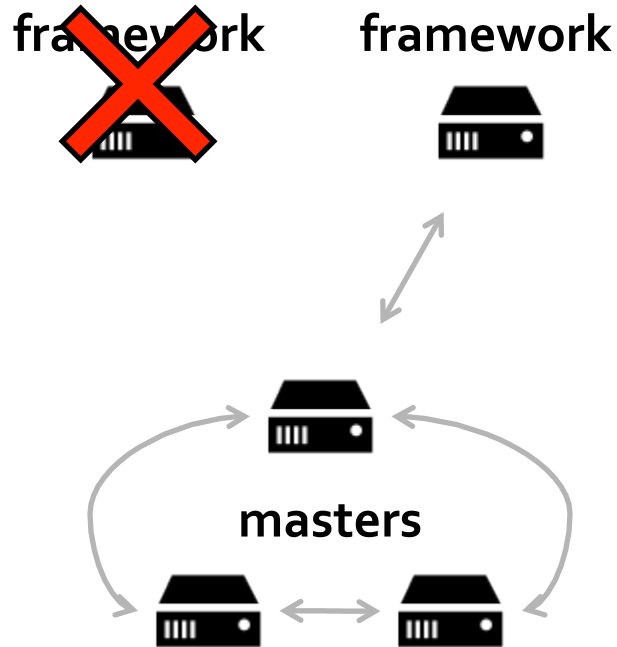
framework ~~framework~~        framework



framework re-registers with
master and resumes operation

*all tasks keep running across
framework failover!*

masters

# agenda

① Mesos

② Spark on Mesos

③ why Mesos?
    ① multi-tenancy
    ② fine-grained sharing
    ③ why not?

④ long-lived services and other frameworks  ⬅
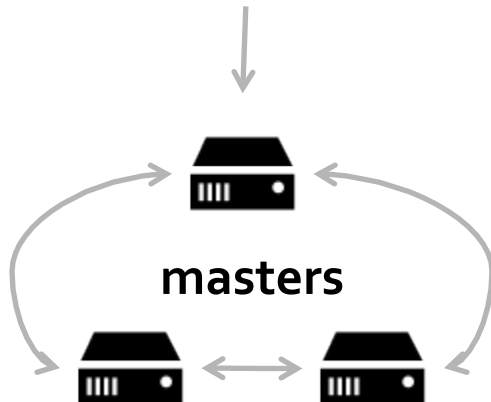
# Apache Aurora (incubating)



Aurora is a Mesos framework that makes it easy to launch services written in Ruby, Java, Scala, Python, Go, etc!

masters

# Marathon (from Mesosphere)

Marathon is a Mesos framework that makes it easy to launch services written in Ruby, Java, Scala, Python, Go, etc!

# Jenkins on Mesos

**ebay** tech blog  Where e-commerce meets world-class technology

## Delivering eBay's CI Solution with Apache Mesos – Part I
*by* THE EBAY PAAS TEAM *on* 04/04/2014
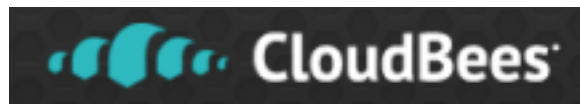*in* CLOUD,DATA INFRASTRUCTURE AND SERVICES,SOFTWARE ENGINEERING

LINKS

eBay Careers

eBay Developers Program

(http://bit.ly/1frLrLf)

**..((((.. CloudBees**

**Apache Mesos and Jenkins - elastic build slaves**

(http://bit.ly/1nHwM3r)

Elastic Mesos: elastic.mesosphere.io

Launch an Apache Mesos Cluster in 3 2 1

Packages:

| Apache Mesos 0.14.2 Release announcement | 04 Nov 2013 |
|---|---|
| Apache Mesos 0.14.2 for Ubuntu 13.04 (AMD 64) and Instructions | SHA 256 |
| Apache Mesos 0.14.2 for Ubuntu 12.10 (AMD 64) and Instructions | SHA 256 |
| Apache Mesos 0.14.2 for Ubuntu 12.04 (AMD 64) and Instructions | SHA 256 |
| Apache Mesos 0.14.2 for Debian 7 (AMD 64) and Instructions | SHA 256 |
| Apache Mesos 0.14.2 for CentOS 6 (x86_64) and Instructions | SHA 256 |
| Apache Mesos 0.14.2 for Red Hat 6 (x86_64) and Instructions | SHA 256 |

# Thank You!

mesos.apache.org

mesos.apache.org/blog

@ApacheMesos