



Fault Tolerant Infrastructure

Building Systems with etcd

@coreoslinux
@brandonphilips

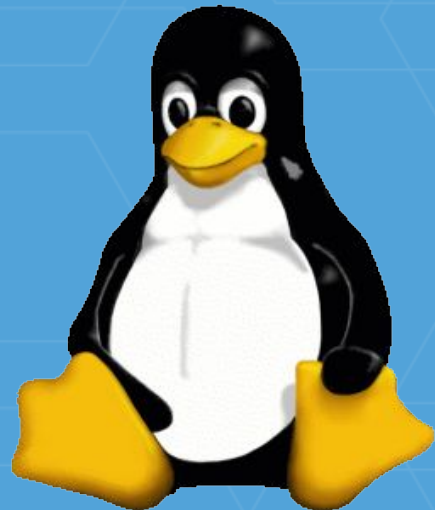
Brandon Philips

CTO, CoreOS

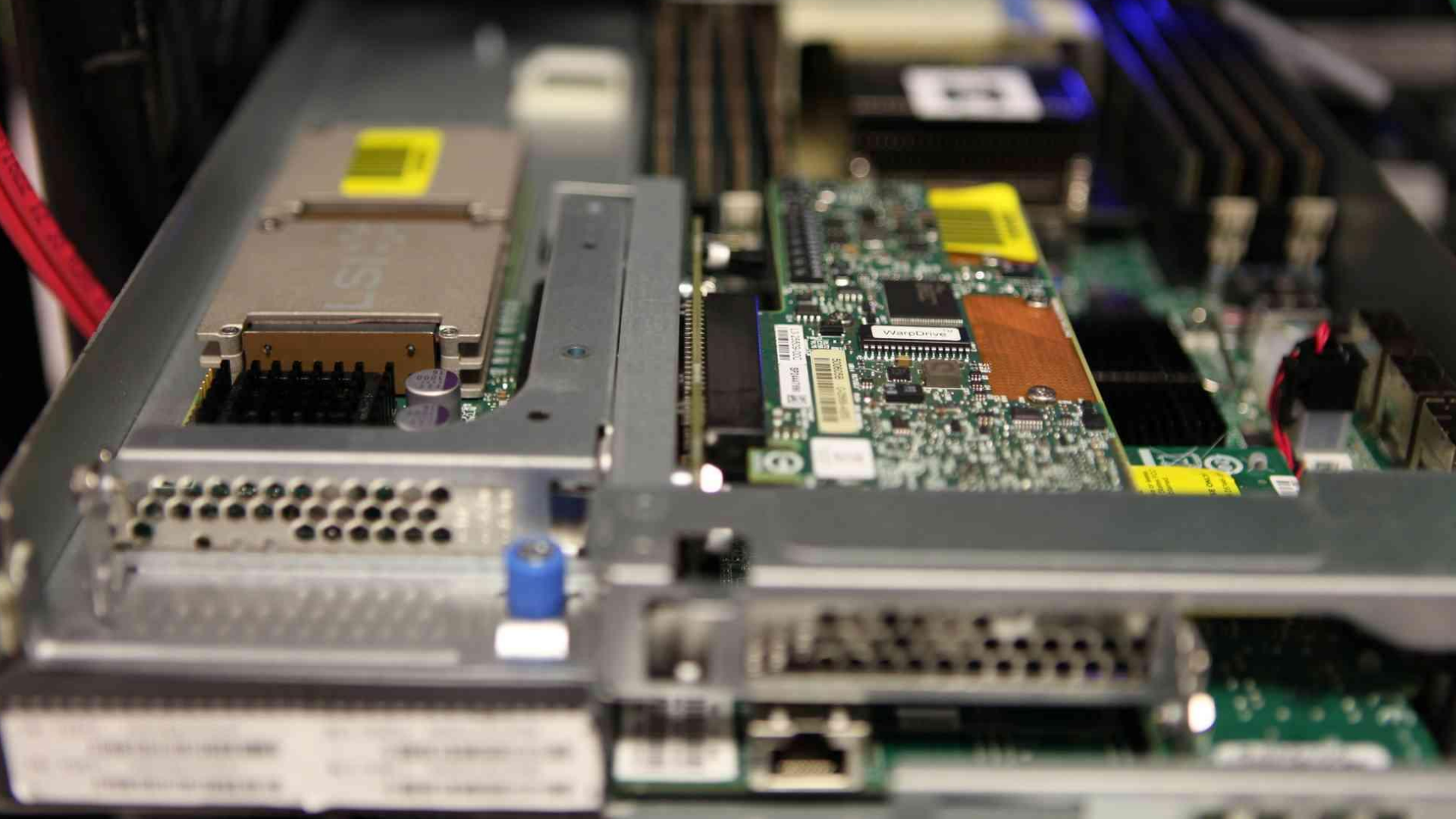
github.com/philips

What is CoreOS?

What is CoreOS?







What is CoreOS?

 etcd

 flannel

 rkt



Top 100 Community Contributors



10k+ total commits
~30 repositories
324 releases
22 CoreOS contributors
605 community contributors

Top 100 Community Contributors





The smartest way to run your container infrastructure.

tectonic.com [@tectonic](https://twitter.com/tectonic)

QUAY

Secure hosting for private Docker repositories

quay.io @quayio

Why build CoreOS?

The Datacenter as a Computer
*An Introduction to the Design of
Warehouse-Scale Machines*

Luiz André Barroso and Urs Hölzle
Google Inc.



you

you as a sw engineer

your

```
with Ada.Text_IO;
```

```
procedure Hello_World is
```

```
  use Ada.Text_IO;
```

```
begin
```

```
  Put_Line("Hello, world!");
```

```
end;
```

```
#include <stdio.h>
```

```
int main()
```

```
{
```

```
  printf("Hello, world!\n");
```

```
}
```

```
package main
```

```
import "fmt"
```

```
func main() {
```

```
  fmt.Println("Hello, world!")
```

```
}
```

your



**container
image**

your



**/bin/java
/opt/app.jar
/lib/libc**

your



```
/bin/python  
/opt/app.py  
/lib/libc
```

your



com.example.app

d474e8c57737625c

your

Signed By: Alice

d474e8c57737625c

you as an ops engineer

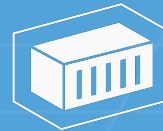
your



`com.example.webapp`
x3



your



`com.example.webapp`
x3



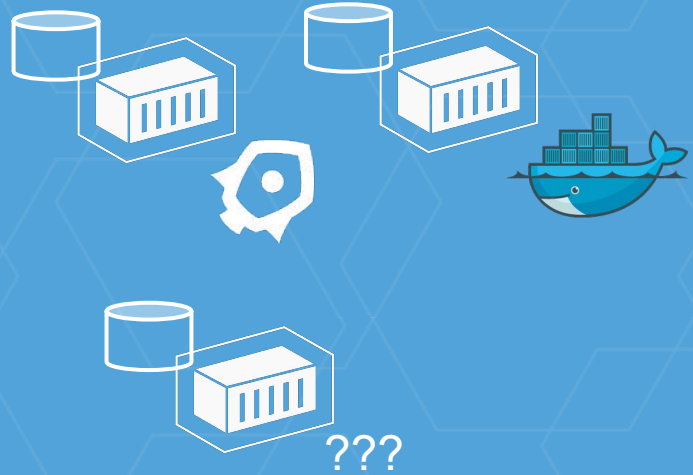
your




com.example.webapp
x3



your

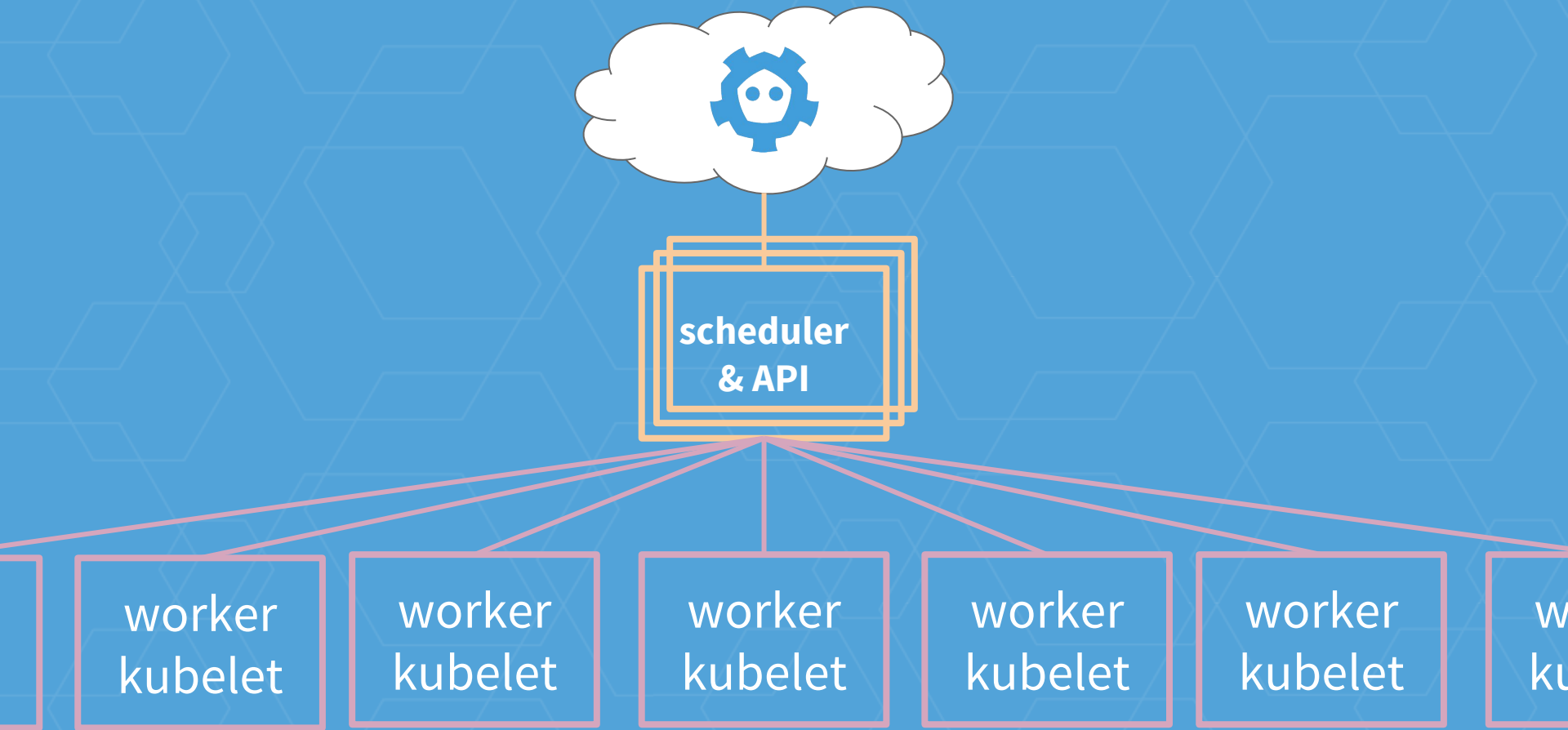


How do we do it?

The background is a solid blue color with a subtle, repeating pattern of light blue hexagons of varying sizes, creating a textured, geometric effect.

cluster operations

architecture in practice



OS operations machine configuration



cluster operations

distributed configuration

github.com/philips/hacks/tree/master/etcd-demos



etcd

/etc

distributed



open source software

failure tolerant

durable

watchable

exposed via HTTP

runtime reconfigurable

Data Store API

-X GET

Get Wait

-X PUT

Put Create CAS

-X DELETE

Delete CAD



etcd basics

clusters

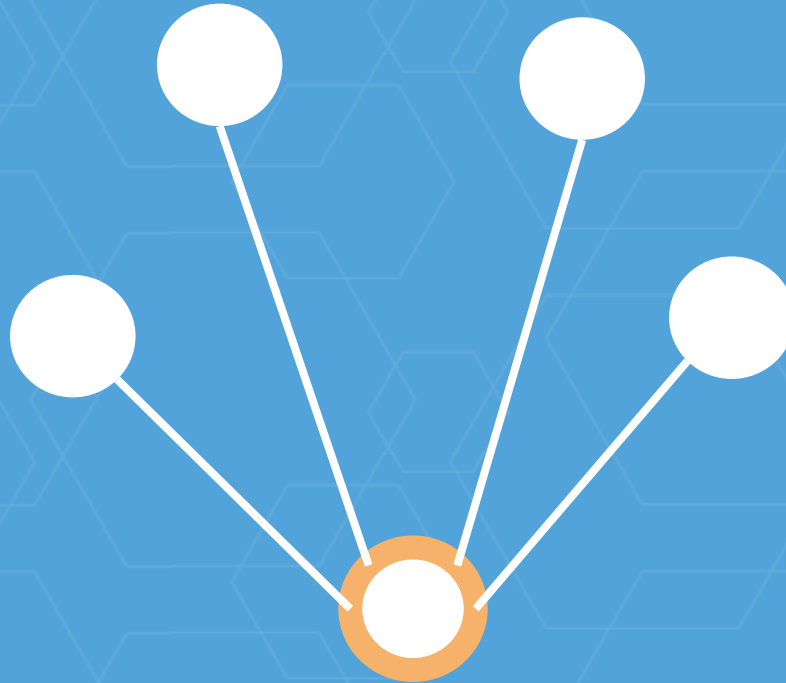
Typical Cluster



Leader



Follower



etcd basics

API



etcd basics

fault tolerance

Available



Leader



Follower



Available



Leader



Follower



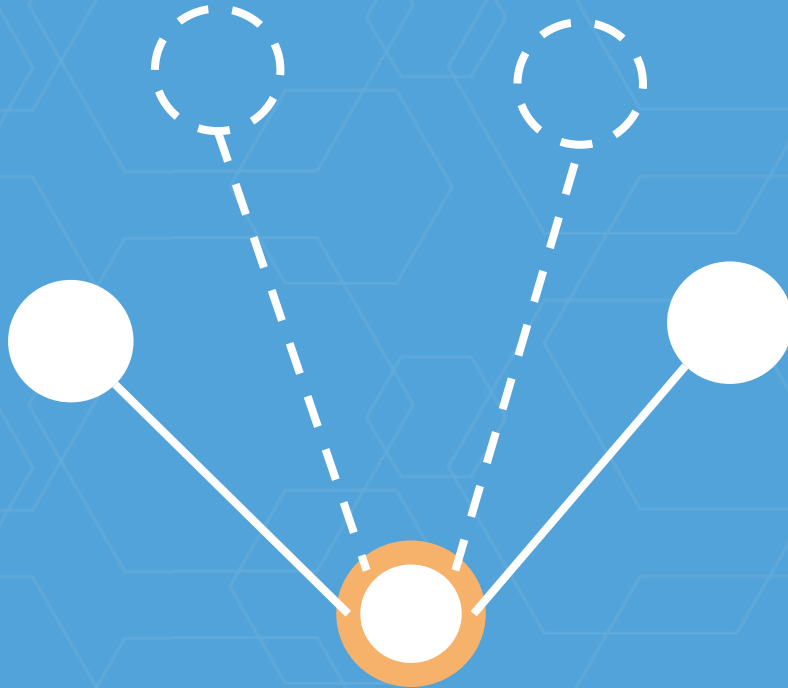
Available



Leader



Follower



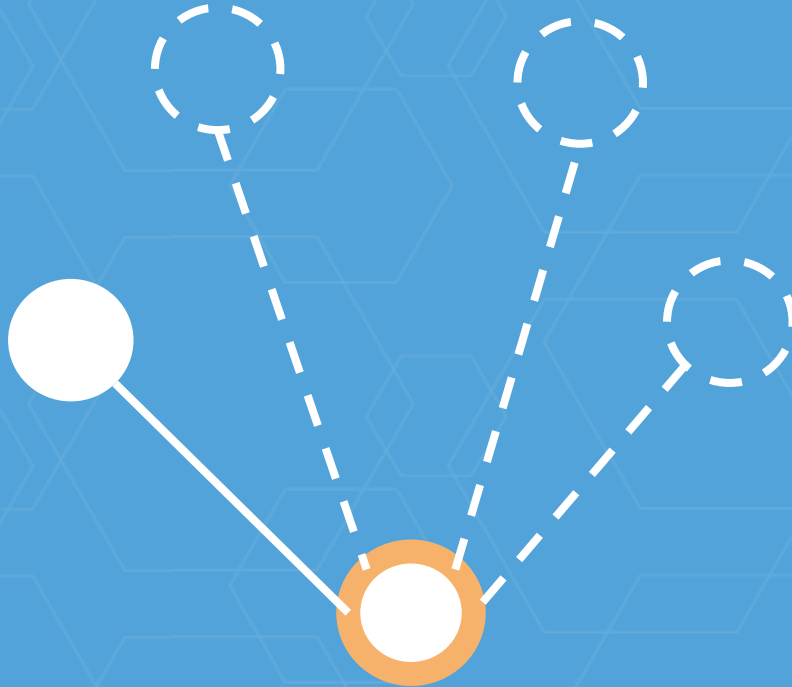
Unavailable



Leader



Follower





etcd basics

leader fault tolerance

Available



Leader



Follower



Available



Leader



Follower



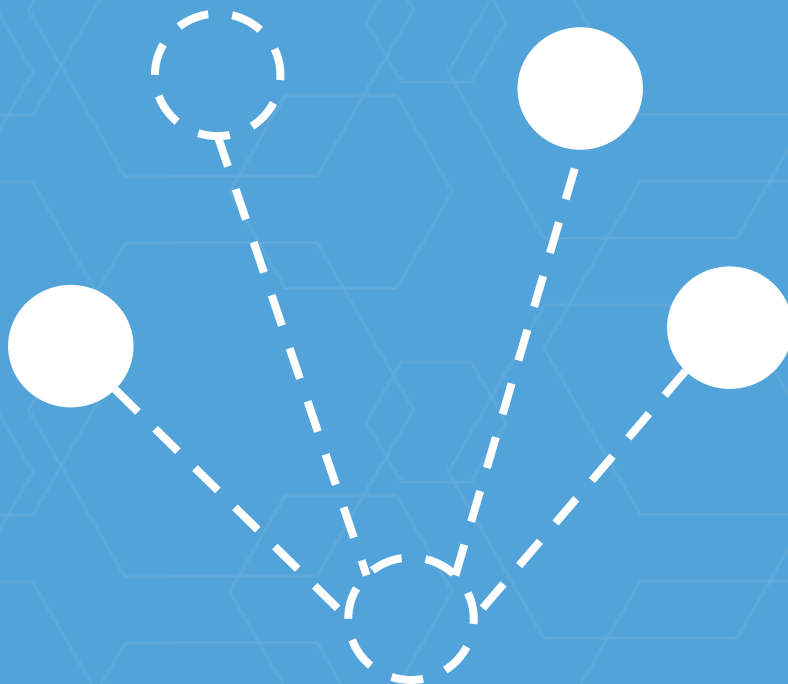
Temporarily Unavailable



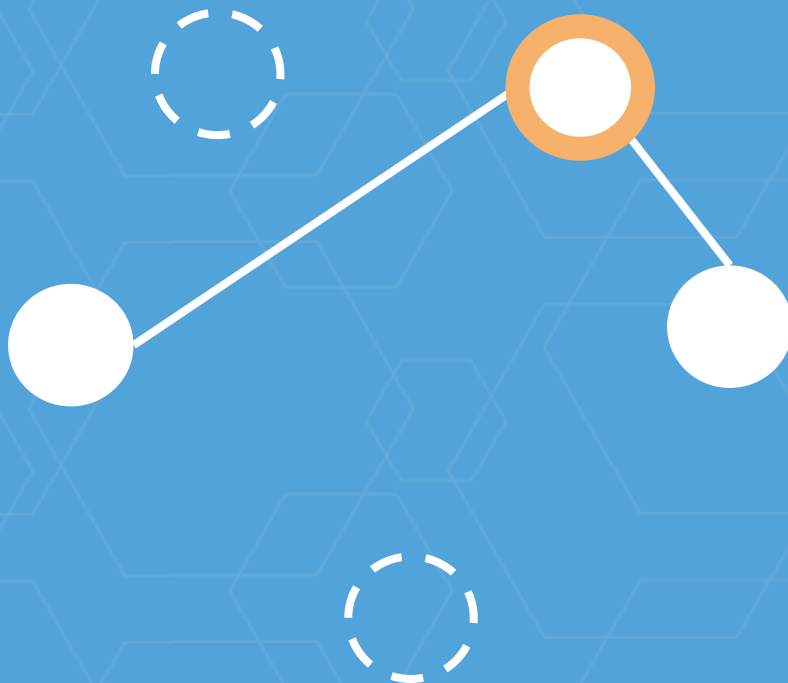
Leader



Follower



Available



Unavailable



etcd durability

wal, snapshots, backups

etcd bootstrap

discovery, static

```
$ curl discovery.etcd.io/new?size=5  
discovery.etcd.io/6eadeac2
```

discovery



discovery



discovery



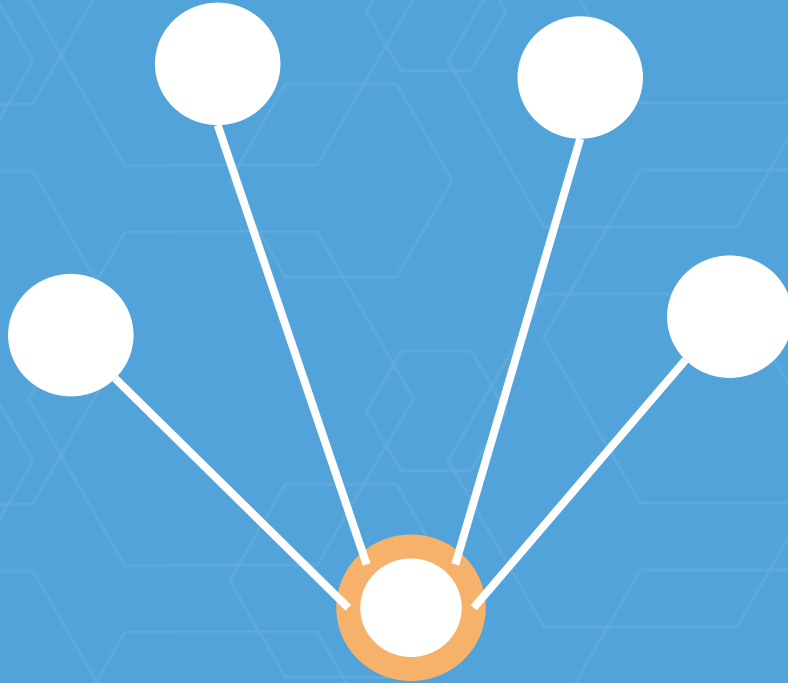
discovery



Leader

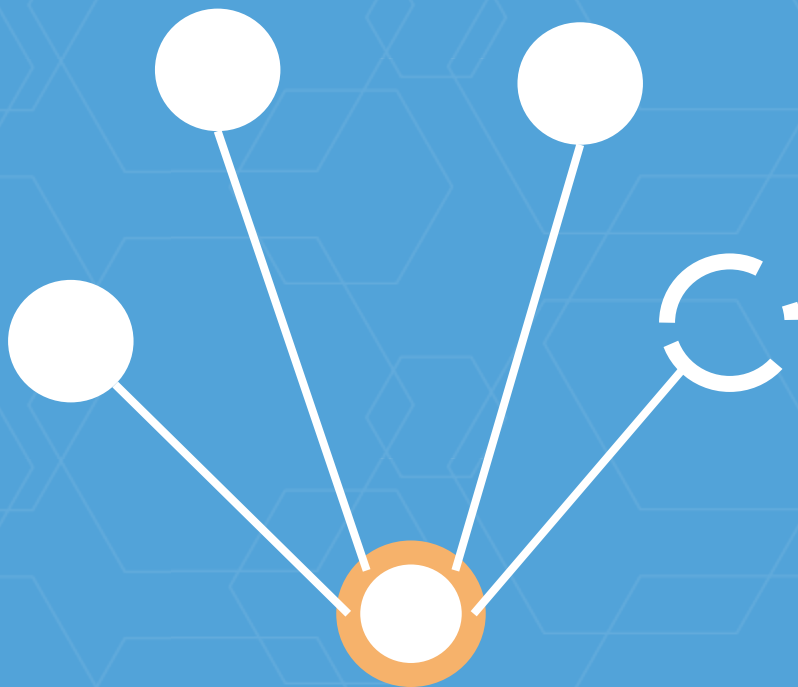


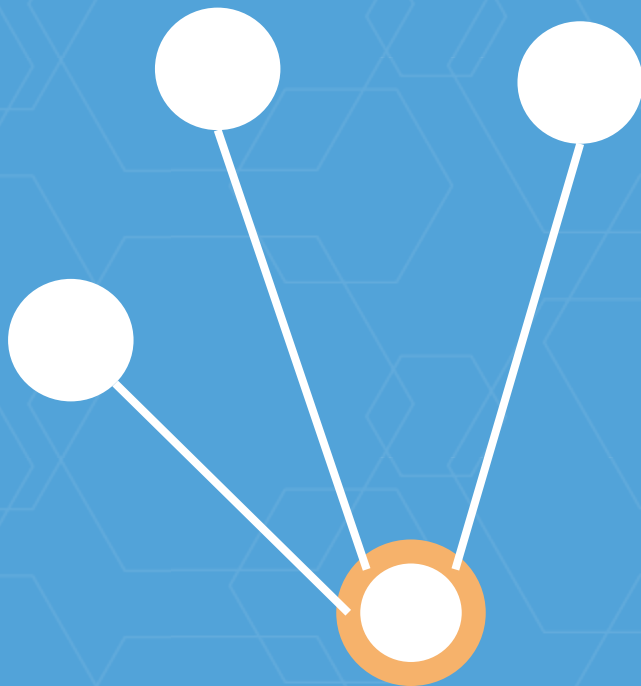
Follower

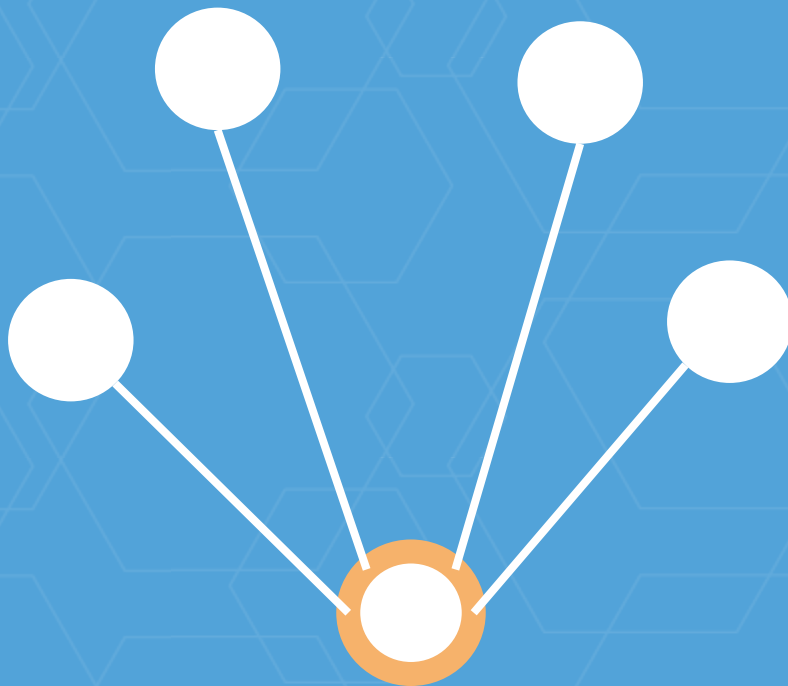


etcd reconfig

live addition and removal



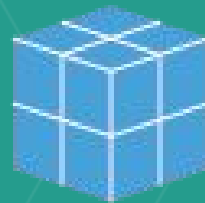




etcd apps

etcd apps
reboot locksmith

Data



Update

A

B



Data

A

B



Cluster Wide Reboot Lock

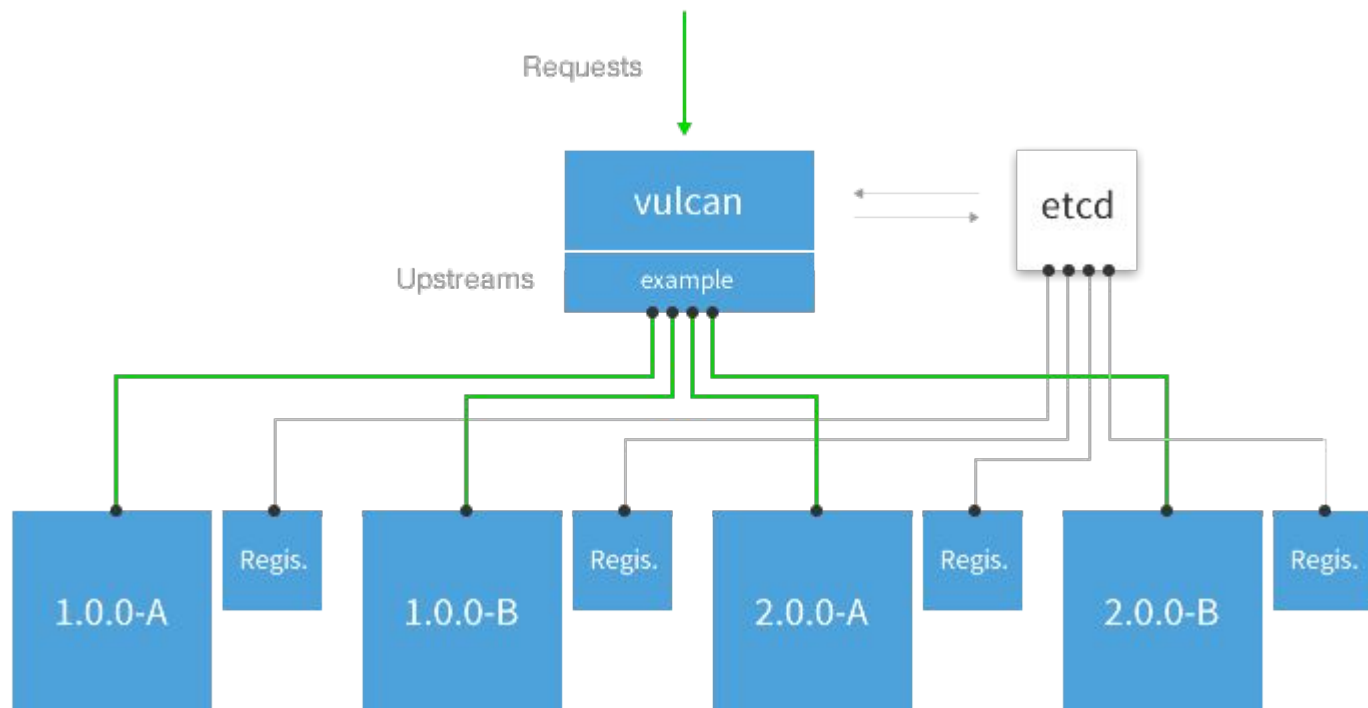
- Need to reboot? Decrement the semaphore key atomically with etcd.
- `manager.Reboot()` and wait...
- After reboot increment the semaphore key in etcd atomically.



etcd apps
skydns

The background is a solid blue color with a pattern of light blue hexagons of varying sizes, some overlapping, creating a geometric texture.

etcd apps
vulcand





etcd apps
confd

The background is a solid blue color with a pattern of light blue, semi-transparent hexagons of varying sizes and orientations, creating a geometric, honeycomb-like texture.

kubernetes
pulling it together

scheduling

k8s/mesos/etc scheduler

scheduling
getting work to servers

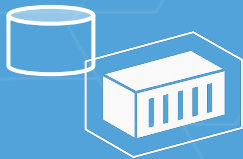
```
$ scp app host:/opt
```

```
$ ssh host systemd-run /opt/app
```



```
$ scp app host:/opt
```

```
$ ssh host systemd-run /opt/app
```



\$ fab deploy:app



\$ fab deploy:app



\$ fab deploy:app



\$ fab deploy:collector-app



\$ fab deploy:collector-app



\$ fab deploy:collector-app



\$ fab deploy deploy:collector-app



\$ fab lowest-loadaverage



\$ fab lowest-loadaverage
host1



```
$ fab lowest-loadaverage
```

```
host1
```

```
$ fab -H host1 deploy:job
```



You



Scheduler API



Scheduler



Machine(s)


```
while true {  
  todo = diff(desState, curState)  
  schedule(todo)  
}
```

```
while true {  
  todo = diff(desState, curState)  
  schedule(todo)  
}
```

```
while true {  
  todo = diff(desState, curState)  
  schedule(todo)  
}
```

```
while true {  
  todo = diff(desState, curState)  
  schedule(todo)  
}
```

services
dns, LBs, k8s labels



k8s labels

flexible service discovery

pod
env=dev
app=web

pod
env=test
app=web

pod
env=prod
app=web

service test.example.com
select(env=dev,app=web)

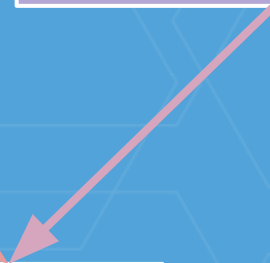
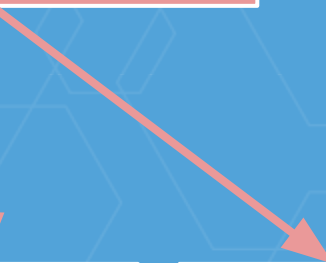
service beta.example.com
select(env=test,app=web)
OR
select(env=prod,app=web)

service example.com
select(env=prod,app=web)

pod
env=dev
app=web

pod
env=test
app=web

pod
env=prod
app=web



github.com/coreos/coreos-kubernetes



etcd.ngrok.io



scheduler
& API

worker
kubelet

worker
kubelet



worker &
API

works on 1 node too

The background is a solid blue color with a pattern of light blue, semi-transparent hexagons of varying sizes and orientations, creating a geometric, honeycomb-like texture.

work with us

coreos.com/careers

thank you

@coreoslinux

@tectonicstack

@brandonphilips