

Kubernetes Architecture



Nigel Poulton

@nigelpoulton www.nigelpoulton.com







Module Outline

Big picture view

Masters

Nodes

Pods

Services

Deployments

Recap



Kubernetes

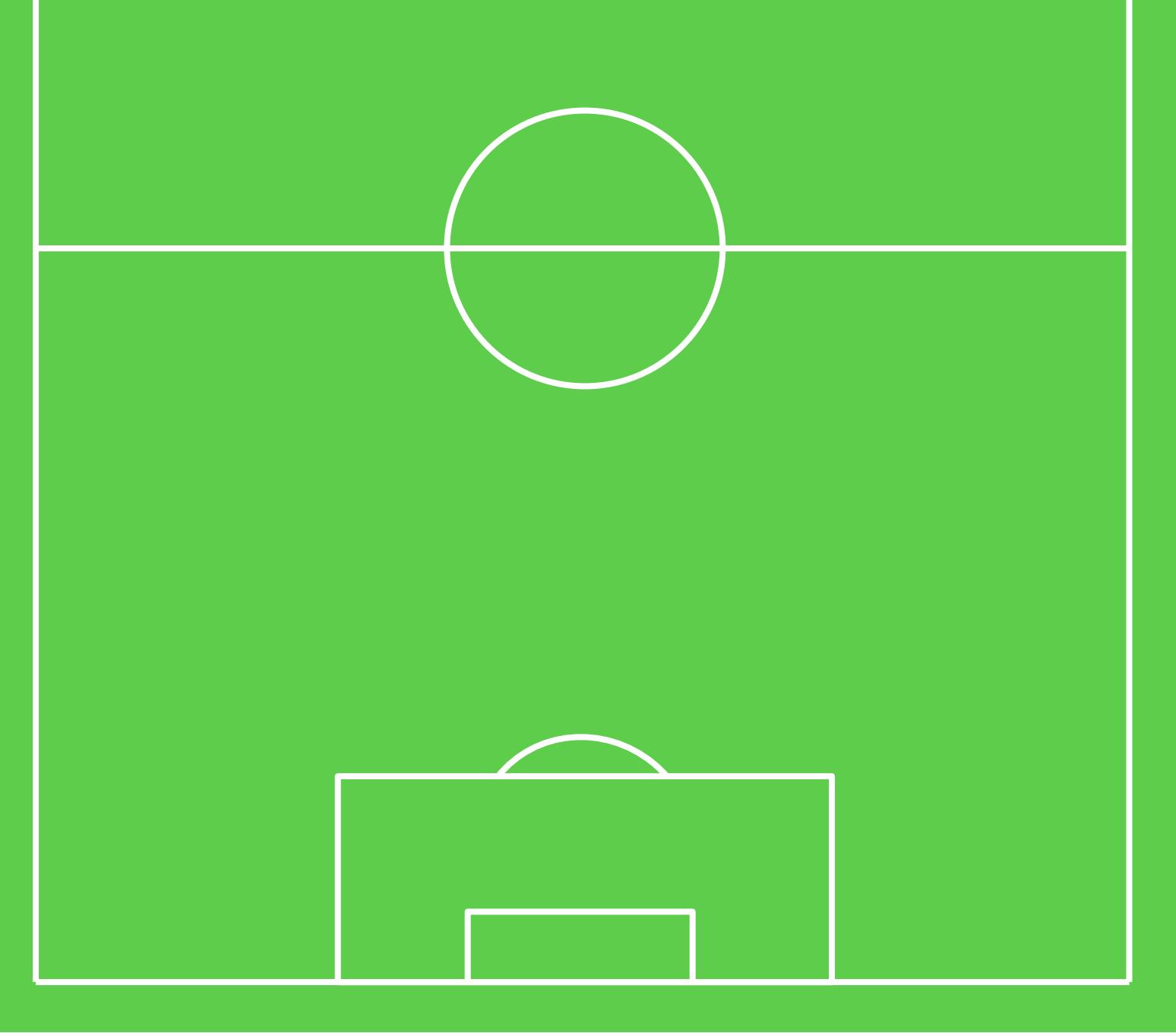
Big Picture View



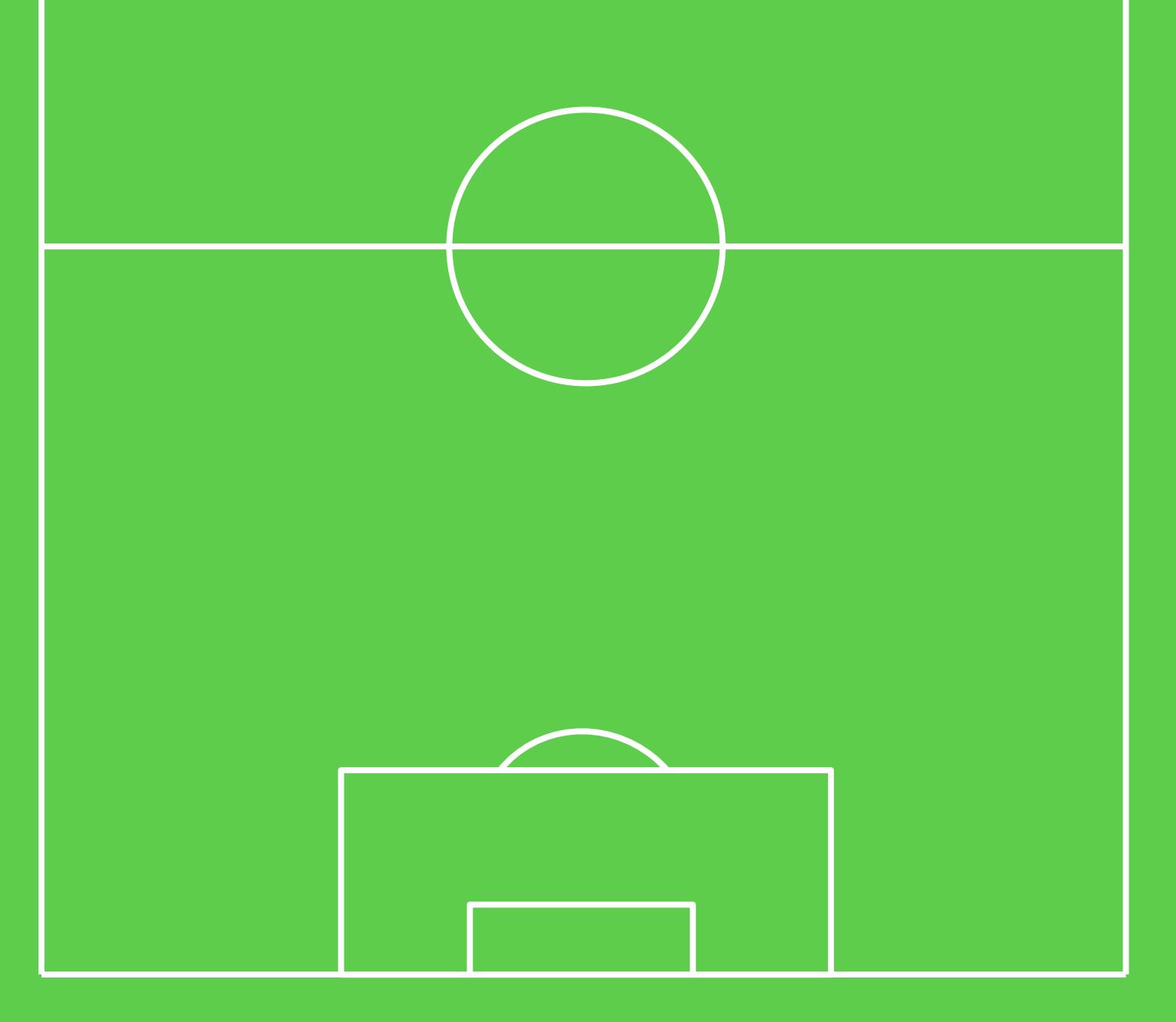
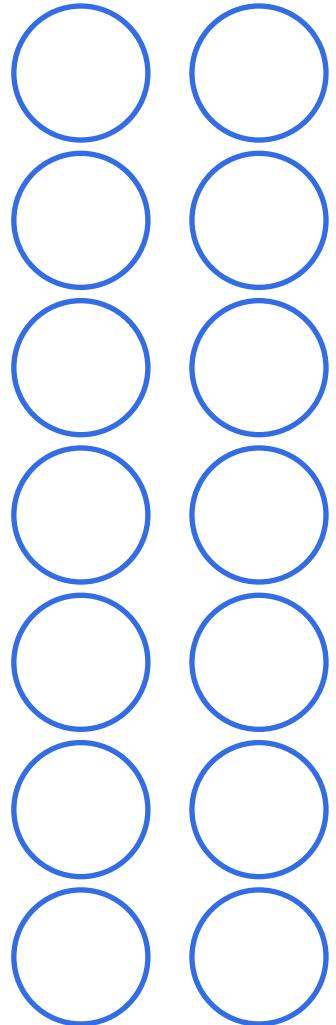
Team



Manager
(coach)



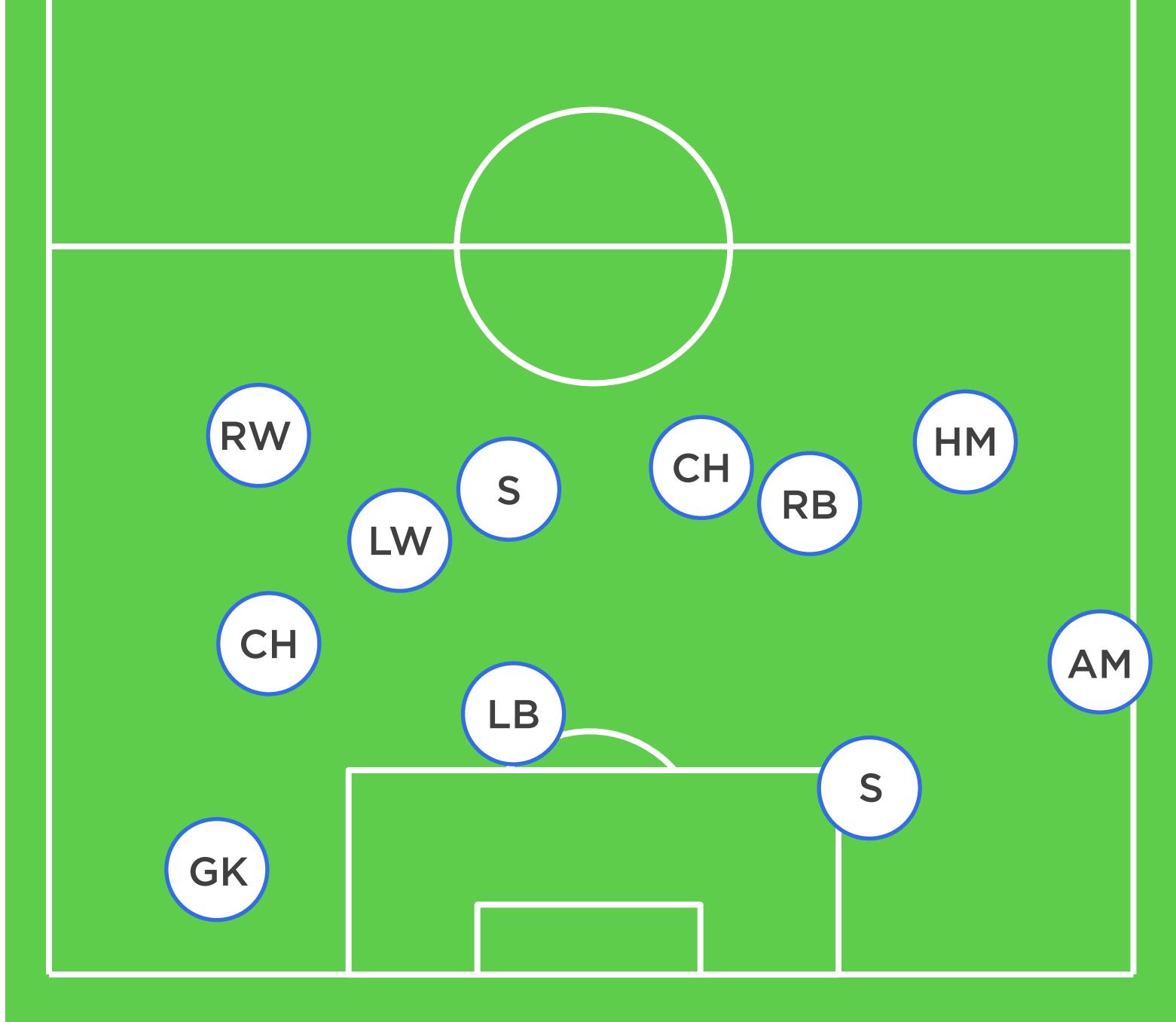
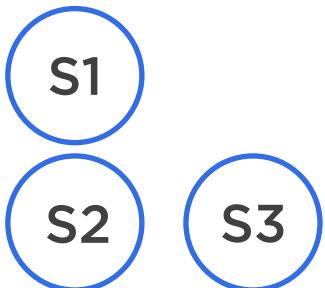
Team

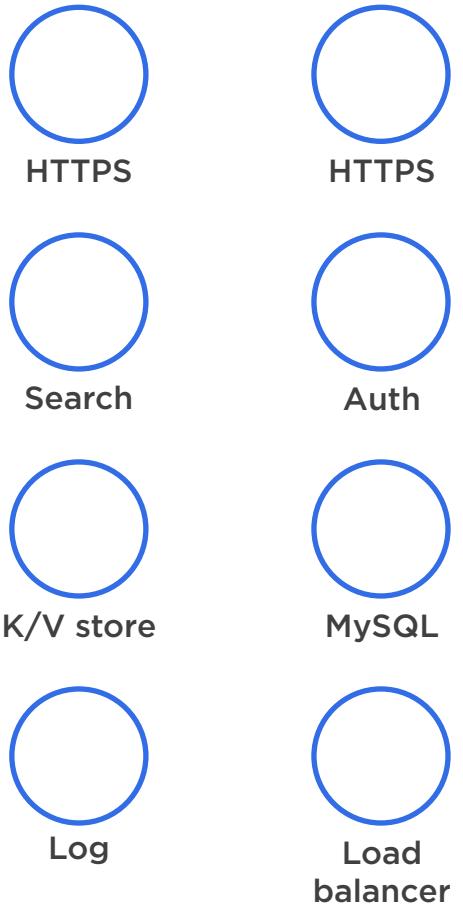


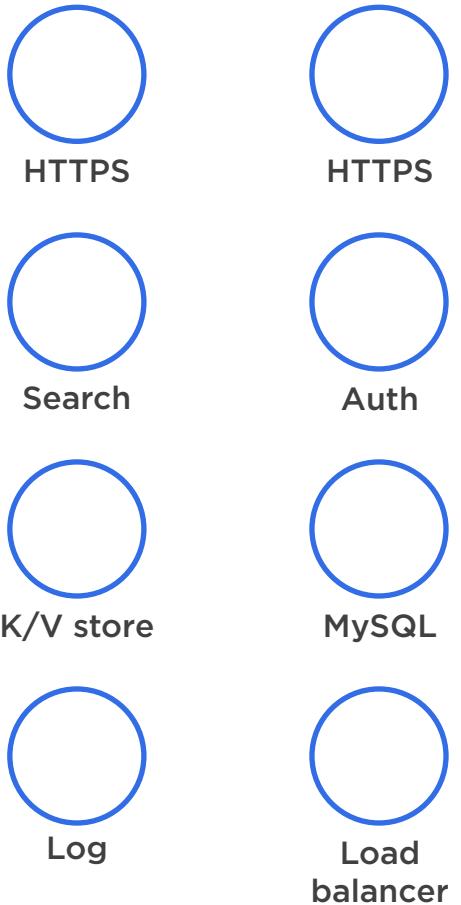
Team



Manager
(coach)

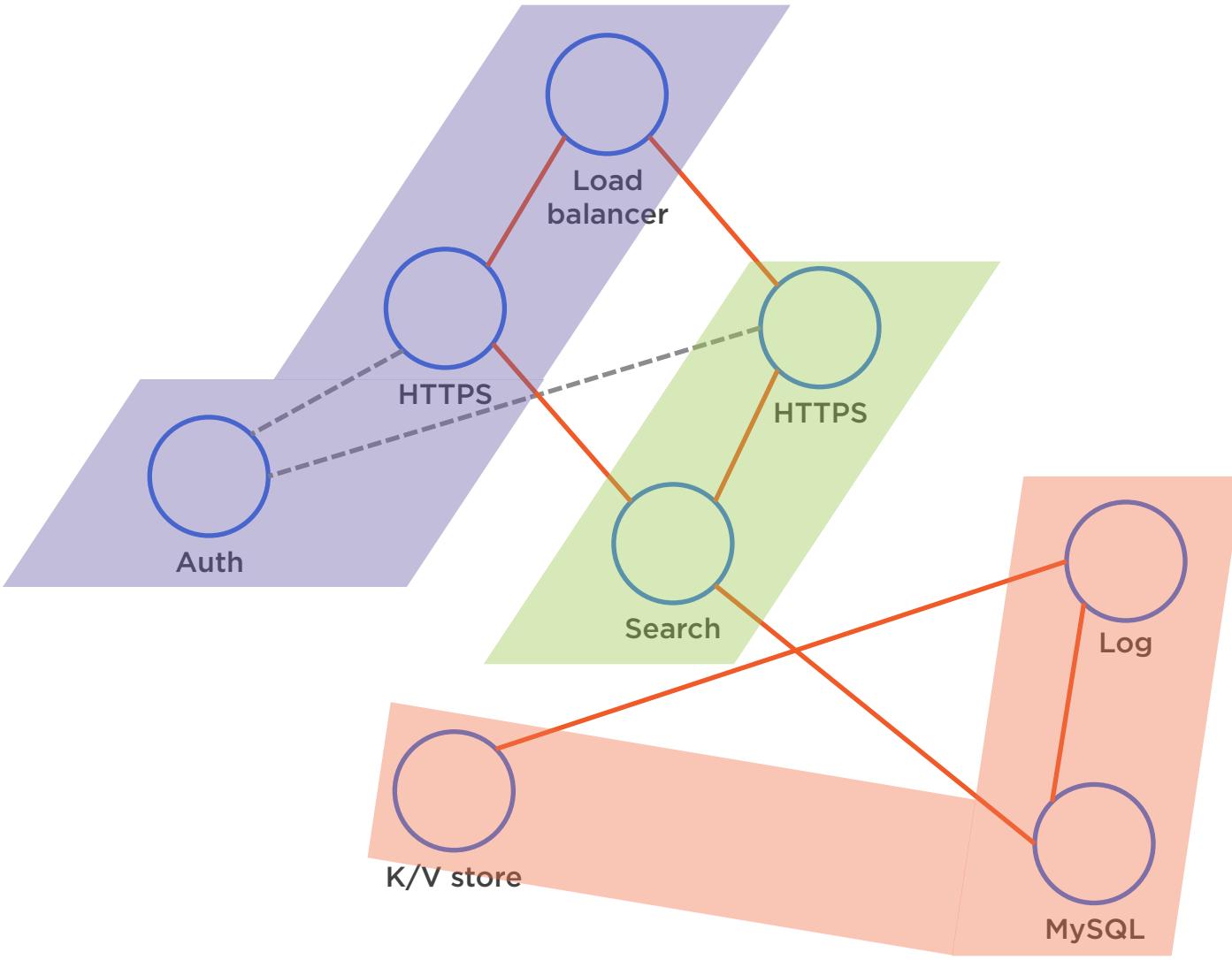


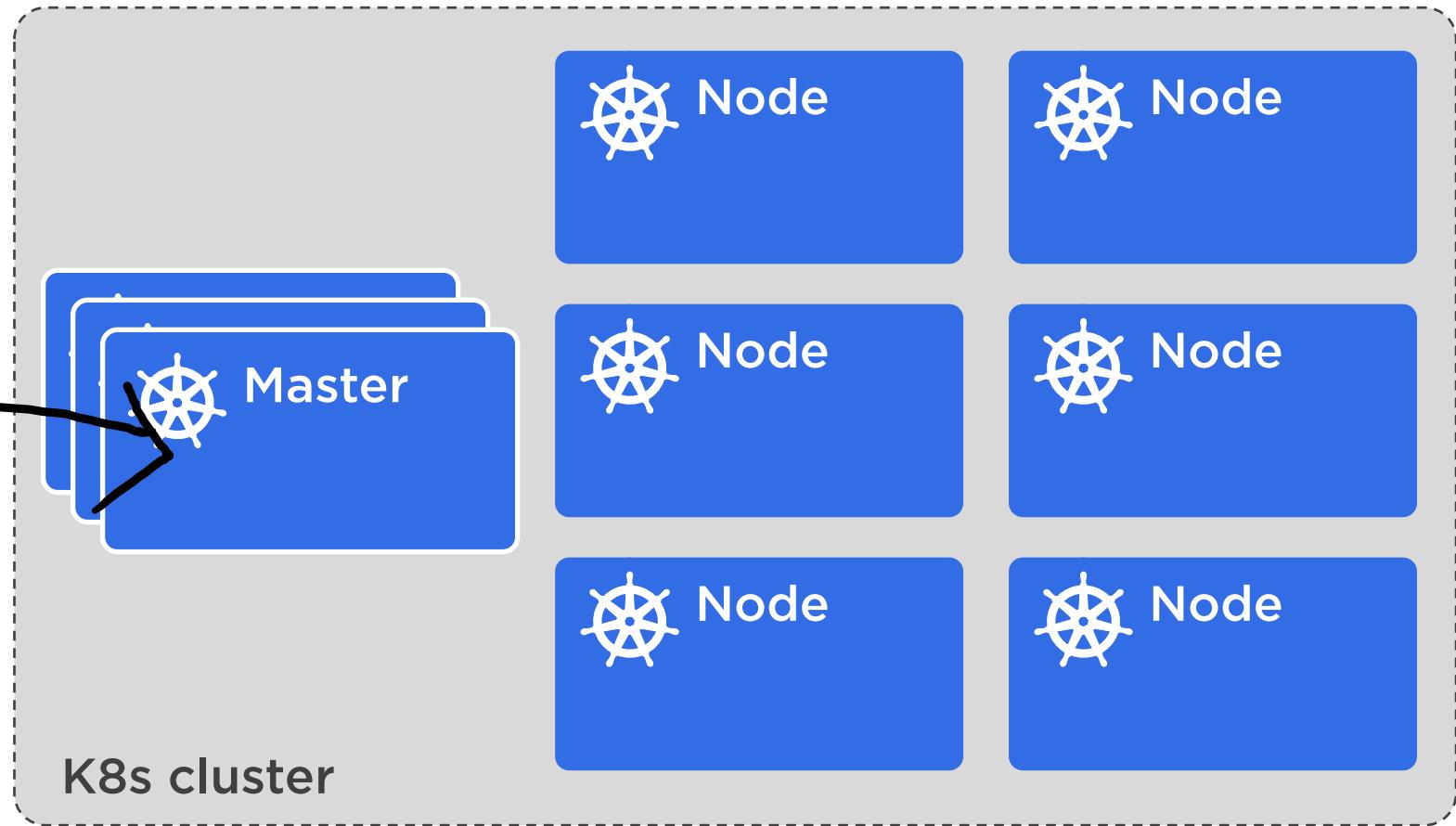
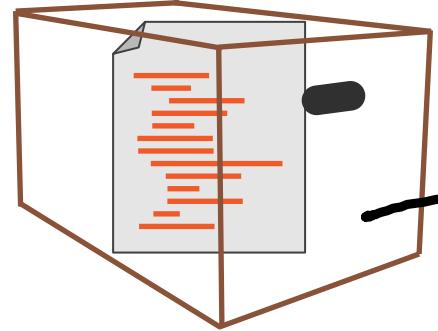






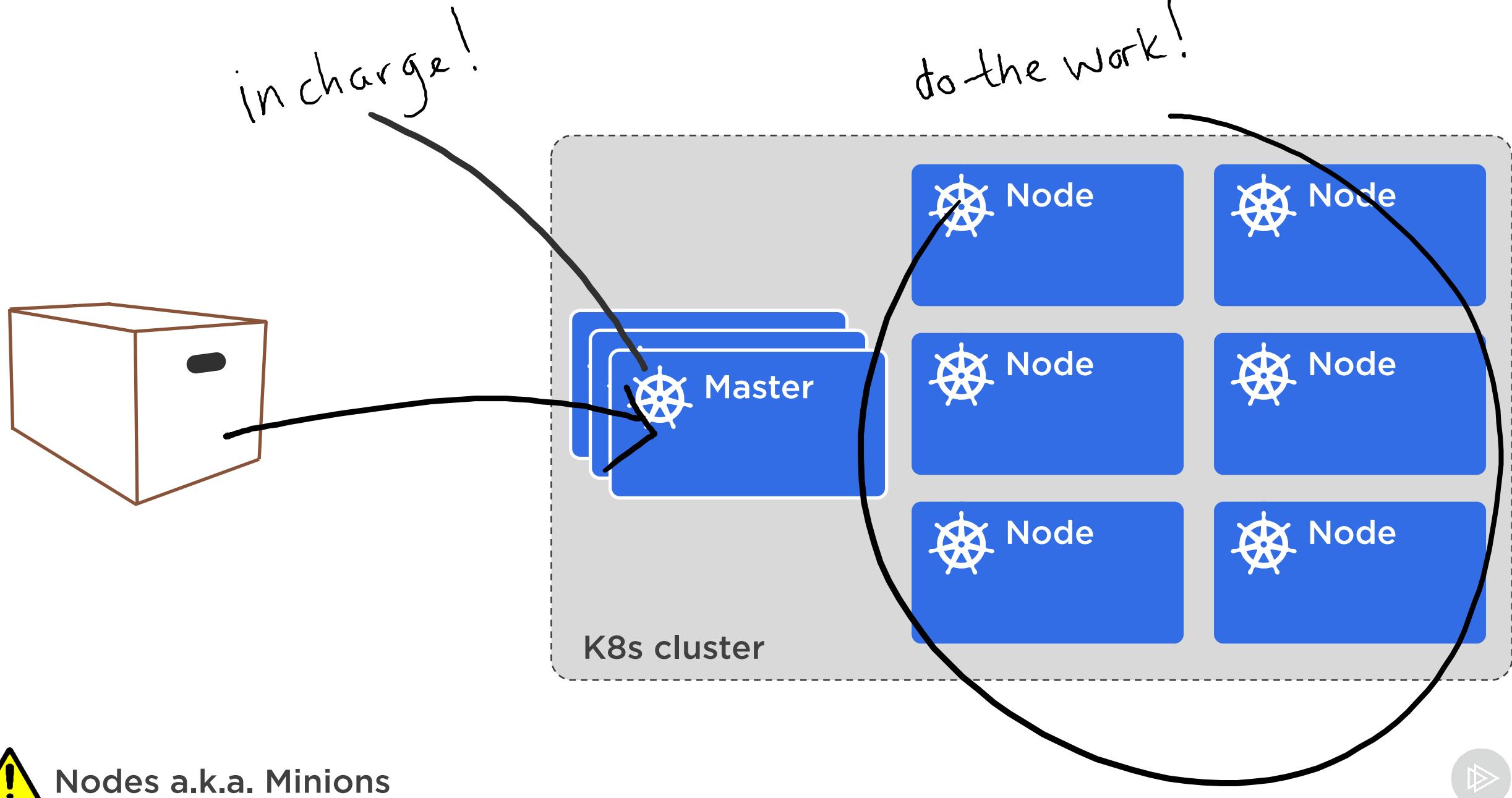
- Node 1
- Node 2
- Node 3

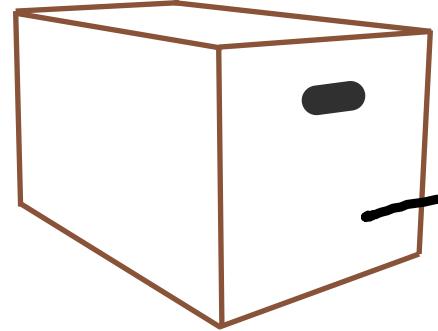




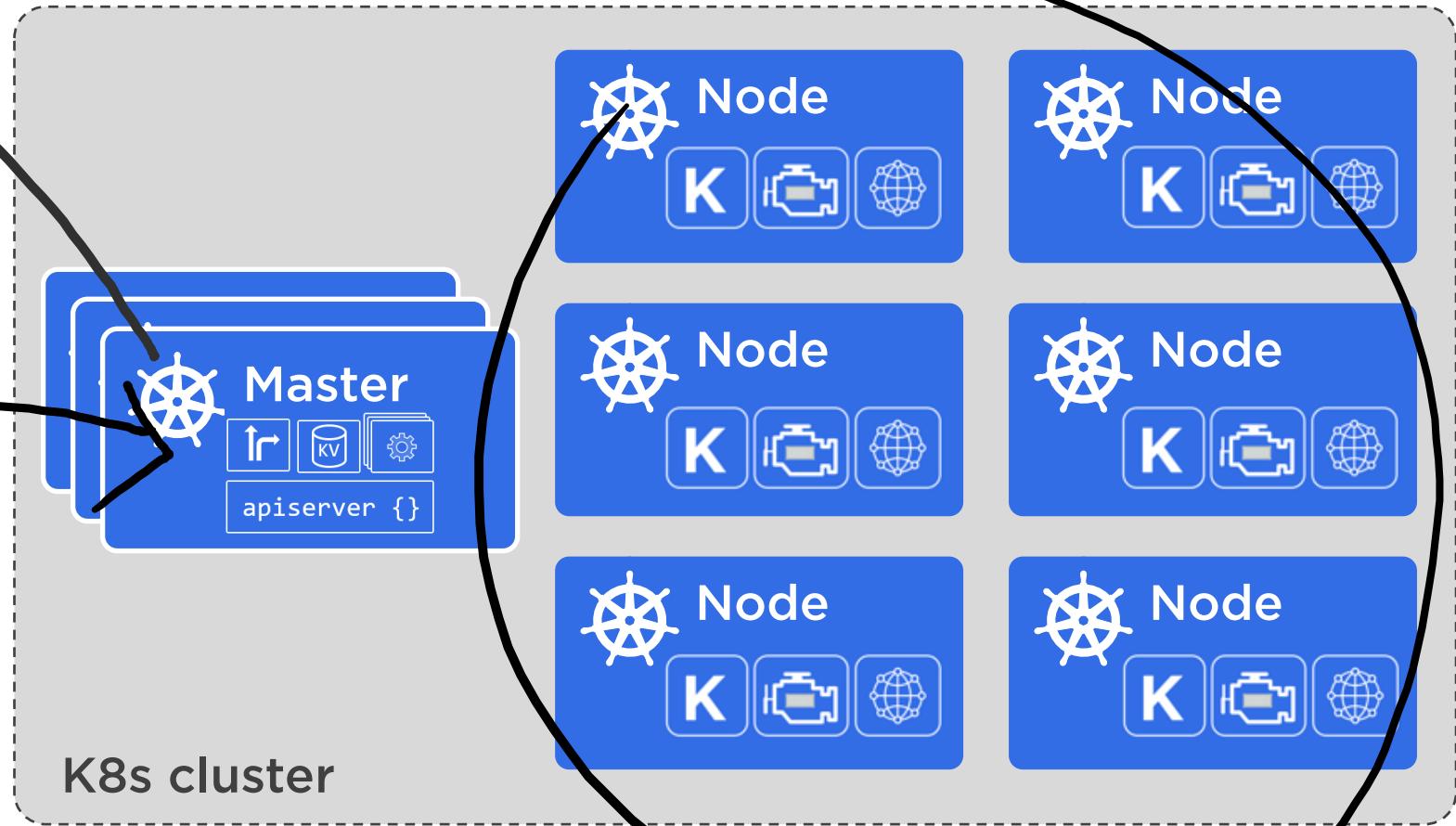
Nodes a.k.a. Minions







in charge!



do the work!



Nodes a.k.a. Minions



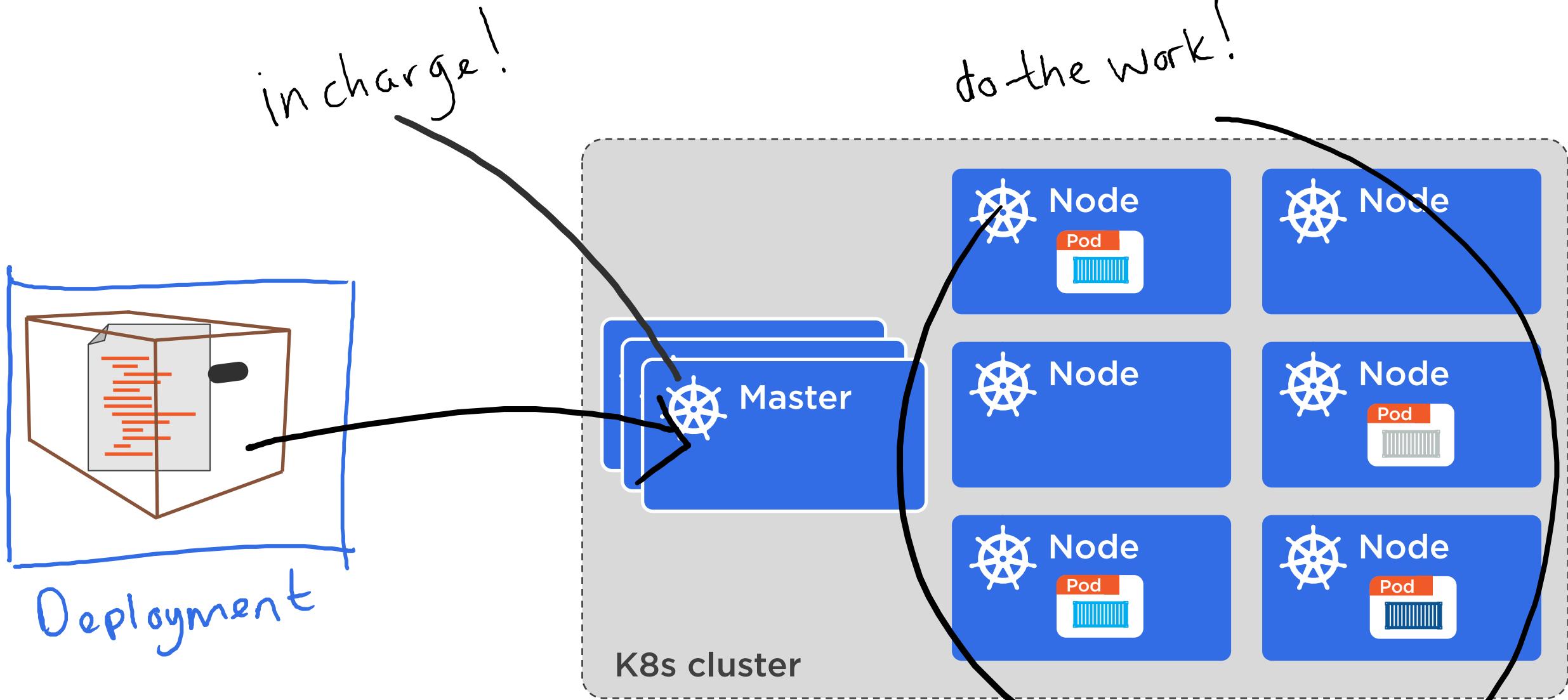


Master



apiserver {}





Nodes a.k.a. Minions



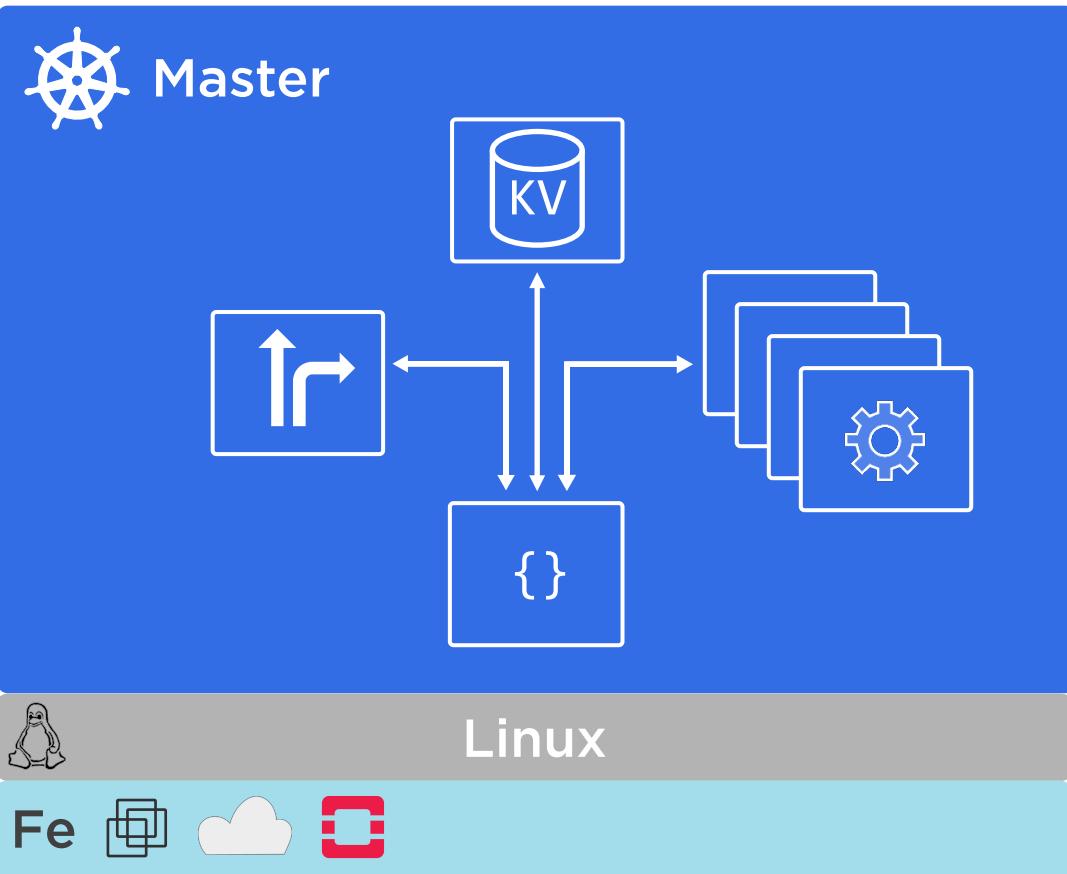
Masters

The Kubernetes Control Plane

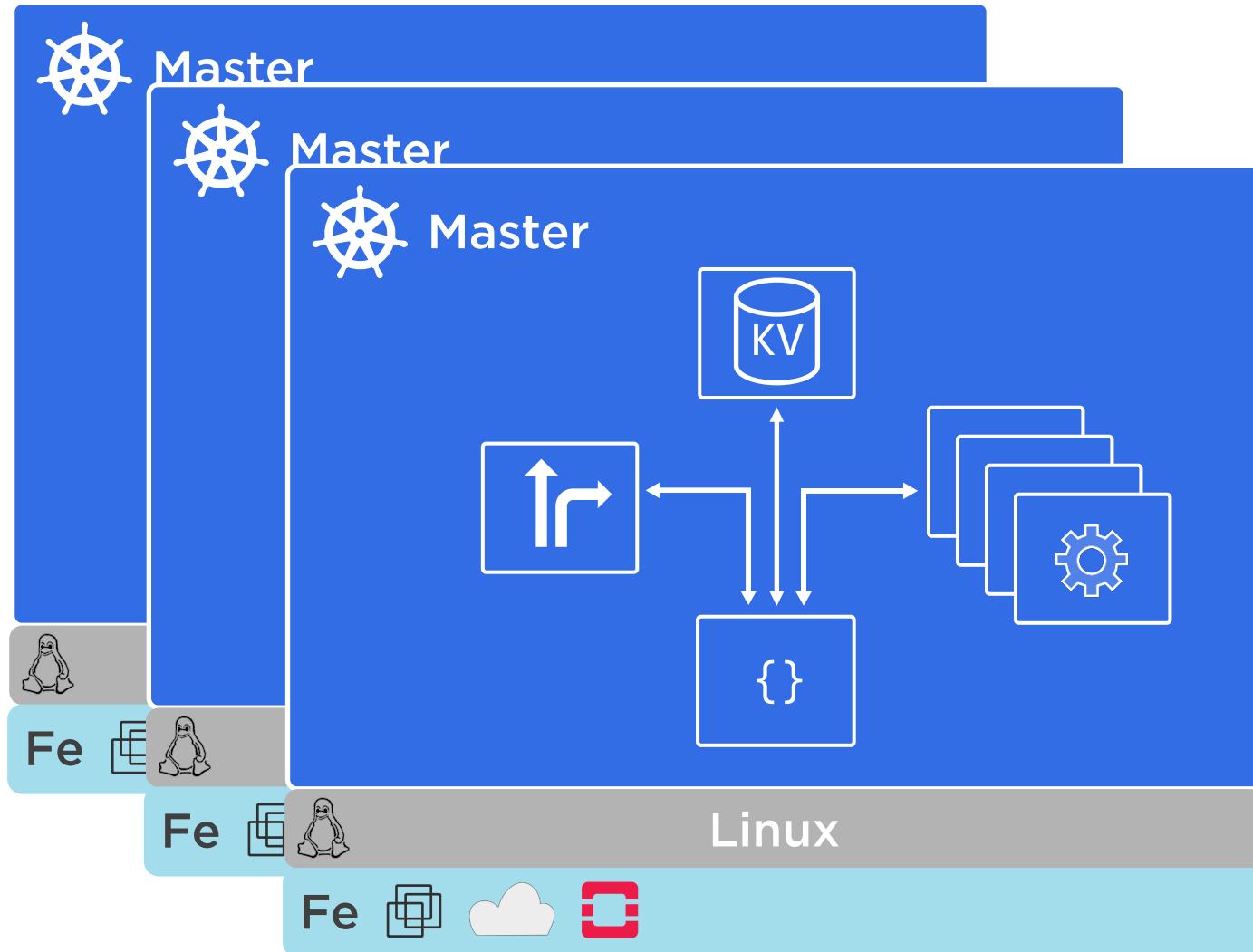


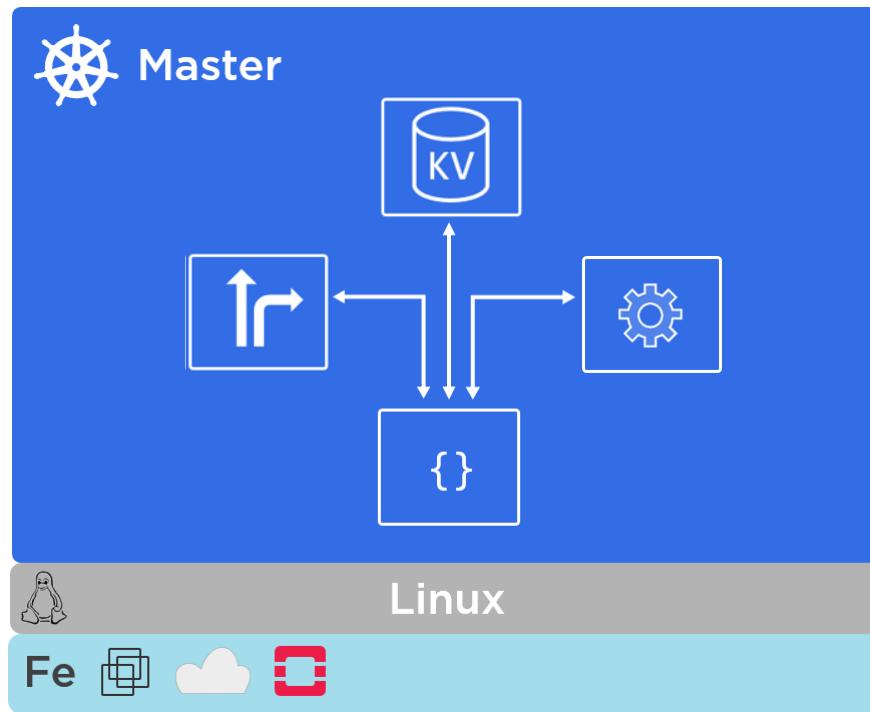
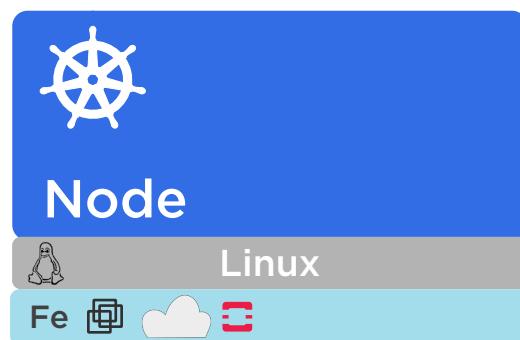
Masters

The Kubernetes Control Plane



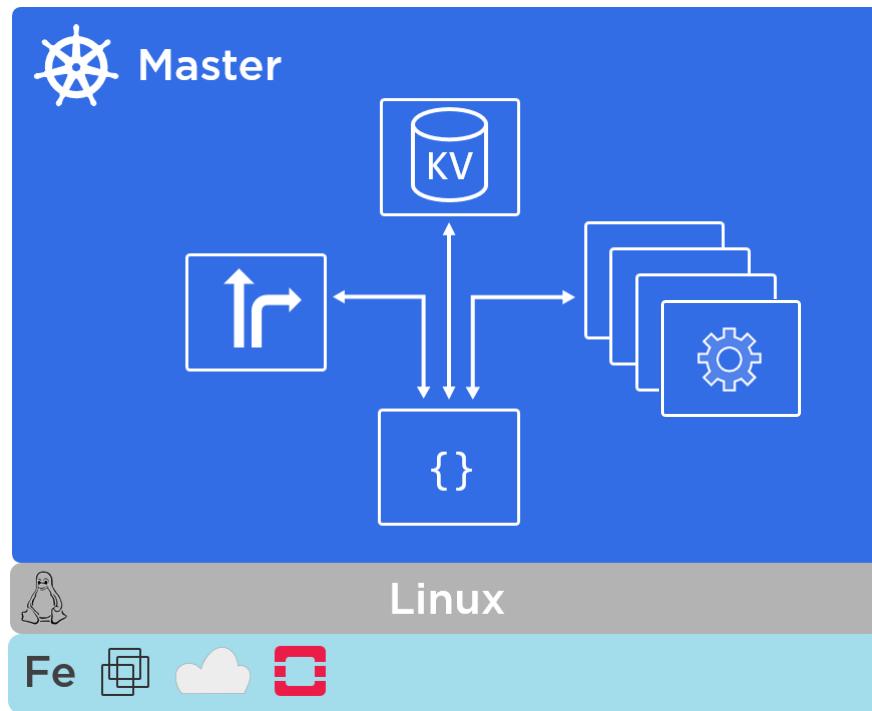
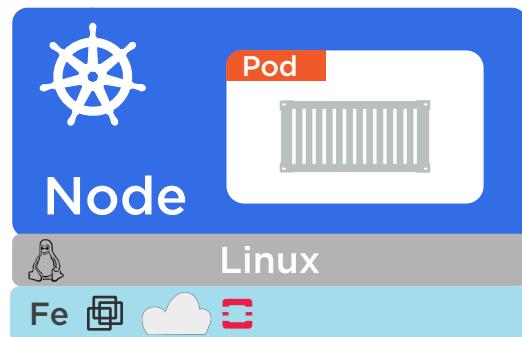
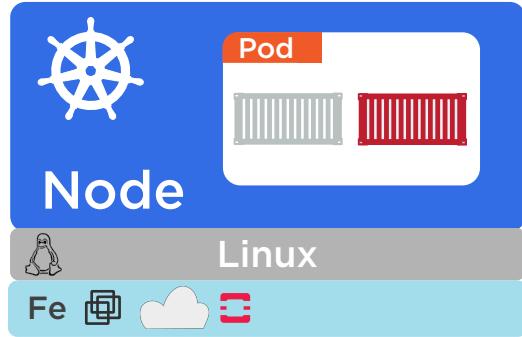
Multi-master HA



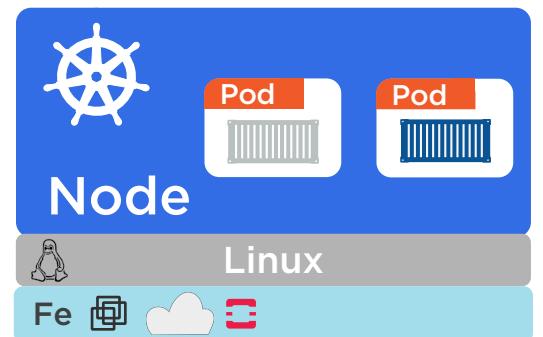
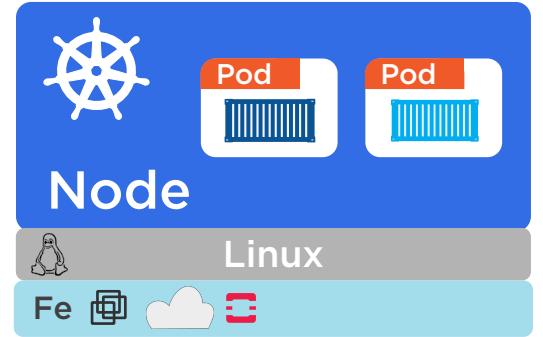


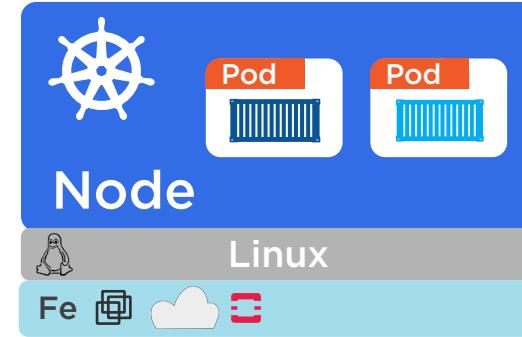
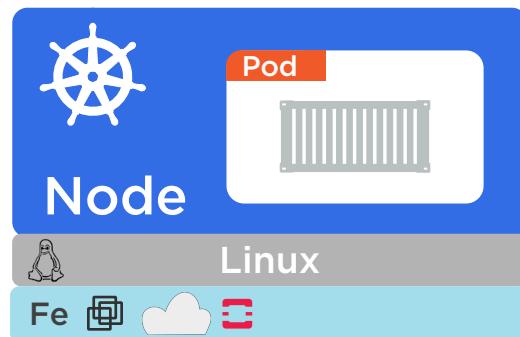
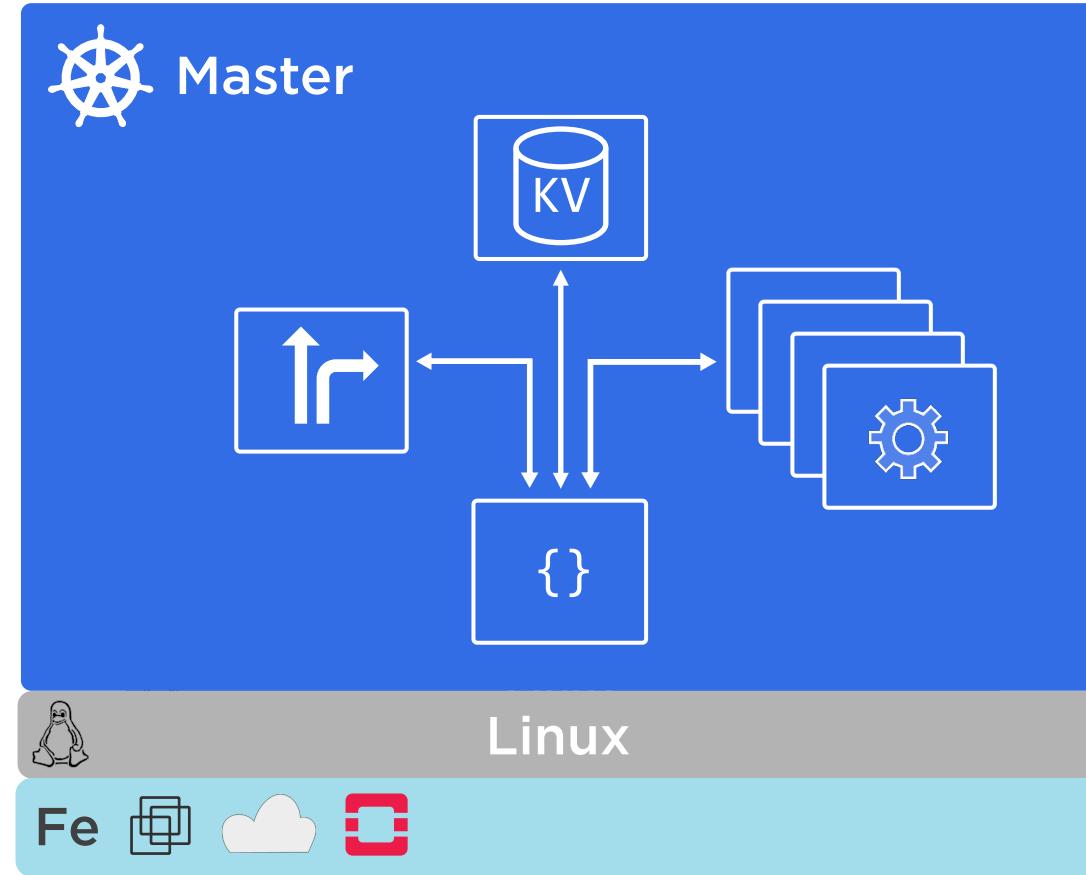
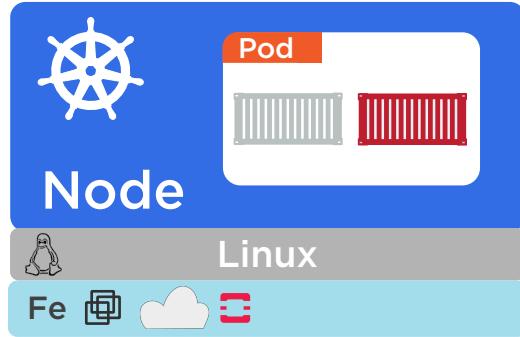
Distributed Control Plane
(future)



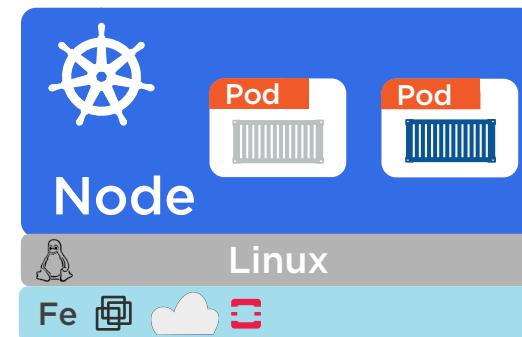


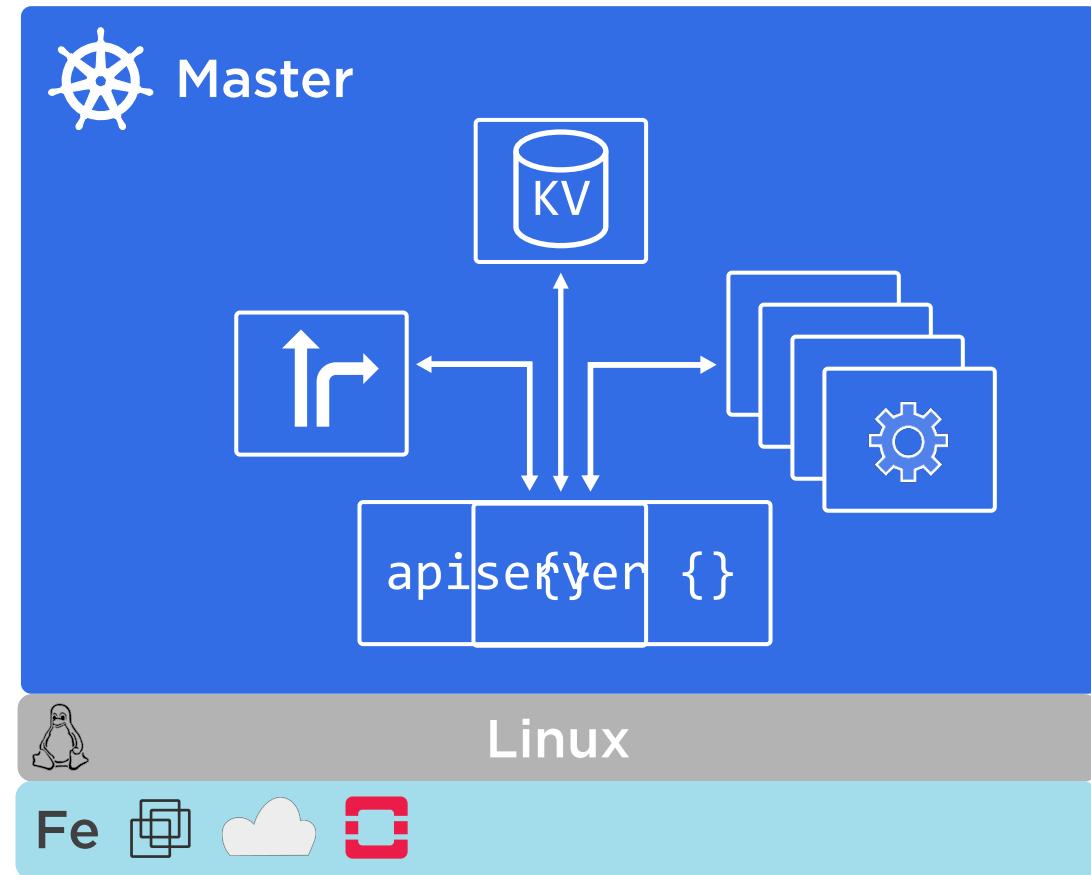
**Don't run user workloads on
“Master”**

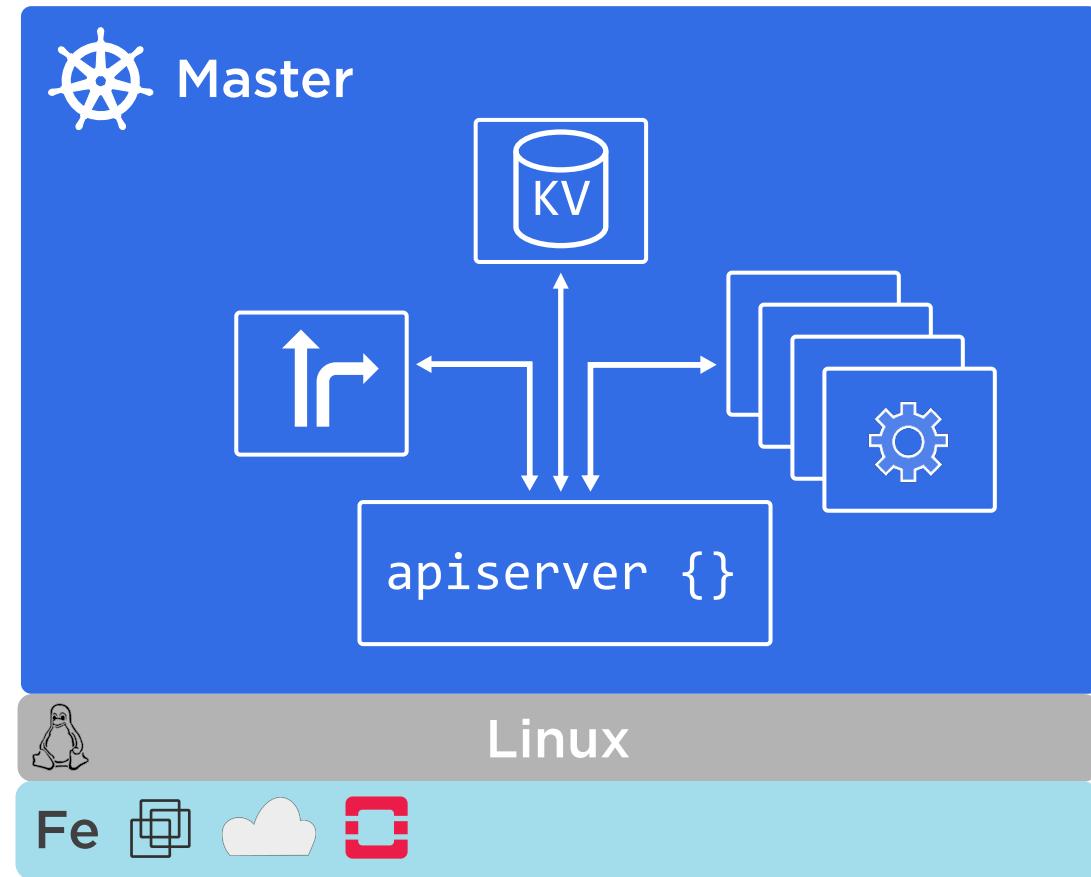




Don't run user workloads on
“Master”







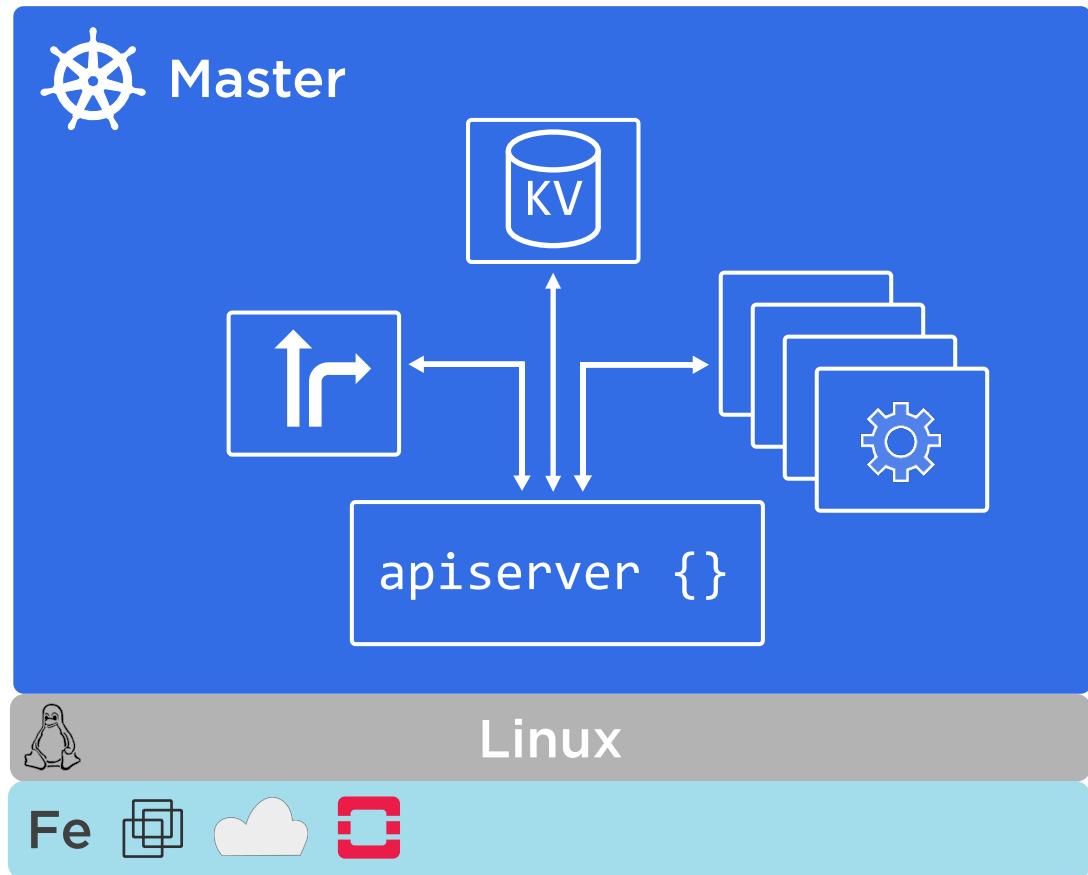
kube-apiserver

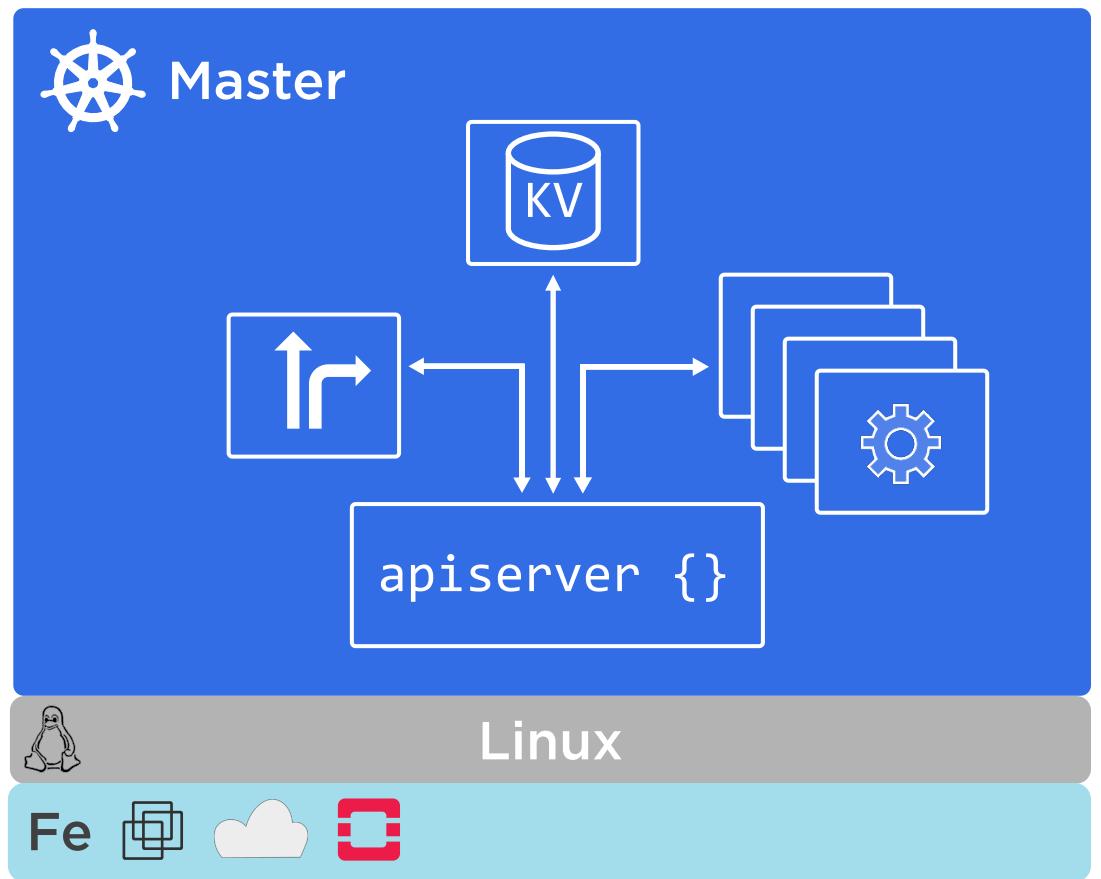
Front-end to the control plane

Exposes the API (REST)

Consumes JSON

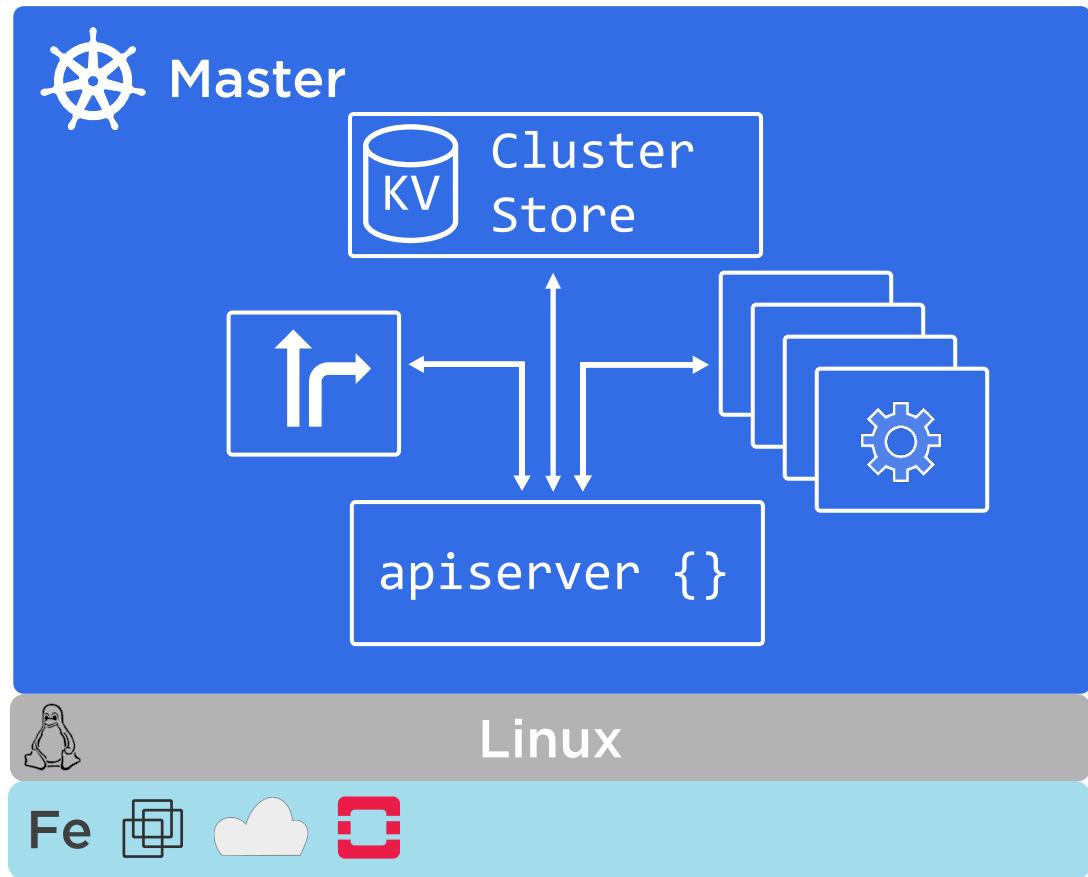
(via manifest files)

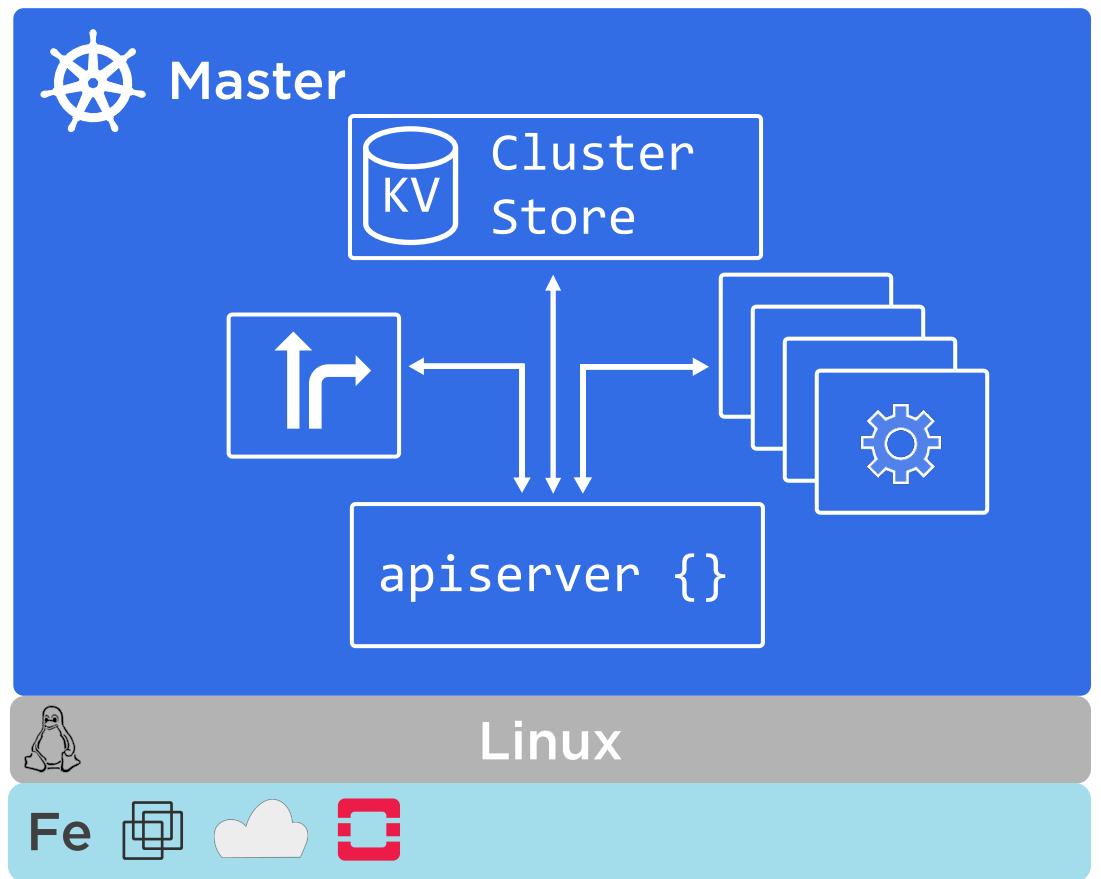




Cluster store

Persistent storage
Cluster state and config
Uses etcd
Distributed, consistent,
watchable...
The “*source of truth*” for
the cluster
Have a backup plan for it!





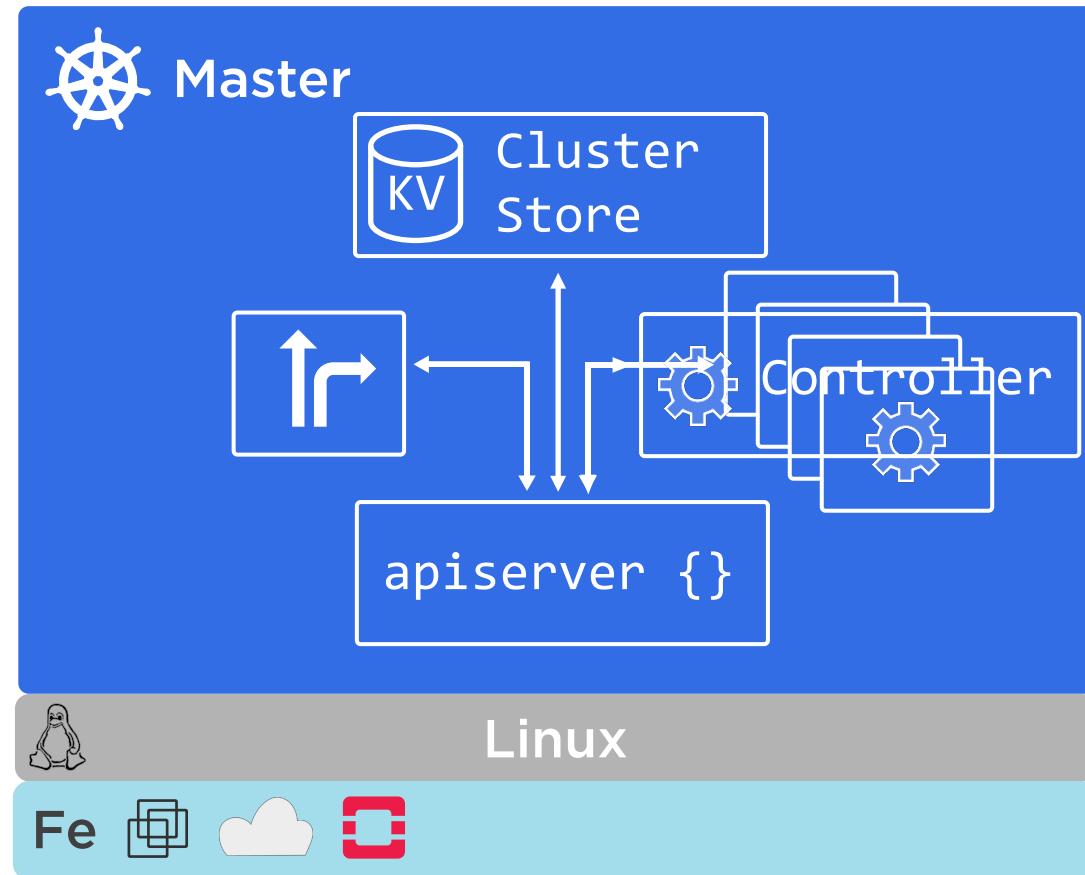
kube-controller-manager

Controller of controllers

- Node controller
- Endpoints controller
- Namespace controller
- ...

Watches for changes

Helps maintain *desired state*



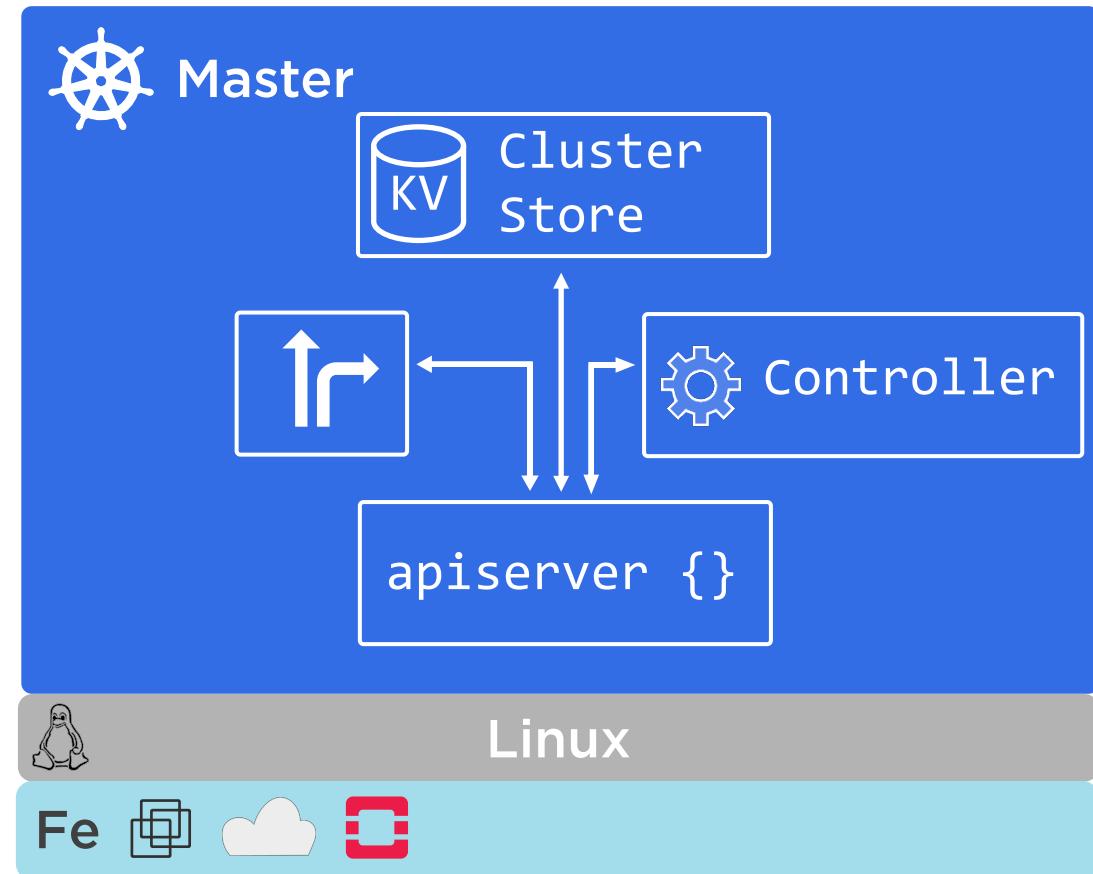
kube-controller-manager

Controller of controllers

- Node controller
- Endpoints controller
- Namespace controller
- ...

Watches for changes

Helps maintain *desired state*

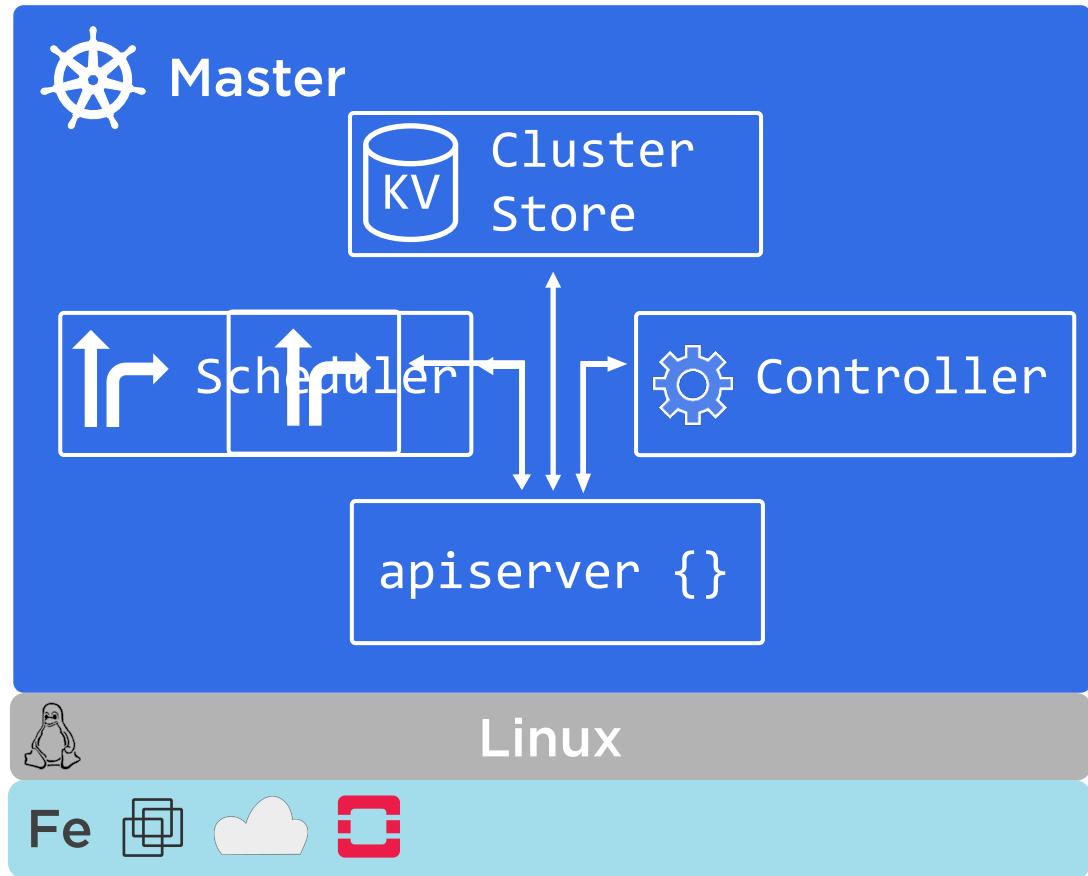


kube-scheduler

Watches apiserver for new pods

Assigns work to nodes

- affinity/anti-affinity
- constraints
- resources
- ...

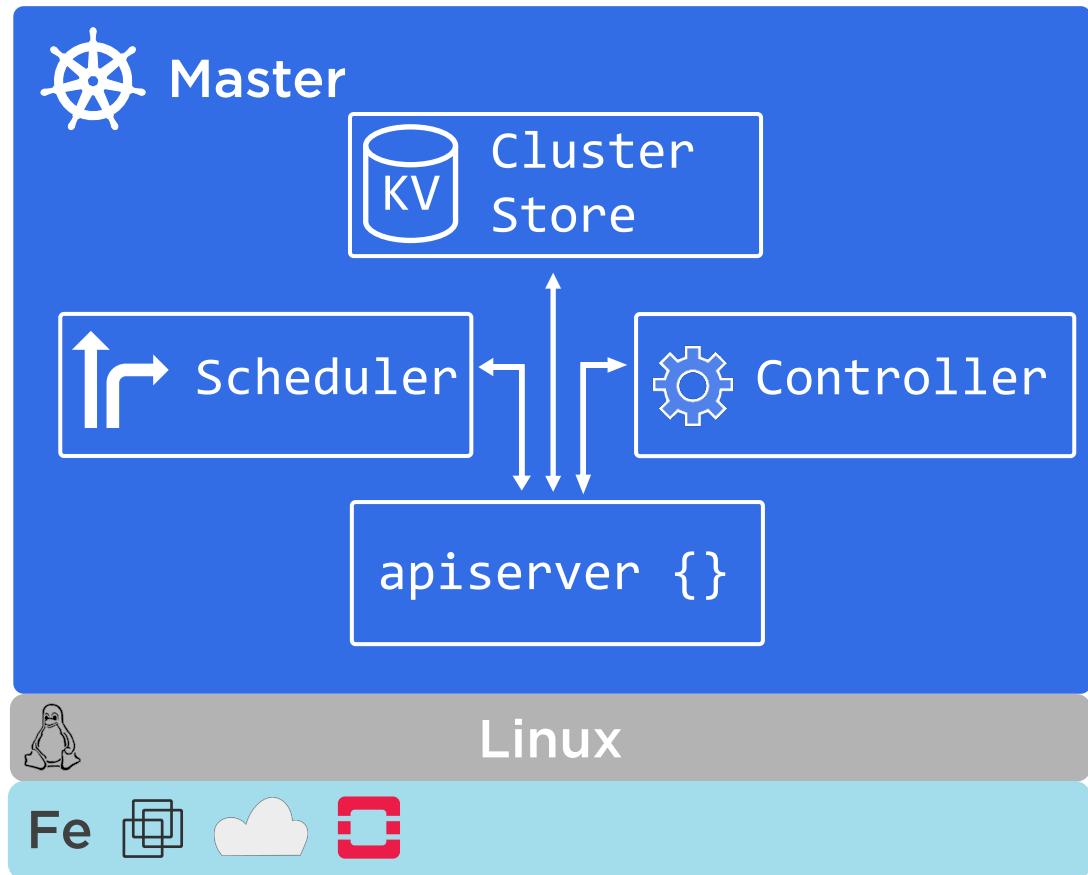


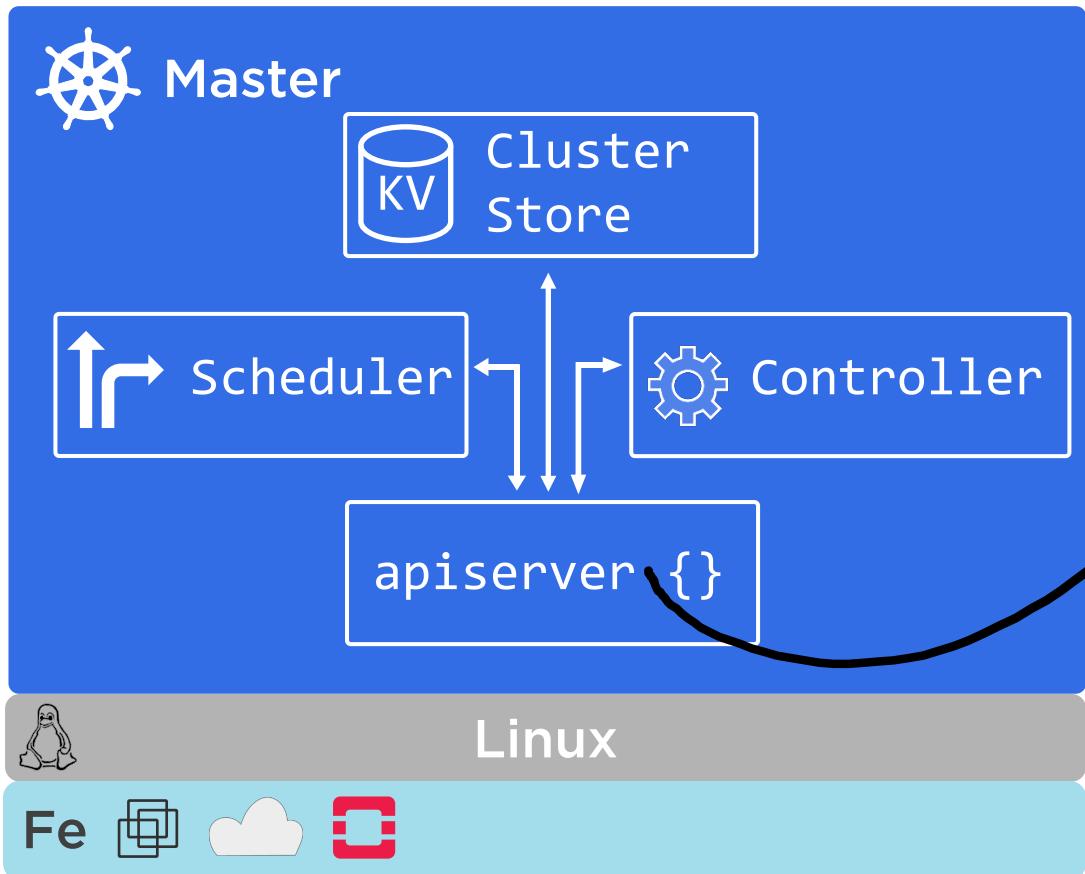
kube-scheduler

Watches apiserver for new pods

Assigns work to nodes

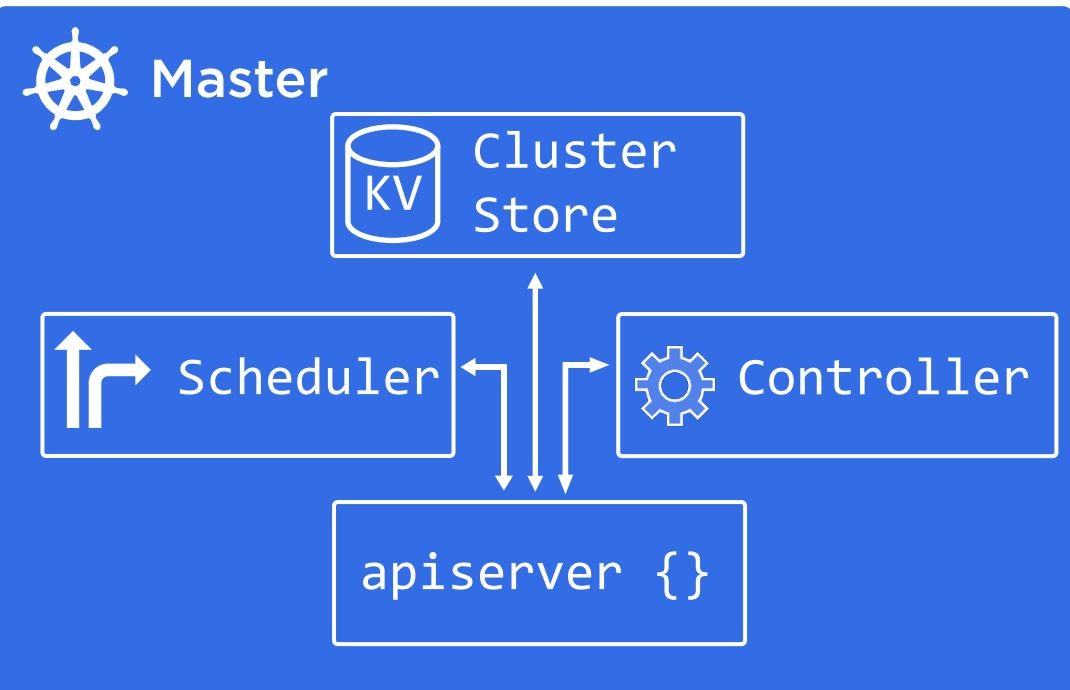
- affinity/anti-affinity
- constraints
- resources
- ...

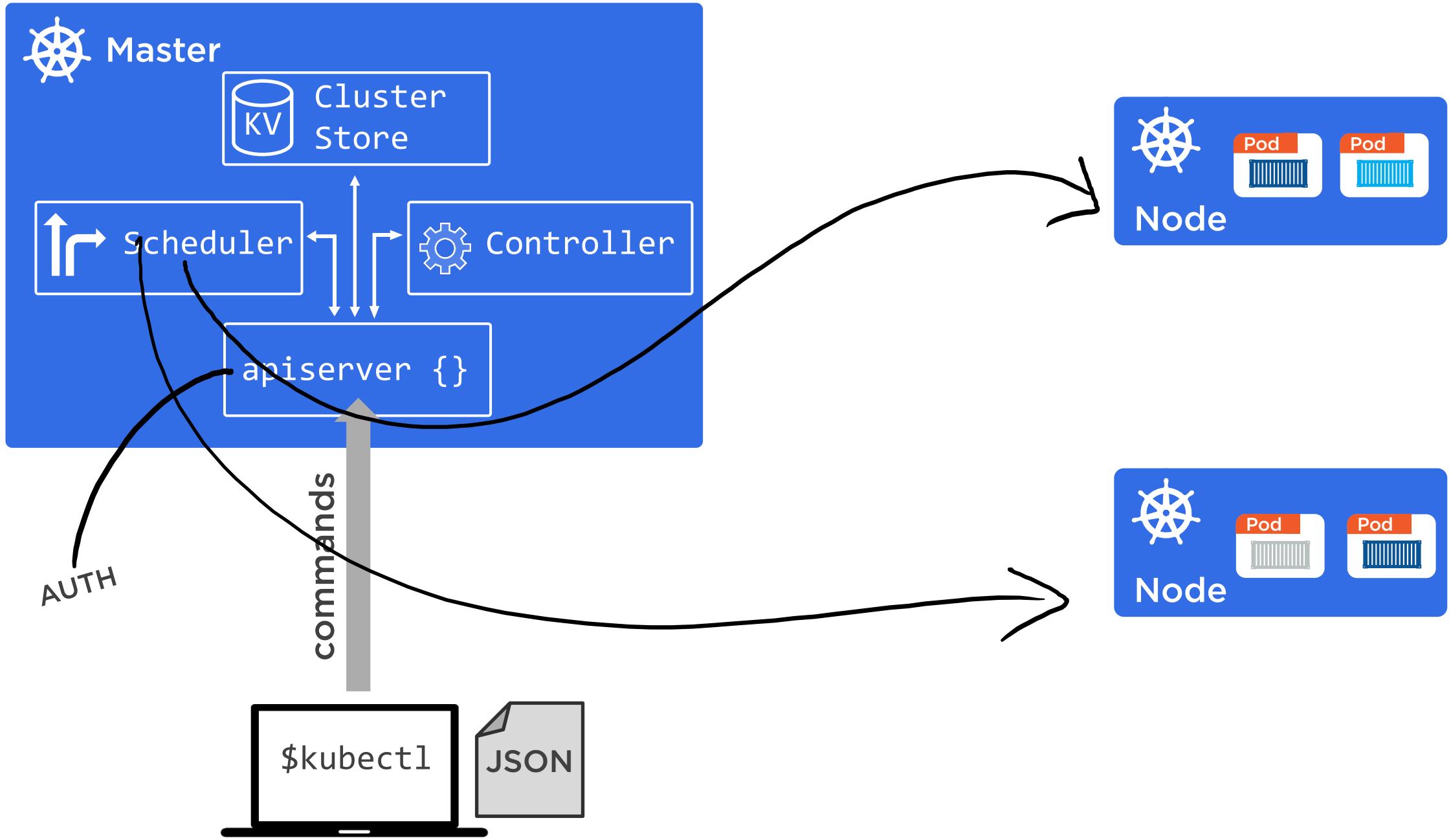




a.k.a
master





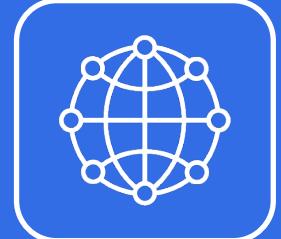
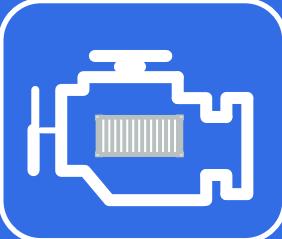


a.k.a. Nodes
“Missions”
The Kubernetes Workers

a.k.a. Nodes
“Missions”
The Kubernetes Workers



Node



Linux

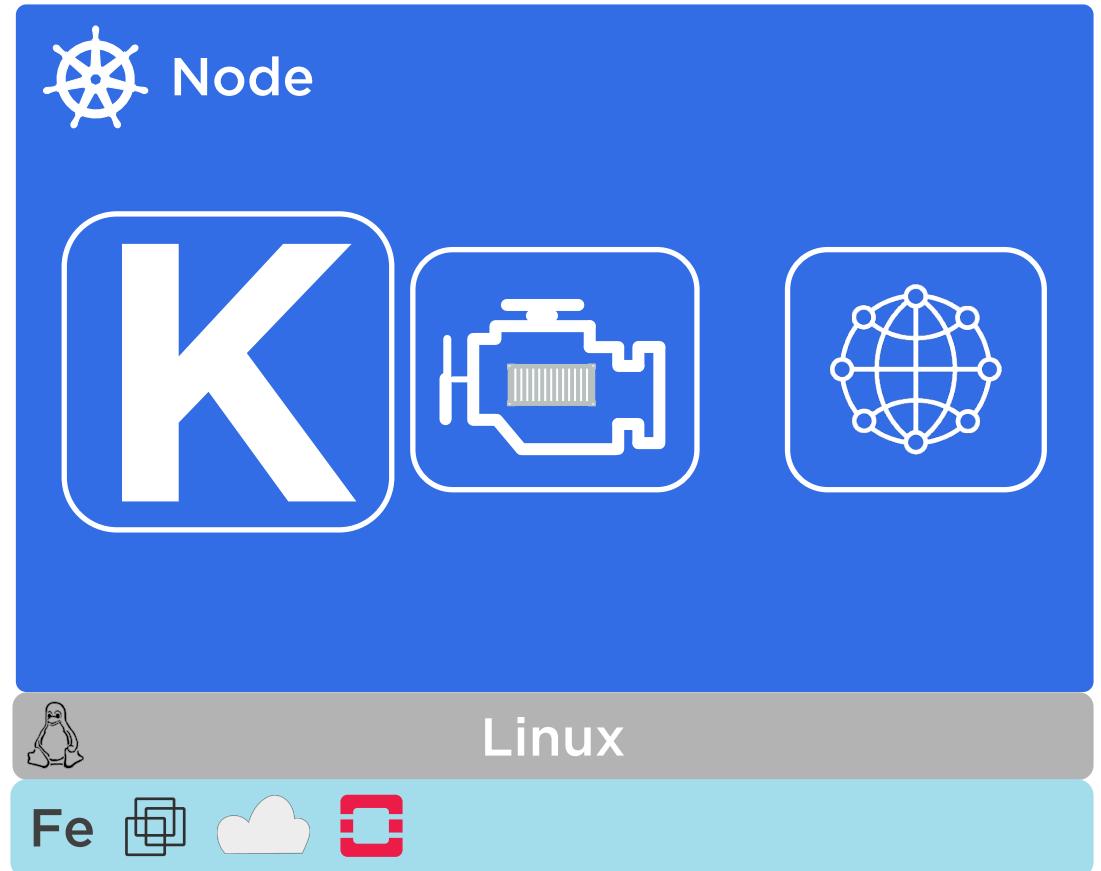
Fe

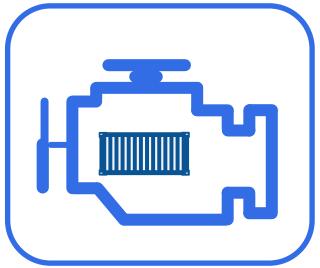




Kubelet

- The **main** Kubernetes agent
- Registers node with cluster
- Watches **apiserver**
- Instantiates **pods**
- Reports back to **master**
- Exposes endpoint on :10255





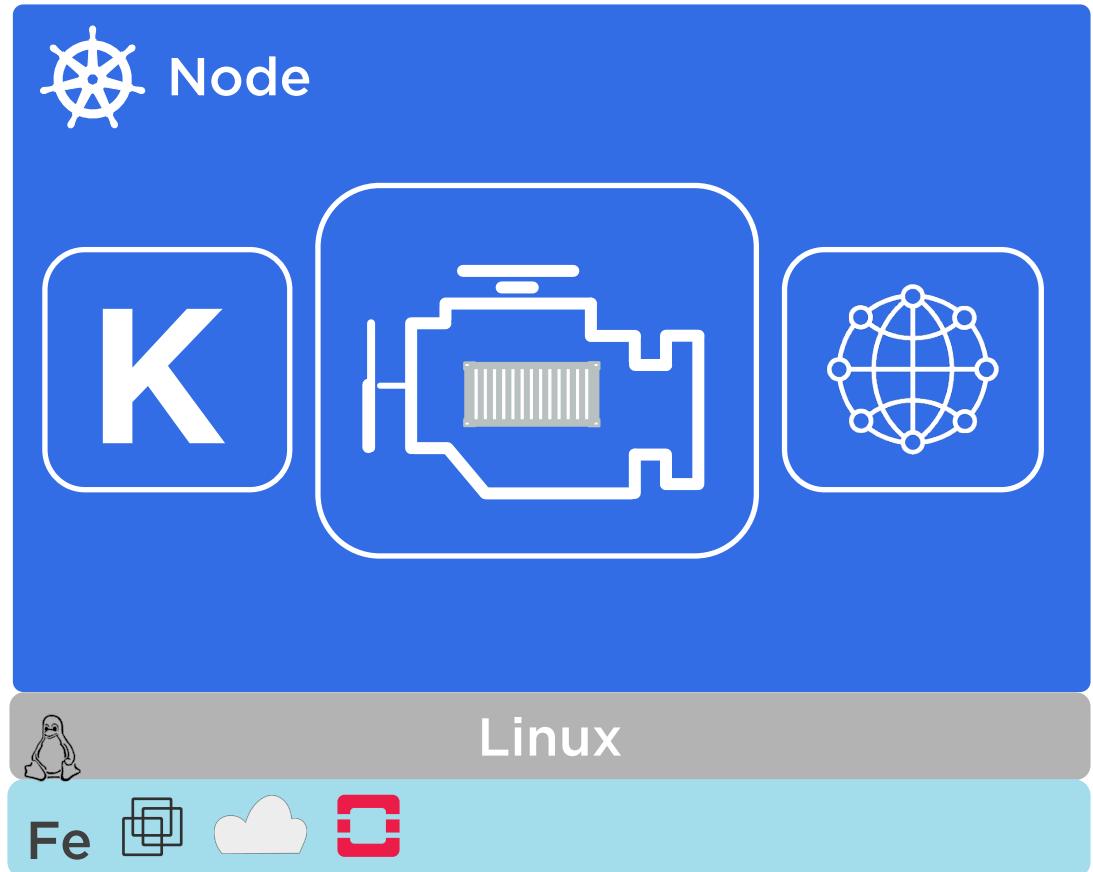
Container Engine

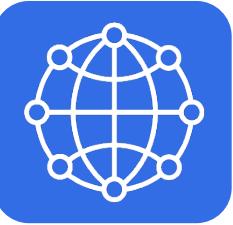
Does container management:

- Pulling images
- Starting/stopping containers
- ...

Pluggable:

- Usually Docker
- Can be rkt

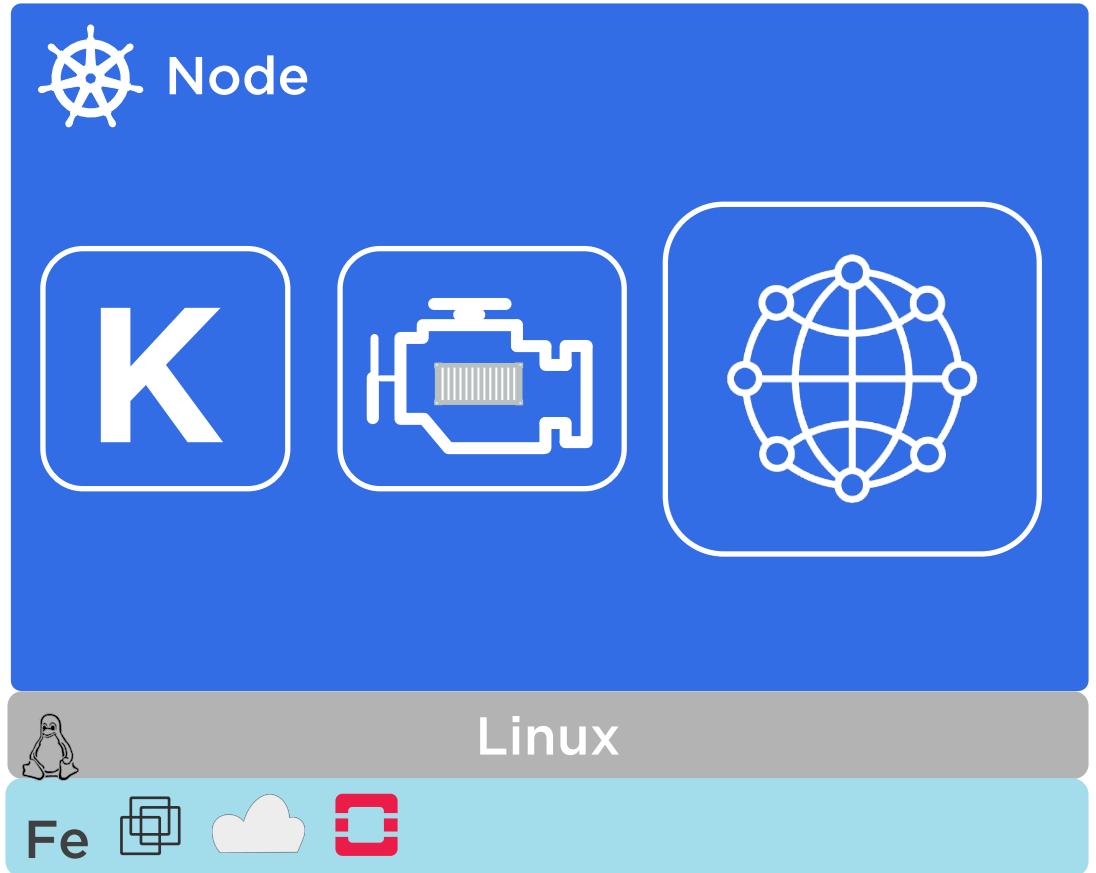




kube-proxy

Kubernetes networking:

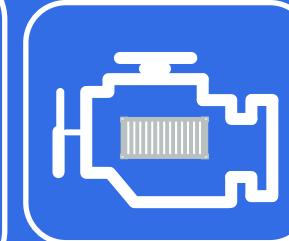
- Pod IP addresses
 - All containers in a pod share a single IP
- Load balances across all pods in a service





Kubelet

Main Kubernetes agent



Linux

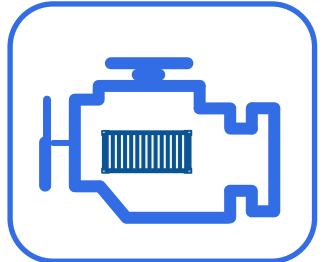
Fe





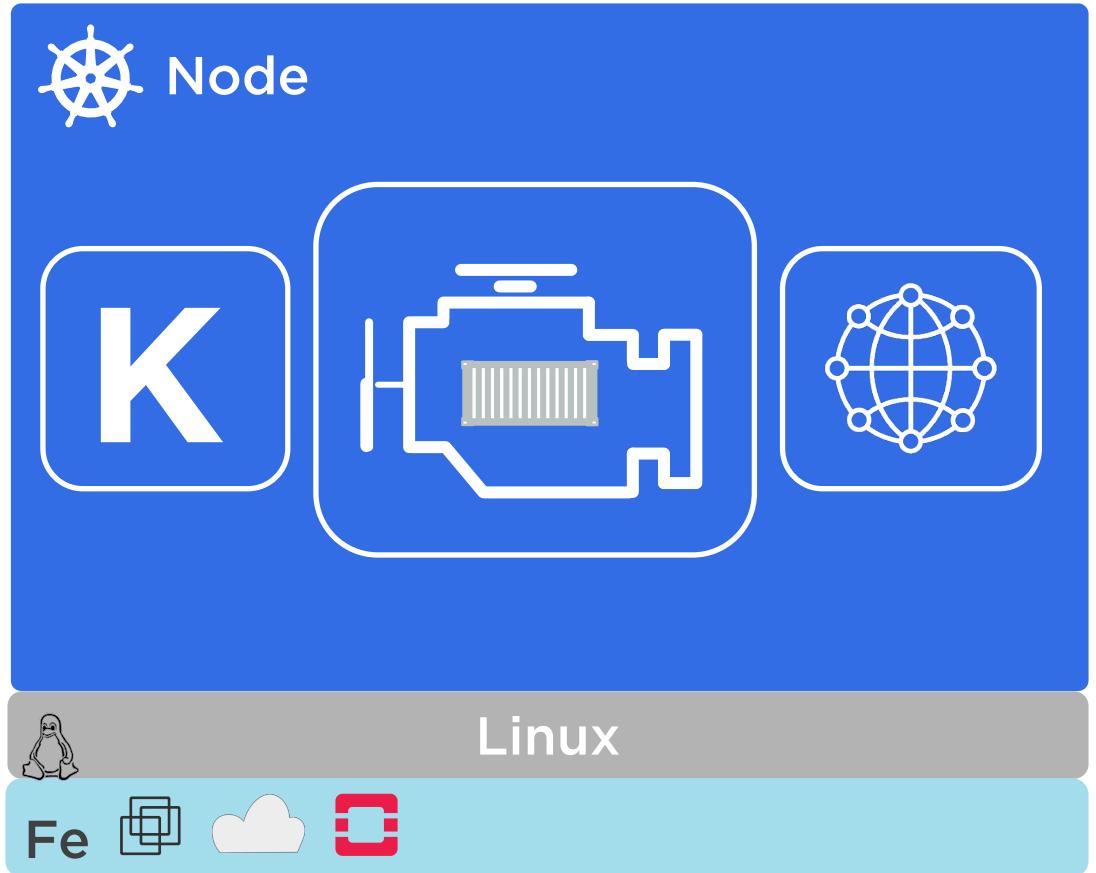
Kubelet

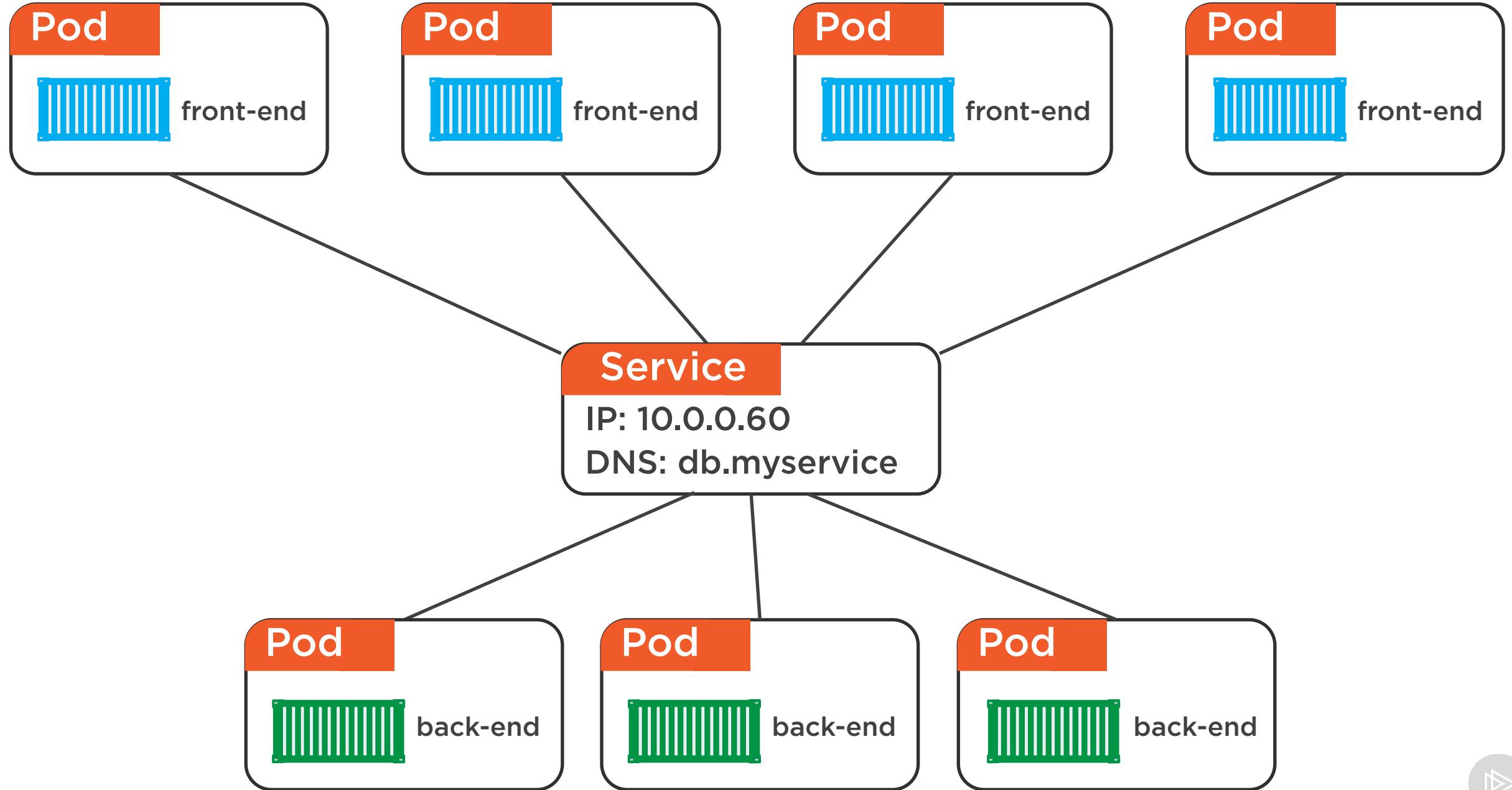
Main Kubernetes agent



Container engine

Docker or rkt

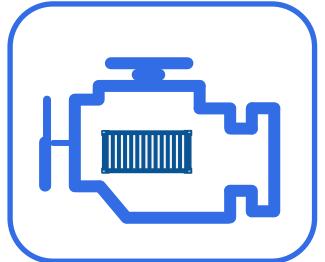






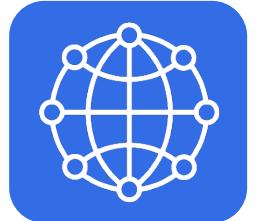
Kubelet

Main Kubernetes agent



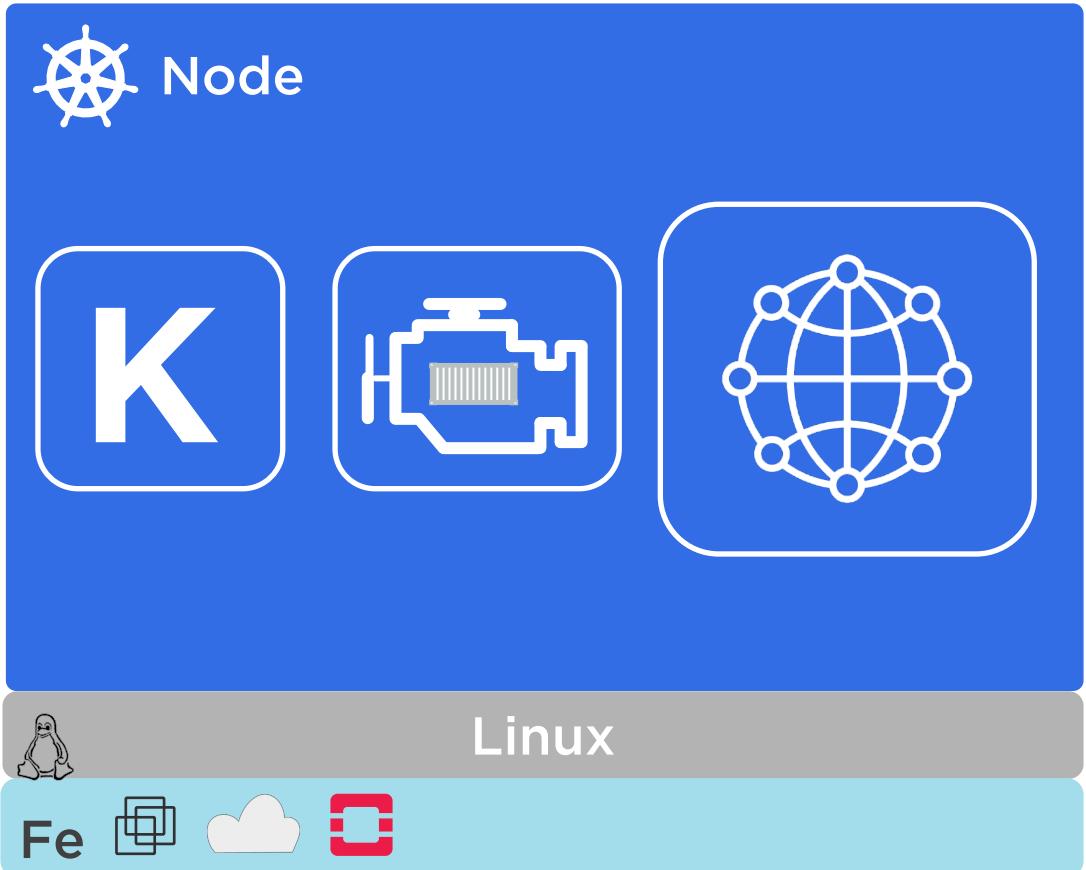
Container engine

Docker or rkt

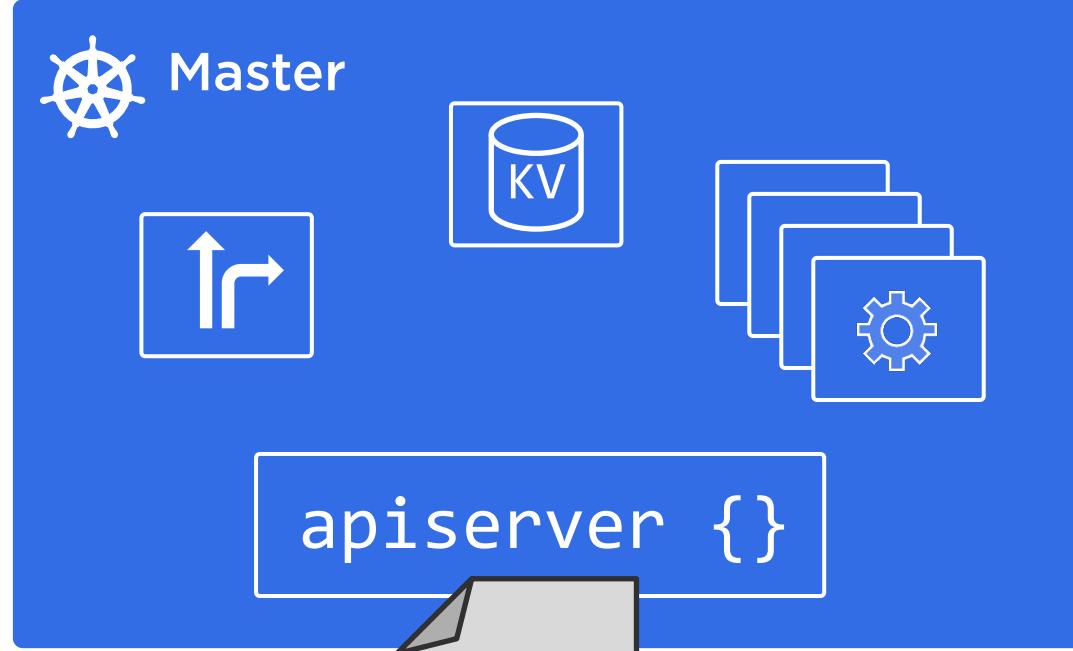


kube-proxy

Kubernetes networking

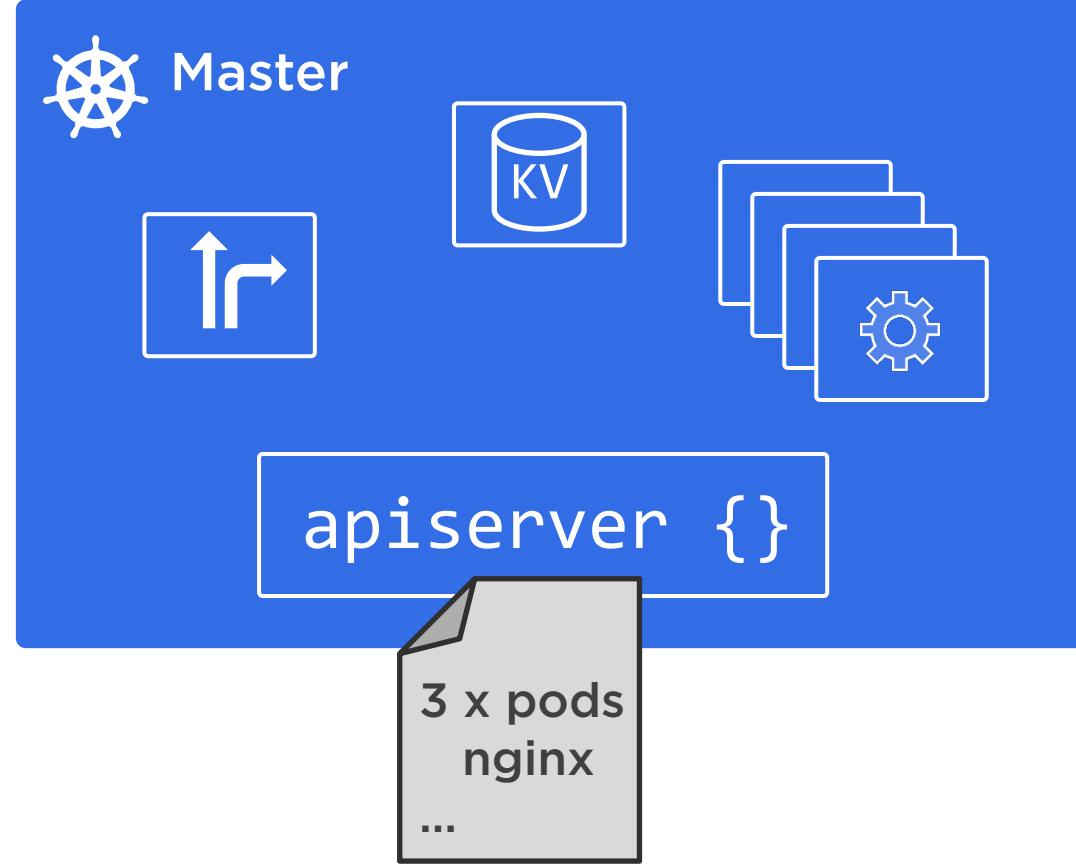


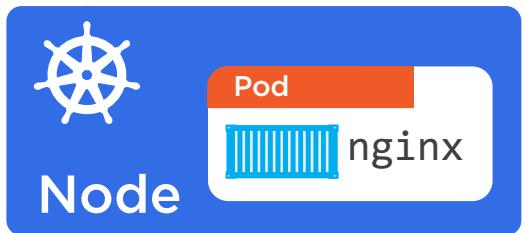
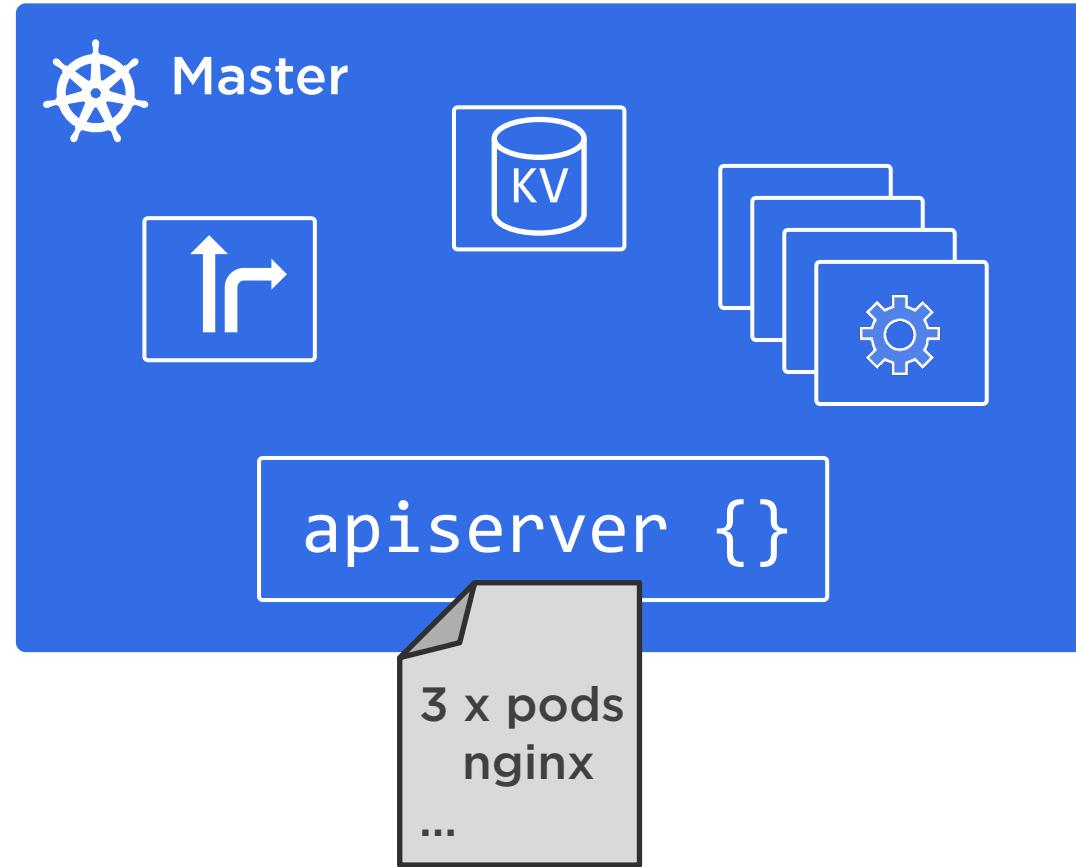
Declarative Model & Desired State

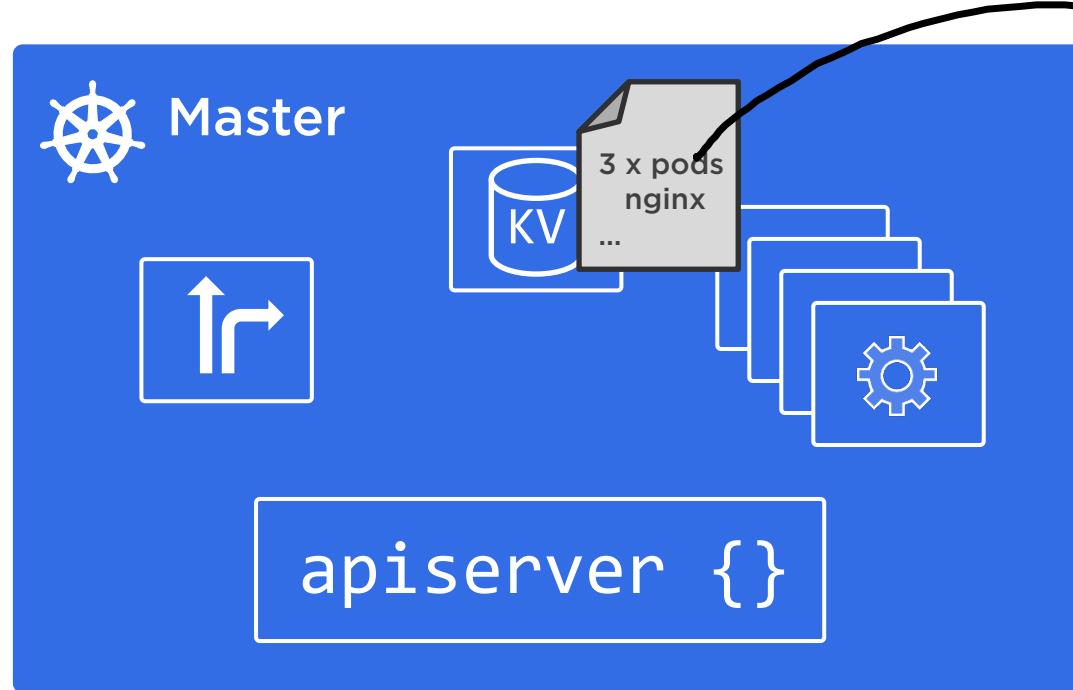


YAML or JSON
Describe **desired**
state

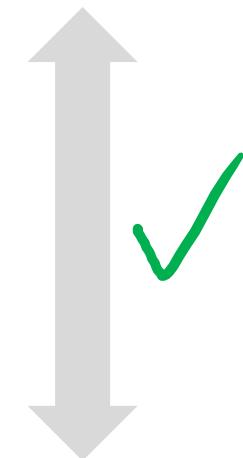




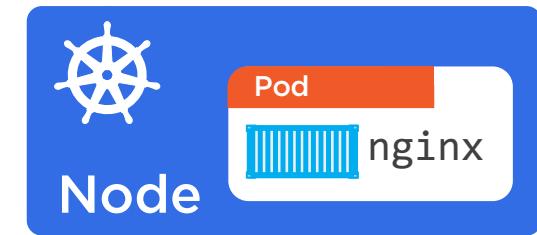
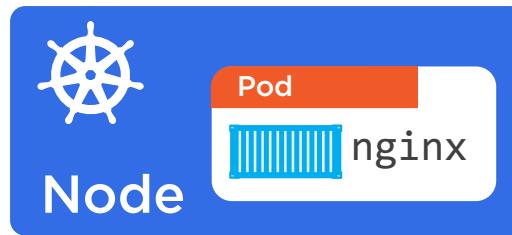
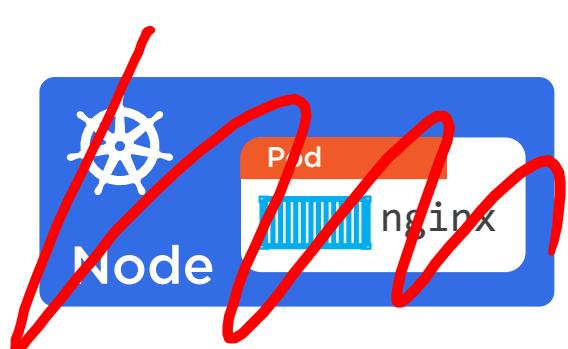


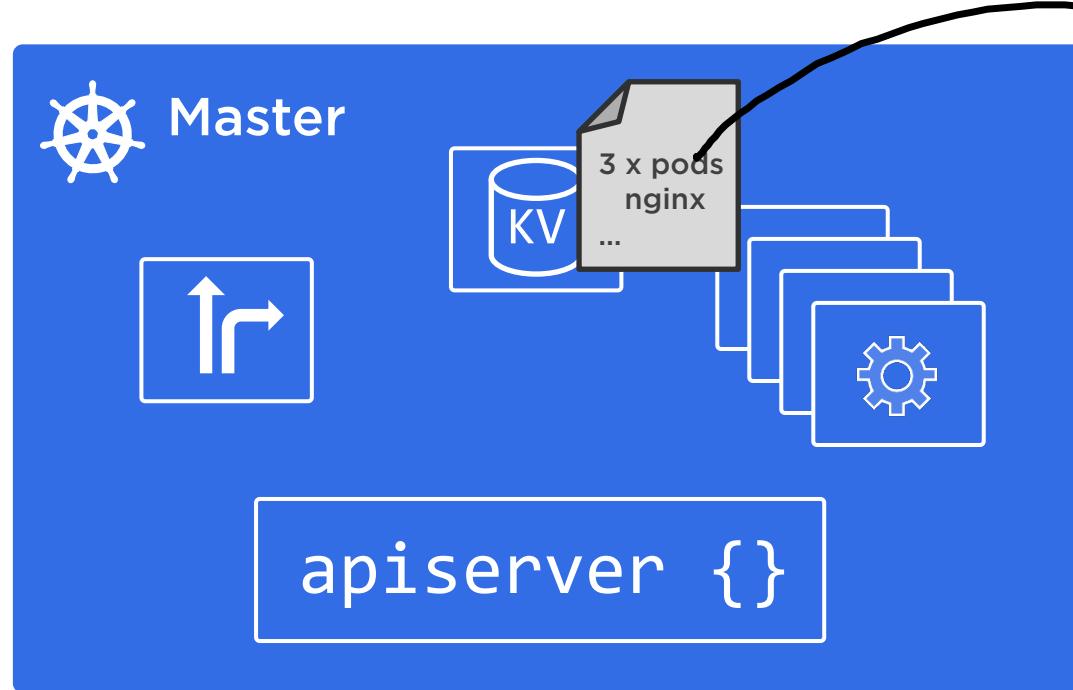


Desired state/
record of intent
• **3 x nginx pods**

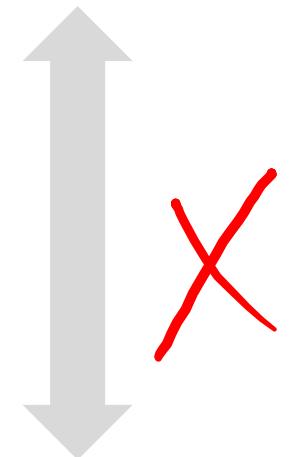


Actual state
• **3 x nginx pods**

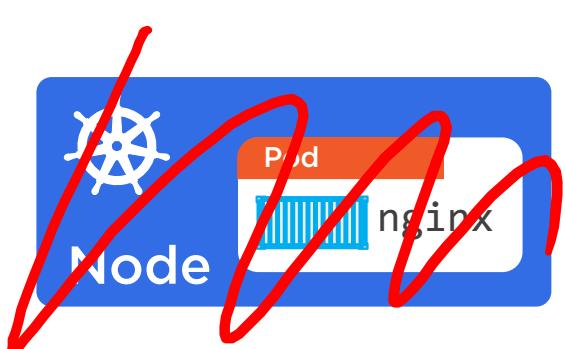


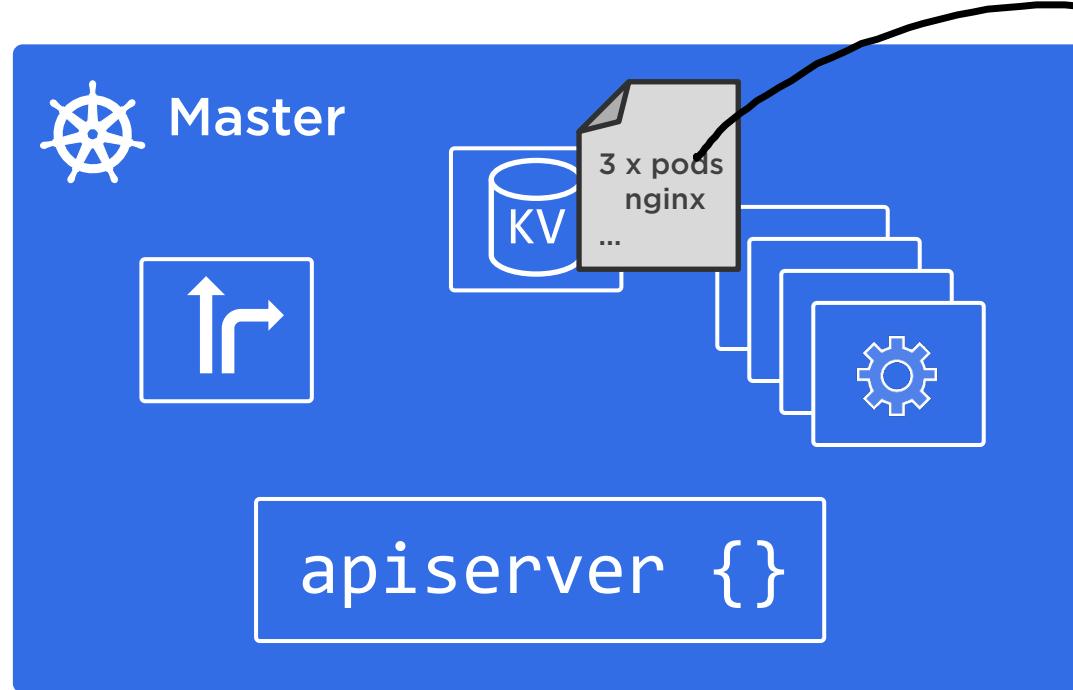


Desired state/
record of intent
• **3 x nginx pods**

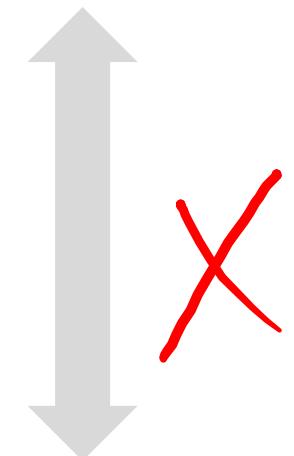


Actual state
• **2 x nginx pods**

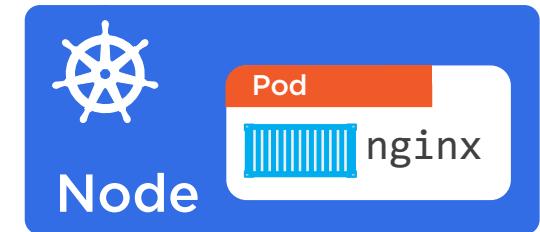
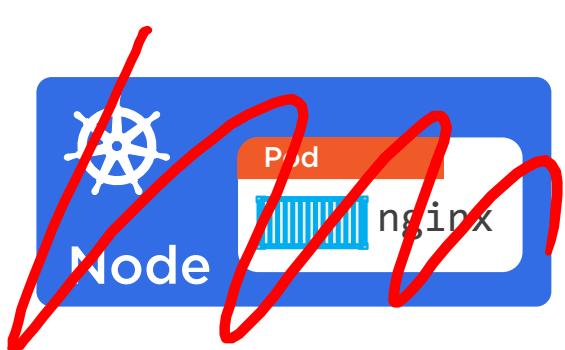


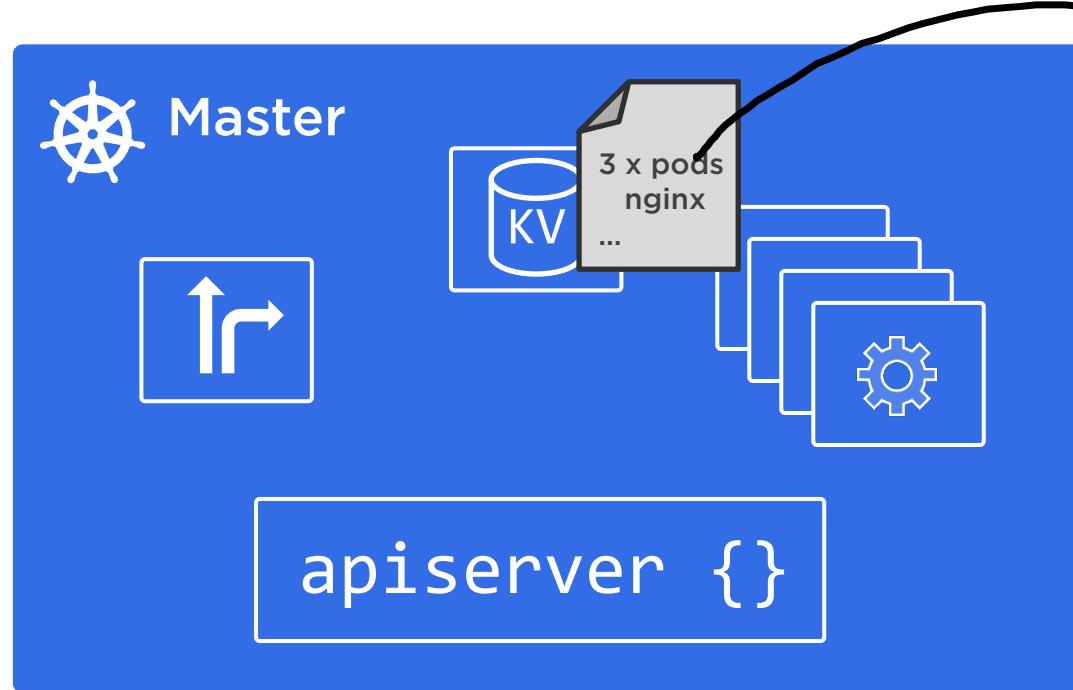


Desired state/
record of intent
• **3 x nginx pods**



Actual state
• **2 x nginx pods**

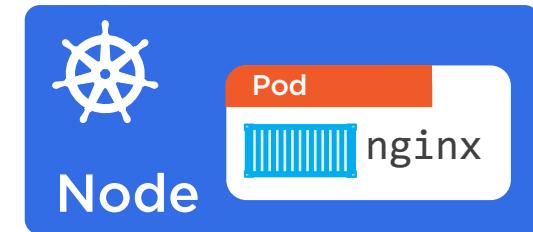




Desired state/
record of intent
• **3 x nginx pods**



Actual state
• **3 x nginx pods**

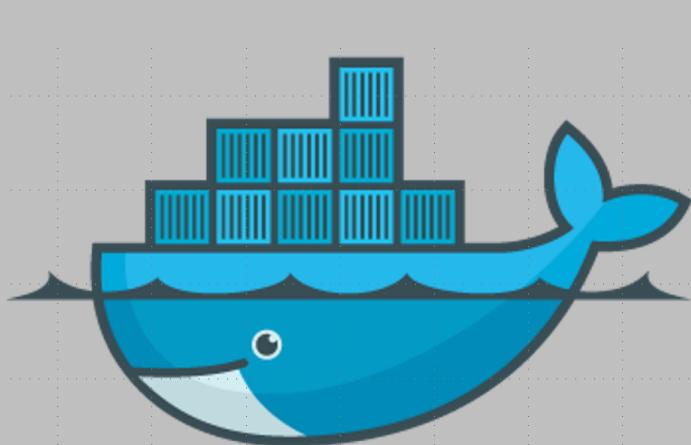




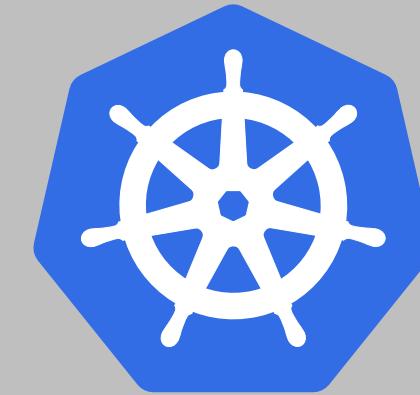
Pods

vm

VM



Container



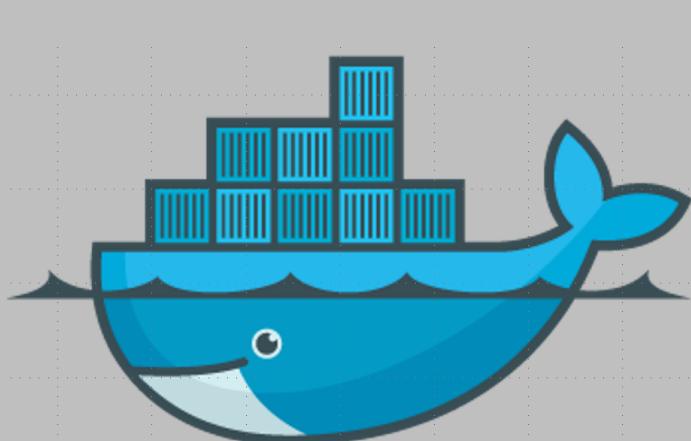
Pod

Atomic units of scheduling



vm

VM



Container



Pod

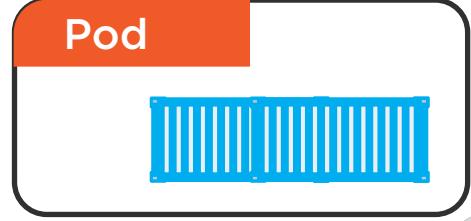
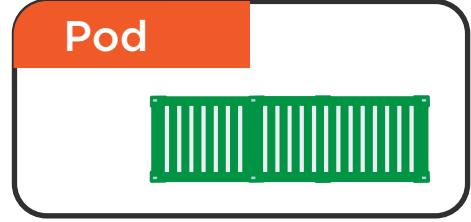
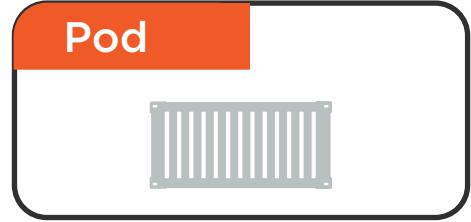
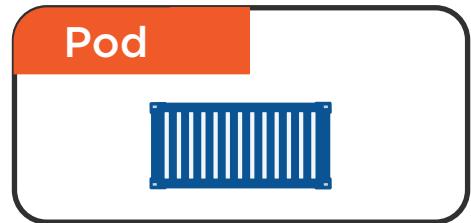
Atomic units of scheduling





Containers always run
inside of pods

Pods can have multiple
containers
(advanced use-case)

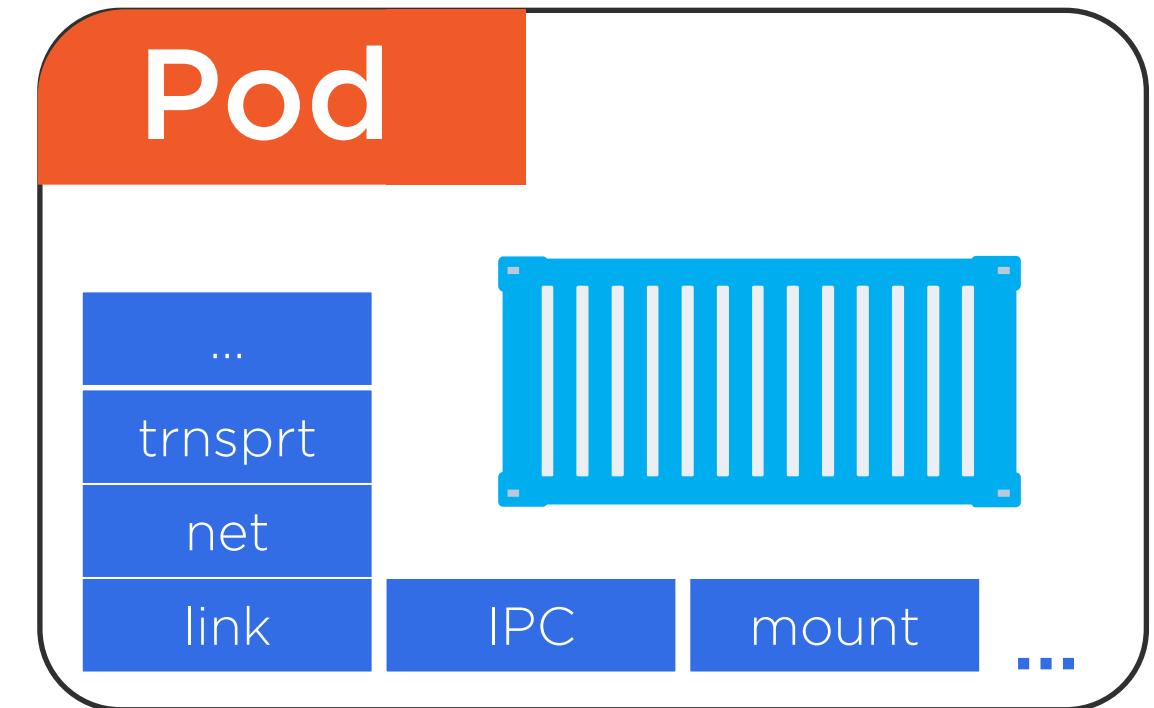


Ring-fenced environment

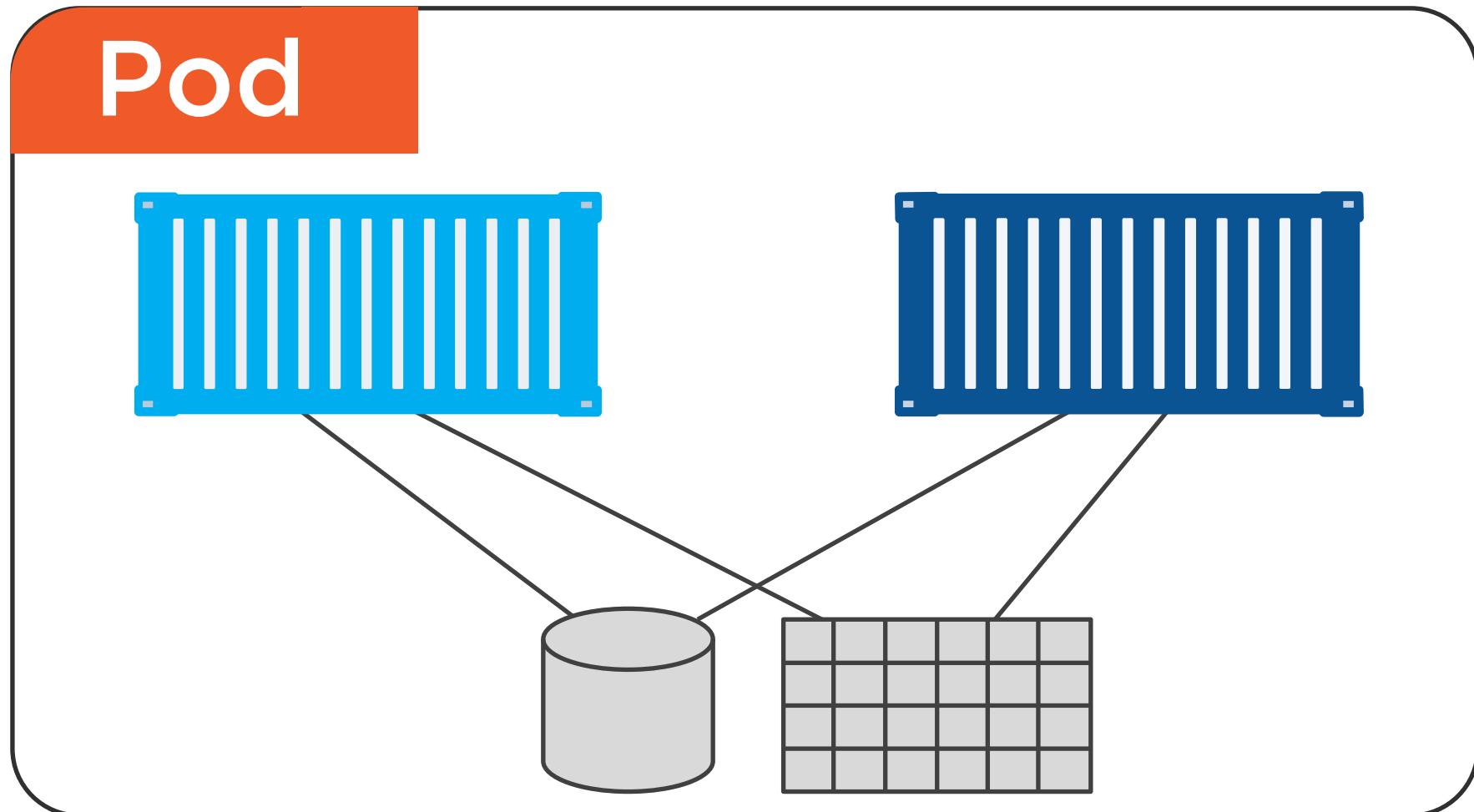
- Network stack
- Kernel namespaces
- ...

n containers

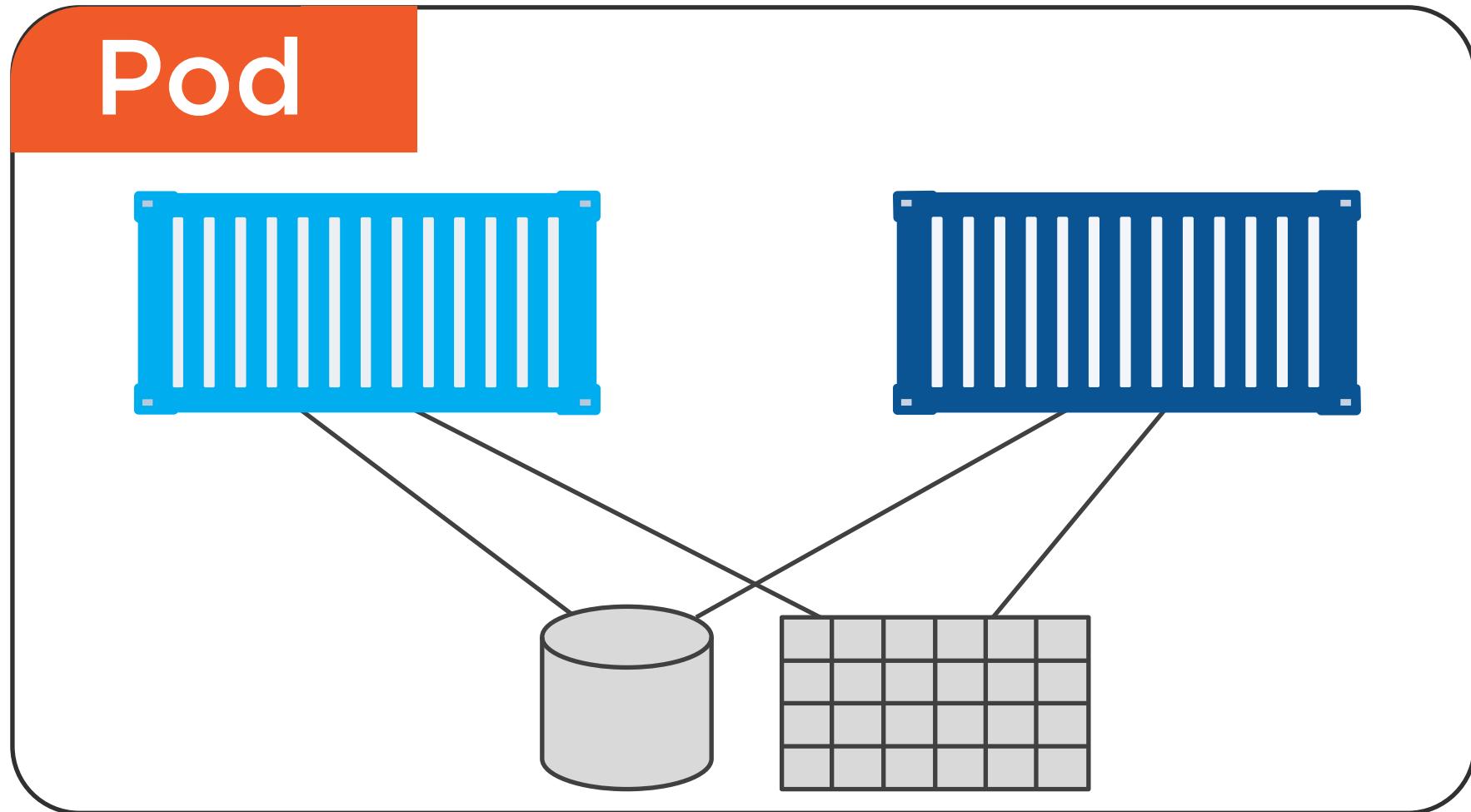
All containers in pod share
the pod environment



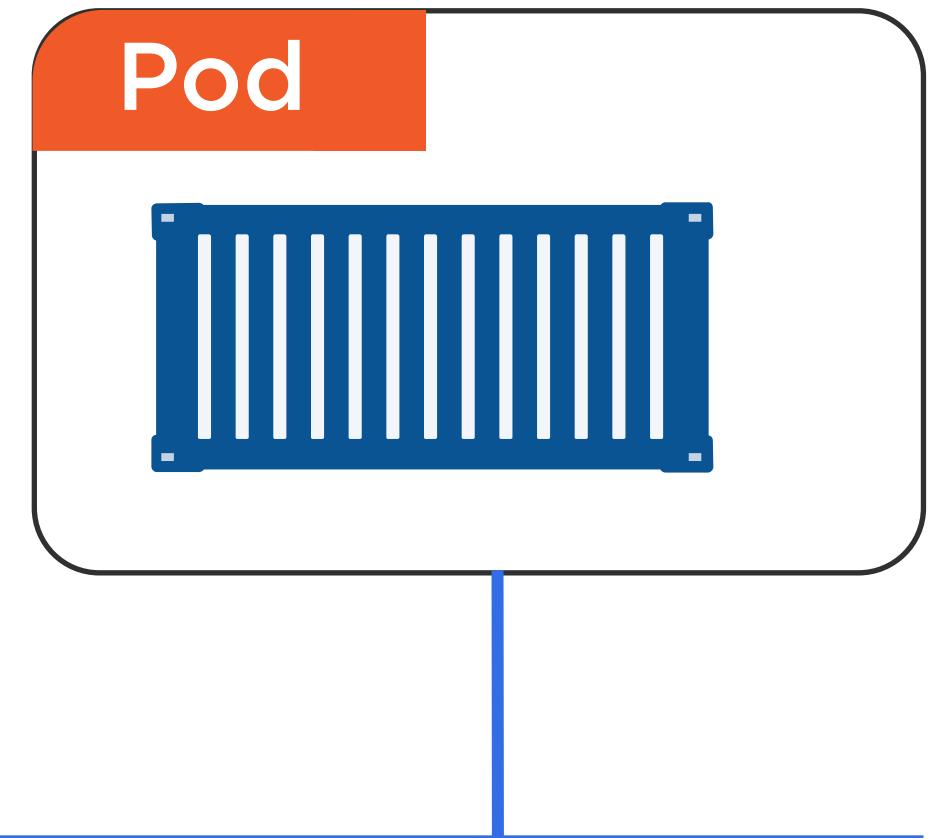
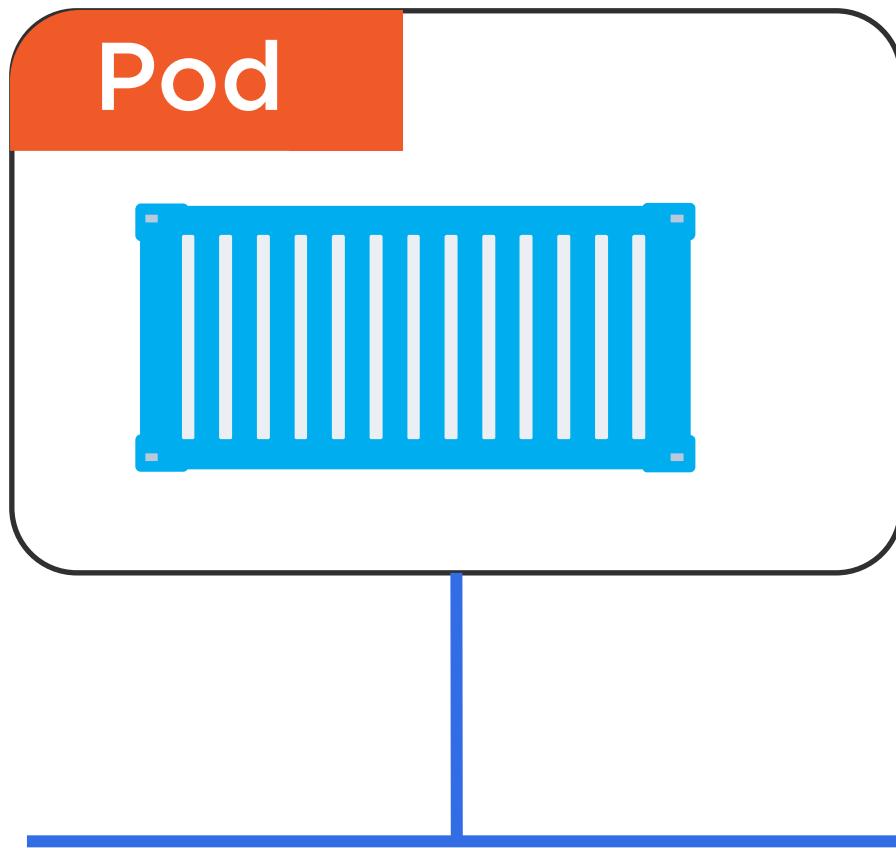
Tight Coupling



Loose Coupling

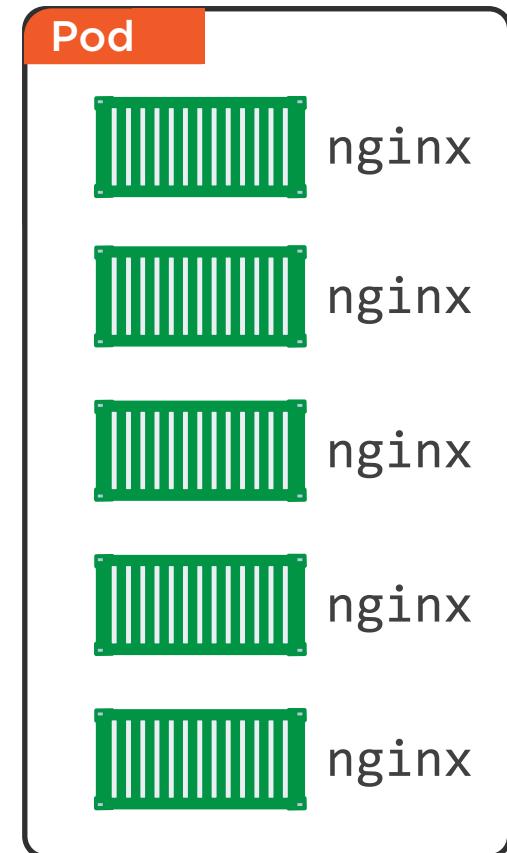
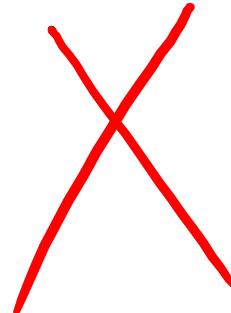
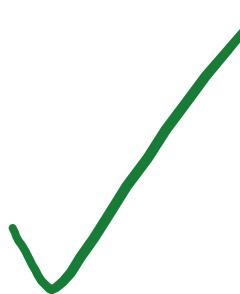


Loose Coupling

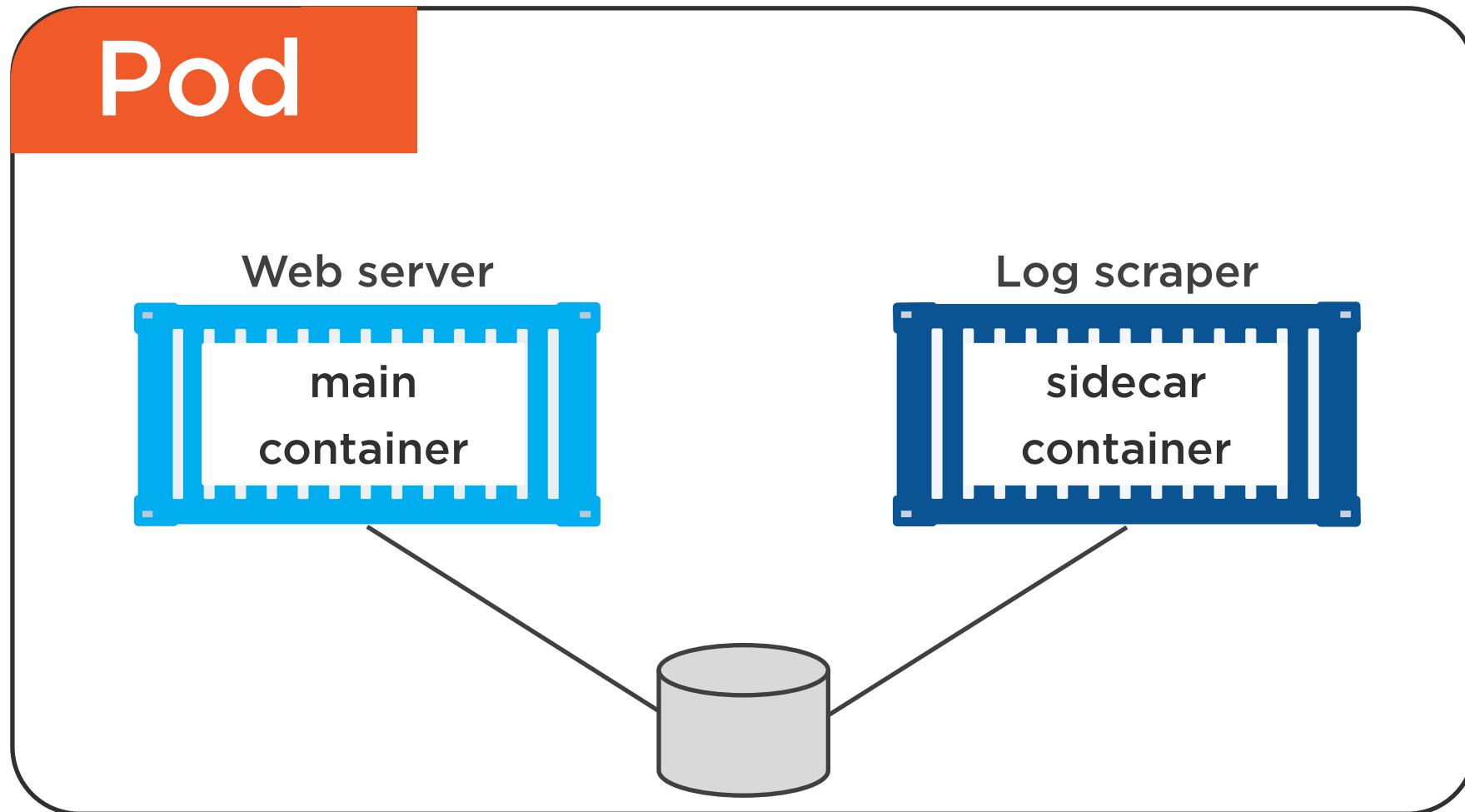




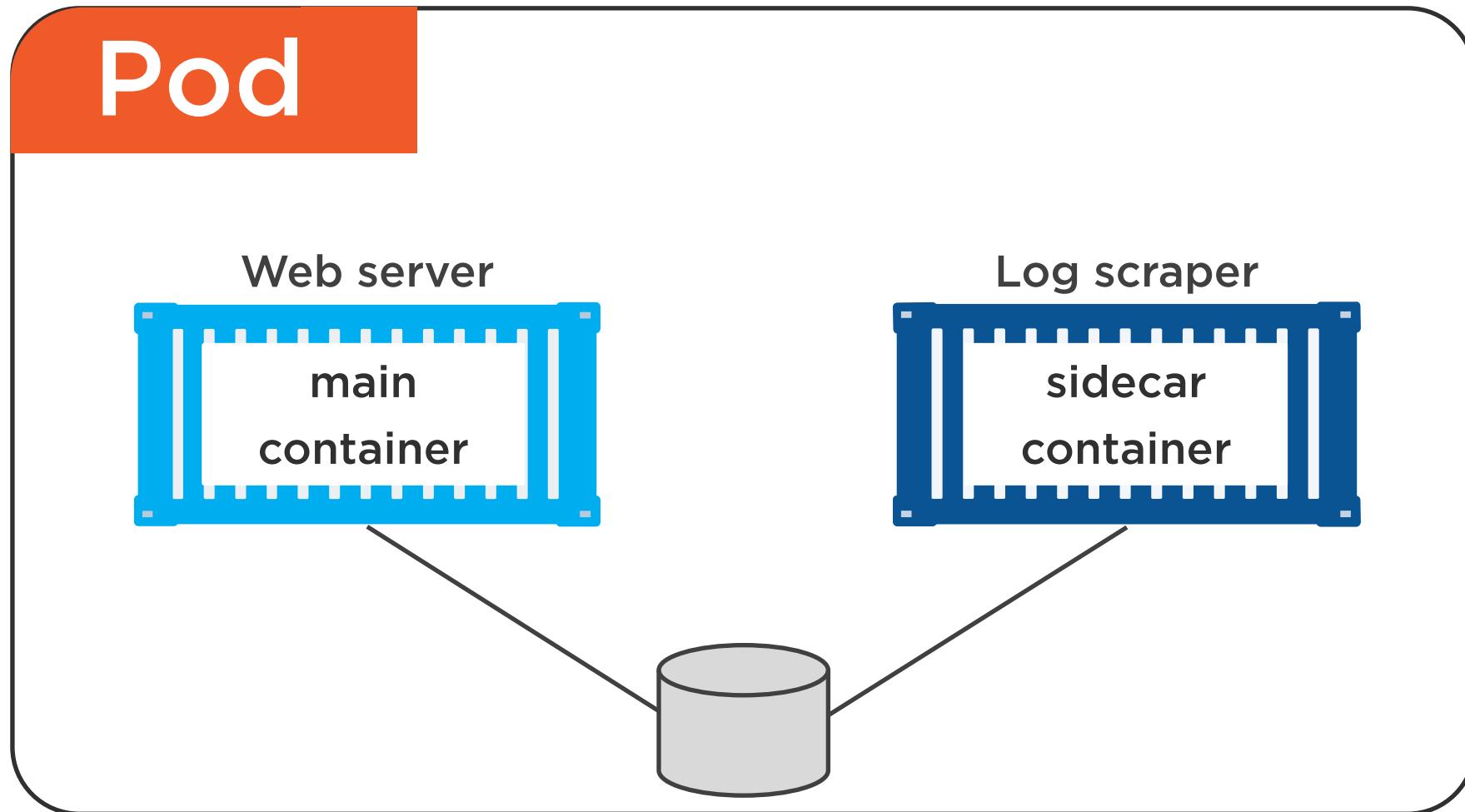
Pods and Scaling



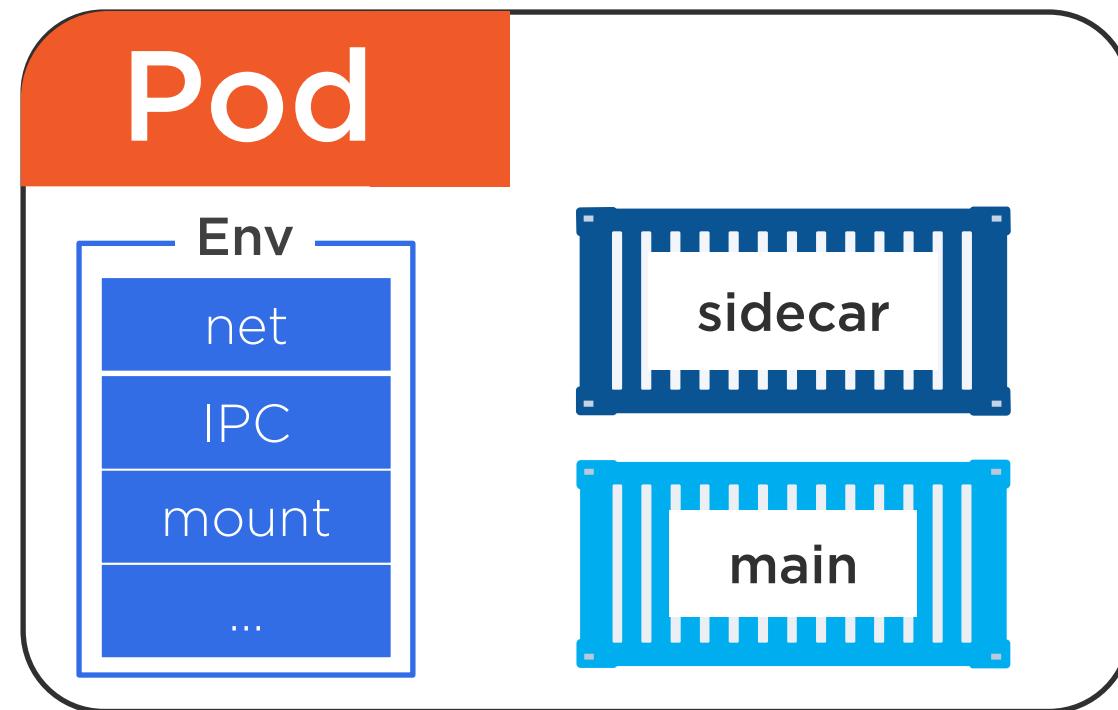
Multi-container Pods



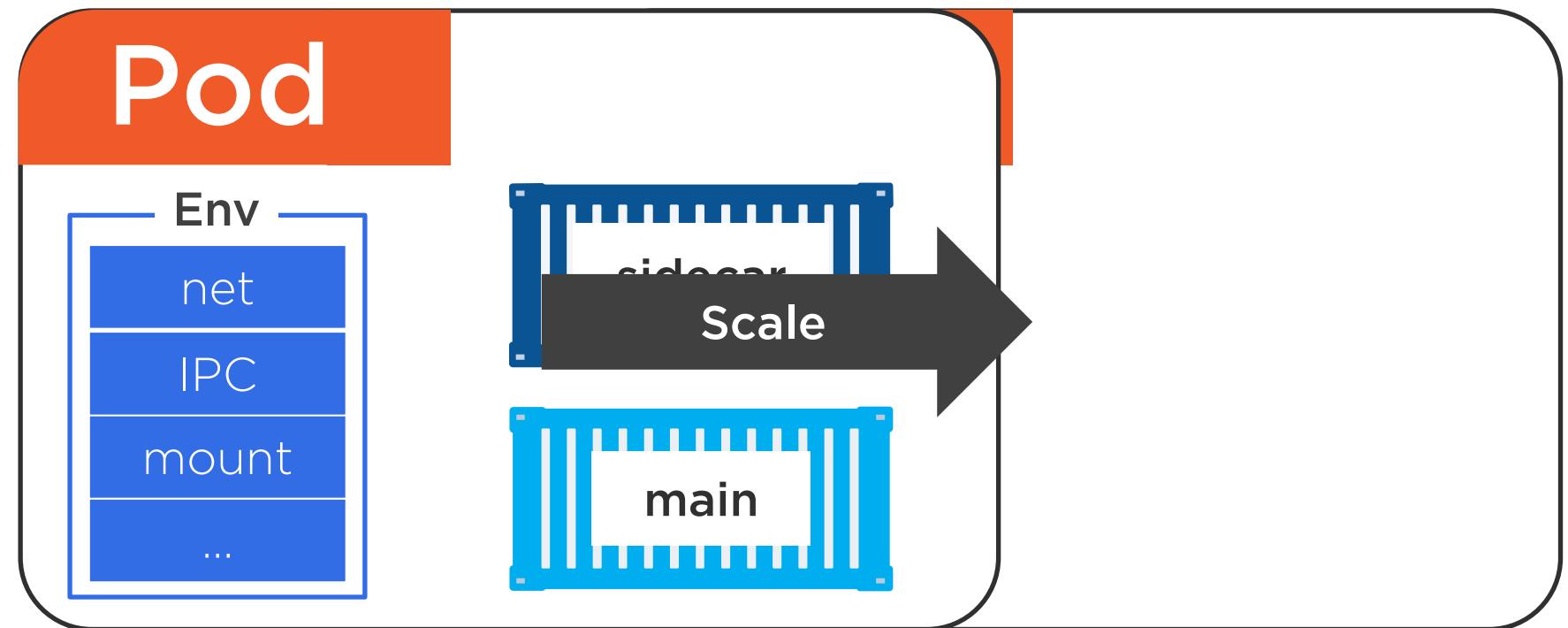
Multi-container Pods



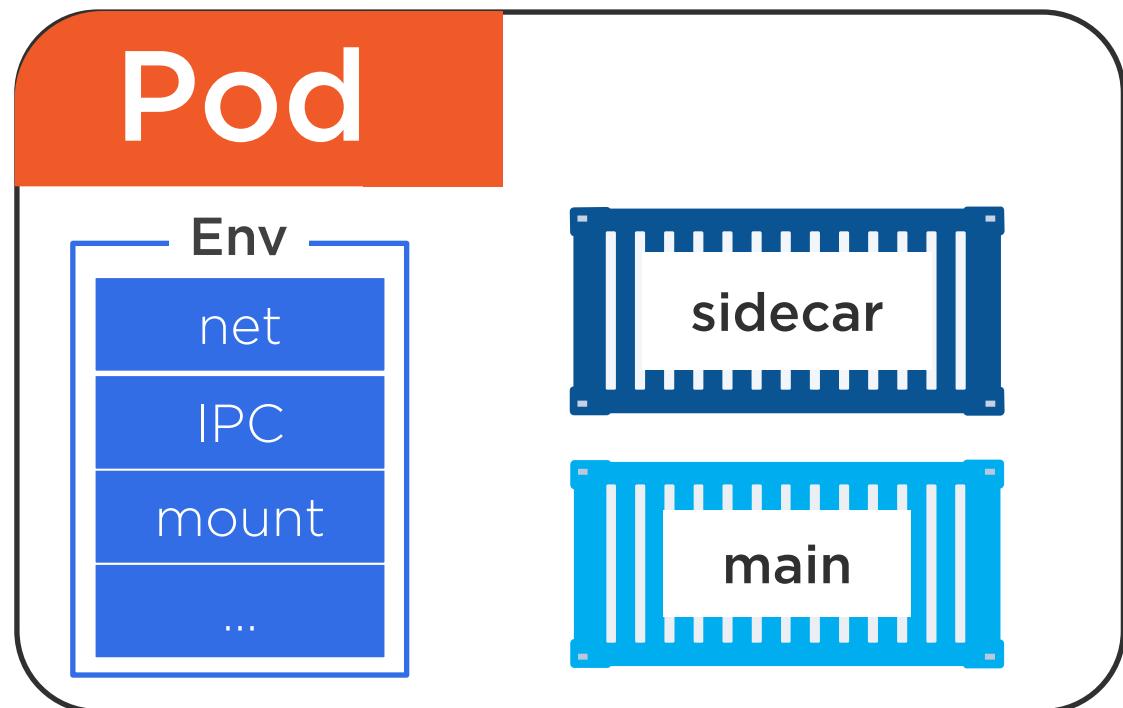
Pods are Atomic



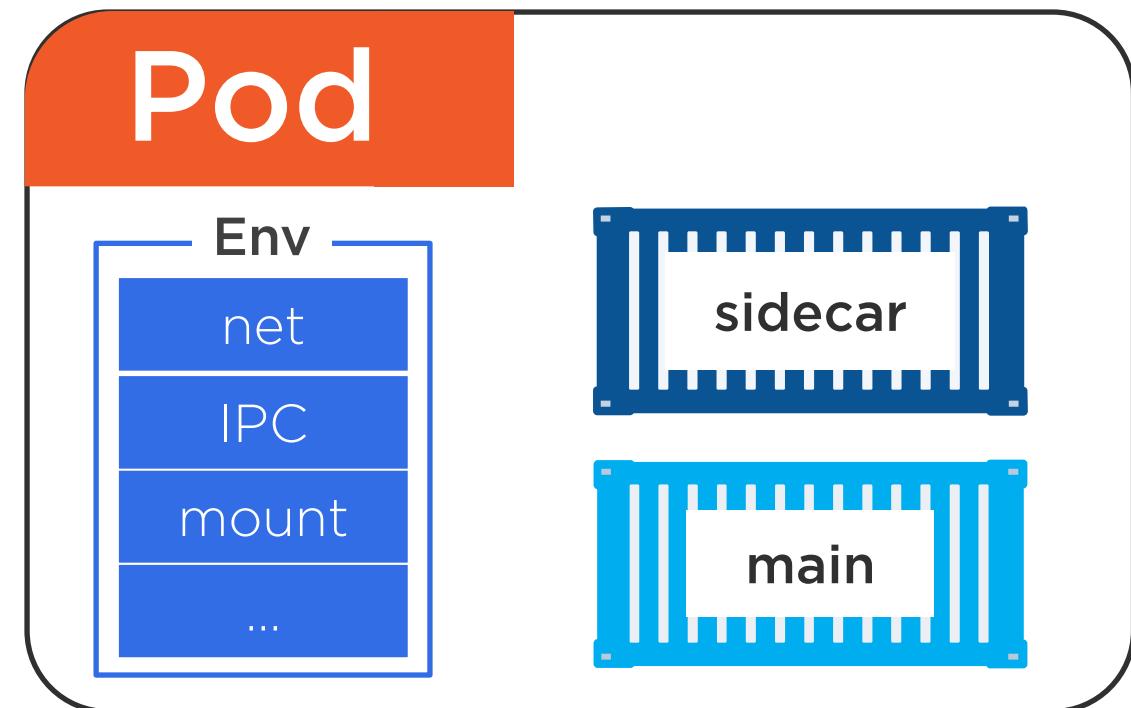
Pods are Atomic



Pods are Atomic



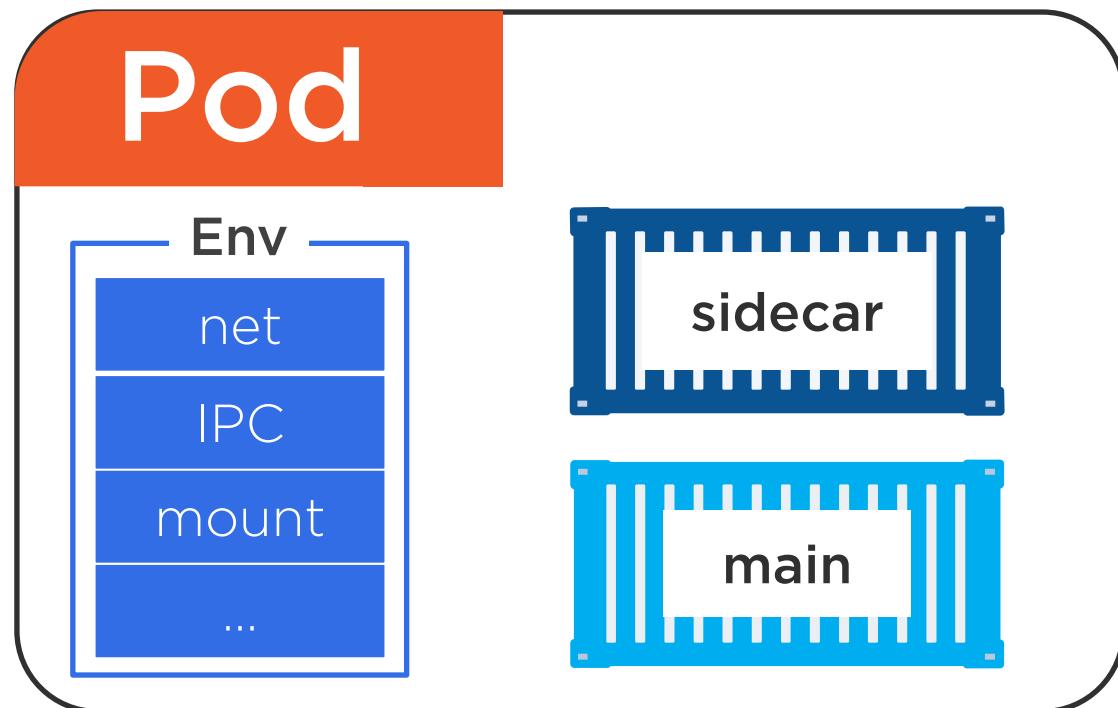
#1 Status:ready



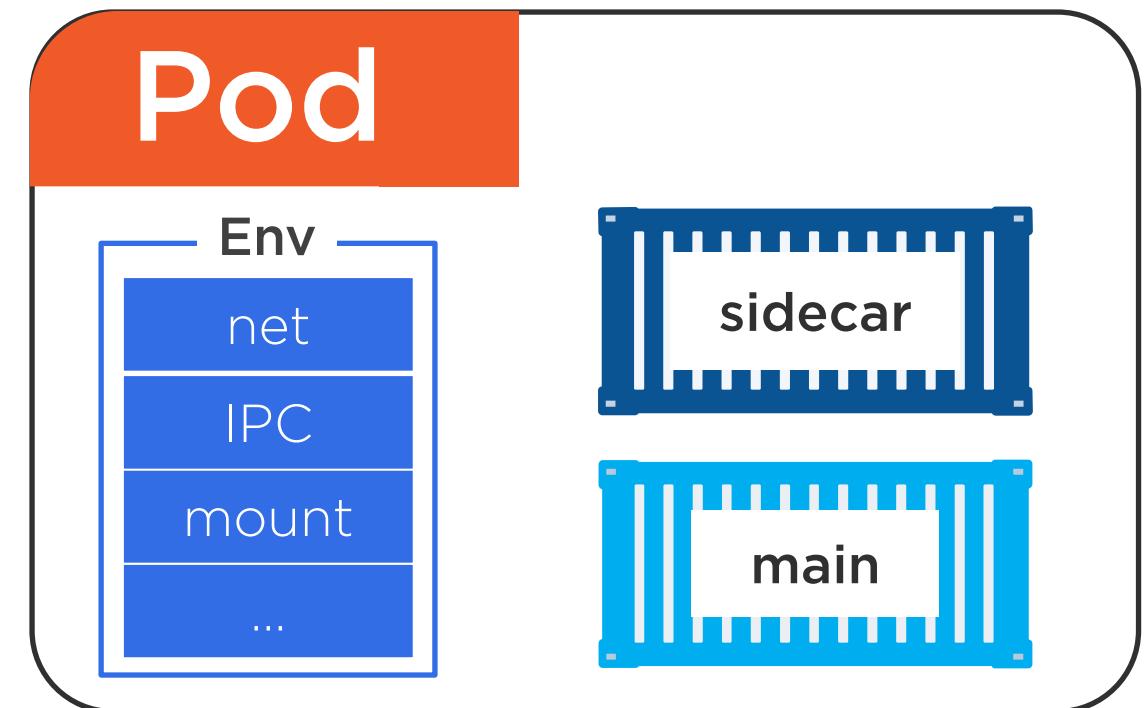
#2 Status:pending



Pods are Atomic



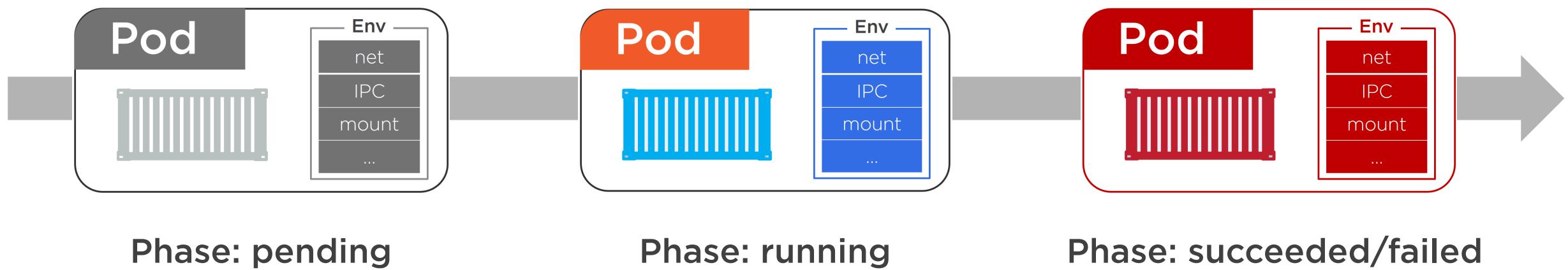
#1 Status:ready



#2 Status:pending



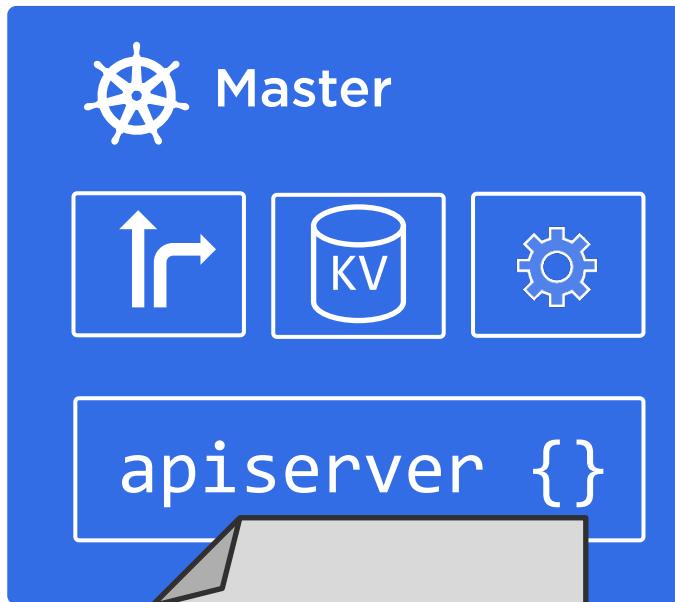
Pod Lifecycle



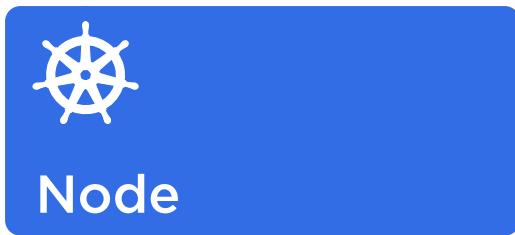
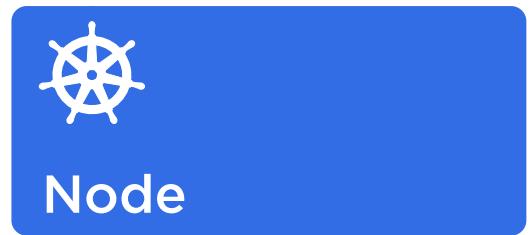
Deploying Pods

Usually via higher level objects



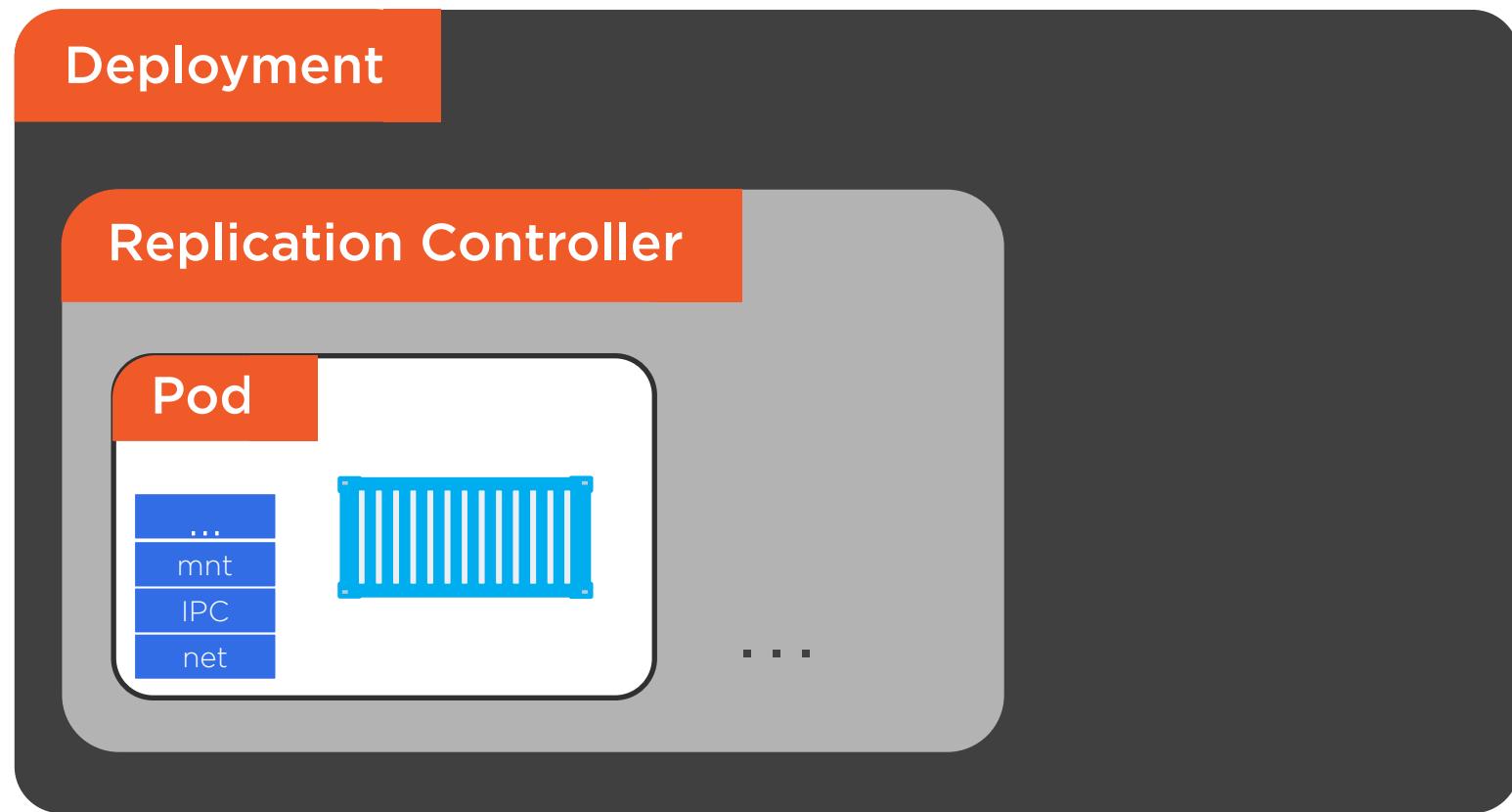


```
apiVersion: v1
kind: Pod
metadata: ✓
```

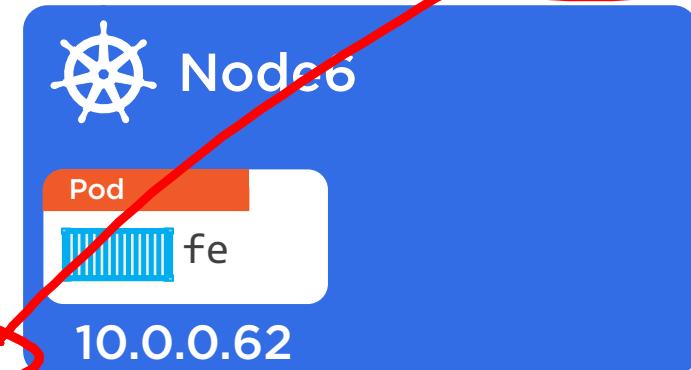
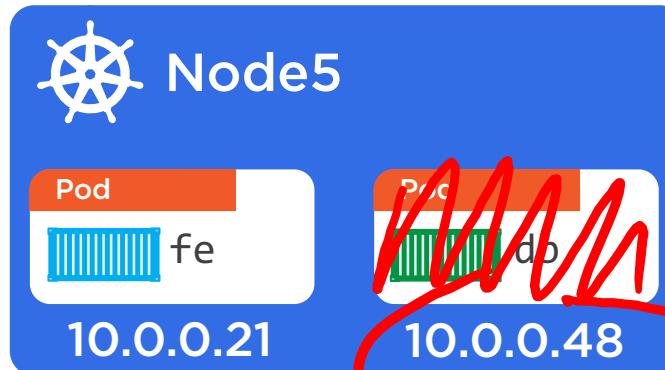
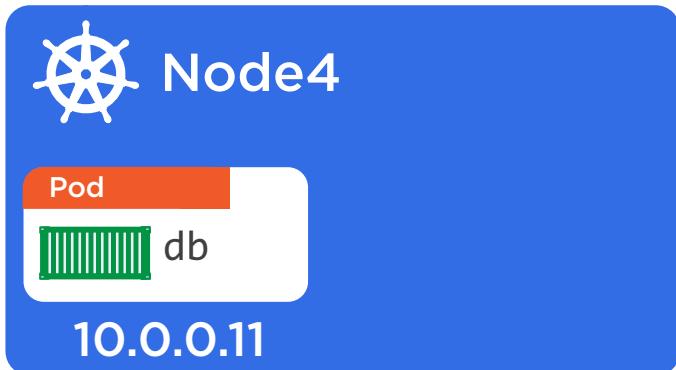
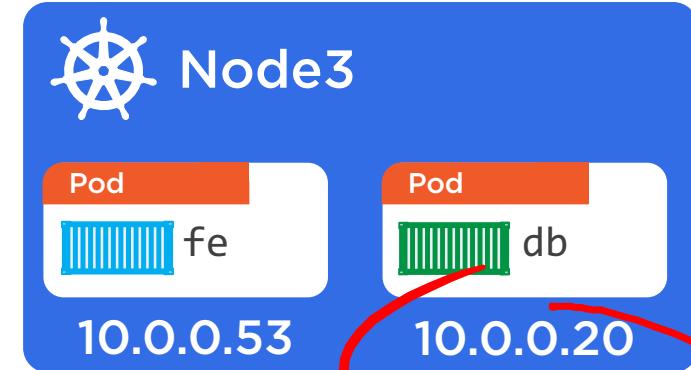
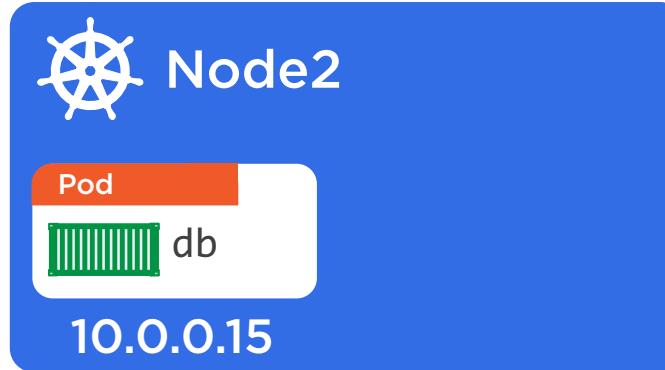
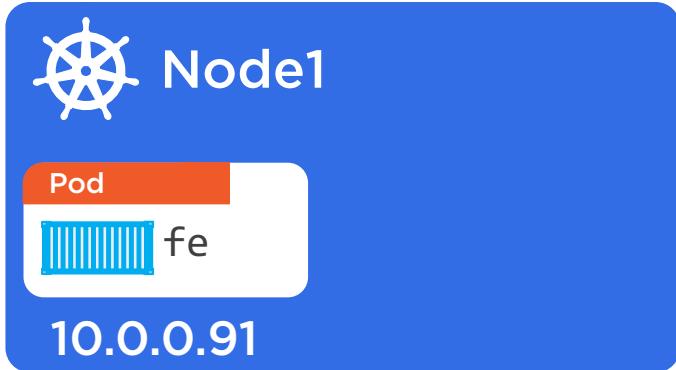


Deploying Pods

Usually via higher level objects

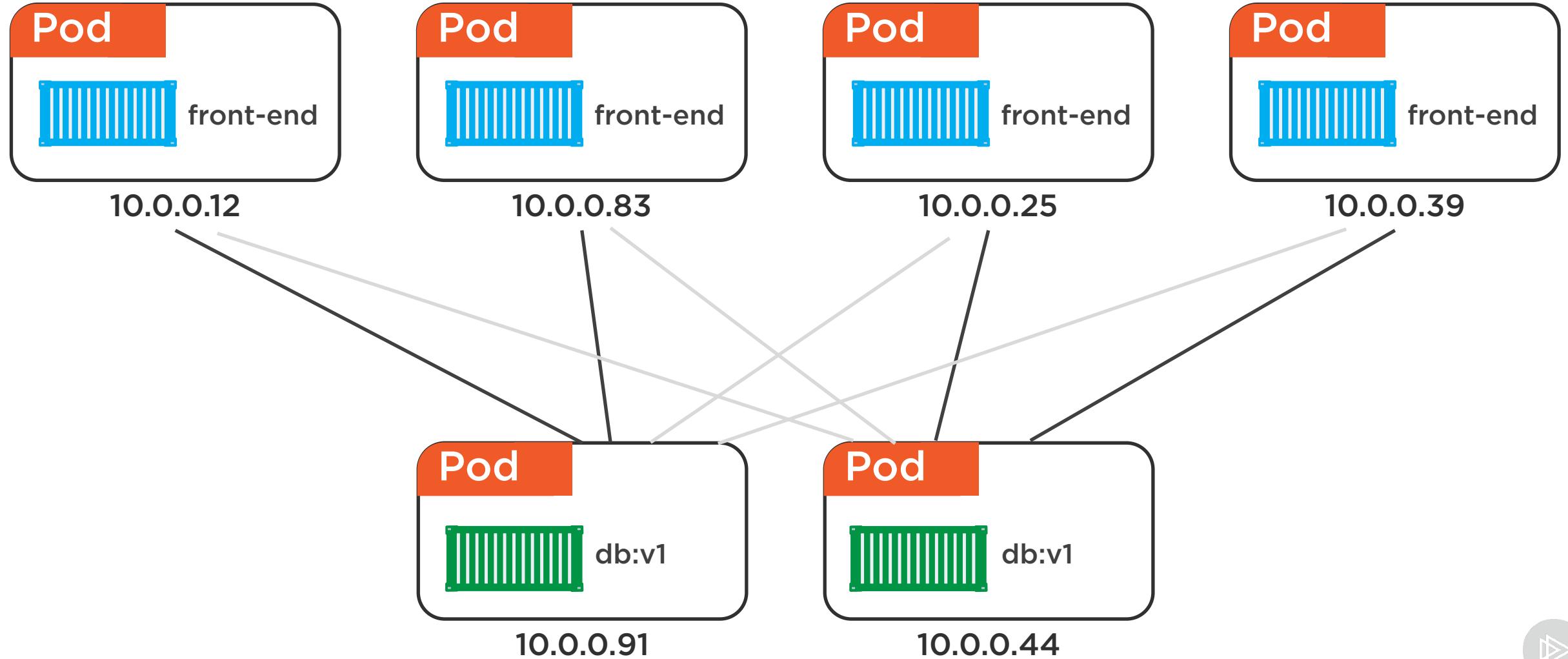


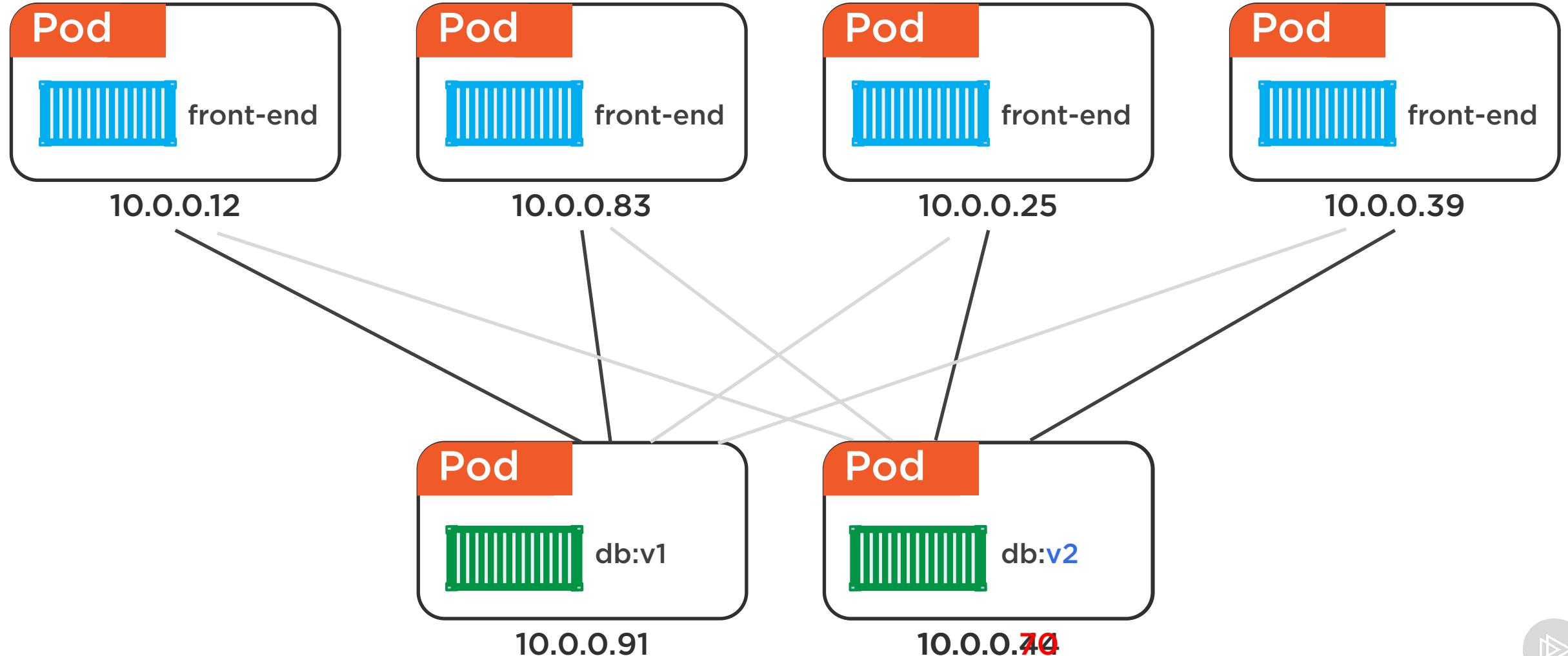
Services

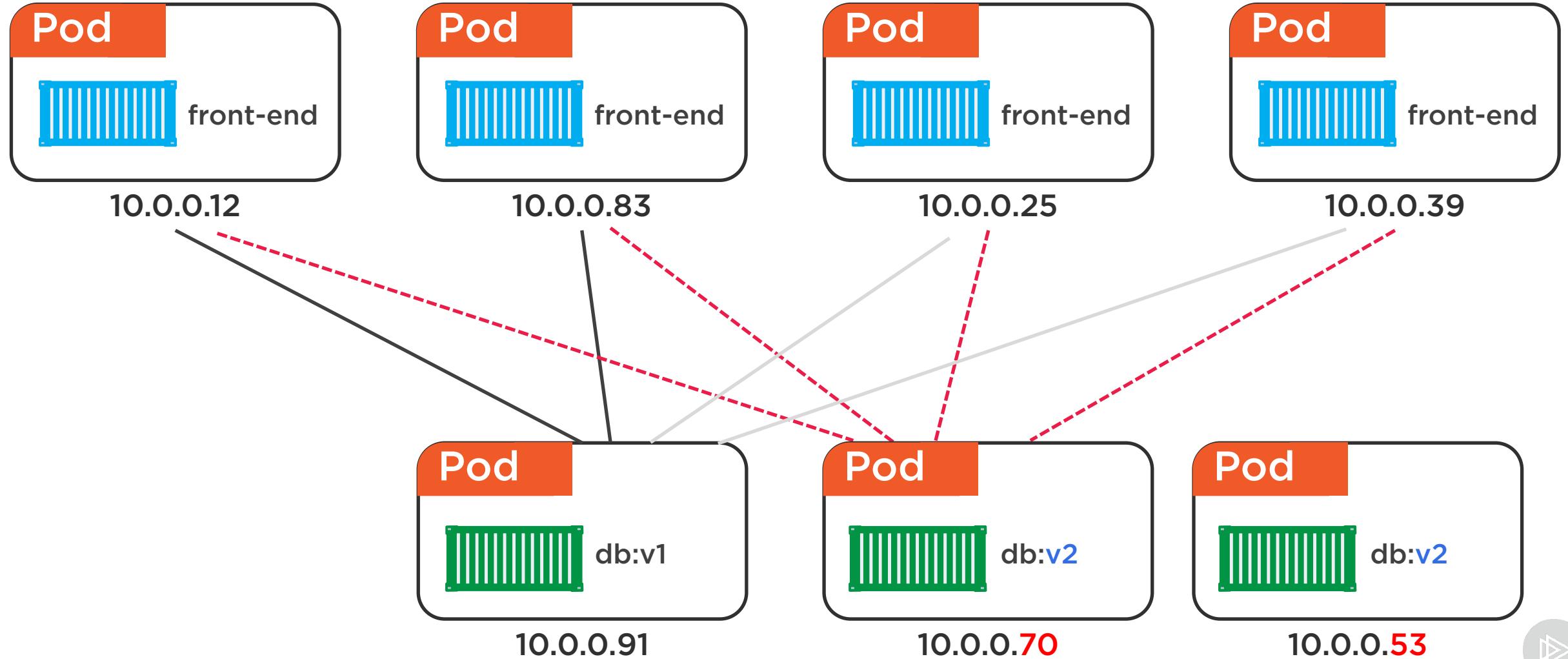


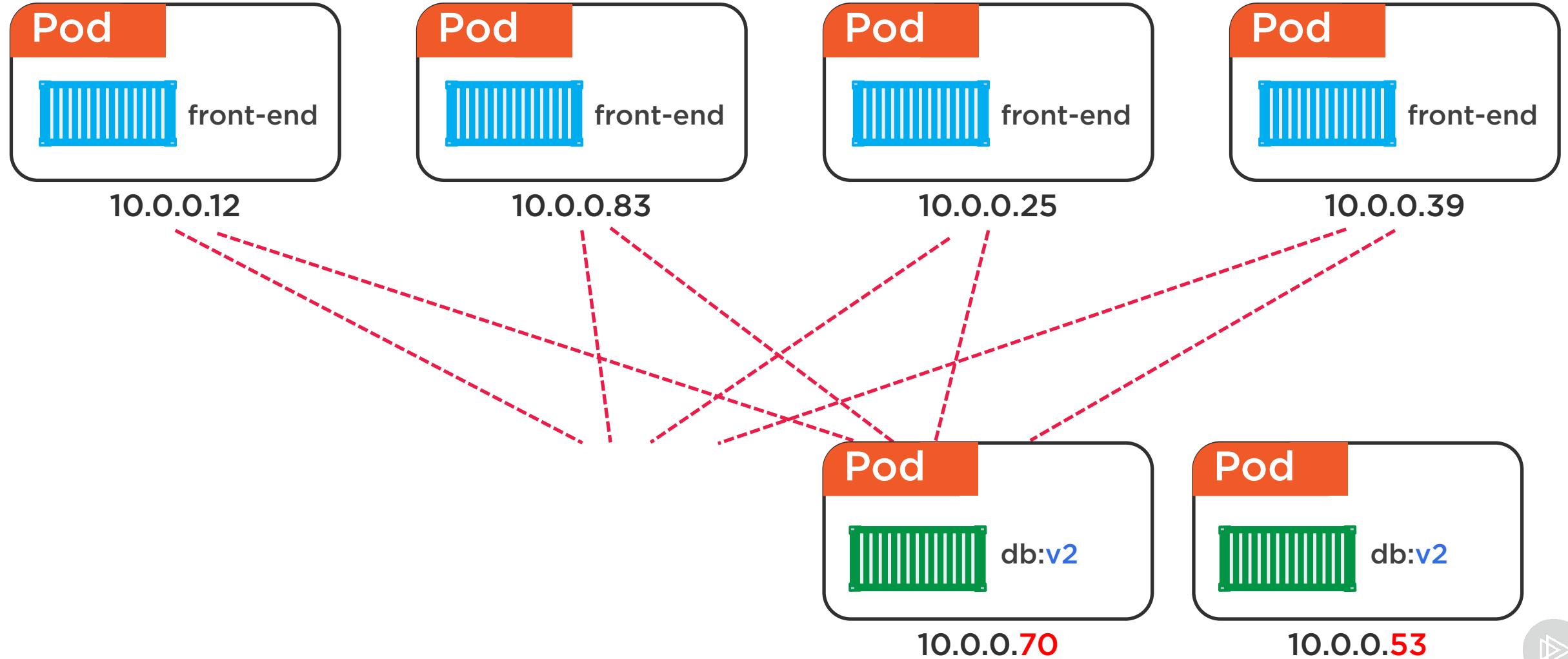
Every new pod gets a new IP = IP churn!

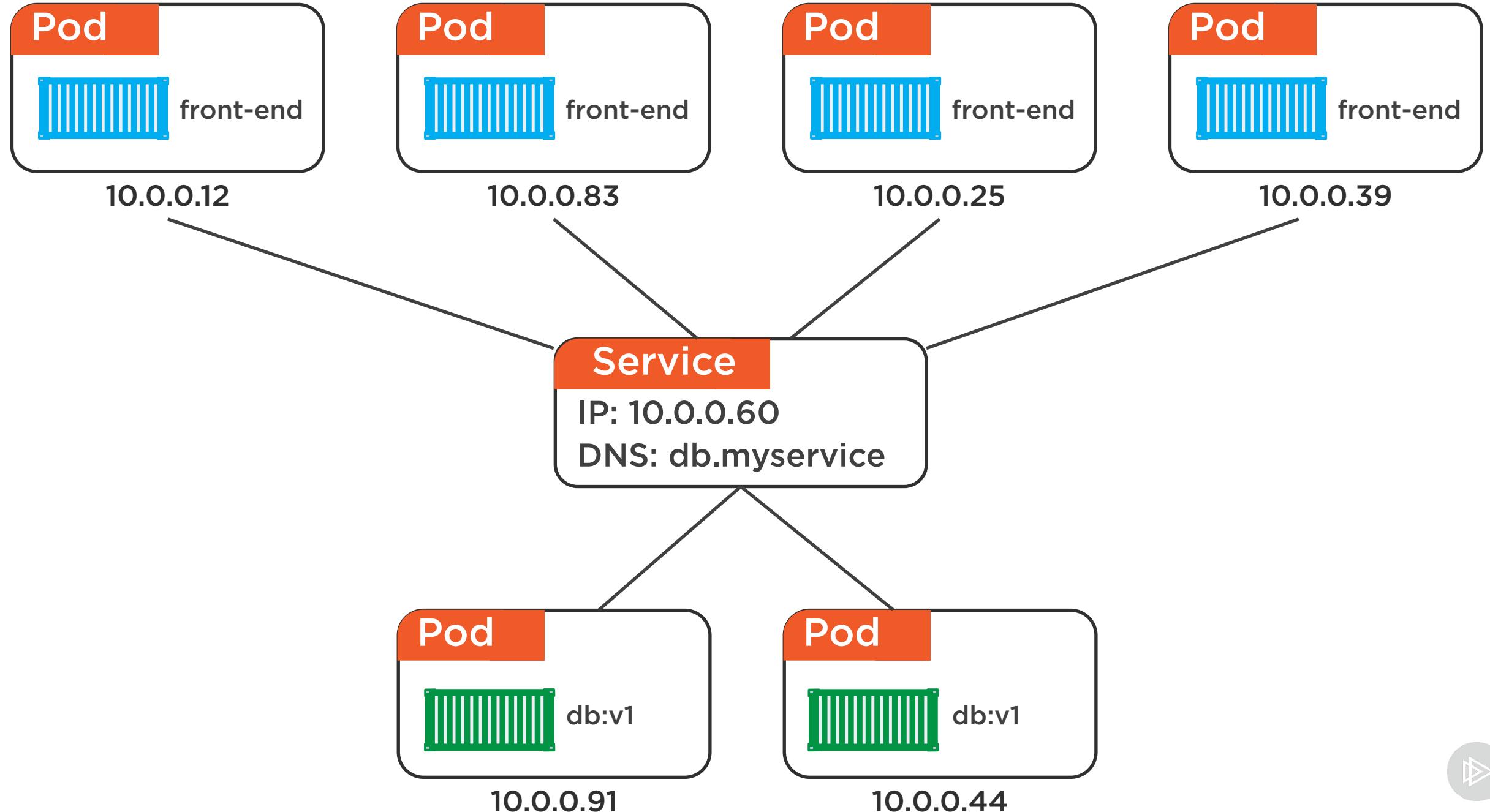


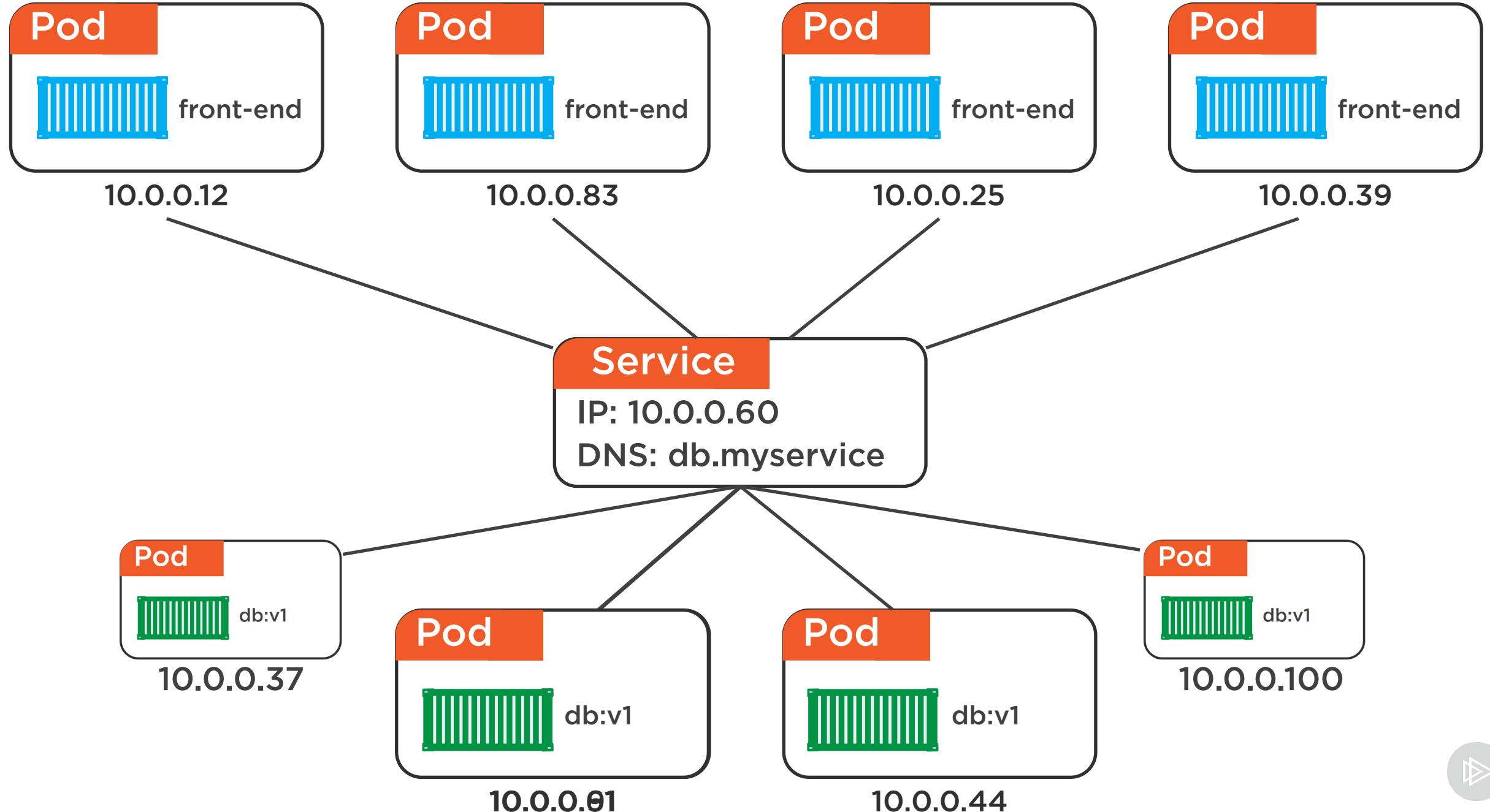


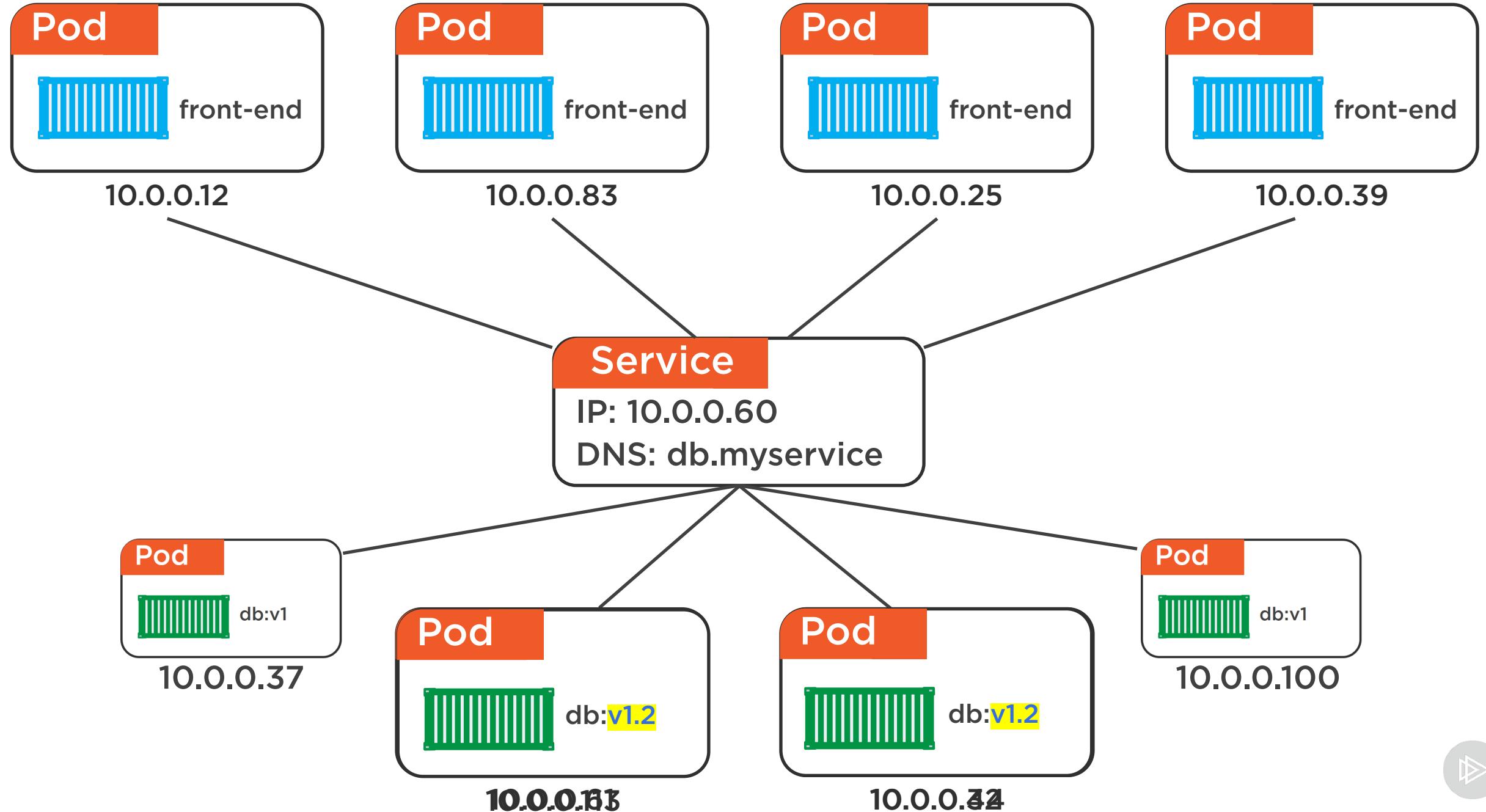




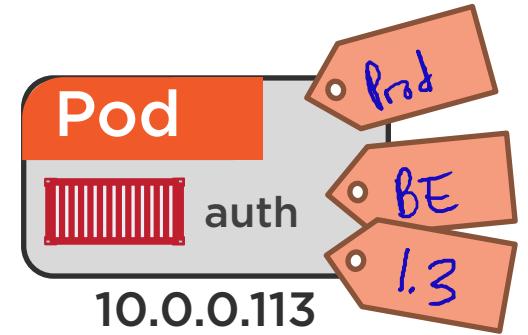
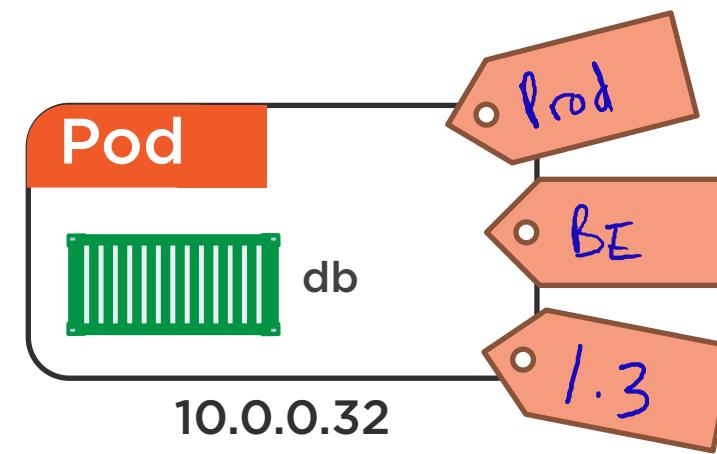
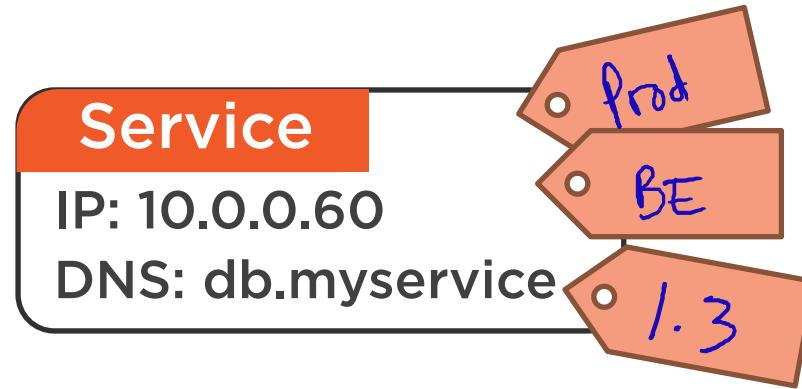
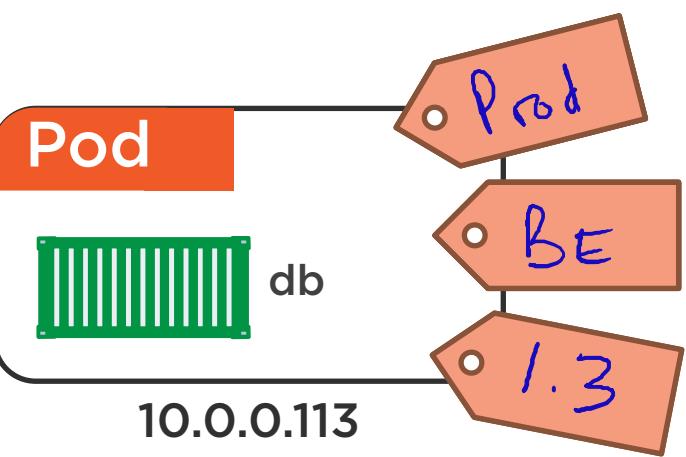


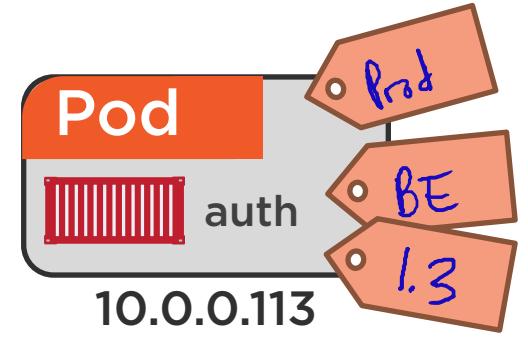
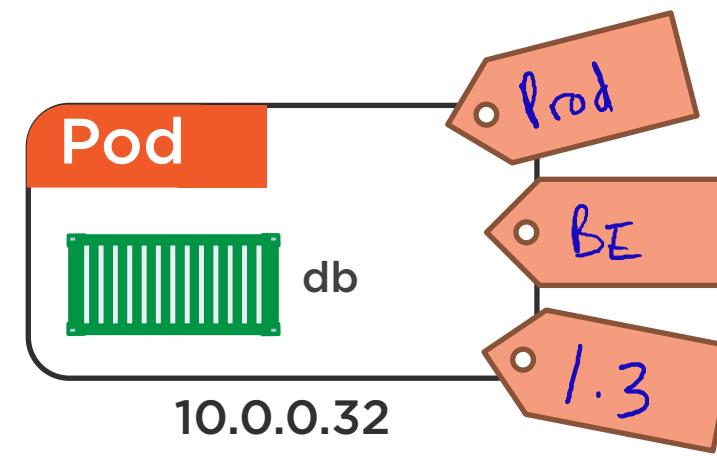
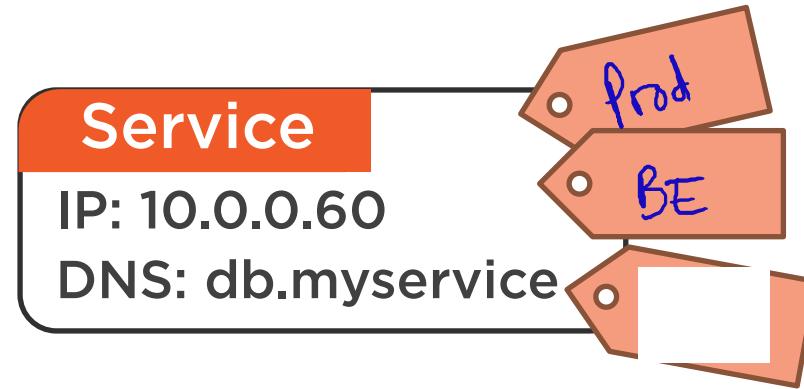
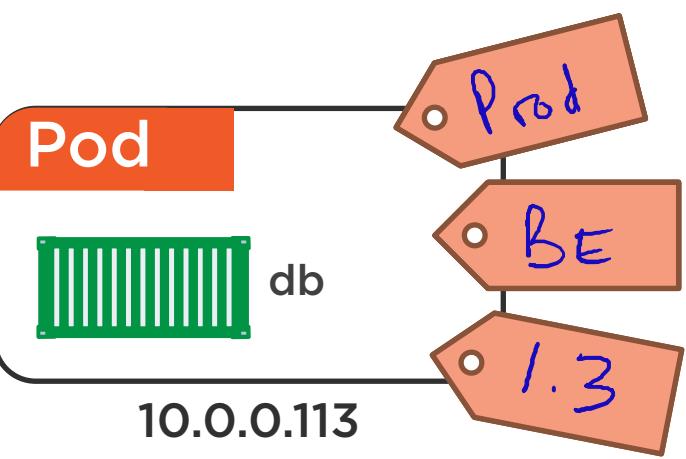


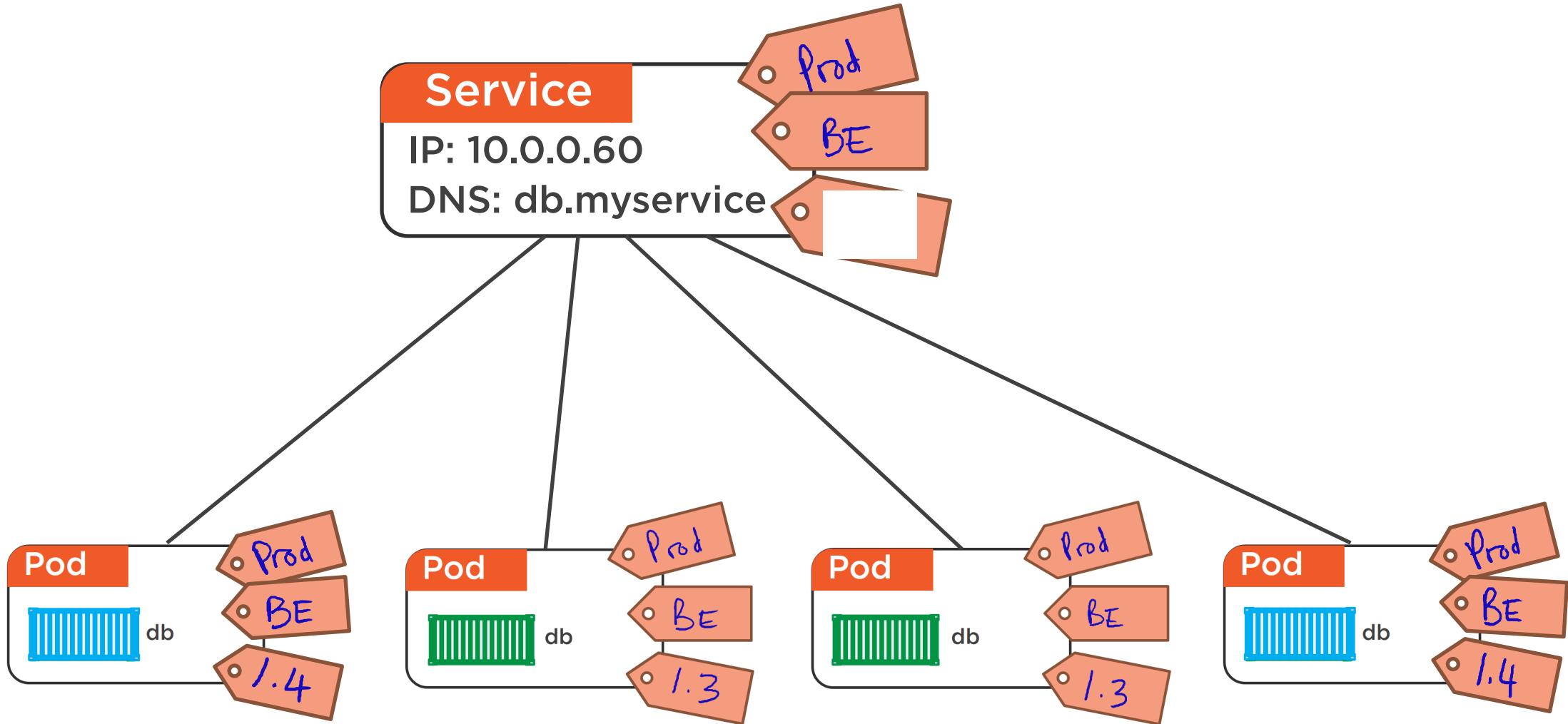


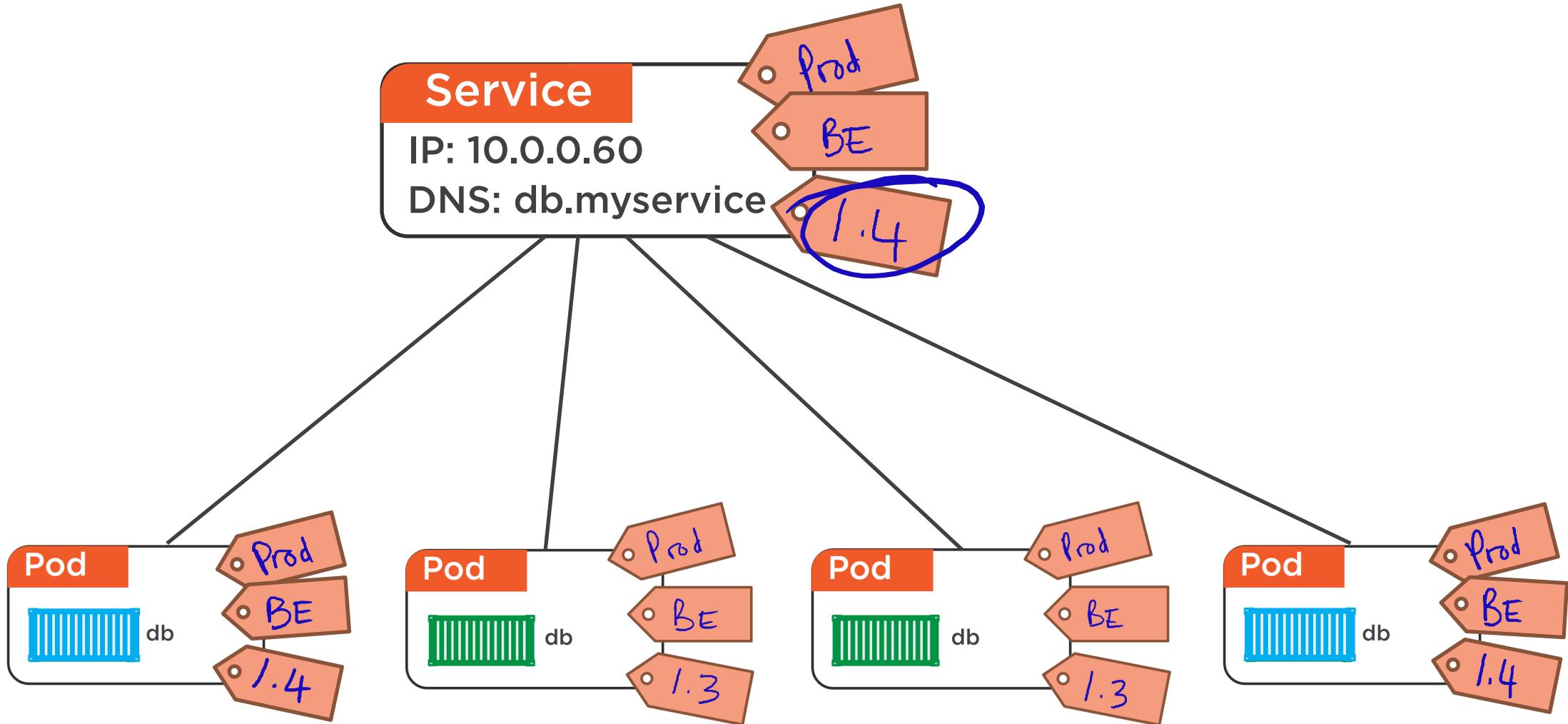




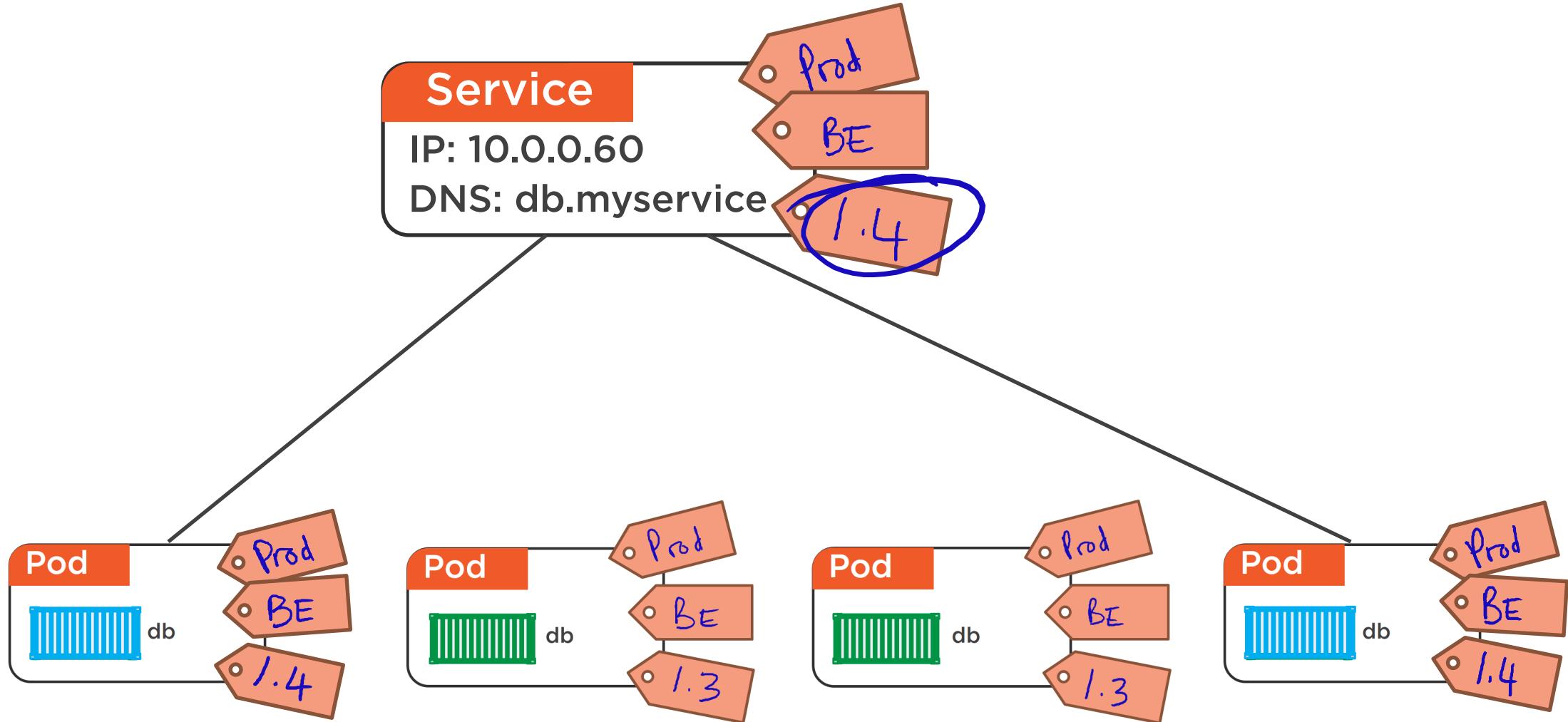




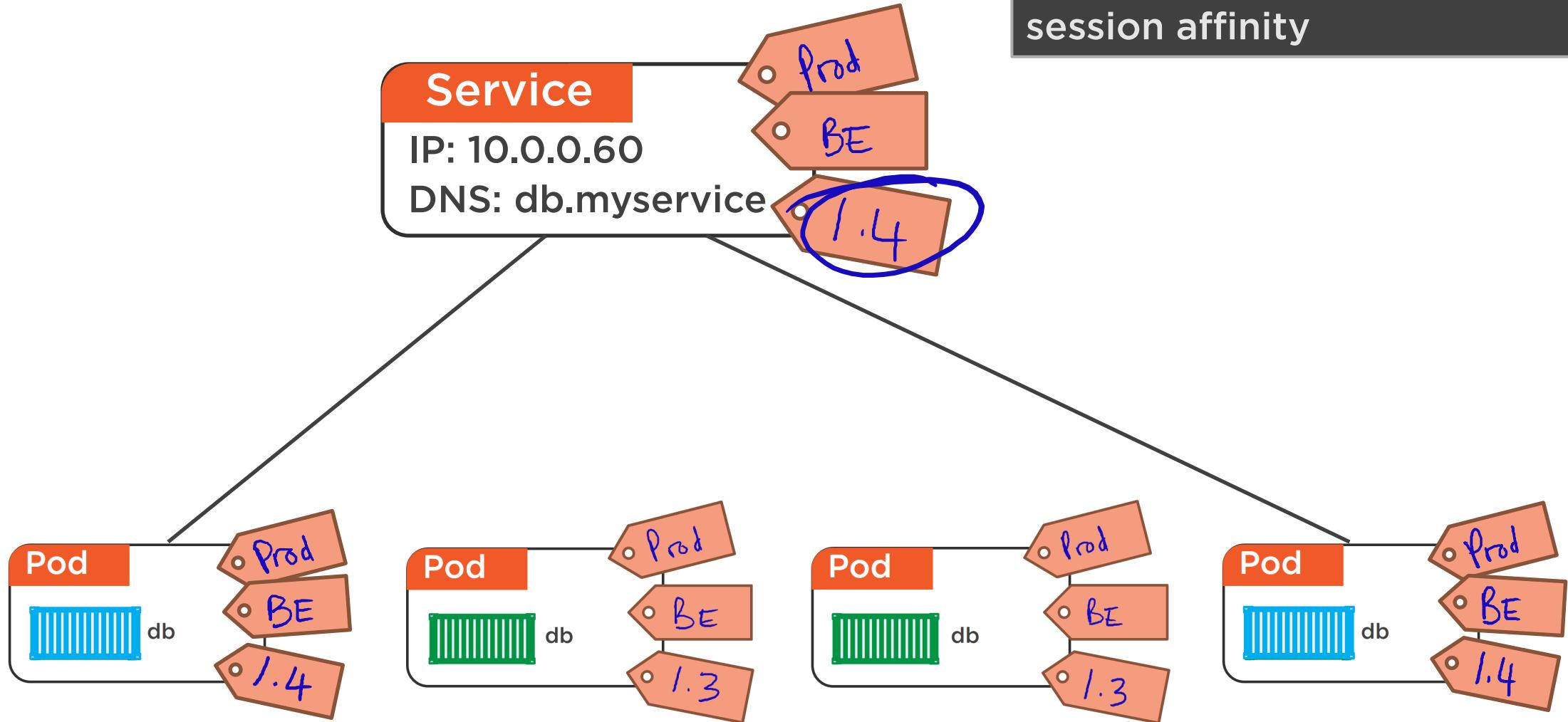




Only send to healthy pods



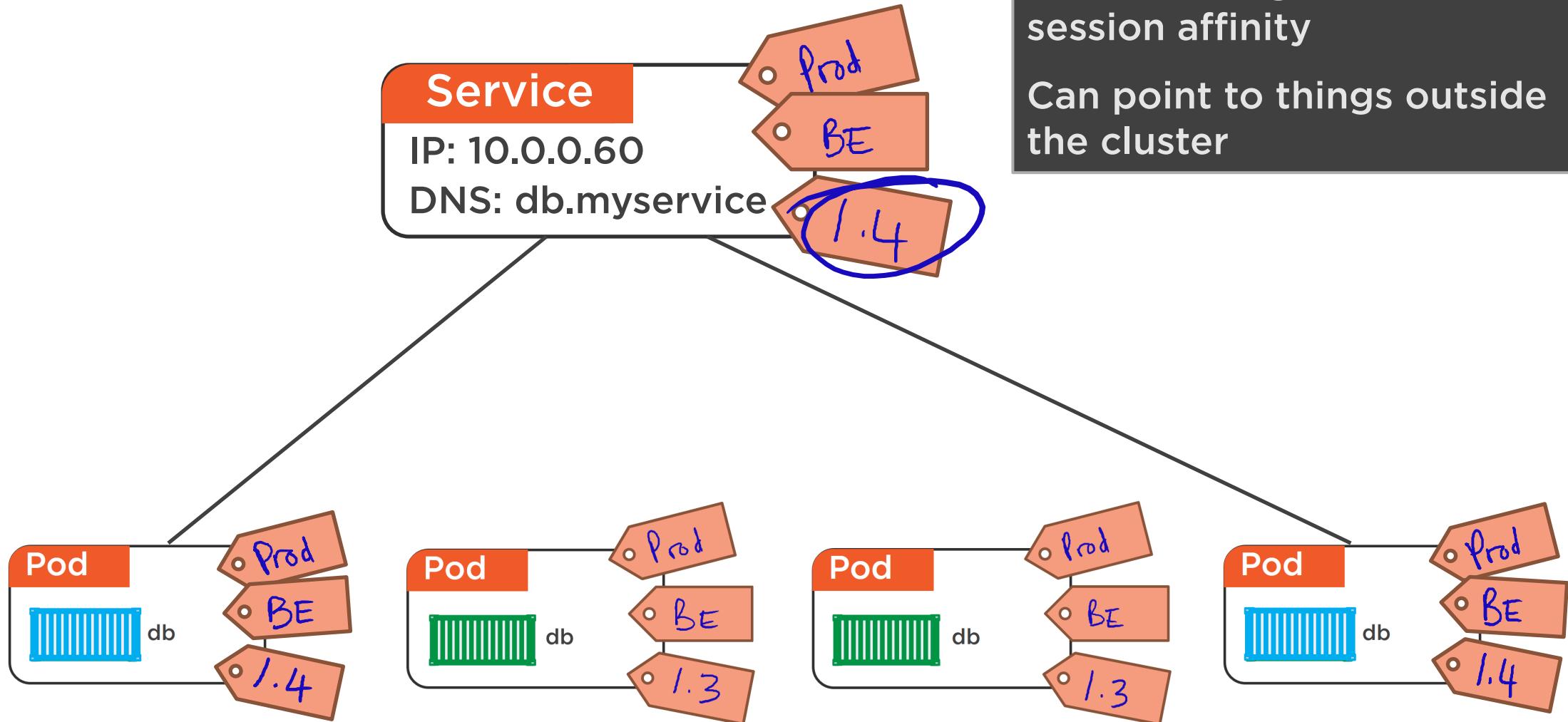
Only send to healthy pods
Can be configured for session affinity

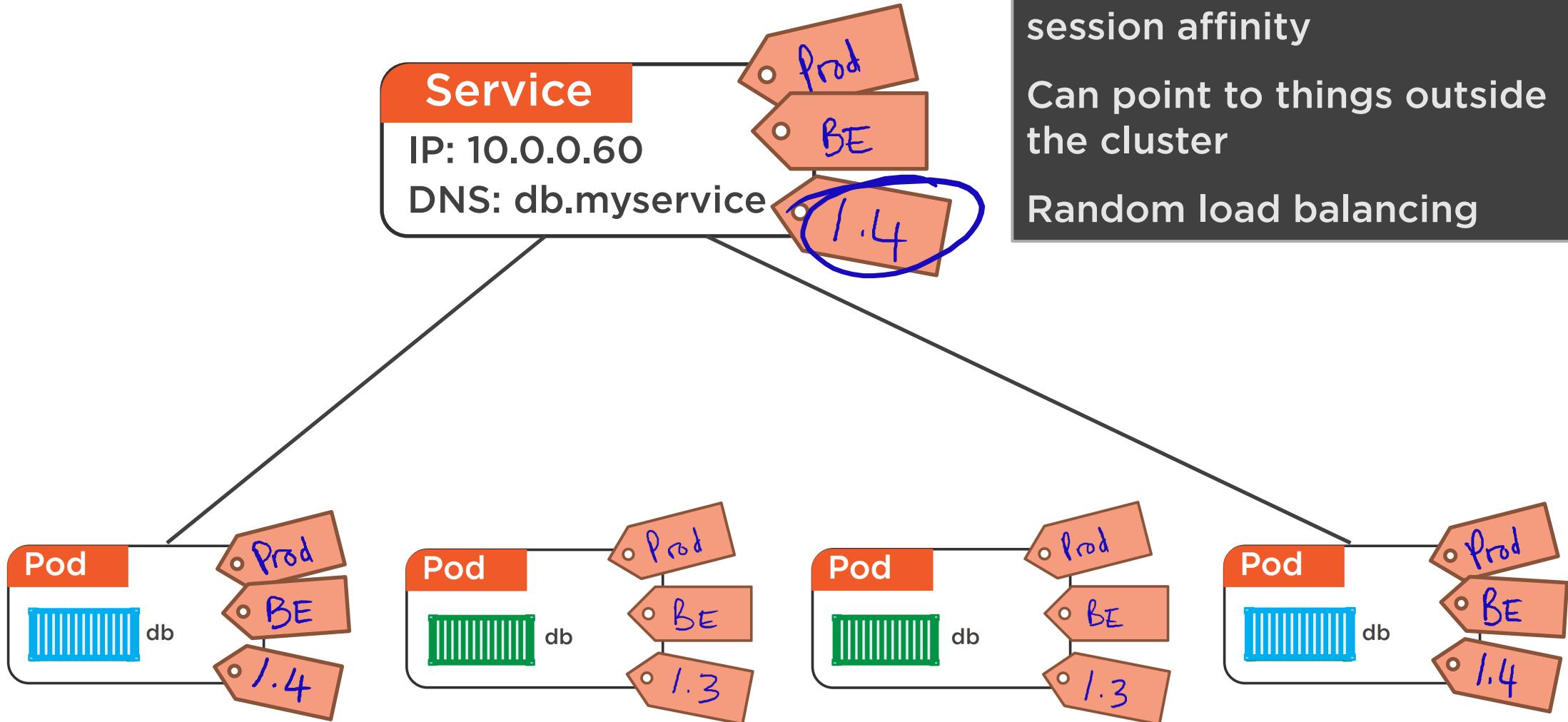


Only send to healthy pods

Can be configured for session affinity

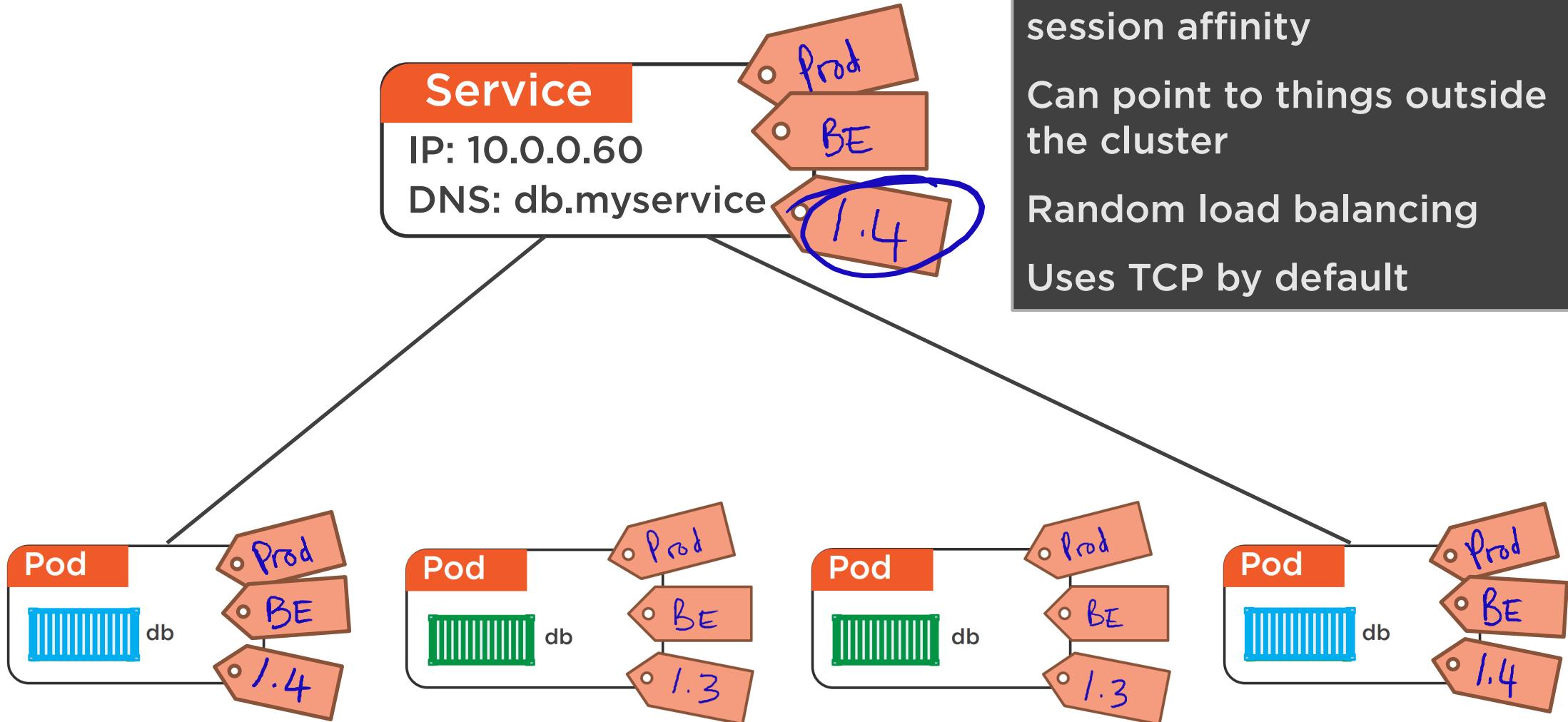
Can point to things outside the cluster





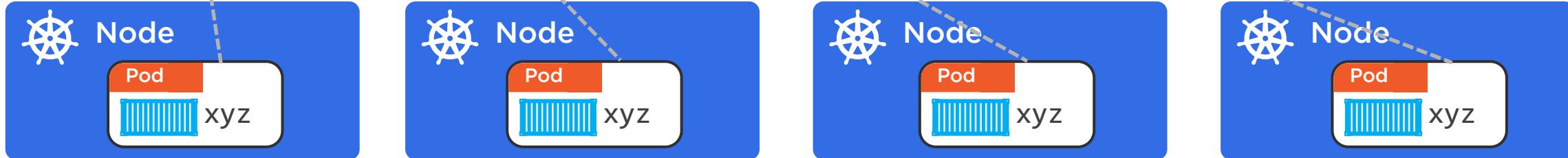
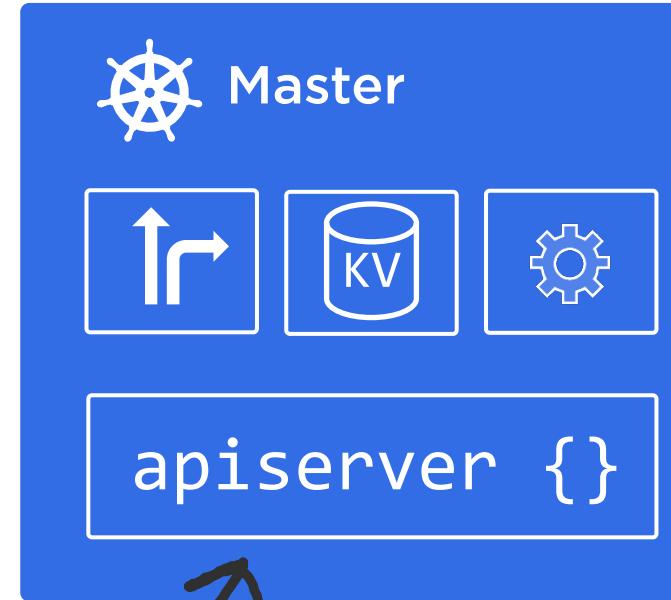
Only send to healthy pods
Can be configured for session affinity
Can point to things outside the cluster
Random load balancing



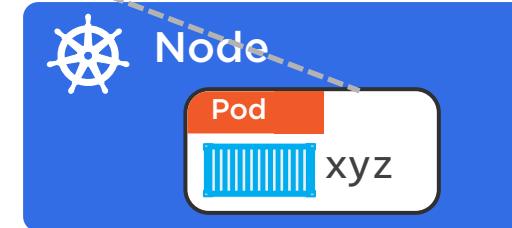
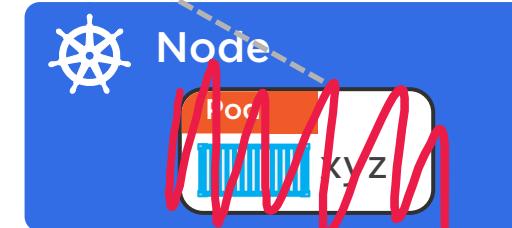
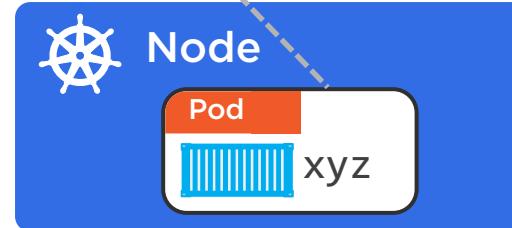
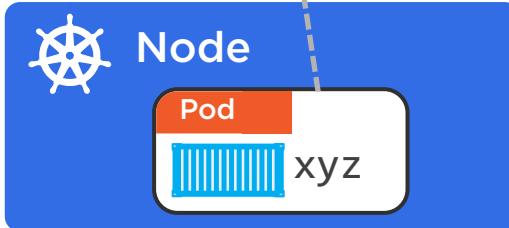
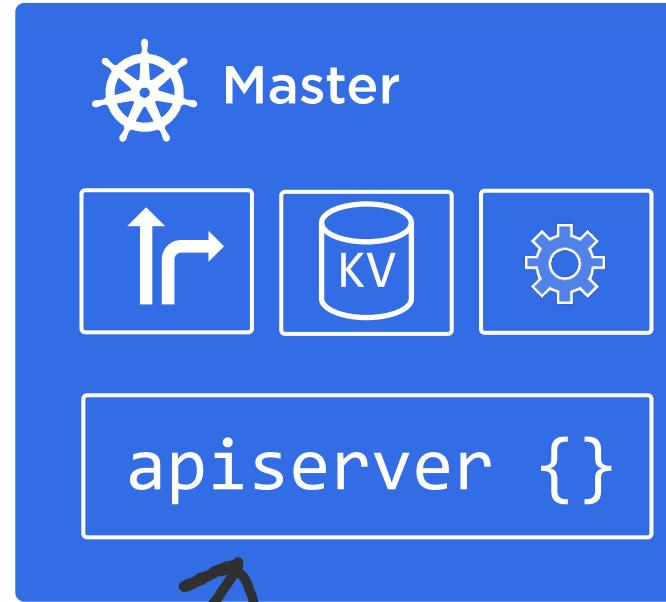


Deployments

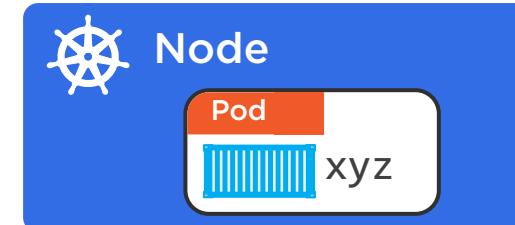
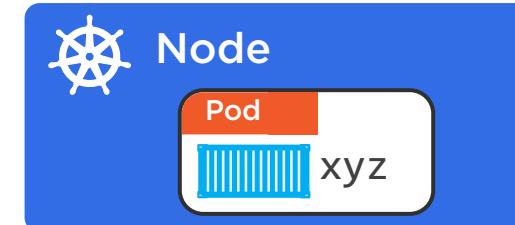
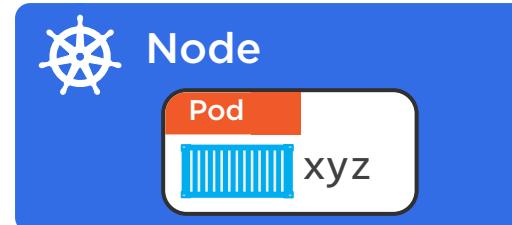
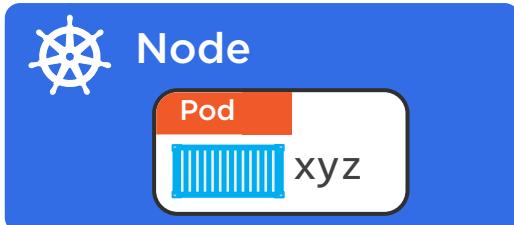
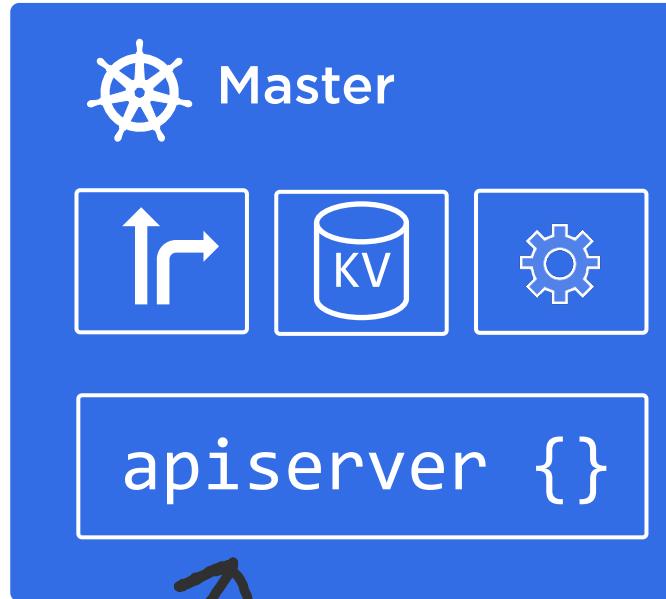
```
apiVersion: v1
kind: ReplicationController
metadata:
  name: xyz
spec:
  replicas: 4
```



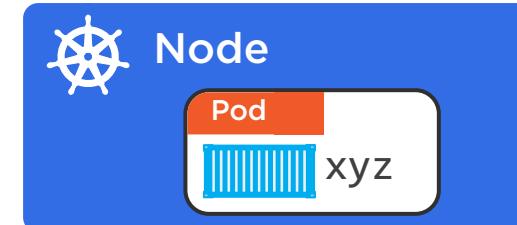
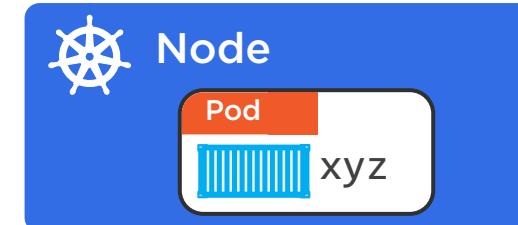
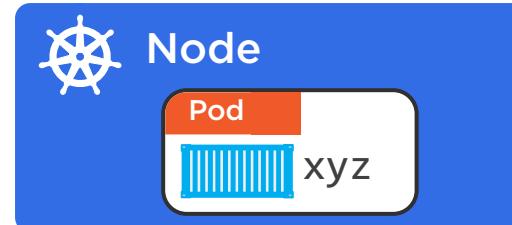
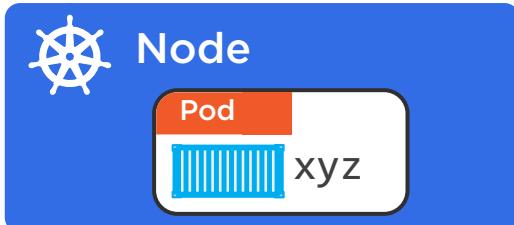
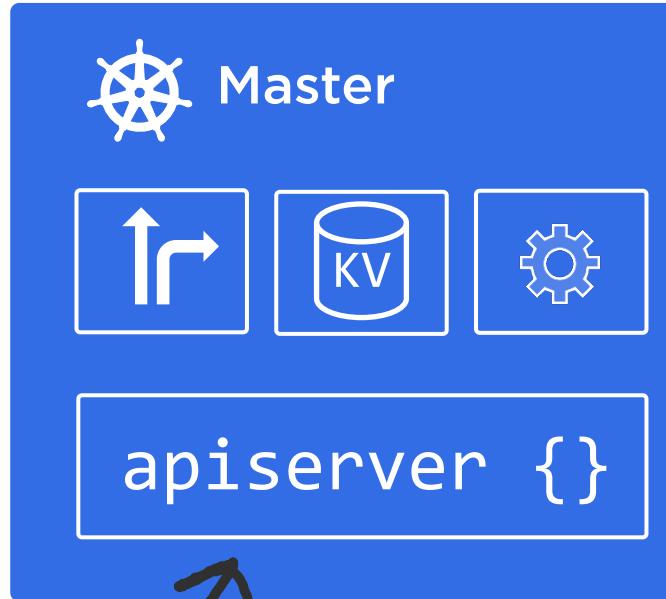
```
apiVersion: v1
kind: ReplicationController
metadata:
  name: xyz
spec:
  replicas: 4
```



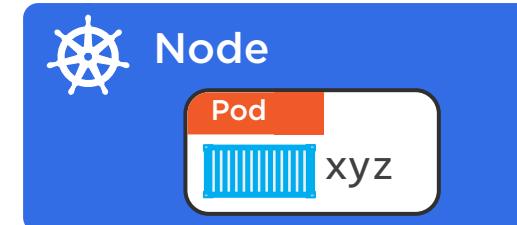
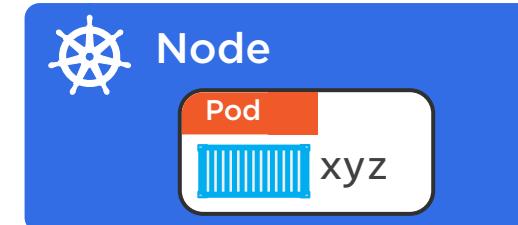
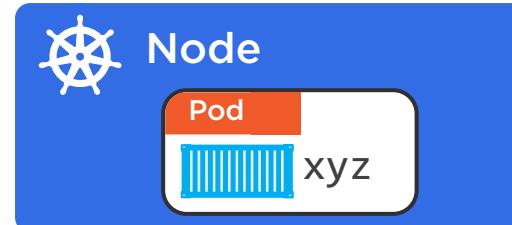
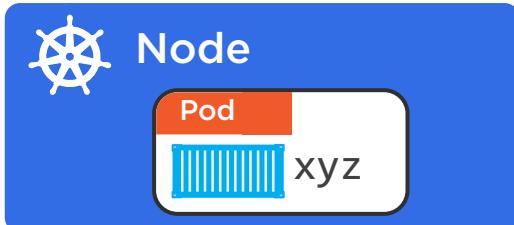
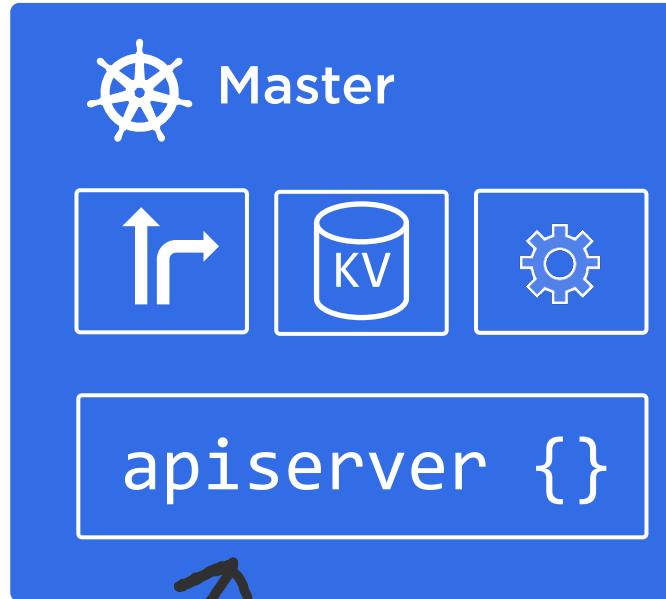
```
apiVersion: v1
kind: ReplicationController
metadata:
  name: xyz
spec:
  replicas: 4
```



```
apiVersion: v1
kind: ReplicationController
metadata:
  name: xyz
spec:
  replicas: 4
```



```
apiVersion: v1
kind: ReplicationController
metadata:
  name: xyz
spec:
  replicas: 4
```



REST objects

Self documenting

Deployed via YAML or
JSON manifests

Spec-once deploy-many

```
apiVersion: extensions/v1beta1
kind: Deployment
metadata:
  name: xyz
spec:
  replicas: 4
```

Simple rolling updates
and rollbacks

Add features to
Replication Controllers
(Replica Sets) *RCv2*

Versioned

Deployed via the
apiserver

REST objects

Self documenting

Deployed via YAML or
JSON manifests

Spec-once deploy-many

```
apiVersion: extensions/v1beta1
kind: Deployment
metadata:
  name: xyz
spec:
  replicas: 4
```

Simple rolling updates
and rollbacks

Add features to
Replication Controllers
(Replica Sets) *RCv2*

Versioned

Deployed via the
apiserver

```
apiVersion: extensions/v1beta1
kind: Deployment
metadata:
  name: xyz
spec:
  replicas: 4
```

Simple rolling updates
and rollbacks

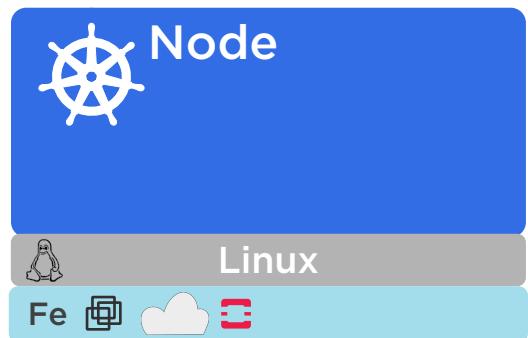
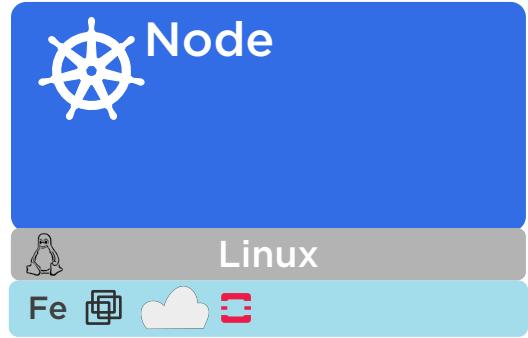
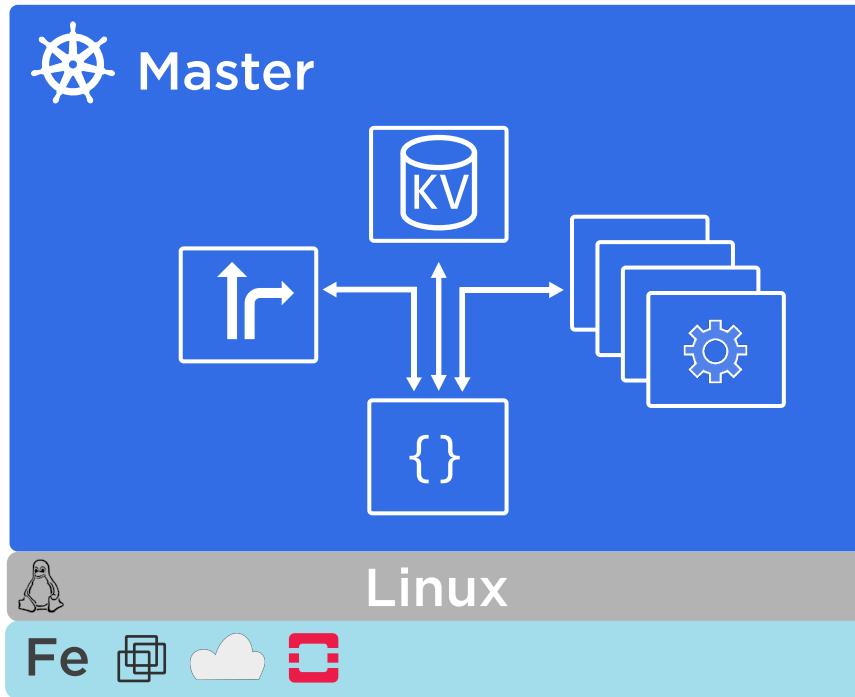
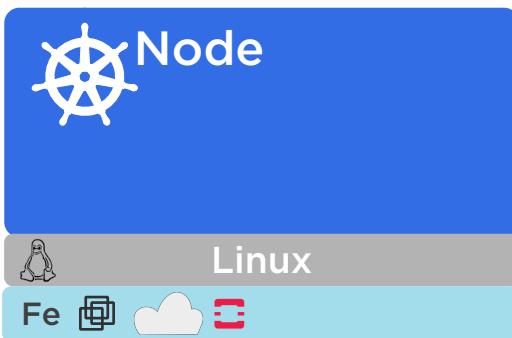
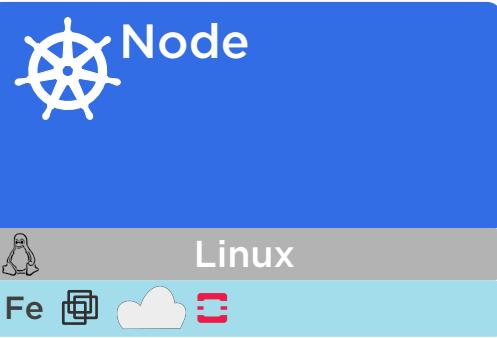
Multiple concurrent versions

- Blue-green deployments
- Canary releases

Simple versioned rollbacks



 Bringing it home!
containerized apps!





Node



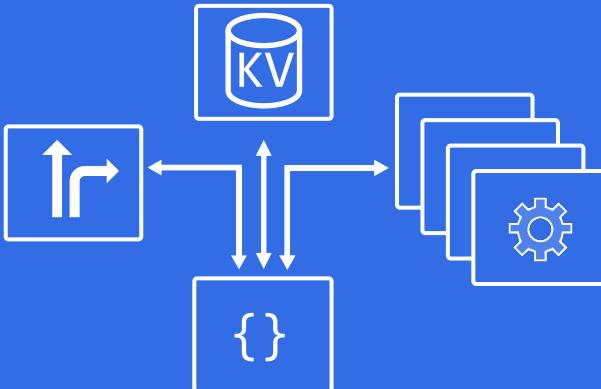
Linux



Linux



Master (control plane)



Linux



Node

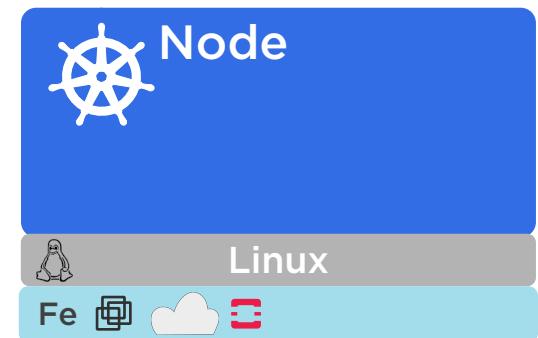
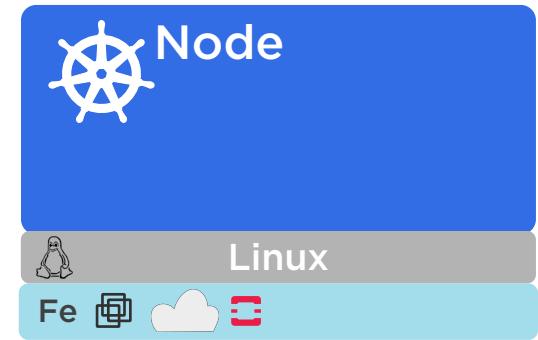
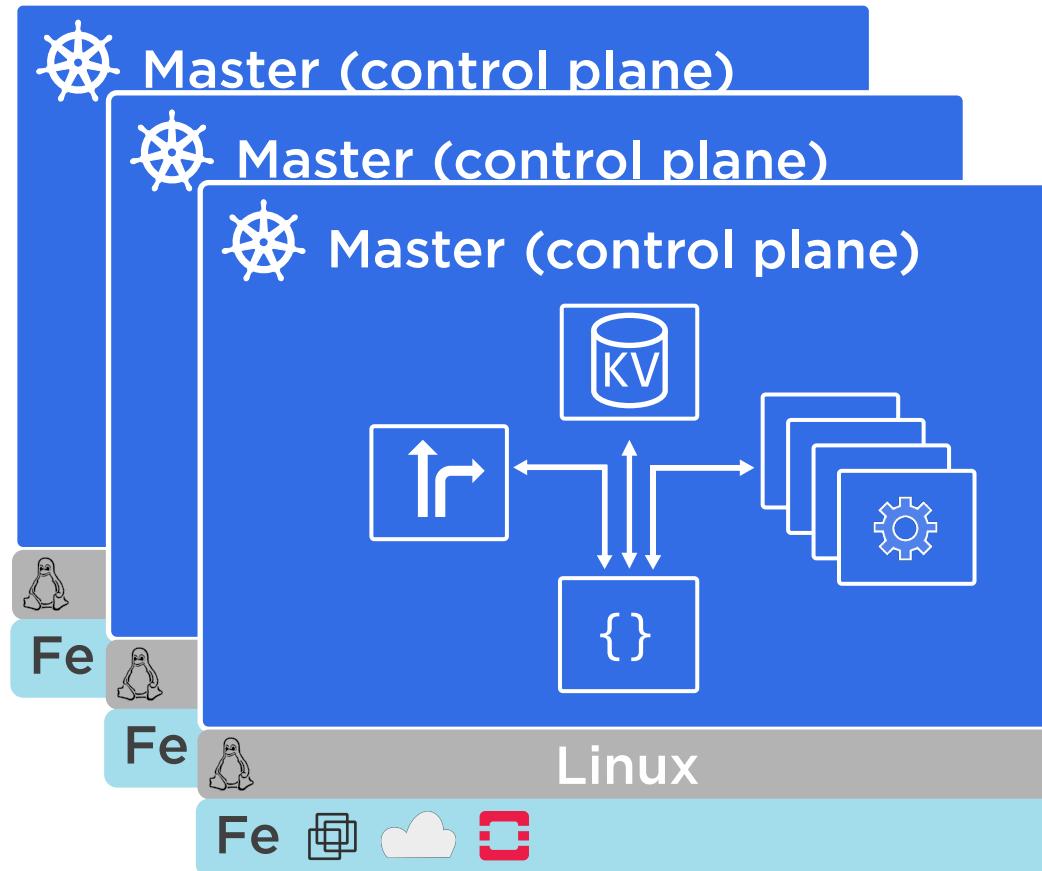
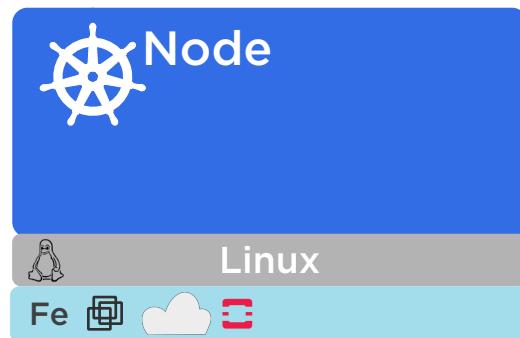
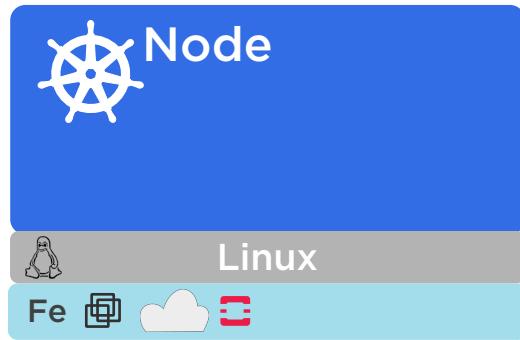


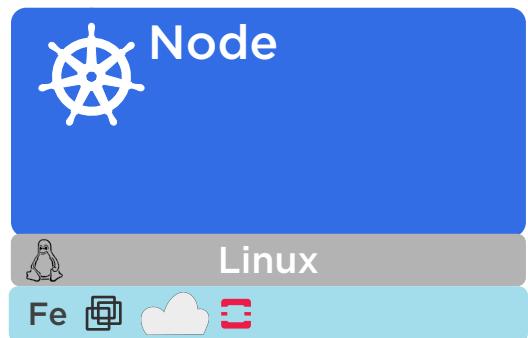
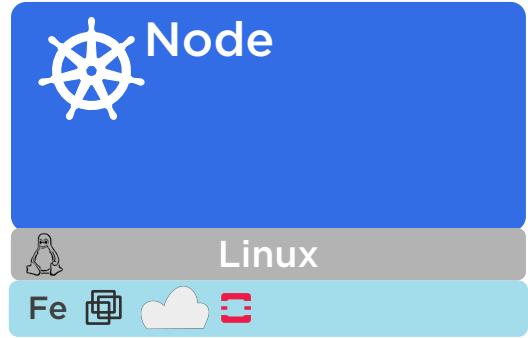
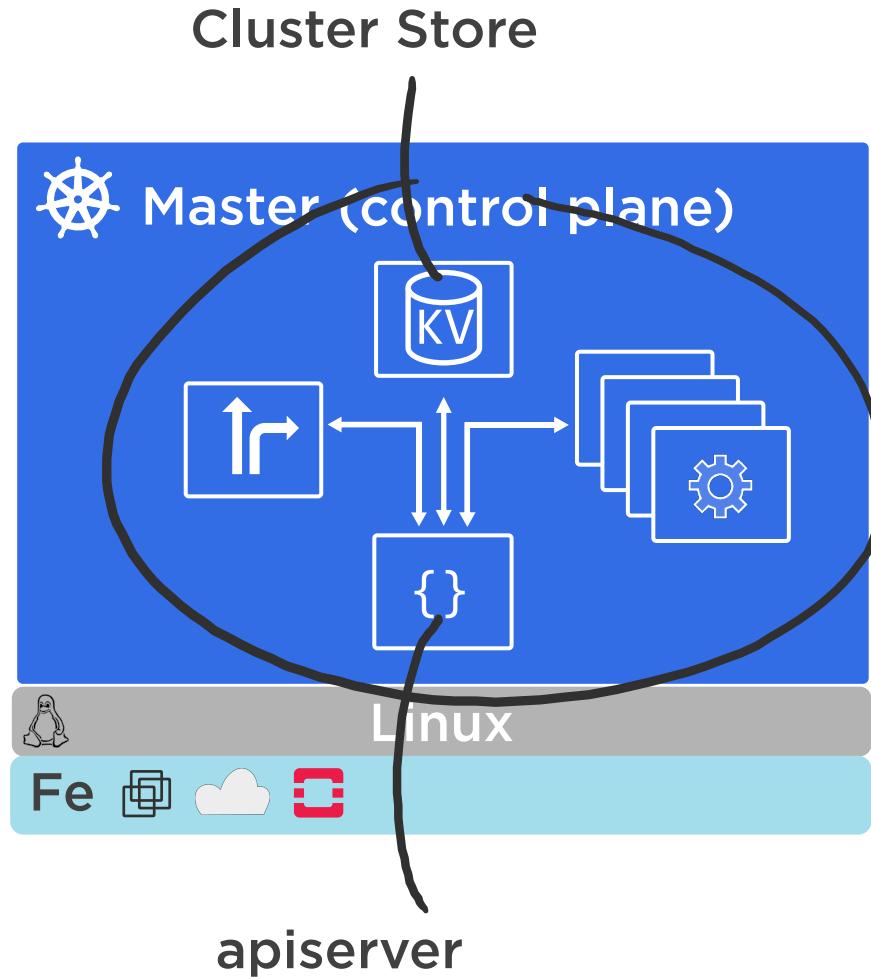
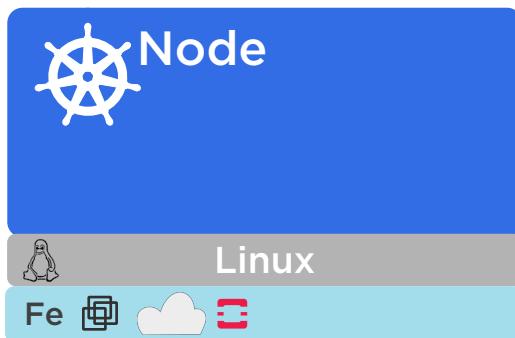
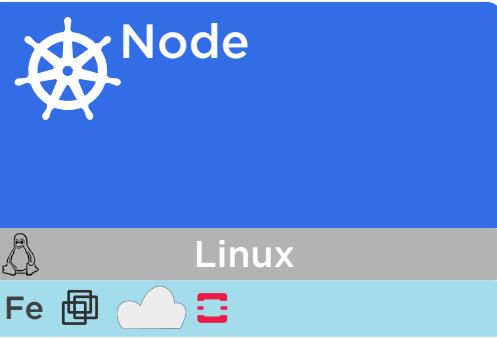
Linux

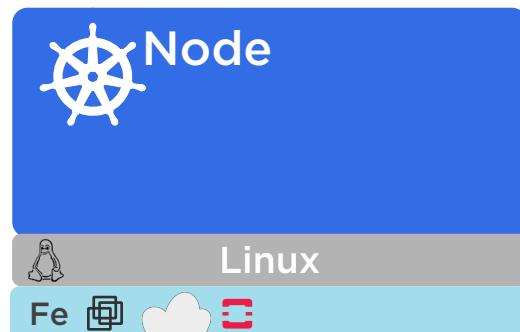
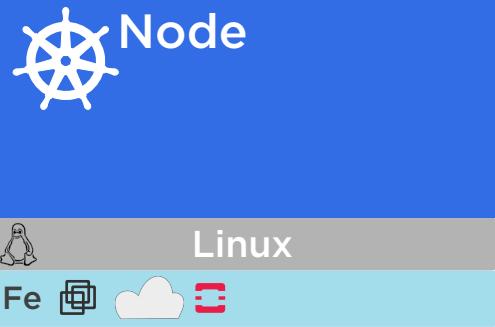


Linux

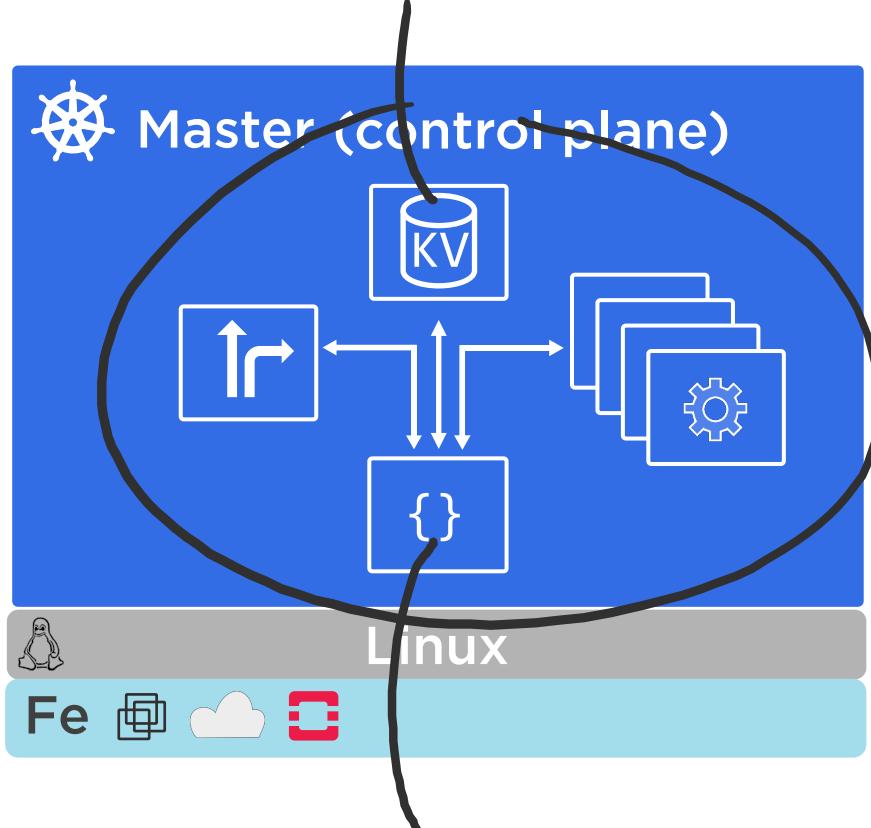




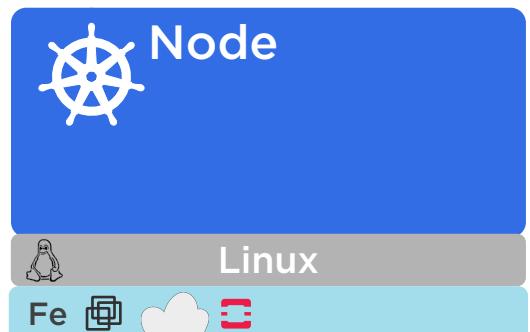
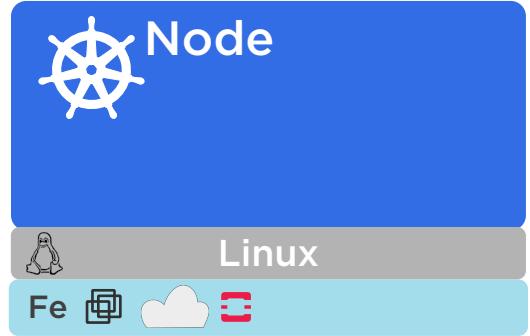


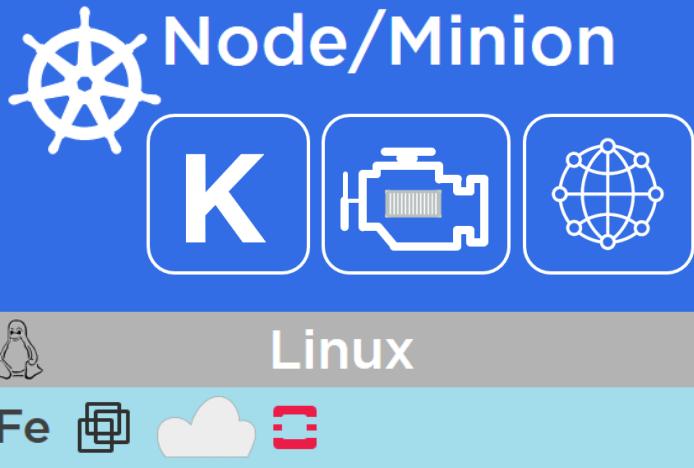


- ## Cluster Store
- Cluster state and config
 - Stateful

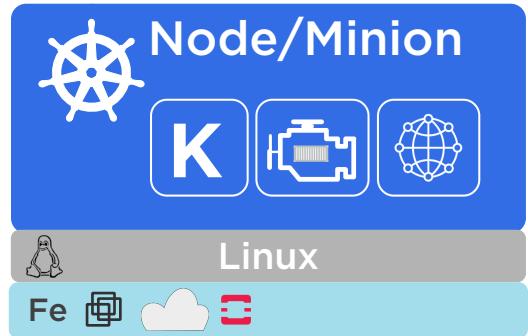
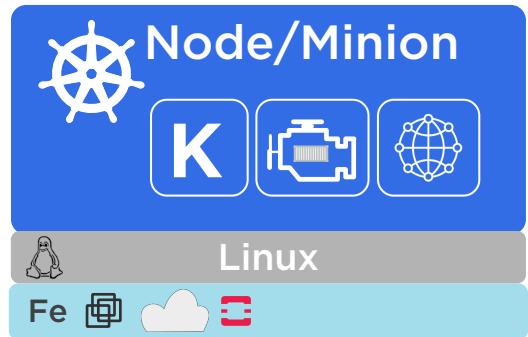
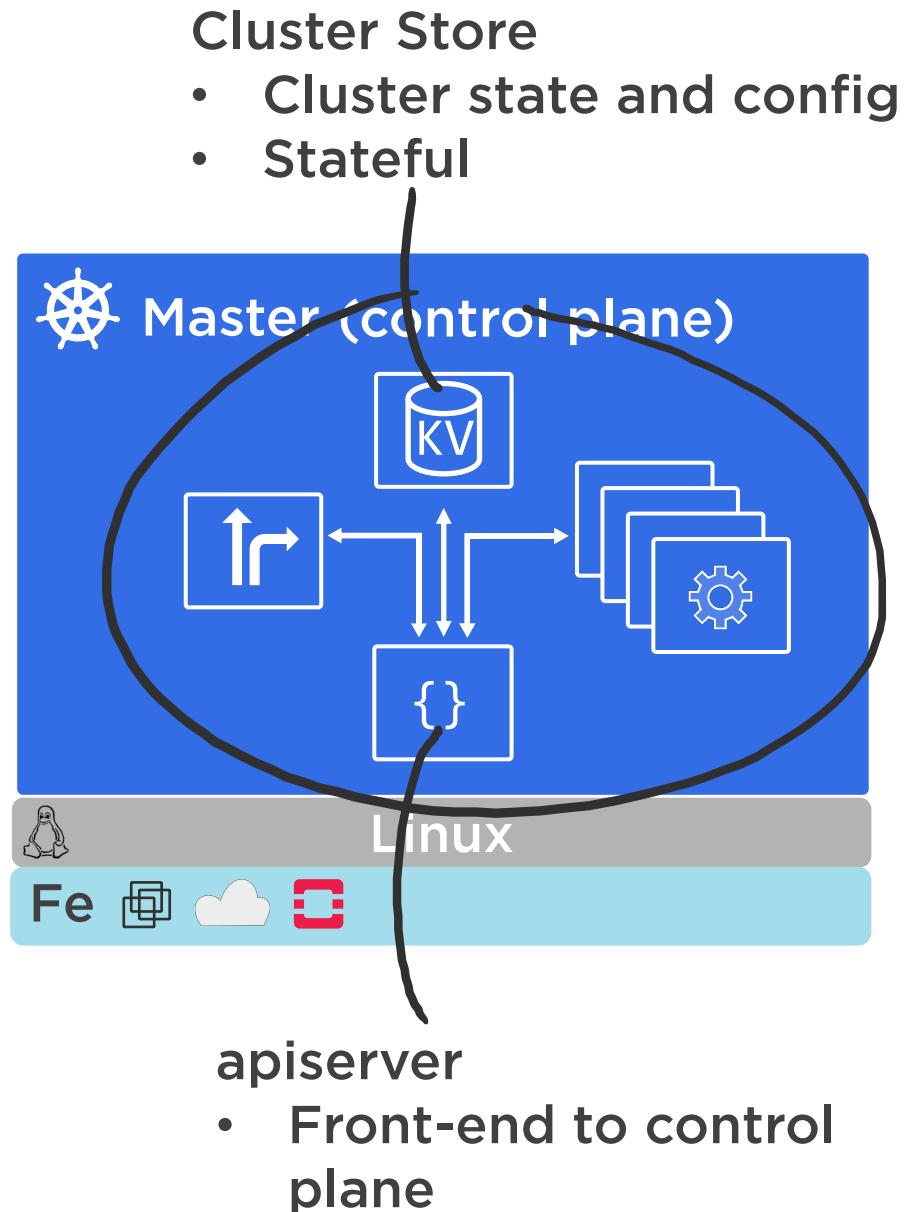


- ## apiserver
- Front-end to control plane





- Kubelet**
Main Kubernetes agent
- Container engine**
Docker or rkt
- kube-proxy**
Kubernetes networking



Objects
in the
K8s API

Pods : Atomic unit of scheduling...

**Replication
Controllers** : Scale pods, desired state etc...

Deployments : RC + rolling updates, rollbacks...

Services : Stable networking...



Coming up next...

Installing Kubernetes