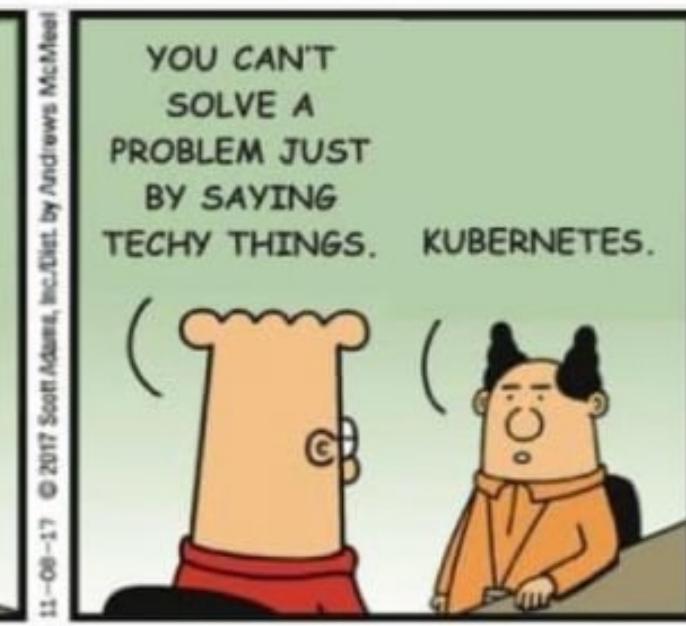
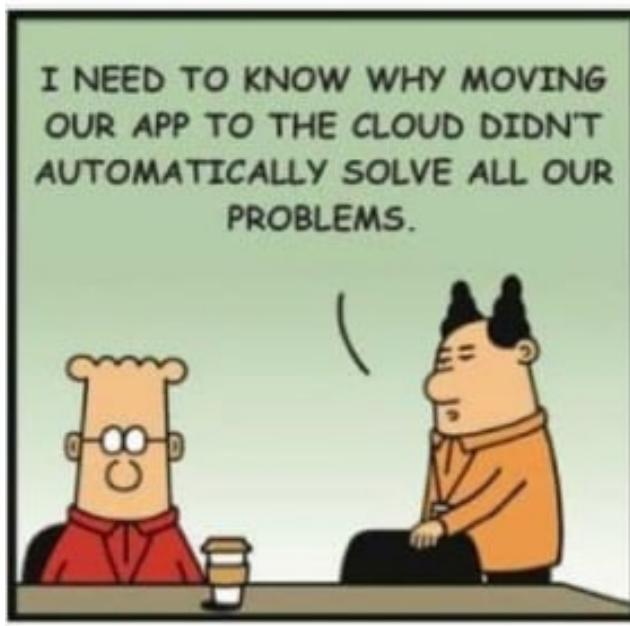


# Kubernetes Internals & Security

Hyunsang Choi

# Kubernetes?

Solved all your problems. You're welcome.

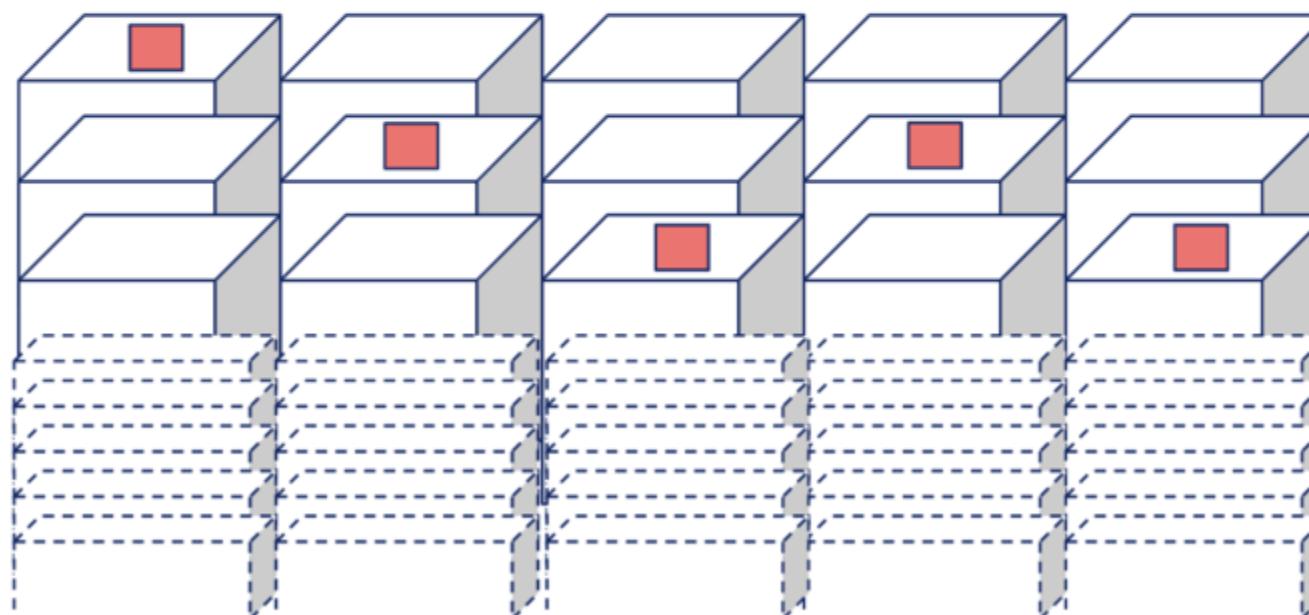


## Containerized App

Kubernetes, run 5 replicas of this  
with 2 CPUs and 3GB RAM each!

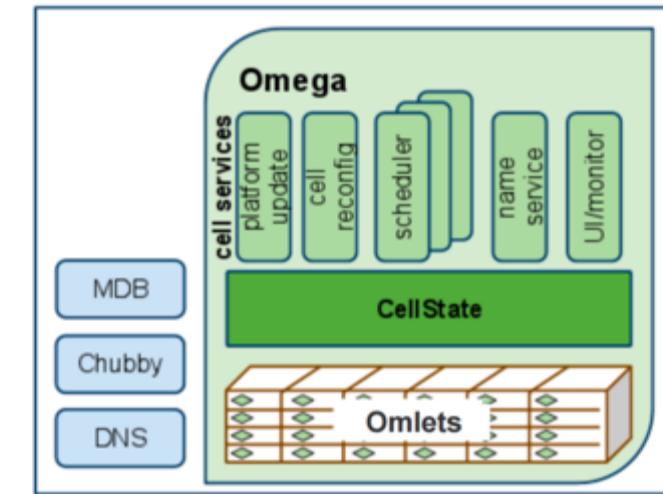
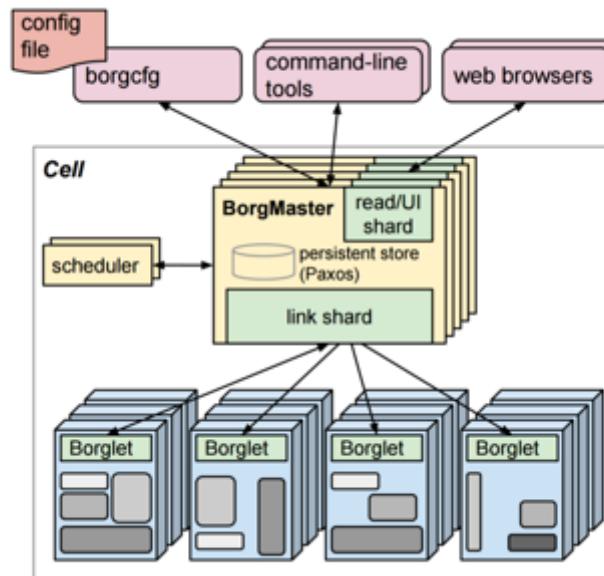
### kubernetes

abstracts away a pool of machines as a single target



# Borg & Omega

- Google's cluster orchestration projects
  - 10/15 years R&D (secret)
  - Borg: Large scale cluster management (written in C++)
  - Omega: Improve performance

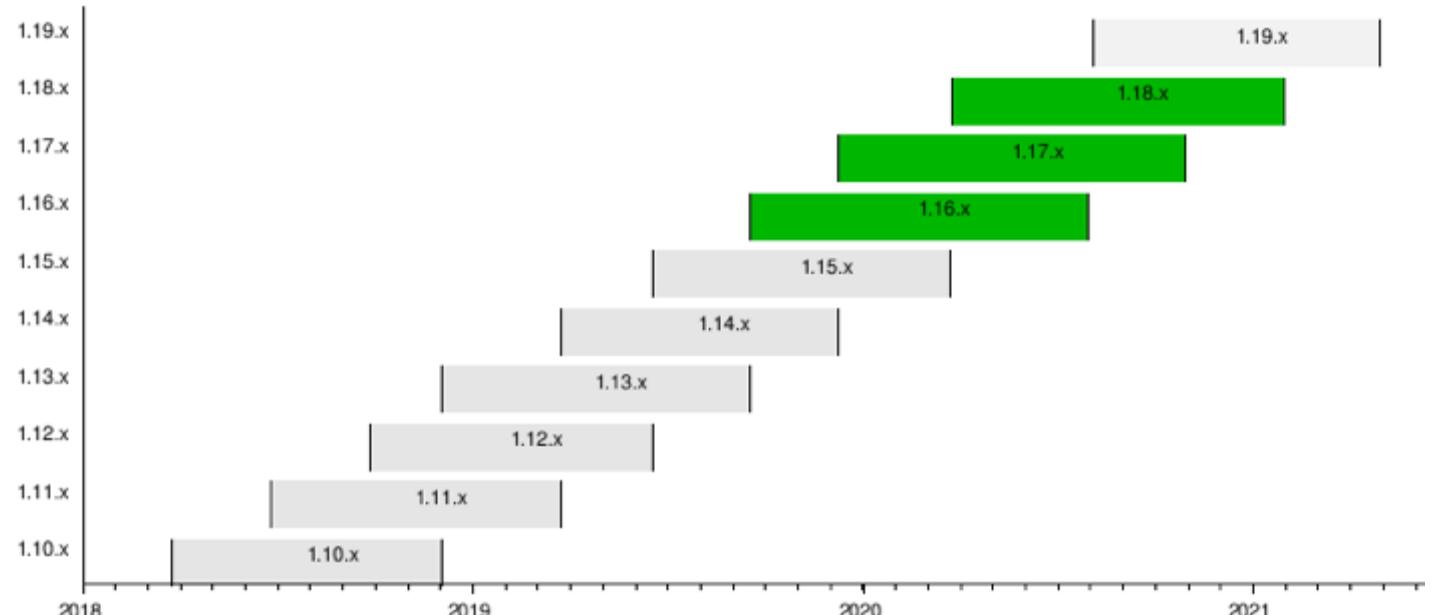


# Borg vs Kubernetes

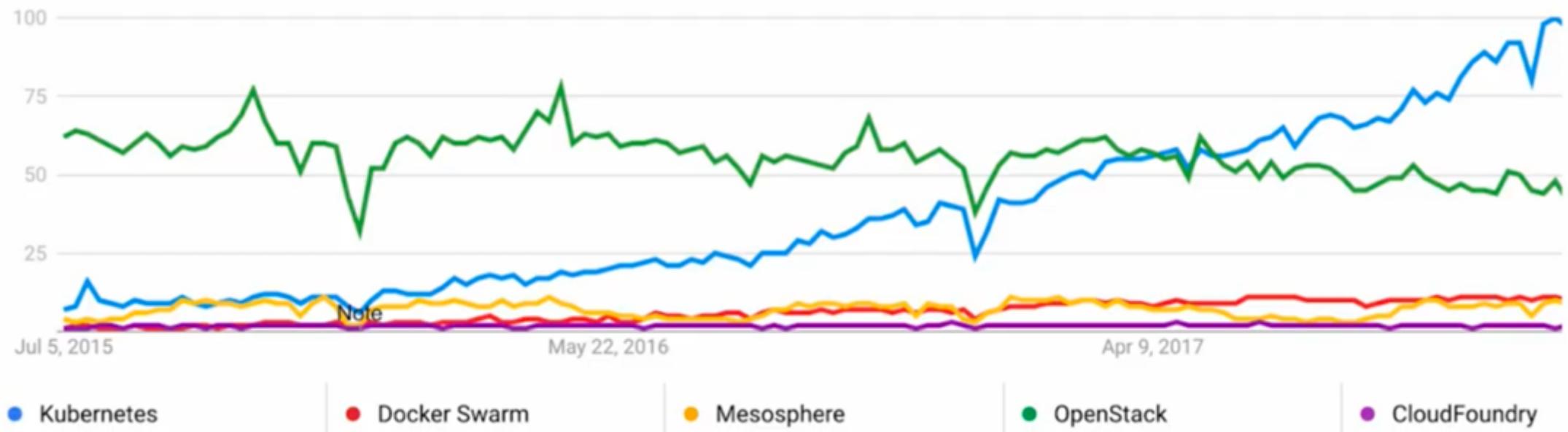
- Borg groups work by ‘job’; Kubernetes adds ‘labels’ for greater flexibility
- Borg uses an IP-per-machine design; Kubernetes uses a network-per-machine and IP-per-Pod design to allow late-binding of ports (letting developers choose ports, not the infrastructure)
- Kubernetes APIs are simpler than Borg's API
- Borg uses [LMCTFY](#) as its container technology; Kubernetes allows the use of Docker, rkt, etc
- Borg is not open source

# Kubernetes

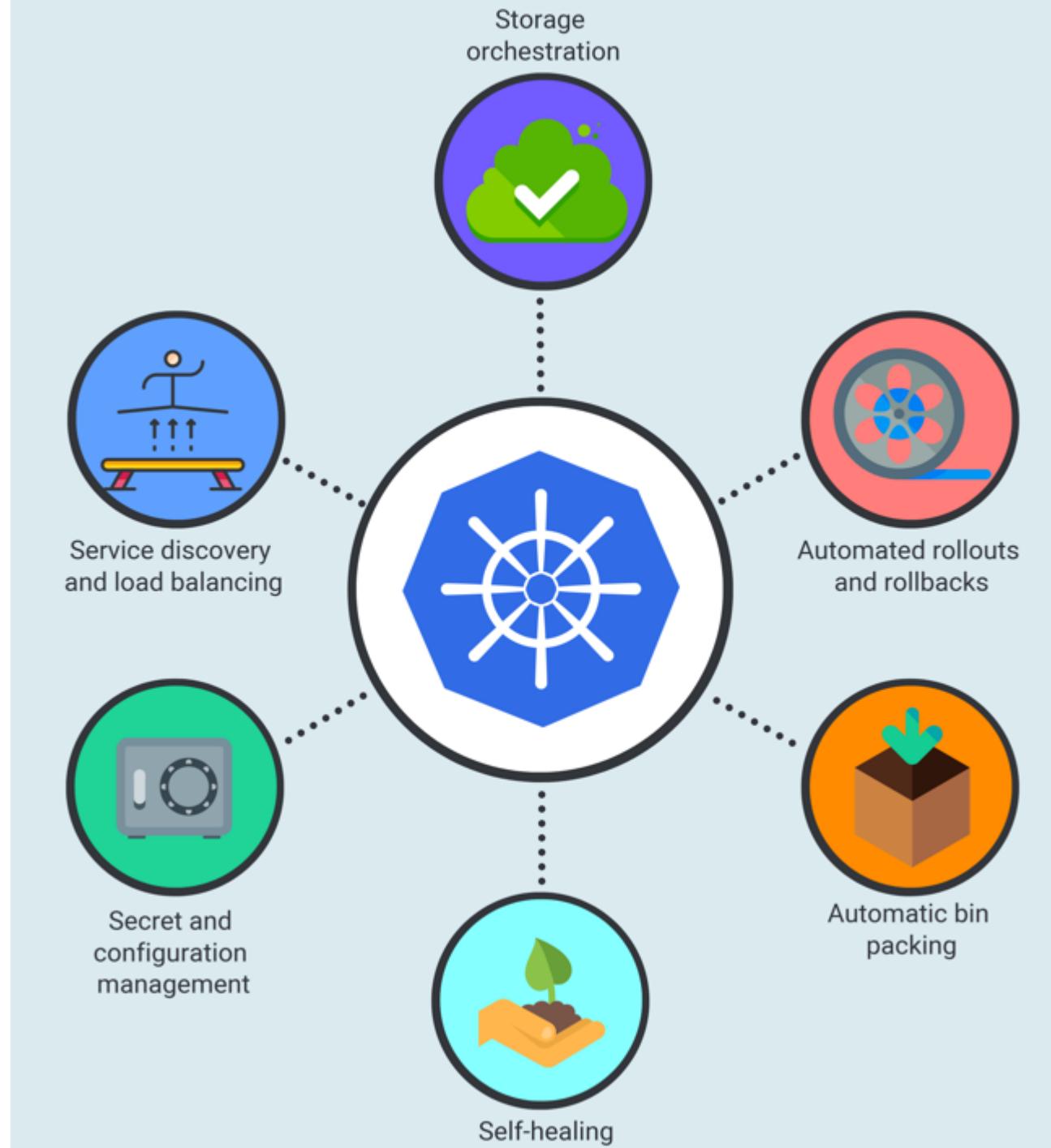
- Open-source container-orchestration system for automating application deployment, scaling, and management
- Greek (κυβερνήτης): helmsman
- Heavily influenced by Borg project
- V1.0: 2015.07.21
- Implemented in Go
- Apache License 2.0
- Maintained by CNCF
- Stylized as k8s



# De-Facto Standard

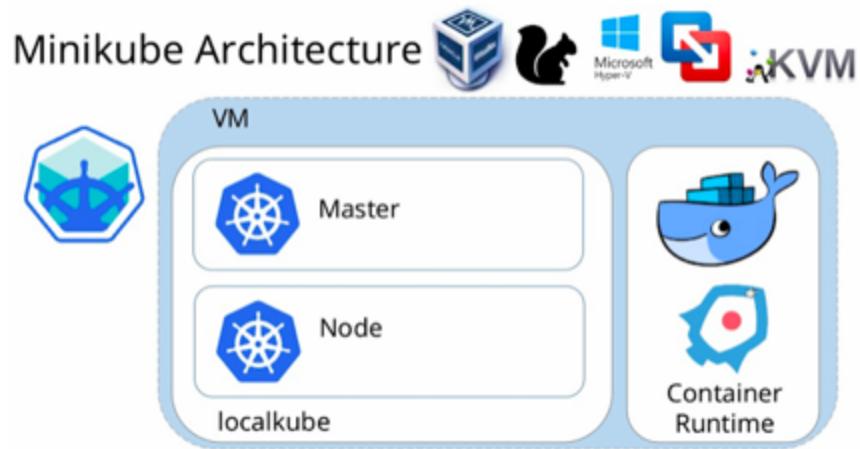


# k8s Features

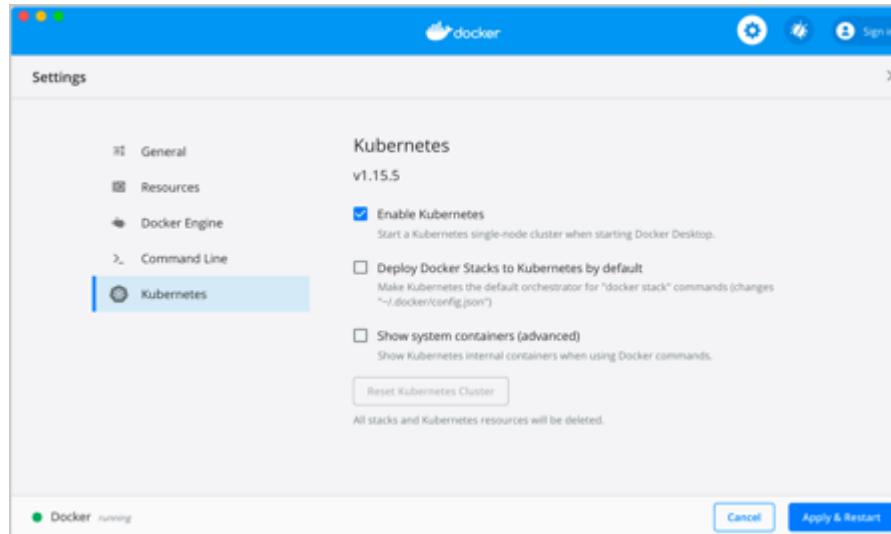


# Kubernetes Practice

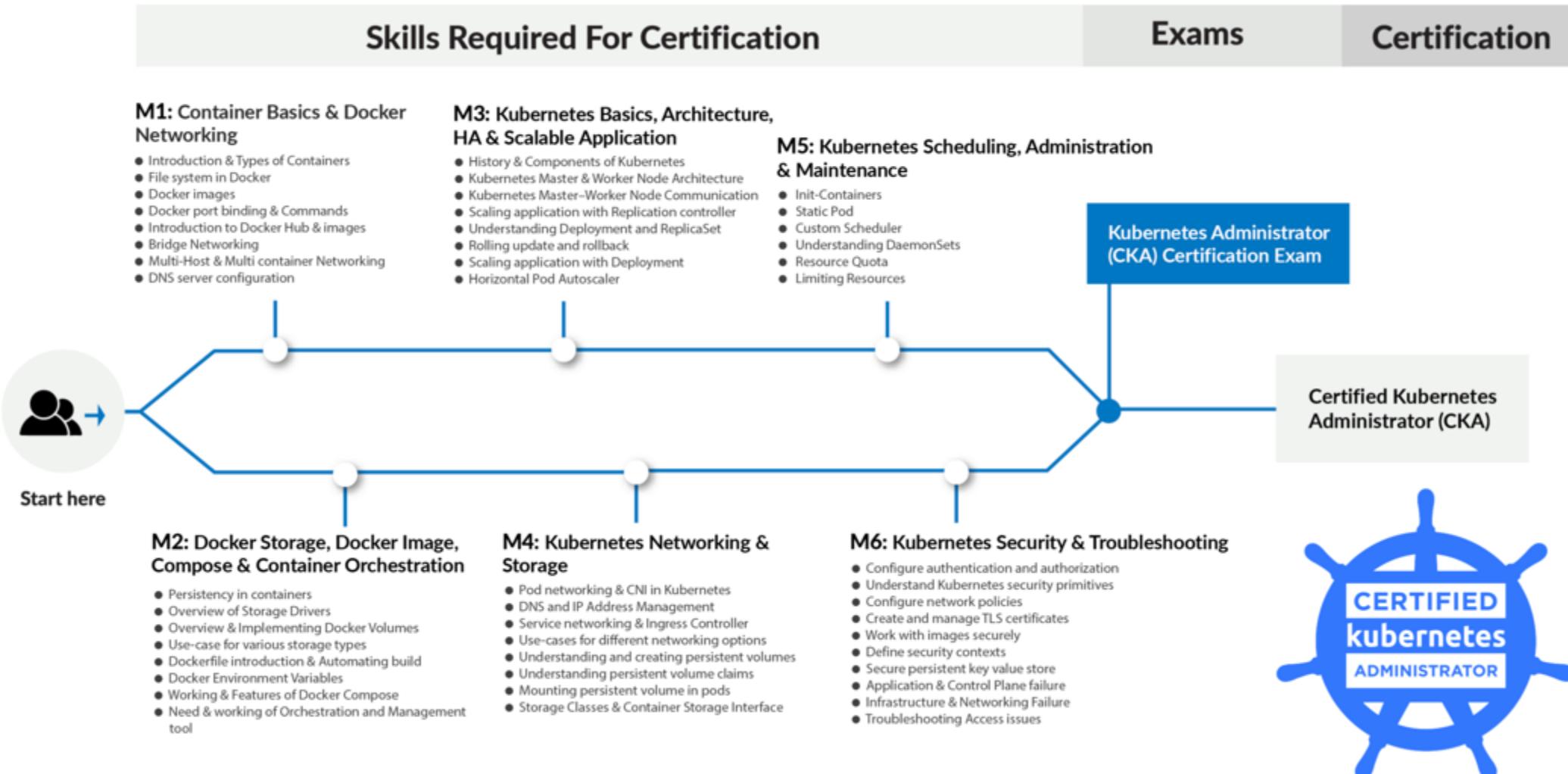
- Minikube



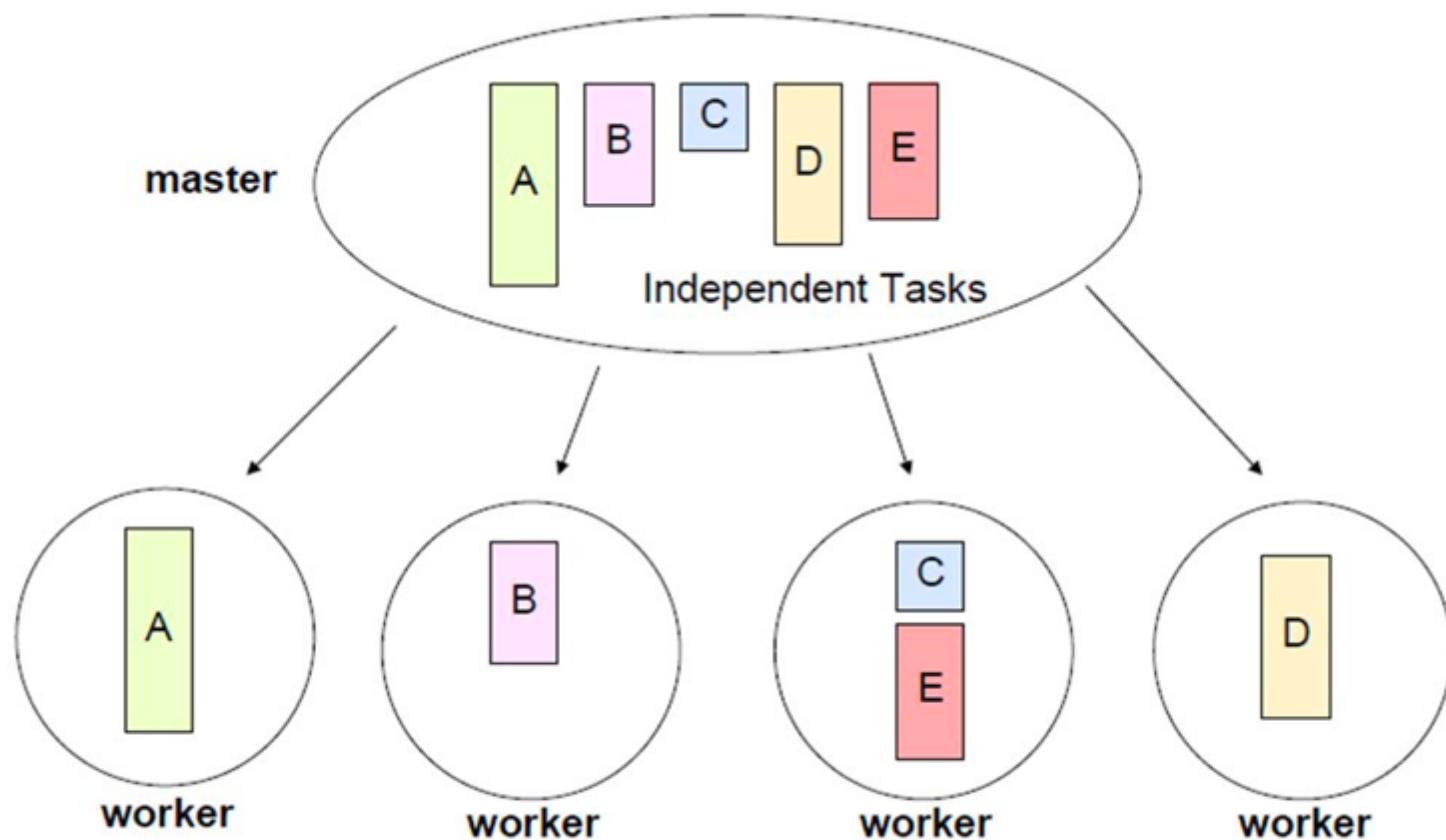
- Katacoda
- Docker desktop
- Cloud
- ...



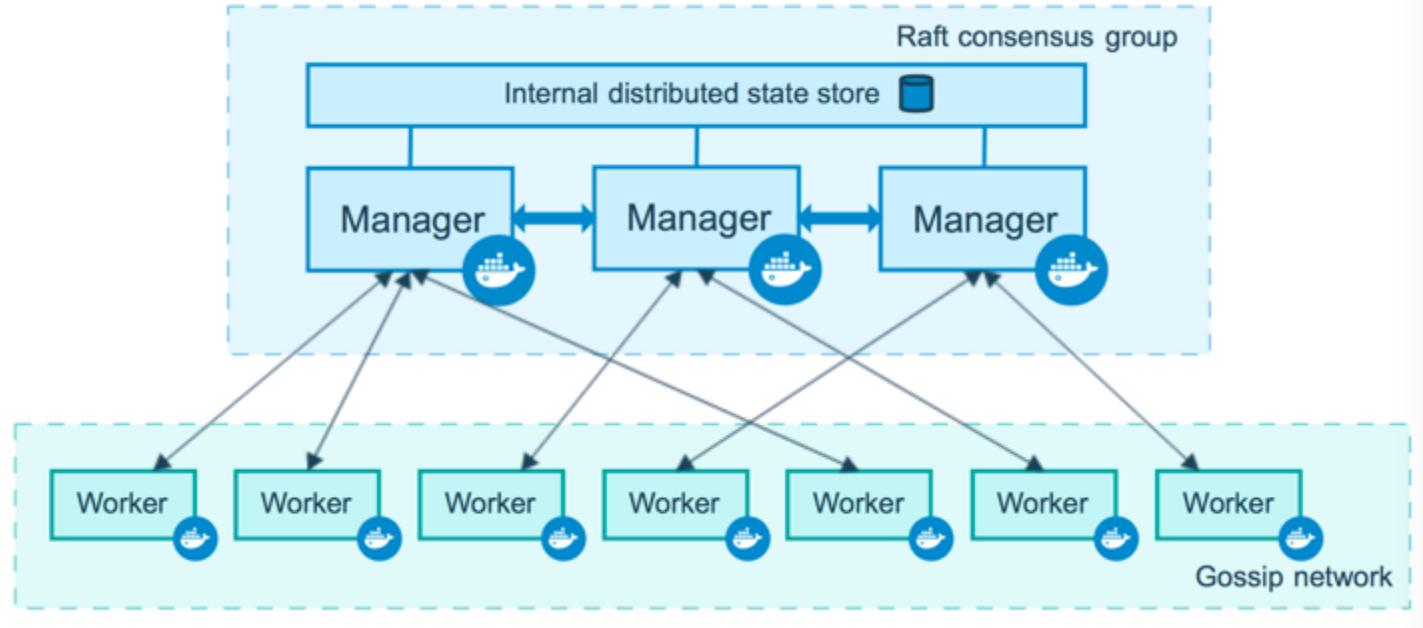
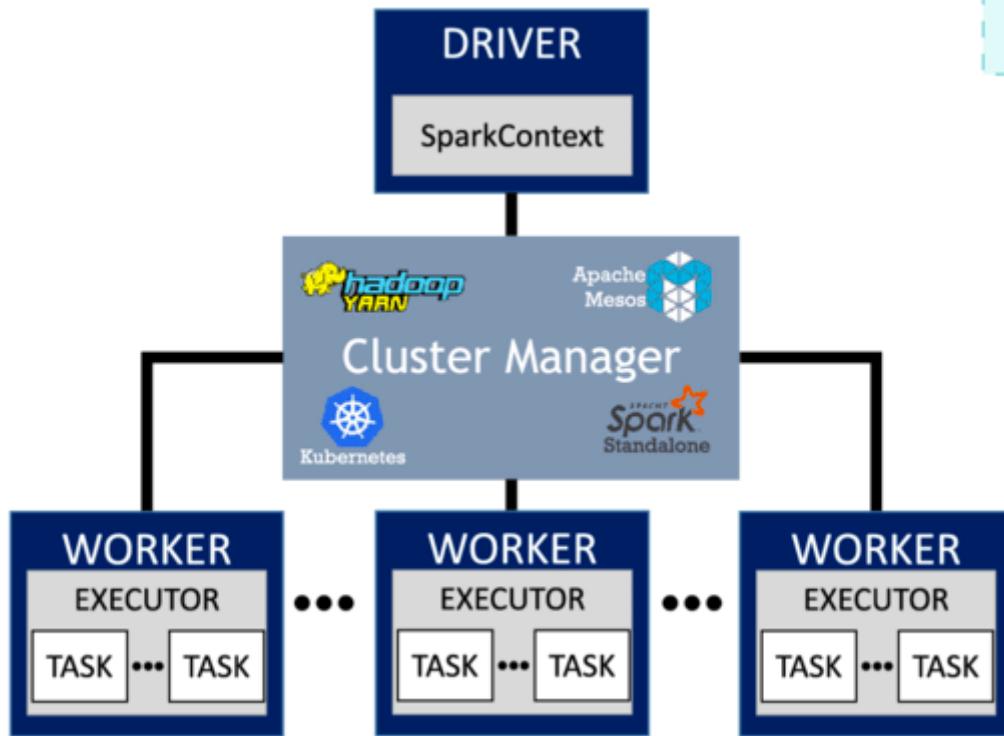
# CKA (Certified Kubernetes Administrator)



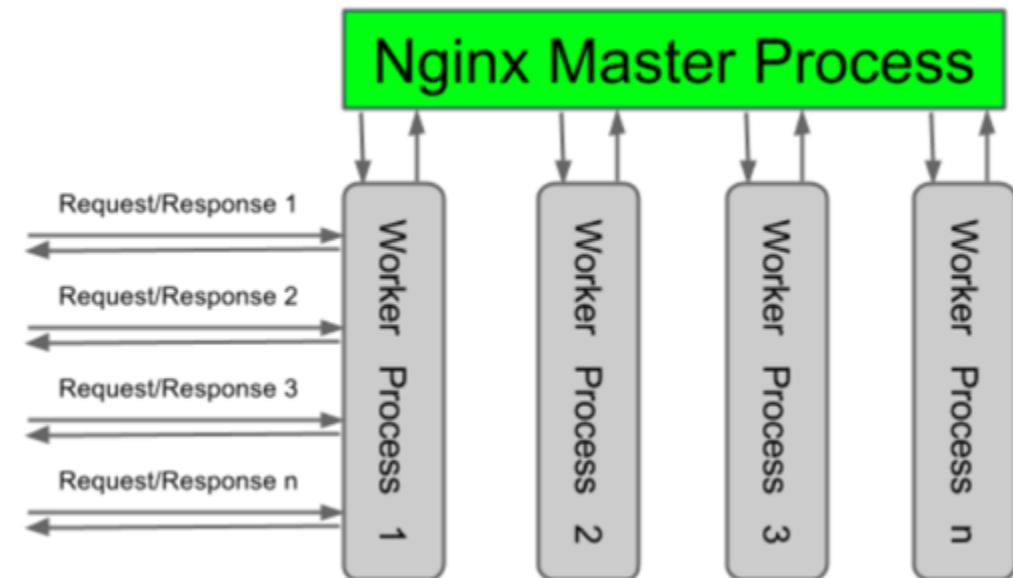
# Master-Worker Pattern



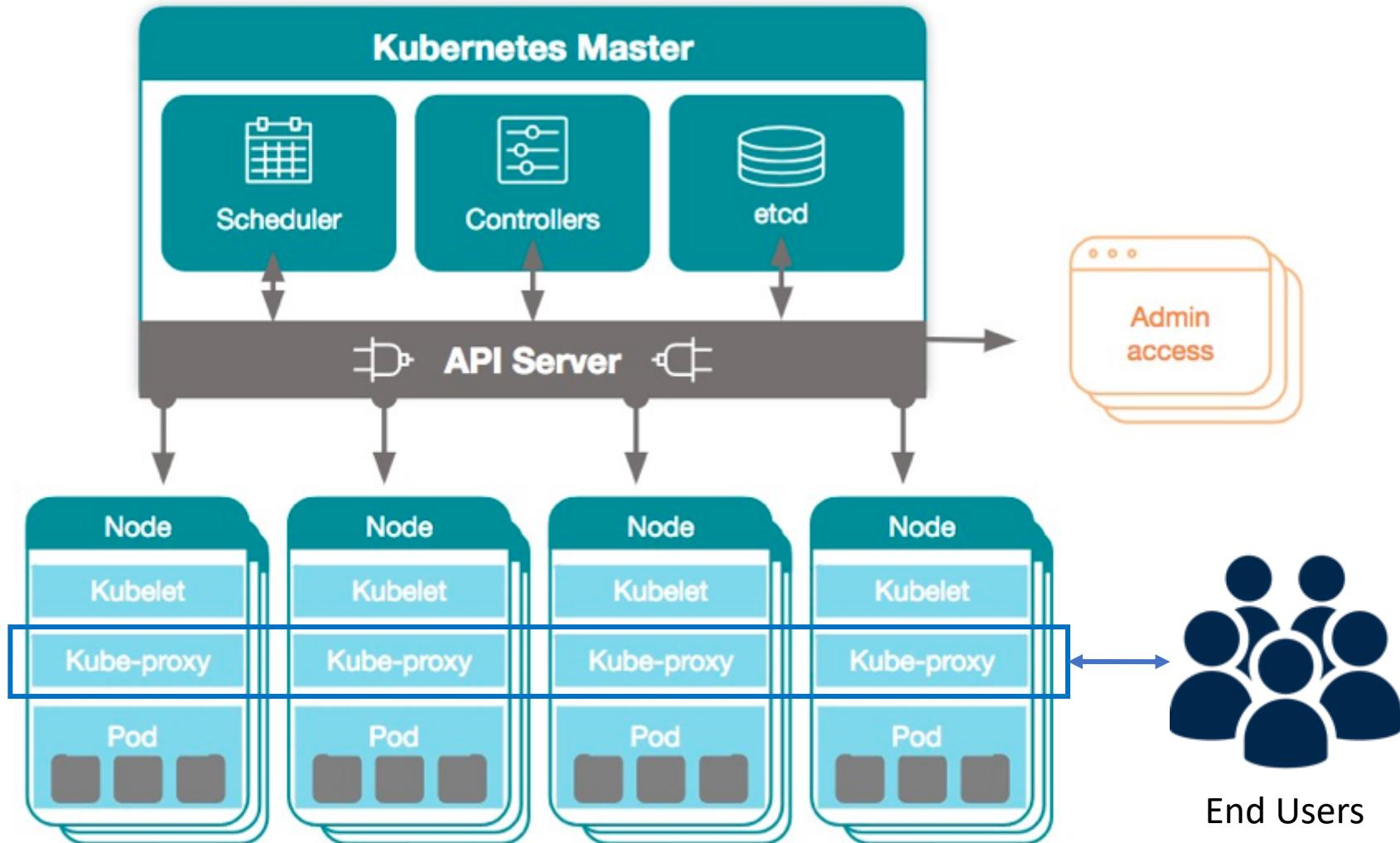
# Examples



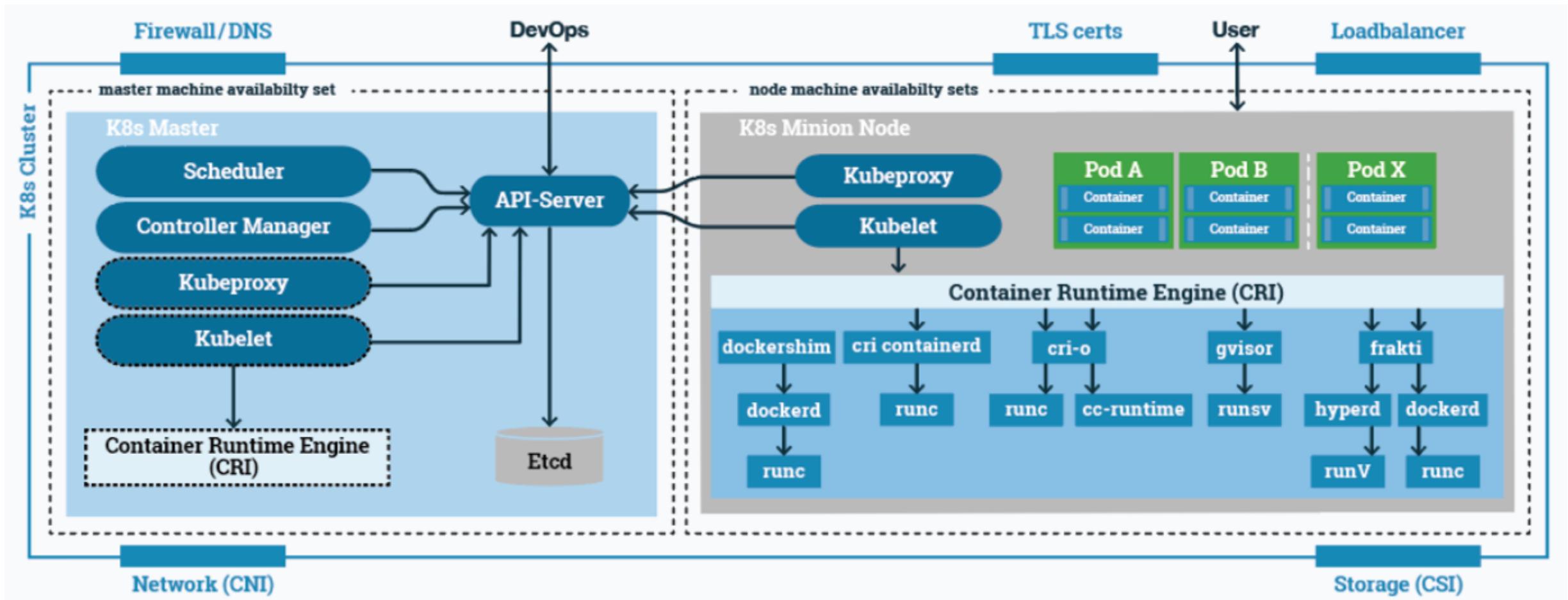
Single master process with "n" number of worker process



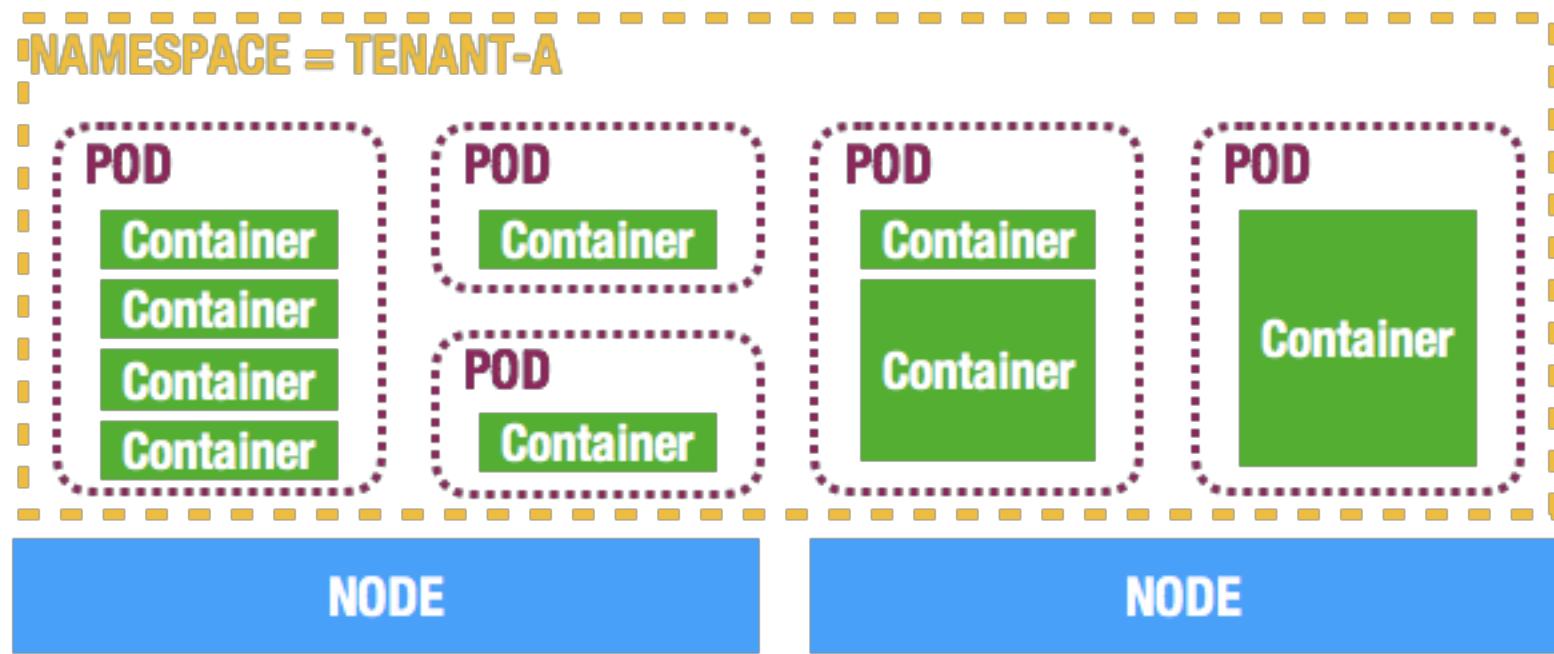
# Kubernetes Architecture



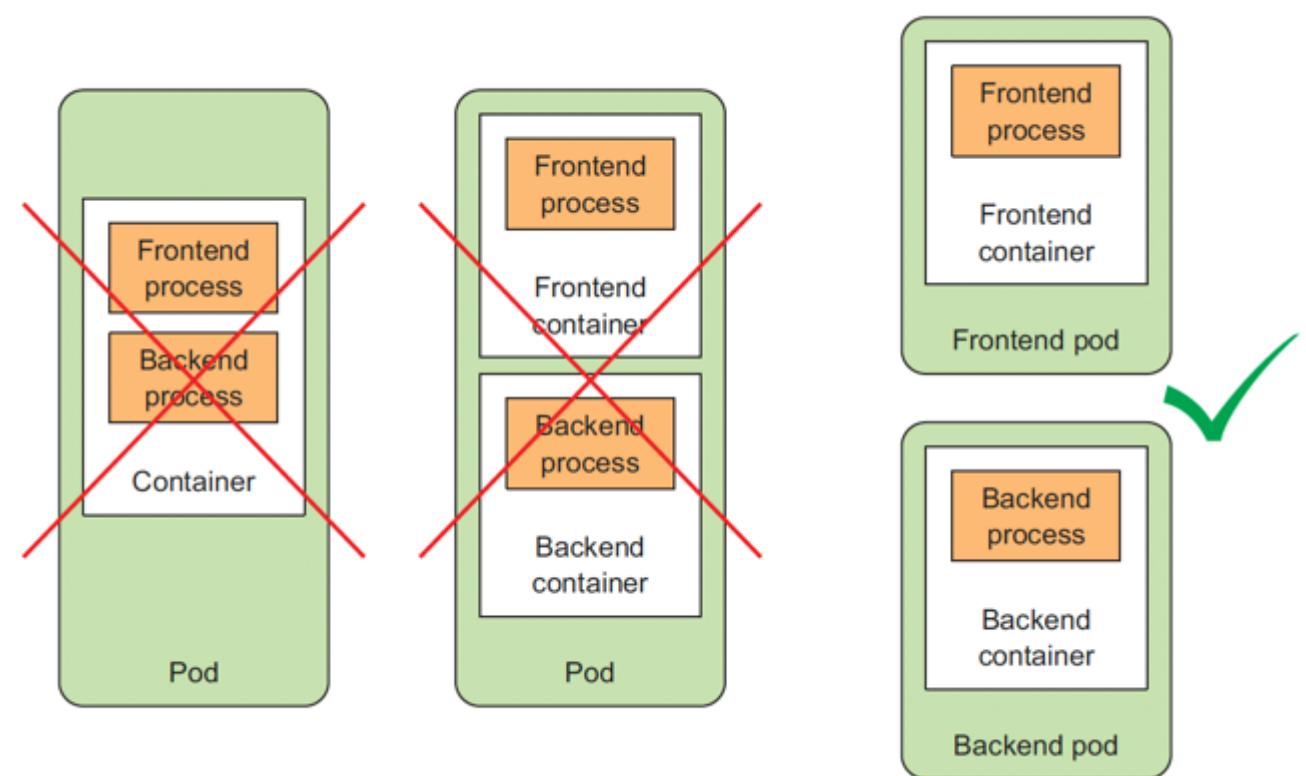
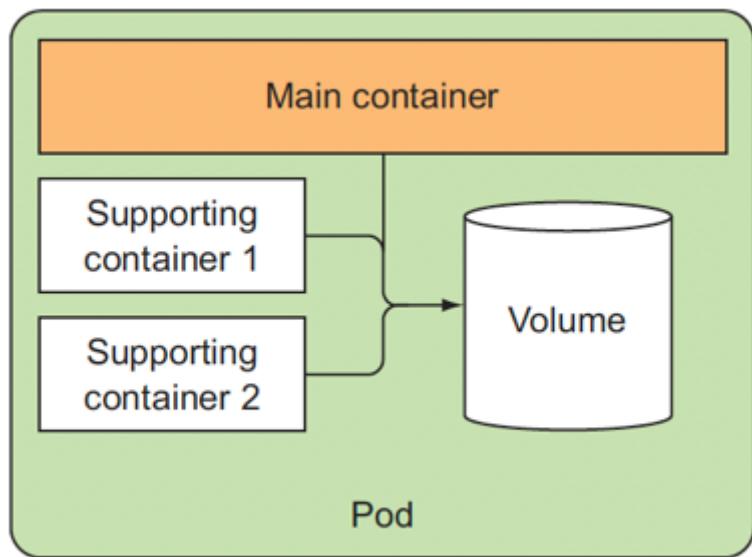
# In-Depth Architecture



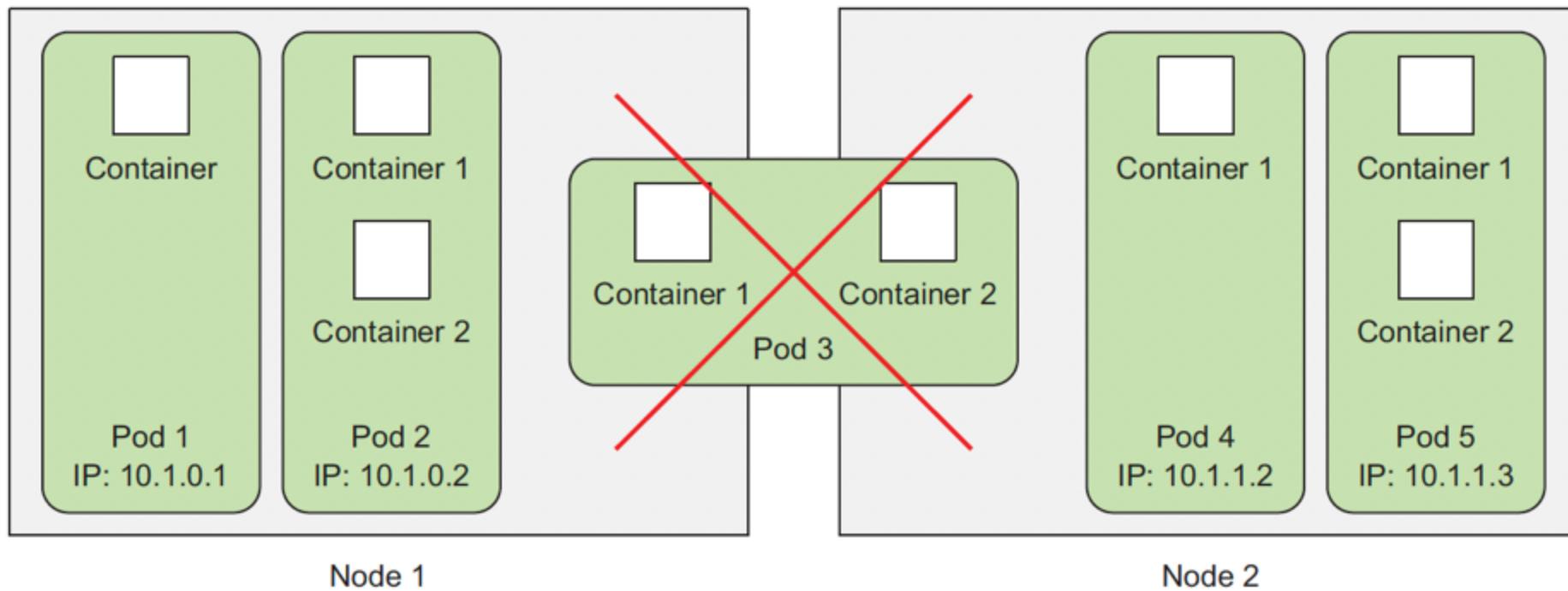
# Pod and Container



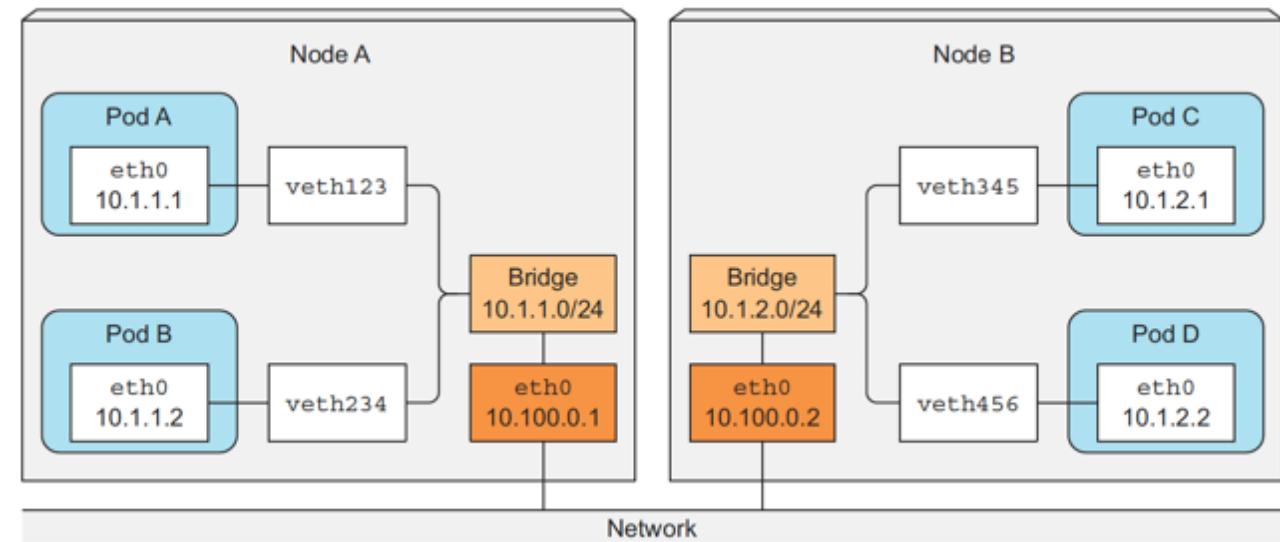
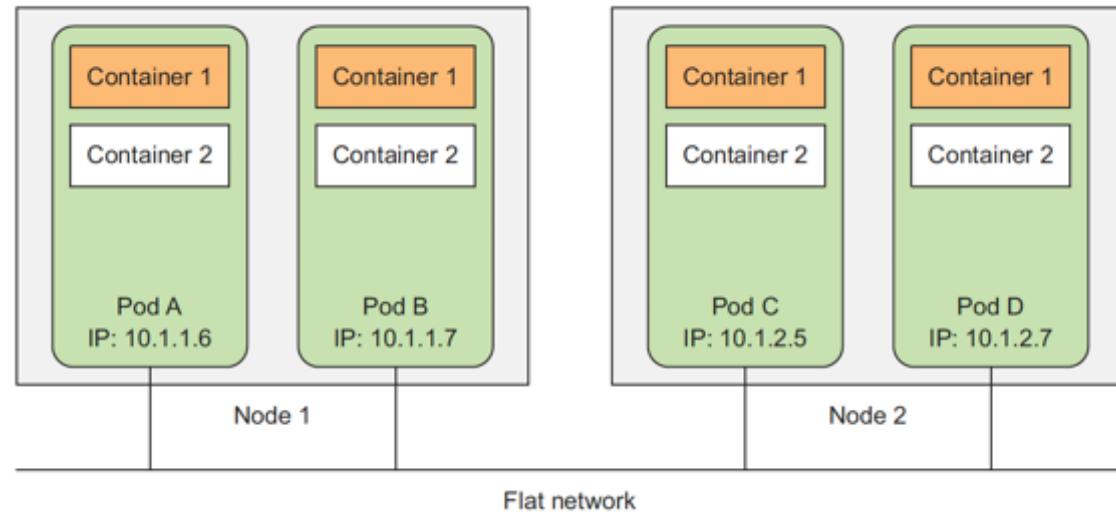
# Pod Design



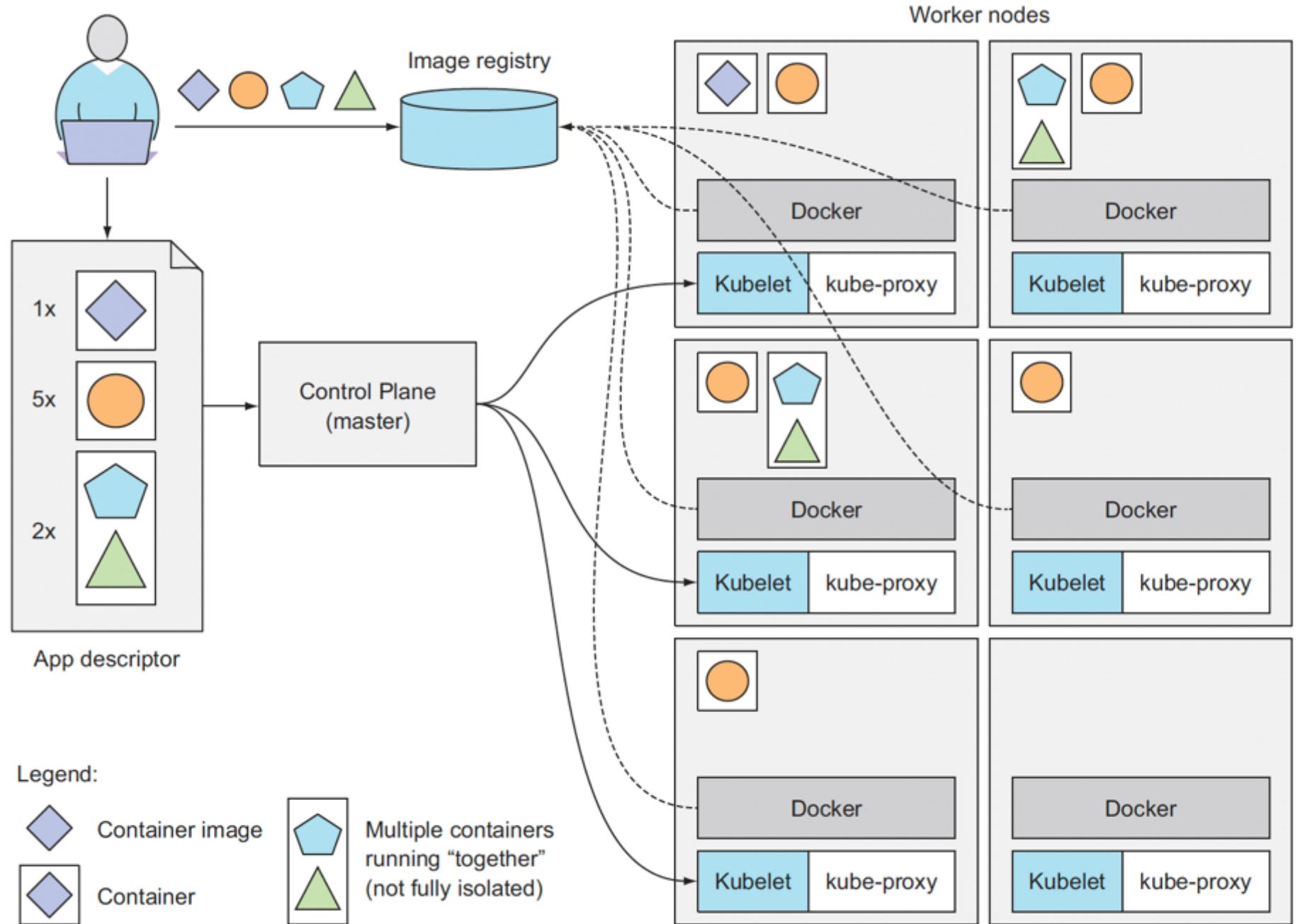
# Pod in a Node



# Pod IP Address

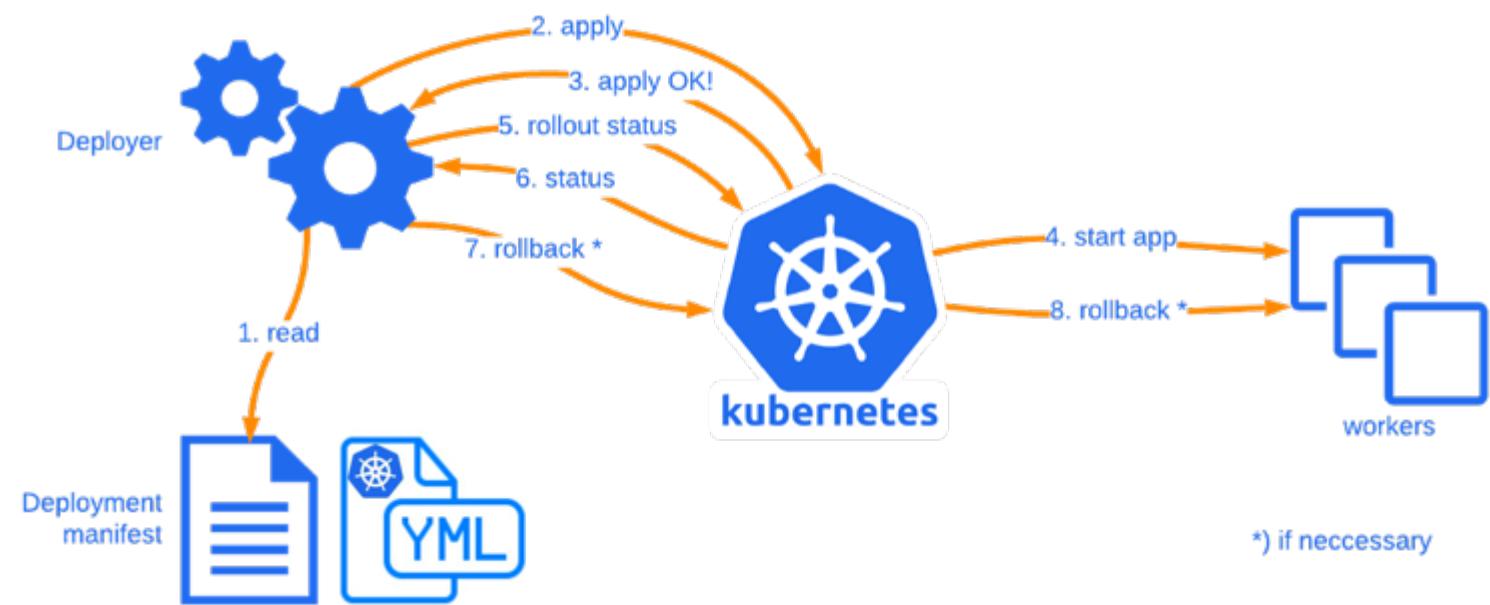
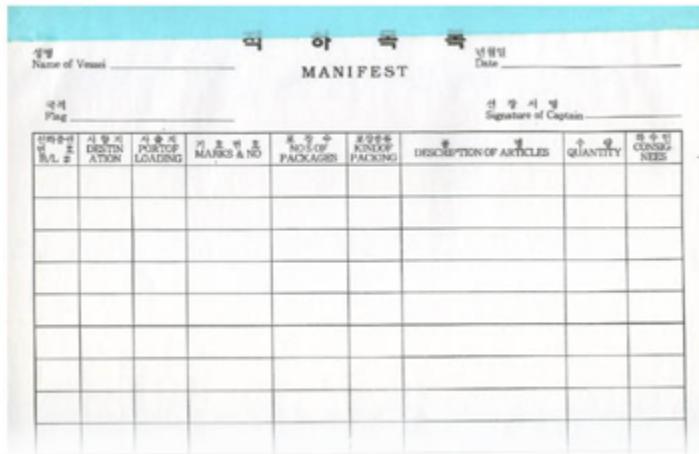


# Overview



# Manifest

- App descriptor = Manifest



# Manifest: YAML File

```
apiVersion: v1
kind: Pod
metadata:
  labels:
    phase: prod
    role: frontend
    name: myfirstpod
  name: myfirstpod
spec:
  containers:
    - name: filepuller
      image: sliranc/filepuller:latest
      volumeMounts:
        - mountPath: /usr/share/nginx/html
          name: static-vol
    - name: webserver
      image: nginx:latest
      ports:
        - containerPort: 80
      volumeMounts:
        - mountPath: /usr/share/nginx/html
          name: static-vol
  volumes:
    - name: static-vol
      emptyDir: {}
```

kind: System object \ resource  
Examples: Pod, RC, Service etc...

Spec: is the specification of the desired state of an **object**.

Containers



# 8 Ways to Create a Pod

Resource	Result	RBAC resource
Pod	Creates Pod	Pods
ReplicationController	Creates Pod with replicas and prevent it from crashing	replicationcontrollers
ReplicaSet		replicasets
Deployment	Creates Pod and replicaset	deployments
DaemonSet	Creates Pod in each node in the cluster	daemonsets
Job	Creates Pod to completion	jobs
CronJob	Creates Pod to completion periodically	cronjobs
StatefulSet	Creates Pod with order	statefulsets

# Label

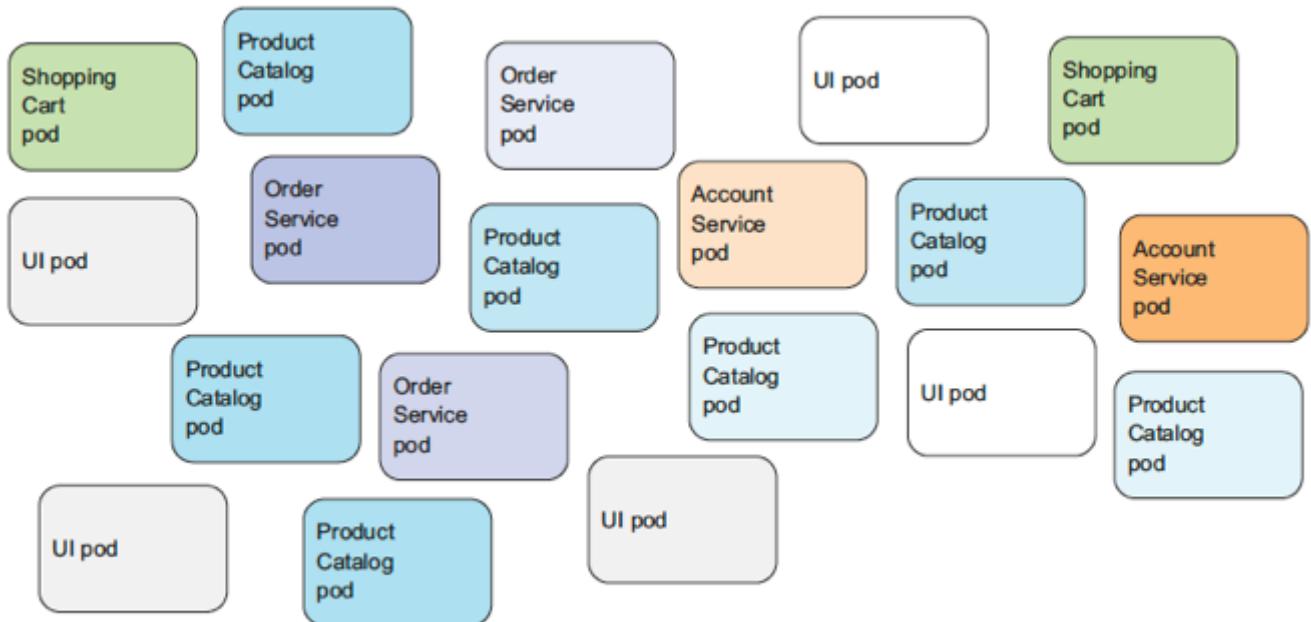
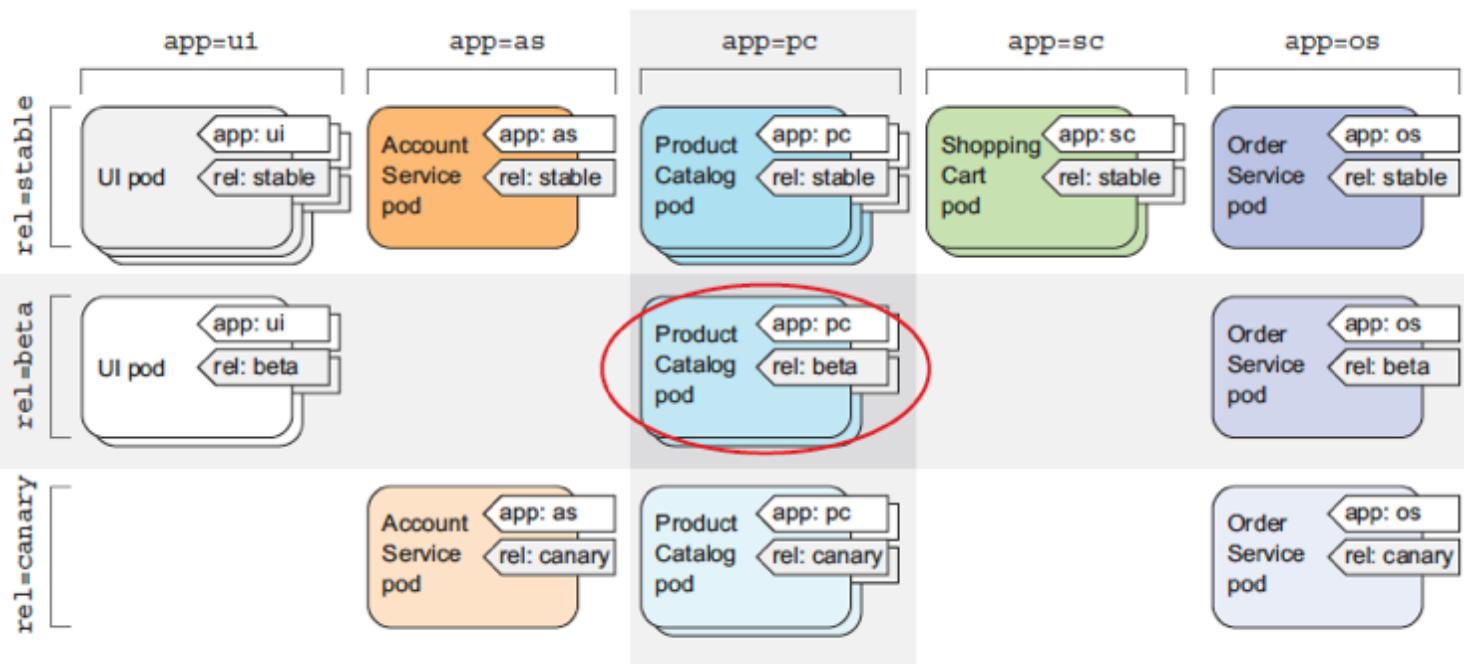


Figure 3.6 Uncategorized pods in a microservices architecture



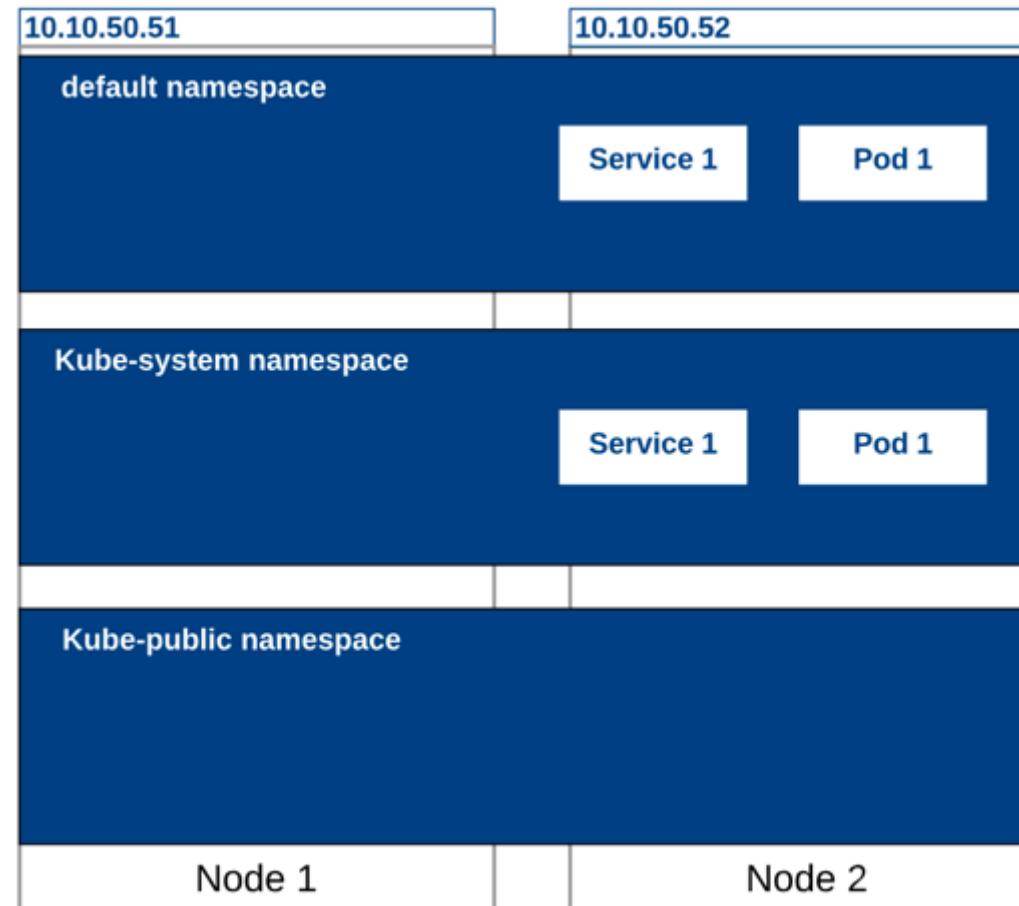
# Label-Selector

Labels can be matched by label selectors

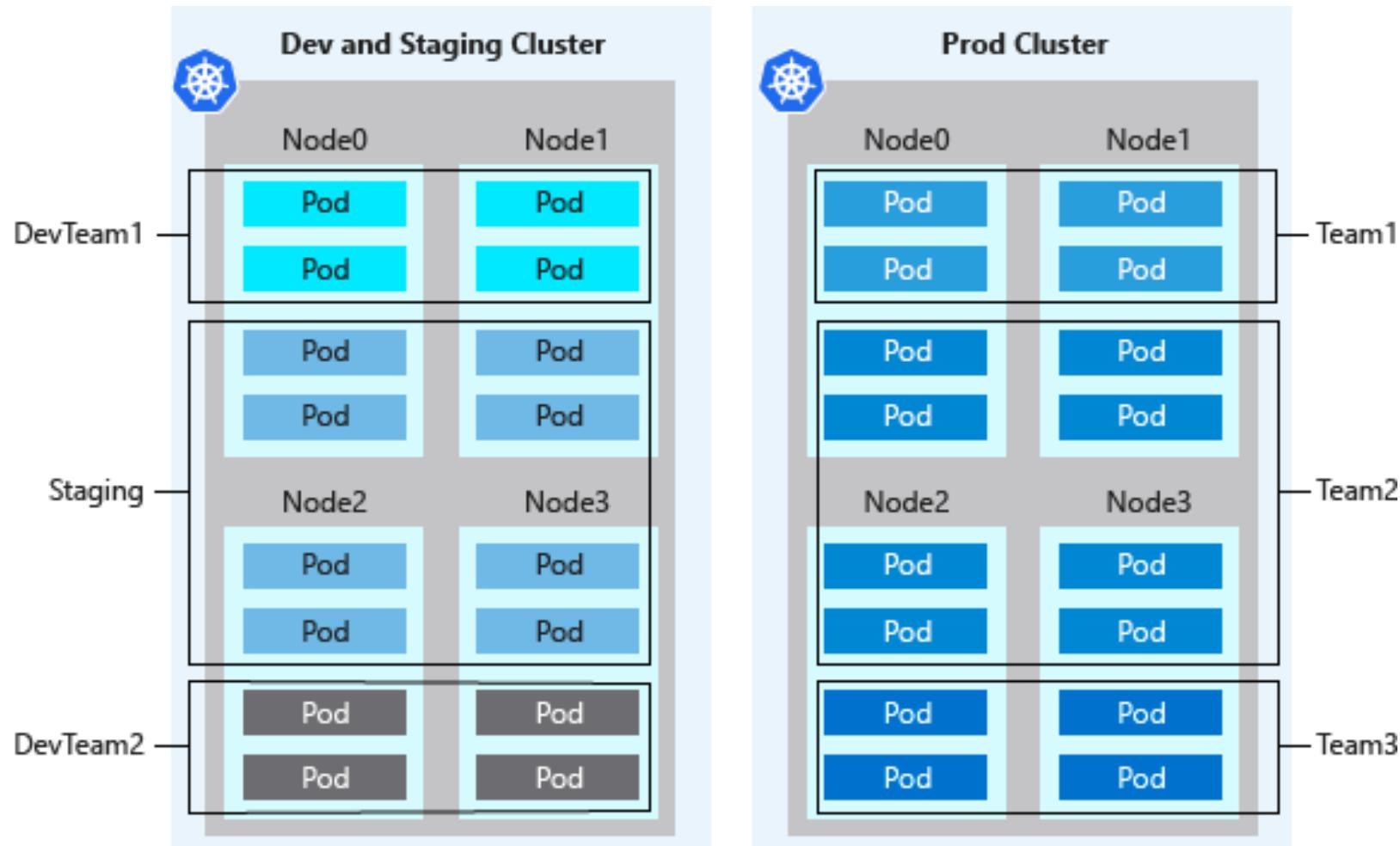
```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: nginx
  labels:
    app: nginx
    env: dev
    stack: frontend
spec:
  replicas: 3
  selector:
    matchLabels
      app: nginx
```

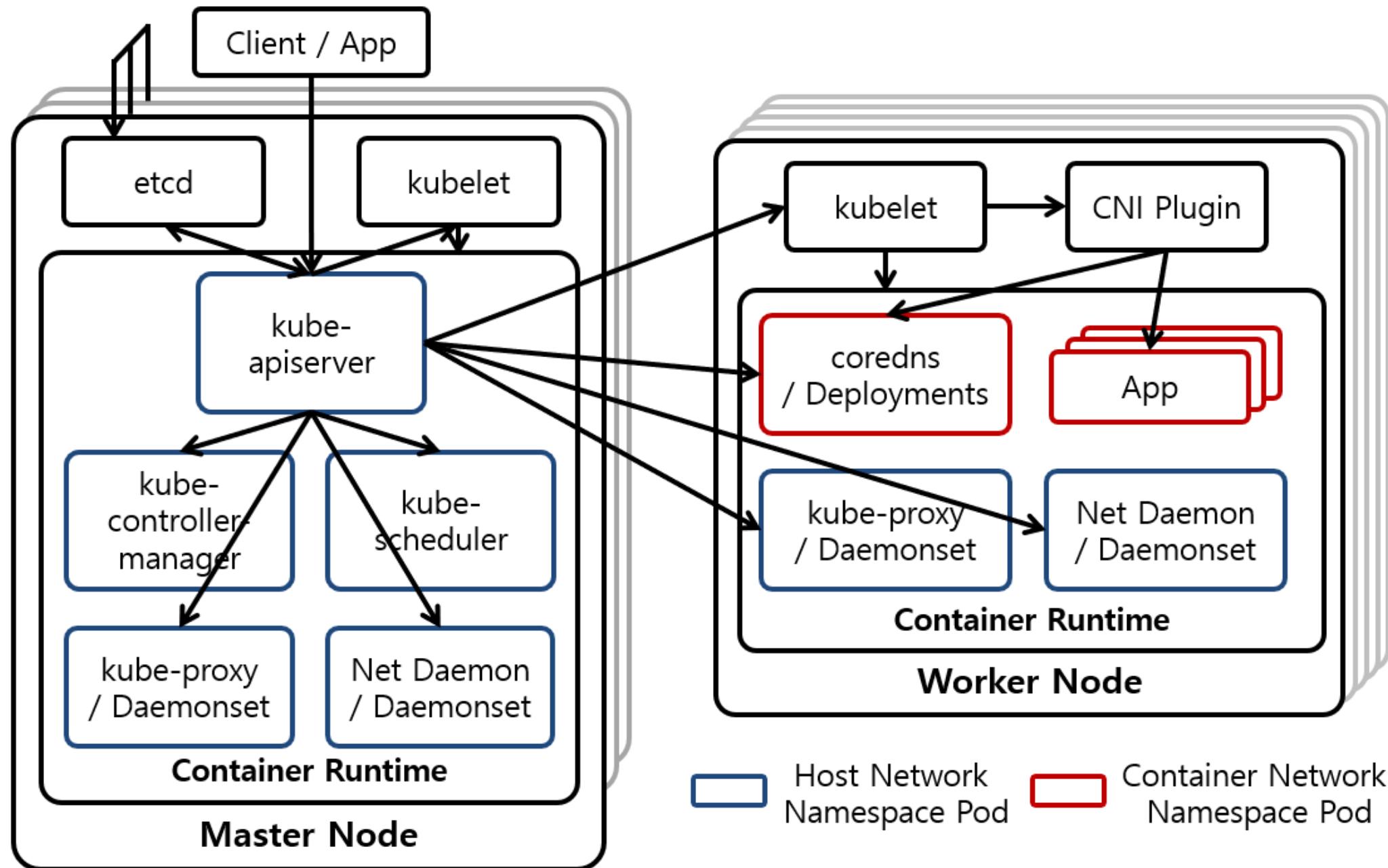


# Namespace, Node, Pod



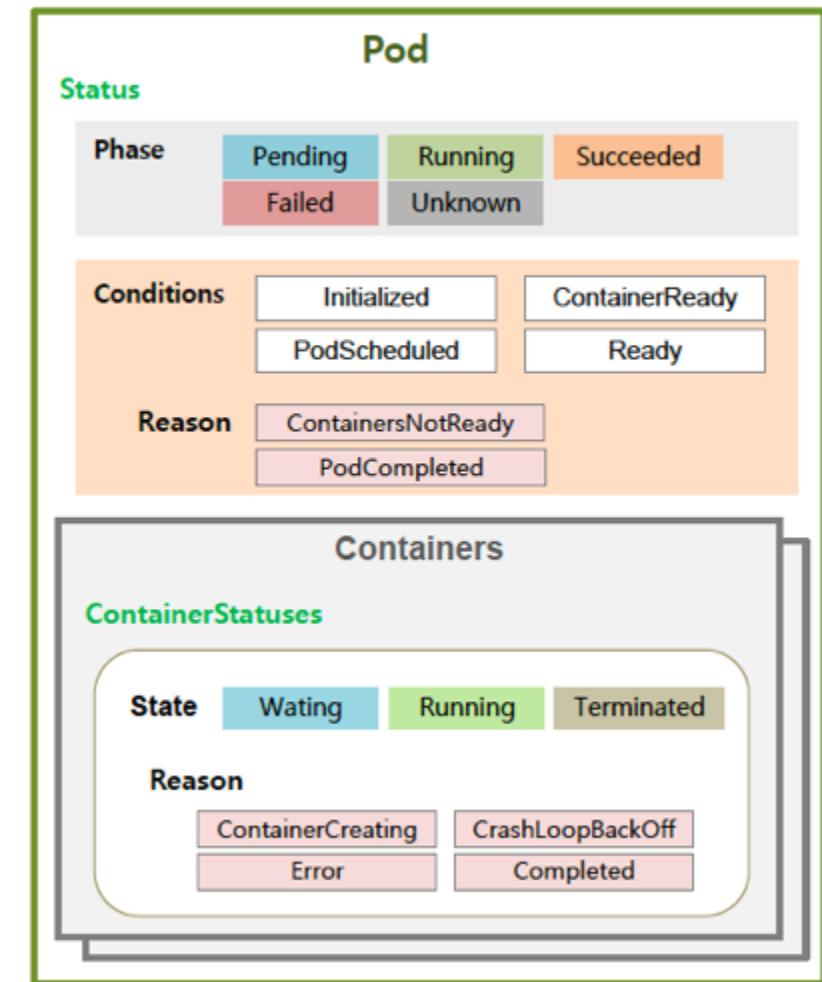
# Namespace Best Practice



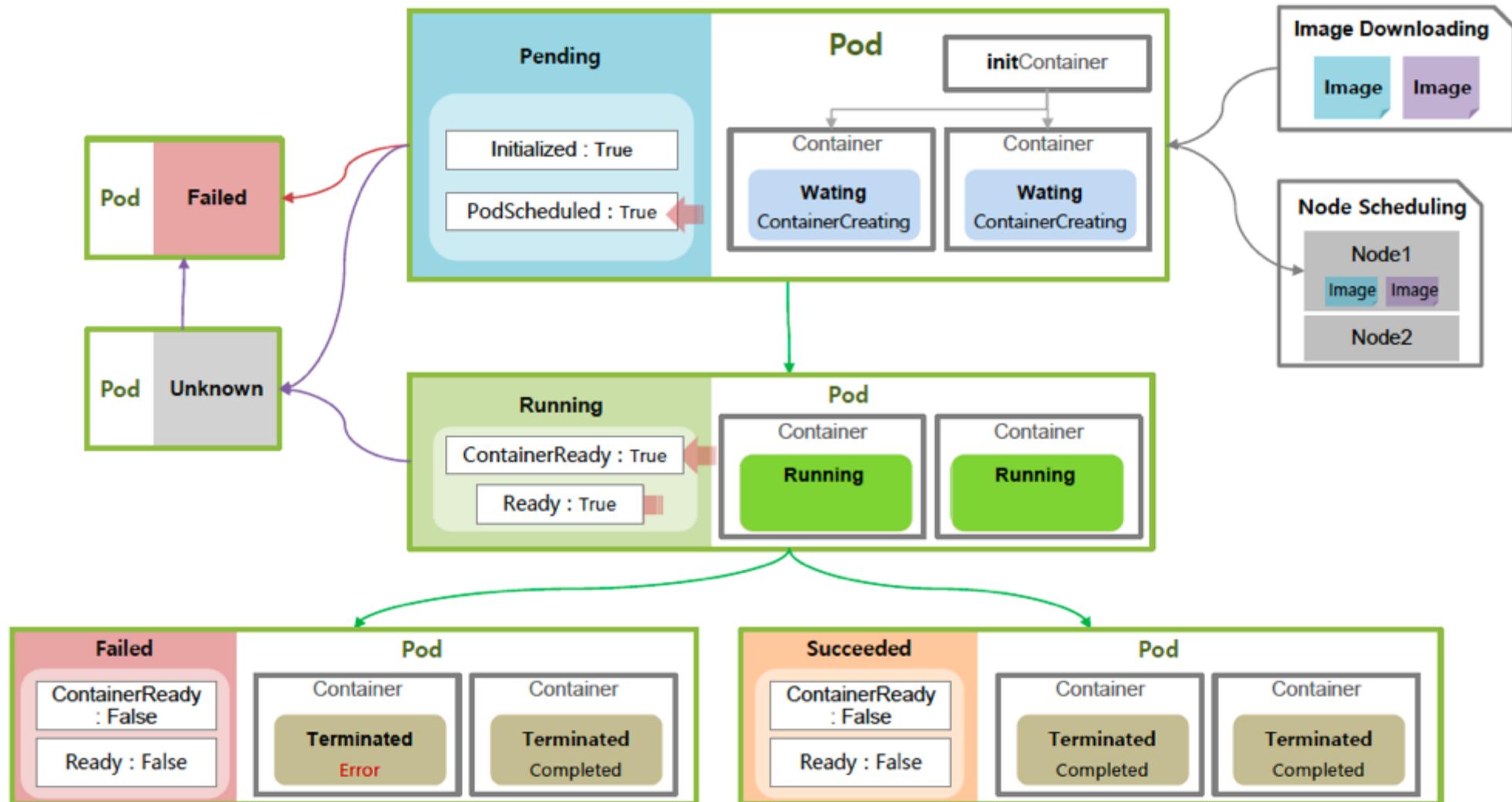


# Pod, Container Status

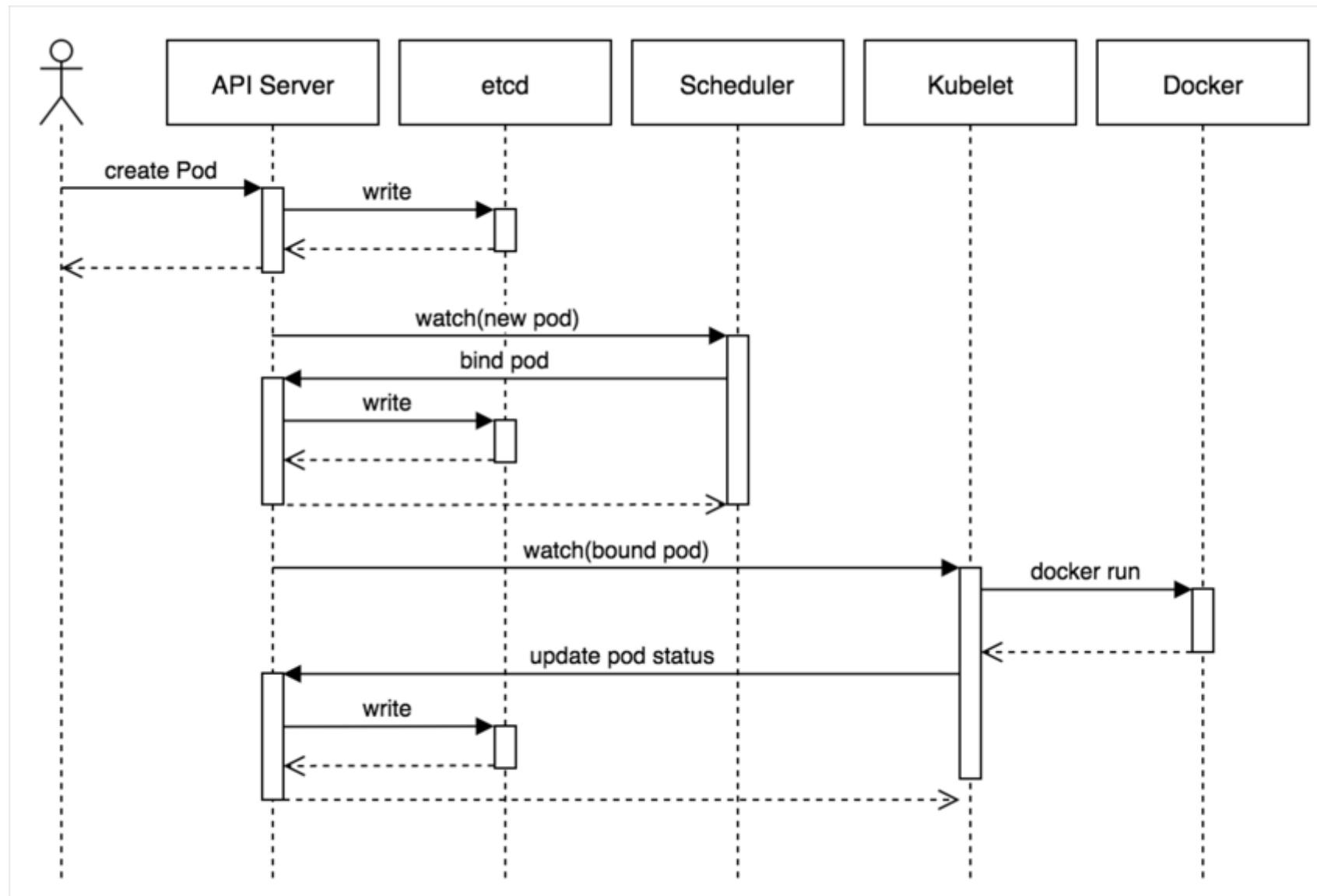
```
status:  
  phase: Pending  
  conditions:  
    - type: Initialized  
      status: 'True'  
      lastProbeTime: null  
      lastTransitionTime: '2019-09-26T22:07:56Z'  
  
    - type: PodScheduled  
      status: 'True'  
      lastProbeTime: null  
      lastTransitionTime: '2019-09-26T22:07:56Z'  
  
  - type: ContainersReady  
    status: 'False'  
    lastProbeTime: null  
    lastTransitionTime: '2019-09-26T22:08:11Z'  
    reason: ContainersNotReady  
  
  - type: Ready  
    status: 'False'  
    lastProbeTime: null  
    lastTransitionTime: '2019-09-26T22:08:11Z'  
    reason: ContainersNotReady  
  
containerStatuses:  
  - name: container  
    state:  
      waiting:  
        reason: ContainerCreating  
        lastState: {}  
        ready: false  
        restartCount: 0  
        image: tmkube/init  
        imageID: ""  
        started: false
```



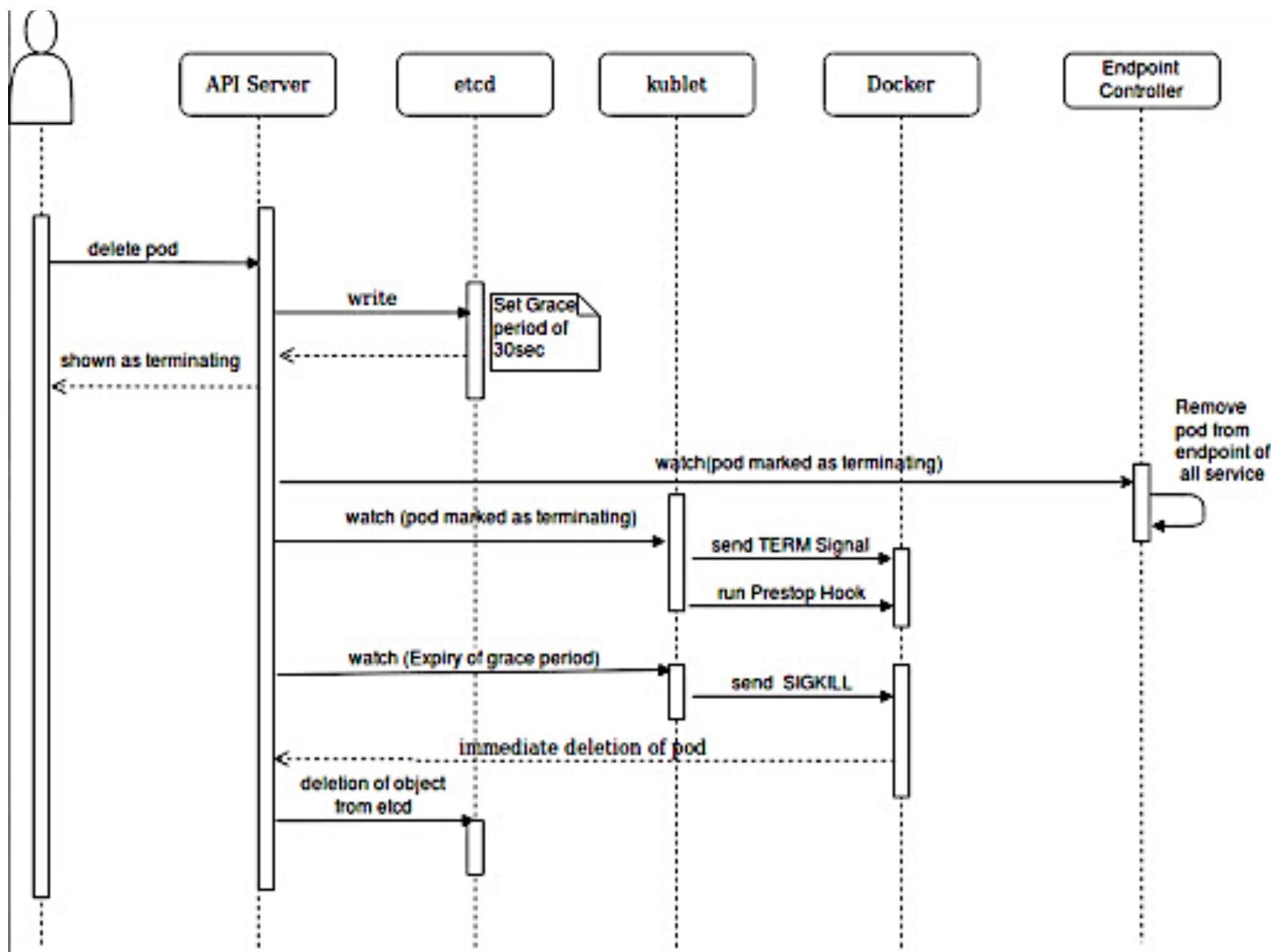
# Pod Lifecycle



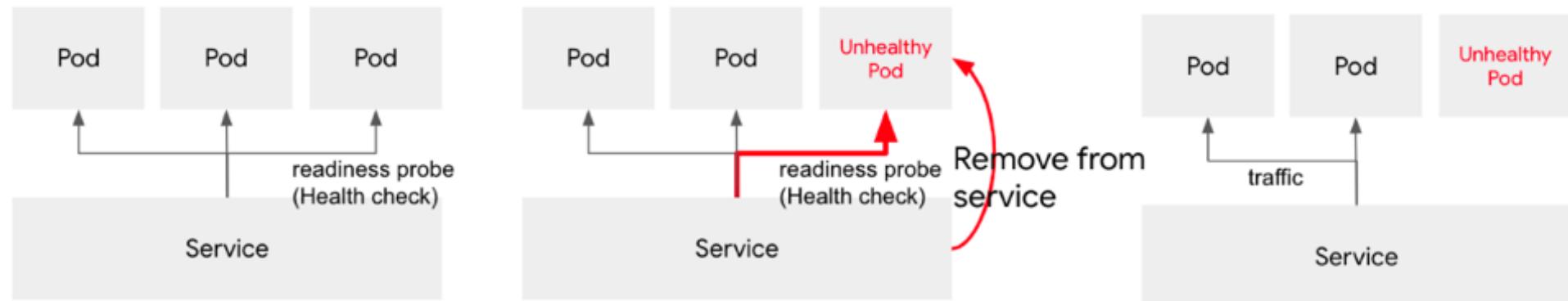
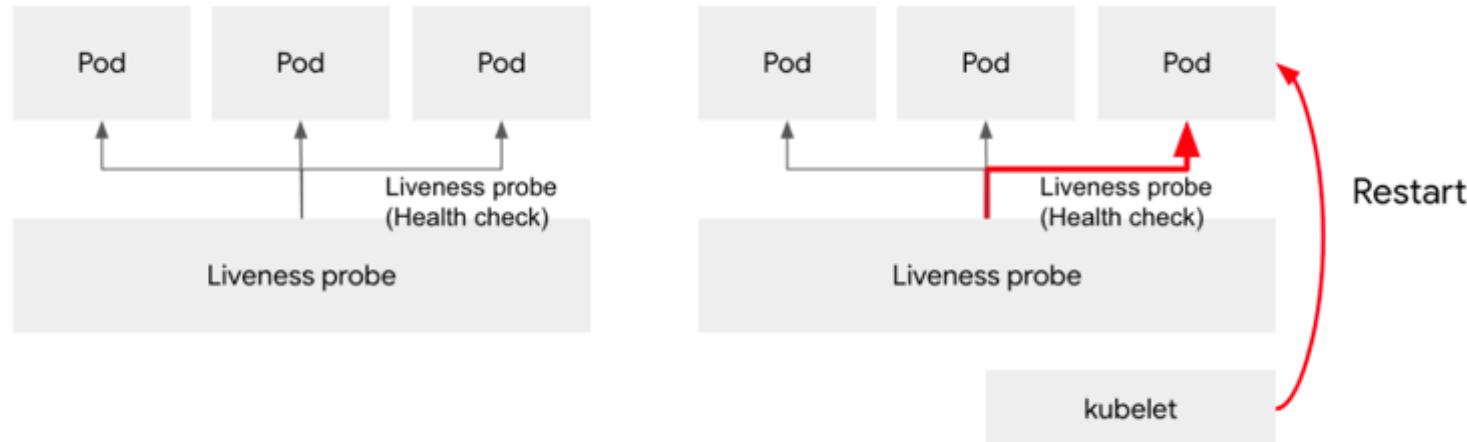
# Create Pod



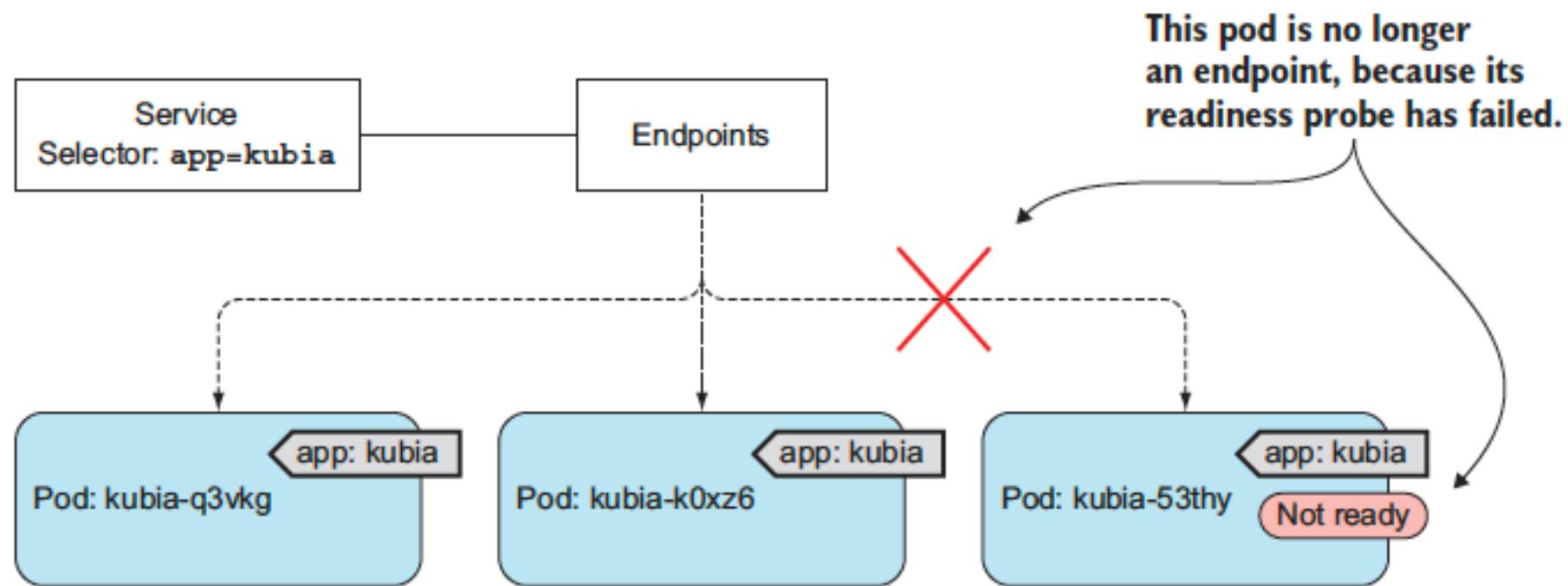
# Delete Pod



# Liveness Probe & Readiness Probe

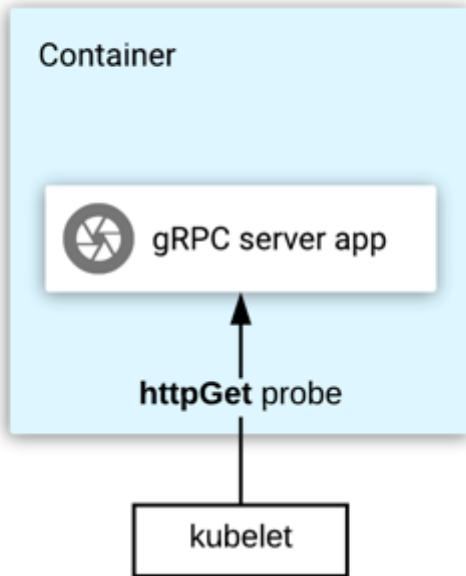


# Readiness, Service, Endpoints

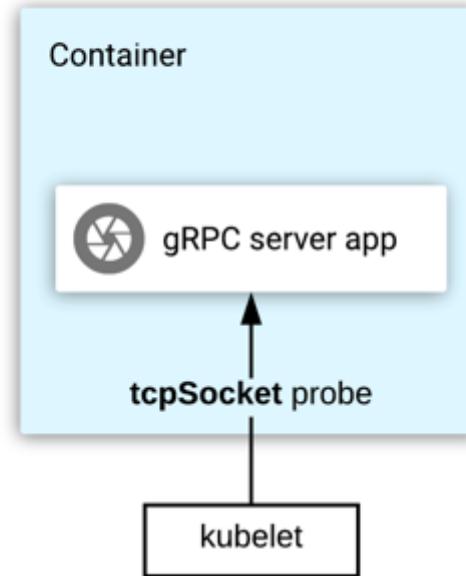


# Probe Methods (httpGet, tcpSocket, exec)

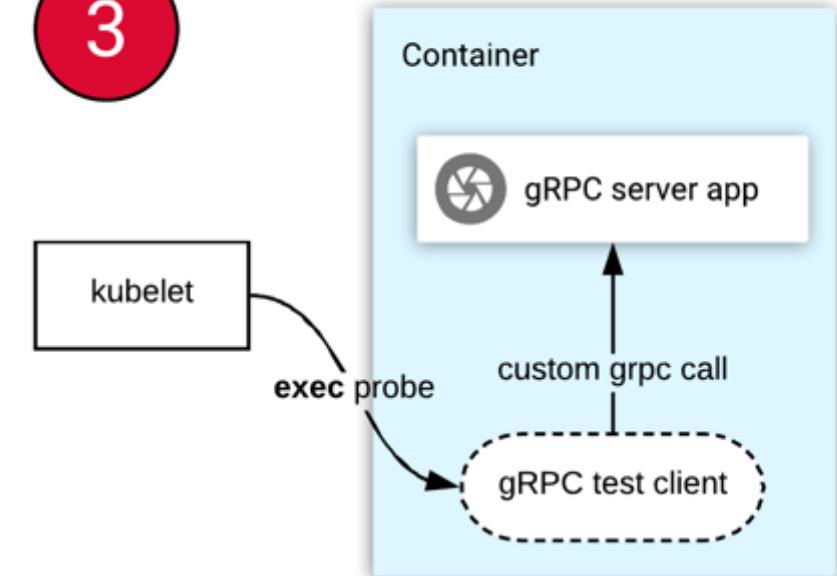
1



2



3

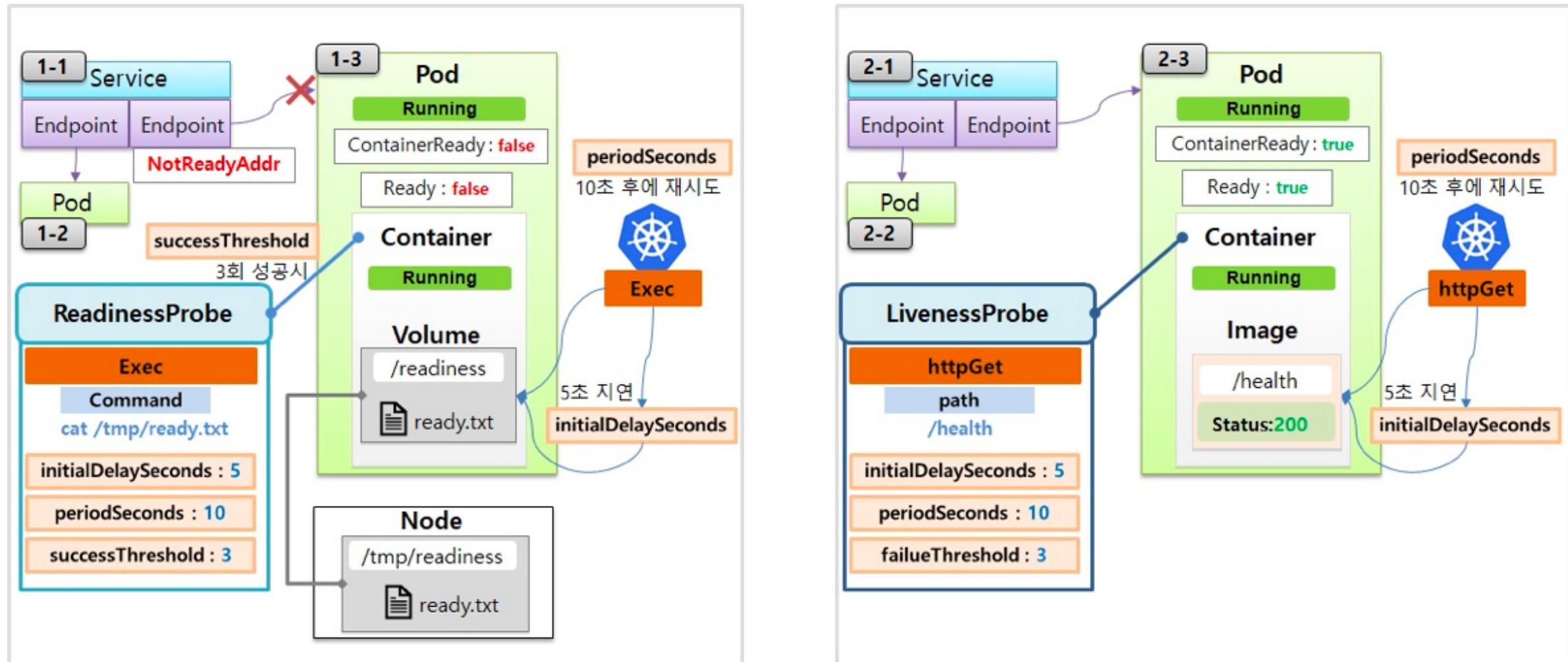


app needs to serve both grpc  
and http/1.1.

opening a tcp socket to a grpc  
server is not meaningful health  
check

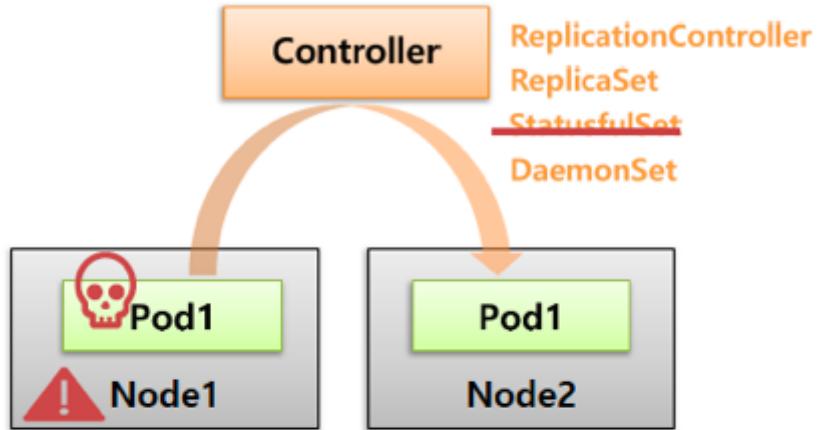
implement a client for your gRPC app, add it  
to container, kubelet invokes the client  
program, it calls a rpc that you implement

# Liveness & Readiness Probe Example

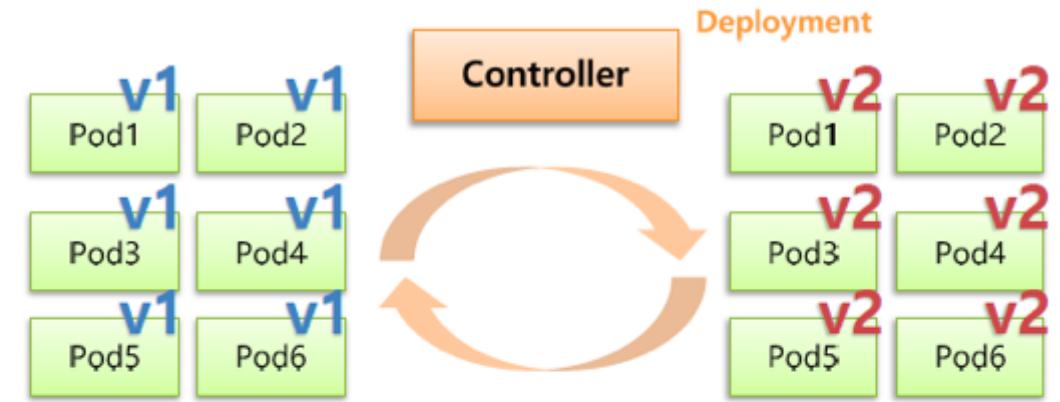


# Controllers (in Controller Manager)

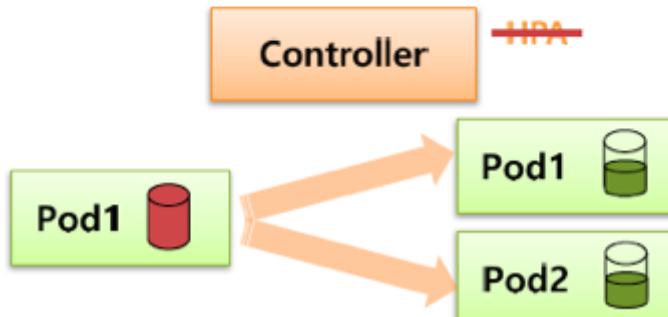
## ✚ Auto Healing



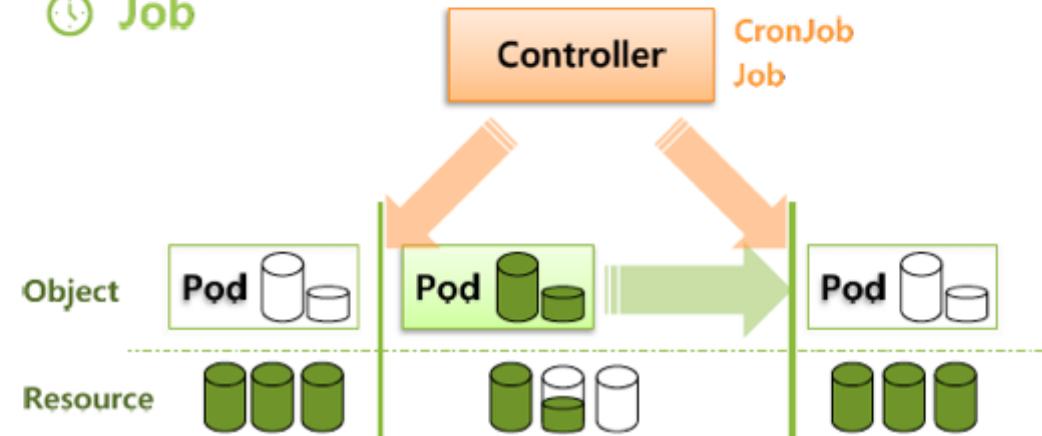
## ⟳ Software Update



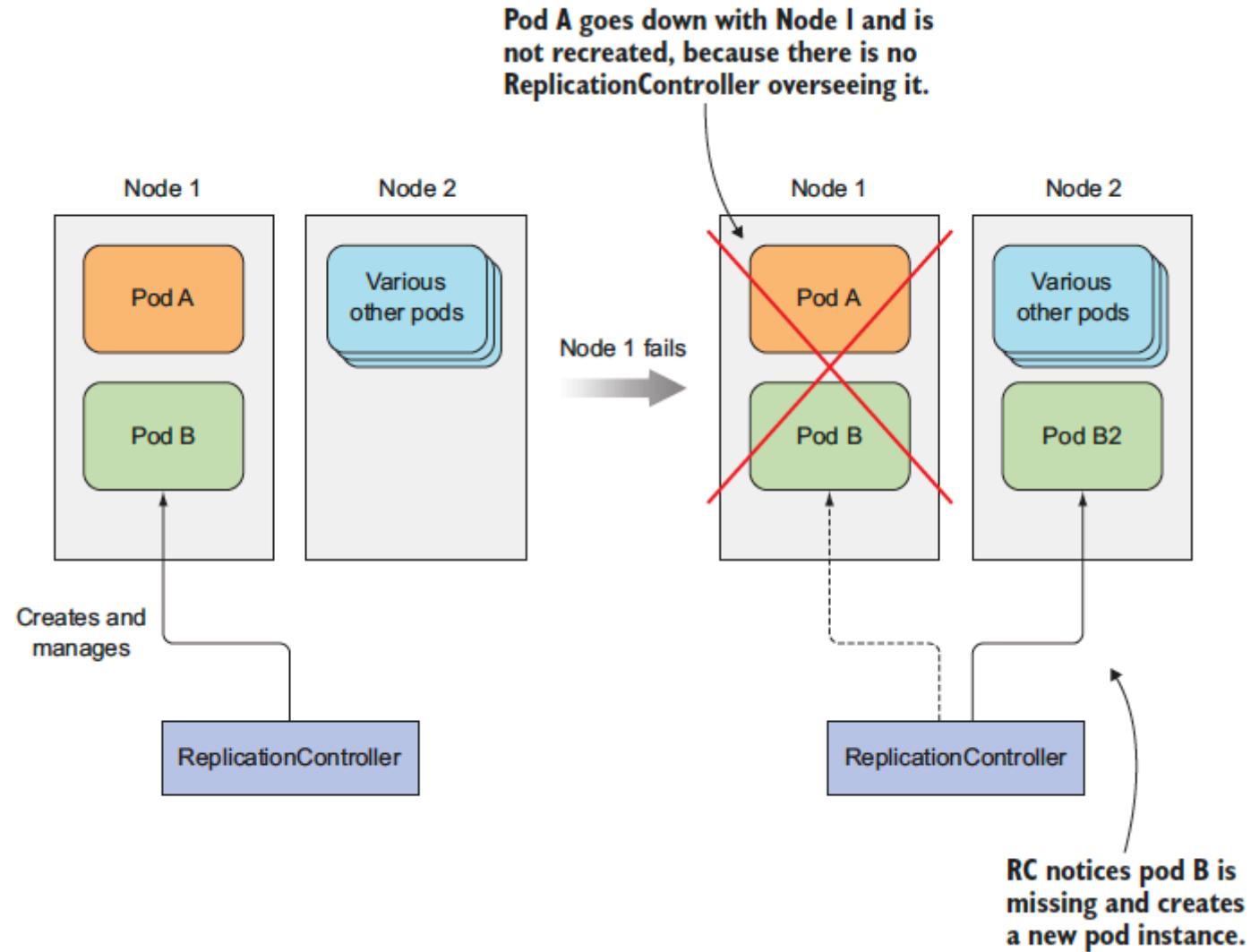
## ✳ Auto Scaling



## ⌚ Job



# ReplicationController



# ReplicationController: YAML File

```
apiVersion: v1
kind: ReplicationController
metadata:
  name: kubia
spec:
  replicas: 3
  selector:
    app: kubia
```

This manifest defines a  
ReplicationController (RC)

The name of this  
ReplicationController

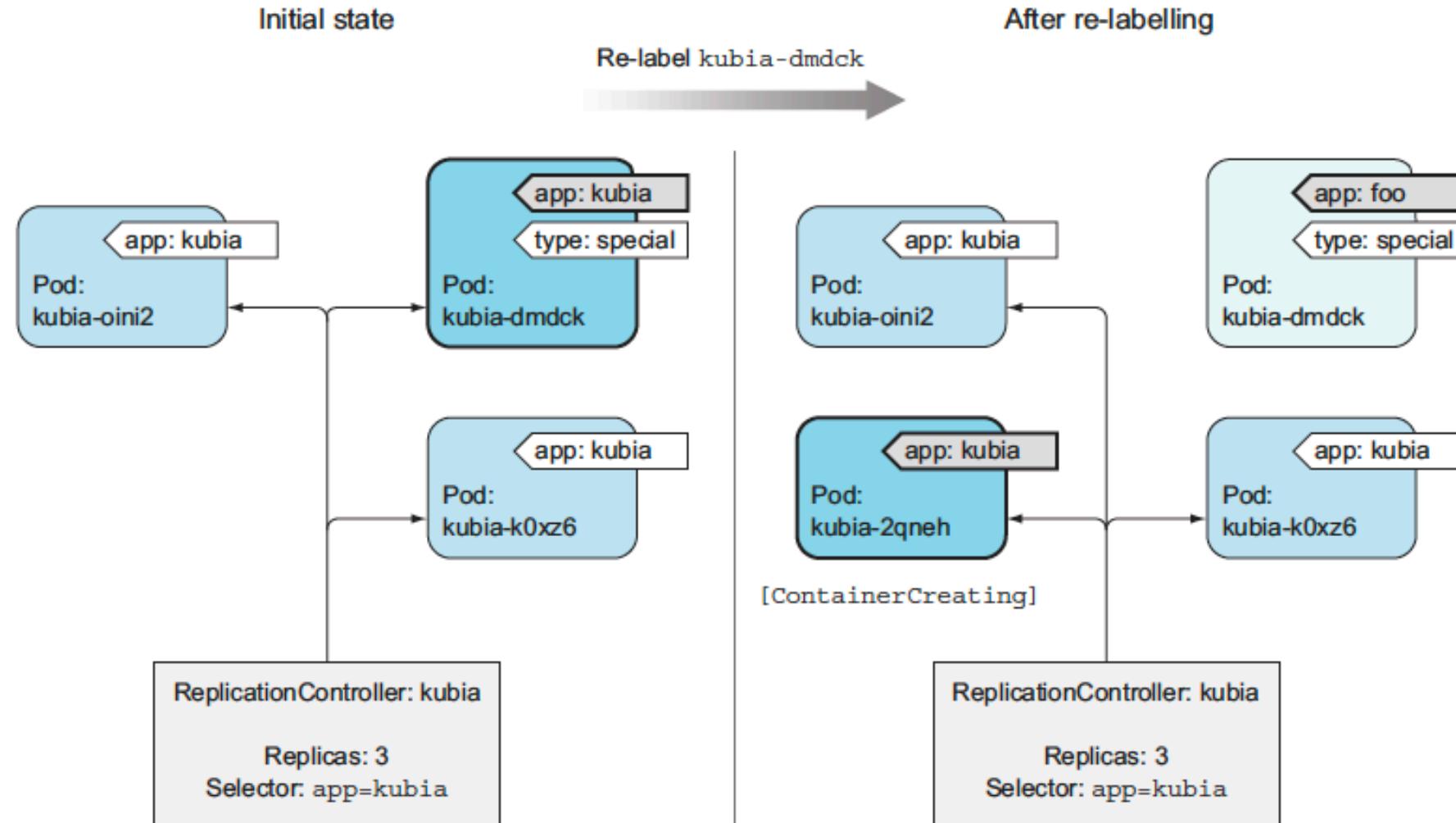
The desired number  
of pod instances

The pod selector determining  
what pods the RC is operating on

```
template:
  metadata:
    labels:
      app: kubia
spec:
  containers:
  - name: kubia
    image: luksa/kubia
    ports:
    - containerPort: 8080
```

The pod template  
for creating new  
pods

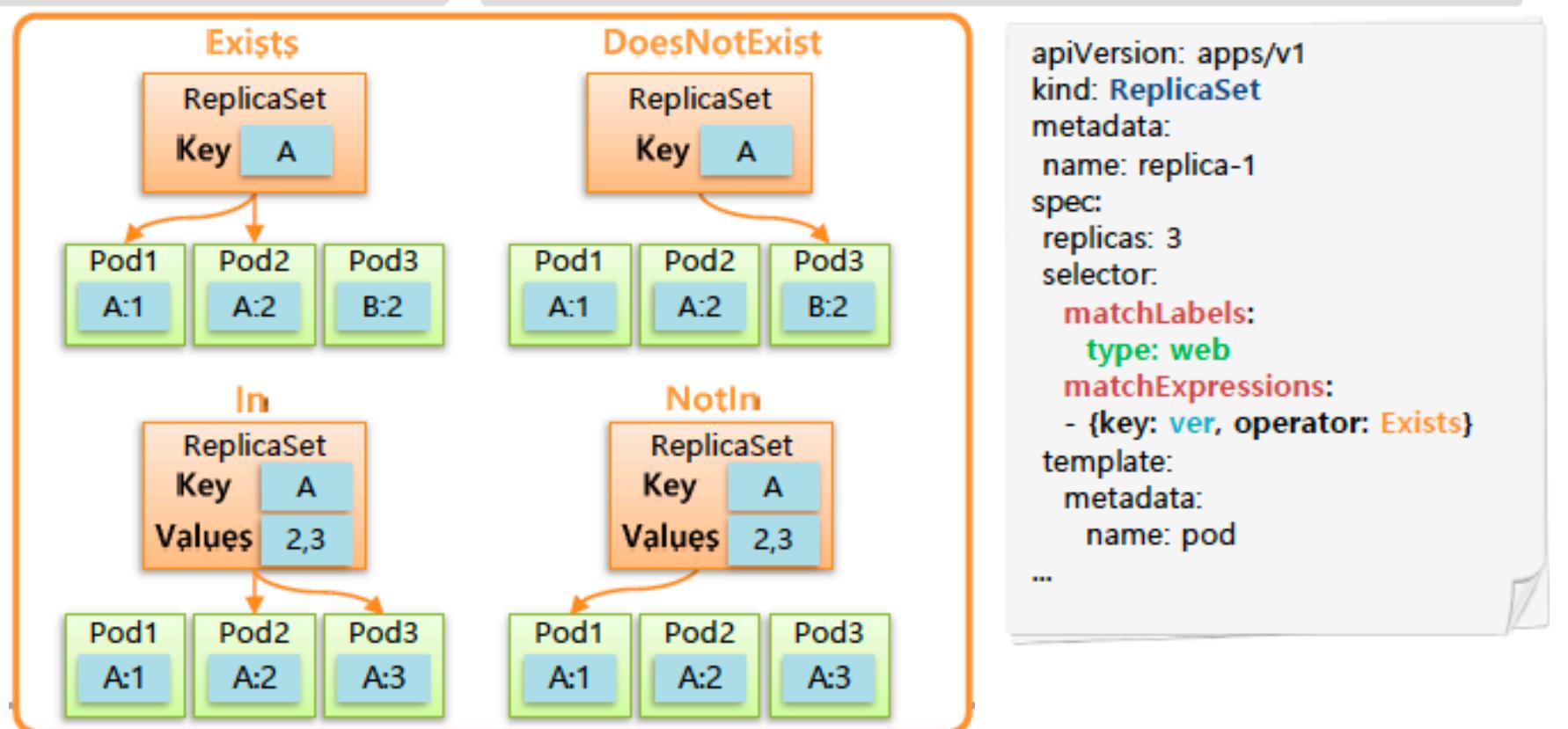
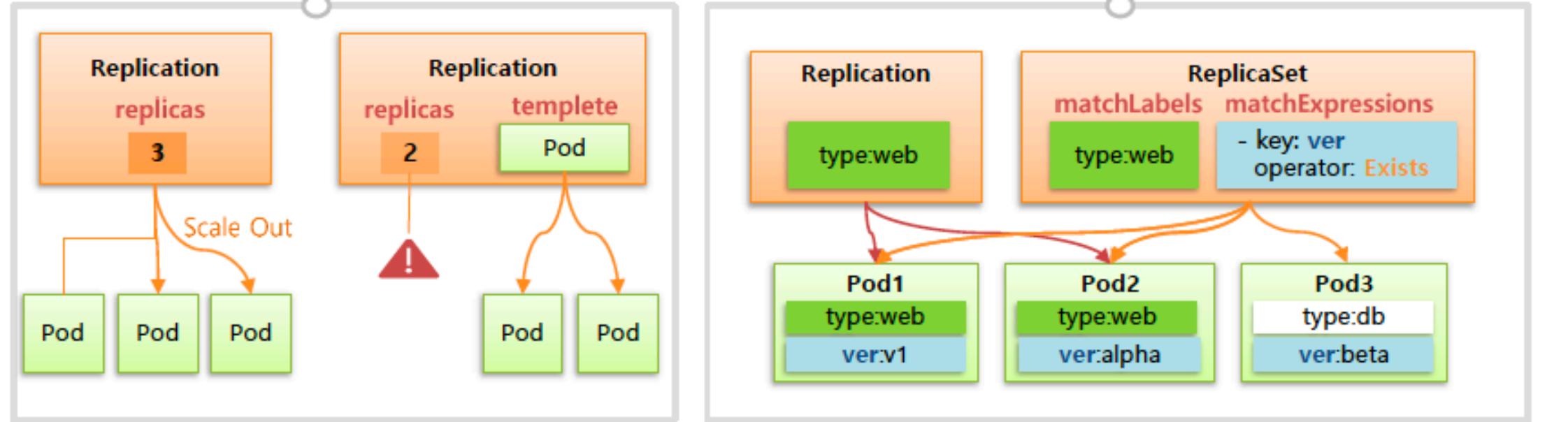
# ReplicationController: Re-Labelling



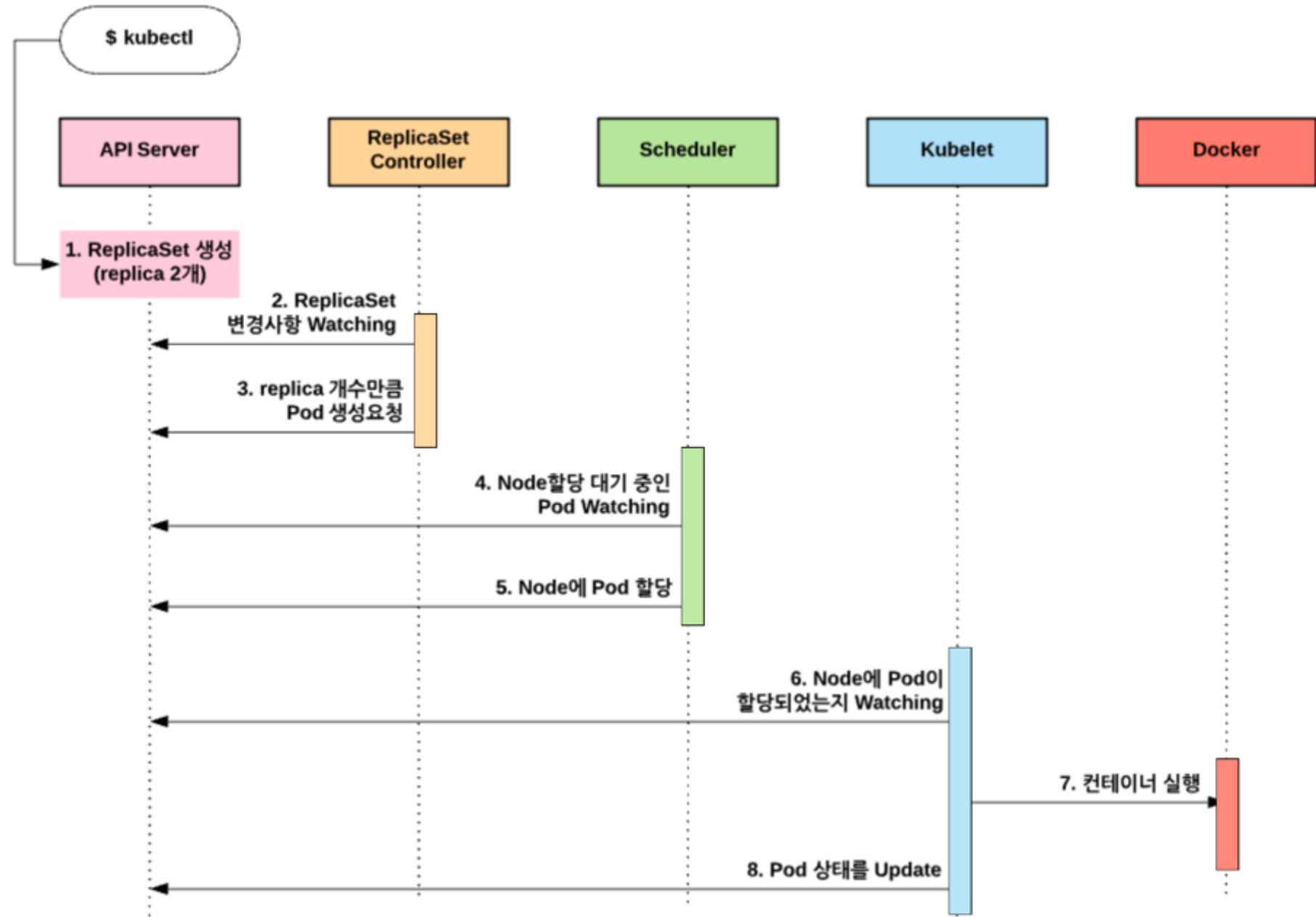
# ReplicaSet

- apiVersion: apps/v1beta2
- ReplicaSet is the next-generation Replication Controller

```
selector:  
  # ReplicationController는 matchLabels만 사용가능  
  matchLabels:  
    component: redis  
  matchExpressions:  
    - {key: tier, operator: In, values: [cache]}  
    - {key: environment, operator: NotIn, values: [dev]}  
    - {key: service, operator: Exists, values: [user]}  
    - {key: service, operator: DoesNotExist, values: [db]}
```

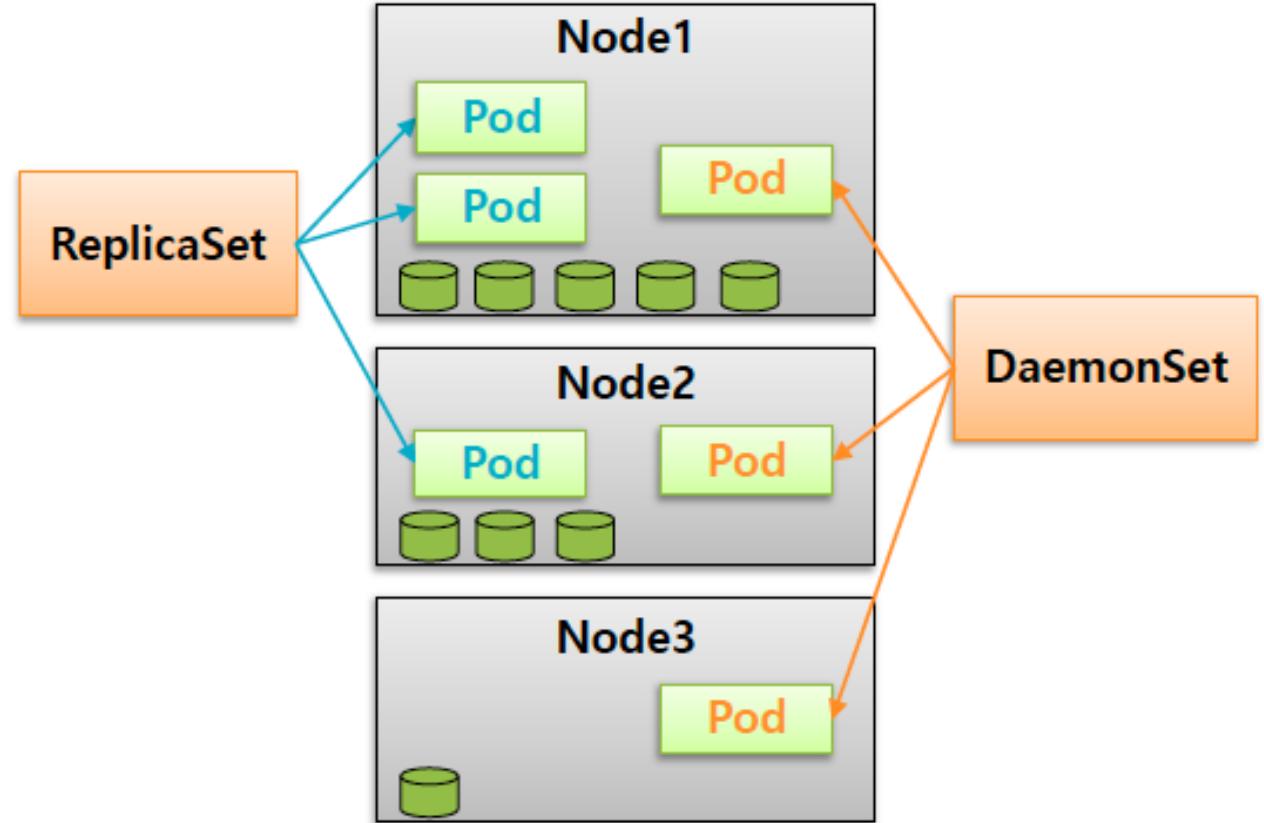


# ReplicaSet



# DaemonSet

- apiVersion: apps/v1beta2



Performance



Prometheus

Logging



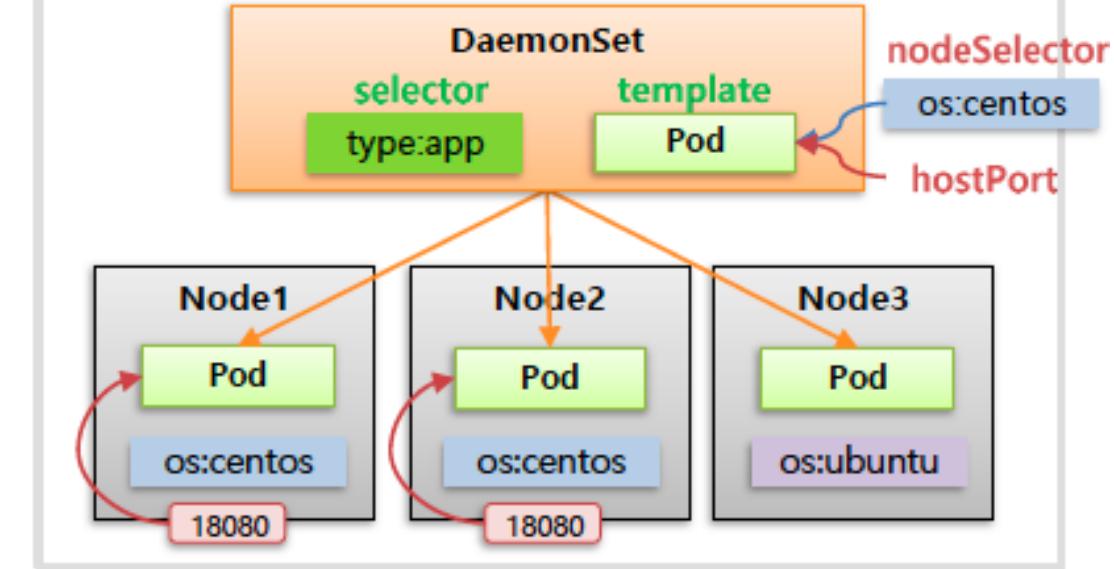
fluentd

Storage

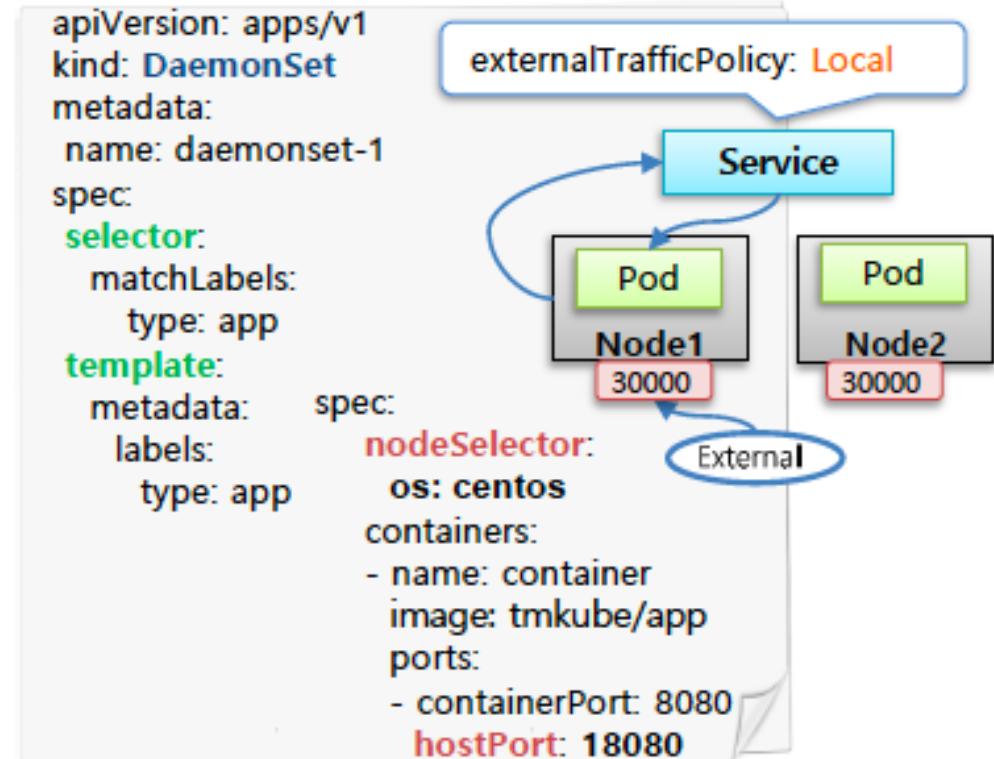


GlusterFS

# DaemonSet: YAML File

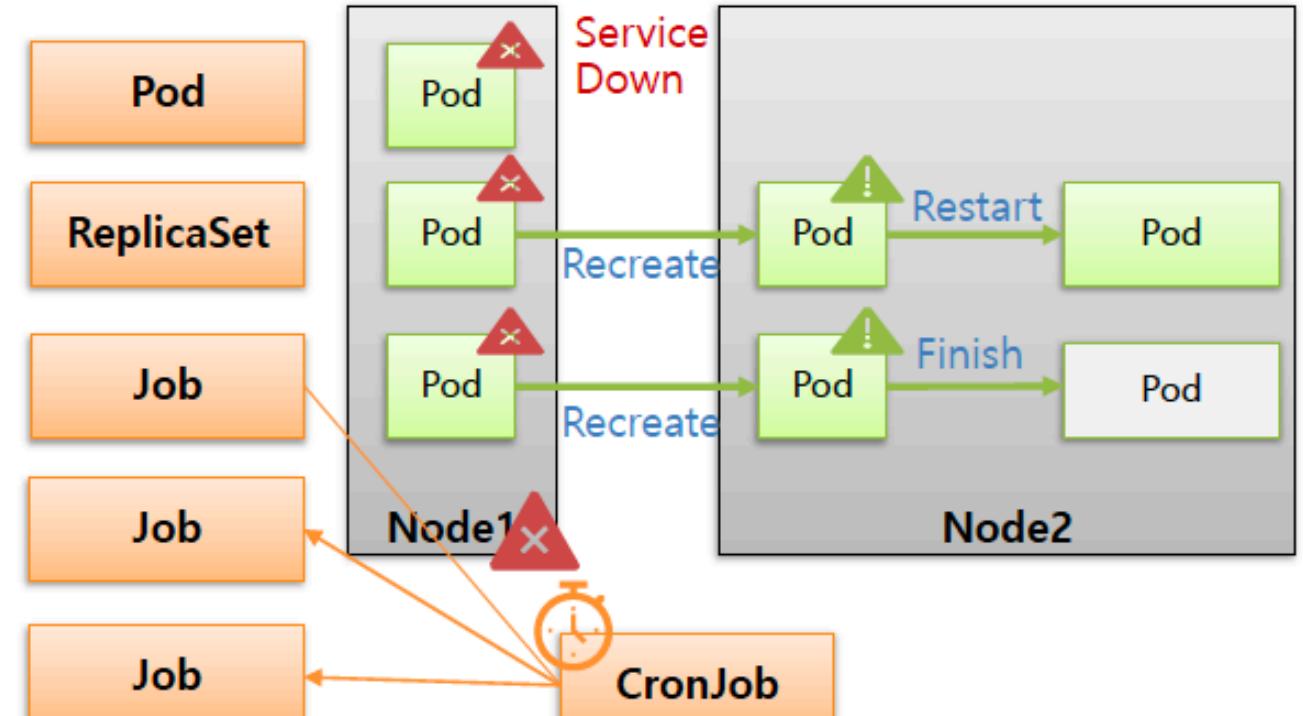


```
apiVersion: apps/v1
kind: DaemonSet
metadata:
  name: daemonset-1
spec:
  selector:
    matchLabels:
      type: app
  template:
    metadata:
      labels:
        type: app
    spec:
      nodeSelector:
        os: centos
      containers:
        - name: container
          image: tmkube/app
          ports:
            - containerPort: 8080
              hostPort: 18080
      hostPort: 18080
```



# Job, CronJob

- apiVersion: batch/v1



Backup



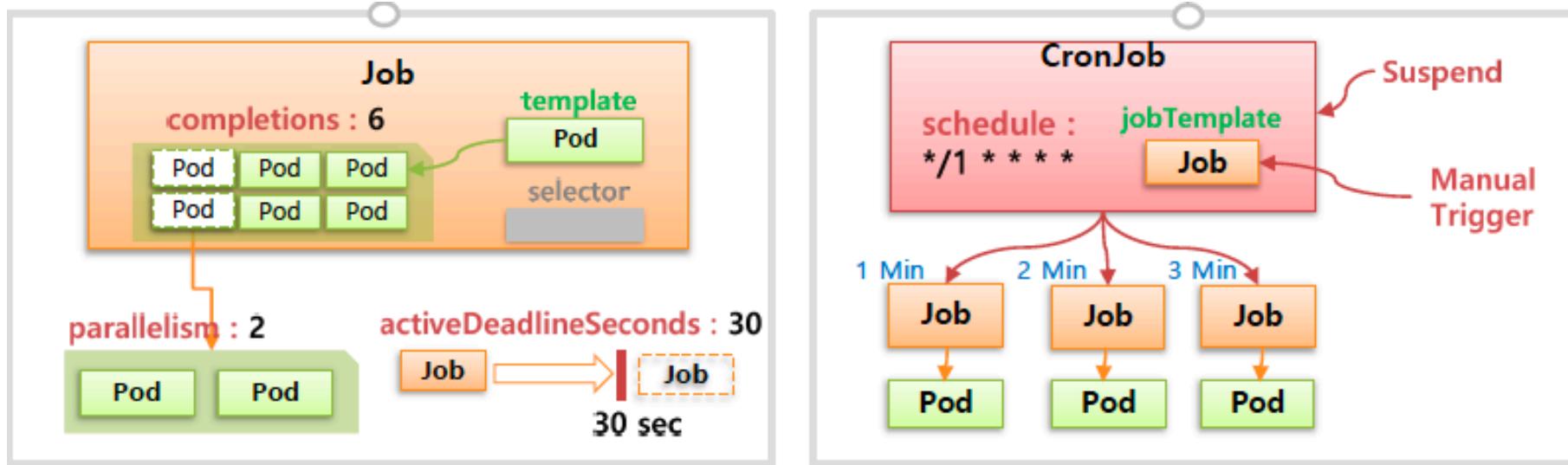
Checking



Messaging



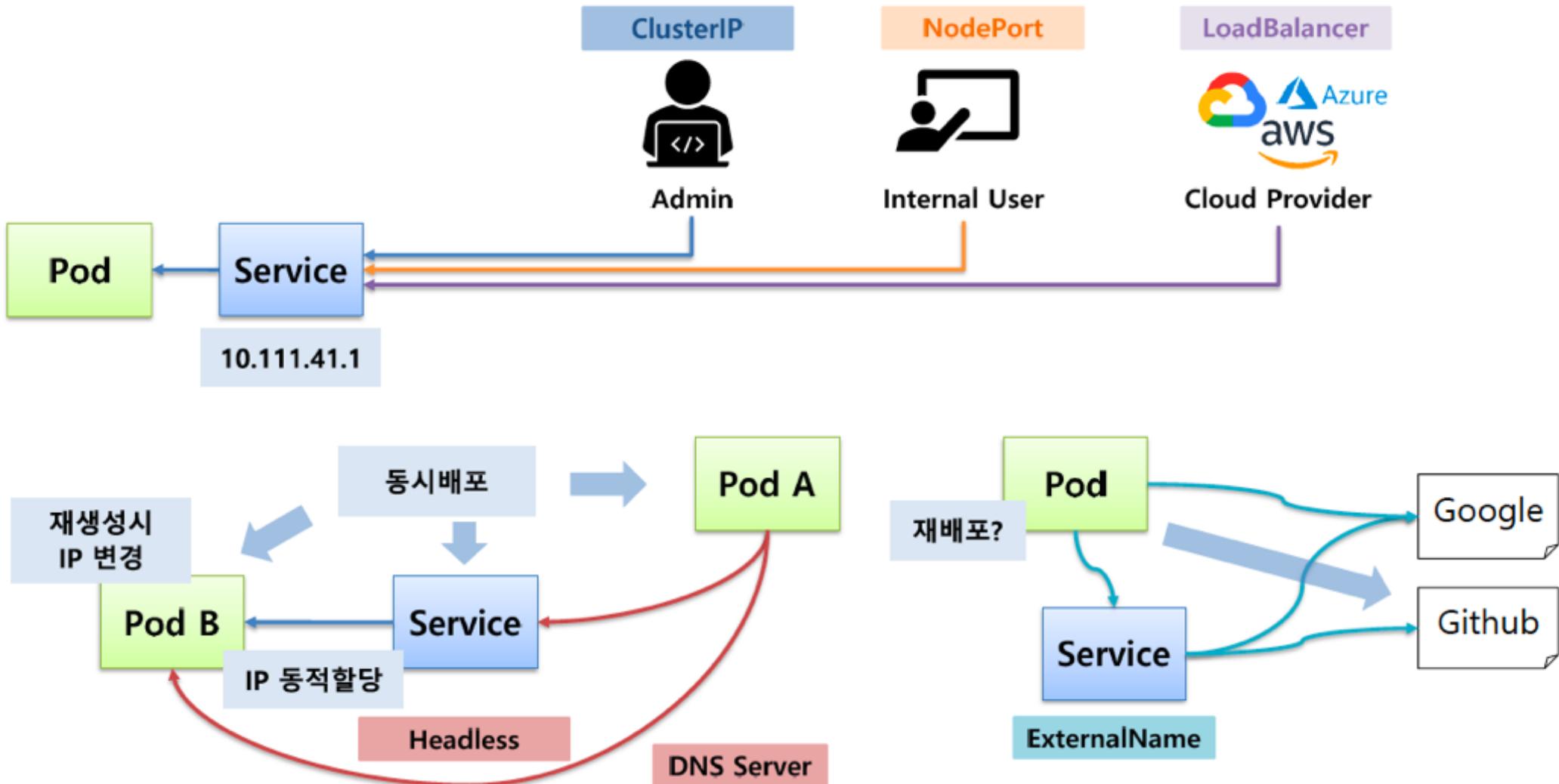
# Job, CronJob: YAML File

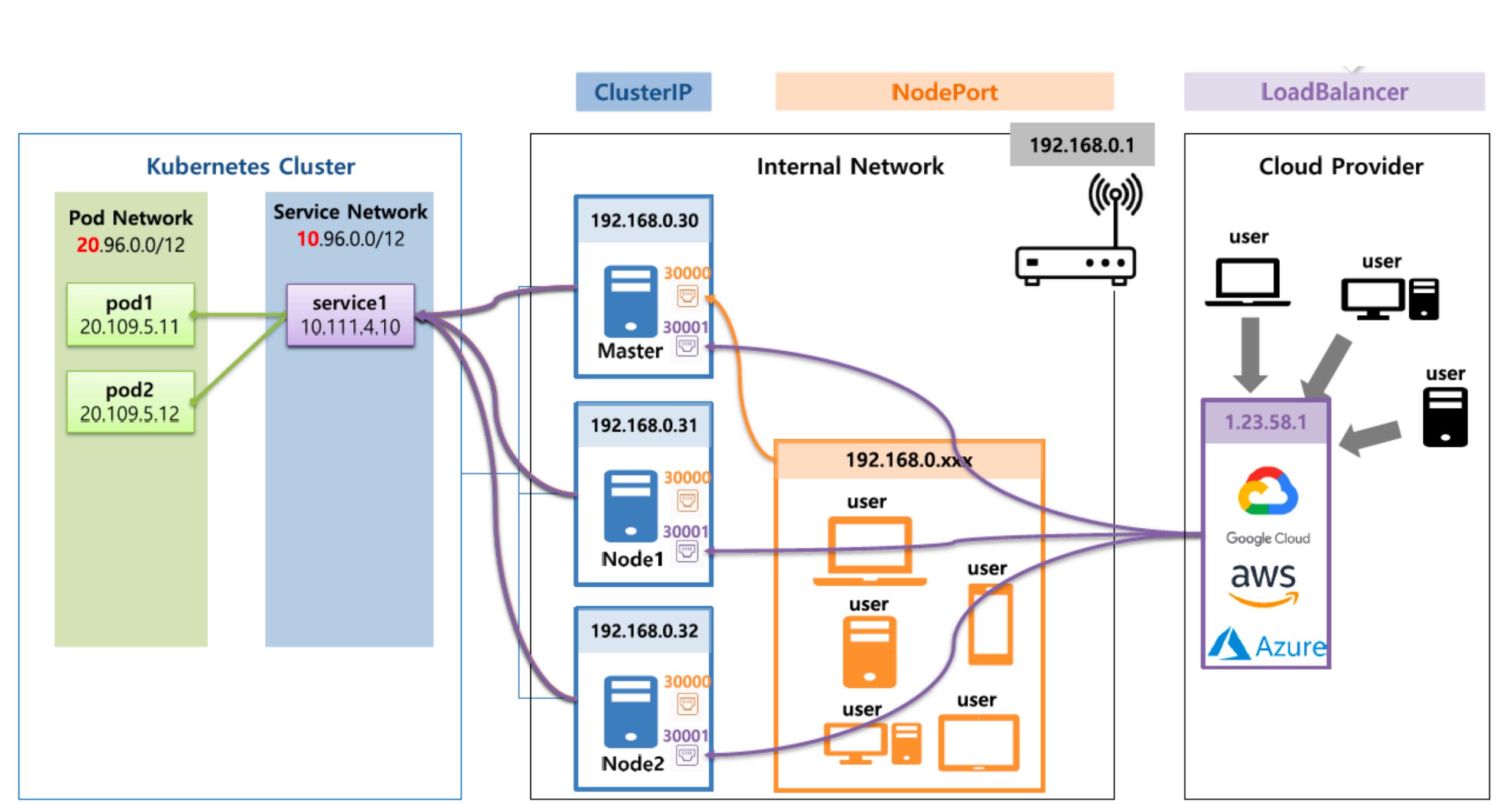


```
apiVersion: batch/v1
kind: Job
metadata:
  name: job-1
spec:
  completions: 6
  parallelism: 2
  activeDeadlineSeconds: 30
  template:
    spec:
      restartPolicy: Never
      containers:
        - name: container
          image: tmkube/init
```

```
apiVersion: batch/v1
kind: CronJob
metadata:
  name: cron-job
spec:
  schedule: "*/1 * * * *"
  concurrencyPolicy: Allow
  jobTemplate:
    spec:
      template:
        spec:
          restartPolicy: Never
          containers:
            - name: container
              image: tmkube/app
```

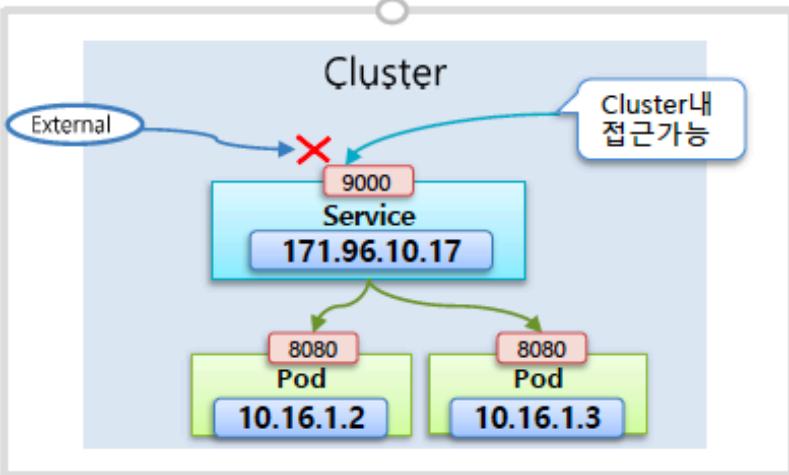
# Service: ClusterIP, NodePort, LoadBalancer, ExternalName, Headless





# Service: ClusterIP, NodePort, LoadBalancer

- 서비스 디버깅
- 내부 트래픽/대쉬보드
- 인증된 사용자 연결



```

apiVersion: v1
kind: Service
metadata:
  name: svc-1
spec:
  selector:
    app: pod
  ports:
    - port: 9000
      targetPort: 8080
      type: ClusterIP

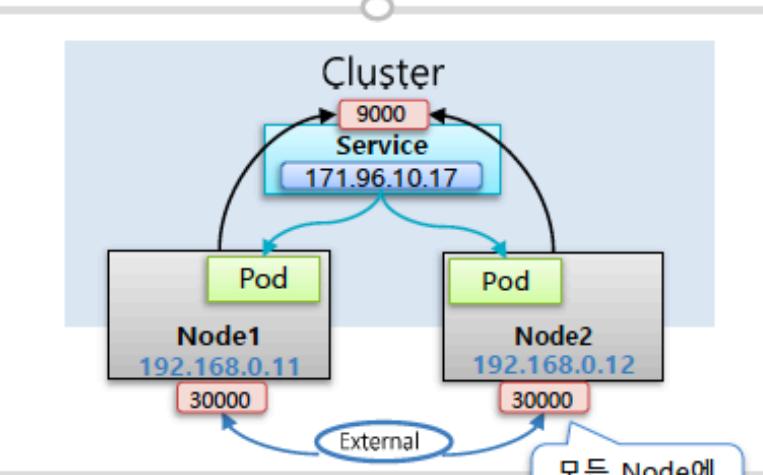
```

```

apiVersion: v1
kind: Pod
metadata:
  name: pod-1
  labels:
    app: pod
spec:
  containers:
    - name: container
      image: tmkube/app
      ports:
        - containerPort: 8080

```

- 데모나 임시 연결용
- 프로덕션 환경(X)

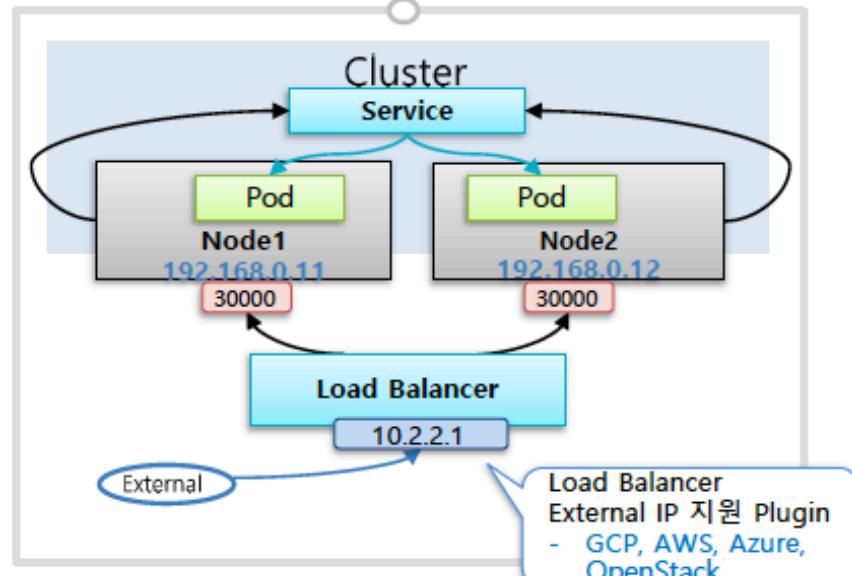


```

apiVersion: v1
kind: Service
metadata:
  name: svc-2
spec:
  selector:
    app: pod
  ports:
    - port: 9000
      targetPort: 8080
      nodePort: 30000
      type: NodePort
  ...
externalTrafficPolicy: Local

```

- 외부 시스템 노출용



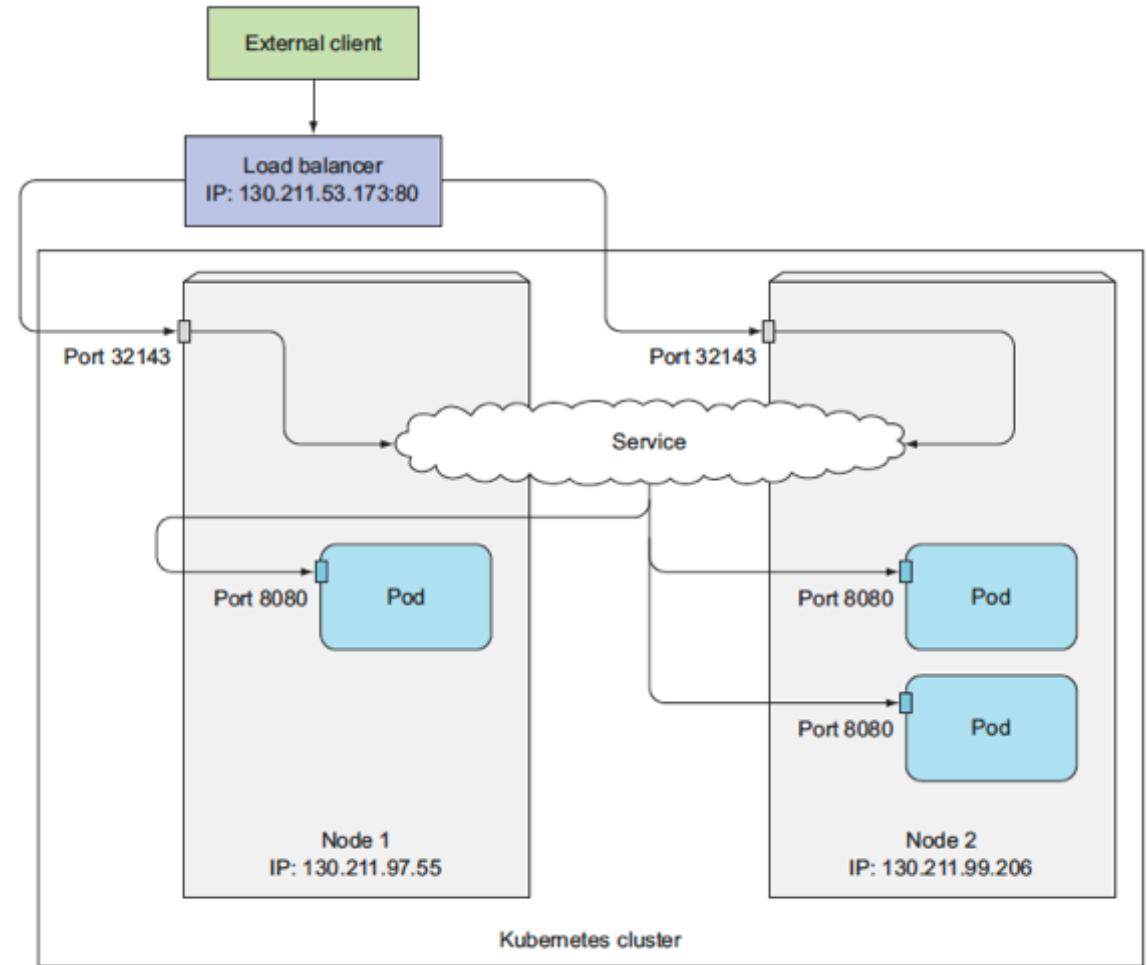
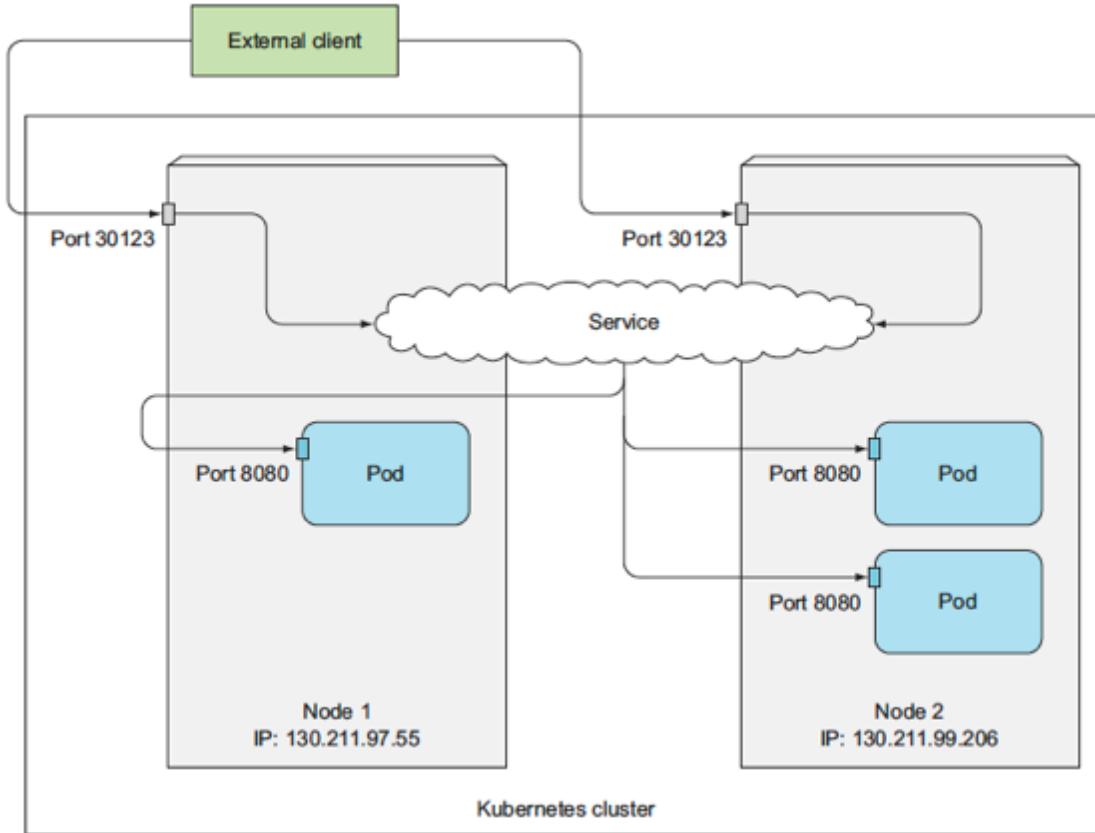
```

apiVersion: v1
kind: Service
metadata:
  name: svc-3
spec:
  selector:
    app: pod
  ports:
    - port: 9000
      targetPort: 8080
      type: LoadBalancer

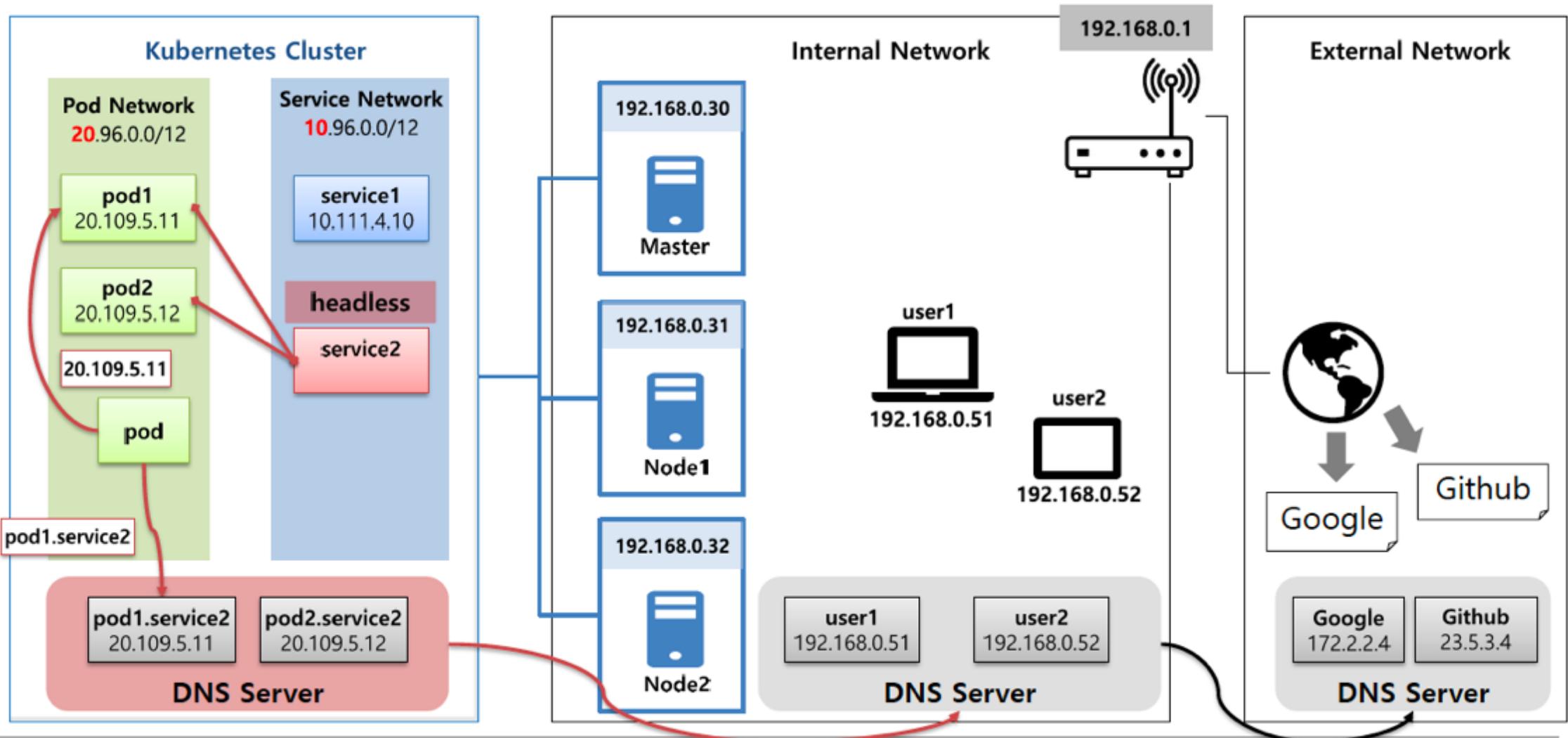
```

Load Balancer External IP 지원 Plugin  
- GCP, AWS, Azure, OpenStack

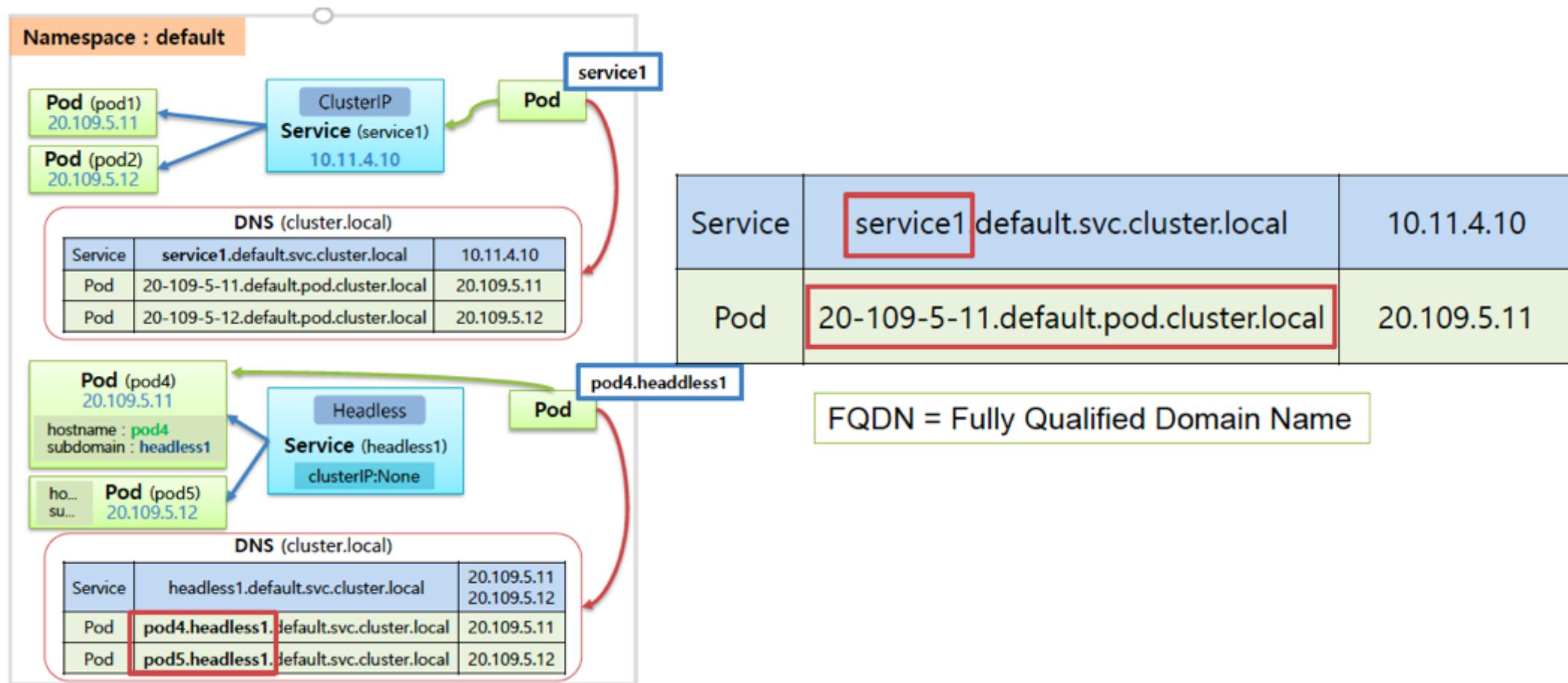
# NodePort vs LoadBalancer



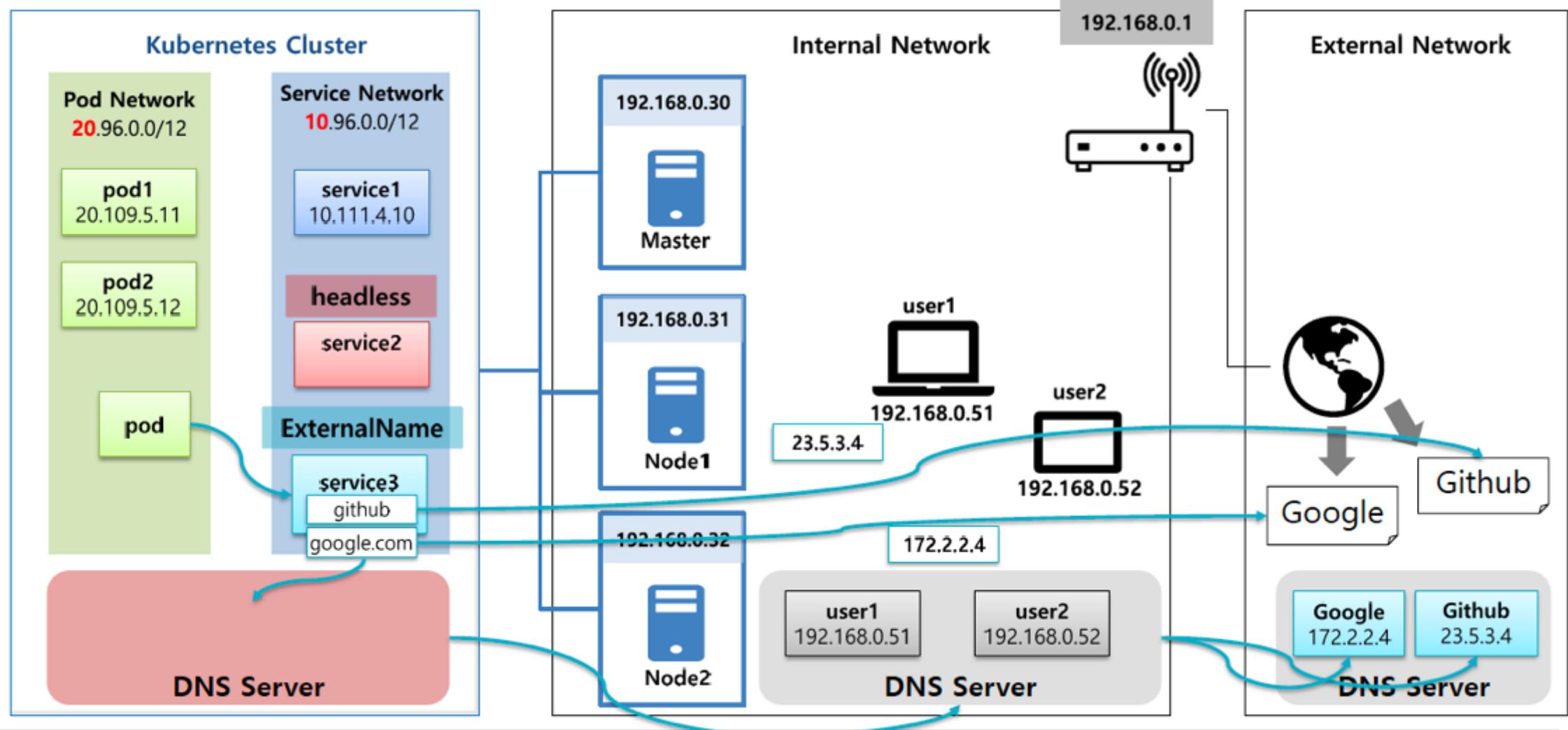
# Headless (1)



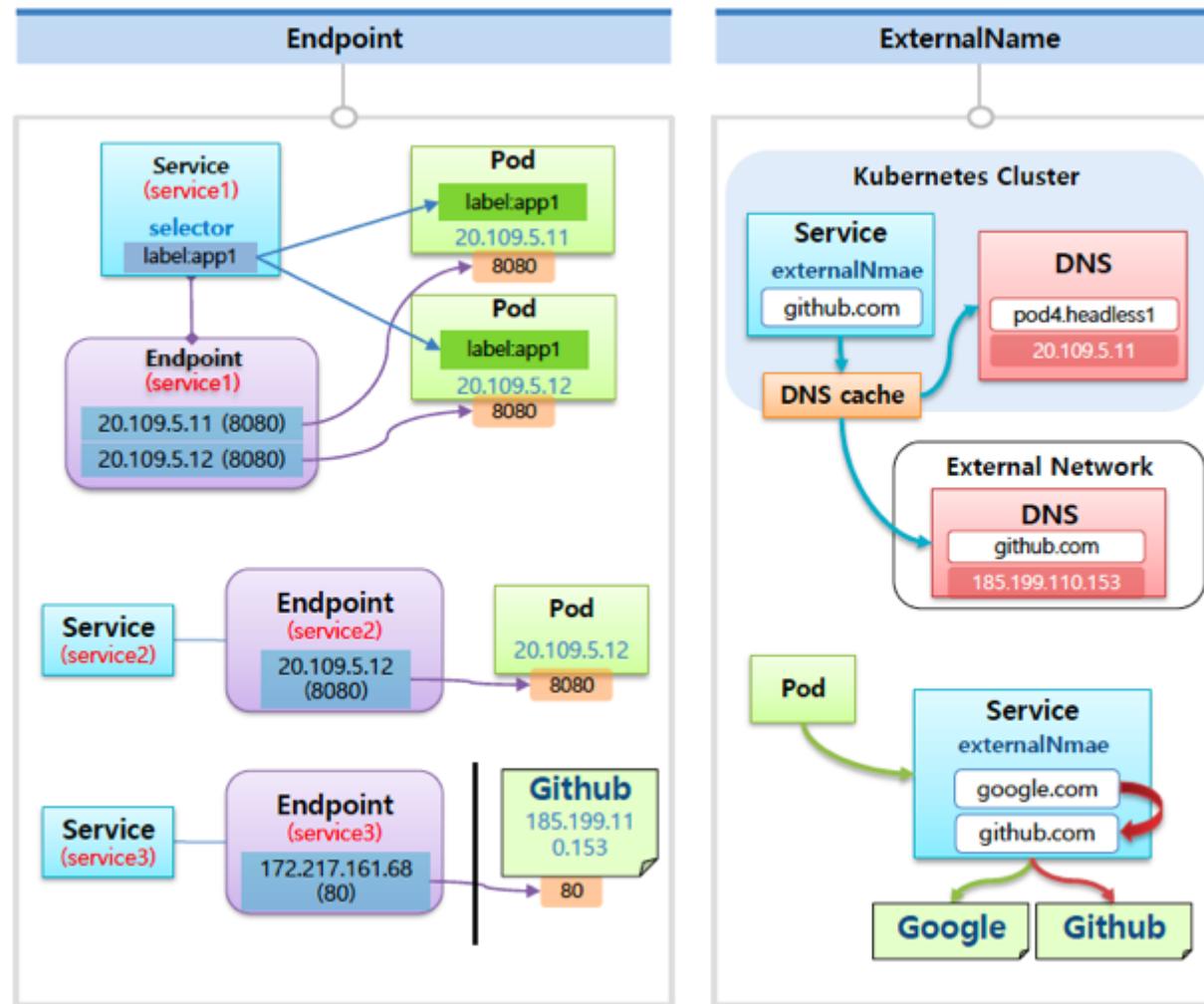
# Headless (2)



# ExternalName

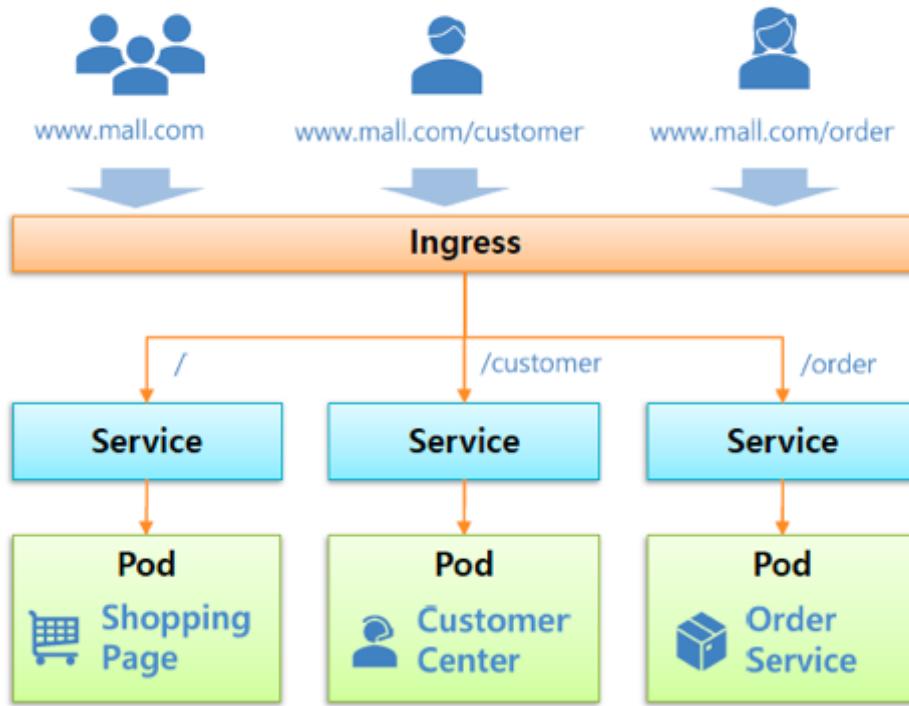


# ExternalName & Endpoint

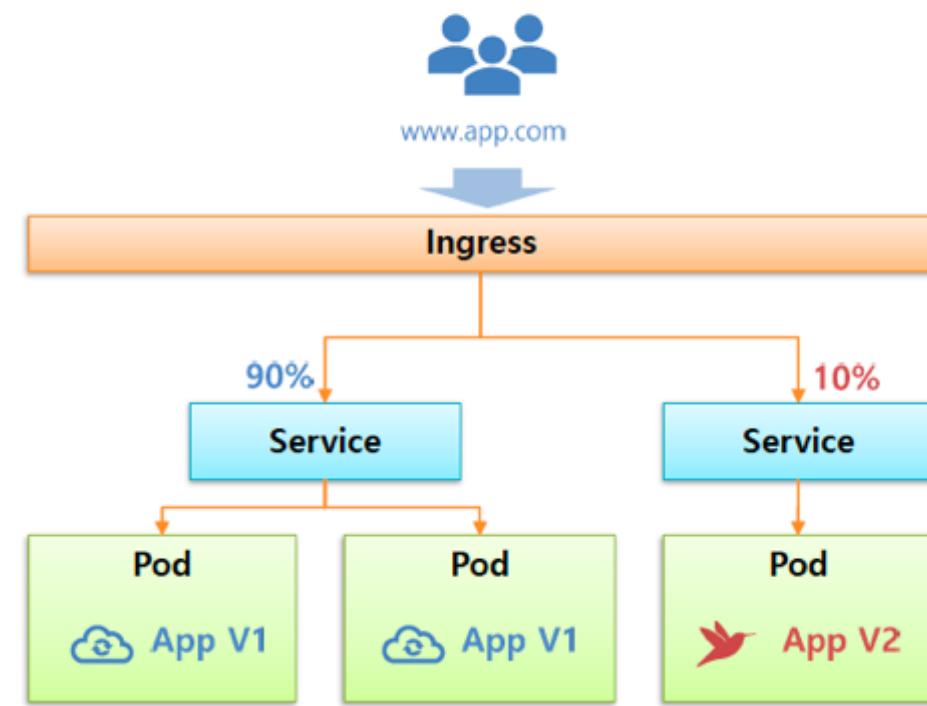


# Ingress

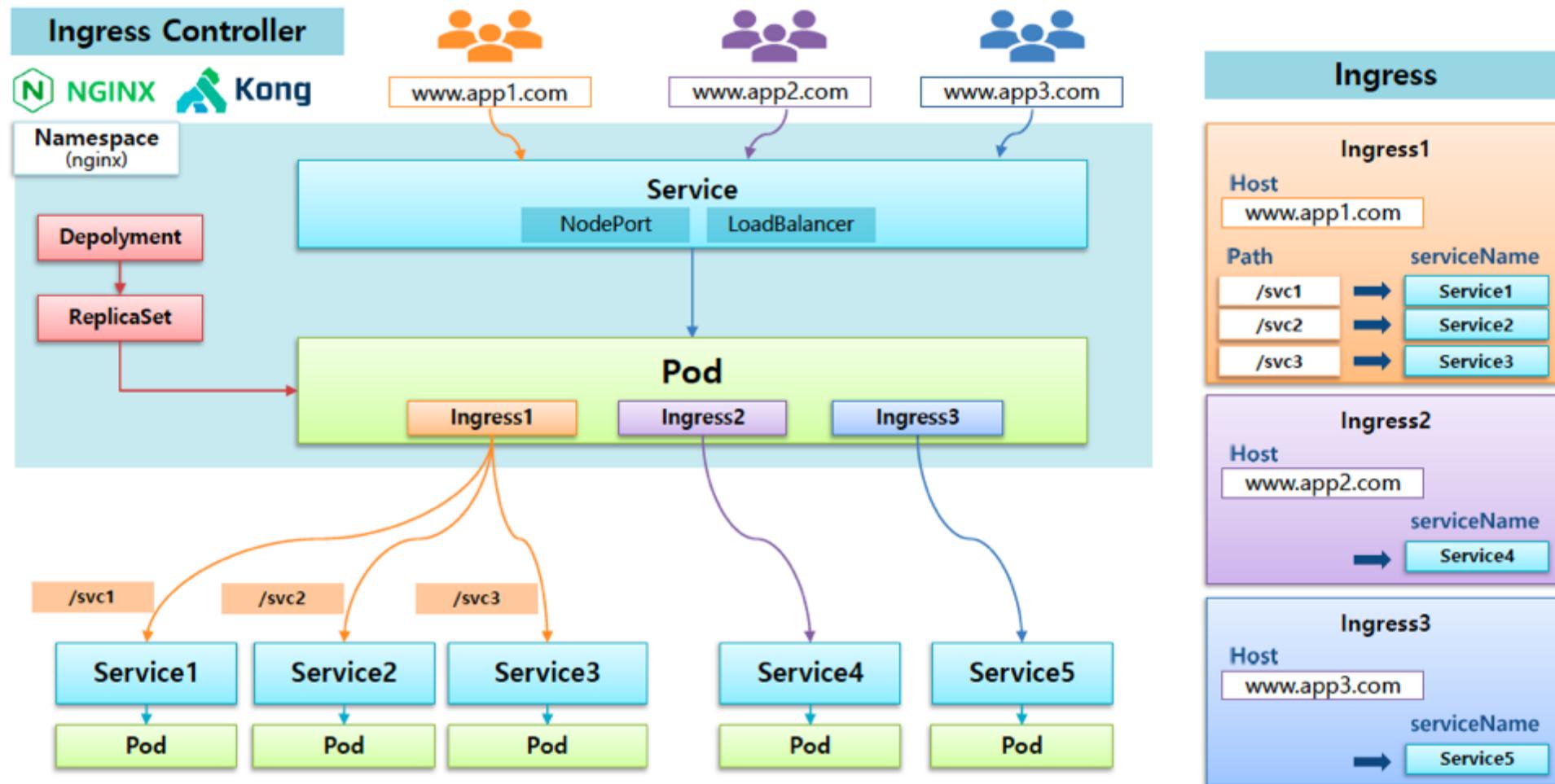
## Service LoadBalancing



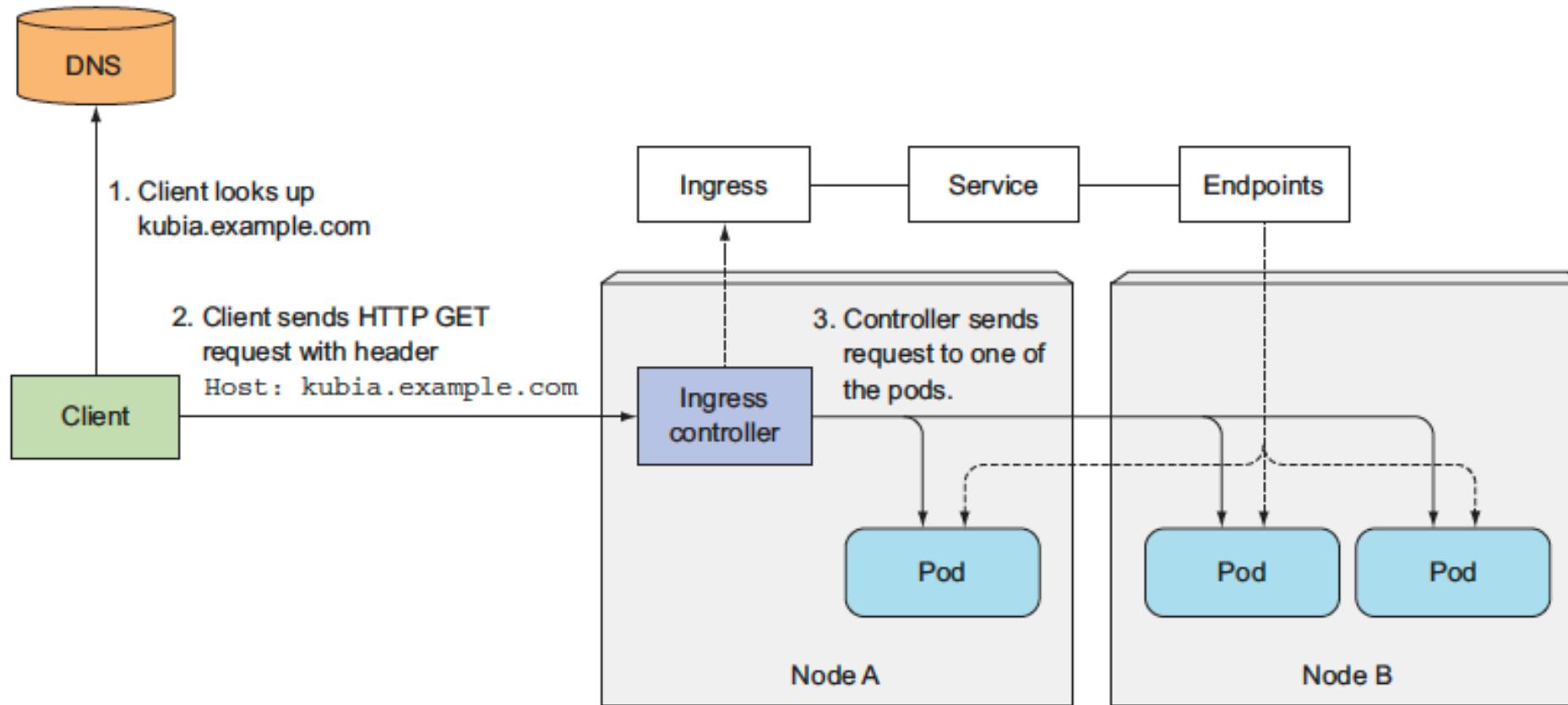
## Canary Upgrade



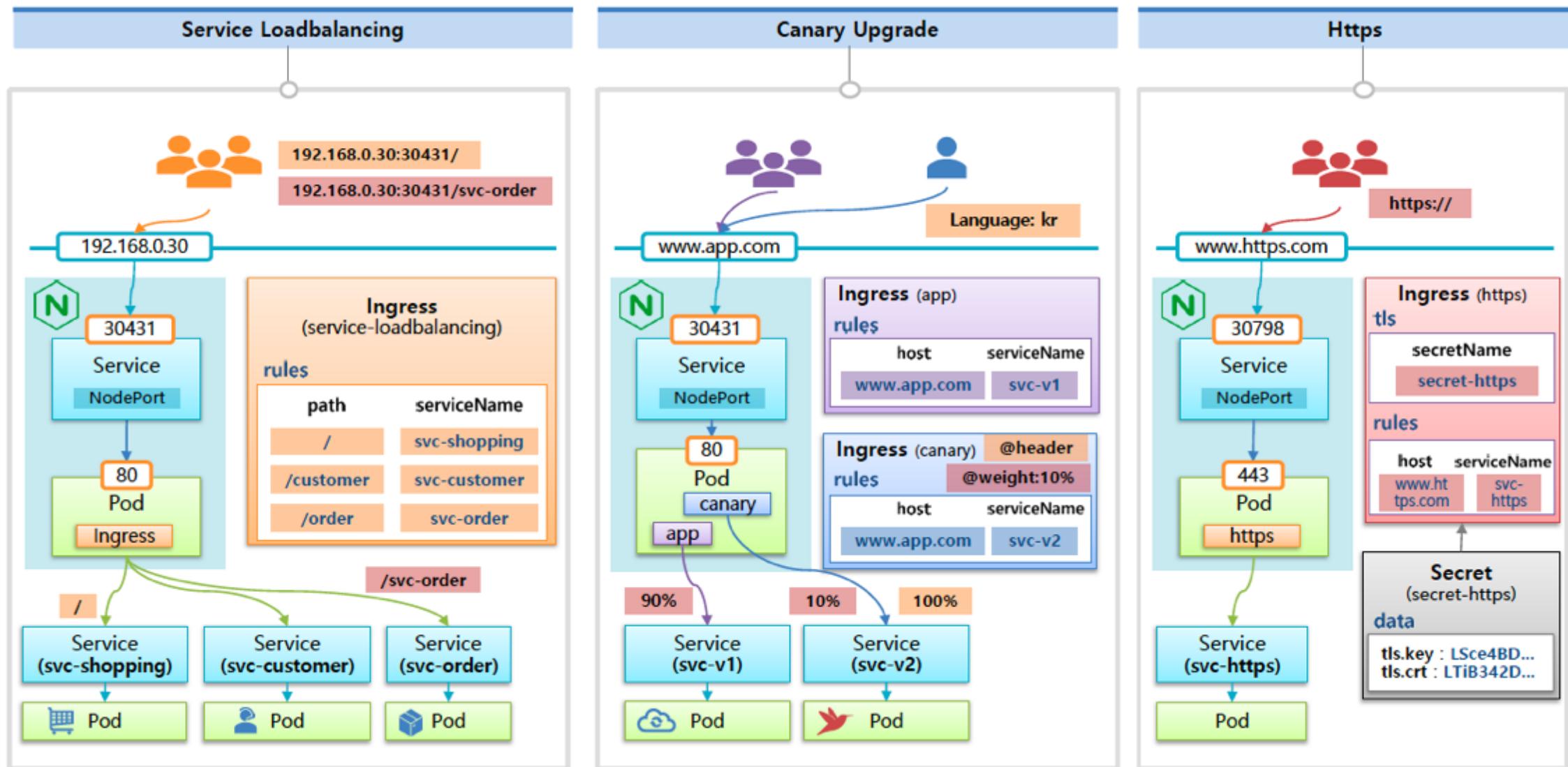
# Ingress Controller



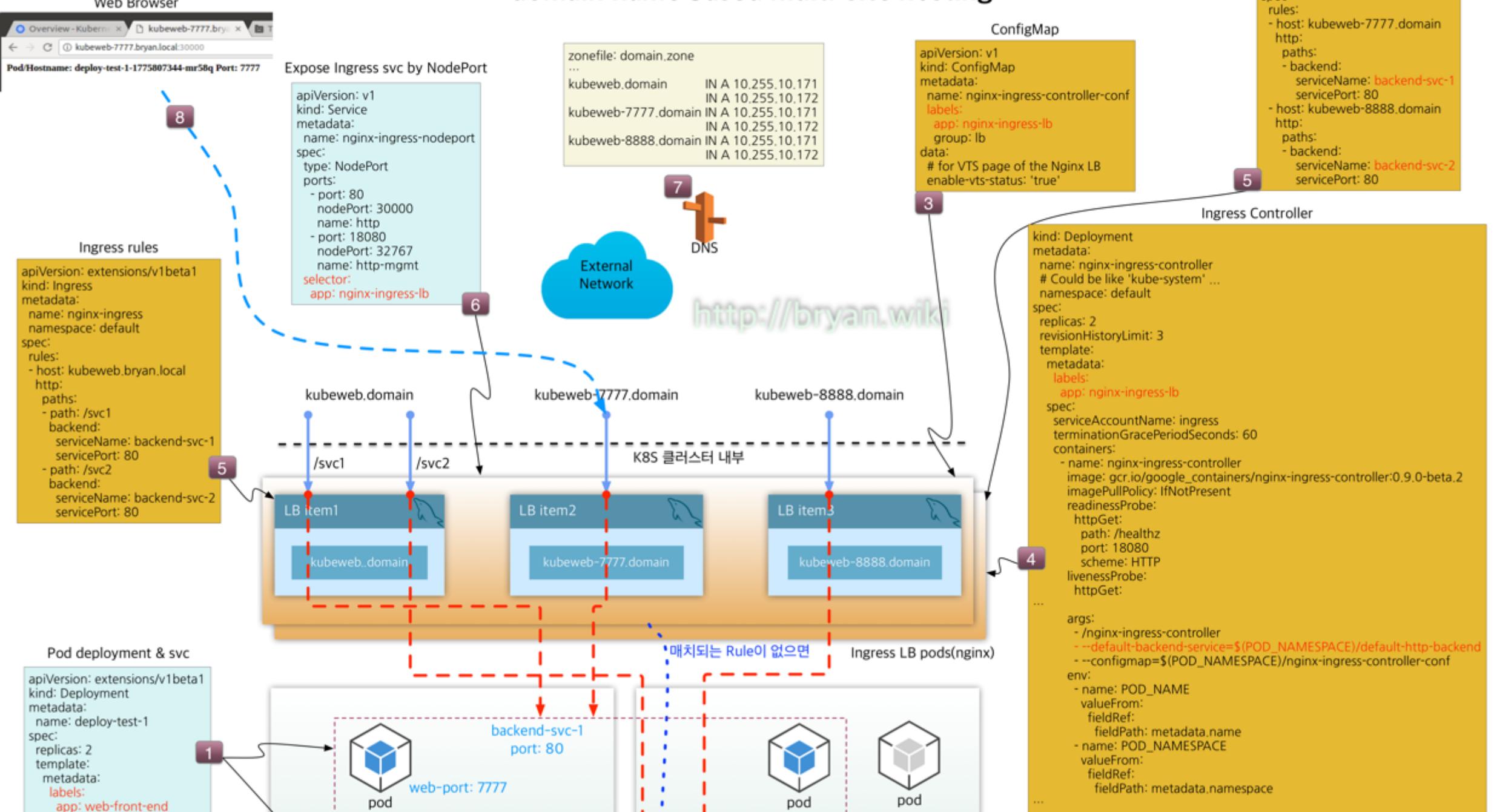
# Ingress, Ingress Controller, Endpoints



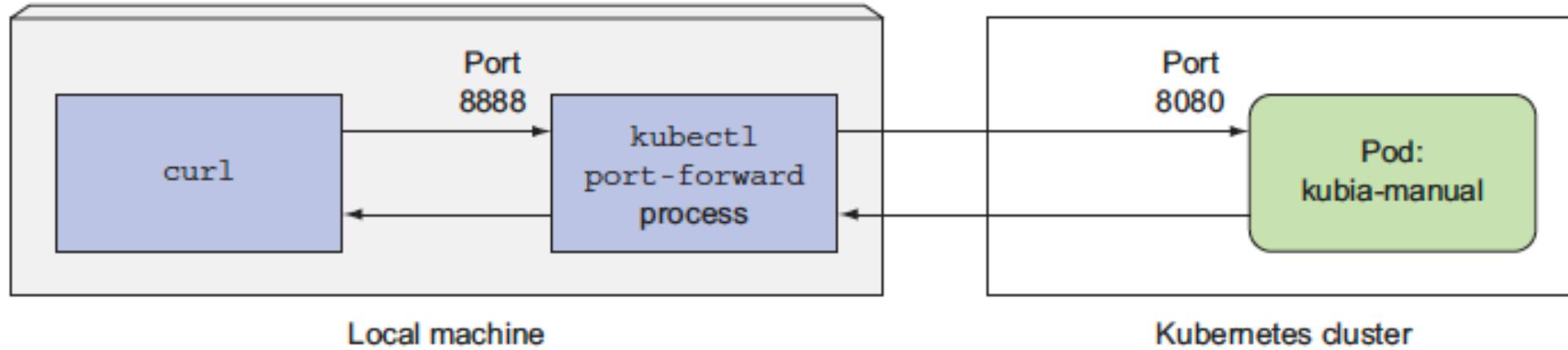
# Ingress Examples



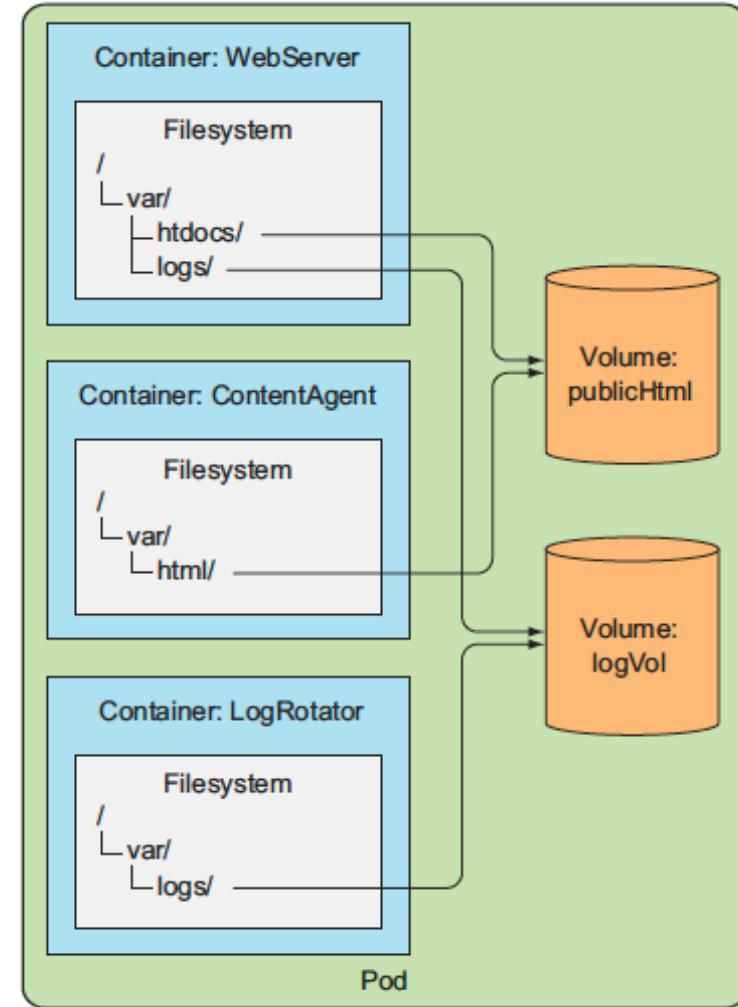
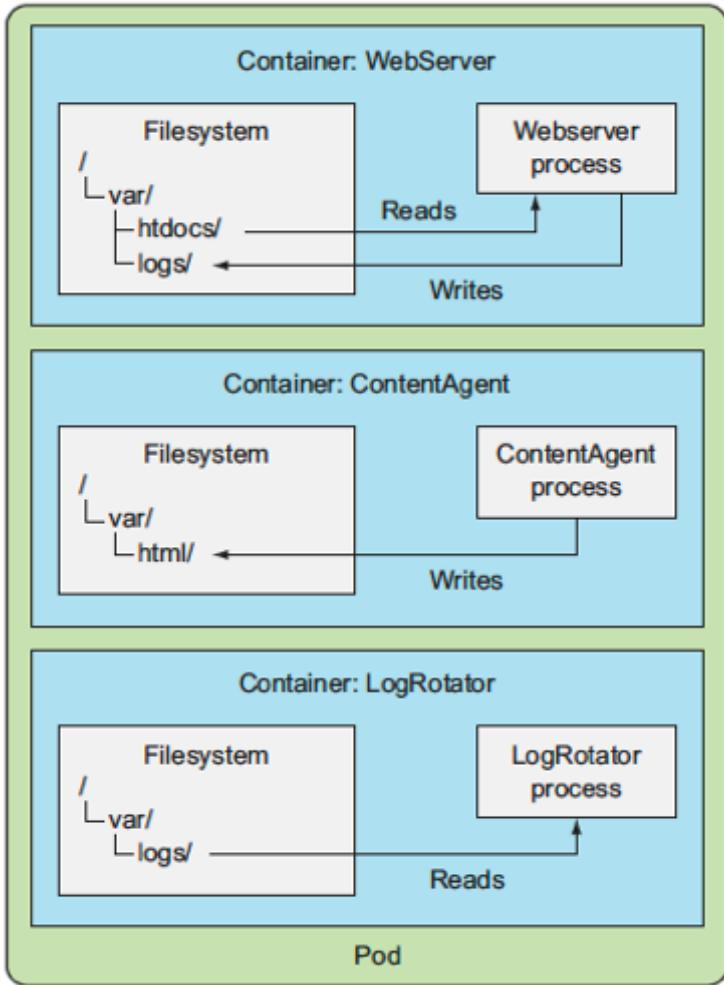
## Ingress Load Balancer 개념도(using nginx container) - domain name based multi-site hosting -



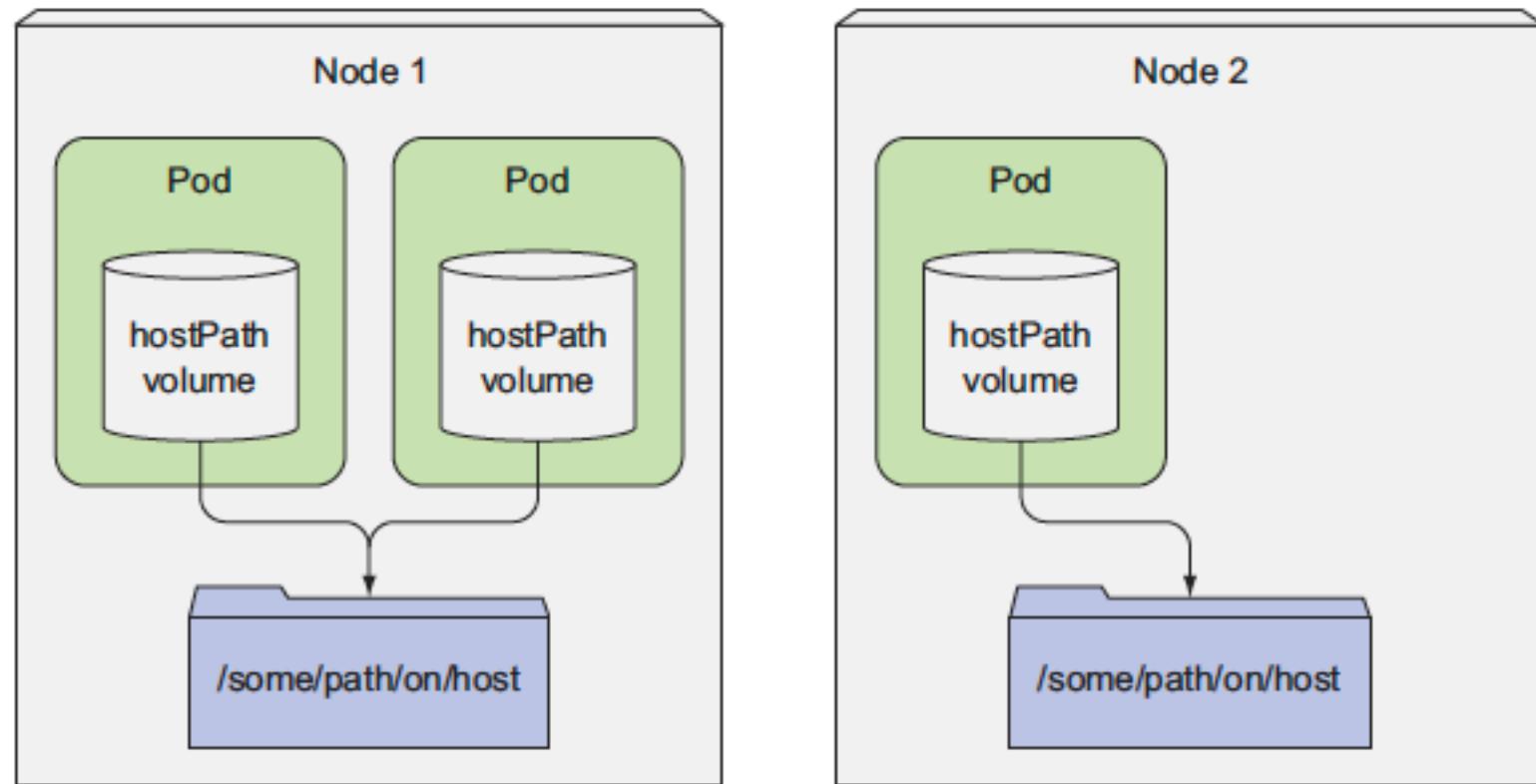
# Port Forward



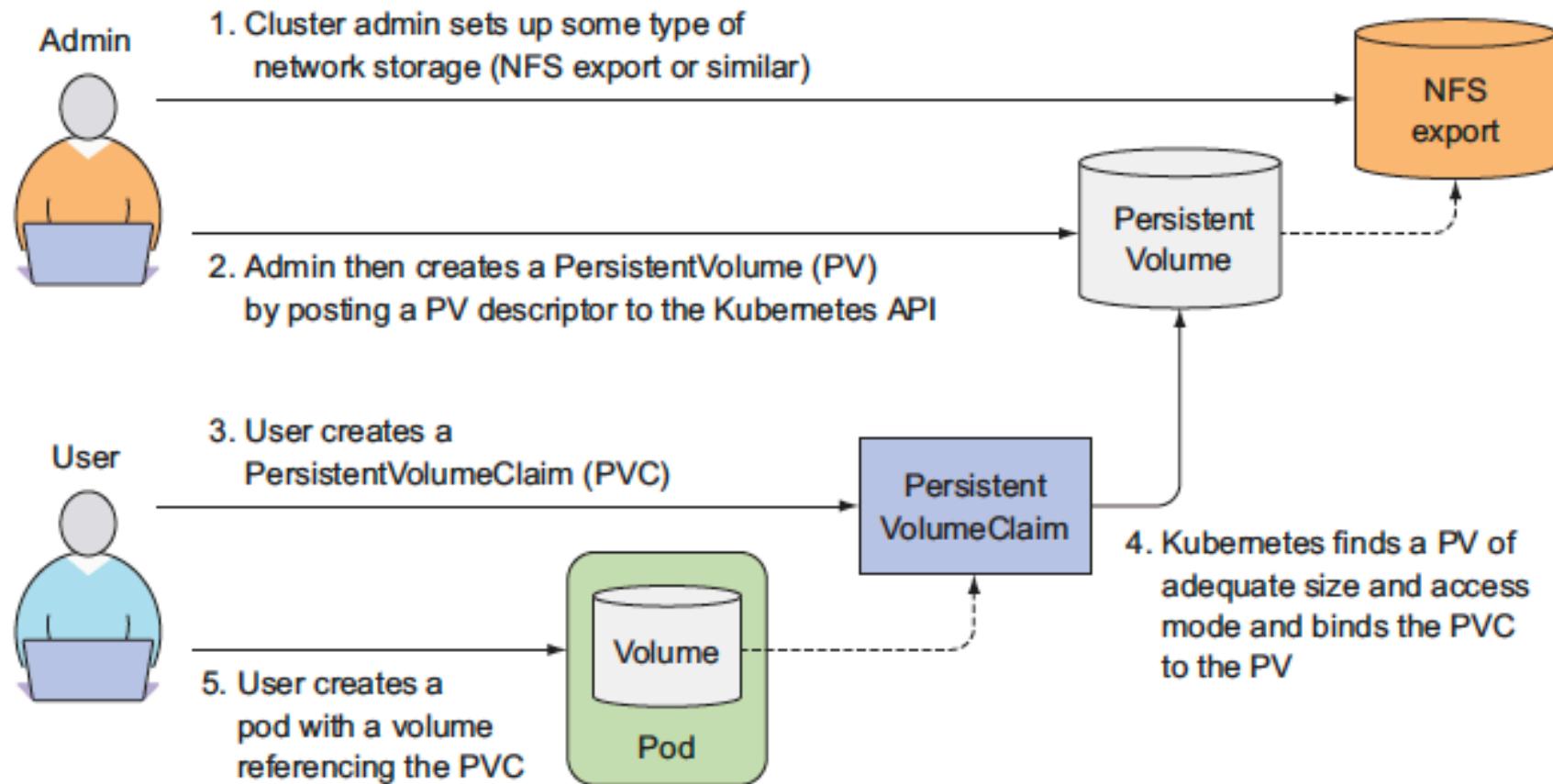
# Volume



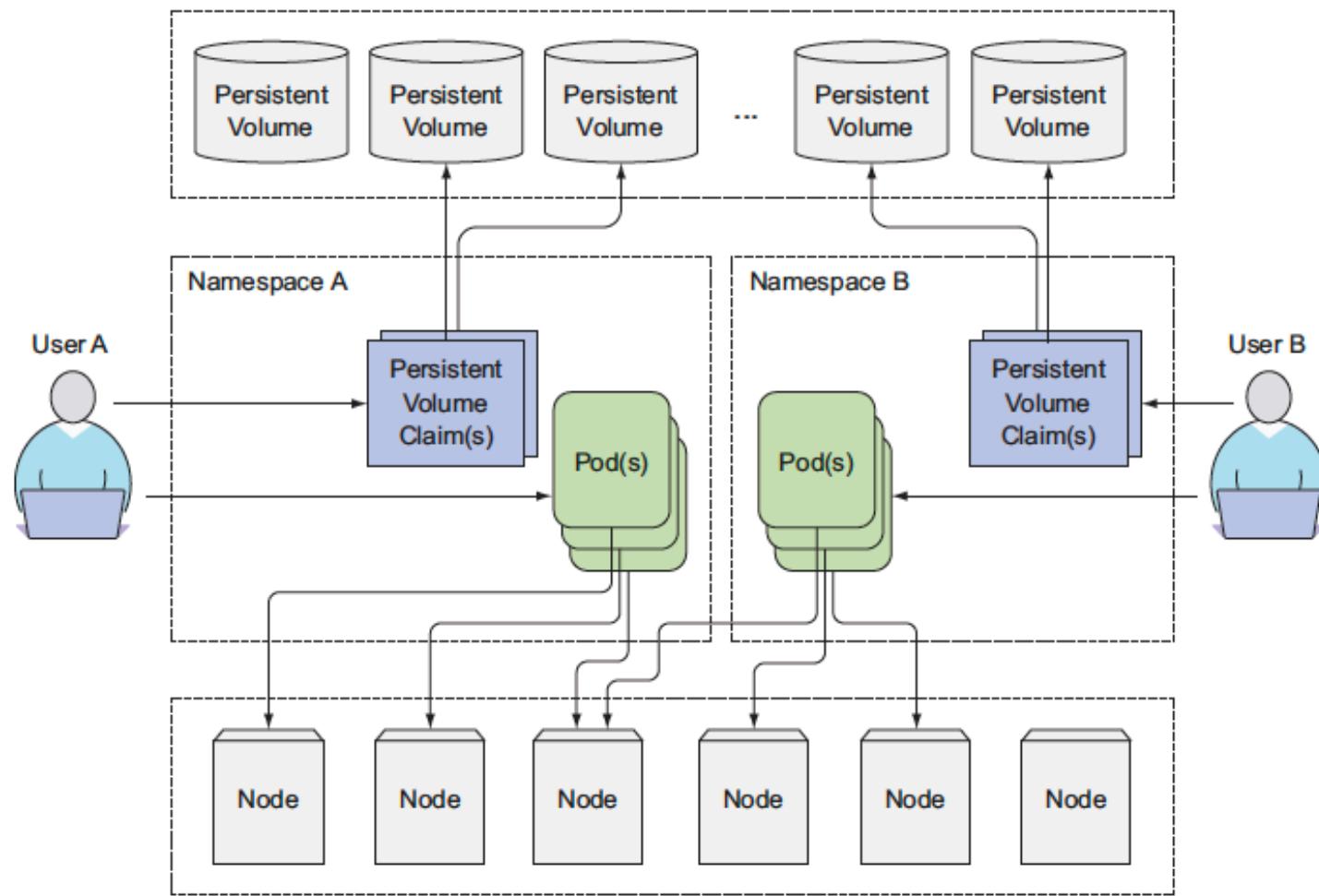
# hostPath



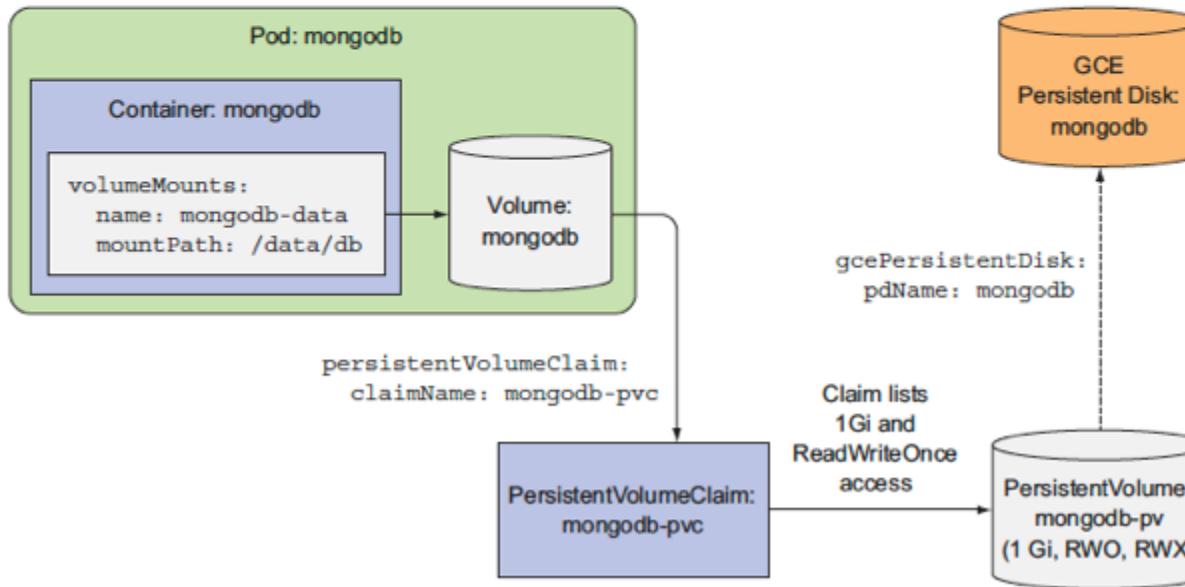
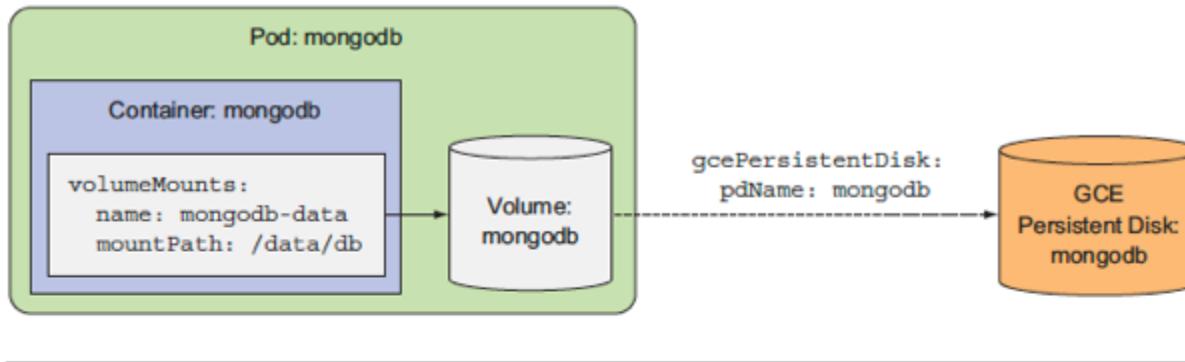
# PersistentVolume & PersistentVolumeClaim



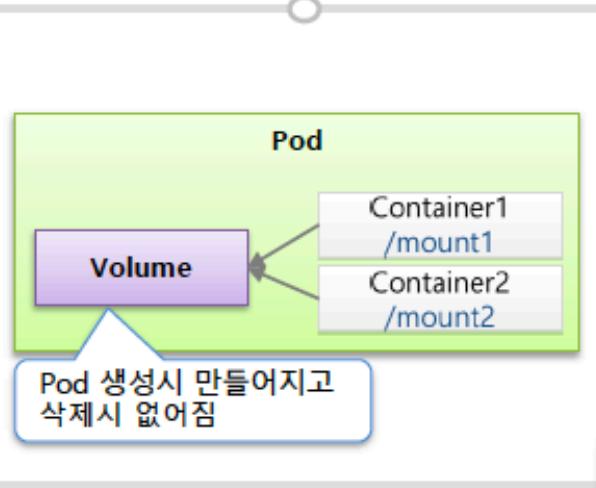
# PersistentVolume, Node, Namespace



# Use Directly vs Through PV and PVC



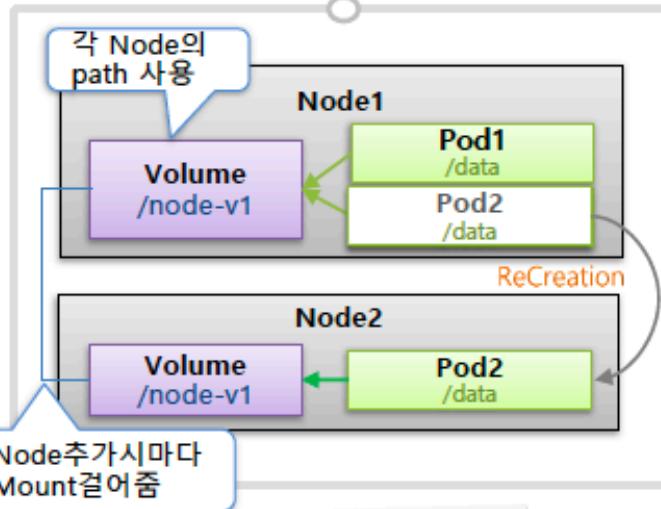
### emptyDir



Pod 생성시 만들어지고  
삭제시 없어짐

```
apiVersion: v1
kind: Pod
metadata:
  name: pod-volume-1
spec:
  containers:
    - name: container1
      image: tmkube/init
    volumes:
      - name: empty-dir
        emptyDir: {}
```

### hostPath

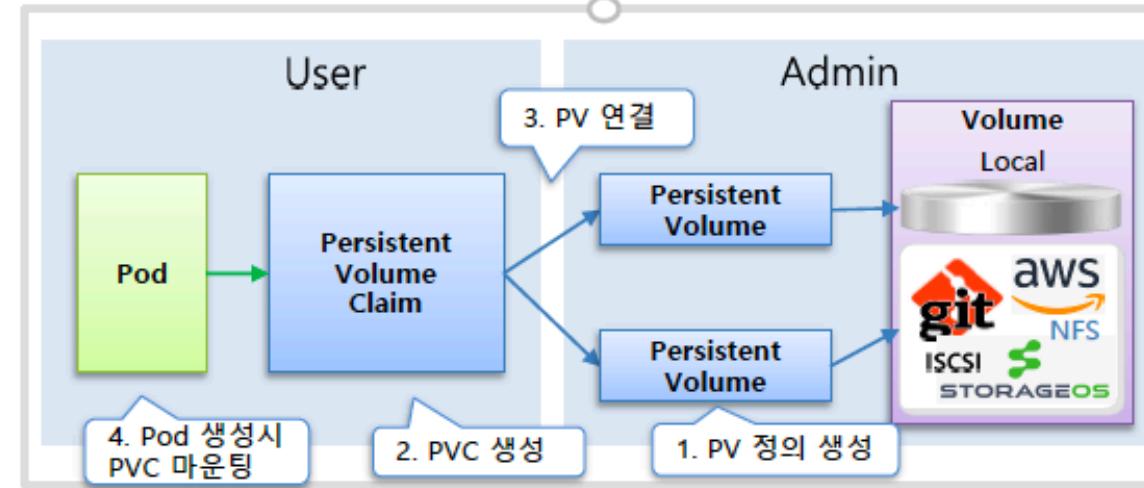


각 Node의  
path 사용  
Node추가시마다  
Mount 걸어줌

```
apiVersion: v1
kind: Pod
metadata:
  name: pod-volume-2
spec:
  containers:
    - name: container
      image: tmkube/init
    volumes:
      - name: host-path
        hostPath:
          path: /node-v
          type: Directory
```

사전에 해당 Node에  
경로가 있어야 함

### PVC / PV



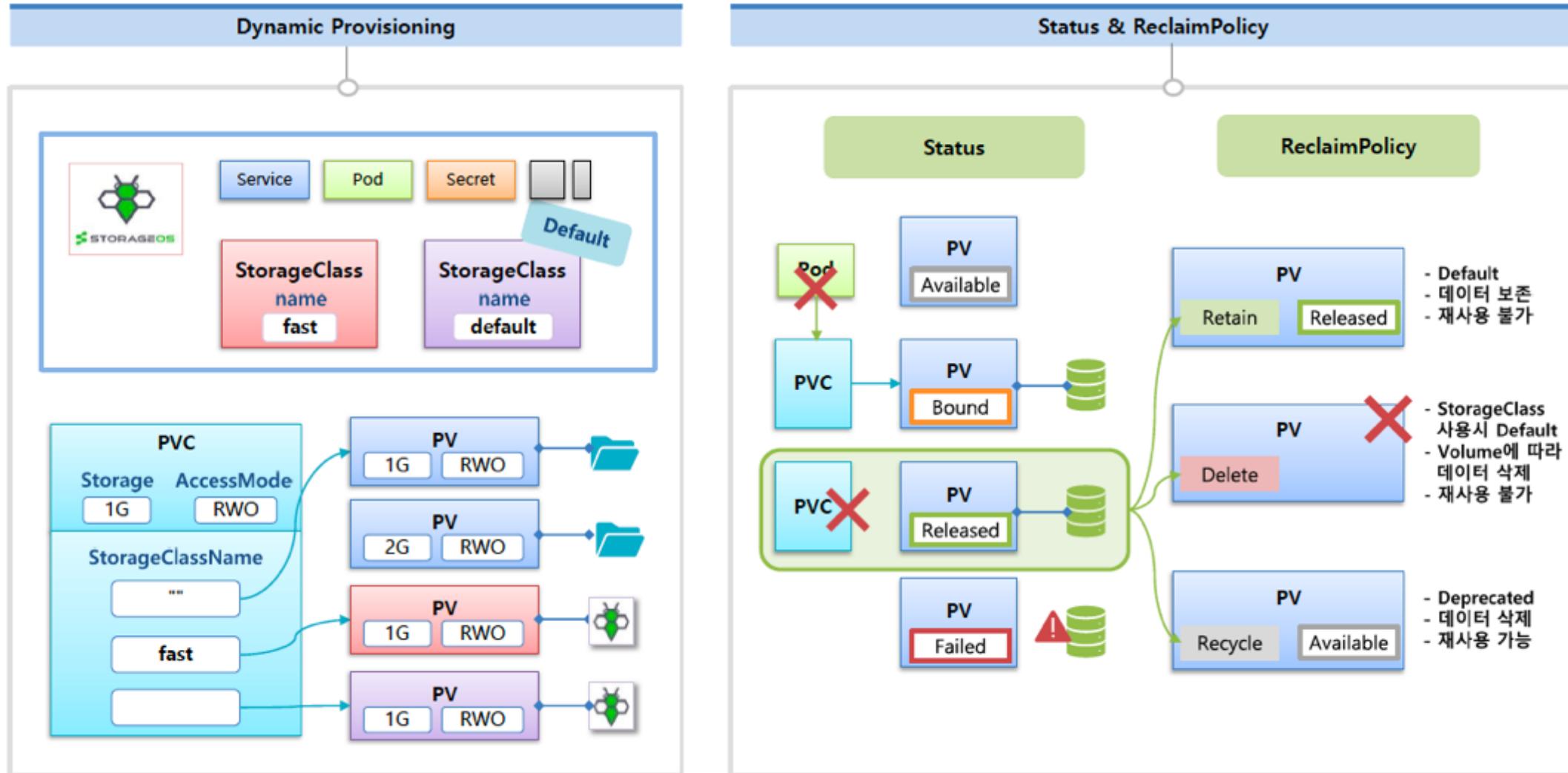
```
apiVersion: v1
kind: Pod
metadata:
  name: pod-volume-3
spec:
  containers:
    - name: container
      image: tmkube/init
    volumes:
      - name: pvc-pv
        persistentVolumeClaim:
          claimName: pvc-01
```

```
apiVersion: v1
kind: PersistentVolumeClaim
metadata:
  name: pvc-01
spec:
  accessModes:
    - ReadWriteOnce
  resources:
    requests:
      storage: 1G
  storageClassName: ""
```

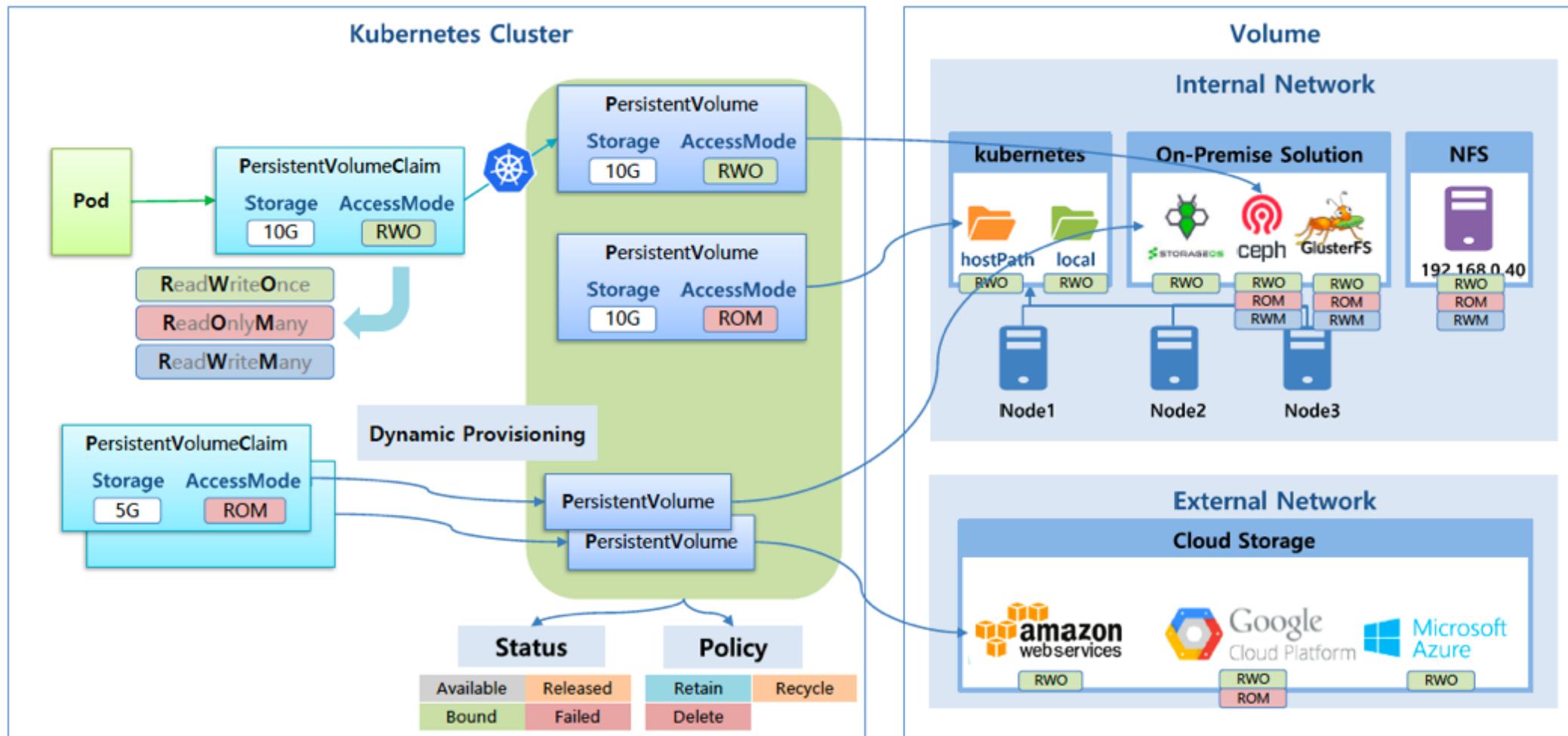
```
apiVersion: v1
kind: PersistentVolume
metadata:
  name: pv-01
spec:
  capacity:
    storage: 1G
  accessModes:
    - ReadWriteOnce
  local:
    path: /node-v
  nodeAffinity:
    required:
      nodeSelectorTerms:
        - matchExpressions:
          - {key: node, operator: In, values: [node1]}
```

Pod가 해당  
Node에 만들어짐

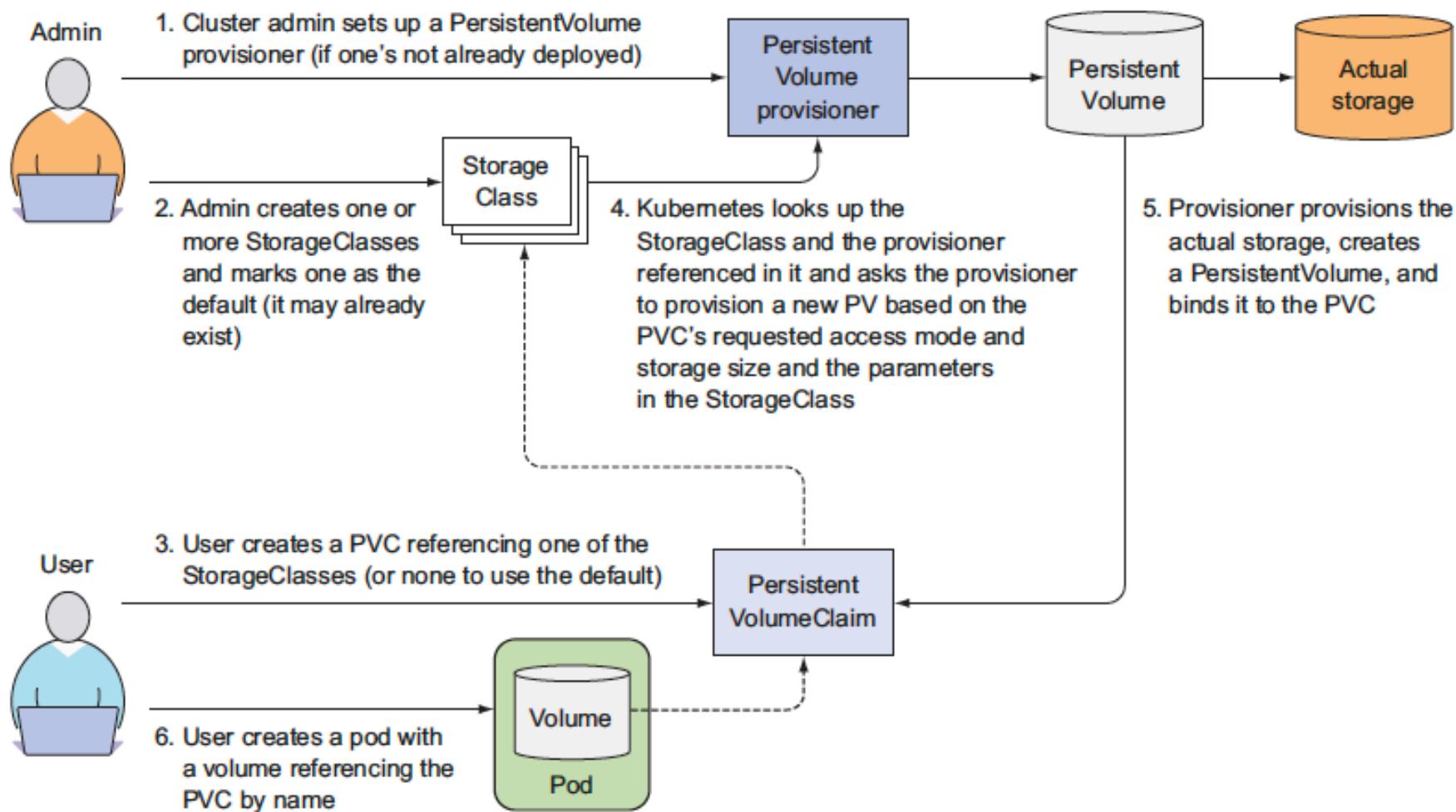
# Dynamic Provisioning

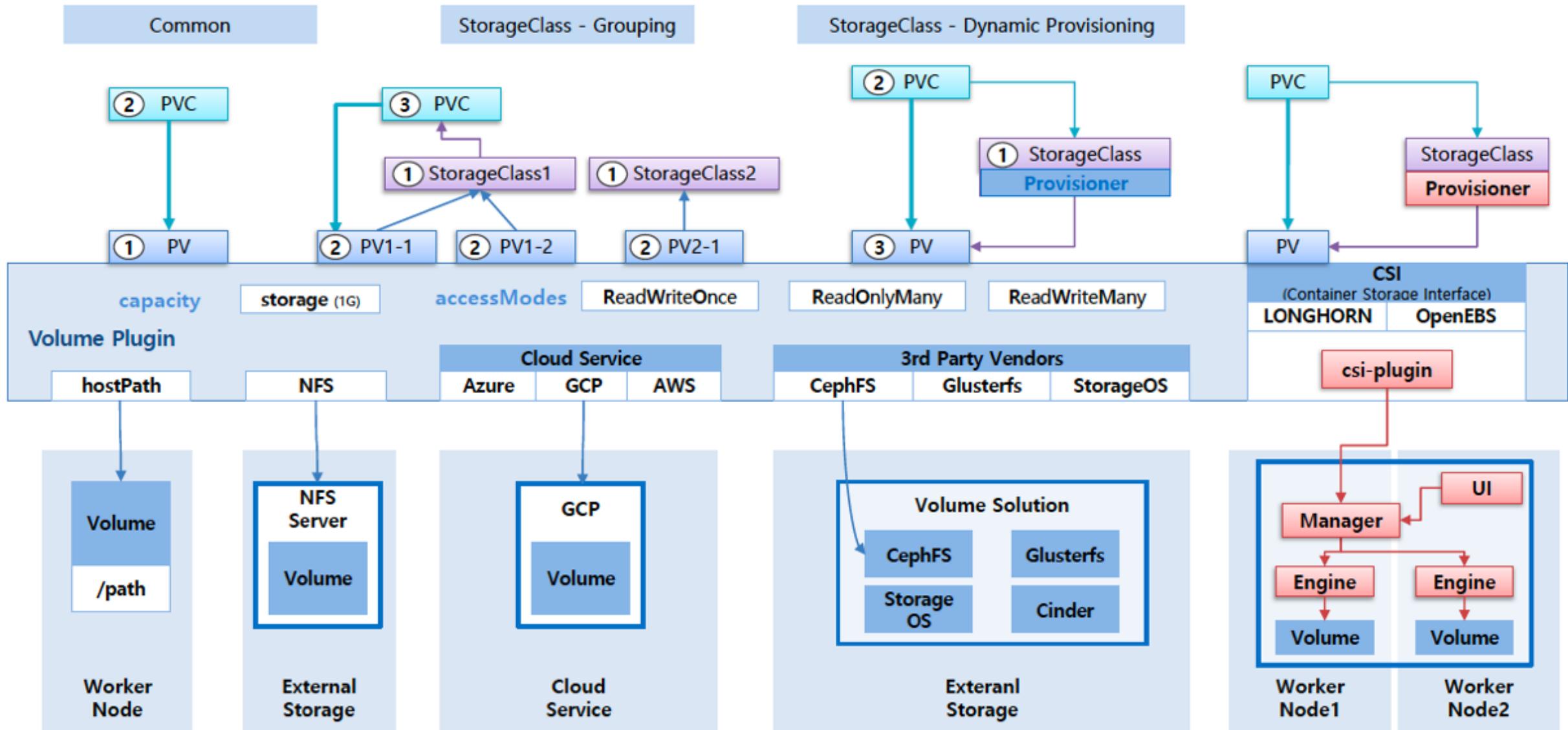


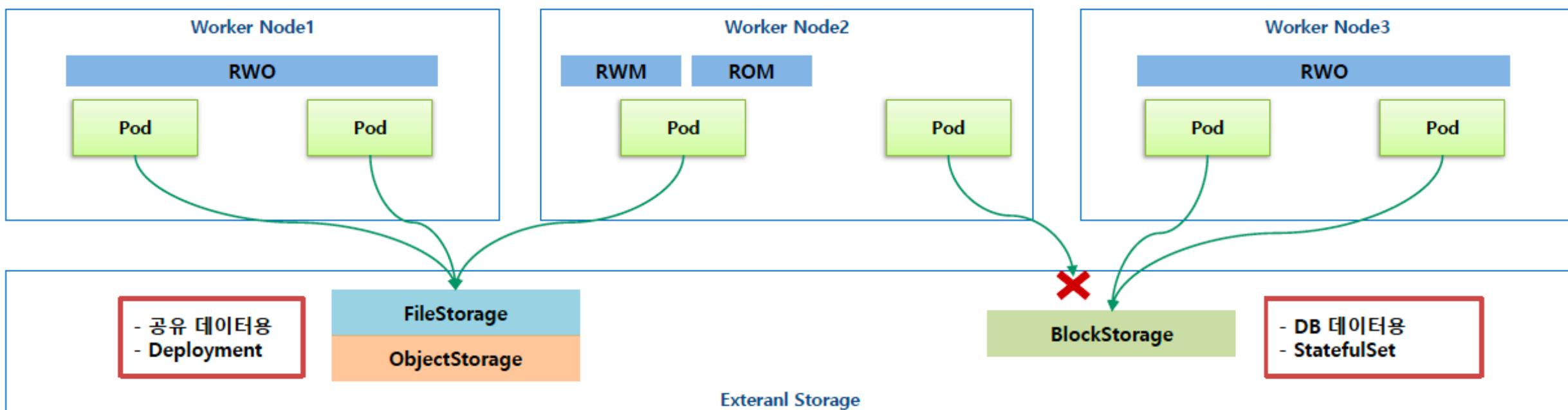
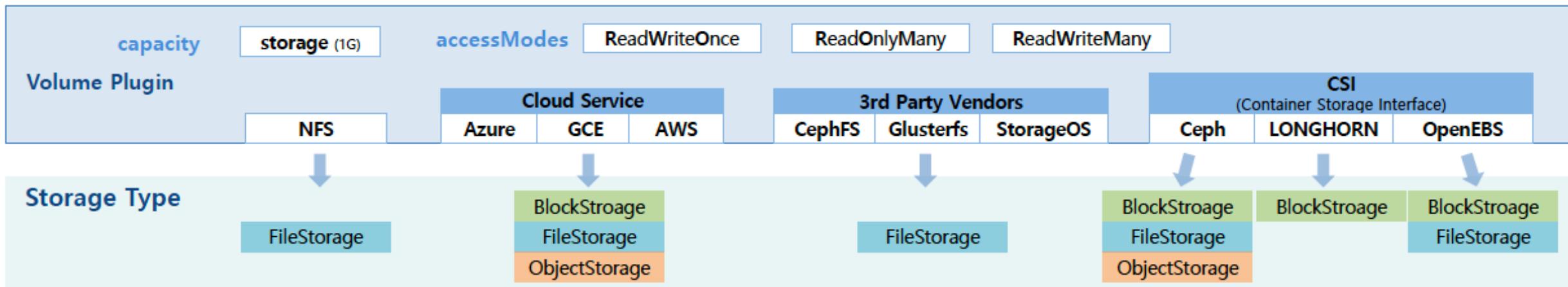
# PV, PVC: Complete Picture

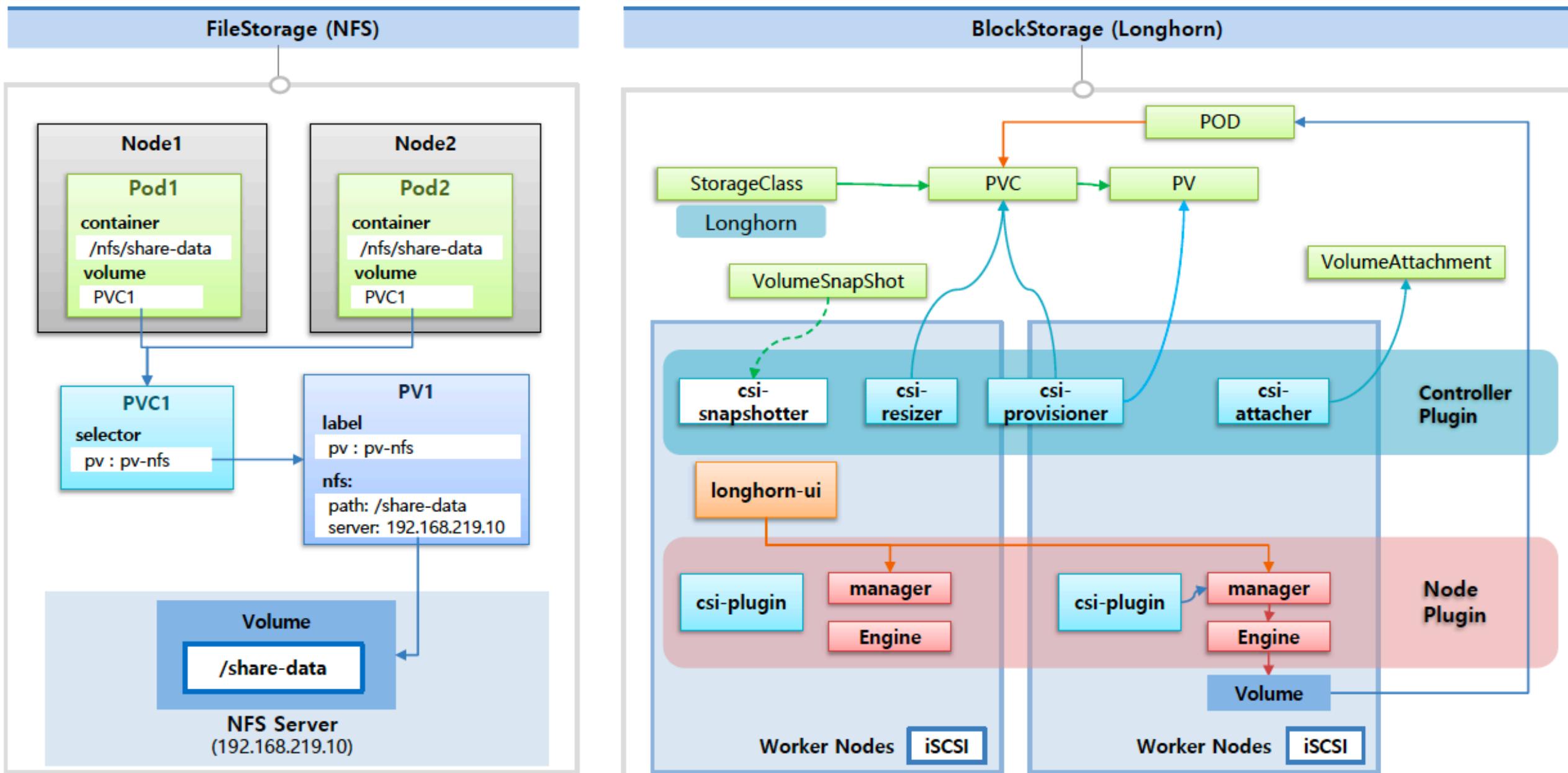


# PV, PVC: Complete Steps

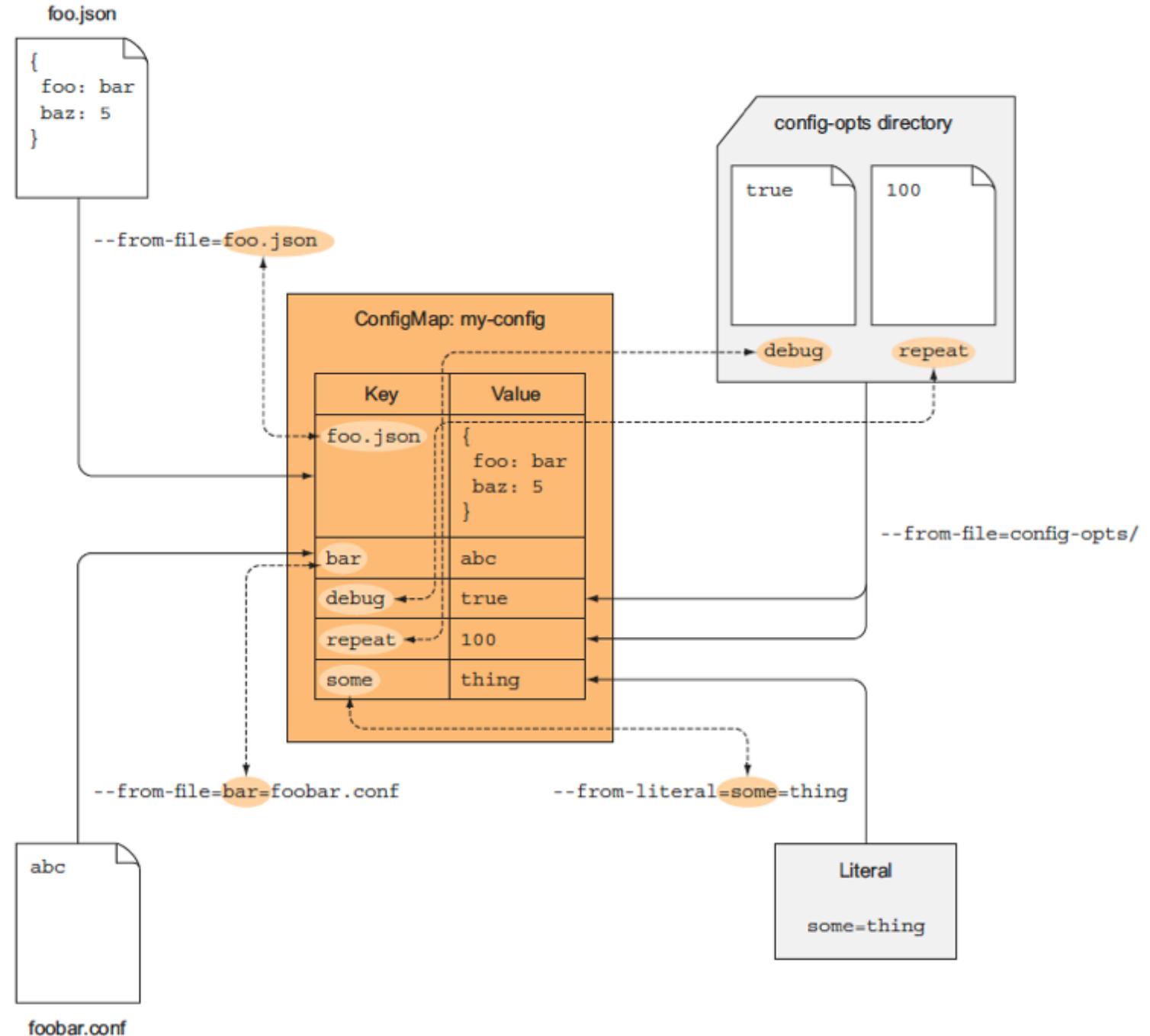




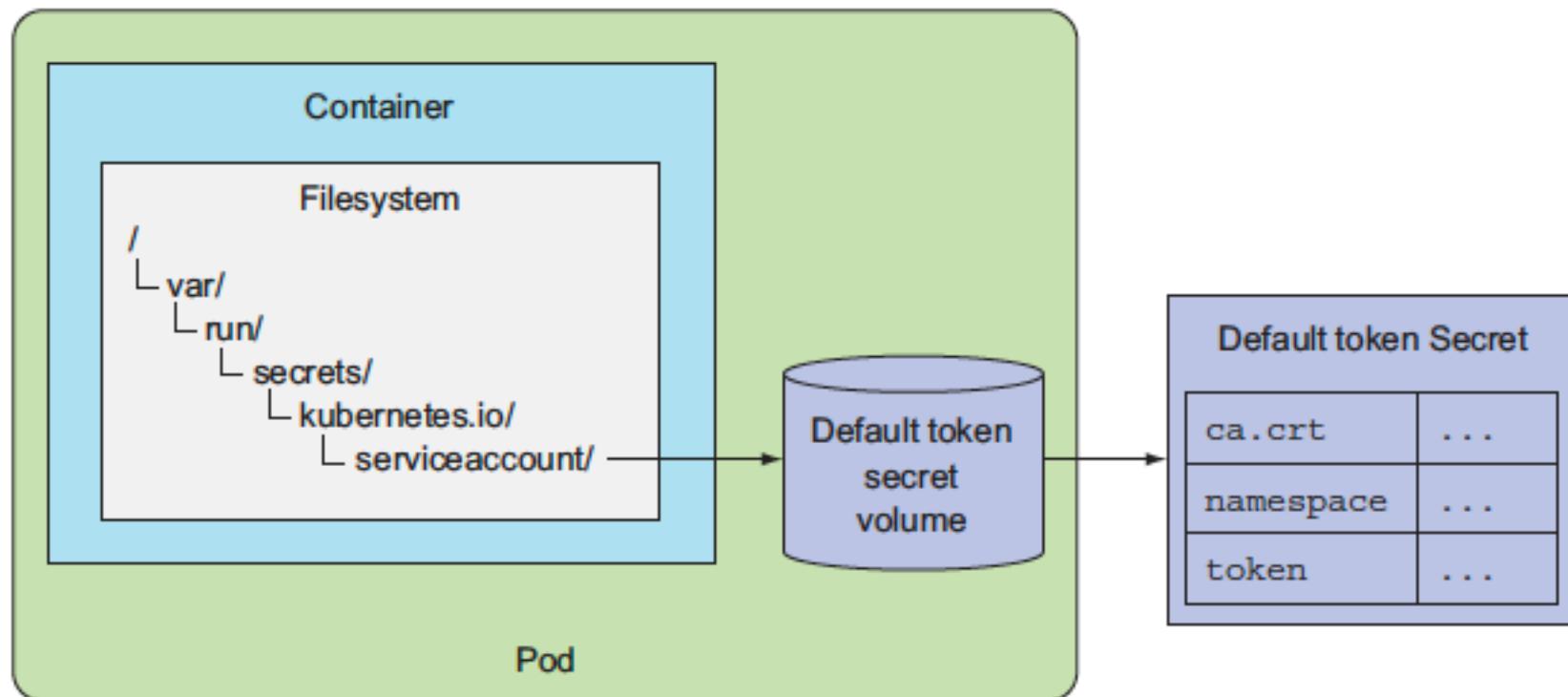




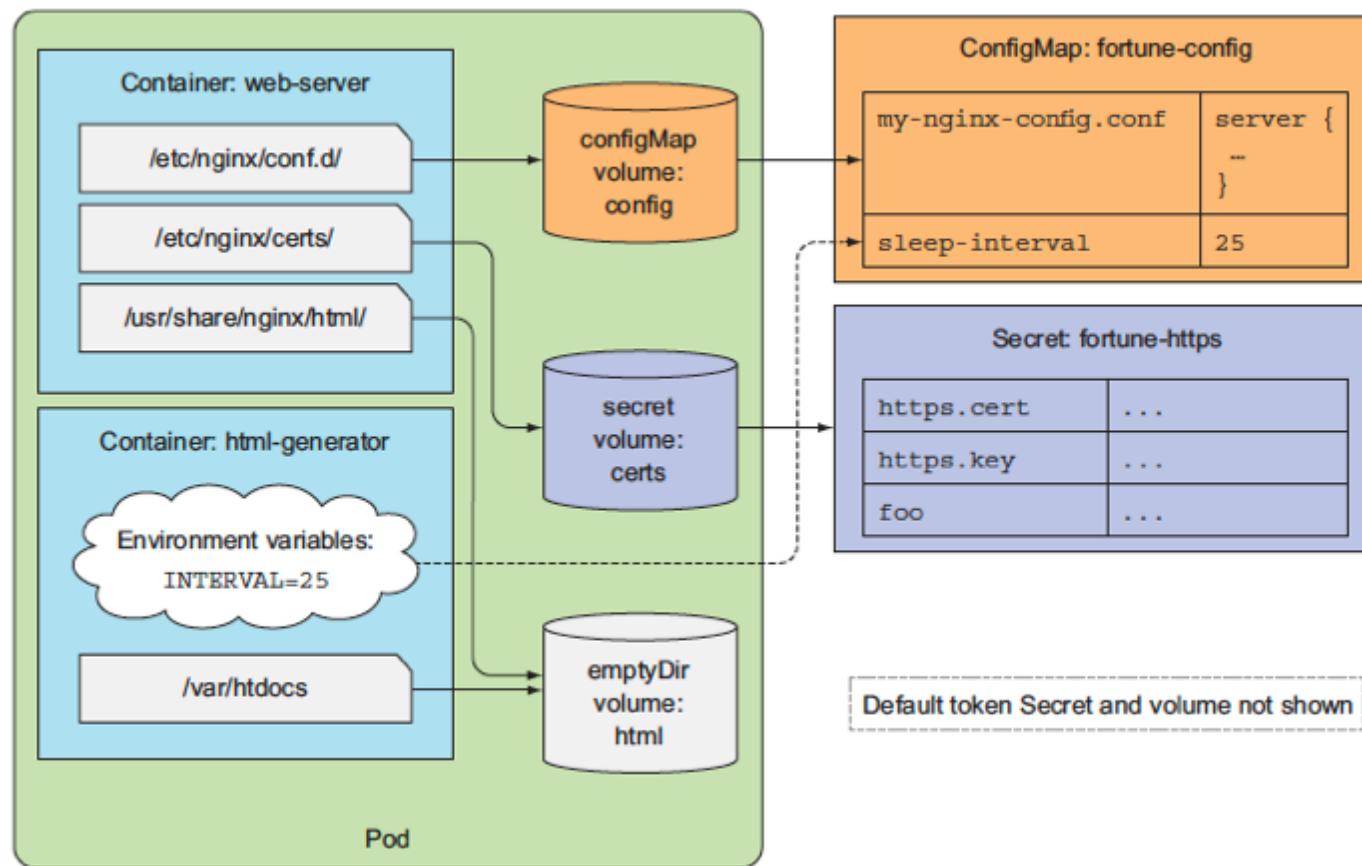
# ConfigMap



# Secret (Default Token)

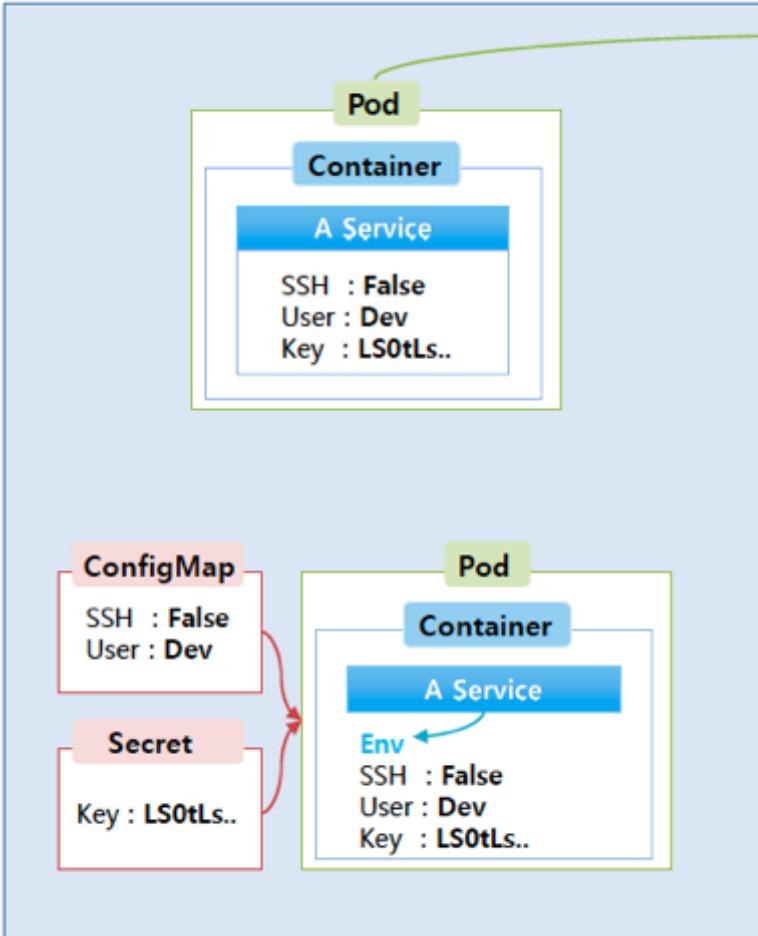


# ConfigMap & Secret

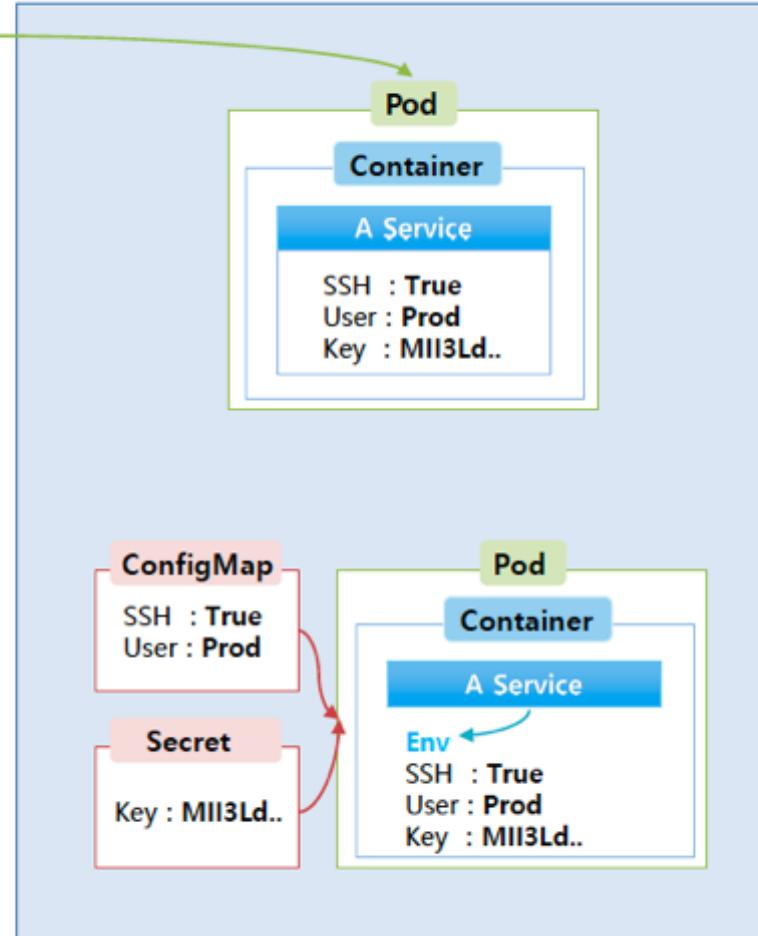


# ConfigMap & Secret

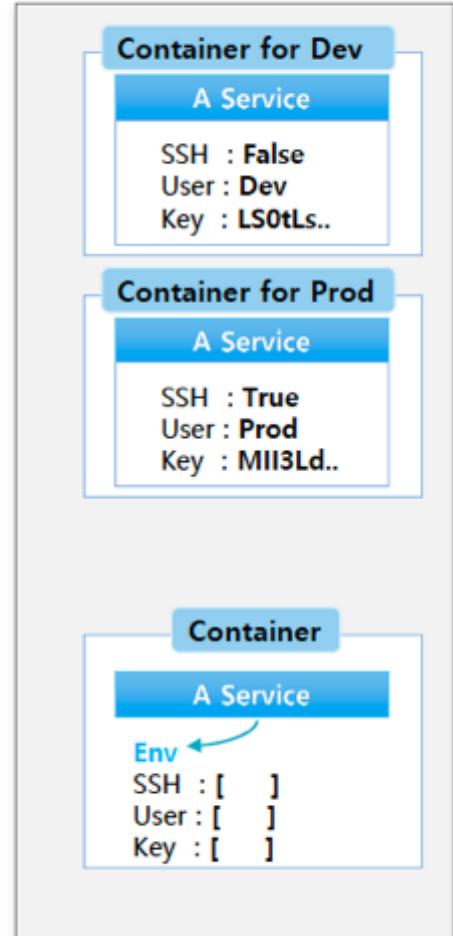
 Dev

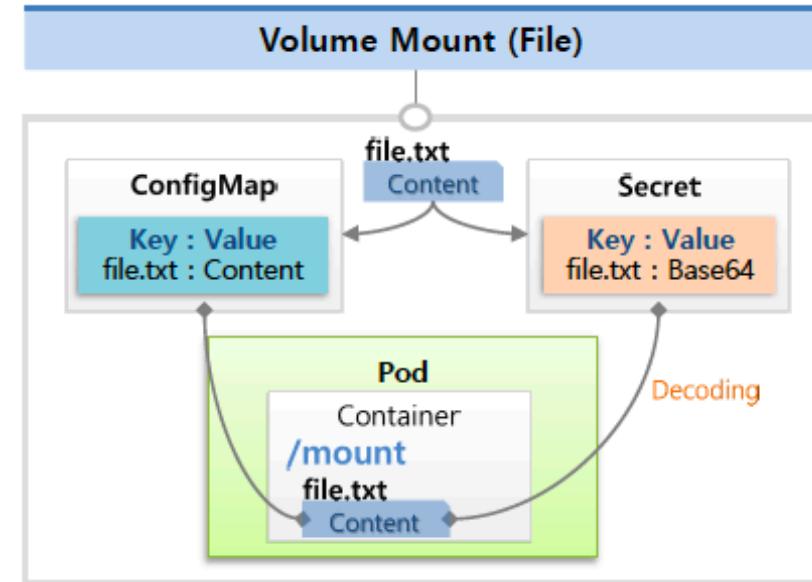
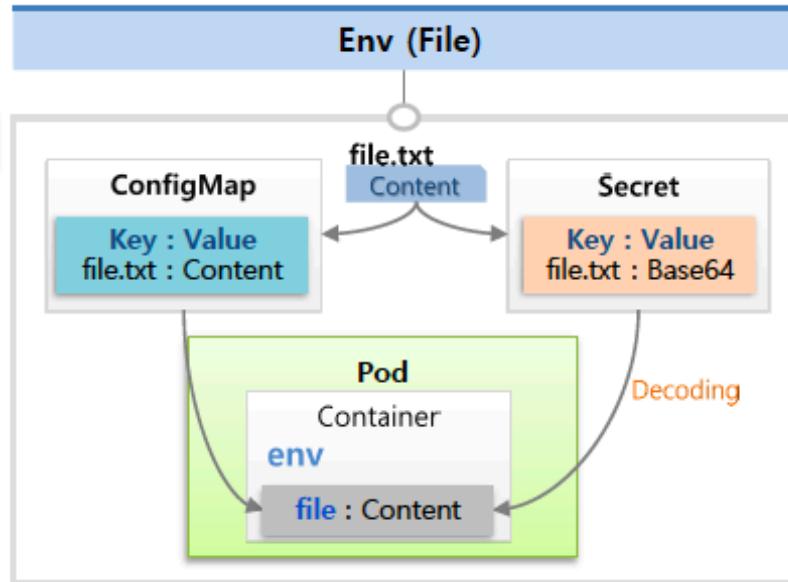
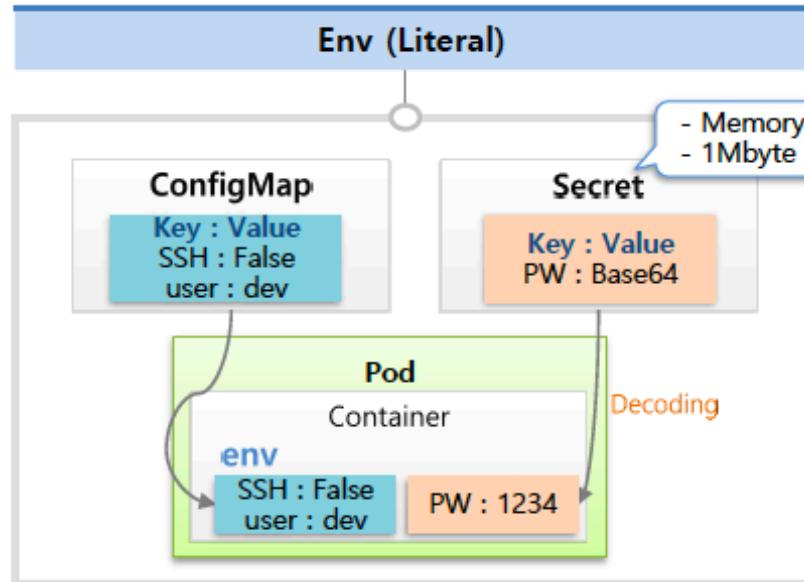


 Production



 Image





```
apiVersion: v1
kind: ConfigMap
metadata:
  name: cm-dev
data:
  SSH: False
  User: dev
```

```
apiVersion: v1
kind: Secret
metadata:
  name: sec-dev
data:
  Key: MTIzNA==
```

```
apiVersion: v1
kind: Pod
metadata:
  name: pod-1
spec:
  containers:
    - name: container
      image: tmkube/init
      envFrom:
        - configMapRef:
            name: cm-dev
        - secretRef:
            name: sec-dev
```

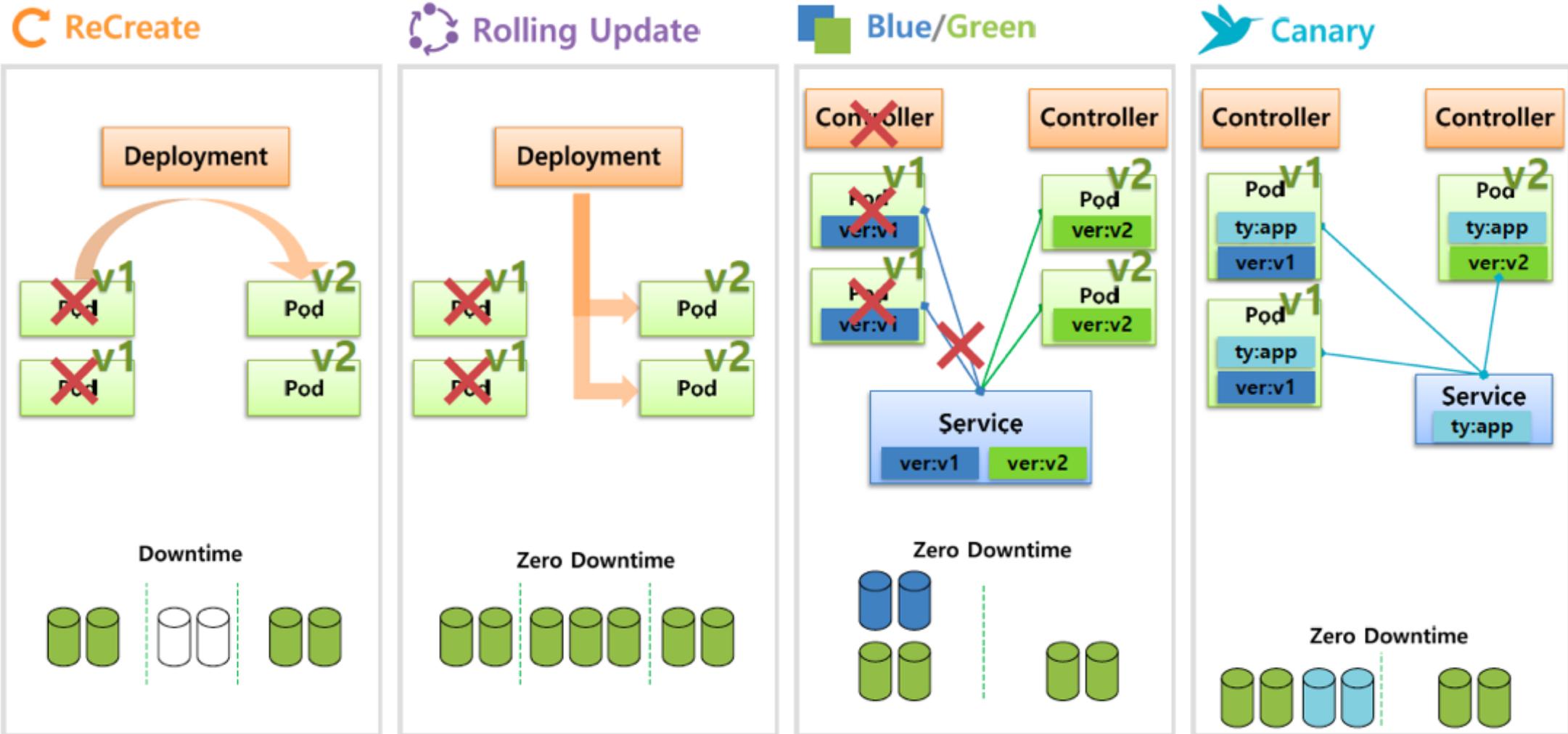
kubectl create configmap cm-file --from-file=./file.txt

kubectl create secret generic sec-file --from-file=./file.txt

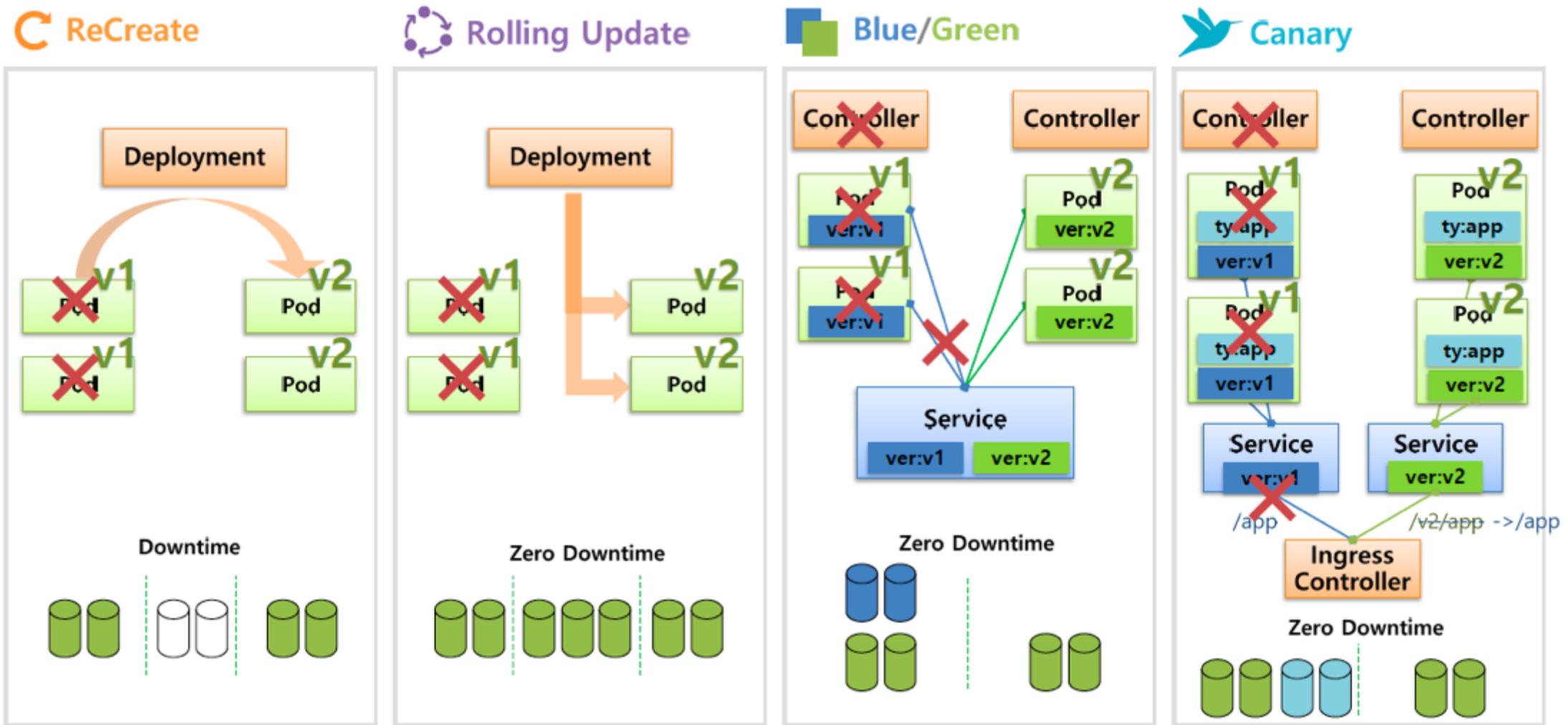
```
apiVersion: v1
kind: Pod
metadata:
  name: file
spec:
  containers:
    - name: container
      image: tmkube/init
      env:
        - name: file
          valueFrom:
            configMapKeyRef:
              name: cm-file
              key: file.txt
```

```
apiVersion: v1
kind: Pod
metadata:
  name: mount
spec:
  containers:
    - name: container
      image: tmkube/init
      volumeMounts:
        - name: file-volume
          mountPath: /mount
  volumes:
    - name: file-volume
      configMap:
        name: cm-file
```

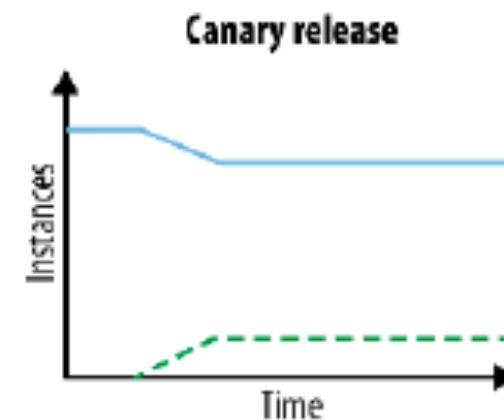
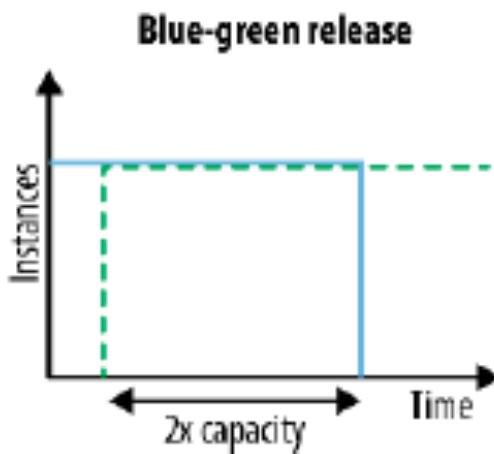
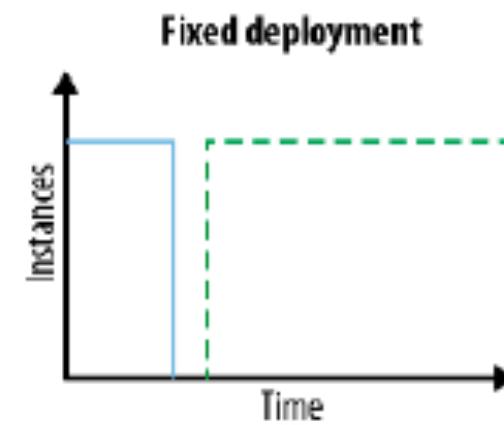
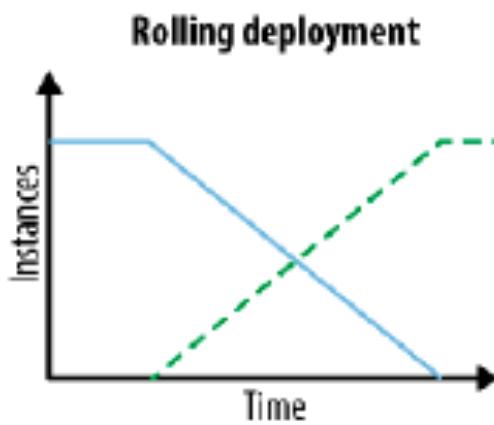
# Deployment



# Deployment

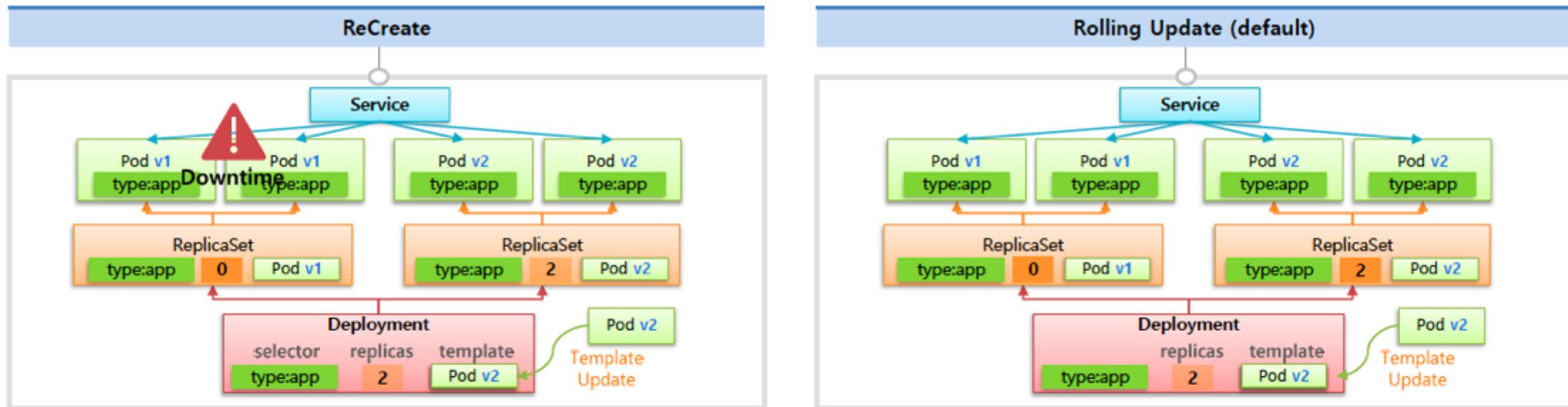


# Deployment Strategies



Strategy	ZERO DOWNTIME	REAL TRAFFIC TESTING	TARGETED USERS	CLOUD COST	ROLLBACK DURATION	NEGATIVE IMPACT ON USER	COMPLEXITY OF SETUP
<b>RECREATE</b> version A is terminated then version B is rolled out	✗	✗	✗	■ ■ ■	■ ■ ■	■ ■ ■	□ □ □
<b>RAMPED</b> version B is slowly rolled out and replacing version A	✓	✗	✗	■ ■ ■	■ ■ ■	■ ■ ■	■ ■ ■
<b>BLUE/GREEN</b> version B is released alongside version A, then the traffic is switched to version B	✓	✗	✗	■ ■ ■	□ □ □	■ ■ ■	■ ■ ■
<b>CANARY</b> version B is released to a subset of users, then proceed to a full rollout	✓	✓	✗	■ ■ ■	■ ■ ■	■ ■ ■	■ ■ ■
<b>A/B TESTING</b> version B is released to a subset of users under specific condition	✓	✓	✓	■ ■ ■	■ ■ ■	■ ■ ■	■ ■ ■
<b>SHADOW</b> version B receives real world traffic alongside version A and doesn't impact the response	✓	✓	✗	■ ■ ■	□ □ □	□ □ □	■ ■ ■

# ReCreate, Rolling Update: YAML



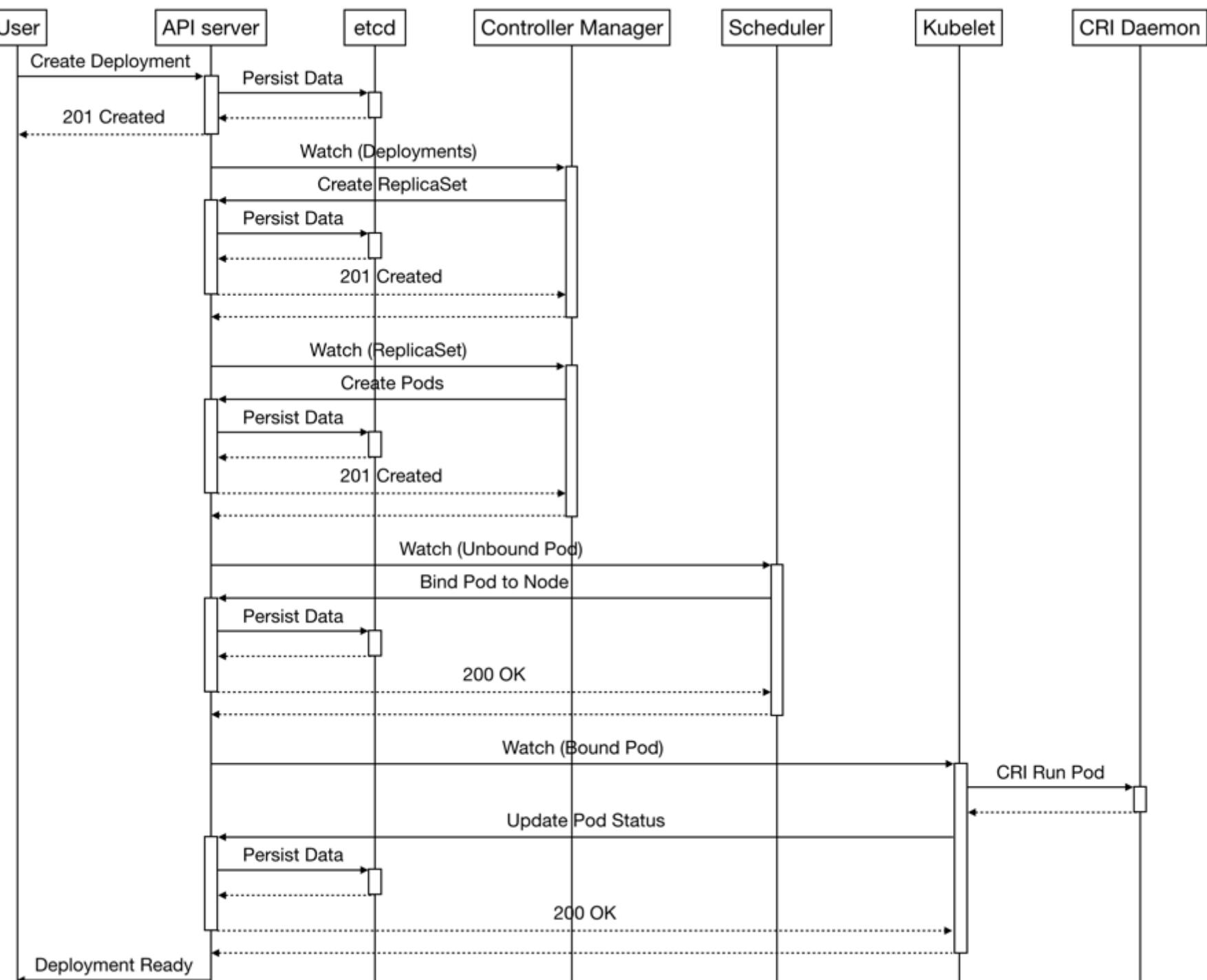
```
apiVersion: v1
kind: Service
metadata:
  name: svc-1
spec:
  selector:
    type: app
  ports:
    - port: 8080
      protocol: TCP
      targetPort: 8080
```

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: deployment-1
spec:
  selector:
    matchLabels:
      type: app
  replicas: 2
  strategy:
    type: Recreate
  revisionHistoryLimit: 1
  template:
    metadata:
      labels:
        type: app
    spec:
      containers:
        - name: container
          image: tmkube/app:v1
```

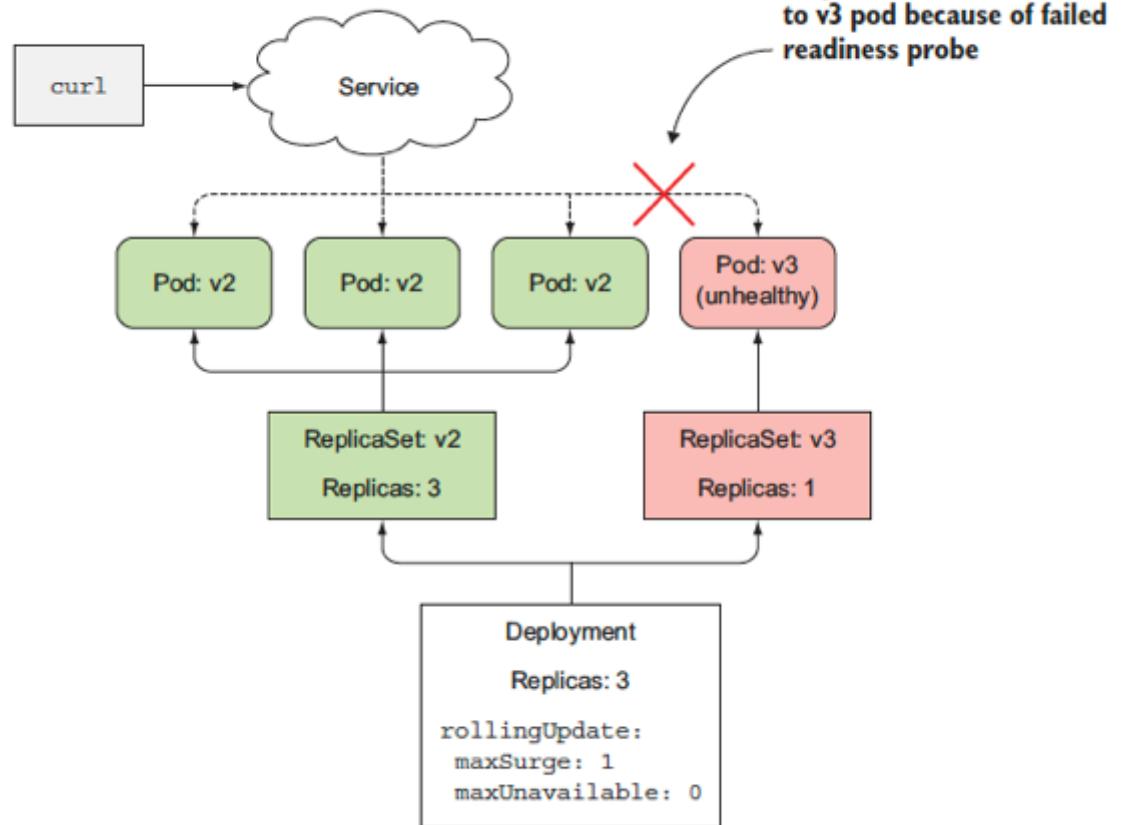
```
apiVersion: v1
kind: Service
metadata:
  name: svc-2
spec:
  selector:
    type: app
  ports:
    - port: 8080
      protocol: TCP
      targetPort: 8080
```

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: deployment-2
spec:
  selector:
    matchLabels:
      type: app
  replicas: 2
  strategy:
    type: RollingUpdate
    minReadySeconds: 10
  template:
    metadata:
      labels:
        type: app
    spec:
      containers:
        - name: container
          image: tmkube/app:v1
```

# Deployment Flow Diagram



# Deployment + readinessProbe



```
apiVersion: apps/v1beta1
kind: Deployment
metadata:
  name: kubia
spec:
  replicas: 3
  minReadySeconds: 10
  strategy:
    rollingUpdate:
      maxSurge: 1
      maxUnavailable: 0
      type: RollingUpdate
  template:
    metadata:
      name: kubia
      labels:
        app: kubia
    spec:
      containers:
        - image: luksa/kubia:v3
  readinessProbe:
    periodSeconds: 1
    httpGet:
      path: /
      port: 8080
```

You're keeping  
minReadySeconds  
set to 10.

You're keeping  
maxUnavailable  
set to 0 to make the deployment  
replace pods one by one

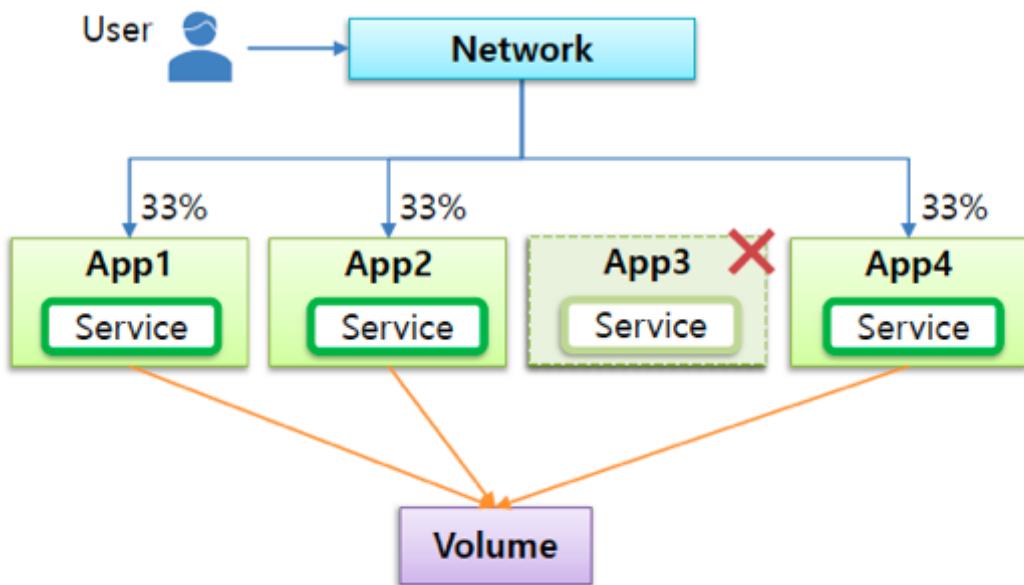
You're defining a readiness probe  
that will be executed every second.

The readiness probe will  
perform an HTTP GET request  
against our container.

# StatefulSet

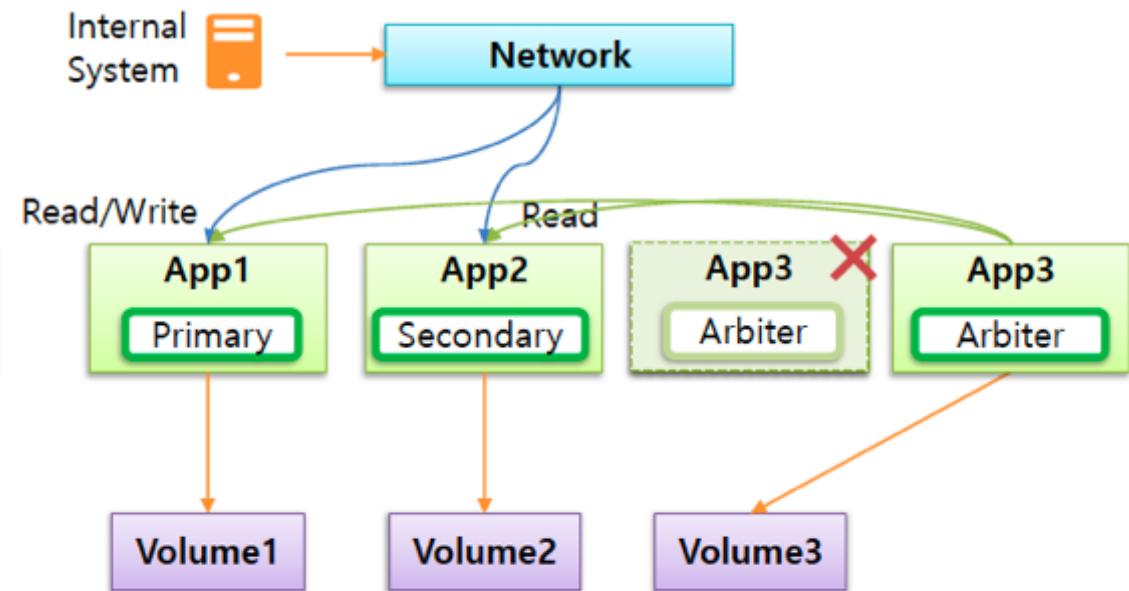
## Stateless Application

### Web Server



## Stateful Application

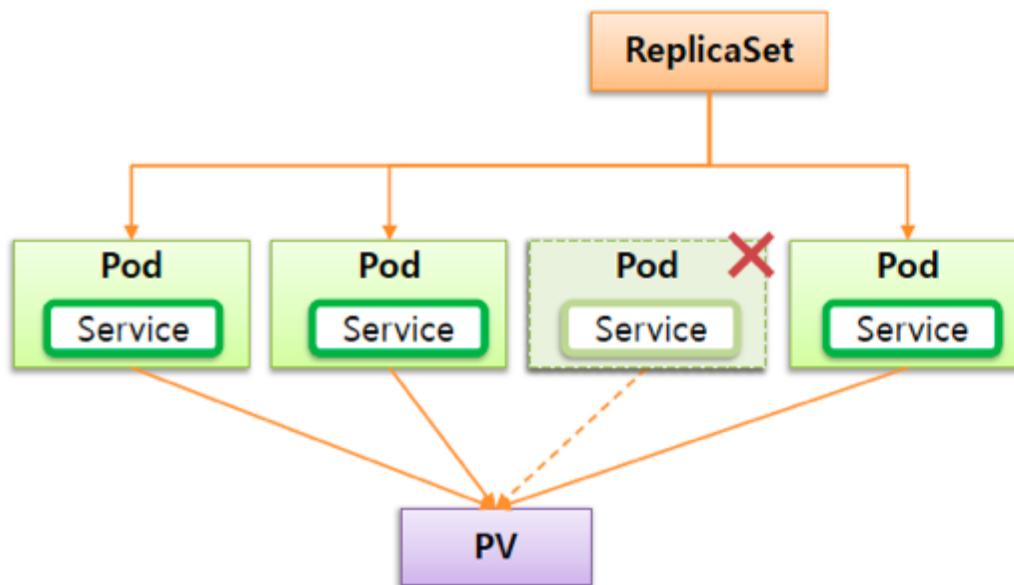
### Database



# StatefulSet

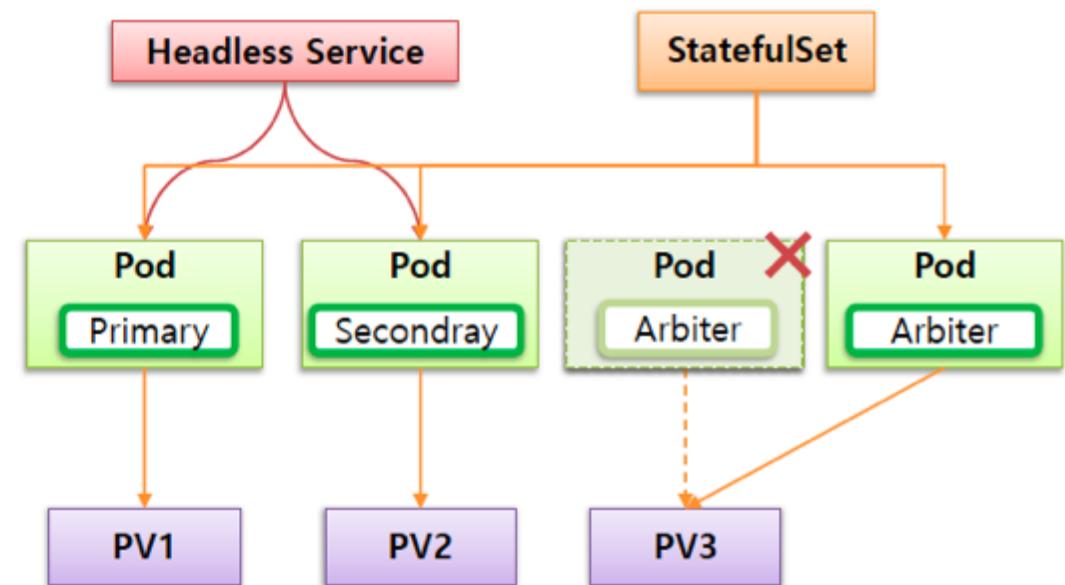
## Stateless Application

Web Server

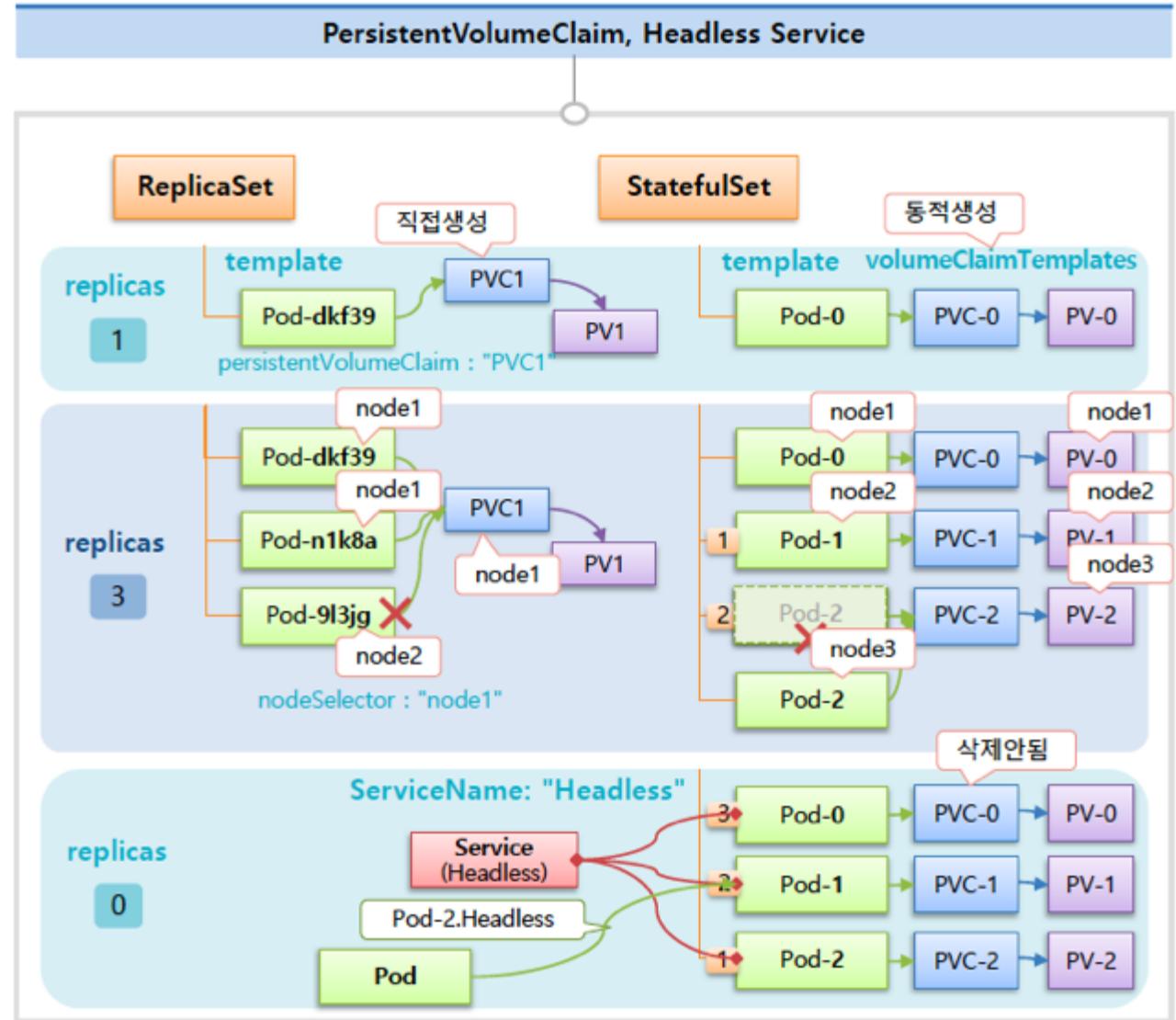
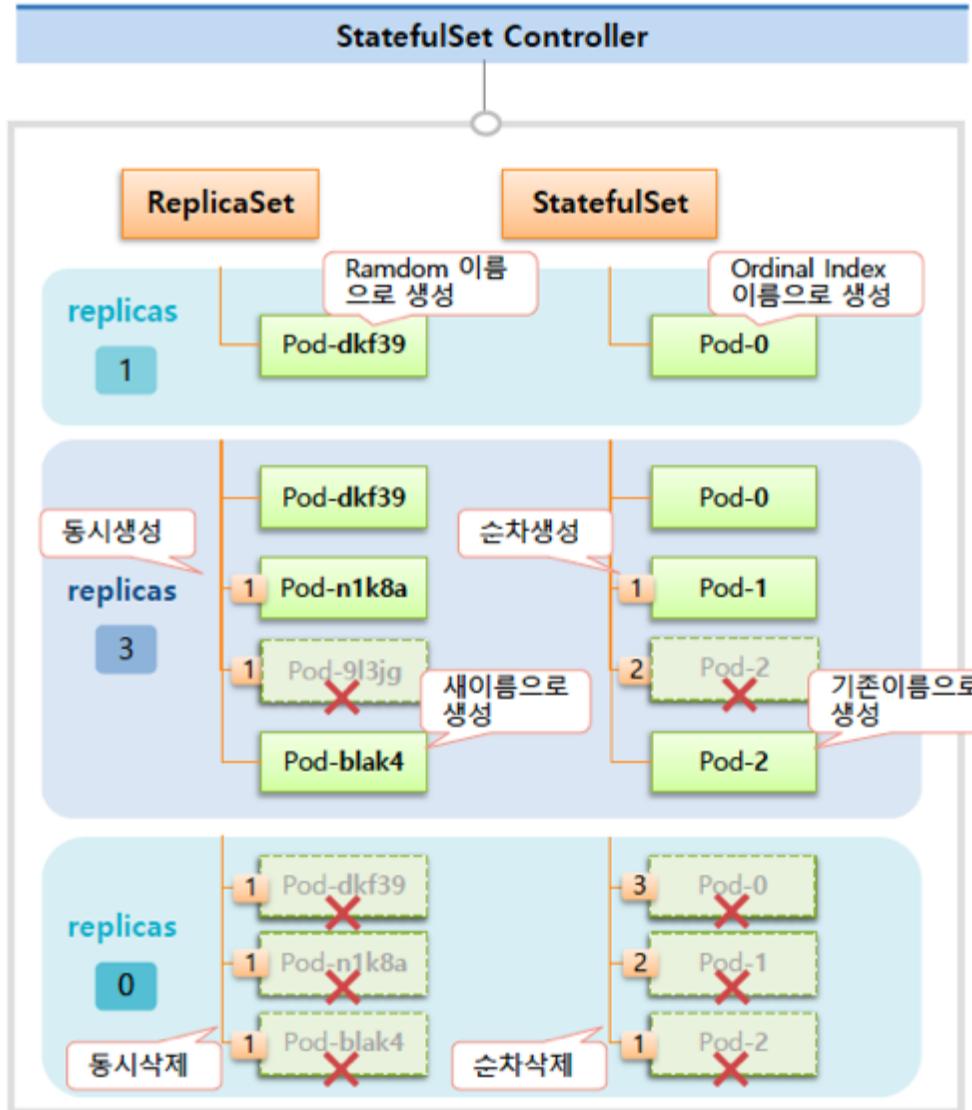


## Stateful Application

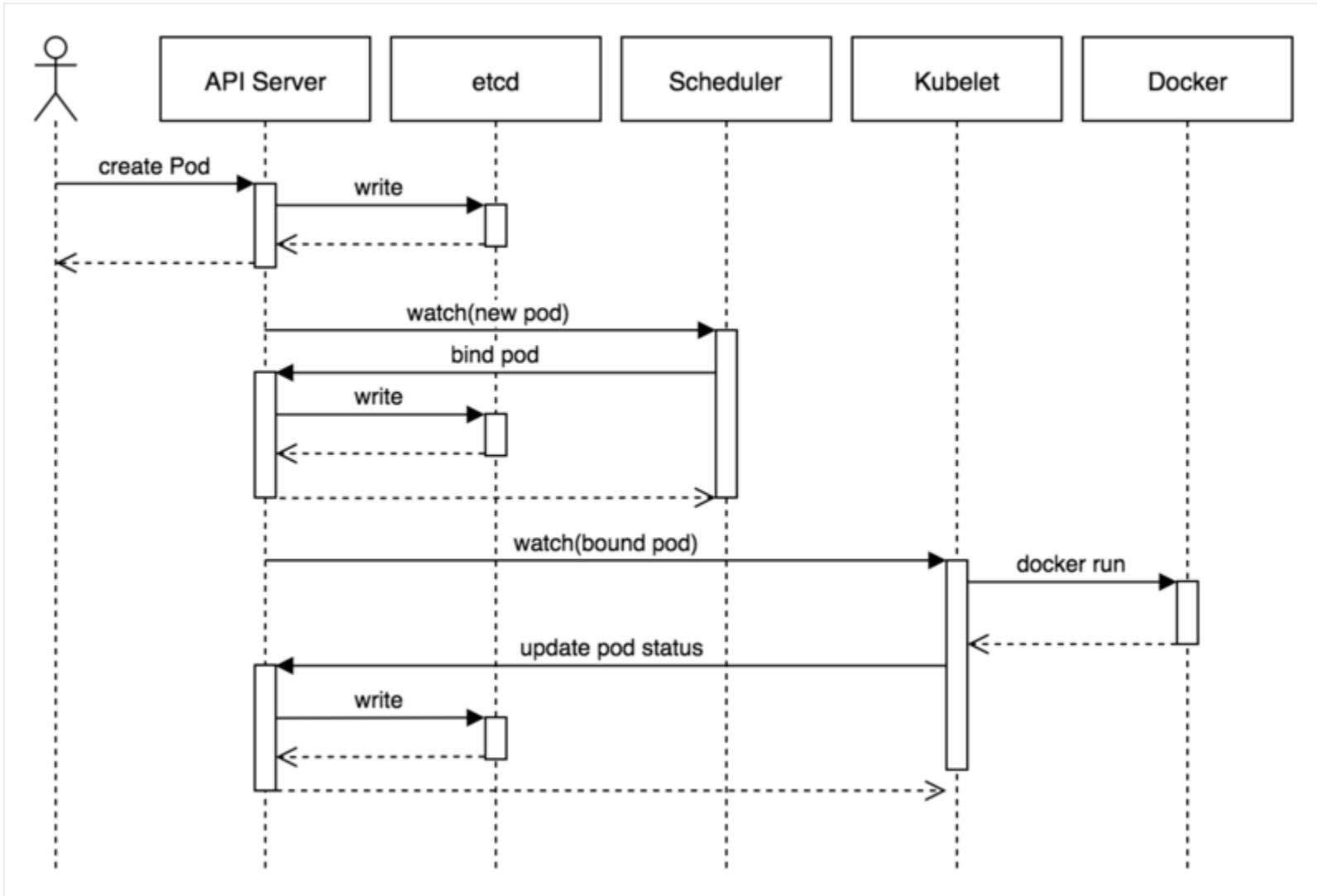
Database



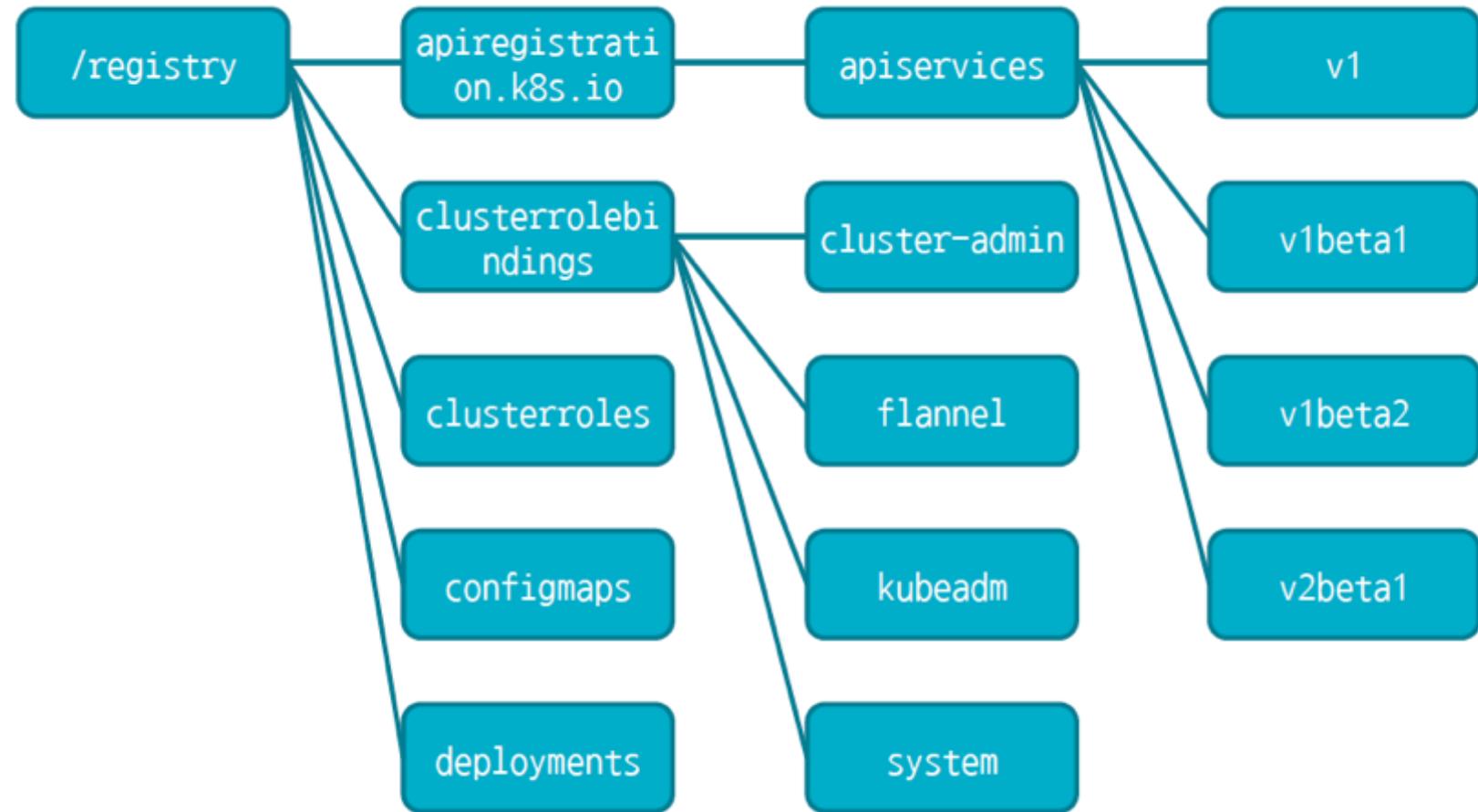
# StatefulSet



# etcd: persist data



# etcd Structure

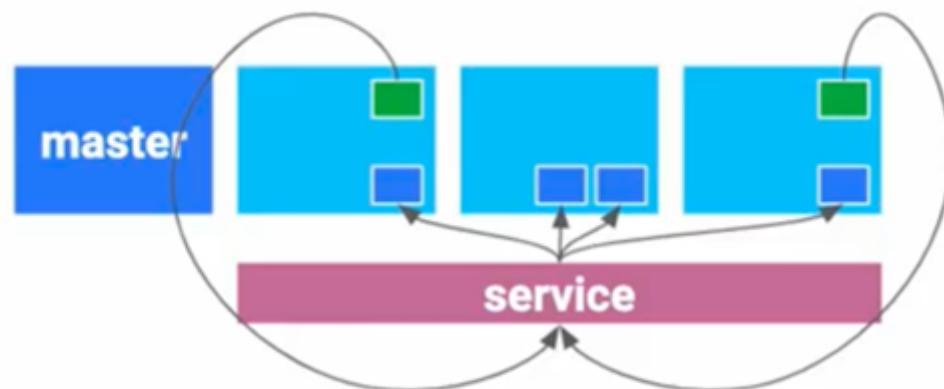


# kube-dns

kube-dns is an addon that resolves  
`\$service.\$namespace(.svc.cluster.local)` to the  
right Service IP. This way you can talk to other apps  
by knowing only the Service name.

Runs as a deployment fronted by a service

When deployed, every container uses the service ip  
for kube-dns in its /etc/resolv.conf



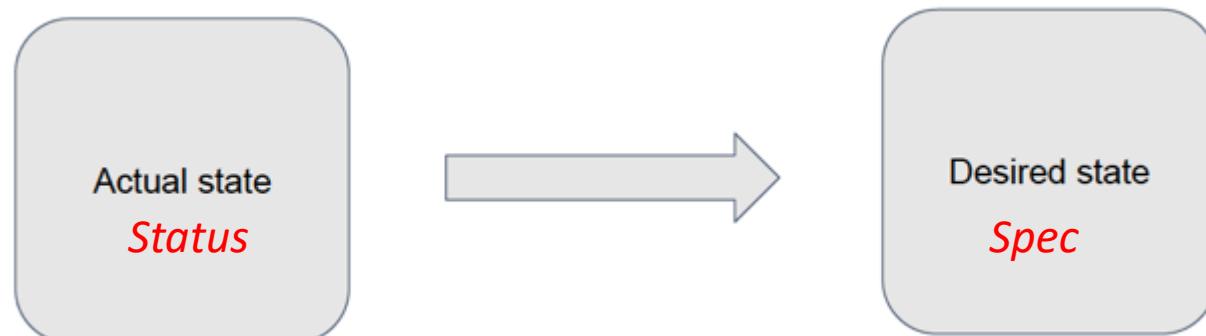
# Declarative Configuration

## Before:

- **You:** provide exact set of instructions to drive to desired state
- **System:** executes instructions
- **You:** monitor system, and provide further instructions if it deviates.

## After:

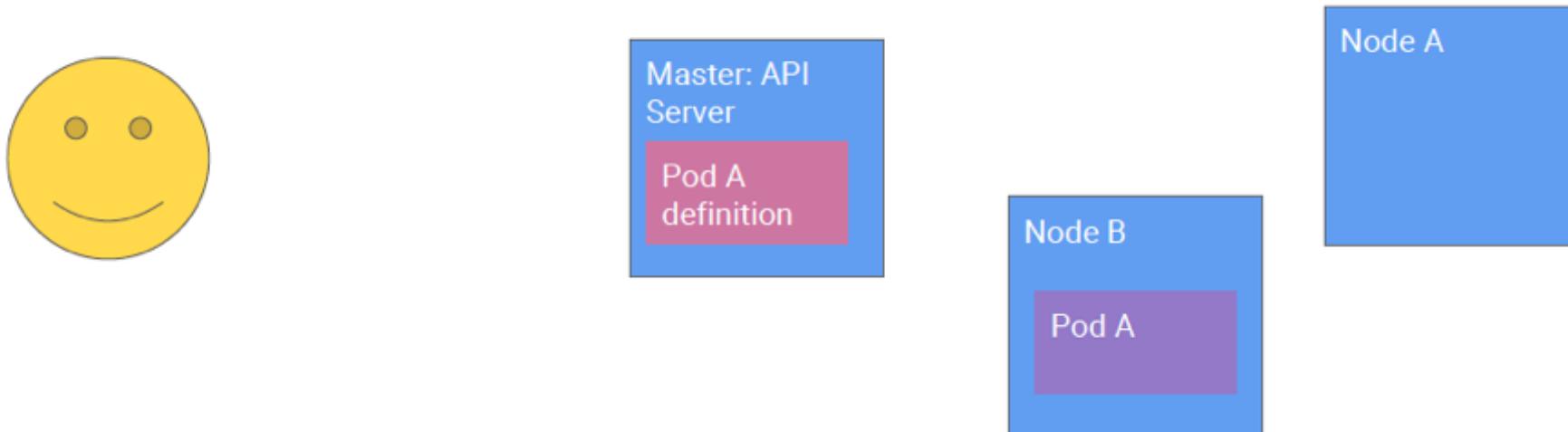
- **You:** define desired state
- **System:** works to drive towards that state



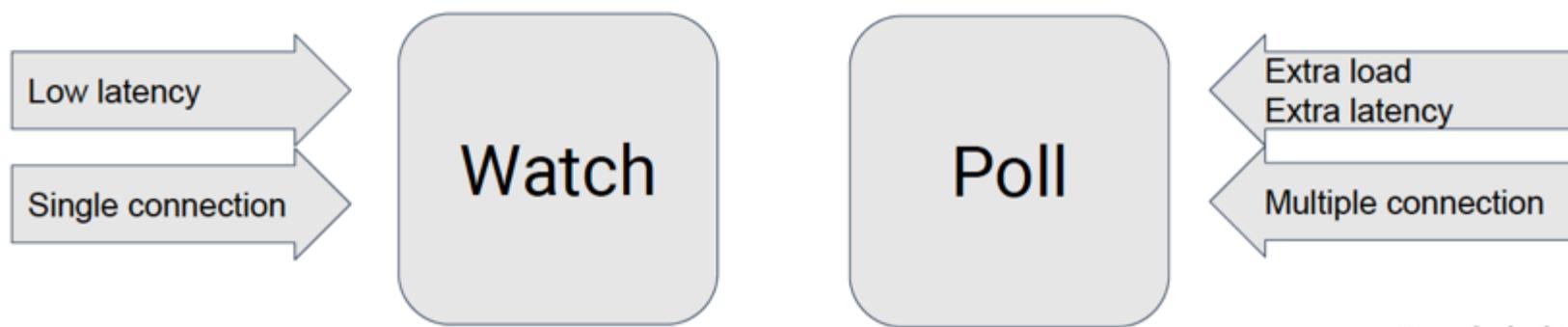
# Kubernetes Way to Deploy Workload

The Kubernetes way!

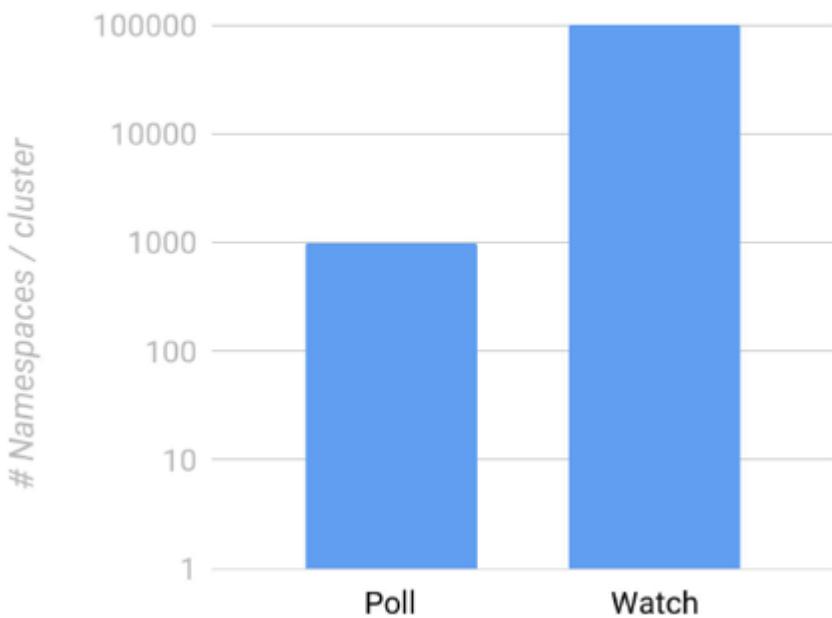
- **You:** create API object that is persisted on kube API server until deletion
- **System:** all components work in parallel to drive to that state



# Kubernetes API server: Watch vs Poll



Scalability of Poll vs Watch

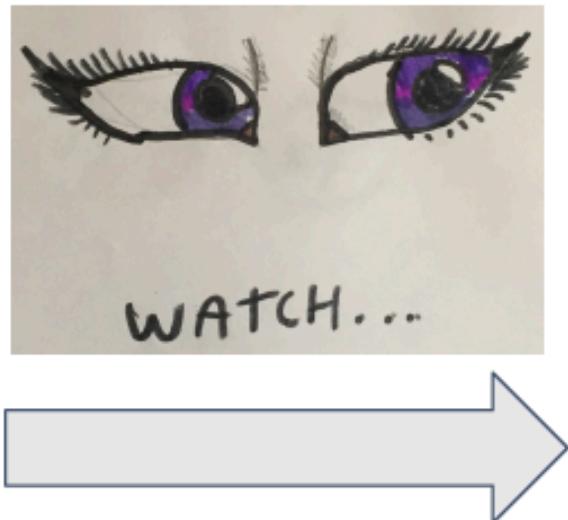


Kubelet on nodes:

- Previous: periodically poll kube-apiserver for secrets and configmaps
- Now: watch individual secrets
- OSS PR: [Kubelet watches necessary secrets/configmaps instead of periodic polling #64752](#)

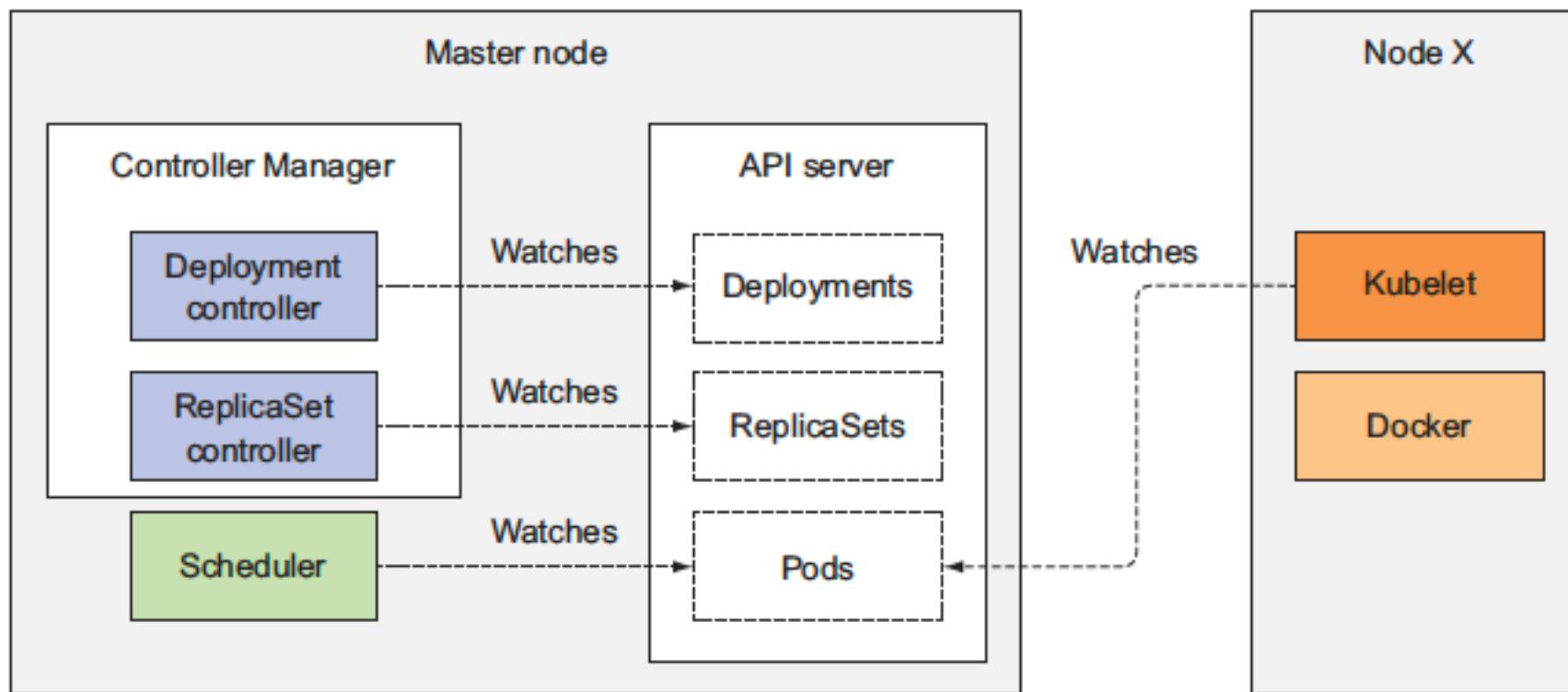
# Watch Event

Actual state

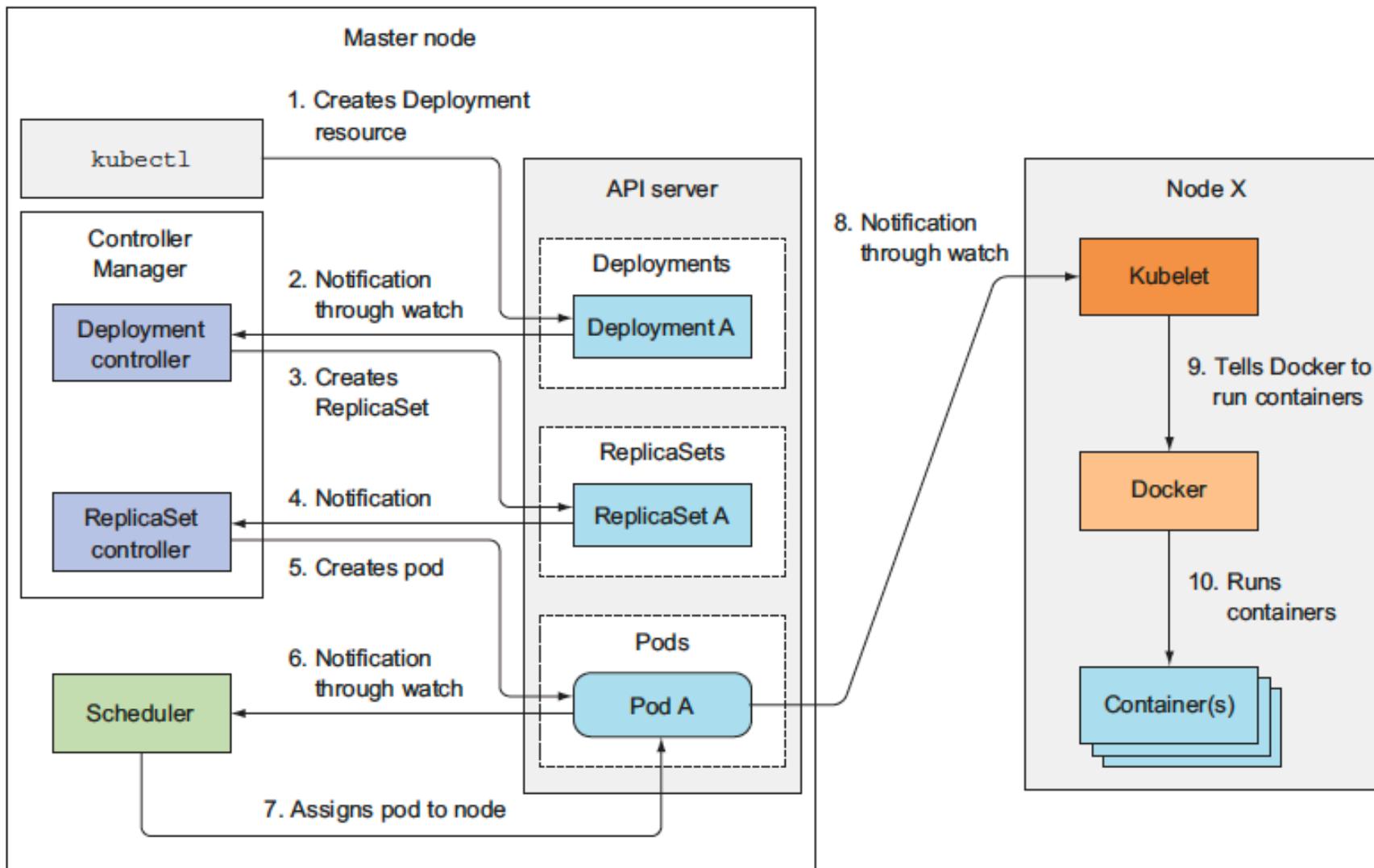


Desired state

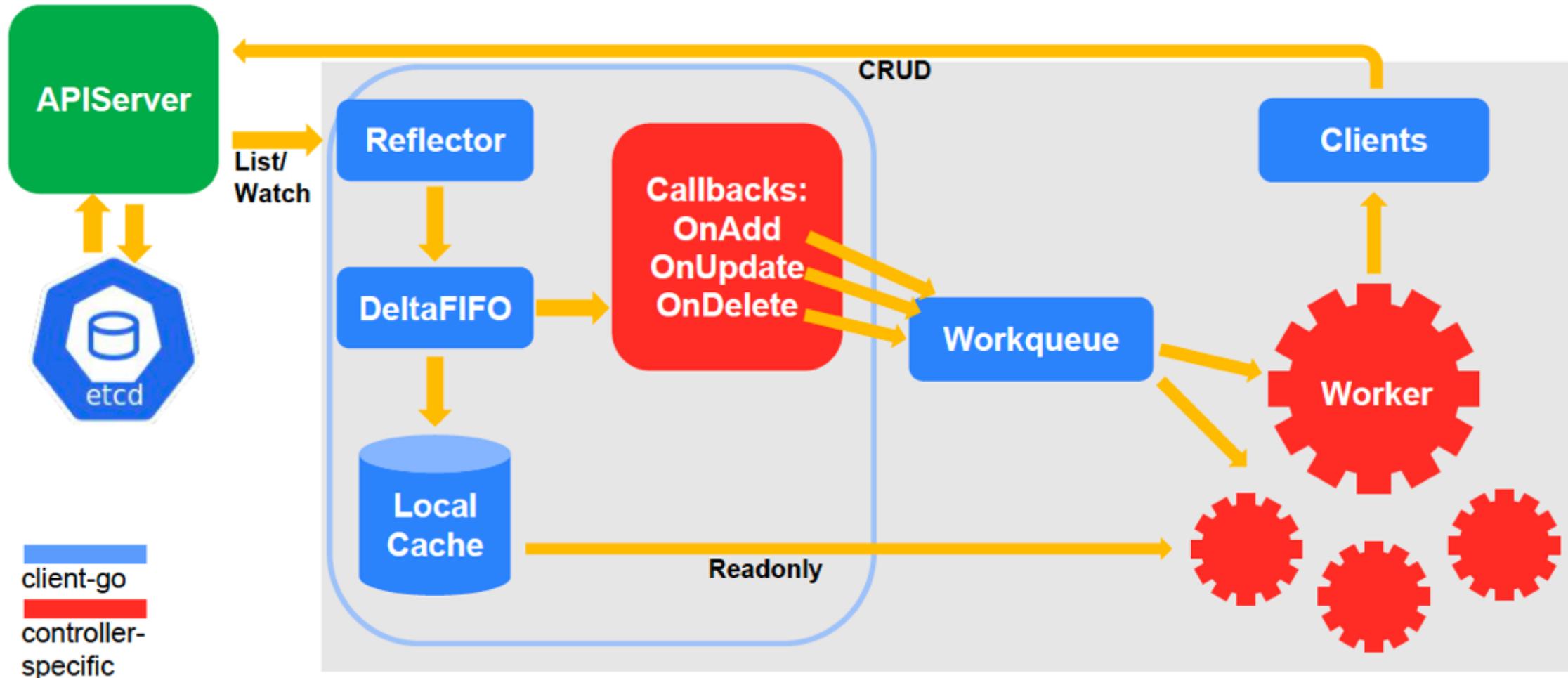
# Watches



# Event Chain



# Kubernetes Controller Workflow

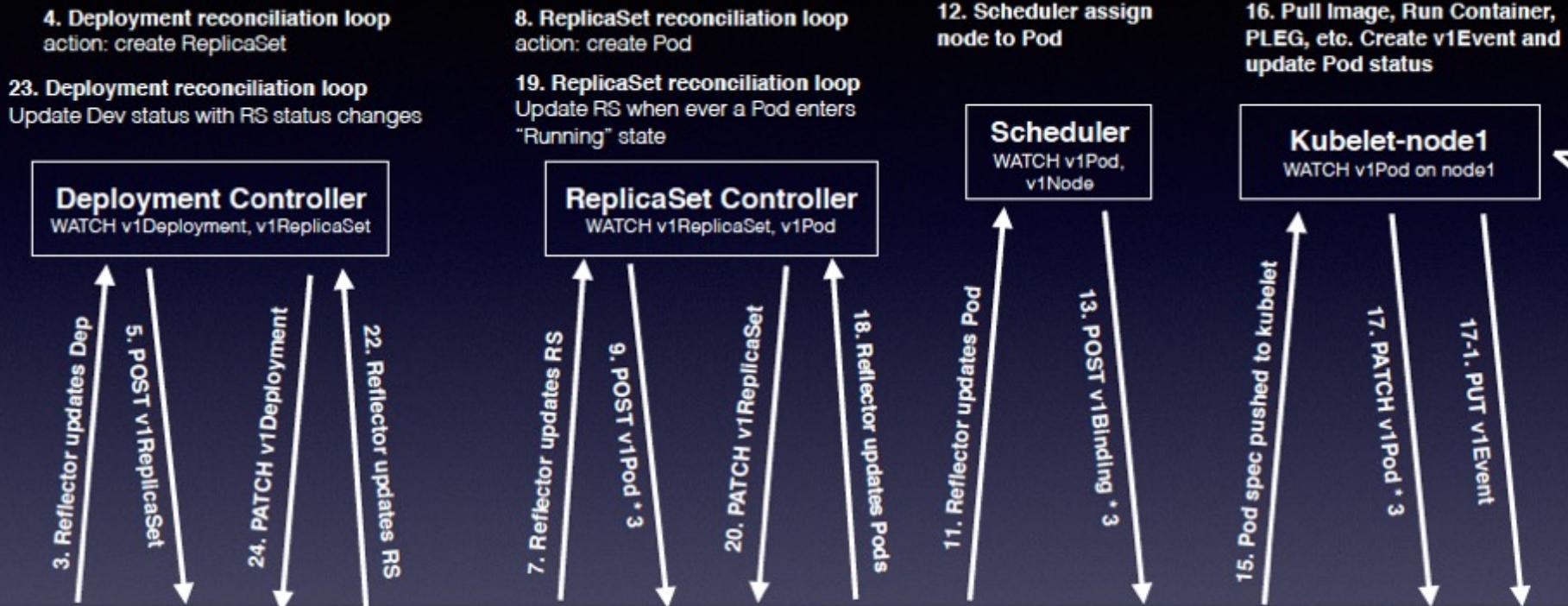


Note: this is a rather high-level idea about Kubernetes reaction chain. It hides some details for illustration purpose and is different than actual behavior



Kubectl

1. POST v1Deployment



## Kubernetes API Server & Etcd

v1Deployment

v1ReplicaSet

v1Pod \* 3

v1Binding \* 3

v1Event \* N

2. Persist v1Deployment API object  
desiredReplica: 3  
availableReplica: 0

25. Update v1Deployment API object  
desiredReplica: 3  
availableReplica: 0 -> 3

6. Persist v1ReplicaSet API object  
desiredReplica: 3  
availableReplica: 0

21. Update v1ReplicaSet API object  
desiredReplica: 3  
availableReplica: 0 -> 3

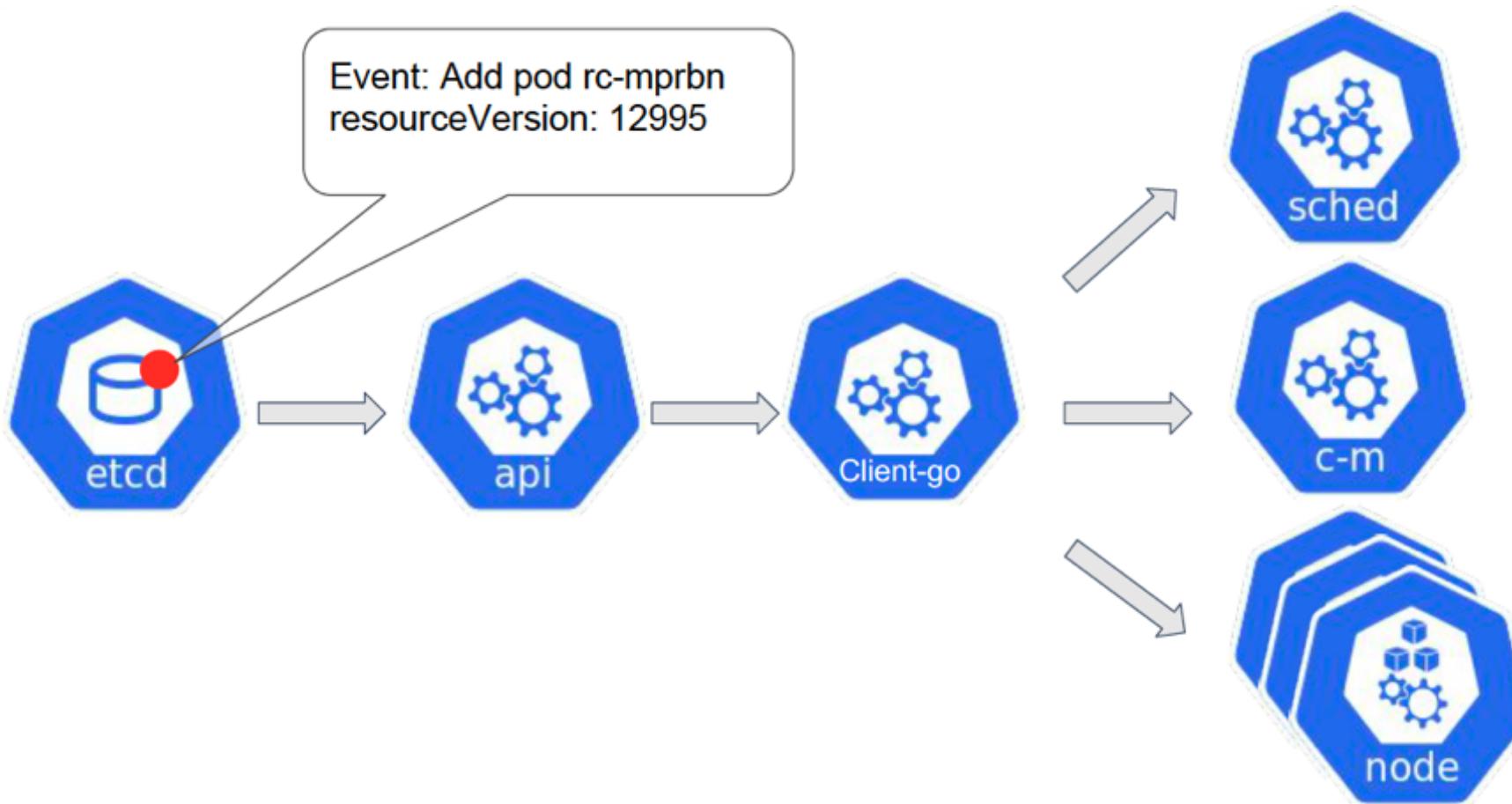
10. Persist 3 v1Pod API objects  
status: Pending  
nodeName: none

17. Update 3 v1Pod API objects  
status: PodInitializing -> Running  
nodeName: node1

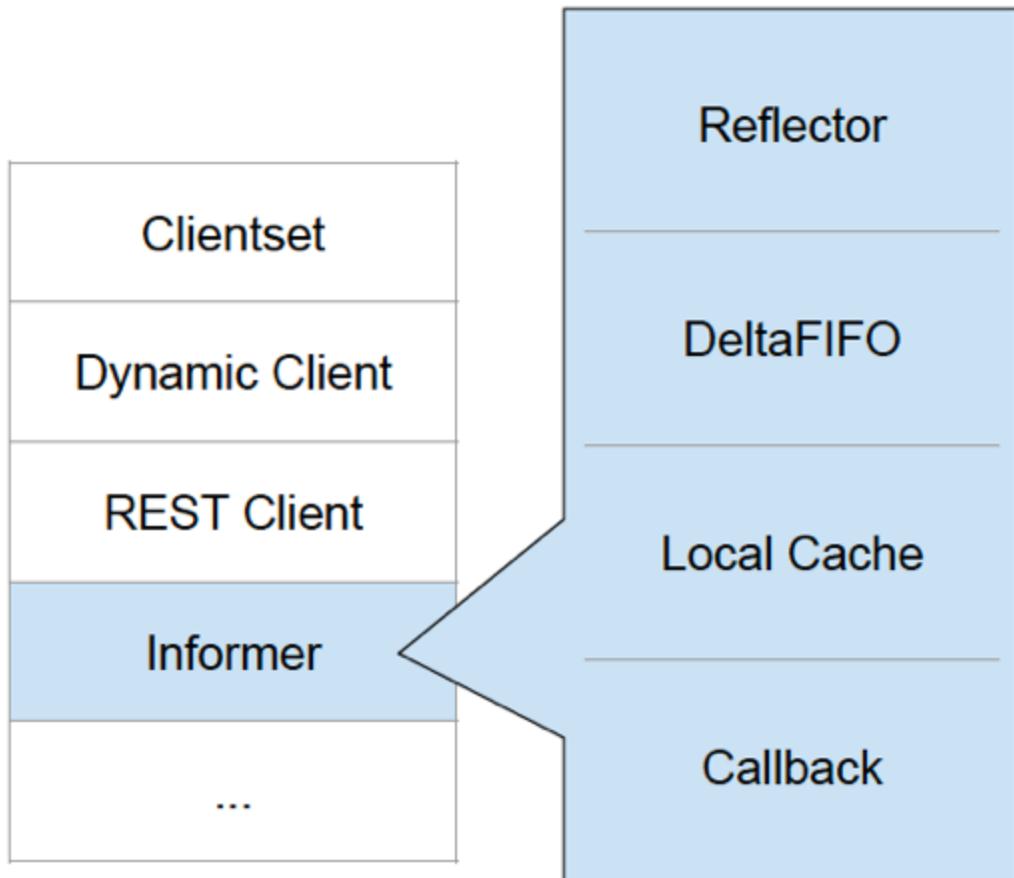
14. Persist 3 v1Binding API objects  
name: PodName  
targetNodeName: node1

17-1. Several v1Event objects created  
e.g. ImagePulled, CreatingContainer,  
StartedContainer, etc., controllers, scheduler  
can also create events throughout the entire  
reaction chain

# Client-go (1)

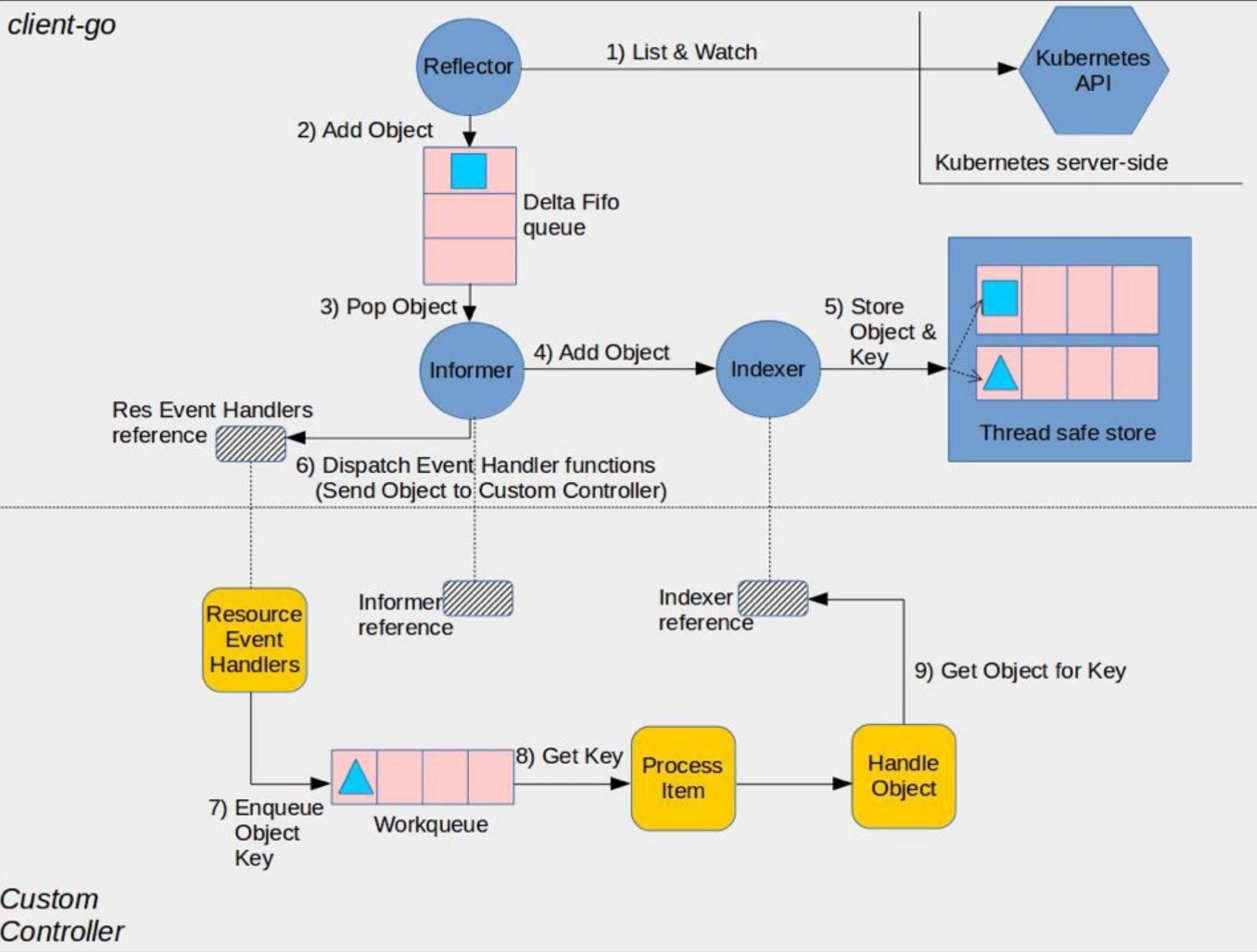


# Client-go (2)



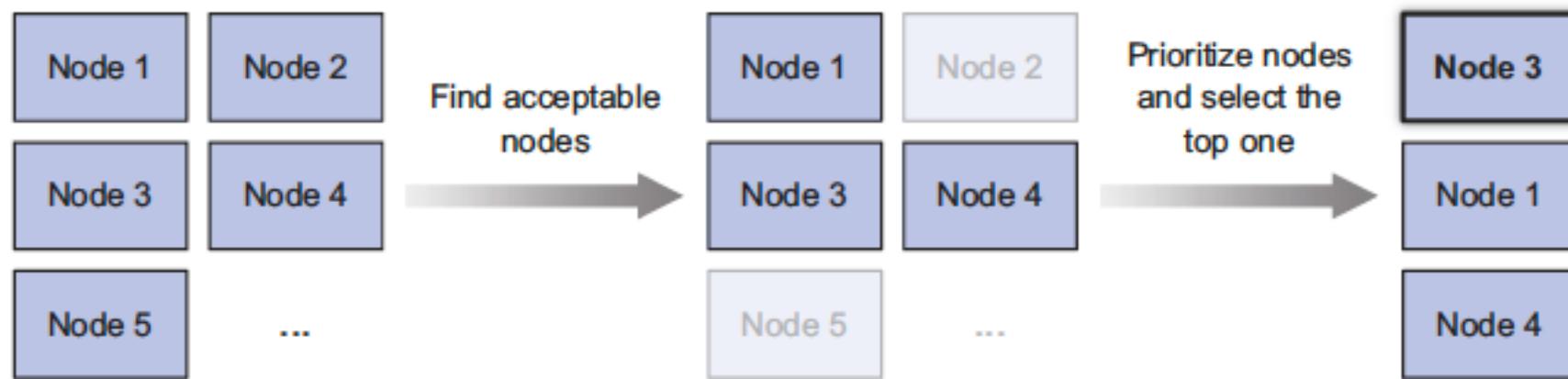
[k8s.io/client-go/tools/cache](https://k8s.io/client-go/tools/cache)  
[k8s.io/client-go/informers](https://k8s.io/client-go/informers)

- Useful component for building event-oriented controllers
- Used by control plane controllers, kubelet, etc.
- Reflector used by kube-apiserver watch cache

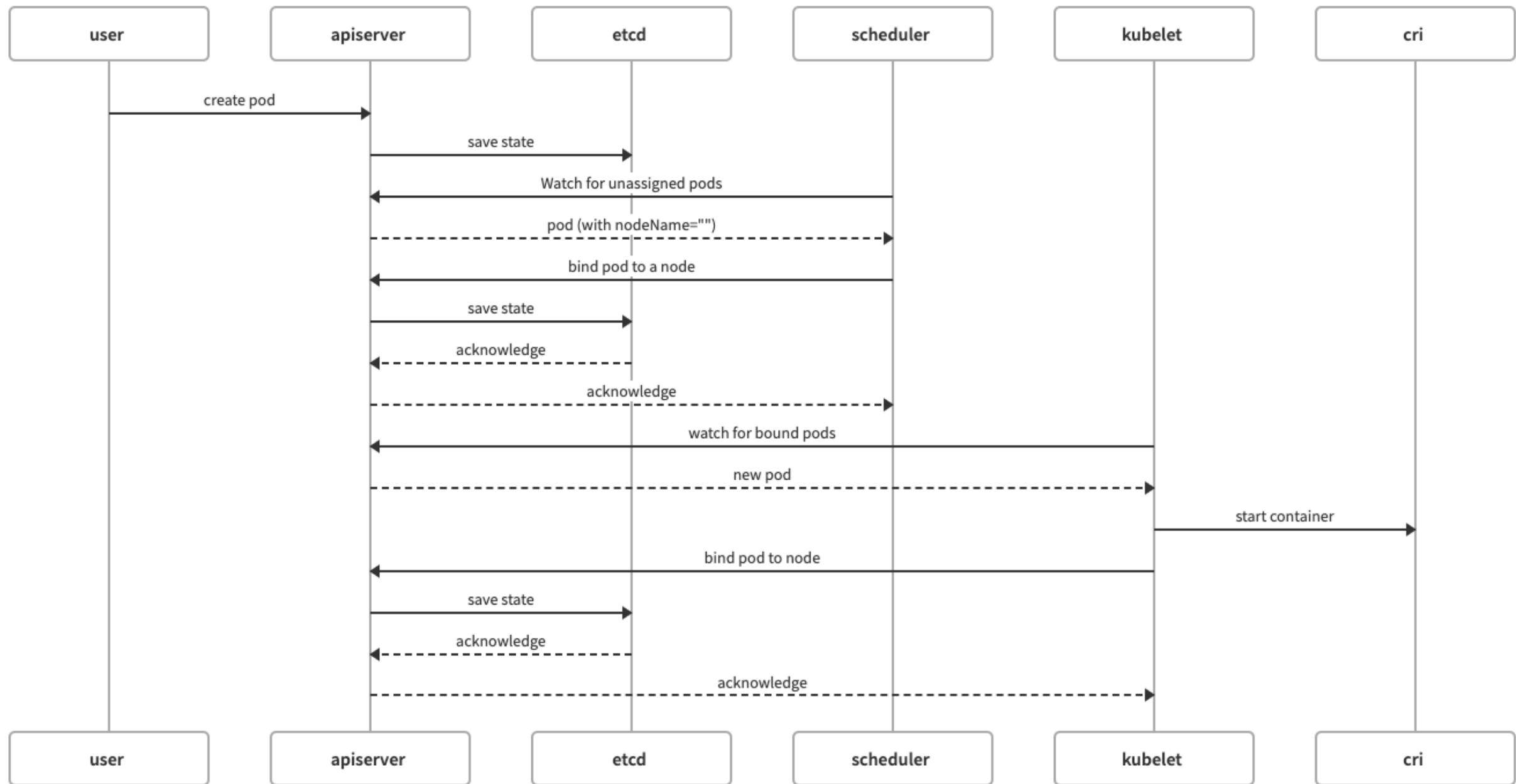


*Custom  
Controller*

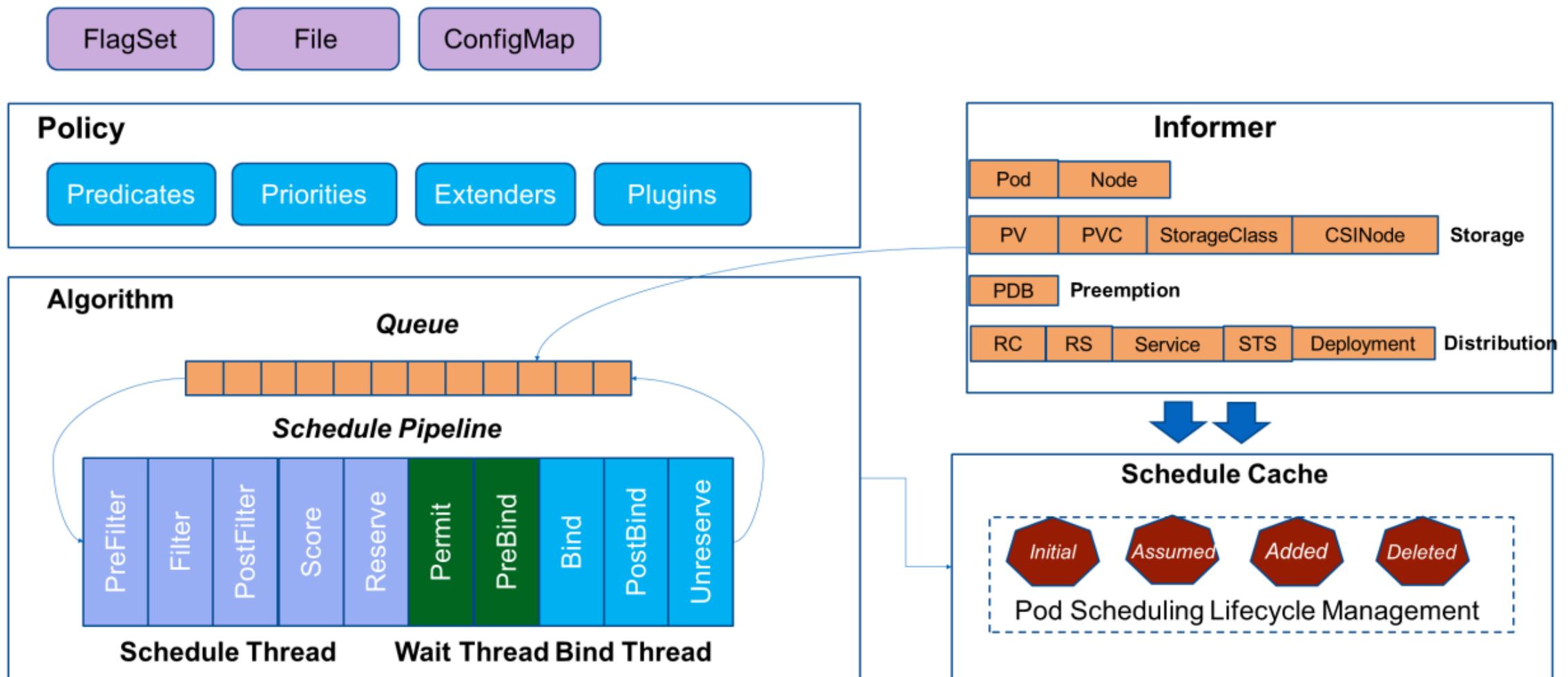
# Kubernetes Scheduler



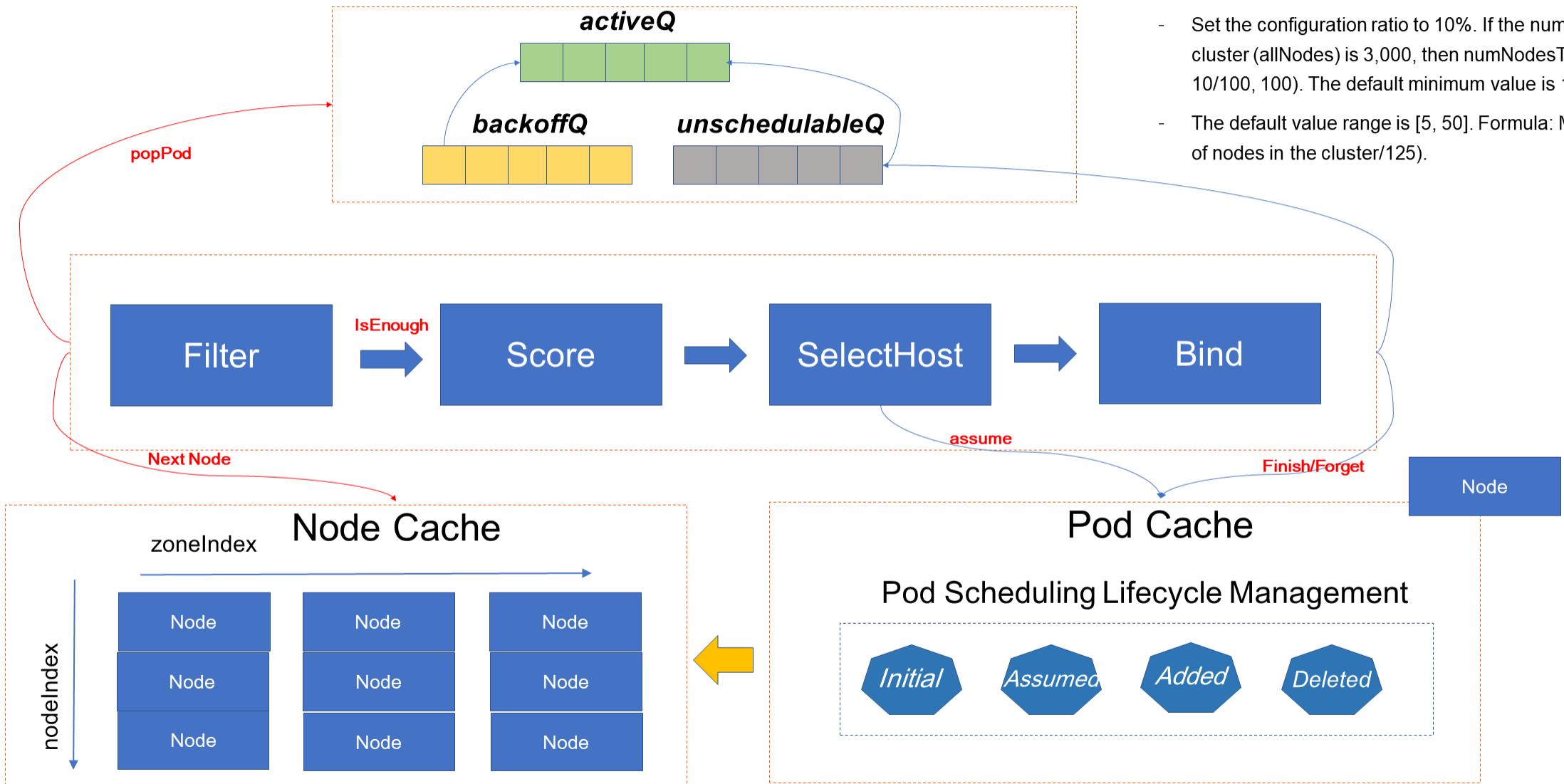
# Kubernetes Scheduler



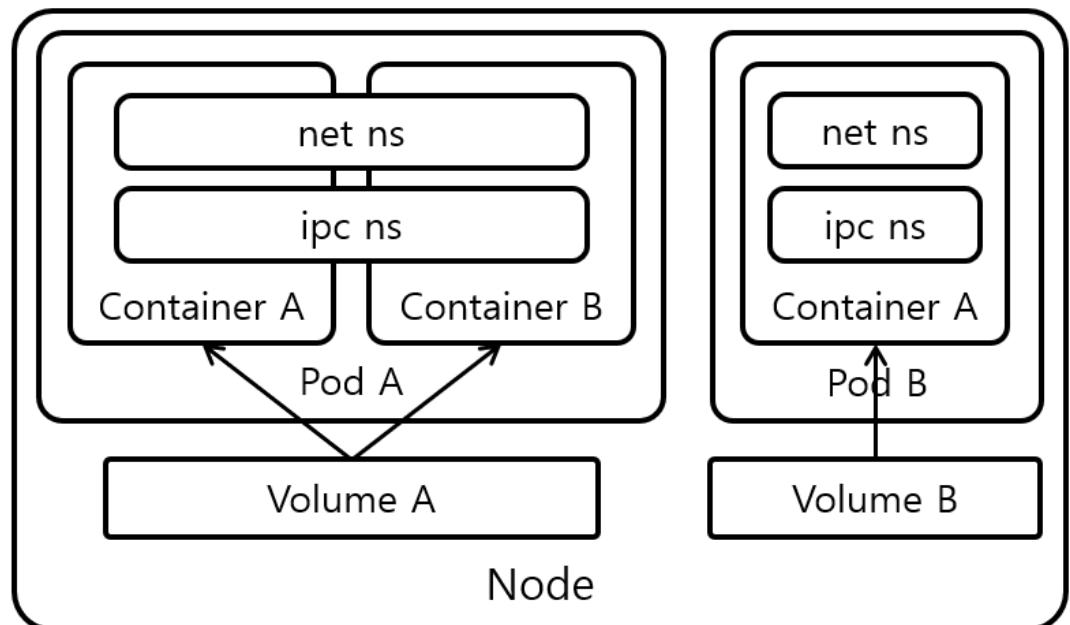
# Scheduler Overview



# Scheduling Process

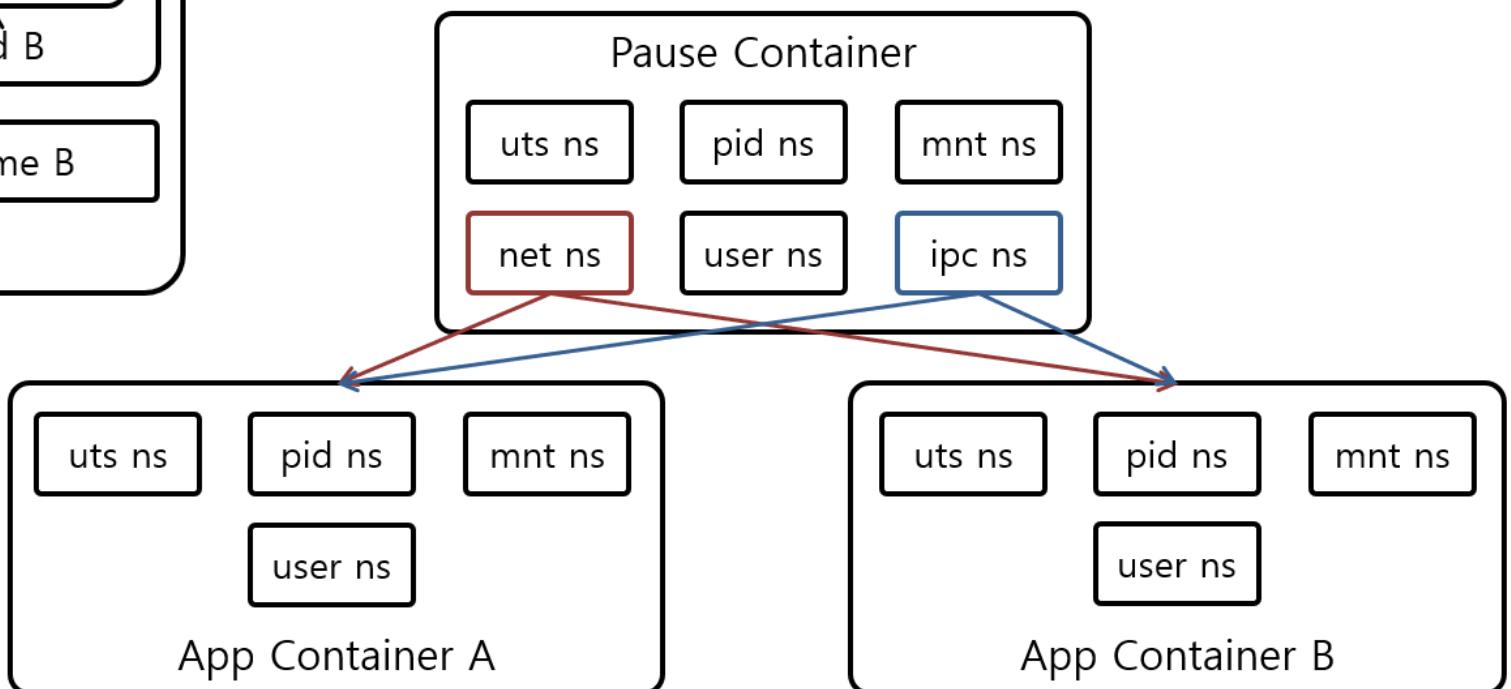


# Pause Container

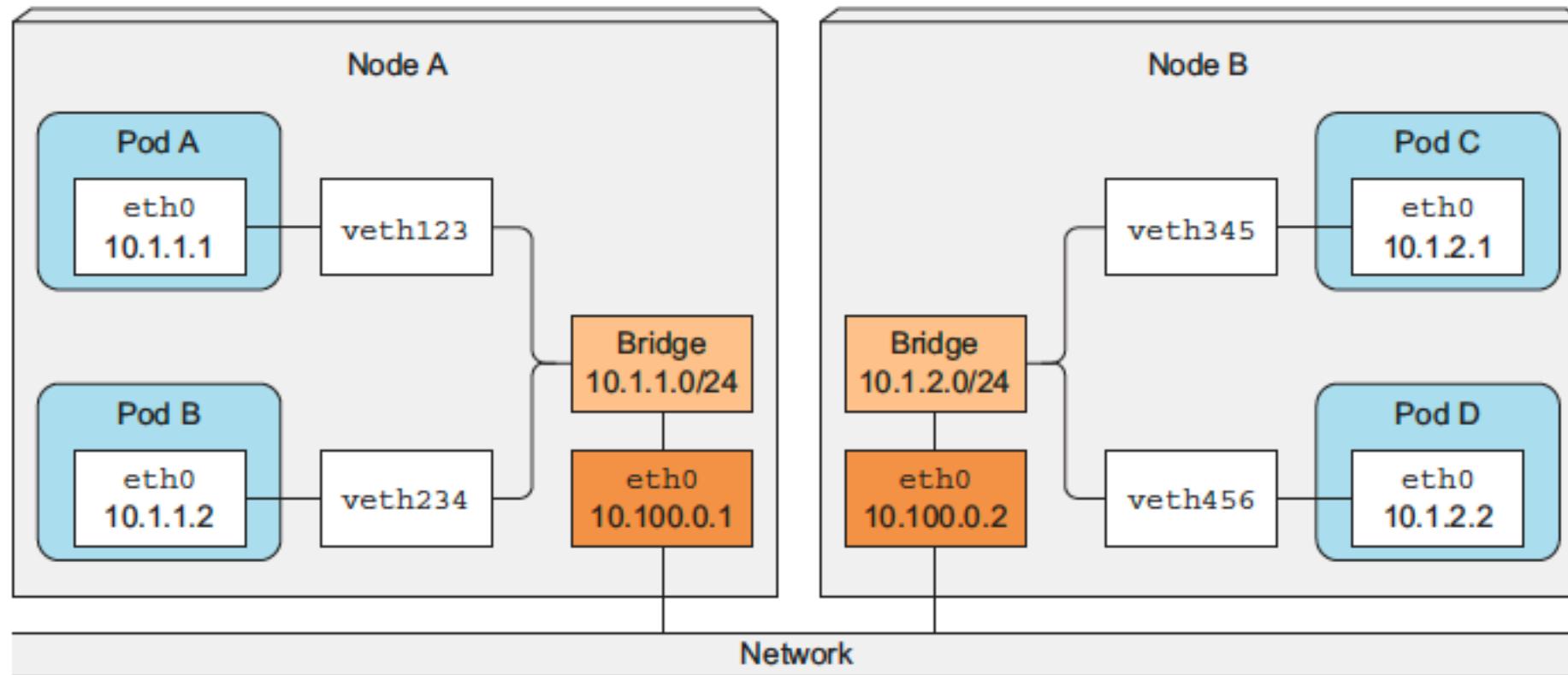


Init Process killed -> SIGKILL to other cloned processes

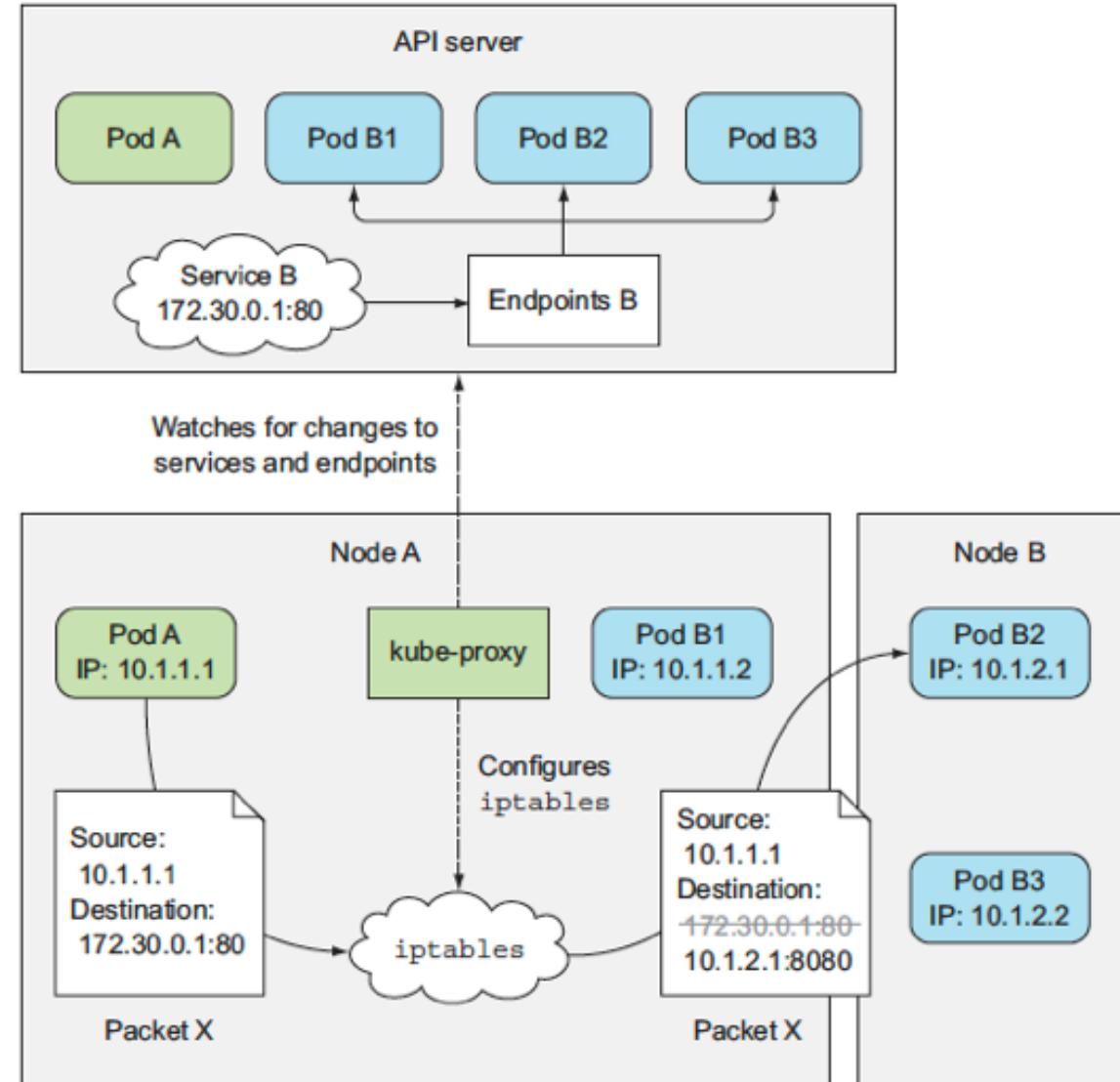
Pause binary Init Process: **pause()** syscall -> blocking until SIG

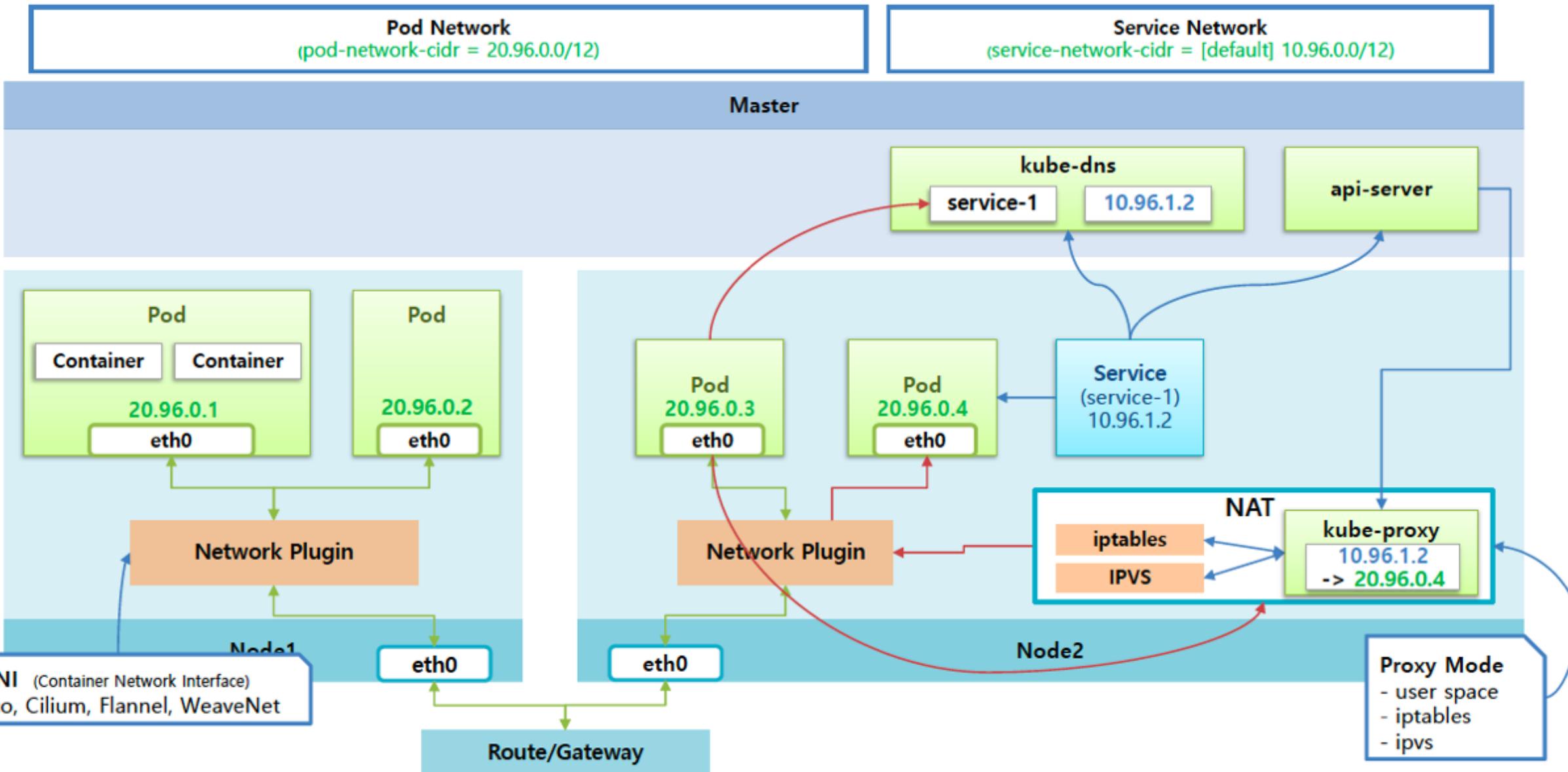


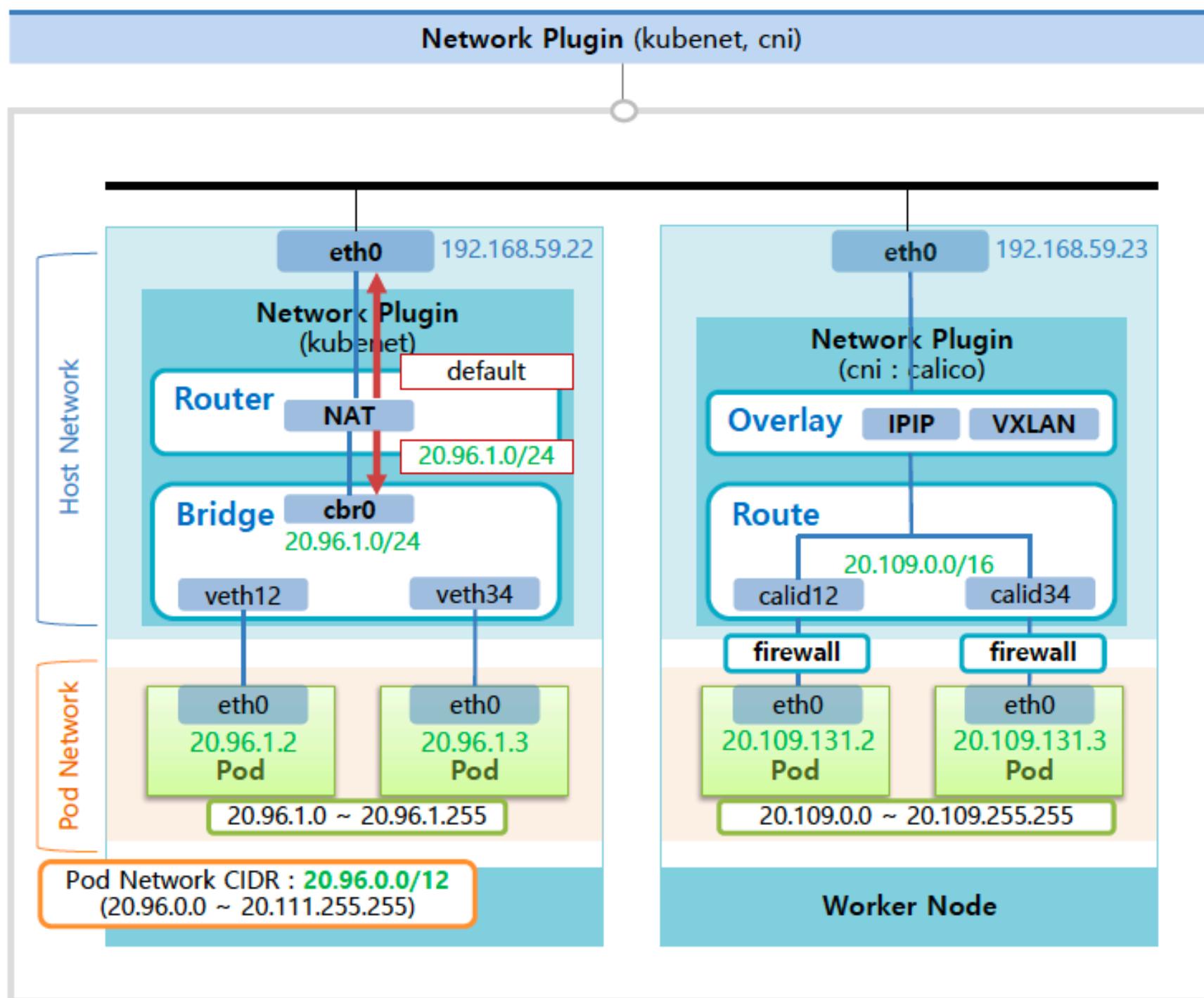
# Communication between Pods and Nodes



# kube-proxy & iptables







## Proxy Mode

Userspace Mode

Iptables/IPVS Mode

### Pod Network Area

Pod  
20.111.156.6

Worker Node

Endpoint  
Service  
10.103.9.116

api-server

Master Node

Service Network CIDR : 10.96.0.0/12  
(10.96.0.0 ~ 10.111.255.255)

kube-proxy

10.103.9.116  
-> 20.111.156.6

iptables

Service CIDR

10.103.9.116

Pod

20.111.156.6

20.111.156.6

iptables

IPVS

10.103.9.116  
-> 20.96.2.5

## Service Type (ClusterIP)

192.168.59.22

eth0

Network Plugin  
(cni : calico)

Overlay

IPIP

Route

NAT

calid12

calid34

eth0

20.111.156.6

Pod A

eth0

20.111.156.7

Pod B

Worker N

Service  
[ClusterIP]  
10.103.9.116

192.168.59.23

eth0

Network  
(cni : calico)

Overlay

IPIP

Route

NAT

calid12

calid34

eth0

20.109.131.2

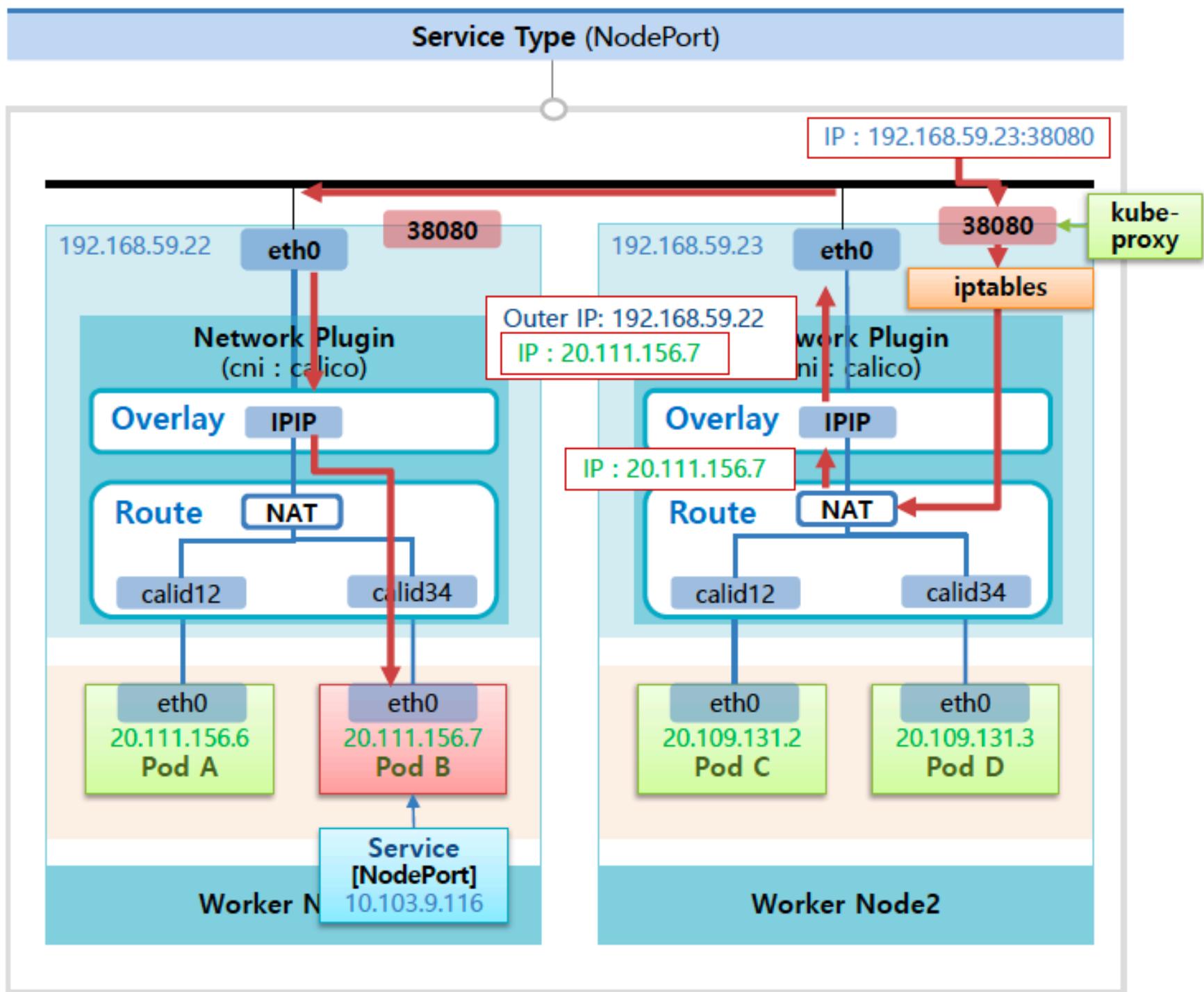
Pod C

eth0

20.109.131.3

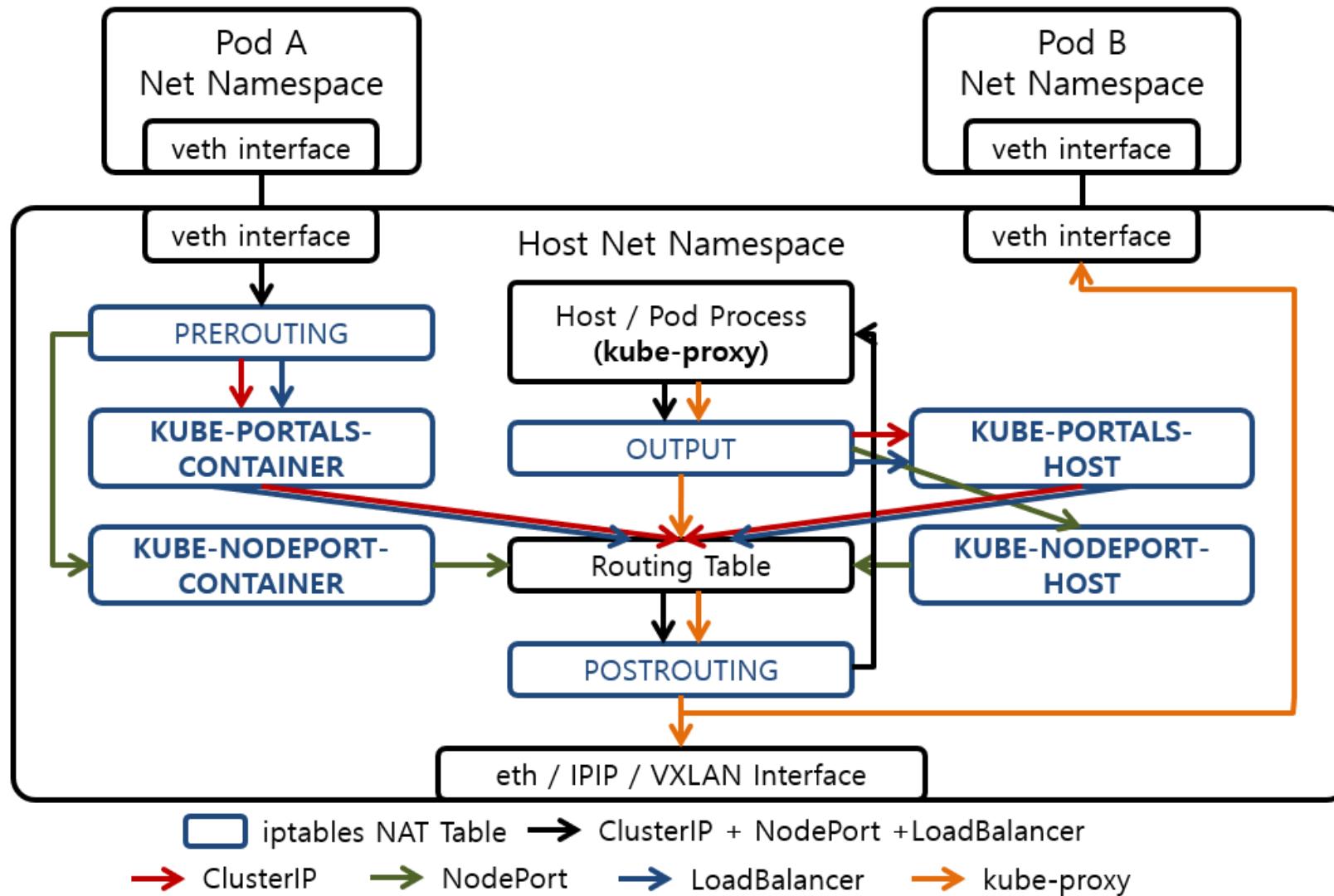
Pod D

Worker Node2

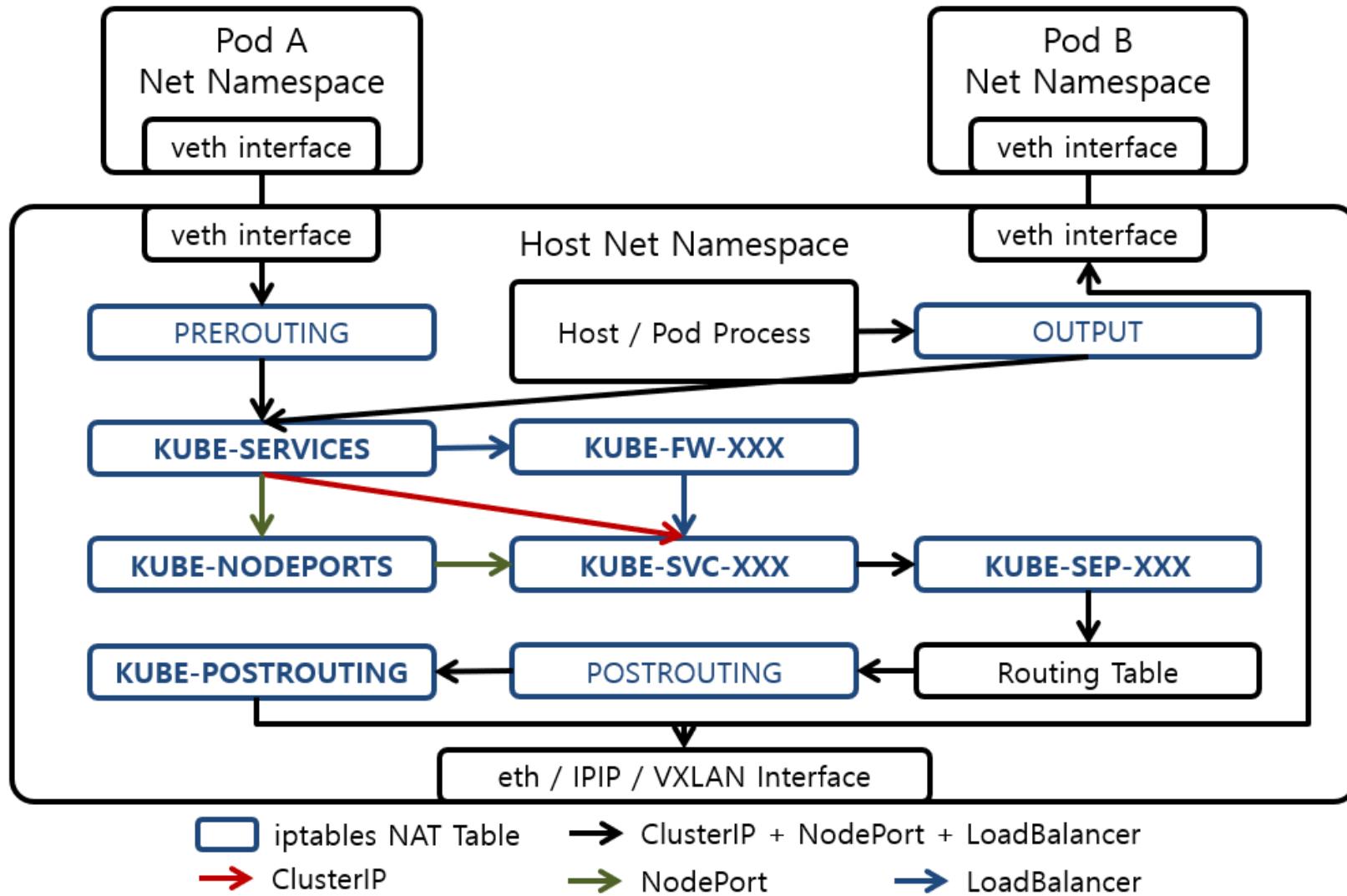


# Userspace Mode

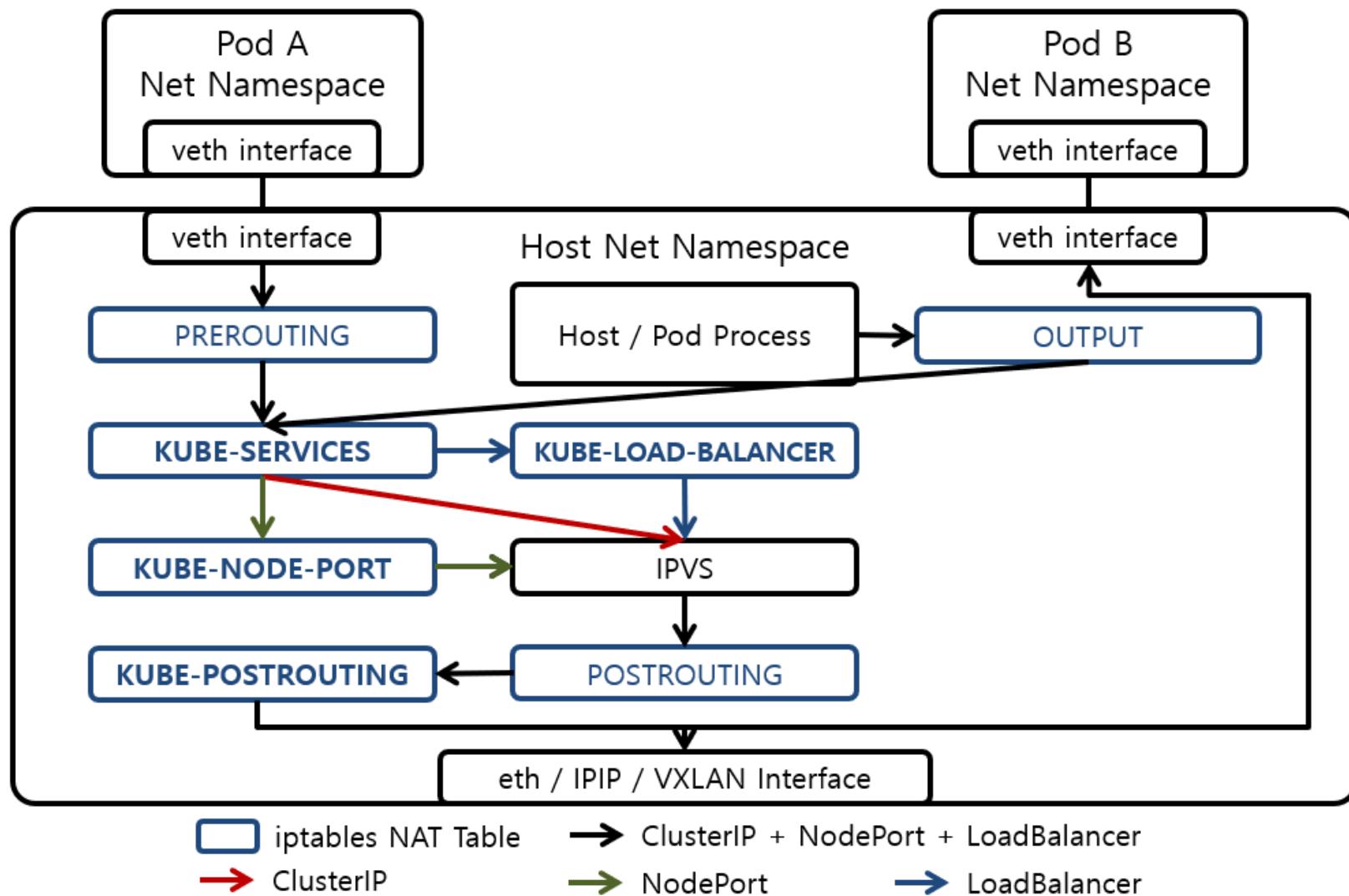
[https://ssup2.github.io/theory\\_analysis/  
Kubernetes Service Proxy/](https://ssup2.github.io/theory_analysis/Kubernetes_Service_Proxy/)



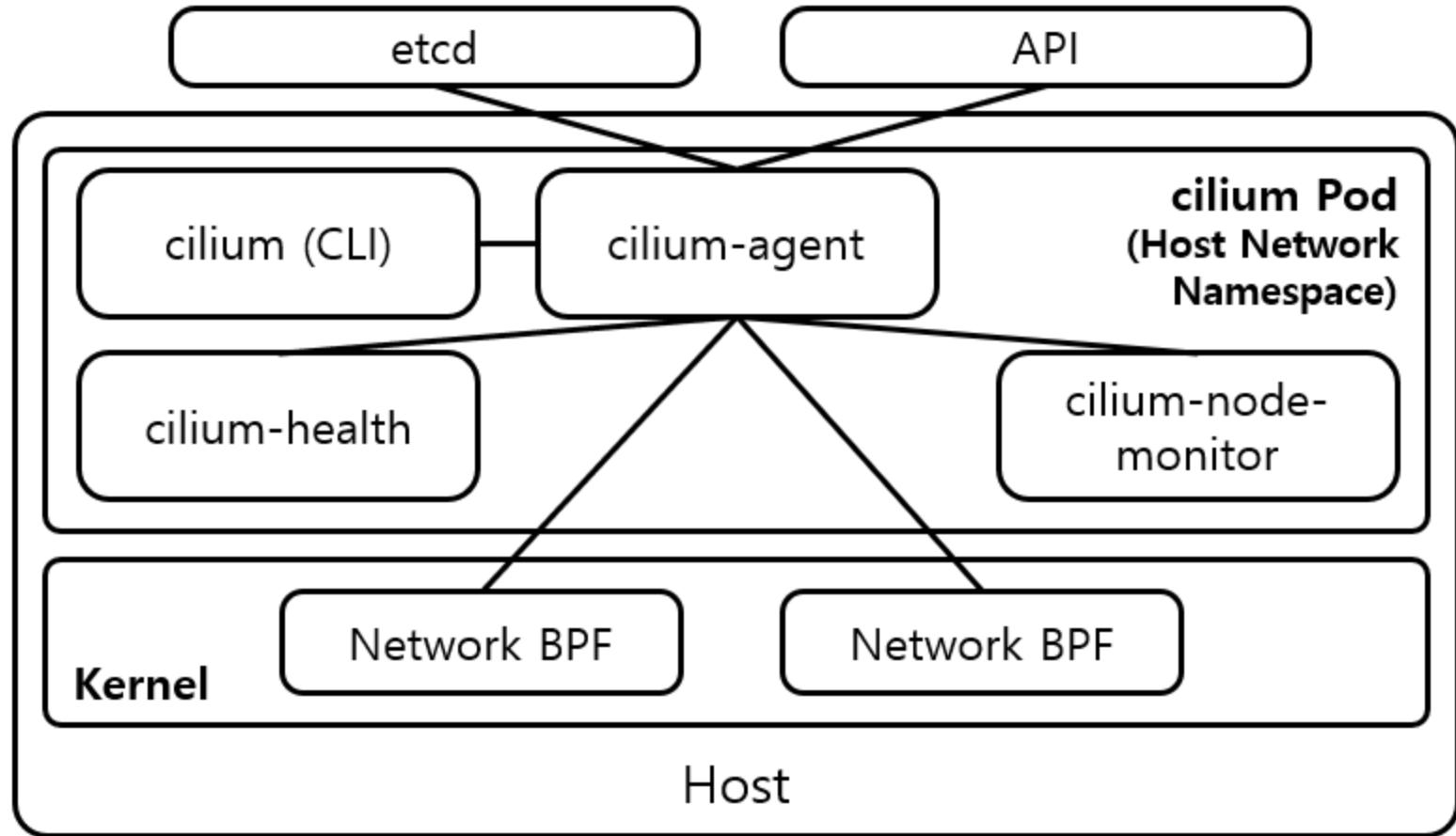
# Iptables Mode



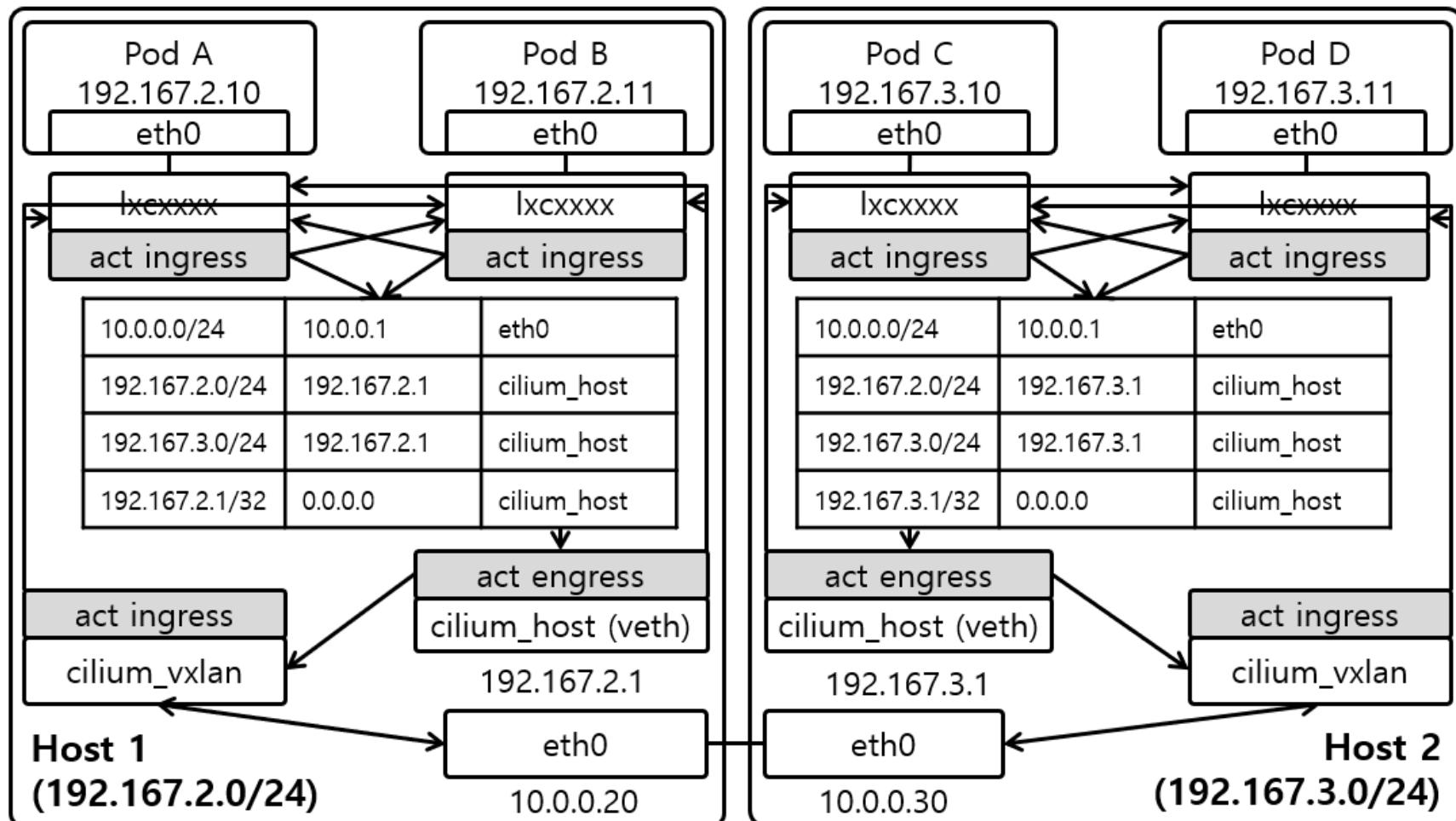
# IPVS Mode



# Cilium Plugin

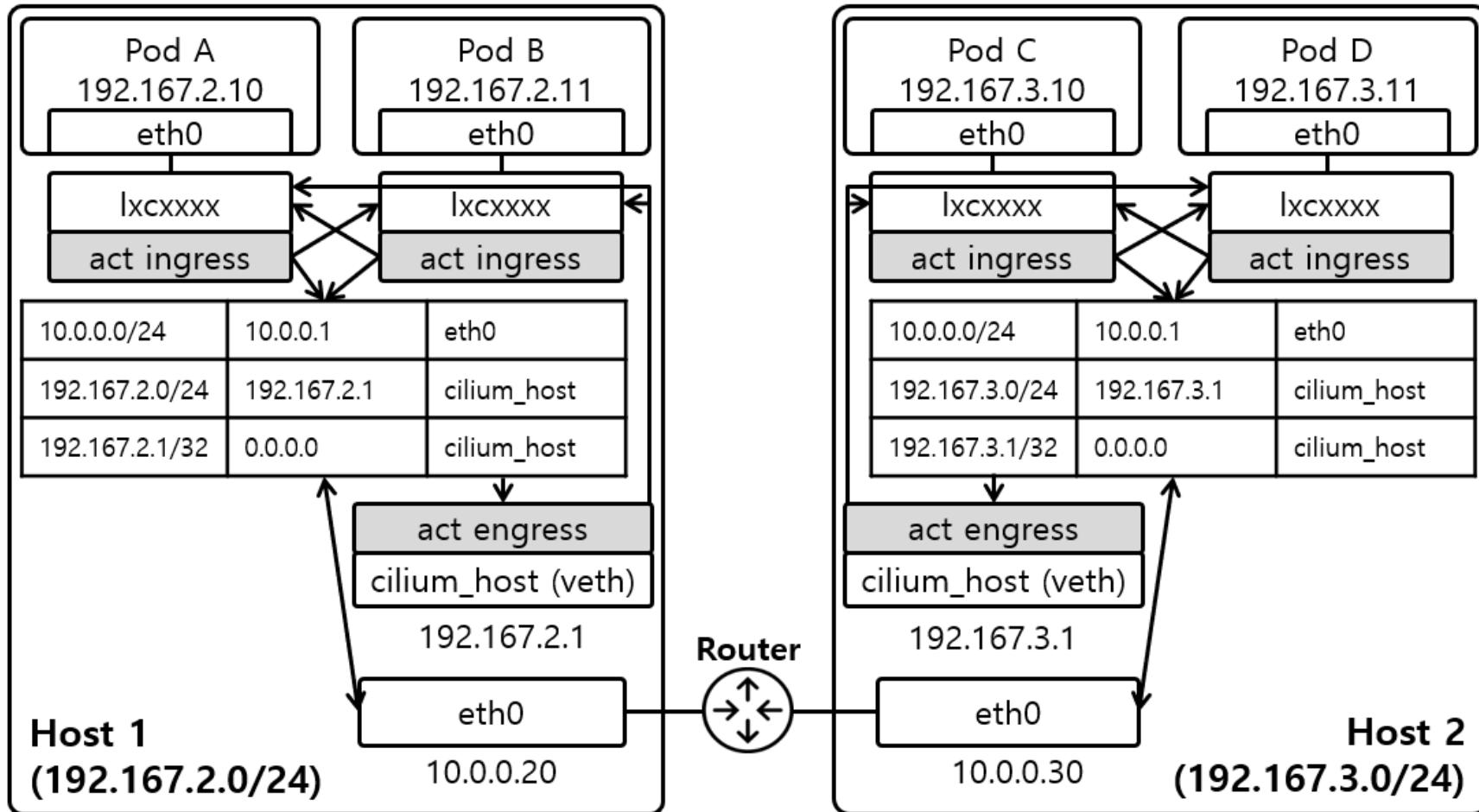


# Cilium Plugin with VXLAN



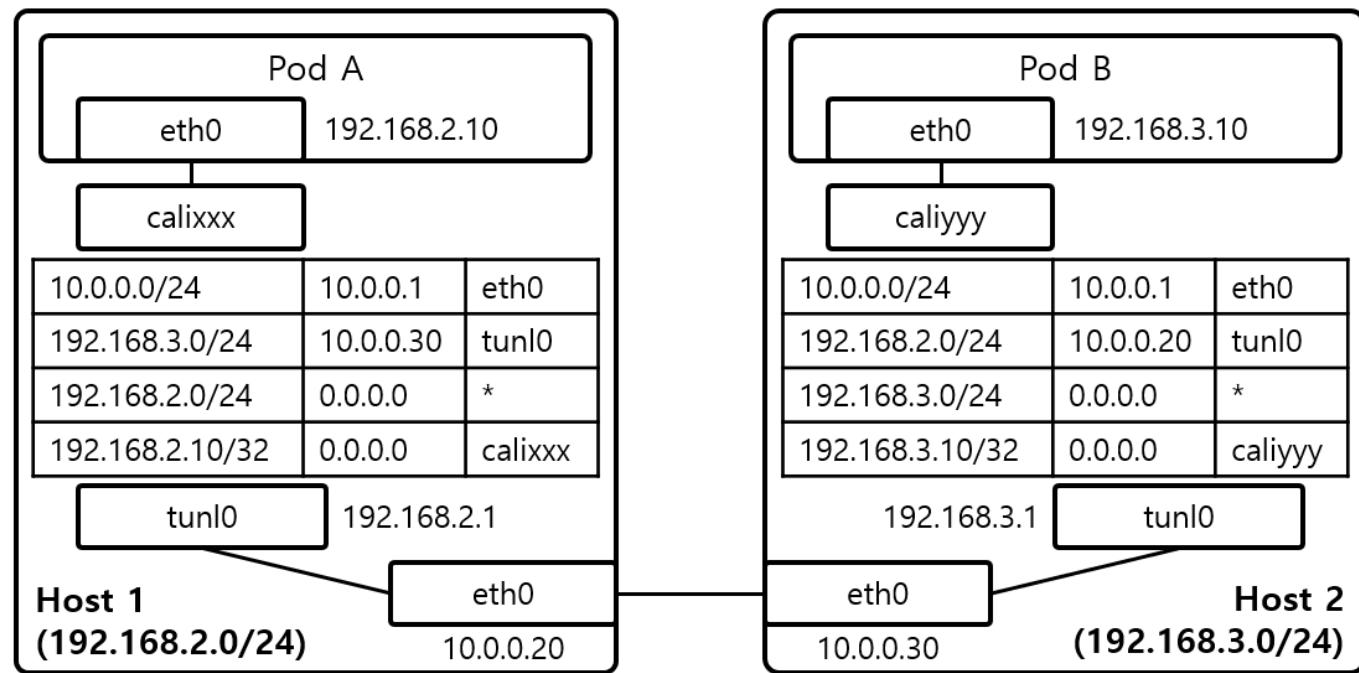
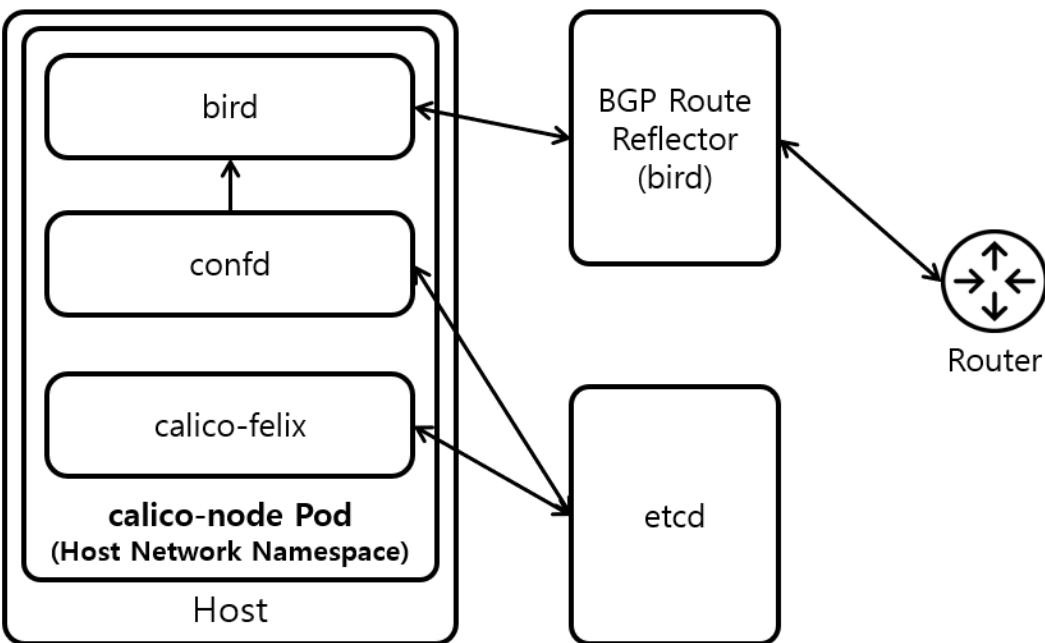
- Host Network : 10.0.0.0/24
- Container Network 192.167.0.0/24

# Cilium Plugin with Host L3



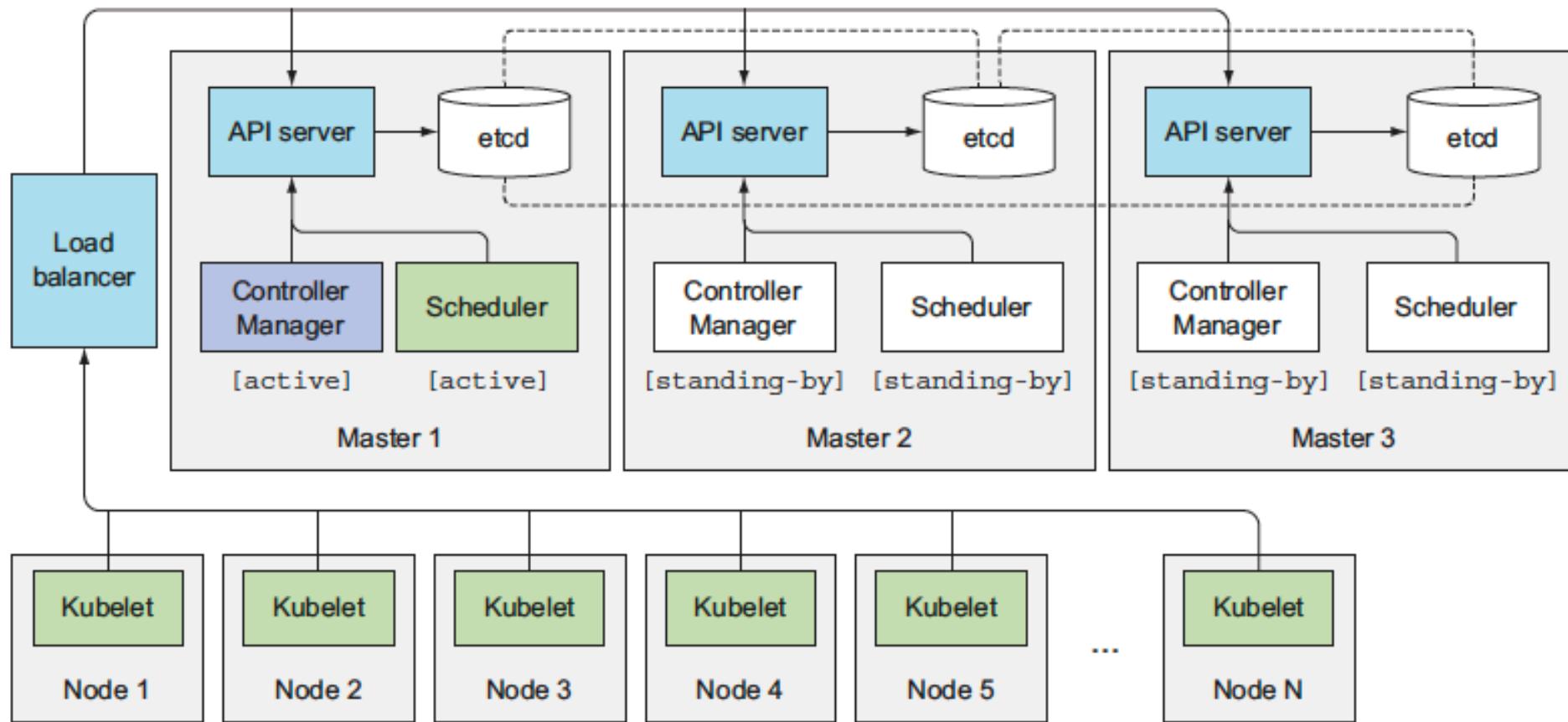
- Host Network : 10.0.0.0/24
- Container Network 192.167.0.0/24

# Calico Plugin

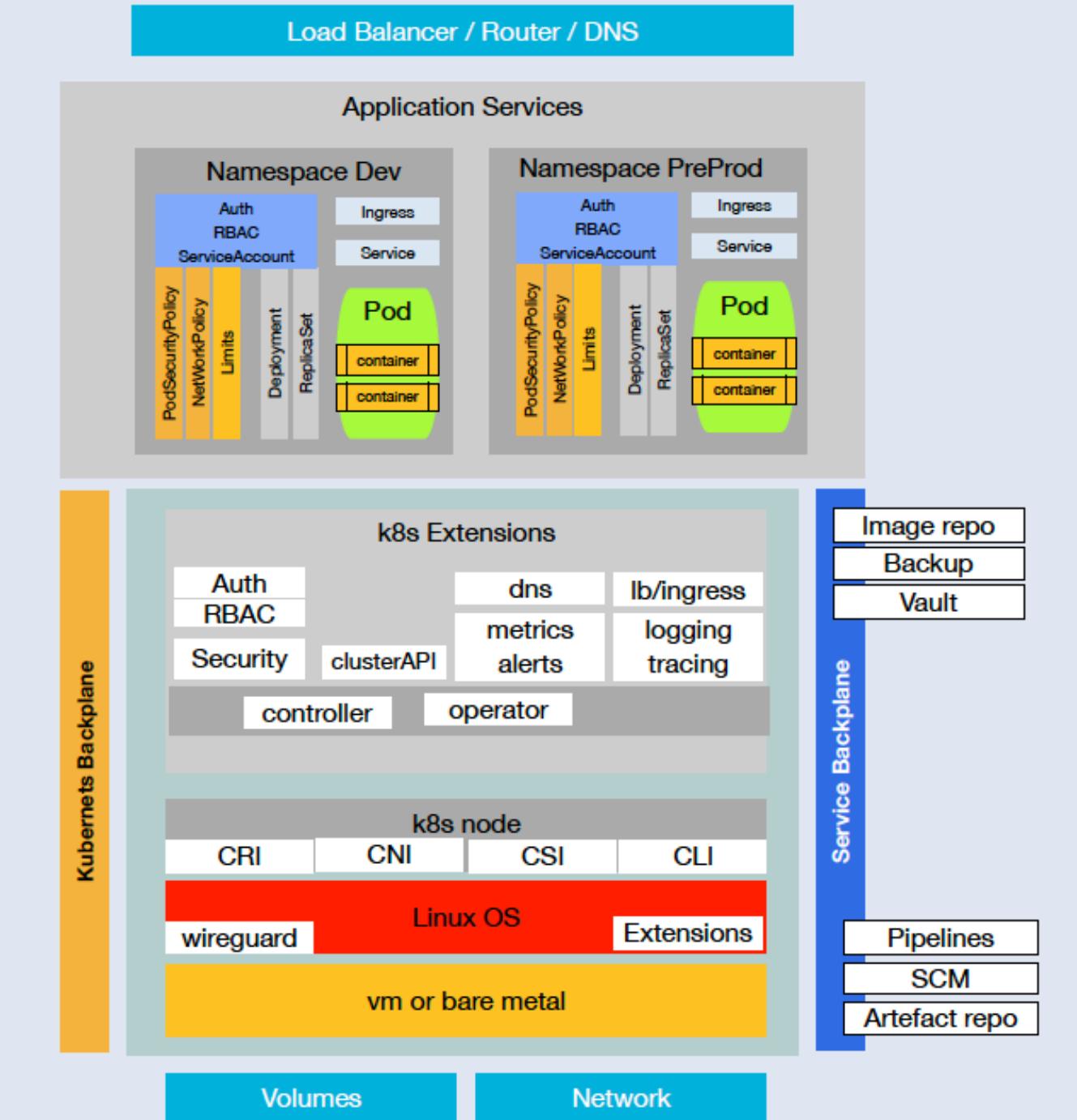


- Host Network : 10.0.0.0/24
- Pod Network 192.168.0.0/24

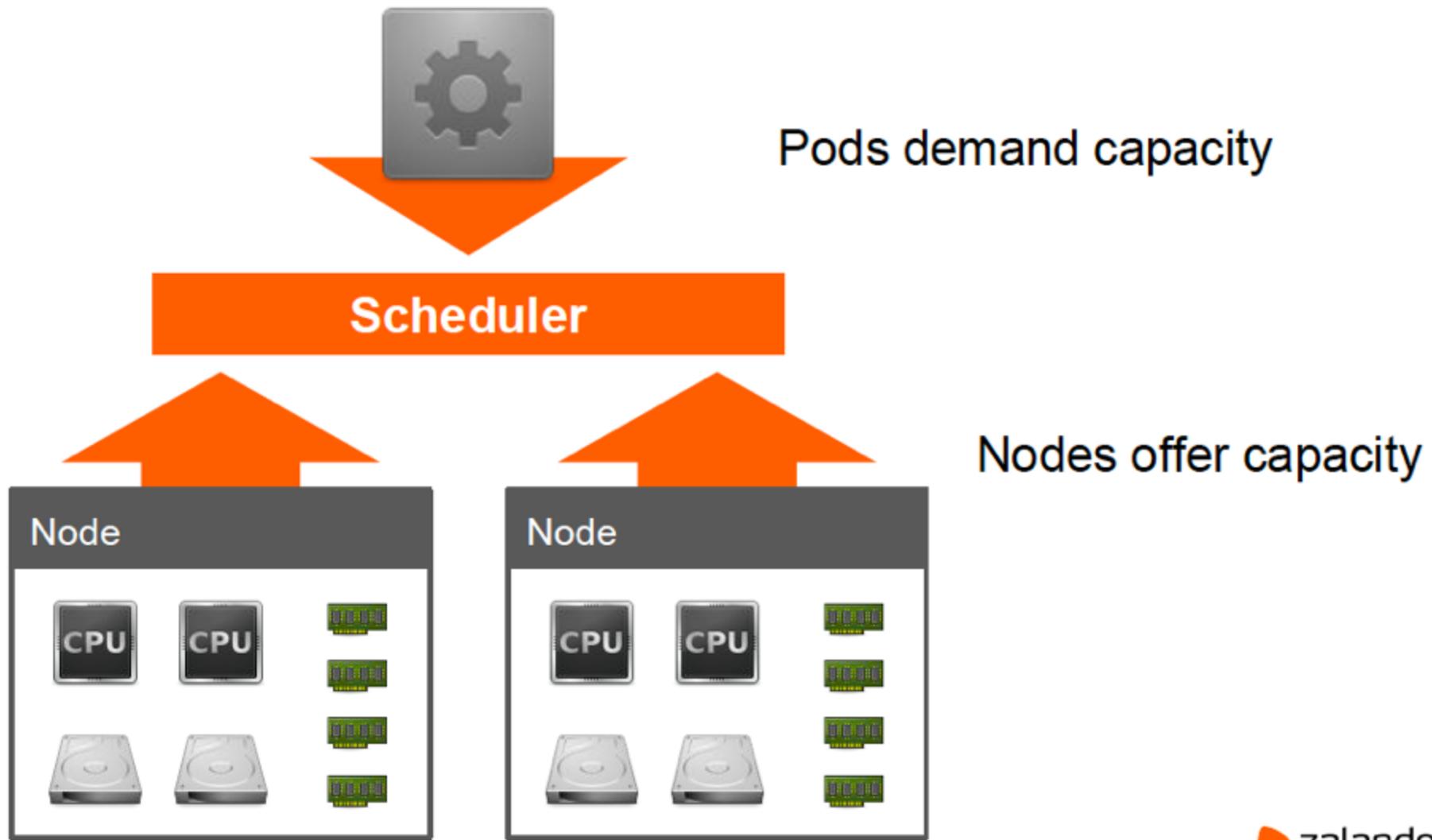
# Highly Available Kubernetes Cluster



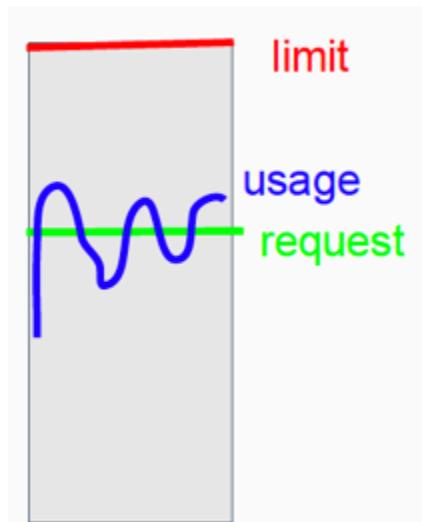
# k8s Stacking



# KUBERNETES: IT'S ALL ABOUT RESOURCES

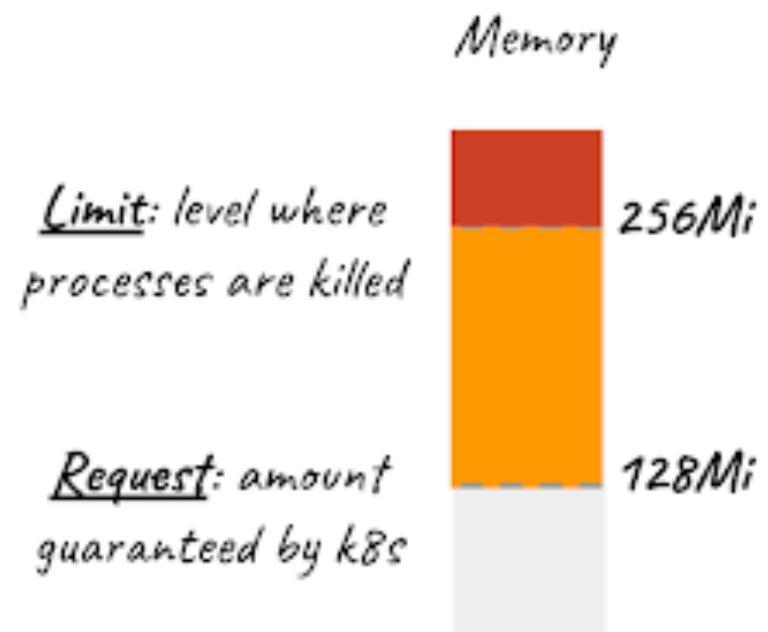
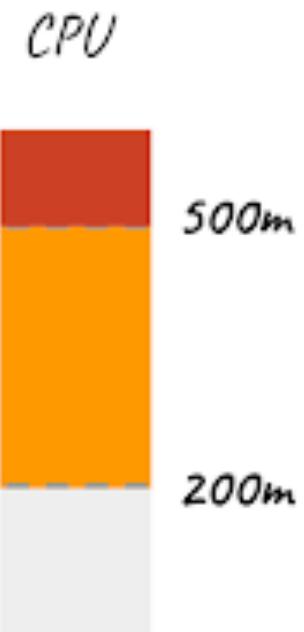


# Request & Limit



Limit: where CPU throttling begins

Request: amount guaranteed by k8s

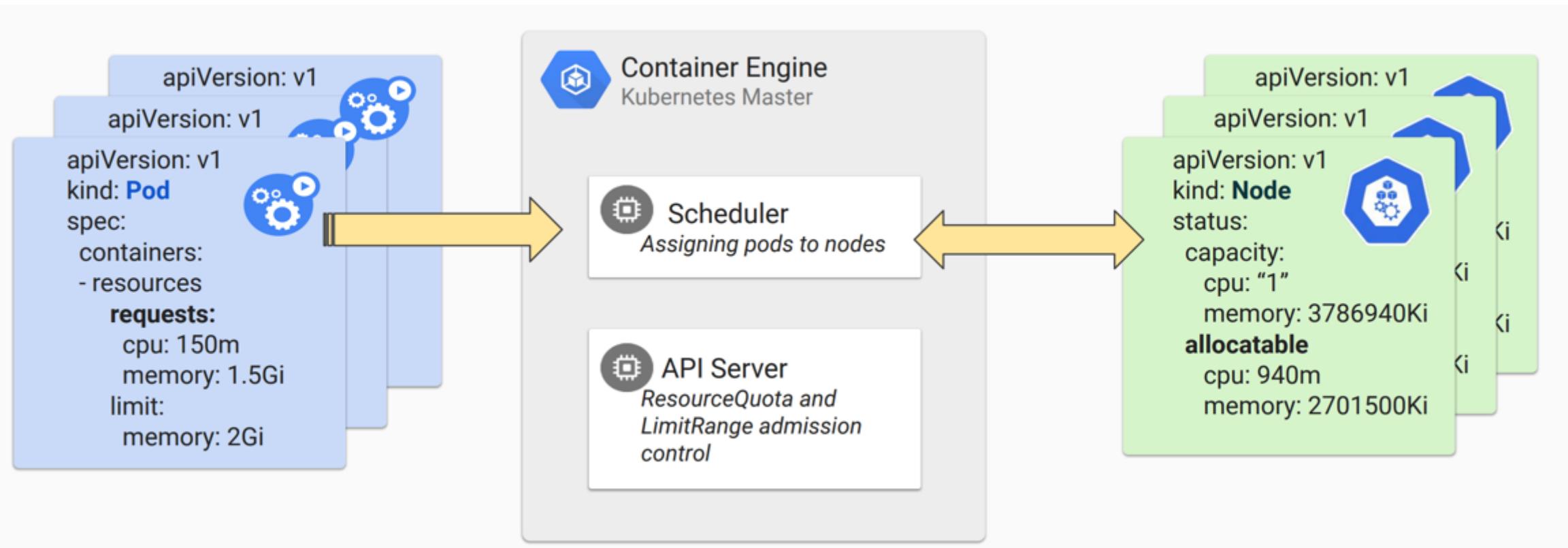


# Requests, Limits by cgroup

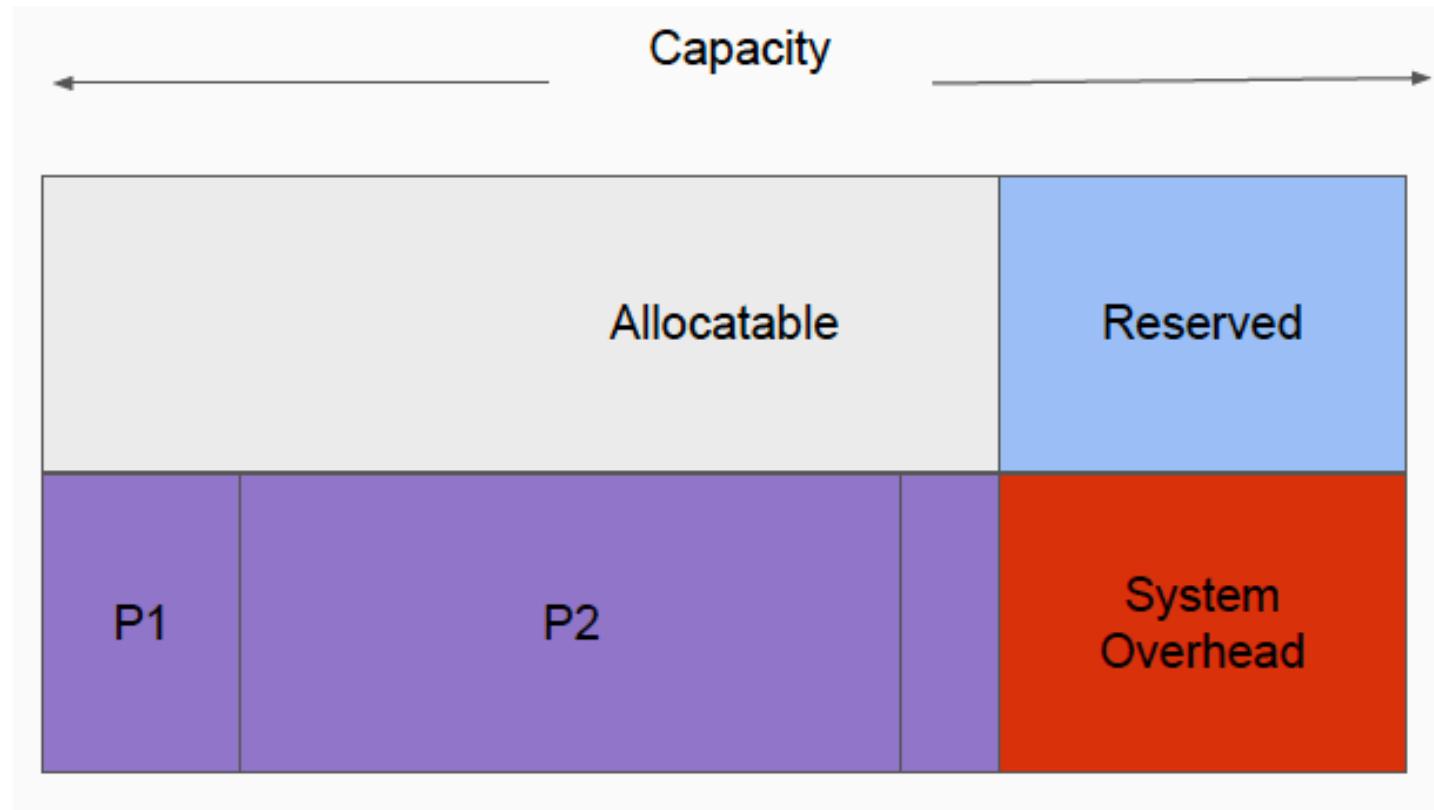
```
docker run --requests(cpu=10m/5m) ..sha512()..  
  
cat /sys/fs/cgroup/cpu/kubepods/burstable/pod5d5..0d/cpu.shares  
10      // relative share of CPU time  
  
cat /sys/fs/cgroup/cpu/kubepods/burstable/pod6e0..0d/cpu.shares  
5       // relative share of CPU time  
  
cat /sys/fs/cgroup/cpuacct/kubepods/burstable/pod5d5..0d/cpuacct.usage  
/sys/fs/cgroup/cpuacct/kubepods/burstable/pod6e0..0d/cpuacct.usage  
13432815283 // total CPU time in nanoseconds  
7528759332  // total CPU time in nanoseconds
```

```
docker run --cpus 1 -m 200m --rm -it busybox  
  
cat /sys/fs/cgroup/cpu/docker/8ab25..1c/cpu.{shares,cfs_*}  
1024    // cpu.shares (default value)  
100000  // cpu.cfs_period_us (100ms period length)  
100000  // cpu.cfs_quota_us (total CPU time in µs consumable per period)  
  
cat /sys/fs/cgroup/memory/docker/8ab25..1c/memory.limit_in_bytes  
209715200
```

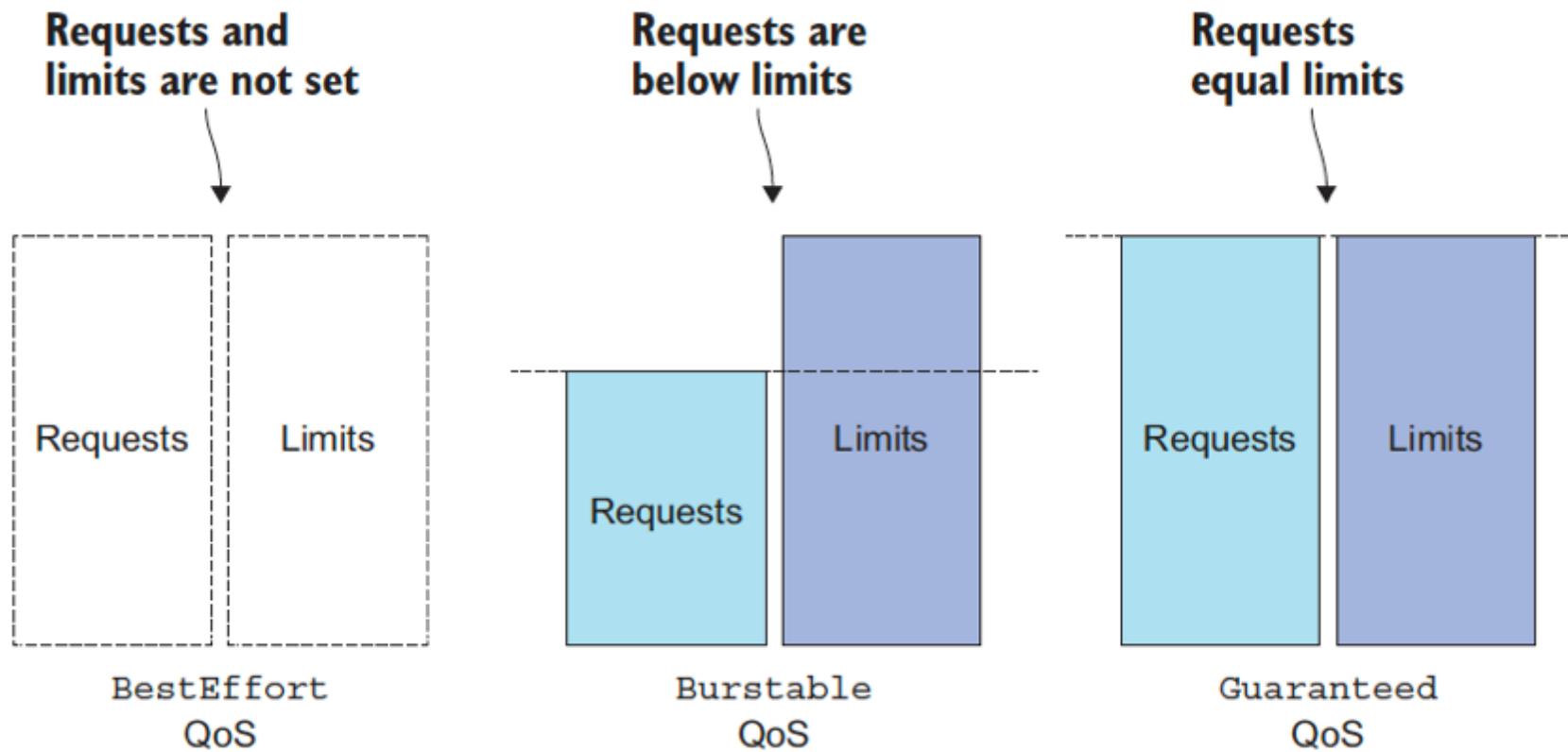
# Scheduler & API Server

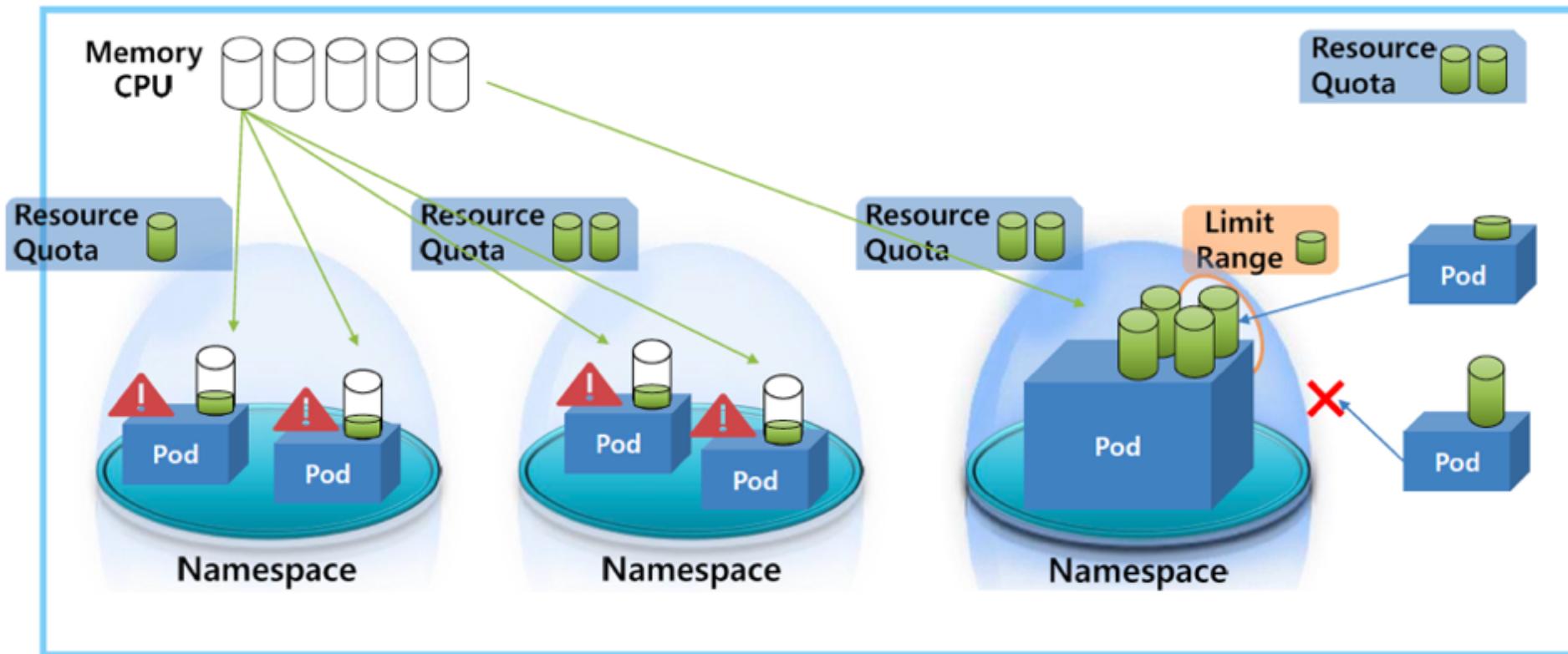
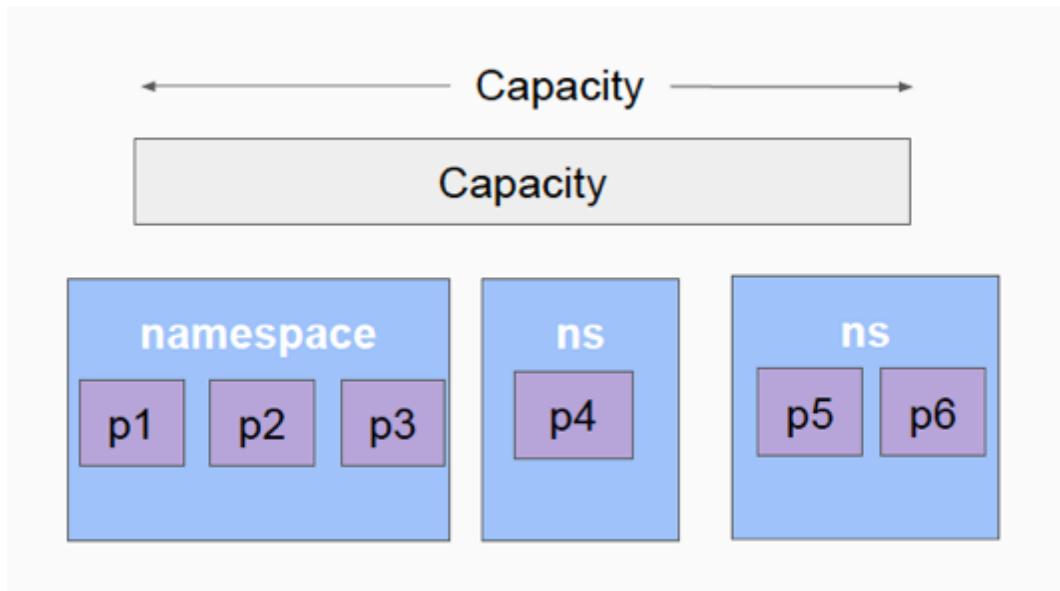


# Capacity & Allocatable

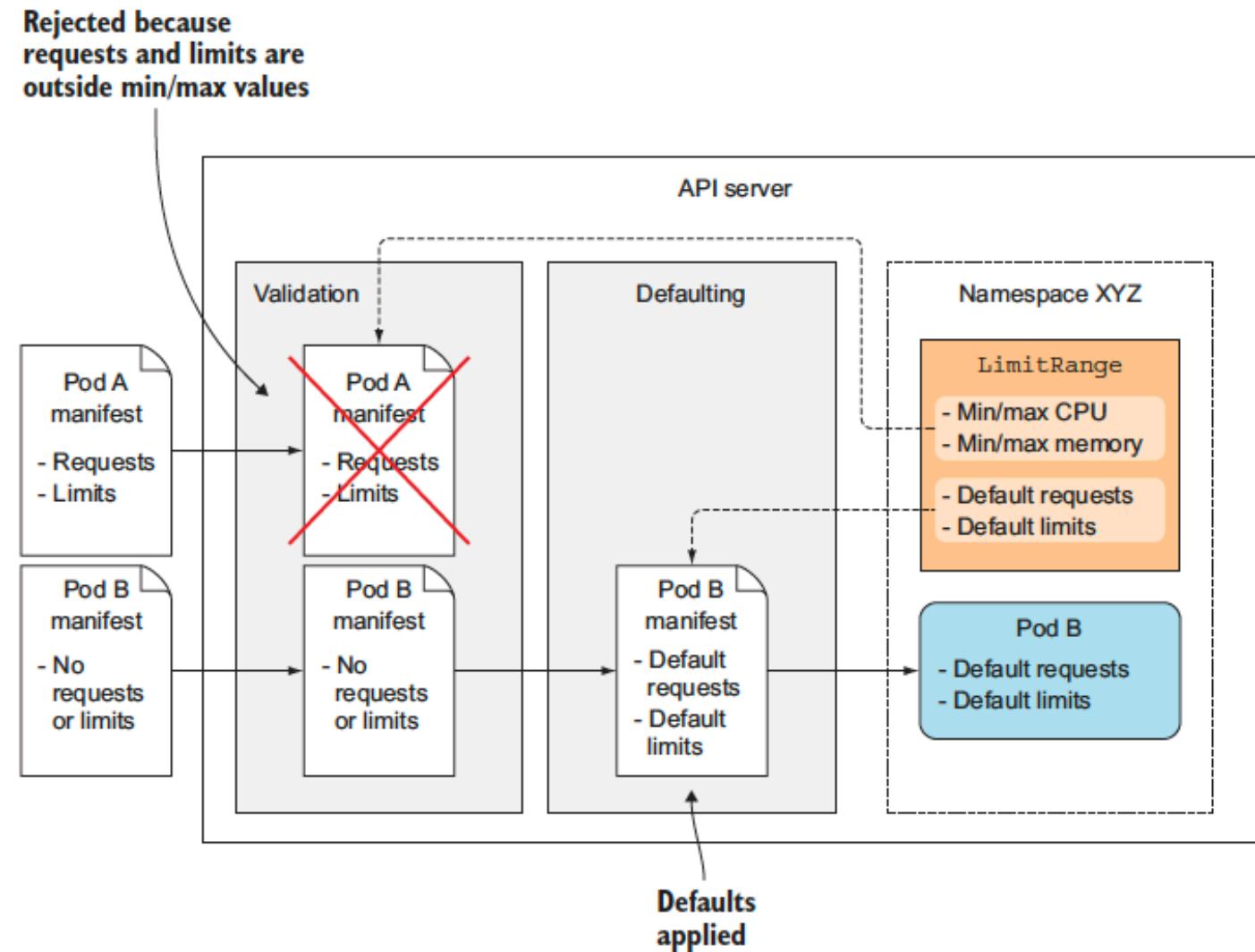


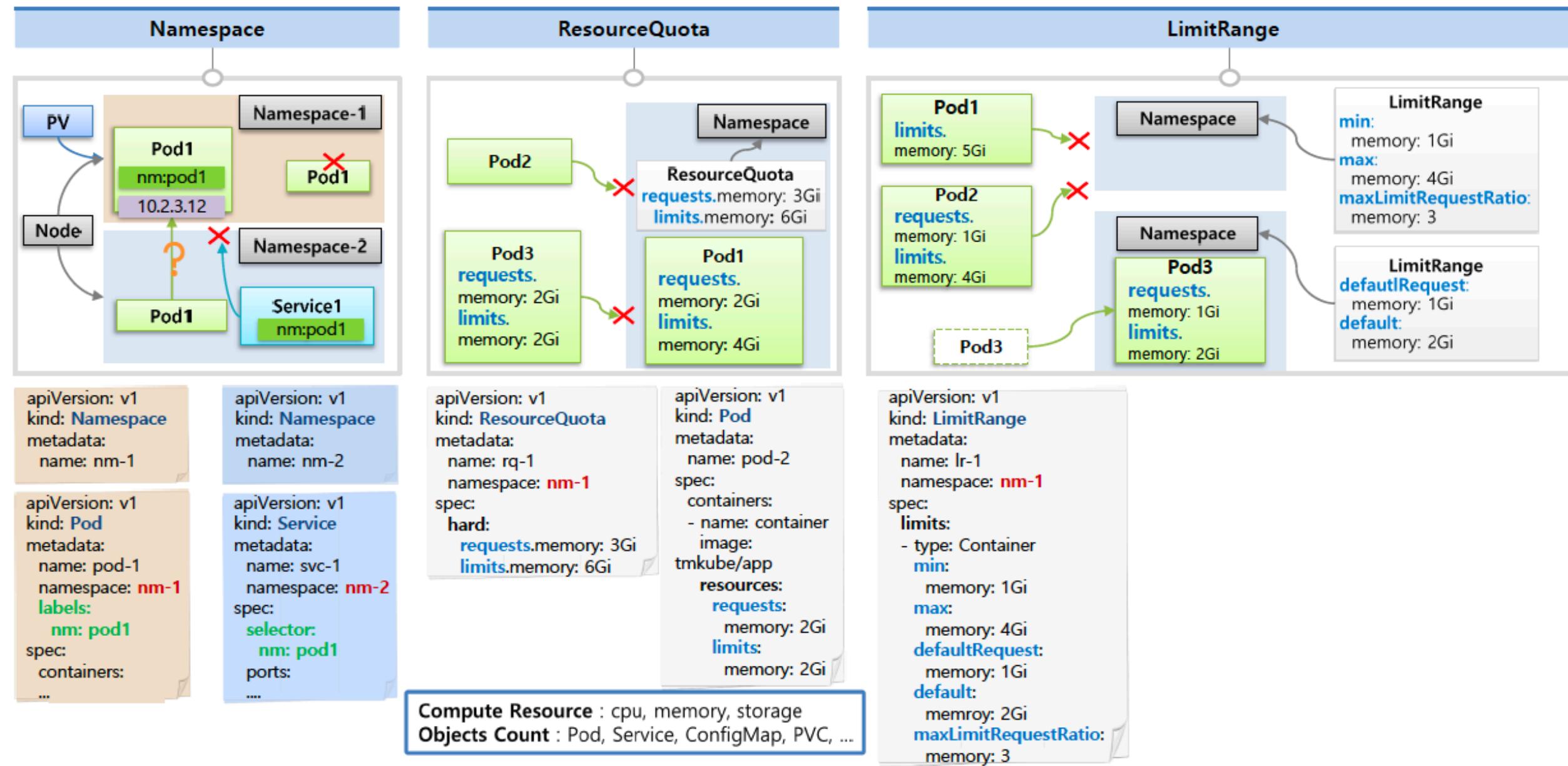
# QoS: BestEffort, Burstable, Guaranteed

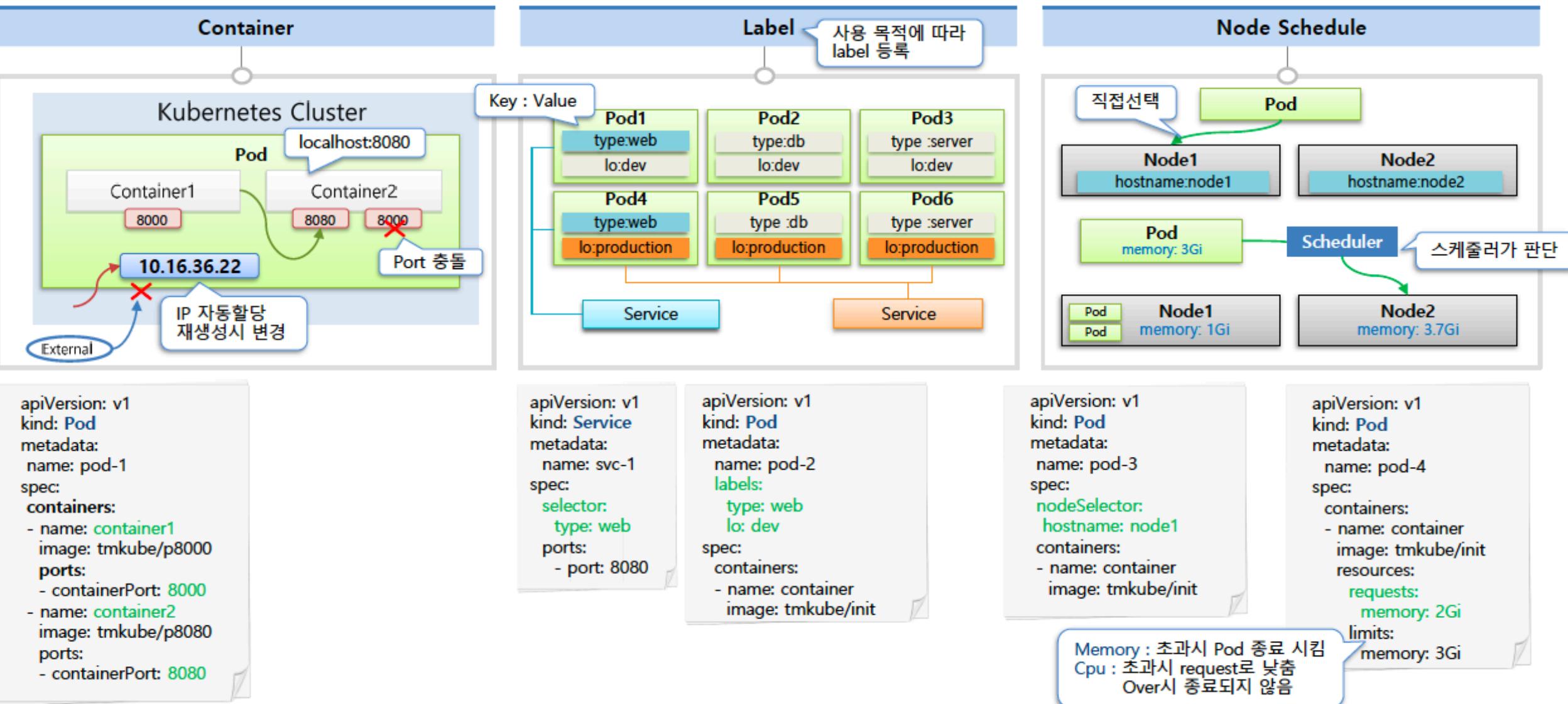




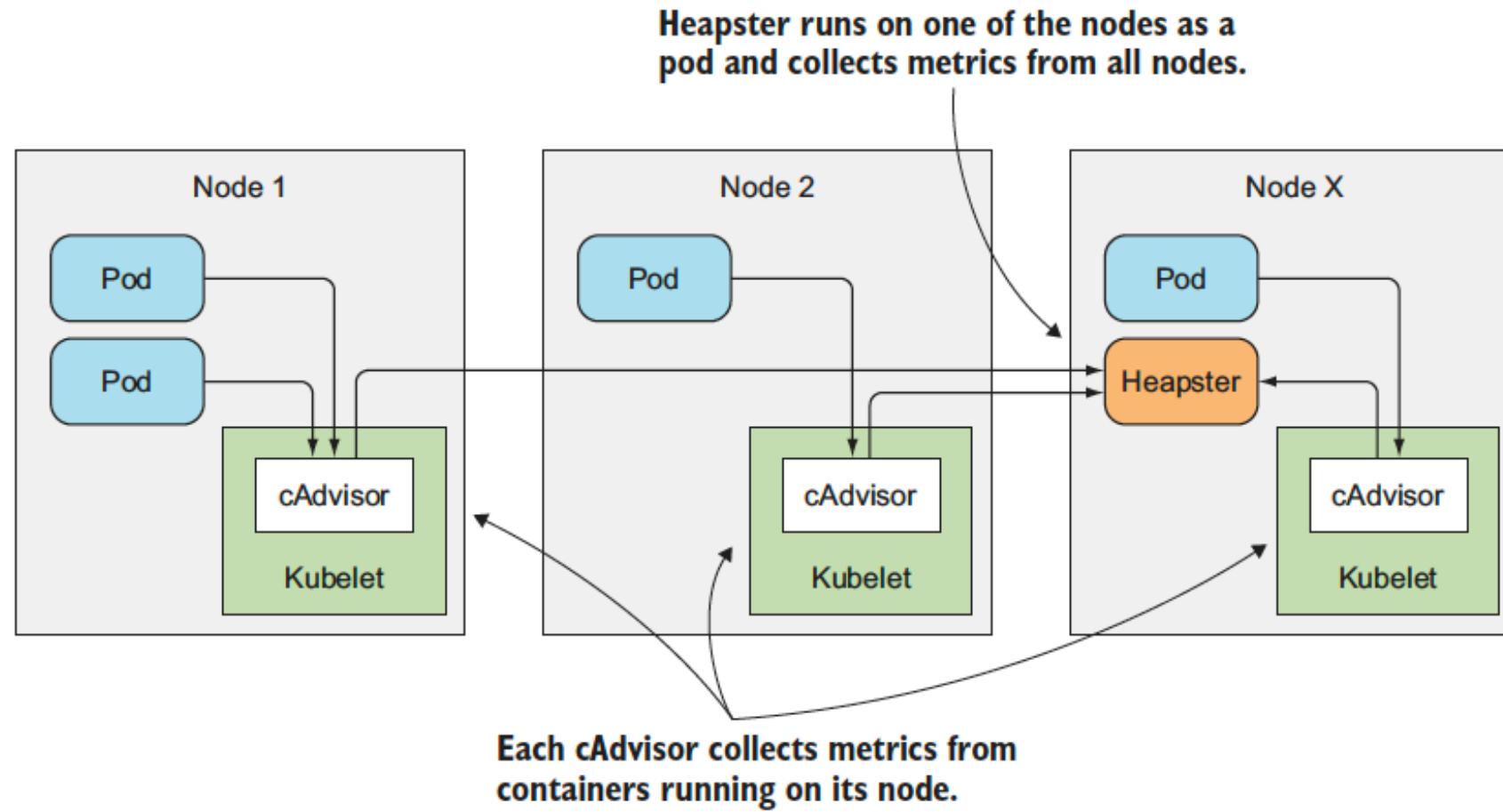
# LimitRange: Admission Control



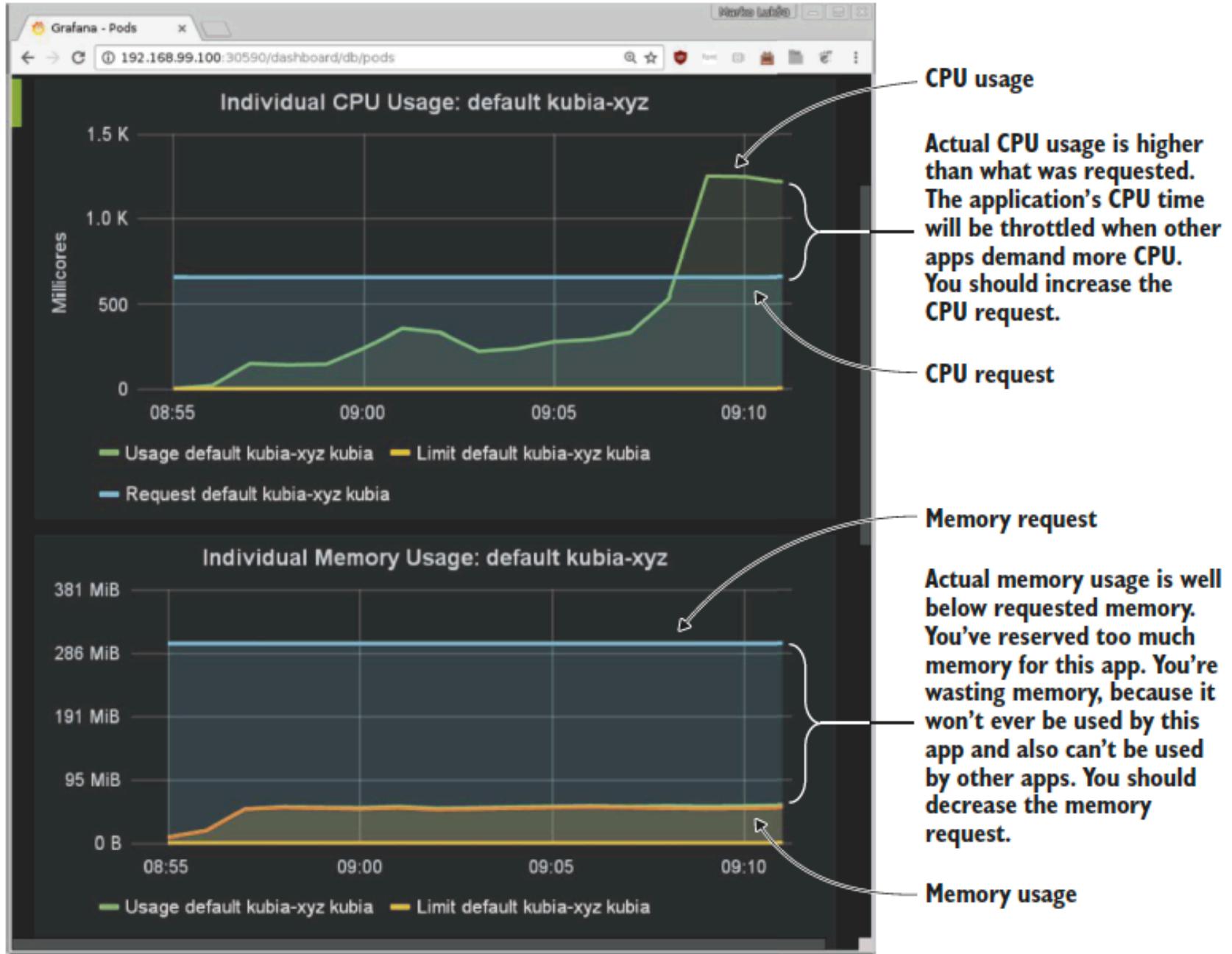




# cAdvisor & Heapster

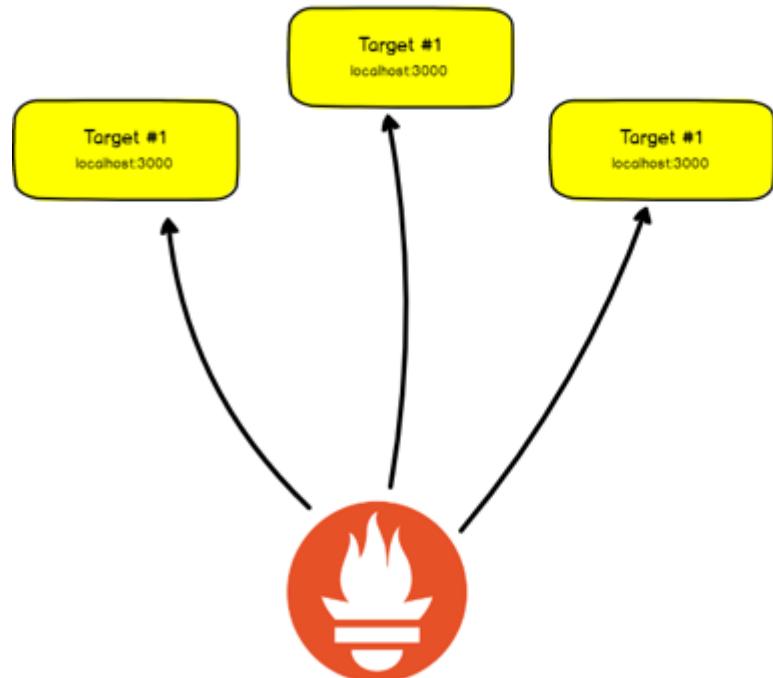


# Grafana

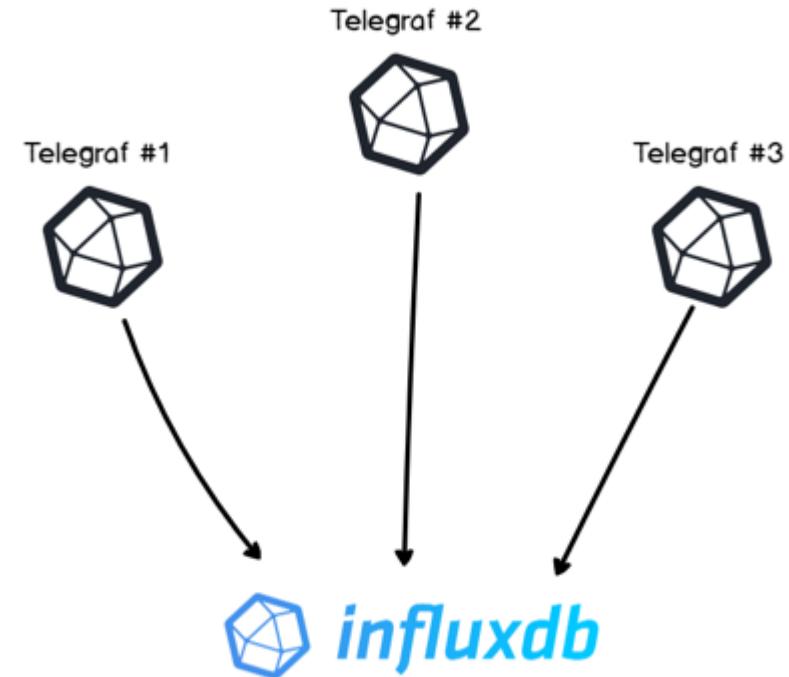


# Prometheus vs influxdb

## Push vs Pull

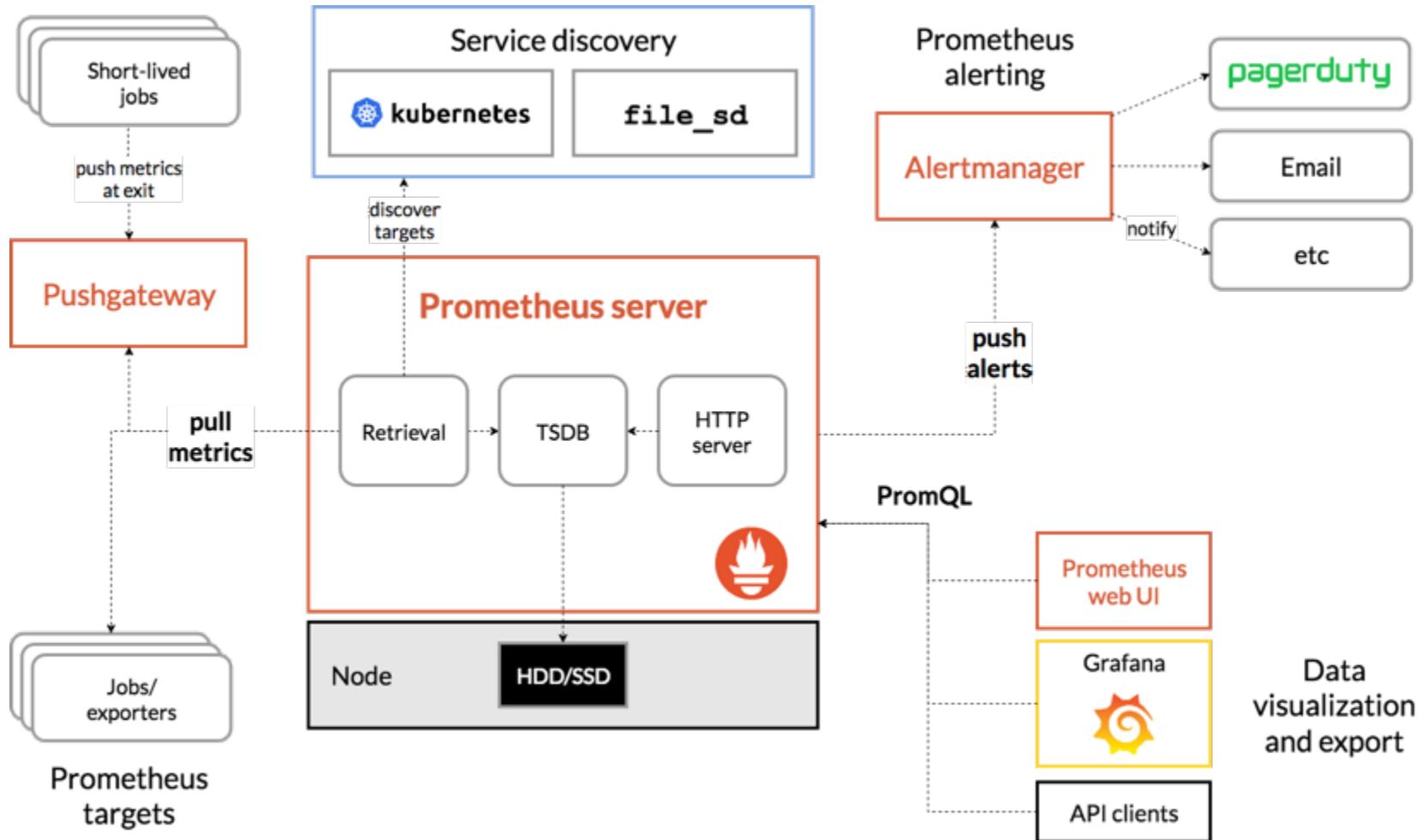


Prometheus pull data every 10 seconds  
Prometheus is **active**



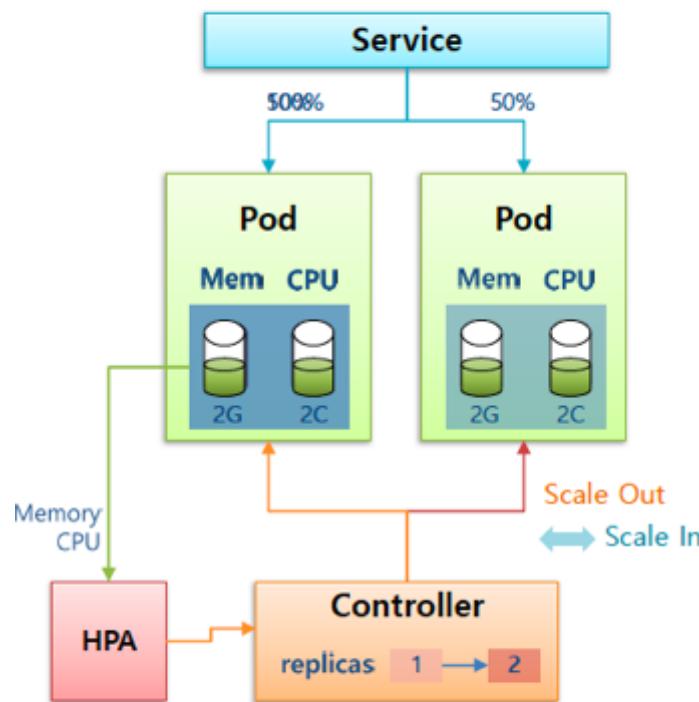
Individual instances push data every 10 seconds  
InfluxDB is **passive**

# Prometheus

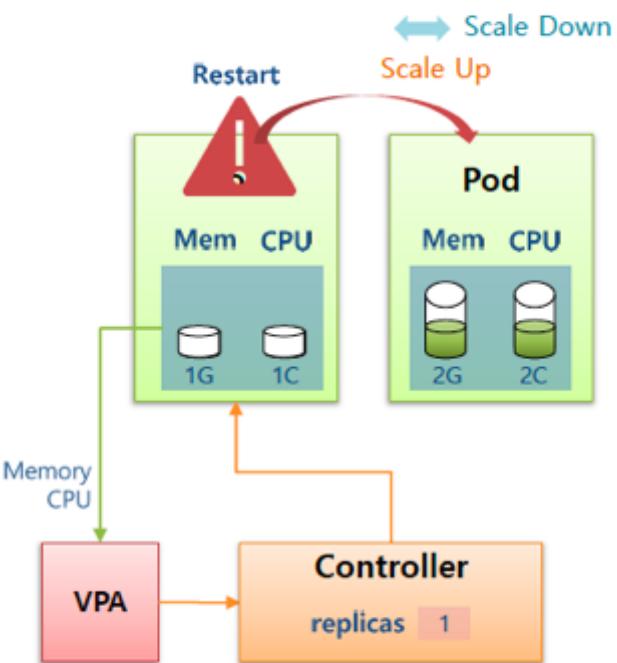




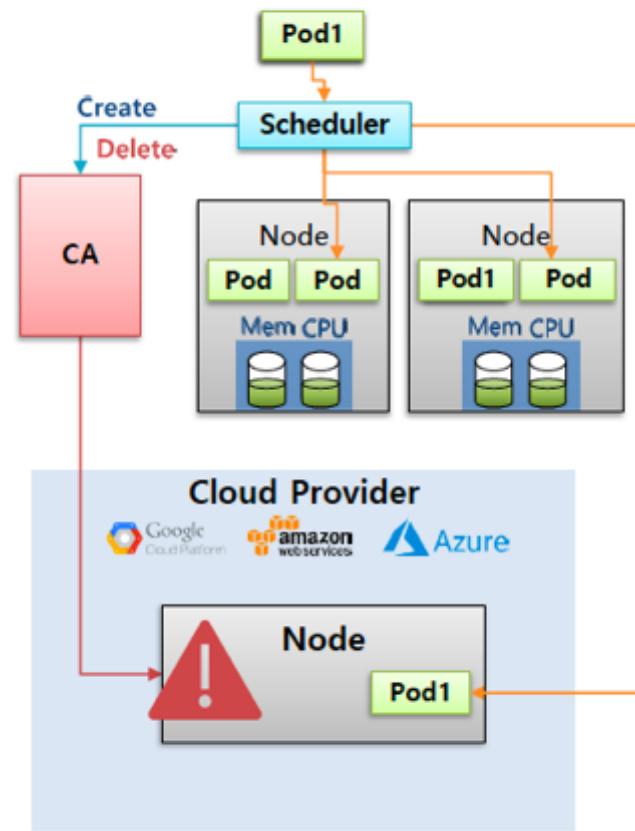
## HPA (Horizontal Pod Autoscaler)



## VPA (Vertical Pod Autoscaler)

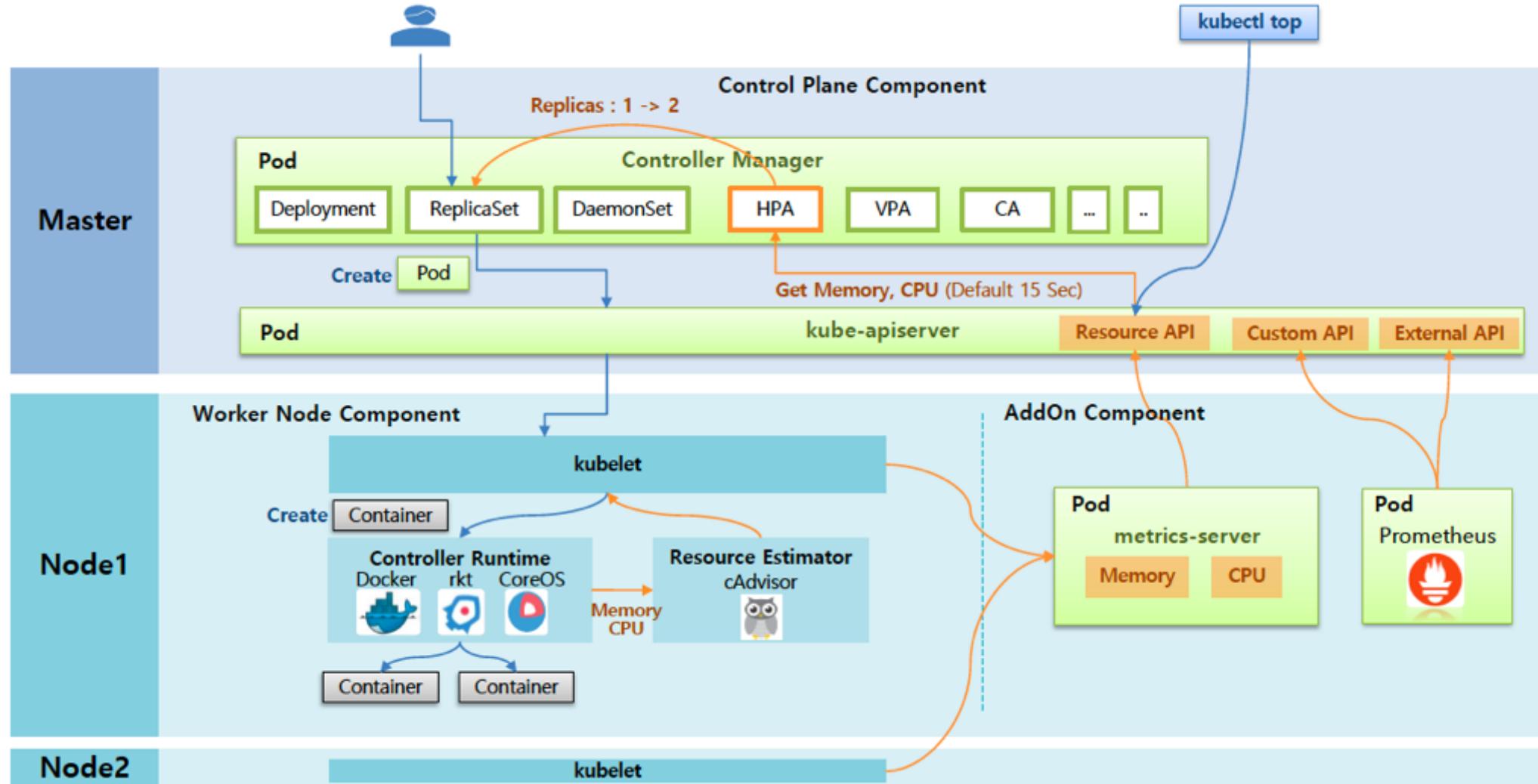


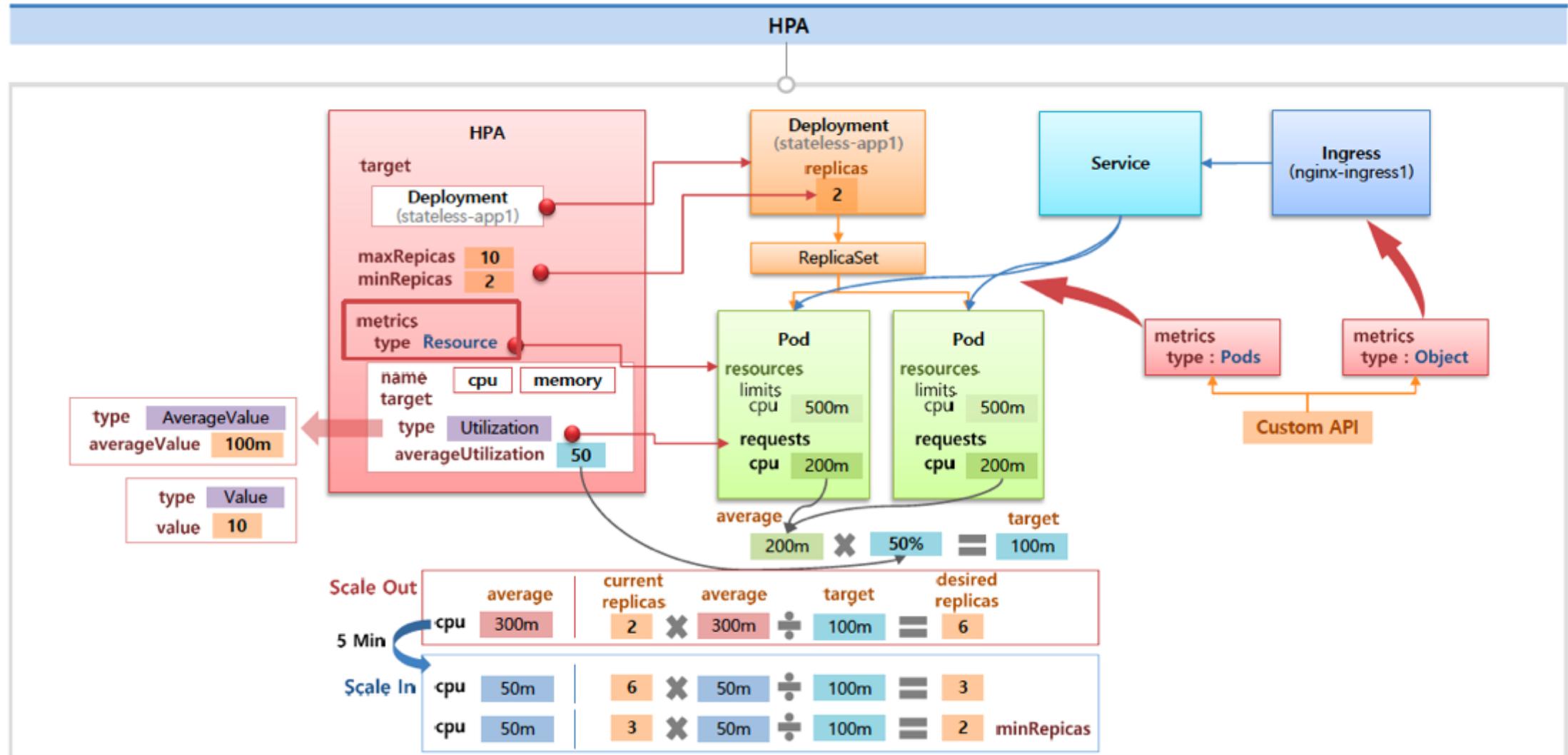
## CA (Cluster Autoscaler)



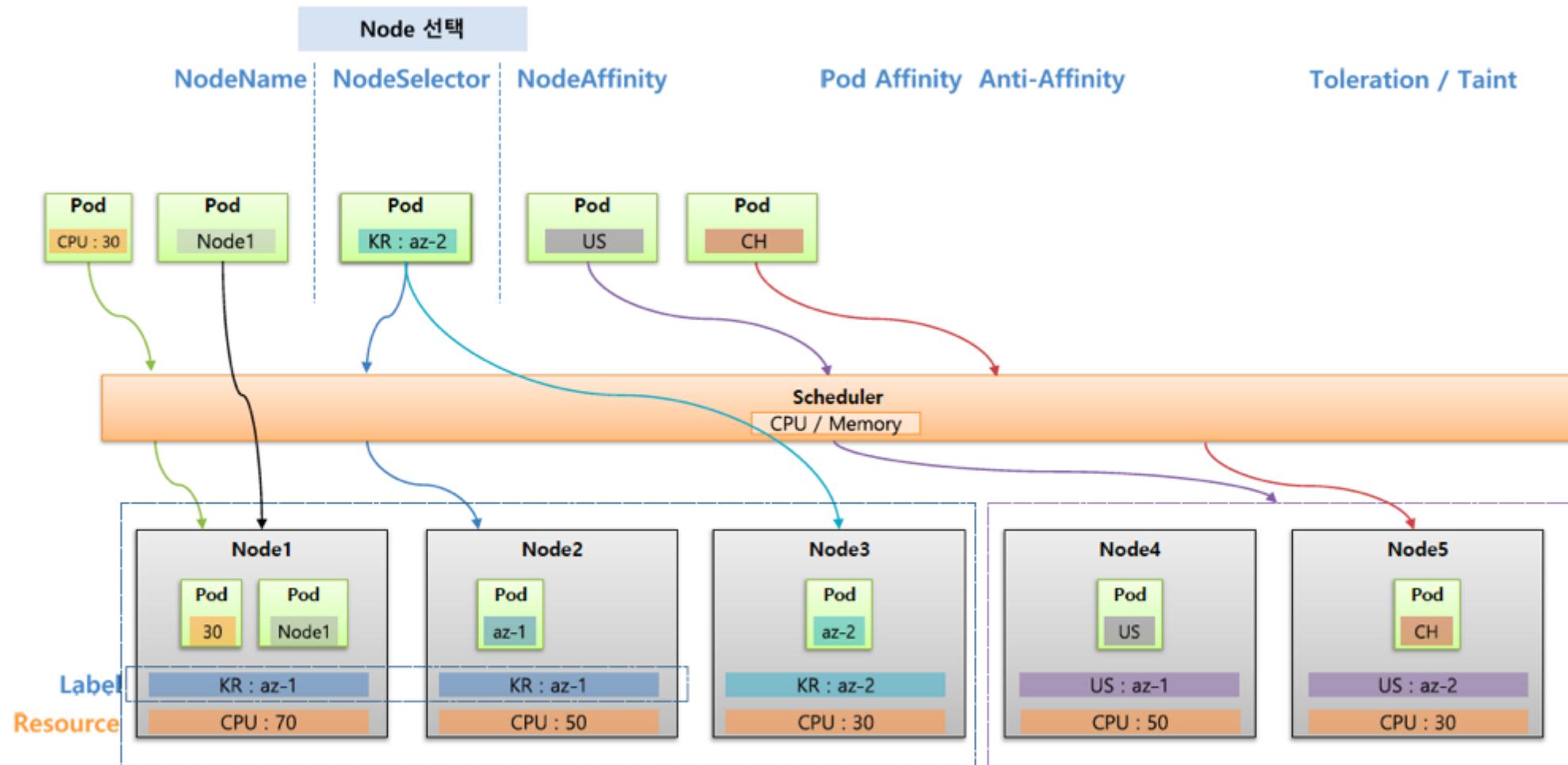
- 기동이 빠르게 되는 App
- Stateless App

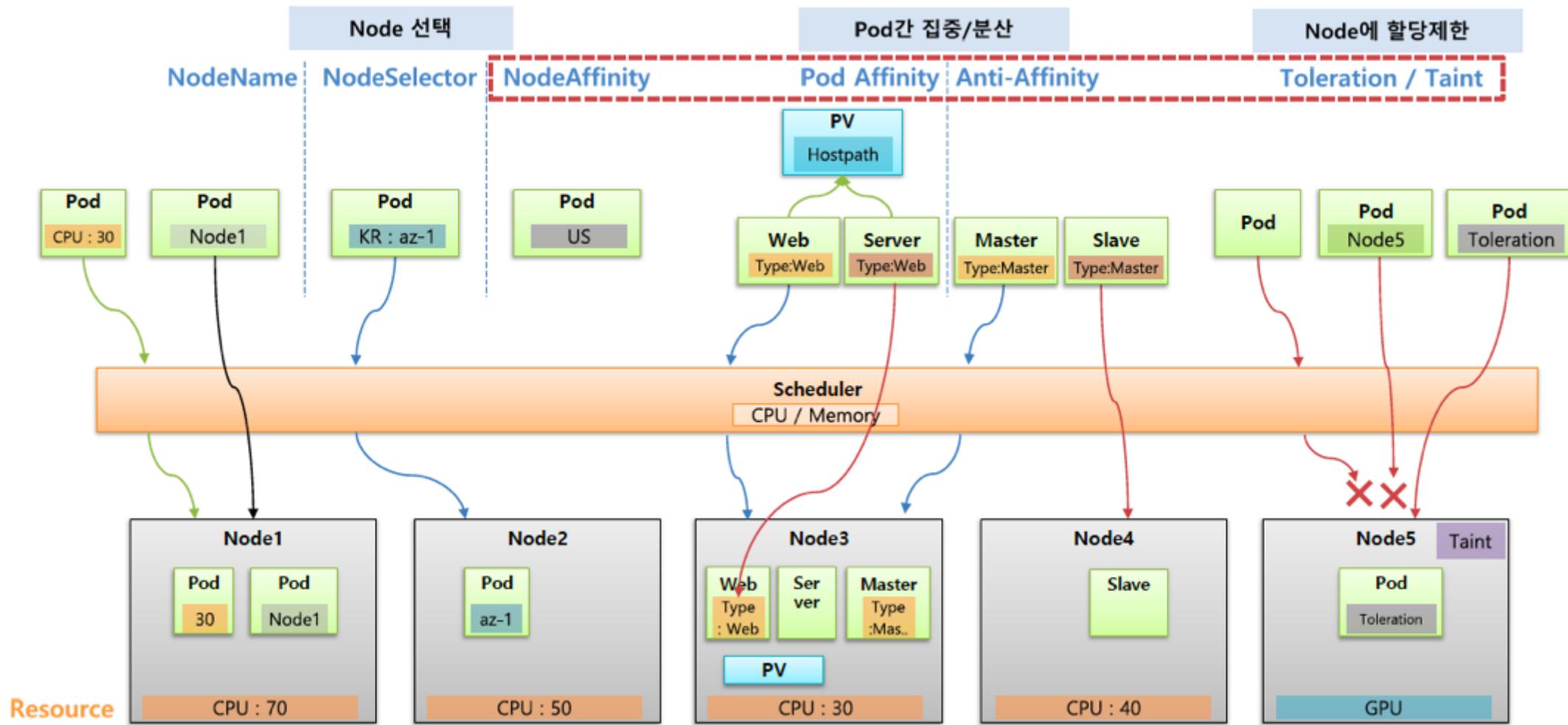
- Stateful App
- 한 Controller에 HPA와 함께 사용안됨





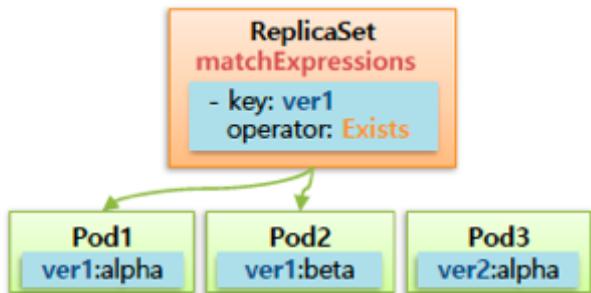




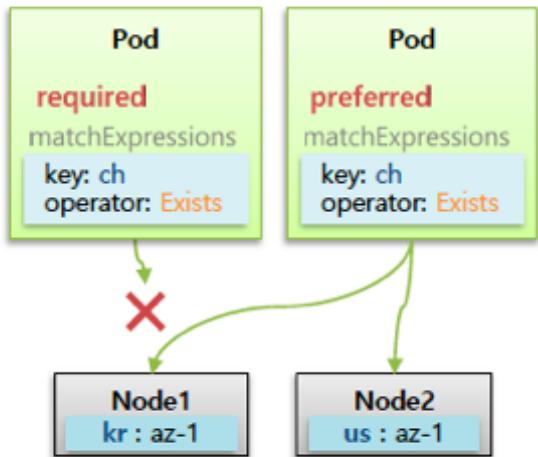


## Node Affinity

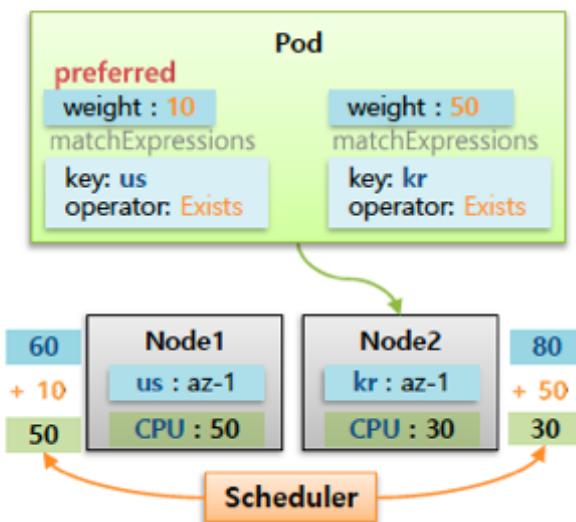
### matchExpressions



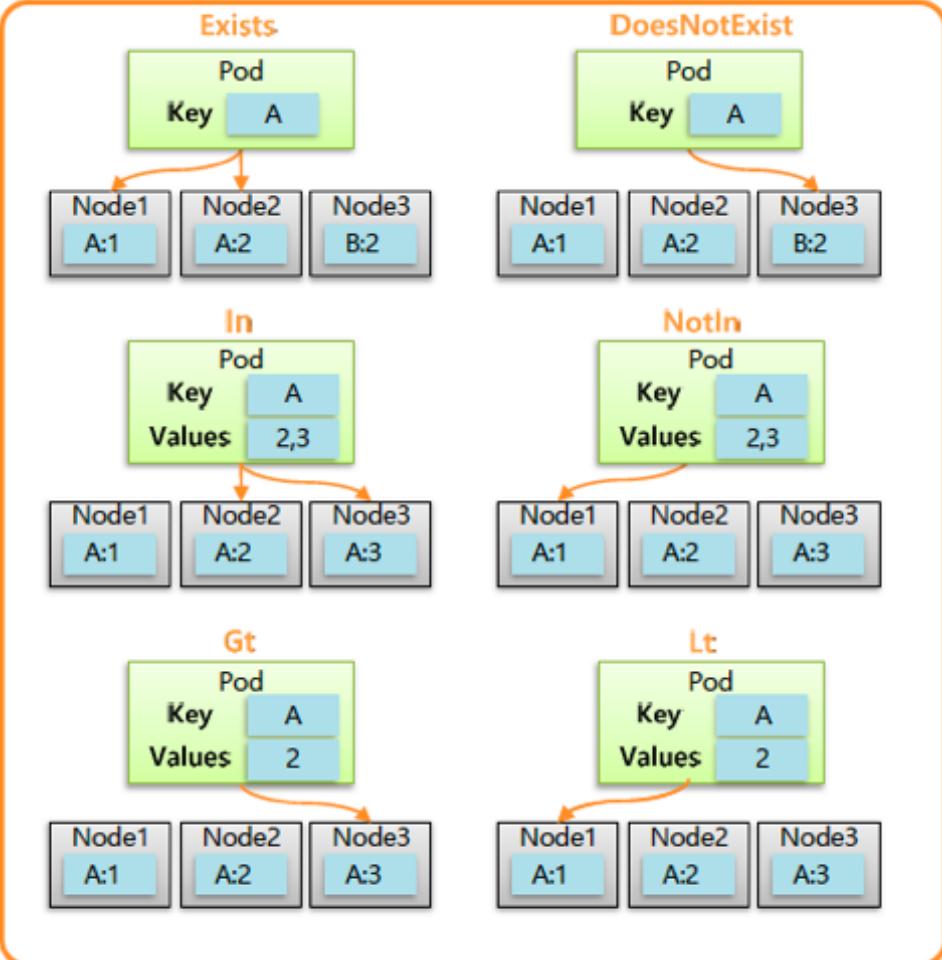
### required vs preferred

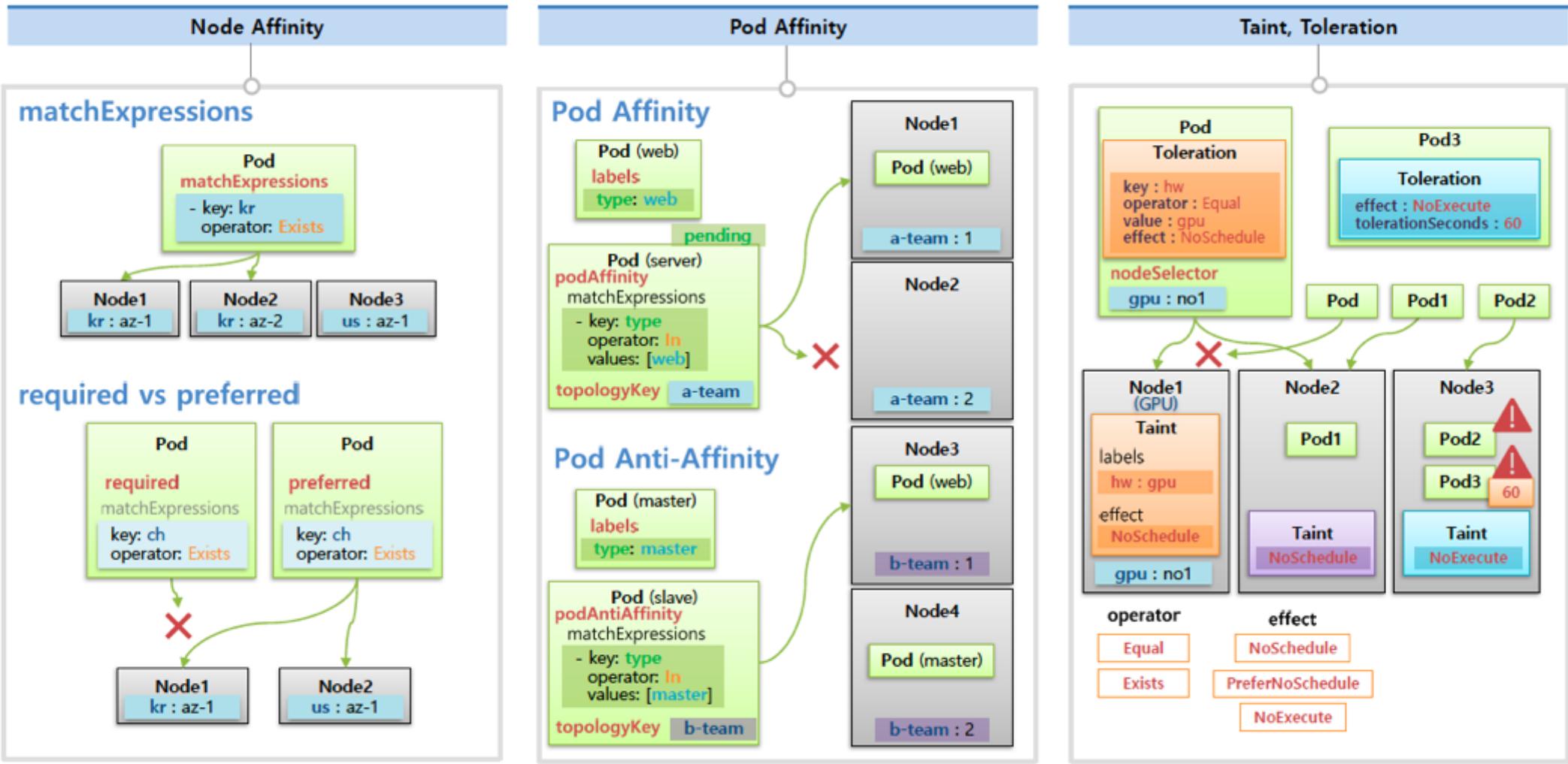


### preferred weight



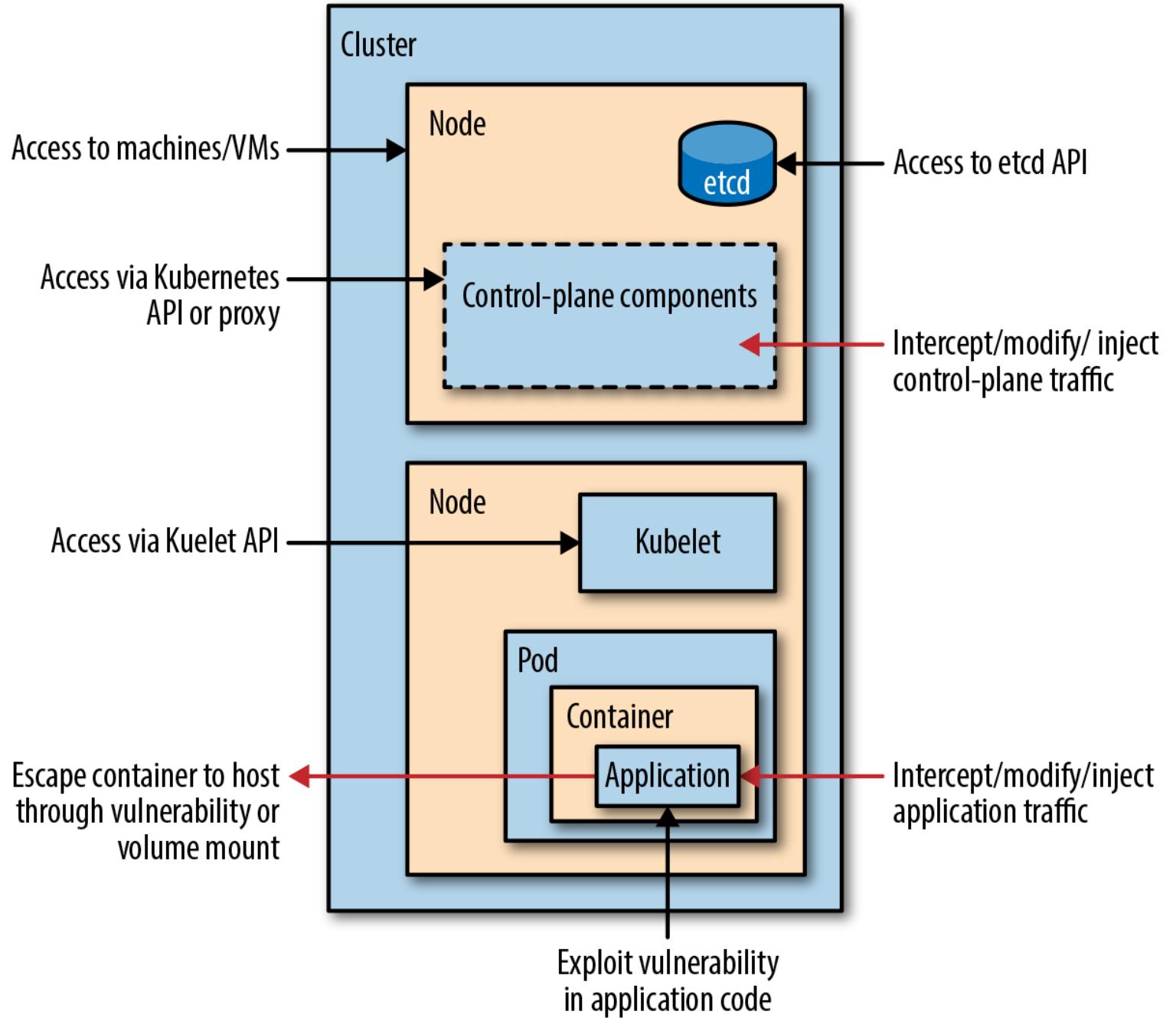
## operator





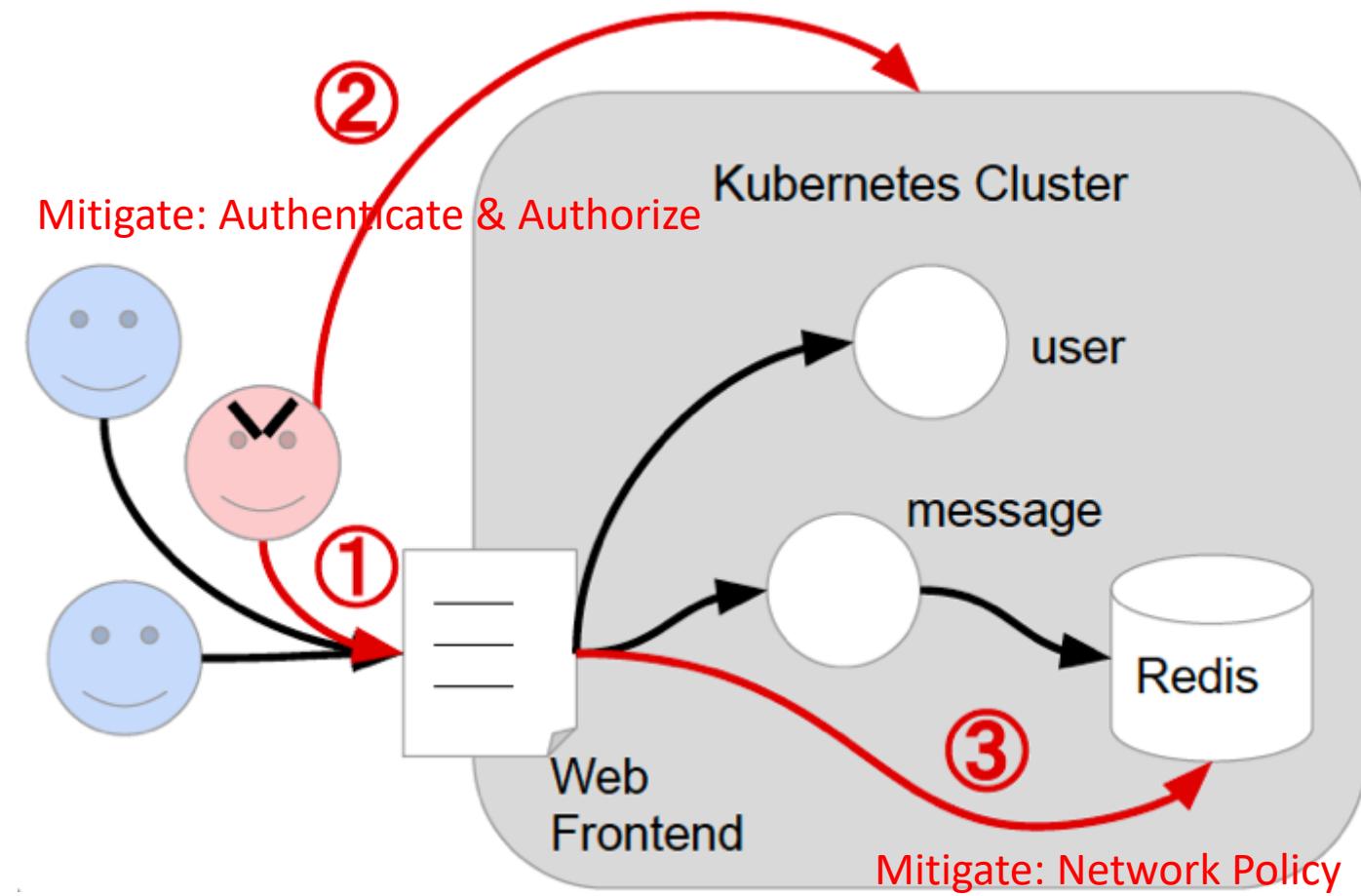


# Attack Surface



# Attack Surface: API Server

1. Get token from frontend Pod
2. Use token to attack cluster API server
3. Get secrets etc. to further attack



# Authentication & Authorization



## Authentication

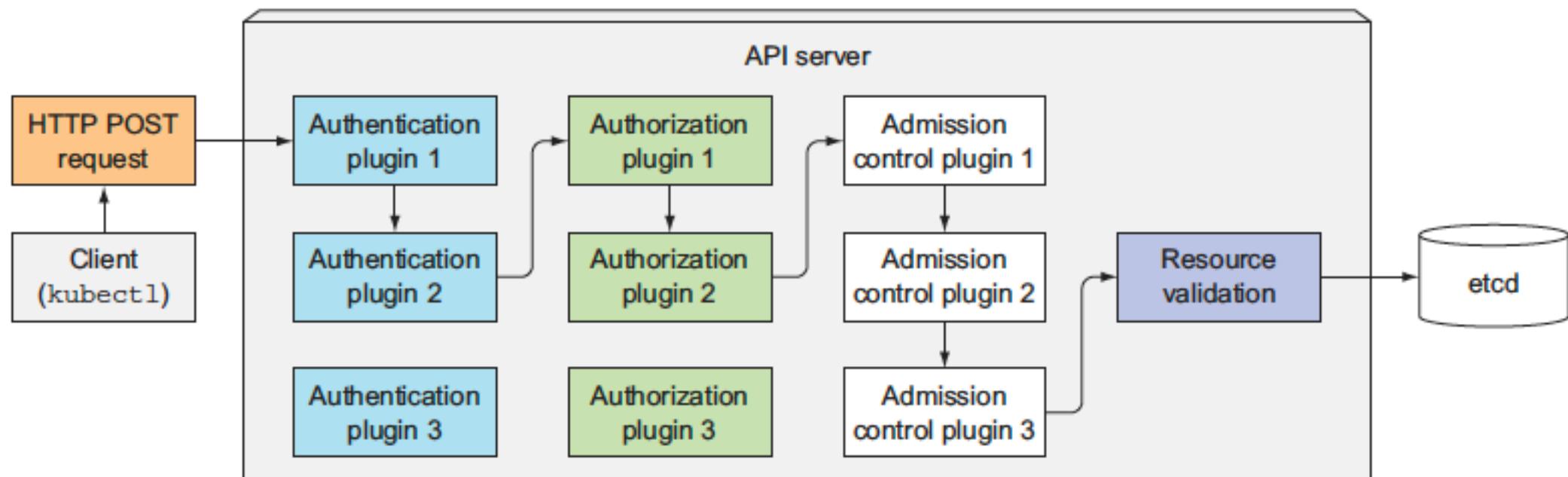
Who you are



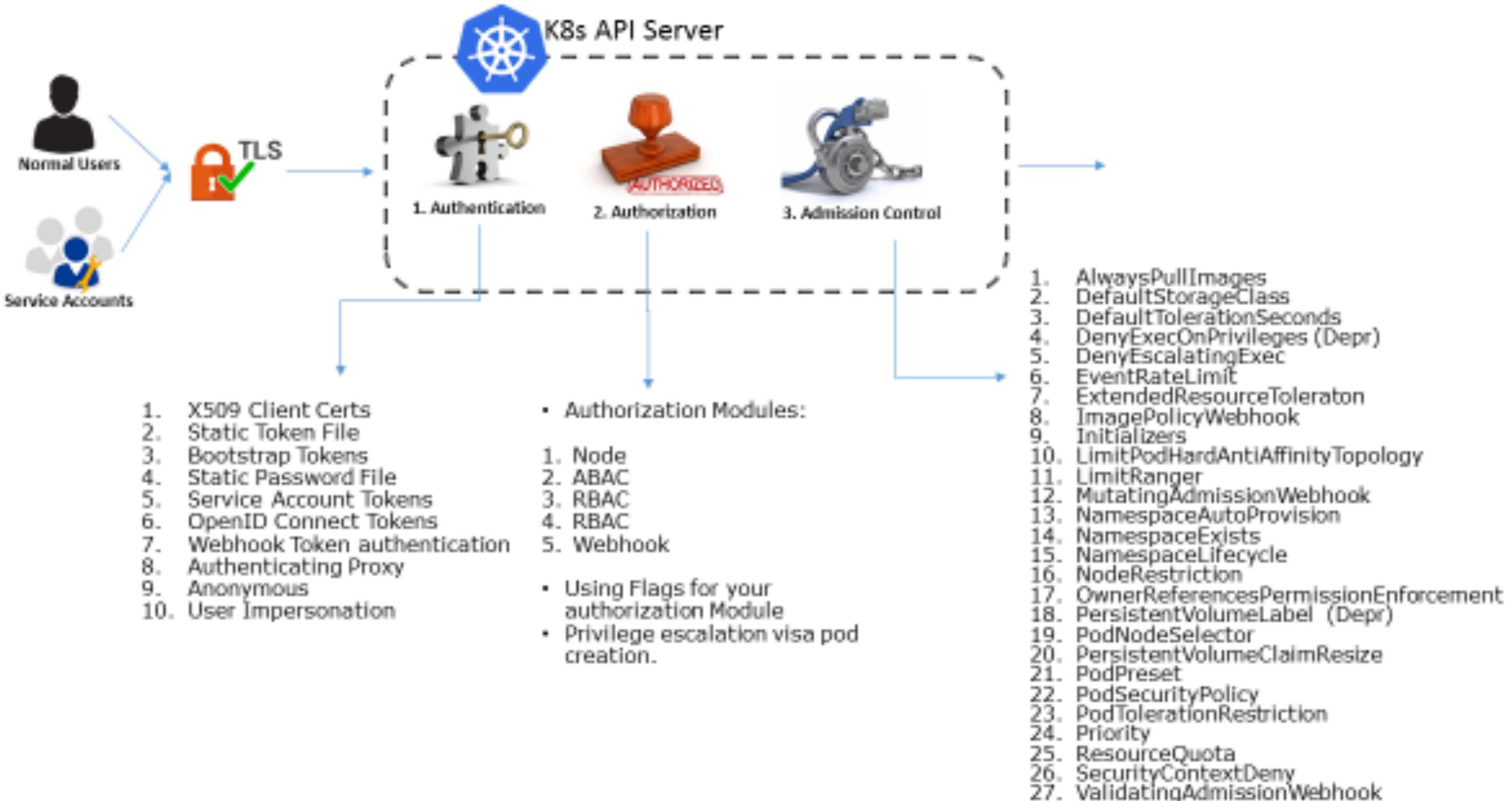
## Authorization

What you can do

# Authentication & Authorization in API Server



# Accessing the Kubernetes API



# Kubernetes Authentication

Mechanism	Secret Source	Usage
X509 Client Certs	CSR generated externally and signed with the cluster CA key	Enterprise CA / PKI
	Via Kubernetes API CertificateSigningRequest	Kubernetes cluster admin
Bearer token	Bootstrap token	Internal use
	Node authentication token	Internal use
	Static token file	Insecure
	ServiceAccount token	Pods, containers, applications, <i>users</i>
	OIDC token	Users
HTTP Basic auth	Static password file	Insecure
Auth proxy	N/A (trust proxy)	Integration
Impersonate	N/A (trust account)	Integration and administration

# Certificate Based Authentication

- Kubernetes is configured with a Certificate Authority (CA)



/etc/kubernetes/pki/ca.crt

Public certificate

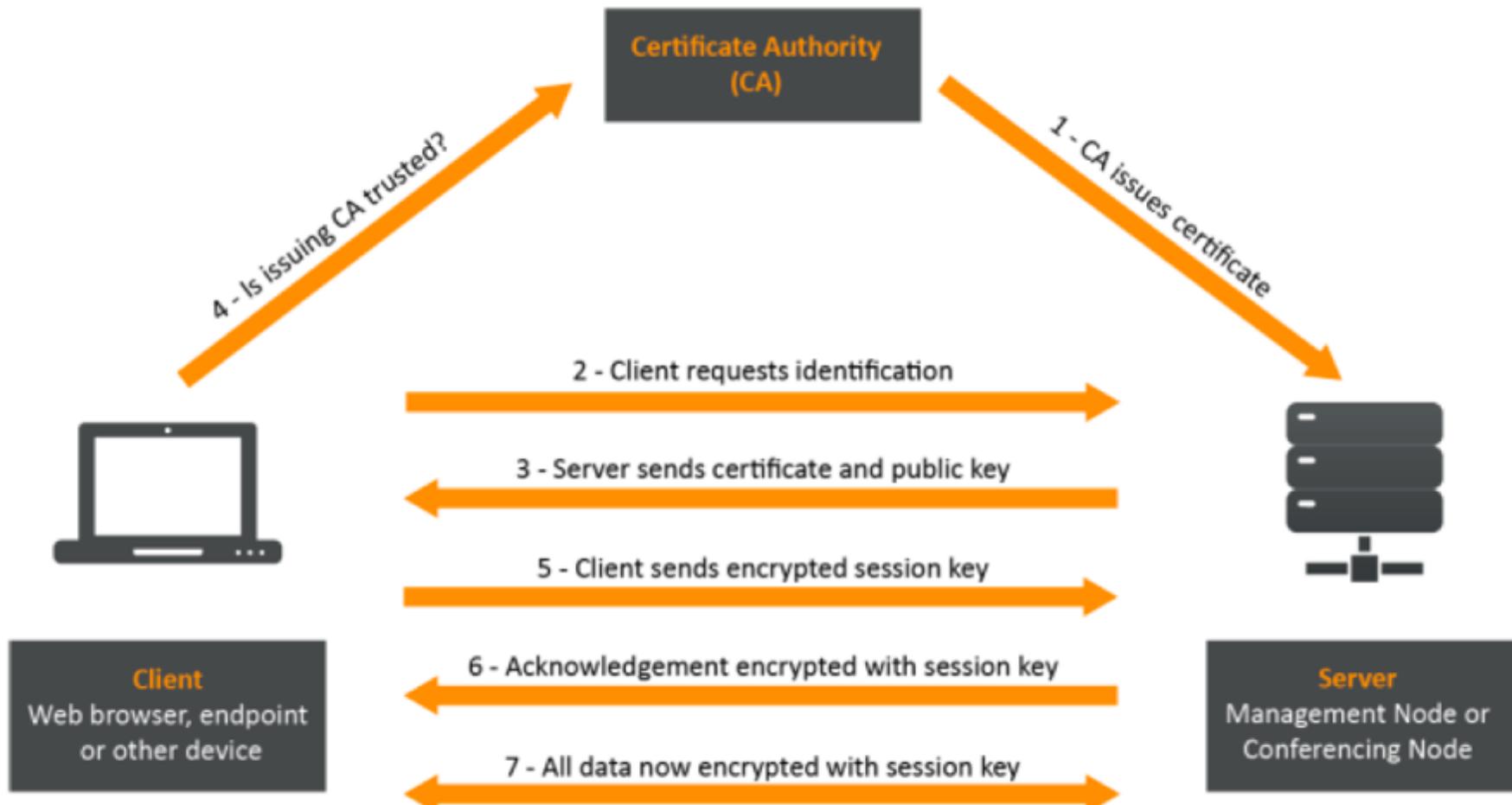
/etc/kubernetes/pki/ca.key

Private key



- Every SSL certificate signed with this CA will be accepted by the Kubernetes API
- Possible options for creating certificates: OpenSSL or CloudFlare's PKI toolkit
- Two important fields in the SSL certificate:
  - Common Name (CN): Kubernetes will interpret this value as the **user**
  - Organization (O): Kubernetes will interpret this value as the **group**

# TLS Certificate



# Certificate & CA in Kubernetes

CERTIFICATE	RESIDUAL TIME	Certificate Path	CERTIFICATE AUTHORITY
admin.conf	365d	/etc/kubernetes	
apiserver	365d	/etc/kubernetes/pki	ca
apiserver-etcd-client	365d	/etc/kubernetes/pki/etcd	etcd-ca
apiserver-kubelet-client	365d	/etc/kubernetes/pki	ca
controller-manager.conf	365d	/etc/kubernetes/pki	
etcd-healthcheck-client	365d	/etc/kubernetes/pki/etcd	etcd-ca
etcd-peer	365d	/etc/kubernetes/pki/etcd	etcd-ca
etcd-server	365d	/etc/kubernetes/pki/etcd	etcd-ca
front-proxy-client	365d	/etc/kubernetes/pki	front-proxy-ca
scheduler.conf	365d	/etc/kubernetes	

CERTIFICATE AUTHORITY	RESIDUAL TIME	Certificate Path
ca	10y	/etc/kubernetes/pki
etcd-ca	10y	/etc/kubernetes/pki/etcd
front-proxy-ca	10y	/etc/kubernetes/pki

참고: front-proxy 인증서는 kube-proxy에서 API 서버 확장을 지원할 때만 kube-proxy에서 필요하다.

# Creating User Certificate (1)



Developer

- Create private key (if it does not exist)

```
openssl genrsa -out juan.key 2048
```

- Create certificate signing request (CSR)

```
openssl req -new -key juan.key -out juan.csr -subj "/CN=juan/O=devs"
```

user

group

- Send the CSR to the administrator



Administrator

- Create certificate from CSR using the cluster authority



```
openssl x509 -req -in juan.csr -CA CA_LOCATION/ca.crt -CAkey  
CA_LOCATION/ca.key -CAcreateserial -out juan.crt -days 500
```

# Creating User Certificate (2)

- To add in your local machine the new configuration:
  - Download the cluster authority and generated certificate
  - Add the new cluster to kubectl

```
kubectl config set-cluster sandbox --certificate-authority=ca.pem  
--embed-certs=true --server=https://<PUBLIC_ADDRESS_OF_YOUR_CLUSTER>:6443
```

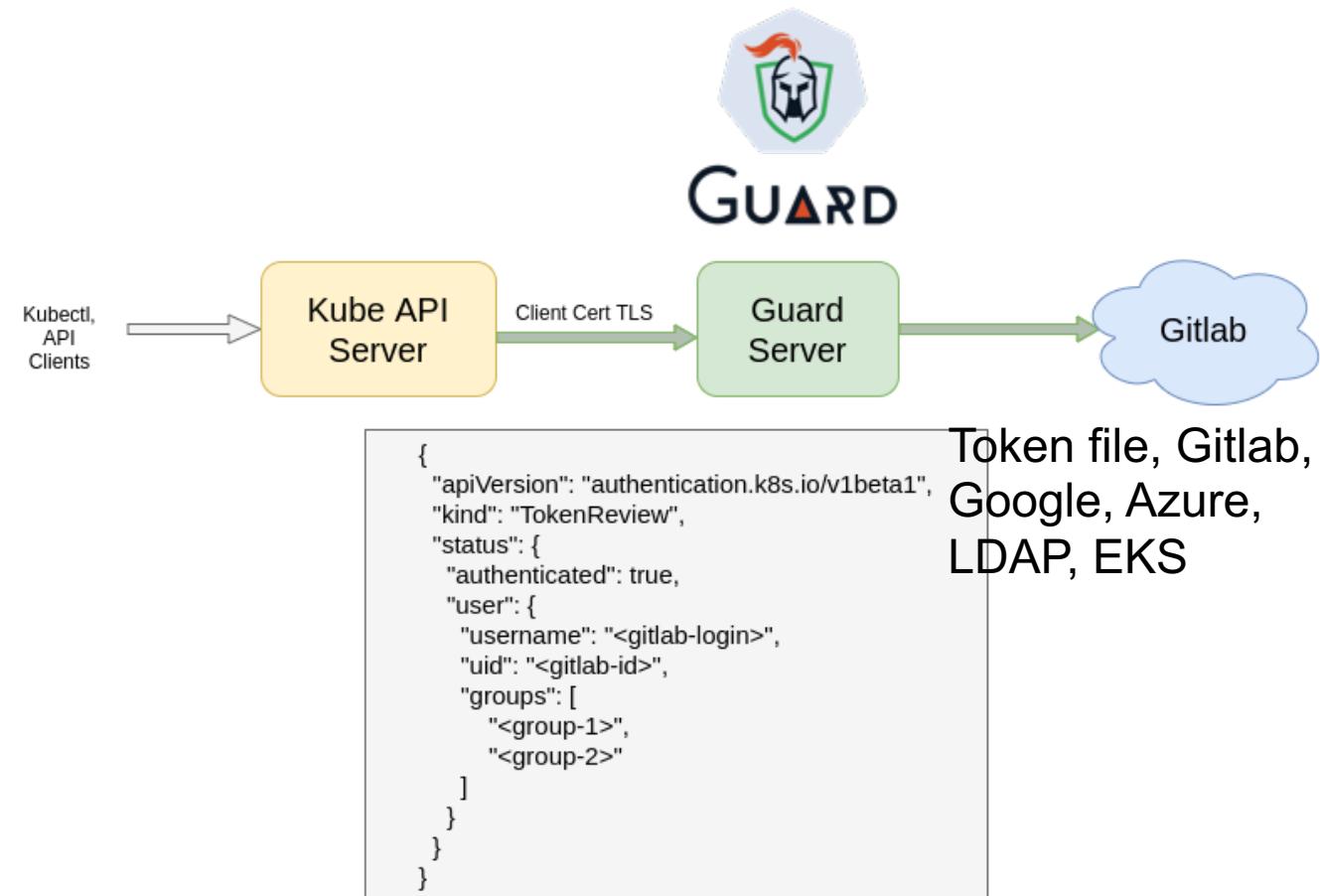
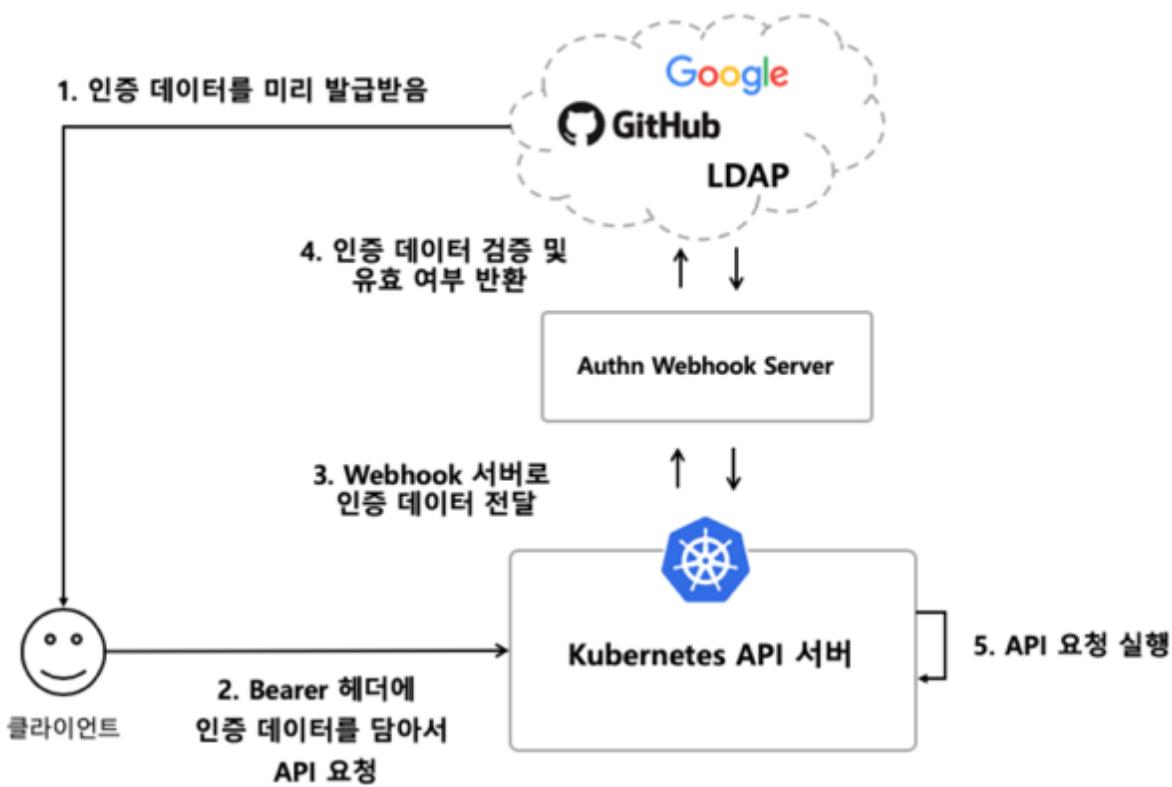
- Add the new credentials to kubectl

```
kubectl config set-credentials juan --client-certificate=juan.crt  
--client-key=juan.key --embed-certs=true
```

- Add the new context to kubectl

```
kubectl config set-context sandbox-juan --cluster=sandbox --user=juan
```

# Token Webhook



# Kubernetes Authorization

Mechanism	Decision source	Usage
Node	API Server built-in	Internal use (kubelets)
ABAC	Static file	Insecure, deprecated
RBAC	API Objects	Users and administrators
WebHook	External services	Integration
AlwaysDeny	API Server built-in	Testing
AlwaysAllow		

# RBAC



**Subjects**



**API Resources**

list  
create  
delete

get  
watch  
patch

**Operations  
(Verbs)**

# Role

```
kind: Role
apiVersion: rbac.authorization.k8s.io/v1beta1
metadata:
  namespace: test
  name: pod-access
rules:
  - apiGroups: [""]
    resources: ["pods"]
    verbs: ["get", "list"]
```

WHICH RESOURCES      WHICH OPERATIONS

- Need to specify:
  - Api group
  - Name

Find it in the [API reference](#), examples

Group	Version	Kind
apps	v1	Deployment

Group	Version	Kind
core	v1	Pod

When it is core, we use an empty string

# API Group

API Group	API Objects
rbac.authorization.k8s.io/v1	ClusterRole Role ClusterRoleBinding RoleBinding
authentication.k8s.io/v1	TokenReview
admissionregistration.k8s.io/v1	MutatingWebhookConfiguration ValidatingWebhookConfiguration
authorization.k8s.io/v1	LocalSubjectAccessReview SelfSubjectAccessReview SelfSubjectRulesReview SubjectAccessReview
certificates.k8s.io/v1beta1	CertificateSigningRequest
policy/v1beta1	PodSecurityPolicy

# Role Binding

```
kind: RoleBinding
apiVersion: rbac.authorization.k8s.io/v1
metadata:
  name: devs-read-pods
  namespace: test
subjects:
- kind: Group
  name: devs
  apiGroup: rbac.authorization.k8s.io
roleRef:
  kind: Role
  name: pod-access
  apiGroup: rbac.authorization.k8s.io
```

- Examples
  - User
  - Group
  - ...

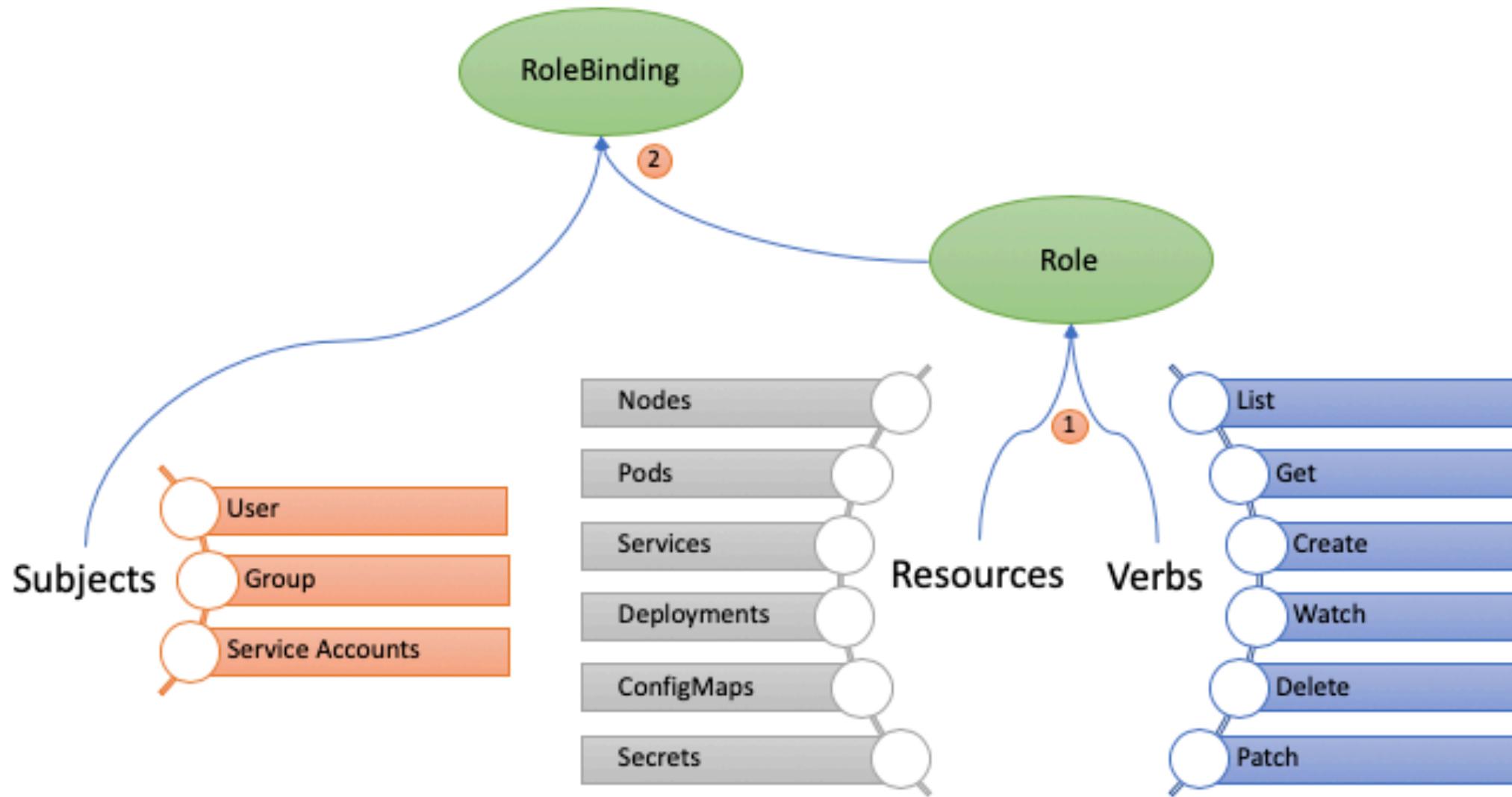
Later we will see another one

WHICH SUBJECTS

Used to specify which api group the kind belongs to

WHICH ROLE (ONLY ONE PER BINDING)

# Role & Role Binding



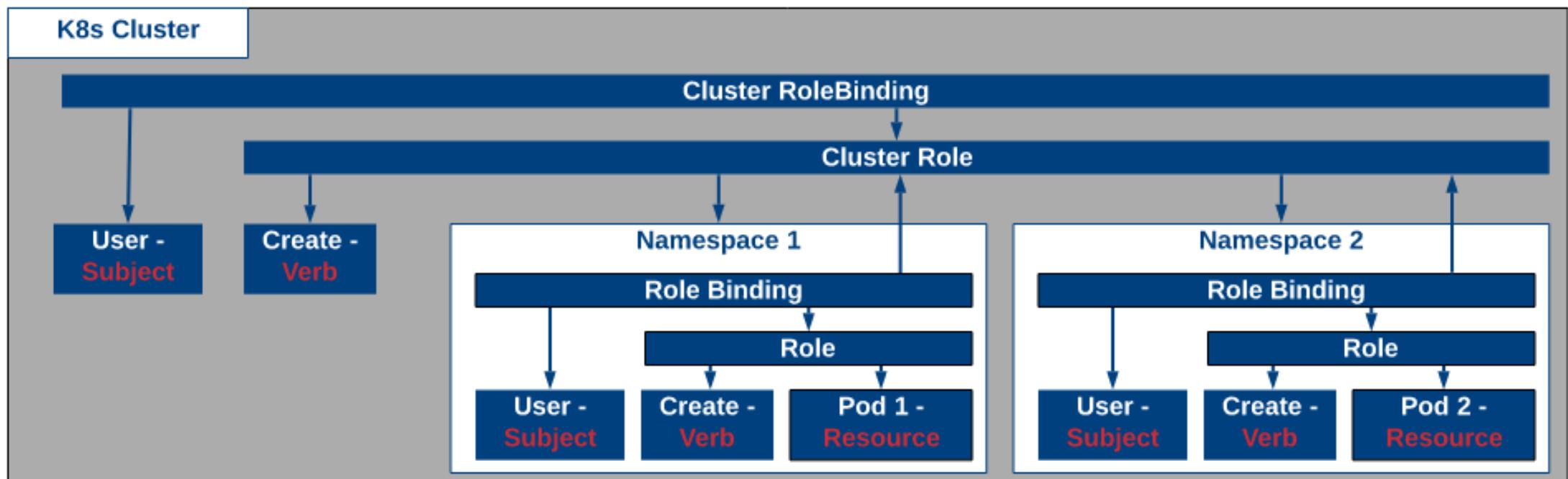
# Cluster Role

```
kind: Role
apiVersion: rbac.authorization.k8s.io/v1beta1
metadata:
  name: pod-access
  namespace: test
rules:
  - apiGroups: [""]
    resources: ["pods"]
    verbs: ["get", "list"]
```

Only difference

```
kind: ClusterRole
apiVersion: rbac.authorization.k8s.io/v1beta1
metadata:
  name: all-pod-access
rules:
  - apiGroups: [""]
    resources: ["pods"]
    verbs: ["get", "list"]
```

# Role & Cluster Role



# Cluster Role Binding

```
kind: RoleBinding
apiVersion: rbac.authorization.k8s.io/v1
metadata:
  name: devs-read-pods
  namespace: test
subjects:
- kind: User
  name: jsalmeron # Name is case sensitive
  apiGroup: rbac.authorization.k8s.io
roleRef:
  kind: Role
  name: ns-admin
  apiGroup: rbac.authorization.k8s.io
```

Only differences

```
kind: ClusterRoleBinding
apiVersion: rbac.authorization.k8s.io/v1
metadata:
  name: salme-reads-all-pods
subjects:
- kind: User
  name: jsalmeron # Name is case sensitive
  apiGroup: rbac.authorization.k8s.io
roleRef:
  kind: ClusterRole
  name: all-pod-access
  apiGroup: rbac.authorization.k8s.io
```

# Default Cluster Role Bindings

- Kubernetes includes some ClusterRoleBindings. For example:
  - **system:basic-user**: For unauthenticated users (group **system:unauthenticated**). No operations are allowed.
  - **cluster-admin**: For members of the **system:masters** group. Can do any operation on the cluster (using **cluster-admin** ClusterRole).



Admin accounts can be created belonging to this group

```
openssl req ... -subj "/CN=dbarranco/O=system:masters"
```

- ClusterRoleBindings for the **different components** of the cluster (kube-controller-manager, kube-scheduler, kube-proxy ...)

# Kubernetes RBAC API Discovery Roles

Default ClusterRole	Default ClusterRoleBinding	Description
<b>system:basic-user</b>	<b>system:authenticated</b> group	Allows a user read-only access to basic information about themselves. Prior to v1.14, this role was also bound to <b>system:unauthenticated</b> by default.
<b>system:discovery</b>	<b>system:authenticated</b> group	Allows read-only access to API discovery endpoints needed to discover and negotiate an API level. Prior to v1.14, this role was also bound to <b>system:unauthenticated</b> by default.
<b>system:public-info-viewer</b>	<b>system:authenticated</b> and <b>system:unauthenticated</b> groups	Allows read-only access to non-sensitive information about the cluster. Introduced in Kubernetes v1.14.

Kubernetes RBAC API discovery roles

Default ClusterRole	Default ClusterRoleBinding	Description
<b>cluster-admin</b>	<b>system:masters</b> group	Allows super-user access to perform any action on any resource. When used in a <b>ClusterRoleBinding</b> , it gives full control over every resource in the cluster and in all namespaces. When used in a <b>RoleBinding</b> , it gives full control over every resource in the role binding's namespace, including the namespace itself.
<b>admin</b>	None	Allows admin access, intended to be granted within a namespace using a <b>RoleBinding</b> . If used in a <b>RoleBinding</b> , allows read/write access to most resources in a namespace, including the ability to create roles and role bindings within the namespace. This role does not allow write access to resource quota or to the namespace itself.
<b>edit</b>	None	Allows read/write access to most objects in a namespace. This role does not allow viewing or modifying roles or role bindings. However, this role allows accessing Secrets and running Pods as any ServiceAccount in the namespace, so it can be used to gain the API access levels of any ServiceAccount in the namespace.
<b>view</b>	None	Allows read-only access to see most objects in a namespace. It does not allow viewing roles or role bindings. This role does not allow viewing Secrets, since reading the contents of Secrets enables access to ServiceAccount credentials in the namespace, which would allow API access as any ServiceAccount in the namespace (a form of privilege escalation).

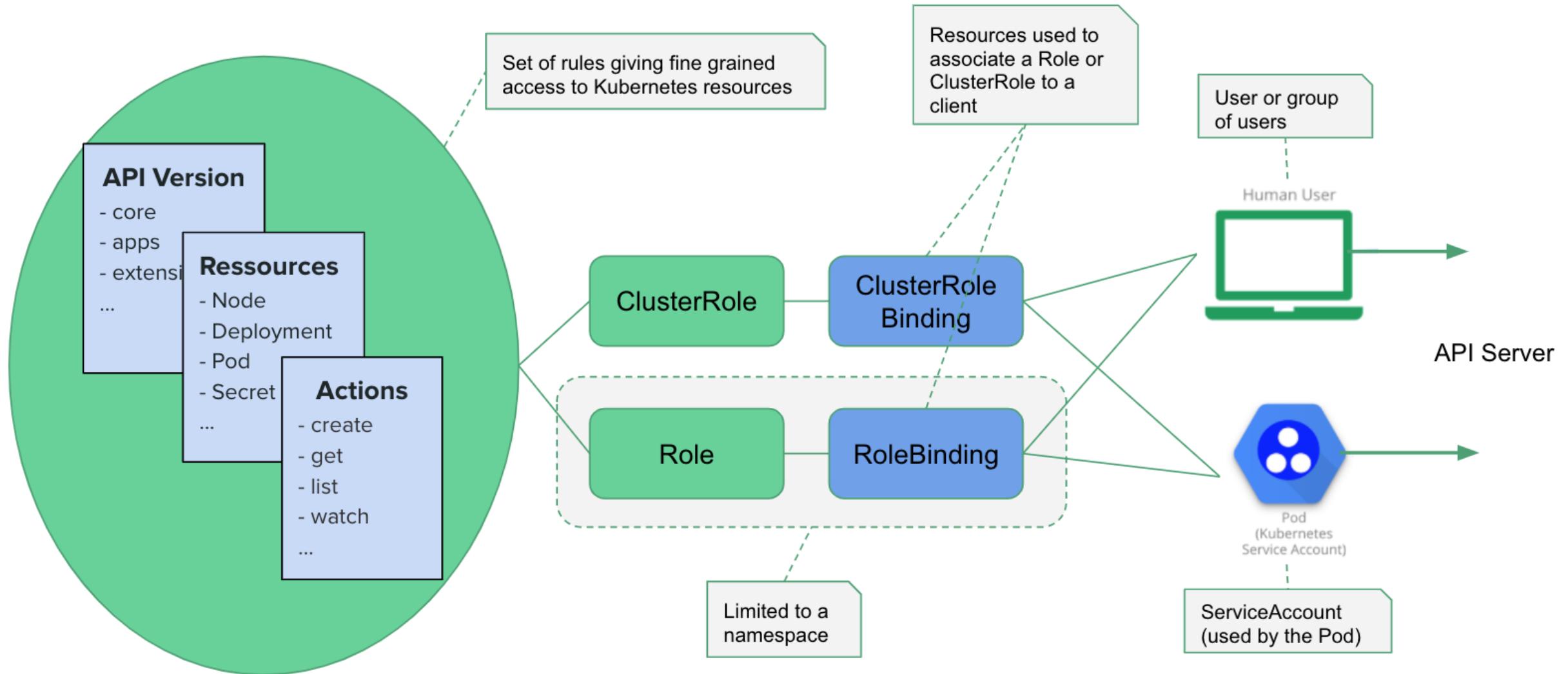
# Core component roles

Default ClusterRole	Default ClusterRoleBinding	Description
<code>system:kube-scheduler</code>	<code>system:kube-scheduler</code> user	Allows access to the resources required by the <a href="#">scheduler</a> component.
<code>system:volume-scheduler</code>	<code>system:kube-scheduler</code> user	Allows access to the volume resources required by the kube-scheduler component.
<code>system:kube-controller-manager</code>	<code>system:kube-controller-manager</code> user	Allows access to the resources required by the <a href="#">controller manager</a> component. The permissions required by individual controllers are detailed in the <a href="#">controller roles</a> .
<code>system:node</code>	None	<p>Allows access to resources required by the kubelet, <b>including read access to all secrets, and write access to all pod status objects</b>. You should use the <a href="#">Node authorizer</a> and <a href="#">NodeRestriction admission plugin</a> instead of the <code>system:node</code> role, and allow granting API access to kubelets based on the Pods scheduled to run on them.</p> <p>The <code>system:node</code> role only exists for compatibility with Kubernetes clusters upgraded from versions prior to v1.8.</p>
<code>system:node-proxier</code>	<code>system:kube-proxy</code> user	Allows access to the resources required by the <a href="#">kube-proxy</a> component.

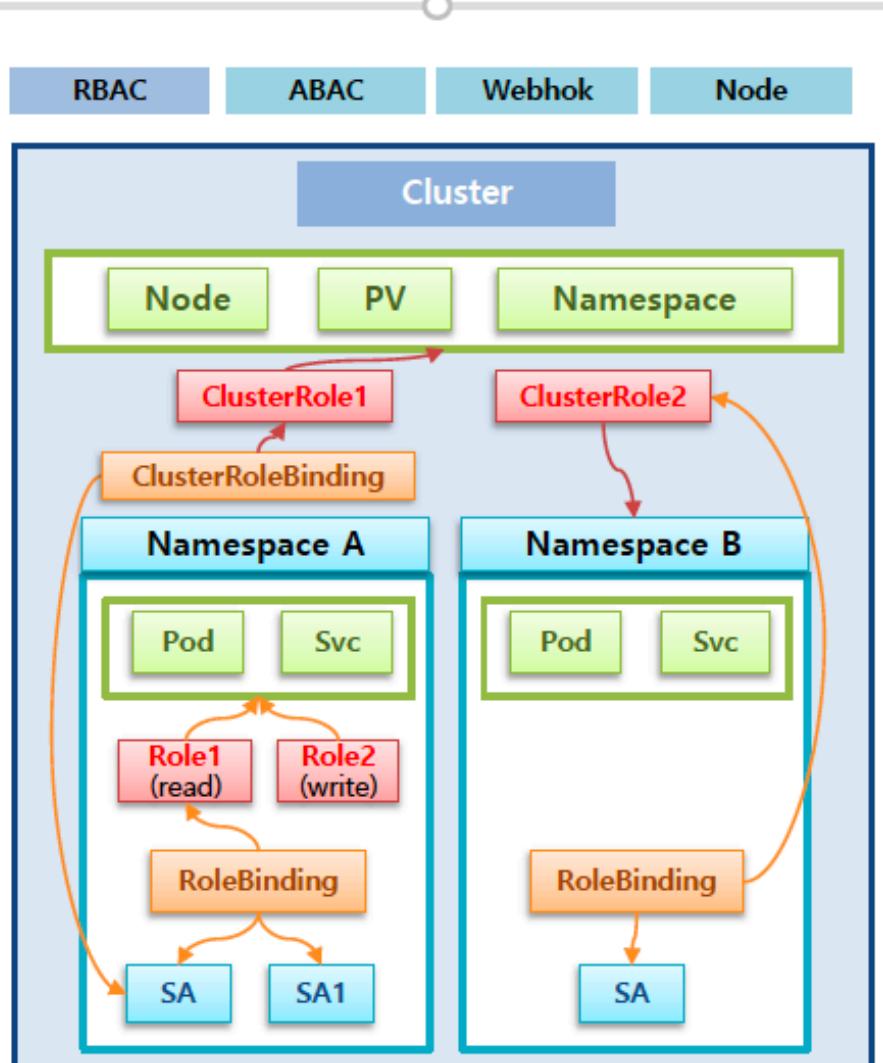
## Other component roles

Default ClusterRole	Default ClusterRoleBinding	Description
<b>system:auth-delegator</b>	None	Allows delegated authentication and authorization checks. This is commonly used by add-on API servers for unified authentication and authorization.
<b>system:heapster</b>	None	Role for the <a href="#">Heapster</a> component (deprecated).
<b>system:kube-aggregator</b>	None	Role for the <a href="#">kube-aggregator</a> component.
<b>system:kube-dns</b>	<b>kube-dns</b> service account in the <b>kube-system</b> namespace	Role for the <a href="#">kube-dns</a> component.
<b>system:kubelet-api-admin</b>	None	Allows full access to the kubelet API.
<b>system:node-bootstrapper</b>	None	Allows access to the resources required to perform <a href="#">kubelet TLS bootstrapping</a> .
<b>system:node-problem-detector</b>	None	Role for the <a href="#">node-problem-detector</a> component.
<b>system:persistent-volume-provisioner</b>	None	Allows access to the resources required by most <a href="#">dynamic volume provisioners</a> .

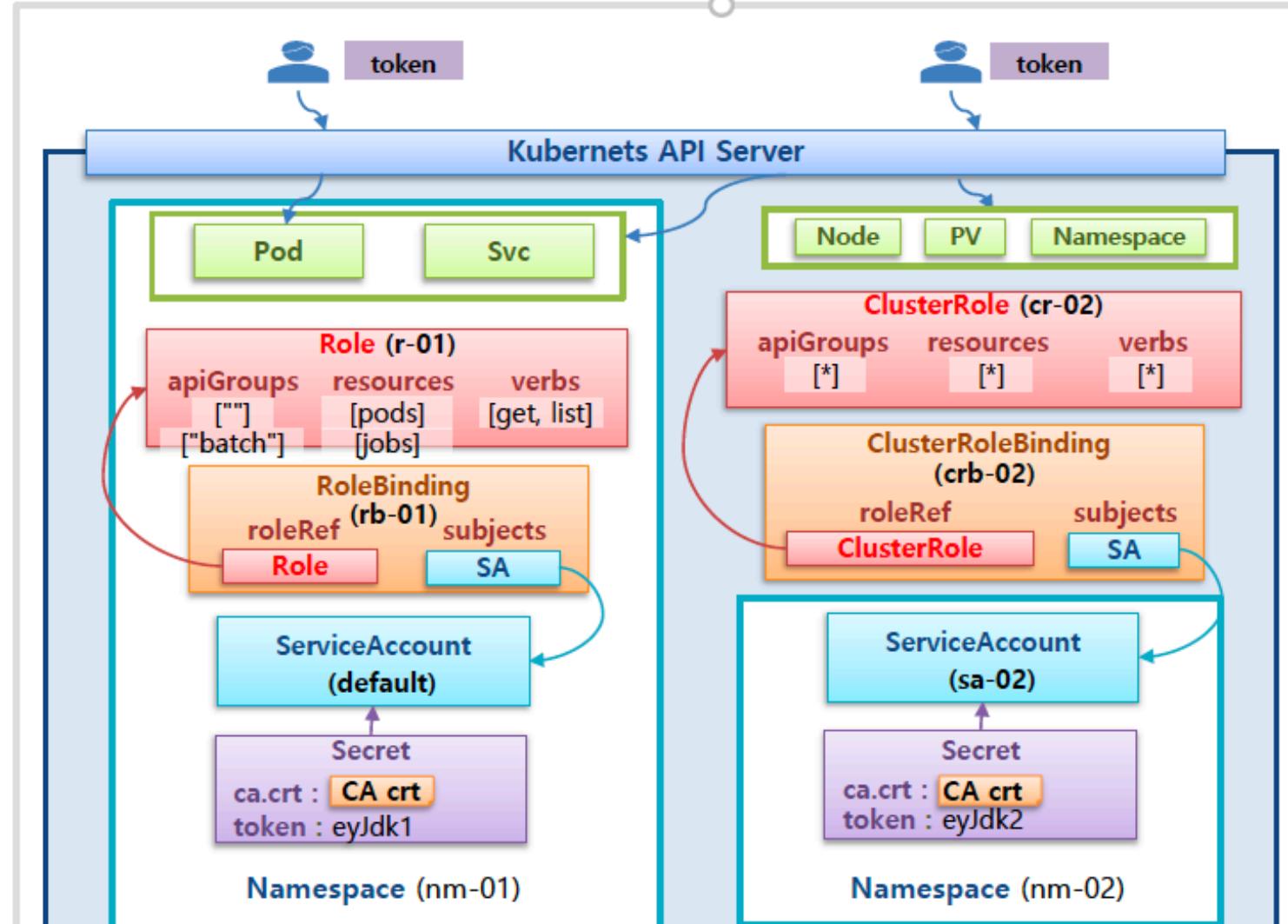
# Kubernetes RBAC



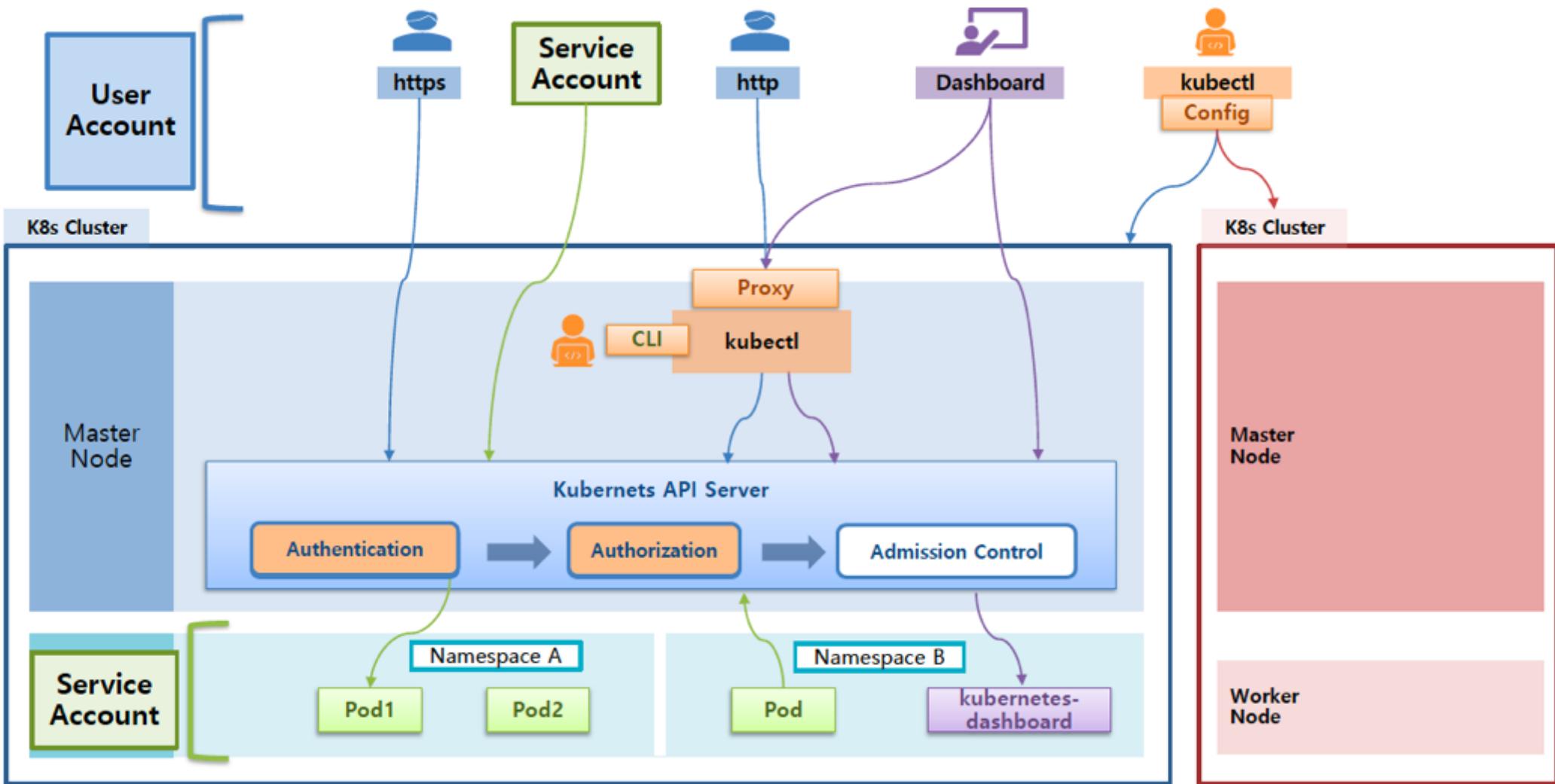
## RBAC (Role, RoleBinding) Overview



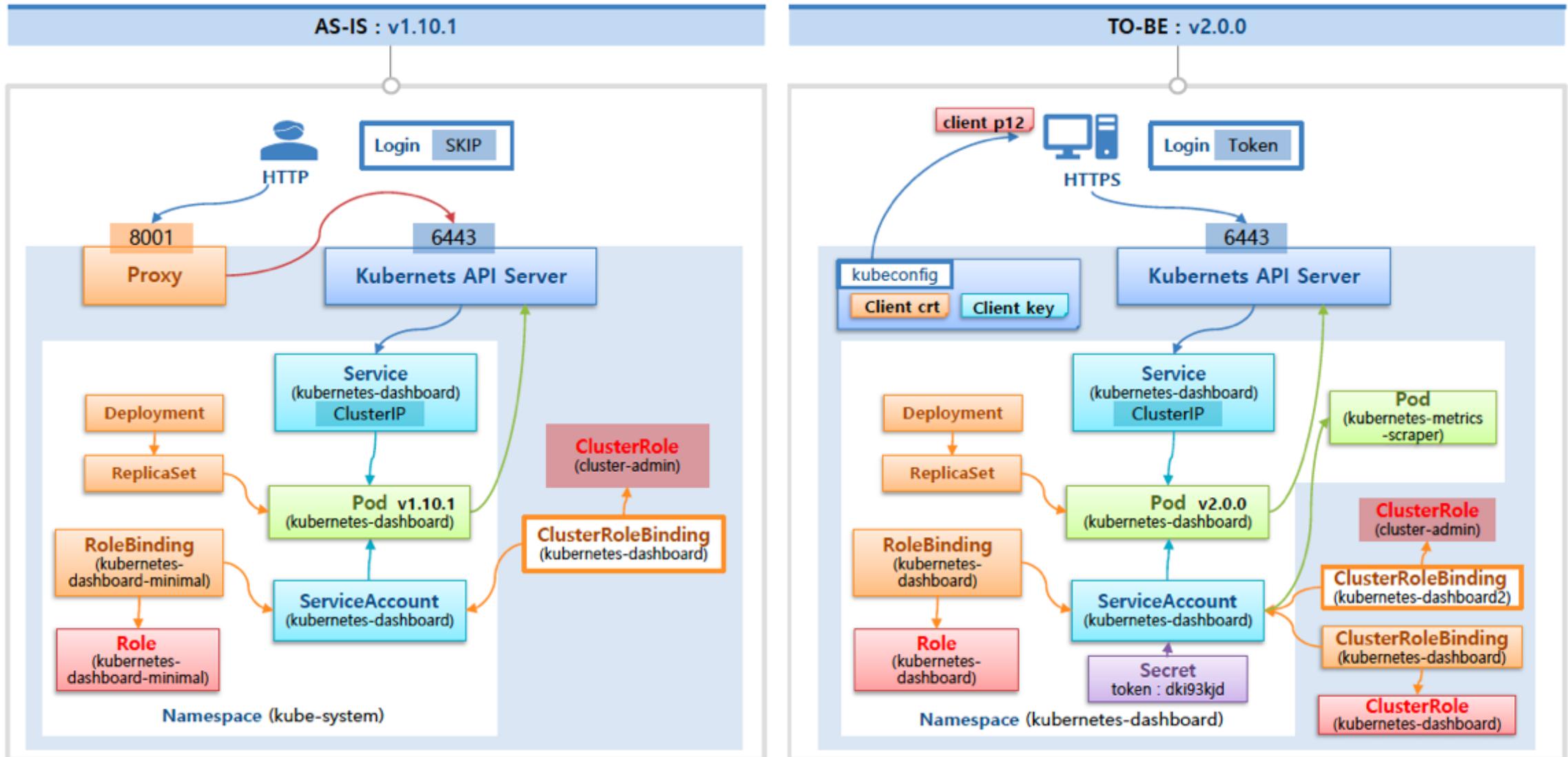
## Role, RoleBinding Detail



# Service Account



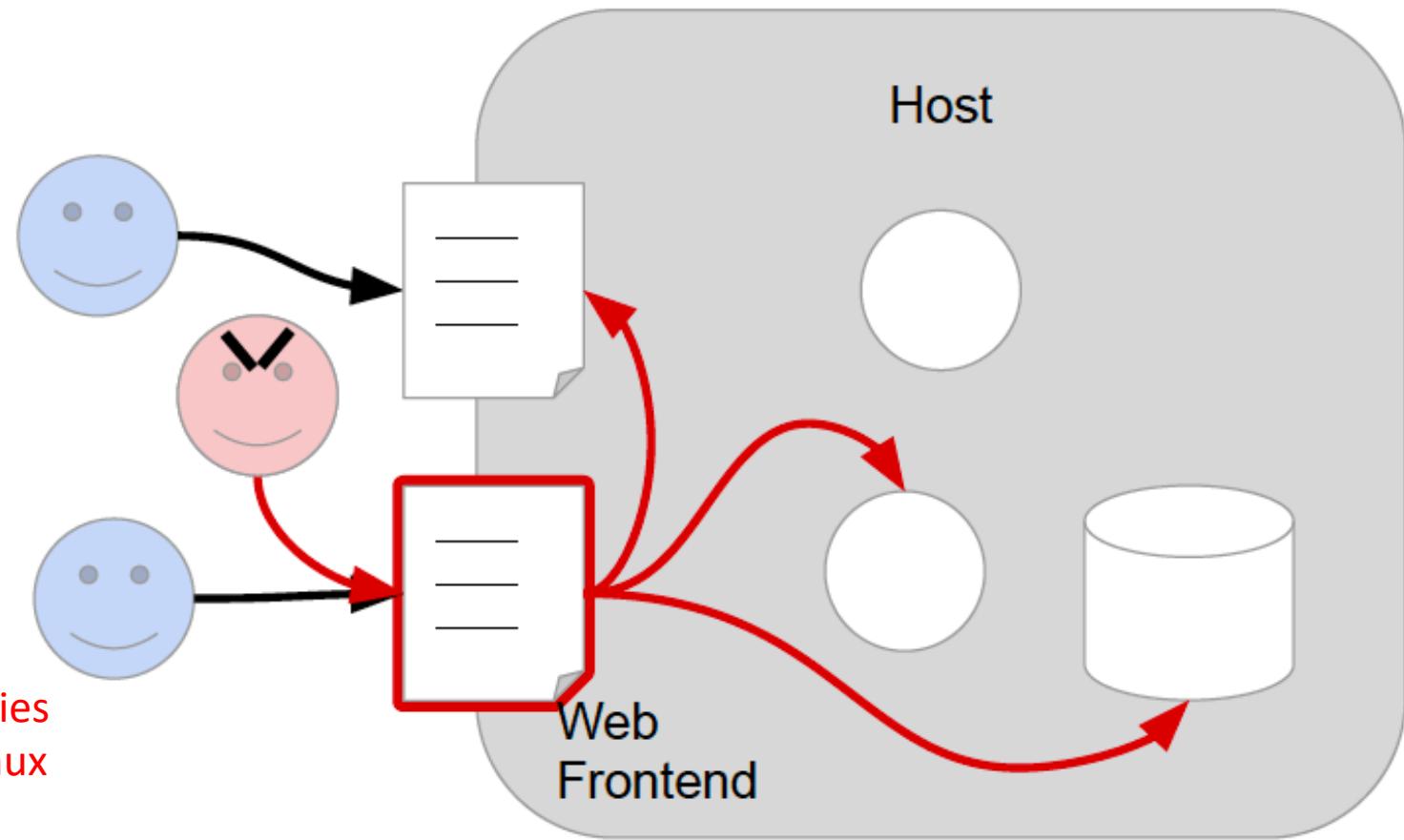
# Kubernetes Dashboard



# Attack Surface: Access to Host

1. Break out of the container using container or kernel exploits
2. Attack the Kubelet
3. Attack other containers running on the same host

Mitigate: Sandboxed Pods, Capabilities Seccomp, AppArmor, SELinux



# Security Context

```
apiVersion: v1
kind: Pod
metadata:
  name: security-context-demo
spec:
  securityContext:
    runAsUser: 1000
    readOnlyRootFilesystem: true
    allowPrivilegeEscalation: false
```

# Seccomp, AppArmor, SELinux

metadata:

annotations:

seccomp.security.alpha.kubernetes.io/pod: runtime/default

container.apparmor.security.beta.kubernetes.io/hello: runtime/default

securityContext:

seLinuxOptions:

level: "s0:c123,c456"

# Security Context (1)

Field	Description
<code>allowPrivilegeEscalation</code> <i>boolean</i>	AllowPrivilegeEscalation controls whether a process can gain more privileges than its parent process. This bool directly controls if the no_new_privs flag will be set on the container process. AllowPrivilegeEscalation is true always when the container is: 1) run as Privileged 2) has CAP_SYS_ADMIN
<code>capabilities</code> <i>Capabilities</i>	The capabilities to add/drop when running containers. Defaults to the default set of capabilities granted by the container runtime.
<code>privileged</code> <i>boolean</i>	Run container in privileged mode. Processes in privileged containers are essentially equivalent to root on the host. Defaults to false.
<code>procMount</code> <i>string</i>	procMount denotes the type of proc mount to use for the containers. The default is DefaultProcMount which uses the container runtime defaults for readonly paths and masked paths. This requires the ProcMountType feature flag to be enabled.
<code>readOnlyRootFilesystem</code> <i>boolean</i>	Whether this container has a read-only root filesystem. Default is false.

# Security Context (2)

<b>runAsGroup</b> <i>integer</i>	The GID to run the entrypoint of the container process. Uses runtime default if unset. May also be set in PodSecurityContext. If set in both SecurityContext and PodSecurityContext, the value specified in SecurityContext takes precedence.
<b>runAsNonRoot</b> <i>boolean</i>	Indicates that the container must run as a non-root user. If true, the Kubelet will validate the image at runtime to ensure that it does not run as UID 0 (root) and fail to start the container if it does. If unset or false, no such validation will be performed. May also be set in PodSecurityContext. If set in both SecurityContext and PodSecurityContext, the value specified in SecurityContext takes precedence.
<b>runAsUser</b> <i>integer</i>	The UID to run the entrypoint of the container process. Defaults to user specified in image metadata if unspecified. May also be set in PodSecurityContext. If set in both SecurityContext and PodSecurityContext, the value specified in SecurityContext takes precedence.
<b>seLinuxOptions</b> <i>SELinuxOptions</i>	The SELinux context to be applied to the container. If unspecified, the container runtime will allocate a random SELinux context for each container. May also be set in PodSecurityContext. If set in both SecurityContext and PodSecurityContext, the value specified in SecurityContext takes precedence.
<b>seccompProfile</b> <i>SeccompProfile</i>	The seccomp options to use by this container. If seccomp options are provided at both the pod & container level, the container options override the pod options.
<b>windowsOptions</b> <i>WindowsSecurityContextOptions</i>	The Windows specific settings applied to all containers. If unspecified, the options from the PodSecurityContext will be used. If set in both SecurityContext and PodSecurityContext, the value specified in SecurityContext takes precedence.

# Capabilities

- Traditional UNIX privilege scheme
  - Process running with EUID 0 (superuser): all privilege actions allowed
  - Process running with EUID  $\neq$  0 : no privilege actions allowed
- Capabilities
  - Since kernel 2.2 (late 90's), 38 capabilities total
  - Collections of distinct privileges that can be enabled per process (thread)
  - CAP\_SYS\_BOOT, CAP\_SYS\_ADMIN, CAP\_CHOWN, CAP\_SYS\_TIME, CAP\_KILL ...
  - capabilities(7), libcap(3)
  - 5 capability sets (effective, permitted, inheritable, ambient, bounding)
  - capset(2), capget(2), prctl
  - systemd: Default, User=..., AmbientCapabilities=X, CapabilityBoundingSet=X

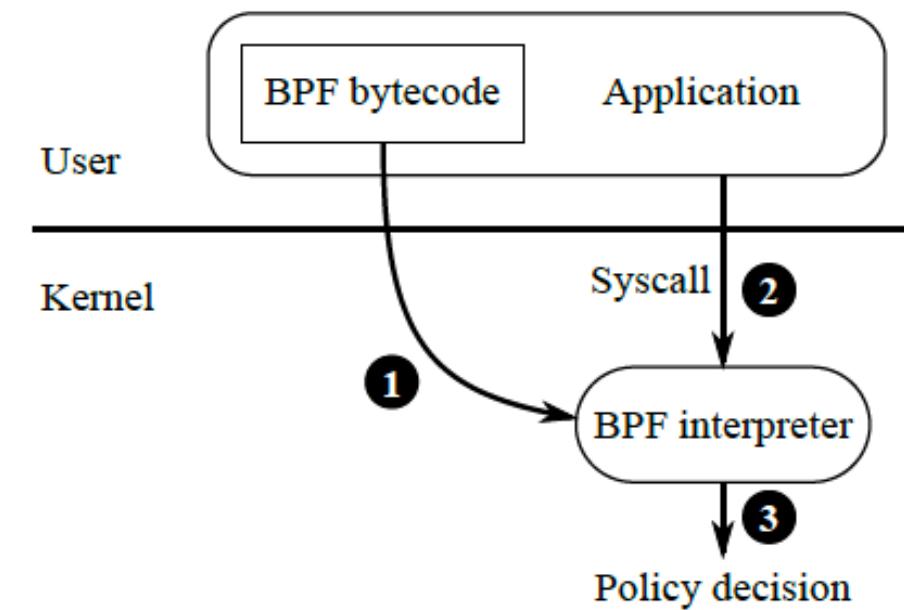
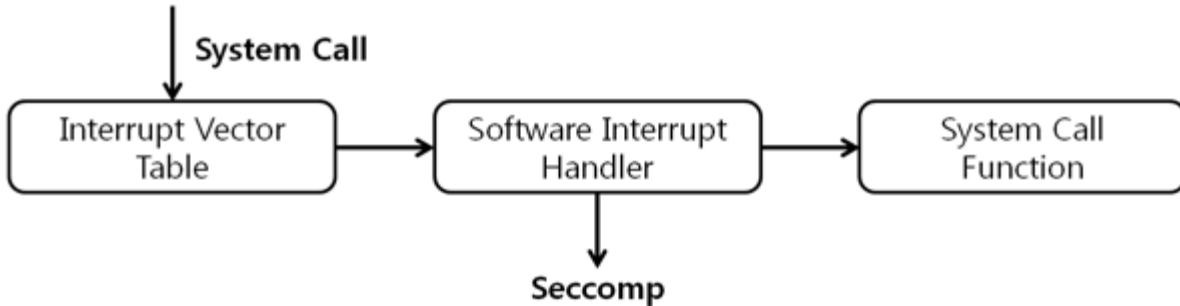
```
$ cat /proc/$$/task/$$/status
...
CapEff: 0000000000000000
CapPrm: 0000000000000000
CapInh: 0000000000000000
CapAmb: 0000000000000000
CapBnd: 0000001fffffffff
                           securebits
```

# Dockerd Capabilities

- By default, Docker starts containers with a restricted set of capabilities
  - SETPCAP, MKNOD, AUDIT\_WRITE, CHOWN, NET\_RAW, DAC\_OVERRIDE, FOWNER, FSETID, KILL, SETGID, SETUID, NET\_BIND\_SERVICE, SYS\_CHROOT, SETFCAP
  - Removed: NET\_ADMIN, SYS\_ADMIN, ...
- Activated command
  - Runs as native process, by default under EUID 0
  - Runs with selected set of capabilities
- Capability-related arguments
  - Add capabilities: `--cap-add` *list*
  - Drop capabilities: `--cap-drop` *list*
  - User: `--user` *name|uid* (no capabilities)

# Seccomp

- Since kernel 2.6.12 (2005)
- Seccomp-bpf
  - Allows filtering of system calls using a configurable policy implemented using Berkeley Packet Filter rules
  - For Google Chrome

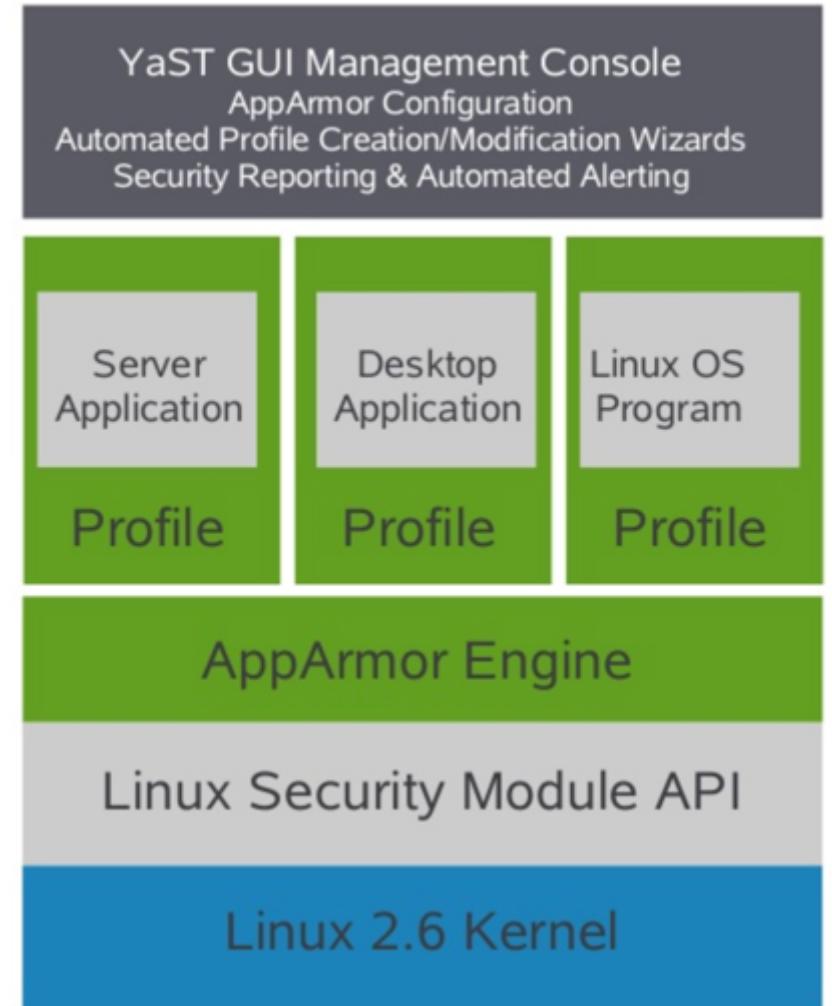


# Mandatory Access Control Systems: SELinux

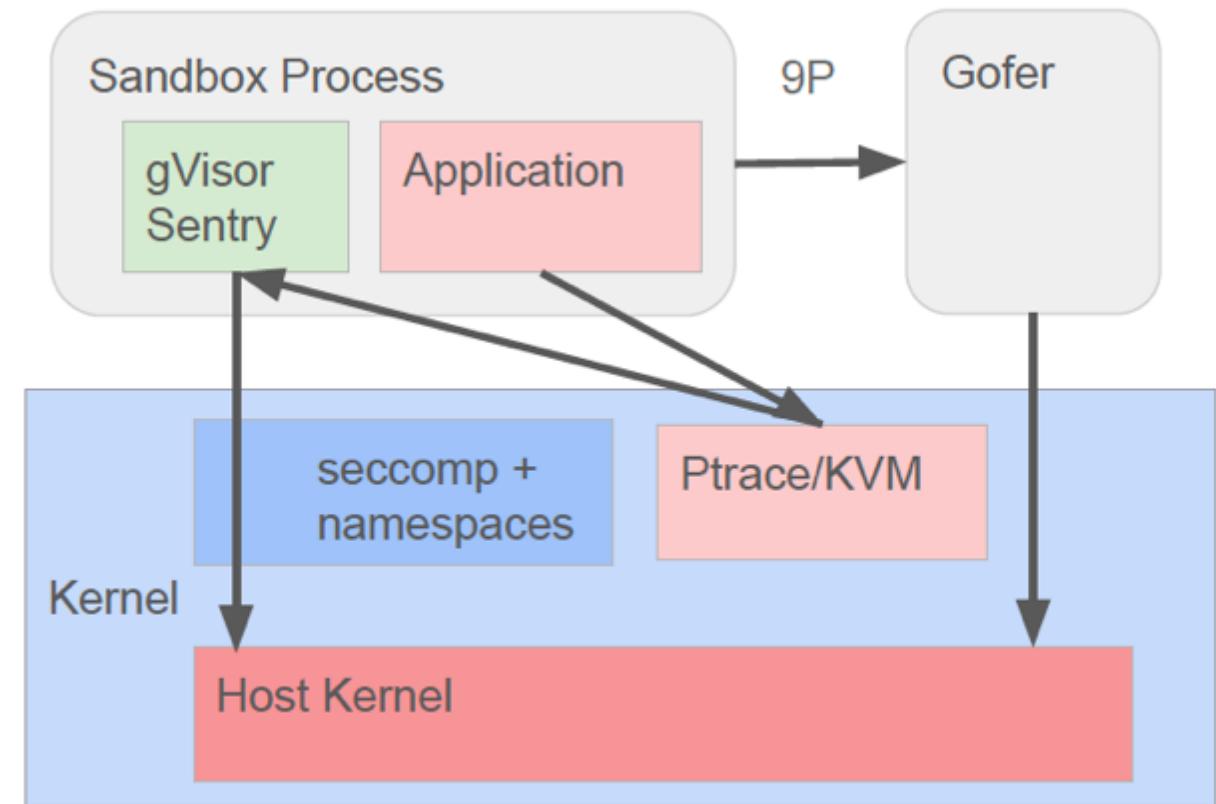
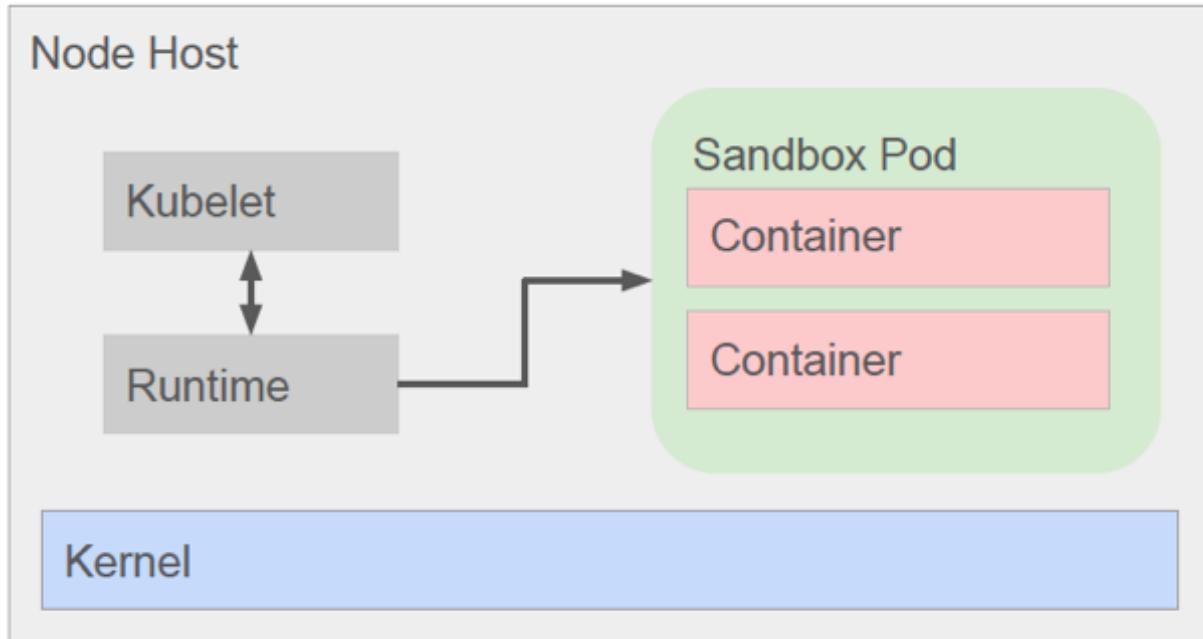
- SELinux (Security Enhanced Linux)
  - sudo setenforce 1 → sudo sestatus
  - Set security labels, 3 modes (enforce, permissive, disable)
  - Even if Apache HTTP Server is compromised, an attacker cannot use that process to read files in user home directories by default, unless a specific SELinux policy rule was added or configured to allow such access
- Protecting host file system from container breakout

# Mandatory Access Control Systems: AppArmor

- AppArmor
  - Restrict programs' capabilities with per-program profiles
  - Profiles: /etc/apparmor.d
  - sudo apparmor\_parser /etc/apparmor.d/...  
→ sudo aa-status
  - 2 modes (enforcement, complain)

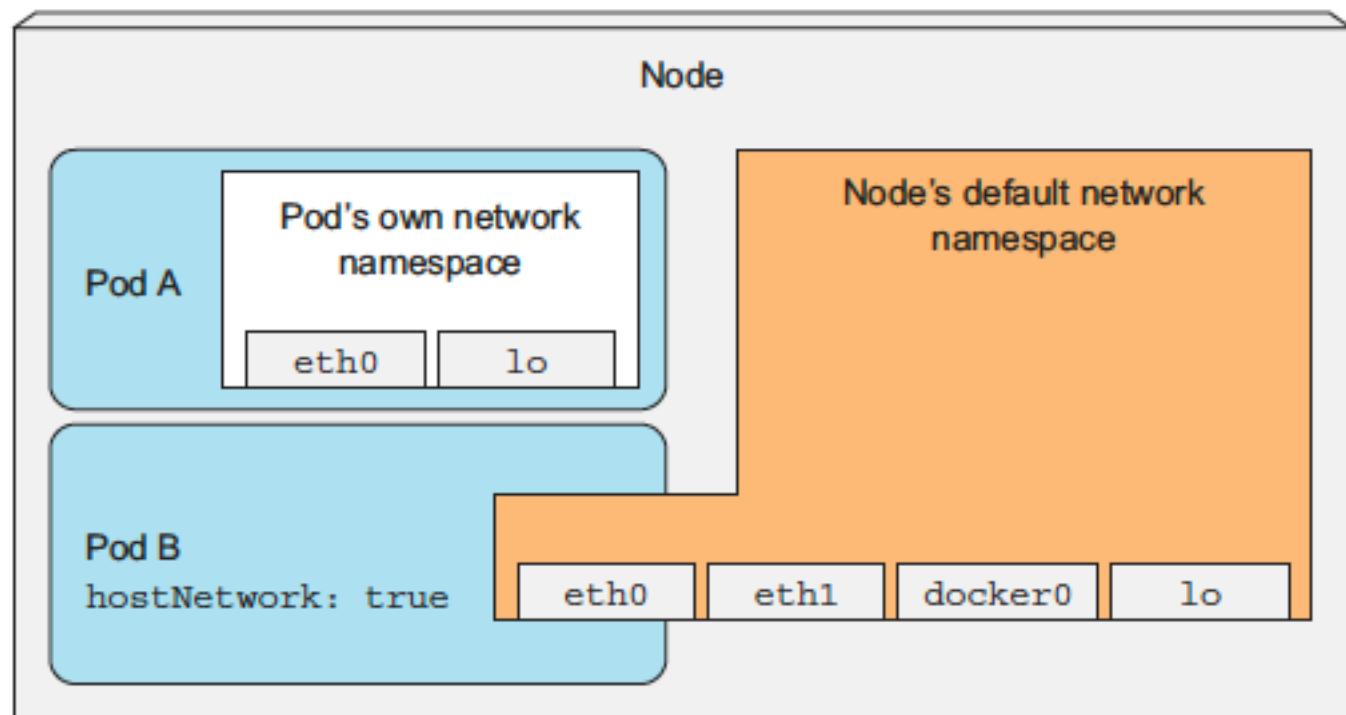


# Sandboxing

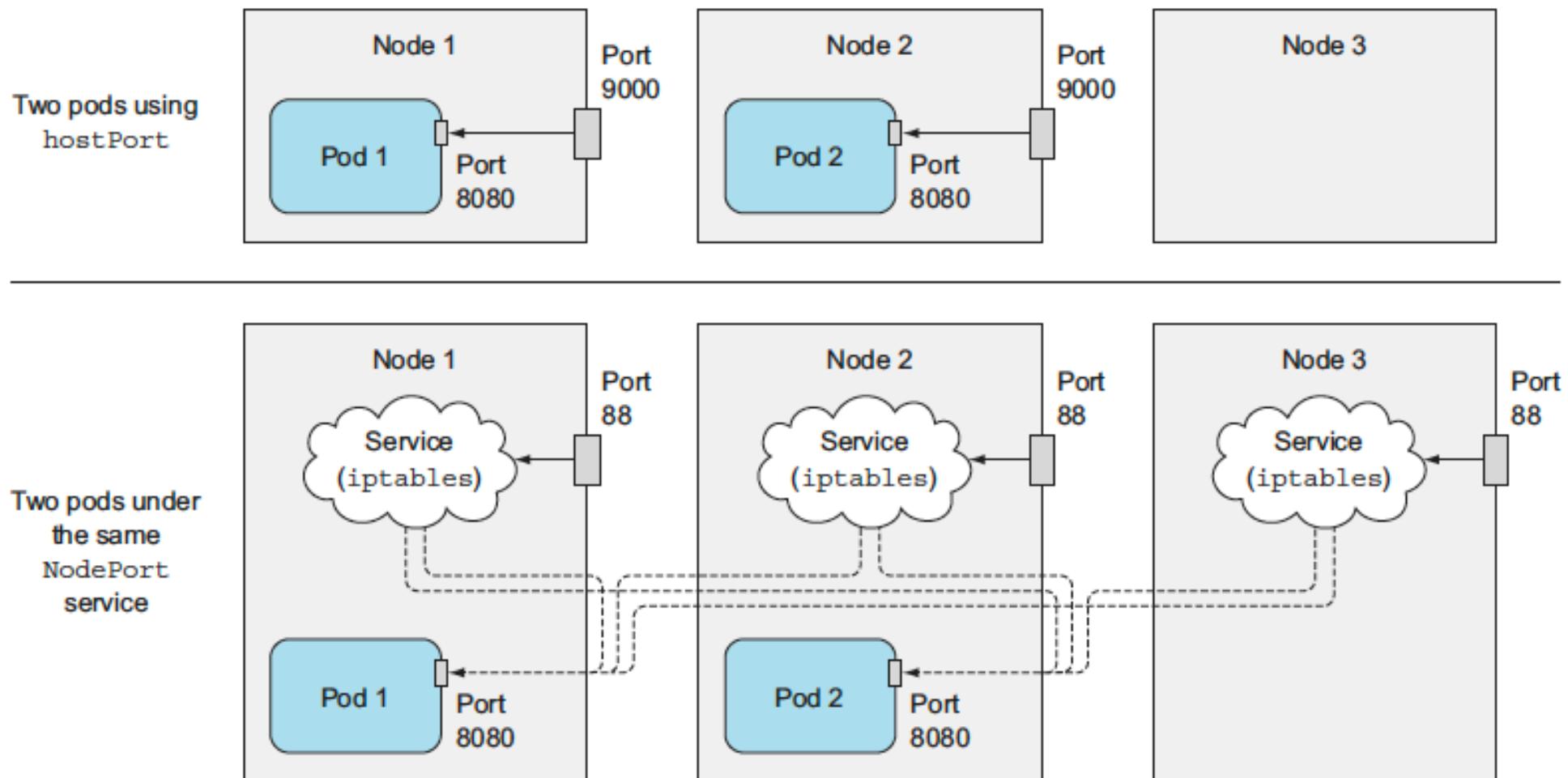


# hostNetwork

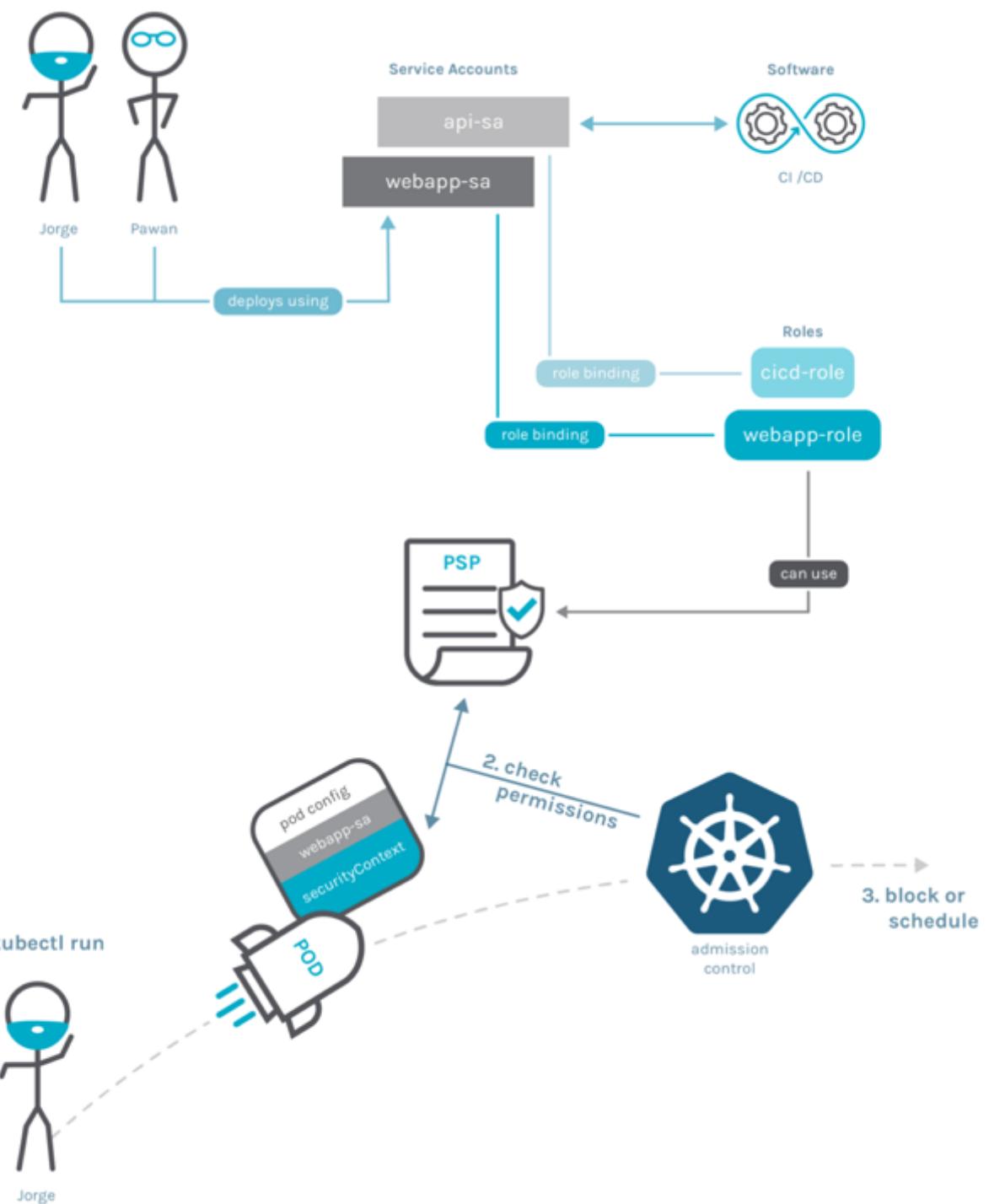
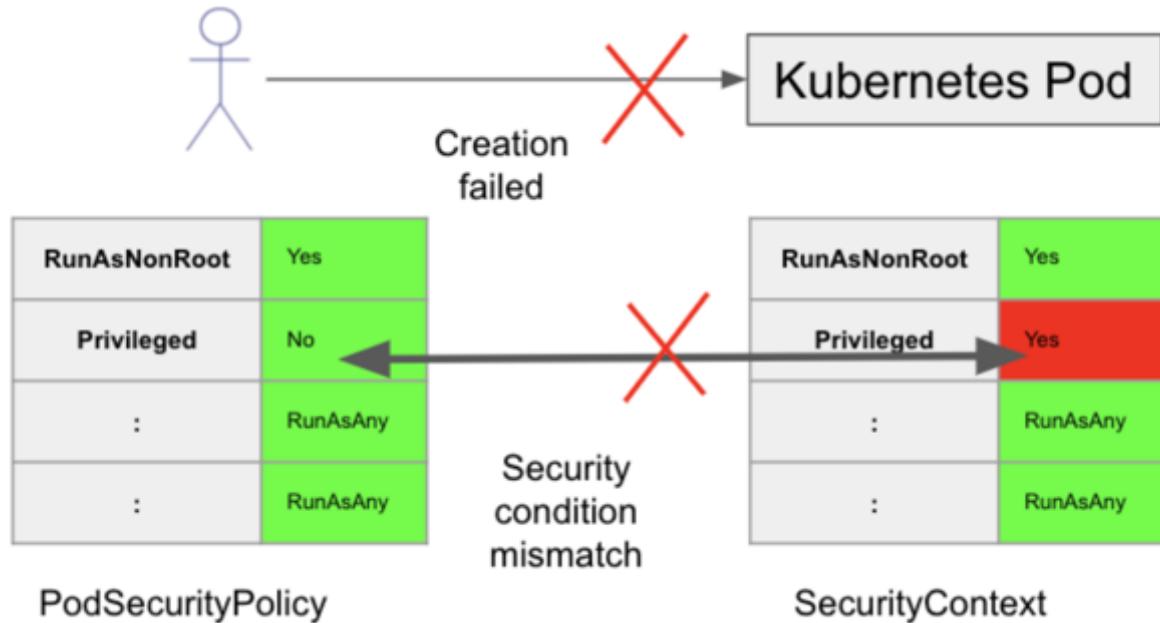
```
spec:  
  hostNetwork: true  
  containers:  
    - name: main  
      image: alpine  
      command: ["/bin/sleep"]
```



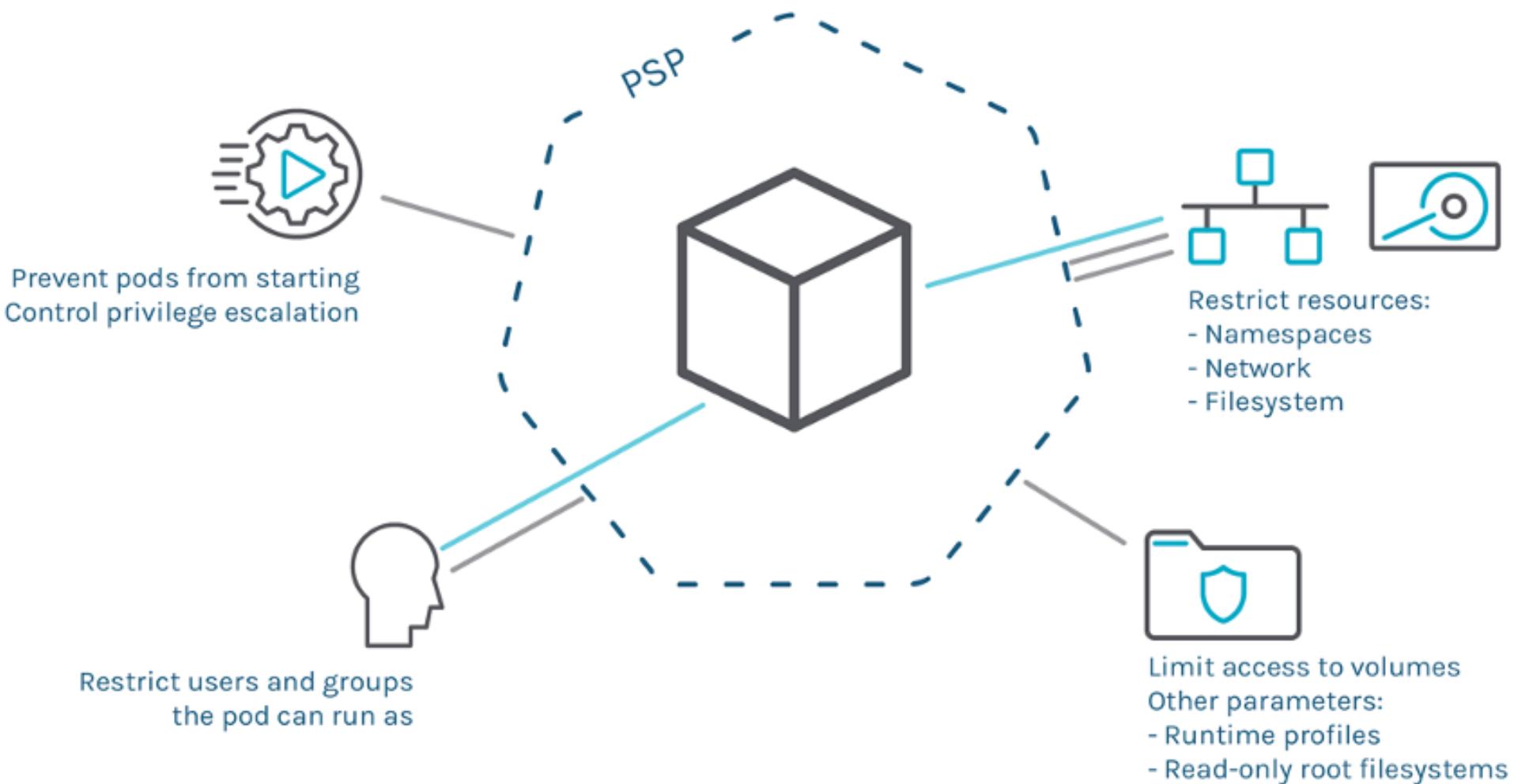
# hostPort



# Pod Security Policy



# Pod Security Policy Benefits



# PSP Examples

```
apiVersion: extensions/v1beta1
kind: PodSecurityPolicy
metadata:
  name: default
spec:
  hostIPC: false
  hostPID: false
  hostNetwork: false
  hostPorts:
    - min: 10000
      max: 11000
    - min: 13000
      max: 14000
  privileged: false
  readOnlyRootFilesystem: true
  runAsUser:
    rule: RunAsAny
  fsGroup:
    rule: RunAsAny
  supplementalGroups:
    rule: RunAsAny
  seLinux:
    rule: RunAsAny
  volumes:
    - '*'
```

Containers are forced to run with a read-only root filesystem.

Containers aren't allowed to use the host's IPC, PID, or network namespace.

They can only bind to host ports 10000 to 11000 (inclusive) or host ports 13000 to 14000.

Containers cannot run in privileged mode.

Containers can run as any user and any group.

They can also use any SELinux groups they want.

All volume types can be used in pods.

```
apiVersion: extensions/v1beta1
kind: PodSecurityPolicy
spec:
  allowedCapabilities:
    - SYS_TIME
  defaultAddCapabilities:
    - CHOWN
  requiredDropCapabilities:
    - SYS_ADMIN
    - SYS_MODULE
  ...
```

Allow containers to add the **SYS\_TIME** capability.

Automatically add the **CHOWN** capability to every container.

Require containers to drop the **SYS\_ADMIN** and **SYS\_MODULE** capabilities.

# PSP Controls (1)

제어 측면	필드 이름
특권을 가진(privileged) 컨테이너의 실행	<code>privileged</code>
호스트 네임스페이스의 사용	<code>hostPID</code> , <code>hostIPC</code>
호스트 네트워킹과 포트의 사용	<code>hostNetwork</code> , <code>hostPorts</code>
볼륨 유형의 사용	<code>volumes</code>
호스트 파일시스템의 사용	<code>allowedHostPaths</code>
특정 FlexVolume 드라이버의 허용	<code>allowedFlexVolumes</code>
파드 볼륨을 소유한 FSGroup 할당	<code>fsGroup</code>
읽기 전용 루트 파일시스템 사용 필요	<code>readOnlyRootFilesystem</code>

# PSP Controls (2)

컨테이너의 사용자 및 그룹 ID	<code>runAsUser</code> , <code>runAsGroup</code> , <code>supplementalGroups</code>
루트 특권으로의 에스컬레이션 제한	<code>allowPrivilegeEscalation</code> , <code>defaultAllowPrivilegeEscalation</code>
리눅스 기능	<code>defaultAddCapabilities</code> , <code>requiredDropCapabilities</code> , <code>allowedCapabilities</code>
컨테이너의 SELinux 컨텍스트	<code>seLinux</code>
컨테이너에 허용된 Proc 마운트 유형	<code>allowedProcMountTypes</code>
컨테이너가 사용하는 AppArmor 프로파일	어노테이션
컨테이너가 사용하는 seccomp 프로파일	어노테이션
컨테이너가 사용하는 sysctl 프로파일	<code>forbiddenSysctls</code> , <code>allowedUnsafeSysctls</code>

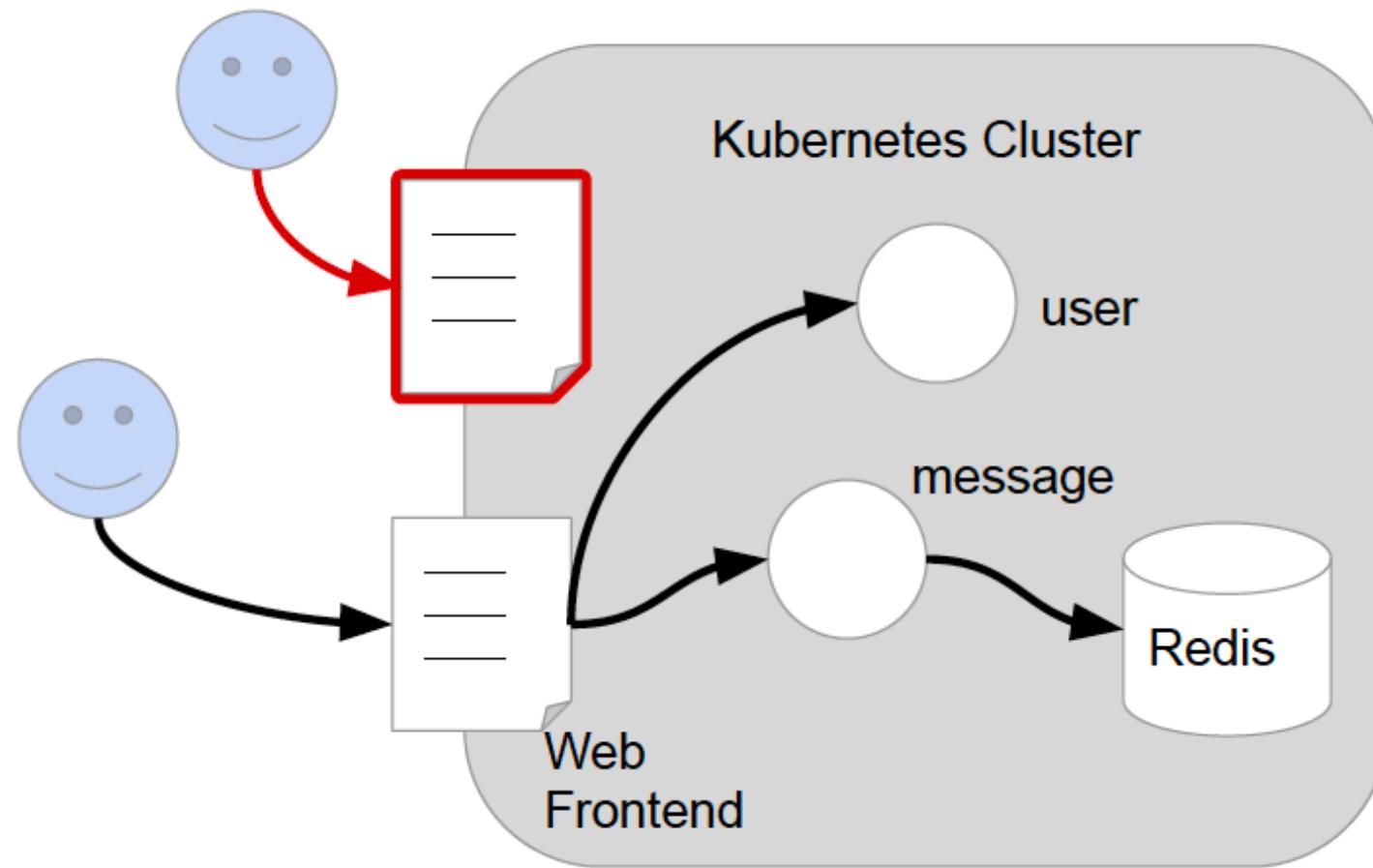
# Sysdig Secure

The screenshot shows the Sysdig Secure web interface. On the left is a dark sidebar with icons and labels: a magnifying glass for 'POLICY EVENTS', a shield for 'POLICIES', a list for 'COMMANDS AUDIT', a camera for 'CAPTURES', a checklist for 'COMPLIANCE', and a circular icon for 'IMAGE SCANNING'. The main area is titled 'KUBERNETES Pod Security Policies > New Simulation'. It features two buttons: 'Import: PSP Policy' and 'Deployment YAML'. A dropdown menu for 'kubernetes.namespace.name' is set to 'all'. Below these, a large text area displays a YAML configuration for a PodSecurityPolicy:

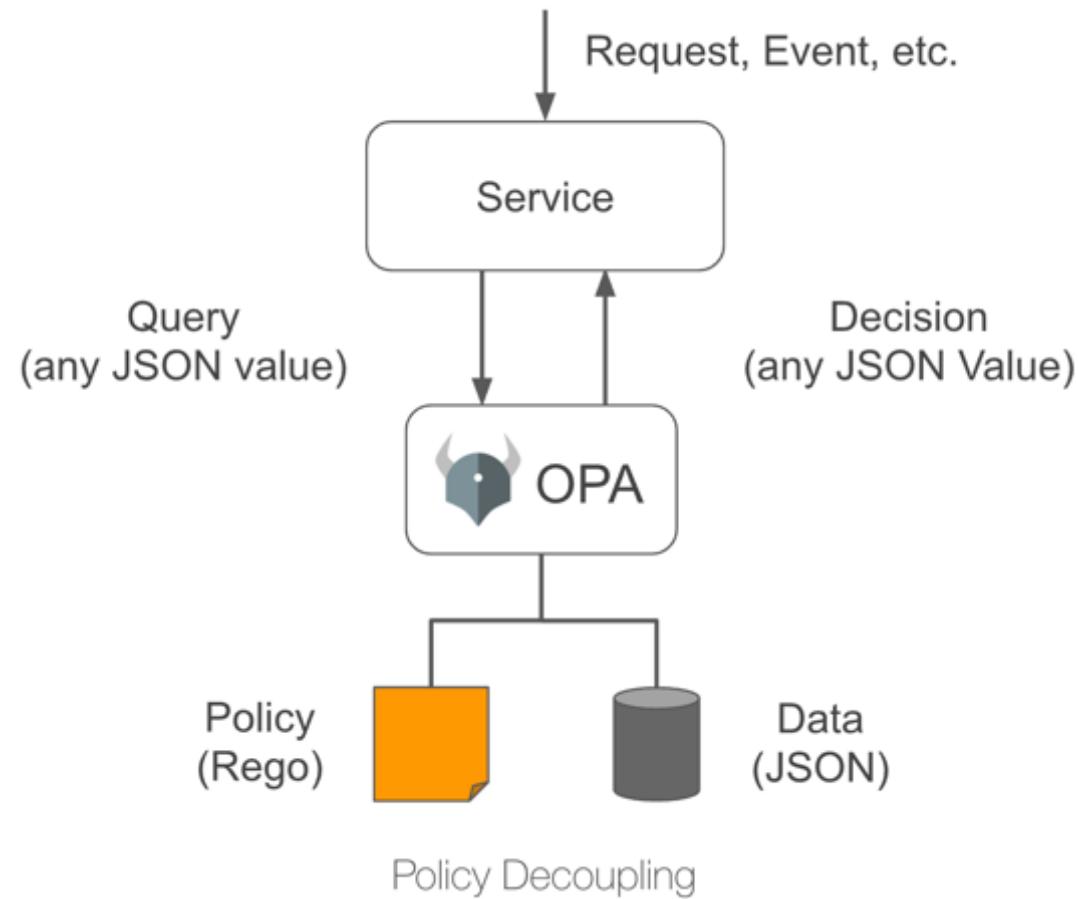
```
apiVersion: policy/v1beta1
kind: PodSecurityPolicy
metadata:
  creationTimestamp: null
  name: pod-security-policy-default-20191110234435
spec:
  allowedHostPaths:
  - pathPrefix: /etc
  fsGroup:
    rule: RunAsAny
  hostNetwork: true
  privileged: true
  runAsUser:
    rule: MustRunAs
  seLinux:
    rule: RunAsAny
  supplementalGroups:
    rule: RunAsAny
  volumes:
  - hostPath
  - secret
```

# Attack Surface: Unsecured Pod

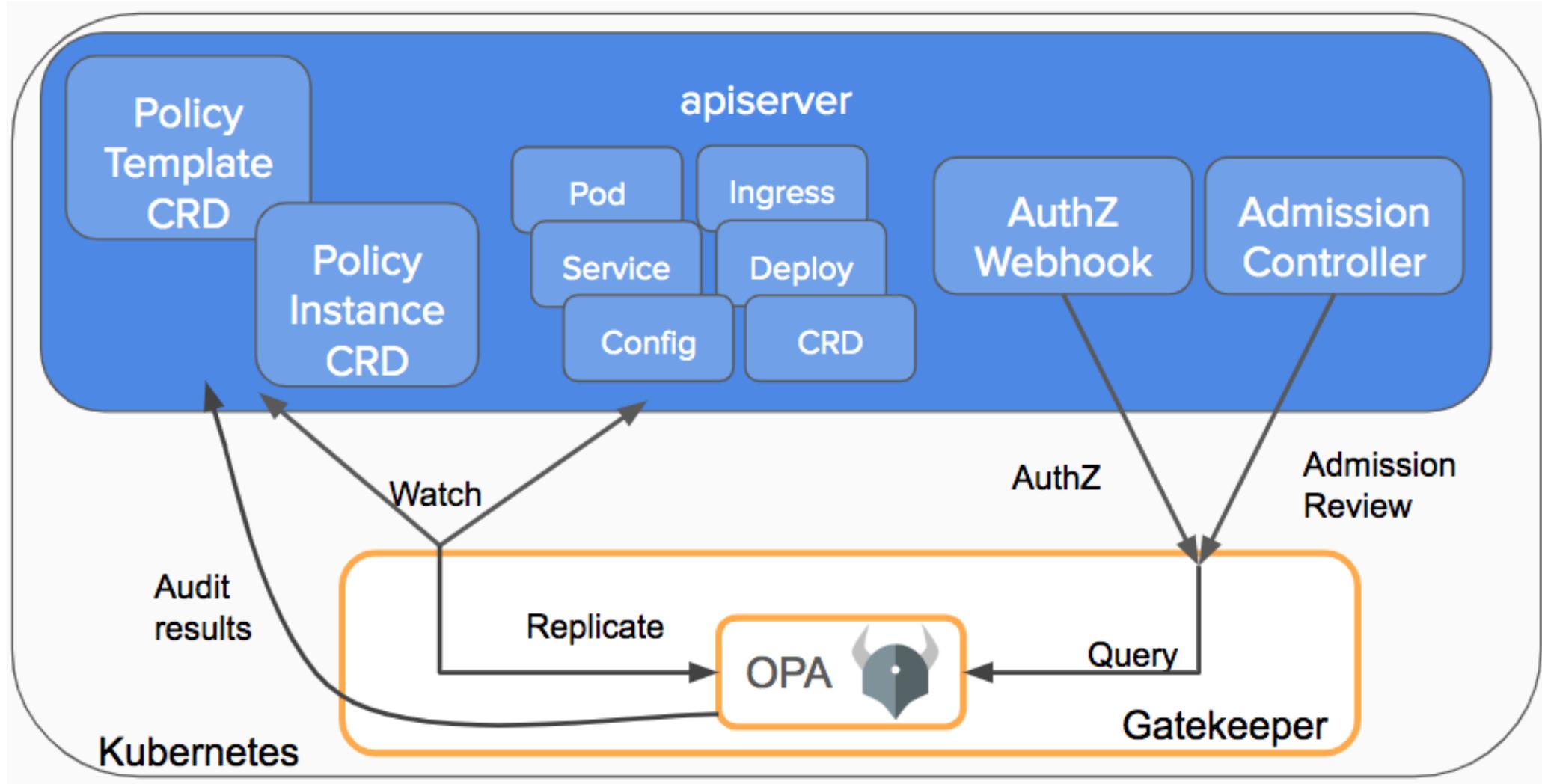
Mitigate: Open Security Policy



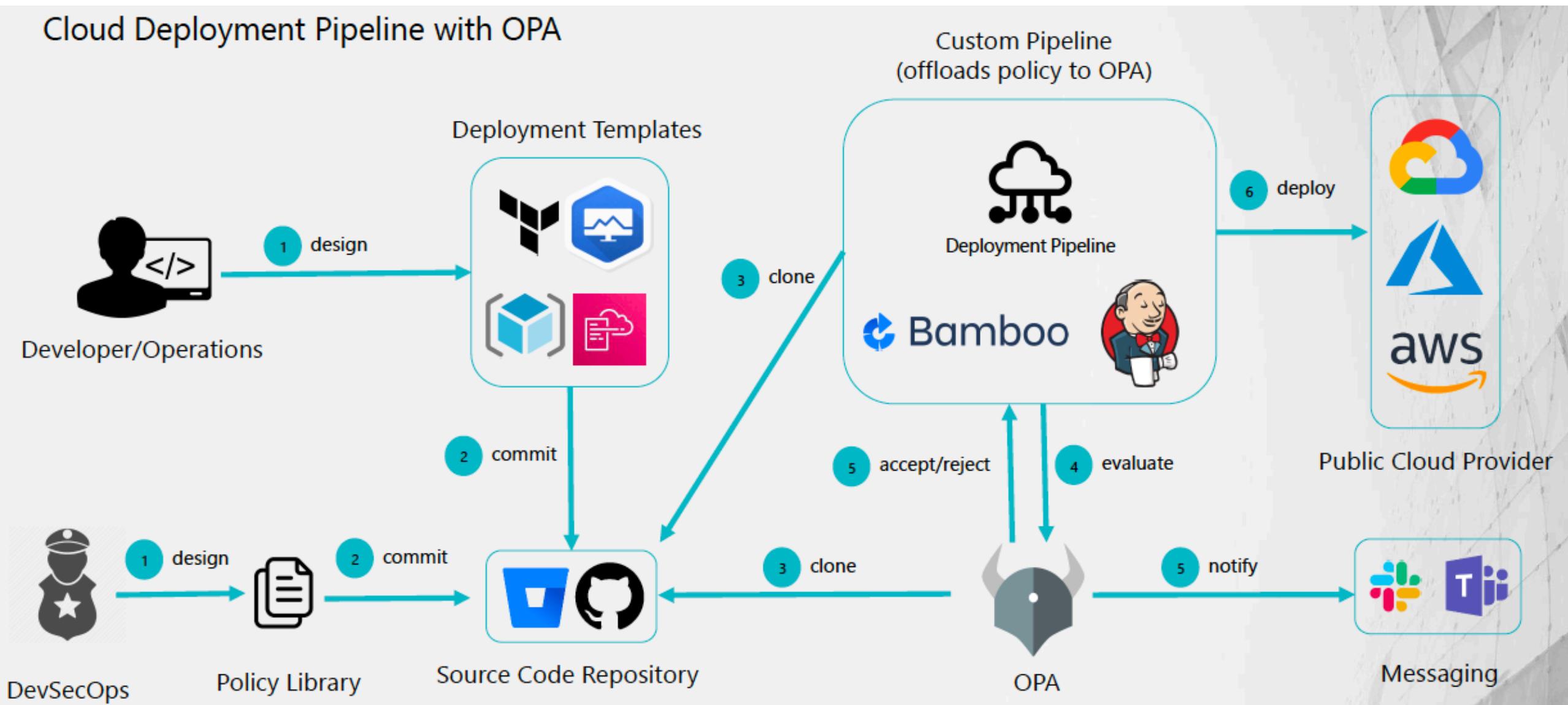
# OPA (Open Policy Agent)



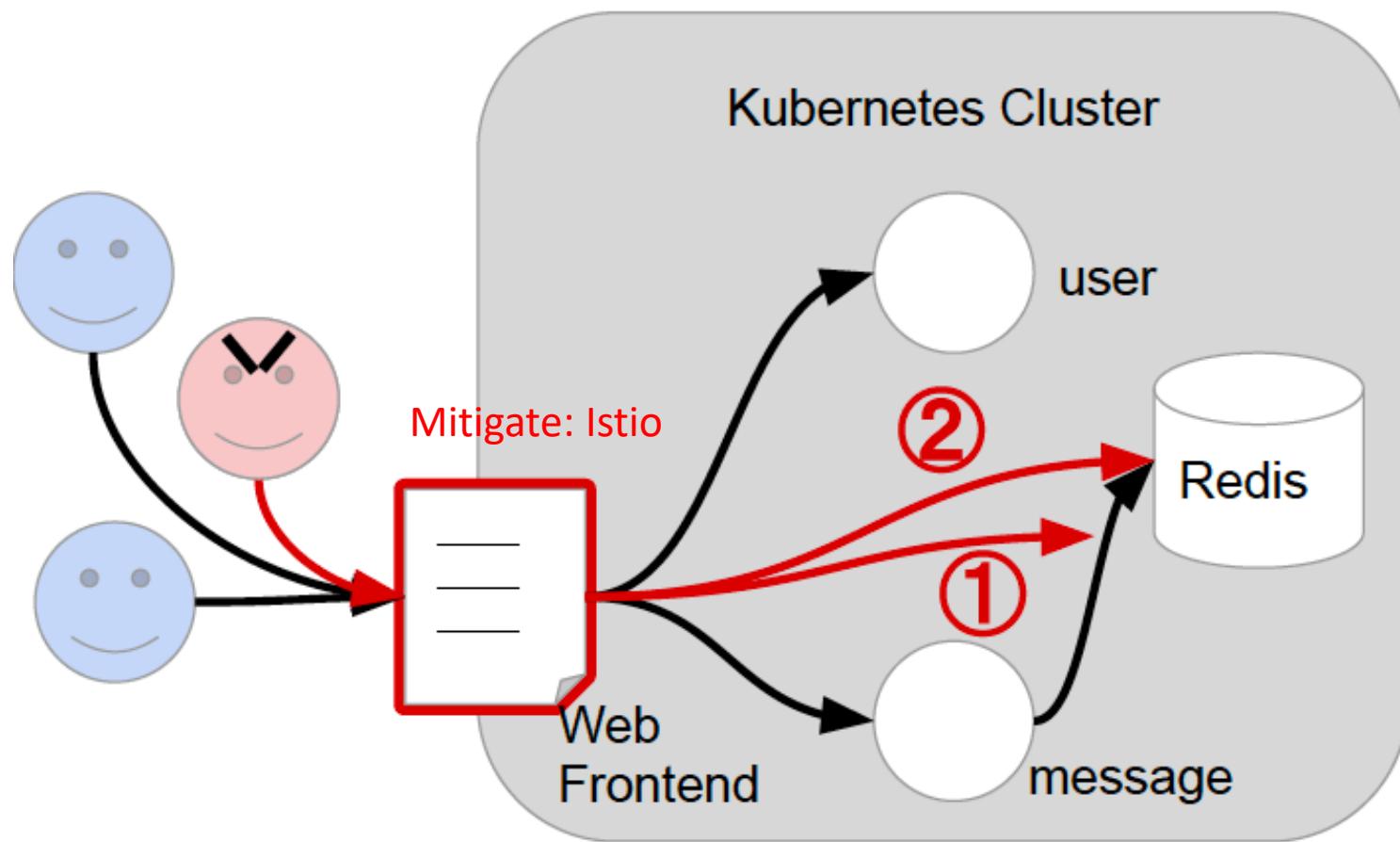
# OPA Gatekeeper



# Cloud Deployment Pipeline with OPA

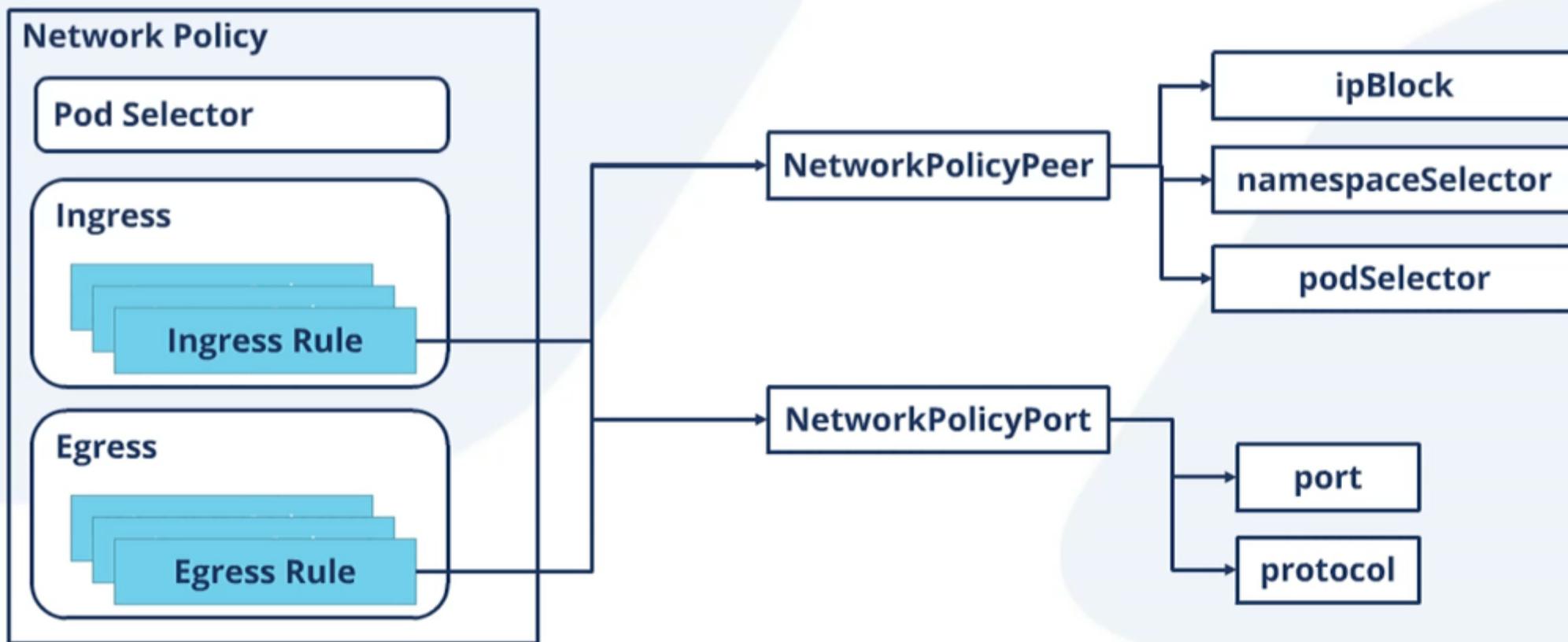


# Attack Surface: Traffic Sniffing

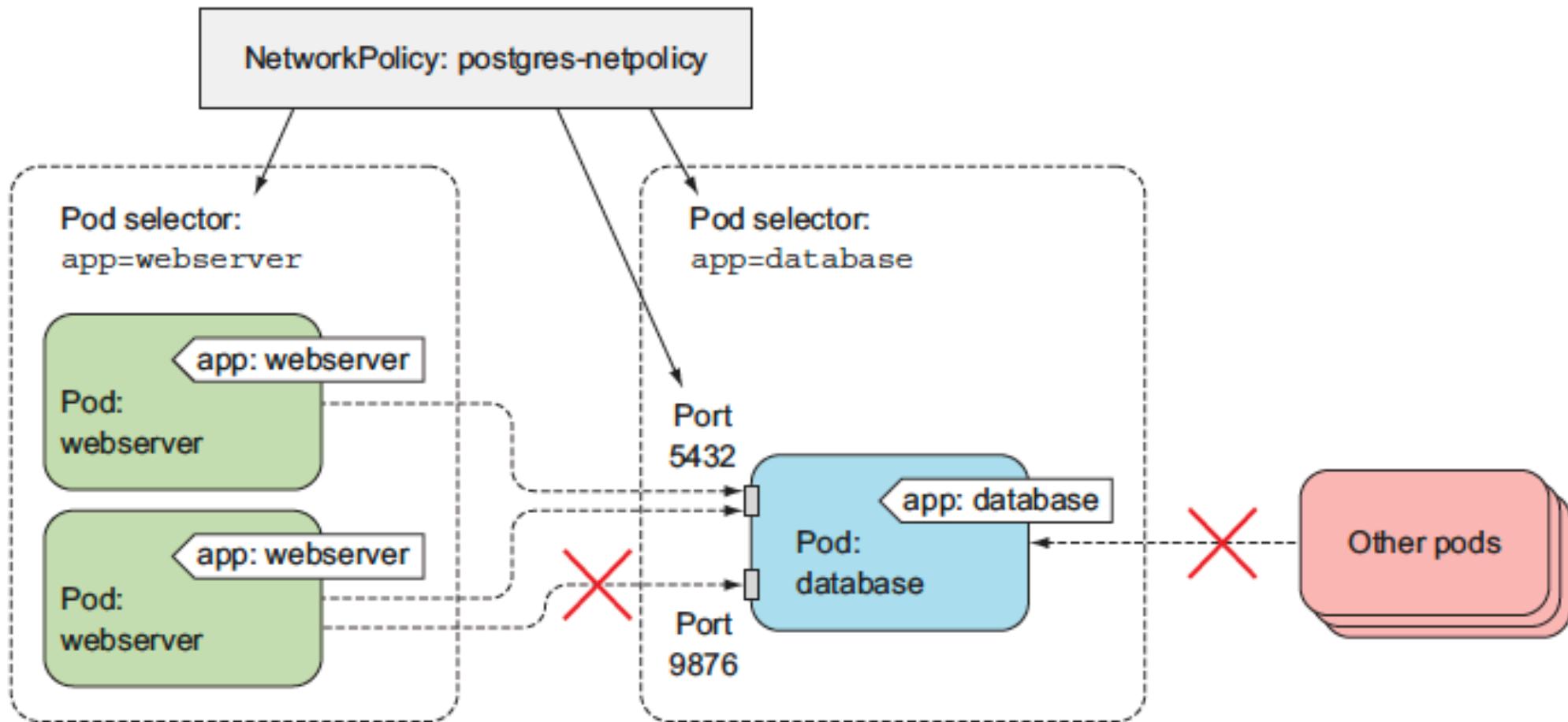


# Network Policy

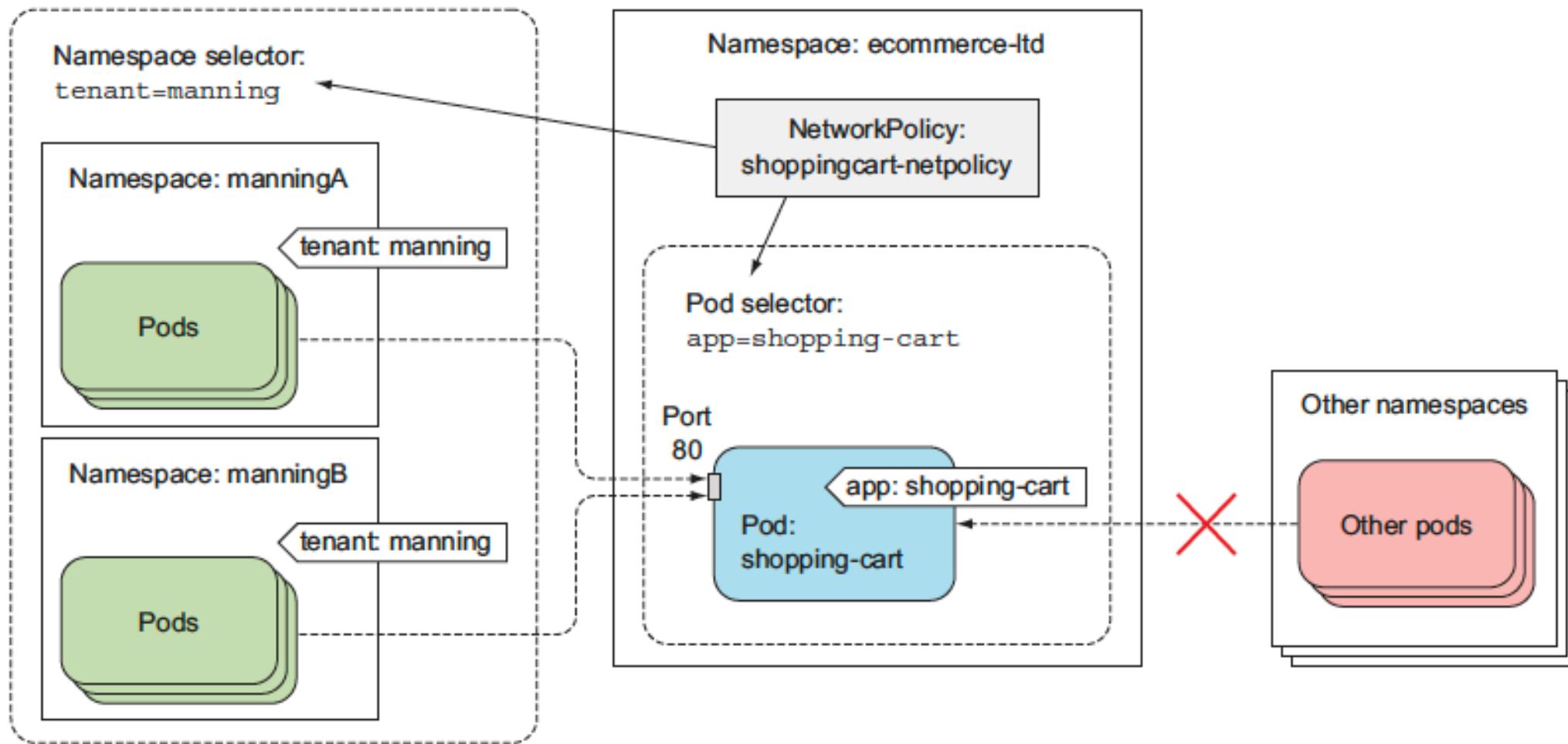
- Rules contain peers (the other side of the connection) and ports



# Network Policy: Pod Selector

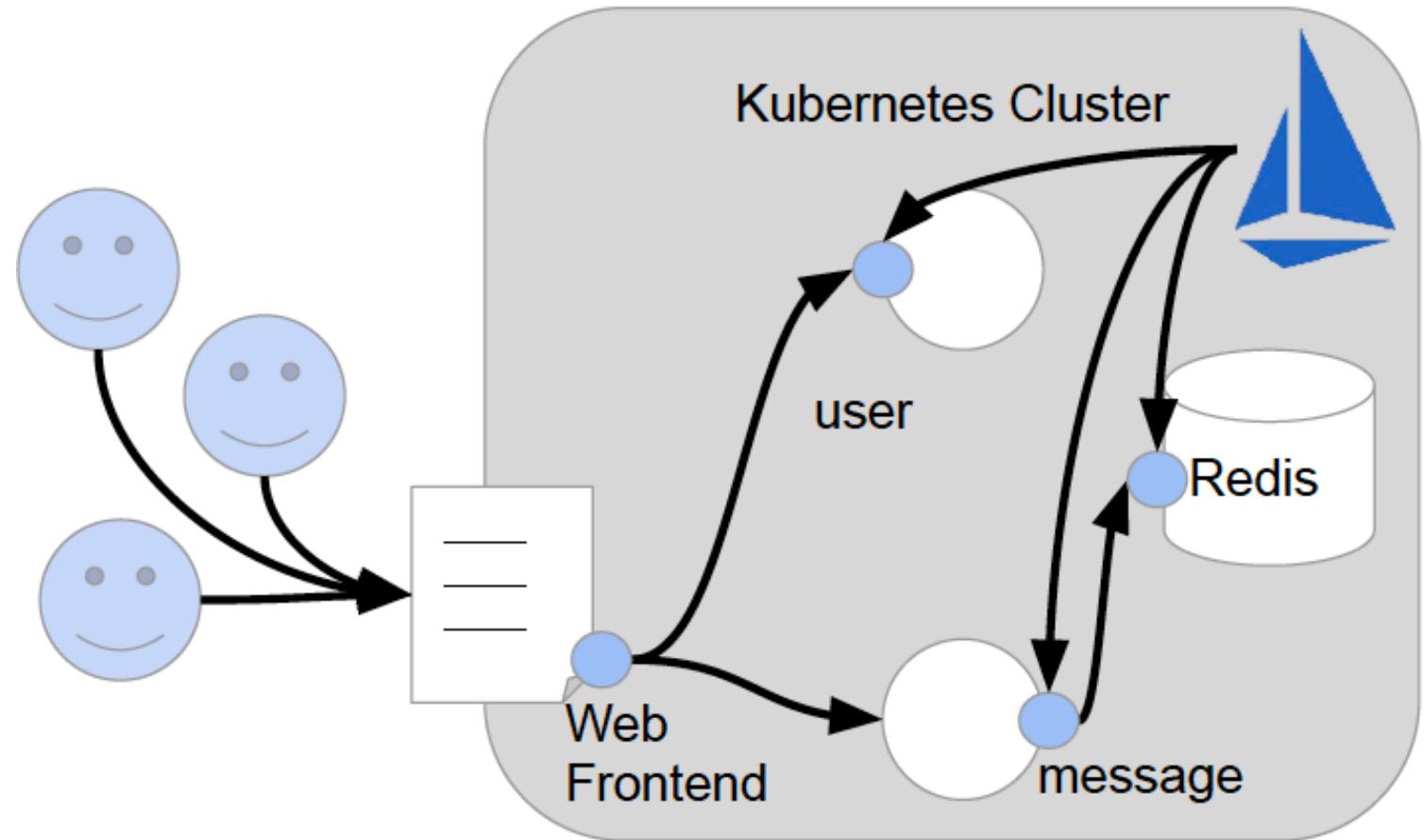


# Network Policy: Namespace Selector

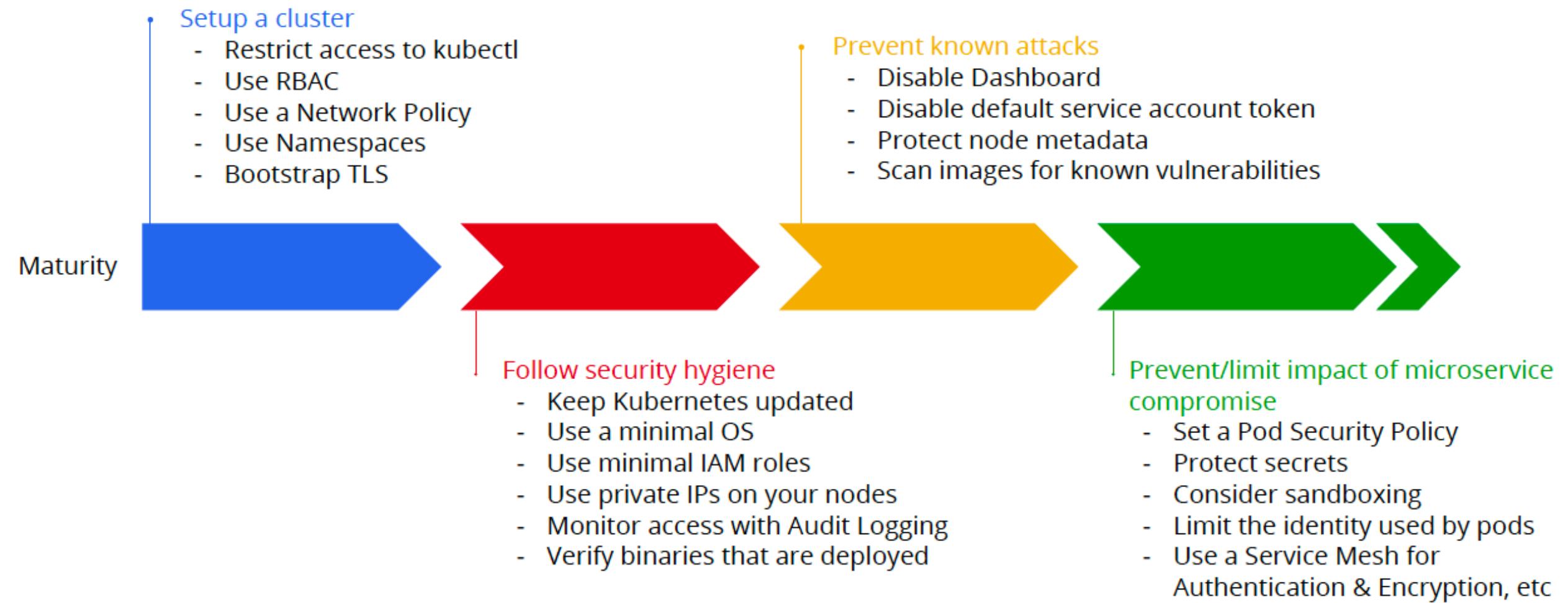


# Istio

1. Proxy data between services
2. End-to-end encryption
3. Rolling certificates
4. Policy managed by central server

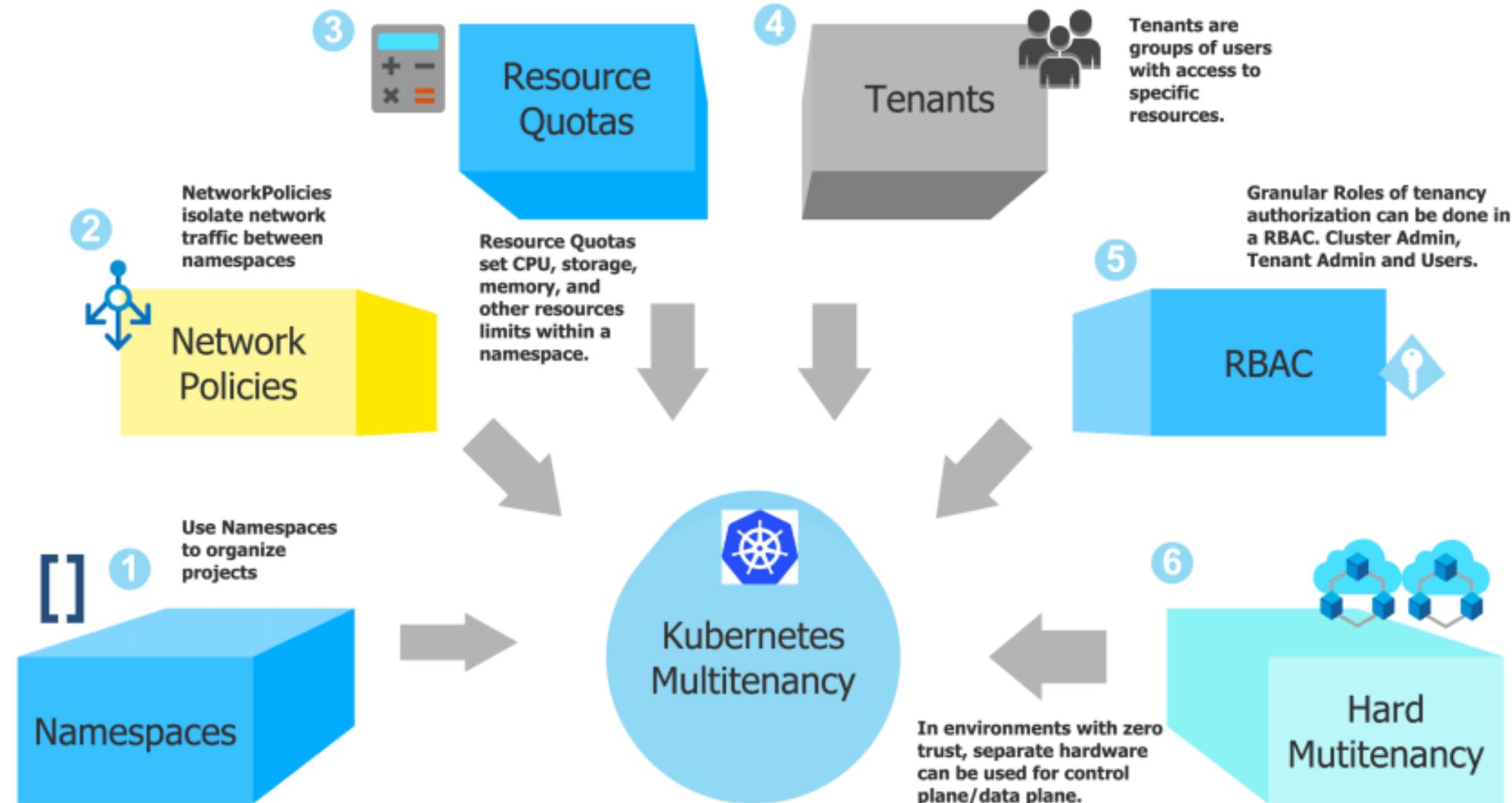


# Kubernetes Security Journey

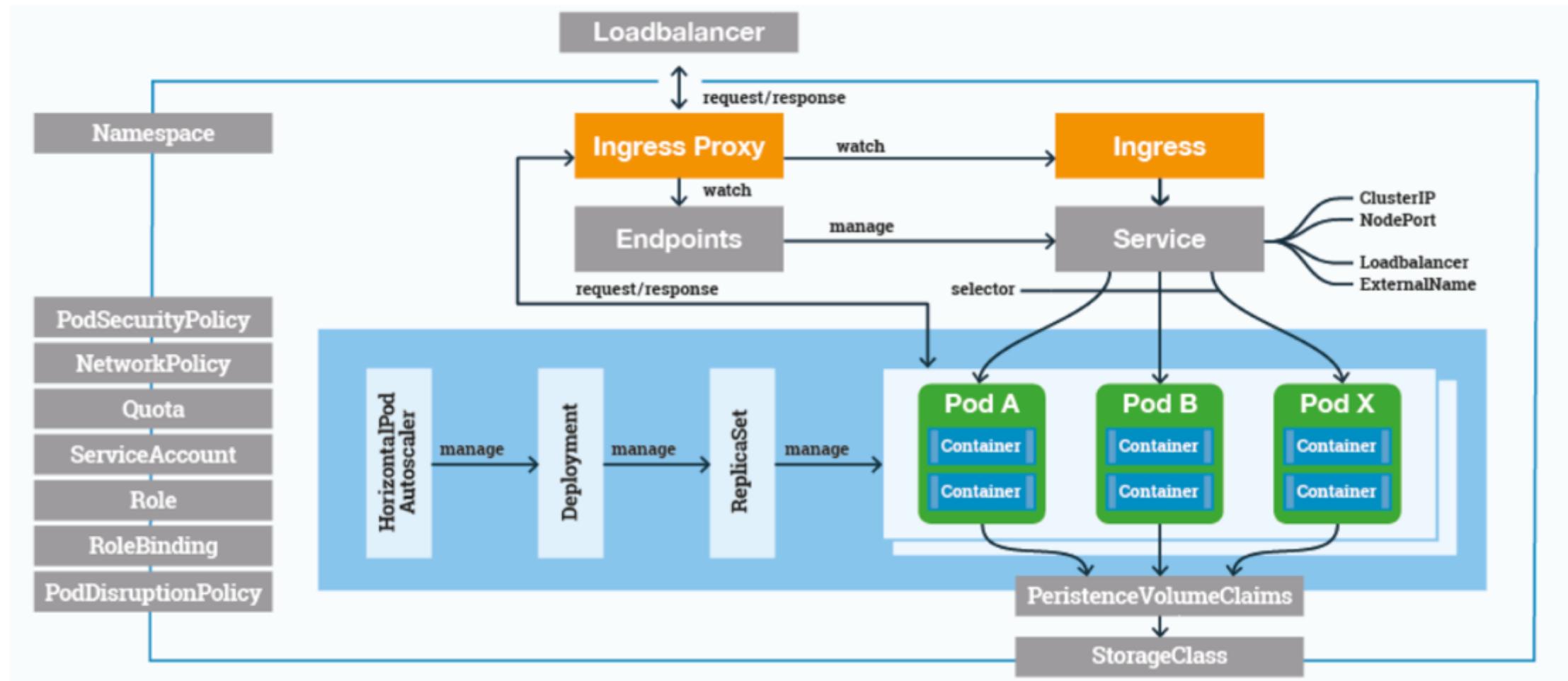


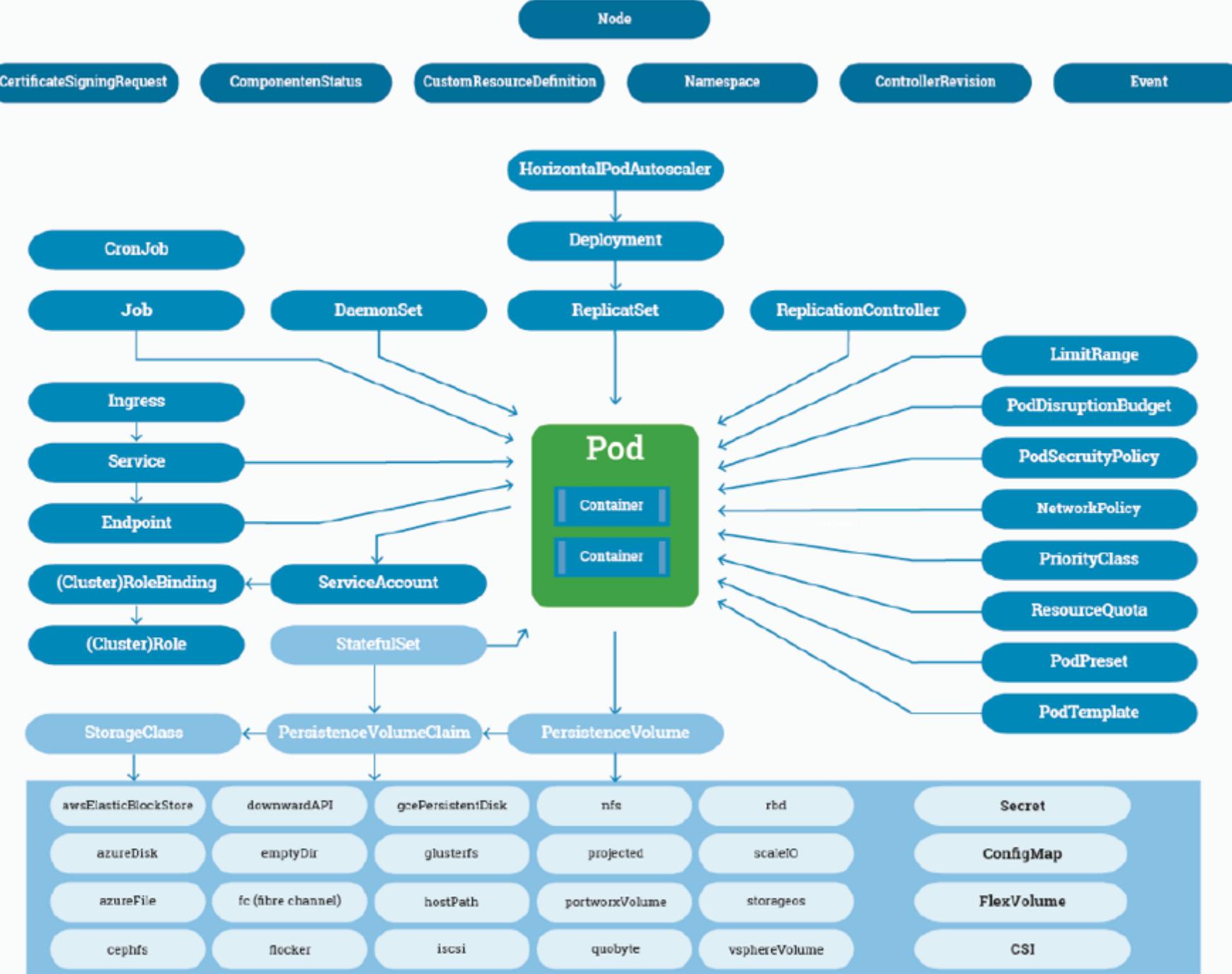


# Kubernetes Multi-Tenancy Challenge



# Kubernetes Service Architecture





# Local Kubernetes Cluster with VirtualBox

