# Building a Kubernetes platform

Chris Nesbitt-Smith

learn**k8s**

# Chris Nesbitt-Smith

UK Gov | LearnK8s | Control Plane | lots of open source

learnk8s

# Multi-tenancy in Kubernetes

# Multi-tenancy in Kubernetes

## Isolation

**Ease of management**

**Cost efficiency**

# Multi-tenancy in Kubernetes

Isolation

## Ease of management

Cost efficiency
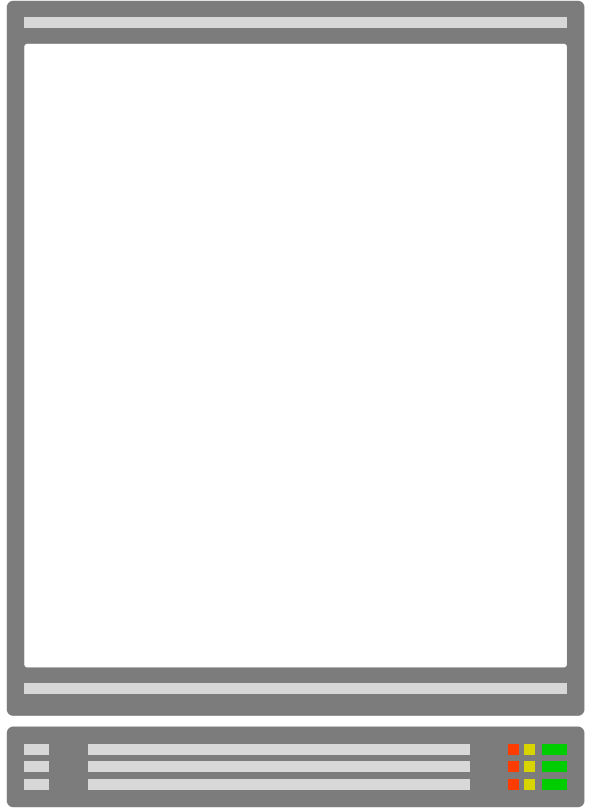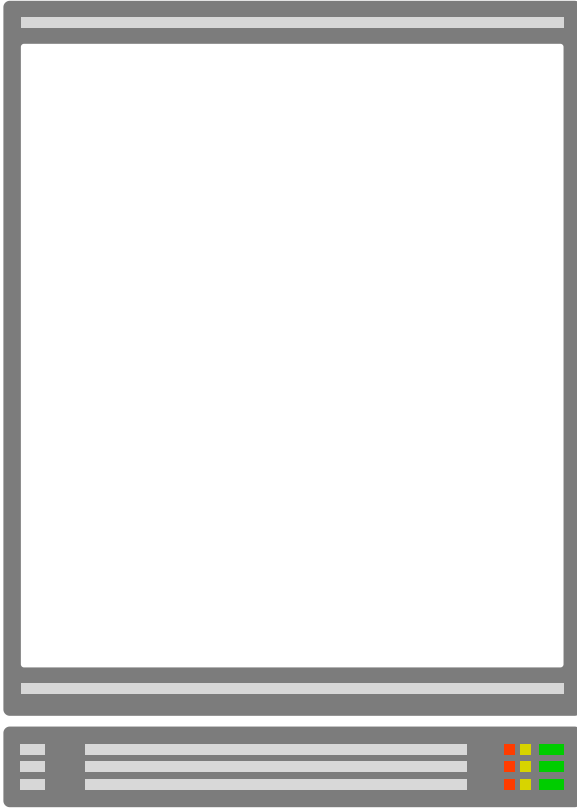
# Multi-tenancy in Kubernetes
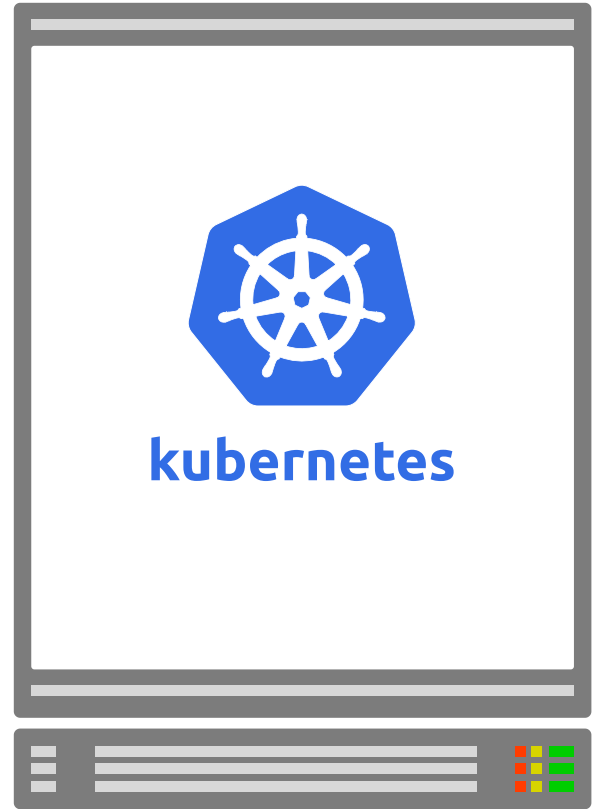
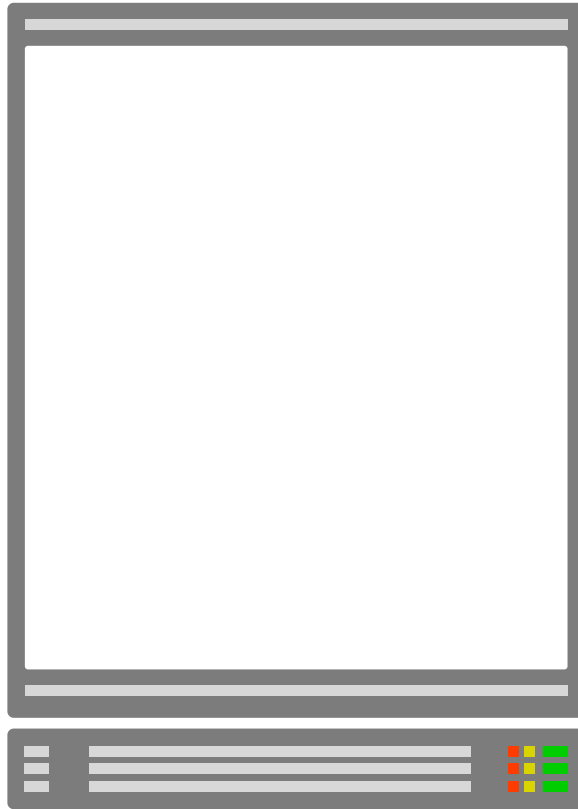Isolation

Ease of management

## Cost efficiency

# Datacentre as a single computer

01

kubernetes

**Worker Node** **Worker Node**

**Worker Node** **Worker Node**

# Namespaces

v1.0.0

v2.0.0

v3.0.0

namespace namespace namespace

v1.0.0 v2.0.0 v3.0.0

# Environments x tenants

dev     test     prod

v1.0.0     v2.0.0     v3.0.0

**Team A**

dev

test

prod

|  | **Team A** | **Team B** |
|------|:----------:|:----------:|
| **dev** | | |
| **test** | | |
| **prod** | | |

|  | Team A | Team B | Team C |
|---|---|---|---|
| **dev** | | | |
| **test** | | | |
| **prod** | | | |

# Environments x tenants *at scale*

**10 TENANTS**

dev

test

prod

**10 TENANTS**

dev

test

prod

**50 TENANTS**

dev
test
prod

dev
test
prod

traffic

NGINX

ClusterIP

ClusterIP

Pod A

Pod A

Pod A
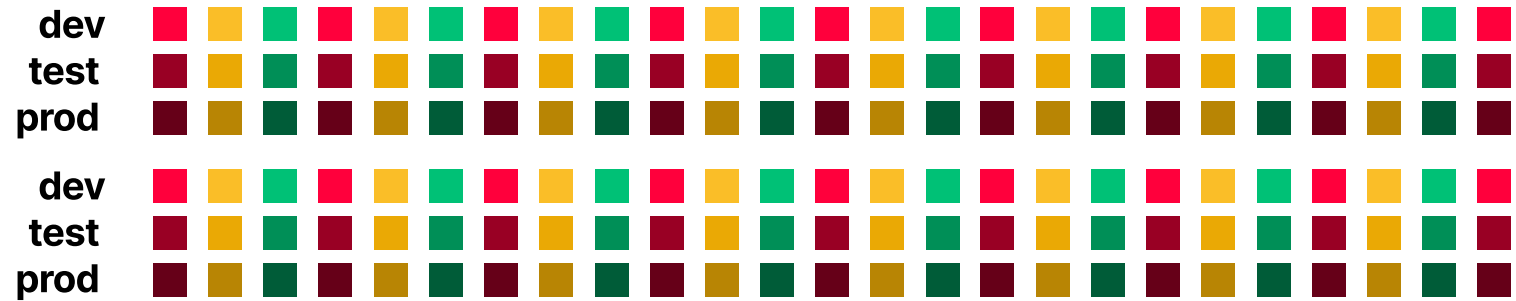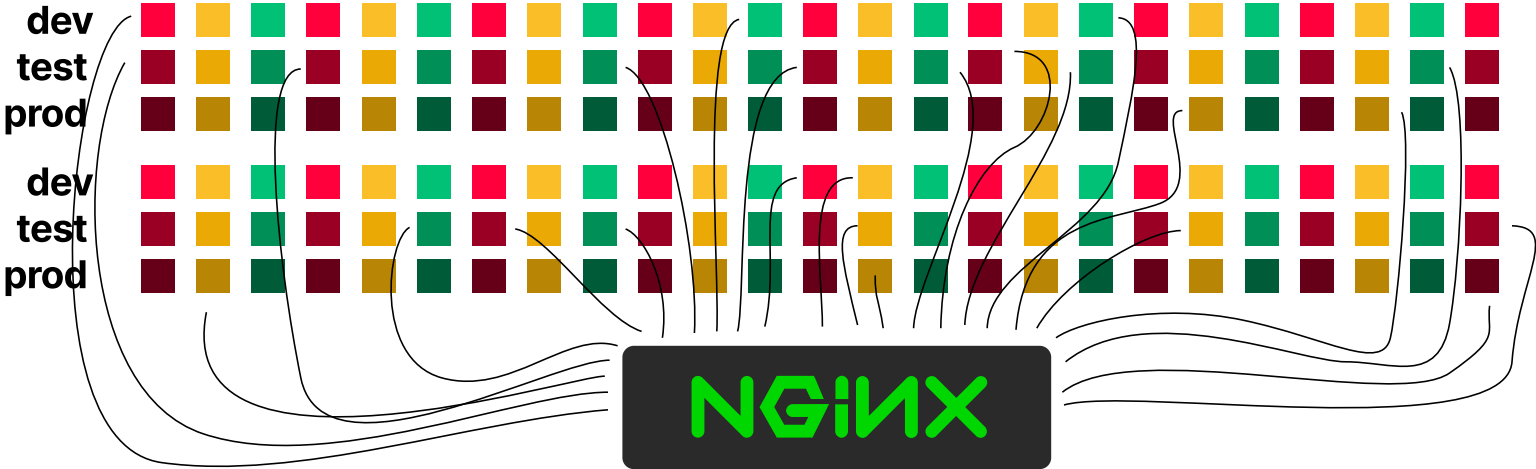
Pod B

cluster

**NGINX**

10 TENANTS

dev test prod

50 TENANTS

dev test prod

dev test prod

**10 TENANTS**

dev
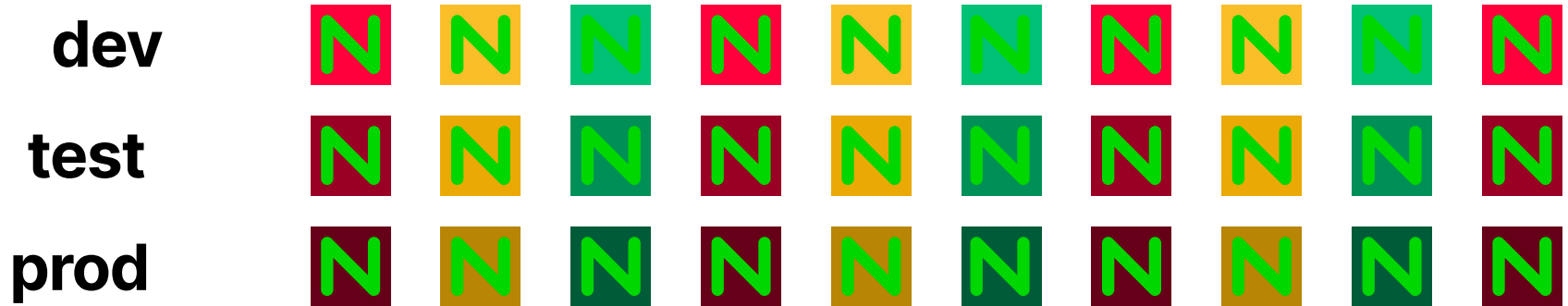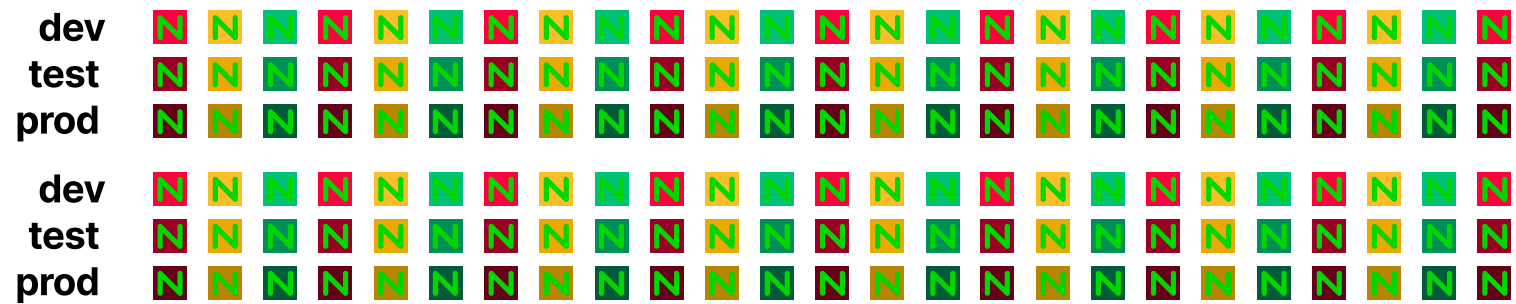
test

prod

**50 TENANTS**

dev
test
prod

dev
test
prod

# 1 vs many: resources

```
~$ cat values.yaml

...
resources:
  requests:
    cpu: 100m
    memory: 90Mi
```

ingress-nginx Helm chart

# Single Ingress

**CPU**

# 100m

**MEMORY**

# 90Mi

**Single Ingress**

**10 × 3**

CPU

**100m**

CPU

**3vCPU**

MEMORY

**90Mi**

MEMORY

**2.7GB**

| Single Ingress | 10 × 3 | 50 × 3 |
| --- | --- | --- |
| **CPU** | **CPU** | **CPU** |
| **100m** | **3vCPU** | **5vCPU** |
| **MEMORY** | **MEMORY** | **MEMORY** |
| **90Mi** | **2.7GB** | **4.5GB** |

| Instance Size | vCPU | Memory (GiB) | Instance Storage (GB) | Network Bandwidth (Gbps)*** | EBS Bandwidth (Gbps) |
|---|---|---|---|---|---|
| c6i.large | 2 | 4 | EBS-Only | Up to 12.5 | Up to 10 |
| c6i.xlarge | 4 | 8 | EBS-Only | Up to 12.5 | Up to 10 |
| c6i.2xlarge | 8 | 16 | EBS-Only | Up to 12.5 | Up to 10 |
| c6i.4xlarge | 16 | 32 | EBS-Only | Up to 12.5 | Up to 10 |
| c6i.8xlarge | 32 | 64 | EBS-Only | 12.5 | 10 |

# $0.34/hr    $248.2/m

# 1 vs many: config

traffic

NGINX

ClusterIP

ClusterIP

Pod A

Pod A

Pod A

Pod B

10s
KEEP ALIVE

cluster

```yaml
kind: ConfigMap
apiVersion: v1
metadata:
  name: nginx-configuration
data:
  keep-alive: "10s"
  proxy-read-timeout: "10s"
  client-max-body-size: "2m"
```
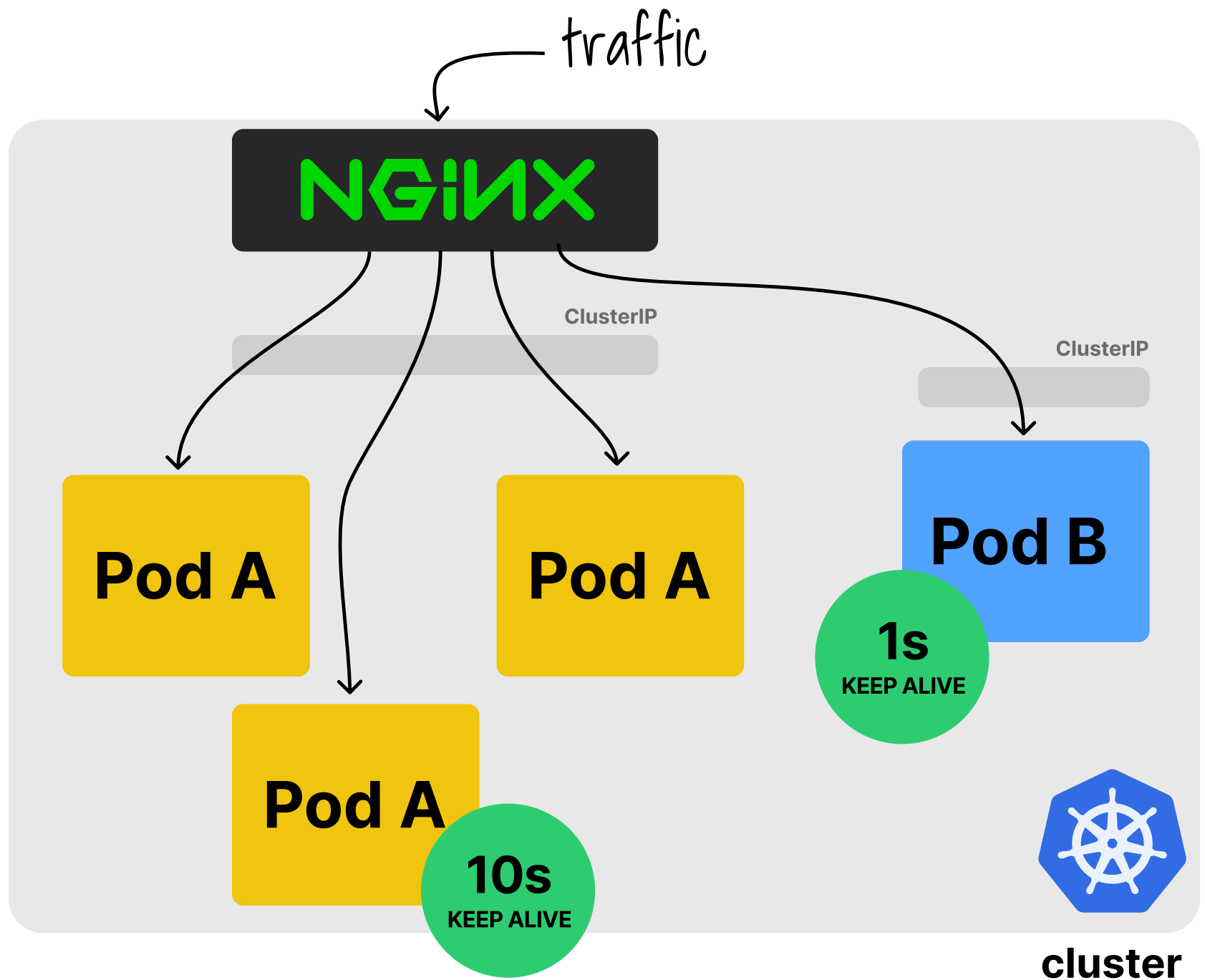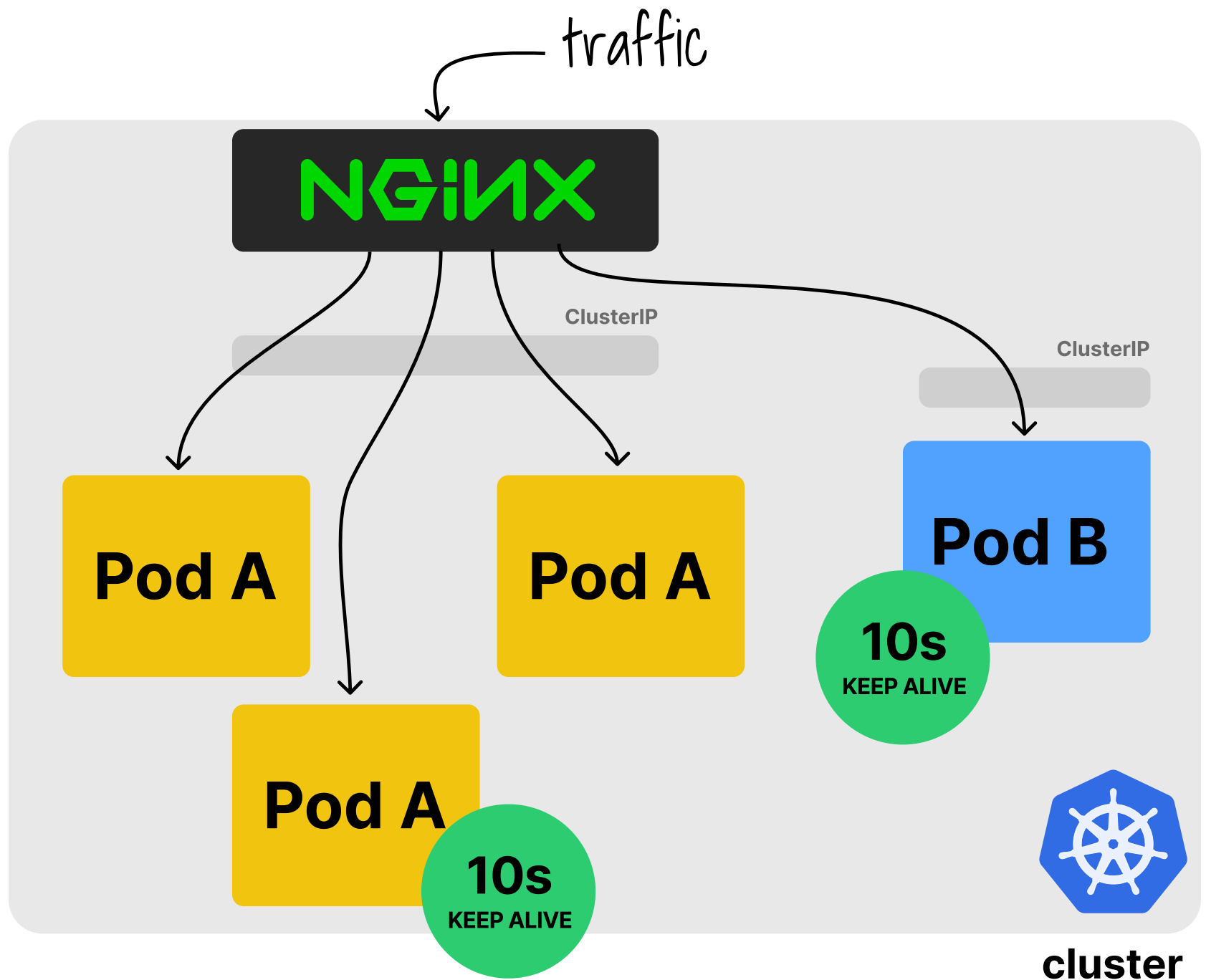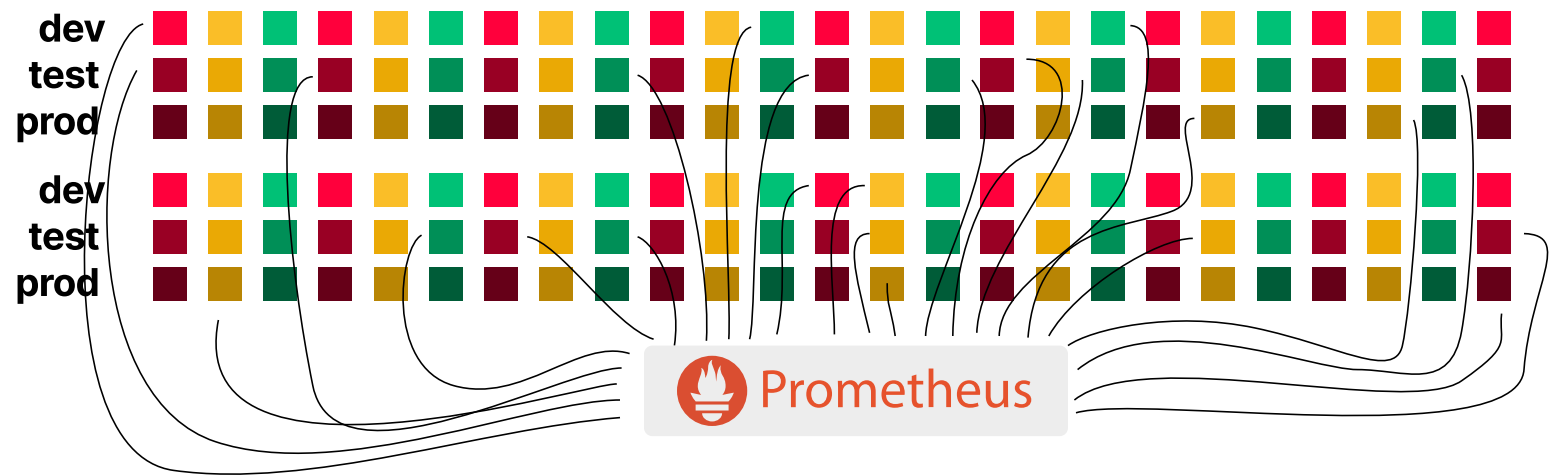
global setting

traffic

NGINX

ClusterIP

ClusterIP

Pod A

Pod A

Pod A

10s
KEEP ALIVE

Pod B

10s
KEEP ALIVE

cluster

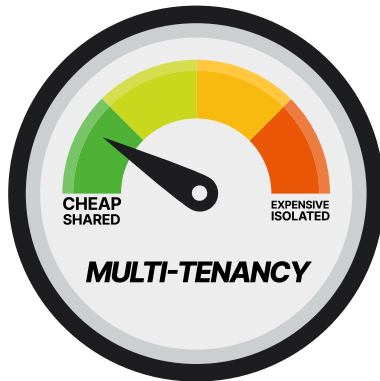CHEAP
SHARED

EXPENSIVE
ISOLATED
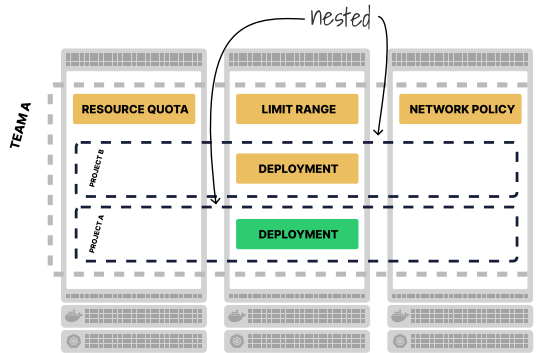
MULTI-TENANCY

# Kubernetes platform tools

**HOON JO**

# Comparing multi-tenancy tools

# Hierarchial Namespace Controller



nested

TEAM A

| RESOURCE QUOTA | LIMIT RANGE | NETWORK POLICY |

PROJECT B

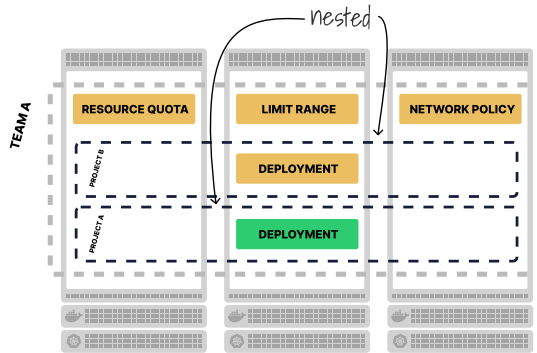DEPLOYMENT

PROJECT A

DEPLOYMENT



CHEAP
SHARED

EXPENSIVE
ISOLATED

**MULTI-TENANCY**

# Hierarchial Namespace Controller

# vCluster

# Karmada



nested

RESOURCE QUOTA | LIMIT RANGE | NETWORK POLICY

DEPLOYMENT

DEPLOYMENT

TEAM A

PROJECT B

PROJECT A

copy pod spec

namespace A

v1.0.0

AGENT

Pod3

API SERVER

CONTROLLER MANAGER

SCHEDULER

CHEAP SHARED — EXPENSIVE ISOLATED

MULTI-TENANCY

CHEAP SHARED — EXPENSIVE ISOLATED

MULTI-TENANCY

CHEAP SHARED — EXPENSIVE ISOLATED

MULTI-TENANCY

# Hierarchical Namespace Controller

<hr />

root namespace

root namespace

RESOURCE QUOTA

ROLE

child 1

child 2

# Demo

# Hierarchial Namespace Controller

# "Nested" namespaces

## Hierarchial Namespace Controller

# "Nested" namespaces
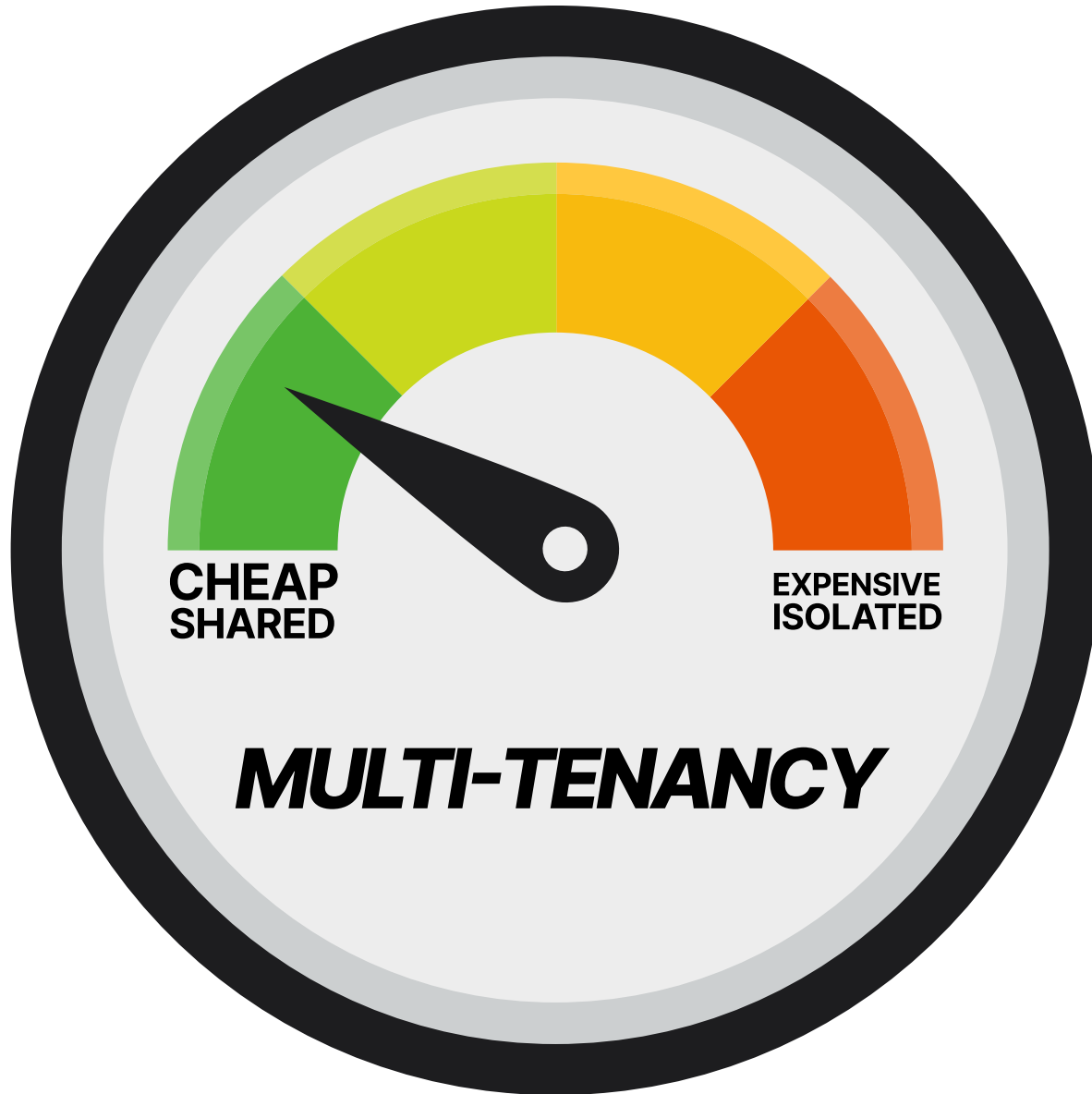# Single controller

# Hierarchial Namespace Controller

# "Nested" namespaces
# Single controller
# Regular namespaces

~$0

# HNC & Roles

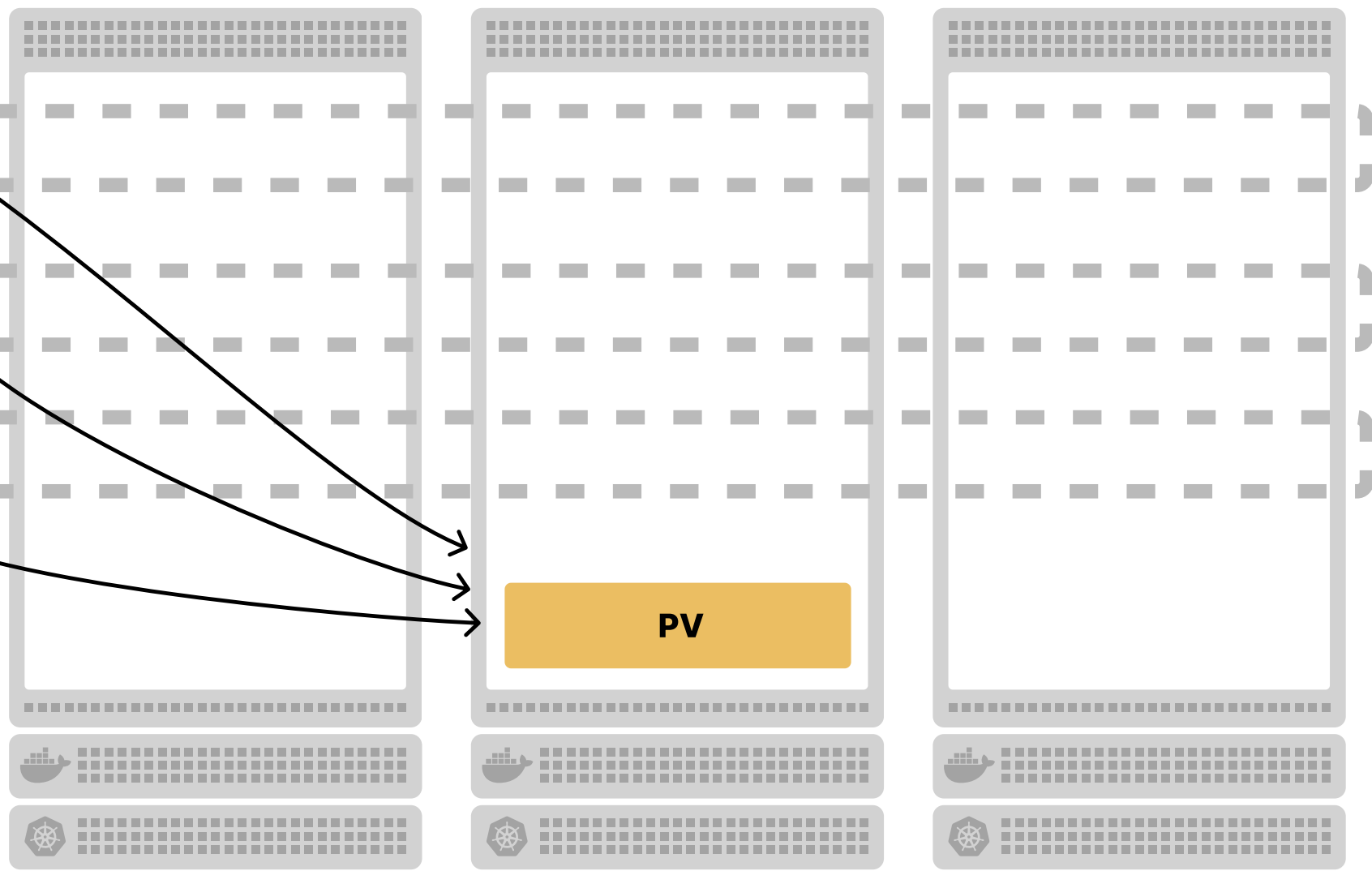| UID | ROLE  | PODS read | write | PVs read | write |
|-----|-------|-----------|-------|----------|-------|
| 1   | teamA | ✓         | ✓     | ✓        | ✓     |

root namespace **ROLE**

TEAM A

TEAM B

TEAM C

**PV**

TEAM A
TEAM B
TEAM C

PV

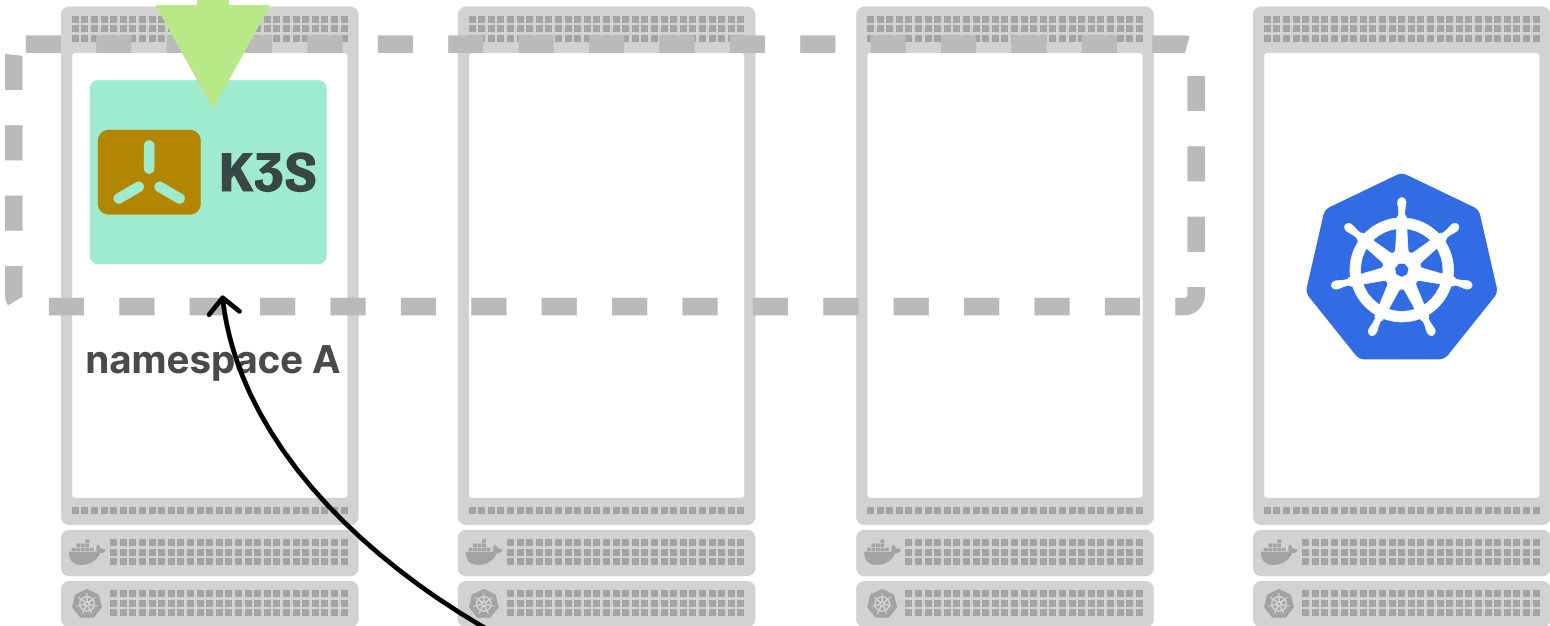# Isolating control planes

master

namespace A

Just a regular pod

master

**K3S**

namespace A

kubectl apply

master

K3S

namespace A

it has no nodes!

kubectl apply

the pods are in (wrong) db!

master

API SERVER

CONTROLLER MANAGER

SCHEDULER

namespace A

it has no nodes!

kubectl apply

copy pod spec

master

API SERVER

CONTROLLER MANAGER    SCHEDULER

namespace A

API SERVER

CONTROLLER MANAGER    SCHEDULER

kubectl apply

master

namespace A

v1.0.0

# vCluster and global resources

kubectl apply -f my-pv.yaml

master

TEAM A NAMESPACE

TEAM B NAMESPACE

API SERVER

CONTROLLER MANAGER

SCHEDULER

API SERVER

CONTROLLER MANAGER

SCHEDULER

vCluster
made by loft

vCluster
made by loft

API SERVER

CONTROLLER MANAGER

SCHEDULER

kubectl apply -f my-pv.yaml

The PV is stored in the tenant control plane
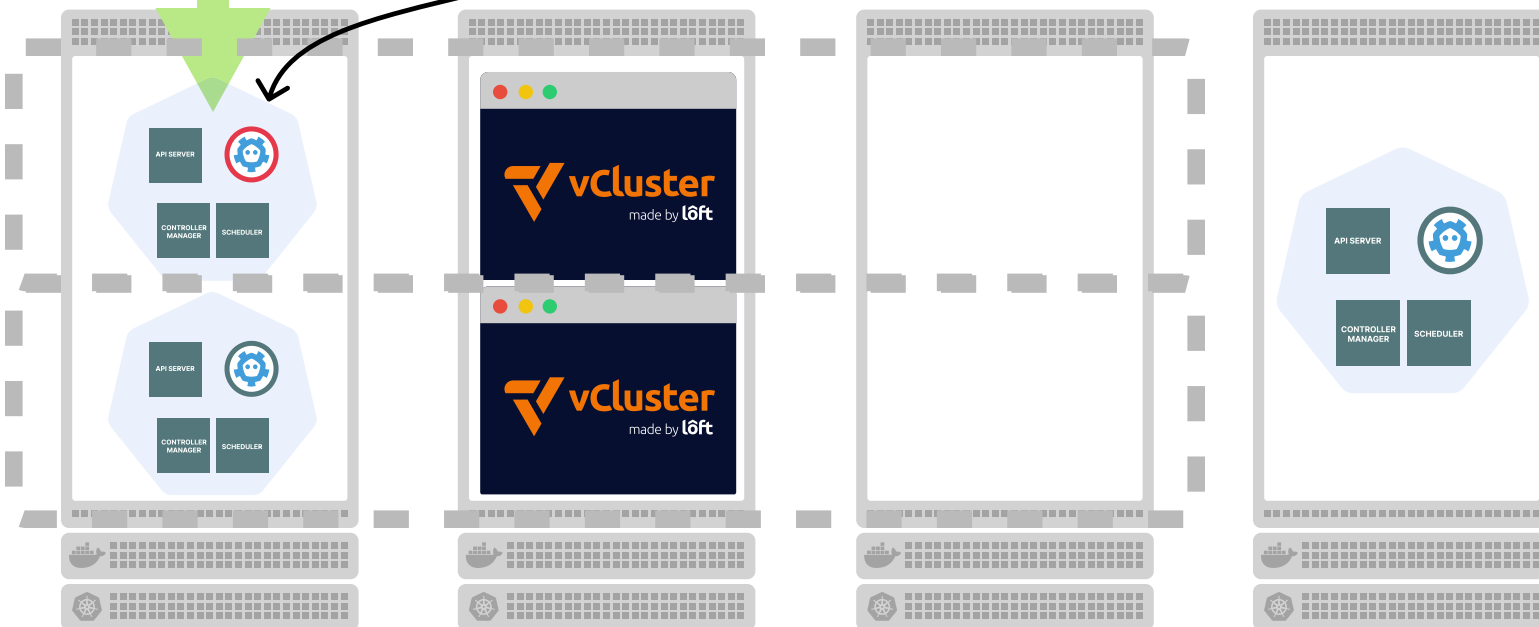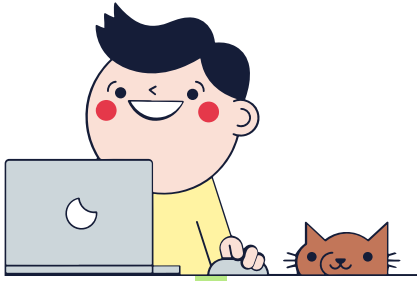
master

TEAM A NAMESPACE

TEAM B NAMESPACE

`kubectl apply -f my-pv.yaml`

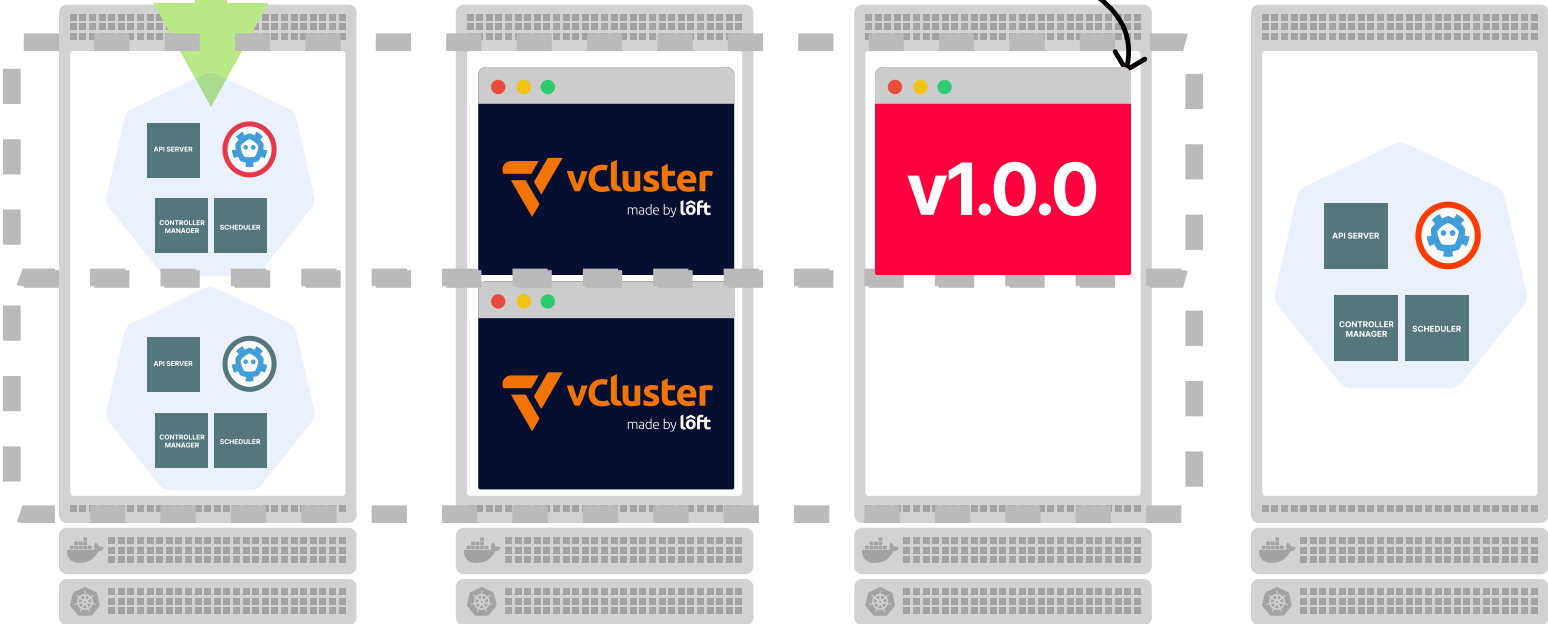the PV is synced

master

TEAM A
NAMESPACE

TEAM B
NAMESPACE

# Demo

# vCluster

# "Nested" control planes

# "Nested" control planes Admin vs tenants

# "Nested" control planes
# Admin vs tenants
# Shared host cluster

CHEAP
SHARED

EXPENSIVE
ISOLATED

MULTI-TENANCY

# + 17 nodes x $12

# + 17 nodes x $12
# + 50 PVs x $1
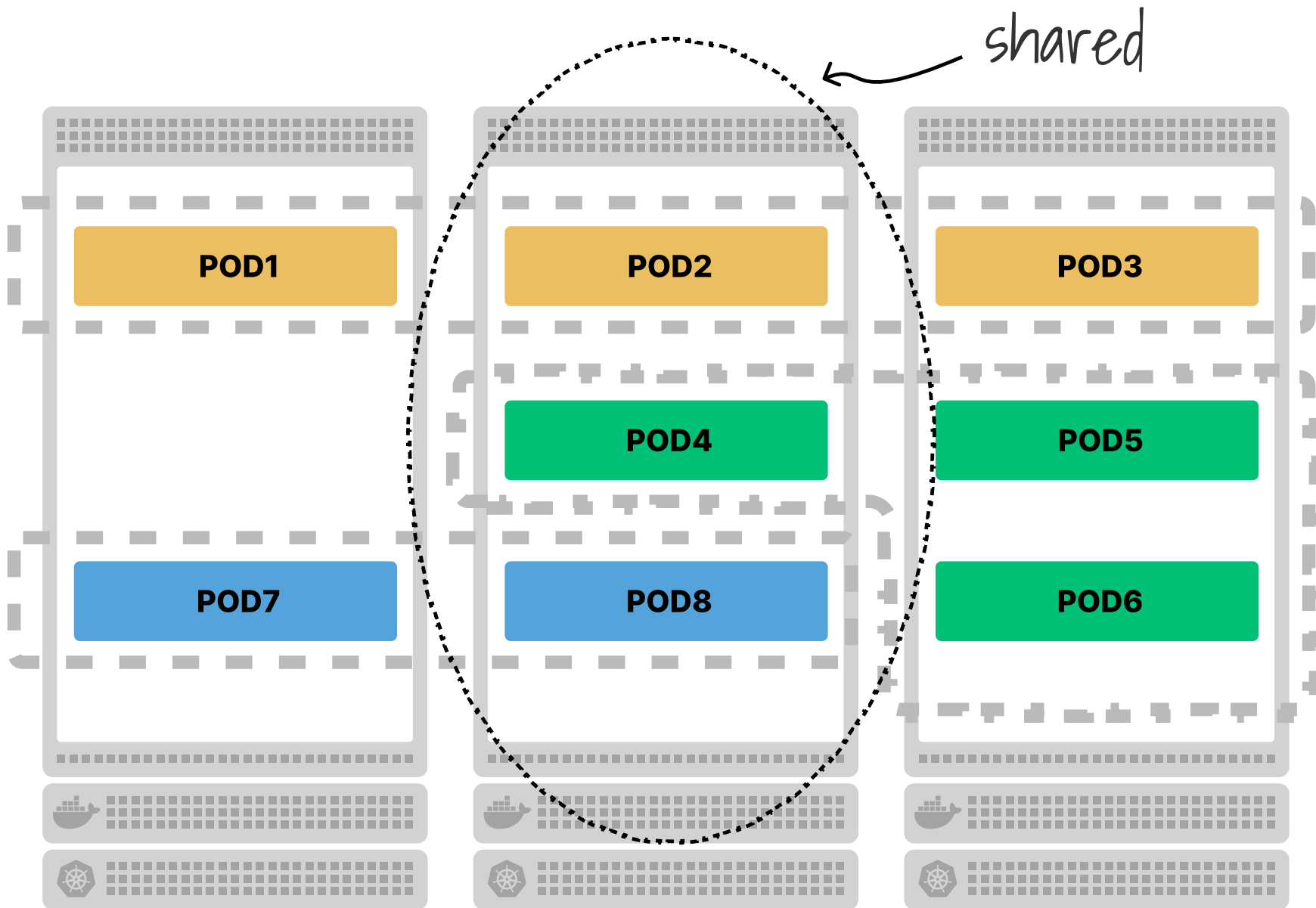
---

+ 17 nodes x $12
+ 50 PVs x $1

_____

= $254 / month

~$5 / month / tenant

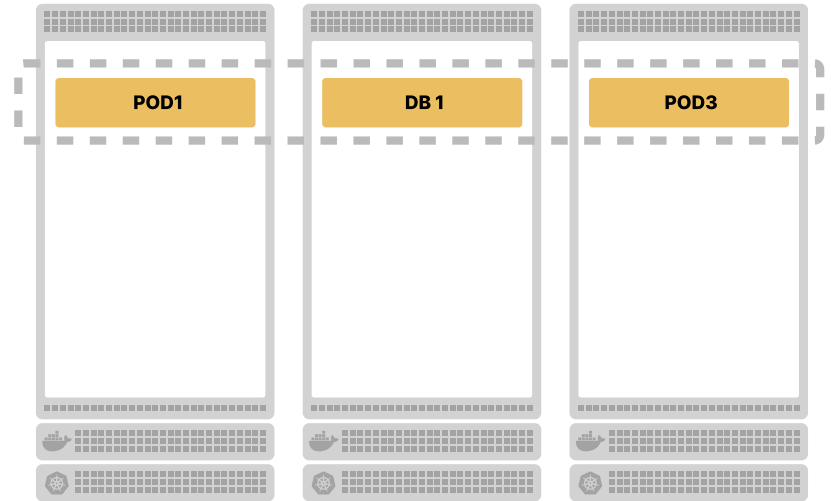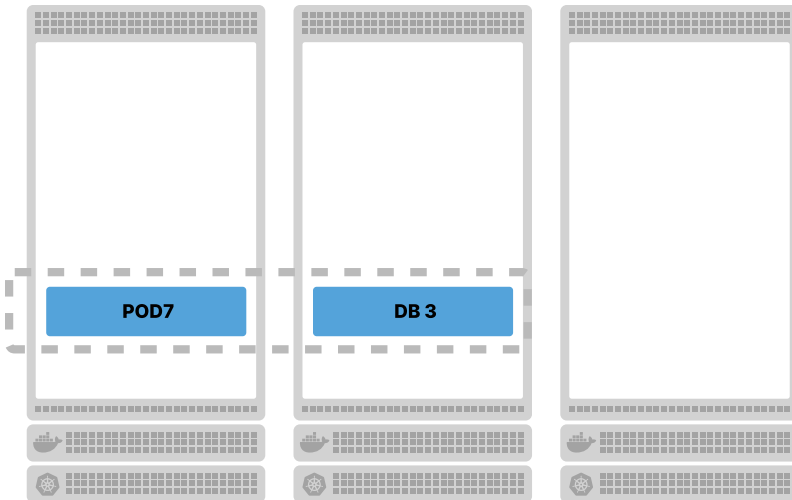# vCluster and shared nodes

noisy neighbours

# POOL 1

**DB 2**

**POD5**

**POD6**

# POOL 2

**POD1**

**DB 1**

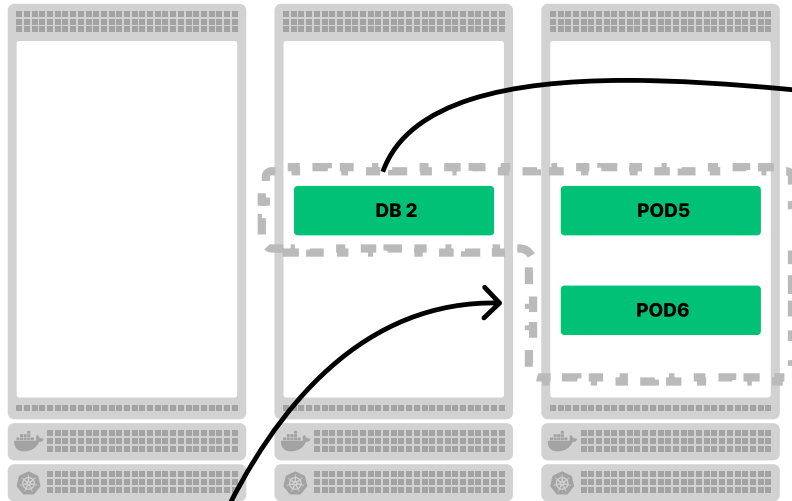**POD3**

# POOL 3
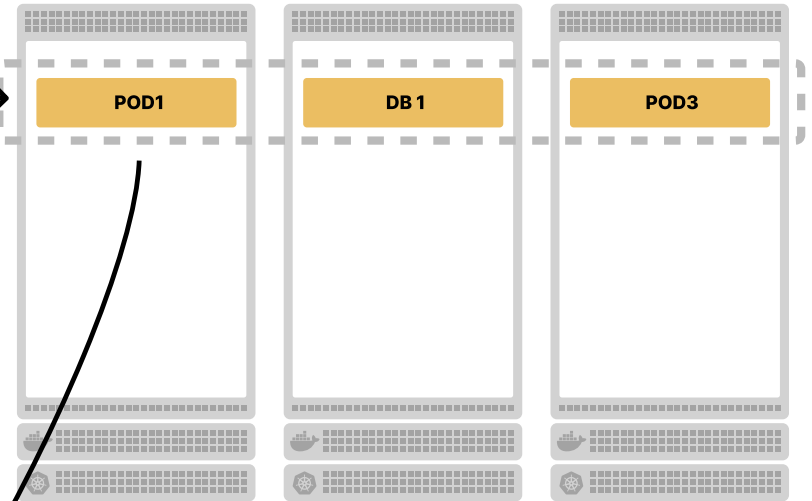
**POD7**

**DB 3**

# CONTROL PLANE

```
--node-selector
--enforce-node-selector
```

# vCluster and shared network

POOL 1

POOL 2

DB 2
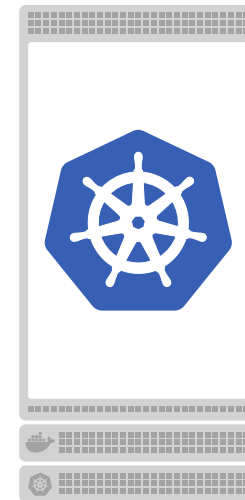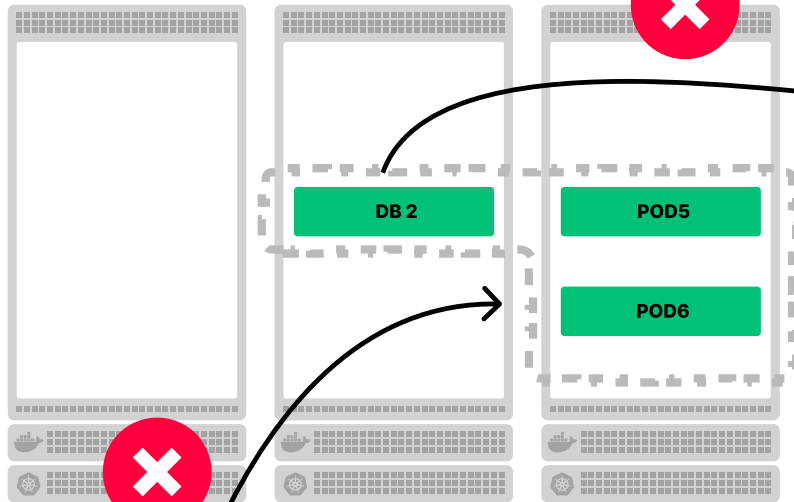
POD5

POD6

POD1

DB 1

POD3

POOL 3

CONTROL PLANE

POD7

DB 3

POOL 1

POOL 2

DB 2

POD5

POD6

POD1

DB 1

POD3

POOL 3

CONTROL PLANE

POD7

DB 3

POOL 1

DB 2
POD5
POD6

POOL 2

POD1
DB 1
POD3

POOL 3

POD7
DB 3

CONTROL PLANE

`--isolate`

# vCluster and shared cluster

POOL 1

CONTROL PLANE

container escape

DB 2

POD

POD6

kubelet take over

control plane escalation

# Dedicated clusters

**Cluster A**

**Cluster B**

**Cluster C**

Pod

Pod

Pod

# Karmada

# Karmada architecture

AGENT

API SERVER

CONTROLLER MANAGER

SCHEDULER

KARMADA API SERVER

KARMADA CONTROLLER MANAGER

KARMADA SCHEDULER

control plane

cluster1

cluster1

AGENT

Pod1

API SERVER

CONTROLLER MANAGER

SCHEDULER

KARMADA
API SERVER

KARMADA
CONTROLLER MANAGER

KARMADA
SCHEDULER

kubectl apply

cluster1

AGENT

KARMADA

Pod1 Pod2

API SERVER

CONTROLLER MANAGER

SCHEDULER

API SERVER

CONTROLLER MANAGER

SCHEDULER

kubectl apply

cluster1

# Independent cluster with central management

kubectl

manager  TEAM A

TEAM B

TEAM C

kubectl

manager   TEAM A

TEAM B

Pod

Pod

TEAM C

Pod

kubectl

manager    TEAM A

TEAM B

Pod

Pod

Pod

Pod

TEAM C

Pod

kubectl

manager TEAM A



TEAM B



TEAM C

# Demo
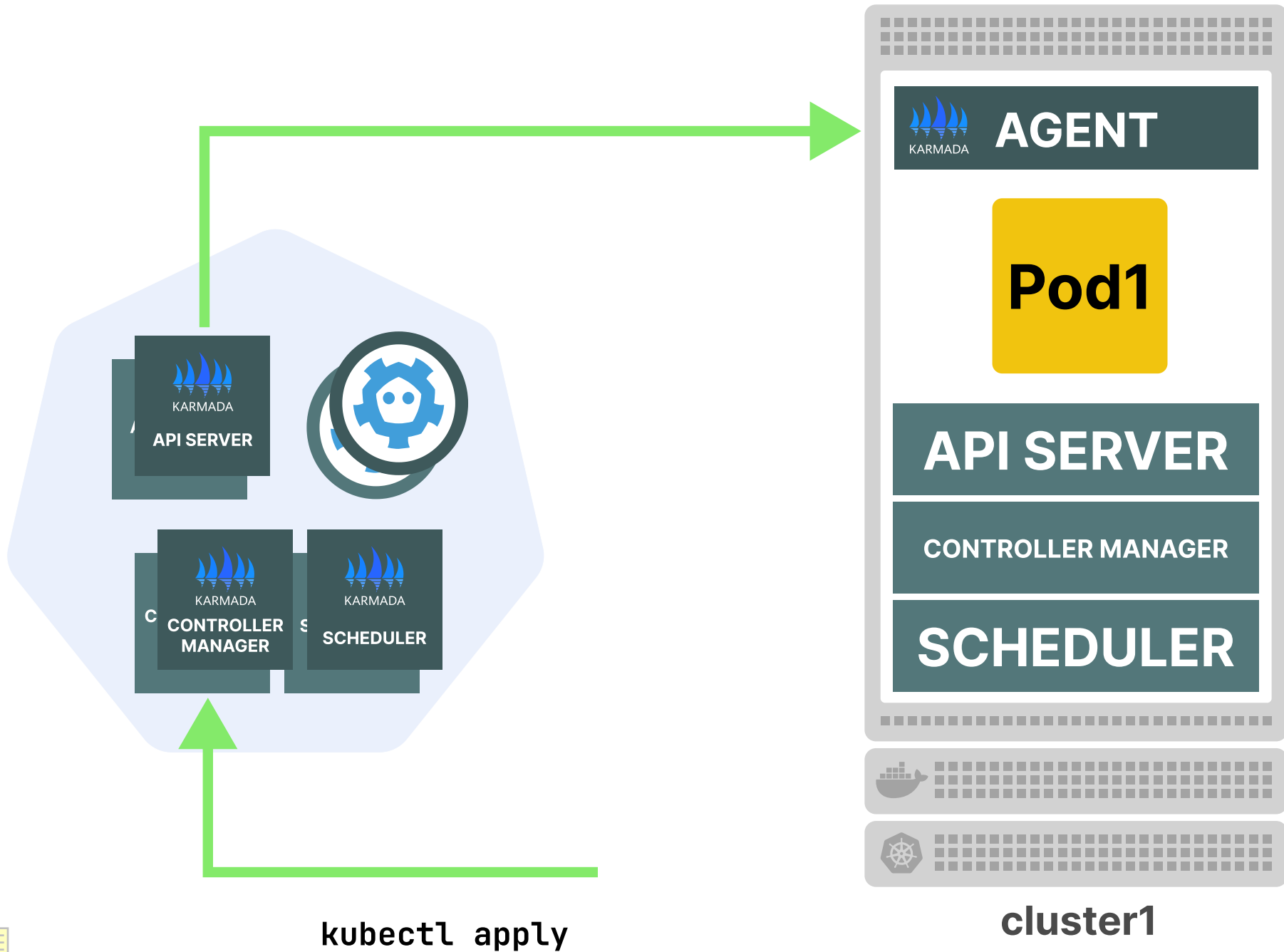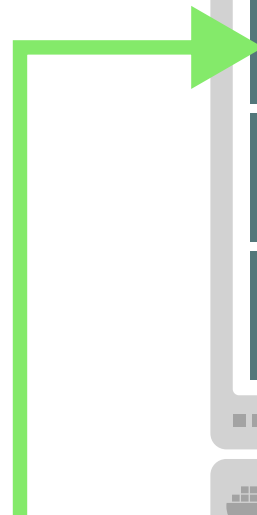
---

# Karmada

# Cluster of clusters

# Cluster of clusters Admin vs tenants

# Cluster of clusters
# Admin vs tenants
# No sharing

# + 51 clusters x $0

# + 51 clusters x $0
# + 51 nodes x $12

_____

**+ 51 clusters x $0**

**+ 51 nodes x $12**

_____

**= $612 / month**

~$12 / month / tenant

# Multi-tenancy baseline

**Multi-tenancy**

**Kubernetes**

monitoring

**Node pools, Sandbox runtime**

# Multi-tenancy

# Kubernetes

monitoring          logging

**Node pools, Sandbox runtime**

# Multi-tenancy

# Kubernetes

monitoring  logging  storage

CI/CD

**Node pools, Sandbox runtime**
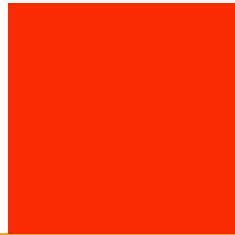
# Multi-tenancy

# Kubernetes

# Costs

# DEDICATED INGRESS FOR 50 TENANTS

**50 × 3**

**CPU**

# 5vCPU

| Instance Size | vCPU | Memory (GiB) | Instance Storage (GB) | Network Bandwidth (Gbps)*** | EBS Bandwidth (Gbps) |
|---|---|---|---|---|---|
| c6i.large | 2 | 4 | EBS-Only | Up to 12.5 | Up to 10 |
| c6i.xlarge | 4 | 8 | EBS-Only | Up to 12.5 | Up to 10 |
| c6i.2xlarge | 8 | 16 | EBS-Only | Up to 12.5 | Up to 10 |
| c6i.4xlarge | 16 | 32 | EBS-Only | Up to 12.5 | Up to 10 |
| c6i.8xlarge | 32 | 64 | EBS-Only | 12.5 | 10 |

**MEMORY**

# 4.5GB

**$0.34/hr**    **$248.2/m**

# Costs*

CHEAP
SHARED

EXPENSIVE
ISOLATED

*MULTI-TENANCY*

# Multi-tenant platform from scratch

## Recap

# Recap

## 1. Isolation VS costs

2. Multi-tenant components (e.g. Ingress)

3. Constant vs linear vs exponential costs

4. HNC and vCluster

5. Karmada

# Recap

1. Isolation VS costs

2. **Multi-tenant components (e.g. Ingress)**

3. Constant vs linear vs exponential costs

4. HNC and vCluster

5. Karmada

# Recap

1. Isolation VS costs

2. Multi-tenant components (e.g. Ingress)

**3. Constant vs linear vs exponential costs**

4. HNC and vCluster

5. Karmada

# Recap

1. Isolation VS costs

2. Multi-tenant components (e.g. Ingress)

3. Constant vs linear vs exponential costs

## 4. HNC and vCluster

5. Karmada

# Recap

1. Isolation VS costs

2. Multi-tenant components (e.g. Ingress)

3. Constant vs linear vs exponential costs

4. HNC and vCluster

5. Karmada

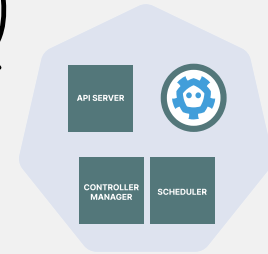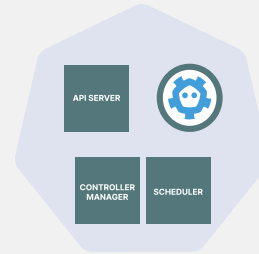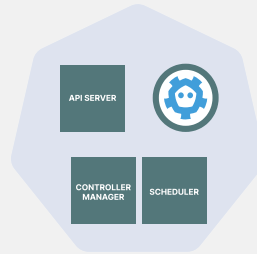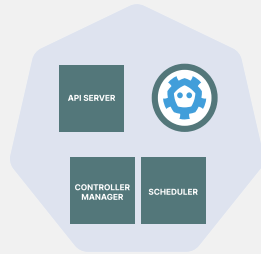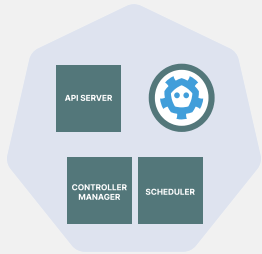# loft

# Thank you!

# Thank you!

---

**in** Chris Nesbitt-Smith

# Hypershift/Kamaji

CLUSTER

POD1    POD2    POD3    POD4    POD5

control plane as a pod

CLUSTER

control plane as a pod

API SERVER

CONTROLLER MANAGER | SCHEDULER

API SERVER

CONTROLLER MANAGER | SCHEDULER

API SERVER

CONTROLLER MANAGER | SCHEDULER

API SERVER

CONTROLLER MANAGER | SCHEDULER

API SERVER

CONTROLLER MANAGER | SCHEDULER

POD1  POD2  POD3  POD4  POD5

your own nodes