# Balancing Isolation and Complexity in a Kubernetes Platform

Chris Nesbitt-Smith

# Chris Nesbitt-Smith

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

learnk8s

# Case study

GitPod

kubernetes / kubernetes

Type / to search

<> Code | ⊙ Issues `1.9k` | ⫚ Pull requests `643` | ⊙ Actions | ⊞ Projects `1` | ⊙ Security | ⌁ Insights

❄ **kubernetes** `Public`

Watch `3216` ⌄ | ⑂ Fork `40.4k` ⌄ | ☆ Star `114k` ⌄

⑂ master ⌄ | ⑂ 56 Branches | ⬚ 1153 Tags

Go to file | t | Add file ⌄ | <> Code ⌄

🤖 **k8s-ci-robot** Merge pull request #130860 from carlory/remove-storage-fea... ••• 6396fa0 · 17 minutes ago ⏱ **128,996 Commits**

| 📁 .github | Add new contribex leads to sig-contribex-approvers | 2 years ago |
|---|---|---|
| 📁 CHANGELOG | Update release notes in changelog-1.30 to fix example cl... | 3 hours ago |
| 📁 LICENSES | Vendor randfill | last week |
| 📁 api | Merge pull request #128622 from jpbetz/admission-polic... | 17 minutes ago |
| 📁 build | Fix KUBE_BUILD_IMAGE_CROSS_TAG mismatch when K... | 9 hours ago |
| 📁 cluster | Fix typo and pass the environment variable required to en... | 5 days ago |
| 📁 cmd | Merge pull request #130354 from siyuanfoundation/forwa... | 4 days ago |
| 📁 docs | Make root approval non-recursive | 3 years ago |
| 📁 hack | Merge pull request #130821 from BenTheElder/revert-procs | 3 days ago |
| 📁 logo | logo: better alignment of layers | 3 years ago |
| 📁 pkg | Merge pull request #128622 from jpbetz/admission-polic... | 17 minutes ago |
| 📁 plugin | address comment | 3 days ago |
| 📁 staging | Merge pull request #128622 from jpbetz/admission-polic... | 17 minutes ago |
| 📁 test | Merge pull request #130860 from carlory/remove-storage... | 17 minutes ago |
| 📁 third_party | Revert "tests: include stdout of failed commands in JUnit" | 2 months ago |
| 📁 vendor | Vendor randfill | last week |
| 📄 .generated_files | remove clearly unnecessary lingering BUILD file references | 3 years ago |

**About**

Production-Grade Container Scheduling and Management

🔗 **kubernetes.io**

`go` `kubernetes` `containers` `cncf`

📖 Readme

⚖ Apache-2.0 license

🛡 Code of conduct

🛡 Security policy

⌁ Activity

⊞ Custom properties

☆ 114k stars

👁 3.2k watching

⑂ 40.4k forks

Report repository

**Releases** 728

🏷 **Kubernetes v1.32.3** `Latest`
last week

+ 727 releases

**Packages**

No packages published

**Contributors** 3,779

github.com/kubernetes/kubernetes

kubernetes / kubernetes

Type / to search

Code  Issues 1.9k  Pull requests 643  Actions  Projects 1  Security  Insights

kubernetes  Public

Watch 3216  Fork 40.4k  Star 114k

master  56 Branches  1153 Tags  Go to file  Add file  Code

k8s-ci-robot  Merge pull request #130860 from carlory/remove-storage-fea...  6396fa0 · 17 minutes ago  128,996 Commits

.github  Add new contribex leads to sig-contribex-approvers  2 years ago
CHANGELOG  Update release notes in changelog-1.30 to fix example cl...  3 hours ago
LICENSES  Vendor randfill  last week
cluster  Fix typo and pass the environment variable required to en...  5 days ago
cmd  Merge pull request #130354 from siyuanfoundation/forwa...  4 days ago
docs  Make root approval non-recursive  3 years ago
hack  Merge pull request #130821 from BenTheElder/revert-procs  3 days ago
logo  logo: better alignment of layers  3 years ago
pkg  Merge pull request #128622 from jpbetz/admission-polic...  17 minutes ago
plugin  address comment  3 days ago
staging  Merge pull request #128622 from jpbetz/admission-polic...  17 minutes ago
test  Merge pull request #130860 from carlory/remove-storage...  17 minutes ago
third_party  Revert "tests: include stdout of failed commands in JUnit"  2 months ago
vendor  Vendor randfill  last week
.generated_files  remove clearly unnecessary lingering BUILD file references  3 years ago

**prefix with https://gitpod.io/#**

About

Production-Grade Container Scheduling and Management

kubernetes.io

go  kubernetes  containers  cncf

Readme
Apache-2.0 license
Code of conduct
Security policy
Activity
Custom properties
114k stars
3.2k watching
40.4k forks

Report repository

Releases 728

Kubernetes v1.32.3  Latest
last week

+ 727 releases

Packages

No packages published

Contributors 3,779

```
  1    # Kubernetes (K8s)
  2
  3    [![CII Best Practices](https://bestpractices.coreinfrastructure.org/projects/569/badge)](https://
       bestpractices.coreinfrastructure.org/projects/569) [![Go Report Card](https://goreportcard.com/badge/
       github.com/kubernetes/kubernetes)](https://goreportcard.com/report/github.com/kubernetes/kubernetes) !
       [GitHub release (latest SemVer)](https://img.shields.io/github/v/release/kubernetes/kubernetes?
       sort=semver)
  4
  5    <img src="https://github.com/kubernetes/kubernetes/raw/master/logo/logo.png" width="100">
  6
  7    ----
  8
  9    Kubernetes, also known as K8s, is an open source system for managing [containerized applications]
 10    across multiple hosts. It provides basic mechanisms for the deployment, maintenance,
 11    and scaling of applications.
 12
 13    Kubernetes builds upon a decade and a half of experience at Google running
 14    production workloads at scale using a system called [Borg],
 15    combined with best-of-breed ideas and practices from the community.
 16
 17    Kubernetes is hosted by the Cloud Native Computing Foundation ([CNCF]).
 18    If your company wants to help shape the evolution of
 19    technologies that are container-packaged, dynamically scheduled,
 20    and microservices-oriented, consider joining the CNCF.
 21    For details about who's involved and how Kubernetes plays a role,
 22    read the CNCF [announcement].
 23
 24    ----
 25
 26    ## To start using K8s
 27
 28    See our documentation on [kubernetes.io].
 29
 30    Take a free course on [Scalable Microservices with Kubernetes].
 31
```

```
 HISTFILE=/workspace/.gitpod/cmd-0 history -r; {
go get && go build ./... && go test ./... && make
} && {
go run .
}
gitpod /workspace/kubernetes (master) $  HISTFILE=/workspace/.gitpod/cmd-0 history -r; {
> go get && go build ./... && go test ./... && make
> } && {
> go run .
> }
go: downloading go1.24.0 (linux/amd64)
go: no package to get in current directory
gitpod /workspace/kubernetes (master) $ █
```

**GitPod**

# Open source

VS Code in the cloud

Dedicated environment

**GitPod**

Open source

**VS Code in the cloud**

Dedicated environment

# GitPod

~~Open source~~
~~VS Code in the cloud~~
**Dedicated environment**

# GitPod in Kubernetes

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

Incoming traffic

pod

Visual Studio Code

docker

PostgreSQL

container with multiple processes

Incoming traffic

microservice 4

**Pod 1**   **Pod 2**   **Pod 3**

microservice 2

microservice 3

microservice 1

single process per container

Incoming traffic

Pod 1

Pod 2

Pod 3

GitPod

GitPod

tenant 2

Pod 1    Pod 2    Pod 3

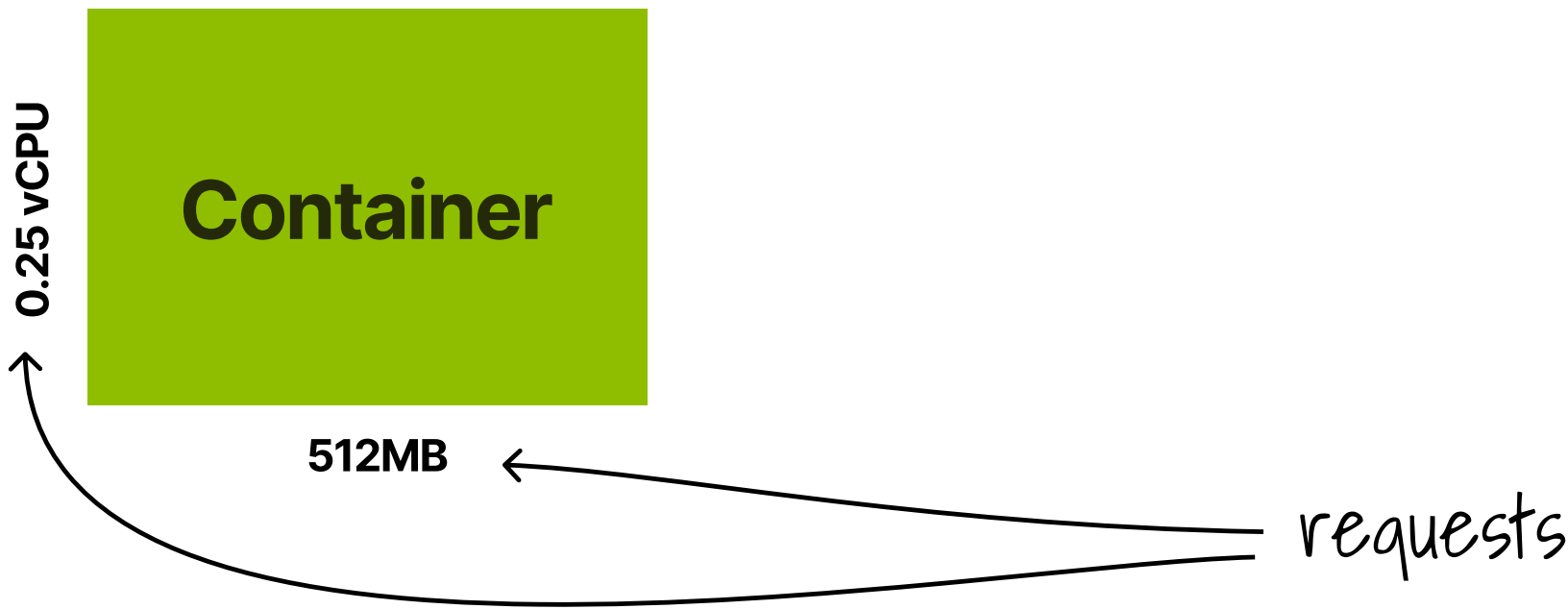tenant 1

Pod 1    Pod 2    Pod 3

tenant 2

# Case study

GitPod

**31 Oct 2024**

**31 Oct 2024**

**31 Oct 2024**

# Resource management

---

## Challenge 1

actual usage

request

limit

actual usage

# The challenge

terrible dx

Memory

Actual memory usage

4GB

Request

0                                    Time

**Slow intellisense
Slow builds**

# CFS hints

**Dynamic resource allocation**

**Swap-space memory**

CFS hints

# Dynamic resource allocation

Swap-space memory

CFS hints

Dynamic resource allocation

# Swap-space memory

# Your options

terrible dx

pay more

**Memory**

Actual memory usage

4GB ━━━━━━━━━━━━━━ Request

0          Time

**Slow intellisense**
**Slow builds**

**Memory**

4GB ━━━━━━━━━━━━━━ Request

Actual memory usage

0          Time

**wasted resourced**

p99 latency
(e.g. in API calls)

frequent scaling
up/down

terrible dx

Pay more

trade-off
between costs
an reliability

**Memory**

Actual memory usage
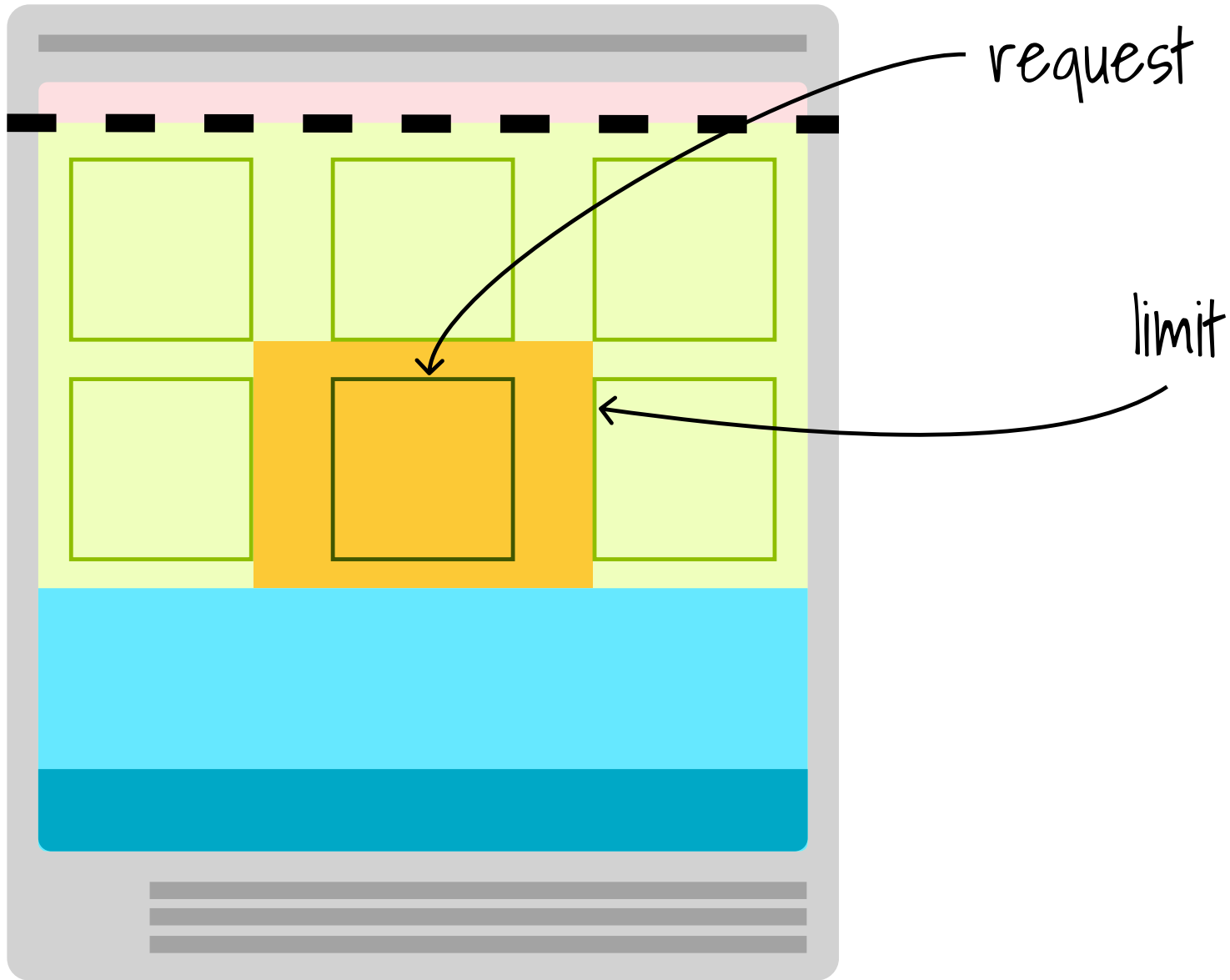
4GB ——— Request

0       **Time**

**Slow intellisense
Slow builds**

**Memory**

Request

4GB

Actual memory usage

0       **Time**

**wasted resourced**

p99 latency
(e.g. in API calls)

frequent scaling
up/down

actual usage

request

**2**

dedicated node

Expensive

resources
oversubscribe

Fully
isolated

Noisy
neighbours

Cheap

dedicated node

# Stormforge Optimize Live



**Optimization Score** ⓘ

**20%**

CPU: 2%
MEM: 79%

⌄ More Information

**Recommendation Summary**

Current Requests → Recommended Requests

$41,161 /mo → $16,987 /mo

4,571.0 Cores → 1,766.8 Cores

1.46 TiB → 1.22 TiB

**Recommendation Impact**

Estimated Savings

$24,174 /mo    59% ⌄

-2804.2 Cores    61% ⌄

-238.25 GiB    16% ⌄

🗑 **Waste**

4,234 Workloads are Overprovisioned

Recommended Adjustments    Estimated Cost Impact

−3,141.4 Cores    $27,640
−609.12 GiB    per month

⊘ **Performance Risk**

391 Workloads are Underprovisioned

Recommended Adjustments    Estimated Cost Impact

+337.2 Cores    $3,466
+371.90 GiB    per month

# Stormforge Optimize Live



## PerfectScale

# Stormforge Optimize Live



Kubex

User-Definable Views and Filters

# Network complexity

---

## Challenge 2

isolated network

Worker Node

Worker Node

GitPod

GitPod

GitPod

kubernetes

PORT 5432     PORT 5005     PORT 8080

# GitPod

SERVICE

PORT 5432    PORT 5005    PORT 8080

GitPod

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

GitPod

GitPod

GitPod

kubernetes

doesn't scale well

# Cost

N N N N N N N N N N N N N N N N N N N N N N N N N N N
N N N N N N N N N N N N N N N N N N N N N N N N N N N
N N N N N N N N N N N N N N N N N N N N N N N N N N N
N N N N N N N N N N N N N N N N N N N N N N N N N N N
N N N N N N N N N N N N N N N N N N N N N N N N N N N
N N N N N N N N N N N N N N N N N N N N N N N N N N N

**tenants**

# Worker Node

# Worker Node

single ingress

kubernetes

# **Too many services**

enableServiceLinks

Failing DNS

Too many services

# enableServiceLinks

Failing DNS

# Network complexity

Too many services
enableServiceLinks
**Failing DNS**

# How Services work

state in etcd



| Pod name | Status | Node | podIP |
|----------|--------|------|-------|
| Pod 1 | RUNNING | worker1 | 10.0.0.1 |
| Pod 2 | RUNNING | worker2 | 10.0.1.1 |

| Service name | IP address | Endpoints |
|--------------|------------|-----------|
| | | |

cluster

when you create a Service...

**SERVICE**
**ClusterIP**

| Pod name | Status | Node | podIP |
|---|---|---|---|
| Pod 1 | RUNNING | worker1 | 10.0.0.1 |
| Pod 2 | RUNNING | worker2 | 10.0.1.1 |

| Service name | IP address | Endpoints |
|---|---|---|

...the endpoint
controller collects
the endpoints

**SERVICE**
**ClusterIP**

| Pod name | Status | Node | podIP |
|----------|--------|------|-------|
| Pod 1 | RUNNING | worker1 | 10.0.0.1 |
| Pod 2 | RUNNING | worker2 | 10.0.1.1 |

| Service name | IP address | Endpoints |
|--------------|-----------|-----------|
| Red | 172.17.0.1 | 10.0.0.1:3000,10.0.1.1:3000 |

kubeproxy

| IP address to intercept | Replace with |
|---|---|

| Pod name | Status | Node | podIP |
|---|---|---|---|
| Pod 1 | RUNNING | worker1 | 10.0.0.1 |
| Pod 2 | RUNNING | worker2 | 10.0.1.1 |

| Service name | IP address | Endpoints |
|---|---|---|
| Red | 172.17.0.1 | 10.0.0.1:3000,10.0.1.1:3000 |

# kubeproxy is notified of the endpoints

| IP address to intercept | Replace with |
|---|---|
| 172.17.0.1 | 10.0.0.1, 10.0.1.1 |

| Pod name | Status | Node | podIP |
|---|---|---|---|
| Pod 1 | RUNNING | worker1 | 10.0.0.1 |
| Pod 2 | RUNNING | worker2 | 10.0.1.1 |

| Service name | IP address | Endpoints |
|---|---|---|
| Red | 172.17.0.1 | 10.0.0.1:3000,10.0.1.1:3000 |

## kubeproxy updates the iptables rules

| IP address to intercept | Replace with |
|---|---|
| 172.17.0.1 | 10.0.0.1, 10.0.1.1 |

| Pod name | Status | Node | podIP |
|---|---|---|---|
| Pod 1 | RUNNING | worker1 | 10.0.0.1 |
| Pod 2 | RUNNING | worker2 | 10.0.1.1 |

| Service name | IP address | Endpoints |
|---|---|---|
| Red | 172.17.0.1 | 10.0.0.1:3000,10.0.1.1:3000 |

| IP address to intercept | Replace with |
|---|---|
| 172.17.0.1 | 10.0.0.1, 10.0.1.1 |

| Service name | IP addresses |
|---|---|
| | |

CoreDNS is also notified!

| Pod name | Status | Node | podIP |
|---|---|---|---|
| Pod 1 | RUNNING | worker1 | 10.0.0.1 |
| Pod 2 | RUNNING | worker2 | 10.0.1.1 |

| Service name | IP address | Endpoints |
|---|---|---|
| Red | 172.17.0.1 | 10.0.0.1:3000,10.0.1.1:3000 |

| IP address to intercept | Replace with |
|---|---|
| 172.17.0.1 | 10.0.0.1, 10.0.1.1 |

| Service name | IP addresses |
|---|---|
| red.namespace.svc.cluster.local | 172.17.0.1 |

a new entry is
added to the DNS

| Pod name | Status | Node | podIP |
|---|---|---|---|
| Pod 1 | RUNNING | worker1 | 10.0.0.1 |
| Pod 2 | RUNNING | worker2 | 10.0.1.1 |

| Service name | IP address | Endpoints |
|---|---|---|
| Red | 172.17.0.1 | 10.0.0.1:3000,10.0.1.1:3000 |

# Services at scale

## Services at scale

# if you have 5,000 services in your Kubernetes cluster, it takes 11 minutes to add a new rule with iptables

```
root@nginx:/# env
KUBERNETES_SERVICE_PORT_HTTPS=443
KUBERNETES_SERVICE_PORT=443
HOSTNAME=nginx
PWD=/
PKG_RELEASE=1~bookworm
HOME=/root
KUBERNETES_PORT_443_TCP=tcp://10.96.0.1:443
DYNPKG_RELEASE=1~bookworm
NJS_VERSION=0.8.9
TERM=xterm
SHLVL=1
KUBERNETES_PORT_443_TCP_PROTO=tcp
KUBERNETES_PORT_443_TCP_ADDR=10.96.0.1
KUBERNETES_SERVICE_HOST=10.96.0.1
KUBERNETES_PORT=tcp://10.96.0.1:443
KUBERNETES_PORT_443_TCP_PORT=443
PATH=/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin
NJS_RELEASE=1~bookworm
```

2 for each service

# CoreDNS at scale

overload

# Your options

Sharing resources

# Your options

Network isolation

default
NAMESPACE

Pod 1

Pod 3

Pod 2

nginx-ingress
NAMESPACE

# Example

**Pod 1**

LABELS

**app: auth**
**role: api**

**Pod 2**

LABELS

**app: backend**
**role: api**

**Pod 3**

LABELS

**app: web**
**role: frontend**

```
~$ cat pod-policy.yaml
kind: NetworkPolicy
apiVersion: networking.k8s.io/v1
metadata:
  name: api-allow
spec:
  podSelector:
    matchLabels:
      app: backend
      role: api
  ingress:
  - from:
      - podSelector:
          matchLabels:
            app: web
```

# Your options

---

Services at scale

```yaml
apiVersion: v1
kind: Pod
metadata:
  name: nginx
spec:
  enableServiceLinks: false
  containers:
  - name: nginx
    image: nginx:1.14.2
    ports:
    - containerPort: 80
```

no more env variables

```
root@nginx:/# env
KUBERNETES_SERVICE_PORT_HTTPS=443
KUBERNETES_SERVICE_PORT=443
HOSTNAME=nginx
PWD=/
PKG_RELEASE=1~boo
HOME=/root
KUBERNETES_PORT_443_TCP=//10.96.0.1:443
DYNPKG_RELEASE=1~bookworm
NJS_VERSION=0.8.9
TERM=xterm
SHLVL=1
KUBERNETES_PORT_443_TCP_PROTO=tcp
KUBERNETES_PORT_443_TCP_ADDR=10.96
KUBERNETES_SERVICE_HOST=10.96.0.1
KUBERNETES_PORT=//10.96.0.1:443
KUBERNETES_PORT_443_TCP_PORT=443
PATH=/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin
NJS_RELEASE=bookworm
```
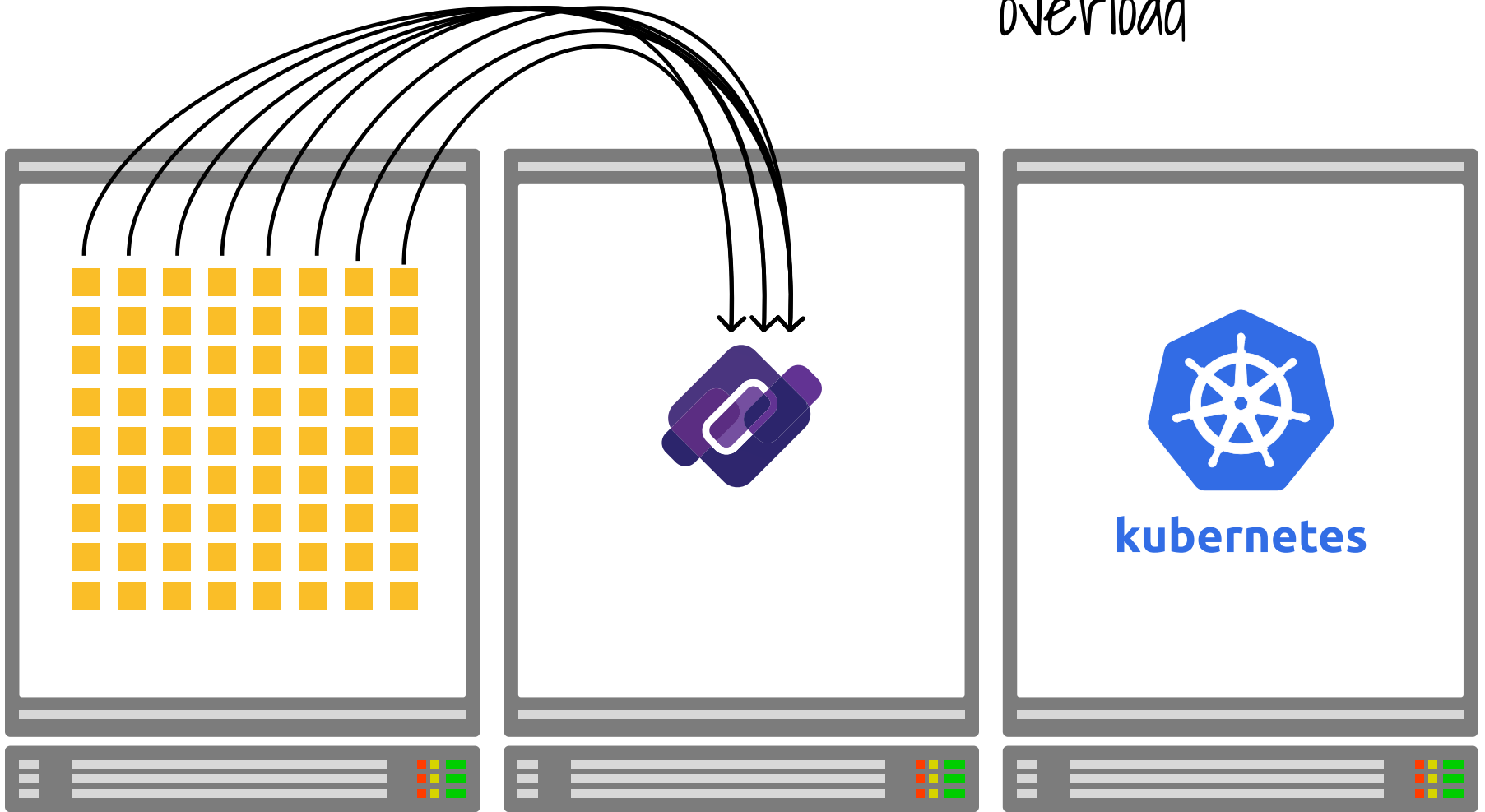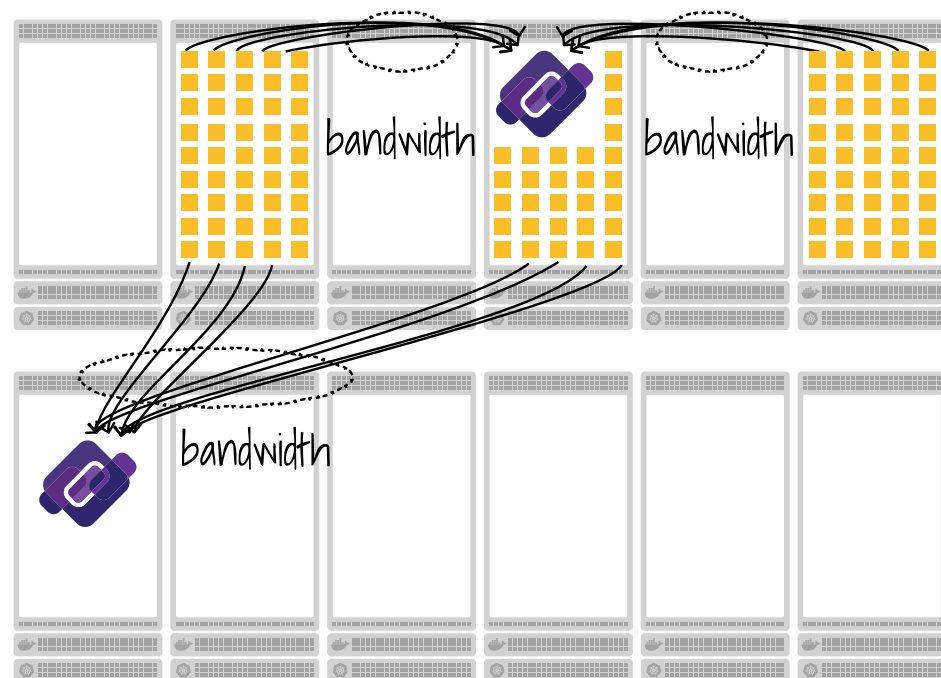
2 for each service

# Your options

iptables limits
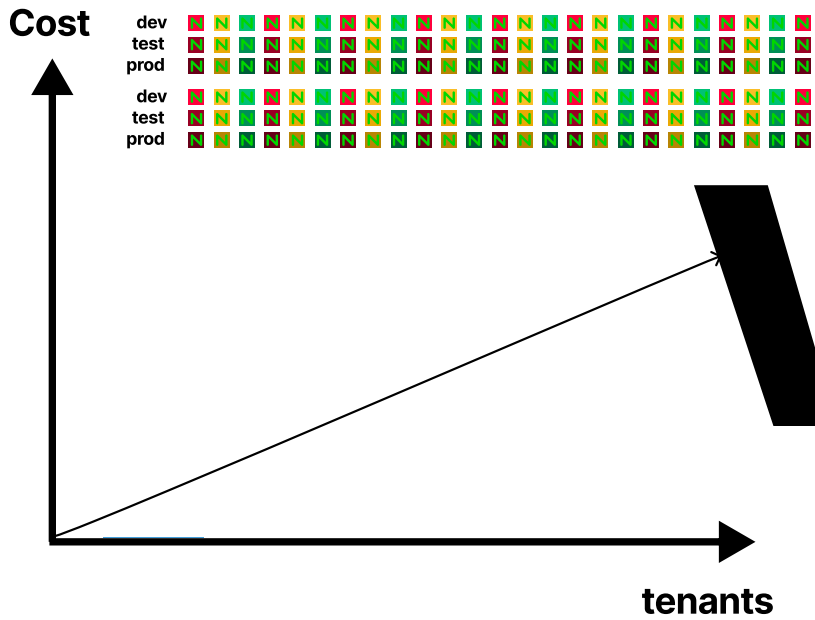
**FROM:**
**Pod** ▉ (10.244.1.7)

**TO:**
**Service** ▮ (10.96.5.81)
**Pod** ▮ (10.244.1.3)

the traffic is rewritten (DNAT)

**eBPF**

10.96.5.81

**10.244.1.7**

**10.244.1.2**

**10.244.1.3**

# Your options

DNS scaling

DNS query

**app**

**CoreDNS**

daemonsets

NodeLocal
DNSCache

NodeLocal
DNSCache

**app**

**CoreDNS**

DNS query

NodeLocal
DNSCache

NodeLocal
DNSCache

**app**

**CoreDNS**

cache hit

NodeLocal
DNSCache

NodeLocal
DNSCache

**app**

**CoreDNS**

# cache miss

NodeLocal DNSCache

NodeLocal DNSCache

**app**

**CoreDNS**

# cache miss



NodeLocal
DNSCache

NodeLocal
DNSCache

**app**

**CoreDNS**

# **Your options**

---

DNS abuse

Tenant 3

Tenant 2

CoreDNS

Tenant 1

Tenant 4

DNS queries

```
cluster.local {
    metadata
    kubernetes {
        pods verified
    }
    firewall query {
        allow [kubernetes/client-namespace] !~ '^tenant-'
        allow [kubernetes/namespace] == [kubernetes/client-namespace]
        allow [kubernetes/namespace] == 'default'
        block true
    }
}
```

# coredns/policy

# Your options

Service meshes

**Ingress Pod**

NGINX

Pod 1    Pod 2    Pod 3

*nested layers*

Ingress Pod

Control plane

Pod 1

Pod 2

Pod 3

Ingress Pod

Control plane

Pod 1

Pod 2

Pod 3

Ingress Pod

Control plane

Pod 1

Pod 2

Pod 3

GET /API/INFO

POST /API/DATA

# Workload isolation

Challenge 3

isolated workspace

Worker Node   Worker Node

GitPod
GitPod
GitPod

kubernetes

isolated workspace

Worker Node

Worker Node

GitPod

GitPod

GitPod

kubernetes

control plane take over

## Workload isolation

# User namespaces*

# Sandboxed runtimes

# Workload isolation

## User namespaces*

## Sandboxed runtimes

# User namespaces

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

GitPod
user: root

kubernetes

the process in the container
runs as root

Worker Node

Worker Node

GitPod
user: root

user: non-root

kubernetes

the process is mapped to a
non-root user in the host

# Sandbox runtimes

CRI

CNI

CSI

kubelet

kubelet delegates to

CRI
CNI
CSI

container**d**

container runtime

kubelet

**Docker image**

kubelet delegates to

CRI

CNI

CSI

containerd

Download from DockerHub

container runtime

kubelet

Docker image

pod

kubelet delegates to

CRI

CNI

CSI

container**d**

Download from
DockerHub

container runtime

kubelet

kubelet delegates to

CRI

CNI

CSI

Firecracker

**Cloud Hypervisor**

kubelet

pod in a (lightweight) VM

kubelet delegates to

**CRI**

**CNI**

**CSI**

Firecracker

**Cloud Hypervisor**

kubelet

# Your options

# Workload isolation

# SecurityContext
gVisor
Sandboxed runtimes

```
~$ cat pod.yaml
apiVersion: v1
kind: Pod
metadata:
  name: security-context-pod
spec:
  securityContext:
    runAsUser: 2500
    fsGroup: 2000
  volumes:
  - name: security-context-vol
    emptyDir: {}
  containers:
  - name: security-context-cont
    image: supergiantkir/k8s-liveliness
    volumeMounts:
    - name: security-context-vol
      mountPath: /data/test
    securityContext:
      allowPrivilegeEscalation: false
```

```
~$ cat pod.yaml
apiVersion: v1
kind: Pod
metadata:
  name: linux-cpb-demo
spec:
  securityContext:
    runAsUser: 3000
  containers:
  - name: linux-cpb-cont
    image: supergiantkir/k8s-liveliness
    securityContext:
      capabilities:
        add: ["NET_ADMIN"]
```

# Workload isolation

SecurityContext

**gVisor**

Sandboxed runtimes

JVM

node js

**System call**

**Kernel**

**Files**  **Network**  **Hardware**  **CPU**

# JVM

# node.js

**System call**

## gVisor

## Kernel

| Files | Network | Hardware | CPU |

# Workload isolation

SecurityContext
gVisor
## Sandboxed runtimes

Firecracker

Cloud Hypervisor

katacontainers

EDERA

WA

varnish/tinykvm*

container escape

POOL 1                                    CONTROL PLANE

dev          prod

Salman Iqbal

19:59 / 58:34 • Environments x Tenants at Scale

# Security as a spectrum

**Less secure** → **More secure**

**coredns policies**

**requests/limits**

**network policies**

**sandboxed runtimes**

**Less secure**

**More secure**

**admission controllers**

**user namespaces**

coredns policies

requests/limits

network policies

sandboxed runtimes

Less secure

More secure

admission controllers

user namespaces

mix and match

All posts > Engineering blog

# We're leaving Kubernetes

31 Oct 2024

Christian Weichel / Co-Founder, CTO at Gitpod

Alejandro de Brito Fontes / Staff Engineer

Kubernetes seems like the obvious choice for building out remote, standardized and automated development environments. We thought so too and have spent six years invested in making the most popular cloud development environment platform at internet scale. That's 1.5 million users, where we regularly see thousands of development environments per day. In that time, we've found that Kubernetes is not the right choice for building development environments.

This is our journey of experiments, failures and dead-ends building development environments on Kubernetes. Over the years, we experimented with many ideas involving SSDs, PVCs, eBPF, seccomp notify, TC and io_uring, shiftfs, FUSE and idmapped mounts, ranging from microVMs, kubevirt to vCluster.

In pursuit of the most optimal infrastructure to balance security, performance and interoperability. All while wrestling with the unique challenges of building a system to scale up, remain secure as it's handling arbitrary code execution, and be stable enough for developers to work in.

This is not a story of whether or not to use Kubernetes for production workloads that's a whole separate

# 31 Oct 2024

# Security as a spectrum

## Understand constraints
Define goals
Tooling

# Security as a spectrum

Understand constraints

# Define goals

Tooling

# Security as a spectrum

Understand constraints

Define goals

# Tooling

# Multi-tenancy spectrum

**Namespaces**

**Namespaces as a Service**

**Kubernetes API as a Service**

**Control plane as a service (internal)**

**Control plane as a service (external)**

**Dedicated clusters**

kØSMOTRON

HNC  capsule  kcp  k3k  vCluster  Kamaji  HYPERSHIFT  KARMADA  sveltos

- Namespaces
- Namespaces a
- Kubernetes A
- Control plane
- Control plane
- Dedicated clusters

**kØSMOTRON**

HNC · capsule · kcp · k3k · vCluster · Kamaji · HYPERSHIFT · KARMADA · sveltos

*lighter options*

🟧 **Namespaces**

🟥 **Namespaces as a Service**

🟧 **Kubernetes API as a Service**

🟨 **Control plane as a service (internal)**

🟨 Control plane as a service (external)

🟩 Dedicated clusters

**CRD isolation**

kØSMOTRON

HNC    kcp    k3k    vCluster    Kamaji    HYPERSHIFT    KARMADA    sveltos

**Namespaces**

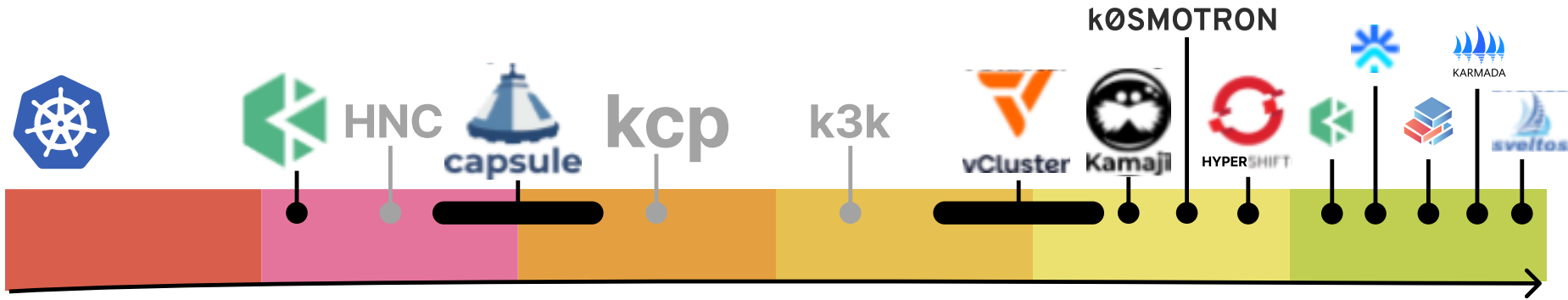**Namespaces as a Service**

**Kubernetes API as a Service**

**Control plane as a service (internal)**

CRD isolaiton

**Control plane as a service (external)**

**Dedicated clusters**

# SINGLE CLUSTER

**SHARED CONTROL PLANE**



**DEDICATED CONTROL PLANE**


vCluster

# SINGLE CLUSTER

# CLUSTER MANAGER

### SHARED CONTROL PLANE

### DEDICATED CONTROL PLANE

vCluster

kØSMOTRON

Kamaji

HYPER SHIFT

# SINGLE CLUSTER

## CLUSTER MANAGER

## DEDICATED CLUSTER

### SHARED CONTROL PLANE



### DEDICATED CONTROL PLANE


vCluster

kØSMOTRON


Kamaji


HYPER SHIFT






KARMADA


sveltos

# SINGLE CLUSTER

## CLUSTER MANAGER

## DEDICATED CLUSTER

**SHARED CONTROL PLANE**

**DEDICATED CONTROL PLANE**

kØSMOTRON

Kamaji

HYPE SHIFT

R

**ADD ONS**

**Network policies**

**Sandboxed runtime**

**CoreDNS Policies**

**Service meshes**

**Kube-proxy alternative**

**ADD ONS**

**Sandboxed runtime**

**Service meshes**

KARMADA

sveltos

**ADD ONS**

**Service meshes**

# Takeaways

Recap

# Recap

**1. Kubernetes is for sharing**

2. Sharing resources is hard

3. Sharing network is hard

4. Securing shared workload is hard

5. Security as spectrum

# Recap

1. Kubernetes is for sharing

2. Sharing resources is hard

3. Sharing network is hard

4. Securing shared workload is hard

5. Security as spectrum

# Recap

1. Kubernetes is for sharing

2. Sharing resources is hard

3. Sharing network is hard

4. Securing shared workload is hard

5. Security as spectrum

# Recap

1. Kubernetes is for sharing

2. Sharing resources is hard

3. Sharing network is hard

4. **Securing shared workload is hard**

5. Security as spectrum

# Recap

1. Kubernetes is for sharing

2.  Sharing resources is hard

3. Sharing network is hard

4. Securing shared workload is hard

5. Security as spectrum

# Thank you!

# Thank you!

---

**in** Chris Nesbitt-Smith | cns.me