



Multitenancy in Kubernetes: Better walls make better tenants

Adrian Ludwin, Senior Software Engineer
Wednesday, June 17, 2020
aludwin@google.com

Google Cloud



Goal of this session

Understand Google's opinionated **best practices** for enterprise multitenancy on Kubernetes, including when and how to apply them.



Topics

- 1 Why use multitenancy?
- 2 What are the principles of multitenancy?
- 3 How do I implement multitenancy?
- 4 Advanced topics
- 5 Conclusion

Why use multitenancy?

What companies care about

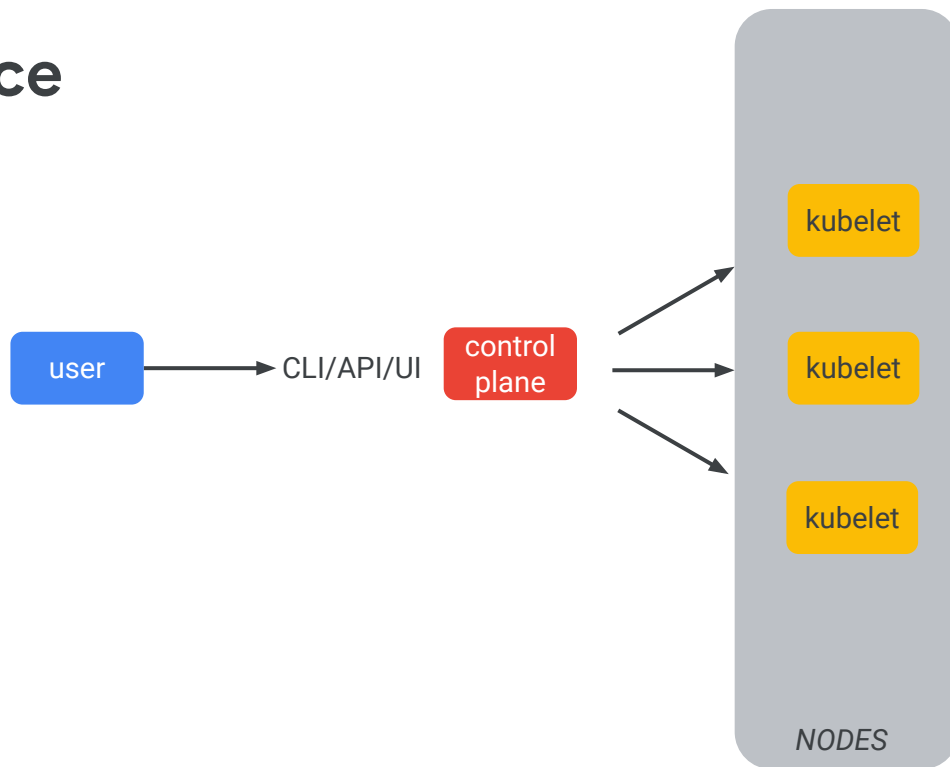


Cost

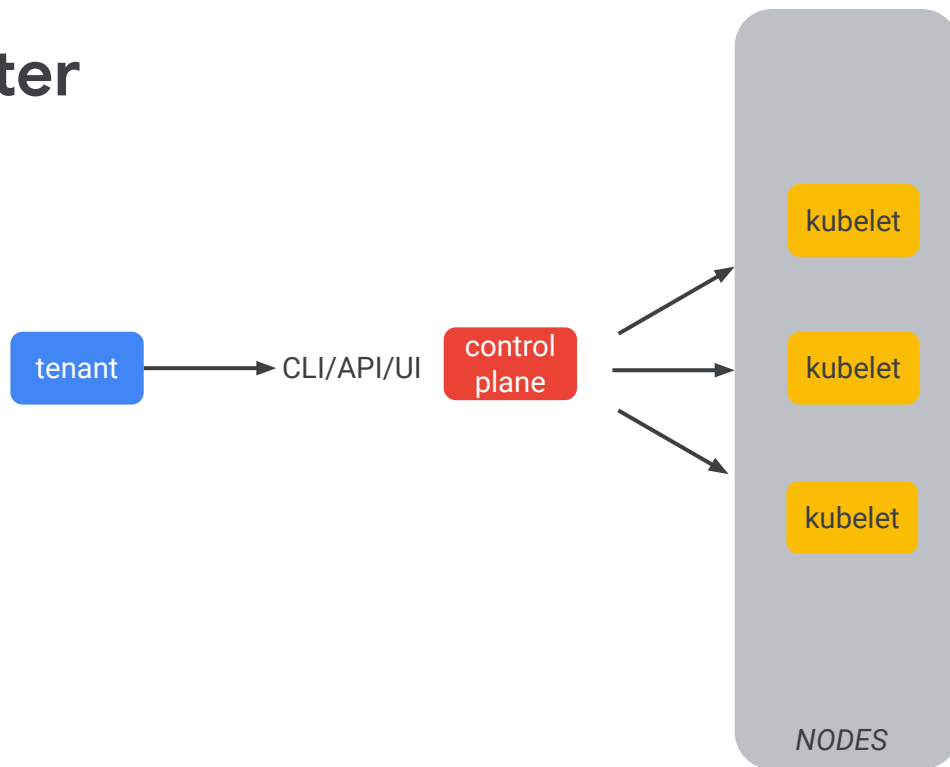


Velocity

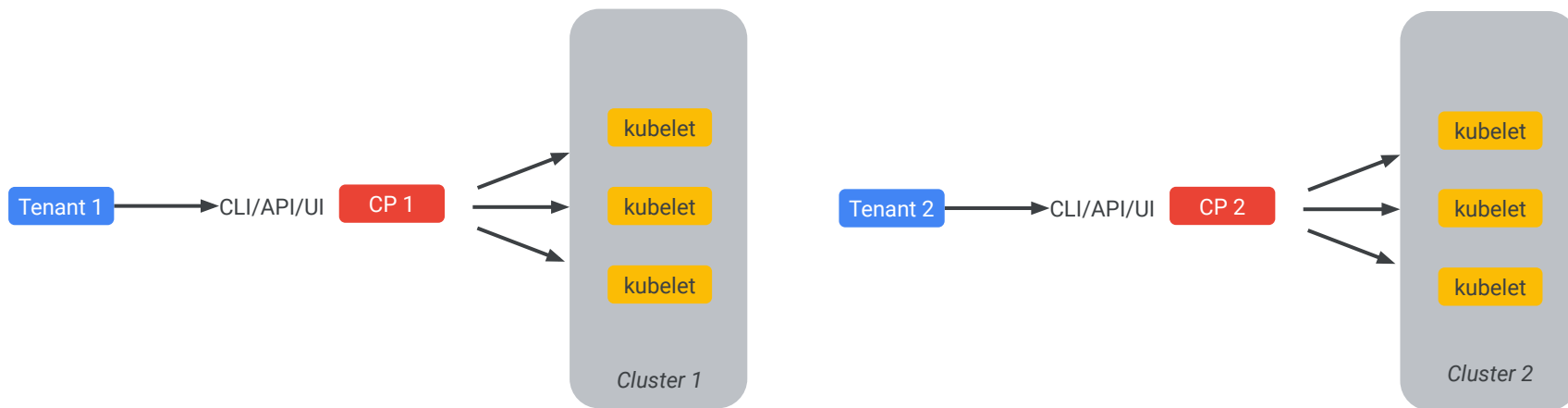
Kubernetes at a glance



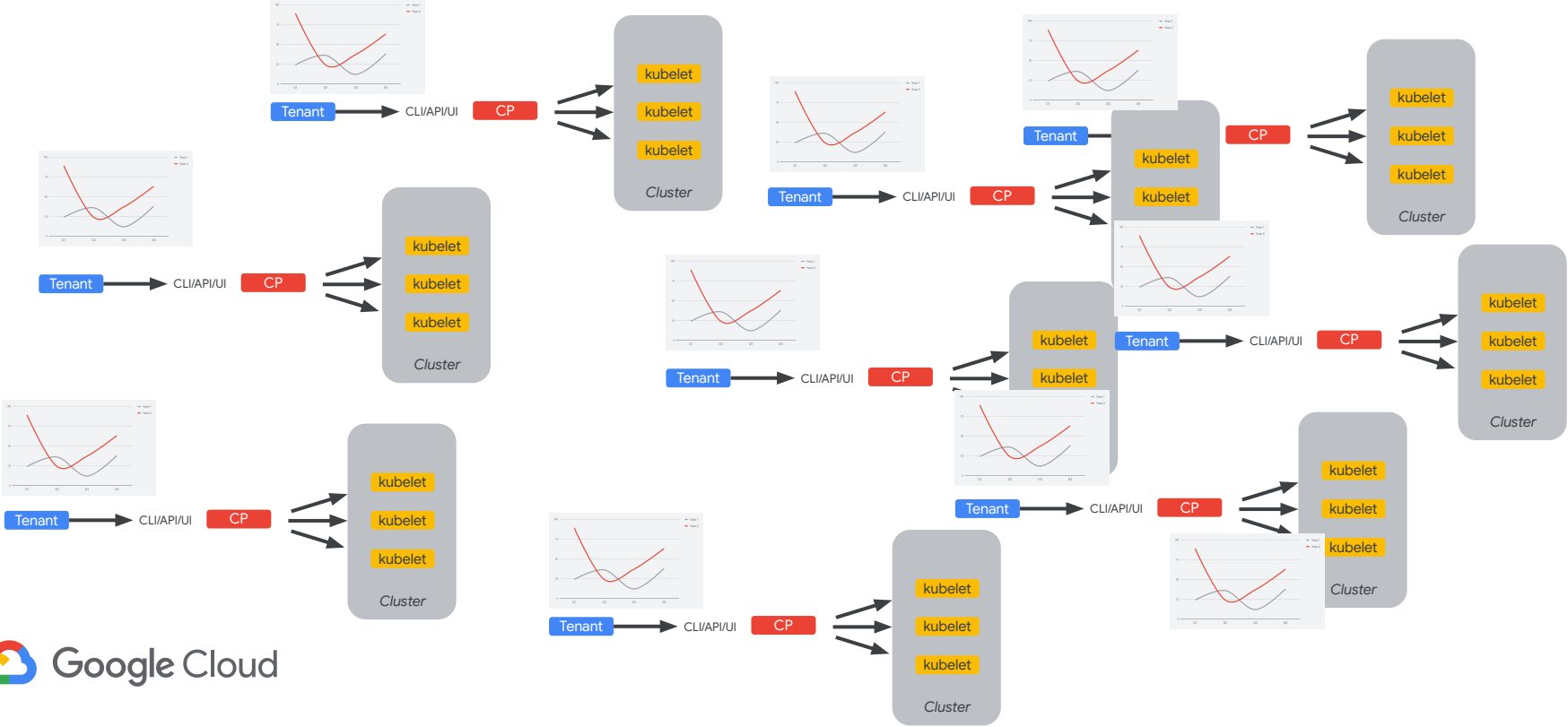
One tenant, one cluster



Multiple tenants, multiple clusters

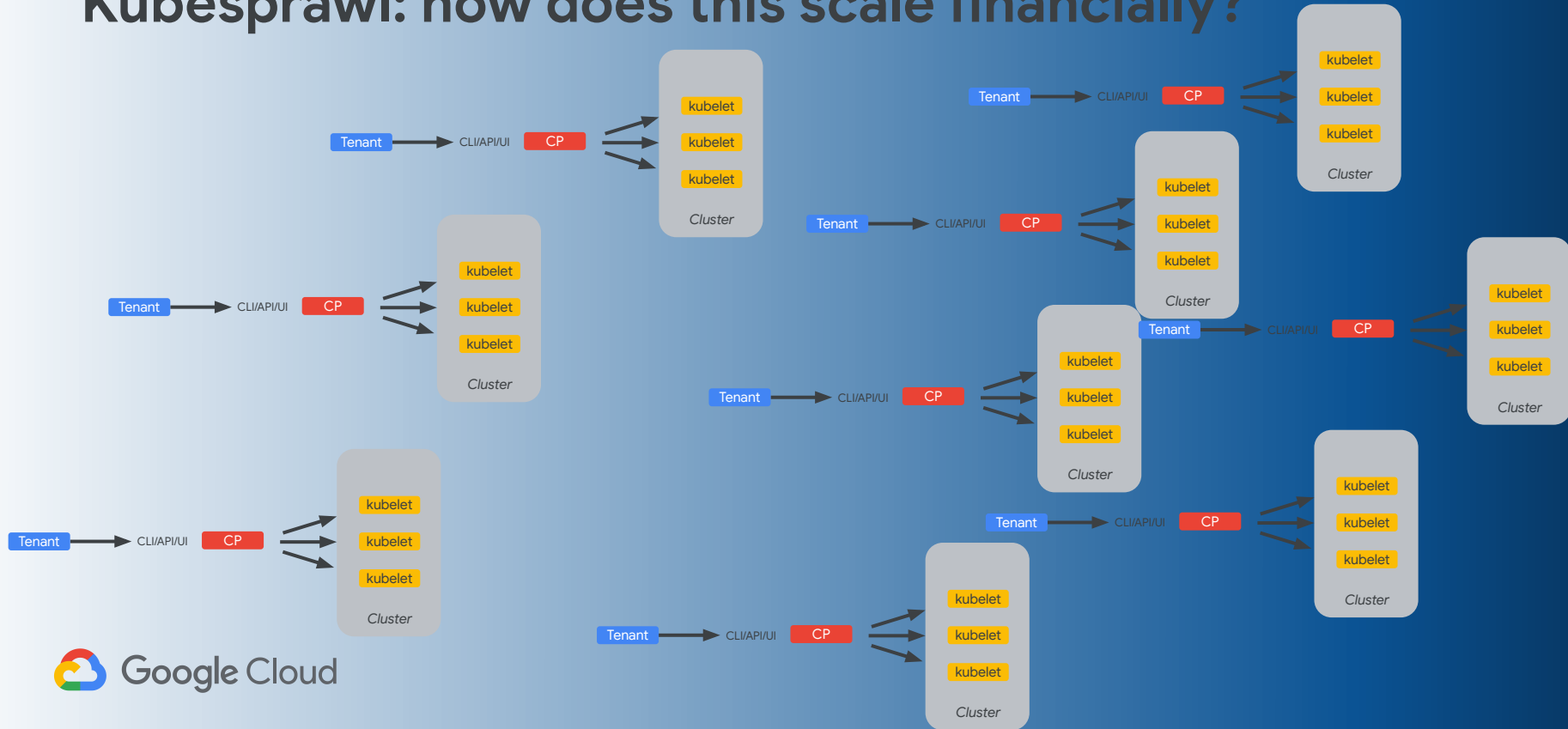


Kubesprawl: how does this scale operationally?

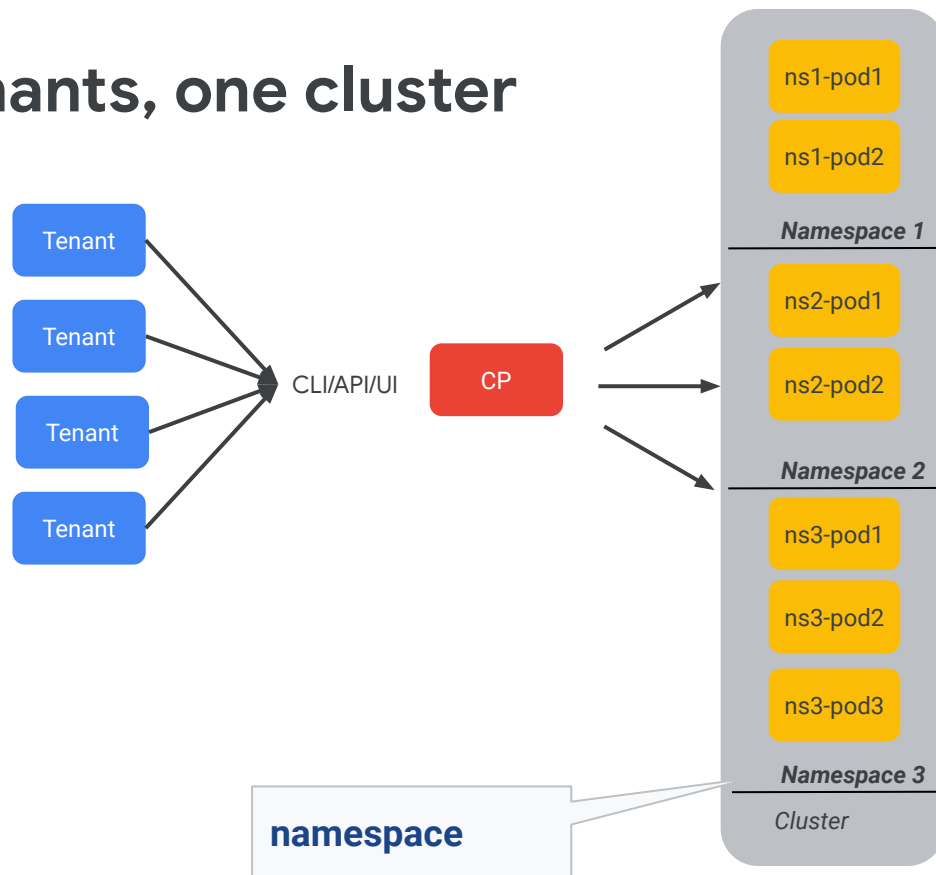




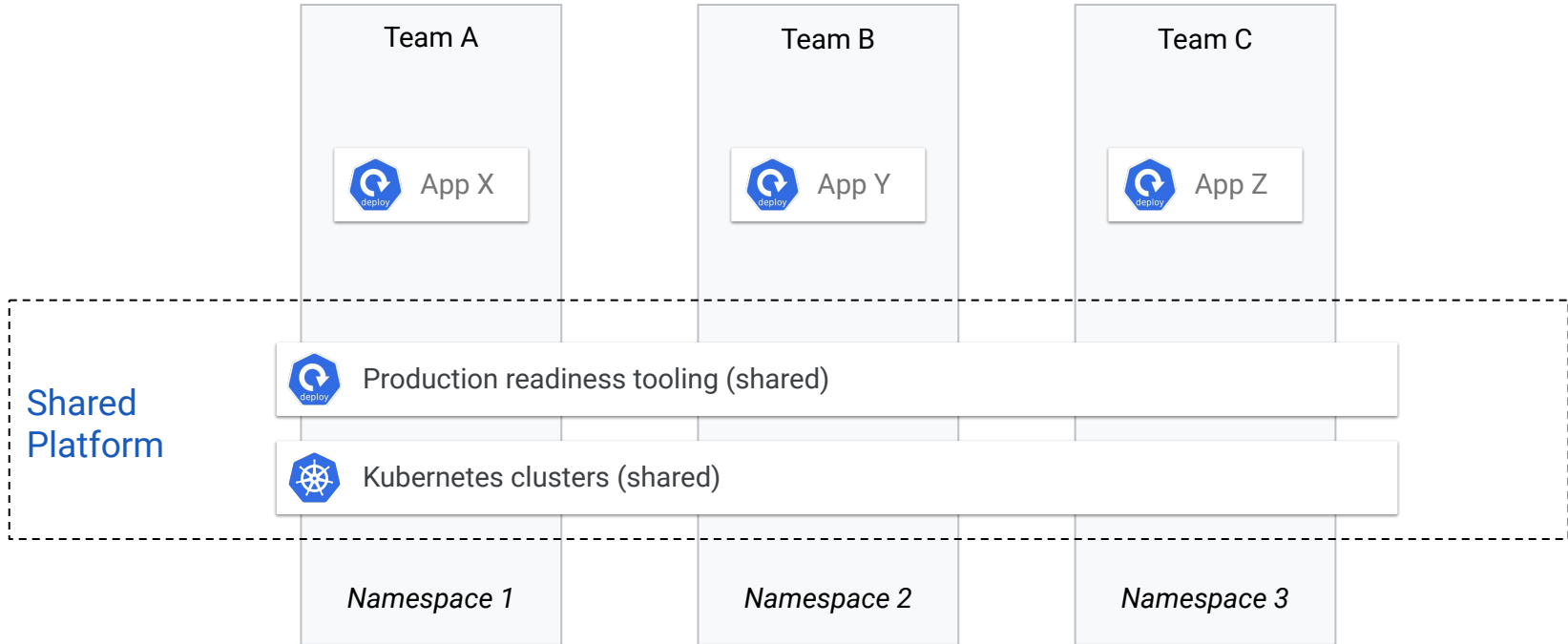
Kubesprawl: how does this scale financially?



Alternative: many tenants, one cluster



Multitenancy reduces overhead and cost



When *should* you use multiple clusters?

Regionalization



Blast radius



Scalability

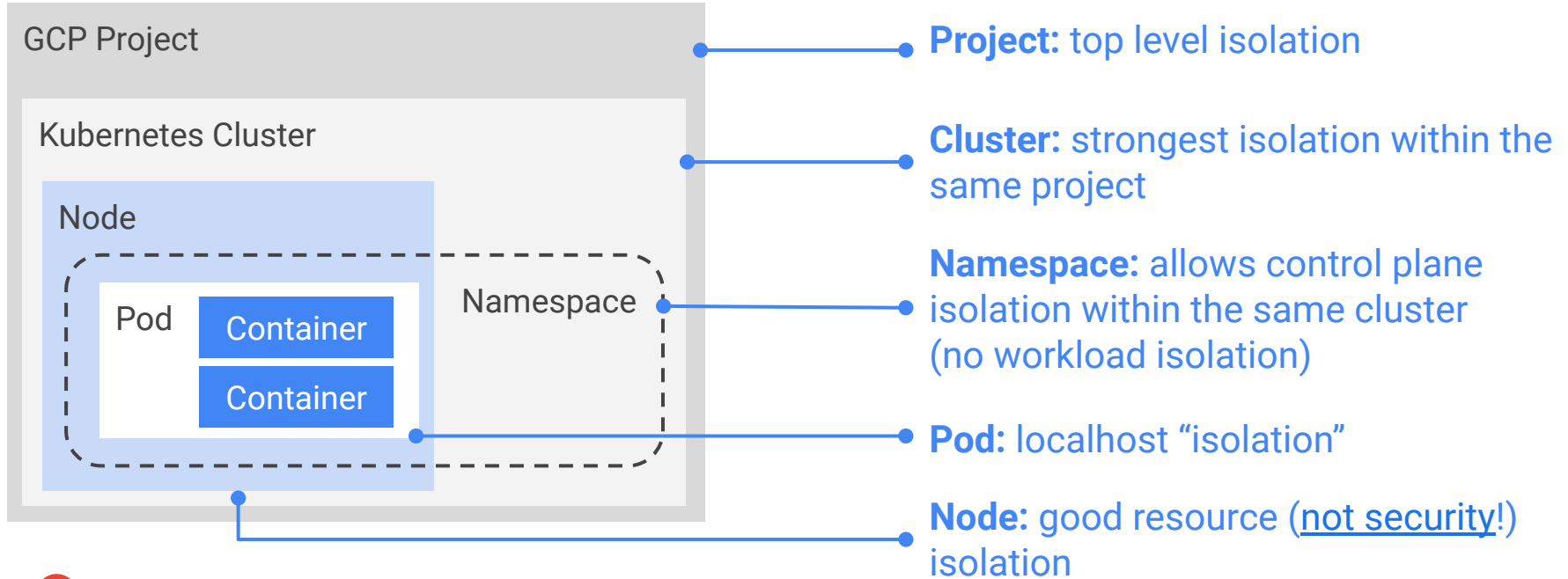


Cluster snowflakes



**What are the
principles of
multitenancy?**

Overview of isolation in Kubernetes



What does “multitenancy in Kubernetes” mean?

Provide **isolation** and
fair **resource sharing** between
multiple users and their workloads
within a single cluster

Best Practices

1

Understand your needs. Not everyone needs everything all at once! Think about cost, overhead, and risk tolerance.

2

Understand the pros/cons of various approaches and technologies to solve your most critical problems.

3

Deploy your solutions, and keep them up-to-date, and iterate to achieve more benefits.

**How do I
implement
multitenancy?**



Namespaces

Namespaces are the primary unit of tenancy in Kubernetes.

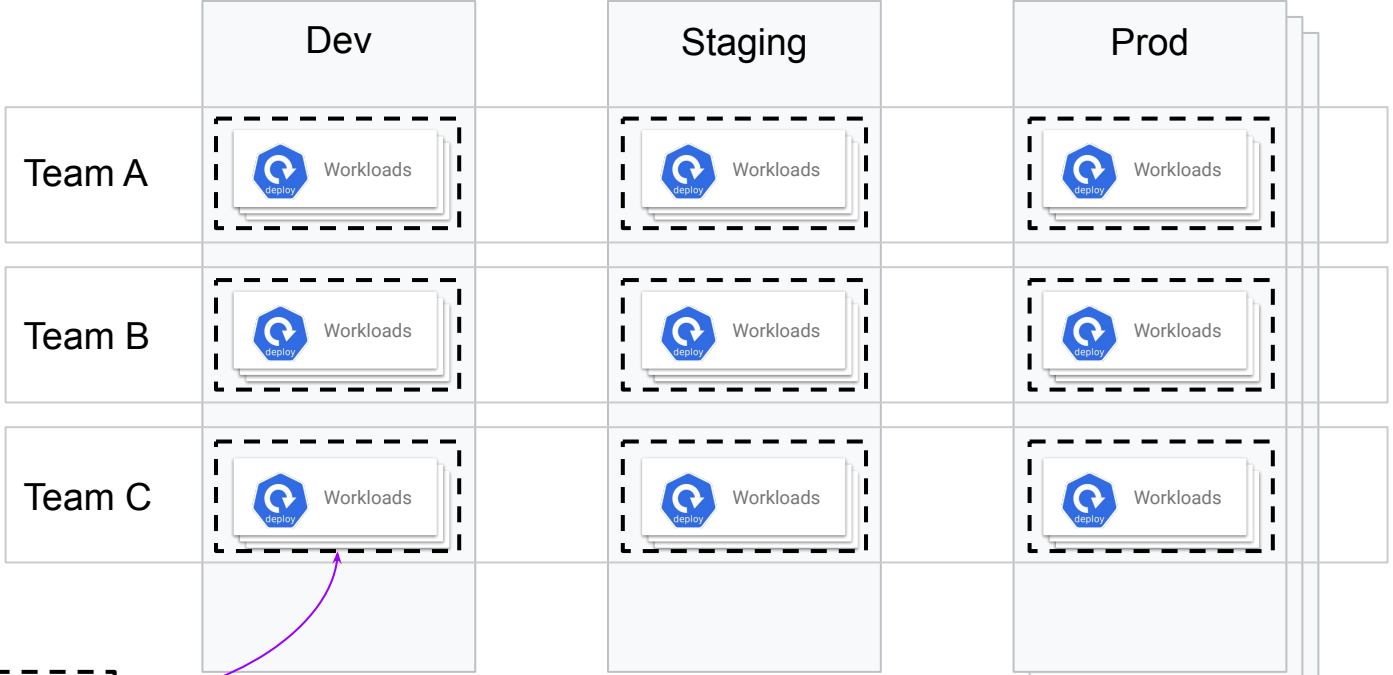
By themselves, they don't do much except organize other objects - but almost all policies support namespaces by default.



Multitenancy across clusters

Clusters

Teams



Namespaces





Properties of namespaces

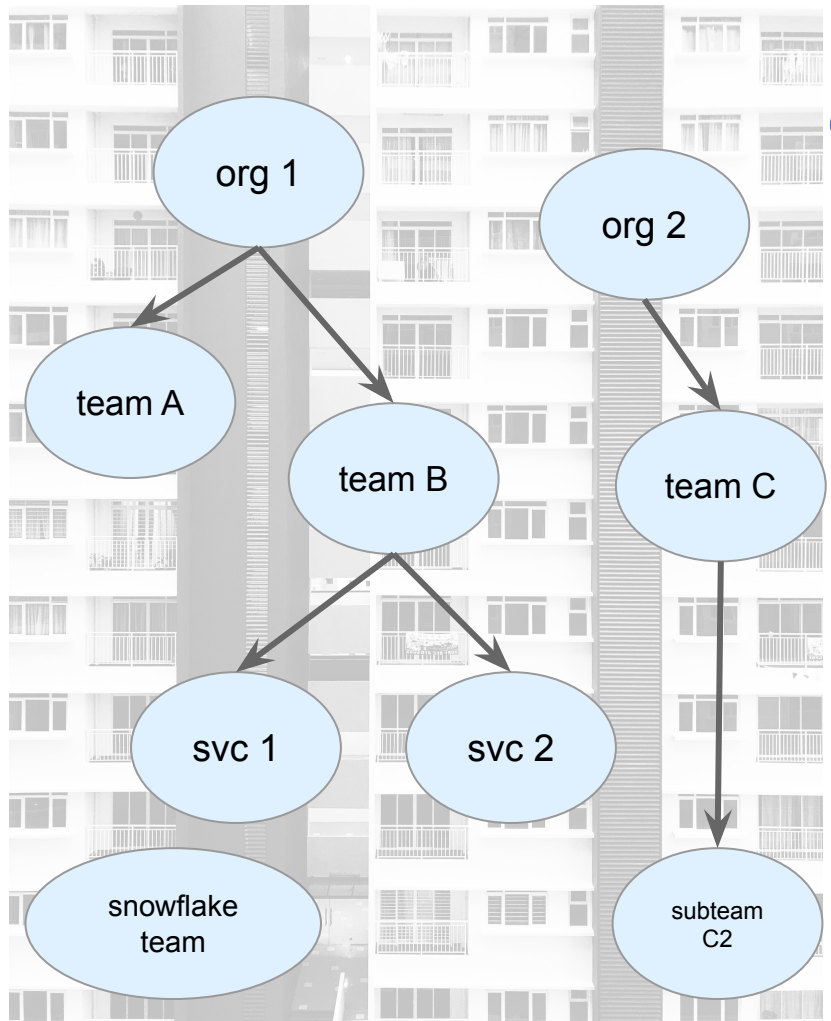
- Require cluster-level permissions to create
- Fully independent policies
- Must be labeled manually to use in policy application (e.g. in Network Policies)
- Included in Kubernetes natively



Hierarchical namespaces

Traditional Kubernetes namespaces are flat, with no relation between them. Hierarchical namespaces express *ownership*, allow for admin delegation and cascading policies.

Hierarchical Namespaces are provided by the [Hierarchical Namespace Controller \(HNC\)](#), a project of wg-multitenancy.





Properties of hierarchical namespaces

- Can use namespace-level “subnamespace” permission to create
- Inherit policies from ancestors
- Can be selected individually or as subtrees
- Provided via the OSS [Hierarchical Namespace Controller](#) (HNC), or as a part of GKE’s [Config Sync](#) (available later in June 2020).

Applying multitenancy

Access
Control

Resource
Sharing

Runtime
Isolation

Insights



Applying multitenancy

Access
Control

Resource
Sharing

Runtime
Isolation

Insights

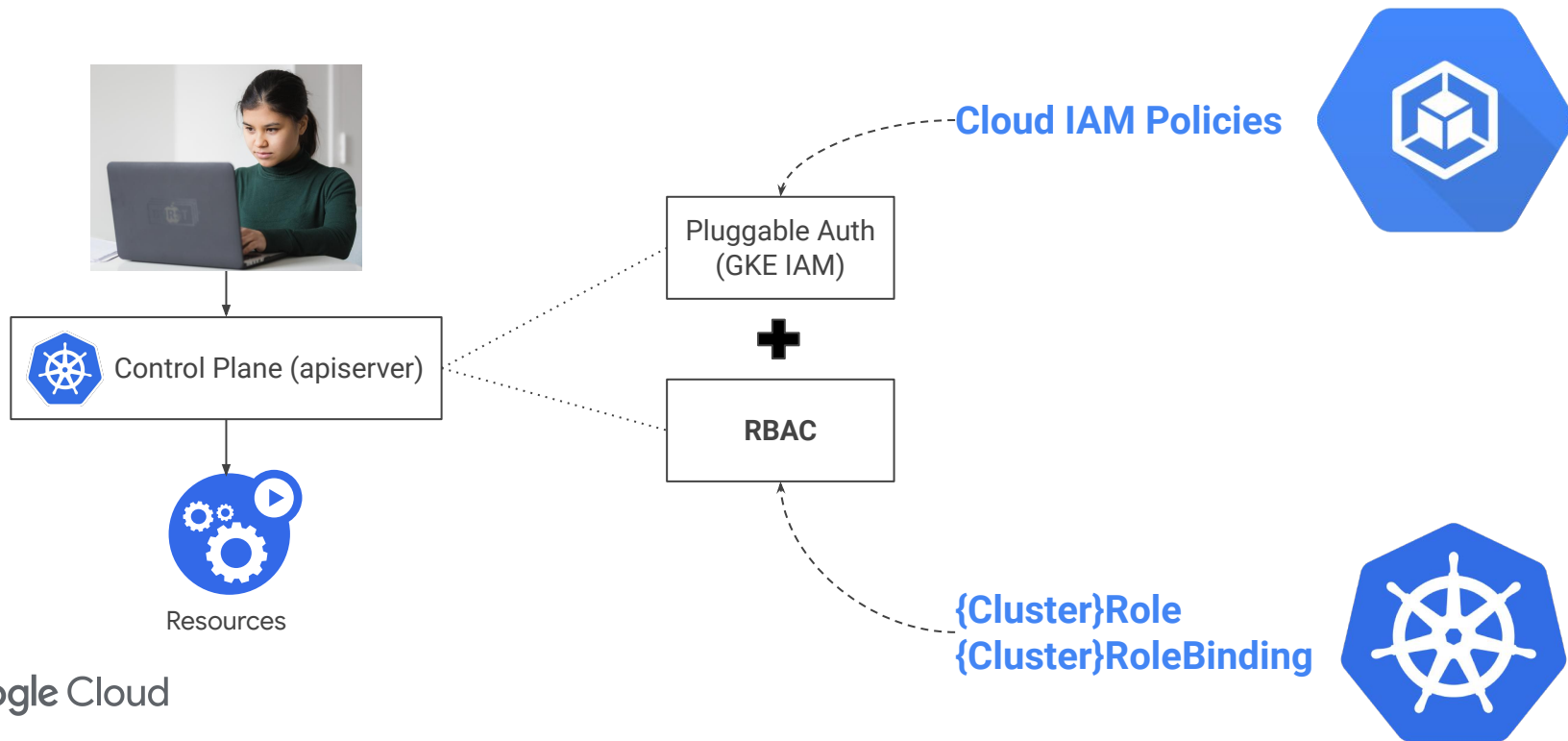


Access control

Tenancy is about *ownership*, and ownership is about control. Job #1 is ensuring that tenants can't control each others' resources.



Authentication and authorization





Role-Based Access Control (RBAC)

RBAC controls access to namespaces in Kubernetes. They're used for:

- Giving **humans** access to Kubernetes resources
 - On GKE, can use Google Groups to give groups of people identical access
- Giving **non-Kubernetes service accounts** access to the Kubernetes API
 - Example: GCP Service Accounts; also work with Google Groups
- Giving access to **pods** calling Kubernetes APIs (with Kubernetes Service Accounts)

Key RBAC concepts:

ClusterRole	A set of cluster-wide permissions. Some useful defaults are preset (e.g. "admin")
Role	Like a ClusterRole, but limited to a single namespace
ClusterRoleBinding	Give a role to one or more <i>subjects</i> (humans, SAs, etc) across the whole cluster.
RoleBinding	Give a role <i>within</i> a single namespace. You can also use ClusterRoles (e.g. "admin") to grant predefined permissions, but limited to that one namespace



GKE Workload Identity

Let Google **manage and rotate the credentials** that are used by your Kubernetes workloads to access GCP services.

- Workloads attaching a Kubernetes SA automatically authenticate as a separate GCP SA when accessing GCP API
- Allows the OAuth scopes and service accounts attached to node pools to follow least privilege

Replaces these workarounds:

- Using node (VM) identity for the pod
- Removes the need for exporting user managed service account keys and embedding them in kubernetes native secrets

Applying multitenancy

Access
Control

Resource
Sharing

Runtime
Isolation

Insights



Resource sharing

Many elements of the cluster are still shared, especially nodes and the apiserver itself. Make sure your tenants can share them fairly.





Resource Sharing

- **Resource Quotas:** no one *namespace* can exceed resource usage. Also used to control dangerous objects (ingress, external services).
- **Limit Range:** no one *pod* can exceed resource usage
- **Pod Affinity/Anti-affinity:** keep pods scheduled together/apart
- **Pod Priority:** pick a winner when there isn't enough to go around

Applying multitenancy

Access
Control

Resource
Sharing

Runtime
Isolation

Insights



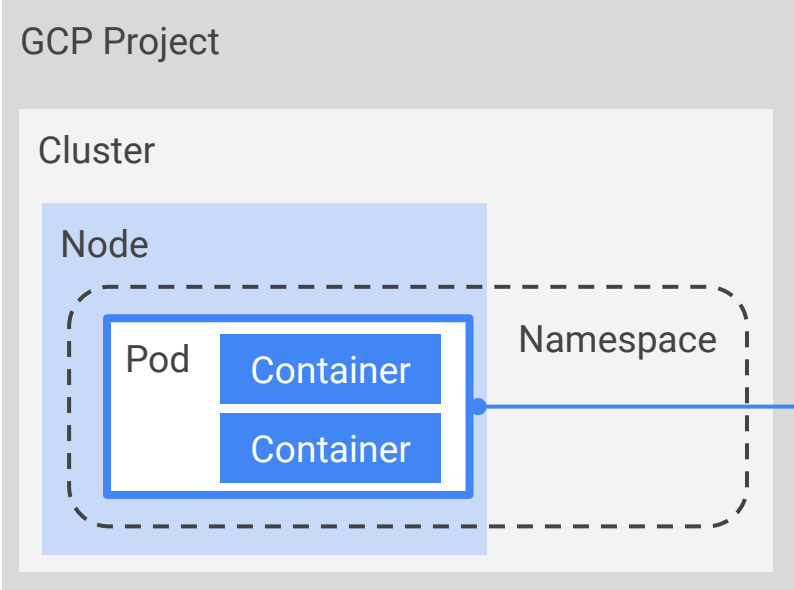
Runtime isolation

Vulnerabilities and attacks are a reality, and containers aren't a security boundary. Consider adding runtime isolation to stop anything getting out of your containers.





Runtime boundaries in Kubernetes



Pod: a unit of scheduling and deploying workloads (no secure isolation by default)



Runtime isolation: overview

- **Pod Security Context:** restrict a pod's workload (eg non-root).
- **Pod Security Policy:** enforce that pods *must* declare an appropriate context. **Hard to enable on a working cluster;** consider alternatives like [OPA Gatekeeper](#) or [Anthos Policy Controller](#).
- **Network Policy:** forbid pods from talking to each other if they have no good reason to.
- **Runtime Class:** run pods in sandboxes like gVisor or Kata Containers.



Runtime isolation: GKE Sandbox

GKE's built-in Runtime Class protects your pods.

Adds a security boundary to containers in GKE based on **gVisor**.

Defense-in-depth security principles **without application changes**, new architecture models, or added complexity.



Applying multitenancy

Access
Control

Resource
Sharing

Runtime
Isolation

Insights



Insights

Tenants need to be able to observe themselves, and you need to be able to observe (and charge!) your tenants.





GKE Usage Metering

- View workloads' resource usage in BigQuery, broken down by namespace and labels
 - Memory, CPU, GPU, PD, network, etc.
- Join usage data with GCP Billing Export data to compute resource costs per tenant

GKE Usage Metering



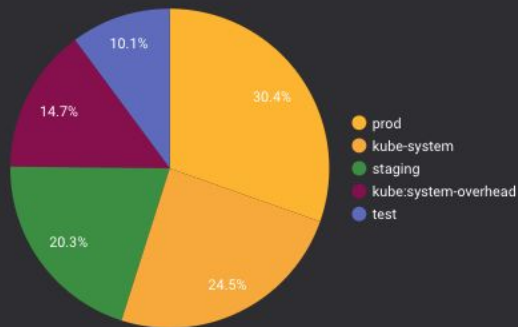
GKE usage metering

Nov 1, 2018 - Nov 13, 2018



Monthly View

Cost breakdown by namespace



namespace	cost
1. prod	32.96
2. kube-system	26.53
3. staging	21.97
4. kube:system-overhead	15.92
5. test	10.99

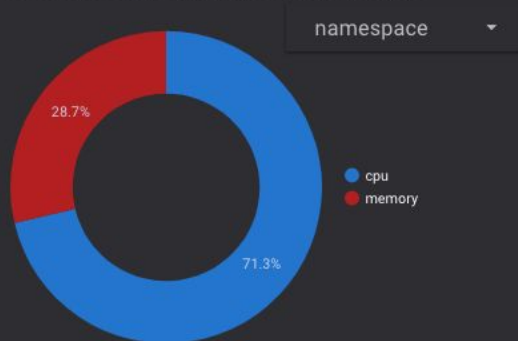
Total GKE cost:

108.37

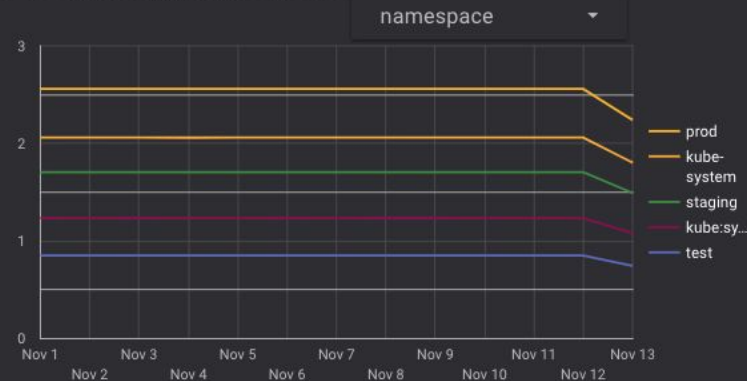
Total GCP cost:

133.25

Cost breakdown by resource type



Cost trends by namespace

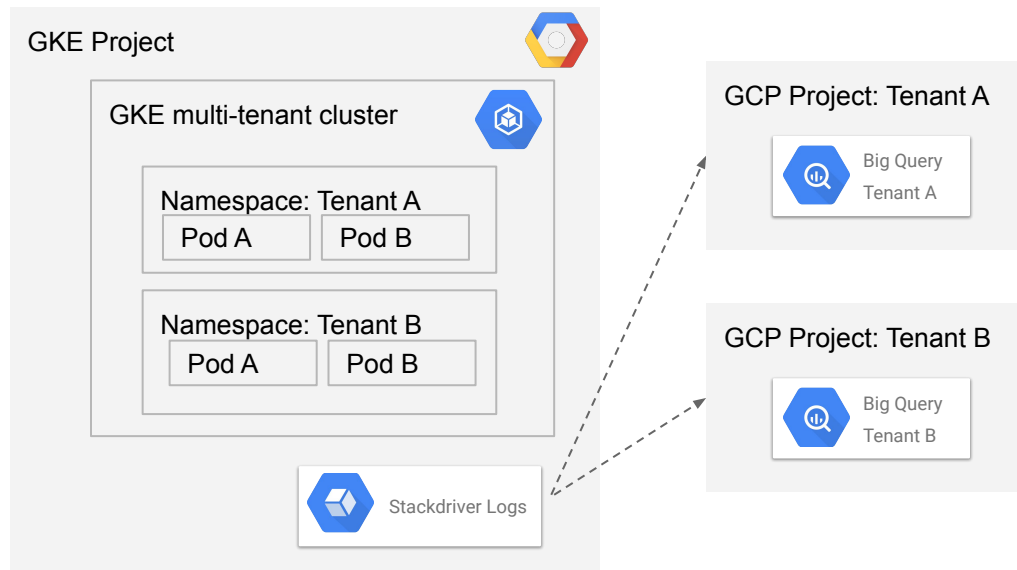




Multitenant logging

Logs often refer to sensitive data, so consider controlling access to them carefully.

- Create a tenant project with a BigQuery dataset for each team.
- Use Log Routing to **filter logs by namespace** and send them to the correct BigQuery dataset.

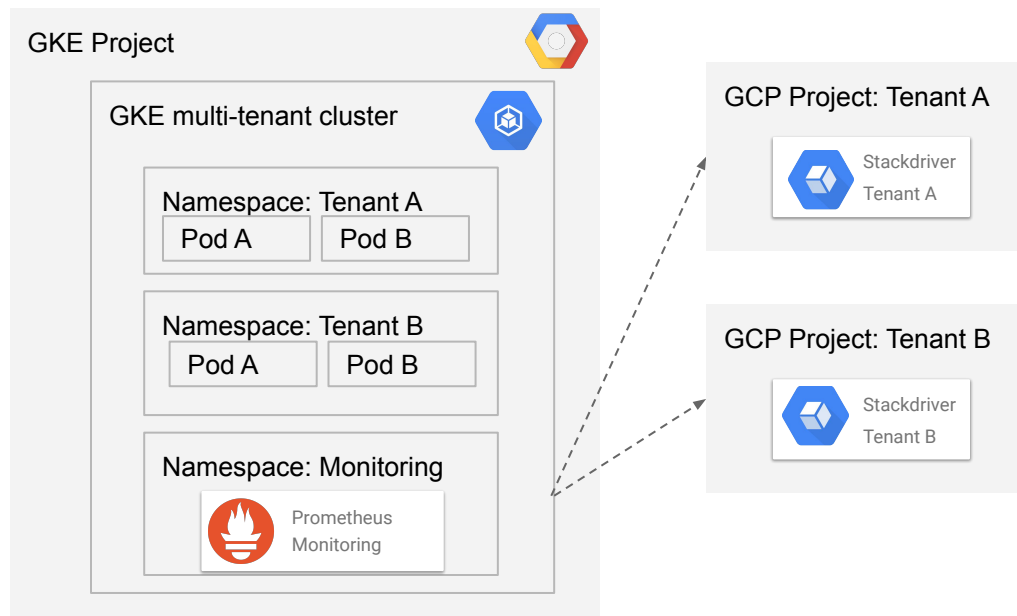




Multitenant monitoring

Metrics are often less sensitive than logs, but you can still control access by tenant if you like.

- Create a tenant project per team.
- As part of the onboarding process, create a **Prometheus / Stackdriver adapter** with a **per namespace config**
- Send each tenant's metrics to the correct project.



Advanced topics





Policy deployment

Be careful about storing your source-of-truth on your cluster.

- Check your policies (e.g. RBAC, Network Policy) into Git as YAML files.
- Have your cluster admin apply them based on Git, or use a CD tool like [GKE Config Sync](#) or [Anthos Config Management](#).
- Test out changes on a canary cluster first, even in prod!
- Separate your policies from your workloads.
 - Possible exception: DaemonSets



Policy enforcement and auditing

Use [OPA Gatekeeper](#) (or [Anthos Policy Controller](#)) to define and apply custom policies.

- Define rules in Rego (Python-like rule language)
- Apply them across your clusters
- Audit violations
- Useful alternative to Pod Security Policies



Further exploration

Hard multitenancy is loosely defined as the condition where tenants are mutually hostile, not relatively co-operative.

- **Virtual Clusters**: an wg-multitenancy project to give each tenant its own control plane while sharing a data plane. Must be combined with sandboxing.
- **SaaS multitenancy**: many different instances of the same application for different **consumers**. Generally requires sandboxing and control plane automation.

Your needs will be very specific to your threat model!

Conclusion

Best Practices

1

Understand your needs. Not everyone needs everything all at once! Think about cost, overhead, and risk tolerance.

2

Understand the pros/cons of various approaches and technologies to solve your most critical problems.

3

Deploy your solutions, and keep them up-to-date, and iterate to achieve more benefits.

Learning more...

The screenshot shows the Google Cloud documentation website. At the top, there is a navigation bar with the Google Cloud logo and links for 'Why Google', 'Solutions', 'Products', 'Pricing', 'Getting Started', 'Docs', and 'Support'. A search icon, a language dropdown set to 'English', and a 'Console' button with a user profile picture are also present. Below the navigation bar, there is a secondary menu with 'Containers', 'Guides', 'Reference', 'Support', and 'Resources'. A 'Contact Sales' button is located on the right side of this menu. The main content area features a left-hand sidebar with a list of topics under 'types of clusters', including 'Private clusters', 'Cluster upgrades', 'Maintenance windows and exclusions', 'Guidelines for creating scalable clusters', 'Node pools', 'Node images', 'Using a node image with containerd', 'Regional clusters', 'Release channels', and 'Alpha clusters'. The main article title is 'Best practices for enterprise multi-tenancy', with a breadcrumb trail: 'Containers > Google Kubernetes Engine (GKE) > Documentation'. To the right of the title are five star icons and a 'Send feedback' button. Below the title is a 'Contents' section with a dropdown arrow, listing 'Assumptions and requirements', 'Setting up folders, projects and clusters', and 'Establish a folder and project hierarchy'.

Google Cloud Why Google Solutions Products Pricing Getting Started Docs Support > English Console

Containers Guides Reference Support Resources Contact Sales

types of clusters

- Private clusters
- Cluster upgrades
- Maintenance windows and exclusions
- Guidelines for creating scalable clusters
- Node pools
- Node images
- Using a node image with containerd
- Regional clusters
- Release channels
- Alpha clusters

Containers > Google Kubernetes Engine (GKE) > Documentation ☆☆☆☆

Best practices for enterprise multi-tenancy

Send feedback

Contents ▾

- Assumptions and requirements
- Setting up folders, projects and clusters
 - Establish a folder and project hierarchy

[Available from the GCP docs website](#)

Learning more...

Kubecon San Diego had some great presentations on building multitenancy systems on Kubernetes. My favourites include:

- [Walls within walls: what if your attacker knows parkour?](#)
- [Kubernetes at Cruise: two years of multitenancy](#)
- Plus [two sessions](#) from the multitenancy working group (wg-multitenancy)

Some other interesting links to follow include:

- [The Multitenancy Working Group](#)
- [Mercari's experience with multitenant Istio](#)

Multitenancy reduces overhead and cost

