



KubeCon

— North America 2017 —

Multitenancy Deep Dive

Thursday, December 7 • 2:00pm - 3:20pm

David Oppenheimer (Google) davidopp@google.com

Quinton Hoole (Huawei) quinton.hoole@huawei.com

Agenda

- Presentations
- Discussion of topics of interest
- Ideas for 2018 (including whether we should create a Working Group)

Presentations

- Quinton Hoole, Huawei
- Jessica Frazelle, Microsoft
- Harry Zhang, Hyper
- David Oppenheimer, Google
- Tim Allclair, Google

Multi-Tenancy Models

Jessie Frazelle - Microsoft
refer to [original doc](#)

Soft Multi-Tenancy

- multiple users within the same organization in the same cluster.
- could have possible bad actors such as people leaving the company, etc.
- Users are not thought to be actively malicious since they are within the same organization, but potential for accidents or “evil leaving employees.”
- A large focus is to prevent accidents.

Hard Multi-Tenancy

- multiple users, from various places, in the same cluster.
- means that anyone on the cluster is thought to be potentially malicious and therefore should not have access to any other tenants resources.

Access to k8s API

For our purposes, we only run untrusted workloads, but we have our own trusted API on top of the kubernetes API

(seems like SaaS from [davidopp's doc.](#))

A different multi-tenancy models would also restrict access to the API and create roles, etc. Refer to that doc for more details.

Host OS

Container Runtime

Network

DNS

AuthN/AuthZ

Isolation of Master and
System nodes.

Isolation of system
services.

Restricting access to
host resources.

Environment Variables

Thoughts about Hard Multi-tenancy in Kubernetes with Hypervisor based Container Runtimes

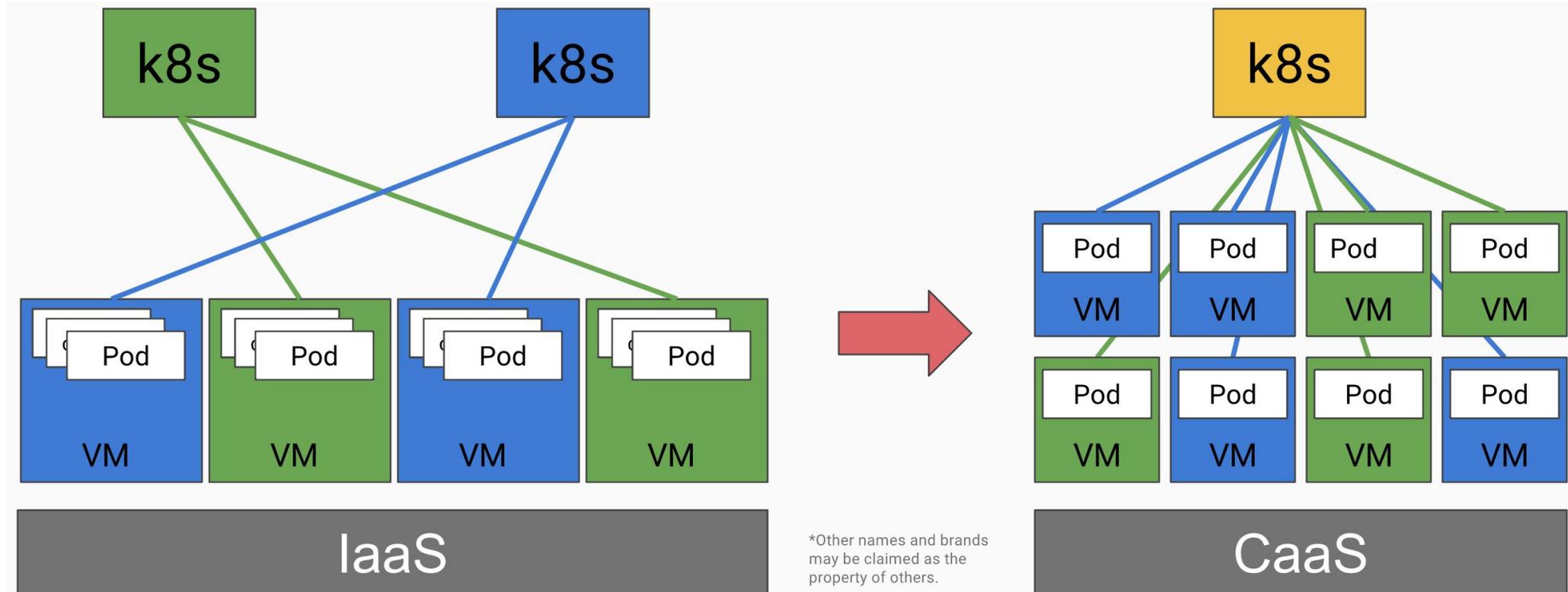
Harry (Lei) Zhang @resouer

Background

- Hypernetes (Stackube):
 - A multi-tenant Kubernetes distro with hypervisor based container runtime
 - runV, now upgrading to KataContainers
 - <https://github.com/openstack/stackube>
 - Upstream Kubernetes + customized plugins
 - The core system behind <https://hyper.sh/>
 - Passed 100% conformance e2e tests

Container Runtime: Isolation & Security

- [KataContainers](#)



Container Runtime: OS Multi-Tenancy

BYOK (Bring Your Own Kernel):

annotations:

```
com.github.katacontainers.KernelPath: "/boot/vmlinuz-custom-myversion"
```

(This has already been implemented)

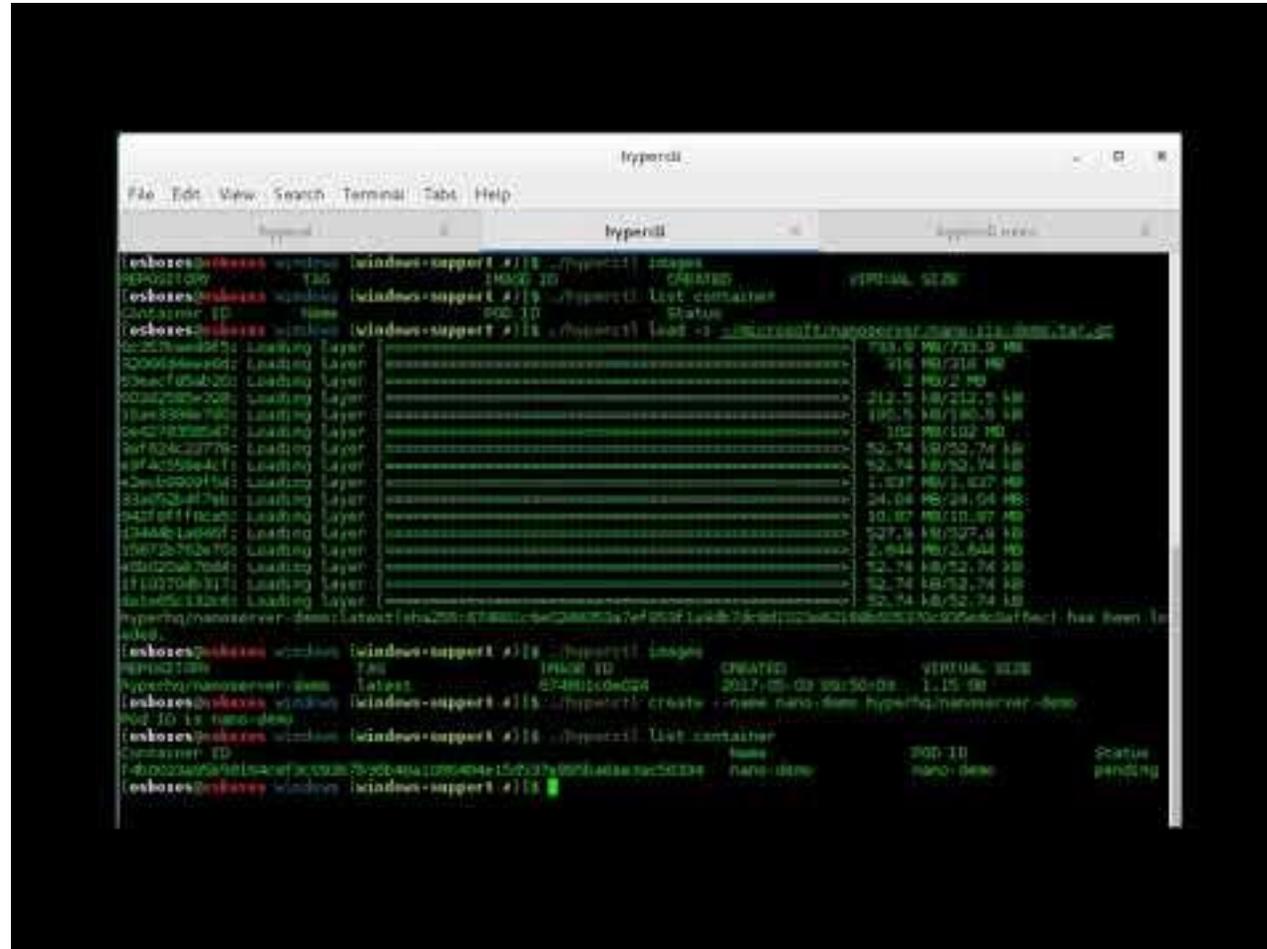
Or even:

annotations:

```
com.github.katacontainers.KernelPath: "/boot/windows-nano-server-myversion"
```

(This has also been concept proved)

Container Runtime: OS Multi-Tenancy



```
Hyper-V
File Edit View Search Terminal Tabs Help
Hyper-V Hyper-V
(esboxes@esboxes ~) windows-support #118 ~/hyperctl images
REPOSITORY TAG IMAGE ID CREATED VIRTUAL SIZE
(esboxes@esboxes ~) windows (windows-support #118 ~/hyperctl) List container
Container ID Name Pool ID Status
(esboxes@esboxes ~) windows (windows-support #118 ~/hyperctl) last -i -p Microsoft.WindowsServer.Manage/WindowsServer-2016-10-14-01
0c2c79a4d3d5f Loading Layer ..... 733.9 MB/733.9 MB
7000634ee4d0 Loading Layer ..... 216 MB/216 MB
03eac705ab200 Loading Layer ..... 2 MB/2 MB
002a702b7008 Loading Layer ..... 212.5 MB/212.5 MB
15aa3399a790 Loading Layer ..... 195.5 MB/195.5 MB
0e4278378547 Loading Layer ..... 102 MB/102 MB
9af824c22776c Loading Layer ..... 52.74 MB/52.74 MB
e9f4c55e44f1 Loading Layer ..... 52.74 MB/52.74 MB
42eb02009f94 Loading Layer ..... 1.037 MB/1.037 MB
83a652b4f790 Loading Layer ..... 24.04 MB/24.04 MB
94278f1f0c40 Loading Layer ..... 10.87 MB/10.87 MB
17446c1a9807 Loading Layer ..... 527.9 MB/527.9 MB
15672c762e70c Loading Layer ..... 2.644 MB/2.644 MB
e01020a705d4 Loading Layer ..... 52.74 MB/52.74 MB
1f102708b317 Loading Layer ..... 52.74 MB/52.74 MB
8ab665c13260 Loading Layer ..... 52.74 MB/52.74 MB
hyperfq/nanoserver-demos/latest/sha256-87880c-5e08800da7ef200f1a4b73c4d12c9a219800370c905e4c4ef4e1 has been loaded.
(esboxes@esboxes ~) windows (windows-support #118 ~/hyperctl) images
REPOSITORY TAG IMAGE ID CREATED VIRTUAL SIZE
hyperfq/nanoserver-demos/latest 87880c0e024 2017-05-03 09:50:08 1.15 MB
(esboxes@esboxes ~) windows (windows-support #118 ~/hyperctl) create --name nano-demo hyperfq/nanoserver-demos/latest is nano-demo
(esboxes@esboxes ~) windows (windows-support #118 ~/hyperctl) List container
Container ID Name Pool ID Status
7ab002549598154ac4f3cc38 26-cb49a1088-494-1552378001488a3c3d334 nano-demo nano-demo pending
(esboxes@esboxes ~) windows (windows-support #118 ~/hyperctl)
```

Tenant

- Stackube:
 - Tenant == Namespace
 - CRD
 - tenant controller with RBAC
 - Keystone
- Q:
 - Do we need nested namespace? (One Tenant with multiple namespaces?)
 - Tenant == Namespace, or Tenant 1:N Namespace

Network

- Stackube:
 - One Network per Tenant
 - CRD
 - network controller
 - Neutron CNI plugin
 - L2 network isolation
 - Pods of same Tenant belong to same subnet
- Q:
 - Do we really need Network API object? Or Network Policy should be the plan?.
 - And what about multiple networks?
 - Is L2 isolation specially preferred for hard multi-tenancy?
 - Do we need to isolate Nodes and Pods by different subnets?

DNS

- Stackube:
 - Per kube-dns per namespace (tenant)
- Q:
 - Discussion: <https://github.com/kubernetes/dns/issues/132>
 - Other approach:
 - Sidecar, we use this in old version of Hypernetes
 - Enforce by CoreDNS (+RBAC)

Summary

- KataContainers can play an important role in hard-multitenant Kubernetes
 - Thanks to CRI
- While other aspects like Tenant, Network, DNS etc still need to be clearly defined or updated to build the whole stack up.
- Then what is the Kubernetes/Cloud Native way to solve them?



KubeCon

— North America 2017 —

Multitenancy taxonomies

David Oppenheimer, *Google*
December 7, 2017

Control plane vs. node multitenancy

All policies are specified through the control plane.

Distinction is whether policy controls sharing of control plane or nodes.

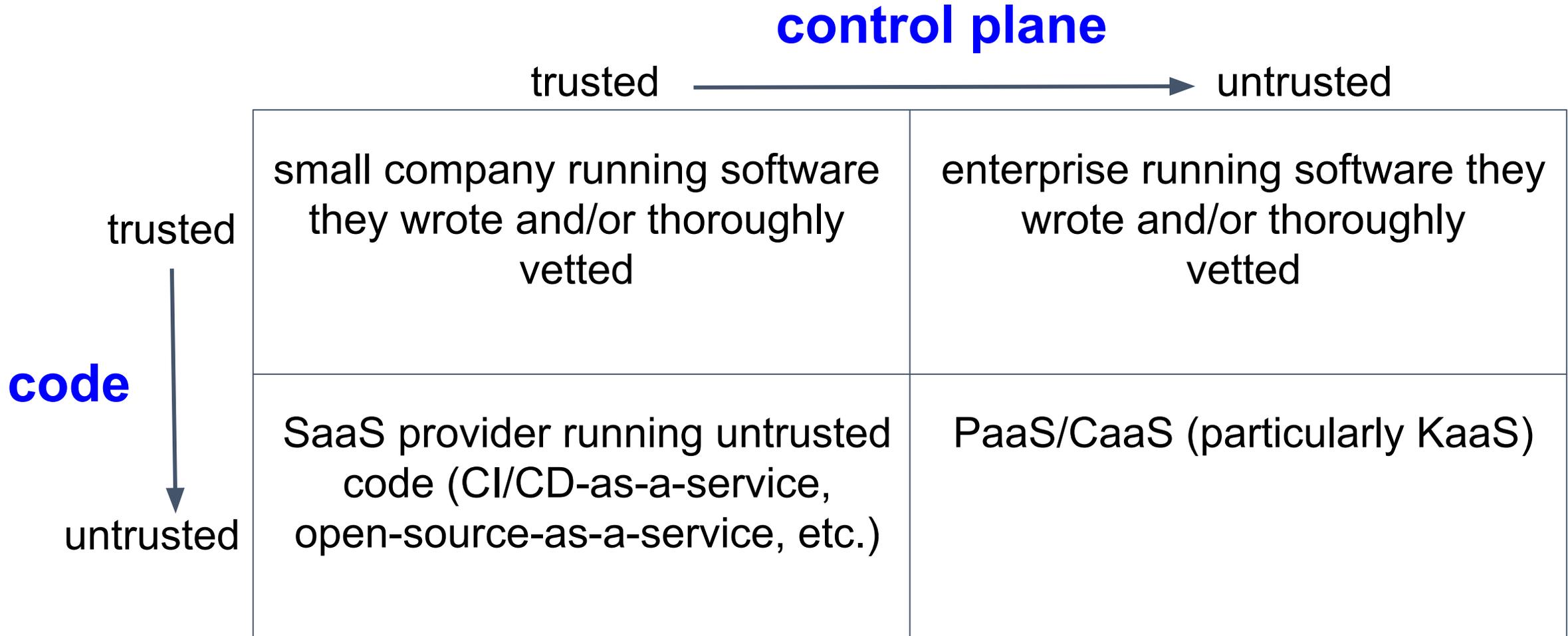
Control plane multitenancy

- RBAC
- EventRateLimit admission controller

Many for node multitenancy

- ResourceQuota / LimitRange / request / limit / priority
- node affinity, pod affinity, taints/tolerations
- PodSecurityPolicy
- NetworkPolicy

How control plane and node are used



Other axes

- What do users see?
 - objects (or subset) in user's namespace(s) or all namespace(s)?
 - nodes?
 - metadata about other tenants (namespace collision, service names in DNS, etc.)?
- Node-level isolation mechanism
 - containers + PodSecurityPolicy, seccomp, AppArmor, SELinux, ...
 - container + hypervisor (nested virtualization)
 - dedicated nodes (taints/tolerations or anti-affinity)



KubeCon

— North America 2017 —

Secure Containers

Tim Allclair, *Google*

Secure Containers

Stronger Isolation

- Sandboxing untrusted code
- VM strength isolation

Work in Progress

- CRI-O with Clear Containers
- Frakti with runV
- (soon!) Kata Containers
- Cloud providers exploring CaaS

It's time to agree on
the abstractions,
before we diverge too
much.

Open Questions

We're kicking off the discussions now.

What are the properties of a sandbox?

- Must it employ full virtualization technology?
 - Or could a sandbox be a very restrictive seccomp profile?
- What does sandboxing imply about networking?
- What does sandboxing imply about auth[nz]?
- What features is it OK to break with a sandbox?
 - E.g. cross-container IPC? host namespaces? etc.

Where is the sandbox boundary?

Pods?

- Easier resource sharing & communication between containers.
- Better for a serverless (nodeless) model
- Much simpler networking

Containers?

- Finer grained control allows for models like trusted sidecars

Or should we consider something else?

Namespaces? Sandbox resource? A combination of pods + containers?

API Design

How do we surface sandboxes to the user?

- Explicit, without choice of backend?
 - Sandbox *bool
- Explicit, with choice of backend?
 - Sandbox string
- Implicit, derived from security attributes?
 - See: Entitlements

Implementation details

Sandboxed & unsandboxed containers should live side-by-side.

Should sandboxing be enforced by the runtime (CRI), or the Kubelet?

Does the kubelet decide which CRI server to talk to, or just pass the sandbox bit on the CreateContainer request?

Stay tuned!

Look for a design proposal soon...

Expect more discussions in sig-node meetings

Thoughts? Questions? Get in touch!

- tallclair@google.com
- **@tallclair** (github, slack, twitter)

Possible group discussion topics (1)

- How are you (and/or your customers) using the existing Kubernetes multitenancy features? What problems/use cases are you solving?
- What would you (and/or your customers) like to do, but can't (or are rolling your own, and would like it supported in upstream)? What problems/use cases/pain points would this address?

(Consider both control plane and node support for multitenancy)

Possible group discussion topics (2)

- Hierararchy vs. labels vs. good enough how it is
 - policies that span namespaces and/or apply within a namespace?
- Need better hiding of tenants from one another?
- Issues with tenants DoSing each other via the control plane?
- Uses cases and missing features for “hard” multitenancy
 - need more isolation in the control plane?
 - need “secure containers”?

2018 planning

- Should we create a multitenancy Working Group?
- Specific multitenancy features you want and/or are interested in working on?

Note: a mailing list has been set up -- please join it!

<https://groups.google.com/forum/#!forum/kubernetes-wg-multitenancy>