**redhat**

# Extending Kubernetes 101

Michael Hausenblas **@mhausenblas**

Developer Advocate, Red Hat

2018-11-15, ContainerConf, Mannheim

# $ whois mhausenblas

- **Developer Advocate @ Red Hat (Go, Kubernetes, OpenShift)**

- Developer Advocate @ Mesosphere (Mesos, DC/OS, Kubernetes)

- Chief Data Engineer @ MapR (HDFS, HBase, Drill, etc.)

- Applied research (4y in Ireland, 7y in Austria)

- Nowadays mainly developing tools in Go (Python, Node, Java, C++)

- Kinda developer turned ops (aka appops)

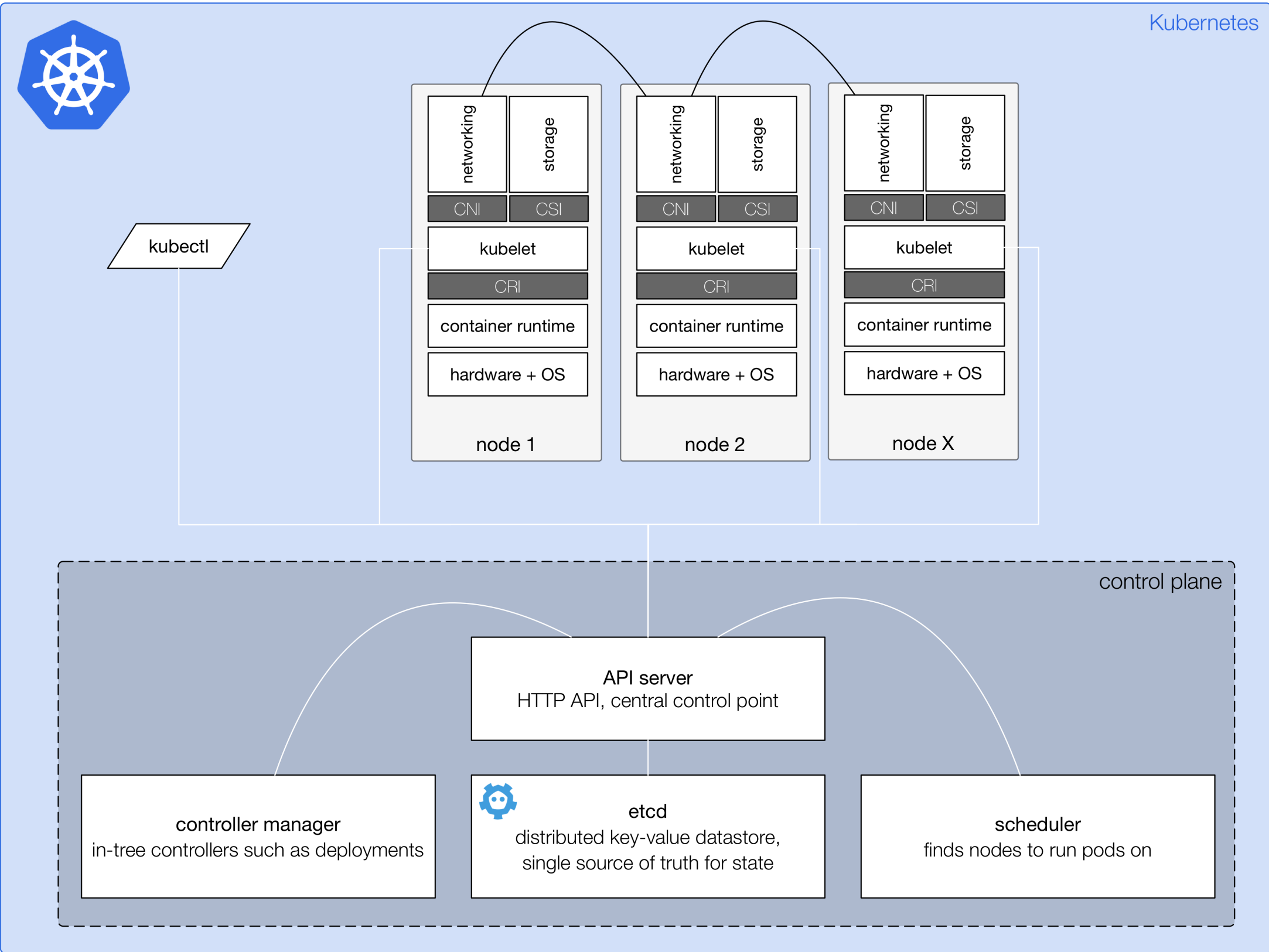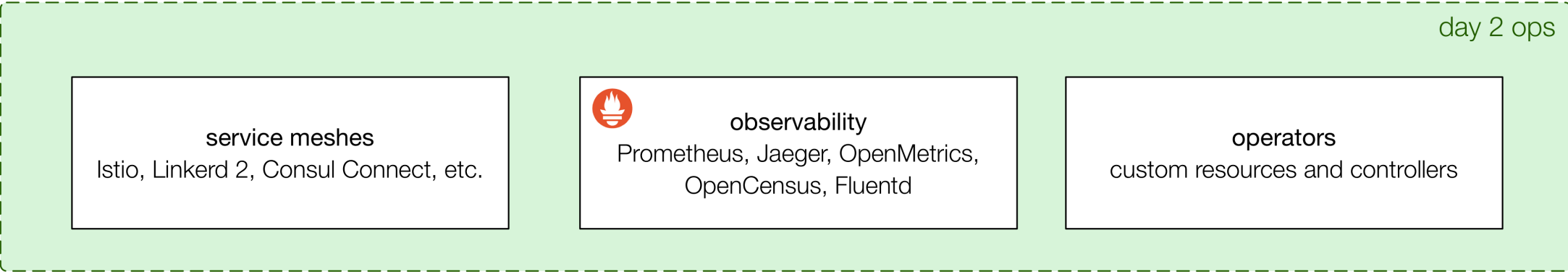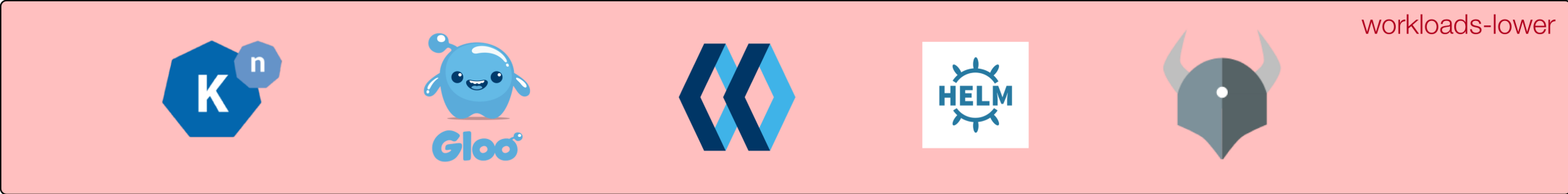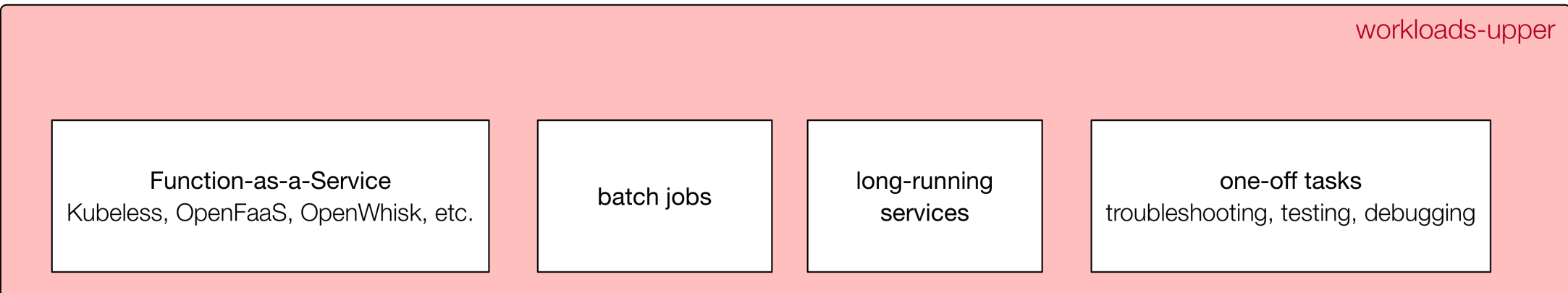Hit me up on Twitter: **@mhausenblas**

**red**hat.

# Kubernetes 101

"Begin at the beginning," the King said, very gravely, "and go on till you come to the end: then stop."

— Lewis Carroll, Alice in Wonderland
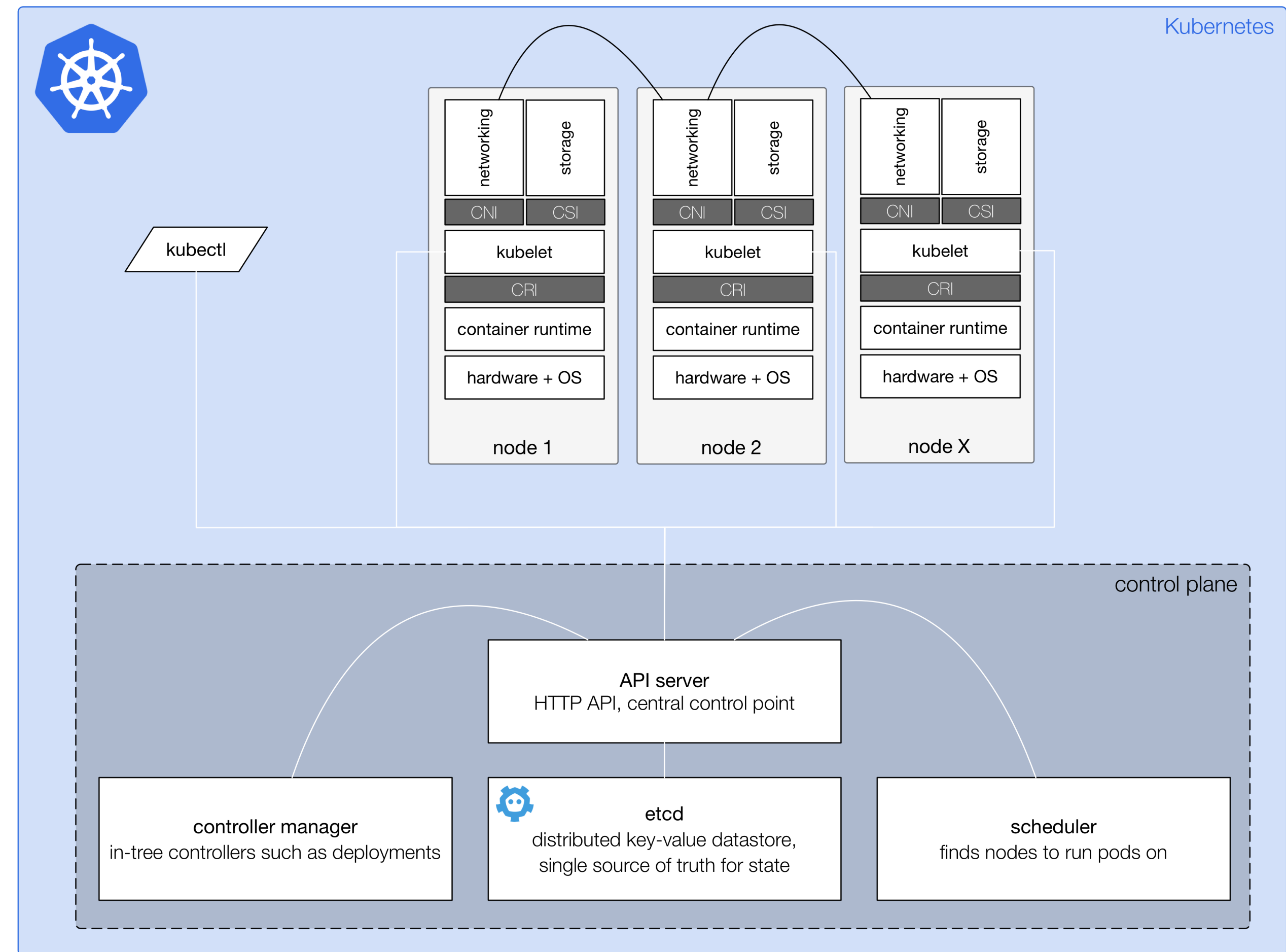
tags: humor

redhat.

workloads-upper

Function-as-a-Service
Kubeless, OpenFaaS, OpenWhisk, etc.

batch jobs

long-running services

one-off tasks
troubleshooting, testing, debugging

workloads-lower

day 2 ops

service meshes
Istio, Linkerd 2, Consul Connect, etc.

observability
Prometheus, Jaeger, OpenMetrics, OpenCensus, Fluentd

operators
custom resources and controllers

Kubernetes

networking    storage
CNI    CSI
kubelet
CRI
container runtime
hardware + OS
node 1

networking    storage
CNI    CSI
kubelet
CRI
container runtime
hardware + OS
node 2

networking    storage
CNI    CSI
kubelet
CRI
container runtime
hardware + OS
node X

kubectl

control plane

API server
HTTP API, central control point

controller manager
in-tree controllers such as deployments

etcd
distributed key-value datastore, single source of truth for state

scheduler
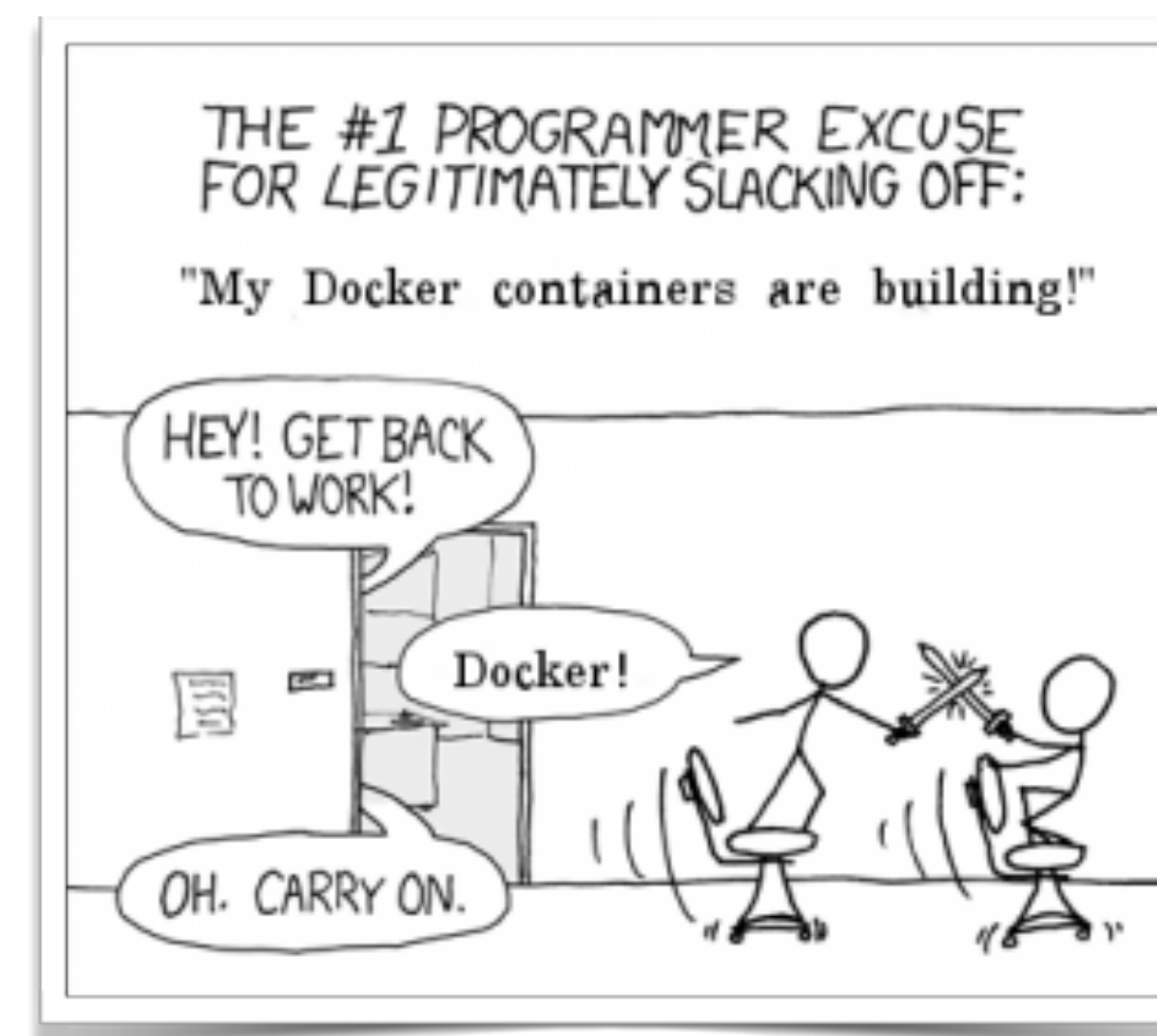finds nodes to run pods on

# Kubernetes

- Container lifecycle management

- Declarative API + control loops

- Robust, flexible, scalable

- Extensible



kubernetes.io

# Roles and responsibilities

- infrastructure admin

- namespace admin

- developer

# How can I customize Kubernetes?

# Customization options in principle

- in-tree (upstream) via SIG or direct PR

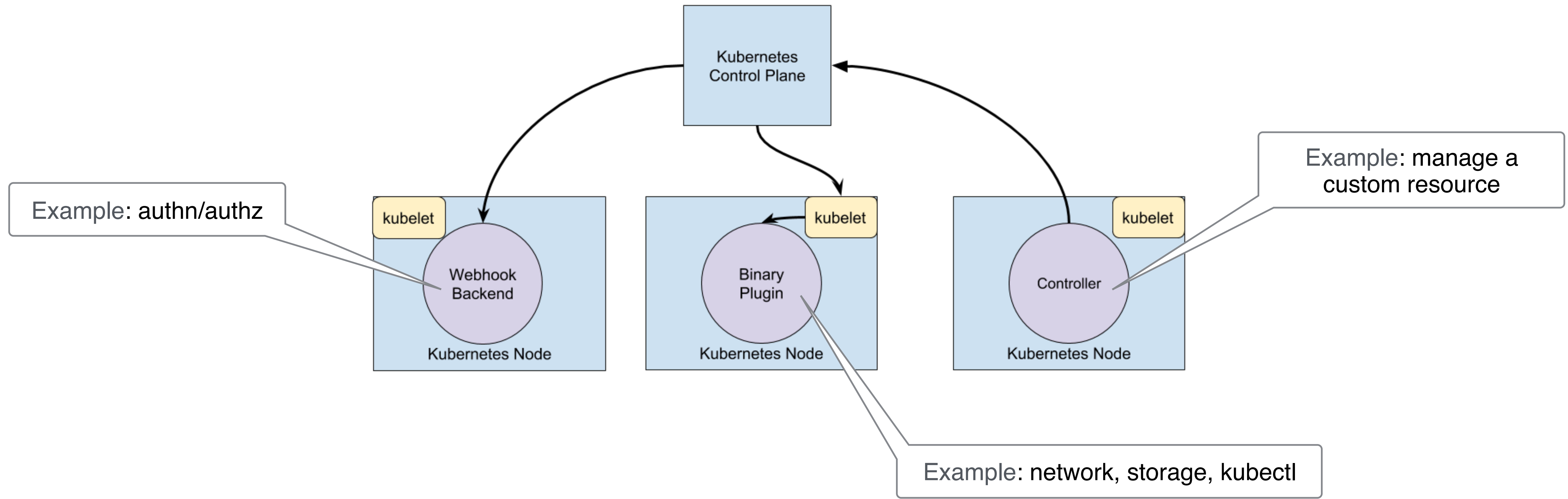- maintain your own fork

- built-in customization approaches



Hit me up on Twitter: @mhausenblas

# Customization approaches

infrastructure (I)

API (A)

- configuration files and flags (`kubelet`, `kube-apiserver`, etc.)
- extension points
  - cloud providers (I)
  - `kubelet` (plugins for network/devices/storage and container runtimes) (I)
  - `kubectl` plugins (I)
  - access extensions in the API server (A)
  - custom resources/controllers (A)
  - extension API servers (A)
  - scheduler extensions (I)

redhat.

# Extension patterns



Kubernetes
Control Plane

kubelet

Webhook
Backend

Kubernetes Node

kubelet

Binary
Plugin

Kubernetes Node

kubelet

Controller

Kubernetes Node

Example: authn/authz

Example: manage a
custom resource

Example: network, storage, kubectl

redhat.

# Cloud providers

- in-tree libraries/controller manager

- interfaces for things like:

    - load balancers

    - network routes

    - nodes/VMs

5 results for repositories matching **cloud-pro**               ❌ Clear filter

**cloud-provider-openstack**
● Go   ★ 63   ⑂ 50   ⚖ Apache-2.0   Updated 2 days ago

**cloud-provider-azure**
● Go   ★ 10   ⑂ 13   ⚖ Apache-2.0   Updated an hour ago

**cloud-provider-gcp**
● Go   ★ 4   ⑂ 6   ⚖ Apache-2.0   Updated 8 hours ago

**cloud-provider-vsphere**
● Go   ★ 8   ⑂ 3   ⚖ Apache-2.0   3 issues need help   Updated 2 days ago

**cloud-provider-aws**
● Go   ★ 7   ⑂ 5   ⚖ Apache-2.0   Updated 7 days ago

**github.com/kubernetes**

redhat.

# `kubelet`: network/device/storage plugins

- Network — standard: CNI

  github.com/containernetworking/cni

  kubernetes.io/docs/concepts/extend-kubernetes/compute-storage-net/network-plugins

- Devices — GPUs, FPGAs, etc.

  kubernetes.io/docs/concepts/extend-kubernetes/compute-storage-net/device-plugins

- Storage — 20+ in-tree, up-and-coming standard: CSI

  kubernetes.io/docs/concepts/storage/volumes/#types-of-volumes

  kubernetes.io/blog/2018/04/10/container-storage-interface-beta

**FEATURE STATE:** `Kubernetes v1.12`   alpha

**FEATURE STATE:** `Kubernetes v1.12`   beta

**FEATURE STATE:** `Kubernetes v1.10`   beta

# `kubelet:` container runtimes

- Container runtime—standard: CRI (since Kubernetes 1.5)

  kubernetes.io/blog/2016/12/container-runtime-interface-cri-in-kubernetes

- Nowadays multiple options:
  - runc
  - containerd
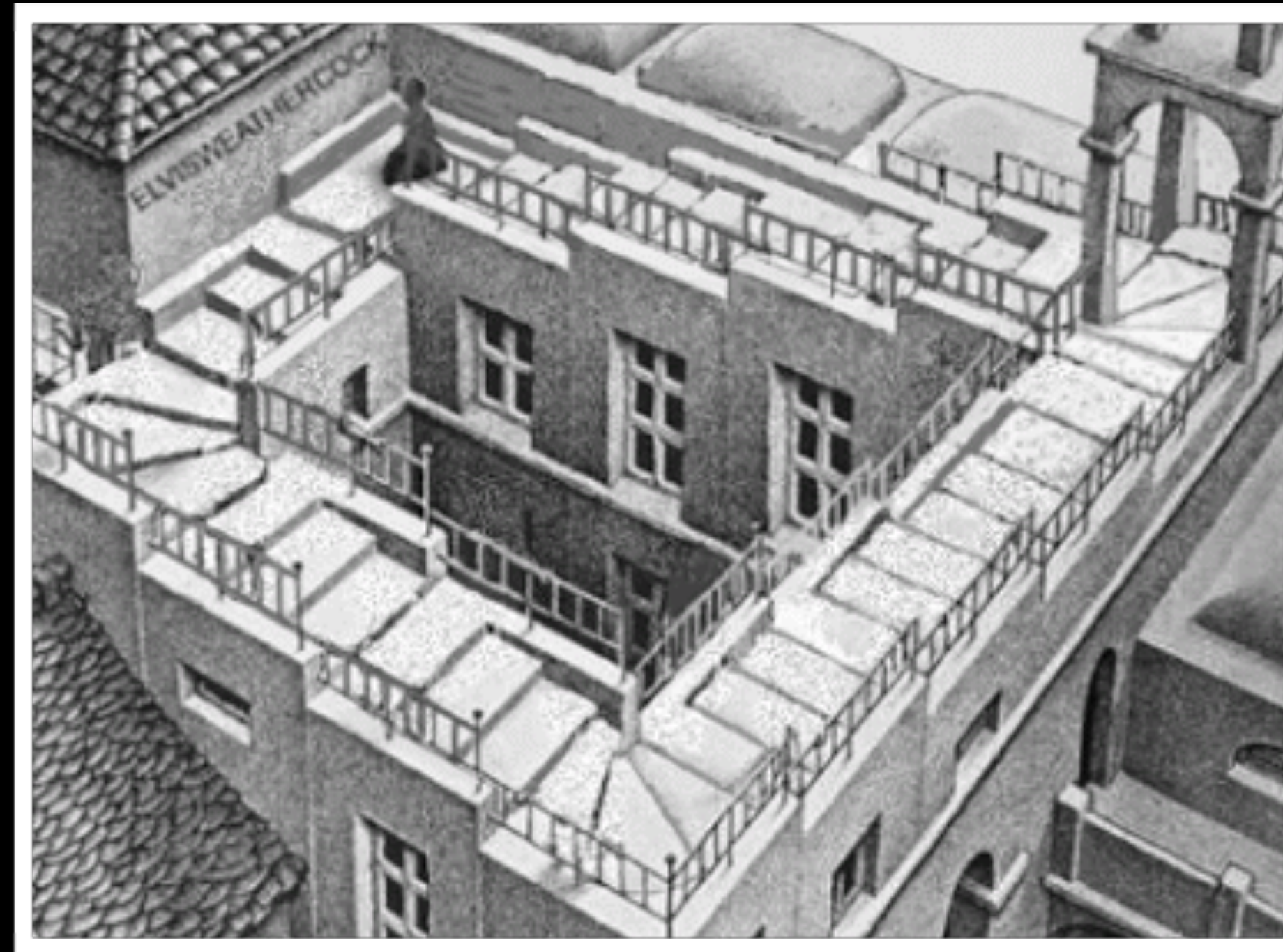  - Kata containers
  - gVisor
  - hyper.sh



cri-o.io

# `kubectl` plugins

- Extend the set of commands

  FEATURE STATE: Kubernetes v1.12    alpha

  kubernetes.io/docs/tasks/extend-kubectl/kubectl-plugins

- Write in any programming language (note: these are binary extensions)

- Examples: context control, service catalog, user verification

as simple plugin in action: `kubectl inspect`

# Extending the Kubernetes API

# Quick control plane refresher



control plane

**API server**
HTTP API, central control point

**controller manager**
in-tree controllers such as deployments

**etcd**
distributed key-value datastore,
single source of truth for state

**scheduler**
finds nodes to run pods on

redhat.

# The life of an API request



Flow diagram is based on Extensible Admission is Beta and Kubernetes deep dive: API Server – part 1.

Hit me up on Twitter: @mhausenblas

# What are (in-tree) core resources?

| Catalog | Kind | 1.5 | | 1.6 | | 1.7 | | 1.8 | | 1.9 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Group | Version | Group | Version | Group | Version | Group | Version | Group | Version |
| Workloads | Container | Core | v1 | Core | v1 | Core | v1 | Core | v1 | Core | v1 |
| | CronJob | Batch | v2alpha1 | Batch | v2alpha1 | Batch | v2alpha1 | Batch | v1beta1 | Batch | v1beta1 |
| | DaemonSet | Extensions | v1beta1 | Extensions | v1beta1 | Extensions | v1beta1 | Apps | v1beta2 | Apps | v1 |
| | Deployment | Extensions | v1beta1 | Apps | v1beta1 | Apps | v1beta1 | Apps | v1beta2 | Apps | v1 |
| | Job | Batch | v1 | Batch | v1 | Batch | v1 | Batch | v1 | Batch | v1 |
| | Pod | Core | v1 | Core | v1 | Core | v1 | Core | v1 | Core | v1 |
| | ReplicaSet | Extensions | v1beta1 | Extensions | v1beta1 | Extensions | v1beta1 | Apps | v1beta2 | Apps | v1 |
| | ReplicationController | Core | v1 | Core | v1 | Core | v1 | Core | v1 | Core | v1 |
| | StatefulSet | Apps | v1beta1 | Apps | v1beta1 | Apps | v1beta1 | Apps | v1beta2 | Apps | v1 |
| Cluster | ClusterRole | RbacAuthorization | v1alpha1 | RBAC | v1beta1 | RBAC | v1beta1 | RBAC | v1 | RBAC | v1 |
| | ClusterRoleBinding | RbacAuthorization | v1alpha1 | RBAC | v1beta1 | RBAC | v1beta1 | RBAC | v1 | RBAC | v1 |
| | ComponentStatus | Core | v1 | Core | v1 | Core | v1 | Core | v1 | Core | v1 |
| | LocalSubjectAccessReview | Authorization | v1beta1 | Authorization | v1 | Authorization | v1 | Authorization | v1 | Authorization | v1 |
| | Namespace | Core | v1 | Core | v1 | Core | v1 | Core | v1 | Core | v1 |
| | Node | Core | v1 | Core | v1 | Core | v1 | Core | v1 | Core | v1 |
| | PersistentVolume | Core | v1 | Core | v1 | Core | v1 | Core | v1 | Core | v1 |
| | ResourceQuota | Core | v1 | Core | v1 | Core | v1 | Core | v1 | Core | v1 |
| | Role | RbacAuthorization | v1alpha1 | RBAC | v1beta1 | RBAC | v1beta1 | RBAC | v1 | RBAC | v1 |
| | RoleBinding | RbacAuthorization | v1alpha1 | RBAC | v1beta1 | RBAC | v1beta1 | RBAC | v1 | RBAC | v1 |
| | SelfSubjectAccessReview | Authorization | v1beta1 | Authorization | v1 | Authorization | v1 | Authorization | v1 | Authorization | v1 |
| | SelfSubjectRulesReview | | | | | | | Authorization | v1 | Authorization | v1 |
| | ServiceAccount | Core | v1 | Core | v1 | Core | v1 | Core | v1 | Core | v1 |
| | SubjectAccessReview | Authorization | v1beta1 | Authorization | v1 | Authorization | v1 | Authorization | v1 | Authorization | v1 |
| | TokenReview | Authentication | v1beta1 | Authorization | v1 | Authorization | v1 | Authorization | v1 | Authorization | v1 |
| | NetworkPolicy | Extensions | v1beta1 | Extensions | v1beta1 | Networking | v1 | Networking | v1 | Networking | v1 |

kubernetes.io/docs/reference/generated/kubernetes-api/v1.12/

redhat.

# Access extensions in the API server

- Admission controllers (in-tree, via configuration of the API server)

  https://kubernetes.io/docs/reference/access-authn-authz/admission-controllers/

- Dynamic Admission Control

  https://kubernetes.io/docs/reference/access-authn-authz/extensible-admission-controllers/

  - Admission Webhooks (beta)

  - *Initializers (alpha)*

Hit me up on Twitter: @mhausenblas

# Custom resources

- Support for "known" resources beyond core resources

  kubernetes.io/docs/concepts/extend-kubernetes/api-extension/custom-resources

  blog.openshift.com/kubernetes-deep-dive-api-server-part-3a

- Use the API server to manage custom resources in `etcd` for you

- Custom resource definition (CRD) and instances

- Use the CLI to interact with custom resources in the usual way,

  for example: `kubectl get mycustomresource`

# Custom resource—example

```
 1   apiVersion: apiextensions.k8s.io/v1beta1
 2   kind:       CustomResourceDefinition
 3   metadata:
 4     name:       databases.example.com
 5   spec:
 6     group:     example.com
 7     version:   v1
 8     names:
 9       kind:    Database
10       plural:  databases
11       scope:   Namespaced
```

# Custom controller

- Implement control loops beyond what thee (in-tree)
  controller manager supports
- Custom controller
  - dealing with core resources
    github.com/kelseyhightower/secrets-controller
  - dealing with custom resources (aka operator)
    github.com/kubernetes/sample-controller

# Custom resources and controllers

| | resource | | controller | |
|---|---|---|---|---|
| | core | custom | in-tree | custom |
| Kubernetes control plane | X | | X | |
| simple controller | X | | | X |
| operator | | X | | X |

Hit me up on Twitter: @mhausenblas

# Operators

# Operators

*operator = custom resource + custom controller*

- Motivation: application lifecycle management

- Use one of 30+ available operators or write your own with:

  - Kubebuilder

  - Kubernetes Operator Kit

  - kutil

  - Metacontroller

  - Operator SDK

redhat.

# Operator use cases

- zero-downtime upgrades of the app the operator supervises

- workflow automations

- policy enforcement

- managing stateful workloads

  - resizing of followers in a distributed datastore

  - backup & restore of a database

  - re-balancing of a distributed message queue

# Operator examples

- etcd

- Prometheus

- Postgres

- Vitess MySQL

- MongoDB

- Couchbase

- Kafka

a simple operator in action: `NoDefaultsPolicy`

# $ operator-sdk new nodefpol-operator

```
$ operator-sdk add api --api-version=nodefpol.k8space.io/v1alpha1 --kind=NoDefaultsPolicy
```

```
$ operator-sdk add controller --api-version=nodefpol.k8space.io/v1alpha1 --kind=NoDefaultsPolicy
```

```
$ kubectl -n ndp-demo apply -f deploy/crds/nodefpol_v1alpha1_nodefaultspolicy_crd.yaml
$ OPERATOR_NAME=nodefpol-operator operator-sdk up local --namespace "ndp-demo"
```

grep '//TODO(user)'

# Extension API servers

- Full control but a lot of effort and responsibility
  kubernetes.io/docs/tasks/access-kubernetes-api/setup-extension-api-server

- Typically more LOC than an controller or operator

- You might end up to manage storage in `etcd` yourself

- And beyond: the Open Service Broker API and the service catalog
  kubernetes.io/docs/concepts/extend-kubernetes/service-catalog

  openservicebrokerapi.org

# Scheduler extensions

A scheduler selects a node to run your pods on, based on resource requirements, QoS, affinity, etc.

jvns.ca/blog/2017/07/27/how-does-the-kubernetes-scheduler-work

- You can modify policies or run multiple schedulers (with pod opt-in)

  kubernetes.io/docs/tasks/administer-cluster/configure-multiple-schedulers

  embano1.github.io/post/sched-reconcile

- You can use a Webhook

  github.com/kubernetes/community/blob/master/contributors/design-proposals/scheduling/scheduler_extender.md

# Other stuff you can customize in Kubernetes

- Monitoring & alerting (Prometheus/Grafana), logging (ELK/EFK stack)

- Secret management (encryption at rest, Vault)

- Ingress
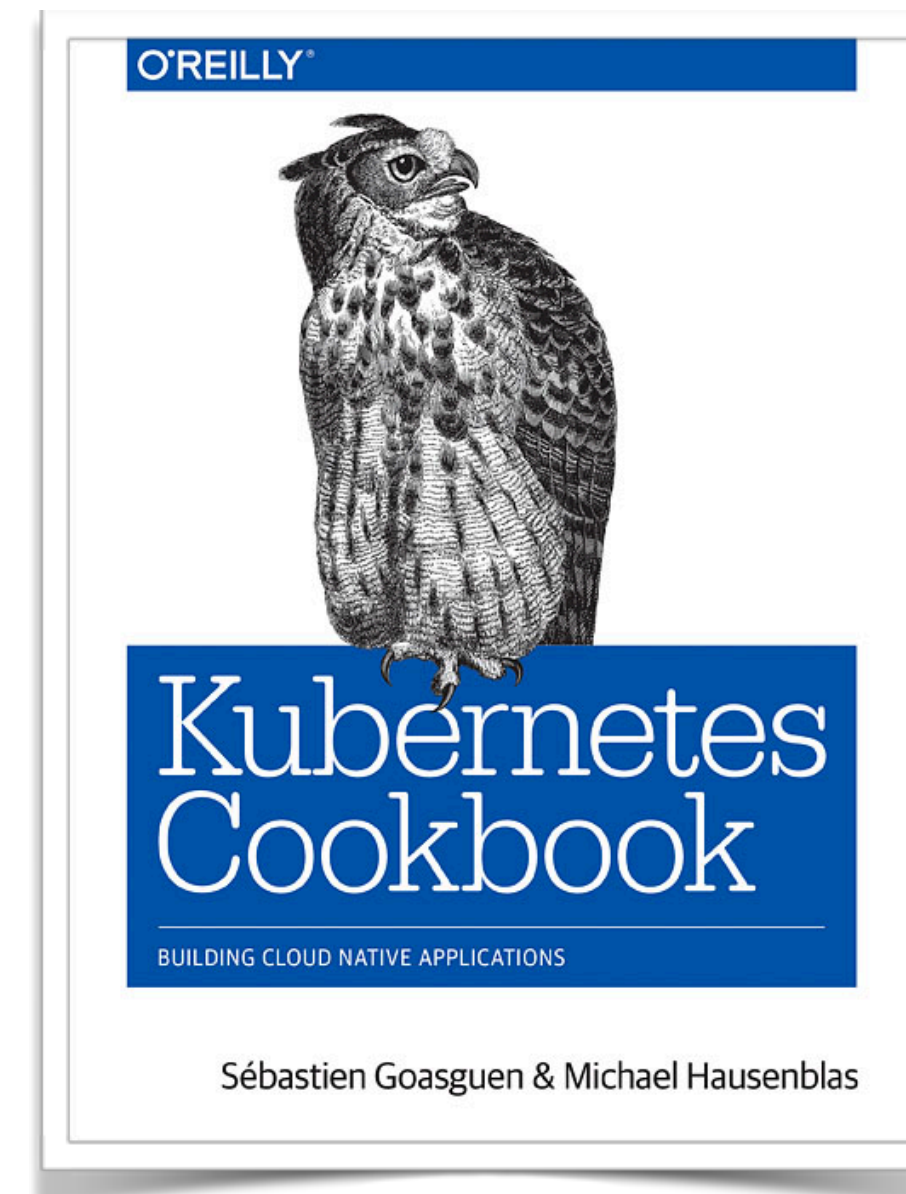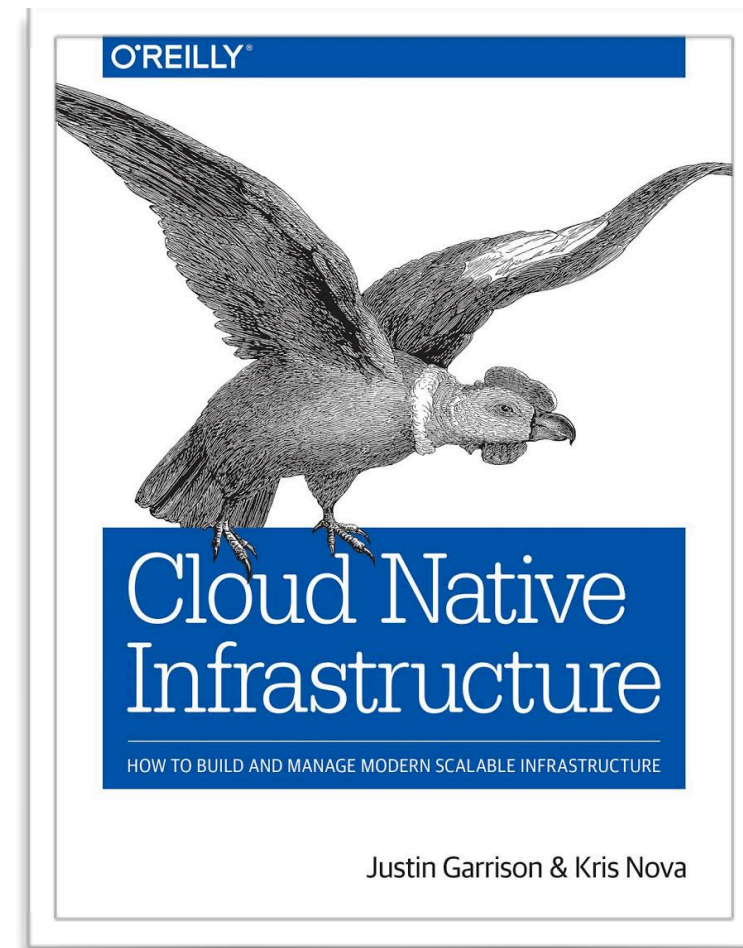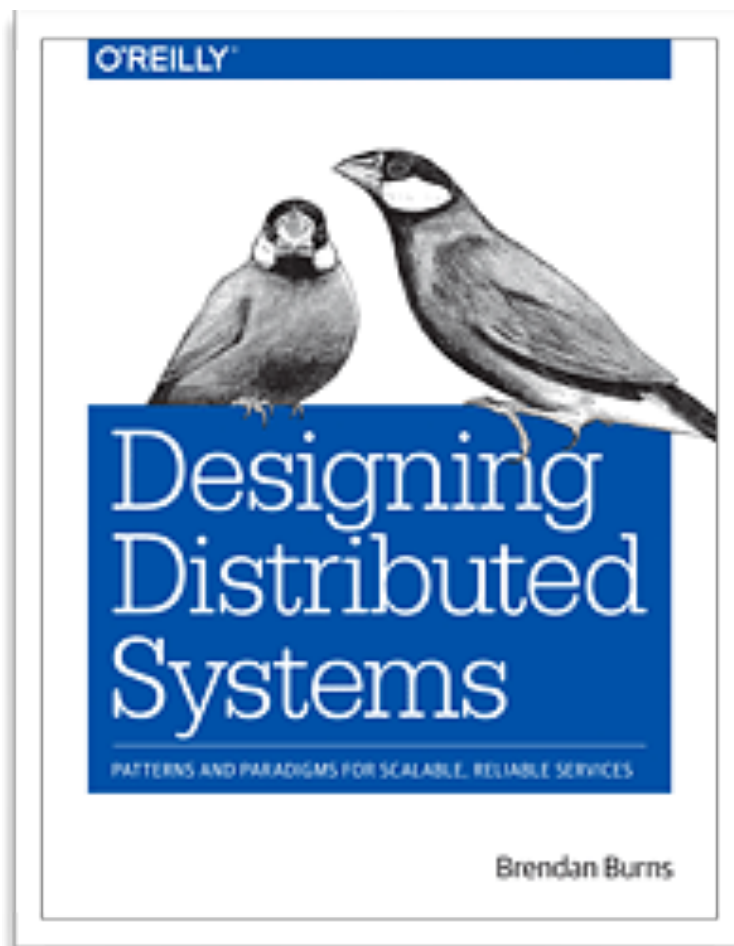  kubernetes.io/docs/concepts/services-networking/ingress

- DNS
  kubernetes.io/docs/tasks/administer-cluster/dns-custom-nameservers

- kube-proxy
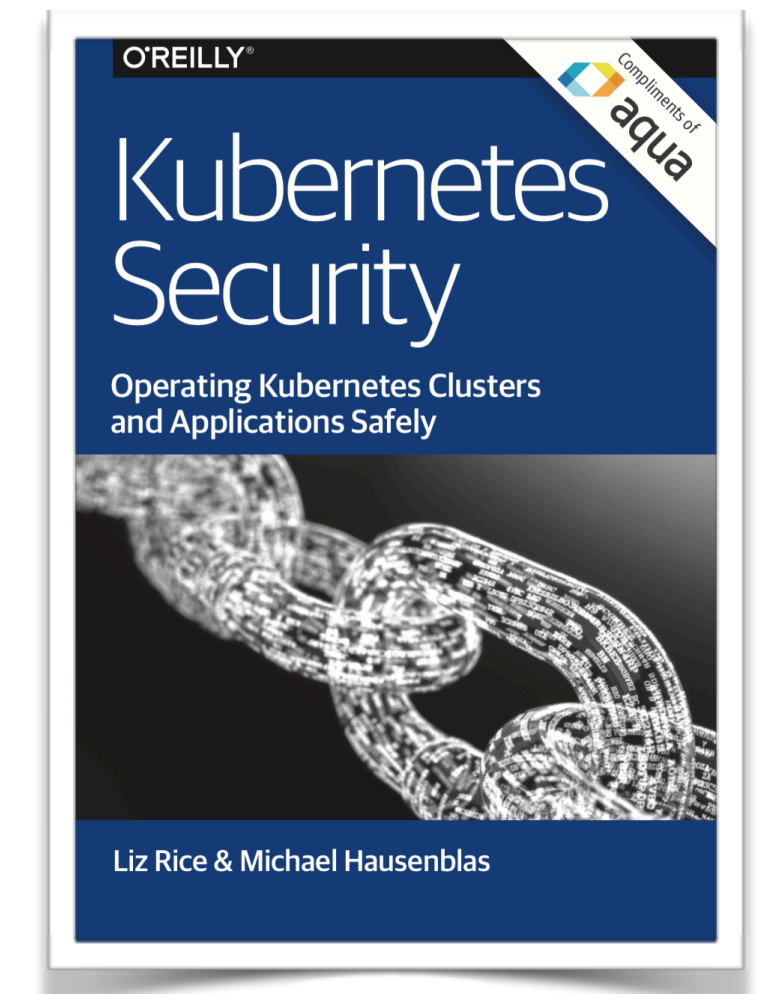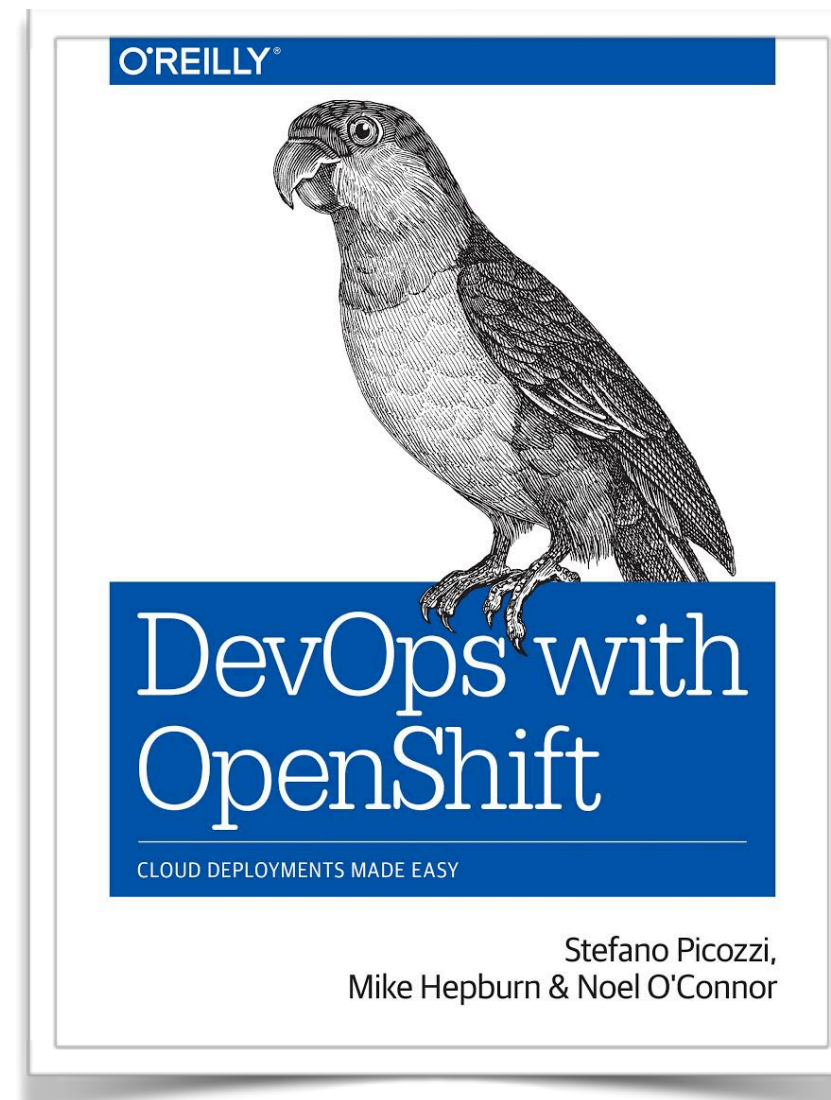  kubernetes.io/blog/2018/07/09/ipvs-based-in-cluster-load-balancing-deep-dive

Hit me up on Twitter: @mhausenblas

redhat.

# Resources

# Articles and slide decks

- *Tim Hockin*—Kubernetes Extensibility

    speakerdeck.com/thockin/kubernetes-extensibility

- *Jonathan Berkhahn & Carolyn Van Slyck*—Kubectl Plugins 101

    kccnceu18.sched.com/event/DqwJ/kubectl-plugins-101-jonathan-berkhahn-ibm-carolyn-van-slyck-microsoft-intermediate-skill-level-slides-attached

- *Adrien Trouillaud*—Kubernetes Custom Resource, Controller & Operator Development Tools

    admiralty.io/kubernetes-custom-resource-controller-and-operator-development-tools.html

- *Toader Sebastian*—A complete guide to Kubernetes Operator SDK

    banzaicloud.com/blog/operator-sdk/

- *Rob Szumski*—Building an Kubernetes Operator for Prometheus and Thanos

    robszumski.com/building-an-operator/

# Repos, examples, tooling

- github.com/kubernetes/kubectl/tree/master/pkg/pluginutils

- github.com/carolynvs/kubectl-flags-plugin

- github.com/jordanwilson230/kubectl-plugins

- github.com/kelseyhightower/denyenv-validating-admission-webhook

- github.com/kubernetes-sigs/controller-tools

- github.com/kubernetes-sigs/kubebuilder

- metacontroller.app

- github.com/yaronha/kube-crd

- github.com/operator-framework/operator-sdk

- github.com/operator-framework/awesome-operators

- reactiveops.github.io/rbac-manager

redhat.

# Kubernetes docs and blog posts

- kubernetes.io/docs/concepts/extend-kubernetes/extend-cluster/

    - kubernetes.io/docs/concepts/extend-kubernetes/api-extension/custom-resources/

    - kubernetes.io/docs/concepts/extend-kubernetes/api-extension/apiserver-aggregation/

    - kubernetes.io/docs/tasks/access-kubernetes-api/setup-extension-api-server/

- kubernetes.io/docs/tasks/extend-kubectl/kubectl-plugins/

- kubernetes.io/docs/reference/access-authn-authz/webhook/

- kubernetes.io/docs/setup/scratch/#cloud-provider

- kubernetes.io/blog/2018/01/extensible-admission-is-beta/

# Videos

- *Tim Hockin & Michael Rubin*—Kubernetes Distributions and 'Kernels'
https://www.youtube.com/watch?v=fXBjA2hH-CQ

- *Stefan Schimanski*:

  - Kubernetes as a API driven platform, Reykjavík Kubernetes Meetup
https://www.youtube.com/watch?v=BiE7oKeEzDU

  - SIG API Machinery Deep Dive
https://www.youtube.com/watch?v=XsFH7OEIIvI

- *James Munnelly*—Extending the Kubernetes API: What the Docs Don't Tell You
https://www.youtube.com/watch?v=PYLFZVv68IM