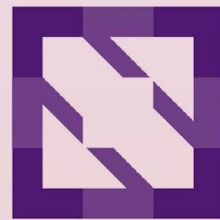




KubeCon



CloudNativeCon

North America 2023





KubeCon



CloudNativeCon

North America 2023

SIG API Machinery

Leila Jalali, Google

Stefan Schimanski, Upbound

Where are we on the journey?



KubeCon



CloudNativeCon

North America 2023

2014
Kubernetes introduced

2015
Kubernetes v1.0
Birth of CNCF

2018
K8s graduated
CNCF

2022
multi-cloud and
hybrid cloud

2023
96% of organizations are
either using or evaluating
Kubernetes



● **kubernetes api server**
Search term



Interest over time

SIG Overview

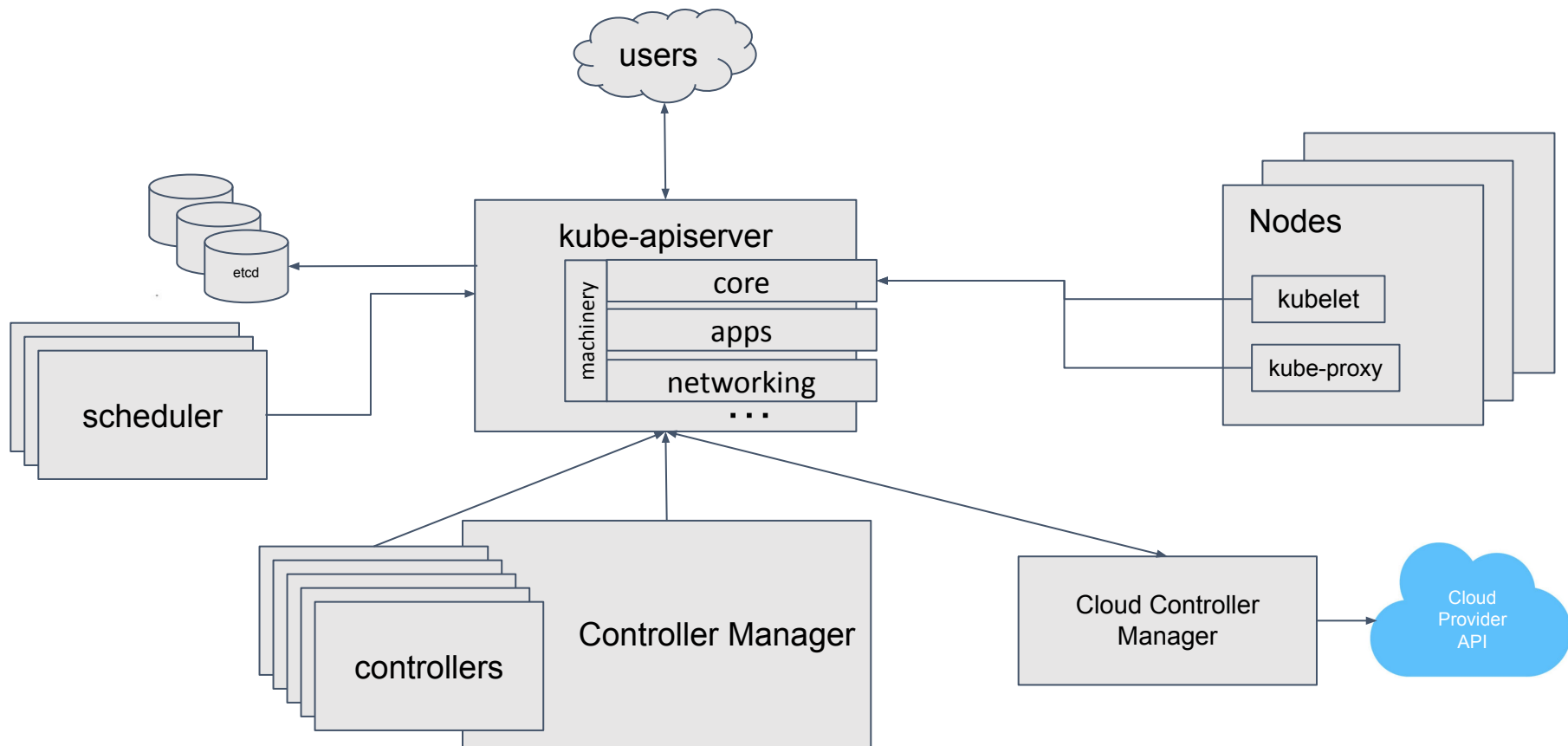


KubeCon



CloudNativeCon

North America 2023



SIG Overview

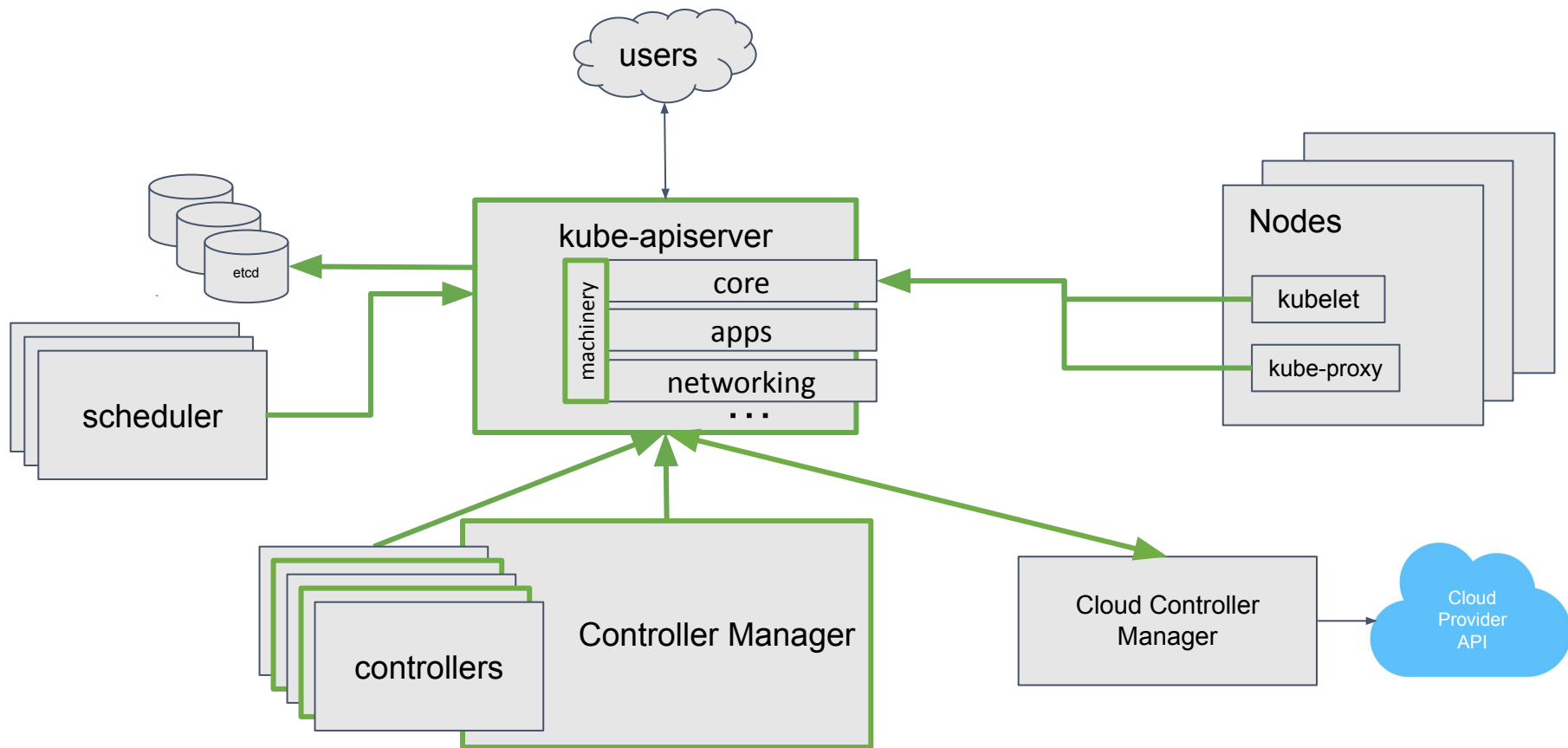


KubeCon



CloudNativeCon

North America 2023



SIG API Machinery is responsible for the development and enhancement of Kubernetes cluster control plane. The scope covers API server, persistence layer (etcd), controller manager, cloud controller manager, CustomResourceDefinition and webhooks. ([SIG API Machinery Charter](#))

- What does the name API Machinery stand for?
- Is API Machinery != All Kubernetes APIs?
- API Machinery == the machinery used by different Kubernetes APIs to interact with the Kubernetes cluster, to be exposed and actuated, and the mechanics to publish, process, and extend them.



What do we own?



KubeCon



CloudNativeCon

North America 2023

- We provide stable core APIs to establish the permanent foundation for the rest of the K8s components to interact with.
- We enable and support other SIGs to be successful through the usage of our machinery

In scope

All aspects of

- How to read, modify, delete objects, including parsing, conversion, defaulting and validation
- Describing and extending the system: OpenAPI, Discovery, CRDs, Webhooks, client/informer libraries
- Maintaining a healthy system: controller-manager, garbage collection, namespace lifecycle
- The persistence layer (etcd, scalability SIGs)

Out of scope

- All the individual Kubernetes APIs

Why is it so important and complex?



KubeCon

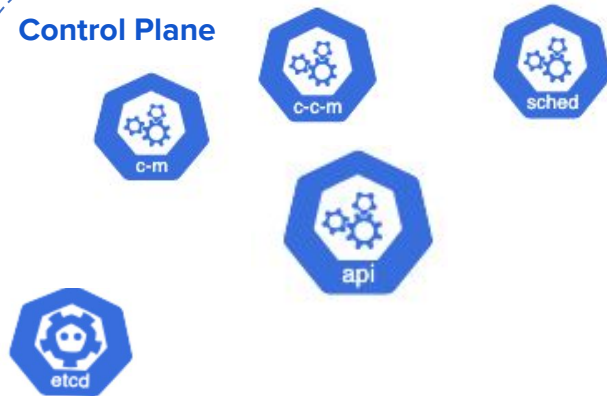


CloudNativeCon

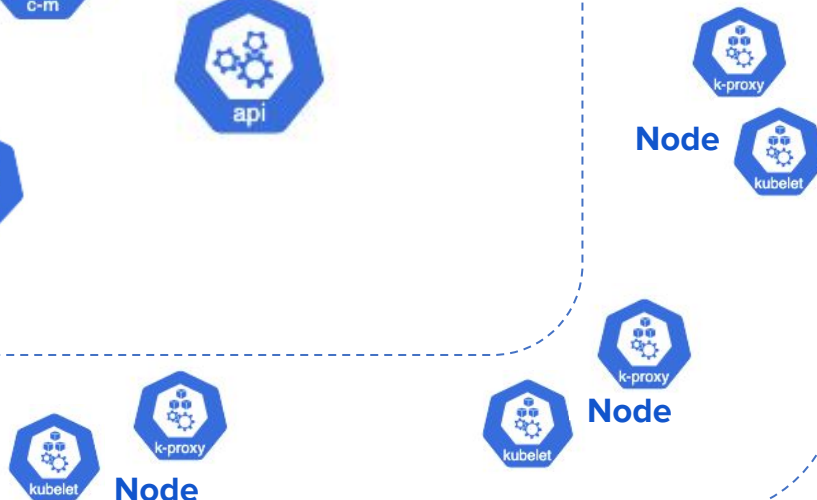
North America 2023

Kubernetes Cluster

Control Plane



Node



Why is it so important and complex?



KubeCon

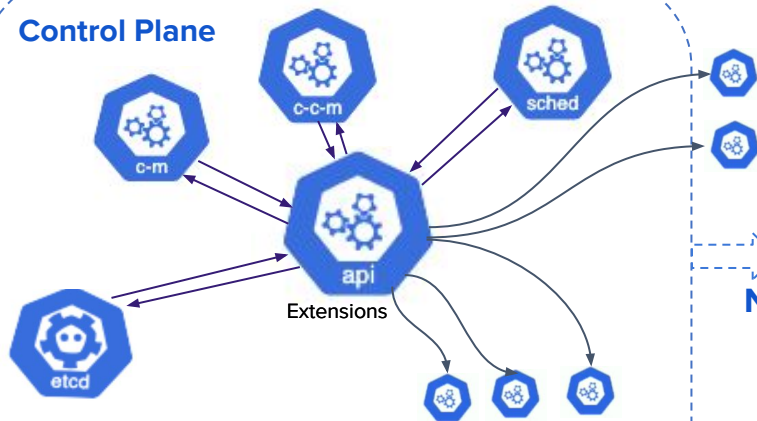


CloudNativeCon

North America 2023

Kubernetes Cluster

Control Plane



Node



Why is it so important and complex?

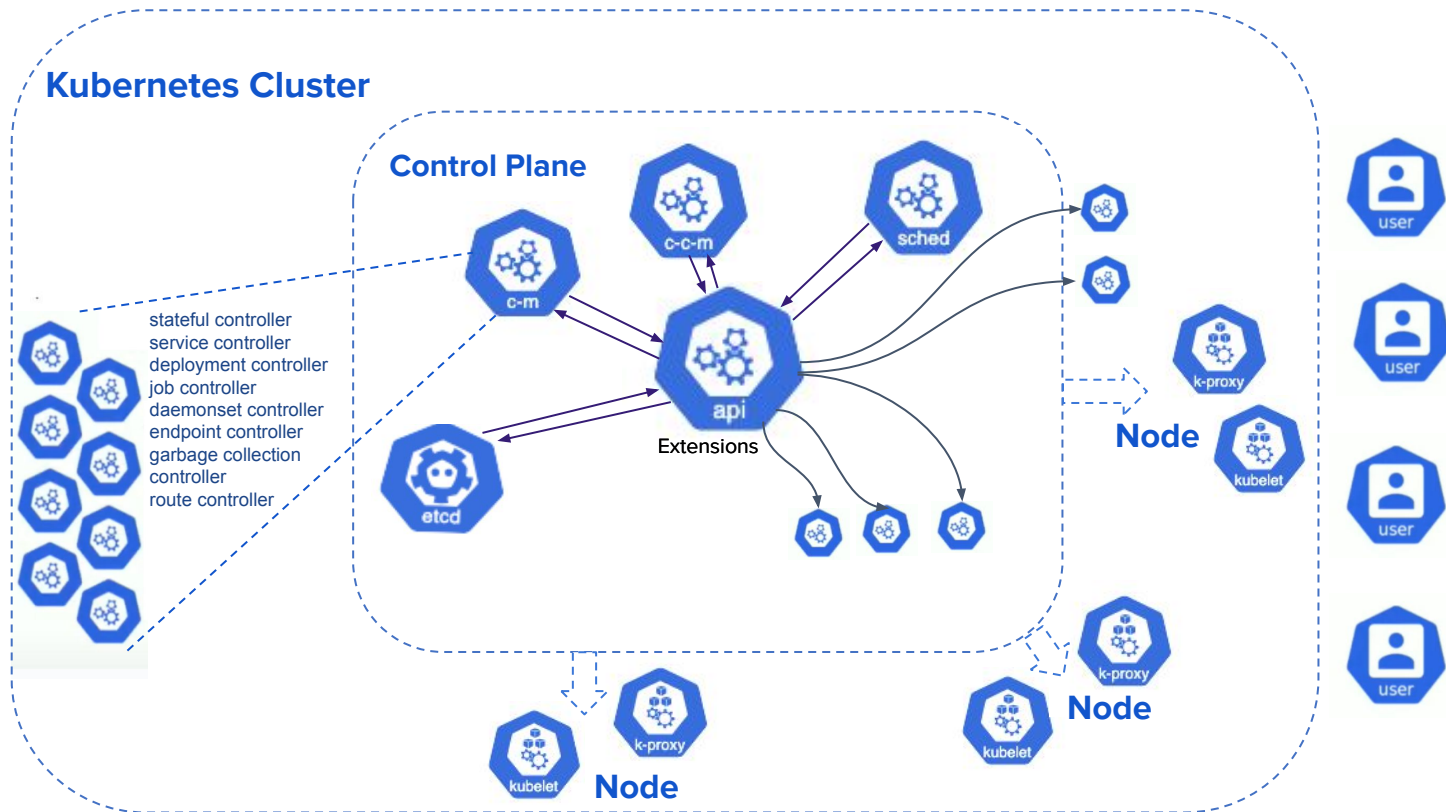


KubeCon



CloudNativeCon

North America 2023



Why is it so important and complex?

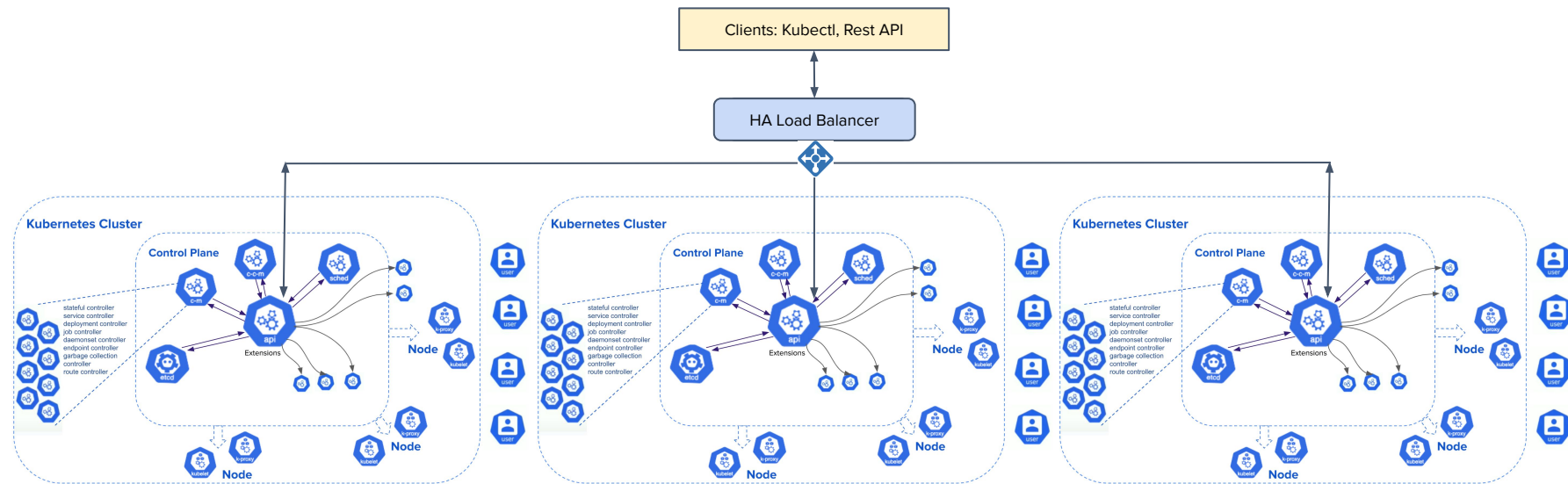


KubeCon



CloudNativeCon

North America 2023



Why is it so important and complex?

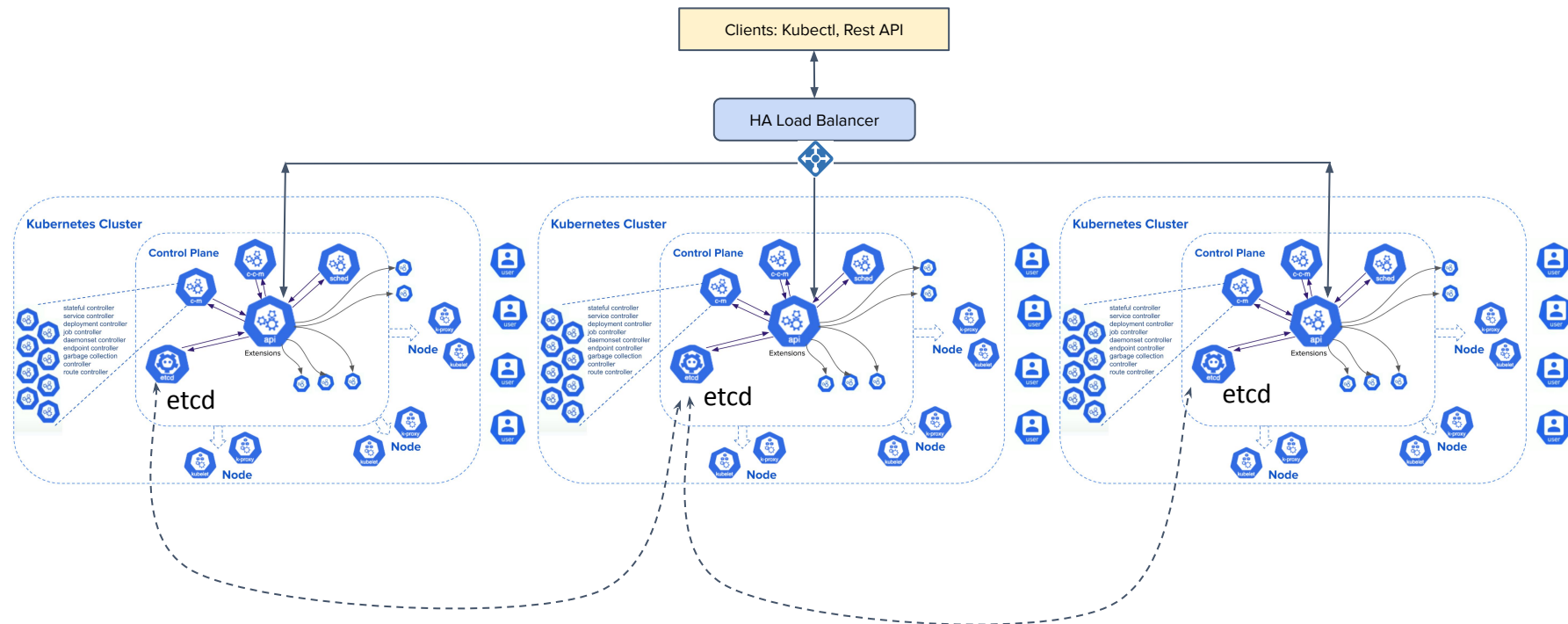


KubeCon



CloudNativeCon

North America 2023



Why is it so important and complex?

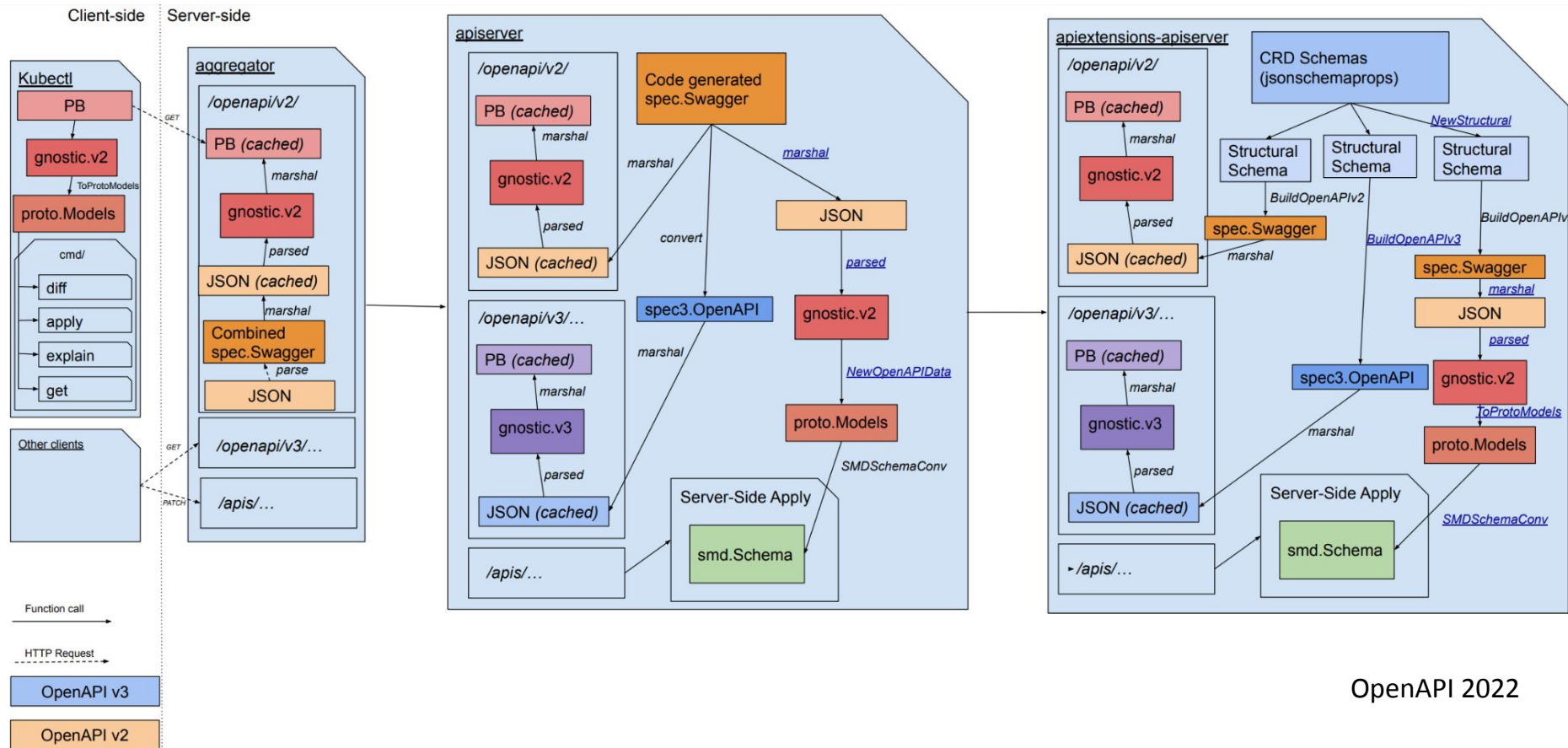


KubeCon



CloudNativeCon

North America 2023



OpenAPI 2022

Why is it so important and complex?

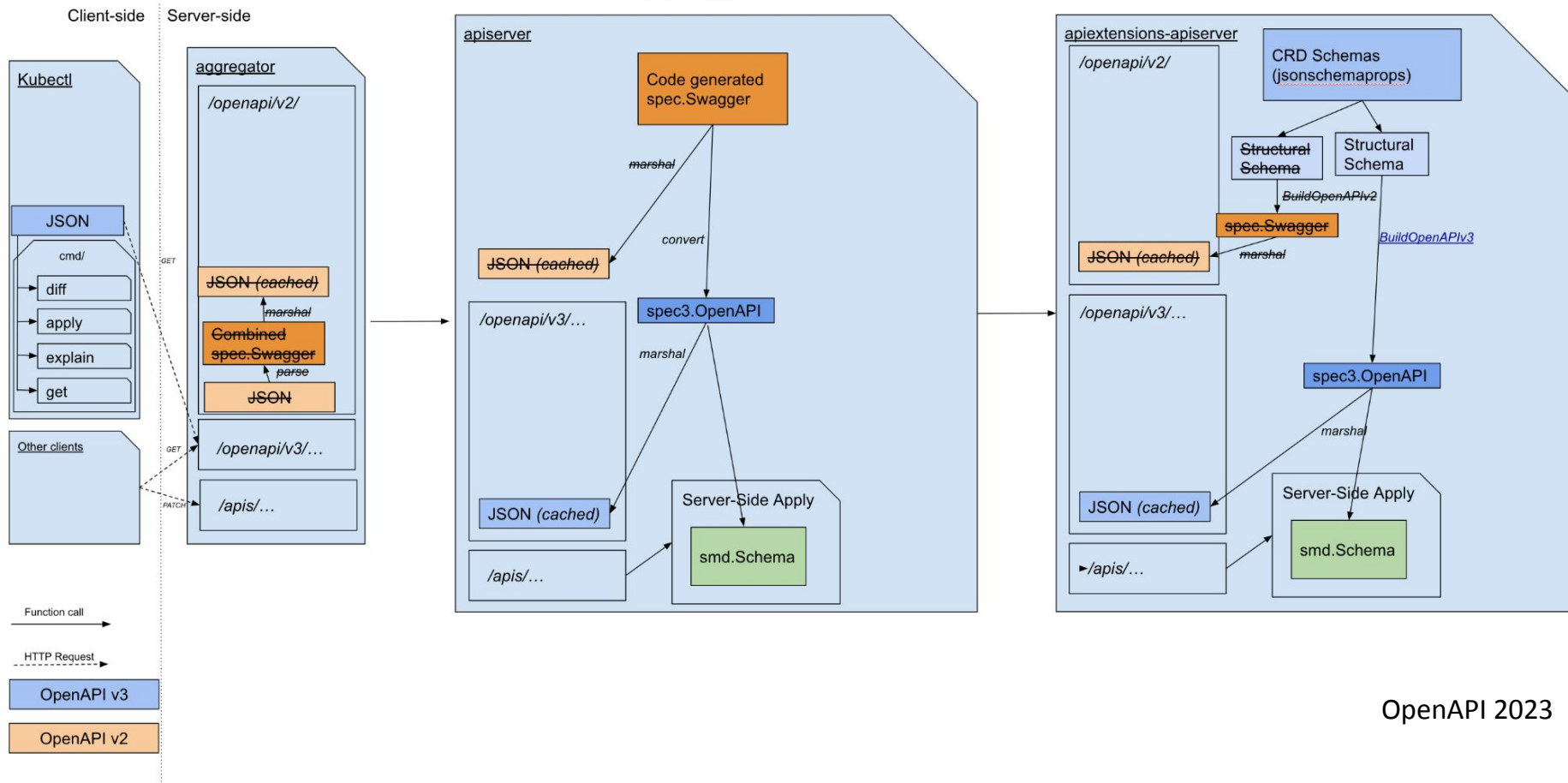


KubeCon



CloudNativeCon

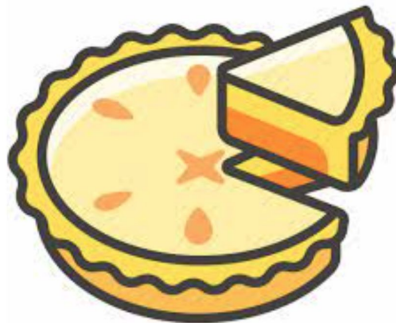
North America 2023



OpenAPI 2023

25-40% code ownership

- 2 approaches with solid proxies:
 - SIG-labelled PRs vs file cataloguing



1.29 Enhancement Tracking



KubeCon



CloudNativeCon

North America 2023

Title	...	Status	...	Stage	...	Type	...
• Priority and Fairness for API Server Requests #1040		Tracked for Enhancements Freeze	▼	Stable	▼	Graduating	▼
• Transition from SPDY to WebSockets #4006		Tracked for Code Freeze	▼	Alpha	▼	Net New	▼
• Support paged LIST queries from the Kubernetes API #365		At Risk for Code Freeze	▼	Stable	▼	Graduating	▼
• Move Storage Version Migrator in-tree #4192		Tracked for Enhancements Freeze	▼	Alpha	▼	Net New	▼
• CRD Validation Expression Language #2876		At Risk for Code Freeze	▼	Stable	▼	Graduating	▼
• CBOR Serializer #4222		Tracked for Enhancements Freeze	▼	Alpha	▼	Net New	▼
• Allow informers for getting a stream of data instead of chunking #3157		Removed from Milestone	▼	Beta	▼	Graduating	▼
• CRD Validation Ratcheting #4008		At Risk for Code Freeze	▼	Beta	▼	Graduating	▼
• Declarative Validation #4153		Tracked for Enhancements Freeze	▼	Alpha	▼	Net New	▼
• CEL for Admission Control #3488		Tracked for Enhancements Freeze	▼	Beta	▼	Graduating	▼

1.29 Enhancement Tracking



KubeCon



CloudNativeCon

North America 2023

Title	Status	Stage	Type
• Priority and Fairness for API Server Requests #1040	Tracked for Enhancements Freeze	Stable	Graduating
• Transition from SPDY to WebSockets #4006	Tracked for Code Freeze	Alpha	Net New
• Support paged LIST queries from the Kubernetes API #365	At Risk for Code Freeze	Stable	Graduating
• Move Storage Version Migrator in-tree #4192	Tracked for Enhancements Freeze	Alpha	Net New
• CRD Validation Expression Language #2876	At Risk for Code Freeze	Stable	Graduating
• CBOR Serializer #4222	Tracked for Enhancements Freeze	Alpha	Net New
• Allow informers for getting a stream of data instead of chunking #3157	Removed from Milestone	Beta	Graduating
• CRD Validation Ratcheting #4008	At Risk for Code Freeze	Beta	Graduating
• Declarative Validation #4153	Tracked for Enhancements Freeze	Alpha	Net New
• CEL for Admission Control #3488	Tracked for Enhancements Freeze	Beta	Graduating

1.29 Enhancement Tracking



KubeCon



CloudNativeCon

North America 2023

Title	...	Status	...	Stage	...	Type	...
Priority and Fairness for API Server Requests #1040		Tracked for Enhancements Freeze	▼	Stable	▼	Graduating	▼
Transition from SPDY to WebSockets #4006		Tracked for Code Freeze	▼	Alpha	▼	Net New	▼
Support paged LIST queries from the Kubernetes API #365		At Risk for Code Freeze	▼	Stable	▼	Graduating	▼
Move Storage Version Migrator in-tree #4192		Tracked for Enhancements Freeze	▼	Alpha	▼	Net New	▼
CRD Validation Expression Language #2876		At Risk for Code Freeze	▼	Stable	▼	Graduating	▼
CBOR Serializer #4222		Tracked for Enhancements Freeze	▼	Alpha	▼	Net New	▼
Allow informers for getting a stream of data instead of chunking #3157		Removed from Milestone	▼	Beta	▼	Graduating	▼
CRD Validation Ratcheting #4008		At Risk for Code Freeze	▼	Beta	▼	Graduating	▼
Declarative Validation #4153		Tracked for Enhancements Freeze	▼	Alpha	▼	Net New	▼
CEL for Admission Control #3488		Tracked for Enhancements Freeze	▼	Beta	▼	Graduating	▼

CRDs, Schema, Declarative APIs

1.29 Enhancement Tracking



KubeCon



CloudNativeCon

North America 2023

Title	Status	Stage	Type
• Priority and Fairness for API Server Requests #1040	Tracked for Enhancements Freeze	Stable	Graduating
• Transition from SPDY to WebSockets #4006	Tracked for Code Freeze	Alpha	Net New
• Support paged LIST queries from the Kubernetes API #365	At Risk for Code Freeze	Stable	Graduating
• Move Storage Version Migrator in-tree #4192	Tracked for Enhancements Freeze	Alpha	Net New
• CRD Validation Expression Language #2876	At Risk for Code Freeze	Stable	Graduating
• CBOR Serializer #4222	Tracked for Enhancements Freeze	Alpha	Net New
• Allow informers for getting a stream of data instead of chunking #3157	Removed from Milestone	Beta	Graduating
• CRD Validation Ratcheting #4008	At Risk for Code Freeze	Beta	Graduating
• Declarative Validation #4153	Tracked for Enhancements Freeze	Alpha	Net New
• CEL for Admission Control #3488	Tracked for Enhancements Freeze	Beta	Graduating



SIG Chairs:

- David Eads ([@deads2k](#)), Red Hat
- Federico Bongiovanni ([@fedebongio](#)), Google

SIG Technical Leads:

- David Eads ([@deads2k](#)), Red Hat
- Joe Betz ([@jpbetz](#)), Google

How to get involved?



KubeCon



CloudNativeCon

North America 2023

- Regular SIG meetings:
 - **SIG Meeting:** 60 min / every 2 weeks (**Wed**, recorded)
 - **PR and Bug triage:** 30 min / twice every week (**Tue & Thu**)
- Regular Working Group meetings (API Expression, Kubebuilder, CEL)
- Upcoming project-based mentorship program

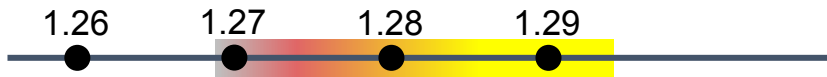
The image is a 3D rendered scene with a soft, pastel color palette of pinks, purples, and oranges. It features various geometric shapes: rectangular blocks of different sizes, a circular platform at the bottom center, and a tall structure on the right made of stacked cubes. Two large spotlights hang from the top, casting a warm, yellowish glow onto the central area. The background is composed of large, flat panels in shades of pink and purple. The overall aesthetic is clean, modern, and artistic.

Some
spotlights on
topics under
the radar.

- <https://github.com/kubernetes/enhancements/issues/4080>

A working kube-based control plane is more than just an apiserver component built on `k/apiserver`. It includes **standard resources** (depending on context namespaces, CRDs, RBAC, secrets, configmaps), and **standard controllers** (think of garbage collection, namespace deletion, etc.). `*kube-apiserver*` today is a bundle of those resources with container orchestration, `*kube-controller-manager*` equally for the corresponding controllers.

Separating the generic parts from container orchestration will allow new use-cases building upon `k/apimachinery` and `k/apiserver`, while **keeping a unified codebase and ecosystem**, and by improving the factoring of `*kube-apiserver*` for easier maintenance due to less complexity by clear layering.



Where in the k/k codebase (eventually)

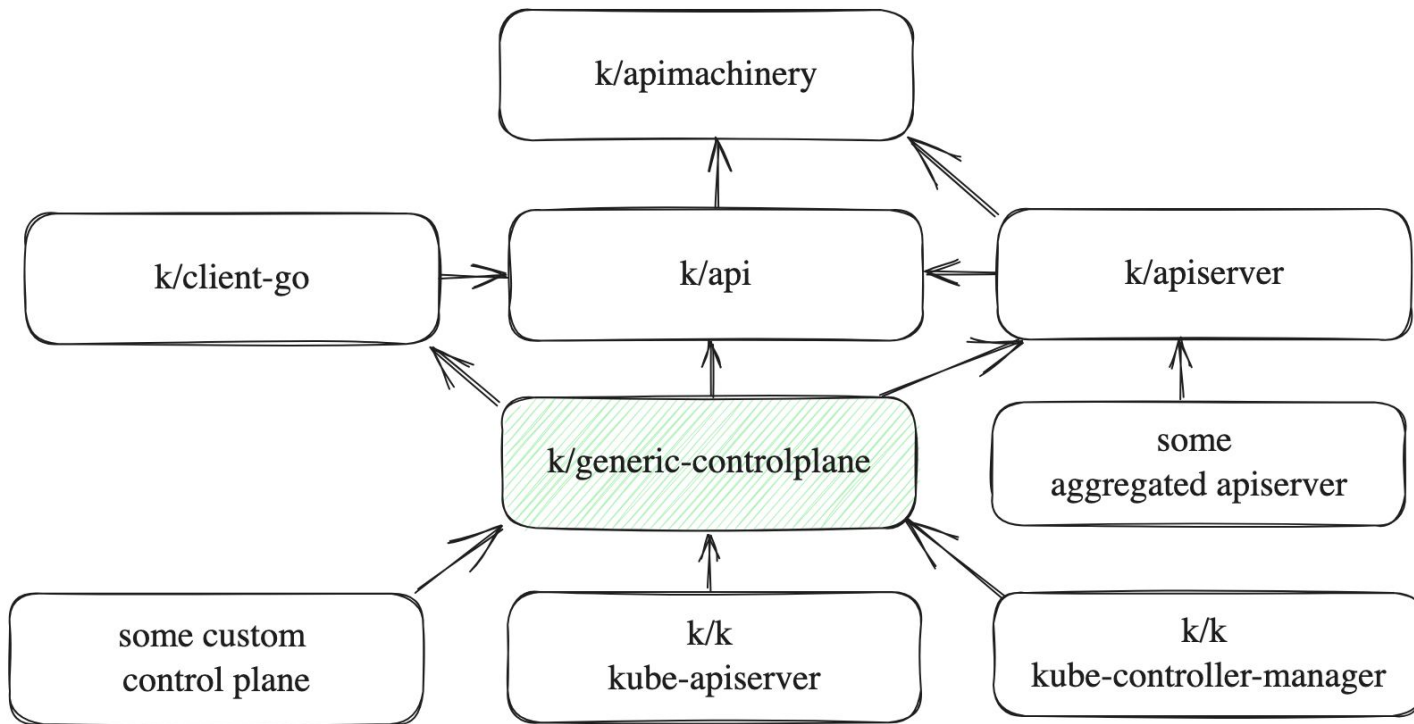


KubeCon



CloudNativeCon

North America 2023



Where in the k/k code base (today)

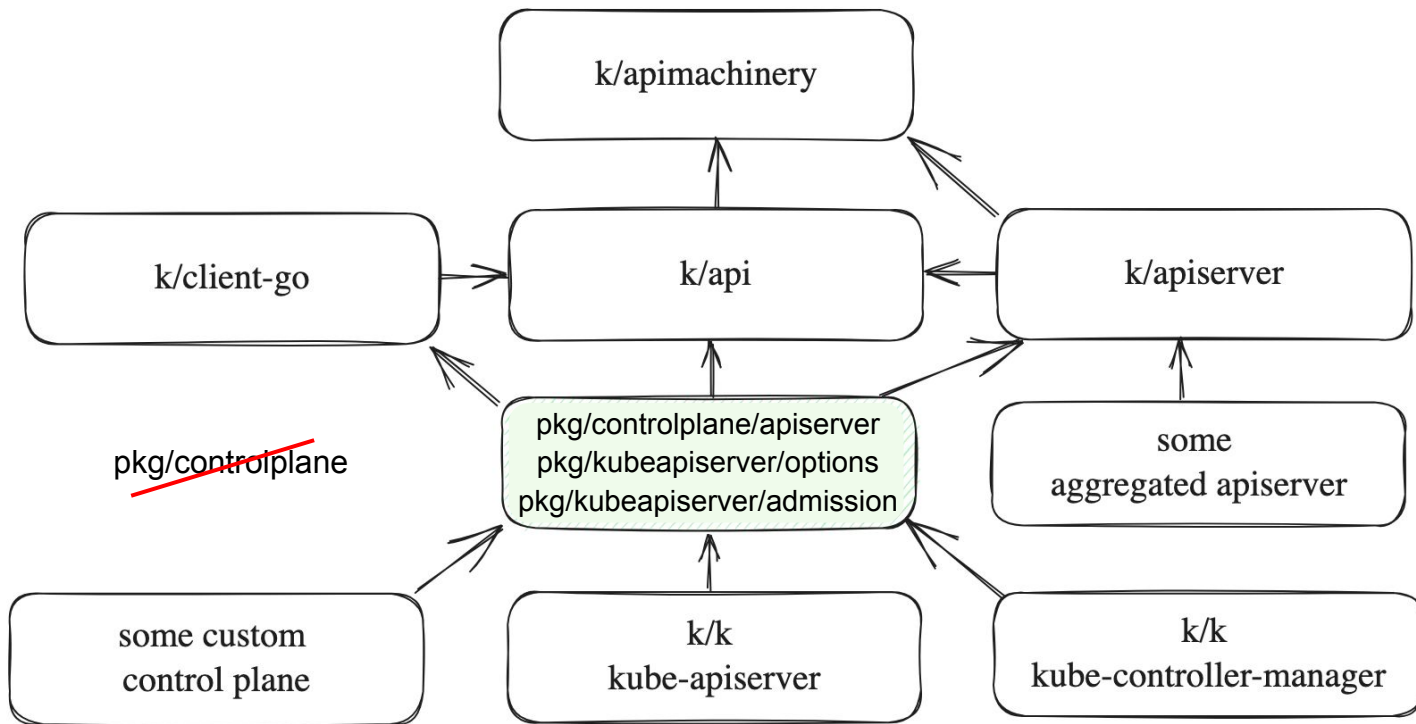


KubeCon



CloudNativeCon

North America 2023



What's in



KubeCon



CloudNativeCon

North America 2023

- CRDs
- namespaces
- secrets optional
- configmaps optional
- RBAC optional
- service accounts optional
- admission webhooks + policies optional
- quota optional
- aggregation, APIServices optional



Garbage
Collection



Namespace
Deletion



Quota

What's in



KubeCon



CloudNativeCon

North America 2023

- CRDs
- namespaces

- secrets



- configmaps



- RBAC



- service accounts



- admission webhooks + policies



- quota



- aggregation, APIServices



Garbage
Collection



Namespace
Deletion



Quota



```
201 }
202 }
203
204 1 usage  Dr. Stefan Schimanski
205 func storageProviders(c controlplaneapiserver.CompletedConfig) ([]controlplaneapiserver.RESTStorageProvider, error) {
206     return []controlplaneapiserver.RESTStorageProvider{
207         withDisabledResources{
208             disabled: map[string][]string{
209                 "v1": {"configmaps", "secrets", "serviceaccounts", "resourcequotas", "resourcequotas/status"},
210             },
211             RESTStorageProvider: &scorerest.GenericConfig{
212                 StorageFactory: c.Extra.StorageFactory,
213                 EventTTL: c.Extra.EventTTL,
214                 LoopbackClientConfig: c.Generic.LoopbackClientConfig,
215                 ServiceAccountIssuer: c.Extra.ServiceAccountIssuer,
216                 ExtendExpiration: c.Extra.ExtendExpiration,
217                 ServiceAccountMaxExpiration: c.Extra.ServiceAccountMaxExpiration,
218                 APIAudiences: c.Generic.Authentication.APIAudiences,
219                 Informers: c.Extra.VersionedInformers,
220             },
221             apiserverinternalrest.StorageProvider: c.Generic.Authentication.Authenticator, c.Generic.Authentication.Authenticator, APIAudiences: c.Generic.Authentication.APIAudiences},
222             FlowControlRestProvider: c.Generic.SharedInformerFactory,
223             coordinationrest.RESTStorageProvider{
224                 eventsrest.RESTStorageProvider{TTL: c.EventTTL},
225             }, nil
226         },
227     }
228
229 2 usages  Dr. Stefan Schimanski
230 type withDisabledResources struct {
231     disabled map[string][]string
232     controlplaneapiserver.RESTStorageProvider
233 }
234
235 storageProviders(c controlplaneapiserver.CompletedConfig) ([]controlplaneapiserver.RESTStorageProvider, error)
```

```
# get https://vbom.ml/util?go-get=1
# get https://vbom.ml/util?go-get=1: Get "https://vbom.ml/util?go-get=1": dial tcp: lookup vbom.ml: no such host
go: github.com/rancher/wrangler@v0.8.5 requires
sigs.k8s.io/cli-utils@v0.16.0 requires
k8s.io/kubectrl@v0.0.0-20191219154910-1528d4eea6dd requires
vbom.ml/util@v0.0.0-20160121211510-db5cfe13f5cc: unrecognized import path "vbom.ml/util": https fetch: Get "https://vbom.ml/util?go-get=1": dial tcp: lookup vbom.ml: no such host
```



KubeCon



CloudNativeCon

North America 2023

What is a minimal Kube Control Plane?

Soon: generic, but extensible



KubeCon



CloudNativeCon

North America 2023

```
$ go run ./cmd/sample-generic-controlplane
```

- [configmaps.v1](#)
- [resourcequotas.v1](#)
- [namespaces.v1](#)
- [secrets.v1](#)
- [serviceaccounts.v1](#)
- [events.v1](#)
- [apiservices.apiregistration.k8s.io/v1](#)
- [events.events.k8s.io/v1](#)
- [selfsubjectreviews.authentication.k8s.io/v1](#)
- [tokenreviews.authentication.k8s.io/v1](#)
- [selfsubjectaccessreviews.authorization.k8s.io/v1](#)
- [selfsubjectrulesreviews.authorization.k8s.io/v1](#)
- [localsubjectaccessreviews.authorization.k8s.io/v1](#)
- [subjectaccessreviews.authorization.k8s.io/v1](#)
- [clusterroles.rbac.authorization.k8s.io/v1](#)
- [rolebindings.rbac.authorization.k8s.io/v1](#)
- [roles.rbac.authorization.k8s.io/v1](#)
- [clusterrolebindings.rbac.authorization.k8s.io/v1](#)
- [certificatesigningrequests.certificates.k8s.io/v1](#)
- [mutatingwebhookconfigurations.admissionregistration.k8s.io/v1](#)
- [validatingwebhookconfigurations.admissionregistration.k8s.io/v1](#)
- [validatingadmissionpolicies.admissionregistration.k8s.io/v1beta1](#)
- [validatingadmissionpolicybindings.admissionregistration.k8s.io/v1beta1](#)
- [customresourcedefinitions.apiextensions.k8s.io/v1](#)
- [leases.coordination.k8s.io/v1](#)
- [prioritylevelconfigurations.flowcontrol.apiserver.k8s.io/v1beta3](#)
- [flowschemas.flowcontrol.apiserver.k8s.io/v1beta3](#)
- [storageversions.internal.apiserver.k8s.io/v1alpha1](#)

WIP: generic, not extensible, no RBAC



KubeCon



CloudNativeCon

North America 2023

```
$ go run ./cmd/sample-minimal-controlplane
```

- **namespaces.v1**
- **events.v1**
- **apiservices.apiregistration.k8s.io/v1**
- **events.events.k8s.io/v1**
- **selfsubjectreviews.authentication.k8s.io/v1**
- **tokenreviews.authentication.k8s.io/v1**
- **leases.coordination.k8s.io/v1**
- **flowschemas.flowcontrol.apiserver.k8s.io/v1beta3**
- **prioritylevelconfigurations.flowcontrol.apiserver.k8s.io/v1beta3**
- **storageversions.internal.apiserver.k8s.io/v1alpha1**
- **your natively implemented APIs**

Tilt: not even a cluster



KubeCon



CloudNativeCon

North America 2023

\$ **tilt** api-resources

- **clusters**.tile.dev/v1alpha1
- **cmdimages**.tile.dev/v1alpha1
- **cmds**.tile.dev/v1alpha1
- **configmaps**.tile.dev/v1alpha1
- **dockercomposelogstreams**.tile.dev/v1alpha1
- **dockercomposeservices**.tile.dev/v1alpha1
- **dockerimages**.tile.dev/v1alpha1
- **extensionrepos**.tile.dev/v1alpha1
- **extensions**.tile.dev/v1alpha1
- **filewatches**.tile.dev/v1alpha1
- **imagemaps**.tile.dev/v1alpha1
- **kubernetesapplies**.tile.dev/v1alpha1
- **kubernetesdiscoveries**.tile.dev/v1alpha1
- **liveupdates**.tile.dev/v1alpha1
- **podlogstreams**.tile.dev/v1alpha1
- **portforwards**.tile.dev/v1alpha1
- **sessions**.tile.dev/v1alpha1
- **tiltfiles**.tile.dev/v1alpha1
- **togglebuttons**.tile.dev/v1alpha1
- **uibuttons**.tile.dev/v1alpha1
- **uiresources**.tile.dev/v1alpha1
- **uisessions**.tile.dev/v1alpha1

Acorn: no cluster, not even kube-apiserver



KubeCon



CloudNativeCon

North America 2023

\$ **acorn** kube kubectl api-resources

- [configmaps.v1](#)
- [secrets.v1](#)
- [subjectaccessreviews.authorization.k8s.io/v1](#)
- [acornimagebuilds.api.acorn.io/v1](#)
- [apps.api.acorn.io/v1](#)
- [builders.api.acorn.io/v1](#)
- [computeclasses.api.acorn.io/v1](#)
- [containerreplicas.api.acorn.io/v1](#)
- [credentials.api.acorn.io/v1](#)
- [devsessions.api.acorn.io/v1](#)
- [events.api.acorn.io/v1](#)
- [imageallowrules.api.acorn.io/v1](#)
- [images.api.acorn.io/v1](#)
- [infos.api.acorn.io/v1](#)
- [jobs.api.acorn.io/v1](#)
- [projects.api.acorn.io/v1](#)
- [regions.api.acorn.io/v1](#)
- [secrets.api.acorn.io/v1](#)
- [volumeclasses.api.acorn.io/v1](#)
- [volumes.api.acorn.io/v1](#)
- [accountevents.account.manager.acorn.io/v1](#)
- [accountquotas.account.manager.acorn.io/v1](#)
- [appcleanuppolicies.account.manager.acorn.io/v1](#)
- [imagerolebindings.account.manager.acorn.io/v1](#)
- [projectquotas.account.manager.acorn.io/v1](#)
- [projectrolebindings.account.manager.acorn.io/v1](#)
- [quotatotals.account.manager.acorn.io/v1](#)
- [quotausages.account.manager.acorn.io/v1](#)
- [workspaces.account.manager.acorn.io/v1](#)
- [quotarequests.admin.acorn.io/v1](#)



KubeCon



CloudNativeCon

North America 2023

What is a minimal Kube Control Plane?

Where do we want to go as Kube project?

Ratcheting Validation



KubeCon



CloudNativeCon

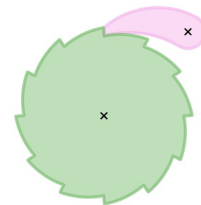
North America 2023

- <https://github.com/kubernetes/enhancements/issues/4008>

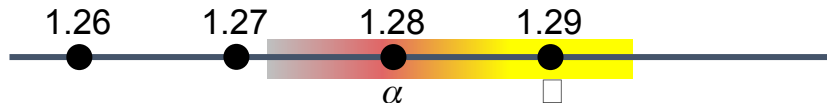
Modifying a value validation on a CRD today means that **you risk breaking the workflow of all your users**, this high price to pay limits adoption, and degrades the Kubernetes user experience: we are prevented from shifting validation logic left.

Goals

- Remove barriers blocking CRD authors from widening value validations
- Remove barriers blocking CRD authors from tightening value validations
- Do this automatically for all CRDs installed into clusters with the feature enabled



Work by  Alex Zielenski, so all the credits to him for starting this important work.



The Challenge



KubeCon



CloudNativeCon

North America 2023

```
apiVersion: apiextensions.k8s.io/v1
kind: CustomResourceDefinition
...
  schema:
    properties:
      spec:
        properties:
          replicas:
            type: number
          ip:
            type: string
```

```
$ cat object.yaml
apiVersion: group/v1
kind: Foo
spec:
  replicas: 2.5
  ip: "1.2.3.4.5"
```

```
$ kubectl apply -f object.yaml
```



The Challenge



KubeCon



CloudNativeCon

North America 2023

```
apiVersion: apiextensions.k8s.io/v1
kind: CustomResourceDefinition
...
  schema:
    properties:
      spec:
        properties:
          replicas:
            type: integer
            minimum: 0
          ip:
            type: string
            format: ipv4
```

```
$ cat object.yaml
apiVersion: group/v1
kind: Foo
spec:
  replicas: 2.5
  ip: "1.2.3.4.5"
```

```
$ kubectl apply -f object.yaml
```



The Challenge



KubeCon



CloudNativeCon

North America 2023

```
apiVersion: apiextensions.k8s.io/v1
kind: CustomResourceDefinition
...
  schema:
    properties:
      spec:
        properties:
          replicas:
            type: integer
            minimum: 0
          ip:
            type: string
            format: ipv4
```

```
$ cat object.yaml
apiVersion: group/v1
kind: Foo
spec:
  replicas: 2.5
  ip: "1.2.3.4.5"
```

```
$ kubectl edit foo
```



```
$ kubectl annotate foo x=y
```



```
$ kubectl edit # remove a finalizer
```



But CEL...



KubeCon



CloudNativeCon

North America 2023

```
apiVersion: apiextensions.k8s.io/v1
kind: CustomResourceDefinition
```

```
...
```

```
  schema:
```

```
    properties:
```

```
      spec:
```

```
        properties:
```

```
          replicas:
```

```
            type: integer
```

```
            x-kubernetes-validations:
```

- rule: "double(int(oldSelf))!=oldSelf || double(int(self))==self"
 message: only integers are allowed
- rule: "oldSelf < 0 || self >= 0"
 message: must be zero or positive

But CEL...



KubeCon



CloudNativeCon

North America 2023

```
apiVersion: apiextensions.k8s.io/v1
kind: CustomResourceDefinition
```

```
...
```

```
  schema:
```

```
    properties:
```

```
      spec:
```

```
        properties:
```

```
          replicas:
```

```
            type: integer
```

```
            x-kubernetes-validations:
```

- rule: "double(int(oldSelf))!=oldSelf || double(int(self))==self"
 message: only integers are allowed
- rule: "oldSelf < 0 || self >= 0"
 message: must be zero or positive



these only apply
on update

New: Automatic Ratcheting Validation



KubeCon



CloudNativeCon

North America 2023

```
apiVersion: apiextensions.k8s.io/v1
kind: CustomResourceDefinition
```

```
...
```

```
schema:
  properties:
    spec:
      properties:
```

```
    replicas:
```

```
      type: integer
      minimum: 0
```

```
    ip:
```

```
      type: string
      format: ipv4
```

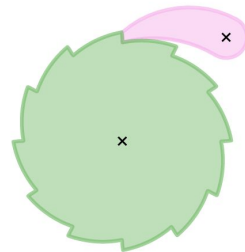
```
apiVersion: group/v1
```

```
kind: Foo
```

```
spec:
```

```
  replicas: 2.5
```

```
  ip: "1.2.3.4.5"
```



Each box is “ratcheted”:

1. Validate **green box** only if **replicas** changes.
2. Validate **blue box** only if **ip** changes.

But CEL...sure – New: OptionalOldSelf



KubeCon



CloudNativeCon

North America 2023

```
apiVersion: apiextensions.k8s.io/v1
kind: CustomResourceDefinition
...
  schema:
    properties:
      spec:
        properties:
          replicas:
            type: integer
            x-kubernetes-validations:
              - rule: >
                  (oldSelf.hasValue() && double(int(oldSelf.value())) != oldSelf.value()) ||
                  double(int(self))==self"
                  message: only integers are allowed
                  optionalOldSelf: true
              - rule: "(oldSelf.hasValue() && oldSelf.value() < 0) || self >= 0"
                  message: must be zero or positive
                  optionalOldSelf: true
```

These also apply
on Create!

Ratcheting Validation



KubeCon



CloudNativeCon


North America 2023

- <https://github.com/kubernetes/enhancements/issues/4008>

Modifying a value validation on a CRD today means that **you risk breaking the workflow of all your users**, this high price to pay limits adoption, and degrades the Kubernetes user experience: we are prevented from shifting validation logic left.

Goals

- Remove barriers blocking CRD authors from widening value validations
- Remove barriers blocking CRD authors from tightening value validations
- Do this automatically for all CRDs installed into clusters with the feature enabled

Work by  Alex Zielenski, so all the credits to him for starting this important work.





KubeCon



CloudNativeCon

North America 2023

Thanks!

Meet us on Thursday, 12pm-3pm during
“Meet the Kubernetes Contributor Community”

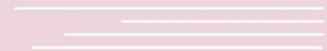


KubeCon



CloudNativeCon

North America 2023



Time for questions!