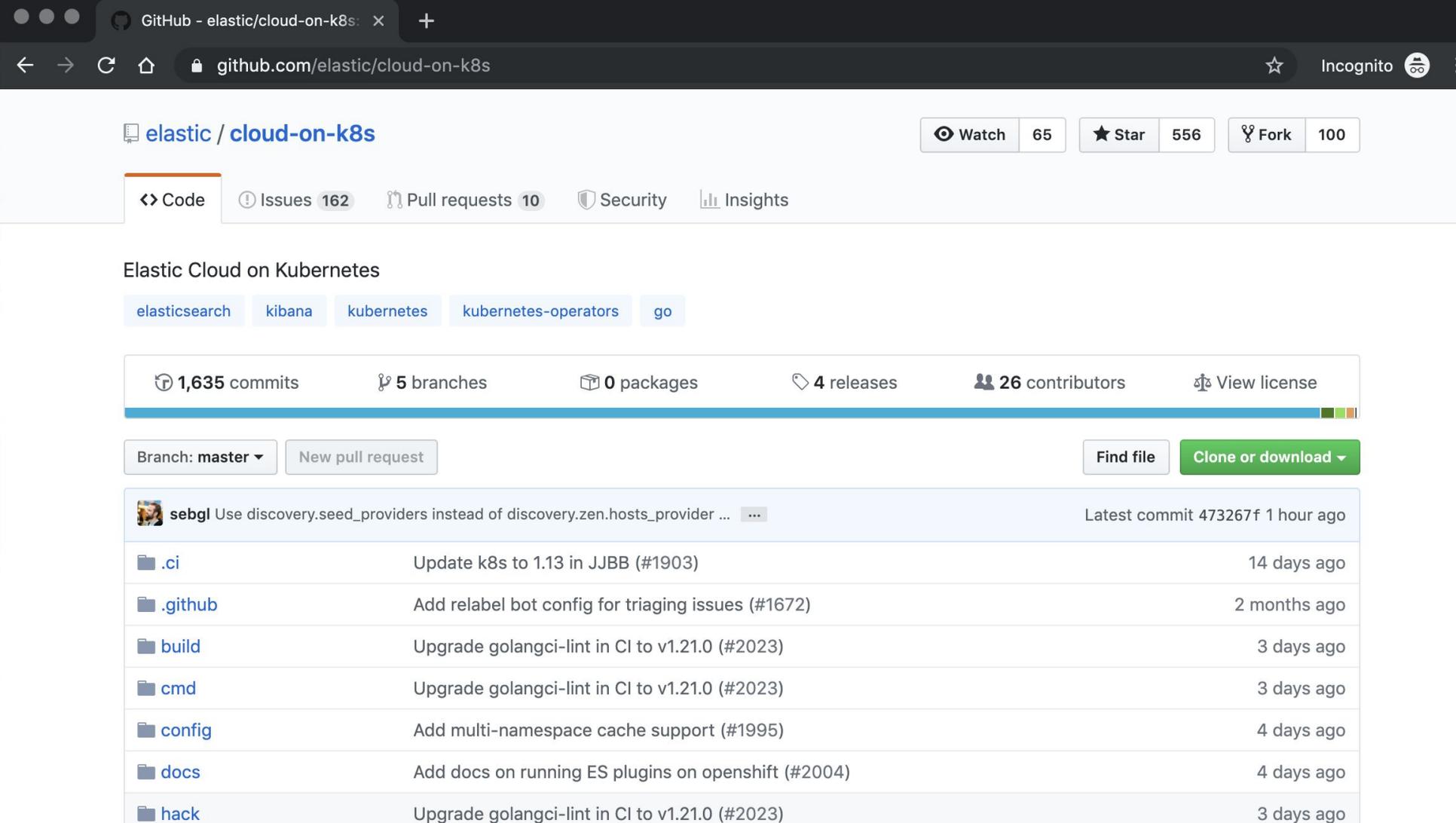# Writing a Kubernetes Operator: the Hard Parts

Sébastien Guilloux  *(@_sebgl)*
KubeCon North America 2019 (San Diego)

github.com/elastic/cloud-on-k8s

Incognito

Watch 65 | Star 556 | Fork 100

<> Code | Issues 162 | Pull requests 10 | Security | Insights

# Elastic Cloud on Kubernetes

elasticsearch | kibana | kubernetes | kubernetes-operators | go

1,635 commits | 5 branches | 0 packages | 4 releases | 26 contributors | View license

Branch: master | New pull request | Find file | Clone or download

sebgl Use discovery.seed_providers instead of discovery.zen.hosts_provider ... ... Latest commit 473267f 1 hour ago

| .ci | Update k8s to 1.13 in JJBB (#1903) | 14 days ago |
| .github | Add relabel bot config for triaging issues (#1672) | 2 months ago |
| build | Upgrade golangci-lint in CI to v1.21.0 (#2023) | 3 days ago |
| cmd | Upgrade golangci-lint in CI to v1.21.0 (#2023) | 3 days ago |
| config | Add multi-namespace cache support (#1995) | 4 days ago |
| docs | Add docs on running ES plugins on openshift (#2004) | 4 days ago |
| hack | Upgrade golangci-lint in CI to v1.21.0 (#2023) | 3 days ago |

# Operators in a nutshell
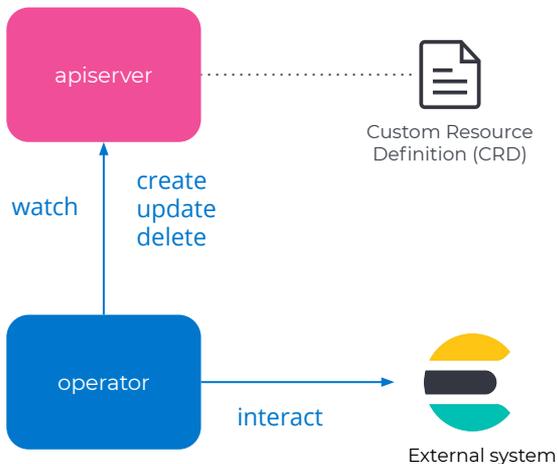## CRDs

apiserver

Custom Resource
Definition (CRD)

```
apiVersion:
elasticsearch.k8s.elastic.co/v1beta1
kind: Elasticsearch
metadata:
  name: elasticsearch-sample
spec:
  version: 7.4.0
  nodeSets:
  - name: master-nodes
    count: 3
    config:
      node.master: true
  - name: data-nodes
    count: 2
    config:
      node.data: true
```

# Operators in a nutshell

## Reconciliation loop



apiserver

Custom Resource
Definition (CRD)

watch

create
update
delete

operator

interact

External system

**New event**
A **watched resource** was created/updated/deleted

**Reconcile!**
**Get** resource spec
Reconcile **Services, Secrets, Pods**, etc.
(maybe) Interact with an **external system**

**Sequential steps**
**Return early**
**Over and over again**

# Tools & libs

# The Hard Parts

(Well, some of them)

elastic

# The operator lives in the past
## Assume you're one step behind reality

```
// retrieve Pods
client.List(&pods)        // [podA, podB]
// we miss podC, create it
client.Create(&podC)
// retrieve Pods again
client.List(&pods)        // [podA, podB]
```

The apiserver client uses a **cached reader** (by default)

elastic

# The operator lives in the past
So what? Examples from real life

**The Infinite Pod Creation Loop**
Pod missing? Create one.
Pod missing? Create one.

**The Split Brain Situation**
3 nodes? Quorum=2.
Add a 4th node. Quorum=3.
3 nodes? Quorum=2.

**The Double Rolling Upgrade Reaction**
Need to upgrade? Delete + Recreate Pods.
Need to upgrade? Delete + Recreate already upgraded Pods.

elastic

# The operator lives in the past
## What can we do?

**Optimistic concurrency** is good enough for most cases

```
> kubectl get pod my-pod -o json | jq .metadata
{
    "namespace": "default",
    "name": "my-pod",
    "uuid": "6210565b-f985-11e9-8ca3-42010a8400bb",
    "resourceVersion": "3721702"
}
```

conflict on creation

conflict on update

elastic

# The operator lives in the past
## What can we do?

**Optimistic concurrency** is good enough for most cases

```go
err := client.Delete(&pod, clientpkg.Preconditions{
    UID: &pod.UID,
    ResourceVersion: &pod.ResourceVersion,
})
if err != nil {
    return err
}
```

conflict on deletion

elastic

# The operator lives in the past
## What if we need some guarantees?

Sometimes this is **not enough**
Especially when dealing with stateful workloads

elastic

# The operator lives in the past
## What if we need some guarantees?

in-memory **expectations**

```
if !expectations.Satisfied() {
    // cache is not up-to-date yet
    return
}
err := client.Create(&pod)
// expect the Pod to be created
expectations.ExpectCreation(pod)
```

In the ReplicaSet controller *github.com/kubernetes/kubernetes/blob/master/pkg/controller/controller_utils.go*
In ECK *github.com/elastic/cloud-on-k8s/tree/master/pkg/controller/common/expectations*

# The operator lives in the past
Best practices

Use **deterministic naming**

Always assume a **stale cache**

The entire reconciliation should be **idempotent**

elastic

# Idempotent reconciliation
## An example

```
if !exists(statefulSet) {
    err := c.Create(statefulSet)
    err := c.Create(headlessSvc)
}
```

The operator may crash here

Or return an error here

*headlessSvc* will never be created

elastic

# Idempotent reconciliation
An example

```
if !exists(statefulSet) {
    err := c.Create(statefulSet)
    err := c.Create(headlessSvc)
}
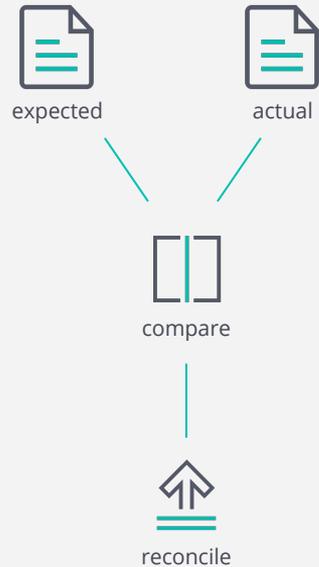```

```
if !exists(statefulSet) {
    err := c.Create(headlessSvc)
    err := c.Create(statefulSet)
}
```

Reorder operations

elastic

# Idempotent reconciliation
## An example

```
if !exists(statefulSet) {
    err := c.Create(statefulSet)
    err := c.Create(headlessSvc)
}
```

```
if !exists(statefulSet) {
    err := c.Create(headlessSvc)
    err := c.Create(statefulSet)
}
```

```
if !exists(statefulSet) {
    err := c.Create(statefulSet)
}
if !exists(headlessSvc) {
    err := c.Create(headlessSvc)
}
```

Decouple reconciliations

elastic

# Reconciling resources

# Reconciling resources

expected    actual

compare

reconcile

## The deep way

```
if !reflect.DeepEqual(expected, actual) {
    // need to update
    // ...
}
```

Not a great fit for:
- metadata
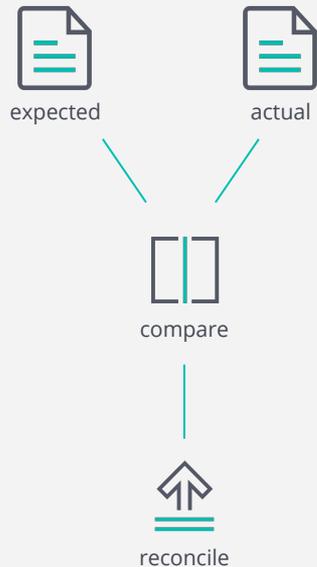- defaulted values

elastic

# Defaulted values

```yaml
apiVersion: v1
kind: Pod
metadata:
  name: mypod
spec:
  containers:
  - name: busybox
    image: busybox
```

1. Create Pod

```yaml
apiVersion: v1
kind: Pod
metadata:
  creationTimestamp: 2019-11-13T10:04:46Z
  namespace: default
  name: mypod
  uid: 052fa624-05fd-11ea-9ab1-42010a84001d
spec:
  containers:
  - name: busybox
    image: busybox
    imagePullPolicy: Always
    env:
    - name: KUBERNETES_PORT_443_TCP_ADDR
      value: c-111-dns-5e14.hcp.westus2.azmk8s.io
    resources:
      requests:
        cpu: 100m
  dnsPolicy: ClusterFirst
  securityContext: {}
```

2. Get Pod

elastic

# Reconciling resources

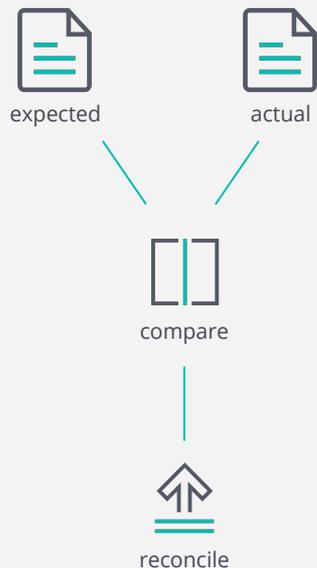

expected    actual

compare

reconcile

## The hard way

```
if sameLabels(expected, actual) &&
   sameAnnotations(expected, actual) &&
   sameReplicas(expected, actual) &&
   sameResourcesLimits(expected, actual) &&
   sameEnvVars(expected, actual) && ...
```

Not a great fit for:
- unit tests
- PR reviewer brain
- real life
- defaulted values

elastic

# Reconciling resources

expected

actual

compare

reconcile

## The smart way

```go
// annotate object with its hash
hash := HashObject(expected)
expected.Annotations[ResourceHash] = hash
// compare expected vs. actual hash
actualHash := actual.Annotations[ResourceHash]
if actualHash != hash {
    // need to update
}
```

Actual hash was **built at creation time**, hence does not include defaulted fields

*github.com/kubernetes/kubernetes/tree/master/pkg/util/hash*
*github.com/elastic/cloud-on-k8s/tree/master/pkg/controller/common/hash*

elastic

```yaml
apiVersion: elasticsearch.k8s.elastic.co/v1beta1
kind: Elasticsearch
metadata:
  name: elasticsearch-sample
spec:
  version: 7.4.0
  nodeSets:
  - name: default
    count: 1
    podTemplate:
      metadata:
        labels: {"foo": "bar"}
      spec:
        affinity:
          nodeAffinity:
            requiredDuringSchedulingIgnoredDuringExecution:
              nodeSelectorTerms:
              - matchExpressions:
                - key: environment
                  operator: In
                  values: ["e2e", "production"]
        containers:
        - name: elasticsearch
          env:
          - name: ES_JAVA_OPTS
            value: "-Xms2g -Xmx4g"
```

# Empower users
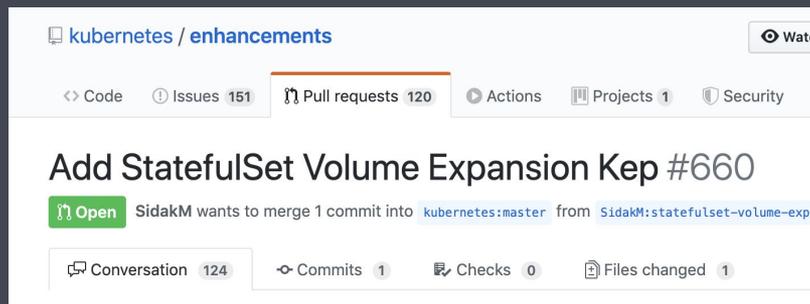## But provide good defaults

optional podTemplate

elastic

# StatefulSets

A few things to know about

elastic

## No volume resize

*KEP in progress*

## No volume resize

*There are workarounds*

## Scheduling conflicts

*Pod vs. PV*

```
kind: StorageClass
metadata:
  name: my-storage-class
volumeBindingMode: WaitForFirstConsumer
```

27

# StatefulSets
## A few things to know about

## Scheduling conflicts

*Local Volumes vs. resources*

StatefulSet **upgrade**:

1.  **Delete** Pod
2.  **Recreate** Pod ← ............ 1.5 **Another Pod** scheduled
3.  No resources available

elastic

## A few things to know about

Stateful workloads
UpdateStrategy

*Pick one*

*RollingUpdate*

*RollingUpdate.Partition*

*OnDelete*

elastic

# StatefulSets
## A few things to know about

You **don't have to** use StatefulSets.

You **can** manage Pods and PVCs directly.

elastic

# StatefulSets

## A few things to know about

You **don't have to** use StatefulSets.

You **can** manage Pods and PVCs directly.

We tried.

elastic

# Backup slides

# Testing
How do you test that monster you ended up with?

- Unit test as much as possible

  - Fake client helps with k8s interactions

- Unit tests for the entire reconciliation are hard

  - Too many code paths to visit & things to mock

- Kubebuilder integration tests

  - Local apiserver + etcd process

  - Our CI had a hard time running ITs in parallel

elastic

# Testing
How do you test that monster you ended up with?

- End-to-end tests
  - 1. Spawn a k8s cluster
  - 2. Deploy the operator
  - 3. Run tests
    - Create an Elasticsearch cluster
    - Verify it's available, with the expected spec
    - Mutate the cluster
    - Verify it eventually has the expected spec
    - Continuously ensure no downtime nor data loss during the mutation

elastic

# Testing
## Multidimensional E2E test matrix

```python
for distribution in ['vanilla', 'gke', 'aks', 'eks', 'openshift']:
    for version in ['1.11', '1.12', '1.13', '1.14', '1.15']:
        for operator in ['0.8', '0.9', '1.0.0-beta1']:
            for elasticsearch in ['6.8.0', '7.1.0', '7.2.0', '7.3.0']:
                for cluster_mutation in ['upscale', 'downscale', 'rolling_upgrade', ...]:
                    runTests(...)
```

elastic

# Namespace management
## Full flexibility

- One operator for the entire cluster
- One operator per namespace
- One operator for [*namespaceA, namespaceB*]
  - Thanks controller-runtime 0.2!
  - Need tooling for RBAC generation

elastic