# ReplicaSets & Deployments

## (The Underrated, OG Operators)

# Why do we care about ReplicaSets (formerly ReplicationControllers)?

redhat.

# Redundancy

Multiple running instances means failure can be tolerated.

# Scale

Multiple running instances mean more requests can be handled.

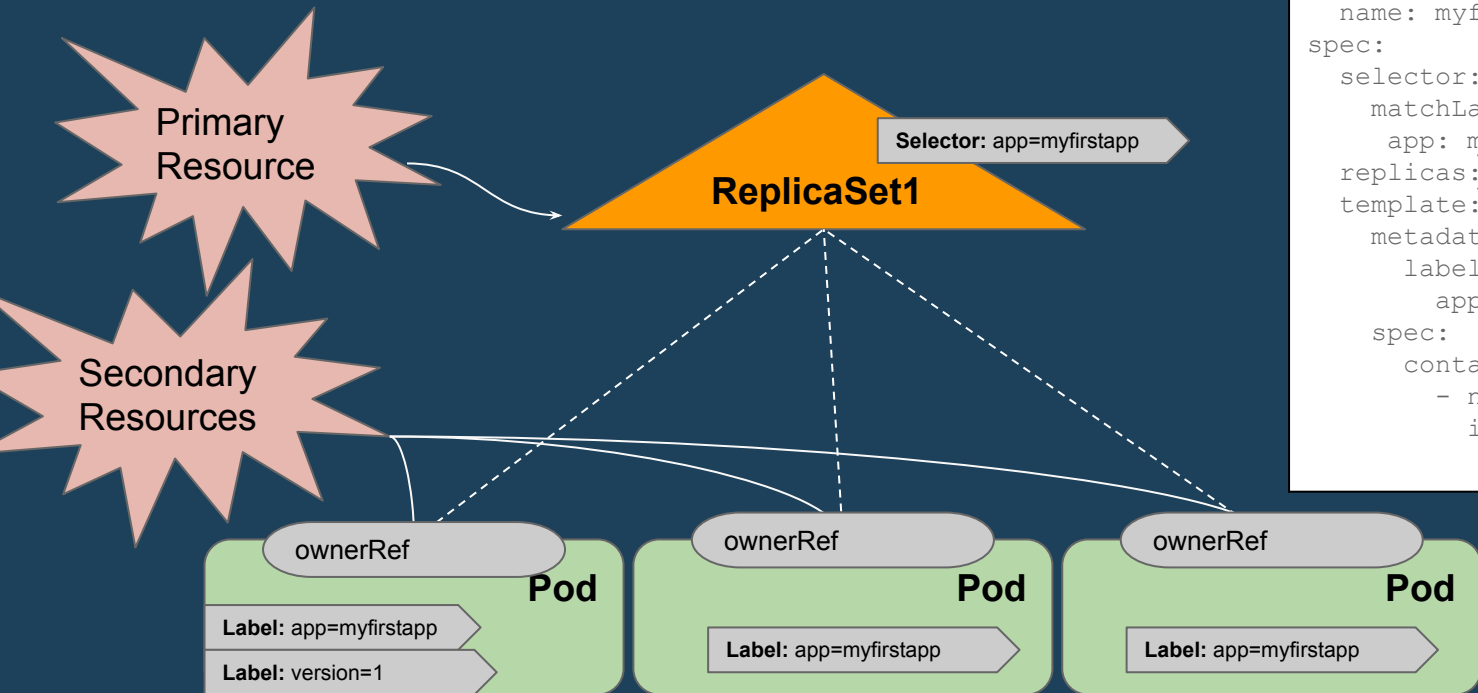| 000 | 001 | 002 | 003 | 004 |
|-----|-----|-----|-----|-----|
| cpu:80% | cpu:75% | cpu:70% | | |
| cpu:20% | cpu:25% | cpu:23% | | |

redhat.

# ReplicaSets in Action!

```
kubectl run myfirstapp --image quay.io/coreostrainme/hello-whoami:2.0.1  --restart=Never -l app=myfirstapp,version=1
```
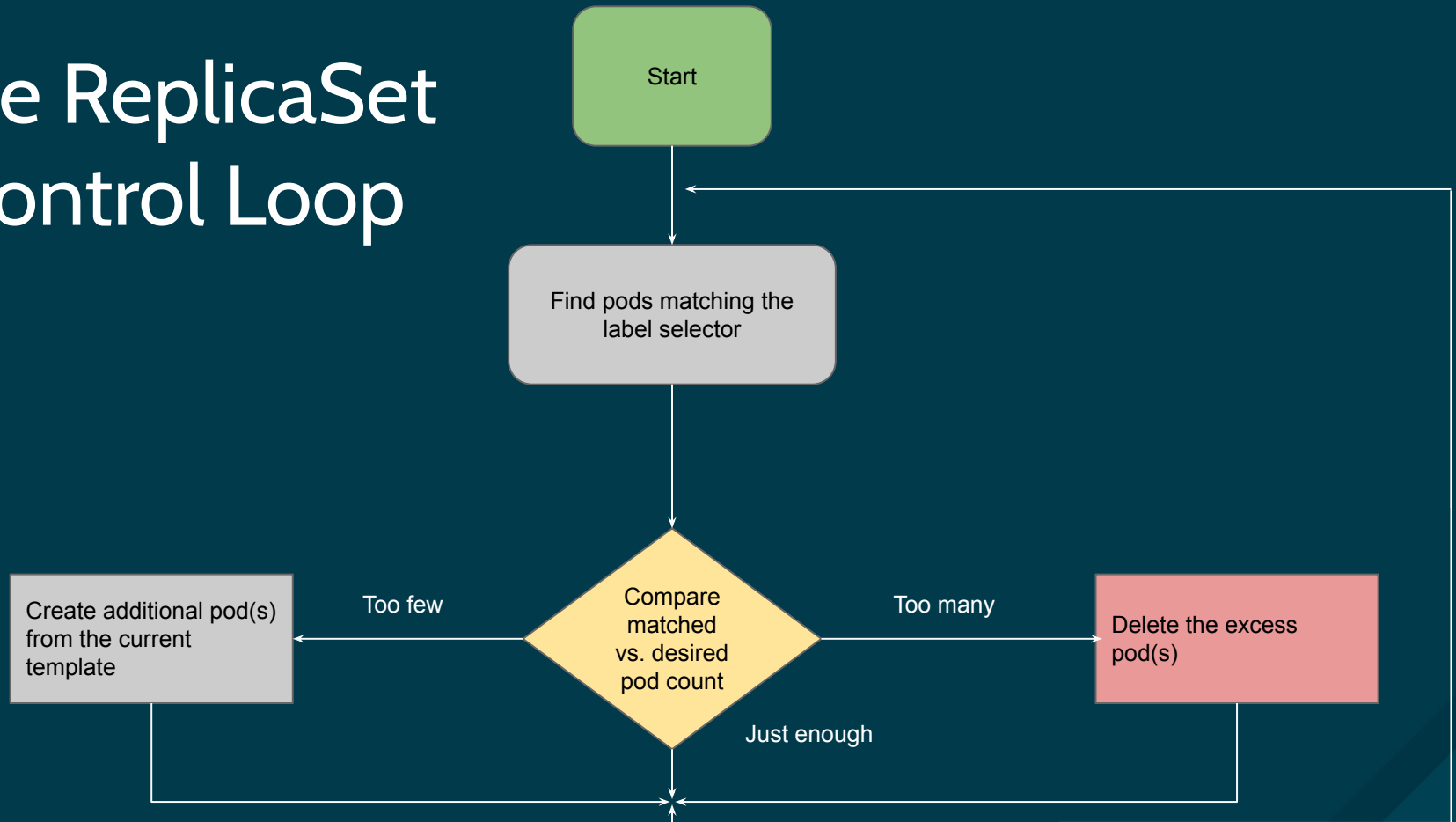
```
kubectl create -f  myfirstreplicaset.yaml
kubectl scale replicaset myfirstreplicaset --replicas=3
```

```
apiVersion: extensions/v1beta1
kind: ReplicaSet
metadata:
  name: myfirstreplicaset
spec:
  selector:
    matchLabels:
      app: myfirstapp
  replicas: ❌ 3
  template:
    metadata:
      labels:
        app: myfirstapp
    spec:
      containers:
        - name: nodejs
          image: myimage
```

Primary Resource

Selector: app=myfirstapp

**ReplicaSet1**

Secondary Resources

ownerRef
**Pod**
**Label:** app=myfirstapp
**Label:** version=1

ownerRef
**Pod**
**Label:** app=myfirstapp

ownerRef
**Pod**
**Label:** app=myfirstapp

redhat.

# The ReplicaSet Control Loop

# How do we accomplish this?

redhat.

# kube-controller-manager

# kube-controller-manager

- Daemon that embeds the core control loops shipped with Kubernetes.

- The control loop is non-terminating loop that regulates the state of the system.

-  In Kubernetes, there are multiple control loops running at all times that watch the shared state of the cluster through the apiserver and make changes attempting to move the current state towards the desired state.

- Examples include: ReplicaSets, Endpoints, Deployments, DaemonSets, etc.

redhat.

# Kubernetes Cluster

```
oc -n kube-system logs kube-controller-manager-localhost -f
```

# Let's head to GitHub.

# The ReplicaSet Controller Initialization.

```
341
342    // NewControllerInitializers is a public map of named controller groups (you can start more than one in an init func)
343    // paired to their InitFunc.  This allows for structured downstream composition and subdivision.
344    func NewControllerInitializers(loopMode ControllerLoopMode) map[string]InitFunc {
345          controllers := map[string]InitFunc{}
346          controllers["endpoint"] = startEndpointController
347          controllers["replicationcontroller"] = startReplicationController
348          controllers["podgc"] = startPodGCController
349          controllers["resourcequota"] = startResourceQuotaController
350          controllers["namespace"] = startNamespaceController
351          controllers["serviceaccount"] = startServiceAccountController
352          controllers["garbagecollector"] = startGarbageCollectorController
353          controllers["daemonset"] = startDaemonSetController
354          controllers["job"] = startJobController
355
       controllers["replicaset"] = startReplicaSetController
356
357          controllers["horizontalpodautoscaling"] = startHPAController
358          controllers["disruption"] = startDisruptionController
359          controllers["statefulset"] = startStatefulSetController
```

redhat.

# startReplicaSetController

```go
59          go statefulset.NewStatefulSetController(
60                  ctx.InformerFactory.Core().V1().Pods(),
61                  ctx.InformerFactory.Apps().V1().StatefulSets(),
62                  ctx.InformerFactory.Core().V1().PersistentVolumeClaims(),
63                  ctx.InformerFactory.Apps().V1().ControllerRevisions(),
64                  ctx.ClientBuilder.ClientOrDie("statefulset-controller"),
65          ).Run(1, ctx.Stop)
66          return nil, true, nil
67      }
68
69      func startReplicaSetController(ctx ControllerContext) (http.Handler, bool, error) {
70          if !ctx.AvailableResources[schema.GroupVersionResource{Group: "apps", Version: "v1", Resource: "replicasets"}] {
71                  return nil, false, nil
72          }
73          go replicaset.NewReplicaSetController(
74                  ctx.InformerFactory.Apps().V1().ReplicaSets(),
75                  ctx.InformerFactory.Core().V1().Pods(),
76                  ctx.ClientBuilder.ClientOrDie("replicaset-controller"),
77                  replicaset.BurstReplicas,
78          ).Run(int(ctx.ComponentConfig.ReplicaSetController.ConcurrentRSSyncs), ctx.Stop)
79          return nil, true, nil
80      }
```

redhat.

# newReplicaSetController

```go
108    // NewReplicaSetController configures a replica set controller with the specified event recorder
109    func NewReplicaSetController(rsInformer appsinformers.ReplicaSetInformer, podInformer coreinformers.PodInformer, kubeClient cl
110            eventBroadcaster := record.NewBroadcaster()
111            eventBroadcaster.StartLogging(glog.Infof)
112            eventBroadcaster.StartRecordingToSink(&v1core.EventSinkImpl{Interface: kubeClient.CoreV1().Events("")})
113            return NewBaseController(rsInformer, podInformer, kubeClient, burstReplicas,
114                    apps.SchemeGroupVersion.WithKind("ReplicaSet"),
115                    "replicaset_controller",
116                    "replicaset",
117                    controller.RealPodControl{
118                            KubeClient: kubeClient,
119                            Recorder:   eventBroadcaster.NewRecorder(scheme.Scheme, v1.EventSource{Component: "replicaset-controlle
120                    },
121            )
122    }
```

redhat.

# newBaseController

# newBaseController

```
141    rsInformer.Informer().AddEventHandler(cache.ResourceEventHandlerFuncs{    2
142            AddFunc:     rsc.enqueueReplicaSet,
143            UpdateFunc: rsc.updateRS,
144            // This will enter the sync loop and no-op, because the replica set has been deleted from the store.
145            // Note that deleting a replica set immediately after scaling it to 0 will not work. The recommended
146            // way of achieving this is by performing a `stop` operation on the replica set.
147            DeleteFunc: rsc.enqueueReplicaSet,
148    })
149    rsc.rsLister = rsInformer.Lister()
150    rsc.rsListerSynced = rsInformer.Informer().HasSynced
```

Watching Primary Resource: kind:ReplicaSet

# newBaseController

```
152         podInformer.Informer().AddEventHandler(cache.ResourceEventHandlerFuncs{  3
153                 AddFunc: rsc.addPod,
154                 // This invokes the ReplicaSet for every pod change, eg: host assignment. Though this might seem like
155                 // overkill the most frequent pod update is status, and the associated ReplicaSet will only list from
156                 // local storage, so it should be ok.
157                 UpdateFunc: rsc.updatePod,
158                 DeleteFunc: rsc.deletePod,
159         })
```

Watching Secondary Resource: kind:Pod

# Creating another ReplicaSet

`kubectl create -f  mysecondreplicaset.yaml`

**ReplicaSet2**

**Selector:** app=myfirstapp

**Selector:** version=1

**ReplicaSet1**

**Selector:** app=myfirstapp

ownerRef

**Pod**

**Label:** app=myfirstapp

**Label:** version=1

ownerRef

**Pod**

**Label:** app=myfirstapp

**Label:** version=1

ownerRef

**Pod**

**Label:** app=myfirstapp

ownerRef

**Pod**

**Label:** app=myfirstapp

```yaml
apiVersion: extensions/v1beta1
kind: ReplicaSet
metadata:
  name: mysecondreplicaset
spec:
  selector:
    matchLabels:
      app: myfirstapp
      version: 1
  replicas: 1
  template:
    metadata:
      labels:
        app: myfirstapp
        version: 1
    spec:
      containers:
        - name: nodejs
          image: myimage
```

redhat.

# Creating Another Orphan Pod

```
kubectl run myfirstapp --image quay.io/coreostrainme/hello-whoami:2.0.1  --restart=Never -l app=myfirstapp,version=1
```

**Selector:** app=myfirstapp

**Selector:** version=1

**ReplicaSet2**

**Pod**

**Label:** app...

**Label:** ...

**Selector:** app=myfirstapp

**ReplicaSet1**

ownerRef

**Pod**

**Label:** app=myfirstapp

**Label:** version=1

ownerRef

**Pod**

**Label:** app=myfirstapp

**Label:** version=1

ownerRef

**Pod**

**Label:** app=myfirstapp

ownerRef

**Pod**

**Label:** app=myfirstapp

```
apiVersion: extensions/v1beta1
kind: ReplicaSet
metadata:
  name: myfirstreplicaset
spec:
  selector:
    matchLabels:
      app: myfirstapp
      version: 1
  replicas: 1
  template:
    metadata:
      labels:
        app: myfirstapp
        version: 1
    spec:
      containers:
        - name: nodejs
          image: myimage
```

**Chapter 4**
**Designing Infrastructure Applications**

*The **reconciler pattern** is a software pattern that can be used or expanded upon for managing cloud native infrastructure. The pattern enforces the idea of having two representations of the infrastructure—the first being the actual state of the infrastructure, and the second being the expected state of the infrastructure.*

*The **reconciler pattern** will force the engineer to have two independent avenues for getting either of these representations, as well as to implement a solution to reconcile the actual state into the expected state.*
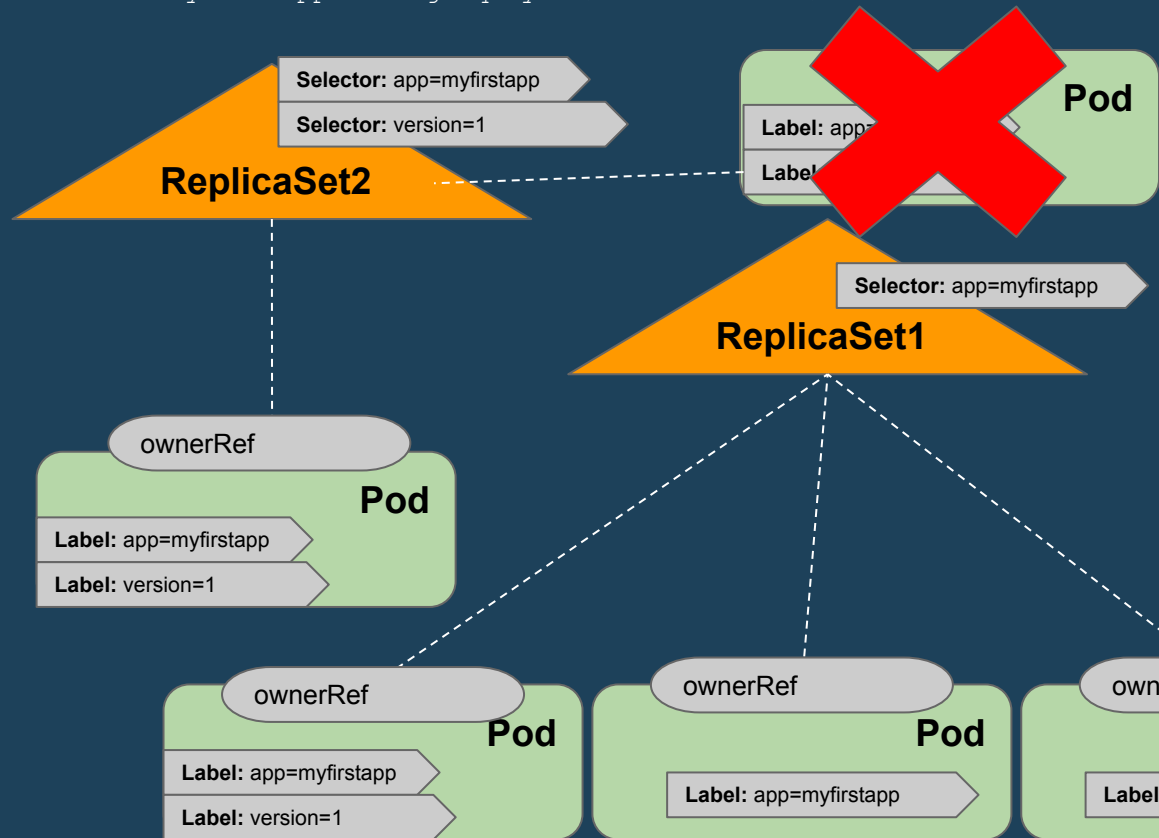
# ReplicaSets in Action!
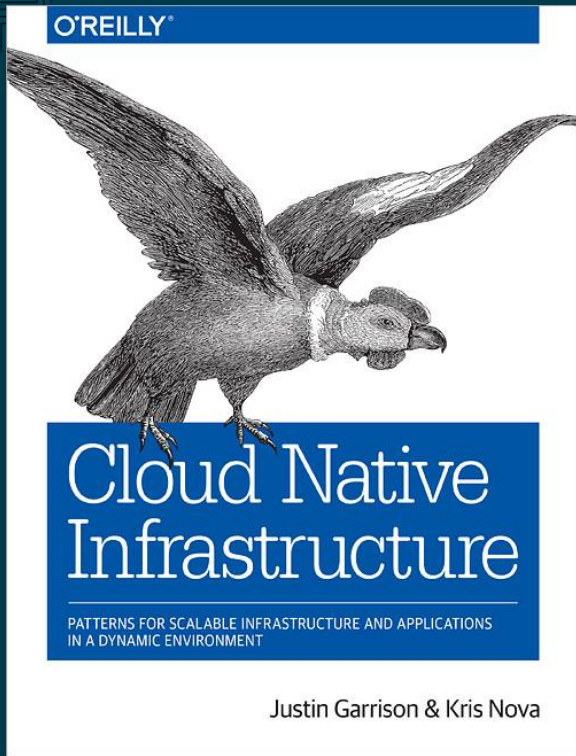
`kubectl create -f myfirstreplicaset.yaml`

```
apiVersion: extensions/v1beta1
kind: ReplicaSet
metadata:
  name: myfirstreplicaset
spec:
  selector:
    matchLabels:
      app: myfirstapp
  replicas: 3
  template:
    metadata:
      labels:
        app: myfirstapp
    spec:
      containers:
        - name: nodejs
          image: myimage
status:
  availableReplicas: 0
  fullyLabeledReplicas: 3
  observedGeneration: 1
  readyReplicas: 0
  replicas: 3
```

status:
  availableReplicas: 3
  fullyLabeledReplicas: 3
  observedGeneration: 1
  readyReplicas: 3
  replicas: 3

**Selector:** app=myfirstapp

**ReplicaSet1**

**Pod**
**Label:** app=myfirstapp

**Pod**
**Label:** app=myfirstapp

**Pod**
**Label:** app=myfirstapp

Kube-API

c.Watch(Replicaset)

c.Watch(Pods, OwnerType: ReplicaSet)

ReplicaSetController

ReplicaSet
Add Event

Func Reconcile

Fetch the rs from cache

r.Client.List Pods by label: rs.metadata.label

r.Client.Create Pod 1

r.Client.Create Pod 2

r.Client.Create Pod 3

r.Client.Update rs Status

Pod 1
Add Event

Pod 2
Add Event

Pod 3
Add Event

Func Reconcile

Fetch the rs from cache

r.Client.List Pods by owner: rs.Name

r.Client.Update rs.Status

redhat.

# ReplicaSets in Action!
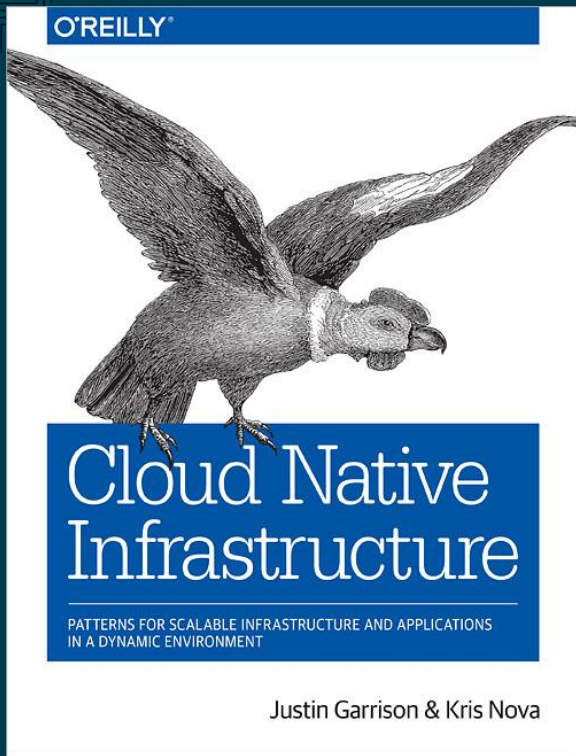
`kubectl create -f  myfirstreplicaset.yaml`

```
apiVersion: extensions/v1beta1
kind: ReplicaSet
metadata:
  name: myfirstreplicaset
spec:
  selector:
    matchLabels:
      app: myfirstapp
  replicas: 3
  template:
    metadata:
      labels:
        app: myfirstapp
    spec:
      containers:
        - name: nodejs
          image: myimage

status:
  availableReplicas: 3
  fullyLabelledReplicas: 3
  observedGeneration: 1
  readyReplicas: 3
  replicas: 3
```

**Selector:** app=myfirstapp

**ReplicaSet1**

**Pod**

**Label:** app=myfirstapp

**Pod**

**Label:** app=myfirstapp

**Pod**

**Label:** app=myfirstapp

Kube-API

c.Watch(Replicaset)

c.Watch(Pods, OwnerType: ReplicaSet)

ReplicaSetController

Pod 1
Delete Event

Func Reconcile

Fetch the rs from cache

r.Client.List Pods by label: rs.metadata.name

r.Client.Update rs.Status

r.Client.Create Pod

Pod 4
Add Event

Func Reconcile

Fetch the rs from cache

r.Client.List Pods by owner: rs.Name

r.Client.Update rs.Status

Pod 4
Update Event
X 4

Func Reconcile

# Let's Identify Primary/Secondary Resources for Existing Kubernetes Controllers!

## (without looking at the code!)

redhat.

**Endpoints**

# EndPoint Controller
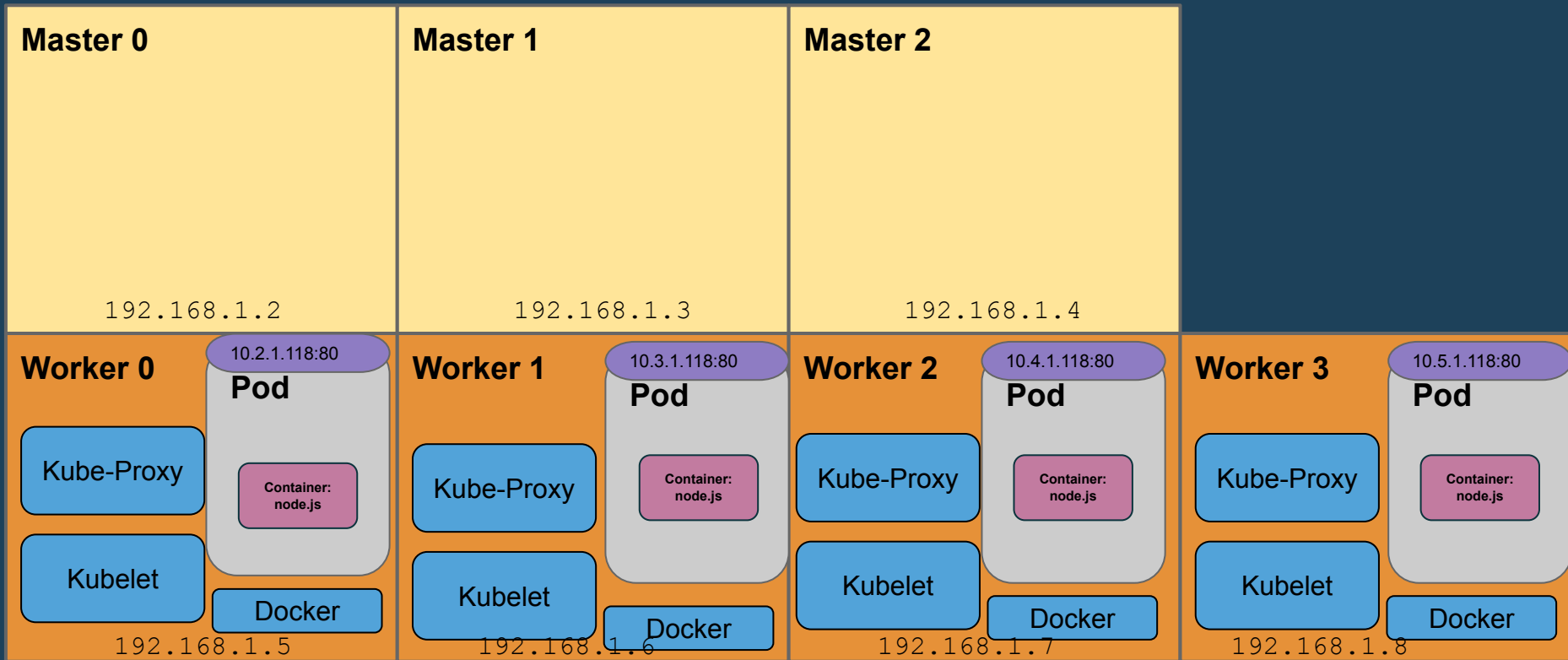
```go
71    // NewEndpointController returns a new *EndpointController.
72    func NewEndpointController(podInformer coreinformers.PodInformer, serviceInformer coreinformers.ServiceInformer,
73        endpointsInformer coreinformers.EndpointsInformer, client clientset.Interface) *EndpointController {
          if client != nil && client.CoreV1().RESTClient().GetRateLimiter() != nil {
              metrics.RegisterMetricAndTrackRateLimiterUsage("endpoint_controller", client.CoreV1().RESTClient().GetRateLimit
          }
77        e := &EndpointController{
78            client:          client,
79            queue:           workqueue.NewNamedRateLimitingQueue(workqueue.DefaultControllerRateLimiter(), "endpoint"),
80            workerLoopPeriod: time.Second,
81        }
82
          serviceInformer.Informer().AddEventHandler(cache.ResourceEventHandlerFuncs{
84            AddFunc: e.enqueueService,
85            UpdateFunc: func(old, cur interface{}) {
                  e.enqueueService(cur)
87            },
88            DeleteFunc: e.enqueueService,
89        })
90        e.serviceLister = serviceInformer.Lister()
91        e.servicesSynced = serviceInformer.Informer().HasSynced
92
          podInformer.Informer().AddEventHandler(cache.ResourceEventHandlerFuncs{
94            AddFunc:    e.addPod,
95            UpdateFunc: e.updatePod,
96            DeleteFunc: e.deletePod,
97        })
98        e.podLister = podInformer.Lister()
99        e.podsSynced = podInformer.Informer().HasSynced
100
101       e.endpointsLister = endpointsInformer.Lister()
102       e.endpointsSynced = endpointsInformer.Informer().HasSynced
103
104       return e
105   }
```

redhat.

DaemonSets

# DaemonSets

| Master 0 | Master 1 | Master 2 |
|---|---|---|
| 192.168.1.2 | 192.168.1.3 | 192.168.1.4 |

**Worker 0**

10.2.1.118:80
**Pod**

Kube-Proxy

Container:
node.js

Kubelet

Docker

192.168.1.5

**Worker 1**

10.3.1.118:80
**Pod**

Kube-Proxy

Container:
node.js

Kubelet

Docker

192.168.1.6

**Worker 2**

10.4.1.118:80
**Pod**

Kube-Proxy

Container:
node.js

Kubelet

Docker

192.168.1.7

**Worker 3**

10.5.1.118:80
**Pod**

Kube-Proxy

Container:
node.js

Kubelet

Docker

192.168.1.8

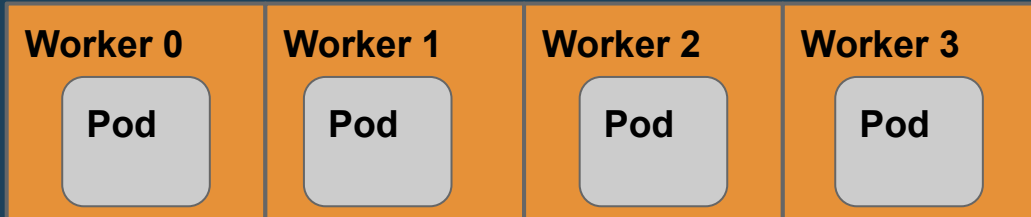redhat.

# DaemonSetController



```
143    func NewDaemonSetsController(
144        daemonSetInformer appsinformers.DaemonSetInformer,
145        historyInformer appsinformers.ControllerRevisionInformer,
146        podInformer coreinformers.PodInformer,
147        nodeInformer coreinformers.NodeInformer,
148        kubeClient clientset.Interface,
149        failedPodsBackoff *flowcontrol.Backoff,
```

Kube-API

c.Watch(DaemonSets)
c.Watch(Pods)
c.Watch(Nodes)

DaemonSetController

Daemon Set

DaemonSet Add Event
r.Client.Create Pod

Pod Add Event
r.Client.Create Pod

Pod Add Event
r.Client.Create Pod

Pod Add Event

Node Add Event
r.Client.Create Pod

| Worker 0 | Worker 1 | Worker 2 | Worker 3 |
|----------|----------|----------|----------|
| Pod | Pod | Pod | Pod |

*Administrator adds a new node!*

# DeploymentControll



```
100  func NewDeploymentController(dInformer appsinformers.DeploymentInformer, rsInformer appsinformers.ReplicaSetInformer, podInfo
101      eventBroadcaster := record.NewBroadcaster()
102      eventBroadcaster.StartLogging(klog.Infof)
103      eventBroadcaster.StartRecordingToSink(&v1core.EventSinkImpl{Interface: client.CoreV1().Events("")})
104
105      if client != nil && client.CoreV1().RESTClient().GetRateLimiter() != nil {
106          if err := metrics.RegisterMetricAndTrackRateLimiterUsage("deployment_controller", client.CoreV1().RESTClient().
107              return nil, err
108          }
109      }
110      dc := &DeploymentController{
111          client:        client,
112          eventRecorder: eventBroadcaster.NewRecorder(scheme.Scheme, v1.EventSource{Component: "deployment-controller"}),
113          queue:         workqueue.NewNamedRateLimitingQueue(workqueue.DefaultControllerRateLimiter(), "deployment"),
114      }
115      dc.rsControl = controller.RealRSControl{
116          KubeClient: client,
117          Recorder:   dc.eventRecorder,
118      }
119
120      dInformer.Informer().AddEventHandler(cache.ResourceEventHandlerFuncs{
121          AddFunc:    dc.addDeployment,
122          UpdateFunc: dc.updateDeployment,
123          // This will enter the sync loop and no-op, because the deployment has been deleted from the store.
124          DeleteFunc: dc.deleteDeployment,
125      })
126      rsInformer.Informer().AddEventHandler(cache.ResourceEventHandlerFuncs{
127          AddFunc:    dc.addReplicaSet,
128          UpdateFunc: dc.updateReplicaSet,
129          DeleteFunc: dc.deleteReplicaSet,
130      })
131      podInformer.Informer().AddEventHandler(cache.ResourceEventHandlerFuncs{
132          DeleteFunc: dc.deletePod,
133      })
134
```

# Let's Identify Primary/Secondary Resources for Existing Kubernetes Controllers!

Garbage Collection assists in deleting objects that have an owner that **no longer** exists.
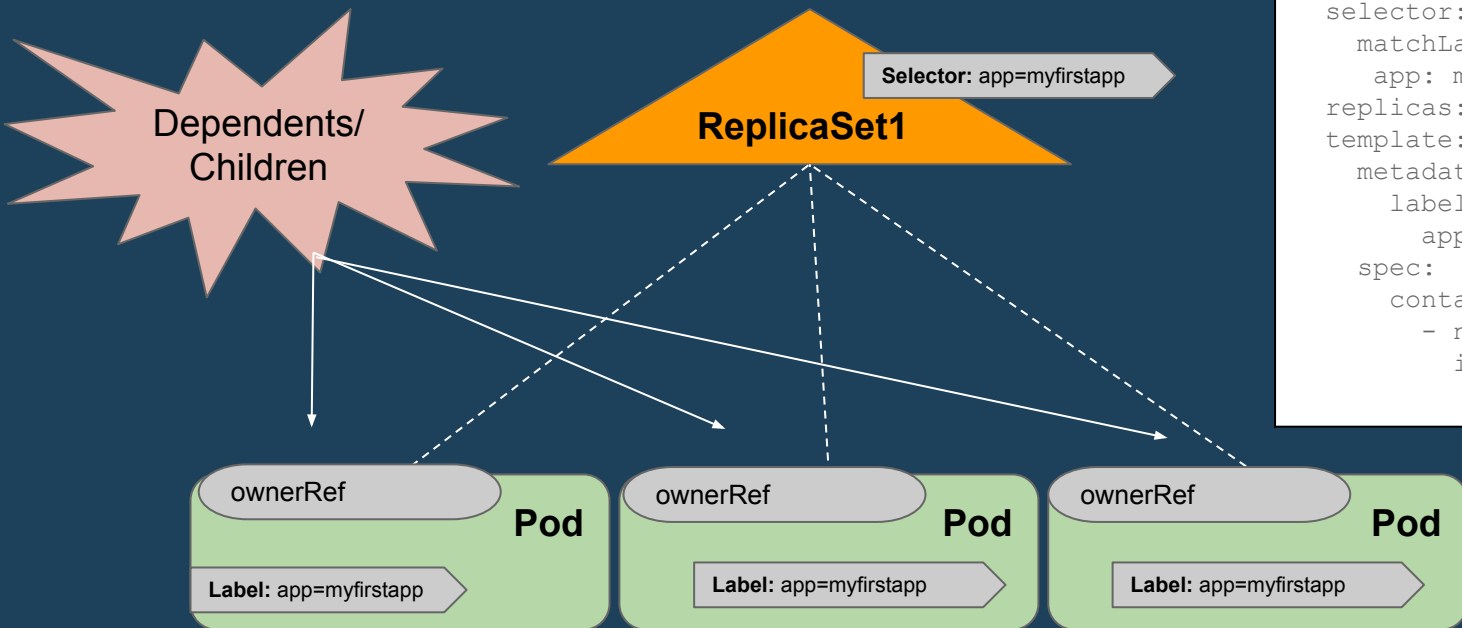
redhat.

# controllermanager.go

```go
336        )
337
338    const (
339            saTokenControllerName = "serviceaccount-token"
340    )
341
342    // NewControllerInitializers is a public map of named controller groups (you can start more than one in an init func)
343    // paired to their InitFunc.  This allows for structured downstream composition and subdivision.
344    func NewControllerInitializers(loopMode ControllerLoopMode) map[string]InitFunc {
345            controllers := map[string]InitFunc{}
346            controllers["endpoint"] = startEndpointController
347            controllers["replicationcontroller"] = startReplicationController
348            controllers["podgc"] = startPodGCController
349            controllers["resourcequota"] = startResourceQuotaController
350            controllers["namespace"] = startNamespaceController
351
352            controllers["garbagecollector"] = startGarbageCollectorController
353                                        = startDaemonSetController
354            controllers["job"] = startJobController
355            controllers["deployment"] = startDeploymentController
356            controllers["replicaset"] = startReplicaSetController
357            controllers["horizontalpodautoscaling"] = startHPAController
```

redhat.

# OwnerReferences

`kubectl create -f  myfirstreplicaset.yaml`



Dependents/
Children

**ReplicaSet1**

Selector: app=myfirstapp

```
apiVersion: extensions/v1beta1
kind: ReplicaSet
metadata:
  name: myfirstreplicaset
spec:
  selector:
    matchLabels:
     app: myfirstapp
  replicas: 3
  template:
    metadata:
      labels:
        app: myfirstapp
    spec:
      containers:
       - name: nodejs
         image: myimage
```

ownerRef
**Pod**
**Label:** app=myfirstapp

ownerRef
**Pod**
**Label:** app=myfirstapp

ownerRef
**Pod**
**Label:** app=myfirstapp

# OwnerReferences

Only applicable when doing "foreground" delete (optional)

GroupVersion of Owner Object (Required)

Kind of Owner Object (Required)

```
ownerReferences:
- apiVersion: apps/v1
  blockOwnerDeletion: true
  controller: true
  kind: ReplicaSet
  name: myfirstreplicaset
  uid: 30c68160-d992-11e8-84d9-e6f5b7702569
```

Strictly informational: shows that a Controller set the ownerReferences (optional).

Name of Owner Object (Required)
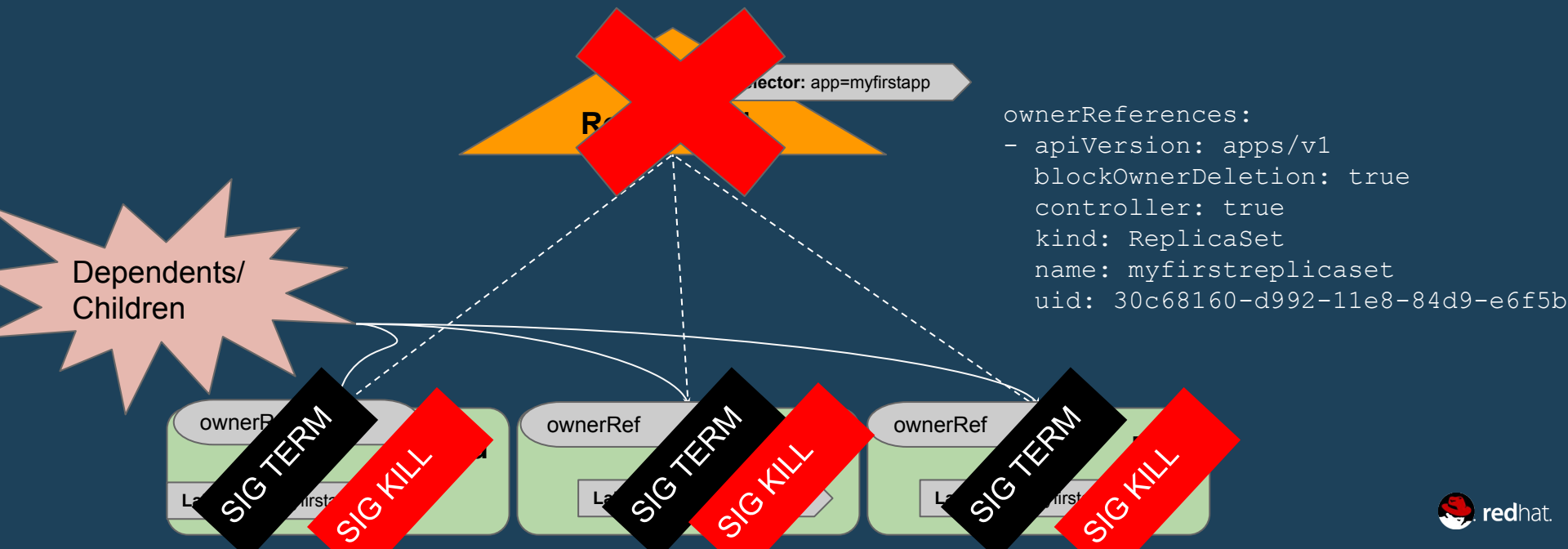
UID of Owner Object (Required)

*querying API for UID not currently supported.

redhat.

# Background Delete

```
oc delete -f myfirstreplicaset
```

```
curl -X DELETE localhost:8080/apis/apps/v1/namespaces/default/replicasets/my-repset \
-d '{"kind":"DeleteOptions","apiVersion":"v1","propagationPolicy":"Background"}' \
-H "Content-Type: application/json"
```

30s..



Dependents/
Children

```
ownerReferences:
- apiVersion: apps/v1
  blockOwnerDeletion: true
  controller: true
  kind: ReplicaSet
  name: myfirstreplicaset
  uid: 30c68160-d992-11e8-84d9-e6f5b
```
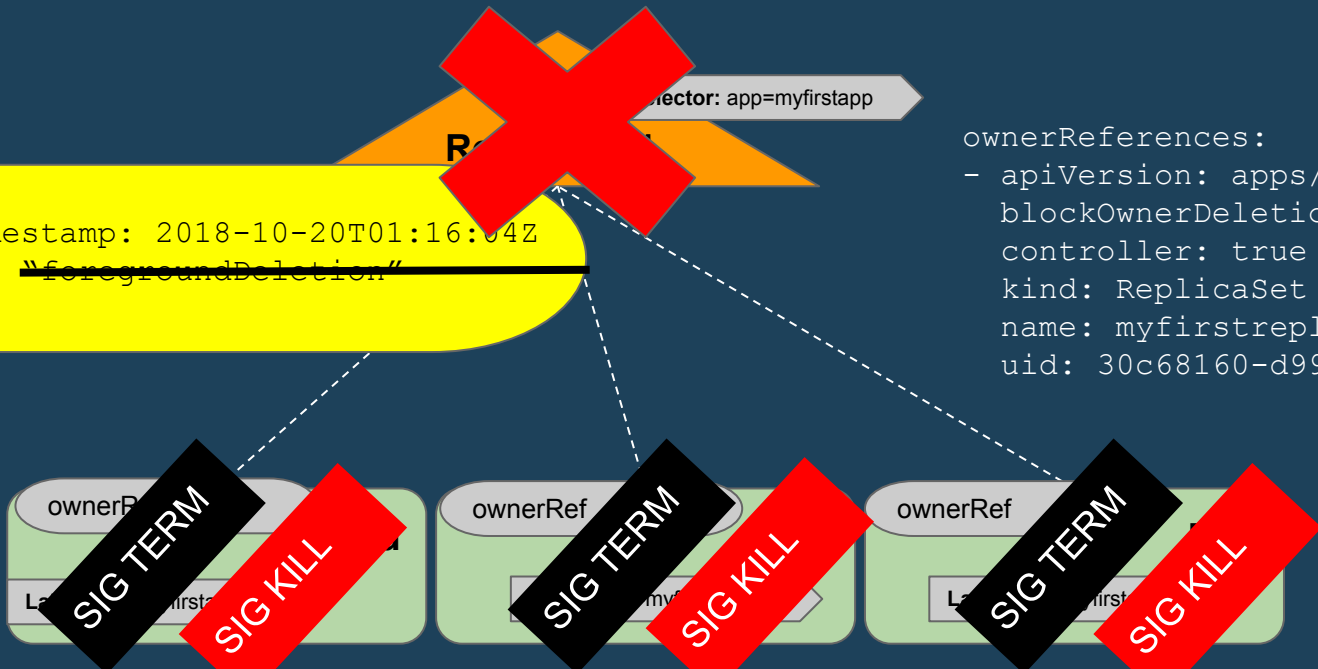
# Foreground Delete

```
oc delete -f myfirstreplicaset
```

```
curl -X DELETE localhost:8080/apis/apps/v1/namespaces/default/replicasets/my-repset \
-d '{"kind":"DeleteOptions","apiVersion":"v1","propagationPolicy":"Foreground"}' \
-H "Content-Type: application/json"
```

30s..

lector: app=myfirstapp

metadata:
  deletionTimestamp: 2018-10-20T01:16:04Z
  Finalizers: "foregroundDeletion"

ownerReferences:
- apiVersion: apps/v1
  blockOwnerDeletion: true
  controller: true
  kind: ReplicaSet
  name: myfirstreplicaset
  uid: 30c68160-d992-11e8-84d9-e6f5b

ownerRef          SIG TERM    SIG KILL
ownerRef          SIG TERM    SIG KILL
ownerRef          SIG TERM    SIG KILL

# Finalizers

Allows controllers to implement conditions that must be completed before the object can be deleted.

```
apiVersion: "stable.example.com/v1"
kind: MySQL
metadata:
  finalizers:
  - finalizer.stable.example.com
```

```
metadata:
  deletionTimestamp: 2018-10-20T01:16:04Z
```

**Pod**

**Pod**

**Pod**

Controller