# Life of a Kubernetes Watch Event

Haowei Cai (roycaihw@github), Google
Wenjia Zhang (wenjiaswe@github), Google

# About us

**Haowei Cai**(@roycaihw)

*Software Engineer for Google Cloud. He is one of the owners of Kubernetes Python client library and an active Kubernetes SIG API Machinery contributor.*


**Wenjia Zhang** (@wenjiaswe)

*Software Engineer on Kubernetes team at Google. Active contributor for Kubernetes and etcd open source projects.*

# Agenda

- **What** is a Kubernetes Watch Event?
- **Why** is Watch Event important for Kubernetes?
- **How** is the life of a Kubernetes Watch Event?
- **Key Takeaways**

What is a Kubernetes Watch Event?

# What is Watch?

Watch is an incremental change notification feed

# Watch vs. Poll

Low latency

Single connection

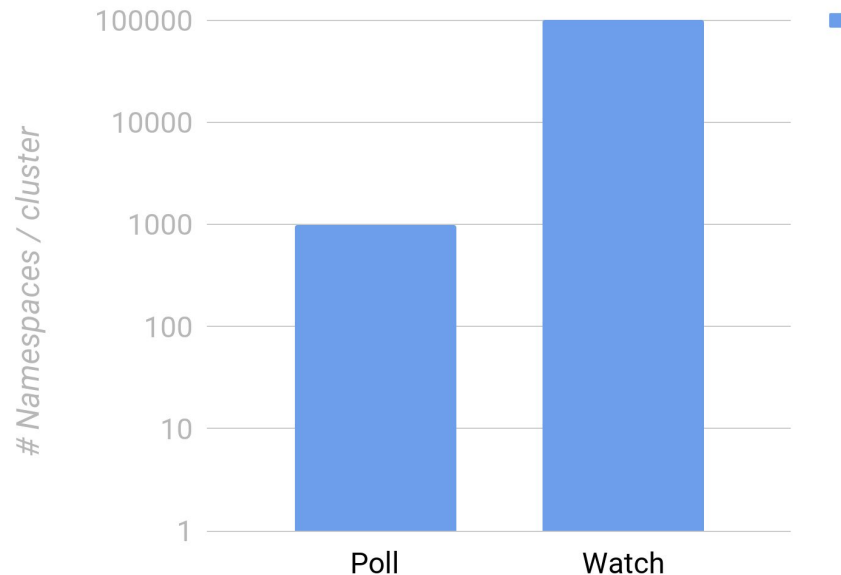## Watch

Extra load
Extra latency

Multiple connection

## Poll

# Watch vs. Poll

Kubelet on nodes:
- Previous: periodically poll kube-apiserver for secrets and configmaps
- Now: watch individual secrets
- OSS PR: Kubelet watches necessary secrets/configmaps instead of periodic polling #64752

## Scalability of Poll vs Watch

# What is Event?

## A Single event to a watched resource

```go
// Event represents a single event to a watched resource.
// +k8s:deepcopy-gen=true
type Event struct {
    Type EventType

    // Object is:
    //  * If Type is Added or Modified: the new state of the object.
    //  * If Type is Deleted: the state of the object immediately before deletion.
    //  * If Type is Error: *api.Status is recommended; other types may make sense
    //    depending on context.
    Object runtime.Object
}
```

**Watched resource:** node      **Event:** add a node

**Watched resource:** pod      **Event**: modify a pod

**Watched resource:** replicaSet      **Event**: delete a replicaSet

**Why** is Watch Event important for Kubernetes?

Level Triggering and Soft Reconciliation

# Declarative configuration

Actual state → Desired state

# Declarative configuration

Actual state

WATCH...

Desired state

**How** is the lifecycle of a K8s watch event?

# Fingerprint of events: resourceVersion

resourceVersion:

a string that identifies the common version of the objects returned by in a list. This value MUST be treated as opaque by clients and passed unmodified back to the server. A resource version is valid on a single kind of resource across namespaces.



Kubernetes object

resourceVersion

# Lifecycle of a K8s watch event

# Lifecycle of a K8s watch event

Event: Add pod rc-mprbn
resourceVersion: 12995

Lifecycle of a K8s watch event

Event: Add pod rc-mprbn
resourceVersion: 12995

# Lifecycle of a K8s watch event



Event: Add pod rc-mprbn
resourceVersion: 12995

etcd → api → Client-go → sched / c-m / node

Watch Event on etcd

# Watch in etcd

Etcd **watch** feature provides an event-based interface for asynchronously monitoring changes to keys.

```
tx01 $ etcdctl --endpoints=$ENDPOINTS
```

```
                        gyuho@tx01: ~ 67x6
tx01 $ etcdctl --endpoints=$ENDPOINTS _
```

# Revision

Revision (etcd) = resourceVersion (apiserver)
- The key space maintains multiple revisions.
- Each atomic mutative operation creates a new revision on the key space.
- All data held by previous revisions remains unchanged.
- If the store is compacted to save space, revisions before the compact revision will be removed.
- Revisions are monotonically increasing over the lifetime of a cluster.

*this is a current implementation detail of the etcd storage layer, and does not guarantee resourceVersion will be numeric or monotonically increasing in the future

# Watch event on etcd



Watches make three guarantees about events:

- Ordered - events are ordered by revision
- Reliable - a sequence of events will never drop any subsequence of events
- Atomic - a list of events is guaranteed to encompass complete revisions

Watch Event on Kube Api-server

# Watch Event on Kube Api-server

```
┌─────────────────┐        ┌─────────────────┐        ┌─────────────────┐
│                 │        │                 │        │                 │
│   etcd watcher  │  ═══>  │   watch_cache   │  ═══>  │   cacheWatcher  │
│                 │        │                 │        │                 │
└─────────────────┘        └─────────────────┘        └─────────────────┘
```

# Watch Event on Kube Api-server

etcd watcher → watch_cache → cacheWatcher

Sending events...

etcd ← GRPC → Client

Client is watching...

# Watch Event on Kube Api-server

# Watch Event on Kube Api-server

# Watch Event on Kube Api-server

# Watch Event on Kube Api-server

# Watch Event in Client-go

# What is Client-go?

| |
|---|
| Clientset |
| Dynamic Client |
| REST Client |
| Informer |
| ... |

https://github.com/kubernetes/client-go

- Go clients for talking to a kubernetes cluster
- Used by Kubernetes itself

# What is Informer?

| Clientset |
| Dynamic Client |
| REST Client |
| **Informer** |
| ... |

| Reflector |
| DeltaFIFO |
| Local Cache |
| Callback |

k8s.io/client-go/tools/cache
k8s.io/client-go/informers

- Useful component for building event-oriented controllers
- Used by control plane controllers, kubelet, etc.
- Reflector used by kube-apiserver watch cache

# Kubernetes controller workflow

# Kubernetes controller workflow

# Kubernetes controller workflow
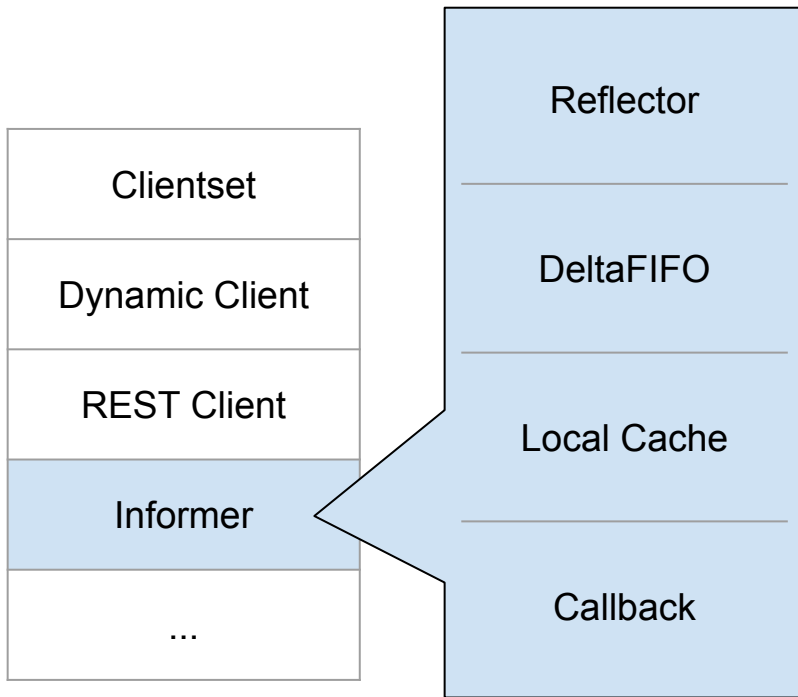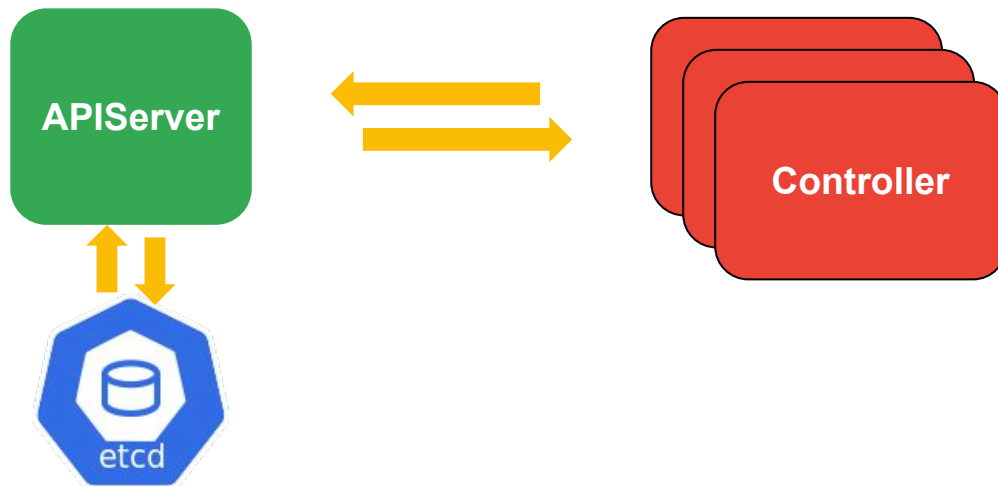
# Kubernetes controller workflow

# List and Watch

```
1    func ListAndWatch() error {
2        // Explicitly set "0" as resource version - it's fine for the List()
3        // to be served from cache and potentially be delayed relative to
4        // etcd contents. Reflector framework will catch up via Watch() eventually.
5        list, err = listWatcher.List(ListOptions{ResourceVersion: "0"})
6
7        // Extract the actual ResourceVersion that we've listed
8        resourceVersion = list.GetResourceVersion()
9
10       for {
11           options = ListOptions{ResourceVersion: resourceVersion}
12           w, err = listWatcher.Watch(options)
13           if err {
14               // most likely apiserver is not responsive. It doesn't make
15               // sense to re-list all objects because most likely we will
16               // be able to restart watch where we ended.
17               if err.IsError("connection refused") {
18                   sleep(time.Second)
19                   continue
20               }
21               // Watch closed normally (EOF) or unexpected error
22               HandleError(err)
23               return nil
24           }
25           // watchHandler watches w and keeps *resourceVersion up to date.
26           watchHandler(w, &resourceVersion)
27       }
28   }
```

Pseudo code

# Recap: resourceVersion

Everything has a ResourceVersion:
- Individual API object (e.g. a Pod) has ResourceVersion when you GET one
- For a list of API objects (e.g. a PodList)
  - The entire list has a ResourceVersion
  - Each API object in items has ResourceVersion

The ResourceVersion of the top-level list is what should be used when starting a watch to observe events occurring after that list was populated.

# Recap: resourceVersion

- ListOption in List Request
  - Unspecified: etcd
  - RV=0: APIServer cache
  - RV>0: the result is at least as fresh as given RV

- ListOption in Watch Request
  - Unspecified: random time point
  - RV=0: the result is an "ADDED" event for every existing object followed by events for changes that occur after the watch was established
    - (main reason: backwards compatibility-- #13910)
  - Best practice: always specify last listed/watched RV

Watch Event on kube-scheduler, kube-controller-manager, kublet… (roycaihw)

# Kubernetes controller workflow

# Scheduler

```
136    // New returns a Scheduler
137    func New(client clientset.Interface,
138            nodeInformer coreinformers.NodeInformer,
139            podInformer coreinformers.PodInformer,
140            pvInformer coreinformers.PersistentVolumeInformer,
141            pvcInformer coreinformers.PersistentVolumeClaimInformer,
142            replicationControllerInformer coreinformers.ReplicationControllerInformer,
143            replicaSetInformer appsinformers.ReplicaSetInformer,
144            statefulSetInformer appsinformers.StatefulSetInformer,
145            serviceInformer coreinformers.ServiceInformer,
146            pdbInformer policyinformers.PodDisruptionBudgetInformer,
147            storageClassInformer storageinformers.StorageClassInformer,
148            recorder record.EventRecorder,
149            schedulerAlgorithmSource kubeschedulerconfig.SchedulerAlgorithmSource,
150            opts ...func(o *schedulerOptions)) (*Scheduler, error) {
151
```

pkg/scheduler/scheduler.go

Watches:
- Node
- Pod
- PV
- PVC
- RC
- RS
- Stateful set
- Service
- PDB
- Storage class

# Scheduler "business" logic

```
1    func InitEventHandlers(PodInformer) {
2         // scheduled pod cache
3         PodInformer.AddEventHandler(
4              FilteringResourceEventHandler{
5                   FilterFunc: func(obj interface{}) bool {
6                        return IsScheduled(obj.Pod())
7                   }
8                   Handler: {
9                        OnAdd: addPodToCache,
10                       OnUpdate: updatePodInCache,
11                       OnDelete: deletePodFromCache,
12                  }
13             }
14        )
15
16        // unscheduled pod queue
17        PodInformer.AddEventHandler(
18             FilteringResourceEventHandler{
19                  FilterFunc: func(obj interface{}) bool {
20                       return !IsScheduled(obj.Pod())
21                  }
22                  Handler: {
23                       OnAdd: addPodToSchedulingQueue,
24                       OnUpdate: updatePodInSchedulingQueue,
25                       OnDelete: deletePodFromSchedulingQueue,
26                  }
27             }
28        )
29   }
50   // SchedulingQueue is an interface for a queue to store pods waiting to be scheduled.
51   // The interface follows a pattern similar to cache.FIFO and cache.Heap and
52   // makes it easy to use those data structures as a SchedulingQueue.
53   type SchedulingQueue interface {
```

**Pseudo code**

Scheduler keeps in-memory:
- PodCache for scheduled pods
- SchedulingQueue for pods waiting to be scheduled

Having knowledge of the state of world enables scheduler to
- calculate resource usage on nodes
- be aware of available PVs that can be binded with pods
- schedule pods onto nodes

**Watch enables on-demand information propagation, and avoids heavy information-polling loads.**

Key Takeaways

Use Watch for your controller!
- It's trustworthy!
- It's efficient!

**KubeCon** | **CloudNativeCon**

North America 2018