



# Rook: Storage Orchestrator for Kubernetes

Kiefer Chang  
SUSE



# Agenda

Kubernetes Recap

Ceph Introduction

Rook Introduction

Rook Architecture

Demo: Deploy a Ceph cluster with Rook

How to Contribute



# About Me

- Kiefer Chang
- Github: <https://github.com/bk201>
- IRC: kiefer\_chang ({#ceph-dashboard,#ceph-devel}@OFTC)
- Experience
  - Storage development
  - OpenStack development (Nova, Cinder, Ironi, etc.)
  - Ceph development (mainly Dashboard and Orchestrator)
- Work for SUSE (Enterprise Storage Team)

# Kubernetes Recap

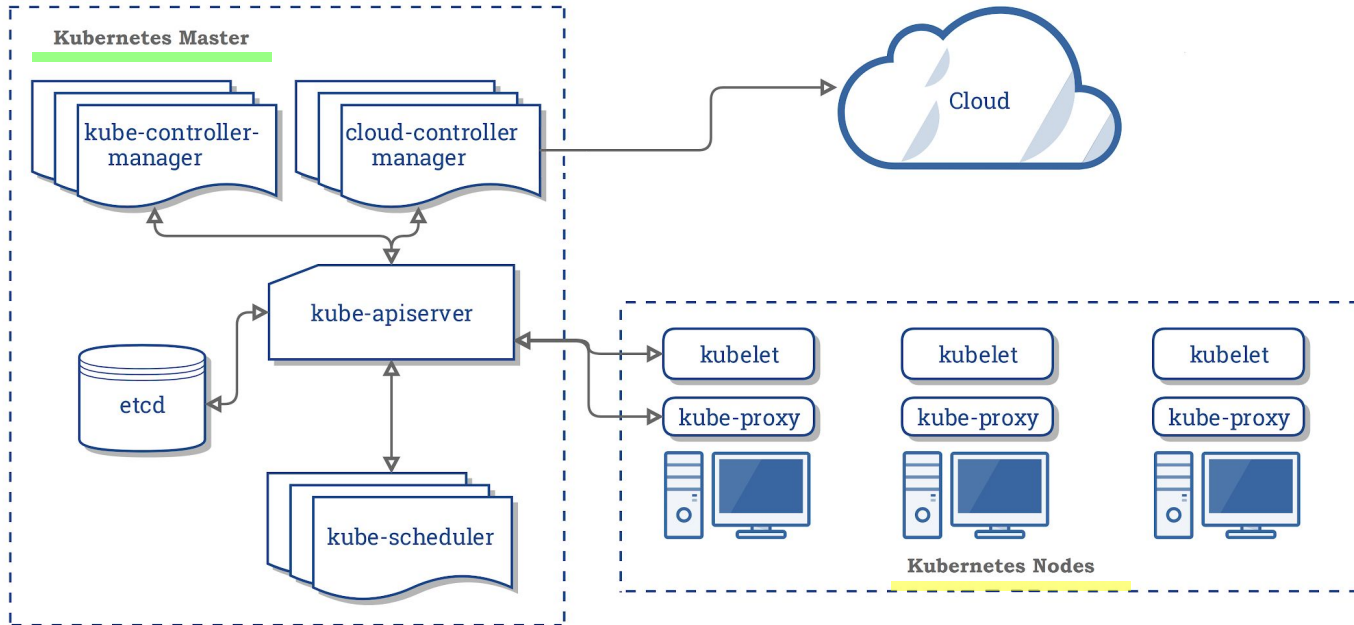
---



# Session Overview

- K8S cluster and components
- Pods and controllers
- Storage
- ConfigMaps
- Secrets
- CRD

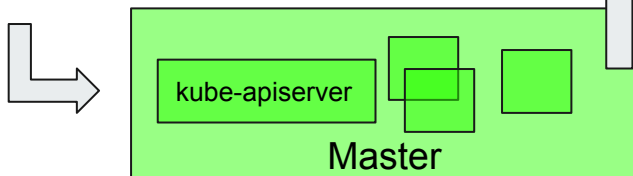
# A Kubernetes Cluster



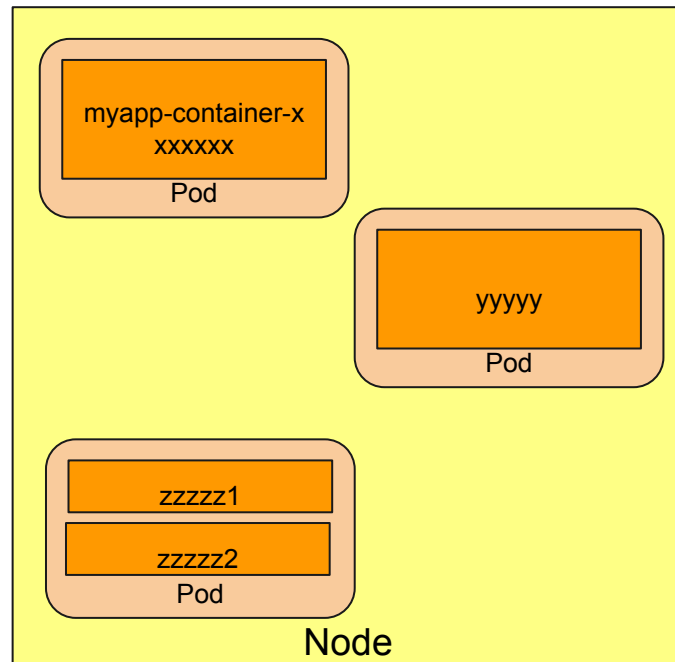
# Manifests and Pods

```
apiVersion: v1
kind: Pod
metadata:
  name: myapp-pod
  labels:
    app: myapp
spec:
  containers:
  - name: myapp-container
    image: busybox
    command: ['sh', '-c', 'echo Hello Kubernetes! && sleep 3600']
```

**Manifests:** Specification of a Kubernetes API object (JSON/YAML)



- **Pods:** smallest unit to be scheduled, contain one or more containers



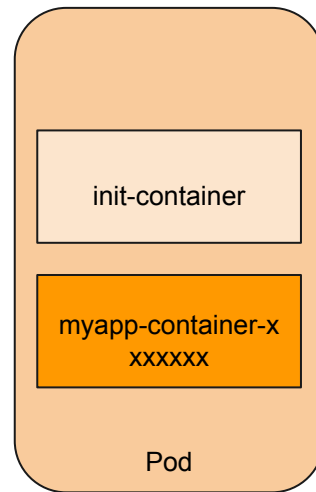
# init Containers

Run before the app containers are started in a Pod

Run to completion and must complete successfully

Useful for

- Register services
- Wait for something
- Utilities







# Pod Controllers

ReplicaSet

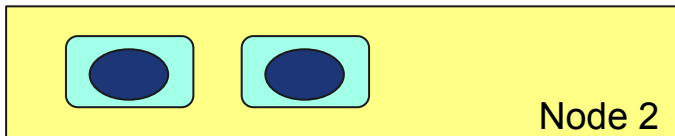
Deployments

DaemonSet

Jobs

# ReplicaSet

maintain a stable set of replica Pods



```
apiVersion: apps/v1
kind: ReplicaSet
metadata:
  name: frontend
  labels:
    app: guestbook
    tier: frontend
spec:
  # modify replicas according to your case
  replicas: 3
  selector:
    matchLabels:
      tier: frontend
  template:
    metadata:
      labels:
        tier: frontend
    spec:
      containers:
        - name: php-redis
          image: gcr.io/google_samples/gb-frontend:v3
```



# Deployment

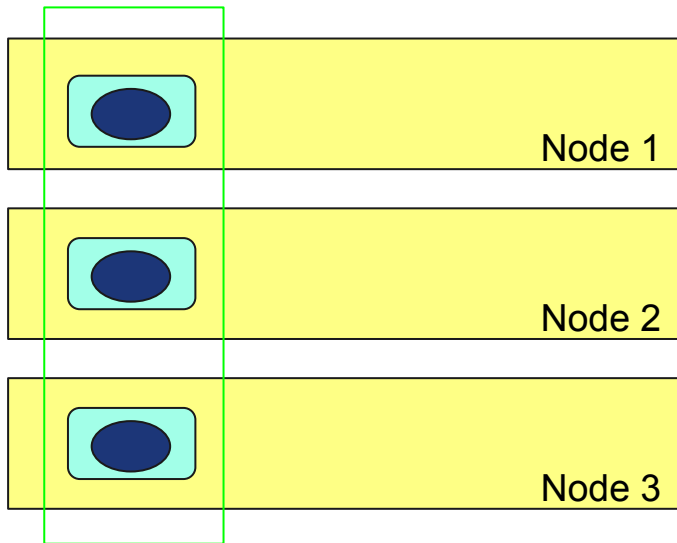
provides declarative updates for [Pods](#) and [ReplicaSets](#)

- Strategy to replace old Pods
  - Recreate
  - Rolling update\*

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: nginx-deployment
  labels:
    app: nginx
spec:
  replicas: 3
  selector:
    matchLabels:
      app: nginx
  template:
    metadata:
      labels:
        app: nginx
    spec:
      containers:
        - name: nginx
          image: nginx:1.7.9
          ports:
            - containerPort: 80
```

# DaemonSet

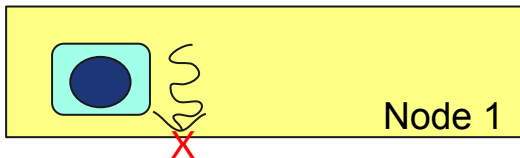
Ensures that all (or some) Nodes run a copy of a Pod



- Storage daemon, such as `glusterd`, `ceph`, on each node.
- Logs collection daemon on every node, such as `fluentd` or `logstash`.
- Monitoring daemon on every node, such as Prometheus Node Exporter..

# Jobs

Creates one or more Pods and ensures that a specified number of them successfully terminate.



```
apiVersion: batch/v1
kind: Job
metadata:
  name: pi
spec:
  template:
    spec:
      containers:
      - name: pi
        image: perl
        command: ["perl", "-Mbignum=bpi", "-wle", "print bpi(2000)"]
        restartPolicy: Never
      backoffLimit: 4
```



# Storage

Volumes

Persistent Volumes

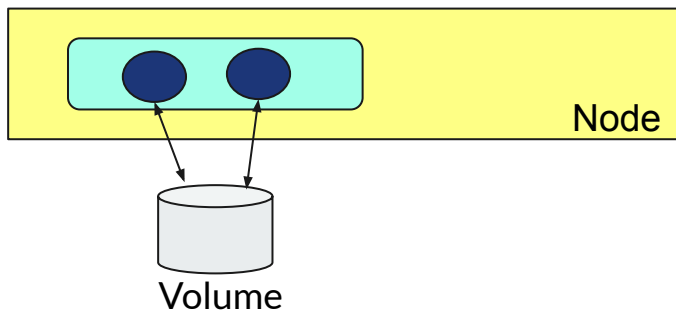
Storage class

Persistent Volume Claim

# Volumes

Data are not lost when Pods are crashed and restarted

Share data between containers in a Pod



## Types of Volumes

Kubernetes supports several types of Volumes:

- [awsElasticBlockStore](#)
- [azureDisk](#)
- [azureFile](#)
- [cephfs](#)
- [cinder](#)
- [configMap](#)
- [csi](#)
- [downwardAPI](#)
- [emptyDir](#)
- [fc \(fibre channel\)](#)
- [flexVolume](#)
- [flocker](#)
- [gcePersistentDisk](#)
- [gitRepo \(deprecated\)](#)
- [glusterfs](#)
- [hostPath](#)
- [iscsi](#)
- [local](#)
- [nfs](#)
- [persistentVolumeClaim](#)
- [projected](#)
- [portworxVolume](#)
- [quobyte](#)
- [rbd](#)
- [scaleIO](#)
- [secret](#)
- [storageos](#)
- [vsphereVolume](#)



# Persistent Volumes

Volume data are not deleted after Pods termination

Provisioned

- by an administrator
- or dynamically

```
apiVersion: v1
kind: PersistentVolume
metadata:
  name: pv0003
spec:
  capacity:
    storage: 5Gi
  volumeMode: Filesystem
  accessModes:
    - ReadWriteOnce
  persistentVolumeReclaimPolicy: Recycle
  storageClassName: slow
  mountOptions:
    - hard
    - nfsvers=4.1
  nfs:
    path: /tmp
    server: 172.17.0.2
```





# Storage Classes

“Classes” of storage

QoS, durability, backup policy..., etc.

```
apiVersion: storage.k8s.io/v1
kind: StorageClass
metadata:
  name: standard
provisioner: kubernetes.io/aws-ebs
parameters:
  type: gp2
reclaimPolicy: Retain
allowVolumeExpansion: true
mountOptions:
  - debug
volumeBindingMode: Immediate
```



# Persistent Volume Claim

A request for storage

```
apiVersion: v1
kind: PersistentVolumeClaim
metadata:
  name: myclaim
spec:
  accessModes:
    - ReadWriteOnce
  volumeMode: Filesystem
  resources:
    requests:
      storage: 8Gi
  storageClassName: slow
  selector:
    matchLabels:
      release: "stable"
    matchExpressions:
      - {key: environment, operator: In, values: [dev]}
```

# How to Claim a PV

```
apiVersion: v1
kind: PersistentVolumeClaim
metadata:
  name: myclaim
spec:
  accessModes:
    - ReadWriteOnce
  volumeMode: Filesystem
  resources:
    requests:
      storage: 8Gi
  storageClassName: standard
  selector:
    matchLabels:
      release: "stable"
    matchExpressions:
      - {key: environment, operator:
In, values: [dev]}
```

PV

```
apiVersion: v1
kind: Pod
metadata:
  name: mypod
spec:
  containers:
    - name: myfrontend
      image: nginx
      volumeMounts:
        - mountPath: "/var/www/html"
          name: mypd
  volumes:
    - name: mypd
      persistentVolumeClaim:
        claimName: myclaim
```

POD

```
apiVersion: storage.k8s.io/v1
kind: StorageClass
metadata:
  name: standard
provisioner: kubernetes.io/aws-ebs
parameters:
  type: gp2
reclaimPolicy: Retain
allowVolumeExpansion: true
mountOptions:
  - debug
volumeBindingMode: Immediate
```

Storage  
Class

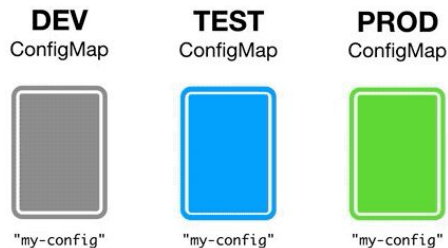
# ConfigMaps

A dictionary of configs (key-value pairs)

Decouple configs and applications

Sources:

- files (ini, JSON, YMAL...), directory
- literals





# Secrets

Store and manage sensitive information, such as passwords, OAuth tokens, and ssh keys

- Outside of Pods
- Stored in etcd

How to use secrets in Pods

- Volumes
- Environment variables

```
apiVersion: v1
kind: Secret
metadata:
  name: mysecret
type: Opaque
data:
  username: YWRtaW4=
  password: MWYyZDFlMmU2N2Rm
```

```
apiVersion: v1
kind: Pod
metadata:
  name: mypod
spec:
  containers:
    - name: mypod
      image: redis
      volumeMounts:
        - name: foo
          mountPath: "/etc/foo"
          readOnly: true
  volumes:
    - name: foo
      secret:
        secretName: mysecret
```

```
apiVersion: v1
kind: Pod
metadata:
  name: secret-env-pod
spec:
  containers:
    - name: mycontainer
      image: redis
      env:
        - name: SECRET_USERNAME
          valueFrom:
            secretKeyRef:
              name: mysecret
              key: username
        - name: SECRET_PASSWORD
          valueFrom:
            secretKeyRef:
              name: mysecret
              key: password
```



# CRD (Custom Resource Definition)

Extends the Kubernetes API or allows you to introduce your own API

```
---
apiVersion: apiextensions.k8s.io/v1beta1
kind: CustomResourceDefinition
metadata:
  name: cephfilesystems.ceph.rook.io
spec:
  group: ceph.rook.io
  names:
    kind: CephFilesystem
    listKind: CephFilesystemList
    plural: cephfilesystems
    singular: cephfilesystem
  scope: Namespaced
  version: v1
  validation:
    openAPIV3Schema:
      properties:
        spec:
          properties:
            metadataServer:
              properties:
                activeCount:
                  minimum: 1
                  maximum: 10
                  type: integer
                activeStandby:
                  type: boolean
              annotations: {}
              placement: {}
              resources: {}
```

# Ceph Introduction

---

Ceph Introduction & How to Contribute

[https://docs.google.com/presentation/d/1-UWfgiXsZjz50j7lYwNwtdAnFzWAbfMb7IAA8JvKwQs/edit#slide=id.g741af3ec10\\_1\\_0](https://docs.google.com/presentation/d/1-UWfgiXsZjz50j7lYwNwtdAnFzWAbfMb7IAA8JvKwQs/edit#slide=id.g741af3ec10_1_0)



# What is Ceph

- An distributed scalable storage
- Open sourced (Mostly LGPL 2.1 & LGPL 3)
- Reliable storage on unreliable hardware
  - No SPoF (Single Point of Failure)
  - Erasure coding or replication for data durability
- Self-healing, auto rebalancing and placement optimization
- Run on commodity hardware
- Enterprise ready





## PREMIER MEMBERS



## GENERAL MEMBERS

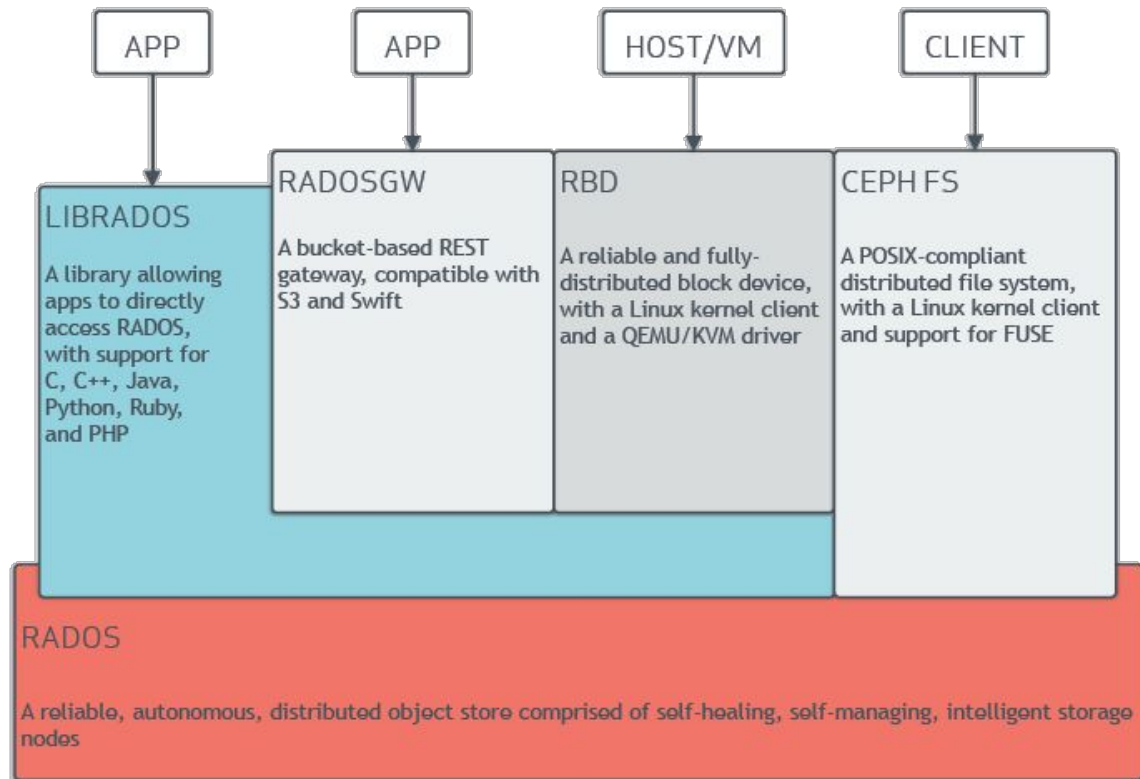


## ASSOCIATE MEMBERS

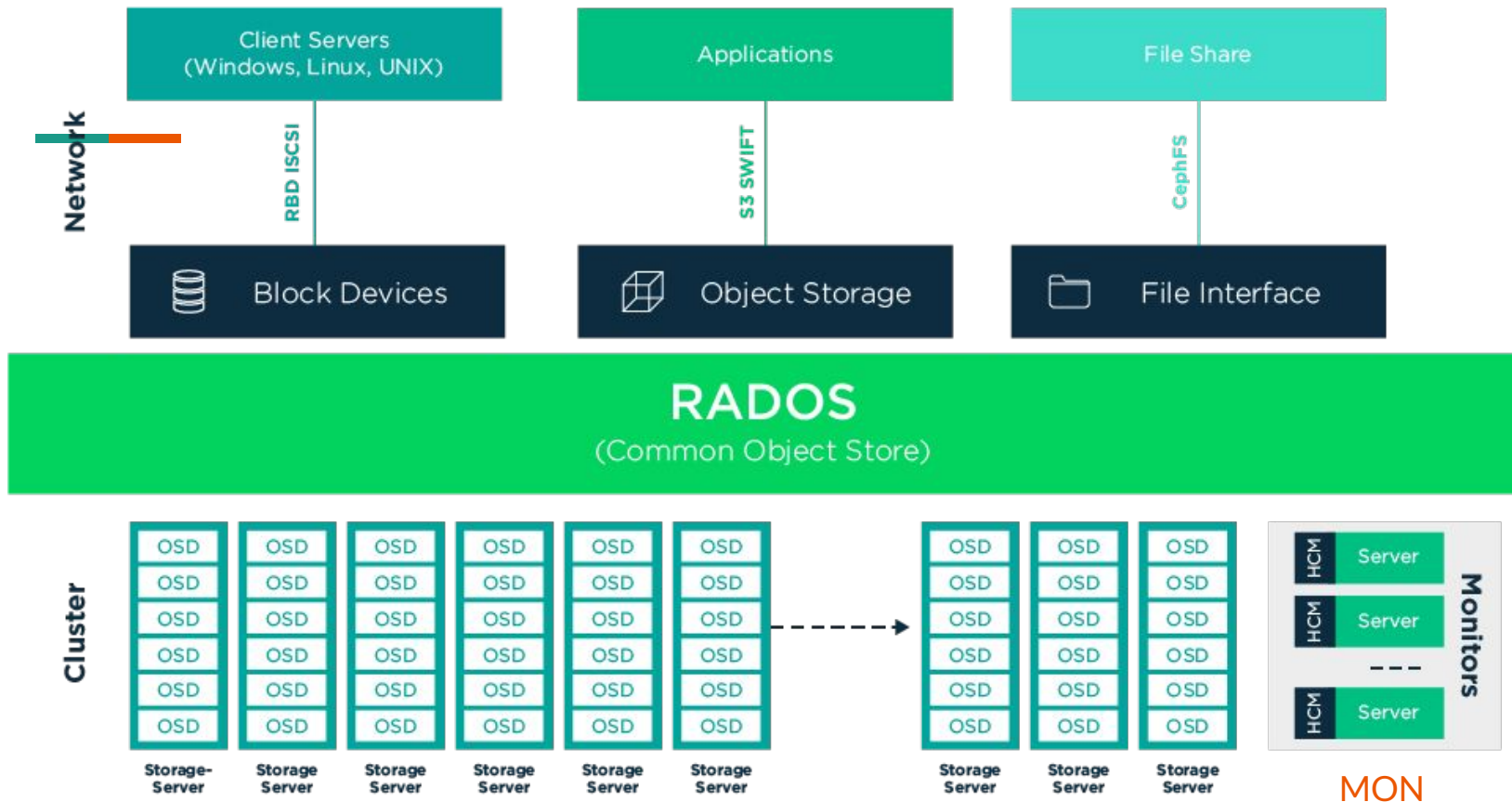


# Ceph is a Unified Storage

- Native protocols
  - Object (radosgw)
  - Block device (RBD)
  - File system (Ceph FS)
- Legacy protocols
  - iSCSI
  - NFS (on CephFS/radosgw)
  - Samba

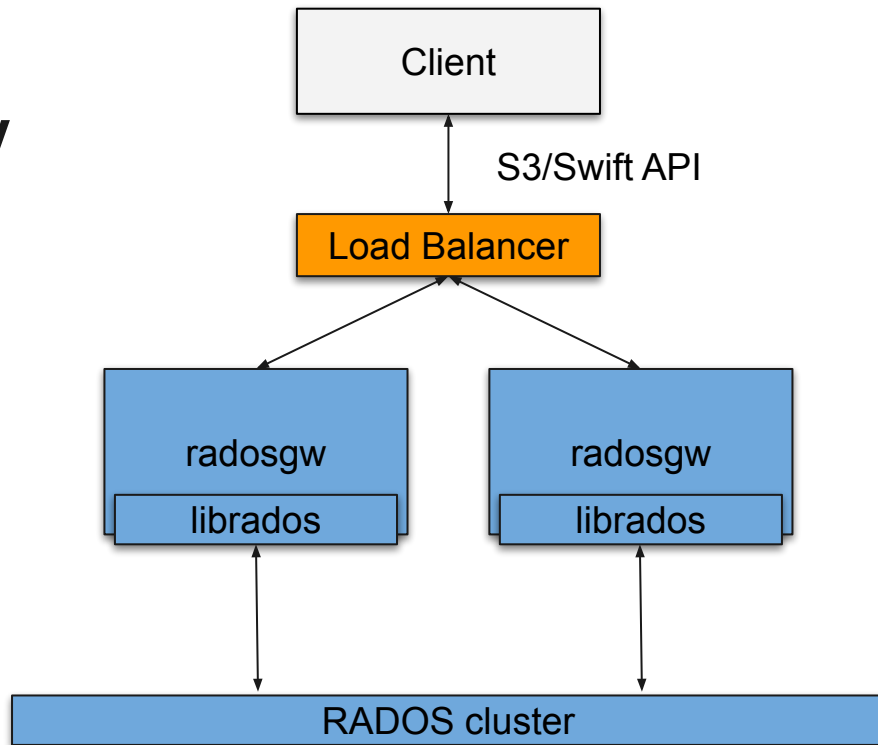


# Ceph Architecture - MONs and OSDs



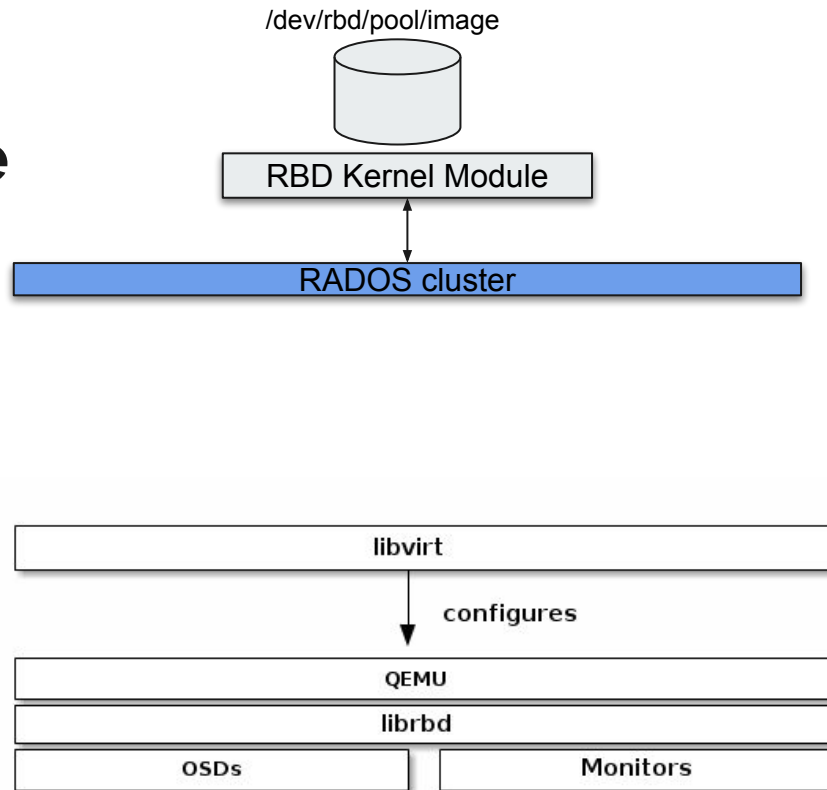
# RGW: RADOS Gateway

- S3/Swift compatible API
  - Users/buckets/objects
  - Encryption
  - Compression
  - Object expiry
  - CORS
- Use cases
  - Public cloud storage
  - Static files
  - VM images
  - Archiving and backup
  - NFS gateway



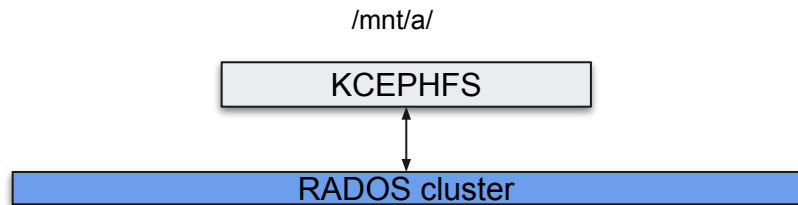
# RBD: RADOS Block Device

- KVM (qemu-kvm) and Linux clients
- Thin-provisioned
- Snapshots and clones
- Cross-sites mirroring
- Use cases
  - VM image/disks
    - For OpenStack components like Nova, Glance, and Cinder
    - Proxmox, Nebula...
  - Kubernetes persistent volume
  - iSCSI gateway



# CephFS: File System

- Shared file system
- Strong consistency
- Scale metadata/data independently
- Snapshots and quotas
- xattrs and lockings features
- Use cases
  - Use as traditional FS, it's nearly POSIX-compatible\*
  - Shared file system for VMs, like OpenStack Manila
  - NFS gateway for windows/others



\* <https://docs.ceph.com/docs/master/cephfs/posix/>

# Dashboard - Summary



English ▾



Dashboard Cluster ▾ Pools Block ▾ NFS Filesystems Object Gateway ▾

Refresh 5 s ▾

## Status

Cluster Status

HEALTH\_OK

Monitors

3 (quorum 0, 1, 2)

OSDs

3 total  
3 up, 3 in

Manager Daemons

1 active  
0 standby

Hosts

2 total

Object Gateways

1 total

Metadata Servers

1 active  
2 standby

ISCSI Gateways

0 total  
0 up, 0 down

## Performance

Client IOPS

0

Client Throughput

0 B/s

Client Read/Write

N/A

Recovery Throughput

0 B/s

Scrub

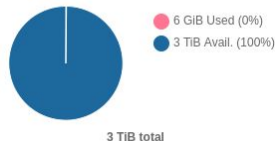
Inactive

## Capacity

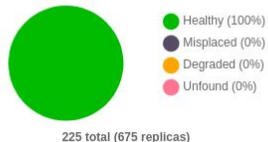
Pools

6

Raw Capacity



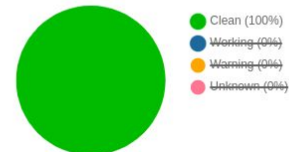
Objects



PGs per OSD

48

PG Status



# Dashboard - Day 2 Operations



English ▾



Dashboard

Cluster ▾

Pools

Block ▾

NFS

Filesystems

Object Gateway ▾

[Pools](#) > Create

## Create Pool

Name \*

data



Pool type \*

erasure



Placement groups \*

128

[Calculation help](#)

The current PGs settings were calculated for you, you should make sure the values suit your needs before submit.

Erasure code profile

default



Flags

☐ EC Overwrites

Applications

No applications added

## Compression

Mode

passive



Algorithm

snappy



Minimum blob size

e.g., 128KiB

Maximum blob size

e.g., 512KiB

Ratio

Compression ratio

## Quotas

Max bytes

e.g., 10GiB

Max objects

0

Create Pool

Cancel



# Dashboard - Monitoring Metrics (Grafana Panels)



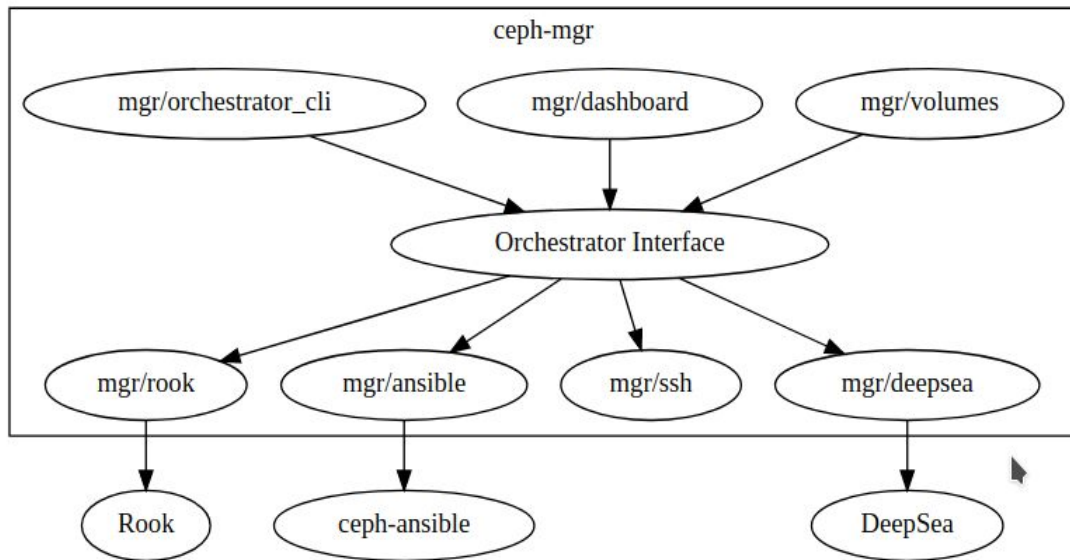


# Deployment Tools

- ceph-deploy
  - <https://docs.ceph.com/docs/master/mgr/ssh/>
  - ssh + commands
  - Likely to be deprecated, moving to ssh orchestrator
- ssh-orchestrator (under active develop, new in Octopus)
  - <https://docs.ceph.com/docs/master/mgr/ssh/>
  - ssh + ceph-daemon (<https://docs.ceph.com/docs/master/bootstrap/#>)
- DeepSea
  - Salt-based
  - <https://github.com/SUSE/DeepSea>
- Ceph-ansible
  - <https://github.com/ceph/ceph-ansible>
  - ansible playbooks
- Rook
  - CRD and operator images to deploy Ceph in Kubernetes
  - <https://github.com/rook/rook/>

# Ceph Orchestrator Module

- Common interface for different orchestrators
  - Let Dashboard/CLI/Utilities work with various different backends



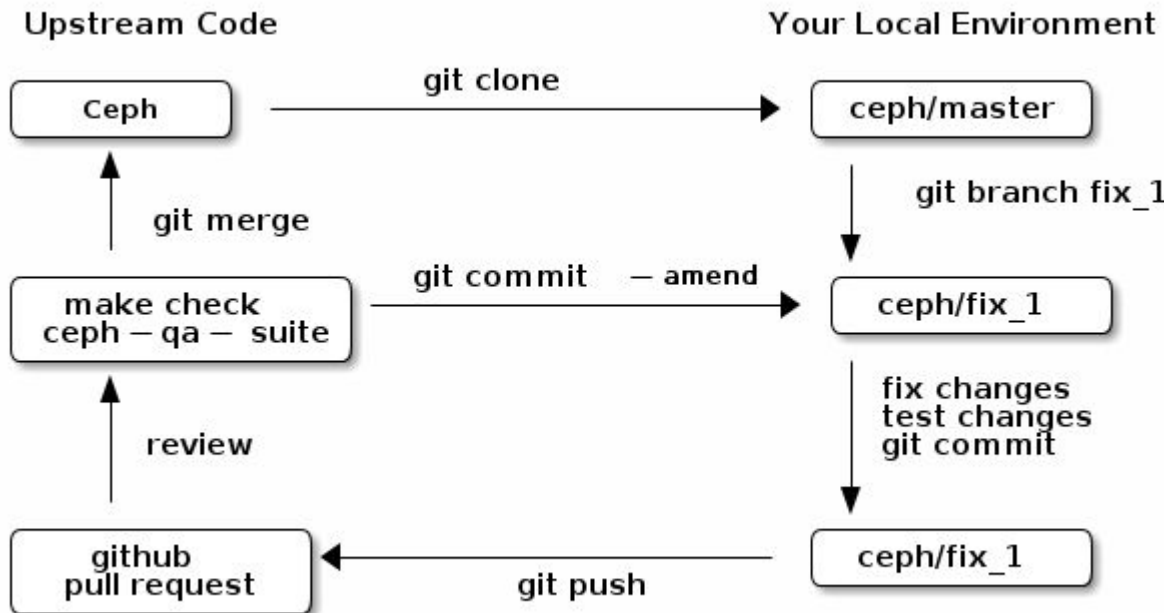
# Contribution Working Flow

## NOTE

*Always create an issue on tracker first if possible*



Jenkins



[https://docs.ceph.com/docs/master/dev/developer\\_guide](https://docs.ceph.com/docs/master/dev/developer_guide)

# Rook Introduction (Ceph-Rook)

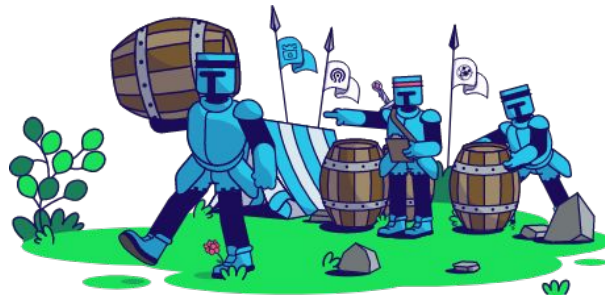
---

# What is Rook

A storage orchestrator for Kubernetes

Leverage Kubernetes' platform power to bootstrap and manage storage in Kubernetes

- Ceph
- CockroachDB
- Yugabyte DB
- EdgeFS
- Cassandra
- Minio
- NFS





# QuickStart: Ceph

TL;DR (If you have a running Kubernetes cluster)

```
git clone https://github.com/rook/rook.git
cd rook
git checkout v1.1.6

cd cluster/examples/kubernetes/ceph
kubectl create -f common.yaml
kubectl create -f operator.yaml
kubectl create -f cluster-test.yaml
```

**\*\* Not for production \*\***



# The Operator Pattern

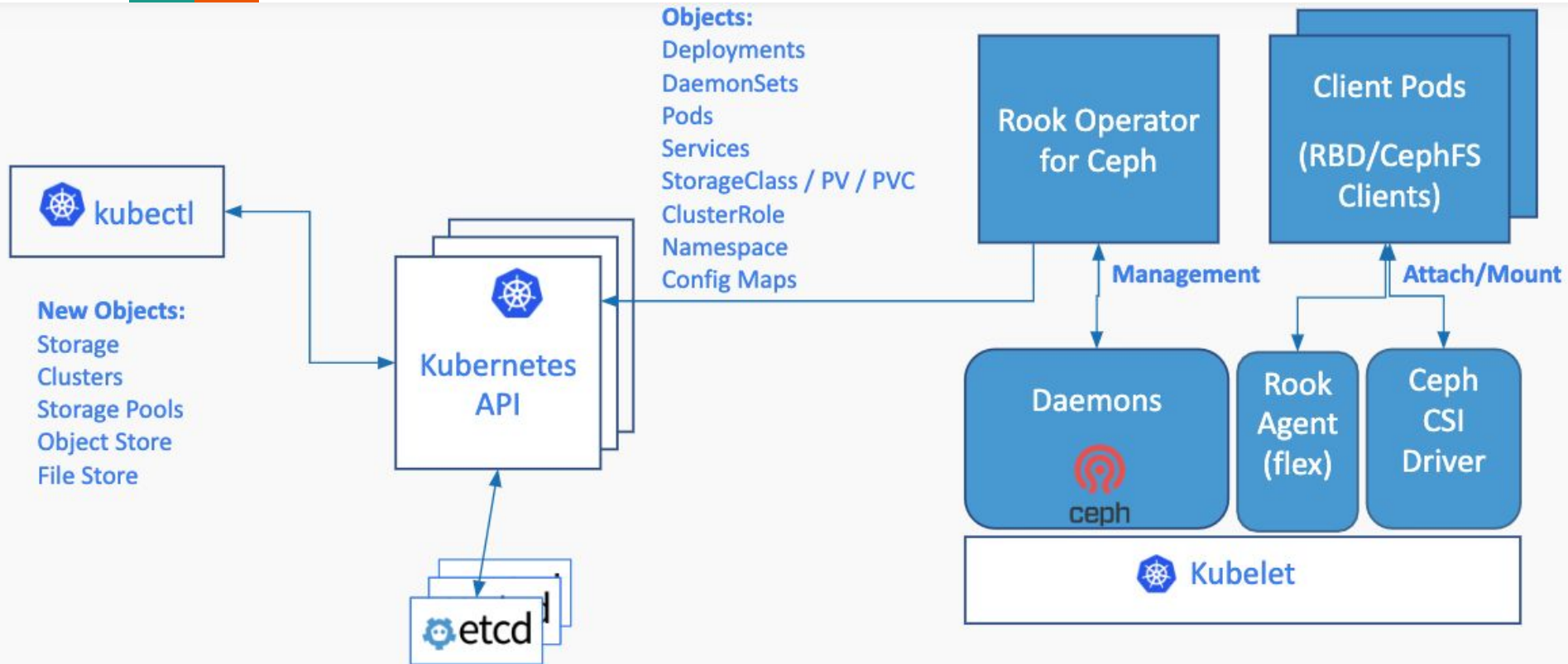
Provision a application in Kubernetes to do the rest of works:

- Bootstrap the Ceph cluster
- Manage Day-2 operations
  - Adding/removing daemons and OSDs
  - Scale out
  - Upgrade

*you still need real operators...*



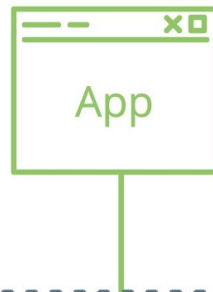
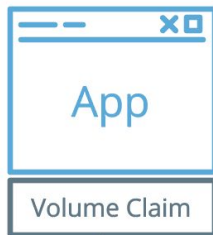
# Rook Architecture



# Rook Architecture



Kubernetes  
Clusters



 ROOK pods

 Agent

 Agent

 Operator

 Agent

 Agent

 Agent

 Mon

 MGR

 MDS

 Mon

 RGW

 RGW

 Mon

 OSD's



 OSD's



 OSD's



 OSD's



 OSD's





# Deploy a Ceph Cluster with Rook

Materials are extracted from Blaine Gardner's presentation: Getting Started as a Rook-Ceph Developer in Cephalocon

Recording: <https://youtu.be/P2ZDIdEPyiw>

Slides: [https://static.sched.com/hosted\\_files/cephalocon2019/39/Intro%20to%20Rook%20Dev.pdf](https://static.sched.com/hosted_files/cephalocon2019/39/Intro%20to%20Rook%20Dev.pdf)



# Create CRDs (cluster/examples/kubernetes/ceph/common.yaml)

```
13 # Namespace where the operator and other rook resources are created
14 apiVersion: v1
15 kind: Namespace
16 metadata:
17   name: rook-ceph
18 # OLM: BEGIN CEPH CRD
19 # The CRD declarations
20 ---
21 apiVersion: apiextensions.k8s.io/v1beta1
22 kind: CustomResourceDefinition
23 metadata:
24   name: cephclusters.ceph.rook.io
25 spec:
26   group: ceph.rook.io
27   names:
28     kind: CephCluster
29     listKind: CephClusterList
30     plural: cephclusters
31     singular: cephcluster
32   scope: Namespaced
33   version: v1
```



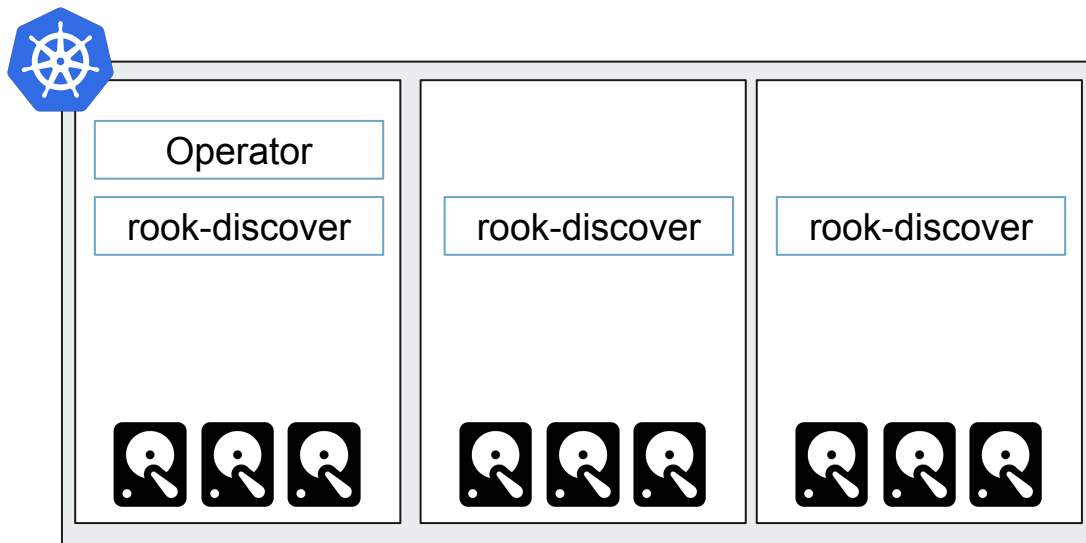
# Deploy the Operator Pod



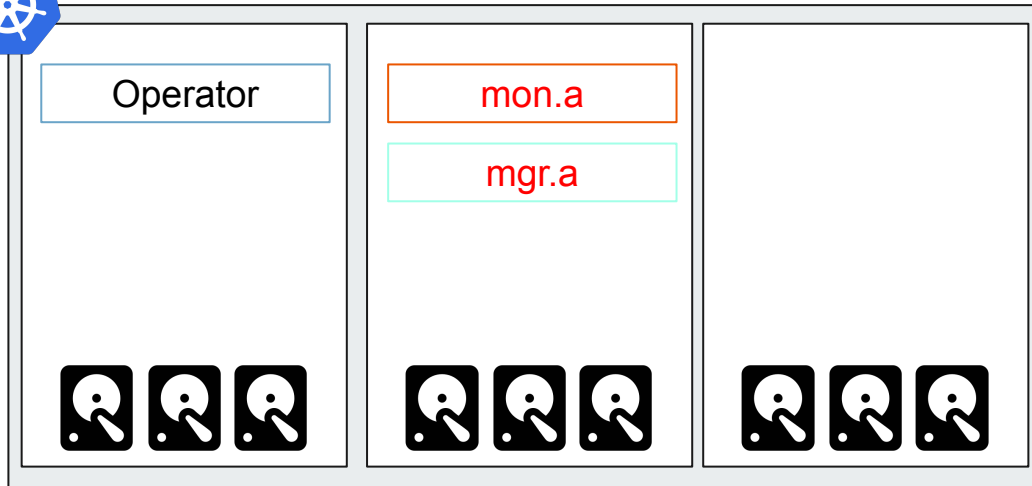
```
12 # OLM: BEGIN OPERATOR DEPLOYMENT
13 apiVersion: apps/v1
14 kind: Deployment
15 metadata:
16   name: rook-ceph-operator
17   namespace: rook-ceph
18   labels:
19     operator: rook
20     storage-backend: ceph
21 spec:
22   selector:
23     matchLabels:
24       app: rook-ceph-operator
25   replicas: 1
26   template:
27     metadata:
28       labels:
29         app: rook-ceph-operator
30   spec:
31     serviceAccountName: rook-ceph-system
32     containers:
33     - name: rook-ceph-operator
34       image: rook/ceph:master
35       args: ["ceph", "operator"]
```

cluster/examples/kubernetes/ceph/operator.yaml

# Operator creates Discover Pods



# Create Ceph Cluster

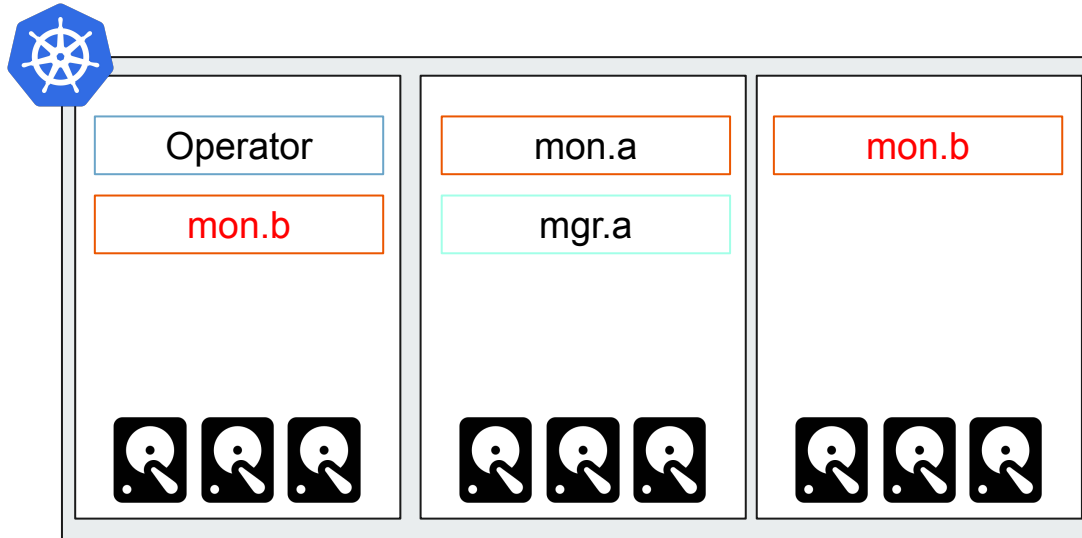


```
12 apiVersion: ceph.rook.io/v1
13 kind: CephCluster
14 metadata:
15   name: rook-ceph
16   namespace: rook-ceph
17 spec:
18   cephVersion:
19     # The container image used to launch the Ceph daemon pods (mon, mgr, osd, mds, rgw).
20     # v13 is mimic, v14 is nautilus, and v15 is octopus.
21     # RECOMMENDATION: In production, use a specific version tag instead of the general v14
22     # versions running within the cluster. See tags available at https://hub.docker.com/r/
23     image: ceph/ceph:v14.2.4-20190917
24     # Whether to allow unsupported versions of Ceph. Currently mimic and nautilus are supp
25     # Octopus is the version allowed when this is set to true.
26     # Do not set to true in production.
27     allowUnsupported: false
28     # The path on the host where configuration files will be persisted. Must be specified.
29     # Important: if you reinstall the cluster, make sure you delete this directory from each
30     # In Minikube, the '/data' directory is configured to persist across reboots. Use "/data
31     dataDirHostPath: /var/lib/rook
32     # Whether or not upgrade should continue even if a check fails
33     # This means Ceph's status could be degraded and we don't recommend upgrading but you mi
34     # Use at your OWN risk
35     # To understand Rook's upgrade process of Ceph, read https://rook.io/docs/rook/master/cep
36     skipUpgradeChecks: false
37     # set the amount of mons to be started
38     mon:
39       count: 3
40       allowMultiplePerNode: false
```

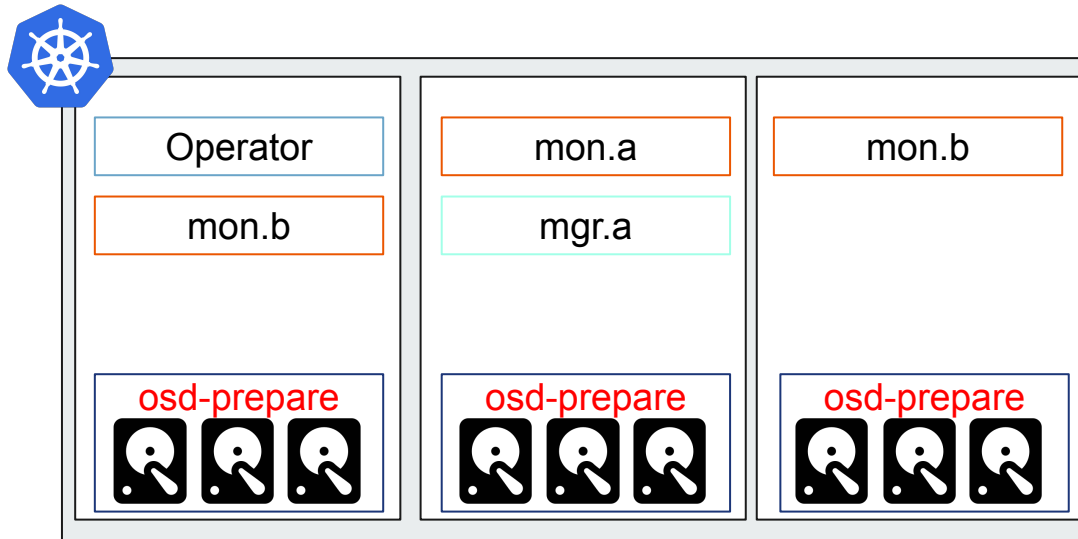
cluster/examples/kubernetes/ceph/cluster.yaml



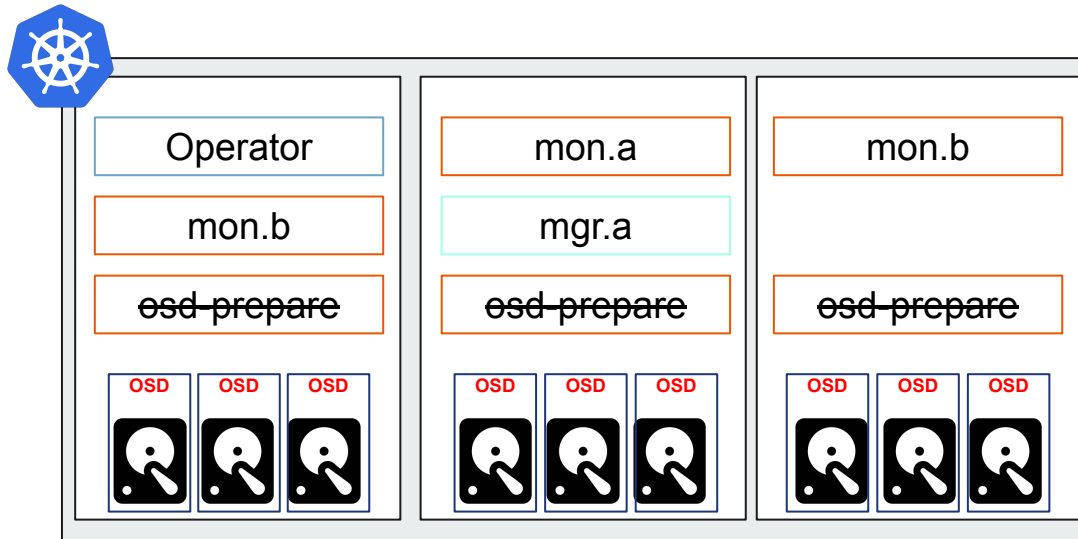
## Create Ceph Cluster (cnt'd)



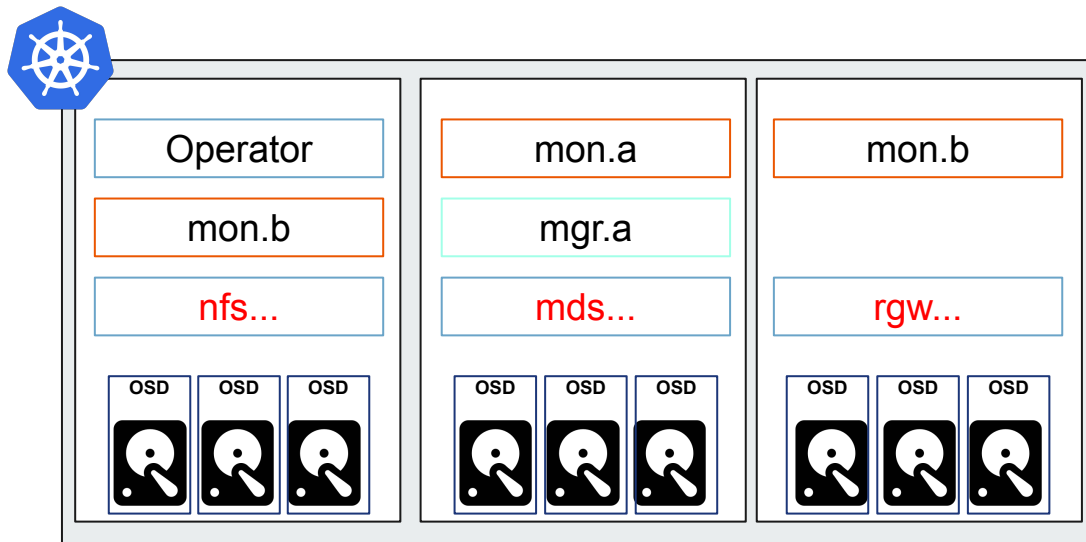
## Create Ceph Cluster (cnt'd)



# Create Ceph Cluster (cnt'd)



# Create Ceph Cluster (cnt'd)





# Claim a RBD Volume

Create a RBD storage class

Claim it in your Pod

```
1  apiVersion: ceph.rook.io/v1
2  kind: CephBlockPool
3  metadata:
4    name: replicapool
5    namespace: rook-ceph
6  spec:
7    failureDomain: host
8    replicated:
9      size: 1
10 ---
11 apiVersion: storage.k8s.io/v1
12 kind: StorageClass
13 metadata:
14   name: rook-ceph-block
15 provisioner: rook-ceph.rbd.csi.ceph.com
16 parameters:
17   # clusterID is the namespace where the rook cluster is running
18   # If you change this namespace, also change the namespace below where the secret namespaces are c
19   clusterID: rook-ceph
20
21   # Ceph pool into which the RBD image shall be created
22   pool: replicapool
23
24   # RBD image format. Defaults to "2".
25   imageFormat: "2"
26
27   # RBD image features. Available for imageFormat: "2". CSI RBD currently supports only `layering`
28   imageFeatures: layering
29
30   # If you change this namespace, also change the namespace below where the secret namespaces are c
```

cluster/examples/kubernetes/ceph/csi/rbd/storageclass-test.yaml

```

14  ---
15  apiVersion: v1
16  kind: PersistentVolumeClaim
17  metadata:
18    name: wp-pv-claim
19    labels:
20      app: wordpress
21  spec:
22    storageClassName: rook-ceph-block
23    accessModes:
24      - ReadWriteOnce
25    resources:
26      requests:
27        storage: 20Gi

```

```

28  ---
29  apiVersion: apps/v1
30  kind: Deployment
31  metadata:
32    name: wordpress
33    labels:
34      app: wordpress
35      tier: frontend
36  spec:
37    selector:
38      matchLabels:
39        app: wordpress
40        tier: frontend
41    strategy:
42      type: Recreate
43    template:
44      metadata:
45        labels:
46          app: wordpress
47          tier: frontend
48      spec:
49        containers:
50          - image: wordpress:4.6.1-apache
51            name: wordpress
52            env:
53              - name: WORDPRESS_DB_HOST
54                value: wordpress-mysql
55              - name: WORDPRESS_DB_PASSWORD
56                value: changeme
57            ports:
58              - containerPort: 80
59              name: wordpress
60            volumeMounts:
61              - name: wordpress-persistent-storage
62                mountPath: /var/www/html
63        volumes:
64          - name: wordpress-persistent-storage
65            persistentVolumeClaim:
66              claimName: wp-pv-claim

```

cluster/examples/kubernetes/wordpress.yaml



# Contribution

Community

 **slack** <https://rook-io.slack.com/> #ceph-dev

 <https://github.com/rook/rook>

Flow: <https://rook.io/docs/rook/v1.1/development-flow.html>

*Fork → Clone → Create local branch → Hack → Push branch → Create PR → Review/CI → Merge*

Most Rook codes are in **Golang**.

Some modules like orchestrator are in Ceph project (**Python**)





# Demo

Deploy a Ceph cluster with Rook

Dashboard & Orchestrator

Persistent Volume Claim with RBD