



Red Hat Ceph Storage 5

Operations Guide

Operational tasks for Red Hat Ceph Storage

Red Hat Ceph Storage 5 Operations Guide

Operational tasks for Red Hat Ceph Storage

Legal Notice

Copyright © 2021 Red Hat, Inc.

The text of and illustrations in this document are licensed by Red Hat under a Creative Commons Attribution–Share Alike 3.0 Unported license ("CC-BY-SA"). An explanation of CC-BY-SA is available at

<http://creativecommons.org/licenses/by-sa/3.0/>

. In accordance with CC-BY-SA, if you distribute this document or an adaptation of it, you must provide the URL for the original version.

Red Hat, as the licensor of this document, waives the right to enforce, and agrees not to assert, Section 4d of CC-BY-SA to the fullest extent permitted by applicable law.

Red Hat, Red Hat Enterprise Linux, the Shadowman logo, the Red Hat logo, JBoss, OpenShift, Fedora, the Infinity logo, and RHCE are trademarks of Red Hat, Inc., registered in the United States and other countries.

Linux[®] is the registered trademark of Linus Torvalds in the United States and other countries.

Java[®] is a registered trademark of Oracle and/or its affiliates.

XFS[®] is a trademark of Silicon Graphics International Corp. or its subsidiaries in the United States and/or other countries.

MySQL[®] is a registered trademark of MySQL AB in the United States, the European Union and other countries.

Node.js[®] is an official trademark of Joyent. Red Hat is not formally related to or endorsed by the official Joyent Node.js open source or commercial project.

The OpenStack[®] Word Mark and OpenStack logo are either registered trademarks/service marks or trademarks/service marks of the OpenStack Foundation, in the United States and other countries and are used with the OpenStack Foundation's permission. We are not affiliated with, endorsed or sponsored by the OpenStack Foundation, or the OpenStack community.

All other trademarks are the property of their respective owners.

Abstract

This document describes how to do operational tasks for Red Hat Ceph Storage. Red Hat is committed to replacing problematic language in our code, documentation, and web properties. We are beginning with these four terms: master, slave, blacklist, and whitelist. Because of the enormity of this endeavor, these changes will be implemented gradually over several upcoming releases. For more details, see our CTO Chris Wright's message

Table of Contents

CHAPTER 1. INTRODUCTION TO THE CEPH ORCHESTRATOR	5
1.1. USE OF THE CEPH ORCHESTRATOR	5
CHAPTER 2. MIGRATING A NON-CONTAINERIZED RED HAT CEPH STORAGE CLUSTER TO A CONTAINERIZED ENVIRONMENT	7
CHAPTER 3. MANAGEMENT OF SERVICES USING THE CEPH ORCHESTRATOR	10
3.1. PLACEMENT SPECIFICATION OF THE CEPH ORCHESTRATOR	10
3.2. DEPLOYING THE CEPH DAEMONS USING THE COMMAND LINE INTERFACE	10
3.3. DEPLOYING THE CEPH DAEMONS ON A SUBSET OF HOSTS USING THE COMMAND LINE INTERFACE	12
3.4. SERVICE SPECIFICATION OF THE CEPH ORCHESTRATOR	13
3.5. DEPLOYING THE CEPH DAEMONS USING THE SERVICE SPECIFICATION	14
CHAPTER 4. MANAGEMENT OF HOSTS USING THE CEPH ORCHESTRATOR	17
4.1. PREREQUISITES	17
4.2. ADDING HOSTS USING THE CEPH ORCHESTRATOR	17
4.3. ADDING MULTIPLE HOSTS USING THE CEPH ORCHESTRATOR	18
4.4. LISTING HOSTS USING THE CEPH ORCHESTRATOR	20
4.5. ADDING LABELS TO HOSTS USING THE CEPH ORCHESTRATOR	20
4.6. REMOVING HOSTS USING THE CEPH ORCHESTRATOR	21
4.7. PLACING HOSTS IN THE MAINTENANCE MODE USING THE CEPH ORCHESTRATOR	24
CHAPTER 5. MANAGEMENT OF MONITORS USING THE CEPH ORCHESTRATOR	26
5.1. CEPH MONITORS	26
5.2. CONFIGURING MONITOR ELECTION STRATEGY	27
5.3. DEPLOYING THE CEPH MONITOR DAEMONS USING THE COMMAND LINE INTERFACE	27
5.4. DEPLOYING THE CEPH MONITOR DAEMONS USING THE SERVICE SPECIFICATION	30
5.5. DEPLOYING THE MONITOR DAEMONS ON SPECIFIC NETWORK USING THE CEPH ORCHESTRATOR	31
5.6. REMOVING THE MONITOR DAEMONS USING THE CEPH ORCHESTRATOR	32
CHAPTER 6. MANAGEMENT OF MANAGERS USING THE CEPH ORCHESTRATOR	34
6.1. PREREQUISITES	34
6.2. DEPLOYING THE MANAGER DAEMONS USING THE CEPH ORCHESTRATOR	34
6.3. REMOVING THE MANAGER DAEMONS USING THE CEPH ORCHESTRATOR	35
CHAPTER 7. MANAGEMENT OF OSDS USING THE CEPH ORCHESTRATOR	37
7.1. CEPH OSDS	37
7.2. CEPH OSD NODE CONFIGURATION	37
7.3. AUTOMATICALLY TUNING OSD MEMORY	37
7.4. LISTING DEVICES FOR CEPH OSD DEPLOYMENT	38
7.5. ZAPPING DEVICES FOR CEPH OSD DEPLOYMENT	40
7.6. DEPLOYING CEPH OSDS ON ALL AVAILABLE DEVICES	41
7.7. DEPLOYING CEPH OSDS ON SPECIFIC DEVICES AND HOSTS	43
7.8. ADVANCED SERVICE SPECIFICATIONS AND FILTERS FOR DEPLOYING OSDS	44
7.9. DEPLOYING CEPH OSDS USING ADVANCED SERVICE SPECIFICATIONS	46
7.10. REMOVING THE OSD DAEMONS USING THE CEPH ORCHESTRATOR	50
7.11. REPLACING THE OSDS USING THE CEPH ORCHESTRATOR	52
7.12. STOPPING THE REMOVAL OF THE OSDS USING THE CEPH ORCHESTRATOR	53
7.13. ACTIVATING THE OSDS USING THE CEPH ORCHESTRATOR	54
7.13.1. Observing the data migration	55
7.14. RECALCULATING THE PLACEMENT GROUPS	55

7.15. USING THE CEPH MANAGER BALANCER MODULE	56
7.16. USING THE CEPH MANAGER CRASH MODULE	59
CHAPTER 8. MANAGEMENT OF MONITORING STACK USING THE CEPH ORCHESTRATOR	64
8.1. DEPLOYING THE MONITORING STACK USING THE CEPH ORCHESTRATOR	64
8.2. REMOVING THE MONITORING STACK USING THE CEPH ORCHESTRATOR	67
CHAPTER 9. BASIC RED HAT CEPH STORAGE CLIENT SETUP	69
9.1. CONFIGURING FILE SETUP ON CLIENT MACHINES	69
9.2. SETTING-UP KEYSRING ON CLIENT MACHINES	69
CHAPTER 10. MANAGEMENT OF MDS SERVICE USING THE CEPH ORCHESTRATOR	71
10.1. PREREQUISITES	71
10.2. DEPLOYING THE MDS SERVICE USING THE COMMAND LINE INTERFACE	71
10.3. DEPLOYING THE MDS SERVICE USING THE SERVICE SPECIFICATION	73
10.4. REMOVING THE MDS SERVICE USING THE CEPH ORCHESTRATOR	75
CHAPTER 11. MANAGEMENT OF CEPH OBJECT GATEWAY USING THE CEPH ORCHESTRATOR	78
11.1. PREREQUISITES	78
11.2. DEPLOYING THE CEPH OBJECT GATEWAY USING THE COMMAND LINE INTERFACE	78
11.3. DEPLOYING THE CEPH OBJECT GATEWAY USING THE SERVICE SPECIFICATION	81
11.4. DEPLOYING A MULTI-SITE CEPH OBJECT GATEWAY USING THE CEPH ORCHESTRATOR	83
11.5. REMOVING THE CEPH OBJECT GATEWAY USING THE CEPH ORCHESTRATOR	88
CHAPTER 12. MANAGEMENT OF NFS-GANESHA GATEWAY USING THE CEPH ORCHESTRATOR	90
12.1. PREREQUISITES	90
12.2. CREATING THE NFS-GANESHA CLUSTER USING THE CEPH ORCHESTRATOR	90
12.3. DEPLOYING THE NFS-GANESHA GATEWAY USING THE COMMAND LINE INTERFACE	92
12.4. DEPLOYING THE NFS-GANESHA GATEWAY USING THE SERVICE SPECIFICATION	93
12.5. UPDATING THE NFS-GANESHA CLUSTER USING THE CEPH ORCHESTRATOR	95
12.6. VIEWING THE NFS-GANESHA CLUSTER INFORMATION USING THE CEPH ORCHESTRATOR	97
12.7. FETCHING THE NFS-GANESHA CLUSTER LOGS USING THE CEPH ORCHESTRATOR	98
12.8. SETTING CUSTOM NFS-GANESHA CONFIGURATION USING THE CEPH ORCHESTRATOR	99
12.9. RESETTING CUSTOM NFS-GANESHA CONFIGURATION USING THE CEPH ORCHESTRATOR	103
12.10. DELETING THE NFS-GANESHA CLUSTER USING THE CEPH ORCHESTRATOR	104
12.11. REMOVING THE NFS-GANESHA GATEWAY USING THE CEPH ORCHESTRATOR	105
CHAPTER 13. MANAGEMENT OF ISCSI GATEWAY USING THE CEPH ORCHESTRATOR	107
13.1. PREREQUISITES	107
13.2. DEPLOYING THE ISCSI GATEWAY USING THE COMMAND LINE INTERFACE	107
13.3. DEPLOYING THE ISCSI GATEWAY USING THE SERVICE SPECIFICATION	108
13.4. REMOVING THE ISCSI GATEWAY USING THE CEPH ORCHESTRATOR	110
CHAPTER 14. HANDLING A NODE FAILURE	112
14.1. PREREQUISITES	113
14.2. CONSIDERATIONS BEFORE ADDING OR REMOVING A NODE	113
14.3. PERFORMANCE CONSIDERATIONS	113
14.4. RECOMMENDATIONS FOR ADDING OR REMOVING NODES	114
14.5. ADDING A CEPH OSD NODE	115
14.6. REMOVING A CEPH OSD NODE	117
14.7. SIMULATING A NODE FAILURE	118
CHAPTER 15. HANDLING A DATA CENTER FAILURE	120
15.1. PREREQUISITES	120
15.2. AVOIDING A DATA CENTER FAILURE	120

15.3. HANDLING A DATA CENTER FAILURE

120

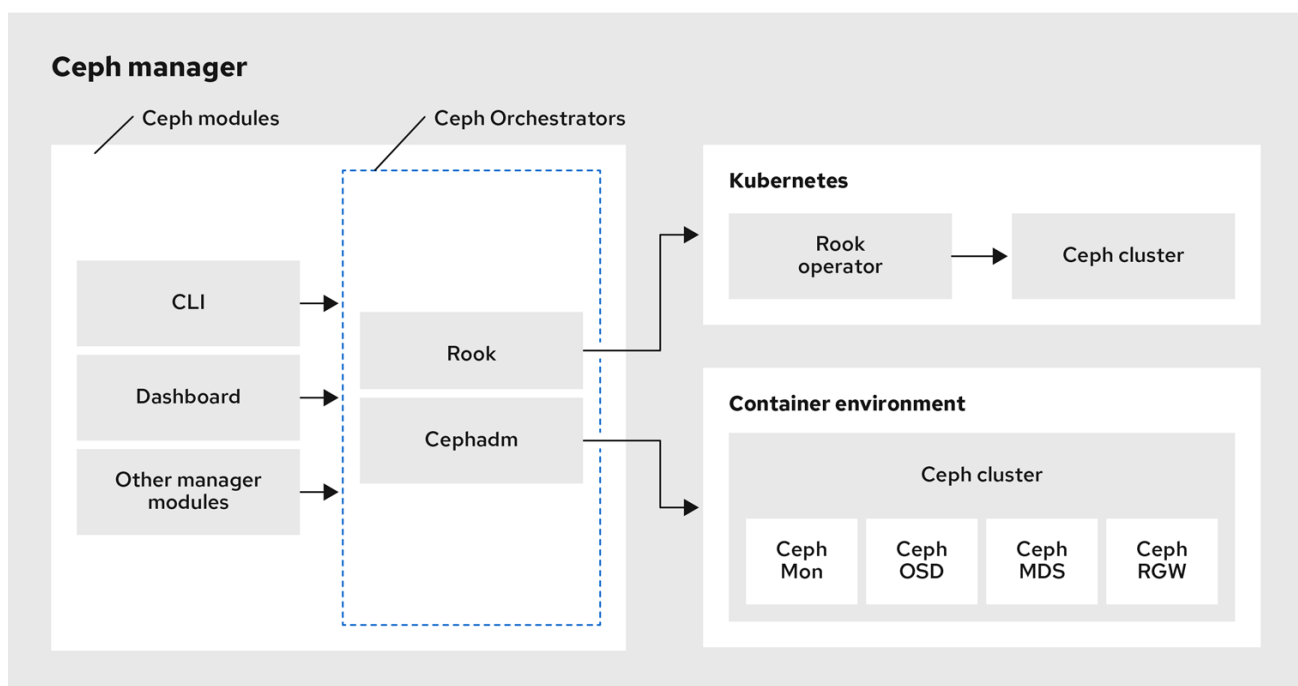
CHAPTER 1. INTRODUCTION TO THE CEPH ORCHESTRATOR

As a storage administrator, you can use the Ceph Orchestrator with Cephadm utility that provides the ability to discover devices and create services in a Red Hat Ceph Storage cluster.

1.1. USE OF THE CEPH ORCHESTRATOR

Red Hat Ceph Storage Orchestrators are manager modules that primarily act as a bridge between a Red Hat Ceph Storage cluster and deployment tools like Rook and Cephadm for a unified experience. They also integrate with the Ceph command line interface and Ceph Dashboard.

The following is a workflow diagram of Ceph Orchestrator:



153_Ceph_0421

Types of Red Hat Ceph Storage Orchestrators

There are three main types of Red Hat Ceph Storage Orchestrators:

- **Orchestrator CLI**: These are common APIs used in Orchestrators and include a set of commands that can be implemented. These APIs also provide a common command line interface (CLI) to orchestrate **ceph-mgr** modules with external orchestration services. The following are the nomenclature used with the Ceph Orchestrator:
 - **Host**: This is the host name of the physical host and not the pod name, DNS name, container name, or host name inside the container.
 - **Service type**: This is the type of the service, such as nfs, mds, osd, mon, rgw, mgr, and iscsi.
 - **Service**: A functional service provided by a Ceph storage cluster such as monitors service, managers service, OSD services, Ceph Object Gateway service, and NFS service.
 - **Daemon**: A specific instance of a service deployed by one or more hosts such as Ceph Object Gateway services can have different Ceph Object Gateway daemons running in three different hosts.

- **Cephadm Orchestrator** - This is a Ceph Orchestrator module that does not rely on an external tool such as Rook or Ansible, but rather manages nodes in a cluster by establishing an SSH connection and issuing explicit management commands. This module is intended for day-one and day-two operations.

Using the Cephadm Orchestrator is the recommended way of installing a Ceph storage cluster without leveraging any deployment frameworks like Ansible. The idea is to provide the manager daemon with access to an SSH configuration and key that is able to connect to all nodes in a cluster to perform any management operations, like creating an inventory of storage devices, deploying and replacing OSDs, or starting and stopping Ceph daemons. In addition, the Cephadm Orchestrator will deploy container images managed by **systemd** in order to allow independent upgrades of co-located services.

This orchestrator will also likely highlight a tool that encapsulates all necessary operations to manage the deployment of container image based services on the current host, including a command that bootstraps a minimal cluster running a Ceph Monitor and a Ceph Manager.

- **Rook Orchestrator** - Rook is an orchestration tool that uses the Kubernetes Rook operator to manage a Ceph storage cluster running inside a Kubernetes cluster. The rook module provides integration between Ceph's Orchestrator framework and Rook. Rook is an open source cloud-native storage operator for Kubernetes.

Rook follows the "operator" model, in which a custom resource definition (CRD) object is defined in Kubernetes to describe a Ceph storage cluster and its desired state, and a rook operator daemon is running in a control loop that compares the current cluster state to desired state and takes steps to make them converge. The main object describing Ceph's desired state is the Ceph storage cluster CRD, which includes information about which devices should be consumed by OSDs, how many monitors should be running, and what version of Ceph should be used. Rook defines several other CRDs to describe RBD pools, CephFS file systems, and so on.

The Rook Orchestrator module is the glue that runs in the **ceph-mgr** daemon and implements the Ceph orchestration API by making changes to the Ceph storage cluster in Kubernetes that describe desired cluster state. A Rook cluster's **ceph-mgr** daemon is running as a Kubernetes pod, and hence, the rook module can connect to the Kubernetes API without any explicit configuration.

CHAPTER 2. MIGRATING A NON-CONTAINERIZED RED HAT CEPH STORAGE CLUSTER TO A CONTAINERIZED ENVIRONMENT

Red Hat Ceph Storage 5 supports only containerized daemons. It does not support non-containerized storage clusters. If you are upgrading a non-containerized Bare-metal storage cluster from Red Hat Ceph Storage 4 to Red Hat Ceph Storage 5, the upgrade process includes the conversion to a containerized deployment. You can manually migrate a non-containerized, Bare-metal, cluster to a containerized cluster using the command line interface.

Prerequisites

- A running Red Hat Ceph Storage non-containerized cluster.
- Root-level access to all the nodes.

Procedure

1. Optional: For two-way RBD mirroring configured using the command-line interface in a bare-metal storage cluster, the cluster does not migrate RBD mirroring. For such a configuration, follow the below steps before migrating the non-containerized storage cluster to a containerized storage cluster:
 - a. Create a user on the Ceph client node:

Syntax

```
ceph auth get client.PRIMARY_CLUSTER_NAME -o
/etc/ceph/ceph.PRIMARY_CLUSTER_NAME.keyring
```

Example

```
[root@rbd-client-site-a ~]# ceph auth get client.rbd-mirror.site-a -o
/etc/ceph/ceph.client.rbd-mirror.site-a.keyring
```

- b. Change the username in the **auth** file in **/etc/ceph** directory:

Example

```
[client.rbd-mirror.rbd-client-site-a]
key = AQCbKbVg+E7POBAA7COSZCodvOrg2LWIFc9+3g==
caps mds = "allow *"
caps mgr = "allow *"
caps mon = "allow *"
caps osd = "allow *"
```

- c. Import the **auth** file to add relevant permissions:

Syntax

```
ceph auth import -i PATH_TO_KEYRING
```

Example

```
[root@rbd-client-site-a ~]# ceph auth import -i /etc/ceph/ceph.client.rbd-mirror.rbd-client-site-a.keyring
```

- d. Check the service name of the RBD mirror node:

Example

```
[root@rbd-client-site-a ~]# systemctl list-units --all

systemctl stop ceph-rbd-mirror@rbd-client-site-a.service
systemctl disable ceph-rbd-mirror@rbd-client-site-a.service
systemctl reset-failed ceph-rbd-mirror@rbd-client-site-a.service
systemctl start ceph-rbd-mirror@rbd-mirror.rbd-client-site-a.service
systemctl enable ceph-rbd-mirror@rbd-mirror.rbd-client-site-a.service
systemctl status ceph-rbd-mirror@rbd-mirror.rbd-client-site-a.service
```

- e. Add the rbd-mirror node to the **/etc/ansible/hosts** file:

Example

```
[rbdmirrors]
ceph.client.rbd-mirror.rbd-client-site-a
```

2. Edit the **group_vars/all.yml** file to include configuration for containers:

```
ceph_docker_image_tag: "latest"
ceph_docker_image: rhceph/rhceph-4-rhel8
containerized_deployment: true
```

3. Navigate to the **/usr/share/ceph-ansible** directory:

```
[ansible@admin ~]$ cd /usr/share/ceph-ansible
```

4. On the Ansible administration node, run the Ansible migration playbook:

Syntax

```
ansible-playbook ./infrastructure-playbooks/switch-from-non-containerized-to-containerized-
ceph-daemons.yml -i INVENTORY_FILE
```

Example

```
[ansible@admin ceph-ansible]$ ansible-playbook ./infrastructure-playbooks/switch-from-non-
containerized-to-containerized-ceph-daemons.yml -i hosts
```

Verify the cluster is switched to containerized environment.

5. On the monitor node, list all running containers:

Example

```
[root@mon ~]$ sudo podman ps
```

Additional Resources

- See the [Configuring two-way mirroring using the command-line interface](#) section in the *Red Hat Ceph Storage Block Device Guide* for more details.
- See the [Red Hat Ceph Storage Installation Guide](#) for more details on installation.

CHAPTER 3. MANAGEMENT OF SERVICES USING THE CEPH ORCHESTRATOR

As a storage administrator, after installing the Red Hat Ceph Storage cluster, you can monitor and manage the services in a storage cluster using the Ceph Orchestrator. A service is a group of daemons that are configured together.

This section covers the following administrative information:

- [Placement specification of the Ceph Orchestrator](#) .
- [Deploying the Ceph daemons using the command line interface](#) .
- [Deploying the Ceph daemons on a subset of hosts using the command line interface](#) .
- [Service specification of the Ceph Orchestrator](#) .
- [Deploying the Ceph daemons using the service specification](#) .

3.1. PLACEMENT SPECIFICATION OF THE CEPH ORCHESTRATOR

You can use the Ceph Orchestrator to deploy **osds**, **mons**, **mgrs**, **mds** and **rgw**, and **iSCSI** services. Red Hat recommends deploying services using placement specifications. You need to know where and how many daemons have to be deployed to deploy a service using the Ceph Orchestrator. Placement specifications can either be passed as command line arguments or as a service specification in a **yaml** file.

There are two ways of deploying the services using the placement specification:

- Using the placement specification directly in the command line interface. For example, if you want to deploy three monitors on the hosts, running the following command deploys three monitors on **host01**, **host02**, and **host03**.

Example

```
[ceph: root@host01 /]# ceph orch apply mon --placement="3 host01 host02 host03"
```

- Using the placement specification in the YAML file. For example, if you want to deploy **node-exporter** on all the hosts, then you can specify the following in the **yaml** file.

Example

```
service_type: mon
placement:
host_pattern: '*'
```

3.2. DEPLOYING THE CEPH DAEMONS USING THE COMMAND LINE INTERFACE

Using the Ceph Orchestrator, you can deploy the daemons such as Ceph Manager, Ceph Monitors, Ceph OSDs, monitoring stack, and others using the **ceph orch** command. Placement specification is passed as **--placement** argument with the Orchestrator commands.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the storage cluster.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Use one of the following methods to deploy the daemons on the hosts:

Method 1

- Specify the number of daemons and the host names:

Syntax

```
ceph orch apply SERVICE_NAME --placement="NUMBER_OF_DAEMONS  
HOST_NAME_1 HOST_NAME_2 HOST_NAME_3"
```

Example

```
[ceph: root@host01 /]# ceph orch apply mon --placement="2 host01 host02 host03"
```

- Add the labels to the hosts and then deploy the daemons using the labels:
 - a. Add the labels to the hosts:

Syntax

```
ceph orch host label add HOSTNAME_1 LABEL_1,LABEL_2
```

Example

```
[ceph: root@host01 /]# ceph orch host label add host01 mon
```

- b. Deploy the daemons with labels:

Syntax

```
ceph orch apply DAEMON_NAME label: _LABEL_1_
```

Example

```
ceph orch apply mon label:mon
```

Method 3

- Add the labels to the hosts and deploy using the **--placement** argument:

- a. Add the labels to the hosts:

Syntax

```
ceph orch host label add HOSTNAME_1 LABEL_1,LABEL_2
```

Example

```
[ceph: root@host01 /]# ceph orch host label add host01 mon
```

- b. Deploy the daemons using the label placement specification:

Syntax

```
ceph orch apply DAEMON_NAME --placement="label:_LABEL_1_"
```

Example

```
ceph orch apply mon --placement="label:mon"
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME  
ceph orch ps --service_name=SERVICE_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=mon  
[ceph: root@host01 /]# ceph orch ps --service_name=mon
```

Additional Resources

- See the [Adding hosts using the Ceph Orchestrator](#) section in the *Red Hat Ceph Storage Operations Guide*.

3.3. DEPLOYING THE CEPH DAEMONS ON A SUBSET OF HOSTS USING THE COMMAND LINE INTERFACE

You can use the **--placement** option to deploy daemons on a subset of hosts. You can specify the number of daemons in the placement specification with the name of the hosts to deploy the daemons.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. List the hosts on which you want to deploy the Ceph daemons:

Example

```
[ceph: root@host01 /]# ceph orch host ls
```

3. Deploy the daemons:

Syntax

```
ceph orch apply SERVICE_NAME --placement="NUMBER_OF_DAEMONS HOST_NAME_1
_HOST_NAME_2 HOST_NAME_3"
```

Example

```
ceph orch apply mgr --placement="2 host01 host02 host03"
```

In this example, the **mgr** daemons are deployed only on two hosts.

Verification

- List the hosts:

Example

```
[ceph: root@host01 /]# ceph orch host ls
```

Additional Resources

- See the [Listing hosts using the Ceph Orchestrator](#) section in the *Red Hat Ceph Storage Operations Guide*.

3.4. SERVICE SPECIFICATION OF THE CEPH ORCHESTRATOR

A service specification is a data structure to specify the service attributes and configuration settings that is used to deploy the Ceph service. The following is an example of the multi-document YAML file, **cluster.yml**, for specifying service specifications:

Example

```
service_type: mon
placement:
  host_pattern: "mon*"
---
service_type: mgr
placement:
  host_pattern: "mgr*"
---
service_type: osd
service_id: default_drive_group
placement:
  host_pattern: "osd*"
data_devices:
  all: true
```

The following list are the parameters where the properties of a service specification are defined as follows:

- **service_type**: The type of service:
 - Ceph services like mon, crash, mds, mgr, osd, rbd, or rbd-mirror.
 - Ceph gateway like nfs or rgw.
 - Monitoring stack like Alertmanager, Prometheus, Grafana or Node-exporter.
 - Container for custom containers.
- **service_id**: A unique name of the service.
- **placement**: This is used to define where and how to deploy the daemons.
- **unmanaged**: If set to **true**, the Orchestrator will neither deploy nor remove any daemon associated with this service.

Stateless service of Orchestrators

A stateless service is a service that does not need information of the state to be available. For example, to start an **rgw** service, additional information is not needed to start or run the service. The **rgw** service does not create information about this state in order to provide the functionality. Regardless of when the **rgw** service starts, the state is the same.

3.5. DEPLOYING THE CEPH DAEMONS USING THE SERVICE SPECIFICATION

Using the Ceph Orchestrator, you can deploy daemons such as ceph Manager, Ceph Monitors, Ceph OSDs, monitoring stack, and others using the service specification in a YAML file.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Navigate to the following directory:

Syntax

```
cd /var/lib/ceph/DAEMON_PATH/
```

Example

```
[ceph: root@host01 /]# cd /var/lib/ceph/mon/
```

3. Create the **yml** file:

Example

```
[ceph: root@host01 mon]# touch mon.yml
```

4. This file can be configured in two different ways:

- Edit the file to include the host details in placement specification:

Syntax

```
service_type: SERVICE_NAME
placement:
  hosts:
    - HOST_NAME_1
    - HOST_NAME_2
```

Example

```
service_type: mon
placement:
  hosts:
    - host01
    - host02
    - host03
```

- Edit the file to include the label details in placement specification:

Syntax

```
■
```

```
service_type: SERVICE_NAME
placement:
  label: "_LABEL_1"
```

Example

```
service_type: mon
placement:
  label: "mon"
```

5. Deploy the Ceph daemons using service specification:

Syntax

```
ceph orch apply -i FILE_NAME.yml
```

Example

```
[ceph: root@host01 mon]# ceph orch apply -i mon.yml
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=mon
```

Additional Resources

- See the [Listing hosts using the Ceph Orchestrator](#) section in the *Red Hat Ceph Storage Operations Guide*.

CHAPTER 4. MANAGEMENT OF HOSTS USING THE CEPH ORCHESTRATOR

As a storage administrator, you can use the Ceph Orchestrator with Cephadm in the backend to add, list, and remove hosts in an existing Red Hat Ceph Storage cluster.

You can also add labels to hosts. Labels are free-form and have no specific meanings. Each host can have multiple labels. For example, apply the **mon** label to all hosts that have monitor daemons deployed, **mgr** for all hosts with manager daemons deployed, **rgw** for Ceph object gateways, and so on.

Labeling all the hosts in the storage cluster helps to simplify system management tasks by allowing you to quickly identify the daemons running on each host. In addition, you can use the Ceph Orchestrator or a YAML file to deploy or remove daemons on hosts that have specific host labels.

This section covers the following administrative tasks:

- [Adding hosts using the Ceph Orchestrator.](#)
- [Adding multiple hosts using the Ceph Orchestrator.](#)
- [Listing hosts using the Ceph Orchestrator.](#)
- [Adding labels to hosts using the Ceph Orchestrator.](#)
- [Removing hosts using the Ceph Orchestrator.](#)
- [Placing hosts in the maintenance mode using the Ceph Orchestrator.](#)

4.1. PREREQUISITES

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.
- The IP addresses of the new hosts should be updated in `/etc/hosts` file.

4.2. ADDING HOSTS USING THE CEPH ORCHESTRATOR

You can use the Ceph Orchestrator with Cephadm in the backend to add hosts to an existing Red Hat Ceph Storage cluster.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Extract the cluster's public SSH keys to a folder:

Syntax

```
ceph cephadm get-pub-key > ~/PATH
```

Example

```
[ceph: root@host01 /]# ceph cephadm get-pub-key > ~/ceph.pub
```

3. Copy Ceph cluster's public SSH keys to the root user's **authorized_keys** file on the new host:

Syntax

```
ssh-copy-id -f -i ~/PATH root@HOST_NAME_2
```

Example

```
[root@host01 ~]# ssh-copy-id -f -i ~/ceph.pub root@host02
```

4. Add hosts to the cluster:

Syntax

```
ceph orch host add HOST_NAME IP_ADDRESS_OF_HOST
```

Example

```
[ceph: root@host01 /]# ceph orch host add host02 10.69.265.25
```

Verification

- List the hosts:

Example

```
[ceph: root@host01 /]# ceph orch host ls
```

Additional Resources

- See the [Listing hosts using the Ceph Orchestrator](#) section in the *Red Hat Ceph Storage Operations Guide*.

4.3. ADDING MULTIPLE HOSTS USING THE CEPH ORCHESTRATOR

You can use the Ceph Orchestrator with Cephadm in the backend to add multiple hosts to a Red Hat Ceph Storage cluster at the same time using the service specification in YAML file format.

Prerequisites

- A running Red Hat Ceph Storage cluster.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Navigate to the following directory:

Syntax

```
cd /var/lib/ceph/DAEMON_PATH/
```

Example

```
[ceph: root@host01 /]# cd /var/lib/ceph/hosts/
```

3. Create the **hosts.yml** file:

Example

```
[ceph: root@host01 hosts]# touch hosts.yml
```

4. Edit the **hosts.yml** file to include the following details:

Example

```
service_type: host
addr: host01
hostname: host01
labels:
- mon
- osd
- mgr
---
service_type: host
addr: host02
hostname: host02
labels:
- mon
- osd
- mgr
---
service_type: host
addr: host03
hostname: host03
labels:
- mon
- osd
```

5. Deploy the hosts using service specification:

Syntax

```
ceph orch apply -i FILE_NAME.yml
```

Example

```
[ceph: root@host01 hosts]# ceph orch apply -i hosts.yml
```

Verification

- List the hosts:

Example

```
[ceph: root@host01 /]# ceph orch host ls
```

Additional Resources

- See the [Listing hosts using the Ceph Orchestrator](#) section in the *Red Hat Ceph Storage Operations Guide*.

4.4. LISTING HOSTS USING THE CEPH ORCHESTRATOR

You can list hosts of a Ceph cluster with Ceph Orchestrators.



NOTE

The *STATUS* of the hosts is blank, in the output of the **ceph orch host ls** command.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the storage cluster.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. List the hosts of the cluster:

Example

```
[ceph: root@host01 /]# ceph orch host ls
```

You will see that the *STATUS* of the hosts is blank which is expected.

4.5. ADDING LABELS TO HOSTS USING THE CEPH ORCHESTRATOR

You can use the Ceph Orchestrator with Cephadm in the backend to add labels to hosts in an existing Red Hat Ceph Storage cluster. Few examples of labels are **mgr**, **mon**, and **osd** based on the service deployed on the hosts.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the storage cluster

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Add labels to the hosts:

Syntax

```
ceph orch host label add HOST_NAME LABEL_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch host label add host02 mon
```

3. Optional: You can add multiple labels to hosts:

Syntax

```
ceph orch host label add HOSTNAME_1 LABEL_1,LABEL_2
```

Example

```
[ceph: root@host01 /]# ceph orch host label add host01 mon,mgr
```

Verification

- List the hosts:

Example

```
[ceph: root@host01 /]# ceph orch host ls
```

4.6. REMOVING HOSTS USING THE CEPH ORCHESTRATOR

You can remove hosts of a Ceph cluster with the Ceph Orchestrators. You need to also remove the **node-exporter** and **crash** services when removing hosts to avoid retaining containers in the host.

**NOTE**

Remove all the services including managers, monitors, OSDs, and others, manually before removing the hosts from the storage cluster.

**IMPORTANT**

If you are removing the bootstrap host, be sure to copy the admin keyring and the configuration file to another host in the storage cluster before you remove the host.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.
- Hosts are added to the storage cluster.
- All the services are deployed.
- Cephadm is deployed on the nodes where the services have to be removed.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. For all Ceph services, except **node-exporter** and **crash**, remove the host from the placement specification files:

Example

```
service_type: rgw
placement:
  hosts:
  - host01
  - host02
```

In this example, **host03** is removed the placement specification of Ceph object gateway service. You need to follow the above step for all the services deployed on the host.

- a. For removing OSDs, deploy the OSDs with **--unmanaged=True**

Example

```
[ceph: root@host01 /]# ceph orch apply osd --all-available-devices --unmanaged=true
```

**NOTE**

This prevents auto-deployment of OSDs on the devices.

- b. Remove the OSDs:

Syntax

```
for osd_id in $(ceph orch ps HOST_NAME --daemon_type osd | grep osd | awk '{print $1}' | cut -c 5-); do ceph orch osd rm $osd_id; done
```

Example

```
[ceph: root@host01 /]# for osd_id in $(ceph orch ps host03 --daemon_type osd | grep osd | awk '{print $1}' | cut -c 5-); do ceph orch osd rm $osd_id; done
```

- c. Remove the hosts:

Syntax

```
ceph orch host rm HOST_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch host rm host03
```

3. On the node from where you have to remove the **node-exporter** and **crash** services, run the following commands:

- a. As a root user, outside the Cephadm shell, fetch details of the **fsid** of the cluster and the name of the service:

Example

```
[root@host03 ~]# cephadm ls
```

- b. Copy the **fsid** and the **name** of the **node-exporter** service.

- c. Remove the service:

Syntax

```
cephadm rm-daemon --fsid CLUSTER_ID --name SERVICE_NAME
```

Example

```
[root@host03 ~]# cephadm rm-daemon --fsid a34c81a0-889b-11eb-af98-001a4a00063d --name node-exporter.host03
```

Verification

- Run **cephadm ls** command to verify service removal:

Example

```
[root@host03 ~]# cephadm ls
```

-
- List the hosts, daemons, and processes:

Example

```
[ceph: root@host01 /]# ceph orch ps
```

Additional Resources

- See the [Adding hosts using the Ceph Orchestrator](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.
- See the [Listing hosts using the Ceph Orchestrator](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

4.7. PLACING HOSTS IN THE MAINTENANCE MODE USING THE CEPH ORCHESTRATOR

You can use the Ceph Orchestrator to place the hosts in and out of the maintenance mode. This stops all the Ceph daemons on the host.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.
- Hosts added to the cluster.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. You can either place the host in maintenance mode or place it out of the maintenance mode:
 - Place the host in maintenance mode:

Syntax

```
ceph orch host maintenance enter HOST_NAME [--force]
```

Example

```
[ceph: root@host01 /]# ceph orch host maintenance enter host02 --force
```

The **--force** flag allows the user to bypass warnings, but not alerts.

- Place the host out of the maintenance mode:

Syntax

```
ceph orch host maintenance exit HOST_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch host maintenance exit host02
```

Verification

- List the hosts:

Example

```
[ceph: root@host01 /]# ceph orch host ls
```

CHAPTER 5. MANAGEMENT OF MONITORS USING THE CEPH ORCHESTRATOR

As a storage administrator, you can deploy additional monitors using placement specification, add monitors using service specification, add monitors to a subnet configuration, and add monitors to specific hosts. Apart from this, you can remove the monitors using the Ceph Orchestrator.

By default, a typical Red Hat Ceph Storage cluster has three or five monitor daemons deployed on different hosts.

Red Hat recommends deploying five monitors if there are five or more nodes in a cluster.

Ceph deploys monitor daemons automatically as the cluster grows, and scales back monitor daemons automatically as the cluster shrinks. The smooth execution of this automatic growing and shrinking depends upon proper subnet configuration.

If your monitor nodes or your entire cluster are located on a single subnet, then Cephadm automatically adds up to five monitor daemons as you add new hosts to the cluster. Cephadm automatically configures the monitor daemons on the new hosts. The new hosts reside on the same subnet as the bootstrapped host in the storage cluster.

Cephadm can also deploy and scale monitors to correspond to changes in the size of the storage cluster.

5.1. CEPH MONITORS

Ceph Monitors are lightweight processes that maintain a master copy of the storage cluster map. All Ceph clients contact a Ceph monitor and retrieve the current copy of the storage cluster map, enabling clients to bind to a pool and read and write data.

Ceph Monitors use a variation of the Paxos protocol to establish consensus about maps and other critical information across the storage cluster. Due to the nature of Paxos, Ceph requires a majority of monitors running to establish a quorum, thus establishing consensus.



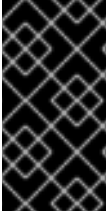
IMPORTANT

Red Hat requires at least three monitors on separate hosts to receive support for a production cluster.

Red Hat recommends deploying an odd number of monitors. An odd number of Ceph Monitors has a higher resiliency to failures than an even number of monitors. For example, to maintain a quorum on a two-monitor deployment, Ceph cannot tolerate any failures; with three monitors, one failure; with four monitors, one failure; with five monitors, two failures. This is why an odd number is advisable. Summarizing, Ceph needs a majority of monitors to be running and to be able to communicate with each other, two out of three, three out of four, and so on.

For an initial deployment of a multi-node Ceph storage cluster, Red Hat requires three monitors, increasing the number two at a time if a valid need for more than three monitors exists.

Since Ceph Monitors are lightweight, it is possible to run them on the same host as OpenStack nodes. However, Red Hat recommends running monitors on separate hosts.



IMPORTANT

Red Hat does NOT support collocating Ceph Monitors and OSDs on the same node. Doing this can have a negative impact to storage cluster performance.

Red Hat ONLY supports collocating Ceph services in containerized environments.

When you remove monitors from a storage cluster, consider that Ceph Monitors use the Paxos protocol to establish a consensus about the master storage cluster map. You must have a sufficient number of Ceph Monitors to establish a quorum.

Additional Resources

- See the Red Hat Ceph Storage [Supported configurations Knowledgebase article](#) for all the supported Ceph configurations.

5.2. CONFIGURING MONITOR ELECTION STRATEGY

The monitor election strategy identifies the net splits and handles failures. You can configure the election monitor strategy in three different modes:

1. **classic** - This is the default mode in which the lowest ranked monitor is voted based on the elector module between the two sites.
2. **disallow** - This mode lets you mark monitors as disallowed, in which case they will participate in the quorum and serve clients, but cannot be an elected leader. This lets you add monitors to a list of disallowed leaders. If a monitor is in the disallowed list, it will always defer to another monitor.
3. **connectivity** - This mode is mainly used to resolve network discrepancies. It evaluates connection scores provided by each monitor for its peers and elects the most connected and reliable monitor to be the leader. This mode is designed to handle net splits, which may happen if your cluster is stretched across multiple data centers or otherwise susceptible. This mode incorporates connection score ratings and elects the monitor with the best score.

Red Hat recommends you to stay in the **classic** mode unless you require features in the other modes.

Before constructing the cluster, change the **election_strategy** to **classic**, **disallow**, or **connectivity** in the following command:

Syntax

```
ceph mon set election_strategy {classic|disallow|connectivity}
```

5.3. DEPLOYING THE CEPH MONITOR DAEMONS USING THE COMMAND LINE INTERFACE

The Ceph Orchestrator deploys one monitor daemon by default. You can deploy additional monitor daemons by using the **placement** specification in the command line interface. To deploy a different number of monitor daemons, specify a different number. If you do not specify the hosts where the monitor daemons should be deployed, the Ceph Orchestrator randomly selects the hosts and deploys the monitor daemons to them.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. There are four different ways of deploying Ceph monitor daemons:

Method 1

- Use placement specification to deploy monitors on hosts:



NOTE

Red Hat recommends that you use the **--placement** option to deploy on specific hosts.

Syntax

```
ceph orch apply mon --placement="HOST_NAME_1 HOST_NAME_2 HOST_NAME_3"
```

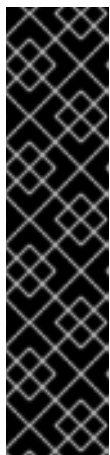
Example

```
[ceph: root@host01 /]# ceph orch apply mon --placement="host01 host02 host03"
```



NOTE

Be sure to include the bootstrap node as the first node in the command.



IMPORTANT

Do not add the monitors individually as **ceph orch apply mon** supersedes and will not add the monitors to all the hosts. For example, if you run the following commands, then the first command creates a monitor on **host01**. Then the second command supersedes the monitor on **host1** and creates a monitor on **host02**. Then the third command supersedes the monitor on **host02** and creates a monitor on **host03**. Eventually, there is a monitor only on the third host.

```
# ceph orch apply mon host01
# ceph orch apply mon host02
# ceph orch apply mon host03
```

Method 2

- Use placement specification to deploy specific number of monitors on specific hosts with labels:

- a. Add the labels to the hosts:

Syntax

```
ceph orch host label add HOSTNAME_1 LABEL_1,LABEL_2
```

Example

```
[ceph: root@host01 /]# ceph orch host label add host01 mon
```

- b. Deploy the daemons:

Syntax

```
ceph orch apply mon --placement="HOST_NAME_1:mon HOST_NAME_2:mon  
HOST_NAME_3:mon"
```

Example

```
[ceph: root@host01 /]# ceph orch apply mon --placement="host01:mon host02:mon  
host03:mon"
```

Method 3

- Use placement specification to deploy specific number of monitors on specific hosts:

Syntax

```
ceph orch apply mon --placement="NUMBER_OF_DAEMONS HOST_NAME_1  
HOST_NAME_2 HOST_NAME_3"
```

Example

```
[ceph: root@host01 /]# ceph orch apply mon --placement="3 host01 host02 host03"
```

Method 4

- Deploy monitor daemons randomly on the hosts in the storage cluster:

Syntax

```
ceph orch apply mon NUMBER_OF_DAEMONS
```

Example

```
[ceph: root@host01 /]# ceph orch apply mon 3
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=mon
```

5.4. DEPLOYING THE CEPH MONITOR DAEMONS USING THE SERVICE SPECIFICATION

The Ceph Orchestrator deploys one monitor daemon by default. You can deploy additional monitor daemons by using the service specification, like a YAML format file.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Navigate to the following directory:

Syntax

```
cd /var/lib/ceph/DAEMON_PATH/
```

Example

```
[ceph: root@host01 /]# cd /var/lib/ceph/mon/
```

3. Create the **mon.yml** file:

Example

```
[ceph: root@host01 mon]# touch mon.yml
```

4. Edit the **mon.yml** file to include the following details:

Syntax

```
service_type: mon
placement:
  hosts:
    - HOST_NAME_1
    - HOST_NAME_2
```

Example

```
service_type: mon
placement:
  hosts:
    - host01
    - host02
```

5. Deploy the monitor daemons:

Syntax

```
ceph orch apply -i FILE_NAME.yml
```

Example

```
[ceph: root@host01 mon]# ceph orch apply -i mon.yml
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=mon
```

5.5. DEPLOYING THE MONITOR DAEMONS ON SPECIFIC NETWORK USING THE CEPH ORCHESTRATOR

The Ceph Orchestrator deploys one monitor daemon by default. You can explicitly specify the IP address or CIDR network for each monitor and control where each monitor is placed.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Disable automated monitor deployment:

Example

```
[ceph: root@host01 /]# ceph orch apply mon --unmanaged
```

3. Deploy monitors on hosts on specific network:

Syntax

```
ceph orch daemon add mon HOST_NAME_1:_IP_OR_NETWORK_
```

Example

```
[ceph: root@host01 /]# ceph orch daemon add mon host03:10.1.2.123
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=mon
```

5.6. REMOVING THE MONITOR DAEMONS USING THE CEPH ORCHESTRATOR

To remove the monitor daemons from the host, you can just redeploy the monitor daemons on other hosts.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- At least one monitor daemon deployed on the hosts.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Run the **ceph orch apply** command to deploy the required monitor daemons:

Syntax

```
ceph orch apply mon "NUMBER_OF_DAEMONS _HOST_NAME_1 HOST_NAME_3"
```

If you want to remove monitor daemons from **host02**, then you can redeploy the monitors on other hosts.

Example

```
[ceph: root@host01 /]# ceph orch apply mon "2 host01 host03"
```

Verification

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=mon
```

Additional Resources

- See [Deploying the Ceph monitor daemons using the command line interface](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.
- See [Deploying the Ceph monitor daemons using the service specification](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

CHAPTER 6. MANAGEMENT OF MANAGERS USING THE CEPH ORCHESTRATOR

As a storage administrator, you can use the Ceph Orchestrator to deploy additional manager daemons. Cephadm automatically installs a manager daemon on the bootstrap node during the bootstrapping process.

This section covers the following administrative tasks:

- [Deploying the manager daemons using the Ceph Orchestrator](#).
- [Removing the manager daemons using the Ceph Orchestrator](#).

6.1. PREREQUISITES

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.
- Hosts are added to the cluster.

6.2. DEPLOYING THE MANAGER DAEMONS USING THE CEPH ORCHESTRATOR

The Ceph Orchestrator deploys two Manager daemons by default. You can deploy additional manager daemons using the **placement** specification in the command line interface. To deploy a different number of Manager daemons, specify a different number. If you do not specify the hosts where the Manager daemons should be deployed, the Ceph Orchestrator randomly selects the hosts and deploys the Manager daemons to them.



NOTE

Ensure your deployment has at least three Ceph Managers in each deployment.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. You can deploy manager daemons in two different ways:

Method 1

- Deploy manager daemons using placement specification on specific set of hosts:

**NOTE**

Red Hat recommends that you use the **--placement** option to deploy on specific hosts.

Syntax

```
ceph orch apply mgr --placement=" HOST_NAME_1 HOST_NAME_2 HOST_NAME_3"
```

Example

```
[ceph: root@host01 /]# ceph orch apply mgr --placement="host01 host02 host03"
```

Method 2

- Deploy manager daemons randomly on the hosts in the storage cluster:

Syntax

```
ceph orch apply mgr NUMBER_OF_DAEMONS
```

Example

```
[ceph: root@host01 /]# ceph orch apply mgr 3
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=mgr
```

6.3. REMOVING THE MANAGER DAEMONS USING THE CEPH ORCHESTRATOR

To remove the manager daemons from the host, you can just redeploy the daemons on other hosts.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.
- Hosts are added to the cluster.
- At least one manager daemon deployed on the hosts.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Run the **ceph orch apply** command to redeploy the required manager daemons:

Syntax

```
ceph orch apply mgr "NUMBER_OF_DAEMONS _HOST_NAME_1 HOST_NAME_3"
```

If you want to remove manager daemons from **host02**, then you can redeploy the manager daemons on other hosts.

Example

```
[ceph: root@host01 /]# ceph orch apply mgr "2 host01 host03"
```

Verification

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=mgr
```

Additional Resources

- See [Deploying the manager daemons using the Ceph Orchestrator](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

CHAPTER 7. MANAGEMENT OF OSDS USING THE CEPH ORCHESTRATOR

As a storage administrator, you can use the Ceph Orchestrators to manage OSDs of a Red Hat Ceph Storage cluster.

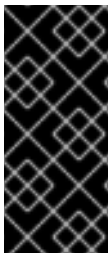
7.1. CEPH OSDS

When a Red Hat Ceph Storage cluster is up and running, you can add OSDs to the storage cluster at runtime.

A Ceph OSD generally consists of one **ceph-osd** daemon for one storage drive and its associated journal within a node. If a node has multiple storage drives, then map one **ceph-osd** daemon for each drive.

Red Hat recommends checking the capacity of a cluster regularly to see if it is reaching the upper end of its storage capacity. As a storage cluster reaches its **near full** ratio, add one or more OSDs to expand the storage cluster's capacity.

When you want to reduce the size of a Red Hat Ceph Storage cluster or replace the hardware, you can also remove an OSD at runtime. If the node has multiple storage drives, you might also need to remove one of the **ceph-osd** daemon for that drive. Generally, it's a good idea to check the capacity of the storage cluster to see if you are reaching the upper end of its capacity. Ensure that when you remove an OSD that the storage cluster is not at its **near full** ratio.



IMPORTANT

Do not let a storage cluster reach the **full** ratio before adding an OSD. OSD failures that occur after the storage cluster reaches the **near full** ratio can cause the storage cluster to exceed the **full** ratio. Ceph blocks write access to protect the data until you resolve the storage capacity issues. Do not remove OSDs without considering the impact on the **full** ratio first.

7.2. CEPH OSD NODE CONFIGURATION

Configure Ceph OSDs and their supporting hardware similarly as a storage strategy for the pool(s) that will use the OSDs. Ceph prefers uniform hardware across pools for a consistent performance profile. For best performance, consider a CRUSH hierarchy with drives of the same type or size.

If you add drives of dissimilar size, adjust their weights accordingly. When you add the OSD to the CRUSH map, consider the weight for the new OSD. Hard drive capacity grows approximately 40% per year, so newer OSD nodes might have larger hard drives than older nodes in the storage cluster, that is, they might have a greater weight.

Before doing a new installation, review the *Requirements for Installing Red Hat Ceph Storage* chapter in the [Installation Guide](#).

7.3. AUTOMATICALLY TUNING OSD MEMORY

The OSD daemons adjust the memory consumption based on the **osd_memory_target** configuration option. The option **osd_memory_target** sets OSD memory based upon the available RAM in the system.

If Red Hat Ceph Storage is deployed on dedicated nodes that do not share memory with other services, Cephadm automatically adjusts the per-OSD consumption based on the total amount of RAM and the number of deployed OSDs.

By default, the **osd_memory_target_autotune** parameter is set to **false** in Red Hat Ceph Storage 5. You can enable this option globally as follows within the Cephadm shell:

Example

```
[ceph: root@host01 /]# ceph config set osd osd_memory_target_autotune true
```

Once the storage cluster is upgraded to Red Hat Ceph Storage 5.0, for cluster maintenance such as addition of OSDs or replacement of OSDs, Red Hat recommends setting **osd_memory_target_autotune** parameter to **true** to autotune osd memory as per system memory.

Cephadm starts with a fraction **mgr/cephadm/autotune_memory_target_ratio**, which defaults to **0.7** of the total RAM in the system, subtract off any memory consumed by non-autotuned daemons such as non-OSDs and for OSDs for which **osd_memory_target_autotune** is false, and then divide by the remaining OSDs.

The **osd_memory_target** parameter is calculated as follows:

Syntax

```
osd_memory_target = TOTAL_RAM_OF_THE_OSD * (1048576) * (0.7) /  
NUMBER_OF OSDS_IN_THE_OSD_NODE
```

For example, if a node has 24 OSDs and has 251G RAM space, then **osd_memory_target** is **7860684936**.

The final targets are reflected in the configuration database with options. You can view the limits and the current memory consumed by each daemon from the **ceph orch ps** output under **MEM LIMIT** column.

You can manually set a specific memory target for an OSD in the storage cluster.

Example

```
[ceph: root@host01 /]# ceph config set osd.123 osd_memory_target 7860684936
```

You can exclude an OSD from memory autotuning.

Example

```
[ceph: root@host01 /]# ceph config set osd.123 osd_memory_target_autotune false
```

7.4. LISTING DEVICES FOR CEPH OSD DEPLOYMENT

You can check the list of available devices before deploying OSDs using the Ceph Orchestrator. The commands are used to print a list of devices discoverable by Cephadm. A storage device is considered available if all of the following conditions are met:

- The device must have no partitions.
- The device must not have any LVM state.

- The device must not be mounted.
- The device must not contain a file system.
- The device must not contain a Ceph BlueStore OSD.
- The device must be larger than 5 GB.

**NOTE**

Ceph will not provision an OSD on a device that is not available.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- All manager and monitor daemons are deployed.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]#cephadm shell
```

2. List the available devices to deploy OSDs:

Syntax

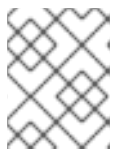
```
ceph orch device ls [--hostname=HOSTNAME_1 HOSTNAME_2] [--wide] [--refresh]
```

Example

```
[ceph: root@host01 /]# ceph orch device ls --wide --refresh
```

Using the **--wide** option provides all details relating to the device, including any reasons that the device might not be eligible for use as an OSD. This option does not support NVMe devices.

3. Optional: To enable *Health*, *Ident*, and *Fault* fields in the output of **ceph orch device ls**, run the following commands:

**NOTE**

These fields are supported by **libstoragegmt** library and currently supports SCSI, SAS, and SATA devices.

- a. As root user, check your hardware's compatibility with **libstoragegmt** library to avoid unplanned interruption to services:

Example

```
[root@host01 ~]# cephadm shell lsmdi ldl
```

In the output, you see the *Health Status* as *Good* with the respective *SCSI VPD 0x83 ID*.



NOTE

If you do not get this information, then enabling the fields might cause erratic behavior of devices.

- b. Enable **libstoragemgmt** support:

Example

```
[ceph: root@host01 /]# ceph config set mgr mgr/cephadm/device_enhanced_scan true
```

Once this is enabled, **ceph orch device ls** will give the output of *Health* field as *Good*.

Verification

- List the devices:

Example

```
[ceph: root@host01 /]# ceph orch device ls
```

7.5. ZAPPING DEVICES FOR CEPH OSD DEPLOYMENT

You need to check the list of available devices before deploying OSDs. If there is no space available on the devices, you can clear the data on the devices by zapping them.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- All manager and monitor daemons are deployed.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. List the available devices to deploy OSDs:

Syntax

```
ceph orch device ls [--hostname=HOSTNAME_1 HOSTNAME_2] [--wide] [--refresh]
```

Example

```
[ceph: root@host01 /]# ceph orch device ls --wide --refresh
```

3. Clear the data of a device:

Syntax

```
ceph orch device zap HOSTNAME FILE_PATH --force
```

Example

```
[ceph: root@host01 /]# ceph orch device zap host02 /dev/sdb --force
```

Verification

- Verify the space is available on the device:

Example

```
[ceph: root@host01 /]# ceph orch device ls
```

You will see that the field under *Available* is *Yes*.

Additional Resources

- See the [Listing devices for Ceph OSD deployment](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

7.6. DEPLOYING CEPH OSDS ON ALL AVAILABLE DEVICES

You can deploy all OSDs on all the available devices. Cephadm allows the Ceph Orchestrator to discover and deploy the OSDs on any available and unused storage device.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- All manager and monitor daemons are deployed.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. List the available devices to deploy OSDs:

Syntax

```
ceph orch device ls [--hostname=HOSTNAME_1 HOSTNAME_2] [--wide] [--refresh]
```

Example

```
[ceph: root@host01 /]# ceph orch device ls --wide --refresh
```

3. Deploy OSDs on all available devices:

Example

```
[ceph: root@host01 /]# ceph orch apply osd --all-available-devices
```

The effect of **ceph orch apply** is persistent which means that the Orchestrator automatically finds the device, adds it to the cluster, and creates new OSDs. This occurs under the following conditions:

- New disks or drives are added to the system.
- Existing disks or drives are zapped.
- An OSD is removed and the devices are zapped.
You can disable automatic creation of OSDs on all the available devices by using the **--unmanaged** parameter.

Example

```
[ceph: root@host01 /]# ceph orch apply osd --all-available-devices --unmanaged=true
```

Setting the parameter **--unmanaged** to **true** disables the creation of OSDs and also there is no change if you apply a new OSD service.



NOTE

The command **ceph orch daemon add** creates new OSDs, but does not add an OSD service.

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- View the details of the node and devices:

Example

```
[ceph: root@host01 /]# ceph osd tree
```

Additional Resources

- See the [Listing devices for Ceph OSD deployment](#) section in the *Red Hat Ceph Storage Operations Guide*.

7.7. DEPLOYING CEPH OSDS ON SPECIFIC DEVICES AND HOSTS

You can deploy all the Ceph OSDs on specific devices and hosts using the Ceph Orchestrator.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- All manager and monitor daemons are deployed.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. List the available devices to deploy OSDs:

Syntax

```
ceph orch device ls [--hostname=HOSTNAME_1 HOSTNAME_2] [--wide] [--refresh]
```

Example

```
[ceph: root@host01 /]# ceph orch device ls --wide --refresh
```

3. Deploy OSDs on specific devices and hosts:

Syntax

```
ceph orch daemon add osd HOSTNAME:_DEVICE_PATH_
```

Example

```
[ceph: root@host01 /]# ceph orch daemon add osd host02:/dev/sdb
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls osd
```

- View the details of the node and devices:

Example

```
[ceph: root@host01 /]# ceph osd tree
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --service_name=SERVICE_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --service_name=osd
```

Additional Resources

- See the [Listing devices for Ceph OSD deployment](#) section in the *Red Hat Ceph Storage Operations Guide*.

7.8. ADVANCED SERVICE SPECIFICATIONS AND FILTERS FOR DEPLOYING OSDS

Service Specification of type OSD is a way to describe a cluster layout using the properties of disks. It gives the user an abstract way to tell Ceph which disks should turn into an OSD with the required configuration without knowing the specifics of device names and paths. For each device and each host, define a **yaml** file or a **json** file.

General settings for OSD specifications

- **service_type:** 'osd': This is mandatory to create OSDS
- **service_id:** Use the service name or identification you prefer. A set of OSDs is created using the specification file. This name is used to manage all the OSDs together and represent an Orchestrator service.
- **placement:** This is used to define the hosts on which the OSDs needs to be deployed. You can use on the following options:
 - **host_pattern:** '*' - A host name pattern used to select hosts.
 - **label:** 'osd_host' - A label used in the hosts where OSD needs to be deployed.
 - **hosts:** 'host01', 'host02' - An explicit list of host names where OSDs needs to be deployed.
- **selection of devices:** The devices where OSDs are created. This allows to separate an OSD from different devices. You can create only BlueStore OSDs which have three components:
 - OSD data: contains all the OSD data
 - WAL: BlueStore internal journal or write-ahead Log
 - DB: BlueStore internal metadata

- **data_devices**: Define the devices to deploy OSD. In this case, OSDs are created in a collocated schema. You can use filters to select devices and folders.
- **wal_devices**: Define the devices used for WAL OSDs. You can use filters to select devices and folders.
- **db_devices**: Define the devices for DB OSDs. You can use the filters to select devices and folders.
- **encrypted**: An optional parameter to encrypt information on the OSD which can set to either **True** or **False**
- **unmanaged**: An optional parameter, set to False by default. You can set it to True if you do not want the Orchestrator to manage the OSD service.
- **block_wal_size**: User-defined value, in bytes.
- **block_db_size**: User-defined value, in bytes.

Filters for specifying devices

Filters are used in conjunction with the **data_devices**, **wal_devices** and **db_devices** parameters.

Name of the filter	Description	Syntax	Example
Model	Target specific disks. You can get details of the model by running lsblk -o NAME,FSTYPE,LABEL,MOUNTPOINT,SIZE,MODEL command or smartctl -i /DEVIVE_PATH	Model: <i>DISK_MODEL_NAME</i>	Model: MC-55-44-XZ
Vendor	Target specific disks	Vendor: <i>DISK_VENDOR_NAME</i>	Vendor: Vendor Cs
Size Specification	Includes disks of an exact size	size: <i>EXACT</i>	size: '10G'
Size Specification	Includes disks size of which is within the range	size: <i>LOW:HIGH</i>	size: '10G:40G'
Size Specification	Includes disks less than or equal to in size	size: <i>:HIGH</i>	size: ':10G'
Size Specification	Includes disks equal to or greater than in size	size: <i>LOW:</i>	size: '40G:'

Rotational	Rotational attribute of the disk. 1 matches all disks that are rotational and 0 matches all the disks that are non-rotational. If rotational=1, then OSD is configured with SSD or NVME. If rotational=0 then the OSD is configured with HDD.	rotational: 0 or 1	rotational: 0
All	Considers all the available disks	all: true	all: true
Limiter	When you specified valid filters but want to limit the amount of matching disks you can use the 'limit' directive. It should be used only as a last resort.	limit: <i>NUMBER</i>	limit: 2

**NOTE**

To create an OSD with non-collocated components in the same host, you have to specify the different type of devices used and the devices should be on the same host.

**NOTE**

The devices used for deploying OSDs must be supported by **libstoragemgmt**.

Additional Resources

- See the [Deploying Ceph OSDs using the advanced specifications](#) section in the *Red Hat Ceph Storage Operations Guide*.
- For more information on **libstoragemgmt**, see the [Listing devices for Ceph OSD deployment](#) section in the *Red Hat Ceph Storage Operations Guide*.

7.9. DEPLOYING CEPH OSDS USING ADVANCED SERVICE SPECIFICATIONS

The service specification of type OSD is a way to describe a cluster layout using the properties of disks. It gives the user an abstract way to tell Ceph which disks should turn into an OSD with the required configuration without knowing the specifics of device names and paths.

You can deploy the OSD for each device and each host by defining a **yml** file or a **json** file.

Prerequisites

- A running Red Hat Ceph Storage cluster.

- Hosts are added to the cluster.
- All manager and monitor daemons are deployed.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. On the monitor node, navigate the following directory:

Syntax

```
cd /var/lib/ceph/DAEMON_PATH/
```

Example

```
[ceph: root@host01 /]# cd /var/lib/ceph/osd/
```

3. Create the **osd_spec.yml** file:

Example

```
[ceph: root@host01 osd]# touch osd_spec.yml
```

4. Edit the **osd_spec.yml** file to include the following details:

- a. Simple scenario: In this case, all the nodes have the same set-up.

Syntax

```
service_type: osd
service_id: SERVICE_ID
placement:
  host_pattern: '*' # optional
data_devices: # optional
  model: DISK_MODEL_NAME # optional
db_devices: # optional
  size: # optional
  all: true # optional
encrypted: true
```

Example

```
service_type: osd
service_id: osd_spec_default
placement:
  host_pattern: '*'
```

```
data_devices:  
  all: true  
  encrypted: true
```

Example

```
service_type: osd  
service_id: osd_spec_default  
placement:  
  host_pattern: '**'  
data_devices:  
  size: '80G'  
db_devices:  
  size: '40G:'
```

- b. Advanced scenario: This would create the desired layout by using all HDDs as **data_devices** with two SSD assigned as dedicated DB or WAL devices. The remaining SSDs are **data_devices** that have the NVMEs vendors assigned as dedicated DB or WAL devices.

Example

```
service_type: osd  
service_id: osd_spec_hdd  
placement:  
  host_pattern: '**'  
data_devices:  
  rotational: 0  
db_devices:  
  model: Model-name  
  limit: 2  
---  
service_type: osd  
service_id: osd_spec_ssd  
placement:  
  host_pattern: '**'  
data_devices:  
  model: Model-name  
db_devices:  
  vendor: Vendor-name
```

- c. Advanced scenario with non-uniform nodes: This applies different OSD specs to different hosts depending on the `host_pattern` key.

Example

```
service_type: osd  
service_id: osd_spec_node_one_to_five  
placement:  
  host_pattern: 'node[1-5]'  
data_devices:  
  rotational: 1  
db_devices:  
  rotational: 0  
---
```

```

service_type: osd
service_id: osd_spec_six_to_ten
placement:
  host_pattern: 'node[6-10]'
data_devices:
  model: Model-name
db_devices:
  model: Model-name

```

- d. Advanced scenario with dedicated WAL and DB devices:

Example

```

service_type: osd
service_id: osd_using_paths
placement:
  hosts:
    - host01
    - host02
data_devices:
  paths:
    - /dev/sdb
db_devices:
  paths:
    - /dev/sdc
wal_devices:
  paths:
    - /dev/sdd

```

5. Before deploying OSDs, do a dry run:



NOTE

This step gives a preview of the deployment, without deploying the daemons.

Example

```
[ceph: root@host01 osd]# ceph orch apply -i osd_spec.yml --dry-run
```

6. Deploy OSDs using service specification:

Syntax

```
ceph orch apply -i FILE_NAME.yml
```

Example

```
[ceph: root@host01 osd]# ceph orch apply -i osd_spec.yml
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls osd
```

- View the details of the node and devices:

Example

```
[ceph: root@host01 /]# ceph osd tree
```

Additional Resources

- See the [Advanced service specifications and filters for deploying OSDs](#) section in the *Red Hat Ceph Storage Operations Guide*.

7.10. REMOVING THE OSD DAEMONS USING THE CEPH ORCHESTRATOR

You can remove the OSD from a cluster by using Cephadm.

Removing an OSD from a cluster involves two steps:

1. Evacuates all placement groups (PGs) from the cluster.
2. Removes the PG-free OSDs from the cluster.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- Ceph Monitor, Ceph Manager and Ceph OSD daemons are deployed on the storage cluster.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Check the device and the node from which the OSD has to be removed:

Example

```
[ceph: root@host01 /]# ceph osd tree
```

3. Remove the OSD:

Syntax

```
ceph orch osd rm OSD_ID [--replace] [--force]
```

Example

```
[ceph: root@host01 /]# ceph orch osd rm 0
```



NOTE

If you remove the OSD from the storage cluster without an option, such as **--replace**, the device is removed from the storage cluster completely. If you want to use the same device for deploying OSDs, you have to first zap the device before adding it to the storage cluster.

- Optional: To remove multiple OSDs from a specific node, run the following command:

Syntax

```
ceph orch osd rm OSD_ID OSD_ID
```

Example

```
[ceph: root@host01 /]# ceph orch osd rm 2 5
```

- Check the status of the OSD removal:

Example

```
[ceph: root@host01 /]# ceph orch osd rm status
OSD_ID HOST STATE PG_COUNT REPLACE FORCE DRAIN_STARTED_AT
2 host01 draining 124 False False 2021-09-07 16:26:07.142980
5 host01 draining 107 False False 2021-09-07 16:26:08.330371
```

When no PGs are left on the OSD, it is decommissioned and removed from the cluster.

Verification

- Verify the details of the devices and the nodes from which the Ceph OSDs are removed:

Example

```
[ceph: root@host01 /]# ceph osd tree
```

Additional Resources

- See the [Deploying Ceph OSDs on all available devices](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.
- See the [Deploying Ceph OSDs on specific devices and hosts](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

- See the [Zapping devices for Ceph OSD deployment](#) section in the *Red Hat Ceph Storage Operations Guide* for more information on clearing space on devices.

7.11. REPLACING THE OSDS USING THE CEPH ORCHESTRATOR

You can replace the OSDs from the cluster by preserving the OSD ID using the **ceph orch rm** command. The OSD is not permanently removed from the CRUSH hierarchy, but is assigned the **destroyed** flag. This flag is used to determine the OSD IDs that can be reused in the next OSD deployment. The 'destroyed' flag is used to determine which OSD ids is reused in the next OSD deployment.

If you use OSD specification for deployment, your newly added disks is assigned the OSD IDs of their replaced counterparts.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- Monitor, Manager and OSD daemons are deployed on the storage cluster.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Check the device and the node from which the OSD has to be replaced:

Example

```
[ceph: root@host01 /]# ceph osd tree
```

3. Replace the OSD:

Syntax

```
ceph orch osd rm OSD_ID --replace
```

Example

```
[ceph: root@host01 /]# ceph orch osd rm 0 --replace
```

4. Check the status of the OSD replacement:

Example

```
[ceph: root@host01 /]# ceph orch osd rm status
```

Verification

verification

- Verify the details of the devices and the nodes from which the Ceph OSDs are replaced:

Example

```
[ceph: root@host01 /]# ceph osd tree
```

Additional Resources

- See the [Deploying Ceph OSDs on all available devices](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.
- See the [Deploying Ceph OSDs on specific devices and hosts](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

7.12. STOPPING THE REMOVAL OF THE OSDS USING THE CEPH ORCHESTRATOR

You can stop the removal of only the OSDs that are queued for removal. This resets the initial state of the OSD and takes it off the removal queue.

If the OSD is in the process of removal, then you cannot stop the process.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- Monitor, Manager and OSD daemons are deployed on the cluster.
- Remove OSD process initiated.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Check the device and the node from which the OSD was initiated to be removed:

Example

```
[ceph: root@host01 /]# ceph osd tree
```

3. Stop the removal of the queued OSD:

Syntax

```
ceph orch osd rm stop OSD_ID
```

Example

```
[ceph: root@host01 /]# ceph orch osd rm stop 0
```

4. Check the status of the OSD removal:

Example

```
[ceph: root@host01 /]# ceph orch osd rm status
```

Verification

- Verify the details of the devices and the nodes from which the Ceph OSDs were queued for removal:

Example

```
[ceph: root@host01 /]# ceph osd tree
```

Additional Resources

- See [Removing the OSD daemons using the Ceph Orchestrator](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

7.13. ACTIVATING THE OSDS USING THE CEPH ORCHESTRATOR

You can activate the OSDs in the cluster in cases where the operating system of the host was reinstalled.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- Monitor, Manager and OSD daemons are deployed on the storage cluster.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. After the operating system of the host is reinstalled, activate the OSDs:

Syntax

```
ceph cephadm osd activate HOSTNAME
```

Example

```
[ceph: root@host01 /]# ceph cephadm osd activate host03
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --service_name=SERVICE_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --service_name=osd
```

7.13.1. Observing the data migration

When you add or remove an OSD to the CRUSH map, Ceph begins rebalancing the data by migrating placement groups to the new or existing OSD(s). You can observe the data migration using **ceph-w** command.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Recently added or removed an OSD.

Procedure

1. To observe the data migration:

Example

```
[ceph: root@host01 /]# ceph -w
```

2. Watch as the placement group states change from **active+clean** to **active, some degraded objects**, and finally **active+clean** when migration completes.
3. To exit the utility, press **Ctrl + C**.

7.14. RECALCULATING THE PLACEMENT GROUPS

Placement groups (PGs) define the spread of any pool data across the available OSDs. A placement group is built upon the given redundancy algorithm to be used. For a 3-way replication, the redundancy is defined to use three different OSDs. For erasure-coded pools, the number of OSDs to use is defined by the number of chunks.

When defining a pool the number of placement groups defines the grade of granularity the data is spread with across all available OSDs. The higher the number the better the equalization of capacity load can be. However, since handling the placement groups is also important in case of reconstruction of data, the number is significant to be carefully chosen upfront. To support calculation a tool is available to produce agile environments.

During lifetime of a storage cluster a pool may grow above the initially anticipated limits. With the growing number of drives a recalculation is recommended. The number of placement groups per OSD should be around 100. When adding more OSDs to the storage cluster the number of PGs per OSD will lower over time. Starting with 120 drives initially in the storage cluster and setting the `pg_num` of the pool to 4000 will end up in 100 PGs per OSD, given with the replication factor of three. Over time, when growing to ten times the number of OSDs, the number of PGs per OSD will go down to ten only. Because small number of PGs per OSD will tend to an unevenly distributed capacity, consider adjusting the PGs per pool.

Adjusting the number of placement groups can be done online. Recalculating is not only a recalculation of the PG numbers, but will involve data relocation, which will be a lengthy process. However, the data availability will be maintained at any time.

Very high numbers of PGs per OSD should be avoided, because reconstruction of all PGs on a failed OSD will start at once. A high number of IOPS is required to perform reconstruction in a timely manner, which might not be available. This would lead to deep I/O queues and high latency rendering the storage cluster unusable or will result in long healing times.

Additional Resources

- See the [PG calculator](#) for calculating the values by a given use case.
- See the [Erasure Code Pools](#) chapter in the *Red Hat Ceph Storage Strategies Guide* for more information.

7.15. USING THE CEPH MANAGER BALANCER MODULE

The balancer is a module for Ceph Manager (`ceph-mgr`) that optimizes the placement of placement groups (PGs) across OSDs in order to achieve a balanced distribution, either automatically or in a supervised fashion.

Currently the balancer module cannot be disabled. It can only be turned off to customize the configuration.

Prerequisites

- A running Red Hat Ceph Storage cluster.

Procedure

1. Ensure the balancer module is enabled:

Example

```
[ceph: root@host01 /]# ceph mgr module enable balancer
```

2. Turn on the balancer module:

Example

```
[ceph: root@host01 /]# ceph balancer on
```

- The default mode is **crush-compat**. The mode can be changed with:

Example

```
[ceph: root@host01 /]# ceph balancer mode upmap
```

or

Example

```
[ceph: root@host01 /]# ceph balancer mode crush-compat
```

Status

The current status of the balancer can be checked at any time with:

Example

```
[ceph: root@host01 /]# ceph balancer status
```

Automatic balancing

By default, when turning on the balancer module, automatic balancing is used:

Example

```
[ceph: root@host01 /]# ceph balancer on
```

The balancer can be turned back off again with:

Example

```
[ceph: root@host01 /]# ceph balancer off
```

This will use the **crush-compat** mode, which is backward compatible with older clients and will make small changes to the data distribution over time to ensure that OSDs are equally utilized.

Throttling

No adjustments will be made to the PG distribution if the cluster is degraded, for example, if an OSD has failed and the system has not yet healed itself.

When the cluster is healthy, the balancer throttles its changes such that the percentage of PGs that are misplaced, or need to be moved, is below a threshold of 5% by default. This percentage can be adjusted using the **target_max_misplaced** setting. For example, to increase the threshold to 7%:

Example

```
[ceph: root@host01 /]# ceph config-key set mgr/balancer/target_max_misplaced .07
```

Supervised optimization

The balancer operation is broken into a few distinct phases:

1. Building a **plan**.
2. Evaluating the quality of the data distribution, either for the current PG distribution, or the PG distribution that would result after executing a **plan**.
3. Executing the **plan**.
 - To evaluate and score the current distribution:

Example

```
[ceph: root@host01 /]# ceph balancer eval
```

- To evaluate the distribution for a single pool:

Syntax

```
ceph balancer eval POOL_NAME
```

Example

```
[ceph: root@host01 /]# ceph balancer eval rbd
```

- To see greater detail for the evaluation:

Example

```
[ceph: root@host01 /]# ceph balancer eval-verbose ...
```

- To generate a plan using the currently configured mode:

Syntax

```
ceph balancer optimize PLAN_NAME
```

Replace *PLAN_NAME* with a custom plan name.

Example

```
[ceph: root@host01 /]# ceph balancer optimize rbd_123
```

- To see the contents of a plan:

Syntax

```
ceph balancer show PLAN_NAME
```

Example

```
[ceph: root@host01 /]# ceph balancer show rbd_123
```

- To discard old plans:

Syntax

```
ceph balancer rm PLAN_NAME
```

Example

```
[ceph: root@host01 /]# ceph balancer rm rbd_123
```

- To see currently recorded plans use the status command:

```
[ceph: root@host01 /]# ceph balancer status
```

- To calculate the quality of the distribution that would result after executing a plan:

Syntax

```
ceph balancer eval PLAN_NAME
```

Example

```
[ceph: root@host01 /]# ceph balancer eval rbd_123
```

- To execute the plan:

Syntax

```
ceph balancer execute PLAN_NAME
```

Example

```
[ceph: root@host01 /]# ceph balancer execute rbd_123
```



NOTE

Only execute the plan if it is expected to improve the distribution. After execution, the plan will be discarded.

7.16. USING THE CEPH MANAGER CRASH MODULE

Using the Ceph manager crash module, you can collect information about daemon crashdumps and store it in the Red Hat Ceph Storage cluster for further analysis.

By default, daemon crashdumps are dumped in `/var/lib/ceph/crash`. You can configure with the option **crash dir**. Crash directories are named by time, date, and a randomly-generated UUID, and contain a metadata file **meta** and a recent log file, with a **crash_id** that is the same.

You can use **ceph-crash.service** to submit these crash automatically and persist in the Ceph Monitors. The **ceph-crash.service** watches watches the crashdump directory and uploads them with **ceph crash post**.

The `RECENT_CRASH` health message is one of the most common health messages in a Ceph cluster. This health message means that one or more Ceph daemons has crashed recently, and the crash has not yet been archived or acknowledged by the administrator. This might indicate a software bug, a hardware problem like a failing disk, or some other problem. The option `mgr/crash/warn_recent_interval` controls the time period of what recent means, which is two weeks by default. You can disable the warnings by running the following command:

Example

```
[ceph: root@host01 /]# ceph config set mgr/crash/warn_recent_interval 0
```

The option `mgr/crash/retain_interval` controls the period for which you want to retain the crash reports before they are automatically purged. The default for this option is one year.

Prerequisites

- A running Red Hat Ceph Storage cluster.

Procedure

1. Ensure the crash module is enabled:

Example

```
[ceph: root@host01 /]# ceph mgr module ls | more
{
  "always_on_modules": [
    "balancer",
    "crash",
    "devicehealth",
    "orchestrator_cli",
    "progress",
    "rbd_support",
    "status",
    "volumes"
  ],
  "enabled_modules": [
    "dashboard",
    "pg_autoscaler",
    "prometheus"
  ]
}
```

2. Save a crash dump: The metadata file is a JSON blob stored in the crash dir as `meta`. You can invoke the ceph command `-i` option, which reads from stdin.

Example

```
[ceph: root@host01 /]# ceph crash post -i meta
```

3. List the timestamp or the UUID crash IDs for all the new and archived crash info:

Example

```
[ceph: root@host01 /]# ceph crash ls
```


-
- 4. List the timestamp or the UUID crash IDs for all the new crash information:

Example

```
[ceph: root@host01 /]# ceph crash ls-new
```

- 5. List the timestamp or the UUID crash IDs for all the new crash information:

Example

```
[ceph: root@host01 /]# ceph crash ls-new
```

- 6. List the summary of saved crash information grouped by age:

Example

```
[ceph: root@host01 /]# ceph crash stat
8 crashes recorded
8 older than 1 days old:
2021-05-20T08:30:14.533316Z_4ea88673-8db6-4959-a8c6-0eea22d305c2
2021-05-20T08:30:14.590789Z_30a8bb92-2147-4e0f-a58b-a12c2c73d4f5
2021-05-20T08:34:42.278648Z_6a91a778-bce6-4ef3-a3fb-84c4276c8297
2021-05-20T08:34:42.801268Z_e5f25c74-c381-46b1-bee3-63d891f9fc2d
2021-05-20T08:34:42.803141Z_96adfc59-be3a-4a38-9981-e71ad3d55e47
2021-05-20T08:34:42.830416Z_e45ed474-550c-44b3-b9bb-283e3f4cc1fe
2021-05-24T19:58:42.549073Z_b2382865-ea89-4be2-b46f-9a59af7b7a2d
2021-05-24T19:58:44.315282Z_1847afbc-f8a9-45da-94e8-5aef0738954e
```

- 7. View the details of the saved crash:

Syntax

```
ceph crash info CRASH_ID
```

Example

```
[ceph: root@host01 /]# ceph crash info 2021-05-24T19:58:42.549073Z_b2382865-ea89-4be2-b46f-9a59af7b7a2d
{
  "assert_condition": "session_map.sessions.empty()",
  "assert_file": "/builddir/build/BUILD/ceph-16.1.0-486-g324d7073/src/mon/Monitor.cc",
  "assert_func": "virtual Monitor::~Monitor()",
  "assert_line": 287,
  "assert_msg": "/builddir/build/BUILD/ceph-16.1.0-486-g324d7073/src/mon/Monitor.cc: In function 'virtual Monitor::~Monitor()' thread 7f67a1aeb700 time 2021-05-24T19:58:42.545485+0000\n/builddir/build/BUILD/ceph-16.1.0-486-g324d7073/src/mon/Monitor.cc: 287: FAILED\nceph_assert(session_map.sessions.empty())\n",
  "assert_thread_name": "ceph-mon",
  "backtrace": [
    "/lib64/libpthread.so.0(+0x12b30) [0x7f679678bb30]",
    "gsignal()",
  ]
}
```

```

    "abort()",
    "(ceph::__ceph_assert_fail(char const*, char const*, int, char const*)+0x1a9)
[0x7f6798c8d37b]",
    "/usr/lib64/ceph/libceph-common.so.2(+0x276544) [0x7f6798c8d544]",
    "(Monitor::~Monitor()+0xe30) [0x561152ed3c80]",
    "(Monitor::~Monitor()+0xd) [0x561152ed3cdd]",
    "main()",
    "__libc_start_main()",
    "_start()"
  ],
  "ceph_version": "16.1.0-486.el8cp",
  "crash_id": "2021-05-24T19:58:42.549073Z_b2382865-ea89-4be2-b46f-9a59af7b7a2d",
  "entity_name": "mon.ceph-adm4",
  "os_id": "rhel",
  "os_name": "Red Hat Enterprise Linux",
  "os_version": "8.3 (Ootpa)",
  "os_version_id": "8.3",
  "process_name": "ceph-mon",
  "stack_sig":
"957c21d558d0cba4cee9e8aaf9227b3b1b09738b8a4d2c9f4dc26d9233b0d511",
  "timestamp": "2021-05-24T19:58:42.549073Z",
  "utsname_hostname": "host02",
  "utsname_machine": "x86_64",
  "utsname_release": "4.18.0-240.15.1.el8_3.x86_64",
  "utsname_sysname": "Linux",
  "utsname_version": "#1 SMP Wed Feb 3 03:12:15 EST 2021"
}

```

- Remove saved crashes older than *KEEP* days: Here, *KEEP* must be an integer.

Syntax

```
ceph crash prune KEEP
```

Example

```
[ceph: root@host01 /]# ceph crash prune 60
```

- Archive a crash report so that it is no longer considered for the **RECENT_CRASH** health check and does not appear in the **crash ls-new** output. It appears in the **crash ls**.

Syntax

```
ceph crash archive CRASH_ID
```

Example

```
[ceph: root@host01 /]# ceph crash archive 2021-05-24T19:58:42.549073Z_b2382865-ea89-4be2-b46f-9a59af7b7a2d
```

- Archive all crash reports:

Example

```
-
```

```
[ceph: root@host01 /]# ceph crash archive-all
```

11. Remove the crash dump:

Syntax

```
ceph crash rm CRASH_ID
```

Example

```
[ceph: root@host01 /]# ceph crash rm 2021-05-24T19:58:42.549073Z_b2382865-ea89-4be2-b46f-9a59af7b7a2d
```

Additional Resources

- See the [Health messages of a Ceph cluster](#) section in the *Red Hat Ceph Storage Troubleshooting Guide* for more information on Ceph health messages.

CHAPTER 8. MANAGEMENT OF MONITORING STACK USING THE CEPH ORCHESTRATOR

As a storage administrator, you can use the Ceph Orchestrator with Cephadm in the backend to deploy monitoring and alerting stack. The monitoring stack consists of Prometheus, Prometheus exporters, Prometheus Alertmanager, and Grafana. Users need to either define these services with Cephadm in a YAML configuration file, or they can use the command line interface to deploy them. When multiple services of the same type are deployed, a highly-available setup is deployed. The node exporter is an exception to this rule.



NOTE

Red Hat Ceph Storage 5.0 does not support custom images for deploying monitoring services such as Prometheus, Grafana, Alertmanager, and node-exporter.

The following monitoring services can be deployed with Cephadm:

- Prometheus is the monitoring and alerting toolkit. It collects the data provided by Prometheus exporters and fires preconfigured alerts if predefined thresholds have been reached. The Prometheus manager module provides a Prometheus exporter to pass on Ceph performance counters from the collection point in **ceph-mgr**.
The Prometheus configuration, including scrape targets, such as metrics providing daemons, is set up automatically by Cephadm. Cephadm also deploys a list of default alerts, for example, health error, 10% OSDs down, or pgs inactive.
- Alertmanager handles alerts sent by the Prometheus server. It deduplicates, groups, and routes the alerts to the correct receiver. By default, the Ceph dashboard is automatically configured as the receiver. The Alertmanager handles alerts sent by the Prometheus server. Alerts can be silenced using the Alertmanager, but silences can also be managed using the Ceph Dashboard.
- Grafana is the visualization and alerting software. The alerting functionality of Grafana is not used by this monitoring stack. For alerting, the Alertmanager is used.
By default, traffic to Grafana is encrypted with TLS. You can either supply your own TLS certificate or use a self-signed one. If no custom certificate has been configured before Grafana has been deployed, then a self-signed certificate is automatically created and configured for Grafana. Custom certificates for Grafana can be configured using the following commands:

Syntax

```
ceph config-key set mgr/cephadm/grafana_key -i
PRESENT_WORKING_DIRECTORY/key.pem
ceph config-key set mgr/cephadm/grafana_cert -i
PRESENT_WORKING_DIRECTORY/certificate.pem
```

Node exporter is an exporter for Prometheus which provides data about the node on which it is installed. It is recommended to install the node exporter on all nodes. This can be done using the `monitoring.yml` file with the `node-exporter` service type.

8.1. DEPLOYING THE MONITORING STACK USING THE CEPH ORCHESTRATOR

The monitoring stack consists of Prometheus, Prometheus exporters, Prometheus Alertmanager, and Grafana. Ceph Dashboard makes use of these components to store and visualize detailed metrics on cluster usage and performance.

You can deploy the monitoring stack using the service specification in YAML file format. All the monitoring services can have the network and port they bind to configured in the **yml** file.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Root-level access to the nodes.

Procedure

1. Enable the prometheus module in the Ceph Manager daemon. This exposes the internal Ceph metrics so that Prometheus can read them:

Example

```
[ceph: root@host01 /]# ceph mgr module enable prometheus
```



IMPORTANT

Ensure this command is run before Prometheus is deployed. If the command was not run before the deployment, you must redeploy Prometheus to update the configuration:

```
ceph orch redeploy prometheus
```

2. Navigate to the following directory:

Syntax

```
cd /var/lib/ceph/DAEMON_PATH/
```

Example

```
[ceph: root@host01 mds/]# cd /var/lib/ceph/monitoring/
```



NOTE

If the directory **monitoring** does not exist, create it.

3. Create the **monitoring.yml** file:

Example

```
[ceph: root@host01 monitoring]# touch monitoring.yml
```

4. Edit the specification file with a content similar to the following example:

Example

```
service_type: prometheus
service_name: prometheus
placement:
  hosts:
  - host01
networks:
- 192.169.142.0/24
---
service_type: node-exporter
---
service_type: alertmanager
service_name: alertmanager
placement:
  hosts:
  - host03
networks:
- 192.169.142.03/24
---
service_type: grafana
service_name: grafana
placement:
  hosts:
  - host02
networks:
- 192.169.142.02/24
```

5. Apply monitoring services:

Example

```
[ceph: root@host01 monitoring]# ceph orch apply -i monitoring.yml
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --service_name=SERVICE_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --service_name=prometheus
```



IMPORTANT

Prometheus, Grafana, and the Ceph dashboard are all automatically configured to talk to each other, resulting in a fully functional Grafana integration in the Ceph dashboard.

8.2. REMOVING THE MONITORING STACK USING THE CEPH ORCHESTRATOR

You can remove the monitoring stack using the **ceph orch rm** command.

Prerequisites

- A running Red Hat Ceph Storage cluster.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Use the **ceph orch rm** command to remove the monitoring stack:

Syntax

```
ceph orch rm SERVICE_NAME --force
```

Example

```
[ceph: root@host01 /]# ceph orch rm grafana
[ceph: root@host01 /]# ceph orch rm prometheus
[ceph: root@host01 /]# ceph orch rm node-exporter
[ceph: root@host01 /]# ceph orch rm alertmanager
[ceph: root@host01 /]# ceph mgr module disable prometheus
```

3. Check the status of the process:

Example

```
[ceph: root@host01 /]# ceph orch status
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps
```

Example

```
[ceph: root@host01 /]# ceph orch ps
```

Additional Resources

- See [Deploying the monitoring stack using the Ceph Orchestrator](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

CHAPTER 9. BASIC RED HAT CEPH STORAGE CLIENT SETUP

As a storage administrator, you have to set up client machines with basic configuration to interact with the storage cluster. Most client machines only need the **ceph-common package** and its dependencies installed. It will supply the basic **ceph** and **rados** commands, as well as other commands like **mount.ceph** and **rbd**.

9.1. CONFIGURING FILE SETUP ON CLIENT MACHINES

Client machines generally need a smaller configuration file than a full-fledged storage cluster member. You can generate a minimal configuration file which can give details to clients to reach the Ceph monitors.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Root access to the nodes.

Procedure

1. On the node where you want to set up the files, create a directory **ceph** in the **/etc** folder:

Example

```
[root@host01 ~]# mkdir /etc/ceph/
```

2. Navigate to **/etc/ceph** directory:

Example

```
[root@host01 ~]# cd /etc/ceph/
```

3. Generate the configuration file in the **ceph** directory:

Example

```
[root@host01 ceph]# ceph config generate-minimal-conf  
  
# minimal ceph.conf for 417b1d7a-a0e6-11eb-b940-001a4a000740  
[global]  
fsid = 417b1d7a-a0e6-11eb-b940-001a4a000740  
mon_host = [v2:10.74.249.41:3300/0,v1:10.74.249.41:6789/0]
```

The contents of this file should be installed in **/etc/ceph/ceph.conf** path. You can use this configuration file to reach the Ceph monitors.

9.2. SETTING-UP KEYRING ON CLIENT MACHINES

Most Ceph clusters are run with the authentication enabled, and the client needs the keys in order to communicate with cluster machines. You can generate the keyring which can give details to clients to reach the Ceph monitors.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Root access to the nodes.

Procedure

1. On the node where you want to set up the keyring, create a directory **ceph** in the **/etc** folder:

Example

```
[root@host01 ~]# mkdir /etc/ceph/
```

2. Navigate to **/etc/ceph** directory in the **ceph** directory:

Example

```
[root@host01 ~]# cd /etc/ceph/
```

3. Generate the keyring for the client:

Syntax

```
ceph auth get-or-create client.CLIENT_NAME -o /etc/ceph/NAME_OF_THE_FILE
```

Example

```
[root@host01 ceph]# ceph auth get-or-create client.fs -o /etc/ceph/ceph.keyring
```

4. Verify the output in the **ceph.keyring** file:

Example

```
[root@host01 ceph]# cat ceph.keyring  
  
[client.fs]  
key = AQAvoH5gkUCsExAATz3xCBLd4n6B6jRv+Z7CVQ==
```

The resulting output should be put into a keyring file, for example **/etc/ceph/ceph.keyring**.

CHAPTER 10. MANAGEMENT OF MDS SERVICE USING THE CEPH ORCHESTRATOR

As a storage administrator, you can use Ceph Orchestrator with Cephadm in the backend to deploy the MDS service. By default, a Ceph File System (CephFS) uses only one active MDS daemon. However, systems with many clients benefit from multiple active MDS daemons.

This section covers the following administrative tasks:

- [Deploying the MDS service using the command line interface](#) .
- [Deploying the MDS service using the service specification](#) .
- [Removing the MDS service using the Ceph Orchestrator](#) .

10.1. PREREQUISITES

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.
- Hosts are added to the cluster.
- All manager, monitor and OSD daemons are deployed.

10.2. DEPLOYING THE MDS SERVICE USING THE COMMAND LINE INTERFACE

Using the Ceph Orchestrator, you can deploy the Metadata Server (MDS) service using the **placement** specification in the command line interface Ceph File System (CephFS) requires one or more MDS.



NOTE

Ensure you have at least two pools, one for Ceph file system (CephFS) data and one for CephFS metadata.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- All manager, monitor and OSD daemons are deployed.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. There are two ways of deploying MDS daemons using placement specification:

Method 1

- Use **ceph fs volume** to create the MDS daemons. This creates the CephFS volume, pools associated to the CephFS, and also starts the MDS service on the hosts.

Syntax

```
ceph fs volume create FILESYSTEM_NAME --placement="NUMBER_OF_DAEMONS
HOST_NAME_1 HOST_NAME_2 HOST_NAME_3"
```



NOTE

By default, replicated pools are created for this command.

Example

```
[ceph: root@host01 /]# ceph fs volume create test --placement="2 host01 host02"
```

Method 2

- Create the pools, CephFS and then deploy MDS service using placement specification:
 - a. Create the pools for CephFS:

Syntax

```
ceph osd pool create DATA_POOL
ceph osd pool create METADATA_POOL
```

Example

```
[ceph: root@host01 /]# ceph osd pool create cephfs_data
[ceph: root@host01 /]# ceph osd pool create cephfs_metadata
```

- b. Create the file system for the data pools and metadata pools:

Syntax

```
ceph fs new FILESYSTEM_NAME DATA_POOL METADATA_POOL
```

Example

```
[ceph: root@host01 /]# ceph fs new test cephfs_data cephfs_metadata
```

- c. Deploy MDS service using the **ceph orch apply** command:

Syntax

```
ceph orch apply mds FILESYSTEM_NAME --placement="NUMBER_OF_DAEMONS
HOST_NAME_1 HOST_NAME_2 HOST_NAME_3"
```

Example

```
[ceph: root@host01 /]# ceph orch apply mds test --placement="2 host01 host02"
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- Check the CephFS status:

Example

```
[ceph: root@host01 /]# ceph fs ls
[ceph: root@host01 /]# ceph fs status
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=mds
```

Additional Resources

- See the [Red Hat Ceph Storage File System guide](#) for more information on creating the Ceph file system (CephFS).

10.3. DEPLOYING THE MDS SERVICE USING THE SERVICE SPECIFICATION

Using the Ceph Orchestrator, you can deploy the MDS service using the service specification.



NOTE

Ensure you have at least two pools, one for Ceph file system (CephFS) data and one for CephFS metadata.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- All manager, monitor and OSD daemons are deployed.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]#cephadm shell
```

2. Navigate to the following directory:

Syntax

```
cd /var/lib/ceph/DAEMON_PATH/
```

Example

```
[ceph: root@host01 mds]# cd /var/lib/ceph/mds/
```



NOTE

If the directory **mds** does not exist, create the directory.

3. Create the **mds.yml** file:

Example

```
[ceph: root@host01 mds/]# touch mds.yml
```

4. Edit the **mds.yml** file to include the following details:

Syntax

```
service_type: mds
service_id: FILESYSTEM_NAME
placement:
  hosts:
    - HOST_NAME_1
    - HOST_NAME_2
    - HOST_NAME_3
```

Example

```
service_type: mds
service_id: fs_name
placement:
  hosts:
    - host01
    - host02
```

5. Deploy MDS service using service specification:

Syntax

```
ceph orch apply -i FILE_NAME.yml
```

Example

```
[ceph: root@host01 mds]# ceph orch apply -i mds.yml
```

- Once the MDS services is deployed and functional, create the CephFS:

Syntax

```
ceph fs new CEPHFS_NAME METADATA_POOL DATA_POOL
```

Example

```
[ceph: root@host01 /]# ceph fs new test metadata_pool data_pool
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=mds
```

Additional Resources

- See the [Red Hat Ceph Storage File System guide](#) for more information on creating the Ceph file system (CephFS).

10.4. REMOVING THE MDS SERVICE USING THE CEPH ORCHESTRATOR

You can remove the service using the **ceph orch rm** command. Alternatively, you can remove the file system and the associated pools.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.

- Hosts are added to the cluster.
- At least one MDS daemon deployed on the hosts.

Procedure

1. There are two ways of removing MDS daemons from the cluster:

Method 1

- Remove the CephFS volume, associated pools and the services:
 - a. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

- b. Set the configuration parameter **mon_allow_pool_delete** to **true**:

Example

```
[ceph: root@host01 /]# ceph config set mon mon_allow_pool_delete true
```

- c. Remove the file system:

Syntax

```
ceph fs volume rm FILESYSTEM_NAME --yes-i-really-mean-it
```

Example

```
[ceph: root@host01 /]# ceph fs volume rm cephfs-new --yes-i-really-mean-it
```

This command will remove the file system, its data, and metadata pools. It also tries to remove MDS using the enabled **ceph-mgr** Orchestrator module.

Method 2

- Use the **ceph orch rm** command to remove the MDS service from the entire cluster:
 - a. List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- b. Remove the service

Syntax

```
ceph orch rm SERVICE_NAME
```


Example

```
[ceph: root@host01 /]# ceph orch rm mds.test
```

Verification

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps
```

Example

```
[ceph: root@host01 /]# ceph orch ps
```

Additional Resources

- See [Deploying the MDS service using the command line interface](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.
- See [Deploying the MDS service using the service specification](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

CHAPTER 11. MANAGEMENT OF CEPH OBJECT GATEWAY USING THE CEPH ORCHESTRATOR

As a storage administrator, you can deploy Ceph object gateway using the command line interface or by using the service specification.

You can also configure multisite object gateways, and remove the Ceph object gateway using the Ceph Orchestrator.

Cephadm deploys Ceph object gateway as a collection of daemons that manages a single-cluster deployment or a particular realm and zone in a multisite deployment.



NOTE

with Cephadm, the object gateway daemons are configured using the monitor configuration database instead of a **ceph.conf** or the command line. If that configuration is not already in the **client.rgw** section, then the object gateway daemons will start up with default settings and binds to the port **80**.

This section covers the following administrative tasks:

- [Deploying the Ceph object gateway using the command line interface](#) .
- [Deploying the Ceph object gateway using the service specification](#) .
- [Deploying a multisite Ceph object gateway using the Ceph Orchestrator](#) .
- [Removing the Ceph object gateway using the Ceph Orchestrator](#) .

11.1. PREREQUISITES

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.
- Hosts are added to the cluster.
- All the managers, monitors, and OSDs are deployed in the storage cluster.

11.2. DEPLOYING THE CEPH OBJECT GATEWAY USING THE COMMAND LINE INTERFACE

Using the Ceph Orchestrator, you can deploy the Ceph object gateway using the **ceph orch** command in the command line interface

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.
- Hosts are added to the cluster.
- All manager, monitor and OSD daemons are deployed.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. You can deploy the Ceph object gateway daemons in three different ways:

Method 1

- Create realm, zone group, zone and then use the placement specification with the host name:
 - a. Create a realm:

Syntax

```
radosgw-admin realm create --rgw-realm=REALM_NAME --default
```

Example

```
[ceph: root@host01 /]# radosgw-admin realm create --rgw-realm=test_realm --default
```

- b. Create a zone group:

Syntax

```
radosgw-admin zonegroup create --rgw-zonegroup=ZONE_GROUP_NAME --master --default
```

Example

```
[ceph: root@host01 /]# radosgw-admin zonegroup create --rgw-zonegroup=default --master --default
```

- c. Create a zone:

Syntax

```
radosgw-admin zone create --rgw-zonegroup=ZONE_GROUP_NAME --rgw-zone=ZONE_NAME --master --default
```

Example

```
[ceph: root@host01 /]# radosgw-admin zone create --rgw-zonegroup=default --rgw-zone=test_zone --master --default
```

- d. Commit the changes:

Syntax

```
radosgw-admin period update --rgw-realm=REALM_NAME --commit
```

Example

```
[ceph: root@host01 /]# radosgw-admin period update --rgw-realm=test_realm --commit
```

- e. Run the **ceph orch apply** command:

Syntax

```
ceph orch apply rgw NAME [--realm=REALM_NAME] [--zone=ZONE_NAME] --
placement="NUMBER_OF_DAEMONS [HOST_NAME_1 HOST_NAME_2]"
```

Example

```
[ceph: root@host01 /]# ceph orch apply rgw test --realm=test_realm --zone=test_zone --
placement="2 host01 host02"
```

Method 2

- Use an arbitrary service name to deploy two Ceph object gateway daemons for a single cluster deployment:

Syntax

```
ceph orch apply rgw SERVICE_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch apply rgw foo
```

Method 3

- Use an arbitrary service name on a labeled set of hosts:

Syntax

```
ceph orch host label add HOST_NAME_1 LABEL_NAME
ceph orch host label add HOSTNAME_2 LABEL_NAME
ceph orch apply rgw SERVICE_NAME '--placement=label:_LABEL_NAME_count-per-
host:_NUMBER_OF_DAEMONS_ ' --port=8000
```

Example

```
[ceph: root@host01 /]# ceph orch host label add host01 rgw # the 'rgw' label can be anything
[ceph: root@host01 /]# ceph orch host label add host02 rgw
[ceph: root@host01 /]# ceph orch apply rgw foo '--placement=label:rgw count-per-host:2' --
port=8000
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=rgw
```

11.3. DEPLOYING THE CEPH OBJECT GATEWAY USING THE SERVICE SPECIFICATION

Using the Ceph Orchestrator, you can deploy the Ceph object gateway using the service specification.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.
- Hosts are added to the cluster.
- All manager, monitor and OSD daemons are deployed.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Navigate to the following directory:

Syntax

```
cd /var/lib/ceph/DAEMON_PATH/
```

Example

```
[ceph: root@host01 nfs/]# cd /var/lib/ceph/rgw/
```

If the **rgw** directory does not exist, create the directory in the path.

3. Create the **rgw.yml** file:

Example

```
[ceph: root@host01 rgw]# touch rgw.yml
```

4. Edit the **rgw.yml** file to include the following details:

Syntax

```
service_type: rgw
service_id: REALM_NAME.ZONE_NAME
placement:
  hosts:
    - HOST_NAME_1
    - HOST_NAME_2
spec:
  rgw_realm: REALM_NAME
  rgw_zone: ZONE_NAME
networks:
  - NETWORK_IP_ADDRESS # RGW service binds to a specific network
```

Example

```
service_type: rgw
service_id: test_realm.test_zone
placement:
  hosts:
    - host01
    - host02
    - host03
spec:
  rgw_realm: test_realm
  rgw_zone: test_zone
networks:
  - 192.169.142.0/24
```

5. Deploy Ceph object gateway using service specification:

Syntax

```
ceph orch apply -i FILE_NAME.yml
```

Example

```
[ceph: root@host01 rgw]# ceph orch apply -i rgw.yml
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=rgw
```

11.4. DEPLOYING A MULTI-SITE CEPH OBJECT GATEWAY USING THE CEPH ORCHESTRATOR

Ceph Orchestrator supports multi-site configurations options for the Ceph Object Gateway.

You can configure each object gateway to work in an active-active zone configuration allowing for writes to a non-primary zone. The multi-site configuration is stored within a container called a realm.

The realm stores zone groups, zones and a time period. The **rgw** daemons handle the synchronization eliminating the need for a separate synchronization agent, thereby operating with an active-active configuration.

You can also deploy multi-site zones using the command line interface (CLI).



NOTE

The following configuration assumes at least two Red Hat Ceph Storage clusters are in geographically separate locations. However, the configuration also works on the same site.

Prerequisites

- At least two running Red Hat Ceph Storage clusters
- At least two Ceph Object Gateway instances, one for each Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.
- Nodes or containers are added to the storage cluster.
- All Ceph Manager, Monitor and OSD daemons are deployed.

Procedure

1. In the **cephadm** shell, configure the primary zone:

- a. Create a realm:

Syntax

```
radosgw-admin realm create --rgw-realm=REALM_NAME --default
```

Example

```
[ceph: root@host01 /]# radosgw-admin realm create --rgw-realm=test_realm --default
```

If the storage cluster has a single realm, then specify the **--default** flag.

- b. Create a primary zone group:

Syntax

```
radosgw-admin zonegroup create --rgw-zonegroup=ZONE_GROUP_NAME --  
endpoints=http://RGW_PRIMARY_HOSTNAME:RGW_PRIMARY_PORT_NUMBER_1_  
--master --default
```

Example

```
[ceph: root@host01 /]# radosgw-admin zonegroup create --rgw-zonegroup=us --  
endpoints=http://rgw1:80 --master --default
```

- c. Create a primary zone:

Syntax

```
radosgw-admin zone create --rgw-zonegroup=PRIMARY_ZONE_GROUP_NAME --rgw-  
zone=PRIMARY_ZONE_NAME --  
endpoints=http://RGW_PRIMARY_HOSTNAME:RGW_PRIMARY_PORT_NUMBER_1_  
--access-key=SYSTEM_ACCESS_KEY --secret=SYSTEM_SECRET_KEY
```

Example

```
[ceph: root@host01 /]# radosgw-admin zone create --rgw-zonegroup=us --rgw-zone=us-  
east-1 --endpoints=http://rgw1:80
```

- d. Optional: Delete the default zone, zone group, and the associated pools:



IMPORTANT

Do not delete the default zone and its pools if you are using the default zone and zone group to store data. Also, removing the default zone group, deletes the system user.

Example

```
[ceph: root@host01 /]# radosgw-admin zonegroup delete --rgw-
zonegroup=default
[ceph: root@host01 /]# ceph osd pool rm default.rgw.log default.rgw.log --
yes-i-really-really-mean-it
[ceph: root@host01 /]# ceph osd pool rm default.rgw.meta
default.rgw.meta --yes-i-really-really-mean-it
[ceph: root@host01 /]# ceph osd pool rm default.rgw.control
default.rgw.control --yes-i-really-really-mean-it
[ceph: root@host01 /]# ceph osd pool rm default.rgw.data.root
default.rgw.data.root --yes-i-really-really-mean-it
[ceph: root@host01 /]# ceph osd pool rm default.rgw.gc default.rgw.gc --
yes-i-really-really-mean-it
```

- e. Create a system user:

Syntax

```
radosgw-admin user create --uid=USER_NAME --display-name="USER_NAME" --
access-key=SYSTEM_ACCESS_KEY --secret=SYSTEM_SECRET_KEY --system
```

Example

```
[ceph: root@host01 /]# radosgw-admin user create --uid=zone.user --display-
name="Zone user" --system
```

Make a note of the **access_key** and **secret_key**.

- f. Add the access key and system key to the primary zone:

Syntax

```
radosgw-admin zone modify --rgw-zone=PRIMARY_ZONE_NAME --access-
key=ACCESS_KEY --secret=SECRET_KEY
```

Example

```
[ceph: root@host01 /]# radosgw-admin zone modify --rgw-zone=us-east-1 --access-
key=NE48APYCAODEPLKBCZVQ--
secret=u24GHQWRE3yxxNBnFBzjM4jn14mF1ckQ4EKL6LoW
```

- g. Commit the changes:

Syntax

```
radosgw-admin period update --commit
```

-

Example

```
[ceph: root@host01 /]# radosgw-admin period update --commit
```

- h. Outside the **cephadm** shell, fetch the **FSID** of the storage cluster and the processes:

Example

```
[root@host01 ~]# systemctl list-units | grep ceph
```

- i. Start the Ceph Object Gateway daemon:

Syntax

```
systemctl start ceph-FSID@DAEMON_NAME
systemctl enable ceph-FSID@DAEMON_NAME
```

Example

```
[root@host01 ~]# systemctl start ceph-62a081a6-88aa-11eb-a367-
001a4a000672@rgw.test_realm.us-east-1.host01.ahdtsw.service
[root@host01 ~]# systemctl enable ceph-62a081a6-88aa-11eb-a367-
001a4a000672@rgw.test_realm.us-east-1.host01.ahdtsw.service
```

2. In the Cephadm shell, configure the secondary zone.

- a. Use the URL path, access key, and secret key of the primary zone and primary zone group, and pull the realm configuration from the host:

Syntax

```
radosgw-admin period pull --url=URL_TO_PRIMARY_ZONE_GATEWAY --access-
key=ACCESS_KEY --secret-key=SECRET_KEY
```

Example

```
[ceph: root@host04 /]# radosgw-admin period pull --url=http://10.74.249.26:80 --access-
key=LIPEYZJLTWXRKXS9LPJC --secret-
key=lsAje0AVDNXNw48LjMAimpCpl7VaxJYSnfD0FFKQ
```

- b. Configure a secondary zone:

Syntax

```
radosgw-admin zone create --rgw-zonegroup=ZONE_GROUP_NAME \
--rgw-zone=SECONDARY_ZONE_NAME --
endpoints=http://RGW_SECONDARY_HOSTNAME:RGW_PRIMARY_PORT_NUMBE
R_1_ \
--access-key=SYSTEM_ACCESS_KEY --secret=SYSTEM_SECRET_KEY \
--endpoints=http://FQDN:80 \
[--read-only]
```

Example

```
[ceph: root@host04 /]# radosgw-admin zone create --rgw-zonegroup=us --rgw-zone=us-east-2 --endpoints=http://rgw2:80 --access-key=LIPEYZJLTWXRKXS9LPJC --secret-key=lsAje0AVDNXNw48LjMAimpCpl7VaxJYSnfD0FFKQ
```

- c. Optional: Delete the default zone:



IMPORTANT

Do not delete the default zone and its pools if you are using the default zone and zone group to store data.

Example

```
[ceph: root@host04 /]# radosgw-admin zone rm --rgw-zone=default
[ceph: root@host04 /]# ceph osd pool rm default.rgw.log default.rgw.log --yes-i-really-really-mean-it
[ceph: root@host04 /]# ceph osd pool rm default.rgw.meta default.rgw.meta --yes-i-really-really-mean-it
[ceph: root@host04 /]# ceph osd pool rm default.rgw.control default.rgw.control --yes-i-really-really-mean-it
[ceph: root@host04 /]# ceph osd pool rm default.rgw.data.root default.rgw.data.root --yes-i-really-really-mean-it
[ceph: root@host04 /]# ceph osd pool rm default.rgw.gc default.rgw.gc --yes-i-really-really-mean-it
```

- d. Update the Ceph configuration database:

Syntax

```
ceph config set _SERVICE_NAME_ rgw_zone _SECONDARY_ZONE_NAME_
```

Example

```
[ceph: root@host04 /]# ceph config set rgw rgw_zone us-east-2
```

- e. Commit the changes:

Syntax

```
radosgw-admin period update --commit
```

Example

```
[ceph: root@host04 /]# radosgw-admin period update --commit
```

- f. Outside the Cephadm shell, fetch the FSID of the storage cluster and the processes:

Example

```
[root@host04 ~]# systemctl list-units | grep ceph
```

- g. Start the Ceph Object Gateway daemon:

Syntax

```
systemctl start ceph-FSID@DAEMON_NAME
systemctl enable ceph-FSID@DAEMON_NAME
```

Example

```
[root@host04 ~]# systemctl start ceph-62a081a6-88aa-11eb-a367-
001a4a000672@rgw.test_realm.us-east-2.host04.ahdtsw.service
[root@host04 ~]# systemctl enable ceph-62a081a6-88aa-11eb-a367-
001a4a000672@rgw.test_realm.us-east-2.host04.ahdtsw.service
```

3. Optional: Deploy multi-site Ceph Object Gateways using the placement specification:

Syntax

```
ceph orch apply rgw _NAME_ --realm=_REALM_NAME_ --
zone=_PRIMARY_ZONE_NAME_ --placement="_NUMBER_OF_DAEMONS_
_HOST_NAME_1_ _HOST_NAME_2_"
```

Example

```
[ceph: root@host04 /]# ceph orch apply rgw east --realm=test_realm --zone=us-east-1 --
placement="2 host01 host02"
```

Verification

- Check the synchronization status to verify the deployment:

Example

```
[ceph: root@host04 /]# radosgw-admin sync status
```

11.5. REMOVING THE CEPH OBJECT GATEWAY USING THE CEPH ORCHESTRATOR

You can remove the Ceph object gateway daemons using the **ceph orch rm** command.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.
- Hosts are added to the cluster.
- At least one Ceph object gateway daemon deployed on the hosts.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- a. Remove the service

Syntax

```
ceph orch rm SERVICE_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch rm rgw.test_realm.test_zone_bb
```

Verification

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps
```

Example

```
[ceph: root@host01 /]# ceph orch ps
```

Additional Resources

- See [Deploying the Ceph object gateway using the placement specification](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.
- See [Deploying the Ceph object gateway using the service specification](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

CHAPTER 12. MANAGEMENT OF NFS-GANESHA GATEWAY USING THE CEPH ORCHESTRATOR

As a storage administrator, you can use the Orchestrator with Cephadm in the backend to deploy the NFS-Ganesha gateway. Cephadm deploys NFS Ganesha using a predefined RADOS pool and optional namespace.



NOTE

Red Hat supports CephFS exports only over the NFS v4.0+ protocol.

This section covers the following administrative tasks:

- [Creating the NFS-Ganesha cluster using the Ceph Orchestrator](#) .
- [Deploying the NFS-Ganesha gateway using the command line interface](#) .
- [Deploying the NFS-Ganesha gateway using the service specification](#) .
- [Viewing the NFS-Ganesha cluster information using the Ceph Orchestrator](#) .
- [Fetching the NFS-Ganesha cluster logs using the Ceph Orchestrator](#) .
- [Setting custom NFS-Ganesha configuration using the Ceph Orchestrator](#) .
- [Resetting custom NFS-Ganesha configuration using the Ceph Orchestrator](#) .
- [Deleting the NFS-Ganesha cluster using the Ceph Orchestrator](#) .
- [Removing the NFS Ganesha gateway using the Ceph Orchestrator](#) .

12.1. PREREQUISITES

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.
- Hosts are added to the cluster.
- All manager, monitor and OSD daemons are deployed.

12.2. CREATING THE NFS-GANESHA CLUSTER USING THE CEPH ORCHESTRATOR

You can create an NFS-Ganesha cluster using the **mgr/nfs** module of the Ceph Orchestrator. This module deploys the NFS cluster using Cephadm in the backend.

This creates a common recovery pool for all NFS-Ganesha daemons, new user based on **clusterid**, and a common NFS-Ganesha config RADOS object.

For each daemon, a new user and a common configuration is created in the pool. Although all the clusters have different namespaces with respect the cluster names, they use the same recovery pool.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- All manager, monitor and OSD daemons are deployed.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Enable the **mgr/nfs** module:

Example

```
[ceph: root@host01 /]# ceph mgr module enable nfs
```

3. Create the cluster:

Syntax

```
ceph nfs cluster create CLUSTER_NAME
["HOST_NAME_1,"HOST_NAME_2,"HOST_NAME_3"]
```

The *CLUSTER_NAME* is an arbitrary string and *HOST_NAME_1* is an optional string signifying the hosts to deploy NFS-Ganesha daemons.

Example

```
[ceph: root@host01 /]# ceph nfs cluster create nfsganesha "host01, host02"
```

This creates an NFS_Ganesha cluster **nfsganesha** with one daemon on **host01** and **host02**.

Verification

- List the cluster details:

Example

```
[ceph: root@host01 /]# ceph nfs cluster ls
```

- Show NFS-Ganesha cluster information:

Syntax

```
ceph nfs cluster info CLUSTER_NAME
```

Example

```
[ceph: root@host01 /]# ceph nfs cluster info nfsganesha
```

Additional Resources

- See [Exporting Ceph File System namespaces over the NFS protocol](#) section in the *Red Hat Ceph Storage File System guide* for more information.
- See [Deploying the Ceph daemons using the placement specification](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

12.3. DEPLOYING THE NFS-GANESHA GATEWAY USING THE COMMAND LINE INTERFACE

You can use the Ceph Orchestrator with Cephadm in the backend to deploy the NFS-Ganesha gateway using the placement specification. In this case, you have to create a RADOS pool and create a namespace before deploying the gateway.



NOTE

Red Hat supports CephFS exports only over the NFS v4.0+ protocol.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- All manager, monitor and OSD daemons are deployed.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Create the RADOS pool namespace, and enable the application:

Syntax

```
ceph osd pool create POOL_NAME _
ceph osd pool application enable POOL_NAME freeform/rgw/rbd/cephfs/nfs
rbd pool init -p POOL_NAME
```

Example

```
[ceph: root@host01 /]# ceph osd pool create nfs-ganesha
[ceph: root@host01 /]# ceph osd pool application enable nfs-ganesha nfs
[ceph: root@host01 /]# rbd pool init -p nfs-ganesha
```

3. Deploy NFS-Ganesha gateway using placement specification in the command line interface:

Syntax

```
ceph orch apply nfs SERVICE_ID --pool POOL_NAME --namespace NAMESPACE --
placement="NUMBER_OF_DAEMONS HOST_NAME_1 HOST_NAME_2 HOST_NAME_3"
```

Example

```
[ceph: root@host01 /]# ceph orch apply nfs foo --pool nfs-ganesha --namespace nfs-ns --
placement="2 host01 host02"
```

This deploys an NFS-Ganesha cluster **nfsganesha** with one daemon on **host01** and **host02**.

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=nfs
```

Additional Resources

- See [Deploying the Ceph daemons using the command line interface](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.
- See the [Creating a block device pool](#) section in the *Red Hat Ceph Storage Block Device Guide* for more information.

12.4. DEPLOYING THE NFS-GANESHA GATEWAY USING THE SERVICE SPECIFICATION

You can use the Ceph Orchestrator with Cephadm in the backend to deploy the NFS-Ganesha gateway using the service specification. In this case, you have to create a RADOS pool and create a namespace before deploying the gateway.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Create the RADOS pool, namespace, and enable RBD:

Syntax

```
ceph osd pool create POOL_NAME _  
ceph osd pool application enable POOL_NAME rbd  
rbd pool init -p POOL_NAME
```

Example

```
[ceph: root@host01 /]# ceph osd pool create nfs-ganesha  
[ceph: root@host01 /]#ceph osd pool application enable nfs-ganesha rbd  
[ceph: root@host01 /]#rbd pool init -p nfs-ganesha
```

3. Navigate to the following directory:

Syntax

```
cd /var/lib/ceph/DAEMON_PATH/
```

Example

```
[ceph: root@host01 nfs/]# cd /var/lib/ceph/nfs/
```

If the **nfs** directory does not exist, create a directory in the path.

4. Create the **nfs.yml** file:

Example

```
[ceph: root@host01 nfs]# touch nfs.yml
```

5. Edit the **nfs.yml** file to include the following details:

Syntax

```
service_type: nfs  
service_id: SERVICE_ID  
placement:  
  hosts:  
    - HOST_NAME_1  
    - HOST_NAME_2  
spec:  
  pool: POOL_NAME  
  namespace: NAMESPACE
```

■

Example

```

service_type: nfs
service_id: foo
placement:
  hosts:
    - host01
    - host02
spec:
  pool: nfs-ganesha
  namespace: nfs-ns

```

6. Deploy NFS-Ganesha gateway using service specification:

Syntax

```
ceph orch apply -i FILE_NAME.yml
```

Example

```
[ceph: root@host01 nfs]# ceph orch apply -i nfs.yml
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=nfs
```

Additional Resources

- See the [Creating a block device pool](#) section in the *Red Hat Ceph Storage Block Device Guide* for more information.

12.5. UPDATING THE NFS-GANESHA CLUSTER USING THE CEPH ORCHESTRATOR

You can update the NFS-Ganesha cluster by changing the placement of the daemons on the hosts using the Ceph Orchestrator with Cephadm in the backend.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- All manager, monitor and OSD daemons are deployed.
- NFS-Ganesha cluster created using the **mgr/nfs** module.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Update the cluster:

Syntax

```
ceph orch apply nfs CLUSTER_NAME [HOST_NAME_1,HOST_NAME_2,HOST_NAME_3]
```

The *CLUSTER_NAME* is an arbitrary string, *HOST_NAME_1* is an optional string signifying the hosts to update the deployed NFS-Ganesha daemons.

Example

```
[ceph: root@host01 /]# ceph orch apply nfs nfsganesha "host02"
```

This updates the **nfsganesha** cluster on **host02**.

Verification

- List the cluster details:

Example

```
[ceph: root@host01 /]# ceph nfs cluster ls
```

- Show NFS-Ganesha cluster information:

Syntax

```
ceph nfs cluster info CLUSTER_NAME
```

Example

```
[ceph: root@host01 /]# ceph nfs cluster info nfsganesha
```

- List the hosts, daemons, and processes:+

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

+ .Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=nfs
```

Additional Resources

- See [Creating the NFS-Ganesha cluster using the Ceph Orchestrator](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

12.6. VIEWING THE NFS-GANESHA CLUSTER INFORMATION USING THE CEPH ORCHESTRATOR

You can view the information of the NFS-Ganesha cluster using the Ceph Orchestrator. You can get the information about all the NFS Ganesha clusters or specific clusters with their port, IP address and the name of the hosts on which the cluster is created.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- All manager, monitor and OSD daemons are deployed.
- NFS-Ganesha cluster created using the **mgr/nfs** module.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. View the NFS-Ganesha cluster information:

Syntax

```
ceph nfs cluster info CLUSTER_NAME
```

Example

```
[ceph: root@host01 /]# ceph nfs cluster info nfsganesha
```

```
{
  "nfsganesha": [
    {
      "hostname": "host02",
```

```

    "ip": [
      "10.74.251.164"
    ],
    "port": 2049
  }
]
}

```

Additional Resources

- See [Creating the NFS-Ganesha cluster using the Ceph Orchestrator](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

12.7. FETCHING THE NFS-GANESHA CLUSTER LOGS USING THE CEPH ORCHESTRATOR

With the Ceph Orchestrator, you can fetch the NFS-Ganesha cluster logs. You need to be on the node where the service is deployed.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Cephadm installed on the nodes where NFS is deployed.
- Root-level access to all the nodes.
- Hosts are added to the cluster.
- NFS-Ganesha cluster created using the **mgr/nfs** module.

Procedure

1. As a root user, fetch the *FSID* of the storage cluster:

Example

```
[root@host03 ~]# cephadm ls
```

Copy the *FSID* and the name of the service.

2. Fetch the logs:

Syntax

```
cephadm logs --fsid FSID --name SERVICE_NAME
```

Example

```
[root@host03 ~]# cephadm logs --fsid 499829b4-832f-11eb-8d6d-001a4a000635 --name
nfs.foo.host03
```

Additional Resources

- See [Deploying the Ceph daemons using the placement specification](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

12.8. SETTING CUSTOM NFS-GANESHA CONFIGURATION USING THE CEPH ORCHESTRATOR

The NFS-Ganesha cluster is defined in default configuration blocks. Using Ceph Orchestrator you can customize the configuration and that will have precedence over the default configuration blocks.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- All manager, monitor and OSD daemons are deployed.
- NFS-Ganesha cluster created using the **mgr/nfs** module.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. The following is an example of the default configuration of NFS-Ganesha cluster:

Example

```
# {{ cephadm_managed }}
NFS_CORE_PARAM {
    Enable_NLM = false;
    Enable_RQUOTA = false;
    Protocols = 4;
}

MDCACHE {
    Dir_Chunk = 0;
}

EXPORT_DEFAULTS {
    Attr_Expiration_Time = 0;
}

NFSv4 {
    Delegations = false;
    RecoveryBackend = 'rados_cluster';
    Minor_Versions = 1, 2;
}
```

```

RADOS_KV {
    UserId = "{{ user }}";
    nodeid = "{{ nodeid}}";
    pool = "{{ pool }}";
    namespace = "{{ namespace }}";
}

RADOS_URLS {
    UserId = "{{ user }}";
    watch_url = "{{ url }}";
}

RGW {
    cluster = "ceph";
    name = "client.{{ rgw_user }}";
}

%url {{ url }}

```

3. Customize the NFS-Ganesha cluster configuration. The following are two examples for customizing the configuration:

- Change the log level:

Example

```

LOG {
    COMPONENTS {
        ALL = FULL_DEBUG;
    }
}

```

- Add custom export block:
 - a. Create the user.



NOTE

User specified in FSAL blocks should have proper caps for NFS-Ganesha daemons to access the Ceph cluster.

Syntax

```

ceph auth get-or-create client.USER_NAME mon 'allow r' osd 'allow rw
pool=POOL_NAME namespace=NFS_CLUSTER_NAME, allow rw tag cephfs
data=FILE_SYSTEM_NAME mds 'allow rw path=EXPORT_PATH'

```

Example

```

[ceph: root@host01 /]# ceph auth get-or-create client.nfstest1 mon 'allow r' osd 'allow
rw pool=nfsganesha namespace=nfs_cluster_name, allow rw tag cephfs
data=filesystem_name' mds 'allow rw path=export_path

```

- b. Navigate to the following directory:

Syntax

```
cd /var/lib/ceph/DAEMON_PATH/
```

Example

```
[ceph: root@host01 /]# cd /var/lib/ceph/nfs/
```

If the **nfs** directory does not exist, create a directory in the path.

- c. Create a new configuration file:

Syntax

```
touch PATH_TO_CONFIG_FILE
```

Example

```
[ceph: root@host01 nfs]# touch nfs-ganesha.conf
```

- d. Edit the configuration file by adding the custom export block. It creates a single export and that is managed by the Ceph NFS export interface.

Syntax

```
EXPORT {
  Export_Id = NUMERICAL_ID;
  Transports = TCP;
  Path = PATH_WITHIN_CEPHFS;
  Pseudo = BINDING;
  Protocols = 4;
  Access_Type = PERMISSIONS;
  Attr_Expiration_Time = 0;
  Squash = None;
  FSAL {
    Name = CEPH;
    Filesystem = "FILE_SYSTEM_NAME";
    User_Id = "USER_NAME";
    Secret_Access_Key = "USER_SECRET_KEY";
  }
}
```

Example

```
EXPORT {
  Export_Id = 100;
  Transports = TCP;
  Path = /;
  Pseudo = /ceph/;
  Protocols = 4;
  Access_Type = RW;
  Attr_Expiration_Time = 0;
  Squash = None;
```

```
FSAL {  
    Name = CEPH;  
    Filesystem = "filesystem name";  
    User_Id = "user id";  
    Secret_Access_Key = "secret key";  
}  
}
```

4. Apply the new configuration the cluster:

Syntax

```
ceph nfs cluster config set _CLUSTER_NAME_ -i _PATH_TO_CONFIG_FILE_
```

Example

```
[ceph: root@host01 nfs]# ceph nfs cluster config set nfs-ganesha -i /root/nfs-ganesha.conf
```

This also restarts the service for the custom configuration.

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=nfs
```

- Verify the custom configuration:

Syntax

```
./bin/rados -p POOL_NAME -N CLUSTER_NAME get userconf-nfs.CLUSTER_NAME -
```

Example

```
[ceph: root@host01 /]# ./bin/rados -p nfs-ganesha -N nfsganesha get userconf-nfs.nfsganesha -
```

Additional Resources

- See the [Resetting custom NFS-Ganesha configuration using the Ceph Orchestrator](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

12.9. RESETTING CUSTOM NFS-GANESHA CONFIGURATION USING THE CEPH ORCHESTRATOR

Using the Ceph Orchestrator, you can reset the user-defined configuration to the default configuration.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- All manager, monitor and OSD daemons are deployed.
- NFS-Ganesha deployed using the **mgr/nfs** module.
- Custom NFS cluster configuration is set-up

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Reset the NFS_Ganesha configuration:

Syntax

```
ceph nfs cluster config reset CLUSTER_NAME
```

Example

```
[ceph: root@host01 /]# ceph nfs cluster config reset nfsganesha
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=nfs
```

- Verify the custom configuration is deleted:

Syntax

```
./bin/rados -p POOL_NAME -N CLUSTER_NAME get userconf-nfs.CLUSTER_NAME -
```

Example

```
[ceph: root@host01 /]# ./bin/rados -p nfs-ganesha -N nfsganesha get userconf-nfs.nfsganesha -
```

Additional Resources

- See [Creating the NFS-Ganesha cluster using the Ceph Orchestrator](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.
- See the [Setting custom NFS-Ganesha configuration using the Ceph Orchestrator](#) section in the *Red Hat Ceph Storage Operations Guide* for more information..

12.10. DELETING THE NFS-GANESHA CLUSTER USING THE CEPH ORCHESTRATOR

You can use the Ceph Orchestrator with Cephadm in the backend to delete the NFS-Ganesha cluster.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- All manager, monitor and OSD daemons are deployed.
- NFS-Ganesha cluster created using the **mgr/nfs** module.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Delete the cluster:

Syntax

```
ceph nfs cluster delete CLUSTER_NAME
```

The `CLUSTER_NAME` is an arbitrary string.

Example

```
[ceph: root@host01 /]# ceph nfs cluster delete nfsganesha
NFS Cluster Deleted Successfully
```

Verification

- List the cluster details:

Example

```
[ceph: root@host01 /]# ceph nfs cluster ls
```

Additional Resources

- See [Exporting Ceph File System namespaces over the NFS protocol](#) section in the *Red Hat Ceph Storage File System guide* for more information.
- See [Creating the NFS-Ganesha cluster using the Ceph Orchestrator](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

12.11. REMOVING THE NFS-GANESHA GATEWAY USING THE CEPH ORCHESTRATOR

You can remove the NFS-Ganesha gateway using the **ceph orch rm** command.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.
- Hosts are added to the cluster.
- At least one NFS-Ganesha gateway deployed on the hosts.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

3. Remove the service

Syntax

```
ceph orch rm SERVICE_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch rm nfs.foo
```

Verification

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps
```

Example

```
[ceph: root@host01 /]# ceph orch ps
```

Additional Resources

- See [Deploying the Ceph daemons using the placement specification](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.
- See [Deploying the NFS-Ganesha gateway using the service specification](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.
- See [Deploying the NFS-Ganesha gateway using the placement specification](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

CHAPTER 13. MANAGEMENT OF ISCSI GATEWAY USING THE CEPH ORCHESTRATOR

As a storage administrator, you can use Ceph Orchestrator to deploy the iSCSI gateway. The iSCSI Gateway presents a Highly Available (HA) iSCSI target that exports RADOS Block Device (RBD) images as SCSI disks.

You can deploy an iSCSI gateway by either using the placement specification or the service specification, like an YAML file.

This section covers the following administrative tasks:

- [Deploying the iSCSI gateway using the placement specification](#) .
- [Deploying the iSCSI gateway using the service specification](#) .
- [Removing the iSCSI gateway using the Ceph Orchestrator](#) .

13.1. PREREQUISITES

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.
- Hosts are added to the cluster.
- All manager, monitor and OSD daemons are deployed.

13.2. DEPLOYING THE ISCSI GATEWAY USING THE COMMAND LINE INTERFACE

Using the Ceph Orchestrator, you can deploy the iSCSI gateway using the **ceph orch** command in the command line interface.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- All manager, monitor and OSD daemons are deployed.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Create the pool:

Syntax

```
ceph osd pool create POOL_NAME
```

Example

```
[ceph: root@host01 /]# ceph osd pool create mypool
```

3. Deploy iSCSI gateway using command line interface:

Syntax

```
ceph orch apply iscsi POOLNAME admin admin --placement="NUMBER_OF_DAEMONS  
HOST_NAME_1 HOST_NAME_2"
```

Example

```
[ceph: root@host01 /]# ceph orch apply iscsi mypool admin admin --placement="1 host01"
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts and process:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=iscsi
```

13.3. DEPLOYING THE ISCSI GATEWAY USING THE SERVICE SPECIFICATION

Using the Ceph Orchestrator, you can deploy the iSCSI gateway using the service specification.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Hosts are added to the cluster.
- All manager, monitor and OSD daemons are deployed.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. Navigate to the following directory:

Syntax

```
cd /var/lib/ceph/DAEMON_PATH/
```

Example

```
[ceph: root@host01 /]# cd /var/lib/ceph/iscsi/
```

Note: If the directory `iscsi` does not exist, create it.

3. Create the **iscsi.yml** file:

Example

```
[ceph: root@host01 iscsi]# touch iscsi.yml
```

4. Edit the **iscsi.yml** file to include the following details:

Syntax

```
service_type: iscsi
service_id: iscsi
placement:
  hosts:
    - HOST_NAME_1
    - HOST_NAME_2
spec:
  pool: POOL_NAME # RADOS pool where ceph-iscsi config data is stored.
  trusted_ip_list: "IP_ADDRESS_1,IP_ADDRESS_2" # optional
  api_port: ... # optional
  api_user: API_USERNAME # optional
  api_password: API_PASSWORD # optional
  api_secure: true/false # optional
  ssl_cert: | # optional
  ...
  ssl_key: | # optional
  ...
```

Example

```
service_type: iscsi
service_id: iscsi
placement:
  hosts:
```

```
- host01
spec:
pool: mypool
```

5. Deploy iSCSI gateway using service specification:

Syntax

```
ceph orch apply -i FILE_NAME.yaml
```

Example

```
[ceph: root@host01 iscsi]# ceph orch apply -i iscsi.yaml
```

Verification

- List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps --daemon_type=DAEMON_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch ps --daemon_type=iscsi
```

13.4. REMOVING THE ISCSI GATEWAY USING THE CEPH ORCHESTRATOR

You can remove the iSCSI gateway daemons using the **ceph orch rm** command.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Root-level access to all the nodes.
- Hosts are added to the cluster.
- At least one iSCSI gateway daemon deployed on the hosts.

Procedure

1. Log into the Cephadm shell:

Example

```
[root@host01 ~]# cephadm shell
```

2. List the service:

Example

```
[ceph: root@host01 /]# ceph orch ls
```

- a. Remove the service

Syntax

```
ceph orch rm SERVICE_NAME
```

Example

```
[ceph: root@host01 /]# ceph orch rm iscsi.iscsi
```

Verification

- List the hosts, daemons, and processes:

Syntax

```
ceph orch ps
```

Example

```
[ceph: root@host01 /]# ceph orch ps
```

Additional Resources

- See [Deploying the iSCSI gateway using the command line interface](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.
- See [Deploying the iSCSI gateway using the service specification](#) section in the *Red Hat Ceph Storage Operations Guide* for more information.

CHAPTER 14. HANDLING A NODE FAILURE

As a storage administrator, you can experience a whole node failing within the storage cluster, and handling a node failure is similar to handling a disk failure. With a node failure, instead of Ceph recovering placement groups (PGs) for only one disk, all PGs on the disks within that node must be recovered. Ceph will detect that the OSDs are all down and automatically start the recovery process, known as self-healing.

There are three node failure scenarios. Here is the high-level workflow for each scenario when replacing a node:

- Replacing the node, but using the root and Ceph OSD disks from the failed node.
 1. Disable backfilling.
 2. Replace the node, taking the disks from old node, and adding them to the new node.
 3. Enable backfilling.
- Replacing the node, reinstalling the operating system, and using the Ceph OSD disks from the failed node.
 1. Disable backfilling.
 2. Create a backup of the Ceph configuration.
 3. Replace the node and add the Ceph OSD disks from failed node.
 4. Configuring disks as JBOD.
 5. Install the operating system.
 6. Restore the Ceph configuration.
 7. Add the new node to the storage cluster using the Ceph Orchestrator commands and Ceph daemons are placed automatically on the respective node.
 8. Enable backfilling.
- Replacing the node, reinstalling the operating system, and using all new Ceph OSDs disks.
 1. Disable backfilling.
 2. Remove all OSDs on the failed node from the storage cluster.
 3. Create a backup of the Ceph configuration.
 4. Replace the node and add the Ceph OSD disks from failed node.
 - a. Configuring disks as JBOD.
 5. Install the operating system.
 6. Add the new node to the storage cluster using the Ceph Orchestrator commands and Ceph daemons are placed automatically on the respective node.
 7. Enable backfilling.

14.1. PREREQUISITES

- A running Red Hat Ceph Storage cluster.
- A failed node.

14.2. CONSIDERATIONS BEFORE ADDING OR REMOVING A NODE

One of the outstanding features of Ceph is the ability to add or remove Ceph OSD nodes at run time. This means that you can resize the storage cluster capacity or replace hardware without taking down the storage cluster.

The ability to serve Ceph clients while the storage cluster is in a **degraded** state also has operational benefits. For example, you can add or remove or replace hardware during regular business hours, rather than working overtime or on weekends. However, adding and removing Ceph OSD nodes can have a significant impact on performance.

Before you add or remove Ceph OSD nodes, consider the effects on storage cluster performance:

- Whether you are expanding or reducing the storage cluster capacity, adding or removing Ceph OSD nodes induces backfilling as the storage cluster rebalances. During that rebalancing time period, Ceph uses additional resources, which can impact storage cluster performance.
- In a production Ceph storage cluster, a Ceph OSD node has a particular hardware configuration that facilitates a particular type of storage strategy.
- Since a Ceph OSD node is part of a CRUSH hierarchy, the performance impact of adding or removing a node typically affects the performance of pools that use the CRUSH ruleset.

Additional Resources

- See the [Red Hat Ceph Storage Storage Strategies Guide](#) for more details.

14.3. PERFORMANCE CONSIDERATIONS

The following factors typically affect a storage cluster's performance when adding or removing Ceph OSD nodes:

- Ceph clients place load on the I/O interface to Ceph; that is, the clients place load on a pool. A pool maps to a CRUSH ruleset. The underlying CRUSH hierarchy allows Ceph to place data across failure domains. If the underlying Ceph OSD node involves a pool that is experiencing high client load, the client load could significantly affect recovery time and reduce performance. Because write operations require data replication for durability, write-intensive client loads in particular can increase the time for the storage cluster to recover.
- Generally, the capacity you are adding or removing affects the storage cluster's time to recover. In addition, the storage density of the node you add or remove might also affect recovery times. For example, a node with 36 OSDs typically takes longer to recover than a node with 12 OSDs.
- When removing nodes, you **MUST** ensure that you have sufficient spare capacity so that you will not reach **full ratio** or **near full ratio**. If the storage cluster reaches **full ratio**, Ceph will suspend write operations to prevent data loss.

- A Ceph OSD node maps to at least one Ceph CRUSH hierarchy, and the hierarchy maps to at least one pool. Each pool that uses a CRUSH ruleset experiences a performance impact when Ceph OSD nodes are added or removed.
- Replication pools tend to use more network bandwidth to replicate deep copies of the data, whereas erasure coded pools tend to use more CPU to calculate **k+m** coding chunks. The more copies that exist of the data, the longer it takes for the storage cluster to recover. For example, a larger pool or one that has a greater number of **k+m** chunks will take longer to recover than a replication pool with fewer copies of the same data.
- Drives, controllers and network interface cards all have throughput characteristics that might impact the recovery time. Generally, nodes with higher throughput characteristics, such as 10 Gbps and SSDs, recover more quickly than nodes with lower throughput characteristics, such as 1 Gbps and SATA drives.

14.4. RECOMMENDATIONS FOR ADDING OR REMOVING NODES

Red Hat recommends adding or removing one OSD at a time within a node and allowing the storage cluster to recover before proceeding to the next OSD. This helps to minimize the impact on storage cluster performance. Note that if a node fails, you might need to change the entire node at once, rather than one OSD at a time.

To remove an OSD:

- Using [Removing the OSD daemons using the Ceph Orchestrator](#).

To add an OSD:

- Using [Deploying Ceph OSDs on all available devices](#).
- Using [Deploying Ceph OSDs using advanced service specification](#).
- Using [Deploying Ceph OSDs on specific devices and hosts](#).

When adding or removing Ceph OSD nodes, consider that other ongoing processes also affect storage cluster performance. To reduce the impact on client I/O, Red Hat recommends the following:

Calculate capacity

Before removing a Ceph OSD node, ensure that the storage cluster can backfill the contents of all its OSDs without reaching the **full ratio**. Reaching the **full ratio** will cause the storage cluster to refuse write operations.

Temporarily disable scrubbing

Scrubbing is essential to ensuring the durability of the storage cluster's data; however, it is resource intensive. Before adding or removing a Ceph OSD node, disable scrubbing and deep-scrubbing and let the current scrubbing operations complete before proceeding.

```
ceph osd set noscrub
ceph osd set nodeep-scrub
```

Once you have added or removed a Ceph OSD node and the storage cluster has returned to an **active+clean** state, unset the **noscrub** and **nodeep-scrub** settings.

```
ceph osd unset noscrub
ceph osd unset nodeep-scrub
```

Limit backfill and recovery

If you have reasonable data durability, there is nothing wrong with operating in a **degraded** state. For example, you can operate the storage cluster with **osd_pool_default_size = 3** and **osd_pool_default_min_size = 2**. You can tune the storage cluster for the fastest possible recovery time, but doing so significantly affects Ceph client I/O performance. To maintain the highest Ceph client I/O performance, limit the backfill and recovery operations and allow them to take longer.

```
osd_max_backfills = 1
osd_recovery_max_active = 1
osd_recovery_op_priority = 1
```

You can also consider setting the sleep and delay parameters such as, **osd_recovery_sleep**.

Increase the number of placement groups

Finally, if you are expanding the size of the storage cluster, you may need to increase the number of placement groups. If you determine that you need to expand the number of placement groups, Red Hat recommends making incremental increases in the number of placement groups. Increasing the number of placement groups by a significant amount will cause a considerable degradation in performance.

14.5. ADDING A CEPH OSD NODE

To expand the capacity of the Red Hat Ceph Storage cluster, you can add an OSD node.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- A provisioned node with a network connection.

Procedure

1. Verify that other nodes in the storage cluster can reach the new node by its short host name.
2. Temporarily disable scrubbing:

Example

```
[ceph: root@host01 /]# ceph osd set noscrub
[ceph: root@host01 /]# ceph osd set nodeep-scrub
```

3. Limit the backfill and recovery features:

Syntax

```
ceph tell DAEMON_TYPE.* injectargs --OPTION_NAME VALUE [--OPTION_NAME VALUE]
```

Example

```
[ceph: root@host01 /]# ceph tell osd.* injectargs --osd-max-backfills 1 --osd-recovery-max-active 1 --osd-recovery-op-priority 1
```

4. Extract the cluster's public SSH keys to a folder:

Syntax

```
ceph cephadm get-pub-key > ~/PATH
```

Example

```
[ceph: root@host01 /]# ceph cephadm get-pub-key > ~/ceph.pub
```

5. Copy ceph cluster's public SSH keys to the root user's **authorized_keys** file on the new host:

Syntax

```
ssh-copy-id -f -i ~/PATH root@HOST_NAME_2
```

Example

```
[ceph: root@host01 /]# ssh-copy-id -f -i ~/ceph.pub root@host02
```

6. Add the new node to the CRUSH map:

Syntax

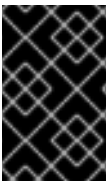
```
ceph orch host add NODE_NAME IP_ADDRESS
```

Example

```
[ceph: root@host01 /]# ceph orch host add host02 168.20.20.2
```

7. Add an OSD for each disk on the node to the storage cluster.

- Using [Deploying Ceph OSDs on all available devices](#).
- Using [Deploying Ceph OSDs using advanced service specification](#).
- Using [Deploying Ceph OSDs on specific devices and hosts](#).



IMPORTANT

When adding an OSD node to a Red Hat Ceph Storage cluster, Red Hat recommends adding one OSD daemon at a time and allowing the cluster to recover to an **active+clean** state before proceeding to the next OSD.

Additional Resources

- See the [Setting a Specific Configuration Setting at Runtime](#) section in the *Red Hat Ceph Storage Configuration Guide* for more details.

- See [Adding a Bucket](#) and [Moving a Bucket](#) sections in the *Red Hat Ceph Storage Storage Strategies Guide* for details on placing the node at an appropriate location in the CRUSH hierarchy,.

14.6. REMOVING A CEPH OSD NODE

To reduce the capacity of a storage cluster, remove an OSD node.



WARNING

Before removing a Ceph OSD node, ensure that the storage cluster can backfill the contents of all OSDs without reaching the **full ratio**. Reaching the **full ratio** will cause the storage cluster to refuse write operations.

Prerequisites

- A running Red Hat Ceph Storage cluster.
- Root-level access to all nodes in the storage cluster.

Procedure

1. Check the storage cluster's capacity:

Syntax

```
ceph df
rados df
ceph osd df
```

2. Temporarily disable scrubbing:

Syntax

```
ceph osd set noscrub
ceph osd set nodeep-scrub
```

3. Limit the backfill and recovery features:

Syntax

```
ceph tell DAEMON_TYPE.* injectargs --OPTION_NAME VALUE [--OPTION_NAME VALUE]
```

Example

```
[ceph: root@host01 /]# ceph tell osd.* injectargs --osd-max-backfills 1 --osd-recovery-max-active 1 --osd-recovery-op-priority 1
```

4. Remove each OSD on the node from the storage cluster:

- Using [Removing the OSD daemons using the Ceph Orchestrator](#).



IMPORTANT

When removing an OSD node from the storage cluster, Red Hat recommends removing one OSD at a time within the node and allowing the cluster to recover to an **active+clean** state before proceeding to remove the next OSD.

- After you remove an OSD, check to verify that the storage cluster is not getting to the **near-full ratio**:

Syntax

```
ceph -s
ceph df
```

- Repeat this step until all OSDs on the node are removed from the storage cluster.
5. Once all OSDs are removed, remove the host:
- Using [Removing hosts using the Ceph Orchestrator](#).

Additional Resources

- See the [Setting a specific configuration setting at runtime](#) section in the *Red Hat Ceph Storage Configuration Guide* for more details.

14.7. SIMULATING A NODE FAILURE

To simulate a hard node failure, power off the node and reinstall the operating system.

Prerequisites

- A healthy running Red Hat Ceph Storage cluster.
- Root-level access to all nodes on the storage cluster.

Procedure

- Check the storage cluster's capacity to understand the impact of removing the node:

Example

```
[ceph: root@host01 /]# ceph df
[ceph: root@host01 /]# rados df
[ceph: root@host01 /]# ceph osd df
```

- Optionally, disable recovery and backfilling:

Example

-

```
[ceph: root@host01 /]# ceph osd set noout
[ceph: root@host01 /]# ceph osd set noscrub
[ceph: root@host01 /]# ceph osd set nodeep-scrub
```

3. Shut down the node.
4. If you are changing the host name, remove the node from CRUSH map:

Example

```
[ceph: root@host01 /]# ceph osd crush rm host03
```

5. Check the status of the storage cluster:

Example

```
[ceph: root@host01 /]# ceph -s
```

6. Reinstall the operating system on the node.
7. Add the new node:
 - Using the [Adding hosts using the Ceph Orchestrator](#).
8. Optionally, enable recovery and backfilling:

Example

```
[ceph: root@host01 /]# ceph osd unset noout
[ceph: root@host01 /]# ceph osd unset noscrub
[ceph: root@host01 /]# ceph osd unset nodeep-scrub
```

9. Check Ceph's health:

Example

```
[ceph: root@host01 /]# ceph -s
```

Additional Resources

- See the [Red Hat Ceph Storage Installation Guide](#) for more details.

CHAPTER 15. HANDLING A DATA CENTER FAILURE

As a storage administrator, you can take preventive measures to avoid a data center failure. These preventive measures include:

- Configuring the data center infrastructure.
- Setting up failure domains within the CRUSH map hierarchy.
- Designating failure nodes within the domains.

15.1. PREREQUISITES

- A healthy running Red Hat Ceph Storage cluster.
- Root-level access to all nodes in the storage cluster.

15.2. AVOIDING A DATA CENTER FAILURE

Configuring the data center infrastructure

Each data center within a stretch cluster can have a different storage cluster configuration to reflect local capabilities and dependencies. Set up replication between the data centers to help preserve the data. If one data center fails, the other data centers in the storage cluster contain copies of the data.

Setting up failure domains within the CRUSH map hierarchy

Failure, or failover, domains are redundant copies of domains within the storage cluster. If an active domain fails, the failure domain becomes the active domain.

By default, the CRUSH map lists all nodes in a storage cluster within a flat hierarchy. However, for best results, create a logical hierarchical structure within the CRUSH map. The hierarchy designates the domains to which each node belongs and the relationships among those domains within the storage cluster, including the failure domains. Defining the failure domains for each domain within the hierarchy improves the reliability of the storage cluster.

When planning a storage cluster that contains multiple data centers, place the nodes within the CRUSH map hierarchy so that if one data center goes down, the rest of the storage cluster stays up and running.

Designating failure nodes within the domains

If you plan to use three-way replication for data within the storage cluster, consider the location of the nodes within the failure domain. If an outage occurs within a data center, it is possible that some data might reside in only one copy. When this scenario happens, there are two options:

- Leave the data in read-only status with the standard settings.
- Live with only one copy for the duration of the outage.

With the standard settings, and because of the randomness of data placement across the nodes, not all the data will be affected, but some data can have only one copy and the storage cluster would revert to read-only mode. However, if some data exist in only one copy, the storage cluster reverts to read-only mode.

15.3. HANDLING A DATA CENTER FAILURE

Red Hat Ceph Storage can withstand catastrophic failures to the infrastructure, such as losing one of the data centers in a stretch cluster. For the standard object store use case, configuring all three data centers can be done independently with replication set up between them. In this scenario, the storage cluster configuration in each of the data centers might be different, reflecting the local capabilities and dependencies.

A logical structure of the placement hierarchy should be considered. A proper CRUSH map can be used, reflecting the hierarchical structure of the failure domains within the infrastructure. Using logical hierarchical definitions improves the reliability of the storage cluster, versus using the standard hierarchical definitions. Failure domains are defined in the CRUSH map. The default CRUSH map contains all nodes in a flat hierarchy. In a three data center environment, such as a stretch cluster, the placement of nodes should be managed in a way that one data center can go down, but the storage cluster stays up and running. Consider which failure domain a node resides in when using 3-way replication for the data.

In the example below, the resulting map is derived from the initial setup of the storage cluster with 6 OSD nodes. In this example, all nodes have only one disk and hence one OSD. All of the nodes are arranged under the default *root*, that is the standard *root* of the hierarchy tree. Because there is a weight assigned to two of the OSDs, these OSDs receive fewer chunks of data than the other OSDs. These nodes were introduced later with bigger disks than the initial OSD disks. This does not affect the data placement to withstand a failure of a group of nodes.

Example

```
[ceph: root@host01 /]# ceph osd tree
ID WEIGHT  TYPE NAME          UP/DOWN REWEIGHT PRIMARY-AFFINITY
-1 0.33554  root default
-2 0.04779  host host03
 0 0.04779  osd.0          up 1.00000    1.00000
-3 0.04779  host host02
 1 0.04779  osd.1          up 1.00000    1.00000
-4 0.04779  host host01
 2 0.04779  osd.2          up 1.00000    1.00000
-5 0.04779  host host04
 3 0.04779  osd.3          up 1.00000    1.00000
-6 0.07219  host host06
 4 0.07219  osd.4          up 0.79999    1.00000
-7 0.07219  host host05
 5 0.07219  osd.5          up 0.79999    1.00000
```

Using logical hierarchical definitions to group the nodes into same data center can achieve data placement maturity. Possible definition types of *root*, *datacenter*, *rack*, *row* and *host* allow the reflection of the failure domains for the three data center stretch cluster:

- Nodes host01 and host02 reside in data center 1 (DC1)
- Nodes host03 and host05 reside in data center 2 (DC2)
- Nodes host04 and host06 reside in data center 3 (DC3)
- All data centers belong to the same structure (allDC)

Since all OSDs in a host belong to the host definition there is no change needed. All the other assignments can be adjusted during runtime of the storage cluster by:

- Defining the *bucket* structure with the following commands:

```
ceph osd crush add-bucket allDC root
ceph osd crush add-bucket DC1 datacenter
ceph osd crush add-bucket DC2 datacenter
ceph osd crush add-bucket DC3 datacenter
```

- Moving the nodes into the appropriate place within this structure by modifying the CRUSH map:

```
ceph osd crush move DC1 root=allDC
ceph osd crush move DC2 root=allDC
ceph osd crush move DC3 root=allDC
ceph osd crush move host01 datacenter=DC1
ceph osd crush move host02 datacenter=DC1
ceph osd crush move host03 datacenter=DC2
ceph osd crush move host05 datacenter=DC2
ceph osd crush move host04 datacenter=DC3
ceph osd crush move host06 datacenter=DC3
```

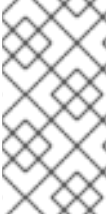
Within this structure any new hosts can be added too, as well as new disks. By placing the OSDs at the right place in the hierarchy the CRUSH algorithm is changed to place redundant pieces into different failure domains within the structure.

The above example results in the following:

Example

```
[ceph: root@host01 /]# ceph osd tree
ID WEIGHT TYPE NAME          UP/DOWN REWEIGHT PRIMARY-AFFINITY
-8 6.00000 root allDC
-9 2.00000 datacenter DC1
-4 1.00000 host host01
 2 1.00000 osd.2 up 1.00000 1.00000
-3 1.00000 host host02
 1 1.00000 osd.1 up 1.00000 1.00000
-10 2.00000 datacenter DC2
-2 1.00000 host host03
 0 1.00000 osd.0 up 1.00000 1.00000
-7 1.00000 host host05
 5 1.00000 osd.5 up 0.79999 1.00000
-11 2.00000 datacenter DC3
-6 1.00000 host host06
 4 1.00000 osd.4 up 0.79999 1.00000
-5 1.00000 host host04
 3 1.00000 osd.3 up 1.00000 1.00000
-1 0 root default
```

The listing from above shows the resulting CRUSH map by displaying the osd tree. Easy to see is now how the hosts belong to a data center and all data centers belong to the same top level structure but clearly distinguishing between locations.

**NOTE**

Placing the data in the proper locations according to the map works only properly within the healthy cluster. Misplacement might happen under circumstances, when some OSDs are not available. Those misplacements will be corrected automatically once it is possible to do so.

Additional Resources

- See the [CRUSH administration](#) chapter in the Red Hat Ceph Storage Storage Strategies Guide for more information.