



Red Hat OpenShift Container Storage 4.8

Recovering a Metro-DR stretch cluster

Disaster recovery tasks for cluster and storage administrators

Red Hat OpenShift Container Storage 4.8 Recovering a Metro-DR stretch cluster

Disaster recovery tasks for cluster and storage administrators

Legal Notice

Copyright © 2022 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 explains how to recover from a metro disaster in Red Hat OpenShift Container Storage. This is a technology preview feature and is available only for deployments using local storage devices. Technology Preview features are not supported with Red Hat production service level agreements (SLAs) and might not be functionally complete. Red Hat does not recommend using them in production. These features provide early access to upcoming product features, enabling customers to test functionality and provide feedback during the development process.

Table of Contents

MAKING OPEN SOURCE MORE INCLUSIVE	3
PROVIDING FEEDBACK ON RED HAT DOCUMENTATION	4
CHAPTER 1. OVERVIEW	5
CHAPTER 2. UNDERSTANDING ZONE FAILURE	6
CHAPTER 3. RECOVERY FOR ZONE-AWARE HA APPLICATIONS WITH RWX STORAGE	7
CHAPTER 4. RECOVERY FOR HA APPLICATIONS WITH RWX STORAGE	8
CHAPTER 5. RECOVERING APPLICATIONS WITH RWO STORAGE	9
CHAPTER 6. RECOVERY FOR STATEFULSET PODS	11

MAKING OPEN SOURCE MORE INCLUSIVE

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](#).

PROVIDING FEEDBACK ON RED HAT DOCUMENTATION

We appreciate your input on our documentation. Do let us know how we can make it better. To give feedback:

- For simple comments on specific passages:
 1. Make sure you are viewing the documentation in the *Multi-page HTML* format. In addition, ensure you see the **Feedback** button in the upper right corner of the document.
 2. Use your mouse cursor to highlight the part of text that you want to comment on.
 3. Click the **Add Feedback** pop-up that appears below the highlighted text.
 4. Follow the displayed instructions.
- For submitting more complex feedback, create a Bugzilla ticket:
 1. Go to the [Bugzilla](#) website.
 2. In the **Component** section, choose **documentation**.
 3. Fill in the **Description** field with your suggestion for improvement. Include a link to the relevant part(s) of documentation.
 4. Click **Submit Bug**.

CHAPTER 1. OVERVIEW

Given the Metro disaster recovery stretch cluster is to provide resiliency in the face of a complete or partial site outage, it is important to understand the different methods of recovery for applications and their storage.

How the application is architected will determine how soon it can be available again on the active zone.

Depending on the site outage, there can be different methods of recovery for applications and their storage. The recovery time depends on the application architecture. The different methods of recovery are as follows:

- Recovery for zone-aware HA applications with RWX storage
- Recovery for HA applications with RWX storage
- Recovery for applications with RWO storage
- Recovery for StatefulSet pods

CHAPTER 2. UNDERSTANDING ZONE FAILURE

For purposes of this section, we will consider a zone failure to be a failure where all OpenShift Container Platform nodes, masters and workers, in a zone are no longer communicating with the resources in the second data zone (e.g., nodes powered down). If communication between data zones is partially still working (intermittent down/up), steps should be taken by cluster, storage, network admins to sever the communication path between the data zones for recovery to succeed.

CHAPTER 3. RECOVERY FOR ZONE-AWARE HA APPLICATIONS WITH RWX STORAGE

Applications that are deployed with **topologyKey: topology.kubernetes.io/zone**, have one or more replicas scheduled in each data zone, and are using shared storage (i.e., RWX cephfs volume) will recover on the active zone within 30-60 seconds for new connections. The short pause is for **HAProxy** to refresh connections if a router pod is now offline in the failed data zone.

An example of this type of application is detailed in the [Install Zone Aware Sample Application](#) section.

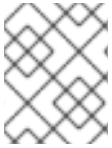


NOTE

When you install the Sample Application, test the failure of a data zone by powering off the OpenShift Container Platform nodes (at least the nodes with OpenShift Container storage devices) to validate that your file-uploader application is available and that new files can be uploaded.

CHAPTER 4. RECOVERY FOR HA APPLICATIONS WITH RWX STORAGE

Applications that are using **topologyKey: kubernetes.io/hostname** or no topology configuration whatsoever, have no protection against all of the application replicas being in the same zone.



NOTE

This can happen even with *podAntiAffinity* and **topologyKey: kubernetes.io/hostname** in the **Pod** spec because this anti-affinity rule is host-based and not zone-based.

If this happens and all replicas are located in the zone that fails, the application using RWX storage will take 6-8 minutes to recover on the active zone. This pause is for the OpenShift Container Platform nodes in the failed zone to become **NotReady** (60 seconds) and then for the default pod eviction timeout to expire (300 seconds).

CHAPTER 5. RECOVERING APPLICATIONS WITH RWO STORAGE

Applications that use RWO storage (ReadWriteOnce) have a known behavior described in this [Kubernetes issue](#). Because of this issue, if there is a data zone failure any application pods in that zone mounting RWO volumes (for example: **cephrbd** based volumes) are stuck with **Terminating** status after 6-8 minutes and will not be recreated on the active zone without manual intervention.

Check the OpenShift Container Platform nodes with a status of **NotReady**. It may have an issue that prevents them from communicating with the OpenShift control plane. They may still be performing IO operations against persistent volumes in-spite of this communication issue.

If two pods are concurrently writing to the same RWO volume, there is a risk of data corruption. Some measure must be taken to ensure that processes on the **NotReady** node are terminated or blocked until they can be terminated.

- Using an out of band management system to power off a node, with confirmation, would be an example of ensuring process termination.
- Withdrawing a network route that is used by nodes at a failed site to communicate with storage would be another solution.



NOTE

Before restoring service to the failed zone or nodes, there must be confirmation that all pods with persistent volumes have terminated successfully.

To get the **Terminating** pods to recreate on the active zone, you can either force delete the pod or delete the finalizer on the associated PV. Once one of these two actions are completed, the application pod should recreate on the active zone and successfully mount its RWO storage.

Force delete the pod

Force deletions do not wait for confirmation from the kubelet that the Pod has been terminated.

```
$ oc delete pod <PODNAME> --grace-period=0 --force --namespace <NAMESPACE>
```

<PODNAME>

Is the name of the pod

<NAMESPACE>

Is the project namespace

Deleting the finalizer on the associated PV

Find the associated PV for the Persistent Volume Claim (PVC) that is mounted by the **Terminating** pod and delete the finalizer using the **oc patch** command.

```
$ oc patch -n openshift-storage pv/<PV_NAME> -p '{"metadata":{"finalizers":[]}}' --type=merge
```

<PV_NAME>

Is the name of the PV

An easy way to find the associated PV is to describe the Terminating pod. If you see a multi-attach warning, it should have the PV names in the warning (for example, pvc-0595a8d2-683f-443b-ae0-6e547f5f5a7c).

```
$ oc describe pod <PODNAME> --namespace <NAMESPACE>
```

<PODNAME>

Is the name of the pod

<NAMESPACE>

Is the project namespace

Example output:

```
[...]
Events:
  Type    Reason          Age    From          Message
  ----    -
  Normal  Scheduled       4m5s   default-scheduler   Successfully assigned openshift-storage/noobaa-db-pg-0 to perf1-mz8bt-worker-d2hdm
  Warning FailedAttachVolume 4m5s   attachdetach-controller Multi-Attach error for volume "pvc-0595a8d2-683f-443b-ae0-6e547f5f5a7c" Volume is already exclusively attached to one node and can't be attached to another
```

CHAPTER 6. RECOVERY FOR STATEFULSET PODS

Pods that are part of a stateful set have a similar issue as Pods mounting RWO volumes. Reference Kubernetes resource [StatefulSet considerations](#) for more information.

To get the Pods part of a StatefulSet to recreate on the active zone after 6-8 minutes, the Pod needs to be force deleted with the same requirements (i.e., OpenShift Container Platform node powered off or communication severed) as Pods with RWO volumes.