Configuring the CRUSH map and creating separate failure domains allows OSDs and cluster nodes to fail without any data loss occurring. The cluster simply operates in a degraded state until the problem is fixed.

Configuring the CRUSH map and creating separate performance domains can reduce performance bottlenecks for clients and applications that use the cluster to store and retrieve data. For example, CRUSH can create one hierarchy for HDDs and another hierarchy for SSDs.

A typical use case for customizing the CRUSH map is to provide additional protection against hardware failures. You can configure the CRUSH map to match the underlying physical infrastructure, which helps mitigate the impact of hardware failures.

By default, the CRUSH algorithm places replicated objects on OSDs on different hosts. You can customize the CRUSH map so that object replicas are placed across OSDs in different shelves, or on hosts in different rooms, or in different racks with distinct power sources.

Another use case is to allocate OSDs with SSD drives to pools used by applications requiring very fast storage, and OSDs with traditional HDDs to pools supporting less demanding workloads.

The CRUSH map can contain multiple hierarchies that you can select through different CRUSH rules. By using separate CRUSH hierarchies, you can establish separate performance domains. Use case examples for configuring separate performance domains are:

- To separate block storage used by VMs from object storage used by applications.

- To separate "cold" storage, containing infrequently accessed data, from "hot" storage, containing frequently accessed data.

If you examine an actual CRUSH map definition, it contains:

- A list of all available physical storage devices.

- A list of all the infrastructure buckets and the IDs of the storage devices or other buckets in each of them. Remember that a bucket is a container, or a branch, in the infrastructure tree. For example, it might represent a location or a piece of physical hardware.

- A list of CRUSH rules to map PGs to OSDs.

- A list of other CRUSH tunables and their settings.

The cluster installation process deploys a default CRUSH map. You can use the `ceph osd crush dump` command to print the CRUSH map in JSON format. You can also export a binary copy of the map and decompile it into a text file:

```
[ceph: root@node /]# ceph osd getcrushmap -o ./map.bin
[ceph: root@node /]# crushtool -d ./map.bin -o ./map.txt
```

## Customizing OSD CRUSH Settings

The CRUSH map contains a list of all the storage devices in the cluster. For each storage device the following information is available:

- The ID of the storage device.

- The name of the storage device.

- The weight of the storage device, normally based on its capacity in terabytes. For example, a 4 TB storage device has a weight of about 4.0. This is the relative amount of data the device can store, which the CRUSH algorithm uses to help ensure uniform object distribution.