

Análisis del impacto de la interrupción de Ceph para VNF de StarOS

Contenido

[Introducción](#)

[Requisito previo](#)

[Requirements](#)

[Componentes Utilizados](#)

[Abreviaturas](#)

[Ceph en Cisco VIM](#)

[Aspectos básicos del mecanismo de supervisión en la Ceph](#)

[Impacto del bloqueo de E/S en StarOS VNF](#)

[Escenarios de E/S de bloqueo prolongado](#)

[Mecanismo de temporizador de carga](#)

[Error de hardware de la tarjeta RAID](#)

[¿Cómo mitigar el impacto?](#)

[Mover al disco local desde el almacenamiento Ceph](#)

[Ajuste De La Configuración De Ceph](#)

[Controlar el problema del hardware de la tarjeta RAID](#)

[Ajuste CEPH OSD RESEREVED PCORES](#)

Introducción

Este documento describe cómo StarOS VNF, que se ejecuta en Cisco Virtualized Infrastructure Manager (VIM), se ve afectado cuando el servicio de almacenamiento Ceph está dañado, y qué se puede hacer para mitigar el impacto. Se explica en el supuesto de que Cisco VIM se utiliza como infraestructura, pero que la misma teoría se puede aplicar a cualquier entorno de Openstack.

Requisito previo

Requirements

Cisco recomienda que tenga conocimiento sobre estos temas:

- Cisco StarOS
- Cisco VIM
- Openstack
- Ceph

Componentes Utilizados

La información que contiene este documento se basa en las siguientes versiones de software y

hardware.

- StarOS: 21.16.c9
- VIM de Cisco: 3.2.2 (Queens de Openstack)

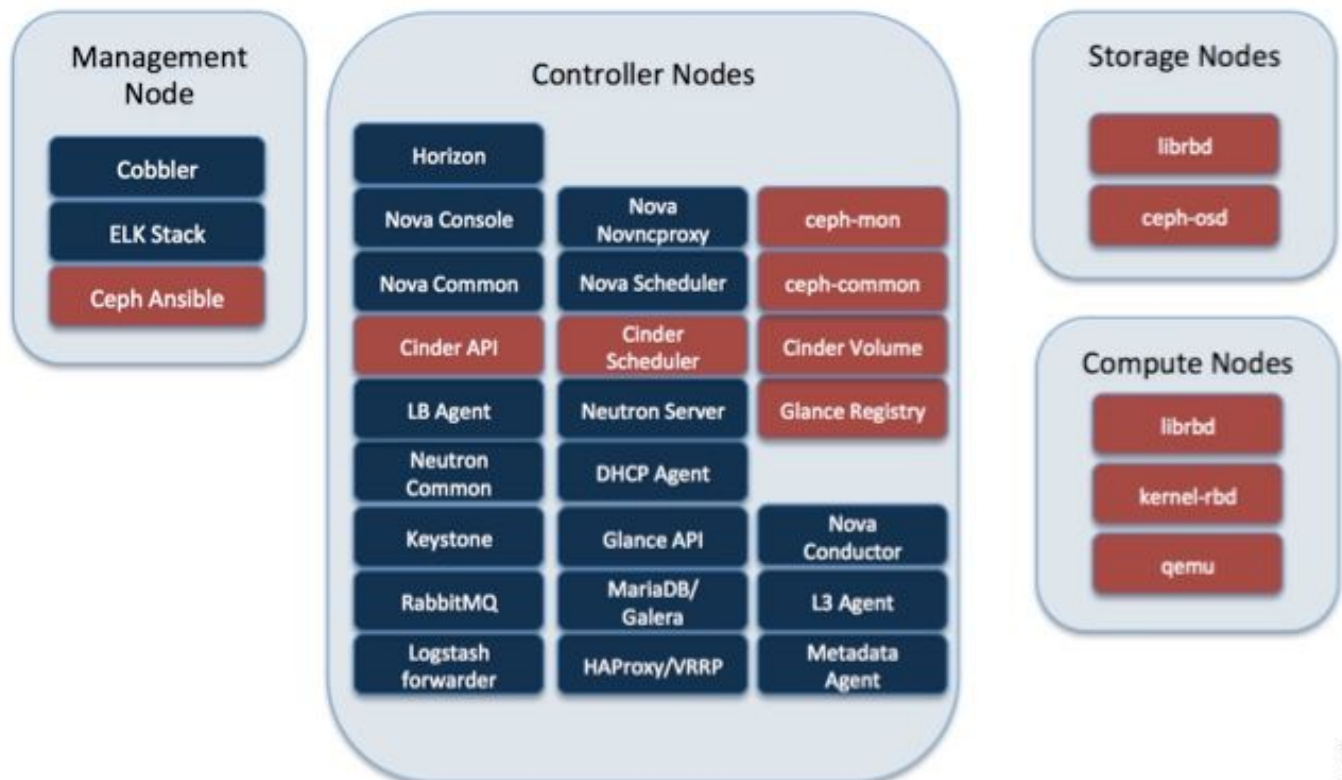
The information in this document was created from the devices in a specific lab environment. All of the devices used in this document started with a cleared (default) configuration. Si tiene una red en vivo, asegúrese de entender el posible impacto de cualquier comando.

Abreviaturas

Cisco VIM	Cisco Virtualized Infrastructure Manager
VNF	Función de red virtual
Ceph OSD	Ceph Object Storage Daemon
StarOS	Sistema operativo para la solución Cisco Mobile Packet Core

Ceph en Cisco VIM

Esta imagen se ha tomado de la Guía del administrador de Cisco VIM. Cisco VIM utiliza Ceph como back end de almacenamiento.



Ceph admite almacenamiento de objetos y bloques y, por lo tanto, se utiliza para almacenar imágenes y volúmenes de VM que se pueden conectar a VM. Entre los diversos servicios OpenStack que dependen del motor de almacenamiento se incluyen:

- Glance (servicio de imágenes OpenStack): utiliza Ceph para almacenar imágenes.
- Cinder (servicio de almacenamiento OpenStack): utiliza Ceph para crear volúmenes que se

pueden conectar a VM.

- Nova (servicio informático OpenStack): utiliza Ceph para conectarse a los volúmenes creados por Cinder.

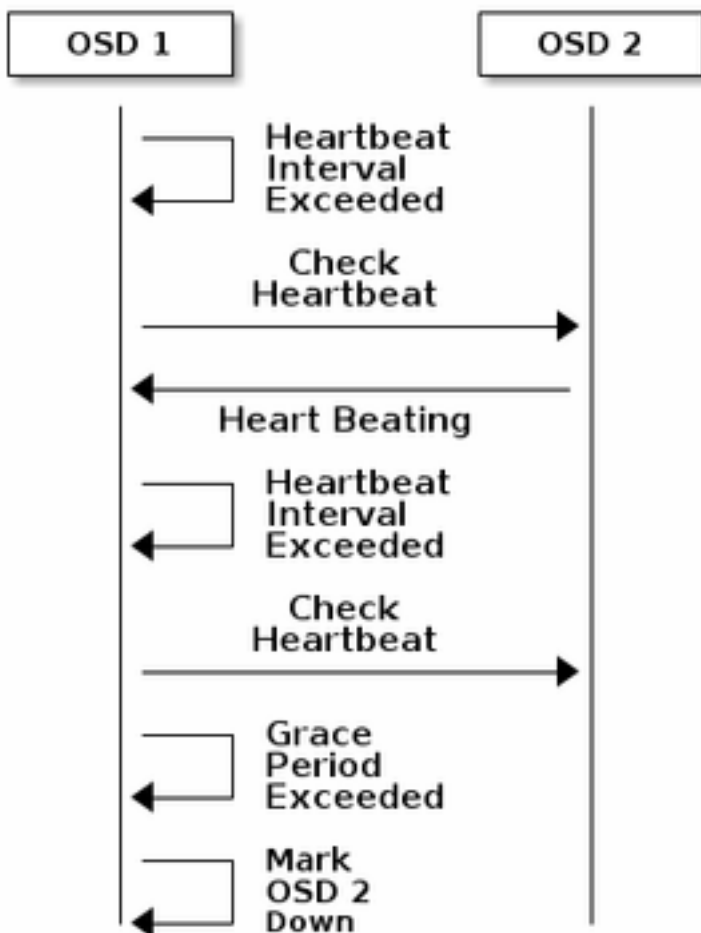
En muchos casos, se crea un volumen en Ceph para **/flash** y **/hd-raid** para StarOS VNF como el ejemplo aquí.

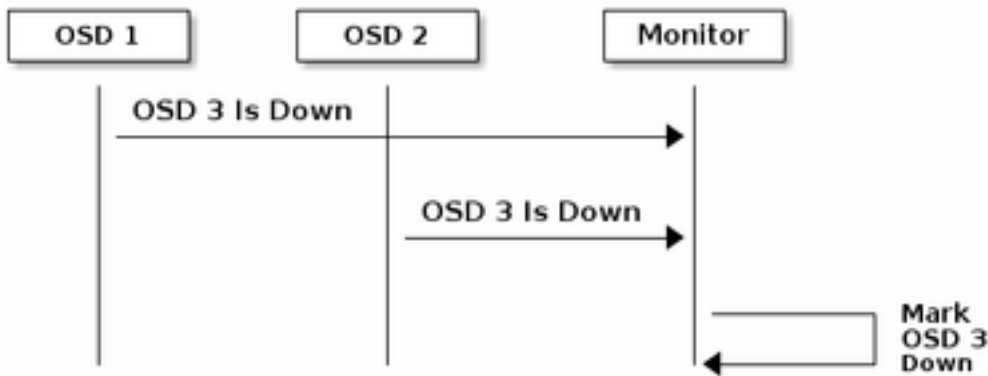
```
openstack volume create --image `glance image-list | grep up-image | awk '{print $2}` --size 16 --type LUKS up1-flash-boot  
openstack volume create --size 20 --type LUKS up1-hd-raid
```

Aspectos básicos del mecanismo de supervisión en la Ceph

Aquí está la explicación del documento de Ceph con respecto al monitoreo:

Cada demonio Ceph OSD verifica el latido del corazón de otros demonios Ceph OSD a intervalos aleatorios inferiores a cada 6 segundos. Si un Ceph OSD Daemon vecino no muestra latidos en un período de gracia de 20 segundos, el Ceph OSD Daemon puede considerar el Ceph OSD Daemon vecino e informarlo a un Ceph Monitor, que actualiza el Ceph Cluster Map. De forma predeterminada, dos demonios Ceph OSD de diferentes hosts deben informar a los Monitores Ceph de que otro demonio Ceph OSD está inactivo antes de que los Monitores Ceph reconozcan que el demonio Ceph OSD notificado está inactivo.





Por lo tanto, en general, toma unos 20 segundos detectar el OSD inactivo y se actualiza el mapa del clúster de Ceph, sólo después de que este VNF pueda utilizar un nuevo OSD. Durante este tiempo, el disco se bloquea.

Impacto del bloqueo de E/S en StarOS VNF

Si la E/S del disco se bloquea durante más de 120 segundos, se reinicia StarOS VNF. Hay una verificación específica para los procesos `xfssyncd/md0` y `xfs_db` que están relacionados con la E/S del disco y StarOS se reinicia intencionalmente cuando detecta un atascado en estos procesos más de 120 segundos.

Registro de consola de depuración de StarOS:

```

[ 1080.859817] INFO: task xfssyncd/md0:25787 blocked for more than 120 seconds.
[ 1080.862844] "echo 0 > /proc/sys/kernel/hung_task_timeout_secs" disables this message.
[ 1080.866184] xfssyncd/md0 D ffff880c036a8290 0 25787 2 0x00000000
[ 1080.869321] ffff880aacf87d30 0000000000000046 00000001000000a9a ffff880a00000000
[ 1080.872665] ffff880aacf87fd8 ffff880c036a8000 ffff880aacf87fd8 ffff880aacf87fd8
[ 1080.876100] ffff880c036a8298 ffff880aacf87fd8 ffff880c0f2f3980 ffff880c036a8000
[ 1080.879443] Call Trace:
[ 1080.880526] [<ffffffffff8123d62e>] ? xfs_trans_commit_iclog+0x28e/0x380
[ 1080.883288] [<ffffffffff810297c9>] ? default_spin_lock_flags+0x9/0x10
[ 1080.886050] [<ffffffffff8157fd7d>] ? _raw_spin_lock_irqsave+0x4d/0x60
[ 1080.888748] [<ffffffffff812301b3>] _xfs_log_force_lsn+0x173/0x2f0
[ 1080.891375] [<ffffffffff8104bae0>] ? default_wake_function+0x0/0x20
[ 1080.894010] [<ffffffffff8123dc15>] _xfs_trans_commit+0x2a5/0x2b0
[ 1080.896588] [<ffffffffff8121ff64>] xfs_fs_log_dummy+0x64/0x90
[ 1080.899079] [<ffffffffff81253cf1>] xfs_sync_worker+0x81/0x90
[ 1080.901446] [<ffffffffff81252871>] xfssyncd+0x141/0x1e0
[ 1080.903670] [<ffffffffff81252730>] ? xfssyncd+0x0/0x1e0
[ 1080.905871] [<ffffffffff81071d5c>] kthread+0x8c/0xa0
[ 1080.908815] [<ffffffffff81003364>] kernel_thread_helper+0x4/0x10
[ 1080.911343] [<ffffffffff81580805>] ? restore_args+0x0/0x30
[ 1080.913668] [<ffffffffff81071cd0>] ? kthread+0x0/0xa0
[ 1080.915808] [<ffffffffff81003360>] ? kernel_thread_helper+0x0/0x10
[ 1080.918411] **** xfssyncd/md0 stuck, resetting card
  
```

Pero no se limita al temporizador de 120 segundos, si la E/S del disco se bloquea durante un tiempo, incluso menos de 120 segundos, VNF puede reiniciarse por diversos motivos. El resultado aquí es un ejemplo que muestra un reinicio debido al problema de E/S del disco, a veces un desperfecto continuo de la tarea de StarOS, etc. Depende de la sincronización de E/S del disco activo frente al problema de almacenamiento.

```
[ 2153.370758] Hangcheck: hangcheck value past margin!  
[ 2153.396850] ata1.01: exception Emask 0x0 SAct 0x0 SErr 0x0 action 0x6 frozen  
[ 2153.396853] ata1.01: failed command: WRITE DMA EXT  
--- skip ---  
SYSLINUX 3.53 0x5d037742 EBIOS Copyright (C) 1994-2007 H. Peter Anvin
```

Básicamente, una E/S de bloqueo larga puede considerarse un problema crítico para StarOS VNF y debe minimizarse tanto como sea posible.

Escenarios de E/S de bloqueo prolongado

Según la investigación de varias implementaciones de clientes y pruebas de laboratorio, se han identificado dos escenarios principales que pueden causar un bloqueo prolongado de E/S en Ceph.

Mecanismo de temporizador de carga

Hay un mecanismo de latido entre los OSD, para detectar el OSD hacia abajo. Según el valor **osd_Hearbeat_Grace**(20 segundos de forma predeterminada), se detecta que el OSD ha fallado. Y hay un mecanismo de temporizador holgado, cuando hay una fluctuación o inestabilidad en el estado OSD el temporizador de gracia se ajusta automáticamente (se hace más largo). Esto puede hacer que el valor **osd_Hearbeat_Grace** sea mayor.

En la situación normal, la gracia del latido es de 20 segundos

```
2019-01-09 16:58:01.715155 mon.ceph-XXXXXX [INF] osd.2 failed (root=default,host=XXXXXX) (2  
reporters from different host after 20.000047 >= grace 20.000000)
```

Sin embargo, después de varias inestabilidad de red de un nodo de almacenamiento, se convierte en un valor mayor.

```
2019-01-10 16:44:15.140433 mon.ceph-XXXXXX [INF] osd.2 failed (root=default,host=XXXXXX) (2  
reporters from different host after 256.588099 >= grace 255.682576)
```

En el ejemplo anterior, se tardan 256 segundos en detectar el OSD como inactivo.

Error de hardware de la tarjeta RAID

Ceph puede no ser capaz de detectar la falla de hardware de la tarjeta RAID de manera oportuna. La falla de la tarjeta RAID termina con una especie de situación de bloqueo OSD. En este caso, el OSD inactivo se detecta después de unos minutos, lo que es suficiente para que el VNF de StarOS se reinicie.

Cuando se cuelga la tarjeta RAID, algunos núcleos de CPU toman el 100% del estado de la guerra.

```
%Cpu20 : 2.6 us, 7.9 sy, 0.0 ni, 0.0 id, 89.4 wa, 0.0 hi, 0.0 si, 0.0 st  
%Cpu21 : 0.0 us, 0.3 sy, 0.0 ni, 99.7 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0 st  
%Cpu22 : 31.3 us, 5.1 sy, 0.0 ni, 63.6 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0 st  
%Cpu23 : 0.0 us, 0.0 sy, 0.0 ni, 28.1 id, 71.9 wa, 0.0 hi, 0.0 si, 0.0 st
```

```
%Cpu24 : 0.0 us, 0.0 sy, 0.0 ni, 0.0 id,100.0 wa, 0.0 hi, 0.0 si, 0.0 st
%Cpu25 : 0.0 us, 0.0 sy, 0.0 ni, 0.0 id,100.0 wa, 0.0 hi, 0.0 si, 0.0 st
```

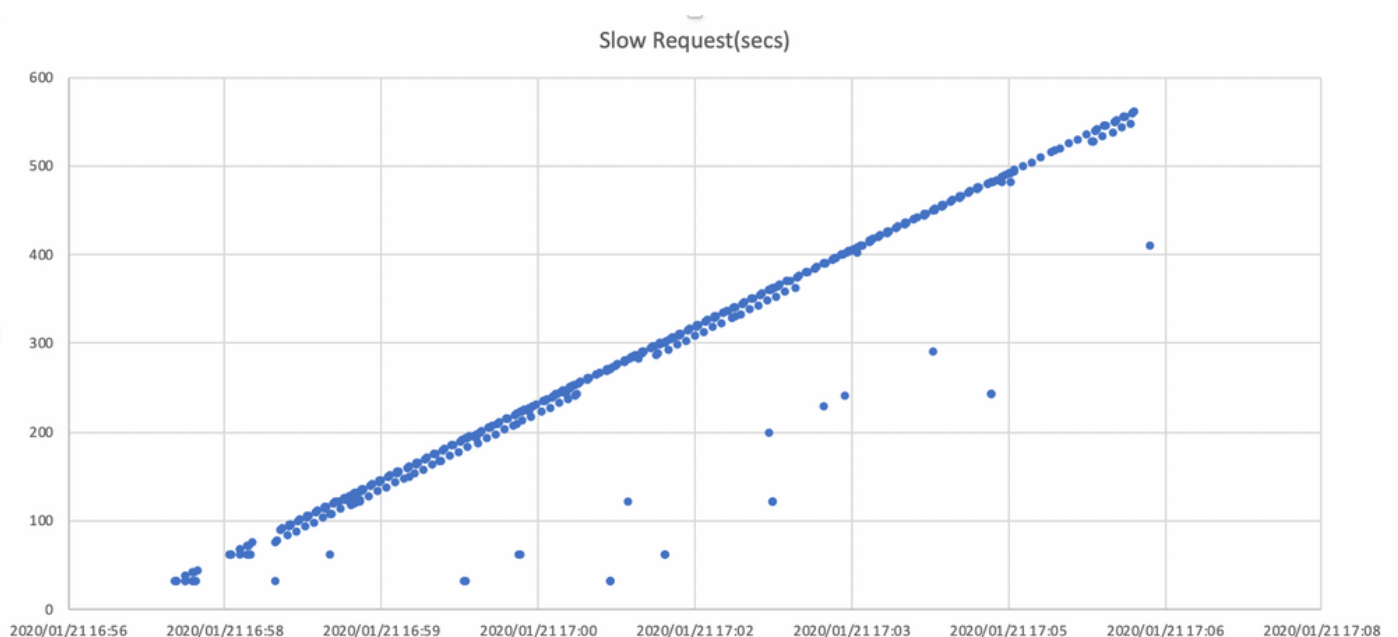
Además, consume todos los núcleos de la CPU gradualmente y OSD también está disminuyendo gradualmente con cierta diferencia de tiempo.

```
2019-01-01 17:08:05.267629 mon.ceph-XXXXX [INF] Marking osd.2 out (has been down for 602
seconds)
2019-01-01 17:09:25.296955 mon.ceph-XXXXX [INF] Marking osd.4 out (has been down for 603
seconds)
2019-01-01 17:11:10.351131 mon.ceph-XXXXX [INF] Marking osd.7 out (has been down for 604
seconds)
2019-01-01 17:16:40.426927 mon.ceph-XXXXX [INF] Marking osd.10 out (has been down for 603
seconds)
```

Paralelamente, las solicitudes lentas se detectan en **ceph.log**.

```
2019-01-01 16:57:26.743372 mon.XXXXX [WRN] Health check failed: 1 slow requests are blocked > 32
sec. Implicated osds 2 (REQUEST_SLOW)
2019-01-01 16:57:35.129229 mon.XXXXX [WRN] Health check update: 3 slow requests are blocked > 32
sec. Implicated osds 2,7,10 (REQUEST_SLOW)
2019-01-01 16:57:38.055976 osd.7 osd.7 [WRN] 1 slow requests, 1 included below; oldest blocked
for > 30.216236 secs
2019-01-01 16:57:39.048591 osd.2 osd.2 [WRN] 1 slow requests, 1 included below; oldest blocked
for > 30.635122 secs
-----skip-----
2019-01-01 17:06:22.124978 osd.7 osd.7 [WRN] 78 slow requests, 1 included below; oldest blocked
for > 554.285311 secs
2019-01-01 17:06:25.114453 osd.4 osd.4 [WRN] 19 slow requests, 1 included below; oldest blocked
for > 546.221508 secs
2019-01-01 17:06:26.125459 osd.7 osd.7 [WRN] 78 slow requests, 1 included below; oldest blocked
for > 558.285789 secs
2019-01-01 17:06:27.125582 osd.7 osd.7 [WRN] 78 slow requests, 1 included below; oldest blocked
for > 559.285915 secs
```

El gráfico de aquí muestra cuánto tiempo se bloquean las solicitudes de E/S con una línea de tiempo. El gráfico se crea trazando los registros de solicitudes lentos en **ceph.log**. Muestra que el tiempo de bloqueo se está alargando con el tiempo.



¿Cómo mitigar el impacto?

Mover al disco local desde el almacenamiento Ceph

La forma más sencilla de mitigar el impacto es pasar a un disco local desde el almacenamiento Ceph. StarOS utiliza 2 discos, /flash y /hd-raid, es posible mover solamente /flash al disco local, lo que hace que StarOS VNF sea más robusto para los problemas de Ceph. El lado negativo del uso del almacenamiento compartido como Ceph es que todos los VNF que lo utilizan se ven afectados al mismo tiempo que ocurre un problema. Al utilizar el disco local, el impacto del problema de almacenamiento se puede minimizar a VNF que se ejecuta sólo en el nodo afectado. Y los escenarios mencionados en la sección anterior son aplicables a Ceph solamente por lo que no son aplicables al disco local. Pero la otra cara del disco local es que el contenido del disco, como la imagen de StarOS, la configuración, el archivo principal, el registro de facturación, no se puede retener cuando se reimplementa la máquina virtual. También puede afectar al mecanismo de reparación automática de VNF.

Ajuste De La Configuración De Ceph

Desde el punto de vista de StarOS VNF, se recomiendan los siguientes parámetros Ceph nuevos para minimizar el tiempo de E/S de bloqueo mencionado anteriormente.

<default settings>

```
"mon_osd_adjust_heartbeat_grace": "true",  
"osd_client_watch_timeout": "30",  
"osd_max_markdown_count": "5",  
"osd_heartbeat_grace": "20",
```

<new settings>

```
"mon_osd_adjust_heartbeat_grace": "false",  
"osd_client_watch_timeout": "10",  
"osd_max_markdown_count": "1",  
"osd_heartbeat_grace": "10",
```

Consta de:

- El mecanismo del temporizador de carga está desactivado, no hay ajuste automático
- El tiempo de gracia del latido se acorta
- El OSD se marca inmediatamente como inactivo (de forma predeterminada 5 veces en los últimos 600 segundos)

Los nuevos parámetros se prueban en un laboratorio, el tiempo de detección para OSD inactivo se reduce a aproximadamente menos de 10 segundos, originalmente era alrededor de 30 segundos con la configuración predeterminada de Ceph.

Controlar el problema del hardware de la tarjeta RAID

Para el escenario de hardware de la tarjeta RAID, puede ser difícil detectar a tiempo la naturaleza del problema, ya que crea una situación en la que OSD funciona intermitentemente mientras se bloquea la E/S. No hay una única solución para esto, pero se recomienda monitorear el registro de hardware del servidor para ver si falla la tarjeta RAID, o el registro lento de la solicitud en ceph.log por alguna secuencia de comandos y tomar alguna acción como hacer que el OSD afectado se desactive de forma proactiva.

Ajuste CEPH_OSD_RESEREVED_PCORES

Esto no está relacionado con los escenarios mencionados, pero si hay un problema con el rendimiento de Ceph debido a la operación de E/S pesada, el aumento del valor CEPH_OSD_RESEREVED_PCORES puede mejorar el rendimiento de E/S de Ceph. De forma predeterminada, CEPH_OSD_RESEREVED_PCORES en Cisco VIM se configura como 2 y se puede aumentar.