

StarOS VNF的Ceph中斷影響分析

目錄

[簡介](#)

[必備條件](#)

[需求](#)

[採用元件](#)

[縮寫](#)

[Cisco VIM中的Ceph](#)

[Ceph中監控機制的基本知識](#)

[阻塞I/O對StarOS VNF的影響](#)

[長遮蔽I/O方案](#)

[滯後計時器機構](#)

[RAID卡硬體故障](#)

[如何減輕影響？](#)

[從Ceph儲存移動到本地磁碟](#)

[Ceph配置調整](#)

[監視RAID卡硬體問題](#)

[CEPH OSD RESERVED PCORES調整](#)

簡介

本文檔介紹當Ceph儲存服務受損時，在Cisco Virtualized Infrastructure Manager(VIM)上運行的StarOS VNF會受到什麼影響，以及可以採取什麼措施來減輕影響。此處的解釋基於以下假設：思科VIM用作基礎設施，但同樣的理論可以應用於任何Openstack環境。

必備條件

需求

思科建議您瞭解以下主題：

- Cisco StarOS
- Cisco VIM
- Openstack
- Ceph

採用元件

本文中的資訊係根據以下軟體和硬體版本：

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

本文中的資訊是根據特定實驗室環境內的裝置所建立。文中使用到的所有裝置皆從已清除（預設

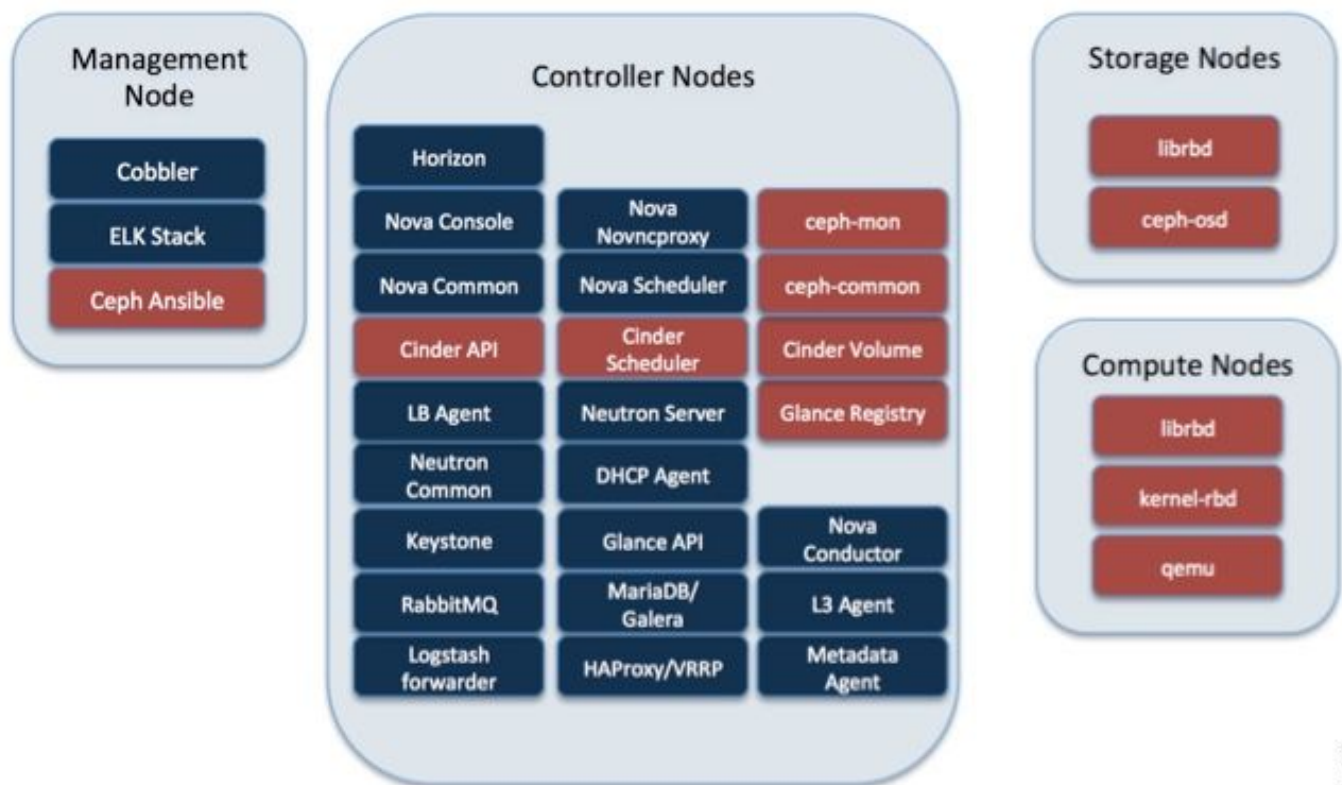
) 的組態來啟動。如果您的網路運作中，請確保您瞭解任何指令可能造成的影響。

縮寫

Cisco VIM	Cisco Virtualized Infrastructure Manager
VNF	虛擬網路功能
Ceph OSD	Ceph對象儲存守護程式
StarOS	適用於思科行動封包核心解決方案的作業系統

Cisco VIM中的Ceph

此圖摘自《Cisco VIM管理員指南》。Cisco VIM使用Ceph作為儲存後端。



Ceph支援塊儲存和對象儲存，因此用於儲存VM映像和可以連線到VM的卷。多個依賴於儲存後端的OpenStack服務包括：

- Glance (OpenStack映像服務) — 使用Ceph儲存映像。
- Cinder (OpenStack儲存服務) — 使用Ceph建立可以連線到VM的卷。
- Nova (OpenStack計算服務) — 使用Ceph連線到Cinder建立的卷。

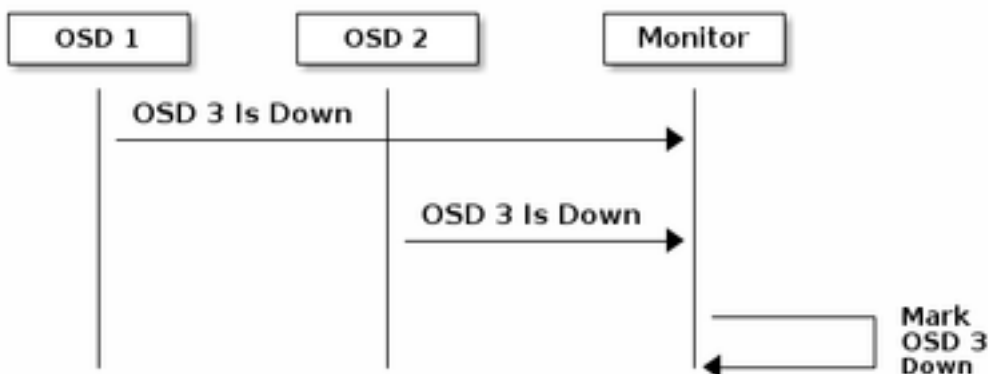
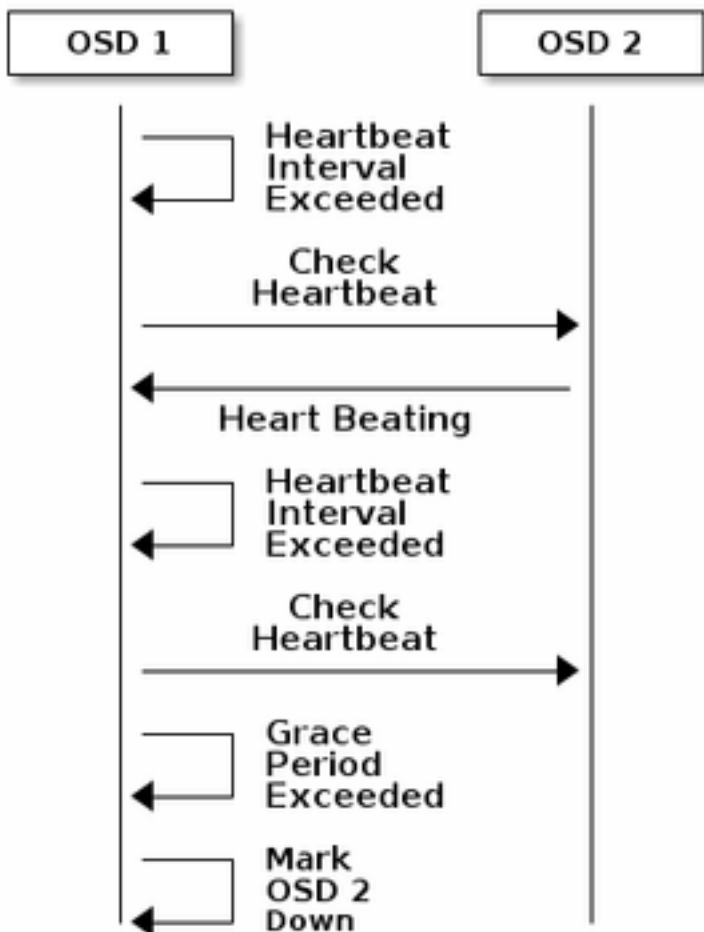
在許多情況下，會像本例一樣在Ceph中為flash和hd-raid為StarOS VNF建立一個卷。

```
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
```

Ceph中監控機制的基本知識

以下是Ceph文檔中有關監控的說明：

每個Ceph OSD守護進程以小於每6秒的隨機間隔檢查其他Ceph OSD守護進程的心跳。如果相鄰Ceph OSD守護進程在20秒的寬限期內未顯示心跳，則Ceph OSD守護進程可能會認為相鄰Ceph OSD守護進程已關閉，並將其報告給Ceph監視器，後者將更新Ceph群集對映。預設情況下，來自不同主機的兩個Ceph OSD守護程式必須向Ceph監控器報告另一個Ceph OSD守護程式已關閉，然後Ceph監控器確認報告的Ceph OSD守護程式已關閉。



因此，一般情況下，檢測OSD關閉大約需要20秒，並且Ceph群集對映會更新，只有在此VNF可以使用新的OSD之後。在此時間期間，會阻止I/O。

阻塞I/O對StarOS VNF的影響

如果磁碟I/O被阻止超過120秒，StarOS VNF將重新啟動。對於與磁碟I/O和StarOS相關的xfssyncd/md0和xfs_db進程，當檢測到這些進程上的停滯時間超過120秒時，會特意重新啟動這些進程。

StarOS調試控制檯日誌：

```
[ 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 0000000100000a9a 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
```

但它不限於120秒計時器，如果磁碟I/O被阻塞一段時間（甚至少於120秒），VNF可能會由於各種原因重新啟動。此處的輸出是一個示例，顯示由於磁碟I/O問題、有時連續的StarOS任務崩潰等原因而重新啟動的情況。這取決於活動磁碟I/O的時間與儲存問題。

```
[ 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

一般來說，長阻塞I/O可以被認為是StarOS VNF的關鍵問題，應儘可能將其降至最低。

長遮蔽I/O方案

根據對多個客戶部署和實驗室測試的研究，確定了2種在Ceph中會導致長阻塞I/O的主要方案。

滯後計時器機構

OSD之間有一個心跳機制，可檢測OSD關閉。根據osd_heartbeat_grace值（預設值為20秒），OSD被檢測為失敗。

此外，還有一個遲滯計時器機構，當OSD狀態有波動或擺動時，寬限計時器自動調整（變長）。這可能會使osd_heartbeat_grace值變大。

在正常情況下，心跳寬限為20秒

```
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)
```

但是在一個儲存節點有多個網路翻動之後，它就變得更有價值了。

```
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)
```

在上面的示例中，檢測OSD為down需要256秒。

RAID卡硬體故障

Ceph可能無法及時檢測RAID卡硬體故障。RAID卡故障最終導致一種OSD掛起情況。在這種情況下，幾分鐘後檢測到OSD關閉，這足以使StarOS VNF重新啟動。

當RAID卡掛起時，某些CPU核心在無線狀態下佔用100%。

```
%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
```

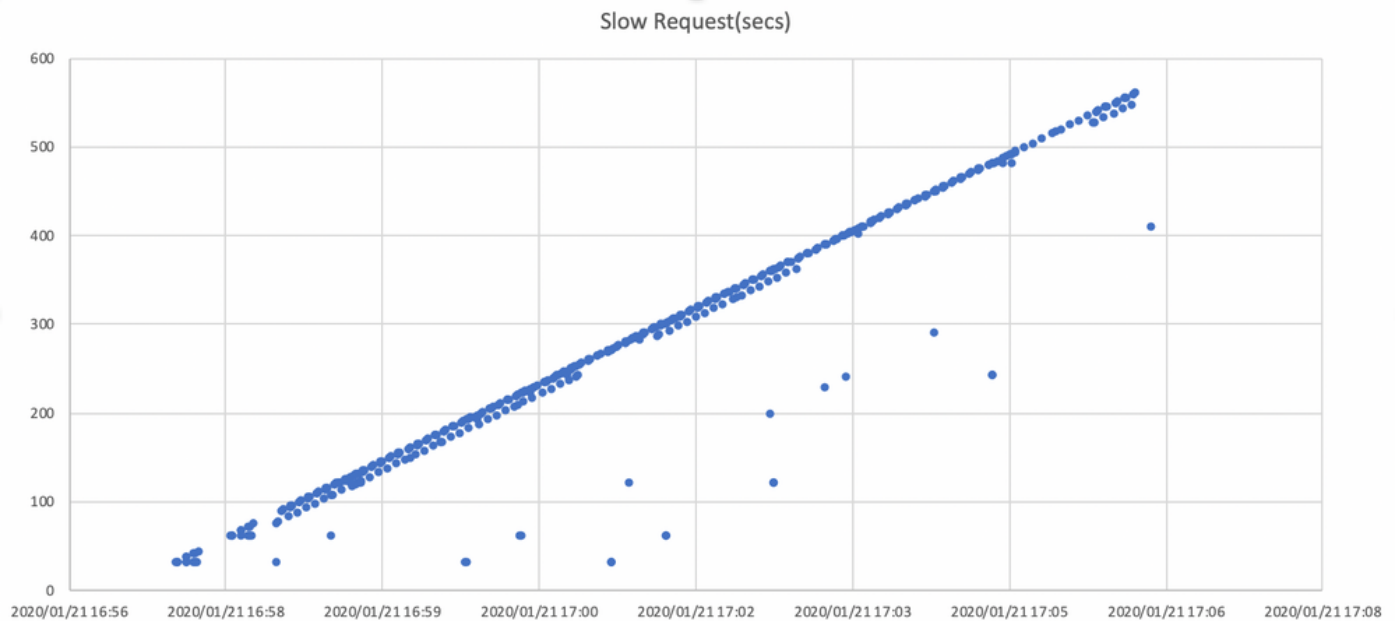
OSD也隨著一定的時間間隔逐漸下降。

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

同時，在ceph.log中檢測到慢速請求。

```
2019-01-01 16:57:26.743372 mon.XXXXXX [WRN] Health check failed: 1 slow requests are blocked > 32
sec. Implicated osds 2 (REQUEST_SLOW)
2019-01-01 16:57:35.129229 mon.XXXXXX [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
```

此處的圖形顯示了使用時間軸阻止I/O請求的時間。該圖形是通過在ceph.log中繪製慢速請求日誌建立的。它顯示阻塞時間隨著時間而變長。



如何減輕影響？

從Ceph儲存移動到本地磁碟

減輕影響的最簡單方法是從Ceph儲存移動到本地磁碟。StarOS使用2個磁碟，/flash和/hd-raid，可以只將/flash移動到本地磁碟，這使StarOS VNF對Ceph問題更加穩健。使用共用儲存（如Ceph）的負面影響是，在發生問題時，所有使用它的VNF都會同時受到影響。通過使用本地磁碟，可以將儲存問題的影響降至最低，使其僅在受影響的節點上運行VNF。並且上一節中提到的方案僅適用於Ceph，因此不適用於本地磁碟。但本地磁碟的另一面是，重新部署虛擬機器時無法保留磁碟的內容，如StarOS映像、配置、核心檔案、計費記錄。也會影響VNF自動修復機制。

Ceph配置調整

從StarOS VNF的角度來看，建議使用以下新的Ceph引數以最小化上述阻塞I/O時間。

<預設設定>

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

<新設定>

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

它包括：

- 滯後計時器機制被禁用，無自動調整
- 心跳寬限期縮短
- OSD立即標籤為關閉（預設情況下在最後600秒內關閉5次）

新引數在實驗室中測試，OSD關閉的檢測時間減少到大約10秒，最初使用Ceph的預設配置為30秒左右。

監視RAID卡硬體問題

對於RAID卡硬體情況，仍可能難以根據問題的性質及時進行檢測，因為它造成了OSD在I/O被阻止時間歇性工作的情況。沒有針對此問題的單一解決方案，但是建議通過某些指令碼監控伺服器硬體日誌以發現RAID卡故障，或監控ceph.log中慢速的請求日誌，並採取一些措施，例如主動關閉受影響的OSD。

CEPH_OSD_RESERVED_PCORES調整

這與上述場景無關，但是如果由於I/O操作繁重導致Ceph效能出現問題，則增加CEPH_OSD_RESERVED_PCORES值可以提高Ceph I/O效能。預設情況下，Cisco VIM上的CEPH_OSD_RESERVED_PCORES配置為2，並可增加。