



# ONTAP 與企業應用程式

## Enterprise applications

NetApp  
May 09, 2024

# 目錄

ONTAP 與企業應用程式	1
Hyper-V	2
部署準則和儲存最佳實務做法	2
Microsoft SQL Server	41
ONTAP 上的 Microsoft SQL Server	41
資料庫組態	41
儲存組態	48
使用 NetApp 管理軟體保護 Microsoft SQL Server 資料	59
使用 ONTAP 進行 Microsoft SQL Server 災難恢復	60
保護 ONTAP 上的 Microsoft SQL Server 安全	61
MySQL	63
ONTAP 上的 MySQL 資料庫	63
資料庫組態	63
主機組態	69
儲存組態	71
Oracle 資料庫	74
ONTAP 上的 Oracle 資料庫	74
組態 ONTAP	74
資料庫組態	83
主機組態	86
網路組態	99
儲存組態	105
Oracle 資料庫虛擬化	119
分層	122
Oracle 資料保護	128
Oracle 災難恢復	148
Oracle 資料庫移轉	170
其他附註	283
PostgreSQL	291
ONTAP 上的 PostgreSQL 資料庫	291
資料庫組態	291
儲存組態	294
資料保護	297
SAP	300
VMware	301
VMware vSphere 搭配 ONTAP	301
使用 ONTAP 的虛擬磁碟區 ( VVols )	335
VMware Site Recovery Manager 搭配 ONTAP	359
vSphere Metro Storage 叢集搭配 ONTAP	377

產品安全性 .....	405
適用於 VMware vSphere 的 ONTAP 工具安全性強化指南 .....	409
法律聲明 .....	423
版權 .....	423
商標 .....	423
專利 .....	423
隱私權政策 .....	423
開放原始碼 .....	423
ONTAP .....	423
適用於 MCC IP 的 ONTAP Mediator .....	424

# ONTAP 與企業應用程式

# Hyper-V

## 部署準則和儲存最佳實務做法

### 總覽

Microsoft Windows Server 是企業級作業系統（OS）、涵蓋網路、安全性、虛擬化、私有雲、混合雲、虛擬桌面基礎架構、存取保護、資訊保護、Web 服務、應用程式平台基礎架構、還有更多。



\* 本文件取代先前發佈的技術報告 \_TR-4568：NetApp 部署指南與 Windows Server\* 儲存最佳實務做法

**NetApp ONTAP (R)** 管理軟體可在 **NetApp** 儲存控制器上執行。有多種格式可供選擇。

- 支援檔案、物件和區塊傳輸協定的統一化架構。如此一來、儲存控制器就能同時做為 NAS 和 SAN 裝置、以及物件儲存區
- All SAN Array (ASA) 僅著重於區塊傳輸協定、並透過為連線主機新增對稱的雙主動式多重路徑功能、將 I/O 恢復時間 (IORT) 最佳化
- 軟體定義的統一化架構
  - 在 VMware vSphere 或 KVM 上執行的 ONTAP Select
  - Cloud Volumes ONTAP 以雲端原生執行個體執行
- 來自超大規模雲端供應商的第一方產品
  - Amazon FSX for NetApp ONTAP 產品
  - Azure NetApp Files
  - Google Cloud NetApp Volumes

ONTAP 提供 NetApp 儲存效率功能、例如 NetApp Snapshot (R) 技術、複製、重複資料刪除、精簡配置、精簡複寫、壓縮、虛擬儲存分層等多項功能、都能提升效能與效率。

Windows Server 和 ONTAP 可在大型環境中共同運作、為資料中心整合和私有雲或混合雲部署帶來巨大價值。這種組合也能有效率地提供不中斷營運的工作負載、並支援無縫擴充。

### 目標對象

本文件適用於設計適用於 Windows Server 的 NetApp 儲存解決方案的系統與儲存架構設計師。

我們在此文件中做出下列假設：

- 讀者對 NetApp 硬體與軟體解決方案有廣泛的瞭解。請參閱 ["叢集管理員系統管理指南"](#) 以取得詳細資料。
- 讀取器一般瞭解區塊存取傳輸協定、例如 iSCSI、FC 和檔案存取傳輸協定 SMB/CIFS。請參閱 ["叢集式 Data ONTAP SAN 管理"](#) 以取得 SAN 相關資訊。請參閱 ["NAS 管理"](#) 以取得 CIFS/SMB 相關資訊。
- 讀者對 Windows Server 作業系統和 Hyper-V 有廣泛的瞭解

如需完整且定期更新的已測試及支援 SAN 和 NAS 組態對照表、請參閱 ["互通性對照表工具IMT \(不含\)"](#) 在

NetApp 支援網站上。有了 IMT、您就能判斷特定環境所支援的確切產品和功能版本。NetApp IMT 定義與 NetApp 支援組態相容的產品元件和版本。具體結果取決於每位客戶依照已發佈規格所安裝的產品。

## NetApp 儲存設備和 Windows Server 環境

如中所述 "總覽" NetApp 儲存控制器提供真正統一化的架構、支援檔案、區塊和物件傳輸協定。這包括 SMB/CIFS、NFS、NVMe/TCP、NVMe/FC、iSCSI、FC (FCP) 和 S3、可建立統一化的用戶端和主機存取。同一個儲存控制器可以同時以 SAN LUN 的形式提供區塊儲存服務、並以 NFS 和 SMB/CIFS 的形式提供檔案服務。ONTAP 也可做為 All SAN Array (ASA) 使用、透過 iSCSI 和 FCP 的對稱雙主動式多重路徑功能來最佳化主機存取、而統一化的 ONTAP 系統則使用非對稱雙主動式多重路徑功能。在這兩種模式中、ONTAP 都會使用 ANA 來管理 NVMe over Fabrics (NVMe of) 多重路徑。

執行 ONTAP 軟體的 NetApp 儲存控制器可在 Windows Server 環境中支援下列工作負載：

- 以持續可用的 SMB 3.0 共享區代管的 VM
- 在 iSCSI 或 FC 上執行的叢集共用磁碟區 (CSV) LUN 上託管的 VM
- SMB 3.0 共用上的 SQL Server 資料庫
- NVMe、iSCSI 或 FC 上的 SQL Server 資料庫
- 其他應用程式工作負載

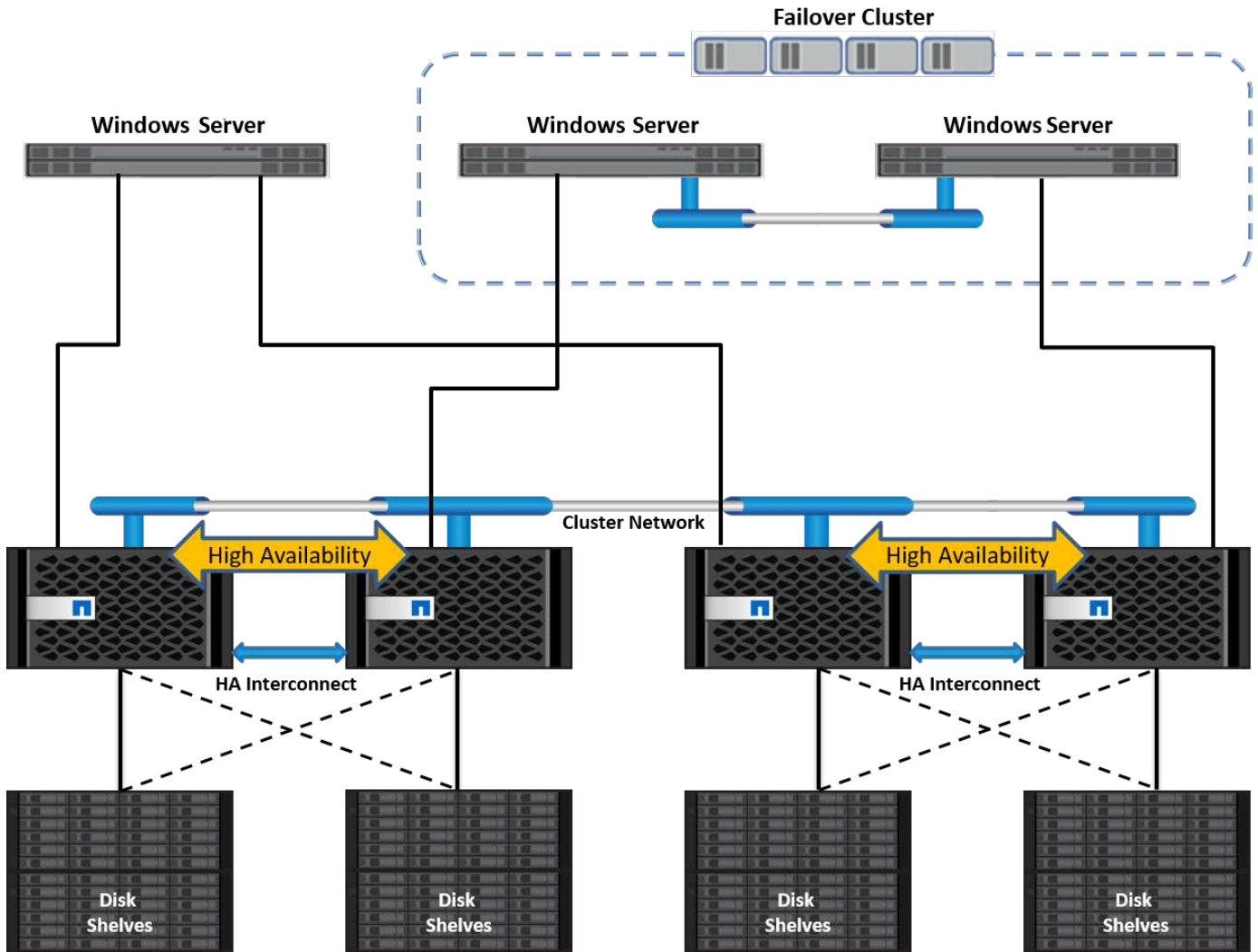
此外、NetApp 儲存效率功能、例如重複資料刪除、NetApp FlexClone (R) 複本、NetApp Snapshot 技術、精簡配置、壓縮、而儲存分層功能可為在 Windows Server 上執行的工作負載提供顯著的價值。

### ONTAP 資料管理

ONTAP 是在 NetApp 儲存控制器上執行的管理軟體。NetApp 儲存控制器稱為節點、是一種硬體裝置、內含處理器、RAM 和 NVRAM。節點可以連接至 SATA、SAS 或 SSD 磁碟機、或是這些磁碟機的組合。

將多個節點彙總至叢集式系統。叢集中的節點會持續彼此通訊、以協調叢集活動。節點也可以使用備援路徑、將資料透明地從節點移至節點、然後移至由兩個 10Gb 乙太網路交換器組成的專用叢集網路。叢集中的節點可以互相接管、在任何容錯移轉案例中提供高可用度。叢集是在整個叢集上進行管理、而非以每個節點為基礎、而且資料是從一或多個儲存虛擬機器 (SVM) 提供。叢集必須至少有一個 SVM 才能提供資料。

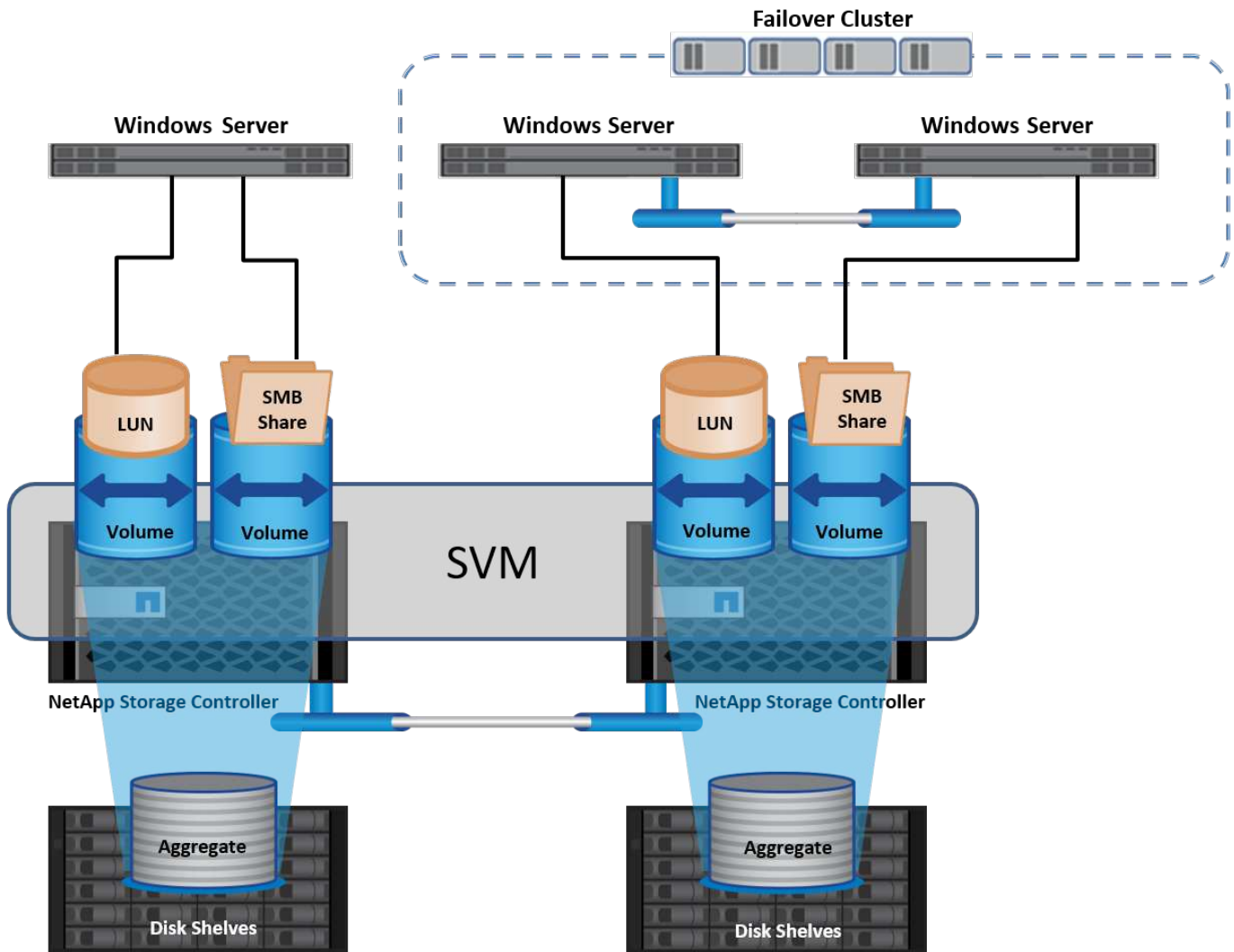
叢集的基本單元是節點、節點會新增至叢集、作為高可用度 (HA) 配對的一部分。HA 配對可透過 HA 互連 (與專用叢集網路分開) 彼此通訊、並維持與 HA 配對磁碟的備援連線、實現高可用度。雖然磁碟櫃可能包含屬於 HA 配對任一成員的磁碟、但磁碟不會在 HA 配對之間共用。下圖說明 Windows Server 環境中的 NetApp 儲存部署。



### 儲存虛擬機器

ONTAP SVM 是一種邏輯儲存伺服器、可從一或多個邏輯介面（生命體）提供 LUN 和 / 或 NAS 命名空間的資料存取。因此 SVM 是儲存區段的基本單元、可在 ONTAP 中實現安全的多租戶共享。每個 SVM 都設定為擁有儲存磁碟區、這些磁碟區是從指派給實體乙太網路或 FC 目標連接埠的實體集合和邏輯介面（生命體）來配置。

邏輯磁碟（LUN）或 CIFS 共用是在 SVM 的磁碟區內建立、並對應至 Windows 主機和叢集、以提供儲存空間、如下圖所示。SVM 不受節點限制、並以叢集為基礎；它們可以使用實體資源、例如叢集中任何位置的磁碟區或網路連接埠。



## 為 Windows Server 佈建 NetApp 儲存設備

儲存設備可同時在 SAN 和 NAS 環境中佈建至 Windows Server。在 SAN 環境中、儲存設備是以 NetApp 磁碟區上 LUN 的磁碟形式提供、作為區塊儲存設備。在 NAS 環境中、儲存設備會在 NetApp 磁碟區上以 CIFS/SMB 共用的形式提供、做為檔案儲存設備。這些磁碟和共用可以在 Windows Server 中套用、如下所示：

- 適用於 Windows Server 主機的儲存設備、適用於應用程式工作負載
- 儲存儲存設備：奈米伺服器和容器
- 用於儲存 VM 的個別 Hyper-V 主機儲存設備
- Hyper-V 叢集的共用儲存設備、以 CSV 的形式儲存 VM
- SQL Server 資料庫的儲存設備

## 管理 NetApp 儲存設備

若要從 Windows Server 2016 連線、設定及管理 NetApp 儲存設備、請使用下列其中一種方法：

- \* 安全 Shell (SSH) 。 \* 使用 Windows Server 上的任何 SSH 用戶端來執行 NetApp CLI 命令。
- \* System Manager 。 \* 這是 NetApp 的 GUI 型管理產品。



- \* NetApp PowerShell Toolkit\* 。這是 NetApp PowerShell 工具套件、用於自動化及實作自訂指令碼和工作流程。

## NetApp PowerShell 工具套件

NetApp PowerShell Toolkit ( PSTK ) 是一個 PowerShell 模組、可提供端點對端自動化、並可讓您管理 NetApp ONTAP 的儲存設備。ONTAP 模組包含超過 2 、 000 個 Cmdlet 、可協助管理 FAS 、 NetApp All Flash FAS ( AFF ) 、市售硬體和雲端資源。

### 值得記住的事項

- NetApp 不支援 Windows Server 儲存空間。儲存空間僅用於 JBOD ( 只有一堆磁碟 ) 、不適用於任何類型的 RAID ( 直接附加儲存 [DAS] 或 SAN ) 。
- ONTAP 不支援 Windows Server 中的叢集式儲存集區。
- NetApp 支援共享虛擬硬碟格式 ( VHDX ) 、用於 Windows SAN 環境中的來賓叢集。
- Windows Server 不支援使用 iSCSI 或 FC LUN 建立儲存池。

### 進一步閱讀

- 如需 NetApp PowerShell 工具組的詳細資訊、請參閱 "[NetApp 支援網站](#)" 。
- 如需 NetApp PowerShell 工具組最佳實務做法的相關資訊、請參閱 "[TR-4475 : NetApp PowerShell 工具套件最佳實務指南](#)" 。

## 網路最佳實務做法

乙太網路可廣泛分為下列群組：

- 虛擬機器的用戶端網路
- 多個儲存網路 ( 連接至儲存系統的 iSCSI 或 SMB )
- 叢集通訊網路 ( 叢集節點之間的活動訊號和其他通訊 )
- 管理網路 ( 用於監控系統並進行疑難排解 )
- 移轉網路 ( 用於主機即時移轉 )
- VM 複寫 ( Hyper-V 複本 )

### 最佳實務做法

- NetApp 建議您針對上述各項功能、使用專用的實體連接埠來隔離網路並提高效率。
- 針對上述每項網路需求 ( 儲存需求除外 ) 、可彙總多個實體網路連接埠以分散負載或提供容錯能力。
- NetApp 建議在 Hyper-V 主機上建立專用的虛擬交換器、以便在 VM 內建立來賓儲存連線。
- 請確定 Hyper-V 主機和來賓 iSCSI 資料路徑使用不同的實體連接埠和虛擬交換器、以確保來賓與主機之間的安全隔離。
- NetApp 建議避免 iSCSI NIC 的 NIC 群組。
- NetApp 建議使用主機上設定的 ONTAP 多重路徑輸入 / 輸出 ( MPIO ) 來進行儲存。
- 如果使用來賓 iSCSI 啟動器、NetApp 建議在來賓 VM 中使用 MPIO 。如果您使用直接移轉磁碟、則必須避免在客體內使用 MPIO 。在這種情況下、在主機上安裝 MPIO 就足夠了。

- NetApp 建議不要將 QoS 原則套用至指派給儲存網路的虛擬交換器。
- NetApp 建議不要在實體 NIC 上使用自動私有 IP 位址（APIPA）、因為 APIPA 不可路由且未在 DNS 中登錄。
- NetApp 建議為 CSV、iSCSI 和即時移轉網路開啟巨型框架、以提高處理量並縮短 CPU 週期。
- NetApp 建議取消勾選允許管理作業系統共用 Hyper-V 虛擬交換器的此網路介面卡選項、以建立虛擬機器專用的網路。
- NetApp 建議建立備援網路路徑（多個交換器）、以進行即時移轉、並建立 iSCSI 網路、以提供恢復能力和 QoS。

## 在 SAN 環境中進行資源配置

ONTAP SVM 支援區塊傳輸協定 iSCSI 和 FC。使用區塊傳輸協定 iSCSI 或 FC 建立 SVM 時、SVM 會分別取得 iSCSI 合格名稱（IQN）或 FC 全球名稱（WWN）。此識別碼為存取 NetApp 區塊儲存設備的主機提供 SCSI 目標。

## 在 Windows Server 上佈建 NetApp LUN

先決條件

在 Windows Server 的 SAN 環境中使用 NetApp 儲存設備有下列需求：

- NetApp 叢集已設定一或多個 NetApp 儲存控制器。
- NetApp 叢集或儲存控制器具有有效的 iSCSI 授權。
- 可使用 iSCSI 和 / 或 FC 組態連接埠。
- FC 分區是在 FC 的 FC 交換器上執行。
- 至少會建立一個 Aggregate。
- SVM 應在每個儲存控制器上、每個乙太網路或光纖通道架構都有一個 LIF、以便使用 iSCSI 或光纖通道來提供資料。

部署

1. 建立新的 SVM、並啟用區塊傳輸協定 iSCSI 和 / 或 FC。您可以使用下列任一方法建立新的 SVM：
  - NetApp 儲存設備上的 CLI 命令
  - 系統管理程式 ONTAP
  - NetApp PowerShell 工具套件
2. 設定 iSCSI 和 / 或 FC 傳輸協定。
3. 在每個叢集節點上指派具有生命的 SVM。
4. 在 SVM 上啟動 iSCSI 和 / 或 FC 服務。
  -
5. 使用 SVM 生命來建立 iSCSI 和 / 或 FC 連接埠集。
6. 使用建立的連接埠集、為 Windows 建立 iSCSI 和 / 或 FC 啟動器群組。

7. 將啟動器新增至啟動器群組。啟動器是 iSCSI 的 IQN 、 FC 的 WWPN 。您可以執行 PowerShell Cmdlet Get-InitiatorPort 從 Windows Server 查詢。

```
# Get the IQN for iSCSI
Get-InitiatorPort | Where \{$_.ConnectionType -eq 'iSCSI'} | Select-Object -Property NodeAddress
```

```
# Get the WWPN for FC
Get-InitiatorPort | Where \{$_.ConnectionType -eq 'Fibre Channel'} | Select-Object -Property PortAddress
```

```
# While adding initiator to the initiator group in case of FC, make sure to provide the initiator(PortAddress) in the standard WWPN format
```

Windows Server 上 iSCSI 的 IQN 也可以在 iSCSI 啟動器內容的組態中進行檢查。

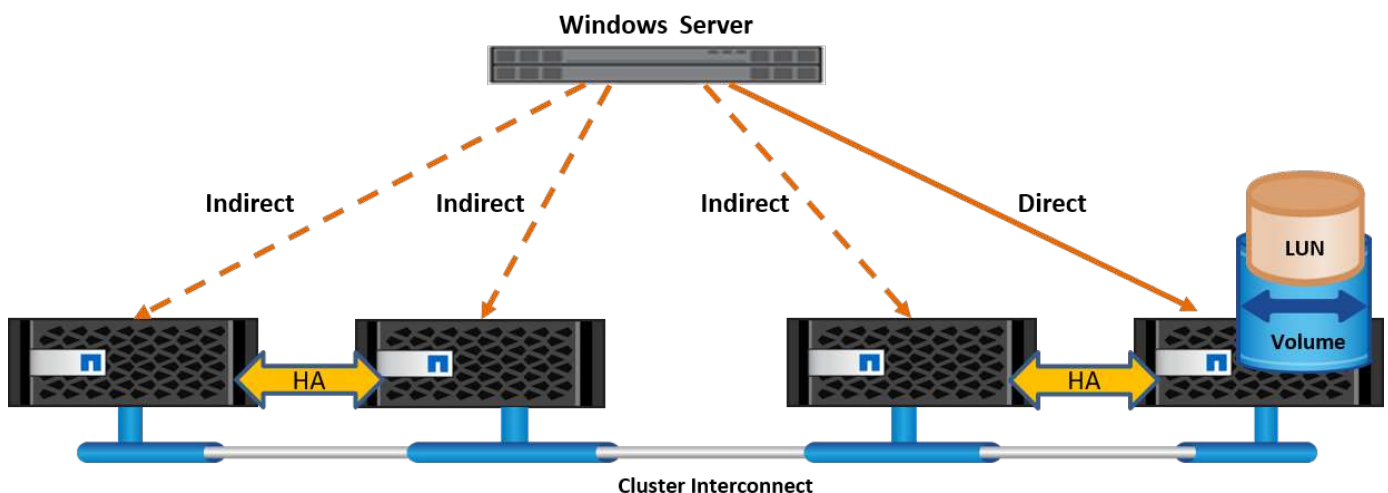
- 使用 Create LUN （創建 LUN ）嚮導創建 LUN ，並將其與創建的啟動器組相關聯。

#### 主機整合

Windows Server 使用非對稱式邏輯單元存取（ALUA）擴充 MPIO 來判斷通往 LUN 的直接和間接路徑。雖然 SVM 擁有的每個 LIF 都接受其 LUN 的讀取 / 寫入要求、但在任何指定時刻、只有其中一個叢集節點實際擁有支援該 LUN 的磁碟。這會將通往 LUN 的可用路徑分為兩種類型、直接或間接、如下圖所示。

LUN 的直接路徑是指 SVM 的生命和要存取的 LUN 位於同一個節點上的路徑。若要從實體目標連接埠移至磁碟、就不需要周遊叢集網路。

間接路徑是指 SVM 的生命與 LUN 位於不同節點上的資料路徑。資料必須穿越叢集網路、才能從實體目標連接埠移至磁碟。



## MPIO

NetApp ONTAP 提供高可用度的儲存設備、其中可存在從儲存控制器到 Windows Server 的多個路徑。多重路徑是指從伺服器到儲存陣列的多重資料路徑。多重路徑可防止硬體故障（纜線切割、交換器和主機匯流排介面卡 [HBA] 故障等）、並可利用多重連線的彙總效能來提供更高的效能限制。當某個路徑或連線無法使用時、多重路徑軟體會自動將負載移至其他可用路徑之一。MPIO 功能將通往儲存設備的多個實體路徑結合為單一邏輯路徑、用於資料存取、以提供儲存恢復能力和負載平衡。若要使用此功能、必須在 Windows Server 上啟用 MPIO 功能。

### 啟用 MPIO

若要在 Windows Server 上啟用 MPIO、請完成下列步驟：

1. 以系統管理員群組成員的身分登入 Windows Server。
2. 啟動 Server Manager。
3. 按一下 [ 管理 ] 區段中的 [ 新增角色和功能 ]。
4. 在「選取功能」頁面中、選取「多重路徑 I/O」。

### 設定 MPIO

使用 iSCSI 傳輸協定時、您必須告知 Windows Server 將多重路徑支援套用至 MPIO 內容中的 iSCSI 裝置。

若要在 Windows Server 上設定 MPIO、請完成下列步驟：

1. 以系統管理員群組成員的身分登入 Windows Server。
2. 啟動 Server Manager。
3. 按一下 [ 工具 ] 區段中的 [ MPIO ]。
4. 在 Discover Multi-Paths 的 MPIO Properties（MPIO 內容）中、選取 Add Support for iSCSI Devices（新增 iSCSI 裝置支援）、然後按一下 Add（新增）。接著會出現提示、要求您重新啟動電腦。
5. 重新啟動 Windows Server、查看 MPIO 內容的「MPIO 裝置」一節中所列的 MPIO 裝置。

### 設定 iSCSI

若要在 Windows Server 上偵測 iSCSI 區塊儲存、請完成下列步驟：

1. 以系統管理員群組成員的身分登入 Windows Server。
2. 啟動 Server Manager。
3. 按一下 [ 工具 ] 區段中的 [ iSCSI 啟動器 ]。
4. 按一下「探索」索引標籤下的「探索入口網站」。
5. 提供與為 SAN 傳輸協定的 NetApp 儲存設備所建立之 SVM 相關聯的生命負載 IP 位址。按一下「進階」、在「一般」索引標籤中設定資訊、然後按一下「確定」。
6. iSCSI 啟動器會自動偵測 iSCSI 目標、並將其列在「目標」索引標籤中。
7. 在探索到的目標中選取 iSCSI 目標。按一下「連線」以開啟「連線至目標」視窗。
8. 您必須在 NetApp 儲存叢集上、從 Windows Server 主機建立多個工作階段至目標 iSCSI 生命期。若要這麼做、請完成下列步驟：

9. 在「連線至目標」視窗中、選取「啟用 MPIO」、然後按一下「進階」。
10. 在「一般」索引標籤下的「進階設定」中、選取本機介面卡做為 Microsoft iSCSI 啟動器、然後選取「啟動器 IP」和「目標入口網站 IP」。
11. 您也必須使用第二個路徑進行連線。因此、請重複步驟 5 至步驟 8、但這次請為第二個路徑選取啟動器 IP 和目標入口網站 IP。
12. 在 iSCSI Properties (iSCSI 屬性) 主窗口的 Discerred Targets (已發現目標) 中選擇 iSCSI 目標、然後單擊 Properties (屬性)。
13. 「內容」視窗顯示已偵測到多個工作階段。選取工作階段、按一下「裝置」、然後按一下 MPIO 以設定負載平衡原則。會顯示為裝置設定的所有路徑、並支援所有負載平衡原則。NetApp 通常建議使用子集循環資源、而此設定是啟用 ALUA 的陣列的預設值。循環配置資源是不支援 ALUA 的主動式陣列的預設值。

#### 偵測區塊儲存

若要在 Windows Server 上偵測 iSCSI 或 FC 區塊儲存、請完成下列步驟：

1. 按一下「伺服器管理員」「工具」區段中的「電腦管理」。
2. 在 [ 電腦管理 ] 中、按一下 [ 儲存設備中的磁碟管理 ] 區段、然後按一下 [ 其他動作及重新掃描磁碟 ]。這樣做會顯示原始 iSCSI LUN。
3. 按一下探索到的 LUN、然後將其上線。然後選取使用 MBR 或 GPT 分割區初始化磁碟。提供磁碟區大小和磁碟機代號、並使用 FAT、FAT32、NTFS 或彈性檔案系統 (Refs) 格式化、以建立新的簡易磁碟區。

#### 最佳實務做法

- NetApp 建議在託管 LUN 的磁碟區上啟用精簡配置。
- 為了避免多重路徑問題、NetApp 建議使用所有 10Gb 工作階段或所有 1Gb 工作階段、連至指定的 LUN。
- NetApp 建議您確認已在儲存系統上啟用 ALUA。ONTAP 預設會啟用 ALUA。
- 在 NetApp LUN 對應至的 Windows Server 主機上、在防火牆設定中、針對輸入和 iSCSI 服務 (TCP 輸出) 啟用 iSCSI 服務 (TCP 輸入)。這些設定可讓 iSCSI 流量進出 Hyper-V 主機和 NetApp 控制器。

#### 在奈米伺服器上佈建 NetApp LUN

##### 先決條件

除了上一節提及的先決條件、儲存角色必須從奈米伺服器端啟用。例如、必須使用 -Storage 選項來部署奈米伺服器。若要部署奈米伺服器、請參閱「[部署奈米伺服器](#)。」

##### 部署

若要在奈米伺服器上配置 NetApp LUN、請完成下列步驟：

1. 請依照「[連線至奈米伺服器](#)。」
2. 若要設定 iSCSI、請在奈米伺服器上執行下列 PowerShell Cmdlet：

```
# Start iSCSI service, if it is not already running
Start-Service msiscsi
```

```
# Create a new iSCSI target portal
New-IscsiTargetPortal -TargetPortalAddress <SVM LIF>
```

```
# View the available iSCSI targets and their node address
Get-IscsiTarget
```

```
# Connect to iSCSI target
Connect-IscsiTarget -NodeAddress <NodeAddress>
```

```
# NodeAddress is retrieved in above cmdlet Get-IscsiTarget
# OR
Get-IscsiTarget | Connect-IscsiTarget
```

```
# View the established iSCSI session
Get-IscsiSession
```

```
# Note the InitiatorNodeAddress retrieved in the above cmdlet Get-
IscsiSession. This is the IQN for Nano server and this needs to be added
in the Initiator group on NetApp Storage
```

```
# Rescan the disks
Update-HostStorageCache
```

### 3. 將啟動器新增至啟動器群組。

```
Add the InitiatorNodeAddress retrieved from the cmdlet Get-IscsiSession
to the Initiator Group on NetApp Controller
```

### 4. 設定 MPIO。

```
# Enable MPIO Feature
Enable-WindowsOptionalFeature -Online -FeatureName MultipathIo
```

```
# Get the Network adapters and their IPs
Get-NetIPAddress -AddressFamily IPv4 -PrefixOrigin <Dhcp or Manual>
```

```
# Create one MPIO-enabled iSCSI connection per network adapter
Connect-IscsiTarget -NodeAddress <NodeAddress> -IsPersistent $True -IsMultipathEnabled $True -InitiatorPortalAddress <IP Address of ethernet adapter>
```

```
# NodeAddress is retrieved from the cmdlet Get-IscsiTarget
# IPs are retrieved in above cmdlet Get-NetIPAddress
```

```
# View the connections
Get-IscsiConnection
```

## 5. 偵測區塊儲存。

```
# Rescan disks
Update-HostStorageCache
```

```
# Get details of disks
Get-Disk
```

```
# Initialize disk
Initialize-Disk -Number <DiskNumber> -PartitionStyle <GPT or MBR>
```

```
# DiskNumber is retrived in the above cmdlet Get-Disk
# Bring the disk online
Set-Disk -Number <DiskNumber> -IsOffline $false
```

```
# Create a volume with maximum size and default drive letter
New-Partition -DiskNumber <DiskNumber> -UseMaximumSize
-AssignDriveLetter
```

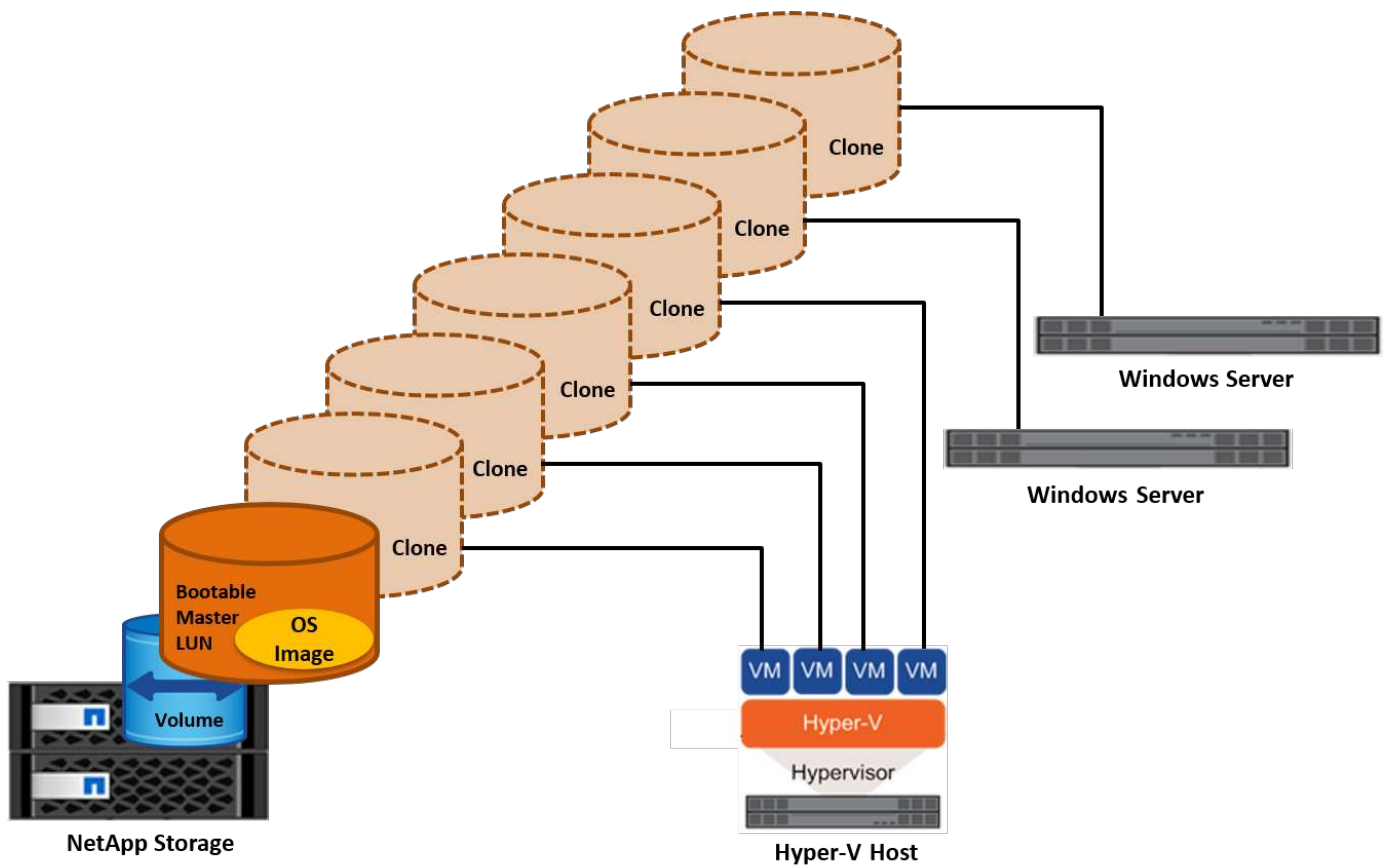
```
# To choose the size and drive letter use -Size and -DriveLetter parameters
# Format the volume
Format-Volume -DriveLetter <DriveLetter> -FileSystem <FAT32 or NTFS or REFS>
```

## 從 SAN 開機

實體主機（伺服器）或 Hyper-V VM 可直接從 NetApp LUN（而非其內部硬碟）啟動 Windows Server 作業系統。在從 SAN 開機的方法中、要從中開機的 OS 映像位於連接至實體主機或 VM 的 NetApp LUN 上。對於實體主機、實體主機的 HBA 會設定為使用 NetApp LUN 進行開機。對於 VM、NetApp LUN 會附加為用於開機的直接移轉磁碟。

### NetApp FlexClone 方法

使用 NetApp FlexClone 技術、可立即複製具有 OS 映像的開機 LUN、並將其附加至伺服器和 VM、以快速提供乾淨的 OS 映像、如下圖所示。



從 SAN 開機以供實體主機使用

### 先決條件

- 實體主機（伺服器）具有適當的 iSCSI 或 FC HBA。



- 您已為支援 Windows Server 的伺服器下載適當的 HBA 裝置驅動程式。
- 伺服器具有適當的 CD/DVD 磁碟機或虛擬媒體來插入 Windows Server ISO 映像、而且已下載 HBA 裝置驅動程式。
- NetApp iSCSI 或 FC LUN 是在 NetApp 儲存控制器上進行佈建。

## 部署

若要設定實體主機從 SAN 開機、請完成下列步驟：

1. 在伺服器 HBA 上啟用 BootBIOS。
2. 對於 iSCSI HBA、請在開機 BIOS 設定中設定啟動器 IP、iSCSI 節點名稱和介面卡開機模式。
3. 在 NetApp 儲存控制器上建立 iSCSI 和 / 或 FC 的啟動器群組時、請將伺服器 HBA 啟動器新增至群組。伺服器的 HBA 啟動器是 FC HBA 的 WWPN 或 iSCSI HBA 的 iSCSI 節點名稱。
4. 在 NetApp 儲存控制器上建立 LUN、LUN ID 為 0、並將其與上一步建立的啟動器群組建立關聯。此 LUN 可作為開機 LUN。
5. 將 HBA 限制為指向開機 LUN 的單一路徑。在開機 LUN 上安裝 Windows Server 之後、可以新增其他路徑、以利用多重路徑功能。
6. 使用 HBA 的 BootBIOS 公用程式、將 LUN 設定為開機裝置。
7. 重新啟動主機並進入主機 BIOS 公用程式。
8. 設定主機 BIOS、使開機 LUN 成為開機順序中的第一個裝置。
9. 從 Windows Server ISO 啟動安裝設定。
10. 當安裝詢問「您要在何處安裝 Windows？」時、請按一下安裝畫面底部的「載入驅動程式」、以啟動「選取要安裝的驅動程式」頁面。提供先前下載的 HBA 裝置驅動程式路徑、並完成驅動程式的安裝。
11. 現在、之前建立的開機 LUN 必須顯示在 Windows 安裝頁面上。選取開機 LUN 以在開機 LUN 上安裝 Windows Server、然後完成安裝。

從 **SAN** 開機以供虛擬機器使用

若要設定從 SAN 開機以供 VM 使用、請完成下列步驟：

## 部署

1. 在 NetApp 儲存控制器上建立 iSCSI 或 FC 的啟動器群組時、請將 iSCSI 的 IQN 或 Hyper-V 伺服器 FC 的 WWN 新增至控制器。
2. 在 NetApp 儲存控制器上建立 LUN 或 LUN 複本、並將它們與上一步建立的啟動器群組建立關聯。這些 LUN 可作為 VM 的開機 LUN。
3. 偵測 Hyper-V 伺服器上的 LUN、將其上線並初始化。
4. 使 LUN 離線。
5. 使用稍後在「Connect Virtual Hard Disk」頁面上的「Attach a Virtual Hard Disk」選項來建立 VM。
6. 將 LUN 新增為傳遞磁碟至 VM。
  - a. 開啟 VM 設定。
  - b. 按一下「IDE 控制器 0」、選取「硬碟」、然後按一下「新增」。選取 IDE 控制器 0、將此磁碟設為

VM 的第一個開機裝置。

- c. 在「硬碟」選項中選取「實體硬碟」、然後從清單中選取一個磁碟做為直接移轉磁碟。磁碟是在前述步驟中設定的 LUN。

7. 在傳遞磁碟上安裝 Windows Server。

#### 最佳實務做法

- 確定 LUN 已離線。否則、磁碟將無法新增為直接移轉磁碟至 VM。
- 當存在多個 LUN 時、請務必在磁碟管理中記下 LUN 的磁碟編號。這是必要的做法、因為列出給 VM 的磁碟會與磁碟編號一起列出。此外、將磁碟選擇為 VM 的直接移轉磁碟也會根據此磁碟編號而定。
- NetApp 建議避免 iSCSI NIC 的 NIC 群組。
- NetApp 建議您使用在主機上設定的 ONTAP MPIO 進行儲存。

## 在 SMB 環境中進行資源配置

ONTAP 使用 SMB3 傳輸協定、為 Hyper-V 虛擬機器提供彈性且高效能的 NAS 儲存設備。

使用 CIFS 通訊協定建立 SVM 時、CIFS 伺服器會在 Windows Active Directory 網域的 SVM 上執行。SMB 共用可用於主目錄、以及裝載 Hyper-V 和 SQL Server 工作負載。ONTAP 支援下列 SMB 3.0 功能：

- 持續處理（持續可用的檔案共用）
- 見證協定
- 叢集式用戶端容錯移轉
- 橫向擴充認知
- ODX
- 遠端 VSS

## 在 Windows Server 上配置 SMB 共享

### 先決條件

在 Windows Server 的 NAS 環境中使用 NetApp 儲存設備有下列需求：

- ONTAP 叢集具有有效的 CIFS 授權。
- 至少會建立一個 Aggregate。
- 系統會建立一個資料邏輯介面（LIF）、而且必須為 CIFS 設定資料 LIF。
- 存在 DNS 設定的 Windows Active Directory 網域伺服器和網域管理員認證。
- NetApp 叢集中的每個節點都會與 Windows 網域控制站同步處理時間。

### Active Directory 網域控制站

NetApp 儲存控制器可以與 Windows Server 類似、在 Active Directory 中加入及操作。在建立 SVM 期間、您可以提供網域名稱和名稱伺服器詳細資料來設定 DNS。SVM 會嘗試搜尋 Active Directory 網域控制站、方法是以類似 Windows Server 的方式、查詢 DNS 中的 Active Directory/ 輕量型目錄存取傳輸協定（LDAP）伺服器。

若要讓 CIFS 設定正常運作、NetApp 儲存控制器必須與 Windows 網域控制器同步時間。NetApp 建議 Windows 網域控制器與 NetApp 儲存控制器之間的時間偏差不超過五分鐘。最佳做法是設定 ONTAP 叢集的網路時間傳輸協定 (NTP) 伺服器、使其與外部時間來源同步。若要將 Windows 網域控制站設定為 NTP 伺服器、請在 ONTAP 叢集上執行下列命令：

```
$domainControllerIP = "<input IP Address of windows domain controller>"
cluster::> system services ntp server create -s "server $domainControllerIP
```

## 部署

1. 在啟用 NAS 傳輸協定 CIFS 的情況下、建立新的 SVM。您可以使用下列任一方法建立新的 SVM：
  - NetApp ONTAP 上的 CLI 命令
  - 系統管理員
  - NetApp PowerShell 工具套件
2. 設定 CIFS 通訊協定
  - a. 提供 CIFS 伺服器名稱。
  - b. 提供必須加入 CIFS 伺服器的 Active Directory。您必須擁有網域管理員認證、才能將 CIFS 伺服器加入 Active Directory。
3. 在每個叢集節點上指派具有生命的 SVM。
4. 在 SVM 上啟動 CIFS 服務。
5. 從 Aggregate 建立具有 NTFS 安全樣式的磁碟區。
6. 在捲上建立 qtree (選用)。
7. 建立對應於 Volume 或 qtree 目錄的共用、以便從 Windows Server 存取。如果共用用於 Hyper-V 儲存設備、請選取在建立共用時啟用 Hyper-V 的持續可用性。這樣做可為檔案共用提供高可用度。
8. 編輯建立的共用、並視需要修改存取共用的權限。SMB 共用的權限必須設定為授與存取此共用之所有伺服器的電腦帳戶存取權。

## 主機整合

NAS 傳輸協定 CIFS 原生整合至 ONTAP。因此、Windows Server 不需要任何額外的用戶端軟體即可存取 NetApp ONTAP 上的資料。NetApp 儲存控制器會在網路上顯示為原生檔案伺服器、並支援 Microsoft Active Directory 驗證。

若要偵測先前使用 Windows Server 建立的 CIFS 共用、請完成下列步驟：

1. 以系統管理員群組成員的身分登入 Windows Server。
2. 請移至 run.exe、輸入為存取共用而建立的 CIFS 共用的完整路徑。
3. 若要將共用永久對應至 Windows Server、請以滑鼠右鍵按一下「此電腦」、按一下「對應網路磁碟機」、然後提供 CIFS 共用的路徑。
4. 某些 CIFS 管理工作可以使用 Microsoft Management Console (MMC) 來執行。在執行這些工作之前、您必須使用 MMC 功能表命令、將 MMC 連接至 NetApp ONTAP 儲存設備。
  - a. 若要在 Windows Server 中開啟 MMC、請按一下 Server Manager 「工具」區段中的「電腦管理」。

- b. 按一下 [ 其他動作 ] 並連線到另一台電腦，這會開啟 [ 選取電腦 ] 對話方塊。
- c. 輸入 CIFS 伺服器名稱或 SVM LIF 的 IP 位址、以連線至 CIFS 伺服器。
- d. 展開「系統工具」和「共用資料夾」、以檢視和管理開啟的檔案、工作階段和共用。

#### 最佳實務做法

- 為了確認當磁碟區從一個節點移至另一個節點時、或在節點故障的情況下、不會發生停機、NetApp 建議您在檔案共用上啟用持續可用性選項。
- 為 Hyper-V over SMB 環境佈建 VM 時、NetApp 建議您在儲存系統上啟用複本卸載功能。如此可縮短 VM 的資源配置時間。
- 如果儲存叢集主控多個 SMB 工作負載、例如 SQL Server、Hyper-V 和 CIFS 伺服器、NetApp 建議將不同的 SMB 工作負載託管在不同的 SVM 上。這項組態非常實用、因為每項工作負載都需要獨特的儲存網路和 Volume 配置。
- NetApp 建議您將 Hyper-V 主機與 NetApp ONTAP 儲存設備連線至 10GB 網路（如果有的話）。如果是 1GB 網路連線、NetApp 建議您建立由多個 1GB 連接埠組成的介面群組。
- 將 VM 從一個 SMB 3.0 共用移轉至另一個 SMB 3.0 共用時、NetApp 建議在儲存系統上啟用 CIFS 複本卸載功能、以便更快移轉。

#### 值得記住的事項

- 為 SMB 環境佈建磁碟區時、必須以 NTFS 安全樣式建立磁碟區。
- 叢集中節點上的時間設定應相應設定。如果 NetApp CIFS 伺服器必須參與 Windows Active Directory 網域、請使用 NTP。
- 持續處理只能在 HA 配對中的節點之間運作。
- 見證通訊協定只能在 HA 配對中的節點之間運作。
- 只有 Hyper-V 和 SQL Server 工作負載才支援持續可用的檔案共用。
- ONTAP 9.4 以上版本支援 SMB 多通道。
- 不支援 RDMA。
- 不支援 Refs。

#### 在奈米伺服器上配置 **SMB** 共享

nano 伺服器不需要額外的用戶端軟體、即可存取 NetApp 儲存控制器上 CIFS 共用區上的資料。

若要將檔案從奈米伺服器複製到 CIFS 共用、請在遠端伺服器上執行下列 Cmdlet：

```
$ip = "<input IP Address of the Nano Server>"
```

```
# Create a New PS Session to the Nano Server  
$session = New-PSSession -ComputerName $ip -Credential ~\Administrator
```

```
Copy-Item -FromSession $s -Path C:\Windows\Logs\DISM\dism.log
-Destination \\cifsshare
```

\* `cifsshare` 是 NetApp 儲存控制器上的 CIFS 共用。

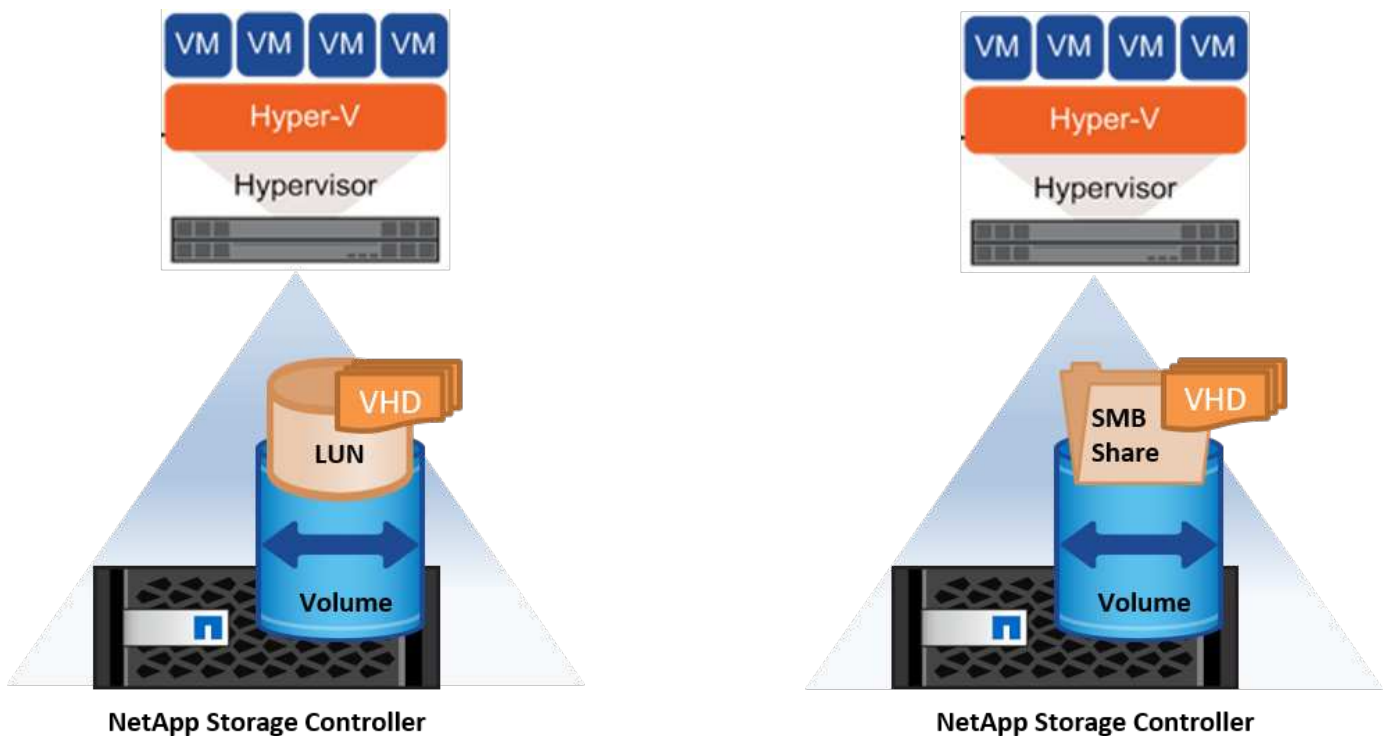
\* 若要將檔案複製到奈米伺服器、請執行下列 Cmdlet ：

```
+
Copy-Item -ToSession $s -Path \\cifsshare\<file> -Destination C:\
```

若要複製資料夾的完整內容、請指定資料夾名稱、並使用 Cmdlet 結尾的 -Recurse 參數。

## NetApp 上的 Hyper-V 儲存基礎架構

Hyper-V 儲存基礎架構可裝載於 ONTAP 儲存系統上。Hyper-V 儲存以儲存 VM 檔案及其磁碟的儲存設備可使用 NetApp LUN 或 NetApp CIFS 共用來提供、如下圖所示。



### NetApp LUN 上的 Hyper-V 儲存設備

- 在 Hyper-V 伺服器機器上佈建 NetApp LUN 。如需詳細資訊、請參閱「[在 SAN 環境中進行資源配置](#)」。
- 從 Server Manager 的「工具」區段開啟 Hyper-V Manager 。
- 選取 Hyper-V 伺服器、然後按一下 Hyper-V 設定。
- 指定將 VM 及其磁碟儲存為 LUN 的預設資料夾。這樣做會將預設路徑設為 Hyper-V 儲存設備的 LUN 。如果您想要明確指定 VM 的路徑、則可以在建立 VM 時指定路徑。

## NetApp CIFS 上的 Hyper-V 儲存設備

在開始執行本節所列的步驟之前、請先檢閱「["在 SMB 環境中進行資源配置"](#)。」 若要在 NetApp CIFS 共用上設定 Hyper-V 儲存設備、請完成下列步驟：

1. 從 Server Manager 的「工具」區段開啟 Hyper-V Manager。
2. 選取 Hyper-V 伺服器、然後按一下 Hyper-V 設定。
3. 指定將 VM 及其磁碟儲存為 CIFS 共用的預設資料夾。這樣做會將預設路徑設為 Hyper-V 儲存設備的 CIFS 共用。如果您想要明確指定 VM 的路徑、則可以在建立 VM 時指定路徑。

Hyper-V 中的每個虛擬機器也可隨附提供給實體主機的 NetApp LUN 和 CIFS 共用。此程序與任何實體主機相同。下列方法可用於將儲存設備配置至 VM：

- 使用 VM 中的 FC 啟動器新增儲存 LUN
- 使用虛擬機器中的 iSCSI 啟動器新增儲存 LUN
- 將直接移轉實體磁碟新增至 VM
- 從主機將 VHD/VHDX 新增至 VM

### 最佳實務做法

- 當 VM 及其資料儲存在 NetApp 儲存設備上時、NetApp 建議定期在 Volume 層級執行 NetApp 重複資料刪除。當相同的 VM 託管在 CSV 或 SMB 共享上時、這種做法可大幅節省空間。重複資料刪除功能會在儲存控制器上執行、不會影響主機系統和 VM 的效能。
- 在 Hyper-V 使用 iSCSI LUN 時、請務必啟用 iSCSI Service (TCP-In) for Inbound 和 iSCSI Service (TCP-Out) for Outbound 在 Hyper-V 主機的防火牆設定中。如此一來、iSCSI 流量就能往返於 Hyper-V 主機和 NetApp 控制器。
- NetApp 建議取消勾選允許管理作業系統共用 Hyper-V 虛擬交換器的此網路介面卡選項。這樣做可為 VM 建立專用網路。

### 值得記住的事項

- 使用虛擬光纖通道來配置 VM 需要啟用 N\_Port ID 虛擬化功能的 FC HBA。最多支援四個 FC 連接埠。
- 如果主機系統已設定多個 FC 連接埠並呈現給 VM、則必須在 VM 中安裝 MPIO 才能啟用多重路徑功能。
- 如果在該主機上使用 MPIO、則無法將直接移轉磁碟配置至主機、因為直接移轉磁碟不支援 MPIO。
- 用於 VHD/VHDX 檔案的磁碟應該使用 64K 格式進行配置。

### 進一步閱讀

- 如需 FC HBA 的相關資訊、請參閱 ["NetApp 互通性對照表"](#)。
- 如需虛擬光纖通道的詳細資訊、請參閱 Microsoft ["Hyper-V 虛擬光纖通道總覽"](#) 頁面。

### 卸載資料傳輸

Microsoft ODX 也稱為複製卸載、可在儲存裝置內或相容儲存裝置之間直接傳輸資料、而無需透過主機電腦傳輸資料。NetApp ONTAP 支援 CIFS 與 SAN 傳輸協定的 ODX 功能。如果複本位於同一個磁碟區內、ODX 可能會改善效能、降低用戶端 CPU 和記憶體的使用率、並降低網路 I/O 頻寬使用率。

有了 ODX、在 SMB 共用區、LUN 內、以及 SMB 共用區與 LUN 之間（如果位於同一個磁碟區）複製檔案、速度更快、效率更高。此方法在相同磁碟區中需要多個作業系統（HDD/VHDX）黃金映像複本的情況下更有幫助。如果複本位於同一個磁碟區內、則可在大幅減少的時間內製作同一個黃金映像的數個複本。ODX 也適用於 Hyper-V 儲存即時移轉、可用於移動 VM 儲存設備。

如果複本是跨磁碟區的、相較於主機型複本、效能可能不會大幅提升。

若要在 CIFS 上啟用 ODX 功能、請在 NetApp 儲存控制器上執行下列 CLI 命令：

1. 啟用適用於 CIFS 的 ODX。  
# 將權限等級設為診斷  
叢集：：> 設定權限診斷

```
#enable the odx feature
cluster::> vserver cifs options modify -vserver <vserver_name> -copy
-offload-enabled true
```

```
#return to admin privilege level
cluster::> set privilege admin
```

2. 若要在 SAN 上啟用 ODX 功能、請在 NetApp 儲存控制器上執行下列 CLI 命令：  
# 將權限等級設為診斷  
叢集：：> 設定權限診斷

```
#enable the odx feature
cluster::> copy-offload modify -vserver <vserver_name> -scsi enabled
```

```
#return to admin privilege level
cluster::> set privilege admin
```

值得記住的事項

- 對於 CIFS、只有當用戶端和儲存伺服器都支援 SMB 3.0 和 ODX 功能時、ODX 才會提供使用。
- 對於 SAN 環境、只有當用戶端和儲存伺服器都支援 ODX 功能時、ODX 才可用。

進一步閱讀

如需 ODX 的相關資訊、請參閱 ["改善 Microsoft 遠端複製效能"](#) 和 ["Microsoft 卸載資料傳輸"](#)。

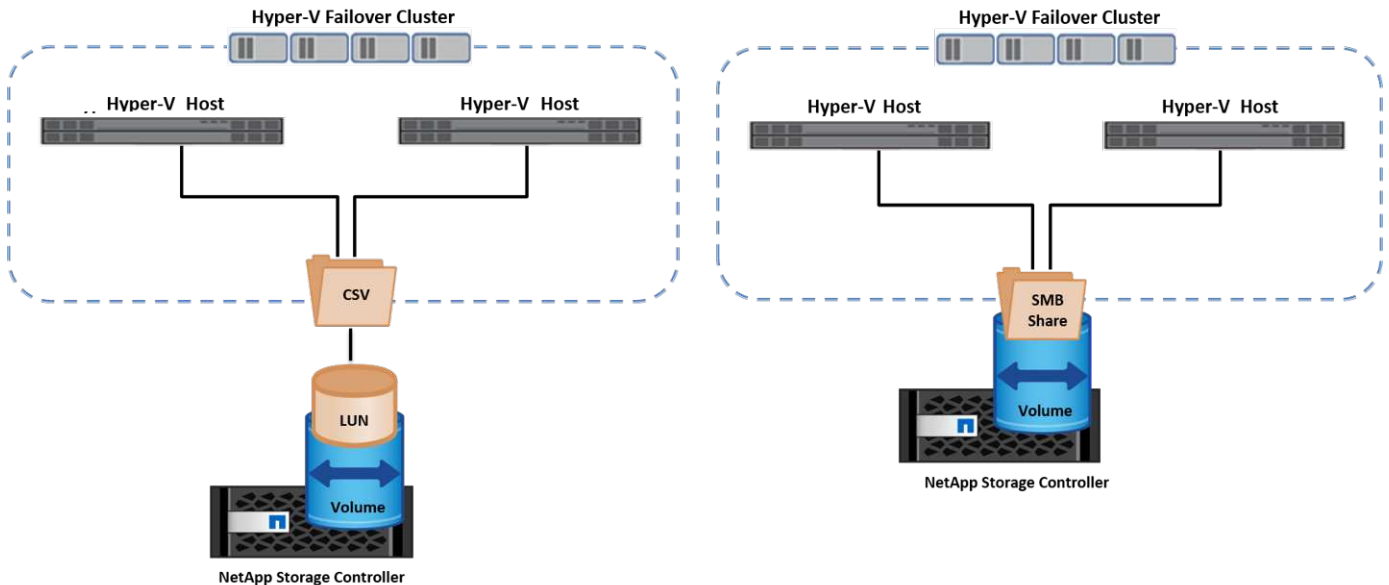
## Hyper-V 叢集：虛擬機器的高可用度與擴充性

容錯移轉叢集可為 Hyper-V 伺服器提供高可用度和擴充性。容錯移轉叢集是一組相互支援的 Hyper-V 伺服器、可一起提高 VM 的可用度和擴充性。

Hyper-V 叢集式伺服器（稱為節點）是由實體網路和叢集軟體所連接。這些節點使用共用儲存設備來儲存 VM 檔案、包括組態、虛擬硬碟（VHD）檔案和 Snapshot 複本。共享儲存設備可以是 NetApp SMB/CIFS 共用、或是 NetApp LUN 上的 CSV、如圖 6 所示。此共享儲存設備提供一致且分散的命名空間、可由叢集中的所有節點同時存取。因此、如果叢集中有一個節點發生故障、另一個節點會透過稱為容錯移轉的程序來提供服務。您可以使用容錯移轉叢集管理單元和容錯移轉叢集 Windows PowerShell Cmdlet 來管理容錯移轉叢集。

#### 叢集共享磁碟區

CSV 可讓容錯移轉叢集中的多個節點同時擁有與 NTFS 或 Refs 磁碟區相同的 NetApp LUN 的讀取 / 寫入存取權。透過 CSV、叢集式角色可以從一個節點快速容錯移轉至另一個節點、而無需變更磁碟機擁有權或卸除及重新掛載磁碟區。CSV 也能簡化容錯移轉叢集中可能大量 LUN 的管理。CSV 提供一般用途的叢集式檔案系統、其分層位於 NTFS 或 Refs 之上。



#### 最佳實務做法

- NetApp 建議關閉 iSCSI 網路上的叢集通訊、以防止內部叢集通訊和 CSV 流量流經同一個網路。
- NetApp 建議使用備援網路路徑（多個交換器）來提供恢復能力和 QoS。

#### 值得記住的事項

- 用於 CSV 的磁碟必須使用 NTFS 或 Refs 進行分割。使用 FAT 或 FAT32 格式化的磁碟無法用於 CSV。
- 用於 CSV 的磁碟應使用 64K 格式進行配置。

#### 進一步閱讀

如需部署 Hyper-V 叢集的相關資訊、請參閱附錄 B：["部署 Hyper-V 叢集"](#)。

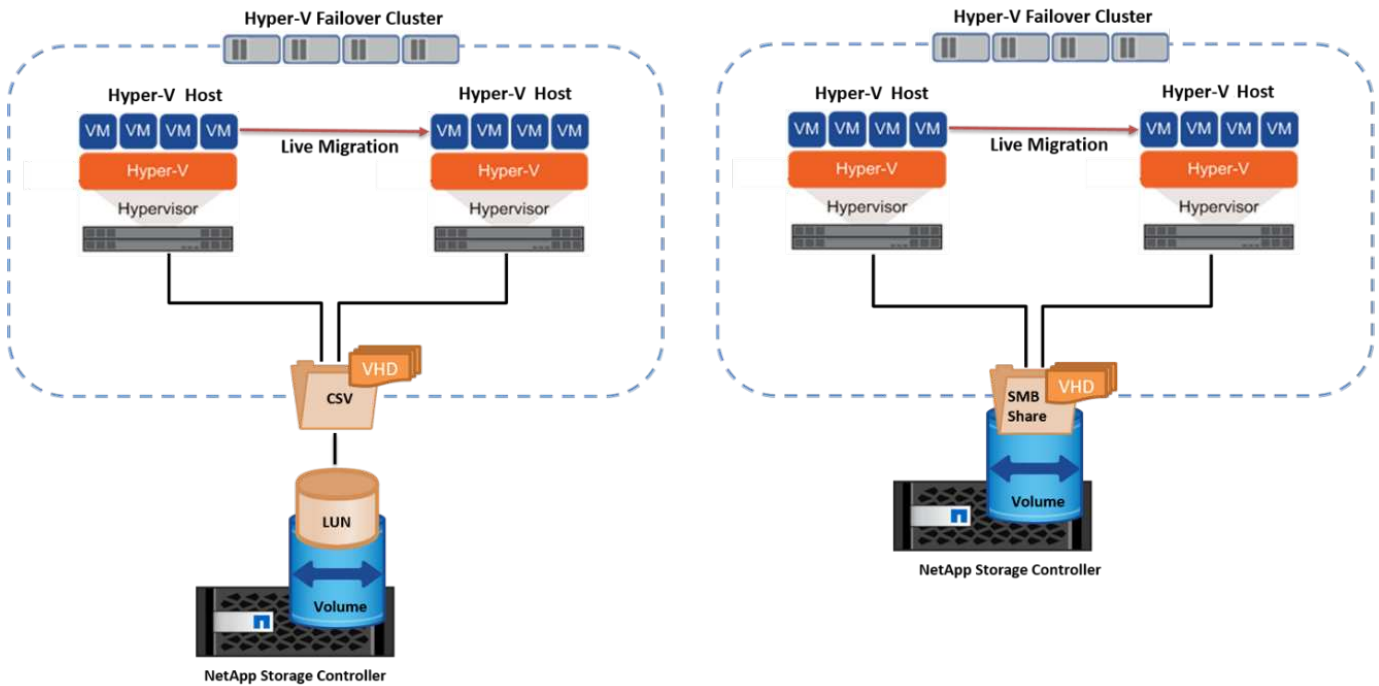
#### Hyper-V 線上即時移轉：VM 移轉

有時在 VM 的生命週期內、必須將其移至 Windows 叢集上的不同主機。如果主機的系統資源不足、或由於維護原因而需要重新開機、則可能需要這麼做。同樣地、可能需要將 VM 移至不同的 LUN 或 SMB 共用區。如果目前的 LUN 或共用區空間不足或效能低於預期、則可能需要這項功能。Hyper-V 線上即時移轉功能可將執行中的 VM 從一部實體 Hyper-V 伺服器移轉至另一部伺服器、對使用者的 VM 可用度沒有影響。您可以在屬於容錯移轉叢集一部分的 Hyper-V 伺服器之間、或是在不屬於任何叢集的單獨 Hyper-V 伺服器之間、即時移轉 VM。



在叢集式環境中進行即時移轉

VM 可在叢集的節點之間無縫移動。VM 移轉是即時的、因為叢集中的所有節點都共用相同的儲存設備、而且可以存取 VM 及其磁碟。下圖說明叢集環境中的即時移轉。



最佳實務做法

- 擁有專屬連接埠、可進行即時移轉流量。
- 擁有專用的主機即時移轉網路、以避免移轉期間發生與網路相關的問題。

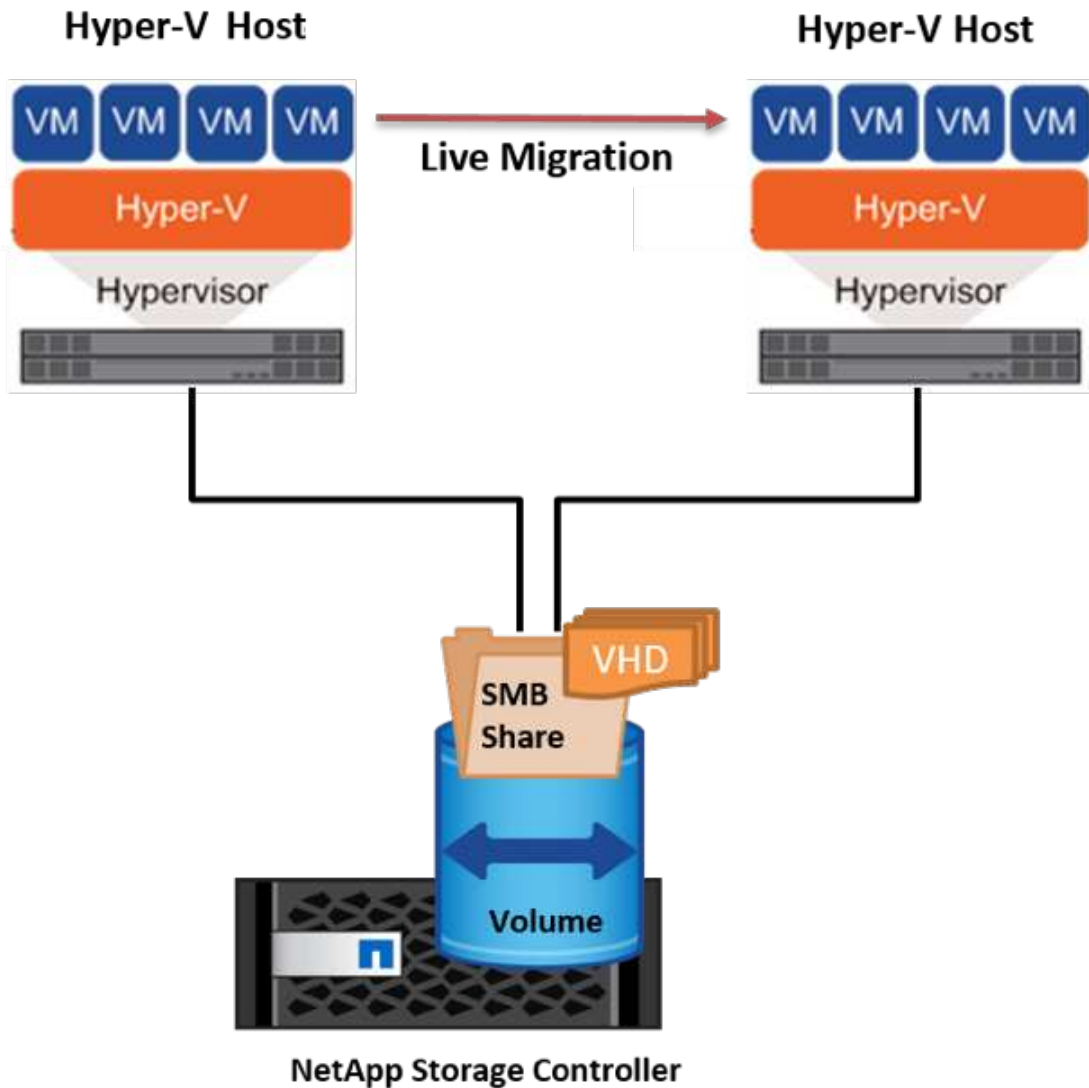
進一步閱讀

如需在叢集環境中部署即時移轉的相關資訊、請參閱 "[附錄 C：在叢集環境中部署 Hyper-V 線上即時移轉](#)"。

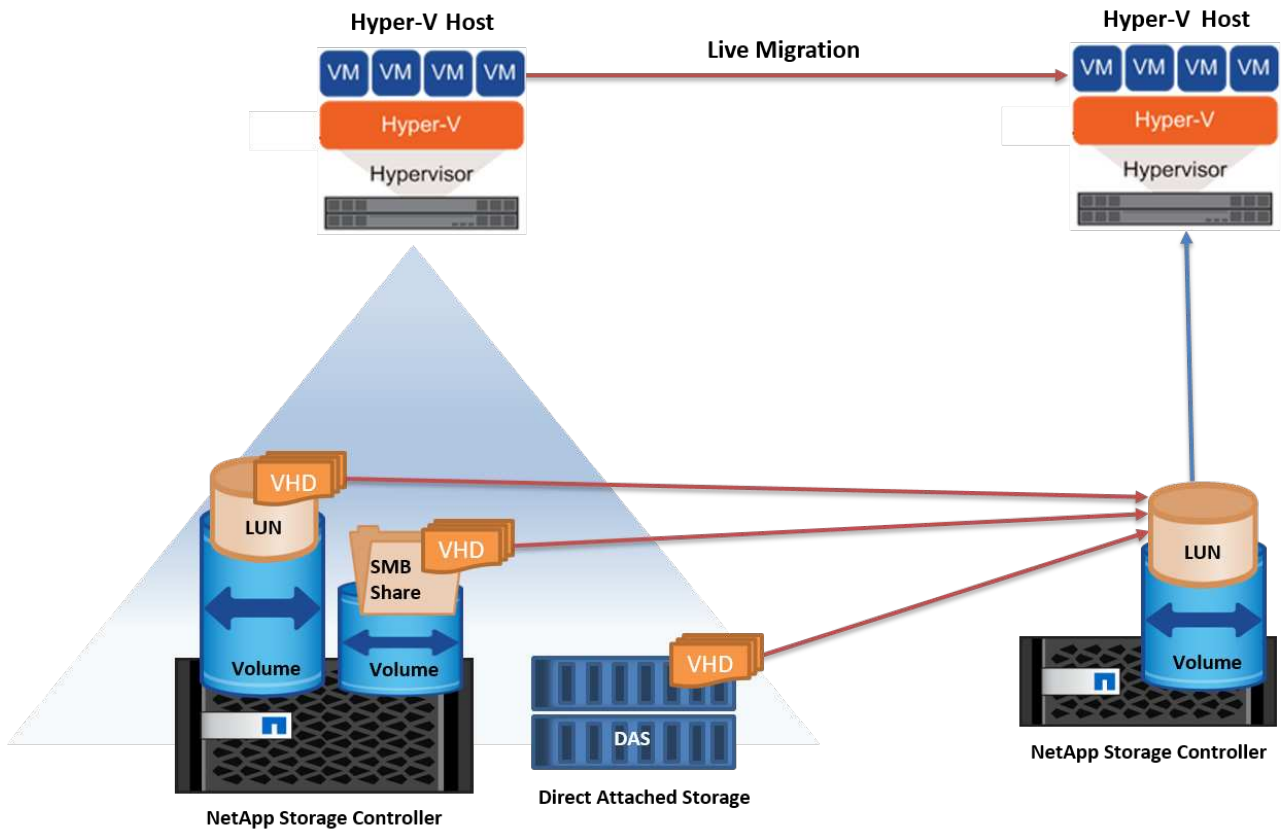
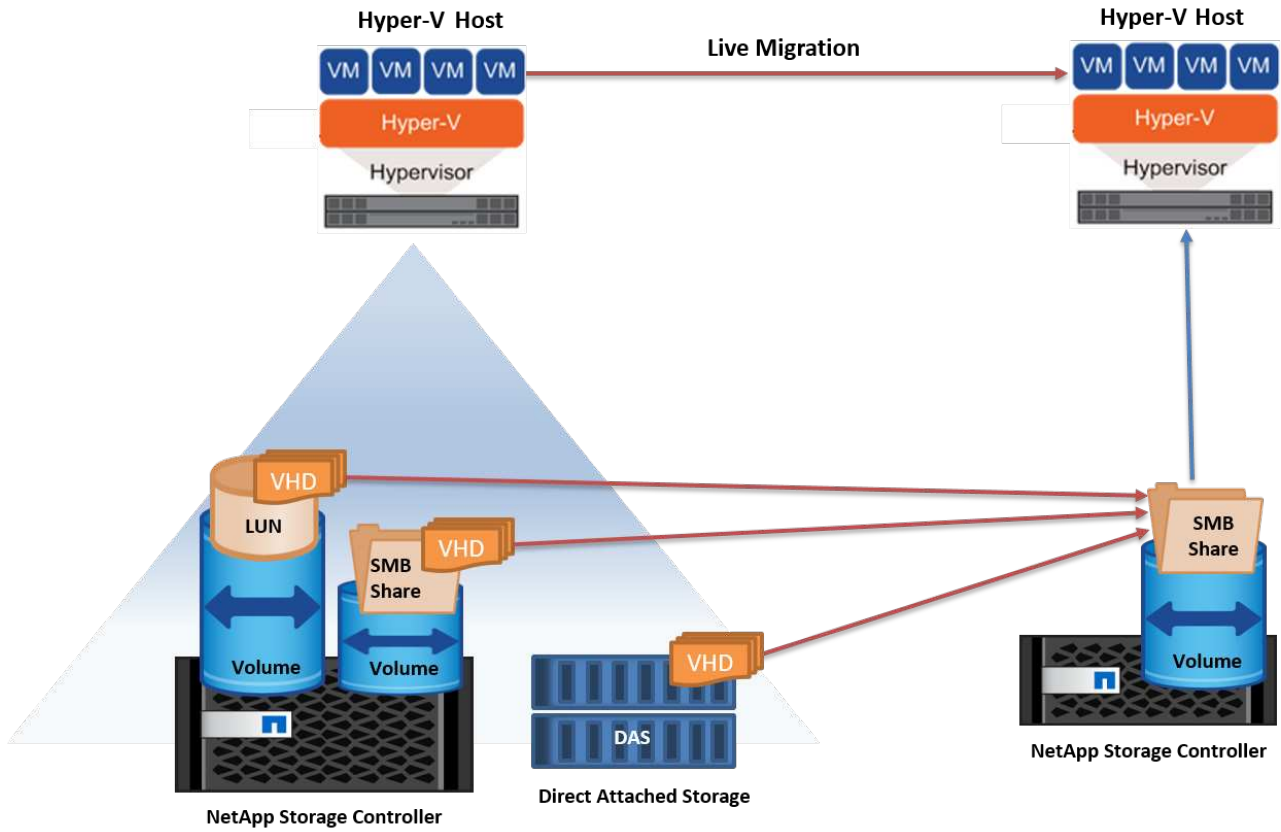
在叢集式環境外進行即時移轉

您可以在兩個非叢集式、不相互關聯的 Hyper-V 伺服器之間即時移轉 VM。此程序可以使用共用或不共用的即時移轉。

- 在共用的即時移轉中、虛擬機器會儲存在 SMB 共用區中。因此、當您即時移轉虛擬機器時、虛擬機器的儲存設備會保留在中央 SMB 共用區上、以便其他節點立即存取、如下圖所示。



- 在「共享無內容即時移轉」中、每部 Hyper-V 伺服器都有自己的本機儲存設備（可以是 SMB 共享區、LUN 或 DAS）、而 VM 的儲存設備則是其 Hyper-V 伺服器的本機儲存設備。VM 在線上即時移轉時、VM 的儲存設備會透過用戶端網路鏡射到目的地伺服器、然後再移轉 VM。儲存在 DAS、LUN 或 SMB/CIFS 共用區上的虛擬機器可移至其他 Hyper-V 伺服器上的 SMB/CIFS 共用區、如下圖所示。也可將其移至 LUN、如第二個圖所示。



進一步閱讀

如需在叢集環境外部部署即時移轉的相關資訊、請參閱 "附錄 D：在叢集環境之外部署 Hyper-V 即時移轉"。

## Hyper-V 儲存即時移轉

在虛擬機器的生命週期內、您可能需要將虛擬機器儲存設備（HDD/VHDX）移至不同的 LUN 或 SMB 共享區。如果目前的 LUN 或共享區空間不足或效能低於預期、則可能需要這項功能。

目前裝載 VM 的 LUN 或共用區可能會用盡空間、重新規劃用途、或是降低效能。在這種情況下、虛擬機器可以在不停機的情況下移至另一個 LUN、或在不同的磁碟區、集合或叢集上共用。如果儲存系統具備複製卸載功能、此程序就會更快完成。NetApp 儲存系統預設為啟用 CIFS 和 SAN 環境的複製卸載。

ODX 功能可在位於遠端伺服器上的兩個目錄之間執行完整檔案或子檔案複本。複本是透過在伺服器之間複製資料來建立（如果來源和目的地檔案都在同一部伺服器上、則複製資料也會複製到同一部伺服器）。建立複本時、用戶端不會從來源讀取資料、也不會寫入目的地。此程序可減少用戶端或伺服器的處理器和記憶體使用量、並將網路 I/O 頻寬降至最低。如果複本位於同一個磁碟區內、則複本速度會更快。如果複本是跨磁碟區的、相較於主機型複本、效能可能不會大幅提升。在繼續主機上的複本作業之前、請確認儲存系統上已設定複本卸載設定。

從主機啟動 VM 儲存即時移轉時、會識別來源和目的地、並將複製活動卸載至儲存系統。由於活動是由儲存系統執行、因此主機 CPU、記憶體或網路的使用率可忽略不計。

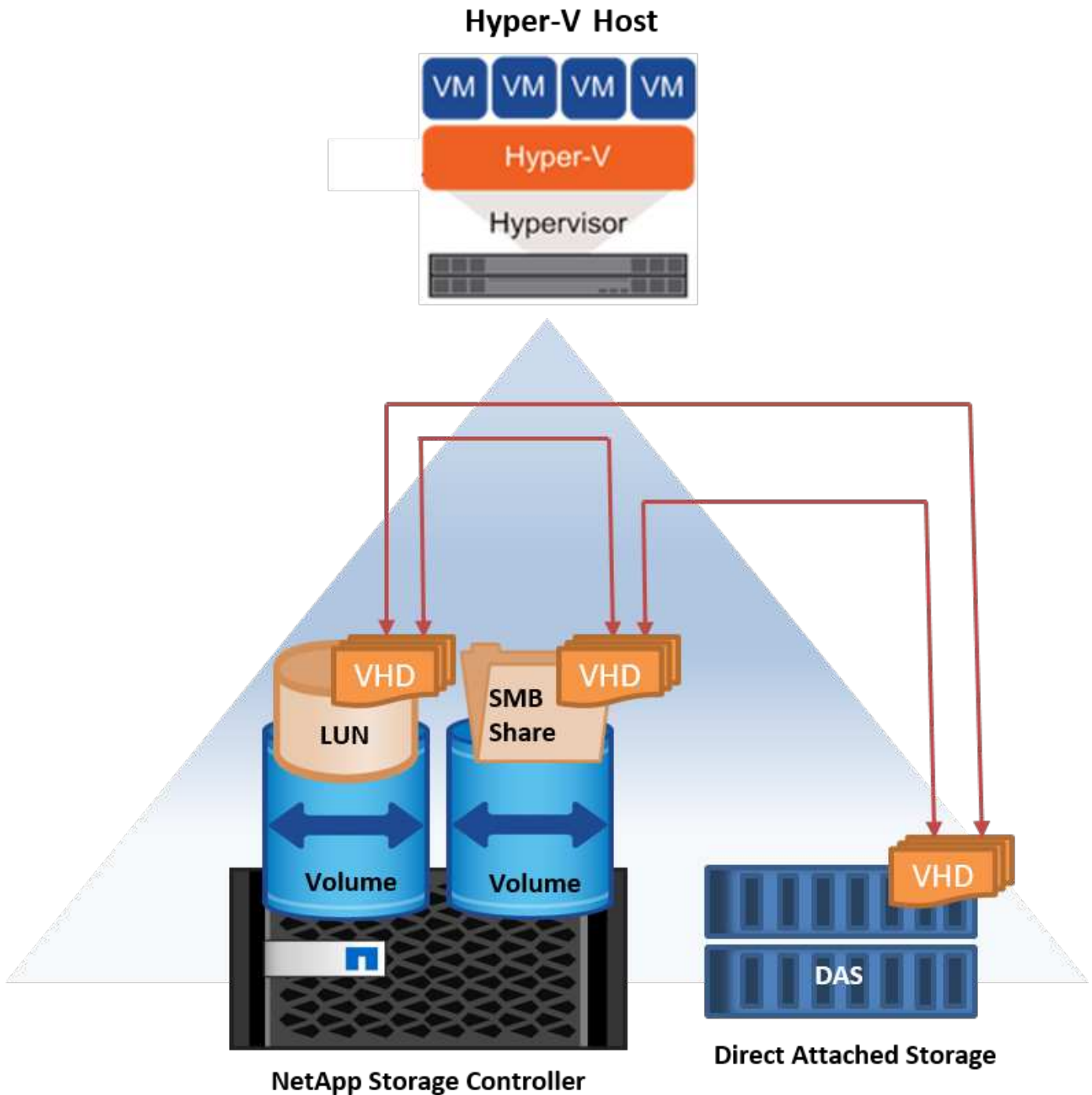
NetApp 儲存控制器支援下列不同的 ODX 情境：

- \* IntraSVM 。 \* 資料歸同一個 SVM 所有：
- \* Intravolume 、 inIntranet 模式 \* 。來源和目的地檔案或 LUN 位於同一個磁碟區內。複本是使用 FlexClone 檔案技術執行、可提供額外的遠端複本效能優勢。
- \* 磁碟區間、內部網路模式 \* 。來源和目的地檔案或 LUN 位於同一個節點上的不同磁碟區上。
- \* 磁碟區間、節點間。 \* 來源和目的地檔案或 LUN 位於不同節點上的不同磁碟區。
- \* InterSVM 。 \* 資料由不同的 SVM 擁有。
- \* 磁碟區間、內部網路模式 \* 。來源和目的地檔案或 LUN 位於同一個節點上的不同磁碟區上。
- \* 磁碟區間、節點間。 \* 來源和目的地檔案或 LUN 位於不同節點上的不同磁碟區。
- \* 叢集間。 \* 從 ONTAP 9.0 開始、ODX 也支援 SAN 環境中的叢集間 LUN 傳輸。叢集間 ODX 僅支援 SAN 通訊協定、不支援 SMB 。

移轉完成後、必須重新設定備份和複寫原則、以反映存放 VM 的新磁碟區。任何先前所進行的備份都無法使用。

VM 儲存設備（HDD/VHDX）可在下列儲存類型之間移轉：

- DAS 和 SMB 共享
- DAS 和 LUN
- SMB 共享區和 LUN
- 在 LUN 之間
- 在 SMB 共享之間

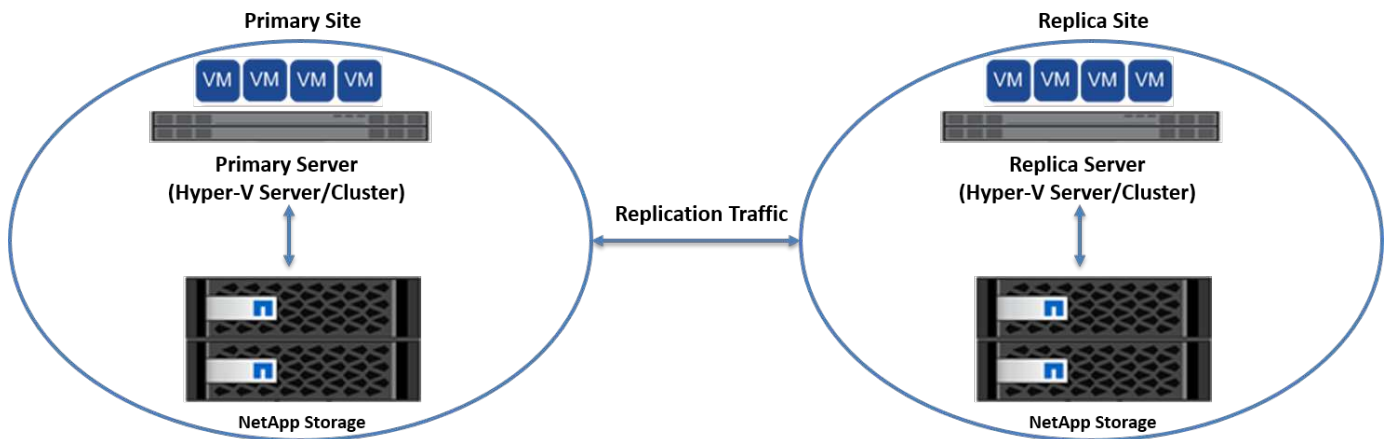


進一步閱讀

如需部署儲存即時移轉的相關資訊、請參閱 ["附錄 E：部署 Hyper-V 儲存即時移轉"](#)。

### Hyper-V 複本：虛擬機器的災難恢復

Hyper-V 複本會將 Hyper-V VM 從主要站台複寫到次要站台上的複本 VM、以非同步方式為 VM 提供災難恢復。主站台上的 Hyper-V 伺服器稱為主伺服器；次要站台上接收複寫 VM 的 Hyper-V 伺服器稱為複本伺服器。下圖顯示 Hyper-V 複本範例案例。您可以在屬於容錯移轉叢集一部分的 Hyper-V 伺服器之間、或是在不屬於任何叢集的單獨 Hyper-V 伺服器之間、使用 Hyper-V 複本來處理 VM。



## 複寫

在主伺服器上為 VM 啟用 Hyper-V 複本之後、初始複寫會在複本伺服器上建立相同的 VM。初始複寫之後、Hyper-V 複本會維護 VM VHD 的記錄檔。根據複寫頻率、以相反順序將記錄檔重新播放至複本 VHD。此記錄和反向順序的使用可確保以非同步方式儲存和複寫最新的變更。如果複寫未與預期頻率一致、就會發出警示。

## 延伸複寫

Hyper-V 複本支援延伸複寫、可在其中設定次要複本伺服器以進行災難恢復。您可以設定次要複本伺服器、讓複本伺服器接收複本 VM 上的變更。在延伸複寫案例中、主要伺服器上主要 VM 上的變更會複寫到複本伺服器。然後將變更複寫到擴充複本伺服器。只有當主要和複本伺服器都停機時、VM 才能容錯移轉至延伸複本伺服器。

## 容錯移轉

容錯移轉不是自動的；程序必須手動觸發。容錯移轉有三種類型：

- \* 測試容錯移轉。\* 此類型用於驗證複本 VM 是否能在複本伺服器上成功啟動、並在複本 VM 上啟動。此程序會在容錯移轉期間建立重複的測試 VM、不會影響正常的正式作業複寫。
- \* 計畫性容錯移轉。\* 此類型用於在計畫性停機或預期停機期間容錯移轉 VM。此程序是在主 VM 上啟動、必須在主伺服器上關閉、然後才會執行規劃的容錯移轉。機器容錯移轉後、Hyper-V 複本會在複本伺服器上啟動複本 VM。
- \* 非計畫性容錯移轉。\* 發生非預期的中斷時、可使用此類型。此程序是在複本 VM 上啟動、只有在主機器故障時才應使用。

## 恢復

當您設定虛擬機器的複寫時、可以指定恢復點的數量。恢復點代表可從複寫機器恢復資料的時間點。

## 進一步閱讀

- 如需在叢集環境外部署 Hyper-V 複本的相關資訊、請參閱["在叢集環境之外部署 Hyper-V 複本"](#)。」
- 如需在叢集環境中部署 Hyper-V 複本的相關資訊、請參閱["在叢集環境中部署 Hyper-V 複本"](#)。」

## 儲存效率

ONTAP 為虛擬化環境（包括 Microsoft Hyper-V）提供領先業界的儲存效率 NetApp 也提供儲存效率保證方案。

## NetApp 重複資料刪除

NetApp 重複資料刪除的運作方式是在儲存磁碟區層級移除重複的區塊、只儲存一個實體複本、無論有多少個邏輯複本。因此、重複資料刪除會產生一種錯覺、認為該區塊有許多複本。重複資料刪除功能會在整個磁碟區的 4KB 區塊層級上自動移除重複的資料區塊。此程序會重新宣告儲存設備、藉由減少實體寫入磁碟的次數、以達到空間和潛在的效能節約。重複資料刪除技術可在 Hyper-V 環境中節省 70% 以上的空間。

## 資源隨需配置

精簡配置是配置儲存設備的有效方法、因為儲存設備並未預先配置。換句話說、當磁碟區或 LUN 是使用精簡配置建立時、儲存系統上的空間就會被閒置。在資料寫入 LUN 或磁碟區之前、該空間會一直保持未使用狀態、而且只會使用儲存資料所需的空間。NetApp 建議在磁碟區上啟用精簡配置、並停用 LUN 保留。

## 服務品質

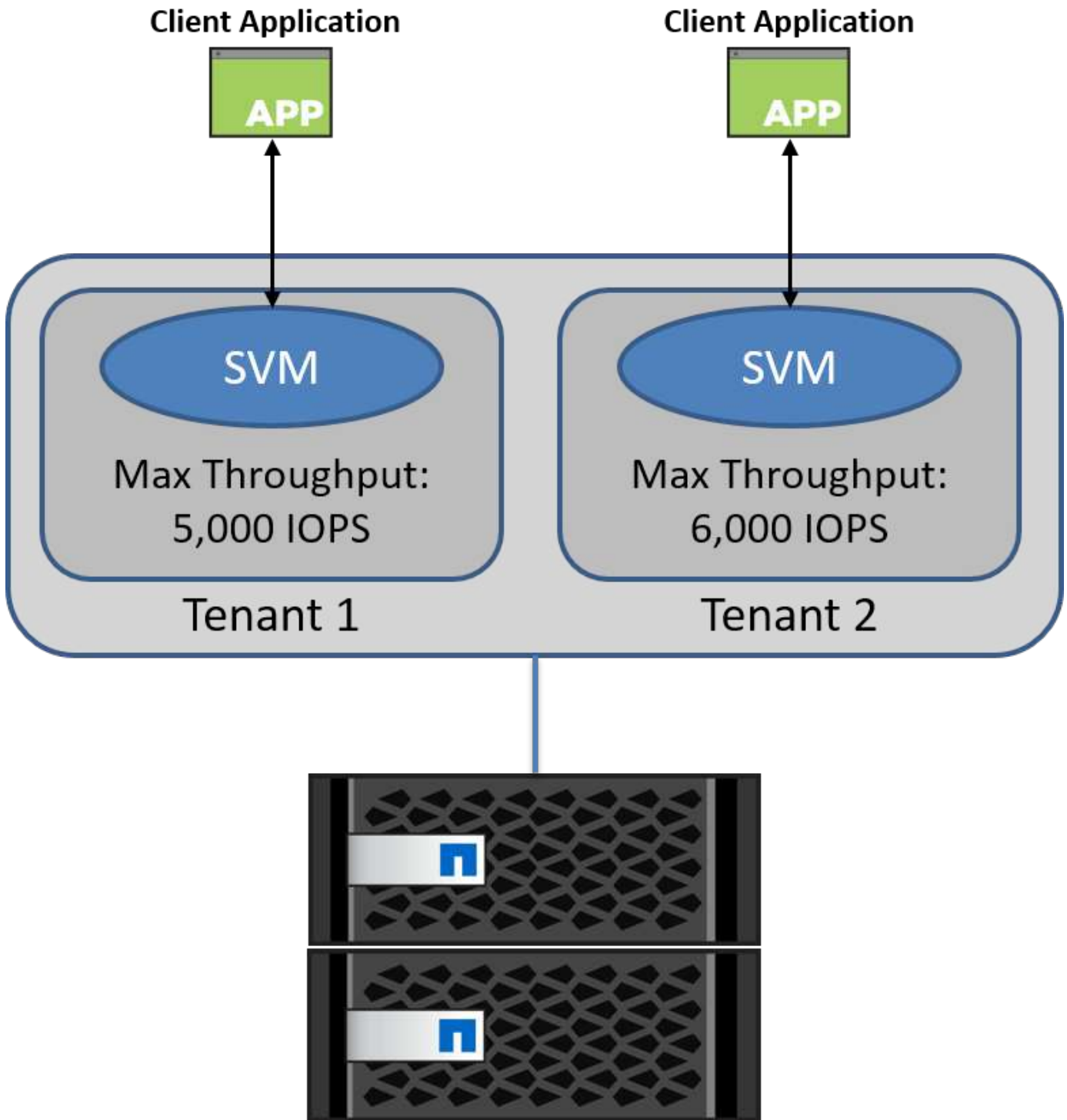
叢集式 ONTAP 中的儲存 QoS 可讓您將儲存物件分組、並設定群組的處理量限制。儲存 QoS 可用於限制工作負載的處理量、並監控工作負載效能。有了這項功能、儲存管理員就能依組織、應用程式、業務單位、或正式作業或開發環境來分隔工作負載。

在企業環境中、儲存 QoS 有助於達成下列目標：

- 防止使用者工作負載彼此影響。
- 保護在 IT 即服務 (ITaaS) 環境中必須符合特定回應時間的關鍵應用程式。
- 防止租戶彼此影響。
- 新增每個租戶、避免效能降低。

QoS 可讓您限制傳送至 SVM、彈性磁碟區、LUN 或檔案的 I/O 量。I/O 可受作業數或原始處理量的限制。

下圖說明 SVM 及其本身的 QoS 原則、可強制執行最大處理量限制。



若要使用自己的 QoS 原則設定 SVM 並監控原則群組、請在 ONTAP 叢集上執行下列命令：

```
# create a new policy group pgl with a maximum throughput of 5,000 IOPS
cluster::> qos policy-group create pgl -vserver vs1 -max-throughput
5000iops
```



```
# create a new policy group pg2 without a maximum throughput
cluster::> qos policy-group create pg2 -vserver vs2
```

```
# monitor policy group performance
cluster::> qos statistics performance show
```

```
# monitor workload performance
cluster::> qos statistics workload performance show
```

## 安全性

ONTAP 為 Windows 作業系統提供安全的儲存系統。

### Windows Defender 防毒軟體

Windows Defender 是預設在 Windows Server 上安裝及啟用的反惡意程式碼軟體。此軟體可主動保護 Windows Server 免於已知惡意軟體的侵害、並可透過 Windows Update 定期更新反惡意程式碼定義。NetApp LUN 和 SMB 共用可以使用 Windows Defender 進行掃描。

進一步閱讀

如需詳細資訊、請參閱 "[Windows Defender 概觀](#)"。

### BitLocker

BitLocker 磁碟機加密是一項資料保護功能、延續自 Windows Server 2012。此加密可保護實體磁碟、LUN 和 CSV。

最佳實務做法

啟用 BitLocker 之前、必須先將 CSV 設為維護模式。因此、NetApp 建議您在 CSV 上建立 VM 之前、先做出有關以 BitLocker 為基礎的安全性的決策、以避免停機。

## 部署奈米伺服器

瞭解如何部署 Microsoft Windows Nano Server。

部署

若要將奈米伺服器部署為 Hyper-V 主機、請完成下列步驟：

1. 以系統管理員群組成員的身分登入 Windows Server。
2. 將 Windows Server ISO 中的 \NanoServer 資料夾中的 NanoServerImageGenerator 資料夾複製到本機硬碟。
3. 若要建立奈米伺服器 VHD/VHDX、請完成下列步驟：

- a. 以系統管理員身分啟動 Windows PowerShell、瀏覽至本機硬碟上複製的 NanoServerImageGenerator 資料夾、然後執行下列 Cmdlet：

```
Set-ExecutionPolicy RemoteSigned
Import-Module .\NanoServerImageGenerator -Verbose
```

- b. 執行下列 PowerShell Cmdlet、為奈米伺服器建立 VHD 做為 Hyper-V 主機。此命令會提示您輸入新 VHD 的系統管理員密碼。

```
New-NanoServerImage -Edition Standard -DeploymentType Guest
-MediaPath <"input the path to the root of the contents of Windows
Server 2016 ISO"> -TargetPath <"input the path, including the
filename and extension where the resulting VHD/VHDX will be created">
-ComputerName <"input the name of the nano server computer you are
about to create"> -Compute
```

.. 在下列範例中、我們建立了支援啟用容錯移轉叢集功能 Hyper-V 主機的 Nano Server VHD。此範例是從安裝於 f:\ 的 ISO 建立 Nano Server VHD。新建立的 VHD 會放在執行 Cmdlet 的資料夾中名為 NanoServer 的資料夾中。電腦名稱為 NanoServer、而產生的 VHD 包含 Windows Server 的標準版本。

```
New-NanoServerImage -Edition Standard -DeploymentType Guest
-MediaPath f:\ -TargetPath .\NanoServer.vhd -ComputerName NanoServer
-Compute -Clustering
```

.. 使用 Cmdlet New-NanoServerImage 設定參數來設定 IP 位址、子網路遮罩、預設閘道、DNS 伺服器、網域名稱、等等。

4. 在 VM 或實體主機中使用 VHD 將 Nano Server 部署為 Hyper-V 主機：

- a. 若要在虛擬機器上部署、請在 Hyper-V Manager 中建立新的虛擬機器、並使用步驟 3 中建立的 VHD。
- b. 若要在實體主機上部署、請將 VHD 複製到實體電腦、並將其設定為從此新 VHD 開機。首先、掛載 VHD、執行 BCDboot e:\windows (VHD 掛載於 E:\ 下)、卸載 VHD、重新啟動實體電腦、然後開機至奈米伺服器。

5. 將 Nano Server 加入網域 (選用)：

- a. 登入網域中的任何電腦、並執行下列 PowerShell Cmdlet 來建立資料 blob：

```
$domain = "<input the domain to which the Nano Server is to be
joined>"
$nanoserver = "<input name of the Nano Server>"
```

```
djoin.exe /provision /domain $domain /machine $nanoserver /savefile
C:\temp\odjblob /reuse
.. 在遠端機器上執行下列 PowerShell Cmdlet 、將 odjblob
檔案複製到奈米伺服器：
```

```
$nanoserver = "<input name of the Nano Server>"
$nanouname = ""<input username of the Nano Server>"
$nanopwd = ""<input password of the Nano Server>"
```

```
$filePath = 'c:\temp\odjblob'
$fileContents = Get-Content -Path $filePath -Encoding Unicode
```

```
$securenanopwd = ConvertTo-SecureString -AsPlainText -Force $nanopwd
$nanosecuredcred = new-object management.automation.pscredential
$nanouname, $securenanopwd
```

```
Invoke-Command -VMName $nanoserver -Credential $nanosecuredcred
-ArgumentList @($filePath,$fileContents) -ScriptBlock \{
    param($filePath,$data)
    New-Item -ItemType directory -Path c:\temp
    Set-Content -Path $filePath -Value $data -Encoding Unicode
    cd C:\temp
    djoin /requestodj /loadfile c:\temp\odjblob /windowspath
    c:\windows /localos
}
```

b. 重新啟動奈米伺服器。

## 連線至奈米伺服器

若要使用 PowerShell 遠端連線至奈米伺服器、請完成下列步驟：

1. 在遠端伺服器上執行下列 Cmdlet 、將 Nano Server 新增為遠端電腦上的信任主機：

```
Set-Item WSMan:\LocalHost\Client\TrustedHosts "<input IP Address of the
Nano Server>"
```

． 如果環境安全、而且您想要將所有主機設定為伺服器上的信任主機、請執行下列命令：

```
Set-Item WSMan:\LocalHost\Client\TrustedHosts *
```

。在遠端伺服器上執行下列 Cmdlet 來啟動遠端工作階段。出現提示時、請提供 Nano 伺服器的密碼。

```
Enter-PSSession -ComputerName "<input IP Address of the Nano Server>"  
-Credential ~\Administrator
```

若要從遠端 Windows Server 使用 GUI 管理工具遠端連線至 Nano Server、請完成下列命令：

1. 以系統管理員群組成員的身分登入 Windows Server。
2. 啟動 Server Manager\*\*
3. 若要從伺服器管理員遠端管理奈米伺服器、請在「所有伺服器」上按一下滑鼠右鍵、按一下「新增伺服器」、提供奈米伺服器的資訊並加以新增。您現在可以在伺服器清單中看到奈米伺服器。選取「奈米伺服器」、按一下滑鼠右鍵、然後開始使用提供的各種選項進行管理。
4. 若要從遠端管理奈米伺服器上的服務、請完成下列步驟：
  - a. 從 Server Manager 的「工具」區段開啟「服務」。
  - b. 以滑鼠右鍵按一下 [ 服務 ( 本機 ) ]。
  - c. 按一下「連線至伺服器」
  - d. 提供 Nano Server 詳細資料、以檢視及管理奈米伺服器上的服務。
5. 如果在奈米伺服器上啟用 Hyper-V 角色、請完成下列步驟、從 Hyper-V Manager 遠端管理：
  - a. 從 Server Manager 的「工具」區段開啟 Hyper-V Manager。
  - b. 以滑鼠右鍵按一下 Hyper-V Manager。
  - c. 按一下「連線至伺服器」、並提供奈米伺服器詳細資料。現在、您可以將奈米伺服器當作 Hyper-V 伺服器來管理、在上面建立和管理 VM。
6. 如果在奈米伺服器上啟用容錯移轉叢集角色、請完成下列步驟、從容錯移轉叢集管理程式進行遠端管理：
  - a. 從 Server Manager 的「工具」區段開啟容錯移轉叢集管理程式。
  - b. 使用奈米伺服器執行叢集相關作業。

## 部署 Hyper-V 叢集

本附錄說明如何部署 Hyper-V 叢集。

### 先決條件

- 至少有兩部 Hyper-V 伺服器彼此連線。
- 每部 Hyper-V 伺服器上至少設定一個虛擬交換器。
- 容錯移轉叢集功能會在每部 Hyper-V 伺服器上啟用。
- SMB 共享區或 CSV 是用來儲存 VM 及其磁碟以用於 Hyper-V 叢集的共用儲存區。

- 不同叢集之間不應共用儲存設備。每個叢集只應有一個 CSV/CIFS 共用。
- 如果 SMB 共用是用作共用儲存設備、則必須設定 SMB 共用的權限、以授予叢集中所有 Hyper-V 伺服器的電腦帳戶存取權。

## 部署

1. 以系統管理員群組成員的身分登入其中一個 Windows Hyper-V 伺服器。
2. 啟動 Server Manager\*\*
3. 按一下 [ 工具 ] 區段中的 [ 容錯移轉叢集管理員 ] 。
4. 按一下「從動作建立叢集」功能表。
5. 提供屬於此叢集一部分的 Hyper-V 伺服器詳細資料。
6. 驗證叢集組態。當系統提示您進行叢集組態驗證時、請選取是、然後選取所需的測試、以驗證 Hyper-V 伺服器是否已通過先決條件、以成為叢集的一部分。
7. 驗證成功後、即會啟動「建立叢集」精靈。在精靈中、提供新叢集的叢集名稱和叢集 IP 位址。接著會為 Hyper-V 伺服器建立新的容錯移轉叢集。
8. 按一下容錯移轉叢集管理程式中新建立的叢集、然後加以管理。
9. 定義要使用的叢集共用儲存設備。可以是 SMB 共享或 CSV 。
10. 將 SMB 共用當成共用儲存設備不需要特殊步驟。
  - 在 NetApp 儲存控制器上設定 CIFS 共用。若要這麼做、請參閱「[在 SMB 環境中進行資源配置](#)」。
11. 若要使用 CSV 做為共用儲存設備、請完成下列步驟：
  - a. 在 NetApp 儲存控制器上設定 LUN 。若要這麼做、請參閱「在 SAN 環境中進行資源配置」一節。
  - b. 請確定容錯移轉叢集中的所有 Hyper-V 伺服器都能看到 NetApp LUN 。若要為容錯移轉叢集的所有 Hyper-V 伺服器執行此動作、請確定其啟動器已新增至 NetApp 儲存設備上的啟動器群組。此外、請務必探索 LUN 、並確定已啟用 MPIO 。
  - c. 在叢集中的任一 Hyper-V 伺服器上、完成下列步驟：
    - i. 將 LUN 連線、初始化磁碟、建立新的簡易磁碟區、並使用 NTFS 或 Refs 格式化。
    - ii. 在容錯移轉叢集管理程式中、展開叢集、展開儲存、在磁碟上按一下滑鼠右鍵、然後按一下新增磁碟。這樣會開啟將 LUN 顯示為磁碟的「將磁碟新增至叢集」精靈。按一下「確定」、將 LUN 新增為磁碟。
    - iii. 現在 LUN 命名為叢集磁碟、並顯示為磁碟下的可用儲存設備。
  - d. 在 LUN (叢集磁碟) 上按一下滑鼠右鍵、然後按一下「Add to Cluster Shared Volumes」現在 LUN 顯示為 CSV 。
  - e. CSV 可從容錯移轉叢集的所有 Hyper-V 伺服器同時顯示、並可從其本機位置 C : \ClusterStorage\ 存取。
12. 建立高可用度 VM :
  - a. 在容錯移轉叢集管理程式中、選取並展開先前建立的叢集。
  - b. 按一下 [ 角色 ] ，然後按一下 [ 動作 ] 中的 [ 虛擬機器 ] 按一下 [ 新增虛擬機器 ] 。
  - c. 從 VM 應位於的叢集中選取節點。
  - d. 在虛擬機器建立精靈中、提供共用儲存設備 ( SMB 共享區或 CSV ) 作為儲存 VM 及其磁碟的路徑。

- e. 使用 Hyper-V Manager 將共享儲存設備（SMB 共享或 CSV）設定為儲存 Hyper-V 伺服器的 VM 及其磁碟的預設路徑。
13. 測試計畫性容錯移轉。使用即時移轉、快速移轉或儲存移轉（移動）、將 VM 移至另一個節點。檢閱 ["叢集環境中的即時移轉"](#) 以取得更多詳細資料。
14. 測試非計畫性容錯移轉停止擁有 VM 的伺服器上的叢集服務。

## 在叢集式環境中部署 Hyper-V 線上即時移轉

本附錄說明如何在叢集式環境中部署即時移轉。

### 先決條件

若要部署即時移轉、您必須在具有共用儲存設備的容錯移轉叢集中設定 Hyper-V 伺服器。檢閱 ["部署 Hyper-V 叢集"](#) 以取得更多詳細資料。

### 部署

若要在叢集環境中使用即時移轉、請完成下列步驟：

1. 在容錯移轉叢集管理程式中、選取並展開叢集。如果叢集不可見、請按一下容錯移轉叢集管理程式、按一下連線至叢集、然後提供叢集名稱。
2. 按一下角色、其中會列出叢集中所有可用的 VM。
3. 在虛擬機器上按一下滑鼠右鍵、然後按一下「移動」。這樣做有三個選項：
  - \* 即時移轉。\* 您可以手動選取節點、或允許叢集選取最佳節點。在即時移轉中、叢集會將 VM 所使用的記憶體從目前節點複製到另一個節點。因此、當 VM 移轉至另一個節點時、VM 所需的記憶體和狀態資訊就已存在、可供 VM 使用。這種移轉方法幾乎是瞬間完成的、但一次只能有一個 VM 進行即時移轉。
  - \* 快速移轉。\* 您可以手動選取節點、或允許叢集選取最佳節點。在快速移轉中、叢集會將 VM 所使用的記憶體複製到儲存設備中的磁碟。因此、當 VM 移轉至另一個節點時、VM 所需的記憶體和狀態資訊可由另一個節點從磁碟中快速讀取。透過快速移轉、可同時移轉多個 VM。
  - \* 虛擬機器儲存移轉。\* 此方法使用「移動虛擬機器儲存設備」精靈。使用此精靈、您可以選取 VM 磁碟及其他要移至其他位置的檔案、這些位置可以是 CSV 或 SMB 共用。

## 在叢集式環境外部署 Hyper-V 即時移轉

本節說明如何在叢集環境外部署 Hyper-V 即時移轉。

### 先決條件

- 獨立式 Hyper-V 伺服器、具備獨立儲存設備或共享的 SMB 儲存設備。
- 安裝在來源伺服器和目的地伺服器上的 Hyper-V 角色。
- 這兩個 Hyper-V 伺服器都屬於同一個網域或彼此信任的網域。

### 部署

若要在非叢集式環境中執行即時移轉、請設定來源和目的地 Hyper-V 伺服器、讓它們能夠傳送和接收即時移轉作業。在兩部 Hyper-V 伺服器上、完成下列步驟：

1. 從 Server Manager 的「工具」區段開啟 Hyper-V Manager 。
2. 按一下 [ 動作 ] 中的 [Hyper-V 設定] 。
3. 按一下「即時移轉」、然後選取「啟用傳入和傳出即時移轉」。
4. 選擇是允許任何可用網路上的即時移轉流量、還是僅允許特定網路上的即時移轉流量。
5. 或者、您也可以從即時移轉的進階區段設定驗證傳輸協定和效能選項。
6. 如果使用 CredSSP 做為驗證傳輸協定、請務必先從目的地 Hyper-V 伺服器登入來源 Hyper-V 伺服器、然後再移動 VM 。
7. 如果使用 Kerberos 做為驗證傳輸協定、請設定受限制的委派。這樣做需要 Active Directory 網域控制站存取。若要設定委派、請完成下列步驟：
  - a. 以系統管理員身分登入 Active Directory 網域控制站。
  - b. 啟動 Server Manager 。
  - c. 按一下 [ 工具 ] 區段中的 [Active Directory 使用者和電腦] 。
  - d. 展開網域、然後按一下「電腦」。
  - e. 從清單中選取來源 Hyper-V 伺服器、以滑鼠右鍵按一下該伺服器、然後按一下「內容」。
  - f. 在 [ 委派 ] 索引標籤中，選取 [ 信任這台電腦只能委派給指定的服務 ] 。
  - g. 選取「僅使用 Kerberos」。
  - h. 按一下「新增」、即可開啟「新增服務」精靈。
  - i. 在 " 新增服務 " 中，按一下 " 使用者和電腦 "，這會開啟 " 選取使用者或電腦 "\*\*
  - j. 提供目的地 Hyper-V 伺服器名稱、然後按一下「確定」。
    - 若要移動 VM 儲存設備、請選取 CIFS 。
    - 若要移動 VM 、請選取 Microsoft Virtual System 移轉服務。
  - k. 在 [ 委派 ] 索引標籤中，按一下 [ 確定 ] 。
  - l. 從「電腦」資料夾中、從清單中選取目的地 Hyper-V 伺服器、然後重複此程序。在 [ 選取使用者或電腦 ] 中，提供來源 Hyper-V 伺服器名稱。
8. 移動 VM 。
  - a. 開啟 Hyper-V Manager 。
  - b. 在 VM 上按一下滑鼠右鍵、然後按一下「移動」。
  - c. 選擇 [ 移動虛擬機器 ] 。
  - d. 指定虛擬機器的目的地 Hyper-V 伺服器。
  - e. 選擇移動選項。若為共用即時移轉、請選擇僅移動虛擬機器。若為「共享無線即時移轉」、請根據您的偏好選擇其他兩個選項中的任何一個。
  - f. 根據您的偏好、在目的地 Hyper-V 伺服器上提供虛擬機器的位置。
  - g. 檢閱摘要、然後按一下「確定」來移動 VM 。

## 部署 Hyper-V Storage 即時移轉

### 瞭解如何設定 Hyper-V 儲存即時移轉

## 先決條件

- 您必須擁有獨立式 Hyper-V 伺服器、並具備獨立儲存設備（DAS 或 LUN）或 SMB 儲存設備（本機或其他 Hyper-V 伺服器共用）。
- Hyper-V 伺服器必須設定為即時移轉。檢閱中的部署一節 "[在叢集環境之外進行即時移轉](#)"。

## 部署

1. 開啟 Hyper-V Manager。
2. 在 VM 上按一下滑鼠右鍵、然後按一下「移動」。
3. 選取 [ 移動虛擬機器的儲存設備 ]。
4. 根據您的偏好選擇儲存設備的移動選項。
5. 提供 VM 項目的新位置。
6. 檢閱摘要、然後按一下「確定」以移動 VM 的儲存設備。

## 在叢集環境外部署 Hyper-V 複本

本附錄說明如何在叢集環境外部署 Hyper-V 複本。

## 先決條件

- 您需要位於相同或不同地理位置的獨立 Hyper-V 伺服器、作為主要和複本伺服器。
- 如果使用不同的站台、則必須設定每個站台的防火牆、以允許主要伺服器和複本伺服器之間的通訊。
- 複本伺服器必須有足夠的空間來儲存複寫的工作負載。

## 部署

1. 設定複本伺服器。
  - a. 為了讓傳入的防火牆規則允許傳入的複寫流量、請執行下列 PowerShell Cmdlet：

```
Enable-Netfirewallrule -displayname "Hyper-V Replica HTTP Listener (TCP-In)"
```

- .. 從 Server Manager 的「工具」區段開啟 Hyper-V Manager。
- .. 按一下「動作」中的「Hyper-V 設定」。
- .. 按一下 [ 複寫組態 ]，然後選取 [ 將此電腦啟用為複本伺服器 ]。
- .. 在驗證和連接埠區段中、選取驗證方法和連接埠。
- .. 在 [ 授權與儲存 ] 區段中，指定儲存複寫 VM 和檔案的位置。

2. 在主伺服器上啟用 VM 複寫。VM 複寫是以每個 VM 為基礎啟用、而非針對整個 Hyper-V 伺服器。
  - a. 在 Hyper-V Manager 中、以滑鼠右鍵按一下虛擬機器、然後按一下「啟用複寫」以開啟「啟用複寫」精靈。
  - b. 提供必須複寫 VM 的複本伺服器名稱。
  - c. 提供驗證類型和複本伺服器連接埠、該伺服器連接埠已設定為在複本伺服器上接收複寫流量。



- d. 選取要複寫的 VHD 。
- e. 選擇將變更傳送至複本伺服器的頻率（持續時間）。
- f. 設定還原點以指定複本伺服器上要維護的還原點數目。
- g. 選擇初始複寫方法、指定將 VM 資料初始複本傳輸到複本伺服器的方法。
- h. 檢閱摘要、然後按一下「完成」。
- i. 此程序會在複本伺服器上建立 VM 複本。

## 複寫

1. 執行測試容錯移轉、確保複本 VM 在複本伺服器上正常運作。測試會在複本伺服器上建立一個暫存 VM 。

  - a. 登入複本伺服器。
  - b. 在 Hyper-V Manager 中、以滑鼠右鍵按一下複本 VM 、按一下複寫、然後按一下測試容錯移轉。
  - c. 選擇要使用的恢復點。
  - d. 此程序會建立名稱相同的 VM 、並附加 -Test 。
  - e. 驗證虛擬機器、確保一切正常運作。
  - f. 容錯移轉之後、如果您為複本測試 VM 選取停止測試容錯移轉、則會刪除該虛擬機器。

2. 執行計畫性容錯移轉、將主要 VM 上的最新變更複寫到複本 VM 。

  - a. 登入主要伺服器。
  - b. 關閉要容錯移轉的 VM 。
  - c. 在 Hyper-V Manager 中、以滑鼠右鍵按一下關閉的 VM 、按一下「複寫」、然後按一下「規劃容錯移轉」。
  - d. 按一下容錯移轉、將最新的 VM 變更傳輸至複本伺服器。

3. 在主 VM 故障時執行非計畫性容錯移轉。

  - a. 登入複本伺服器。
  - b. 在 Hyper-V Manager 中、以滑鼠右鍵按一下複本 VM 、按一下複寫、然後按一下容錯移轉。
  - c. 選擇要使用的恢復點。
  - d. 按一下「容錯移轉」以容錯移轉 VM 。

## 在叢集環境中部署 Hyper-V 複本

瞭解如何使用 Windows Server 容錯移轉叢集來部署及設定 Hyper-V 複本。

### 先決條件

- 您必須將 Hyper-V 叢集放在相同或不同地理位置的同一位置、做為主要叢集和複本叢集。檢閱 ["部署 Hyper-V 叢集"](#) 以取得更多詳細資料。
- 如果使用不同的站台、則必須設定每個站台的防火牆、以允許主要叢集和複本叢集之間的通訊。
- 複本叢集必須有足夠的空間來儲存複寫的工作負載。

## 部署

1. 在叢集的所有節點上啟用防火牆規則。在主要叢集和複本叢集中的所有節點上、以管理員權限執行下列 PowerShell Cmdlet 。

```
# For Kerberos authentication
get-clusternode | ForEach-Object \{Invoke-command -computername $_.name
-scripblock \{Enable-Netfirewallrule -displayname "Hyper-V Replica HTTP
Listener (TCP-In)"}\}
```

```
# For Certificate authentication
get-clusternode | ForEach-Object \{Invoke-command -computername $_.name
-scripblock \{Enable-Netfirewallrule -displayname "Hyper-V Replica
HTTPS Listener (TCP-In)"}\}
```

2. 設定複本叢集。
  - a. 使用 NetBIOS 名稱和 IP 位址來設定 Hyper-V 複本代理程式、以作為用於複本叢集的叢集連線點。
    - i. 開啟容錯移轉叢集管理程式。
    - ii. 展開叢集、按一下「角色」、然後按一下「從動作設定角色」窗格。
    - iii. 在「選取角色」頁面中選取 Hyper-V 複本代理人。
    - iv. 提供要用作叢集連線點的 NetBIOS 名稱和 IP 位址（用戶端存取點）。
    - v. 此程序會建立 Hyper-V 複本代理程式角色。確認已成功上線。
  - b. 設定複寫設定。
    - i. 以滑鼠右鍵按一下在前述步驟中建立的複本代理程式、然後按一下複寫設定。
    - ii. 選取 " 將此叢集啟用為複本伺服器 "。
    - iii. 在驗證和連接埠區段中、選取驗證方法和連接埠。
    - iv. 在授權與儲存區段中、選取允許將 VM 複寫至此叢集的伺服器。此外、請指定儲存複寫虛擬機器的預設位置。

## 複寫

複寫類似於一節中所述的程序 ["叢集環境外部的複本"](#)。

## 何處可找到其他資訊

### Microsoft Windows 和 Hyper-V 的其他資源

- 概念ONTAP  
<https://docs.netapp.com/us-en/ontap/concepts/introducing-data-management-software-concept.html>
- 現代 SAN 的最佳實務做法  
<https://www.netapp.com/media/10680-tr4080.pdf>

- NetApp All SAN Array Data Availability and Integrity with the NetApp ASA  
<https://www.netapp.com/pdf.html?item=/media/85671-tr-4968.pdf>
- SMB 文件  
<https://docs.netapp.com/us-en/ontap/smb-admin/index.html>
- 開始使用 Nano Server  
<https://technet.microsoft.com/library/mt126167.aspx>
- Windows Server 上 Hyper-V 的新功能  
<https://technet.microsoft.com/windows-server-docs/compute/hyper-v/what-s-new-in-hyper-v-on-windows>

# Microsoft SQL Server

## ONTAP 上的 Microsoft SQL Server

ONTAP 為您的 Microsoft SQL Server 資料庫提供企業級的安全性與效能解決方案、同時也提供世界級的工具來管理您的環境。



本文件取代先前發佈的技術報告 \_TR-4590 : Microsoft SQL Server 與 ONTAP 的最佳實務做法指南

NetApp 假設讀者具備下列工作知識：

- 軟件ONTAP
- NetApp SnapCenter 做為備份軟體、包括：
  - 適用於Microsoft Windows的解決方案SnapCenter
  - 適用於 SQL Server 的 SnapCenter 外掛程式
- Microsoft SQL Server 架構與管理

本最佳實務做法一節的範圍僅限於根據 NetApp 建議的儲存基礎架構設計原則和偏好的標準進行技術設計。端點對端點實作超出範圍。

如需 NetApp 產品的組態相容性、請參閱 "[NetApp互通性對照表工具IMT \(不含\)](#)"。

### Microsoft SQL Server 工作負載

部署 SQL Server 之前、您必須先瞭解 SQL Server 執行個體所支援應用程式的資料庫工作負載需求。每個應用程式對於容量、效能和可用度的需求各不相同、因此每個資料庫都應該設計成能以最佳方式支援這些需求。許多組織會使用應用程式需求來定義 SLA、將資料庫分類為多個管理層。SQL Server 工作負載的說明如下：

- OLTP 資料庫通常也是組織中最重要的資料庫。這些資料庫通常會支援面對客戶的應用程式、因此被視為是公司核心營運不可或缺的一環。關鍵任務 OLTP 資料庫及其支援的應用程式通常都有需要高效能的 SLA、而且對效能降級和可用度很敏感。他們也可能是永遠在容錯移轉叢集或永遠在可用度群組的候選對象。這些類型資料庫的 I/O 組合通常以 75% 至 90% 隨機讀取和 25% 至 10% 寫入為特徵。
- 決策支援系統 (DSS) 資料庫也可稱為資料倉儲。這些資料庫在許多仰賴分析技術的企業中、都是關鍵任務。執行查詢時、這些資料庫會對 CPU 使用率和從磁碟讀取作業敏感。在許多組織中、DSS 資料庫在月底、每季和每年都是最重要的此工作負載通常有 100% 的讀取 I/O 混合。

## 資料庫組態

### Microsoft SQL Server CPU 組態

若要改善系統效能、您需要修改 SQL Server 設定和伺服器組態、才能使用適當數量的處理器來執行。

## 超執行緒

超執行緒是 Intel 專屬的同步多執行緒（SMT）實作、可改善在 x86 微處理器上執行的運算（多工）平行化。

使用超執行緒的硬體可讓邏輯超執行緒 CPU 在作業系統中顯示為實體 CPU。SQL Server 接著會看到作業系統所呈現的實體 CPU、並可使用超執行緒處理器。如此一來、可提高平行處理能力、進而提升效能。

此處的注意點是、每個 SQL Server 版本都有自己的運算能力限制。如需詳細資訊、請參閱依 SQL Server 版本計算容量限制。

SQL Server 授權有兩個選項。第一個稱為伺服器 + 用戶端存取授權（CAL）模式、第二個是每個處理器核心模式。雖然您可以使用伺服器 + CAL 策略存取 SQL Server 中所有可用的產品功能、但每個插槽的硬體限制為 20 個 CPU 核心。即使您的伺服器的 SQL Server Enterprise Edition + CAL 每個插槽有超過 20 個 CPU 核心、應用程式也無法在該執行個體上同時使用所有這些核心。

下圖顯示啟動後的 SQL Server 記錄訊息、指出核心限制的強制執行。

記錄項目指出 **SQL Server** 啟動後所使用的核心數。

```
2017-01-11 07:16:30.71 Server      Microsoft SQL Server 2016
(RTM) - 13.0.1601.5 (X64)
Apr 29 2016 23:23:58
Copyright (c) Microsoft Corporation
Enterprise Edition (64-bit) on Windows Server 2016
Datacenter 6.3 <X64> (Build 14393: )

2017-01-11 07:16:30.71 Server      UTC adjustment: -8:00
2017-01-11 07:16:30.71 Server      (c) Microsoft Corporation.
2017-01-11 07:16:30.71 Server      All rights reserved.
2017-01-11 07:16:30.71 Server      Server process ID is 10176.
2017-01-11 07:16:30.71 Server      System Manufacturer:
'FUJITSU', System Model: 'PRIMERGY RX2540 M1'.
2017-01-11 07:16:30.71 Server      Authentication mode is MIXED.
2017-01-11 07:16:30.71 Server      Logging SQL Server messages
in file 'C:\Program Files\Microsoft SQL Server
\MSSQL13.MSSQLSERVER\MSSQL\Log\ERRORLOG'.
2017-01-11 07:16:30.71 Server      The service account is 'SEA-
TM\FUJIA2R30$'. This is an informational message; no user action
is required.
2017-01-11 07:16:30.71 Server      Registry startup parameters:
-d C:\Program Files\Microsoft SQL Server
\MSSQL13.MSSQLSERVER\MSSQL\DATA\master.mdf
-e C:\Program Files\Microsoft SQL Server
\MSSQL13.MSSQLSERVER\MSSQL\Log\ERRORLOG
-l C:\Program Files\Microsoft SQL Server
\MSSQL13.MSSQLSERVER\MSSQL\DATA\mastlog.ldf
-T 3502
-T 834
2017-01-11 07:16:30.71 Server      Command Line Startup
Parameters:
-s "MSSQLSERVER"
2017-01-11 07:16:30.72 Server      SQL Server detected 2 sockets
with 18 cores per socket and 36 logical processors per socket,
72 total logical processors; using 40 logical processors based
on SQL Server licensing. This is an informational message; no
user action is required.
2017-01-11 07:16:30.72 Server      SQL Server is starting at
```

因此、若要使用所有 CPU、您應該使用每個處理器核心授權。如需 SQL Server 授權的詳細資訊、請參閱 ["SQL Server 2022：現代化的資料平台"](#)。

## CPU 親和性

除非您遇到效能問題、否則您不太可能需要變更處理器親和性預設值、但仍值得您瞭解它們是什麼、以及它們的運作方式。

SQL Server 可透過兩個選項來支援處理器關聯性：

- CPU 親和性遮罩
- 關聯性 I/O 遮罩

SQL Server 會使用作業系統提供的所有 CPU（如果選擇了每個處理器核心授權）。它會在所有 CPU 上建立排程器、以便為任何指定的工作負載充分利用資源。當多工處理時、伺服器上的作業系統或其他應用程式可以將處理執行緒從一個處理器切換至另一個處理器。SQL Server 是一種資源密集的應用程式、發生這種情況時、效能可能會受到影響。為了將影響降至最低、您可以設定處理器、以便將所有 SQL Server 負載導向預先選定的處理器群組。這是透過使用 CPU 關聯性遮罩來達成的。

關聯性 I/O 遮罩選項會將 SQL Server 磁碟 I/O 繫結至 CPU 子集。在 SQL Server OLTP 環境中、此延伸功能可提升 SQL Server 執行緒的效能、以執行 I/O 作業。

## 平行度上限（MAXDOP）

根據預設、如果選擇個別處理器核心授權、SQL Server 會在查詢執行期間使用所有可用的 CPU。

雖然這對大型查詢很有幫助、但可能會造成效能問題並限制並行處理。更好的方法是將平行度限制在單一 CPU 插槽中的實體核心數量。例如、在具有兩個實體 CPU 插槽、每個插槽有 12 個核心的伺服器上、無論超執行緒為何、MAXDOP 都應設為 12。MAXDOP 無法限制或規定要使用的 CPU。而是限制單一批次查詢可以使用的 CPU 數量。



\* NetApp 建議 \* 適用於資料倉儲等 DSS、從 50 開始使用 MAXDOP、並視需要探索調校或調校。進行變更時、請務必測量應用程式中的關鍵查詢。

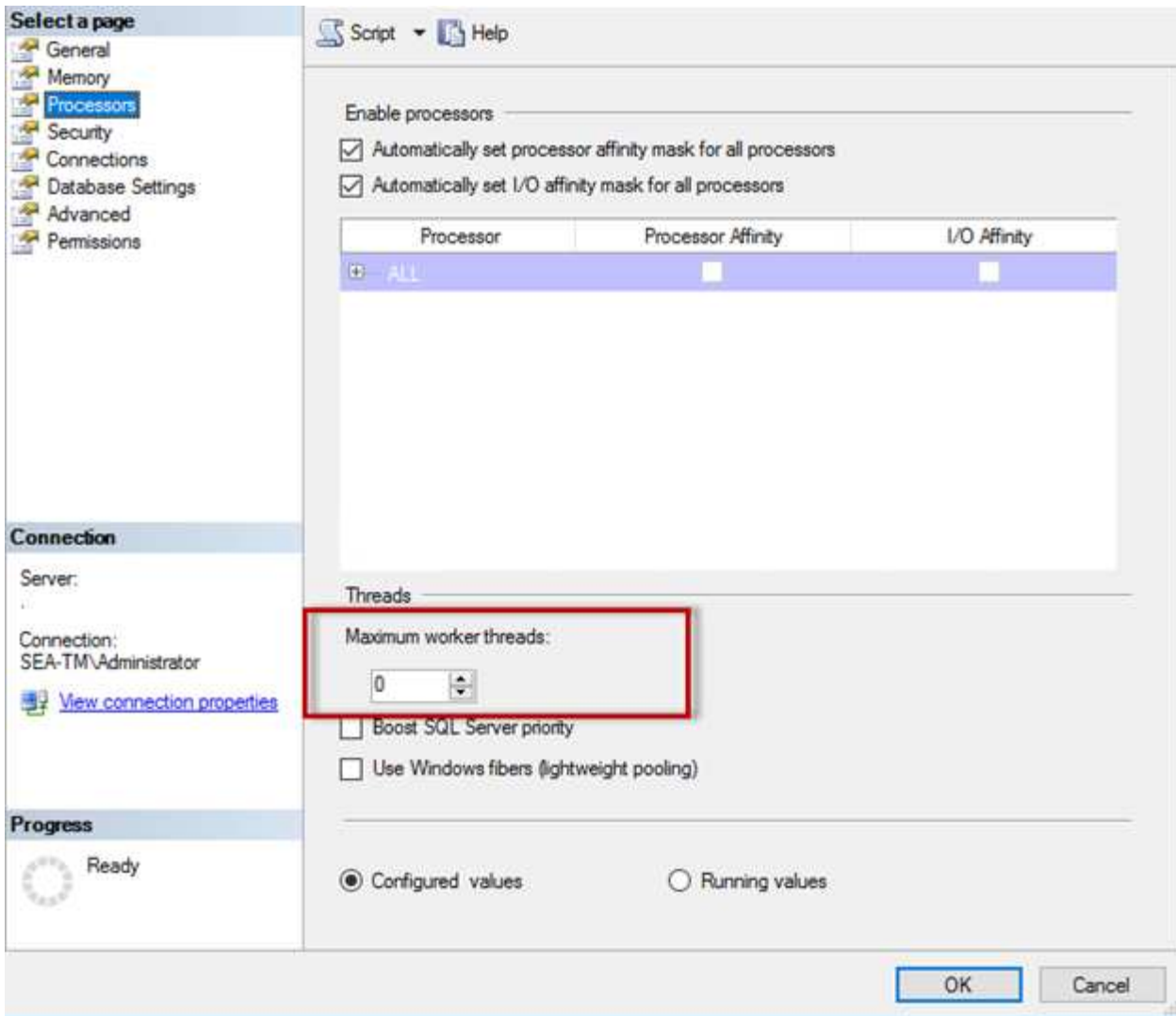
## 最大工作執行緒數

當大量用戶端連線至 SQL Server 時、最大工作執行緒選項可協助最佳化效能。

通常會為每個查詢要求建立個別的作業系統執行緒。如果數百個同時連線到 SQL Server、則每個查詢要求一個執行緒會消耗大量的系統資源。「最大工作執行緒」選項可讓 SQL Server 建立工作執行緒集區、以服務更多查詢要求、進而協助改善效能。

預設值為 0、可讓 SQL Server 在啟動時自動設定工作執行緒數量。這適用於大多數系統。Max Worker 執行緒是進階選項、如果沒有經驗豐富的資料庫管理員（DBA）的協助、就不應變更。

何時應設定 SQL Server 使用更多工作執行緒？如果每個排程器的平均工作佇列長度超過 1、您可能會因為在系統中新增更多執行緒而受益、但前提是負載未受 CPU 限制或遇到任何其他繁重的等待。如果發生上述任一種情況、新增更多執行緒並不會有幫助、因為它們最終會等待其他系統瓶頸。如需工作者執行緒上限的詳細資訊、請參閱 ["設定最大工作執行緒伺服器組態選項"](#)。



使用 SQL Server Management Studio 設定最大工作執行緒。

The following example shows how to configure the max work threads option using T-SQL.

```
EXEC sp_configure 'show advanced options', 1;
GO
RECONFIGURE ;
GO
EXEC sp_configure 'max worker threads', 900 ;
GO
RECONFIGURE;
GO
```

## Microsoft SQL Server 記憶體組態

下節說明如何設定 SQL Server 記憶體設定、以最佳化資料庫效能。

## 最大伺服器記憶體

最大伺服器記憶體選項可設定 SQL Server 執行個體可使用的最大記憶體容量。

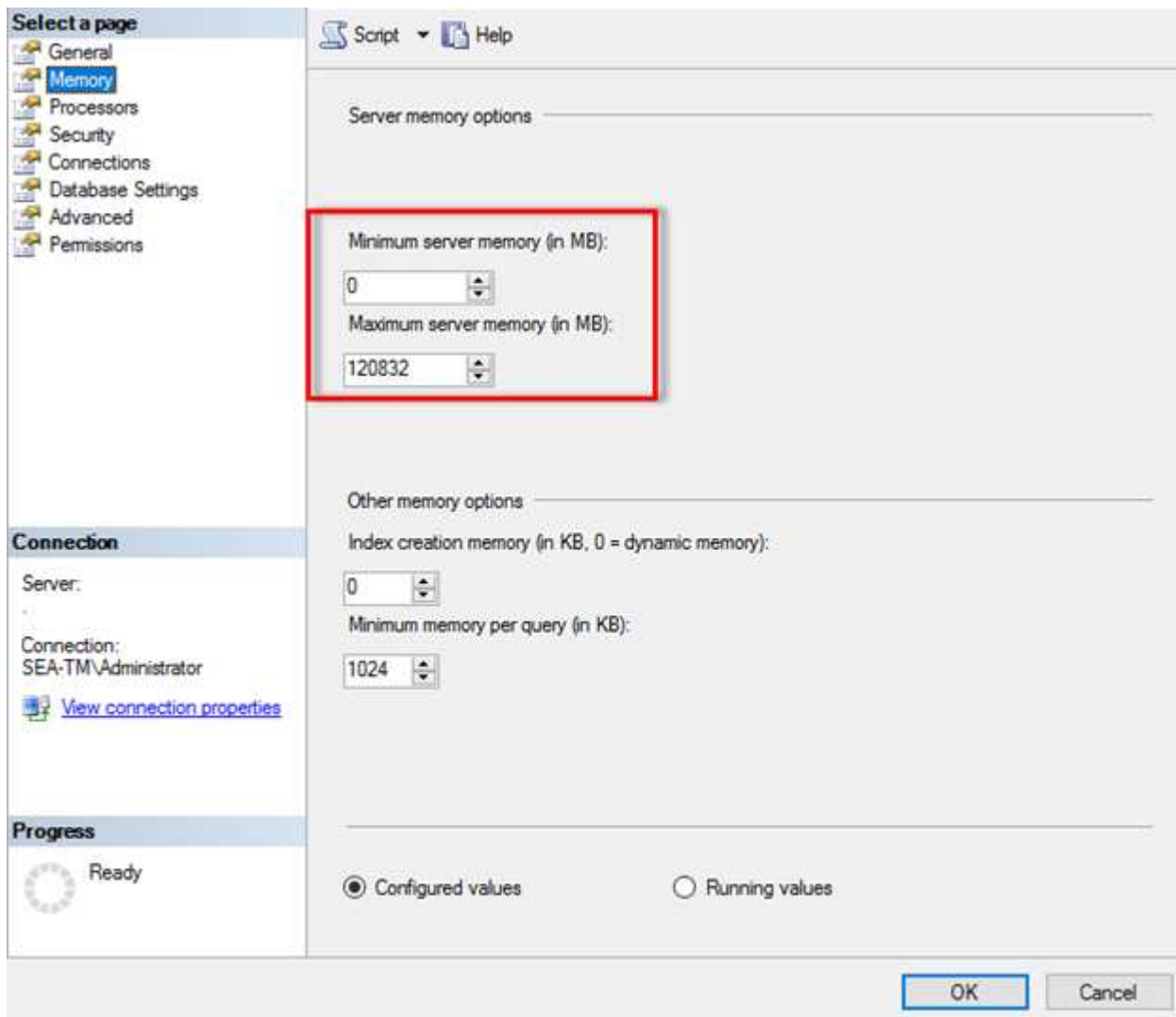
如果在執行 SQL Server 的同一部伺服器上執行多個應用程式、而且您想要保證這些應用程式有足夠的記憶體可以正常運作、通常會使用此功能。

有些應用程式只會在啟動時使用任何可用的記憶體、即使需要也不會要求更多。這就是最大伺服器記憶體設定的作用所在。

在具有多個 SQL Server 執行個體的 SQL Server 叢集上、每個執行個體可能會爭用資源。為每個 SQL Server 執行個體設定記憶體限制、有助於保證每個執行個體的最佳效能。



\* NetApp 建議 \* 為作業系統保留至少 4GB 至 6GB 的 RAM 、以避免效能問題。



使用 **SQL Server Management Studio** 調整最小和最大伺服器記憶體。

若要使用 SQL Server Management Studio 調整最小或最大伺服器記憶體、必須重新啟動 SQL Server 服務。您可以使用以下代碼、使用 TransAct SQL ( T-SQL ) 來調整伺服器記憶體：



```
EXECUTE sp_configure 'show advanced options', 1
GO
EXECUTE sp_configure 'min server memory (MB)', 2048
GO
EXEC sp_configure 'max server memory (MB)', 120832
GO
RECONFIGURE WITH OVERRIDE
```

## 不一致的記憶體存取

非一致性記憶體存取（NUMA）是一種記憶體存取最佳化方法、可協助提高處理器速度、而不會增加處理器匯流排的負載。

如果在安裝 SQL Server 的伺服器上設定 NUMA、則不需要額外的組態、因為 SQL Server 可感知 NUMA 並在 NUMA 硬體上執行良好。

## 索引會建立記憶體

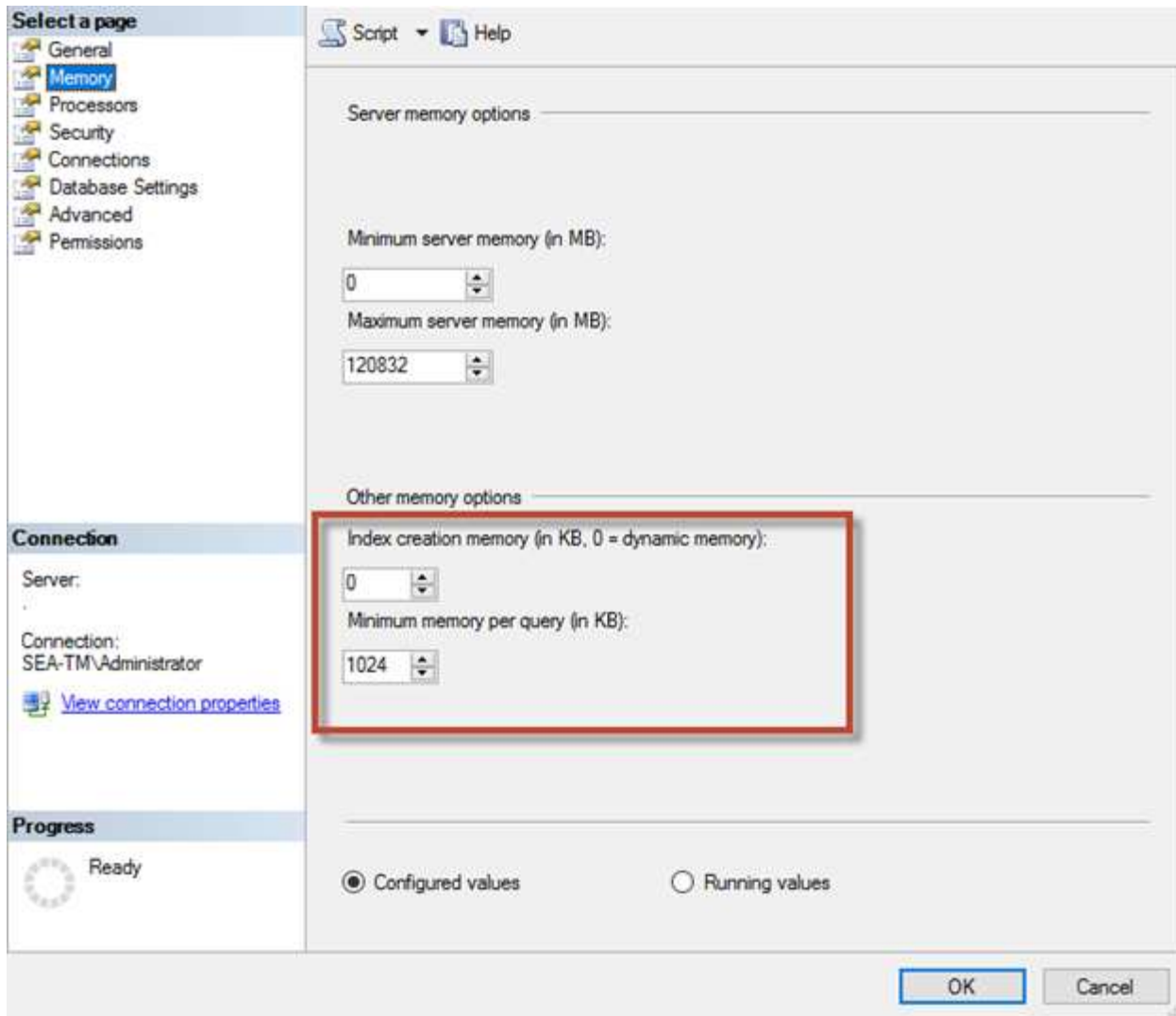
索引建立記憶體選項是另一個進階選項、您通常不應變更。

它控制最初分配給建立索引的最大 RAM 容量。此選項的預設值為 0、這表示它是由 SQL Server 自動管理。不過、如果您在建立索引時遇到困難、請考慮增加此選項的值。

## 每個查詢的最小記憶體

執行查詢時、SQL Server 會嘗試分配最適當的記憶體量、以有效執行。

根據預設、[每個查詢的最小記憶體] 設定會為每個要執行的查詢分配  $\geq$  至 1024KB。最佳做法是將此設定保留為預設值 0、以允許 SQL Server 動態管理分配給索引建立作業的記憶體量。不過、如果 SQL Server 的 RAM 超過有效執行所需的容量、則如果您增加此設定、某些查詢的效能可能會大幅提升。因此、只要 SQL Server、任何其他應用程式或作業系統未使用的伺服器上有可用的記憶體、則提升此設定有助於整體 SQL Server 效能。如果沒有可用的可用記憶體、增加此設定可能會影響整體效能。



## 緩衝區集區副檔名

緩衝區集區擴充可將 NVRAM 擴充與資料庫引擎緩衝區集區無縫整合、大幅改善 I/O 處理量。

並非每個 SQL Server 版本都提供緩衝區集區擴充功能。僅適用於 64 位元 SQL Server Standard、Business Intelligence 和 Enterprise 版本。

緩衝區集區擴充功能可利用非揮發性儲存設備（通常是 SSD）來擴充緩衝區集區快取。此擴充功能可讓緩衝區集區容納更大的資料庫工作集、強制在 RAM 和 SSD 之間分頁 I/O、並有效地將小型隨機 I/O 從機械磁碟卸載到 SSD。由於 SSD 延遲較低、隨機 I/O 效能較佳、因此緩衝區集區擴充功能可大幅改善 I/O 處理量。

緩衝區集區擴充功能提供下列優點：

- 增加隨機 I/O 處理量
- 降低 I/O 延遲
- 提高交易處理量
- 利用較大的混合式緩衝區集區來改善讀取效能
- 一種快取架構、可充分利用現有和未來的低成本記憶體



- NetApp 建議 \* 將緩衝區集區延伸設定為：
- 請確定將 SSD 支援的 LUN (例如 NetApp AFF) 呈現給 SQL Server 主機、以便將其用作緩衝區集區擴充目標磁碟。
- 副檔名必須與緩衝區集區大小相同或大於該檔案。

下列範例顯示 T-SQL 命令、可設定 32GB 的緩衝區集區擴充。

```
USE master
GO
ALTER SERVER CONFIGURATION
SET BUFFER POOL EXTENSION ON
(FILENAME = 'P:\BUFFER POOL EXTENSION\SQLServerCache.BUFFER POOL
EXTENSION', SIZE = 32 GB);
GO
```

## Microsoft SQL Server 共用執行個體與專用執行個體的比較

可將多個 SQL Server 設定為每部伺服器的單一執行個體、或設定為多個執行個體。正確的決策通常取決於各種因素、例如同伺服器是用於正式作業或開發、無論執行個體是否對業務營運和效能目標至關重要。

共享執行個體組態一開始可能比較容易設定、但可能會導致資源被分割或鎖定的問題、進而導致在共享 SQL Server 執行個體上主控資料庫的其他應用程式效能問題。

疑難排解效能問題可能很複雜、因為您必須找出哪個執行個體是根本原因。此問題與作業系統授權和 SQL Server 授權的成本相比較。如果應用程式效能至關重要、則強烈建議使用專用執行個體。

Microsoft 在伺服器層級的每個核心授權 SQL Server、而非每個執行個體授權。因此、資料庫管理員想要安裝伺服器能處理的 SQL Server 執行個體數量、以節省授權成本、這可能會導致日後發生重大效能問題。



- \* NetApp 建議 \* 盡可能選擇專屬的 SQL Server 執行個體、以獲得最佳效能。

## 儲存組態

### Microsoft SQL Server 儲存考量

結合 ONTAP 儲存解決方案與 Microsoft SQL Server、可建立企業級資料庫儲存設計、滿足現今最嚴苛的應用程式需求。

若要最佳化這兩種技術、瞭解 SQL Server I/O 模式和特性是非常重要的。精心設計的 SQL Server 資料庫儲存配置可支援 SQL Server 的效能及 SQL Server 基礎架構的管理。良好的儲存配置也能讓初始部署成功、並隨著業務成長、環境隨時間而順利成長。

## 資料儲存設計

對於不使用 SnapCenter 支援功能執行備份的 SQL Server 資料庫、Microsoft 建議將資料和記錄檔放在不同的磁碟機上。對於同時更新和要求資料的應用程式、記錄檔會密集寫入、而且資料檔（視應用程式而定）會密集讀寫。對於資料擷取、不需要記錄檔。因此、您可以從放在自己磁碟機上的資料檔案來滿足資料要求。

當您建立新資料庫時、Microsoft 建議您為資料和記錄指定個別的磁碟機。若要在資料庫建立之後移動檔案、資料庫必須離線。如需更多 Microsoft 建議、請參閱 "[將資料和記錄檔放在不同的磁碟機上](#)"。

## 集合體

Aggregate 是 NetApp 儲存組態的最低層級儲存容器。網際網路上存在一些舊版文件、建議將 IO 分隔到不同的基礎磁碟機集。ONTAP 不建議這麼做。NetApp 已使用資料檔案和交易記錄檔分離的共用和專用集合體、執行各種 I/O 工作負載特性分析測試。測試結果顯示、一個大型集合體含有更多 RAID 群組和磁碟機、可最佳化和改善儲存效能、而且管理員更容易管理、原因有兩個：

- 一個大型集合體可讓所有檔案都能使用所有磁碟機的 I/O 功能。
- 一個大型 Aggregate 可讓您以最有效率的方式使用磁碟空間。

若要獲得高可用度（HA）、請將 SQL Server Always On Availability Group 次要同步複本放在 Aggregate 中的獨立儲存虛擬機器（SVM）上。為了進行災難恢復、請將非同步複本放在 DR 站台中屬於獨立儲存叢集的集合上、並使用 NetApp SnapMirror 技術複寫內容。NetApp 建議在集合體中至少有 10% 的可用空間、以獲得最佳的儲存效能。

## 磁碟區

NetApp FlexVol 磁碟區會建立並位於集合體內。此術語有時會造成混淆、因為 ONTAP 磁碟區不是 LUN。ONTAP Volume 是資料的管理容器。磁碟區可能包含檔案、LUN、甚至 S3 物件。磁碟區不會佔用空間、只會用於管理內含的資料。

## Volume 設計考量

在您建立資料庫 Volume 設計之前、請務必瞭解 SQL Server I/O 模式和特性的差異、視工作負載及備份與還原需求而定。請參閱下列 NetApp 建議的彈性磁碟區：

- 避免在主機之間共用磁碟區。例如、雖然可以在單一磁碟區中建立 2 個 LUN、並將每個 LUN 共用至不同的主機、但這應該避免、因為這樣可能會使管理複雜化。
- 請使用 NTFS 掛載點而非磁碟機代號、以超越 Windows 中 26 個磁碟機代號的限制。使用 Volume 掛載點時、一般建議將 Volume 標籤命名為與掛載點相同的名稱。
- 適當時、請設定 Volume 自動調整大小原則、以協助避免空間不足的情況。17 Microsoft SQL Server with ONTAP 最佳實務指南 © 2022 NetApp、Inc. 版權所有。
- 如果您在 SMB 共用上安裝 SQL Server、請確定已在 SMB/CIFS 磁碟區上啟用 Unicode 以建立資料夾。
- 將磁碟區中的快照保留值設為零、以便從作業角度進行監控。
- 停用快照排程和保留原則。而是使用 SnapCenter 來協調 SQL Server 資料磁碟區的 Snapshot 複本。
- 將 SQL Server 系統資料庫放在專用磁碟區上。
- Tempdb 是 SQL Server 用來做為暫用工作區的系統資料庫、特別是用於 I/O 密集型 DBCC CHECKDB 作業。因此、請將此資料庫放在具有獨立磁碟集的專用磁碟區上。在磁碟區數是一項挑戰的大型環境中、您可以將 Tempdb 整合為較少的磁碟區、並在經過仔細規劃之後、將其儲存在與其他系統資料庫相同的磁碟區中。對 tempdb 的資料保護不是高優先順序、因為每次重新啟動 SQL Server 時都會重新建立此資料庫。

- 將使用者資料檔 (.mdf) 放在不同的磁碟區上、因為它們是隨機讀取/寫入工作負載。建立交易記錄備份的頻率通常高於資料庫備份。因此、請將交易記錄檔 (.ldf) 放在不同的磁碟區、或將 VMDK 放在資料檔案中、以便為每個檔案建立獨立的備份排程。這種分隔方式也能將記錄檔的連續寫入 I/O 與資料檔案的隨機讀寫 I/O 隔離、大幅提升 SQL Server 效能。

## LUN

- 請確定使用者資料庫檔案和用於儲存記錄備份的記錄目錄位於不同的磁碟區、以防止保留原則在與 SnapVault 技術搭配使用時覆寫快照。
- 請確定 SQL Server 資料庫位於與具有非資料庫檔案 (例如全文搜尋相關檔案) 的 LUN 分開的 LUN 上。
- 將資料庫次要檔案 (作為檔案群組的一部分) 放在不同的磁碟區、可改善 SQL Server 資料庫的效能。只有當資料庫的 .mdf 檔案未與任何其他 .mdf 檔案共用其 LUN 時、此分隔才有效。
- 如果您使用 DiskManager 或其他工具建立 LUN、請確保在格式化 LUN 時、將分割區的分配單元大小設為 64K。
- 請參閱 ["適用於現代 SAN 的 ONTAP 最佳實務做法下的 Microsoft Windows 和原生 MPIO"](#) 將 Windows 上的多重路徑支援套用至 MPIO 內容中的 iSCSI 裝置。

## Microsoft SQL Server 資料庫檔案和檔案群組

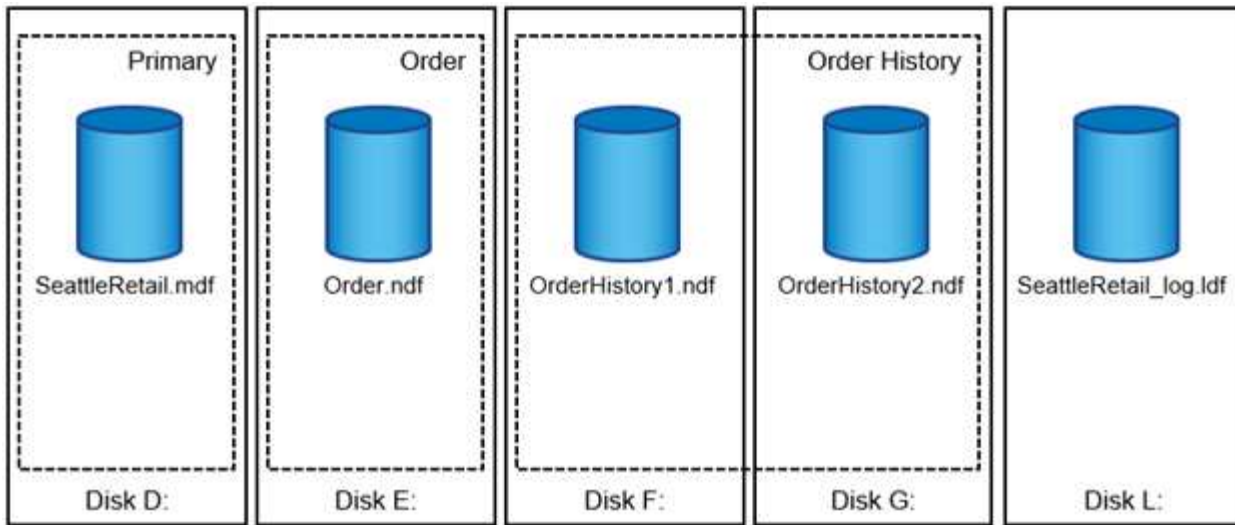
在 ONTAP 上正確放置 SQL Server 資料庫檔案是初始部署階段的關鍵。如此可確保最佳效能、空間管理、備份與還原時間、並可根據您的業務需求進行設定。

理論上、SQL Server (64 位元) 支援每個執行個體 32,767 個資料庫、以及 524272TB 的資料庫大小、雖然一般安裝通常有多個資料庫。不過 SQL Server 可以處理的資料庫數量取決於負載和硬體。看到 SQL Server 執行個體託管數十個、數百個甚至數千個小型資料庫、並不罕見。

每個資料庫都包含一或多個資料檔案、以及一或多個交易記錄檔。交易記錄會儲存資料庫交易的相關資訊、以及每個工作階段所做的所有資料修改。每次修改資料時、SQL Server 都會在交易記錄中儲存足夠的資訊、以復原 (復原) 或重做 (重新執行) 動作。SQL Server 交易記錄是 SQL Server 在資料完整性和健全性方面聲譽的重要一環。交易記錄對於 SQL Server 的原子性、一致性、隔離和耐用性 (ACID) 功能至關重要。一旦資料頁發生任何變更、SQL Server 就會立即寫入交易記錄檔。每個 Data 操縱語言 (DML) 陳述式 (例如、SELECT、INSERT、UPDATE 或 DELETE) 都是完整的交易、而且交易記錄會確保整個以 Set 為基礎的作業都能進行、確保交易的完整性。

每個資料庫都有一個主要資料檔案、預設會有 .mdf 副檔名。此外、每個資料庫都可以有次要資料庫檔案。根據預設、這些檔案的副檔名為 .NDF。

所有資料庫檔案都會分組為檔案群組。檔案群組是邏輯單元、可簡化資料庫管理。它們允許在邏輯物件放置和實體資料庫檔案之間進行分隔。當您建立資料庫物件表格時、您可以在檔案群組中指定它們應該放置的位置、而無需擔心基礎資料檔案組態。



將多個資料檔案放入檔案群組的功能可讓您將負載分散到不同的儲存裝置、有助於改善系統的 I/O 效能。與此相反的是、由於 SQL Server 會循序寫入交易記錄檔、因此無法從多個檔案中獲益。

檔案群組中的邏輯物件放置與實體資料庫檔案之間的分隔、可讓您微調資料庫檔案配置、充分發揮儲存子系統的效益。例如、將產品部署給不同客戶的獨立軟體廠商 (ISV)、可以根據基礎 I/O 組態和部署階段的預期資料量、調整資料庫檔案數量。這些變更對應用程式開發人員來說是透明的、他們將資料庫物件放置在檔案群組中、而非資料庫檔案中。



\* NetApp 建議 \* 避免將主要檔案群組用於系統物件以外的任何項目。為使用者物件建立個別的檔案群組或一組檔案群組、可簡化資料庫管理和災難恢復、尤其是大型資料庫。

您可以在建立資料庫或將新檔案新增至現有資料庫時指定初始檔案大小和自動成長參數。SQL Server 在選擇要將資料寫入哪個資料檔案時、會使用比例填滿演算法。它會將大量資料按比例寫入檔案中的可用空間。檔案中的可用空間越大、其處理的寫入次數就越多。



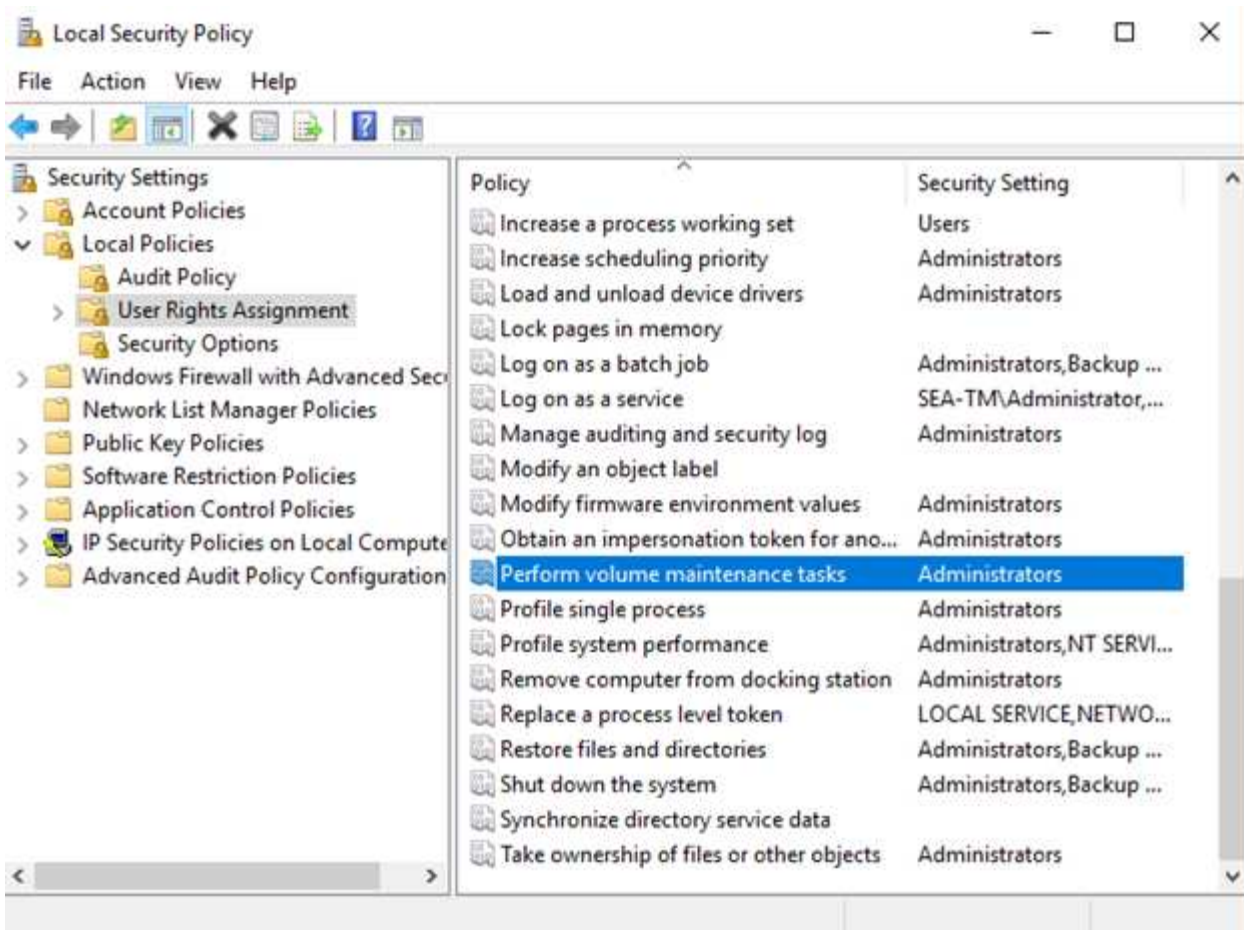
\* NetApp 建議 \* 單一檔案群組中的所有檔案都具有相同的初始大小和自動成長參數、其成長大小以 MB 為單位定義、而非百分比。這有助於比例填滿演算法在資料檔案之間平均平衡寫入活動。

每次 SQL Server 增加檔案時、都會以零填滿新分配的空間。該程序會封鎖所有需要寫入對應檔案的工作階段、或在交易記錄增加時產生交易記錄。

SQL Server 一律會將交易記錄檔歸零、而且該行為無法變更。不過、您可以啟用或停用即時檔案初始化來控制資料檔案是否正在歸零。啟用即時檔案初始化有助於加速資料檔案成長、並縮短建立或還原資料庫所需的時間。

與即時檔案初始化有關的安全風險較小。啟用此選項時、資料檔案的未分配部分可能會包含先前刪除的作業系統檔案資訊。資料庫管理員可以檢查這類資料。

您可以將 SA\_SA\_SAM\_VOLUM\_NAME 權限 (也稱為「執行 Volume 維護工作」) 新增至 SQL Server 啟動帳戶、以啟用即時檔案初始化。您可以在本機安全性原則管理應用程式 (secpol.msc) 下執行此動作、如下圖所示。開啟「執行 Volume 維護工作」權限的內容、並將 SQL Server 啟動帳戶新增至該處的使用者清單。



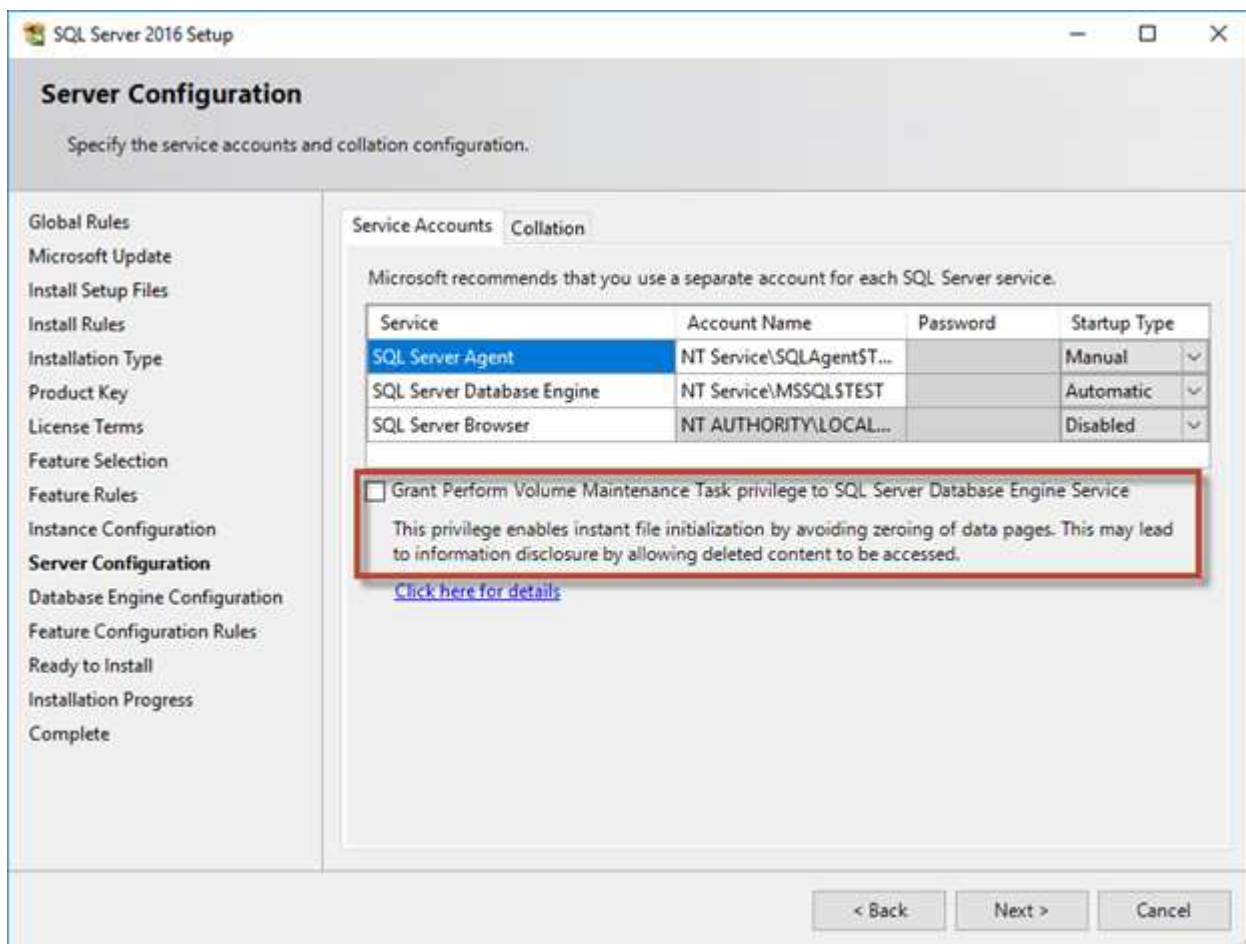
若要檢查權限是否已啟用、您可以使用下列範例中的程式碼。此程式碼會設定兩個追蹤旗標、強制 SQL Server 將其他資訊寫入錯誤記錄檔、建立小型資料庫、以及讀取記錄內容。

```
DBCC TRACEON(3004,3605,-1)
GO
CREATE DATABASE DelMe
GO
EXECUTE sp_readerrorlog
GO
DROP DATABASE DelMe
GO
DBCC TRACEOFF(3004,3605,-1)
GO
```

如果未啟用即時檔案初始化、SQL Server 錯誤記錄會顯示 SQL Server 除了將 ldf 記錄檔歸零之外、還會將 MDF 資料檔案歸零、如下例所示。當啟用即時檔案初始化時、它只會顯示記錄檔的零位。

	LogDate	ProcessInfo	Text
365	2017-02-09 08:10:07.660	spid53	Ckpt dbid 3 flush delta counts.
366	2017-02-09 08:10:07.660	spid53	Ckpt dbid 3 logging active xact info.
367	2017-02-09 08:10:07.750	spid53	Ckpt dbid 3 phase 1 ended (8)
368	2017-02-09 08:10:07.750	spid53	About to log Checkpoint end.
369	2017-02-09 08:10:07.880	spid53	Ckpt dbid 3 complete
370	2017-02-09 08:10:08.130	spid53	Starting up database 'DelMe'.
371	2017-02-09 08:10:08.150	spid53	FixupLog Tail(progress) zeroing C:\Program Files\Micros
372	2017-02-09 08:10:08.160	spid53	Zeroing C:\Program Files\Microsoft SQL Server\MSSQ
373	2017-02-09 08:10:08.170	spid53	Zeroing completed on C:\Program Files\Microsoft SQL
374	2017-02-09 08:10:08.710	spid53	Ckpt dbid 6 started
375	2017-02-09 08:10:08.710	spid53	About to log Checkpoint begin.

執行 Volume 維護工作在 SQL Server 2016 中已簡化、稍後會在安裝過程中提供選項。此圖顯示選項、可授予 SQL Server 資料庫引擎服務執行 Volume 維護工作的權限。



控制資料庫檔案大小的另一個重要資料庫選項是自動壓縮。啟用此選項時、SQL Server 會定期縮減資料庫檔案、減少檔案大小、並釋出空間給作業系統。這項作業需要大量資源、而且很少有用、因為當新資料進入系統時、資料庫檔案會在一段時間後再次增加。永遠不要在資料庫上啟用自動壓縮。

## Microsoft SQL Server 記錄目錄

記錄目錄是在 SQL Server 中指定、用於在主機層級儲存交易記錄備份資料。如果您使用



SnapCenter 來備份記錄檔、則 SnapCenter 使用的每個 SQL Server 主機都必須設定一個主機記錄目錄、才能執行記錄備份。由於包含資料庫儲存庫、因此與備份、還原或複製作業相關的中繼資料會儲存在中央資料庫儲存庫中 SnapCenter 。

主機記錄目錄的大小計算方式如下：

主機記錄目錄大小 = ( (最大 DB LDF 大小 x 每日記錄變更率 %) ) x (快照保留) ÷ ( 1 - LUN 額外負荷空間 % )

主機記錄目錄大小調整公式假設 LUN 負荷空間為 10%

將記錄目錄放在專用磁碟區或 LUN 上。主機記錄目錄中的資料量取決於備份的大小和保留備份的天數。SnapCenter 每個 SQL Server 主機只允許一個主機記錄目錄。您可以在 SnapCenter → 主機 → 組態外掛程式中設定主機記錄目錄。



- NetApp 建議 \* 下列主機記錄目錄：
- 請確定主機記錄目錄未被任何其他可能毀損備份快照資料的資料類型共用。
- 請勿將使用者資料庫或系統資料庫放置在裝載點的 LUN 上。
- 在 SnapCenter 複製交易記錄的專用 FlexVol 磁碟區上建立主機記錄目錄。
- 使用 SnapCenter 精靈將資料庫移轉至 NetApp 儲存設備、以便將資料庫儲存在有效位置、進而成功執行 SnapCenter 備份與還原作業。請記住、移轉程序會中斷運作、並可能導致資料庫在移轉進行中時離線。
- SQL Server 的容錯移轉叢集執行個體 ( FCI ) 必須符合下列條件：
  - 如果您使用容錯移轉叢集執行個體、則主機記錄目錄 LUN 必須是與要備份 SnapCenter 的 SQL Server 執行個體位於同一個叢集群組中的叢集磁碟資源。
  - 如果您使用容錯移轉叢集執行個體、則使用者資料庫必須放置在共用 LUN 上、這些 LUN 是指派給與 SQL Server 執行個體相關聯的叢集群組的實體磁碟叢集資源。

## Microsoft SQL Server tempdb 檔案

Tempdb 資料庫的使用率可能很高。除了在 ONTAP 上最佳放置使用者資料庫檔案之外、也可以變更 tempdb 資料檔案以減少分配爭用

當 SQL Server 必須寫入特殊系統頁面以分配新物件時、網頁爭用可能會發生在全域分配對應 ( GAM )、共用全域分配對應 ( SGAM ) 或頁面可用空間 ( PFS ) 頁面上。鎖條可保護 (鎖定) 記憶體中的這些頁面。在忙碌的 SQL Server 執行個體上、在 tempdb 的系統頁面上取得鎖定可能需要很長時間。這會導致查詢執行時間變慢、也稱為鎖定爭用。請參閱下列建立 tempdb 資料檔案的最佳實務做法：

- 對於 < 或 = 至 8 核心：tempdb 資料檔案 = 核心數
- 若為 > 8 核心：8 個 tempdb 資料檔案

下列範例指令碼會建立八個 tempdb 檔案、並將 tempdb 移至掛載點、以修改 tempdb C:\MSSQL\tempdb 適用於 SQL Server 2012 及更新版本。

```
use master  
  
go
```

```

-- Change logical tempdb file name first since SQL Server shipped with
logical file name called tempdev

alter database tempdb modify file (name = 'tempdev', newname =
'tempdev01');

-- Change location of tempdev01 and log file

alter database tempdb modify file (name = 'tempdev01', filename =
'C:\MSSQL\tempdb\tempdev01.mdf');

alter database tempdb modify file (name = 'templog', filename =
'C:\MSSQL\tempdb\templog.ldf');

GO

-- Assign proper size for tempdev01

ALTER DATABASE [tempdb] MODIFY FILE ( NAME = N'tempdev01', SIZE = 10GB );

ALTER DATABASE [tempdb] MODIFY FILE ( NAME = N'templog', SIZE = 10GB );

GO

-- Add more tempdb files

ALTER DATABASE [tempdb] ADD FILE ( NAME = N'tempdev02', FILENAME =
N'C:\MSSQL\tempdb\tempdev02.ndf' , SIZE = 10GB , FILEGROWTH = 10%);

ALTER DATABASE [tempdb] ADD FILE ( NAME = N'tempdev03', FILENAME =
N'C:\MSSQL\tempdb\tempdev03.ndf' , SIZE = 10GB , FILEGROWTH = 10%);

ALTER DATABASE [tempdb] ADD FILE ( NAME = N'tempdev04', FILENAME =
N'C:\MSSQL\tempdb\tempdev04.ndf' , SIZE = 10GB , FILEGROWTH = 10%);

ALTER DATABASE [tempdb] ADD FILE ( NAME = N'tempdev05', FILENAME =
N'C:\MSSQL\tempdb\tempdev05.ndf' , SIZE = 10GB , FILEGROWTH = 10%);

ALTER DATABASE [tempdb] ADD FILE ( NAME = N'tempdev06', FILENAME =
N'C:\MSSQL\tempdb\tempdev06.ndf' , SIZE = 10GB , FILEGROWTH = 10%);

ALTER DATABASE [tempdb] ADD FILE ( NAME = N'tempdev07', FILENAME =
N'C:\MSSQL\tempdb\tempdev07.ndf' , SIZE = 10GB , FILEGROWTH = 10%);

ALTER DATABASE [tempdb] ADD FILE ( NAME = N'tempdev08', FILENAME =

```

```
N'C:\MSSQL\tempdb\tempdev08.ndf' , SIZE = 10GB , FILEGROWTH = 10%);
```

```
GO
```

從 SQL Server 2016 開始、在安裝期間會自動偵測作業系統可見的 CPU 核心數量、並根據該數量、SQL Server 會計算和設定最佳效能所需的 tempdb 檔案數量。

## Microsoft SQL Server 與儲存效率

ONTAP 儲存效率經過最佳化、可儲存及管理 SQL Server 資料、但所需的儲存空間最少、對系統整體效能幾乎沒有影響或完全沒有影響。

儲存效率是 RAID、資源配置（整體配置和使用率）、鏡像和其他資料保護技術的組合。NetApp 技術包括快照、精簡配置和複製、可最佳化基礎架構中的現有儲存設備、並可延後或避免未來的儲存支出。這些技術搭配使用越多、節省的成本就越大。

空間效率功能（例如壓縮、壓縮和重複資料刪除）的設計、是為了增加符合特定實體儲存量的邏輯資料量。結果是降低成本和管理成本。

在高層級、壓縮是一種數學程序、可偵測及編碼資料模式、以減少空間需求。相反地、重複資料刪除功能會偵測實際重複的資料區塊、並移除額外的複本。資料實作可讓多個邏輯區塊在媒體上共用相同的實體區塊。



請參閱以下關於精簡配置的章節、以瞭解儲存效率與部分保留之間互動的說明。

### 壓縮

在提供 All Flash 儲存系統之前、以陣列為基礎的壓縮價值有限、因為大多數 I/O 密集的工作負載都需要大量磁碟來提供可接受的效能。儲存系統的容量總是比所需的容量大得多、這是大量磁碟機的副作用。固態儲存設備的興起、改變了這種情況。不再需要純粹為了獲得良好效能而大幅過度配置磁碟機。儲存系統中的磁碟機空間可與實際容量需求相符。

固態硬碟機（SSD）的 IOPS 容量增加、幾乎總是比旋轉硬碟機節省成本、但壓縮技術可以增加固態媒體的有效容量、進而進一步節省成本。

壓縮資料的方法有好幾種。許多資料庫都包含自己的壓縮功能、但在客戶環境中很少會發現這種情況。其原因通常是對壓縮資料的 \* 變更 \* 效能會受到影響、而對於某些應用程式而言、資料庫層級壓縮的授權成本較高。最後、對資料庫作業的整體效能影響。對於執行資料壓縮與解壓縮的 CPU、而非實際的資料庫工作、支付高昂的每 CPU 授權成本是不合理的。更好的選擇是將壓縮工作卸載到儲存系統。

### 自適應壓縮

即使在以微秒為單位測量延遲的 All Flash 環境中、主動式壓縮也已針對企業工作負載進行徹底測試、且未對效能產生任何影響。有些客戶甚至報告使用壓縮技術時效能會提高、因為資料會保持在快取中的壓縮、有效增加控制器中可用的快取數量。

ONTAP 以 4KB 單位管理實體區塊。自適應壓縮使用 8KB 的預設壓縮區塊大小、也就是以 8KB 為單位壓縮資料。這與關係式資料庫最常使用的 8KB 區塊大小相符。隨著將更多資料壓縮成單一單元、壓縮演算法就會變得更有效率。32 KB 壓縮區塊大小比 8 KB 壓縮區塊單元更具空間效率。這表示使用預設 8KB 區塊大小的調適式壓縮確實會導致效率稍微降低、但使用較小的壓縮區塊大小也有很大的好處。資料庫工作負載包括大量的覆寫活動。若要覆寫 32 KB 壓縮資料區塊的 8KB 資料、必須讀回整個 32 KB 的邏輯資料、將其解壓縮、更新所需的 8 KB 區域、重新壓縮、然後將整個 32 KB 寫入磁碟機。這對儲存系統來說是非常昂貴的作業、也是因為某些競爭

儲存陣列以較大的壓縮區塊大小為基礎、也會對資料庫工作負載造成重大效能損失的原因。



調適式壓縮所使用的區塊大小最多可增加至 32KB。這可能會改善儲存效率、而且當大量的這類資料儲存在陣列上時、應該考慮用於靜態檔案、例如交易記錄檔和備份檔案。在某些情況下、使用 16KB 或 32KB 區塊大小的作用中資料庫、也可能因為增加適應式壓縮的區塊大小而受惠。請洽詢 NetApp 或合作夥伴代表、瞭解這是否適合您的工作負載。



在串流備份目的地上、不應將大於 8KB 的壓縮區塊大小與重複資料刪除一起使用。原因是備份資料的細微變更會影響 32KB 壓縮時間。如果視窗移動、則產生的壓縮資料會在整個檔案中有所不同。重複資料刪除是在壓縮之後進行、這表示重複資料刪除引擎會以不同的方式檢視每個壓縮備份。如果需要重複資料刪除串流備份、則只應使用 8KB 區塊調適性壓縮。調適性壓縮較為理想、因為它的區塊大小較小、不會中斷重複資料刪除的效率。由於類似的原因、主機端壓縮也會影響重複資料刪除的效率。

#### 壓縮對齊

資料庫環境中的調適性壓縮需要考量壓縮區塊對齊。這樣做只是對隨機覆寫非常特定區塊的資料的考量。這種方法的概念與整體檔案系統對齊方式類似、檔案系統的開始必須與 4K 裝置邊界對齊、檔案系統的區塊大小必須是 4K 的倍數。

例如、只有在檔案與檔案系統本身的 8KB 邊界對齊時、才會壓縮寫入 8KB 檔案。這表示它必須落在檔案的前 8KB、檔案的第二 8KB 等。確保正確對齊的最簡單方法是使用正確的 LUN 類型、建立的任何分割區都應該與 8K 的倍數裝置開始偏移、並使用資料庫區塊大小的倍數檔案系統區塊大小。

備份或交易記錄等資料會循序寫入跨越多個區塊的作業、所有這些區塊都會被壓縮。因此、不需要考慮對齊。唯一令人擔憂的 I/O 模式是隨機覆寫檔案。

#### 資料壓縮

資料壓縮技術可改善壓縮效率。如前所述、僅有調適式壓縮功能、最多可節省 2 : 1、因為它僅限於在 4KB WAFL 區塊中儲存 8KB I/O。較大區塊大小的壓縮方法可提供更好的效率。不過、這些資料不適合受到小型區塊覆寫的資料。解壓縮 32KB 的資料單元、更新 8KB 部分、重新壓縮及回寫磁碟機、都會產生額外的負荷。

資料壓縮的運作方式是允許將多個邏輯區塊儲存在實體區塊內。例如、含有高度壓縮資料（例如文字或部分完整區塊）的資料庫、可能會從 8KB 壓縮至 1KB。如果沒有壓縮、1KB 的資料仍會佔用整個 4KB 區塊。即時資料壓縮功能可將 1KB 的壓縮資料與其他壓縮資料一起儲存在 1KB 的實體空間中。這不是一項壓縮技術、只是在磁碟機上分配空間的一種更有效率的方法、因此不應產生任何可偵測的效能影響。

節省的程度各不相同。已壓縮或加密的資料通常無法進一步壓縮、因此資料集無法從資料壓縮中獲益。相反地、新初始化的資料檔案僅包含區塊中繼資料和零、最多可壓縮至 80 : 1。

#### 對溫度敏感的儲存效率

溫度敏感儲存效率（TSSE）是 ONTAP 9.8 及更新版本中提供的產品、它仰賴區塊存取熱圖來識別不常存取的區塊、並以更高的效率加以壓縮。

#### 重複資料刪除

重複資料刪除是從資料集移除重複的區塊大小。例如、如果 10 個不同的檔案中存在相同的 4KB 區塊、重複資料刪除會將所有 10 個檔案中的 4KB 區塊重新導向至相同的 4KB 實體區塊。結果是該資料的效率提升 10 : 1。

VMware 來賓開機 LUN 等資料通常會極好地刪除重複資料、因為這些資料包含相同作業系統檔案的多個複本。效率達到 100 : 1 以上。

部分資料不包含重複資料。例如、Oracle 區塊包含資料庫的全域唯一標頭、以及近乎唯一的標尾。因此、Oracle 資料庫的重複資料刪除功能很少能節省 1% 以上的成本。使用 MS SQL 資料庫進行重複資料刪除的效果稍微好一些、但區塊層級的獨特中繼資料仍是一項限制。

在某些情況下、使用 16KB 和大型區塊大小的資料庫可節省高達 15% 的空間。每個區塊的初始 4KB 包含全域唯一的標頭、最後 4KB 區塊則包含近乎獨特的標尾。內部區塊是重複資料刪除的候選項目、但實際上、這幾乎完全歸功於重複資料刪除零位資料。

許多競爭陣列都宣稱能夠根據資料庫複製多次的假設來刪除重複的資料庫。在這方面，也可以使用 NetApp 重複資料刪除技術，但 ONTAP 提供更好的選擇：NetApp FlexClone 技術。最終結果相同；資料庫的多個複本會建立共用大部分基礎實體區塊。使用 FlexClone 比花時間複製資料庫檔案然後刪除複製檔案更有效率。實際上，它是不重複數據刪除，而不是重複數據刪除，因為從一開始就不會創建重複數據。

### 效率與精簡配置

效率功能是精簡配置的形式。例如、佔用 100GB 磁碟區的 100GB LUN 可能會壓縮至 50GB。由於磁碟區仍為 100GB、因此尚未實現實際節省。必須先縮小磁碟區的大小、才能將儲存的空間用於系統的其他位置。如果稍後變更為 100GB LUN、則資料的壓縮性會降低、LUN 的大小會增加、而且磁碟區可能會填滿。

強烈建議採用精簡配置、因為它可以簡化管理、同時大幅改善可用容量、並節省相關成本。原因很簡單：資料庫環境通常包含大量的空空間、大量的磁碟區和 LUN、以及可壓縮的資料。如果磁碟區和 LUN 的儲存空間有一天 100% 滿、而且包含 100% 不可壓縮的資料、則大量資源配置會導致保留空間。這種情況不太可能發生。精簡配置可回收空間並在其他地方使用、並可讓容量管理以儲存系統本身為基礎、而非許多較小的磁碟區和 LUN。

有些客戶偏好針對特定工作負載使用完整資源配置、或是根據既定的營運和採購實務做法。

- 注意：\* 如果磁碟區是完整配置的磁碟區、則必須小心將該磁碟區的所有效率功能完全停用、包括使用解壓縮和移除重複資料刪除 `sis undo` 命令。Volume 不應出現在 `volume efficiency show` 輸出。如果有、則磁碟區仍會部分設定為使用效率功能。因此、覆寫保證會以不同的方式運作、這會增加組態超視導致磁碟區意外用盡空間的機會、進而導致資料庫 I/O 錯誤。

### 效率最佳實務做法

NetApp 建議：

#### AFF 預設值

在 All Flash AFF 系統上執行的 ONTAP 上建立的磁碟區會自動精簡佈建、並啟用所有的內嵌效率功能。雖然資料庫通常無法從重複資料刪除中獲益、而且可能包含不可壓縮的資料、但預設設定仍適用於幾乎所有的工作負載。ONTAP 旨在有效處理所有類型的資料和 I/O 模式、無論是否能節省成本。只有在充分瞭解理由且有偏離的好處時、才應變更預設值。

#### 一般建議

- 如果磁碟區和（或）LUN 並未精簡配置、您必須停用所有效率設定、因為使用這些功能並不會節省成本、而將複雜資源配置與啟用空間效率的組合、可能會導致非預期的行為、包括空間不足的錯誤。
- 如果資料不需要覆寫、例如備份或資料庫交易記錄檔、您可以在冷卻週期較短的情況下啟用 TSSE、以達到更高的效率。

- 某些檔案可能包含大量不可壓縮的資料、例如、當檔案的應用程式層級已啟用壓縮時、就會進行加密。如果上述任何情況屬實、請考慮停用壓縮、以便在包含可壓縮資料的其他磁碟區上執行更有效率的作業。
- 請勿將 32KB 壓縮和重複資料刪除同時用於資料庫備份。請參閱一節 [\[自適應壓縮\]](#) 以取得詳細資料。

## 資料庫壓縮

SQL Server 本身也具備可壓縮及有效管理資料的功能。SQL Server 目前支援兩種類型的資料壓縮：資料列壓縮和頁面壓縮。

資料列壓縮會變更資料儲存格式。例如、它會將整數和小數位數變更為可變長度格式、而非原生固定長度格式。它也會消除空白、將固定長度字元字串變更為可變長度格式。頁面壓縮會實作列壓縮及其他兩種壓縮策略（前置壓縮和字典壓縮）。您可以在中找到有關頁面壓縮的詳細資料 "[頁面壓縮實作](#)"。

SQL Server 2008 及更新版本的 Enterprise、Developer 及 Evaluation 版本目前支援資料壓縮。雖然壓縮可以由資料庫本身執行、但在 SQL Server 環境中很少會發生這種情況。

以下是管理 SQL Server 資料檔案空間的建議

- 在 SQL Server 環境中使用自動精簡配置、以提高空間使用率、並在使用空間保證功能時、降低整體儲存需求。
- 對於大多數常見的部署組態、請使用自動擴充、因為儲存管理員只需要監控集合體中的空間使用量。
- 建議不要在包含 SQL Server 資料檔案的任何磁碟區上啟用重複資料刪除功能、除非已知該磁碟區包含相同資料的多個複本、例如將資料庫從備份還原至單一磁碟區。

## 空間回收

空間回收可定期啟動、以恢復 LUN 中未使用的空間。有了 SnapCenter、您可以使用下列 PowerShell 命令來啟動空間回收。

```
Invoke-SdHostVolumeSpaceReclaim -Path drive_path
```

如果您需要執行空間回收、則此程序應在低活動期間執行、因為它最初會在主機上使用週期。

## 使用 NetApp 管理軟體保護 Microsoft SQL Server 資料

規劃資料庫備份是根據業務需求而定。透過結合 ONTAP 的 NetApp Snapshot 技術並運用 Microsoft SQL Server API、無論使用者資料庫的大小為何、您都能快速進行應用程式一致的備份。如需更進階或橫向擴充的資料管理需求、NetApp 提供 SnapCenter。

### SnapCenter

SnapCenter 是適用於企業應用程式的 NetApp 資料保護軟體。使用適用於 SQL Server 的 SnapCenter 外掛程式、以及由 SnapCenter Plug-in for Microsoft Windows 管理的作業系統作業、即可快速輕鬆地保護 SQL Server 資料庫。

SQL Server 執行個體可以是獨立式安裝、容錯移轉叢集執行個體、也可以永遠位於可用性群組。結果是、從單一窗口即可保護、複製及還原主要或次要複本的資料庫。SnapCenter 可以同時管理內部部署、雲端和混合式組態的 SQL Server 資料庫。資料庫複本也可以在幾分鐘內在原始主機或替代主機上建立、以供開發或報告之用。



\* NetApp 建議 \* 使用 SnapCenter 來建立 Snapshot 複本。以下所述的 T-SQL 方法也能正常運作、但 SnapCenter 提供完整的備份、還原及複製程序自動化功能。它也會執行探索、以確保建立正確的快照。不需要預先設定。

...

SQL Server 也需要作業系統與儲存設備之間的協調、以確保建立時快照中有正確的資料。在大多數情況下、唯一安全的方法是使用 SnapCenter 或 T-SQL。在沒有這種額外協調的情況下建立的快照、可能無法可靠地恢復。

如需適用於 SnapCenter 的 SQL Server 外掛程式的詳細資訊、請參閱 ["TR-4714：使用 NetApp SnapCenter 的 SQL Server 最佳實務指南"](#)。

## 使用 T-SQL 快照保護資料庫

在 SQL Server 2022 中、Microsoft 推出了 T-SQL 快照、可提供建立指令碼和自動化備份作業的路徑。您可以為資料庫準備快照、而非執行完整大小的複本。資料庫準備好備份之後、您就可以利用 ONTAP REST API 來建立快照。

以下是備份工作流範例：

1. 使用 ALTER 命令凍結資料庫。如此一來、資料庫就能在基礎儲存設備上準備一致的快照。凍結之後、您可以使用備份命令來解凍資料庫並記錄快照。
2. 使用新的備份群組和備份伺服器命令、同時在儲存磁碟區上執行多個資料庫的快照。
3. 執行完整備份或僅複製完整備份。這些備份也會記錄在 msdb 中。
4. 使用快照完整備份後以正常串流方式進行的記錄備份、執行時間點還原。如果需要、也支援串流差異備份。

若要深入瞭解、請參閱 ["瞭解 T-SQL 快照的 Microsoft 文件"](#)。

## 使用 ONTAP 進行 Microsoft SQL Server 災難恢復

企業資料庫和應用程式基礎架構通常需要複寫、才能在最短的停機時間內、避免自然災難或非預期的業務中斷。

SQL Server 全年無休可用性群組複寫功能是絕佳的選擇、而 NetApp 提供的選項可將資料保護與全年無休整合。不過、在某些情況下、您可能會想要考慮使用 ONTAP 複寫技術。ONTAP 複寫選項（包括 MetroCluster 和 SnapMirror）可在最小化效能影響的情況下更好地擴充、保護非 SQL 資料、並通常提供完整的基礎架構複寫和災難恢復解決方案。

### SnapMirror 非同步

SnapMirror 技術提供快速靈活的非同步企業解決方案、可在 LAN 和 WAN 上複寫資料。SnapMirror 技術只會在建立初始鏡像之後、將變更的資料區塊傳輸到目的地、大幅降低網路頻寬需求。

以下是 SnapMirror for SQL Server 的建議：

- 如果使用 CIFS、目的地 SVM 必須是來源 SVM 所屬 Active Directory 網域的成員、如此一來、在從災難恢復期間、NAS 檔案中儲存的存取控制清單（ACL）就不會中斷。
- 不需要使用與來源 Volume 名稱相同的目的地 Volume 名稱、但可讓將目的地 Volume 掛載至目的地的程序更容易管理。如果使用 CIFS、您必須在來源命名空間的路徑和目錄結構中、使目的地 NAS 命名空間相同。

- 為了一致性的目的、請勿從控制器排程 SnapMirror 更新。而是在完整備份或記錄備份完成後、從 SnapCenter 啟用 SnapMirror 更新以更新 SnapMirror。
- 將包含 SQL Server 資料的磁碟區分散到叢集中的不同節點、以允許所有叢集節點共用 SnapMirror 複寫活動。此套裝作業系統可最佳化節點資源的使用。

如需 SnapMirror 的詳細資訊、請參閱 ["TR-4015：ONTAP 9 的 SnapMirror 組態與最佳實務做法指南"](#)。

## 保護 ONTAP 上的 Microsoft SQL Server 安全

確保 SQL Server 資料庫環境的安全性是一項不只是管理資料庫本身的多維工作。ONTAP 提供數項獨特功能、旨在保護資料庫基礎架構的儲存層面。

### Snapshot 複本

儲存快照是目標資料的時間點複本。ONTAP 的實作功能包括設定各種原則、每個磁碟區最多可儲存 1024 個快照。ONTAP 中的快照具有極高的空間效率。空間只會在原始資料集變更時使用。它們也是唯讀的。快照可以刪除、但無法變更。

在某些情況下、可以直接在 ONTAP 上排程快照。在其他情況下、在建立快照之前、可能需要 SnapCenter 等軟體來協調應用程式或作業系統的作業。無論哪種方法最適合您的工作負載、積極的快照策略都能透過頻繁且容易存取的方式、來提供資料安全性、從開機 LUN 到關鍵任務資料庫、都能備份所有資料。

- 注意 \*：ONTAP 彈性 Volume（或更簡單的）磁碟區與 LUN 並不同。Volume 是管理容器、用於儲存檔案或 LUN 等資料。例如、資料庫可能放置在 8-LUN 等量磁碟區集上、而所有 LUN 都包含在單一磁碟區中。

如需快照的詳細資訊、請按一下 ["請按這裡。"](#)

### 防竄改快照

從 ONTAP 9.12.1 開始、快照不只是唯讀的、也可以防止意外或刻意刪除。此功能稱為防竄改快照。您可以透過快照原則設定及強制執行保留期間。產生的快照必須等到到期日才會刪除。沒有系統管理或支援中心覆寫。

如此可確保入侵者、惡意內部人員、甚至勒索軟體攻擊都無法入侵備份、即使備份導致存取 ONTAP 系統本身也是如此。如果結合頻繁的快照排程、就能以極低的 RPO 提供極強大的資料保護功能。

如需防竄改快照的詳細資訊、請按一下 ["請按這裡。"](#)

### SnapMirror 複寫

快照也可以複寫到遠端系統。這包括防竄改快照、在遠端系統上套用及強制執行保留期間。因此資料保護效益與本機快照相同、但資料位於第二個儲存陣列上。如此可確保原始陣列的毀損不會影響備份。

第二個系統也會開啟新的管理安全選項。例如、某些 NetApp 客戶會將主要和次要儲存系統的驗證認證資料加以分隔。沒有單一管理使用者可以存取這兩個系統、這表示惡意系統管理員無法刪除所有資料複本。

如需 SnapMirror 的詳細資訊、請按一下 ["請按這裡。"](#)



## 儲存虛擬機器

新設定的 ONTAP 儲存系統類似於新佈建的 VMware ESX 伺服器、因為在建立虛擬機器之前、兩者都無法支援任何使用者。透過 ONTAP、您可以建立儲存虛擬機器（SVM）、成為最基本的儲存管理單元。每個 SVM 都有自己的儲存資源、傳輸協定組態、IP 位址和 FCP WWN。這是 ONTAP 多租戶的基礎。

例如、您可以為關鍵的正式作業工作負載設定一個 SVM、並在不同的網路區段上設定第二個 SVM 以進行開發活動。然後、您可以限制特定管理員存取正式作業 SVM、同時讓開發人員更廣泛地控制開發 SVM 中的儲存資源。您可能也需要為財務和人力資源團隊提供第三個 SVM、以儲存特別重要的純眼資料。

如需有關 SVM 的詳細資訊、請按一下 ["請按這裡。"](#)

## 管理 RBAC

ONTAP 提供強大的角色型存取控制（RBAC）功能、可用於管理登入。某些管理員可能需要完整的叢集存取權、而其他管理員可能只需要存取特定的 SVM。進階服務台人員可能需要增加磁碟區大小的能力。因此、您可以授予系統管理使用者執行工作職責所需的存取權限、而無需再提供其他權限。此外、您可以使用來自不同廠商的 PKI 來保護這些登入安全、僅限制對 ssh 金鑰的存取、並強制執行失敗的登入嘗試鎖定。

如需管理存取控制的詳細資訊、請按一下 ["請按這裡。"](#)

## Multifactor 驗證

ONTAP 和其他某些 NetApp 產品現在支援使用各種方法的多因素驗證（MFA）。因此、光是使用者名稱 / 密碼就不再是安全執行緒、而沒有第二個因素的資料、例如 FOB 或智慧型手機應用程式。

如需詳細資訊、請按一下 ["請按這裡。"](#)

## API RBAC

自動化需要 API 呼叫、但並非所有工具都需要完整的管理存取權。為了協助保護自動化系統的安全、API 層級也提供 RBAC。您可以將自動化使用者帳戶限制為所需的 API 呼叫。例如、監控軟體不需要變更存取權、只需要讀取存取權。配置儲存設備的工作流程不需要刪除儲存設備的能力。

若要深入瞭解、請啟動 [here](#).

## 多重管理驗證（MAV）

若要進一步進行多重「因素」驗證、需要兩位不同的管理員（每位管理員都有自己的認證）來核准某些活動。這包括變更登入權限、執行診斷命令和刪除資料。

如需多管理驗證（MAV）的詳細資訊、請按一下 ["請按這裡"](#)

# MySQL

## ONTAP 上的 MySQL 資料庫

MySQL 及其變種、包括 MariaDB 和 Percona MySQL、是全球最受歡迎的資料庫。



ONTAP 和 MySQL 資料庫上的本文件取代先前發佈的 [\\_TR-4722](#)：ONTAP 最佳實務做法的 MySQL 資料庫。

ONTAP 是適用於 MySQL 資料庫的理想平台、因為 ONTAP 確實是專為資料庫所設計。為了滿足資料庫工作負載的需求、我們特別建立了許多功能、例如隨機 IO 延遲最佳化、以提供進階服務品質（QoS）到基本 FlexClone 功能。

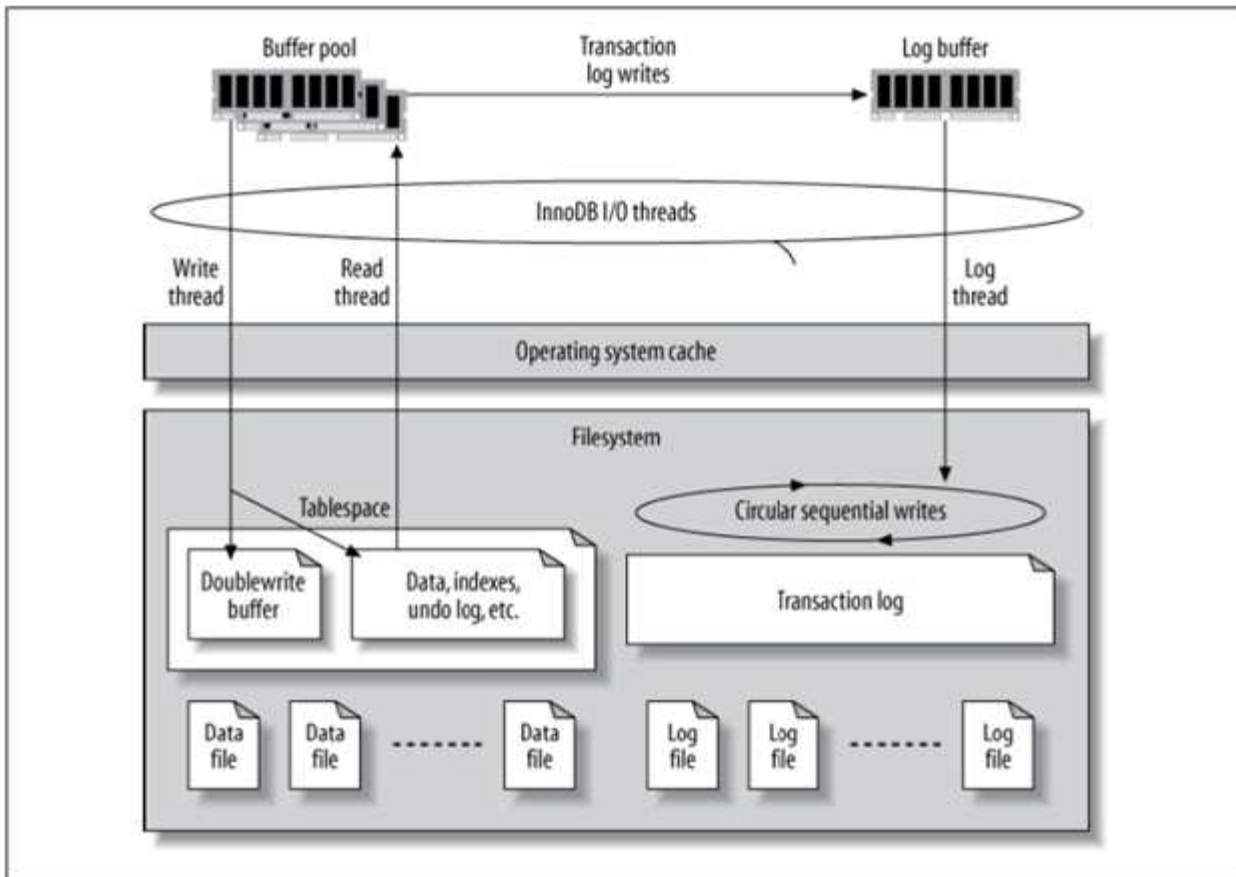
其他功能（例如不中斷升級）（包括儲存設備更換）、可確保關鍵資料庫仍可使用。您也可以透過 MetroCluster 為大型環境進行即時災難恢復、或是使用 SnapMirror 主動式同步來選擇資料庫。

最重要的是、ONTAP 提供無與倫比的效能、並能根據您的獨特需求調整解決方案的規模。我們的高階系統可提供超過 1M IOPS、延遲時間以微秒為測量單位、但如果您只需要 10 萬次 IOPS、您可以使用更小的控制器來調整儲存解決方案的大小、而該控制器仍可執行完全相同的儲存作業系統。

## 資料庫組態

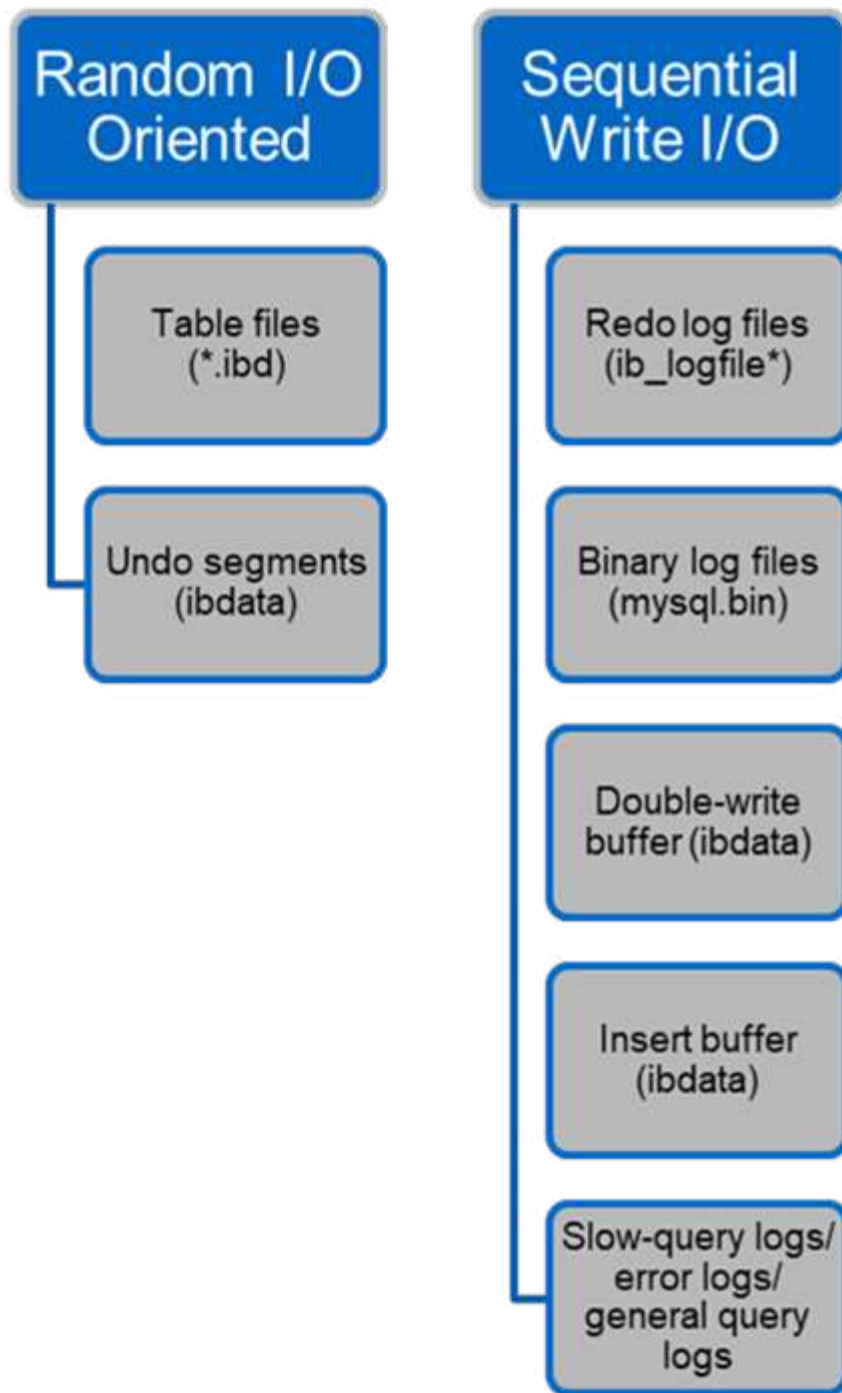
### MySQL 和 InnoDB

InnoDB 是儲存設備和 MySQL 伺服器之間的中間層、可將資料儲存到磁碟機中。



MySQL I/O 分為兩種類型：

- 隨機檔案 I/O
- 連續檔案 I/O



資料檔案會隨機讀取及覆寫、因此 IOPS 較高。因此、建議使用 SSD 儲存設備。

重做記錄檔和二進位記錄檔是交易記錄檔。它們會依序寫入、因此您可以在 HDD 上使用寫入快取獲得良好效能。恢復時會進行連續讀取、但很少會造成效能問題、因為記錄檔大小通常比資料檔案小、而連續讀取比隨機讀取快（發生在資料檔案上）。

雙寫入緩衝區是 InnoDB 的一項特殊功能。InnoDB 會先將排清的頁面寫入雙寫入緩衝區、然後將頁面寫入資料檔案的正確位置。此程序可防止頁面毀損。如果沒有雙寫入緩衝區、如果在寫入磁碟機程序期間發生電源故障、頁面可能會毀損。由於寫入雙寫入緩衝區是循序的、因此對 HDD 進行了高度最佳化。恢復時會進行連續讀取。

由於 ONTAP NVRAM 已提供寫入保護、因此不需要雙重寫入緩衝。MySQL 有一個參數、

`skip_innodb_doublewrite`，以停用雙寫入緩衝區。此功能可大幅提升效能。

插入緩衝區也是 InnoDB 的一項特殊功能。如果非唯一的次要索引區塊不在記憶體中、InnoDB 會將項目插入插入緩衝區、以避免隨機 I/O 作業。插入緩衝區會定期合併到資料庫的次要索引樹中。插入緩衝區可將 I/O 要求合併至同一個區塊、以減少 I/O 作業的數量；隨機 I/O 作業可以是連續的。插入緩衝區也針對 HDD 進行了高度最佳化。連續寫入和讀取都會在正常作業期間執行。

復原區段以隨機 I/O 為導向。為了保證多重版本併發（MVCC）、InnoDB 必須在復原區段中登錄舊影像。從復原區段讀取先前影像需要隨機讀取。如果您執行具有可重複讀取的長交易（例如 `mysqldump` — 單一交易）或執行長查詢、可能會發生隨機讀取。因此、在這種情況下、將復原區段儲存在 SSD 上會更好。如果您只執行簡短的交易或查詢、隨機讀取並不是問題。



- 由於 InnoDB I/O 特性、NetApp 建議 \* 下列儲存設計配置。
- 一個用於儲存 MySQL 的隨機和連續 I/O 導向檔案的磁碟區
- 另一個用於儲存 MySQL 純粹循序 I/O 導向檔案的磁碟區

此配置也能協助您設計資料保護原則與策略。

## MySQL 組態參數

NetApp 建議使用一些重要的 MySQL 組態參數、以獲得最佳效能。

參數	價值
<code>InnoDB_log_file_size</code>	256M
<code>InnoDB_Flush 記錄_AT_TRx_Commit</code>	2.
<code>InnoDB_doublewrite</code>	0%
<code>InnoDB_Flush 方法</code>	<code>fsync</code>
<code>InnoDB_緩衝區_Pool_size</code>	11g
<code>InnoDB_IO_capAC</code>	8192
<code>InnoDB_緩衝區集區執行個體</code>	8.
<code>InnoDB_LRU_SCAN_depth</code>	8192
<code>open_file_limit</code>	65535

若要設定本節所述的參數、您必須在 MySQL 組態檔（`my.cnf`）中變更這些參數。NetApp 最佳實務做法是在內部執行測試的結果。

## InnoDB\_log\_file\_size

為 InnoDB 記錄檔大小選取適當的大小、對於寫入作業以及在伺服器當機後擁有適當的還原時間都很重要。

由於有這麼多交易登入檔案、因此記錄檔大小對於寫入作業非常重要。修改記錄時、變更不會立即回寫到表格區。相反地、變更會記錄在記錄檔的結尾、而且頁面會標示為「髒污」。InnoDB 使用其記錄檔將隨機 I/O 轉換成連續 I/O

當記錄檔已滿時、會依序將不完整頁面寫入資料表空間、以釋放記錄檔中的空間。例如、假設某個伺服器在交易過程中當機、而且寫入作業只會記錄在記錄檔中。伺服器必須經過復原階段、記錄檔中記錄的變更會重新播放、才能重新上線。記錄檔中的項目越多、伺服器恢復所需的時間就越長。

在此範例中、記錄檔大小會同時影響還原時間和寫入效能。為記錄檔大小選擇正確的數字時、請平衡恢復時間與寫入效能。一般而言、128M 和 512M 之間的任何項目都是物超所值的。

## InnoDB\_Flush 記錄\_AT\_TRx\_Commit

當資料發生變更時、變更不會立即寫入儲存設備。

而是記錄在記錄緩衝區中、這是 InnoDB 分配給記錄在記錄檔中的緩衝區變更的記憶體部分。InnoDB 會在交易提交、緩衝區滿時、或每秒一次（以先發生的事件為準）、將緩衝區排清至記錄檔。控制此程序的組態變數是 InnoDB\_Flush 記錄\_AT\_TRx\_Commit。價值選項包括：

- 當您設定時 `innodb_flush_log_trx_at_commit=0`、InnoDB 會將修改過的資料（在 InnoDB 緩衝區集中）寫入記錄檔（`IB_logfile`）、並每秒清除記錄檔（寫入儲存區）。不過、提交交易時、它不會執行任何動作。如果發生電源故障或系統當機、則沒有任何未排清的資料可恢復、因為它不會寫入記錄檔或磁碟機。
- 當您設定時 `innodb_flush_log_trx_commit=1`、InnoDB 會將記錄緩衝區寫入交易記錄檔、並在每筆交易中將記錄檔排清至持久儲存區。例如、對於所有交易認可、InnoDB 會寫入記錄、然後寫入儲存設備。儲存速度變慢會對效能造成負面影響、例如每秒 InnoDB 交易數會減少。
- 當您設定時 `innodb_flush_log_trx_commit=2`、InnoDB 會在每次提交時將記錄緩衝區寫入記錄檔、但不會將資料寫入儲存區。InnoDB 每秒會清除一次資料。即使發生電源故障或系統當機、記錄檔中也有選項 2 資料可供使用、而且可恢復。

如果效能是主要目標、請將值設為 2。由於 InnoDB 每秒一次寫入磁碟機、而非每次提交交易時、效能大幅提升。如果發生停電或當機、可從交易記錄中恢復資料。

如果資料安全是主要目標、請將值設為 1、以便每次提交交易時、InnoDB 都會將資料清出磁碟機。不過、效能可能會受到影響。



\* NetApp 建議 \* 將 InnoDB\_Flush 日誌\_AT\_TRx\_Commit 值設為 2、以獲得更好的效能。

## InnoDB\_doublewrite

何時 `innodb_doublewrite` 啟用（預設）時、InnoDB 會將所有資料儲存兩次：先儲存至雙寫入緩衝區、然後儲存至實際資料檔案。

您可以使用關閉此參數 `--skip-innodb_doublewrite` 針對效能標竿、或是當您比較在意最高效能、而非資料完整性或可能的故障時。InnoDB 使用稱為雙寫入的檔案清除技術。InnoDB 在將頁面寫入資料檔案之前、會將其寫入稱為雙寫入緩衝區的鄰近區域。寫入和清除雙寫入緩衝區完成後、InnoDB 會將頁面寫入資料檔案中的適當位置。如果作業系統或 `mysqld` 程序在頁面寫入期間當機、InnoDB 之後可以在當機恢復期間、從雙寫入緩衝區找到良好的頁面複本。



\* NetApp 建議 \* 停用雙寫入緩衝區。ONTAP NVRAM 的功能相同。雙重緩衝會不必要地損害效能。

## InnoDB\_緩衝區\_Pool\_size

InnoDB 緩衝資源池是任何調校活動中最重要的部分。

InnoDB 在很大程度上仰賴緩衝區集區來快取索引和資料、調適性雜湊索引、插入緩衝區、以及許多內部使用的其他資料結構。緩衝區集區也會緩衝資料的變更、這樣就不需要立即對儲存設備執行寫入作業、進而改善效能。緩衝區集區是 InnoDB 的一部分、必須據此調整其大小。設定緩衝區集區大小時、請考量下列因素：

- 若為僅 InnoDB 的專用機器、請將緩衝區集區大小設為 80% 以上的可用 RAM 。
- 如果不是 MySQL 專用伺服器、請將 RAM 大小設為 50% 。

## InnoDB\_Flush 方法

InnoDB\_flush\_method 參數指定 InnoDB 如何開啟及排清記錄檔和資料檔。

最佳化

在 InnoDB 最佳化中、如果適用、設定此參數會調整資料庫效能。

下列選項用於透過 InnoDB 排清檔案：

- `fsync`。InnoDB 使用 `fsync()` 系統呼叫以清除資料和記錄檔。此選項為預設設定。
- `O_DSYNC`。InnoDB 使用 `O_DSYNC` 用於打開和刷新日誌文件和 `fsync()` 以刷新數據文件的選項。InnoDB 不使用 `O_DSYNC` 直接來說、因為 UNIX 的許多種類都有問題。
- `O_DIRECT`。InnoDB 使用 `O_DIRECT` 選項（或 `directio()` 在 Solaris 上）開啟資料檔案及使用 `fsync()` 清除資料和記錄檔。此選項可在某些版本的 GNU/Linux、FreeBSD 和 Solaris 上使用。
- `O_DIRECT_NO_FSYNC`。InnoDB 使用 `O_DIRECT` 排清 I/O 時的選項；不過、它會跳過 `fsync()` 之後進行系統通話。此選項不適用於某些類型的檔案系統（例如 XFS）。如果您不確定檔案系統是否需要 `fsync()` 系統呼叫（例如為了保留所有檔案中繼資料）使用 `O_DIRECT` 選項。

觀察

在 NetApp 實驗室測試中、`fsync` 預設選項用於 NFS 和 SAN、與相較之下、這是一項很棒的效能改進工具 `O_DIRECT`。使用「齊平」方法時為 `O_DIRECT` 使用 ONTAP 時、我們觀察到用戶端會以序列方式、在 4096 區塊的邊界寫入大量的單位位元組寫入資料。這些寫入會增加網路延遲並降低效能。

## InnoDB\_IO\_capAC

在 InnoDB 外掛程式中、從 MySQL 5.7 新增名為 `InnoDB_IO_capACure`。

它可控制 InnoDB 執行的 IOPS 上限（包括髒頁的排清率、以及插入緩衝區 [ibuf] 批次大小）。`InnoDB_IO_capAC` 容量參數會根據 InnoDB 背景工作來設定 IOPS 上限、例如從緩衝區排清頁面、以及從變更緩衝區合併資料。

將 `InnoDB_IO_capAC` 容量參數設定為系統每秒可執行的 I/O 作業大約數目。理想情況下、請盡可能將設定保持在低的位置、但不要太低、讓背景活動變慢。如果設定太高、資料會從緩衝區移除、並太快插入緩衝區以供快取、以提供顯著效益。



\* NetApp 建議 \* 如果在 NFS 上使用此設定、請分析 IOPS ( Sys台 / Fio ) 的測試結果、並據此設定參數。除非您在 InnoDB 緩衝資源池中看到比您想要的更多修改或不乾淨頁面、否則請使用最小的值來進行排清和清除。



除非您證明較低的值不足以應付工作負載、否則請勿使用極端值、例如 20,000 或更多。

InnoDB\_IO\_capACure. 參數可規範排清率和相關 I/O



您可以將此參數或 InnoDB\_IO\_capACure\_max 參數設定得太高、並以提早排清的方式浪費 I/O 作業、進而嚴重損害效能。

## InnoDB\_LRU\_SCAN\_depth

◦ innodb\_lru\_scan\_depth 參數會影響 InnoDB 緩衝區集區之排清作業的演算法和啟發性。

此參數主要是效能專家調校 I/O 密集工作負載的興趣所在。對於每個緩衝區集區執行個體、此參數會指定最少使用 (LRU) 頁面清單中、頁面清理程式執行緒應繼續掃描的程度、以尋找要清除的髒頁面。此背景作業每秒執行一次。

您可以上下調整值、將可用頁數降至最低。請勿將此值設定得比所需值高得多、因為掃描可能會產生重大的效能成本。此外、請考慮在變更緩衝區集區執行個體數目時調整此參數、因為 `innodb_lru_scan_depth * innodb_buffer_pool_instances` 定義頁面清理程式執行緒每秒執行的工作量。

小於預設值的設定適用於大部分的工作負載。只有在典型工作負載下有備用 I/O 容量時、才考慮增加此值。相反地、如果寫入密集的工作負載使 I/O 容量飽和、請降低該值、尤其是當您擁有大型緩衝區集區時。

## open\_file\_limits

◦ open\_file\_limits 參數決定作業系統允許 mysqld 開啟的檔案數目。

此參數在執行階段的值是系統允許的實際值、可能與您在伺服器啟動時指定的值不同。在 MySQL 無法變更開啟檔案數量的系統上、此值為 0。有效 `open_files_limit` 此值是根據系統啟動時指定的值 (如果有) 和的值而定 `max_connections` 和 `table_open_cache` 使用這些公式：

- $10 + \text{max\_connections} + (\text{table\_open\_cache} \times 2)$
- $\text{max\_connections} \times 5$ .
- 如果為正、則作業系統限制
- 如果作業系統限制為無限：`open_files_limit` 在啟動時指定值；若無、則指定值為 5、000

伺服器會嘗試使用這四個值的最大值來取得檔案描述元數目。如果無法取得這麼多描述元、伺服器會嘗試取得系統允許的數量。

## 主機組態



## MySQL 容器化

MySQL 資料庫的容器化日漸普及。

低層級的容器管理幾乎總是透過 Docker 來執行。OpenShift 和 Kubernetes 等容器管理平台讓大型容器環境的管理變得更簡單。容器化的優點包括成本較低、因為不需要授權 Hypervisor。此外、容器可讓多個資料庫彼此隔離執行、同時共用相同的基礎核心和作業系統。容器可在微秒內完成佈建。

NetApp 提供 Astra Trident、可提供進階的儲存管理功能。例如、Astra Trident 可讓在 Kubernetes 中建立的容器自動在適當的層級上配置儲存設備、套用匯出原則、設定快照原則、甚至將一個容器複製到另一個容器。如需其他資訊、請參閱 "[Astra Trident文件](#)"。

## MySQL 和 NFSv3 插槽表

Linux 上的 NFSv3 效能取決於所呼叫的參數 `tcp_max_slot_table_entries`。

TCP 插槽表是與主機匯流排介面卡（HBA）佇列深度相當的 NFSv3。這些表格可控制任何時間都可以處理的 NFS 作業數量。預設值通常為 16、這對於最佳效能而言太低。相反的問題發生在較新的 Linux 核心上、這會自動將 TCP 插槽表格限制增加到要求使 NFS 伺服器飽和的層級。

為了達到最佳效能並避免效能問題、請調整控制 TCP 插槽表的核心參數。

執行 `sysctl -a | grep tcp.*.slot_table` 並觀察下列參數：

```
# sysctl -a | grep tcp.*.slot_table
sunrpc.tcp_max_slot_table_entries = 128
sunrpc.tcp_slot_table_entries = 128
```

所有 Linux 系統都應該包括在內 `sunrpc.tcp_slot_table_entries`、但只有部分包含在內 `sunrpc.tcp_max_slot_table_entries`。兩者都應設為 128。

### 注意

若未設定這些參數、可能會對效能造成重大影響。在某些情況下、效能會受到限制、因為 Linux 作業系統沒有發出足夠的 I/O。在其他情況下、隨著 Linux 作業系統嘗試發出的 I/O 數量超過可服務的數量、I/O 延遲也會增加。

## I/O 排程器和 MySQL

Linux 核心可讓您以低層級控制 I/O 排程封鎖裝置的方式。

Linux 各版本的預設值差異極大。MySQL 建議您使用 `NOOP` 或是 `deadline` 在 Linux 上具有原生非同步 I/O（AIO）的 I/O 排程器。一般而言、NetApp 客戶和內部測試都能透過 `NoOps` 獲得更好的結果。

MySQL 的 InnoDB 儲存引擎使用 Linux 上的非同步 I/O 子系統（原生 AIO）來執行預先讀取和寫入資料檔案頁面的要求。此行為由控制 `innodb_use_native_aio` 組態選項、預設為啟用。有了原生的整合式全功能電腦、I/O 排程器的類型對 I/O 效能有更大的影響。執行效能標竿、判斷哪一個 I/O 排程器可為您的工作負載和環境提供最佳結果。

如需設定 I/O 排程器的指示、請參閱相關的 Linux 和 MySQL 文件。

## MySQL 檔案描述元

若要執行、MySQL 伺服器需要檔案描述元、而且預設值不足。

它會使用它們來開啟新的連線、將資料表儲存在快取中、建立暫時資料表來解決複雜的查詢、以及存取持續性的查詢。如果 `mysqld` 在需要時無法開啟新檔案、它就能停止正常運作。此問題的常見症狀是錯誤 24 「開啟的檔案太多」。 `mysqld` 可以同時開啟的檔案描述元數量是由定義的 `open_files_limit` 設定在組態檔案中的選項 (`/etc/my.cnf`)。但是 `open_files_limit` 也取決於作業系統的限制。這種相依性會使設定變數變得更複雜。

MySQL 無法設定 `open_files_limit` 選項高於下所指定的選項 `ulimit 'open files'`。因此、您必須在作業系統層級明確設定這些限制、才能讓 MySQL 視需要開啟檔案。檢查 Linux 檔案限制的方法有兩種：

- `ulimit` 命令會快速為您提供所允許或鎖定參數的詳細說明。執行此命令所做的變更並非永久性變更、將會在系統重新開機後清除。
- 變更為 `/etc/security/limit.conf` 檔案是永久性的、不受系統重新開機影響。

請務必同時變更使用者 `mysql` 的硬限制和軟限制。以下摘錄來自組態：

```
mysql hard nofile 65535
mysql soft nofile 65353
```

同時、請在中更新相同的組態 `my.cnf` 以完全使用開放式檔案限制。

## 儲存組態

### 使用 NFS 的 MySQL

MySQL 文件建議您將 NFSv4 用於 NAS 部署。

#### ONTAP NFS 傳輸大小

根據預設、ONTAP 將 NFS IO 大小限制為 64K。使用 MySQL 資料庫的隨機 IO 使用的區塊大小要小得多、遠低於 64K 上限。大型區塊 IO 通常是平行處理的、因此 64K 最大值也不是限制。

有些工作負載的上限為 64K、因此會造成限制。尤其是、如果資料庫執行的 IO 數量較少、但容量較大、例如完整表格掃描備份作業等單執行緒作業、將會更快、更有效率地執行。ONTAP 搭配資料庫工作負載的最佳 IO 處理大小為 256k。針對下列特定作業系統所列出的 NFS 裝載選項已相應從 64K 更新至 256k。

指定 ONTAP SVM 的最大傳輸大小可變更如下：

```
Cluster01::> set advanced
```

```
Warning: These advanced commands are potentially dangerous; use them only  
when directed to do so by NetApp personnel.
```

```
Do you want to continue? {y|n}: y
```

```
Cluster01::*> nfs server modify -vserver vserver1 -tcp-max-xfer-size  
262144
```



切勿將 ONTAP 上允許的傳輸大小上限降至低於目前掛載之 NFS 檔案系統的 rsize/wsize 值。這可能會在某些作業系統中造成當機或甚至資料毀損。例如、如果 NFS 用戶端目前設定為 rsize/wsize 65536、則 ONTAP 最大傳輸大小可在 65536 到 1048576 之間調整、因為用戶端本身受到限制、因此沒有任何影響。將傳輸大小上限降至 65536 以下可能會損害可用度或資料。

- NetApp 推薦 \*



設定下列 NFSv4 Fstab ( /etc/fstab ) 設定：

```
nfs4 rw,  
hard,nointr,bg,vers=4,proto=tcp,noatime,rsize=262144,wsiz=262144
```



NFSv3 的常見問題是停電後鎖定的 InnoDB 記錄檔。使用時間或切換記錄檔可解決此問題。不過、NFSv4 有鎖定作業、並追蹤開啟的檔案和委派。

## MySQL 搭配 SAN

有兩個選項可以使用常見的雙磁碟區模型、來設定 MySQL 與 SAN 。

只要 I/O 和容量需求在單一 LUN 檔案系統的限制範圍內、就可以將較小的資料庫放在一對標準 LUN 上。例如、需要大約 2K 隨機 IOPS 的資料庫、可以裝載在單一 LUN 上的單一檔案系統上。同樣地、只有 100GB 大小的資料庫也能放在單一 LUN 上、而不會造成管理問題。

大型資料庫需要多個 LUN 。例如、需要 10 萬 IOPS 的資料庫最可能需要至少八個 LUN 。由於磁碟機的 SCSI 通道數量不足、單一 LUN 將成為瓶頸。10TB 資料庫同樣難以在單一 10TB LUN 上進行管理。邏輯磁碟區管理程式的設計旨在將多個 LUN 的效能和容量功能結合在一起、以改善效能和管理能力。

在這兩種情況下、一對 ONTAP 磁碟區應該足夠。透過簡單的組態、資料檔案 LUN 會與記錄 LUN 一樣置於專用磁碟區中。在邏輯 Volume Manager 組態下、資料檔案 Volume 群組中的所有 LUN 都將位於專用磁碟區、而記錄 Volume 群組的 LUN 則位於第二個專用磁碟區。

- NetApp 建議 \* 在 SAN 上使用兩個檔案系統進行 MySQL 部署：
- 第一個檔案系統會儲存所有 MySQL 資料、包括資料表空間、資料和索引。
- 第二個檔案系統會儲存所有記錄（二進位記錄、慢速記錄和交易記錄）。



以這種方式分隔資料有多種原因、包括：

- 資料檔案和記錄檔的 I/O 模式各不相同。將它們區隔、就能透過 QoS 控制提供更多選項。
- 最佳化使用 Snapshot 技術需要能夠個別還原資料檔案。將資料檔案與記錄檔混合會干擾資料檔案還原。
- NetApp SnapMirror 技術可用於為資料庫提供簡單、低 RPO 的災難恢復功能、但是資料檔案和記錄檔需要不同的複寫排程。



使用這種基本的兩個磁碟區配置、以符合未來需求的解決方案、以便在需要時使用所有 ONTAP 功能。

- NetApp 建議 \* 使用 ext4 檔案系統格式化磁碟機、因為有下列功能：
- 延伸的區塊管理方法、用於日誌檔案系統（JFS）中的區塊管理功能、以及延伸檔案系統（XFS）的延遲分配功能。
- ext4 允許檔案系統最多 1 個 exbibyte（ $2^{60}$  位元組）、檔案最多 16 個 tebibytes（ $16 * 2^{40}$  位元組）。相反地、ext3 檔案系統僅支援 16TB 的最大檔案系統大小、最大檔案大小為 2TB。
- 在 ext4 檔案系統中、多區塊分配（mballoc）會在單一作業中為檔案分配多個區塊、而非如 ext3 中逐一分配。此組態可減少多次呼叫區塊分配器的成本、並最佳化記憶體配置。
- 雖然 XFS 是許多 Linux 套裝作業系統的預設值、但它以不同方式管理中繼資料、不適合某些 MySQL 組態。



- NetApp 建議 \* 將 4K 區塊大小選項與 mkfs 公用程式搭配使用、以符合現有區塊 LUN 大小。



```
mkfs.ext4 -b 4096
```

NetApp LUN 以 4KB 實體區塊儲存資料、產生八個 512 位元組的邏輯區塊。

如果您未設定相同的區塊大小、I/O 將無法正確對齊實體區塊、而且可能會在 RAID 群組中的兩個不同磁碟機中寫入、導致延遲。



請務必調整 I/O、以順利進行讀寫作業。但是、當 I/O 從不在實體區塊開頭的邏輯區塊開始時、I/O 會未對齊。I/O 作業只有在邏輯區塊（實體區塊中的第一個邏輯區塊）開始時才會對齊。

# Oracle 資料庫

## ONTAP 上的 Oracle 資料庫

ONTAP 專為 Oracle 資料庫所設計。數十年來、ONTAP 已針對關聯式資料庫 I/O 的獨特需求進行最佳化、並特別建立多項 ONTAP 功能、以滿足 Oracle 資料庫的需求、甚至是 Oracle Inc. 本身的要求。



本文件取代先前發佈的技術報告 TR-3633：ONTAP 上的 Oracle 資料庫；TR-4591：Oracle 資料保護：備份、還原、複寫；TR-4592：MetroCluster 上的 Oracle；以及 TR-4534：將 Oracle 資料庫移轉至 NetApp 儲存系統

除了 ONTAP 為您的資料庫環境帶來價值的許多可能方式之外、也有許多使用者需求、包括資料庫大小、效能需求和資料保護需求。NetApp 儲存設備的已知部署包括從虛擬化環境（約 6、000 個在 VMware ESX 下執行的資料庫）到目前規模為 996TB 且不斷成長的單一執行個體資料倉儲、應有盡有。因此、在 NetApp 儲存設備上設定 Oracle 資料庫的最佳實務做法很少。

在 NetApp 儲存設備上操作 Oracle 資料庫的需求有兩種解決方法。首先、當有明確的最佳實務做法時、將會特別加以說明。我們將根據 Oracle 儲存解決方案架構設計師的特定業務需求、詳細說明許多設計考量。

## 組態 ONTAP

### RAID 和 Oracle 資料庫

RAID 是指使用備援功能來保護資料、避免磁碟機遺失。

在設定用於 Oracle 資料庫和其他企業應用程式的 NetApp 儲存設備時、偶爾會出現 RAID 層級的問題。許多舊版 Oracle 儲存陣列組態的最佳實務做法都包含使用 RAID 鏡射和 / 或避免使用特定類型 RAID 的警告。雖然這些來源提出有效點、但這些來源不適用於 RAID 4、以及 ONTAP 中使用的 NetApp RAID DP 和 RAID-TEC 技術。

RAID 4、RAID 5、RAID 6、RAID DP 和 RAID-TEC 都使用同位元檢查來確保磁碟機故障不會導致資料遺失。與鏡像相比、這些 RAID 選項可提供更好的儲存使用率、但大多數 RAID 實作都有影響寫入作業的缺點。在其他 RAID 實作上完成寫入作業可能需要多個磁碟機讀取才能重新產生同位元資料、這是一種通常稱為 RAID 懲罰的程序。

不過、ONTAP 並不會因此而受到此 RAID 處罰。這是因為 NetApp WAFL（隨處寫入檔案配置）與 RAID 層整合。寫入作業會整合在 RAM 中、並準備為完整的 RAID 等量磁碟區、包括同位元檢查產生。ONTAP 不需要執行讀取即可完成寫入、這表示 ONTAP 和 WAFL 可以避免 RAID 的損失。對於延遲關鍵作業（例如重作記錄）的效能不受阻礙、隨機的資料檔案寫入不會因重新產生同位元檢查而導致任何 RAID 損失。

在統計可靠性方面、即使是 RAID DP 也能提供比 RAID 鏡射更好的保護。主要問題是 RAID 重建期間對磁碟機的需求。有了鏡射 RAID 集、當磁碟機在重建時發生故障、而在 RAID 組中重建其合作夥伴時、資料遺失的風險遠高於 RAID DP 組中發生三重磁碟機故障的風險。

### Oracle 資料庫與儲存容量管理

使用可預測、可管理的高效能企業儲存設備來管理資料庫或其他企業應用程式、需要磁碟

機上的一些可用空間來進行資料和中繼資料管理。所需的可用空間量取決於使用的磁碟機類型和業務程序。

可用空間定義為任何不用於實際資料的空間、並包括集合體本身的未分配空間、以及組成磁碟區內的未使用空間。也必須考慮精簡配置。例如、某個磁碟區可能包含 1TB LUN、其中只有 50% 被實際資料使用。在精簡佈建的環境中、這似乎是消耗 500GB 的空間。不過、在完全佈建的環境中、1TB 的完整容量似乎正在使用中。500GB 的未分配空間會隱藏起來。實際資料未使用此空間、因此應納入總可用空間的計算。

NetApp 對於企業應用程式所使用的儲存系統建議如下：

### SSD 集合體、包括 AFF 系統



\* NetApp 建議 \* 至少有 10% 的可用空間。這包括所有未使用的空間、包括集合體或磁碟區內的可用空間、以及因使用完整資源配置而分配但實際資料未使用的任何可用空間。邏輯空間並不重要、問題在於實際可用的實體空間可用於資料儲存。

10% 可用空間的建議非常保守。SSD 集合體可支援使用率更高的工作負載、而不會對效能造成任何影響。不過、隨著 Aggregate 的使用率增加、如果未仔細監控使用率、則用盡空間的風險也會增加。此外、當系統的容量達到 99% 時、可能不會導致效能降低、但在訂購額外硬體時、可能需要管理人員努力避免系統完全裝滿、而且可能需要一些時間來採購和安裝額外的磁碟機。

### HDD 集合體、包括 Flash Pool 集合體



\* 使用旋轉磁碟機時、NetApp 建議 \* 至少有 15% 的可用空間。這包括所有未使用的空間、包括集合體或磁碟區內的可用空間、以及因使用完整資源配置而分配但實際資料未使用的任何可用空間。在免費語音方法中、效能將受到 10% 的影響。

## Oracle 資料庫與儲存虛擬機器

### Oracle 資料庫儲存管理集中在儲存虛擬機器（SVM）上

SVM 在 ONTAP CLI 中稱為 vserver、是基本的儲存功能單元、比較 SVM 與 VMware ESX 伺服器上的客體是很有用的。

首次安裝時、ESX 沒有預先設定的功能、例如代管來賓作業系統或支援終端使用者應用程式。它是空容器、直到定義虛擬機器（VM）為止。ONTAP 類似。第一次安裝 ONTAP 時、只有建立 SVM、它才具備資料服務功能。定義資料服務的是 SVM 特性設定。

與儲存架構的其他層面一樣、SVM 和邏輯介面（LIF）設計的最佳選項、在很大程度上取決於擴充需求和業務需求。

### SVM

我們並未正式提供 ONTAP 的 SVM 資源配置最佳實務做法。正確的方法取決於管理和安全要求。

大多數客戶只需操作一部主要 SVM、即可滿足大部分的日常需求、然後針對特殊需求建立少量 SVM。例如、您可能想要建立：

- 由專業團隊管理的關鍵業務資料庫 SVM
- 開發群組的 SVM 已獲得完整的管理控制權、可讓他們獨立管理自己的儲存設備

- 必須限制管理團隊的 SVM、用於處理敏感業務資料、例如人力資源或財務報告資料

在多租戶環境中、每個租戶的資料都可以獲得專用的 SVM。每個叢集、HA 配對和節點的 SVM 和生命量限制取決於所使用的傳輸協定、節點模型和 ONTAP 版本。請參閱 "[NetApp Hardware Universe](#)" 針對這些限制。

## ONTAP QoS 的 Oracle 資料庫效能管理

安全有效地管理多個 Oracle 資料庫、需要有效的 QoS 策略。原因在於現代儲存系統的效能功能不斷提升。

具體而言、由於採用 All Flash 儲存設備的情況增加、因此能夠整合工作負載。依賴旋轉媒體的儲存陣列往往只支援有限數量的 I/O 密集工作負載、因為舊版旋轉式磁碟機技術的 IOPS 功能有限。一或兩個高作用中的資料庫會在儲存控制器達到限制之前、使基礎磁碟機飽和。這種情況已經改變。即使是功能最強大的儲存控制器、相對少數 SSD 磁碟機的效能也能達到飽和。這意味著控制器的完整功能可以充分發揮、而不會因為旋轉媒體延遲尖峰而突然降低效能。

舉例來說、簡單的雙節點 HA AFF A800 系統能夠在延遲超過 1 毫秒之前、提供高達 100 萬次的隨機 IOPS 服務。只有很少單一工作負載會達到這類層級。充分利用此 AFF A800 系統陣列、將需要託管多個工作負載、同時確保可預測性、這需要 QoS 控制。

ONTAP 中有兩種服務品質 (QoS)：IOPS 和頻寬。QoS 控制可套用至 SVM、磁碟區、LUN 和檔案。

### IOPS QoS

IOPS QoS 控制顯然是以指定資源的 IOPS 總計為基礎、但 IOPS QoS 的許多層面可能並不符合直覺。剛開始有幾位客戶對於達到 IOPS 臨界值時、延遲明顯增加感到困惑。延遲增加是限制 IOPS 的自然結果。從邏輯上講、它的運作方式與權杖系統類似。例如、如果包含資料檔案的特定磁碟區有 10K IOPS 限制、則每個到達的 I/O 都必須先接收權杖才能繼續處理。只要在指定的秒數內使用的權杖不超過 10K、就不會有延遲。如果 IO 作業必須等待接收其權杖、則此等待會顯示為額外延遲。工作負載相對於 QoS 限制的推動越大、每個 IO 在佇列中等待處理的時間就越長、使用者認為延遲越高。



將 QoS 控制套用至資料庫交易 / 重做記錄資料時、請務必謹慎。雖然重做記錄的效能需求通常比資料檔案低很多、但重做記錄活動卻很繁瑣。IO 會以簡短的脈衝形式發生、而顯示適合平均重做 IO 層級的 QoS 限制、對於實際需求而言可能太低。結果可能會造成嚴重的效能限制、因為每次重做記錄突增時、QoS 都會啟動。一般而言、重作和歸檔記錄不應受到 QoS 的限制。

### 頻寬 QoS

並非所有 I/O 大小都相同。例如、資料庫可能會執行大量的小區塊讀取、導致達到 IOPS 臨界值、但資料庫也可能執行完整的資料表掃描作業、這項作業會由極少數的大量區塊讀取所組成、佔用大量頻寬、但 IOPS 相對較少。

同樣地、VMware 環境在開機期間可能會產生極高的隨機 IOPS、但在外部備份期間執行的 IO 會較少、但會較大。

有時有效管理效能需要 IOPS 或頻寬 QoS 限制、甚至兩者都需要。

### 最低 / 保證的 QoS

許多客戶尋求的解決方案都包含保證的 QoS、這比看起來更難達成、而且可能相當浪費。例如、如果要放置 10 個具有 10K IOPS 保證的資料庫、就必須針對所有 10 個資料庫同時以 10K IOPS 執行的情況來調整系統規模、總共需要 10 萬個。

最適合用於最低 QoS 控制的是保護關鍵工作負載。例如、假設 ONTAP 控制器的 IOPS 最高可達 50 萬、同時混合了生產與開發工作負載。您應該將 QoS 原則上限套用至開發工作負載、以防止任何指定的資料庫壟斷控制器。然後、您可以將最低 QoS 原則套用至正式作業工作負載、以確保它們在需要時隨時都能使用所需的 IOPS。

## 調適性 QoS

調適性 QoS 是指 ONTAP 功能、其中 QoS 限制是根據儲存物件的容量而定。它很少用於資料庫、因為資料庫的大小與其效能需求之間通常沒有任何連結。大型資料庫可能幾乎無法運作、而較小的資料庫則可能是 IOPS 最密集的資料庫。

Adaptive QoS 對於虛擬化資料存放區非常有用、因為這類資料集的 IOPS 需求往往與資料庫的總大小相關。較新的資料存放區包含 1TB 的 VMDK 檔案、可能需要的效能約為 2TB 資料存放區的一半。Adaptive QoS 可讓您在資料存放區填入資料時、自動增加 QoS 限制。

## Oracle 資料庫與 ONTAP 效率功能

ONTAP 空間效率功能已針對 Oracle 資料庫進行最佳化。在幾乎所有情況下、最佳方法是在啟用所有效率功能的情況下、保留預設值。

空間效率功能（例如壓縮、壓縮和重複資料刪除）的設計、是為了增加符合特定實體儲存量的邏輯資料量。結果是降低成本和管理成本。

在高層級、壓縮是一種數學程序、可偵測及編碼資料模式、以減少空間需求。相反地、重複資料刪除功能會偵測實際重複的資料區塊、並移除額外的複本。資料實作可讓多個邏輯區塊在媒體上共用相同的實體區塊。



請參閱以下關於精簡配置的章節、以瞭解儲存效率與部分保留之間互動的說明。

## 壓縮

在提供 All Flash 儲存系統之前、以陣列為基礎的壓縮價值有限、因為大多數 I/O 密集的工作負載都需要大量磁碟來提供可接受的效能。儲存系統的容量總是比所需的容量大得多、這是大量磁碟機的副作用。固態儲存設備的興起、改變了這種情況。不再需要純粹為了獲得良好效能而大幅過度配置磁碟機。儲存系統中的磁碟機空間可與實際容量需求相符。

固態硬碟機（SSD）的 IOPS 容量增加、幾乎總是比旋轉硬碟機節省成本、但壓縮技術可以增加固態媒體的有效容量、進而進一步節省成本。

壓縮資料的方法有好幾種。許多資料庫都包含自己的壓縮功能、但在客戶環境中很少會發現這種情況。其原因通常是對壓縮資料的 \* 變更 \* 效能會受到影響、而對於某些應用程式而言、資料庫層級壓縮的授權成本較高。最後、對資料庫作業的整體效能影響。對於執行資料壓縮與解壓縮的 CPU、而非實際的資料庫工作、支付高昂的每 CPU 授權成本是不合理的。更好的選擇是將壓縮工作卸載到儲存系統。

## 自適應壓縮

即使在以微秒為單位測量延遲的 All Flash 環境中、主動式壓縮也已針對企業工作負載進行徹底測試、且未對效能產生任何影響。有些客戶甚至報告使用壓縮技術時效能會提高、因為資料會保持在快取中的壓縮、有效增加控制器中可用的快取數量。

ONTAP 以 4KB 單位管理實體區塊。自適應壓縮使用 8KB 的預設壓縮區塊大小、也就是以 8KB 為單位壓縮資料。這與關係式資料庫最常使用的 8KB 區塊大小相符。隨著將更多資料壓縮成單一單元、壓縮演算法就會變得更有效率。32 KB 壓縮區塊大小比 8 KB 壓縮區塊單元更具空間效率。這表示使用預設 8KB 區塊大小的調適式



壓縮確實會導致效率稍微降低、但使用較小的壓縮區塊大小也有很大的好處。資料庫工作負載包括大量的覆寫活動。若要覆寫 32 KB 壓縮資料區塊的 8KB 資料、必須讀回整個 32 KB 的邏輯資料、將其解壓縮、更新所需的 8 KB 區域、重新壓縮、然後將整個 32 KB 寫入磁碟機。這對儲存系統來說是非常昂貴的作業、也是因為某些競爭儲存陣列以較大的壓縮區塊大小為基礎、也會對資料庫工作負載造成重大效能損失的原因。



調適式壓縮所使用的區塊大小最多可增加至 32KB。這可能會改善儲存效率、而且當大量的這類資料儲存在陣列上時、應該考慮用於靜態檔案、例如交易記錄檔和備份檔案。在某些情況下、使用 16KB 或 32KB 區塊大小的作用中資料庫、也可能因為增加適應式壓縮的區塊大小而受惠。請洽詢 NetApp 或合作夥伴代表、瞭解這是否適合您的工作負載。



在串流備份目的地上、不應將大於 8KB 的壓縮區塊大小與重複資料刪除一起使用。原因是備份資料的細微變更會影響 32KB 壓縮時間。如果視窗移動、則產生的壓縮資料會在整個檔案中有所不同。重複資料刪除是在壓縮之後進行、這表示重複資料刪除引擎會以不同的方式檢視每個壓縮備份。如果需要重複資料刪除串流備份、則只應使用 8KB 區塊調適性壓縮。調適性壓縮較為理想、因為它的區塊大小較小、不會中斷重複資料刪除的效率。由於類似的原因、主機端壓縮也會影響重複資料刪除的效率。

### 壓縮對齊

資料庫環境中的調適性壓縮需要考量壓縮區塊對齊。這樣做只是對隨機覆寫非常特定區塊的資料的考量。這種方法的概念與整體檔案系統對齊方式類似、檔案系統的開始必須與 4K 裝置邊界對齊、檔案系統的區塊大小必須是 4K 的倍數。

例如、只有在檔案與檔案系統本身的 8KB 邊界對齊時、才會壓縮寫入 8KB 檔案。這表示它必須落在檔案的前 8KB、檔案的第二 8KB 等。確保正確對齊的最簡單方法是使用正確的 LUN 類型、建立的任何分割區都應該與 8K 的倍數裝置開始偏移、並使用資料庫區塊大小的倍數檔案系統區塊大小。

備份或交易記錄等資料會循序寫入跨越多個區塊的作業、所有這些區塊都會被壓縮。因此、不需要考慮對齊。唯一令人擔憂的 I/O 模式是隨機覆寫檔案。

### 資料壓縮

資料壓縮技術可改善壓縮效率。如前所述、僅有調適式壓縮功能、最多可節省 2 : 1、因為它僅限於在 4KB WAFL 區塊中儲存 8KB I/O。較大區塊大小的壓縮方法可提供更好的效率。不過、這些資料不適合受到小型區塊覆寫的資料。解壓縮 32KB 的資料單元、更新 8KB 部分、重新壓縮及回寫磁碟機、都會產生額外的負荷。

資料壓縮的運作方式是允許將多個邏輯區塊儲存在實體區塊內。例如、含有高度壓縮資料（例如文字或部分完整區塊）的資料庫、可能會從 8KB 壓縮至 1KB。如果沒有壓縮、1KB 的資料仍會佔用整個 4KB 區塊。即時資料壓縮功能可將 1KB 的壓縮資料與其他壓縮資料一起儲存在 1KB 的實體空間中。這不是一項壓縮技術、只是在磁碟機上分配空間的一種更有效率的方法、因此不應產生任何可偵測的效能影響。

節省的程度各不相同。已壓縮或加密的資料通常無法進一步壓縮、因此資料集無法從資料壓縮中獲益。相反地、新初始化的資料檔案僅包含區塊中繼資料和零、最多可壓縮至 80 : 1。

### 對溫度敏感的儲存效率

溫度敏感儲存效率（TSSE）是 ONTAP 9.8 及更新版本中提供的產品、它仰賴區塊存取熱圖來識別不常存取的區塊、並以更高的效率加以壓縮。

### 重複資料刪除

重複資料刪除是從資料集移除重複的區塊大小。例如、如果 10 個不同的檔案中存在相同的 4KB 區塊、重複資

料刪除會將所有 10 個檔案中的 4KB 區塊重新導向至相同的 4KB 實體區塊。結果是該資料的效率提升 10 : 1。

VMware 來賓開機 LUN 等資料通常會極好地刪除重複資料、因為這些資料包含相同作業系統檔案的多個複本。效率達到 100 : 1 以上。

部分資料不包含重複資料。例如、Oracle 區塊包含資料庫的全域唯一標頭、以及近乎唯一的標尾。因此、Oracle 資料庫的重複資料刪除功能很少能節省 1% 以上的成本。使用 MS SQL 資料庫進行重複資料刪除的效果稍微好一些、但區塊層級的獨特中繼資料仍是一項限制。

在某些情況下、使用 16KB 和大型區塊大小的資料庫可節省高達 15% 的空間。每個區塊的初始 4KB 包含全域唯一的標頭、最後 4KB 區塊則包含近乎獨特的標尾。內部區塊是重複資料刪除的候選項目、但實際上、這幾乎完全歸功於重複資料刪除零位資料。

許多競爭陣列都宣稱能夠根據資料庫複製多次的假設來刪除重複的資料庫。在這方面，也可以使用 NetApp 重複資料刪除技術，但 ONTAP 提供更好的選擇：NetApp FlexClone 技術。最終結果相同；資料庫的多個複本會建立共用大部分基礎實體區塊。使用 FlexClone 比花時間複製資料庫檔案然後刪除複製檔案更有效率。實際上，它是不重複數據刪除，而不是重複數據刪除，因為從一開始就不會創建重複數據。

### 效率與精簡配置

效率功能是精簡配置的形式。例如、佔用 100GB 磁碟區的 100GB LUN 可能會壓縮至 50GB。由於磁碟區仍為 100GB、因此尚未實現實際節省。必須先縮小磁碟區的大小、才能將儲存的空間用於系統的其他位置。如果稍後變更為 100GB LUN、則資料的壓縮性會降低、LUN 的大小會增加、而且磁碟區可能會填滿。

強烈建議採用精簡配置、因為它可以簡化管理、同時大幅改善可用容量、並節省相關成本。原因很簡單：資料庫環境通常包含大量的空空間、大量的磁碟區和 LUN、以及可壓縮的資料。如果磁碟區和 LUN 的儲存空間有一天 100% 滿、而且包含 100% 不可壓縮的資料、則大量資源配置會導致保留空間。這種情況不太可能發生。精簡配置可回收空間並在其他地方使用、並可讓容量管理以儲存系統本身為基礎、而非許多較小的磁碟區和 LUN。

有些客戶偏好針對特定工作負載使用完整資源配置、或是根據既定的營運和採購實務做法。

- 注意：\* 如果磁碟區是完整配置的磁碟區、則必須小心將該磁碟區的所有效率功能完全停用、包括使用解壓縮和移除重複資料刪除 `sis undo` 命令。Volume 不應出現在 `volume efficiency show` 輸出。如果有、則磁碟區仍會部分設定為使用效率功能。因此、覆寫保證會以不同的方式運作、這會增加組態超視導致磁碟區意外用盡空間的機會、進而導致資料庫 I/O 錯誤。

### 效率最佳實務做法

NetApp 建議：

#### AFF 預設值

在 All Flash AFF 系統上執行的 ONTAP 上建立的磁碟區會自動精簡佈建、並啟用所有的內嵌效率功能。雖然資料庫通常無法從重複資料刪除中獲益、而且可能包含不可壓縮的資料、但預設設定仍適用於幾乎所有的工作負載。ONTAP 旨在有效處理所有類型的資料和 I/O 模式、無論是否能節省成本。只有在充分瞭解理由且有偏離的好處時、才應變更預設值。

#### 一般建議

- 如果磁碟區和（或）LUN 並未精簡配置、您必須停用所有效率設定、因為使用這些功能並不會節省成本、而將複雜資源配置與啟用空間效率的組合、可能會導致非預期的行為、包括空間不足的錯誤。

- 如果資料不需要覆寫、例如備份或資料庫交易記錄檔、您可以在冷卻週期較短的情況下啟用 TSSE、以達到更高的效率。
- 某些檔案可能包含大量不可壓縮的資料、例如、當檔案的應用程式層級已啟用壓縮時、就會進行加密。如果上述任何情況屬實、請考慮停用壓縮、以便在包含可壓縮資料的其他磁碟區上執行更有效率的作業。
- 請勿將 32KB 壓縮和重複資料刪除同時用於資料庫備份。請參閱一節 [\[自適應壓縮\]](#) 以取得詳細資料。

## 使用 Oracle 資料庫進行精簡配置

簡化 Oracle 資料庫的資源配置需要仔細規劃、因為結果是在儲存系統上設定的空間比實際可用的空間更多。這項工作非常值得、因為正確完成後、可大幅節省成本、並改善管理能力。

精簡配置有多種形式、是 ONTAP 為企業應用程式環境提供的許多功能不可或缺的一部分。精簡配置也與效率技術密切相關、原因相同：效率功能可儲存比儲存系統技術更多的邏輯資料。

幾乎任何快照的使用都需要精簡配置。例如、NetApp 儲存設備上的典型 10TB 資料庫、包含約 30 天的快照。這種配置可在作用中的檔案系統中看到大約 10TB 的資料、而在快照中則有 300TB 的資料。總儲存容量為 310TB、通常位於大約 12TB 到 15TB 的空間上。作用中資料庫消耗 10TB、而其餘 300TB 的資料僅需要 2TB 到 5TB 的空間、因為只會儲存對原始資料所做的變更。

複製也是精簡配置的範例。一位主要 NetApp 客戶建立 40 個 80 TB 資料庫的複本、供開發人員使用。如果使用這些複本的 40 位開發人員都在每個資料檔案中覆寫每個區塊、則需要超過 3.2PB 的儲存空間。實際上、週轉率較低、而且由於磁碟機上只儲存變更、因此集體空間需求接近 40 TB。

### 空間管理

由於資料變更率可能會意外增加、因此精簡配置應用程式環境時必須謹慎處理。例如、如果資料庫資料表重新編製索引、或是將大規模的修補套用至 VMware 來賓系統、快照所造成的空間使用量就會迅速增加。錯誤的備份可能會在很短的時間內寫入大量資料。最後、如果檔案系統在非預期的情況下用盡可用空間、可能很難恢復某些應用程式。

幸運的是、這些風險可以透過仔細設定來解決 volume-autogrow 和 snapshot-autodelete 原則。如同其名稱所暗示、這些選項可讓使用者建立原則、以自動清除快照佔用的空間、或是增加磁碟區以容納額外資料。有許多選項可供選擇、需求因客戶而異。

請參閱 "[邏輯儲存管理文件](#)" 以完整討論這些功能。

### 部分保留

「部分保留」是指磁碟區中 LUN 在空間效率方面的行為。選項 fractional-reserve 設為 100%、磁碟區中的所有資料在任何資料模式下都能達到 100% 的營業額、而不會耗盡磁碟區上的空間。

例如、假設資料庫位於 1TB 磁碟區中的單一 250GB LUN 上。建立快照將會立即在磁碟區中保留額外的 250GB 空間、以確保磁碟區不會因任何原因而用盡空間。使用分數保留通常是浪費、因為資料庫磁碟區中的每個位元組都不太可能需要覆寫。沒有理由為永遠不會發生的事件預留空間。不過、如果客戶無法監控儲存系統中的空間使用量、而且必須確定空間永遠不會用盡、則使用快照需要 100% 的部分保留。

### 壓縮與重複資料刪除

壓縮和重複資料刪除都是精簡配置的形式。例如、50TB 的資料佔用空間可能會壓縮至 30TB、因此可節省 20TB。為了讓壓縮產生任何效益、其中某些 20TB 必須用於其他資料、或是儲存系統必須購買的容量低於

50TB。因此、儲存的資料量比儲存系統技術上的資料量還多。從資料觀點來看、即使磁碟機僅佔用 30TB、資料仍有 50TB。

資料集的可壓縮性隨時都會變更、這會導致實際空間的使用量增加。這種使用量的增加意味著、在監控和使用方面、必須像其他形式的精簡配置一樣管理壓縮 `volume-autogrow` 和 `snapshot-autodelete`。

有關壓縮和重複資料刪除的詳細資訊、請參閱連結：[efficiency.html](#) 一節

#### 壓縮與部分保留

壓縮是一種精簡配置形式。部分保留會影響壓縮的使用、並附有一個重要附註；在建立快照之前、會保留空間。通常、只有存在快照時、部分保留才會很重要。如果沒有快照、則部分保留並不重要。這不是壓縮的情況。如果在具有壓縮功能的磁碟區上建立 LUN、ONTAP 會保留空間以容納快照。在組態期間、這種行為可能會令人困惑、但這是預期的。

舉例來說、請考慮使用 5GB LUN 的 10GB 磁碟區、該磁碟區已壓縮至 2.5GB、但沒有快照。請考慮以下兩種情況：

- 分數保留 = 100 會導致 7.5 GB 使用率
- 部分保留量 = 0 會導致使用率為 2.5GB

第一個案例包括目前資料使用 2.5 GB 的空間、以及 5 GB 的空間、可在預期使用快照時、讓來源的營業額達到 100%。第二個案例不會保留額外空間。

雖然這種情況可能令人困惑、但實際上並不可能發生。壓縮意味著精簡配置、而 LUN 環境中的精簡配置則需要部分保留。壓縮資料永遠可以被無法壓縮的東西覆寫、這表示必須精簡配置磁碟區以進行壓縮、以節省任何成本。



- NetApp 建議 \* 下列保留組態：
- 設定 `fractional-reserve` 與一起進行基本容量監控時為 0 `volume-autogrow` 和 `snapshot-autodelete`。
- 設定 `fractional-reserve` 如果沒有監控能力、或在任何情況下都無法排放空間、則達到 100。

#### 可用空間和 LVM 空間分配

在檔案系統環境中自動精簡配置作用中 LUN 的效率、可能會隨著資料刪除而隨時間而喪失。除非刪除的資料會以零覆寫（另請參閱 "[ASMRU](#)" 或是隨著修剪 / 取消對應空間回收而釋放空間、「清除」資料會佔用檔案系統中越來越多的未分配空白空間。此外、在許多資料庫環境中、主動式 LUN 的精簡配置功能有限、因為資料檔案會在建立時初始化為全尺寸。

仔細規劃 LVM 組態可提高效率、並將儲存資源配置和 LUN 調整大小的需求降至最低。當使用 Veritas VxVM 或 Oracle ASM 等 LVM 時、基礎 LUN 會分割成僅在需要時才使用的範圍。例如、如果資料集的大小從 2TB 開始、但隨時間而成長至 10TB、則此資料集可放置在配置在 LVM 磁碟群組中的 10TB 精簡配置 LUN 上。在建立時、它只會佔用 2TB 的空間、而且只會在分配範圍以容納資料成長時、才會要求額外的空間。只要監控空間、此程序就安全無虞。

#### Oracle 資料庫和 ONTAP 控制器容錯移轉 / 切換

需要瞭解儲存設備接管和切入功能、才能確保 Oracle 資料庫作業不會因這些作業而中斷。

此外、如果不正確使用、接管和切入作業所使用的引數可能會影響資料完整性。

- 在正常情況下、傳入的寫入資料會同步鏡射至指定的控制器、以供其合作夥伴使用。在 NetApp MetroCluster 環境中、寫入也會鏡射到遠端控制器。除非寫入儲存在所有位置的非揮發性媒體中、否則不會對主機應用程式進行確認。
- 儲存寫入資料的媒體稱為非揮發性記憶體或 NVMEM。它有時也稱為非揮發性隨機存取記憶體（NVRAM）、雖然它是日誌、但仍可視為寫入快取。在正常作業中、不會讀取來自 NVMEM 的資料；只有在軟體或硬體故障時、才會用來保護資料。當資料寫入磁碟機時、資料會從系統的 RAM 傳輸、而非從 NVMEM 傳輸。
- 在接管作業期間、高可用度（HA）配對中的一個節點會接管其合作夥伴的作業。切換基本上相同、但適用於遠端節點接管本機節點功能的 MetroCluster 組態。

在例行維護作業期間、儲存設備接管或切換作業應該是透明的、但網路路徑變更時、操作可能會短暫暫停。然而、網路連線可能很複雜、而且容易出錯、因此 NetApp 強烈建議您在將儲存系統投入生產之前、先徹底測試接管和轉換作業。這樣做是確保正確設定所有網路路徑的唯一方法。在 SAN 環境中、請仔細檢查命令的輸出 `sanlun lun show -p` 以確保所有預期的主要和次要路徑都可用。

發出強制接管或關機時、請務必小心。使用這些選項強制變更儲存組態、表示會忽略擁有磁碟機的控制器狀態、而替代節點則強制控制磁碟機。不正確地強制接管可能會導致資料遺失或毀損。這是因為強制接管或變更會捨棄 NVMEM 的內容。在接管或切換完成後、資料遺失表示儲存在磁碟機上的資料可能會從資料庫的角度還原到稍微舊的狀態。

很少需要強制接管正常的 HA 配對。在幾乎所有故障情況下、節點都會關機並通知合作夥伴、以便進行自動容錯移轉。有些邊緣情況、例如發生滾動故障、節點之間的互連中斷、然後一個控制器遺失、需要強制接管。在這種情況下、節點之間的鏡像會在控制器故障之前遺失、這表示當機的控制器將不再擁有正在進行的寫入複本。然後需要強制接管、這表示資料可能會遺失。

同樣的邏輯也適用於 MetroCluster 轉換。在正常情況下、可進行的作業幾乎透明化。然而、災難可能會導致仍在運作的站台和災難站台之間的連線中斷。從仍在運作的站台觀點來看、問題可能只是站台之間的連線中斷、而原始站台可能仍在處理資料。如果節點無法驗證主控制器的狀態、則只能強制進行移轉。

- NetApp 建議 \* 採取下列預防措施：
- 請務必小心、避免意外強制接管或切入。一般而言、不應強制、強制變更可能會導致資料遺失。
- 如果需要強制接管或移除、請確定應用程式已關機、所有檔案系統均已卸除、且邏輯 Volume Manager（LVM）磁碟區群組已移除。必須卸載 ASM 磁碟群組。
- 在強制 MetroCluster 轉換的情況下、請將故障節點從所有仍在運作的儲存資源中隔離。如需詳細資訊、請參閱 MetroCluster 管理與災難恢復指南、以取得相關版本的 ONTAP。



## MetroCluster 和多個集合體

MetroCluster 是一種同步複寫技術、可在連線中斷時切換至非同步模式。這是客戶最常提出的要求、因為保證同步複寫意味著站台連線中斷會導致資料庫 I/O 完全停止、使資料庫停止運作。

透過 MetroCluster、集合體在連線恢復後會快速重新同步。與其他儲存技術不同、MetroCluster 在站台故障後絕不應要求完整的重新鏡射。只能運送差異變更。

在跨集合體的資料集中、在循環災難案例中需要額外的資料恢復步驟、風險很小。具體而言、如果（a）站台之間的連線中斷、（b）連線恢復、（c）集合體會達到某種狀態、其中有些是同步的、有些則不是同步的、然後（d）主站台會遺失、結果是無法運作的站台、而集合體彼此之間不會同步。如果發生這種情況、資料集的某些部分會彼此同步、因此無法在沒有恢復的情況下啟動應用程式、資料庫或資料存放區。如果資料集橫跨整

個集合體、NetApp 強烈建議您利用快照式備份、搭配眾多可用工具之一、在這種不尋常的情況下驗證快速的恢復性。

## 資料庫組態

### Oracle 資料庫區塊大小

ONTAP 內部使用可變的區塊大小、這表示 Oracle 資料庫可以設定任何所需的區塊大小。不過、檔案系統區塊大小可能會影響效能、在某些情況下、較大的重做區塊大小可能會改善效能。

#### 資料檔案區塊大小

部分作業系統提供多種檔案系統區塊大小選擇。對於支援 Oracle 資料檔案檔案的檔案系統、使用壓縮時區塊大小應為 8KB。不需要壓縮時、可以使用 8KB 或 4KB 的區塊大小。

如果資料檔案放在具有 512 位元組區塊的檔案系統上、則可能會有未對齊的檔案。根據 NetApp 建議、LUN 和檔案系統可能已正確對齊、但檔案 I/O 可能未對齊。這種錯誤的調整會導致嚴重的效能問題。

支援重做記錄的檔案系統必須使用重做區塊大小的倍數。這通常需要重做記錄檔系統和重做記錄本身都使用 512 位元組的區塊大小。

#### 重做區塊大小

重做率極高時、4KB 區塊大小可能會執行得更好、因為重做率較高、可在較少且更有效率的作業中執行 I/O。如果重做速率大於 50Mbps、請考慮測試 4KB 區塊大小。

在具有 4KB 區塊大小和許多極小型交易的檔案系統上、使用 512 位元組區塊大小的重做記錄檔來識別資料庫中的一些客戶問題。將多個 512 位元組的變更套用到單一 4KB 檔案系統區塊所涉及的額外負荷、導致效能問題、而這些問題已透過將檔案系統變更為使用 512 位元組的區塊大小來解決。



\* NetApp 建議 \* 除非相關客戶支援或專業服務組織建議您變更重做區塊大小、否則請勿變更重做區塊大小、否則變更將以正式產品文件為基礎。

### Oracle 資料庫參數：DB\_FILE\_Multifblock\_read\_count

◦ `db_file_multiblock_read_count` 參數控制 Oracle 在連續 I/O 期間讀取為單一作業的 Oracle 資料庫區塊數量上限

不過、此參數不會影響 Oracle 在任何及所有讀取作業期間讀取的區塊數、也不會影響隨機 I/O 只有連續 I/O 的區塊大小會受到影響。

Oracle 建議使用者不要設定此參數。如此可讓資料庫軟體自動設定最佳值。這通常表示此參數設為可產生 1MB I/O 大小的值。例如、1MB 讀取 8KB 區塊需要 128 個區塊才能讀取、因此此參數的預設值為 128。

NetApp 在客戶站台上觀察到的大多數資料庫效能問題、都涉及此參數的設定不正確。使用 Oracle 版本 8 和 9 變更此值的理由是正確的。因此、參數可能會在不知情的情況下出現在 `init.ora` 檔案、因為資料庫已就地升級至 Oracle 10 及更新版本。傳統設定為 8 或 16、而預設值為 128、會大幅損害連續 I/O 效能。



\* NetApp 建議 \* 設定 `db_file_multiblock_read_count` 參數不應出現在 `init.ora` 檔案：NetApp 從未遇到過變更此參數可改善效能的情況、但在許多情況下、它會對連續 I/O 處理量造成明顯損害。

## Oracle 資料庫參數：`filesystemio_options`

### Oracle 初始化參數 `filesystemio_options` 控制非同步和直接 I/O 的使用

與一般的看法相反、非同步和直接 I/O 並不相互排斥。NetApp 發現、在客戶環境中、此參數經常設定錯誤、而這種錯誤設定直接導致許多效能問題。

非同步 I/O 表示 Oracle I/O 作業可以平行化。在各種作業系統上均可使用非同步 I/O 之前、使用者已設定數個 `dbwriter` 程序、並變更伺服器程序組態。透過非同步 I/O、作業系統本身就能以高效率且平行的方式代表資料庫軟體執行 I/O。此程序不會讓資料面臨風險、而且關鍵作業（例如 Oracle 重做記錄）仍會同步執行。

直接 I/O 會略過作業系統緩衝區快取。UNIX 系統上的 I/O 通常會流經作業系統緩衝區快取。這對不維護內部快取的應用程式很有用、但 Oracle 在 SGA 中擁有自己的緩衝區快取。在幾乎所有情況下、最好是啟用直接 I/O 並將伺服器 RAM 分配給 SGA、而非仰賴作業系統緩衝區快取。Oracle SGA 更有效率地使用記憶體。此外、當 I/O 流經作業系統緩衝區時、它會受到額外處理、因此會增加延遲。當低延遲是關鍵需求時、在大量寫入 I/O 時、延遲特別明顯。

的選項 `filesystemio_options` 是：

- \* 非同步 \*。Oracle 將 I/O 要求提交給作業系統以進行處理。此程序可讓 Oracle 執行其他工作、而非等待 I/O 完成、進而增加 I/O 平行化。
- **directio**。Oracle 直接針對實體檔案執行 I/O、而非透過主機作業系統快取來路由 I/O。
- \* 無 \*。Oracle 使用同步和緩衝 I/O 在此組態中、選擇共享與專用伺服器程序與 `dbWriters` 數量更為重要。
- **setall**。Oracle 同時使用非同步和直接 I/O 在幾乎所有情況下、都是使用 `setall` 是最佳的。



◦ `filesystemio_options` 參數在 DNFS 和 ASM 環境中無效。使用 DNFS 或 ASM 時、會自動同時使用非同步和直接 I/O

有些客戶過去曾遇到非同步 I/O 問題、尤其是先前的 Red Hat Enterprise Linux 4 (RHEL4) 版本。網際網路上有些過時的建議仍建議避免非同步 IO、因為資訊過時。在所有目前的作業系統上、非同步 I/O 都是穩定的。沒有理由停用它、作業系統沒有已知的錯誤。

如果資料庫已使用緩衝 I/O、則直接 I/O 的交換器也可能需要變更 SGA 大小。停用緩衝 I/O 可消除主機作業系統快取為資料庫提供的效能優勢。將 RAM 新增回 SGA 可修復此問題。最終結果應該是 I/O 效能的改善。

雖然 Oracle SGA 使用 RAM 幾乎比使用 OS 緩衝區快取更好、但可能無法判斷最佳值。例如、最好在資料庫伺服器上使用具有極小型 SGA 大小的緩衝 I/O、其中有許多間歇性作用中的 Oracle 執行個體。這種配置可讓所有執行中的資料庫執行個體靈活使用作業系統上的剩餘可用 RAM。這是非常不尋常的情況、但在某些客戶據點已發現這種情況。



\* NetApp 建議 \* 設定 `filesystemio_options` 至 `'setall'` 但請注意、在某些情況下、遺失主機緩衝區快取可能需要增加 Oracle SGA。

## Oracle Real Application Clusters ( RAC ) 逾時

Oracle RAC 是一款叢集軟體產品、內含多種類型的內部活動訊號處理程序、可監控叢集的健全狀況。



中的資訊 "遺失計數" 本節包含使用網路儲存設備的 Oracle RAC 環境的重要資訊、在許多情況下、需要變更預設 Oracle RAC 設定、以確保 RAC 叢集在網路路徑變更和儲存設備容錯移轉 / 切換作業之後仍能順利運作。

### 磁碟逾時

主要儲存相關 RAC 參數為 `disktimeout`。此參數控制投票檔案 I/O 必須完成的臨界值。如果是 `disktimeout` 超過參數、RAC 節點就會從叢集中移出。此參數的預設值為 200。此值應足以用於標準儲存設備接管和恢復程序。

NetApp 強烈建議您在將 RAC 組態投入生產之前、先徹底測試這些組態、因為許多因素會影響接管或恢復作業。除了完成儲存容錯移轉所需的時間之外、連結集合化控制傳輸協定 (LACP) 變更也需要額外的時間才能傳播。此外、SAN 多重路徑軟體必須偵測 I/O 逾時、然後在替代路徑上重試。如果資料庫處於極活躍狀態、則必須在處理投票磁碟 I/O 之前、先佇列並重新嘗試大量 I/O。

如果無法執行實際的儲存接管或恢復、則可以在資料庫伺服器上進行纜線拉出測試來模擬影響。



- NetApp 建議 \* 下列事項：
- 離開 `disktimeout` 預設值為 200 的參數。
- 務必徹底測試 RAC 組態。

### 遺失計數

。 `misscount` 參數通常只會影響 RAC 節點之間的網路心跳。預設值為 30 秒。如果網格二進位檔位於儲存陣列上、或作業系統開機磁碟機不是本機磁碟機、此參數可能會變得很重要。這包括 FC SAN 上具有開機磁碟機的主機、NFS 開機作業系統、以及位於虛擬化資料存放區 (例如 VMDK 檔案) 上的開機磁碟機。

如果因儲存接管或恢復而中斷開機磁碟機的存取、網格二進位位置或整個作業系統可能會暫時停止運作。ONTAP 完成儲存作業所需的時間、以及作業系統變更路徑和恢復 I/O 所需的時間、可能會超過 `misscount` 臨界值。因此、節點會在連線到開機 LUN 或網格二進位檔恢復後立即停止。在大多數情況下、會發生遷離和後續重新開機、而不會出現記錄訊息來指出重新開機的原因。並非所有組態都會受到影響、因此請在 RAC 環境中測試任何 SAN 開機、NFS 開機或資料存放區型主機、以便在與開機磁碟機的通訊中斷時、RAC 保持穩定。

非本機開機磁碟機或非本機檔案系統代管的情況 `grid` 二進位檔案 `misscount` 需要變更以符合 `disktimeout`。如果變更此參數、請進行進一步測試、以識別對 RAC 行為的任何影響、例如節點容錯移轉時間。





- NetApp 建議 \* 下列事項：
- 離開 `misscount` 參數的預設值為 30、除非符合下列其中一項條件：
  - `grid` 二進位檔位於網路附加磁碟機上、包括 NFS、iSCSI、FC 和資料存放區型磁碟機。
  - 作業系統是 SAN 開機。
- 在這種情況下、請評估網路中斷對 OS 或的存取造成的影響 `GRID_HOME` 檔案系統。在某些情況下、這類中斷會導致 Oracle RAC 精靈停止運作、進而導致 `misscount` 根據的逾時和遷離。逾時預設為 27 秒、即的值 `misscount` 減號 `reboottime`。在這種情況下、請增加 `misscount` 200 比對 `disktimeout`。

## 主機組態

### 使用 IBM AIX 的 Oracle 資料庫

IBM AIX 與 ONTAP 上 Oracle 資料庫的組態主題。

#### 並行 I/O

若要在 IBM AIX 上達到最佳效能、必須同時使用 I/O 如果沒有並行 I/O、效能限制很可能是因為 AIX 執行序列化的原子 I/O、這會產生重大的負荷。

NetApp 最初建議使用 `cio` 掛載選項可強制在檔案系統上使用並行 I/O、但此程序有缺點、不再需要。自從推出 AIX 5.2 和 Oracle 10gR1 之後、AIX 上的 Oracle 就可以開啟個別檔案來同時執行 IO、而不是強制在整個檔案系統上同時執行 I/O。

啟用並行 I/O 的最佳方法是設定 `init.ora` 參數 `filesystemio_options` 至 `setall`。這樣做可讓 Oracle 開啟特定檔案、以與並行 I/O 搭配使用

使用 `cio` 作為掛載選項，強制使用並行 I/O，這可能會產生負面影響。例如、強制並行 I/O 會停用檔案系統上的預先讀取、這可能會損害 Oracle 資料庫軟體以外的 I/O 效能、例如複製檔案和執行磁帶備份。此外、Oracle GoldenGate 和 SAP BR\* Tools 等產品與使用不相容 `cio` 裝載選項搭配特定版本的 Oracle。



- NetApp 建議 \* 下列事項：
- 請勿使用 `cio` 檔案系統層級的掛載選項。而是透過使用來啟用並行 I/O `filesystemio_options=setall`。
- 請僅使用 `cio` 如果無法設定掛載選項、則應選擇掛載選項 `filesystemio_options=setall`。

### AIX NFS 裝載選項

下表列出 Oracle 單一執行個體資料庫的 AIX NFS 掛載選項。

檔案類型	掛載選項
ADR 首頁	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsiz=262144,wsiz=262144</code>

檔案類型	掛載選項
控制檔 資料檔案 重作記錄	rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144
ORACLE_HOME	rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,intr

下表列出 RAC 的 AIX NFS 掛載選項。

檔案類型	掛載選項
ADR 首頁	rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144
控制檔 資料檔案 重作記錄	rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,nointr,noac
CRS/Voting	rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,nointr,noac
專屬 ORACLE_HOME	rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144
共享 ORACLE_HOME	rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,nointr

單一執行個體與 RAC 掛載選項之間的主要差異在於新增 noac 至掛載選項。此新增功能可停用主機作業系統快取、使 RAC 叢集中的所有執行個體都能一致地檢視資料狀態。

雖然使用 cio 裝載選項和 init.ora 參數 filesystemio\_options=setall 停用主機快取的效果相同、仍需使用 noac。noac 為共享的必要項目 ORACLE\_HOME 部署以促進 Oracle 密碼檔案和等檔案的一致性 spfile 參數檔。如果 RAC 叢集中的每個執行個體都有專用的 ORACLE\_HOME，則不需要此參數。

### AIX jfs/JFS2 掛載選項

下表列出 AIX jfs/JFS2 掛載選項。

檔案類型	掛載選項
ADR 首頁	預設值
控制檔 資料檔案 重作記錄	預設值
Oracle_Home	預設值

使用 AIX 之前 hdisk 任何環境中的裝置（包括資料庫）都請檢查參數 queue\_depth。此參數不是 HBA 佇列深度、而是與個別主機的 SCSI 佇列深度相關 hdisk device。Depending on how the LUNs are configured, the value for `queue\_depth` 效能可能太低。測試顯示最佳值為 64。

## 使用 HP-UX 的 Oracle 資料庫

適用於 HP-UX with ONTAP 上 Oracle 資料庫的組態主題。

### HP-UX NFS 掛載選項

下表列出單一執行個體的 HP-UX NFS 掛載選項。

檔案類型	掛載選項
ADR 首頁	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,suid</code>
控制檔 資料檔案 重作記錄	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,forcedirectio, nointr,suid</code>
ORACLE_HOME	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,suid</code>

下表列出適用於 RAC 的 HP-UX NFS 掛載選項。

檔案類型	掛載選項
ADR 首頁	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,noac,suid</code>
控制檔 資料檔案 重作記錄	<code>rw, bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,nointr,noac,forcedirectio,suid</code>
CRS/ 投票	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,nointr,noac,forcedirectio,suid</code>
專屬 ORACLE_HOME	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,suid</code>
共享 ORACLE_HOME	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,nointr,noac,suid</code>

單一執行個體與 RAC 掛載選項之間的主要差異在於新增 `noac` 和 `forcedirectio` 至掛載選項。此新增功能可停用主機作業系統快取、使 RAC 叢集中的所有執行個體都能一致地檢視資料狀態。雖然使用 `init.ora` 參數 `filesystemio_options=setall` 停用主機快取的效果相同、仍需使用 `noac` 和 `forcedirectio`。

原因 `noac` 為共享的必要項目 `ORACLE_HOME` 部署是為了促進檔案的一致性、例如 Oracle 密碼檔案和 `spfiles`。如果 RAC 叢集中的每個執行個體都有專用的 `ORACLE_HOME`，不需要此參數。

## HP-UX VxFS 掛載選項

對於託管 Oracle 二進位檔的檔案系統、請使用下列掛載選項：

```
delaylog,nodatainlog
```

對於包含資料檔案、重做記錄檔、歸檔記錄檔和控制檔的檔案系統、若 HP-UX 版本不支援並行 I/O、請使用下列掛載選項：

```
nodatainlog,mincache=direct,convosync=direct
```

支援並行 I/O（VxFS 5.0.1 及更新版本、或 ServiceGuard Storage Management Suite）時、請針對包含資料檔案、重作記錄檔、歸檔記錄檔和控制檔的檔案系統、使用這些掛載選項：

```
delaylog,cio
```



參數 `db_file_multiblock_read_count` 在 VxFS 環境中尤其重要。Oracle 建議在 Oracle 10g R1 及更新版本中保留此參數、除非另有特別指示。Oracle 8KB 區塊大小的預設值為 128。如果此參數的值強制為 16 或更少、請移除 `convosync=direct` 裝載選項、因為它可能會損害連續 I/O 效能。此步驟會損害其他效能層面、只有在的價值下才應採取 `db_file_multiblock_read_count` 必須從預設值變更。

## 使用 Linux 的 Oracle 資料庫

Linux 作業系統專屬的組態主題。

### Linux NFSv3 TCP 插槽表

TCP 插槽表是與主機匯流排介面卡（HBA）佇列深度相當的 NFSv3。這些表格可控制任何時間都可以處理的 NFS 作業數量。預設值通常為 16、這對於最佳效能而言太低。相反的問題發生在較新的 Linux 核心上、這會自動將 TCP 插槽表格限制增加到要求使 NFS 伺服器飽和的層級。

為了達到最佳效能並避免效能問題、請調整控制 TCP 插槽表的核心參數。

執行 `sysctl -a | grep tcp.*.slot_table` 並觀察下列參數：

```
# sysctl -a | grep tcp.*.slot_table
sunrpc.tcp_max_slot_table_entries = 128
sunrpc.tcp_slot_table_entries = 128
```

所有 Linux 系統都應該包括在內 `sunrpc.tcp_slot_table_entries`、但只有部分包含在內 `sunrpc.tcp_max_slot_table_entries`。兩者都應設為 128。

## 注意

若未設定這些參數、可能會對效能造成重大影響。在某些情況下、效能會受到限制、因為 Linux 作業系統沒有發出足夠的 I/O 在其他情況下、隨著 Linux 作業系統嘗試發出的 I/O 數量超過可服務的數量、I/O 延遲也會增加。

## Linux NFS 裝載選項

下表列出單一執行個體的 Linux NFS 掛載選項。

檔案類型	掛載選項
ADR 首頁	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144</code>
控制檔 資料檔案 重作記錄	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,nointr</code>
Oracle_Home	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,nointr</code>

下表列出 RAC 的 Linux NFS 掛載選項。

檔案類型	掛載選項
ADR 首頁	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,actimeo=0</code>
控制檔 資料檔案 重作記錄	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,nointr,actimeo=0</code>
CRS/ 投票	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,nointr,noac,actimeo=0</code>
專屬 ORACLE_HOME	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144</code>
共享 ORACLE_HOME	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,nointr,actimeo=0</code>

單一執行個體與 RAC 掛載選項之間的主要差異在於新增 `actimeo=0` 至掛載選項。此新增功能可停用主機作業系統快取、使 RAC 叢集中的所有執行個體都能一致地檢視資料狀態。雖然使用 `init.ora` 參數 `filesystemio_options=setall` 停用主機快取的效果相同、仍需使用 `actimeo=0`。

原因 `actimeo=0` 為共享的必要項目 `ORACLE_HOME` 部署是為了促進檔案的一致性、例如 Oracle 密碼檔案和 `spfiles`。如果 RAC 叢集中的每個執行個體都有專用的 `ORACLE_HOME`，則不需要此參數。

一般而言、非資料庫檔案的掛載應該與單一執行個體資料檔案所使用的選項相同、不過特定應用程式可能有不同的需求。避免使用掛載選項 `noac` 和 `actimeo=0` 如果可能、因為這些選項會停用檔案系統層級的預先讀取和緩

衝處理。這可能會對擷取、轉譯和載入等程序造成嚴重的效能問題。

### 存取與 GetAttr

有些客戶指出、存取和 GetAttr 等極高層級的其他 IOPS、可能會主導他們的工作負載。在極端情況下、讀取和寫入等作業可低於總作業量的 10%。這是任何包含使用的資料庫的正常行為 `actimeo=0` 和/或 `noac` 在 Linux 上、因為這些選項會導致 Linux 作業系統持續從儲存系統重新載入檔案中繼資料。存取和 GetAttr 等作業是從資料庫環境中的 ONTAP 快取提供服務的低影響作業。它們不應被視為真正的 IOPS、例如讀寫、而會對儲存系統產生真正的需求。不過、這些其他 IOPS 確實會產生一些負載、尤其是在 RAC 環境中。若要解決這種情況、請啟用 DNFS、以略過作業系統緩衝區快取、避免這些不必要的中繼資料作業。

### Linux Direct NFS

另一個掛載選項、稱為 `nosharecache` (a) DNFS 啟用時、需要使用、(b) 來源磁碟區多次裝載於具有巢狀 NFS 裝載的單一伺服器 (c) 上。此組態主要出現在支援 SAP 應用程式的環境中。例如、NetApp 系統上的單一磁碟區可能有位於的目錄 `/vol/oracle/base` 再來一次 `/vol/oracle/home`。如果 `/vol/oracle/base` 安裝於 `/oracle` 和 `/vol/oracle/home` 安裝於 `/oracle/home`，結果是來自相同來源的巢狀 NFS 掛載。

作業系統可以偵測到這個事實 `/oracle` 和 `/oracle/home` 位於同一個磁碟區、即相同的來源檔案系統。作業系統接著會使用相同的裝置控制代碼來存取資料。這樣做可以改善 OS 快取和某些其他作業的使用、但會干擾 DNFS。如果 DNFS 必須存取檔案、例如 `spfile`、開啟 `/oracle/home`，它可能會錯誤地嘗試使用錯誤的資料路徑。結果是 I/O 作業失敗。在這些組態中、新增 `nosharecache` 裝載選項至任何與該主機上其他 NFS 檔案系統共用來源 FlexVol 磁碟區的 NFS 檔案系統。這樣做會強制 Linux 作業系統為該檔案系統分配獨立的裝置控制代碼。

### Linux Direct NFS 和 Oracle RAC

使用 DNFS 對 Linux 作業系統上的 Oracle RAC 有特殊的效能優勢、因為 Linux 沒有強制直接 I/O 的方法、而 RAC 需要此方法才能在節點之間保持一致。因應措施是 Linux 需要使用 `actimeo=0` 掛載選項、會使檔案資料從作業系統快取立即過期。此選項會強制 Linux NFS 用戶端持續重新讀取屬性資料、進而損害延遲並增加儲存控制器的負載。

啟用 DNFS 會略過主機 NFS 用戶端、避免此損害。多家客戶在啟用 DNFS 時、報告 RAC 叢集的效能大幅提升、ONTAP 負載大幅降低 (特別是其他 IOPS)。

### Linux Direct NFS 和 `oranzstab` 檔案

在 Linux 上搭配多重路徑選項使用 DNFS 時、必須使用多個子網路。在其他作業系統上、可使用建立多個 DNFS 通道 `LOCAL` 和 `DONTROUTE` 可在單一子網路上設定多個 DNFS 通道的選項。不過、這在 Linux 上無法正常運作、因此可能會產生非預期的效能問題。在 Linux 中、用於 DNFS 流量的每個 NIC 都必須位於不同的子網路上。

### I/O 排程器

Linux 核心可讓您以低層級控制 I/O 排程封鎖裝置的方式。Linux 各版本的預設值差異極大。測試顯示、截止日期通常會提供最佳結果、但有時 `NOOP` 會稍微好一點。效能差異最小、但如果需要從資料庫組態擷取最大可能效能、請測試這兩個選項。在許多組態中、`CFQ` 是預設值、而且已證明資料庫工作負載的效能有重大問題。

如需設定 I/O 排程器的指示、請參閱相關的 Linux 廠商文件。

## 多重路徑

部分客戶在網路中斷期間遭遇當機、因為多重路徑常駐程式未在其系統上執行。在最新版本的 Linux 上、作業系統的安裝程序和多重路徑常駐程式可能會讓這些作業系統容易受到此問題的影響。套件已正確安裝、但未設定為在重新開機後自動啟動。

例如、RHEL5.5 上多重路徑常駐程式的預設值可能如下所示：

```
[root@host1 iscsi]# chkconfig --list | grep multipath
multipathd      0:off  1:off  2:off  3:off  4:off  5:off  6:off
```

您可以使用下列命令來修正此問題：

```
[root@host1 iscsi]# chkconfig multipathd on
[root@host1 iscsi]# chkconfig --list | grep multipath
multipathd      0:off  1:off  2:on   3:on   4:on   5:on   6:off
```

## ASM 鏡像

ASM鏡射可能需要變更Linux多重路徑設定、以允許ASM辨識問題並切換至其他故障群組。大部分關於「不完整」的ASM組態ONTAP 都使用外部備援、這表示資料保護是由外部陣列提供、而ASM不會鏡射資料。某些站台使用具有一般備援的ASM來提供雙向鏡像、通常是跨不同站台。

中顯示的 Linux 設定 "[NetApp 主機公用程式文件](#)" 包含會導致 I/O 無限期佇列的多重路徑參數這表示 LUN 裝置上沒有作用中路徑的 I/O 會在 I/O 完成所需的時間內等待。這通常是很理想的做法、因為 Linux 主機會在 SAN 路徑變更完成、FC 交換器重新開機或儲存系統完成容錯移轉所需的時間內等待。

這種不受限制的佇列行為會導致 ASM 鏡像發生問題、因為 ASM 必須收到 I/O 故障、才能在替代 LUN 上重試 I/O。

在 Linux 中設定下列參數 `multipath.conf` 用於 ASM 鏡像的 ASM LUN 檔案：

```
polling_interval 5
no_path_retry 24
```

這些設定會為 ASM 裝置建立 120 秒的逾時。逾時會計算為 `polling_interval * no_path_retry` 秒。在某些情況下可能需要調整確切的值、但 120 秒的逾時時間應足以滿足大部分的使用需求。具體而言、120 秒的時間應該能讓控制器接管或恢復、而不會產生 I/O 錯誤、導致故障群組離線。

較低 `no_path_retry` 此值可縮短 ASM 切換至替代故障群組所需的時間、但這也會增加在維護活動（例如控制器接管）期間不必要的容錯移轉風險。仔細監控 ASM 鏡像狀態、即可降低風險。如果發生不必要的容錯移轉、只要執行重新同步的速度相對較快、鏡像就能快速重新同步。如需更多資訊、請參閱 ASM Fast Mirror Resync 上的 Oracle 說明文件、以瞭解所使用的 Oracle 軟體版本。

## Linux xfs、ext3 和 ext4 掛載選項



\* NetApp 建議 \* 使用預設掛載選項。

## 使用 ASMLib/AFD 的 Oracle 資料庫（ASM 篩選器驅動程式）

### 使用 AFD 和 ASMLib 的 Linux 作業系統專屬組態主題

#### ASMLib 區塊大小

ASMLib 是選用的 ASM 管理程式庫和相關公用程式。其主要值是將 LUN 或 NFS 型檔案標記為具有人類可讀標籤的 ASM 資源。

ASMLib 的最新版本會偵測稱為每個實體區塊指數（LBPPBE）的邏輯區塊的 LUN 參數。ONTAP SCSI 目標直到最近才回報此值。現在會傳回一個值、表示偏好 4KB 區塊大小。這不是區塊大小的定義、但它是使用 LBPPBE 的任何應用程式的提示、可能會更有效率地處理特定大小的 I/O。不過、ASMLib 會將 LBPPBE 解譯為區塊大小、並在建立 ASM 裝置時持續標記 ASM 標頭。

此程序可能會以多種方式造成升級和移轉問題、全部是因為無法在同一個 ASM 磁碟群組中混合使用不同區塊大小的 ASMLib 裝置。

例如、較舊的陣列通常回報 LBPPBE 值為 0、或根本沒有回報此值。ASMLib 會將此解譯為 512 位元組的區塊大小。較新的陣列會被解譯為具有 4KB 區塊大小。無法在同一個 ASM 磁碟群組中混合使用 512 位元組和 4KB 的裝置。這樣做會阻止用戶使用兩個陣列中的 LUN 或使用 ASM 作為遷移工具來增加 ASM 磁盤組的大小。在其他情況下、RMAN 可能不允許在具有 512 位元組區塊大小的 ASM 磁碟群組和具有 4KB 區塊大小的 ASM 磁碟群組之間複製檔案。

首選的解決方案是修補 ASMLib。Oracle 錯誤 ID 為 13999609、而 Oracle 錯誤 ID 則存在於 oracleas-support-2.1.8-1 及更高版本中。此修補程式可讓使用者設定參數 `ORACLEASM_USE_LOGICAL_BLOCK_SIZE` 至 `true` 在中 `/etc/sysconfig/oracleasm` 組態檔。這樣做會阻止 ASMLib 使用 LBPPBE 參數、這表示新陣列上的 LUN 現在會被辨識為 512 位元組區塊裝置。

此選項不會變更先前由 ASMLib 戳記的 LUN 區塊大小。例如、如果具有 512 位元組區塊的 ASM 磁碟群組必須移轉至回報 4KB 區塊的新儲存系統、則選項



`ORACLEASM_USE_LOGICAL_BLOCK_SIZE` 必須先設定、才能使用 ASMLib 標記新的 LUN。如果裝置已被 `oracleasm` 戳記、則必須先重新格式化、然後再重新設定新的區塊大小。首先、請使用取消設定裝置 `oracleasm deletedisk`、然後使用清除裝置的前 1GB `dd if=/dev/zero of=/dev/mapper/device bs=1048576 count=1024`。最後、如果裝置先前已分割、請使用 `kpartx` 命令移除過時的分割區、或只是重新開機作業系統。

如果無法修補 ASMLib、可以從組態中移除 ASMLib。這項變更會造成中斷、需要在 ASM 磁碟上加蓋戳記、並確定 `asm_diskstring` 參數設定正確。不過、這項變更並不需要移轉資料。

#### ASM Filter Drive（AFD）區塊大小

AFD 是選用的 ASM 管理程式庫、正在取代 ASMLib。從儲存角度來看、它與 ASMLib 非常類似、但它還包含其他功能、例如能夠封鎖非 Oracle I/O、以降低使用者或應用程式錯誤可能毀損資料的機會。

#### 裝置區塊大小

如同 ASMLib、AFD 也會讀取 LUN 參數每個實體區塊指數（LBPPBE）的邏輯區塊、並依預設使用實體區塊大小、而非邏輯區塊大小。

如果將 AFD 新增至現有組態、而 ASM 裝置已格式化為 512 位元組區塊裝置、則可能會造成問題。AFD 驅動程式會將 LUN 辨識為 4K 裝置、而 ASM 標籤與實體裝置之間的不符將會妨礙存取。同樣地、移轉也會受到影響、因為無法在同一個 ASM 磁碟群組中混合使用 512 位元組和 4KB 的裝置。這樣做會阻止用戶使用兩個陣列中的



LUN 或使用 ASM 作為遷移工具來增加 ASM 磁盤組的大小。在其他情況下、RMAN 可能不允許在具有 512 位元組區塊大小的 ASM 磁碟群組和具有 4KB 區塊大小的 ASM 磁碟群組之間複製檔案。

解決方案很簡單 - AFD 包含一個參數、可控制它是否使用邏輯區塊或實體區塊大小。這是影響系統上所有裝置的全域參數。若要強制 AFD 使用邏輯區塊大小、請設定 `options oracleafd oracleafd_use_logical_block_size=1` 在中 `/etc/modprobe.d/oracleafd.conf` 檔案：

多重路徑傳輸大小

最近的 Linux 核心變更會強制執行傳送至多重路徑裝置的 I/O 大小限制、而 AFD 則不遵守這些限制。然後會拒絕 I/O、導致 LUN 路徑離線。結果是無法安裝 Oracle Grid、設定 ASM 或建立資料庫。

解決方案是在 ONTAP LUN 的 `multipath.conf` 檔案中手動指定傳輸長度上限：

```
devices {
    device {
        vendor "NETAPP"
        product "LUN.*"
        max_sectors_kb 4096
    }
}
```



即使目前沒有問題、如果使用 AFD 來確保未來的 Linux 升級不會意外造成問題、也應設定此參數。

## 使用 Microsoft Windows 的 Oracle 資料庫

Microsoft Windows with ONTAP 上 Oracle 資料庫的組態主題。

### NFS

Oracle 支援搭配直接 NFS 用戶端使用 Microsoft Windows。這項功能提供 NFS 管理效益的途徑、包括跨環境檢視檔案、動態調整磁碟區大小、以及使用較便宜的 IP 傳輸協定。如需在 Microsoft Windows 上使用 DNFS 安裝及設定資料庫的詳細資訊、請參閱正式的 Oracle 文件。不存在任何特殊的最佳實務做法。

### SAN

為達到最佳壓縮效率、請確保 NTFS 檔案系統使用 8K 或更大的分配單元。使用 4K 分配單元（通常是預設）會對壓縮效率造成負面影響。

## Oracle 資料庫與 Solaris

特定於 Solaris OS 的組態主題。

### Solaris NFS 掛載選項

下表列出單一執行個體的 Solaris NFS 掛載選項。

檔案類型	掛載選項
ADR 首頁	<code>rw,bg,hard,[vers=3,vers=4.1],roto=tcp,timeo=600,rsize=262144,wsiz=262144</code>
控制檔 資料檔案 重作記錄	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,nointr,llock,suid</code>
ORACLE_HOME	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,suid</code>

的用途 `llock` 已獲證實、可移除在儲存系統上取得和釋放鎖定的相關延遲、大幅提升客戶環境的效能。在將多個伺服器設定為掛載相同檔案系統的環境中、請謹慎使用此選項、並將 Oracle 設定為掛載這些資料庫。雖然這是非常不尋常的組態、但只有少數客戶使用。如果第二次意外啟動某個執行個體、可能會因為 Oracle 無法偵測到外部伺服器上的鎖定檔案而導致資料毀損。NFS 鎖定不會提供保護；如同 NFS 第 3 版一樣、它們只是建議事項。

因為 `llock` 和 `forcedirectio` 參數是互斥的、這一點很重要 `filesystemio_options=setall` 存在於 `init.ora` 檔案就是這樣 `directio` 已使用。如果沒有此參數、就會使用主機作業系統緩衝區快取、而且效能可能會受到負面影響。

下表列出了 Solaris NFS RAC 掛載選項。

檔案類型	掛載選項
ADR 首頁	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,noac</code>
控制檔 資料檔案 重作記錄	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,nointr,noac,forcedirectio</code>
CRS/ 投票	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,nointr,noac,forcedirectio</code>
專屬 ORACLE_HOME	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,suid</code>
共享 ORACLE_HOME	<code>rw,bg,hard,[vers=3,vers=4.1],proto=tcp,timeo=600,rsize=262144,wsiz=262144,nointr,noac,suid</code>

單一執行個體與 RAC 掛載選項之間的主要差異在於新增 `noac` 和 `forcedirectio` 至掛載選項。此新增功能可停用主機作業系統快取、使 RAC 叢集中的所有執行個體都能一致地檢視資料狀態。雖然使用 `init.ora` 參數 `filesystemio_options=setall` 停用主機快取的效果相同、仍需使用 `noac` 和 `forcedirectio`。

原因 `actimeo=0` 為共享的必要項目 `ORACLE_HOME` 部署是為了促進檔案的一致性、例如 Oracle 密碼檔案和 `spfiles`。如果 RAC 叢集中的每個執行個體都有專用的 `ORACLE_HOME`，不需要此參數。

## Solaris UFS 掛載選項

NetApp 強烈建議您使用記錄掛載選項、以便在 Solaris 主機當機或 FC 連線中斷時保留資料完整性。記錄掛載選項也可保留 Snapshot 備份的使用性。

## Solaris ZFS

必須仔細安裝和設定 Solaris ZFS、才能提供最佳效能。

### mvector

Solaris 11 變更了 IT 處理大型 I/O 作業的方式、可能會在 SAN 儲存陣列上造成嚴重的效能問題。NetApp 錯誤報告 630173 「Solaris 11 ZFS 效能回歸」中詳細說明了此問題。" 解決方案是變更為的 OS 參數 `zfs_mvvector_max_size`。

以 root 執行下列命令：

```
[root@host1 ~]# echo "zfs_mvvector_max_size/W 0t131072" |mdb -kw
```

如果這項變更發生任何非預期的問題、您可以以 root 執行下列命令、輕鬆地將其還原：

```
[root@host1 ~]# echo "zfs_mvvector_max_size/W 0t1048576" |mdb -kw
```

### 核心

可靠的 ZFS 效能需要修補 Solaris 核心、以因應 LUN 對齊問題。此修正程式是隨 Solaris 10 中的修補程式 147440-19 和適用於 Solaris 11 的 SRU 10.5 一起推出的。只能將 Solaris 10 及更新版本與 ZFS 搭配使用。

### LUN 組態

若要設定 LUN、請完成下列步驟：

1. 建立類型的 LUN `solaris`。
2. 安裝所指定的適當主機公用程式套件 (Huk) ["NetApp互通性對照表工具IMT \(不含\)"](#)。
3. 請依照 Huk 中的說明進行操作、完全符合上述說明。以下概述基本步驟、但請參閱 ["最新文件"](#) 以瞭解正確的程序。
  - a. 執行 `host_config` 更新的公用程式 `sd.conf/sdd.conf` 檔案：這樣做可讓 SCSI 磁碟機正確探索 ONTAP LUN。
  - b. 請遵循所提供的指示 `host_config` 啟用多重路徑輸入 / 輸出 (MPIO) 的公用程式。
  - c. 重新開機。此步驟是必要步驟、以便在整個系統中辨識任何變更。
4. 分割 LUN 並確認它們已正確對齊。請參閱「附錄 B：WAFL 校準驗證」、瞭解如何直接測試及確認校準。

### zPools

只有在中的步驟之後才應建立 zpool ["LUN 組態"](#) 執行。如果程序未正確執行、可能會因為 I/O 對齊而導致嚴重的效能降低。ONTAP 的最佳效能要求 I/O 必須與磁碟機上的 4K 邊界對齊。在 zpool 上建立的檔案系統使用有效

的區塊大小、並透過稱為的參數加以控制 `ashift`，您可以執行命令來檢視 `zdb -C`。

的價值 `ashift` 預設為 9、表示  $2^9$  或 512 位元組。為了獲得最佳效能 `ashift` 值必須為 12 ( $2^{12}=4K$ )。此值是在創建 `zpool` 時設置的，不能更改，這意味着 `zpool` 中的數據 `ashift` 除 12 個以外、應將資料複製到新建立的 `zPool`、以進行移轉。

建立 `zPool` 之後、請驗證的值 `ashift` 繼續之前。如果值不是 12、則表示未正確探索到 LUN。銷毀 `zpool`、確認相關主機公用程式文件中顯示的所有步驟均已正確執行、然後重新建立 `zPool`。

## zPools 和 Solaris LDoms

Solaris LDoms 還需要確保 I/O 對齊正確無誤。雖然 LUN 可能會被正確發現為 4K 裝置、但 LDOM 上的虛擬 `vdsk` 裝置不會繼承 I/O 網域的組態。以該 LUN 為基礎的 `vdsk` 預設為 512 位元組區塊。

需要額外的組態檔案。首先、必須針對 Oracle 錯誤 15824910 修補個別的 LDOM、才能啟用其他組態選項。此修補程式已移轉至所有目前使用的 Solaris 版本。一旦 LDOM 獲得修補、就可以依照下列方式設定新的正確對齊 LUN：

1. 識別要在新的 `zPool` 中使用的 LUN 或 LUN。在此範例中、它是 `c2d1` 裝置。

```
[root@LDM1 ~]# echo | format
Searching for disks...done
AVAILABLE DISK SELECTIONS:
  0. c2d0 <Unknown-Unknown-0001-100.00GB>
     /virtual-devices@100/channel-devices@200/disk@0
  1. c2d1 <SUN-ZFS Storage 7330-1.0 cyl 1623 alt 2 hd 254 sec 254>
     /virtual-devices@100/channel-devices@200/disk@1
```

2. 擷取要用於 ZFS Pool 的裝置之 VDC 執行個體：

```
[root@LDOM1 ~]# cat /etc/path_to_inst
#
# Caution! This file contains critical kernel state
#
"/fcoe" 0 "fcoe"
"/iscsi" 0 "iscsi"
"/pseudo" 0 "pseudo"
"/scsi_vhci" 0 "scsi_vhci"
"/options" 0 "options"
"/virtual-devices@100" 0 "vnex"
"/virtual-devices@100/channel-devices@200" 0 "cnex"
"/virtual-devices@100/channel-devices@200/disk@0" 0 "vdc"
"/virtual-devices@100/channel-devices@200/pciv-communication@0" 0 "vpci"
"/virtual-devices@100/channel-devices@200/network@0" 0 "vnet"
"/virtual-devices@100/channel-devices@200/network@1" 1 "vnet"
"/virtual-devices@100/channel-devices@200/network@2" 2 "vnet"
"/virtual-devices@100/channel-devices@200/network@3" 3 "vnet"
"/virtual-devices@100/channel-devices@200/disk@1" 1 "vdc" << We want
this one
```

### 3. 編輯 /platform/sun4v/kernel/drv/vdc.conf :

```
block-size-list="1:4096";
```

這表示裝置執行個體 1 的區塊大小為 4096 。

另一個範例是假設需要將 vdisk 執行個體 1 至 6 設定為 4K 區塊大小和 /etc/path\_to\_inst 內容如下：

```
"/virtual-devices@100/channel-devices@200/disk@1" 1 "vdc"
"/virtual-devices@100/channel-devices@200/disk@2" 2 "vdc"
"/virtual-devices@100/channel-devices@200/disk@3" 3 "vdc"
"/virtual-devices@100/channel-devices@200/disk@4" 4 "vdc"
"/virtual-devices@100/channel-devices@200/disk@5" 5 "vdc"
"/virtual-devices@100/channel-devices@200/disk@6" 6 "vdc"
```

### 4. 最終結果 vdc.conf 檔案應包含下列項目：

```
block-size-list="1:8192","2:8192","3:8192","4:8192","5:8192","6:8192";
```

## 注意

設定 VC.conf 並建立 vdsk 之後、必須重新啟動 LDOM。無法避免此步驟。區塊大小變更只會在重新開機後生效。繼續使用 zpool 組態、並確保如前所述、移位已正確設定為 12。

## ZFS Intent Log (ZIL)

一般而言、沒有理由在不同的裝置上找到 ZFS Intent Log (ZIL)。記錄檔可以與主集區共用空間。獨立 ZIL 的主要用途是使用缺乏現代儲存陣列寫入快取功能的實體磁碟機。

## logbias

設定 logbias 託管 Oracle 資料的 ZFS 檔案系統參數。

```
zfs set logbias=throughput <filesystem>
```

使用此參數可降低整體寫入層級。根據預設值、寫入的資料會先提交至 ZIL、然後再提交至主儲存池。此方法適用於使用純磁碟機組態的組態、包括 SSD 型 ZIL 裝置和主儲存池的旋轉媒體。這是因為它允許在可用的最低延遲媒體上、在單一 I/O 交易中進行認可。

使用包含其快取功能的現代化儲存陣列時、通常不需要使用此方法。在極少數情況下、可能需要在單一交易中寫入記錄檔、例如由高度集中、對延遲敏感的隨機寫入所組成的工作負載。寫入放大的形式會產生影響、因為記錄的資料最終會寫入主儲存池、導致寫入活動加倍。

## 直接 I/O

許多應用程式 (包括 Oracle 產品) 都可以啟用直接 I/O、藉此略過主機緩衝區快取此策略無法在 ZFS 檔案系統中正常運作。雖然會略過主機緩衝區快取、但 ZFS 本身仍會繼續快取資料。使用 Fio 或 Sio 等工具執行效能測試時、這項動作可能會產生誤導性的結果、因為很難預測 I/O 是否到達儲存系統、或是是否在作業系統中本機快取。此動作也會讓使用此類模擬測試來比較 ZFS 效能與其他檔案系統的情況變得非常困難。實際上、在真實使用者工作負載下、檔案系統效能幾乎沒有任何差異。

## 多個 zPools

必須在 zpool 層級執行快照型備份、還原、複製及歸檔 ZFS 型資料、而且通常需要多個 zPools。zpool 類似於 LVM 磁碟群組、應使用相同的規則進行設定。例如、資料庫的配置最好是存放在資料檔案上 zpool1 以及駐留在上的歸檔記錄、控制檔和重做記錄 zpool2。此方法允許標準熱備份、將資料庫置於熱備份模式、然後是的快照 zpool1。接著會從熱備份模式移除資料庫、強制進行記錄歸檔、並建立快照 zpool2 已建立。還原作業需要卸載 zfs 檔案系統、並在執行 SnapRestore 還原作業之後、將 zPool 完全離線。然後可以重新上線並恢復資料庫。

## filesystemio\_options

Oracle 參數 filesystemio\_options 使用 ZFS 的方式不同。如果 setall 或 directio 使用時、寫入作業會同步並略過 OS 緩衝區快取、但讀取會由 ZFS 進行緩衝。此動作會導致效能分析方面的困難、因為有時會被 ZFS 快取攔截和服務 I/O、使儲存延遲和總 I/O 比預期的要少。

## 網路組態

## Oracle 資料庫的邏輯介面設計

Oracle 資料庫需要存取儲存設備。邏輯介面（Lifs）是將儲存虛擬機器（SVM）連接到網路、然後再連接到資料庫的網路配送。需要適當的 LIF 設計、才能確保每個資料庫工作負載都有足夠的頻寬、而且容錯移轉不會導致儲存服務遺失。

本節概述 LIF 的主要設計原則。如需更完整的文件、請參閱 "[ONTAP 網路管理文件](#)"。與資料庫架構的其他層面一樣、儲存虛擬機器（SVM、在 CLI 稱為 vserver）和邏輯介面（LIF）設計的最佳選項、在很大程度上取決於擴充需求和業務需求。

建置 LIF 策略時、請考量下列主要主題：

- \* 效能。\* 網路頻寬是否足夠？
- \* 恢復能力。\* 設計中是否有任何單點故障？
- \* 管理能力。\* 網路是否能不中斷地擴充？

這些主題適用於端點對端點解決方案、從主機到交換器、再到儲存系統。

### LIF 類型

有多種 LIF 類型。"[LIF 類型的 ONTAP 文件](#)" 提供更完整的本主題資訊、但從功能觀點來看、生命可分為下列群組：

- \* 用於管理儲存叢集的叢集與節點管理生命期 \*。
- \* SVM 管理階層。\* 允許透過 REST API 或 ONTAPI（也稱為 ZAPI）存取 SVM 的介面、以執行快照建立或磁碟區調整大小等功能。SnapManager for Oracle（SMO）等產品必須能夠存取 SVM 管理 LIF。
- \* 資料生命。\* FC、iSCSI、NVMe/FC、NVMe/TCP、NFS、或 SMB/CIFS 資料。



用於 NFS 流量的資料 LIF 也可透過變更防火牆原則來進行管理 data 至 mgmt 或其他允許 HTTP、HTTPS 或 SSH 的原則。這項變更可避免每部主機的組態設定、以同時存取 NFS 資料 LIF 和個別的管理 LIF、進而簡化網路組態。儘管這兩者都使用 IP 傳輸協定、但無法同時為 iSCSI 和管理流量設定介面。iSCSI 環境需要個別的管理 LIF。

### SAN LIF 設計

SAN 環境中的 LIF 設計相對簡單、原因有一：多重路徑。所有現代化的 SAN 實作均可讓用戶端透過多個、不受限制的網路路徑存取資料、並選擇最佳的存取路徑或路徑。因此、LIF 設計的效能更容易因應、因為 SAN 用戶端會在最佳可用路徑之間自動平衡 I/O 負載。

如果路徑無法使用、用戶端會自動選取不同的路徑。因此設計簡易性讓 SAN 的工作更容易管理。這並不表示 SAN 環境總是更容易管理、因為 SAN 儲存設備還有許多其他層面比 NFS 複雜得多。這只是表示 SAN LIF 設計更簡單。

#### 效能

在 SAN 環境中、LIF 效能的最重要考量是頻寬。例如、雙節點 ONTAP AFF 叢集每個節點有兩個 16GB FC 連接埠、可在每個節點之間提供高達 32GB 的頻寬。

## 恢復能力

AFF 儲存系統上的 SAN Lifs 不會容錯移轉。如果 SAN LIF 因控制器容錯移轉而失敗、則用戶端的多重路徑軟體會偵測路徑遺失、並將 I/O 重新導向至不同的 LIF。使用 ASA 儲存系統時、在短暫延遲後將會容錯移轉、但這不會中斷 IO、因為其他控制器上已經有作用中的路徑。發生容錯移轉程序是為了在所有定義的連接埠上還原主機存取。

## 管理能力

LIF 移轉是 NFS 環境中較常見的工作、因為 LIF 移轉通常與在叢集周圍重新放置磁碟區有關。當磁碟區移轉至 HA 配對內時、無需在 SAN 環境中移轉 LIF。這是因為在磁碟區移動完成之後、ONTAP 會傳送路徑變更通知給 SAN、而 SAN 用戶端會自動重新最佳化。與 SAN 的 LIF 移轉主要與重大實體硬體變更有關。例如、如果需要不中斷營運的控制器升級、則 SAN LIF 會移轉至新硬體。如果發現 FC 連接埠故障、LIF 就可以移轉至未使用的連接埠。

## 設計建議

NetApp 提出下列建議：

- 請勿建立超過所需的路徑。過多的路徑會使整體管理更為複雜、並可能導致部分主機上路徑容錯移轉的問題。此外、有些主機對 SAN 開機等組態有非預期的路徑限制。
- 極少數組態需要四條以上的路徑才能連接到 LUN。如果擁有 LUN 的節點及其 HA 合作夥伴故障、則無法存取主控 LUN 的集合、因此限制將超過兩個節點的路徑通告至 LUN 的價值。在非主要 HA 配對的節點上建立路徑、在這種情況下並無幫助。
- 雖然可視 LUN 路徑的數量可以透過選擇 FC 區域中包含哪些連接埠來進行管理、但通常較容易在 FC 區域中包含所有潛在目標點、並控制 ONTAP 層級的 LUN 可見度。
- 在 ONTAP 8.3 及更新版本中、選擇性 LUN 對應 (SLM) 功能為預設功能。透過 SLM、任何新的 LUN 都會自動從擁有基礎 Aggregate 的節點和節點的 HA 合作夥伴通告。這種安排可避免建立連接埠集或設定分區、以限制連接埠存取。每個 LUN 都可在最佳效能和恢復能力所需的最小節點數上使用。
- 如果必須在兩個控制器之外移轉 LUN、則可以使用新增額外的節點 `lun mapping add-reporting-nodes` 命令、以便在新節點上通告 LUN。這樣做會建立通往 LUN 的額外 SAN 路徑、以進行 LUN 移轉。但是、主機必須執行探索作業、才能使用新路徑。
- 不要過度擔心間接流量。最好在 I/O 密集環境中避免間接流量、因為每微秒的延遲都是關鍵、但對於一般工作負載而言、可見的效能影響卻微不足道。

## NFS LIF 設計

與 SAN 通訊協定不同、NFS 定義多個資料路徑的能力有限。NFSv4 的平行 NFS (pNFS) 擴充解決了這項限制、但由於乙太網路速度已達 100GB、而且在新增其他路徑時、很少會有價值。

## 效能與恢復能力

雖然測量 SAN LIF 效能主要是從所有主要路徑計算總頻寬、但判斷 NFS LIF 效能需要仔細瞭解確切的網路組態。例如、兩個 10Gb 連接埠可設定為原始實體連接埠、或是可設定為連結集合體控制傳輸協定 (LACP) 介面群組。如果將它們設定為介面群組、則可根據流量是交換還是路由、使用不同的多個負載平衡原則。最後、Oracle Direct NFS (DNFS) 提供目前不存在於任何 OS NFS 用戶端的負載平衡組態。

與 SAN 通訊協定不同的是、NFS 檔案系統需要在通訊協定層恢復能力。例如、LUN 一律設定為啟用多重路徑、表示儲存系統可使用多個備援通道、每個通道都使用 FC 傳輸協定。另一方面、NFS 檔案系統則取決於單一 TCP/IP 通道的可用度、而該通道只能在實體層受到保護。這種配置是為何存在連接埠容錯移轉和 LACP 連接埠集合等選項。



在 NFS 環境中、網路傳輸協定層會同時提供效能和恢復能力。因此、這兩個主題彼此交織在一起、必須一起討論。

### 將生命體繫結至連接埠群組

若要將 LIF 繫結至連接埠群組、請將 LIF IP 位址與一組實體連接埠建立關聯。將實體連接埠集合在一起的主要方法是 LACP。LACP 的容錯功能相當簡單；LACP 群組中的每個連接埠都會受到監控、並在發生故障時從連接埠群組中移除。不過、對於 LACP 在效能方面的運作方式、有許多誤解：

- LACP 不需要在交換器上進行組態以符合端點。例如、ONTAP 可設定 IP 型負載平衡、而交換器則可使用 MAC 型負載平衡。
- 使用 LACP 連線的每個端點可以個別選擇封包傳輸連接埠、但無法選擇用於接收的連接埠。這表示從 ONTAP 到特定目的地的流量會連結到特定連接埠、而傳回流量可能會到達不同的介面。但這不會造成問題。
- LACP 不會一直平均分配流量。在擁有許多 NFS 用戶端的大型環境中、通常甚至會使用 LACP 集合中的所有連接埠。不過、環境中的任何一個 NFS 檔案系統都只能使用一個連接埠的頻寬、而非整個集合。
- 雖然 ONTAP 上有資源配置資源配置資源 LACP 原則、但這些原則並不會解決從交換器到主機的連線問題。例如、主機上有四埠 LACP 主幹的組態、ONTAP 上有四埠 LACP 主幹的組態、仍只能使用單一連接埠讀取檔案系統。雖然 ONTAP 可以透過所有四個連接埠傳輸資料、但目前沒有任何交換器技術可以透過所有四個連接埠從交換器傳送到主機。僅使用一個。

在包含許多資料庫主機的大型環境中、最常見的方法是使用 IP 負載平衡、建立一個包含適當數量 10Gb（或更快）介面的 LACP 集合體。只要有足夠的用戶端、這種方法就能讓 ONTAP 提供所有連接埠的均勻使用。當組態中的用戶端較少時、負載平衡會中斷、因為 LACP 主幹不會動態重新分配負載。

建立連線後、特定方向的流量只會放置在一個連接埠上。例如、對透過四埠 LACP 主幹連接的 NFS 檔案系統執行完整表格掃描的資料庫、只會透過一個網路介面卡（NIC）讀取資料。如果只有三個資料庫伺服器在這種環境中、則可能所有三個都從同一個連接埠讀取、而其他三個連接埠則處於閒置狀態。

### 將生命與實體連接埠繫結

將 LIF 繫結至實體連接埠、可更精細地控制網路組態、因為 ONTAP 系統上的指定 IP 位址一次只與一個網路連接埠相關聯。然後、可透過設定容錯移轉群組和容錯移轉原則來實現恢復能力。

### 容錯移轉原則和容錯移轉群組

網路中斷期間的生命行為是由容錯移轉原則和容錯移轉群組所控制。不同版本的 ONTAP 已變更組態選項。請參閱 ["適用於容錯移轉群組和原則的 ONTAP 網路管理文件"](#) 以取得所部署 ONTAP 版本的特定詳細資料。

ONTAP 8.3 及更高版本可根據廣播網域來管理 LIF 容錯移轉。因此、系統管理員可以定義所有可存取指定子網路的連接埠、並允許 ONTAP 選取適當的容錯移轉 LIF。這種方法可由部分客戶使用、但由於缺乏可預測性、因此在高速儲存網路環境中有限制。例如、環境可同時包含 1Gb 連接埠、以供例行檔案系統存取、而 10Gb 連接埠則可用於資料檔案 I/O。如果兩種連接埠都存在於同一個廣播網域中、LIF 容錯移轉可能會導致資料檔案 I/O 從 10Gb 連接埠移至 1Gb 連接埠。

總而言之、請考慮下列實務做法：

1. 將容錯移轉群組設定為使用者定義。
2. 將儲存容錯移轉（SFO）合作夥伴控制器上的連接埠填入容錯移轉群組、以便在儲存容錯移轉期間、生命體跟隨集合體。如此可避免產生間接流量。

3. 使用效能特性與原始 LIF 相符的容錯移轉連接埠。例如、單一實體 10Gb 連接埠上的 LIF 應包含單一 10Gb 連接埠的容錯移轉群組。四埠 LACP LIF 應容錯移轉至另一個四埠 LACP LIF。這些連接埠將是廣播網域中定義的連接埠子集。
4. 將容錯移轉原則設為僅限 SFO 合作夥伴。這樣做可確保 LIF 在容錯移轉期間跟隨集合體。

## 自動還原

設定 `auto-revert` 視需要設定參數。大多數客戶偏好將此參數設為 `true` 讓 LIF 還原至其主連接埠。不過、在某些情況下、客戶將此設定為「假」、表示在將 LIF 傳回其主連接埠之前、可以調查非預期的容錯移轉。

## LIF 與 Volume 比率

常見的誤解是、磁碟區和 NFS 生命體之間必須有一對一的關係。雖然在叢集中的任何位置移動磁碟區都需要此組態、但絕不會產生額外的互連流量、但絕對不需要此組態。必須考慮叢集間流量、但僅存在叢集間流量並不會造成問題。為 ONTAP 所發佈的許多基準測試主要包括間接 I/O

例如、資料庫專案中包含相對少數的效能關鍵資料庫、只需要總共 40 個磁碟區、可能需要將 1 : 1 磁碟區轉換為 LIF 策略、這種安排需要 40 個 IP 位址。然後、任何磁碟區都可以連同相關的 LIF 一起移至叢集中的任何位置、而且流量永遠是直接的、即使在微秒層級、也能將每個延遲來源減至最低。

舉例來說、大型託管環境的管理可能更容易、因為客戶與生命的關係是一對一。隨著時間的推移、可能需要將磁碟區移轉至不同的節點、這會造成一些間接流量。但是、除非互連交換器上的網路連接埠飽和、否則效能影響應該無法偵測。如果有疑慮、可以在其他節點上建立新的 LIF、並在下一個維護時段更新主機、以移除組態中的間接流量。

## 用於 Oracle 資料庫的 TCP/IP 和乙太網路組態

ONTAP 上的許多 Oracle 客戶都使用乙太網路、NFS、iSCSI、NVMe / TCP 的網路傳輸協定、尤其是雲端。

### 主機作業系統設定

大多數應用程式廠商文件都包含特定的 TCP 和乙太網路設定、以確保應用程式能以最佳方式運作。這些相同的設定通常足以提供最佳的 IP 型儲存效能。

### 乙太網路流量控制

這項技術可讓用戶端要求傳送者暫時停止資料傳輸。這通常是因為接收者無法快速處理傳入的資料。一次、要求傳送者停止傳輸的中斷程度比接收者丟棄封包的中斷程度低、因為緩衝區已滿。現今作業系統中使用的 TCP 堆疊已不再如此。事實上、流量控制所造成的問題比解決的問題還多。

近年來、乙太網路流量控制所造成的效能問題不斷增加。這是因為乙太網路流量控制是在實體層運作。如果網路組態允許任何主機作業系統將乙太網路流量控制要求傳送至儲存系統、則所有連線的用戶端都會暫停 I/O。由於單一儲存控制器服務的用戶端數量不斷增加、因此其中一或多個用戶端傳送流量控制要求的可能性會增加。在擁有廣泛作業系統虛擬化的客戶據點、經常會發現這個問題。

NetApp 系統上的 NIC 不應接收流量控制要求。實現此結果的方法因網路交換器製造商而異。在大多數情況下、可將乙太網路交換器上的流量控制設定為 `receive desired` 或 `receive on`，這意味着流控制請求不會轉發到儲存控制器。在其他情況下、儲存控制器上的網路連線可能不允許停用流程控制。在這些情況下、用戶端必須設定為永遠不要傳送流量控制要求、方法是變更至主機伺服器本身的 NIC 組態、或是變更主機伺服器所連接的交換器連接埠。



\* NetApp 建議 \* 確保 NetApp 儲存控制器不會接收以太網路流量控制封包。這通常可以透過設定控制器所連接的交換器連接埠來完成、但有些交換器硬體有限制、可能需要改用用戶端變更。

## MTU 大小

使用巨型框架的結果顯示、透過降低 CPU 和網路成本、可在速度較低的網路中提供一些效能改善、但效益通常並不顯著。



\* NetApp 建議 \* 盡可能實作巨型框架、以實現任何可能的效能效益、並確保解決方案符合未來需求。

在 10Gb 網路中使用巨型框架幾乎是強制性的。這是因為大多數的 10Gb 實作都達到每秒封包數的限制、而不需要巨型框架、就能達到 10Gb 標誌。使用巨型框架可改善 TCP/IP 處理效率、因為它可讓作業系統、伺服器、NIC 和儲存系統處理較少但較大的封包。效能的改善因 NIC 而異、但成效相當顯著。

對於巨型框架實作、通常但不正確的看法是、所有連線的裝置都必須支援巨型框架、而且 MTU 大小必須與端點對端點相符而是在建立連線時、兩個網路端點會協商最高的雙方可接受的框架大小。在一般環境中、網路交換器的 MTU 大小設為 9216、NetApp 控制器設為 9000、用戶端則設為 9000 和 1514 的混合。支援 9000 MTU 的用戶端可以使用巨型框架、而只支援 1514 的用戶端可以協商較低的值。

在完全交換的環境中、這種配置的問題很少發生。不過、在沒有中繼路由器被迫分割巨型框架的路由環境中、請務必小心。



- NetApp 建議 \* 設定下列項目：
- 使用 1 GB 以太網路 (GbE) 時、巨型框架是理想的選擇、但不是必要的。
- 使用 10GbE 及更快的速度、需要巨型框架才能達到最佳效能。

## TCP 參數

三項設定通常設定錯誤：TCP 時間戳記、選擇性認可 (SACK) 和 TCP 視窗縮放。網際網路上的許多過時文件建議停用一或多個這些參數、以改善效能。這項建議在多年前就有一些優點、因為 CPU 功能較低、因此有助於盡可能降低 TCP 處理的成本。

然而、在現代化的作業系統中、停用任何這些 TCP 功能通常會導致無法偵測的效益、同時也可能造成效能受損。在虛擬化網路環境中、效能受損的可能性特別大、因為這些功能是有有效處理封包遺失和網路品質變更所必需的。



\* NetApp 建議 \* 在主機上啟用 TCP 時間戳記、SACK 和 TCP 視窗縮放功能、而且在任何目前的作業系統中、這三個參數都應該預設為開啟。

## Oracle 資料庫的 FC 組態

為 Oracle 資料庫設定 FC SAN 主要是為了遵循日常的 SAN 最佳實務做法。

這包括典型的規劃措施、例如確保主機和儲存系統之間的 SAN 上有足夠的頻寬、使用 FC 交換器廠商所需的 FC 連接埠設定、檢查所有必要裝置之間是否存在所有 SAN 路徑、避免 ISL 爭用、並使用適當的 SAN 架構監控。

## 分區

FC 區域不得包含多個啟動器。這種安排一開始可能會運作、但啟動器之間的串擾最終會影響效能和穩定性。

雖然在極少數情況下、來自不同廠商的 FC 目標連接埠行為造成問題、但多目標區域通常被視為安全區域。例如、避免將 NetApp 和非 NetApp 儲存陣列的目標連接埠同時納入同一區域。此外、將 NetApp 儲存系統和磁帶裝置置於同一個區域、更有可能造成問題。

## Oracle 資料庫與直接連線 ONTAP 連線

儲存管理員有時偏好從組態中移除網路交換器、以簡化其基礎架構。在某些情況下可能會支援這項功能。

### iSCSI 和 NVMe / TCP

使用 iSCSI 或 NVMe / TCP 的主機可以直接連線至儲存系統、並正常運作。原因是路徑。直接連線至兩個不同的儲存控制器、會產生兩個不同的資料流路徑。遺失路徑、連接埠或控制器並不會妨礙其他路徑的使用。

### NFS

可以使用直接連線的 NFS 儲存設備、但有很大的限制：如果沒有大量的指令碼工作、容錯移轉將無法運作、這是客戶的責任。

直接連線的 NFS 儲存設備會造成不中斷的容錯移轉複雜化、這是因為本機作業系統上會發生路由。例如、假設主機的 IP 位址為 192.168.1.1/24、並直接連線至 IP 位址為 192.168.1.50/24 的 ONTAP 控制器。在容錯移轉期間、該位址 192.168.1.50 可以容錯移轉至其他控制器、而且主機可以使用該位址、但主機如何偵測其存在？原來的 192.168.1.1 位址仍然存在於不再連線至作業系統的主機 NIC 上。目的地為 192.168.1.5 的流量將繼續傳送至無法運作的網路連接埠。

第二個 OS NIC 可設定為 19 可以與故障的 over 192.168.1.50 位址進行通訊、但本機路由表預設會使用一個 \* 且只有一個 \* 位址來與 192.168.1.0/24 子網路通訊。系統管理員可以建立指令碼架構、以偵測失敗的網路連線、並變更本機路由表或使介面正常運作。具體程序取決於所使用的作業系統。

實際上、NetApp 客戶確實有直接連線的 NFS、但通常僅適用於容錯移轉期間 IO 暫停的工作負載。使用硬掛載時、在這類暫停期間不應有任何 IO 錯誤。IO 應該會暫停運作、直到服務還原為止、無論是透過容錯回復或手動介入、在主機上的 NIC 之間移動 IP 位址。

### FC 直接連線

無法使用 FC 傳輸協定將主機直接連接至 ONTAP 儲存系統。原因是使用 NPIV。用於識別 FC 網路的 ONTAP FC 連接埠的 WWN 使用稱為 NPIV 的虛擬化類型。任何連接至 ONTAP 系統的裝置都必須能夠辨識 NPIV WWN。目前沒有任何 HBA 廠商提供可安裝在能夠支援 NPIV 目標的主機上的 HBA。

## 儲存組態

### FC SAN

#### Oracle 資料庫 I/O 的 LUN 對齊

LUN 對齊是指針對基礎檔案系統配置最佳化 I/O。

在 ONTAP 系統上、儲存設備是以 4KB 為單位進行組織。資料庫或檔案系統 8KB 區塊應對應至兩個 4KB 區塊。如果 LUN 組態發生錯誤、在任一方向將對齊移至 1KB、則每個 8KB 區塊會存在於三個不同的 4KB 儲存區塊、而非兩個。這種安排會導致延遲增加、並導致在儲存系統中執行額外的 I/O。

對齊也會影響 LVM 架構。如果在整個磁碟機裝置上定義邏輯磁碟區群組內的實體磁碟區（不建立分割區）、LUN 上的前 4KB 區塊會與儲存系統上的前 4KB 區塊對齊。這是正確的對齊方式。磁碟分割發生問題、因為它們會移轉作業系統使用 LUN 的起始位置。只要偏移量以 4KB 的整體單位移動、LUN 就會對齊。

在 Linux 環境中、在整個磁碟機裝置上建立邏輯磁碟區群組。當需要磁碟分割時、請執行檢查對齊 `fdisk -u` 並驗證每個分割區的開始時間為八個之倍數。這表示分割區從八個 512 位元組磁區的倍數開始、即 4KB。

另請參閱一節中有關壓縮區塊對齊的討論 "效率"。任何與 8KB 壓縮區塊邊界對齊的配置、也會與 4KB 邊界對齊。

#### 錯誤對齊警告

資料庫重做 / 交易記錄通常會產生未對齊的 I/O、導致 ONTAP 上未對齊 LUN 的錯誤警告。

記錄會以不同大小的寫入方式、連續寫入記錄檔。不符合 4KB 界限的記錄寫入作業通常不會造成效能問題、因為下一個記錄寫入作業會完成區塊。結果是 ONTAP 幾乎能將所有寫入作業視為完整的 4KB 區塊來處理、即使某些 4KB 區塊中的資料是以兩個不同的作業來寫入。

使用公用程式（例如）來驗證對齊 `sio` 或 `dd` 可在定義的區塊大小下產生 I/O。您可以使用檢視儲存系統上的 I/O 對齊統計資料 `stats` 命令。請參閱 "WAFS 對齊驗證" 以取得更多資訊。

在 Solaris 環境中進行對齊更為複雜。請參閱 "SAN 主機組態 ONTAP" 以取得更多資訊。

#### 注意

在 Solaris x86 環境中、由於大多數組態都有多層分割區、因此請格外注意正確的對齊方式。Solaris x86 分割區磁碟片通常位於標準主開機記錄分割區表格的上方。

## Oracle 資料庫 LUN 規模調整和 LUN 數量

選擇最佳 LUN 大小和要使用的 LUN 數量、對於 Oracle 資料庫的最佳效能和管理性至關重要。

LUN 是 ONTAP 上的虛擬化物件、存在於託管集合體中的所有磁碟機中。因此、LUN 的效能不受其大小影響、因為無論選擇何種大小、LUN 都會充分發揮彙總的效能潛力。

為了方便起見、客戶可能想要使用特定大小的 LUN。例如、如果資料庫建置在由兩個 LUN 組成的 LVM 或 Oracle ASM 磁碟群組上、每個 LUN 均為 1TB、則該磁碟群組必須以 1TB 為增量來擴充。最好是從八個 LUN（每個 LUN 為 500GB）構建磁盤組、以便可以以更小的增量來增加磁盤組。

我們不鼓勵建立通用標準 LUN 大小的做法、因為這樣做可能會使管理變得複雜。例如、當資料庫或資料存放區的範圍介於 1TB 到 2TB 時、100GB 的標準 LUN 大小可能運作良好、但大小為 20TB 的資料庫或資料存放區需要 200 個 LUN。這表示伺服器重新開機時間較長、不同 UI 中需要管理的物件較多、而 SnapCenter 等產品必須在許多物件上執行探索。使用較少、較大的 LUN 可避免此類問題。

- LUN 數量比 LUN 大小更重要。
- LUN 大小大多由 LUN 數需求控制。
- 避免建立超過所需數量的 LUN。

## LUN 計數

與 LUN 大小不同、LUN 數量確實會影響效能。應用程式效能通常取決於透過 SCSI 層執行平行 I/O 的能力。因此、兩個 LUN 的效能優於單一 LUN。使用 LVM（例如 Veritas VxVM、Linux LVM2 或 Oracle ASM）是提高平行度的最簡單方法。

NetApp 客戶通常從 LUN 數量增加到 16 個以上獲得最小的效益、不過測試 100% SSD 環境時、隨機 I/O 非常繁重、這已證實可進一步改善至 64 個 LUN。

- NetApp 建議 \* 下列事項：



一般而言、四到十六個 LUN 足以支援任何特定資料庫工作負載的 I/O 需求。由於主機 SCSI 實作的限制、少於四個 LUN 可能會造成效能限制。

## Oracle 資料庫 LUN 放置

資料庫 LUN 在 ONTAP 磁碟區內的最佳放置方式、主要取決於如何使用各種 ONTAP 功能。

### 磁碟區

與剛接觸 ONTAP 的客戶混淆的一個常見點是使用 FlexVols、通常稱為「Volume」。

磁碟區不是 LUN。這些詞彙與許多其他廠商產品（包括雲端供應商）同義。ONTAP Volume 是簡單的管理容器。它們本身不會提供資料、也不會佔用空間。它們是檔案或 LUN 的容器、可改善及簡化管理、尤其是大規模管理。

### 磁碟區和 LUN

相關 LUN 通常位於單一磁碟區中。例如、需要 10 個 LUN 的資料庫通常會將所有 10 個 LUN 放在同一個磁碟區上。



- 使用 LUN 對磁碟區的比例 1 : 1 表示每個磁碟區有一個 LUN、這是 \* 非 \* 正式最佳實務做法。
- 而是應將磁碟區視為工作負載或資料集的容器。每個磁碟區可能只有一個 LUN、或者可能有許多 LUN。正確的答案取決於管理需求。
- 在不必要數量的磁碟區之間分散 LUN、可能會導致額外的額外負荷和排程問題、例如快照作業、UI 中顯示的物件過多、並導致在達到 LUN 限制之前達到平台磁碟區限制。

### 磁碟區、LUN 和快照

Snapshot 原則和排程會放置在磁碟區上、而非 LUN 上。如果資料集由 10 個 LUN 組成、則當這些 LUN 位於同一個磁碟區中時、只需要單一快照原則。

此外、在單一磁碟區中共同定位給定資料集的所有相關 LUN、可提供原子快照作業。例如、如果基礎 LUN 全部放在單一磁碟區上、則位於 10 個 LUN 上的資料庫、或是由 10 個不同作業系統組成的 VMware 應用程式環境、都可以作為單一且一致的物件加以保護。如果將快照放在不同的磁碟區上、則即使同時排程、快照仍可能保持 100% 同步。

在某些情況下、由於恢復需求、相關的 LUN 集可能需要分割成兩個不同的磁碟區。例如、資料庫可能有四個 LUN 用於資料檔案、兩個 LUN 用於記錄。在這種情況下、具有 4 個 LUN 的資料檔案磁碟區和具有 2 個 LUN

的記錄磁碟區可能是最佳選擇。原因在於可進行的可恢復性是不相關的。例如、資料檔案磁碟區可以選擇性地還原為較早的狀態、這表示所有四個 LUN 都會還原為快照狀態、而記錄磁碟區與其重要資料則不會受到影響。

### Volume、LUN 和 SnapMirror

SnapMirror 原則和作業就像快照作業一樣、是在磁碟區上執行、而不是在 LUN 上執行。

在單一磁碟區中共同定位相關 LUN、可讓您建立單一 SnapMirror 關係、並透過單一更新來更新所有包含的資料。與快照一樣、更新也將是一項原子作業。SnapMirror 目的地將保證擁有來源 LUN 的單一時間點複本。如果 LUN 分散在多個磁碟區、則複本可能彼此一致、也可能不一致。

### 磁碟區、LUN 和 QoS

雖然 QoS 可以選擇性地套用至個別 LUN、但通常在磁碟區層級設定 QoS 會比較容易。例如、指定 ESX 伺服器中的來賓所使用的所有 LUN 都可以放置在單一磁碟區上、然後就可以套用 ONTAP 調適性 QoS 原則。結果是將每 TB IOPS 的自我擴充限制套用至所有 LUN。

同樣地、如果資料庫需要 10 萬次 IOPS、而且佔用 10 個 LUN、則在單一磁碟區上設定單一的 10 萬次 IOPS 限制、比在每個 LUN 上設定 10 個個別的 10K IOPS 限制更容易。

### 多重 Volume 配置

在某些情況下、跨多個磁碟區散佈 LUN 可能會有幫助。主要原因是控制器分段。例如、HA 儲存系統可能會裝載單一資料庫、其中需要每個控制器的完整處理與快取潛力。在這種情況下、典型的設計是將一半的 LUN 放在控制器 1 的單一磁碟區、而另一半的 LUN 則放在控制器 2 的單一磁碟區中。

同樣地、控制器分段也可用於負載平衡。HA 系統託管 100 個資料庫、每個資料庫各有 10 個 LUN、每個資料庫可在兩個控制器上接收 5 個 LUN 磁碟區。如此一來、每個控制器就能以對稱的方式進行對稱載入、同時還能配置額外的資料庫。

不過、這些範例都不涉及 1 : 1 的磁碟區對 LUN 比率。目標仍然是透過在磁碟區中共同定位相關 LUN 來最佳化管理性。

其中一個例子是、1 : 1 LUN 對磁碟區比率非常合理、其中每個 LUN 可能真正代表單一工作負載、而且每個工作負載都需要個別管理。在這種情況下、1 : 1 的比率可能是最佳的。

### Oracle 資料庫 LUN 調整大小和以 LVM 為基礎的調整大小

當 SAN 型檔案系統達到容量上限時、有兩個選項可以增加可用空間：

- 增加 LUN 的大小
- 將 LUN 新增至現有的磁碟區群組、並擴充內含的邏輯磁碟區

雖然 LUN 調整大小是增加容量的選項、但通常最好使用 LVM、包括 Oracle ASM。LVM 存在的主要原因之一、是為了避免需要調整 LUN 大小。使用 LVM 時、多個 LUN 會結合在一個虛擬儲存池中。從該池中切出的邏輯卷由 LVM 管理、可以輕鬆調整大小。另一項優點是在所有可用 LUN 之間分配給定的邏輯磁碟區、以避免在特定磁碟機上出現熱點。通常可以使用 Volume Manager 將邏輯磁碟區的基礎範圍重新放置到新的 LUN、以執行透明移轉。

### 使用 Oracle 資料庫的 LVM 分拆

LVM 分拆是指在多個 LUN 之間分配資料。如此一來、許多資料庫的效能大幅提升。

在快閃磁碟機時代之前、使用區塊延展來協助克服旋轉磁碟機的效能限制。例如、如果作業系統需要執行 1MB 讀取作業、則從單一磁碟機讀取 1MB 的資料時、需要大量的磁碟機磁頭搜尋和讀取、因為 1MB 會緩慢傳輸。如果將 1MB 的資料分散在 8 個 LUN 上、則作業系統可能會同時執行 8 個 128K 讀取作業、並縮短完成 1MB 傳輸所需的時間。

由於必須事先知道 I/O 模式、因此使用旋轉磁碟機進行分拆會更困難。如果串列區塊延展未針對真正的 I/O 模式正確調整、則等量區塊配置可能會損害效能。使用 Oracle 資料庫、特別是搭配 All Flash 組態、分拆作業更容易設定、並經證實可大幅提升效能。

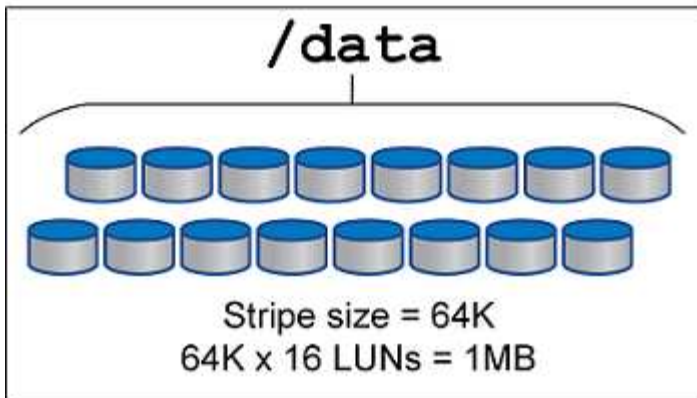
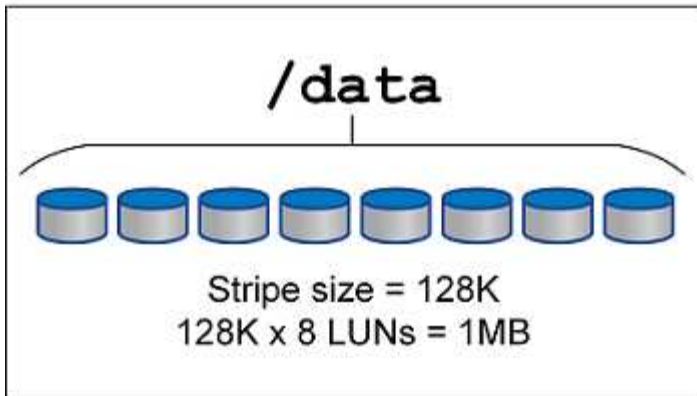
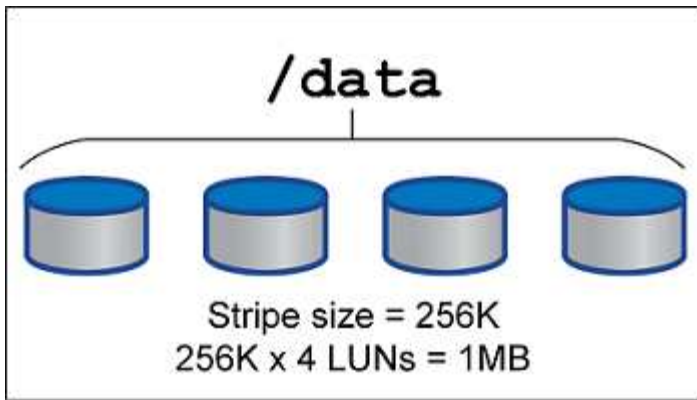
依預設、邏輯磁碟區管理程式（例如 Oracle ASM 等量磁碟區）不屬於原生 OS LVM。其中有些 LUN 會將多個 LUN 連結在一起、成為串連的裝置、導致資料檔案存在於一台 LUN 裝置上、而只存在於一台 LUN 裝置上。這會造成熱點。其他 LVM 實作預設為分散式擴充。這與分拆類似、但卻是比較粗糙的。磁碟區群組中的 LUN 會切成大型片段、稱為區段、通常以百萬位元組為單位測量、然後邏輯磁碟區會分佈在這些區段中。結果是對檔案進行隨機 I/O、應該能在 LUN 之間妥善分配、但連續 I/O 作業的效率卻不如以前那麼高。

效能密集的應用程式 I/O 幾乎總是（a）以基本區塊大小為單位、或（b）1 MB。

等量分配組態的主要目標是確保單一檔案 I/O 可作為單一單元執行、而多區塊 I/O 的大小應為 1MB、可在等量磁碟區中的所有 LUN 之間平均平行處理。這表示等量磁碟區大小不得小於資料庫區塊大小、且等量磁碟區大小乘以 LUN 數量應為 1MB。

下圖顯示等量磁碟區大小和寬度調校的三個可能選項。選擇 LUN 數量以滿足上述效能需求、但在所有情況下、單一等量磁碟區內的總資料為 1MB。





## NFS

### Oracle 資料庫的 NFS 組態

NetApp 已提供企業級 NFS 儲存設備超過 30 年、由於其簡易性、隨著雲端型基礎架構的推向、其使用量也不斷增加。

NFS 傳輸協定包含多個不同需求的版本。如需 ONTAP 的 NFS 組態完整說明、請參閱 ["TR-4067 ONTAP 最佳實務做法上的 NFS"](#)。下列各節涵蓋一些較重要的需求和一般使用者錯誤。

### NFS 版本

NetApp 必須支援作業系統 NFS 用戶端。

- NFSv3 支援符合 NFSv3 標準的作業系統。
- Oracle DNFS 用戶端支援 NFSv3。

- 所有遵循 NFSv4 標準的作業系統都支援 NFSv4 。
- NFSv4.1 和 NFSv4.2 需要特定的作業系統支援。請參閱 "NetApp IMT" 適用於支援的作業系統。
- Oracle DNFS 支援 NFSv4.1 需要 Oracle 12.2.0.2 或更高版本。



◦ "NetApp 支援對照表" 對於 NFSv3 和 NFSv4 、不包含特定的作業系統。一般支援所有遵守 RFC 的作業系統。搜尋線上 IMT 以取得 NFSv3 或 NFSv4 支援時、請勿選取特定的作業系統、因為不會顯示任何相符項目。一般原則隱含支援所有作業系統。

### Linux NFSv3 TCP 插槽表

TCP 插槽表是與主機匯流排介面卡（HBA）佇列深度相當的 NFSv3 。這些表格可控制任何時間都可以處理的 NFS 作業數量。預設值通常為 16、這對於最佳效能而言太低。相反的問題發生在較新的 Linux 核心上、這會自動將 TCP 插槽表格限制增加到要求使 NFS 伺服器飽和的層級。

為了達到最佳效能並避免效能問題、請調整控制 TCP 插槽表的核心參數。

執行 `sysctl -a | grep tcp.*.slot_table` 並觀察下列參數：

```
# sysctl -a | grep tcp.*.slot_table
sunrpc.tcp_max_slot_table_entries = 128
sunrpc.tcp_slot_table_entries = 128
```

所有 Linux 系統都應該包括在內 `sunrpc.tcp_slot_table_entries`、但只有部分包含在內 `sunrpc.tcp_max_slot_table_entries`。兩者都應設為 128。

#### 注意

若未設定這些參數、可能會對效能造成重大影響。在某些情況下、效能會受到限制、因為 Linux 作業系統沒有發出足夠的 I/O 在其他情況下、隨著 Linux 作業系統嘗試發出的 I/O 數量超過可服務的數量、I/O 延遲也會增加。

### ADR 和 NFS

部分客戶回報的效能問題是由於中的資料 I/O 過多所造成 ADR 位置。問題通常不會在累積許多效能資料之前發生。I/O 過多的原因不明、但此問題似乎是 Oracle 處理程序重複掃描目標目錄以進行變更所致。

移除 `noac` 和/或 `actimeo=0` 掛載選項可執行主機作業系統快取、並降低儲存 I/O 層級。



\* NetApp 建議 \* 不要放置 ADR 檔案系統上的資料 `noac` 或 `actimeo=0` 因為效能問題很可能會發生。獨立 ADR 如有必要、請將資料移至不同的掛載點。

### NFS-rootonly 和 mount-rootonly

ONTAP 包含一個稱為的 NFS 選項 `nfs-rootonly` 控制伺服器是否接受來自高連接埠的 NFS 流量連線。為了安全起見、只有 root 使用者可以使用低於 1024 的來源連接埠來開啟 TCP/IP 連線、因為這類連接埠通常是保留給作業系統使用、而非使用者處理程序。此限制有助於確保 NFS 流量來自實際的作業系統 NFS 用戶端、而非模擬 NFS 用戶端的惡意程序。Oracle DNFS 用戶端是 `userspace` 驅動程式、但程序是以 root 執行、因此通常不需要變更的值 `nfs-rootonly`。連線是從低連接埠建立。

◦ `mount-rootonly` 選項僅適用於 NFSv3。它控制是否從大於 1024 的連接埠接受 RPC 掛載呼叫。使用 DNFS 時、用戶端會再次以 root 執行、因此能夠開啟低於 1024 的連接埠。此參數無效。

透過 NFS 4.0 及更高版本開啟與 DNFS 連線的程序不會以 root 執行、因此需要 1024 以上的連接埠。◦ `nfs-rootonly` 參數必須設為停用、DNFS 才能完成連線。

如果 `nfs-rootonly` 啟用時、結果會在掛載階段開啟 DNFS 連線時暫停。sqlplus 輸出類似於以下內容：

```
SQL>startup
ORACLE instance started.
Total System Global Area 4294963272 bytes
Fixed Size                  8904776 bytes
Variable Size               822083584 bytes
Database Buffers           3456106496 bytes
Redo Buffers                 7868416 bytes
```

參數可變更如下：

```
Cluster01::> nfs server modify -nfs-rootonly disabled
```



在極少數情況下、您可能需要將 `NFS-rootonly` 和 `mount-rootonly` 都變更為停用。如果伺服器正在管理大量的 TCP 連線、則可能沒有低於 1024 的連接埠可用、而且作業系統必須使用較高的連接埠。需要變更這兩個 ONTAP 參數、才能完成連線。

#### NFS 匯出原則：超級使用者和 `setuid`

如果 Oracle 二進位檔位於 NFS 共用區、則匯出原則必須包含超級使用者和 `setuid` 權限。

用於一般檔案服務（例如使用者主目錄）的共享 NFS 匯出通常會佔用 root 使用者。這表示已掛載檔案系統的主機上 root 使用者的要求、會重新對應為具有較低權限的不同使用者。這有助於防止特定伺服器上的根使用者存取共用伺服器上的資料、進而保護資料安全。在共享環境中、`setuid` 位元也可能是安全性風險。`setuid` 位元可讓處理程序以不同於使用者的身分執行、而非以使用者的身分執行指令。例如、root 擁有的 Shell 指令碼搭配 `setuid` 位元、會以 root 執行。如果其他使用者可以變更該 Shell 指令碼、任何非 root 使用者都可以透過更新指令碼、以 root 身分發出命令。

Oracle 二進位檔包含 root 擁有的檔案、並使用 `setuid` 位元。如果在 NFS 共用上安裝 Oracle 二進位檔、匯出原則必須包含適當的超級使用者和 `setuid` 權限。在以下範例中、這兩項規則都包含在內 `allow-suid` 及許可 `superuser`（root）使用系統驗證來存取 NFS 用戶端。

```
Cluster01::> export-policy rule show -vserver vserver1 -policyname orabin
-fields allow-suid,superuser
vserver  policyname ruleindex superuser allow-suid
-----
vserver1 orabin      1          sys      true
```

## NFSv4/4.1 組態

對於大多數應用程式、NFSv3 和 NFSv4 之間的差異很小。應用程式 I/O 通常是非常簡單的 I/O、而且不會從 NFSv4 中提供的某些進階功能中獲得顯著效益。較高版本的 NFS 不應從資料庫儲存的角度視為「升級」、而應視為包含其他功能的 NFS 版本。例如、如果需要 Kerberos 隱私模式 (krb5p) 的端點對端安全性、則需要 NFSv4。



\* 如果需要 NFSv4 功能、NetApp 建議 \* 使用 NFSv4.1。NFSv4.1 中的 NFSv4 傳輸協定有一些功能性增強功能、可改善某些邊緣情況的恢復能力。

切換至 NFSv4 比單純將掛載選項從 `ves=3` 變更為 `ves=4.1` 更複雜。如需更完整的 NFSv4 組態與 ONTAP 說明、包括作業系統設定指南、請參閱 ["TR-4067 ONTAP 最佳實務做法上的 NFS"](#)。本 TR 的下列各節說明使用 NFSv4 的一些基本要求。

## NFSv4 網域

NFSv4/4.1 組態的完整說明已超出本文件的範圍、但常見的問題之一是網域對應不相符。從系統管理員的角度來看、NFS 檔案系統的行為似乎正常、但應用程式會報告某些檔案的權限和 / 或 `setuid` 錯誤。在某些情況下、系統管理員不正確地判斷應用程式二進位檔的權限已受損、並在實際問題是網域名稱時執行 `chown` 或 `chmod` 命令。

NFSv4 網域名稱是在 ONTAP SVM 上設定：

```
Cluster01::> nfs server show -fields v4-id-domain
vserver    v4-id-domain
-----
vserver1   my.lab
```

主機上的 NFSv4 網域名稱是在中設定 `/etc/idmap.cfg`

```
[root@host1 etc]# head /etc/idmapd.conf
[General]
#Verbosity = 0
# The following should be set to the local NFSv4 domain name
# The default is the host's DNS domain name.
Domain = my.lab
```

網域名稱必須相符。如果沒有、則會在中顯示類似下列的對應錯誤 `/var/log/messages`：

```
Apr 12 11:43:08 host1 nfsidmap[16298]: nss_getpwnam: name 'root@my.lab'
does not map into domain 'default.com'
```

應用程式二進位檔（例如 Oracle 資料庫二進位檔）包含 `root` 擁有的具有 `setuid` 位元的檔案、這表示 NFSv4 網域名稱不相符會導致 Oracle 啟動失敗、並會發出呼叫檔案擁有權或權限的警告 `oradism`、位於 `$ORACLE_HOME/bin` 目錄。其內容應如下所示：

```
[root@host1 etc]# ls -l /orabin/product/19.3.0.0/dbhome_1/bin/oradism
-rwsr-x--- 1 root oinstall 147848 Apr 17 2019
/orabin/product/19.3.0.0/dbhome_1/bin/oradism
```

如果此檔案的擁有權為 nobody、則可能是 NFSv4 網域對應問題。

```
[root@host1 bin]# ls -l oradism
-rwsr-x--- 1 nobody oinstall 147848 Apr 17 2019 oradism
```

若要修正此問題、請參閱 `/etc/idmap.cfg` 根據 ONTAP 上的 vv4 識別碼網域設定來建立檔案、並確保檔案一致。如果沒有、請進行必要的變更、然後執行 `nfsidmap -c`、然後等待一段時間讓變更傳播。接著、檔案擁有權應正確辨識為 root。如果使用者嘗試執行 `chown root` 在 NFS 網域設定修正之前、可能需要在這個檔案上執行 `chown root` 再一次。

## Oracle directNFS

Oracle 資料庫可使用 NFS 的方式有兩種。

首先、它可以使用以作業系統一部分的原生 NFS 用戶端所掛載的檔案系統。這有時稱為核心 NFS 或 kNFS。NFS 檔案系統是由 Oracle 資料庫安裝及使用、與任何其他應用程式使用 NFS 檔案系統的方式完全相同。

第二種方法是 Oracle Direct NFS (DNFS)。這是在 Oracle 資料庫軟體中實作 NFS 標準。它不會改變 DBA 設定或管理 Oracle 資料庫的方式。只要儲存系統本身有正確的設定、就應該對 DBA 團隊和終端使用者透明使用 DNFS。

啟用 DNFS 功能的資料庫仍會掛載一般的 NFS 檔案系統。資料庫開啟後、Oracle 資料庫會開啟一組 TCP/IP 工作階段、並直接執行 NFS 作業。

### Direct NFS

Oracle Direct NFS 的主要值是略過主機 NFS 用戶端、並直接在 NFS 伺服器上執行 NFS 檔案作業。啟用它只需要變更 Oracle 磁碟管理程式 (ODM) 程式庫。Oracle 說明文件中提供此程序的說明。

使用 DNFS 可大幅提升 I/O 效能、並降低主機和儲存系統的負載、因為 I/O 是以最有效率的方式執行。

此外、Oracle DNFS 還包含 \* 選項 \*、可用於網路介面多重路徑和容錯功能。例如、兩個 10Gb 介面可以結合在一起、以提供 20Gb 的頻寬。如果某個介面發生故障、則會在另一個介面上重試 I/O。整體作業與 FC 多重路徑非常類似。多重路徑在數年前是最常見的標準、那就是 1 個乙太網路。10Gb NIC 足以應付大多數 Oracle 工作負載、但如果需要更多 10Gb NIC、則可加以連結。

使用 DNFS 時、必須安裝 Oracle Doc 1495104.1 中所述的所有修補程式。如果無法安裝修補程式、則必須評估環境、確保該文件中所述的錯誤不會造成問題。在某些情況下、無法安裝所需的修補程式會導致無法使用 DNFS。

請勿將 DNFS 用於任何類型的循環名稱解析、包括 DNS、DDNS、NIS 或任何其他方法。這包括 ONTAP 中可用的 DNS 負載平衡功能。當使用 DNFS 的 Oracle 資料庫將主機名稱解析為 IP 位址時、後續查詢時不得變更。這可能會導致 Oracle 資料庫當機、並可能導致資料毀損。

## 直接 NFS 和主機檔案系統存取

使用 DNFS 有時會導致依賴掛載在主機上的可見檔案系統的應用程式或使用者活動發生問題、因為 DNFS 用戶端會從主機作業系統不定期存取檔案系統。DNFS 用戶端可以在不瞭解作業系統的情況下建立、刪除及修改檔案。

使用單一執行個體資料庫的掛載選項時、會啟用檔案和目錄屬性的快取、這也表示目錄內容會快取。因此、DNFS 可以建立檔案、而且在作業系統重新讀取目錄內容和讓使用者看到檔案之前、會有短暫的延遲。這通常不是問題、但在極少數情況下、SAP BR\*Tools 等公用程式可能會發生問題。如果發生這種情況、請變更掛載選項、以使用 Oracle RAC 的建議來解決此問題。這項變更會導致停用所有主機快取。

只有在使用 (a) DNFS 時才變更掛載選項、且 (b) 檔案可見度延遲所造成的問題。如果未使用 DNFS、在單一執行個體資料庫上使用 Oracle RAC 掛載選項會導致效能降低。



請參閱附註 nosharecache 在中 "[Linux NFS 裝載選項](#)" 針對可能產生異常結果的 Linux 特定 DNFS 問題。

## Oracle 資料庫和 NFS 會租用和鎖定

NFSv3 無狀態。這實際上意味著 NFS 伺服器 (ONTAP) 無法追蹤哪些檔案系統是掛載的、由誰掛載、或哪些鎖定是真的就位。

ONTAP 確實有一些功能會記錄掛載嘗試、因此您可以知道哪些用戶端可能正在存取資料、而且可能會出現諮詢鎖定、但這項資訊並不保證 100% 完整。無法完成、因為追蹤 NFS 用戶端狀態並非 NFSv3 標準的一部分。

### NFSv4 狀態

相反地、NFSv4 是有狀態的。NFSv4 伺服器會追蹤哪些用戶端正在使用哪些檔案系統、哪些檔案存在、哪些檔案和 / 或檔案區域被鎖定等 這表示 NFSv4 伺服器之間需要定期通訊、才能保持狀態資料最新。

NFS 伺服器所管理的最重要狀態是 NFSv4 鎖定和 NFSv4 租賃、它們彼此之間有很大的關聯。您必須瞭解每個項目本身的運作方式、以及它們彼此之間的關係。

### NFSv4 鎖定

有了 NFSv3、鎖定是建議事項。NFS 用戶端仍可修改或刪除「鎖定」檔案。NFSv3 鎖本身不會過期、必須將其移除。這會造成問題。例如、如果叢集式應用程式會建立 NFSv3 鎖定、而其中一個節點發生故障、您該怎麼做？您可以在仍在運作的節點上對應用程式進行編碼、以移除鎖定、但您如何知道這是安全的？可能是「故障」節點可以運作、但無法與叢集的其他部分通訊？

有了 NFSv4、鎖定的持續時間有限。只要持有鎖定的用戶端繼續與 NFSv4 伺服器簽入、就不允許其他用戶端取得這些鎖定。如果用戶端無法使用 NFSv4 進行存回、伺服器最終會撤銷鎖定、而其他用戶端則能要求並取得鎖定。

### NFSv4 租賃

NFSv4 鎖定與 NFSv4 租用相關聯。當 NFSv4 用戶端與 NFSv4 伺服器建立連線時、它會取得租用。如果用戶端取得鎖定 (鎖定類型眾多)、則鎖定會與租用相關聯。

此租用具有定義的逾時時間。根據預設、ONTAP 會將逾時值設為 30 秒：

```
Cluster01::*> nfs server show -vserver vserver1 -fields v4-lease-seconds

vserver    v4-lease-seconds
-----
vserver1   30
```

這表示 NFSv4 用戶端需要每 30 秒與 NFSv4 伺服器簽入一次、才能續約。

租賃會自動由任何活動續約、因此如果用戶端正在工作、就不需要執行額外作業。如果某個應用程式變得很安靜、而且沒有真正的工作、則需要改為執行某種保持活動狀態的作業（稱為順序）。基本上只是說：「我還在這裏、請重新整理我的租約。」

```
*Question:* What happens if you lose network connectivity for 31 seconds?
NFSv3 無狀態。這並不需要用戶端的通訊。NFSv4
可設定狀態、一旦租用期間結束、租用即會過期、鎖定會被撤銷、而鎖定的檔案會提供給其他用戶端
使用。
```

有了 NFSv3、您可以四處移動網路纜線、重新啟動網路交換器、進行組態變更、並確保不會發生任何問題。應用程式通常只會耐心等待網路連線再次運作。

有了 NFSv4、您有 30 秒的時間（除非您已在 ONTAP 中增加該參數的值）來完成工作。如果您超過此上限、您的租約將會逾時。這通常會導致應用程式當機。

舉例來說、如果您有 Oracle 資料庫、而且網路連線中斷（有時稱為「網路分割區」）超過租用逾時、您就會使資料庫當機。

以下是 Oracle 警示記錄中發生這種情況的範例：

```
2022-10-11T15:52:55.206231-04:00
Errors in file /orabin/diag/rdbms/ntap/NTAP/trace/NTAP_ckpt_25444.trc:
ORA-00202: control file: '/redo0/NTAP/ctrl/control01.ctl'
ORA-27072: File I/O error
Linux-x86_64 Error: 5: Input/output error
Additional information: 4
Additional information: 1
Additional information: 4294967295
2022-10-11T15:52:59.842508-04:00
Errors in file /orabin/diag/rdbms/ntap/NTAP/trace/NTAP_ckpt_25444.trc:
ORA-00206: error in writing (block 3, # blocks 1) of control file
ORA-00202: control file: '/redo1/NTAP/ctrl/control02.ctl'
ORA-27061: waiting for async I/Os failed
```

如果您查看系統記錄檔、您應該會看到以下幾個錯誤：

```
Oct 11 15:52:55 host1 kernel: NFS: nfs4_reclaim_open_state: Lock reclaim failed!
Oct 11 15:52:55 host1 kernel: NFS: nfs4_reclaim_open_state: Lock reclaim failed!
Oct 11 15:52:55 host1 kernel: NFS: nfs4_reclaim_open_state: Lock reclaim failed!
```

記錄訊息通常是問題的第一個徵象、而非應用程式凍結。通常、在網路中斷期間、您完全看不到任何內容、因為程序和作業系統本身都遭到封鎖、無法嘗試存取 NFS 檔案系統。

網路重新運作後、就會出現錯誤。在上述範例中、一旦重新建立連線、作業系統就會嘗試重新取得鎖定、但時間太晚了。租約已到期、鎖定已移除。這會導致一個錯誤、該錯誤會傳播到 Oracle 層、並導致警示記錄中出現訊息。根據資料庫的版本和組態、您可能會看到這些模式的變化。

總之、NFSv3 可容忍網路中斷、但 NFSv4 更敏感、並規定了一段定義的租用期。

如果無法接受 30 秒的逾時、該怎麼辦？如果您管理一個動態變化的網路、在其中重新啟動交換器或重新放置纜線、導致網路偶爾中斷、該怎麼辦？您可以選擇延長租用期、但是否需要說明 NFSv4 寬限期。

#### NFSv4 寬限期

如果 NFSv3 伺服器重新開機、幾乎可以立即為 IO 服務。它並未維持任何形式的用戶端狀態。結果是、ONTAP 接管作業通常似乎接近即時。當控制器準備好開始提供資料時、就會傳送 ARP 給網路、以表示拓撲的變化。客戶端通常幾乎立即檢測到這種情況、數據恢復流動。

不過 NFSv4 會短暫暫停。這只是 NFSv4 運作方式的一部分。

NFSv4 伺服器需要追蹤租用、鎖定、以及使用何種資料的人員。如果 NFS 伺服器出現問題並重新開機、或停電一段時間、或在維護活動期間重新啟動、則會導致租約 / 鎖定、而其他用戶端資訊也會遺失。伺服器需要先找出哪個用戶端正在使用哪些資料、才能恢復作業。這就是寬限期的開始。

如果您突然關閉 NFSv4 伺服器的電源、當恢復 IO 時、嘗試恢復 IO 的用戶端會收到回應、基本上說：「我遺失了租用 / 鎖定資訊。您是否要重新登錄鎖定？」這就是寬限期的開始。ONTAP 預設為 45 秒：

```
Cluster01::> nfs server show -vserver vserver1 -fields v4-grace-seconds

vserver    v4-grace-seconds
-----
vserver1   45
```

結果是、在重新啟動之後、控制器會暫停 IO、而所有用戶端都會回收租用和鎖定。寬限期結束後、伺服器將恢復 IO 作業。

#### 租用逾時與寬限期比較

寬限期與租用期間已連線。如上所述、預設的租用逾時為 30 秒、這表示 NFSv4 用戶端必須至少每 30 秒與伺服器簽入一次、否則就會遺失租約、進而導致鎖定。存在寬限期、可讓 NFS 伺服器重建租用 / 鎖定資料、預設為 45 秒。ONTAP 要求寬限期比租用期長 15 秒。如此可確保設計為至少每 30 秒續約的 NFS 用戶端環境、在重新



啟動後能夠與伺服器簽入。45 秒的寬限期可確保所有預期至少每 30 秒續約一次的客戶都有機會續約。

如果無法接受 30 秒的逾時、您可以選擇延長租用期。如果您想要將租用逾時延長至 60 秒、以便承受 60 秒的網路中斷、您必須將寬限期延長至至少 75 秒。ONTAP 要求比租用期高 15 秒。這表示您將會在控制器容錯移轉期間經歷更長的 IO 暫停時間。

這通常不是問題。一般使用者每年只會更新 ONTAP 控制器一或兩次、而且由於硬體故障而造成的非計畫性容錯移轉極少。此外、如果您的網路發生 60 秒網路中斷的可能性、而您需要將租用逾時時間延長至 60 秒、那麼您可能不會反對罕見的儲存系統容錯移轉、導致暫停時間也達 75 秒。您已確認網路暫停超過 60 秒、而且速度較快。

使用 **Oracle** 資料庫進行 **NFS** 快取

如果存在下列任一掛載選項、則會停用主機快取：

```
cio, actimeo=0, noac, forcedirectio
```

這些設定可能會嚴重影響軟體安裝、修補及備份 / 還原作業的速度。在某些情況下、尤其是叢集式應用程式、這些選項是必要的、因為必須在叢集中的所有節點之間提供快取一致性。在其他情況下、客戶誤用這些參數、結果是不必要的效能損害。

許多客戶在安裝或修補應用程式二進位檔時、會暫時移除這些掛載選項。如果使用者在安裝或修補程序過程中確認沒有其他處理程序正在使用目標目錄、則可安全地執行此移除。

**Oracle** 資料庫的 **NFS** 傳輸大小

根據預設、ONTAP 將 NFS I/O 大小限制為 64K。

大多數應用程式和資料庫的隨機 I/O 使用的區塊大小要小得多、遠低於 64K 的最大值。大型區塊 I/O 通常是平行處理的、因此 64K 的最大值也不是取得最大頻寬的限制。

有些工作負載的上限為 64K、因此會造成限制。特別是、如果資料庫執行的 I/O 數量較少、但容量較大、則備份或還原作業或資料庫完整表格掃描等單執行緒作業、會更快、更有效率地執行。ONTAP 的最佳 I/O 處理大小為 256k。

指定 ONTAP SVM 的最大傳輸大小可變更如下：

```
Cluster01::> set advanced
Warning: These advanced commands are potentially dangerous; use them only
when directed to do so by NetApp personnel.
Do you want to continue? {y|n}: y
Cluster01::*> nfs server modify -vserver vserver1 -tcp-max-xfer-size
262144
Cluster01::*>
```

## 注意

切勿將 ONTAP 上允許的傳輸大小上限降至低於目前掛載之 NFS 檔案系統的 rsize/wsize 值。這可能會在某些作業系統中造成當機或甚至資料毀損。例如、如果 NFS 用戶端目前設定為 rsize/wsize 65536、則 ONTAP 最大傳輸大小可在 65536 到 1048576 之間調整、因為用戶端本身受到限制、因此沒有任何影響。將傳輸大小上限降至 65536 以下可能會損害可用度或資料。

## Oracle 資料庫與 NVFAIL

NVFAIL 是 ONTAP 中的一項功能、可確保災難性容錯移轉案例期間的完整性。

資料庫在儲存設備容錯移轉事件期間容易受損、因為它們會維持大型內部快取。如果災難性事件需要強制 ONTAP 容錯移轉或強制 MetroCluster 切換、無論整體組態的健全狀況為何、都可能有效捨棄先前確認的變更。儲存陣列的內容會及時向後跳轉、而且資料庫快取的狀態不再反映磁碟上資料的狀態。這種不一致會導致資料毀損。

快取可能發生在應用程式或伺服器層。例如、Oracle Real Application Cluster (RAC) 組態、主站台和遠端站台上的伺服器都處於作用中狀態、可在 Oracle SGA 中快取資料。強制切入作業會導致資料遺失、因此資料庫可能會發生毀損、因為儲存在 SGA 中的區塊可能與磁碟上的區塊不符。

較不明顯的快取用途是在作業系統檔案系統層。來自掛載 NFS 檔案系統的區塊可能會快取到作業系統中。或者、以位於主要站台上的 LUN 為基礎的叢集式檔案系統、可以掛載到遠端站台的伺服器上、然後再次快取資料。在這些情況下、NVRAM 故障或強制接管或強制性的作業系統、可能會導致檔案系統毀損。

ONTAP 使用 NVFAIL 及其相關設定、保護資料庫和作業系統不受此案例影響。

## ASM 回收公用程式和 ONTAP 零區塊偵測

啟用即時壓縮時、ONTAP 可有效移除寫入檔案或 LUN 的歸零區塊。Oracle ASM 回收公用程式 (ASRU) 等公用程式的運作方式是將零寫入未使用的 ASM 範圍。

這可讓 DBA 在資料刪除後回收儲存陣列上的空間。ONTAP 會攔截零並取消分配 LUN 的空間。回收程序非常快速、因為儲存系統中沒有寫入資料。

從資料庫的角度來看、ASM 磁碟群組包含零、讀取 LUN 的這些區域會產生零串流、但 ONTAP 不會將零儲存在磁碟機上。而是進行簡單的中繼資料變更、在內部將 LUN 的歸零區域標記為任何資料的空白。

由於類似的原因、涉及零位資料的效能測試無效、因為零區塊實際上並未在儲存陣列內以寫入方式處理。



使用 ASRU 時、請確定已安裝所有 Oracle 建議的修補程式。

## Oracle 資料庫虛擬化

使用 VMware、Oracle OLVM 或 KVM 來虛擬化資料庫、對於選擇虛擬化技術的 NetApp 客戶而言、這是越來越常見的選擇、即使是他們最關鍵的關鍵任務資料庫也一樣。

### 支援能力

對於 Oracle 虛擬化支援政策、尤其是 VMware 產品、存在許多誤解。聽說 Oracle 完全不支援虛擬化並不罕

見。這個概念不正確、導致錯失從虛擬化中獲益的機會。Oracle Doc ID 249212.1 討論實際需求、客戶很少會將此視為疑慮。

如果虛擬化伺服器發生問題、而 Oracle Support 先前不知道該問題、可能會要求客戶在實體硬體上重現問題。執行產品尖端版本的 Oracle 客戶可能不想使用虛擬化技術、因為可能會發生支援問題、但這種情況對於虛擬化客戶而言、並不是一個真正的世界、因為他們使用的是 Oracle 產品版本。

## 儲存簡報

考慮將資料庫虛擬化的客戶應根據其業務需求做出儲存決策。雖然這是所有 IT 決策的一般陳述、但資料庫專案尤其重要、因為需求的大小和範圍差異極大。

儲存簡報有三個基本選項：

- Hypervisor 資料存放區上的虛擬化 LUN
- iSCSI LUN 由虛擬機器上的 iSCSI 啟動器管理、而非 Hypervisor
- 由虛擬機器掛載的 NFS 檔案系統（非從 NFS 型資料存放區）
- 直接裝置對應。客戶不喜歡 VMware RDM、但實體裝置通常與 KVM 和 OLVM 虛擬化類似。

## 效能

將儲存設備呈現給虛擬化來賓作業系統的方法通常不會影響效能。主機作業系統、虛擬化網路驅動程式和 Hypervisor 資料存放區實作均經過高度最佳化、只要遵循基本最佳實務做法、通常都能在 Hypervisor 和儲存系統之間使用所有可用的 FC 或 IP 網路頻寬。在某些情況下、使用一種儲存呈現方法比使用另一種方法來獲得最佳效能可能會稍微容易一些、但最終結果應該是可比較的。

## 管理能力

決定如何將儲存設備呈現給虛擬化來賓作業系統的關鍵因素是可管理性。沒有正確或錯誤的方法。最佳方法取決於 IT 營運需求、技能和偏好。

考量因素包括：

- \* 透明度。\* VM 管理其檔案系統時、資料庫管理員或系統管理員更容易識別其資料的檔案系統來源。檔案系統和 LUN 的存取方式與實體伺服器的存取方式完全相同。
- \* 一致性。\* VM 擁有其檔案系統時、使用或不使用 Hypervisor 層會影響管理能力。資源配置、監控、資料保護等程序也可在整個資產中使用、包括虛擬化和非虛擬化環境。

另一方面、在另一個 100% 虛擬化的資料中心中、最好還是在整個佔用空間中使用資料存放區型儲存設備、但前提是上述相同的理由（一致性）、也就是能夠使用相同的程序來進行資源配置、保護、監控和資料保護。

- \* 穩定性與疑難排解。\* VM 擁有其檔案系統時、由於整個儲存堆疊都存在於 VM 上、因此提供良好、穩定的效能與疑難排解問題變得更簡單。Hypervisor 唯一的角色是傳輸 FC 或 IP 框架。當資料存放區包含在組態中時、它會引入另一組逾時、參數、記錄檔和潛在錯誤、使組態複雜化。
- \* 可攜性。\* VM 擁有其檔案系統時、移動 Oracle 環境的程序會變得更簡單。檔案系統可在虛擬化與非虛擬化的來賓作業系統之間輕鬆移動。
- \* 廠商鎖定 \* 將資料放入資料存放區後、使用不同的 Hypervisor 或將資料從虛擬化環境中移出、將變得非常困難。

- \* 啟用 Snapshot : \* 虛擬化環境中的傳統備份程序可能會因為頻寬相對有限而成為問題。例如、四埠 10GbE 主幹可能足以支援許多虛擬化資料庫的日常效能需求、但這類主幹可能不足以使用 RMAN 或其他需要串流完整資料複本的備份產品來執行備份。因此、日益整合的虛擬化環境需要透過儲存快照來執行備份。如此可避免僅為了支援備份時間內的頻寬和 CPU 需求而需要過度建置 Hypervisor 組態。

使用來賓擁有的檔案系統有時會讓您更輕鬆地利用快照型備份和還原、因為需要保護的儲存物件可以更輕鬆地鎖定目標。然而、越來越多的虛擬化資料保護產品能夠與資料存放區和快照完美整合。在決定如何將儲存設備呈現給虛擬化主機之前、應充分考慮備份策略。

## 半虛擬化驅動程式

為了達到最佳效能、使用半虛擬化網路驅動程式至關重要。使用資料存放區時、需要半虛擬化 SCSI 驅動程式。半虛擬化的裝置驅動程式可讓來賓更深入地整合至 Hypervisor 、而非模擬的驅動程式、在該驅動程式中、Hypervisor 會花費更多的 CPU 時間來模擬實體硬體的行為。

## 過度使用 RAM

過度使用 RAM 意味著在不同主機上設定的虛擬化 RAM 多於實體硬體上的虛擬化 RAM 。否則可能會造成非預期的效能問題。虛擬化資料庫時、Oracle SGA 的基礎區塊不得由 Hypervisor 交換至儲存設備。這樣做會導致效能結果極不穩定。

## 資料存放區等量分割

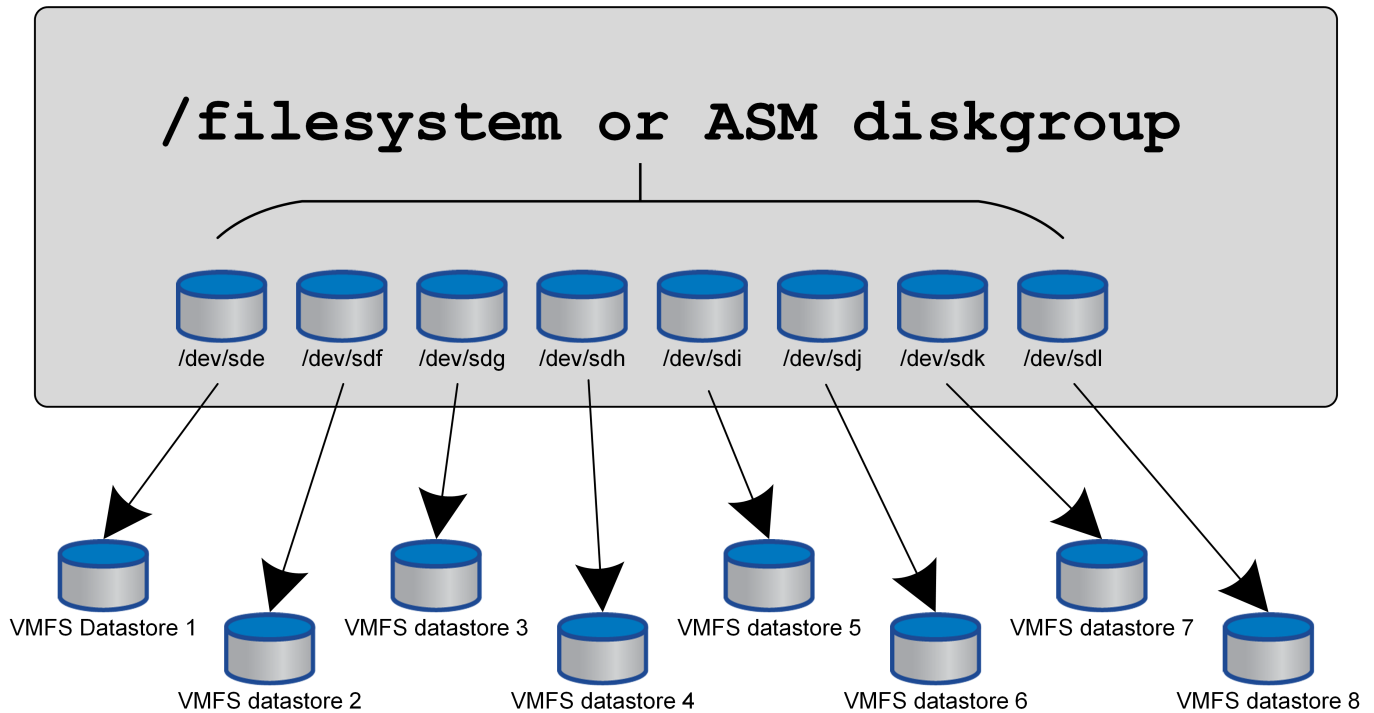
在使用資料庫搭配資料存放區時、效能方面有一個關鍵因素需要考量、那就是分段。

VMFS 等資料存放區技術可以跨越多個 LUN 、但它們不是等量磁區的裝置。LUN 會串聯。最終結果可能是 LUN 熱點。例如、典型的 Oracle 資料庫可能有 8-LUN ASM 磁碟群組。所有 8 個虛擬化 LUN 均可在 8-LUN VMFS 資料存放區上進行佈建、但無法保證資料所在的 LUN 。產生的組態可能是全部 8 個虛擬化 LUN 、佔用 VMFS 資料存放區內的單一 LUN 。這會成為效能瓶頸。

通常需要分拆。有些 Hypervisor (包括 KVM) 可以使用 LVM 等量分拆來建置資料存放區、如前所述 ["請按這裡"](#)。有了 VMware 、架構看起來有點不同。每個虛擬化 LUN 都必須放置在不同的 VMFS 資料存放區上。

例如：

## Virtualized host



這種方法的主要驅動因素不是 ONTAP、而是因為單一 VM 或 Hypervisor LUN 可平行處理的作業數量、有固有的限制。單一 ONTAP LUN 通常支援的 IOPS 遠高於主機所能要求的 IOPS。單一 LUN 效能限制幾乎是主機作業系統的普遍結果。結果是大多數資料庫需要 4 到 8 個 LUN 才能滿足效能需求。

VMware 架構需要仔細規劃其架構、以確保此方法不會遇到資料存放區和 / 或 LUN 路徑最大化問題。此外、每個資料庫都不需要一組唯一的 VMFS 資料存放區。主要需求是確保每個主機都有一組乾淨的 4 到 8 個 IO 路徑、從虛擬化 LUN 到儲存系統本身的後端 LUN。在極少數情況下、更多的資料存取器可能對真正極致的效能需求有所助益、但 4-8 個 LUN 通常足以滿足 95% 的資料庫需求。單一 ONTAP 磁碟區包含 8 個 LUN、可透過典型的 OS/ONTAP/ 網路組態、支援多達 250,000 個隨機 Oracle 區塊 IOPS。

## 分層

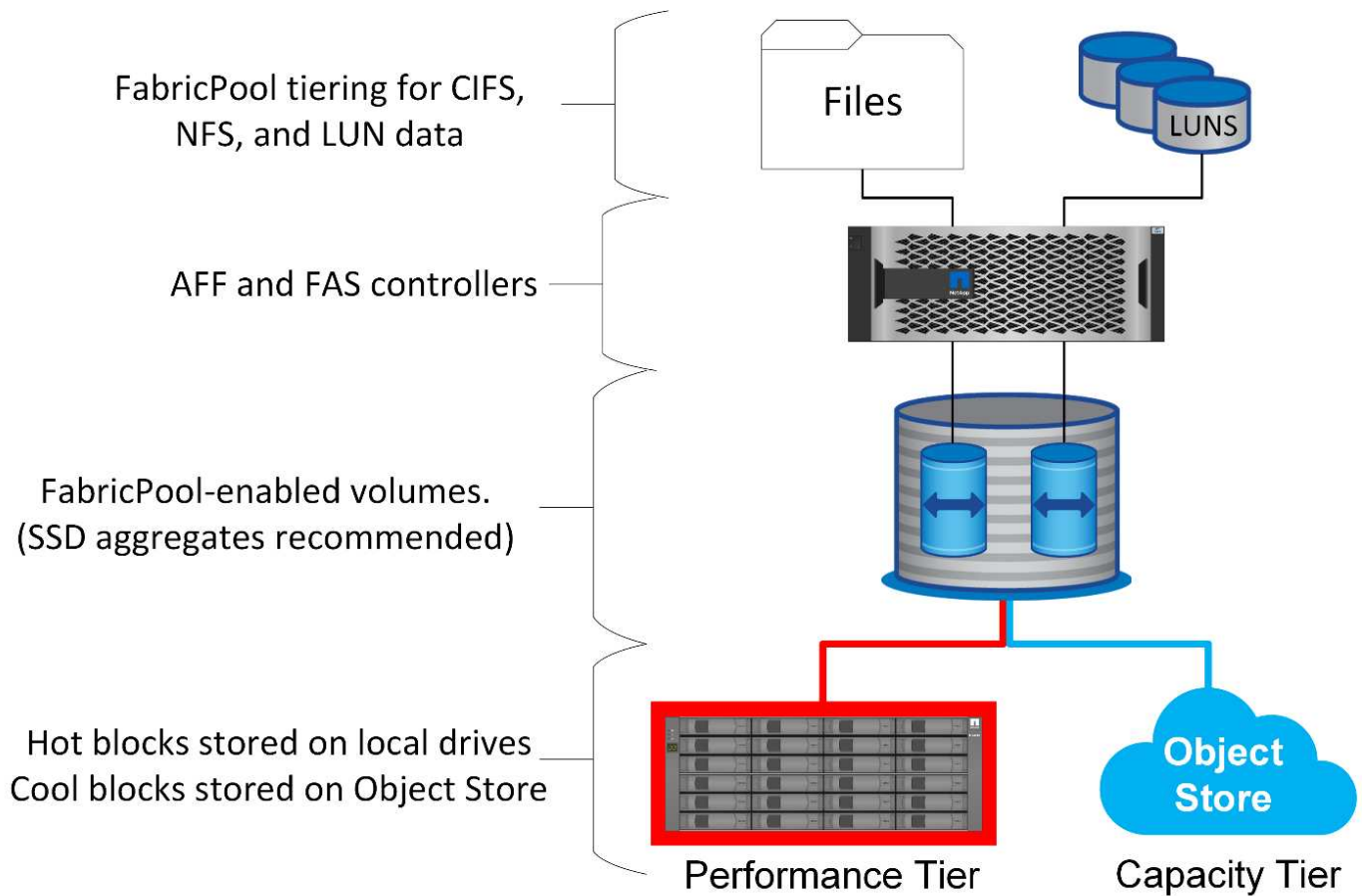
### Oracle 資料庫 FabricPool 分層概述

瞭解 FabricPool 分層對 Oracle 和其他資料庫的影響、需要瞭解低階 FabricPool 架構。

#### 架構

FabricPool 是一項分層技術、將區塊歸類為熱區塊或冷卻區塊、並將其置於最適當的儲存層。效能層最常位於 SSD 儲存設備上、並裝載熱資料區塊。容量層位於物件存放區、並裝載 Cool 資料區塊。物件儲存支援包括 NetApp StorageGRID、ONTAP S3、Microsoft Azure Blob 儲存設備、Alibaba 雲端物件儲存服務、IBM 雲端物件儲存設備、Google 雲端儲存設備和 Amazon AWS S3。

有多個分層原則可供使用、可控制區塊如何分類為熱或冷區、並可針對每個磁碟區設定原則、並視需要進行變更。只有資料區塊會在效能層和容量層之間移動。定義 LUN 和檔案系統結構的中繼資料一律保留在效能層上。因此、管理作業會集中在 ONTAP 上。檔案和 LUN 與儲存在任何其他 ONTAP 組態上的資料並無不同。NetApp AFF 或 FAS 控制器會套用定義的原則、將資料移至適當的層級。



### 物件存放區提供者

物件儲存傳輸協定使用簡單的 HTTP 或 HTTPS 要求來儲存大量的資料物件。物件儲存設備的存取必須可靠、因為 ONTAP 的資料存取取決於要求的即時服務。選項包括 Amazon S3 Standard 和 In常用 存取選項、以及 Microsoft Azure Hot 和 Cool Blob Storage、IBM Cloud 和 Google Cloud。不支援 Amazon Glacier 和 Amazon Archive 等歸檔選項、因為擷取資料所需的時間可能超過主機作業系統和應用程式的容許值。

NetApp StorageGRID 也受到支援、是最佳的企業級解決方案。它是高效能、可擴充且高度安全的物件儲存系統、可為 FabricPool 資料及其他物件儲存應用程式提供地理上的備援、這些應用程式越來越可能成為企業應用程式環境的一部分。

StorageGRID 也能避免許多公有雲供應商為了從服務讀取資料而收取的出口費用、進而降低成本。

### 資料與中繼資料

請注意、此處的「資料」一詞適用於實際的資料區塊、而非中繼資料。只有資料區塊是階層式的、而中繼資料則保留在效能層級中。此外、只有讀取實際的資料區塊、才能影響區塊的「熱」或「冷」狀態。只要讀取檔案的名稱、時間戳記或所有權中繼資料、就不會影響基礎資料區塊的位置。

### 備份

雖然 FabricPool 可以大幅減少儲存佔用空間、但它本身並不是備份解決方案。NetApp WAFL 中繼資料一律保持在效能層級。如果災難性災難破壞效能層、則無法使用容量層上的資料來建立新環境、因為該環境不包含 WAFL 中繼資料。

不過、FabricPool 可以成為備份策略的一部分。例如、FabricPool 可以使用 NetApp SnapMirror 複寫技術進行

設定。每一半的鏡像都可以有自己與物件儲存目標的連線。結果是兩個資料的複本。主要複本包含效能層上的區塊、以及容量層中的相關區塊、而複本則是第二組效能和容量區塊。

## 分層原則

### Oracle 資料庫 FabricPool 分層原則

ONTAP 提供四項原則、可控制效能層上的 Oracle 資料如何成為移轉至容量層的候選對象。

#### 僅限 Snapshot

- `snapshot-only tiering-policy` 僅適用於未與作用中檔案系統共用的區塊。它基本上會導致資料庫備份分層。在建立快照之後、區塊會成為分層的候選項目、然後區塊會被覆寫、導致區塊只存在於快照中。A 之前的延遲 `snapshot-only` 區塊視為冷區、由控制 `tiering-minimum-cooling-days` 音量設定。ONTAP 9.8 的範圍為 2 至 183 天。

許多資料集的變更率都很低、因此這項原則可節省的成本極低。例如、在 ONTAP 上觀察到的典型資料庫每週變更率低於 5%。資料庫歸檔記錄檔可能佔用大量空間、但通常會繼續存在於作用中的檔案系統中、因此不會成為根據此原則分層的候選項目。

#### 自動

- `auto` 分層原則可將分層延伸至快照專用區塊、以及作用中檔案系統內的區塊。區塊冷卻前的延遲由控制 `tiering-minimum-cooling-days` 音量設定。ONTAP 9.8 的範圍為 2 至 183 天。

此方法可啟用無法與搭配使用的分層選項 `snapshot-only` 原則。例如、資料保護原則可能需要保留 90 天的特定記錄檔。如果將冷卻期設定為 3 天、則任何超過 3 天的記錄檔都會從效能層中分層移出。此動作可釋放效能層級上的大量空間、同時仍可讓您檢視及管理完整的 90 天資料。

#### 無

- `none` 分層原則可防止任何額外的區塊從儲存層分層、但容量層中的任何資料仍會保留在容量層中、直到讀取為止。如果接著讀取區塊、則會將其拉回並放置在效能層上。

使用的主要原因 `none` 分層原則旨在防止區塊分層、但隨著時間推移而變更原則可能會很有用。例如、假設某個特定資料集已廣泛分層至容量層、但卻產生對完整效能功能的非預期需求。可變更原則以防止任何額外的分層、並確認隨 IO 增加而讀取的任何區塊仍保留在效能層中。

#### 全部

- `all` 分層原則取代了 `backup` 原則自 ONTAP 9.6 起。◦ `backup` 原則僅套用至資料保護磁碟區、意指 SnapMirror 或 NetApp SnapVault 目的地。◦ `all` 原則的功能相同、但不限於資料保護磁碟區。

有了這項原則、就能立即將區塊視為酷炫、並立即分層至容量層。

此原則特別適用於長期備份。它也可以用作階層式儲存管理 (HSM) 的形式。過去、HSM 通常用於將檔案的資料區塊分層至磁帶、同時讓檔案本身在檔案系統上保持可見。具有的 FabricPool Volume `all` 原則可讓您將檔案儲存在可見且可管理的環境中、但幾乎無需佔用本機儲存層的空間。

## Oracle 資料庫和 FabricPool 擷取原則

分層原則可控制哪些 Oracle 資料庫區塊從效能層分層分層到容量層。擷取原則可控制已階層的區塊讀取時所發生的情況。

### 預設

所有 FabricPool 磁碟區的初始設定為 `default` 這表示行為是由「雲端擷取原則」所控制。「確切的行為取決於所使用的分層原則。

- `auto`—僅擷取隨機讀取的資料
- `snapshot-only`—擷取所有依序或隨機讀取的資料
- `none`—擷取所有依序或隨機讀取的資料
- `all`—請勿從容量層擷取資料

### 讀取中

設定 `cloud-retrieval-policy` 對讀取會覆寫預設行為、因此讀取任何階層資料會導致資料傳回效能層。

例如、在的下、一個 Volume 可能已被輕度使用了很長時間 `auto` 分層原則和大多數區塊現在都是分層的。

如果業務發生非預期的變更、需要重複掃描部分資料以準備特定報告、則可能需要變更 `cloud-retrieval-policy` 至 `on-read` 確保讀取的所有資料都會傳回效能層、包括依序和隨機讀取資料。這將改善連續 I/O 相對於磁碟區的效能。

### 推廣

升級原則的行為取決於分層原則。如果分層原則是 `auto`、然後設定 `cloud-retrieval-policy` 到 `promote` 在下次分層掃描時、從容量層移回所有區塊。

如果分層原則是 `snapshot-only`，則只會傳回與作用中檔案系統相關聯的區塊。通常這不會有任何影響、因為唯一的區塊會在底下分層 `snapshot-only` 原則將是與快照完全相關的區塊。作用中檔案系統中不會有階層式區塊。

然而、如果磁碟區上的資料是由 Volume SnapRestore 或快照的檔案複製作業還原、則目前使用中檔案系統可能需要一些因為只與快照相關而分層的區塊。您可能需要暫時變更 `cloud-retrieval-policy` 原則目標 `promote` 快速擷取所有本機需要的區塊。

### 永不

請勿從容量層擷取區塊。

## 分層策略

### Oracle 資料庫完整檔案 FabricPool 分層

雖然 FabricPool 分層是在區塊層級運作、但在某些情況下、它可用於提供檔案層級的分層。



許多應用程式資料集都是按日期組織、而且隨著時間的演進、存取這些資料的可能性通常會越來越小。例如、某家銀行可能會有包含五年客戶對帳單的 PDF 檔案儲存庫、但只有最近幾個月是有效的。FabricPool 可用於將較舊的資料檔案重新放置到容量層。14 天的冷卻期可確保最近 14 天的 PDF 檔案仍保留在效能層級上。此外、至少每 14 天讀取一次的檔案仍會很熱、因此會保留在效能層級上。

#### 原則

若要實作檔案型分層方法、您必須有已寫入且未隨後修改的檔案。tiering-minimum-cooling-days 原則應該設定得足夠高、以便您可能需要的檔案保留在效能層。例如、需要最近 60 天資料的資料集、且需要設定最佳效能 tiering-minimum-cooling-days 期間為 60。也可以根據檔案存取模式來達成類似的結果。例如、如果需要最近 90 天的資料、而應用程式正在存取該 90 天的資料範圍、則資料會保留在效能層。設定 tiering-minimum-cooling-days 期間為 2、資料變少後、系統會提示您分層。

- auto 由於只有 auto 原則會影響作用中檔案系統中的區塊。



任何類型的資料存取都會重設熱圖資料。病毒掃描、索引甚至是讀取來源檔案的備份活動、都會因為必要而防止分層 tiering-minimum-cooling-days 從未達到臨界值。

#### Oracle 部分檔案 FabricPool 分層

由於 FabricPool 是在區塊層級運作、因此可能變更的檔案可以部分分層化為物件儲存、而部分保留在效能層級上。

這在資料庫中很常見。已知包含非作用中區塊的資料庫也可用於 FabricPool 分層。例如、供應鏈管理資料庫可能包含歷史資訊、這些資訊必須在必要時提供、但在正常作業期間無法存取。FabricPool 可用於選擇性地重新定位非使用中的區塊。

例如、在 FabricPool 磁碟區上執行的資料檔案 tiering-minimum-cooling-days 90 天的期間會保留過去 90 天內在效能層級上存取的任何區塊。然而、90 天內未存取的任何項目都會重新移至容量層。在其他情況下、正常的應用程式活動會將正確的區塊保留在正確的層級上。例如、如果資料庫通常用於定期處理過去 60 天的資料、則會降低許多 tiering-minimum-cooling-days 可以設定期間、因為應用程式的自然活動可確保區塊不會提早重新定位。

- auto 原則應與資料庫一起使用。許多資料庫都有定期活動、例如季末流程或重新編製索引作業。如果這些作業的期間大於 tiering-minimum-cooling-days 可能會發生效能問題。例如、如果季度末處理需要 1TB 的資料、而這些資料原本沒有受到影響、則該資料現在可能會出現在容量層。從容量層讀取的速度通常極快、可能不會造成效能問題、但實際結果將取決於物件儲存區組態。

#### 原則

◦ tiering-minimum-cooling-days 原則的設定應足夠高、以保留效能層上可能需要的檔案。例如、如果資料庫需要最新 60 天的資料、而且效能最佳、則需要設定 tiering-minimum-cooling-days 期間為 60 天。也可以根據檔案的存取模式來達成類似的結果。例如、如果需要最近 90 天的資料、而應用程式正在存取該 90 天的資料範圍、則資料會保留在效能層。設定 tiering-minimum-cooling-days 在資料變得不活躍之後、將會立即將資料分級至 2 天。

- auto 由於只有 auto 原則會影響作用中檔案系統中的區塊。



任何類型的資料存取都會重設熱圖資料。因此、資料庫完整表格掃描、甚至是讀取來源檔案的備份活動、都會因為需要而防止分層 tiering-minimum-cooling-days 從未達到臨界值。

FabricPool 最重要的用途可能是提高已知冷資料的效率、例如資料庫交易記錄。

大部分的關聯式資料庫都是以交易記錄歸檔模式運作、以提供時間點還原。對資料庫所做的變更會記錄交易記錄中的變更、並保留交易記錄而不會被覆寫。結果可能需要保留大量歸檔的交易記錄檔。許多其他應用程式工作流程也有類似的例子、這些工作流程會產生必須保留的資料、但很難存取。

FabricPool 透過提供整合式分層的單一解決方案來解決這些問題。檔案會儲存並保留在其一般位置、但實際上不會佔用主要陣列上的任何空間。

### 原則

使用 `tiering-minimum-cooling-days` 幾天的原則會導致在效能層上保留最近建立的檔案（近期最可能需要的檔案）中的區塊。之後、舊檔案的資料區塊會移至容量層。

- `auto` 無論主要檔案系統中的記錄是否已刪除或繼續存在、都會在達到冷卻臨界值時強制執行提示分層。將所有可能需要的記錄儲存在作用中檔案系統的單一位置、也能簡化管理。沒有理由搜尋快照以找出需要還原的檔案。

某些應用程式（例如 Microsoft SQL Server）會在備份作業期間截斷交易記錄檔、使記錄不再位於作用中的檔案系統中。使用可節省容量 `snapshot-only` 分層原則、但 `auto` 原則對記錄資料並不實用、因為作用中檔案系統中應該很少會冷卻記錄資料。

## Oracle 與 FabricPool 快照分層

FabricPool 的初始版本以備份使用案例為目標。唯一可以分層的區塊類型是不再與作用中檔案系統中的資料相關聯的區塊。因此、只能將快照資料區塊移至容量層。當您需要確保效能不受影響時、這仍然是最安全的分層選項之一。

### 原則 - 本機快照

有兩個選項可將非作用中的快照區塊分層到容量層。首先 `snapshot-only` 原則僅針對快照區塊。儘管如此 `auto` 原則包括 `snapshot-only` 區塊、也會將區塊分層、從作用中的檔案系統中移出。這可能不理想。

- `tiering-minimum-cooling-days` 值應設為時間週期、以便在效能層上提供還原期間所需的資料。例如、重要正式作業資料庫的大多數還原案例、都會在過去幾天的某個時間加入還原點。設定 `tiering-minimum-cooling-days` 值 3 可確保檔案的任何還原都會產生可立即提供最大效能的檔案。作用中檔案中的所有區塊仍會顯示在快速儲存設備上、而無需從容量層恢復。

### 原則 - 複寫的快照

使用 SnapMirror 或 SnapVault 複寫的僅用於恢復的快照通常應使用 FabricPool `all` 原則。使用此原則、會複寫中繼資料、但所有資料區塊都會立即傳送至容量層、以獲得最大效能。大部分的恢復程序都涉及循序 I/O、這是固有的效率。應該評估物件存放區目的地的恢復時間、但在設計完善的架構中、此恢復程序不需要比從本機資料恢復慢很多。

如果複製的資料也要用於複製、則會使用 `auto` 原則比較適當、請使用 `tiering-minimum-cooling-days` 包含預期在複製環境中經常使用的資料的值。例如、資料庫的作用中工作集可能包含前三天讀取或寫入的資料、但也可能包含另外 6 個月的歷史資料。如果是、則是 `auto` SnapMirror 目的地的原則可讓工作集在效能層上使用。

傳統應用程式備份包括 Oracle Recovery Manager 等產品、可在原始資料庫之外建立檔案型備份。

```
`tiering-minimum-cooling-days` policy of a few days preserves the most recent backups, and therefore the backups most likely to be required for an urgent recovery situation, on the performance tier. The data blocks of the older files are then moved to the capacity tier.
```

- ``auto``

原則是最適合備份資料的原則。如此可確保在達到冷卻臨界值時、無論主要檔案系統中的檔案是否已刪除或繼續存在、都能立即分層。將所有可能需要的檔案儲存在作用中檔案系統的單一位置、也能簡化管理。沒有理由搜尋快照以找出需要還原的檔案。

- `snapshot-only` 原則可以生效、但該原則僅適用於不在作用中檔案系統中的區塊。因此、必須先刪除 NFS 或 SMB 共用上的檔案、才能分層化資料。

使用 LUN 組態時、此原則的效率會更低、因為從 LUN 刪除檔案只會從檔案系統中繼資料中移除檔案參照。LUN 上的實際區塊會一直保留到位、直到被覆寫為止。這種情況可能會造成從刪除檔案到覆寫區塊並成為分層候選項目之間的長時間延遲。移動有一些好處 `snapshot-only` 區塊到容量層、但整體而言、備份資料的 FabricPool 管理最適合搭配使用 `auto` 原則。



此方法可協助使用者更有效率地管理備份所需的空間、但 FabricPool 本身並不是備份技術。將備份檔案分層至物件存放區可簡化管理、因為檔案仍可在原始儲存系統上看到、但物件存放區目的地中的資料區塊則取決於原始儲存系統。如果來源磁碟區遺失、物件儲存區資料將無法再使用。

## Oracle 資料庫和物件儲存區存取中斷

使用 FabricPool 分層資料集會導致主要儲存陣列與物件存放區層之間的相依性。有許多物件儲存選項可提供不同層級的可用度。請務必瞭解主要儲存陣列與物件儲存層之間可能中斷連線的影響。

如果發給 ONTAP 的 I/O 需要容量層的資料、而 ONTAP 無法到達容量層來擷取區塊、則 I/O 最終會逾時。此逾時的影響取決於所使用的傳輸協定。在 NFS 環境中、ONTAP 會根據傳輸協定回應 EJUKEBOX 或 EDELAY 回應。有些較舊的作業系統可能會將此視為錯誤、但 Oracle Direct NFS 用戶端目前的作業系統和修補程式層級會將此視為可重試的錯誤、並繼續等待 I/O 完成。

較短的逾時時間適用於 SAN 環境。如果需要物件存放區環境中的區塊、而且無法連線兩分鐘、則會將讀取錯誤傳回主機。ONTAP 磁碟區和 LUN 保持連線、但主機作業系統可能會將檔案系統標示為處於錯誤狀態。

物件儲存連線問題 `snapshot-only` 由於只有備份資料是階層式的、因此原則並不令人擔心。通訊問題會拖慢資料恢復速度、但不會影響使用中的資料。◦ `auto` 和 `all` 原則允許從作用中 LUN 分層處理冷資料、這表示物件儲存區資料擷取期間發生的錯誤可能會影響資料庫可用度。採用這些原則的 SAN 部署只能搭配專為高可用度所設計的企業級物件儲存設備和網路連線使用。NetApp StorageGRID 是絕佳的選擇。

## Oracle 資料保護

## 使用 ONTAP 保護 Oracle 資料

NetApp 知道資料庫中有最重要的關鍵任務資料。

企業無法在沒有資料存取權的情況下運作、有時資料會定義業務。此資料必須受到保護；然而、資料保護不只是確保可用的備份、還要快速可靠地執行備份、而且還要安全地儲存。

資料保護的另一面是資料恢復。當資料無法存取時、企業就會受到影響、而且可能無法運作、直到資料還原為止。此程序必須快速且可靠。最後、大多數資料庫都必須防範災難、這表示必須維護資料庫的複本。複本必須是最新的。複本也必須快速且簡單、才能使其成為完全運作的資料庫。



本文件取代先前發佈的技術報告 [\\_TR-4591](#)：Oracle 資料保護：備份、還原及複寫。

### 規劃

正確的企業資料保護架構、取決於資料保留、可恢復性、以及各種事件中斷的容忍度等業務需求。

例如、請考慮範圍內的應用程式、資料庫和重要資料集數量。為單一資料集建立備份策略、以確保符合典型 SLA 的要求相當簡單、因為沒有太多物件需要管理。隨著資料集數量增加、監控作業變得更複雜、系統管理員可能不得不花費更多時間來解決備份故障。當環境達到雲端與服務供應商的規模時、需要完全不同的方法。

資料集大小也會影響策略。例如、由於資料集太小、因此有許多選項可用於 100GB 資料庫的備份與還原。只要使用傳統工具從備份媒體複製資料、通常就能提供足夠的恢復 RTO。100TB 資料庫通常需要完全不同的策略、除非 RTO 允許多天中斷、否則傳統的複製備份與還原程序可能是可以接受的。

最後、除了備份與還原程序本身之外、還有其他因素。例如、是否有支援關鍵生產活動的資料庫、使恢復成為僅由熟練的 DBA 執行的罕見事件？或者、資料庫是否是大型開發環境的一部分、而在大型開發環境中、恢復是經常發生的事、並由一群通才的 IT 團隊負責管理？

## Oracle 資料庫 RTO、RPO 和 SLA 規劃

ONTAP 可讓您輕鬆根據業務需求量身打造 Oracle 資料庫資料保護策略。

這些需求包括恢復速度、最大允許資料遺失量、以及備份保留需求等因素。資料保護計畫也必須考量資料保留與還原的各種法規要求。最後、必須考量不同的資料恢復情境、從使用者或應用程式錯誤所造成的典型和可預見的恢復、到包括站點完全遺失的災難恢復情境。

資料保護與還原原則的細微變更、可能會對儲存、備份與還原的整體架構造成重大影響。在開始設計工作之前、務必先定義並記錄標準、以免資料保護架構複雜化。不必要的功能或保護層級會導致不必要的成本和管理成本、而最初忽略的需求可能導致專案方向錯誤、或是需要最後一分鐘的設計變更。

### 恢復時間目標

恢復時間目標（RTO）定義了恢復服務所允許的最長時間。例如、人力資源資料庫的 RTO 可能為 24 小時、因為雖然在工作天內無法存取這些資料、但企業仍可繼續營運。相反地、支援銀行總分類帳的資料庫、其 RTO 會以分鐘甚至秒為單位來衡量。RTO 為零是不可能的、因為必須有辦法區分實際服務中斷和例行事件、例如遺失的網路封包。然而、典型的要求是接近零的 RTO。

### 恢復點目標

恢復點目標（RPO）定義了最大可容忍的資料遺失量。在許多情況下、RPO 完全取決於快照或 SnapMirror 更

新的頻率。

在某些情況下、RPO 可能會變得更具侵略性、因此可選擇性地更頻繁地保護某些資料。在資料庫內容中、RPO 通常是在特定情況下可能遺失多少記錄資料的問題。在資料庫因產品錯誤或使用者錯誤而受損的典型還原案例中、RPO 應為零、表示不應有資料遺失。恢復程序包括還原較早的資料庫檔案複本、然後重新播放記錄檔、將資料庫狀態提升至所需的時間點。此作業所需的記錄檔應已在原始位置中。

在不尋常的情況下、記錄資料可能會遺失。例如、意外或惡意 `rm -rf *` 資料庫檔案可能會導致刪除所有資料。唯一的選擇是從備份還原、包括記錄檔、有些資料必然會遺失。在傳統備份環境中改善 RPO 的唯一選項是執行記錄資料的重複備份。然而、由於資料持續移動、而且難以將備份系統維持為持續運作的服務、因此這項功能也有其侷限性。進階儲存系統的優點之一、就是能夠保護資料、避免意外或惡意損壞檔案、進而在不移動資料的情況下提供更好的 RPO。

## 災難恢復

災難恢復包括在發生實體災難時恢復服務所需的 IT 架構、原則和程序。這可能包括洪水、火災或惡意或疏忽意圖的人。

災難恢復不只是一套恢復程序。這是識別各種風險、定義資料恢復和服務持續性需求、以及提供適當架構及相關程序的完整程序。

在建立資料保護需求時、必須區分典型的 RPO 和 RTO 需求、以及災難恢復所需的 RPO 和 RTO 需求。有些應用程式環境需要零 RPO 和近乎零的 RTO、才能因應資料遺失的情況、從相對正常的使用者錯誤到破壞資料中心的火災。然而、這些高層級的保護措施會產生成本和管理上的後果。

一般而言、非災難性資料恢復需求應嚴格、原因有兩個。首先、應用程式錯誤和使用者錯誤會導致資料受損、幾乎是不可避免的。其次、只要儲存系統未遭銷毀、就不難設計出可提供零 RPO 和低 RTO 的備份策略。沒有理由不解決容易補救的重大風險、這就是為何用於本機恢復的 RPO 和 RTO 目標應該積極主動的原因。

災難恢復 RTO 和 RPO 需求因災難可能性及相關資料遺失或業務中斷所造成的後果而異。RPO 和 RTO 需求應以實際業務需求為基礎、而非一般原則。他們必須考慮多種邏輯和實體災難案例。

## 邏輯災難

邏輯災難包括使用者、應用程式或作業系統錯誤所造成的資料毀損、以及軟體故障。邏輯災難也可能包括由外部人士利用病毒或蠕蟲或利用應用程式弱點發動的惡意攻擊。在這些情況下、實體基礎架構並未受損、但基礎資料已不再有效。

越來越常見的邏輯災難類型稱為勒索軟體、其中攻擊模式是用來加密資料。加密不會損壞資料、但會在付款給第三方之前無法使用。越來越多的企業被勒索軟體攻擊的目標鎖定。針對此威脅、NetApp 提供防竄改快照、即使儲存管理員也無法在設定的到期日之前變更受保護的資料。

## 實體災難

實體災難包括基礎架構元件故障、其超過備援功能、導致資料遺失或服務延長中斷。例如、RAID 保護可提供磁碟機備援、而使用 HBA 則可提供 FC 連接埠和 FC 纜線備援。這類元件的硬體故障是可預見的、不會影響可用度。

在企業環境中、通常可以使用備援元件來保護整個站台的基礎架構、直到唯一可預見的實體災難案例完全遺失站台為止。災難恢復規劃則取決於站台對站台的複寫。

在理想的環境中、所有資料都會在地理上分散的站台之間同步複寫。這類複寫不一定可行、甚至可能、原因有幾個：

- 同步複寫不可避免地會增加寫入延遲、因為所有變更都必須複寫到兩個位置、應用程式 / 資料庫才能繼續處理。因此產生的效能影響有時是不可接受的、無法使用同步鏡射。
- 100% SSD 儲存設備的採用率增加、意味著更有可能注意到額外的寫入延遲、因為效能期望包括數十萬 IOPS 和低於毫秒的延遲。若要充分發揮 100% SSD 的效益、可能需要重新規劃災難恢復策略。
- 資料集會繼續以位元組為單位增加、因此必須確保足夠的頻寬來維持同步複寫、因此會產生挑戰。
- 資料集的複雜度也隨之增加、管理大規模同步複寫也會帶來挑戰。
- 雲端型策略通常需要較長的複寫距離和延遲、進一步排除同步鏡射的使用。

NetApp 提供的解決方案包括同步複寫、可滿足最嚴苛的資料恢復需求、並可提供更優異的效能與靈活度的非同步解決方案。此外、NetApp 技術也能與許多第三方複寫解決方案無縫整合、例如 Oracle DataGuard

### 保留時間

資料保護策略的最後一個層面是資料保留時間、這可能會大幅改變。

- 一般要求是在主要站台上進行 14 天夜間備份、以及儲存在次要站台上的 90 天備份。
- 許多客戶會建立獨立的季度歸檔、儲存在不同的媒體上。
- 持續更新的資料庫可能不需要歷史資料、而備份只需保留幾天。
- 法規要求可能需要在 365 天內恢復任何任意交易的點。

## ONTAP 提供 Oracle 資料庫可用度

ONTAP 旨在提供最大的 Oracle 資料庫可用度。ONTAP 高可用度功能的完整說明不在本文件的範圍之內。然而、與資料保護一樣、在設計資料庫基礎架構時、對此功能的基本瞭解非常重要。

### HA 配對

高可用度的基本單位是 HA 配對。每對都包含備援連結、可支援將資料複寫到 NVRAM。NVRAM 不是寫入快取。控制器內的 RAM 會做為寫入快取。NVRAM 的用途是暫時記錄資料、以防止發生非預期的系統故障。在這方面、它與資料庫重做記錄類似。

NVRAM 和資料庫重做記錄都可用來快速儲存資料、讓資料的變更能夠儘快提交。磁碟機（或資料檔案）上的持續資料更新直到稍後在 ONTAP 和大多數資料庫平台上稱為檢查點的程序期間才會進行。正常作業期間不會讀取 NVRAM 資料或資料庫重做記錄。

如果控制器突然故障、可能會有擱置中的變更、這些變更儲存在 NVRAM 中、但尚未寫入磁碟機。合作夥伴控制器會偵測故障、控制磁碟機、並套用儲存在 NVRAM 中的必要變更。

### 接管與恢復

接管與恢復是指在 HA 配對中的節點之間轉移儲存資源責任的程序。接管和恢復有兩個層面：

- 管理可存取磁碟機的網路連線能力
- 管理磁碟機本身

支援 CIFS 和 NFS 流量的網路介面、都是設定在主位置和容錯移轉位置。接管包括將網路介面移至與原始位置位於同一子網路的實體介面上的暫存主目錄。贈品包括將網路介面移回其原始位置。您可以視需要調整確切行為。

支援 SAN 區塊傳輸協定（例如 iSCSI 和 FC）的網路介面不會在接管和恢復期間重新定位。而是應使用包含完整 HA 配對的路徑來佈建 LUN、以產生主要路徑和次要路徑。



您也可以設定其他控制器的路徑、以支援在較大叢集中的節點之間重新放置資料、但這並不屬於 HA 程序的一部分。

接管與恢復的第二個層面是磁碟擁有權的轉移。確切的程序取決於多種因素、包括接管 / 恢復的原因、以及發出的命令列選項。目標是盡可能有效率地執行作業。雖然整體程序可能需要幾分鐘的時間、但磁碟機的實際擁有權從節點移轉至節點的時間通常只需幾秒鐘。

### 接管時間

主機 I/O 在接管和恢復作業期間會短暫暫停 I/O、但在正確設定的環境中不應發生應用程式中斷。I/O 延遲的實際轉換程序通常是以秒為單位來測量、但主機可能需要額外的時間來識別資料路徑的變更並重新提交 I/O 作業。

中斷的性質取決於傳輸協定：

- 支援 NFS 和 CIFS 流量的網路介面會在移轉至新實體位置後、向網路發出位址解析傳輸協定（ARP）要求。這會導致網路交換器更新其媒體存取控制（MAC）位址表、並繼續處理 I/O 在計畫性接管和恢復的情況下、中斷通常以秒為單位來衡量、在許多情況下都無法偵測到。有些網路可能會較慢、無法完全辨識網路路徑的變更、有些作業系統可能會在很短的時間內排入大量 I/O、因此必須重新嘗試。這可能會延長恢復 I/O 所需的時間
- 支援 SAN 通訊協定的網路介面不會轉換到新位置。主機作業系統必須變更使用中的路徑。主機觀察到 I/O 暫停的情形取決於多種因素。從儲存系統的角度來看、無法提供 I/O 的時間只有幾秒鐘。不過、不同的主機作業系統可能需要額外的時間、才能讓 I/O 逾時、再重試。較新的作業系統更能更快辨識路徑變更、但較舊的作業系統通常需要 30 秒才能辨識變更。

下表顯示儲存系統無法將資料提供給應用程式環境的預期接管時間。在任何應用程式環境中都不應發生任何錯誤、而是在 IO 處理過程中、接管應該會顯示為短暫的暫停。

	NFS	AFF	ASA
計畫性接管	15 秒	6-10 秒	2-3 秒
非計畫性接管	30 秒	6-10 秒	2-3 秒

## Checksum 與 Oracle 資料庫完整性

ONTAP 及其支援的通訊協定包含多項保護 Oracle 資料庫完整性的功能、包括靜態資料和透過網路傳輸的資料。

ONTAP 內的邏輯資料保護包含三項關鍵需求：

- 資料必須受到保護、以免資料毀損。

- 資料必須受到保護、避免磁碟機故障。
- 資料變更必須受到保護、以免遺失。

以下各節將討論這三項需求。

### 網路毀損：Checksum

最基本的資料保護層級是 Checksum、這是儲存在資料旁的特殊錯誤偵測程式碼。在網路傳輸期間、會使用 Checksum 和（在某些情況下）多個 Checksum 來偵測資料毀損。

例如、FC 框架包含稱為循環備援檢查（CRC）的校驗和形式、以確保有效負載不會在傳輸過程中毀損。傳輸器會同時傳送資料和資料的 CRC。FC 訊框的接收器會重新計算接收資料的 CRC、以確保其符合傳輸的 CRC。如果新計算的 CRC 與附加至框架的 CRC 不相符、則資料會毀損、FC 框架會遭到捨棄或拒絕。iSCSI I/O 作業包括 TCP/IP 層和以太網路層的校驗和、此外、為了提供額外的保護、也可在 SCSI 層提供選用的 CRC 保護。TCP 層或 IP 層偵測到線路上的任何位元毀損、導致封包重新傳輸。與 FC 一樣、SCSI CRC 中的錯誤也會導致作業遭到捨棄或拒絕。

### 磁碟機毀損：Checksum

Checksum 也可用來驗證儲存在磁碟機上的資料完整性。寫入磁碟機的資料區塊會以 Checksum 功能儲存、產生與原始資料相關的不可預測數字。從磁碟機讀取資料時、會重新計算總和檢查碼、並與儲存的總和檢查碼進行比較。如果不相符、則資料已毀損、必須由 RAID 層還原。

### 資料毀損：寫入遺失

最難偵測的毀損類型之一是遺失或錯誤寫入。確認寫入後、必須將其寫入正確位置的媒體。就地資料毀損使用儲存在資料中的簡單檢查碼、相當容易偵測。但是、如果寫入資料只是遺失、則先前版本的資料可能仍然存在、而且總和檢查碼是正確的。如果寫入放置在錯誤的實體位置、則相關的 Checksum 將再次對儲存的資料有效、即使寫入已銷毀其他資料。

這項挑戰的解決方案如下：

- 寫入作業必須包含中繼資料、以指出寫入的預期位置。
- 寫入作業必須包含某種版本識別碼。

ONTAP 寫入區塊時、會包含區塊所屬的資料。如果後續讀取識別出某個區塊、但中繼資料指出該區塊在位置 456 找到時屬於位置 123、則表示該寫入已放錯位置。

更難偵測完全遺失的寫入。這項說明非常複雜、但基本上 ONTAP 是以寫入作業導致磁碟機上兩個不同位置的更新方式來儲存中繼資料。如果寫入遺失、後續的資料讀取和相關中繼資料會顯示兩個不同的版本識別。這表示磁碟機未完成寫入。

遺失或放錯位置的寫入毀損極少發生、但隨著磁碟機持續成長、資料集也逐漸擴充至 EB 規模、風險也會增加。任何支援資料庫工作負載的儲存系統都應包含遺失寫入偵測。

### 磁碟機故障：RAID、RAID DP 和 RAID-TEC

如果發現磁碟機上的資料區塊毀損、或整個磁碟機故障且完全無法使用、則必須重新建立資料。這是在 ONTAP 中使用同位元磁碟機來完成的。資料會在多個資料磁碟機之間進行等量分割、然後產生同位元檢查資料。這會與原始資料分開儲存。

ONTAP 最初使用 RAID 4、每組資料磁碟機使用單一同位元檢查磁碟機。結果是群組中的任何一個磁碟機都可



能發生故障、而不會導致資料遺失。如果同位元磁碟機故障、則沒有資料受損、也可以建構新的同位元磁碟機。如果單一資料磁碟機故障、其餘磁碟機可與同位元磁碟機搭配使用、以重新產生遺失的資料。

當磁碟機很小時、同時發生兩個磁碟機故障的統計機率可忽略不計。隨著磁碟機容量的增加、磁碟機故障後重建資料所需的時間也隨之增加。這增加了第二個磁碟機故障會導致資料遺失的時間。此外、重建程序也會在正常運作的磁碟機上建立許多額外的 I/O。隨著磁碟機老化、導致第二個磁碟機故障的額外負載風險也會增加。最後、即使持續使用 RAID 4、資料遺失的風險並未增加、資料遺失的後果也會更加嚴重。在 RAID 群組故障時遺失的資料越多、還原資料所需的時間就越長、業務中斷也就越長。

這些問題導致 NetApp 開發 NetApp RAID DP 技術、這是 RAID 6 的變體。此解決方案包含兩個同位元檢查磁碟機、表示 RAID 群組中的任何兩個磁碟機都可能發生故障、而不會造成資料遺失。磁碟機的大小持續成長、最終導致 NetApp 開發 NetApp RAID-TEC 技術、引進第三個同位元磁碟機。

某些歷史資料庫最佳實務做法建議使用 RAID-10、也稱為等量鏡射。因為有多個雙磁碟故障案例、因此資料保護功能比 RAID DP 更少、而在 RAID DP 中則沒有。

由於效能考量、有些歷史資料庫最佳實務做法表示 RAID-10 較 RAID-4/5/6 選項更為偏好。這些建議有時是指 RAID 罰款。雖然這些建議一般都是正確的、但不適用於在 ONTAP 中實作 RAID。效能考量與同位元重生有關。在傳統的 RAID 實作中、處理資料庫執行的例行隨機寫入作業需要多個磁碟讀取才能重新產生同位元資料並完成寫入。其懲罰定義為執行寫入作業所需的額外讀取 IOPS。

ONTAP 不會發生 RAID 損失、因為寫入會分段在記憶體中產生同位元檢查、然後以單一 RAID 等量磁碟寫入磁碟。完成寫入作業不需要讀取。

總而言之、相較於 RAID 10、RAID DP 和 RAID-TEC 可提供更多可用容量、更好的磁碟機故障防護、而且不會犧牲效能。

### 硬體故障保護：**NVRAM**

任何服務資料庫工作負載的儲存陣列、都必須儘快服務寫入作業。此外、寫入作業也必須受到保護、避免因電源故障等非預期事件而遺失。這表示任何寫入作業都必須安全地儲存在至少兩個位置。

AFF 和 FAS 系統仰賴 NVRAM 來滿足這些需求。寫入程序的運作方式如下：

1. 傳入寫入資料儲存在 RAM 中。
2. 必須對磁碟上的資料所做的變更、會同時記入本機節點和合作夥伴節點上的 NVRAM。NVRAM 不是寫入快取、而是類似資料庫重做記錄的日誌。在正常情況下、系統不會讀取。它僅用於恢復、例如在 I/O 處理期間發生電源故障後。
3. 然後寫入會被確認給主機。

此階段的寫入程序從應用程式的角度來看已完成、而且資料會受到保護、不會遺失、因為資料會儲存在兩個不同的位置。最後、變更會寫入磁碟、但此程序會從應用程式的觀點超出頻外、因為它會在寫入確認之後發生、因此不會影響延遲。此程序再次類似於資料庫記錄。對資料庫的變更會盡快記錄在重做記錄檔中、然後將變更確認為已認可。資料檔案的更新會在稍後進行、不會直接影響處理速度。

如果控制器發生故障、合作夥伴控制器會取得所需磁碟的所有權、並重新執行 NVRAM 中記錄的資料、以恢復發生故障時正在執行的任何 I/O 作業。

### 硬體故障保護：**NVFAIL**

如前所述、寫入必須先登入本機 NVRAM 及至少一個其他控制器上的 NVRAM、才會被確認。此方法可確保硬體故障或停電不會導致機內 I/O 遺失如果本機 NVRAM 故障或連線至 HA 合作夥伴失敗、則無法再鏡射此傳輸中

的資料。

如果本機 NVRAM 回報錯誤、節點會關機。此關機會導致容錯移轉至 HA 合作夥伴控制器。由於發生故障的控制器尚未確認寫入作業、因此不會遺失任何資料。

除非強制容錯移轉、否則 ONTAP 不允許在資料不同步時進行容錯移轉。以這種方式強制變更條件、即表示資料可能會留在原始控制器中、而且資料遺失是可以接受的。

如果強制容錯移轉、則資料庫特別容易受損、因為資料庫會在磁碟上保留大量的內部資料快取。如果發生強制容錯移轉、先前確認的變更將會有效捨棄。儲存陣列的內容會有效地及時向後跳轉、而且資料庫快取的狀態不再反映磁碟上資料的狀態。

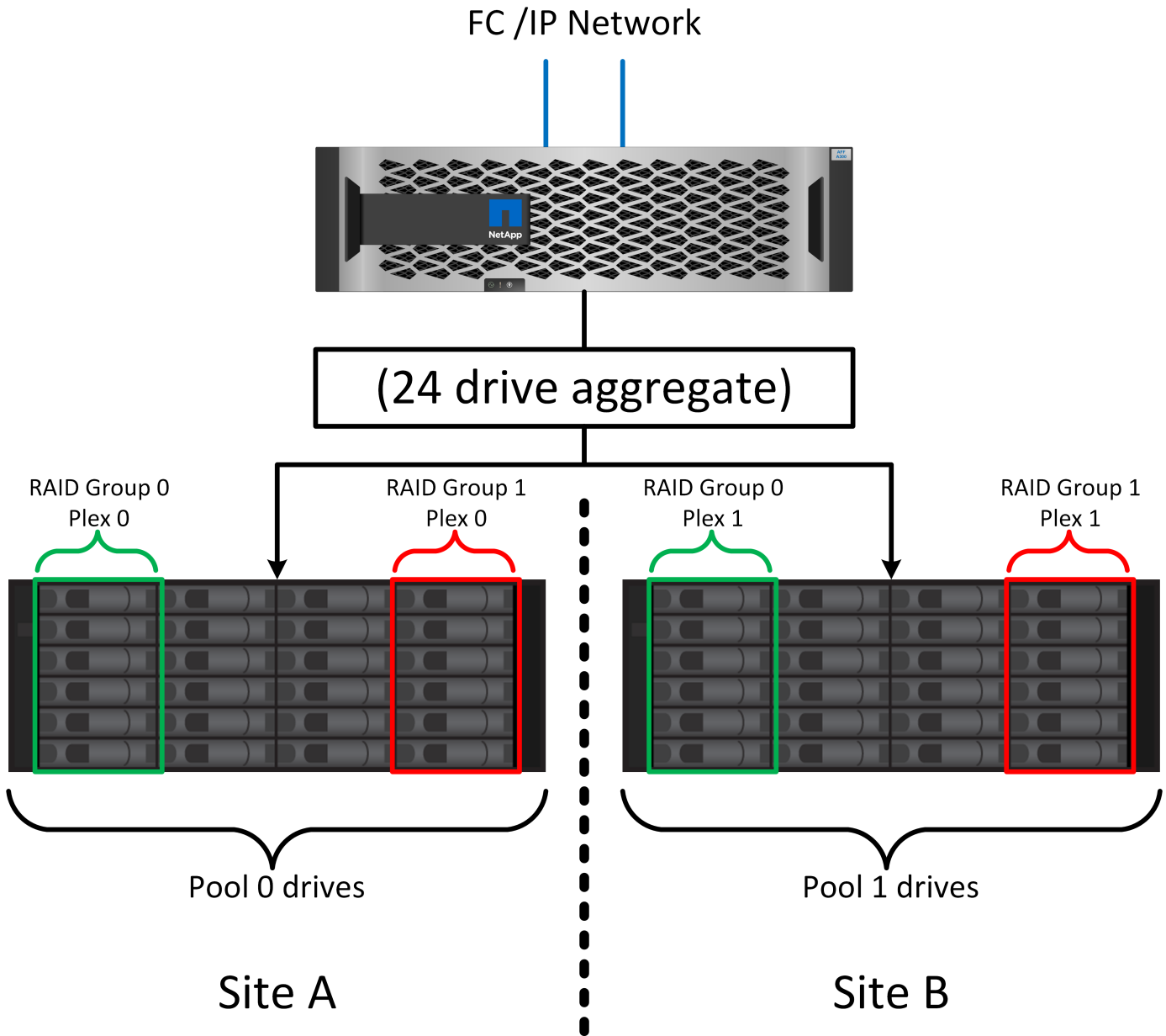
為了保護資料不受這種情況影響、ONTAP 允許設定磁碟區、以針對 NVRAM 故障提供特殊保護。觸發時、此保護機制會導致磁碟區進入稱為 NVFAIL 的狀態。此狀態會導致 I/O 錯誤、導致應用程式關機、使其不使用過時的資料。資料不應遺失、因為儲存陣列上應該存在任何已確認的寫入資料。

通常的後續步驟是讓系統管理員在手動將 LUN 和磁碟區重新上線之前、先完全關閉主機。雖然這些步驟可能涉及一些工作、但這種方法是確保資料完整性的最安全方法。並非所有資料都需要這項保護、因此 NVFAIL 行為可依每個磁碟區設定。

站台和機櫃故障保護：**SyncMirror** 和叢

SyncMirror 是一項鏡射技術、可增強但不取代 RAID DP 或 RAID-TEC。它會鏡射兩個不同 RAID 群組的內容。邏輯組態如下：

- 磁碟機會根據位置設定成兩個集區。一個集區由站台 A 上的所有磁碟機組成、第二個集區由站台 B 上的所有磁碟機組成
- 接著會根據鏡射的 RAID 群組集建立通用儲存池（稱為 Aggregate）。從每個站台擷取的磁碟機數量相等。例如、20 個磁碟機的 SyncMirror Aggregate 將由站台 A 的 10 個磁碟機和站台 B 的 10 個磁碟機組成
- 指定站台上的每組磁碟機都會自動設定為一個或多個完全備援的 RAID-DP 或 RAID-TEC 群組、而不受鏡像的使用影響。即使在站台遺失之後、也能提供持續的資料保護。



上圖說明 SyncMirror 組態範例。在控制器上建立了 24 個磁碟機的集合體、其中 12 個磁碟機來自於站台 A 上配置的機櫃、12 個磁碟機來自於站台 B 上配置的機櫃磁碟機分為兩個鏡射 RAID 群組。RAID 群組 0 包括站台 A 上的 6 磁碟機叢、鏡射到站台 B 上的 6 磁碟機叢同樣地、RAID Group 1 也包含站台 A 的 6 磁碟叢、鏡射到站台 B 的 6 磁碟叢

SyncMirror 通常用於提供 MetroCluster 系統的遠端鏡射、每個站台都有一份資料複本。有時候、它是用來在單一系統中提供額外的備援層級。特別是提供機架層級的備援。磁碟機櫃已包含雙電源供應器和控制器、整體上比金屬板稍多、但在某些情況下、可能需要額外的保護。例如、有一位 NetApp 客戶部署 SyncMirror、用於汽車測試期間使用的行動即時分析平台。系統由獨立的 UPS 系統提供獨立的電源供應器、分為兩個實體機架。

=checksum

對於習慣使用 Oracle RMAN 串流備份的 DBA 來說、檢查和主題特別重要、因為它們會移轉至快照型備份。RMAN 的一項功能是在備份作業期間執行完整性檢查。雖然這項功能有一些價值、但其主要優點是資料庫並未用於現代化的儲存陣列。當實體磁碟機用於 Oracle 資料庫時、幾乎可以確定磁碟機老化時最終會發生毀損、這是由真正儲存陣列中的陣列型校驗和所解決的問題。

使用真正的儲存陣列、資料完整性可在多個層級上使用校驗和加以保護。如果 IP 型網路中的資料毀損、傳輸控制傳輸協定 (TCP) 層會拒絕封包資料並要求重新傳輸。FC 傳輸協定包括校驗和、封裝的 SCSI 資料也一樣。在陣列上之後、ONTAP 就有 RAID 和 Checksum 保護。可能會發生毀損、但如同大多數企業陣列一樣、系統會偵測並修正毀損。一般而言、整個磁碟機都會故障、導致 RAID 重建、資料庫完整性也不會受到影響。ONTAP 偵測到 Checksum 錯誤的頻率較低、表示磁碟機上的資料已損壞。然後磁碟機故障、RAID 重建就會開始。資料完整性再次不受影響。

Oracle 資料檔案和重做記錄架構也設計成即使在極端情況下、也能提供最高程度的資料完整性。在最基本的層級、Oracle 區塊幾乎包含每個 I/O 的 Checksum 和基本邏輯檢查如果 Oracle 尚未當機或將資料表空間離線、則資料會保持不變。資料完整性檢查的程度可調整、Oracle 也可設定為確認寫入。因此、幾乎所有的當機和故障情況都可以恢復、而且在極罕見的不可恢復情況下、系統會立即偵測到毀損。

大多數使用 Oracle 資料庫的 NetApp 客戶在移轉至快照型備份後、都會停止使用 RMAN 和其他備份產品。仍有一些選項可讓 RMAN 使用 SnapCenter 執行區塊層級的還原。然而、日常使用的 RMAN、NetBackup 及其他產品只會偶爾用於建立每月或每季的歸檔複本。

有些客戶選擇執行 dbv 定期對現有資料庫執行完整性檢查。NetApp 不鼓勵這種做法、因為它會產生不必要的 I/O 負載。如上所述、如果資料庫先前沒有遇到問題、就有可能發生 dbv 偵測問題接近零、此公用程式會在網路和儲存系統上產生非常高的連續 I/O 負載。除非有理由相信存在毀損、例如暴露於已知的 Oracle 錯誤、否則沒有理由執行 dbv。

## 備份與還原基礎知識

### Oracle 資料庫和快照型備份

NetApp Snapshot 技術是 ONTAP 上 Oracle 資料庫資料保護的基礎。

關鍵值如下：

- \* 簡易性。\* 快照是特定時間點資料容器內容的唯讀複本。
- \* 效率。\* 快照在建立時不需要任何空間。只有在資料變更時才會使用空間。
- \* 管理能力。\* 由於快照是儲存作業系統的原生部分、因此以快照為基礎的備份策略很容易設定和管理。如果儲存系統已開機、就可以開始建立備份。
- \* 擴充性。\* 最多可保留 1024 個檔案和 LUN 的單一容器備份。對於複雜的資料集、可透過一組一致的快照來保護多個資料容器。
- 無論磁碟區包含 1024 個快照或無、效能都不會受到影響。

雖然許多儲存廠商都提供快照技術、但 ONTAP 中的 Snapshot 技術是獨一無二的、可為企業應用程式和資料庫環境帶來顯著效益：

- Snapshot 複本是基礎 Write -Anywhere File Layout (WAFL) 的一部分。它們不是附加技術或外部技術。這簡化了管理、因為儲存系統是備份系統。
- 快照複本不會影響效能、但某些邊緣情況除外、例如當基礎儲存系統填滿的快照中儲存了大量資料時。
- 「一致性群組」一詞通常是指一組儲存物件、這些物件是以一致的資料集合來管理。特定 ONTAP 磁碟區的快照構成一致性群組備份。

ONTAP 快照的擴充能力也優於競爭技術。客戶可以儲存 5、50 或 500 個快照、而不會影響效能。磁碟區目前允許的最大快照數為 1024。如果需要額外的快照保留、則有選項可將快照串聯至其他磁碟區。

因此、保護託管在 ONTAP 上的資料集非常簡單且具有高度擴充性。備份不需要移動資料、因此備份策略可以根據業務需求量身打造、而非限制網路傳輸率、大量磁帶機或磁碟接移區域。

快照是否為備份？

將快照當作資料保護策略使用的常見問題之一、就是「真實」資料和快照資料位於同一個磁碟機上。遺失這些磁碟機將會導致主要資料和備份遺失。

這是一項有效的考量。本機快照是用於日常備份與還原需求、在這方面、快照是備份。在 NetApp 環境中、將近 99% 的還原案例都仰賴快照來滿足最嚴苛的 RTO 需求。

然而、本機快照不應是唯一的備份策略、因此 NetApp 提供 SnapMirror 和 SnapVault 複寫等技術、可快速有效地將快照複寫至一組不同的磁碟機。在架構正確的解決方案中、快照加上快照複寫功能可將磁帶的使用降至最低、甚至是每季歸檔、或完全消除。

快照型備份

使用 ONTAP Snapshot 複本保護資料有許多選項、快照是許多其他 ONTAP 功能的基礎、包括複寫、災難恢復和複製。快照技術的完整說明不在本文件的範圍內、但以下各節提供一般概觀。

建立資料集快照的主要方法有兩種：

- 損毀一致的備份
- 應用程式一致的備份

資料集的損毀一致備份是指在單一時間點擷取整個資料集結構。如果資料集儲存在單一 NetApp FlexVol Volume 中、則程序很簡單；您可以隨時建立 Snapshot。如果資料集橫跨多個磁碟區、則必須建立一致性群組（CG）快照。建立 CG 快照有多種選項、包括 NetApp SnapCenter 軟體、原生 ONTAP 一致性群組功能、以及使用者維護的指令碼。

當備份點還原足夠時、主要會使用損毀一致的備份。當需要更精細的恢復時、通常需要應用程式一致的備份。

「應用程式一致性」一詞通常是錯誤的。例如、將 Oracle 資料庫置於備份模式稱為應用程式一致的備份、但資料並未以任何方式保持一致或停止。資料會在整個備份過程中持續變更。相反地、大部分的 MySQL 和 Microsoft SQL Server 備份確實會在執行備份之前先將資料關閉。VMware 可能會或可能不會使某些檔案一致。

一致性群組

術語「一致性群組」是指儲存陣列將多個儲存資源視為單一映像來管理的能力。例如、資料庫可能包含 10 個 LUN。陣列必須能夠以一致的方式備份、還原及複寫這 10 個 LUN。如果 LUN 的映像備份時不一致、則無法還原。複寫這 10 個 LUN 需要所有複本彼此完全同步。

討論 ONTAP 時、「一致性群組」一詞並不常使用、因為一致性一向是 ONTAP 內 Volume 和 Aggregate 架構的基本功能。許多其他儲存陣列會將 LUN 或檔案系統視為個別單元進行管理。接著可選擇性地將它們設定為「一致性群組」、以保護資料、但這是組態中的額外步驟。

ONTAP 一向能夠擷取一致的本機和複寫資料映像。雖然 ONTAP 系統上的各種磁碟區通常並未正式描述為一致性群組、但這就是它們的名稱。該磁碟區的快照是一致性群組映像、該快照的還原是一致性群組還原、SnapMirror 和 SnapVault 都提供一致性群組複寫。

一致性群組快照

一致性群組快照（CG 快照）是基本 ONTAP Snapshot 技術的延伸。標準快照作業可在單一磁碟區內建立所有

資料的一致映像、但有時必須在多個磁碟區甚至跨多個儲存系統建立一致的快照集。結果是一組快照、其使用方式與只有一個個別磁碟區的快照相同。它們可用於本機資料還原、複寫以進行災難恢復、或複製為單一一致的單元。

CG 快照的最大已知用途是資料庫環境、其大小約為 1PB、跨越 12 個控制器。在此系統上建立的 CG 快照已用於備份、恢復和複製。

大多數情況下、當資料集跨越磁碟區且必須保留寫入順序時、所選管理軟體會自動使用 CG 快照。在此情況下、無需瞭解 CG 快照的技術詳細資料。然而、在某些情況下、複雜的資料保護需求需要對資料保護和複寫程序進行詳細控制。自動化工作流程或使用自訂指令碼來呼叫 CG 快照 API 是其中的一些選項。若要瞭解最佳選項和 CG-snapshot 的角色、需要更詳細的技術說明。

建立一組 CG 快照是兩個步驟：

1. 在所有目標磁碟區上建立寫入屏障。
2. 在圍籬狀態下建立這些磁碟區的快照。

寫入隔離是連續建立的。這表示當隔離程序在多個磁碟區之間設定時、寫入 I/O 會依序凍結在第一個磁碟區上、因為它會繼續提交到稍後出現的磁碟區。這可能一開始就違反了保留寫入順序的要求、但這僅適用於非同步在主機上發出的 I/O、不需仰賴任何其他寫入。

例如、資料庫可能會發出許多非同步資料檔案更新、並允許作業系統重新排序 I/O、並根據其本身的排程器組態完成這些更新。這類 I/O 的順序無法保證、因為應用程式和作業系統已釋出保留寫入順序的要求。

以計數器為例、大部分的資料庫記錄活動都是同步的。在確認 I/O 之前、資料庫不會繼續進行記錄寫入、而且必須保留這些寫入的順序。如果記錄 I/O 到達圍籬式磁碟區、則不會予以確認、應用程式會在進一步寫入時加以封鎖。同樣地、檔案系統中繼資料 I/O 通常是同步的。例如、檔案刪除作業不得遺失。如果具有 xfs 檔案系統的作業系統刪除了檔案、而更新 xfs 檔案系統中繼資料的 I/O 則會移除位於圍籬磁碟區上的該檔案參照、則檔案系統活動就會暫停。這可確保 CG 快照作業期間檔案系統的完整性。

在目標磁碟區之間設定寫入屏障之後、就可以開始建立快照。由於磁碟區的狀態會從相關寫入點凍結、因此不需要同時精確建立快照。為了防範建立 CG 快照的應用程式中的瑕疵、初始寫入屏障包含可設定的逾時時間、ONTAP 會在定義的秒數後自動釋放隔離功能並繼續寫入處理。如果所有快照都是在逾時期間發生之前建立的、則產生的一組快照是有效的一致性群組。

## 相關寫入順序

從技術觀點來看、一致性群組的關鍵在於保留寫入順序、特別是根據寫入順序。例如、寫入 10 個 LUN 的資料庫會同時寫入所有 LUN。許多寫入都是以非同步方式發出、這表示完成的順序不重要、實際完成的順序會因作業系統和網路行為而異。

某些寫入作業必須存在於磁碟上、資料庫才能繼續進行其他寫入作業。這些關鍵寫入作業稱為「相關寫入」。後續寫入 I/O 則取決於磁碟上是否有這些寫入資料。這 10 個 LUN 的任何快照、恢復或複寫都必須確保相關寫入順序受到保證。檔案系統更新是寫入順序相關寫入的另一個範例。必須保留檔案系統變更的順序、否則整個檔案系統可能會毀損。

## 策略

以快照為基礎的備份主要有兩種方法：

- 損毀一致的備份
- 快照保護的熱備份

資料庫的損毀一致備份是指在單一時間點擷取整個資料庫結構、包括資料檔案、重做記錄和控制檔。如果資料庫儲存在單一 NetApp FlexVol Volume 中、則程序很簡單；您可以隨時建立 Snapshot。如果資料庫橫跨磁碟區、則必須建立一致性群組（CG）快照。建立 CG 快照有多種選項、包括 NetApp SnapCenter 軟體、原生 ONTAP 一致性群組功能、以及使用者維護的指令碼。

當備份點還原足夠時、主要會使用損毀一致的 Snapshot 備份。在某些情況下可以套用歸檔記錄檔、但如果需要更精細的時間點還原、則最好使用線上備份。

快照型線上備份的基本程序如下：

1. 將資料庫放入 backup 模式。
2. 建立所有託管資料檔案的磁碟區快照。
3. 結束 backup 模式。
4. 執行命令 `alter system archive log current` 強制記錄歸檔。
5. 為所有託管歸檔記錄的磁碟區建立快照。

此程序會產生一組快照、其中包含備份模式中的資料檔案、以及在備份模式中產生的重要歸檔記錄。這是恢復資料庫的兩項需求。控制檔等檔案也應受到保護、以方便使用、但唯一的絕對需求是保護資料檔案和歸檔記錄。

雖然不同的客戶可能有非常不同的策略、但幾乎所有這些策略最終都是以下列相同原則為基礎。

#### 快照型還原

在設計 Oracle 資料庫的 Volume 配置時、第一個決定是是否使用 Volume NetApp SnapRestore（VBSR）技術。

Volume 型 SnapRestore 可讓磁碟區立即還原至較早的時間點。由於磁碟區上的所有資料都已還原、因此 VBSR 可能不適用於所有使用案例。例如、如果整個資料庫（包括資料檔案、重做記錄和歸檔記錄）儲存在單一磁碟區上、且此磁碟區使用 VBSR 還原、則資料會遺失、因為較新的歸檔記錄和重做資料會被捨棄。

還原不需要 VSR。許多資料庫都可以使用檔案型單一檔案 SnapRestore（SFSR）來還原、或只是將檔案從快照複製回作用中的檔案系統。

當資料庫非常大或必須盡快恢復時、最好使用 VBSR、而使用 VSR 需要隔離資料檔案。在 NFS 環境中、指定資料庫的資料檔案必須儲存在專用的磁碟區中、而這些磁碟區不會受到任何其他類型的檔案污染。在 SAN 環境中、資料檔案必須儲存在專用 FlexVol 磁碟區上的專用 LUN 中。如果使用 Volume Manager（包括 Oracle 自動儲存管理 [AS]）、則磁碟群組也必須專用於資料檔案。

以這種方式隔離資料檔案、可讓檔案還原至較早的狀態、而不會損壞其他檔案系統。

#### Snapshot保留

對於 SAN 環境中具有 Oracle 資料的每個 Volume `percent-snapshot-space` 應設為零、因為在 LUN 環境中保留快照空間並不實用。如果百分比保留設為 100、則具有 LUN 的磁碟區快照需要在磁碟區中有足夠的可用空間、但不包括快照保留空間、以吸收所有資料 100% 的營業額。如果將百分比保留設為較低的值、則需要相對較小的可用空間、但它一律會排除快照保留。這表示 LUN 環境中的快照保留空間會被浪費。

在 NFS 環境中、有兩個選項：

- 設定 `percent-snapshot-space` 根據預期的快照空間使用量。

- 設定 `percent-snapshot-space` 以歸零並統整管理作用中和快照空間使用量。

使用第一個選項、`percent-snapshot-space` 設為非零值、通常約 20%。然後、使用者就會隱藏此空間。不過、此值並不會限制使用率。如果具有 20% 保留的資料庫擁有 30% 的營業額、則快照空間可能會超出 20% 保留空間的範圍、並佔用無保留空間。

將保留設定為 20% 等值的主要優點是驗證某些空間永遠可供快照使用。例如、保留 20% 的 1TB 磁碟區只允許資料庫管理員 (DBA) 儲存 800GB 的資料。此組態保證至少有 200GB 的空間可供快照使用。

何時 `percent-snapshot-space` 設為零、則使用者可以使用磁碟區中的所有空間、以提供更好的可見度。DBA 必須瞭解、如果他 / 她看到 1TB 的磁碟區運用快照、則這 1TB 的空間會在使用中資料和 Snapshot 週轉之間共享。

終端使用者之間的選項 1 和選項 2 之間沒有明確的偏好設定。

### ONTAP 和第三方快照

Oracle Doc ID 604683.1 說明第三方快照支援的需求、以及備份與還原作業的多種選項。

第三方廠商必須保證公司的快照符合下列要求：

- 快照必須與 Oracle 建議的還原與還原作業整合。
- 快照必須在快照點保持一致的資料庫損毀。
- 快照中的每個檔案都會保留寫入順序。

ONTAP 和 NetApp Oracle 管理產品符合這些要求。

### 使用 SnapRestore 快速恢復 Oracle 資料庫

NetApp SnapRestore 技術可從快照快速還原 ONTAP 中的資料。

當關鍵資料集無法使用時、關鍵業務營運就會中斷。磁帶可能會中斷、甚至是從磁碟型備份還原、在網路上傳輸速度可能會很慢。SnapRestore 可提供近乎即時的資料集還原功能、避免這些問題。即使是 PB 規模的資料庫、只要花幾分鐘的時間就能完全還原。

SnapRestore 有兩種形式：檔案 /LUN 型和磁碟區型。

- 無論是 2TB LUN 或 4KB 檔案、個別檔案或 LUN 都能在數秒內還原。
- 無論是 10GB 或 100TB 的資料、檔案或 LUN 的容器都能在數秒內還原。

「檔案或 LUN 的容器」通常指的是 FlexVol Volume。例如、您可以在單一磁碟區中擁有 10 個組成 LVM 磁碟群組的 LUN、或是一個磁碟區可以儲存 1000 位使用者的 NFS 主目錄。您可以將整個磁碟區還原為單一作業、而不需為每個個別檔案或 LUN 執行還原作業。此程序也適用於包含多個磁碟區的橫向擴充容器、例如 FlexGroup 或 ONTAP 一致性群組。

SnapRestore 之所以能如此快速有效地運作、是因為快照的本質、基本上是在特定時間點對磁碟區內容進行平行唯讀檢視。作用中區塊是可以變更的實際區塊、而快照則是建立快照時構成檔案和 LUN 之區塊狀態的唯讀檢視。

ONTAP 僅允許唯讀存取快照資料、但可透過 SnapRestore 重新啟動資料。快照會重新啟用為資料的讀寫檢視、並將資料恢復至先前的狀態。SnapRestore 可以在磁碟區或檔案層級運作。這項技術基本上相同、但行為上有



一些小差異。

## Volume SnapRestore

Volume 型 SnapRestore 會將整個資料量傳回至較早的狀態。這項作業不需要資料移動、也就是說還原程序基本上是即時的、雖然 API 或 CLI 作業可能需要幾秒鐘的時間才能處理。還原 1GB 的資料並不比還原 1PB 的資料複雜或耗時。這項功能是許多企業客戶移轉至 ONTAP 儲存系統的主要原因。即使是最大的資料集、也能以秒為單位提供 RTO。

Volume 型 SnapRestore 的缺點之一、是因為磁碟區內的變更會隨時間累積。因此、每個快照和作用中檔案資料都取決於到該點之前的變更。將磁碟區還原為較早的狀態、表示捨棄所有後續對資料所做的變更。然而、不太明顯的是、這包括後續建立的快照。這並不總是理想的。

例如、資料保留 SLA 可能會指定每晚備份 30 天。將資料集還原至五天前以 Volume SnapRestore 建立的快照、將會捨棄前五天建立的所有快照、違反 SLA。

有許多選項可解決此限制：

1. 資料可從先前的快照複製、而非執行整個 Volume 的 SnapRestore。此方法最適合較小的資料集。
2. 您可以複製快照、而非還原快照。此方法的限制在於來源快照是複本的相依性。因此、除非也刪除複本、或將其分割成一個不同的 Volume、否則無法將其刪除。
3. 使用檔案型 SnapRestore。

## File SnapRestore

檔案型 SnapRestore 是更精細的快照型還原程序。個別檔案或 LUN 的狀態會還原、而非還原整個磁碟區的狀態。不需要刪除快照、此作業也不需要對先前的快照建立任何相依性。檔案或 LUN 會立即在作用中磁碟區中可用。

SnapRestore 還原檔案或 LUN 時不需要移動資料。不過、需要進行一些內部中繼資料更新、以反映檔案或 LUN 中的基礎區塊現在同時存在於快照和作用中磁碟區中。不應影響效能、但此程序會封鎖快照的建立、直到快照完成為止。根據還原的檔案總大小、處理速度約為 5Gbps (18TB/小時)。

## Oracle 資料庫線上備份

在備份模式中保護和恢復 Oracle 資料庫需要兩組資料。請注意、這不是唯一的 Oracle 備份選項、而是最常見的選項。

- 備份模式中資料檔案的快照
- 資料檔案處於備份模式時所建立的歸檔記錄

如果需要完整恢復 (包括所有已提交的交易)、則需要第三個項目：

- 一組目前的重做記錄

有許多方法可以推動線上備份的還原。許多客戶使用 ONTAP CLI 還原快照、然後使用 Oracle RMAN 或 sqlplus 來完成還原。這在大型正式作業環境中尤其常見、因為資料庫還原的可能性和頻率極低、而且任何還原程序都是由熟練的 DBA 來處理。為了實現完整的自動化、NetApp SnapCenter 等解決方案包含 Oracle 外掛程式、其中包含命令列和圖形介面。

有些大型客戶已在主機上設定基本指令碼、以便在特定時間將資料庫置於備份模式、以準備排程的快照、藉此採

取更簡單的方法。例如、排程命令 `alter database begin backup 23 時 58 分`、`alter database end backup` 於 00 : 02、然後於午夜直接在儲存系統上排程快照。如此一來、就能實現簡單易用、擴充性極高的備份策略、無需外部軟體或授權。

#### 資料配置

最簡單的配置是將資料檔案隔離到一個或多個專用磁碟區。它們必須不受任何其他檔案類型的污染。這是為了確保資料檔案磁碟區可以透過 SnapRestore 作業快速還原、而不會破壞重要的重做記錄檔、控制檔或歸檔記錄。

SAN 對專用磁碟區內的資料檔案隔離有類似的需求。在 Microsoft Windows 等作業系統中、單一磁碟區可能包含多個資料檔案 LUN、每個 LUN 都有 NTFS 檔案系統。在其他作業系統中、通常會有邏輯 Volume Manager。例如、使用 Oracle ASM 時、最簡單的選項是將 ASM 磁碟群組的 LUN 限制在單一磁碟區、以便作為一個單元進行備份和還原。如果基於效能或容量管理的理由而需要額外的磁碟區、則在新磁碟區上建立額外的磁碟群組、將可簡化管理。

如果遵循這些準則、則可直接在儲存系統上排程快照、而無需執行一致性群組快照。原因是 Oracle 備份不需要同時備份資料檔案。線上備份程序旨在讓資料檔案在數小時內緩慢串流至磁帶、因此能夠繼續更新。

在使用分佈於不同磁碟區的 ASM 磁碟群組等情況下、會產生複雜性。在這種情況下、必須執行 CG 快照、以確保 ASM 中繼資料在所有組成磁碟區之間一致。

- 注意：\* 驗證 ASM `spfile` 和 `passwd` 檔案不在主控資料檔案的磁碟群組中。這會影響選擇性還原資料檔案和僅還原資料檔案的能力。

#### 本機恢復程序— NFS

此程序可以手動或透過 SnapCenter 等應用程式來驅動。基本程序如下：

1. 關閉資料庫。
2. 在所需還原點之前、立即將資料檔案磁碟區復原至快照。
3. 將歸檔記錄重播至所需的點。
4. 如果需要完整還原、請重新播放目前的重做記錄。

此程序假設所需的歸檔記錄檔仍存在於作用中的檔案系統中。如果沒有、則必須還原歸檔記錄、或將 RMAN/sqlplus 導向快照目錄中的資料。

此外、對於較小的資料庫、終端使用者可以直接從中復原資料檔案 `.snapshot` 目錄、無需自動化工具或儲存管理員協助執行 `snaprestore` 命令。

#### 本機恢復程序— SAN

此程序可以手動或透過 SnapCenter 等應用程式來驅動。基本程序如下：

1. 關閉資料庫。
2. 將託管資料檔案的磁碟群組置於系統中。此程序會因所選的邏輯磁碟區管理程式而異。使用 ASM 時、此程序需要卸除磁碟群組。在 Linux 中、必須卸除檔案系統、且必須停用邏輯磁碟區和磁碟區群組。目標是停止要還原之目標 Volume 群組上的所有更新。
3. 在所需還原點之前、立即將資料檔案磁碟群組還原至快照。
4. 重新啟動新還原的磁碟群組。

5. 將歸檔記錄重播至所需的點。
6. 如果需要完整還原、請重新播放所有重做記錄。

此程序假設所需的歸檔記錄檔仍存在於作用中的檔案系統中。如果沒有、則必須將歸檔記錄 LUN 離線並執行還原、以還原歸檔記錄。這也是將歸檔記錄分割成專用磁碟區的範例。如果歸檔記錄與重做記錄共用一個磁碟區群組、則必須先將重做記錄複製到其他位置、才能還原整個 LUN 集。此步驟可防止這些最終記錄的交易遺失。

### Oracle 資料庫儲存快照最佳化備份

當 Oracle 12c 發行時、快照型備份與還原變得更簡單、因為不需要將資料庫置於熱備份模式。結果是能夠直接在儲存系統上排程快照式備份、同時仍保留執行完整或時間點還原的能力。

雖然 DBA 較熟悉熱備份還原程序、但很長一段時間以來、它仍可使用資料庫處於熱備份模式時未建立的快照。恢復期間、Oracle 10g 和 11g 需要額外的步驟、才能使資料庫保持一致。使用 Oracle 12c、sqlplus 和 rman 包含額外的邏輯、可在非熱備份模式的資料檔案備份上重播歸檔記錄。

如前所述、復原快照型熱備份需要兩組資料：

- 在備份模式下建立的資料檔案快照
- 資料檔案處於熱備份模式時所產生的歸檔記錄

在還原期間、資料庫會從資料檔案讀取中繼資料、以選取所需的歸檔記錄進行還原。

儲存快照最佳化的還原需要稍微不同的資料集、才能達到相同的結果：

- 資料檔案的快照、加上一種識別快照建立時間的方法
- 從最新資料檔案檢查點的時間到快照的確切時間、都會歸檔記錄檔

在還原期間、資料庫會從資料檔案讀取中繼資料、以識別所需的最早歸檔記錄。可以執行完整或時間點恢復。執行時間點還原時、必須知道資料檔案快照的時間。指定的恢復點必須在快照建立時間之後。NetApp 建議您在快照時間中加入至少幾分鐘、以因應時鐘變化。

如需完整的詳細資料、請參閱 Oracle 各版本的 Oracle 12c 說明文件中有關「使用儲存 Snapshot 最佳化進行恢復」主題的 Oracle 文件。此外、請參閱 Oracle 文件 ID 文件 ID 604683.1、瞭解 Oracle 協力廠商快照支援。

### 資料配置

最簡單的配置是將資料檔案隔離為一個或多個專用磁碟區。它們必須不受任何其他檔案類型的污染。這是為了確保資料檔案磁碟區可以透過 SnapRestore 作業快速還原、而不會破壞重要的重做記錄檔、控制檔或歸檔記錄檔。

SAN 對專用磁碟區內的資料檔案隔離有類似的需求。在 Microsoft Windows 等作業系統中、單一磁碟區可能包含多個資料檔案 LUN、每個 LUN 都有 NTFS 檔案系統。在其他作業系統中、通常也會有邏輯 Volume Manager。例如、使用 Oracle ASM 時、最簡單的選項是將磁碟群組限制在單一磁碟區、以便作為一個單元進行備份和還原。如果基於效能或容量管理的理由而需要額外的磁碟區、則在新磁碟區上建立額外的磁碟群組、將可更輕鬆地進行管理。

如果遵循這些準則、則可直接在 ONTAP 上排程快照、而無需執行一致性群組快照。原因是快照最佳化備份不需要同時備份資料檔案。

在 ASM 磁碟群組等情況下、會發生複雜的情況、而 ASM 磁碟群組會分散在不同的磁碟區中。在這種情況下、必須執行 CG 快照、以確保 ASM 中繼資料在所有組成磁碟區之間一致。

[ 注意 ] 確認 ASM spfile 和 passwd 檔案不在主控資料檔案的磁碟群組中。這會影響選擇性還原資料檔案和僅還原資料檔案的能力。

#### 本機恢復程序— NFS

此程序可以手動或透過 SnapCenter 等應用程式來驅動。基本程序如下：

1. 關閉資料庫。
2. 在所需還原點之前、立即將資料檔案磁碟區復原至快照。
3. 將歸檔記錄重播至所需的點。

此程序假設所需的歸檔記錄檔仍存在於作用中的檔案系統中。如果沒有、則必須還原歸檔記錄、或 rman 或 sqlplus 可導向至中的資料 .snapshot 目錄。

此外、對於較小的資料庫、終端使用者可以直接從中復原資料檔案 .snapshot 無需自動化工具或儲存管理員協助執行 SnapRestore 命令的目錄。

#### 本機恢復程序— SAN

此程序可以手動或透過 SnapCenter 等應用程式來驅動。基本程序如下：

1. 關閉資料庫。
2. 將託管資料檔案的磁碟群組置於系統中。此程序會因所選的邏輯磁碟區管理程式而異。使用 ASM 時、此程序需要卸除磁碟群組。在 Linux 中、必須卸除檔案系統、並停用邏輯磁碟區和磁碟區群組。目標是停止要還原之目標 Volume 群組上的所有更新。
3. 在所需還原點之前、立即將資料檔案磁碟群組還原至快照。
4. 重新啟動新還原的磁碟群組。
5. 將歸檔記錄重播至所需的點。

此程序假設所需的歸檔記錄檔仍存在於作用中的檔案系統中。如果沒有、則必須將歸檔記錄 LUN 離線並執行還原、以還原歸檔記錄。這也是將歸檔記錄分割成專用磁碟區的範例。如果歸檔記錄與重做記錄共用磁碟區群組、則必須在還原整體 LUN 組之前、將重做記錄複製到其他位置、以免遺失最終記錄的交易。

#### 完整恢復範例

假設資料檔案已毀損或毀損、且需要完整還原。執执行程序如下：

```

[oracle@host1 ~]$ sqlplus / as sysdba
Connected to an idle instance.
SQL> startup mount;
ORACLE instance started.
Total System Global Area 1610612736 bytes
Fixed Size                2924928 bytes
Variable Size             1040191104 bytes
Database Buffers         553648128 bytes
Redo Buffers              13848576 bytes
Database mounted.
SQL> recover automatic;
Media recovery complete.
SQL> alter database open;
Database altered.
SQL>

```

#### 時間點恢復範例

整個恢復過程只需一個命令：recover automatic。

如果需要時間點恢復、則必須知道快照的時間戳記、並可識別如下：

```

Cluster01::> snapshot show -vserver vserver1 -volume NTAP_oradata -fields
create-time
vserver   volume           snapshot         create-time
-----
vserver1  NTAP_oradata    my-backup       Thu Mar 09 10:10:06 2017

```

快照建立時間列於 3 月 9 日和 10 : 10 : 06 。為了安全起見、快照時間會增加一分鐘：

```

[oracle@host1 ~]$ sqlplus / as sysdba
Connected to an idle instance.
SQL> startup mount;
ORACLE instance started.
Total System Global Area 1610612736 bytes
Fixed Size                2924928 bytes
Variable Size             1040191104 bytes
Database Buffers         553648128 bytes
Redo Buffers              13848576 bytes
Database mounted.
SQL> recover database until time '09-MAR-2017 10:44:15' snapshot time '09-
MAR-2017 10:11:00';

```

恢復作業現在已啟動。它指定的快照時間為 10 : 11 : 00、記錄時間後一分鐘、以計算可能的時鐘差異、目標恢復時間為 10 : 44。接下來、sqlplus 會要求所需的歸檔記錄檔、以達到所需的 10 : 44 恢復時間。

```
ORA-00279: change 551760 generated at 03/09/2017 05:06:07 needed for
thread 1
ORA-00289: suggestion : /orlogs_nfs/arch/1_31_930813377.dbf
ORA-00280: change 551760 for thread 1 is in sequence #31
Specify log: {<RET>=suggested | filename | AUTO | CANCEL}
ORA-00279: change 552566 generated at 03/09/2017 05:08:09 needed for
thread 1
ORA-00289: suggestion : /orlogs_nfs/arch/1_32_930813377.dbf
ORA-00280: change 552566 for thread 1 is in sequence #32
Specify log: {<RET>=suggested | filename | AUTO | CANCEL}
ORA-00279: change 553045 generated at 03/09/2017 05:10:12 needed for
thread 1
ORA-00289: suggestion : /orlogs_nfs/arch/1_33_930813377.dbf
ORA-00280: change 553045 for thread 1 is in sequence #33
Specify log: {<RET>=suggested | filename | AUTO | CANCEL}
ORA-00279: change 753229 generated at 03/09/2017 05:15:58 needed for
thread 1
ORA-00289: suggestion : /orlogs_nfs/arch/1_34_930813377.dbf
ORA-00280: change 753229 for thread 1 is in sequence #34
Specify log: {<RET>=suggested | filename | AUTO | CANCEL}
Log applied.
Media recovery complete.
SQL> alter database open resetlogs;
Database altered.
SQL>
```



使用來完成資料庫的快照還原 `recover automatic` 命令不需要特定授權、而是使用時間點還原 `snapshot time` 需要 Oracle 進階壓縮授權。

## Oracle 資料庫管理與自動化工具

ONTAP 在 Oracle 資料庫環境中的主要價值來自核心 ONTAP 技術、例如即時 Snapshot 複本、簡單的 SnapMirror 複寫、以及快速建立 FlexClone Volume。

在某些情況下、直接在 ONTAP 上簡單設定這些核心功能即可滿足需求、但更複雜的需求則需要協調層。

### SnapCenter

SnapCenter 是 NetApp 資料保護的旗艦產品。在極低的層級上、它與 SnapManager 產品在執行資料庫備份的方式上類似、但它是從頭開始打造、提供單一窗口來管理 NetApp 儲存系統上的資料保護。

SnapCenter 包括快照式備份與還原、SnapMirror 與 SnapVault 複寫等基本功能、以及大型企業大規模營運所需的其他功能。這些進階功能包括擴充的角色型存取控制 (RBAC) 功能、可與協力廠商協調化產品整合的 RESTful API、資料庫主機上 SnapCenter 外掛程式的不中斷中央管理、以及專為雲端規模環境設計的使用者介

面。

休息

ONTAP 也包含豐富的 RESTful API 集。這可讓協力廠商建立資料保護及其他管理應用程式、並與 ONTAP 進行深度整合。此外、想要建立自己的自動化工作流程和公用程式的客戶也能輕鬆使用 RESTful API。

## Oracle 災難恢復

### 使用 ONTAP 進行 Oracle 資料庫災難恢復

災難恢復是指在發生災難性事件（例如破壞儲存系統甚至整個站台的火災）之後還原資料服務。



本文件取代先前發佈的技術報告 [\\_TR-4591 : Oracle Data Protection](#) 和 [\\_TR-4592 : Oracle on MetroCluster](#)。

當然、災難恢復可以透過使用 SnapMirror 簡單複寫資料來完成、許多客戶會在每小時更新鏡射複本。

對於大多數客戶而言、DR 不只需要擁有遠端資料複本、還需要能夠快速使用該資料。NetApp 提供兩種技術來滿足這種需求：MetroCluster 和 SnapMirror 主動同步

MetroCluster 指的是硬體組態中的 ONTAP、其中包括低階同步鏡射儲存設備和許多其他功能。MetroCluster 等整合式解決方案可簡化現今複雜的橫向擴充資料庫、應用程式及虛擬化基礎架構。它以一個簡單的中央儲存陣列取代多種外部資料保護產品和策略。它也能在單一叢集式儲存系統中提供整合式備份、還原、災難恢復和高可用性（HA）。

SnapMirror 主動同步是以 SnapMirror Synchronous 為基礎。使用 MetroCluster、每個 ONTAP 控制器都負責將其磁碟機資料複寫到遠端位置。有了 SnapMirror 主動式同步、您基本上擁有兩個不同的 ONTAP 系統、可維護 LUN 資料的獨立複本、但可以合作呈現該 LUN 的單一執行個體。從主機的角度來看、這是單一 LUN 實體。

雖然 SnapMirror 主動式同步和 MetroCluster 在內部的運作方式截然不同、但對於主機而言、結果卻非常相似。主要差異在於精細度。如果您只需要選取要同步複寫的工作負載、SnapMirror 主動同步是更好的選擇。如果您需要複寫整個環境、甚至是資料中心、MetroCluster 是更好的選擇。此外、SnapMirror 主動式同步目前僅適用於 SAN、而 MetroCluster 則是多重傳輸協定、包括 SAN、NFS 和 SMB。

## MetroCluster

### MetroCluster 實體架構和 Oracle 資料庫

瞭解 Oracle 資料庫在 MetroCluster 環境中的運作方式、需要對 MetroCluster 系統的實體設計進行一些說明。



本文件取代先前發佈的技術報告 [\\_TR-4592 : Oracle on MetroCluster](#)。

MetroCluster 可在 3 種不同組態中使用

- HA 可與 IP 連線配對
- HA 可與 FC 連線配對

- 單一控制器、具備 FC 連線能力

[ 注意 ] 「連線」一詞是指用於跨站台複寫的叢集連線。它並不指主機協定。無論叢集間通訊所使用的連線類型為何、MetroCluster 組態中的所有主機端通訊協定都會如常支援。

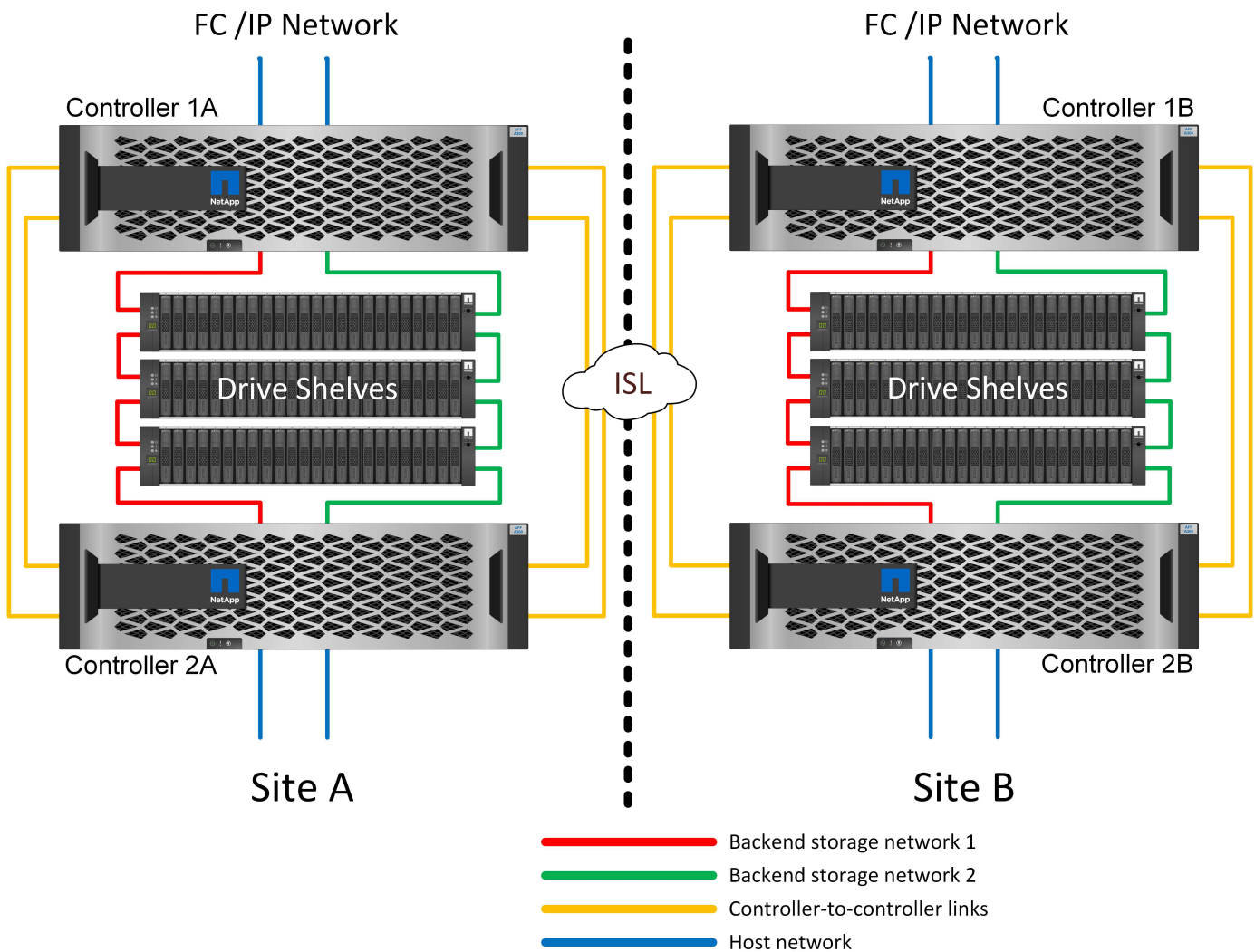
#### 知識產權MetroCluster

HA 配對 MetroCluster IP 組態每個站台使用兩或四個節點。此組態選項可增加與雙節點選項相關的複雜度和成本、但它提供重要的優點：站台內備援。簡單的控制器故障不需要透過 WAN 存取資料。透過替代本機控制器、資料存取仍保持在本機狀態。

大多數客戶都選擇 IP 連線、因為基礎架構需求較為簡單。過去、高速跨站台連線通常較容易使用深色光纖和 FC 交換器進行配置、但如今、高速、低延遲的 IP 電路更容易使用。

由於唯一的跨站台連線適用於控制器、因此架構也更簡單。在 FC SAN 附加 MetroCluster 中、控制器會直接寫入另一個站台上的磁碟機、因此需要額外的 SAN 連線、交換器和橋接器。相反地、IP 組態中的控制器會透過控制器寫入相對的磁碟機。

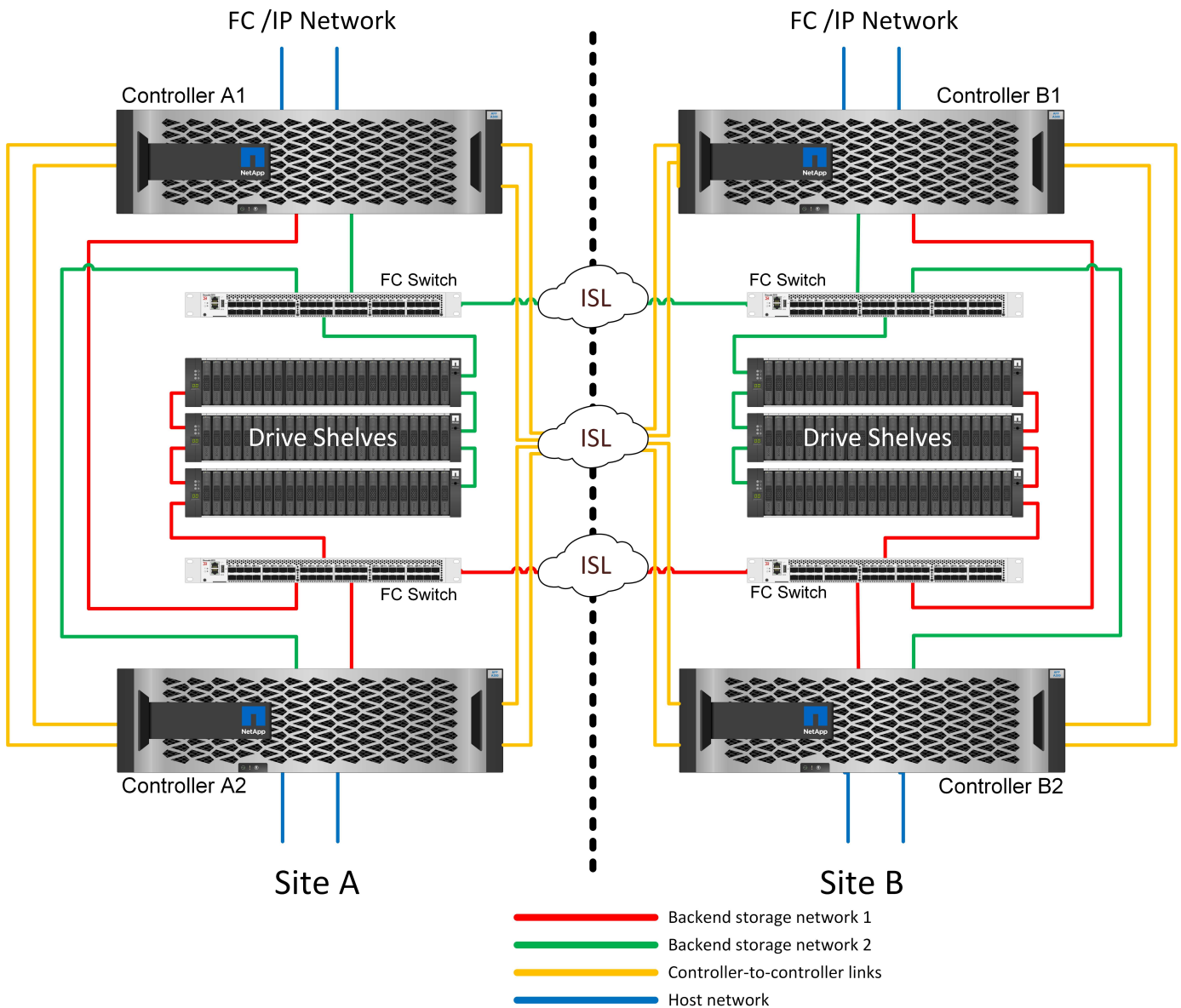
如需其他資訊、請參閱 ONTAP 正式文件和 "[SIP解決方案架構與設計MetroCluster](#)"。





## HA 配對 FC SAN 附加 MetroCluster

HA 配對 MetroCluster FC 組態每個站台使用兩個或四個節點。此組態選項可增加與雙節點選項相關的複雜度和成本、但它提供重要的優點：站台內備援。簡單的控制器故障不需要透過 WAN 存取資料。透過替代本機控制器、資料存取仍保持在本機狀態。



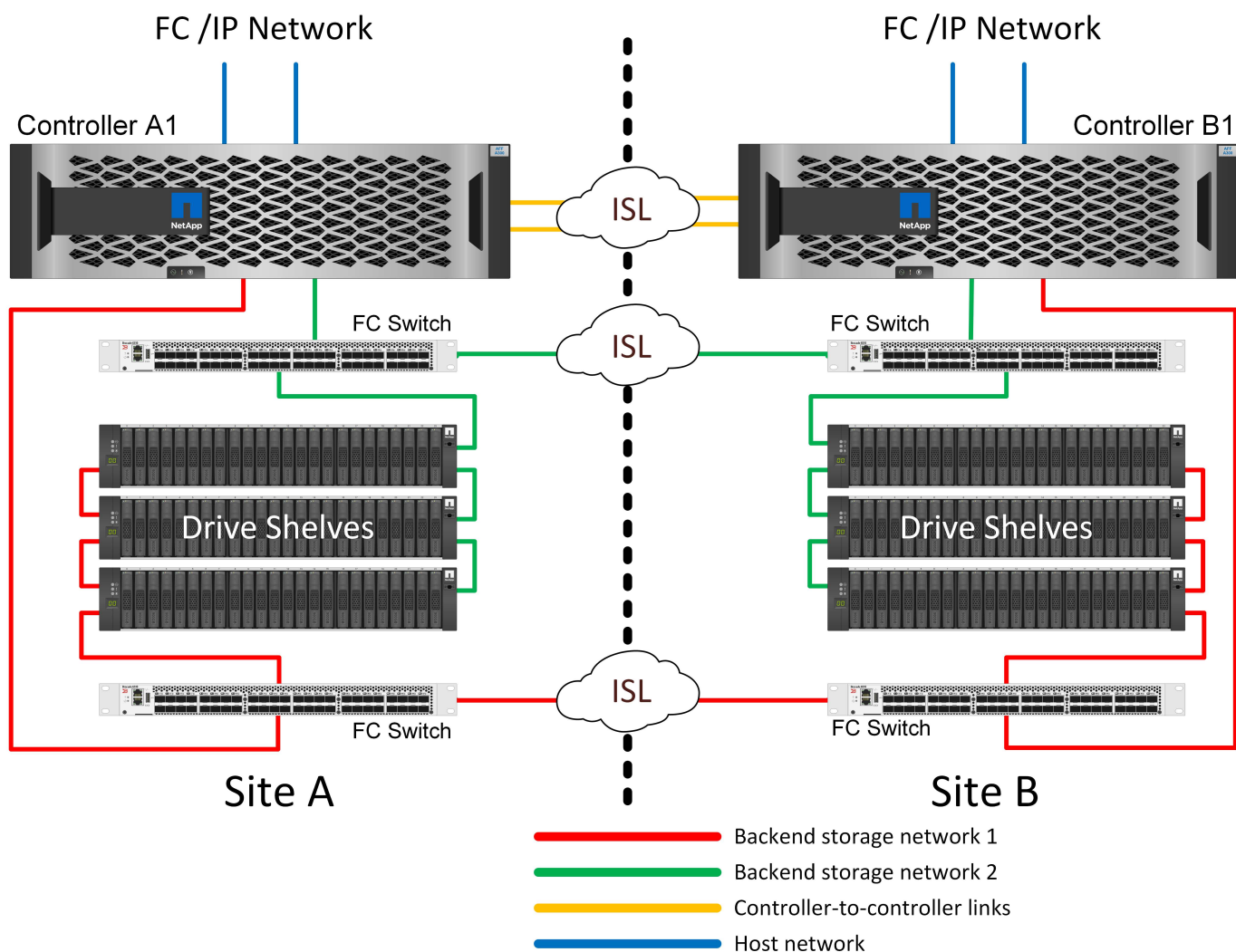
有些多站台基礎架構並非設計用於主動式作業、而是更多用於主要站台和災難恢復站台。在這種情況下、HA 配對 MetroCluster 選項通常較為理想、原因如下：

- 雖然雙節點 MetroCluster 叢集是 HA 系統、但控制器意外故障或規劃的維護作業需要資料服務必須在相反的站台上線。如果站台之間的網路連線能力不支援所需的頻寬、效能就會受到影響。唯一的選項是將各種主機作業系統和相關服務容錯移轉至替代站台。HA 配對 MetroCluster 叢集可消除此問題、因為遺失控制器會導致同一個站台內的簡單容錯移轉。
- 有些網路拓撲並非設計用於跨站台存取、而是使用不同的子網路或隔離的 FC SAN。在這種情況下、雙節點 MetroCluster 叢集不再作為 HA 系統運作、因為替代控制器無法將資料提供給位於相反站台的伺服器。HA 配對 MetroCluster 選項是提供完整備援的必要條件。
- 如果將雙站台基礎架構視為單一的高可用度基礎架構、則雙節點 MetroCluster 組態很適合。不過、如果系統

在站台故障後必須長時間運作、則最好使用 HA 配對、因為它會繼續在單一站台內提供 HA。

### 雙節點 FC SAN 附加 MetroCluster

雙節點 MetroCluster 組態每個站台僅使用一個節點。此設計比 HA 配對選項簡單、因為要設定和維護的元件較少。此外、它也降低了佈線和 FC 交換方面的基礎架構需求。最後、它能降低成本。



這項設計的明顯影響是、控制器在單一站台上故障、表示資料可從另一個站台取得。這種限制不一定是個問題。許多企業都有多站台資料中心作業、並有延伸、高速、低延遲的網路、基本上是一個基礎架構。在這些情況下、MetroCluster 的雙節點版本是慣用的組態。多家服務供應商目前以 PB 規模使用雙節點系統。

### MetroCluster 恢復功能

MetroCluster 解決方案沒有單點故障：

- 每個控制器都有兩條通往本機站台磁碟櫃的路徑。
- 每個控制器都有兩條通往遠端站台磁碟機櫃的路徑。
- 每個控制器都有兩條通往另一個站台上控制器的路徑。
- 在 HA 配對組態中、每個控制器都有兩條路徑通往本機合作夥伴。

總而言之、您可以移除組態中的任何一個元件、而不會影響 MetroCluster 提供資料的能力。這兩個選項之間恢復能力的唯一差異是 HA 配對版本在站台故障後仍是整個 HA 儲存系統。

## MetroCluster 邏輯架構和 Oracle 資料庫

瞭解 Oracle 資料庫在 MetroCluster 環境中的運作方式需要對 MetroCluster 系統的邏輯功能進行一些說明。

### 站台故障保護：NVRAM 和 MetroCluster

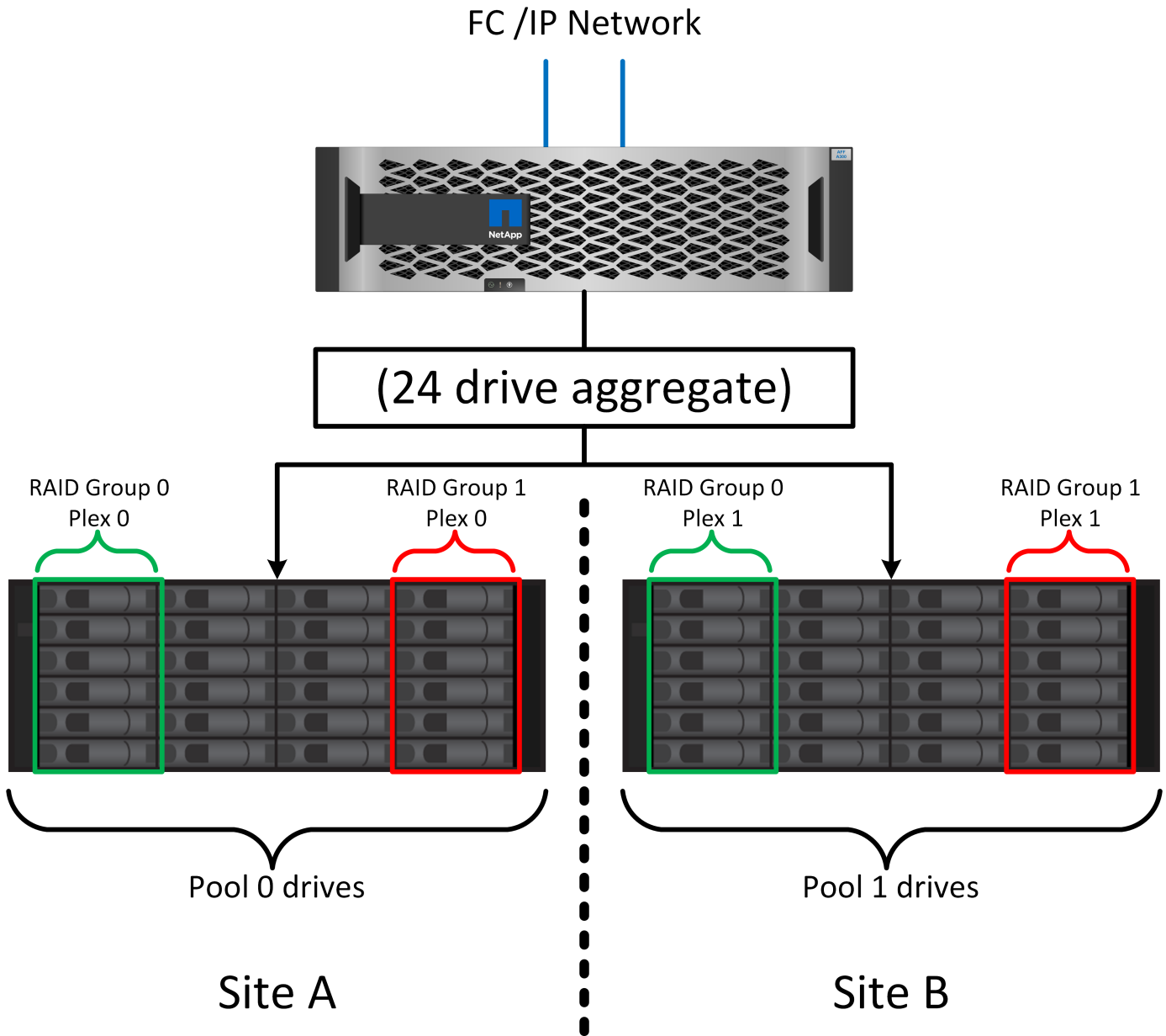
MetroCluster 以下列方式擴充 NVRAM 資料保護：

- 在雙節點組態中、NVRAM 資料會使用交換器間連結（ISL）複寫到遠端合作夥伴。
- 在 HA 配對組態中、NVRAM 資料會同時複寫到本機合作夥伴和遠端合作夥伴。
- 寫入內容必須複寫到所有合作夥伴、才能予以確認。此架構可將 NVRAM 資料複寫至遠端合作夥伴、保護機上 I/O 不受站台故障影響。此程序不涉及磁碟機層級的資料複寫。擁有該集合體的控制器負責將資料複寫至集合體中的兩個叢集、但在站台遺失時仍必須保護資料、避免在執行中遺失 I/O。只有當合作夥伴控制器必須接管故障控制器時、才會使用複寫的 NVRAM 資料。

### 站台和機櫃故障保護：SyncMirror 和叢

SyncMirror 是一項鏡射技術、可增強但不取代 RAID DP 或 RAID-TEC。它會鏡射兩個不同 RAID 群組的內容。邏輯組態如下：

1. 磁碟機會根據位置設定成兩個集區。一個集區由站台 A 上的所有磁碟機組成、第二個集區由站台 B 上的所有磁碟機組成
2. 接著會根據鏡射的 RAID 群組集建立通用儲存池（稱為 Aggregate）。從每個站台擷取的磁碟機數量相等。例如、20 個磁碟機的 SyncMirror Aggregate 將由站台 A 的 10 個磁碟機和站台 B 的 10 個磁碟機組成
3. 指定站台上的每組磁碟機都會自動設定為一個或多個完全備援的 RAID DP 或 RAID-TEC 群組、而不受鏡像的使用影響。在鏡射下使用 RAID、即使在站台遺失之後、也能提供資料保護。



上圖說明 SyncMirror 組態範例。在控制器上建立了 24 個磁碟機的集合體、其中 12 個磁碟機來自於站台 A 上配置的機櫃、12 個磁碟機來自站台 B 上配置的機櫃磁碟機分為兩個鏡射 RAID 群組。RAID 群組 0 包含站台 A 的 6 磁碟機叢、鏡射到站台 B 的 6 磁碟機叢同樣地、RAID 群組 1 也包含站台 A 的 6 磁碟機叢、鏡射到站台 B 的 6 磁碟機叢

SyncMirror 通常用於提供 MetroCluster 系統的遠端鏡射、每個站台都有一份資料複本。有時候、它是用來在單一系統中提供額外的備援層級。特別是提供機架層級的備援。磁碟機櫃已包含雙電源供應器和控制器、整體上比金屬板稍多、但在某些情況下、可能需要額外的保護。例如、有一位 NetApp 客戶部署 SyncMirror、用於汽車測試期間使用的行動即時分析平台。系統分為兩個實體機架、分別隨附獨立的電源饋送和獨立的 UPS 系統。

#### 備援故障：NVFAIL

如前所述、寫入必須先登入本機 NVRAM 及至少一個其他控制器上的 NVRAM、才會被確認。此方法可確保硬體故障或停電不會導致機內 I/O 遺失如果本機 NVRAM 故障或連線至其他節點失敗、則資料將不再鏡射。

如果本機 NVRAM 回報錯誤、節點會關機。當使用 HA 配對時、此關機會導致容錯移轉至合作夥伴控制器。使

用 MetroCluster 時、行為取決於所選的整體組態、但可能會導致自動容錯移轉至遠端記事。無論如何、由於發生故障的控制器尚未確認寫入作業、因此不會遺失任何資料。

站台對站台連線故障會封鎖 NVRAM 複寫至遠端節點、這種情況更為複雜。寫入不再複寫到遠端節點、因此如果控制器發生災難性錯誤、可能會導致資料遺失。更重要的是、在這些情況下、嘗試容錯移轉至其他節點會導致資料遺失。

控制因素是 NVRAM 是否同步。如果 NVRAM 已同步、則節點對節點容錯移轉可安全地繼續進行、不會有資料遺失的風險。在 MetroCluster 組態中、如果 NVRAM 和基礎 Aggregate plex 同步、則可以安全地繼續進行轉換、而不會有資料遺失的風險。

除非強制進行容錯移轉或切換、否則 ONTAP 不允許在資料不同步時進行容錯移轉或切換。以這種方式強制變更條件、即表示資料可能會留在原始控制器中、而且資料遺失是可以接受的。

如果強制進行容錯移轉或切換、資料庫和其他應用程式尤其容易毀損、因為它們會在磁碟上保留較大的內部資料快取。如果發生強制容錯移轉或切換、先前確認的變更將會有效捨棄。儲存陣列的內容會有效地及時向後跳轉、而且快取狀態不再反映磁碟上資料的狀態。

為了避免這種情況發生、ONTAP 允許設定磁碟區、以針對 NVRAM 故障提供特殊保護。觸發時、此保護機制會導致磁碟區進入稱為 NVFAIL 的狀態。此狀態會導致 I/O 錯誤、導致應用程式當機。這項當機會導致應用程式關機、使其不使用過時的資料。資料不應遺失、因為記錄中應存在任何已認可的交易資料。通常的後續步驟是讓系統管理員在手動將 LUN 和磁碟區重新上線之前、先完全關閉主機。雖然這些步驟可能涉及一些工作、但這種方法是確保資料完整性的最安全方法。並非所有資料都需要這項保護、因此 NVFAIL 行為可依每個磁碟區設定。

## HA 配對與 MetroCluster

MetroCluster 提供兩種組態：雙節點和 HA 配對。雙節點組態在 NVRAM 上的運作方式與 HA 配對相同。如果發生突然故障、合作夥伴節點可以重新執行 NVRAM 資料、以確保磁碟機一致、並確保沒有遺失任何已確認的寫入資料。

HA 配對組態也會將 NVRAM 複寫到本機合作夥伴節點。簡單的控制器故障會在合作夥伴節點上重新執行 NVRAM、而獨立 HA 配對則不使用 MetroCluster。萬一突然完全遺失站台、遠端站台也需要 NVRAM、才能讓磁碟機保持一致、開始提供資料。

MetroCluster 的一個重要層面是、在正常作業條件下、遠端節點無法存取合作夥伴資料。每個站台基本上都是一個可假設對方站台特性的個別系統。此程序稱為「轉換」、包含計畫性的轉換、可在不中斷營運的情況下、將站合作業移轉至另一個站台。它也包括站台遺失的非計畫性情況、以及災難恢復需要手動或自動切換。

### 切換與切換

術語切換和切換是指在 MetroCluster 組態中、在遠端控制器之間轉換磁碟區的程序。此程序僅適用於遠端節點。在四個磁碟區組態中使用 MetroCluster 時、本機節點容錯移轉是先前所述的相同接管和恢復程序。

### 計畫性切換與切換

規劃的切換或切換類似於節點之間的接管或恢復。此程序有多個步驟、可能需要幾分鐘的時間、但實際發生的是儲存設備和網路資源的多階段順暢轉換。控制傳輸的速度比執行完整命令所需的時間快得多。

接管 / 恢復與切換 / 切換回復之間的主要差異在於對 FC SAN 連線能力的影響。使用本機接管 / 恢復功能、主機會遺失通往本機節點的所有 FC 路徑、並仰賴其原生 MPIO 來切換至可用的替代路徑。連接埠不會重新定位。透過切換和切換、控制器上的虛擬 FC 目標連接埠會轉換到另一個站台。它們在 SAN 上實際上已經停用一段時間、然後重新出現在替代控制器上。

## SyncMirror 逾時

SyncMirror 是一項 ONTAP 鏡射技術、可針對機櫃故障提供保護。當機櫃之間相隔一段距離時、就能獲得遠端資料保護。

SyncMirror 無法提供通用同步鏡像。因此、可用度更高。有些儲存系統使用固定的全或全自動鏡射、有時稱為 Domino 模式。這種形式的鏡像在應用程式中受到限制、因為如果與遠端站台的連線中斷、所有寫入活動都必須停止。否則、寫字會存在於某個站台、但不會存在於另一個站台。一般而言、如果站台對站台連線中斷超過一段短時間（例如 30 秒）、這類環境就會設定為使 LUN 離線。

這種行為是小型環境子集的理想選擇。不過、大多數應用程式都需要一套解決方案、能夠在正常作業條件下提供保證同步複寫、但能夠暫停複寫。站台對站台連線能力完全中斷通常被視為近乎災難的情況。一般而言、這類環境會保持在線上狀態並提供資料、直到連線能力修復或正式決定關閉環境以保護資料為止。純粹因為遠端複寫失敗而需要自動關閉應用程式、這是不尋常的。

SyncMirror 支援同步鏡射需求、並可靈活調整逾時時間。如果與遠端控制器和 / 或叢的連線中斷、30 秒定時器就會開始倒數。當計數器達到 0 時、會使用本機資料繼續寫入 I/O 處理。資料的遠端複本可以使用、但會在連線恢復之前、及時凍結。重新同步利用 Aggregate 層級快照、將系統儘快恢復至同步模式。

值得注意的是、在許多情況下、這種通用的「全或全無」Domino 模式複寫功能更適合在應用程式層上實作。例如、Oracle DataGuard 包括最大保護模式、可在任何情況下保證執行個體的長時間複寫。如果複寫連結失敗超過可設定的逾時時間、資料庫就會關閉。

### 使用 Fabric 附加 MetroCluster 自動進行無人值守切換

自動無人值守切換（AUSO）是一項 Fabric 附加 MetroCluster 功能、可提供一種跨站台 HA 的形式。如前所述、MetroCluster 有兩種類型：每個站台上只有一個控制器、或每個站台上有一個 HA 配對。HA 選項的主要優點是、計畫性或非計畫性控制器關機仍可讓所有 I/O 成為本機。單一節點選項的優勢在於降低成本、複雜度和基礎架構。

AUSO 的主要價值在於改善 Fabric 附加 MetroCluster 系統的 HA 功能。每個站台都會監控相對站台的健全狀況、如果沒有節點仍可提供資料、AUSO 就會導致快速的轉換。這種方法在每個站台只有一個節點的 MetroCluster 組態中特別有用、因為在可用度方面、它使組態更接近 HA 配對。

AUSO 無法在 HA 配對層級提供全方位監控。HA 配對可提供極高的可用度、因為它包含兩條備援實體纜線、可用於直接節點對節點通訊。此外、HA 配對中的兩個節點都能存取備援迴圈上的同一組磁碟、為一個節點提供另一條路由來監控另一個節點的健全狀況。

MetroCluster 叢集存在於站台之間、節點對節點通訊和磁碟存取都仰賴站台對站台網路連線。監控叢集其餘部分的活動訊號的能力有限。AUSO 必須區分其他站台實際停機、而非因為網路問題而無法使用的情況。

因此、如果 HA 配對中的控制器偵測到因特定原因（例如系統異常）而發生的控制器故障、就會提示接管。如果連線完全中斷、也可能會提示接管、有時也稱為「失去心跳」。

只有在原始站台偵測到特定故障時、MetroCluster 系統才能安全地執行自動切換。此外、擁有儲存系統所有權的控制器必須能夠保證磁碟和 NVRAM 資料同步。控制器無法保證進行變更的安全性、因為它與來源站台失去接觸、而該站台仍可運作。如需將交換作業自動化的其他選項、請參閱下一節中的 MetroCluster tiebreaker（MCTB）解決方案資訊。

### MetroCluster tiebreaker 搭配網路附加 MetroCluster

- ["NetApp MetroCluster tiebreaker"](#) 軟體可在第三個站台上執行、以監控 MetroCluster 環境的健全狀況、傳送通知、並在災難情況下強制切換。如需有關斷路器的完整說明、請參閱 ["NetApp 支援網站"](#)但 MetroCluster 斷路

器的主要用途是偵測站台遺失。它還必須區分站台遺失和連線中斷。例如、不應因為斷路器無法到達主要站台而進行切入、這就是為什麼斷路器也會監控遠端站台與主要站台聯絡的能力。

與 AUSO 的自動切換功能也相容於 MCTB。AUSO 反應非常迅速、因為它的設計是偵測特定故障事件、然後只有在 NVRAM 和 SyncMirror 叢同步時才叫用切入。

相反地、斷路器位於遠端位置、因此必須等到定時器結束後才會宣告站台停機。tiebreaker 最終會偵測 AUSO 涵蓋的控制器故障類型、但一般而言、AUSO 已經開始進行開關作業、而且可能會在 tiebreaker 運作之前完成開關作業。產生的第二個來自 tiebreaker 的切換命令將會遭到拒絕。

- 注意：\* 強制切入時、MCTB 軟體無法驗證 NVRAM 是否與 / 或叢同步。如果已設定自動切換、則應在維護活動期間停用、導致 NVRAM 或 SyncMirror 叢同步中斷。

此外、MCTB 可能無法因應導致下列事件順序的滾動災難：

1. 站台之間的連線中斷超過 30 秒。
2. SyncMirror 複寫逾時、且作業會繼續在主要站台上執行、使遠端複本過時。
3. 主站台會遺失。結果是主站台上存在未複寫的變更。因此、由於下列幾個原因、可能不希望進行任何一次的重新操作：
  - 關鍵資料可能會出現在主要站台上、而且該資料最終可能會恢復。允許應用程式繼續作業的轉換作業、將會有效捨棄該關鍵資料。
  - 當站台遺失時、使用主要站台上儲存資源的仍在運作中站台上的應用程式可能已快取資料。切入會導致資料的過時版本與快取不相符。
  - 當發生站台遺失時、使用主要站台上儲存資源的仍在運作中站台上的作業系統、可能已快取資料。切入會導致資料的過時版本與快取不相符。最安全的選項是將斷路器設定為在偵測到站台故障時傳送警示、然後讓人員決定是否強制進行轉換。應用程式和（或）作業系統可能需要先關機、才能清除任何快取資料。此外、NVFAIL 設定也可用於新增進一步的保護、並協助簡化容錯移轉程序。

## ONTAP Mediator 搭配 MetroCluster IP

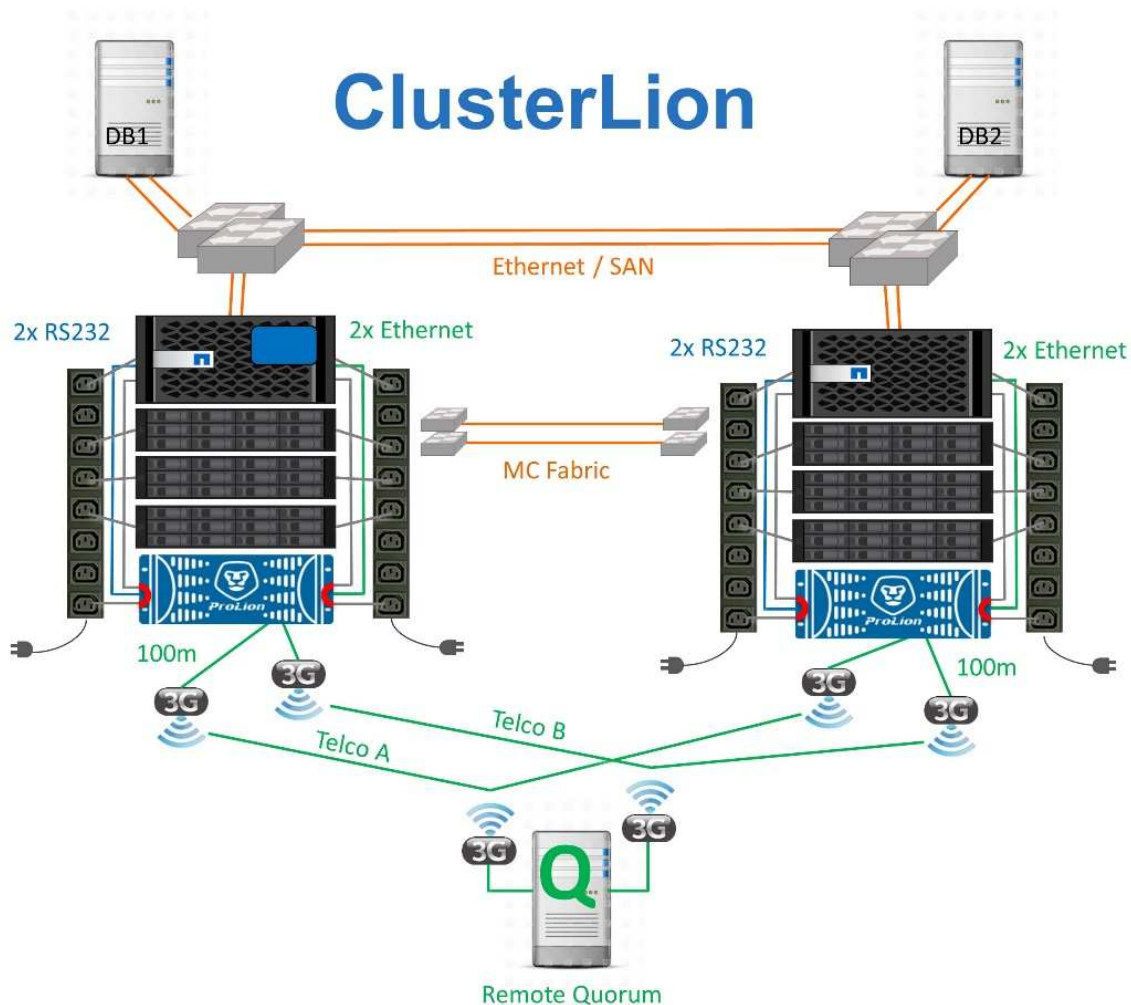
ONTAP Mediator 可搭配 MetroCluster IP 和某些其他 ONTAP 解決方案使用。它是一項傳統的斷路器服務、就像上述的 MetroCluster tiebreaker 軟體一樣、但也包含一項重要功能、即執行自動無人值守的移除。

光纖連接的 MetroCluster 可直接存取位於相對站台的儲存裝置。這可讓一個 MetroCluster 控制器從磁碟機讀取心跳資料、以監控其他控制器的健全狀況。這可讓一個控制器辨識另一個控制器的故障、並執行切換。

相反地、MetroCluster IP 架構只會透過控制器控制器連線路由所有 I/O、而無法直接存取遠端站台上的儲存裝置。這會限制控制器偵測故障和執行轉換的能力。因此、ONTAP Mediator 必須作為斷路器裝置、才能偵測站台遺失並自動執行轉換。

## 使用 ClusterLion 的虛擬第三站點

ClusterLion 是一款先進的 MetroCluster 監控設備、可作為虛擬第三站點使用。此方法可讓 MetroCluster 安全部署在雙站台組態中、並具備全自動的轉換功能。此外、ClusterLion 還能執行額外的網路層級監控、並執行後置作業。完整文件可從 ProLion 取得。



- ClusterLion 設備會使用直接連接的乙太網路和序列纜線來監控控制器的健全狀況。
- 這兩台設備透過備援的 3G 無線連線彼此連線。
- ONTAP 控制器的電源會透過內部中繼路由傳送。發生站台故障時、包含內部 UPS 系統的 ClusterLion 會先切斷電源連線、然後再啟動切入。此程序可確保不會發生任何大腦分割狀況。
- ClusterLion 會在 30 秒 SyncMirror 逾時內執行切換、或完全不執行。
- 除非 NVRAM 和 SyncMirror 叢集的狀態同步、否則 ClusterLion 不會執行切入。
- 由於 ClusterLion 只會在 MetroCluster 完全同步時執行切入、因此不需要 NVFAIL。此組態可讓擴充 Oracle RAC 等站台跨距環境保持連線、即使在非計畫性的轉換期間亦然。
- 支援包括光纖連接的 MetroCluster 和 MetroCluster IP

## SyncMirror 的 Oracle 資料庫

SyncMirror 是 MetroCluster 系統的 Oracle 資料保護基礎、是最大效能的橫向擴充同步鏡射技術。

### 使用 SyncMirror 保護資料

在最簡單的層級上、同步複寫表示必須先對鏡射儲存設備的兩側進行任何變更、然後才會被確認。例如、如果資料庫正在寫入記錄檔、或是正在修補 VMware 來賓作業系統、則寫入作業絕不能遺失。作為一種協議級別、在



兩個站點上的非易失性介質被認可之前，存儲系統不得確認寫入內容。只有這樣、在不遺失資料的風險下繼續作業是安全的。

使用同步複寫技術是設計和管理同步複寫解決方案的第一步。最重要的考量是瞭解在各種計畫性和非計畫性失敗案例中可能發生的情況。並非所有同步複寫解決方案都提供相同的功能。如果您需要提供零恢復點目標（RPO）的解決方案、亦即零資料遺失、則必須考慮所有故障情況。特別是、當站台之間的連線中斷而無法進行複寫時、預期會產生什麼結果？

#### SyncMirror 資料可用度

MetroCluster 複寫是以 NetApp SyncMirror 技術為基礎、其設計旨在有效率地切換至同步模式及從同步模式切換到同步模式。這項功能符合要求同步複寫、但也需要高可用度資料服務的客戶需求。例如、如果中斷與遠端站台的連線、通常最好讓儲存系統繼續以非複寫狀態運作。

許多同步複寫解決方案只能以同步模式運作。這種類型的全或全無複寫有時稱為 Domino 模式。這類儲存系統會停止提供資料、而不允許資料的本機和遠端複本進行非同步處理。如果複寫被強制中斷、重新同步可能會非常耗時、而且可能會讓客戶在重新建立鏡像期間暴露在完全資料遺失的風險中。

SyncMirror 不僅可以在無法連線到遠端站台時、無縫切換至同步模式、也可以在連線恢復時、快速重新同步至 RPO = 0 狀態。遠端站台的資料過時複本也可在重新同步期間保留為可用狀態、以確保資料的本機和遠端複本隨時都存在。

在需要 Domino 模式的情況下、NetApp 提供 SnapMirror 同步（SM-S）。應用程式層級選項也存在、例如 Oracle DataGuard 或主機端磁碟鏡射的延長逾時。如需其他資訊和選項、請洽詢您的 NetApp 或合作夥伴客戶團隊。

#### 使用 MetroCluster 進行 Oracle 資料庫容錯移轉

Metrocluster is an ONTAP feature that can protect your Oracle databases with RPO=0 synchronous mirroring across sites, and it scales up to support hundreds of databases on a single MetroCluster system. It's also simple to use. The use of MetroCluster does not necessarily add to or change any best practices for operating a enterprise applications and databases. 通常的最佳實務做法仍適用、如果您的需求只需要 RPO = 0 資料保護、則 MetroCluster 會滿足您的需求。然而、大多數客戶不僅使用 MetroCluster 來保護 RPO = 0 資料、還能在災難期間改善 RTO、並在站台維護活動中提供透明的容錯移轉。

#### 使用預先設定的作業系統進行容錯移轉

SyncMirror 在災難恢復站點上提供資料的同步複本、但要讓資料可用、則需要作業系統和相關應用程式。基本自動化可大幅改善整體環境的容錯移轉時間。Oracle RAC、Veritas 叢集伺服器（VCS）或 VMware HA 等叢集式產品通常用於在站台之間建立叢集、在許多情況下、容錯移轉程序可以使用簡單的指令碼來驅動。

如果主節點遺失、叢集軟體（或指令碼）會設定為在替代站台上線應用程式。其中一個選項是建立預先針對構成應用程式的 NFS 或 SAN 資源所預先設定的待命伺服器。如果主站台發生故障、叢集軟體或指令碼替代方案會執行類似下列的一系列動作：

1. 強制 MetroCluster 進行重新操作
2. 執行 FC LUN 探索（僅限 SAN）

### 3. 掛載檔案系統

### 4. 啟動應用程式

此方法的主要需求是在遠端站台上執行作業系統。它必須預先設定應用程式二進位檔、也就是說、修補等工作必須在主要站台和待命站台上執行。或者、應用程式二進位檔可鏡射至遠端站台、並在宣告災難時掛載。

實際的啟動程序很簡單。LUN 探索等命令每個 FC 連接埠只需要幾個命令。檔案系統掛載只不過是 mount 只需一個命令、即可在 CLI 上啟動和停止資料庫和 ASM。如果在切換之前、磁碟區和檔案系統並未在災難恢復站台上使用、則無需設定 `dr-force- nvfail` 在磁碟區上。

使用虛擬化作業系統進行容錯移轉

資料庫環境的容錯移轉可延伸至包含作業系統本身。理論上、此容錯移轉可以使用開機 LUN 來完成、但通常是使用虛擬化的作業系統來完成。此程序類似於下列步驟：

1. 強制 MetroCluster 進行重新操作
2. 裝載託管資料庫伺服器虛擬機器的資料存放區
3. 啟動虛擬機器
4. 手動啟動資料庫、或將虛擬機器設定為自動啟動資料庫

例如、ESX 叢集可以跨越站台。在發生災難時、虛擬機器可在移至災難恢復站台後上線。只要主控虛擬化資料庫伺服器的資料存放區在災難發生時並未使用、就不需要設定 `dr-force- nvfail` 在相關的磁碟區上。

## Oracle 資料庫、MetroCluster 和 NVFAIL

NVFAIL 是 ONTAP 中的一般資料完整性功能、其設計可讓資料庫發揮最大的資料完整性保護。



本節將進一步說明基本的 ONTAP NVFAIL、以涵蓋 MetroCluster 專屬主題。

使用 MetroCluster 時、寫入必須登入至少一個其他控制器的本機 NVRAM 和 NVRAM、才能被確認。此方法可確保硬體故障或停電不會導致機內 I/O 遺失如果本機 NVRAM 故障或連線至其他節點失敗、則資料將不再鏡射。

如果本機 NVRAM 回報錯誤、節點會關機。當使用 HA 配對時、此關機會導致容錯移轉至合作夥伴控制器。使用 MetroCluster 時、行為取決於所選的整體組態、但可能會導致自動容錯移轉至遠端記事。無論如何、由於發生故障的控制器尚未確認寫入作業、因此不會遺失任何資料。

站台對站台連線故障會封鎖 NVRAM 複寫至遠端節點、這種情況更為複雜。寫入不再複寫到遠端節點、因此如果控制器發生災難性錯誤、可能會導致資料遺失。更重要的是、在這些情況下、嘗試容錯移轉至其他節點會導致資料遺失。

控制因素是 NVRAM 是否同步。如果 NVRAM 已同步、則節點對節點容錯移轉可安全地繼續進行、而不會有資料遺失的風險。在 MetroCluster 組態中、如果 NVRAM 和基礎 Aggregate plex 同步、則在不遺失資料的情況下繼續進行轉換是安全的。

除非強制進行容錯移轉或切換、否則 ONTAP 不允許在資料不同步時進行容錯移轉或切換。以這種方式強制變更條件、即表示資料可能會留在原始控制器中、而且資料遺失是可以接受的。

如果強制進行容錯移轉或切換、則資料庫特別容易遭到毀損、因為資料庫會在磁碟上保留較大的內部資料快取。如果發生強制容錯移轉或切換、先前確認的變更將會有效捨棄。儲存陣列的內容會有效地及時向後跳轉、而且資

料庫快取的狀態不再反映磁碟上資料的狀態。

為了保護應用程式不受這種情況影響、ONTAP 允許設定磁碟區、以針對 NVRAM 故障提供特殊保護。觸發時、此保護機制會導致磁碟區進入稱為 NVFAIL 的狀態。此狀態會導致 I/O 錯誤、導致應用程式關機、使其不使用過時的資料。資料不應遺失、因為儲存系統上仍有任何已確認的寫入資料、而資料庫則應在記錄中顯示任何已認可的交易資料。

通常的後續步驟是讓系統管理員在手動將 LUN 和磁碟區重新上線之前、先完全關閉主機。雖然這些步驟可能涉及一些工作、但這種方法是確保資料完整性的最安全方法。並非所有資料都需要這項保護、因此 NVFAIL 行為可依每個磁碟區設定。

#### 手動強制 NVFAIL

最安全的選項是透過指定來強制轉換跨站台散佈的應用程式叢集（包括 VMware、Oracle RAC 及其他）`-force-nvfail-all` 在命令列。此選項可作為緊急措施使用、以確保所有快取資料均已清除。如果主機使用的儲存資源原本位於災難性站台上、則會收到 I/O 錯誤或過時的檔案處理 (ESTALE) 錯誤。Oracle 資料庫當機、檔案系統可能完全離線、或切換至唯讀模式。

在完成重新操作之後、`in-nvfailed-state` 需要清除旗標、且 LUN 必須置於線上。完成此活動後、即可重新啟動資料庫。這些工作可以自動化、以降低 RTO。

#### `dr-force-nvfail`

作為一般安全措施、請設定 `dr-force-nvfail` 在所有可能在正常作業期間從遠端站台存取的磁碟區上加上旗標、表示這些磁碟區是在容錯移轉之前使用的活動。此設定的結果是、選取的遠端磁碟區在進入時無法使用 `in-nvfailed-state` 在進行重新操作時。在完成重新操作之後、`in-nvfailed-state` 旗標必須清除、且 LUN 必須置於線上。這些活動完成後、即可重新啟動應用程式。這些工作可以自動化、以降低 RTO。

結果就像使用 `-force-nvfail-all` 手動切換的旗標。然而、受影響的磁碟區數量可能僅限於必須受到保護的磁碟區、不受具有過時快取的應用程式或作業系統的影響。

對於不使用的環境、有兩項關鍵需求 `dr-force-nvfail` 在應用程式磁碟區上：

- 在主站台遺失後、強制進行的重新操作不得超過 30 秒。
- 在維護工作期間、或是在 SyncMirror 叢或 NVRAM 複寫不同步的任何其他情況下、切勿進行切入。第一項需求可以透過使用已設定為在站台故障 30 秒內執行轉換的斷路器軟體來達成。這並不表示切入作業必須在偵測站台故障的 30 秒內執行。這表示、如果站台確認運作已過 30 秒、就不再安全地強制進行轉換。

第二項需求可在已知 MetroCluster 組態不同步時停用所有自動切換功能、以部分滿足。更好的選擇是擁有可監控 NVRAM 複寫和 SyncMirror 叢的健全狀況的斷路器解決方案。如果叢集未完全同步、則斷路器不應觸發切入。

NetApp MCTB 軟體無法監控同步處理狀態、因此當 MetroCluster 因任何原因而未同步時、應該停用同步處理狀態。ClusterLion 確實包含 NVRAM 監控和叢監視功能、除非 MetroCluster 系統確認完全同步、否則可將其設定為不觸發切入。

#### MetroCluster 上的 Oracle 單一執行個體

如前所述、MetroCluster 系統的存在並不一定會新增或變更任何操作資料庫的最佳實務做法。目前在客戶 MetroCluster 系統上執行的大多數資料庫都是單一執行個體、並遵循 Oracle on ONTAP 文件中的建議。

使用預先設定的作業系統進行容錯移轉

SyncMirror 在災難恢復站點上提供資料的同步複本、但要讓資料可用、則需要作業系統和相關應用程式。基本自動化可大幅改善整體環境的容錯移轉時間。例如 Veritas Cluster Server (VCS) 等叢集產品通常用於在站台之間建立叢集、而且在許多情況下、容錯移轉程序可以使用簡單的指令碼來驅動。

如果主節點遺失、叢集軟體（或指令碼）會設定為在替代站台上線資料庫。其中一個選項是建立預先設定為 NFS 或 SAN 資源的備用伺服器、以供組成資料庫。如果主站台發生故障、叢集軟體或指令碼替代方案會執行類似下列的一系列動作：

1. 強制 MetroCluster 進行重新操作
2. 執行 FC LUN 探索（僅限 SAN）
3. 掛載檔案系統和 / 或掛載 ASM 磁碟群組
4. 啟動資料庫

此方法的主要需求是在遠端站台上執行作業系統。它必須預先設定 Oracle 二進位檔、這也表示 Oracle 修補等工作必須在主要站台和待命站台上執行。或者、Oracle 二進位檔可鏡射至遠端站台、並在宣告災難時掛載。

實際的啟動程序很簡單。LUN 探索等命令每個 FC 連接埠只需要幾個命令。檔案系統掛載只不過是 mount 只需一個命令、即可在 CLI 上啟動和停止資料庫和 ASM。如果在切換之前、磁碟區和檔案系統並未在災難恢復站台上使用、則無需設定 `dr-force-nvfail` 在磁碟區上。

使用虛擬化作業系統進行容錯移轉

資料庫環境的容錯移轉可延伸至包含作業系統本身。理論上、此容錯移轉可以使用開機 LUN 來完成、但通常是使用虛擬化的作業系統來完成。此程序類似於下列步驟：

1. 強制 MetroCluster 進行重新操作
2. 裝載託管資料庫伺服器虛擬機器的資料存放區
3. 啟動虛擬機器
4. 手動啟動資料庫、或將虛擬機器設定為自動啟動資料庫、例如 ESX 叢集可能跨越站台。在發生災難時、虛擬機器可在移至災難恢復站台後上線。只要主控虛擬化資料庫伺服器的資料存放區在災難發生時並未使用、就不需要設定 `dr-force-nvfail` 在相關的磁碟區上。

## MetroCluster 上的延伸 Oracle RAC

許多客戶透過在各個站台之間延伸 Oracle RAC 叢集來最佳化 RTO、進而實現完全主動式的組態。整體設計變得更複雜、因為它必須包含 Oracle RAC 的仲裁管理。此外、從兩個站台存取資料、這表示強制轉換可能會導致使用過時的資料複本。

雖然兩個站台上都有資料複本、但只有目前擁有 Aggregate 的控制器才能提供資料。因此、使用擴充的 RAC 叢集時、遠端節點必須透過站台對站台連線來執行 I/O。結果會增加 I/O 延遲、但這種延遲通常不是問題。RAC 互連網路也必須延伸至站台、這表示無論如何都需要高速、低延遲的網路。如果增加的延遲確實造成問題、則叢集可以主動被動方式運作。接著、需要將 I/O 密集作業導向至擁有該集合體的控制器本機的 RAC 節點。然後、遠端節點會執行較輕的 I/O 作業、或純粹作為暖待機伺服器使用。

如果需要雙主動式擴充 RAC、則應考慮使用 ASM 鏡像來取代 MetroCluster。ASM 鏡像可讓您偏好資料的特定複本。因此、可以內建擴充 RAC 叢集、讓所有讀取作業都在本機進行。讀取 I/O 永遠不會跨越網站、因此可提供最低的延遲。所有寫入活動仍必須傳輸站台間連線、但任何同步鏡射解決方案都無法避免此類流量。



如果開機 LUN（包括虛擬化開機磁碟）與 Oracle RAC 搭配使用、請使用 `misscount` 可能需要變更參數。如需 RAC 逾時參數的詳細資訊、請參閱 "[Oracle RAC 搭配 ONTAP](#)"。

## 雙站台組態

雙站台擴充 RAC 組態可提供雙主動式資料庫服務、可在不中斷營運的情況下、在許多（但並非全部）災難案例中順利運作。

## RAC 投票檔案

在 MetroCluster 上部署擴充 RAC 時、首先應考慮仲裁管理。Oracle RAC 有兩種機制可管理仲裁：磁碟心跳和網路心跳。磁碟心跳會使用投票檔案來監控儲存設備存取。只要基礎儲存系統提供 HA 功能、單一投票資源就足以搭配單一站台 RAC 組態。

在早期版本的 Oracle 中、投票檔案會放置在實體儲存裝置上、但在目前版本的 Oracle 中、投票檔案會儲存在 ASM 磁碟群組中。



NFS 支援 Oracle RAC。在網格安裝程序期間、會建立一組 ASM 程序、將用於網格檔案的 NFS 位置顯示為 ASM 磁碟群組。此程序對終端使用者來說幾乎透明、安裝完成後不需要持續進行 ASM 管理。

雙站台組態的第一項需求是確保每個站台都能以保證不中斷災難恢復程序的方式存取超過半數的投票檔案。這項工作在投票檔案儲存在 ASM 磁碟群組之前很簡單、但現在管理員必須瞭解 ASM 備援的基本原則。

ASM 磁碟群組有三種備援選項 `external`、`normal` 和 `high`。換句話說、非鏡射、鏡射和 3 向鏡射。名為的較新選項 `flex` 也可以使用、但很少使用。備援裝置的備援層級和放置位置可控制故障情況發生的情況。例如：

- 將投票檔案放在上 `diskgroup` 與 `external` 如果站台間連線中斷、備援資源保證可收回一個站台。
- 將投票檔案放在上 `diskgroup` 與 `normal` 如果站台間連線中斷、每個站台只有一個 ASM 磁碟的備援功能可確保兩個站台的節點遷離、因為兩個站台都不會有大部分的仲裁。
- 將投票檔案放在上 `diskgroup` 與 `high` 當兩個站台都可以運作且彼此可連線時、一個站台上有一個磁碟和另一個站台上的單一磁碟的備援功能可讓雙主動式作業運作。但是、如果單一磁碟站台與網路隔離、則該站台會被逐出。

## RAC 網路心跳

Oracle RAC 網路活動訊號可監控叢集互連中的節點可連性。若要保留在叢集中、節點必須能夠連絡其他節點的一半以上。在雙站台架構中、此需求會為 RAC 節點數建立下列選項：

- 如果每個站台放置相同數量的節點、則會在網路連線中斷時、在某個站台上造成遷離。
- 在另一個站台上放置 N 個節點、在另一個站台上放置 N+1 個節點、可確保站台之間的連線中斷、導致站台的網路仲裁中剩餘節點數量較多、而節點移出數量較少的站台。

在 Oracle 12cR2 之前、無法控制哪一方在站台遺失時會發生遷離。當每個站台的節點數量相等時、會由主要節點控制遷離、這通常是第一個要開機的 RAC 節點。

Oracle 12cR2 引進節點加權功能。這項功能可讓管理員更有效地控制 Oracle 如何解決大腦分裂狀況。例如、下列命令可設定 RAC 中特定節點的偏好設定：

```
[root@host-a ~]# /grid/bin/crsctl set server css_critical yes
CRS-4416: Server attribute 'CSS_CRITICAL' successfully changed. Restart
Oracle High Availability Services for new value to take effect.
```

重新啟動 Oracle 高可用度服務後、組態如下所示：

```
[root@host-a lib]# /grid/bin/crsctl status server -f | egrep
'^NAME|CSS_CRITICAL='
NAME=host-a
CSS_CRITICAL=yes
NAME=host-b
CSS_CRITICAL=no
```

節點 host-a 現已指定為關鍵伺服器。如果兩個 RAC 節點是隔離的、host-a 生存、和 host-b 被逐出。



如需完整詳細資料、請參閱 Oracle 白皮書《Oracle Clusterware 12c Release 2 Technical Overview》。

對於 12cR2 之前的 Oracle RAC 版本、可透過檢查 CRS 記錄來識別主節點、如下所示：

```
[root@host-a ~]# /grid/bin/crsctl status server -f | egrep
'^NAME|CSS_CRITICAL='
NAME=host-a
CSS_CRITICAL=yes
NAME=host-b
CSS_CRITICAL=no
[root@host-a ~]# grep -i 'master node' /grid/diag/crs/host-
a/crs/trace/crsd.trc
2017-05-04 04:46:12.261525 : CRSSE:2130671360: {1:16377:2} Master Change
Event; New Master Node ID:1 This Node's ID:1
2017-05-04 05:01:24.979716 : CRSSE:2031576832: {1:13237:2} Master Change
Event; New Master Node ID:2 This Node's ID:1
2017-05-04 05:11:22.995707 : CRSSE:2031576832: {1:13237:221} Master
Change Event; New Master Node ID:1 This Node's ID:1
2017-05-04 05:28:25.797860 : CRSSE:3336529664: {1:8557:2} Master Change
Event; New Master Node ID:2 This Node's ID:1
```

此記錄表示主節點為 2 和節點 host-a ID 為 1。這意味著 host-a 不是主節點。您可以使用命令確認主節點的身分識別 `olsnodes -n`。

```
[root@host-a ~]# /grid/bin/olsnodes -n
host-a 1
host-b 2
```

識別碼為的節點 2 是 host-b，這是主節點。在每個站台上節點數量相等的組態中、站台為 host-b 如果這兩組因為任何原因而失去網路連線、則該站台仍可生存。

識別主節點的記錄項目可能會超出系統的使用期限。在這種情況下、可以使用 Oracle 叢集登錄（OCR）備份的時間戳記。

```
[root@host-a ~]# /grid/bin/ocrconfig -showbackup
host-b      2017/05/05 05:39:53      /grid/cdata/host-cluster/backup00.ocr
0
host-b      2017/05/05 01:39:53      /grid/cdata/host-cluster/backup01.ocr
0
host-b      2017/05/04 21:39:52      /grid/cdata/host-cluster/backup02.ocr
0
host-a      2017/05/04 02:05:36      /grid/cdata/host-cluster/day.ocr      0
host-a      2017/04/22 02:05:17      /grid/cdata/host-cluster/week.ocr    0
```

此範例顯示主節點是 host-b。它也表示主節點的變更來源 host-a 至 host-b 5 月 4 日下午 2：05 至 21：39 之間。這種識別主節點的方法只有在也檢查了 CRS 記錄檔時才安全使用、因為主節點可能自上一次的 OCR 備份後變更。如果發生此變更、則應可在 OCR 記錄中看到。

大多數客戶選擇單一投票磁碟群組來服務整個環境、以及每個站台上相同數量的 RAC 節點。磁碟群組應放置在包含資料庫的網站上。結果是連線中斷會導致遠端站台被逐出。遠端站台將不再擁有仲裁、也無法存取資料庫檔案、但本機站台會繼續如常運作。連線恢復後、遠端執行個體即可重新上線。

發生災難時、需要進行轉換、才能讓資料庫檔案和投票磁碟群組在正常運作的網站上線。如果災難允許 AUSO 觸發切換、則不會觸發 NVFAIL、因為已知叢集處於同步狀態、且儲存資源正常上線。AUSO 是一項非常快速的作業、應在完成之前完成 disktimeout 期間過期。

由於只有兩個站台、因此無法使用任何類型的自動外部中斷軟體、這表示強制切換必須是手動操作。

### 三站台組態

擴充的 RAC 叢集可更輕鬆地建構三個站台。裝載 MetroCluster 系統每一半的兩個站台也支援資料庫工作負載、而第三個站台則是資料庫和 MetroCluster 系統的斷路器。Oracle tiebreaker 組態可能只需在第三站台上放置用於投票的 ASM 磁碟群組成員、也可能在第三站台上加入作業執行個體、以確保 RAC 叢集中有奇數個節點。



有關在擴展 RAC 配置中使用 NFS 的重要信息，請參閱 Oracle 文檔中的“quorum failure group（仲裁故障組）”。總而言之、NFS 掛載選項可能需要修改以包含軟選項、以確保主仲裁資源所在的第三站台連線中斷、不會使主 Oracle 伺服器或 Oracle RAC 程序掛起。

## SnapMirror 主動同步

## 採用 SnapMirror 主動同步的 Oracle 資料庫

SnapMirror 主動同步可針對個別 Oracle 資料庫和應用程式環境、啟用選擇性的 RPO = 0 同步鏡射。

SnapMirror 主動同步基本上是 SAN 的強化 SnapMirror 功能、可讓主機從主控 LUN 的系統以及主控其複本的系統存取 LUN。

SnapMirror 主動式同步和 SnapMirror 同步可共用複寫引擎、不過 SnapMirror 主動式同步則包含其他功能、例如企業應用程式的透明應用程式容錯移轉和容錯回復。

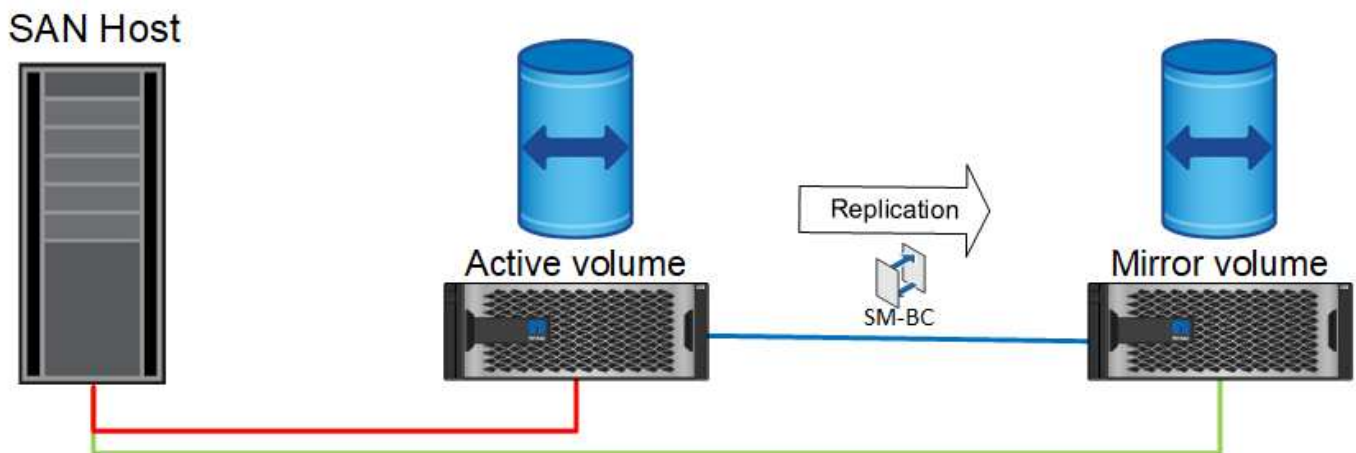
實際上、它的運作方式類似於精細的 MetroCluster 版本、可針對個別工作負載啟用選擇性且精細的 RPO = 0 同步複寫。低階路徑行為與 MetroCluster 非常不同、但主機觀點的最終結果卻相似。

### 路徑存取

透過 SnapMirror 主動式同步、儲存裝置可從主要和遠端儲存陣列的主機作業系統中看到。路徑是透過非對稱式邏輯單元存取 (ALUA) 進行管理、這是一種業界標準傳輸協定、用於識別儲存系統與主機之間的最佳化路徑。

存取 I/O 最短的裝置路徑被視為主動 / 最佳化路徑、其餘路徑則被視為主動 / 非最佳化路徑。

SnapMirror 主動同步關係是位於不同叢集上的一對 SVM 之間的關係。兩個 SVM 都能提供資料、但 ALUA 會優先使用 SVM、而 SVM 目前擁有 LUN 所在磁碟機的所有權。透過 SnapMirror 主動同步互連、將 IO 透過代理至遠端 SVM。



### 同步複寫

在正常作業中、遠端複本始終為 RPO = 0 同步複本、但有一個例外。如果無法複寫資料、SnapMirror 主動式同步將會釋放複寫資料和恢復服務 IO 的需求。如果客戶認為複寫連結遺失的情況近乎災難、或是不想在無法複寫資料時停止業務作業、則偏好使用此選項。

### 儲存硬體

與其他儲存災難恢復解決方案不同、SnapMirror 主動式同步提供非對稱式平台靈活度。每個站台的硬體不一定相同。此功能可讓您調整支援 SnapMirror 主動同步所用硬體的大小。如果遠端儲存系統需要支援完整的正式作業工作負載、則它可以與主要站台相同、但如果災難導致 I/O 減少、遠端站台上較小的系統可能會更具成本效益。



## 中間器ONTAP

ONTAP Mediator 是從 NetApp 支援下載的軟體應用程式。Mediator 可自動執行主與遠端站台儲存叢集的容錯移轉作業。它可以部署在內部部署或雲端的小型虛擬機器（VM）上。設定之後、它會成為第三個站台、用來監控兩個站台的容錯移轉案例。

### Oracle 資料庫容錯移轉搭配 SnapMirror 主動式同步

在 SnapMirror 主動式同步上託管 Oracle 資料庫的主要原因、是在計畫性和非計畫性儲存事件期間提供透明的容錯移轉。

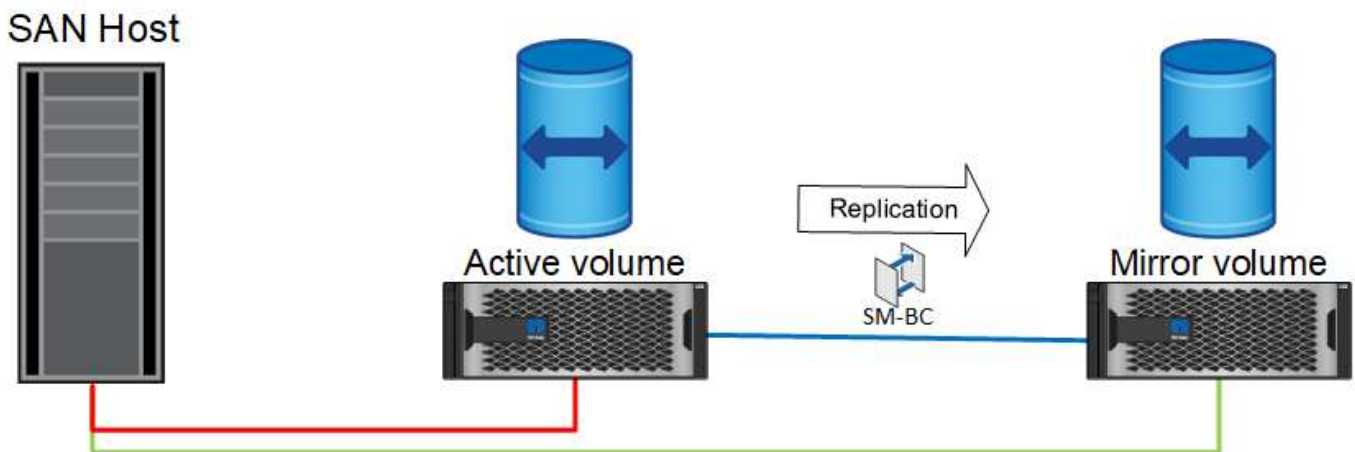
SnapMirror 主動式同步支援兩種儲存容錯移轉作業：計畫性和非計畫性、運作方式略有不同。系統管理員會手動啟動計畫性容錯移轉、以快速切換至遠端站台、而非計畫性容錯移轉則由第三站台的協調員自動啟動。計畫性容錯移轉的主要目的是執行漸進式修補與升級、執行災難恢復測試、或是採用正式的原則、在一年內在站台之間切換作業、以證明完整的主動式同步功能。

下圖顯示正常、容錯移轉及容錯回復作業期間的情況。為了便於說明、它們描述了複寫的 LUN。在實際的 SnapMirror 主動式同步組態中、複寫是以磁碟區為基礎、其中每個磁碟區都包含一個或多個 LUN、但為了讓圖片更簡單、磁碟區層已經移除。

#### 正常運作

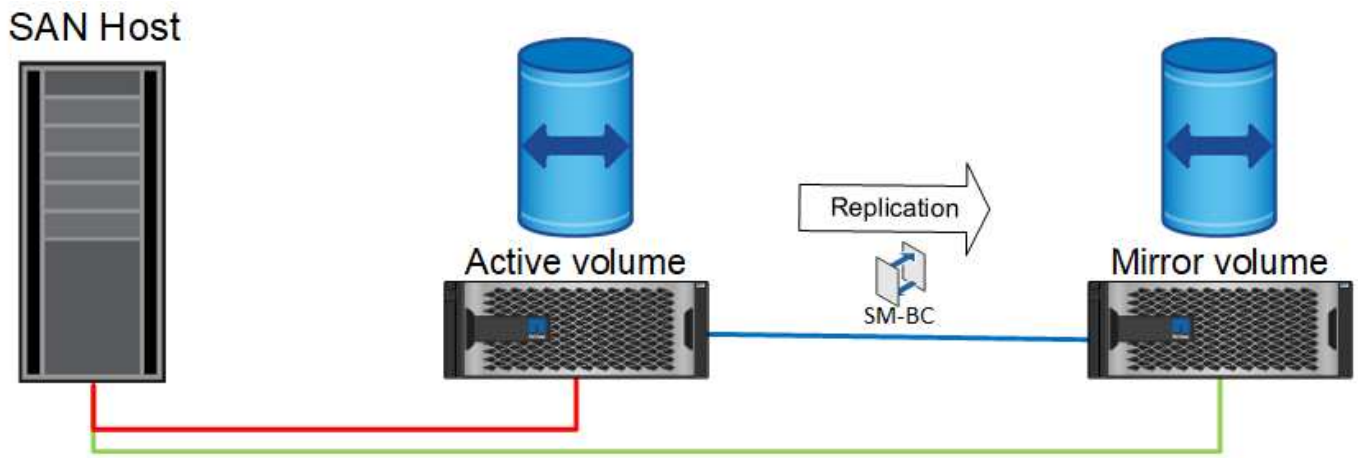
在正常作業中、可以從本機或遠端複本存取 LUN。紅色線表示 ALUA 所通告的最佳化路徑、結果應該是 IO 優先傳送至此路徑。

綠線是一條活動路徑，但由於該路徑上的 IO 需要通過 SnapMirror 活動同步路徑傳遞，因此會產生更多延遲。額外的延遲時間取決於用於 SnapMirror 主動同步的站台之間互連的速度。



#### 故障

如果主動鏡像複本因為計畫性或非計畫性容錯移轉而無法使用、則顯然無法再使用。然而、遠端系統已擁有通往遠端站台的同步複本和 SAN 路徑。遠端系統能夠為該 LUN 提供 IO 服務。



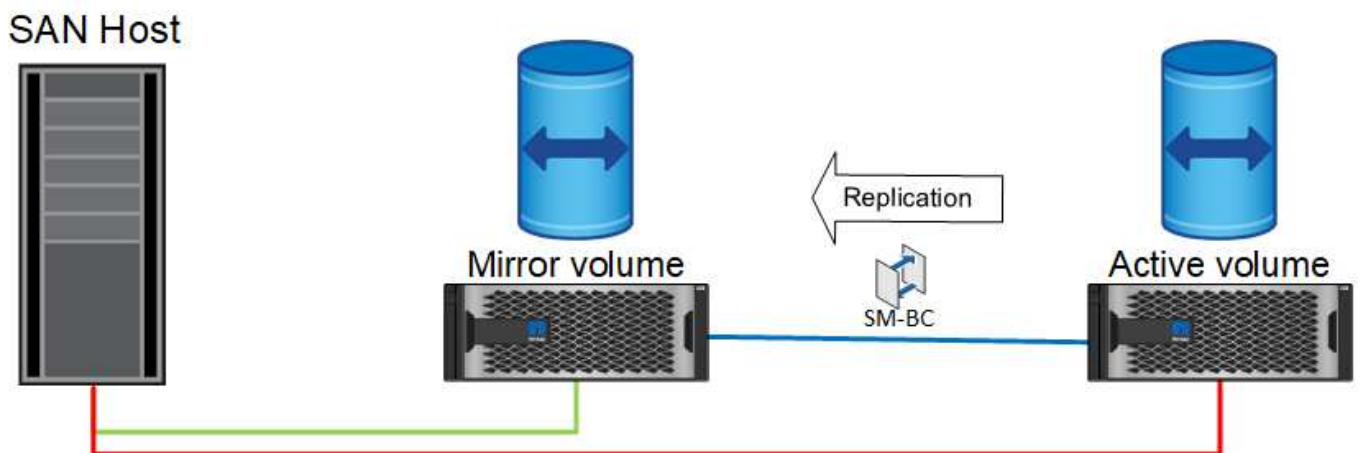
### 容錯移轉

容錯移轉會導致遠端複本變成使用中複本。路徑會從「Active」變更為「Active/Optimized」、IO 也會持續提供服務、不會遺失資料。



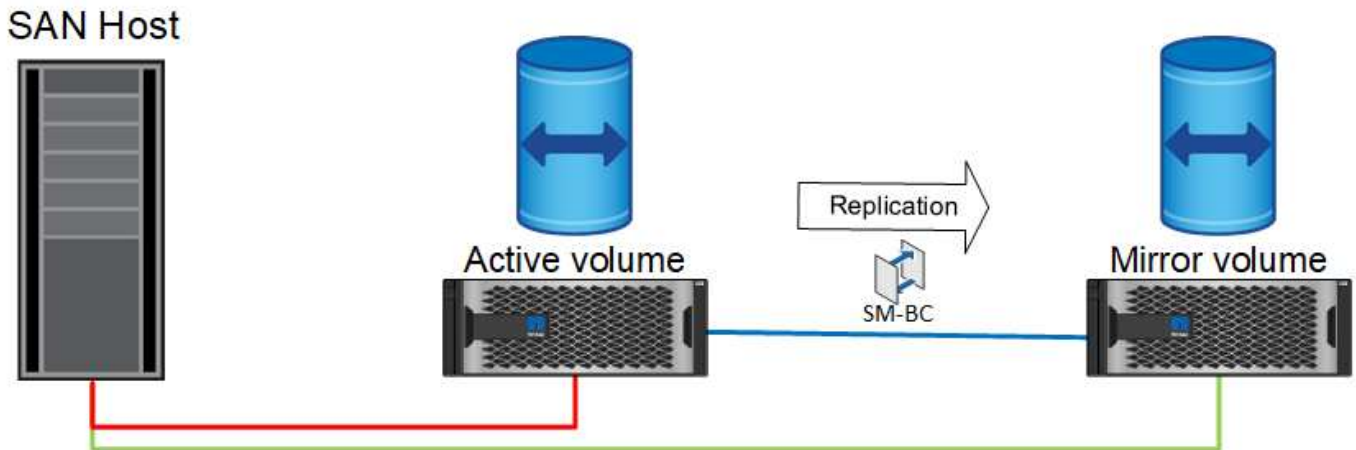
### 修復

一旦來源系統恢復服務、SnapMirror 主動式同步就能重新同步複寫、但會執行其他方向。現在的組態基本上與起點相同、只是主動鏡射站台已經翻轉。



## 容錯回復

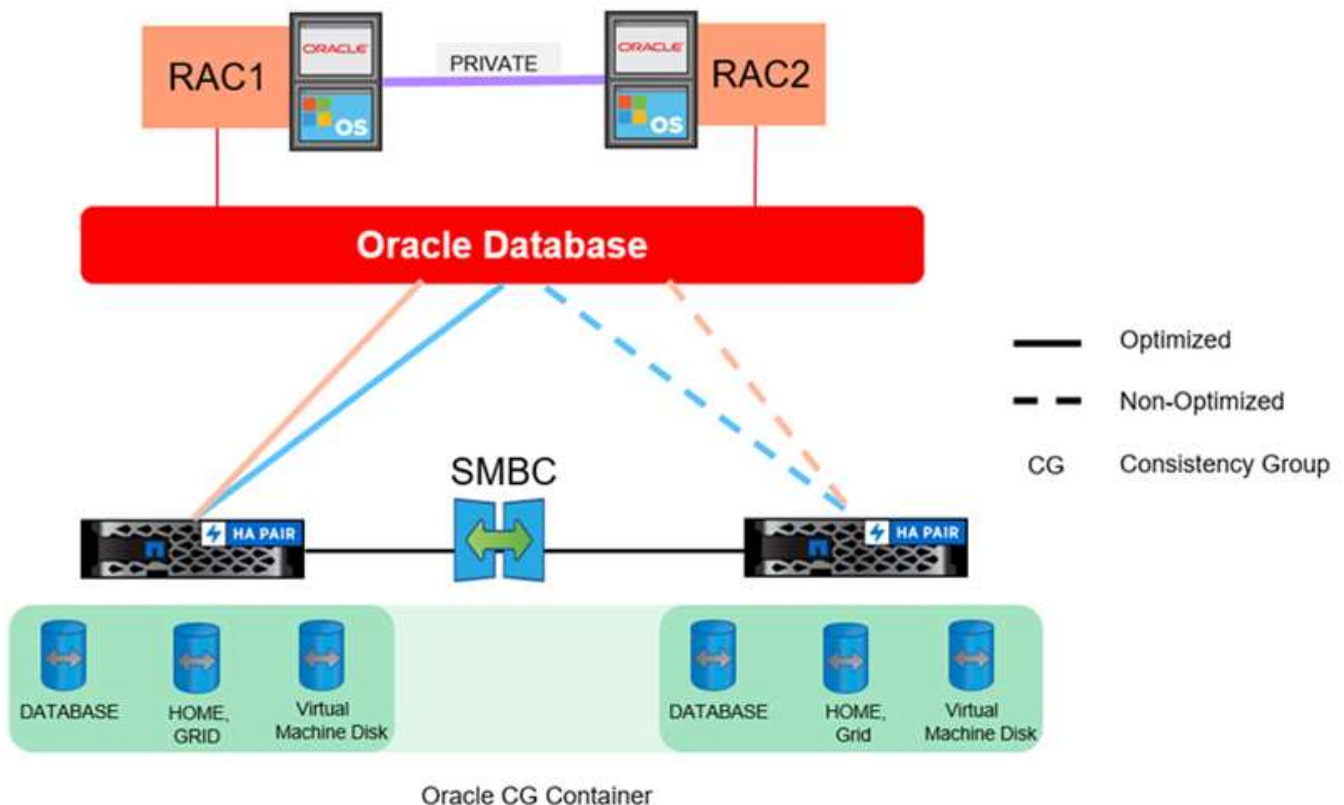
如果需要、管理員可以執行容錯回復、並將 LUN 的作用中複本移回原始控制器。



## 單一執行個體 Oracle 資料庫搭配 SnapMirror 主動式同步

下圖顯示一個簡單的部署模式、您可以將儲存裝置分區或從主要和遠端儲存叢集連線至 Oracle 資料庫。

Oracle 僅在主要系統上設定。此模式可因應儲存端災難時的無縫儲存容錯移轉、在沒有任何應用程式停機的情況下、不會造成資料遺失。然而、此模型無法在站台故障期間提供資料庫環境的高可用度。這類架構對於想要尋找零資料遺失解決方案、且儲存服務可用度高的客戶而言非常實用、但請接受、資料庫叢集的整體遺失需要手動作業。



這種方法也能節省 Oracle 授權成本。在遠端站台上預先設定 Oracle 資料庫節點時、所有核心都必須根據大多數 Oracle 授權合約獲得授權。如果安裝 Oracle 資料庫伺服器 and 掛載仍在運作的資料複本所需的時間所造成的延遲是可以接受的、則此設計可能非常具成本效益。

## Oracle RAC 搭配 SnapMirror 主動式同步

SnapMirror 主動式同步可針對資料集複寫提供精細的控制、例如負載平衡或個別應用程式容錯移轉。整體架構看起來像延伸 RAC 叢集、但有些資料庫是專供特定站台使用、整體負載則分散。

例如、您可以建置一個 Oracle RAC 叢集、主控六個個別資料庫。其中三個資料庫的儲存設備主要託管於站台 A、其他三個資料庫的儲存設備則裝載於站台 B 此組態可將跨網站流量降至最低、以確保最佳效能。此外、應用程式也會設定為使用儲存系統本機的資料庫執行個體、並使用作用中路徑。如此可將 RAC 互連流量降至最低。最後、這項整體設計可確保所有運算資源均能平均使用。隨著工作負載改變、資料庫可以在站台之間選擇性地來回容錯、以確保負載均勻。

除了精細度之外、使用 SnapMirror 主動式 SynCare 的 Oracle RAC 基本原則和選項與相同 "[MetroCluster 上的 Oracle RAC](#)"

## Oracle 資料庫和 SnapMirror 主動式同步失敗案例

有多種 SnapMirror 主動式同步 (SM-AS) 故障案例、每種案例的結果各不相同。

案例	結果
複寫連結失敗	中介程序可辨識這種分割腦案例、並在保留主複本的節點上恢復 I/O。當站台之間的連線恢復上線時、替代站台會執行自動重新同步。
主要站台儲存設備故障	自動非計畫性容錯移轉是由 Mediator 啟動。  無 I/O 中斷。
遠端站台儲存設備故障	沒有 I/O 中斷。由於網路造成同步複寫中斷、主機確認其擁有者是合法的 I/O 服務者 (共識)、因此暫時暫停。因此、I/O 暫停數秒、然後 I/O 就會恢復。  網站上線時會自動重新同步。
Mediator 或 Mediator 與儲存陣列之間的連結遺失	I/O 會繼續並與遠端叢集保持同步、但如果沒有 Mediator、則無法自動進行非計畫性 / 計畫性容錯移轉和容錯回復。
HA 叢集中的其中一個儲存控制器遺失	HA 叢集中的合作夥伴節點會嘗試接管 (n)。如果接管失敗、Mediator 會注意到儲存設備中的兩個節點都已關閉、並自動執行非計畫性容錯移轉至遠端叢集。
磁碟遺失	IO 會持續連續三次發生磁碟故障。這是 RAID-TEC 的一部分。

案例	結果
在一般部署中遺失整個站台	<p>故障站台上的伺服器顯然無法再使用。支援叢集的應用程式可設定為在兩個站台上執行、並在其他站台上繼續作業、不過大多數的應用程式都需要類似於 SM 要求中介者的第三站台斷路器。</p> <p>如果沒有應用程式層級的叢集、應用程式就必須在正常運作的站台上啟動。這會影響可用性、但會保留 RPO =0。不會遺失任何資料。</p>

## Oracle 資料庫移轉

### 將 Oracle 資料庫移轉至 ONTAP 儲存系統

利用新儲存平台的功能、有一項不可避免的需求；資料必須放在新的儲存系統上。ONTAP 讓移轉程序變得簡單、包括 ONTAP 到 ONTAP 的移轉與升級、外部 LUN 匯入、以及直接使用主機作業系統或 Oracle 資料庫軟體的程序。



本文件取代先前發佈的技術報告 [\\_TR-4534](#)：將 Oracle 資料庫移轉至 NetApp 儲存系統

若是新的資料庫專案、這並不是問題、因為資料庫和應用程式環境都已建置就緒。然而、移轉對於業務中斷、完成移轉所需的時間、所需的技能組合、以及將風險降至最低等方面、都會帶來特殊挑戰。

#### 指令碼

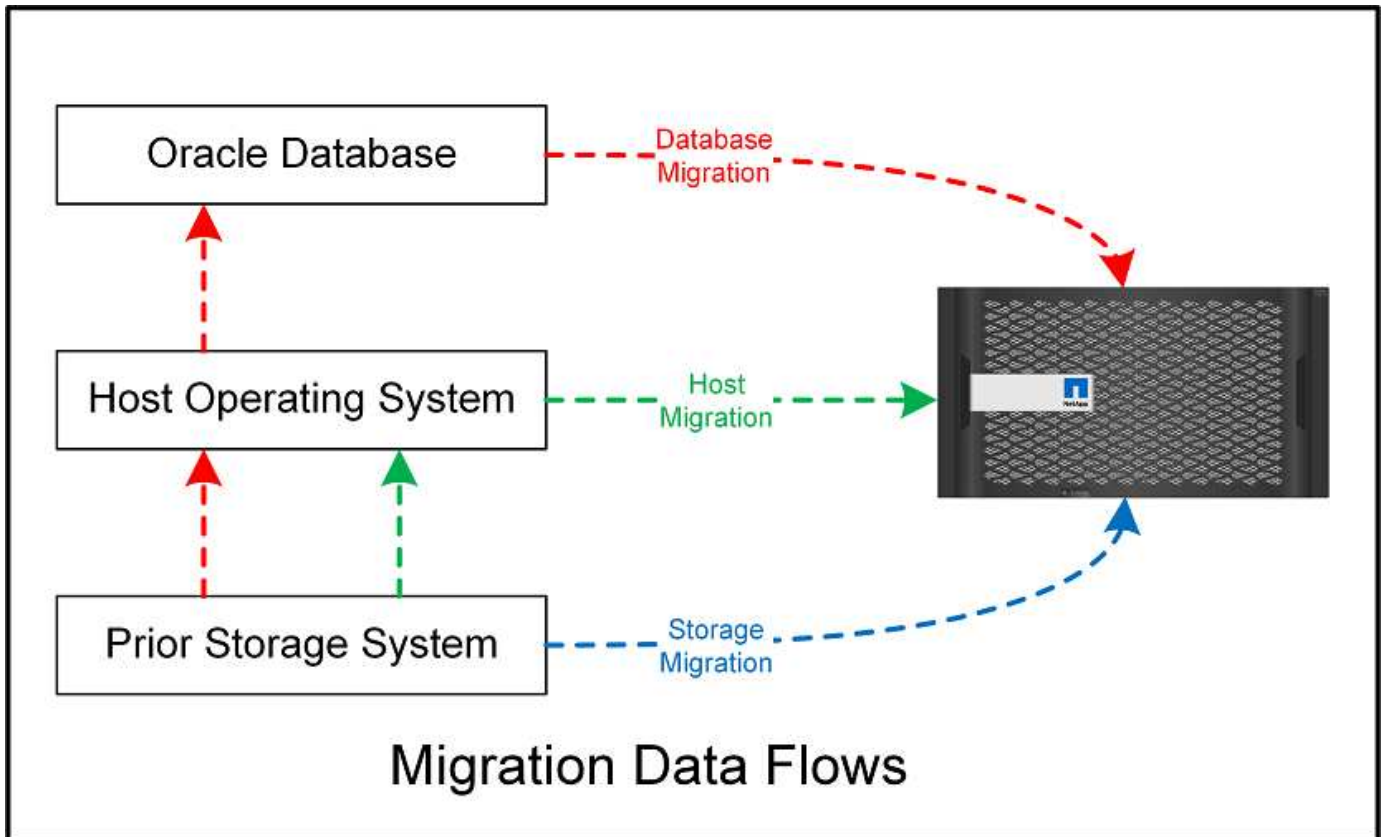
本文件提供範例指令碼。這些指令碼提供自動化移轉各個層面的範例方法、以降低使用者錯誤的機率。這些指令碼可降低 IT 人員對於移轉作業的整體需求、並加速整體程序。這些指令碼都是從 NetApp 專業服務和 NetApp 合作夥伴執行的實際移轉專案中擷取而來。本文件中會顯示使用範例。

### Oracle 資料庫移轉規劃

Oracle 資料移轉可在下列三個層級中進行：資料庫、主機或儲存陣列。

不同之處在於整體解決方案的哪個元件負責移動資料：資料庫、主機作業系統或儲存系統。

下圖顯示移轉層級和資料流的範例。在資料庫層級移轉的情況下、資料會從原始儲存系統透過主機和資料庫層移至新環境。主機層級移轉類似、但資料不會通過應用程式層、而是使用主機程序寫入新位置。最後、隨著儲存層級移轉、NetApp FAS 系統等陣列也會負責資料移動。



資料庫層級移轉通常指透過待命資料庫傳送 Oracle 記錄以完成 Oracle 層級的移轉。主機層級的移轉是使用主機作業系統組態的原生功能來執行。此組態包括檔案複製作業、使用 CP、tar 和 Oracle Recovery Manager (RMAN) 等命令、或使用邏輯 Volume Manager (LVM) 來重新定位檔案系統的基礎位元組。Oracle 自動儲存管理 (ASM) 被歸類為主機層級功能、因為它的執行層級低於資料庫應用程式層級。ASM 取代主機上一般的邏輯磁碟區管理程式。最後、資料可以在儲存陣列層級上移轉、也就是在作業系統層級之下。

#### 規劃考量

最佳移轉選項取決於多種因素、包括要移轉的環境規模、避免停機的需求、以及執行移轉所需的整體工作。大型資料庫顯然需要更多時間和精力來進行移轉、但這類移轉的複雜度極低。小型資料庫可以快速移轉、但如果要移轉數千個資料庫、則工作規模可能會造成複雜問題。最後、資料庫越大、業務關鍵的可能性就越大、因此需要將停機時間降至最低、同時保留一條後端路徑。

此處將討論規劃移轉策略的一些考量事項。

#### 資料大小

要移轉的資料庫大小顯然會影響移轉規劃、但大小不一定會影響轉換時間。當必須移轉大量資料時、主要考量的是頻寬。複製作業通常是以高效率的連續 I/O 來執行保守估計、假設複製作業的可用網路頻寬使用率為 50%。例如、8GB FC 連接埠理論上可傳輸約 800MBps。假設使用率為 50%、則資料庫的複製速度約為 400Mbps。因此、10TB 資料庫可在此速率下在大約七小時內複製。

遠距離移轉通常需要更具創意的辦法、例如中所說明的記錄傳送程序 "[線上資料檔案移動](#)"。遠距 IP 網路在任何接近 LAN 或 SAN 速度的地方、都很少會有頻寬。在某種情況下、NetApp 協助 220 TB 資料庫進行遠距移轉、而且產生的歸檔記錄率非常高。所選的資料傳輸方法是每天運送磁帶、因為這種方法提供最大可能的頻寬。

## 資料庫計數

在許多情況下、移動大量資料的問題並不是資料大小、而是支援資料庫的組態複雜度。只要知道必須移轉 50TB 的資料庫、就無法獲得足夠的資訊。它可以是單一 50TB 關鍵任務資料庫、4、000 個舊資料庫的集合、或是正式作業和非正式作業資料的混合。在某些情況下、大部分資料是由來源資料庫的複本所組成。這些複本完全不需要移轉、因為它們可以輕鬆地重新建立、特別是當新架構設計為使用 NetApp FlexClone Volume 時。

在移轉規劃方面、您必須瞭解範圍內有多少資料庫、以及它們的優先順序。隨著資料庫數量的增加、偏好的移轉選項在堆疊中會較低和較低。例如、在 RMAN 和短暫停機的情況下、複製單一資料庫可能很容易。這是主機層級的複寫。

如果有 50 個資料庫、可能會更容易避免設定新的檔案系統結構來接收 RMAN 複本、而改為將資料移到適當位置。此程序可透過利用主機型 LVM 移轉來將資料從舊 LUN 重新放置到新的 LUN 來完成。這樣做會將責任從資料庫管理員 (DBA) 團隊移轉至作業系統團隊、因此資料會以透明方式移轉至資料庫。檔案系統組態不變。

最後、如果必須移轉 200 部伺服器上的 500 個資料庫、則可使用 ONTAP 外部 LUN 匯入 (FLI) 功能等儲存型選項來執行 LUN 的直接移轉。

## 重新架構需求

一般而言、必須變更資料庫檔案配置才能運用新儲存陣列的功能、但情況並非總是如此。例如、EF 系列 All Flash 陣列的功能主要是針對 SAN 效能和 SAN 可靠性。在大多數情況下、資料庫可以移轉至 EF 系列陣列、而無需特別考量資料配置。唯一的要求是高 IOPS、低延遲和強大的可靠性。雖然 RAID 組態或動態磁碟集區等因素都有最佳實務做法、但 EF 系列專案很少需要對整體儲存架構進行任何重大變更、才能運用這些功能。

相反地、移轉至 ONTAP 通常需要更多考量資料庫配置、以確保最終組態能提供最大價值。ONTAP 本身可為資料庫環境提供許多功能、即使沒有任何特定的架構工作也沒問題。最重要的是、當目前的硬體達到使用壽命時、它能夠不中斷地移轉至新的硬體。一般而言、移轉至 ONTAP 是您最後需要執行的移轉作業。隨後的硬體即會就地升級、資料也不會中斷營運地移轉至新媒體。

有了一些規劃、就能獲得更多效益。使用快照時最重要的考量事項。快照是執行近乎即時的備份、還原及複製作業的基礎。作為快照功能的範例、已知最大的用途是在 6 個控制器上的約 250 個 LUN 上執行單一資料庫 996TB。此資料庫可在 2 分鐘內備份、2 分鐘內還原、15 分鐘內複製完成。其他優點包括：能夠在叢集內移動資料、以因應工作負載的變化、以及應用服務品質 (QoS) 控制來在多資料庫環境中提供良好且一致的效能。

QoS 控制、資料重新配置、快照和複製等技術幾乎可在任何組態中運作。不過、一般需要考慮一些方法才能最大化效益。在某些情況下、資料庫儲存配置可能需要變更設計、以最大化對新儲存陣列的投資。這類設計變更可能會影響移轉策略、因為主機型或儲存型移轉會複寫原始資料配置。完成移轉並提供針對 ONTAP 最佳化的資料配置、可能需要其他步驟。中所示的程序 "[Oracle 移轉程序概述](#)" 稍後、我們將示範一些方法、讓您不只是移轉資料庫、還能以最少的心力將資料庫移轉至最佳的最終配置。

## 轉換時間

應決定轉換期間允許的服務中斷上限。假設整個移轉程序會造成中斷、這是常見的錯誤。許多工作都可以在服務中斷開始之前完成、許多選項都能在不中斷或中斷的情況下完成移轉。即使無法避免中斷、您仍必須定義允許的服務中斷上限、因為轉換時間的持續時間因程序而異。

例如、複製 10TB 資料庫通常需要大約七小時才能完成。如果企業需要中斷七小時、檔案複製是輕鬆安全的移轉選項。如果五小時不可接受、則只需簡單的記錄傳送程序 (請參閱 "[Oracle 記錄傳送](#)") 只需最少的努力就能設定、將轉換時間縮短至約 15 分鐘。在此期間、資料庫管理員可以完成此程序。如果 15 分鐘是不可接受的、則可透過指令碼將最後的轉換程序自動化、將轉換時間縮短至幾分鐘。您可以隨時加快移轉速度、但這樣做的代價是時間和精力。轉換時間目標應以企業可接受的內容為基礎。

## 回溯路徑

沒有移轉作業完全沒有風險。即使技術運作正常、使用者也永遠可能發生錯誤。與所選移轉路徑相關的風險、必須與移轉失敗的後果一併考量。例如、Oracle ASM 的透明線上儲存移轉功能是其重要功能之一、這是最可靠的方法之一。然而、資料正以這種方法進行無法扭轉的複製。在 ASM 發生問題的極不可能發生的事件中、沒有簡單的回傳路徑。唯一的選項是還原原始環境、或使用 ASM 將移轉回復至原始 LUN。如果系統能夠執行此類作業、則可在原始儲存系統上執行快照類型備份、將風險降至最低、但不會消除。

## 排練

某些移轉程序必須在執行前經過完整驗證。移轉和排練轉換程序是關鍵任務資料庫的常見要求、移轉必須成功、停機時間必須降至最低。此外、使用者驗收測試通常是移轉後工作的一部分、只有在完成這些測試之後、才能將整個系統恢復正常運作。

如果需要排練、有幾項 ONTAP 功能可以讓流程更輕鬆。特別是、快照可以重設測試環境、並快速建立資料庫環境的多個具空間效益的複本。

## 程序

### Oracle 移轉程序概述

Oracle 移轉資料庫有許多可用的程序。正確的選擇取決於您的業務需求。

在許多情況下、系統管理員和 DBA 都有自己偏好的方法來重新定位實體磁碟區資料、鏡像和磁碟鏡射、或是利用 Oracle RMAN 來複製資料。

這些程序主要是為不熟悉某些可用選項的 IT 人員提供指引。此外、這些程序還說明每種移轉方法的工作、時間需求和專長類別需求。如此一來、NetApp 和合作夥伴專業服務或 IT 管理等其他方就能更充分瞭解每個程序的要求。

建立移轉策略沒有單一最佳實務做法。建立計畫需要先瞭解可用度選項、然後選擇最符合業務需求的方法。下圖說明客戶所做的基本考量和典型結論、但並不適用於所有情況。

例如、一個步驟會引發資料庫總大小的問題。下一步取決於資料庫是否大於或小於 1TB。建議的步驟就是根據一般客戶實務做法提出建議。大多數客戶不會使用 DataGuard 複製小型資料庫、但有些客戶可能會使用。大多數客戶不會因為所需時間而嘗試複製 50TB 資料庫、但有些客戶可能會有足夠大的維護時間來允許這類作業。

您可以找到最適合移轉路徑的考量類型流程圖 ["請按這裡"](#)。

### 線上資料檔案移動

Oracle 12cR1 及更高版本可在資料庫保持連線時移動資料檔案。此外、它還能在不同的檔案系統類型之間運作。例如、資料檔案可從 xfs 檔案系統重新定位至 ASM。由於需要個別資料檔案移動作業的數量、因此通常不會大規模使用此方法、但這是一個值得考慮的選項、因為較小的資料庫資料檔案數量較少。

此外、單純移動資料檔案是移轉部分現有資料庫的好選項。例如、較不活躍的資料檔案可重新放置到更具成本效益的儲存設備、例如可在物件儲存區中儲存閒置區塊的 FabricPool Volume。

### 資料庫層級移轉

資料庫層級的移轉意味著允許資料庫重新配置資料。具體而言、這表示記錄傳送。RMAN 和 ASM 等技術是 Oracle 產品、但為了進行移轉、它們會在主機層級運作、在主機層級複製檔案和管理磁碟區。



## 記錄傳送

資料庫層級移轉的基礎是 Oracle 歸檔記錄檔、其中包含資料庫變更的記錄檔。歸檔記錄通常是備份與還原策略的一部分。恢復程序從還原資料庫開始、然後重新播放一或多個歸檔記錄檔、將資料庫恢復到所需的狀態。這項相同的基本技術可用於執行移轉、幾乎不會中斷營運。更重要的是、這項技術可在不影響原始資料庫的情況下進行移轉、並保留一條向外移轉的路徑。

移轉程序從將資料庫備份還原至次要伺服器開始。您可以透過多種方式進行、但大多數客戶都會使用正常的備份應用程式來還原資料檔案。資料檔案還原後、使用者便可建立記錄傳送方法。目標是建立由主要資料庫所產生的持續歸檔記錄摘要、並在還原的資料庫上重新播放、使兩者都接近相同的狀態。轉換時間到達時、來源資料庫會完全關閉、最後的歸檔記錄會複製並重新播放、有時則會複製重做記錄。重做記錄也必須列入考量、因為它們可能包含已提交的部分最終交易。

在傳輸和重播這些記錄之後、兩個資料庫彼此之間的一致性。此時、大多數客戶都會執行一些基本測試。如果在移轉過程中發生任何錯誤、則記錄重新執行應會報告錯誤並失敗。還是建議您根據已知的查詢或應用程式導向的活動來執行一些快速測試、以驗證組態是否最佳化。在關閉原始資料庫之前建立一個最終測試表格、以確認該表格是否存在於移轉的資料庫中、這也是常見的做法。此步驟可確保在最終記錄同步期間不會發生任何錯誤。

簡單的記錄傳送移轉作業可針對原始資料庫進行額外設定、這對關鍵任務資料庫特別有用。來源資料庫不需要變更組態、移轉環境的還原和初始組態也不會影響正式作業。設定記錄傳送之後、它會對正式作業伺服器提出一些 I/O 需求。不過、記錄傳送是由簡單的歸檔記錄循序讀取所組成、這對正式作業資料庫效能不太可能有任何影響。

已證實、記錄傳送對於長期、高變更率的移轉專案特別有用。在某個案例中、單一 220TB 資料庫已移轉至距離約 500 英哩的新位置。變更率極高、安全性限制無法使用網路連線。記錄傳送是使用磁帶和快遞業者來執行。原始資料庫的複本最初是使用下列程序來還原。然後、快遞業者每週寄送記錄、直到最後一組磁帶送達時、記錄才會套用至複本資料庫。

## Oracle DataGuard

在某些情況下、保證提供完整的 DataGuard 環境。使用術語 DataGuard 來指稱任何記錄傳送或待命資料庫組態是不正確的。Oracle DataGuard 是管理資料庫複寫的全方位架構、但它不是複寫技術。在移轉工作中、完整的 DataGuard 環境的主要優點是能從一個資料庫透明切換到另一個資料庫。如果發現問題（例如新環境的效能或網路連線問題）、Dataguard 也能將切換回原始資料庫。完整設定的 DataGuard 環境不僅需要設定資料庫層、也需要設定應用程式、以便應用程式能夠偵測主要資料庫位置的變更。一般而言、不需要使用 DataGuard 來完成移轉、但有些客戶內部擁有豐富的 DataGuard 專業知識、而且已經仰賴它來進行移轉工作。

## 重新架構

如前所述、運用儲存陣列的進階功能有時需要變更資料庫配置。此外、從 ASM 移轉至 NFS 檔案系統等儲存傳輸協定變更、也必然會改變檔案系統配置。

記錄傳送方法（包括 DataGuard）的主要優點之一是複寫目的地不需要與來源相符。使用記錄傳送方法從 ASM 移轉至一般檔案系統時沒有任何問題、反之亦然。您可以在目的地變更資料檔案的精確配置、以最佳化易插拔資料庫（PDB）技術的使用、或選擇性地設定特定檔案的 QoS 控制。換句話說、以記錄傳送為基礎的移轉程序可讓您輕鬆安全地最佳化資料庫儲存配置。

## 伺服器資源

資料庫層級移轉的一項限制是需要第二部伺服器。使用第二部伺服器的方法有兩種：

1. 您可以使用第二部伺服器作為資料庫的永久新主目錄。
2. 您可以使用第二部伺服器做為暫存伺服器。在資料移轉至新儲存陣列完成並測試之後、LUN 或 NFS 檔案系

統會從暫存伺服器中斷連線、然後重新連線至原始伺服器。

第一個選項是最簡單的、但在需要非常強大伺服器的大型環境中、使用它可能不可行。第二個選項需要額外的工作、才能將檔案系統重新放置回原始位置。這可以是一項簡單的作業、使用 NFS 做為儲存傳輸協定、因為檔案系統可以從暫存伺服器卸載、然後重新掛載到原始伺服器上。

區塊型檔案系統需要額外的工作、才能更新 FC 分區或 iSCSI 啟動器。對於大多數邏輯磁碟區管理員（包括 ASM）、LUN 會在原始伺服器上可用後自動偵測並上線。不過、某些檔案系統和 LVM 實作可能需要更多工作才能匯出和匯入資料。確切的程序可能會有所不同、但通常很容易建立簡單且可重複的程序、以完成移轉並將資料重新存放在原始伺服器上。

雖然可以在單一伺服器環境中設定記錄傳送和複寫資料庫、但新執行個體必須具有不同的處理程序 SID、才能重新播放記錄。您可以使用不同的 SID、在不同的處理序識別碼集下暫時開啟資料庫、稍後再變更。然而、這樣做可能會導致許多複雜的管理活動、並使資料庫環境面臨使用者錯誤的風險。

#### 主機層級移轉

在主機層級移轉資料是指使用主機作業系統和相關公用程式來完成移轉。此程序包括複製資料的任何公用程式、包括 Oracle RMAN 和 Oracle ASM。

#### 資料複製

不應低估簡單複製作業的價值。現代化的網路基礎架構可以以每秒 GB 的速度來移動資料、而檔案複製作業則是以高效率的連續讀寫 I/O 為基礎相較於記錄傳送、主機複本作業無法避免造成更多中斷、但移轉不只是資料移動而已。通常包括網路變更、資料庫重新啟動時間和移轉後測試。

複製資料所需的實際時間可能並不重要。此外、複製作業會保留保證的回傳路徑、因為原始資料不會受到影響。如果在移轉過程中遇到任何問題、可以重新啟動原始資料的原始檔案系統。

#### 重新建立平台

重組是指 CPU 類型的變更。當資料庫從傳統的 Solaris、AIX 或 HP-UX 平台移轉至 x86 Linux 時、由於 CPU 架構的變更、資料必須重新格式化。SPARC、IA64 和 Power CPU 稱為 Big endian 處理器、而 x86 和 x86\_64 架構則稱為小 endian。因此、Oracle 資料檔案中的某些資料會根據使用中的處理器而有不同的訂購方式。

傳統上、客戶都使用 DataPump 跨平台複寫資料。datapump 是一種公用程式、可建立特殊類型的邏輯資料匯出、以便更快地匯入目的地資料庫。因為它會建立資料的邏輯複本、所以 DataPump 會將處理器位準的相依性留在背後。有些客戶仍使用資料平台來重新建立平台、但 Oracle 11g 提供更快速的選項：跨平台可攜式表格空間。這項進階功能可將資料表空間轉換成不同的 endian 格式。這是一種實體轉型、效能優於 DataPump 匯出、它必須將實體位元組轉換為邏輯資料、然後再轉換回實體位元組。

關於 DataPump 和可攜式資料表空間的完整討論不在 NetApp 文件的範圍之內、但 NetApp 根據我們協助客戶移轉至具有新 CPU 架構的新儲存陣列記錄的經驗、提供一些建議：

- 如果使用 DataPump、則應在測試環境中測量完成移轉所需的時間。客戶有時會對完成移轉所需的時間感到驚訝。這種非預期的額外停機可能會造成中斷。
- 許多客戶誤以為跨平台可攜式資料表空間不需要資料轉換。當使用具有不同序位元組的 CPU 時、會使用 RMAN convert 必須事先對資料檔案執行作業。這不是即時操作。在某些情況下、轉換程序可以透過在不同資料檔案上執行多個執行緒來加速、但無法避免轉換程序。

## 邏輯 Volume Manager 導向的移轉

LVMS 的運作方式是將一組或多個 LUN 拆分為一般稱為擴充的小型單元。然後將擴充集區用作建立邏輯磁碟區的來源、這些邏輯磁碟區基本上是虛擬化的。此虛擬化層以各種方式提供價值：

- 邏輯磁碟區可以使用從多個 LUN 擷取的範圍。在邏輯磁碟區上建立檔案系統時、它可以使用所有 LUN 的完整效能功能。此外、它也能提升磁碟區群組中所有 LUN 的平均載入速度、提供更可預測的效能。
- 您可以新增邏輯磁碟區、並在某些情況下移除範圍、以調整其大小。在邏輯磁碟區上調整檔案系統大小通常不會中斷營運。
- 透過移動基礎範圍、邏輯磁碟區可以不中斷地移轉。

使用 LVM 移轉的運作方式有兩種：移動範圍或鏡射 / 去除範圍。LVM 移轉使用高效率的大型區塊連續 I/O、而且很少會造成任何效能問題。如果這確實是問題、通常有節流 I/O 速率的選項。如此可增加完成移轉所需的時間、同時減輕主機和儲存系統的 I/O 負擔。

### 鏡射與鏡射

某些 Volume 管理程式（例如 AIX LVM）可讓使用者指定每個範圍的複本數量、並控制裝載每個複本的裝置。移轉作業是透過取得現有的邏輯磁碟區、將基礎範圍鏡射到新磁碟區、等待複本同步、然後丟棄舊複本來完成。如果需要返回路徑、可以在放置鏡射複本之前建立原始資料的快照。或者、您也可以強制刪除內含的鏡像複本之前、暫時關閉伺服器以遮罩原始 LUN。這樣做會在資料的原始位置保留可恢復的資料複本。

### 擴展移轉

幾乎所有的 Volume 管理程式都允許移轉擴充、有時也有多個選項。例如、某些 Volume 管理程式可讓管理員將特定邏輯磁碟區的個別擴充區從舊儲存區重新定位到新儲存區。Volume 管理程式（例如 Linux LVM2）提供 `pvmove` 命令、可將指定 LUN 裝置上的所有延伸重新定位至新 LUN。移除舊 LUN 之後、即可將其移除。



作業的主要風險是從組態中移除舊的、未使用的 LUN。變更 FC 分區和移除過時的 LUN 裝置時、必須格外小心。

### Oracle 自動儲存管理

Oracle ASM 是結合邏輯 Volume Manager 與檔案系統的產品。在較高層級、Oracle ASM 會將 LUN 集合起來、分成小的分配單元、並將其呈現為稱為 ASM 磁碟群組的單一磁碟區。ASM 也能透過設定備援層級來鏡射磁碟群組。磁碟區可以是無鏡射（外部備援）、鏡射（正常備援）或三向鏡射（高備援）。設定備援層級時、請務必謹慎、因為建立後無法變更。

ASM 也提供檔案系統功能。雖然檔案系統無法直接從主機看到、但 Oracle 資料庫仍可在 ASM 磁碟群組上建立、移動及刪除檔案與目錄。此外、您也可以使用 `asmcmd` 公用程式來瀏覽結構。

與其他 LVM 實作一樣、Oracle ASM 也會在所有可用 LUN 之間、對每個檔案的 I/O 進行分拆和負載平衡、以最佳化 I/O 效能。其次、基礎擴充可重新定位、以便同時調整 ASM 磁碟群組的大小和移轉。Oracle ASM 會透過重新平衡作業來自動化程序。新的 LUN 會新增至 ASM 磁碟群組、而舊的 LUN 會被丟棄、這會觸發磁碟群組中的磁碟區重新配置及後續刪除已清空的 LUN。此程序是最獲證實的移轉方法之一、而 ASM 提供透明移轉的可靠性、可能是最重要的功能。



由於 Oracle ASM 的鏡射層級是固定的、因此無法搭配鏡射和鏡射移轉方法使用。

## 儲存層級移轉

儲存層級移轉是指在應用程式和作業系統層級以下執行移轉。過去、這有時是指使用專門的裝置來複製網路層級的 LUN、但現在這些功能在 ONTAP 中是原生的。

## SnapMirror

使用 NetApp SnapMirror 資料複製軟體、幾乎可以通用地從 NetApp 系統之間移轉資料庫。此程序包括為要移轉的磁碟區設定鏡射關係、允許它們進行同步處理、然後等待轉換時間。當來源資料庫到達時、即會關閉、執行最後一個鏡像更新、而且鏡像也會中斷。然後、複本磁碟區就可以開始使用、方法是掛載包含的 NFS 檔案系統目錄、或是探索包含的 LUN 並啟動資料庫。

在單一 ONTAP 叢集中重新放置磁碟區並不視為移轉作業、而是例行作業 `volume move` 營運。SnapMirror 用作叢集中的資料複製引擎。此程序完全自動化。當磁碟區的屬性（例如 LUN 對應或 NFS 匯出權限）與磁碟區本身一起移動時、無需執行其他移轉步驟。重新配置不會中斷主機作業。在某些情況下、必須更新網路存取、以確保以最有效率的方式存取新重新部署的資料、但這些工作也不會中斷營運。

## 外部 LUN 匯入 (FLI)

FLI 是一項功能、可讓執行 8.3 或更高版本的 Data ONTAP 系統從另一個儲存陣列移轉現有 LUN。此程序很簡單：ONTAP 系統會分區到現有的儲存陣列、就像是任何其他 SAN 主機一樣。然後 Data ONTAP 控制所需的舊版 LUN、並移轉基礎資料。此外、匯入程序會在資料移轉時使用新 Volume 的效率設定、也就是說、資料可以在移轉過程中內嵌進行壓縮及刪除重複資料。

Data ONTAP 8.3 中首次實作的 FLI 僅允許離線移轉。這是非常快速的傳輸、但仍表示在移轉完成之前、LUN 資料無法使用。線上移轉是在 Data ONTAP 8.3.1 中推出。這類移轉可讓 ONTAP 在傳輸過程中提供 LUN 資料、將中斷情形減至最低。當主機重新分區以透過 ONTAP 使用 LUN 時、會發生短暫的中斷。不過、一旦進行這些變更、資料就會再次存取、並在整個移轉程序中保持可存取的狀態。

讀取 I/O 會透過 ONTAP 代理、直到複製作業完成為止、而寫入 I/O 會同步寫入外部和 ONTAP LUN。這兩個 LUN 複本會以這種方式保持同步、直到系統管理員執行完整的轉換程式來釋放外部 LUN、而不再複製寫入內容。

FLI 的設計可與 FC 搭配使用、但如果您想要變更為 iSCSI、則可在移轉完成後、輕鬆將移轉的 LUN 重新對應為 iSCSI LUN。

FLI 的功能包括自動對齊偵測與調整。在這種情況下、「對齊」一詞是指 LUN 裝置上的分割區。最佳效能需要將 I/O 與 4K 區塊對齊。如果分割區的偏移量不是 4K 的倍數、效能就會受到影響。

第二個對齊層面無法透過調整分割區偏移（檔案系統區塊大小）來修正。例如、ZFS 檔案系統通常預設為 512 位元組的內部區塊大小。其他使用 AIX 的客戶偶爾會建立具有 512 或 1、024 位元組區塊大小的 JFS2 檔案系統。雖然檔案系統可能會與 4K 邊界對齊、但在該檔案系統中建立的檔案不會受到影響、效能也會受到影響。

在此情況下不應使用 FLI。雖然資料在移轉後仍可存取、但結果是檔案系統的效能嚴重限制。一般而言、任何支援 ONTAP 上隨機覆寫工作負載的檔案系統、都應該使用 4K 區塊大小。這主要適用於資料庫資料檔案和 VDI 部署等工作負載。區塊大小可使用相關的主機作業系統命令來識別。

例如、在 AIX 上、可以使用檢視區塊大小 `lsfs -q`。使用 Linux、`xfs_info` 和 `tune2fs` 可用於 `xfs` 和 `ext3/ext4`。與 `zfs`、命令是 `zdb -C`。

控制區塊大小的參數為 `ashift` 而且通常預設值為 9，即  $2^9$  或 512 位元組。為了獲得最佳效能 `ashift` 值必須為 12（ $2^{12}=4K$ ）。此值是在創建 `zpool` 時設置的，不能更改，這意味着使用的數據 `zpool` `ashift` 除 12 個以外、應將資料複製到新建立的 `zPool`、以進行移轉。

Oracle ASM 沒有基本區塊大小。唯一的要求是必須正確對齊 ASM 磁碟所在的磁碟分割區。

## 7-Mode Transition Tool

7-Mode Transition Tool (7MTT) 是一種自動化公用程式、用於將大型 7-Mode 組態移轉至 ONTAP。大多數資料庫客戶發現其他方法都比較容易、部分原因是他們通常會依資料庫來移轉環境資料庫、而非重新配置整個儲存設備佔用空間。此外、資料庫通常只是較大型儲存環境的一部分。因此、資料庫通常會個別移轉、其餘的環境則可以使用 7MTT 進行移轉。

有少數客戶擁有專為複雜資料庫環境設計的儲存系統。這些環境可能包含許多磁碟區、快照和許多組態詳細資料、例如匯出權限、LUN 啟動器群組、使用者權限和輕量型目錄存取傳輸協定組態。在這種情況下、7MTT 的自動化功能可簡化移轉作業。

7MTT 可在下列兩種模式中的其中一種運作：

- 複製型轉換 (CBT) \* 7MTT 搭配 CBT、可從新環境中現有的 7 模式系統設定 SnapMirror 磁碟區。資料同步後、7MTT 會協調轉換程序。
- 複製 - 自由轉換 (CFT) \* 採用 CFT 的 7MTT 是根據現有 7-Mode 磁碟櫃的原位轉換而定。不會複製任何資料、也可以重複使用現有的磁碟櫃。保留現有的資料保護與儲存效率組態。

這兩種選項的主要差異在於：無複製轉換是一種非常有效的方法、其中所有連接至原始 7-Mode HA 配對的磁碟櫃都必須重新放置到新環境中。沒有選項可以移動一部分機櫃。複製型方法可讓選取的磁碟區移動。此外、由於可重新儲存磁碟櫃和轉換中繼資料所需的連結、因此也可能會有較長的轉換時間、且無需複製。根據現場經驗、NetApp 建議允許 1 小時重新配置及重新配置磁碟櫃、15 分鐘至 2 小時的中繼資料轉換時間。

## Oracle 資料檔案移轉

單一命令即可移動個別的 Oracle 資料檔案。

例如、下列命令會將資料檔案 IOPST.dbf 從檔案系統中移出 /oradata2 至檔案系統 /oradata3。

```
SQL> alter database move datafile '/oradata2/NTAP/IOPS002.dbf' to
'/oradata3/NTAP/IOPS002.dbf';
Database altered.
```

使用此方法移動資料檔案可能會很慢、但通常不應產生足夠的 I/O、而會干擾日常資料庫工作負載。相反地、透過 ASM 重新平衡移轉可以更快執行、但卻會在資料移動時降低整體資料庫的速度。

您可以建立測試資料檔、然後移動資料檔案、輕鬆測量移動資料檔案所需的時間。操作所耗用的時間會記錄在 v\$ 工作階段資料中：

```

SQL> set linesize 300;
SQL> select elapsed_seconds||': '||message from v$session_longops;
ELAPSED_SECONDS||': '||MESSAGE
-----
-----
351:Online data file move: data file 8: 22548578304 out of 22548578304
bytes done
SQL> select bytes / 1024 / 1024 /1024 as GB from dba_data_files where
FILE_ID = 8;
        GB
-----
        21

```

在此範例中、移動的檔案是 datafile 8、大小為 21 GB、需要約 6 分鐘才能移轉。所需時間顯然取決於儲存系統、儲存網路的功能、以及移轉時發生的整體資料庫活動。

### 透過記錄傳送進行 Oracle 資料庫移轉

使用記錄傳送進行移轉的目標是在新位置建立原始資料檔案的複本、然後建立將變更傳送到新環境的方法。

一旦建立、記錄傳送和重播就能自動進行、使複本資料庫與來源保持大致同步。例如、cron 工作可排程至 (a) 將最近的記錄複製到新位置、並 (b) 每 15 分鐘重播一次。這樣做可在轉換時將中斷次數降至最低、因為必須重播不超過 15 分鐘的歸檔記錄。

以下程序基本上也是資料庫複製作業。所示邏輯類似於 NetApp SnapManager for Oracle (SMO) 和 NetApp SnapCenter Oracle Plug-in 內的引擎。有些客戶已使用指令碼或 WFA 工作流程中所示的程序來進行自訂的複製作業。雖然此程序比使用 SnapCenter 或 SMO 更為手冊化、但仍可輕鬆撰寫指令碼、而 ONTAP 中的資料管理 API 則可進一步簡化程序。

### 記錄傳送 - 檔案系統至檔案系統

本範例示範如何將名為華夫餅的資料庫從一般檔案系統移轉至位於不同伺服器上的其他一般檔案系統。它也說明 SnapMirror 可用來快速複製資料檔案、但這並不是整體程序不可或缺的一部分。

### 建立資料庫備份

第一步是建立資料庫備份。具體而言、此程序需要一組資料檔案、可用於歸檔記錄重新執行。

### 環境

在此範例中、來源資料庫位於 ONTAP 系統上。建立資料庫備份最簡單的方法是使用快照。將資料庫置於熱備份模式幾秒鐘 snapshot create 在託管資料檔案的磁碟區上執行作業。

```

SQL> alter database begin backup;
Database altered.

```

```
Cluster01::*> snapshot create -vserver vserver1 -volume jfsc1_oradata
hotbackup
Cluster01::*>
```

```
SQL> alter database end backup;
Database altered.
```

結果是磁碟上的快照稱為 hotbackup 包含處於熱備份模式時的資料檔案映像。如果結合適當的歸檔記錄以使資料檔案一致、則此快照中的資料可作為還原或複製的基礎。在這種情況下、它會複寫到新的伺服器。

### 還原至新環境

現在必須在新環境中還原備份。這可以透過多種方式完成、包括 Oracle RMAN、從備份應用程式（如 NetBackup）還原、或是簡單複製置於熱備份模式的資料檔案。

在此範例中、SnapMirror 用於將快照熱備份複寫到新位置。

1. 建立新的磁碟區以接收快照資料。從初始化鏡像 jfsc1\_oradata 至 vol\_oradata。

```
Cluster01::*> volume create -vserver vserver1 -volume vol_oradata
-aggregate data_01 -size 20g -state online -type DP -snapshot-policy
none -policy jfsc3
[Job 833] Job succeeded: Successful
```

```
Cluster01::*> snapmirror initialize -source-path vserver1:jfsc1_oradata
-destination-path vserver1:vol_oradata
Operation is queued: snapmirror initialize of destination
"vserver1:vol_oradata".
Cluster01::*> volume mount -vserver vserver1 -volume vol_oradata
-junction-path /vol_oradata
Cluster01::*>
```

2. 在 SnapMirror 設定狀態後、表示同步已完成、請特別根據所需的快照來更新鏡像。

```
Cluster01::*> snapmirror show -destination-path vserver1:vol_oradata
-fields state
source-path          destination-path      state
-----
vserver1:jfsc1_oradata vserver1:vol_oradata SnapMirrored
```

```
Cluster01::*> snapmirror update -destination-path vserver1:vol_oradata
-source-snapshot hotbackup
Operation is queued: snapmirror update of destination
"vserver1:vol_oradata".
```

3. 您可以透過檢視來驗證同步成功與否 newest-snapshot 鏡射磁碟區上的欄位。

```
Cluster01::*> snapmirror show -destination-path vserver1:vol_oradata
-fields newest-snapshot
source-path          destination-path      newest-snapshot
-----
vserver1:jfsc1_oradata vserver1:vol_oradata hotbackup
```

4. 然後鏡射可能會中斷。

```
Cluster01::> snapmirror break -destination-path vserver1:vol_oradata
Operation succeeded: snapmirror break for destination
"vserver1:vol_oradata".
Cluster01::>
```

5. 掛載新的檔案系統。使用區塊型檔案系統時、精確的程序會因使用中的 LVM 而異。必須設定 FC 分區或 iSCSI 連線。建立與 LUN 的連線後、命令如 Linux pvscan 可能需要探索哪些磁碟區群組或 LUN 需要正確設定、才能讓 ASM 發現。

在此範例中、使用簡單的 NFS 檔案系統。此檔案系統可直接掛載。

```
fas8060-nfs1:/vol_oradata          19922944    1639360    18283584    9%
/oradata
fas8060-nfs1:/vol_logs              9961472     128        9961344     1%
/logs
```

## 建立控制檔建立範本

接下來必須建立控制檔範本。 backup controlfile to trace 命令會建立文字命令以重新建立控制檔。在某些情況下、此功能可用於從備份還原資料庫、而且通常用於執行資料庫複製等工作的指令碼。

1. 以下命令的輸出用於為遷移的數據庫重新創建控制文件。

```
SQL> alter database backup controlfile to trace as '/tmp/waffle.ctrl';
Database altered.
```



## 2. 建立控制檔之後、請將檔案複製到新伺服器。

```
[oracle@jpsc3 tmp]$ scp oracle@jpsc1:/tmp/waffle.ctl /tmp/  
oracle@jpsc1's password:  
waffle.ctl                                100% 5199  
5.1KB/s   00:00
```

### 備份參數檔案

在新環境中也需要一個參數檔。最簡單的方法是從目前的 spfile 或 pfile 建立 pfile。在此範例中、來源資料庫使用的是 spfile。

```
SQL> create pfile='/tmp/waffle.tmp.pfile' from spfile;  
File created.
```

### 建立 oratab 項目

需要建立 oratab 項目、才能正常運作如 oraenv 等公用程式。若要建立 oratab 項目、請完成下列步驟。

```
WAFFLE:/orabin/product/12.1.0/dbhome_1:N
```

### 準備目錄結構

如果所需目錄尚未存在、您必須建立它們、否則資料庫啟動程序會失敗。若要準備目錄結構、請完成下列最低需求。

```
[oracle@jpsc3 ~]$ . oraenv  
ORACLE_SID = [oracle] ? WAFFLE  
The Oracle base has been set to /orabin  
[oracle@jpsc3 ~]$ cd $ORACLE_BASE  
[oracle@jpsc3 orabin]$ cd admin  
[oracle@jpsc3 admin]$ mkdir WAFFLE  
[oracle@jpsc3 admin]$ cd WAFFLE  
[oracle@jpsc3 WAFFLE]$ mkdir adump dpdump pfile scripts xdb_wallet
```

### 參數檔案更新

1. 若要將參數檔複製到新伺服器、請執行下列命令。預設位置為 \$ORACLE\_HOME/dbs 目錄。在這種情況下、pfile 可以放在任何地方。它只是移轉程序中的中間步驟。

```

[oracle@jfsc3 admin]$ scp oracle@jfsc1:/tmp/waffle.tmp.pfile
$ORACLE_HOME/dbs/waffle.tmp.pfile
oracle@jfsc1's password:
waffle.pfile                                100%  916
0.9KB/s   00:00

```

1. 視需要編輯檔案。例如、如果歸檔記錄位置已變更、則必須變更 pfile 以反映新位置。在此範例中、只有控制檔正在重新定位、部分是為了在記錄檔和資料檔案系統之間散佈。

```

[root@jfsc1 tmp]# cat waffle.pfile
WAFFLE.__data_transfer_cache_size=0
WAFFLE.__db_cache_size=507510784
WAFFLE.__java_pool_size=4194304
WAFFLE.__large_pool_size=20971520
WAFFLE.__oracle_base='/orabin'#ORACLE_BASE set from environment
WAFFLE.__pga_aggregate_target=268435456
WAFFLE.__sga_target=805306368
WAFFLE.__shared_io_pool_size=29360128
WAFFLE.__shared_pool_size=234881024
WAFFLE.__streams_pool_size=0
*.audit_file_dest='/orabin/admin/WAFFLE/adump'
*.audit_trail='db'
*.compatible='12.1.0.2.0'
*.control_files='/oradata//WAFFLE/control01.ctl','/oradata//WAFFLE/control02.ctl'
*.control_files='/oradata/WAFFLE/control01.ctl','/logs/WAFFLE/control02.ctl'
*.db_block_size=8192
*.db_domain=''
*.db_name='WAFFLE'
*.diagnostic_dest='/orabin'
*.dispatchers='(PROTOCOL=TCP) (SERVICE=WAFFLEXDB)'
*.log_archive_dest_1='LOCATION=/logs/WAFFLE/arch'
*.log_archive_format='%t_%s_%r.dbf'
*.open_cursors=300
*.pga_aggregate_target=256m
*.processes=300
*.remote_login_passwordfile='EXCLUSIVE'
*.sga_target=768m
*.undo_tablespace='UNDOTBS1'

```

2. 編輯完成後、請根據此 pfile 建立 spfile。

```
SQL> create spfile from pfile='waffle.tmp.pfile';
File created.
```

## 重新建立控制檔

在前一個步驟中、的輸出 `backup controlfile to trace` 已複製到新伺服器。所需輸出的特定部分是 `controlfile recreation` 命令。此資訊可在檔案中標記的區段下找到 Set #1. `NORESETLOGS`。從這條線開始 `create controlfile reuse database` 並應包含這個字 `noresetlogs`。結尾是分號 ( ; ) 字元。

1. 在此範例程序中、檔案會讀取如下內容。

```
CREATE CONTROLFILE REUSE DATABASE "WAFFLE" NORESETLOGS ARCHIVELOG
  MAXLOGFILES 16
  MAXLOGMEMBERS 3
  MAXDATAFILES 100
  MAXINSTANCES 8
  MAXLOGHISTORY 292
LOGFILE
  GROUP 1 '/logs/WAFFLE/redo/redo01.log' SIZE 50M BLOCKSIZE 512,
  GROUP 2 '/logs/WAFFLE/redo/redo02.log' SIZE 50M BLOCKSIZE 512,
  GROUP 3 '/logs/WAFFLE/redo/redo03.log' SIZE 50M BLOCKSIZE 512
-- STANDBY LOGFILE
DATAFILE
  '/oradata/WAFFLE/system01.dbf',
  '/oradata/WAFFLE/sysaux01.dbf',
  '/oradata/WAFFLE/undotbs01.dbf',
  '/oradata/WAFFLE/users01.dbf'
CHARACTER SET WE8MSWIN1252
;
```

2. 視需要編輯此指令碼、以反映各種檔案的新位置。例如、已知可支援高 I/O 的某些資料檔案、可能會重新導向至高效能儲存層上的檔案系統。在其他情況下、這些變更可能純粹是因為系統管理員的理由、例如在專用磁碟區中隔離指定的 PDB 資料檔案。
3. 在此範例中 `DATAFILE stanza` 保持不變、但重做記錄會移至中的新位置 `/redo` 而非與歸檔登入共用空間 `/logs`。

```
CREATE CONTROLFILE REUSE DATABASE "WAFFLE" NORESETLOGS ARCHIVELOG
  MAXLOGFILES 16
  MAXLOGMEMBERS 3
  MAXDATAFILES 100
  MAXINSTANCES 8
  MAXLOGHISTORY 292
LOGFILE
  GROUP 1 '/redo/redo01.log' SIZE 50M BLOCKSIZE 512,
  GROUP 2 '/redo/redo02.log' SIZE 50M BLOCKSIZE 512,
  GROUP 3 '/redo/redo03.log' SIZE 50M BLOCKSIZE 512
-- STANDBY LOGFILE
DATAFILE
  '/oradata/WAFFLE/system01.dbf',
  '/oradata/WAFFLE/sysaux01.dbf',
  '/oradata/WAFFLE/undotbs01.dbf',
  '/oradata/WAFFLE/users01.dbf'
CHARACTER SET WE8MSWIN1252
;
```

```

SQL> startup nomount;
ORACLE instance started.
Total System Global Area  805306368 bytes
Fixed Size                  2929552 bytes
Variable Size               331353200 bytes
Database Buffers           465567744 bytes
Redo Buffers                 5455872 bytes
SQL> CREATE CONTROLFILE REUSE DATABASE "WAFFLE" NORESETLOGS  ARCHIVELOG
 2      MAXLOGFILES 16
 3      MAXLOGMEMBERS 3
 4      MAXDATAFILES 100
 5      MAXINSTANCES 8
 6      MAXLOGHISTORY 292
 7 LOGFILE
 8   GROUP 1 '/redo/redo01.log'  SIZE 50M BLOCKSIZE 512,
 9   GROUP 2 '/redo/redo02.log'  SIZE 50M BLOCKSIZE 512,
10   GROUP 3 '/redo/redo03.log'  SIZE 50M BLOCKSIZE 512
11  -- STANDBY LOGFILE
12  DATAFILE
13   '/oradata/WAFFLE/system01.dbf',
14   '/oradata/WAFFLE/sysaux01.dbf',
15   '/oradata/WAFFLE/undotbs01.dbf',
16   '/oradata/WAFFLE/users01.dbf'
17  CHARACTER SET WE8MSWIN1252
18  ;
Control file created.
SQL>

```

如果有任何檔案放錯位置或參數設定錯誤、就會產生錯誤、指出必須修正的項目。資料庫已掛載、但尚未開啟且無法開啟、因為使用中的資料檔案仍標示為處於熱備份模式。必須先套用歸檔記錄檔、才能使資料庫一致。

#### 初始記錄複寫

為了使資料檔案一致、至少需要執行一項記錄回覆作業。有許多選項可供重播記錄。在某些情況下、原始伺服器上的原始歸檔記錄檔位置可以透過 NFS 共用、而且記錄回覆可以直接完成。在其他情況下、必須複製歸檔記錄。

例如、簡單 scp 作業可將所有目前記錄從來源伺服器複製到移轉伺服器：

```
[oracle@jpsc3 arch]$ scp jpsc1:/logs/WAFFLE/arch/* ./
oracle@jpsc1's password:
1_22_912662036.dbf          100%   47MB
47.0MB/s   00:01
1_23_912662036.dbf          100%   40MB
40.4MB/s   00:00
1_24_912662036.dbf          100%   45MB
45.4MB/s   00:00
1_25_912662036.dbf          100%   41MB
40.9MB/s   00:01
1_26_912662036.dbf          100%   39MB
39.4MB/s   00:00
1_27_912662036.dbf          100%   39MB
38.7MB/s   00:00
1_28_912662036.dbf          100%   40MB
40.1MB/s   00:01
1_29_912662036.dbf          100%   17MB
16.9MB/s   00:00
1_30_912662036.dbf          100%   636KB
636.0KB/s   00:00
```

#### 初始記錄重新播放

檔案在歸檔記錄位置後、可以發出命令來重新播放 `recover database until cancel` 接著是回應 `AUTO` 自動重播所有可用的記錄。

```

SQL> recover database until cancel;
ORA-00279: change 382713 generated at 05/24/2016 09:00:54 needed for
thread 1
ORA-00289: suggestion : /logs/WAFFLE/arch/1_23_912662036.dbf
ORA-00280: change 382713 for thread 1 is in sequence #23
Specify log: {<RET>=suggested | filename | AUTO | CANCEL}
AUTO
ORA-00279: change 405712 generated at 05/24/2016 15:01:05 needed for
thread 1
ORA-00289: suggestion : /logs/WAFFLE/arch/1_24_912662036.dbf
ORA-00280: change 405712 for thread 1 is in sequence #24
ORA-00278: log file '/logs/WAFFLE/arch/1_23_912662036.dbf' no longer
needed for
this recovery
...
ORA-00279: change 713874 generated at 05/26/2016 04:26:43 needed for
thread 1
ORA-00289: suggestion : /logs/WAFFLE/arch/1_31_912662036.dbf
ORA-00280: change 713874 for thread 1 is in sequence #31
ORA-00278: log file '/logs/WAFFLE/arch/1_30_912662036.dbf' no longer
needed for
this recovery
ORA-00308: cannot open archived log '/logs/WAFFLE/arch/1_31_912662036.dbf'
ORA-27037: unable to obtain file status
Linux-x86_64 Error: 2: No such file or directory
Additional information: 3

```

最終的歸檔記錄回覆會回報錯誤、但這是正常現象。記錄會指出這一點 sqlplus 正在尋找特定的記錄檔、但找不到該檔案。原因很可能是記錄檔尚未存在。

如果在複製歸檔記錄之前可以關閉來源資料庫、則此步驟只能執行一次。歸檔記錄會複製並重新播放、然後程序會直接繼續進行轉換程序、以複寫重要的重作記錄。

### 遞增記錄複寫及重新播放

在大多數情況下、移轉作業不會立即執行。移轉程序可能在幾天甚至幾週前完成、這表示記錄必須持續運送至複本資料庫並重新執行。因此、當轉換程式到達時、必須傳輸和重播最少的資料。

這樣做有許多方式可以撰寫指令碼、但其中最受歡迎的方法之一是使用 rsync、這是通用的檔案複寫公用程式。使用此公用程式最安全的方法是將其設定為常駐程式。例如、rsyncd.conf 下列檔案顯示如何建立名為的資源 waffle.arch 使用 Oracle 使用者認證存取、並對應至 /logs/WAFFLE/arch。最重要的是、資源設為唯讀、可讀取正式作業資料、但不變更。

```
[root@jfscl arch]# cat /etc/rsyncd.conf
[waffle.arch]
  uid=oracle
  gid=dba
  path=/logs/WAFFLE/arch
  read only = true
[root@jfscl arch]# rsync --daemon
```

下列命令會將新伺服器的保存檔記錄目的地與 `rsync` 資源同步 `waffle.arch` 在原始伺服器上。◦ `t` 引數 `rsync -potg` 根據時間戳記比較檔案清單、只複製新檔案。此程序提供新伺服器的遞增更新。此命令也可在 `cron` 中排程為定期執行。



```

[oracle@jfsc3 arch]$ rsync -potg --stats --progress jfsc1::waffle.arch/*
/logs/WAFFLE/arch/
1_31_912662036.dbf
    650240 100% 124.02MB/s    0:00:00 (xfer#1, to-check=8/18)
1_32_912662036.dbf
    4873728 100% 110.67MB/s    0:00:00 (xfer#2, to-check=7/18)
1_33_912662036.dbf
    4088832 100%  50.64MB/s    0:00:00 (xfer#3, to-check=6/18)
1_34_912662036.dbf
    8196096 100%  54.66MB/s    0:00:00 (xfer#4, to-check=5/18)
1_35_912662036.dbf
    19376128 100%  57.75MB/s    0:00:00 (xfer#5, to-check=4/18)
1_36_912662036.dbf
     71680 100% 201.15kB/s    0:00:00 (xfer#6, to-check=3/18)
1_37_912662036.dbf
    1144320 100%   3.06MB/s    0:00:00 (xfer#7, to-check=2/18)
1_38_912662036.dbf
    35757568 100%  63.74MB/s    0:00:00 (xfer#8, to-check=1/18)
1_39_912662036.dbf
    984576 100%   1.63MB/s    0:00:00 (xfer#9, to-check=0/18)
Number of files: 18
Number of files transferred: 9
Total file size: 399653376 bytes
Total transferred file size: 75143168 bytes
Literal data: 75143168 bytes
Matched data: 0 bytes
File list size: 474
File list generation time: 0.001 seconds
File list transfer time: 0.000 seconds
Total bytes sent: 204
Total bytes received: 75153219
sent 204 bytes  received 75153219 bytes  150306846.00 bytes/sec
total size is 399653376  speedup is 5.32

```

在收到記錄之後、必須重新播放記錄。前面的範例顯示使用 sqlplus 來手動執行 recover database until cancel，這是一種可以輕鬆自動化的程序。此處顯示的範例使用中所述的指令碼 "重播資料庫上的記錄"。指令碼會接受指定需要重新執行作業之資料庫的引數。如此可在多資料庫移轉作業中使用相同的指令碼。

```
[oracle@jfsc3 logs]$ ./replay.logs.pl WAFFLE
ORACLE_SID = [WAFFLE] ? The Oracle base remains unchanged with value
/orabin
SQL*Plus: Release 12.1.0.2.0 Production on Thu May 26 10:47:16 2016
Copyright (c) 1982, 2014, Oracle. All rights reserved.
Connected to:
Oracle Database 12c Enterprise Edition Release 12.1.0.2.0 - 64bit
Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
SQL> ORA-00279: change 713874 generated at 05/26/2016 04:26:43 needed for
thread 1
ORA-00289: suggestion : /logs/WAFFLE/arch/1_31_912662036.dbf
ORA-00280: change 713874 for thread 1 is in sequence #31
Specify log: {<RET>=suggested | filename | AUTO | CANCEL}
ORA-00279: change 814256 generated at 05/26/2016 04:52:30 needed for
thread 1
ORA-00289: suggestion : /logs/WAFFLE/arch/1_32_912662036.dbf
ORA-00280: change 814256 for thread 1 is in sequence #32
ORA-00278: log file '/logs/WAFFLE/arch/1_31_912662036.dbf' no longer
needed for
this recovery
ORA-00279: change 814780 generated at 05/26/2016 04:53:04 needed for
thread 1
ORA-00289: suggestion : /logs/WAFFLE/arch/1_33_912662036.dbf
ORA-00280: change 814780 for thread 1 is in sequence #33
ORA-00278: log file '/logs/WAFFLE/arch/1_32_912662036.dbf' no longer
needed for
this recovery
...
ORA-00279: change 1120099 generated at 05/26/2016 09:59:21 needed for
thread 1
ORA-00289: suggestion : /logs/WAFFLE/arch/1_40_912662036.dbf
ORA-00280: change 1120099 for thread 1 is in sequence #40
ORA-00278: log file '/logs/WAFFLE/arch/1_39_912662036.dbf' no longer
needed for
this recovery
ORA-00308: cannot open archived log '/logs/WAFFLE/arch/1_40_912662036.dbf'
ORA-27037: unable to obtain file status
Linux-x86_64 Error: 2: No such file or directory
Additional information: 3
SQL> Disconnected from Oracle Database 12c Enterprise Edition Release
12.1.0.2.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
```

## 轉換

當您準備好切換至新環境時、必須執行最後一次同步、其中包括歸檔記錄和重做記錄。如果尚未知道原始的重做記錄位置、可以如下所示識別：

```
SQL> select member from v$logfile;
MEMBER
-----
-----
/logs/WAFFLE/redo/redo01.log
/logs/WAFFLE/redo/redo02.log
/logs/WAFFLE/redo/redo03.log
```

1. 關閉來源資料庫。
2. 使用所需的方法、在新伺服器上執行歸檔記錄的最後一次同步。
3. 來源重做記錄檔必須複製到新伺服器。在此範例中、重做記錄會重新定位到新的目錄 /redo。

```
[oracle@jpsc3 logs]$ scp jpsc1:/logs/WAFFLE/redo/* /redo/
oracle@jpsc1's password:
redo01.log
100% 50MB 50.0MB/s 00:01
redo02.log
100% 50MB 50.0MB/s 00:00
redo03.log
100% 50MB 50.0MB/s 00:00
```

4. 在此階段、新的資料庫環境包含所有必要的檔案、使其與來源完全相同。歸檔記錄必須最後重播一次。

```

SQL> recover database until cancel;
ORA-00279: change 1120099 generated at 05/26/2016 09:59:21 needed for
thread 1
ORA-00289: suggestion : /logs/WAFFLE/arch/1_40_912662036.dbf
ORA-00280: change 1120099 for thread 1 is in sequence #40
Specify log: {<RET>=suggested | filename | AUTO | CANCEL}
AUTO
ORA-00308: cannot open archived log
'/logs/WAFFLE/arch/1_40_912662036.dbf'
ORA-27037: unable to obtain file status
Linux-x86_64 Error: 2: No such file or directory
Additional information: 3
ORA-00308: cannot open archived log
'/logs/WAFFLE/arch/1_40_912662036.dbf'
ORA-27037: unable to obtain file status
Linux-x86_64 Error: 2: No such file or directory
Additional information: 3

```

5. 完成後、必須重新執行重作記錄。如果出現此訊息 Media recovery complete 會傳回、程序成功、資料庫會同步、並可開啟。

```

SQL> recover database;
Media recovery complete.
SQL> alter database open;
Database altered.

```

#### 記錄傳送 - ASM 至檔案系統

本範例說明如何使用 Oracle RMAN 移轉資料庫。這與先前的檔案系統傳送檔案系統記錄檔範例非常類似、但主機看不到 ASM 上的檔案。唯一用於移轉位於 ASM 裝置上的資料的選項是重新放置 ASM LUN、或使用 Oracle RMAN 來執行複製作業。

雖然 RMAN 是從 Oracle ASM 複製檔案的必要條件、但 RMAN 的使用不限於 ASM。RMAN 可用於從任何類型的儲存設備移轉至任何其他類型。

此範例顯示將名為 pake 的資料庫從 ASM 儲存設備重新放置到位於路徑上不同伺服器上的一般檔案系統 /oradata 和 /logs。

#### 建立資料庫備份

第一步是建立要移轉到替代伺服器的資料庫備份。由於來源使用 Oracle ASM、因此必須使用 RMAN。簡單的 RMAN 備份可執行如下。此方法會建立標記備份、可在稍後的程序中由 RMAN 輕鬆識別。

第一個命令定義備份的目的地類型和要使用的位置。第二個只會啟動資料檔案的備份。

```

RMAN> configure channel device type disk format '/rman/pancake/%U';
using target database control file instead of recovery catalog
old RMAN configuration parameters:
CONFIGURE CHANNEL DEVICE TYPE DISK FORMAT    '/rman/pancake/%U';
new RMAN configuration parameters:
CONFIGURE CHANNEL DEVICE TYPE DISK FORMAT    '/rman/pancake/%U';
new RMAN configuration parameters are successfully stored
RMAN> backup database tag 'ONTAP_MIGRATION';
Starting backup at 24-MAY-16
allocated channel: ORA_DISK_1
channel ORA_DISK_1: SID=251 device type=DISK
channel ORA_DISK_1: starting full datafile backup set
channel ORA_DISK_1: specifying datafile(s) in backup set
input datafile file number=00001 name=+ASM0/PANCAKE/system01.dbf
input datafile file number=00002 name=+ASM0/PANCAKE/sysaux01.dbf
input datafile file number=00003 name=+ASM0/PANCAKE/undotbs101.dbf
input datafile file number=00004 name=+ASM0/PANCAKE/users01.dbf
channel ORA_DISK_1: starting piece 1 at 24-MAY-16
channel ORA_DISK_1: finished piece 1 at 24-MAY-16
piece handle=/rman/pancake/lgr6c161_1_1 tag=ONTAP_MIGRATION comment=NONE
channel ORA_DISK_1: backup set complete, elapsed time: 00:00:03
channel ORA_DISK_1: starting full datafile backup set
channel ORA_DISK_1: specifying datafile(s) in backup set
including current control file in backup set
including current SPFILE in backup set
channel ORA_DISK_1: starting piece 1 at 24-MAY-16
channel ORA_DISK_1: finished piece 1 at 24-MAY-16
piece handle=/rman/pancake/lhr6c164_1_1 tag=ONTAP_MIGRATION comment=NONE
channel ORA_DISK_1: backup set complete, elapsed time: 00:00:01
Finished backup at 24-MAY-16

```

## 備份控制檔

稍後的程序中需要備份控制檔 duplicate database 營運。

```
RMAN> backup current controlfile format '/rman/pancake/ctrl.bkp';
Starting backup at 24-MAY-16
using channel ORA_DISK_1
channel ORA_DISK_1: starting full datafile backup set
channel ORA_DISK_1: specifying datafile(s) in backup set
including current control file in backup set
channel ORA_DISK_1: starting piece 1 at 24-MAY-16
channel ORA_DISK_1: finished piece 1 at 24-MAY-16
piece handle=/rman/pancake/ctrl.bkp tag=TAG20160524T032651 comment=NONE
channel ORA_DISK_1: backup set complete, elapsed time: 00:00:01
Finished backup at 24-MAY-16
```

### 備份參數檔案

在新環境中也需要一個參數檔。最簡單的方法是從目前的 spfile 或 pfile 建立 pfile。在此範例中、來源資料庫使用 spfile。

```
RMAN> create pfile='/rman/pancake/pfile' from spfile;
Statement processed
```

### ASM 檔案重新命名指令碼

移動資料庫時、控制檔中目前定義的數個檔案位置會變更。下列指令碼會建立 RMAN 指令碼、以簡化程序。此範例顯示的資料庫資料檔案數量極少、但資料庫通常包含數百個甚至數千個資料檔案。

此指令碼位於 ["ASM 至檔案系統名稱轉換"](#) 它有兩件事。

首先、它會建立一個參數、重新定義稱為的重做記錄位置 `log_file_name_convert`。基本上是交替欄位清單。第一個欄位是目前重做記錄檔的位置、第二個欄位是新伺服器上的位置。然後重複該模式。

第二個功能是提供資料檔案重新命名的範本。指令碼會循環瀏覽資料檔案、擷取名稱和檔案編號資訊、並將其格式化為 RMAN 指令碼。然後、它會對暫存檔案執行相同的操作。結果是一個簡單的 RMAN 指令碼、可視需要加以編輯、以確保檔案還原至所需的位置。

```

SQL> @/rman/mk.rename.scripts.sql
Parameters for log file conversion:
*.log_file_name_convert = '+ASM0/PANCAKE/redo01.log',
'/NEW_PATH/redo01.log','+ASM0/PANCAKE/redo02.log',
'/NEW_PATH/redo02.log','+ASM0/PANCAKE/redo03.log', '/NEW_PATH/redo03.log'
rman duplication script:
run
{
set newname for datafile 1 to '+ASM0/PANCAKE/system01.dbf';
set newname for datafile 2 to '+ASM0/PANCAKE/sysaux01.dbf';
set newname for datafile 3 to '+ASM0/PANCAKE/undotbs101.dbf';
set newname for datafile 4 to '+ASM0/PANCAKE/users01.dbf';
set newname for tempfile 1 to '+ASM0/PANCAKE/temp01.dbf';
duplicate target database for standby backup location INSERT_PATH_HERE;
}
PL/SQL procedure successfully completed.

```

擷取此畫面的輸出。◦ `log_file_name_convert` 參數會如下所述放置在 `pfile` 中。RMAN 資料檔案重新命名和重複指令碼必須據此編輯、才能將資料檔案放置在所需的位置。在此範例中、所有的項目都放在 `/oradata/pancake` ◦

```

run
{
set newname for datafile 1 to '/oradata/pancake/pancake.dbf';
set newname for datafile 2 to '/oradata/pancake/sysaux.dbf';
set newname for datafile 3 to '/oradata/pancake/undotbs1.dbf';
set newname for datafile 4 to '/oradata/pancake/users.dbf';
set newname for tempfile 1 to '/oradata/pancake/temp.dbf';
duplicate target database for standby backup location '/rman/pancake';
}

```

### 準備目錄結構

指令碼幾乎可以執行、但首先必須有目錄結構。如果所需目錄尚未存在、則必須建立這些目錄、否則資料庫啟動程序會失敗。以下範例反映最低需求。

```

[oracle@jfsc2 ~]$ mkdir /oradata/pancake
[oracle@jfsc2 ~]$ mkdir /logs/pancake
[oracle@jfsc2 ~]$ cd /orabin/admin
[oracle@jfsc2 admin]$ mkdir PANCAKE
[oracle@jfsc2 admin]$ cd PANCAKE
[oracle@jfsc2 PANCAKE]$ mkdir adump dpdump pfile scripts xdb_wallet

```

## 建立 oratab 項目

下列命令是 oraenv 等公用程式正常運作所需的命令。

```
PANCAKE:/orabin/product/12.1.0/dbhome_1:N
```

## 參數更新

必須更新儲存的 pfile、以反映新伺服器上的任何路徑變更。資料檔案路徑變更是由 RMAN 複製指令碼所變更、幾乎所有資料庫都需要變更 control\_files 和 log\_archive\_dest 參數。也可能有必須變更的稽核檔案位置和參數、例如 db\_create\_file\_dest 在 ASM 之外可能無關緊要。經驗豐富的 DBA 應仔細審查建議的變更、然後再繼續。

在此範例中、主要變更為控制檔位置、記錄歸檔目的地、以及新增 log\_file\_name\_convert 參數。



```

PANCAKE.__data_transfer_cache_size=0
PANCAKE.__db_cache_size=545259520
PANCAKE.__java_pool_size=4194304
PANCAKE.__large_pool_size=25165824
PANCAKE.__oracle_base='/orabin'#ORACLE_BASE set from environment
PANCAKE.__pga_aggregate_target=268435456
PANCAKE.__sga_target=805306368
PANCAKE.__shared_io_pool_size=29360128
PANCAKE.__shared_pool_size=192937984
PANCAKE.__streams_pool_size=0
*.audit_file_dest='/orabin/admin/PANCAKE/adump'
*.audit_trail='db'
*.compatible='12.1.0.2.0'
*.control_files='+ASM0/PANCAKE/control01.ctl','+ASM0/PANCAKE/control02.ctl'
*.control_files='/oradata/pancake/control01.ctl','/logs/pancake/control02.ctl'
*.db_block_size=8192
*.db_domain=''
*.db_name='PANCAKE'
*.diagnostic_dest='/orabin'
*.dispatchers='(PROTOCOL=TCP) (SERVICE=PANCAKEXDB)'
*.log_archive_dest_1='LOCATION=+ASM1'
*.log_archive_dest_1='LOCATION=/logs/pancake'
*.log_archive_format='%t_%s_%r.dbf'
'/logs/path/redo02.log'
*.log_file_name_convert = '+ASM0/PANCAKE/redo01.log',
'/logs/pancake/redo01.log', '+ASM0/PANCAKE/redo02.log',
'/logs/pancake/redo02.log', '+ASM0/PANCAKE/redo03.log',
'/logs/pancake/redo03.log'
*.open_cursors=300
*.pga_aggregate_target=256m
*.processes=300
*.remote_login_passwordfile='EXCLUSIVE'
*.sga_target=768m
*.undo_tablespace='UNDOTBS1'

```

確認新參數之後、必須使參數生效。存在多個選項、但大多數客戶會根據文字 pfile 建立 spfile。

```
bash-4.1$ sqlplus / as sysdba
SQL*Plus: Release 12.1.0.2.0 Production on Fri Jan 8 11:17:40 2016
Copyright (c) 1982, 2014, Oracle. All rights reserved.
Connected to an idle instance.
SQL> create spfile from pfile='/rman/pancake/pfile';
File created.
```

## 啟動 nomount

複寫資料庫之前的最後一個步驟是啟動資料庫程序、但不要掛載檔案。在此步驟中、spfile 可能會出現問題。如果是 startup nomount 命令因參數錯誤而失敗、關機很簡單、請修正 pfile 範本、將其重新載入為 spfile、然後再試一次。

```
SQL> startup nomount;
ORACLE instance started.
Total System Global Area 805306368 bytes
Fixed Size 2929552 bytes
Variable Size 373296240 bytes
Database Buffers 423624704 bytes
Redo Buffers 5455872 bytes
```

## 複製資料庫

將先前的 RMAN 備份還原至新位置、比此程序中的其他步驟花費更多時間。必須複製資料庫、而不需變更資料庫 ID (DBID) 或重新設定記錄。這可防止套用記錄、這是完全同步複本的必要步驟。

使用 RMAN AS aux 連線至資料庫、並使用在前一個步驟中建立的指令碼發出重複資料庫命令。

```
[oracle@jpsc2 pancake]$ rman auxiliary /
Recovery Manager: Release 12.1.0.2.0 - Production on Tue May 24 03:04:56
2016
Copyright (c) 1982, 2014, Oracle and/or its affiliates. All rights
reserved.
connected to auxiliary database: PANCAKE (not mounted)
RMAN> run
2> {
3> set newname for datafile 1 to '/oradata/pancake/pancake.dbf';
4> set newname for datafile 2 to '/oradata/pancake/sysaux.dbf';
5> set newname for datafile 3 to '/oradata/pancake/undotbs1.dbf';
6> set newname for datafile 4 to '/oradata/pancake/users.dbf';
7> set newname for tempfile 1 to '/oradata/pancake/temp.dbf';
8> duplicate target database for standby backup location '/rman/pancake';
9> }
executing command: SET NEWNAME
```

```

executing command: SET NEWNAME
executing command: SET NEWNAME
executing command: SET NEWNAME
executing command: SET NEWNAME
Starting Duplicate Db at 24-MAY-16
contents of Memory Script:
{
  restore clone standby controlfile from  '/rman/pancake/ctrl.bkp';
}
executing Memory Script
Starting restore at 24-MAY-16
allocated channel: ORA_AUX_DISK_1
channel ORA_AUX_DISK_1: SID=243 device type=DISK
channel ORA_AUX_DISK_1: restoring control file
channel ORA_AUX_DISK_1: restore complete, elapsed time: 00:00:01
output file name=/oradata/pancake/control01.ctl
output file name=/logs/pancake/control02.ctl
Finished restore at 24-MAY-16
contents of Memory Script:
{
  sql clone 'alter database mount standby database';
}
executing Memory Script
sql statement: alter database mount standby database
released channel: ORA_AUX_DISK_1
allocated channel: ORA_AUX_DISK_1
channel ORA_AUX_DISK_1: SID=243 device type=DISK
contents of Memory Script:
{
  set newname for tempfile  1 to
"/oradata/pancake/temp.dbf";
  switch clone tempfile all;
  set newname for datafile  1 to
"/oradata/pancake/pancake.dbf";
  set newname for datafile  2 to
"/oradata/pancake/sysaux.dbf";
  set newname for datafile  3 to
"/oradata/pancake/undotbs1.dbf";
  set newname for datafile  4 to
"/oradata/pancake/users.dbf";
  restore
  clone database
  ;
}
executing Memory Script
executing command: SET NEWNAME

```

```

renamed tempfile 1 to /oradata/pancake/temp.dbf in control file
executing command: SET NEWNAME
executing command: SET NEWNAME
executing command: SET NEWNAME
executing command: SET NEWNAME
Starting restore at 24-MAY-16
using channel ORA_AUX_DISK_1
channel ORA_AUX_DISK_1: starting datafile backup set restore
channel ORA_AUX_DISK_1: specifying datafile(s) to restore from backup set
channel ORA_AUX_DISK_1: restoring datafile 00001 to
/oradata/pancake/pancake.dbf
channel ORA_AUX_DISK_1: restoring datafile 00002 to
/oradata/pancake/sysaux.dbf
channel ORA_AUX_DISK_1: restoring datafile 00003 to
/oradata/pancake/undotbs1.dbf
channel ORA_AUX_DISK_1: restoring datafile 00004 to
/oradata/pancake/users.dbf
channel ORA_AUX_DISK_1: reading from backup piece
/rman/pancake/1gr6c161_1_1
channel ORA_AUX_DISK_1: piece handle=/rman/pancake/1gr6c161_1_1
tag=ONTAP_MIGRATION
channel ORA_AUX_DISK_1: restored backup piece 1
channel ORA_AUX_DISK_1: restore complete, elapsed time: 00:00:07
Finished restore at 24-MAY-16
contents of Memory Script:
{
  switch clone datafile all;
}
executing Memory Script
datafile 1 switched to datafile copy
input datafile copy RECID=5 STAMP=912655725 file
name=/oradata/pancake/pancake.dbf
datafile 2 switched to datafile copy
input datafile copy RECID=6 STAMP=912655725 file
name=/oradata/pancake/sysaux.dbf
datafile 3 switched to datafile copy
input datafile copy RECID=7 STAMP=912655725 file
name=/oradata/pancake/undotbs1.dbf
datafile 4 switched to datafile copy
input datafile copy RECID=8 STAMP=912655725 file
name=/oradata/pancake/users.dbf
Finished Duplicate Db at 24-MAY-16

```

## 初始記錄複寫

您現在必須將變更從來源資料庫傳送至新位置。這樣做可能需要多個步驟的組合。最簡單的方法是讓來源資料庫

上的 RMAN 將歸檔記錄寫入共用網路連線。如果無法使用共用位置、則另一種方法是使用 RMAN 寫入本機檔案系統、然後使用 rcp 或 rsync 複製檔案。

在此範例中 /rman 目錄是一種 NFS 共用、可同時用於原始和移轉的資料庫。

此處的一個重要問題是 disk format 條款。備份的磁碟格式為 %h\_%e\_%a.dbf，這表示您必須使用資料庫的執行緒編號、序號和啟動 ID 格式。雖然字母不同、但這與相符 log\_archive\_format='%t %s %r.dbf' pfile 中的參數。此參數也會以執行緒編號、序號和啟動 ID 的格式來指定封存記錄。最終結果是來源上的記錄檔備份使用資料庫預期的命名慣例。如此一來、就能執行像這樣的作業 recover database 更簡單、因為 sqlplus 能正確預測要重新播放的歸檔記錄名稱。

```

RMAN> configure channel device type disk format
'/rman/pancake/logship/%h_%e_%a.dbf';
old RMAN configuration parameters:
CONFIGURE CHANNEL DEVICE TYPE DISK FORMAT
'/rman/pancake/arch/%h_%e_%a.dbf';
new RMAN configuration parameters:
CONFIGURE CHANNEL DEVICE TYPE DISK FORMAT
'/rman/pancake/logship/%h_%e_%a.dbf';
new RMAN configuration parameters are successfully stored
released channel: ORA_DISK_1
RMAN> backup as copy archivelog from time 'sysdate-2';
Starting backup at 24-MAY-16
current log archived
allocated channel: ORA_DISK_1
channel ORA_DISK_1: SID=373 device type=DISK
channel ORA_DISK_1: starting archived log copy
input archived log thread=1 sequence=54 RECID=70 STAMP=912658508
output file name=/rman/pancake/logship/1_54_912576125.dbf RECID=123
STAMP=912659482
channel ORA_DISK_1: archived log copy complete, elapsed time: 00:00:01
channel ORA_DISK_1: starting archived log copy
input archived log thread=1 sequence=41 RECID=29 STAMP=912654101
output file name=/rman/pancake/logship/1_41_912576125.dbf RECID=124
STAMP=912659483
channel ORA_DISK_1: archived log copy complete, elapsed time: 00:00:01
...
channel ORA_DISK_1: starting archived log copy
input archived log thread=1 sequence=45 RECID=33 STAMP=912654688
output file name=/rman/pancake/logship/1_45_912576125.dbf RECID=152
STAMP=912659514
channel ORA_DISK_1: archived log copy complete, elapsed time: 00:00:01
channel ORA_DISK_1: starting archived log copy
input archived log thread=1 sequence=47 RECID=36 STAMP=912654809
output file name=/rman/pancake/logship/1_47_912576125.dbf RECID=153
STAMP=912659515
channel ORA_DISK_1: archived log copy complete, elapsed time: 00:00:01
Finished backup at 24-MAY-16

```

## 初始記錄重新播放

檔案在歸檔記錄位置後、可以發出命令來重新播放 `recover database until cancel` 接著是回應 `AUTO` 自動重播所有可用的記錄。參數檔目前正在將歸檔記錄導向 `/logs/archive` 但這與 `RMAN` 用於保存日誌的位置不匹配。在恢復資料庫之前、可依下列方式暫時重新導向位置。

```

SQL> alter system set log_archive_dest_1='LOCATION=/rman/pancake/logship'
scope=memory;
System altered.
SQL> recover standby database until cancel;
ORA-00279: change 560224 generated at 05/24/2016 03:25:53 needed for
thread 1
ORA-00289: suggestion : /rman/pancake/logship/1_49_912576125.dbf
ORA-00280: change 560224 for thread 1 is in sequence #49
Specify log: {<RET>=suggested | filename | AUTO | CANCEL}
AUTO
ORA-00279: change 560353 generated at 05/24/2016 03:29:17 needed for
thread 1
ORA-00289: suggestion : /rman/pancake/logship/1_50_912576125.dbf
ORA-00280: change 560353 for thread 1 is in sequence #50
ORA-00278: log file '/rman/pancake/logship/1_49_912576125.dbf' no longer
needed
for this recovery
...
ORA-00279: change 560591 generated at 05/24/2016 03:33:56 needed for
thread 1
ORA-00289: suggestion : /rman/pancake/logship/1_54_912576125.dbf
ORA-00280: change 560591 for thread 1 is in sequence #54
ORA-00278: log file '/rman/pancake/logship/1_53_912576125.dbf' no longer
needed
for this recovery
ORA-00308: cannot open archived log
'/rman/pancake/logship/1_54_912576125.dbf'
ORA-27037: unable to obtain file status
Linux-x86_64 Error: 2: No such file or directory
Additional information: 3

```

最終的歸檔記錄回覆會回報錯誤、但這是正常現象。此錯誤表示 sqlplus 正在尋找特定的記錄檔、但找不到該檔案。原因很可能是記錄檔尚未存在。

如果在複製歸檔記錄之前可以關閉來源資料庫、則此步驟只能執行一次。歸檔記錄會複製並重新播放、然後程序會直接繼續進行轉換程序、以複寫重要的重作記錄。

#### 遞增記錄複寫及重新播放

在大多數情況下、移轉作業不會立即執行。移轉程序可能在幾天甚至幾週前完成、這表示記錄必須持續運送至複本資料庫並重新執行。這樣做可確保轉換程序到達時、必須傳輸和重播最少的資料。

此程序很容易撰寫指令碼。例如、您可以在原始資料庫上排程下列命令、以確保用於記錄傳送的位置持續更新。

```
[oracle@jfscl pancake]$ cat copylogs.rman
configure channel device type disk format
'/rman/pancake/logship/%h_%e_%a.dbf';
backup as copy archivelog from time 'sysdate-2';
```

```
[oracle@jfscl pancake]$ rman target / cmdfile=copylogs.rman
Recovery Manager: Release 12.1.0.2.0 - Production on Tue May 24 04:36:19
2016
Copyright (c) 1982, 2014, Oracle and/or its affiliates. All rights
reserved.
connected to target database: PANCAKE (DBID=3574534589)
RMAN> configure channel device type disk format
'/rman/pancake/logship/%h_%e_%a.dbf';
2> backup as copy archivelog from time 'sysdate-2';
3>
4>
using target database control file instead of recovery catalog
old RMAN configuration parameters:
CONFIGURE CHANNEL DEVICE TYPE DISK FORMAT
'/rman/pancake/logship/%h_%e_%a.dbf';
new RMAN configuration parameters:
CONFIGURE CHANNEL DEVICE TYPE DISK FORMAT
'/rman/pancake/logship/%h_%e_%a.dbf';
new RMAN configuration parameters are successfully stored
Starting backup at 24-MAY-16
current log archived
allocated channel: ORA_DISK_1
channel ORA_DISK_1: SID=369 device type=DISK
channel ORA_DISK_1: starting archived log copy
input archived log thread=1 sequence=54 RECID=123 STAMP=912659482
RMAN-03009: failure of backup command on ORA_DISK_1 channel at 05/24/2016
04:36:22
ORA-19635: input and output file names are identical:
/rman/pancake/logship/1_54_912576125.dbf
continuing other job steps, job failed will not be re-run
channel ORA_DISK_1: starting archived log copy
input archived log thread=1 sequence=41 RECID=124 STAMP=912659483
RMAN-03009: failure of backup command on ORA_DISK_1 channel at 05/24/2016
04:36:23
ORA-19635: input and output file names are identical:
/rman/pancake/logship/1_41_912576125.dbf
continuing other job steps, job failed will not be re-run
...
channel ORA_DISK_1: starting archived log copy
```



```
input archived log thread=1 sequence=45 RECID=152 STAMP=912659514
RMAN-03009: failure of backup command on ORA_DISK_1 channel at 05/24/2016
04:36:55
ORA-19635: input and output file names are identical:
/rman/pancake/logship/1_45_912576125.dbf
continuing other job steps, job failed will not be re-run
channel ORA_DISK_1: starting archived log copy
input archived log thread=1 sequence=47 RECID=153 STAMP=912659515
RMAN-00571: =====
RMAN-00569: ===== ERROR MESSAGE STACK FOLLOWS =====
RMAN-00571: =====
RMAN-03009: failure of backup command on ORA_DISK_1 channel at 05/24/2016
04:36:57
ORA-19635: input and output file names are identical:
/rman/pancake/logship/1_47_912576125.dbf
Recovery Manager complete.
```

在收到記錄之後、必須重新播放記錄。先前的範例顯示使用 sqlplus 來手動執行 `recover database until cancel` 可輕鬆自動化。此處顯示的範例使用中所說的指令碼 ["重播待命資料庫上的記錄"](#)。指令碼會接受一個引數、指定需要重新執行作業的資料庫。此程序允許在多資料庫移轉工作中使用相同的指令碼。

```

[root@jpsc2 pancake]# ./replaylogs.pl PANCAKE
ORACLE_SID = [oracle] ? The Oracle base has been set to /orabin
SQL*Plus: Release 12.1.0.2.0 Production on Tue May 24 04:47:10 2016
Copyright (c) 1982, 2014, Oracle. All rights reserved.
Connected to:
Oracle Database 12c Enterprise Edition Release 12.1.0.2.0 - 64bit
Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
SQL> ORA-00279: change 560591 generated at 05/24/2016 03:33:56 needed for
thread 1
ORA-00289: suggestion : /rman/pancake/logship/1_54_912576125.dbf
ORA-00280: change 560591 for thread 1 is in sequence #54
Specify log: {<RET>=suggested | filename | AUTO | CANCEL}
ORA-00279: change 562219 generated at 05/24/2016 04:15:08 needed for
thread 1
ORA-00289: suggestion : /rman/pancake/logship/1_55_912576125.dbf
ORA-00280: change 562219 for thread 1 is in sequence #55
ORA-00278: log file '/rman/pancake/logship/1_54_912576125.dbf' no longer
needed for this recovery
ORA-00279: change 562370 generated at 05/24/2016 04:19:18 needed for
thread 1
ORA-00289: suggestion : /rman/pancake/logship/1_56_912576125.dbf
ORA-00280: change 562370 for thread 1 is in sequence #56
ORA-00278: log file '/rman/pancake/logship/1_55_912576125.dbf' no longer
needed for this recovery
...
ORA-00279: change 563137 generated at 05/24/2016 04:36:20 needed for
thread 1
ORA-00289: suggestion : /rman/pancake/logship/1_65_912576125.dbf
ORA-00280: change 563137 for thread 1 is in sequence #65
ORA-00278: log file '/rman/pancake/logship/1_64_912576125.dbf' no longer
needed for this recovery
ORA-00308: cannot open archived log
'/rman/pancake/logship/1_65_912576125.dbf'
ORA-27037: unable to obtain file status
Linux-x86_64 Error: 2: No such file or directory
Additional information: 3
SQL> Disconnected from Oracle Database 12c Enterprise Edition Release
12.1.0.2.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options

```

## 轉換

準備好切換至新環境時、您必須執行最後一次同步。使用一般檔案系統時、由於原始的重作記錄會複製並重新播放、因此很容易確保移轉的資料庫與原始資料庫 100% 同步。使用 ASM 執行此作業的方法並不理想。只有歸檔日誌可以輕鬆地重新記錄。為了確保不會遺失任何資料、必須謹慎執行原始資料庫的最終關機。

1. 首先、必須將資料庫暫時禁用、確保不會進行任何變更。這種停止可能包括停用排程作業、關閉接聽程式及 / 或關閉應用程式。
2. 執行此步驟後、大多數 DBA 會建立一個虛擬表格、做為關機的標記。
3. 強制記錄歸檔、以確保在歸檔記錄檔中記錄建立虛擬表格。若要這麼做、請執行下列命令：

```
SQL> create table cutovercheck as select * from dba_users;
Table created.
SQL> alter system archive log current;
System altered.
SQL> shutdown immediate;
Database closed.
Database dismounted.
ORACLE instance shut down.
```

4. 若要複製最後一個歸檔記錄檔、請執行下列命令。資料庫必須可用、但不可開啟。

```
SQL> startup mount;
ORACLE instance started.
Total System Global Area  805306368 bytes
Fixed Size                  2929552 bytes
Variable Size               331353200 bytes
Database Buffers            465567744 bytes
Redo Buffers                 5455872 bytes
Database mounted.
```

5. 若要複製歸檔記錄檔、請執行下列命令：

```

RMAN> configure channel device type disk format
'/rman/pancake/logship/%h_%e_%a.dbf';
2> backup as copy archivelog from time 'sysdate-2';
3>
4>
using target database control file instead of recovery catalog
old RMAN configuration parameters:
CONFIGURE CHANNEL DEVICE TYPE DISK FORMAT
'/rman/pancake/logship/%h_%e_%a.dbf';
new RMAN configuration parameters:
CONFIGURE CHANNEL DEVICE TYPE DISK FORMAT
'/rman/pancake/logship/%h_%e_%a.dbf';
new RMAN configuration parameters are successfully stored
Starting backup at 24-MAY-16
allocated channel: ORA_DISK_1
channel ORA_DISK_1: SID=8 device type=DISK
channel ORA_DISK_1: starting archived log copy
input archived log thread=1 sequence=54 RECID=123 STAMP=912659482
RMAN-03009: failure of backup command on ORA_DISK_1 channel at
05/24/2016 04:58:24
ORA-19635: input and output file names are identical:
/rman/pancake/logship/1_54_912576125.dbf
continuing other job steps, job failed will not be re-run
...
channel ORA_DISK_1: starting archived log copy
input archived log thread=1 sequence=45 RECID=152 STAMP=912659514
RMAN-03009: failure of backup command on ORA_DISK_1 channel at
05/24/2016 04:58:58
ORA-19635: input and output file names are identical:
/rman/pancake/logship/1_45_912576125.dbf
continuing other job steps, job failed will not be re-run
channel ORA_DISK_1: starting archived log copy
input archived log thread=1 sequence=47 RECID=153 STAMP=912659515
RMAN-00571: =====
RMAN-00569: ===== ERROR MESSAGE STACK FOLLOWS =====
RMAN-00571: =====
RMAN-03009: failure of backup command on ORA_DISK_1 channel at
05/24/2016 04:59:00
ORA-19635: input and output file names are identical:
/rman/pancake/logship/1_47_912576125.dbf

```

6. 最後、在新伺服器上重播剩餘的歸檔記錄。

```

[root@jpsc2 pancake]# ./replaylogs.pl PANCAKE
ORACLE_SID = [oracle] ? The Oracle base has been set to /orabin
SQL*Plus: Release 12.1.0.2.0 Production on Tue May 24 05:00:53 2016
Copyright (c) 1982, 2014, Oracle. All rights reserved.
Connected to:
Oracle Database 12c Enterprise Edition Release 12.1.0.2.0 - 64bit
Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
SQL> ORA-00279: change 563137 generated at 05/24/2016 04:36:20 needed
for thread 1
ORA-00289: suggestion : /rman/pancake/logship/1_65_912576125.dbf
ORA-00280: change 563137 for thread 1 is in sequence #65
Specify log: {<RET>=suggested | filename | AUTO | CANCEL}
ORA-00279: change 563629 generated at 05/24/2016 04:55:20 needed for
thread 1
ORA-00289: suggestion : /rman/pancake/logship/1_66_912576125.dbf
ORA-00280: change 563629 for thread 1 is in sequence #66
ORA-00278: log file '/rman/pancake/logship/1_65_912576125.dbf' no longer
needed
for this recovery
ORA-00308: cannot open archived log
'/rman/pancake/logship/1_66_912576125.dbf'
ORA-27037: unable to obtain file status
Linux-x86_64 Error: 2: No such file or directory
Additional information: 3
SQL> Disconnected from Oracle Database 12c Enterprise Edition Release
12.1.0.2.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options

```

7. 在此階段、複寫所有資料。資料庫已準備好從待命資料庫轉換為作用中的作業資料庫、然後開啟。

```

SQL> alter database activate standby database;
Database altered.
SQL> alter database open;
Database altered.

```

8. 確認虛擬表格是否存在、然後將其丟棄。

```

SQL> desc cutovercheck
Name                                                    Null?    Type
-----
-----
USERNAME                                                NOT NULL VARCHAR2 (128)
USER_ID                                                  NOT NULL NUMBER
PASSWORD                                                VARCHAR2 (4000)
ACCOUNT_STATUS                                          NOT NULL VARCHAR2 (32)
LOCK_DATE                                               DATE
EXPIRY_DATE                                             DATE
DEFAULT_TABLESPACE                                     NOT NULL VARCHAR2 (30)
TEMPORARY_TABLESPACE                                  NOT NULL VARCHAR2 (30)
CREATED                                                 NOT NULL DATE
PROFILE                                                 NOT NULL VARCHAR2 (128)
INITIAL_RSRC_CONSUMER_GROUP                            VARCHAR2 (128)
EXTERNAL_NAME                                           VARCHAR2 (4000)
PASSWORD_VERSIONS                                       VARCHAR2 (12)
EDITIONS_ENABLED                                       VARCHAR2 (1)
AUTHENTICATION_TYPE                                    VARCHAR2 (8)
PROXY_ONLY_CONNECT                                    VARCHAR2 (1)
COMMON                                                  VARCHAR2 (3)
LAST_LOGIN                                              TIMESTAMP (9) WITH
TIME_ZONE
ORACLE_MAINTAINED                                       VARCHAR2 (1)
SQL> drop table cutovercheck;
Table dropped.

```

#### 不中斷的重作記錄移轉

有時資料庫會在整體上正確組織、但重做記錄除外。這可能是因為許多原因、其中最常見的原因與快照有關。SnapManager for Oracle、SnapCenter 和 NetApp Snap Creator 儲存管理架構等產品可讓您近乎即時地恢復資料庫、但前提是您必須還原資料檔案磁碟區的狀態。如果重做記錄檔與資料檔案共用空間、則無法安全執行還原、因為還原會導致重做記錄檔毀損、這可能表示資料遺失。因此、重做記錄必須重新定位。

此程序很簡單、可在不中斷營運的情況下執行。

#### 目前的重做記錄組態

1. 識別重做記錄群組的數目及其各自的群組編號。

```

SQL> select group#||' '||member from v$logfile;
GROUP#||' '||MEMBER
-----
-----
1 /redo0/NTAP/redo01a.log
1 /redo1/NTAP/redo01b.log
2 /redo0/NTAP/redo02a.log
2 /redo1/NTAP/redo02b.log
3 /redo0/NTAP/redo03a.log
3 /redo1/NTAP/redo03b.log
rows selected.

```

## 2. 輸入重做記錄檔的大小。

```

SQL> select group#||' '||bytes from v$log;
GROUP#||' '||BYTES
-----
-----
1 524288000
2 524288000
3 524288000

```

## 建立新記錄

### 1. 針對每個重做記錄、建立一個大小和成員數目相符的新群組。

```

SQL> alter database add logfile ('/newredo0/redo01a.log',
'/newredo1/redo01b.log') size 500M;
Database altered.
SQL> alter database add logfile ('/newredo0/redo02a.log',
'/newredo1/redo02b.log') size 500M;
Database altered.
SQL> alter database add logfile ('/newredo0/redo03a.log',
'/newredo1/redo03b.log') size 500M;
Database altered.
SQL>

```

### 2. 驗證新組態。

```

SQL> select group#||' '||member from v$logfile;
GROUP#||' '||MEMBER
-----
1 /redo0/NTAP/redo01a.log
1 /redo1/NTAP/redo01b.log
2 /redo0/NTAP/redo02a.log
2 /redo1/NTAP/redo02b.log
3 /redo0/NTAP/redo03a.log
3 /redo1/NTAP/redo03b.log
4 /newredo0/redo01a.log
4 /newredo1/redo01b.log
5 /newredo0/redo02a.log
5 /newredo1/redo02b.log
6 /newredo0/redo03a.log
6 /newredo1/redo03b.log
12 rows selected.

```

## 刪除舊記錄

1. 刪除舊記錄（群組 1、2 和 3）。

```

SQL> alter database drop logfile group 1;
Database altered.
SQL> alter database drop logfile group 2;
Database altered.
SQL> alter database drop logfile group 3;
Database altered.

```

2. 如果您遇到錯誤、導致無法刪除作用中記錄、請強制切換至下一個記錄檔、以釋放鎖定並強制建立全域檢查點。請參閱下列此程序範例。由於此記錄檔中仍有作用中的資料、因此拒絕嘗試丟棄位於舊位置的記錄檔群組 2。

```

SQL> alter database drop logfile group 2;
alter database drop logfile group 2
*
ERROR at line 1:
ORA-01623: log 2 is current log for instance NTAP (thread 1) - cannot
drop
ORA-00312: online log 2 thread 1: '/redo0/NTAP/redo02a.log'
ORA-00312: online log 2 thread 1: '/redo1/NTAP/redo02b.log'

```

3. 記錄歸檔之後再加上檢查點、可讓您捨棄記錄檔。



```
SQL> alter system archive log current;
System altered.
SQL> alter system checkpoint;
System altered.
SQL> alter database drop logfile group 2;
Database altered.
```

4. 然後從檔案系統刪除記錄。您應該非常小心地執行此程序。

## Oracle 資料庫主機資料複本

如同資料庫層級的移轉、主機層的移轉也提供儲存設備廠商的不受侷連的方法。

換句話說、有時候「只複製檔案」是最佳選擇。

雖然這種低技術方法似乎過於基本、但它確實提供了顯著的效益、因為不需要特殊軟體、而且在程序期間、原始資料仍保持安全不變。主要的限制是檔案複製資料移轉是一項破壞性程序、因為必須在複製作業開始之前關閉資料庫。沒有適當的方法可以同步處理檔案中的變更、因此檔案必須在開始複製之前完全處於禁用狀態。

如果複製作業所需的關機不理想、則下一個最佳的主機型選項是使用邏輯 Volume Manager (LVM)。包括 Oracle ASM 在內的許多 LVM 選項都具有類似的功能、但也有一些必須考量的限制。在大多數情況下、可在不中斷或停機的情況下完成移轉。

### 檔案系統複製到檔案系統

不應低估簡單複製作業的效用。這項作業需要在複製程序期間停機、但這是一個非常可靠的程序、不需要操作系統、資料庫或儲存系統的專門知識。此外、它也非常安全、因為它不會影響原始資料。通常、系統管理員會將來源檔案系統變更為唯讀安裝、然後重新啟動伺服器、以保證沒有任何東西會損壞目前的資料。複製程序可以撰寫指令碼、確保能以最快的速度執行、而不會發生使用者錯誤的風險。由於 I/O 類型是簡單的資料循序傳輸、因此具有極高的頻寬效率。

下列範例示範安全快速移轉的一個選項。

### 環境

要移轉的環境如下：

- 目前的檔案系統

```
ontap-nfs1:/host1_oradata      52428800  16196928  36231872  31%
/oradata
ontap-nfs1:/host1_logs        49807360   548032   49259328  2% /logs
```

- 新檔案系統

```
ontap-nfs1:/host1_logs_new      49807360      128  49807232      1%
/new/logs
ontap-nfs1:/host1_oradata_new    49807360      128  49807232      1%
/new/oradata
```

## 總覽

資料庫可由 DBA 移轉、只需關閉資料庫並複製檔案即可、但如果必須移轉許多資料庫、或是將停機時間降至最低、則此程序很容易撰寫指令碼。使用指令碼也能降低使用者錯誤的機率。

所示範例指令碼可自動化下列作業：

- 關閉資料庫
- 將現有檔案系統轉換為唯讀狀態
- 將所有資料從來源複製到目標檔案系統、以保留所有檔案權限
- 卸載舊的和新的檔案系統
- 將新檔案系統重新掛載到與先前檔案系統相同的路徑

## 程序

### 1. 關閉資料庫。

```
[root@host1 current]# ./dbshut.pl NTAP
ORACLE_SID = [oracle] ? The Oracle base has been set to /orabin
SQL*Plus: Release 12.1.0.2.0 Production on Thu Dec 3 15:58:48 2015
Copyright (c) 1982, 2014, Oracle. All rights reserved.
Connected to:
Oracle Database 12c Enterprise Edition Release 12.1.0.2.0 - 64bit
Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
SQL> Database closed.
Database dismounted.
ORACLE instance shut down.
SQL> Disconnected from Oracle Database 12c Enterprise Edition Release
12.1.0.2.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
NTAP shut down
```

### 2. 將檔案系統轉換為唯讀。如所示、使用指令碼可以更快完成這項工作 "將檔案系統轉換為唯讀"。

```
[root@host1 current]# ./mk.fs.readonly.pl /oradata
/oradata unmounted
/oradata mounted read-only
[root@host1 current]# ./mk.fs.readonly.pl /logs
/logs unmounted
/logs mounted read-only
```

### 3. 確認檔案系統現在為唯讀。

```
ontap-nfs1:/host1_oradata on /oradata type nfs
(ro,bg,vers=3,rsz=65536,wsz=65536,addr=172.20.101.10)
ontap-nfs1:/host1_logs on /logs type nfs
(ro,bg,vers=3,rsz=65536,wsz=65536,addr=172.20.101.10)
```

### 4. 將檔案系統內容與同步 rsync 命令。

```
[root@host1 current]# rsync -rlpogt --stats --progress
--exclude=.snapshot /oradata/ /new/oradata/
sending incremental file list
./
NTAP/
NTAP/IOPS.dbf
 10737426432 100% 153.50MB/s   0:01:06 (xfer#1, to-check=10/13)
NTAP/iops.dbf.zip
  22823573 100%  12.09MB/s   0:00:01 (xfer#2, to-check=9/13)
...
NTAP/undotbs02.dbf
 1073750016 100% 131.60MB/s   0:00:07 (xfer#10, to-check=1/13)
NTAP/users01.dbf
  5251072 100%   3.95MB/s   0:00:01 (xfer#11, to-check=0/13)
Number of files: 13
Number of files transferred: 11
Total file size: 18570092218 bytes
Total transferred file size: 18570092218 bytes
Literal data: 18570092218 bytes
Matched data: 0 bytes
File list size: 277
File list generation time: 0.001 seconds
File list transfer time: 0.000 seconds
Total bytes sent: 18572359828
Total bytes received: 228
sent 18572359828 bytes  received 228 bytes  162204017.96 bytes/sec
total size is 18570092218  speedup is 1.00
```

```

[root@host1 current]# rsync -rlpogt --stats --progress
--exclude=.snapshot /logs/ /new/logs/
sending incremental file list
./
NTAP/
NTAP/1_22_897068759.dbf
    45523968 100%  95.98MB/s    0:00:00 (xfer#1, to-check=15/18)
NTAP/1_23_897068759.dbf
    40601088 100%  49.45MB/s    0:00:00 (xfer#2, to-check=14/18)
...
NTAP/redo/redo02.log
    52429312 100%  44.68MB/s    0:00:01 (xfer#12, to-check=1/18)
NTAP/redo/redo03.log
    52429312 100%  68.03MB/s    0:00:00 (xfer#13, to-check=0/18)
Number of files: 18
Number of files transferred: 13
Total file size: 527032832 bytes
Total transferred file size: 527032832 bytes
Literal data: 527032832 bytes
Matched data: 0 bytes
File list size: 413
File list generation time: 0.001 seconds
File list transfer time: 0.000 seconds
Total bytes sent: 527098156
Total bytes received: 278
sent 527098156 bytes  received 278 bytes  95836078.91 bytes/sec
total size is 527032832  speedup is 1.00

```

5. 卸載舊檔案系統、並重新放置複製的資料。如所示、使用指令碼可以更快完成這項工作 "取代檔案系統"。

```

[root@host1 current]# ./swap.fs.pl /logs,/new/logs
/new/logs unmounted
/logs unmounted
Updated /logs mounted
[root@host1 current]# ./swap.fs.pl /oradata,/new/oradata
/new/oradata unmounted
/oradata unmounted
Updated /oradata mounted

```

6. 確認新檔案系統已就位。

```
ontap-nfs1:/host1_logs_new on /logs type nfs
(rw,bg,vers=3,rsz=65536,wsz=65536,addr=172.20.101.10)
ontap-nfs1:/host1_oradata_new on /oradata type nfs
(rw,bg,vers=3,rsz=65536,wsz=65536,addr=172.20.101.10)
```

## 7. 啟動資料庫。

```
[root@host1 current]# ./dbstart.pl NTAP
ORACLE_SID = [oracle] ? The Oracle base has been set to /orabin
SQL*Plus: Release 12.1.0.2.0 Production on Thu Dec 3 16:10:07 2015
Copyright (c) 1982, 2014, Oracle. All rights reserved.
Connected to an idle instance.
SQL> ORACLE instance started.
Total System Global Area 805306368 bytes
Fixed Size 2929552 bytes
Variable Size 390073456 bytes
Database Buffers 406847488 bytes
Redo Buffers 5455872 bytes
Database mounted.
Database opened.
SQL> Disconnected from Oracle Database 12c Enterprise Edition Release
12.1.0.2.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
NTAP started
```

### 全自動轉換

此範例指令碼接受資料庫 SID 的引數、後面接著通用分隔的檔案系統配對。如前所示、命令發出方式如下：

```
[root@host1 current]# ./migrate.oracle.fs.pl NTAP /logs,/new/logs
/oradata,/new/oradata
```

執行時、範例指令碼會嘗試執行下列順序。如果在任何步驟中遇到錯誤、它都會終止：

1. 關閉資料庫。
2. 將目前的檔案系統轉換為唯讀狀態。
3. 使用每個以逗號分隔的檔案系統引數配對、並將第一個檔案系統同步到第二個檔案系統。
4. 卸除先前的檔案系統。
5. 更新 /etc/fstab 檔案如下：
  - a. 請在下列位置建立備份 /etc/fstab.bak。

- b. 註解先前和新檔案系統的先前項目。
  - c. 為使用舊掛載點的新檔案系統建立新項目。
6. 掛載檔案系統。
7. 啟動資料庫。

下列文字提供此指令碼的執行範例：

```
[root@host1 current]# ./migrate.oracle.fs.pl NTAP /logs,/new/logs
/oradata,/new/oradata
ORACLE_SID = [oracle] ? The Oracle base has been set to /orabin
SQL*Plus: Release 12.1.0.2.0 Production on Thu Dec 3 17:05:50 2015
Copyright (c) 1982, 2014, Oracle. All rights reserved.
Connected to:
Oracle Database 12c Enterprise Edition Release 12.1.0.2.0 - 64bit
Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
SQL> Database closed.
Database dismounted.
ORACLE instance shut down.
SQL> Disconnected from Oracle Database 12c Enterprise Edition Release
12.1.0.2.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
NTAP shut down
sending incremental file list
./
NTAP/
NTAP/1_22_897068759.dbf
    45523968 100% 185.40MB/s   0:00:00 (xfer#1, to-check=15/18)
NTAP/1_23_897068759.dbf
    40601088 100%  81.34MB/s   0:00:00 (xfer#2, to-check=14/18)
...
NTAP/redo/redo02.log
    52429312 100%  70.42MB/s   0:00:00 (xfer#12, to-check=1/18)
NTAP/redo/redo03.log
    52429312 100%  47.08MB/s   0:00:01 (xfer#13, to-check=0/18)
Number of files: 18
Number of files transferred: 13
Total file size: 527032832 bytes
Total transferred file size: 527032832 bytes
Literal data: 527032832 bytes
Matched data: 0 bytes
File list size: 413
File list generation time: 0.001 seconds
```

```

File list transfer time: 0.000 seconds
Total bytes sent: 527098156
Total bytes received: 278
sent 527098156 bytes received 278 bytes 150599552.57 bytes/sec
total size is 527032832 speedup is 1.00
Succesfully replicated filesystem /logs to /new/logs
sending incremental file list
./
NTAP/
NTAP/IOPS.dbf
  10737426432 100% 176.55MB/s 0:00:58 (xfer#1, to-check=10/13)
NTAP/iops.dbf.zip
  22823573 100% 9.48MB/s 0:00:02 (xfer#2, to-check=9/13)
... NTAP/undotbs01.dbf
  309338112 100% 70.76MB/s 0:00:04 (xfer#9, to-check=2/13)
NTAP/undotbs02.dbf
  1073750016 100% 187.65MB/s 0:00:05 (xfer#10, to-check=1/13)
NTAP/users01.dbf
  5251072 100% 5.09MB/s 0:00:00 (xfer#11, to-check=0/13)
Number of files: 13
Number of files transferred: 11
Total file size: 18570092218 bytes
Total transferred file size: 18570092218 bytes
Literal data: 18570092218 bytes
Matched data: 0 bytes
File list size: 277
File list generation time: 0.001 seconds
File list transfer time: 0.000 seconds
Total bytes sent: 18572359828
Total bytes received: 228
sent 18572359828 bytes received 228 bytes 177725933.55 bytes/sec
total size is 18570092218 speedup is 1.00
Succesfully replicated filesystem /oradata to /new/oradata
swap 0 /logs /new/logs
/new/logs unmounted
/logs unmounted
Mounted updated /logs
Swapped filesystem /logs for /new/logs
swap 1 /oradata /new/oradata
/new/oradata unmounted
/oradata unmounted
Mounted updated /oradata
Swapped filesystem /oradata for /new/oradata
ORACLE_SID = [oracle] ? The Oracle base has been set to /orabin
SQL*Plus: Release 12.1.0.2.0 Production on Thu Dec 3 17:08:59 2015
Copyright (c) 1982, 2014, Oracle. All rights reserved.

```

```
Connected to an idle instance.
SQL> ORACLE instance started.
Total System Global Area  805306368 bytes
Fixed Size                 2929552 bytes
Variable Size             390073456 bytes
Database Buffers          406847488 bytes
Redo Buffers               5455872 bytes
Database mounted.
Database opened.
SQL> Disconnected from Oracle Database 12c Enterprise Edition Release
12.1.0.2.0 - 64bit Production
With the Partitioning, OLAP, Advanced Analytics and Real Application
Testing options
NTAP started
[root@host1 current]#
```

### Oracle ASM spfile 和 passwd 移轉

在完成涉及 ASM 的移轉時、有一個困難是 ASM 專屬的 spfile 和密碼檔案。根據預設、這些關鍵中繼資料檔案會建立在定義的第一個 ASM 磁碟群組上。如果必須撤出和移除特定的 ASM 磁碟群組、則必須重新放置管理該 ASM 執行個體的 spfile 和密碼檔案。

另一個需要重新放置這些檔案的使用案例是在部署資料庫管理軟體時、例如 SnapManager for Oracle 或 SnapCenter Oracle 外掛程式。這些產品的其中一項功能是透過還原代管資料檔案的 ASM LUN 狀態、快速還原資料庫。這樣做需要在執行還原之前將 ASM 磁碟群組離線。只要指定資料庫的資料檔案隔離在專用的 ASM 磁碟群組中、這不是問題。

當該磁碟群組也包含 ASM spfile/passwd 檔案時、唯一可以將磁碟群組離線的方法是關閉整個 ASM 執行個體。這是一項破壞性程序、也就是說、spfile/passwd 檔案必須重新放置。

### 環境

1. 資料庫 SID = Toast
2. 目前的資料檔案位於 +DATA
3. 上目前的記錄檔和控制檔 +LOGS
4. 建立為的新 ASM 磁碟群組 +NEWDATA 和 +NEWLOGS

### ASM spfile/passwd 檔案位置

您可以不中斷地重新放置這些檔案。不過、為了安全起見、NetApp 建議您關閉資料庫環境、以便確定檔案已重新放置、且組態已正確更新。如果伺服器上有多個 ASM 執行個體、則必須重複此程序。

### 識別 ASM 執行個體

根據中記錄的資料來識別 ASM 執行個體 oratab 檔案：ASM 執行個體以 + 符號表示。



```
-bash-4.1$ cat /etc/oratab | grep '^+'
+ASM:/orabin/grid:N          # line added by Agent
```

此伺服器上有一個稱為 +ASM 的 ASM 執行個體。

確定所有資料庫都已關閉

唯一可見的 SMON 程序應該是使用中 ASM 執行個體的 SMON。另一個 SMON 程序的存在表示資料庫仍在執行中。

```
-bash-4.1$ ps -ef | grep smon
oracle      857      1  0 18:26 ?          00:00:00 asm_smon_+ASM
```

唯一的 SMON 程序是 ASM 執行個體本身。這表示沒有其他資料庫正在執行中、而且在不中斷資料庫作業的風險下繼續作業是安全的。

尋找檔案

使用識別 ASM spfile 和密碼檔案的目前位置 spget 和 pwget 命令。

```
bash-4.1$ asmcmd
ASMCMDB> spget
+DATA/spfile.ora
```

```
ASMCMDB> pwget --asm
+DATA/orapwasm
```

這些檔案都位於的基礎上 +DATA 磁碟群組。

複製檔案

使用將檔案複製到新的 ASM 磁碟群組 spcopy 和 pwcopu 命令。如果新磁碟群組是最近建立的、而且目前是空的、則可能需要先掛載。

```
ASMCMDB> mount NEWDATA
```

```
ASMCMDB> spcopy +DATA/spfile.ora +NEWDATA/spfile.ora
copying +DATA/spfile.ora -> +NEWDATA/spfilea.ora
```

```
ASMCMD> pwcopy +DATA/orapwasm +NEWDATA/orapwasm
copying +DATA/orapwasm -> +NEWDATA/orapwasm
```

檔案現已從複製 +DATA 至 +NEWDATA 。

### 更新 ASM 執行個體

現在必須更新 ASM 執行個體、以反映位置變更。◦ spset 和 pwset 命令會更新啟動 ASM 磁碟群組所需的 ASM 中繼資料。

```
ASMCMD> spset +NEWDATA/spfile.ora
ASMCMD> pwset --asm +NEWDATA/orapwasm
```

### 使用更新的檔案啟動 ASM

此時、ASM 執行個體仍會使用這些檔案的先前位置。必須重新啟動執行個體、以強制重新讀取新位置的檔案、並釋放先前檔案上的鎖定。

```
-bash-4.1$ sqlplus / as sysasm
SQL> shutdown immediate;
ASM diskgroups volume disabled
ASM diskgroups dismounted
ASM instance shutdown
```

```
SQL> startup
ASM instance started
Total System Global Area 1140850688 bytes
Fixed Size 2933400 bytes
Variable Size 1112751464 bytes
ASM Cache 25165824 bytes
ORA-15032: not all alterations performed
ORA-15017: diskgroup "NEWDATA" cannot be mounted
ORA-15013: diskgroup "NEWDATA" is already mounted
```

### 移除舊的 spfile 和密碼檔案

如果程序已成功執行、先前的檔案將不再鎖定、現在可以移除。

```
-bash-4.1$ asmcmd
ASMCMD> rm +DATA/spfile.ora
ASMCMD> rm +DATA/orapwasm
```

## Oracle ASM 至 ASM 複本

Oracle ASM 本質上是輕量的組合 Volume Manager 和檔案系統。由於檔案系統並不容易看到、因此 RMAN 必須用於執行複製作業。雖然複製型移轉程序既安全又簡單、但會造成部分中斷。可以將中斷降至最低、但不能完全消除。

如果您想要不中斷地移轉 ASM 型資料庫、最好的方法是利用 ASM 的功能、在移轉舊 LUN 的同時、重新平衡 ASM 擴充至新 LUN 的平衡。這樣做通常是安全且不中斷營運的、但它不提供回溯路徑。如果遇到功能或效能問題、唯一的選項是將資料移回來源。

您可以將資料庫複製到新位置而非移動資料、以避免此風險、避免原始資料受到影響。資料庫可以在新位置進行完整測試後再上線運作、如果發現問題、原始資料庫則可作為回復選項使用。

此程序是 RMAN 的眾多選項之一。其設計允許建立初始備份的兩個步驟程序、然後透過記錄重播進行同步處理。這項程序最適合將停機時間降至最低、因為它可讓資料庫在初始基準複本期間維持運作並提供資料。

### 複製資料庫

Oracle RMAN 會建立目前位於 ASM 磁碟群組的來源資料庫層級 0（完整）複本 +DATA 移至新位置 +NEWDATA。

```

-bash-4.1$ rman target /
Recovery Manager: Release 12.1.0.2.0 - Production on Sun Dec 6 17:40:03
2015
Copyright (c) 1982, 2014, Oracle and/or its affiliates. All rights
reserved.
connected to target database: TOAST (DBID=2084313411)
RMAN> backup as copy incremental level 0 database format '+NEWDATA' tag
'ONTAP_MIGRATION';
Starting backup at 06-DEC-15
using target database control file instead of recovery catalog
allocated channel: ORA_DISK_1
channel ORA_DISK_1: SID=302 device type=DISK
channel ORA_DISK_1: starting datafile copy
input datafile file number=00001
name=+DATA/TOAST/DATAFILE/system.262.897683141
...
input datafile file number=00004
name=+DATA/TOAST/DATAFILE/users.264.897683151
output file name=+NEWDATA/TOAST/DATAFILE/users.258.897759623
tag=ONTAP_MIGRATION RECID=5 STAMP=897759622
channel ORA_DISK_1: datafile copy complete, elapsed time: 00:00:01
channel ORA_DISK_1: starting incremental level 0 datafile backup set
channel ORA_DISK_1: specifying datafile(s) in backup set
including current SPFILE in backup set
channel ORA_DISK_1: starting piece 1 at 06-DEC-15
channel ORA_DISK_1: finished piece 1 at 06-DEC-15
piece
handle=+NEWDATA/TOAST/BACKUPSET/2015_12_06/nnsnn0_ontap_migration_0.262.89
7759623 tag=ONTAP_MIGRATION comment=NONE
channel ORA_DISK_1: backup set complete, elapsed time: 00:00:01
Finished backup at 06-DEC-15

```

### 強制歸檔記錄切換

您必須強制使用歸檔記錄切換、以確保歸檔記錄包含所有必要資料、使複本完全一致。如果沒有此命令、重做記錄檔中可能仍會有關鍵資料。

```

RMAN> sql 'alter system archive log current';
sql statement: alter system archive log current

```

### 關閉來源資料庫

由於資料庫已關機、並處於有限存取、唯讀模式、因此在此步驟中就會開始中斷。若要關閉來源資料庫、請執行下列命令：

```

RMAN> shutdown immediate;
using target database control file instead of recovery catalog
database closed
database dismounted
Oracle instance shut down
RMAN> startup mount;
connected to target database (not started)
Oracle instance started
database mounted
Total System Global Area      805306368 bytes
Fixed Size                     2929552 bytes
Variable Size                 390073456 bytes
Database Buffers              406847488 bytes
Redo Buffers                   5455872 bytes

```

## 控制檔備份

您必須備份控制檔、以防您必須中止移轉並還原至原始儲存位置。備份控制檔的複本並非 100% 必要、但它確實讓將資料庫檔案位置重設回原始位置的程序變得更簡單。

```

RMAN> backup as copy current controlfile format '/tmp/TOAST.ctrl';
Starting backup at 06-DEC-15
allocated channel: ORA_DISK_1
channel ORA_DISK_1: SID=358 device type=DISK
channel ORA_DISK_1: starting datafile copy
copying current control file
output file name=/tmp/TOAST.ctrl tag=TAG20151206T174753 RECID=6
STAMP=897760073
channel ORA_DISK_1: datafile copy complete, elapsed time: 00:00:01
Finished backup at 06-DEC-15

```

## 參數更新

目前的 spfile 包含對舊 ASM 磁碟群組內控制檔目前位置的參照。您必須編輯此檔案、只要編輯中繼 pfile 版本即可輕鬆完成。

```

RMAN> create pfile='/tmp/pfile' from spfile;
Statement processed

```

## 更新 pfile

更新任何參照舊 ASM 磁碟群組的參數、以反映新的 ASM 磁碟群組名稱。然後儲存更新的 pfile。請確定 db\_create 有參數存在。

在以下範例中、請參考 +DATA 變更為 +NEWDATA 以黃色反白顯示。兩個主要參數是 db\_create 在正確位置建立任何新檔案的參數。

```
*.compatible='12.1.0.2.0'  
*.control_files='+NEWLOGS/TOAST/CONTROLFILE/current.258.897683139'  
*.db_block_size=8192  
*. db_create_file_dest='+NEWDATA'  
*. db_create_online_log_dest_1='+NEWLOGS'  
*.db_domain=''   
*.db_name='TOAST'  
*.diagnostic_dest='/orabin'  
*.dispatchers='(PROTOCOL=TCP) (SERVICE=TOASTXDB)'  
*.log_archive_dest_1='LOCATION='+NEWLOGS'  
*.log_archive_format='%t_%s_%r.dbf'
```

### 更新 init.ora 檔案

大多數以 ASM 為基礎的資料庫都使用 init.ora 檔案位於 \$ORACLE\_HOME/dbs 目錄、指向 ASM 磁碟群組上的 spfile。此檔案必須重新導向至新 ASM 磁碟群組上的位置。

```
-bash-4.1$ cd $ORACLE_HOME/dbs  
-bash-4.1$ cat initTOAST.ora  
SPFILE='+DATA/TOAST/spfileTOAST.ora'
```

變更此檔案的方式如下：

```
SPFILE='+NEWLOGS/TOAST/spfileTOAST.ora
```

### 參數檔案重新建立

spfile 現在已準備好由編輯的 pfile 中的資料填入。

```
RMAN> create spfile from pfile='/tmp/pfile';  
Statement processed
```

### 啟動資料庫以開始使用新的 spfile

啟動資料庫以確保它現在使用新建立的 spfile、並正確記錄對系統參數的任何進一步變更。

```

RMAN> startup nomount;
connected to target database (not started)
Oracle instance started
Total System Global Area      805306368 bytes
Fixed Size                    2929552 bytes
Variable Size                 373296240 bytes
Database Buffers              423624704 bytes
Redo Buffers                   5455872 bytes

```

## 還原控制檔

RMAN 所建立的備份控制檔也可直接還原至新 spfile 中指定的位置。

```

RMAN> restore controlfile from
'+DATA/TOAST/CONTROLFILE/current.258.897683139';
Starting restore at 06-DEC-15
using target database control file instead of recovery catalog
allocated channel: ORA_DISK_1
channel ORA_DISK_1: SID=417 device type=DISK
channel ORA_DISK_1: copied control file copy
output file name=+NEWLOGS/TOAST/CONTROLFILE/current.273.897761061
Finished restore at 06-DEC-15

```

裝入資料庫並驗證新控制檔的使用。

```

RMAN> alter database mount;
using target database control file instead of recovery catalog
Statement processed

```

```

SQL> show parameter control_files;
NAME                                TYPE                                VALUE
-----                                -
control_files                       string
+NEWLOGS/TOAST/CONTROLFILE/cur
                                         rent.273.897761061

```

## 記錄重新播放

資料庫目前使用舊位置的資料檔案。在使用複本之前、必須先進行同步處理。初始複製程序已經過時間、變更主要記錄在歸檔記錄中。這些變更會複寫如下：

1. 執行包含歸檔記錄的 RMAN 遞增備份。

```
RMAN> backup incremental level 1 format '+NEWLOGS' for recover of copy
with tag 'ONTAP_MIGRATION' database;
Starting backup at 06-DEC-15
allocated channel: ORA_DISK_1
channel ORA_DISK_1: SID=62 device type=DISK
channel ORA_DISK_1: starting incremental level 1 datafile backup set
channel ORA_DISK_1: specifying datafile(s) in backup set
input datafile file number=00001
name=+DATA/TOAST/DATAFILE/system.262.897683141
input datafile file number=00002
name=+DATA/TOAST/DATAFILE/sysaux.260.897683143
input datafile file number=00003
name=+DATA/TOAST/DATAFILE/undotbs1.257.897683145
input datafile file number=00004
name=+DATA/TOAST/DATAFILE/users.264.897683151
channel ORA_DISK_1: starting piece 1 at 06-DEC-15
channel ORA_DISK_1: finished piece 1 at 06-DEC-15
piece
handle=+NEWLOGS/TOAST/BACKUPSET/2015_12_06/nnndn1_ontap_migration_0.268.
897762693 tag=ONTAP_MIGRATION comment=NONE
channel ORA_DISK_1: backup set complete, elapsed time: 00:00:01
channel ORA_DISK_1: starting incremental level 1 datafile backup set
channel ORA_DISK_1: specifying datafile(s) in backup set
including current control file in backup set
including current SPFILE in backup set
channel ORA_DISK_1: starting piece 1 at 06-DEC-15
channel ORA_DISK_1: finished piece 1 at 06-DEC-15
piece
handle=+NEWLOGS/TOAST/BACKUPSET/2015_12_06/ncsnn1_ontap_migration_0.267.
897762697 tag=ONTAP_MIGRATION comment=NONE
channel ORA_DISK_1: backup set complete, elapsed time: 00:00:01
Finished backup at 06-DEC-15
```

2. 重新播放記錄。



```

RMAN> recover copy of database with tag 'ONTAP_MIGRATION';
Starting recover at 06-DEC-15
using channel ORA_DISK_1
channel ORA_DISK_1: starting incremental datafile backup set restore
channel ORA_DISK_1: specifying datafile copies to recover
recovering datafile copy file number=00001
name=+NEWDATA/TOAST/DATAFILE/system.259.897759609
recovering datafile copy file number=00002
name=+NEWDATA/TOAST/DATAFILE/sysaux.263.897759615
recovering datafile copy file number=00003
name=+NEWDATA/TOAST/DATAFILE/undotbs1.264.897759619
recovering datafile copy file number=00004
name=+NEWDATA/TOAST/DATAFILE/users.258.897759623
channel ORA_DISK_1: reading from backup piece
+NEWLOGS/TOAST/BACKUPSET/2015_12_06/nnndn1_ontap_migration_0.268.8977626
93
channel ORA_DISK_1: piece
handle=+NEWLOGS/TOAST/BACKUPSET/2015_12_06/nnndn1_ontap_migration_0.268.
897762693 tag=ONTAP_MIGRATION
channel ORA_DISK_1: restored backup piece 1
channel ORA_DISK_1: restore complete, elapsed time: 00:00:01
Finished recover at 06-DEC-15

```

## 啟動

還原的控制檔仍會參照原始位置的資料檔案、也會包含複製資料檔案的路徑資訊。

1. 若要變更使用中的資料檔案、請執行 `switch database to copy` 命令。

```

RMAN> switch database to copy;
datafile 1 switched to datafile copy
"+NEWDATA/TOAST/DATAFILE/system.259.897759609"
datafile 2 switched to datafile copy
"+NEWDATA/TOAST/DATAFILE/sysaux.263.897759615"
datafile 3 switched to datafile copy
"+NEWDATA/TOAST/DATAFILE/undotbs1.264.897759619"
datafile 4 switched to datafile copy
"+NEWDATA/TOAST/DATAFILE/users.258.897759623"

```

使用中的資料檔案現在是複製的資料檔案、但最終的重做記錄檔中可能仍有變更。

2. 若要重播所有剩餘記錄、請執行 `recover database` 命令。如果出現此訊息 `media recovery complete` 出現時、程序成功。

```

RMAN> recover database;
Starting recover at 06-DEC-15
using channel ORA_DISK_1
starting media recovery
media recovery complete, elapsed time: 00:00:01
Finished recover at 06-DEC-15

```

此程序只會變更一般資料檔案的位置。必須重新命名暫存資料檔案、但不需要複製、因為它們只是暫時性的。資料庫目前關閉、因此暫存資料檔案中沒有作用中的資料。

3. 若要重新放置暫存資料檔案、請先識別其位置。

```

RMAN> select file#||' '||name from v$tempfile;
FILE#||' '||NAME
-----
-----
1 +DATA/TOAST/TEMPFILE/temp.263.897683145

```

4. 使用 RMAN 命令重新定位暫存資料檔案、為每個資料檔案設定新名稱。使用 Oracle 託管檔案（OMF）時、不需要完整名稱；ASM 磁碟群組已足夠。開啟資料庫時、OMF 會連結至 ASM 磁碟群組上的適當位置。若要重新定位檔案、請執行下列命令：

```

run {
set newname for tempfile 1 to '+NEWDATA';
switch tempfile all;
}

```

```

RMAN> run {
2> set newname for tempfile 1 to '+NEWDATA';
3> switch tempfile all;
4> }
executing command: SET NEWNAME
renamed tempfile 1 to +NEWDATA in control file

```

## 重做記錄移轉

移轉程序即將完成、但重做記錄仍位於原始 ASM 磁碟群組中。重作記錄無法直接重新定位。而是會建立新的重做記錄集、並將其新增至組態、然後刪除舊的記錄。

1. 識別重做記錄群組的數目及其各自的群組編號。

```

RMAN> select group#||' '||member from v$logfile;
GROUP#||' '||MEMBER
-----
-----
1 +DATA/TOAST/ONLINELOG/group_1.261.897683139
2 +DATA/TOAST/ONLINELOG/group_2.259.897683139
3 +DATA/TOAST/ONLINELOG/group_3.256.897683139

```

2. 輸入重做記錄檔的大小。

```

RMAN> select group#||' '||bytes from v$log;
GROUP#||' '||BYTES
-----
-----
1 52428800
2 52428800
3 52428800

```

3. 針對每個重做記錄、建立具有相符組態的新群組。如果您未使用 OMF、則必須指定完整路徑。這也是使用的範例 `db_create_online_log` 參數。如先前所示、此參數設為 `+NEWLOGS`。此組態可讓您使用下列命令來建立新的線上記錄檔、而無需指定檔案位置、甚至是特定的 ASM 磁碟群組。

```

RMAN> alter database add logfile size 52428800;
Statement processed
RMAN> alter database add logfile size 52428800;
Statement processed
RMAN> alter database add logfile size 52428800;
Statement processed

```

4. 開啟資料庫。

```

SQL> alter database open;
Database altered.

```

5. 刪除舊記錄。

```

RMAN> alter database drop logfile group 1;
Statement processed

```

6. 如果您遇到錯誤、導致無法刪除作用中記錄、請強制切換至下一個記錄檔、以釋放鎖定並強制建立全域檢查點。範例如下所示。嘗試丟棄位於舊位置的記錄檔群組 3、因為此記錄檔中仍有作用中資料、因此遭到拒

絕。檢查點之後的記錄封存可讓您刪除記錄檔。

```

RMAN> alter database drop logfile group 3;
RMAN-00571: =====
RMAN-00569: ===== ERROR MESSAGE STACK FOLLOWS =====
RMAN-00571: =====
RMAN-03002: failure of sql statement command at 12/08/2015 20:23:51
ORA-01623: log 3 is current log for instance TOAST (thread 4) - cannot
drop
ORA-00312: online log 3 thread 1:
'+LOGS/TOAST/ONLINELOG/group_3.259.897563549'
RMAN> alter system switch logfile;
Statement processed
RMAN> alter system checkpoint;
Statement processed
RMAN> alter database drop logfile group 3;
Statement processed

```

7. 檢閱環境、確定所有位置型參數都已更新。

```

SQL> select name from v$datafile;
SQL> select member from v$logfile;
SQL> select name from v$tempfile;
SQL> show parameter spfile;
SQL> select name, value from v$parameter where value is not null;

```

8. 下列指令碼示範如何簡化此程序：

```

[root@host1 current]# ./checkdbdata.pl TOAST
TOAST datafiles:
+NEWDATA/TOAST/DATAFILE/system.259.897759609
+NEWDATA/TOAST/DATAFILE/sysaux.263.897759615
+NEWDATA/TOAST/DATAFILE/undotbs1.264.897759619
+NEWDATA/TOAST/DATAFILE/users.258.897759623
TOAST redo logs:
+NEWLOGS/TOAST/ONLINELOG/group_4.266.897763123
+NEWLOGS/TOAST/ONLINELOG/group_5.265.897763125
+NEWLOGS/TOAST/ONLINELOG/group_6.264.897763125
TOAST temp datafiles:
+NEWDATA/TOAST/TEMPFILE/temp.260.897763165
TOAST spfile
spfile                                string
+NEWDATA/spfiletoast.ora
TOAST key parameters
control_files +NEWLOGS/TOAST/CONTROLFILE/current.273.897761061
log_archive_dest_1 LOCATION=+NEWLOGS
db_create_file_dest +NEWDATA
db_create_online_log_dest_1 +NEWLOGS

```

9. 如果 ASM 磁碟群組已完全撤出、現在可以使用卸載 `asmcmd`。不過、在許多情況下、屬於其他資料庫或 ASM spfile/passwd 檔案的檔案可能仍存在。

```

-bash-4.1$ . oraenv
ORACLE_SID = [TOAST] ? +ASM
The Oracle base remains unchanged with value /orabin
-bash-4.1$ asmcmd
ASMCMDB> umount DATA
ASMCMDB>

```

### Oracle ASM 至檔案系統複本

Oracle ASM 至檔案系統複製程序與 ASM 至 ASM 複製程序非常類似、具有類似的優點和限制。主要差異在於使用可見檔案系統時、不同命令和組態參數的語法、而非使用 ASM 磁碟群組。

### 複製資料庫

Oracle RMAN 用於建立目前位於 ASM 磁碟群組的來源資料庫層級 0（完整）複本 +DATA 移至新位置 /oradata。

```

RMAN> backup as copy incremental level 0 database format
'/oradata/TOAST/%U' tag 'ONTAP_MIGRATION';
Starting backup at 13-MAY-16
using target database control file instead of recovery catalog
allocated channel: ORA_DISK_1
channel ORA_DISK_1: SID=377 device type=DISK
channel ORA_DISK_1: starting datafile copy
input datafile file number=00001 name=+ASM0/TOAST/system01.dbf
output file name=/oradata/TOAST/data_D-TOAST_I-2098173325_TS-SYSTEM_FNO-
1_01r5fhjg tag=ONTAP_MIGRATION RECID=1 STAMP=911722099
channel ORA_DISK_1: datafile copy complete, elapsed time: 00:00:07
channel ORA_DISK_1: starting datafile copy
input datafile file number=00002 name=+ASM0/TOAST/sysaux01.dbf
output file name=/oradata/TOAST/data_D-TOAST_I-2098173325_TS-SYSAUX_FNO-
2_02r5fhjo tag=ONTAP_MIGRATION RECID=2 STAMP=911722106
channel ORA_DISK_1: datafile copy complete, elapsed time: 00:00:07
channel ORA_DISK_1: starting datafile copy
input datafile file number=00003 name=+ASM0/TOAST/undotbs101.dbf
output file name=/oradata/TOAST/data_D-TOAST_I-2098173325_TS-UNDOTBS1_FNO-
3_03r5fhjt tag=ONTAP_MIGRATION RECID=3 STAMP=911722113
channel ORA_DISK_1: datafile copy complete, elapsed time: 00:00:07
channel ORA_DISK_1: starting datafile copy
copying current control file
output file name=/oradata/TOAST/cf_D-TOAST_id-2098173325_04r5fhk5
tag=ONTAP_MIGRATION RECID=4 STAMP=911722118
channel ORA_DISK_1: datafile copy complete, elapsed time: 00:00:01
channel ORA_DISK_1: starting datafile copy
input datafile file number=00004 name=+ASM0/TOAST/users01.dbf
output file name=/oradata/TOAST/data_D-TOAST_I-2098173325_TS-USERS_FNO-
4_05r5fhk6 tag=ONTAP_MIGRATION RECID=5 STAMP=911722118
channel ORA_DISK_1: datafile copy complete, elapsed time: 00:00:01
channel ORA_DISK_1: starting incremental level 0 datafile backup set
channel ORA_DISK_1: specifying datafile(s) in backup set
including current SPFILE in backup set
channel ORA_DISK_1: starting piece 1 at 13-MAY-16
channel ORA_DISK_1: finished piece 1 at 13-MAY-16
piece handle=/oradata/TOAST/06r5fhk7_1_1 tag=ONTAP_MIGRATION comment=NONE
channel ORA_DISK_1: backup set complete, elapsed time: 00:00:01
Finished backup at 13-MAY-16

```

### 強制歸檔記錄切換

必須強制使用歸檔記錄交換器、才能確保歸檔記錄包含所有必要資料、使複本完全一致。如果沒有此命令、重做記錄檔中可能仍會有關鍵資料。若要強制使用歸檔記錄交換器、請執行下列命令：

```
RMAN> sql 'alter system archive log current';
sql statement: alter system archive log current
```

### 關閉來源資料庫

由於資料庫已關機、並處於有限存取的唯一讀模式、因此此步驟開始造成中斷。若要關閉來源資料庫、請執行下列命令：

```
RMAN> shutdown immediate;
using target database control file instead of recovery catalog
database closed
database dismounted
Oracle instance shut down
RMAN> startup mount;
connected to target database (not started)
Oracle instance started
database mounted
Total System Global Area      805306368 bytes
Fixed Size                    2929552 bytes
Variable Size                 331353200 bytes
Database Buffers              465567744 bytes
Redo Buffers                   5455872 bytes
```

### 控制檔備份

備份控制檔、以防您必須中止移轉並還原至原始儲存位置。備份控制檔的複本並非 100% 必要、但它確實讓將資料庫檔案位置重設回原始位置的程序變得更簡單。

```
RMAN> backup as copy current controlfile format '/tmp/TOAST.ctrl';
Starting backup at 08-DEC-15
using channel ORA_DISK_1
channel ORA_DISK_1: starting datafile copy
copying current control file
output file name=/tmp/TOAST.ctrl tag=TAG20151208T194540 RECID=30
STAMP=897939940
channel ORA_DISK_1: datafile copy complete, elapsed time: 00:00:01
Finished backup at 08-DEC-15
```

### 參數更新

```
RMAN> create pfile='/tmp/pfile' from spfile;
Statement processed
```

## 更新 pfile

任何參照舊 ASM 磁碟群組的參數都應該更新、在某些情況下、當不再相關時、就會刪除。更新它們以反映新的檔案系統路徑、並儲存更新的 pfile。請確定已列出完整的目標路徑。若要更新這些參數、請執行下列命令：

```
*.audit_file_dest='/orabin/admin/TOAST/adump'  
*.audit_trail='db'  
*.compatible='12.1.0.2.0'  
*.control_files='/logs/TOAST/arch/control01.ctl','/logs/TOAST/redo/control  
02.ctl'  
*.db_block_size=8192  
*.db_domain=''  
*.db_name='TOAST'  
*.diagnostic_dest='/orabin'  
*.dispatchers='(PROTOCOL=TCP) (SERVICE=TOASTXDB)'  
*.log_archive_dest_1='LOCATION=/logs/TOAST/arch'  
*.log_archive_format='%t_%s_%r.dbf'  
*.open_cursors=300  
*.pga_aggregate_target=256m  
*.processes=300  
*.remote_login_passwordfile='EXCLUSIVE'  
*.sga_target=768m  
*.undo_tablespace='UNDOTBS1'
```

## 停用原始的 init.ora 檔案

此檔案位於 \$ORACLE\_HOME/dbs 目錄和通常位於 pfile 中、作為指向 ASM 磁碟群組上 spfile 的指標。若要確定不再使用原始 spfile、請重新命名。不過、請勿刪除它、因為如果必須中止移轉、就需要此檔案。

```
[oracle@jfscl ~]$ cd $ORACLE_HOME/dbs  
[oracle@jfscl dbs]$ cat initTOAST.ora  
SPFILE='+ASM0/TOAST/spfileTOAST.ora'  
[oracle@jfscl dbs]$ mv initTOAST.ora initTOAST.ora.prev  
[oracle@jfscl dbs]$
```

## 參數檔案重新建立

這是重新定位 spfile 的最後一步。原始 spfile 不再使用、而且資料庫目前是使用中繼檔案啟動（但未掛載）。此檔案的內容可以寫入新的 spfile 位置、如下所示：

```
RMAN> create spfile from pfile='/tmp/pfile';  
Statement processed
```



## 啟動資料庫以開始使用新的 spfile

您必須啟動資料庫以釋放中繼檔案上的鎖定、並只使用新的 spfile 檔案來啟動資料庫。啟動資料庫也能證明新的 spfile 位置正確、而且其資料有效。

```
RMAN> shutdown immediate;
Oracle instance shut down
RMAN> startup nomount;
connected to target database (not started)
Oracle instance started
Total System Global Area      805306368 bytes
Fixed Size                     2929552 bytes
Variable Size                  331353200 bytes
Database Buffers               465567744 bytes
Redo Buffers                    5455872 bytes
```

## 還原控制檔

已在路徑上建立備份控制檔 /tmp/TOAST.ctrl 請稍早在程序中進行。新的 spfile 將控制檔位置定義為 /logfs/TOAST/ctrl/ctrlfile1.ctrl 和 /logfs/TOAST/redo/ctrlfile2.ctrl。不過、這些檔案尚不存在。

1. 此命令會將控制檔資料還原至 spfile 中定義的路徑。

```
RMAN> restore controlfile from '/tmp/TOAST.ctrl';
Starting restore at 13-MAY-16
using channel ORA_DISK_1
channel ORA_DISK_1: copied control file copy
output file name=/logs/TOAST/arch/control01.ctrl
output file name=/logs/TOAST/redo/control02.ctrl
Finished restore at 13-MAY-16
```

2. 發出 mount 命令、以便正確探索控制檔並包含有效資料。

```
RMAN> alter database mount;
Statement processed
released channel: ORA_DISK_1
```

驗證 control\_files 參數、請執行下列命令：

```
SQL> show parameter control_files;
NAME                                TYPE                                VALUE
-----                                -
control_files                       string
/logs/TOAST/arch/control01.ctl
,
/logs/TOAST/redo/control02.c
tl
```

## 記錄重新播放

資料庫目前正在使用舊位置的資料檔案。在使用複本之前、必須先同步資料檔案。在初始複製程序期間已經過時間、變更主要記錄在歸檔記錄中。以下兩個步驟會複寫這些變更。

1. 執行包含歸檔記錄的 RMAN 遞增備份。

```
RMAN> backup incremental level 1 format '/logs/TOAST/arch/%U' for
recover of copy with tag 'ONTAP_MIGRATION' database;
Starting backup at 13-MAY-16
using target database control file instead of recovery catalog
allocated channel: ORA_DISK_1
channel ORA_DISK_1: SID=124 device type=DISK
channel ORA_DISK_1: starting incremental level 1 datafile backup set
channel ORA_DISK_1: specifying datafile(s) in backup set
input datafile file number=00001 name=+ASM0/TOAST/system01.dbf
input datafile file number=00002 name=+ASM0/TOAST/sysaux01.dbf
input datafile file number=00003 name=+ASM0/TOAST/undotbs101.dbf
input datafile file number=00004 name=+ASM0/TOAST/users01.dbf
channel ORA_DISK_1: starting piece 1 at 13-MAY-16
channel ORA_DISK_1: finished piece 1 at 13-MAY-16
piece handle=/logs/TOAST/arch/09r5fj8i_1_1 tag=ONTAP_MIGRATION
comment=NONE
channel ORA_DISK_1: backup set complete, elapsed time: 00:00:01
Finished backup at 13-MAY-16
RMAN-06497: WARNING: control file is not current, control file
AUTOBACKUP skipped
```

2. 重播記錄。

```

RMAN> recover copy of database with tag 'ONTAP_MIGRATION';
Starting recover at 13-MAY-16
using channel ORA_DISK_1
channel ORA_DISK_1: starting incremental datafile backup set restore
channel ORA_DISK_1: specifying datafile copies to recover
recovering datafile copy file number=00001 name=/oradata/TOAST/data_D-
TOAST_I-2098173325_TS-SYSTEM_FNO-1_01r5fhjg
recovering datafile copy file number=00002 name=/oradata/TOAST/data_D-
TOAST_I-2098173325_TS-SYSAUX_FNO-2_02r5fhjo
recovering datafile copy file number=00003 name=/oradata/TOAST/data_D-
TOAST_I-2098173325_TS-UNDOTBS1_FNO-3_03r5fhjt
recovering datafile copy file number=00004 name=/oradata/TOAST/data_D-
TOAST_I-2098173325_TS-USERS_FNO-4_05r5fhk6
channel ORA_DISK_1: reading from backup piece
/logs/TOAST/arch/09r5fj8i_1_1
channel ORA_DISK_1: piece handle=/logs/TOAST/arch/09r5fj8i_1_1
tag=ONTAP_MIGRATION
channel ORA_DISK_1: restored backup piece 1
channel ORA_DISK_1: restore complete, elapsed time: 00:00:01
Finished recover at 13-MAY-16
RMAN-06497: WARNING: control file is not current, control file
AUTOBACKUP skipped

```

## 啟動

還原的控制檔仍會參照原始位置的資料檔案、也會包含複製資料檔案的路徑資訊。

1. 若要變更使用中的資料檔案、請執行 `switch database to copy` 命令：

```

RMAN> switch database to copy;
datafile 1 switched to datafile copy "/oradata/TOAST/data_D-TOAST_I-
2098173325_TS-SYSTEM_FNO-1_01r5fhjg"
datafile 2 switched to datafile copy "/oradata/TOAST/data_D-TOAST_I-
2098173325_TS-SYSAUX_FNO-2_02r5fhjo"
datafile 3 switched to datafile copy "/oradata/TOAST/data_D-TOAST_I-
2098173325_TS-UNDOTBS1_FNO-3_03r5fhjt"
datafile 4 switched to datafile copy "/oradata/TOAST/data_D-TOAST_I-
2098173325_TS-USERS_FNO-4_05r5fhk6"

```

2. 雖然資料檔案應完全一致、但仍需最後一步才能重播線上重作記錄中記錄的其餘變更。使用 `recover database` 命令重播這些變更、並使複本 100% 與原始版本相同。不過、複本尚未開啟。

```

RMAN> recover database;
Starting recover at 13-MAY-16
using channel ORA_DISK_1
starting media recovery
archived log for thread 1 with sequence 28 is already on disk as file
+ASM0/TOAST/redo01.log
archived log file name=+ASM0/TOAST/redo01.log thread=1 sequence=28
media recovery complete, elapsed time: 00:00:00
Finished recover at 13-MAY-16

```

## 重新部署暫存資料檔案

1. 識別仍在原始磁碟群組中使用的暫存資料檔案位置。

```

RMAN> select file#||' '||name from v$tempfile;
FILE#||' '||NAME
-----
-----
1 +ASM0/TOAST/temp01.dbf

```

2. 若要重新放置資料檔案、請執行下列命令。如果有許多 tempfiles、請使用文字編輯器建立 RMAN 命令、然後剪下並貼上。

```

RMAN> run {
2> set newname for tempfile 1 to '/oradata/TOAST/temp01.dbf';
3> switch tempfile all;
4> }
executing command: SET NEWNAME
renamed tempfile 1 to /oradata/TOAST/temp01.dbf in control file

```

## 重做記錄移轉

移轉程序即將完成、但重做記錄仍位於原始 ASM 磁碟群組中。重作記錄無法直接重新定位。而是建立新的重做記錄集、並在刪除舊記錄之後新增至組態。

1. 識別重做記錄群組的數目及其各自的群組編號。

```

RMAN> select group#||' '||member from v$logfile;
GROUP#||' '||MEMBER
-----
-----
1 +ASM0/TOAST/redo01.log
2 +ASM0/TOAST/redo02.log
3 +ASM0/TOAST/redo03.log

```

2. 輸入重做記錄檔的大小。

```

RMAN> select group#||' '||bytes from v$log;
GROUP#||' '||BYTES
-----
-----
1 52428800
2 52428800
3 52428800

```

3. 對於每個重做記錄、請使用與目前重做記錄群組相同的大小、使用新的檔案系統位置來建立新群組。

```

RMAN> alter database add logfile '/logs/TOAST/redo/log00.rdo' size
52428800;
Statement processed
RMAN> alter database add logfile '/logs/TOAST/redo/log01.rdo' size
52428800;
Statement processed
RMAN> alter database add logfile '/logs/TOAST/redo/log02.rdo' size
52428800;
Statement processed

```

4. 移除仍位於先前儲存設備上的舊記錄檔群組。

```

RMAN> alter database drop logfile group 4;
Statement processed
RMAN> alter database drop logfile group 5;
Statement processed
RMAN> alter database drop logfile group 6;
Statement processed

```

5. 如果遇到阻止刪除作用中記錄的錯誤、請強制切換至下一個記錄檔、以釋放鎖定並強制建立全域檢查點。範例如下所示。嘗試丟棄位於舊位置的記錄檔群組 3、因為此記錄檔中仍有作用中資料、因此遭到拒絕。記錄歸檔之後再加上檢查點、即可刪除記錄檔。

```

RMAN> alter database drop logfile group 4;
RMAN-00571: =====
RMAN-00569: ===== ERROR MESSAGE STACK FOLLOWS =====
RMAN-00571: =====
RMAN-03002: failure of sql statement command at 12/08/2015 20:23:51
ORA-01623: log 4 is current log for instance TOAST (thread 4) - cannot
drop
ORA-00312: online log 4 thread 1:
'+NEWLOGS/TOAST/ONLINELOG/group_4.266.897763123'
RMAN> alter system switch logfile;
Statement processed
RMAN> alter system checkpoint;
Statement processed
RMAN> alter database drop logfile group 4;
Statement processed

```

6. 檢閱環境、確定所有位置型參數都已更新。

```

SQL> select name from v$datafile;
SQL> select member from v$logfile;
SQL> select name from v$tempfile;
SQL> show parameter spfile;
SQL> select name, value from v$parameter where value is not null;

```

7. 下列指令碼示範如何簡化此程序。

```

[root@jfscl current]# ./checkdbdata.pl TOAST
TOAST datafiles:
/oradata/TOAST/data_D-TOAST_I-2098173325_TS-SYSTEM_FNO-1_01r5fhjg
/oradata/TOAST/data_D-TOAST_I-2098173325_TS-SYSAUX_FNO-2_02r5fhjo
/oradata/TOAST/data_D-TOAST_I-2098173325_TS-UNDOTBS1_FNO-3_03r5fhjt
/oradata/TOAST/data_D-TOAST_I-2098173325_TS-USERS_FNO-4_05r5fhk6
TOAST redo logs:
/logs/TOAST/redo/log00.rdo
/logs/TOAST/redo/log01.rdo
/logs/TOAST/redo/log02.rdo
TOAST temp datafiles:
/oradata/TOAST/temp01.dbf
TOAST spfile
spfile                                string
/orabin/product/12.1.0/dbhome_
                                         1/dbs/spfileTOAST.ora
TOAST key parameters
control_files /logs/TOAST/arch/control01.ctl,
/logs/TOAST/redo/control02.ctl
log_archive_dest_1 LOCATION=/logs/TOAST/arch

```

8. 如果 ASM 磁碟群組已完全撤出、現在可以使用卸載 `asmcmd`。在許多情況下、屬於其他資料庫或 ASM `spfile/passwd` 檔案的檔案仍會存在。

```

-bash-4.1$ . oraenv
ORACLE_SID = [TOAST] ? +ASM
The Oracle base remains unchanged with value /orabin
-bash-4.1$ asmcmd
ASMCMDB> umount DATA
ASMCMDB>

```

### 資料檔案清理程序

根據 Oracle RMAN 的使用方式而定、移轉程序可能會導致資料檔案的語法較長或較隱密。在此所示範例中、備份是以的檔案格式執行 `/oradata/TOAST/%U`。%U 表示 RMAN 應為每個資料檔案建立預設的唯一名稱。結果與下列文字所示類似。資料檔案的傳統名稱會內嵌在名稱中。您可以使用中所指的指令碼方法來清除此問題 **"ASM 移轉清理"**。

```
[root@jfscl current]# ./fixuniquenames.pl TOAST
#sqlplus Commands
shutdown immediate;
startup mount;
host mv /oradata/TOAST/data_D-TOAST_I-2098173325_TS-SYSTEM_FNO-1_01r5fhjg
/oradata/TOAST/system.dbf
host mv /oradata/TOAST/data_D-TOAST_I-2098173325_TS-SYSAUX_FNO-2_02r5fhjo
/oradata/TOAST/sysaux.dbf
host mv /oradata/TOAST/data_D-TOAST_I-2098173325_TS-UNDOTBS1_FNO-
3_03r5fhjt /oradata/TOAST/undotbs1.dbf
host mv /oradata/TOAST/data_D-TOAST_I-2098173325_TS-USERS_FNO-4_05r5fhk6
/oradata/TOAST/users.dbf
alter database rename file '/oradata/TOAST/data_D-TOAST_I-2098173325_TS-
SYSTEM_FNO-1_01r5fhjg' to '/oradata/TOAST/system.dbf';
alter database rename file '/oradata/TOAST/data_D-TOAST_I-2098173325_TS-
SYSAUX_FNO-2_02r5fhjo' to '/oradata/TOAST/sysaux.dbf';
alter database rename file '/oradata/TOAST/data_D-TOAST_I-2098173325_TS-
UNDOTBS1_FNO-3_03r5fhjt' to '/oradata/TOAST/undotbs1.dbf';
alter database rename file '/oradata/TOAST/data_D-TOAST_I-2098173325_TS-
USERS_FNO-4_05r5fhk6' to '/oradata/TOAST/users.dbf';
alter database open;
```

### Oracle ASM 重新平衡

如前所述、Oracle ASM 磁碟群組可透過重新平衡程序、以透明方式移轉至新的儲存系統。總而言之、重新平衡程序需要在現有的 LUN 群組中新增大小相同的 LUN、然後再中斷先前 LUN 的作業。Oracle ASM 會以最佳配置自動將基礎資料重新定位至新儲存設備、然後在完成時釋出舊的 LUN。

移轉程序使用高效率的循序 I/O、通常不會造成任何效能中斷、但可視需要調整移轉率。

### 識別要移轉的資料

```
SQL> select name||' '||group_number||' '||total_mb||' '||path||'
' ||header_status from v$asm_disk;
NEWDATA_0003 1 10240 /dev/mapper/3600a098038303537762b47594c315864 MEMBER
NEWDATA_0002 1 10240 /dev/mapper/3600a098038303537762b47594c315863 MEMBER
NEWDATA_0000 1 10240 /dev/mapper/3600a098038303537762b47594c315861 MEMBER
NEWDATA_0001 1 10240 /dev/mapper/3600a098038303537762b47594c315862 MEMBER
SQL> select group_number||' '||name from v$asm_diskgroup;
1 NEWDATA
```

### 建立新的 LUN

建立大小相同的新 LUN、並視需要設定使用者和群組成員資格。LUN 應顯示為 CANDIDATE 磁碟。



```
SQL> select name||' '||group_number||' '||total_mb||' '||path||'
' ||header_status from v$asm_disk;
0 0 /dev/mapper/3600a098038303537762b47594c31586b CANDIDATE
0 0 /dev/mapper/3600a098038303537762b47594c315869 CANDIDATE
0 0 /dev/mapper/3600a098038303537762b47594c315858 CANDIDATE
0 0 /dev/mapper/3600a098038303537762b47594c31586a CANDIDATE
NEWDATA_0003 1 10240 /dev/mapper/3600a098038303537762b47594c315864 MEMBER
NEWDATA_0002 1 10240 /dev/mapper/3600a098038303537762b47594c315863 MEMBER
NEWDATA_0000 1 10240 /dev/mapper/3600a098038303537762b47594c315861 MEMBER
NEWDATA_0001 1 10240 /dev/mapper/3600a098038303537762b47594c315862 MEMBER
```

## 新增 LUN

雖然可以同時執行新增和刪除作業、但通常只需兩個步驟即可輕鬆新增 LUN。首先、將新 LUN 新增至磁碟群組。此步驟會將一半的擴充從目前的 ASM LUN 移轉至新的 LUN。

重新平衡的力量代表資料傳輸的速度。資料傳輸的平行度越高、資料傳輸的數量就越多。執行移轉時、必須執行有效率的連續 I/O 作業、而這些作業不太可能造成效能問題。不過、若有需要、可利用調整進行中移轉的重新平衡能力 `alter diskgroup [name] rebalance power [level]` 命令。典型移轉使用 5 個值。

```
SQL> alter diskgroup NEWDATA add disk
'/dev/mapper/3600a098038303537762b47594c31586b' rebalance power 5;
Diskgroup altered.
SQL> alter diskgroup NEWDATA add disk
'/dev/mapper/3600a098038303537762b47594c315869' rebalance power 5;
Diskgroup altered.
SQL> alter diskgroup NEWDATA add disk
'/dev/mapper/3600a098038303537762b47594c315858' rebalance power 5;
Diskgroup altered.
SQL> alter diskgroup NEWDATA add disk
'/dev/mapper/3600a098038303537762b47594c31586a' rebalance power 5;
Diskgroup altered.
```

## 監控作業

可透過多種方式監控和管理重新平衡作業。在此範例中、我們使用下列命令。

```
SQL> select group_number,operation,state from v$asm_operation;
GROUP_NUMBER OPERA STAT
-----
1 REBAL RUN
1 REBAL WAIT
```

移轉完成時、不會回報任何重新平衡作業。

```
SQL> select group_number,operation,state from v$asm_operation;
no rows selected
```

## 丟棄舊的 LUN

移轉作業現在已完成一半。您可能需要執行一些基本效能測試、以確保環境健全。確認之後、可藉由丟棄舊的 LUN 來重新放置其餘的資料。請注意、這不會導致 LUN 立即發行。此中斷作業會先發出 Oracle ASM 重新定位延伸、然後再釋放 LUN。

```
sqlplus / as sysasm
SQL> alter diskgroup NEWDATA drop disk NEWDATA_0000 rebalance power 5;
Diskgroup altered.
SQL> alter diskgroup NEWDATA drop disk NEWDATA_0001 rebalance power 5;
Diskgroup altered.
SQL> alter diskgroup newdata drop disk NEWDATA_0002 rebalance power 5;
Diskgroup altered.
SQL> alter diskgroup newdata drop disk NEWDATA_0003 rebalance power 5;
Diskgroup altered.
```

## 監控作業

可透過多種方式監控和管理重新平衡作業。在此範例中、我們使用下列命令：

```
SQL> select group_number,operation,state from v$asm_operation;
GROUP_NUMBER OPERA STAT
-----
1 REBAL RUN
1 REBAL WAIT
```

移轉完成時、不會回報任何重新平衡作業。

```
SQL> select group_number,operation,state from v$asm_operation;
no rows selected
```

## 移除舊的 LUN

從磁碟群組移除舊 LUN 之前、您應該先對標頭狀態執行一次最後檢查。從 ASM 發佈 LUN 後、它不再列出名稱、而且標頭狀態會列為 FORMER。這表示這些 LUN 可以安全地從系統中移除。

```

SQL> select name||' '||group_number||' '||total_mb||' '||path||'
' ||header_status from v$asm_disk;
NAME||' '||GROUP_NUMBER||' '||TOTAL_MB||' '||PATH||' '||HEADER_STATUS
-----
-----
0 0 /dev/mapper/3600a098038303537762b47594c315863 FORMER
0 0 /dev/mapper/3600a098038303537762b47594c315864 FORMER
0 0 /dev/mapper/3600a098038303537762b47594c315861 FORMER
0 0 /dev/mapper/3600a098038303537762b47594c315862 FORMER
NEWDATA_0005 1 10240 /dev/mapper/3600a098038303537762b47594c315869 MEMBER
NEWDATA_0007 1 10240 /dev/mapper/3600a098038303537762b47594c31586a MEMBER
NEWDATA_0004 1 10240 /dev/mapper/3600a098038303537762b47594c31586b MEMBER
NEWDATA_0006 1 10240 /dev/mapper/3600a098038303537762b47594c315858 MEMBER
8 rows selected.

```

## LVM 移轉

此處介紹的程序顯示了以 LVM 為基礎的磁碟區群組移轉原則、稱為 `datavg`。這些範例來自 Linux LVM、但這些原則同樣適用於 AIX、HP-UX 和 VxVM。精確命令可能會有所不同。

1. 識別目前在中的 LUN `datavg` Volume 群組。

```

[root@host1 ~]# pvdisplay -C | grep datavg
/dev/mapper/3600a098038303537762b47594c31582f datavg lvm2 a-- 10.00g
10.00g
/dev/mapper/3600a098038303537762b47594c31585a datavg lvm2 a-- 10.00g
10.00g
/dev/mapper/3600a098038303537762b47594c315859 datavg lvm2 a-- 10.00g
10.00g
/dev/mapper/3600a098038303537762b47594c31586c datavg lvm2 a-- 10.00g
10.00g

```

2. 建立相同或稍大實體大小的新 LUN、並將其定義為實體磁碟區。

```
[root@host1 ~]# pvcreate /dev/mapper/3600a098038303537762b47594c315864
Physical volume "/dev/mapper/3600a098038303537762b47594c315864"
successfully created
[root@host1 ~]# pvcreate /dev/mapper/3600a098038303537762b47594c315863
Physical volume "/dev/mapper/3600a098038303537762b47594c315863"
successfully created
[root@host1 ~]# pvcreate /dev/mapper/3600a098038303537762b47594c315862
Physical volume "/dev/mapper/3600a098038303537762b47594c315862"
successfully created
[root@host1 ~]# pvcreate /dev/mapper/3600a098038303537762b47594c315861
Physical volume "/dev/mapper/3600a098038303537762b47594c315861"
successfully created
```

3. 將新的磁碟區新增至磁碟區群組。

```
[root@host1 tmp]# vgextend datavg
/dev/mapper/3600a098038303537762b47594c315864
Volume group "datavg" successfully extended
[root@host1 tmp]# vgextend datavg
/dev/mapper/3600a098038303537762b47594c315863
Volume group "datavg" successfully extended
[root@host1 tmp]# vgextend datavg
/dev/mapper/3600a098038303537762b47594c315862
Volume group "datavg" successfully extended
[root@host1 tmp]# vgextend datavg
/dev/mapper/3600a098038303537762b47594c315861
Volume group "datavg" successfully extended
```

4. 發行 `pvmove` 命令將每個目前 LUN 的範圍重新放置到新 LUN。 - `i [seconds]` 引數會監控作業的進度。

```

[root@host1 tmp]# pvmove -i 10
/dev/mapper/3600a098038303537762b47594c31582f
/dev/mapper/3600a098038303537762b47594c315864
  /dev/mapper/3600a098038303537762b47594c31582f: Moved: 0.0%
  /dev/mapper/3600a098038303537762b47594c31582f: Moved: 14.2%
  /dev/mapper/3600a098038303537762b47594c31582f: Moved: 28.4%
  /dev/mapper/3600a098038303537762b47594c31582f: Moved: 42.5%
  /dev/mapper/3600a098038303537762b47594c31582f: Moved: 57.1%
  /dev/mapper/3600a098038303537762b47594c31582f: Moved: 72.3%
  /dev/mapper/3600a098038303537762b47594c31582f: Moved: 87.3%
  /dev/mapper/3600a098038303537762b47594c31582f: Moved: 100.0%
[root@host1 tmp]# pvmove -i 10
/dev/mapper/3600a098038303537762b47594c31585a
/dev/mapper/3600a098038303537762b47594c315863
  /dev/mapper/3600a098038303537762b47594c31585a: Moved: 0.0%
  /dev/mapper/3600a098038303537762b47594c31585a: Moved: 14.9%
  /dev/mapper/3600a098038303537762b47594c31585a: Moved: 29.9%
  /dev/mapper/3600a098038303537762b47594c31585a: Moved: 44.8%
  /dev/mapper/3600a098038303537762b47594c31585a: Moved: 60.1%
  /dev/mapper/3600a098038303537762b47594c31585a: Moved: 75.8%
  /dev/mapper/3600a098038303537762b47594c31585a: Moved: 90.9%
  /dev/mapper/3600a098038303537762b47594c31585a: Moved: 100.0%
[root@host1 tmp]# pvmove -i 10
/dev/mapper/3600a098038303537762b47594c315859
/dev/mapper/3600a098038303537762b47594c315862
  /dev/mapper/3600a098038303537762b47594c315859: Moved: 0.0%
  /dev/mapper/3600a098038303537762b47594c315859: Moved: 14.8%
  /dev/mapper/3600a098038303537762b47594c315859: Moved: 29.8%
  /dev/mapper/3600a098038303537762b47594c315859: Moved: 45.5%
  /dev/mapper/3600a098038303537762b47594c315859: Moved: 61.1%
  /dev/mapper/3600a098038303537762b47594c315859: Moved: 76.6%
  /dev/mapper/3600a098038303537762b47594c315859: Moved: 91.7%
  /dev/mapper/3600a098038303537762b47594c315859: Moved: 100.0%
[root@host1 tmp]# pvmove -i 10
/dev/mapper/3600a098038303537762b47594c31586c
/dev/mapper/3600a098038303537762b47594c315861
  /dev/mapper/3600a098038303537762b47594c31586c: Moved: 0.0%
  /dev/mapper/3600a098038303537762b47594c31586c: Moved: 15.0%
  /dev/mapper/3600a098038303537762b47594c31586c: Moved: 30.4%
  /dev/mapper/3600a098038303537762b47594c31586c: Moved: 46.0%
  /dev/mapper/3600a098038303537762b47594c31586c: Moved: 61.4%
  /dev/mapper/3600a098038303537762b47594c31586c: Moved: 77.2%
  /dev/mapper/3600a098038303537762b47594c31586c: Moved: 92.3%
  /dev/mapper/3600a098038303537762b47594c31586c: Moved: 100.0%

```

5. 完成此程序後、請使用從磁碟區群組中刪除舊的 LUN `vgreduce` 命令。如果成功、現在即可安全地從系統移除 LUN。

```
[root@host1 tmp]# vgreduce datavg
/dev/mapper/3600a098038303537762b47594c31582f
Removed "/dev/mapper/3600a098038303537762b47594c31582f" from volume
group "datavg"
[root@host1 tmp]# vgreduce datavg
/dev/mapper/3600a098038303537762b47594c31585a
Removed "/dev/mapper/3600a098038303537762b47594c31585a" from volume
group "datavg"
[root@host1 tmp]# vgreduce datavg
/dev/mapper/3600a098038303537762b47594c315859
Removed "/dev/mapper/3600a098038303537762b47594c315859" from volume
group "datavg"
[root@host1 tmp]# vgreduce datavg
/dev/mapper/3600a098038303537762b47594c31586c
Removed "/dev/mapper/3600a098038303537762b47594c31586c" from volume
group "datavg"
```

## 外部 LUN 匯入

使用 FLI 進行 Oracle 移轉：規劃

NetApp 中記錄了使用 FLI 移轉 SAN 資源的程序 ["TR-4380：使用外部 LUN Import 進行 SAN 移轉"](#)。

從資料庫和主機的觀點來看、不需要採取任何特殊步驟。更新 FC 區域並在 ONTAP 上提供 LUN 之後、LVM 應該能夠從 LUN 讀取 LVM 中繼資料。此外、這些磁碟區群組也可以開始使用、無需進一步的組態步驟。在極少數情況下、環境可能會包含硬編碼的組態檔案、其中包含先前儲存陣列的參考資料。例如、內含的 Linux 系統 `/etc/multipath.conf` 參考指定裝置 WWN 的規則必須更新、以反映 FLI 所做的變更。



如需支援組態的相關資訊、請參閱 NetApp 相容性對照表。如果您的環境未包含在內、請聯絡 NetApp 代表以取得協助。

此範例顯示 Linux 伺服器上代管的 ASM 和 LVM LUN 移轉。其他作業系統支援 FLI、雖然主機端命令可能不同、但原則相同、ONTAP 程序相同。

## 識別 LVM LUN

準備的第一步是識別要移轉的 LUN。在此所示範例中、會在裝載兩個 SAN 型檔案系統 `/orabin` 和 `/backups`。

```
[root@host1 ~]# df -k
Filesystem                1K-blocks      Used Available Use%
Mounted on
/dev/mapper/rhel-root      52403200    8811464  43591736  17% /
devtmpfs                   65882776         0  65882776   0% /dev
...
fas8060-nfs-public:/install 199229440 119368128  79861312  60%
/install
/dev/mapper/sanvg-lvorabin  20961280  12348476   8612804  59%
/orabin
/dev/mapper/sanvg-lvbackups 73364480  62947536  10416944  86%
/backups
```

Volume 群組的名稱可以從裝置名稱中擷取、該名稱使用格式（Volume 群組名稱） - （邏輯磁碟區名稱）。在這種情況下、會呼叫 Volume 群組 sanvg。

◦ pvdisplay 命令可用於識別支援此 Volume 群組的 LUN、如下所示。在這種情況下、共有 10 個 LUN 組成 sanvg Volume 群組。

```
[root@host1 ~]# pvdisplay -C -o pv_name,pv_size,pv_fmt,vg_name
PV                               PSize  VG
/dev/mapper/3600a0980383030445424487556574266 10.00g sanvg
/dev/mapper/3600a0980383030445424487556574267 10.00g sanvg
/dev/mapper/3600a0980383030445424487556574268 10.00g sanvg
/dev/mapper/3600a0980383030445424487556574269 10.00g sanvg
/dev/mapper/3600a098038303044542448755657426a 10.00g sanvg
/dev/mapper/3600a098038303044542448755657426b 10.00g sanvg
/dev/mapper/3600a098038303044542448755657426c 10.00g sanvg
/dev/mapper/3600a098038303044542448755657426d 10.00g sanvg
/dev/mapper/3600a098038303044542448755657426e 10.00g sanvg
/dev/mapper/3600a098038303044542448755657426f 10.00g sanvg
/dev/sda2                          278.38g rhel
```

## 識別 ASM LUN

ASM LUN 也必須移轉。若要以 sysasm 使用者的身分從 sqlplus 取得 LUN 和 LUN 路徑的數目、請執行下列命令：

```

SQL> select path||' '||os_mb from v$asm_disk;
PATH||' '||OS_MB
-----
-----
/dev/oracleasm/disks/ASM0 10240
/dev/oracleasm/disks/ASM9 10240
/dev/oracleasm/disks/ASM8 10240
/dev/oracleasm/disks/ASM7 10240
/dev/oracleasm/disks/ASM6 10240
/dev/oracleasm/disks/ASM5 10240
/dev/oracleasm/disks/ASM4 10240
/dev/oracleasm/disks/ASM1 10240
/dev/oracleasm/disks/ASM3 10240
/dev/oracleasm/disks/ASM2 10240
10 rows selected.
SQL>

```

## FC 網路變更

目前環境包含 20 個要移轉的 LUN。更新目前的 SAN、讓 ONTAP 能夠存取目前的 LUN。資料尚未移轉、但 ONTAP 必須從目前的 LUN 讀取組態資訊、才能為該資料建立新的主目錄。

AFF/FAS 系統上至少必須將一個 HBA 連接埠設定為啟動器連接埠。此外、必須更新 FC 區域、讓 ONTAP 能夠存取外部儲存陣列上的 LUN。某些儲存陣列已設定 LUN 遮罩、限制哪些 WWN 可以存取指定的 LUN。在這種情況下、LUN 遮罩也必須更新、才能授予 ONTAP WWN 存取權。

完成此步驟後、ONTAP 應能使用檢視外部儲存陣列 `storage array show` 命令。它傳回的關鍵欄位是用來識別系統上外部 LUN 的首碼。在以下範例中、為外部陣列上的 LUN `FOREIGN_1` 在 ONTAP 中使用前置碼顯示 `FOR-1`。

### 識別外部陣列

```

Cluster01::> storage array show -fields name,prefix
name          prefix
-----
FOREIGN_1     FOR-1
Cluster01::>

```

### 識別外部 LUN

可以通過傳送來列出 LUN `array-name` 至 `storage disk show` 命令。移轉程序期間會多次參照傳回的資料。



```

Cluster01::> storage disk show -array-name FOREIGN_1 -fields disk,serial
disk      serial-number
-----
FOR-1.1   800DT$HuVWBX
FOR-1.2   800DT$HuVWBZ
FOR-1.3   800DT$HuVWBW
FOR-1.4   800DT$HuVWB Y
FOR-1.5   800DT$HuVWB/
FOR-1.6   800DT$HuVWBa
FOR-1.7   800DT$HuVWBd
FOR-1.8   800DT$HuVWBb
FOR-1.9   800DT$HuVWBc
FOR-1.10  800DT$HuVWB e
FOR-1.11  800DT$HuVWBf
FOR-1.12  800DT$HuVWBg
FOR-1.13  800DT$HuVWB i
FOR-1.14  800DT$HuVWBh
FOR-1.15  800DT$HuVWBj
FOR-1.16  800DT$HuVWBk
FOR-1.17  800DT$HuVWBm
FOR-1.18  800DT$HuVWB l
FOR-1.19  800DT$HuVWB o
FOR-1.20  800DT$HuVWB n
20 entries were displayed.
Cluster01::>

```

### 將外部陣列 LUN 登錄為匯入候選項目

外部 LUN 一開始會歸類為任何特定的 LUN 類型。在匯入資料之前、必須將 LUN 標記為外部、因此是匯入程序的候選項目。將序號傳送至即可完成此步驟 `storage disk modify` 命令、如下列範例所示。請注意、此程序只會將 LUN 標記為 ONTAP 中的外部。不會將任何資料寫入外部 LUN 本身。

```

Cluster01::*> storage disk modify {-serial-number 800DT$HuVWBW} -is
-foreign true
Cluster01::*> storage disk modify {-serial-number 800DT$HuVWBX} -is
-foreign true
...
Cluster01::*> storage disk modify {-serial-number 800DT$HuVWBn} -is
-foreign true
Cluster01::*> storage disk modify {-serial-number 800DT$HuVWB o} -is
-foreign true
Cluster01::*>

```

## 建立磁碟區以裝載移轉的 LUN

需要一個磁碟區來裝載移轉的 LUN。確切的 Volume 組態取決於運用 ONTAP 功能的整體計畫。在此範例中、ASM LUN 會放置在一個磁碟區中、而 LVM LUN 則放置在第二個磁碟區中。這樣做可讓您將 LUN 當作個別群組來管理、例如分層、建立快照或設定 QoS 控制。

設定 `snapshot-policy`to`none`。移轉程序可能包括大量資料流動。因此、如果快照是意外建立的、可能會大幅增加空間使用量、因為快照中會擷取不需要的資料。

```
Cluster01::> volume create -volume new_asm -aggregate data_02 -size 120G
-snapshot-policy none
[Job 1152] Job succeeded: Successful
Cluster01::> volume create -volume new_lvm -aggregate data_02 -size 120G
-snapshot-policy none
[Job 1153] Job succeeded: Successful
Cluster01::>
```

## 建立 ONTAP LUN

建立磁碟區之後、必須建立新的 LUN。一般而言、建立 LUN 需要使用者指定 LUN 大小之類的資訊、但在此情況下、外部磁碟引數會傳遞給命令。因此、ONTAP 會從指定的序號複寫目前的 LUN 組態資料。它也會使用 LUN 幾何資料和分割表格資料來調整 LUN 對齊、並建立最佳效能。

在此步驟中、序號必須與外部陣列交叉參照、以確保正確的外部 LUN 與正確的新 LUN 相符。

```
Cluster01::*> lun create -vserver vserver1 -path /vol/new_asm/LUN0 -ostype
linux -foreign-disk 800DT$HuVWBW
Created a LUN of size 10g (10737418240)
Cluster01::*> lun create -vserver vserver1 -path /vol/new_asm/LUN1 -ostype
linux -foreign-disk 800DT$HuVWBX
Created a LUN of size 10g (10737418240)
...
Created a LUN of size 10g (10737418240)
Cluster01::*> lun create -vserver vserver1 -path /vol/new_lvm/LUN8 -ostype
linux -foreign-disk 800DT$HuVWBn
Created a LUN of size 10g (10737418240)
Cluster01::*> lun create -vserver vserver1 -path /vol/new_lvm/LUN9 -ostype
linux -foreign-disk 800DT$HuVWBo
Created a LUN of size 10g (10737418240)
```

## 建立匯入關係

LUN 現已建立、但尚未設定為複寫目的地。在執行此步驟之前、必須先將 LUN 離線。這項額外步驟旨在保護資料不受使用者錯誤影響。如果 ONTAP 允許在線上 LUN 上執行移轉、可能會造成打字錯誤、導致覆寫作用中資料。強制使用者先將 LUN 離線的額外步驟、有助於確認使用正確的目標 LUN 做為移轉目的地。

```

Cluster01::*> lun offline -vserver vserver1 -path /vol/new_asm/LUN0
Warning: This command will take LUN "/vol/new_asm/LUN0" in Vserver
        "vserver1" offline.
Do you want to continue? {y|n}: y
Cluster01::*> lun offline -vserver vserver1 -path /vol/new_asm/LUN1
Warning: This command will take LUN "/vol/new_asm/LUN1" in Vserver
        "vserver1" offline.
Do you want to continue? {y|n}: y
...
Warning: This command will take LUN "/vol/new_lvm/LUN8" in Vserver
        "vserver1" offline.
Do you want to continue? {y|n}: y
Cluster01::*> lun offline -vserver vserver1 -path /vol/new_lvm/LUN9
Warning: This command will take LUN "/vol/new_lvm/LUN9" in Vserver
        "vserver1" offline.
Do you want to continue? {y|n}: y

```

LUN 離線後、您可以將外部 LUN 序號傳送至、以建立匯入關係 `lun import create` 命令。

```

Cluster01::*> lun import create -vserver vserver1 -path /vol/new_asm/LUN0
-foreign-disk 800DT$HuVWBW
Cluster01::*> lun import create -vserver vserver1 -path /vol/new_asm/LUN1
-foreign-disk 800DT$HuVWBX
...
Cluster01::*> lun import create -vserver vserver1 -path /vol/new_lvm/LUN8
-foreign-disk 800DT$HuVWBn
Cluster01::*> lun import create -vserver vserver1 -path /vol/new_lvm/LUN9
-foreign-disk 800DT$HuVWBo
Cluster01::*>

```

建立所有匯入關係之後、即可將 LUN 重新上線。

```

Cluster01::*> lun online -vserver vserver1 -path /vol/new_asm/LUN0
Cluster01::*> lun online -vserver vserver1 -path /vol/new_asm/LUN1
...
Cluster01::*> lun online -vserver vserver1 -path /vol/new_lvm/LUN8
Cluster01::*> lun online -vserver vserver1 -path /vol/new_lvm/LUN9
Cluster01::*>

```

## 建立啟動器群組

啟動器群組 (igroup) 是 ONTAP LUN 遮罩架構的一部分。除非先授予主機存取權、否則無法存取新建立的 LUN。這是透過建立一個 igroup、列出應授予存取權的 FC WWN 或 iSCSI 啟動器名稱來完成。在撰寫本報告

時、僅 FC LUN 支援 FLI。不過、轉換為 iSCSI 後移轉是一項簡單的工作、如所示 "傳輸協定轉換"。

在此範例中、會建立一個 igroup、其中包含兩個 WWN、對應於主機 HBA 上可用的兩個連接埠。

```
Cluster01::*> igroup create linuxhost -protocol fcp -ostype linux
-initiator 21:00:00:0e:1e:16:63:50 21:00:00:0e:1e:16:63:51
```

### 將新 LUN 對應至主機

在建立 igroup 之後、LUN 會對應至定義的 igroup。這些 LUN 僅適用於此 igroup 中包含的 WWN。NetApp 假設移轉程序目前階段主機尚未分區至 ONTAP。這一點很重要、因為如果主機同時分區到外部陣列和新的 ONTAP 系統、則可能會在每個陣列上發現具有相同序號的 LUN。這種情況可能導致多重路徑故障或資料受損。

```
Cluster01::*> lun map -vserver vserver1 -path /vol/new_asm/LUN0 -igroup
linuxhost
Cluster01::*> lun map -vserver vserver1 -path /vol/new_asm/LUN1 -igroup
linuxhost
...
Cluster01::*> lun map -vserver vserver1 -path /vol/new_lvm/LUN8 -igroup
linuxhost
Cluster01::*> lun map -vserver vserver1 -path /vol/new_lvm/LUN9 -igroup
linuxhost
Cluster01::*>
```

### 使用 FLI 進行 Oracle 移轉：轉換

由於需要變更 FC 網路組態、因此無法避免在外部 LUN 匯入期間發生中斷。不過、中斷時間不一定比重新啟動資料庫環境和更新 FC 分區所需的時間長、以便將主機 FC 連線能力從外部 LUN 切換至 ONTAP。

此程序可歸納如下：

1. 在外部 LUN 上執行所有 LUN 活動。
2. 將主機 FC 連線重新導向至新的 ONTAP 系統。
3. 觸發匯入程序。
4. 重新探索 LUN。
5. 重新啟動資料庫。

您不需要等待移轉程序完成。一旦開始移轉給定的 LUN、就可以在 ONTAP 上使用、並在資料複製程序繼續進行時提供資料。所有讀取都會傳送到外部 LUN、而且所有寫入都會同步寫入兩個陣列。複製作業非常快速、重新導向 FC 流量的負荷也很小、因此對效能的任何影響都應該是暫時性的、而且最小的。如果有疑慮、您可以延遲重新啟動環境、直到移轉程序完成、匯入關係已刪除為止。

## 關閉資料庫

在本範例中、停止環境的第一步是關閉資料庫。

```
[oracle@host1 bin]$ . oraenv
ORACLE_SID = [oracle] ? FLIDB
The Oracle base remains unchanged with value /orabin
[oracle@host1 bin]$ sqlplus / as sysdba
SQL*Plus: Release 12.1.0.2.0
Copyright (c) 1982, 2014, Oracle. All rights reserved.
Connected to:
Oracle Database 12c Enterprise Edition Release 12.1.0.2.0 - 64bit
Production
With the Partitioning, Automatic Storage Management, OLAP, Advanced
Analytics
and Real Application Testing options
SQL> shutdown immediate;
Database closed.
Database dismounted.
ORACLE instance shut down.
SQL>
```

## 關閉網格服務

其中一個要移轉的 SAN 型檔案系統也包含 Oracle ASM 服務。若要停止基礎 LUN、則需要卸除檔案系統、這也意味著在此檔案系統上停止任何開啟檔案的處理程序。

```
[oracle@host1 bin]$ ./crsctl stop has -f
CRS-2791: Starting shutdown of Oracle High Availability Services-managed
resources on 'host1'
CRS-2673: Attempting to stop 'ora.evmd' on 'host1'
CRS-2673: Attempting to stop 'ora.DATA.dg' on 'host1'
CRS-2673: Attempting to stop 'ora.LISTENER.lsnr' on 'host1'
CRS-2677: Stop of 'ora.DATA.dg' on 'host1' succeeded
CRS-2673: Attempting to stop 'ora.asm' on 'host1'
CRS-2677: Stop of 'ora.LISTENER.lsnr' on 'host1' succeeded
CRS-2677: Stop of 'ora.evmd' on 'host1' succeeded
CRS-2677: Stop of 'ora.asm' on 'host1' succeeded
CRS-2673: Attempting to stop 'ora.cssd' on 'host1'
CRS-2677: Stop of 'ora.cssd' on 'host1' succeeded
CRS-2793: Shutdown of Oracle High Availability Services-managed resources
on 'host1' has completed
CRS-4133: Oracle High Availability Services has been stopped.
[oracle@host1 bin]$
```

## 卸除檔案系統

如果所有程序都關閉、`umount` 作業就會成功。如果權限遭拒、檔案系統上必須有鎖定的程序。◦ `fuser` 命令可協助識別這些程序。

```
[root@host1 ~]# umount /orabin
[root@host1 ~]# umount /backups
```

## 停用 Volume 群組

卸除指定 Volume 群組中的所有檔案系統後、即可停用該 Volume 群組。

```
[root@host1 ~]# vgchange --activate n sanvg
  0 logical volume(s) in volume group "sanvg" now active
[root@host1 ~]#
```

## FC 網路變更

現在可以更新 FC 區域、以移除主機對外部陣列的所有存取權、並建立對 ONTAP 的存取權。

## 開始匯入程序

若要啟動 LUN 匯入程序、請執行 `lun import start` 命令。

```
Cluster01::lun import*> lun import start -vserver vserver1 -path
/vol/new_asm/LUN0
Cluster01::lun import*> lun import start -vserver vserver1 -path
/vol/new_asm/LUN1
...
Cluster01::lun import*> lun import start -vserver vserver1 -path
/vol/new_lvm/LUN8
Cluster01::lun import*> lun import start -vserver vserver1 -path
/vol/new_lvm/LUN9
Cluster01::lun import*>
```

## 監控匯入進度

您可以使用監控匯入作業 `lun import show` 命令。如下所示、目前正在匯入所有 20 個 LUN、這表示即使資料複製作業仍在進行中、仍可透過 ONTAP 存取資料。

```

Cluster01::lun import*> lun import show -fields path,percent-complete
vserver    foreign-disk path                               percent-complete
-----
vserver1   800DT$HuVWB/ /vol/new_asm/LUN4 5
vserver1   800DT$HuVWBW /vol/new_asm/LUN0 5
vserver1   800DT$HuVWBX /vol/new_asm/LUN1 6
vserver1   800DT$HuVWBZ /vol/new_asm/LUN2 6
vserver1   800DT$HuVWBa /vol/new_asm/LUN3 5
vserver1   800DT$HuVWBb /vol/new_asm/LUN5 4
vserver1   800DT$HuVWBc /vol/new_asm/LUN6 4
vserver1   800DT$HuVWBd /vol/new_asm/LUN7 4
vserver1   800DT$HuVWBd /vol/new_asm/LUN8 4
vserver1   800DT$HuVWBe /vol/new_asm/LUN9 4
vserver1   800DT$HuVWBf /vol/new_lvm/LUN0 5
vserver1   800DT$HuVWBg /vol/new_lvm/LUN1 4
vserver1   800DT$HuVWBh /vol/new_lvm/LUN2 4
vserver1   800DT$HuVWBh /vol/new_lvm/LUN3 3
vserver1   800DT$HuVWBj /vol/new_lvm/LUN4 3
vserver1   800DT$HuVWBk /vol/new_lvm/LUN5 3
vserver1   800DT$HuVWBk /vol/new_lvm/LUN6 4
vserver1   800DT$HuVWBm /vol/new_lvm/LUN7 3
vserver1   800DT$HuVWBn /vol/new_lvm/LUN8 2
vserver1   800DT$HuVWBn /vol/new_lvm/LUN9 2
20 entries were displayed.

```

如果您需要離線程序、請延遲重新探索或重新啟動服務、直到 `lun import show` 命令表示所有移轉均已成功完成。接著您可以依照中所述、完成移轉程序 "外部 LUN 匯入：完成"。

如果您需要線上移轉、請繼續在新的主目錄中重新探索 LUN、並啟動服務。

### 掃描 SCSI 裝置變更

在大多數情況下、重新探索新 LUN 最簡單的選項是重新啟動主機。這樣做會自動移除舊的過時裝置、正確探索所有新的 LUN、並建置相關的裝置、例如多重路徑裝置。以下範例顯示出完全線上的示範程序。

注意：在重新啟動主機之前、請確定中的所有項目都已存在 `/etc/fstab` 這項參照移轉的 SAN 資源會被註解出來。如果未執行此操作、且 LUN 存取有問題、作業系統可能無法開機。這種情況不會損害資料。不過、開機進入救援模式或類似模式並修正可能非常不方便 `/etc/fstab` 如此一來、就能開機作業系統以進行疑難排解。

本範例所使用 Linux 版本上的 LUN 可與重新掃描 `rescan-scsi-bus.sh` 命令。如果命令成功、每個 LUN 路徑都會出現在輸出中。輸出可能很難解譯、但如果分區和 `igroup` 組態正確、許多 LUN 應該會顯示為包含 NETAPP 廠商字串。

```

[root@host1 /]# rescan-scsi-bus.sh
Scanning SCSI subsystem for new devices
Scanning host 0 for SCSI target IDs 0 1 2 3 4 5 6 7, all LUNs
  Scanning for device 0 2 0 0 ...
OLD: Host: scsi0 Channel: 02 Id: 00 Lun: 00
      Vendor: LSI      Model: RAID SAS 6G 0/1  Rev: 2.13
      Type:   Direct-Access                    ANSI SCSI revision: 05
Scanning host 1 for SCSI target IDs 0 1 2 3 4 5 6 7, all LUNs
  Scanning for device 1 0 0 0 ...
OLD: Host: scsi1 Channel: 00 Id: 00 Lun: 00
      Vendor: Optiarc  Model: DVD RW AD-7760H  Rev: 1.41
      Type:   CD-ROM                      ANSI SCSI revision: 05
Scanning host 2 for SCSI target IDs 0 1 2 3 4 5 6 7, all LUNs
Scanning host 3 for SCSI target IDs 0 1 2 3 4 5 6 7, all LUNs
Scanning host 4 for SCSI target IDs 0 1 2 3 4 5 6 7, all LUNs
Scanning host 5 for SCSI target IDs 0 1 2 3 4 5 6 7, all LUNs
Scanning host 6 for SCSI target IDs 0 1 2 3 4 5 6 7, all LUNs
Scanning host 7 for all SCSI target IDs, all LUNs
  Scanning for device 7 0 0 10 ...
OLD: Host: scsi7 Channel: 00 Id: 00 Lun: 10
      Vendor: NETAPP   Model: LUN C-Mode          Rev: 8300
      Type:   Direct-Access                    ANSI SCSI revision: 05
  Scanning for device 7 0 0 11 ...
OLD: Host: scsi7 Channel: 00 Id: 00 Lun: 11
      Vendor: NETAPP   Model: LUN C-Mode          Rev: 8300
      Type:   Direct-Access                    ANSI SCSI revision: 05
  Scanning for device 7 0 0 12 ...
...
OLD: Host: scsi9 Channel: 00 Id: 01 Lun: 18
      Vendor: NETAPP   Model: LUN C-Mode          Rev: 8300
      Type:   Direct-Access                    ANSI SCSI revision: 05
  Scanning for device 9 0 1 19 ...
OLD: Host: scsi9 Channel: 00 Id: 01 Lun: 19
      Vendor: NETAPP   Model: LUN C-Mode          Rev: 8300
      Type:   Direct-Access                    ANSI SCSI revision: 05
0 new or changed device(s) found.
0 remapped or resized device(s) found.
0 device(s) removed.

```

## 檢查多重路徑裝置

LUN 探索程序也會觸發多重路徑裝置的重新開發、但已知 Linux 多重路徑驅動程式偶爾會發生問題。的輸出 `multipath - ll` 應檢查以驗證輸出是否如預期。例如、下列輸出顯示與相關的多重路徑裝置 NETAPP 廠商字串。每個裝置有四條路徑、其中兩條優先順序為 50、兩條優先順序為 10。雖然確切的輸出可能會因 Linux 的不同版本而有所不同、但此輸出的外觀與預期相同。





請參閱您用來驗證的 Linux 版本的主機公用程式文件 `/etc/multipath.conf` 設定正確。

```
[root@host1 /]# multipath -ll
3600a098038303558735d493762504b36 dm-5 NETAPP ,LUN C-Mode
size=10G features='4 queue_if_no_path pg_init_retries 50
retain_attached_hw_handle' hwhandler='1 alua' wp=rw
|+- policy='service-time 0' prio=50 status=active
| |- 7:0:1:4 sdat 66:208 active ready running
| `-- 9:0:1:4 sdbn 68:16 active ready running
`-+- policy='service-time 0' prio=10 status=enabled
   |- 7:0:0:4 sdf 8:80 active ready running
   `-- 9:0:0:4 sdz 65:144 active ready running
3600a098038303558735d493762504b2d dm-10 NETAPP ,LUN C-Mode
size=10G features='4 queue_if_no_path pg_init_retries 50
retain_attached_hw_handle' hwhandler='1 alua' wp=rw
|+- policy='service-time 0' prio=50 status=active
| |- 7:0:1:8 sdax 67:16 active ready running
| `-- 9:0:1:8 sdbr 68:80 active ready running
`-+- policy='service-time 0' prio=10 status=enabled
   |- 7:0:0:8 sdj 8:144 active ready running
   `-- 9:0:0:8 sdad 65:208 active ready running
...
3600a098038303558735d493762504b37 dm-8 NETAPP ,LUN C-Mode
size=10G features='4 queue_if_no_path pg_init_retries 50
retain_attached_hw_handle' hwhandler='1 alua' wp=rw
|+- policy='service-time 0' prio=50 status=active
| |- 7:0:1:5 sdau 66:224 active ready running
| `-- 9:0:1:5 sdbo 68:32 active ready running
`-+- policy='service-time 0' prio=10 status=enabled
   |- 7:0:0:5 sdg 8:96 active ready running
   `-- 9:0:0:5 sdaa 65:160 active ready running
3600a098038303558735d493762504b4b dm-22 NETAPP ,LUN C-Mode
size=10G features='4 queue_if_no_path pg_init_retries 50
retain_attached_hw_handle' hwhandler='1 alua' wp=rw
|+- policy='service-time 0' prio=50 status=active
| |- 7:0:1:19 sdbi 67:192 active ready running
| `-- 9:0:1:19 sdcc 69:0 active ready running
`-+- policy='service-time 0' prio=10 status=enabled
   |- 7:0:0:19 sdu 65:64 active ready running
   `-- 9:0:0:19 sdao 66:128 active ready running
```

## 重新啟動 LVM Volume 群組

如果正確探索到 LVM LUN、則會發現 `vgchange --activate y` 命令應該成功。這是邏輯 Volume Manager 的價值範例。由於磁碟區群組中繼資料是寫入 LUN 本身、因此 LUN 的 WWN 變更甚至是序列號都不重要。

作業系統掃描 LUN、並發現 LUN 上寫入的少量資料、可將其識別為屬於的實體磁碟區 sanvg volumegroup。然後、它會建置所有必要的裝置。只需重新啟動 Volume 群組即可。

```
[root@host1 ~]# vgchange --activate y sanvg
Found duplicate PV fpCzdLTuKfy2xDZjailNliJh3TjLUBiT: using
/dev/mapper/3600a098038303558735d493762504b46 not /dev/sdp
Using duplicate PV /dev/mapper/3600a098038303558735d493762504b46 from
subsystem DM, ignoring /dev/sdp
2 logical volume(s) in volume group "sanvg" now active
```

### 重新掛載檔案系統

磁碟區群組重新啟動後、檔案系統可以裝入、所有原始資料均完整無缺。如前所述、即使資料複寫仍在後端群組中作用中、檔案系統仍可完全運作。

```
[root@host1 ~]# mount /orabin
[root@host1 ~]# mount /backups
[root@host1 ~]# df -k
```

Filesystem	1K-blocks	Used	Available	Use%
Mounted on				
/dev/mapper/rhel-root	52403200	8837100	43566100	17% /
devtmpfs	65882776	0	65882776	0% /dev
tmpfs	6291456	84	6291372	1%
/dev/shm				
tmpfs	65898668	9884	65888784	1% /run
tmpfs	65898668	0	65898668	0%
/sys/fs/cgroup				
/dev/sda1	505580	224828	280752	45% /boot
fas8060-nfs-public:/install	199229440	119368256	79861184	60%
/install				
fas8040-nfs-routable:/snapomatic	9961472	30528	9930944	1%
/snapomatic				
tmpfs	13179736	16	13179720	1%
/run/user/42				
tmpfs	13179736	0	13179736	0%
/run/user/0				
/dev/mapper/sanvg-lvorabin	20961280	12357456	8603824	59%
/orabin				
/dev/mapper/sanvg-lvbackups	73364480	62947536	10416944	86%
/backups				

### 重新掃描 ASM 設備

重新掃描 SCSI 裝置時、應已重新探索 ASMLib 裝置。重新探索可透過重新啟動 ASMLib、然後掃描磁碟來線上驗證。



此步驟僅與使用 ASMLib 的 ASM 組態相關。

注意：若未使用 ASMLib、請使用 `/dev/mapper` 裝置應已自動重新建立。不過、權限可能不正確。在 ASMLib 不存在的情況下、您必須為基礎裝置設定特殊權限。這樣做通常是透過中的特殊項目來完成 `/etc/multipath.conf` 或 `udev` 規則、或可能同時在兩個規則集中。這些檔案可能需要更新、以反映環境中的 WWN 或序號變更、以確保 ASM 裝置仍擁有正確的權限。

在此範例中、重新啟動 ASMLib 並掃描磁碟時、會顯示與原始環境相同的 10 個 ASM LUN。

```
[root@host1 ~]# oracleasm exit
Unmounting ASMLib driver filesystem: /dev/oracleasm
Unloading module "oracleasm": oracleasm
[root@host1 ~]# oracleasm init
Loading module "oracleasm": oracleasm
Configuring "oracleasm" to use device physical block size
Mounting ASMLib driver filesystem: /dev/oracleasm
[root@host1 ~]# oracleasm scandisks
Reloading disk partitions: done
Cleaning any stale ASM disks...
Scanning system for ASM disks...
Instantiating disk "ASM0"
Instantiating disk "ASM1"
Instantiating disk "ASM2"
Instantiating disk "ASM3"
Instantiating disk "ASM4"
Instantiating disk "ASM5"
Instantiating disk "ASM6"
Instantiating disk "ASM7"
Instantiating disk "ASM8"
Instantiating disk "ASM9"
```

### 重新啟動網絡服務

現在、LVM 和 ASM 裝置已上線且可供使用、可以重新啟動網絡服務。

```
[root@host1 ~]# cd /orabin/product/12.1.0/grid/bin
[root@host1 bin]# ./crsctl start has
```

### 重新啟動資料庫

網絡服務重新啟動後、即可啟動資料庫。在嘗試啟動資料庫之前、可能需要等待幾分鐘、ASM 服務才能完全可用。

```
[root@host1 bin]# su - oracle
[oracle@host1 ~]$ . oraenv
ORACLE_SID = [oracle] ? FLIDB
The Oracle base has been set to /orabin
[oracle@host1 ~]$ sqlplus / as sysdba
SQL*Plus: Release 12.1.0.2.0
Copyright (c) 1982, 2014, Oracle. All rights reserved.
Connected to an idle instance.
SQL> startup
ORACLE instance started.
Total System Global Area 3221225472 bytes
Fixed Size 4502416 bytes
Variable Size 1207962736 bytes
Database Buffers 1996488704 bytes
Redo Buffers 12271616 bytes
Database mounted.
Database opened.
SQL>
```

**Oracle 移轉與 FLI**：完成

從主機的角度來看、移轉已完成、但在匯入關係刪除之前、仍會從外部陣列提供 I/O。

刪除關係之前、您必須確認所有 LUN 的移轉程序已完成。

```

Cluster01::*> lun import show -vserver vserver1 -fields foreign-
disk,path,operational-state
vserver    foreign-disk path                operational-state
-----
vserver1 800DT$HuVWB/ /vol/new_asm/LUN4 completed
vserver1 800DT$HuVWBW /vol/new_asm/LUN0 completed
vserver1 800DT$HuVWBX /vol/new_asm/LUN1 completed
vserver1 800DT$HuVWBZ /vol/new_asm/LUN2 completed
vserver1 800DT$HuVWBa /vol/new_asm/LUN5 completed
vserver1 800DT$HuVWBb /vol/new_asm/LUN6 completed
vserver1 800DT$HuVWBc /vol/new_asm/LUN7 completed
vserver1 800DT$HuVWBd /vol/new_asm/LUN8 completed
vserver1 800DT$HuVWB e /vol/new_asm/LUN9 completed
vserver1 800DT$HuVWBf /vol/new_lvm/LUN0 completed
vserver1 800DT$HuVWBg /vol/new_lvm/LUN1 completed
vserver1 800DT$HuVWBh /vol/new_lvm/LUN2 completed
vserver1 800DT$HuVWB i /vol/new_lvm/LUN3 completed
vserver1 800DT$HuVWBj /vol/new_lvm/LUN4 completed
vserver1 800DT$HuVWBk /vol/new_lvm/LUN5 completed
vserver1 800DT$HuVWB l /vol/new_lvm/LUN6 completed
vserver1 800DT$HuVWBm /vol/new_lvm/LUN7 completed
vserver1 800DT$HuVWBn /vol/new_lvm/LUN8 completed
vserver1 800DT$HuVWB o /vol/new_lvm/LUN9 completed
20 entries were displayed.

```

## 刪除匯入關係

移轉程序完成後、請刪除移轉關係。完成後、I/O 將由 ONTAP 上的磁碟機獨家提供。

```

Cluster01::*> lun import delete -vserver vserver1 -path /vol/new_asm/LUN0
Cluster01::*> lun import delete -vserver vserver1 -path /vol/new_asm/LUN1
...
Cluster01::*> lun import delete -vserver vserver1 -path /vol/new_lvm/LUN8
Cluster01::*> lun import delete -vserver vserver1 -path /vol/new_lvm/LUN9

```

## 取消註冊外部 LUN

最後、修改磁碟以移除 is-foreign 指定。

```

Cluster01::*> storage disk modify {-serial-number 800DT$HuVWBW} -is
-foreign false
Cluster01::*> storage disk modify {-serial-number 800DT$HuVWBX} -is
-foreign false
...
Cluster01::*> storage disk modify {-serial-number 800DT$HuVWBn} -is
-foreign false
Cluster01::*> storage disk modify {-serial-number 800DT$HuVWBo} -is
-foreign false
Cluster01::*>

```

使用 **FLI** 進行 **Oracle** 移轉：傳輸協定轉換

變更用於存取 LUN 的傳輸協定是常見的需求。

在某些情況下、它是將資料移轉至雲端的整體策略的一部分。TCP/IP 是雲端的傳輸協定、從 FC 變更為 iSCSI 可讓您更輕鬆地移轉至各種雲端環境。在其他情況下、iSCSI 可能需要善用 IP SAN 降低的成本。有時候、移轉可能會使用不同的傳輸協定作為臨時措施。例如、如果外部陣列和 ONTAP 型 LUN 無法共存於同一個 HBA 上、您可以使用足夠長的 iSCSI LUN、從舊陣列複製資料。然後、您可以在從系統移除舊 LUN 之後、將其轉換回 FC。

下列程序示範從 FC 轉換至 iSCSI 的過程、但整體原則適用於從 iSCSI 轉換至 FC 的反轉過程。

### 安裝 **iSCSI** 啟動器

大多數作業系統預設都包含軟體 iSCSI 啟動器、但如果不包含軟體 iSCSI 啟動器、則可輕鬆安裝。

```

[root@host1 /]# yum install -y iscsi-initiator-utils
Loaded plugins: langpacks, product-id, search-disabled-repos,
subscription-
           : manager
Resolving Dependencies
--> Running transaction check
---> Package iscsi-initiator-utils.x86_64 0:6.2.0.873-32.el7 will be
updated
--> Processing Dependency: iscsi-initiator-utils = 6.2.0.873-32.el7 for
package: iscsi-initiator-utils-iscsiuio-6.2.0.873-32.el7.x86_64
---> Package iscsi-initiator-utils.x86_64 0:6.2.0.873-32.0.2.el7 will be
an update
--> Running transaction check
---> Package iscsi-initiator-utils-iscsiuio.x86_64 0:6.2.0.873-32.el7 will
be updated
---> Package iscsi-initiator-utils-iscsiuio.x86_64 0:6.2.0.873-32.0.2.el7
will be an update
--> Finished Dependency Resolution
Dependencies Resolved

```

```

=====
===
Package                Arch    Version                Repository
Size
=====
===
Updating:
iscsi-initiator-utils  x86_64 6.2.0.873-32.0.2.el7 ol7_latest 416
k
Updating for dependencies:
iscsi-initiator-utils-iscsiuio x86_64 6.2.0.873-32.0.2.el7 ol7_latest 84
k
Transaction Summary
=====
===
Upgrade 1 Package (+1 Dependent package)
Total download size: 501 k
Downloading packages:
No Presto metadata available for ol7_latest
(1/2): iscsi-initiator-utils-6.2.0.873-32.0.2.el7.x86_6 | 416 kB 00:00
(2/2): iscsi-initiator-utils-iscsiuio-6.2.0.873-32.0.2. | 84 kB 00:00
-----
---
Total                2.8 MB/s | 501 kB
00:00Cluster01
Running transaction check
Running transaction test
Transaction test succeeded
Running transaction
  Updating   : iscsi-initiator-utils-iscsiuio-6.2.0.873-32.0.2.el7.x86
1/4
  Updating   : iscsi-initiator-utils-6.2.0.873-32.0.2.el7.x86_64
2/4
  Cleanup    : iscsi-initiator-utils-iscsiuio-6.2.0.873-32.el7.x86_64
3/4
  Cleanup    : iscsi-initiator-utils-6.2.0.873-32.el7.x86_64
4/4
rhel-7-server-eus-rpms/7Server/x86_64/productid | 1.7 kB 00:00
rhel-7-server-rpms/7Server/x86_64/productid | 1.7 kB 00:00
  Verifying  : iscsi-initiator-utils-6.2.0.873-32.0.2.el7.x86_64
1/4
  Verifying  : iscsi-initiator-utils-iscsiuio-6.2.0.873-32.0.2.el7.x86
2/4
  Verifying  : iscsi-initiator-utils-iscsiuio-6.2.0.873-32.el7.x86_64
3/4
  Verifying  : iscsi-initiator-utils-6.2.0.873-32.el7.x86_64

```

4/4

Updated:

```
iscsi-initiator-utils.x86_64 0:6.2.0.873-32.0.2.e17
```

Dependency Updated:

```
iscsi-initiator-utils-iscsiuio.x86_64 0:6.2.0.873-32.0.2.e17
```

Complete!

```
[root@host1 ~]#
```

## 識別 iSCSI 啟動器名稱

在安裝過程中會產生唯一的 iSCSI 啟動器名稱。在 Linux 上、它位於 `/etc/iscsi/initiatorname.iscsi` 檔案：此名稱用於識別 IP SAN 上的主機。

```
[root@host1 ~]# cat /etc/iscsi/initiatorname.iscsi
InitiatorName=iqn.1992-05.com.redhat:497bd66ca0
```

## 建立新的啟動器群組

啟動器群組 (igroup) 是 ONTAP LUN 遮罩架構的一部分。除非先授予主機存取權、否則無法存取新建立的 LUN。此步驟的完成方法是建立一個 igroup、列出需要存取的 FC WWN 或 iSCSI 啟動器名稱。

在此範例中、會建立包含 Linux 主機 iSCSI 啟動器的 igroup。

```
Cluster01::*> igroup create -igroup linuxiscsi -protocol iscsi -ostype
linux -initiator iqn.1994-05.com.redhat:497bd66ca0
```

## 關閉環境

變更 LUN 傳輸協定之前、必須完全禁用 LUN。任何要轉換的 LUN 上的資料庫都必須關機、檔案系統必須卸載、而且必須停用磁碟區群組。使用 ASM 時、請確定已卸除 ASM 磁碟群組、並關閉所有網絡服務。

## 從 FC 網路取消對應 LUN

LUN 完全禁用後、請從原始 FC igroup 移除對應。

```
Cluster01::*> lun unmap -vserver vserver1 -path /vol/new_asm/LUN0 -igroup
linuxhost
Cluster01::*> lun unmap -vserver vserver1 -path /vol/new_asm/LUN1 -igroup
linuxhost
...
Cluster01::*> lun unmap -vserver vserver1 -path /vol/new_lvm/LUN8 -igroup
linuxhost
Cluster01::*> lun unmap -vserver vserver1 -path /vol/new_lvm/LUN9 -igroup
linuxhost
```



## 將 LUN 重新對應至 IP 網路

將每個 LUN 的存取權授予新的 iSCSI 型啟動器群組。

```
Cluster01::*> lun map -vserver vserver1 -path /vol/new_asm/LUN0 -igroup linuxiscsi
Cluster01::*> lun map -vserver vserver1 -path /vol/new_asm/LUN1 -igroup linuxiscsi
...
Cluster01::*> lun map -vserver vserver1 -path /vol/new_lvm/LUN8 -igroup linuxiscsi
Cluster01::*> lun map -vserver vserver1 -path /vol/new_lvm/LUN9 -igroup linuxiscsi
Cluster01::*>
```

## 探索 iSCSI 目標

iSCSI 探索分為兩個階段。第一是探索目標、這與探索 LUN 不同。iscsiadm 下列命令會探查指定的入口網站群組 -p argument 並儲存提供 iSCSI 服務的所有 IP 位址和連接埠清單。在這種情況下、預設連接埠 3260 上有四個 iSCSI 服務的 IP 位址。



如果無法到達任何目標 IP 位址、此命令可能需要幾分鐘的時間才能完成。

```
[root@host1 ~]# iscsiadm -m discovery -t st -p fas8060-iscsi-public1
10.63.147.197:3260,1033 iqn.1992-
08.com.netapp:sn.807615e9ef6111e5a5ae90e2ba5b9464:vs.3
10.63.147.198:3260,1034 iqn.1992-
08.com.netapp:sn.807615e9ef6111e5a5ae90e2ba5b9464:vs.3
172.20.108.203:3260,1030 iqn.1992-
08.com.netapp:sn.807615e9ef6111e5a5ae90e2ba5b9464:vs.3
172.20.108.202:3260,1029 iqn.1992-
08.com.netapp:sn.807615e9ef6111e5a5ae90e2ba5b9464:vs.3
```

## 探索 iSCSI LUN

發現 iSCSI 目標後、請重新啟動 iSCSI 服務以探索可用的 iSCSI LUN、並建置相關裝置、例如多重路徑或 ASMLib 裝置。

```
[root@host1 ~]# service iscsi restart
Redirecting to /bin/systemctl restart iscsi.service
```

## 重新啟動環境

重新啟動 Volume 群組、重新掛載檔案系統、重新啟動 RAC 服務等、以重新啟動環境。為了預防這種情況、

NetApp 建議您在轉換程序完成後重新啟動伺服器、以確保所有組態檔案都正確無誤、並移除所有過時的裝置。

注意：在重新啟動主機之前、請確定中的所有項目都已存在 `/etc/fstab` 這項參照移轉的 SAN 資源會被註解出來。如果未執行此步驟、且 LUN 存取有問題、則可能是無法開機的作業系統。此問題不會損壞資料。不過、開機進入救援模式或類似模式進行修正可能會非常不方便 `/etc/fstab` 這樣就能啟動作業系統、開始進行疑難排解工作。

## Oracle 移轉程序範例指令碼

提供的指令碼是如何為各種作業系統和資料庫工作撰寫指令碼的範例。這些都是依現狀供應。如果特定程序需要支援、請聯絡 NetApp 或 NetApp 經銷商。

### 資料庫關機

下列 Perl 指令碼會採用 Oracle SID 的單一引數、並關閉資料庫。它可以以 Oracle 使用者或 root 身分執行。

```

#!/usr/bin/perl
use strict;
use warnings;
my $oraclesid=$ARGV[0];
my $oracleuser='oracle';
my @out;
my $uid=$<;
if ($uid == 0) {
@out=`su - $oracleuser -c '. oraenv << EOF1
77 Migration of Oracle Databases to NetApp Storage Systems © 2021 NetApp,
Inc. All rights reserved
$oraclesid
EOF1
sqlplus / as sysdba << EOF2
shutdown immediate;
EOF2
`
`;}
else {
@out=`. oraenv << EOF1
$oraclesid
EOF4
sqlplus / as sysdba << EOF2
shutdown immediate;
EOF2
`;};
print @out;
if ("@out" =~ /ORACLE instance shut down/) {
print "$oraclesid shut down\n";
exit 0;}
elsif ("@out" =~ /Connected to an idle instance/) {
print "$oraclesid already shut down\n";
exit 0;}
else {
print "$oraclesid failed to shut down\n";
exit 1;}

```

## 資料庫啟動

下列 Perl 指令碼會採用 Oracle SID 的單一引數、並關閉資料庫。它可以以 Oracle 使用者或 root 身分執行。

```

#!/usr/bin/perl
use strict;
use warnings;
my $oraclesid=$ARGV[0];
my $oracleuser='oracle';
my @out;
my $uid=$<;
if ($uid == 0) {
@out=`su - $oracleuser -c '. oraenv << EOF1
$oraclesid
EOF1
sqlplus / as sysdba << EOF2
startup;
EOF2
`
`;}
else {
@out=`. oraenv << EOF3
$oraclesid
EOF1
sqlplus / as sysdba << EOF2
startup;
EOF2
`;};
print @out;
if ("@out" =~ /Database opened/) {
print "$oraclesid started\n";
exit 0;}
elsif ("@out" =~ /cannot start already-running ORACLE/) {
print "$oraclesid already started\n";
exit 1;}
else {
78 Migration of Oracle Databases to NetApp Storage Systems © 2021 NetApp,
Inc. All rights reserved
print "$oraclesid failed to start\n";
exit 1;}

```

### 將檔案系統轉換為唯讀

下列指令碼會採用檔案系統引數、並嘗試將其卸除並重新掛載為唯讀。在移轉過程中、這樣做非常有用、因為必須將檔案系統保留在可複寫資料的位置、但必須保護其免於意外損壞。

```

#!/usr/bin/perl
use strict;
#use warnings;
my $filesystem=$ARGV[0];
my @out=`umount '$filesystem'`;
if ($? == 0) {
    print "$filesystem unmounted\n";
    @out = `mount -o ro '$filesystem'`;
    if ($? == 0) {
        print "$filesystem mounted read-only\n";
        exit 0;}}
else {
    print "Unable to unmount $filesystem\n";
    exit 1;}
print @out;

```

## 取代檔案系統

下列指令碼範例用於將一個檔案系統取代為另一個檔案系統。因為它編輯了 /etc/fstab 文件，所以它必須以 root 身份運行。它接受新舊檔案系統的單一逗號分隔引數。

1. 若要取代檔案系統、請執行下列指令碼：

```

#!/usr/bin/perl
use strict;
#use warnings;
my $oldfs;
my $newfs;
my @oldfstab;
my @newfstab;
my $source;
my $mountpoint;
my $leftover;
my $oldfstabentry='';
my $newfstabentry='';
my $migratedfstabentry='';
($oldfs, $newfs) = split (',', $ARGV[0]);
open(my $filehandle, '<', '/etc/fstab') or die "Could not open
/etc/fstab\n";
while (my $line = <$filehandle>) {
    chomp $line;
    ($source, $mountpoint, $leftover) = split(/[ , ]/, $line, 3);
    if ($mountpoint eq $oldfs) {
        $oldfstabentry = "#Removed by swap script $source $oldfs $leftover";}
    elsif ($mountpoint eq $newfs) {

```

```

$newfstabentry = "#Removed by swap script $source $newfs $leftover";
$migratedfstabentry = "$source $oldfs $leftover";
else {
push (@newfstab, "$line\n")}
79 Migration of Oracle Databases to NetApp Storage Systems © 2021
NetApp, Inc. All rights reserved
push (@newfstab, "$oldfstabentry\n");
push (@newfstab, "$newfstabentry\n");
push (@newfstab, "$migratedfstabentry\n");
close($filehandle);
if ($oldfstabentry eq ''){
die "Could not find $oldfs in /etc/fstab\n";}
if ($newfstabentry eq ''){
die "Could not find $newfs in /etc/fstab\n";}
my @out=`umount '$newfs'`;
if ($? == 0) {
print "$newfs unmounted\n";}
else {
print "Unable to unmount $newfs\n";
exit 1;}
@out=`umount '$oldfs'`;
if ($? == 0) {
print "$oldfs unmounted\n";}
else {
print "Unable to unmount $oldfs\n";
exit 1;}
system("cp /etc/fstab /etc/fstab.bak");
open ($filehandle, ">", '/etc/fstab') or die "Could not open /etc/fstab
for writing\n";
for my $line (@newfstab) {
print $filehandle $line;}
close($filehandle);
@out=`mount '$oldfs'`;
if ($? == 0) {
print "Mounted updated $oldfs\n";
exit 0;}
else{
print "Unable to mount updated $oldfs\n";
exit 1;}
exit 0;

```

以本指令碼的使用範例為例、假設中的資料 /oradata 移轉至 /neworadata 和 /logs 移轉至 /newlogs。執行此工作最簡單的方法之一、就是使用簡單的檔案複製作業、將新裝置重新放置回原始安裝點。

2. 假設舊的和新的檔案系統存在於中 /etc/fstab 檔案如下：

```

cluster01:/vol_oradata /oradata nfs rw,bg,vers=3,rsize=65536,wsiz=65536
0 0
cluster01:/vol_logs /logs nfs rw,bg,vers=3,rsize=65536,wsiz=65536 0 0
cluster01:/vol_neworadata /neworadata nfs
rw,bg,vers=3,rsize=65536,wsiz=65536 0 0
cluster01:/vol_newlogs /newlogs nfs rw,bg,vers=3,rsize=65536,wsiz=65536
0 0

```

3. 執行時、此指令碼會卸載目前的檔案系統、並以新的：

```

[root@jpsc3 scripts]# ./swap.fs.pl /oradata,/neworadata
/neworadata unmounted
/oradata unmounted
Mounted updated /oradata
[root@jpsc3 scripts]# ./swap.fs.pl /logs,/newlogs
/newlogs unmounted
/logs unmounted
Mounted updated /logs

```

4. 指令碼也會更新 /etc/fstab 請據此歸檔。在此處所示範例中、包含下列變更：

```

#Removed by swap script cluster01:/vol_oradata /oradata nfs
rw,bg,vers=3,rsize=65536,wsiz=65536 0 0
#Removed by swap script cluster01:/vol_neworadata /neworadata nfs
rw,bg,vers=3,rsize=65536,wsiz=65536 0 0
cluster01:/vol_neworadata /oradata nfs
rw,bg,vers=3,rsize=65536,wsiz=65536 0 0
#Removed by swap script cluster01:/vol_logs /logs nfs
rw,bg,vers=3,rsize=65536,wsiz=65536 0 0
#Removed by swap script cluster01:/vol_newlogs /newlogs nfs
rw,bg,vers=3,rsize=65536,wsiz=65536 0 0
cluster01:/vol_newlogs /logs nfs rw,bg,vers=3,rsize=65536,wsiz=65536 0
0

```

## 自動化資料庫移轉

此範例示範如何使用關機、啟動及檔案系統置換指令碼來完全自動化移轉。

```

#!/usr/bin/perl
use strict;
#use warnings;
my $oraclesid=$ARGV[0];

```

```

my @oldfs;
my @newfs;
my $x=1;
while ($x < scalar(@ARGV)) {
    ($oldfs[$x-1], $newfs[$x-1]) = split (',', $ARGV[$x]);
    $x+=1;}
my @out=`./dbshut.pl '$oraclesid'`;
print @out;
if ($? ne 0) {
    print "Failed to shut down database\n";
    exit 0;}
$x=0;
while ($x < scalar(@oldfs)) {
    my @out=`./mk.fs.readonly.pl '$oldfs[$x]'`;
    if ($? ne 0) {
        print "Failed to make filesystem $oldfs[$x] readonly\n";
        exit 0;}
    $x+=1;}
$x=0;
while ($x < scalar(@oldfs)) {
    my @out=`rsync -rlpogt --stats --progress --exclude='.snapshot'
'$oldfs[$x]/' '/$newfs[$x]/'`;
    print @out;
    if ($? ne 0) {
        print "Failed to copy filesystem $oldfs[$x] to $newfs[$x]\n";
        exit 0;}
    else {
        print "Succesfully replicated filesystem $oldfs[$x] to
$newfs[$x]\n";}
    $x+=1;}
$x=0;
while ($x < scalar(@oldfs)) {
    print "swap $x $oldfs[$x] $newfs[$x]\n";
    my @out=`./swap.fs.pl '$oldfs[$x],$newfs[$x]'`;
    print @out;
    if ($? ne 0) {
        print "Failed to swap filesystem $oldfs[$x] for $newfs[$x]\n";
        exit 1;}
    else {
        print "Swapped filesystem $oldfs[$x] for $newfs[$x]\n";}
    $x+=1;}
my @out=`./dbstart.pl '$oraclesid'`;
print @out;

```



## 顯示檔案位置

此指令碼會收集許多重要的資料庫參數、並以易讀的格式列印。此指令碼在檢閱資料配置時非常實用。此外、指令碼也可以修改以供其他用途使用。

```
#!/usr/bin/perl
#use strict;
#use warnings;
my $oraclesid=$ARGV[0];
my $oracleuser='oracle';
my @out;
sub dosql{
    my $command = @_ [0];
    my @lines;
    my $uid=$<;
    if ($uid == 0) {
        @lines=`su - $oracleuser -c "export ORAENV_ASK=NO;export
ORACLE_SID=$oraclesid;. oraenv -s << EOF1
EOF1
sqlplus -S / as sysdba << EOF2
set heading off
$command
EOF2
"
        `; }
    else {
        $command=~s/\\\\\\\\/\\/g;
        @lines=`export ORAENV_ASK=NO;export ORACLE_SID=$oraclesid;. oraenv
-s << EOF1
EOF1
sqlplus -S / as sysdba << EOF2
set heading off
$command
EOF2
        `; };
    return @lines;
}
print "\n";
@out=dosql('select name from v\\\\\\\\$datafile;');
print "$oraclesid datafiles:\n";
for $line (@out) {
    chomp($line);
    if (length($line)>0) {print "$line\n";}}
print "\n";
@out=dosql('select member from v\\\\\\\\$logfile;');
print "$oraclesid redo logs:\n";
for $line (@out) {
```

```

        chomp($line);
        if (length($line)>0) {print "$line\n";}}
print "\n";
@out=dosql('select name from v\\\\\\\\$tempfile;');
print "$oraclesid temp datafiles:\n";
for $line (@out) {
    chomp($line);
    if (length($line)>0) {print "$line\n";}}
print "\n";
@out=dosql('show parameter spfile;');
print "$oraclesid spfile\n";
for $line (@out) {
    chomp($line);
    if (length($line)>0) {print "$line\n";}}
print "\n";
@out=dosql('select name||\'' \'|value from v\\\\\\\\$parameter where
isdefault=\'FALSE\';');
print "$oraclesid key parameters\n";
for $line (@out) {
    chomp($line);
    if ($line =~ /control_files/) {print "$line\n";}
    if ($line =~ /db_create/) {print "$line\n";}
    if ($line =~ /db_file_name_convert/) {print "$line\n";}
    if ($line =~ /log_archive_dest/) {print "$line\n";}}
    if ($line =~ /log_file_name_convert/) {print "$line\n";}
    if ($line =~ /pdb_file_name_convert/) {print "$line\n";}
    if ($line =~ /spfile/) {print "$line\n";}
print "\n";

```

## ASM 移轉清理

```

#!/usr/bin/perl
#use strict;
#use warnings;
my $oraclesid=$ARGV[0];
my $oracleuser='oracle';
my @out;
sub dosql{
    my $command = @_ [0];
    my @lines;
    my $uid=$<;
    if ($uid == 0) {
        @lines=`su - $oracleuser -c "export ORAENV_ASK=NO;export
ORACLE_SID=$oraclesid;. oraenv -s << EOF1
EOF1

```

```

sqlplus -S / as sysdba << EOF2
set heading off
$command
EOF2
"
    `; }
    else {
        $command=~s/\\\\\\\\/\\/g;
        @lines=`export ORAENV_ASK=NO;export ORACLE_SID=$oraclesid;. oraenv
-s << EOF1
EOF1
sqlplus -S / as sysdba << EOF2
set heading off
$command
EOF2
    `; }
return @lines}
print "\n";
@out=dosql('select name from v\\\\\\\\$datafile;');
print @out;
print "shutdown immediate;\n";
print "startup mount;\n";
print "\n";
for $line (@out) {
    if (length($line) > 1) {
        chomp($line);
        ($first, $second, $third, $fourth)=split('_', $line);
        $fourth =~ s/^TS-//;
        $newname=lc("$fourth.dbf");
        $path2file=$line;
        $path2file=~ /(^.*\\.\/)/;
        print "host mv $line $1$newname\n";}}
print "\n";
for $line (@out) {
    if (length($line) > 1) {
        chomp($line);
        ($first, $second, $third, $fourth)=split('_', $line);
        $fourth =~ s/^TS-//;
        $newname=lc("$fourth.dbf");
        $path2file=$line;
        $path2file=~ /(^.*\\.\/)/;
        print "alter database rename file '$line' to
'$1$newname';\n";}}
print "alter database open;\n";
print "\n";

```

## ASM 至檔案系統名稱轉換

```
set serveroutput on;
set wrap off;
declare
    cursor df is select file#, name from v$datafile;
    cursor tf is select file#, name from v$tempfile;
    cursor lf is select member from v$logfile;
    firstline boolean := true;
begin
    dbms_output.put_line(CHR(13));
    dbms_output.put_line('Parameters for log file conversion:');
    dbms_output.put_line(CHR(13));
    dbms_output.put('*.log_file_name_convert = ');
    for lfrec in lf loop
        if (firstline = true) then
            dbms_output.put('''' || lfrec.member || ''', ');
            dbms_output.put(''''/NEW_PATH/' ||
regexp_replace(lfrec.member, '^.*./', '') || ''');
        else
            dbms_output.put(', ''' || lfrec.member || ''', ');
            dbms_output.put(''''/NEW_PATH/' ||
regexp_replace(lfrec.member, '^.*./', '') || ''');
        end if;
        firstline:=false;
    end loop;
    dbms_output.put_line(CHR(13));
    dbms_output.put_line(CHR(13));
    dbms_output.put_line('rman duplication script:');
    dbms_output.put_line(CHR(13));
    dbms_output.put_line('run');
    dbms_output.put_line('{');
    for dfrec in df loop
        dbms_output.put_line('set newname for datafile ' ||
dfrec.file# || ' to ''' || dfrec.name || ''';');
    end loop;
    for tfrec in tf loop
        dbms_output.put_line('set newname for tempfile ' ||
tfrec.file# || ' to ''' || tfrec.name || ''';');
    end loop;
    dbms_output.put_line('duplicate target database for standby backup
location INSERT_PATH_HERE;');
    dbms_output.put_line('}');
end;
/
```

## 在資料庫上重新播放記錄

此指令碼接受 Oracle SID 的單一引數、用於處於掛載模式的資料庫、並嘗試重新播放所有目前可用的歸檔記錄。

```
#!/usr/bin/perl
use strict;
my $oraclesid=$ARGV[0];
my $oracleuser='oracle';
84 Migration of Oracle Databases to NetApp Storage Systems © 2021 NetApp,
Inc. All rights reserved
my $uid = $<;
my @out;
if ($uid == 0) {
@out=`su - $oracleuser -c '. oraenv << EOF1
$oraclesid
EOF1
sqlplus / as sysdba << EOF2
recover database until cancel;
auto
EOF2
`
`;}
else {
@out=`. oraenv << EOF1
$oraclesid
EOF1
sqlplus / as sysdba << EOF2
recover database until cancel;
auto
EOF2
`;
}
print @out;
```

## 在待命資料庫上重新播放記錄

此指令碼與前述指令碼相同、但其設計用於待命資料庫。

```

#!/usr/bin/perl
use strict;
my $oraclesid=$ARGV[0];
my $oracleuser='oracle';
my $uid = $<;
my @out;
if ($uid == 0) {
@out=`su - $oracleuser -c '. oraenv << EOF1
$oraclesid
EOF1
sqlplus / as sysdba << EOF2
recover standby database until cancel;
auto
EOF2
';}
else {
@out=`. oraenv << EOF1
$oraclesid
EOF1
sqlplus / as sysdba << EOF2
recover standby database until cancel;
auto
EOF2
';}
}
print @out;

```

## 其他附註

### Oracle 資料庫效能最佳化與基準測試程序

準確測試資料庫儲存效能是極為複雜的主題。需要瞭解下列問題：

- IOPS 與處理量
- 前景與背景 I/O 作業之間的差異
- 延遲對資料庫的影響
- 許多作業系統和網路設定也會影響儲存效能

此外、還有非儲存資料庫工作要考量。最佳化儲存效能並不會帶來實用效益、因為儲存效能不再是效能的限制因素。

大多數資料庫客戶現在都選擇 All Flash Array、這會造成一些額外考量。例如、請考慮在雙節點 AFF A900 系統上進行效能測試：

- 有了 80/20 讀取 / 寫入比率、兩個 A900 節點可在延遲甚至超過 150  $\mu$ s 標記之前、提供超過 1M 的隨機資料庫 IOPS。這遠遠超出了大多數資料庫目前的效能需求、很難預測預期的改善。儲存設備將會大幅清除、成為瓶頸。
- 網路頻寬是效能限制的常見來源。例如、旋轉式磁碟解決方案通常是資料庫效能的瓶頸、因為 I/O 延遲非常高。當 All Flash 陣列移除延遲限制時、障礙會頻繁移轉至網路。在虛擬化環境和刀鋒系統中、這一點尤其顯著、因為這些環境和刀鋒系統的真正網路連線能力難以視覺化。如果由於頻寬限制而無法充分利用儲存系統本身、這可能會使效能測試變得複雜。
- 由於 All Flash 陣列的延遲大幅改善、因此一般無法將 All Flash 陣列與含有旋轉磁碟的陣列進行效能比較。測試結果通常沒有意義。
- 將尖峰 IOPS 效能與 All Flash 陣列進行比較通常並不實用、因為資料庫不受儲存 I/O 限制例如、假設一個陣列可維持 500K 的隨機 IOPS、而另一個陣列則可維持 300K。如果資料庫花費 99% 的時間處理 CPU、這種差異在現實世界中是不相關的。工作負載永遠不會使用儲存陣列的完整功能。相反地、尖峰 IOPS 功能在整合平台中可能非常重要、而在整合平台中、儲存陣列預期會載入至尖峰容量。
- 請務必在任何儲存測試中考慮延遲和 IOPS。市面上許多儲存陣列都宣稱 IOPS 極高、但延遲卻使這些 IOPS 在這類層級上無法使用。使用 All Flash Array 的典型目標是 1 毫秒標記。更好的測試方法不是測量最大可能的 IOPS、而是判斷儲存陣列在平均延遲大於 1 毫秒之前可以維持多少 IOPS。

## Oracle 自動工作負載儲存庫與基準測試

Oracle 效能比較的黃金標準是 Oracle 自動工作負載儲存庫（AWR）報告。

有多種類型的 AWR 報告。從儲存點來看、執行所產生的報告 `awrrpt.sql Command` 是最全面且最有價值的命令、因為它針對特定資料庫執行個體、並包含一些詳細的分佈圖、可根據延遲來中斷儲存 I/O 事件。

比較兩個效能陣列的理想方法是在每個陣列上執行相同的工作負載、並產生精確鎖定工作負載的 AWR 報告。在執行時間極長的工作負載中、可以使用包含開始和停止時間的單一 AWR 報告、但最好將 AWR 資料分成多份報告。例如、如果批次工作從午夜執行至上午 6 點、請建立一系列從午夜-1 點、上午 1 點-2 點開始的一小時 AWR 報告。

在其他情況下、應最佳化非常簡短的查詢。最佳選項是以查詢開始時建立的 AWR 快照為基礎的 AWR 報告、以及在查詢結束時建立的第二個 AWR 快照。否則資料庫伺服器應保持安靜、以將會使分析中查詢活動模糊的背景活動降至最低。



如果無法取得 AWR 報告、Oracle 狀態報告是一個很好的替代方案。它們包含與 AWR 報告相同的大部分 I/O 統計資料。

## Oracle AWR 與疑難排解

AWR 報告也是分析效能問題的最重要工具。

與基準測試一樣、效能疑難排解也需要您精確測量特定工作負載。如果可能、請在向 NetApp 支援中心回報效能問題、或與 NetApp 或合作夥伴客戶團隊合作、討論新解決方案時提供 AWR 資料。

提供 AWR 資料時、請考量下列需求：

- 執行 `awrrpt.sql` 產生報告的命令。輸出可以是文字或 HTML。
- 如果使用 Oracle Real Application Clusters（RAC）、請為叢集中的每個執行個體產生 AWR 報告。
- 鎖定問題存在的特定時間。AWR 報告的最長可接受使用時間通常為一小時。如果問題持續數小時、或涉及多小時作業、例如批次工作、請提供多個一小時的 AWR 報告、涵蓋整個分析期間。

- 如有可能、請將 AWR 快照時間間隔調整為 15 分鐘。此設定可執行更詳細的分析。這也需要額外執行 `awrrpt.sql` 提供每 15 分鐘間隔的報告。
- 如果問題是非常短的執行查詢、請根據作業開始時建立的 AWR 快照、以及作業結束時建立的第二個 AWR 快照、提供 AWR 報告。否則、資料庫伺服器應保持安靜、以將會使分析中作業的活動受到影響的背景活動減至最低。
- 如果在特定時間回報效能問題、但未在其他時間回報、請提供額外的 AWR 資料、以展現良好的效能來進行比較。

## calibr\_IO

◦ `calibrate_io` 絕對不可使用命令來測試、比較或基準測試儲存系統。如 Oracle 文件所述、本程序會校正儲存設備的 I/O 功能。

校準與基準測試不同。此命令的目的是發佈 I/O、藉由最佳化發行給主機的 I/O 層級、協助校正資料庫作業並改善其效率。因為執行的 I/O 類型 `calibrate_io` 作業並不代表實際的資料庫使用者 I/O、結果無法預測、而且經常無法重現。

## SLOB2

SLOB2 是愚蠢的小 Oracle 基準測試工具、已成為評估資料庫效能的首選工具。這是由 Kevin Closson 開發的、可在取得 "<https://kevinclosson.net/slob/>"。安裝和設定需要幾分鐘的時間、它使用實際的 Oracle 資料庫來在使用者可定義的資料表空間上產生 I/O 模式。這是少數幾種可用的測試選項之一、可將全快閃陣列與 I/O 飽和它也有助於產生更低層級的 I/O、以模擬低 IOPS 但對延遲敏感的儲存工作負載。

## 交換台工作台

交換基準台可用於測試資料庫效能、但使用交換基準台的方式會對儲存造成壓力、這是非常困難的。NetApp 尚未從 SwingWorkbench 中看到任何測試結果、這些測試產生足夠的 I/O、使其在任何 AFF 陣列上都成為重大負載。在有限的情況下、訂單輸入測試 (OET) 可用於從延遲點評估儲存設備。這在資料庫對於特定查詢具有已知延遲相依性的情況下很有用。必須注意確保主機和網路已正確設定、以實現 All Flash 陣列的延遲潛力。

## HammerDB

HammerDB 是一種資料庫測試工具、可模擬 TPC-C 和 TPC-H 基準測試等。建構足夠大的資料集以正確執行測試可能需要很長時間、但它可以是評估 OLTP 和資料倉儲應用程式效能的有效工具。

## Orion

Oracle Orion 工具通常與 Oracle 9 搭配使用、但並未加以維護、以確保與各種主機作業系統的變更相容。由於與作業系統和儲存組態不相容、因此很少與 Oracle 10 或 Oracle 11 搭配使用。

Oracle 重新編寫了該工具，默認情況下，該工具與 Oracle 12c 一起安裝。雖然本產品已經過改良、並使用許多與實際 Oracle 資料庫相同的呼叫、但它並未使用 Oracle 所使用的相同程式碼路徑或 I/O 行為。例如、大部分的 Oracle I/O 都是同步執行、這表示資料庫會暫停、直到 I/O 作業在前景完成為止。只是以隨機 I/O 淹沒儲存系統、並不是真正的 Oracle I/O 複製、也無法提供直接的儲存陣列比較方法、也無法測量組態變更的影響。

也就是說、Orion 有一些使用案例、例如一般測量特定主機網路儲存組態的最大可能效能、或是測量儲存系統的健全狀況。仔細測試後、只要參數包括 IOPS、處理量和延遲的考量、並嘗試忠實複製真實的工作負載、就能設計出可用的 Orion 測試來比較儲存陣列或評估組態變更的影響。



## 過時的 NFSv3 鎖定和 Oracle 資料庫

如果 Oracle 資料庫伺服器當機、則重新啟動時可能會發生過時的 NFS 鎖定問題。請仔細注意伺服器上的名稱解析設定、以避免此問題。

產生此問題的原因是建立鎖定和清除鎖定會使用兩種稍微不同的名稱解析方法。涉及兩個過程：Network Lock Manager (NLM) 和 NFS 用戶端。NLM 使用 `uname -n` 以決定主機名稱、而 `rpc.statd` 程序用途 `gethostbyname()`。這些主機名稱必須相符、作業系統才能正確清除過時的鎖定。例如、主機可能正在尋找擁有的鎖定 `dbserver5`、但鎖定已由主機登錄為 `dbserver5.mydomain.org`。如果 `gethostbyname()` 不會傳回與相同的值 `uname -a`，則鎖定釋放程序未成功。

下列範例指令碼會驗證名稱解析是否完全一致：

```
#!/usr/bin/perl
$uname=`uname -n`;
chomp($uname);
($name, $aliases, $addrtype, $length, @addrs) = gethostbyname $uname;
print "uname -n yields: $uname\n";
print "gethostbyname yields: $name\n";
```

如果 `gethostbyname` 不符 `uname`、可能是過時的鎖定。例如、此結果顯示潛在問題：

```
uname -n yields: dbserver5
gethostbyname yields: dbserver5.mydomain.org
```

解決方案通常是透過變更主機出現在中的順序來找到 `/etc/hosts`。例如、假設 `hosts` 檔案包含下列項目：

```
10.156.110.201 dbserver5.mydomain.org dbserver5 loghost
```

若要解決此問題、請變更完整網域名稱和簡短主機名稱出現的順序：

```
10.156.110.201 dbserver5 dbserver5.mydomain.org loghost
```

`gethostbyname()` 現在傳回短 `dbserver5` 主機名稱、符合的輸出 `uname`。因此、鎖定會在伺服器當機後自動清除。

## Oracle 資料庫的 WAFL 對齊驗證

正確的 WAFL 對齊對於良好的效能至關重要。雖然 ONTAP 以 4KB 單位管理區塊、但這並不表示 ONTAP 以 4KB 單位執行所有作業。事實上、ONTAP 支援不同大小的區塊作業、但基礎會計是由 WAFL 以 4KB 為單位進行管理。

「對齊」一詞是指 Oracle I/O 與這些 4KB 單元的相對應方式。最佳效能要求 Oracle 8KB 區塊位於磁碟機上的

兩個 4KB WAFL 實體區塊上。如果區塊偏移 2KB、則此區塊會位於一半的 4KB 區塊、一個獨立的完整 4KB 區塊、然後是第三個 4KB 區塊的一半。這種安排會導致效能降低。

對齊並不涉及 NAS 檔案系統。Oracle 資料檔案會根據 Oracle 區塊的大小、與檔案的開頭對齊。因此、8KB、16KB 和 32KB 的區塊大小一律會對齊。所有區塊作業都會從檔案開頭偏移、單位為 4 KB。

相反地、LUN 在啟動時通常會包含某種驅動程式標頭或檔案系統中繼資料、以建立偏移。對齊在現代作業系統中很少是個問題、因為這些作業系統是專為可能使用原生 4KB 磁碟區的實體磁碟機所設計、因此也需要將 I/O 與 4KB 邊界對齊才能獲得最佳效能。

不過、有一些例外情況。資料庫可能已從未針對 4KB I/O 最佳化的舊版作業系統移轉、或是在建立分割區時發生使用者錯誤、可能導致偏移量、而大小單位不是 4KB。

下列範例僅適用於 Linux、但程序可適用於任何作業系統。

一致

以下範例顯示單一磁碟分割的單一 LUN 對齊檢查。

首先、建立使用磁碟機上所有可用分割區的分割區。

```
[root@host0 iscsi]# fdisk /dev/sdb
Device contains neither a valid DOS partition table, nor Sun, SGI or OSF
disklabel
Building a new DOS disklabel with disk identifier 0xb97f94c1.
Changes will remain in memory only, until you decide to write them.
After that, of course, the previous content won't be recoverable.
The device presents a logical sector size that is smaller than
the physical sector size. Aligning to a physical sector (or optimal
I/O) size boundary is recommended, or performance may be impacted.
Command (m for help): n
Command action
   e   extended
   p   primary partition (1-4)
p
Partition number (1-4): 1
First cylinder (1-10240, default 1):
Using default value 1
Last cylinder, +cylinders or +size{K,M,G} (1-10240, default 10240):
Using default value 10240
Command (m for help): w
The partition table has been altered!
Calling ioctl() to re-read partition table.
Syncing disks.
[root@host0 iscsi]#
```

您可以使用下列命令以數學方式檢查對齊方式：

```
[root@host0 iscsi]# fdisk -u -l /dev/sdb
Disk /dev/sdb: 10.7 GB, 10737418240 bytes
64 heads, 32 sectors/track, 10240 cylinders, total 20971520 sectors
Units = sectors of 1 * 512 = 512 bytes
Sector size (logical/physical): 512 bytes / 4096 bytes
I/O size (minimum/optimal): 4096 bytes / 65536 bytes
Disk identifier: 0xb97f94c1

   Device Boot      Start         End      Blocks   Id  System
/dev/sdb1            32      20971519    10485744    83   Linux
```

輸出顯示單位為 512 位元組、且分割區的開頭為 32 個單位。總共  $32 \times 512 = 16,384$  位元組、這是 4KB WAFL 區塊的整數倍數。此分割區已正確對齊。

若要驗證正確的對齊方式、請完成下列步驟：

1. 識別 LUN 的通用唯一識別碼 (UUID)。

```
FAS8040SAP::> lun show -v /vol/jfs_luns/lun0
      Vserver Name: jfs
      LUN UUID: ed95d953-1560-4f74-9006-85b352f58fcd
      Mapped: mapped`
```

2. 進入 ONTAP 控制器上的節點 Shell。

```
FAS8040SAP::> node run -node FAS8040SAP-02
Type 'exit' or 'Ctrl-D' to return to the CLI
FAS8040SAP-02> set advanced
set not found. Type '?' for a list of commands
FAS8040SAP-02> priv set advanced
Warning: These advanced commands are potentially dangerous; use
them only when directed to do so by NetApp
personnel.
```

3. 在第一步中識別的目標 UUID 上開始收集統計資料。

```
FAS8040SAP-02*> stats start lun:ed95d953-1560-4f74-9006-85b352f58fcd
Stats identifier name is 'Ind0xffffffff08b9536188'
FAS8040SAP-02*>
```

4. 執行一些 I/O 請務必使用 `iflag` 用於確保 I/O 同步且無緩衝的引數。



請務必小心使用此命令。反轉 `if` 和 `of` 引數會破壞資料。

```
[root@host0 iscsi]# dd if=/dev/sdb1 of=/dev/null iflag=dsync count=1000
bs=4096
1000+0 records in
1000+0 records out
4096000 bytes (4.1 MB) copied, 0.0186706 s, 219 MB/s
```

5. 停止統計資料並檢視對齊分佈圖。所有 I/O 都應位於 .0 貯體、表示 I/O 與 4KB 區塊邊界對齊。

```
FAS8040SAP-02*> stats stop
StatisticsID: Ind0xffffffff08b9536188
lun:ed95d953-1560-4f74-9006-85b352f58fcd:instance_uuid:ed95d953-1560-
4f74-9006-85b352f58fcd
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_align_histo.0:186%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_align_histo.1:0%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_align_histo.2:0%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_align_histo.3:0%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_align_histo.4:0%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_align_histo.5:0%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_align_histo.6:0%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_align_histo.7:0%
```

## 未對齊

以下範例顯示 I/O 未對齊：

1. 建立不符合 4KB 邊界的分割區。這不是現代作業系統的預設行為。

```
[root@host0 iscsi]# fdisk -u /dev/sdb
Command (m for help): n
Command action
  e   extended
  p   primary partition (1-4)
p
Partition number (1-4): 1
First sector (32-20971519, default 32): 33
Last sector, +sectors or +size{K,M,G} (33-20971519, default 20971519):
Using default value 20971519
Command (m for help): w
The partition table has been altered!
Calling ioctl() to re-read partition table.
Syncing disks.
```

2. 已建立磁碟分割、並使用 33 磁區偏移值、而非預設的 32。重複中所述的程序 "一致"。直方圖顯示如下：

```

FAS8040SAP-02*> stats stop
StatisticsID: Ind0xffffffff0468242e78
lun:ed95d953-1560-4f74-9006-85b352f58fcd:instance_uuid:ed95d953-1560-4f74-9006-85b352f58fcd
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_align_histo.0:0%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_align_histo.1:136%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_align_histo.2:4%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_align_histo.3:0%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_align_histo.4:0%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_align_histo.5:0%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_align_histo.6:0%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_align_histo.7:0%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:read_partial_blocks:31%

```

錯誤的對齊是顯而易見的。I/O 大多落在 \* 之中 \* .1 符合預期偏移的貯體。建立分割區時、它會比最佳化的預設值更進一步移入 512 個位元組、這表示長條圖偏移 512 個位元組。

此外 read\_partial\_blocks 統計資料為非零、這表示執行的 I/O 並未填滿整個 4KB 區塊。

## 重作記錄

此處說明的程序適用於資料檔案。Oracle 重做記錄和歸檔記錄檔有不同的 I/O 模式。例如、重做記錄是單一檔案的循環覆寫。如果使用預設的 512 位元組區塊大小、寫入統計資料看起來會像這樣：

```

FAS8040SAP-02*> stats stop
StatisticsID: Ind0xffffffff0468242e78
lun:ed95d953-1560-4f74-9006-85b352f58fcd:instance_uuid:ed95d953-1560-4f74-9006-85b352f58fcd
lun:ed95d953-1560-4f74-9006-85b352f58fcd:write_align_histo.0:12%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:write_align_histo.1:8%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:write_align_histo.2:4%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:write_align_histo.3:10%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:write_align_histo.4:13%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:write_align_histo.5:6%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:write_align_histo.6:8%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:write_align_histo.7:10%
lun:ed95d953-1560-4f74-9006-85b352f58fcd:write_partial_blocks:85%

```

I/O 會分散到所有分佈式分佈區、但這並不是效能考量。不過、重做記錄率極高可能會因為使用 4KB 區塊大小而受惠。在這種情況下、最好確定重做記錄 LUN 已正確對齊。不過、這對於資料檔案對齊的良好效能並不重要。

# PostgreSQL

## ONTAP 上的 PostgreSQL 資料庫

PostgreSQL 隨附的變種包括 PostgreSQL、PostgreSQL Plus 和 EDBS PostgreSQL 進階伺服器（EPAS）。PostgreSQL 通常部署為多層應用程式的後端資料庫。它受一般中介軟體套件的支援（例如 PHP、Java、Python、Tcl/Tk、ODBC、和 JDBC），過去一直是開放原始碼資料庫管理系統的熱門選擇。ONTAP 是執行 PostgreSQL 資料庫的絕佳選擇、可確保其可靠性、高效能及高效率的資料管理功能。



ONTAP 和 PostgreSQL 資料庫上的這份文件取代先前發佈的 \_TR-4770：ONTAP 最佳實務做法的 PostgreSQL 資料庫。 \_

隨著資料呈指數成長、企業的資料管理變得更加複雜。這種複雜性會增加授權、營運、支援和維護成本。若要降低整體 TCO、請考慮使用可靠、高效能的後端儲存設備、從商業資料庫切換為開放原始碼資料庫。

ONTAP 是理想的平台、因為 ONTAP 是專為資料庫所設計。為了滿足資料庫工作負載的需求、我們特別建立了許多功能、例如隨機 IO 延遲最佳化、以提供進階服務品質（QoS）到基本 FlexClone 功能。

其他功能（例如不中斷升級）（包括儲存設備更換）、可確保關鍵資料庫仍可使用。您也可以透過 MetroCluster 為大型環境進行即時災難恢復、或是使用 SnapMirror 主動式同步來選擇資料庫。

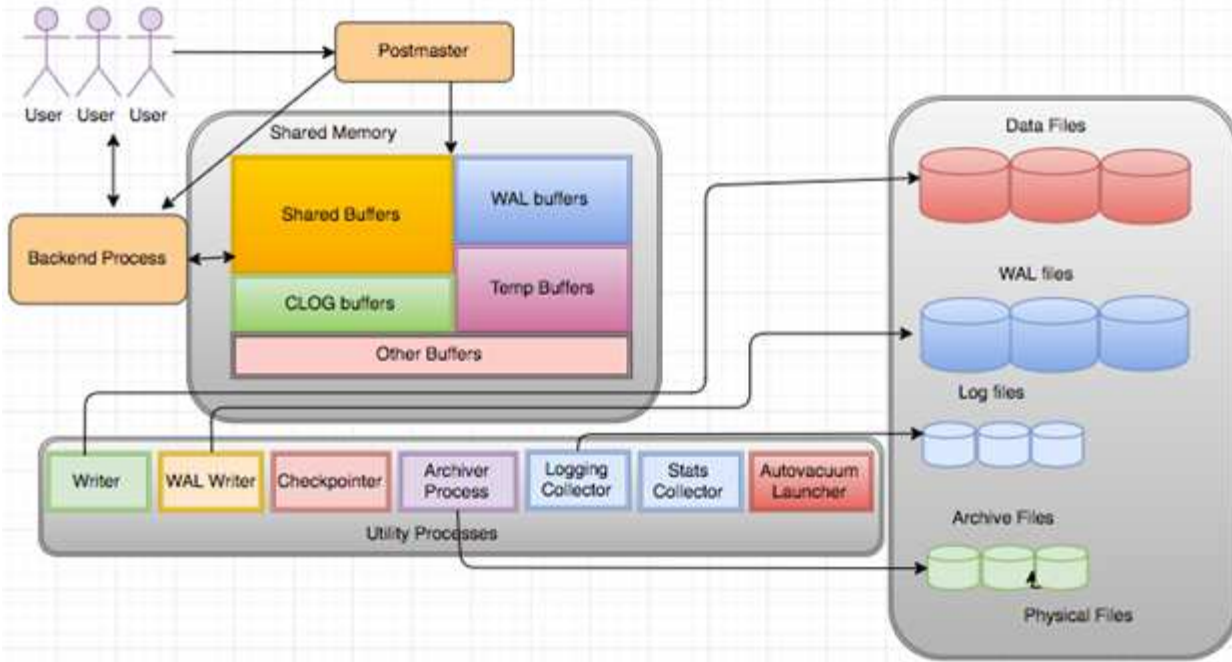
最重要的是、ONTAP 提供無與倫比的效能、並能根據您的獨特需求調整解決方案的規模。我們的高階系統可提供超過 1M IOPS、延遲時間以微秒為測量單位、但如果您只需要 10 萬次 IOPS、您可以使用仍執行相同儲存作業系統的較小控制器來調整儲存解決方案的大小。

## 資料庫組態

### PostgreSQL 架構

PostgreSQL 是以用戶端和伺服器架構為基礎的 RDBMS。PostgreSQL 執行個體稱為資料庫叢集、是資料庫的集合、而非伺服器集合。

## PostgreSQL Basic Architecture



PostgreSQL 資料庫中有三個主要元素：郵政局長、前端（用戶端）和後端用戶端會傳送要求給郵政局長、其中包含 IP 傳輸協定等資訊、以及要連線的資料庫。郵政局長會驗證連線、並將其傳送至後端程序以進行進一步的通訊。後端程序會執行查詢、並將結果直接傳送至前端（用戶端）。

PostgreSQL 執行個體是以多重處理模式為基礎、而非多重執行緒模式。它會為不同的工作產生多個處理程序、而且每個處理程序都有自己的功能。主要程序包括用戶端程序、Wal 寫入程序、背景寫入程序及檢查指標程序：

- 當用戶端（前景）程序傳送讀取或寫入要求至 PostgreSQL 執行個體時、它不會直接讀取或寫入資料至磁碟。它首先緩衝共享緩衝區和預先寫入記錄（Wal）緩衝區中的資料。
- Wal 寫入器程序會操控共用緩衝區和 Wal 緩衝區的內容、以寫入 Wal 記錄檔。Wal 記錄檔通常是 PostgreSQL 的交易記錄檔、並依序寫入。因此、為了改善資料庫的回應時間、PostgreSQL 會先寫入交易記錄檔、並確認用戶端。
- 若要将資料庫置於一致的狀態、背景寫入器程序會定期檢查共用緩衝區是否有髒頁。然後將資料排清至儲存在 NetApp 磁碟區或 LUN 上的資料檔案。
- checkpointer 程序也會定期執行（比背景程序更少）、並防止對緩衝區進行任何修改。它會向 Wal 寫入器程序發出訊號、將檢查點記錄寫入並清除至儲存在 NetApp 磁碟上的 Wal 記錄檔結尾。它也會向背景寫入器程序發出訊號、要求將所有髒頁寫入磁碟並清除。

## PostgreSQL 初始化參數

您可以使用建立新的資料庫叢集 `initdb` 方案。一 `initdb` 指令碼會建立定義叢集的資料檔案、系統表格和範本資料庫（`template0` 和 `template1`）。

範本資料庫代表常用資料庫。其中包含系統表格、標準檢視、功能和資料類型的定義。 `pgdata` 做為的引數 `initdb` 指定資料庫叢集位置的指令碼。

PostgreSQL 中的所有資料庫物件都是由各自的 OID 在內部管理。表格和索引也由個別的 OID 管理。資料庫物件與其各自的 OID 之間的關係會儲存在適當的系統目錄表格中、視物件類型而定。例如、資料庫和堆積表格的

OID 會儲存在中 `pg_database` 和 `"pg_class"`。您可以在 PostgreSQL 用戶端上發佈查詢來判斷 OID。

每個資料庫都有自己的個別資料表和索引檔案、限制為 1GB。每個表格都有兩個相關的檔案、分別以後綴表示 `_fsm` 和 `_vm`。它們稱為可用空間地圖和可見度地圖。這些檔案會儲存可用空間容量的相關資訊、並在表格檔案中的每個頁面上都有可見度。索引只有個別的可用空間地圖、而且沒有可見度地圖。

◦ `pg_xlog/pg_wal` 目錄包含預先寫入記錄。預先寫入記錄可用來改善資料庫的可靠性和效能。每當您更新表格中的列時、PostgreSQL 會先將變更寫入預先寫入記錄、然後將修改寫入實際的資料頁面到磁碟。

`pg_xlog` 目錄通常包含數個檔案、但 `initdb` 只會建立第一個檔案。視需要新增額外檔案。每個 `xlog` 檔案長度為 16MB。

## ONTAP 的 PostgreSQL 資料庫組態

有幾種 PostgreSQL 調校組態可以改善效能。

最常用的參數如下：

- `max_connections = <num>`：一次擁有的最大資料庫連線數。使用此參數可限制磁碟交換並終止效能。根據應用程式需求、您也可以針對連線集區設定調整此參數。
- `shared_buffers = <num>`：提高資料庫伺服器效能的最簡單方法。對於大多數現代硬體而言、預設值為低。在部署期間、系統上的可用 RAM 約為 25%。此參數設定會因其與特定資料庫執行個體的運作方式而異；您可能必須根據試用和錯誤來增加和減少值。不過、將其設為高可能會降低效能。
- `effective_cache_size = <num>`：此值告訴 PostgreSQL 的最佳化程式 PostgreSQL 有多少記憶體可用於快取資料、並有助於判斷是否使用索引。較大的值會增加使用索引的可能性。此參數應設為分配給的記憶體容量 `shared_buffers` 加上可用的作業系統快取容量。此值通常超過系統總記憶體的 50%。
- `work_mem = <num>`：此參數控制用於排序作業和雜湊表的記憶體容量。如果您在應用程式中進行大量排序、可能需要增加記憶體容量、但請謹慎。它不是系統範圍的參數、而是每次操作的參數。如果複雜查詢中有多個排序作業、則會使用多個 `work_mem` 記憶體單元、而多個後端也可能同時執行此作業。如果值太大、此查詢通常會引導您的資料庫伺服器進行切換。此選項先前在舊版 PostgreSQL 中稱為 `sort_mem`。
- `fsync = <boolean> (on or off)`：此參數確定在提交事務之前是否應使用 `fsync ()` 將所有 Wal 頁面同步到磁碟。關閉它有時會改善寫入效能、並將其開啟、以提高系統當機時避免毀損的風險。
- `checkpoint_timeout`：檢查點處理程序會將已提交的資料清除至磁碟。這涉及磁碟上的大量讀寫作業。此值以秒為單位設定、較低的值可減少損毀恢復時間、而增加的值則可減少檢查點呼叫、進而降低系統資源的負載。根據應用程式的關鍵程度、使用量、資料庫可用度、設定 `checkpoint` 逾時的值。
- `commit_delay = <num>` 和 `commit_siblings = <num>`：這些選項可同時用於撰寫多筆同時提交的交易、以協助改善效能。如果交易提交時有多個 `command_siblings` 物件處於作用中狀態、伺服器會等待 `command_delay` 微秒、嘗試一次提交多個交易。
- `max_worker_processes / max_parallel_workers`：配置流程的最佳工作人員數量。`max_parallel_workers` 對應可用的 CPU 數量。視應用程式設計而定、查詢可能需要較少的工作人員來執行平行作業。最好保持兩個參數的值相同、但在測試後調整值。
- `random_page_cost = <num>`：此值控制 PostgreSQL 檢視非連續磁碟讀取的方式。較高的值表示 PostgreSQL 較可能使用連續掃描、而非索引掃描、表示伺服器有快速磁碟。請在評估其他選項（例如計畫型最佳化、吸塵、索引以變更查詢或架構）之後、修改此設定。
- `effective_io_concurrency = <num>`：此參數設置 PostgreSQL 嘗試同時執行的並行磁碟 I/O 操作數。提高此值會增加任何個別 PostgreSQL 工作階段嘗試平行啟動的 I/O 作業數。允許範圍為 1 到 1,000、或為零、以停用非同步 I/O 要求的發出。目前、此設定只會影響點陣圖堆疊掃描。固態硬碟（SSD）和其他記憶體型儲存設備（NVMe）通常可以處理許多並行要求、因此最佳價值可能在數百種環境中。



如需 PostgreSQL 組態參數的完整清單、請參閱 PostgreSQL 文件。

## 吐司

Toast 代表「超大型屬性儲存技術」。PostgreSQL 使用固定的頁面大小（通常為 8KB）、不允許 Tuple 跨越多個頁面。因此、無法直接儲存大欄位值。當您嘗試儲存超過此大小的資料列時、Toast 會將大型資料欄的資料分成較小的「片段」、並將其儲存在 Toast 資料表中。

只有在將結果集傳送至用戶端時、才會拔出（如果完全選取）已烘烤的屬性大值。表格本身比沒有任何離線儲存設備（Toast）時更小、可容納更多資料列到共用緩衝區快取中。

## 真空

在正常的 PostgreSQL 作業中、因為更新而刪除或過時的 Tuple 不會從其表格中實際移除、直到執行真空為止。因此、您必須定期執行吸氣、尤其是在經常更新的表格上。接著必須回收它所佔用的空間、以便新列重複使用、以避免磁碟空間中斷。不過、它不會將空間傳回作業系統。

頁面內的可用空間不會分散。真空會重新寫入整個區塊、有效地將剩餘的資料列封裝起來、並在頁面中留下一個連續的可用空間區塊。

相反地、使用「完全真空」技術、可以撰寫完全新版的表格檔案、而不會有任何空間。此動作可將表格大小減至最小、但可能需要很長時間。在作業完成之前、它也需要額外的磁碟空間來容納表格的新複本。例行真空的目標是避免完全真空。此程序不僅能將資料表保持在最小大小、還能維持磁碟空間的穩定狀態使用量。

## PostgreSQL 表格空間

初始化資料庫叢集時會自動建立兩個資料表空間。

◦ `pg_global` 表空間用於共享系統目錄。◦ `pg_default` 表空間是 `template1` 和 `template0` 資料庫的預設資料表空間。如果初始化叢集的磁碟分割或磁碟區空間不足且無法擴充、則可在不同的磁碟分割上建立資料表空間、並在重新設定系統之前使用。

大量使用的索引可以放在快速、高可用度的磁碟上、例如固態裝置。此外、儲存歸檔資料的表格、無論是極少使用或非效能關鍵、都可以儲存在成本較低、速度較慢的磁碟系統上、例如 SAS 或 SATA 磁碟機。

資料表空間是資料庫叢集的一部分、無法視為資料檔案的獨立集合。它們取決於主資料目錄中包含的中繼資料、因此無法附加至不同的資料庫叢集或個別備份。同樣地、如果您遺失資料表空間（透過檔案刪除、磁碟故障等方式）、資料庫叢集可能會變得無法讀取或無法啟動。將資料表空間放在像 RAM 磁碟一樣的暫存檔案系統上、會危及整個叢集的可靠性。

建立表格後、如果要求的使用者擁有足夠的權限、則可以從任何資料庫使用表格區。PostgreSQL 使用符號連結來簡化資料表空間的實作。PostgreSQL 會在中新增一列 `pg_tablespace` 表（叢集範圍表格）、並將新的物件識別碼（OID）指派給該列。最後、伺服器會使用 OID 在叢集與指定目錄之間建立符號連結。目錄 `$PGDATA/pg_tblspc` 包含指向叢集中定義的每個非內建表格空間的符號連結。

## 儲存組態

### 具有 NFS 檔案系統的 PostgreSQL 資料庫

PostgreSQL 資料庫可以裝載在 NFSv3 或 NFSv4 檔案系統上。最佳選項取決於資料庫以外的因素。

例如、在某些叢集式環境中、NFSv4 鎖定行為可能較為理想。(請參閱 ["請按這裡"](#) 如需其他詳細資料)

否則、資料庫功能應接近相同、包括效能。唯一的要求是使用 `hard` 掛載選項。這是為了確保軟性逾時不會產生無法恢復的 IO 錯誤。

如果選擇 NFSv4 作為傳輸協定、NetApp 建議使用 NFSv4.1。NFSv4.1 中的 NFSv4 傳輸協定有一些功能性增強功能、可改善 NFSv4.0 的恢復能力。

針對一般資料庫工作負載、請使用下列掛載選項：

```
rw,hard,nointr,bg,vers=[3|4],proto=tcp,rsize=65536,wsiz=65536
```

如果需要大量連續 IO、則可依照下節所述增加 NFS 傳輸大小。

## NFS 傳輸大小

根據預設、ONTAP 將 NFS I/O 大小限制為 64K。

大多數應用程式和資料庫的隨機 I/O 使用的區塊大小要小得多、遠低於 64K 的最大值。大型區塊 I/O 通常是平行處理的、因此 64K 的最大值也不是取得最大頻寬的限制。

有些工作負載的上限為 64K、因此會造成限制。特別是、如果資料庫執行的 I/O 數量較少、但容量較大、則備份或還原作業或資料庫完整表格掃描等單執行緒作業、會更快、更有效率地執行。ONTAP 的最佳 I/O 處理大小為 256k。

指定 ONTAP SVM 的最大傳輸大小可變更如下：

```
Cluster01::> set advanced
Warning: These advanced commands are potentially dangerous; use them only
when directed to do so by NetApp personnel.
Do you want to continue? {y|n}: y
Cluster01::*> nfs server modify -vserver vserver1 -tcp-max-xfer-size
262144
Cluster01::*>
```

### 注意

切勿將 ONTAP 上允許的傳輸大小上限降至低於目前掛載之 NFS 檔案系統的 `rsiz/wsiz` 值。這可能會在某些作業系統中造成當機或甚至資料毀損。例如、如果 NFS 用戶端目前設定為 `rsiz/wsiz 65536`、則 ONTAP 最大傳輸大小可在 65536 到 1048576 之間調整、因為用戶端本身受到限制、因此沒有任何影響。將傳輸大小上限降至 65536 以下可能會損害可用度或資料。

在 ONTAP 層級增加傳輸大小後、將會使用下列掛載選項：

```
rw,hard,nointr,bg,vers=[3|4],proto=tcp,rsiz=262144,wsiz=262144
```

## NFSv3 TCP 插槽表

如果 NFSv3 與 Linux 搭配使用、則正確設定 TCP 插槽表至關重要。

TCP 插槽表是與主機匯流排介面卡（HBA）佇列深度相當的 NFSv3。這些表格可控制任何時間都可以處理的 NFS 作業數量。預設值通常為 16、這對於最佳效能而言太低。相反的問題發生在較新的 Linux 核心上、這會自動將 TCP 插槽表格限制增加到要求使 NFS 伺服器飽和的層級。

為了達到最佳效能並避免效能問題、請調整控制 TCP 插槽表的核心參數。

執行 `sysctl -a | grep tcp.*.slot_table` 並觀察下列參數：

```
# sysctl -a | grep tcp.*.slot_table
sunrpc.tcp_max_slot_table_entries = 128
sunrpc.tcp_slot_table_entries = 128
```

所有 Linux 系統都應該包括在內 `sunrpc.tcp_slot_table_entries`、但只有部分包含在內 `sunrpc.tcp_max_slot_table_entries`。兩者都應設為 128。

### 注意

若未設定這些參數、可能會對效能造成重大影響。在某些情況下、效能會受到限制、因為 Linux 作業系統沒有發出足夠的 I/O 在其他情況下、隨著 Linux 作業系統嘗試發出的 I/O 數量超過可服務的數量、I/O 延遲也會增加。

## 具有 SAN 檔案系統的 PostgreSQL

具有 SAN 的 PostgreSQL 資料庫通常裝載在 xfs 檔案系統上、但如果作業系統廠商支援、也可以使用其他資料庫

雖然單一 LUN 通常可支援高達十萬 IOPS、但 IO 密集的資料庫通常需要使用含分段的 LVM。

### LVM 分拆

在快閃磁碟機時代之前、使用區塊延展來協助克服旋轉磁碟機的效能限制。例如、如果作業系統需要執行 1MB 讀取作業、則從單一磁碟機讀取 1MB 的資料時、需要大量的磁碟機磁頭搜尋和讀取、因為 1MB 會緩慢傳輸。如果將 1MB 的資料分散在 8 個 LUN 上、則作業系統可能會同時執行 8 個 128K 讀取作業、並縮短完成 1MB 傳輸所需的時間。

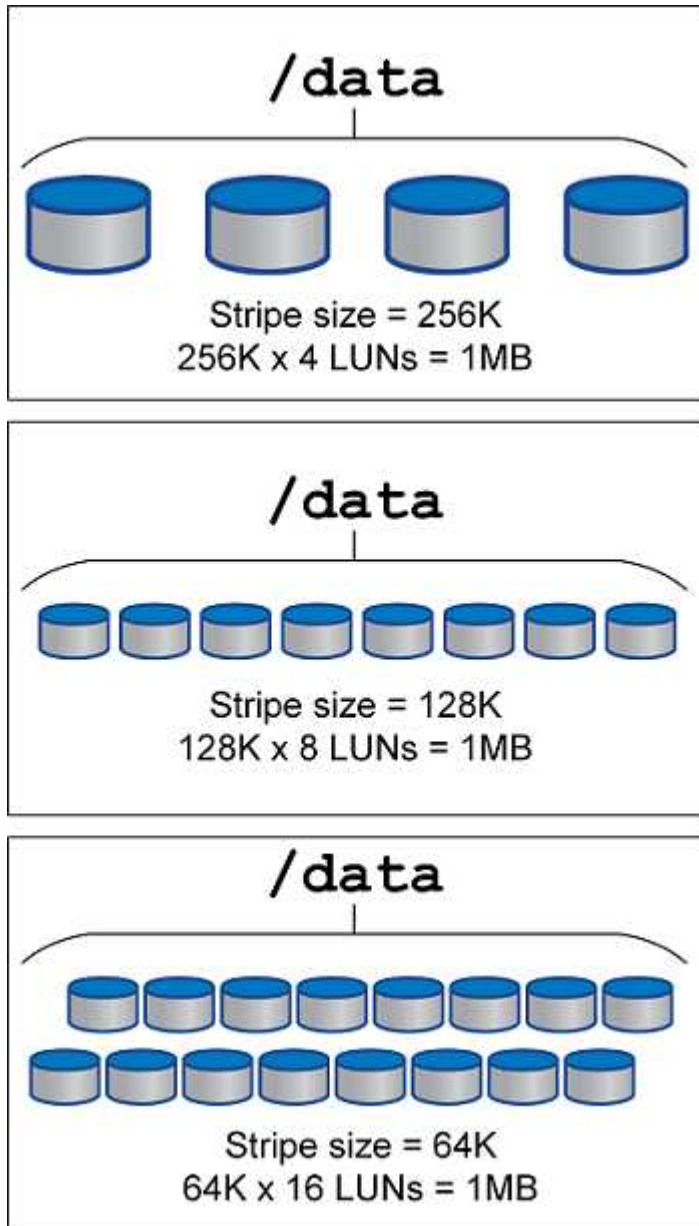
由於必須事先知道 I/O 模式、因此使用旋轉磁碟機進行分拆會更困難。如果串列區塊延展未針對真正的 I/O 模式正確調整、則等量區塊配置可能會損害效能。使用 Oracle 資料庫、特別是搭配 All Flash 組態、分拆作業更容易設定、並經證實可大幅提升效能。

依預設、邏輯磁碟區管理程式（例如 Oracle ASM 等量磁碟區）不屬於原生 OS LVM。其中有些 LUN 會將多個 LUN 連結在一起、成為串連的裝置、導致資料檔案存在於一台 LUN 裝置上、而只存在於一台 LUN 裝置上。這會造成熱點。其他 LVM 實作預設為分散式擴充。這與分拆類似、但卻是比較粗糙的。磁碟區群組中的 LUN 會切成大型片段、稱為區段、通常以百萬位元組為單位測量、然後邏輯磁碟區會分佈在這些區段中。結果是對檔案進行隨機 I/O、應該能在 LUN 之間妥善分配、但連續 I/O 作業的效率卻不如以前那麼高。

效能密集的應用程式 I/O 幾乎總是（a）以基本區塊大小為單位、或（b）1 MB。

等量分配組態的主要目標是確保單一檔案 I/O 可作為單一單元執行、而多區塊 I/O 的大小應為 1MB、可在等量磁碟區中的所有 LUN 之間平均平行處理。這表示等量磁碟區大小不得小於資料庫區塊大小、且等量磁碟區大小乘以 LUN 數量應為 1MB。

下圖顯示等量磁碟區大小和寬度調校的三個可能選項。選擇 LUN 數量以滿足上述效能需求、但在所有情況下、單一等量磁碟區內的總資料為 1MB。



## 資料保護

### PostgreSQL 資料保護

儲存設計的主要層面之一、是為 PostgreSQL Volume 提供保護。客戶可以使用傾印方法或使用檔案系統備份來保護 PostgreSQL 資料庫。本節說明備份個別資料庫或整個叢集的不同方法。

備份 PostgreSQL 資料有三種方法：

- SQL Server 傾印
- 檔案系統層級備份
- 持續歸檔

SQL Server 傾印方法背後的想法是使用 SQL Server 命令產生檔案、當檔案傳回伺服器時、就能像在傾印時一樣重新建立資料庫。PostgreSQL 提供公用程式 `pg_dump` 和 `pg_dump_all` 用於建立個別和叢集層級的備份。這些傾印都是邏輯的、而且沒有足夠的資訊可供 Wal 重新執行使用。

另一種備份策略是使用檔案系統層級的備份、系統管理員可在其中直接複製 PostgreSQL 用來將資料儲存在資料庫中的檔案。此方法是在離線模式下執行：必須關閉資料庫或叢集。另一種方法是使用 `pg_basebackup` 執行 PostgreSQL 資料庫的熱串流備份。

## PostgreSQL 資料庫和儲存快照

PostgreSQL 的 Snapshot 型備份需要為資料檔案、Wal 檔案和歸檔的 Wal 檔案組態快照、才能提供完整或時間點還原。

對於 PostgreSQL 資料庫而言、快照的平均備份時間介於數秒到數分鐘之間。此備份速度比快 60 至 100 倍 `pg_basebackup` 以及其他檔案系統備份方法。

NetApp 儲存設備上的快照既可一致當機、也可一致應用程式。當機一致的快照是在儲存設備上建立、而不會停止資料庫、而在資料庫處於備份模式時則會建立應用程式一致的快照。NetApp 也可確保後續快照是永久遞增備份、以節省儲存成本並提高網路效率。

由於快照速度很快、不會影響系統效能、因此您可以每天排程多個快照、而不必像其他串流備份技術一樣建立單一的每日備份。當需要還原與還原作業時、系統停機時間會減少兩項重要功能：

- NetApp SnapRestore 資料還原技術代表還原作業只需數秒即可完成。
- 積極的恢復點目標（RPO）意味著必須套用較少的資料庫記錄、同時也加速轉送恢復。

若要備份 PostgreSQL、您必須使用（一致性群組）Wal 和歸檔記錄、確保資料磁碟區同時受到保護。使用 Snapshot 技術複製 Wal 檔案時、請務必執行 `pg_stop` 清除所有必須歸檔的 Wal 項目。如果您在還原期間清除 Wal 項目、則只需要停止資料庫、卸載或刪除現有的資料目錄、並在儲存設備上執行 SnapRestore 作業。還原完成後、您可以掛載系統並將其恢復至目前狀態。對於時間點恢復、您也可以還原 Wal 和歸檔記錄檔、然後 PostgreSQL 會決定最一致的點並自動恢復。

一致性群組是 ONTAP 中的一項功能、建議在有多個磁碟區掛載到單一執行個體或具有多個資料表空間的資料庫時使用。一致性群組快照可確保將所有磁碟區分組並加以保護。一致性群組可從 ONTAP 系統管理員有效管理、您甚至可以複製它來建立資料庫的執行個體複本、以供測試或開發之用。

如需一致性群組的詳細資訊、請參閱 ["NetApp 一致性群組總覽"](#)。

## PostgreSQL 資料保護軟體

適用於 PostgreSQL 資料庫的 NetApp SnapCenter 外掛程式、結合 Snapshot 和 NetApp FlexClone 技術、可為您提供下列優點：

- 快速備份與還原。

- 節省空間的複本。
- 能夠建置快速有效的災難恢復系統。

在下列情況下、您可能偏好選擇 NetApp 的優質備份合作夥伴、例如 Veeam Software 和 CommVault：



- 在異質環境中管理工作負載
- 將備份儲存至雲端或磁帶、以供長期保留
- 支援多種作業系統版本和類型

適用於 PostgreSQL 的 SnapCenter 外掛程式是社群支援的外掛程式、可在 NetApp Automation Store 取得設定與文件。透過 SnapCenter、使用者可以遠端備份資料庫、複製及還原資料。



# VMware

## VMware vSphere 搭配 ONTAP

### VMware vSphere 搭配 ONTAP

ONTAP 在將近 20 年來一直是 VMware vSphere 環境的領先儲存解決方案、並持續新增創新功能來簡化管理、同時降低成本。本文件介紹 ONTAP vSphere 的解決方案、包括最新的產品資訊和最佳實務做法、以簡化部署、降低風險及簡化管理。



本文件取代先前發佈的技術報告 [\\_TR-4597 : VMware vSphere for ONTAP](#) \_

最佳實務做法是輔助其他文件、例如指南和相容性清單。這些技術是根據實驗室測試和 NetApp 工程師與客戶廣泛的現場經驗所開發。它們可能不是每個環境中唯一能運作的支援實務做法、但通常是最簡單的解決方案、能滿足大多數客戶的需求。

本文件著重於在 vSphere 7.0 或更新版本上執行的 ONTAP (9.x) 最新版本中的功能。請參閱 ["NetApp 互通性對照表工具"](#) 和 ["VMware 相容性指南"](#) 以取得與特定版本相關的詳細資料。

為何 ONTAP 選擇適用於 vSphere 的呢？

有許多理由讓成千上萬的客戶選擇 ONTAP 作為 vSphere 的儲存解決方案、例如支援 SAN 和 NAS 傳輸協定的統一儲存系統、使用節省空間的快照功能提供強大的資料保護功能、以及豐富的工具來協助您管理應用程式資料。使用與 Hypervisor 分開的儲存系統、您可以卸載許多功能、並將 vSphere 主機系統的投資效益最大化。這種方法不僅能確保主機資源專注於應用程式工作負載、也能避免儲存作業對應用程式造成隨機效能影響。

搭配 vSphere 使用 VMware 是一項絕佳組合、可降低主機硬體與 VMware 軟體的費用。ONTAP 您也可以透過一致的高效能、以較低的成本保護資料。由於虛擬化工作負載是行動工作負載、因此您可以使用 Storage VMotion、在 VMFS、NFS 或 vVols 資料存放區之間移動 VM、探索不同的方法、所有這些都在同一個儲存系統上。

以下是客戶今日重視的關鍵因素：

- \*統一化儲存設備。\* 執行 ONTAP 此功能的系統以多種重要方式統一化。這種方法原本是指 NAS 和 SAN 兩種傳輸協定、ONTAP 而除了 NAS 的原始優勢之外、它仍是 SAN 的領導平台。在 vSphere 環境中、這種方法也可能代表虛擬桌面基礎架構 (VDI) 的統一化系統、以及虛擬伺服器基礎架構 (VSI)。執行 ONTAP VMware 軟體的系統通常比傳統企業陣列便宜、但在同一個系統中擁有進階的儲存效率功能來處理 VDI。此外、從 SSD 到 SATA、還能統一化各種儲存媒體、並將這些媒體輕鬆延伸到雲端。ONTAP 無需購買單一 Flash 陣列即可獲得效能、SATA 陣列可用於歸檔、而獨立的系統則可用於雲端。將它們緊密連結在一起。ONTAP
- \* 虛擬磁碟區和儲存原則型管理。\* NetApp 是 VMware 早期開發 vSphere 虛擬磁碟區 (vVols) 的設計合作夥伴、為 vVols 和 VMware vSphere API for Storage Aware (VASA) 提供架構輸入和早期支援。這種方法不僅能為 VMFS 帶來精細的 VM 儲存管理、也支援透過儲存原則型管理來自動化儲存資源配置。此方法可讓儲存架構設計師設計具有不同功能的儲存資源池、讓 VM 管理員輕鬆使用。這個解決方案是 vVol 擴充儲存產業的領導廠商、可在單一叢集內支援數十萬個 vVols、而企業陣列和小型 Flash 陣列廠商則可支援每個陣列數千個 vVols。ONTAP NetApp 也透過即將推出的 vVols 3.0 支援功能、推動精細 VM 管理的演進。
- \* 儲存效率。\* 雖然 NetApp 是第一批為正式作業工作負載提供重複資料刪除技術的公司、但這項創新技術並不是這方面的第一項或最後一項。它從快照開始、這是一種不具效能影響的空間效率資料保護機制、搭配 FlexClone 技術、可立即製作 VM 的讀取 / 寫入複本、以供正式作業和備份使用。NetApp 繼續提供內嵌功能、包括重複資料刪除、壓縮及零區塊重複資料刪除、讓昂貴的 SSD 發揮最大的儲存容量。最近、利用壓縮



技術、將較小的I/O作業和檔案封裝到磁碟區塊中的功能更為豐富。ONTAP這些功能的結合、讓客戶看到VSI的節約效益高達5：1、VDI的節約效益高達30：1。

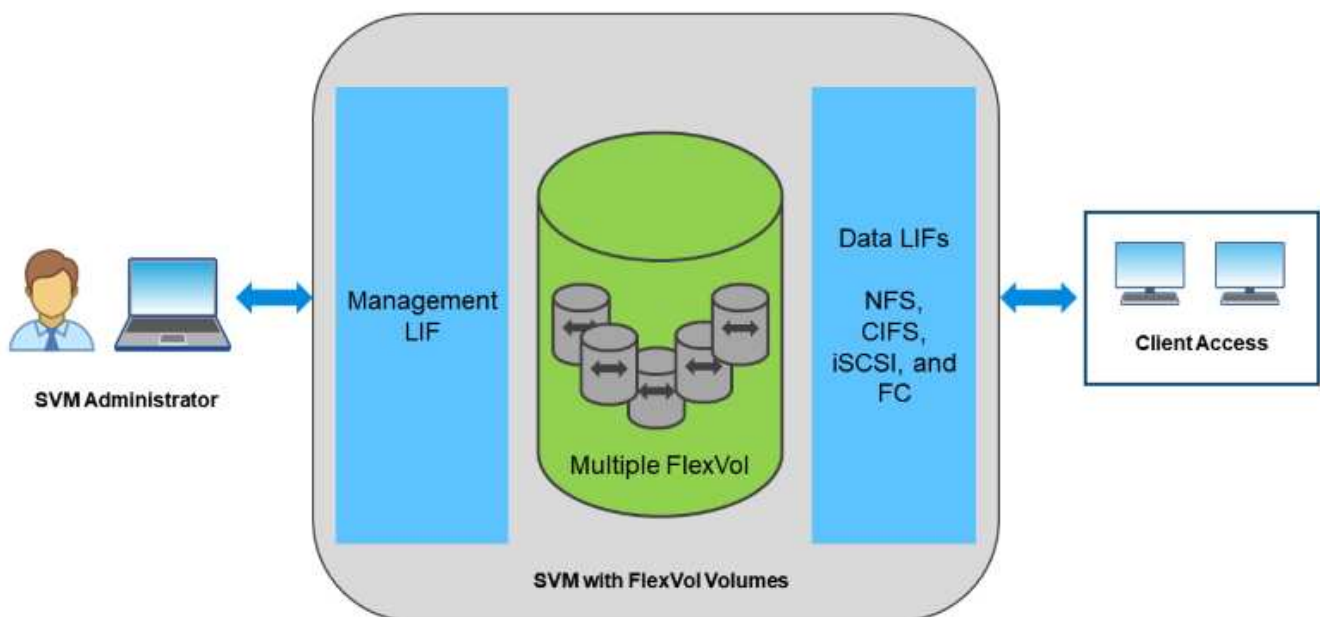
- \*混合雲\*無論是用於內部部署的私有雲、公有雲基礎架構、或是結合兩者最佳功能的混合雲、ONTAP 均可利用此解決方案協助您打造資料架構、以簡化及最佳化資料管理。從高效能All Flash系統開始著手、然後將它們與磁碟或雲端儲存系統結合、以提供資料保護和雲端運算。您可以選擇Azure、AWS、IBM或Google雲端、以最佳化成本並避免受限。視需要運用OpenStack和Container技術的進階支援。NetApp也提供雲端型備份 (SnapMirror Cloud、Cloud Backup Service VMware及Cloud Sync VMware)、以及FabricPool 適用於VMware的儲存分層與歸檔工具 (VMware®) ONTAP、協助降低營運成本、並充分運用雲端的廣泛應用。
- \*及更多資訊。\*善用NetApp AFF Sa系列陣列的極致效能、加速虛擬化基礎架構、同時管理成本。使用橫向擴充ONTAP 的叢集、享受完全不中斷營運的體驗、從維護到升級、到儲存系統的完整更換。使用NetApp加密功能保護閒置資料、無需額外成本。透過精細的服務品質功能、確保效能符合業務服務層級。它們都是業界領先的企業資料管理軟體 ONTAP 所提供的各種功能的一部分。

## 統一化儲存設備

NetApp ONTAP 透過簡化的軟體定義方法、統一化儲存設備、實現安全高效的管理、更高的效能、以及無縫的擴充性。這種方法可加強資料保護、並有效運用雲端資源。

這種統一化方法原本是指在單一儲存系統上同時支援 NAS 和 SAN 傳輸協定、而 ONTAP 則是 SAN 的領先平台、同時也是 NAS 的原始優勢。ONTAP 現在也提供 S3 物件傳輸協定支援。雖然 S3 不用於資料存放區、但您可以將它用於來賓應用程式。您可以在中深入瞭解 ONTAP 中的 S3 傳輸協定支援 "[S3組態總覽](#)"。

儲存虛擬機器 (SVM) 是 ONTAP 中安全的多租戶共享單元。這是一種邏輯結構、可讓用戶端存取執行 ONTAP 軟體的系統。SVM可透過邏輯介面 (LIF)、透過多種資料存取傳輸協定同時提供資料。SVM透過NAS 傳輸協定 (例如CIFS和NFS) 提供檔案層級的資料存取、並透過SAN傳輸協定 (例如iSCSI、FC/FCoE和NVMe) 提供區塊層級的資料存取。SVM 可以同時將資料單獨提供給 SAN 和 NAS 用戶端、也可以搭配 S3 使用。



在vSphere環境中、這種方法也可能代表虛擬桌面基礎架構 (VDI) 的統一化系統、以及虛擬伺服器基礎架構

(VSI)。執行ONTAP VMware軟體的系統通常比傳統企業陣列便宜、但在同一個系統中擁有進階的儲存效率功能來處理VDI。此外、從SSD到SATA、還能統一化各種儲存媒體、並將這些媒體輕鬆延伸到雲端。ONTAP無需購買單一 Flash 陣列即可獲得效能、SATA 陣列可用於歸檔、而獨立的系統則可用於雲端。將它們緊密連結在一起。ONTAP

- 注意：\* 如需更多有關 SVM、統一化儲存設備和用戶端存取的資訊、請參閱 "儲存虛擬化" 在VMware的VMware®文件中心。ONTAP

## 適用於VMware的虛擬化工具ONTAP

NetApp提供數種獨立式軟體工具、可搭配ONTAP 使用以管理您的虛擬化環境。

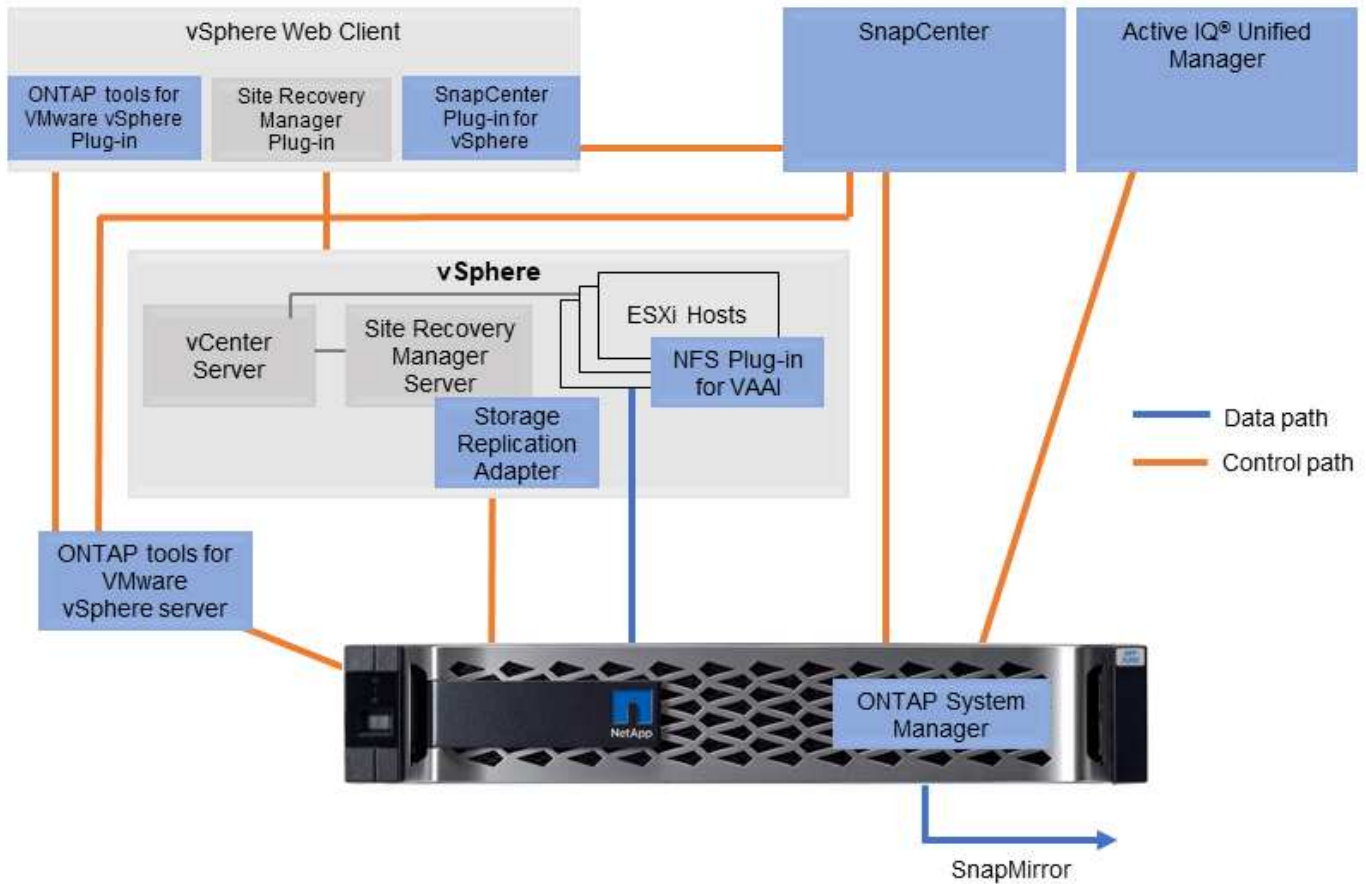
下列工具隨ONTAP 附於VMware授權中、不需額外付費。請參閱圖1、瞭解這些工具如何在vSphere環境中共同運作。

### VMware vSphere適用的工具ONTAP

VMware vSphere的VMware vSphere工具是一套搭配vSphere使用VMware vCenter儲存設備的工具。ONTAP ONTAPvCenter外掛程式先前稱為虛擬儲存主控台 (VSC)、可簡化儲存管理與效率功能、增強可用度、並降低儲存成本與營運成本、無論您使用SAN或NAS。它採用最佳實務做法來配置資料存放區、並針對NFS和區塊儲存環境最佳化ESXi主機設定。為獲得所有這些效益、NetApp建議您在ONTAP 使用vSphere搭配執行ONTAP VMware軟體的系統時、將這些VMware不完整的工具當作最佳實務做法。其中包括伺服器應用裝置、vCenter的使用者介面擴充、VASA Provider及儲存複寫介面卡。幾乎所有的功能都可以透過簡單的REST API來自動化、而大多數的現代化自動化工具都能耗用這些API ONTAP 。

- \* vCenter UI延伸功能\* ONTAP 此功能可簡化營運團隊和vCenter管理員的工作。\*此功能可在vCenter UI 中、直接使用易於使用的內容相關功能表來管理主機和儲存設備、資訊入口小程式和原生警示功能、以利簡化工作流程。
- \*適用於ONTAP VMware的VASA Provider for Sfor Sfor the。\*適用於ONTAP VMware的VASA Provider支援VMware vStorage API for Storage Aware (VASA) 架構。它是ONTAP VMware vSphere的VMware vSphere配套工具、可做為單一虛擬應用裝置、方便部署。VASA Provider將vCenter Server與ONTAP VMware連線、以協助資源配置及監控VM儲存設備。它可支援VMware虛擬磁碟區 (vVols)、管理儲存功能設定檔和個別VM vVols效能、以及監控容量和設定檔法規遵循的警示。
- 儲存複寫介面卡。SRA與VMware Site Recovery Manager (SRM) 搭配使用、可在正式作業與災難恢復站台之間管理資料複寫、並在不中斷營運的情況下測試DR複本。它有助於自動化探索、還原及重新保護等工作。其中包括適用於Windows SRM伺服器和SRM應用裝置的SRA伺服器應用裝置和SRA介面卡。

下圖說明ONTAP vSphere的各項功能。



### 適用於VMware VAAI的NFS外掛程式

適用於VMware VAAI的NetApp NFS外掛程式是ESXi主機的外掛程式、可讓ESXi主機在ONTAP 支援VMware的NFS資料存放區時、使用VAAI功能。它支援複製卸載、以進行複製作業、保留大型虛擬磁碟檔案的空間、以及快照卸載。將複本作業卸載到儲存設備並不一定能更快完成、但它確實能降低網路頻寬需求、並卸載CPU週期、緩衝區和佇列等主機資源。您可以使用ONTAP VMware vSphere的支援功能、在ESXi主機上安裝外掛程式、或是在支援的情況下安裝vSphere Lifecycle Manager (vLCM)。

### 虛擬磁碟區 (vVols) 與儲存原則型管理 (SPBM)

NetApp是VMware在開發vSphere虛擬磁碟區 (vVols) 方面的早期設計合作夥伴、為vVols和VMware vSphere API提供架構輸入和早期支援、以提高儲存感知 (VASA) 的能力。這種方法不僅將 VM 精細儲存管理帶入 VMFS、也支援透過儲存原則型管理 (SPBM) 來自動化儲存資源配置。

SPBM提供一個架構、可在虛擬化環境可用的儲存服務與透過原則配置的儲存元素之間、做為抽象層。此方法可讓儲存架構設計師設計具有不同功能的儲存資源池、讓VM管理員輕鬆使用。接著系統管理員可以根據已配置的儲存資源池來比對虛擬機器工作負載需求、以便精細控制每個VM或虛擬磁碟層級的各種設定。

這個解決方案是vVols擴充儲存產業的領導廠商、可在單一叢集內支援數十萬個vVols、而企業陣列和小型Flash陣列廠商則支援每個陣列數千個vVols。ONTAPNetApp也透過即將推出的vVols 3.0支援功能、推動VM精細管理的演進。



如需VMware vSphere虛擬磁碟區、SPBM及ONTAP VMware的詳細資訊、請參閱 ["TR-4400: VMware vSphere虛擬Volume ONTAP with VMware"](#)。

## 資料存放區與傳輸協定

### vSphere 資料存放區和傳輸協定功能概觀

使用七種傳輸協定將VMware vSphere連接至執行ONTAP VMware軟體的系統上的資料存放區：

- FCP
- FCoE
- NVMe / FC
- NVMe / TCP
- iSCSI
- NFS v3
- NFS v4.1

FCP、FCoE、NVMe/FC、NVMe/TCP和iSCSI是區塊傳輸協定、使用vSphere虛擬機器檔案系統（VMFS）將VM儲存在ONTAP 包含ONTAP FlexVol 在一個EesfVolume中的SesfLUN或NVMe命名空間內。請注意、從vSphere 7.0開始、VMware不再支援正式作業環境中的軟體FCoE。NFS是一種檔案傳輸協定、可將VM放入資料存放區（只是ONTAP 指不需要VMFS的功能）。SMB（CIFS）、iSCSI、NVMe / TCP或NFS也可直接從客體作業系統使用到ONTAP 不支援。

下表提供 vSphere 支援的傳統資料存放區功能與 ONTAP。此資訊不適用於vVols資料存放區、但通常適用於使用支援ONTAP 的版本資訊的vSphere 6.x及更新版本。您也可以諮詢 ["VMware 組態上限"](#) 以確認特定vSphere版本的特定限制。

功能/特色	FC/FCoE	iSCSI	NVMe	NFS
格式	VMFS或原始裝置對應（RDM）	VMFS或RDM	VMFS	不適用
資料存放區或LUN的最大數量	每個主機有1024個LUN	每個伺服器有1024個LUN	每個伺服器256個Namespaces	256 個掛載預設 NFS。最大磁碟區為8。使用VMware vSphere的VMware vSphere功能將其提升至256個。ONTAP
最大資料存放區大小	64TB	64TB	64TB	100TB FlexVol 以上的穩定區塊（含FlexGroup 不穩定區）
最大資料存放區檔案大小	62TB	62TB	62TB	62TB、含 ONTAP 9.12.1P2 及更新版本
每個LUN或檔案系統的最佳佇列深度	64-256	64-256	自動協商	請參閱中的NFS.MaxQuestue深度 <a href="#">"建議的ESXi主機和其他ONTAP 功能設定"</a> 。

下表列出支援的VMware儲存相關功能。

容量/功能	FC/FCoE	iSCSI	NVMe	NFS
vMotion	是的	是的	是的	是的
Storage VMotion	是的	是的	是的	是的
VMware HA	是的	是的	是的	是的
儲存分散式資源排程器 (SDR)	是的	是的	是的	是的
啟用資料保護 (VADP) 的VMware vStorage API備份軟體	是的	是的	是的	是的
VM內的Microsoft叢集服務 (Microsoft Cluster Service、英文) 或容錯移轉叢集	是的	是*	是*	不支援
容錯能力	是的	是的	是的	是的
Site Recovery Manager	是的	是的	否*	僅適用於v3 *
精簡配置的VM (虛擬磁碟)	是的	是的	是的	是的 此設定是不使用 VAAI 時、NFS 上所有 VM 的預設值。
VMware原生多重路徑	是的	是的	是、使用新的高效能外掛程式 (HPP)	NFS v4.1 工作階段主幹需要 ONTAP 9.14.1 及更新版本

下表列出支援ONTAP 的功能不完整的儲存管理功能。

功能/特色	FC/FCoE	iSCSI	NVMe	NFS
重複資料刪除技術	節省陣列成本	節省陣列成本	節省陣列成本	資料存放區的節約效益
資源隨需配置	資料存放區或RDM	資料存放區或RDM	資料存放區	資料存放區
調整資料存放區大小	僅成長	僅成長	僅成長	擴充、自動擴充及縮減
適用於Windows、Linux應用程式 (客體內) 的列舉外掛程式SnapCenter	是的	是的	否	是的
使用ONTAP VMware vSphere的VMware VMware vSphere的支援工具來監控及主機組態	是的	是的	否	是的

功能/特色	FC/FCoE	iSCSI	NVMe	NFS
使用ONTAP VMware vSphere的VMware vSphere的VMware 工具進行資源配置	是的	是的	否	是的

下表列出支援的備份功能。

功能/特色	FC/FCoE	iSCSI	NVMe	NFS
ONTAP 快照	是的	是的	是的	是的
SRM支援複寫備份	是的	是的	否*	僅適用於v3 *
Volume SnapMirror	是的	是的	是的	是的
VMDK映像存取	啟用VADP的備份軟體	啟用VADP的備份軟體	啟用VADP的備份軟體	啟用VADP的備份軟體、vSphere Client和vSphere Web Client資料存放區瀏覽器
VMDK檔案層級存取	啟用VADP的備份軟體、僅限Windows	啟用VADP的備份軟體、僅限Windows	啟用VADP的備份軟體、僅限Windows	啟用VADP的備份軟體和協力廠商應用程式
NDMP精細度	資料存放區	資料存放區	資料存放區	資料存放區或 VM

- NetApp建議將來賓iSCSI用於Microsoft叢集、而非在VMFS資料存放區中使用支援多寫入器的VMDK。Microsoft和VMware完全支援這種方法、ONTAP 提供優異的靈活度搭配使用VMware (SnapMirror至ONTAP 內部部署或雲端的等化系統)、易於設定和自動化、SnapCenter 並可透過VMware加以保護。vSphere 7新增叢集式VMDK選項。這與啟用多寫入器的VMDK不同、因為VMDK需要透過FC傳輸協定呈現資料存放區、而且此傳輸協定已啟用叢集式VMDK支援。其他限制也適用。請參閱 VMware 的 ["Windows Server容錯移轉叢集的設定"](#) 組態準則文件。

\*使用NVMe與NFS v4.1的資料存放區需要vSphere複寫。SRM不支援陣列型複寫。

#### 選擇儲存傳輸協定

執行ONTAP 支援所有主要儲存傳輸協定的系統、因此客戶可以根據現有和規劃的網路基礎架構和員工技能、選擇最適合自己環境的系統。NetApp測試通常顯示以類似線路速度執行的傳輸協定之間沒有什麼差異、因此最好將重點放在網路基礎架構和員工能力上、而不只是原始傳輸協定效能。

下列因素可能有助於考量選擇傳輸協定：

- \*目前的客戶環境。\*雖然IT團隊通常擅長管理乙太網路IP基礎架構、但並非所有人都擅長管理FC SAN架構。但是、使用非專為儲存流量設計的通用 IP 網路可能無法正常運作。請考量您所擁有的網路基礎架構、任何計畫性的改善、以及員工管理這些基礎架構的技能和可用度。
- \*易於設定。\*除了FC架構的初始組態設定 (額外的交換器和纜線、分區、以及HBA和韌體的互通性驗證) 之外、區塊傳輸協定也需要建立及對應LUN、以及由客體作業系統探索及格式化。NFS磁碟區建立及匯出之後、便會由ESXi主機掛載並準備好使用。NFS沒有特殊的硬體限制或韌體可管理。
- \*易於管理。\*有了SAN傳輸協定、如果需要更多空間、就必須採取幾個步驟、包括擴充LUN、重新掃描以探索新的大小、然後擴充檔案系統)。雖然可以擴充LUN、但減少LUN的大小並不可行、而且恢復未使用的空間可能需要額外的心力。NFS可輕鬆調整規模或縮減規模、儲存系統也能自動調整大小。SAN透過客體作業

系統修剪/取消對應命令提供空間回收、讓刪除檔案的空間可以傳回陣列。NFS資料存放區的這類空間回收較為困難。

- \*儲存空間的透明度。\*在NFS環境中、儲存使用率通常比較容易查看、因為精簡配置可立即回收節約效益。同樣地、相同資料存放區中的其他VM或其他儲存系統磁碟區也可立即使用重複資料刪除和複製的節約效益。NFS資料存放區的VM密度通常也較高、可減少資料存放區的管理數量、進而改善重複資料刪除的節約效益、並降低管理成本。

#### 資料存放區配置

可靈活建立VM和虛擬磁碟的資料存放區。ONTAP雖然ONTAP 使用VSC來配置vSphere的資料存放區時會套用許多功能不實的最佳實務做法（請參閱一節 "[建議的ESXi主機和其他ONTAP 功能設定](#)"）、以下是一些額外的考量準則：

- 部署vSphere搭配ONTAP 使用不間斷的NFS資料存放區、可實現高效能且易於管理的實作、提供VM對資料存放區的比率、而這些比率無法透過區塊型儲存傳輸協定取得。此架構可減少相關資料存放區數量、使資料存放區密度增加十倍。雖然較大的資料存放區可提升儲存效率並提供營運效益、但請考慮使用至少四個資料存放區FlexVol（VMware Volume）、將VM儲存在單ONTAP 一的VMware控制器上、以從硬體資源中獲得最大效能。此方法也可讓您建立具有不同恢復原則的資料存放區。根據業務需求、部分備份或複製的頻率可能會比其他更高。由於資料存放區FlexGroup 是依設計進行擴充、因此不需要使用多個資料存放區來提升效能。
- NetApp 建議對大多數 NFS 資料存放區使用 FlexVol Volume。從 ONTAP 9.8 FlexGroup 磁碟區開始、也支援作為資料存放區使用、一般建議在某些使用案例中使用。一般不建議使用其他 ONTAP 儲存容器、例如 qtree、因為目前 VMware vSphere 的 ONTAP 工具或 VMware vSphere 的 NetApp SnapCenter 外掛程式都不支援這些容器。也就是說、將資料存放區部署為單一磁碟區中的多個 qtree、對於高度自動化的環境來說可能很有用、因為資料存放區層級配額或 VM 檔案複本可以讓環境受益。
- 適用於不只FlexVol 4TB、更能滿足8TB的需求。這種規模對於效能、管理簡易性和資料保護來說、是一個很好的平衡點。從小規模開始（例如4TB）、視需要擴充資料存放區（最高100TB）。較小的資料存放區可更快從備份或災難後恢復、並可在叢集之間快速移動。請考慮使用ONTAP 不同步自動調整大小、以便在使用空間變更時自動擴充及縮小磁碟區。VMware vSphere資料存放區資源配置精靈的「VMware vSphere資料存放區資源配置精靈」預設會針對新的資料存放區使用自動調整大小。ONTAP您可以使用System Manager或命令列、進一步自訂「成長」和「縮減」臨界值、以及最大和最小大小。
- 此外、VMFS資料存放區也可以設定LUN、以供FC、iSCSI或FCoE存取。VMFS可讓叢集中的每個ESX伺服器同時存取傳統LUN。VMFS資料存放區的大小最多可達64TB、最多可包含32個2TB LUN（VMFS 3）或單一64TB LUN（VMFS 5）。大部分系統的LUN大小僅為16TB、ONTAP All SAN陣列系統的LUN大小上限為12TB。因此、在ONTAP 大多數的不實系統上、可使用四個16TB LUN來建立最大大小的VMFS 5資料存放區。雖然多個LUN的高I/O工作負載（使用高階FAS 的功能或AFF 功能性系統）可獲得效能優勢、但由於建立、管理及保護資料存放區LUN的管理複雜度增加、以及提高可用度風險、因此這項優勢已被抵銷。NetApp 一般建議針對每個資料存放區使用單一大型LUN、而且只有在需要超越16TB資料存放區的情況下才需要跨距。與NFS一樣、請考慮使用多個資料存放區（Volume）、在單ONTAP 一的VMware控制器上發揮最大效能。
- 老舊的客體作業系統（OS）需要與儲存系統一致、才能獲得最佳效能和儲存效率。然而、Microsoft和Linux 經銷商（例如Red Hat）所支援的現代化作業系統不再需要調整、以使檔案系統分割區與虛擬環境中基礎儲存系統的區塊保持一致。如果您使用的是可能需要調整的舊作業系統、請在NetApp支援知識庫中搜尋文章、使用「VM對齊」、或向NetApp銷售或合作夥伴聯絡人索取TR-3747的複本。
- 避免在來賓作業系統中使用重組公用程式、因為這不會帶來效能效益、也會影響儲存效率和快照空間使用量。也請考慮在客體作業系統中關閉虛擬桌面的搜尋索引。
- 以創新的儲存效率功能引領業界、讓您充分發揮可用磁碟空間的效益。ONTAP利用預設的即時重複資料刪除與壓縮技術、支援更高的效率。AFF資料會在集合體中的所有磁碟區中進行重複資料刪除、因此您不再需要將類似的作業系統和類似的應用程式群組在單一資料存放區中、以達到最大的節約效益。
- 在某些情況下、您甚至不需要資料存放區。為獲得最佳效能與管理能力、請避免將資料存放區用於高I/O應用

程式、例如資料庫和某些應用程式。相反地、請考慮使用來賓擁有的檔案系統、例如NFS或iSCSI檔案系統、由來賓或RDM管理。如需特定的應用程式指南、請參閱適用於您應用程式的NetApp技術報告。例如、["Oracle資料庫ONTAP"](#) 提供虛擬化的相關章節、並提供實用的詳細資料。

- 一流磁碟（或改良的虛擬磁碟）可讓vCenter管理的磁碟獨立於vSphere 6.5及更新版本的VM。雖然主要是由API管理、但在vVols上也很實用、尤其是由OpenStack或Kubernetes工具管理時。支援的項目包括ONTAP VMware ONTAP vSphere的VMware vSphere的支援功能和VMware vSphere的支援功能。

#### 資料存放區與VM移轉

將VM從另一個儲存系統上的現有資料存放區移轉至ONTAP 支援區時、請謹記以下幾項實務做法：

- 使用Storage VMotion將大部分虛擬機器移至ONTAP VMware。這種方法不僅不中斷虛ONTAP 擬機器的執行、還能讓諸如即時重複資料刪除和壓縮等儲存效率功能、在資料移轉時處理資料。請考慮使用vCenter功能從清單清單清單中選取多個VM、然後在適當的時間排程移轉（按一下「Actions」（動作）時使用Ctrl鍵）。
- 雖然您可以仔細規劃移轉至適當的目的地資料存放區、但通常較容易大量移轉、然後視需要組織。如果您有特定的資料保護需求、例如不同的 Snapshot 排程、您可能會想要使用此方法來引導您移轉至不同的資料存放區。
- 大多數VM及其儲存設備可能會在執行（Hot）時移轉、但從其他儲存系統移轉附加（非資料存放區）儲存設備（例如ISO、LUN或NFS磁碟區）可能需要冷移轉。
- 需要更謹慎移轉的虛擬機器包括使用附加儲存設備的資料庫和應用程式。一般而言、請考慮使用應用程式的工具來管理移轉。對於Oracle、請考慮使用Oracle工具（例如RMAN或ASM）來移轉資料庫檔案。請參閱["TR-4534"](#) 以取得更多資訊。同樣地、對於SQL Server、請考慮使用SQL Server Management Studio或NetApp工具、例如SnapManager 適用於SQL Server或SnapCenter VMware。

#### VMware vSphere適用的工具ONTAP

搭配執行ONTAP VMware vCenter軟體的系統使用vSphere時、最重要的最佳實務做法是安裝及使用ONTAP VMware vSphere外掛程式（前身為虛擬儲存主控台）的VMware VMware vSphere資訊工具。無論使用SAN或NAS、此vCenter外掛程式都能簡化儲存管理、提升可用度、並降低儲存成本和營運成本。它採用最佳實務做法來配置資料存放區、並針對多重路徑和HBA逾時最佳化ESXi主機設定（如附錄B所述）。由於它是 vCenter 外掛程式、因此可用於所有連線至 vCenter 伺服器的 vSphere Web 用戶端。

外掛程式也能協助您在ONTAP vSphere環境中使用其他的功能。它可讓您安裝適用於 VMware VAAI 的 NFS 外掛程式、以便將複本卸載至 ONTAP 進行虛擬機器複製作業、保留大型虛擬磁碟檔案的空間、以及 ONTAP 快照卸載。

外掛程式也是VASA Provider許多功能的管理介面、ONTAP 可支援vVols的儲存原則型管理。在登錄VMware vSphere的VMware vSphere基礎架構工具之後ONTAP、請使用它來建立儲存功能設定檔、將其對應至儲存設備、並確保資料存放區在一段時間內符合設定檔的要求。VASA Provider也提供一個介面、可用來建立及管理VVOL資料存放區。

一般而言、NetApp建議在ONTAP vCenter內使用VMware vSphere的VMware vCenter功能的VMware vCenter功能、來配置傳統和vVols資料存放區、以確保遵循最佳實務做法。

#### 一般網路

使用vSphere搭配執行ONTAP VMware軟體的系統時、設定網路設定很簡單、而且類似於其他網路組態。以下是幾點需要考量的事項：

- 將儲存網路流量與其他網路區隔。使用專屬的VLAN或獨立的交換器來儲存、即可建立獨立的網路。如果儲存網路共用實體路徑（例如上行鏈路）、您可能需要QoS或額外的上行鏈路連接埠、以確保有足夠的頻寬。



請勿將主機直接連線至儲存設備；請使用交換器來建立備援路徑、並讓 VMware HA 在不需介入的情況下運作。請參閱 ["直接連線網路"](#) 以取得更多資訊。

- 如果您的網路需要並支援巨型框架、尤其是使用iSCSI時、可以使用巨型框架。如果使用、請確定在儲存設備和ESXi主機之間的路徑中、所有網路裝置、VLAN等上的設定都相同。否則、您可能會看到效能或連線問題。MTU也必須在ESXi虛擬交換器、VMkernel連接埠、以及每ONTAP 個節點的實體連接埠或介面群組上設定相同。
- NetApp僅建議停用ONTAP 叢集內叢集網路連接埠上的網路流量控制。對於用於資料流量的其餘網路連接埠、NetApp並未提出其他最佳實務做法建議。您應視需要啟用或停用。請參閱 ["TR-4182"](#) 以取得流程控制的更多背景資訊。
- 當ESXi和ONTAP VMware ESXi儲存陣列連接至乙太網路儲存網路時、NetApp建議將這些系統連接的乙太網路連接埠設定為快速擴充樹狀傳輸協定（RSTP）邊緣連接埠、或使用Cisco PortFast功能。NetApp建議在使用Cisco PortFast功能的環境中、啟用跨距樹狀結構PortFast主幹功能、並在ESXi伺服器或ONTAP VMware®儲存陣列上啟用802.1Q VLAN主幹連線。
- NetApp建議下列連結集合最佳實務做法：
  - 使用交換器、透過 Cisco 的 Virtual PortChannel（VPC）等多機箱連結集合群組方法、在兩個獨立的交換器機箱上支援連接埠的連結集合。
  - 除非您使用已設定LACP的DVSwitches 5.1或更新版本、否則請停用連接至ESXi的交換器連接埠LACP。
  - 使用 LACP 為具有連接埠或 IP 雜湊的動態多重模式介面群組的 ONTAP 儲存系統建立連結集合體。請參閱 ["網路管理"](#) 以取得進一步指引。
  - 在 ESXi 上使用靜態連結集合（例如、EtherChannel）和標準 vSwitch、或是搭配 vSphere Distributed Switch 使用 LACP 型連結集合時、請使用 IP 雜湊成組原則。如果未使用連結集合、請改用「根據來源虛擬連接埠 ID 建立路由」。

下表提供網路組態項目的摘要、並指出套用設定的位置。

項目	ESXi	交換器	節點	SVM
IP 位址	VMkernel	否*	否*	是的
連結集合體	虛擬交換器	是的	是的	否*
VLAN	VMkernel和VM連接埠群組	是的	是的	否*
流程控制	NIC	是的	是的	否*
跨距樹狀結構	否	是的	否	否
MTU（用於巨型框架）	虛擬交換器與VMkernel連接埠（9000）	是（設為上限）	有（9000）	否*
容錯移轉群組	否	否	是（建立）	是（選取）

- SVM lifs連接到具有VLAN、MTU及其他設定的連接埠、介面群組或VLAN介面。不過、這些設定不會在SVM層級進行管理。

這些裝置擁有自己的IP位址進行管理、但這些位址並未用於ESXi儲存網路環境。

## SAN (FC、FCoE、NVMe/FC、iSCSI)、RDM

NetApp ONTAP 使用 iSCSI、光纖通道傳輸協定 (FCP 或 FC 簡稱) 和 NVMe over Fabrics (NVMe of)、為 VMware vSphere 提供企業級區塊儲存。以下是使用 vSphere 和 ONTAP 實作 VM 儲存區塊傳輸協定的最佳實務做法。

在 vSphere 中、有三種使用區塊儲存 LUN 的方法：

- 使用 VMFS 資料存放區
- 使用原始裝置對應 (RDM)
- 由軟體啟動器從 VM 客體作業系統存取及控制的 LUN

VMFS 是高效能的叢集式檔案系統、可提供共用儲存資源池的資料存放區。VMFS 資料存放區可設定為使用 FC、iSCSI、FCoE 或 NVMe 命名空間存取 LUN、使用 NVMe / FC 或 NVMe / TCP 通訊協定存取。VMFS 可讓叢集中的每個 ESX 伺服器同時存取儲存設備。從 ONTAP 9.12.1P2 開始、LUN 大小上限通常為 128TB (ASA 系統則為較早版本)；因此、使用單一 LUN 可建立最大大小為 64TB 的 VMFS 5 或 6 資料存放區。

vSphere 內建多個儲存裝置路徑的支援功能、稱為原生多重路徑 (NMP)。NMP 可偵測支援儲存系統的儲存類型、並自動設定 NMP 堆疊以支援使用中儲存系統的功能。

NMP 和 ONTAP 都支援非對稱邏輯單元存取 (ALUA)、可協調最佳化和非最佳化的路徑。在本功能中、ALUA 最佳化路徑會使用主控所存取 LUN 的節點上的目標連接埠、遵循直接資料路徑。ONTAP 預設會在 vSphere 和 ONTAP VMware 中同時開啟 ALUA。NMP 將 ONTAP 叢集識別為 ALUA、並使用 ALUA 儲存陣列類型外掛程式 (VMW\_SATP\_ALUA) 並選取循環路徑選擇外掛程式 (VMW\_PSP\_RR)。

ESXi 6 最多可支援 256 個 LUN、並可支援多達 1、024 條通往 LUN 的總路徑。ESXi 不會看到超出這些限制的任何 LUN 或路徑。假設 LUN 數量上限、則路徑限制允許每個 LUN 有四個路徑。在更大 ONTAP 的實體叢集中、可以在 LUN 限制之前達到路徑限制。為了解決此限制、ONTAP 支援 8.3 版及更新版本中的選擇性 LUN 對應 (SLM)、

對於向指定 LUN 通告路徑的節點、SLM 會有限制。NetApp 最佳實務做法是每個 SVM 每個節點至少有一個 LIF、並使用 SLM 來限制通告給裝載 LUN 及其 HA 合作夥伴之節點的路徑。雖然存在其他路徑、但預設不會通告這些路徑。您可以使用新增和移除在 SLMs 中的報告節點引數來修改通告的路徑。請注意、在 8.3 之前的版本中建立的 LUN 會通告所有路徑、而且必須加以修改、才能只向主機 HA 配對通告路徑。如需更多關於 SLM、請參閱第 5.9 節 "TR-4080"。先前的連接埠集方法也可用於進一步減少 LUN 的可用路徑。PortSets 可減少 igroup 中的啟動器可透過哪些可見路徑來查看 LUN、進而提供協助。

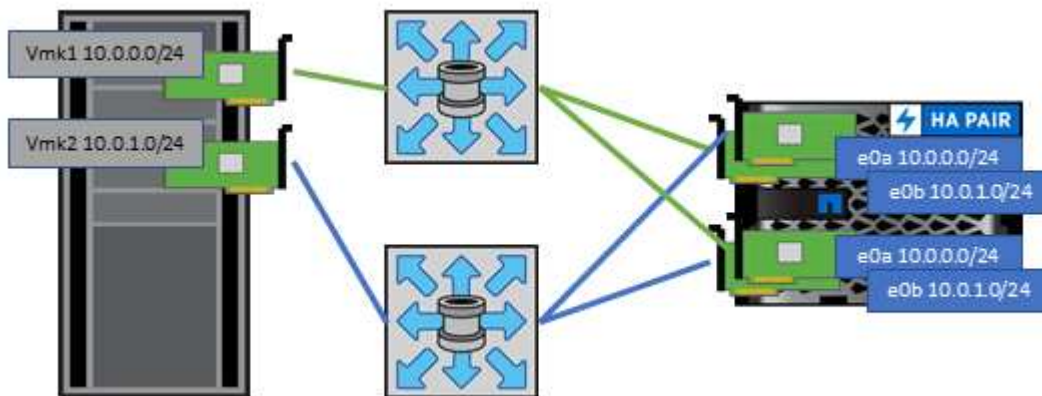
- 根據預設、會啟用 SLM。除非您使用連接埠集、否則不需要額外的組態。
- 對於在 Data ONTAP 8.3 之前建立的 LUN、請執行手動套用 SLM `lun mapping remove-reporting-nodes` 用於移除 LUN 報告節點、並限制 LUN 存取 LUN 所屬節點及其 HA 合作夥伴的命令。

區塊傳輸協定 (iSCSI、FC 和 FCoE) 使用 LUN ID 和序號以及唯一名稱來存取 LUN。FC 和 FCoE 使用全球名稱 (WWNN 和 WWPN)、iSCSI 則使用 iSCSI 合格名稱 (IQN)。儲存設備內部的 LUN 路徑對區塊傳輸協定毫無意義、而且不會出現在傳輸協定的任何位置。因此、只包含 LUN 的磁碟區根本不需要內部掛載、而包含資料存放區所用 LUN 的磁碟區則不需要使用交會路徑。NVMe 子系統 ONTAP 的運作方式類似。

其他應考慮的最佳實務做法：

- 請確定 ONTAP 已為叢集中每個節點上的每個 SVM 建立邏輯介面 (LIF)、以達到最大可用度和行動性。最佳實務做法是每個節點使用兩個實體連接埠和 LIF、每個光纖使用一個連接埠。ONTAP ALUA 可用來剖析路徑、識別作用中最佳化 (直接) 路徑、以及作用中未最佳化路徑。ALUA 用於 FC、FCoE 和 iSCSI。

- 對於iSCSI網路、當存在多個虛擬交換器時、請在不同的網路子網路上使用多個VMkernel網路介面搭配NIC群組。您也可以使用多個實體NIC來連接至多個實體交換器、以提供HA並提高處理量。下圖提供多重路徑連線的範例。在靜態中ONTAP、設定單一模式介面群組以容錯移轉兩個或多個連結連接至兩個或多個交換器、或使用LACP或其他連結集合技術搭配多重模式介面群組、以提供HA及連結集合的優點。
- 如果 ESXi 中使用挑戰握手驗證傳輸協定 ( CHAP ) 進行目標驗證、則也必須使用 CLI 在 ONTAP 中進行設定 (vserver iscsi security create) 或使用 System Manager (在 Storage (儲存) > SVM (SVM) > SVM Settings (SVM 設定) > Protocols (傳輸協定) > iSCSI (iSCSI) 下編輯啟動器安全性)。
- 使用VMware vSphere的VMware vCenter工具來建立及管理LUN和群組。ONTAP外掛程式會自動決定伺服器的WWPN、並建立適當的igroup。它也會根據最佳實務做法來設定LUN、並將其對應至正確的igroup。
- 請謹慎使用 RDM 、因為它們可能較難管理、而且也會使用路徑、而路徑的限制如前所述。支援這兩種LUN ONTAP "實體與虛擬相容模式" RDM。
- 如需更多關於將NVMe/FC搭配vSphere 7.0使用的資訊、請參閱此 "NVMe / FC主機組態指南ONTAP" 和 "TR-4684"下圖說明從vSphere主機到ONTAP VMware LUN的多重路徑連線能力。



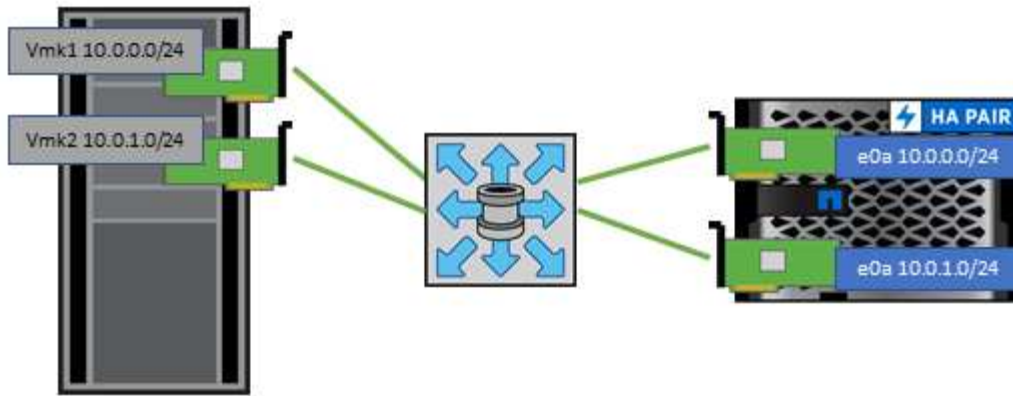
## NFS

NetApp ONTAP 是一種企業級橫向擴充 NAS 陣列、其中包括許多其他功能。ONTAP 讓 VMware vSphere 能夠從許多 ESXi 主機同時存取 NFS 連線的資料存放區、遠超出對 VMFS 檔案系統的限制。如中所述、將 NFS 搭配 vSphere 使用可提供一些易用性和儲存效率可見度效益 "資料存放區" 區段。

搭配ONTAP vSphere使用VMware NFS時、建議採用下列最佳實務做法：

- 在叢集中的每個節點上、為每個SVM使用單一邏輯介面 (LIF) ONTAP。不再需要過去針對每個資料存放區的LIF建議。雖然直接存取 (同一個節點上的 LIF 和資料存放區) 是最佳的選擇、但別擔心間接存取、因為效能影響通常很小 (微秒)。
- VMware自VMware Infrastructure 3以來就一直支援NFSv3。vSphere 6.0新增對NFSv4.1的支援、可啟用某些進階功能、例如Kerberos安全性。NFSv3使用用戶端鎖定功能時、NFSv4.1會使用伺服器端鎖定功能。雖然可以透過這兩種傳輸協定匯出一個資料區、但ESXi只能透過一個傳輸協定掛載。ONTAP此單一傳輸協定掛載並不排除其他ESXi主機透過不同版本掛載相同的資料存放區。請務必指定要在掛載時使用的傳輸協定版本、以便所有主機使用相同版本、因此鎖定樣式相同。請勿在主機之間混合使用NFS版本。如有可能、請使用主機設定檔檢查是否符合規定。
  - 由於NFSv3與NFSv4.1之間沒有自動資料存放區轉換、因此請建立新的NFSv4.1資料存放區、並使用Storage VMotion將VM移轉至新的資料存放區。

- 請參閱中的NFS v4.1互通性表附註 "[NetApp互通性對照表工具](#)" 支援所需的特定ESXi修補程式層級。
- VMware 支援從 vSphere 8.0U2 開始與 NFSv3 進行 nconnect 。如需 nconnect 的詳細資訊、請參閱 "[NFSv3 nConnect 功能搭配 NetApp 和 VMware](#)"
- NFS匯出原則用於控制vSphere主機的存取。您可以將一個原則與多個磁碟區（資料存放區）搭配使用。使用NFSv3時、ESXi會使用sys（UNIX）安全樣式、並需要root掛載選項來執行VM。在現階段、此選項稱為超級使用者、使用超級使用者選項時、不需要指定匿名使用者ID。ONTAP請注意、匯出具有不同值的原則規則 -anon 和 -allow-suid 可能會導致SVM發現ONTAP 問題、因為使用這些功能。以下是範例原則：
  - 存取傳輸協定：NFS（包括nfs3和nfs4）
  - 用戶端配對規格：192 · 168 · 42 · 21
  - RO存取規則：系統
  - RW 存取規則
  - 匿名 UID
  - 超級使用者：sys
- 如果使用 NetApp NFS 外掛程式 for VMware VAAI 、則通訊協定應設為 nfs 而非 nfs3 匯出原則規則建立或修改時。即使資料傳輸協定是 NFSv3 、VAAI 複製卸載功能仍需要 NFSv4 傳輸協定才能運作。將通訊協定指定為 nfs 包括 NFSv3 和 NFSv4 版本。
- NFS資料存放區磁碟區是從SVM的根磁碟區連結而來、因此ESXi也必須能夠存取根磁碟區、才能瀏覽及掛載資料存放區磁碟區。根 Volume 和資料存放區磁碟區交會嵌套的任何其他磁碟區的匯出原則、都必須包含ESXi 伺服器授予其唯讀存取權的規則或規則。以下是根 Volume 的範例原則、也使用 VAAI 外掛程式：
  - 存取傳輸協定：NFS（包括nfs3和nfs4）
  - 用戶端配對規格：192 · 168 · 42 · 21
  - RO存取規則：系統
  - RW存取規則：Never（root Volume的最佳安全性）
  - 匿名 UID
  - 超級使用者：sys（也適用於採用VAAI的根Volume）
- 使用ONTAP VMware vSphere的VMware Infrastructure（最重要的最佳實務做法）：
  - 使用VMware vSphere的VMware VMware VMware vSphere功能來配置資料存放區、因為它能自動簡化匯出原則的管理。ONTAP
  - 使用外掛程式為VMware叢集建立資料存放區時、請選取叢集而非單一ESX伺服器。此選項會觸發IT自動將資料存放區掛載至叢集中的所有主機。
  - 使用外掛程式掛載功能、將現有的資料存放區套用至新的伺服器。
  - 如果不使用ONTAP VMware vSphere的VMware vSphere功能、請針對所有伺服器或需要額外存取控制的每個伺服器叢集、使用單一匯出原則。
- 雖然供應彈性的Volume命名空間結構、可利用交會在樹狀結構中排列磁碟區、但這種方法對vSphere沒有任何價值。ONTAP無論儲存設備的命名空間階層為何、它都會在資料存放區根目錄中為每個VM建立一個目錄。因此、最佳實務做法是將vSphere磁碟區的交會路徑掛載到SVM的根磁碟區、這就是ONTAP VMware vSphere的VMware vSphere功能如何配置資料存放區。沒有巢狀結點路徑也表示除了根磁碟區之外、沒有任何磁碟區相依於任何磁碟區、即使是刻意將磁碟區離線或銷毀、也不會影響其他磁碟區的路徑。
- 對於NFS資料存放區上的NTFS分割區、4K區塊大小是可以的。下圖說明從vSphere主機連線至ONTAP VMware NFS資料存放區的能力。



下表列出NFS版本及支援的功能。

vSphere功能	NFSv3.	NFSv4.1
vMotion與Storage vMotion	是的	是的
高可用度	是的	是的
容錯能力	是的	是的
DRS	是的	是的
主機設定檔	是的	是的
儲存DRS	是的	否
儲存I/O控制	是的	否
SRM	是的	否
虛擬磁碟區	是的	否
硬體加速 (VAAI)	是的	是的
Kerberos驗證	否	是 (vSphere 6.5及更新版本增強支援AES、krb5i)
多重路徑支援	否	是的

### 資料量FlexGroup

將 ONTAP 和 FlexGroup Volume 搭配 VMware vSphere 使用、即可建立簡單且可擴充的資料存放區、充分發揮整個 ONTAP 叢集的完整功能。

ONTAP 9.8 搭配適用於 VMware vSphere 9.8 的 ONTAP 工具、以及適用於 VMware 4.4 版本的 SnapCenter 外掛程式、新增了對 vSphere 中 FlexGroup Volume 備份資料存放區的支援。FlexGroup Volume 可簡化大型資料存放區的建立、並在 ONTAP 叢集上自動建立必要的分散式組成磁碟區、讓 ONTAP 系統發揮最大效能。

如需 FlexGroup Volume 的詳細資訊、請參閱 ["FlexCache 與 FlexGroup Volume 技術報告"](#)。

如果您需要具備完整 ONTAP 叢集功能的單一可擴充 vSphere 資料存放區、或是擁有可從全新 FlexGroup 複製機制獲益的大型複製工作負載、請將 FlexGroup Volume 搭配 vSphere 一起使用。

## 複本卸載

除了針對 vSphere 工作負載進行廣泛的系統測試之外、ONTAP 9.8 也為 FlexGroup 資料存放區新增了複本卸載機制。這套新系統使用改良的複製引擎、在背景中的成員之間複製檔案、同時允許存取來源和目的地。然後使用此本機快取、根據需要快速產生 VM 複本。

若要啟用 FlexGroup 最佳化複本卸載、請參閱 ["如何設定 ONTAP FlexGroups 以允許 VAAI 複本卸載"](#)

您可能會發現、如果您使用的是 VAAI 複製、但複製的複製量不足以保持快取溫暖、則您的複本可能不會比主機複本快。如果是這種情況、您可以調整快取逾時、以更符合您的需求。

請考慮下列案例：

- 您已經建立了 8 個組成要素的新 FlexGroup
- 新 FlexGroup 的快取逾時設定為 160 分鐘

在此案例中、要完成的前 8 個複本將為完整複本、而非本機檔案複本。在 160 秒逾時過期之前、對該 VM 進行任何額外的複製、都會以循環方式使用每個組成要素內部的檔案複製引擎、以建立幾乎立即的複本、並在組成的磁碟區之間平均分配。

Volume 接收的每個新複製工作都會重設逾時。如果範例 FlexGroup 中的組成磁碟區在逾時之前未收到複製要求、則會清除該特定 VM 的快取、而且需要再次填入該磁碟區。此外、如果原始複本的來源變更（例如、您已更新範本）、則每個成分的本機快取都會失效、以避免任何衝突。如前所述、快取是可調整的、可設定以符合您環境的需求。

如需搭配 VAAI 使用 FlexGroups 的詳細資訊、請參閱知識庫文章：["VAAI：快取如何與 FlexGroup 磁碟區搭配運作？"](#)

在無法充分利用 FlexGroup 快取、但仍需要快速跨磁碟區複製的環境中、請考慮使用 vVols。使用 vVols 進行跨磁碟區複製的速度比使用傳統資料存放區快得多、而且不依賴快取。

## QoS 設定

支援使用 ONTAP 系統管理員或叢集 Shell 在 FlexGroup 層級設定 QoS、但無法提供 VM 認知或 vCenter 整合。

QoS（最大 / 最小 IOPS）可在 vCenter UI 中的個別虛擬機器或當時資料存放區中的所有虛擬機器上設定、或使用 ONTAP 工具透過 REST API 設定。在所有 VM 上設定 QoS 會取代任何個別 VM 設定。設定未來不會延伸至新的或移轉的 VM；您可以在新的 VM 上設定 QoS、或是將 QoS 重新套用至資料存放區中的所有 VM。

請注意、VMware vSphere 會將 NFS 資料存放區的所有 IO 視為每個主機的單一佇列、而一個 VM 上的 QoS 節流會影響同一個資料存放區中其他 VM 的效能。這與 vVols 形成對照、vVols 可在移轉至其他資料存放區時維持其 QoS 原則設定、並在節流時不會影響其他 VM 的 IO。

## 指標

ONTAP 9.8 也為 FlexGroup 檔案新增了檔案型效能指標（IOPS、處理量和延遲）、這些指標可在適用於 VMware vSphere 儀表板和 VM 報告的 ONTAP 工具中檢視。VMware vSphere 外掛程式的支援功能也可讓您使用最大和/或最小 IOPS 的組合來設定服務品質（QoS）規則。ONTAP 這些設定可以跨資料存放區中的所有 VM 進行設定、也可以針對特定 VM 個別設定。

- 使用 ONTAP 工具來建立 FlexGroup 資料存放區、以確保以最佳方式建立 FlexGroup、並將匯出原則設定為符合您的 vSphere 環境。不過、使用 ONTAP 工具建立 FlexGroup Volume 之後、您會發現 vSphere 叢集中的所有節點都使用單一 IP 位址來掛載資料存放區。這可能導致網路連接埠出現瓶頸。若要避免此問題、請卸載資料存放區、然後使用標準 vSphere 資料存放區精靈、使用在 SVM 上的整個生命體之間進行負載平衡的循環 DNS 名稱來重新掛載資料存放區。重新掛載之後、ONTAP 工具將再次能夠管理資料存放區。如果 ONTAP 工具無法使用、請使用 FlexGroup 預設值、並依照中的準則建立匯出原則 "[資料存放區和傳輸協定 - NFS](#)"。
- 調整 FlexGroup VMware 資料存放區規模時、請記住 FlexGroup、此功能包含 FlexVol 多個較小的、可建立較大命名空間的支援區。因此、將資料存放區大小調整為至少 8 倍（假設預設為 8 個組成要素）、即最大 VMDK 檔案的大小、加上 10-20% 的未使用保留空間、以便靈活地重新平衡。例如、如果您的環境中有 6TB VMDK、請將 FlexGroup 資料存放區大小調整為不小於 52.8 TB（6X8+10%）。
- VMware 和 NetApp 支援從 ONTAP 9.14.1 開始的 NFSv4.1 工作階段主幹。如需特定版本詳細資料、請參閱 NetApp NFS 4.1 互通性對照表附註。NFSv3 不支援多個實體路徑到一個 Volume、但支援從 vSphere 8.0U2 開始的 nconnect。如需 nconnect 的詳細資訊、請參閱 "[NFSv3 nConnect 功能搭配 NetApp 和 VMware](#)"。
- 使用適用於 VMware VAAI 的 NFS 外掛程式進行複本卸載。請注意、雖然如前所述、在 FlexGroup 資料存放區內增強複製功能、但在 FlexVol 和 / 或 FlexGroup 磁碟區之間複製 VM 時、ONTAP 並未提供與 ESXi 主機複本相比的顯著效能優勢。因此、在決定使用 VAAI 或 FlexGroups 時、請考慮您的複製工作負載。修改組成磁碟區數量是最佳化 FlexGroup 型複製的一種方法。如同調整先前提到的快取逾時。
- 使用適用於 VMware vSphere 9.8 或更新版本的 ONTAP 工具、使用 ONTAP 度量（儀表板和 VM 報告）來監控 FlexGroup VM 的效能、並在個別 VM 上管理 QoS。目前無法透過 ONTAP REST 指令或 API 取得這些指標。
- SnapCenter Plug-in for VMware vSphere 4.4 版及更新版本支援在主要儲存系統上的 FlexGroup 資料存放區中備份及還原 VM。4.6 號選擇控制閥增加了對基於 FlexGroup 的數據存儲的 SnapMirror 支持。使用陣列型快照和複寫是保護資料的最有效方法。

## 網路組態

使用 vSphere 搭配執行 ONTAP VMware 軟體的系統時、設定網路設定很簡單、而且類似於其他網路組態。

以下是幾點需要考量的事項：

- 將儲存網路流量與其他網路區隔。使用專屬的 VLAN 或獨立的交換器來儲存、即可建立獨立的網路。如果儲存網路共用實體路徑（例如上行鏈路）、您可能需要 QoS 或額外的上行鏈路連接埠、以確保有足夠的頻寬。請勿將主機直接連線至儲存設備；請使用交換器來建立備援路徑、並讓 VMware HA 在不需介入的情況下運作。請參閱 "[直接連線網路](#)" 以取得更多資訊。
- 如果您的網路需要並支援巨型框架、尤其是使用 iSCSI 時、可以使用巨型框架。如果使用、請確定在儲存設備和 ESXi 主機之間的路徑中、所有網路裝置、VLAN 等上的設定都相同。否則、您可能看到效能或連線問題。MTU 也必須在 ESXi 虛擬交換器、VMkernel 連接埠、以及每 ONTAP 個節點的實體連接埠或介面群組上設定相同。
- NetApp 僅建議停用 ONTAP 叢集內叢集網路連接埠上的網路流量控制。對於用於資料流量的其餘網路連接埠、NetApp 並未提出其他最佳實務做法建議。您應該視需要啟用或停用它。請參閱 "[TR-4182](#)" 以取得流程控制的更多背景資訊。
- 當 ESXi 和 ONTAP VMware ESXi 儲存陣列連接至乙太網路儲存網路時、NetApp 建議將這些系統連接的乙太網路連接埠設定為快速擴充樹狀傳輸協定 (RSTP) 邊緣連接埠、或使用 Cisco PortFast 功能。NetApp 建議在使用 Cisco PortFast 功能的環境中、啟用跨距樹狀結構 PortFast 主幹功能、並在 ESXi 伺服器或 ONTAP

VMware®儲存陣列上啟用802.1Q VLAN主幹連線。

• NetApp建議下列連結集合最佳實務做法：

- 使用交換器、透過 Cisco 的 Virtual PortChannel ( VPC ) 等多機箱連結集合群組方法、在兩個獨立的交換器機箱上支援連接埠的連結集合。
- 除非您使用已設定LACP的DVSwitches 5.1或更新版本、否則請停用連接至ESXi的交換器連接埠LACP。
- 使用LACP建立鏈路集合體、以動態ONTAP 多重模式介面群組搭配IP雜湊、以利支援靜態儲存系統。
- 在ESXi上使用IP雜湊群組原則。

下表提供網路組態項目的摘要、並指出套用設定的位置。

項目	ESXi	交換器	節點	SVM
IP 位址	VMkernel	否*	否*	是的
連結集合體	虛擬交換器	是的	是的	否*
VLAN	VMkernel和VM連接埠群組	是的	是的	否*
流程控制	NIC	是的	是的	否*
跨距樹狀結構	否	是的	否	否
MTU (用於巨型框架)	虛擬交換器與VMkernel連接埠(9000)	是 (設為上限)	有 (9000)	否*
容錯移轉群組	否	否	是 (建立)	是 (選取)

- SVM lifs連接到具有VLAN、MTU及其他設定的連接埠、介面群組或VLAN介面。不過、這些設定不會在SVM層級進行管理。

這些裝置擁有自己的IP位址進行管理、但這些位址並未用於ESXi儲存網路環境。

## SAN (FC、FCoE、NVMe/FC、iSCSI)、RDM

在vSphere中、有三種使用區塊儲存LUN的方法：

- 使用VMFS資料存放區
- 使用原始裝置對應 (RDM)
- 由軟體啟動器從VM客體作業系統存取及控制的LUN

VMFS是高效能的叢集式檔案系統、可提供共用儲存資源池的資料存放區。VMFS資料存放區可設定LUN、使用NVMe / FC傳輸協定存取的FC、iSCSI、FCoE或NVMe命名空間來存取。VMFS可讓叢集中的每個ESX伺服器同時存取傳統LUN。支援的最大LUN大小通常為16TB；因此、使用四個16TB LUN (所有SAN陣列系統支援的最大VMFS LUN大小為64TB)、即可建立最大大小為64TB的VMFS 5資料存放區 (請參閱本節的第一個表格) ONTAP。由於VMware不具備小型的個別佇列深度、所以在VMware中、VMFS資料存放區的擴充程度比傳統陣列架構的擴充程度更高、而且相對簡單。ONTAP ONTAP

vSphere內建多個儲存裝置路徑的支援功能、稱為原生多重路徑 (NMP)。NMP可偵測支援儲存系統的儲存類型、並自動設定NMP堆疊以支援使用中儲存系統的功能。



NMP 和 ONTAP 都支援非對稱邏輯單元存取 (ALUA)、可協調最佳化和非最佳化的路徑。在本功能中、ALUA最佳化路徑會使用主控所存取LUN的節點上的目標連接埠、遵循直接資料路徑。ONTAP預設會在vSphere和ONTAP VMware中同時開啟ALUA。NMP 將 ONTAP 叢集識別為 ALUA、並使用 ALUA 儲存陣列類型外掛程式 (VMW\_SATP\_ALUA) 並選取循環路徑選擇外掛程式 (VMW\_PSP\_RR)。

ESXi 6最多可支援256個LUN、並可支援多達1、024條通往LUN的總路徑。ESXi不會看到任何超出這些限制的LUN或路徑。假設LUN數量上限、則路徑限制允許每個LUN有四個路徑。在更大ONTAP的實體叢集中、可以在LUN限制之前達到路徑限制。為了解決此限制、ONTAP 支援8.3版及更新版本中的選擇性LUN對應 (SLM),

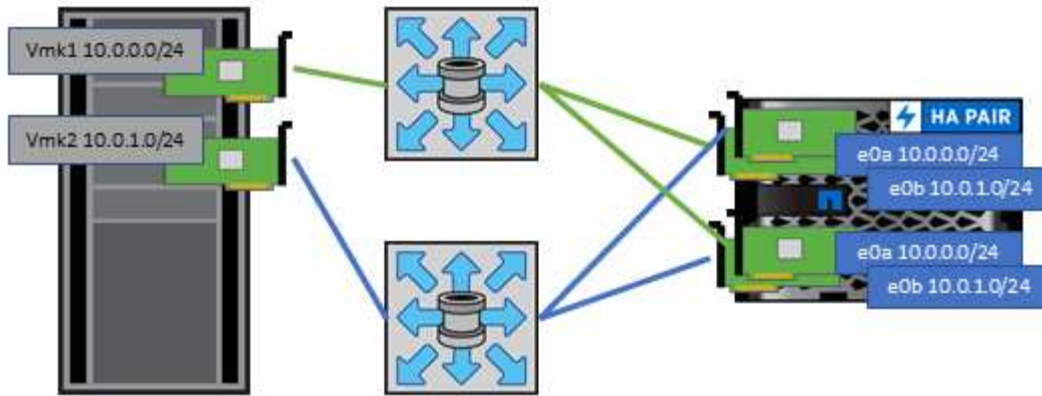
對於向指定LUN通告路徑的節點、SLM會有限制。NetApp最佳實務做法是每個SVM每個節點至少有一個LIF、並使用SLM來限制通告給裝載LUN及其HA合作夥伴之節點的路徑。雖然存在其他路徑、但預設不會通告這些路徑。您可以使用新增和移除在SLMs中的報告節點引數來修改通告的路徑。請注意、在 8.3 之前版本中建立的LUN 會通告所有路徑、而且需要修改以僅通告主控 HA 配對的路徑。如需更多關於SLM,請參閱第5.9節 "[TR-4080](#)"。先前的連接埠集方法也可用於進一步減少LUN的可用路徑。PortSets可減少igroup中的啟動器可透過哪些可見路徑來查看LUN、進而提供協助。

- 根據預設、會啟用SLM.除非您使用連接埠集、否則不需要額外的組態。
- 對於在 Data ONTAP 8.3 之前建立的 LUN、請執行手動套用 SLM lun mapping remove-reporting-nodes 用於移除 LUN 報告節點、並限制 LUN 存取 LUN 所屬節點及其 HA 合作夥伴的命令。

區塊傳輸協定 (iSCSI、FC和FCoE) 使用LUN ID和序號以及唯一名稱來存取LUN。FC和FCoE使用全球名稱 (WWNN和WWPN)、iSCSI則使用iSCSI合格名稱 (IQN)。儲存設備內部的LUN路徑對區塊傳輸協定毫無意義、而且不會出現在傳輸協定的任何位置。因此、只包含LUN的磁碟區根本不需要內部掛載、而包含資料存放區所用LUN的磁碟區則不需要使用交會路徑。NVMe子系統ONTAP 的運作方式類似。

其他應考慮的最佳實務做法：

- 請確定ONTAP 已為叢集中每個節點上的每個SVM建立邏輯介面 (LIF)、以達到最大可用度和行動性。最佳實務做法是每個節點使用兩個實體連接埠和LIF、每個光纖使用一個連接埠。ONTAPALUA可用來剖析路徑、識別作用中最佳化 (直接) 路徑、以及作用中未最佳化路徑。ALUA用於FC、FCoE和iSCSI。
- 對於iSCSI網路、當存在多個虛擬交換器時、請在不同的網路子網路上使用多個VMkernel網路介面搭配NIC 群組。您也可以使用多個實體NIC來連接至多個實體交換器、以提供HA並提高處理量。下圖提供多重路徑連線的範例。在 ONTAP 中、請使用單一模式介面群組、其中包含多個連結至不同交換器或具有多重模式介面群組的 LACP、以獲得高可用度和連結集合效益。
- 如果 ESXi 中使用挑戰握手驗證傳輸協定 (CHAP) 進行目標驗證、則也必須使用 CLI 在 ONTAP 中進行設定 (vserver iscsi security create) 或使用 System Manager (在 Storage (儲存) > SVM (SVM) > SVM Settings (SVM 設定) > Protocols (傳輸協定) > iSCSI (iSCSI) 下編輯啟動器安全性)。
- 使用VMware vSphere的VMware vCenter工具來建立及管理LUN和群組。ONTAP外掛程式會自動決定伺服器的WWPN、並建立適當的igroup。它也會根據最佳實務做法來設定LUN、並將其對應至正確的igroup。
- 請謹慎使用 RDM、因為它們可能較難管理、而且也會使用路徑、而路徑的限制如前所述。支援這兩種LUN ONTAP "[實體與虛擬相容模式](#)" RDM。
- 如需更多關於將NVMe/FC搭配vSphere 7.0使用的資訊、請參閱此 "[NVMe / FC主機組態指南ONTAP](#)" 和 "[TR-4684](#)"。下圖說明從 vSphere 主機到 ONTAP LUN 的多重路徑連線能力。



## NFS

vSphere可讓客戶使用企業級NFS陣列、同時存取ESXi叢集中所有節點的資料存放區。如資料存放區一節所述、在使用NFS搭配vSphere時、會有一些易於使用和儲存效率可見度的優點。

搭配ONTAP vSphere使用VMware NFS時、建議採用下列最佳實務做法：

- 在叢集中的每個節點上、為每個SVM使用單一邏輯介面（LIF）ONTAP。不再需要過去針對每個資料存放區的LIF建議。雖然直接存取（同一個節點上的LIF和資料存放區）是最佳選擇、但別擔心間接存取、因為效能影響通常很小（微秒）。
- 目前支援的所有VMware vSphere版本均可使用NFS v3和v4.1。正式支援nconnect已新增至vSphere 8.0更新2 for NFS v3。對於NFS v4.1、vSphere持續支援工作階段主幹、Kerberos驗證及完整性Kerberos驗證。請務必注意、工作階段主幹需要ONTAP 9.14.1或更新版本。您可以深入瞭解nconnect功能、以及它如何改善的效能"[NFSv3 nConnect 功能搭配 NetApp 和 VMware](#)"。

值得注意的是、NFSv3和NFSv4.1使用不同的鎖定機制。NFSv3使用用戶端端鎖定、而NFSv4.1則使用伺服器端鎖定。雖然ONTAP磁碟區可以透過兩種傳輸協定匯出、但ESXi只能透過一種傳輸協定掛載資料存放區。不過、這並不表示其他ESXi主機無法透過不同版本掛載相同的資料存放區。為了避免任何問題、請務必指定要在掛載時使用的通訊協定版本、確保所有主機都使用相同版本、因此使用相同的鎖定樣式。避免在主機之間混合使用NFS版本是非常重要的。如有可能、請使用主機設定檔來檢查法規遵循狀況。

由於**NFSv3**和**NFSv4.1**之間沒有自動資料存放區轉換、因此請建立新的**NFSv4.1**資料存放區、並使用**Storage VMotion**將VM移轉至新的資料存放區。

請參閱中的NFS v4.1互通性表附註"[NetApp互通性對照表工具](#)"支援所需的特定ESXi修補程式層級。

\* NFS匯出原則用於控制vSphere主機的存取。您可以將一個原則與多個磁碟區（資料存放區）搭配使用。使用NFSv3時、ESXi會使用sys（UNIX）安全樣式、並需要root掛載選項來執行VM。在現階段、此選項稱為超級使用者、使用超級使用者選項時、不需要指定匿名使用者ID。ONTAP請注意、匯出具有不同值的原則規則-anon和-allow-suid可能會導致SVM發現ONTAP問題、因為使用這些功能。以下是範例原則：

存取傳輸協定：**nfs3**。

用戶端比對規格：192.168.42.21

RO存取規則：系統

RW存取規則：系統

匿名UID

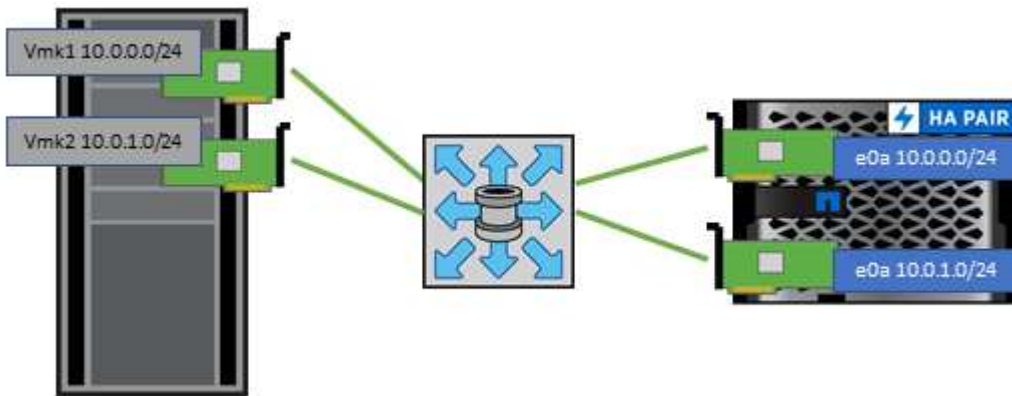
超級用戶：sys

\* 如果使用NetApp NFS外掛程式for VMware VAAI、則通訊協定應設為nfs匯出原則規則建立或修改時。需要NFSv4傳輸協定、VAAI複本卸載才能運作、並將傳輸協定指定為nfs自動同時包含NFSv3和NFSv4版本。

\* NFS資料存放區磁碟區是從SVM的根磁碟區連接而成、因此ESXi也必須擁有根磁碟區的存取權、才能瀏覽及裝載資料存放區磁碟區。根Volume和資料存放區磁碟區交會嵌套的任何其他磁碟區的匯出原則、都必須包含ESXi伺服器授予其唯讀存取權的規則或規則。以下是根Volume的範例原則、也使用VAAI外掛程式：

存取傳輸協定：**NFS**（包括 **nfs3** 和 **nfs4**）  
 用戶端比對規格：192.168.42.21  
**RO** 存取規則：系統  
**RW** 存取規則：永不（root Volume 的最佳安全性）  
**匿名 UID**  
 超級用戶：sys（VAAI 的 root Volume 也需要）

\* 使用適用於 VMware vSphere 的 ONTAP 工具（最重要的最佳實務做法）：  
 使用適用於 **VMware vSphere** 的 **ONTAP** 工具來配置資料存放區、因為它能自動簡化匯出原則的管理。  
 使用外掛程式為 VMware 叢集建立資料存放區時、請選取叢集、而非單一 ESX 伺服器。此選項會觸發 IT 自動將資料存放區掛載至叢集中的所有主機。  
 使用外掛裝載功能將現有的資料存放區套用至新伺服器。  
 如果不使用 ONTAP 工具來執行 VMware vSphere、請針對所有伺服器或需要額外存取控制的每個伺服器叢集、使用單一匯出原則。  
 \* 雖然 ONTAP 提供彈性的 Volume 命名空間結構、可利用交叉路口在樹狀結構中排列磁碟區、但這種方法對 vSphere 沒有價值。無論儲存設備的命名空間階層為何、它都會在資料存放區根目錄中為每個 VM 建立一個目錄。因此、最佳實務做法是將 vSphere 磁碟區的交會路徑掛載到 SVM 的根磁碟區、這就是 ONTAP VMware vSphere 的 VMware vSphere 功能如何配置資料存放區。沒有巢狀結點路徑也表示除了根磁碟區之外、沒有任何磁碟區相依於任何磁碟區、即使是刻意將磁碟區離線或銷毀、也不會影響其他磁碟區的路徑。  
 \* 對於 NFS 資料存放區上的 NTFS 分割區、4K 區塊大小是很好的。下圖說明從 vSphere 主機連線至 ONTAP VMware NFS 資料存放區的能力。



下表列出 NFS 版本及支援的功能。

vSphere 功能	NFSv3.	NFSv4.1
vMotion 與 Storage vMotion	是的	是的
高可用性	是的	是的
容錯能力	是的	是的
DRS	是的	是的
主機設定檔	是的	是的
儲存 DRS	是的	否
儲存 I/O 控制	是的	否
SRM	是的	否
虛擬磁碟區	是的	否
硬體加速 (VAAI)	是的	是的

vSphere功能	NFSv3.	NFSv4.1
Kerberos驗證	否	是 (vSphere 6.5及更新版本增強支援AES、krb5i)
多重路徑支援	否	有 (ONTAP 9.14.1)

## 直接連線網路

儲存管理員有時偏好從組態中移除網路交換器、以簡化其基礎架構。在某些情況下可能會支援這項功能。

## iSCSI 和 NVMe / TCP

使用 iSCSI 或 NVMe / TCP 的主機可以直接連線至儲存系統、並正常運作。原因是路徑。直接連線至兩個不同的儲存控制器、可產生兩個不同的資料流路徑。遺失路徑、連接埠或控制器並不會妨礙其他路徑的使用。

## NFS

可以使用直接連線的 NFS 儲存設備、但有很大的限制：如果沒有大量的指令碼工作、容錯移轉將無法運作、這是客戶的責任。

直接連線的 NFS 儲存設備會造成不中斷的容錯移轉複雜化、這是因為本機作業系統上會發生路由。例如、假設主機的 IP 位址為 192.168.1.1/24、並直接連線至 IP 位址為 192.168.1.50/24 的 ONTAP 控制器。在容錯移轉期間、該位址 192.168.1.50 可以容錯移轉至其他控制器、而且主機可以使用該位址、但主機如何偵測其存在？原來的 192.168.1.1 位址仍然存在於不再連線至作業系統的主機 NIC 上。目的地為 192.168.1.5 的流量將繼續傳送至無法運作的網路連接埠。

第二個 OS NIC 可設定為 19 可以與故障的 over 192.168.1.50 位址進行通訊、但本機路由表預設會使用一個 \* 且只有一個 \* 位址來與 192.168.1.0/24 子網路通訊。系統管理員可以建立指令碼架構、以偵測失敗的網路連線、並變更本機路由表或使介面正常運作。具體程序取決於所使用的作業系統。

實際上、NetApp 客戶確實有直接連線的 NFS、但通常僅適用於容錯移轉期間 IO 暫停的工作負載。使用硬掛載時、在這類暫停期間不應有任何 IO 錯誤。IO 應該會暫停運作、直到服務還原為止、無論是透過容錯回復或手動介入、在主機上的 NIC 之間移動 IP 位址。

## FC Direct Connect

無法使用 FC 傳輸協定將主機直接連接至 ONTAP 儲存系統。原因是使用 NPIV。用於識別 FC 網路的 ONTAP FC 連接埠的 WWN 使用稱為 NPIV 的虛擬化類型。任何連接至 ONTAP 系統的裝置都必須能夠辨識 NPIV WWN。目前沒有任何 HBA 廠商提供可安裝在能夠支援 NPIV 目標的主機上的 HBA。

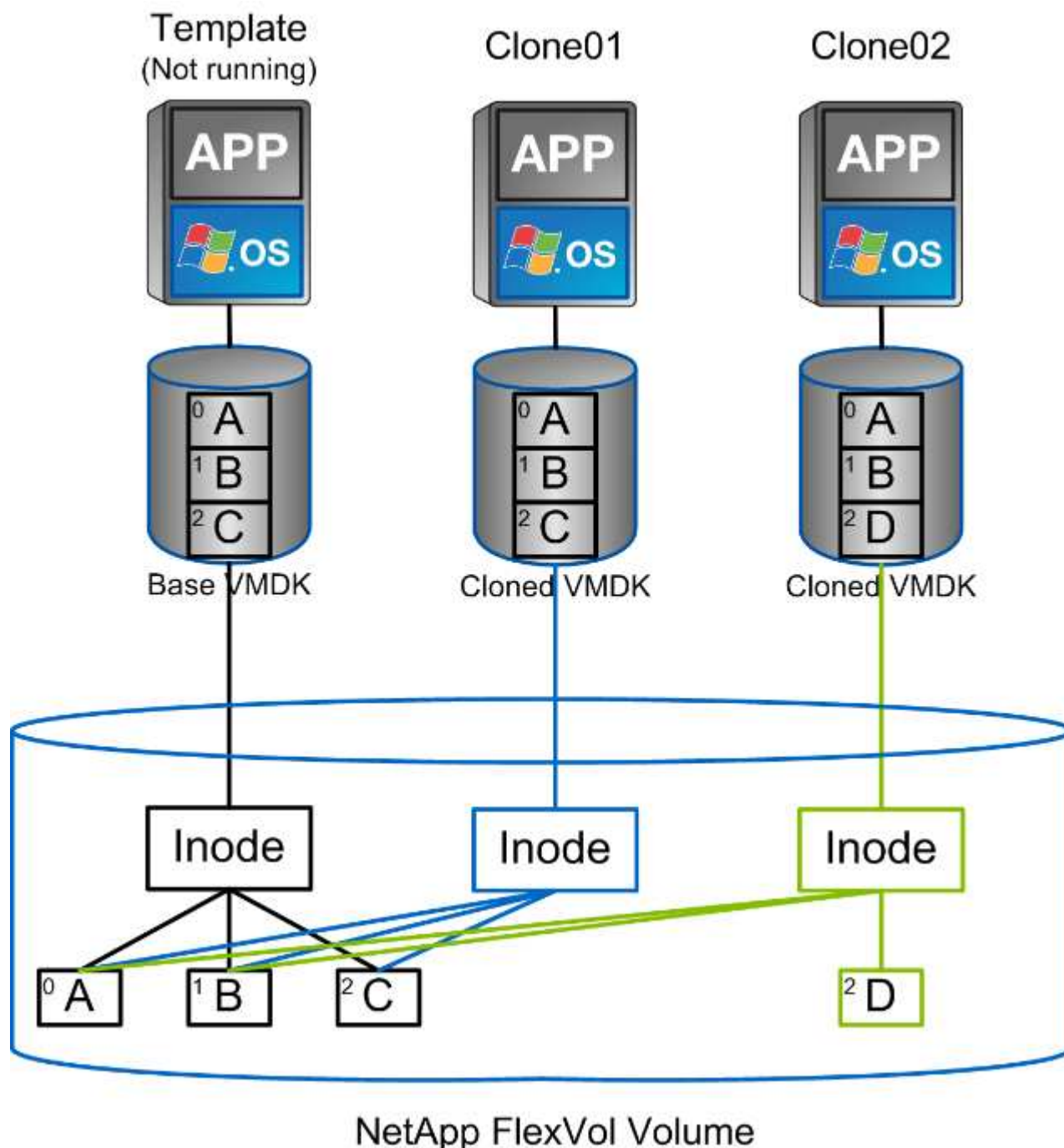
## 虛擬機器與資料存放區複製

複製儲存物件可讓您快速建立複本以供進一步使用、例如配置額外的 VM、備份/還原作業等。

在 vSphere 中、您可以複製 VM、虛擬磁碟、vVol 或資料存放區。複製之後、通常可透過自動化程序、進一步自訂物件。vSphere 同時支援完整複本複本複本、以及連結的複本、可在其中分別追蹤原始物件的變更。

連結的複本非常適合用來節省空間、但會增加 vSphere 處理 VM 的 I/O 量、進而影響該 VM 的效能、甚至可能影響整個主機的效能。這就是 NetApp 客戶經常使用儲存系統型複本來充分發揮兩者優勢的原因：有效使用儲存設備並提高效率。

下圖說明ONTAP 了還原複製。



複製可ONTAP 透過多種機制卸載至執行支援軟體的系統、通常是VM、vVol或資料存放區層級。其中包括：

- VVols使用NetApp vSphere API for Storage Aware (VASA) Provider。ONTAP 複本可用於支援 vCenter 管理的 vVol 快照、這些快照的空間效率極低、而且 I/O 效果極低、可用來建立和刪除這些快照。也可以使用vCenter複製VM、這些VM也會卸載到ONTAP VMware、無論是在單一資料存放區/ Volume內、或是在資料存放區/磁碟區之間。
- 使用vSphere API進行vSphere複製與移轉-陣列整合 (VAAI)。VM複製作業可在ONTAP SAN和NAS環境中卸載至VMware (NetApp提供ESXi外掛程式來啟用VAAI for NFS)。vSphere僅卸載NAS資料存放區中的冷 (關機) VM作業、而熱VM (複製與儲存vMotion) 上的作業也會卸載用於SAN。根據來源、目的地及安裝的產品授權、使用最有效率的方法。ONTAPVMware Horizon View也會使用此功能。

- SRA（與VMware Site Recovery Manager搭配使用）。在此、複本可用來在不中斷營運的情況下測試災難恢復複本的還原。
- 使用SnapCenter NetApp工具（如VMware）進行備份與恢復。VM複製可用來驗證備份作業、以及掛載VM備份、以便複製個別檔案。

VMware、NetApp及協力廠商工具可叫用不需載入的複製。ONTAP卸載至ONTAP 不完整的複本有多項優點。在大多數情況下、它們都能節省空間、只需要儲存設備來變更物件；讀取和寫入時不會產生額外的效能影響、在某些情況下、透過在高速快取中共用區塊來改善效能。此外、也會從ESXi伺服器卸載CPU週期和網路I/O。使用FlexClone授權、在傳統資料存放區內使用FlexVol 實體磁碟區的複本卸載作業既快速又有效率、但FlexVol 在各個實體磁碟區之間的複本可能會較慢。如果您將VM範本保留為複製來源、請考慮將其放在資料存放區磁碟區內（使用資料夾或內容程式庫來組織它們）、以獲得快速且具空間效益的複本。

您也可以直接複製ONTAP 實體內部的磁碟區或LUN、以複製資料存放區。有了NFS資料存放區、FlexClone技術就能複製整個Volume、而且可從ONTAP VMware匯出該實體複本、並由ESXi作為另一個資料存放區來掛載。對於VMFS資料存放區、ONTAP VMware可以複製一個或整個Volume內的LUN、包括其中的一個或多個LUN。包含VMFS的LUN必須對應至ESXi啟動器群組（igroup）、然後由ESXi重新簽名、才能掛載並作為一般資料存放區使用。在某些臨時使用案例中、可以掛載複製的VMFS、而無需重新簽名。複製資料存放區之後、就可以登錄、重新設定及自訂其中的VM、就像是個別複製的VM一樣。

在某些情況下、您可以使用額外的授權功能來強化複製功能、例如SnapRestore 針對備份或FlexClone的功能。這些授權通常包含在授權套裝組合中、不需額外付費。vVol 複製作業需要 FlexClone 授權、也必須支援 vVol 的託管快照（從 Hypervisor 卸載至 ONTAP）。FlexClone授權也能在資料存放區/磁碟區內使用時、改善特定的VAAI型複本（建立即時、節省空間的複本、而非區塊複本）。SRA也會在測試災難恢復複本的恢復時使用此複本、SnapCenter 而使用此複本來執行複製作業、並瀏覽備份複本來還原個別檔案。

## 資料保護

備份虛擬機器並快速恢復這些虛擬機器、是ONTAP vSphere的絕佳優勢之

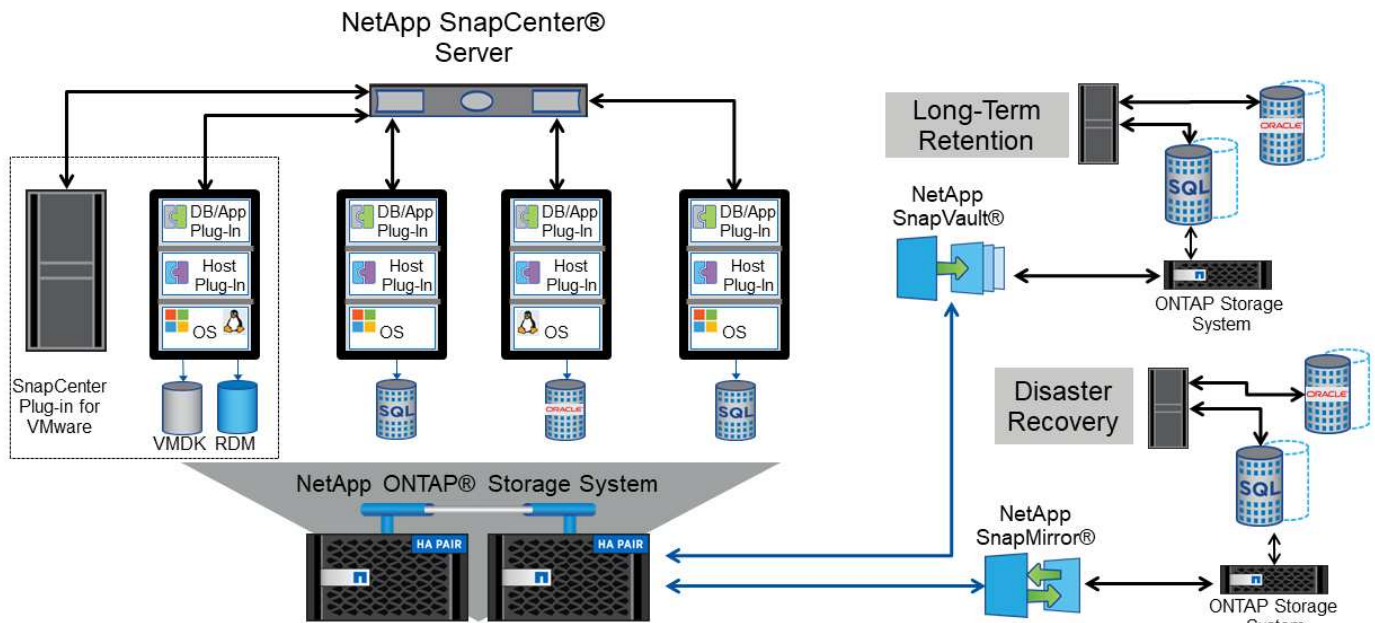
一、SnapCenter 而使用VMware vSphere的VMware vCenter外掛程式、可輕鬆管理vCenter內部的這項功能。

使用快照快速複製 VM 或資料存放區、而不會影響效能、然後使用 SnapMirror 將其傳送至次要系統、以提供長期的異地資料保護。這種方法只儲存變更的資訊、可將儲存空間和網路頻寬減至最低。

利用此功能、您可以建立可套用至多個工作的備份原則。SnapCenter這些原則可定義排程、保留、複寫及其他功能。它們持續允許選用 VM 一致的快照、利用 Hypervisor 的能力來在拍攝 VMware 快照之前、先關閉 I/O。不過、由於VMware快照的效能影響、除非您需要靜止客體檔案系統、否則一般不建議使用這些快照。請改用快照來提供一般保護、並使用 SnapCenter 外掛程式等應用程式工具來保護交易資料、例如 SQL Server 或 Oracle。這些快照與 VMware（一致性）快照不同、適合長期保護。VMware快照僅供參考 ["建議"](#) 因效能和其他影響而短期使用。

這些外掛程式提供延伸功能、可在實體和虛擬環境中保護資料庫。有了vSphere、您可以使用它們來保護SQL Server或Oracle資料庫、其中資料儲存在RDM LUN、iSCSI LUN直接連線至來賓作業系統、或VMDK或NFS資料存放區上的VMDK檔案。外掛程式可指定不同類型的資料庫備份、支援線上或離線備份、以及保護資料庫檔案和記錄檔。除了備份與還原之外、外掛程式也支援複製資料庫以供開發或測試之用。

下圖說明SnapCenter 了一套功能性部署的範例。



如需增強的災難恢復功能、請考慮搭配ONTAP VMware Site Recovery Manager使用NetApp SRA for VMware。除了支援將資料存放區複製寫至DR站台之外、它也能透過複製複製的資料存放區、在DR環境中進行不中斷營運的測試。SRA內建自動化功能、可在故障解決後、輕鬆從災難中恢復並重新保護正式作業。

最後、若要獲得最高層級的資料保護、請考慮使用NetApp MetroCluster VMware vSphere Metro Storage Cluster (VMSC) 組態。VMSC是VMware認證的解決方案、結合了同步複製與陣列式叢集、提供高可用性叢集的同樣優勢、但分散在不同站台、以防止站台災難發生。NetApp MetroCluster 解決方案提供具成本效益的組態、可讓您以透明的方式從任何單一儲存元件故障中恢復同步複製、並在發生站台災難時提供單一命令恢復功能。VMSC的詳細說明請參閱 "TR-4128"。

## 服務品質 (QoS)

執行ONTAP 支援功能的系統可使用ONTAP 「支援服務品質」功能、限制檔案、LUN、磁碟區或整個SVM等不同儲存物件的每秒處理量 (以Mbps和/或I/O (IOPS) 為單位)。

處理量限制可在部署前控制未知工作負載或測試工作負載、以確保不會影響其他工作負載。它們也可用於在識別出高性能工作負載之後加以限制。也支援以IOPS為基礎的最低服務層級、以提供一致的效能、適用於ONTAP VMware的SAN物件、以及ONTAP 支援VMware的NAS物件。

有了NFS資料存放區、QoS原則可套用至FlexVol 整個VMware磁碟區或其中的個別VMDK檔案。使用ONTAP 使用VMware LUN的VMFS資料存放區時、QoS原則可套用至FlexVol 包含LUN或個別LUN的VMware磁碟區、但不能套用個別VMDK檔案、因為ONTAP VMware對VMFS檔案系統沒有任何認知。使用vVols時、可使用儲存功能設定檔和VM儲存原則、在個別VM上設定最低和/或最高QoS。

物件的QoS最大處理量限制可設定為Mbps和/或IOPS。如果兩者皆使用、ONTAP 則由支援執行第一個上限。工作負載可以包含多個物件、QoS原則也可以套用至一或多個工作負載。當原則套用至多個工作負載時、工作負載會共用原則的總限制。不支援巢狀物件 (例如、磁碟區內的檔案無法各自擁有自己的原則)。QoS最低值只能以IOPS設定。

下列工具目前可用於管理ONTAP 不實的QoS原則、並將其套用至物件：

- CLI ONTAP

- 系統管理程式ONTAP
- OnCommand Workflow Automation
- Active IQ Unified Manager
- NetApp PowerShell Toolkit for ONTAP
- VMware vSphere VASA Provider適用的工具ONTAP

若要將QoS原則指派給NFS上的VMDK、請注意下列準則：

- 原則必須套用至 `vmname-flat.vmdk` 其中包含實際的虛擬磁碟映像、而非 `vmname.vmdk`（虛擬磁碟描述元檔案）或 `vmname.vmx`（VM 描述元檔案）。
- 請勿將原則套用至其他 VM 檔案、例如虛擬交換檔案 (`vmname.vswp`)。
- 使用 vSphere Web 用戶端尋找檔案路徑（資料存放區 > 檔案）時、請注意它會結合的資訊 - `flat.vmdk` 和 `.vmdk` 只要顯示一個檔案、其中包含的名稱即可 `.vmdk` 但是的大小 - `flat.vmdk`。新增 `-flat` 輸入檔案名稱以取得正確路徑。

若要將QoS原則指派給LUN、包括VMFS和RDM、ONTAP 顯示為Vserver的SVM、LUN路徑和序號、可從ONTAP VMware vSphere的「VMware vSphere的VMware vSphere」（VMware vSphere）「VMware vCenter工具」首頁上的「儲存系統」功能表取得。選取儲存系統（SVM）、然後選取相關物件 > SAN。使用ONTAP 其中一項功能來指定QoS時、請使用此方法。

利用ONTAP VMware vSphere或Virtual Storage Console 7.1及更新版本的VMware vSphere或Virtual Storage Console 7.1工具、可輕鬆將最大和最小QoS指派給VVol型VM。為 vVol 容器建立儲存功能設定檔時、請在效能功能下指定最大和 / 或最小 IOPS 值、然後使用 VM 的儲存原則參考此 SCP。建立VM或將原則套用至現有VM時、請使用此原則。

使用VMware vSphere 9.8及更新版本的VMware vSphere 9.8版的VMware VMware vCenter資料存放區提供增強的QoS功能。FlexGroup ONTAP您可以輕鬆地在資料存放區或特定VM的所有VM上設定QoS。如FlexGroup 需詳細資訊、請參閱本報告的「參考資料」一節。

## QoS和VMware SIOC ONTAP

VMware vSphere儲存I/O控制（SIOC）是相輔相成的技術、vSphere和儲存管理員可以搭配使用、來管理執行VMware軟體之系統上所託管vSphere VM的效能。ONTAP ONTAP每個工具都有自己的優點、如下表所示。由於VMware vCenter和ONTAP VMware vCenter的範圍不同、有些物件可由一個系統來查看和管理、而非由另一個系統來管理。

屬性	QoS ONTAP	VMware SIOC
當作用中時	原則永遠處於作用中狀態	存在爭用時作用中（超過臨界值的資料存放區延遲）
單位類型	IOPS、Mbps	IOPS、共享
vCenter或應用程式範圍	多個vCenter環境、其他Hypervisor和應用程式	單一vCenter伺服器
在VM上設定QoS？	僅NFS上的VMDK	NFS或VMFS上的VMDK
設定LUN上的QoS（RDM）？	是的	否
在LUN（VMFS）上設定QoS？	是的	否



屬性	QoS ONTAP	VMware SIOC
在Volume (NFS資料存放區) 上設定QoS?	是的	否
在SVM (租戶) 上設定QoS?	是的	否
以原則為基礎的方法?	是; 可由原則中的所有工作負載共用、或完全套用至原則中的每個工作負載。	是、使用vSphere 6.5及更新版本。
需要授權	隨附ONTAP 於此功能	Enterprise Plus

## VMware Storage Distributed Resource Scheduler

VMware儲存分散式資源排程器 (SDR) 是vSphere功能、可根據目前的I/O延遲和空間使用量、將VM放置在儲存設備上。接著、它會在資料存放區叢集中的資料存放區之間 (也稱為Pod)、在不中斷營運的情況下移動VM或VMDK、並選取將VM或VMDK置於資料存放區叢集中的最佳資料存放區。資料存放區叢集是類似資料存放區的集合、從 vSphere 管理員的觀點來看、這些資料存放區會彙總成單一使用量單位。

搭配 ONTAP 工具使用適用於 VMware vSphere 的 SDR 時、您必須先使用外掛程式建立資料存放區、使用 vCenter 建立資料存放區叢集、然後將資料存放區新增至該叢集。建立資料存放區叢集之後、可直接從「詳細資料」頁面上的資源配置精靈、將其他資料存放區新增至資料存放區叢集。

SDR的ONTAP 其他最佳實務做法包括：

- 叢集中的所有資料存放區都應該使用相同類型的儲存設備 (例如SAS、SATA或SSD)、無論是所有VMFS或NFS資料存放區、都具有相同的複寫和保護設定。
- 請考慮在預設 (手動) 模式下使用SDR。此方法可讓您檢閱建議、並決定是否要套用建議。請注意VMDK移轉的下列影響：
  - 當SDR在資料存放區之間移動VMDK時、ONTAP 任何從還原複製或重複資料刪除所節省的空间都會遺失。您可以重新執行重複資料刪除、以重新獲得這些節約效益。
  - 在 SDR 移動 VMDK 之後、NetApp 建議在來源資料存放區重新建立快照、因為其他情況下空間會被移動的 VM 鎖定。
  - 在同一個集合體上的資料存放區之間移動VMDK並沒有什麼好處、而且SDR無法看到可能共用該集合體的其他工作負載。

## 儲存原則型管理和 vVols

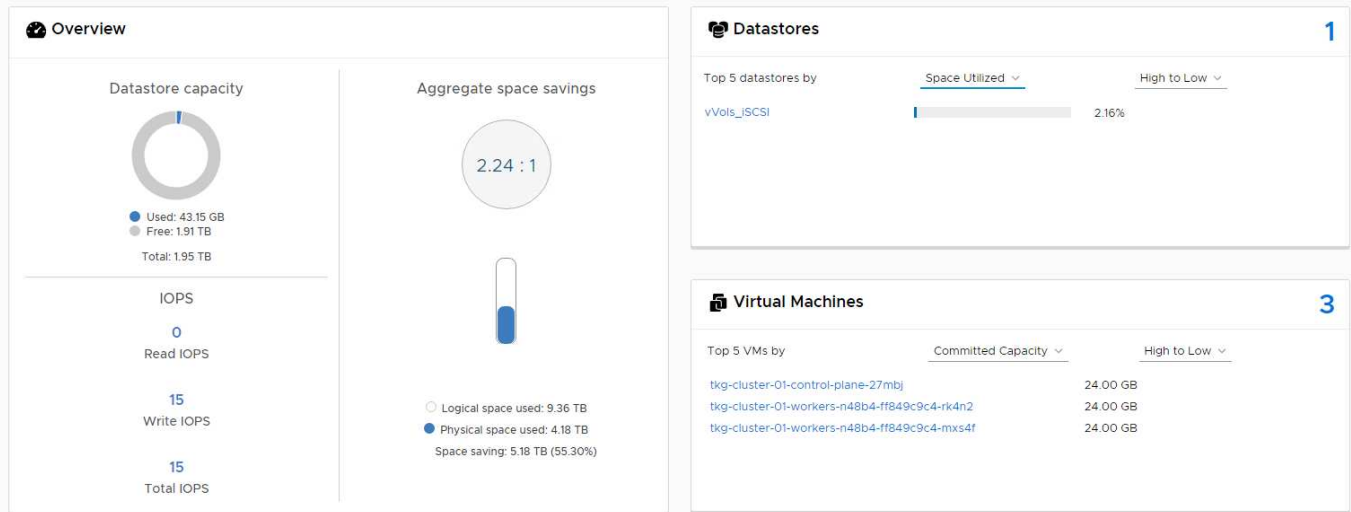
VMware vSphere API for Storage感知 (VASA) 可讓儲存管理員輕鬆設定具有明確定義功能的資料存放區、並讓VM管理員在需要時使用這些功能來配置VM、而不需要彼此互動。請看一下這種方法、瞭解它如何簡化虛擬化儲存作業、並避免許多瑣碎的工作。

在VASA之前、VM管理員可以定義VM儲存原則、但他們必須與儲存管理員合作、以識別適當的資料存放區、通常是使用文件或命名慣例。有了VASA、儲存管理員可以定義一系列的儲存功能、包括效能、分層、加密及複寫。一組磁碟區或一組磁碟區的功能稱為儲存功能設定檔 (scp)。

SCP 支援虛擬機器資料 VVols 的最低和 / 或最高 QoS。只AFF 有在不支援的系統上才支援最低QoS。VMware vSphere的VMware vSphere工具包含儀表板、可顯示VM精細的效能、以及在VMware系統上用於vVols的邏輯容量。ONTAP ONTAP

下圖說明ONTAP VMware vSphere 9.8 vVols儀表板的各項功能。

The dashboard displays IOPS, latency, throughput, and logical space values obtained from ONTAP.



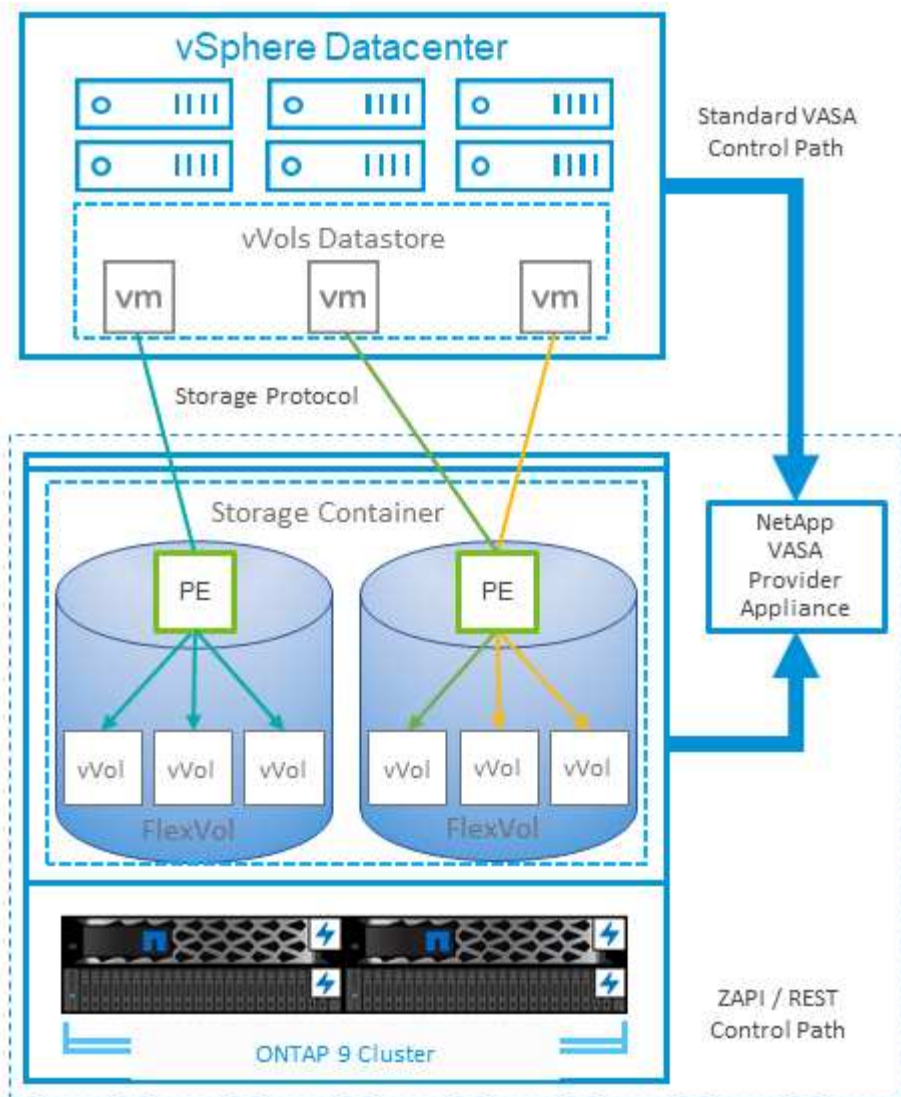
定義儲存功能設定檔之後、就可以使用識別其需求的儲存原則來配置VM。VM儲存原則與資料存放區儲存功能設定檔之間的對應、可讓vCenter顯示相容資料存放區清單以供選擇。這種方法稱為儲存原則型管理。

VASA提供查詢儲存設備的技術、並將一組儲存功能傳回vCenter。VASA廠商供應商會提供儲存系統API與架構之間的轉譯、以及vCenter所瞭解的VMware API。NetApp的 VASA Provider for ONTAP 是 ONTAP 工具的一部分、適用於 VMware vSphere 應用裝置 VM、vCenter 外掛程式則提供介面、可配置及管理 vVol 資料存放區、並可定義儲存功能設定檔 (CDP)。

支援VMFS和NFS vVol資料存放區。ONTAP將vVols與SAN資料存放區搭配使用、可帶來NFS的部分效益、例如VM層級的精細度。以下是一些最佳實務做法、您可以在中找到更多資訊 "[TR-4400](#)"：

- VVol資料存放區可由FlexVol 多個叢集節點上的多個支援功能區所組成。最簡單的方法是單一資料存放區、即使磁碟區具有不同的功能也一樣。SPBM可確保VM使用相容的Volume。然而、這些磁碟區必須全部屬於ONTAP 單一的一套功能、並使用單一傳輸協定來存取。每個節點的每個傳輸協定只需一個LIF就足夠了。避免在ONTAP 單一vVol資料存放區中使用多個版本的支援、因為儲存功能可能因版本而異。
- 使用ONTAP VMware vSphere外掛程式的VMware vCenter工具來建立及管理VVol資料存放區。除了管理資料存放區及其設定檔之外、它還會自動建立傳輸協定端點、以便在需要時存取vVols。如果使用LUN、請注意LUN PE是使用LUN ID 300以上的LUN來對應。確認 ESXi 主機進階系統設定 `Disk.MaxLUN` 允許大於 300 的 LUN ID 號碼 (預設值為 1,024)。若要執行此步驟、請在 vCenter 中選取 ESXi 主機、然後選取「設定」索引標籤、再選取「尋找」 `Disk.MaxLUN` 在進階系統設定清單中。
- 請勿安裝或移轉VASA Provider、vCenter Server (應用裝置或Windows) 或ONTAP VMware vSphere的各種支援工具到vVols資料存放區、因為這些工具彼此相依、因此在停電或其他資料中心中斷時、您無法管理這些工具。
- 定期備份VASA Provider VM。至少要為包含 VASA Provider 的傳統資料存放區建立每小時快照。如需保護及恢復VASA Provider的詳細資訊、請參閱此 "[知識庫文章](#)"。

下圖顯示vVols元件。



## 雲端移轉與備份

另一ONTAP 項優勢是廣泛支援混合雲、將內部部署私有雲中的系統與公有雲功能合併在一起。以下是一些可搭配vSphere使用的NetApp雲端解決方案：

- \* 雲端 Volume 。 \* NetApp Cloud Volumes Service for Amazon Web Services 或 Google Cloud Platform 和 Azure NetApp Files for anf 可在領先業界的公有雲環境中、提供高效能、多重傳輸協定的託管儲存服務。VMware Cloud VM來賓可以直接使用。
- \* 《NetApp》資料管理軟體可在您選擇的雲端上、為您的資料提供控制、保護、靈活度及效率。Cloud Volumes ONTAP Cloud Volumes ONTAPCloud Volumes ONTAP 是以 ONTAP 儲存設備為基礎的雲端原生資料管理軟體。搭配Cloud Manager一起使用、即可部署Cloud Volumes ONTAP 及管理包含ONTAP 內部部署的各種系統的不二執行個體。利用先進的 NAS 和 iSCSI SAN 功能、搭配整合式資料管理、包括快照和 SnapMirror 複寫。
- \* Cloud Services 。 \*使用Cloud Backup Service NetApp或SnapMirror Cloud、利用公有雲儲存設備保護內部部署系統的資料。可協助您在NAS、物件儲存區和物件儲存區之間移轉及保持資料同步。Cloud Sync Cloud Volumes Service
- \* FabricPool 《》 **FabricPool** \* 《》 《》 提供快速且簡單的ONTAP 資料分層功能。冷區塊可移轉至公有雲或私有 StorageGRID 物件存放區中的物件存放區、並在再次存取 ONTAP 資料時自動重新叫用。或是將物件層用作SnapVault 已由效益管理的資料的第三層保護。您可以使用這種方法 "[儲存更多 VM 快照](#)" 在一

線ONTAP 和/或二線的不二元儲存系統上。

- 《》。\*使用NetApp軟體定義的儲存設備、將您的私有雲端延伸至遠端設施和辦公室、您可以使用《》來支援區塊和檔案服務、以及您在企業資料中心擁有的相同vSphere資料管理功能。ONTAP Select ONTAP Select

在設計VM型應用程式時、請考慮未來的雲端行動力。例如、應用程式和資料檔案不會放在一起、而是使用個別LUN或NFS匯出來匯出資料。這可讓您將VM和資料分別移轉至雲端服務。

## vSphere資料加密

如今、透過加密保護閒置資料的需求與日俱增。雖然最初的重點是財務和醫療資訊、但無論資訊儲存在檔案、資料庫或其他資料類型中、都越來越有興趣保護所有資訊。

執行ONTAP 此軟體的系統可透過閒置加密、輕鬆保護任何資料。NetApp儲存加密 (NSE) 使用自我加密的磁碟機ONTAP 搭配使用、以保護SAN和NAS資料。NetApp也提供NetApp Volume Encryption和NetApp Aggregate Encryption、這是一種簡單、以軟體為基礎的方法、可加密任何磁碟機上的磁碟區。此軟體加密不需要特殊的磁碟機或外部金鑰管理員、ONTAP 客戶可免費使用。您可以在不中斷用戶端或應用程式的情況下升級及開始使用、而且它們已通過FIPS 140-2第1級標準驗證、包括內建金鑰管理程式。

有幾種方法可以保護在VMware vSphere上執行的虛擬化應用程式資料。其中一種方法是在客體作業系統層級使用VM內部的軟體來保護資料。vSphere 6.5等較新的Hypervisor現在也支援VM層級的加密、這是另一種替代方案。不過、NetApp軟體加密既簡單又簡單、而且具有下列優點：

- \*對虛擬伺服器CPU沒有影響。\*某些虛擬伺服器環境需要其應用程式的每個可用CPU週期、但測試顯示、Hypervisor層級加密需要高達5倍的CPU資源。即使加密軟體支援 Intel 的 AES-NI 指令集來卸載加密工作負載 (如同 NetApp 軟體加密一樣)、由於新的 CPU 與舊版伺服器不相容、因此這種方法可能不可行。
- 隨附機上金鑰管理程式。NetApp軟體加密包含內建金鑰管理程式、不需額外付費、因此無需購買和使用複雜的高可用度金鑰管理伺服器、即可輕鬆開始使用。
- \*對儲存效率沒有影響。\*目前廣泛使用重複資料刪除與壓縮等儲存效率技術、是以具成本效益的方式使用Flash磁碟媒體的關鍵。不過、加密資料通常無法進行重複資料刪除或壓縮。NetApp硬體與儲存加密的運作層級較低、可充分運用領先業界的NetApp儲存效率功能、不像其他方法。
- \*輕鬆進行資料存放區精細加密。\*有了NetApp Volume Encryption、每個磁碟區都能獲得自己的AES 256位元金鑰。如果您需要變更、只要使用一個命令即可。如果您有多個租戶、或需要證明不同部門或應用程式的獨立加密、這種方法非常適合。這種加密是在資料存放區層級進行管理、比管理個別VM容易得多。

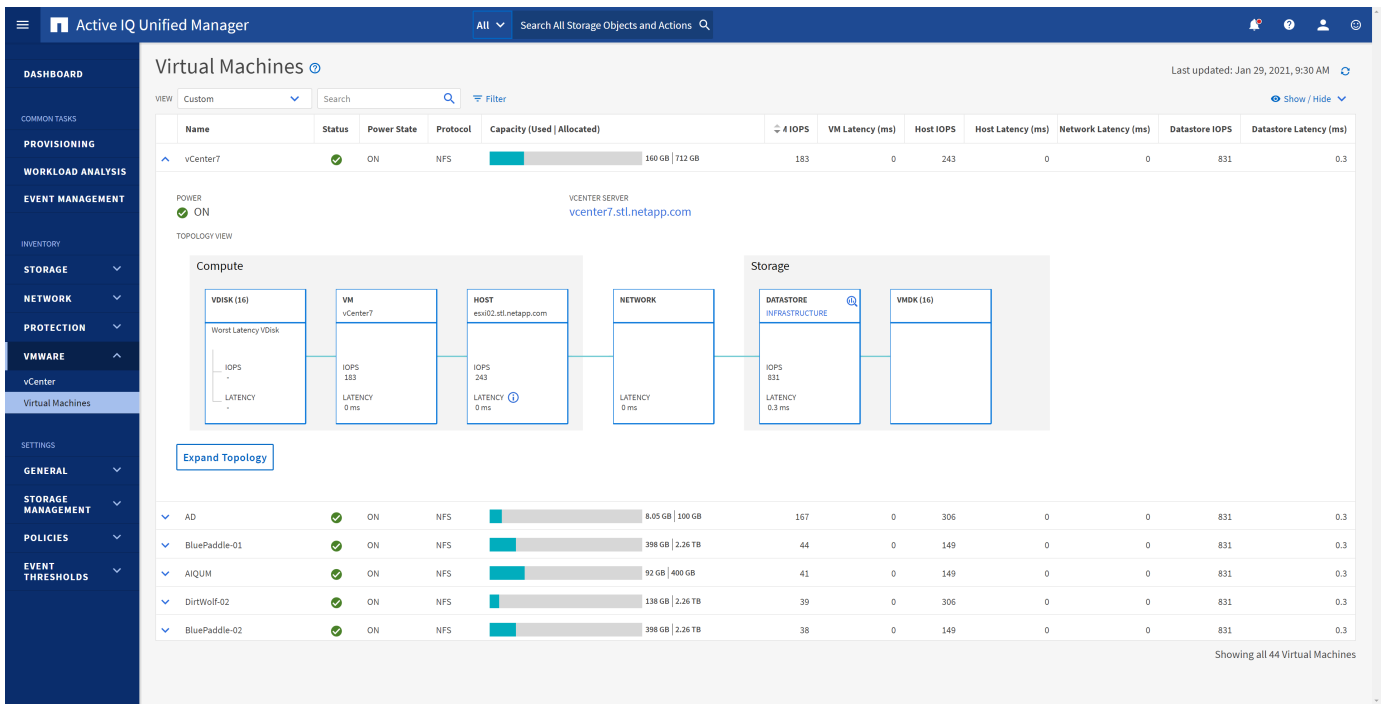
開始使用軟體加密非常簡單。安裝授權之後、只要指定通關密碼、即可設定內建金鑰管理程式、然後建立新的磁碟區、或是執行儲存端磁碟區移轉、即可啟用加密功能。NetApp正致力於在未來的VMware工具版本中、為加密功能提供更多整合式支援。

## Active IQ Unified Manager

利用VMware Infrastructure、您可以清楚掌握虛擬基礎架構中的虛擬機器、並監控及疑難排解虛擬環境中的儲存與效能問題。Active IQ Unified Manager

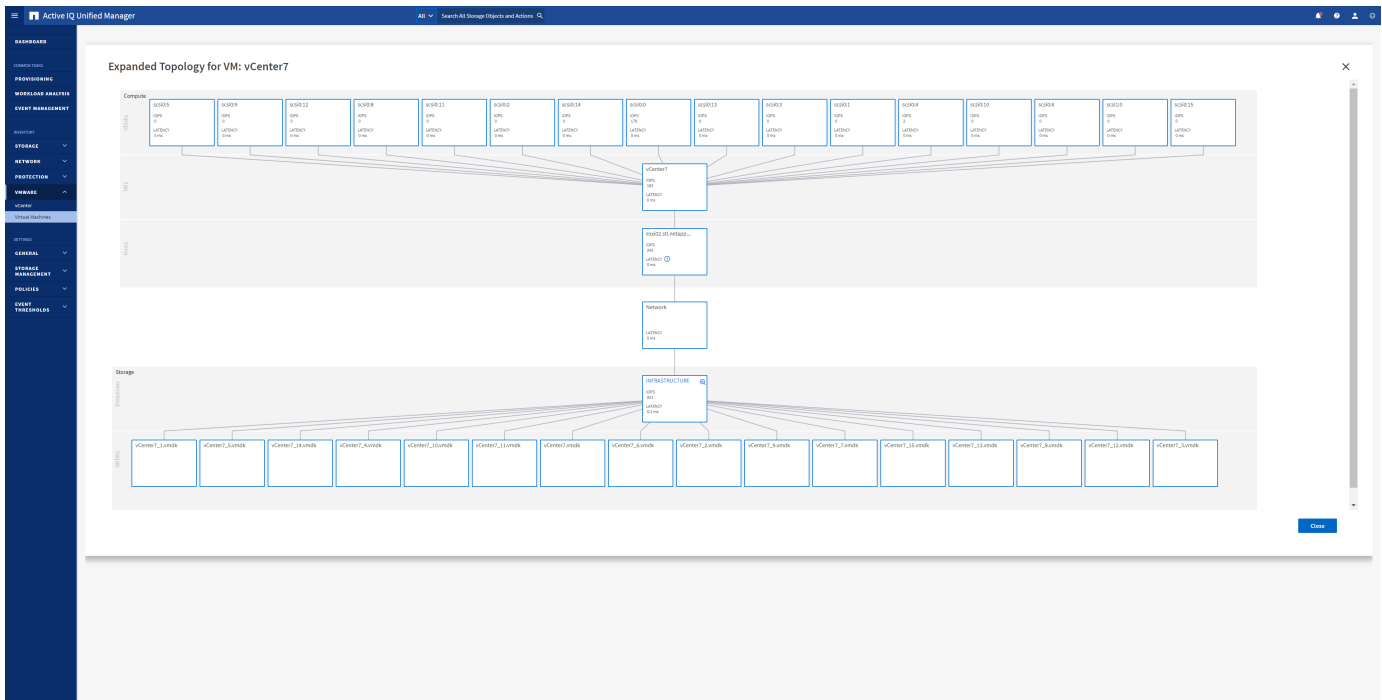
典型的虛擬基礎架構部署ONTAP 在整個運算、網路和儲存層之間、有許多不同的元件。VM應用程式中的任何效能延遲都可能是因為各個元件在各個層面上所面臨的延遲問題。

下列螢幕快照顯示Active IQ Unified Manager 「VMware虛擬機器」檢視畫面。



Unified Manager會在拓撲檢視中呈現虛擬環境的底層子系統、以判斷運算節點、網路或儲存設備是否發生延遲問題。此檢視也會強調導致效能延遲的特定物件、以便採取補救步驟並解決根本問題。

下列螢幕快照顯示AIQUM擴充拓撲。



## 儲存原則型管理和 vVols

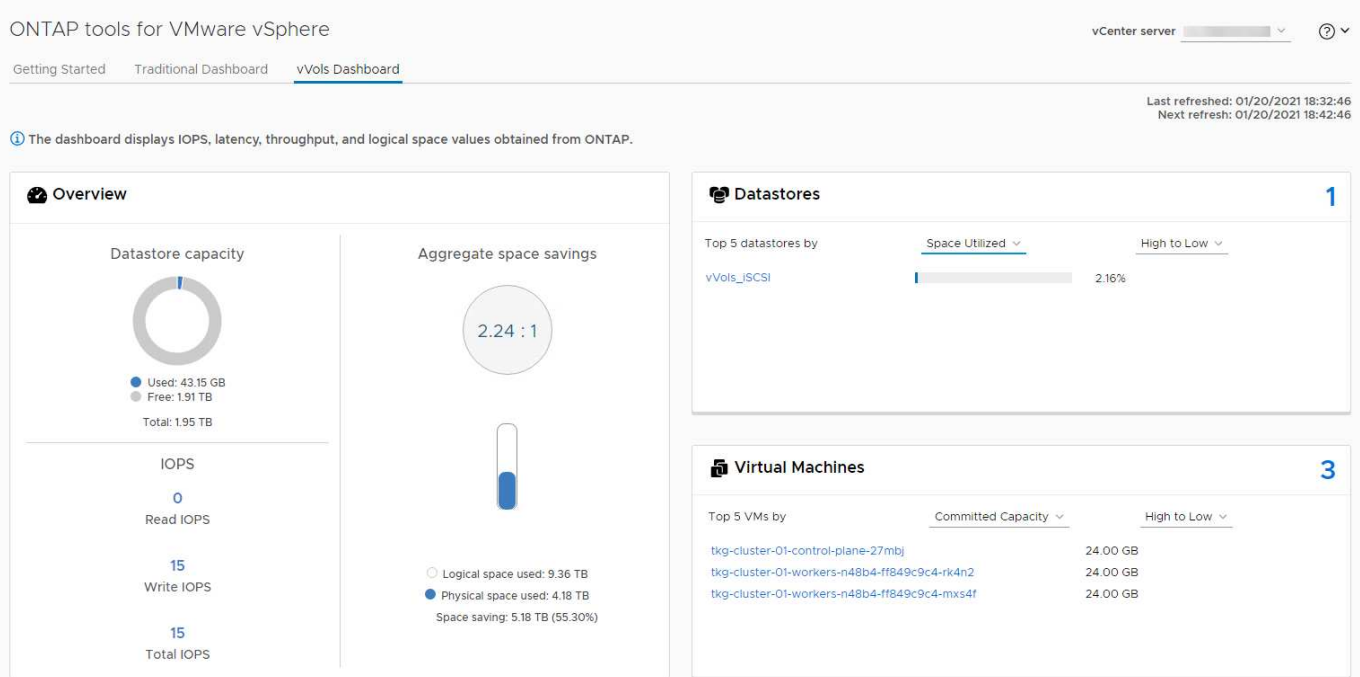
VMware vSphere API for Storage感知 (VASA) 可讓儲存管理員輕鬆設定具有明確定義功能的資料存放區、並讓VM管理員在需要時使用這些功能來配置VM、而不需要彼此互動。

請看一下這種方法、瞭解它如何簡化虛擬化儲存作業、並避免許多瑣碎的工作。

在VASA之前、VM管理員可以定義VM儲存原則、但他們必須與儲存管理員合作、以識別適當的資料存放區、通常是使用文件或命名慣例。有了VASA、儲存管理員可以定義一系列的儲存功能、包括效能、分層、加密及複寫。一組磁碟區或一組磁碟區的功能稱為儲存功能設定檔 (scp)。

SCP 支援虛擬機器資料 VVols 的最低和 / 或最高 QoS。只AFF 有在不支援的系統上才支援最低QoS。VMware vSphere的VMware vSphere工具包含儀表板、可顯示VM精細的效能、以及在VMware系統上用於vVols的邏輯容量。ONTAP ONTAP

下圖說明ONTAP VMware vSphere 9.8 vVols儀表板的各項功能。



定義儲存功能設定檔之後、就可以使用識別其需求的儲存原則來配置VM。VM儲存原則與資料存放區儲存功能設定檔之間的對應、可讓vCenter顯示相容資料存放區清單以供選擇。這種方法稱為儲存原則型管理。

VASA提供查詢儲存設備的技術、並將一組儲存功能傳回vCenter。VASA廠商供應商會提供儲存系統API與架構之間的轉譯、以及vCenter所瞭解的VMware API。NetApp 的 VASA Provider for ONTAP 是 ONTAP 工具的一部分、適用於 VMware vSphere 應用裝置 VM、vCenter 外掛程式則提供介面、可配置及管理 vVol 資料存放區、並可定義儲存功能設定檔 ( CDP )。

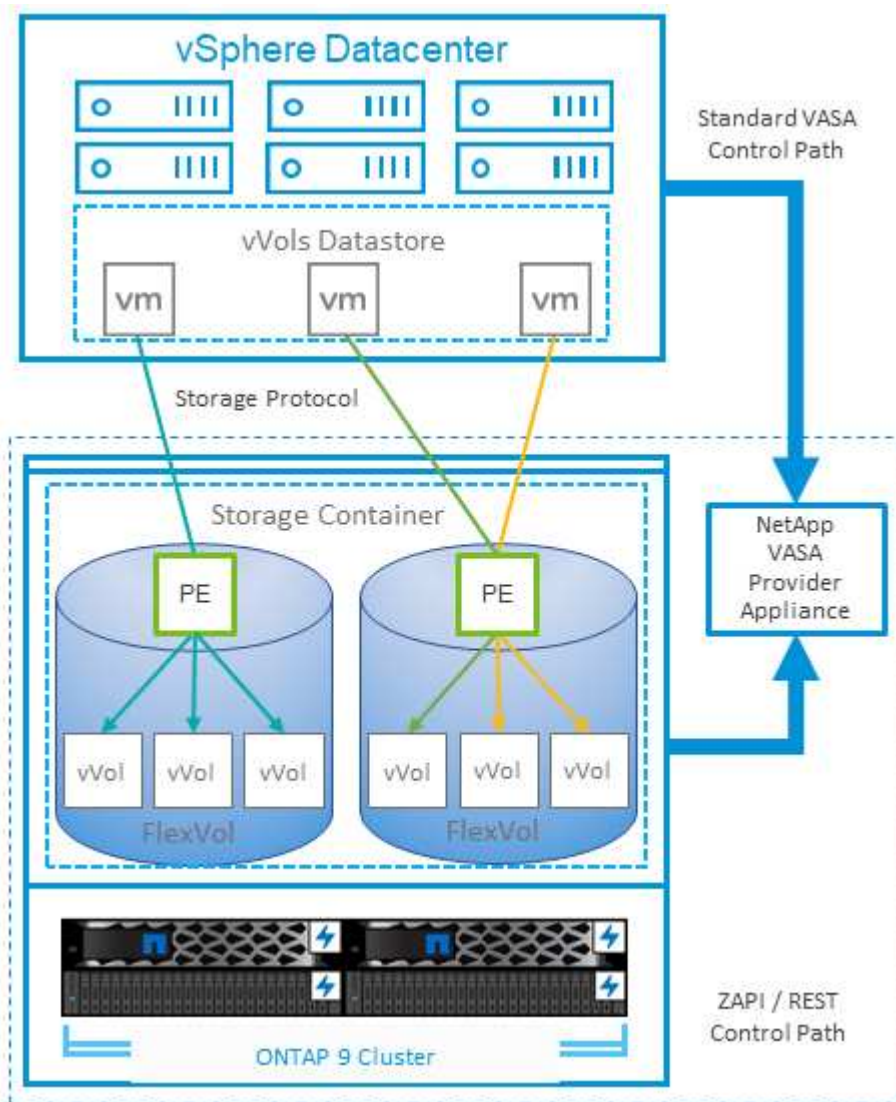
支援VMFS和NFS vVol資料存放區。ONTAP將vVols與SAN資料存放區搭配使用、可帶來NFS的部分效益、例如VM層級的精細度。以下是一些最佳實務做法、您可以在中找到更多資訊 "[TR-4400](#)"：

- VVol資料存放區可由FlexVol 多個叢集節點上的多個支援功能區所組成。最簡單的方法是單一資料存放區、即使磁碟區具有不同的功能也一樣。SPBM可確保VM使用相容的Volume。然而、這些磁碟區必須全部屬於ONTAP 單一的一套功能、並使用單一傳輸協定來存取。每個節點的每個傳輸協定只需一個LIF就足夠了。避免在ONTAP 單一vVol資料存放區中使用多個版本的支援、因為儲存功能可能因版本而異。
- 使用ONTAP VMware vSphere外掛程式的VMware vCenter工具來建立及管理VVol資料存放區。除了管理資料存放區及其設定檔之外、它還會自動建立傳輸協定端點、以便在需要時存取vVols。如果使用LUN、請注意LUN PE是使用LUN ID 300以上的LUN來對應。確認 ESXi 主機進階系統設定 `Disk.MaxLUN` 允許大於 300 的 LUN ID 號碼 (預設值為 1,024)。若要執行此步驟、請在 vCenter 中選取 ESXi 主機、然後選取「設定」索引標籤、再選取「尋找」 `Disk.MaxLUN` 在進階系統設定清單中。
- 請勿安裝或移轉VASA Provider、vCenter Server (應用裝置或Windows) 或ONTAP VMware vSphere的各

種支援工具到vVols資料存放區、因為這些工具彼此相依、因此在停電或其他資料中心中斷時、您無法管理這些工具。

- 定期備份VASA Provider VM。至少要為包含 VASA Provider 的傳統資料存放區建立每小時快照。如需保護及恢復VASA Provider的詳細資訊、請參閱此 ["知識庫文章"](#)。

下圖顯示vVols元件。



## VMware Storage Distributed Resource Scheduler

VMware儲存分散式資源排程器（SDR）是vSphere功能、可根據目前的I/O延遲和空間使用量、將VM放置在儲存設備上。

接著、它會在資料存放區叢集中的資料存放區之間（也稱為Pod）、在不中斷營運的情況下移動VM或VMDK、並選取將VM或VMDK置於資料存放區叢集中的最佳資料存放區。資料存放區叢集是類似資料存放區的集合、從vSphere 管理員的觀點來看、這些資料存放區會彙總成單一使用量單位。

搭配 ONTAP 工具使用適用於 VMware vSphere 的 SDR 時、您必須先使用外掛程式建立資料存放區、使用 vCenter 建立資料存放區叢集、然後將資料存放區新增至該叢集。建立資料存放區叢集之後、可直接從「詳細資料」頁面上的資源配置精靈、將其他資料存放區新增至資料存放區叢集。

SDR的ONTAP 其他最佳實務做法包括：

- 叢集中的所有資料存放區都應該使用相同類型的儲存設備（例如SAS、SATA或SSD）、無論是所有VMFS或NFS資料存放區、都具有相同的複寫和保護設定。
- 請考慮在預設（手動）模式下使用SDR。此方法可讓您檢閱建議、並決定是否要套用建議。請注意VMDK移轉的下列影響：
  - 當SDR在資料存放區之間移動VMDK時、ONTAP 任何從還原複製或重複資料刪除所節省的空间都會遺失。您可以重新執行重複資料刪除、以重新獲得這些節約效益。
  - 在 SDR 移動 VMDK 之後、NetApp 建議在來源資料存放區重新建立快照、因為其他情況下空間會被移動的 VM 鎖定。
  - 在同一個集合體上的資料存放區之間移動VMDK並沒有什麼好處、而且SDR無法看到可能共用該集合體的其他工作負載。

## 建議的ESXi主機和其他ONTAP 功能設定

NetApp 針對 NFS 和區塊傳輸協定開發了一組最佳 ESXi 主機設定。此外、我們也針對多重路徑和 HBA 逾時設定提供具體指引、以根據 NetApp 和 VMware 內部測試、讓 ONTAP 正常運作。

使用適用於 VMware vSphere 的 ONTAP 工具可輕鬆設定這些值：從摘要儀表板、按一下主機系統 Portlet 中的「編輯設定」、或在 vCenter 中的主機上按一下滑鼠右鍵、然後瀏覽至 ONTAP 工具 > 設定建議值。

以下是 9.8~9.13 版本目前建議的主機設定。

主機設定	* NetApp推薦價值*	需要重新開機
* ESXi進階組態*		
VMFS3.HardwareAcceleratedLocking	保留預設值（1）	否
VMFS3.EnableBlockDelete	保留預設值（0）、但可視需要變更。 如需詳細資訊、請參閱 " <a href="#">VMware KB 2007427</a> "	否
VMFS3.EnableVMFS6 取消對應	保留預設值（1） 如需詳細資訊、請參閱 " <a href="#">VMware vSphere API：陣列整合（VAAI）</a> "	否
* NFS 設定 *		
net.TcpipHeapSize.	vSphere 6.0或更新版本、設定為32。 所有其他NFS組態、設定為30	是的
net.TcpipHeapMax	大多數vSphere 6.x版本的設定為512MB。 若為6.5U3、6.7U3及7.0或更新版本、則設為1024MB。	是的



NFS.MaxVolumes	vSphere 6.0 或更新版本、設定為 256 所有其他 NFS 組態都設為 64 。	否
NFS41.MaxVolumes	vSphere 6.0 或更新版本、設定為 256 。	否
NFS.MaxQuesteepeth1^	vSphere 6.0或更新版本、設定為128	是的
NFS.HeartbeatMaxFailures	所有NFS組態皆設為10	否
NFS.Heartbeat頻率	針對所有 NFS 組態設定為 12	否
NFS.Heartbeattimeout	針對所有 NFS 組態設定為 5 。	否
SUNPC。MaxConnPerIP	vSphere 7.0 或更新版本、設定為 128 。	否
* FC/FCoE設定*		
路徑選擇原則	使用具有ALUA的FC路徑時、設定為RR（循環配置資源）。所有其他組態皆設為「固定」。將此值設為RR有助於在所有主動/最佳化路徑之間提供負載平衡。固定值適用於舊版非ALUA組態、有助於防止Proxy I/O換句話說、它有助於在Data ONTAP 以7-Mode運作的環境中、讓I/O不再移至高可用度（HA）配對的其他節點	否
disk.QFullSampleSize.	所有組態皆設為32。 設定此值有助於避免I/O錯誤。	否
disk.QFullThreshold	所有組態均設為8。 設定此值有助於避免I/O錯誤。	否
Emulex FC HBA逾時	使用預設值。	否
QLogic FC HBA逾時	使用預設值。	否
* iSCSI設定*		
路徑選擇原則	所有iSCSI路徑均設為RR（循環配置資源）。將此值設為RR有助於在所有主動/最佳化路徑之間提供負載平衡。	否
disk.QFullSampleSize.	所有組態皆設為32。 設定此值有助於避免I/O錯誤	否
disk.QFullThreshold	所有組態均設為8。 設定此值有助於避免I/O錯誤。	否



1 -使用VMware vSphere ESXi 7.0.1和VMware vSphere ESXi 7.0.2時、NFS進階組態選項MaxQuesteDepth可能無法如預期運作。請參閱 "[VMware KB 86331](#)" 以取得更多資訊。

建立完等量磁碟區和LUN時、也會指定特定的預設設定：ONTAP ONTAP FlexVol

《》 工具* ONTAP	預設設定
Snapshot保留 (-%快照空間)	0%
部分保留 (-分數保留)	0%
存取時間更新 (-atime-update)	錯
最小預先讀取 (-min-readahead)	錯
排程快照	無
儲存效率	已啟用
Volume保證	無 (精簡配置)
Volume自動調整大小	大幅縮減
LUN空間保留	已停用
LUN空間配置	已啟用

## 效能的多重路徑設定

雖然目前未由可用的 ONTAP 工具進行設定、NetApp 仍會建議下列組態選項：

- 在高效能環境中或使用單一LUN資料存放區測試效能時、請考慮將循環配置資源 (VMW\_PSP\_RR) 路徑選擇原則 (PSP) 的負載平衡設定、從預設的IOPS設定1000變更為1。請參閱VMware KB "[2069356](#)" 以取得更多資訊。
- 在vSphere 6.7 Update 1中、VMware為Round Robin PSP引進新的延遲負載平衡機制。新選項會在選取I/O最佳路徑時、考量I/O頻寬和路徑延遲您可以在具有非對等路徑連線能力的環境中使用它、例如在一條路徑上有比另一條路徑上有更多網路躍點的情況、或是在使用 NetApp All SAN Array 系統時使用。請參閱 "[路徑選取外掛程式和原則](#)" 以取得更多資訊。

## 其他文件

對於帶有 vSphere 7 的 FCP 和 iSCSI、詳細資料請參閱 "[搭配 ONTAP 使用 VMware vSphere 7.x](#)"

對於帶有 vSphere 8 的 FCP 和 iSCSI、詳細資料請參閱 "[搭配 ONTAP 使用 VMware vSphere 8.x](#)"

如需使用 vSphere 7 的 NVMe、請參閱 "[如需更多詳細資料、請參閱適用於 ESXi 7.x with ONTAP 的 NVMe 主機組態](#)"

如需使用 vSphere 8 的 NVMe、請參閱 "[如需更多詳細資料、請參閱適用於 ESXi 8.x 與 ONTAP 的 NVMe 主機組態](#)"

# 使用 ONTAP 的虛擬磁碟區 ( VVols )

## 總覽

ONTAP 在過去二十多年來一直是 VMware vSphere 環境的領先儲存解決方案、並持續新增創新功能來簡化管理、同時降低成本。

本文件涵蓋適用於 VMware vSphere 虛擬磁碟區 ( VVols ) 的 ONTAP 功能、包括最新的產品資訊和使用案例、以及最佳實務做法和其他資訊、可簡化部署並減少錯誤。



本文件將先前發佈的技術報告 [\\_TR-4400: VMware vSphere 虛擬磁碟區 \(vVols\) 取代為 ONTAP](#)

最佳實務做法是輔助其他文件、例如指南和相容性清單。這些技術是根據實驗室測試和NetApp工程師與客戶廣泛的現場經驗所開發。它們可能不是唯一可行或受支援的實務做法、但通常是最簡單的解決方案、可滿足大多數客戶的需求。



本文件已更新、納入 vSphere 8.0 更新 1 中新增的 vVols 功能、ONTAP 工具 9.12 版本支援這些功能。

## 虛擬磁碟區 (vVols) 概觀

NetApp 於 2012 年開始與 VMware 合作、為 vSphere 5 支援 vSphere API for Storage Aware (VASA)。此早期的 VASA Provider 允許定義設定檔中的儲存功能、以便在資源配置時用來篩選資料存放區、並在之後檢查原則是否符合規定。隨著時間演進、新增新功能以實現更多自動化資源配置、並新增虛擬磁碟區或 vVols、其中個別的儲存物件用於虛擬機器檔案和虛擬磁碟。這些物件可以是 LUN、檔案、現在也可以是 vSphere 8 : NVMe namespaces。NetApp 與 VMware 密切合作、成為 2015 年 vSphere 6 隨附的 vVols 參考合作夥伴、再次成為 vSphere 8 中使用 NVMe over Fabric 的 vVols 設計合作夥伴。NetApp 持續強化 vVols、以充分利用 ONTAP 的最新功能。

有幾個元件需要注意：

<p>* VASA Provider*</p>
<p>這是處理 VMware vSphere 與儲存系統之間通訊的軟體元件。對於 ONTAP、VASA Provider 會在名為 ONTAP tools for VMware vSphere 的應用裝置 (簡稱 ONTAP 工具) 中執行。ONTAP 工具也包括 vCenter 外掛程式、適用於 VMware Site Recovery Manager 的儲存複寫介面卡 (SRA)、以及用於建置自己自動化的 REST API 伺服器。ONTAP 工具在 vCenter 中設定及登錄後、就不再需要直接與 ONTAP 系統互動、因為幾乎所有的儲存需求都可以直接從 vCenter UI 中進行管理、或透過 REST API 自動化來進行管理。</p>
<p>* 傳輸協定端點 (PE) *</p>
<p>傳輸協定端點是 ESXi 主機和 vVols 資料存放區之間 I/O 的 Proxy。ONTAP VASA Provider 會自動建立這些 LUN、每個 vVols 資料存放區的 FlexVol 磁碟區一個傳輸協定端點 LUN (大小為 4 MB)、或是在資料存放區中主控 FlexVol 磁碟區的儲存節點上、每個 NFS 介面 (LIF) 一個 NFS 裝載點。ESXi 主機直接裝載這些傳輸協定端點、而非個別的 vVol LUN 和虛擬磁碟檔案。不需要管理 VASA Provider 自動建立、掛載、卸載及刪除的傳輸協定端點、以及任何必要的介面群組或匯出原則。</p>
<p>* 虛擬傳輸協定端點 (VPE) *</p>
<p>vSphere 8 的新功能是、搭配 vVols 使用 NVMe over Fabrics (NVMe of) 時、ONTAP 中不再有傳輸協定端點的概念。相反地、只要第一個虛擬機器開機、ESXi 主機就會為每個 ANA 群組自動產生一個虛擬 PE。ONTAP 會自動為資料存放區所使用的每個 FlexVol 磁碟區建立全日空群組。</p> <p>使用 NVMe for vVols 的另一個優點是、VASA Provider 不需要任何繫結要求。而是由 ESXi 主機根據 VPE 在內部處理 vVol 繫結功能。這可減少 vVol 綁定風暴對服務造成影響的機會。</p>
<p>如需詳細資訊、請參閱 <a href="#">"NVMe 和虛擬磁碟區"</a> 開啟 <a href="#">"vmware.com"</a></p>
<p>* 虛擬磁碟區資料存放區 *</p>

虛擬 Volume 資料存放區是 VVols 容器的邏輯資料存放區呈現、由 VASA Provider 建立和維護。該容器代表從 VASA Provider 管理的儲存系統所配置的儲存容量集區。ONTAP 工具支援將多個 FlexVol 磁碟區（稱為備份磁碟區）分配至單一 VVols 資料存放區、這些 VVols 資料存放區可跨越 ONTAP 叢集中的多個節點、將 Flash 和混合式系統結合在一起、提供不同的功能。管理員可以使用資源配置精靈或 REST API 來建立新的 FlexVol Volume、或選擇預先建立的 FlexVol Volume 來備份儲存區（如果有的話）。

\* 虛擬磁碟區（vVols） \*

VVols 是儲存在 vVols 資料存放區中的實際虛擬機器檔案和磁碟。使用術語 vVol（奇異）是指單一特定檔案、LUN 或命名空間。ONTAP 會根據資料存放區使用的傳輸協定、建立 NVMe 命名空間、LUN 或檔案。有幾種不同類型的 vVols；最常見的是組態（中繼資料檔案）、資料（虛擬磁碟或 VMDK）和交換（在 VM 開機時建立）。受 VMware VM 加密保護的 VVols 屬於其他類型。VMware VM 加密不應與 ONTAP Volume 或 Aggregate 加密混淆。

## 原則型管理

VMware vSphere API for Storage Aware（VASA）可讓 VM 管理員輕鬆使用所需的任何儲存功能來配置 VM、而無需與儲存團隊互動。在 VASA 之前、VM 管理員可以定義 VM 儲存原則、但必須與儲存管理員合作、以識別適當的資料存放區、通常是使用文件或命名慣例。有了 VASA、具有適當權限的 vCenter 管理員就能定義一系列儲存功能、vCenter 使用者隨後可以使用這些功能來配置 VM。VM 儲存原則與資料存放區儲存功能設定檔之間的對應可讓 vCenter 顯示相容資料存放區清單以供選擇、並可啟用其他技術、例如 Aria（前身為 vRealize）Automation 或 Tanzu Kubernetes Grid、以自動從指派的原則中選取儲存區。這種方法稱為儲存原則型管理。雖然儲存功能設定檔和原則也可用於傳統的資料存放區、但我們的重點是 vVols 資料存放區。

有兩個要素：

\* 儲存功能設定檔（SCP） \*

儲存功能設定檔（SCP）是一種儲存範本形式、可讓 vCenter 管理員定義所需的儲存功能、而無需實際瞭解如何在 ONTAP 中管理這些功能。採用範本樣式的方法、讓管理員能夠以一致且可預測的方式輕鬆提供儲存服務。SCP 中說明的功能包括效能、傳輸協定、儲存效率及其他功能。具體功能因版本而異。使用 vCenter UI 中的 ONTAP Tools for VMware vSphere 功能表來建立這些工具。您也可以使用 REST API 來建立 CDP。您可以選擇個別功能來手動建立這些功能、或是從現有（傳統）資料存放區自動產生。

\* VM 儲存原則 \*

VM 儲存原則是在 vCenter 的原則和設定檔下建立。對於 vVols、請使用 NetApp vVols 儲存類型供應商的規則來建立規則集。ONTAP 工具提供簡化的方法、讓您只需選取一個 SCP、而非強制您指定個別規則。

如上所述、使用原則有助於簡化資源配置磁碟區的工作。只要選取適當的原則、VASA Provider 就會顯示支援該原則的 VVols 資料存放區、並將 VVol 放入符合法規的個別 FlexVol 磁碟區（圖 1）。

## 使用儲存原則部署 VM

## New Virtual Machine

- ✓ 1 Select a creation type
- ✓ 2 Select a name and folder
- ✓ 3 Select a compute resource
- 4 Select storage**
- 5 Select compatibility
- 6 Select a guest OS
- 7 Customize hardware
- 8 Ready to complete

### Select storage

Select the storage for the configuration and disk files

Encrypt this virtual machine (Requires Key Management Server)

VM Storage Policy

Platinum

Disable Storage DRS for this virtual machine

	Name	Storage Compatibility	Capacity	Provisioned	Free	Type	Clu
<input checked="" type="radio"/>	vVolsiSCSI	Compatible	100 GB	40.74 GB	64.88 GB	vVol	
<input type="radio"/>	vVolsNFS2202...	Compatible	2 TB	36.88 GB	1.96 TB	vVol	
<input type="radio"/>	local-esx01	Incompatible	3.63 TB	1.46 GB	3.63 TB	VMFS 6	
<input type="radio"/>	local-esx07	Incompatible	1.81 TB	3.85 GB	1.81 TB	VMFS 6	
<input type="radio"/>	local-esx08	Incompatible	1.69 TB	1.43 GB	1.69 TB	VMFS 6	
<input type="radio"/>	local-esx09	Incompatible	1.81 TB	3.85 GB	1.81 TB	VMFS 6	
<input type="radio"/>	local-esx15	Incompatible	3.63 TB	1.46 GB	3.63 TB	VMFS 6	
<input type="radio"/>	tier001_ds	Incompatible	22 TB	23.73 TB	18.09 TB	NFS v3	

CANCEL

BACK

NEXT

VM 佈建完成後、VASA Provider 將繼續檢查法規遵循狀況、並在備用磁碟區不再符合原則時、在 vCenter 中警示 VM 管理員 (圖 2)。

VM 儲存原則法規遵循

## Storage Policies

### VM Storage Policies

AFF\_VASA10

### VM Storage Policy Compliance

Noncompliant

### Last Checked Date

5/20/2022, 12:59:35 PM

### VM Replication Groups

[CHECK COMPLIANCE](#)

## NetApp VVols 支援

ONTAP 自 2012 年首次推出 VASA 規格以來、就一直提供支援。雖然其他 NetApp 儲存系統可能支援 VASA、但本文件著重於目前支援的 ONTAP 9 版本。

### ONTAP

除了 AFF、ASA 和 FAS 系統上的 ONTAP 9 之外、NetApp 還支援 ONTAP Select 上的 VMware 工作負載、適用於 NetApp 的 Amazon FSX 搭配 AWS 上的 VMware Cloud、Azure NetApp Files 搭配 Azure VMware 解決方案、Cloud Volumes Service 搭配 Google Cloud VMware Engine、以及 Equinix 中的 NetApp 私有儲存設備、但具體功能可能會因服務供應商和可用的網路連線而異。也可從 vSphere 來賓存取儲存在這些組態中的資料、以及 Cloud Volumes ONTAP。

在發佈時、超大規模環境僅限於傳統的 NFS v3 資料存放區、因此 VVols 僅適用於內部部署 ONTAP 系統、或雲端連線系統、這些系統提供內部部署系統的完整功能、例如由全球各地的 NetApp 合作夥伴和服務供應商代管的系統。

\_ 如需 ONTAP 的詳細資訊、請參閱 "[產品文件ONTAP](#)" \_

\_ 如需 ONTAP 和 VMware vSphere 最佳實務做法的詳細資訊、請參閱 "[TR-4597](#)" \_

### 搭配 ONTAP 使用 vVols 的優點

當 VMware 在 2015 年推出 VVols 支援 VASA 2.0 時、他們將其描述為「整合與管理架構、為外部儲存設備（SAN/NAS）提供全新的作業模式」。此作業模式可提供多項優點、搭配 ONTAP 儲存設備使用。

#### 原則型管理

如第 1.2 節所述、原則型管理可讓 VM 使用預先定義的原則進行佈建及後續管理。這有助於 IT 作業的多種方式：

- \* 提高速度。\* ONTAP 工具不需要 vCenter 管理員與儲存團隊一起開啟儲存資源配置活動的問題單。不過、vCenter 和 ONTAP 系統上的 ONTAP 工具 RBAC 角色仍可允許個別的團隊（例如儲存團隊）、或是由同一個團隊進行個別活動、只要有需要、就能限制特定功能的存取。
- \* 更聰明的資源配置。\* 儲存系統功能可透過 VASA API 公開、讓資源配置工作流程能夠充分利用進階功能、而無需 VM 管理員瞭解如何管理儲存系統。
- \* 更快的資源配置。\* 可在單一資料存放區中支援不同的儲存功能、並根據 VM 原則自動選擇適合的 VM。
- \* 避免錯誤。\* 儲存和 VM 原則是事先開發的、並可視需要套用、而無需每次佈建 VM 時都自訂儲存設備。當儲存功能從定義的原則中移出時、就會發出法規遵循警報。如前所述、SCP 可讓初始資源配置可預測且可重複執行、而 VM 儲存原則則以 SCP 為基礎、則可確保正確放置。
- \* 更好的容量管理。\* VASA 和 ONTAP 工具可讓您視需要將儲存容量向下檢視至大量的彙總層級、並在容量開始不足時提供多層警示。

#### 現代化 SAN 上的 VM 精細管理

使用光纖通道和 iSCSI 的 SAN 儲存系統是 VMware 首次支援 ESX 的系統、但它們缺乏從儲存系統管理個別 VM 檔案和磁碟的能力。而是配置 LUN 並由 VMFS 管理個別檔案。這使得儲存系統難以直接管理個別 VM 儲存效能、複製和保護。VVols 提供 ONTAP 強大、高效能的 SAN 功能、讓使用 NFS 儲存設備的客戶能夠享有更精細的儲存空間。

現在、使用適用於 VMware vSphere 9.12 及更新版本的 vSphere 8 和 ONTAP 工具、Vols 對於舊版 SCSI 型傳

輸協定所使用的相同精細控制功能現在也可在採用 NVMe over Fabrics 的現代化光纖通道 SAN 中使用、以及在規模上獲得更高的效能。有了 vSphere 8.0 更新 1、現在可以使用 vVols 部署完整的端點對端 NVMe 解決方案、而無需在 Hypervisor 儲存堆疊中進行任何 I/O 轉譯。

#### 更強大的儲存卸載功能

雖然 VAAI 提供多種卸載至儲存設備的作業、但 VASA Provider 仍會解決一些落差。SAN VAAI 無法將 VMware 託管的快照卸載至儲存系統。NFS VAAI 可以卸載 VM 託管的快照、但儲存原生快照對 VM 有限制。由於 VVols 使用個別 LUN、命名空間或檔案來儲存虛擬機器磁碟、因此 ONTAP 可以快速有效地複製檔案或 LUN、以建立不再需要差異檔案的 VM 精細快照。NFS VAAI 也不支援卸載熱（開啟電源）Storage VMotion 移轉的複製作業。當使用 VAAI 搭配傳統 NFS 資料存放區時、必須關閉虛擬機器電源、以允許移轉卸載。ONTAP 工具中的 VASA Provider 可提供近乎即時且具儲存效率的複本、以進行熱移轉和冷移轉、也支援近乎即時的 vVols 跨磁碟區移轉複本。由於這些顯著的儲存效率效益、您可能可以在中充分利用 vVols 工作負載 "效率保證" 方案。同樣地、如果使用 VAAI 的跨磁碟區複製無法滿足您的需求、您可能會因為 vVols 複製體驗的改善而解決您的業務挑戰。

#### vVols 的常見使用案例

除了這些優點之外、我們也會看到 vVol 儲存設備的常見使用案例：

- \* 隨選虛擬機器資源配置 \*
  - 私有雲或服務供應商 IaaS。
  - 透過 Aria（前身為 vRealize）套件、OpenStack 等、充分運用自動化與協調功能
- \* 一流磁碟（FCD）\*
  - VMware Tanzu Kubernetes Grid [TKG] 持續磁碟區。
  - 透過與 VMDK 生命週期管理功能相隨的方式、提供 Amazon EBS 般的服務。
- \* 隨需提供暫存虛擬機器 \*
  - 測試 / 開發實驗室
  - 訓練環境

#### vVols 的常見優點

在充分發揮其優勢時（例如在上述使用案例中）、vVols 提供下列具體改善：

- 在單一磁碟區內或 ONTAP 叢集中的多個磁碟區之間快速建立複本、相較於傳統的 VAAI 複本、這是一項優勢。而且儲存效率也很高。磁碟區內的複製作業會使用 ONTAP 檔案複製、就像 FlexClone 磁碟區一樣、而且只會儲存來源 vVol 檔案 /LUN/ 命名空間的變更。因此、為了生產或其他應用程式的目的而建立的長期虛擬機器會迅速建立、佔用最少空間、並可從虛擬機器層級保護（使用適用於 VMware vSphere 的 NetApp SnapCenter 外掛程式、VMware 託管快照或 VADP 備份）和效能管理（搭配 ONTAP QoS）中獲益。
- VVols 是搭配 vSphere CSI 使用 TKG 時的理想儲存技術、可提供由 vCenter 管理員管理的獨立儲存類別和容量。
- Amazon EBS 類似的服務可透過 FCD 提供、因為 FCD VMDK 就像名稱所示、是 vSphere 中的一流公民、生命週期可獨立管理、與可能附加的虛擬機器分開管理。

## 搭配 ONTAP 使用 vVols

將 VVols 搭配 ONTAP 使用的關鍵在於 VASA Provider 軟體、此軟體是 VMware vSphere

虛擬應用裝置 ONTAP 工具的一部分。

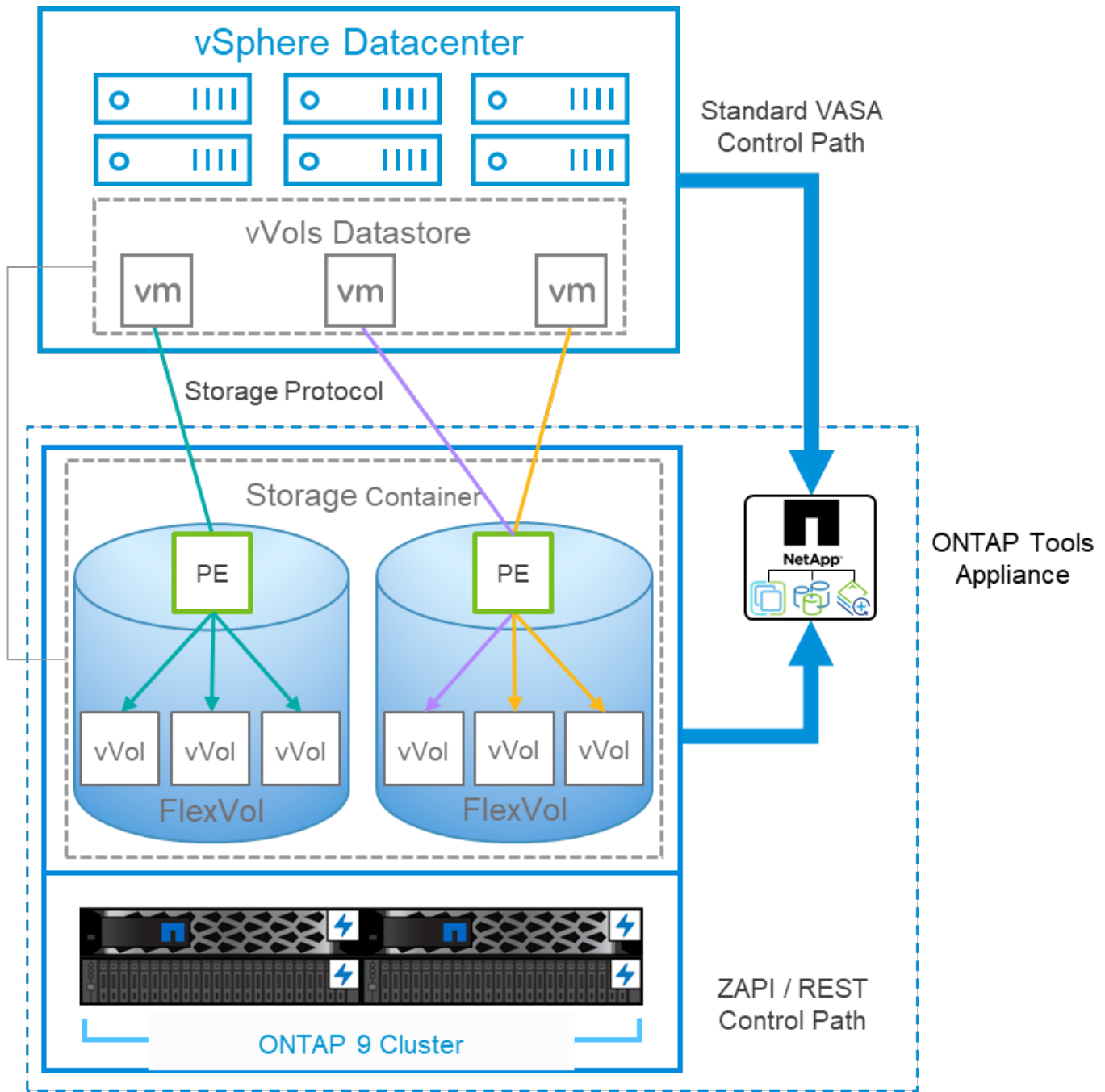
ONTAP 工具也包括 vCenter UI 擴充、REST API 伺服器、適用於 VMware Site Recovery Manager 的儲存複寫介面卡、監控和主機組態工具、以及一系列可協助您更妥善管理 VMware 環境的報告。

#### 產品與文件

ONTAP FlexClone 授權（隨附於 ONTAP One）和 ONTAP 工具應用裝置是唯一需要搭配 ONTAP 使用 vVols 的其他產品。ONTAP 工具的最新版本是以單一統一應用裝置的形式提供、可在 ESXi 上執行、提供以前三種不同應用裝置和伺服器的功能。對於 vVols、使用 ONTAP 工具 vCenter UI 延伸或 REST API 作為 ONTAP 與 vSphere 功能的一般管理工具和使用界面、以及提供特定 vVols 功能的 VASA 供應商、是非常重要的。SRA 元件包含在傳統資料存放區中、但 VMware Site Recovery Manager 不會將 SRA 用於 vVols、而是在 SRM 8.3 及更新版本中實作新服務、利用 vVols 複寫的 VASA 提供者。

使用 iSCSI 或 FCP 時的 ONTAP 工具 VASA Provider 架構





## 產品安裝

若要進行新的安裝、請將虛擬應用裝置部署到 vSphere 環境中。目前版本的 ONTAP 工具會自動向 vCenter 註冊、並預設啟用 VASA Provider。除了 ESXi 主機和 vCenter Server 資訊外、您還需要應用裝置的 IP 位址組態詳細資料。如前所述、VASA Provider 要求 ONTAP FlexClone 授權已安裝在您計畫用於 vVols 的任何 ONTAP 叢集上。此應用裝置內建監控程式、可確保可用性、最佳實務做法是設定 VMware High Availability 及選擇性的容錯功能。如需其他詳細資料、請參閱第 4.1 節。請勿將 ONTAP 工具應用裝置或 vCenter Server 應用裝置 (VCSA) 安裝或移至 vVols 儲存設備、因為這可能會導致應用裝置無法重新啟動。

使用 NetApp 支援網站 (NSS) 上可供下載的升級 ISO 檔案、即可支援 ONTAP 工具的就地升級。請依照部署與設定指南的指示來升級應用裝置。

如需調整虛擬應用裝置規模及瞭解組態限制、請參閱本知識庫文章：["VMware vSphere ONTAP 工具規模調整指南"](#)

## 產品文件

下列文件可協助您部署 ONTAP 工具。

"如需完整的文件儲存庫與擴大機； #44 ；請造訪 [docs.netapp.com](https://docs.netapp.com) 的連結"

## 開始使用

- "版本資訊"
- "瞭解適用於 VMware vSphere 的 ONTAP 工具"
- "VMware 工具快速入門 ONTAP"
- "部署 ONTAP 各種工具"
- "升級 ONTAP 功能"

## 使用 ONTAP VMware 工具

- "配置傳統資料存放區"
- "配置 vVols 資料存放區"
- "設定角色型存取控制"
- "設定遠端診斷"
- "設定高可用度"

## 保護及管理資料存放區

- "保護傳統的資料存放區" 使用 SRM
- "保護 vVols 型虛擬機器" 使用 SRM
- "監控傳統的資料存放區和虛擬機器"
- "監控 vVols 資料存放區和虛擬機器"

除了產品文件之外、還有支援知識庫文章、這些文章可能很有用。

- "如何執行 VASA Provider 災難恢復解決方案指南"

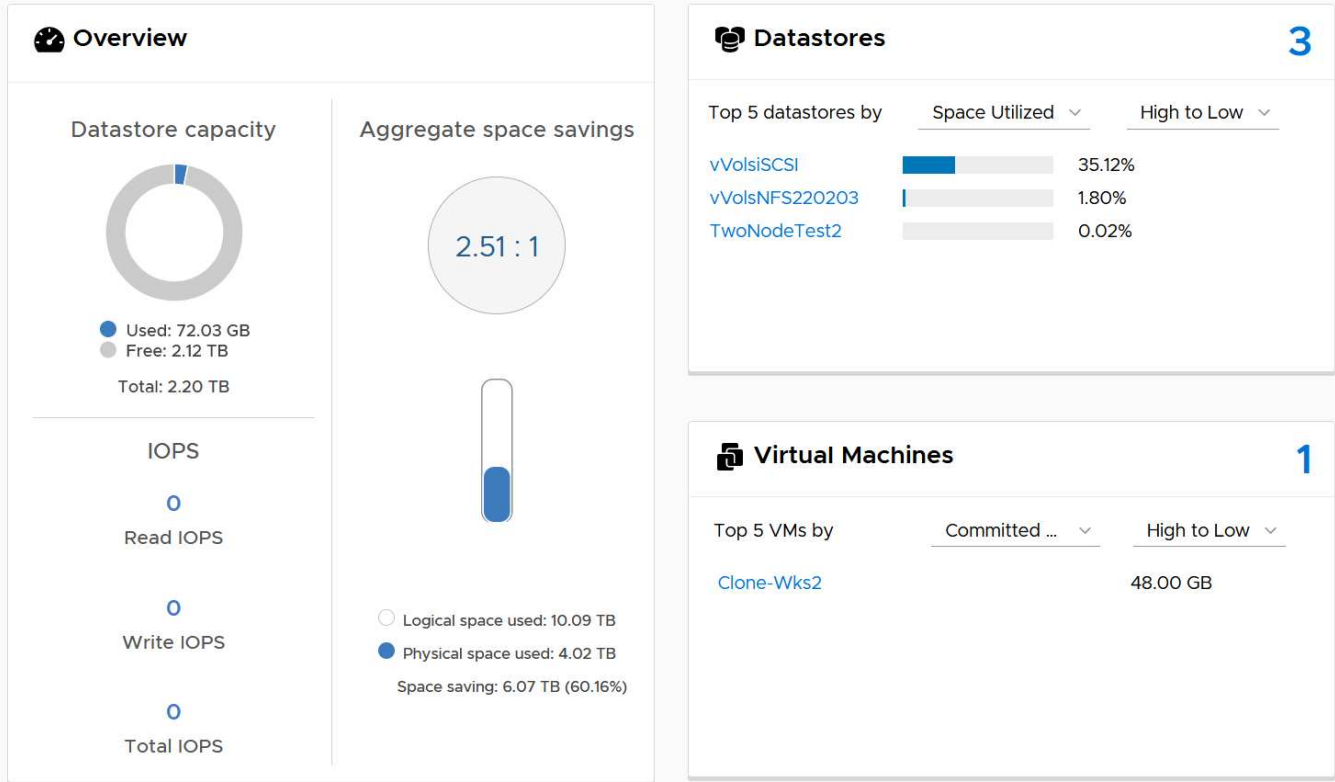
## Vasa Provider Dashboard

VASA Provider 包含儀表板、內含個別 vVols VM 的效能與容量資訊。這些資訊直接來自 VVol 檔案和 LUN 的 ONTAP、包括前 5 名虛擬機器的延遲、IOPS、處理量和正常運作時間、以及前 5 名資料存放區的延遲和 IOPS。使用 ONTAP 9.7 或更新版本時、預設會啟用此功能。擷取並顯示在儀表板中的初始資料可能需要 30 分鐘。

## ONTAP 工具 vVols 儀表板

Last refreshed: 05/20/2022 15:00:57  
Next refresh: 05/20/2022 15:10:57

? The dashboard displays IOPS, latency, throughput, and logical space values obtained from ONTAP.



**最佳實務做法**

搭配 vSphere 使用 ONTAP vVols 非常簡單、並遵循已發佈的 vSphere 方法（請參閱您的 ESXi 版本的 VMware 文件中的「在 vSphere 儲存環境下使用虛擬磁碟區」）。以下是一些與 ONTAP 一起考量的額外實務做法。

- 限制 \*

一般而言、ONTAP 支援 VMware 所定義的 VVols 限制（請參閱已發佈 "組態最大值"）。下表摘要列出 VVols 大小和數量的特定 ONTAP 限制。請務必檢查 "NetApp Hardware Universe" 以獲得 LUN 和檔案數量和大小的更新限制。

- ONTAP VVols 限制 \*

容量/功能	SAN (SCSI 或 NVMe of)	NFS
VVols 大小上限	62 TiB*	62 TiB*
每個 FlexVol Volume 的最大 VVols 數	1024	20 億
每個 ONTAP 節點的 VVols 數目上限	最多 12288**	500 億
每個 ONTAP 配對的最大 VVols 數	高達 24576 個 **	500 億

容量/功能	SAN (SCSI 或 NVMe of)	NFS
每個 ONTAP 叢集的 VVols 數量上限	高達 98,304**	沒有特定的叢集限制
QoS 物件上限 (共享原則群組和個別 VVols 服務層級)	12、000 至 ONTAP 9.3 ; 40、000、含 ONTAP 9.4 及更新版本	

- 大小限制是根據執行 ONTAP 9.12.1P2 及更新版本的 ASA 系統或 AFF 和 FAS 系統而定。
  - SAN vVols (NVMe 命名空間或 LUN) 的數量會因平台而異。請務必檢查 ["NetApp Hardware Universe"](#) 以獲得 LUN 和檔案數量和大小的更新限制。
- 將 ONTAP 工具用於 VMware vSphere 的 UI 延伸或 REST API、以佈建 vVols 資料存放區 \*\* 和傳輸協定端點 \*

雖然可以使用一般 vSphere 介面建立 vVols 資料存放區、但使用 ONTAP 工具會視需要自動建立傳輸協定端點、並使用 ONTAP 最佳實務做法並符合您定義的儲存功能設定檔來建立 FlexVol 磁碟區。只要在主機 / 叢集 / 資料中心上按一下滑鼠右鍵、然後選取 ONTAP tools\_ 和 *Provision datastortity* 即可。您只需在精靈中選擇所需的 vVols 選項即可。

- 切勿將 ONTAP 工具應用裝置或 vCenter Server Appliance (VCSA) 儲存在他們正在管理的 VVols 資料存放區。 \*

如果您需要重新開機設備、這可能會導致「雞和蛋的情況」、因為它們在重新開機時無法重新連結自己的 vVols。您可以將它們儲存在由不同 ONTAP 工具和 vCenter 部署所管理的 vVols 資料存放區。

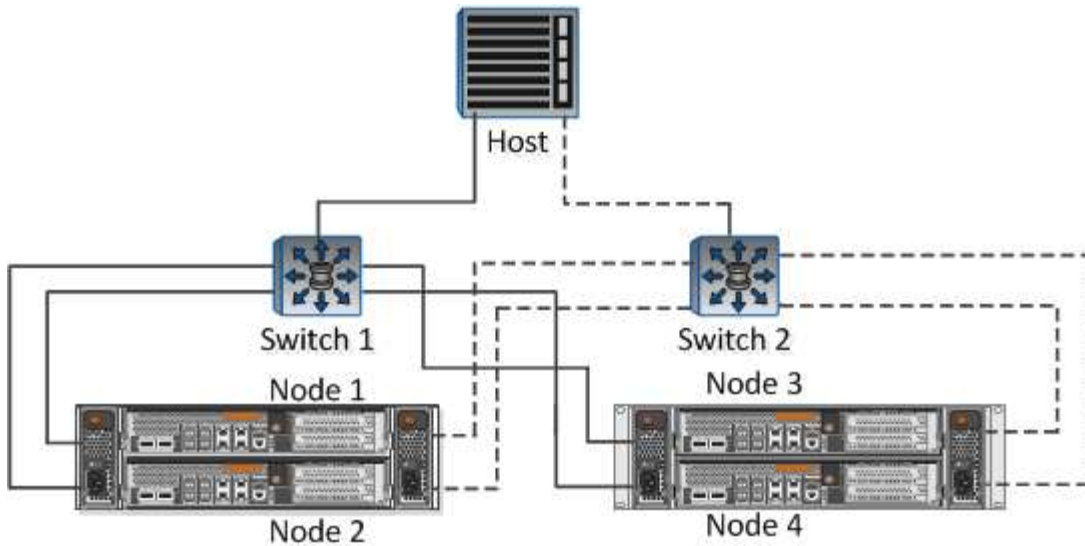
- 避免在不同的 ONTAP 版本中執行 vVols 作業。 \*

支援的儲存功能 (例如 QoS、特性設定等) 已在 VASA Provider 的不同版本中有所變更、有些則視 ONTAP 版本而定。在 ONTAP 叢集中使用不同版本、或在具有不同版本的叢集之間移動 vVols、可能會導致非預期行為或法規遵循警示。

- 使用 NVMe / FC 或 FCP for vVols 之前、請先將光纖通道架構分區。 \*

ONTAP 工具 VASA 供應商負責管理 FCP、iSCSI 群組、以及 ONTAP 中的 NVMe 子系統、這些子系統是以受管理 ESXi 主機的探索啟動器為基礎。不過、它並未與光纖通道交換器整合以管理分區。在進行任何資源配置之前、必須根據最佳實務做法進行分區。以下是單一啟動器分區至四個 ONTAP 系統的範例：

單一啟動器分區：



如需更多最佳實務做法、請參閱下列文件：

["TR-4080 現代 SAN ONTAP 9 的最佳實務做法"](#)

["\\_TR-4684 使用 NVMe 來實作和設定現代化 SAN"](#)

- 根據您的需求規劃您的支援 FlexVols 。 \*

您可以將多個備份磁碟區新增至 vVols 資料存放區、以便在 ONTAP 叢集上分散工作負載、支援不同的原則選項、或增加允許的 LUN 或檔案數量。不過、如果需要最高的儲存效率、請將所有的備份磁碟區放在單一集合體上。或者、如果需要最大的複製效能、請考慮使用單一 FlexVol 磁碟區、並將範本或內容庫保留在相同的磁碟區中。VASA Provider 將許多 VVols 儲存作業卸載至 ONTAP、包括移轉、複製和快照。在單一 FlexVol 磁碟區內完成此作業時、會使用節省空間的檔案複本、而且幾乎可以立即使用。當跨 FlexVol 磁碟區執行此作業時、複本會快速可用、並使用即時重複資料刪除和壓縮功能、但在背景工作使用背景重複資料刪除和壓縮在磁碟區上執行之前、最大的儲存效率可能無法恢復。視來源和目的地而定、部分效率可能會降低。

- 讓儲存功能設定檔（SCP）保持簡單。 \*

避免將功能設定為任何、以指定不需要的功能。這可將選擇或建立 FlexVol 磁碟區時發生的問題減至最低。例如、在 VASA Provider 7.1 及更早版本中、如果將壓縮保留在預設的 SCP 設定「否」、則會嘗試停用壓縮、即使在 AFF 系統上也一樣。

- 使用預設的 SCP 做為範例範本來建立您自己的範本。 \*

隨附的 SCP 適用於大多數一般用途、但您的需求可能有所不同。

- 請考慮使用最大 IOPS 來控制未知虛擬機器或測試虛擬機器。 \*

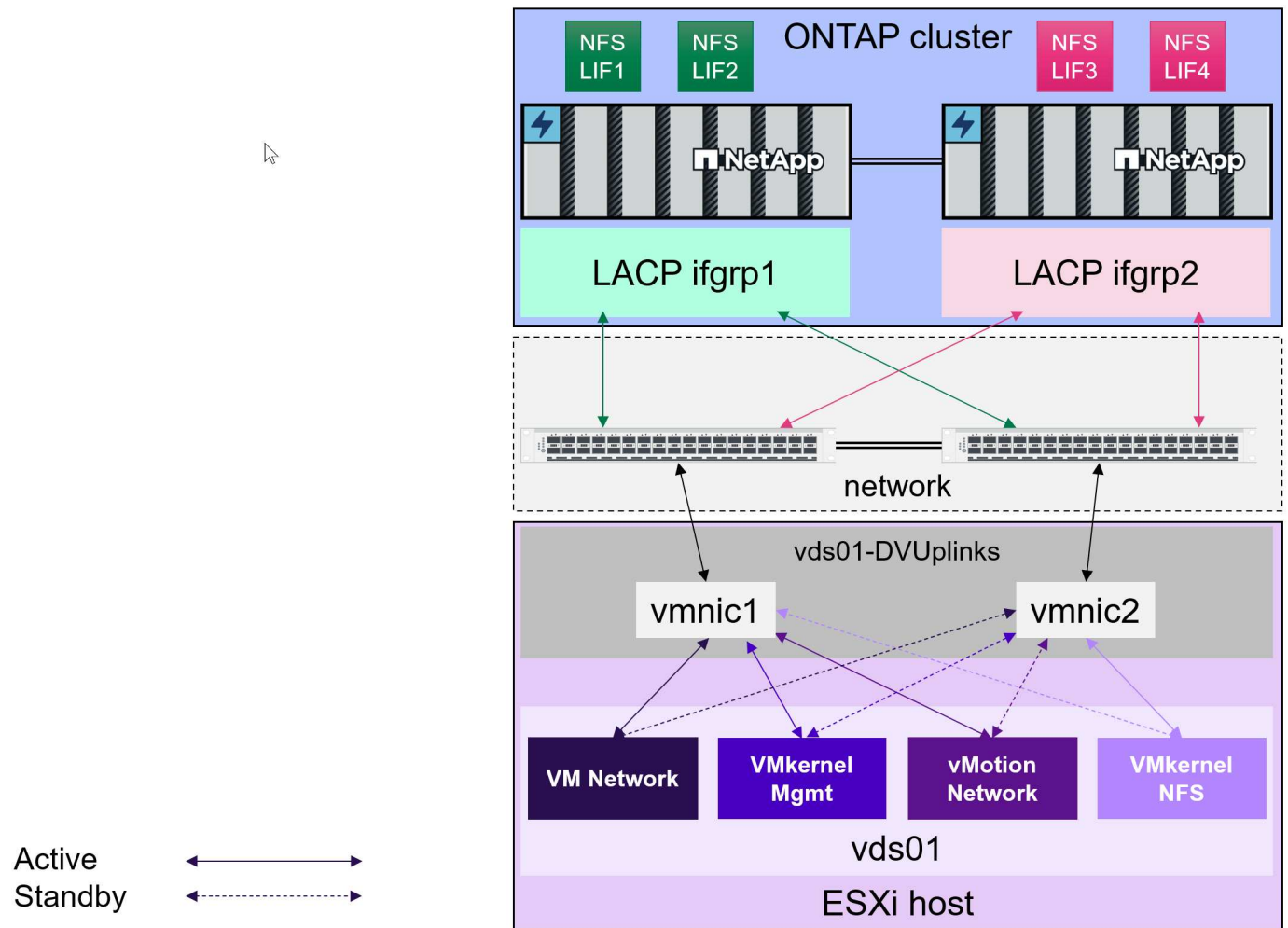
VASA Provider 7.1 首度提供最大 IOPS、可將未知工作負載的 IOPS 限制在特定的 vVol、以避免對其他更重要的工作負載造成影響。如需效能管理的詳細資訊、請參閱表 4。

- 確保您擁有足夠的資料生命。 \*  
每個 HA 配對每個節點至少建立兩個生命期。根據您的工作負載、可能需要更多資源。
- 遵循所有通訊協定最佳實務做法。 \*

請參閱 NetApp 和 VMware 針對您所選傳輸協定的其他最佳實務做法指南。一般而言、除了上述變更之外、沒

有其他變更。

- 使用 vVols over NFS v2\* 的網路組態範例



## 部署 VVols Storage

建立 VM 的 vVols 儲存設備有幾個步驟。

對於使用 ONTAP 用於傳統資料存放區的現有 vSphere 環境、可能不需要前兩個步驟。您可能已經使用 ONTAP 工具來管理、自動化及報告 VMFS 或傳統 NFS 型儲存設備。以下章節將詳細說明這些步驟。

1. 建立儲存虛擬機器 (SVM) 及其傳輸協定組態。您將選擇 NVMe / FC、NFSv3、NFSv4.1、iSCSI、FCP、或是混合使用這些選項。您可以使用 ONTAP 系統管理員精靈或叢集 Shell 命令列。
  - 每個節點至少有一個 LIF 用於每個交換器 / 架構連線。最佳做法是為 FCP、iSCSI 或 NVMe 型傳輸協定、每個節點建立兩個以上的資料傳輸協定。
  - 現在可以建立磁碟區、但讓 `_資源存放區_` 精靈建立磁碟區比較簡單。此規則的唯一例外是、您打算將 vVols 複寫搭配 VMware Site Recovery Manager 使用。使用現有的 FlexVol 磁碟區與現有的 SnapMirror 關係、更容易設定。請注意、不要在任何磁碟區上啟用用於 vVols 的 QoS、因為這是由 SPBM 和 ONTAP 工具管理的。
2. 使用從 NetApp 支援網站下載的 OVA 部署適用於 VMware vSphere 的 ONTAP 工具。
3. 為您的環境設定 ONTAP 工具。

- 將 ONTAP 叢集新增至 *Storage Systems* 下的 ONTAP 工具
    - 雖然 ONTAP 工具和 SRA 同時支援叢集層級和 SVM 層級認證、但 VASA Provider 僅支援儲存系統的叢集層級認證。這是因為許多用於 vVols 的 API 只能在叢集層級使用。因此、如果您打算使用 vVols、則必須使用叢集範圍認證來新增 ONTAP 叢集。
  - 如果 ONTAP 資料生命期與 VMkernel 介面卡位於不同的子網路上、則必須將 VMkernel 介面卡子網路新增至 ONTAP 工具的「設定」功能表中的「所選子網路」清單。根據預設、ONTAP 工具僅允許本機子網路存取、以保護您的儲存流量。
  - ONTAP 工具隨附數個預先定義的原則、這些原則可以使用或查看 [使用原則管理 VM](#) 以取得建立 SCP 的指引。
4. 使用 vCenter 中的 ONTAP tools\_ 功能表來啟動 *Provision datastortory* 精靈。
  5. 提供有意義的名稱、並選取所需的傳輸協定。您也可以提供資料存放區的說明。
  6. 選取一個或多個要由 VVols 資料存放區支援的 SCP。這會篩選出任何無法符合設定檔的 ONTAP 系統。從產生的清單中、選取所需的叢集和 SVM。
  7. 使用精靈為每個指定的 FlexVol 建立新的磁碟區、或選取適當的選項按鈕來使用現有的磁碟區。
  8. 從 vCenter UI 的 *Policies and Profiles* 功能表、為每個要在資料存放區中使用的 SCP 建立 VM 原則。
  9. 選擇「NetApp.Cluster.Data.ONTAP.VP.VVOL」儲存規則集。「NetApp.Cluster.Data.ONTAP.VP.VASA10」儲存規則集適用於對非 VVols 資料存放區的 SPBM 支援
  10. 建立 VM 儲存原則時、您將依名稱指定儲存功能設定檔。在此步驟中、您也可以使用複寫索引標籤來設定 SnapMirror 原則比對、並使用標記索引標籤來設定標籤型比對。請注意、標記必須已建立、才能選取。
  11. 建立 VM、在 Select Storage 下選取 VM Storage Policy 和相容的資料存放區。

### 將 VM 從傳統資料存放區移轉至 vVols

將 VM 從傳統資料存放區移轉至 vVols 資料存放區、就像在傳統資料存放區之間移動 VM 一樣簡單。只要選取虛擬機器、然後從「動作」清單中選取「移轉」、然後選取移轉類型\_僅變更儲存設備\_即可。移轉複本作業將會隨 vSphere 6.0 及更新版本卸載、以便將 SAN VMFS 移轉至 vVols、但不會從 NAS VMDK 移轉至 vVols。

### 使用原則管理 VM

若要利用原則型管理來自動化儲存資源配置、我們需要：

- 使用儲存功能設定檔（FlexVol）定義儲存設備的功能（ONTAP 節點和 Volume）。
- 建立對應至定義的 CDP 的 VM 儲存原則。

NetApp 已從 VASA Provider 7.2 開始簡化功能與對應、並在更新版本中持續改善。本節著重於這種新方法。早期版本支援更多功能、並允許將其個別對應至儲存原則、但不再支援此方法。

### ONTAP 工具版本的儲存功能設定檔功能

* SCP 功能 *	* 能力價值 *	* 支援的版本 *	附註
壓縮	是、否、任何	全部	AFF 在 7.2 及更新版本中為必填項目。
重複資料刪除	是、否、任何	全部	7.2 及更新版本的 AFF 適用的調查。





## Create Storage Capability Profile

1 General

2 Platform

3 Protocol

4 Performance

5 Storage attributes

6 Summary

### General

Specify a name and description for the storage capability profile. ?

Name:

New\_SCP

Description:

CANCEL

NEXT

## Create Storage Capability Profile

1 General

2 Platform

3 Protocol

4 Performance

5 Storage attributes

6 Summary

### Platform

Platform:

All Flash FAS (AFF)

CANCEL

BACK

NEXT

## Create Storage Capability Profile

- 1 General
- 2 Platform
- 3 Protocol**
- 4 Performance
- 5 Storage attributes
- 6 Summary

### Protocol

Protocol:

Any

- Any
- FCP
- NFS
- NFS 4.1
- iSCSI
- NVMe/FC

CANCEL

BACK

NEXT

## Create Storage Capability Profile

- 1 General
- 2 Platform
- 3 Protocol
- 4 Performance**
- 5 Storage attributes
- 6 Summary

### Performance

None ⓘ

QoS policy group ⓘ

Min IOPS: 1000

Max IOPS:

Unlimited

CANCEL

BACK

NEXT

### Create Storage Capability Profile

- 1 General
- 2 Platform
- 3 Protocol
- 4 Performance
- 5 Storage attributes
- 6 Summary

### Storage attributes

Deduplication:	Yes	▼
Compression:	Yes	▼
Space reserve:	Thin	▼
Encryption:	Yes	▼
Tiering policy (FabricPool):	Snapshot	▼

CANCEL
BACK
NEXT

### Create Storage Capability Profile

- 1 General
- 2 Platform
- 3 Protocol
- 4 Performance
- 5 Storage attributes
- 6 Summary

### Summary

Name:	New_SCP	
Description:	N/A	
Platform:	All Flash FAS (AFF)	
Protocol:	Any	
Min IOPS:	1000 IOPS	
Max IOPS:	Unlimited	
Space reserve:	Thin	
Deduplication:	Yes	
Compression:	Yes	
Encryption:	Yes	
Tiering policy (FabricPool):	Snapshot	

CANCEL
BACK
FINISH

• 建立 VVols 資料存放區 \*

建立必要的 SCP 之後、就可以使用它們來建立 vVols 資料存放區（也可以選用資料存放區的 FlexVol 磁碟區）。以滑鼠右鍵按一下您要建立 VVols 資料存放區的主機、叢集或資料中心、然後選取 ONTAP tools\_ > *Provision Datastore*。選取一個或多個要由資料存放區支援的 FlexVol、然後從現有的 FlexVol Volume 和 / 或為資料存放區配置新的 Volume 中進行選取。最後、為資料存放區指定預設的 SCP、用於未依原則指定 SCP 的 VM、以及交換 VVols（這些不需要高效能儲存）。

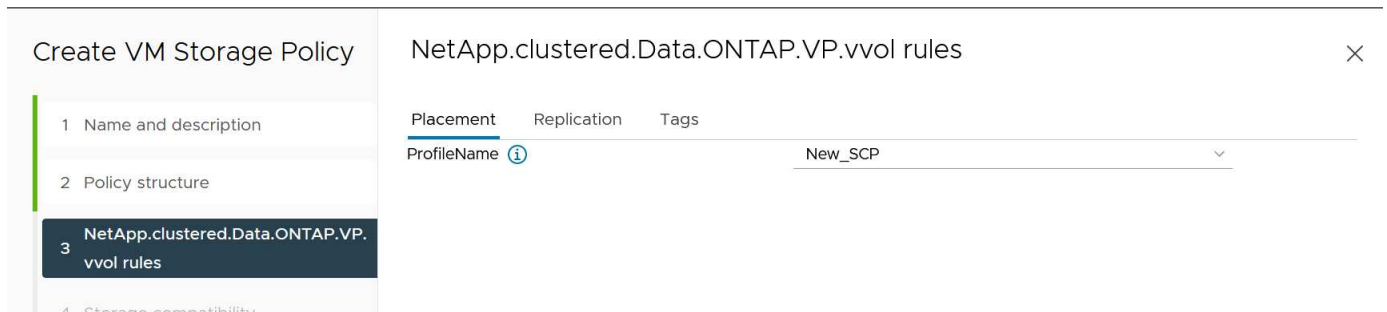
建立 VM 儲存原則

vSphere 中使用 VM 儲存原則來管理儲存 I/O 控制或 vSphere 加密等選用功能。它們也可與 vVols 搭配使用、將特定的儲存功能套用至 VM。使用「NetApp.Cluster.Data.ONTAP.VP.VVOL」儲存類型和「profilename」規則、透過使用原則將特定的 SCP 套用至 VM。請參閱連結：[vmware-vols-ontap.html#Best Practices \[ 透過 NFS v3 使用 vVols 的範例網路組態 \]](#)、以瞭解 ONTAP 工具 VASA Provider 的範例。「NetApp.Cluster.Data.ONTAP.VP.VASA10」儲存設備的規則將用於非 vVols 型資料存放區。

較早的版本類似、但如所述 [ONTAP 工具版本的儲存功能設定檔功能](#)、您的選項可能有所不同。

一旦建立儲存原則、就可以在佈建新 VM 時使用、如所示 "[使用儲存原則部署 VM](#)"。有關搭配 VASA Provider 7.2 使用效能管理功能的準則、請參考 [使用 ONTAP 工具 9.10 及更新版本進行效能管理](#)。

## 使用 ONTAP 工具 VASA Provider 9.10 建立 VM 儲存原則



## 使用 ONTAP 工具 9.10 及更新版本進行效能管理

- ONTAP 工具 9.10 使用自己的平衡放置演算法、將新的 vVol 置於 vVols 資料存放區內的最佳 FlexVol Volume 中。放置方式是根據指定的 SCP 和相符的 FlexVol 磁碟區。如此可確保資料存放區和備份儲存設備符合指定的效能需求。
- 變更效能功能（例如最小和最大 IOPS）需要特別注意特定組態。
  - \* 可以在 SCP 中指定最小和最大 IOPS \*、並在 VM 原則中使用。
    - 變更 SCP 中的 IOPS 不會變更 vVols 上的 QoS、直到編輯 VM 原則、然後重新套用至使用它的 VM 為止（請參閱 [適用於 ONTAP 工具 9.10 及更新版本的儲存功能](#)）。或是使用所需的 IOPS 建立新的 SCP、並變更原則以使用它（並重新套用至 VM）。一般而言、建議您只為不同的服務層級定義個別的 SCP 和 VM 儲存原則、並只需變更 VM 上的 VM 儲存原則即可。
    - AFF 和 FAS 身分具有不同的 IOPs 設定。最小值和最大值均可在 AFF 上使用。不過、非 AFF 系統只能使用最大 IOPs 設定。
- 在某些情況下、可能需要在原則變更後移轉 vVol（手動或由 VASA Provider 和 ONTAP 自動移轉）：
  - 有些變更不需要移轉（例如變更最大 IOPS、可立即套用至 VM、如上所述）。
  - 如果儲存 vVol 的目前 FlexVol Volume 不支援原則變更（例如、平台不支援要求的加密或分層原則）、您將需要在 vCenter 中手動移轉 VM。
- ONTAP 工具會使用目前支援的 ONTAP 版本來建立個別的非共用 QoS 原則。因此、每個個別的 VMDK 都會收到自己的 IOP 分配。

## 重新套用 VM 儲存原則

## VM Storage Policies

CREATE CHECK EDIT CLONE REAPPLY DELETE

Filter

<input type="checkbox"/>	Name	VC
<input type="checkbox"/>	Management Storage Policy - Large	vm-is-vcenter01.vtme.netapp.com
<input type="checkbox"/>	VVol No Requirements Policy	vm-is-vcenter01.vtme.netapp.com
<input type="checkbox"/>	Management Storage Policy - Stretched Lite	vm-is-vcenter01.vtme.netapp.com
<input type="checkbox"/>	VM Encryption Policy	vm-is-vcenter01.vtme.netapp.com
<input type="checkbox"/>	Management Storage policy - Encryption	vm-is-vcenter01.vtme.netapp.com
<input type="checkbox"/>	Management Storage Policy - Single Node	vm-is-vcenter01.vtme.netapp.com
<input type="checkbox"/>	Management Storage policy - Thin	vm-is-vcenter01.vtme.netapp.com
<input checked="" type="checkbox"/>	AFF_ISCSI_VMSP	vm-is-vcenter01.vtme.netapp.com
<input type="checkbox"/>	Host-local PMem Default Storage Policy	vm-is-vcenter01.vtme.netapp.com

1 14 items

## 保護 vVols

以下各節概述將 VMware vVols 與 ONTAP 儲存設備搭配使用的程序和最佳實務做法。

### Vasa Provider 高可用度

NetApp VASA Provider 與 vCenter 外掛程式和 REST API 伺服器（先前稱為虛擬儲存主控台 [VSC]）及儲存複寫介面卡一起執行、是虛擬應用裝置的一部分。如果 VASA Provider 不可用、則使用 vVols 的 VM 將繼續執行。但是、無法建立新的 vVols 資料存放區、而且 VVols 無法由 vSphere 建立或繫結。這表示無法開啟使用 vVols 的虛擬機器、因為 vCenter 將無法要求建立交換 vVol。而且執行中的 VM 無法使用 VMotion 移轉至其他主機、因為 VVols 無法繫結至新主機。

Vasa Provider 7.1 及更新版本支援新功能、確保服務在需要時可用。其中包括監控 VASA Provider 和整合式資料庫服務的新看門狗程序。如果偵測到故障、它會更新記錄檔、然後自動重新啟動服務。

vSphere 管理員必須使用相同的可用度功能來設定進一步的保護、以保護其他關鍵任務 VM 免於軟體、主機硬體和網路中的故障。虛擬應用裝置不需要額外的組態即可使用這些功能、只要使用標準 vSphere 方法進行設定即可。NetApp 已測試並支援這些解決方案。

vSphere High Availability 可輕鬆設定為在發生故障時、在主機叢集中的其他主機上重新啟動 VM。vSphere 容錯功能可建立次要 VM、以持續複寫、並可隨時接管、進而提供更高的可用度。如需這些功能的詳細資訊、請參閱 ["適用於 VMware vSphere 的 ONTAP 工具文件 \(設定 ONTAP 工具的高可用度\)"](#) 以及 VMware vSphere 文件（請在 ESXi 和 vCenter Server 下尋找 vSphere 可用度）。

ONTAP 工具 VASA Provider 會即時自動備份 VVols 組態至託管 ONTAP 系統、其中 VVols 資訊儲存在 FlexVol Volume 中繼資料中。如果 ONTAP 工具應用裝置因任何原因而無法使用、您可以輕鬆快速地部署新的工具並匯入組態。如需 VASA Provider 恢復步驟的詳細資訊、請參閱本知識庫文章：

["如何執行 VASA Provider 災難恢復解決方案指南"](#)

### VVols 複寫

許多 ONTAP 客戶使用 NetApp SnapMirror 將傳統的資料存放區複寫到次要儲存系統、然後在發生災難時使用次

要系統來恢復個別 VM 或整個站台。在大多數情況下、客戶都會使用軟體工具來管理這項作業、例如 NetApp SnapCenter 外掛程式 for VMware vSphere 等備份軟體產品、或是 VMware 的 Site Recovery Manager (搭配 ONTAP 工具中的儲存複寫介面卡) 等災難恢復解決方案。

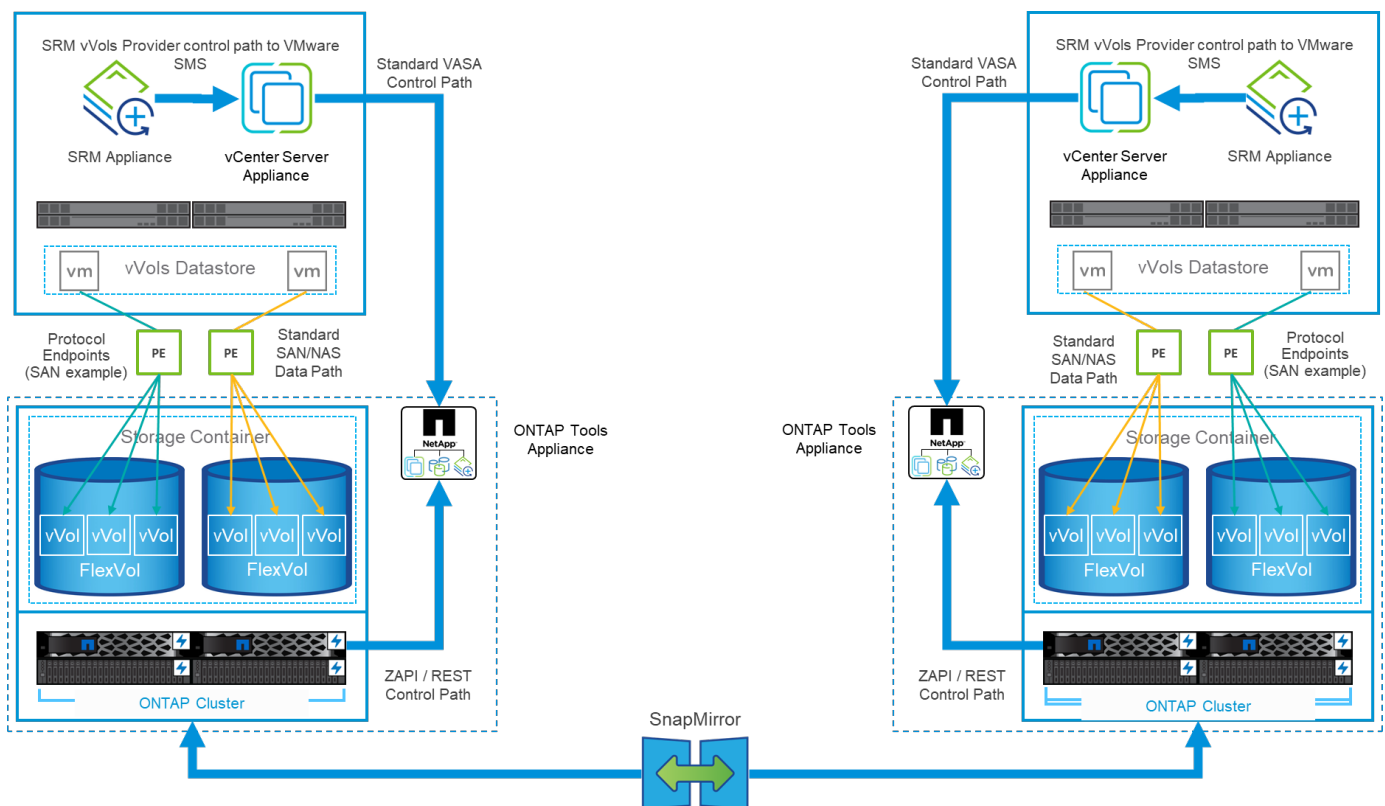
這項軟體工具需求對於管理 vVols 複寫更為重要。雖然有些方面可以由原生功能進行管理 (例如、VMware 託管的 vVols 快照會卸載到 ONTAP、而會使用快速、有效率的檔案或 LUN 複本)、但一般而言、管理複寫和還原需要協調。VVols 的中繼資料受到 ONTAP 和 VASA Provider 的保護、但在次要站台使用時需要額外處理。

ONTAP 工具 9.7.1 搭配 VMware Site Recovery Manager (SRM) 8.3 版本、新增了災難恢復與移轉工作流程協調的支援、充分利用 NetApp SnapMirror 技術。

在 ONTAP 工具 9.7.1 的初始版 SRM 支援中、必須預先建立 FlexVols 並啟用 SnapMirror 保護、才能將它們作為 VVols 資料存放區的備份磁碟區。從 ONTAP 工具 9.10 開始、不再需要此程序。您現在可以將 SnapMirror 保護新增至現有的備份磁碟區、並更新您的 VM 儲存原則、以利用與 SRM 整合的災難恢復、移轉協調和自動化功能、來善用原則型管理。

目前、VMware SRM 是 NetApp 唯一支援的 vVols 災難恢復與移轉自動化解決方案、ONTAP 工具會先檢查是否存在已在 vCenter 註冊的 SRM 8.3 或更新版本伺服器、然後再允許您啟用 vVols 複寫、雖然您可以利用 ONTAP 工具 REST API 來建立自己的服務。

### 使用 SRM 複寫 VVols



### MetroCluster 支援

雖然 ONTAP 工具無法觸發 MetroCluster 切轉、但它確實支援採用統一 vSphere 都會儲存叢集 (VMSC) 組態的 NetApp MetroCluster 系統來支援 VVols 備份磁碟區。MetroCluster 系統的移轉是以正常方式處理。

雖然 NetApp SnapMirror Business Continuity (SM-BC) 也可作為 VMSC 組態的基礎、但 VVols 目前不支援這項功能。

如需 NetApp MetroCluster 的詳細資訊、請參閱以下指南：

["\\_TR-4689 MetroCluster IP 解決方案架構與設計\\_"](#)

["\\_TR-4705 NetApp MetroCluster 解決方案架構與設計\\_"](#)

["VMware KB 2031038 VMware vSphere 支援 NetApp MetroCluster\\_"](#)

## **VVols 備份總覽**

有幾種方法可以保護 VM、例如使用客體內備份代理程式、將 VM 資料檔案附加至備份 Proxy、或使用定義的 API（例如 VMware VADP）。VVols 可以使用相同的機制來保護、許多 NetApp 合作夥伴也支援 VM 備份、包括 vVols。

如前所述、VMware vCenter 託管快照會卸載至節省空間且快速的 ONTAP 檔案 /LUN 複本。這些資料可用於快速手動備份、但 vCenter 最多可限制為 32 個快照。您可以使用 vCenter 拍攝快照、並視需要進行還原。

從 SnapCenter Plugin for VMware vSphere（SCV）4.6 開始、搭配 ONTAP 工具 9.10 及更新版本使用時、新增了對使用 ONTAP FlexVol Volume 快照的 vVols 型虛擬機器的損毀一致備份與還原支援、並支援 SnapMirror 和 SnapVault 複寫。每個磁碟區最多可支援 1023 個快照。此外、選擇控制閥也可以使用 SnapMirror 搭配鏡射儲存庫原則、在次要磁碟區上儲存更多快照、並保留更長的時間。

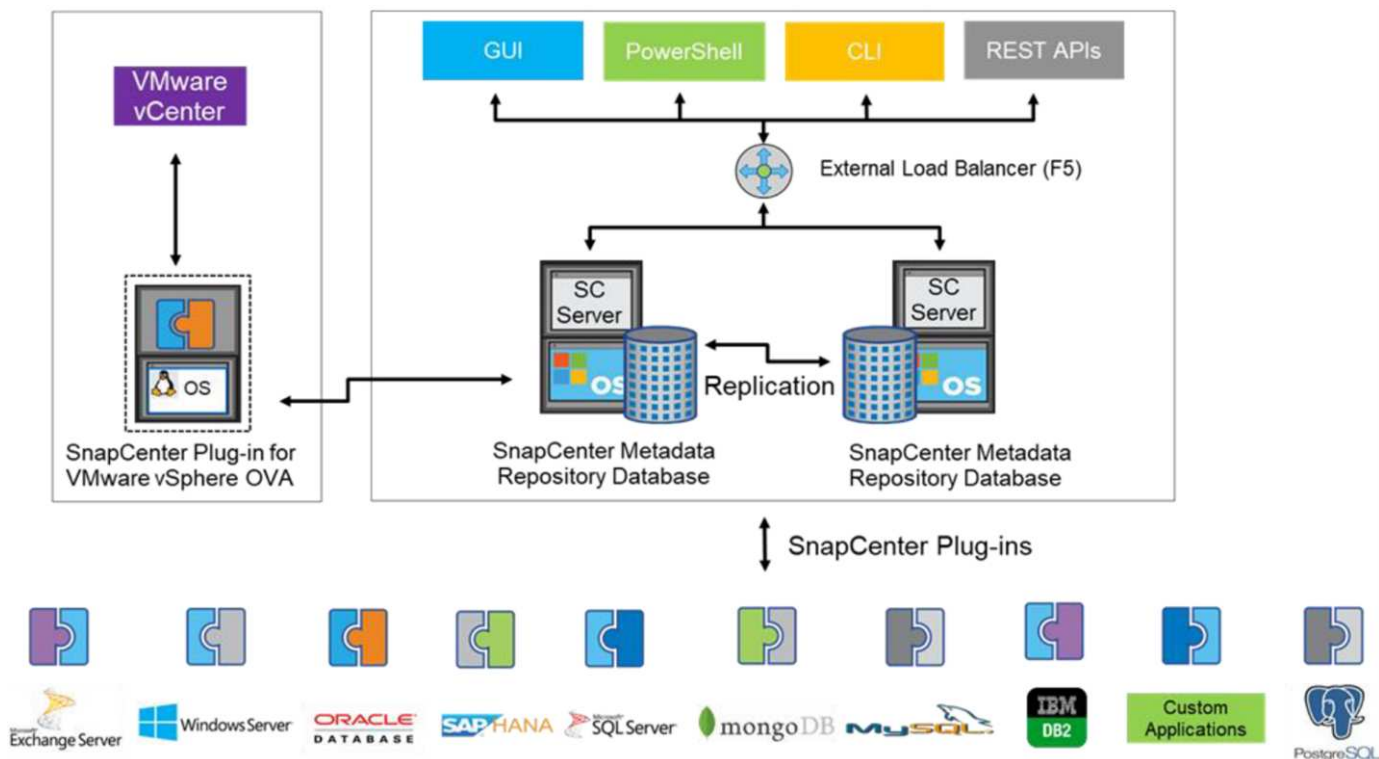
vSphere 8.0 支援是採用分離式本機外掛程式架構的 4.7 號選擇控制閥推出的。vSphere 8.0U1 支援已新增至 4 號選擇控制閥 4.8、該選擇控制閥已完全移轉至新的遠端外掛程式架構。

## **適用於 VMware vSphere 的 VVols Backup 搭配 SnapCenter 外掛程式**

有了 NetApp SnapCenter、您現在可以根據標籤和 / 或資料夾、為 vVols 建立資源群組、以自動利用 ONTAP 的 FlexVol 快照來執行 vVols 型 VM。這可讓您定義備份與還原服務、以便在虛擬機器在環境中動態佈建時、自動保護這些虛擬機器。

適用於 VMware vSphere 的 SnapCenter 外掛程式部署為登錄為 vCenter 擴充功能的獨立應用裝置、可透過 vCenter UI 或 REST API 進行管理、以實現備份與還原服務自動化。

## **架構SnapCenter**



由於其他 SnapCenter 外掛程式在撰寫本文時尚未支援 vVols，因此我們將著重於本文中的獨立部署模型。

由於 SnapCenter 使用 ONTAP FlexVol 快照，因此 vSphere 上不會產生任何額外負荷，也不會像傳統 VM 使用 vCenter 託管快照時所見到的效能損失。此外，由於選擇控制閥的功能是透過 REST API 公開，因此使用 VMware Aria Automation、Ansible、Terraform 等工具，以及幾乎任何其他能夠使用標準 REST API 的自動化工具，都能輕鬆建立自動化工作流程。

如需 SnapCenter 有關靜態 API 的資訊，請參閱 ["REST API 總覽"](#)

如需 SnapCenter VMware vSphere REST API 的靜態外掛程式資訊，請參閱 ["VMware vSphere REST API 的 VMware 外掛程式 SnapCenter"](#)

#### 最佳實務做法

下列最佳實務做法可協助您充分發揮 SnapCenter 部署的效益。

- 選擇控制閥同時支援 vCenter Server RBAC 和 ONTAP RBAC，並包含預先定義的 vCenter 角色，這些角色會在登錄外掛程式時自動為您建立。您可以深入瞭解受支援的 RBAC 類型 ["請按這裡"](#)。
  - 使用 vCenter UI，使用所述的預先定義角色指派最低權限的帳戶存取權 ["請按這裡"](#)。
  - 如果您在 SnapCenter 伺服器上使用選擇控制閥，則必須指派 *SnapCenterAdmin* 角色。
  - ONTAP RBAC 是指用於新增及管理選擇控制閥所使用儲存系統的使用者帳戶。ONTAP RBAC 不適用於 vVols 型備份。深入瞭解 ONTAP RBAC 和選擇控制閥 ["請按這裡"](#)。
- 使用 SnapMirror 將備份資料集複寫到第二個系統，以取得來源磁碟區的完整複本。如前所述，您也可以使用鏡射資料保險箱原則，以長期保留備份資料，而不受來源磁碟區快照保留設定影響。vVols 支援這兩種機制。



- 由於選擇控制閥也需要 VMware vSphere 的 ONTAP 工具才能使用 vVols 功能、因此請務必檢查 NetApp 互通性對照表工具 (IMT) 以瞭解特定版本的相容性
- 如果您在 VMware SRM 中使用 vVols 複寫、請注意您的原則 RPO 和備份排程
- 使用保留設定來設計備份原則、以符合組織定義的還原點目標 (RPO)
- 在資源群組上設定通知設定、以便在執行備份時收到狀態通知 (請參閱下方圖 10)

## 資源群組通知選項

### Edit Resource Group

#### 1. General info & notification

#### 2. Resource

#### 3. Spanning disks

#### 4. Policies

#### 5. Schedules

#### 6. Summary

**vCenter Server:**

**Name:**

**Description:**

**Notification:**

**Email send from:**

**Email send to:**

**Email subject:**

**Latest Snapshot name**  Enable \_recent suffix for latest Snapshot Copy ⓘ

**Custom snapshot format:**  Use custom name format for Snapshot copy

Note that the Plug-in for VMware vSphere cannot do the following:

使用這些文件開始使用選擇控制閥

["深入瞭解 SnapCenter 解 VMware vSphere 的功能"](#)

["部署 SnapCenter VMware vSphere 的 VMware vCenter 外掛程式"](#)

## 疑難排解

有幾種疑難排解資源可提供額外資訊。

## NetApp 支援網站

除了 NetApp 虛擬化產品的各種知識庫文章之外、NetApp 支援網站 也提供方便的登陸頁面 ["VMware vSphere 適用的工具 ONTAP"](#) 產品。此入口網站提供 NetApp 社群文章、下載、技術報告及 VMware 解決方案討論的連

結。您可以在以下網址取得：

["NetApp 支援網站"](#)

其他解決方案文件可於此處取得：

["\\_NetApp 虛擬化解決方案\\_"](#)

### 產品疑難排解

ONTAP 工具의 各種元件（例如 vCenter 外掛程式、VASA Provider 和儲存複寫介面卡）都會記錄在 NetApp 文件儲存庫中。不過、每個都有一個知識庫的個別小節、可能有特定的疑難排解程序。這些解決了 VASA Provider 可能遇到的最常見問題。

### Vasa Provider UI 問題

vCenter vSphere Web Client 偶爾會遇到與 Serenity 元件有關的問題、導致無法顯示 VASA Provider for ONTAP 功能表項目。請參閱部署指南或本知識庫中的解決 VASA Provider 註冊問題 ["文章"](#)。

### VVols 資料存放區資源配置失敗

有時 vCenter 服務在建立 vVols 資料存放區時可能會逾時。若要修正此問題、請重新啟動 VMware-SPS 服務、然後使用 vCenter 功能表（Storage > New Datastore）重新掛載 vVols 資料存放區。本節內容涵蓋《管理指南》中 vCenter Server 6.5 的 VVols 資料存放區資源配置失敗。

### 升級 Unified Appliance 無法掛載 ISO

由於 vCenter 發生錯誤、用於將 Unified Appliance 從一個版本升級至下一個版本的 ISO 可能無法掛載。如果 ISO 能夠附加至 vCenter 中的應用裝置、請遵循此知識庫中的程序 ["文章"](#) 以解決此問題。

## VMware Site Recovery Manager 搭配 ONTAP

### VMware Site Recovery Manager 搭配 ONTAP

ONTAP 自 2002 年引進現代化資料中心以來、一直是 VMware vSphere 環境的領先儲存解決方案、並持續新增創新功能、以簡化管理、同時降低成本。

本文件介紹 VMware 領先業界的災難恢復（DR）軟體 VMware Site Recovery Manager（SRM）ONTAP 解決方案、包括最新的產品資訊和最佳實務做法、可簡化部署、降低風險並簡化後續管理。



本文件取代先前發佈的技術報告 [\\_TR-4900: VMware Site Recovery Manager with ONTAP](#)

最佳實務做法是輔助其他文件、例如指南和相容性工具。這些技術是根據實驗室測試和 NetApp 工程師與客戶廣泛的現場經驗所開發。在某些情況下、建議的最佳實務做法可能不適合您的環境；不過、它們通常是最簡單的解決方案、能滿足大多數客戶的需求。

本文件著重於最新版 ONTAP 9 的功能、這些功能搭配適用於 VMware vSphere 9.12 的 ONTAP 工具（包括 NetApp 儲存複寫介面卡 [SRA] 和 VASA Provider [VP]）、以及 VMware Site Recovery Manager 8.7。

為何ONTAP 要搭配SRM使用此功能？

NetApp資料管理平台採用ONTAP VMware軟體、是SRM最廣泛採用的儲存解決方案之一。理由十分豐富：安全、高性能、統一化的傳輸協定（NAS與SAN一起）資料管理平台、提供業界定義的儲存效率、多租戶、服務控制品質、節省空間的快照資料保護、以及使用SnapMirror進行複寫。所有這些工具都能運用原生混合式多雲端整合技術來保護VMware工作負載、並在彈指之間提供大量的自動化與協調工具。

當您使用SnapMirror進行陣列型複寫時、您將會充分利用ONTAP最成熟的技術之一。SnapMirror可讓您享有安全且高效率的資料傳輸優勢、只複製變更的檔案系統區塊、而非複製整個VM或資料存放區。即使是這些區塊、也能善用空間節約效益、例如重複資料刪除、壓縮及解壓縮等。現代化ONTAP的支援系統現在使用獨立於版本的SnapMirror、讓您能夠靈活選擇來源和目的地叢集。SnapMirror已真正成為最強大的災難恢復工具之一。

無論您使用傳統NFS、iSCSI或光纖通道附加資料存放區（現在支援vVols資料存放區）、SRM都能提供強大的第一方產品、充分發揮ONTAP災難恢復或資料中心移轉規劃與協調的最佳功能。

### SRM如何運用ONTAP VMware®

SRM利用ONTAP VMware ONTAP vSphere的進階資料管理技術與VMware vSphere的VMware vSphere工具整合、此虛擬應用裝置包含三個主要元件：

- vCenter外掛程式先前稱為虛擬儲存主控台（VSC）、可簡化儲存管理與效率功能、增強可用度、並降低儲存成本與營運成本、無論您使用SAN或NAS。它採用最佳實務做法來配置資料存放區、並針對NFS和區塊儲存環境最佳化ESXi主機設定。為獲得所有這些好處、NetApp建議您在使用vSphere搭配執行ONTAP VMware軟體的系統時、使用此外掛程式。
- VASA Provider for ONTAP VMware vStorage API for Storage感知（VASA）架構。VASA Provider將vCenter Server與ONTAP VMware連線、以協助資源配置及監控VM儲存設備。它可支援VMware虛擬磁碟區（vVols）、並管理儲存功能設定檔（包括vVols複寫功能）和個別VM vVols效能。它也提供監控容量和設定檔法規遵循的警示。搭配SRM使用時、VASA Provider ONTAP for VMware可支援vVols型虛擬機器、而不需要在SRM伺服器上安裝SRA介面卡。
- SRA與SRM搭配使用、可管理傳統VMFS與NFS資料存放區的正式作業與災難恢復站台之間的VM資料複寫、也可用於災難恢復複本的不中斷測試。它有助於自動化探索、還原及重新保護等工作。其中包括適用於Windows SRM伺服器和SRM應用裝置的SRA伺服器應用裝置和SRA介面卡。

在SRM伺服器上安裝並設定SRA介面卡、以保護VASA Provider設定中的非vVols資料存放區和/或啟用vVols複寫之後、即可開始設定vSphere環境進行災難恢復。

SRA與VASA Provider提供SRM伺服器的命令與控制介面、可管理ONTAP包含VMware虛擬機器（VM）的VMware FlexVols、以及保護它們的SnapMirror複寫。

從SRM 8.3開始、SRM伺服器引進新的SRM vVols Provider控制路徑、讓它能與vCenter伺服器通訊、並透過該路徑與VASA Provider通訊、而不需要SRA。如此一來、SRM伺服器就能比ONTAP以往更深入地控制這個叢集、因為VASA提供完整的API來進行緊密耦合的整合。

SRM可以使用NetApp專屬的FlexClone技術、在不中斷營運的情況下測試您的災難恢復計畫、在災難恢復站台上為受保護的資料存放區建立近乎即時的複本。SRM會建立沙箱以安全地進行測試、讓組織和客戶在發生真正災難時受到保護、讓您對組織在災難期間執行容錯移轉的能力充滿信心。

如果發生真正的災難、甚至是計畫性的移轉、SRM可讓您透過最終的SnapMirror更新（如果您選擇這樣做）、將任何最後一分鐘的變更傳送至資料集。然後中斷鏡射、並將資料存放區掛載至DR主機。此時、您的VM可以根據預先規劃的策略、以任何順序自動開機。

## SRM搭配ONTAP 功能完善和其他使用案例：混合雲和移轉

與ONTAP 本地儲存選項相比、將SRM部署與豐富的資料管理功能整合、可大幅提升擴充性與效能。除此之外、它還能帶來混合雲的靈活度。混合雲可將未使用的資料區塊從高效能陣列分層、到您偏好的超大規模擴充系統 (FabricPool 例如NetApp Sirse-)、藉此節省成本StorageGRID。您也可以使用ONTAP Select SnapMirror來搭配使用Cloud Volumes ONTAP 以軟體定義的功能、或是使用功能 (CVO) 或雲端型DR的邊緣型系統 "Equinix的NetApp私有儲存設備" 適用於Amazon Web Services (AWS)、Microsoft Azure和Google Cloud Platform (GCP)、可在雲端中建立完全整合的儲存設備、網路和運算服務堆疊。

接著、您可以在雲端服務供應商的資料中心內執行測試容錯移轉、因為 FlexClone 的儲存佔用空間接近零。保護組織的成本現在比以往更低。

SRM也可利用SnapMirror、將VM從一個資料中心有效地傳輸到另一個資料中心、甚至是同一個資料中心內、無論您是自己、或是透過任何數量的NetApp合作夥伴服務供應商、來執行計畫性的移轉作業。

## 部署最佳實務做法

以下各節概述 ONTAP 和 VMware SRM 的部署最佳實務做法。

### 適用於SMT的SVM配置與區隔

利用NetApp技術、儲存虛擬機器 (SVM) 的概念可在安全的多租戶環境中提供嚴格的區隔。ONTAP某個SVM上的SVM使用者無法從另一個SVM存取或管理資源。如此一ONTAP 來、您就能為不同的業務單位建立獨立的SVM、以便在同一個叢集上管理自己的SRM工作流程、進而提升整體儲存效率、進而充分運用此項技術。

考慮ONTAP 使用SVM範圍的帳戶和SVM管理生命體來管理功能、不僅能改善安全控管、也能提升效能。使用SVM範圍的連線時、效能自然會更高、因為SRA不需要處理整個叢集中的所有資源、包括實體資源。而是只需要瞭解抽象化至特定SVM的邏輯資產。

僅使用NAS傳輸協定 (無法存取SAN) 時、您甚至可以設定下列參數來使用新的NAS最佳化模式 (請注意、此名稱是如此、因為SRA和VASA在應用裝置中使用相同的後端服務)：

1. 登入控制面板、網址為 `https://<IP address>:9083` 然後按一下「網路型 CLI 介面」。
2. 執行命令 `vp updateconfig -key=enable.qtree.discovery -value=true`。
3. 執行命令 `vp updateconfig -key=enable.optimised.sra -value=true`。
4. 執行命令 `vp reloadconfig`。

### 部署ONTAP vVols的各種功能與考量

如果您打算搭配vVols使用SRM、則必須使用叢集範圍認證和叢集管理LIF來管理儲存設備。這是因為VASA供應商必須瞭解基礎實體架構、才能滿足VM儲存原則所需的原則。例如、如果您的原則需要All Flash儲存設備、則VASA Provider必須能夠查看哪些系統是Flash。

另一個部署最佳實務做法是、切勿將ONTAP 您的VMware Tools應用裝置儲存在其所管理的vVols資料存放區。這可能會導致您無法開啟VASA供應商的電源、因為您無法為應用裝置建立切換VVOL、因為應用裝置離線。

### 管理ONTAP 功能的最佳實務做法

如前所述、您可以ONTAP 使用叢集或SVM範圍內的認證和管理生命體來管理等叢集。為了達到最佳效能、您可能想要在不使用 vVols 時、考慮使用 SVM 範圍的認證。不過、在這樣做的過程中、您應該瞭解某些需求、而且

您確實會失去某些功能。

- 預設的vsadmin SVM帳戶沒有執行ONTAP 各項功能工作所需的存取層級。因此、您需要建立新的SVM帳戶。
- 如果您使用的是 ONTAP 9.8 或更新版本、NetApp 建議您使用 ONTAP 系統管理員的使用者功能表、以及 ONTAP 工具應用裝置上的 JSON 檔案來建立 RBAC 最低權限使用者帳戶 `https://<IP address>:9083/vsc/config/`。使用您的系統管理員密碼下載Json檔案。這可用於SVM或叢集範圍內的帳戶。

如果您使用ONTAP 的是32個以上版本的版本、則應使用中提供的RBAC使用者建立工具 (RUC) "[NetApp 支援網站工具箱](#)"。

- 由於vCenter UI外掛程式、VASA Provider和SRA伺服器都是完全整合的服務、因此您必須以在vCenter UI中新增儲存設備以供ONTAP 支援VMware工具的相同方式、將儲存設備新增至SRM的SRA介面卡。否則、SRA伺服器可能無法辨識透過SRA介面卡從SRM傳送的要求。
- 使用SVM範圍認證時、不會執行NFS路徑檢查。這是因為實體位置從SVM邏輯上抽象化。不過這並不是令人擔心的問題、因為使用ONTAP 間接路徑時、現代的功能不再受到明顯的效能下降影響。
- 可能不會報告儲存效率所節省的Aggregate空間。
- 在支援的情況下、無法更新負載共用鏡像。
- 可能無法在ONTAP 以SVM範圍認證來管理的各種系統上執行EMS記錄。

## 營運最佳實務做法

以下各節概述 VMware SRM 和 ONTAP 儲存設備的最佳作業實務做法。

### 資料存放區與傳輸協定

- 如有可能、請務必使用ONTAP 「資訊工具」來配置資料存放區和磁碟區。這可確保磁碟區、交會路徑、LUN、igroup、匯出原則、和其他設定均以相容的方式進行設定。
- 當ONTAP 透過SRA使用陣列型複寫時、SRM支援iSCSI、Fibre Channel及NFS版本3 with VMware®9。SRM不支援使用傳統或vVols資料存放區的NFS 4.1版陣列型複寫。
- 若要確認連線能力、請務必確認您可以從目的地ONTAP 叢集掛載並卸載DR站台上的新測試資料存放區。測試您要用於資料存放區連線的每個傳輸協定。最佳實務做法是使用ONTAP 「VMware工具」來建立測試資料存放區、因為它是依照SRM的指示來執行所有資料存放區自動化作業。
- 每個站台的SAN傳輸協定應該是同質的。您可以混合使用NFS和SAN、但SAN傳輸協定不應在站台內混合使用。例如、您可以在站台A中使用FCP、而在站台B中使用iSCSI您不應在站台A同時使用FCP和iSCSI原因是SRA不會在恢復站台建立混合式igroup、而且SRM不會篩選指派給SRA的啟動器清單。
- 先前的指南建議建立 LIF 至資料位置。也就是說、務必使用實體擁有磁碟區的節點上的LIF來掛載資料存放區。這已不再是ONTAP 現今版本的更新要求。只要可能、而且如果提供叢集範圍的認證、ONTAP 工具仍會選擇在資料的本機生命體之間平衡負載、但這並不是高可用度或高效能的需求。
- ONTAP 9 可設定為自動移除快照、以在自動調整大小無法提供足夠的緊急容量時、在空間不足的情況下保留正常運作時間。此功能的預設設定不會自動刪除 SnapMirror 所建立的快照。如果刪除 SnapMirror 快照、則 NetApp 無法針對受影響的磁碟區進行反向和重新同步複寫。若要防止 ONTAP 刪除 SnapMirror 快照、請設定 Snapshot 自動刪除功能以嘗試。

```
snap autodelete modify -volume -commitment try
```

- Volume 自動調整大小應設為 `grow` 適用於包含 SAN 資料存放區和的磁碟區 `grow_shrink` 適用於 NFS 資料存放區。深入瞭解 "設定磁碟區以自動擴充或縮減"。
- 當恢復計畫中的資料存放區數量和保護群組數量減至最少時、SRM 會發揮最佳效能。因此、您應該考慮在受 SRM 保護的環境中最佳化虛擬機器密度、以因應 RTO 至關重要的環境。
- 使用 Distributed Resource Scheduler (DRS) 協助平衡受保護和恢復 ESXi 叢集的負載。請記住、如果您計畫進行容錯回復、當您執行重新保護之前受保護的叢集時、就會成為新的恢復叢集。DRS 有助於平衡兩個方向的放置。
- 儘可能避免將 IP 自訂功能與 SRM 搭配使用、因為這樣可能會增加 RTO。

## 儲存原則型管理 (SPBM) 和 vVols

從 SRM 8.3 開始、支援使用 vVols 資料存放區的 VM 保護。當在 ONTAP 「VMware Tools 設定」功能表中啟用 vVols 複寫時、VASA 供應商會將 SnapMirror 排程公開給 VM 儲存原則、如下列螢幕擷取畫面所示。

以下範例顯示啟用 vVols 複寫。

### Manage Capabilities

- Enable VASA Provider**  
vStorage APIs for Storage Awareness (VASA) is a set of application program interfaces (APIs) that enables vSphere vCenter to recognize the capabilities of storage arrays.
- Enable vVols replication**  
Enables replication of vVols when used with VMware Site Recovery Manager 8.3 or later.
- Enable Storage Replication Adapter (SRA)**  
Storage Replication Adapter (SRA) allows VMware Site Recovery Manager (SRM) to integrate with third party storage array technology.

Enter authentication details for VASA Provider and SRA server:

IP address or hostname: 192.168.64.7  
 Username: Administrator  
 Password: \_\_\_\_\_

CANCEL

APPLY

以下螢幕快照提供「建立VM儲存原則」精靈中所顯示的 SnapMirror 排程範例。

## Create VM Storage Policy

- 1 Name and description
- 2 Policy structure
- 3 NetApp.clustered.Data.ONTAP.VP...
- 4 Storage compatibility
- 5 Review and finish

## NetApp.clustered.Data.ONTAP.VP.vvol rules

Placement   **Replication**   Tags

Disabled  
 Custom

Provider: NetApp.clustered.Data.ONTAP.VP.vvolReplication

Replication ⓘ Asynchronous REMOVE

Replication Schedule ⓘ [Select Value] REMOVE  
[Select Value]  
hourly

CANCEL   BACK   NEXT

支援將故障切換至不同儲存設備的功能。ONTAP例如、系統可能會從ONTAP Select 位於邊緣位置的停止執行、到AFF 核心資料中心的故障轉移。無論儲存設備的相似性為何、您都必須針對啟用複寫的VM儲存原則、設定儲存原則對應和反轉對應、以確保恢復站台所提供的服務符合期望和要求。下列螢幕擷取畫面會強調顯示原則對應範例。

## New Storage Policy Mappings

- 1 Creation mode
- 2 Recovery storage policies
- 3 Reverse mappings
- 4 Ready to complete

## Recovery storage policies

Configure recovery storage policy mappings for one or more storage policies.

Search...

vc1.demo.netapp.com

- Host-local PMem Default Storage Policy
- VC1 Storage Policy \*
- VM Encryption Policy
- vSAN Default Storage Policy
- VVol No Requirements Policy

ADD MAPPINGS

vc1.demo.netapp.com	vc2.demo.netapp.com
<input checked="" type="radio"/> VC1 Storage Policy	<input checked="" type="radio"/> VC2 Storage Policy

1 mapping(s)

CANCEL   BACK   NEXT

## 為vVols資料存放區建立複寫的磁碟區

與先前的vVols資料存放區不同、複寫的vVols資料存放區必須從啟用複寫的開始建立、而且必須使用ONTAP 在具有SnapMirror關係的SnapMirror系統上預先建立的磁碟區。這需要預先設定叢集對等和SVM對等等項目。這些

活動應由 ONTAP 管理員執行、因為這有助於在多個站台之間管理 ONTAP 系統的人員與主要負責 vSphere 作業的人員之間、嚴格區分責任。

這是vSphere管理員的新要求。由於建立的磁碟區超出ONTAP 了功能性測試工具的範圍、因此在ONTAP 定期排程的重新探索期間之前、系統管理員不會察覺到您所做的變更。因此、當您建立要與vVols搭配使用的Volume或SnapMirror關係時、一律執行重新探索是最佳實務做法。只要在主機或叢集上按一下滑鼠右鍵、然後選取「ONTAP 工具」 > 「更新主機和儲存資料」、如下面的螢幕擷取畫面所示。

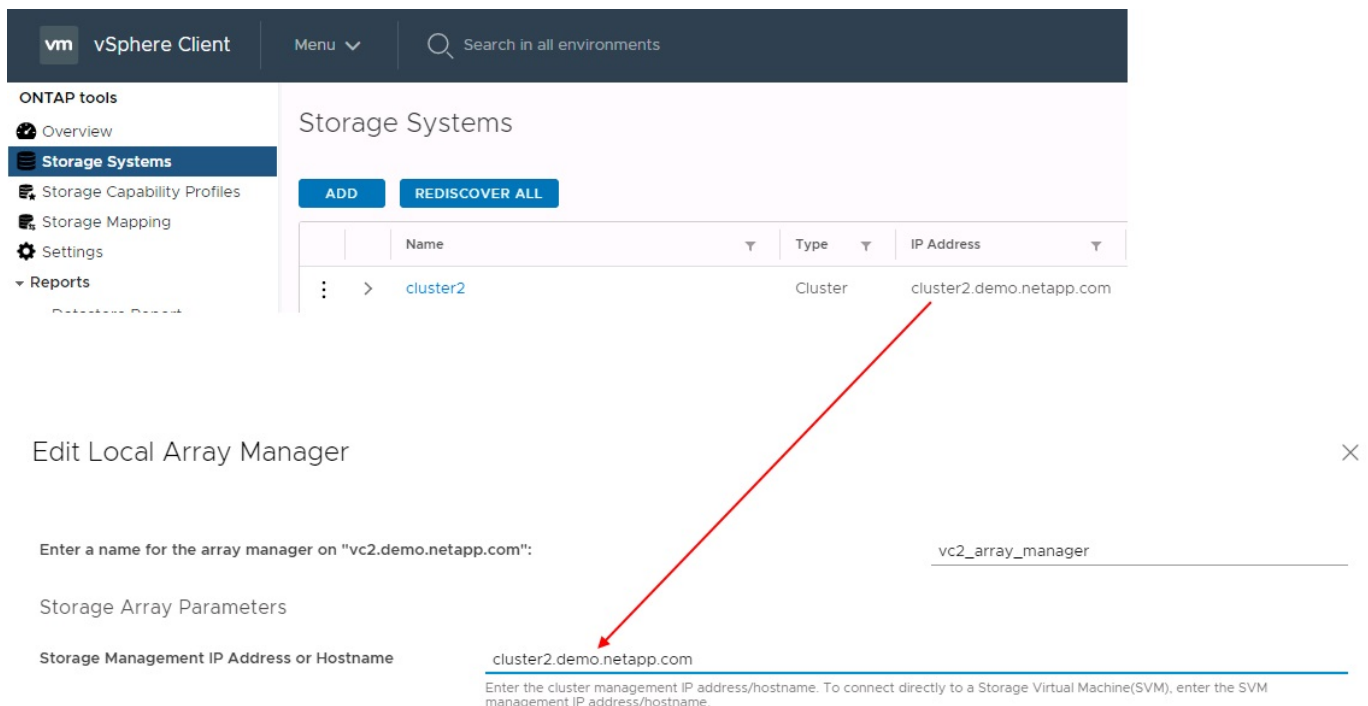


在vVols和SRM方面、請務必謹慎處理。請勿在同一個vVols資料存放區中混用受保護和未受保護的VM。原因是當您使用SRM容錯移轉至DR站台時、只有屬於保護群組的VM才會在DR中上線。因此、當您重新保護（將SnapMirror從災難恢復還原至正式作業）時、可能會覆寫未容錯移轉的VM、並可能包含寶貴的資料。

### 關於陣列配對

系統會為每個陣列配對建立陣列管理程式。有了SRM和ONTAP VMware等工具、每個陣列配對都是以SVM的範圍來完成、即使您使用叢集認證資料也是如此。這可讓您根據指派給租戶的SVM進行管理、在租戶之間分割DR工作流程。您可以為指定的叢集建立多個陣列管理員、而且這些管理員可以是非對稱的。您可以在不同ONTAP的叢集之間進行扇出或扇入。例如、您可以將叢集1上的SVM-A和SVM-B複製到叢集2上的SVM-C、叢集3上的SVM-D、或反之。

在SRM中設定陣列配對時、您應該一律以新增至ONTAP「VMware工具」的相同方式、在SRM中新增這些配對、也就是說、它們必須使用相同的使用者名稱、密碼和管理LIF。這項需求可確保SRA與陣列正常通訊。下列螢幕快照說明ONTAP 叢集在「叢集工具」中的顯示方式、以及如何將其新增至陣列管理程式。





## 關於複寫群組

複寫群組包含一起還原的虛擬機器邏輯集合。這個功能可讓 VASA Provider 自動為您建立複寫群組。ONTAP 由於 SnapMirror 複寫是在磁碟區層級進行、因此一個磁碟區中的所有 VM 都位於相同的複寫群組中。ONTAP

複寫群組的考量因素有多種、以及如何將 VM 分散到 FlexVol 整個流程區。將類似的 VM 分組在同一個磁碟區中、可以提高儲存效率、因為較舊的 ONTAP 系統缺乏 Aggregate 層級的重複資料刪除功能、但群組會增加磁碟區的大小、並減少磁碟區 I/O 並行處理。在現代 ONTAP 系統中、透過在同一個集合體中跨 FlexVol 磁碟區散佈 VM、以達到最佳的效能與儲存效率平衡、進而運用彙總層級的重複資料刪除技術、並在多個磁碟區之間取得更高的 I/O 平行化。您可以將磁碟區中的虛擬機器一起還原、因為保護群組（如下所述）可以包含多個複寫群組。此配置的缺點是、磁碟區 SnapMirror 不會將 Aggregate 重複資料刪除納入考量、因此可能會多次透過線路傳輸區塊。

複寫群組的最後一個考量是、每個群組的本質都是一個邏輯一致性群組（請勿與 SRM 一致性群組混淆）。這是因為磁碟區中的所有 VM 都會使用相同的快照一起傳輸。因此、如果您有必須彼此一致的 VM、請考慮將它們儲存在同 FlexVol 一個地方。

## 關於保護群組

保護群組會將虛擬機器和資料存放區定義為群組、這些群組會從受保護的站台一起還原。受保護站台是指在正常穩定狀態作業期間、保護群組中設定的 VM 存在的位置。請務必注意、雖然 SRM 可能會針對保護群組顯示多個陣列管理程式、但保護群組無法跨越多個陣列管理程式。因此、您不應該跨不同 SVM 上的資料存放區跨 VM 檔案。

## 關於恢復計畫

恢復計畫會定義在相同程序中恢復哪些保護群組。您可以在相同的恢復計畫中設定多個保護群組。此外、若要啟用更多執行恢復計畫的選項、可在多個恢復計畫中加入單一保護群組。

恢復計畫可讓 SRM 管理員定義恢復工作流程、將 VM 指派給優先順序群組、從 1（最高）指派至 5（最低）、預設值為 3（中）。在優先順序群組中、可設定 VM 以因應相依性。

例如、您的公司可能擁有第 1 層關鍵業務應用程式、而該應用程式則仰賴 Microsoft SQL Server 來執行其資料庫。因此、您決定將虛擬機器置於優先順序群組 1。在優先順序群組 1 中、您開始規劃訂單以啟動服務。您可能希望 Microsoft Windows 網域控制器在 Microsoft SQL 伺服器之前開機、而 Microsoft SQL 伺服器必須在應用程式伺服器之前上線、依此類推。您可以將所有這些 VM 新增至優先順序群組、然後設定相依性、因為相依性僅適用於指定的優先順序群組。

NetApp 強烈建議您與應用程式團隊合作、瞭解容錯移轉案例中所需的作業順序、並據此建構您的恢復計畫。

## 測試容錯移轉

最佳實務做法是、只要對受保護的 VM 儲存設備組態進行變更、就必須執行測試容錯移轉。如此可確保在發生災難時、Site Recovery Manager 能夠在預期的 RTO 目標內還原服務。

NetApp 也建議偶爾確認來賓應用程式功能、尤其是在重新設定 VM 儲存設備之後。

執行測試還原作業時、會在 ESXi 主機上為 VM 建立私有測試球型網路。不過、此網路不會自動連線至任何實體網路介面卡、因此無法在 ESXi 主機之間提供連線功能。為了在 DR 測試期間允許在不同 ESXi 主機上執行的 VM 之間進行通訊、會在 DR 站台的 ESXi 主機之間建立實體私有網路。若要驗證測試網路是否為私有網路、可以實體分隔測試網路、或使用 VLAN 或 VLAN 標記來分隔測試網路。此網路必須與正式作業網路隔離、因為在恢復 VM 時、無法將其置於可能與實際正式作業系統衝突的 IP 位址正式作業網路上。在 SRM 中建立恢復計畫時、所建立的測試網路可選取為私有網路、以便在測試期間連接 VM。

在測試通過驗證且不再需要之後、請執行清除作業。執行清除功能會將受保護的VM恢復至初始狀態、並將恢復計畫重設為「就緒」狀態。

## 容錯移轉考量

除了本指南所述的作業順序之外、還有其他幾個考量因素是站台容錯移轉。

您可能必須面對的一個問題是站台之間的網路差異。某些環境可能會在主要站台和DR站台上使用相同的網路IP位址。這項功能稱為「延伸虛擬LAN (VLAN)」或「延伸網路設定」。其他環境可能需要在主要站台使用不同的網路IP位址 (例如不同的VLAN)、相對於DR站台。

VMware提供多種方法來解決此問題。例如VMware NSS-T Data Center等網路虛擬化技術、會從作業環境的第2層到第7層、將整個網路堆疊抽象化、以提供更多可攜的解決方案。深入瞭解 ["支援 SRM 的 NSX-T 選項"](#)。

SRM也可讓您在VM恢復時變更其網路組態。此重新設定包括 IP 位址、閘道位址和 DNS 伺服器設定等設定。不同的網路設定會在個別 VM 恢復時套用到它們、您可以在恢復計畫中的 VM 內容設定中指定。

若要設定SRM將不同的網路設定套用到多個VM、而不需要編輯恢復計畫中每個VM的內容、VMware提供一種稱為DR-IP-customizer的工具。如需瞭解如何使用此公用程式、請參閱 ["VMware 文件"](#)。

## 重新保護

恢復之後、恢復站台將成為新的正式作業站台。由於恢復作業中斷了SnapMirror複寫、因此新的正式作業站台不會受到任何未來災難的保護。最佳實務做法是在恢復後立即將新的正式作業站台保護到另一個站台。如果原始正式作業站台可運作、VMware管理員可以將原始正式作業站台當作新的恢復站台、以保護新正式作業站台、有效反轉保護方向。只有在非災難性故障時、才能使用重新保護功能。因此、原始vCenter Server、ESXi伺服器、SRM伺服器及對應的資料庫最終必須可還原。如果無法使用、則必須建立新的保護群組和新的恢復計畫。

## 容錯回復

容錯回復作業基本上是以不同於以往的方向進行容錯移轉。最佳實務做法是在嘗試容錯回復之前、或是在容錯移轉至原始站台之前、先確認原始站台是否恢復為可接受的功能層級。如果原始站台仍遭入侵、您應該延遲容錯回復、直到故障獲得充分補救為止。

另一個容錯回復最佳做法是在完成重新保護後、在執行最終容錯回復之前、一律執行測試容錯移轉。如此可驗證原始站台上的系統是否能夠完成作業。

## 重新保護原始網站

在容錯回復之後、您應該向所有相關人員確認他們的服務已恢復正常、然後再重新執行「重新保護」、

在容錯回復後執行重新保護、基本上會使環境回到最初的狀態、並再次從正式作業站台執行SnapMirror複寫至還原站台。

## 複寫拓撲

在流程9中ONTAP、叢集管理員可以看到叢集的實體元件、但使用叢集的應用程式和主機無法直接看到這些元件。實體元件提供一個共享資源集區、用於建構邏輯叢集資源。應用程式和主機只能透過含有磁碟區和LIF的SVM存取資料。

在VMware vCenter Site Recovery Manager中、每個NetApp SVM都被視為陣列。SRM支援特定的陣列對陣列 (或SVM對SVM) 複寫配置。

單一VM無法在多個SRM陣列上擁有資料（虛擬機器磁碟（VMDK）或RDM）、原因如下：

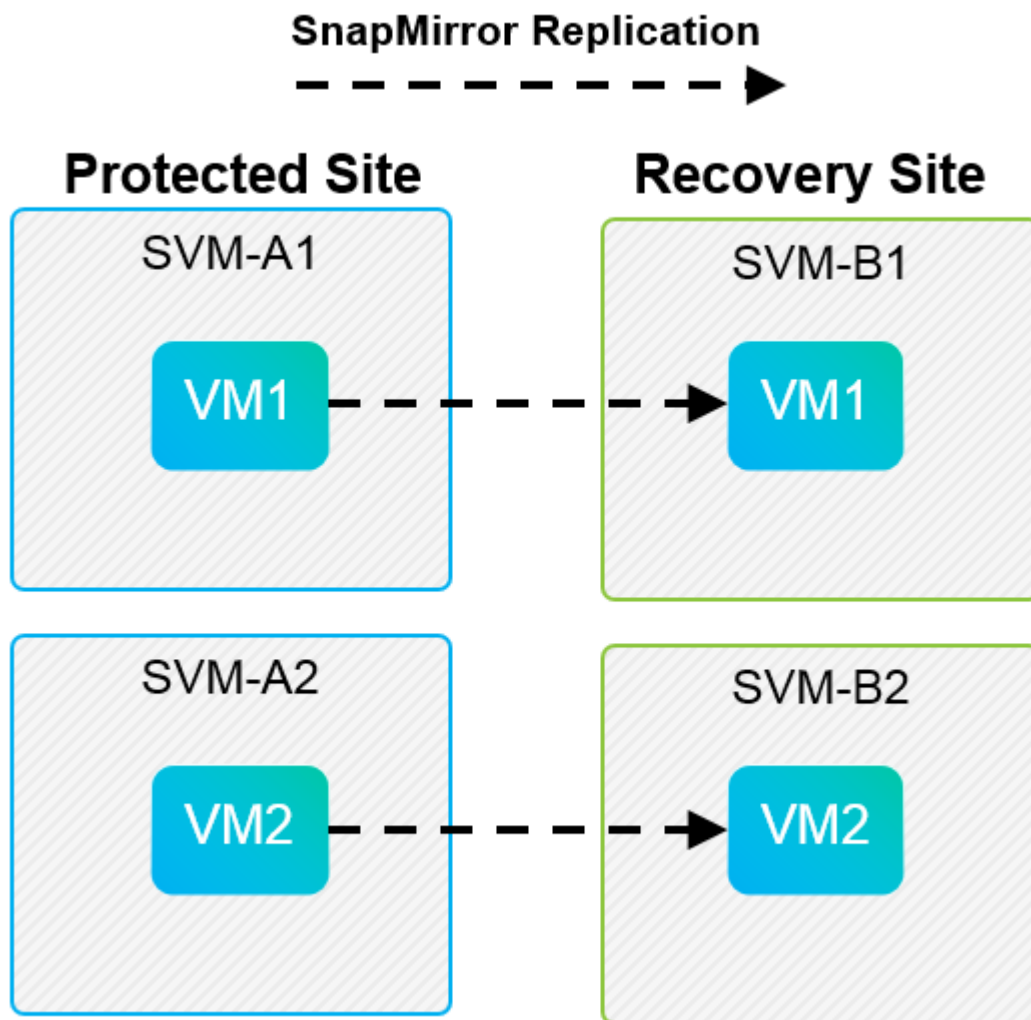
- SRM只能看到SVM、而非個別的實體控制器。
- SVM可控制橫跨叢集中多個節點的LUN和磁碟區。

#### 最佳實務做法

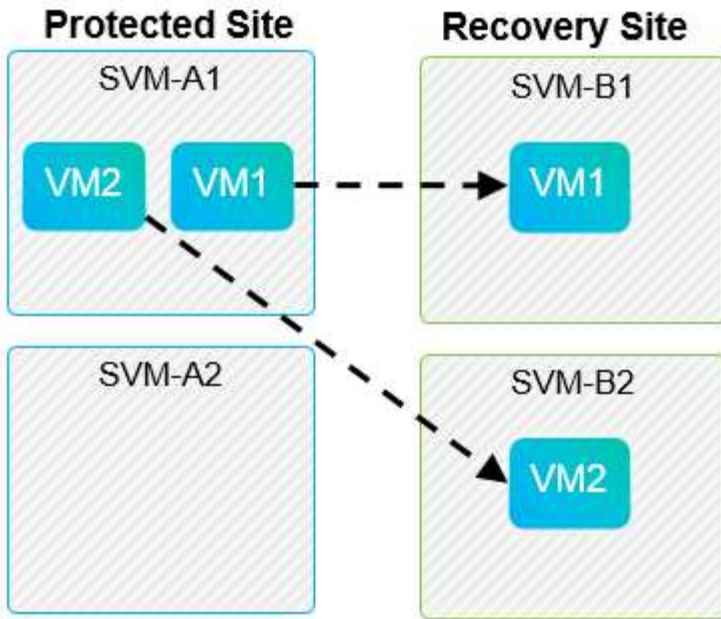
若要判斷可支援性、請謹記此規則：若要使用SRM和NetApp SRA來保護VM、VM的所有部分都必須只存在於一個SVM上。此規則同時適用於受保護的站台和恢復站台。

#### 支援的SnapMirror配置

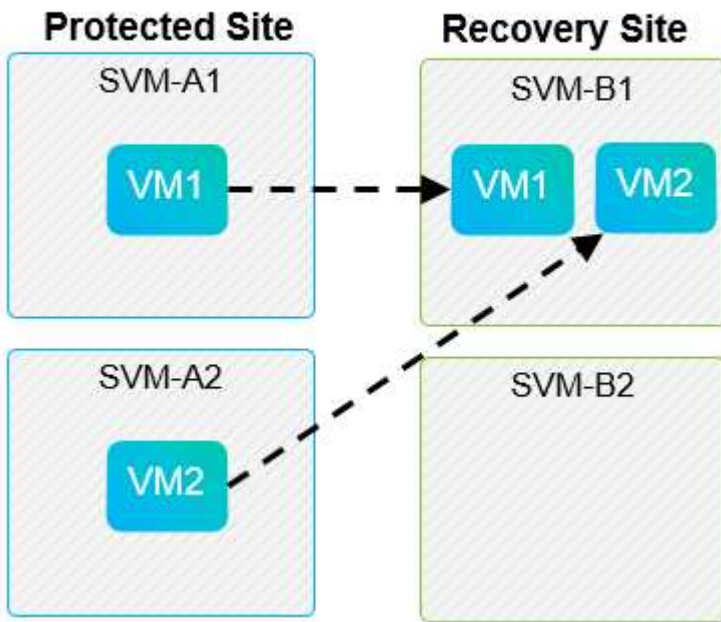
下圖顯示了SRM和SRA支援的SnapMirror關係配置案例。複寫磁碟區中的每個VM在每個站台只擁有一個SRM陣列（SVM）上的資料。

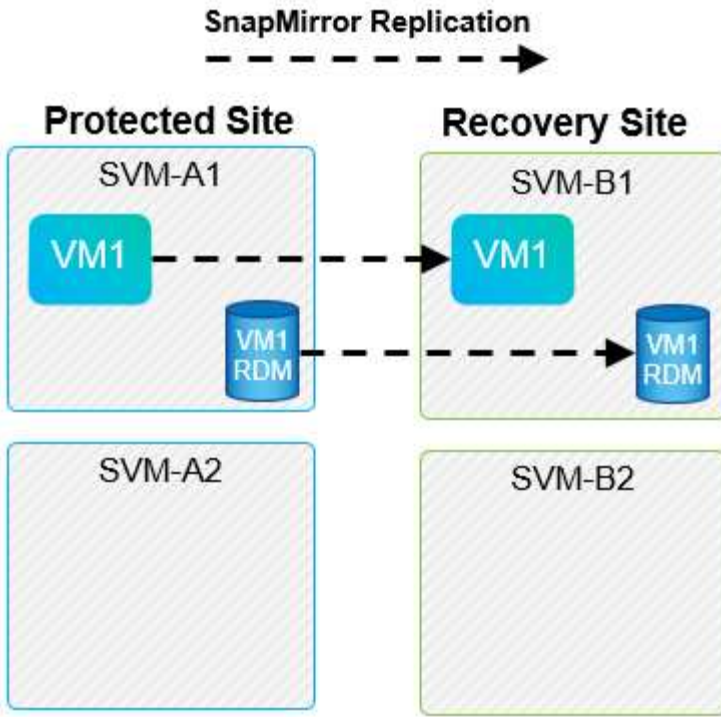


SnapMirror Replication



SnapMirror Replication





### 支援的Array Manager配置

當您在SRM中使用陣列型複寫（ABR）時、保護群組會隔離為單一陣列配對、如下面的快照所示。在此案例中、SVM1 和 SVM2 與我們合作 SVM3 和 SVM4 在恢復站點上。不過、您只能在建立保護群組時、從兩個陣列配對中選取一個。

#### New Protection Group

- 1 Name and direction
- 2 Type
- 3 Datastore groups
- 4 Recovery plan
- 5 Ready to complete

#### Type ×

Select the type of protection group you want to create:

- Datastore groups (array-based replication)**  
Protect all virtual machines which are on specific datastores.
- Individual VMs (vSphere Replication)  
Protect specific virtual machines, regardless of the datastores.
- Virtual Volumes (vVol replication)  
Protect virtual machines which are on replicated vVol storage.
- Storage policies (array-based replication)  
Protect virtual machines with specific storage policies.

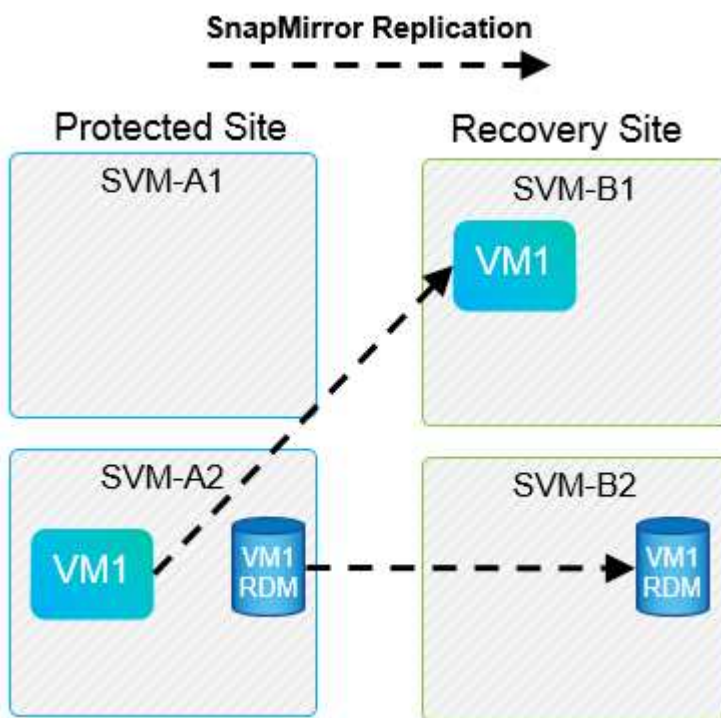
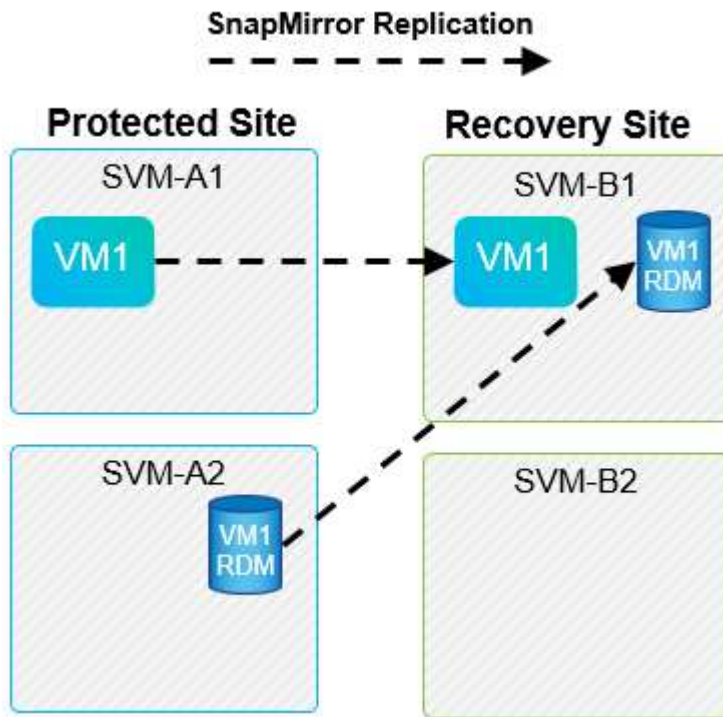
Select array pair

	Array Pair	Array Manager Pair
<input type="radio"/>	✓ cluster1:svm1 ↔ cluster2:svm2	vc1 array manager ↔ vc2 array manager
<input type="radio"/>	✓ cluster1:svm3 ↔ cluster2:svm4	vc1 trad datastores ↔ vc2 trad datastores

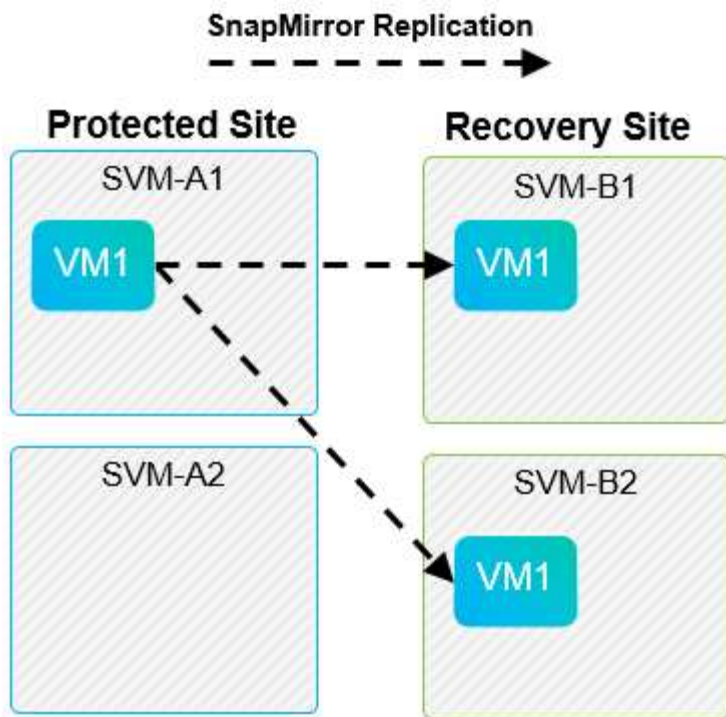
CANCEL
BACK
NEXT

## 不支援的配置

不受支援的組態在個別VM擁有的多個SVM上有資料（VMDK或RDM）。在下列圖中所示的範例中、vm1 無法使用 SRM 進行保護設定、原因是 vm1 在兩個 SVM 上有資料。

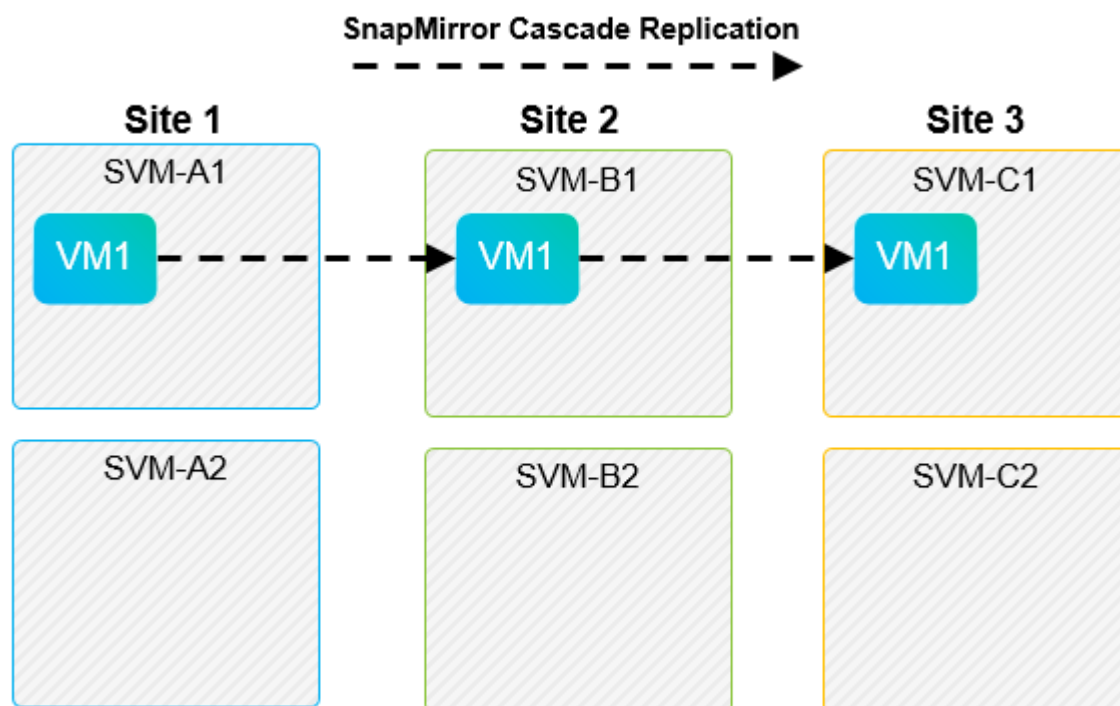


任何將個別NetApp磁碟區從一個來源SVM複寫到同一個SVM或不同SVM中的多個目的地的複寫關係、都稱為SnapMirror連出。SRM不支援連出。在下圖所示範例中、vm1 無法在 SRM 中進行保護設定、因為它會與 SnapMirror 一起複寫到兩個不同的位置。



### SnapMirror串聯

SRM不支援SnapMirror關係的串聯、在這種關係中、來源磁碟區會複寫到目的地磁碟區、而目的地磁碟區也會使用SnapMirror複寫到另一個目的地磁碟區。在下圖所示的案例中、SRM無法用於任何站台之間的容錯移轉。



### SnapMirror與SnapVault

NetApp SnapVault 解決方案軟體可在NetApp儲存系統之間、以磁碟形式備份企業資料。可在同一個環境中共存的VMware vCenter和SnapMirror、不過SRM僅支援SnapMirror關係的容錯移轉。SnapVault



## NetApp 支援 mirror-vault 原則類型。

為了執行效能提升8.2、從一開始就重建了這個系統。SnapVault ONTAP儘管以前Data ONTAP 的使用者應該會發現相似點、SnapVault 但本版的VMware已經做出重大的改善。其中一項重大進展是SnapVault、能夠在傳輸過程中維持主要資料的儲存效率。

架構上的一項重要變更是SnapVault、在ONTAP Volume層級進行的不只是qtree層級的不完整複寫、7-Mode SnapVault 的情況就是如此。這項設定表示SnapVault、來源的不景點必須是一個Volume、而且該Volume必須複寫到SnapVault 自己的Volume上的不二系統。

在使用 SnapVault 的環境中、會在主要儲存系統上建立特別命名的快照。根據實作的組態而定、命名快照可由 SnapVault 排程或 NetApp Active IQ Unified Manager 等應用程式在主要系統上建立。然後，在主系統上創建的命名快照將被複製到 SnapMirror 目標，並從該目的地將其保存到 SnapVault 目的地。

您可以在串聯組態中建立來源Volume、將磁碟區複寫到DR站台的SnapMirror目的地、然後從該磁碟區保存到SnapVault 一個目的地。來源Volume也可建立在連出關係中、其中一個目的地是SnapMirror目的地、另一個目的地SnapVault 是一個目的地。不過、SRA不會在SnapVault 發生SRM容錯移轉或複寫反轉時、自動重新設定「還原」關係、以使用SnapMirror目的地Volume作為資料庫的來源。

如需SnapMirror和SnapVault 適用於ONTAP SnapMirror的更新資訊、請參閱 ["TR-4015 《SnapMirror組態最佳實務指南ONTAP》 \(英文\) 。](#)"

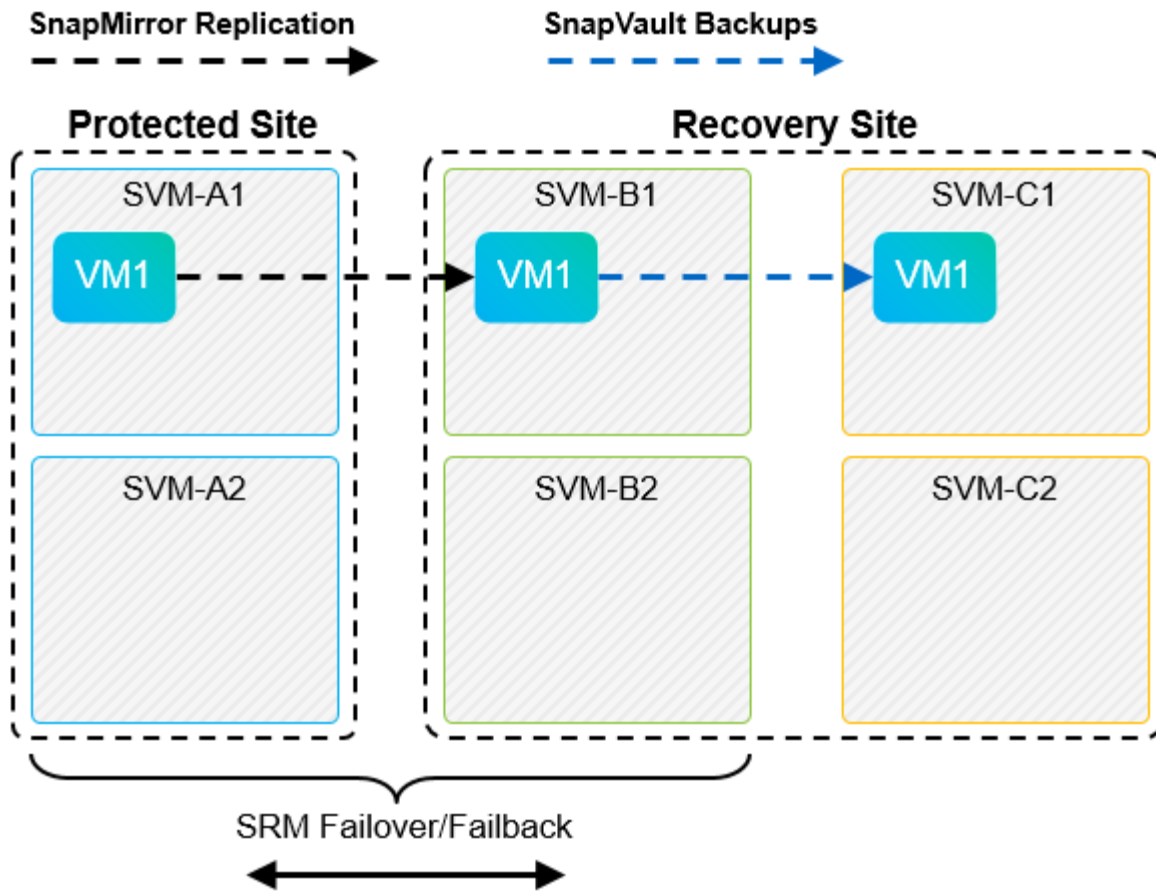
### 最佳實務做法

如果SnapVault 在同一個環境中使用了VMware vCenter和SRM、NetApp建議使用SnapMirror SnapVault 來進行還原串聯組態、SnapVault 以便從DR站台的SnapMirror目的地執行還原備份。發生災難時、此組態會使主要站台無法存取。將SnapVault 還原目的地保留在恢復站台、可在SnapVault 容錯移轉後重新設定還原功能、SnapVault 以便在恢復站台上操作時繼續執行還原備份。

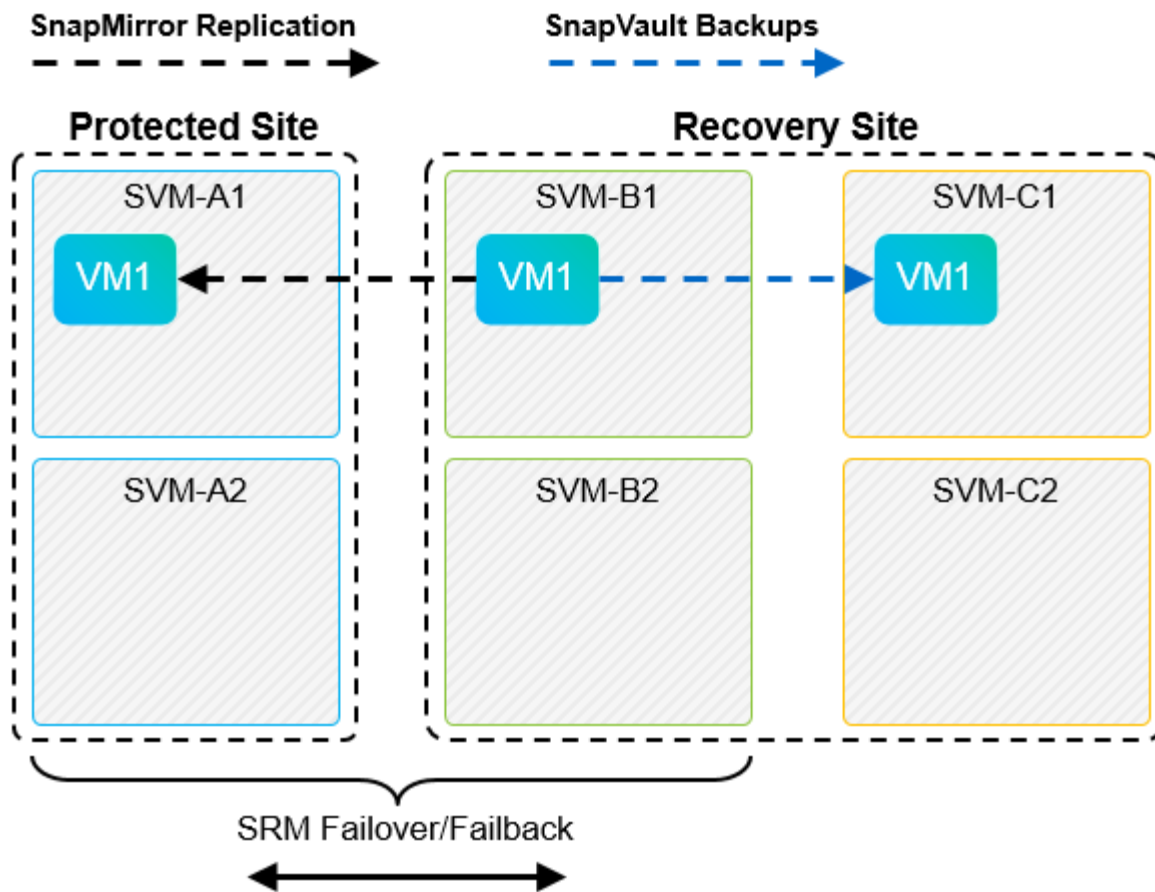
在VMware環境中、每個資料存放區都有通用唯一識別碼 (UUID)、而且每個VM都有唯一的託管物件ID (MOID)。在容錯移轉或容錯回復期間、SRM不會維護這些ID。由於SRM在容錯移轉期間不會維護資料存放區UUID和VM MOID、因此在SRM容錯移轉之後、任何依賴這些ID的應用程式都必須重新設定。例如NetApp Active IQ Unified Manager 解決方案就是應用程式、它可協調SnapVault vSphere環境中的功能複寫。

下圖說明SnapMirror至SnapVault SnapMirror串聯組態。如果該站台位於DR站台或第三站台、但不受主站台中斷影響、則可重新設定環境、以便在容錯移轉後繼續備份。SnapVault





下圖說明使用SRM將SnapMirror複寫還原回主要站台之後的組態。環境也經過重新設定、SnapVault 使目前的SnapMirror來源產生了不支援的資料。此設定為SnapMirror SnapVault 的橫向風扇組態。



在SRM執行容錯回復並第二次反轉SnapMirror關係之後、正式作業資料就會回到主要站台。此資料現在的保護方式與容錯移轉至DR站台之前相同、透過SnapMirror和SnapVault 還原備份。

### 在Site Recovery Manager環境中使用qtree

qtree是允許應用NAS檔案系統配額的特殊目錄。利用SnapMirror複寫的磁碟區中、能夠建立qtree和qtree。ONTAP不過、SnapMirror不允許複寫個別qtree或qtree層級的複寫。所有SnapMirror複寫僅位於磁碟區層級。因此、NetApp不建議搭配SRM使用qtree。

### 混合式FC與iSCSI環境

藉由支援的SAN傳輸協定 (FC、FCoE和iSCSI) ONTAP、支援的LUN服務、也就是能夠建立LUN並將其對應至連接的主機。由於叢集由多個控制器組成、因此有多個邏輯路徑是由多重路徑I/O管理、可通往任何個別LUN。主機上使用非對稱邏輯單元存取 (ALUA)、以便選取LUN的最佳化路徑、並使其成為資料傳輸的作用中路徑。如果任何LUN的最佳化路徑有所變更 (例如、因為包含的磁碟區已移動)、ONTAP 則針對此變更、支援不中斷地自動辨識及調整。如果最佳化路徑無法使用、ONTAP 則不中斷營運地切換至任何其他可用路徑。

VMware SRM和NetApp SRA支援在一個站台使用FC傳輸協定、在另一個站台使用iSCSI傳輸協定。不過、它不支援在同一個ESXi主機或同一個叢集中的不同主機上混合使用FC附加資料存放區和iSCSI附加資料存放區。SRM不支援此組態、因為在SRM容錯移轉或測試容錯移轉期間、SRM會在要求中包含ESXi主機中的所有FC和iSCSI啟動器。

## 最佳實務做法

SRM和SRA支援受保護站台與恢復站台之間的混合FC和iSCSI傳輸協定。不過、每個站台只能設定一個FC或iSCSI傳輸協定、而非在同一個站台設定兩個傳輸協定。如果要求在同一個站台同時設定FC和iSCSI傳輸協定、NetApp建議某些主機使用iSCSI、而其他主機則使用FC。在此情況下、NetApp也建議設定SRM資源對應、以便將VM設定為容錯移轉至一組主機或另一組主機。

## 使用vVols複寫時疑難排解SRM

使用vVols複寫時、SRM內部的工作流程與SRA和傳統資料存放區使用的工作流程大不相同。例如、沒有Array Manager概念。因此、discoverarrays 和 discoverdevices 從未見過命令。

疑難排解時、瞭解下列新工作流程會有所助益：

1. 查詢複製對等方：探索兩個故障網域之間的複寫合約。
2. 查詢FaultDomain：探索故障網域階層。
3. 查詢複製群組：探索來源或目標網域中的複寫群組。
4. SyncReplicationGroup：在來源與目標之間同步資料。
5. 查詢點時間複本：探索目標上的時間點複本。
6. testFailoverReplicationGroupStart：開始測試容錯移轉。
7. testFailoverReplicationGroupStop：結束測試容錯移轉。
8. 促銷複製群組：將目前正在測試的群組推廣至正式作業。
9. PrepareFailoverReplicationGroup：準備災難恢復。
10. 容錯移轉複製群組：執行災難恢復。
11. 混響複寫群組：啟動反轉複寫。
12. queryMatchingContainer：尋找容器（連同主機或複寫群組）、以特定原則來滿足資源配置要求。
13. 查詢資源中繼資料：從VASA提供者探索所有資源的中繼資料、可傳回資源使用率做為查詢配對Container功能的答案。

設定vVols複寫時最常見的錯誤是無法發現SnapMirror關係。這是因為磁碟區和SnapMirror關係是在ONTAP 不屬於「需求工具」範圍的情況下建立。因此、最佳實務做法是在ONTAP 嘗試建立複寫的vVols資料存放區之前、務必確認SnapMirror關係已完全初始化、並在兩個站台上執行「ReDiscovery工具」中的重新探索。

## 其他資訊

若要深入瞭解本文所述資訊、請檢閱下列文件和 / 或網站：

- TR-4597：VMware vSphere ONTAP for VMware  
["https://docs.netapp.com/us-en/ontap-apps-dbs/vmware/vmware-vsphere-overview.html"](https://docs.netapp.com/us-en/ontap-apps-dbs/vmware/vmware-vsphere-overview.html)
- TR-4400: VMware vSphere 虛擬 Volume ONTAP with VMware  
["https://docs.netapp.com/us-en/ontap-apps-dbs/vmware/vmware-vmvols-overview.html"](https://docs.netapp.com/us-en/ontap-apps-dbs/vmware/vmware-vmvols-overview.html)
- TR-4015 《SnapMirror組態最佳實務指南ONTAP》（英文）  
<https://www.netapp.com/media/17229-tr4015.pdf?v=127202175503P>

- RBAC使用者建立工具ONTAP 以供參考  
["https://mysupport.netapp.com/site/tools/tool-eula/rbac"](https://mysupport.netapp.com/site/tools/tool-eula/rbac)
- VMware vSphere資源的相關工具ONTAP  
["https://mysupport.netapp.com/site/products/all/details/otv/docsandkb-tab"](https://mysupport.netapp.com/site/products/all/details/otv/docsandkb-tab)
- VMware Site Recovery Manager文件  
["https://docs.vmware.com/en/Site-Recovery-Manager/index.html"](https://docs.vmware.com/en/Site-Recovery-Manager/index.html)

請參閱 ["互通性對照表工具IMT \(不含\)"](#) 在 NetApp 支援網站上，驗證您的特有環境是否支援本文件中所述的明確產品與功能版本。NetApp IMT 解決方案定義了可用於建構NetApp支援組態的產品元件和版本。具體結果取決於每位客戶依照已發佈規格所安裝的產品。

## vSphere Metro Storage 叢集搭配 ONTAP

### vSphere Metro Storage 叢集搭配 ONTAP

VMware 領先業界的 vSphere Hypervisor 可部署為稱為 vSphere Metro Storage Cluster (VMSC) 的延伸叢集。

NetApp® MetroCluster™ 和 SnapMirror 主動同步 (以前稱為 SnapMirror 業務連續性或 SMBC) 均支持 VMSC 解決方案，如果一個或多個故障域發生整體中斷，則可提供高級業務連續性。不同故障模式的恢復能力取決於您選擇的組態選項。

適用於 **vSphere** 環境的持續可用度解決方案

ONTAP 架構是靈活且可擴充的儲存平台、可為資料存放區提供 SAN (FCP、iSCSI 和 NVMe of) 和 NAS (NFS v3 和 v4.1) 服務。NetApp AFF、ASA 和 FAS 儲存系統使用 ONTAP 作業系統來提供額外的通訊協定、以供 S3 和 SMB/CIFS 等來賓儲存設備存取。

NetApp MetroCluster 使用 NetApp 的 HA (控制器容錯移轉或 CFO) 功能來防範控制器故障。它還包括本機 SyncMirror 技術、災難時的叢集容錯移轉 (隨需控制器容錯移轉或 CFOD)、硬體備援、以及地理區隔、以達到高可用性。SyncMirror 會將資料寫入兩個叢中、以同步鏡射 MetroCluster 組態的兩個部份資料：本機叢 (位於本機櫃上) 主動提供資料、而遠端叢 (位於遠端機櫃上) 通常不會提供資料。所有 MetroCluster 元件 (例如控制器、儲存設備、纜線、交換器 (與 Fabric MetroCluster 搭配使用) 和介面卡) 均具備硬體備援功能。

NetApp SnapMirror 主動式同步可透過 FCP 和 iSCSI SAN 傳輸協定提供資料存放區精細保護、讓您只能選擇性地保護高優先順序的工作負載。它提供對本機和遠端站台的主動式存取、而 NetApp MetroCluster 則是主動式待命解決方案。目前、主動式同步是一種非對稱式解決方案、其中一端較另一端更偏好、提供更好的效能。這是使用 ALUA (非對稱邏輯單元存取) 功能來達成的、此功能會自動通知 ESXi 主機偏好的控制器。不過、NetApp 已宣佈啟用主動式同步功能、即將啟用完全對稱的存取。

若要跨兩個站台建立 VMware HA/DRS 叢集、ESXi 主機會由 vCenter Server Appliance (VCSA) 使用和管理。vSphere 管理、VMotion® 和虛擬機器網路是透過兩個站台之間的備援網路連線。管理 HA/DRS 叢集的 vCenter Server 可連線至兩個站台的 ESXi 主機、並應使用 vCenter HA 進行設定。

請參閱 ["如何在 vSphere Client 中建立和設定叢集"](#) 設定 vCenter HA。

您也應該參閱 ["VMware vSphere Metro儲存叢集建議實務做法"](#)。

## 什麼是 vSphere Metro Storage Cluster ？

vSphere Metro Storage Cluster (VMSC) 是經過認證的組態、可保護虛擬機器 (VM) 和容器免於故障。這是透過使用延伸儲存概念和 ESXi 主機叢集來達成的、這些主機分佈在不同的故障網域、例如機架、建築物、校園或甚至城市。NetApp MetroCluster 和 SnapMirror 主動同步儲存技術可分別為主機叢集提供 RPO = 0 或近乎 RPO = 0 的保護。VMSC 組態的設計是為了確保即使完整的實體或邏輯「站台」故障、資料仍可隨時使用。在成功通過 VMSC 認證程序之後、必須通過 VMSC 組態一部分的儲存裝置認證。所有支援的儲存裝置都可以在中找到 "[VMware 儲存相容性指南](#)"。

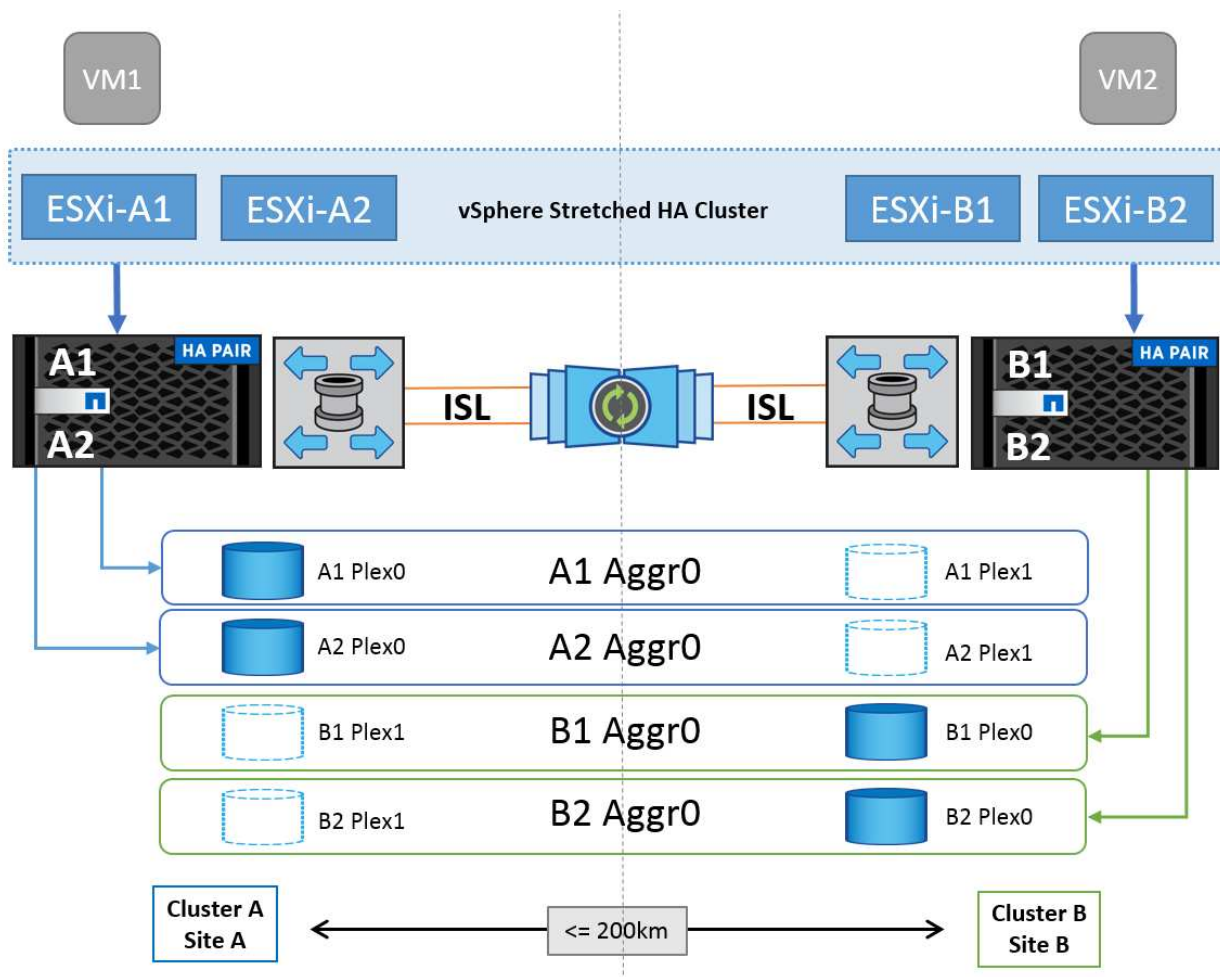
如果您想要更多有關 vSphere Metro Storage 叢集設計準則的資訊、請參閱下列文件：

- "[VMware vSphere 支援 NetApp MetroCluster](#)"
- "[VMware vSphere 支援 NetApp SnapMirror 業務持續運作](#)" (現在稱為 SnapMirror 主動同步)

視延遲考量因素而定、NetApp MetroCluster 可部署在兩種不同的組態中、以搭配 vSphere 使用：

- Stretch MetroCluster
- Fabric MetroCluster

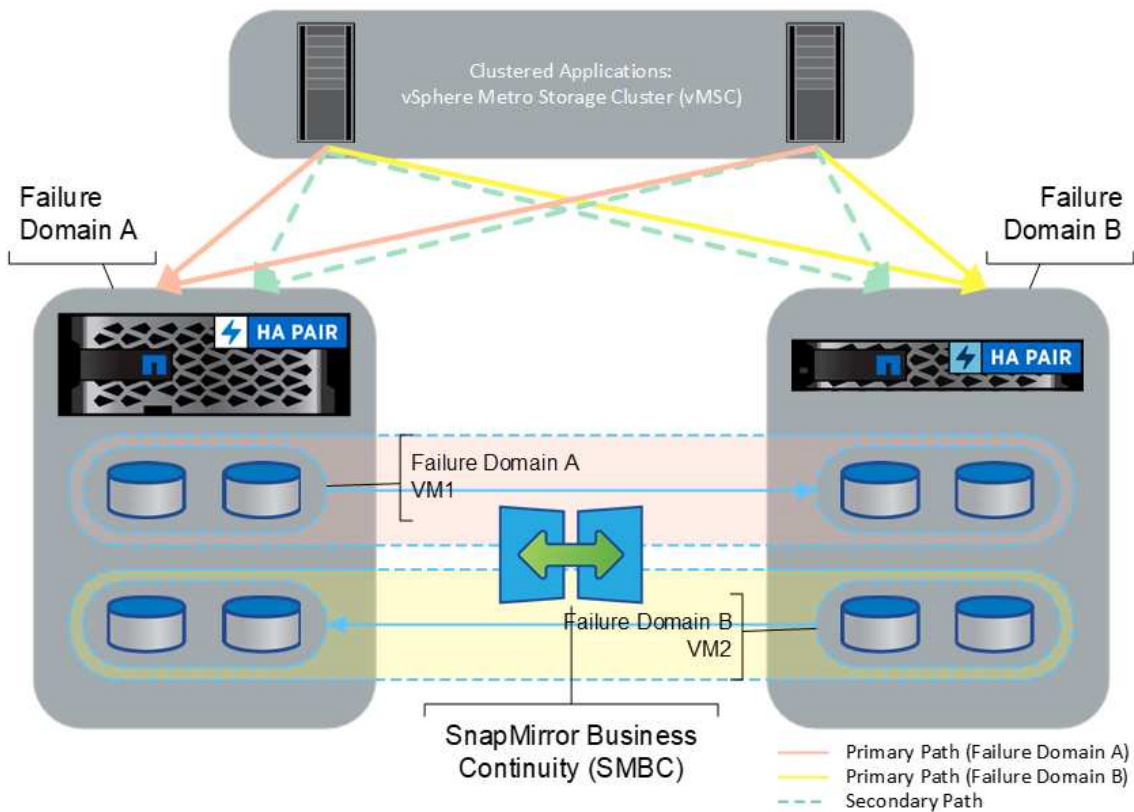
以下說明 Stretch MetroCluster 的高階拓撲圖。



請參閱 "[本文檔MetroCluster](#)" 取得 MetroCluster 的特定設計與部署資訊。

SnapMirror 主動式同步也可透過兩種不同方式部署。

- 非對稱
- 對稱 ( ONTAP 9.14.1 中的私有預覽)



請參閱 "NetApp文件" 取得 SnapMirror 主動同步的特定設計與部署資訊。

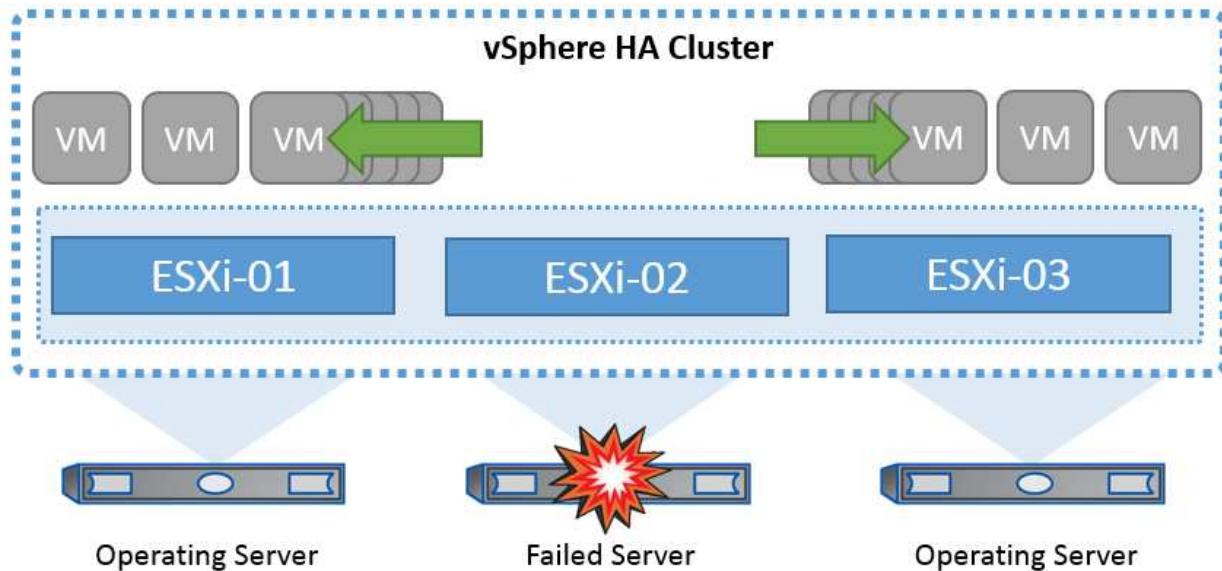
## VMware vSphere 解決方案概觀

vCenter Server Appliance (VCSA) 是強大的集中式管理系統、也是 vSphere 的單一窗口、可讓管理員有效操作 ESXi 叢集。它有助於執行重要功能、例如 VM 資源配置、VMotion 作業、高可用度 (HA)、分散式資源排程器 (DRS)、Tanzu Kubernetes Grid 等。它是 VMware 雲端環境中的重要元件、設計時應考慮到服務可用度。

### vSphere 高可用度

VMware 的叢集技術可將 ESXi 伺服器分組到虛擬機器的共用資源集區中、並提供 vSphere High Availability (HA)。vSphere HA 可為在虛擬機器中執行的應用程式提供易於使用的高可用度。當叢集上啟用 HA 功能時、每部 ESXi 伺服器都會與其他主機保持通訊、以便在任何 ESXi 主機無回應或隔離時、HA 叢集可在叢集中的未運作主機之間、協調在該 ESXi 主機上執行的虛擬機器的還原作業。萬一來賓作業系統發生故障、vSphere HA 會在同一部實體伺服器上重新啟動受影響的虛擬機器。vSphere HA 可減少計畫性停機、避免非計畫性停機、並快速從停機中恢復。

vSphere HA 叢集可從故障伺服器還原 VM。



請務必瞭解 VMware vSphere 不知道 NetApp MetroCluster 或 SnapMirror 主動同步、並視主機和 VM 群組關聯性組態而定、將 vSphere 叢集中的所有 ESXi 主機視為 HA 叢集作業的合格主機。

#### 主機故障偵測

建立 HA 叢集之後、叢集中的所有主機都會參與選舉、其中一部主機即成為主主機。每個從屬設備都會對主主機執行網路活動訊號、而主設備則會在所有從屬主機上執行網路活動訊號。vSphere HA 叢集的主要主機負責偵測從屬主機的故障。

視偵測到的故障類型而定、在主機上執行的虛擬機器可能需要容錯移轉。

在 vSphere HA 叢集中、偵測到三種類型的主機故障：

- 故障：主機停止運作。
- 隔離：主機會變成網路隔離。
- 分割區 - 主機失去與主主機的網路連線。

主機會監控叢集中的從屬主機。這種通訊是透過每秒交換網路訊號來完成。當主主機停止從從屬主機接收這些心跳時、它會在宣告主機故障之前先檢查主機的活動性。主要主機執行的活性檢查是判斷從屬主機是否與其中一個資料存放區交換活動訊號。此外、主機會檢查主機是否回應傳送至其管理 IP 位址的 ICMP Ping、以偵測其是否只是與主節點隔離、或完全與網路隔離。它會透過 ping 預設閾值來執行此作業。您可以手動指定一或多個隔離位址、以增強隔離驗證的可靠性。

#### 最佳實務做法

NetApp 建議指定至少兩個額外的隔離位址、而且每個位址都是站台本機位址。這將提高隔離驗證的可靠性。

#### 主機隔離回應

隔離回應是 vSphere HA 中的一項設定、可決定當 vSphere HA 叢集中的主機失去管理網路連線但仍繼續執行時、在虛擬機器上觸發的動作。此設定有三個選項：「已停用」、「關機並重新啟動 VM」和「關機並重新啟動 VM」。

「關機」比「關機」好、因為「關機」無法清除磁碟或認可交易的最新變更。如果虛擬機器在 300 秒內未關

機、則會關閉電源。若要變更等待時間、請使用進階選項 `das.isolationshutdowntimeout`。

在 HA 起始隔離回應之前、會先檢查 vSphere HA 主要代理程式是否擁有包含 VM 組態檔案的資料存放區。如果沒有、則主機不會觸發隔離回應、因為沒有主節點可重新啟動 VM。主機會定期檢查資料存放區狀態、以判斷是否由擁有主角色的 vSphere HA 代理程式宣告。

#### 最佳實務做法

NetApp 建議將「主機隔離回應」設定為「已停用」。

如果主機與 vSphere HA 主主機隔離或分割、而主主機無法透過心跳資料存放區或 ping 進行通訊、就可能發生分割腦部狀況。主機會宣告隔離的主機當機、並在叢集中的其他主機上重新啟動 VM。現在存在分割腦狀況、因為有兩個執行中的虛擬機器執行個體、只有其中一個執行個體可以讀取或寫入虛擬磁碟。現在可以透過設定 VM 元件保護 (VMCP) 來避免發生大腦分裂的情況。

### VM 元件保護 (VMCP)

與 HA 相關的 vSphere 6 功能增強功能之一是 VMCP。VMCP 針對區塊 (FC、iSCSI、FCoE) 和檔案儲存 (NFS)、提供增強的保護、防止所有路徑中斷 (APD) 和永久裝置遺失 (PDL) 情況。

#### 永久裝置遺失 (PDL)

當儲存設備永久故障或被管理性移除、且不預期返回時、會發生 PDL 狀況。NetApp 儲存陣列會向 ESXi 發出 SCSI Sense 程式碼、聲明該裝置已永久遺失。在 vSphere HA 的「故障條件和 VM 回應」區段中、您可以設定在偵測到 PDL 條件後應回應的內容。

#### 最佳實務做法

NetApp 建議將「使用 PDL 的資料存放區回應」設定為「\* 關閉並重新啟動 VM\*」。偵測到這種情況時、將會在 vSphere HA 叢集中的健全主機上立即重新啟動 VM。

#### 所有下行路徑 (APD)

當主機無法存取儲存裝置、且沒有通往陣列的路徑可用時、便會發生 APD 狀況。ESXi 認為這是裝置的暫時性問題、因此預期裝置會再次出現。

偵測到 APD 狀況時、會啟動定時器。140 秒後、APD 條件會正式宣告、且裝置會標示為 APD 逾時。140 秒過後、HA 會開始計算 VM 容錯移轉 APD 延遲中指定的分鐘數。指定時間過後、HA 會重新啟動受影響的虛擬機器。您可以設定 VMCP 在需要時以不同的方式回應 (停用、問題事件、或關機和重新啟動 VM)。

#### 最佳實務做法

NetApp 建議將「使用 APD 的資料存放區回應」設定為「\* 關閉並重新啟動 VM (保守) \*」。

保守是指 HA 能夠重新啟動 VM 的可能性。如果設為保守、HA 只會重新啟動受 APD 影響的 VM、前提是它知道其他主機可以重新啟動。在積極的情況下、HA 會嘗試重新啟動 VM、即使它不知道其他主機的狀態。如果沒有可存取其所在資料存放區的主機、這可能導致 VM 無法重新啟動。

如果 APD 狀態已解決、且在逾時之前已還原對儲存設備的存取、則 HA 不會不必要地重新啟動虛擬機器、除非您明確將其設定為如此。如果即使環境已從 APD 條件恢復、仍需要回應、則 APD 逾時後的 APD 恢復回應應設定為重設虛擬機器。



## 最佳實務做法 \_

NetApp 建議將 APD 逾時後的 APD 恢復回應設定為停用。

### 適用於 NetApp MetroCluster 的 VMware DRS 實作

VMware DRS 是一項功能、可將叢集中的主機資源集合在一起、主要用於在虛擬基礎架構中的叢集內進行負載平衡。VMware DRS 主要會計算 CPU 和記憶體資源、以便在叢集中執行負載平衡。由於 vSphere 不知道延伸叢集、因此在負載平衡時會考慮兩個站台中的所有主機。為了避免跨站台流量、NetApp 建議您設定 DRS 關聯性規則、以管理虛擬機器的邏輯分隔。這可確保除非發生完整的站台故障、否則 HA 和 DRS 只會使用本機主機。

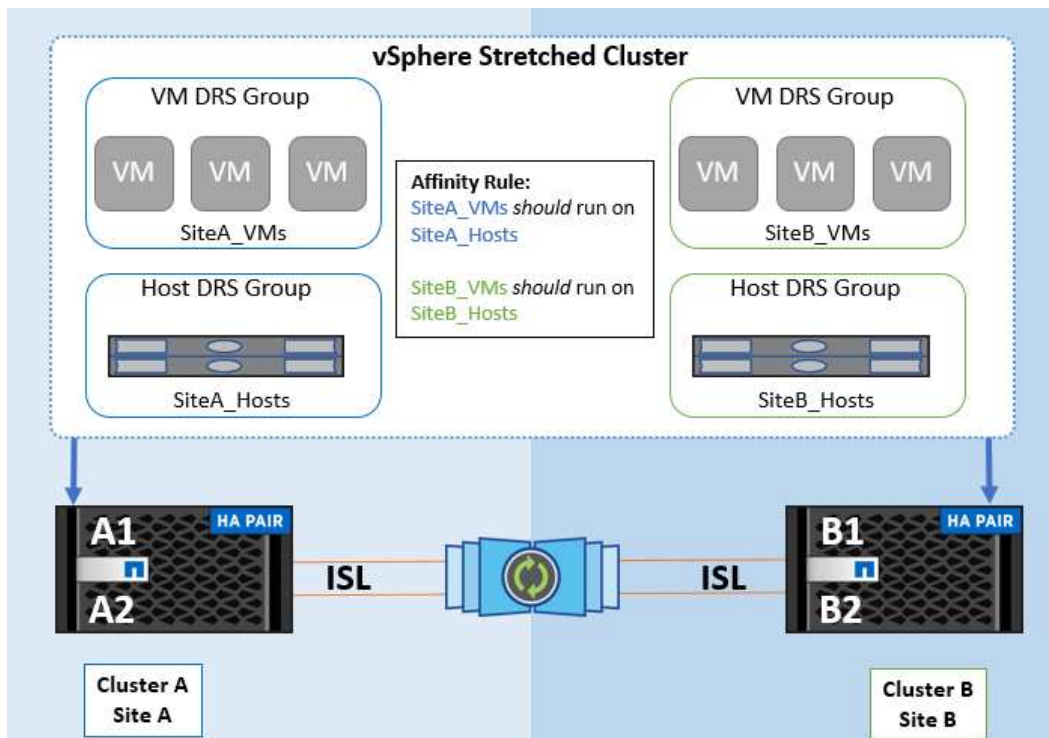
如果您為叢集建立 DRS 關聯性規則、您可以指定 vSphere 如何在虛擬機器容錯移轉期間套用該規則。

您可以指定 vSphere HA 容錯移轉行為的規則有兩種類型：

- VM 反關聯性規則會強制指定的虛擬機器在容錯移轉動作期間保持分離。
- VM 主機關聯性規則會在容錯移轉動作期間、將指定的虛擬機器放置在特定主機或已定義主機群組的成員上。

使用 VMware DRS 中的 VM 主機關聯性規則、可以在站台 A 和站台 B 之間有邏輯分隔、以便 VM 在主機上執行、而該主機與陣列是設定為指定資料存放區的主要讀取 / 寫入控制器。此外、VM 主機關聯性規則可讓虛擬機器保持儲存設備的本機狀態、進而在站台之間發生網路故障時確定虛擬機器連線。

以下是 VM 主機群組和關聯規則的範例。



## 最佳實務做法 \_

NetApp 建議實作「應該」規則、而非「必須」規則、因為在發生故障時、vSphere HA 會違反這些規則。使用「必須」規則可能導致服務中斷。

服務的可用度應永遠高於效能。在完整資料中心故障的情況下、「必須」規則必須從 VM 主機關聯群組中選擇主機、而當資料中心無法使用時、虛擬機器將不會重新啟動。

## 使用 NetApp MetroCluster 實作 VMware Storage DRS

VMware Storage DRS 功能可將資料存放區集合至單一單元、並在超過儲存 I/O 控制臨界值時平衡虛擬機器磁碟。

依預設、啟用 Storage DRS 的 DRS 叢集會啟用儲存 I/O 控制。儲存 I/O 控制功能可讓管理員控制 I/O 壅塞期間分配給虛擬機器的儲存 I/O 數量、讓更重要的虛擬機器能夠優先選擇較不重要的虛擬機器來分配 I/O 資源。

Storage DRS 使用 Storage VMotion 將虛擬機器移轉至資料存放區叢集中的不同資料存放區。在 NetApp MetroCluster 環境中、必須在該站台的資料存放區內控制虛擬機器移轉。例如、在站台 A 的主機上執行的虛擬機器 A、最好能在站台 A 的 SVM 資料存放區內移轉如果無法這麼做、虛擬機器將繼續運作、但效能降低、因為虛擬磁碟讀取 / 寫入將透過站台間連結來自站台 B。

### 最佳實務做法

NetApp 建議針對儲存站台親和性建立資料存放區叢集、也就是說、站台 A 的站台親和性資料存放區不應與站台 B 具有站台親和性的資料存放區叢集混合使用

每當使用 Storage VMotion 新佈建或移轉虛擬機器時、NetApp 建議手動更新這些虛擬機器的所有 VMware DRS 規則。這將確定主機和資料存放區在站台層級的虛擬機器關聯性、進而降低網路和儲存負荷。

## VMSC 設計與實作準則

本文件概述 VMSC 搭配 ONTAP 儲存系統的設計與實作準則。

### NetApp 儲存組態

NetApp MetroCluster 的設定指示（稱為 MCC 組態）可在以下網址取得：["資訊文件MetroCluster"](#)。SnapMirror 主動同步的說明也可在以下網址取得：["SnapMirror營運不中斷總覽"](#)。

設定 MetroCluster 之後、管理就像管理傳統的 ONTAP 環境一樣。您可以使用命令列介面（CLI）、系統管理員或 Ansible 等各種工具來設定儲存虛擬機器（SVM）。設定 SVM 後、在叢集上建立邏輯介面（生命體）、磁碟區和邏輯單元編號（LUN）、以用於正常作業。這些物件將會使用叢集對等網路自動複寫到其他叢集。

如果不使用 MetroCluster、您可以使用 SnapMirror 主動式同步功能、在不同故障網域中的多個 ONTAP 叢集之間提供資料存放區精細保護和主動式存取。SnapMirror 主動式同步會使用一致性群組、確保一或多個資料存放區之間的寫入順序一致性、您可以根據應用程式和資料存放區需求、建立多個一致性群組。一致性群組對於需要在多個資料存放區之間進行資料同步的應用程式特別有用。SnapMirror 主動式同步也支援原始裝置對應（RDM）和來賓 iSCSI 啟動器的來賓連線儲存設備。如需更多關於一致性群組的資訊、請參閱["一致性群組總覽"](#)。

與 MetroCluster 相比、使用 SnapMirror 主動式同步管理 VMSC 組態有一些差異。首先、這是僅限 SAN 的組態、沒有 NFS 資料存放區可以使用 SnapMirror 主動式同步進行保護。其次、您必須將兩個 LUN 複本對應到 ESXi 主機、以便它們存取兩個故障網域中的複寫資料存放區。

## VMware vSphere HA

### 建立 vSphere HA 叢集

建立 vSphere HA 叢集是一個多步驟程序、完整記錄於["如何在 docs.vmware.com 上的 vSphere Client 中建立和設定叢集"](#)。簡言之、您必須先建立空叢集、然後使用 vCenter 新增主機、並指定叢集的 vSphere HA 和其他

設定。

- 附註：\* 本文件並無取代之處 "[VMware vSphere Metro儲存叢集建議實務做法](#)"

若要設定 HA 叢集、請完成下列步驟：

1. 連線至 vCenter UI 。
2. 在主機和叢集中、瀏覽至您要建立 HA 叢集的資料中心。
3. 以滑鼠右鍵按一下資料中心物件、然後選取新叢集。在基礎知識之下、確保您已啟用 vSphere DRS 和 vSphere HA 。

The screenshot shows the 'New Cluster' configuration wizard in vSphere. The 'Basics' tab is selected, and the cluster name is 'MCC Cluster'. The location is set to 'Raleigh'. The 'vSphere DRS' and 'vSphere HA' options are both enabled (checked). The 'vSAN' option is disabled. Below the table, there is a checkbox for 'Manage all hosts in the cluster with a single image' which is checked. Underneath, there is a section 'Choose how to set up the cluster's image' with three radio button options: 'Compose a new image' (selected), 'Import image from an existing host in the vCenter inventory', and 'Import image from a new host'. At the bottom, there is a checkbox for 'Manage configuration at a cluster level' which is unchecked.

Name	MCC Cluster
Location	Raleigh
vSphere DRS	<input checked="" type="checkbox"/>
vSphere HA	<input checked="" type="checkbox"/>
vSAN	<input type="checkbox"/> Enable vSAN ESA

Manage all hosts in the cluster with a single image

Choose how to set up the cluster's image

Compose a new image

Import image from an existing host in the vCenter inventory

Import image from a new host

Manage configuration at a cluster level

1. 選取叢集、然後移至「組態」標籤。選取 vSphere HA 、然後按一下編輯。
2. 在 [ 主機監控 ] 下，選取 [ 啟用主機監控 ] 選項。

vSphere HA



Failures and responses | Admission Control | Heartbeat Datastores | Advanced Options

You can configure how vSphere HA responds to the failure conditions on this cluster. The following failure conditions are supported: host, host isolation, VM component protection (datastore with PDL and APD), VM and application.

Enable Host Monitoring

> Host Failure Response	Restart VMs ▾
> Response for Host Isolation	Disabled ▾
> Datastore with PDL	Power off and restart VMs ▾
> Datastore with APD	Power off and restart VMs - Conservative restart policy ▾
> VM Monitoring	Disabled ▾

CANCEL

OK

1. 在「故障與回應」標籤上、於「VM 監控」下、選取「僅限 VM 監控」選項或「VM 與應用程式監控」選項。

> Response for Host Isolation Disabled ▾

---

> Datastore with PDL Power off and restart VMs ▾

---

> Datastore with APD Power off and restart VMs - Conservative restart policy ▾

---

▾ VM Monitoring

**Enable heartbeat monitoring**

VM monitoring resets individual VMs if their VMware tools heartbeats are not received within a set time. Application monitoring resets individual VMs if their in-guest heartbeats are not received within a set time.

Disabled

VM Monitoring Only

Turns on VMware tools heartbeats. When heartbeats are not received within a set time, the VM is reset.

**VM and Application Monitoring**

Turns on application heartbeats. When heartbeats are not received within a set time, the VM is reset.

CANCEL
OK

1. 在 [ 許可控制 ] 下，將 HA 接入控制選項設定為叢集資源保留；使用 50% 的 CPU/ MEM 。

vSphere HA

Failures and responses | Admission Control | Heartbeat Datastores | Advanced Options

Admission control is a policy used by vSphere HA to ensure failover capacity within a cluster. Raising the number of potential host failures will increase the availability constraints and capacity reserved.

Host failures cluster tolerates: 1  
Maximum is one less than number of hosts in cluster.

Define host failover capacity by: Cluster resource Percentage

Override calculated failover capacity.

Reserved failover CPU capacity: 50 % CPU

Reserved failover Memory capacity: 50 % Memory

Reserve Persistent Memory failover capacity

Override calculated Persistent Memory failover capacity

CANCEL OK

1. 按一下「確定」。
2. 選取 DRS、然後按一下編輯。
3. 除非應用程式要求、否則請將自動化層級設為手動。

vSphere DRS

Automation | Additional Options | Power Management | Advanced Options

Automation Level: Manual  
DRS generates both power-on placement recommendations, and migration recommendations for virtual machines. Recommendations need to be manually applied or ignored.

Migration Threshold: Conservative (Less Frequent vMotions) | Aggressive (More Frequent vMotions)

Predictive DRS:  Enable

Virtual Machine Automation:  Enable

1. 啟用 VM 元件保護、請參閱 "[docs.vmware.com](https://docs.vmware.com)"。
2. 建議使用 MCC 的 VMSC 使用下列其他 vSphere HA 設定：

故障	回應
主機故障	重新啟動 VM
主機隔離	已停用
永久裝置遺失（PDL）的資料存放區	關閉並重新啟動 VM
All Paths Down（APD）資料存放區	關閉並重新啟動 VM
客人不會心碎	重設 VM
VM 重新啟動原則	由虛擬機器的重要性決定
主機隔離的回應	關閉並重新啟動 VM
使用 PDL 的資料存放區回應	關閉並重新啟動 VM
對具有 APD 的資料存放區的回應	關閉並重新啟動 VM（保守）
APD 的 VM 容錯移轉延遲	3 分鐘
APD 逾時的 APD 恢復回應	已停用
VM 監控靈敏度	預設為高

#### 設定資料存放區以進行心跳

當管理網路故障時、vSphere HA 會使用資料存放區來監控主機和虛擬機器。您可以設定 vCenter 如何選取心跳資料存放區。若要設定資料存放區以進行心跳、請完成下列步驟：

1. 在資料存放區心跳區段中、從指定清單中選取使用資料存放區、並在需要時自動補充資料。
2. 選取您要 vCenter 從兩個站台使用的資料存放區、然後按下 OK。

vSphere HA








Failures and responses   Admission Control   **Heartbeat Datastores**   Advanced Options

vSphere HA uses datastores to monitor hosts and virtual machines when the HA network has failed. vCenter Server selects 4 datastores for each host using the policy and datastore preferences specified below.

Heartbeat datastore selection policy:

- Automatically select datastores accessible from the hosts
- Use datastores only from the specified list
- Use datastores from the specified list and complement automatically if needed

Available heartbeat datastores

	Name ↑	Datastore Cluster	Hosts Mounting Datastore
<input checked="" type="checkbox"/>	 d11	N/A	2
<input checked="" type="checkbox"/>	 d12	N/A	2
<input checked="" type="checkbox"/>	 d21	N/A	2
<input checked="" type="checkbox"/>	 d22	N/A	2
<input type="checkbox"/>	 d31	N/A	2
<input type="checkbox"/>	 d32	N/A	2
<input type="checkbox"/>	 d41	N/A	2
<input type="checkbox"/>	 d42	N/A	2

11 items

CANCEL OK

### 設定進階選項

- 主機故障偵測 \*

當 HA 叢集內的主機無法連線至網路或叢集中的其他主機時、就會發生隔離事件。根據預設、vSphere HA 會使用其管理網路的預設閘道做為預設隔離位址。不過、您可以為主機指定其他隔離位址來執行 ping、以判斷是否應該觸發隔離回應。新增兩個可 ping 的隔離 IP、每個站台一個。請勿使用閘道 IP。使用的 vSphere HA 進階設定為 `das.isolationaddress`。您可以將 ONTAP 或 Mediator IP 位址用於此用途。

請參閱 "[core.vmware.com](https://core.vmware.com)" 以取得更多資訊



vSphere HA

Failures and responses   Admission Control   Heartbeat Datastores   **Advanced Options**

You can set advanced options that affect the behavior of your vSphere HA cluster.

+ Add   ✕ Delete

Option	Value
das.IgnoreRedundantNetWarning	true
das.Isolationaddress0	10.61.99.100
das.Isolationaddress1	10.61.99.110
das.heartbeatDsPerHost	4

4 items

CANCEL   OK

新增稱為 das.心跳 DsPerHost 的進階設定、可能會增加心跳資料存放區的數量。使用四個心跳資料存放區（HB DSS）、每個站台兩個。使用「從清單中選取但輔助」選項。這是必要的、因為如果某個站台發生故障、您仍需要兩個 HB DSS。但是、這些不需要透過 MCC 或 SnapMirror 主動同步來保護。

請參閱 "[core.vmware.com](http://core.vmware.com)" 以取得更多資訊

### 適用於 NetApp MetroCluster 的 VMware DRS 關聯性

在本節中、我們會為 MetroCluster 環境中的每個站台 \ 叢集、建立 VM 和主機 DRS 群組。然後我們設定 VM\Host 規則、使 VM 主機與本機儲存資源的關聯性一致。例如、站台 A VM 屬於 VM 群組 sitea\_vms、站台 A 主機屬於主機群組 sitea\_hosts。接下來、在 VM\Host 規則中、我們指出 sitea\_vms 應該在 sitea\_hosts 中的主機上執行。

#### 最佳實務做法

- NetApp 強烈建議在組 \* 中的主機上運行規範 \*，而不是規範 \* 必須在組 \* 中的主機上運行。萬一站台 A 主機故障、站台 A 的 VM 需要透過 vSphere HA 在站台 B 的主機上重新啟動、但後者的規格不允許 HA 在站台 B 上重新啟動 VM、因為這是硬規則。以前的規格是軟性規則、在 HA 發生時會違反、因此可提供可用度而非效能。
- 附註：\* 您可以建立事件型警示、在虛擬機器違反 VM-Host 關聯性規則時觸發。在 vSphere Client 中、新增虛擬機器的警示、並選取「VM 正在違反 VM-Host Affinity Rule」作為事件觸發程序。如需建立及編輯警

示的詳細資訊、請參閱 "[vSphere 監控與效能](#)" 文件。

#### 建立 DRS 主機群組

若要建立站台 A 和站台 B 專屬的 DRS 主機群組、請完成下列步驟：

1. 在 vSphere Web Client 中、以滑鼠右鍵按一下資源清冊中的叢集、然後選取「設定」。
2. 按一下 VM\Host Groups。
3. 按一下「新增」
4. 輸入群組的名稱（例如、sitea\_hosts）。
5. 從「類型」功能表中、選取「主機群組」。
6. 按一下「新增」、然後從站台 A 選取所需的主機、再按一下「確定」。
7. 重複這些步驟、為站台 B 新增另一個主機群組
8. 按一下「確定」。

#### 建立 DRS VM 群組

若要建立站台 A 和站台 B 專屬的 DRS VM 群組、請完成下列步驟：

1. 在 vSphere Web Client 中、以滑鼠右鍵按一下資源清冊中的叢集、然後選取「設定」。
2. 按一下 VM\Host Groups。
3. 按一下「新增」
4. 輸入群組的名稱（例如、sitea\_vms）。
5. 從 Type（類型）功能表中、選取 VM Group（VM 群組）。
6. 按一下「新增」、然後從站台 A 選取所需的 VM、再按一下「確定」。
7. 重複這些步驟、為站台 B 新增另一個主機群組
8. 按一下「確定」。

#### 建立 VM Host 規則

若要建立站台 A 和站台 B 特有的 DRS 關聯性規則、請完成下列步驟：

1. 在 vSphere Web Client 中、以滑鼠右鍵按一下資源清冊中的叢集、然後選取「設定」。
2. 按一下 VM\Host Rules。
3. 按一下「新增」
4. 輸入規則的名稱（例如、sitea\_fit射）。
5. 確認已核取「啟用規則」選項。
6. 從 Type（類型）功能表中、選取 Virtual Machines to Hosts（虛擬機器至主機）。
7. 選取 VM 群組（例如、sitea\_vms）。
8. 選取主機群組（例如、sitea\_hosts）。
9. 重複這些步驟、為站台 B 新增另一個 VM\ 主機規則

10. 按一下「確定」。

## Create VM/Host Rule | Cluster-01 ×

Name	sitea_affinity <input checked="" type="checkbox"/> Enable rule.
Type	Virtual Machines to Hosts <span>▼</span>

Virtual machines that are members of the Cluster VM Group sitea\_vms should run on host group sitea\_hosts.

VM Group:

sitea_vms <span>▼</span>
Should run on hosts in group <span>▼</span>

Host Group:

sitea_hosts <span>▼</span>
----------------------------

CANCEL	OK
--------	----

### VMware vSphere Storage DRS for NetApp MetroCluster

建立資料存放區叢集

若要為每個站台設定資料存放區叢集、請完成下列步驟：

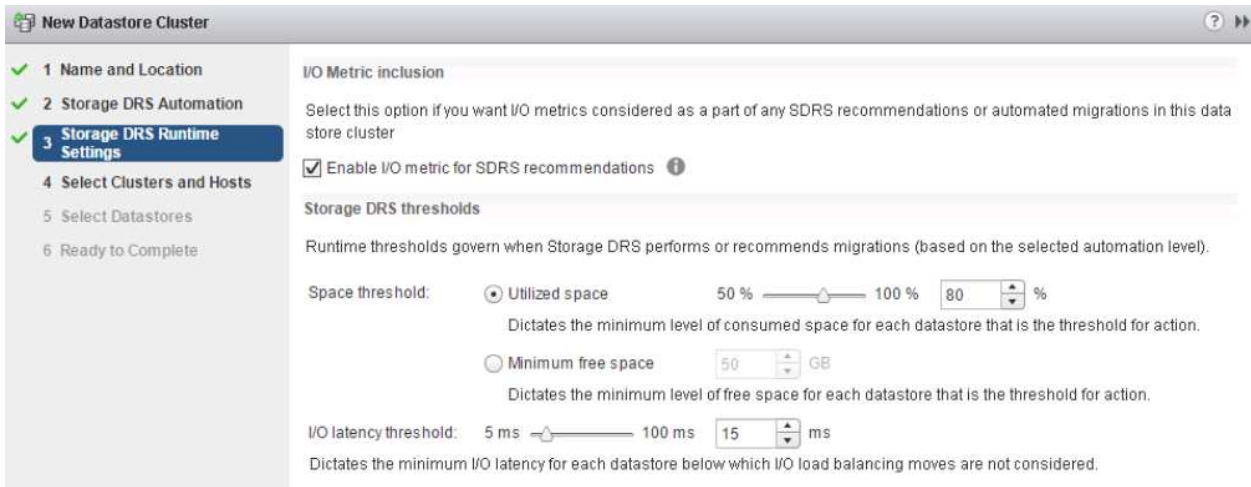
1. 使用 vSphere Web 用戶端、瀏覽至 HA 叢集位於 Storage 下的資料中心。
2. 以滑鼠右鍵按一下資料中心物件、然後選取儲存 > 新資料存放區叢集。
3. 選取「開啟 Storage DRS」選項、然後按一下「下一步」。
4. 將所有選項設定為「無自動化（手動模式）」、然後按一下「下一步」。

最佳實務做法 \_

- NetApp 建議您將儲存 DRS 設定為手動模式、以便系統管理員決定並控制何時需要移轉。

Storage DRS automation	
Cluster automation level	<input checked="" type="radio"/> <b>No Automation (Manual Mode)</b> vCenter Server will make migration recommendations for virtual machine storage, but will not perform automatic migrations.
	<input type="radio"/> <b>Fully Automated</b> Files will be migrated automatically to optimize resource usage.

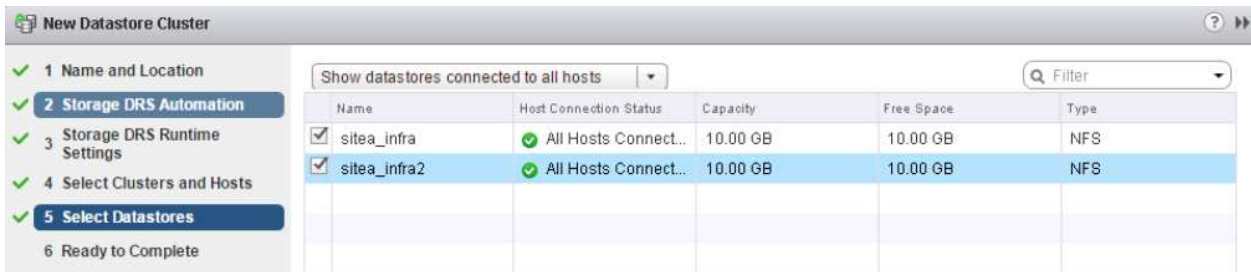
1. 確認已核取「啟用 SDR 建議的 I/O 度量」核取方塊；度量設定可以保留預設值。



1. 選取 HA 叢集、然後按一下「下一步」。



1. 選取屬於站台 A 的資料存放區、然後按一下「下一步」。



1. 檢閱選項、然後按一下「完成」。
2. 重複這些步驟以建立站台 B 資料存放區叢集、並確認只選取站台 B 的資料存放區。

## vCenter Server 可用性

您的 vCenter Server 應用裝置（VCSA）應使用 vCenter HA 加以保護。vCenter HA 可讓您在主動式被動式 HA 配對中部署兩個 VCSA。每個故障網域各有一個。您可以在上閱讀更多有關 vCenter HA 的資訊 "[docs.vmware.com](https://docs.vmware.com)"。

## 計畫性和非計畫性事件的恢復能力

NetApp MetroCluster 和 SnapMirror 主動同步是強大的工具、可增強 NetApp 硬體和 ONTAP® 軟體的高可用度和不中斷營運。

這些工具可為整個儲存環境提供全站台保護、確保資料永遠可用。無論您是使用獨立式伺服器、高可用度伺服器

叢集、 Docker 容器或虛擬化伺服器、 NetApp 技術都能在電力中斷、冷卻或網路連線中斷、儲存陣列關機或作業錯誤等情況下、無縫維持儲存可用度。

MetroCluster 和 SnapMirror 主動式同步提供三種基本方法、可在發生計畫性或非計畫性事件時維持資料連續性：

- 備援元件、可防止單一元件故障
- 本機 HA 接管、用於影響單一控制器的事件
- 完整的站台保護：將儲存設備和用戶端存取從來源叢集移至目的地叢集、以快速恢復服務

這表示在單一元件故障時、作業會順暢地繼續、並在更換故障元件時自動恢復至備援作業。

除了單節點叢集（例如 ONTAP Select 等軟體定義版本）之外、所有 ONTAP 叢集都具有稱為接管和恢復的內建 HA 功能。叢集中的每個控制器都會與另一個控制器配對、形成 HA 配對。這些配對可確保每個節點都在本機上連線至儲存設備。

接管是一種自動化程序、其中一個節點會接管另一個節點的儲存設備、以維護資料服務。GiveBack 是還原正常作業的反向程序。可以規劃接管、例如執行硬體維護或 ONTAP 升級、或是因節點緊急或硬體故障而非計畫性地進行。

在接管期間、 MetroCluster 組態中的網路附加儲存邏輯介面（ NAS 生命期）會自動容錯移轉。但是、儲存區域網路生命（ SAN 生命）不會容錯移轉；它們會繼續使用邏輯單元編號（ LUN ）的直接路徑。

如需 HA 接管與恢復的詳細資訊、請參閱 "[HA配對管理總覽](#)"。值得一提的是、這項功能並非 MetroCluster 或 SnapMirror 主動式同步的專屬功能。

當某個站台離線、或是作為整個站台維護的計畫活動時、就會使用 MetroCluster 進行站台切換。其餘站台則假設擁有離線叢集的儲存資源（磁碟和集合體）、而故障站台上的 SVM 則會在災難站台上線並重新啟動、保留其完整身分以供用戶端和主機存取。

有了 SnapMirror 主動式同步、由於兩個複本都是同時使用的、因此您現有的主機將繼續運作。NetApp Mediator 是確保站台容錯移轉正確進行所需的工具。

## 使用 MCC 的 VMSC 的失敗案例

以下各節概述 VMSC 和 NetApp MetroCluster 系統各種故障情況的預期結果。

### 單一儲存路徑故障

在這種情況下、如果 HBA 連接埠、網路連接埠、前端資料交換器連接埠或 FC 或乙太網路纜線等元件故障、 ESXi 主機會將該儲存裝置的特定路徑標記為已停用。如果在 HBA/ 網路 / 交換器連接埠上提供恢復功能、就能為儲存裝置設定多個路徑、 ESXi 理想情況下會執行路徑切換。在這段期間內、虛擬機器會持續執行而不會受到影響、因為提供多條路徑可通往儲存設備、因此可確保儲存設備的可用度。

- 附註： \* 在此案例中、 MetroCluster 行為並無變更、所有資料存放區仍會保留在各自站台內。

### 最佳實務做法 \_

在使用 NFS/iSCSI 磁碟區的環境中、 NetApp 建議在標準 vSwitch 中、至少為 NFS vmkernel 連接埠設定兩個網路上行鏈路、而在對應 NFS vmkernel 介面的連接埠群組中、則必須設定相同的上行鏈路。 NIC 群組可在雙主動式或雙主動式待命模式中進行設定。

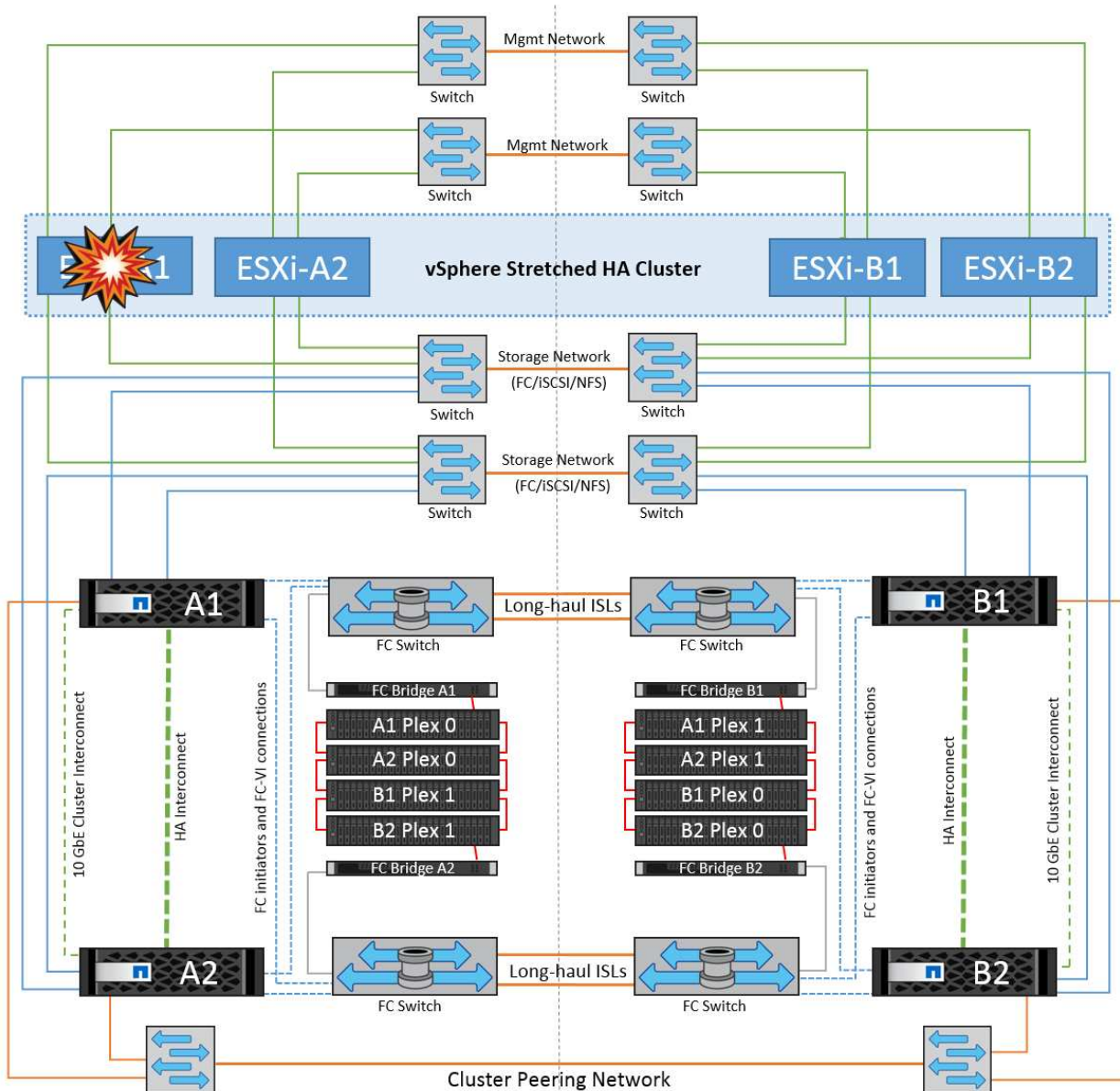
此外、對於 iSCSI LUN 、必須將 vmkernel 介面繫結至 iSCSI 網路介面卡、以設定多重路徑。如需詳細資訊、請參閱 vSphere 儲存文件。

最佳實務做法 \_

在使用光纖通道 LUN 的環境中、NetApp 建議至少有兩個 HBA 、以保證 HBA/ 連接埠層級的恢復能力。NetApp 也建議將單一啟動器分區至單一目標分區、以做為設定分區的最佳實務做法。

應使用虛擬儲存主控台（VSC）來設定多重路徑原則、因為它會為所有新的和現有的 NetApp 儲存裝置設定原則。

### 單一 ESXi 主機故障



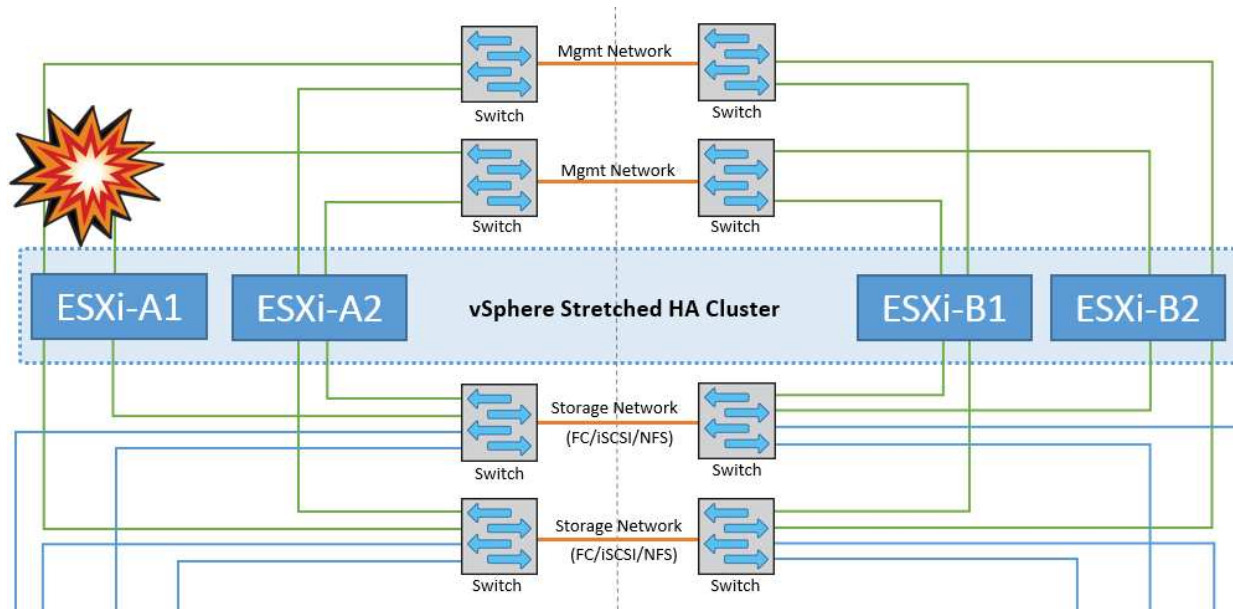
在這種情況下、如果 ESXi 主機發生故障、VMware HA 叢集中的主節點會偵測主機故障、因為主機不再接收到網路心跳。若要判斷主機是否真的停機或只是網路分割區、主節點會監控資料存放區的訊次、如果沒有、則會 ping 失敗主機的管理 IP 位址、以執行最終檢查。如果所有這些檢查都是負數、則主節點會將此主機宣告為故障主機、而在該故障主機上執行的所有虛擬機器都會在叢集中的正常主機上重新開機。

如果已設定 DRS VM 和主機關聯性規則（VM 群組 sitea\_vms 中的 VM 應在主機群組 sitea\_hosts 中執行主機）、則 HA 主機會先檢查站台 A 的可用資源。如果站台 A 沒有可用的主機、則主機會嘗試在站台 B 的主機上重新啟動 VM。

如果本機台有資源限制、則可能會在其他站台的 ESXi 主機上啟動虛擬機器。不過、如果將虛擬機器移轉回本機台中任何仍在運作的 ESXi 主機、而違反任何規則、則定義的 DRS VM 和主機關聯性規則將會修正。如果 DRS 設定為手動、NetApp 建議您啟動 DRS、並套用建議來修正虛擬機器的放置位置。

在此案例中、MetroCluster 行為並無任何變更、所有資料存放區仍會保持不變、不受其個別站台影響。

### ESXi 主機隔離

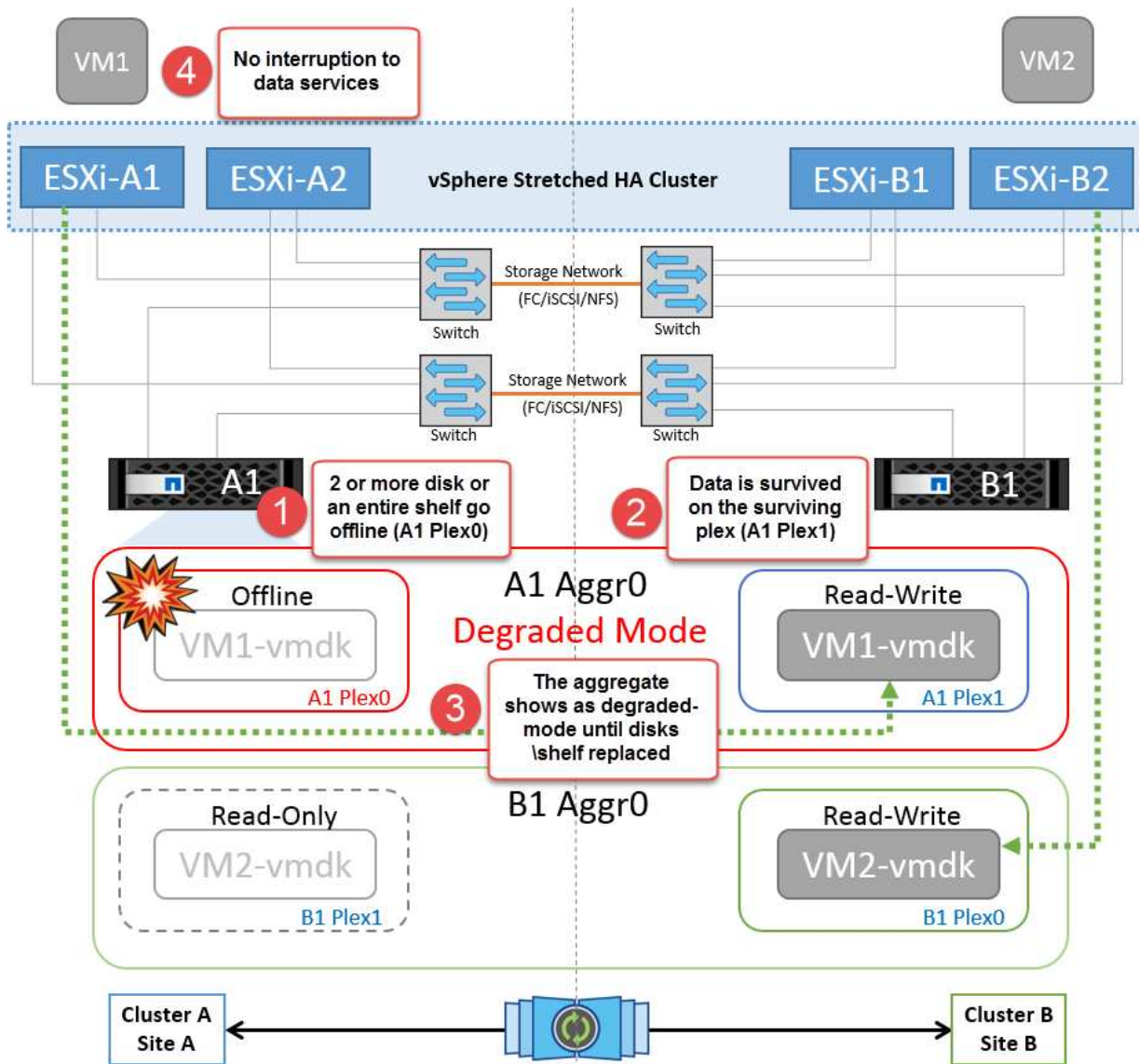


在此案例中、如果 ESXi 主機的管理網路中斷、HA 叢集中的主節點將不會接收任何訊息、因此此主機會在網路中隔離。若要判斷它是否發生故障或只是隔離、主節點會開始監控資料存放區心跳。如果主機存在、則主機會由主節點宣告為隔離。根據設定的隔離回應、主機可能會選擇關閉、關閉虛擬機器、甚至讓虛擬機器保持開機。隔離回應的預設時間間隔為 30 秒。

在此案例中、MetroCluster 行為並無任何變更、所有資料存放區仍會保持不變、不受其個別站台影響。

### 磁碟機櫃故障

在此案例中、有兩個以上的磁碟或整個機櫃發生故障。資料是從仍在運作的複合環境提供、不會中斷資料服務。磁碟故障可能會影響本機或遠端叢。由於只有一個叢處於作用中狀態、因此集合體會顯示為降級模式。更換故障磁碟後、受影響的集合體將自動重新同步以重建資料。重新同步後、集合體將自動返回正常的鏡射模式。如果單一 RAID 群組中有兩個以上的磁碟發生故障、則必須從頭重建叢。

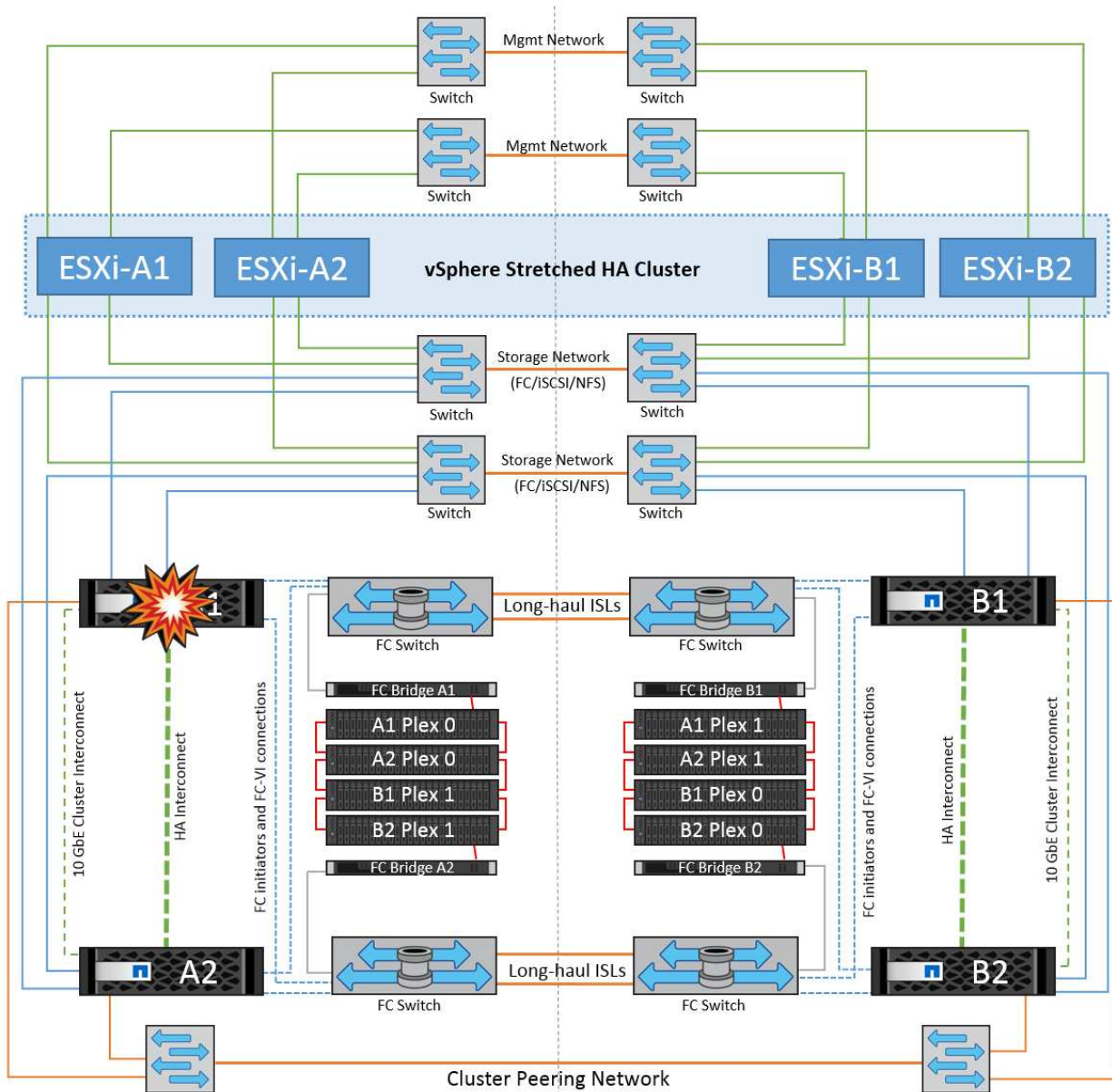


- 注意：\* 在此期間、虛擬機器 I/O 作業不會受到影響、但效能會降低、因為資料是透過 ISL 連結從遠端磁碟機櫃存取。

#### 單一儲存控制器故障

在這種情況下、兩個儲存控制器中的其中一個會在一個站台發生故障。由於每個站台都有 HA 配對、因此一個節點的故障會以透明方式自動觸發容錯移轉至另一個節點。例如、如果節點 A1 故障、其儲存設備和工作負載會自動傳輸至節點 A2。虛擬機器將不會受到影響、因為所有的叢集都仍然可用。第二個站台節點（B1 和 B2）不受影響。此外、vSphere HA 將不會採取任何行動、因為叢集中的主節點仍會接收到網路心跳。

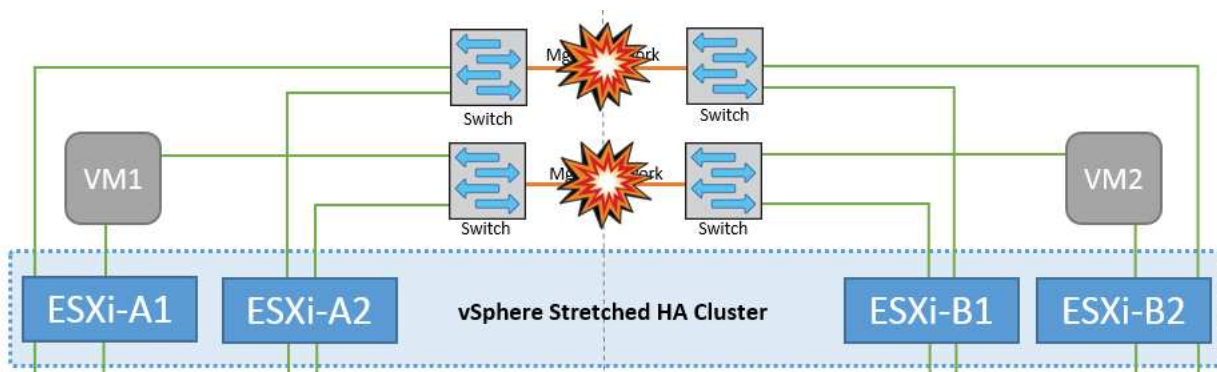




如果容錯移轉是循環災難的一部分（節點 A1 容錯移轉至 A2）、而且之後發生 A2 故障、或是站台 A 完全故障、則災難後的切換可能會發生在站台 B

交換器間連結故障

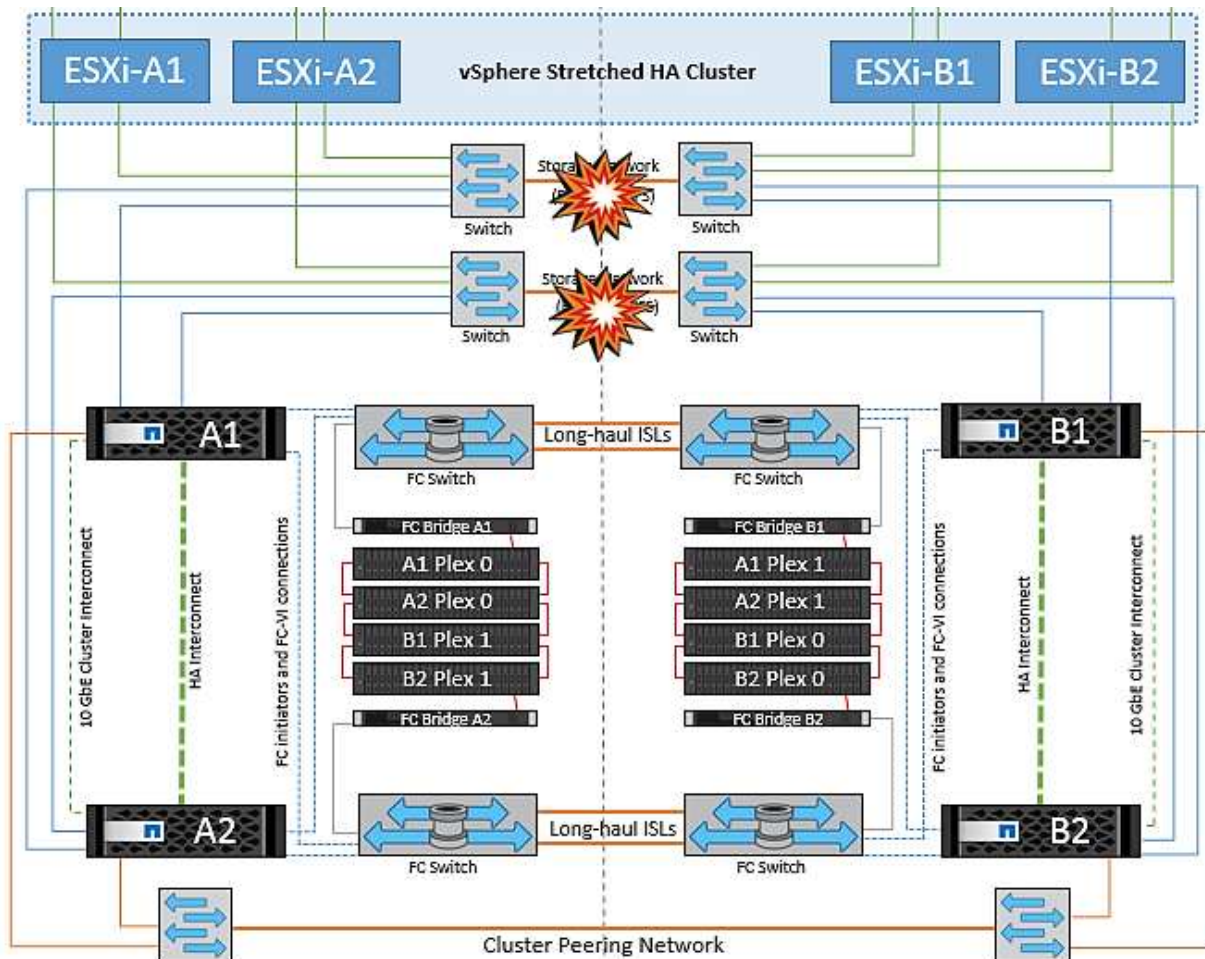
管理網路的交換器間連結故障



在此案例中、如果前端主機管理網路的 ISL 連結失敗、站台 A 的 ESXi 主機將無法與站台 B 的 ESXi 主機通訊這會導致網路分割區、因為特定站台的 ESXi 主機將無法將網路心跳傳送至 HA 叢集中的主節點。因此、由於分割區的緣故、將會有兩個網路區段、每個區段中都會有一個主節點、可保護 VM 免於特定站台內的主機故障。

- 附註：\* 在此期間、虛擬機器仍在執行中、在此案例中、MetroCluster 行為並無變更。所有的資料存放區都會繼續保持不變、不受其個別站台影響。

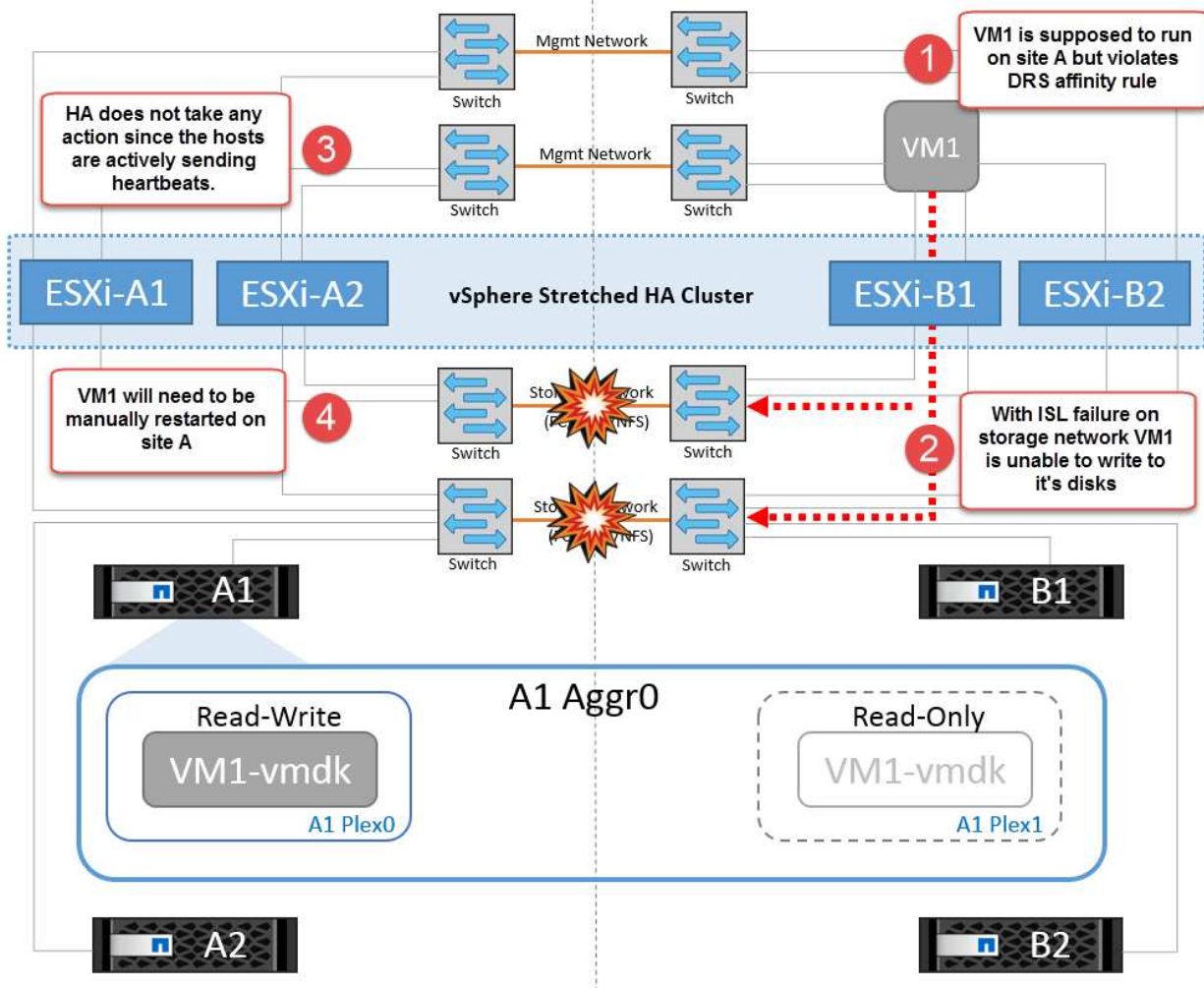
儲存網路的交換器間連結故障



在此案例中、如果後端儲存網路的 ISL 連結故障、站台 A 的主機將無法存取站台 B 的儲存磁碟區或叢集 B 的 LUN、反之亦然。VMware DRS 規則的定義、是為了讓主機儲存站台的關聯性能讓虛擬機器在不影響站台的情況下執行。

在此期間、虛擬機器會繼續在各自的站台上執行、在此案例中、MetroCluster 行為不會有任何變更。所有的資料存放區都會繼續保持不變、不受其個別站台影響。

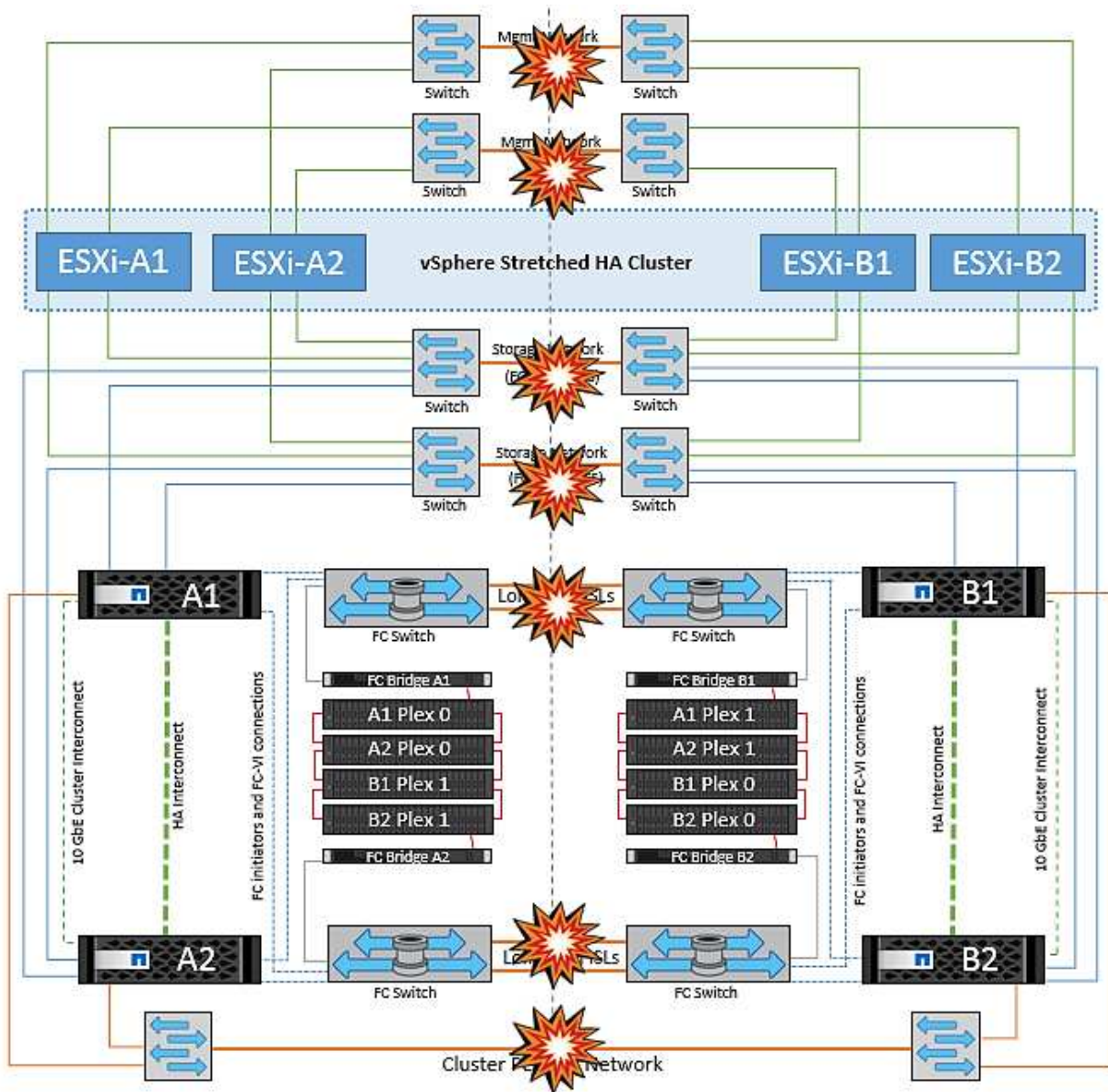
如果因為某種原因違反關聯規則（例如、VM1 原本應從站台 A 執行、其磁碟位於本機叢集 A 節點上、而 VM1 則是在站台 B 的主機上執行）、則虛擬機器的磁碟將透過 ISL 連結遠端存取。由於 ISL 連結故障、在站台 B 執行的 VM1 將無法寫入其磁碟、因為通往儲存磁碟區的路徑已關閉、且該特定虛擬機器已關閉。在這些情況下、VMware HA 不會採取任何行動、因為主機正在主動傳送心跳。這些虛擬機器必須在各自的站台手動關閉並開啟電源。下圖說明違反 DRS 關聯性規則的虛擬機器。



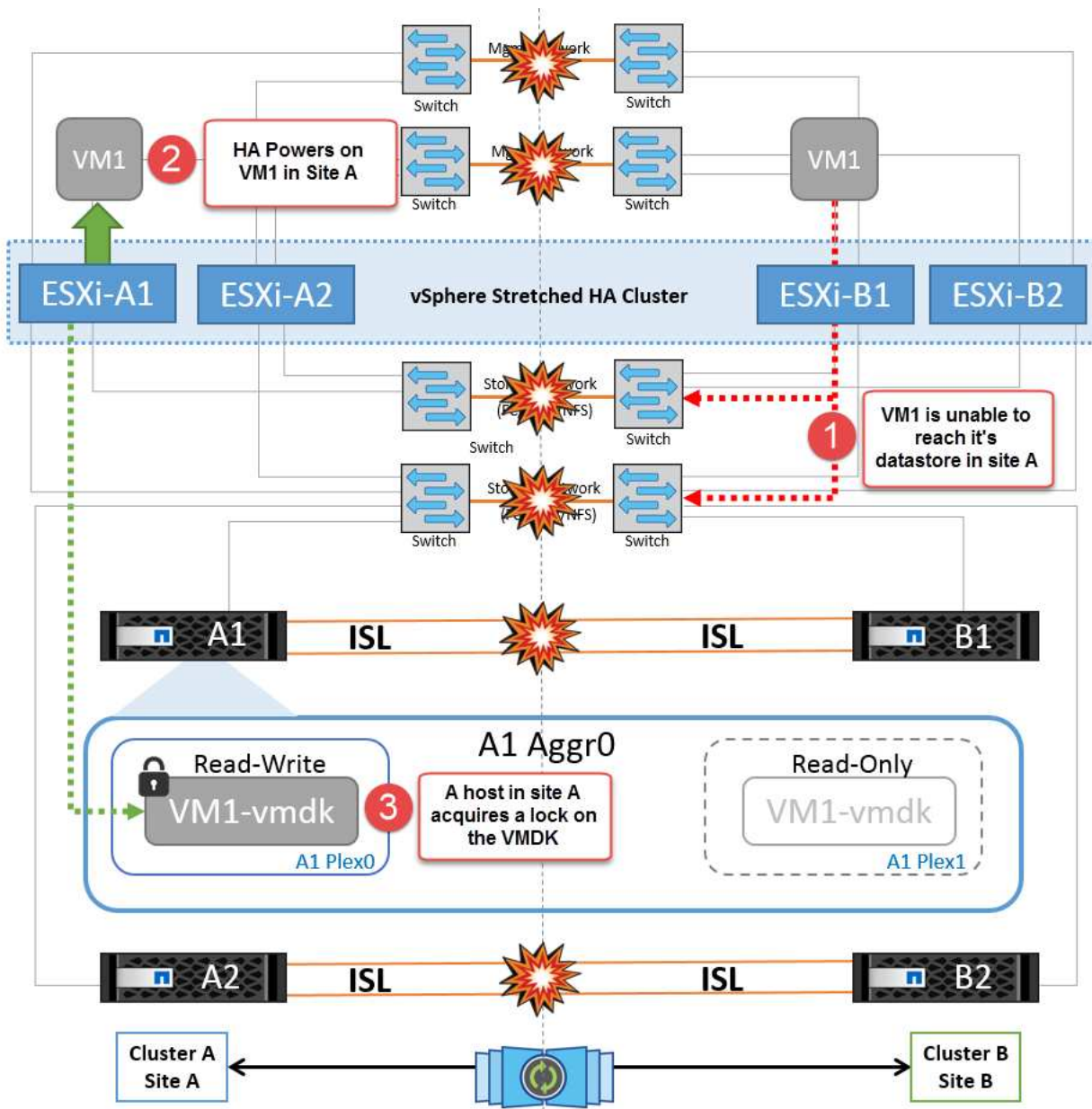
所有交換器間故障或完整資料中心分割區

在這種情況下、站台之間的所有 ISL 連結都會中斷、而且兩個站台彼此之間會隔離。如先前的案例所述、例如管理網路和儲存網路的 ISL 故障、虛擬機器在完全 ISL 故障時不會受到影響。

在站台之間分割 ESXi 主機之後、vSphere HA 代理程式會檢查資料存放區心跳、而且在每個站台中、本機 ESXi 主機將能夠將資料存放區心跳更新至各自的讀寫磁碟區 /LUN。站台 A 中的主機會假設站台 B 中的其他 ESXi 主機故障、因為沒有網路 / 資料存放區檢測信號。站台 A 的 vSphere HA 會嘗試重新啟動站台 B 的虛擬機器、最終會失敗、因為站台 B 的資料存放區因為儲存 ISL 故障而無法存取。站台 B 也會再次出現類似的情況



NetApp 建議判斷是否有任何虛擬機器違反 DRS 規則。從遠端站台執行的任何虛擬機器都會停機、因為它們將無法存取資料存放區、vSphere HA 會在本機站台上重新啟動該虛擬機器。當 ISL 連結恢復上線後、在遠端站台上執行的虛擬機器將會停止運作、因為無法有兩個執行個體使用相同的 MAC 位址執行虛擬機器。

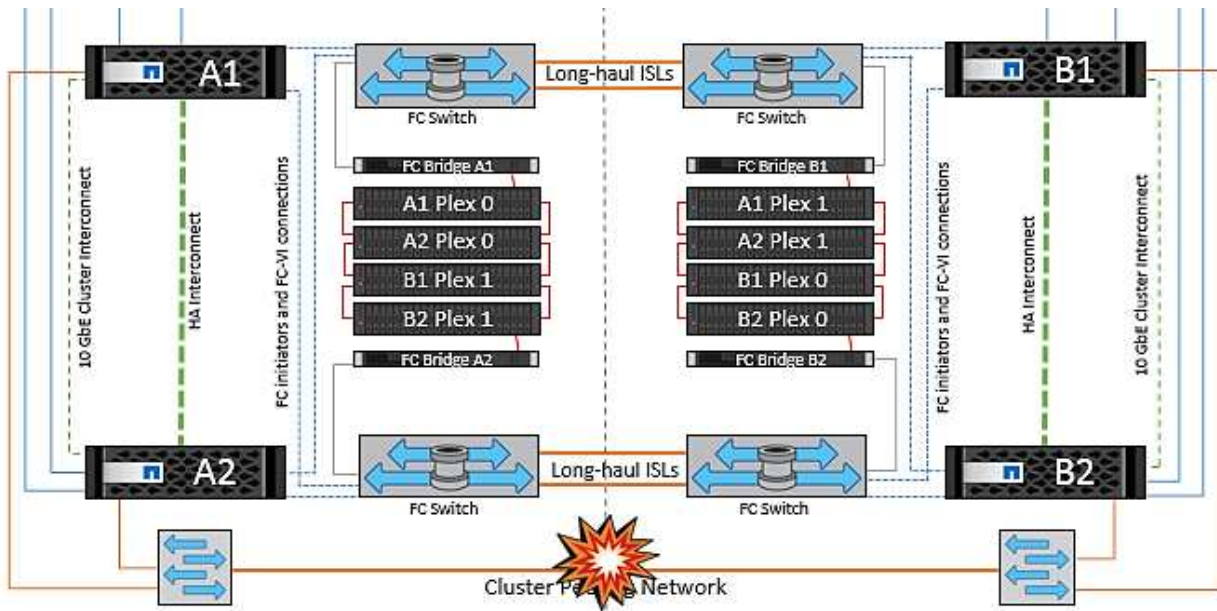


### NetApp MetroCluster 中兩個 Fabric 上的交換器間連結故障

在一個或多個 ISL 故障的情況下、流量會繼續流經其餘的連結。如果兩個架構上的所有 ISL 都發生故障、使得儲存和 NVRAM 複寫站台之間沒有連結、則每個控制器都會繼續提供其本機資料。還原至少一個 ISL 時、所有的叢會自動重新同步。

在所有 ISL 停機之後所發生的任何寫入動作、都不會鏡射到另一個站台。當組態處於此狀態時、發生災難時的切入將會遺失尚未同步的資料。在這種情況下、需要手動介入才能在進行重新操作後恢復。如果很可能在較長的時間內沒有可用的 ISL、系統管理員可以選擇關閉所有資料服務、以避免在發生災難時發生資料遺失的風險。在至少有一個 ISL 可供使用之前、應將執行此動作的可能性與需要進行重新操作的災難可能性進行權衡。或者、如果 ISL 在串聯案例中發生故障、系統管理員可能會在所有連結失敗之前、觸發已規劃的切換至其中一個站台。





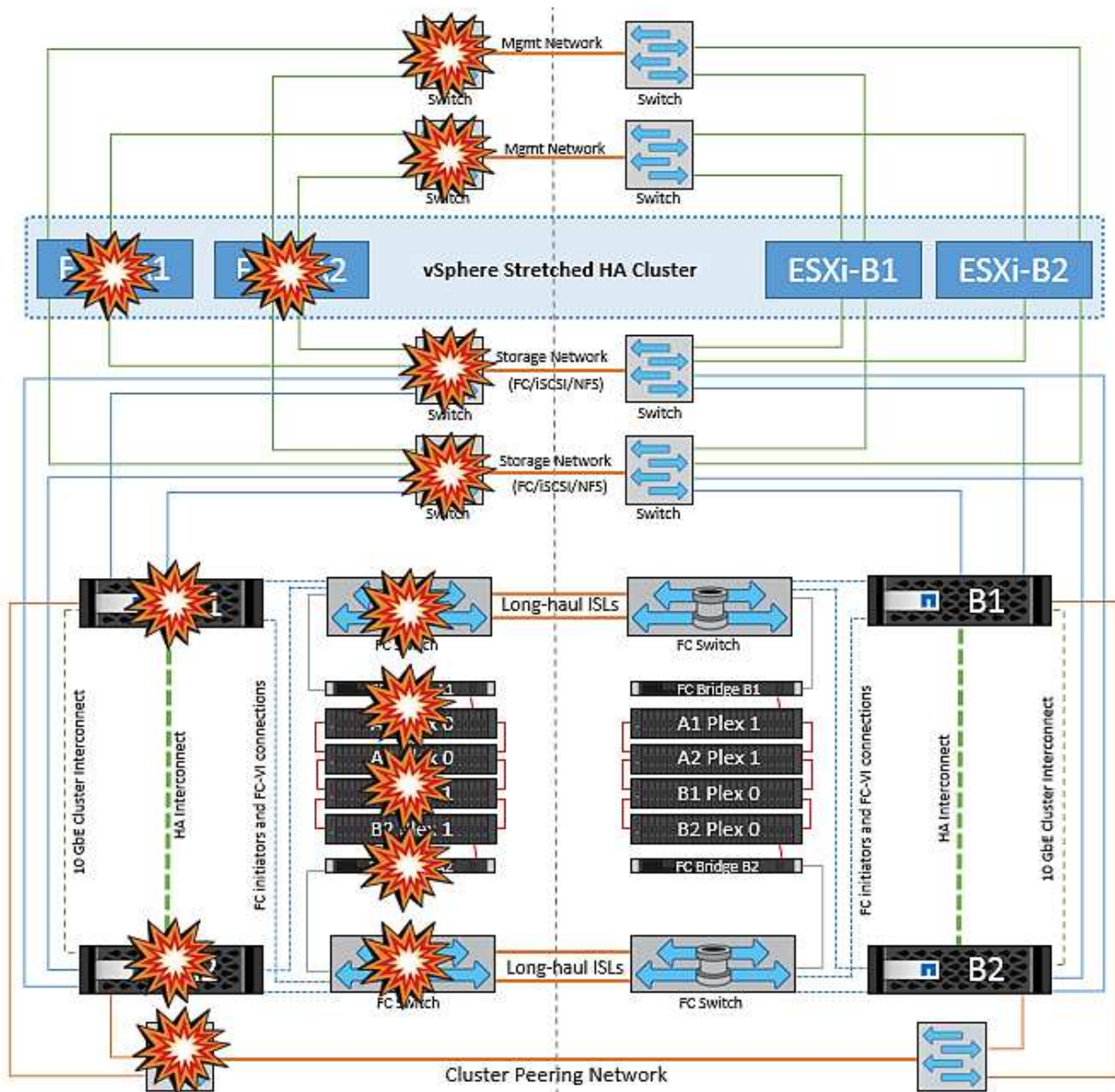
#### 完成站台故障

在完整站台 A 故障案例中、站台 B 的 ESXi 主機因為故障而無法從站台 A 的 ESXi 主機取得網路心跳。站台 B 的 HA 主機會驗證資料存放區心跳不存在、宣告站台 A 的主機故障、並嘗試重新啟動站台 B 中的站台 A 虛擬機器在此期間、儲存管理員會執行一次轉換、以恢復仍在運作的站台上故障節點的服務、該站台將還原站台 B 上站台 A 的所有儲存服務站台 A 磁碟區或 LUN 在站台 B 上可用後、HA 主代理程式會嘗試重新啟動站台 B 中的站台 A 虛擬機器

如果 vSphere HA 主要代理程式嘗試重新啟動虛擬機器（包括登錄及開機）失敗、則會在延遲後重試重新啟動。重新啟動之間的延遲時間最多可設定為 30 分鐘。vSphere HA 會嘗試重新啟動這些項目、最多嘗試次數（預設為六次）。

- 附註：\* 在放置管理程式找到適當的儲存設備之前、HA 主機不會開始重新啟動嘗試、因此在整個站台發生故障的情況下、這將是在執行切入之後。

如果站台 A 已切換、則可透過容錯移轉至正常運作的節點、無縫地處理其中一個仍在運作的站台 B 節點的後續故障。在這種情況下、四個節點的工作現在僅由一個節點執行。在這種情況下、恢復將包括執行恢復到本機節點的贈品。然後、當站台 A 還原時、會執行切換作業、以還原組態的穩定狀態作業。



## 產品安全性

### VMware vSphere適用的工具ONTAP

採用 ONTAP Tools for VMware vSphere 的軟體工程採用下列安全開發活動：

- \*威脅建模。\*威脅建模的目的是在軟體開發生命週期初期、發現某項功能、元件或產品的安全性瑕疵。威脅模式是對影響應用程式安全性的所有資訊的結構化呈現。基本上、它是透過安全性觀點來檢視應用程式及其環境。
- \*動態應用程式安全性測試 (dast)。\*這項技術的設計、是為了偵測應用程式在執行狀態下的易受影響狀況。Dast會測試開放Web應用程式的公開HTTP和HTML介面。
- \*協力廠商程式碼貨幣。\*在開放原始碼軟體 (開放原始碼軟體) 的軟體開發過程中、您必須解決與產品內建的任何開放原始碼軟體相關的安全性弱點。這是一項持續努力、因為新的開放源碼版本可能隨時都有新發現的弱點報告。
- \*弱點掃描。\*弱點掃描的目的是在NetApp產品中發現常見且已知的安全性弱點之後、再將弱點發佈給客戶。



- \*滲透測試。\*滲透測試是評估系統、Web應用程式或網路以找出攻擊者可能利用的安全性弱點的程序。NetApp的滲透測試（筆測試）是由一群獲核准且值得信賴的第三方公司進行。其測試範圍包括利用精密的利用方法或工具、對類似於惡意入侵者或駭客的應用程式或軟體發動攻擊。

## 產品安全功能

適用於 VMware vSphere 的 ONTAP 工具在每個版本中都包含下列安全功能。

- 登入橫幅。SSH預設為停用、如果從VM主控台啟用、則僅允許一次性登入。使用者在登入提示中輸入使用者名稱後、會顯示下列登入橫幅：

\*警告：\*禁止未經授權存取本系統、並依法律予以起訴。存取本系統即表示您同意、若懷疑有未獲授權的使用情形、您的行動可能受到監控。

使用者透過 SSH 通道完成登入後、會顯示下列文字：

```
Linux vsc1 4.19.0-12-amd64 #1 SMP Debian 4.19.152-1 (2020-10-18) x86_64
The programs included with the Debian GNU/Linux system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/*/copyright.
Debian GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent
permitted by applicable law.
```

- 角色型存取控制（**RBAC**）。ONTAP 有兩種RBAC控制項與VMware Tools相關聯：
  - 原生vCenter Server權限
  - vCenter外掛程式特定權限。如需詳細資訊、請參閱 ["此連結"](#)。
- \*加密的通訊通道。\*所有外部通訊都是透過使用TLS 1.2版的HTTPS進行。
- \*最小的連接埠曝光。\*只有必要的連接埠會在防火牆上開啟。

下表說明開放連接埠的詳細資料。

TCP v4/v6連接埠#	方向	功能
8143.	傳入	用於REST API的HTTPS連線
8043.	傳入	HTTPS連線
9060	傳入	HTTPS連線 用於透過 https 連線的 SOAP 此連接埠必須開啟、才能讓用戶端 連線至 ONTAP 工具 API 伺服器。
22	傳入	SSH（預設為停用）
9080	傳入	HTTPS連線- VP和SRA -僅從回送進行內部連線
9083.	傳入	HTTPS 連線： VP 與 SRA 用於透過 https 連線的 SOAP

TCP v4/v6連接埠#	方向	功能
1162	傳入	VP SNMP設陷封包
1527.	僅限內部使用	僅在此電腦與本身之間的外部連線不接受（僅限內部連線）
443..	雙向	用於連線ONTAP 至叢集

- 支援憑證授權單位（CA）簽署的憑證。ONTAP VMware vSphere的各種工具支援CA簽署的憑證。請參閱 "[知識庫文章](#)" 以取得更多資訊。
- \*稽核記錄。\*您可以下載支援套裝組合、而且內容極為詳細。使用者登入和登出活動會記錄在個別的記錄檔中。ONTAP VASA API呼叫會記錄在專屬的VASA稽核記錄（本機CXF.log）中。
- \*密碼原則。\*遵循下列密碼原則：
  - 密碼不會記錄在任何記錄檔中。
  - 密碼不會以純文字形式傳達。
  - 密碼是在安裝程序本身期間設定的。
  - 密碼歷程記錄是可設定的參數。
  - 密碼最短使用期限設為24小時。
  - 密碼欄位的自動完成功能已停用。
  - 利用SHA256雜湊功能、將所有儲存的認證資訊加密。ONTAP

## SnapCenter 外掛程式 VMware vSphere

適用於VMware vSphere軟體工程的NetApp SnapCenter 支援外掛程式使用下列安全開發活動：

- \*威脅建模。\*威脅建模的目的是在軟體開發生命週期初期、發現某項功能、元件或產品的安全性瑕疵。威脅模式是對影響應用程式安全性的所有資訊的結構化呈現。基本上、它是透過安全性觀點來檢視應用程式及其環境。
- \*動態應用程式安全性測試（dast）。\*專為偵測應用程式執行狀態中的易受影響狀況而設計的技術。Dast會測試開放Web應用程式的公開HTTP和HTML介面。
- \*協力廠商程式碼貨幣。\*在開發軟體及使用開放原始碼軟體（開放原始碼軟體）的過程中、必須解決與產品整合的開放原始碼軟體（開放原始碼軟體）相關的安全性弱點。這是一項持續努力、因為開放源碼軟體會元件的版本可能隨時報告新發現的弱點。
- \*弱點掃描。\*弱點掃描的目的是在NetApp產品中發現常見且已知的安全性弱點之後、再將弱點發佈給客戶。
- \*滲透測試。\*滲透測試是評估系統、Web應用程式或網路以找出攻擊者可能利用的安全性弱點的程序。NetApp的滲透測試（筆測試）是由一群獲核准且值得信賴的第三方公司進行。其測試範圍包括利用精密的利用方法或工具、對惡意入侵者或駭客等應用程式或軟體發動攻擊。
- \*產品安全性事件回應活動。\*公司內部和外部都發現安全性弱點、如果未及時解決、可能會對 NetApp 的聲譽造成嚴重風險。為了推動此程序、產品安全性事件回應團隊（PSIRT）會報告並追蹤弱點。

### 產品安全功能

適用於VMware vSphere的NetApp SnapCenter VMware vCenter外掛程式在每個版本中都包含下列安全功能：

- 受限的Shell存取。SSH預設為停用、且只有在從VM主控台啟用時、才允許一次性登入。
- \*登入橫幅中的存取警告。\*使用者在登入提示中輸入使用者名稱後、會顯示下列登入橫幅：

\*警告：\*禁止未經授權存取本系統、並依法律予以起訴。存取本系統即表示您同意、若懷疑有未獲授權的使用情形、您的行動可能受到監控。

使用者透過SSH通道完成登入後、會顯示下列輸出：

```
Linux vsc1 4.19.0-12-amd64 #1 SMP Debian 4.19.152-1 (2020-10-18) x86_64
The programs included with the Debian GNU/Linux system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/*/copyright.
Debian GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent
permitted by applicable law.
```

- 角色型存取控制（RBAC）。ONTAP 有兩種RBAC控制項與VMware Tools相關聯：
  - 原生vCenter Server權限。
  - VMware vCenter外掛程式特定權限。如需詳細資訊、請參閱 "[角色型存取控制（RBAC）](#)"。
- \*加密的通訊通道。\*所有外部通訊都是使用TLS透過HTTPS進行。
- \*最小的連接埠曝光。\*只有必要的連接埠會在防火牆上開啟。

下表提供開放連接埠詳細資料。

TCP v4/v6連接埠號碼	功能
8144.	用於REST API的HTTPS連線
8080	用於OVA GUI的HTTPS連線
22	SSH（預設為停用）
3306.	MySQL（僅限內部連線；預設為停用外部連線）
443..	Ngin像（資料保護服務）

- 支援憑證授權單位（CA）簽署的憑證。SnapCenter VMware vSphere的支援外掛程式支援CA簽署憑證的功能。請參閱 "[如何建立及/或將SSL憑證匯入SnapCenter VMware vSphere（選擇控制器）的VMware外掛程式](#)"。
- \*密碼原則。\*下列密碼原則有效：
  - 密碼不會記錄在任何記錄檔中。
  - 密碼不會以純文字形式傳達。
  - 密碼是在安裝程序本身期間設定的。
  - 所有認證資訊均使用SHA256雜湊來儲存。
- 基本作業系統映像。\*本產品隨附適用於OVA的Debian基礎作業系統、存取受限且停用Shell存取。如此可減少攻擊佔用空間。每SnapCenter 個發行版基礎作業系統都會更新最新的安全修補程式、以達到最大的安全覆蓋範圍。

NetApp針對SnapCenter VMware vSphere應用裝置開發有關VMware vSphere外掛程式的軟體功能與安全性修補程式、然後將其作為套裝軟體平台發佈給客戶。由於這些應用裝置包括特定的Linux子作業系統相依性、以及我們的專屬軟體、因此NetApp建議您不要變更子作業系統、因為這會對NetApp應用裝置造成重大影響。這可能會影響NetApp支援應用裝置的能力。NetApp建議測試及部署我們最新的應用裝置程式碼版本、因為這些版本已發行以修補任何與安全性相關的問題。

## 適用於 VMware vSphere 的 ONTAP 工具安全性強化指南

### 適用於 VMware vSphere 的 ONTAP 工具安全性強化指南

適用於 VMware vSphere 的 ONTAP 工具安全性強化指南提供一套完整的指示、可協助您設定最安全的設定。

這些指南同時適用於應用程式和應用裝置本身的客體作業系統。

### 驗證 VMware vSphere 安裝套件的 ONTAP 工具完整性

有兩種方法可供客戶驗證其 ONTAP 工具安裝套件的完整性。

1. 驗證校驗和
2. 驗證簽名

OTV 安裝套件的下載頁面提供校驗和。使用者必須根據下載頁面所提供的 Checksum 來驗證下載套件的總和。

### 驗證 ONTAP 工具 OVA 的簽名

vApp 安裝套件以 tarball 的形式提供。此 tarball 包含虛擬應用裝置的中繼和根憑證、以及 README 檔案和 OVA 套件。README 檔案可引導使用者驗證 vApp OVA 套件的完整性。

客戶也必須在 vCenter 7.0U3E 版及更新版本上傳所提供的根憑證和中介憑證。對於 7.0.1 與 7.0.U3E 之間的 vCenter 版本、VMware 不支援驗證憑證的功能。客戶不需要上傳任何 vCenter 6.x 版的憑證

### 將信任的根憑證上傳至 vCenter

1. 使用 VMware vSphere Client 登入 vCenter Server 。
2. 指定管理員 @vspece.pengil 或 vCenter 單一登入管理員群組的其他成員的使用者名稱和密碼。如果您在安裝期間指定不同的網域、請以管理員 @ mydomain.
3. 瀏覽至「憑證管理」使用者介面： a.從主選單中、選取管理。B.按一下 [ 憑證 ] 底下的 [ 憑證管理 ] 。
4. 如果系統提示您、請輸入 vCenter Server 的認證。
5. 按一下 [ 信任的根憑證 ] 底下的 [ 新增 ] 。
6. 按一下瀏覽並選取憑證 .pem 檔案 ( OTV\_OVa\_INT\_ROOT\_CERT\_CHERC.pem ) 的位置。
7. 按一下「新增」憑證即會新增至儲存區。

請參閱 ["將信任的根憑證新增至憑證存放區"](#) 以取得更多資訊。部署 VApp (使用 OVA 檔案) 時、可在「Review details」(檢閱詳細資料) 頁面上驗證 vApp 套件的數位簽章。如果下載的 VApp 套件為正版、「發行者」欄會顯示「信任的憑證」(如下面的螢幕擷取畫面所示)。

## Deploy OVF Template

- ✓ 1 Select an OVF template
- ✓ 2 Select a name and folder
- ✓ 3 Select a compute resource
- 4 Review details**
- 5 License agreements
- 6 Select storage
- 7 Select networks
- 8 Customize template
- 9 Ready to complete

### Review details

Verify the template details.

Publisher	Entrust Code Signing CA - OVCS2 (Trusted certificate)
Product	Virtual Appliance - NetApp Inc. ONTAP tools for VMware vSphere
Version	See appliance for version
Vendor	NetApp Inc.
Description	Virtual Appliance - NetApp Inc. ONTAP tools for VMware vSphere for netapp storage systems. For more information or support please visit <a href="https://www.netapp.com/">https://www.netapp.com/</a>
Download size	2.2 GB
Size on disk	3.9 GB (thin provisioned) 53.0 GB (thick provisioned)

CANCEL

BACK

NEXT

Activate  
Go to Sys

### 驗證工具 ISO 和 ONTAP tar.gz 的簽名

NetApp 會在產品下載頁面上與客戶共用程式碼簽署憑證、以及適用於 OTV-ISO 和 SRA.tgz 的產品 zip 檔案。

從程式碼簽署憑證中、使用者可以擷取公開金鑰、如下所示：

```
#> openssl x509 -in <code-sign-cert, pem file> -pubkey -noout > <public-key name>
```

接著應使用公開金鑰來驗證 ISO 和 tgz 產品 zip 的簽名、如下所示：

```
#> openssl dgst -sha256 -verify <public-key> -signature <signature-file>  
<binary-name>
```

範例：

```
#> openssl x509 -in OTV_ISO_CERT.pem -pubkey -noout > OTV_ISO.pub
#> openssl dgst -sha256 -verify OTV_ISO.pub -signature netapp-ontap-tools-
for-vmware-vmware-9.12-upgrade-iso.sig netapp-ontap-tools-for-vmware-
vsphere-9.12-upgrade.iso
Verified OK => response
```

## 連接埠與傳輸協定

此處列出的必要連接埠和通訊協定、可讓 VMware vSphere 伺服器的 ONTAP 工具與其他實體（例如託管儲存系統、伺服器和其他元件）之間進行通訊。

### OTV 所需的傳入和傳出連接埠

請注意下表列出正確運作 ONTAP 工具所需的輸入和輸出連接埠。請務必確保只開啟表中所述的連接埠、以進行遠端機器的連線、而所有其他連接埠則應封鎖、以進行遠端機器的連線。這將有助於確保系統的安全性。

下表說明開放連接埠的詳細資料。

TCP v4/v6 連接埠 #	方向	* 功能 *
8143.	傳入	用於REST API的HTTPS連線
8043.	傳入	HTTPS連線
9060	傳入	HTTPS 連線 用於透過 HTTPS 連線的 SOAP 此連接埠必須開啟、才能讓用戶端連線至 ONTAP 工具 API 伺服器。
22	傳入	SSH（預設為停用）
9080	傳入	HTTPS連線- VP和SRA -僅從回送進行內部連線
9083.	傳入	HTTPS 連線： VP 與 SRA 用於透過 HTTPS 連線的 SOAP
1162	傳入	VP SNMP設陷封包
8443	傳入	遠端外掛程式
1527.	僅限內部使用	只有在此電腦和本身之間才有 Derby 資料庫連接埠、不接受外部連線—僅限內部連線
8150	僅限內部使用	記錄完整性服務會在連接埠上執行
443..	雙向	用於連線ONTAP 至叢集

### 控制對 Derby 資料庫的遠端存取

系統管理員可以使用下列命令來存取 derby 資料庫。您可以透過 ONTAP 工具本機 VM 以及遠端伺服器來存取它、步驟如下：

```
java -classpath "/opt/netapp/vpserver/lib/*" org.apache.derby.tools.ij;
connect 'jdbc:derby://<OTV-
IP>:1527//opt/netapp/vpserver/vvoldb;user=<user>;password=<password>';
```

**[.Underline] example:**

```
root@UnifiedVSC:~# java -classpath "/opt/netapp/vpserver/lib/*" org.apache.derby.tools.ij;
ij version 10.15
ij> connect 'jdbc:derby://localhost:1527//opt/netapp/vpserver/vvoldb;user=app;password=
ij> show tables;
TABLE_SCHEM      |TABLE_NAME      |REMARKS
-----|-----|-----
SYS              |SYSALIASES      |
SYS              |SYSCHECKS       |
SYS              |SYSCOLPERMS     |
SYS              |SYSCOLUMNS     |
SYS              |SYSCONGLOMERATES|
SYS              |SYSCONSTRAINTS |
SYS              |SYSDEPENDS      |
SYS              |SYSFILES        |
SYS              |SYSFOREIGNKEYS  |
SYS              |SYSKEYS         |
SYS              |SYSPERMS       |
```

## 適用於 VMware vSphere 存取點的 ONTAP 工具（使用者）

ONTAP Tools for VMware vSphere 安裝會建立並使用三種類型的使用者：

1. 系統使用者：root 使用者帳戶
2. 應用程式使用者：系統管理員使用者、主要使用者及資料庫使用者帳戶
3. 支援使用者：診斷使用者帳戶

### 1. 系統使用者

系統（root）使用者是由安裝在基礎作業系統（Debian）上的 ONTAP 工具所建立。

- 預設的系統使用者「root」是由 ONTAP 工具安裝在 Debian 上建立的。其預設值為停用、可透過「Maint」主控台在特定的基礎上啟用。

### 2. 應用程式使用者

應用程式使用者在 ONTAP 工具中會命名為本機使用者。這些是在 ONTAP 工具應用程式中建立的使用者。下表列出應用程式使用者的類型：

使用者	說明
系統管理員使用者	它是在 ONTAP 工具安裝期間建立、使用者在部署 ONTAP 工具時提供認證。使用者可以在「Maint」主控台中變更「密碼」。密碼將在 90 天內過期、使用者預期會變更相同的密碼。
維護使用者	它是在 ONTAP 工具安裝期間建立、使用者在部署 ONTAP 工具時提供認證。使用者可以在「Maint」主控台中變更「密碼」。這是維護使用者、是為了執行維護主控台作業而建立的。

使用者	說明
資料庫使用者	它是在 ONTAP 工具安裝期間建立、使用者在部署 ONTAP 工具時提供認證。使用者可以在「Maint」主控台中變更「密碼」。密碼將在 90 天內過期、使用者預期會變更相同的密碼。

### 3. 支援使用者（診斷使用者）

在 ONTAP 工具安裝期間、系統會建立支援使用者。此使用者可在伺服器發生任何問題或中斷時、用來存取 ONTAP 工具、並收集記錄。根據預設、此使用者已停用、但可透過「Maint」主控台臨時啟用。請務必注意、此使用者將在一段時間後自動停用。

### 相互 TLS（憑證型驗證）

ONTAP 9.7 版及更新版本支援相互 TLS 通訊。從適用於 VMware 的 ONTAP 工具和 vSphere 9.12 開始、系統會使用相互 TLS 與新增的叢集進行通訊（視 ONTAP 版本而定）。

### ONTAP

對於所有先前新增的儲存系統：在升級期間、所有新增的儲存系統都會自動受到信任、而且會設定憑證型驗證機制。

如下面的螢幕擷取畫面所示、叢集設定頁面會顯示為每個叢集設定的相互 TLS（憑證型驗證）狀態。

Name	Type	IP Address	ONTAP Release	Status	Capacity	NFS VAAI	Supported Protocols
CL_st121-vmim-ucs501m_1670870260	Cluster	10.224.05.142	9.12.0	Normal	20.42%		

#### \* 叢集新增 \*

在叢集新增工作流程期間、如果所新增的叢集支援 MTLS、則預設會設定 MTLS。使用者不需要為此進行任何組態。下列螢幕擷取畫面會顯示在叢集新增期間顯示給使用者的畫面。



## Add Storage System

 Any communication between ONTAP tools plug-in and the storage system should be mutually authenticated.

vCenter server 10.224.58.52 ▼

Name or IP address:

Username:

Password:

Port:

443

Advanced options ▲

ONTAP Cluster  
Certificate:

Automatically fetch  Manually upload

CANCEL

ADD

## Add Storage System

 Any communication between ONTAP tools plug-in and the storage system should be mutually authenticated.

vCenter server	10.224.58.52 ▾
Name or IP address:	10.234.85.142
Username:	admin
Password:	.....
Port:	443
Advanced options	>

CANCEL

ADD

## Add Storage System

 Any communication between ONTAP tools plug-in and the storage system should be mutually authenticated.

vCenter server

10.234.85.52 ▼

### Authorize Cluster Certificate

Host 10.234.85.142 has identified itself with a self-signed certificate.

[Show certificate](#)

Do you want to trust this certificate?

NO

YES

CANCEL

ADD

## Authorize Cluster Certificate

Host 10.234.85.142 has identified itself with a self-signed certificate.

[Hide certificate](#)

### Certificate Information

This certificate identifies the 10.234.85.142 host.

#### Issued By

**Name (CN or DN):** C1\_sti21-vsimg-ucs581m\_1678878260

#### Issued To

**Name (CN or DN):** C1\_sti21-vsimg-ucs581m\_1678878260

#### Validity

**Issued On:** 03/15/2023 11:16:06

**Expires On:** 03/14/2024 11:16:06

#### Fingerprint Information

**SHA-1 Fingerprint:** 2C:38:E3:5C:4B:F3:5D:3F:39:C8:CE:4A:8  
2:C1:A6:EE:34:53:A0:F3

**SHA-256 Fingerprint:** 05:0F:FE:CD:B0:C6:FC:6F:EB:8A:FC:86:F  
7:E3:EF:D4:8D:CA:02:92:9B:E1:A4:70:84:  
52:F8:76:98:64:FA:23

Do you want to trust this certificate?

NO

YES

### 叢集編輯

在叢集編輯作業期間、有兩種情況：

- 如果 ONTAP 憑證過期、則使用者必須取得新憑證並上傳憑證。
- 如果 OTV 憑證過期、則使用者可以勾選核取方塊來重新產生該憑證。
  - 產生 ONTAP 的新用戶端憑證 \_

# Modify Storage System

Settings   Provisioning Options

IP address or hostname:  ▼

Port:

Username:

Password:

Upload Certificate (Optional)  [BROWSE](#)

Skip monitoring of this storage system

Generate a new client certificate for ONTAP

CANCEL

OK



## ONTAP 工具 HTTPS 憑證

根據預設、ONTAP 工具會使用在安裝期間自動建立的自我簽署憑證、以確保 HTTPS 存取安全無虞。ONTAP 工具提供下列功能：

### 1. 重新產生 HTTPS 憑證

在 ONTAP 工具安裝期間、會安裝 HTTPS CA 憑證、並將憑證儲存在金鑰庫中。使用者可以選擇透過維護主控台重新產生 HTTPS 憑證。

您可以在 *main* 主控台中存取上述選項、方法是瀏覽至「應用程式組態」→「重新產生憑證」。

## 登入橫幅

使用者在登入提示中輸入使用者名稱後、會顯示下列登入橫幅。請注意、SSH 預設為停用、從 VM 主控台啟用時僅允許一次性登入。

```
WARNING: Unauthorized access to this system is forbidden and will be
prosecuted by law. By accessing this system, you agree that your actions
may be monitored if unauthorized usage is suspected.
```

使用者透過SSH通道完成登入後、會顯示下列文字：

```
Linux UnifiedVSC 5.10.0-21-amd64 #1 SMP Debian 5.10.162-1 (2023-01-21)
x86_64
```

```
The programs included with the Debian GNU/Linux system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/*/copyright.
```

```
Debian GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent
permitted by applicable law.
```

## 閒置逾時

為了防止未經授權的存取、系統會設定閒置逾時、自動登出在使用授權資源期間處於非使用中狀態的使用者。如此可確保只有授權使用者才能存取資源、並協助維護安全性。

- 根據預設、vSphere Client 工作階段會在閒置 120 分鐘後關閉、要求使用者再次登入才能繼續使用用戶端。您可以編輯 `webclient.properties` 檔案來變更逾時值。您可以設定 vSphere Client 的逾時時間 "[設定 vSphere Client 逾時值](#)"
- ONTAP 工具的網路 CLI 工作階段登出時間為 30 分鐘。

## 每位使用者的並行要求上限（網路安全保護：DOS 攻擊）

依預設、每位使用者的並行要求上限為 48 個。ONTAP 工具中的根使用者可以根據其環境需求變更此值。\* 此值不應設為非常高的值、因為它提供了一種機制來防範拒絕服務（DOS）攻擊。\*

使用者可以在 `/opt/NetApp/vscserver/etc/dosfilterParams.json` 檔案中變更並行工作階段的最大數量及其他支援參數。

我們可以使用下列參數來設定篩選器：

- **delayMs**：在考慮所有請求之前，為其提供的延遲（以毫秒為單位）超過了速率限制。給予 -1 即可拒絕要求。
- **THROLMS\_**：異步等待信號量的時間。
- **maxRequestM**：允許執行此要求的時間。

- **ipWhitelist**：以逗號分隔的 IP 位址清單、不會受到速率限制。（這可以是 vCenter、ESXi 和 SRA IP）
- **maxRequestsPerSec**：每秒來自連線的最大要求數。
- 在 `_dosfilterParams` 檔案中的預設值：`*`

```
{
  "delayMs": "-1",
  "throttleMs": "1800000",
  "maxRequestMs": "300000",
  "ipWhitelist": "10.224.58.52",
  "maxRequestsPerSec": "48"
}
```

## 網路時間傳輸協定（NTP）組態

有時、網路時間組態不一致、可能會發生安全問題。請務必確保網路中的所有裝置都有正確的時間設定、以避免發生此類問題。

### \* 虛擬應用裝置 \*

您可以從虛擬應用裝置的維護主控台設定 NTP 伺服器。使用者可以在 `_系統組態_ => _新增 NTP 伺服器_` 選項下新增 NTP 伺服器詳細資料

根據預設、NTP 的服務為 `ntpd`。這是一項舊版服務、在某些情況下、虛擬機器無法順利運作。

### \* Debian\*

在 Debian 上、使用者可以存取 `/etc/ntp.conf` 檔案來取得 NTP 伺服器的詳細資料。

## 密碼原則

首次部署 ONTAP 工具或升級至 9.12 版或更新版本的使用者、必須同時遵循系統管理員和資料庫使用者的強式密碼原則。在部署過程中、系統會提示新使用者輸入密碼。對於升級至 9.12 版或更新版本的瀏覽欄位使用者、維護主控台將提供遵循強式密碼原則的選項。

- 一旦使用者登入主控台、就會對照複雜的規則集來檢查密碼、如果發現未遵循、則會要求使用者重設相同的密碼。
- 密碼預設有效時間為 90 天、75 天之後、使用者會開始收到變更密碼的通知。
- 每個週期都需要設定新密碼、系統不會將最後一個密碼當作新密碼。
- 每當使用者登入主控台時、會在載入主功能表之前、先檢查密碼原則、例如下列螢幕擷取畫面：



```
Maintenance Console : "Netapp ONTAP tools for VMware vSphere"  
Discovered interfaces: eth0 (ENABLED)  
validating password policies
```

- 如果發現未遵循密碼原則或 ONTAP 工具 9.11 或更早版本的升級設定、然後使用者會看到下列畫面來重設密碼：

```
Your Administrator and Database password is expired or does not match password policy:  
-----  
1 ) Change 'administrator' user password  
2 ) Change database password  
  
x ) Exit  
  
Enter your choice: _
```

- 如果使用者嘗試設定弱密碼或再次輸入上一個密碼、則使用者將會看到下列錯誤：

```
Changing password for administrator.  
  
User: administrator  
Enter new password:  
Retype new password:  
  
Password doesn't matches the password policy.  
For security reasons, it is recommended to use a password that is of eight to thirty characters and  
contains a minimum of one upper, one lower, one digit, and one special character.  
  
Enter new password:  
Retype new password:  
Check if new decoder works ?  
New decoder worked successfully  
08-02/23 13:36:53 Your new password must be different  
  
Error updating sra credential file  
  
Press ENTER to continue._
```

# 法律聲明

法律聲明提供版權聲明、商標、專利等存取權限。

## 版權

["https://www.netapp.com/company/legal/copyright/"](https://www.netapp.com/company/legal/copyright/)

## 商標

NetApp、NetApp 標誌及 NetApp 商標頁面上列出的標章均為 NetApp、Inc. 的商標。其他公司與產品名稱可能為其各自所有者的商標。

["https://www.netapp.com/company/legal/trademarks/"](https://www.netapp.com/company/legal/trademarks/)

## 專利

如需最新的 NetApp 擁有專利清單、請參閱：

<https://www.netapp.com/pdf.html?item=/media/11887-patentspage.pdf>

## 隱私權政策

["https://www.netapp.com/company/legal/privacy-policy/"](https://www.netapp.com/company/legal/privacy-policy/)

## 開放原始碼

通知檔案提供有關 NetApp 軟體所使用之協力廠商版權與授權的資訊。

## ONTAP

"ONTAP 9.13.1 注意事項"

"關於此功能的注意事項ONTAP 9.12.1.1"

"關於此功能的注意事項ONTAP 9.12.0"

"ONTAP 9.11.1 注意事項"

"關於本產品的注意事項ONTAP 9.10.1"

"ONTAP 9.10.0 注意事項"

"關於此功能的注意事項ONTAP"

"關於本產品的注意事項ONTAP 9.8"

"關於產品的注意ONTAP 事項9.7"

"關於此功能的注意事項ONTAP"

"關於本產品的注意事項ONTAP"

"關於產品的注意ONTAP 9.4"

"關於本產品的注意事項ONTAP"

"關於此功能的注意事項ONTAP 9.2"

"關於此產品的注意事項ONTAP"

## 適用於 MCC IP 的 ONTAP Mediator

"9.9.1 ONTAP Mediator for MCC IP 通知"

"9.8 ONTAP Mediator for MCC IP 通知"

"9.7 適用於 MCC IP 的 ONTAP Mediator 通知"

## 版權資訊

Copyright © 2024 NetApp, Inc. 版權所有。台灣印製。非經版權所有人事先書面同意，不得將本受版權保護文件的任何部分以任何形式或任何方法（圖形、電子或機械）重製，包括影印、錄影、錄音或儲存至電子檢索系統中。

由 NetApp 版權資料衍伸之軟體必須遵守下列授權和免責聲明：

此軟體以 NETAPP「原樣」提供，不含任何明示或暗示的擔保，包括但不限於有關適售性或特定目的適用性之擔保，特此聲明。於任何情況下，就任何已造成或基於任何理論上責任之直接性、間接性、附隨性、特殊性、懲罰性或衍生性損害（包括但不限於替代商品或服務之採購；使用、資料或利潤上的損失；或企業營運中斷），無論是在使用此軟體時以任何方式所產生的契約、嚴格責任或侵權行為（包括疏忽或其他）等方面，NetApp 概不負責，即使已被告知有前述損害存在之可能性亦然。

NetApp 保留隨時變更本文所述之任何產品的權利，恕不另行通知。NetApp 不承擔因使用本文所述之產品而產生的責任或義務，除非明確經過 NetApp 書面同意。使用或購買此產品並不會在依據任何專利權、商標權或任何其他 NetApp 智慧財產權的情況下轉讓授權。

本手冊所述之產品受到一項（含）以上的美國專利、國外專利或申請中專利所保障。

有限權利說明：政府機關的使用、複製或公開揭露須受 DFARS 252.227-7013（2014 年 2 月）和 FAR 52.227-19（2007 年 12 月）中的「技術資料權利 - 非商業項目」條款 (b)(3) 小段所述之限制。

此處所含屬於商業產品和 / 或商業服務（如 FAR 2.101 所定義）的資料均為 NetApp, Inc. 所有。根據本協議提供的所有 NetApp 技術資料和電腦軟體皆屬於商業性質，並且完全由私人出資開發。美國政府對於該資料具有非專屬、非轉讓、非轉授權、全球性、有限且不可撤銷的使用權限，僅限於美國政府為傳輸此資料所訂合約所允許之範圍，並基於履行該合約之目的方可使用。除非本文另有規定，否則未經 NetApp Inc. 事前書面許可，不得逕行使用、揭露、重製、修改、履行或展示該資料。美國政府授予國防部之許可權利，僅適用於 DFARS 條款 252.227-7015(b)（2014 年 2 月）所述權利。

## 商標資訊

NETAPP、NETAPP 標誌及 <http://www.netapp.com/TM> 所列之標章均為 NetApp, Inc. 的商標。文中所涉及的所有其他公司或產品名稱，均為其各自所有者的商標，不得侵犯。