



NVA-1173 NetApp AIPod 搭配 NVIDIA DGX 系統

NetApp Solutions

NetApp
September 26, 2024

目錄

NVA-1173 NetApp AIPOd 搭配 NVIDIA DGX 系統	1
NVA-1173 NetApp AIPOd 搭配 NVIDIA DGX 系統 - 簡介	1
NVA-1173 NetApp AIPOd 搭配 NVIDIA DGX 系統 - 硬體元件	2
NVA-1173 NetApp AIPOd 搭配 NVIDIA DGX 系統 - 軟體元件	5
NVA-1173 NetApp AIPOd 搭配 NVIDIA DGX H100 系統 - 解決方案架構	8
NVA-1173 NetApp AIPOd 搭配 NVIDIA DGX 系統 - 部署詳細資料	10
NVA-1173 NetApp AIPOd 搭配 NVIDIA DGX 系統 - 解決方案驗證與規模調整指南	18
NVA-1173 NetApp AIPOd 搭配 NVIDIA DGX 系統 - 結論與其他資訊	19

NVA-1173 NetApp AI Pod 搭配 NVIDIA DGX 系統

NVA-1173 NetApp AI Pod 搭配 NVIDIA DGX 系統 - 簡介

POWERED BY



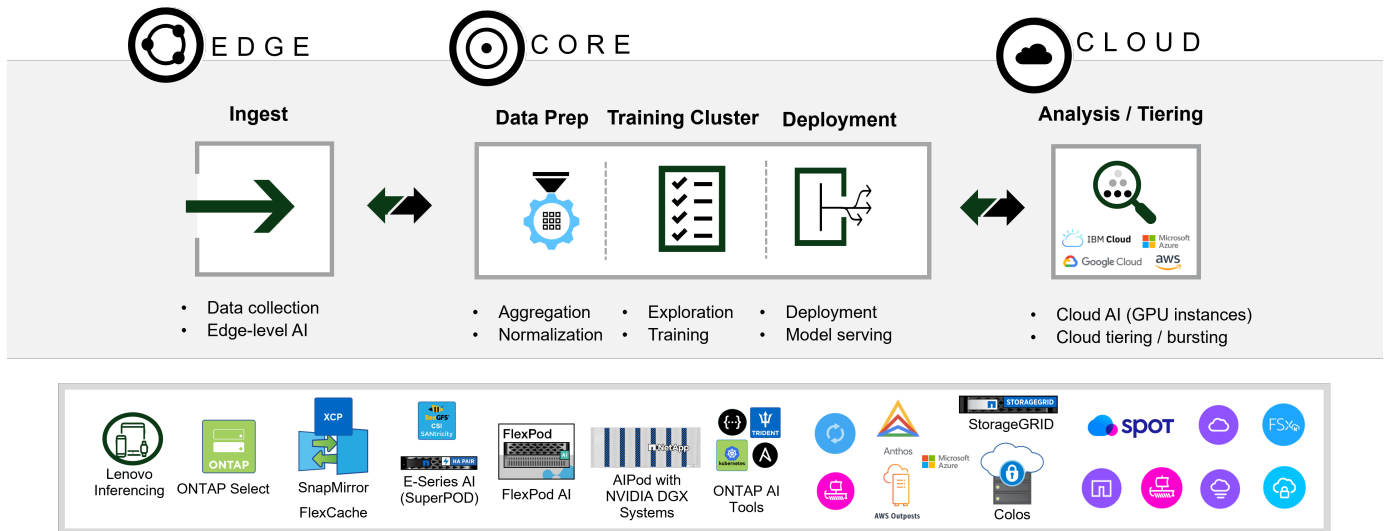
NVIDIA®

NetApp 解決方案工程

執行摘要

NetApp 採用 NVIDIA DGX 的 AI Pod ；系統和 NetApp 雲端連線儲存系統、可消除設計複雜度和猜測、簡化機器學習（ML）和人工智慧（AI）工作負載的基礎架構部署。以 NVIDIA DGX BasePOD 和 NetApp AFF 儲存系統、讓客戶能夠從小規模開始、不中斷地成長、同時智慧地管理從邊緣到核心、再到雲端再回來的資料。NetApp AI Pod 是 NetApp AI 解決方案較大型產品組合的一部分、如下圖所示。

NetApp AI 解決方案產品組合



本文件說明 AI Pod 參考架構的關鍵元件、系統連線能力與組態資訊、驗證測試結果、以及解決方案規模調整指南。本文件適用於有興趣為 ML/DL 和分析工作負載部署高效能基礎架構的 NetApp 和合作夥伴解決方案工程師及客戶策略決策者。

NVA-1173 NetApp AI Pod 搭配 NVIDIA DGX 系統 - 硬體元件

本節重點介紹 NetApp AI Pod 搭配 NVIDIA DGX 系統的硬體元件。

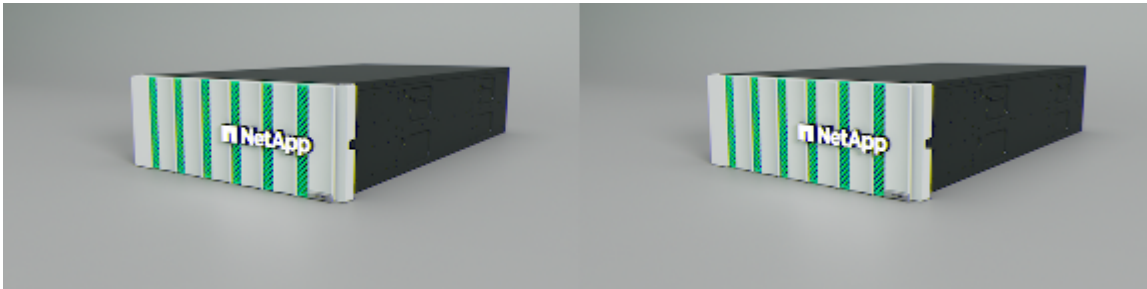
NetApp AFF 儲存系統

NetApp AFF 最先進的儲存系統、讓 IT 部門能夠以領先業界的效能、優異的靈活度、雲端整合、以及同級最佳的資料管理、滿足企業儲存需求。專為Flash而設計AFF 的支援功能、可協助加速、管理及保護業務關鍵資料。

AFF A90 儲存系統

搭載 NetApp ONTAP 資料管理軟體的 NetApp AFF A90 提供內建資料保護、選購的反勒索軟體功能、以及支援最關鍵業務工作負載所需的高效能與恢復能力。它可避免對關鍵任務作業造成中斷、將效能調校降至最低、並保護資料免受勒索軟體攻擊。它提供：
•領先業界的效能
•毫不妥協的資料安全性
•簡化不中斷升級

NetApp AFF A90 儲存系統 _



領先業界的效能

AFF A90 可輕鬆管理深度學習、AI 和高速分析等新一代工作負載、以及 Oracle、SAP HANA、Microsoft SQL Server 和虛擬化應用程式等傳統企業資料庫。它可讓業務關鍵應用程式以最高速度執行、每個 HA 配對最多可達 240 萬 IOPS、延遲低至 100 μ s、效能比先前的 NetApp 機型高達 50%。有了 NFS over RDMA、pNFS 和工作階段 Trunking、客戶就能使用現有的資料中心網路基礎架構、達到新一代應用程式所需的高網路效能。客戶也可以透過統一化的多重傳輸協定支援來擴充 SAN、NAS 和物件儲存設備、並透過統一化的單一 ONTAP 資料管理軟體、在內部部署或雲端中提供最大的靈活度。此外、Active IQ 和 Cloud Insights 也提供 AI 型預測分析功能、可最佳化系統健全狀況。

毫不妥協的資料安全性

AFF A90 系統包含一套完整的 NetApp 整合式與應用程式一致的資料保護軟體。它提供內建的資料保護功能、以及先進的反勒索軟體解決方案、可用於搶佔和攻擊後恢復。惡意檔案可能遭到封鎖、無法寫入磁碟、而且儲存異常狀況也可輕鬆監控、以獲得深入見解。

簡化不中斷升級

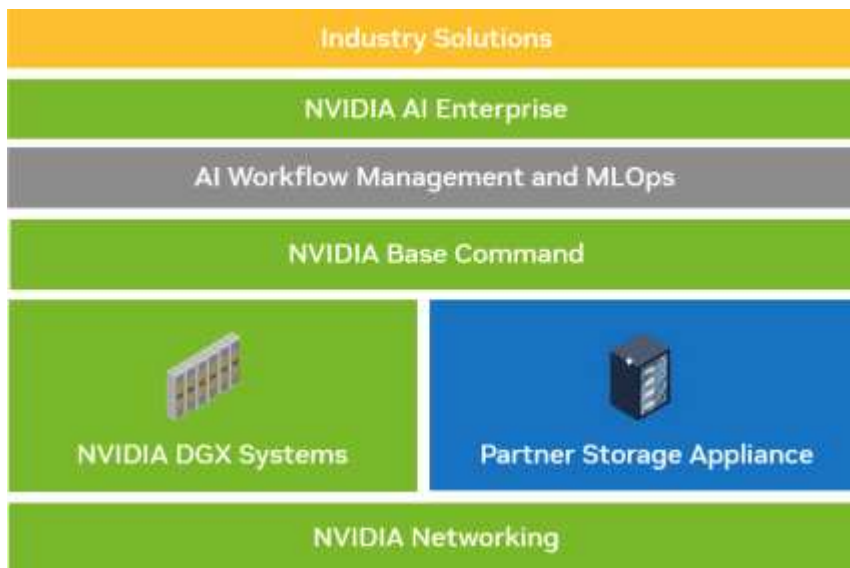
AFF A90 可作為不中斷營運的機箱內升級、升級至現有 A800 客戶。NetApp 透過我們先進的可靠性、可用度、可維修性和管理性（RASM）功能、讓您輕鬆更新並消除對關鍵任務作業的中斷。此外、NetApp 還能進一步提升營運效率、並簡化 IT 團隊的日常活動、因為 ONTAP 軟體會自動為所有系統元件套用韌體更新。

對於規模最大的部署、AFF A1K 系統提供最高的效能和容量選項、而其他 NetApp 儲存系統（例如 AFF A70 和 AFF C800）則提供更低成本的小型部署選項。

NVIDIA DGX基礎POD

NVIDIA DGX BasePOD 是整合式解決方案、包含 NVIDIA 硬體和軟體元件、MLOps 解決方案和協力廠商儲存設備。客戶可運用 NVIDIA 產品和驗證的合作夥伴解決方案的橫向擴充系統設計最佳實務做法、為 AI 開發建置高效率且可管理的平台。圖 1 重點介紹 NVIDIA DGX BasePOD 的各種元件。

NVIDIA DGX BasePOD 解決方案 _



NVIDIA DGX H100 系統

NVIDIA DGX H100 系統；系統是 AI 的強大功能、可透過 NVIDIA H100 Tensor Core GPU 的突破性效能加速。

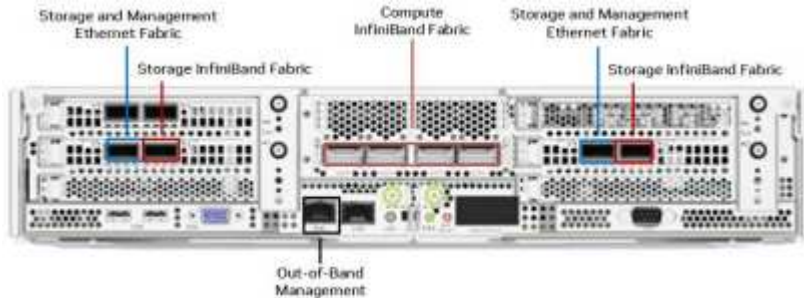
NVIDIA DGX H100 系統 _



DGX H100 系統的主要規格為：•八個 NVIDIA H100 GPU 。•每個 GPU 80 GB GPU 記憶體、總計 640GB 。•四個 NVIDIA NVSwitch™ 晶片。•支援 PCIe 5.0 的雙 56 核心 Intel® Xeon® Platinum 8480 處理器。•2 TB 的

DDR5 系統記憶體。•四個 OSFP 連接埠、可服務八個單埠 NVIDIA ConnectX®-7 (InfiniBand / 乙太網路) 介面卡、以及兩個雙埠 NVIDIA ConnectX-7 (InfiniBand / 乙太網路) 介面卡。•兩個適用於 DGX 作業系統的 1.92 TB M.2 NVMe 磁碟機、八個用於儲存 / 快取的 3.84 TB U.2 NVMe 磁碟機。•最大 10.2 kw 功率。DGX H100 CPU 匣的後端連接埠如下所示。四個 OFP 連接埠可為 InfiniBand 運算架構提供八個 ConnectX-7 介面卡。每對雙連接埠 ConnectX-7 介面卡都提供平行路徑、可通往儲存和管理架構。額外連接埠用於 BMC 存取。

_NVIDIA DGX H100 後面板 _



NVIDIA Networking

NVIDIA Quantum-2 QM9700 交換器

_NVIDIA Quantum-2 QM9700 InfiniBand 交換器 _



NVIDIA Quantum-2 QM9700 交換器搭配 400GB / 秒 InfiniBand 連線能力、可為 NVIDIA Quantum-2 InfiniBand BasePOD 組態中的運算架構提供強大動力。ConnectX-7 單埠介面卡用於 InfiniBand 運算架構。每個 NVIDIA DGX 系統都有兩個連線至每個 QM9700 交換器、可在系統之間提供多個高頻寬、低延遲的路徑。

NVIDIA Spectrum 3 SN4600 交換器

_NVIDIA Spectrum 3 SN4600 交換器 _



NVIDIA Spectrum 和 #8482;-3 SN4600 交換器總共提供 128 個連接埠 (每個交換器 64 個)、以提供備援連線功能、以利 DGX BasePOD 的頻內管理。NVIDIA SN4600 交換器可提供介於 1 GbE 和 200 GbE 之間的速度。對於透過乙太網路連線的儲存設備、也會使用 NVIDIA SN4600 交換器。NVIDIA DGX 雙連接埠 ConnectX-7 介面卡上的連接埠可用於頻內管理和儲存連線。

NVIDIA Spectrum SN2201 交換器

_NVIDIA Spectrum SN2201 交換器 _



NVIDIA Spectrum SN2201 交換器提供 48 個連接埠、可提供額外管理的連線功能。額外管理可為 DGX BasePOD 中的所有元件提供整合式管理連線。

NVIDIA ConnectX-7 介面卡

[_NVIDIA ConnectX-7 介面卡_](#)



NVIDIA ConnectX-7 介面卡可提供 25/50/100/200/400G 的處理量。NVIDIA DGX 系統同時使用單連接埠和雙連接埠 ConnectX-7 介面卡、以使用 400GB InfiniBand 和乙太網路、提供 DGX BasePOD 部署的靈活性。

NVA-1173 NetApp AI Pod 搭配 NVIDIA DGX 系統 - 軟體元件

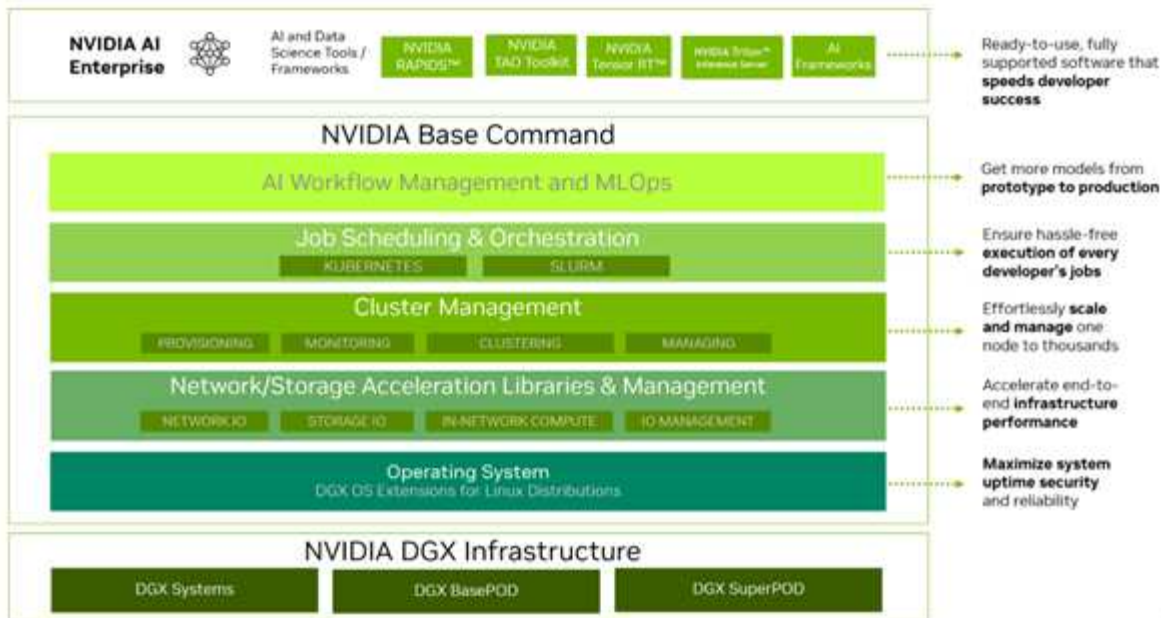
本節重點介紹 NetApp AI Pod 搭配 NVIDIA DGX 系統的軟體元件。

NVIDIA 軟體

NVIDIA Base Command

NVIDIA Base Command[®]；為每個 DGX BasePOD 提供強大動力、讓組織能夠充分發揮 NVIDIA 軟體創新的最大效益。企業可以透過備受肯定的平台、充分發揮投資潛力、包括企業級協調與叢集管理、加速運算、儲存與網路基礎架構的程式庫、以及針對 AI 工作負載最佳化的作業系統（OS）。

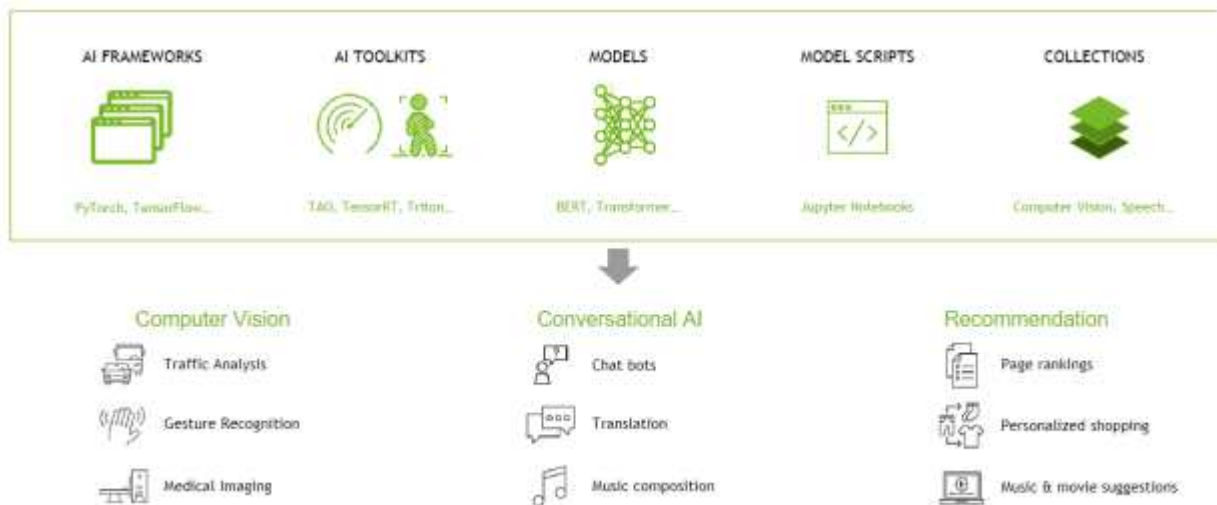
[_NVIDIA BaseCommand 解決方案_](#)



NVIDIA GPU雲端 (NGC)

NVIDIA NGC™ 提供軟體、可滿足具有各種 AI 專業水準的資料科學家、開發人員和研究人員的需求。NGC 軟體會掃描一組彙總的常見弱點與曝險 (CVE)、加密和私密金鑰。它經過測試與設計、可擴充至多個 GPU、在許多情況下、可擴充至多節點、確保使用者在 DGX 系統上的投資發揮最大效益。

NVIDIA GPU Cloud _



NVIDIA AI Enterprise

NVIDIA AI Enterprise 是端點對端軟體平台、可為每個企業提供泛用 AI、為專為在 NVIDIA DGX 平台上執行而最佳化的泛用 AI 基礎模型提供最快且最有效率的執行時間。憑藉生產級的安全性、穩定性和管理能力、它簡化了泛型 AI 解決方案的開發。DGX BasePOD 隨附 NVIDIA AI Enterprise、可讓企業開發人員存取預先訓練的模型、最佳化架構、微服務、加速程式庫、以及企業支援。

NetApp 軟體

NetApp ONTAP

NetApp最新一代的儲存管理軟體、即支援企業將基礎架構現代化、並移轉至雲端就緒的資料中心。ONTAP利用領先業界的資料管理功能ONTAP、無論資料位於何處、只要使用一組工具、即可管理及保護資料。您也可以自由地將資料移至任何需要的位置：邊緣、核心或雲端。支援眾多功能、可簡化資料管理、加速及保護關鍵資料、並在混合雲架構中提供新一代基礎架構功能。ONTAP

加速並保護資料

提供優異的效能與資料保護、並以下列方式擴充這些功能：ONTAP

- 效能與較低的延遲。ONTAP 以最低可能延遲提供最高的處理量、包括支援使用 NFS over RDMA、平行 NFS (pNFS) 和 NFS 工作階段主幹的 NVIDIA GPUDirect 儲存設備 (GDS)。
- 資料保護：ONTAP 提供內建的資料保護功能、以及業界最強大的反勒索軟體保證、並可在所有平台上進行通用管理。
- NetApp Volume Encryption (NVE)。支援內建和外部金鑰管理、提供原生Volume層級的加密功能。ONTAP
- 儲存設備多租戶和多因素驗證。支援以最高安全等級共享基礎架構資源。ONTAP

簡化資料管理

資料管理對於企業IT營運和資料科學家而言至關重要、因此可將適當的資源用於AI應用程式和訓練AI/ML資料集。下列關於NetApp技術的其他資訊超出此驗證範圍、但可能會因您的部署而有所差異。

包含下列功能的資料管理軟體、可簡化及簡化作業、並降低您的總營運成本：ONTAP

- 快照和複製功能可讓 ML/DL 工作流程進行協同作業、平行實驗和強化資料管理。
- SnapMirror 可在混合雲和多站台環境中順暢地移動資料、隨時隨地提供所需的資料。
- 即時資料精簡與擴充重複資料刪除技術。資料壓縮可減少儲存區塊內的空間浪費、重複資料刪除技術可大幅提升有效容量。這適用於本機儲存的資料、以及分層至雲端的資料。
- 最低、最大及可調適的服務品質 (AQO)。精細的服務品質 (QoS) 控制有助於維持高共享環境中關鍵應用程式的效能等級。
- NetApp FlexGroups 可在儲存叢集中的所有節點上散佈資料、為極大型的資料集提供大容量和更高效能。
- NetApp FabricPool自動將冷資料分層至公有和私有雲端儲存選項、包括Amazon Web Services (AWS)、Azure和NetApp StorageGRID 等儲存解決方案。如需FabricPool 更多有關資訊、請參閱 "[TR-4598 : FabricPool 最佳實務做法](#)"。
- NetApp FlexCache。提供遠端磁碟區快取功能、可簡化檔案發佈、減少 WAN 延遲、並降低 WAN 頻寬成本。FlexCache 可在多個站台之間進行分散式產品開發、並可從遠端位置加速存取公司資料集。

符合未來需求的基礎架構

下列功能可協助滿足嚴苛且不斷變化的業務需求：ONTAP

- 無縫擴充和不中斷營運。ONTAP 支援在線上新增現有控制器和橫向擴充叢集的容量。客戶可以升級至最新技術、例如NVMe和32GB FC、而不需進行昂貴的資料移轉或中斷運作。

- 雲端連線：NetApp是最具雲端連線能力的儲存管理軟體、可在所有公有雲中選擇軟體定義儲存（AI）和雲端原生執行個體（NetApp） ONTAP ONTAP Select Cloud Volumes Service 。
- 與新興應用程式整合。利用支援現有企業應用程式的相同基礎架構、為新一代平台和應用程式提供企業級資料服務、例如自動駕駛車輛、智慧城市和產業4.0。ONTAP

NetApp DataOps工具套件

NetApp DataOps Toolkit是一款以Python為基礎的工具、可簡化開發/訓練工作區和推斷伺服器的管理、這些工作區都以高效能橫向擴充的NetApp儲存設備為後盾。DataOps Toolkit 可作為獨立公用程式運作、在 Kubernetes 環境中更有效、利用 NetApp Astra Trident 來自動化儲存作業。主要功能包括：

- 快速配置以高效能橫向擴充NetApp儲存設備為後盾的新高容量JupyterLab工作區。
- 快速配置以企業級NetApp儲存設備為後盾的全新NVIDIA Triton Inference Server執行個體。
- 近乎即時的高容量 JupyterLab 工作區複製、以便進行實驗或快速迭代。
- 高容量 JupyterLab 工作區的近乎即時的快照、用於備份及 / 或可追蹤性 / 基準化。
- 近乎即時的高容量高效能資料磁碟區資源配置、複製及快照。

NetApp Astra Trident

Astra Trident是完全受支援的開放原始碼儲存協調工具、適用於容器和Kubernetes配送、包括Anthos。Trident可搭配整個 NetApp 儲存產品組合使用、包括 NetApp ONTAP、也支援 NFS、NVMe / TCP 和 iSCSI 連線。Trident可讓終端使用者從NetApp儲存系統配置及管理儲存設備、而無需儲存管理員介入、進而加速DevOps 工作流程。

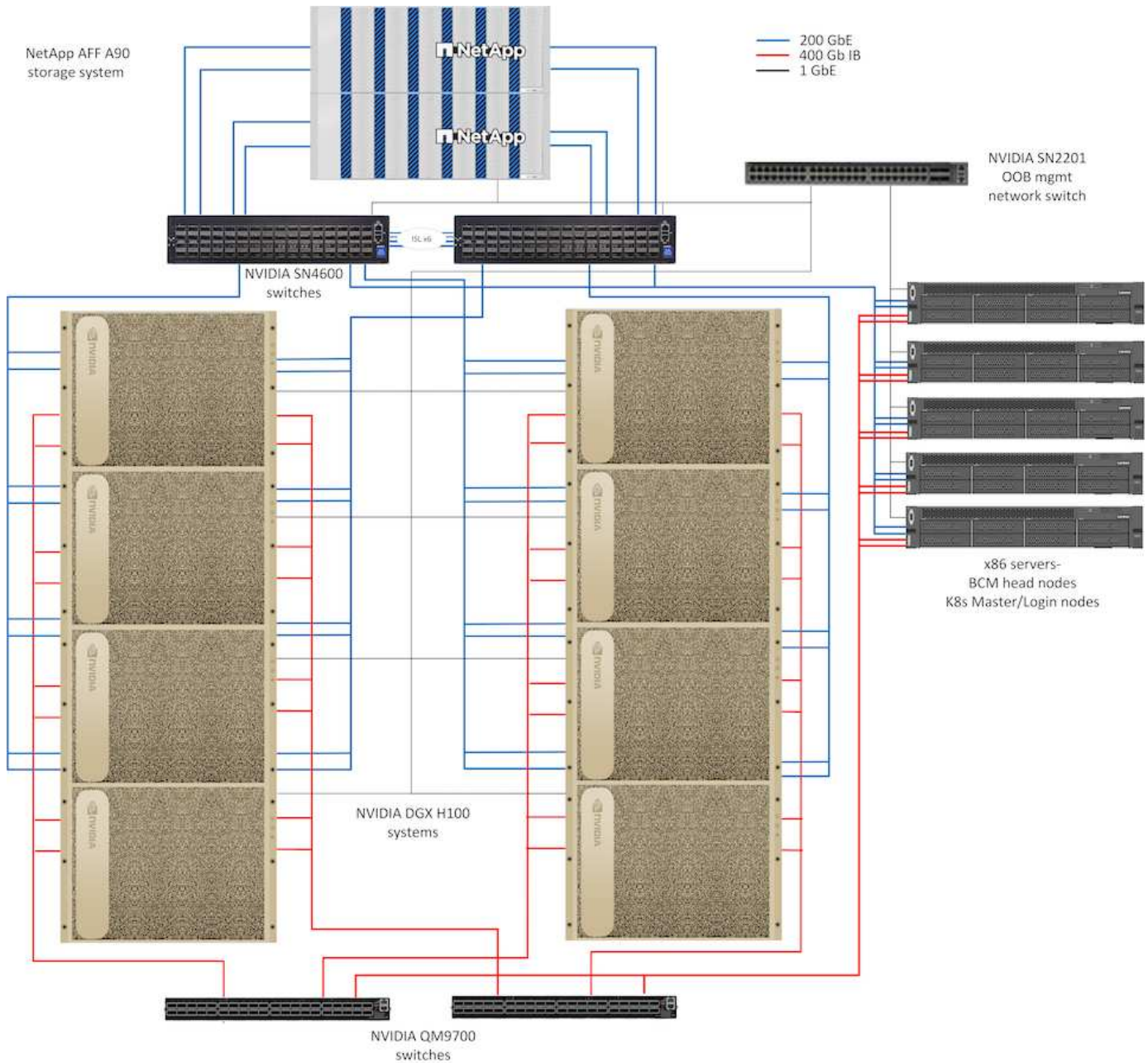
NVA-1173 NetApp AIPod 搭配 NVIDIA DGX H100 系統 - 解決方案架構

本節著重於採用 NVIDIA DGX 系統的 NetApp AIPod 架構。

NetApp AIPod 搭配 DGX 系統

此參考架構運用獨立的架構來進行運算叢集互連和儲存存取、並在運算節點之間提供 400GB / 秒 InfiniBand（IB）連線能力。下圖顯示 NetApp AIPod 搭配 DGX H100 系統的整體解決方案拓撲。

NetApp AIPod 解決方案拓撲 _



網路設計

在此組態中、運算叢集架構使用一對 QM9700 400GB / 秒 IB 交換器、這些交換器會連接在一起以獲得高可用度。每個 DGX H100 系統都使用八個連線連接至交換器、其中偶數連接埠連接至一台交換器、而奇數連接埠則連接至另一台交換器。

對於儲存系統存取、頻內管理和用戶端存取、會使用一對 SN4600 乙太網路交換器。交換器是透過交換器間的連結連接、並設定多個 VLAN 來隔離各種流量類型。在特定 VLAN 之間啟用基本 L3 路由、可在同一台交換器的用戶端和儲存介面之間、以及在交換器之間啟用多條路徑、以實現高可用度。對於較大型的部署、可將乙太網路擴充至葉脊組態、為脊椎交換器新增額外的交換器配對、並視需要新增額外的葉片。

除了運算互連和高速乙太網路之外、所有實體裝置也會連線至一或多個 SN2201 乙太網路交換器、以進行頻外管理。["部署詳細資料"](#)如需網路組態的詳細資訊、請參閱頁面。

DGX H100 系統的儲存存取總覽

每個 DGX H100 系統均配置兩個雙連接埠 ConnectX-7 介面卡、用於管理和儲存流量、而此解決方案則將每個介面卡上的兩個連接埠連接到同一個交換器。接著、每個卡的一個連接埠會設定為 LACP MLAG 連結、每個交換器都有一個連接埠、而用於頻內管理、用戶端存取和使用者層級儲存存取的 VLAN 則會裝載在此連結上。

每張卡上的另一個連接埠用於連線至 AFF A90 儲存系統、並可用於數種組態、視工作負載需求而定。對於使用 NFS over RDMA 來支援 NVIDIA Magnum IO GPUDirect 儲存設備的組態、連接埠會個別用於不同 VLAN 中的 IP 位址。對於不需要 RDMA 的部署、儲存介面也可以設定為 LACP 繫結、以提供高可用度和額外頻寬。無論是否使用 RDMA、用戶端都可以使用 NFS v4.1 pNFS 和工作階段主幹來掛載儲存系統、以便平行存取叢集中的所有儲存節點。["部署詳細資料"](#)如需用戶端組態的詳細資訊、請參閱頁面。

如需 DGX H100 系統連線的詳細資訊["NVIDIA BasePOD 文件"](#)、請參閱。

儲存系統設計

每個 AFF A90 儲存系統都使用每個控制器的六個 200 GbE 連接埠進行連線。每個控制器有四個連接埠用於從 DGX 系統存取工作負載資料、每個控制器有兩個連接埠則設定為 LACP 介面群組、以支援從管理層伺服器存取叢集管理成品和使用者主目錄。儲存系統的所有資料存取都是透過 NFS 提供、其中儲存虛擬機器 (SVM) 專用於 AI 工作負載存取、另有專用於叢集管理用途的獨立 SVM。

["部署詳細資料"](#)如需儲存系統組態的詳細資訊、請參閱頁面。

管理層伺服器

此參考架構也包含五部以 CPU 為基礎的伺服器、供管理層使用。其中兩個系統是 NVIDIA Base Command Manager 的主要節點、用於叢集部署和管理。其餘三個系統則用於提供額外的叢集服務、例如 Kubernetes 主節點或登入節點、以便使用 Slurm 進行工作排程。使用 Kubernetes 的部署可運用 NetApp Astra Trident CSI 驅動程式、為 AFF A900 儲存系統上的管理和 AI 工作負載提供自動化的資源配置和資料服務、

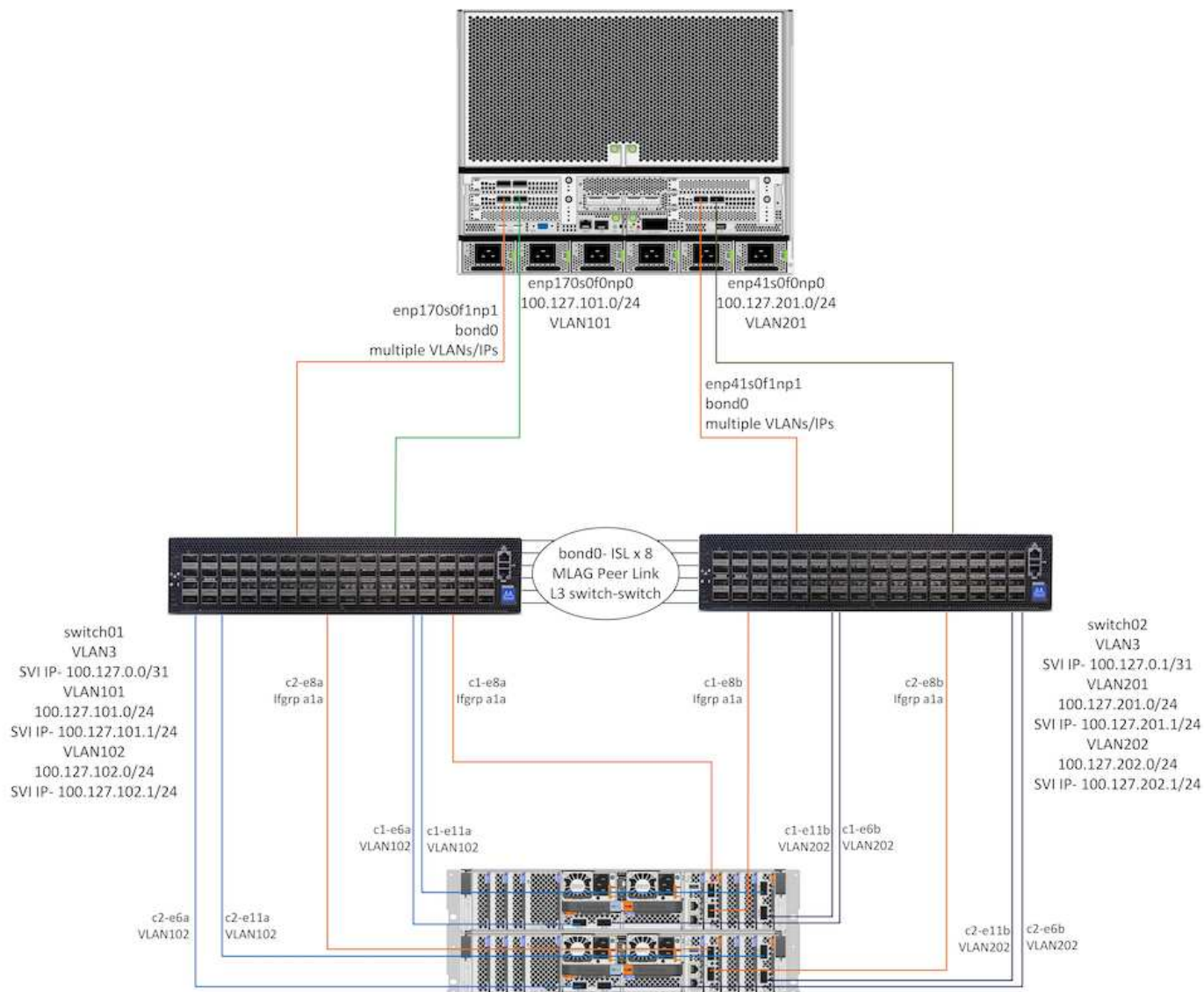
每部伺服器都會實體連接至 IB 交換器和乙太網路交換器、以啟用叢集部署和管理、並透過管理 SVM 將 NFS 裝載至儲存系統、以儲存叢集管理產出工件、如前所述。

NVA-1173 NetApp AIPod 搭配 NVIDIA DGX 系統 - 部署詳細資料

本節說明驗證此解決方案時所使用的部署詳細資料。所使用的 IP 位址是範例、應根據部署環境進行修改。如需實作此組態時所使用之特定命令的詳細資訊、請參閱適當的產品文件。

下圖顯示 1 個 DGX H100 系統和 1 個 HA AFF A90 控制器配對的詳細網路和連線資訊。以下各節中的部署指南是根據此圖表中的詳細資料而定。

NetApp AIPod 網路組態 _



下表顯示最多 16 個 DGX 系統和 2 個 AFF A90 HA 配對的佈線範例。

交換器與連接埠	裝置	裝置連接埠
Switch1 連接埠 1-16	DGX-H100-01 至 -16	enp170s0f0np0 、 SLOT1 連接埠 1
Switch1 連接埠 17-32	DGX-H100-01 至 -16	enp170s0f1np1 、 SLOT1 連接埠 2
Switch1 連接埠 33-36	AFF A90-01 到 -04	連接埠 e6a
Switch1 連接埠 37-40	AFF A90-01 到 -04	連接埠 e11a
Switch1 連接埠 41-44	AFF A90-01 到 -04	連接埠 e8a
Switch1 連接埠 57-64	ISL 到交換器 2	連接埠 57-64
Switch2 連接埠 1-16	DGX-H100-01 至 -16	enp41s0f0np0 、 插槽 2 連接埠 1
Switch2 連接埠 17-32	DGX-H100-01 至 -16	enp41s0f1np1 、 插槽 2 連接埠 2
Switch2 連接埠 33-36	AFF A90-01 到 -04	連接埠 e6b.
Switch2 連接埠 37-40	AFF A90-01 到 -04	連接埠 e11b.

交換器與連接埠	裝置	裝置連接埠
Switch2 連接埠 41-44	AFF A90-01 到 -04	連接埠 e8b.
Switch2 連接埠 57-64	ISL 到交換器 1	連接埠 57-64

下表顯示此驗證所使用之各種元件的軟體版本。

裝置	軟體版本
NVIDIA SN4600交換器	Cumulus Linux v5.9.1
NVIDIA DGX 系統	DGX OS 6.2.1 版 (Ubuntu 22.04 LTS)
Mellanox OFED	24.01
NetApp AFF 產品系列A90	NetApp ONTAP 9.14.1

儲存網路組態

本節概述乙太網路儲存網路組態的主要詳細資料。如需設定 InfiniBand 運算網路的相關資訊，請參閱"[NVIDIA BasePOD 文件](#)"。如需更多關於交換器組態"[NVIDIA Cumulus Linux 文件](#)"的詳細資訊、請參閱。

下面概述了用於配置 SN4600 交換機的基本步驟。此程序假設纜線和基本交換器設定（管理 IP 位址、授權等）已完成。

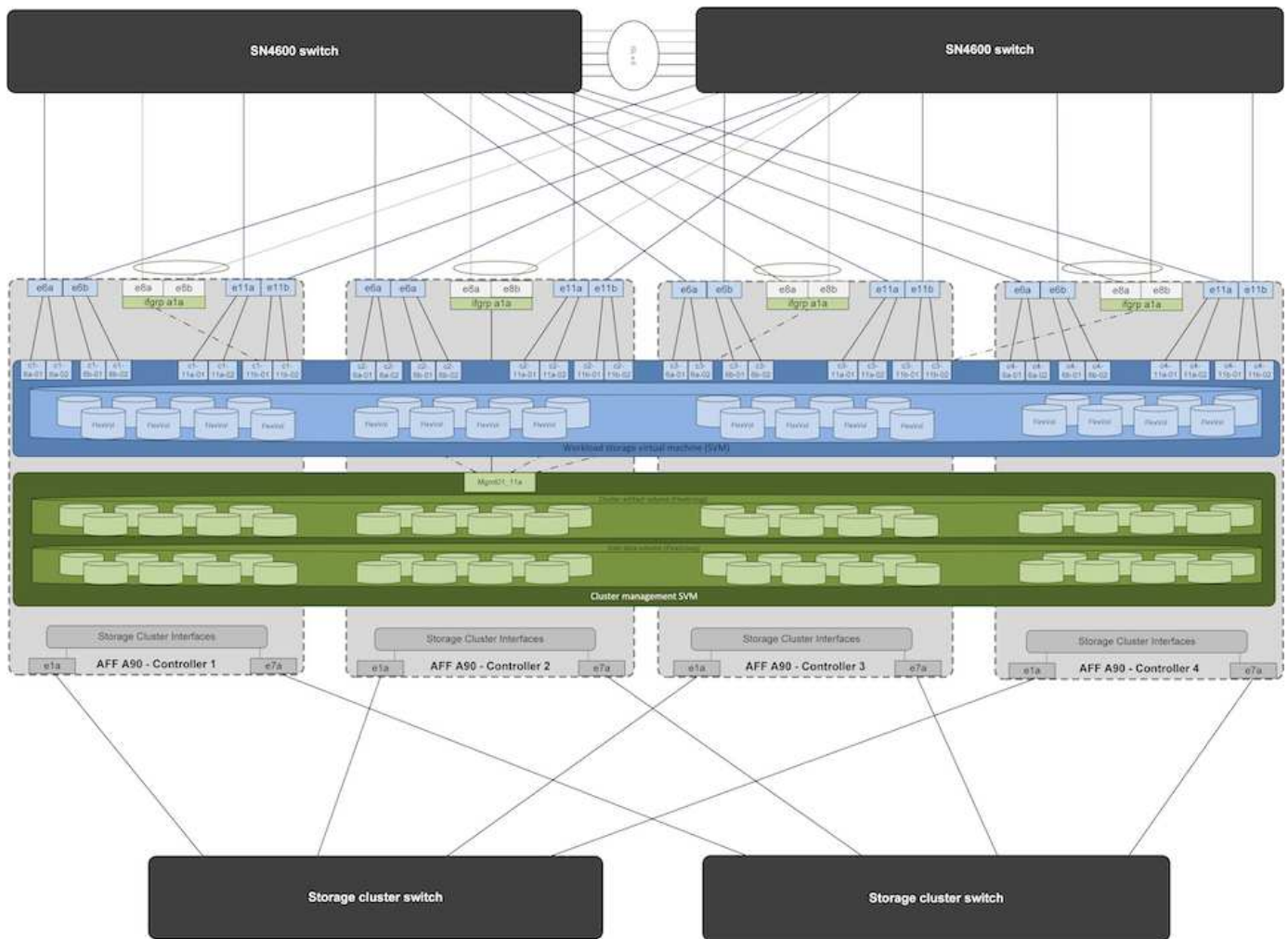
- 設定交換器之間的 ISL 連結、以啟用多連結集合體（ MLAG ）和容錯移轉流量
 - 這項驗證使用 8 個連結、為測試中的儲存組態提供足夠的頻寬
 - 如需啟用 MLAG 的特定指示、請參閱 Cumulus Linux 文件。
- 為兩台交換器上的每對用戶端連接埠和儲存連接埠設定 LACP MLAG
 - DGX-H100-01 （ enp170s0f1np1 和 enp41s0f1np1 ） 、 DGX-H100-02 等各交換器的連接埠 swp17 （ bond1-16 ）
 - 每個交換器上的連接埠 swp41 適用於 AFF A90-01 （ e8a 和 e8b ） 、連接埠 swp42 適用於 AFF A90-02 等 （ bond17-20 ）
 - NV Set 介面 bondX bond 成員 swpX
 - NV Set 介面 bondx bond MLAG id X
- 將所有連接埠和 MLAG 連結新增至預設橋接網域
 - NV 設定 int swp1-1633-40 橋接網域 br_default
 - NV 設定 int bond1-20 橋接網域 br_default
- 在每台交換器上啟用 roce
 - NV 設定無損模式
- 設定 VLANs - 2 用於用戶端連接埠、 2 用於儲存連接埠、 1 用於管理、 1 用於 L3 交換器至交換器
 - 交換器 1-
 - 當用戶端 NIC 發生故障時、用於 L3 交換器的 VLAN 3 交換器路由
 - 每個 DGX 系統上儲存連接埠 1 的 VLAN 101 （ enp170s0f0np0 、 SLOT1 連接埠 1 ）

- 每個 AFF A90 儲存控制器上連接埠 e6a 和 e11a 的 VLAN 102
 - 使用 MLAG 介面管理每個 DGX 系統和儲存控制器的 VLAN 301
 - 交換器 2-
 - 當用戶端 NIC 發生故障時、用於 L3 交換器的 VLAN 3 交換器路由
 - VLAN 201 用於每個 DGX 系統上的儲存連接埠 2 (enp41s0f0np0 、 slot2 連接埠 1)
 - 每個 AFF A90 儲存控制器上的 VLAN 202 連接埠 e6b 和 e11b
 - 使用 MLAG 介面管理每個 DGX 系統和儲存控制器的 VLAN 301
6. 視情況將實體連接埠指派給每個 VLAN 、例如用戶端 VLAN 中的用戶端連接埠、以及儲存 VLAN 中的儲存連接埠
- NV 設定 int <swpX> 橋接網域 br_default 存取 <vlan id>
 - MLAG 連接埠應保留為主幹連接埠、以視需要在連結的介面上啟用多個 VLAN 。
7. 在每個 VLAN 上設定交換器虛擬介面 (SVI) 、以做為閘道並啟用 L3 路由
- 交換器 1-
 - NV 設定 int VLAN3 IP 位址 100.127.0/31
 - NV 設定 int vlan101 IP 位址 100.127.101.1/24
 - NV 設定 int vlan102 IP 位址 100.127.102.1/24
 - 交換器 2-
 - NV 設定 int VLAN3 IP 位址 100.127.0.0.1/31
 - NV 設定 int vlan201 IP 位址 100.127.201.1/24
 - NV 設定 int vlan202 IP 位址 100.127.202.1/24
8. 建立靜態路由
- 靜態路由會自動為同一台交換器上的子網路建立
 - 當用戶端連結故障時、交換器到交換器的路由需要額外的靜態路由
 - 交換器 1-
 - NV 將 VRF 預設路由器靜態設為 100.127.128.0/17 、透過 100.127.0.1
 - 交換器 2-
 - NV 將 VRF 預設路由器靜態設定為 100.127.0/17 、透過 100.127.0.0

儲存系統組態

本節說明此解決方案 A90 儲存系統組態的重要詳細資料。如需 ONTAP 系統組態的詳細資訊、請參閱 [ONTAP 說明文件] 。下圖顯示儲存系統的邏輯組態。

NetApp A90 儲存叢集邏輯組態 _



以下概述設定儲存系統的基本步驟。此程序假設已完成基本儲存叢集安裝。

1. 在每個控制器上設定 1 個 Aggregate、所有可用分割區減 1 個備援磁碟區
 - Aggr create -node <node> -Aggregate <node> 資料 a01 -diskcount <47>
2. 在每個控制器上設定 ifgrp
 - NET 連接埠 ifgrp create -node <node> -ifgrp A1A -mode imody_lacp -distr-function 連接埠
 - net 連接埠 ifgrp add-port -node <node> -ifgrp <ifgrp> -ports <node> : e8a、<node> : e8b
3. 在每個控制器上的 ifgrp 上設定管理 VLAN 連接埠
 - net port VLAN create -node AFF — a90-01 - 連接埠 A1A -vlan-id 31
 - net port VLAN create -node AFF — a90-02 — port A1A — vlan — id 31
 - net port VLAN create -node AFF — a90-03 — port A1A — vlan — id 31
 - net port VLAN create -node AFF — a90-04 — port A1A — vlan-id 31
4. 建立廣播網域
 - 廣播網域 create -broadcast-domain VLAN21 -MTU 9000 連接埠 AFF a90-01:e6a、AFF AFF a90-01:e11a、AFF a90-02:e6a、AFF a90-02:e11a、AFF a90-03:e6a、AFF a90-03:e11a、AFF a90-04:e04-e90:e11a
 - 廣播網域 create -broadcast-domain VLAN22 -MTU 9000 連接埠 aaa 穎 90-01:e6b、AFF AFF a90-01

: e11b 、 AFF a90-02 : e6b 、 AFF a90-02 : e11b 、 AFF a90-03 : e6b 、 AFF a90-03 : e11b
、 AFF a90-04-e90

- 廣播網域 create -broadcast-domain vlan31 -MTU 9000 連接埠 AFF a90-01 : A1A-31 、 AFF a90-02 : A1A-31 、 AFF a90-03 : A1A-31 、 AFF a90-04 : A1A-31

5. 建立管理 SVM *

6. 設定管理 SVM

- 建立 LIF
 - net int create -vserver baspoe-mgmt -lif vlan31-01 -home-node-a90-01 AFF -home-port A1A-31 -address 192.168.31.X -netmask 255 · 255 · 255 · 255 · 255 · 0
- 建立 FlexGroup Volume -
 - Vol create -vserver baspoe-mgmt -volume home -size 10T -auto-Provision -as FlexGroup -jite-path /home
 - Vol create -vserver baspoe-mgmt -volume cm -size 10T -auto-Providence-as FlexGroup -jite-path /cm
- 建立匯出原則
 - 匯出原則規則 create -vserver baspoe-mgmt -policy default -client-match 192.168.31.0/24 -rorule sys -rwrule sys -Superuser sys

7. 建立資料 SVM *

8. 設定資料 SVM

- 設定 SVM 以支援 RDMA
 - vserver modify -vserver baspoe-data -RDMA 已啟用
- 建立生命
 - net int create -vserver baspoe-data -lif c1-6a-lif1 -home-node AFF — a90-01 -home-port e6a -address 100.127.102.101 -netmask 255 · 255 · 255 · 255 · 255 · 255 · 255 · 0
 - net int create -vserver baspoe-data -lif c1-6a-lif2 -home-node AFF — a90-01 -home-port e6a -address 100.127.102.102 -netmask 255 · 255 · 255 · 255 · 255 · 255 · 255 · 0
 - net int create -vserver baspoe-data -lif c1-6B-lif1 -home-node-a90-01 AFF -home-port e6b -address 100.127.202.101 -netmask 255 · 255 · 255 · 255 · 255 · 255 · 255 · 0
 - net int create -vserver baspoe-data -lif c1-6B-lif2 -home-node-a90-01 -home-port AFF e6b -address 100.127.202.102 -netmask 255 · 255 · 255 · 255 · 255 · 255 · 255 · 0
 - net int create -vserver baspoe-data -lif c1-11a-lif1 -home-node-a90-01 -home-port AFF e11a -address 100.127.102.103 -netmask 255 · 255 · 255 · 255 · 255 · 255 · 255 · 0
 - net int create -vserver baspoe-data -lif c1-11a-lif2 -home-node-a90-01 -home-port AFF e11a -address 100.127.102.104 -netmask 255 · 255 · 255 · 255 · 255 · 255 · 255 · 0
 - net int create -vserver baspoe-data -lif c1-11b-lif1 -home-node-a90-01 AFF -home-port e11b -address 100.127.202.103 -netmask 255 · 255 · 255 · 255 · 255 · 255 · 255 · 0
 - net int create -vserver baspoe-data -lif c1-11b-lif2 -home-node-a90-01 -home-port AFF e11b -address 100.127.202.104 -netmask 255 · 255 · 255 · 255 · 255 · 255 · 255 · 0
 - net int create -vserver baspoe-data -lif c2-6a-lif1 -home-node-a90-02 AFF -home-port e6a -address 100.127.102.105 -netmask 255 · 255 · 255 · 255 · 255 · 255 · 255 · 0
 - net int create -vserver baspoe-data -lif c2-6a-lif2 -home-node-a90-02 AFF -home-port e6a -address

```
100.127.102.106 -netmask 255 · 255 · 255 · 255 · 255 · 255 · 255 · 0
```

- net int create -vserver baspoe-data -lif c2-6B-lif1 -home-node-a90-02 AFF -home-port e6b -address 100.127.202.105 -netmask 255 · 255 · 255 · 255 · 255 · 255 · 255 · 0
- net int create -vserver baspoe-data -lif c2-6B-lif2 -home-node-a90-02 -home-port AFF e6b -address 100.127.202.106 -netmask 255 · 255 · 255 · 255 · 255 · 255 · 255 · 0
- net int create -vserver baspoe-data -lif c2-11a-lif1 -home-node-a90-02 -home-port AFF e11a -address 100.127.102.107 -netmask 255 · 255 · 255 · 255 · 255 · 255 · 255 · 0
- net int create -vserver baspoe-data -lif c2-11a-lif2 -home-node-a90-02 -home-port AFF e11a -address 100.127.102.108 -netmask 255 · 255 · 255 · 255 · 255 · 255 · 255 · 0
- net int create -vserver baspoe-data -lif c2-11b-lif1 -home-node-a90-02 AFF -home-port e11b -address 100.127.202.107 -netmask 255 · 255 · 255 · 255 · 255 · 255 · 255 · 0
- net int create -vserver baspoe-data -lif c2-11b-lif2 -home-node-a90-02 AFF -home-port e11b -address 100.127.202.108 -netmask 255 · 255 · 255 · 255 · 255 · 255 · 255 · 0

9. 設定 RDMA 存取的生命

- 對於使用 ONTAP 9 · 15.1 的部署、實體資訊的 Roce QoS 組態需要 ONTAP CLI 中無法使用的 OS 層級命令。如需連接埠組態的協助、請聯絡 NetApp 支援部門以取得 ROCE 支援。NFS over RDMA 功能完全沒有問題
- 從 ONTAP 9 · 16.1 開始、實體介面會自動設定適當的設定、以支援端點對端點的設備。
- net int modify -vserver baspoed-data -lif * -rdma-protocols roce

10. 在資料 SVM 上設定 NFS 參數

- NFS modify -vserver baspox-data -v4.1 已啟用 -v4.1-pNFS 已啟用 -v4.1-trunking -tcp-max-transfer-size 262144

11. 建立 FlexGroup Volume -

- Vol create -vserver baspode-data -volume data -size 100T -auto-Provision -as FlexGroup -jite-path /data

12. 建立匯出原則

- 匯出原則規則 create -vserver baspoe-data -policy default -client-match 100.127.101.0/24 -rorule sys -rwRule sys -Superuser sys
- 匯出原則規則 create -vserver baspoe-data -policy default -client-match 100.127.201.0/24 -rorule sys -rwRule sys -Superuser sys

13. 建立航線

- Route add -vserver baspo_data -destination 100.127.0/17 - gateway 100.127.102.1 metric 20
- Route add -vserver baspo_data -destination 100.127.0/17 - 閘道 100.127.202.1 metric 30
- Route add -vserver baspo_data -destination 100.127.128.0/17 - 閘道 100.127.202.1 metric 20
- Route add -vserver baspo_data -destination 100.127.128.0/17 - 閘道 100.127.102.1 metric 30

用於存取 **ROCE** 儲存設備的 **DGX H100** 組態

本節說明 DGX H100 系統組態的重要詳細資料。其中許多組態項目可包含在部署至 DGX 系統的 OS 映像中、或在開機時由 Base Command Manager 實作。此處列出這些項目以供參考、如需在 BCM 中設定節點和軟體映像"BCM 文件"的詳細資訊、請參閱。

1. 安裝其他套件
 - IPMItool
 - python3-pip
2. 安裝 Python 套件
 - 輔助子
 - matplotlib
3. 在套件安裝後重新設定 dpkg
 - dpkg --configure -a
4. 安裝 MOFED
5. 設定效能調校的 mst 值
 - mstconfig -y -d <aa:00.0,29:00.0> set advanced_pci_settings=1 NUM_OF_VFS=0
MAX_ACC_Out_read=44
6. 修改設定後重設介面卡
 - mlxfwreset -d <aa:00.0,29:00.0> -y 重設
7. 在 PCI 裝置上設定 MaxReadReq
 - setpci -s <aa:00.0,29:00.0> 68.W=5957
8. 設定 RX 和 TX 環狀緩衝區大小
 - Ethtool -G <enp170s0f0np0,enp41s0f0np0> Rx 8192 Tx 8192
9. 使用 mlx_QoS 設定 PFC 和 DSCP
 - mlx_qos -i <enp170s0f0np0,enp41s0f0np0> --fc 0 、 0 、 0 、 0 、 -trust = dscp --cable_len=3
10. 將 ToS 設為在網路連接埠上傳輸流量
 - ECHO 106 > /sys/class/InfiniBand / <mlx5_7,mlx5_1> / tc/1/tra流量 類別
11. 在適當的子網路上、使用 IP 位址設定每個儲存 NIC
 - 100 、 127.101.0/24 適用於儲存 NIC 1
 - 100100127.201.0/24 適用於儲存 NIC 2
12. 設定 LACP 繫結的頻內網路連接埠 (enp170s0f1np1 、 enp41s0f1np1)
13. 為通往每個儲存子網路的主要和次要路徑設定靜態路由
 - 新增路由– net 100.127.0/17 GW 100.127.101.1 公制 20
 - 新增路由– net 100.127.0/17 GW 100.127.201.1 公制 30
 - 新增路由– net 100.127.128.0/17 GW 100.127.201.1 公制 20
 - 新增路由– net 100.127.128.0/17 GW 100.127.101.1 公制 30
14. 裝載 /home Volume
 - 掛載 -o ves=3 、 nconnect =16 、 rsize=262144 、 wsize=262144 192.168.31.X : /home /home
15. 裝載 / 資料磁碟區
 - 下列掛載選項是在安裝資料 Volume 時使用的 -

- ves=4.1# 可讓 pNFS 平行存取多個儲存節點
- proto=RDMA # 會將傳輸通訊協定設定為 RDMA、而非預設 TCP
- max_connect = 16# 可讓 NFS 工作階段主幹聚合儲存連接埠頻寬
- 寫入 = 熱切 # 可改善緩衝寫入的寫入效能
- rsize=262144、wsize=262144 # 將 I/O 傳輸大小設為 256k

NVA-1173 NetApp AIPOd 搭配 NVIDIA DGX 系統 - 解決方案驗證與規模調整指南

本節著重於 NetApp AIPOd 搭配 NVIDIA DGX 系統的解決方案驗證與規模調整指南。

解決方案驗證

此解決方案中的儲存組態已使用一系列使用開放原始碼工具 FIO 的綜合工作負載進行驗證。這些測試包括讀寫 I/O 模式、用於模擬 DGX 系統執行深度學習訓練工作所產生的儲存工作負載。儲存組態已通過驗證、使用雙插槽 CPU 伺服器叢集同時執行 FIO 工作負載、以模擬 DGX 系統叢集。每個用戶端都設定了先前所述的相同網路組態、並加入下列詳細資料。

此驗證使用下列掛載選項：

ves=4.1	啟用 pNFS 以平行存取多個儲存節點
proto=RDMA	將傳輸通訊協定設定為 RDMA、而非預設 TCP
連接埠 = 20049	為 RDMA NFS 服務指定正確的連接埠
max_connect = 16	啟用 NFS 工作階段主幹以彙總儲存連接埠頻寬
寫入 = 渴望	改善緩衝寫入的寫入效能
rsize=262144、wsize=262144	將 I/O 傳輸大小設為 256k

此外、用戶端設定的 NFS max_Session_插槽值為 1024。在透過 RDMA 使用 NFS 測試解決方案時、儲存網路連接埠已設定為主動 / 被動連結。此驗證使用下列連結參數：

mode=active-backup	將連結設定為主動 / 被動模式
primary = <interface name>	所有用戶端的主要介面都分散在交換器上
MII-monitor-interval = 100	指定 100ms 的監控時間間隔
容錯移轉 -Mac-policy=active	指定主動鏈的 MAC 地址是綁定的 MAC 地址。這是在連結介面上正確操作 RDMA 所需的。

儲存系統的設定方式如前所述、為兩對 A900 HA（4 個控制器）搭配兩個 NS224 磁碟櫃（每對 HA 連接 24 個 1.9TB NVMe 磁碟機）。如架構一節所述、所有控制器的儲存容量都是使用 FlexGroup 磁碟區進行組合、而來自所有用戶端的資料則分散在叢集中的所有控制器上。

儲存系統規模調整指南

NetApp 已成功完成 DGX BasePOD 認證、兩對通過測試的 A90 HA 可輕鬆支援 16 個 DGX H100 系統的叢

集。對於儲存效能需求較高的大型部署、可在單一叢集中新增最多 12 個 HA 配對（24 個節點）的額外 AFF 系統至 NetApp ONTAP 叢集。使用本解決方案所述的 FlexGroup 技術、24 節點叢集可在單一命名空間中提供 40 PB 以上的資料傳輸量、以及高達 300 Gbps 的傳輸量。其他 NetApp 儲存系統（例如 AFF A400、A250 和 C800）則提供較低的效能和 / 或較高的容量選項、以較低的成本進行較小型的部署。由於 ONTAP 9 支援混合模式叢集、因此客戶可以從較小的初始佔用空間開始、並在容量和效能需求增加時、將更多或更大的儲存系統新增至叢集。下表顯示每個 AFF 機型所支援的 A100 和 H100 GPU 數量的粗略估計值。

NetApp 儲存系統規模調整指南

		Throughput ²	Raw capacity (typical / max)	Connectivity	# NVIDIA A100 GPUs supported ³	# NVIDIA H100 GPUs supported ⁴
NetApp® AFF A900	1 HA pair ¹	28GB/s	182TB / 14.7PB	100 GbE	1 - 64	1-32
	12 HA pairs	336GB/s	2.1PB / 176.4PB		768	384
AFF A800	1 HA pair	25GB/s	368TB / 3.6PB	100 GbE	1 - 64	1-32
	12 HA pairs	300GB/s	4.4PB / 43.2PB		768	384
AFF C800	1 HA pair	21GB/s	368TB / 3.6PB	100 GbE	1-48	1-24
	12 HA pairs	252GB/s	4.4PB / 43.2PB		576	288
AFF A400	1 HA pair	11GB/s	182TB / 14.7PB	40/100 GbE	1 - 32	1-16
	12 HA pairs	132GB/s	2.1PB / 176.4PB		384	192
AFF C400	1 HA pair	8GB/s	182TB / 14.7PB	40/100 GbE	1 - 16	1-8
	12 HA pairs	128GB/s	2.1PB / 176.4PB		192	96
AFF A250	1 HA pair	7.4GB/s	91.2TB / 4.4PB	25 GbE 40/100GbE	1 - 16	1-8
	4 HA pairs	29.6GB/s	364.8TB / 17.6PB		64	32
AFF C250	1 HA pair	5 GB/s	91.2TB / 4.4PB	25 GbE 40/100GbE	1-8	1-4
	4 HA pairs	20 GB/s	364.8TB / 17.6PB		32	8

1 – 1 AFF = 1 HA pair = 2 Nodes. 12 HA pairs = 24 nodes
2 – 100% sequential read

3 – Based on workload testing in NVA-1153
4 – Based on BasePOD validation test results

NVA-1173 NetApp AIPod 搭配 NVIDIA DGX 系統 - 結論與其他資訊

本節包含 NetApp AIPod 搭配 NVIDIA DGX 系統的其他資訊參考資料。

結論

DGX BasePOD 架構是新一代的深度學習平台、需要同樣進階的儲存與資料管理功能。透過將 DGX BasePOD 與 NetApp AFF 系統結合、NetApp AIPod 與 DGX 系統架構幾乎可以在任何規模上實作。AFF 結合 NetApp ONTAP 優異的雲端整合功能和軟體定義功能、提供跨越邊緣、核心和雲端的完整資料傳輸管道、讓 DL 專案成功完成。

其他資訊

若要深入瞭解本文件所述資訊、請參閱下列文件和 / 或網站：

- NetApp ONTAP 數據管理軟體 ONTAP — 資訊庫

["https://docs.netapp.com/us-en/ontap-family/"](https://docs.netapp.com/us-en/ontap-family/)

- NetApp AFF A900 儲存系統 -

["https://www.netapp.com/data-storage/aff-a-series/aff-a900/"](https://www.netapp.com/data-storage/aff-a-series/aff-a900/)

- NetApp ONTAP RDMA 資訊 -

["https://docs.netapp.com/us-en/ontap/nfs-rdma/index.html"](https://docs.netapp.com/us-en/ontap/nfs-rdma/index.html)

- NetApp DataOps工具套件

["https://github.com/NetApp/netapp-dataops-toolkit"](https://github.com/NetApp/netapp-dataops-toolkit)

- NetApp Astra Trident

"總覽"

- NetApp GPUDirect 儲存部落格 -

["https://www.netapp.com/blog/ontap-reaches-171-gpudirect-storage/"](https://www.netapp.com/blog/ontap-reaches-171-gpudirect-storage/)

- NVIDIA DGX基礎POD

["https://www.nvidia.com/en-us/data-center/dgx-basepod/"](https://www.nvidia.com/en-us/data-center/dgx-basepod/)

- NVIDIA DGX H100 系統

["https://www.nvidia.com/en-us/data-center/dgx-h100/"](https://www.nvidia.com/en-us/data-center/dgx-h100/)

- NVIDIA Networking

["https://www.nvidia.com/en-us/networking/"](https://www.nvidia.com/en-us/networking/)

- NVIDIA Magnum IO™ ; GPUDirect® ; 儲存設備

["https://docs.nvidia.com/gpudirect-storage"](https://docs.nvidia.com/gpudirect-storage)

- NVIDIA Base Command

["https://www.nvidia.com/en-us/data-center/base-command/"](https://www.nvidia.com/en-us/data-center/base-command/)

- NVIDIA Base Command Manager

["https://www.nvidia.com/en-us/data-center/base-command/manager"](https://www.nvidia.com/en-us/data-center/base-command/manager)

- NVIDIA AI Enterprise

["https://www.nvidia.com/en-us/data-center/products/ai-enterprise/"](https://www.nvidia.com/en-us/data-center/products/ai-enterprise/)

感謝

本文檔是 NetApp 解決方案與 ONTAP 工程團隊的工作成果：David Arnette、Olga Kornievskaia、Dustin Fischer、Srikanth Kaligotla、Mohit Kumar 和 Raghuram Sudhaakar。作者也想感謝 NVIDIA 和 NVIDIA

DGX BasePOD 工程團隊持續提供支援。

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。