



搭载NVIDIA DGX 系统的NetApp AI Pod NetApp artificial intelligence solutions

NetApp
February 12, 2026

目录

搭载NVIDIA DGX 系统的NetApp AIPOd	1
NVA-1173 NetApp AIPOd与NVIDIA DGX 系统 - 简介	1
内容提要	1
搭载NVIDIA DGX 系统的 NVA-1173 NetApp AIPOd - 硬件组件	2
NetApp AFF存储系统	2
NVIDIA DGX BasePOD	3
NVA-1173 NetApp AIPOd与NVIDIA DGX 系统 - 软件组件	5
NVIDIA软件	5
NetApp 软件	7
NVA-1173 NetApp AIPOd与NVIDIA DGX H100 系统 - 解决方案架构	8
搭载 DGX 系统的NetApp AIPOd	8
网络设计	9
DGX H100 系统的存储访问概述	10
存储系统设计	10
管理平面服务器	11
NVA-1173 NetApp AIPOd与NVIDIA DGX 系统 - 部署详情	11
存储网络配置	13
存储系统配置	14
NVA-1173 NetApp AIPOd与NVIDIA DGX 系统 - 解决方案验证和规模调整指南	19
解决方案验证	19
存储系统规模指南	19
NVA-1173 NetApp AIPOd与NVIDIA DGX 系统 - 结论及其他信息	20
结束语	20
追加信息	20
致谢	21

搭载NVIDIA DGX 系统的NetApp AI Pod

NVA-1173 NetApp AI Pod与NVIDIA DGX 系统 - 简介

POWERED BY



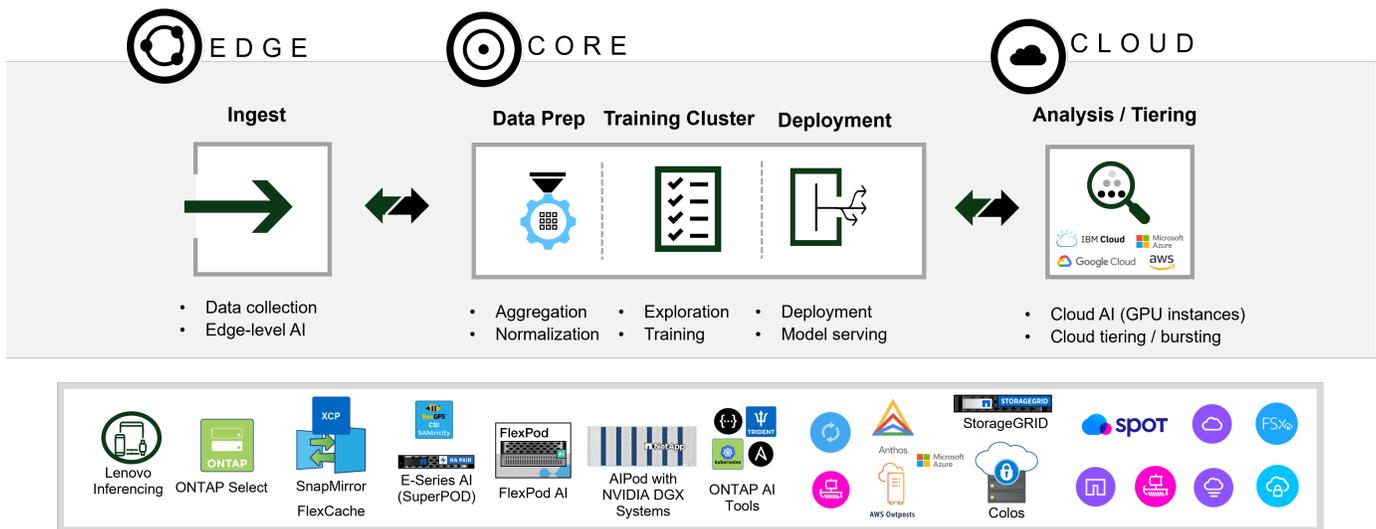
NVIDIA

NetApp解决方案工程

内容提要

NetApp™ AI Pod配备NVIDIA DGX™ 系统和NetApp云连接存储系统，通过消除设计复杂性和猜测，简化了机器学习 (ML) 和人工智能 (AI) 工作负载的基础设施部署。基于NVIDIA DGX BasePOD™ 设计，旨在为下一代工作负载提供卓越的计算性能，搭载NVIDIA DGX 系统的AI Pod增加了NetApp AFF存储系统，使客户能够从小规模开始并无中断地发展，同时智能地管理从边缘到核心再到云端的数据。NetApp AI Pod是NetApp AI 解决方案产品组合的一部分，如下图所示。

NetApp 人工智能解决方案组合



本文档描述了AI Pod参考架构的关键组件、系统连接和配置信息、验证测试结果和解决方案规模指导。本文档适用于有兴趣为 ML/DL 和分析工作负载部署高性能基础架构的NetApp和合作伙伴解决方案工程师以及客户战略决策者。

搭载NVIDIA DGX 系统的 NVA-1173 NetApp AI Pod - 硬件组件

本节重点介绍带有NVIDIA DGX 系统的NetApp AI Pod的硬件组件。

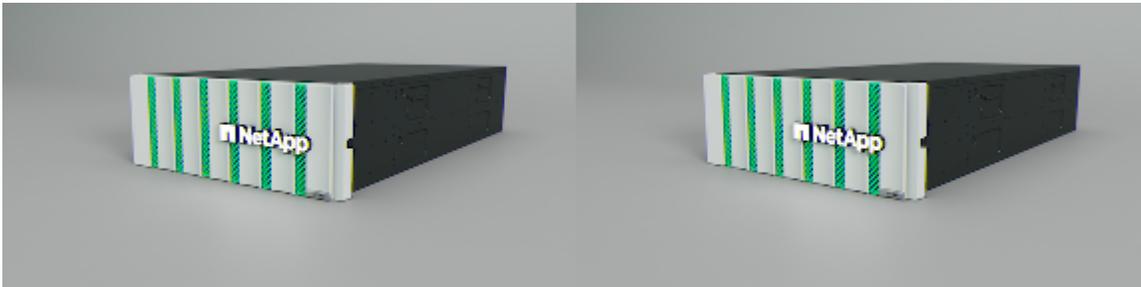
NetApp AFF存储系统

NetApp AFF最先进的存储系统使 IT 部门能够通过业界领先的性能、卓越的灵活性、云集成和一流的数据管理来满足企业存储需求。AFF系统专为闪存设计，有助于加速、管理和保护关键业务数据。

AFF A90存储系统

由NetApp ONTAP数据管理软件提供支持的NetApp AFF A90提供内置数据保护、可选的反勒索软件功能以及支持最关键业务工作负载所需的高性能和弹性。它消除了对关键任务操作的中断，最大限度地减少了性能调整，并保护您的数据免受勒索软件攻击。它提供：
• 行业领先的性能
• 不折不扣的数据安全性
• 简化的无中断升级

NetApp AFF A90存储系统



行业领先的性能

AFF A90可轻松管理深度学习、人工智能和高速分析等下一代工作负载以及 Oracle、SAP HANA、Microsoft SQL Server 和虚拟化应用程序等传统企业数据库。它使关键业务应用程序保持最高速度运行，每个 HA 对高达 2.4M IOPS，延迟低至 100 μ s，并且性能比以前的NetApp型号提高高达 50%。借助 NFS over RDMA、pNFS 和会话中继，客户可以使用现有的数据中心网络基础设施实现下一代应用程序所需的高水平网络性能。客户还可以通过对 SAN、NAS 和对象存储的统一多协议支持进行扩展和增长，并通过统一的单一ONTAP数据管理软件为本地或云端数据提供最大的灵活性。此外，还可以通过Active IQ和Cloud Insights提供的基于 AI 的预测分析来优化系统健康状况。

不妥协的数据安全

AFF A90系统包含一整套NetApp集成和应用程序一致的数据保护软件。它提供内置数据保护和尖端反勒索软件解决方案，用于预防和攻击后恢复。可以阻止恶意文件写入磁盘，并且可以轻松监控存储异常以获取洞察。

简化的无中断升级

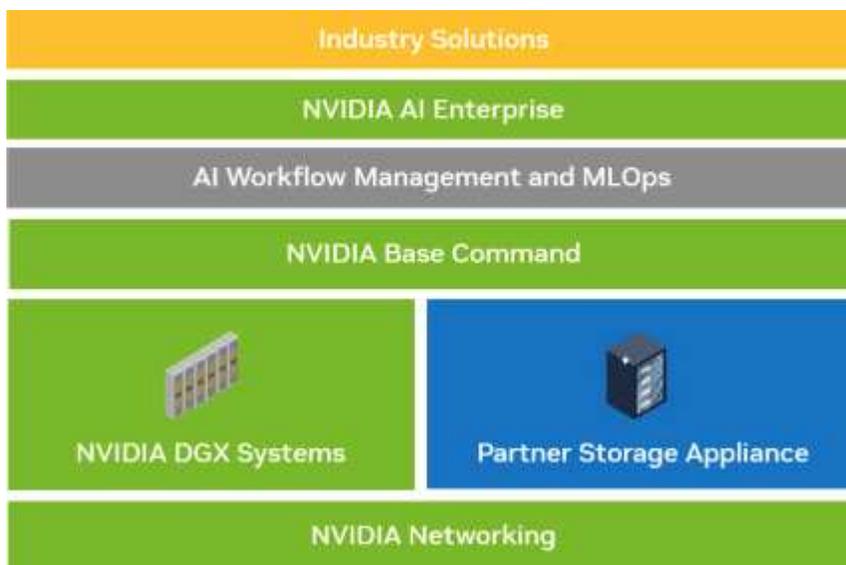
对于现有的 A800 客户来说，AFF A90 可以作为无中断机箱内升级。NetApp 凭借其先进的可靠性、可用性、可维护性和可管理性 (RASM) 功能，可以轻松更新并消除关键任务操作的中断。此外，由于 ONTAP 软件会自动为所有系统组件应用固件更新，NetApp 进一步提高了运营效率并简化了 IT 团队的日常活动。

对于最大的部署，AFF A1K 系统提供最高的性能和容量选项，而其他 NetApp 存储系统（如 AFF A70 和 AFF C800）则以较低的成本为较小的部署提供选项。

NVIDIA DGX BasePOD

NVIDIA DGX BasePOD 是由 NVIDIA 硬件和软件组件、MLOps 解决方案以及第三方存储组成的集成解决方案。利用 NVIDIA 产品和经过验证的合作伙伴解决方案的横向扩展系统设计最佳实践，客户可以实现高效且易于管理的 AI 开发平台。图 1 突出显示了 NVIDIA DGX BasePOD 的各个组件。

NVIDIA DGX BasePOD 解决方案



NVIDIA DGX H100 系统

NVIDIA DGX H100™ 系统是 AI 的强大引擎，由 NVIDIA H100 Tensor Core GPU 的突破性性能加速。

NVIDIA DGX H100 系统

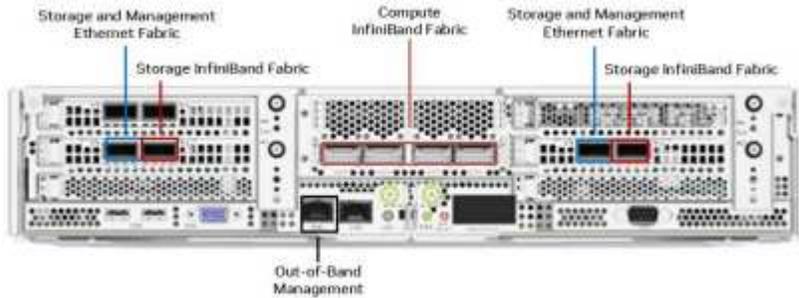


DGX H100 系统的主要规格如下：

- 八个 NVIDIA H100 GPU。
- 每个 GPU 配备 80 GB GPU 内存，总计 640GB。
- 四个 NVIDIA NVSwitch 芯片。
- 双 56 核 Intel Xeon Platinum 8480 处理器，支持 PCIe 5.0。
- 2 TB DDR5 系统内存。
- 四个 OSFP 端口，服务于八个单端口 NVIDIA ConnectX™-7 (InfiniBand/以太网) 适配器和

两个双端口NVIDIA ConnectX-7 (InfiniBand/以太网) 适配器。•两个 1.92 TB M.2 NVMe 驱动器用于 DGX OS, 八个 3.84 TB U.2 NVMe 驱动器用于存储/缓存。•最大功率10.2 kW。 DGX H100 CPU 托盘的后端口如下所示。四个 OSFP 端口为 InfiniBand 计算结构的八个 ConnectX-7 适配器提供服务。每对双端口 ConnectX-7 适配器为存储和管理结构提供并行路径。带外端口用于BMC访问。

NVIDIA DGX H100 后面板



NVIDIA 网络连接

NVIDIA Quantum-2 QM9700 交换机

NVIDIA Quantum-2 QM9700 InfiniBand 交换机



具有 400Gb/s InfiniBand 连接的NVIDIA Quantum-2 QM9700 交换机为NVIDIA Quantum-2 InfiniBand BasePOD 配置中的计算结构提供动力。 ConnectX-7 单端口适配器用于 InfiniBand 计算结构。每个NVIDIA DGX 系统与每个 QM9700 交换机都有双重连接，从而在系统之间提供多条高带宽、低延迟路径。

NVIDIA Spectrum-3 SN4600 交换机

NVIDIA Spectrum-3 SN4600 交换机



NVIDIA Spectrum™-3 SN4600 交换机总共提供 128 个端口（每个交换机 64 个），为 DGX BasePOD 的带内管理提供冗余连接。 NVIDIA SN4600 交换机可以提供 1 GbE 到 200 GbE 之间的速度。对于通过以太网连接的存储设备，也使用NVIDIA SN4600 交换机。 NVIDIA DGX 双端口 ConnectX-7 适配器上的端口用于带内管理和存储连接。

NVIDIA Spectrum SN2201 交换机

NVIDIA Spectrum SN2201 交换机



NVIDIA Spectrum SN2201 交换机提供 48 个端口，可为带外管理提供连接。带外管理为 DGX BasePOD 中的所有组件提供整合的管理连接。

NVIDIA ConnectX-7 适配器

NVIDIA ConnectX-7 适配器



NVIDIA ConnectX-7 适配器可提供 25/50/100/200/400G 的吞吐量。NVIDIA DGX 系统使用单端口和双端口 ConnectX-7 适配器，为具有 400Gb/s InfiniBand 和以太网的 DGX BasePOD 部署提供灵活性。

NVA-1173 NetApp AI Pod与NVIDIA DGX 系统 - 软件组件

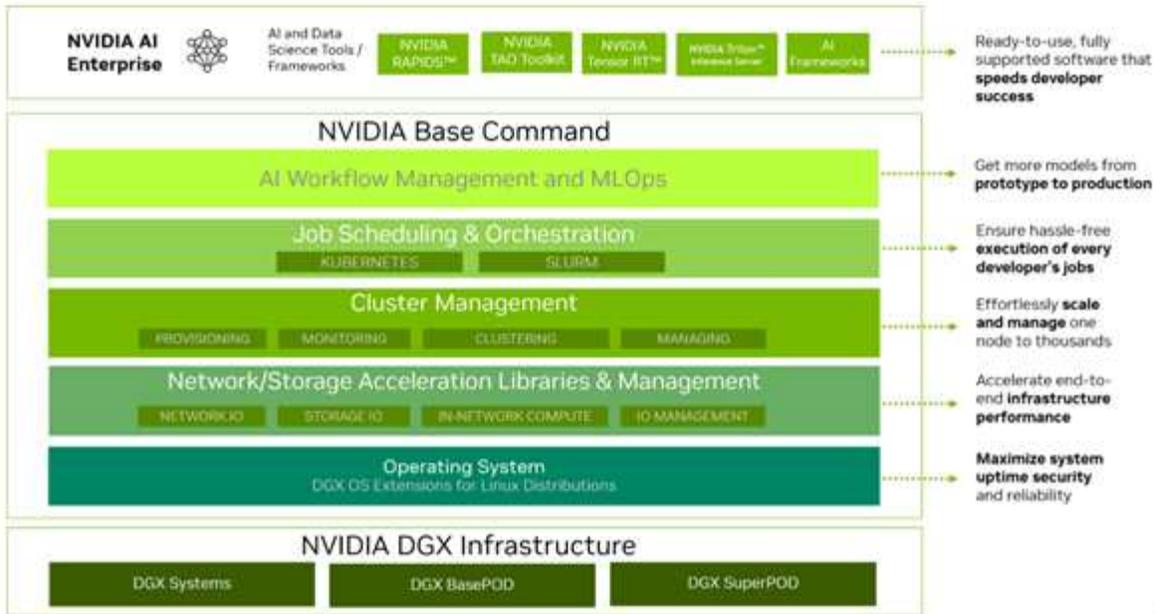
本节重点介绍带有NVIDIA DGX 系统的NetApp AI Pod的软件组件。

NVIDIA 软件

NVIDIA 基本命令

NVIDIA Base Command™ 为每个 DGX BasePOD 提供支持，使组织能够充分利用NVIDIA软件创新的最佳成果。企业可以通过经过验证的平台充分发挥其投资潜力，该平台包括企业级编排和集群管理、加速计算、存储和网络基础设施的库以及针对 AI 工作负载优化的操作系统 (OS)。

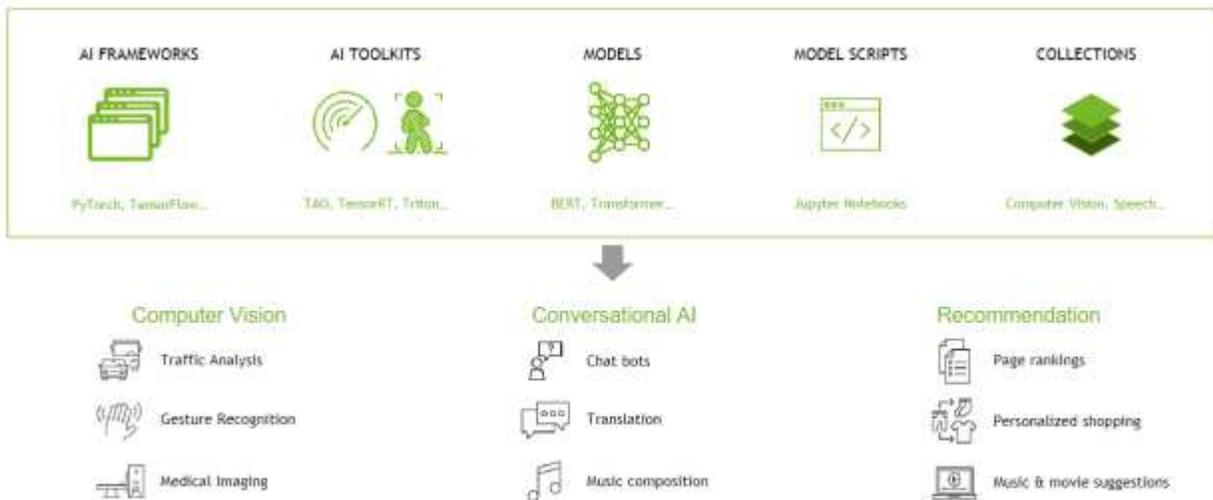
NVIDIA BaseCommand 解决方案



NVIDIA GPU 云 (NGC)

NVIDIA NGC 提供的软件可以满足具有不同 AI 专业水平的数据科学家、开发人员和研究人员的需求。NGC 上托管的软件会针对一组常见漏洞和暴露 (CVE)、加密和私钥进行扫描。它经过测试和设计，可扩展到多个 GPU，在许多情况下，可扩展到多节点，确保用户最大限度地利用其在 DGX 系统上的投资。

NVIDIA GPU 云



NVIDIA AI 企业版

NVIDIA AI Enterprise 是一个端到端软件平台，可让每个企业都能够使用生成式 AI，为在 NVIDIA DGX 平台上优化的生成式 AI 基础模型提供最快、最高效的运行。凭借生产级的安全性、稳定性和可管理性，它简化了生成式 AI 解决方案的开发。NVIDIA AI Enterprise 包含在 DGX BasePOD 中，企业开发人员可以访问预训练模型、优化框架、微服务、加速库和企业支持。

NetApp 软件

NetApp ONTAP

ONTAP 9 是NetApp最新一代存储管理软件，它支持企业实现基础架构现代化并过渡到云就绪数据中心。ONTAP利用业界领先的数据管理功能，只需一套工具即可管理和保护数据，无论数据位于何处。您还可以将数据自由移动到任何需要的地方：边缘、核心或云端。ONTAP 9 包含众多功能，可简化数据管理、加速和保护关键数据，并支持跨混合云架构的下一代基础架构功能。

加速并保护数据

ONTAP提供卓越级别的性能和数据保护，并通过以下方式扩展这些功能：

- 性能和更低的延迟。ONTAP以最低的延迟提供最高的吞吐量，包括支持使用 NFS over RDMA、并行 NFS (pNFS) 和 NFS 会话中继的NVIDIA GPUDirect Storage (GDS)。
- 数据保护。ONTAP提供内置数据保护功能和业界最强大的反勒索软件保障，并在所有平台上实现通用管理。
- NetApp卷加密 (NVE)。ONTAP提供原生卷级加密，同时支持板载和外部密钥管理。
- 存储多租户和多因素身份验证。ONTAP支持以最高级别的安全性共享基础设施资源。

简化数据管理

数据管理对于企业 IT 运营和数据科学家至关重要，以便将适当的资源用于 AI 应用程序和训练 AI/ML 数据集。以下有关NetApp技术的附加信息超出了本次验证的范围，但可能与您的部署相关。

ONTAP数据管理软件包括以下功能，可简化操作并降低总运营成本：

- 快照和克隆支持 ML/DL 工作流的协作、并行实验和增强数据治理。
- SnapMirror可在混合云和多站点环境中实现无缝数据移动，并在需要的时间和地点提供数据。
- 内联数据压缩和扩展重复数据删除。数据压缩减少了存储块内部浪费的空间，重复数据删除显著增加了有效容量。这适用于本地存储的数据和分层到云的数据。
- 最小、最大和自适应服务质量 (AQoS)。细粒度的服务质量 (QoS) 控制有助于维持高度共享环境中关键应用程序的性能水平。
- NetApp FlexGroups 支持在存储集群中的所有节点上分布数据，为超大数据集提供巨大的容量和更高的性能。
- NetApp FabricPool。提供冷数据自动分层到公共和私有云存储选项，包括 Amazon Web Services (AWS)、Azure 和NetApp StorageGRID存储解决方案。有关FabricPool的更多信息，请参阅 "[TR-4598: FabricPool最佳实践](#)"。
- NetApp FlexCache。提供远程卷缓存功能，可简化文件分发、减少 WAN 延迟并降低 WAN 带宽成本。FlexCache支持跨多个站点的分布式产品开发，以及从远程位置加速访问公司数据集。

面向未来的基础设施

ONTAP具有以下功能，可帮助满足苛刻且不断变化的业务需求：

- 无缝扩展和无中断操作。ONTAP支持在线向现有控制器和横向扩展集群添加容量。客户可以升级到最新技术，例如 NVMe 和 32Gb FC，而无需昂贵的数据迁移或中断。
- 云连接。ONTAP是与云连接最紧密的存储管理软件，在所有公共云中均提供软件定义存储 (ONTAP Select) 和云原生实例 (Google Cloud NetApp Volumes) 的选项。

- 与新兴应用程序的集成。ONTAP使用支持现有企业应用的相同基础架构，为下一代平台和应用（如自动驾驶汽车、智能城市和工业 4.0）提供企业级数据服务。

NetApp DataOps 工具包

NetApp DataOps Toolkit 是一款基于 Python 的工具，可简化由高性能、横向扩展NetApp存储支持的开发/培训工作区和推理服务器的管理。DataOps Toolkit 可以作为独立实用程序运行，并且在利用NetApp Trident自动化存储操作的 Kubernetes 环境中更加有效。主要功能包括：

- 快速配置由高性能、横向扩展NetApp存储支持的新的容量 JupyterLab 工作区。
- 快速配置由企业级NetApp存储支持的全新NVIDIA Triton 推理服务器实例。
- 近乎即时地克隆容量的 JupyterLab 工作区，以实现实验或快速迭代。
- 用于备份和/或可追溯性/基准的高容量 JupyterLab 工作区的近乎即时的快照。
- 近乎即时地配置、克隆和快照容量、高性能数据卷。

NetApp Trident

Trident是一个完全受支持的开源存储编排器，适用于容器和 Kubernetes 发行版（包括 Anthos）。Trident可与整个NetApp存储产品组合配合使用，包括NetApp ONTAP，并且还支持 NFS、NVMe/TCP 和 iSCSI 连接。Trident允许最终用户从其NetApp存储系统配置和管理存储，而无需存储管理员的干预，从而加速 DevOps 工作流程。

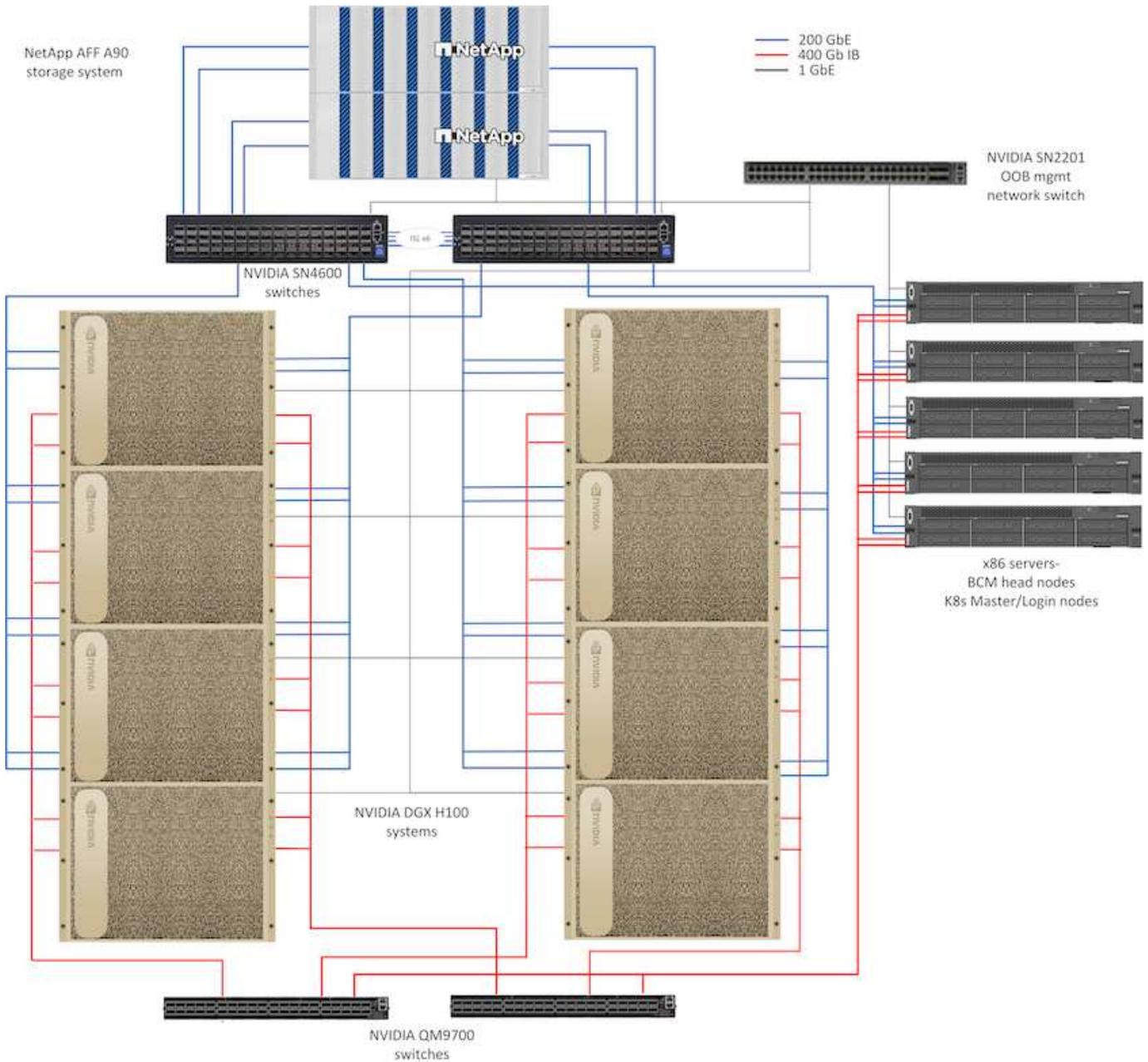
NVA-1173 NetApp AIPOd与NVIDIA DGX H100 系统 - 解决方案架构

本节重点介绍采用NVIDIA DGX 系统的NetApp AIPOd的架构。

搭载 DGX 系统的NetApp AIPOd

该参考架构利用单独的结构进行计算集群互连和存储访问，并在计算节点之间实现 400Gb/s InfiniBand (IB) 连接。下图展示了NetApp AIPOd与 DGX H100 系统的整体解决方案拓扑。

NetApp AIPOd 解决方案拓扑



网络设计

在此配置中，计算集群结构使用一对 QM9700 400Gb/s IB 交换机，它们连接在一起以实现高可用性。每个 DGX H100 系统使用八个连接连接到交换机，其中偶数端口连接到一个交换机，奇数端口连接到另一个交换机。

对于存储系统访问、带内管理和客户端访问，使用一对 SN4600 以太网交换机。交换机之间通过交换机间链路连接，并配置多个 VLAN 来隔离各种流量类型。在特定 VLAN 之间启用基本 L3 路由，以在同一交换机上的客户端和存储接口之间以及交换机之间启用多条路径，从而实现高可用性。对于更大的部署，可以通过根据需要为主干交换机添加额外的交换机对以及为其他叶子交换机添加额外的交换机对，将以太网网络扩展为叶子-主干配置。

除了计算互连和高速以太网网络之外，所有物理设备还连接到一个或多个 SN2201 以太网交换机，以进行带外管理。请参阅[部署详细信息](#)页面以获取有关网络配置的更多信息。

DGX H100 系统的存储访问概述

每个 DGX H100 系统都配备了两个双端口 ConnectX-7 适配器用于管理和存储流量，并且对于此解决方案，每个卡上的两个端口都连接到同一个交换机。然后将每个卡的一个端口配置为 LACP MLAG 绑定，并将一个端口连接到每个交换机，并且带内管理、客户端访问和用户级存储访问的 VLAN 都托管在此绑定上。

每张卡上的另一个端口用于连接 AFF A90 存储系统，并且可以根据工作负载要求以多种配置使用。对于使用 NFS over RDMA 来支持 NVIDIA Magnum IO GPUDirect Storage 的配置，端口单独使用，并且 IP 地址位于单独的 VLAN 中。对于不需要 RDMA 的部署，存储接口也可以配置 LACP 绑定，以提供高可用性和额外的带宽。无论是否使用 RDMA，客户端都可以使用 NFS v4.1 pNFS 和会话中继挂载存储系统，以实现集群中所有存储节点的并行访问。请参阅["部署详细信息"](#)页面以获取有关客户端配置的更多信息。

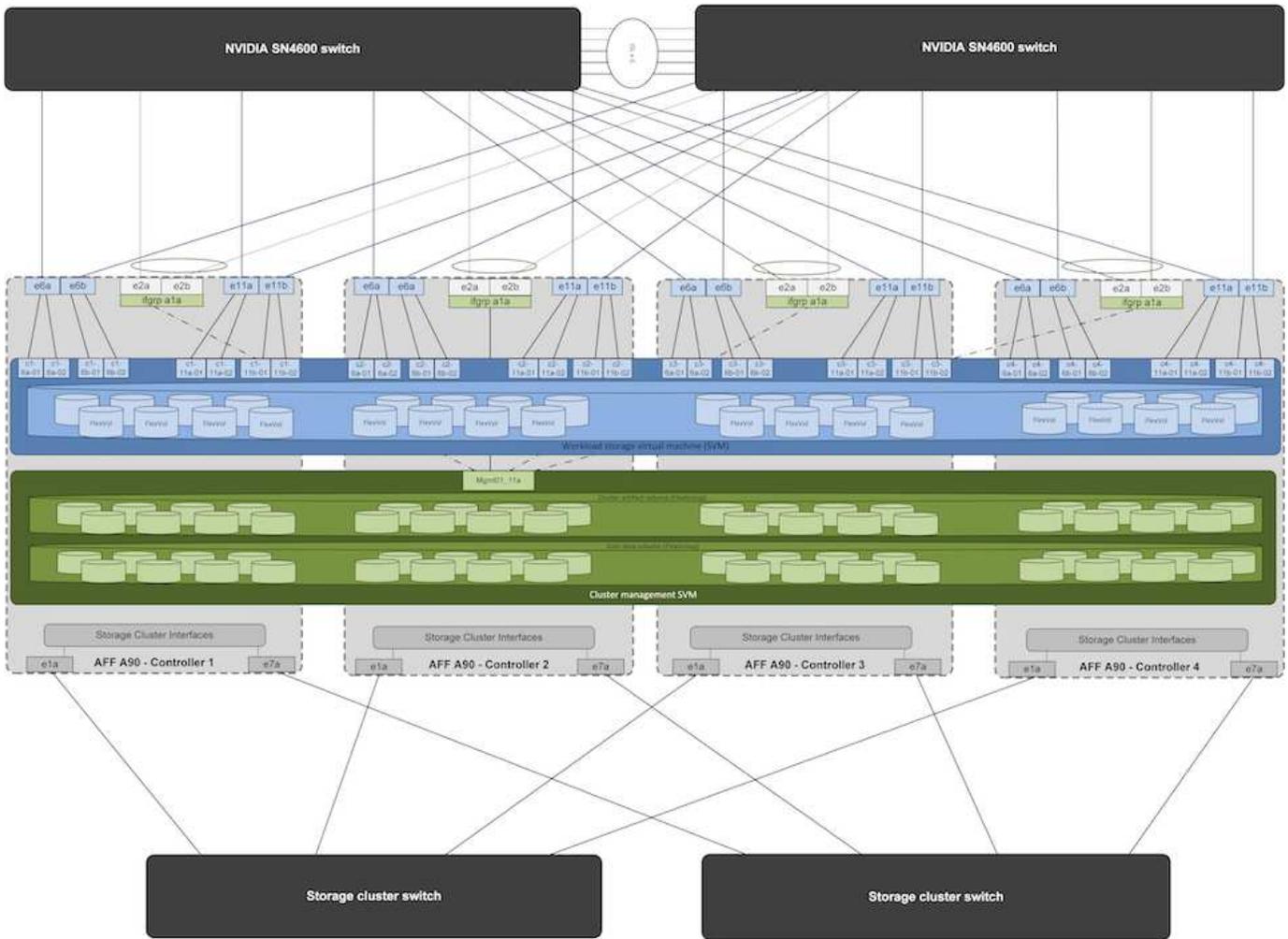
有关 DGX H100 系统连接的详细信息，请参阅["NVIDIA BasePOD 文档"](#)。

存储系统设计

每个 AFF A90 存储系统使用每个控制器的六个 200 GbE 端口进行连接。每个控制器的四个端口用于从 DGX 系统访问工作负载数据，每个控制器的两个端口配置为 LACP 接口组，以支持从管理平面服务器访问集群管理工作和用户主目录。存储系统的所有数据访问均通过 NFS 提供，其中有一个专用于 AI 工作负载访问的存储虚拟机 (SVM) 和一个专用于集群管理用途的单独 SVM。

管理 SVM 只需要一个 LIF，该 LIF 托管在每个控制器上配置的 2 端口接口组上。其他 FlexGroup 卷在管理 SVM 上进行配置，以容纳集群管理构件，如集群节点映像、系统监控历史数据和最终用户主目录。下图显示了存储系统的逻辑配置。

NetApp A90 存储集群逻辑配置



管理平面服务器

该参考架构还包括五个基于 CPU 的服务器，用于管理平面。其中两个系统用作 NVIDIA Base Command Manager 的头节点，用于集群部署和管理。其他三个系统用于提供额外的集群服务，例如 Kubernetes 主节点或利用 Slurm 进行作业调度的部署的登录节点。利用 Kubernetes 的部署可以利用 NetApp Trident CSI 驱动程序为 AFF A900 存储系统上的管理和 AI 工作负载提供具有持久存储的自动配置和数据服务。

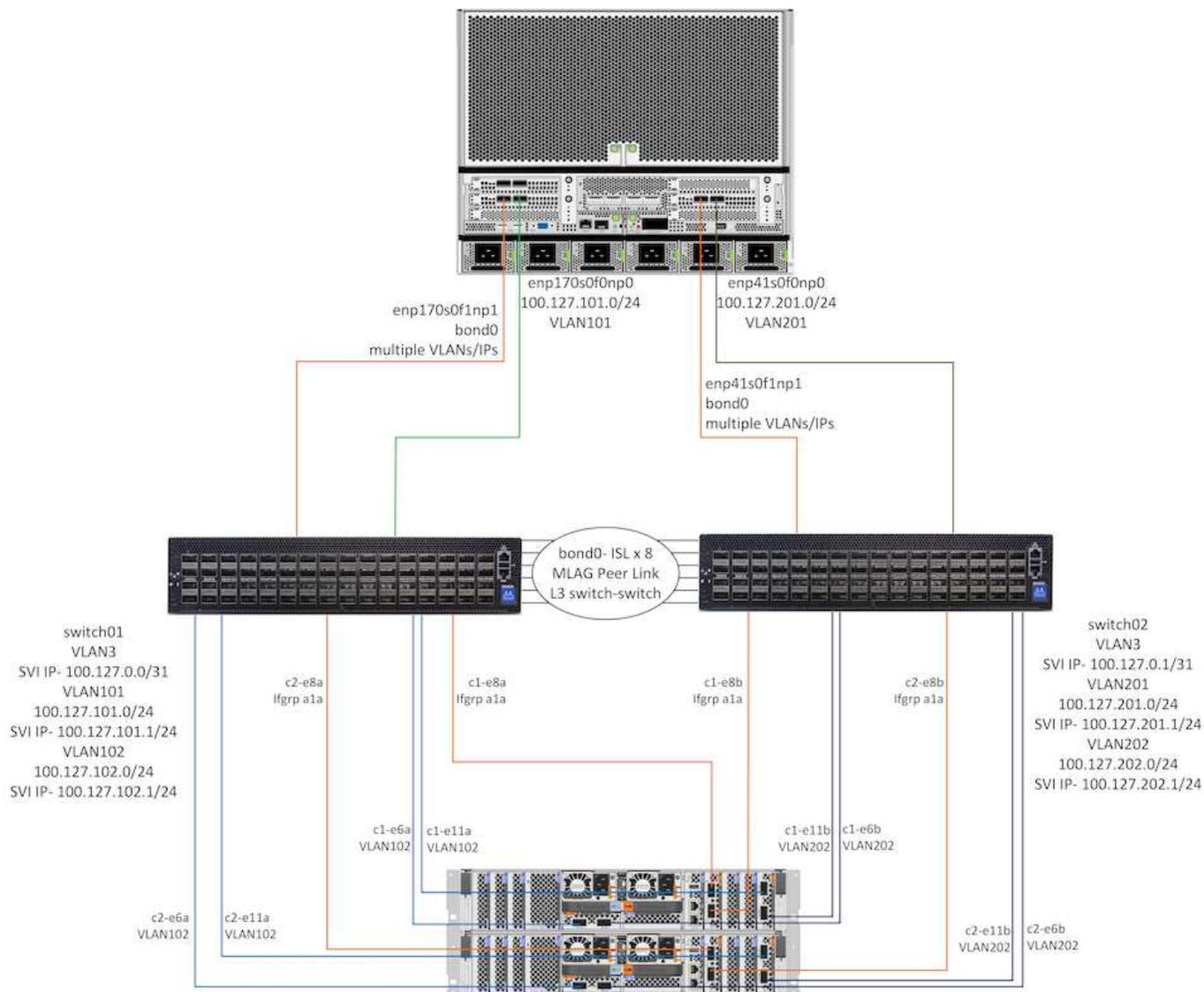
每台服务器都物理连接到 IB 交换机和以太网交换机，以实现集群部署和管理，并通过管理 SVM 配置 NFS 挂载到存储系统，以存储前面所述的集群管理工作。

NVA-1173 NetApp AI Pod 与 NVIDIA DGX 系统 - 部署详情

本节介绍验证此解决方案期间使用的部署细节。使用的 IP 地址仅供参考，请根据部署环境进行修改。有关此配置的实现中使用的特定命令的更多信息，请参阅相应的产品文档。

下图显示了 1 个 DGX H100 系统和 1 个 HA 对 AFF A90 控制器的详细网络和连接信息。以下部分中的部署指南基于此图中的详细信息。

NetApp AI pod 网络配置



下表显示了最多 16 个 DGX 系统和 2 个 AFF A90 HA 对的示例布线分配。

交换机和端口	设备	设备端口
交换机1端口1-16	DGX-H100-01 至 -16	enp170s0f0np0, 插槽1端口1
交换机1端口17-32	DGX-H100-01 至 -16	enp170s0f1np1, 插槽1端口2
交换机1端口33-36	AFF-A90-01 至 -04	端口 e6a
交换机1端口37-40	AFF-A90-01 至 -04	端口 e11a
交换机1端口41-44	AFF-A90-01 至 -04	端口 e2a
交换机1端口57-64	ISL 到交换机 2	端口 57-64
交换机2端口1-16	DGX-H100-01 至 -16	enp41s0f0np0, 插槽2端口1
交换机2端口17-32	DGX-H100-01 至 -16	enp41s0f1np1, 插槽 2 端口 2
交换机2端口33-36	AFF-A90-01 至 -04	端口 e6b
交换机2端口37-40	AFF-A90-01 至 -04	端口 e11b

交换机和端口	设备	设备端口
交换机2端口41-44	AFF-A90-01 至 -04	端口 e2b
交换机2端口57-64	ISL 到交换机 1	端口 57-64

下表显示了本次验证中使用的各个组件的软件版本。

设备	软件版本
NVIDIA SN4600 交换机	Cumulus Linux v5.9.1
NVIDIA DGX 系统	DGX 操作系统 v6.2.1 (Ubuntu 22.04 LTS)
Mellanox OFED	24.01
NetApp AFF A90	NetApp ONTAP 9.14.1

存储网络配置

本节概述以太网存储网络配置的关键细节。有关配置 InfiniBand 计算网络的信息，请参阅["NVIDIA BasePOD 文档"](#)。有关交换机配置的详细信息，请参阅["NVIDIA Cumulus Linux 文档"](#)。

配置 SN4600 交换机的基本步骤概述如下。此过程假定布线和基本交换机设置（管理 IP 地址、许可等）已完成。

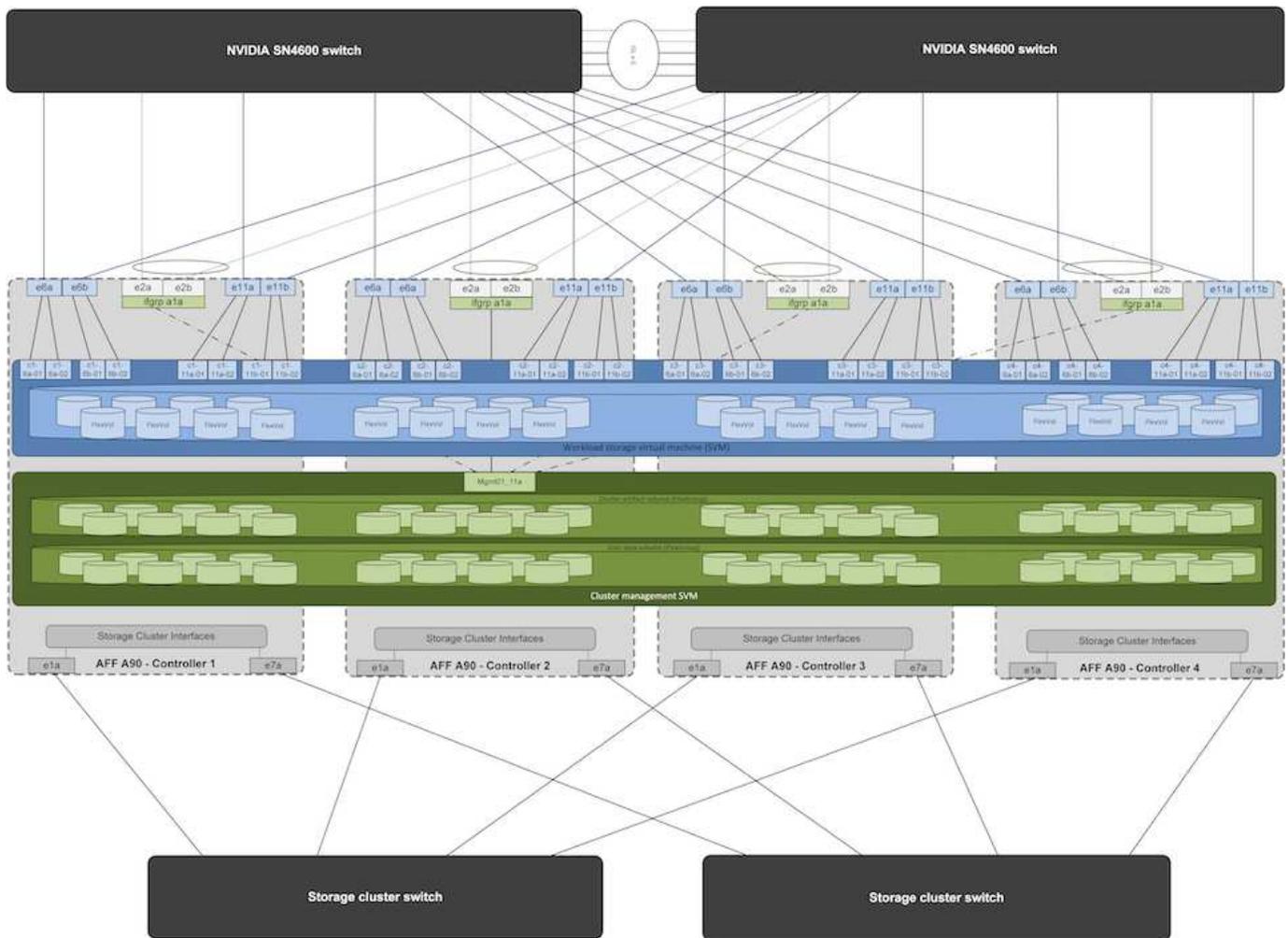
- 配置交换机之间的 ISL 绑定以启用多链路聚合 (MLAG) 和故障转移流量
 - 本次验证使用了 8 条链路，为测试的存储配置提供了足够的带宽
 - 有关启用 MLAG 的具体说明，请参阅 Cumulus Linux 文档。
- 为两台交换机上的每对客户端端口和存储端口配置 LACP MLAG
 - 每个交换机上的端口 swp17 用于 DGX-H100-01 (enp170s0f1np1 和 enp41s0f1np1)，端口 swp18 用于 DGX-H100-02，等等 (bond1-16)
 - 每个交换机上的端口 swp41 用于 AFF-A90-01 (e2a 和 e2b)，端口 swp42 用于 AFF-A90-02，等等 (bond17-20)
 - nv 设置接口 bondX 债券成员 swpX
 - nv 设置接口 bondx 绑定 mlag id X
- 将所有端口和 MLAG 绑定添加到默认桥接域
 - nv 设置 int swp1-16,33-40 桥接域 br_default
 - nv 设置 int bond1-20 桥接域 br_default
- 在每台交换机上启用 RoCE
 - nv 设置 roce 模式无损
- 配置 VLAN - 2 个用于客户端端口，2 个用于存储端口，1 个用于管理，1 个用于 L3 交换机到交换机
 - 开关 1-
 - VLAN 3 用于在客户端 NIC 发生故障时进行 L3 交换机到交换机的路由
 - 每个 DGX 系统上的存储端口 1 的 VLAN 101 (enp170s0f0np0, slot1 端口 1)

- 每个AFF A90存储控制器上的端口 e6a 和 e11a 的 VLAN 102
 - VLAN 301 用于使用 MLAG 接口对每个 DGX 系统和存储控制器进行管理
 - 开关 2-
 - VLAN 3 用于在客户端 NIC 发生故障时进行 L3 交换机到交换机的路由
 - 每个 DGX 系统上的存储端口 2 的 VLAN 201 (enp41s0f0np0, slot2 端口 1)
 - 每个AFF A90存储控制器上的端口 e6b 和 e11b 的 VLAN 202
 - VLAN 301 用于使用 MLAG 接口对每个 DGX 系统和存储控制器进行管理
6. 根据需要将物理端口分配给每个 VLAN，例如客户端 VLAN 中的客户端端口和存储 VLAN 中的存储端口
- nv 设置 int <swpX> 桥接域 br_default 访问 <vlan id>
 - MLAG 端口应保持为中继端口，以根据需要在绑定接口上启用多个 VLAN。
7. 在每个 VLAN 上配置交换机虚拟接口 (SVI) 以充当网关并启用 L3 路由
- 开关 1-
 - nv 设置 int vlan3 ip 地址 100.127.0.0/31
 - nv 设置 int vlan101 ip 地址 100.127.101.1/24
 - nv 设置 int vlan102 ip 地址 100.127.102.1/24
 - 开关 2-
 - nv 设置 int vlan3 ip 地址 100.127.0.1/31
 - nv 设置 int vlan201 ip 地址 100.127.201.1/24
 - nv 设置 int vlan202 ip 地址 100.127.202.1/24
8. 创建静态路由
- 同一交换机上的子网将自动创建静态路由
 - 当客户端链路发生故障时，交换机到交换机的路由需要额外的静态路由
 - 开关 1-
 - nv 设置 VRF 默认路由器静态 100.127.128.0/17 通过 100.127.0.1
 - 开关 2-
 - nv 设置 VRF 默认路由器静态 100.127.0.0/17 通过 100.127.0.0

存储系统配置

本节介绍此解决方案的 A90 存储系统配置的关键细节。有关ONTAP系统配置的更多详细信息，请参阅["ONTAP 文档"](#)。下图显示了存储系统的逻辑配置。

NetApp A90 存储集群逻辑配置



配置存储系统的基本步骤概述如下。此过程假设基本存储集群安装已经完成。

1. 在每个控制器上配置 1 个聚合，所有可用分区减去 1 个备用分区
 - `aggr create -node <节点> -aggregate <节点>_data01 -diskcount <47>`
2. 在每个控制器上配置 ifgrps
 - 网络端口 `ifgrp create -node <节点> -ifgrp a1a -mode multimode_lacp -distr-function port`
 - 网络端口 `ifgrp add-port -node <节点> -ifgrp <ifgrp> -ports <节点>:e2a,<节点>:e2b`
3. 在每个控制器上的 ifgrp 上配置 mgmt vlan 端口
 - 网络端口 vlan 创建 -节点 aff-a90-01 -端口 a1a -vlan-id 31
 - 网络端口 vlan 创建 -节点 aff-a90-02 -端口 a1a -vlan-id 31
 - 网络端口 vlan 创建 -节点 aff-a90-03 -端口 a1a -vlan-id 31
 - 网络端口 vlan 创建 -节点 aff-a90-04 -端口 a1a -vlan-id 31
4. 创建广播域
 - 广播域创建-广播域vlan21-mtu 9000-端口aff-a90-01: e6a, aff-a90-01: e11a, aff-a90-02: e6a, aff-a90-02: e11a, aff-a90-03: e6a, aff-a90-03: e11a, aff-a90-04: e6a, aff-a90-04: e11a
 - 广播域创建-广播域vlan22-mtu 9000-端口aaff-a90-01: e6b, aff-a90-01: e11b, aff-a90-02: e6b, aff-a90-02: e11b, aff-a90-03: e6b, aff-a90-03: e11b, aff-a90-04: e6b, aff-a90-04: e11b

- 广播域创建-广播域vlan31-mtu 9000-端口aff-a90-01:a1a-31, aff-a90-02:a1a-31, aff-a90-03:a1a-31, aff-a90-04:a1a-31

5. 创建管理 SVM *

6. 配置管理 SVM

- 创建 LIF
 - net int create -vserver basepod-mgmt -lif vlan31-01 -home-node aff-a90-01 -home-port a1a-31 -address 192.168.31.X -netmask 255.255.255.0
- 创建FlexGroup卷-
 - 卷创建-vserver basepod-mgmt-volume home-size 10T-auto-provision-as flexgroup-junction-path /home
 - 卷创建-vserver basepod-mgmt-volume cm-size 10T-auto-provision-as flexgroup-junction-path /cm
- 制定出口政策
 - 导出策略规则创建-vserver basepod-mgmt-policy default-client-match 192.168.31.0/24-rerule sys-rerule sys-superuser sys

7. 创建数据 SVM *

8. 配置数据 SVM

- 配置 SVM 以支持 RDMA
 - vserver nfs 修改-vserver basepod-data -rdma 已启用
- 创建 LIF
 - net int create -vserver basepod-data -lif c1-6a-lif1 -home-node aff-a90-01 -home-port e6a -address 100.127.102.101 -netmask 255.255.255.0
 - net int create -vserver basepod-data -lif c1-6a-lif2 -home-node aff-a90-01 -home-port e6a -address 100.127.102.102 -netmask 255.255.255.0
 - net int create -vserver basepod-data -lif c1-6b-lif1 -home-node aff-a90-01 -home-port e6b -address 100.127.202.101 -netmask 255.255.255.0
 - net int create -vserver basepod-data -lif c1-6b-lif2 -home-node aff-a90-01 -home-port e6b -address 100.127.202.102 -netmask 255.255.255.0
 - net int create -vserver basepod-data -lif c1-11a-lif1 -home-node aff-a90-01 -home-port e11a -address 100.127.102.103 -netmask 255.255.255.0
 - net int create -vserver basepod-data -lif c1-11a-lif2 -home-node aff-a90-01 -home-port e11a -address 100.127.102.104 -netmask 255.255.255.0
 - net int create -vserver basepod-data -lif c1-11b-lif1 -home-node aff-a90-01 -home-port e11b -address 100.127.202.103 -netmask 255.255.255.0
 - net int create -vserver basepod-data -lif c1-11b-lif2 -home-node aff-a90-01 -home-port e11b -address 100.127.202.104 -netmask 255.255.255.0
 - net int create -vserver basepod-data -lif c2-6a-lif1 -home-node aff-a90-02 -home-port e6a -address 100.127.102.105 -netmask 255.255.255.0
 - net int create -vserver basepod-data -lif c2-6a-lif2 -home-node aff-a90-02 -home-port e6a -address 100.127.102.106 -netmask 255.255.255.0
 - net int create -vserver basepod-data -lif c2-6b-lif1 -home-node aff-a90-02 -home-port e6b -address 100.127.202.105 -netmask 255.255.255.0

- net int create -vserver basepod-data -lif c2-6b-lif2 -home-node aff-a90-02 -home-port e6b -address 100.127.202.106 -netmask 255.255.255.0
- net int create -vserver basepod-data -lif c2-11a-lif1 -home-node aff-a90-02 -home-port e11a -address 100.127.102.107 -netmask 255.255.255.0
- net int create-vserver basepod-data-lif c2-11a-lif2-home-node aff-a90-02-home-port e11a-address 100.127.102.108-netmask 255.255.255.0
- net int create -vserver basepod-data -lif c2-11b-lif1 -home-node aff-a90-02 -home-port e11b -address 100.127.202.107 -netmask 255.255.255.0
- net int create-vserver basepod-data-lif c2-11b-lif2-home-node aff-a90-02-home-port e11b-address 100.127.202.108-netmask 255.255.255.0

9. 配置 LIF 以进行 RDMA 访问

- 对于使用 ONTAP 9.15.1 的部署，物理信息的 RoCE QoS 配置需要 ONTAP CLI 中不可用的操作系统级命令。请联系 NetApp 支持以获取有关 RoCE 支持端口配置的帮助。NFS over RDMA 功能正常
- 从 ONTAP 9.16.1 开始，物理接口将自动配置适当的设置以实现端到端 RoCE 支持。
- net int 修改-vserver basepod-data -lif * -rdma-protocols roce

10. 在数据 SVM 上配置 NFS 参数

- nfs 修改 -vserver basepod-data -v4.1 已启用 -v4.1-pnfs 已启用 -v4.1-trunking 已启用 -tcp-max-transfer-size 262144

11. 创建 FlexGroup 卷

- 卷创建-vserver basepod-data-volume 数据-size 100T-auto-provision-as flexgroup-junction-path /data

12. 创建导出策略

- 导出策略规则创建-vserver basepod-data-policy default-client-match 100.127.101.0/24-rorule sys-rwrule sys-superuser sys
- 导出策略规则创建-vserver basepod-data-policy default-client-match 100.127.201.0/24-rorule sys-rwrule sys-superuser sys

13. 创建路线

- 路由添加-vserver basepod_data-目的地100.127.0.0/17-网关100.127.102.1度量20
- 路由添加-vserver basepod_data-目的地100.127.0.0/17-网关100.127.202.1度量30
- 路由添加-vserver basepod_data-目的地100.127.128.0/17-网关100.127.202.1度量20
- 路由添加-vserver basepod_data-目的地100.127.128.0/17-网关100.127.102.1度量30

用于 RoCE 存储访问的 DGX H100 配置

本节介绍 DGX H100 系统配置的关键细节。许多配置项可以包含在部署到 DGX 系统的 OS 映像中，或者在启动时由 Base Command Manager 实现。这里列出它们以供参考，有关在 BCM 中配置节点和软件映像的更多信息，请参阅["BCM 文档"](#)。

1. 安装其他软件包

- ipmitool
- python3-pip

2. 安装 Python 包

- 波罗米科
 - matplotlib
3. 软件包安装后重新配置 dpkg
 - dpkg——配置-a
 4. 安装 MOFED
 5. 设置 mst 值以进行性能调整
 - mstconfig -y -d <aa:00.0,29:00.0> 设置 ADVANCED_PCI_SETTINGS=1 NUM_OF_VFS=0
MAX_ACC_OUT_READ=44
 6. 修改设置后重置适配器
 - mlxfwreset -d <aa:00.0,29:00.0> -y 重置
 7. 在 PCI 设备上设置 MaxReadReq
 - setpci -s <aa:00.0,29:00.0> 68.W=5957
 8. 设置 RX 和 TX 环形缓冲区大小
 - ethtool -G <enp170s0f0np0,enp41s0f0np0> rx 8192 tx 8192
 9. 使用 mlx_qos 设置 PFC 和 DSCP
 - mlx_qos -i <enp170s0f0np0,enp41s0f0np0> --pfc 0,0,0,1,0,0,0,0 --trust=dscp --cable_len=3
 10. 为网络端口上的 RoCE 流量设置 ToS
 - echo 106 > /sys/class/infiniband/<mlx5_7,mlx5_1>/tc/1/traffic_class
 11. 在适当的子网上为每个存储 NIC 配置一个 IP 地址
 - 100.127.101.0/24 用于存储 NIC 1
 - 100.127.201.0/24 用于存储 NIC 2
 12. 配置带内网络端口进行 LACP 绑定 (enp170s0f1np1、enp41s0f1np1)
 13. 为每个存储子网的主路径和次路径配置静态路由
 - 路由添加 -net 100.127.0.0/17 gw 100.127.101.1 metric 20
 - 路由添加 -net 100.127.0.0/17 gw 100.127.201.1 公制 30
 - 路由添加 -net 100.127.128.0/17 gw 100.127.201.1 公制 20
 - 路由添加 -net 100.127.128.0/17 gw 100.127.101.1 公制 30
 14. 挂载 /home 卷
 - 安装-o vers = 3, nconnect = 16, rsize = 262144, wsize = 262144 192.168.31.X: /home /home
 15. 挂载/数据卷
 - 安装数据卷时使用了以下安装选项-
 - vers=4.1 # 启用 pNFS 来并行访问多个存储节点
 - proto=rdma # 将传输协议设置为 RDMA，而不是默认的 TCP
 - max_connect=16 #启用 NFS 会话中继来聚合存储端口带宽
 - write=eager # 提高缓冲写入的写入性能

- rsize=262144,wsiz=262144 # 将 I/O 传输大小设置为 256k

NVA-1173 NetApp AIPod与NVIDIA DGX 系统 - 解决方案验证和规模调整指南

本节重点介绍采用NVIDIA DGX 系统的NetApp AIPod的解决方案验证和尺寸调整指导。

解决方案验证

使用开源工具 FIO 通过一系列合成工作负载验证了此解决方案中的存储配置。这些测试包括读写 I/O 模式，旨在模拟执行深度学习训练作业的 DGX 系统产生的存储工作负载。使用同时运行 FIO 工作负载的 2 插槽 CPU 服务器集群来验证存储配置，以模拟 DGX 系统集群。每个客户端都配置了前面描述的不同网络配置，并添加了以下详细信息。

以下安装选项用于此验证：

版本=4.1	启用 pNFS 来并行访问多个存储节点
原型=rdma	将传输协议设置为 RDMA，而不是默认的 TCP
端口=20049	为 RDMA NFS 服务指定正确的端口
最大连接数=16	启用 NFS 会话中继来聚合存储端口带宽
写=渴望	提高缓冲写入的写入性能
rsize=262144,wsiz=262144	将 I/O 传输大小设置为 256k

此外，客户端的 NFS max_session_slots 值配置为 1024。由于该解决方案是使用 NFS over RDMA 进行测试的，因此存储网络端口配置了主动/被动结合。本次验证使用了以下债券参数：

模式=主动备份	将绑定设置为主动/被动模式
primary=<接口名称>	所有客户端的主接口分布在交换机上
mii-监控间隔=100	指定监控间隔为100ms
故障转移 mac 策略=活动	指定活动链路的 MAC 地址是绑定的 MAC。这是 RDMA 通过绑定接口正确运行所必需的。

存储系统配置如下，包括两个 A900 HA 对（4 个控制器），每个 HA 对连接两个 NS224 磁盘架，每个磁盘架有 24 个 1.9TB NVMe 磁盘驱动器。如架构部分所述，所有控制器的存储容量使用FlexGroup卷进行组合，并且所有客户端的数据分布在集群中的所有控制器上。

存储系统规模指南

NetApp已成功完成 DGX BasePOD 认证，经测试的两个 A90 HA 对可以轻松支持由 16 个 DGX H100 系统组成的集群。对于具有更高存储性能要求的大型部署，可以将额外的AFF系统添加到NetApp ONTAP集群中，单个集群中最多可包含 12 个 HA 对（24 个节点）。使用本解决方案中描述的FlexGroup技术，24 节点集群可以在单个命名空间中提供超过 79 PB 和高达 552 GBps 的吞吐量。其他NetApp存储系统（例如AFF A400、A250 和 C800）以较低的成本为较小规模的部署提供较低的性能和/或更高的容量选项。由于ONTAP 9 支持混合模型集群，客户可以从较小的初始占用空间开始，并随着容量和性能需求的增长向集群添加更多或更大的存储系统。下表粗略估计了每个AFF型号支持的 A100 和 H100 GPU 的数量。

		Throughput ²	Raw capacity (typical ³ / max)	Connectivity	# NVIDIA A100 GPUs supported ⁴	# NVIDIA H100 GPUs supported ⁵
NetApp [®] AFF A1K	1 HA pair ¹	56 GB/s	368TB / 14.7PB	200 GbE	1-160	1-80
	12 HA pairs	672 GB/s	4.4PB / 176.4PB		1920	960
AFF A90	1 HA pair	46 GB/s	368TB / 6.6PB	200 GbE	1 – 128	1-64
	12 HA pairs	552 GB/s	4.4PB / 79.2PB		1536	768
AFF A70	1 HA pair	21 GB/s	368TB / 6.6PB	200 GbE	1-48	1-24
	12 HA pairs	252 GB/s	4.4PB / 79.2PB		576	288

NVA-1173 NetApp AIPod与NVIDIA DGX 系统 - 结论及其他信息

本节包含有关带有NVIDIA DGX 系统的NetApp AIPod的更多信息的参考。

结束语

DGX BasePOD 架构是下一代深度学习平台，需要同样先进的存储和数据管理功能。通过将 DGX BasePOD 与NetApp AFF系统相结合， NetApp AIPod与 DGX 系统架构几乎可以在任何规模上实现。结合NetApp ONTAP卓越的云集成和软件定义功能， AFF可为成功的 DL 项目提供涵盖边缘、核心和云的全方位数据管道。

追加信息

要了解有关本文档中描述的信息的更多信息，请参阅以下文档和/或网站：

- NetApp ONTAP数据管理软件 — ONTAP信息库
["https://docs.netapp.com/us-en/ontap-family/"](https://docs.netapp.com/us-en/ontap-family/)
- NetApp AFF A90存储系统-
<https://www.netapp.com/pdf.html?item=/media/7828-ds-3582-aff-a-series-ai-era.pdf>
- NetApp ONTAP RDMA 信息-
["https://docs.netapp.com/us-en/ontap/nfs-rdma/index.html"](https://docs.netapp.com/us-en/ontap/nfs-rdma/index.html)
- NetApp DataOps 工具包
["https://github.com/NetApp/netapp-dataops-toolkit"](https://github.com/NetApp/netapp-dataops-toolkit)
- NetApp Trident
["概述"](#)
- NetApp GPUDirect 存储博客-

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

- NVIDIA DGX BasePOD

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

- NVIDIA DGX H100 系统

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

- NVIDIA 网络连接

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

- NVIDIA Magnum IO™ GPUDirect® 存储

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

- NVIDIA基本命令

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

- NVIDIA基础命令管理器

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

- NVIDIA AI 企业版

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

致谢

本文档由NetApp解决方案和ONTAP工程团队（David Arnette、Olga Kornievskaia、Dustin Fischer、Srikanth Kaligotla、Mohit Kumar 和 Raghuram Sudhaakar）编写。作者还要感谢NVIDIA和NVIDIA DGX BasePOD工程团队的持续支持。

版权信息

版权所有 © 2026 NetApp, Inc.。保留所有权利。中国印刷。未经版权所有者事先书面许可，本档中受版权保护的任何部分不得以任何形式或通过任何手段（图片、电子或机械方式，包括影印、录音、录像或存储在电子检索系统中）进行复制。

从受版权保护的 NetApp 资料派生的软件受以下许可和免责声明的约束：

本软件由 NetApp 按“原样”提供，不含任何明示或暗示担保，包括但不限于适销性以及针对特定用途的适用性的隐含担保，特此声明不承担任何责任。在任何情况下，对于因使用本软件而以任何方式造成的任何直接性、间接性、偶然性、特殊性、惩罚性或后果性损失（包括但不限于购买替代商品或服务；使用、数据或利润方面的损失；或者业务中断），无论原因如何以及基于何种责任理论，无论出于合同、严格责任或侵权行为（包括疏忽或其他行为），NetApp 均不承担责任，即使已被告知存在上述损失的可能性。

NetApp 保留在不另行通知的情况下随时对本文档所述的任何产品进行更改的权利。除非 NetApp 以书面形式明确同意，否则 NetApp 不承担因使用本文档所述产品而产生的任何责任或义务。使用或购买本产品不表示获得 NetApp 的任何专利权、商标权或任何其他知识产权许可。

本手册中描述的产品可能受一项或多项美国专利、外国专利或正在申请的专利的保护。

有限权利说明：政府使用、复制或公开本文档受 DFARS 252.227-7013（2014 年 2 月）和 FAR 52.227-19（2007 年 12 月）中“技术数据权利 — 非商用”条款第 (b)(3) 条规定的限制条件的约束。

本文档中所含数据与商业产品和/或商业服务（定义见 FAR 2.101）相关，属于 NetApp, Inc. 的专有信息。根据本协议提供的所有 NetApp 技术数据和计算机软件具有商业性质，并完全由私人出资开发。美国政府对这些数据的使用权具有非排他性、全球性、受限且不可撤销的许可，该许可既不可转让，也不可再许可，但仅限在与交付数据所依据的美国政府合同有关且受合同支持的情况下使用。除本文档规定的情形外，未经 NetApp, Inc. 事先书面批准，不得使用、披露、复制、修改、操作或显示这些数据。美国政府对国防部的授权仅限于 DFARS 的第 252.227-7015(b)（2014 年 2 月）条款中明确的权利。

商标信息

NetApp、NetApp 标识和 <http://www.netapp.com/TM> 上所列的商标是 NetApp, Inc. 的商标。其他公司和产品名称可能是其各自所有者的商标。