# NetApp

NetApp Verified Architecture

# NVIDIA Certified Storage Program

Reference Design for NetApp

NVA-1181-DESIGN

By NetApp ®

In partnership with

**NetApp**     **NVIDIA**

## Abstract

The NVIDIA-Certified™ Storage program supports AI factories with rigorous performance testing, ensuring high-quality, efficient data solutions. This complements NVIDIA-Certified Systems and NVIDIA Enterprise Reference Architectures, empowering partners and customers to deploy AI infrastructures confidently, accelerating innovation and reducing deployment risks..

TABLE OF CONTENTS

# Executive summary

Data is the foundation of AI applications. It started with perception AI — understanding images, words and sounds. Then generative AI — creating text, images and sound. And now, we're entering the era of agentic AI and physical AI — and data is more important than ever as the pace of the AI revolution continues to accelerate.

The AI surge has created unprecedented demand for accelerated AI data centers, or AI Factories. The growth of computing power has been driven by increasing model complexity, multi-modal techniques, and more complex reasoning. This demand is compounded by an exponential growth in data. According to IDC's 2024 Global DataSphere Forecast, enterprises will generate 317 zettabytes of data annually by 2028 — including the creation of 29 zettabytes of unique data — of which 78% will be unstructured data and 44% of that will be audio and video.

As enterprises build AI factories to address these needs, the importance of this data cannot be overstated: access to high quality data directly impacts the performance and reliability of AI models. Data is essential for developing and optimizing AI applications, and it must be fed across all stages of the AI pipeline, from model building to training, tuning, and inference, with varied storage requirements at each stage. It's the fuel for the AI factory.

NVIDIA-Certified™ Systems and NVIDIA Enterprise Reference Architectures (RAs) are the building blocks for AI factories, bringing together accelerated computing, networking, and software for high-performance enterprise AI deployments at scale. The NVIDIA-Certified Storage program is designed to support the massive data demands of enterprise AI factories by offering a comprehensive storage certification that plugs into these programs. This empowers partners and customers to deploy to build AI factories that efficiently leverage data for faster, more accurate, and reliable AI models, unlocking new possibilities that were previously unimaginable.

The NVIDIA-Certified Storage program represents a significant step forward in ensuring that enterprise AI factories are built on a foundation of high-performance, reliable storage solutions. By rigorously validating storage systems against real-world AI workloads and integrating them seamlessly with NVIDIA's accelerated computing, networking, and software, customers can confidently deploy data hungry AI factories at scale. As AI continues to transform industries, the comprehensive NVIDIA-Certified program for accelerated computing systems and enterprise storage will play a crucial role in driving business innovation, reducing deployment risks, and achieving the potential of AI across the enterprise landscape.


The NVIDIA- Certified Storage program offers two levels of certification: Foundation and Enterprise. These storage certifications then integrate seamlessly with corresponding NVIDIA Enterprise RAs, which follow a prescriptive design pattern and scalable units of clustered NVIDIA-Certified systems to ensure optimal performance. The Foundation level storage certification certifies storage partners for NVIDIA's PCIe Optimized Reference Configurations for design patterns like 2-4-3-200 and 2-8-5-200, which specify 2 CPU sockets, 4/8 GPUs, 3/5 network adapters, and 200 GbE total bandwidth per GPU. The larger scale Enterprise level storage certification certifies storage partners for the HGX Reference Configuration following the 2-8-9-400 design pattern for larger-scale systems with 2 CPUs, 8 GPUs, 9 network adapters, and 400 GbE total bandwidth per GPU.

# NVIDIA-Certified Storage overview with NetApp

## Target audience

This solution is intended for organizations with AI/ML workloads that require deeper integration into broad data estates and traditional IT infrastructure tools and processes.

The target audience for the solution includes the following groups:

- IT and line of business decision makers planning for the most efficient infrastructure to deliver on AI/ML initiatives with the fastest time to market and ROI.

- Data scientists and data engineers who are interested in maximizing efficiency for critical data focused portions of the AI/ML workflow.

- IT architects and AI engineers who need to deliver a reliable and secure infrastructure that enables automated data workflows and compliance with existing data and process governance standards.

## NVIDIA-NetApp relationship

NetApp and NVIDIA have collaborated to provide enterprise customers with validated architectures for AI model training, fine tuning, and inferencing for over seven years with hundreds of joint customers. By certifying NetApp storage and management capabilities,  enterprises are empowered to use intelligent data infrastructure for AI innovation and effective data strategies. The intelligent data infrastructure from NetApp and NVIDIA's accelerated computing together deliver a comprehensive foundation for enterprise AI, covering edge, core, and cloud.

# NetApp AFF A90 with NetApp ONTAP

The NetApp AFF A90 powered by NetApp ONTAP® data management software provides built-in data protection, anti-ransomware capabilities, and high performance, scalability and resiliency required to support the most critical business workloads. It eliminates disruptions to mission-critical operations, minimizes performance tuning, and safeguards your data from ransomware attacks. For detailed information about NetApp AFF A-series systems please visit - https://www.netapp.com/aff-a-series/

## NetApp Storage Reference Architecture

NetApp storage reference architecture supports high performance, scalability, and flexibility for AI workloads by optimizing data pipelines from ingestion to model training, fine tuning, and inference using NetApp data fabric. NetApp ONTAP and AI solutions manage large datasets across training, fine-tuning, and inference. NetApp provides a unified, scalable platform to accelerate AI deployment with pre-configured infrastructure, ensuring seamless integration with NVIDIA Enterprise Reference Architectures ( Enterprise RAs).

# NVIDIA Enterprise RA with NetApp AFF A90

NVIDIA Enterprise RAs provide recommendations for deploying AI and data analytics solutions in enterprise environments. It integrates NVIDIA accelerated computing, networking, software stacks, and partner ecosystems to offer a complete AI infrastructure stack. NetApp All Flash FAS (AFF) systems meet enterprise storage needs with superior data management, enhancing business speed without sacrificing efficiency, reliability, or flexibility. NetApp recommends AFF A90 for scaling with NVIDIA Enterprise RAs, allowing enterprises to grow from small setups to large GPU farms. This paper covers scaling storage along with GPUs as outlined in the NVIDIA Enterprise RA Whitepaper. Four compute nodes form a scale unit (SU), and NVIDIA recommends scaling by SU rather than by individual compute nodes. From a connectivity side, each NetApp AFF A90 HA pair will be equipped with X50131A dual port adapters. The

quantity of these adapters is detailed in tables for each RA based on the number of scale units (SUs). The breakout cables used in the NVIDIA Enterprise RA will transition from 800GbE on the switch port to 2x 400GbE transceiver ports on the same switch port, then to 2x 200GbE storage ports.

## NVIDIA Enterprise RA 2-4-3-200

2-4-3-200 is a PCIe optimized NVIDIA-certified compute node supporting four GPUs (NVIDIA Blackwell RTX™ PRO 6000 Server Edition, L40S, H100 NVL, H200 NVL), three network adapters (1x NVIDIA ConnectX®-7 NIC for N-S and 2x NVIDIA BlueField®-3 SuperNICs), and two CPUs. This can scale from 1 SU up to 8 SUs in a cluster. Table 1 shows storage scaling requirements for up to eight scale units and 128 GPUs. Each scalable unit provides the following connectivity components:

o  For the E-W fabric: On 4 servers (SU), each equipped with 2x B3140H BlueField-3 SuperNICs (1x 400GbE port), providing a total of 8x 400Gb/s connections and an aggregate bandwidth of 3.2Tb/s

o  For the Converged N-S fabric for storage: On 4 servers, each equipped with 1 dual port ConnectX-7 NIC (2x 200GbE ports per server), providing a total of 8x 200Gb/s connections and an aggregate bandwidth of 1.6Tb/s

o  For the storage side: Each NetApp AFF A90 HA pair will be equipped with 4x X50131A dual port adapters

o  Eight storage ports are used for 2-4-3-200 configuration, each port providing 200Gb/s bandwidth

**Note**: There will be instances where both links on the breakout cable will not be used to get best redundancy for the N-S (storage) traffic

**Table 1) NVIDIA Enterprise RA 2-4-3-200 with NetApp AFF A90**

| Compute Counts | | | Storage Counts | Switch counts | | Converged Network Allocated Ports | | Storage - Transceiver counts | | Cable Counts |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | HA Pair to Leaf | | Storage |
| Nodes | GPUs | SUs | A90 HA Pair | Leaf | Spine | Switch To Storage - 800G | Ports Per HA pair - 200G | Storage - 200G | Switch - 800G | Node-to-leaf - 4 * 200G Breakouts |
| 4 | 16 | 1 | 1 | Converged in E-W Net | | 4 | 8 | 8 | 2 | 4 |
| 8 | 32 | 2 | 1 | Converged in E-W Net | | 4 | 8 | 8 | 2 | 4 |
| 12 | 48 | 3 | 1 | Converged in E-W Net | | 4 | 8 | 8 | 2 | 4 |
| 16 | 64 | 4 | 1 | Converged in E-W Net | | 4 | 8 | 8 | 2 | 4 |
| 20 | 80 | 5 | 1 | Converged in E-W Net | | 4 | 8 | 8 | 2 | 4 |
| 24 | 96 | 6 | 1 | Converged in E-W Net | | 4 | 8 | 8 | 2 | 4 |
| 28 | 112 | 7 | 1 | Converged in E-W Net | | 4 | 8 | 8 | 2 | 4 |
| 32 | 128 | 8 | 1 | Converged in E-W Net | | 4 | 8 | 8 | 2 | 4 |

## NVIDIA Enterprise  RA 2-8-5-200

2-8-5-200 is a PCIe optimized, NVIDIA-Certified compute node supporting eight PCIe GPUs (such as NVIDIA Blackwell RTX™ PRO 6000 Server Edition, L40S, H100 NVL, H200 NVL), five network adapters (1x NVIDIA ConnectX-7 NIC and 4x BlueField-3 SuperNICs), and two CPUs. It scales from 1 SU to 8 SUs in a cluster, doubling the GPU count compared to the 2-4-3 configuration. Supported devices are listed in the NVIDIA Enterprise Reference Architecture Whitepaper. Table 2 shows storage scaling for up to eight scale units and 256 GPUs. Each scalable unit provides the following connectivity building blocks:

o  For the E-W fabric: On 4 servers (SU), each equipped with 4x B3140H BlueField-3 SuperNICs (1x 400GbE port), providing a total of 16x 400Gb/s connections and an aggregate bandwidth of 6.4Tb/s

- For the Converged N-S fabric for storage: On 4 servers, each equipped with 1 dual port ConnectX-7 NIC (2x 200GbE ports per server), providing a total of 8x 200Gb/s connections and an aggregate bandwidth of 1.6Tb/s

- For the storage side, each NetApp AFF A90 HA pair will have X50131A dual port adapters, and the number of adapters are listed in Table 2 below based on the amount of SUs

- Eight storage ports are used for 2-4-5-200 configuration, each port providing 200Gb/s bandwidth

**Note**: There will be instances where both links on the breakout cable will not be used to get best redundancy for the N-S (storage) traffic

**Table 2) NVIDIA Enterprise RA 2-8-5-200 with NetApp AFF A90**

| Compute counts | | | Storage counts | Switch counts | | Converged Network Allocated Ports | | Storage - Transceiver counts | | Cable Counts |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | | | | | HA Pair to Leaf | | Storage Node-to-leaf - 4 * 200G Breakouts |
| Nodes | GPUs | SUs | A90 HA Pair | Leaf | Spine | Switch To Storage - 800G | Ports Per HA pair - 200G | Storage - 200G | Switch - 800G | |
| 4 | 32 | 1 | 1 | Converged in E-W Net | | 4 | 8 | 8 | 2 | 4 |
| 8 | 64 | 2 | 1 | Converged in E-W Net | | 4 | 8 | 8 | 2 | 4 |
| 12 | 96 | 3 | 1 | Converged in E-W Net | | 4 | 8 | 8 | 2 | 4 |
| 16 | 128 | 4 | 1 | Converged in E-W Net | | 4 | 8 | 8 | 2 | 4 |
| 20 | 160 | 5 | 2 | Converged in E-W Net | | 8 | 8 | 8 | 4 | 8 |
| 24 | 192 | 6 | 2 | 2 | N/A | 8 | 8 | 8 | 4 | 8 |
| 28 | 224 | 7 | 2 | 2 | N/A | 8 | 8 | 8 | 4 | 8 |
| 32 | 256 | 8 | 2 | 2 | N/A | 8 | 8 | 8 | 4 | 8 |

## NVIDIA Enterprise RA 2-8-9-400

2-8-9-400 is a PCIe Optimized 2-8-9-200 (CPU-GPU-Network Adaptor-Network Bandwidth) NVIDIA-Certified scale-out compute nodes that support eight PCIe GPUs (such as NVIDIA HGX™ H100, H200, or B200) and nine network adapters (1x ConnectX-7 NIC for N-S and 8x BlueField-3 SuperNICs) and two CPUs. This configuration supports scaling from 1 SU to 32 SUs in a cluster. Table 3 shows the storage scaling requirement for up to 32 scale units and 256 GPUs    ). Each scalable unit provides the following connectivity building blocks:

- On 4 servers, each with two dual port B3140H BlueField-3 SuperNICs (8 ports per server), providing a total of 32x 400Gb/s connections with an aggregate bandwidth of 12.8Tb/s

- For the Converged North-South fabric for storage: Each of the four servers is equipped with 1x dual-port ConnectX-7 adapter (each server has 2 ports of 200Gb). This configuration provides a total of 16 connections at 200Gb/s each, resulting in an aggregate bandwidth of 1.6Tb/s

- Eight storage ports are used for 2-4-5-200 configuration, each port providing 200Gb/s bandwidth

**Note**: There will be instances where both links on the breakout cable will not be used to get best redundancy for the N-S (storage) traffic

**Table 3) NVIDIA Enterprise RA 2-8-9-400 with NetApp AFF A90**

| Compute counts | | | Storage Counts | Switch counts | | Converged Network Allocated Ports | | Storage Transceiver counts | | Cable Counts |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | | | Max Ports Switch to Storage (800G) | Ports Per HA pair (200G) | HA Pair to Leaf | | 4 * 200G Breakouts - Storage to switch |
| Nodes | GPUs | SUs | A90 HA Pair | Leaf | Spine | | | Storage (200G) | Switch (800G) | |

| 4 | 32 | 1 | 1 | Converged in E-W Net | | 4 | 8 | 8 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|
| 8 | 64 | 2 | 1 | Converged in E-W Net | | 4 | 8 | 8 | 2 | 4 |
| 12 | 96 | 3 | 1 | Converged in E-W Net | | 4 | 8 | 8 | 2 | 4 |
| 16 | 128 | 3 | 1 | 2 | N/A | 4 | 8 | 8 | 2 | 4 |
| 24 | 192 | 3 | 2 | 2 | N/A | 8 | 8 | 8 | 4 | 8 |
| 32 | 256 | 8 | 2 | 2 | N/A | 8 | 8 | 8 | 4 | 8 |
| 48 | 384 | 12 | 3 | 2 | N/A | 12 | 8 | 8 | 6 | 12 |
| 64 | 512 | 16 | 3 | 2 | N/A | 16 | 8 | 8 | 6 | 12 |
| 96 | 768 | 24 | 4 | 5 | 3 | 24 | 8 | 8 | 8 | 16 |
| 128 | 1024 | 32 | 4 | 6 | 3 | 32 | 8 | 8 | 8 | 16 |

# Conclusion

The NVIDIA-Certified storage program with NetApp data management provides a scalable solution for enterprise AI workloads, addressing security, data management, and resource utilization. It ensures high availability without sacrificing storage efficiency, removing bottlenecks in AI pipelines and allowing data scientists and ML engineers to focus on innovation for quicker business results.

# Where to find additional information

To learn more about the information that is described in this document, review the following documents and/or websites:

NVA-1175 NVIDIA DGX SuperPOD with NetApp AFF A90 Storage Systems Deployment Guide

NetApp Documentation

NetApp AI Solutions Documentation

NetApp Install and Maintain AFF Storage Systems

NFS over RDMA

What is pNFS

# Version history

As an option, use the NetApp Table style to create a Version History table. Do not add a table number or caption.

| Version | Date | Document version history |
|---|---|---|
| Version 1.0 | June 2025 | New document |

Refer to the Interoperability Matrix Tool (IMT) on the NetApp Support site to validate that the exact product and feature versions described in this document are supported for your specific environment. The NetApp IMT defines the product components and versions that can be used to construct configurations that are supported by NetApp. Specific results depend on each customer's installation in accordance with published specifications.