

REFERENCE ARCHITECTURE

AIRI Pure Storage NVIDIA DGX BasePOD Reference Architecture and Configuration Guide



Contents

Abstract	3
Introduction	3
Network Architecture	3
Compute Fabric	5
Storage Fabric and In-Band Management Network	6
Out-of-Band Management Network	6
AIRI—An NVIDIA DGX BasePOD Certified Reference Architecture	7
NVIDIA Base Command	8
NVIDIA AI Enterprise	8
NVIDIA Base Command Manager	9
Major Components (Bill of Material)	10
DGX with NVIDIA Quantum InfiniBand Compute Fabric	10
DGX with RoCE Compute Fabric	11
Learn more about Pure Storage solutions for AI	13



Abstract

AIRI® is a certified NVIDIA DGX BasePOD reference architecture, representing the next evolution of the innovative and efficient, industry-first AI Ready Infrastructure (AIRI) from Pure Storage and NVIDIA. AIRI combines Pure Storage FlashBlade//S™ with the NVIDIA DGX BasePOD platform to provide the underlying infrastructure and software to accelerate deployment and execution of enterprise AI workloads.

This guide provides a prescriptive reference architecture for IT administrators to deploy a validated solution designed for inference, data exploration, and other computationally intensive algorithms that enable AI today. An AIRI system is built on the NVIDIA DGX BasePOD reference architecture and includes NVIDIA® DGX™ A100 or DGX H100 systems, Pure Storage® FlashBlade//S storage, NVIDIA networking, NVIDIA Base Command and NVIDIA AI Enterprise software. This guide is designed to support organizations looking to satisfy the needs of multiple workloads, from providing as-a-service access for small interactive jobs to supporting cluster-wide jobs that make full use of multi-GPU and multi-node resources.

Introduction

AI and other computationally intensive workloads have seen an unprecedented growth in the market. The need to scale out GPU-accelerated computing infrastructure to support this growth has become a significant challenge that IT administrators must address. AI is a fundamentally different workload than traditional enterprise applications running on CPU-based servers. AI necessitates consideration of specific networking, storage, and infrastructure management approaches proven to enable improved scalability, performance, and cost-effective manageability for AI workloads.

The AIRI DGX BasePOD reference architecture is an optimized AI infrastructure design consisting of NVIDIA DGX A100 or DGX H100 systems, [Pure Storage FlashBlade//S storage](#), and NVIDIA networking switches to support various job types. The AIRI design enables execution of multiple concurrent interactive jobs, multinode AI model training, inference, and data exploration work, all using NVIDIA AI software.

Because of its integrated nature, the AIRI solution is exceptionally scalable, allowing for flexible deployments based on workload, environmental, and budget constraints. In this Configuration Guide, we demonstrate four-node systems as an example.

NOTE: *We used two DGX A100 configurations as example configurations for validating this guide. DGX H100 is also certified for an AIRI solution. Contact an NVIDIA and Pure Storage partner for all possible configurations.*

Network Architecture

The AIRI system has three networks:

- **Compute fabric:** Connects the eight 200Gb/s NVIDIA ConnectX®-7 NICs from each DGX A100 through separate network planes for inter-node communication. You can configure them in InfiniBand or Ethernet mode.
- **Ethernet fabric:** Used for system management and storage. Uses one port from each dual-port ConnectX-7 NIC to connect with FlashBlade, cabled for 200GbE.
- **Out-of-band Ethernet network:** Connects the 1GbE RJ45 BMC port of each DGX A100 system to an additional Ethernet switch for node management and monitoring.

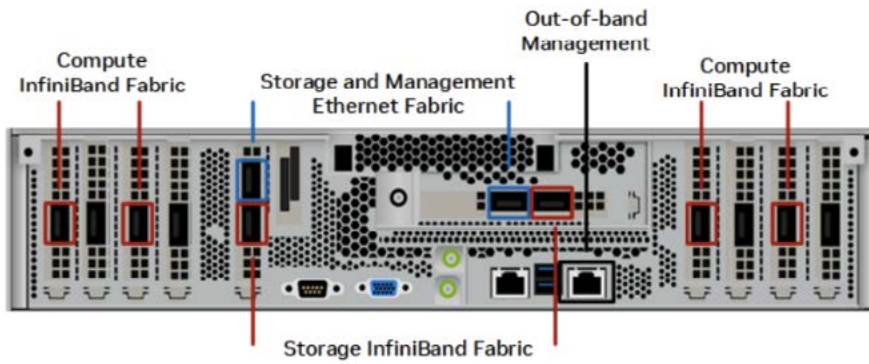


FIGURE 1 Network connections for DGX A100

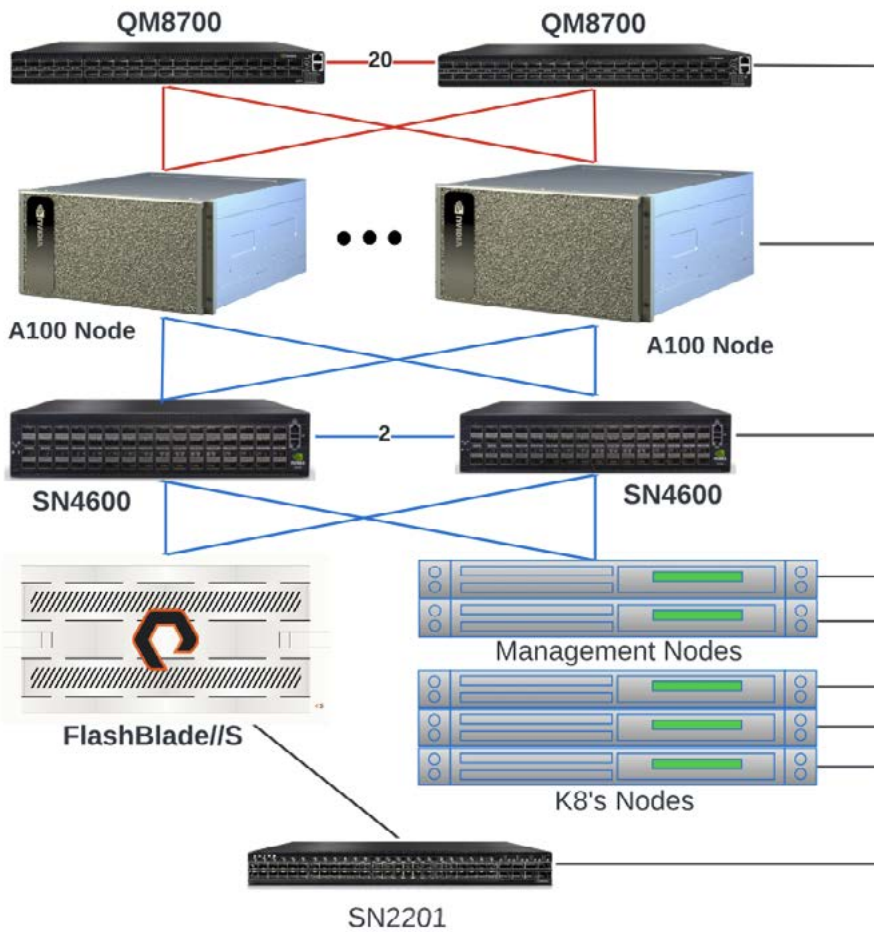


FIGURE 2 AIRI Network connectivity overview



Compute Fabric

For a DGX A100 system, the compute fabric may be NVIDIA InfiniBand (IB) or RDMA over Converged Ethernet (RoCE). This configuration is entirely independent of the DGX-to-FlashBlade fabric configuration. Each DGX A100 system has up to eight single port CX-7 connections for the compute fabric (Figure 1). The fabric design maximizes performance for typical communications traffic of AI workloads and provides some redundancy in the event of hardware failures and minimizing cost.

Choosing which type of compute fabric depends on the use case and architecture requirements.

This guide details both network configurations in an AIRI solution.

InfiniBand

If using IB mode, the compute fabric may utilize NVIDIA Quantum QM8700 200Gb/s InfiniBand switches. Each switch includes 40 QSFP56 ports used for communication to each DGX A100 system and between switches in the compute fabric. All connections are 200Gb/s, maximizing the bandwidth between network elements. No InfiniBand partitioning or other segmentation is used, with the QM8700 switch providing the subnet manager for the compute fabric. Connection to the out-of-band management ports on the switch to the out-of-band management fabric may be made if needed, but it is not critical to the AIRI system's operation.



FIGURE 3 NVIDIA Quantum QM8700 200 Gb/s InfiniBand Switch

RoCE

If using RoCE, the compute fabric may utilize NVIDIA Spectrum SN4600 Open Ethernet Switches. The 32 QSFP56 ports on compute fabric switches are configured for 200Gb/s operation and used for connecting to the DGX A100 systems. The switch is configured for RoCE traffic and uses VLANs for each discrete subnet configured on each DGX A100 NIC. Support is not available for All-All communication between GPUs in the cluster as it requires Layer 3 routing. This may be supported in a future update.



FIGURE 4 NVIDIA Spectrum SN4600 Open Ethernet Switch



Storage Fabric and In-Band Management Network

The storage fabric employs an Ethernet network fabric. In this example, the connections from the DGX compute nodes are NVIDIA ConnectX-7 200GbE connections.

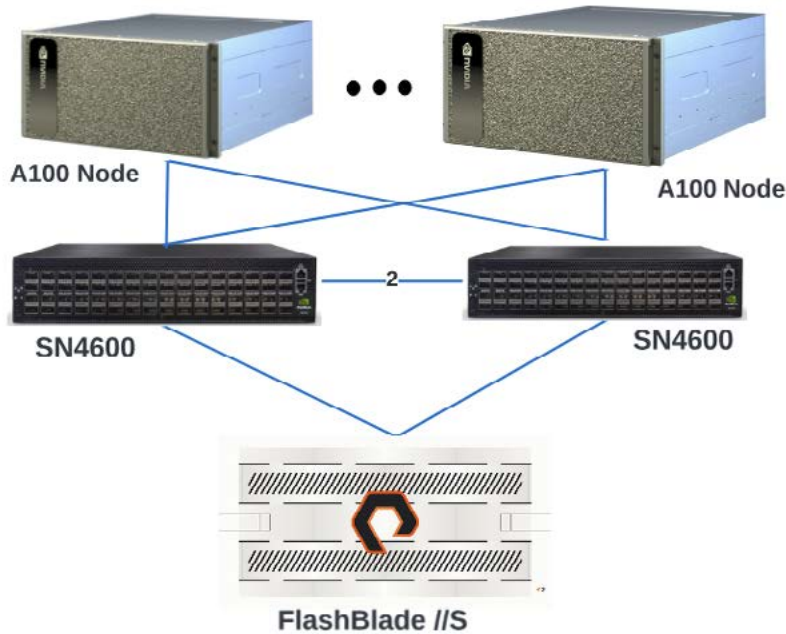


FIGURE 5 Storage fabric topology for a two-node AIRI BasePOD deployment

Each node has two connections originating on one port on each of the dual-port NVIDIA ConnectX[®]-7 NICs. For DGX, the recommended number of Inter-Peer Links is 4×200G, although 2×200G is also acceptable.

NOTE: *The AIRI configuration guide supports using Top-of-Rack switches from NVIDIA, Cisco, or Arista.*

The Storage and in-band Ethernet network provides connectivity for cluster services such as [Slurm](#) and [Kubernetes](#) and other services outside of the cluster, such as the NVIDIA NGC registry, code repositories, and data sources.

The in-band network may use the SN4600 switches (Figure 5) from the storage fabric. For most configurations, two SN4600 switches will be sufficient to support all the in-band and storage network connections. In-band management and storage networks are VLAN-separated in typical deployments. These switches may be connected directly to the data center core switch, with traffic routed to the storage or management network as needed to integrate with existing infrastructure and cluster access mechanisms.

Out-of-Band Management Network

The out-of-band Ethernet network is used for system management via the Baseboard Management Controller (BMC) and provides connectivity to manage all networking equipment. Out-of-band management is critical to the cluster's operation by providing low usage paths that ensure management traffic does not conflict with other cluster services.



The out-of-band management network is based on the Spectrum SN2201 switches (Figure 7), connecting up to 48 x 1G/100M/10M Base-T host-ports running the Cumulus Linux network operating system. These switches are usually also connected to the data center management network. It is recommended that the serial console of the out-of-band management network switches be connected to a console server in the data center to facilitate reconfiguration and connectivity in the event of a network failure.



FIGURE 6 1GbE NVIDIA SN2201 Data Center Open Ethernet Switch

AIRI—An NVIDIA DGX BasePOD Certified Reference Architecture

AIRI™, a DGX BasePOD certified reference architecture, is an evolution of the industry’s first complete AI-ready infrastructure. Architected by Pure Storage® and NVIDIA, and consisting of the latest NVIDIA DGX systems, NVIDIA networking and Pure Storage FlashBlade//S, AIRI is a complete hardware and software solution built on NVIDIA DGX BasePOD that maximizes AI results with the industry’s most efficient, disaggregated and modular scale-out storage platform.

NVIDIA DGX BasePOD is a prescriptive AI infrastructure for enterprises, eliminating the design challenges, lengthy deployment cycle and management complexity traditionally associated with scaling AI infrastructure. Certification ensures a validated, proven full-stack solution that is simple-to-use, enabling you to reap the benefits of AI right away with a fast, efficient, and future-proof infrastructure to meet your AI demands at enterprise scale. Powered by NVIDIA Base Command software, DGX BasePOD provides the essential foundation for AI development optimized for enterprise businesses.

FlashBlade//S is the ideal data storage platform for AI, as it was purpose-built from the ground up for modern, unstructured workloads and accelerates AI processes with the most efficient storage platform at every step of your data pipeline. A centralized data platform in a deep learning architecture increases the productivity of data scientists and makes scaling and operations simpler and more agile for the data architect.



FIGURE 7 AIRI stack



NVIDIA Base Command

NVIDIA Base Command is the operating system of the DGX data center, helping organizations speed the ROI of AI. Base Command powers the DGX platform, enabling organizations to leverage the best of NVIDIA software innovation. Enterprises can tap into the full potential of their DGX infrastructure with a proven platform that includes AI workflow management, enterprise-grade cluster management, libraries that accelerate compute, storage and network infrastructure, and system software optimized for running AI workloads.

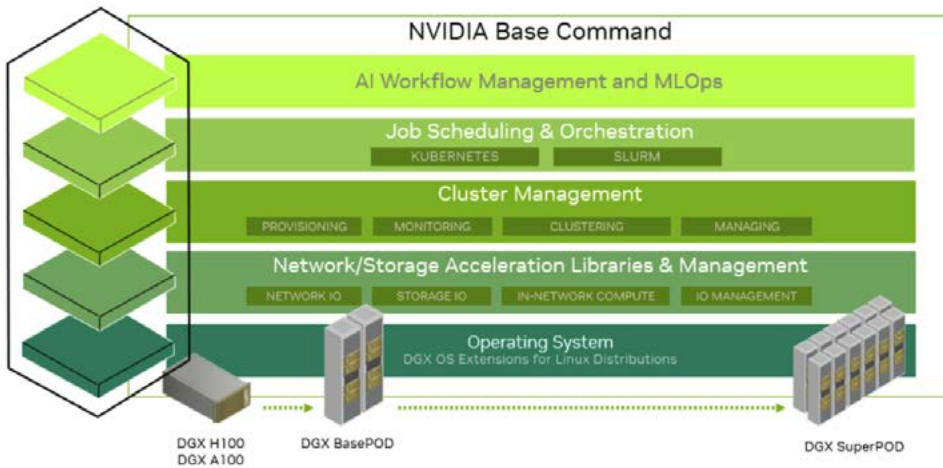


FIGURE 8 NVIDIA Base Command stack

NVIDIA AI Enterprise

NVIDIA AI Enterprise addresses the complexities of organizations trying to build and maintain their own high-performance, secure, cloud-native AI software platform. It includes the full AI software stack used for accelerating the data science pipeline and streamlining development and deployment of production AI, including generative AI, computer vision, speech AI, and more.

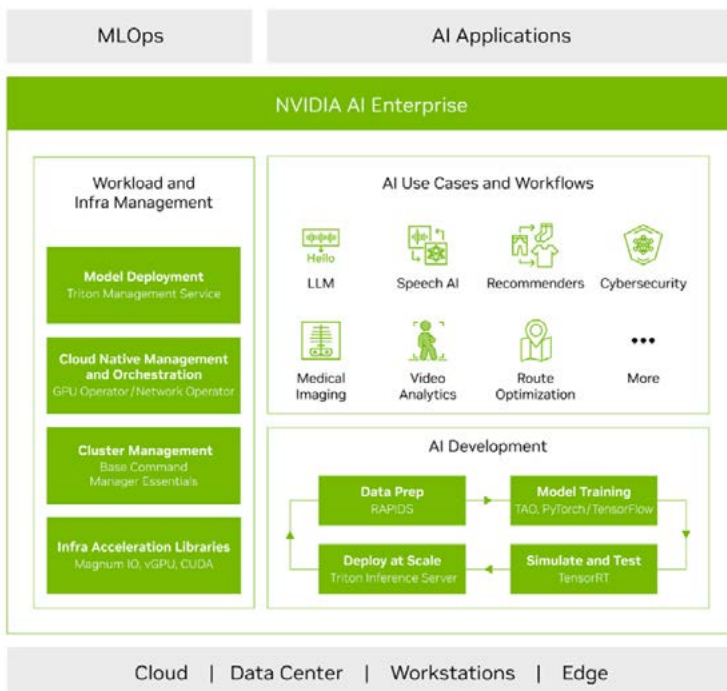


FIGURE 9 NVIDIA AI Enterprise



NVIDIA Base Command Manager

NVIDIA Base Command Manager is an integral part of the DGX BasePOD architecture providing cluster management, alerting, image deployment, job scheduling, allocation, among many other capabilities that help manage and maintain a cluster. The architecture described in this document was architected around a Base Command Manager configuration Type 1.

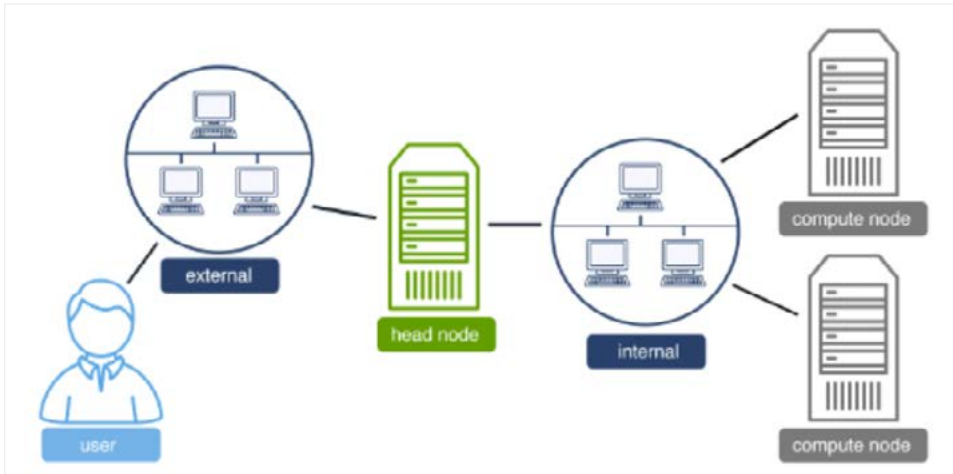


FIGURE 10 Base Command Manager Network Type 1

Storage network mount points to the FlashBlade were created using the Base Command Manager FS Mount function for each image category. In this example a value of `nconnect=8` was used to increase NFS connectivity to the FlashBlade from each compute node.

Example:

FS Mount list

node016 > FSMount list

FILESYSTEM	DEVICE	MOUNTPOINT	MOUNTOPTIO...
devpts	none	/dev/pts	gid=5,mode=620
proc	none	/proc	defaults,nosuid
sysfs	none	/sys	defaults
tmpfs	none	/dev/shm	defaults
nfs	\$localnfsserver:/c...	/cm/shared	rsize=32768,wsiz...
nfs	\$localnfsserver:/h...	/home	rsize=32768,wsiz...
nfs	192.168.125.50:/b...	/mnt/basepod	defaults,nconnect=8

FIGURE 11 Base Command Manager FS Mount list example



Major Components (Bill of Material)

When planning an AIRI deployment, you need to consider several factors to determine the total floorspace and resources needed. The design explicitly described in this configuration guide uses at minimum two data center racks: one as a utility rack for login, management, and storage servers, plus additional racks for the DGX A100 systems that comprise the compute muscle.

DGX with NVIDIA Quantum InfiniBand Compute Fabric

This section outlines the major components of a four-node AIRI system with NVIDIA Quantum InfiniBand compute fabric.

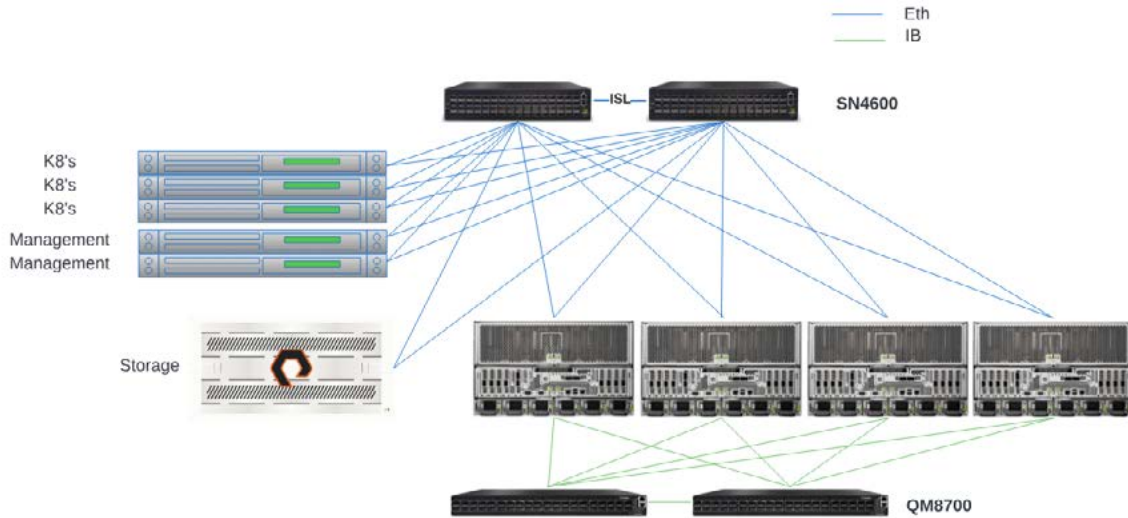


FIGURE 12 Logical topology for a four-node AIRI system with NVIDIA Quantum InfiniBand compute fabric

Compute

Quantity	SKU	Description
4	DGXA-G327H+P2CMI36	NVIDIA DGX A100 system with eight 40 or 80 GB A100 GPUs, with Nvidia ConnectX-7 200 Gb/s InfiniBand, Up to 2x Dual-Port NVIDIA ConnectX-7 VPI 10/25/50/100/200Gb/s Ethernet
Min 1 BCM	Varies	BasePOD cluster management servers are commonly configured in an HA pair



Compute Networking

Quantity	SKU	Description
2	920-9B110-00FH-0MD	QM8700 InfiniBand switches, compute fabric
16	MFS1S00-HxxV	AOC from DGX A100 to compute fabric leaf Switch
8	MCA1J00-HxxxE	200Gb/s InfiniBand DAC from leaf to leaf of compute fabric

Storage and In-Band Networking

Quantity	SKU	Description
1	STOR-ENDPOINT	Pure Storage FlashBlade //S500*
2	920-9N302-00FA-0C0	SN4600 switches for in-band management and storage
2	MCP1650-VxxxE26	200 GbE DAC for in-band fabric ISL
8	MFS1S00-Cxxx	200 GbE AOC for DGX A100 systems for in-band and storage
4	MFA1A00-Cxxx	100 GbE AOC for management servers to in-band management switches

Out of Band Network

Quantity	SKU	Description
1	920-9N110-00F1-0C0	SN2201 switches for OOB management
27	No specific requirement	1 GbE Cat 6 cables for OOB management system to switch
2	MFA1A00-Cxxx	100 GbE cables OOB to in-band switches

DGX with RoCE Compute Fabric

This section lists the major components of a four-node AIRI system with dedicated Ethernet (RoCE) compute fabric.

In this example SN4600 are used, however an alternate configuration using the NVIDIA Spectrum SN5000 series switches could also be used to take advantage of its 51.2Tb/s switching and routing capability, 30 percent lower power consumption and a flexible port configuration supporting: 64× 800GbE, 128× 400GbE, and 256× 200/100GbE ports.

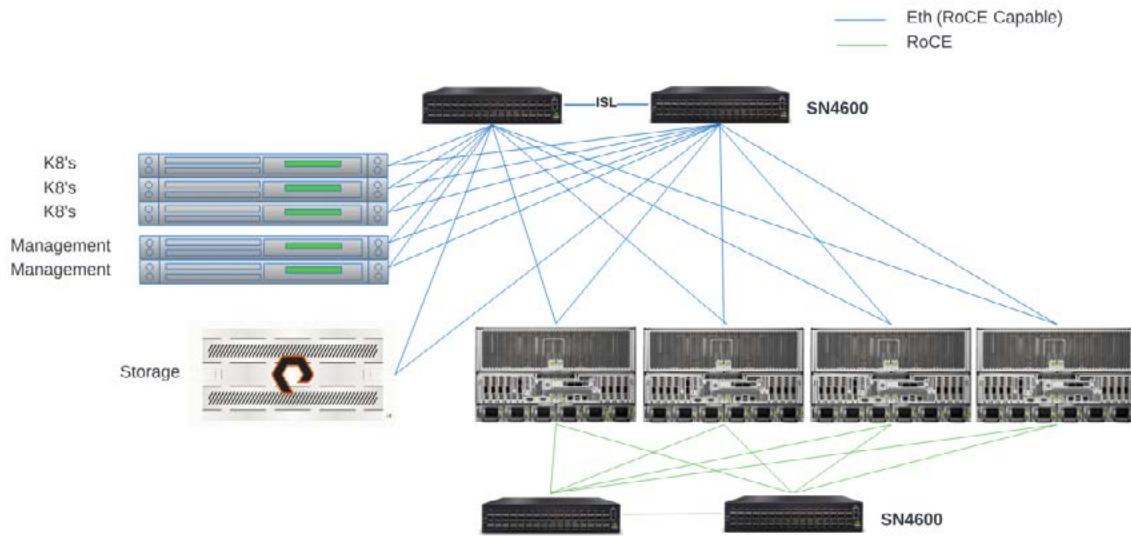


FIGURE 13 Logical topology for a four-node AIRI system with dedicated RoCE compute fabric

Compute

Quantity	SKU	Description
4	DGXA-G327H+P2CMI36	NVIDIA DGX A100 system with eight 40 or 80 GB A100 GPUs with Nvidia ConnectX-7 200 Gb/s InfiniBand, Up to 2x Dual-Port NVIDIA ConnectX-7 VPI 10/25/50/100/200Gb/s Ethernet
Min 1 BCM	Varies	BasePOD cluster management servers are commonly configured in an HA pair

Compute Networking

Quantity	SKU	Description
2	920-9N302-00FA-0C0	SN4600 switches for in-band management and storage
2	MCP1650-VxxxE26	200 GbE DAC for in-band fabric ISL
16	MFS1S00-Cxxx	200 GbE AOC for DGX A100 systems for in-band and storage



Storage and In-Band Networking

Quantity	SKU	Description
1	STOR-ENDPOINT	Pure Storage FlashBlade//S500*
2	920-9N302-00FA-0C0	SN4600 switches for in-band management and storage
2	MCP1650-VxxxE26	200 GbE DAC for in-band fabric ISL
8	MFS1S00-Cxxx	200 GbE AOC for DGX A100 systems for in- band and storage
4	MFA1A00-Cxxx	100 GbE AOC for management servers to in-band management switches

Out of Band Networking

Quantity	SKU	Description
1	920-9N110-00F1-0C0	SN2201 switches for OOB management
27	No specific requirement	1 GbE Cat 6 cables for OOB management system to switch
2	MFA1A00-Cxxx	100 GbE cables OOB to in-band switches

Learn more about Pure Storage solutions for AI

- [Pure Storage AIRI](#)
- [Pure Storage AI Solutions](#)

* Both the //S500 and S//200 platforms are available as solutions for BasePod. Please contact your Pure Storage sales team to determine the best match for your workload requirements.

purestorage.com

800.379.PURE

