



NVIDIA DGX SuperPOD: Next Generation Scalable Infrastructure for AI Leadership

Reference Architecture

Featuring NVIDIA DGX B200 Systems

Abstract

NVIDIA DGX SuperPOD™ with NVIDIA DGX™ B200 systems is the next generation of data center architecture for artificial intelligence (AI). Designed to provide the levels of computing performance required to solve advanced computational challenges in AI, high performance computing (HPC), and hybrid applications where the two are combined to improve prediction performance and time-to-solution. DGX SuperPOD is based upon the infrastructure built at NVIDIA for internal research purposes and is designed to solve the most challenging computational problems of today. Systems based on the DGX SuperPOD architecture have been deployed at customer data centers and cloud-service providers around the world.

Two key tenants of DGX SuperPOD are to embody the best combination of technologies available to be the premiere platform for AI computing, and to be designed in a manner that allows predictable scaling to fit workloads of different sizes. To make DGX SuperPOD the platform of choice for AI computing, DGX SuperPOD is powered by several key NVIDIA technologies, including:

- > NVIDIA DGX B200 system—to provide the most powerful computational building block for AI and HPC.
- > NVIDIA NDR (400 Gbps) InfiniBand—bringing the highest performance, lowest latency, and most scalable network interconnect.



- > NVIDIA NVLink® technology—networking technologies that connect GPUs at the NVLink layer to provide unprecedented performance for most demanding communication patterns.

The DGX SuperPOD architecture integrates NVIDIA software solutions including NVIDIA Base Command™, NVIDIA AI Enterprise, CUDA, and NVIDIA Magnum IO™. These technologies help keep the system running at the highest levels of availability, performance, and with NVIDIA Enterprise Support (NVEX), keeps all components and applications running smoothly.

This reference architecture (RA) discusses the components that define the scalable and modular architecture of DGX SuperPOD. The system is built on the concept of scalable units (SU), each containing 32 DGX B200 systems, which provides for rapid deployment of systems of multiple sizes. This RA includes details regarding the SU design and specifics of InfiniBand, NVLink network, Ethernet fabric topologies, storage system specifications, recommended rack layouts, and wiring guides.

Contents

Key Components of DGX SuperPOD.....	2
NVIDIA DGX B200 System.....	2
NVIDIA InfiniBand Technology	3
Runtime and System Management.....	3
Components.....	4
Design Requirements	5
System Design	5
Compute Fabric.....	6
Storage Fabric (High Speed Storage).....	6
In-Band Management Network.....	6
Out-of-Band Management Network.....	6
Storage Requirements.....	7
High-Performance Storage	7
User Storage.....	7
DGX SuperPOD Architecture	8
Network Fabrics	11
Compute Fabric	12
InfiniBand Storage Fabric.....	13
Ethernet Storage Fabric.....	14
In-Band Management Network.....	15
Out-of-Band Management Network.....	16
Storage Architecture.....	17
DGX SuperPOD Software.....	20
NVIDIA Base Command	20
NVIDIA NGC	21
NVIDIA AI Enterprise.....	21
Run:ai.....	21
Summary	22
Appendix A. Major Components.....	iii

Key Components of DGX SuperPOD

The DGX SuperPOD architecture has been designed to maximize performance for state-of-the-art model training, scale to exaflops of performance, provide the highest performance to storage and support all customers in the enterprise, higher education, research, and the public sector. It is a digital twin of the main NVIDIA research and development system, meaning the company's software, applications, and support structure are first tested and vetted on the same architecture. By using SUs, system deployment times are reduced from months to weeks. Leveraging the DGX SuperPOD design reduces time-to-solution and time-to-market of next generation models and applications.

DGX SuperPOD is the integration of key NVIDIA components, as well as storage solutions from partners certified to work in the DGX SuperPOD environment.

NVIDIA DGX B200 System

The NVIDIA DGX B200 system (Figure 1) is an AI powerhouse that enables enterprises to expand the frontiers of business innovation and optimization. The DGX B200 system delivers breakthrough AI performance with the most powerful chips ever built, in an eight GPU configuration. The NVIDIA Blackwell GPU architecture provides the latest technologies that brings months of computational effort down to days and hours, on some of the largest AI/ML workloads.

Figure 1. DGX B200 system



Some of the key highlights of the DGX B200 system when compared to the DGX H200 system include:

- 72 petaFLOPS FP8 training and 144 petaFLOPS FP4 inference
- Fifth generation of NVIDIA NVLink.
- 1,440 GB of aggregated HBM3 memory

NVIDIA InfiniBand Technology

InfiniBand is a high-performance, low latency, RDMA capable networking technology, proven over 20 years in the harshest compute environments to provide the best inter-node network performance. **InfiniBand** continues to evolve and lead data center network performance.

The NDR generation InfiniBand, NDR, has a peak speed of 400 Gbps per direction with an extremely low port-to-port latency, and is backwards compatible with the previous generations of InfiniBand specifications. InfiniBand is more than just peak bandwidth and **low latency**. InfiniBand provides additional features to optimize performance including adaptive routing (AR), collective communication with SHARP™, dynamic network healing with SHIELD™, and supports several network topologies including fat-tree, Dragonfly, and multi-dimensional Torus to build the largest fabrics and compute systems possible.

Runtime and System Management

The DGX SuperPOD RA represents the best practices for building high-performance data centers. There is flexibility in how these systems can be presented to customers and users. NVIDIA Base Command Manager software is used to manage all DGX SuperPOD deployments.

DGX SuperPOD can be deployed on-premises, meaning the customer owns and manages the hardware as a traditional system. This can be within a customer's data center or co-located at a commercial data center, but the customer owns the hardware.

Components

The hardware components of DGX SuperPOD are described in Table 1. The software components are shown in Table 2.

Table 1. DGX SuperPOD hardware components by NVIDIA

Component	NVIDIA Technology	Description
Compute nodes	NVIDIA DGX B200 system with eight B200 GPUs	The world's premier purpose-built AI systems featuring NVIDIA B200 Tensor Core GPUs, fifth-generation NVIDIA NVLink, and fourth-generation NVIDIA NVSwitch™ technologies.
Compute fabric	NVIDIA Quantum QM9700 NDR 400 Gbps InfiniBand	Rail-optimized, non-blocking, full fat-tree network with eight NDR400 connections per system
InfiniBand Storage fabric	NVIDIA Quantum QM9700 NDR 400 Gbps InfiniBand	The fabric is optimized to match peak performance of the configured storage array
Ethernet Storage fabric	NVIDIA Spectrum-X SN5600 800 Gbps Ethernet	Optional storage fabric for ethernet based storage solutions
Compute/storage InfiniBand fabric management	NVIDIA Unified Fabric Manager Appliance, Enterprise Edition	NVIDIA UFM combines enhanced, real-time network telemetry with AI powered cyber intelligence and analytics to manage scale-out InfiniBand data centers
In-band management network	NVIDIA SN5600 switch	64 port 800 Gbps and up to 256 ports of 200 Gbps Ethernet switch providing high port density with high performance
In-band and Out-of-band (OOB) management network	NVIDIA SN2201 switch	48 port 1 Gbps Ethernet and 4 x 100 Gbps switch leveraging copper ports to minimize complexity

Table 2. DGX SuperPOD software components

Component	Description
NVIDIA Base Command Manager	Comprehensive AI infrastructure management for AI clusters. It automates provisioning and administration and supports cluster sizes into the thousands of nodes.
NVIDIA AI Enterprise	NVIDIA AI Enterprise is an end-to-end, cloud-native software platform that accelerates data science pipelines and streamlines development and deployment of production-grade co-pilots and other generative AI applications.
Magnum IO	Enables increased performance for AI and HPC
NVIDIA NGC	The NGC catalog provides a collection of GPU-optimized containers for AI and HPC
Slurm	A classic workload manager used to manage complex workloads in a multi-node, batch-style, compute environment
Run:ai	Cloud-native AI workload and GPU orchestration platform enabling fractional, full, and multi-node support for the entire enterprise AI lifecycle including interactive development environments, training and inference



Note: DGX SuperPOD only supports NVIDIA Base Command Manager with Slurm, or Run:ai

Design Requirements

DGX SuperPOD is designed to minimize system bottlenecks throughout the tightly coupled configuration to provide the best performance and application scalability. Each subsystem has been thoughtfully designed to meet this goal. In addition, the overall design remains flexible so that data center requirements can be tailored to better integrate into existing data centers.

System Design

DGX SuperPOD is optimized for a customers' particular workload of multi-node AI and HPC applications:

- > A modular architecture based on SUs of 32 DGX B200 systems each.
- > A fully tested system scales to four SUs, but larger deployments can be built based on customer requirements.
- > Single rack that can support two DGX B200 systems per rack, so that the rack layout can be modified to accommodate different data center requirements.

- > Storage partner equipment that has been certified to work in DGX SuperPOD environments.
- > Full system support—including compute, storage, network, and software—is provided by NVIDIA Enterprise Experience (NVEX).

Compute Fabric

- > The compute fabric is rail-optimized, balanced, full-fat tree topology
- > Managed NDR switches are used throughout the design to provide better management of the fabric.
- > The fabric is designed to support the latest SHaRP features.

Storage Fabric (High Speed Storage)

The storage fabric provides high bandwidth to shared storage. It also has the following characteristics:

- > It is independent of the compute fabric to maximize performance of both storage and application performance.
- > Provides single-node bandwidth of at least 40 GBps to each DGX B200 system.
- > Storage is provided over InfiniBand or RDMA over Converged Ethernet to provide maximum performance and minimize CPU overhead.
- > It is flexible and can scale to meet specific capacity and bandwidth requirements.
- > Connectivity to management nodes required to provide storage access independent of compute nodes.

In-Band Management Network

- > The in-band management network fabric is Ethernet-based and is used for node provisioning, data movement, Internet access, and other services that must be accessible by the users.
- > The in-band management network connections for compute and management nodes operate at 200 Gbps and are bonded for resiliency.

Out-of-Band Management Network

The OOB management network connects all the base management controller (BMC) ports, as well as other devices that should be physically isolated from users. The Switch Management Network is a subset of the Out-Of-Band Network that provides additional security and resiliency.

Storage Requirements

The DGX SuperPOD compute architecture must be paired with a high-performance, balanced, storage system to maximize overall system performance. DGX SuperPOD is designed to use two separate storage systems, high-performance storage (HPS) and user storage, optimized for key operations of throughput, parallel I/O, as well as higher IOPS and metadata workloads.

High-Performance Storage

High-Performance Storage is provided via InfiniBand connected storage from a DGX SuperPOD certified storage partner, and is engineered and tested with the following attributes in mind:

- > High-performance, resilient, POSIX-style file system optimized for multi-threaded read and write operations across multiple nodes.
- > RDMA on InfiniBand or Ethernet support
- > Local system RAM for transparent caching of data.
- > Leverage local flash device transparently for read and write caching.

The specific storage fabric topology, capacity, and components are determined by the DGX SuperPOD certified storage partner as part of the DGX SuperPOD design process.

User Storage

User Storage differs from High-Performance storage in that it exposes an NFS share on the in-band management fabric for multiple uses. It is typically used for “home directory” type usage (especially with clusters deployed with Slurm), administrative scratch space, and shared storage as needed by DGX SuperPOD components in a High Availability configuration (e.g., Base Command Manager), and log files.

With that in mind, User Storage has the following requirements:

- > 100 Gb/s connectivity is required.
- > Designed for high metadata performance, IOPS, and key enterprise features such as checkpointing. This is different than the HPS, which is optimized for parallel I/O and large capacity.
- > Communicate over Ethernet, using NFS.

User storage in a DGX SuperPOD is often satisfied with existing NFS servers already deployed, such that a new export is created and made accessible to the DGX SuperPOD’s in-band management network. User Storage is therefore not described in detail in this DGX SuperPOD reference architecture. However, we require 100 Gb/s minimum bandwidth for the user storage.

DGX SuperPOD Architecture

The DGX SuperPOD architecture is a combination of DGX systems, InfiniBand and Ethernet networking, management nodes, and storage. Figure 2 shows the rack layout of a single SU. In this example, power consumption per rack exceeds 25 kW. The rack layout can be adjusted to meet local data center requirements, such as maximum power per rack and rack layout between DGX systems and supporting equipment to meet local needs for power and cooling distribution.

Figure 2. Complete single SU rack layout

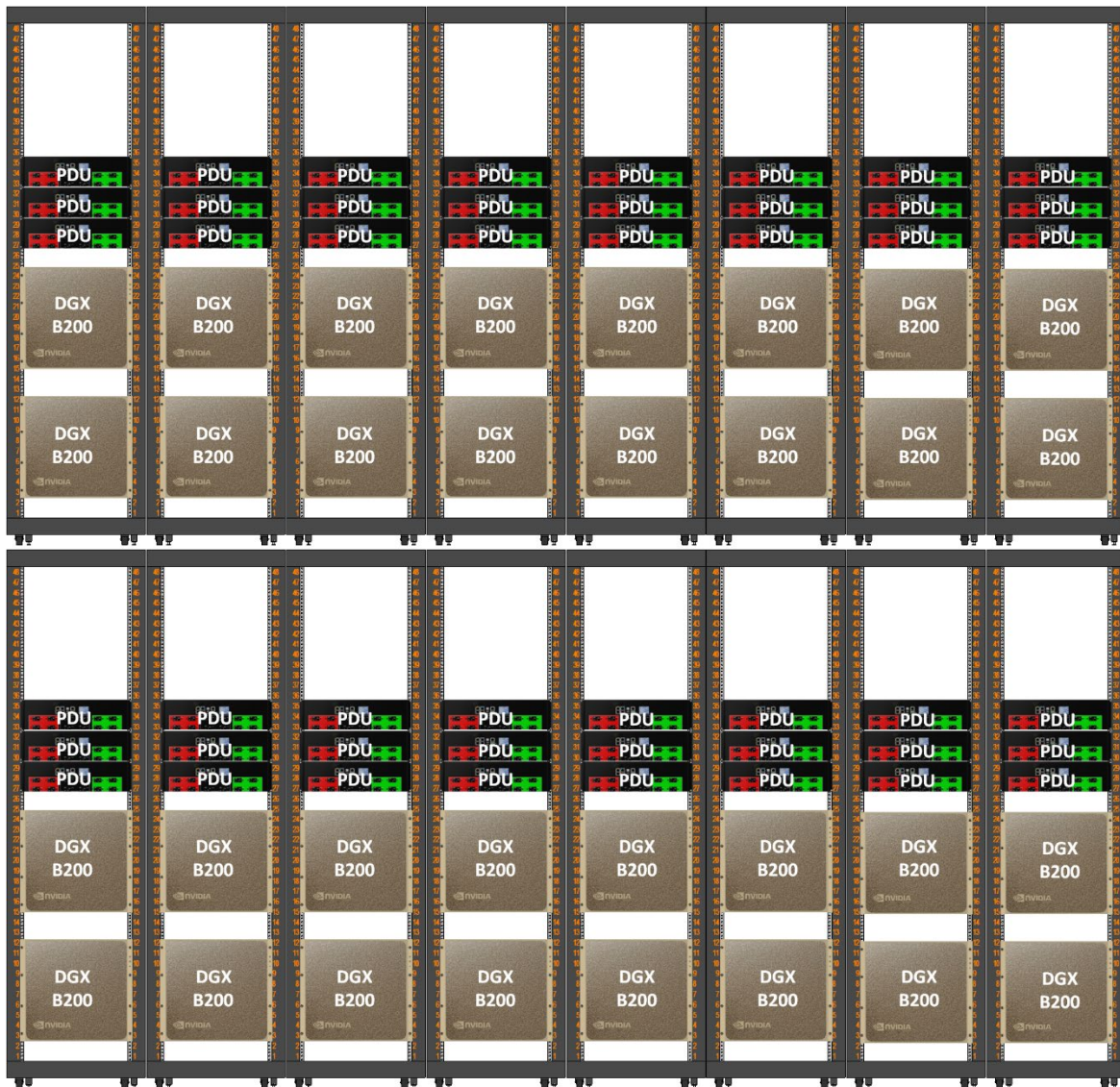
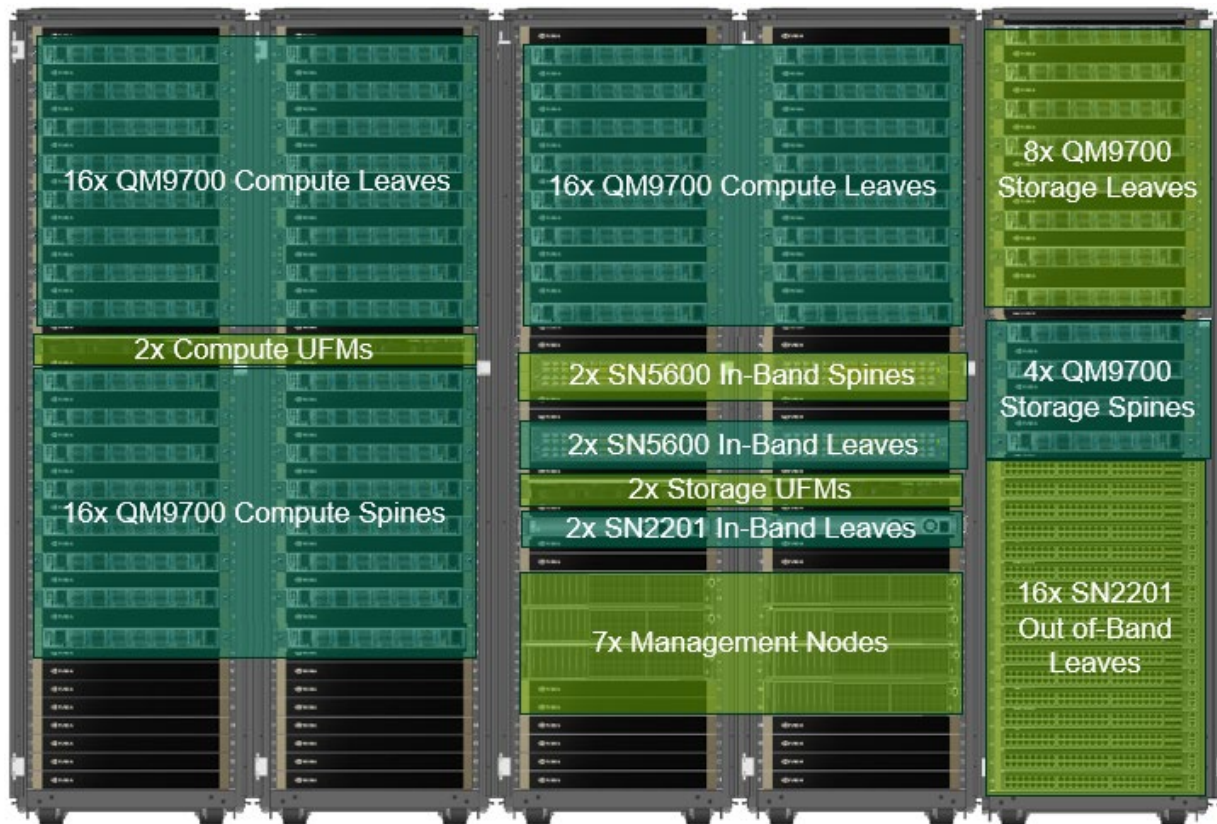


Figure 3 shows an example management rack configuration with networking switches, management servers, storage arrays, and UFM appliances. Sizes and quantities will vary depending upon models used.

Figure 3. Management rack configuration



This reference architecture is focused on 4 SU units with 128 DGX nodes. DGX SuperPOD can scale to much larger configurations up to and beyond 64 SU with 2000+ DGX B200 nodes. See Table 3 for more information.

Table 3. Larger DGX SuperPOD component counts

SU Count	Node Count	GPU Count	InfiniBand Switch count			Cable Count		
			Leaf	Spine	Core	Node-Leaf	Leaf-Spine	Spine-Core
2	64	512	16	8	--	512	512	--
4	128	1024	32	16	--	1024	1024	--
8	256	2048	64	32	--	2048	2048	--
16	512	4096	128	128	64	4096	4096	4096
32	1024	8192	256	256	128	8192	8192	8192
64	2048	16384	512	512	256	16384	16384	16384

Contact NVIDIA for information regarding DGX SuperPOD solutions of four scalable units or more.

Network Fabrics

Building systems by SU provides the most efficient designs. However, if a different node count is required due to budgetary constraints, data center constraints, or other needs, the fabric should be designed to support the full SU, including leaf switches and leaf-spine cables, and leave the portion of the fabric unused where these nodes would be located. This will ensure optimal traffic routing and ensure that performance is consistent across all portions of the fabric.

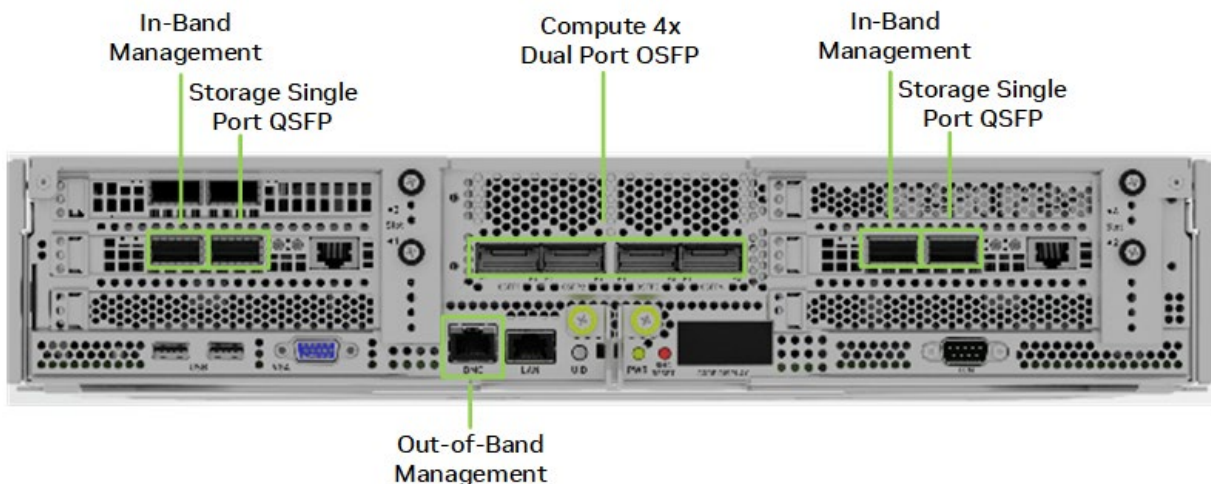
DGX SuperPOD configurations utilize four network fabrics:

- > Compute Fabric
- > Storage Fabric
- > In-Band Management Network
- > Out-of-Band Management Network

Each network is detailed in this section.

Figure 4 shows the ports on the back of the DGX B200 CPU tray and the connectivity provided. The compute fabric ports in the middle use a two-port transceiver to access all eight GPUs. Each pair of in-band management and storage ports provide parallel pathways into the DGX B200 system for increased performance. The OOB port is used for BMC access. (The LAN port next to the BMC port is not used in DGX SuperPOD configurations.)

Figure 4. DGX B200 network ports



Compute Fabric

Figure 5 shows the compute fabric layout for the full 127-node DGX SuperPOD. Each group of 32 nodes is rail-aligned. Traffic per rail of the DGX B200 systems is always one hop away from the other 31 nodes in a SU. Traffic between nodes, or between rails, traverses the spine layer.

Figure 5. Compute fabric for full 127-node DGX SuperPOD

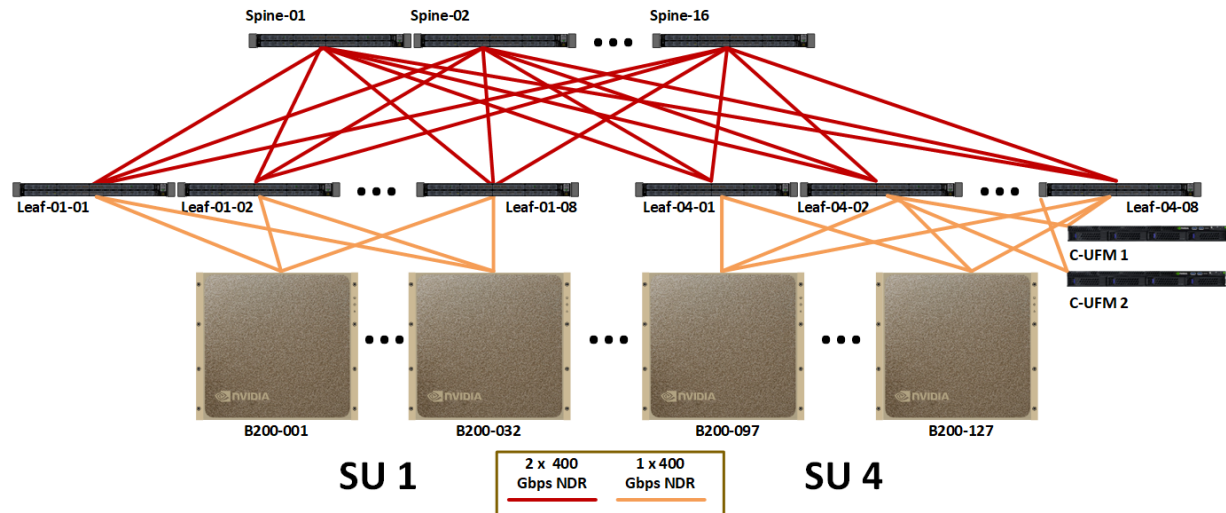


Table 4 shows the number of cables and switches required for the compute fabric for different SU sizes.

Table 4. Compute fabric component count

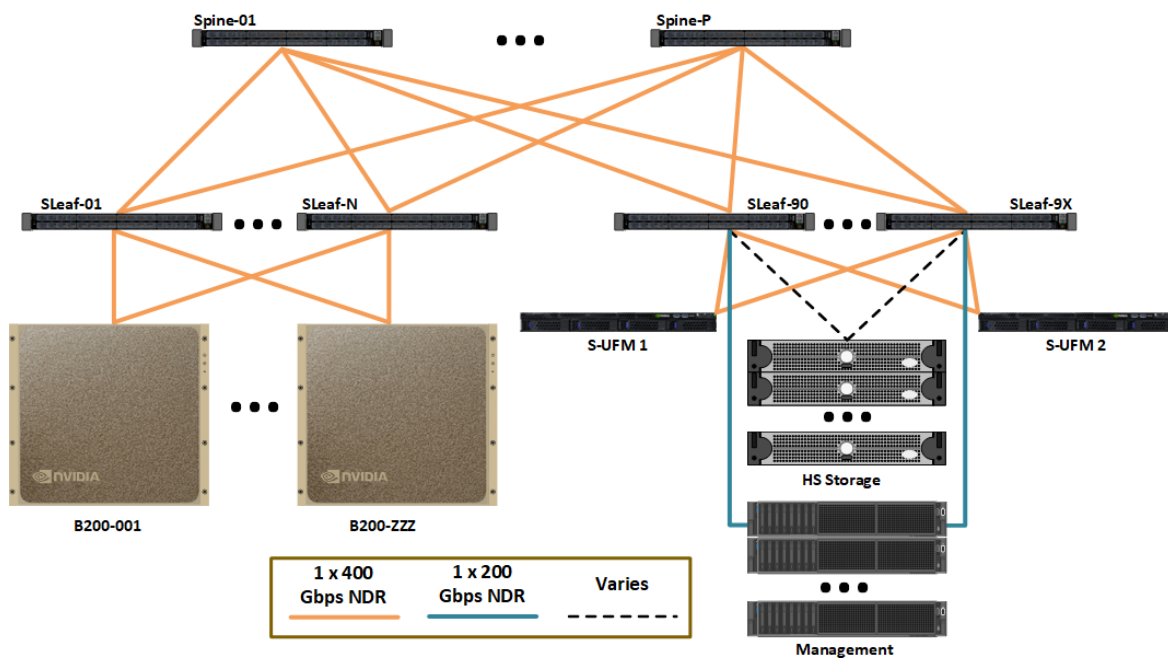
SU Count	Node Count	GPU Count	InfiniBand Switch Count		Cable Count	
			Leaf	Spine	Compute + UFM	Spine-Leaf
1	31 ¹	248	8	4	252	256
2	63	504	16	8	508	512
3	95	760	24	16	764	768
4	127	1016	32	16	1020	1024

1. This is a 32 node per SU design, however a DGX system must be removed to accommodate for UFM connectivity.

InfiniBand Storage Fabric

The storage fabric employs an InfiniBand network fabric that is essential to maximum bandwidth (Figure 6). This is because the I/O per-node for the DGX SuperPOD must exceed 40 GBps. High bandwidth- requirements with advanced fabric management features, such as congestion control and AR, provide significant benefits for the storage fabric.

Figure 6. Storage fabric logical design



The InfiniBand storage fabric uses [MQM9700-NS2F](#) switches (Figure 7). The high-speed storage devices are connected at a 1:1 port to uplink ratio. The DGX B200 system connections are slightly oversubscribed with a ratio near 4:3 with adjustments as needed to enable more storage flexibility regarding cost and performance.

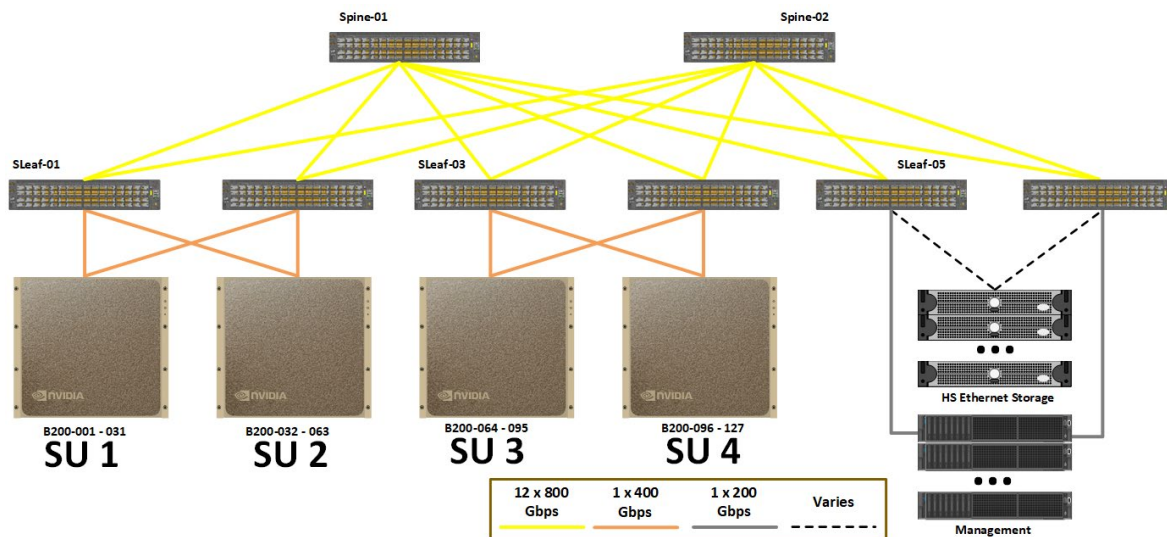
Figure 7. MQM9700-NS2F switch



Ethernet Storage Fabric

The Ethernet storage fabric employs a high-speed Ethernet network fabric that is essential to maximum bandwidth (Figure 8). This is because the I/O per-node for the DGX SuperPOD must exceed 40 GBps. High bandwidth- requirements with advanced fabric management features, provide significant benefits for the storage fabric. Supported ethernet storage appliance leverages RoCE to provide best performance and minimizes CPU usage.

Figure 8. Storage fabric logical design



The storage fabric uses SN5600 switches (Figure 9). The high-speed storage devices are connected at a 1:1 port to uplink ratio. The DGX B200 system connections are slightly oversubscribed with a ratio near 4:3 with adjustments as needed to enable more storage flexibility regarding cost and performance.

Figure 9. NVIDIA Spectrum SN5600 Ethernet Switch



In-Band Management Network

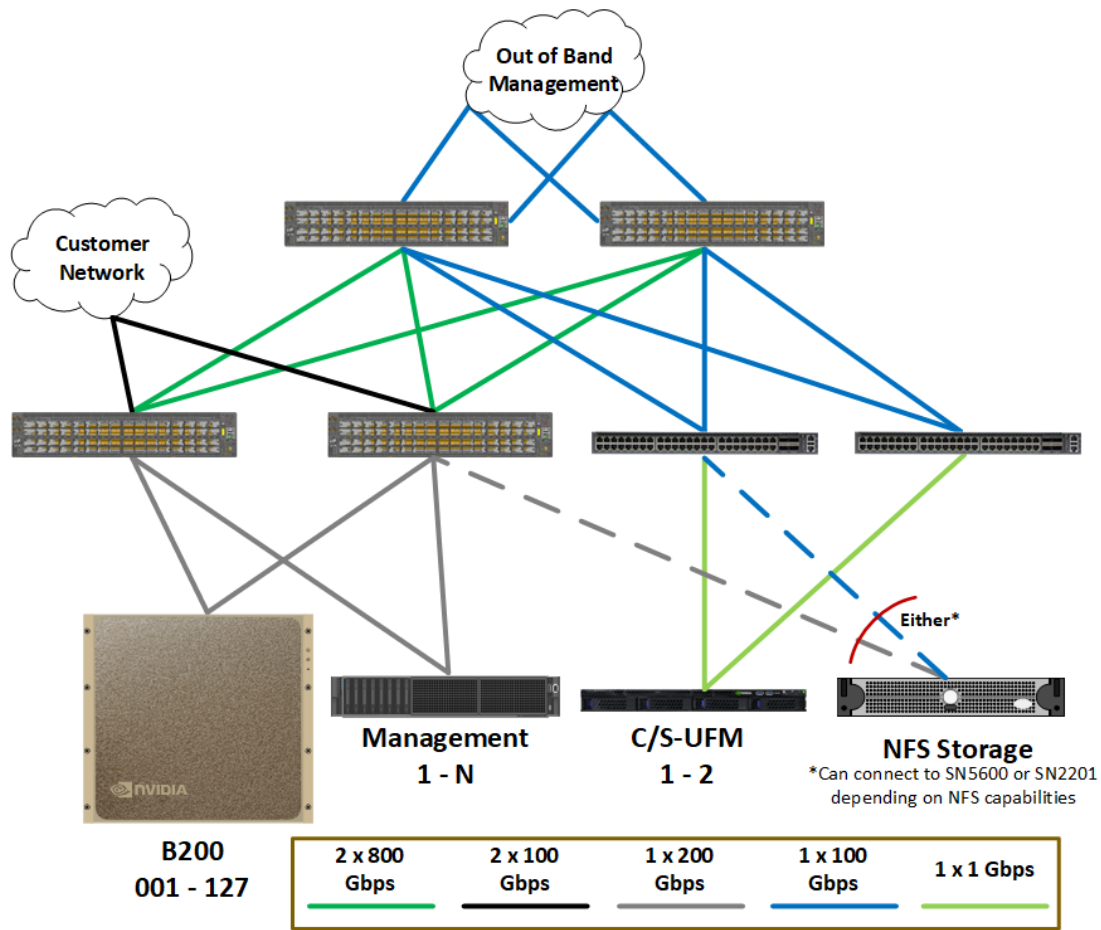
The in-band management network provides several key functions:

- > Connects all the services that manage the cluster.
- > Enables access to the data NFS tier.
- > Provides connectivity for the in-cluster services such as Base Command Manager, Slurm, Run:ai and to other services outside of the cluster such as the NGC registry, code repositories, and data sources.

Figure 10 shows the logical layout of the in-band Ethernet network. The in-band network connects the compute nodes and management nodes. In addition, the OOB network is connected to the in-band network to provide high-speed interfaces from the management nodes to support parallel operations to devices connected to the OOB storage fabric, such as storage.

The OOB fabric and the In-Band fabric are logically separated on the spine layer to ensure secure isolation for these networks.

Figure 10. In-band Ethernet network

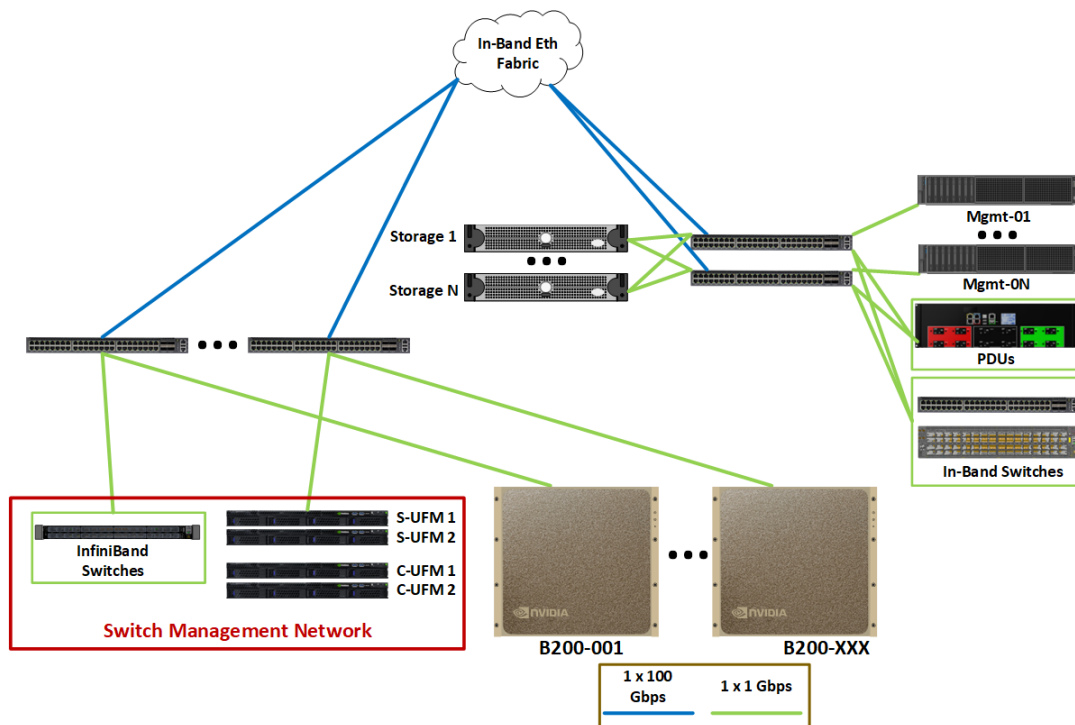


The in-band management network uses SN5600 and SN2201 switches (Figure 9 and 13)

Out-of-Band Management Network

Figure 12 shows the OOB Ethernet fabric. It connects the management ports of all devices including DGX and management servers, storage, networking gear, rack PDUs, and all other devices. These are separated onto their own fabric because there is no use-case where users need access to these ports and are secured using logical network separation. Figure 12 shows the Switch Management Network is a subset of the Out-Of-Band Network that provides additional security and resiliency.

Figure 12. Logical OOB management network layout



The OOB management network uses SN2201 switches (Figure 13).

Figure 13. SN2201 switch



Storage Architecture

Data, lots of data, is the key to development of accurate deep learning (DL) models. Data volume continues to grow exponentially, and data used to train individual models continues to grow as well. Data format, not just volume can play a key factor in the rate at which data is accessed so storage system performance must scale commensurately.

The key I/O operation in DL training is re-read. It is not just that data is read, but it must be reused again and again due to the iterative nature of DL training. Pure read performance still is important as some model types can train in a fraction of an epoch (ex: some recommender models) and inference of existing can be highly I/O intensive, much more so than training. Write performance can also be important. As DL models grow and time-to-train, writing checkpoints is necessary for fault tolerance. The size of checkpoint files can be terabytes in size and while not written frequently are typically written synchronously that blocks forward progress of DL models.

Ideally, data is cached during the first read of the dataset, so data does not have to be retrieved across the network. Shared filesystems typically use RAM as the first layer of cache. Reading files from cache can be an order of magnitude faster than from remote storage. In addition, the DGX B200 system provides local NVMe storage that can also be used for caching or staging data.

DGX SuperPOD is designed to support all workloads, but the storage performance required to maximize training performance can vary depending on the type of model and dataset. The guidelines in Table 5 and Table 6 are provided to help determine the I/O levels required for different types of models.

Table 5. Storage performance requirements

Level	Work Description	Dataset Size
Standard	Multiple concurrent LLM or fine-tuning training jobs and periodic checkpoints, where the compute requirements dominate the data I/O requirements significantly.	Most datasets can fit within the local compute systems' memory cache during training. The datasets are single modality, and models have millions of parameters.
Enhanced	Multiple concurrent multimodal training jobs and periodic checkpoints, where the data I/O performance is an important factor for end-to-end training time.	Datasets are too large to fit into local compute systems' memory cache requiring more I/O during training, not enough to obviate the need for frequent I/O. The datasets have multiple modalities and models have billions (or higher) of parameters.

Table 6. Guidelines for storage performance

Performance Characteristic	Good (GBps)	Better (GBps)
Single SU aggregate system read	40	125
Single SU aggregate system write	20	62
4 SU aggregate system read	160	500
4 SU aggregate system write	80	250

High-speed storage provides a shared view of an organization's data to all nodes. It must be optimized for small, random I/O patterns, and provide high peak node performance and high aggregate filesystem performance to meet the variety of workloads an organization may encounter. High-speed storage should support both efficient multi-threaded reads and writes from a single system, but most DL workloads will be read-dominant.

Use cases in automotive and other computer vision-related tasks, where high-resolution images are used for training (and in some cases are uncompressed) involve datasets that easily exceed 30 TB in size. In these cases, 4 GBps per GPU for read performance is needed.

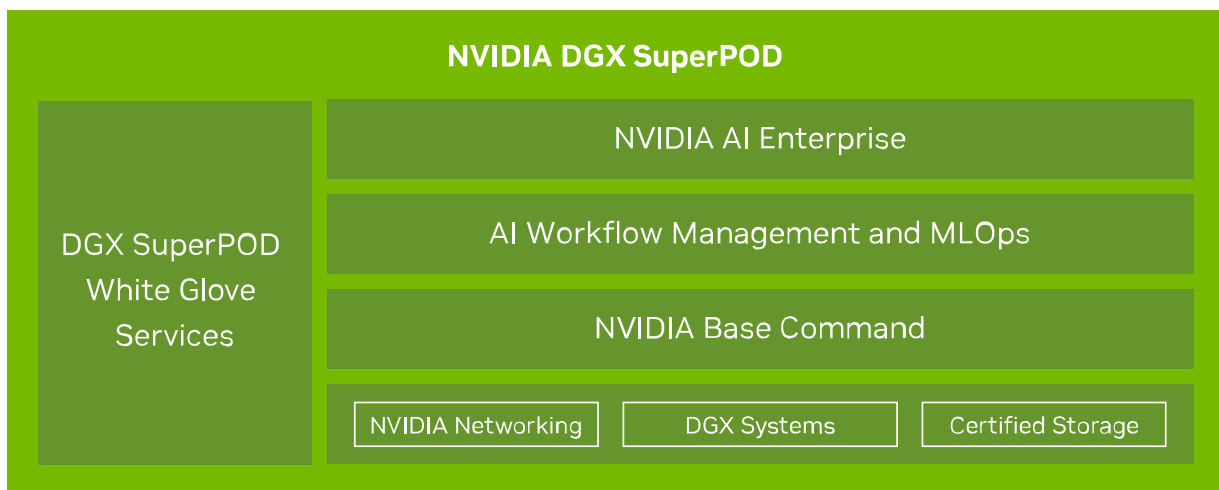
While NLP and LLM cases often do not require as much read performance for training, peak performance for reads and writes are needed for creating and reading checkpoint files. This is a synchronous operation and training stops during this phase. If you are looking for best end-to-end training performance, do not ignore I/O operations for checkpoints. Consider at least ½ of the read performance as recommended write performance for LLM and large model use cases.

The preceding metrics assume a variety of workloads, datasets, and need for training locally and directly from the high-speed storage system. It is best to characterize workloads and organizational needs before finalizing performance and capacity requirements.

DGX SuperPOD Software

DGX SuperPOD is an integrated hardware and software solution. The included software (Figure 14) is optimized for AI from top to bottom. From the accelerated frameworks and workflow management through to system management and low-level operating system (OS) optimizations, every part of the stack is designed to maximize the performance and value of DGX SuperPOD.

Figure 14. DGX SuperPOD high-level software architecture



NVIDIA Base Command

[NVIDIA Base Command](#) powers every DGX SuperPOD, enabling organizations to leverage the best of NVIDIA software innovation. Enterprises can unleash the full potential of their investment with a proven platform that includes enterprise-grade orchestration and cluster management, libraries that accelerate compute, storage and network infrastructure, and an OS optimized for AI workloads.

NVIDIA NGC

NGC provides software to meet the needs of data scientists, developers, and researchers with various levels of AI expertise.

Software hosted on NGC undergoes scans against an aggregated set of common vulnerabilities and exposures (CVEs), crypto, and private keys.

Software from the NGC catalog is tested and ensured to scale to multiple GPUs and in some cases, to scale to multi-node, ensuring users maximize the use of their DGX SuperPOD.

NVIDIA AI Enterprise

NVIDIA AI Enterprise is the end-to-end software platform that brings generative AI into reach for every enterprise, providing the fastest and most efficient runtime for generative AI foundation models developed with the NVIDIA DGX platform. With production-grade security, stability, and manageability, it streamlines the development of generative AI solutions. NVIDIA AI Enterprise is included with DGX SuperPOD for enterprise developers to access pretrained models, optimized frameworks, microservices, accelerated libraries, and enterprise support.

Run:ai

Run:ai is cloud native AI workload and GPU orchestration platform that simplifies and accelerates AI and machine learning with DGX SuperPOD through dynamic resource allocation, comprehensive AI lifecycle support, strategic resource management and advanced scheduling. Run:ai maximizes GPU efficiency and workload capacity. Its policy engine, open architecture, and visibility into AI workloads foster strategic alignment with business objectives. This results in increases in cluster efficiency and utilization, all with zero manual resource intervention, accelerating innovation and providing a scalable, agile, and cost-effective solution for enterprises.

Summary

DGX SuperPOD with NVIDIA DGX B200 systems is the next generation of data center scale architecture to meet the demanding and growing needs of AI training. This RA document for DGX SuperPOD represents the architecture used by NVIDIA for our own AI model and HPC research and development. DGX SuperPOD continues to build upon its high-performance roots to enable training of the largest NLP models, support the expansive needs of training models for automotive applications, and scaling-up recommender models for greater accuracy and faster turn-around-time.

DGX SuperPOD represents a complete system of not just hardware but all the necessary software to accelerate time-to-deployment, streamline system management, proactively identify system issues. The combination of all these components keeps systems running reliably, with maximum performance, and enables users to push the bounds of state-of-the-art. The platform is designed to both support the workloads of today and grow to support tomorrow's applications.

Appendix A. Major Components

Major components for the DGX SuperPOD configuration are listed in Table 7. These are representative of the configuration and must be finalized based on actual design.

Table 7. Major components of the 4 SU, 127-node DGX SuperPOD

Count	Component	Recommended Model
Racks		
70	Rack (Legrand)	NVIDPD13
Nodes		
127	DGX nodes	NVIDIA DGX systems
4	UFM appliance	NVIDIA Unified Fabric Manager Appliance
7	Management servers	Intel based x86 2 × Socket, 24 core or greater, 384 GB RAM, OS (2x480GB M.2 or SATA/SAS SSD in RAID 1), NVME 7.68 TB (raw), 4x NDR VPI Ports, TPM
Management Network		
4	In-band management	NVIDIA SN5600 Spectrum-4 based 800GbE 2U Open Ethernet switch with Cumulus Linux Authentication, 64 OSFP ports and 1 SFP28 port, 2 power supplies (AC), x86 CPU, Secure-boot, standard depth, C2P airflow, Tool-less Rail Kit, 920-9N42F-00RI-7C0
2	In-band management	NVIDIA SN2201 switch with Cumulus Linux, 48 RJ45 ports, P2C, 920-9N110-00F1-0C0
17	OOB management	NVIDIA SN2201 switch with Cumulus Linux, 48 RJ45 ports, P2C, 920-9N110-00F1-0C0
Compute Fabric		
48	Fabric switches	NVIDIA Quantum QM9700 switch, 920-9B210-00FN-0M0
Storage Fabric		
16	Fabric switches	NVIDIA Quantum QM9700 switch, 920-9B210-00FN-0M0
PDUs		
192	Rack PDUs	Raritan PX3-5878I2R-P1Q2R1A15D5
12	Rack PDUs	Raritan PX3-5747V-V2

Associated cables and transceivers are listed in Table 8f. All networking components are multi-mode fiber.

Table 8. Estimate of cables required for a 4 SU, 127-node DGX SuperPOD

Count	Component	Connection	Recommended Model
In-Band Ethernet Cables			
68	Ethernet 800Gb/s (2x400Gb/sTwin-port OSFP, DR8 multimode, parallel, 8-channel transceiver	Leaf and spine transceivers	980-9I510-F4NS00
2	DR1 Splitter cable 1x 400Gb/s to 4x 100Gb/s	Spine to SN2201 Leaf Leaf to NFS	Off-the-shelf, POT EFALU-PA2S1Q-005M or similar
4	100gb/s single mode 1-lane (DR1), QSFP28 optical transceiver	Spine to SN2201 on Leaf Transceivers	980-9I042-00C000
4	2x400GbE Twin-port OSFP 100-meter Single mode Ethernet transceiver	Spine to SN2201 on Spine NFS Connectivity to 5600	980-9I30H-F4NM00
254	DGX System 400G QSFP112 Multimode Transceivers	QSFP112 transceivers on DGX Systems,	980-9I693-00NS00
134	MMF MPO12 APC to 2xMPO12 APC 10m	DGX systems, Management nodes to leaf	980-9I570-00N030
14	Ethernet (ETH) 400Gb/s, Single-port, OSFP, multimode parallel transceiver	OSFP transceivers on SLURM and Management Nodes	980-9I51S-F4NS00
Varies	Varies	Customer NFS Storage	Varies
8	NVIDIA passive Copper cable, IB twin port NDR, up to 800Gb/s, OSFP, 1.5m	Leaf – Spine layer	980-9IA0Q-00N01A
4	Cat5e for UFM to Inband	UFM to Inband	Cat5e
OOB Ethernet Cables			
381	1 Gbps	DGX systems	Cat5e
64	1 Gbps	InfiniBand Switches	Cat5e
11	1 Gbps	Management/UFM nodes	Cat5e
6	1 Gbps	In-band Ethernet switches	Cat5e
2	1 Gbps	UFM Back-to-Back	Cat5e
204	1 Gbps	PDUs	Cat5e
34	100gb/s single mode 1-lane (DR1), QSFP28 optical transceiver	Spine to out of band SN2201	980-9I042-00C000
10	2x400GbE Twin-port OSFP 100-meter Single mode Ethernet transceiver	Spine to SN2201 on Spine	980-9I30H-F4NM00

20	DR1 Splitter cable 1x 400Gb/s to 4x 100Gb/s	Spine to SN2201 Leaf	Off-the-shelf, POT EFALU-PA2S1Q-005M or similar
Varies	1 Gbps	Storage	Cat5e
Compute InfiniBand Cabling			
2044	NDR Fiber Cables ¹ , 400 Gbps	DGX systems to leaf, leaf to spine, UFM to leaf ports	980-9I570-00N030
1536	Switch 2x400G OSFP Finned-top Multimode Transceivers	Leaf and spine transceivers	980-9I510-00NS00
508	System 2x400G OSFP Flat-top Multimode Transceivers	Transceivers in the DGX B200 systems	980-9I51A-00NS00
4	UFM System 400G OSFP Multimode Transceivers	UFM to leaf connections	980-9I51S-00NS00
InfiniBand Storage Cables ^{1,2}			
498	NDR Fiber Cables, 400 Gbps	DGX systems to leaf, leaf to spine, UFM to leaf connections	980-9I570-00N030
48	NDR AOC Cables, 2x 200 Gbps QSFP56-QSFP56	Storage	980-9I117-00H030
4	UFM System 400G OSFP Multimode Transceivers	UFM to leaf connections	980-9I51S-00NS00
369	Switch 2x400G OSFP Finned-top Multimode Transceivers	Leaf and spine transceivers	980-9I510-00NS00
254	DGX System 400G QSFP112 Multimode Transceivers	QSFP112 transceivers	980-9I693-00NS00
4	HDR 400 Gbps to 2x200 Gbps AOC Cables	Slurm management	980-9I117-00H030
Varies	Storage Cables, 400 Gbps to 2x200 Gbps AOC Cables	Varies	980-9I117-00H030
Ethernet Storage Cables ^{1,2}			
514	MMF MPO12 APC to 2xMPO12 APC 10m	DGX systems to leaf, leaf to spine, to SLURM nodes	980-9I570-00N030
386	2x400GbE Twin-port OSFP 50-meter Multimode Ethernet transceiver	Leaf and spine transceivers	980-9I510-F4NS00
8	400Gbs Single-port OSFP, 400Gbs Multimode SR4 50m	OSFP transceivers on SLURM Management Nodes	980-9I51S-00NS00
254	400GbE Single-port, QSFP112 50-meter Multimode Ethernet transceiver	QSFP112 transceivers on DGX Systems,	980-9I693-F4NS00
Varies	100gb/s single mode 1-lane (DR1), QSFP28 optical transceiver	Leaf transceivers for Storage	980-9I042-00C000

Varies	800gb/s to 4x 100Gb/s splitter cable	Leaf to Storage Cables	Varies
1. Part number will depend on exact cable lengths needed based on data center requirements. 2. Count and cable type required depend on specific storage selected.			

Notice

This document is provided for information purposes only and shall not be regarded as a warranty of a certain functionality, condition, or quality of a product. NVIDIA Corporation ("NVIDIA") makes no representations or warranties, expressed or implied, as to the accuracy or completeness of the information contained in this document and assumes no responsibility for any errors contained herein. NVIDIA shall have no liability for the consequences or use of such information or for any infringement of patents or other rights of third parties that may result from its use. This document is not a commitment to develop, release, or deliver any Material (defined below), code, or functionality.

NVIDIA reserves the right to make corrections, modifications, enhancements, improvements, and any other changes to this document, at any time without notice.

Customer should obtain the latest relevant information before placing orders and should verify that such information is current and complete.

NVIDIA products are sold subject to the NVIDIA standard terms and conditions of sale supplied at the time of order acknowledgement, unless otherwise agreed in an individual sales agreement signed by authorized representatives of NVIDIA and customer ("Terms of Sale"). NVIDIA hereby expressly objects to applying any customer general terms and conditions with regards to the purchase of the NVIDIA product referenced in this document. No contractual obligations are formed either directly or indirectly by this document.

NVIDIA products are not designed, authorized, or warranted to be suitable for use in medical, military, aircraft, space, or life support equipment, nor in applications where failure or malfunction of the NVIDIA product can reasonably be expected to result in personal injury, death, or property or environmental damage. NVIDIA accepts no liability for inclusion and/or use of NVIDIA products in such equipment or applications and therefore such inclusion and/or use is at customer's own risk.

NVIDIA makes no representation or warranty that products based on this document will be suitable for any specified use. Testing of all parameters of each product is not necessarily performed by NVIDIA. It is customer's sole responsibility to evaluate and determine the applicability of any information contained in this document, ensure that the product is suitable and fit for the application planned by customer, and perform the necessary testing for the application to avoid a default of the application or the product. Weaknesses in customer's product designs may affect the quality and reliability of the NVIDIA product and may result in additional or different conditions and/or requirements beyond those contained in this document. NVIDIA accepts no liability related to any default, damage, costs, or problem that may be based on or attributable to: (i) the use of the NVIDIA product in any manner that is contrary to this document or (ii) customer product designs.

No license, either expressed or implied, is granted under any NVIDIA patent right, copyright, or other NVIDIA intellectual property right under this document. Information published by NVIDIA regarding third-party products or services does not constitute a license from NVIDIA to use such products or services or a warranty or endorsement thereof. Use of such information may require a license from a third party under the patents or other intellectual property rights of the third party, or a license from NVIDIA under the patents or other intellectual property rights of NVIDIA.

Reproduction of information in this document is permissible only if approved in advance by NVIDIA in writing, reproduced without alteration and in full compliance with all applicable export laws and regulations, and accompanied by all associated conditions, limitations, and notices.

THIS DOCUMENT AND ALL NVIDIA DESIGN SPECIFICATIONS, REFERENCE BOARDS, FILES, DRAWINGS, DIAGNOSTICS, LISTS, AND OTHER DOCUMENTS (TOGETHER AND SEPARATELY, "MATERIALS") ARE BEING PROVIDED "AS IS." NVIDIA MAKES NO WARRANTIES, EXPRESSED, IMPLIED, STATUTORY, OR OTHERWISE WITH RESPECT TO THE MATERIALS, AND EXPRESSLY DISCLAIMS ALL IMPLIED WARRANTIES OF NONINFRINGEMENT, MERCHANTABILITY, AND FITNESS FOR A PARTICULAR PURPOSE. TO THE EXTENT NOT PROHIBITED BY LAW, IN NO EVENT WILL NVIDIA BE LIABLE FOR ANY DAMAGES, INCLUDING WITHOUT LIMITATION ANY DIRECT, INDIRECT, SPECIAL, INCIDENTAL, PUNITIVE, OR CONSEQUENTIAL DAMAGES, HOWEVER CAUSED AND REGARDLESS OF THE THEORY OF LIABILITY, ARISING OUT OF ANY USE OF THIS DOCUMENT, EVEN IF NVIDIA HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES. Notwithstanding any damages that customer might incur for any reason whatsoever, NVIDIA's aggregate and cumulative liability towards customer for the products described herein shall be limited in accordance with the Terms of Sale for the product.

Trademarks

NVIDIA, the NVIDIA logo, NVIDIA DGX, NVIDIA DGX SuperPOD, NVIDIA Base Command, are trademarks and/or registered trademarks of NVIDIA Corporation in the U.S. and other countries. Other company and product names may be trademarks of the respective companies with which they are associated.

Copyright

© 2024 NVIDIA Corporation. All rights reserved.