



NVIDIA DGX SuperPOD: Next Generation Scalable Infrastructure for AI Leadership

Reference Architecture

Featuring NVIDIA DGX H100 Systems

Document History

RA-11333-001

Version	Date	Authors	Description of Change
01	2022-11-16	Craig Tierney, Chris Kawalek, Premal Savla, Robert Sohigian, Shawn Kaiser, and Tony Paikeday	Initial Release
02	2022-12-7	Shawn Kaiser and Craig Tierney	Minor updates
03	2022-12-16	Craig Tierney, Robert Sohigian, and Shawn Kaiser	Updates to the BOM
04	2023-1-13	Craig Tierney, Robert Sohigian, and Shawn Kaiser	Update DGX H100 networking for QSFP and OSFP
05	2023-03-22	Robert Sohigian	New template and minor updates

Abstract

The NVIDIA DGX SuperPOD™ with NVIDIA DGX™ H100 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. The 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.

To achieve the most scalability, DGX SuperPOD is powered by several key NVIDIA technologies, including:

- ▶ NVIDIA DGX H100 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—networking technologies that connect GPUs at the NVLink layer to provide unprecedented performance for most demanding communication patterns.

The DGX SuperPOD architecture is managed by NVIDIA solutions including NVIDIA Base Command™, NVIDIA AI Enterprise, CUDA, and 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 the DGX SuperPOD. The system is built upon building blocks of scalable units (SU), each containing 32 DGX H100 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 the DGX SuperPOD.....	1
NVIDIA DGX H100 System	1
NVIDIA InfiniBand Technology	2
Runtime and System Management	2
Components.....	3
Design Requirements	4
System Design	4
InfiniBand Fabrics.....	4
Compute Fabric.....	4
Storage Fabric	4
Ethernet Fabrics.....	5
In-Band Management Network.....	5
Out-of-Band Management Network.....	5
Storage Requirements.....	5
High-Performance Storage	5
User Storage.....	5
DGX SuperPOD Architecture	6
Network Fabrics	8
Compute—InfiniBand Fabric	9
Storage—InfiniBand Fabric	10
In-Band Management Network.....	11
Out-of-Band Management Network.....	12
Storage Architecture.....	13
DGX SuperPOD Software	16
NVIDIA Base Command	16
NVIDIA NGC.....	17
NVIDIA AI Enterprise.....	17
Summary	18
Appendix A. Major Components	iv

Key Components of the 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. Using SUs, system deployment times are reduced from months to weeks. Leveraging the DGX SuperPOD designs reduces time-to-solution and time-to-market of next generation models and applications.

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

NVIDIA DGX H100 System

The NVIDIA DGX H100 system (Figure 1) is an AI powerhouse that enables enterprises to expand the frontiers of business innovation and optimization. The DGX H100 system, which is the fourth-generation NVIDIA DGX system, delivers AI excellence in an eight GPU configuration. The NVIDIA Hopper GPU architecture provides latest technologies such as the transformer engines and fourth-generation NVLink technology that brings months of computational effort down to days and hours, on some of the largest AI/ML workloads.

Figure 1. DGX H100 system



Some of the key highlights of the DGX H100 system over the DGX A100 system include:

- ▶ Up to 9X more performance with 32 petaFLOPS at FP8 precision.
- ▶ Dual 56-core 4th Gen Intel® Xeon® capable processors with PCIe 5.0 support and DDR5 memory.
- ▶ 2X faster networking and storage @ 400 Gbps InfiniBand/Ethernet with NVIDIA ConnectX®-7 smart network interface cards (SmartNICs).
- ▶ 1.5X higher bandwidth per GPU @ 900 GBps with fourth generation of NVIDIA NVLink.
- ▶ 640 GB of aggregated HBM3 memory with 24 TB/s of aggregate memory bandwidth, 1.5X higher than DGX A100 system.

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. Driven by the InfiniBand Trade Association (IBTA), it continues to evolve and lead data center network performance.

The latest generation InfiniBand, NDR, has a peak speed of 400 Gbps per direction. It is backwards compatible with the previous generations of InfiniBand specifications. InfiniBand is more than just peak performance. 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 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. For on-premises solutions, the customer has the option to operate the system with a secure, cloud-native interface through NVIDIA NGC™.

Components

The components of the DGX SuperPOD are described in Table 1.

Table 1. Four SU, 127-node DGX SuperPOD components

Component	Technology	Description
Compute nodes	127 × NVIDIA DGX H100 system with eight 80 GB H100 GPUs	Fourth generation of the world's premier purpose-built AI systems featuring NVIDIA H100 Tensor Core GPUs, 4 th generation NVIDIA NVLink® and 3 rd generation NVIDIA NVSwitch™ technologies.
Compute fabric	NVIDIA Quantum QM9700 NDR 400 Gbps InfiniBand	Rail-optimized, full fat-tree network with eight NDR400 connections per system
Storage fabric	NVIDIA Quantum QM9700 NDR 400 Gb/s InfiniBand	The fabric is optimized to match peak performance of the configured storage array
Compute/storage fabric management	NVIDIA Unified Fabric Manager, 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 SN4600C switch	64 port 100 Gbps Ethernet switch providing high port density with high performance
Out-of-band (OOB) management network	NVIDIA SN2201 switch	48 port 1 Gbps Ethernet switch leveraging copper ports to minimize complexity
DGX SuperPOD software stack	NVIDIA Base Command Manager	Cluster management for DGX SuperPOD
	NVIDIA AI Enterprise	Best-in-class development tools and frameworks for the AI practitioner and reliable management and orchestration for IT professionals
	Magnum IO	The NVIDIA 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
User environment	Slurm	Slurm is a classic workload manager used to manage complex workloads in a multi-node, batch-style, compute environment

Design Requirements

The 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

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

- ▶ A modular architecture based on SUs of 32 DGX H100 systems each.
- ▶ A fully tested system scales to four SUs, but larger deployments can be built based on customer requirements.
- ▶ Rack design can support one, two, or four DGX H100 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 Support NVES).

InfiniBand Fabrics

Compute Fabric

- ▶ The InfiniBand compute fabric is rail-optimized to the top layer of the fabric.
- ▶ The InfiniBand fabric is a balanced, full-fat tree.
- ▶ Managed NDR switches are used throughout the design to provide better management of the fabric.
- ▶ The fabric is designed to support the latest SHaRPv3 features.

Storage Fabric

The storage fabric provides high bandwidth to shared storage. It also has these 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 H100 system.
- ▶ Storage is provided over InfiniBand and leverages RDMA to provide maximum performance and minimize CPU overhead.
- ▶ It is flexible and can scaled to meet specific capacity and bandwidth requirements.
- ▶ User-accessible management nodes provide access to shared storage.

Ethernet Fabrics

Multiple Ethernet fabrics are used to support management communications, Ethernet-based storage targets, Internet access, and other traditional TCP/IP based services.

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 servers operate at 100 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 system users.

Storage Requirements

The DGX SuperPOD compute architecture must be paired with a high-performance, balanced, storage system to maximize overall system performance. The 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

HPS must provide:

- ▶ High-performance, resilient, POSIX-style file system optimized for multi-threaded read and write operations across multiple nodes.
- ▶ Native InfiniBand support.
- ▶ Local system RAM for transparent caching of data.
- ▶ Leverage local disk transparently for caching of larger datasets.

User Storage

User storage must:

- ▶ Be 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 to provide a secondary path to storage so, that in the event of a failure of the storage fabric or HPS, nodes can still be accessed and managed by administrators in parallel.

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 40 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

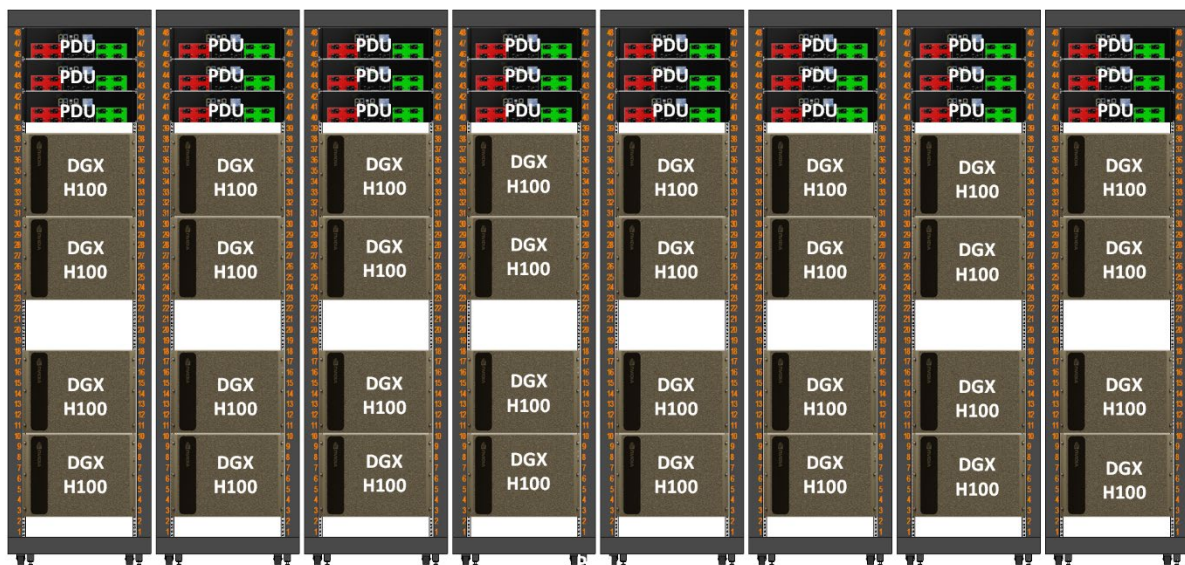
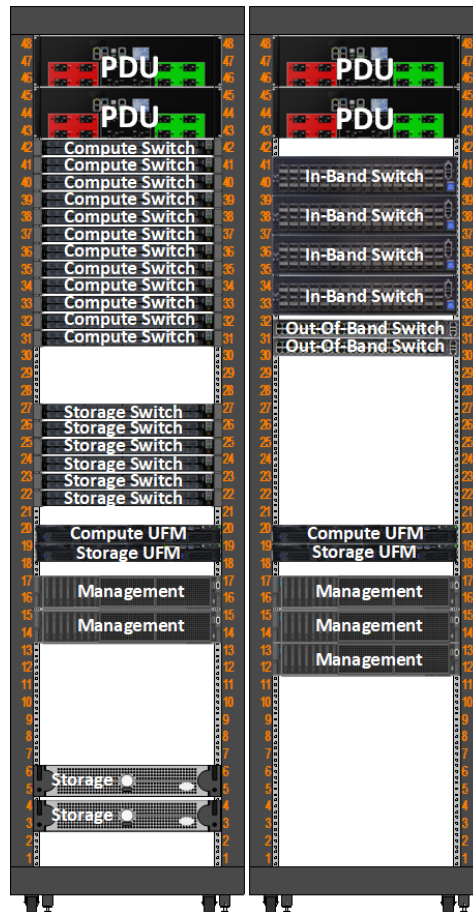


Figure 3 shows a typical management rack configuration with InfiniBand and Ethernet switches, management servers, storage arrays, and UFM appliances.

Figure 3. Typical management rack

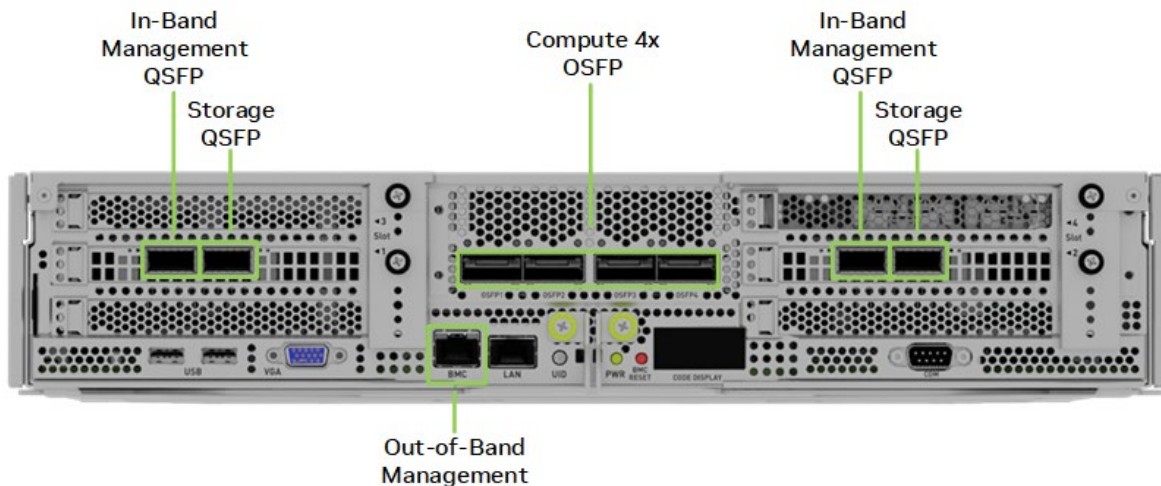


Network Fabrics

Several networks are deployed on the DGX SuperPOD. The compute fabric is used for inter-node communication through the applications. A separate storage fabric is used to isolate storage traffic. There are two Ethernet fabrics for in-band and OOB management. Requirements for each section are detailed below. In addition, designs for the network are provided after the requirements.

Figure 4 shows the different ports on the back of the DGX H100 CPU tray and the connectivity provided. The InfiniBand compute fabric ports in the middle use a two-port transceiver to access all eight GPUs. Each pair of in-band Ethernet management and InfiniBand storage ports provide parallel pathways into the DGX H100 system for increased performance. The OOB port is used for BMC access. In addition, there is an additional LAN port next to the BMC but is not used in the DGX SuperPOD.

Figure 4. DGX H100 network ports



Compute—InfiniBand 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 H100 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 InfiniBand fabric for full 127 node DGX SuperPOD

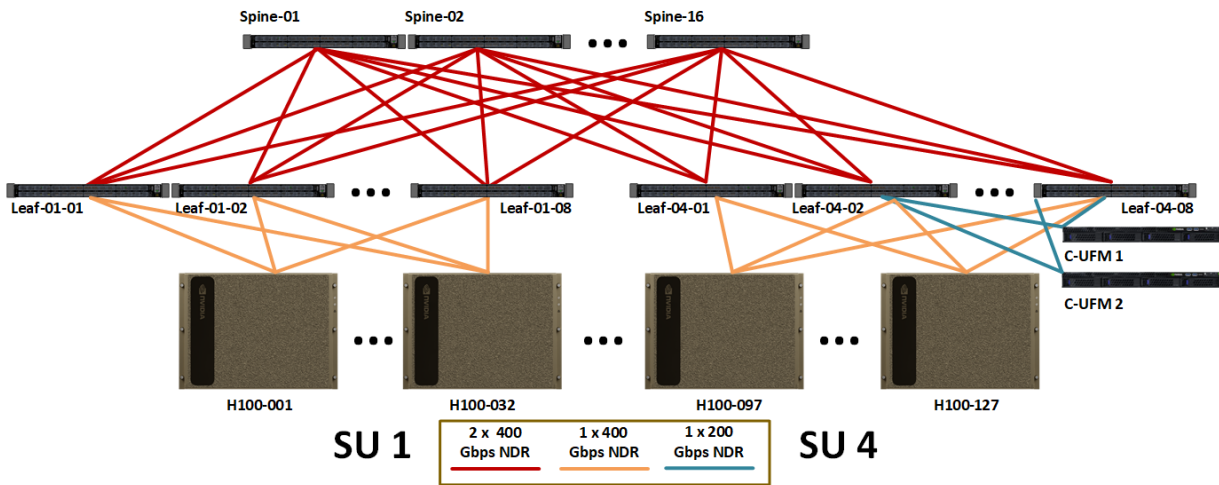


Table 2 shows the number of cables and switches required for the compute fabric for different SU sizes.

Table 2. Compute fabric component count

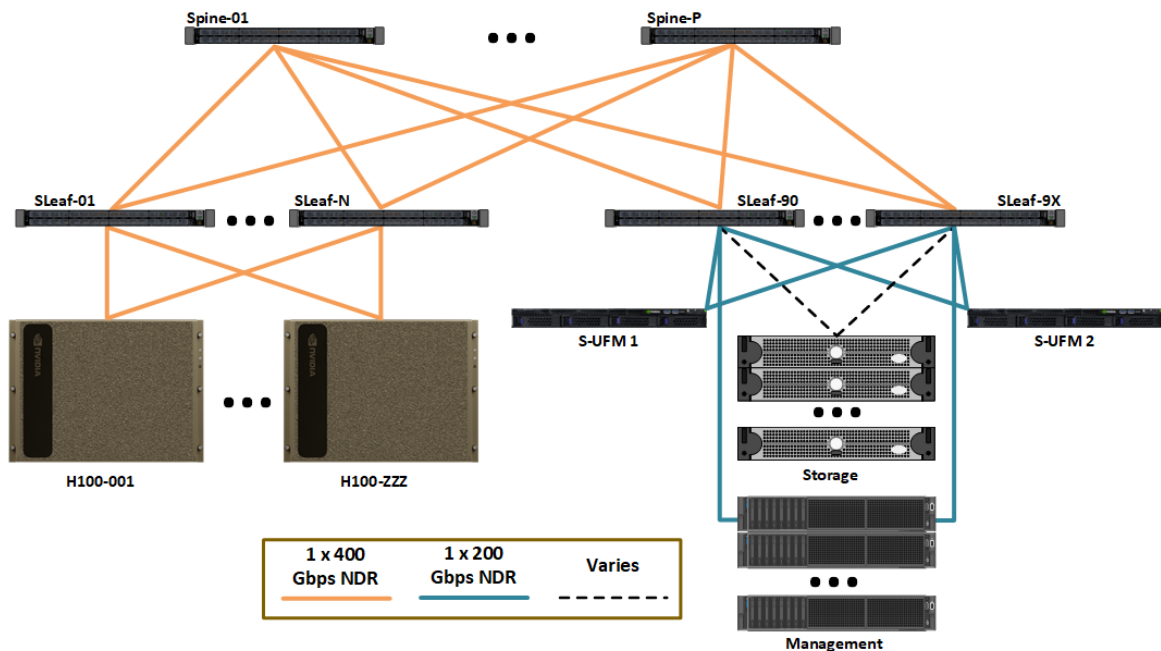
SU Count	Cluster Size # Nodes	Cluster Size # GPUs	Leaf Switch Count	Spine Switch Count	Compute + UFM Node Cable Count	Spine-Leaf Cable Count
1	31 ¹	248	8	4	252	256
2	63	504	16	8	508	512
3	95	760	24	16	764	768
4	127	1016	32	16	1020	1024
1. This is a 32 node per SU design, however a DGX Node must be removed to accommodate for UFM connectivity.						

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.

Storage—InfiniBand 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. InfiniBand storage fabric logical design



The storage fabric uses [MQM9700-NS2F](#) switches (Figure 7). The storage devices are connected at a 1:1 port to uplink ratio. The DGX H100 system connections are slightly oversubscribed with a ratio near 4:3 with adjustments as needed to allow for more storage flexibility regarding cost and performance.

Figure 7. MQM9700-NS2F switch



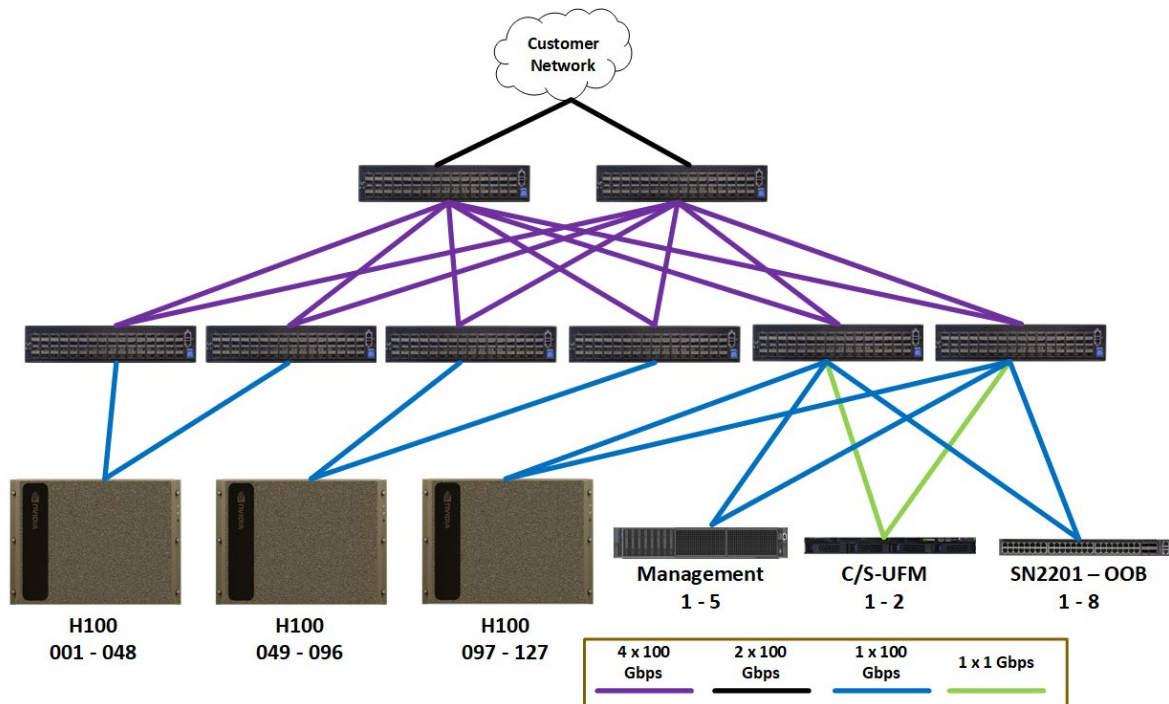
In-Band Management Network

The in-band management network provides several key functions:

- ▶ Connects all the services that manage the cluster.
- ▶ Enables access to the home filesystem and storage pool.
- ▶ Provides connectivity for the in-cluster services such as Base Command Manager, Slurm and to other services outside of the cluster such as the NGC registry, code repositories, and data sources.

Figure 8 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.

Figure 8. In-band Ethernet network



The in-band management network uses [SN4600C](#) switches (Figure 9).

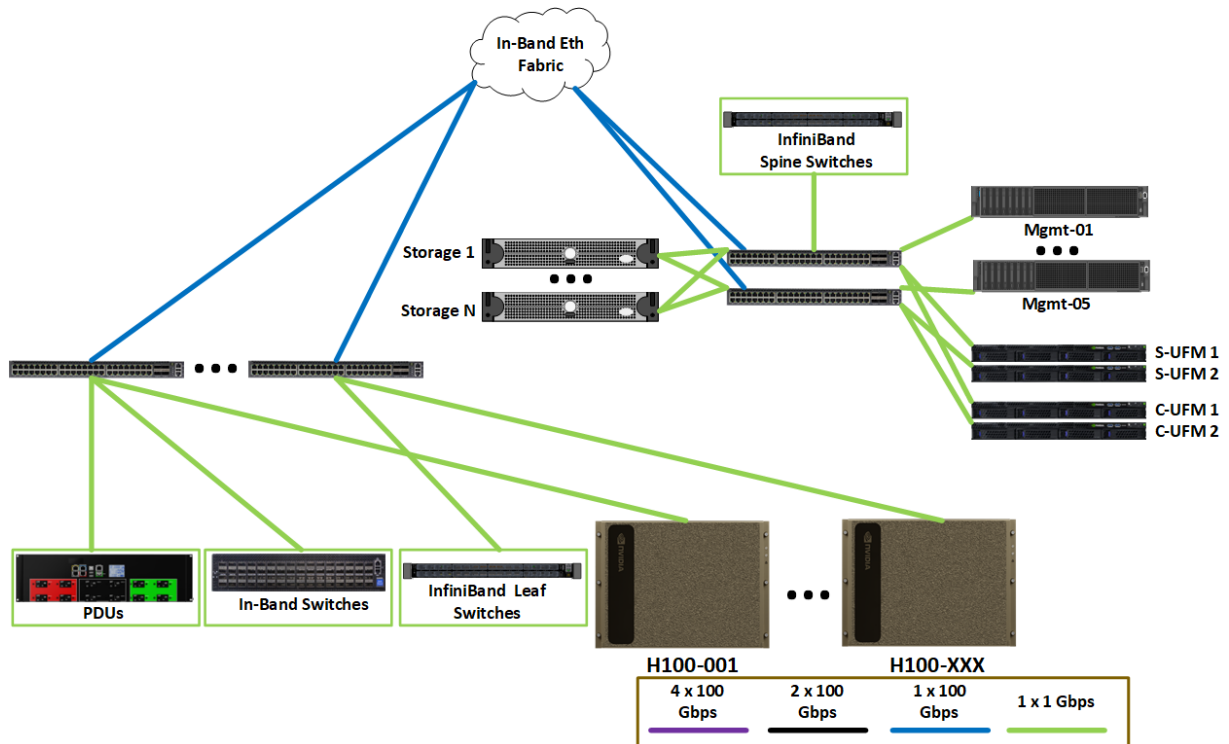
Figure 9. SN4600C switch



Out-of-Band Management Network

Figure 10 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 separate onto their own fabric since there is no use-case where users need access to these ports and are secured using logical network separation.

Figure 10. Logical OOB management network layout



The OOB management network uses SN2201 switches (Figure 11).

Figure 11. 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. The performance of the DGX H100 system is up to nine times faster than its predecessor. To achieve this in practice, 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 in size 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 H100 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 3 and Table 4 are provided to help determine the I/O levels required for different types of models.


Table 3. Storage performance requirements

Performance Level	Work Description	Dataset Size
Good	Natural Language Processing (NLP)	Datasets generally fit within local cache
Better	Image processing with compressed images (ex: ImageNet)	Many to most datasets can fit within the local system's cache
Best	Training with 1080p, 4K, or uncompressed images, offline inference, ETL,	Datasets are too large to fit into cache, massive first epoch I/O requirements, workflows that only read the dataset once

Table 4. Guidelines for storage performance

Performance Characteristic	Good (GBps)	Better (GBps)	Best (GBps)
Single-node read	4	8	40
Single-node write	2	4	20
Single SU aggregate system read	15	40	125
Single SU aggregate system write	7	20	62
4 SU aggregate system read	60	160	500
4 SU aggregate system write	30	80	250

Even for the best category above, it is desirable that the single node read performance is closer to the maximum network performance of 80 GBps.

	<p>Note: As datasets get larger, they may no longer fit in cache on the local system. Pairing large datasets that do not fit in cache with very fast GPUs can create a situation where it is difficult to achieve maximum training performance. NVIDIA GPUDirect Storage® (GDS) provides a way to read data from the remote filesystem or local NVMe directly into GPU memory providing higher sustained I/O performance with lower latency. Using the storage fabric on the DGX SuperPOD, a GDS-enabled application should be able to read data at over 40 GBps directly into the GPUs.</p>
-------------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

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 1080p 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.

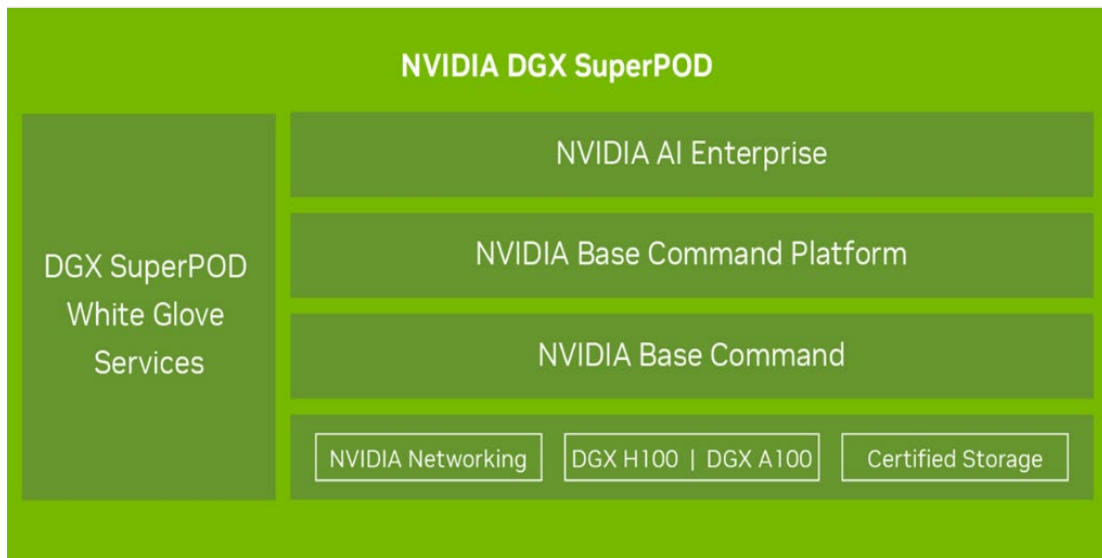
While NLP cases often do not require as much read performance for training, peak performance for reads and writes are needed for creating and reading checkpoint files. This is a synchronous operation and training stops during this phase. If you are looking for best end-to-end training performance, do not ignore I/O operations for checkpoints.

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 12) 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 12. DGX SuperPOD high-level architecture



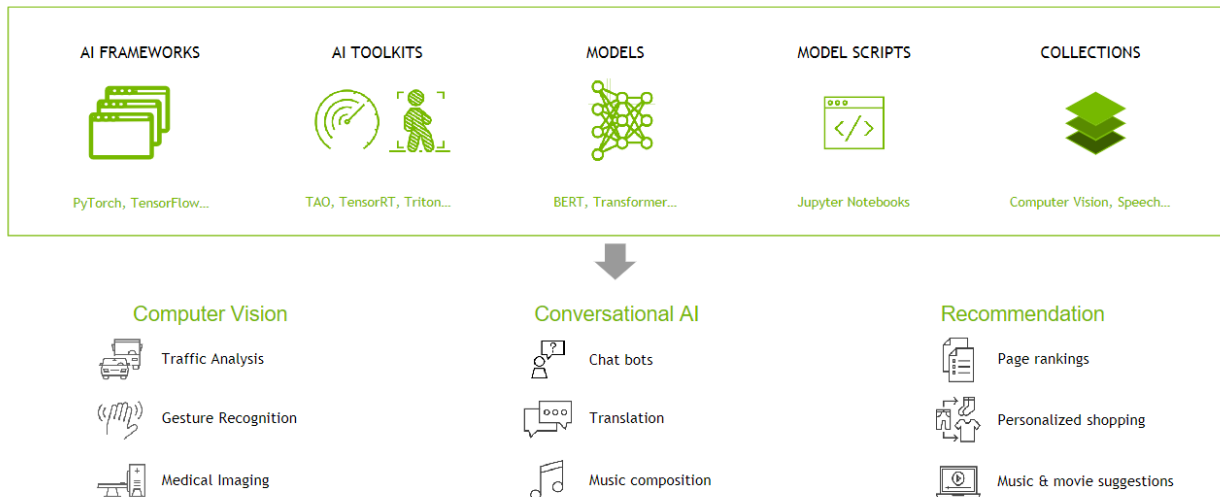
NVIDIA Base Command

[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 (Figure 13) provides software to meet the needs of data scientists, developers, and researchers with various levels of AI expertise.

Figure 13. NGC catalog overview



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 a suite of AI and data analytics software optimized for the development and deployment of AI. NVIDIA AI Enterprise includes proven, open-sourced containers and frameworks such as NVIDIA RAPIDS, NVIDIA TAO Toolkit, NVIDIA TensorRT™ and NVIDIA Triton Inference Server, which are certified and supported to run on DGX SuperPOD. NVIDIA AI Enterprise is included with DGX SuperPOD and is used in combination with NVIDIA Base Command and NVIDIA NGC.

Summary

DGX SuperPOD with NVIDIA DGX H100 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, and support the same accelerated software that you leverage across DGX SuperPOD, on laptops, or other NVIDIA GPU-based systems. 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 5. These are representative of the configuration and must be finalized based on actual design.

Table 5. Major components of the 4 SU, 127-node DGX SuperPOD

Count	Component	Recommended Model
Racks		
38	Rack (Legrand)	NVIDPD13
Nodes		
127	GPU nodes	NVIDIA DGX H100 systems
4	UFM appliance	NVIDIA Unified Fabric Manager Appliance 3.1
5	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 HDR200 VPI Ports, TPM 2.0
Ethernet Network		
8	In-band management	NVIDIA SN4600C switch with Cumulus Linux
8	OOB management	NVIDIA SN2201 switch with Cumulus Linux
Compute InfiniBand Fabric		
48	Fabric switches	NVIDIA Quantum QM9700 switch, 920-9B210-00FN-OM0
Storage InfiniBand Fabric		
16	Fabric switches	NVIDIA Quantum QM9700 switch, 920-9B210-00FN-OM0
PDUs		
96	Rack PDUs	Raritan PX3-5878I2R-P1Q2R1A15D5
12	Rack PDUs	Raritan PX3-5747V-V2

Associated cables and transceivers are listed in Table 6.

Table 6. Estimate of cables required for a 4 SU, 127-node DGX SuperPOD

Count	Component	Connection	Recommended Model
In-Band Ethernet Cables			
254	100 Gbps	DGX H100 system	Varies
32	100 Gbps QSFP to QSFP AOC	Management nodes	Varies
6	100 Gbps	ISL Cables	Varies
Varies	Ethernet (perf varies)	Storage	Varies
Varies	Varies	Core DC	Varies
OOB Ethernet Cables			
127	1 Gbps	DGX H100 systems	Cat5e
64	1 Gbps	IB Switches	Cat5e
11	1 Gbps	Mgmt/UFM nodes	Cat5e
8	1 Gbps	In-band Eth Switches	Cat5e
Varies	1 Gbps	Storage	Cat5e
108	1 Gbps	PDUs	Cat5e
16	100 Gbps	Two uplinks per OOB	Varies
Compute InfiniBand Cabling			
2040	NDR Cables1, 400 Gbps	DGX H100 systems to leaf, leaf to spine	NVIDIA 980-9I57X-00N010
2	NDR Cables, 200 Gbps	UFM to leaf ports	NVIDIA 980-9I1111-00H010
1536	Switch OSFP Transceivers	Leaf and spine transceivers	NVIDIA 980-9IA20-00NS00
508	System OSFP Transceivers	Transceivers in the DGX H100 Systems	NVIDIA 980-9I89P-00N000
4	UFM System Transceivers	UFM to leaf connections	NVIDIA 980-9I89R-00NS00
Storage InfiniBand Cables1			
494	NDR Cables, 400 Gbps	DGX H100 systems to leaf, leaf to spine	NVIDIA 980-9I57X-00N010
482	NDR Cables, 200 Gbps	Storage	NVIDIA 980-9I1111-00H010
6	NDR Cables, 200 Gbps	Management/UFM nodes	NVIDIA 980-9I1111-00H010
395	Switch Transceivers	Leaf and spine transceivers	NVIDIA 980-9IA20-00NS00
254	DGX System Transceivers	QSFP112 Transceivers	NVIDIA 980-9I693-00NS00
8	Mgmt/UFM Transceivers	NDR200 Transceivers	NVIDIA 980-9I693-00NS00
Varies	Storage Cables, NDR200	Varies	Varies
Varies	Storage Transceivers	Varies	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. 3. All networking components are multi-mode fiber.			

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 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

© 2023 NVIDIA Corporation. All rights reserved.