# DDN A³I® SOLUTIONS FOR NVIDIA® CLOUD PARTNERS

**Fully-validated and optimized AI High-Performance Storage (HPS) solutions for cloud service partners featuring NVIDIA HGX™ H100 and H200 8-GPU based platforms.**

## Executive Summary

DDN A³I Solutions are proven at-scale to deliver optimal data performance for Artificial Intelligence (AI), Data Analytics, and High-Performance Computing (HPC) applications running on GPUs in NVIDIA HGX ™ H100 and H200 platforms. DDN data platforms are used in some of the largest NVIDIA DGX SuperPOD™ deployments, including NVIDIA Eos, Selene, and Cambridge-1.

This document describes fully validated reference architectures for NVIDIA Cloud Partners (NCPs). The solutions integrate DDN AI400X2 Turbo appliances and DDN Insight software with NVIDIA HGX platforms for High-Performance Storage (HPS).

# 1. DDN A³I End-To-End Enablement for NVIDIA Cloud Partners

DDN A³I solutions (Accelerated, Any-Scale AI) are architected to achieve the most from at-scale AI, Data Analytics and HPC applications running on HGX systems. They provide predictable performance, capacity, and capability through a tight integration between DDN and NVIDIA systems. Every layer of hardware and software engaged in delivering and storing data is optimized for fast, responsive, and reliable access.

DDN A³I solutions are designed, developed, and optimized in close collaboration with NVIDIA. The deep integration of DDN AI appliances with HGX systems ensures a reliable experience. DDN A³I solutions are configured for flexible deployment in a wide range of environments and scale seamlessly in capacity and capability to match evolving workload needs.

DDN brings the same advanced technologies used to power the world's largest supercomputers in a fully integrated package for HGX systems that is easy to deploy and manage. DDN A³I solutions are proven to maximize benefits for at-scale AI, Analytics, and HPC workloads on HGX systems.

This section describes the advanced features of DDN A³I Solutions for NCPs.

## 1.1. DDN A³I Shared Parallel Architecture

The DDN A³I shared parallel architecture and client protocol ensures high levels of performance, scalability, security, and reliability for HGX systems. Multiple parallel data paths extend from the drives all the way to containerized applications running on the GPUs in the HGX system. With DDN's true end-to-end parallelism, data is delivered with high-throughput, low-latency, and massive concurrency in transactions. This ensures applications achieve the most from HGX systems with all GPU cycles put to productive use. Optimized parallel data-delivery directly translates to increased application performance and faster completion times. The DDN A³I shared parallel architecture also contains redundancy and automatic failover capability to ensure high reliability, resiliency, and data availability in case a network connection or server becomes unavailable.

## 1.2. DDN A³I Streamlined Deep Learning Workflows

DDN A³I solutions enable and accelerate end-to-end data pipelines for deep learning (DL) workflows of all scale running on HGX systems. The DDN shared parallel architecture enables concurrent and continuous execution of all phases of DL workflows across multiple HGX systems. This eliminates the management overhead and risks of moving data between storage locations. At the application level, data is accessed through a standard and highly operable file interface, for a familiar and intuitive user experience.

Significant acceleration can be achieved by executing an application across multiple HGX systems simultaneously and engaging parallel training efforts of candidate neural network variants. These advanced optimizations maximize the potential of DL frameworks. DDN works closely with NVIDIA and its customers to develop solutions and technologies that allow widely used DL frameworks to run reliably on HGX systems.

## 1.3. DDN A³I Multi-rail Networking

DDN A³I solutions integrate a wide range of networking technologies and topologies to ensure streamlined deployment and optimal performance for AI infrastructure. The latest generation NVIDIA Quantum InfiniBand and Spectrum-X Ethernet technology provide both high-bandwidth and low-latency data transfers between applications, compute servers and storage appliances.

DDN A³I Multi-rail enables the grouping of multiple network interfaces on an HGX system to achieve faster aggregate data transfer capabilities without any switch configuration such as channel groups or bonding. The feature balances traffic dynamically across all the interfaces, and actively monitors link health for rapid failure detection and automatic recovery. DDN A³I Multi-rail makes designing, deploying, and managing high-performance networks very simple, and is proven to deliver complete connectivity for at-scale infrastructure for NCPs.

## 1.4. DDN A³I Advanced Optimizations for HGX System Architecture

The DDN A³I client's NUMA-aware capabilities enable strong optimization for HGX systems. It automatically pins threads to ensure I/O activity across the HGX system is optimally localized, reducing latencies, and increasing the utilization efficiency of the whole environment. Further enhancements reduce overhead when reclaiming memory pages from page cache to accelerate buffered operations to storage. The DDN A³I client software for HGX systems has been validated at-scale.

## 1.5. DDN A³I Hot Nodes

DDN Hot Nodes is a powerful software enhancement that enables the use of the NVME devices in an HGX system as a local cache for read-only operations. This method significantly improves the performance of applications if a data set is accessed multiple times during a particular workflow.

This is typical with DL training, where the same input data set or portions of the same input data set are accessed repeatedly over multiple training iterations. Traditionally, the application on the HGX system reads the input data set from shared storage directly, thereby continuously consuming shared storage resources. With Hot Nodes, as the input data is read during the first training iteration, the DDN software automatically writes a copy of the data on the local NVME devices. During subsequent reads, data is delivered to the application from the local cache rather than the shared storage. This entire process is managed by the DDN client software running on the HGX system. Data access is seamless, and the cache is fully transparent to users and applications. The use of the local cache eliminates network traffic and reduces the load on the shared storage system. This allows other critical DL training operations like checkpointing to complete faster by engaging the full capabilities of the shared storage system.

DDN Hot Nodes includes extensive data management tools and performance monitoring facilities. These tools enable user-driven local cache management and make integration simple with task schedulers. For example, training input data can be loaded to the local cache on an HGX system as a pre-flight task before the AI training application is engaged. As well, the metrics expose insightful information about cache utilization and performance, enabling system administrators to further optimize their data loading and maximize application and infrastructure efficiency gains.

## 1.6. DDN A$^3$I Container Client

Containers encapsulate applications and their dependencies to provide simple, reliable, and consistent execution. DDN enables a direct high-performance connection between the application containers on the HGX system and the DDN parallel filesystem. This brings significant application performance benefits by enabling low-latency, high-throughput parallel data access directly from a container. Additionally, the limitations of sharing a single host-level connection to storage between multiple containers disappear. The DDN in-container filesystem mounting capability is added at runtime through a universal wrapper that does not require any modification to the application or container.

Containerized versions of popular DL frameworks specially optimized for HGX systems are available from NVIDIA. They provide a solid foundation that enables data scientists to rapidly develop and deploy applications on HGX systems. In some cases, open-source versions of the containers are available, further enabling access and integration for developers. The DDN A$^3$I container client provides high-performance parallelized data access directly from containerized applications on an HGX system. This provides containerized DL frameworks with the most efficient dataset access possible, eliminating all latencies introduced by other layers of the computing stack.

## 1.7. DDN A$^3$I S3 Data Services

DDN S3 Data Services provide hybrid file and object data access to the shared namespace. The multi-protocol access to the unified namespace provides tremendous workflow flexibility and simple end-to-end integration. Data can be captured directly to storage through the S3 interface and accessed immediately by containerized applications on an HGX system through a file interface. The shared namespace can also be presented through an S3 interface, for easy collaboration with multisite and multicloud deployments. The DDN S3 Data Services architecture delivers robust performance, scalability, security, and reliability features.

## 1.8. DDN A³I Multitenancy

DDN A³I makes it very simple to operate a secure multitenant environment at-scale through its native client and comprehensive digital security framework. DDN A³I multitenancy makes it simple to share HGX systems across a large pool of users and still maintain secure data segregation. Multitenancy provides quick, seamless, dynamic HGX system resource provisioning for users. It eliminates resource silos, complex software release management, and unnecessary data movement between data storage locations.

DDN A³I brings a very powerful multitenancy capability to HGX systems and makes it very simple for customers to deliver a secure, shared innovation space, for at-scale data-intensive applications. Containers bring security challenges and are vulnerable to unauthorized privilege escalation and data access. The DDN A³I digital security framework provides extensive controls, including a global root_squash to prevent unauthorized data access or modification from a malicious user, and even if a node or container is compromised.

DDN A³I also provides secure client authentication to prevent rogue client impersonation through Shared-Secret-Key (SSK) or Kerberos authentication. Tenant directories can be allocated a fixed capacity of space and inodes they can consume through the Project quotas features. If required, DDN A³I provides client-side encryption through the fscrypt interface and supports multiple encryptions keys and encryption per-file or directory level, providing strong privacy of the data.

## 2. DDN A³I Solutions with NVIDIA HGX Systems

The DDN A³I scalable architecture integrates HGX systems with DDN AI shared parallel file storage appliances and delivers fully optimized end-to-end AI, Analytics, and HPC workflow acceleration on GPUs. DDN A³I solutions greatly simplify the deployment of HGX systems, while also delivering performance and efficiency for maximum GPU saturation, and high levels of scalability.

This section describes the components integrated in DDN A³I Solutions for NCPs.

## 2.1. DDN AI400X2T Appliance

The AI400X2T appliance (Figure 1) is a fully integrated and optimized shared data platform with predictable capacity, capability, and performance. Every AI400X2T appliance delivers over 110 GB/s and 3M IOPS directly to HGX systems. Shared performance scales linearly as additional AI400X2T appliances are integrated to the storage system. The all-NVMe configuration provides optimal performance for a wide variety of workload and data types and ensures that NCPs achieve the most from at-scale GPU applications, while maintaining a single, shared, centralized data platform.



*Figure 1. DDN AI400X2T all-NVME storage appliance.*

The AI400X2T appliance integrates the DDN A$^3$I shared parallel architecture and includes a wide range of capabilities described in section 1, including automated data management, digital security, and data protection, as well as extensive monitoring. The AI400X2T appliances enable NCPs to go beyond basic infrastructure and implement complete data governance pipelines at-scale.

For NCP deployments, DDN uses two distinct appliance configurations to deploy a shared data platform. The AI400X2T-OSS appliance provides data storage through four OSS and eight OST appliances and is available in 120, 250 and 500 TB useable capacity options. The AI400X2T-MDS appliance provides metadata storage through four MDSs and four MDT appliances, each appliance provides 9.2 billion inodes. Both appliance configurations must be used jointly to provide a filesystem and must be connected to HGX systems through RDMA over Converged Ethernet (RoCE) using ConnectX®-7 HCAs. Each appliance provides eight interfaces, two per OSS/MDS, to connect to the storage fabric.

## 2.2. DDN Insight

DDN Insight (Figure 2) is a centralized management and monitoring software suite for AI400X2T appliances. It provides extensive performance and health monitoring of all DDN storage systems in a cloud infrastructure from a single web-based user interface. DDN Insight greatly simplifies IT operations and enables automated and proactive storage platform management guided by analytics and intelligent software.
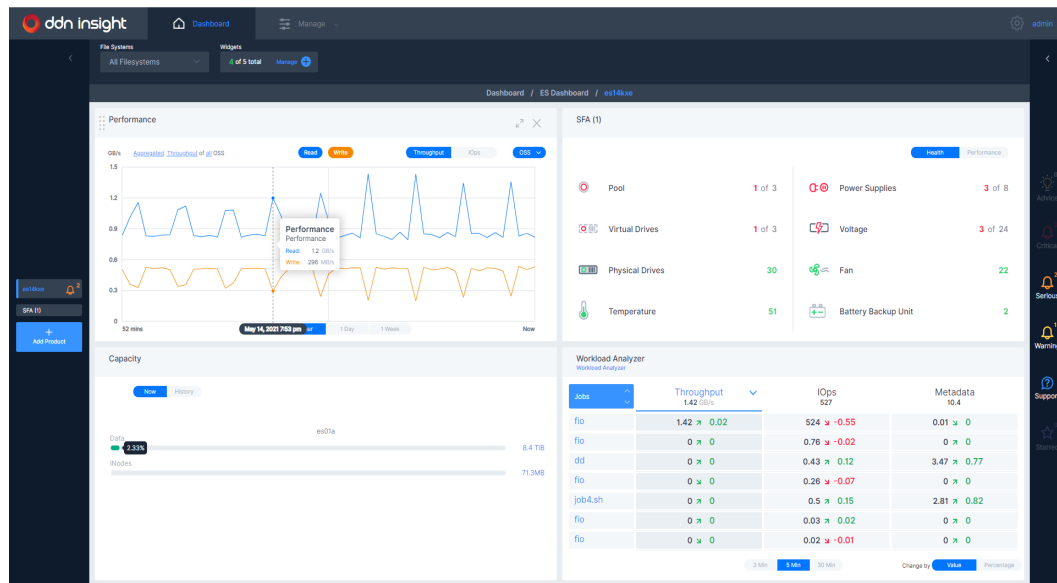


*Figure 2. DDN Insight Workload Analyzer Tool.*

Performance monitoring is an important aspect of operating a cloud infrastructure efficiently. Provided the several variables that affect data I/O performance, the identification of bottlenecks and degradation is crucial while production workloads are engaged. DDN Insight provides deep real-time analysis across the entire cloud infrastructure, tracking I/O transactions from applications running on HGX systems all the way through individual drives in the AI400X2T appliances. The embedded analytics engine makes it simple for cloud operators to visualize I/O performance across their entire infrastructure through intuitive user interfaces. These include extensive logging, trending, and comparison tools, for analyzing I/O performance of specific applications and users over time. The open backend database also makes it simple to extend the benefits of DDN Insight and integrate other AI infrastructure components within the engine, or export data to third party monitoring systems.

DDN Insight is provided as a turnkey server appliance from DDN.

## 3. DDN A³I Reference Architectures for NCP Deployments

The DDN AI400X2T appliance is a turnkey shared storage system. The AI400X2T appliances deliver optimal GPU performance for every workload and data type in a dense, power efficient 2RU chassis. The AI400X2T appliance simplifies the design, deployment, and management of a cloud deployment and provides predictable performance, capacity, and scaling. The appliance is designed for seamless integration with HGX systems and enables customers to move rapidly from test to production. In addition, DDN provides complete expert design, deployment, and support services globally. The DDN field engineering organization has already deployed dozens of solutions for customers based on the A³I reference architectures.

DDN proposes the following recommended architectures for NCP deployments.

## 3.1. Recommended Storage Sizing

The storage performance required to maximize training performance can vary depending on the type of model and dataset. The guidelines in Table 1 and Table 2 provide guidance and targets to determine the I/O levels required for different types of models.

| Performance Level | Workload Examples | Dataset Size |
|---|---|---|
| Good | Classical Natural Language Processing (NLP) workload, such as BERT | Datasets generally fit within local cache |
| Better | Training compressed images, compressed audio and text data, such as Large Language Model (LLM) Training | Many to most datasets can fit within the local system's cache, or datasets are distributed across different nodes |
| Best | Extract, transform, and load (ETL), training with large video and image files such as AV replay, generative networks such as stable diffusion, 3D images such as Medical uNet, genomics workload and protein prediction such as AlphaFold | Datasets are too large to fit into cache, massive first epoch I/O requirements, workflows that only read the dataset once, or in-place data processing |

*Table 1. Guidelines for storage performance requirements.*

| Performance Characteristics | Good (GBps) | Better (GBps) | Best (GBps) |
|---|---|---|---|
| 4 SU (1k SU group) aggregate system read | 60 | 160 | 500 |
| 4 SU (1k SU group) aggregate system write | 30 | 80 | 250 |
| 8 SU (2k SU group) aggregate system read | 120 | 320 | 1000 |
| 8 SU (2k SU group) aggregate system write | 60 | 160 | 500 |
| 32 SU (8k SU group) aggregate system read | 480 | 1280 | 4000 |
| 32 SU (8k SU group) aggregate system write | 240 | 640 | 2000 |
| 64 SU (16k SU group) aggregate system read | 960 | 2560 | 8000 |
| 64 SU (16k SU group) aggregate system write | 480 | 1280 | 4000 |

*Table 2. Guidelines for storage performance scaling.*

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. When looking for best end-to-end training performance, I/O operations for checkpoints are important to the HPS design. Modern LLM workload has a significant requirement in the write performance to not consume too much time in writing checkpoints.

As a reference, Table 3 has hypothetical calculations for required write rate for a LLM training. It has these constants:

- Number of bytes per parameter: 14.
- Total write time percentage of total training time: 1%.
- Seconds per checkpointing interval: 3,600s

| Number of Parameters (in Billions) | Size (TB) | Tensor Parallel Domain Size | Pipeline Parallel Domain Size | Number of Files | Total Write Rate (GB/s) | Per Node Write Rate (GB/s) | Per GPU Write Rate (GB/s) |
|---|---|---|---|---|---|---|---|
| 175 | 2.4 | 8 | 16 | 128 | 68.06 | 4.25 | 0.53 |
| 530 | 7.2 | 8 | 32 | 280 | 206.11 | 5.89 | 0.74 |
| 1,000 | 13.7 | 8 | 64 | 512 | 388.89 | 6.08 | 0.76 |

*Table 3. Estimates of LLM checkpoint size.*

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

As general guidance, DDN recommends to base sizing on the "Better" guideline, so that that the shared storage be sized to ensure close to 1 GB/s per second of read and 500 MB/s write throughput for every NVIDIA GPU (Table 4). This ensures the minimum performance required to operate the GPU infrastructure, based on NVIDIA reference design for NCP deployments.

| | | Scalable Units (SUs) | | | |
|---|---|---|---|---|---|
| | | 4 | 8 | 32 | 64 |
| **Compute components** | NVIDIA HGX Systems | 127 | 255 | 1023 | 2047 |
| | NVIDIA GPUs | 1016 | 2040 | 8184 | 16376 |
| **DDN Storage components** | DDN Metadata appliances | 2 | 4 | 12 | 23 |
| | DDN Data appliances | 10 | 20 | 68 | 137 |
| **DDN Storage specification** | Aggregate read throughput | 1.1 TB/s | 2.2 TB/s | 7.5 TB/s | 15 TB/s |
| | Aggregate write throughput | 700 GB/s | 1.4 TB/s | 4.8 TB/s | 9.6 TB/s |
| | Per GPU read throughput | 1.1 GB/s | 1.1 GB/s | 0.91 GB/s | 0.92 GB/s |
| | Per GPU write throughput | 0.7 GB/s | 0.7 GB/s | 0.58 GB/s | 0.59 GB/s |
| | Number of namespaces | 1 | 1 | 2 | 3 |
| | Minimum aggregated useable capacity | 1.2 PB | 2.4 PB | 8.2 PB | 16.4 PB |
| | Aggregate useable inodes | 18.4 billion | 37 billion | 110 billion | 212 billion |
| | Physical, rack units | 24 | 48 | 160 | 320 |
| | Power, nominal | 25 KW | 50 KW | 173 kW | 346 kW |
| | Cooling, nominal | 87 kBTU/hr | 174 kBTU/hr | 581 kBTU/hr | 1,163 kBTU/hr |

Table 4. DDN recommended storage sizing for NCP deployments.

For configurations not listed in Table 4, a single namespace can be used up to 25 SUs and scales up to 64 AI400X2T-OSS (Data) paired with 12 AI400X2T-MDS (Metadata) appliances. Each namespace can have up to 128 tenants configured. The minimum number of appliances for a single namespace is one AI400X2T-OSS and one AI400X2T-MDS, when using more than 6 AI400X2T-OSS appliances, DDN requires to pair a AI400X2T-MDS in a 6:1 ratio (one AI400X2T-MDS for up to six AI400X2T-OSS). This ratio is done with the assumption of 80KB and 300KB per inode for respectively the smallest (120TB) and the biggest (500TB) capacity available.

The sizing shown in Table 4 is based on "Peak Single-GPU" read and write performance guidance from Table 2.

DDN recommends this level of performance to ensure that the storage can satisfy a wide range of data types, data access patterns as required by at-scale AI training applications. This recommendation is based on DDN's experience as a widely deployed data storage platform for NVIDA DGX SuperPOD.

For convenience, DDN details configurations based on "aggregate system" read and write performance guidance from Table 2 at the end of this document in Appendix 6.

## 3.2. Recommended Storage Networking

The NCP reference design includes several networks, one of which is used for storage traffic. The storage network provides connectivity between the AI400X2T appliances, the compute nodes, and management nodes. This network is designed to meet the high-throughput, low-latency, and scalability requirements of NCP deployments.

The NVIDIA SN5600 switch (Figure 3) is recommended for NCP storage connectivity. It provides 128 ports of 400GbE over 64 OSFP ports in a 2RU form factor. The cables listed in Section 3.2.3 are validated with the SN5600 switch.



*Figure 3. NVIDIA SN5600 switch.*

An overview of the NCP storage network architecture is shown in Figure 4.

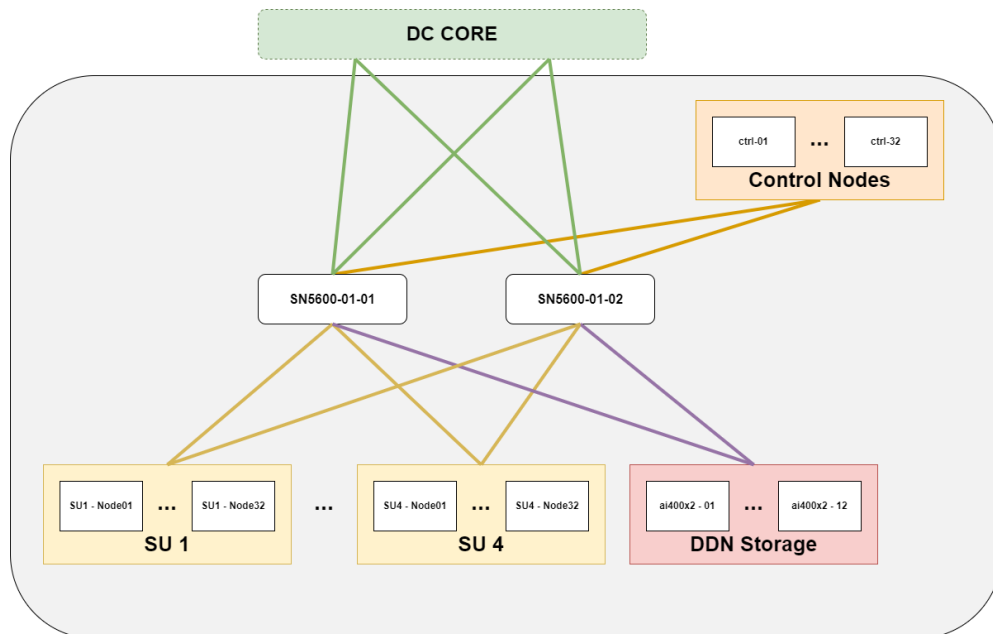DDN requires the storage network to run in RDMA over Converged Ethernet (RoCE) mode.



*Figure 4. Overview of the NCP storage network reference architecture.*

## 3.2.1. HGX System Network Connectivity

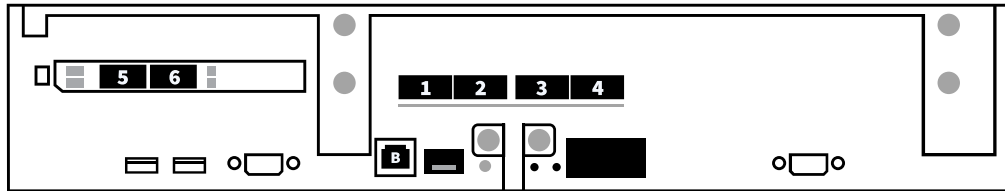Figures 5 and 6 detail recommended port and network connections for each HGX system.



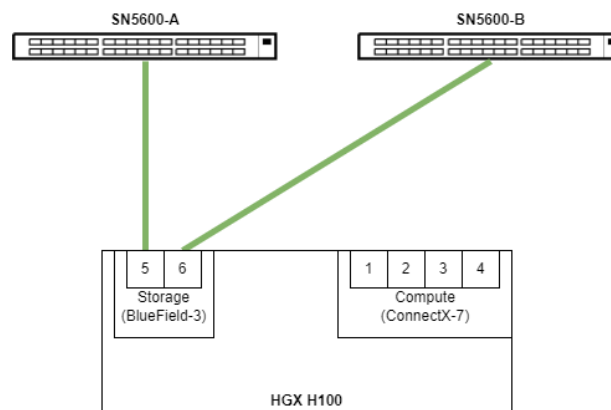*Figure 5. Recommended HGX system network port connections.*



*Figure 6. Recommended HGX system network connection diagram.*

DDN requires that OSFP ports 1 to 4 (ConnectX-7 HCAs) on the HGX systems be connected to the compute network. Ports 5 and 6 (Bluefield-3 DPU HCA) should be connected to the in-band management network, which also serves as the storage network using RDMA Over Converged Ethernet (RoCE). DDN does not recommend the use of bonding, link redundancy is handled by DDN A³I Multi-rail. As well, the management BMC ("B") port should be connected to the out-of-band management network.

### 3.2.2. AI400X2T Appliance Network Connectivity

DDN recommends using two splitter cables to connect the AI400X2 appliance to the storage network, see Section 3.2.3 for splitter cables reference. One splitter cable must be connected to ports 2, 6, 10 and 14 of the appliances and going to a switch "A." The other splitter cable must be connected to ports 4, 8, 12 and 16 of the appliances and going to another switch "B" to allow link redundancy, see Figure 7 and Figure 8.

As well, the management ("M") and BMC ("B") ports for both controllers should be connected to the out-of-band management network. Note that each AI400X2T appliance requires one inter-controller network port connection ("I") using short Ethernet cable supplied.
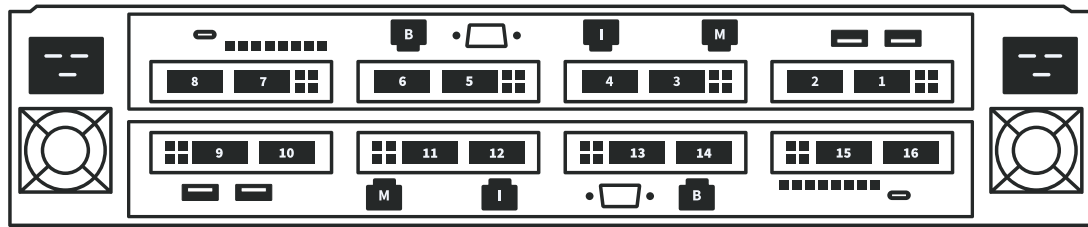


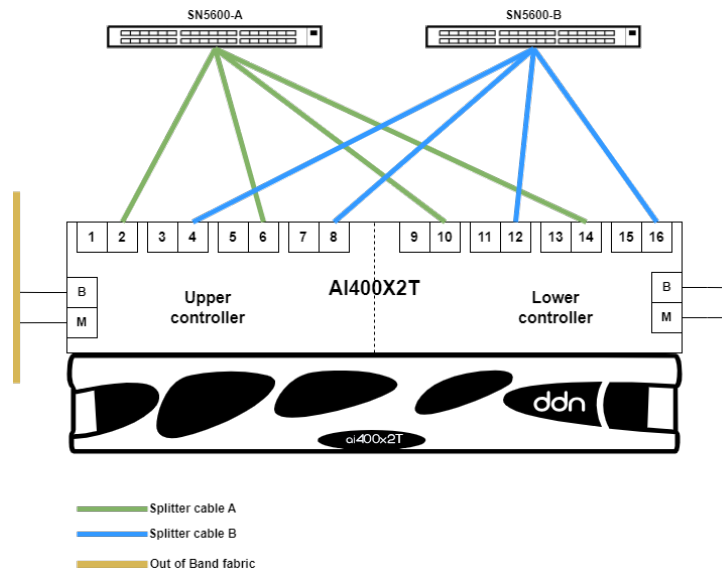*Figure 7. Recommended AI400X2T appliance network port connections.*



*Figure 8. Recommended AI400X2T network connection diagram with splitter cables.*

### 3.2.3. AI400X2T Appliance Cabling with SN5600 Switches

The AI400X2T appliance connects to the storage network with eight 200GbE interfaces. Particular attention must be given to the cabling selection to ensure compatibility between different Ethernet connectivity and data rates. DDN and NVIDIA have validated the cables in Table 5 to connect AI400X2T appliances with SN5600 switches. The use of splitter cables ensures the most efficient use of switch ports. Two splitter cables are required per each AI400X2T appliance.

| Cable Type | Part Number | Description |
|---|---|---|
| Direct Attach Copper | MCP7Y40-Nxxx | NVIDIA 800Gb/s Twin-port OSFP to 4x200Gb/s QSFP112 DAC Splitter Cable<br><br>xxx indicates length in meters: 001, 01A, 002, 02A, 003 |
| Active Copper Cable | MCA7J75-Nxxx | NVIDIA 800Gb/s Twin-port OSFP to 4x200Gb/s QSFP112 ACC Splitter<br><br>xxx indicates the length in meters: 004, 005 |

Table 5. DAC and ACC splitter cables.

## 3.3. NCP Deployments with 127 HGX Systems

Figure 9 illustrates the DDN A[3]I reference architecture for NCP deployments with 127 HGX systems, ten DDN AI400X2T-OSS, and two DDN AI400X2T-MDS appliances and a DDN Insight server. Every HGX system connects to the storage network with two 200GbE links. Each AI400X2T appliance connects to the storage network with eight 200GbE links using the appropriate cable type. The DDN Insight server connects to the AI400X2T appliances over the 1GbE out-of-band management network. It does not require a connection to the storage network.
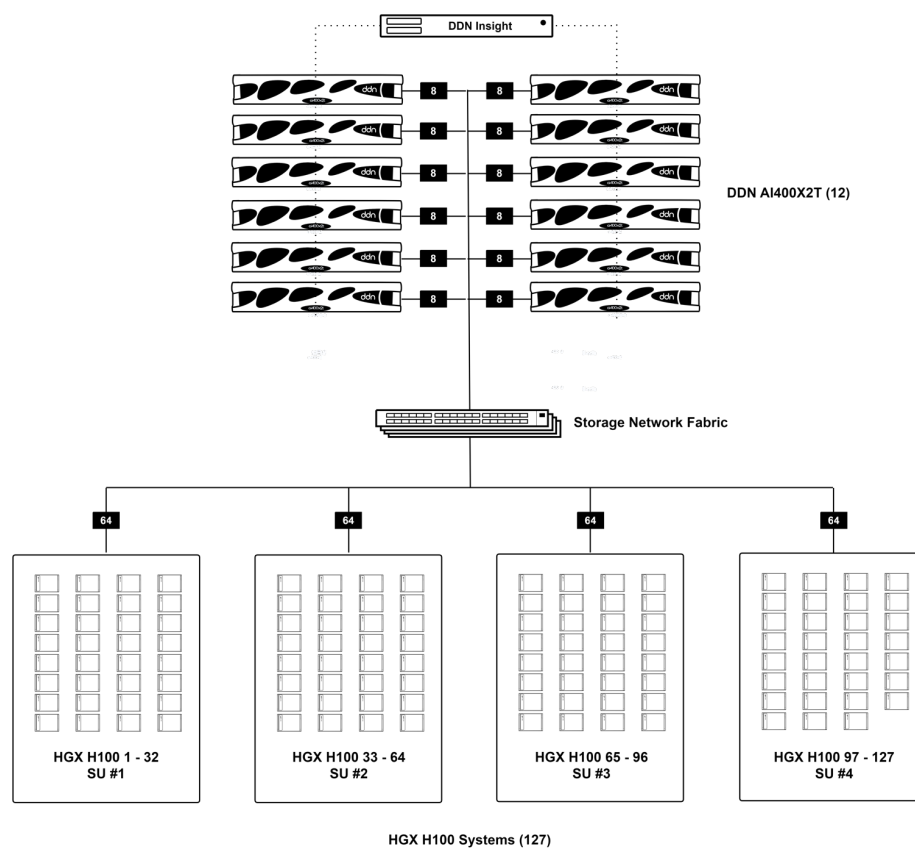


*Figure 9. DDN A[3]I reference architecture for NCP deployments with 127 HGX systems.*

| Description | Count |
|---|---|
| 200GbE links (Storage) | 96 |
| OSPF ports (Switch side) / Splitter cables (MCP7Y40-Nxxx) | 24 |
| 1GbE links (Management) | 48 |

*Table 6. Cable counts for NCP deployments with 127 HGX systems.*

## 3.4. NCP Deployments with 255 HGX Systems

Figure 10 illustrates the DDN A[3]I reference architecture for NCP deployments with 255 HGX systems, 20 DDN AI400X2T-OSS and four DDN AI400X2T-MDS appliances and a DDN Insight server. Every HGX system connects to the storage network with two 200GbE links. Each AI400X2T appliance connects to the storage network with eight 200GbE links using the appropriate cable type. The DDN Insight server connects to the AI400X2T appliances over the 1GbE out-of-band management network. It does not require a connection to the storage network.
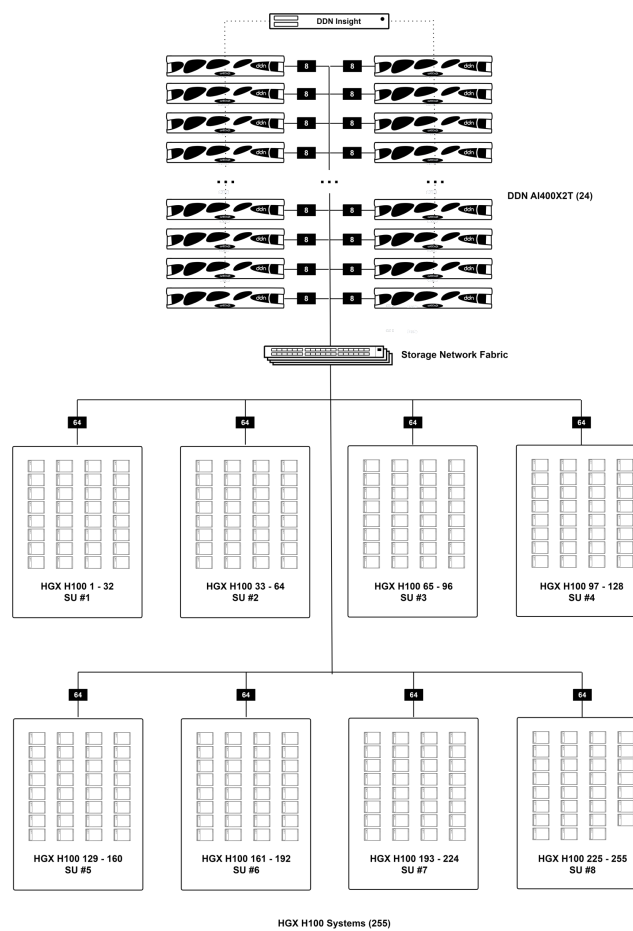


*Figure 10. DDN A[3]I reference architecture for NCP deployments with 255 HGX systems.*

| Description | Count |
|---|---|
| 200GbE links (Storage) | 192 |
| OSPF ports (Switch side) / Splitter cables (MCP7Y40-Nxxx) | 48 |
| 1GbE links (Management) | 96 |

*Table 7. Cable counts for NCP deployments with 255 systems.*

## 3.5. NCP Deployments with 1023 HGX Systems

Figure 11 illustrates the DDN A³I reference architecture for NCP deployments 1023 HGX systems, 68 DDN AI400X2T-OSS appliances and 12 DDN AI400X2T-MDS and a DDN Insight server. Every HGX system connects to the storage network with two 200GbE links. Each AI400X2T appliance connects to the storage network with eight 200GbE links using the appropriate cable type. The DDN Insight server connects to the AI400X2T appliances over the 1GbE out-of-band management network. It does not require a connection to the storage network.
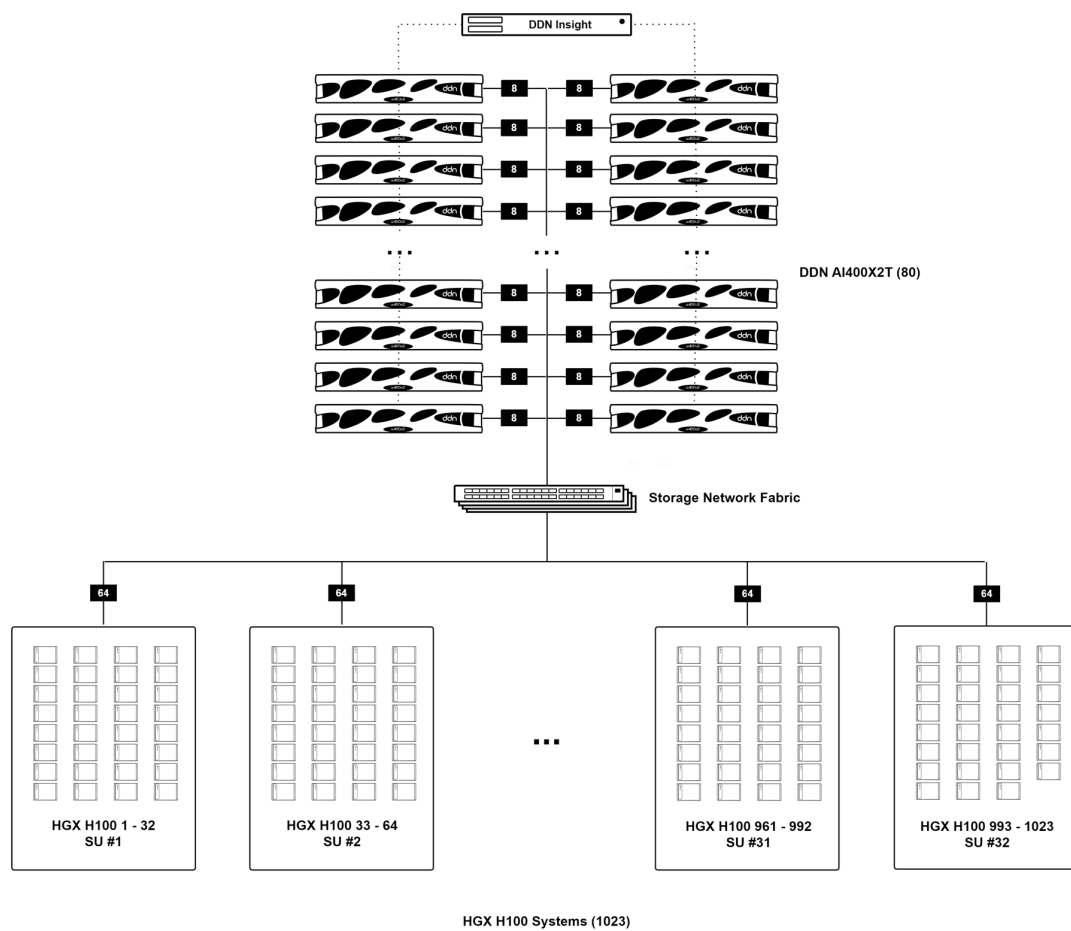


Figure 11. DDN A³I reference architecture for NCP deployments with 1023 HGX systems.

| Description | Count |
|---|---|
| 200GbE links (Storage) | 640 |
| OSPF ports (Switch side) / Splitter cables (MCP7Y40-Nxxx) | 160 |
| 1GbE links (Management) | 320 |

*Table 8. Cable counts for NCP deployments with 1023 HGX systems.*

## 3.6. NCP Deployments with 2047 HGX Systems

Figure 12 illustrates the DDN A[3]I reference architecture for NCP deployments with 2047 HGX systems, 137 DDN AI400X2T-OSS appliances and 23 DDN AI400X2T-MDS and a DDN Insight server. Every HGX system connects to the storage network with two 200GbE links. Each AI400X2T appliance connects to the storage network with eight 200GbE links using the appropriate cable type. The DDN Insight server connects to the AI400X2T appliances over the 1GbE out-of-band management network. It does not require a connection to the storage network.
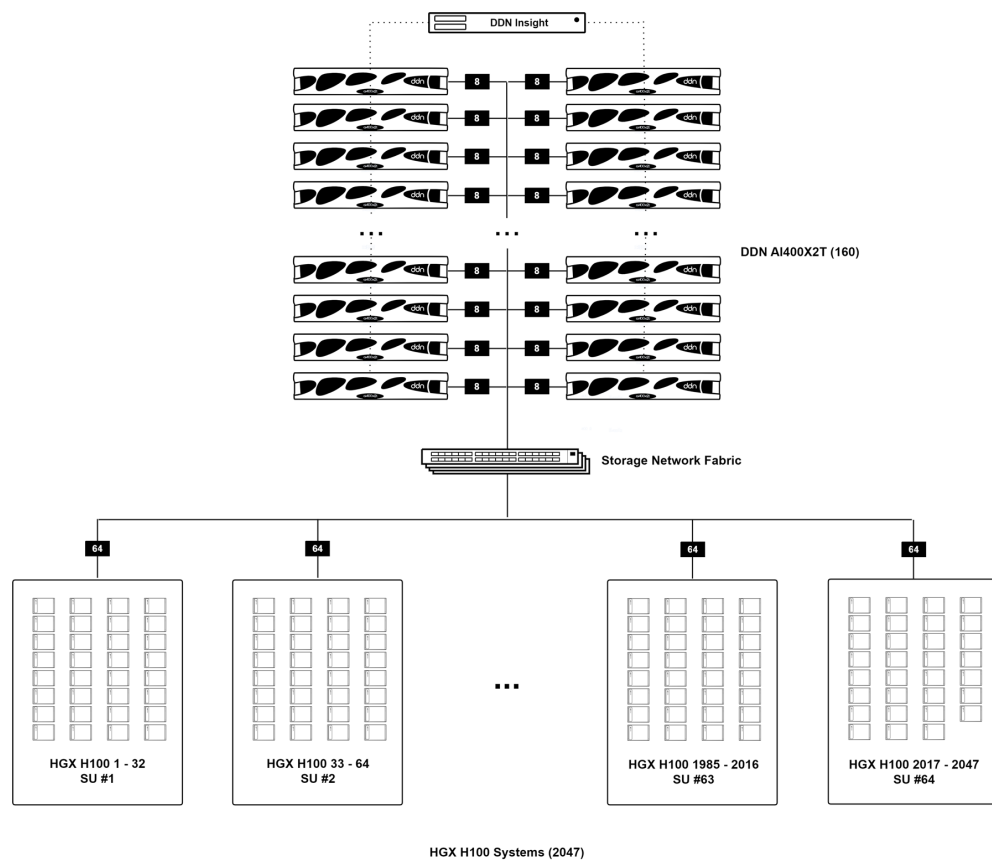


*Figure 12. DDN A[3]I reference architecture for NCP deployments with 2047 HGX systems.*

| Description | Count |
|---|---|
| 200GbE links (Storage) | 1280 |
| OSPF ports (Switch side) / Splitter cables (MCP7Y40-Nxxx) | 320 |
| 1GbE links (Management) | 640 |

*Table 9. Cable counts for NCP deployments with 2047 HGX systems.*

## 4. DDN A³I Solutions Performance Validation

DDN conducts extensive engineering integration, optimization, and validation efforts in close collaboration with NVIDIA and customers to ensure the best possible end-user experience using the reference designs in this document. The joint validation confirms functional integration, and optimal performance out-of-the-box for NCP deployments.

Performance testing on the DDN A³I architecture is conducted with industry standard synthetic throughput and IOPS applications, as well as widely used DL frameworks and data types. The results demonstrate that with the DDN A³I shared parallel architecture, GPU-accelerated applications can engage the full capabilities of the data infrastructure and the HGX systems. Performance is distributed evenly across all the HGX systems in the NCP deployment and scales linearly as more HGX systems are engaged.

This section details some of the results from recent at-scale testing integrating AI400X2T appliances with HGX systems in an NCP deployment.

## 4.1. HGX System Performance Validation

The tests described in this section were executed on one HGX H100 system equipped with eight H100 GPUs and one AI400X2T appliances running DDN EXAScaler 6.3.1.

For the storage network, the HGX H100 system is connected to a NVIDIA SN4700 switch with two 200 GbE links (see recommendation in section 3.2.1). The AI400X2T appliance is connected to the same network with eight 200 GbE links each (see recommendation in section 3.2.2).

This test environment allows us to demonstrate and project storage performance for a scalable unit configuration used as part of the NCP reference architecture.

This test demonstrates the peak performance of the DDN NCP reference architecture using the FIO open-source synthetic benchmark tool. The tool is set to simulate a general-purpose workload without any performance-enhancing optimizations. Separate tests were run to measure both 100% read and 100% write workload scenarios.

The FIO configuration parameters used for these tests were:

- blocksize = 1024k
- direct = 1
- iodepth= 8
- ioengine = libaio
- bw-threads = 128

Figure 13 demonstrates that the DDN solution can deliver close to 50 GB/s of read throughput and close to 50 GB/s of write throughput to a single HGX H100 system. In this test, data is accessed through a single posix mount provided by the DDN shared parallel filesystem client installed on the HGX H100 system. Storage is accessed through two 200GbE links on a Bluefield DPU.
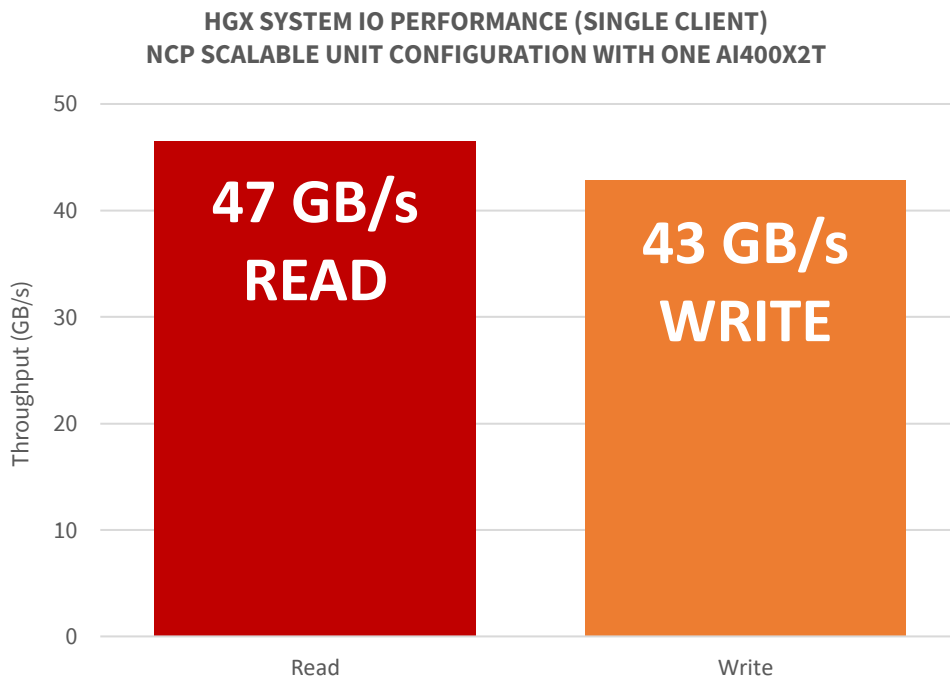
**HGX SYSTEM IO PERFORMANCE (SINGLE CLIENT)**
**NCP SCALABLE UNIT CONFIGURATION WITH ONE AI400X2T**



*Figure 13. FIO throughput using a single HGX H100 system.*

Figure 14 demonstrates the projection that the DDN software evenly distributes the full read and write performance of the AI400X2T appliances with all 127 HGX systems engaged simultaneously. The DDN solution utilizes the network links on every HGX system, ensuring optimal performance for a very wide range of data access patterns and data types.

This test demonstrates that the storage solution is provisioned to deliver at least 1.1 GB/s of read and 700 MB/s of write throughput for every H100 GPU in the SU.
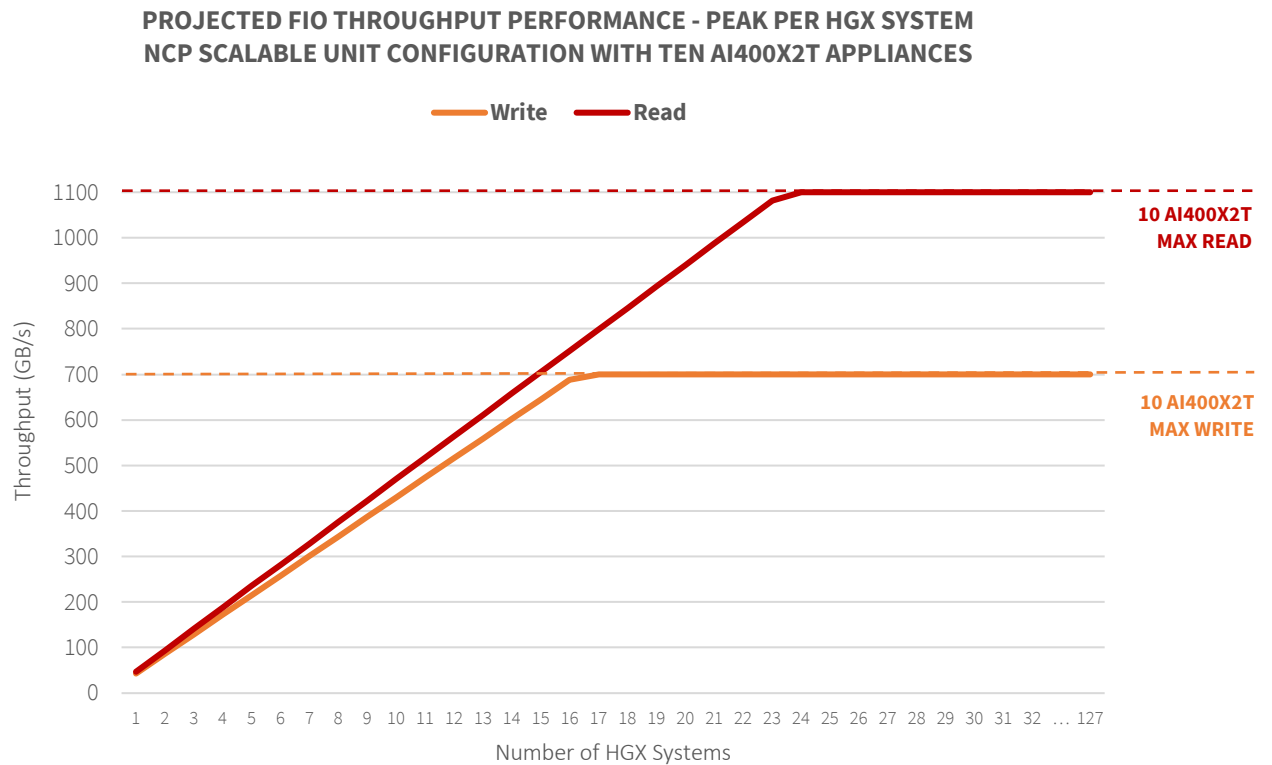
**PROJECTED FIO THROUGHPUT PERFORMANCE - PEAK PER HGX SYSTEM
NCP SCALABLE UNIT CONFIGURATION WITH TEN AI400X2T APPLIANCES**



*Figure 14. Projected FIO throughput performance with 32 HGX systems.*

## 5. Contact DDN to Unleash the Power of Your Cloud Deployment

DDN has long been a partner of choice for organizations pursuing at-scale data-driven projects. Beyond technology platforms with proven capabilities, DDN provides significant technical expertise through its global research and development and field technical organizations.

A worldwide team with hundreds of engineers and technical experts can be called upon to optimize every phase of a customer project: initial inception, solution architecture, systems deployment, customer support, and future scaling needs.

Strong customer focus coupled with technical excellence and deep field experience ensures that DDN delivers the best possible solution to any challenge. Taking a consultative approach, DDN experts will perform an in-depth evaluation of requirements and provide application-level optimization of data workflows for a project. They will then design and propose an optimized, highly reliable, and easy to use solution that best enables and accelerates the customer effort.

Drawing from the company's rich history in successfully deploying large-scale projects, DDN experts will create a structured program to define and execute a testing protocol that reflects the customer environment and meets project objectives. DDN has equipped its laboratories with leading GPU compute platforms to provide unique benchmarking and testing capabilities for AI and DL applications.

Contact DDN today and engage our team of experts to unleash the power of your NCP deployment.

## 6. Appendix

For convenience, DDN details configurations based on "aggregate system" read and write performance guidance from Table 2 in Table 10 below.

| | | Scalable Units (SUs) | | | |
|---|---|---|---|---|---|
| | | **4** | **8** | **32** | **64** |
| **Compute components** | NVIDIA HGX Systems | 127 | 255 | 1023 | 2047 |
| | NVIDIA GPUs | 1016 | 2040 | 8184 | 16376 |
| **DDN Storage components** | DDN Metadata appliances | 1 | 1 | 2 | 4 |
| | DDN Data appliances | 2 | 3 | 12 | 24 |
| **DDN Storage specification** | Aggregate read throughput | 160 GB/s | 320 GB/s | 1280 GB/s | 2560 GB/s |
| | Aggregate write throughput | 80 GB/s | 160 GB/s | 640 GB/s | 1280 GB/s |
| | Per GPU read throughput | 157.48 MB/s | 156.86 MB/s | 156.40 MB/s | 156.33 MB/s |
| | Per GPU write throughput | 78.74 MB/s | 78.43 MB/s | 78.20 MB/s | 78.16 MB/s |
| | Minimum aggregated useable capacity | 240 TB | 360 TB | 1440 TB | 2880 TB |
| | Aggregate useable inodes | 9 billion | 9 billion | 18 billion | 37 billion |
| | Physical, rack units | 6 | 8 | 28 | 56 |
| | Power, nominal | 7 kW | 9 kW | 31 kW | 61 kW |
| | Cooling, nominal | 21 kBTU/hr | 28 kBTU/hr | 101 kBTU/hr | 203 kBTU/hr |

*Table 10. DDN storage sizing for NCP deployments based on "aggregate system" read and write performance guidance.*