# Lenovo Hybrid AI 289 Platform Guide
## Product Guide

The evolution from Generative AI to Agentic AI has revolutionized the landscape of business and enterprise operations globally. By leveraging the capabilities of intelligent agents, companies can now streamline processes, enhance efficiency, and maintain a competitive edge.

These AI agents are adept at handling routine tasks, allowing skilled employees to focus on strategic initiatives and areas where their expertise truly adds value. This symbiotic relationship between AI agents and human employees fosters a collaborative environment that drives innovation and success.

Enterprises must proactively identify opportunities where AI agents can be integrated to support their operations, ensuring they remain agile and effective in an ever-evolving market. This new foundation of AI-driven optimization not only boosts productivity but also empowers employees to contribute more meaningfully to the organization's vision and goals.

Lenovo Hybrid AI 289 is a platform that enables enterprises of all sizes to quickly deploy hybrid AI factory infrastructure, supporting Enterprise AI use cases as either a new, greenfield environment or an extension of their existing IT infrastructure.
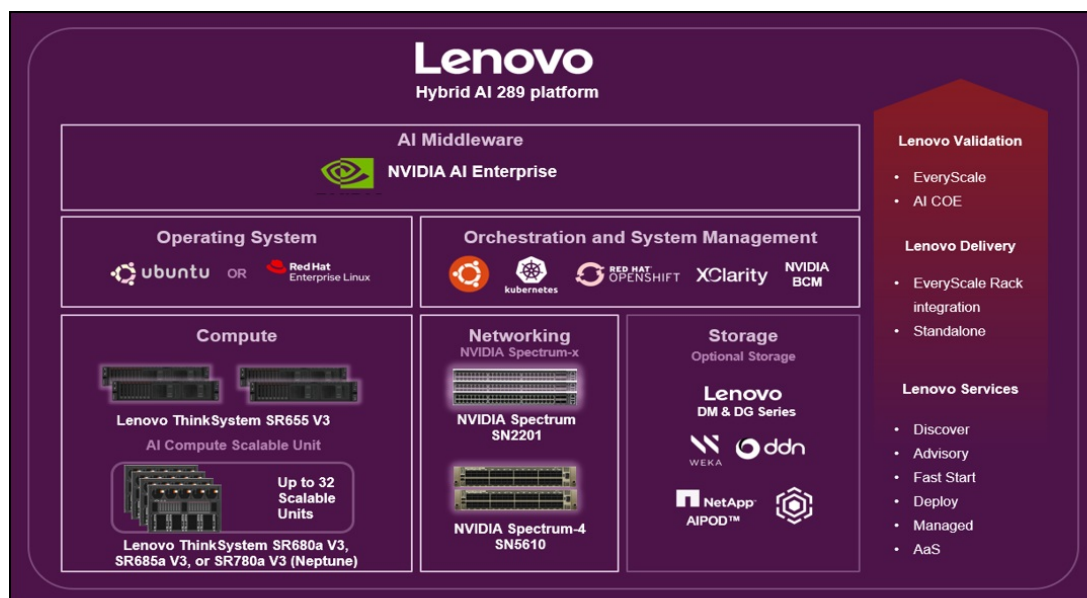


Figure 1. Lenovo Hybrid AI 289 platform overview

The offering is based on the NVIDIA 2-8-9 HGX configuration— 2x CPUs, 8x GPUs, and 9x network adapters — and is ideally suited for model training, fine-tuning, and large-scale inference use cases. It combines market leading Lenovo ThinkSystem GPU-rich servers with NVIDIA Hopper GPUs, NVIDIA Spectrum X networking and enables the use of the NVIDIA AI Enterprise software stack with NVIDIA Blueprints.
Following the principle of *From Exascale to EveryScale™* , Lenovo, widely recognized as a leader in High Performance Computing, leverages its expertise and capabilities from Supercomputing to create tailored enterprise-class hybrid AI factories.

Ideally utilizing Lenovo EveryScale Infrastructure (LESI) it comes with the EveryScale Solution verified interoperability for the tested Best Recipe hardware and software stack. Additionally, EveryScale allows Lenovo Hybrid AI platform deployments to be delivered as fully pre-built, rack-integrated systems that are ready for immediate use.

## Use Cases

The 32 petaflops of processing power per NVIDIA HGX H200 8-GPU baseboard and the additional bandwidth in the East-West compute network makes the 289 platform best suited for LLM fine-tuning and training. The Lenovo Hybrid AI 289 platform supports up to 32 scalable units, totaling 128 AI compute nodes and 1,024 GPUs, with high network bandwidth—ideal for even the most demanding AI training workloads. LLM training times that may have taken weeks with previous AI infrastructures can be reduced to a matter of days using the 289 platform.

Although inference and large-scale RAG applications can be implemented on the 289 platform, medium-sized LLMs with 58B parameters at 16-bit precision can fit on a single H200/B200 SXM GPU, meaning that the added network bandwidth in 289 is unnecessary in medium inference cases. Furthermore, by reducing parameter precision to 4 bits, 4 58B parameter models can fit on one SXM GPU. If per-GPU or per-node inference will be the main use case of your enterprise, consider the Lenovo Hybrid AI 285 platform.

As some of the leading open-source LLMs such as DeepSeek-R1-0528 and Llama 4 Maverick have hundreds of billions of parameters as opposed to tens of billions, inference is becoming more compute intensive. The addition of LLM tooling and RAG only adds to the compute needs of large-scale AI systems. Note that at the time of this document's writing Kimi-K2 is the largest and most performant open-source model, which has 1 trillion parameters at FP8 precision. At that size, Kimi-K2 consumes the entire memory of one HGX H200 8-GPU board, or it can be quantized to 4-bit precision to fit two Kimi-K2 models on one HGX H200. It's reasonable to expect model sizes to grow in the future, requiring more and more compute. Consider whether your enterprise values model accuracy, speed, and/or concurrency when performing model and platform selection. If utilizing the most powerful LLMs for inference and providing those models access to large amounts of data for RAG and tooling, 289 may provide the compute needed, or the 285 platform with B200 GPUs would be required to run per-node inference.

To discuss your enterprise's exact needs and decide on the optimal solution for your use case, please contact your Lenovo representative.

## Overview

The Lenovo 289 Platform comes in 3 main configurations: **3 Scalable Units**, **8 Scalable Units**, and **32 Scalable units**. All 3 configurations provide 100G or 200G connectivity to storage, and 200G connectivity to the enterprise network and support servers.

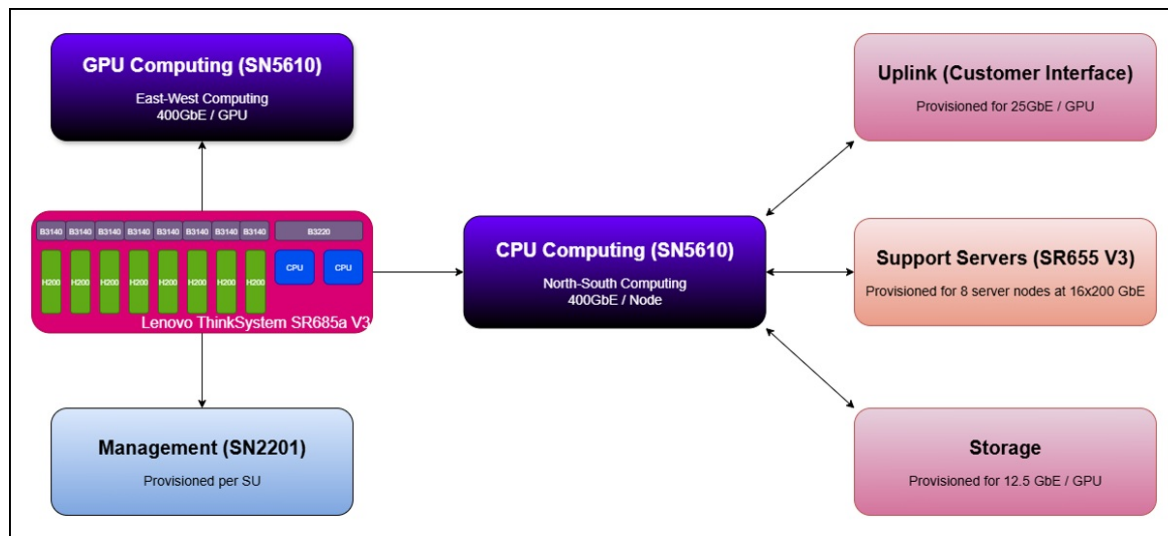See figures below for a sizing overview.

Figure 2. Lenovo Hybrid AI 289 Converged Network Architecture Diagram
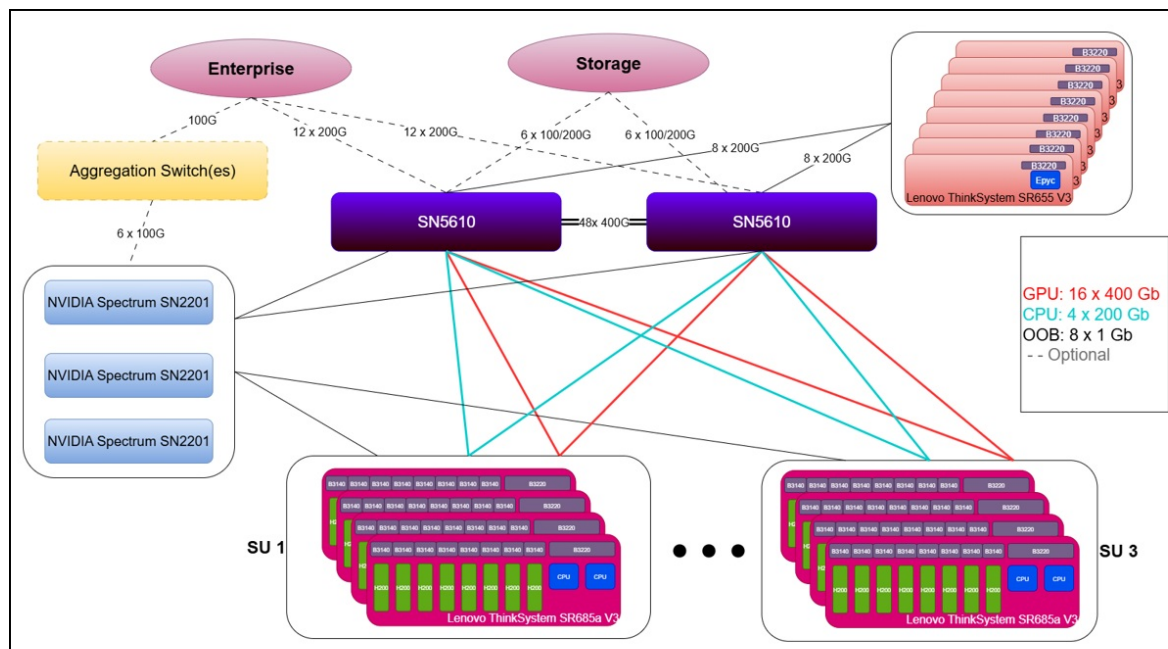


Figure 3. Lenovo Hybrid AI 289 platform with 3 Scalable Units

It can be deployed even for larger sizes with 8 Scalable Units with 32 Servers and 256 GPUs by breaking out the network using four more SN5610s to create a dedicated E/W network for GPU-to-GPU communication.
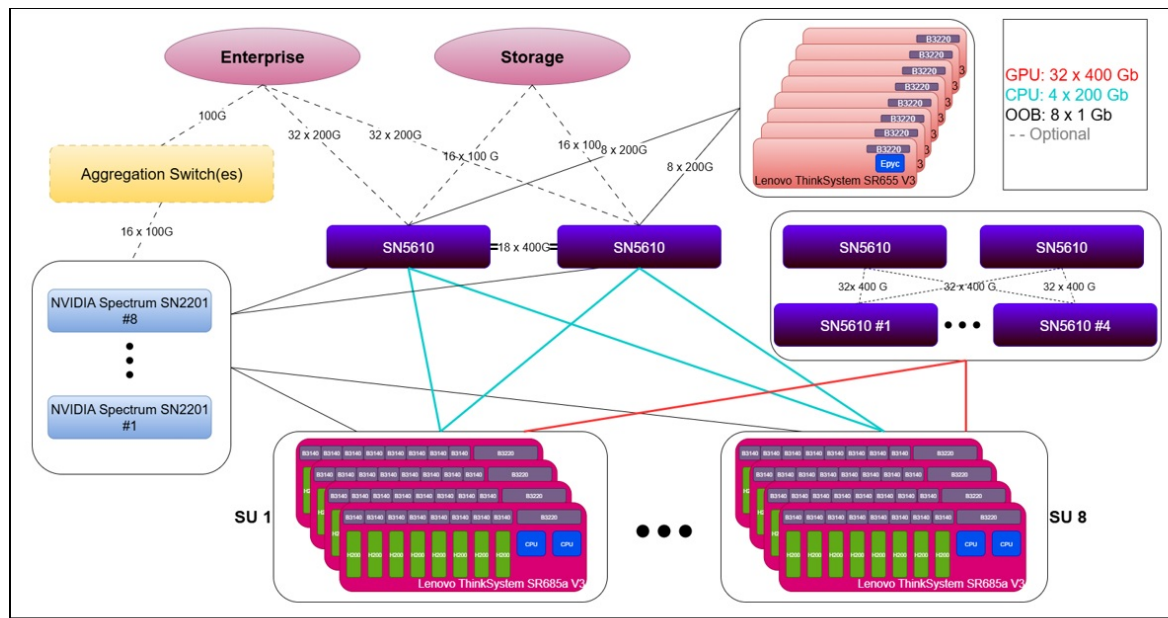
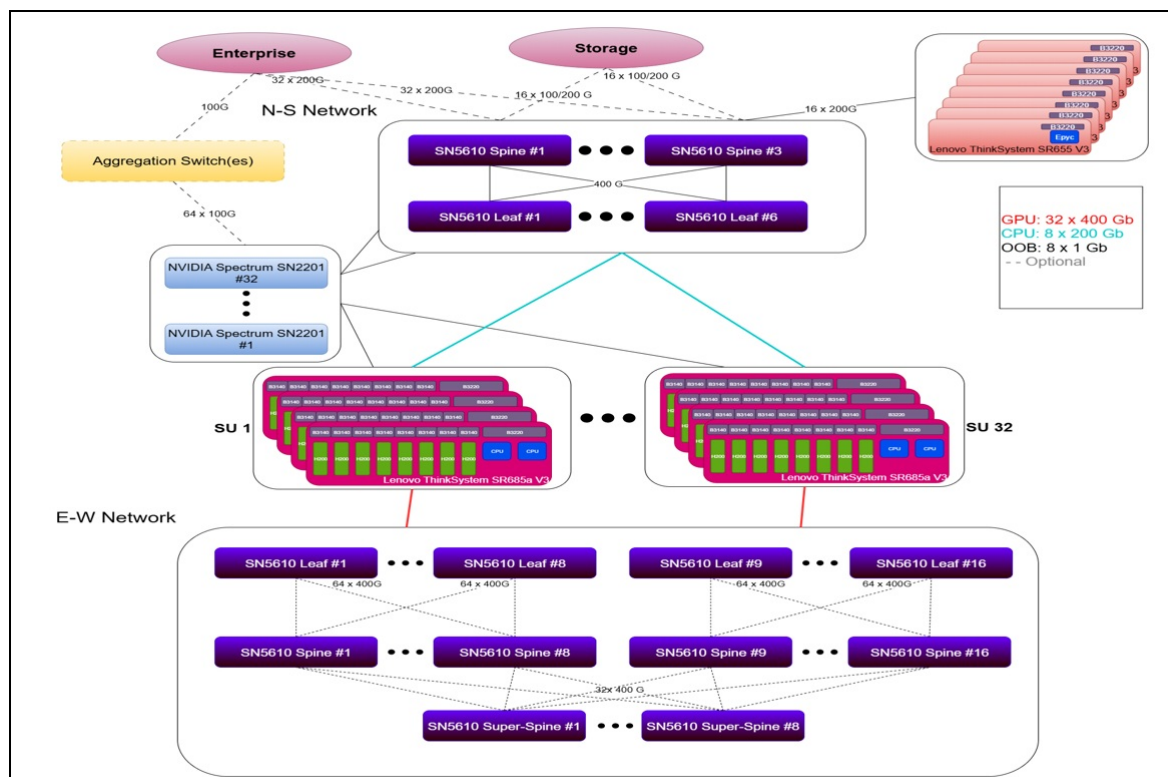Figure 4. Lenovo Hybrid AI 289 platform with 8 Scalable Units



Figure 5. Lenovo Hybrid AI 289 platform with 32 Scalable Units

Note: The implementation of a Super-Spine layer enables ease of scaling beyond 32 SU's. If there is no intention of scaling, the network can be implemented without a Super-Spine layer, which only requires 24 switches as opposed to 40 and leads to reduced costs.

# Low Level Architecture Details

Properly understanding and constructing the network fabric is crucial to getting maximum performance and efficiency out of the 289 platform. The table below shows what resources are needed as the system scales up in a per-SU basis. The high-level architectures in the previous section show the main components and connections between them, while the table below describes exact number of connections between layers of the network.

Table 1. Components for different node counts

| Nodes | GPUs | SUs | Leaf switch count | Spine switch count | Super-spine | Node to leaf Compute Transceivers | Node to leaf Switch Transceivers | Switch to Switch Transceivers |
|---|---|---|---|---|---|---|---|---|
| 4 | 32 | 1 | 2 | 0 | 0 | 32 | 16 | 16 |
| 8 | 64 | 2 | 2 | 0 | 0 | 64 | 32 | 32 |
| 12 | 96 | 3 | 2 | 0 | 0 | 96 | 48 | 48 |
| 16 | 128 | 4 | 2 | 0 | 0 | 128 | 64 | 64 |
| 24 | 192 | 6 | 3 | 2 | 0 | 192 | 96 | 192 |
| 32 | 256 | 8 | 4 | 2 | 0 | 256 | 128 | 256 |
| 48 | 384 | 12 | 6 | 4 | 0 | 384 | 192 | 384 |
| 64 | 512 | 16 | 8 | 8 | 0 | 512 | 256 | 1024 |
| 96 | 768 | 24 | 16 | 16 | 8 | 768 | 384 | 2048 |
| 128 | 1024 | 32 | 16 | 16 | 8 | 1024 | 512 | 2048 |

To maximally utilize the network bandwidth and compute resources when going above 3 SUs, the compute network topology should follow a rail-optimized configuration. In a rail-optimized configuration, the $i^{th}$ NIC on every node is connected to the same switch, which allows the system to take advantage of the NVSwitch technology present on every HGX GPU to reduce interference between flows. This works because on every HGX system, the data can be quickly moved to a GPU connected to the same switch through the NICs as the destination GPU. Below we will go into more detail on the network fabric for these configurations, and permutations can be made upon the specified fabrics when deploying a different amount of SU's by consulting the above table. Note that for the cases where the leaf switch count does not evenly divide 8, rails will have to connect to different NICs/GPUs depending on the node.

In a rail-optimized configuration, a rail refers to the group of NICs/GPUs that are all connected to the same switch. So, to describe the fabric, we will map NICs/GPUs 1 through 8 to one of the available switches in that architecture.

## 8 SU GPU Network Fabric

Table 2.Network fabric for 8 Scalable Units

| Rail | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| GPU / NIC | 1, 5 | 2, 6 | 3, 7 | 4, 8 |

## 32 SU GPU Network Fabric

Recall from the earlier table that above 24 SU's, a Super-Spine layer can be introduced. In the case where a Super Spine is introduced, we will break the SU's into two groups: SU's 1-16 and SU's 17-32. Note that the leaf and spine layers are fully connected from 1 to 8 and 9 to 16, and the Spine and Super-Spine layers are all fully connected.

1-16

For this group of SU's, the GPU/NIC to rail connection is  $i$-to-$i$, meaning that GPU 1 is connected to rail 1, GPU 2 is connected to rail 2, etc.

Table 3. Network fabric for 32 Scalable units, nodes 1 through 16

| Rail | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| GPU / NIC | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |

17-32

For this group of SU's, the GPU/NIC to rail connection is $i$-to-$i + 8$

Table 4. Network fabric for 32 Scalable units, nodes 17 through 32

| Rail | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
|---|---|---|---|---|---|---|---|---|
| GPU / NIC | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |

**Out-Of-Band (OOB) Management**

On a per-node basis there are 2 OOB Management connections, one to the B3220 DPU RJ45 Connector and one to the BMC. While this is the recommended fabric for OOB Management, all B3140 SuperNICs also have RJ45 Connectors that could be added to the management fabric.

## Components

The main hardware components of Lenovo Hybrid AI platforms are Compute nodes and the Networking infrastructure. As an integrated solution they can come together in either a Lenovo EveryScale Rack (Machine Type 1410) or Lenovo EveryScale Client Site Integration Kit (Machine Type 7X74).

Topics in this section:

- AI Compute Nodes and Service Nodes
- Networking
- Lenovo EveryScale Solution

### AI Compute Nodes and Service Nodes

The Lenovo Hybrid AI 289 platform provides three possible GPU server nodes and service nodes to use as the AI Compute Node:

- Lenovo ThinkSystem SR680a V3
- Lenovo ThinkSystem SR685a V3
- Lenovo ThinkSystem SR780a V3 Water Cooled
- Service Nodes – SR655 V3
- Configuration

### Lenovo ThinkSystem SR680a V3

The AI Compute Node leverages the SR680a V3

Figure 6. List of components and connectors in Lenovo ThinkSystem SR680a V3

The SR680a V3 is a 2-socket 5th Gen AMD EPYC 9005 server that supports the NVIDIA HGX H200 system with 9 network adapters in a 8U rack server chassis.



Figure 7. SR680a V3 AI Compute Node Block Diagram

For the CPU, the AI Compute node is configured with two **Intel Xeon Platinum 8570 64 Core 2.4 GHz** processors with an all-core boost frequency of 4GHz. The Xeon Platinum provides 56 cores and 112 threads for each of the Multi Instance GPUs (MIGs). With a 300 MB L3 Cache, this processor excels at accessing frequently used data compared to its AMD EPYC 9535 counterpart, which only has a 256 MB L3 Cache.

**Lenovo ThinkSystem SR685a V3**
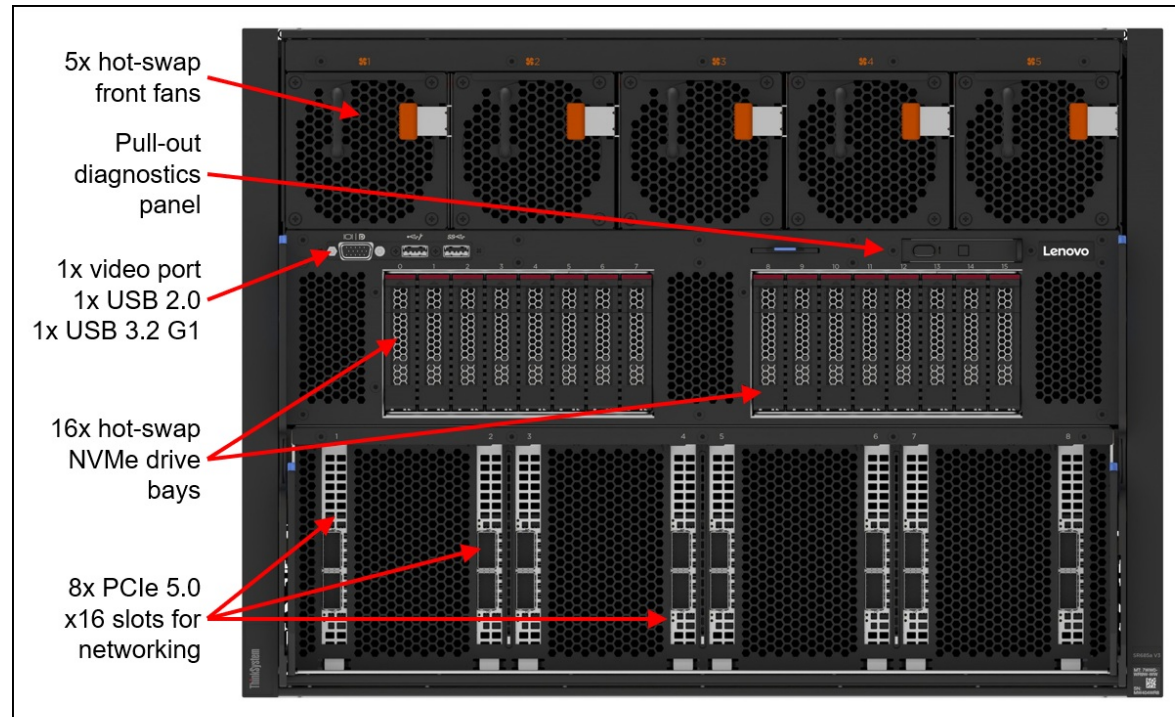
The AI Compute Node leverages the SR685a V3



Figure 8. List of components and connectors in Lenovo ThinkSystem SR685a V3

The SR685a V3 is a 2-socket 5th Gen AMD EPYC 9005 server that supports the NVIDIA HGX H200 system with 9 network adapters in a 8U rack server chassis.
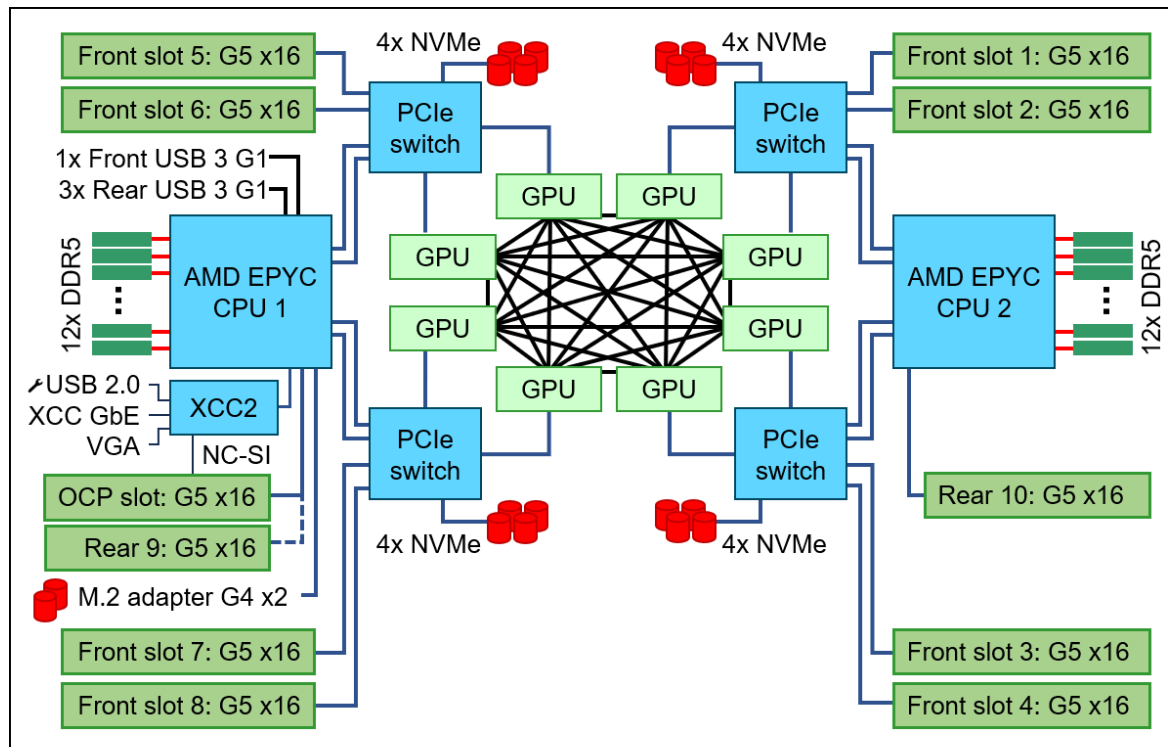
Figure 9. SR680a V3 AI Compute Node Block Diagram

For the CPU, this AI Compute node is configured with two **AMD EPYC 9535 64 Core 2.4 GHz** processors with an all-core boost frequency of 3.5GHz. Besides providing consistently more than 2GHz frequency this ensures that with 7 Multi Instance GPUs (MIG) on 8 physical GPUs there are 2 Cores available per MIG plus a few additional Cores for Operating System and other operations. With 12 Memory Channels per processor socket the AMD based server provides superior Memory bandwidth versus computing Intel-based platforms ensuring highest performance. Leveraging 64GB 6400MHz Memory DIMMs for a total of 1.5TB of main memory providing 192GB memory per GPU or a minimum of 1.5X the HGX H200 GPU memory.

**Lenovo ThinkSystem SR780a V3 Water Cooled**

The AI Compute Node leverages the SR780a V3

Figure 10. List of components and connectors in Lenovo ThinkSystem SR780a V3

The SR780a V3 is a liquid-cooled 2-socket 5th Gen AMD EPYC 9005 server that supports the NVIDIA HGX H200 system with 9 network adapters in a 5U rack server chassis.



Figure 11. AI Compute Node Block Diagram for the SR780a V3

For the CPU, the AI Compute node is configured with two **Intel Xeon Platinum 8570 64 Core 2.4 GHz** processors with an all-core boost frequency of 4GHz (Same as the SR680a V3 CPU). The Xeon Platinum provides 56 cores and 112 threads for each of the Multi Instance GPUs (MIGs). With a 300 MB L3 Cache, this processor excels at accessing frequently used data compared to its AMD EPYC 9535 counterpart which only has a 256 MB L3 Cache.

- **Ports**

  On all 3 servers (SR680a V3, SR685a V3, and SR780a V3), front slots 1 through 8 are connected to NICs 1 through 8 for the B3140 SuperNIC. The B3220 is present on the Rear 9 slot.

- **Ethernet**

  For Converged (North-South) Network an Ethernet adapter with redundant 200Gb/s connections provides ample bandwidth to storage, service nodes and the Enterprise network. The **NVIDIA BlueField-3 B3220 P-Series FHHL DPU** provides the two 200Gb/s Ethernet ports, a 1Gb/s Management board and a 16-Core ARM chip enabling Cloud Orchestration, Storage Acceleration, Secure Infrastructure and Tenant Networking.

  The Ethernet adapters for the Compute (East-West) Network are directly connected to the GPUs via PCIe switches minimizing latency and enabling NVIDIA GPUDirect and GPUDirect Storage operation. For pure Inference workload they are optional, but for training and fine-tuning operation they should provide at least 200Gb/s per GPU.

  By using the **NVIDIA BlueField-3 B3140H E-Series HHHL DPU** or alternatively a ConnectX8 400Gbit Ethernet adapter in combination with NVIDIA Spectrum 4 networking the East-West traffic can utilize Spectrum X operation.

- **Storage**

  The system is completed by local storage with two 960GB Read Intensive M.2 in RAID1 configuration for the operating system and four 3.84TB Read Intensive E3.S drives for local application data.

- **GPU Selection**

  **NVIDIA HGX H200**

  The NVIDIA HGX H200 is a powerful GPU designed to accelerate both generative AI and high-performance computing (HPC) workloads. It boasts a massive 141GB of HBM3e memory, which is nearly double the capacity of its predecessor, the H100. This increased memory, coupled with a 4.8 terabytes per second (TB/s) memory bandwidth, enables the HGX H200 to handle larger and more complex AI models, like large language models (LLMs), with significantly improved performance. The HGX H200 is built with maximum power and scalability with up to 4 Petaflops of FP8 performance. The HGX H200 uses 8 SXM based GPUs that are directly connected rather than PCIe which allows for faster communication with 900 GB/s bandwidth.

  **NVIDIA HGX B200**

  The NVIDIA HGX B200 is designed to supercharge AI and High Performance (HPC) workloads. This GPU is built off the Blackwell architecture, with 180GB of HBM3e memory per GPU. The HGX B200 has 1.8 terabytes per second (TB/s) of GPU-GPU interconnection bandwidth which ensures rapid data transfer and highly efficient processing. Additionally, the HGX B200 uses an advanced liquid-cooling system to effectively dissipate heat and reduce data center electricity costs.

**Service Nodes – SR655 V3**

The AI Service node leverages the SR655 V3

When deploying beyond two AI Compute nodes additional Service nodes are needed to manage the overall AI cluster environment.

Two **Management Nodes** provide a high-availability for the System Management and Monitoring provided through NVIDIA Base Command Manager (BMC) as described further in the AI Software Stack chapter.

For the Container operations three **Scheduling Nodes** build the Kubernetes control plane providing redundant operations and quorum capability.



Figure 12. Lenovo ThinkSystem SR655 V3

The Lenovo ThinkSystem SR655 V3 is an optimal choice for a homogeneous host environment, featuring a single socket AMD EPYC 9335 with 32 cores operating at 3.0 GHz base with an all-core boost frequency of 4.0GHz. The system is fully equipped with twelve 32GB 6400MHz Memory DIMMs, two 960GB Read Intensive M.2 drives in RAID1 configuration for the operating system, and two 3.84TB Read Intensive U.2 drives for local data storage. Additionally, it includes a **NVIDIA BlueField-3 B3220L E-Series FHHL DPU** adapter to connect the Service Nodes to the Converged Network.

**Configuration**

The following table lists the configuration of a single Service Node.

Table 5. ThinkSystem SR655 V3 Service Node

| Part Number | Description | Quantity |
|---|---|---|
| 7D9E-CTOLWW | ThinkSystem SR655 V3 | |
| BLKK | ThinkSystem V3 2U 24x2.5" Chassis | 1 |
| C2AC | ThinkSystem SR655 V3 MB w/IO+PIB+FB,2U | 1 |
| C2AQ | ThinkSystem AMD EPYC 9335 32C 210W 3.0GHz Processor | 1 |
| C0CJ | ThinkSystem 32GB TruDDR5 6400MHz (2Rx8) RDIMM-A | 12 |
| BPQV | ThinkSystem V3 2U x16/x16/E PCIe Gen5 Riser1 or 2 | 1 |
| BVBG | ThinkSystem NVIDIA BlueField-3 B3220 VPI QSFP112 2P 200G PCIe Gen5 x16 Adapter | 1 |
| B8P9 | ThinkSystem M.2 NVMe 2-Bay RAID Adapter | 1 |
| BXMH | ThinkSystem M.2 PM9A3 960GB Read Intensive NVMe PCIe 4.0 x4 NHS SSD | 2 |
| BS7Y | ThinkSystem V3 2U 8x2.5" NVMe Gen5 Backplane | 1 |
| C0ZU | ThinkSystem 2.5" U.2 Multi Vendor 3.84TB Read Intensive NVMe PCIe 5.0 x4 HS SSD | 2 |
| BLL6 | ThinkSystem 2U V3 Performance Fan Module | 6 |
| 6400 | 2.8m, 13A/100-250V, C13 to C14 Jumper Cord | 2 |
| BLKH | ThinkSystem 1100W 230V Titanium Hot-Swap Gen2 Power Supply | 2 |
| B8LA | ThinkSystem Toolless Slide Rail Kit v2 | 1 |
| C1PT | ThinkSystem SR635 V3/SR655 V3 Root of Trust Module Low Voltage-RoW V2 | 1 |
| BQQ6 | ThinkSystem 2U V3 EIA right with FIO | 1 |
| 5PS7B08762 | 5Yr Premier NBD Resp + KYD SR655 V3 | 1 |

**Networking**

The default setup of the Lenovo Hybrid AI 289 platform leverages NVIDIA Networking with the NVIDIA Spectrum-4 SN5610 for the Converged and Compute Network and the NVIDIA SN2201 for the Management Network.

**NVIDIA SN5610**

The SN5610 smart-leaf, spine, and super-spine switch offers 64 ports of 800GbE in a dense 2U form factor. The SN5610 is ideal for NVIDIA Spectrum-X deployments and enables both standard leaf/spine designs with top-of-rack (ToR) switches as well as end-of-row (EoR) topologies. The SN5610 offers diverse connectivity in combinations of 1 to 800GbE and boasts an industry-leading total throughput of 51.2Tb/s.


Figure 13. NVIDIA SN5610 Switch

The **Converged (North-South) Network** handles storage and in-band management, linking the Enterprise IT environment to the Agentic AI platform. Built on Ethernet with RDMA over Converged Ethernet (RoCE), it supports current and new cloud and storage services as outlined in the AI Compute node configuration.

The Converged Network connects to the Enterprise IT network with up to 40 Ethernet connections at 200Gb/s for up to five Scalable Units (SU) or 64 Ethernet connections at 200Gb/s for up to eight SUs. This setup guarantees a minimum bandwidth of 25Gb/s per GPU.

In addition to providing access to the AI agents and functions of the AI platform, this connection is utilized for all data ingestion from the Enterprise IT data during indexing and embedding into the Retrieval-Augmented Generation (RAG) process. It is also used for data retrieval during AI operations.

The Storage connectivity is exactly half that and described in the Storage Connectivity chapter.

The **Compute (East-West) Network** facilitates application communication between the GPUs across the Compute nodes of the AI platform. It is designed to achieve minimal latency and maximal performance using a rail-optimized, fully non-blocking fat tree topology with NVIDIA Spectrum-X.

Spectrum X reduces latency and increases bandwidth for Ethernet by using advanced functionality that splits network packets, tags them, and allows the switch to balance them across network lanes. The receiving node then reassembles the packets regardless of the order they were received in. This process helps avoid hash collusion and congestion, which often lead to suboptimal performance in low-entropy networks.

> **Tip:** In a pure Inference use case, the Compute Network is typically not necessary, but for training and fine-tuning operations it is a crucial component of the solution.

For configurations of up to five Scalable Units, the Compute and Converged Network are integrated utilizing the same switches. When deploying more than five units, it is necessary to separate the fabric.

The following table lists the configuration of the NVIDIA Spectrum-4 SN5610 (link to BOM)

Table 6. NVIDIA Spectrum-4 SN5600 configuration

| Part Number | Description | Quantity |
|---|---|---|
| 7D5FCTONWW | NVIDIA SN5600 800GbE Managed Switch with Cumulus | |
| C0Q5 | NVIDIA SN5600 800GbE Managed Switch with Cumulus | 2 |
| 6252 | 2.5m, 16A/100-250V, C19 to C20 Jumper Cord | 4 |
| 5WS7B98401 | 5Yr Premier NBD Resp NVID SN5600 oPSE | 2 |

**NVIDIA Spectrum SN2201**

The SN2201 is ideal as an out-of-band (OOB) management switch or as a ToR switch connecting up to 48 1G Base-T host ports with non-blocking 100GbE spine uplinks. Featuring highly advanced hardware and software along with ASIC-level telemetry and a 16 megabyte (MB) fully shared buffer, the SN2201 delivers unique and innovative features to 1G switching.



Figure 14. NVIDIA Spectrum SN2201

The **Out-of-Band (Management) Network** encompasses all AI Compute node and BlueField-3 DPU base management controllers (BMC) as well as the network infrastructure management.

The following table lists the configuration of the NVIDIA Spectrum SN2201.

Table 7. NVIDIA Spectrum SN2201 configuration

| Part Number | Description | Quantity |
|---|---|---|
| 7D5FCTOFWW | NVIDIA SN2201 1GbE Managed Switch with Cumulus (PSE) | |
| BPC7 | NVIDIA SN2201 1GbE Managed Switch with Cumulus (PSE) | 1 |
| 6201 | 1.5m, 10A/100-250V, C13 to C14 Jumper Cord | 2 |
| 5WS7B98268 | 5Yr Premier NBD Resp NVID SN2201 PSE | 1 |

**Lenovo EveryScale Solution**

The Server and Networking components and Operating System can come together as a Lenovo EveryScale Solution. It is a framework for designing, manufacturing, integrating and delivering data center solutions, with a focus on High Performance Computing (HPC), Technical Computing, and Artificial Intelligence (AI) environments.

Lenovo EveryScale provides Best Recipe guides to warrant interoperability of hardware, software and firmware among a variety of Lenovo and third-party components.

Addressing specific needs in the data center, while also optimizing the solution design for application performance, requires a significant level of effort and expertise. Customers need to choose the right hardware and software components, solve interoperability challenges across multiple vendors, and determine optimal firmware levels across the entire solution to ensure operational excellence, maximize performance, and drive best total cost of ownership.

Lenovo EveryScale reduces this burden on the customer by pre-testing and validating a large selection of Lenovo and third-party components, to create a "Best Recipe" of components and firmware levels that work seamlessly together as a solution. From this testing, customers can be confident that such a best practice solution will run optimally for their workloads, tailored to the client's needs.

In addition to interoperability testing, Lenovo EveryScale hardware is pre-integrated, pre-cabled, pre-loaded with the best recipe and optionally an OS-image and tested at the rack level in manufacturing, to ensure a reliable delivery and minimize installation time in the customer data center.

## Configurations

**Configurations**

The 2-8-9 platform can come in 3 main variants: 3 Scalable Units, 8 Scalable Units, and 32 Scalable Units, depending on the use case. Each node requires at a minimum 8,555 W, so 4 nodes in a rack would require ~34,000W. If your rack cannot support these power requirements, the nodes can be split into two racks per SU

Power requirements:

- SR680a: 8996 W
- SR685a: 8978.2 W
- SR780a: 8555 W

The networking decision depends on whether the platform is designed to support up to three, eight, or thirty-two Scalable Units in total. Subsequently, the solution can be expanded seamlessly without downtime by incorporating additional Scalable Units, ultimately reaching a total of eight or 32 as needed.

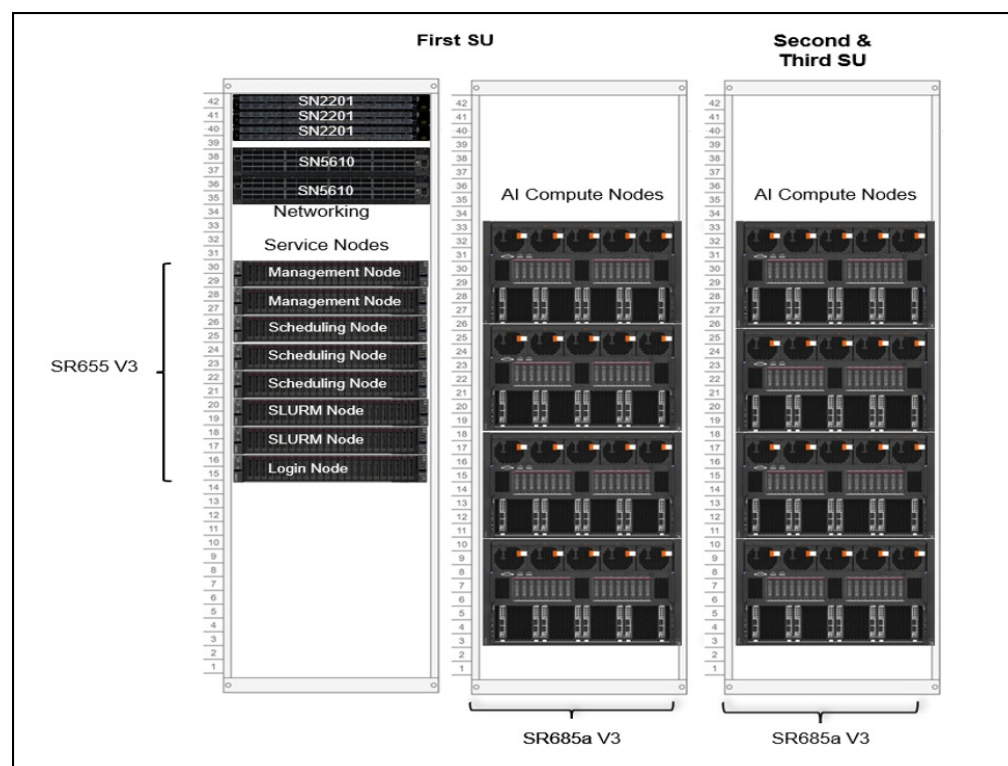The following describe the three and eight Scalable Units deployments:
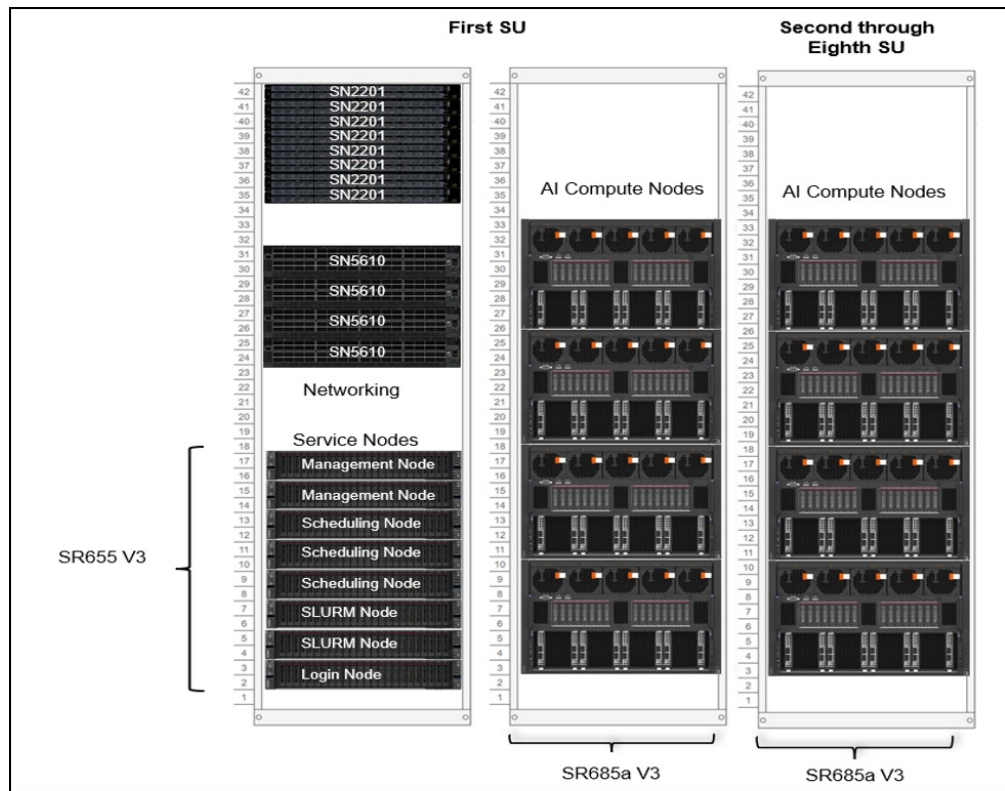


Figure 15. 3 SU Deployment

Figure 16. 8 SU Deployment

## Scalability

### Custom Deployment

For high-end scenarios requiring more than eight scalable units, the network can be custom designed to any required size. Lenovo will develop a fully bespoke solution tailored to match the workflow and workload requirements in that case.

Note: For the 3SU configuration, Direct Attach Copper (DAC) cabling can be used to reduce cost while maintaining the necessary bandwidth. This cabling will reduce costs when compared to using fiber-optic cabling.

## AI Software Stack

Deploying AI to production involves implementing multiple layers of software. The process begins with the system management and operating system of the compute nodes, progresses through a workload or container scheduling and cluster management environment, and culminates in the AI software stack that enables delivering AI agents to users.

The following table provides all of the recommended software layers and their roles in the Lenovo Hybrid AI 289 platform.

Note that not all software is required to function; however, this is the recommended stack.

Table 8. AI Software Stack

| Software Role | Software Package |
|---|---|
| Operating System | Ubuntu |

| Software Role | Software Package |
|---|---|
| Orchestration | Base Command Manager Essentials and XClarity |
| Container Runtime | Containerd |
| Container Orchestration | Kubernetes |
| Container Network Interface (CNI) | Calico |
| Load Balancer - API Service Gateway | Nginx |
| Load Balancer – Network Services | MetalLB |
| Ingress Controller | Nginx |
| Package Manager | Helm |
| Operator | GPU Operator |
| Operator | Network Operator |
| Operator | NIM Operator |
| RBAC | Permission Manager |
| Observability | Prometheus |
| Observability | Grafana |
| Observability | NVIDIA NetQ |
| Storage | NFS Provisioner |

In the following sections, we take a deeper dive into the software elements:

- XClarity System Management
- Linux Operating System
- Base Command Manager - Provisioning (NVAIE)
- Kubernetes Container Orchestration
- Data and AI Applications
- NVIDIA AI Enterprise

## XClarity System Management

Lenovo XClarity Administrator is a centralized, resource-management solution that simplifies infrastructure management, speeds responses, and enhances the availability of Lenovo server systems and solutions. It runs as a virtual appliance that automates discovery, inventory, tracking, monitoring, and provisioning for server, network, and storage hardware in a secure environment.

Table 9. XClarity System Management

| Part Number | Description |
|---|---|
| 00MT203 | Lenovo XClarity Pro, Per Managed Endpoint w/5 Yr SW S&S |

## Linux Operating System

The AI Compute nodes are deployed with Linux. Traditionally  Canonical Ubuntu is the default choice for AI environments with its optimizations across all prominent AI hardware platforms and up to 12 years of security maintenance, but Lenovo Hybrid AI platforms support Red Hat Enterprise Linux (RHEL) as alternative choice.

Table 10. Linux Operating System

| Part number | Description |
|---|---|
| 7S1B000YWW | Canonical Ubuntu Pro 5Yr w/Canonical weekday Support |

### Base Command Manager - Provisioning (NVAIE)

Base Command Manager (BCM) provisions the AI environment, incorporating the components such as the Operating System, Vanilla Kubernetes (K8S), GPU Operator, and Network Operator to manage the AI workloads. BCM Supports 3 types of network topologies depending on how the user wants nodes to be accessed.

- Type 1: All communication is centralized through the head node, providing a controlled and secured gateway
- Type 2: Worker nodes can be accessed directly via a router, so that traffic does not need to go through the head node.
- Type 3: A routed public network is used, where regular nodes are on Internalnet and the head node is on Managementnet.

The following table displays the recommended BCM Networks for a Type 2 network.

Table 11. BCM Networks

| BCM Network | Traffic Type | Purpose | Creation Phase |
|---|---|---|---|
| managementnet | IPMI | OOB Management | BCM Deployment |
| internalnet | Management/Data | In-band management | BCM deployment |
| storagenet | Storage | Backend Storage | Storage Configuration |
| failovernet | High Availability | HA Heartbeat | HA Configuration, setup by deployment wizard |
| Kube-cluster-pod | Pod network | Pod-to-pod traffic | K8s deployment, setup by deployment wizard |
| Kube-cluster-service | Service Network | App traffic | K8s deployment, setup by deployment wizard |

### Kubernetes Container Orchestration

Canonical Ubuntu Pro comes with Kubernetes, the leading AI container deployment and workload management tool in the market, which can be used for edge and centralized data center deployments.

Canonical Kubernetes is used across industries for mission critical workloads, and uniquely offers up to 12 years of security for those customers who cannot, or choose not to upgrade their Kubernetes versions.

When implementing a Scalable Unit Deployment and above Ubuntu Charmed Kubernetes is used.

When choosing Red Hat for the Operating System, Red Hat OpenShift as the matching Kubernetes implementation is required.

### Data and AI Applications

Canonical's Ubuntu Pro includes a portfolio of open source applications in the data and AI space including leading projects for ML space with Kubeblow and MLFlow, big data and database with Spark, Kafka, PostgreSQL, Mongo and others. Ubuntu Pro enables customers on their open source AI journey to simplify deployment and maintenance of these applications and provides security maintenance.

### NVIDIA AI Enterprise

The Lenovo Hybrid AI 289 platform is designed for NVIDIA AI Enterprise, which is a comprehensive suite of artificial intelligence and data analytics software designed for optimized development and deployment in enterprise settings.

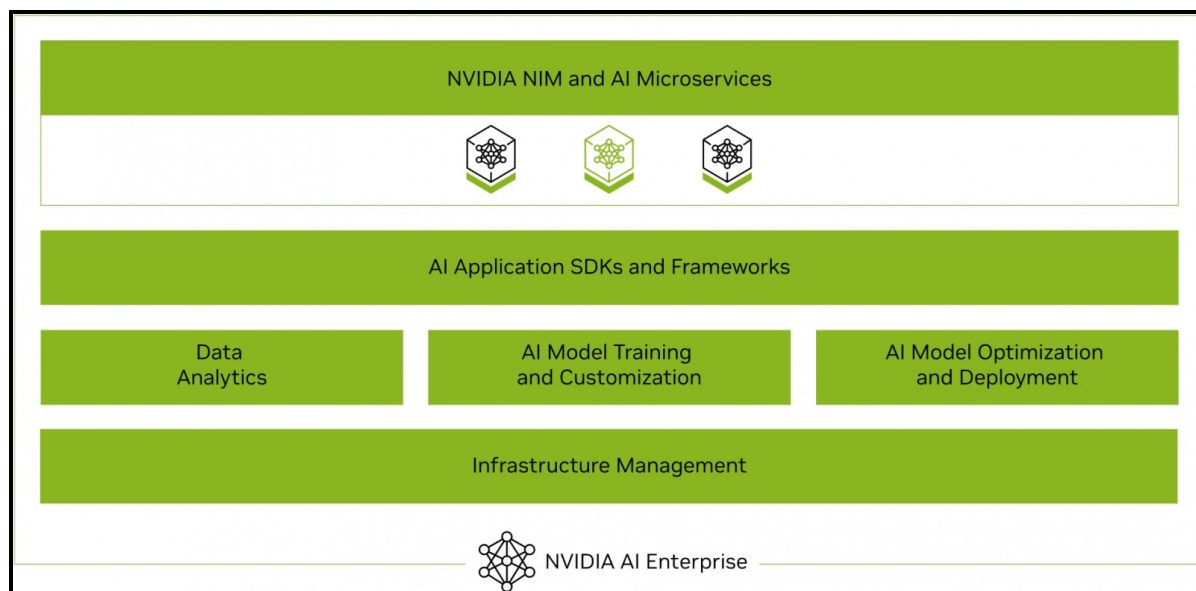**Tip**: Entitlement to NVIDIA AI Enterprise for 5 years is included with the NVIDIA H200 NVL PCIe GPU.

Figure 17. NVIDIA AI Enterprise software stack

NVIDIA AI Enterprise includes workload and infrastructure management software known as Base Command Manager. This software provisions the AI environment, incorporating the components such as the Operating System, Kubernetes (K8S), GPU Operator, and Network Operator to manage the AI workloads.

Additionally, NVIDIA AI Enterprise provides access to ready-to-use open-sourced containers and frameworks from NVIDIA like NVIDIA NeMo, NVIDIA RAPIDS, NVIDIA TAO Toolkit, NVIDIA TensorRT and NVIDIA Triton Inference Server.

- **NVIDIA NeMo** is an end-to-end framework for building, customizing, and deploying enterprise-grade generative AI models; NeMo lets organizations easily customize pretrained foundation models from NVIDIA and select community models for domain-specific use cases.

- **NVIDIA RAPIDS** is an open-source suite of GPU-accelerated data science and AI libraries with APIs that match the most popular open-source data tools. It accelerates performance by orders of magnitude at scale across data pipelines.

- **NVIDIA TAO Toolkit** simplifies model creation, training, and optimization with TensorFlow and PyTorch and it enables creating custom, production-ready AI models by fine-tuning NVIDIA pretrained models and large training datasets.

- **NVIDIA TensorRT**, an SDK for high-performance deep learning inference, includes a deep learning inference optimizer and runtime that delivers low latency and high throughput for inference applications. TensorRT is built on the NVIDIA CUDA parallel programming model and enables you to optimize inference using techniques such as quantization, layer and tensor fusion, kernel tuning, and others on NVIDIA GPUs. https://developer.nvidia.com/tensorrt-getting-started

- **NVIDIA TensorRT-LLM** is an open-source library that accelerates and optimizes inference performance of the latest large language models (LLMs). TensorRT-LLM wraps TensorRT's deep learning compiler and includes optimized kernels from FasterTransformer, pre- and post-processing, and multi-GPU and multi-node communication. https://developer.nvidia.com/tensorrt

- **NVIDIA Triton Inference Server** optimizes the deployment of AI models at scale and in production for both neural networks and tree-based models on GPUs.

It also provides full access to the NVIDIA NGC catalogue, a collection of tested enterprise software, services and tools supporting end-to-end AI and digital twin workflows and can be integrated with MLOps platforms such as ClearML, Domino Data Lab, Run:ai, UbiOps, and Weights & Biases.

Finally, NVIDIA AI Enterprise introduced NVIDIA Inference Microservices (NIM), a set of performance-optimized, portable microservices designed to accelerate and simplify the deployment of AI models. Those containerized GPU-accelerated pretrained, fine-tuned, and customized models are ideally suited to be self-hosted and deployed on the Lenovo Hybrid AI 289 platform.

The ever-growing catalog of NIM microservices contains models for a wide range of AI use cases, from chatbot assistants to computer vision models for video processing. The image below shows some of the NIM microservices, organized by use case.
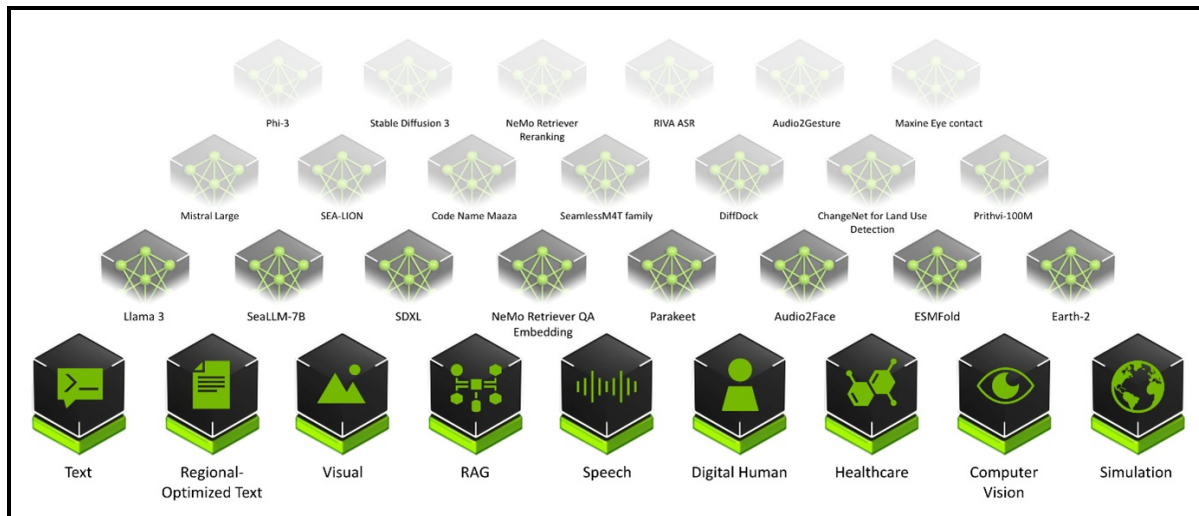


Figure 18. NVIDIA Inference Microservice catalog

## Storage Connectivity

Lenovo Hybrid AI platforms do not include storage but do interface with any storage technology that is validated by NVIDIA for NVIDIA OVX or certified by the NVIDIA Certified program.

Lenovo Storage validated by NVIDIA includes: Lenovo DM Series, Lenovo DG Series, Lenovo ThinkAgile HX series with Nutanix and ThinkAgile VX series with VMware, and Lenovo DSS-G for IBM Storage Scale. Lenovo is a qualified hardware platform for Cloudian, DDN, and WEKA.

The storage directly attached to the AI platform primarily hosts the vector database supporting data retrieval for Retrieval-Augmented Generation (RAG) applications. Additionally, it functions as high-performance storage for retraining or fine-tuning models.

The Converged Network connects to the Storage with up to 20 Ethernet connections at 200Gb/s for up to five Scalable Units (SU) or 32 Ethernet connections at 200Gb/s for up to eight SUs. This setup guarantees a minimum bandwidth of 12.5Gb/s per GPU. The connected storage system should be configured to support this bandwidth requirement per GPU accordingly.

### ThinkSystem DM & DG Series Storage

For enterprise organizations AI applications are treated as any other workload and require the same data management and enterprise data security features. The new DM7200F and DG5200 platforms provide all flash storage that prepares your infrastructure and data for AI workloads.

Benefits of DM & DG storage for GenAI & RAG:

- Enterprise Security features including autonomous ransomware protection
- Deduplication and compression
- All flash performance
- Flexible scaling
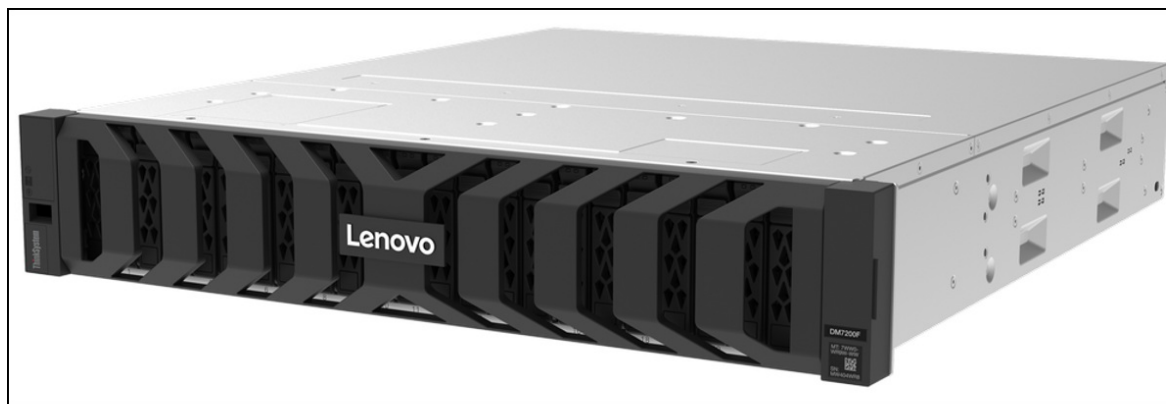- Unified file, object, and block eliminates data silos



Figure 19. Lenovo ThinkSystem DM7200F Unified Storage Arrays

## Lenovo AI Center of Excellence

In addition to the choice of utilizing Lenovo EveryScale Infrastructure framework for the Enterprise AI platform to ensure tested and warranted interoperability, Lenovo operates an AI Lab and CoE at the headquarters in Morrisville, North Carolina, USA to test and enable AI applications and use cases on the Lenovo EveryScale AI platform.

The AI Lab environment allows customers and partners to execute proof of concepts for their use cases or test their AI middleware or applications. It is configured as a diverse AI platform with a range of systems and GPU options, including NVIDIA L40S and NVIDIA HGX8 H200.

The software environment utilizes Canonical Ubuntu Linux along with Canonical MicroK8s to offer a multi-tenant Kubernetes environment. This setup allows customers and partners to schedule their respective test containers effectively.

### Lenovo AI Innovators

Lenovo Hybrid AI platforms offer the necessary infrastructure for a customer's hybrid AI factory. To fully leverage the potential of AI integration within business processes and operations, software providers, both large and small, are developing specialized AI applications tailored to a wide array of use cases.

To support the adoption of those AI applications, Lenovo continues to invest in and extend its AI Innovators Program to help organizations gain access to enterprise AI by partnering with more than 50 of the industry's leading software providers.

Partners of the Lenovo AI Innovators Program get access to our AI Discover Labs, where they validate their solutions and jointly support Proof of Concepts and Customer engagements.

LAII provides customers and channel partners with a range of validated solutions across various vertical use cases, such as for Retail or Public Security. These solutions are designed to facilitate the quick and safe deployment of AI solutions that optimally address the business requirements.

The following is a selection of case studies involving Lenovo customers implementing an AI solution:

- Kroeger (Retail) – Reducing Customer friction and loss prevention
- Peak (Logistics) – Streamlining supply chain ops for fast and efficient deliveries
- Bikal (AI at Scale) – Delivering shared AI platform for education
- VSAAS (Smart Cities) – Enabling accurate and effective public security

### Lenovo Validated Designs

Lenovo Validated Designs (LVDs) are pre-tested, optimized solution designs enabling reliability, scalability, and efficiency in specific workloads or industries. These solutions integrate Lenovo hardware like ThinkSystem servers, storage, and networking with software and best practices to solve common IT challenges. Developed with technology partners such as VMware, Intel, and Red Hat, LVDs ensure performance, compatibility, and easy deployment through rigorous validation.

Lenovo Validated Designs are intended to simplify the planning, implementation, and management of complex IT infrastructures. They provide detailed guidance, including architectural overviews, component models, deployment considerations, and bills of materials, tailored to specific use cases such as artificial intelligence (AI), big data analytics, cloud computing, virtualization, retail, or smart manufacturing. By offering a pretested solution, LVDs aim to reduce risk, accelerate deployment, and assist organizations in achieving faster time-to-value for their IT investments.

Lenovo Hybrid AI platforms act as infrastructure frameworks for LVDs addressing data center-based AI solutions. They provide the hardware/software reference architecture, optionally Lenovo EveryScale integrated solution delivery method, and general sizing guidelines.

## AI Services

Lenovo Hybrid AI platforms solutions are specifically designed to enable broad adoption in the Enterprise supported by Lenovo's powerful IT partner ecosystem.

In addition to custom deployments, as the foundation of NVIDIA's Enterprise Reference Architecture, it is fully compatible with NVIDIA Blueprints for agentic and generative AI use cases.

This enables both Lenovo AI Partners and Lenovo Professional Services to accelerate deployment and provide enterprises with the fastest time to production.

Lenovo Hybrid AI Advantage enables customers to overcome the barriers they face in realizing ROI from AI investments by providing critical expertise needed to accelerate business outcomes. Leveraging a responsible approach to AI, Lenovo AI expertise, Lenovo's advanced partner ecosystem and industry leading technology we help customers realize the benefits of AI faster.

### AI Discover Workshop

Lenovo AI Discover Workshops help customers visualize and map out their strategy and resources for AI adoption to rapidly unlock real business value. Lenovo's experts assess the organization's AI readiness across security, people, technology, and process – a proven methodology – with recommendations that put customers on a path to AI success. With a focus on real outcomes, AI Discover leverage proven frameworks, processes and policies to deliver a technology roadmap that charts the path to AI success.

### AI Fast Start

With customers looking to unlock the transformative power of AI, Lenovo AI Fast Start empowers customers to rapidly build and deploy production-ready AI solutions tailored to their needs. Optimized for NVIDIA AI Enterprise and leveraging accelerators like NVIDIA NIMs, Lenovo AI Fast Start accelerates use case development and platform readiness for AI deployment at scale allowing customers to go from concept to production ready deployment in just weeks. Easy to use containerized and optimized inference engines for popular NVIDIA AI Foundation models empower developers to deliver results. AI Fast Start provides access to AI Experts, platforms and technologies supporting onsite and remote models to achieve business objectives.

## Lenovo TruScale

Lenovo TruScale XaaS is your set of flexible IT services that makes everything easier. Streamline IT procurement, simplify infrastructure and device management, and pay only for what you use – so your business is free to grow and go anywhere.

Lenovo TruScale is the unified solution that gives you simplified access to:

- The industry's broadest portfolio – from pocket to cloud – all delivered as a service
- A single-contract framework for full visibility and accountability
- The global scale to rapidly and securely build teams from anywhere
- Flexible fixed and metered pay-as-you-go models with minimal upfront cost
- The growth-driving combination of hardware, software, infrastructure, and solutions – all from one single provider with one point of accountability.

For information about Lenovo TruScale offerings that are available in your region, contact your local Lenovo sales representative or business partner.

## Lenovo Financial Services

Why wait to obtain the technology you need now? No payments for 90 days and predictable, low monthly payments make it easy to budget for your Lenovo solution.

- **Flexible**
  Our in-depth knowledge of the products, services and various market segments allows us to offer greater flexibility in structures, documentation and end of lease options.

- **100% Solution Financing**
  Financing your entire solution including hardware, software, and services, ensures more predictability in your project planning with fixed, manageable payments and low monthly payments.

- **Device as a Service (DaaS)**
  Leverage latest technology to advance your business. Customized solutions aligned to your needs. Flexibility to add equipment to support growth. Protect your technology with Lenovo's Premier Support service.

- **24/7 Asset management**
  Manage your financed solutions with electronic access to your lease documents, payment histories, invoices and asset information.

- **Fair Market Value (FMV) and $1 Purchase Option Leases**
  Maximize your purchasing power with our lowest cost option. An FMV lease offers lower monthly payments than loans or lease-to-own financing. Think of an FMV lease as a rental. You have the flexibility at the end of the lease term to return the equipment, continue leasing it, or purchase it for the fair market value. In a $1 Out Purchase Option lease, you own the equipment. It is a good option when you are confident you will use the equipment for an extended period beyond the finance term. Both lease types have merits depending on your needs. We can help you determine which option will best meet your technological and budgetary goals.

Ask your Lenovo Financial Services representative about this promotion and how to submit a credit application. For the majority of credit applicants, we have enough information to deliver an instant decision and send a notification within minutes.

## Bill of Materials

**Example BoMs for Deployment Sizes**

Consult the full BoM for per-system to determine the exact quantities for each deployment size. Outlined below is a high level overview of each deployment size will require.

**3 SU Example Bill of Materials**

- 12x Lenovo ThinkSystem SR685a V3 with HGX H200 GPU
- 8x Lenovo ThinkSystem SR655 V3
- 2x NVIDIA SN5610 Switches
- 3x NVIDIA SN2201 Switches
- 4x Onyx Heavy Duty Rack Cabinet

**8 SU Example Bill of Materials**

- 32x Lenovo ThinkSystem SR685a V3 with HGX H200 GPU
- 8x Lenovo ThinkSystem SR655 V3
- 8x NVIDIA SN5610 Switches
- 8x NVIDIA SN2201 Switches
- 9x Onyx Heavy Duty Rack Cabinet

**32 SU Example Bill of Materials**

- 128x Lenovo ThinkSystem SR685a V3 with HGX H200 GPU
- 8x Lenovo ThinkSystem SR655 V3
- 49x NVIDIA SN5610 Switches
- 32x NVIDIA SN2201 Switches
- 37x Onyx Heavy Duty Rack Cabinet

Storage is optional and not included in this BoM.

**Per System Bill of Materials**

In this section:

- ThinkSystem SR685a V3 BOM
- ThinkSystem SR680a V3 BOM
- ThinkSystem SR780a V3
- ThinkSystem SR780a V3 B200 BoM
- ThinkSystem SR655 V3 BoM
- NVIDIA SN5610 Switch BoM
- NVIDIA SN2201 Switch BoM
- Power Distribution Unit (PDU) BoM
- Rack Cabinet BoM
- XClarity Software BoM

**ThinkSystem SR685a V3 BOM**

Table 12. ThinkSystem SR685a V3 BoM

| Part Number | Product Description | Qty per System | Total Qty |
|---|---|---|---|
| 7DHCCTO1WW | Server : ThinkSystem SR685a V3 - 3 Year Warranty with Controlled GPU | 1 | 12 |
| C1EX | ThinkSystem SR685a V3 H200 GPU Base | 1 | 12 |
| C2AL | ThinkSystem AMD EPYC 9535 64C 300W 2.4GHz Processor | 2 | 24 |
| C0CP | ThinkSystem 96GB TruDDR5 6400MHz (2Rx4) RDIMM-A | 24 | 288 |
| BTQ1 | ThinkSystem 2.5" U.3 PM1743 7.68TB Read Intensive NVMe PCIe 5.0 x4 HS SSD | 2 | 24 |
| BT3B | ThinkSystem SR850 V3/SR860 V3 8x 2.5" AnyBay Backplane | 2 | 24 |
| B8P9 | ThinkSystem M.2 NVMe 2-Bay RAID Adapter | 1 | 12 |
| BXMH | ThinkSystem M.2 PM9A3 960GB Read Intensive NVMe PCIe 4.0 x4 NHS SSD | 2 | 24 |
| C0Q4 | ThinkSystem NVIDIA BlueField-3 B3140H VPI QSFP112 1P 400G PCIe Gen5 x16 Adapter | 8 | 96 |
| BVBG | ThinkSystem NVIDIA BlueField-3 B3220 VPI QSFP112 2P 200G PCIe Gen5 x16 Adapter | 1 | 12 |
| C1HM | ThinkSystem NVIDIA HGX H200 141GB 700W 8-GPU Board | 1 | 12 |
| C1EP | ThinkSystem SR685a V3 x16 PCIe Rear IO Riser | 1 | 12 |
| C8VK | ThinkSystem SR685a V3 8x x16 PCIe GPU Switch Riser for Turin | 1 | 12 |
| C4HK | ThinkSystem 2600W 230V Titanium Hot-Swap Gen2 Power Supply v4 | 8 | 96 |
| 6252 | 2.5m, 16A/100-250V, C19 to C20 Jumper Cord | 8 | 96 |
| C84K | ThinkSystem SR685a V3 System Board for Turin | 1 | 12 |
| C1EN | ThinkSystem SR685a V3 Midplane Cable Assembly Switch Side | 1 | 12 |
| BAVU | Front Operator Panel with Quad Line LCD Display | 1 | 12 |
| C1EV | ThinkSystem SR685a V3 Midplane Cable Assembly CPU Side | 1 | 12 |
| C1ET | ThinkSystem SR685a V3 Root of Trust Module | 1 | 12 |
| C1EQ | ThinkSystem SR685a V3 CPU Complex Power Interface Board | 1 | 12 |
| BQ27 | ThinkSystem SR665 V3/SR655 V3 Standard Heatsink | 2 | 24 |
| BABV | ThinkSystem Screw for fix M.2 Adapter | 1 | 12 |
| C1FH | ThinkSystem SR685a V3 / SR680a V3 Front Fan | 5 | 60 |
| C1FG | ThinkSystem SR685a V3 / SR680a V3 Rear Fan | 10 | 120 |
| C1FF | ThinkSystem 8GPU Server Front Fan Control Board | 1 | 12 |
| C6UG | ThinkSystem SR685a V3 8GPU Server Power Distribution Board | 1 | 12 |
| C1EY | ThinkSystem 8GPU Server CFFV4 Power Interface Board | 1 | 12 |
| 5641PX3 | XClarity Pro, Per Endpoint w/3 Yr SW S&S | 1 | 12 |
| 1340 | Lenovo XClarity Pro, Per Managed Endpoint w/3 Yr SW S&S | 1 | 12 |
| 7Q01CTSAWW | SERVER KEEP YOUR DRIVE ADD-ON | 1 | 12 |
| 7Q01CTS4WW | SERVER PREMIER 24X7 4HR RESP | 1 | 12 |

**ThinkSystem SR680a V3 BOM**

Table 13. ThinkSystem SR680a V3 BoM

| Part Number | Product Description | Qty per System | Total Qty |
|---|---|---|---|
| 7DHECTO1WW | Server : ThinkSystem SR680a V3 - 3 Year Warranty with Controlled GPU | 1 | 12 |
| C1EL | ThinkSystem SR680a V3 H200 GPU Base | 1 | 12 |
| BYWG | Intel Xeon Platinum 8570 56C 350W 2.1GHz Processor | 2 | 24 |
| C5H9 | ThinkSystem 64GB TruDDR5 5600MHz (2Rx4) 10x4 RDIMM | 32 | 384 |
| BTQ1 | ThinkSystem 2.5" U.3 PM1743 7.68TB Read Intensive NVMe PCIe 5.0 x4 HS SSD | 2 | 24 |
| BT3B | ThinkSystem SR850 V3/SR860 V3 8x 2.5" AnyBay Backplane | 2 | 24 |
| BXMH | ThinkSystem M.2 PM9A3 960GB Read Intensive NVMe PCIe 4.0 x4 NHS SSD | 1 | 12 |
| BS7M | Intel VROC (VMD NVMe RAID) Standard for M.2 | 1 | 12 |
| C0Q4 | ThinkSystem NVIDIA BlueField-3 B3140H VPI QSFP112 1P 400G PCIe Gen5 x16 Adapter | 8 | 96 |
| BVBG | ThinkSystem NVIDIA BlueField-3 B3220 VPI QSFP112 2P 200G PCIe Gen5 x16 Adapter | 1 | 12 |
| C1HM | ThinkSystem NVIDIA HGX H200 141GB 700W 8-GPU Board | 1 | 12 |
| C1F3 | ThinkSystem SR680a V3 x16 PCIe Rear IO Riser | 2 | 24 |
| C1F2 | ThinkSystem SR685a V3 / SR680a V3 8x x16 PCIe GPU Switch Riser | 1 | 12 |
| C4HK | ThinkSystem 2600W 230V Titanium Hot-Swap Gen2 Power Supply v4 | 8 | 96 |
| B4L2 | 2.0m, 16A/100-250V, C19 to C20 Jumper Cord | 8 | 96 |
| AUMY | ThinkSystem Lift Handles | 1 | 12 |
| A4AA | ThinkSystem Toolless Fixed Rail Kit | 1 | 12 |
| C1FJ | ThinkSystem SR680a V3 System Board | 1 | 12 |
| C1F4 | ThinkSystem SR680a V3 Midplane Cable Assembly Switch side | 1 | 12 |
| C1FG | ThinkSystem SR685a V3 / SR680a V3 Rear Fan | 10 | 120 |
| BAVU | Front Operator Panel with Quad Line LCD Display | 1 | 12 |
| C1FF | ThinkSystem 8GPU Server Front Fan Control Board | 1 | 12 |
| C1FE | ThinkSystem SR680a V3 Midplane Cable Assembly CPU side | 1 | 12 |
| BPDR | ThinkSystem V4 2U Standard Heatsink | 2 | 24 |
| C1EY | ThinkSystem 8GPU Server CFFV4 Power Interface Board | 1 | 12 |
| C1FH | ThinkSystem SR685a V3 / SR680a V3 Front Fan | 5 | 60 |
| C1F8 | ThinkSystem 8GPU Server Power Distribution Board | 1 | 12 |
| 5641PX3 | XClarity Pro, Per Endpoint w/3 Yr SW S&S | 1 | 12 |
| 1340 | Lenovo XClarity Pro, Per Managed Endpoint w/3 Yr SW S&S | 1 | 12 |
| 7Q01CTSAWW | SERVER KEEP YOUR DRIVE ADD-ON | 1 | 12 |
| 7Q01CTS4WW | SERVER PREMIER 24X7 4HR RESP | 1 | 12 |

**ThinkSystem SR780a V3**

Table 14. ThinkSystem SR780a V3 BoM

| Part Number | Product Description | Qty per System | Total Qty |
|---|---|---|---|
| 7DJ5CTO1WW | Server : ThinkSystem SR780a V3 - 3 Year Warranty with Controlled GPU | 1 | 12 |
| C2F0 | ThinkSystem SR780a V3 H200 HGX GPU Base | 1 | 12 |
| BYWG | Intel Xeon Platinum 8570 56C 350W 2.1GHz Processor | 2 | 24 |
| C5H9 | ThinkSystem 64GB TruDDR5 5600MHz (2Rx4) 10x4 RDIMM | 32 | 384 |
| BTQ1 | ThinkSystem 2.5" U.3 PM1743 7.68TB Read Intensive NVMe PCIe 5.0 x4 HS SSD | 2 | 24 |
| C2ET | ThinkSystem 8x2.5" PCIe Gen5 AnyBay Storage Backplane | 1 | 12 |
| BXMH | ThinkSystem M.2 PM9A3 960GB Read Intensive NVMe PCIe 4.0 x4 NHS SSD | 2 | 24 |
| C0Q4 | ThinkSystem NVIDIA BlueField-3 B3140H VPI QSFP112 1P 400G PCIe Gen5 x16 Adapter | 8 | 96 |
| BVBG | ThinkSystem NVIDIA BlueField-3 B3220 VPI QSFP112 2P 200G PCIe Gen5 x16 Adapter | 1 | 12 |
| C2ER | ThinkSystem NVIDIA HGX H200 141GB 700W 8-GPU Liquid Cooled Board | 1 | 12 |
| C2EX | ThinkSystem SR780a V3 x16 PCIe Rear IO Riser | 2 | 24 |
| C2EU | ThinkSystem SR780a V3 8x x16 PCIe GPU Switch Riser | 1 | 12 |
| C4HK | ThinkSystem 2600W 230V Titanium Hot-Swap Gen2 Power Supply v4 | 8 | 96 |
| 6252 | 2.5m, 16A/100-250V, C19 to C20 Jumper Cord | 8 | 96 |
| AUMY | ThinkSystem Lift Handles | 1 | 12 |
| C2P8 | ThinkSystem SR780a V3 Rail Kit | 1 | 12 |
| C2EZ | ThinkSystem SR780a V3 System Board | 1 | 21 |
| BAVU | Front Operator Panel with Quad Line LCD Display | 1 | 12 |
| C2RF | ThinkSystem Rear Fans | 5 | 60 |
| C1EY | ThinkSystem 8GPU Server CFFV4 Power Interface Board | 1 | 12 |
| C2RE | ThinkSystem Internal Fans | 6 | 72 |
| C2EW | ThinkSystem SR780a V3 Fan Control Boards | 1 | 12 |
| C4E4 | ThinkSystem SR780a V3 Two Rear Riser Assembly | 1 | 12 |
| C1F8 | ThinkSystem 8GPU Server Power Distribution Board | 1 | 12 |
| 5641PX3 | XClarity Pro, Per Endpoint w/3 Yr SW S&S | 1 | 12 |
| 1340 | Lenovo XClarity Pro, Per Managed Endpoint w/3 Yr SW S&S | 1 | 12 |
| 7Q01CTS4WW | SERVER PREMIER 24X7 4HR RESP | 1 | 12 |
| 7Q01CTSAWW | SERVER KEEP YOUR DRIVE ADD-ON | 1 | 12 |

**ThinkSystem SR780a V3 B200 BoM**

Table 15. ThinkSystem SR780a V3 B200 BoM

| Part Number | Product Description | Qty per System | Total Qty |
|---|---|---|---|
| 7DJ5CTO1WW | Server : ThinkSystem SR780a V3 - 3 Year Warranty for AI | 1 | 12 |
| C695 | ThinkSystem SR780a V3 B200 HDX GPU Base | 1 | 12 |
| BYWG | Intel Xeon Platinum 8570 56C 350W 2.1GHz Processor | 2 | 24 |
| C5H9 | ThinkSystem 64GB TruDDR5 5600MHz (2Rx4) 10x4 RDIMM | 32 | 384 |
| BTQ1 | ThinkSystem 2.5" U.3 PM1743 7.68TB Read Intensive NVMe PCIe 5.0 x4 HS SSD | 2 | 24 |
| C2ET | ThinkSystem 8x2.5" PCIe Gen5 AnyBay Storage Backplane | 1 | 12 |
| BXMH | ThinkSystem M.2 PM9A3 960GB Read Intensive NVMe PCIe 4.0 x4 NHS SSD | 2 | 24 |
| C0Q4 | ThinkSystem NVIDIA BlueField-3 B3140H VPI QSFP112 1P 400G PCIe Gen5 x16 Adapter | 8 | 96 |
| BVBG | ThinkSystem NVIDIA BlueField-3 B3220 VPI QSFP112 2P 200G PCIe Gen5 x16 Adapter | 1 | 12 |
| C696 | ThinkSystem NVIDIA HGX B200 1000W 180GB 8-GPU Liquid-Cooled Board | 1 | 12 |
| C2EX | ThinkSystem SR780a V3 x16 PCIe Rear IO Riser | 1 | 12 |
| C2EU | ThinkSystem SR780a V3 8x x16 PCIe GPU Switch Riser | 1 | 12 |
| C4HK | ThinkSystem 2600W 230V Titanium Hot-Swap Gen2 Power Supply v4 | 8 | 96 |
| 6252 | 2.5m, 16A/100-250V, C19 to C20 Jumper Cord | 8 | 96 |
| AUMY | ThinkSystem Lift Handles | 1 | 12 |
| C2P8 | ThinkSystem SR780a V3 Rail Kit | 1 | 12 |
| C2EZ | ThinkSystem SR780a V3 System Board | 1 | 12 |
| BAVU | Front Operator Panel with Quad Line LCD Display | 1 | 12 |
| AVEN | ThinkSystem 1x1 2.5" HDD Filler | 6 | 72 |
| C517 | ThinkSystem SR680a V3 Server Power Distribution Board 2 | 1 | 12 |
| C2RF | ThinkSystem Rear Fans | 5 | 60 |
| C1EY | ThinkSystem 8GPU Server CFFV4 Power Interface Board | 1 | 12 |
| C2RE | ThinkSystem Internal Fans | 6 | 72 |
| C2EW | ThinkSystem SR780a V3 Fan Control Boards | 1 | 12 |
| C4E4 | ThinkSystem SR780a V3 Two Rear Riser Assembly | 1 | 12 |
| 5641PX3 | XClarity Pro, Per Endpoint w/3 Yr SW S&S | 1 | 12 |
| 1340 | Lenovo XClarity Pro, Per Managed Endpoint w/3 Yr SW S&S | 1 | 12 |

## ThinkSystem SR655 V3 BoM

Table 16. ThinkSystem SR655 V3 BoM

| Part Number | Product Description | Qty per System | Total Qty |
|---|---|---|---|
| 7D9ECTOLWW | ThinkSystem SR655 V3 | | 5 |
| BLKK | ThinkSystem V3 2U 24x2.5" Chassis | 1 | 5 |
| C2AC | ThinkSystem SR655 V3 MB w/IO+PIB+FB,2U | 1 | 5 |
| C2AQ | ThinkSystem AMD EPYC 9335 32C 210W 3.0GHz Processor | 1 | 5 |
| C0CJ | ThinkSystem 32GB TruDDR5 6400MHz (2Rx8) RDIMM-A | 12 | 60 |
| BPQV | ThinkSystem V3 2U x16/x16/E PCIe Gen5 Riser1 or 2 | 1 | 5 |
| BVBG | ThinkSystem NVIDIA BlueField-3 B3220 VPI QSFP112 2P 200G PCIe Gen5 x16 Adapter | 1 | 5 |
| B8P9 | ThinkSystem M.2 NVMe 2-Bay RAID Adapter | 1 | 5 |
| BXMH | ThinkSystem M.2 PM9A3 960GB Read Intensive NVMe PCIe 4.0 x4 NHS SSD | 2 | 10 |
| BS7Y | ThinkSystem V3 2U 8x2.5" NVMe Gen5 Backplane | 1 | 5 |
| C0ZU | ThinkSystem 2.5" U.2 Multi Vendor 3.84TB Read Intensive NVMe PCIe 5.0 x4 HS SSD | 2 | 10 |
| BLL6 | ThinkSystem 2U V3 Performance Fan Module | 6 | 30 |
| 6400 | 2.8m, 13A/100-250V, C13 to C14 Jumper Cord | 2 | 10 |
| BLKH | ThinkSystem 1100W 230V Titanium Hot-Swap Gen2 Power Supply | 2 | 10 |
| B8LA | ThinkSystem Toolless Slide Rail Kit v2 | 1 | 5 |
| C1PT | ThinkSystem SR635 V3/SR655 V3 Root of Trust Module Low Voltage-RoW V2 | 1 | 5 |
| BQQ6 | ThinkSystem 2U V3 EIA right with FIO | 1 | 5 |
| 5PS7B08762 | 5Yr Premier NBD Resp + KYD SR655 V3 | 1 | 5 |

## NVIDIA SN5610 Switch BoM

Table 17. NVIDIA SN5610 Switch BoM

| Part Number | Product Description | Qty per System | Total Qty |
|---|---|---|---|
| NVIDIA SKU: 920-9N42F-00RI-3C1 | NVIDIA SN5610 | 1 | 2 |

## NVIDIA SN2201 Switch BoM

Table 18. NVIDIA SN2201 Switch BoM

| Part Number | Description | Qty per System | Total Qty |
|---|---|---|---|
| 7D5FCTOFWW | NVIDIA SN2201 1GbE Managed Switch with Cumulus (PSE) | | 2 |
| BPC7 | NVIDIA SN2201 1GbE Managed Switch with Cumulus (PSE) | 1 | 2 |
| 6201 | 1.5m, 10A/100-250V, C13 to C14 Jumper Cord | 2 | 4 |
| 5WS7B98268 | 5Yr Premier NBD Resp NVID SN2201 PSE | 1 | 2 |

## Power Distribution Unit (PDU) BoM

Table 19. Power Distribution Unit (PDU) BoM

| Part Number | Product Description | Qty per System | Total Qty |
|---|---|---|---|
| 7DGMCTO1WW | -SB- 0U 18 C13/C15 and 18 C13/C15/C19 Switched and Monitored 63A 3 Phase WYE PDU v2 | | 2 |

## Rack Cabinet BoM

Table 20. Rack Cabinet BoM

| Part Number | Product Description | Qty per System | Total Qty |
|---|---|---|---|
| 1410O42 | Lenovo EveryScale 42U Onyx Heavy Duty Rack Cabinet | | 1 |
| BHC4 | Lenovo EveryScale 42U Onyx Heavy Duty Rack Cabinet | 1 | 1 |
| BJPD | 21U Front Cable Management Bracket | 2 | 2 |
| BHC7 | ThinkSystem 42U Onyx Heavy Duty Rack Side Panel | 2 | 2 |
| BJPA | ThinkSystem 42U Onyx Heavy Duty Rack Rear Door | 1 | 1 |
| 5AS7B07693 | Lenovo EveryScale Rack Setup Services | 1 | 1 |

## XClarity Software BoM

Table 21. XClarity Software BoM

| Part Number | Product Description | Qty per System | Total Qty |
|---|---|---|---|
| SBCV | Lenovo XClarity XCC2 Platinum Upgrade (FOD) | | 3 |
| 00MT203 | Lenovo XClarity Pro, Per Managed Endpoint w/5 Yr SW S&S | | 5 |

## Seller training courses

The following sales training courses are offered for employees and partners (login required). Courses are listed in date order.

1. **VTT AI: Introducing The Lenovo Hybrid AI 285 Platform with Cisco Networking**
   2025-08-26 | 54 minutes | Employees Only

   Please view this session as Pierce Beary, Sr. AI Solution Manager, ISG ESMB Segment and AI explains:

   - Value propositions for the Hybrid AI 285 platform
   - Updates for the Hybrid AI 285 platform
   - Leveraging Cisco networking with the 285 platform
   - Future plans for the 285 platform

   Tags: Artificial Intelligence (AI), Technical Sales, ThinkSystem

   Published: 2025-08-26
   Length: 54 minutes

   > **Start the training:**
   > Employee link: Grow@Lenovo

   Course code: DVAI220

2. **VTT AI: Introducing the Lenovo Hybrid AI 285 Platform April 2025**
   2025-04-30 | 60 minutes | Employees Only

   The Lenovo Hybrid AI 285 Platform enables enterprises of all sizes to quickly deploy AI infrastructures supporting use cases as either new greenfield environments or as an extension to current infrastructures. The 285 Platform enables the use of the NVIDIA AI Enterprise software stack. The AI Hybrid 285 platform is the perfect foundation supporting Lenovo Validated Designs.
   • Technical overview of the Hybrid AI 285 platform
   • AI Hybrid platforms as infrastructure frameworks for LVDs addressing data center-based AI solutions.
   • Accelerate AI adoption and reduce deployment risks

   Tags: Artificial Intelligence (AI), Nvidia, Technical Sales, Lenovo Hybrid AI 285

   Published: 2025-04-30
   Length: 60 minutes

   > **Start the training:**
   > Employee link: Grow@Lenovo

   Course code: DVAI215

## Related publications and links

For more information, see these resources:

- Lenovo EveryScale support page:
  https://datacentersupport.lenovo.com/us/en/solutions/ht505184
- x-config configurator:
  https://lesc.lenovo.com/products/hardware/configurator/worldwide/bhui/asit/x-config.jnlp
- Implementing AI Workloads using NVIDIA GPUs on ThinkSystem Servers:
  https://lenovopress.lenovo.com/lp1928-implementing-ai-workloads-using-nvidia-gpus-on-thinksystem-servers
- Making LLMs Work for Enterprise Part 3: GPT Fine-Tuning for RAG:
  https://lenovopress.lenovo.com/lp1955-making-llms-work-for-enterprise-part-3-gpt-fine-tuning-for-rag
- Lenovo to Deliver Enterprise AI Compute for NetApp AIPod Through Collaboration with NetApp and NVIDIA
  https://lenovopress.lenovo.com/lp1962-lenovo-to-deliver-enterprise-ai-compute-for-netapp-aipod-nvidia

## Related product families

Product families related to this document are the following:

- AI Servers
- ThinkSystem SR655 V3 Server
- ThinkSystem SR680a V3 Server
- ThinkSystem SR685a V3 Server
- ThinkSystem SR780a V3 Server

# Notices

Lenovo may not offer the products, services, or features discussed in this document in all countries. Consult your local Lenovo representative for information on the products and services currently available in your area. Any reference to a Lenovo product, program, or service is not intended to state or imply that only that Lenovo product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe any Lenovo intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any other product, program, or service. Lenovo may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not give you any license to these patents. You can send license inquiries, in writing, to:

Lenovo (United States), Inc.
8001 Development Drive
Morrisville, NC 27560
U.S.A.
Attention: Lenovo Director of Licensing

LENOVO PROVIDES THIS PUBLICATION "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. Some jurisdictions do not allow disclaimer of express or implied warranties in certain transactions, therefore, this statement may not apply to you.

This information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. Lenovo may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice.

The products described in this document are not intended for use in implantation or other life support applications where malfunction may result in injury or death to persons. The information contained in this document does not affect or change Lenovo product specifications or warranties. Nothing in this document shall operate as an express or implied license or indemnity under the intellectual property rights of Lenovo or third parties. All information contained in this document was obtained in specific environments and is presented as an illustration. The result obtained in other operating environments may vary. Lenovo may use or distribute any of the information you supply in any way it believes appropriate without incurring any obligation to you.

Any references in this publication to non-Lenovo Web sites are provided for convenience only and do not in any manner serve as an endorsement of those Web sites. The materials at those Web sites are not part of the materials for this Lenovo product, and use of those Web sites is at your own risk. Any performance data contained herein was determined in a controlled environment. Therefore, the result obtained in other operating environments may vary significantly. Some measurements may have been made on development-level systems and there is no guarantee that these measurements will be the same on generally available systems. Furthermore, some measurements may have been estimated through extrapolation. Actual results may vary. Users of this document should verify the applicable data for their specific environment.

This document, LP2286, was created or updated on August 28, 2025.

Send us your comments in one of the following ways:

- Use the online Contact us review form found at:
  https://lenovopress.lenovo.com/LP2286
- Send your comments in an e-mail to:
  comments@lenovopress.com

This document is available online at  https://lenovopress.lenovo.com/LP2286.

## Trademarks

Lenovo and the Lenovo logo are trademarks or registered trademarks of Lenovo in the United States, other countries, or both. A current list of Lenovo trademarks is available on the Web at https://www.lenovo.com/us/en/legal/copytrade/.

The following terms are trademarks of Lenovo in the United States, other countries, or both:
Lenovo®
AnyBay®
From Exascale to EveryScale®
Lenovo Hybrid AI Advantage
ThinkAgile®
ThinkSystem®
XClarity®

The following terms are trademarks of other companies:

AMD and AMD EPYC™ are trademarks of Advanced Micro Devices, Inc.

Intel® and Xeon® are trademarks of Intel Corporation or its subsidiaries.

Linux® is the trademark of Linus Torvalds in the U.S. and other countries.

IBM® is a trademark of IBM in the United States, other countries, or both.

Other company, product, or service names may be trademarks or service marks of others.