# Cisco Nexus Hyperfabric AI

## Contents

# Overview and definitions

**Q: What is Cisco [Nexus® Hyperfabric](#) AI?**

A: Cisco Nexus Hyperfabric AI is a cloud-managed full stack AI infrastructure solution delivered as hardware + software + service. Using a cloud controller managed by Cisco, customers easily design, deploy, and manage their network fabric, GPU servers, and storage servers. It reinvents the IT operations lifecycle of the AI infrastructure in the data center by simplifying every step of the process and ensures repeatable and predictable outcomes by IT generalist, application, and DevOps teams. The vertical stack solution consists of purpose-built hardware, software, cloud management, Day-2 automation, and Cisco support. Cisco Nexus Hyperfabric AI is best suited for customers looking to build out their private cloud AI infrastructure.

**Q: Why are we bringing Cisco Nexus Hyperfabric AI to market in 2025?**

A: The adoption of AI is increasingly viewed by organizations as a key enabler to drive innovation: "By 2027, 40% of enterprises will deploy GenAI network fabrics to enable cost and performance-optimized support for AI workloads in their own data centers."

(IDC Perspective, March 2024.) Our findings show us that 95% of customers are aware AI will increase workloads, but only 17% are equipped to handle this complexity, with 23% having limited or no capacity to meet the AI demand with current infrastructures. Given the significant business value and growth in AI adoption across enterprises, Cisco is working closely with NVIDIA and VAST Data to ensure that customers can rapidly and reliably deploy AI wherever it is needed. This is in addition to simplifying the IT lifecycle enabling IT generalists, data science teams, and DevOps teams to easily design, deploy, and operate AI and non-AI data center fabrics. The solution is based on a converged Ethernet network, so organizations easily support and scale by leveraging existing skills and processes.

**Q: What are the components of Cisco Nexus Hyperfabric AI?**

A: Cisco Nexus Hyperfabric AI consists of five primary components:

- **Cloud Controller:** a scalable, globally distributed multitenant cloud service that is used to design, plan, control, upgrade, and monitor fabrics using a browser or APIs known as Cisco Nexus Hyperfabric.

- **Cisco 6000 Series Switches:** installed with Cisco Nexus Hyperfabric AI–managed software, they connect to the cloud for centralized real-time visibility and control.

- **Cisco UCS® NVIDIA GPU/DPU Servers:** Cisco UCS-C885A-M8-CN1 server that packs 8 NVIDIA H200 GPUs, with BlueField-3 DPU (3240H) and SuperNIC (3140H) connected using Cisco® optics that can efficiently run training/inferencing/fine-tuning jobs.

- **UCS-VAST Storage (Optional):** Cisco UCS-C225-M8N-1P servers pack 8 drives with each of the servers with a cluster consisting of 11 servers that can be easily expanded. The specs start with 1PB of storage that can be used for training, fine-tuning, inferencing, Retrieval-Augmented Generation (RAG), and other data engineering work.

- **NVAIE:** Integration and access to the NVIDIA AI for Enterprise (NVAIE) software stack from get-go allows access to different model training software and data catalogues so that engineers and scientists can get started immediately.

Cisco Nexus Hyperfabric AI is a single integrated solution **from Cisco** that can be designed, ordered, deployed, and operated as a cohesive solution.

**Q: What are the capabilities of Cisco Nexus Hyperfabric AI?**

A: Cisco Nexus Hyperfabric AI integrates these components into a single solution:

- Its cloud controller, operated by Cisco, serves as a single point for configuration, monitoring, and maintenance of all tenant customer fabrics, and provides health monitoring and visibility of GPU servers and storage servers, using real-time connection to switches and servers deployed either on-premises or in colocation facilities.

- Network fabrics consist of cloud-managed Cisco 6000 Series Switches that offer automated zero-touch provisioning, and a "Helping Hands assistant," which provides step-by-step cabling instructions combined with real-time verification.

- Network fabric spans the GPU servers as well as storage servers and can be set as a mesh or leaf-spine configuration. With RDMA over Converged Ethernet (RoCE) Version 2, your investment in network hardware and expertise is protected.

- The GPU cluster consisting of UCS-C885A-M8-CN1 servers is highly scalable and flexible so that you can adapt the cluster size to your needs.

- Optional storage cluster, consisting of Cisco UCS-C225-M8N-1P high-density storage servers built on the latest NVMe drives ensure fast Input/Output Operations Per Second (IOPS), that feeds data to GPUs using Cisco optics so that you get most value from your GPU investments.

- In a shared-responsibility model, automation and operations from Cisco support manage the cloud controller, the fabric underlay and overlay networks, and the software upgrade process, while customers maintain direct control of all interconnections to their applications, hosts, and the rest of their network.

- Assertion-based monitoring continuously verifies the availability and reliability of the fabric and connected resources, and making it easier to detect the root cause of any issue. A built-in designer helps customers construct a validated fabric design based on desired compute, storage, host and port capacity, oversubscription, and environmental considerations including cabling and power, and then creates an accurate Bill of Materials (BoM).

- Self-service fabric tenancies empower host and application teams to monitor and manage the fabric services they have been allocated, removing the need to depend on IT for most support services.

**Q: Which partnerships are part of the Cisco Nexus Hyperfabric AI solution?**

A: The Nexus Hyperfabric AI cluster solution leverages the following partners:

1. NVIDIA AI for Enterprise (NVAIE).

2. NVIDIA for GPUs and network interface cards (NICs) based on Hyperscale Graphics eXtension (HGX) reference architectures.

3. VAST storage software on Cisco UCS storage servers.

**Q: What is Cisco Hyperfabric AI Certification with NVIDIA Enterprise Reference Architecture?**

A: The [Cisco Nexus Hyperfabric AI infrastructure](#) leverages the NVIDIA AI Enterprise reference architecture to establish best practices and guidelines for building high-performance, scalable, and secure infrastructures in enterprise environments. This architecture integrates processors, servers, networking devices, storage, management tools, and software into an NVIDIA-validated design. It provides customers with a vertically integrated solution, ensuring optimized performance and seamless scalability.

**Q: How will Cisco Nexus Hyperfabric AI with NVIDIA ERA be managed?**

A: The primary management tool for ERA-certified Cisco Nexus Hyperfabric AI solutions is the Hyperfabric Controller.

- Hyperfabric Controller: provides assertion-based, high-level management and control across all three layers: network, GPUs, and storage.

- Cisco Intersight®:  provides detailed management of GPU servers and storage servers.

- NVIDIA AI Enterprise software: provides AI/ML tools tailored for data scientists.

- VAST VMS: performs data- and storage-management tasks.

These management capabilities will include comprehensive fabric control and leverage NVIDIA's proprietary insights, directly integrated into Cisco Silicon One®, to enhance visibility into workload connectivity. The collaboration and integration between Cisco and NVIDIA products are expected to deepen further over time.

**Q: What switches does Cisco Nexus Hyperfabric AI support?**

A. Three models of Cisco 6000 switches will be offered with the first release of Cisco Nexus Hyperfabric AI. Two are Cisco Silicon One Q200 Processor 1RU platforms, and one is Cisco Silicon One G200 processor.

- Cisco HF6100-60L4D with 60x 10/25/50GbE SFP56, 4x 100/400GbE QSFP56-DD (16x through 100GbE breakout)

- Cisco HF6100-32D with 32x 100/400GbE QSFP56-DD (128x 100GbE through 400:100 breakout)

- Cisco HF6100-64E with 64x 800G OSFP

Additional [Cisco 6000](#) Series Switches will be offered in the future.

**Q: Why is As a Service so important? The hardware is all on premises.**

A: AI infrastructure is very expensive. Acquisition costs are high, as are operational costs for space, power, and cooling, so customers need to see a fast return on their investment. This means many enterprise customers need to get their AI jobs running as fast as possible. Bringing up a new AI cluster

will require the coordination of MLOps, NetOps, IT, facilities, and SecOps skills that might be in short supply after years of moving workloads to the cloud. Working with a trusted partner such as Cisco will enable customers to act quickly and with greater confidence in this rapidly evolving area.

**Q: Is Cisco Nexus Hyperfabric AI a closed system? How does it integrate with existing data centers?**

A. Cisco Nexus Hyperfabric AI is a complete, turnkey solution designed for ease of deployment and accelerated day-2 value. It is standards-based and will interoperate with existing data-center fabrics connected, for example, through border gateways. As a turnkey solution, Cisco Nexus Hyperfabric AI will run only on specific Cisco hardware, and only that hardware can be part of a Nexus Hyperfabric AI cluster.

**Q: What about other storage partners? I don't use VAST.**

A. VAST storage is our partner at launch. Customers are welcome to use their own storage vendor through the north/south (N/S) interface into the cluster; however, it will not be managed by Hyperfabric.

# Deployment guidelines

**Q: Is the Cisco 6000 Series Switch the same as the old Cisco Nexus 6000 switch family?**

A. They are not the same. The old Cisco Nexus 6000 switches have reached end-of-life and end-of-support and are not supported or compatible with Cisco Nexus Hyperfabric AI.

**Q: Can Nexus Hyperfabric manage Cisco Nexus 9000 Series Switches?**

A. The initial release of Cisco Nexus Hyperfabric does not support Cisco Nexus 9000 Series Switches. The new Cisco 6000 switches are required. Support for Cisco Nexus 9000 switches is being evaluated.

**Q: At a high level, how does Cisco Nexus Hyperfabric AI work?**

A. Customers begin by logging into the Nexus Hyperfabric AI designer, where they can start designing their AI infrastructure by selecting a pre-validated template that suits their requirements. Once these specifications are entered, the designer automatically creates a validated, complete fabric topology, including the necessary UCS compute servers, storage servers, optics, cabling, and switches. Nexus Hyperfabric AI seamlessly integrates with Cisco's ordering tools to ensure there are no

errors when converting the design into a bill of materials. The bill of materials also includes Cisco's Day 0/1 installation and deployment services, as well as software subscription licenses for Cisco Hyperfabric, Cisco Intersight, VAST storage software, and NVIDIA NVAIE. When the Cisco 6000 switches, Cisco UCS 885 compute servers, and Cisco UCS 225 storage servers arrive on site, the Cisco Customer Experience team or partners will handle the rack and stack (Day 0/1) to complete the installation and deployment process. Once deployed, the switches automatically connect to the cloud for claiming and provisioning by the cloud controller through zero-touch plug-and-play. This process quickly sets up a fully operational AI infrastructure, including the network fabric, in just minutes. Assertion-based monitoring continuously checks the availability and reliability of the AI infrastructure and its connected resources, immediately identifying the root cause of any issues. If changes are needed to adjust the capacity or design, customers can easily modify the in-progress design, approve the updates, and follow the same process again. The product offers guidance on all necessary physical changes to transition from the old topology to the new, including cabling adjustments, and automatically reconfigures itself.

**Q: How can I connect multiple fabrics with Cisco Nexus Hyperfabric AI?**

A: Cisco Nexus Hyperfabric AI is designed to interoperate with existing Ethernet fabrics and Layer-3 networks, thereby integrating north-south traffic with east-west traffic. It delivers support for:

1. Northbound-routed peering through Border Gateway Protocol (BGP) and static routes for both IPv4 and IPv6 address families. Nexus Hyperfabric supports multi-Virtual Routing and Forwarding (VRF) routing for concurrently peering multiple virtual-routing instances with upstream routed networks.

2. East/westbound peering with existing Layer-2 Ethernet-based networks and fabrics using multichassis links.

**Q: What deployment use cases does the first release support?**

A: The initial release is designed to support the following use cases:

1. Initial deployment of Gen AI infrastructure, Small Language Model (SLM), fine-tuning, and RAG-type deployments are supported Day 1.

2. Limited Gen AI infrastructure expertise, i.e. both depth and breadth, to build an infrastructure expertise to support your Gen AI workloads.

3. Where capacity needs are unknown upfront, scale as you know more, rather than overbuilding upfront.

4. Gen AI deployment that needs more than 8 GPUs and, optionally, 1PB storage capacity.

5. Share the infrastructure among multiple groups for training, fine-tuning, RAG, and data engineering works.

## Q: Where is the Nexus Hyperfabric Cloud Controller hosted?

A: The Nexus Hyperfabric Cloud Controller is maintained by Cisco and is hosted in the public cloud. The Cloud Controller includes global scalability, as well as multi-region redundancy, without any additional configuration for end users.

The Nexus Hyperfabric Cloud Controller is reachable through a single URL ([https://nexushyperfabric.cisco.com](https://nexushyperfabric.cisco.com)) and is used for all communications (including switch-management cloud connectivity, primary user interfaces, and RESTful API endpoints).

## Q: What are the plans for Cloud Controller hosting outside the United States?

A: Cisco plans to offer local cloud controllers hosted natively in Europe and Asia in future releases.

## Q: Is Cisco Nexus Hyperfabric strictly cloud-based, or are there any plans for air-gapped environments?

A: The Cloud Controller is managed by Cisco operations staff. There are no plans to hand off that responsibility to entities not owned or operated by Cisco.

## Q: What are the plans for government cloud compliance (under the Federal Risk and Authorization Management Program [FedRAMP])?

A: Cisco is investigating offering local cloud controllers hosted in FedRAMP in a future release and will consider other government cloud-compliance environments.

## Q: Will Nexus Hyperfabric AI have integrations with other products?

A: In the first release, HashiCorp Terraform and Red Hat Ansible could be used to fully automate provisioning, and telemetry may be sent to a variety of data collectors, including Amazon S3, Cisco Splunk, and ServiceNow. The Nexus Hyperfabric controller is designed to be operated API-first, and includes a native RESTful API designed to allow flexible integrations both for provisioning tools such as Ansible or Terraform and for external network- and incident-management tools.

## Q: How can I export logging to my own logging platforms?

A: The Nexus Hyperfabric Cloud Controller will support external cloud-to-cloud integrations for delivering logging information through external Amazon S3 buckets to logging files in syslog format. In addition, the Cloud Controller will also allow API integration for querying internal states, alerts, and other telemetry data from the on-premises fabrics.

## Q: What are the bandwidth requirements for cloud connectivity?

A: Target bandwidth per switch is less than 2Mb/sec at steady state. Certain types of real-time monitoring may increase upstream bandwidth when initiated by end users through the Cloud Controller UI.

## Q: How is real-time telemetry uplifted to the Cloud Controller?

A: The Cisco 6000 Series Switches include a telemetry agent that establishes an outbound TLS connection to the cloud controller. This outbound Transport Layer Security (TLS) (TCP/443) session is designed to work with existing network security controls for standard web clients and provides all connectivity required for configuration, monitoring, and real-time telemetry for each switch managed by Nexus Hyperfabric.

**Q: What happens if I lose cloud connectivity?**

A: The cloud connection from a deployed fabric to the Cloud Controller is used for management and telemetry purposes only. All stateful protocols are maintained independently on the switches residing on premises. Any fabric disconnected from the cloud will continue to operate normally (including local underlay and overlay fabric-management protocols, external peering protocols such as BGP, and all standard bridging functions).

When disconnected from the cloud, a fabric will not be able to accept configuration updates, and telemetry will not be received by the cloud controller. The on-switch agent will intelligently buffer telemetry if disconnected and will upload this information and resynchronize any configuration updates once the fabric is reconnected to the cloud controller.

**Q: How can I interact with the Nexus Hyperfabric API?**

A: The Nexus Hyperfabric Cloud Controller implements a RESTful API using the same URL as standard management functions do. This API is extensively documented in an online API programming guide and is designed to be compatible with OpenAPI specifications.

In addition, Nexus Hyperfabric will support providers for both Ansible and Terraform to help users exercise the Nexus Hyperfabric API according to their existing management tooling.

**Q: Who will a customer contact for Technical Assistance Center (TAC) support for the cluster?**

A: Cisco.

**Q: What are the power requirements?**

A: In the range of 10-16kW/server per GPU Server. Storage servers, optics. and switches will require additional power.

**Q: How are the components of Hyperfabric AI Connected?**

A: AI clusters are generally built around multiple logical network fabrics. The primary fabrics are distinct in their optimization based on which function they provide. Logically, there are multiple network fabrics that form the AI cluster network:

- A backend network interconnecting GPUs across the servers.

- A frontend network connecting the servers to the rest of the applications and end users.

- A storage network connecting the GPU servers to storage servers.

- A management network connecting to each switch and server in the network and used by Cisco Hyperfabric and Cisco Intersight consoles.

These networks make up a single cluster, controlled by a single instance of Hyperfabric AI.

**Q: Is the cluster air cooled?**

A: Yes. Future versions may require liquid cooling for large clusters and faster GPUs.

## How to purchase

**Q: When will Cisco Nexus Hyperfabric AI be available?**

A: Cisco Nexus Hyperfabric AI is planned to be released in Q3, 2025. Many of the components of the product such as Cisco Nexus Hyperfabric and Cisco UCS compute servers are orderable now.

**Q: How do I buy Cisco Nexus Hyperfabric AI?**

A: Cisco Nexus Hyperfabric AI may be purchased from a certified Cisco Nexus Hyperfabric AI reseller. Organizations able to purchase products directly from Cisco may also purchase the solution. In addition you can design your data center fabric before buying it.

**Q: How do I design a data center fabric before buying it?**

A: Anyone with a Cisco Connection Online (CCO) account or with an identity provider (IdP) tied to CCO may request and receive an organization tenancy in the cloud controller. Within the tenancy, customers may design blueprints that contain all the details needed to order, deploy, configure, and operate the fabric, including:

1. Physical components including switches, optics, compute servers, storage servers, and connectors, air-flow direction, power consumption, and a cabling plan.

2. A bill of materials for all the Cisco and partner components that is the source-of-truth and integrated with Cisco Commerce Workspace (CCW) to automate the process of obtaining an accurate quote.

3. The logical network including the entire network overlay and underlay and upstream route peering.

4. Host-port assignments to server infrastructure teams.

5. API integration to automate the provisioning of the blueprint.

Additional accounts may then be added to that tenancy so teams can collaborate. Once the equipment is deployed, it is automatically provisioned according to the blueprint.

**Q: Can Cisco Nexus Hyperfabric AI be sold only by certified resellers?**

A: Yes, only Cisco Nexus Hyperfabric AI certified resellers may sell the product.

**Q: What is the Cisco Secure AI Factory with NVIDIA?**

A: Cisco Secure AI Factory with NVIDIA represents a scalable, high-performance, and secure AI infrastructure solution offered by Cisco in collaboration with NVIDIA and their strategic partners. This solution aims to accelerate the adoption of AI for enterprises by providing a security-first architecture that integrates Cisco security solutions such as Cisco AI Defense and Cisco Hypershield. It leverages high-performance AI infrastructure, including Cisco UCS AI servers and Cisco switches, to enable efficient model training, customization, and inferencing. It offers a pre-validated AI full-stack solution, Cisco's Nexus Hyperfabric AI, with flexible deployment options to simplify the deployment and implementation of AI Infrastructure.

**Q: Where can I get more information?**

A: More information may be found at https://www.cisco.com/site/us/en/products/networking/data-center-networking/nexus-hyperfabric/hyperfabric-ai/index.html.