

TECHNICAL WHITE PAPER  
March 2025



# VMware Private AI Foundation with NVIDIA on HGX Servers

Reference Design for Inference

**vmware**<sup>®</sup>  
by **Broadcom**



## Table of contents

<b>Document history</b> .....	<b>4</b>
<b>Executive summary</b> .....	<b>5</b>
<b>Introduction</b> .....	<b>6</b>
Intended audience .....	7
Terminology .....	7
<b>Core components</b> .....	<b>8</b>
NVIDIA AI Enterprise .....	9
NVIDIA vGPU (C-Series) .....	9
NGC .....	10
VMware Private AI Foundation with NVIDIA.....	12
VMware Cloud Foundation .....	15
HGX systems .....	17
Ethernet networking.....	18
<b>Reference architecture</b> .....	<b>19</b>
Physical architecture .....	19
Virtual architecture.....	23
Management domain .....	24
Management domain architecture .....	24
Management domain settings .....	25
Workload domain .....	28
Workload domain architecture .....	28
Workload domain settings .....	29
VMware Private AI Foundation with NVIDIA prerequisites .....	31
<b>Validation</b> .....	<b>32</b>
<b>Performance</b> .....	<b>33</b>
Benchmarking GenAI-Perf .....	33
Deploy a DLVM .....	33
Launch NIM in a DLVM.....	34
Launch GenAI-Perf .....	36
Performance comparison of virtual vs. bare metal .....	38
Inference sizing guidance.....	41

Conclusion ..... 43

Additional information ..... 43

About the authors ..... 43

Acknowledgments ..... 43

## Document history

Version	Date	Authors	Change Summary
1.0	2025-03-17	Yuankun Fu, Agustin Malanco, Ramesh Radhakrishnan	Initial release based on VCF 5.2.1

## Executive summary

As AI continues to revolutionize industries, organizations—driven by concerns over cost efficiency, security, and agility—are increasingly adopting private cloud solutions to power inference workloads. VMware® Private AI Foundation with NVIDIA is a generative AI (GenAI) platform that enables AI professionals to run RAG workflows, fine-tune and customize LLM models, and run inference workloads in their on-prem data centers. The platform addresses critical issues related to privacy, choice, cost, performance, and compliance. VMware Private AI Foundation with NVIDIA comprises VMware Cloud Foundation® and NVIDIA AI Enterprise (featuring NVIDIA vGPU, NVIDIA NIM™ and NVIDIA NeMo™ microservices, and NVIDIA AI Blueprints).

This reference design details our recommended architecture. With the products and guidelines listed here, IT teams can easily deploy a robust, future-proof infrastructure from which data scientists can easily deploy AI inference applications. The foundation of this solution is made up of NVIDIA-certified HGX servers with 8x H100 or H200 GPUs, NVSwitches and NVLinks for high-speed inter-GPU communication, and NVIDIA Spectrum™-X (Ethernet-based) networking.

These components are reliable, easy to manage, performant (high throughput and low latency), and provide exceptional AI inference capabilities.

This paper provides AI professionals with

- A list of core components and infrastructure choices
- Deployment considerations
- Performance validation of VMware Private AI Foundation with NVIDIA, offering organizations a comprehensive guide to optimizing AI inference workloads in a private cloud environment.

## Introduction

The rapid evolution of artificial intelligence (AI) has driven enterprises to prioritize scalable, cost-efficient infrastructure for managing AI workloads. While cloud platforms initially provided agility for experimentation, the costs of large-scale AI inference in the cloud—along with risks of data exposure and governance—are compelling organizations to shift towards on-prem solutions. However, this transition introduces several key challenges:

1. **GPU underutilization:** On-prem GPUs are often underutilized in enterprise data centers, with scenarios such as assigning GPUs to a less frequently used model and over-provisioning for peak loads, leading to wasted resources and reduced return on investment (ROI). Optimizing GPU usage is crucial to prevent resource hoarding and enhance efficiency, similar to the optimization challenges faced in the early stages of CPU utilization.
2. **Giving data scientists a cloud-like interface:** The rapid pace of AI models and toolkit updates creates challenges for data scientists who need a flexible, cloud-like interface. At the same time, infrastructure provisioning remains largely an IT responsibility, while data scientists require the freedom to focus on model development.
3. **Model governance:** As AI models increasingly use sensitive data, effective governance is critical. Enterprises must enforce security policies, prevent model drift, and ensure compliance. Private AI frameworks provide the necessary control to secure proprietary data and maintain model reliability.
4. **Familiar management interface:** VMware's user-friendly infrastructure management interface is widely adopted in enterprise IT environments, enabling IT teams to manage AI workloads efficiently without the steep learning curve associated with new or open-source platforms. This familiarity enhances operational efficiency and reduces the potential for errors.

VMware Private AI Foundation with NVIDIA—a joint GenAI platform by Broadcom and NVIDIA—tackles these challenges by providing on-demand GPU allocation, enabling GPU sharing, and automating infrastructure management. It gives data scientists the resource flexibility they need to focus on actual development while allowing IT teams to streamline deployment, enforce governance policies, secure proprietary data, and ensure model compliance over time.

This paper presents a reference design for deploying VMware Private AI Foundation with NVIDIA to support AI inference workloads on NVIDIA-Certified HGX Systems equipped with 8x H100/H200 GPUs, NVSwitches, and NVLinks over Ethernet networking. The proposed design offers a prescriptive solution to ensure efficient, secure, and agile AI development, deployment, and operation.

## Intended audience

This reference design is for the following groups:

- **New users:** Organizations that own or plan to procure NVIDIA-Certified HGX Systems (8x H100/H200 GPUs with NVSwitches and NVLinks) and plan to deploy VMware Private AI Foundation with NVIDIA for AI inference workloads on a VCF-based private cloud infrastructure.
- **Infrastructure architects:** VCF admins, DevOps, and SRE teams responsible for deploying and managing both physical and virtual infrastructures. This paper provides guidance on hardware, networking, and system configurations to integrate VMware Private AI Foundation with NVIDIA.
- **Software architects:** Data scientists and AI developers who will develop inference workloads on the platform. This paper offers architecture and performance insights to help optimize their AI models and workloads.

## Terminology

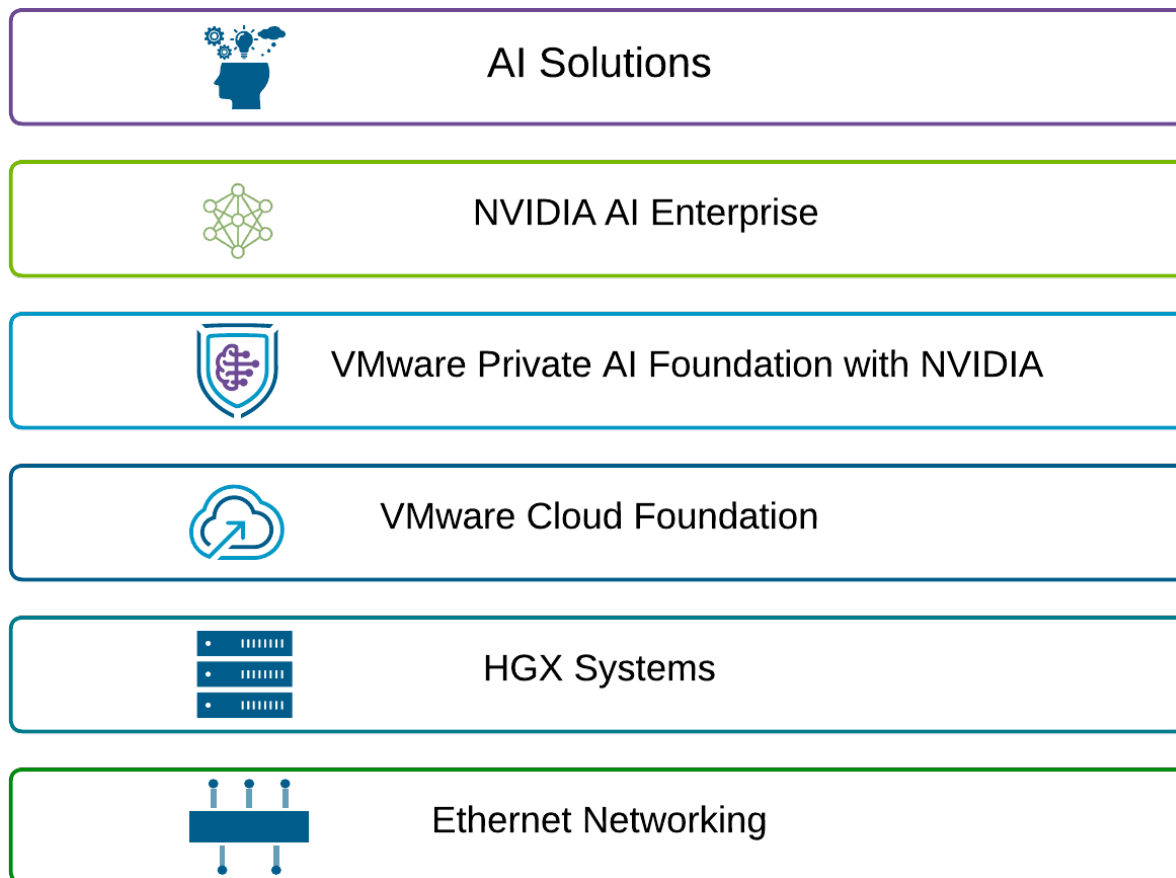
Table 1. Terms and acronyms used throughout this paper

Term	Definition
CNI	Container Network Interface
DLVM	Deep Learning Virtual Machines
DRS	Distributed Resource Scheduler
DSM	Data Service Manager
ESA	Express Storage Architecture for vSAN
HA	High Availability
LLM	Large Language Model
LCM	Life Cycle Management
NGC	NVIDIA GPU Cloud
OOB	Out-Of-Band
OSA	Original Storage Architecture for vSAN
RAG	Retrieval Augmented Generation
RDMA	Remote Direct Memory Access
SRE	Site Reliability Engineers
TEP	Tunnel End Point
VCF	VMware Cloud Foundation
VDS	Virtual Distributed Switch
vGPU	NVIDIA vGPU (C-Series)
VIB	vSphere Installation Bundles
VKS	vSphere Kubernetes Service
vLCM	vSphere Lifecycle Manager
VM	Virtual Machine
VI WLD	Virtual Infrastructure Workload Domain

## Core components

Figure 1 depicts the layered architecture of the solution, where each layer represents a crucial integration point that often requires manual setup and configuration for the efficient execution of GenAI solutions. VMware Private AI Foundation with NVIDIA is an advanced service for VMware Cloud Foundation (VCF) that is available for purchase. NVIDIA AI Enterprise, which is purchased directly from NVIDIA, is also required to deliver essential software packages such as NVIDIA vGPU (C-Series), GPU Operator, NVIDIA Network Operator, and NVIDIA NIM. This reference design simplifies the deployment and optimization of each layer by providing validated, prescriptive guidance, ensuring a smoother and more efficient implementation.

Figure 1. VMware Private AI Foundation with NVIDIA on HGX servers



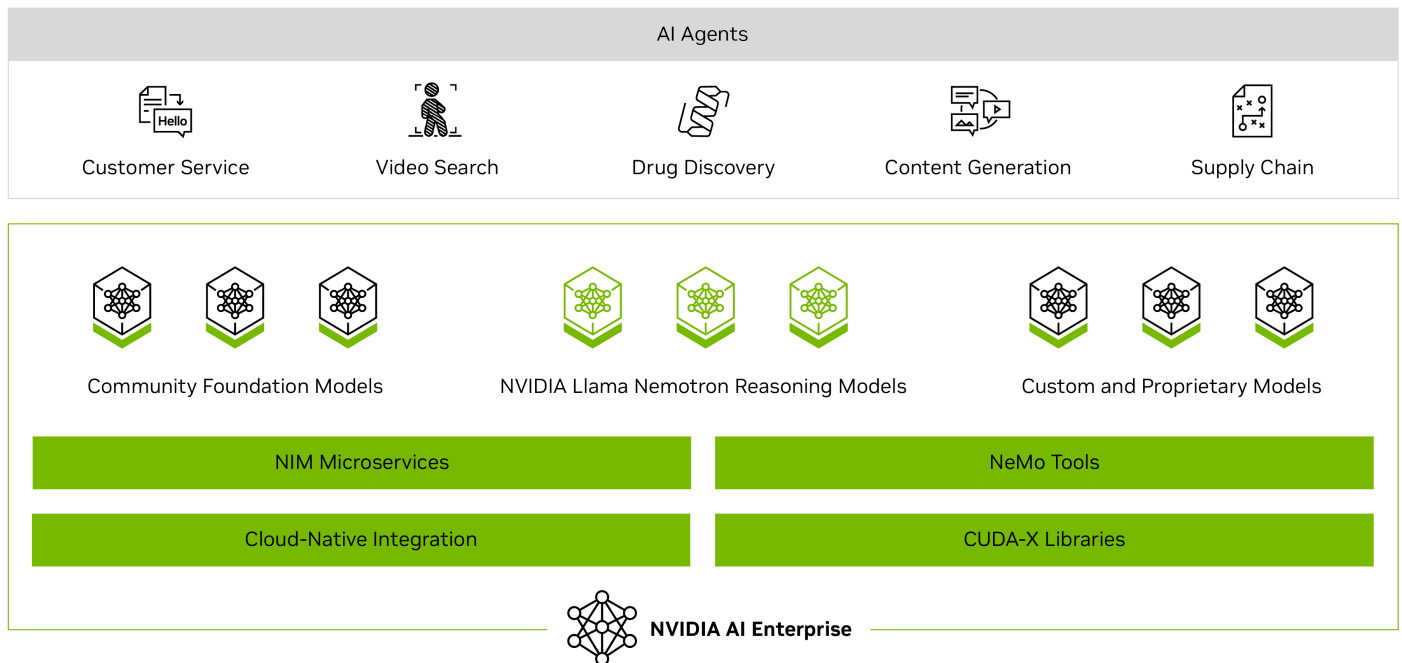
In the following sections, a technical overview and key considerations for each core component are provided, shown in Figure 1.



## NVIDIA AI Enterprise

[NVIDIA AI Enterprise](#) is a cloud-native software platform that streamlines the development and deployment of production-grade, end-to-end generative AI pipelines and helps organizations build data flywheels for the next era of agentic AI. The product versions for each release are shown in the [NVIDIA AI Enterprise release notes](#).

Figure 2. NVIDIA AI Enterprise overview



## NVIDIA vGPU (C-Series)

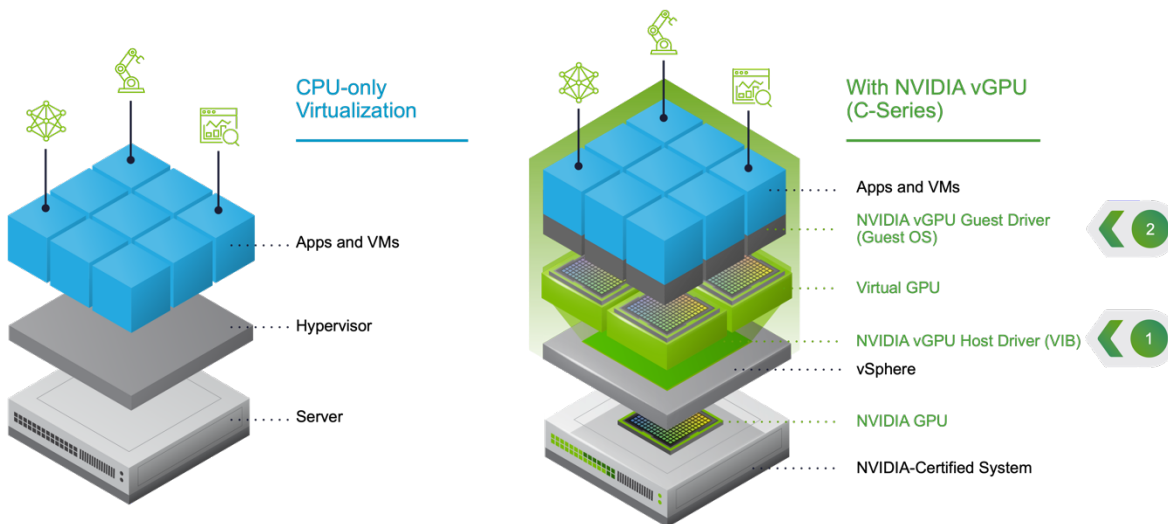
NVIDIA virtual GPU (vGPU) technology enables efficient GPU resource pooling and sharing of one or multiple physical GPUs in a virtualized environment. It is built on a variation of passthrough, providing direct access to host physical GPUs for optimal performance while retaining the flexibility of virtualization with features like snapshots, vMotion, and more. This approach aligns with VMware software's ability to abstract physical infrastructure—compute, storage, and network—into a software layer (hypervisor) for resource pooling and consumption, allowing multiple VMs to leverage underlying hardware effectively. Additionally, all 8x GPUs (or a subset) connected via high-speed NVSwitch and NVIDIA NVLink™ in an HGX server can be allocated to a single VM with the vSphere [device-group](#) capability.

# VMware Private AI Foundation with NVIDIA on HGX Servers: Reference Design for Inference

The [NVIDIA vGPU \(C-Series\)](#) is designed for compute-intensive workloads. As shown in Figure 3, unlike CPU-only virtualization, NVIDIA vGPU (C-Series) requires drivers to be installed at both the host and guest levels. At the host level, a vSphere installation bundle (VIB) is installed, which can be automated through vCenter Lifecycle Manager (vLCM). On the guest side, it is essential to install an NVIDIA vGPU (C-Series) guest driver version that matches the host VIB to ensure compatibility and optimal performance. Manually performing this installation can be error-prone and time-consuming. **Deploying a deep learning virtual machine (DLVM) simplifies this process by automatically installing the appropriate NVIDIA vGPU (C-Series) guest driver based on the host's VIB version.** This approach improves hardware utilization and simplifies the deployment and management of GPU resources in virtualized environments.

The NVIDIA vGPU (C-Series) host driver VIB can be downloaded from the [NVIDIA NGC](#) catalog. The `NVD-AIE-xx.zip` package also includes a `mgmt-daemon` VIB for monitoring GPU metrics in VMware Aria Operations. These two VIBs can be integrated into a vLCM image for use during the creation of a virtual infrastructure workload domain (VI WLD) or remediation of a cluster. The NVIDIA vGPU (C-Series) guest driver can also be downloaded from the [NVIDIA NGC](#) catalog. However, before obtaining these drivers, **an NVIDIA AI Enterprise license must first be acquired.** Additionally, an NVIDIA license server instance must be configured within the [NVIDIA Application Hub](#), which generates a licensing portal API key and a client configuration token file. The key and token file are then used to enable the full capabilities of the guest vGPU driver in the AI workstation or the AI Kubernetes cluster.

Figure 3. vGPU technology diagram

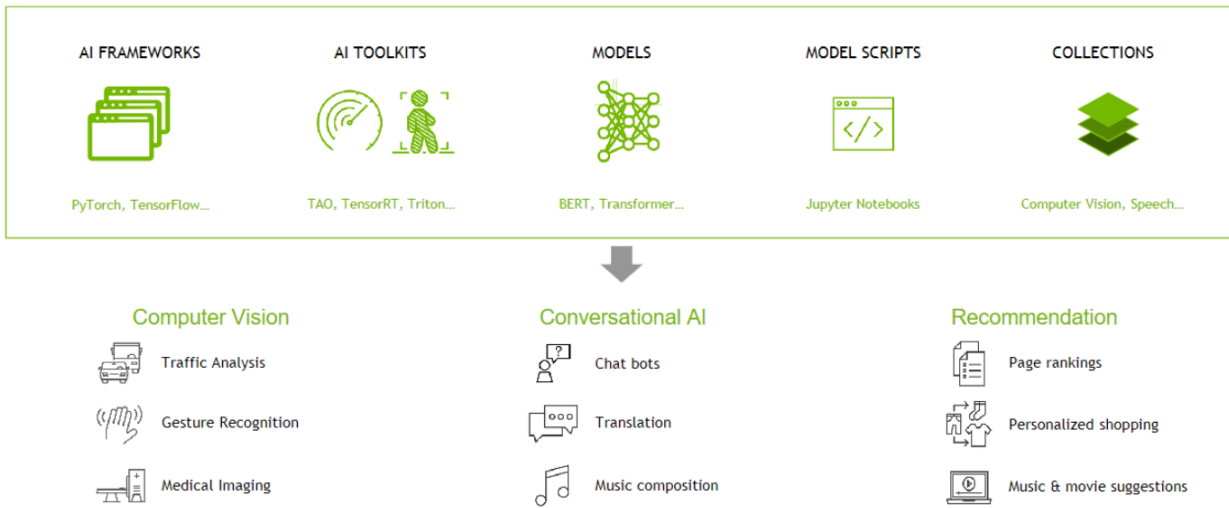


## NGC

NVIDIA AI Enterprise includes [NVIDIA NIM](#), other AI microservices, GPU Operator, [NVIDIA NIM Operator](#), and [NeMo Retriever](#). All of these components are available from [NVIDIA NGC](#). Downloading these resources requires an [NGC API Key](#).

# VMware Private AI Foundation with NVIDIA on HGX Servers: Reference Design for Inference

Figure 4. NVIDIA NGC catalog overview

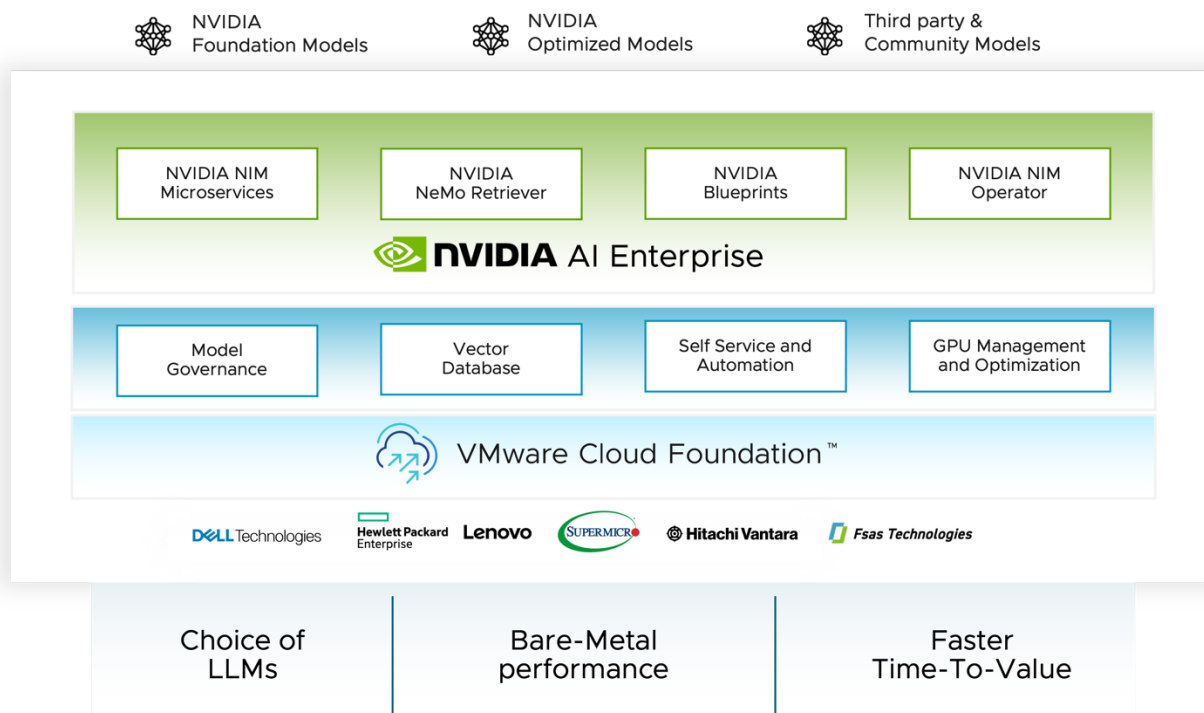


Software and containers hosted on NGC undergo scans against an aggregated set of common vulnerabilities and exposures (CVEs), crypto, and private keys. They are tested and designed to scale up to multiple GPUs and, in many cases, to multi-node. For more information, refer to the [NVIDIA NGC User Guide](#).

## VMware Private AI Foundation with NVIDIA

VMware Private AI Foundation with NVIDIA is an add-on advanced service to VCF for provisioning AI workloads based on NVIDIA AI Enterprise. It is a jointly engineered product between NVIDIA and Broadcom that enables enterprises to deploy generative AI workflows, such as retrieval augmented generation (RAG), using their proprietary data. The solution allows organizations to run inference workloads on models hosted in their private stores, addressing concerns related to privacy, choice, cost, performance, and compliance.

Figure 5. VMware Private AI Foundation with NVIDIA

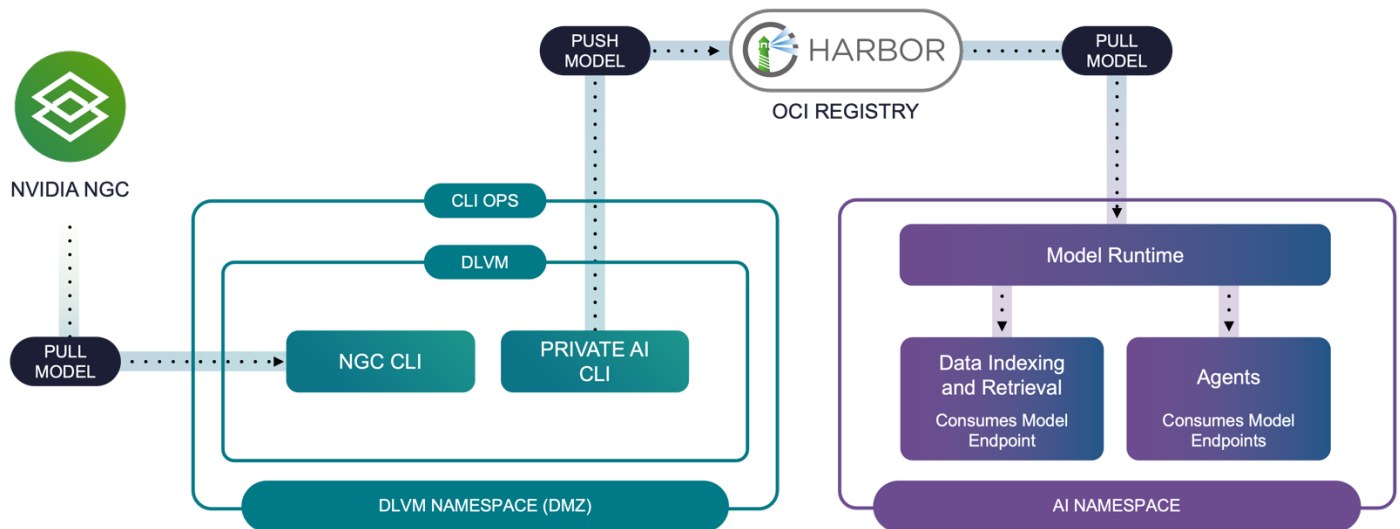


# VMware Private AI Foundation with NVIDIA on HGX Servers: Reference Design for Inference

VMware Private AI Foundation with NVIDIA offers a comprehensive solution for enterprises to develop and deploy AI applications securely and efficiently with the following key components (shown in Figure 5):

1. **Model governance** allows data scientists to test, evaluate, and store pre-trained LLMs or containers that are deemed safe and suitable for business use. The workflow, illustrated in Figure 6, begins with model testing and validation within a deep learning VM (DLVM) in an isolated environment to ensure safety and control. Once validated, these artifacts are stored in a model gallery hosted on the Harbor Registry, with evaluation processes customized to meet each enterprise's specific requirements. After passing evaluations, the models are promoted to a Kubernetes-based deployment, making them accessible to developers in a scalable environment.

Figure 6. Model Governance Journey in VMware Private AI Foundation with NVIDIA



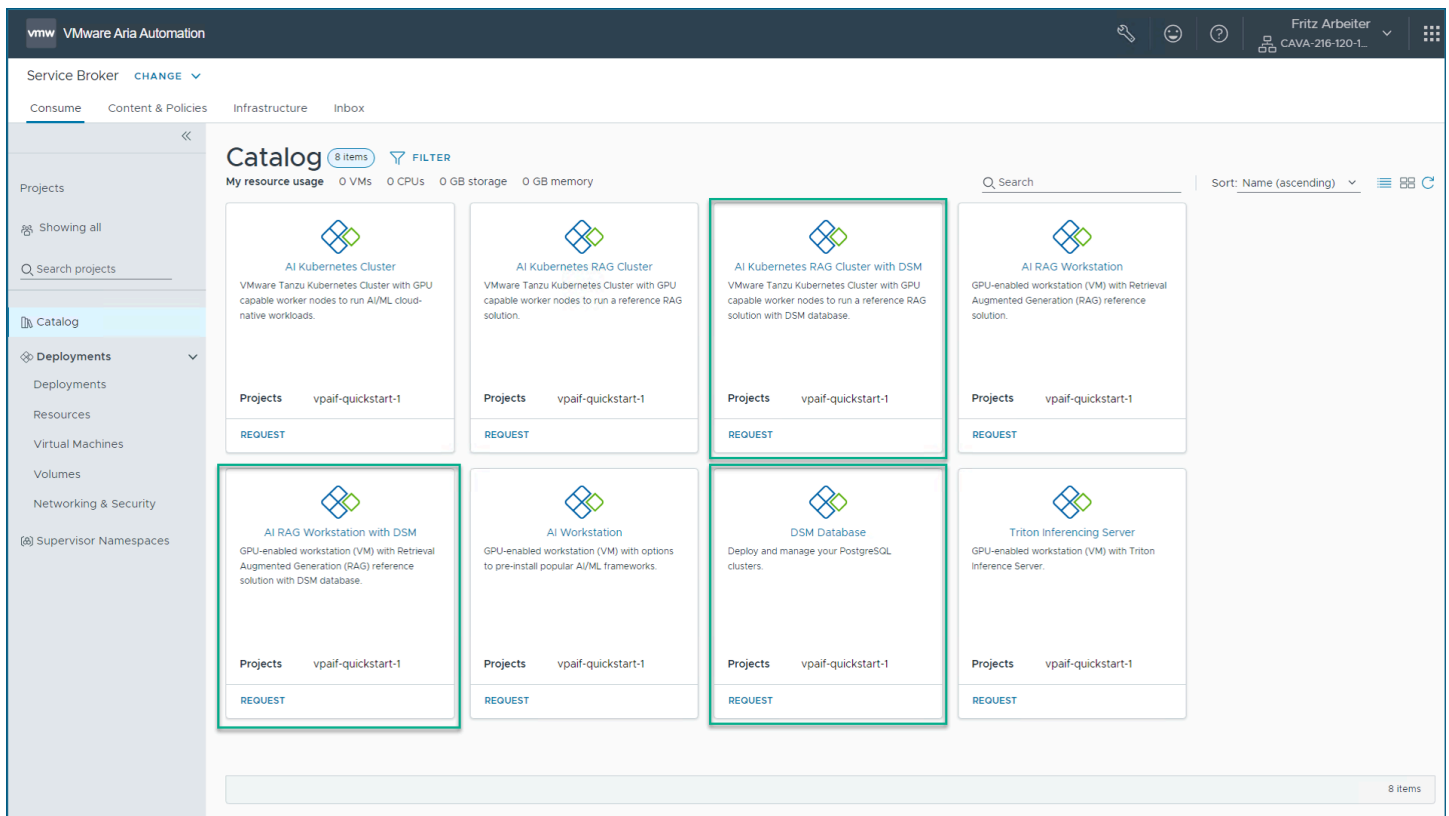
2. **Vector database** functionality is provided via PostgreSQL with the pgVector extension and deployed as a containerized service running in a Kubernetes cluster. Deployment can be automated through VMware Aria Automation Data Service Manager (DSM) integration. In addition to vector database services, DSM can automate provisioning databases of choice (for example, MySQL) for various use cases, such as [storing conversation history](#) in a real-world RAG application. This solution ensures the secure storage and retrieval of vector embeddings for private business data, which is essential for RAG applications, while enforcing RBAC and encryption. PostgreSQL's proven reliability and enterprise-grade features guarantee that private business data remains both protected and accessible for AI workflows.
3. **Self-service automation**, powered by VMware Aria Automation Service Broker, enables the provisioning of DLVMs for model development and Kubernetes clusters for production scaling. This solution benefits both the IT team and data scientists. Data scientists gain access to a cloud-like experience through a self-service catalog, equipped with the necessary software bundles. The [Quickstart Wizard](#) creates five basic catalog items that can be generally categorized into two use cases, and IT Ops can easily create or customize additional catalog items (shown in Figure 7) using [Automation Assembler](#). With this approach, there is no need for tickets or lengthy interactions between the IT team and data scientists. Instead, there is a smooth, streamlined selection process that empowers both teams to quickly get what they need.

# VMware Private AI Foundation with NVIDIA on HGX Servers: Reference Design for Inference

Table 2. Two fundamental use cases in self-service automation

Use Case	Target Users	Environment	Key Features
Development	Data scientists / Infrastructure Architects	Deep Learning VM	A GPU-enabled VM image with curated VM settings and software bundles, such as vGPU drivers, AI development tools, and LLM microservices (for example, NIM, NeMo Retriever, etc.)
Production	DevOps / Infrastructure Architects	Kubernetes cluster	Automate 50+ steps to simplify deployment of AI applications on vSphere Kubernetes Service once they are production-ready

Figure 7. Find catalog items related to VMware Private AI Foundation with NVIDIA in Aria Automation Service Broker



4. **GPU Monitoring** in Aria Operations provides real-time visibility into hardware performance, tracking metrics like GPU utilization, memory consumption, and temperature across VMs on the system as a whole.

## VMware Cloud Foundation

VMware Cloud Foundation (VCF) is a unified full-stack private cloud platform designed to streamline the deployment, automation, and management of software-defined infrastructure. The platform is optimized for modern workloads such as AI and container-based applications. It integrates core technologies—VMware vSphere® (compute), VMware® vSAN™ (storage), VMware NSX® (networking), and the VMware Aria suite (VMware Aria® Automation™ and VMware Aria Operations)—into a single, efficient, and consistent solution. By adopting a standardized, automated approach, VCF simplifies resource management, reduces technical debt, and enhances operational agility, enabling organizations to accelerate application delivery and embrace cloud-native technologies. The private cloud platform supports GPU virtualization for AI tasks and offers the flexibility of on-prem and edge deployments, providing a consistent cloud operating model for the scalable, efficient management of diverse workloads.

Table 3 lists the [bill of materials \(BOM\)](#) and corresponding versions in VCF 5.2.1 used in this design. Table 4 lists each component’s function and key features.

Table 3. VCF 5.2.1 BOM used in this reference design

Software Component	Version
Cloud Builder VM	5.2.1
VMware SDDC Manager	5.2.1
VMware vCenter® Server Appliance	8.0 Update 3c
VMware ESX®	8.0 Update 3b
VMware NSX	4.2.1
VMware Aria Suite Lifecycle	8.18
VMware Private AI Foundation with NVIDIA	5.2.1

Table 4. VCF components used in this reference design

Component	Function	Key Features
<a href="#">SDDC Manager</a>	Centralized management for VCF	<ul style="list-style-type: none"> <li>Automates deployment/lifecycle of vSphere, vSAN, NSX</li> <li>Manages workload domains</li> <li>Enforces pre-validated architectures</li> <li>Provisions with an API-driven infrastructure</li> </ul>
<a href="#">vSphere</a>	Enterprise virtualization platform	<ul style="list-style-type: none"> <li>Abstracts physical compute resources</li> <li>Features foundation products ESX hypervisor and vCenter appliance</li> <li>Integrates vGPU and Fabric Manager (to configure NVSwitch memory fabrics) into ESX</li> <li>Uses virtual distributed switch (VDS) to simplify networking</li> <li>Uses VMware vSphere® Distributed Resource Scheduler™ (DRS) for workload balancing</li> <li>Includes VMware vSphere® High Availability (HA) for automated failover</li> </ul>

# VMware Private AI Foundation with NVIDIA on HGX Servers: Reference Design for Inference

Component	Function	Key Features
↳ <a href="#">vSAN</a>	Software-defined storage solution	<ul style="list-style-type: none"> <li>• Creates a shared storage pool from local or direct-attached storage</li> <li>• Supports original storage architecture (OSA) and express storage architecture (ESA)</li> <li>• Includes vSAN file services for enterprise-grade NFS shares</li> </ul>
↳ vSphere Kubernetes Service (VKS)	Kubernetes native infrastructure on vSphere	<ul style="list-style-type: none"> <li>• Transforms vSphere into a Kubernetes platform by enabling a supervisor cluster that runs on ESX</li> <li>• Provides a VM service that manages VMs in a declarative manner through the Kubernetes API</li> <li>• Uses Harbor Registry for AI model/container governance</li> <li>• Uses <a href="#">vSphere CSI</a> for container-native storage capabilities</li> <li>• Integrates NSX and <a href="#">Antrea CNI</a> for container networking</li> </ul>
<a href="#">NSX</a>	Software-defined networking and security solution	<ul style="list-style-type: none"> <li>• Supports multi-tenancy through virtual private clouds (VPCs)</li> <li>• Provides logical switching using GENEVE encapsulation</li> <li>• Uses a two-tier routing model (Tier-0 and Tier-1 gateways) with network services such as NAT, VPN, and more</li> <li>• Provides microsegmentation via distributed firewall (DFW)</li> <li>• Integrates load-balancing services</li> </ul>
<a href="#">Workspace ONE Access</a>	Identity and access management component within VCF	<ul style="list-style-type: none"> <li>• Provides single sign-on capabilities with multi-factor authentication</li> <li>• Centralizes identity governance across all VCF services and workloads</li> <li>• Integrates with existing identity providers</li> </ul>
<a href="#">Aria Lifecycle</a>	Deploy, configure, manage, and upgrade Aria Suite products	<ul style="list-style-type: none"> <li>• Deploys in the management domain</li> <li>• Integrates with SDDC Manager, allowing users to initiate deployments from the SDDC Manager UI and complete them in the Aria Suite Lifecycle interface</li> </ul>
<a href="#">Aria Automation</a> *	Automates infrastructure provisioning, configuration, and management	<ul style="list-style-type: none"> <li>• Enables self-service provisioning and management of infrastructure as a service (IaaS)</li> <li>• Automates Kubernetes, network, data services, and private AI workloads</li> </ul>
<a href="#">Aria Operations</a> *	Monitor, optimize, plan, and scale applications and infrastructure	<ul style="list-style-type: none"> <li>• Provides comprehensive monitoring, troubleshooting, and capacity management for VCF</li> <li>• Includes operations management: Optimizes the cost and performance of compute, storage, and network resources</li> <li>• Includes security management: Provides visibility into hardening, governance, and compliance</li> </ul>

\* Beginning with VCF 9.0, Aria Automation and Aria Operations will be rebranded as VCF Automation and VCF Operations, respectively.

↳ This symbol signifies the component is part of the larger component above it.



## HGX systems

NVIDIA [HGX series](#) servers have become very popular. The HGX platform, inspired by the NVIDIA DGX series, provides a flexible and configurable GPU solution from your trusted OEM provider. Server vendors can directly buy the standard NVIDIA HGX baseboard. This baseboard features 8 SXM form-factor GPUs linked through 4 NVSwitches and NVLinks. OEM vendors can customize other components like CPU, RAM, storage, and NICs around the baseboard while keeping a consistent GPU setup. This allows them to submit their systems for HGX certification under the [NVIDIA-Certified Systems program](#). This approach makes it easier to integrate multiple GPUs and advanced technologies, leading to more flexible, efficient, and powerful systems for AI and HPC.

Table 5 provides a filtered list of HGX systems [retrieved from the NVIDIA qualified system catalog](#). These systems include 8x H100/H200 SXM GPUs, and partners have certified them (in March 2025) for the data center used in this reference design.

Table 5. NVIDIA-Certified HGX systems with 8x-H100/H200-SXM-GPU

Partner	System	NVIDIA GPU	CPU
ASUS	ESC N8-E11	8x H100-80GB GPU	2x Intel Sapphire Rapids
	ESC N8-E11	8x H200-141GB GPU	2x Intel Sapphire Rapids
	ESC N8A-E12	8x H100-80GB GPU	2x AMD Genoa
Dell Technologies	PowerEdge XE9680	8x H100-80GB GPU	2x Intel Emerald Rapids
	PowerEdge XE9680	8x H200-141GB GPU	2x Intel Emerald Rapids
	PowerEdge XE9680	8x H100-80GB GPU	2x Intel Sapphire Rapids
Gigabyte	G593-ZD2-A	8x H100-80GB GPU	2x AMD Genoa
	G593-SD2-A	8x H100-80GB GPU	2x Intel Sapphire Rapids
	G593-SD1-A	8x H100-80GB GPU	2x Intel Sapphire Rapids
	G593-SD0-A	8x H100-80GB GPU	2x Intel Sapphire Rapids
HPE	Cray XD670 Gen 11	8x H100-80GB GPU	2x Intel Sapphire Rapids
	Cray XD670 Gen 11	8x H200-141GB GPU	2x Intel Emerald Rapids
Lenovo ISG	SR685a V3	8x H200-141GB GPU	2x AMD Genoa
	SR685a V3	8x H100-80GB GPU	2x AMD Genoa
Quanta QCT	QuantaGrid D74H-7U	8x H100-80GB GPU	2x Intel Sapphire Rapids
Supermicro	AS-4125GS-TNHT2-LCC	8x H100-80GB GPU	2x AMD Genoa
	AS-8125GS-TNHR	8x H100-80GB GPU	2x AMD Genoa
	CAS-IH828	8x H100-80GB GPU	2x Intel Sapphire Rapids
	SYS-421GE-TNHR2-LCC	8x H100-80GB GPU	2x Intel Sapphire Rapids
	SYS-821GE-TNHR (Rear I/O)	8x H100-80GB GPU	2x Intel Sapphire Rapids
	SYS-821GE-TNHR (Rear I/O)	8x H200-141GB GPU	2x Intel Sapphire Rapids
	SYS-821GE-TNHR-LCC	8x H100-80GB GPU	2x Intel Sapphire Rapids

**Important:** To support an LLM such as [DeepSeek-R1 NIM](#) as a NIM microservice, a single 8x H200 HGX server can host the model, as outlined in the NVIDIA blog post: [DeepSeek-R1 Now Live With NVIDIA NIM](#). However, for multi-node inference (such as deploying a single DeepSeek-R1 NIM across two H100 HGX servers), the paper does not cover cases where the vendor restricts enabling address translation services (ATS) or access control services (ACS), which may render the server incompatible with a hypervisor.

## Ethernet networking

Ethernet's simplicity, scalability, and cost-effectiveness make it an ideal choice for private AI inference deployments. Moreover, for future-proofing infrastructure, particularly for model customization or fine-tuning, Ethernet offers efficient scalability and flexibility.

Proper networking is critical to ensuring that VMware Private AI Foundation with NVIDIA on HGX servers does not have any bottlenecks or suffer performance degradation for AI workloads. Advancements in Ethernet technology, such as RoCE v2 and lossless fabrics, combined with its open ecosystem and scalability, make it a good fit for both AI inference and training.

For more information, refer to [Broadcom Networking](#) and [NVIDIA Spectrum Ethernet](#).

## Reference architecture

This reference architecture for VMware Private AI Foundation with NVIDIA on HGX servers acts as a blueprint that provides prescriptive architectures designed to meet the changing needs of AI workloads.

### Physical architecture

The components of VMware Private AI Foundation with NVIDIA on HGX servers are described in Table 7. For larger-scale deployments, please refer to Section 5.2 for detailed sizing guidance or reach out to VMware Professional Services for further assistance.

Table 6. Physical architecture components

Component	Technology
Inference servers in a Virtual Infrastructure Workload Domain (4–16)	<ul style="list-style-type: none"> <li>• HGX System with 8x H100/H200 SXM GPUs and 4x NVSwitches interconnected by Gen4 NVLink</li> <li>• Min 2x 25Gbps NICs for VCF infrastructure service network (for example, management, vSAN storage, etc.)</li> <li>• Min 2x 100GbE NICs for inference workload network</li> <li>• Although the UI of VMware Private AI Foundation with NVIDIA requires a minimum of 3 ESXi hosts, a 4-node (N+1) configuration is recommended for resilience, ensuring redundancy for replicas and RAID-5 erasure coding while mitigating failure risks during maintenance or unexpected issues.</li> </ul>
Inference fabric	Min 100 GbE Ethernet Switch
Management and storage fabric	Min 25 GbE Ethernet Switch
Out-of-band management fabric	Min 1 GbE Ethernet Switch
Management servers in the Management Domain	<ul style="list-style-type: none"> <li>• Min 4x vSAN Ready nodes certified for vSAN OSA or ESA</li> <li>• Memory and compute should be sized based on the VCF components planned to be deployed in the Mgmt Domain</li> <li>• Use the <a href="#">vSAN Sizer Tool</a> for sizing guidance</li> </ul>

Figure 8 depicts the physical architecture for a minimum of 4-node vSAN-ready setup to create the VCF management domain, alongside up to 16 HGX servers to establish the VI WLD with the Ethernet networking switch’s radix (ports per switch) equal to 32. Further scale-out is possible with a network redesign. Each HGX H100/H200 system uses 8 connections to link the workload networks. The full architecture includes three distinct networks: an Ethernet-based workload network backed by two switches, each with at least 100 GbE; an Ethernet fabric for VCF management and storage backed by two switches, each with at least 2x 25GbE switches; and an out-of-band Ethernet network.

**Note:** Figure 8 represents an example setup. For VI WLD, the minimum configuration requires 4 HGX servers and 1 workload network switch. If the workload Ethernet switch has a higher radix (more than 32) or if multiple switches are used in a multi-layer network design, additional HGX servers can be incorporated into the deployment. Any OEM HGX server that meets the minimum requirements outlined in Table 5 can be used.

# VMware Private AI Foundation with NVIDIA on HGX Servers: Reference Design for Inference

Figure 8. Physical architecture with a minimum of 4 and up to 16 HGX servers based on ethernet switch radix = 32.

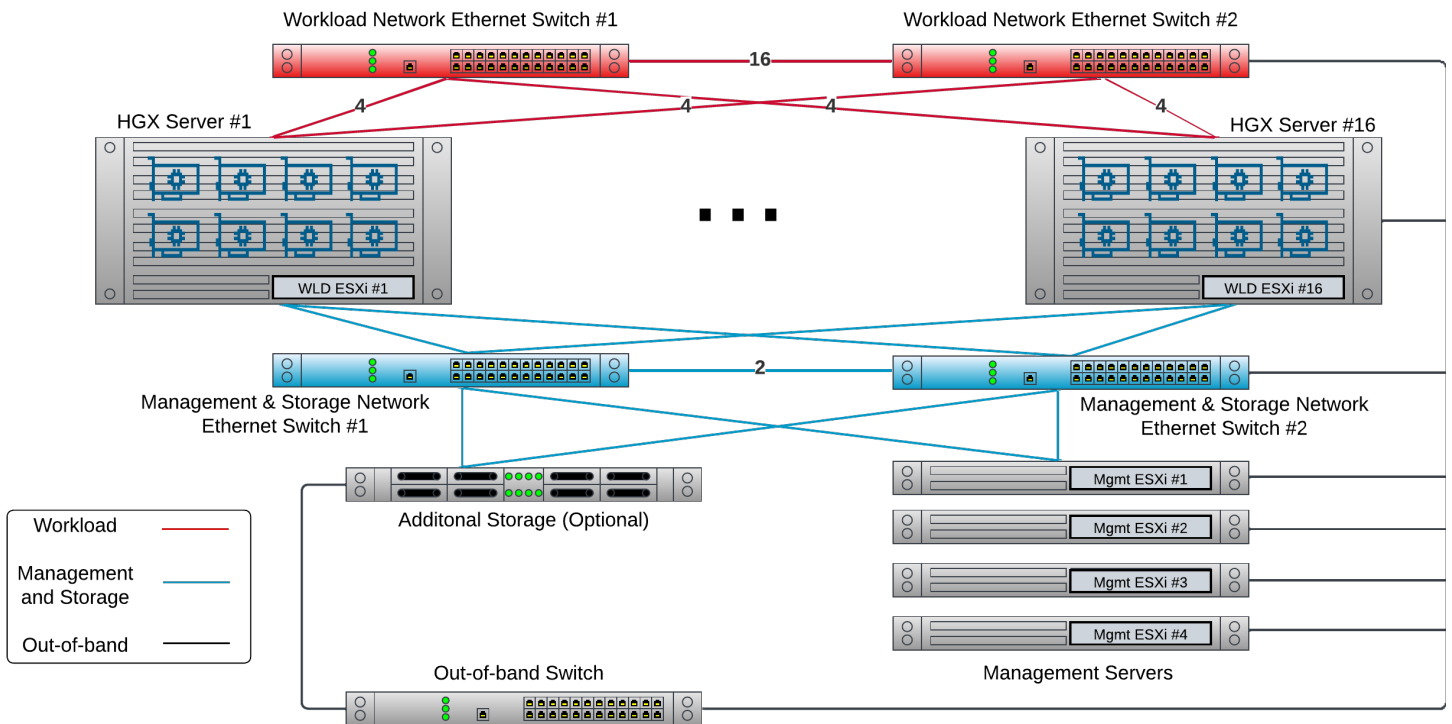
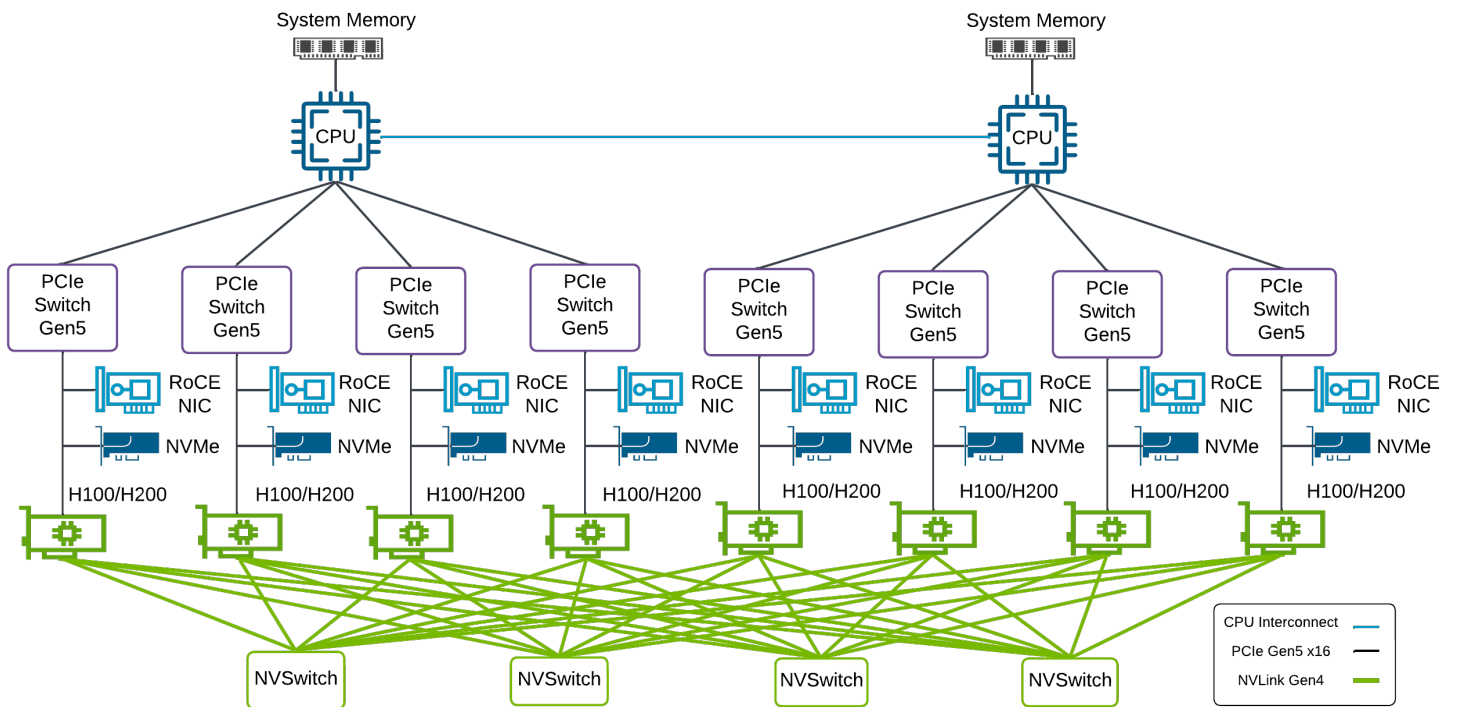


Figure 9 presents the balanced topology diagram of a typical HGX server. Table 6 provides the list of the network prerequisites for the enterprise networking team to prepare accordingly.

Figure 9. Topology diagram of a typical HGX server



# VMware Private AI Foundation with NVIDIA on HGX Servers: Reference Design for Inference

Table 7. Typical HGX server system configuration for GenAI inference

Parameter	Inference Server Configuration
GPU	H100 / H200 SXM
GPU Configuration	8× GPUs connected by 4× NVSwitch chips and NVLinks within a server GPUs should be balanced across CPU sockets and root ports.
CPU	x86 PCIe Gen5-capable CPUs are recommended, such as Intel Xeon scalable processor (Sapphire Rapids) or AMD Genoa.
CPU Sockets	Min 2 CPU sockets
CPU Speed	Min 2.1 GHz base clock
CPU Cores	Min 6× physical CPU cores per GPU
System Memory	Min 1.5× (or 2.0×) of the total GPU memory size of DDR5 is recommended. Evenly spread across all CPU sockets and memory channels.
PCIe	Min 1× Gen5 ×16 link per Gen5 GPU is recommended.
PCIe Topology	For balanced PCIe architecture, GPUs should be evenly distributed across CPU sockets and PCIe root ports. NICs and NVMe drives should be placed within the same PCIe switch or root complex as the GPUs.
PCIe Switches	Direct CPU attached is preferred.
VCF networks NICs	RDMA capability is not necessarily required for AI inference within an HGX server. <a href="#">NVIDIA ConnectX NICs</a> or <a href="#">Broadcom Ethernet NICs</a>
VCF networks NICs Speed	vSAN ESA requires a minimum of 25Gbps. The network for AI inference requires a minimum of 25 GbE per GPU to meet text data's bandwidth AI inference demands. To future-proof the infrastructure, a minimum of 2× 100 GbE NICs per HGX server can be optionally used.
Out-of-band network NIC	1GbE Ethernet NIC.
Storage	At least 1× 3.84TB Gen4 NVMe per CPU socket. To future-proof the infrastructure, 1× 3.84 TB U.2 NVMe drive for each GPU under the same PCIe switch of a GPU is preferred. For this inference reference architecture, all the NVMe disks available will be pooled and controlled by vSAN ESA.

# VMware Private AI Foundation with NVIDIA on HGX Servers: Reference Design for Inference

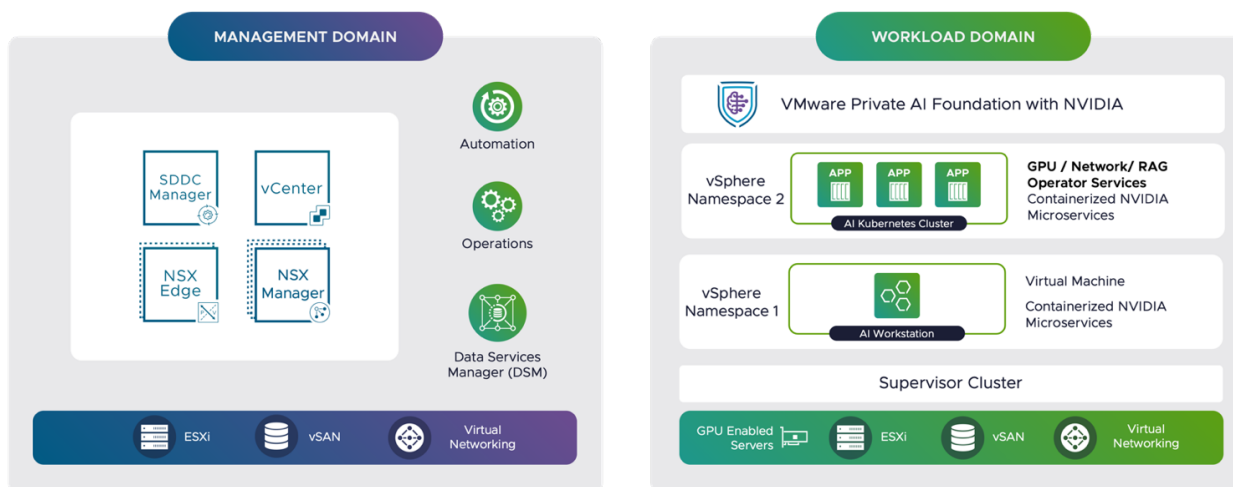
Table 8. Network to be prepared by the enterprise networking team.

Network	Function	VLAN/Overlay	Routed or NAT	WLD
VM management	IP Connectivity for "VM Management" traffic.	VLAN	Routed	Mgmt WLD & VI WLD
ESXi Management	IP connectivity for ESXi	VLAN	Routed	Mgmt WLD & VI WLD
vSAN	vSAN traffic	VLAN	Routed or isolated based on L2 boundary	Mgmt WLD & VI WLD
vSphere vMotion	vMotion traffic	VLAN	Routed or isolated based on L2 boundary	Mgmt WLD & VI WLD
Host Overlay	NSX TEP traffic	VLAN	Routed	Mgmt WLD & VI WLD
Edge Uplink01	Edge Cluster Northbound connectivity via eBGP	VLAN	Routed	Mgmt WLD & VI WLD
Edge Uplink02	Edge Cluster Northbound connectivity via eBGP	VLAN	Routed	Mgmt WLD & VI WLD
Edge Overlay	Edge Node TEP VLAN	VLAN	Routed	Mgmt WLD & VI WLD
Supervisor Cluster Management Network	Used by the Supervisor control plane nodes. This VLAN can be shared with the "VM Management" VLAN	VLAN	Routed	VI WLD
Supervisor Namespace Networks	Used to allocate IP addresses to the workload network in the Supervisor Namespace	Overlay	NAT	VI WLD
Service IP Pool Network	Used by Kubernetes applications that need a service IP address	Overlay	NAT	VI WLD
Egress IP Pool Network	Used by NSX to create an IP pool for load balancing	VLAN	Routed	VI WLD
Ingress IP Pool Network	Used by NSX to create an IP pool for NAT endpoint use	VLAN	Routed	VI WLD
Kubernetes Cluster Service Pool Network	Internal CIDR block from which IPs for Kubernetes ClusterIP Services will be allocated	Overlay	NAT	VI WLD

## Virtual architecture

The virtual architecture in VCF is structured into two primary domains: the **Management (Mgmt) Domain** and the **Workload Domain (VI WLD)**, as shown in Figure 10. This separation ensures isolation between infrastructure management and AI workloads, enhancing security, scalability, operational efficiency, and separate lifecycle management.

Figure 10. Conceptual view of virtual architecture



The **Management Domain** is a dedicated set of infrastructure resources in the form of a vSphere cluster that hosts all the core management components and services required to operate the VCF environment. It requires **vSAN** and hosts critical components such as **SDDC Manager**, **vCenter**, **NSX**, **Automation**, **Operations**, and **Data Services Manager**. While these tools reside in the Mgmt Domain, they extend their functionality to VI WLDs, enabling provisioning, orchestration, policy enforcement, and governance.

**Workload Domains** are dedicated to running AI workloads and user applications. Best practices recommend creating separate VI WLDs for each business subsidiary in large organizations, ensuring resource isolation between different teams while operating within the same VCF infrastructure. Each VI WLD functions independently, with its own **vCenter** instance residing in the Mgmt Domain to manage its virtualized resources.

Within a VI WLD, AI workloads are provisioned as **AI Workstations (DLVMs)** within vSphere namespaces. As workload demands increase, **AI Kubernetes clusters** can be dynamically deployed or decommissioned based on priorities. Each cluster consists of nodes implemented as GPU-enabled VMs, allowing seamless scalability. **vSphere Namespaces** not only serve as resource pool boundaries but also enforce permissions and apply policies (e.g., storage policies), further enhancing isolation by segregating different users and projects. Additionally, other VMware Private AI Foundation with NVIDIA components and capabilities are deployed within vSphere namespaces.

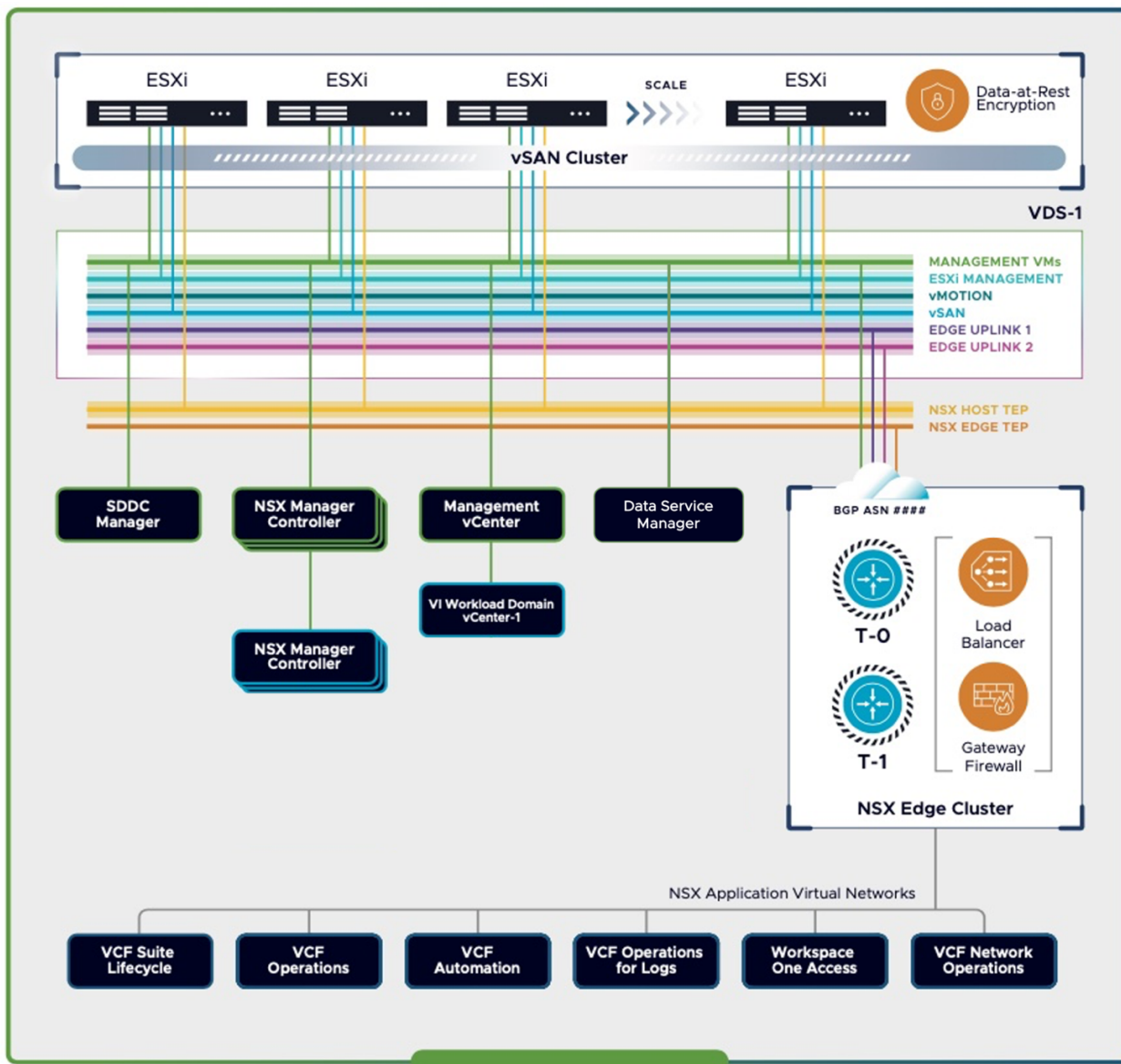
At the core of each VI WLD, layered on top of the vSphere infrastructure, is the **Supervisor Cluster**—a Kubernetes cluster responsible for provisioning vSphere namespaces and managing the resources within them. The Supervisor serves as the critical control plane component, enabling Kubernetes orchestration within vSphere by running directly on ESXi hosts. Acting as a translation layer between Kubernetes and vSphere, the Supervisor manages resource allocation and the full lifecycle of containerized workloads, including deployment, storage provisioning, and networking configuration, all while ensuring enterprise-grade security and control. New AI Kubernetes clusters can be deployed either through **Aria Automation** or a single `kubectl` command using the Kubernetes API, significantly streamlining production AI and data science operations.

## Management domain

### Management domain architecture

Figure 11 depicts the essential components in the VCF management domain, which requires a minimum of four ESX hosts to enable the vSAN service. The network segments are color-coded to represent different PortGroups in the vSphere Distributed Switch (VDS-1); detailed settings for VDS-1 and port groups are provided in Tables 10 and 11, respectively. The vCenter Server and NSX Managers with green borders manage the Mgmt Domain, while those with blue borders oversee a workload domain (VI WLD), as discussed in Section 3.4. Additionally, the NSX Edge cluster manages ingress and egress traffic to the Mgmt Domain, providing security and load balancing. The Aria Suite components are connected to NSX Application Virtual Networks.

Figure 11. Management Domain virtual architecture





# VMware Private AI Foundation with NVIDIA on HGX Servers: Reference Design for Inference

## Management domain settings

Table 9 shows the function and settings of each component in the management domain.

Table 9. Management Domain settings

Category	Component	Function	Number	Settings / Notes	Deployment Reference
SDDC	SDDC manager	Central management UI for VCF	1	<ul style="list-style-type: none"> <li>Use Cloud Builder's default configuration of SDDC Manager to deploy</li> <li>Protected by vSphere HA</li> </ul>	<a href="#">SDDC Manager Design for VMware Cloud Foundation</a>
vSphere	vCenter	Central management UI for configuring and monitoring vSphere infrastructure	1 for Mgmt Domain 1 for each VI WLD	<ul style="list-style-type: none"> <li>Mgmt Domain's vCenter is deployed by Cloud Builder VM</li> <li>VI WLD vCenter is deployed by SDDC Manager during VI WLD creation</li> </ul>	<a href="#">vCenter Server Design for VMware Cloud Foundation</a>
	vSAN	Create a single, shared storage pool across ESXi hosts in a domain	Require at least 4 ESXi hosts	<ul style="list-style-type: none"> <li>Required by Mgmt Domain</li> <li>Support OSA or ESA</li> </ul>	<a href="#">vSAN Design for VMware Cloud Foundation</a>
VMware Private AI Foundation with NVIDIA	Data Service Manger (DSM)	Provision and lifecycle of PostgreSQL-based vector database	1	<ul style="list-style-type: none"> <li>Requires a 1:1 relationship between each vCenter Server and a DSM appliance to manage a VI WLD</li> </ul>	<a href="#">VMware Data Services Manager Design for Private AI Ready Infrastructure for</a>
Aria Suite	Aria Automation Appliance	Automate IT processes and service delivery	A cluster of 3 appliances	<ul style="list-style-type: none"> <li>Enable Self-service catalog</li> <li>Deployed and LCM controlled by Aria Lifecycle</li> </ul>	<a href="#">Private Cloud Automation for VMware Cloud Foundation</a>
	Aria Operations Appliance	Monitor and optimize performance, capacity, and configuration	A cluster of 3 appliances (primary, data, replica)	<ul style="list-style-type: none"> <li>Provide the scale capacity required for monitoring up to 12,000 VMs or objects</li> <li>Support scale-out with additional data nodes.</li> </ul>	<a href="#">Intelligent Operations Management for VMware Cloud Foundation</a>
	↳ Aria Operations Cloud Proxies	Enhance scalability by collecting data from each VCF instance and sending it to the Aria Operations cluster	2 appliances deployed on the local-instance NSX segment	<ul style="list-style-type: none"> <li>Deployed and LCM controlled by Aria Lifecycle</li> </ul>	<a href="#">Configuring Cloud Proxies in VMware Aria Operations</a>
	Aria Suite Lifecycle	Life cycle management for Aria Suite and vIDM	1 on the cross-instance NSX segment	<ul style="list-style-type: none"> <li>Automate deployment, configuration, patching, upgrade, and content management across Aria Suite</li> <li>Enable VCF mode</li> <li>Protected by vSphere HA</li> </ul>	<a href="#">VMware Aria Suite Lifecycle Design for VMware Cloud Foundation</a>

# VMware Private AI Foundation with NVIDIA on HGX Servers: Reference Design for Inference

Category	Component	Function	Number	Settings / Notes	Deployment Reference
Workspace ONE	VMware Identity Manager (vIDM)	Central identity and access management solution for SSO, authentication, and user directory integration	a cluster of medium-size appliances	<ul style="list-style-type: none"> <li>• Deployed and LCM controlled by Aria Lifecycle</li> <li>• Clustered vIDM to support clustered deployment of Aria Automation</li> <li>• Protected by vSphere HA</li> </ul>	<a href="#">Workspace ONE Access Design for VMware Cloud Foundation</a>
NSX	NSX Manager	Central management UI of NSX for configuring and monitoring network components in Mgmt Domain	A cluster of 3 appliances	<ul style="list-style-type: none"> <li>• Deployed by Cloud Builder with a Virtual IP (VIP) address and an anti-affinity rule to ensure NSX managers running on different ESXi hosts</li> </ul>	<a href="#">Logical Design for NSX for VMware Cloud Foundation</a>
	NSX Edge Cluster	<ul style="list-style-type: none"> <li>• Provides HA, scalability, and distributed firewall services</li> <li>• Support Application Virtual Networks (AVNs)</li> </ul>	1	<ul style="list-style-type: none"> <li>• Profile Type: Set to Default.</li> <li>• Network Configuration: Layer 2 (L2) Uniform.</li> </ul>	<a href="#">Installing NSX Edge</a>
	↳ Tier-0 Router	North-south connectivity between NSX logical network and physical infrastructure	1	<ul style="list-style-type: none"> <li>• Active/Active (ECMP) routing for redundancy and load balancing.</li> <li>• BGP configured with unique ASNs for local and remote peers.</li> <li>• Large form factor to handle high throughput.</li> </ul>	<a href="#">NSX Design for VMware Cloud Foundation</a>
	↳ Tier-1 Router	East-west routing for internal network communication and stateful services	1 or more.	<ul style="list-style-type: none"> <li>• Active/Passive setup for stateful services like NAT and load balancing.</li> </ul>	<a href="#">NSX Design for VMware Cloud Foundation</a>

↳ This symbol signifies the component is part of the larger component above it.

# VMware Private AI Foundation with NVIDIA on HGX Servers: Reference Design for Inference

Tables 10 and 11 present the settings of VDS and its PortGroups (Network) in the Mgmt Domain.

Table 10. Management Domain VDS settings

VDS Configuration	Setting	Considerations	Deployment Reference
Single VDS	Single VDS with minimum 2 uplinks (NICs)	<ul style="list-style-type: none"> <li>A single VDS prepared for NSX enhances operational simplicity.</li> <li>If using only two uplinks, enable traffic sharing on the same NICs.</li> <li>Using 4+ uplinks is recommended to ensure proper traffic segmentation and bandwidth allocation.</li> </ul>	<a href="#">vSphere Networking Design for VMware Cloud Foundation</a>
Network I/O Control	Enabled	<ul style="list-style-type: none"> <li>Allow for proper bandwidth prioritization, increasing resiliency and performance of the network.</li> </ul>	<a href="#">vSphere Networking Design for VMware Cloud Foundation</a>
Data Path Mode	Enhanced Data Path	<ul style="list-style-type: none"> <li>Recommended mode for vSphere clusters running NSX Edge nodes</li> </ul>	<a href="#">vSphere Networking Design for VMware Cloud Foundation</a>
NSX Transport Zones	VLAN and Overlay Transport Zones	<ul style="list-style-type: none"> <li>Allows for flexible network segmentation and supports both traditional VLAN-based and overlay-(Geneve)-based networking.</li> </ul>	<a href="#">vSphere Networking Design for VMware Cloud Foundation</a>

Table 11. Management Domain VDS portgroup settings

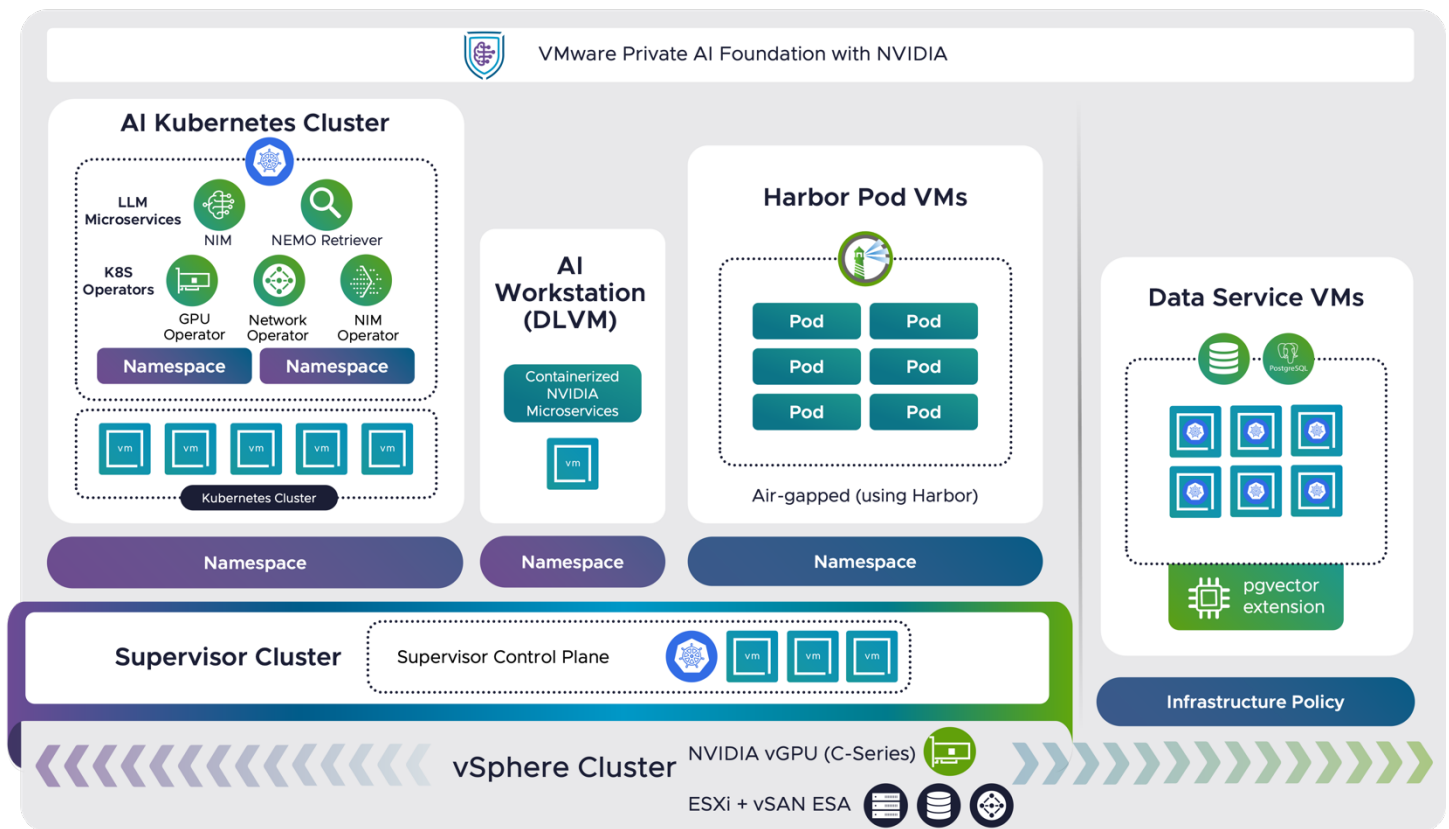
Portgroup Name	Function	VMkernel adapter Required?	MTU	Teaming Policy
VM management	IP connectivity for VM Management traffic	No	1500	Route based on physical NIC
ESXi Management	IP connectivity for ESXi	Yes	1500	Route based on physical NIC
vSAN	vSAN traffic	Yes	9000	Route based on physical NIC
NFS (Optional)	NFS traffic for supplemental storage	Yes	9000	Route based on physical NIC
Host Overlay	NSX Fabric TEP traffic	Yes	9000	Not Applicable
Edge uplinks and overlay	NSX Edge Nodes TEP traffic and uplink fabric traffic	No	9000	Explicit Failover

## Workload domain

### Workload domain architecture

Figure 12 provides an expanded view of the VI workload domain architecture compared to the conceptual overview in Figure 10. At its base, the vSphere cluster enables the Supervisor Cluster, supporting vSphere namespaces that facilitate resource isolation, governance, and policy enforcement. Within these namespaces, users can deploy diverse workloads such as AI Kubernetes clusters, AI workstations, and RAG applications—provisioned and managed by the Supervisor. The Supervisor can also integrate Harbor, a unified repository for container images and AI models, supporting OCI-compatible formats for offline or air-gapped environments. [Data Service Manager \(DSM\)](#) running in the Management Domain uses the vSphere infrastructure policy to deploy Database VMs to the registered VI WLD. These Database VMs then initiate a Kubernetes cluster and run vector database containers within pods. All components are pre-configured or customizable as service catalogs within Aria Automation, ensuring seamless operation in VMware Private AI Foundation with NVIDIA deployments.

Figure 12. Workload Domain virtual architecture



# VMware Private AI Foundation with NVIDIA on HGX Servers: Reference Design for Inference

## Workload domain settings

Table 12 details the functions and configurations of each component within the workload domain. Notably, the settings of the two vSphere Distributed Switches (VDS) and PortGroups in the workload domain mirror those utilized in the management domain.

Table 12. Workload Domain settings

Category	Component	Function	Number	Settings / Notes	Deployment Reference
vSphere	vCenter in Mgmt Domain	Central management UI for configuring and monitoring vSphere infrastructure	1 per VI WLD	<ul style="list-style-type: none"> <li>Deployed by SDDC Manager during VI WLD creation</li> <li>Located within Mgmt Domain</li> <li>Integrated into pre-existing SSO domain in the Mgmt Domain</li> </ul>	<a href="#">VMware Cloud Foundation 5.2 Design Guide</a>
	vSAN	Create a single, shared storage pool across ESXi hosts in a domain	Require at least 4 ESXi hosts	<ul style="list-style-type: none"> <li>Enable ESA with RAID5 policy for performance and storage space efficiency</li> </ul>	<a href="#">vSAN Design for VMware Cloud Foundation</a>
	VDS	Provide virtual networking across multiple ESXi hosts	2 per VI WLD	<ul style="list-style-type: none"> <li>Each VDS configured with 2 uplinks (NICs) to ensure traffic isolation between overlay-based networks and GPU/Data networks on non-overlay networks.</li> <li>The second VDS is prepared for NSX integration</li> </ul>	<a href="#">NSX Reference Design Guide 4.2</a>  <a href="#">vSphere Networking Design for VMware Cloud Foundation</a>
	Host VIBs	Enable NVIDIA vGPU (C-Series) and Fabric Management for NVSwitch and NVlink based systems	2 VIBs per ESXi host	<ul style="list-style-type: none"> <li>Recommend using vLCM images to install in the cluster of VI WLD</li> <li>NVIDIA vGPU (C-Series) Host VIB contains vGPU host driver and NVSwitch fabric management</li> <li>mgmt-daemon VIB for monitoring GPU metrics in Aria Ops</li> </ul>	<ul style="list-style-type: none"> <li>- <a href="#">Installing and configuring the NVIDIA vGPU Manager VIB</a></li> <li>- <a href="#">Deploy a GPU-Accelerated VI Workload Domain for VMware Private AI Foundation with NVIDIA</a></li> </ul>
VKS	Content Library	Centralized repository to store and manage DLVMs' templates, Kubernetes Releases, and OVAs	1 or more	<ul style="list-style-type: none"> <li>Import DLVM template from <a href="#">link</a></li> <li>Use the check-in/check-out feature for version control of VM templates</li> <li>Configured as a subscribed library to synchronize <a href="#">DLVM templates</a> published by VMware</li> </ul>	<a href="#">Create a Content Library with Deep Learning VM Images for VMware Private AI Foundation with NVIDIA</a>
	VM Classes	Define resource allocations for Kubernetes cluster nodes or VMs	As needed	<ul style="list-style-type: none"> <li>Logical entities that specify the resources assigned to VMs or Kubernetes nodes, including CPU, memory, PCI devices, and</li> </ul>	<a href="#">Configure vGPU-Based VM Classes for AI Workloads for VMware Private AI</a>

# VMware Private AI Foundation with NVIDIA on HGX Servers: Reference Design for Inference

Category	Component	Function	Number	Settings / Notes	Deployment Reference
		provisioned by VM Service		Device Groups (vGPUs with NVSwitches, NVLinks, and Virtual Functions on NICs).	<a href="#">Foundation with NVIDIA</a>
	Supervisor Contour Service	Ingress Service for Supervisor Services	Per Supervisor Cluster	<ul style="list-style-type: none"> <li>Default values for Contour supervisor Service.</li> <li>Contour is required to enable the Harbor and Model Store Supervisor Services</li> </ul>	<a href="#">Install Contour as a Supervisor Service</a>
	Supervisor Harbor Service	Serves as a private registry for storing and managing OCI-compliant artifacts, such as container images and AI models.	Per Supervisor Cluster	<ul style="list-style-type: none"> <li>Update harbor-data-values.yml file to specify settings like FQDN and Storage Class.</li> </ul>	<a href="#">Setting Up a Private Harbor Registry in VMware Private AI Foundation with NVIDIA</a>
NSX	NSX Manager in Mgmt Domain	Central management UI of NSX for configuring and monitoring network components in VI WLD	A cluster of 3 appliances per VI WLD	<ul style="list-style-type: none"> <li>Deployed by SDDC Manager during VI WLD creation</li> <li>Separate NSX Manager deployment for isolating permissions, policies, services, etc.</li> </ul>	<a href="#">NSX Design for VMware Cloud Foundation</a>
	NSX Edge Cluster	Grouping of one or more edge nodes, representing a pool of capacity for NSX services	1 or more per VI WLD	<ul style="list-style-type: none"> <li>Profile Type: Set to Default.</li> <li>Network Configuration: Layer 2 (L2) Uniform.</li> </ul>	
	↳ Tier-0 Router	Provides north-south connectivity between the NSX logical network and the physical infrastructure	1	<ul style="list-style-type: none"> <li>Active/Active mode with ECMP routing</li> <li>BGP Configured</li> <li>ASNs defined (local &amp; remote)</li> <li>Large Form Factor</li> </ul>	
	↳ Tier-1 Router	Offers east-west routing functionality, connecting to logical switches for internal network communication and providing stateful network services.	1 or more	<ul style="list-style-type: none"> <li>Configured in Active/Passive mode for VKS enablement</li> </ul>	

## VMware Private AI Foundation with NVIDIA prerequisites

After provisioning the above VI WLD, infrastructure architects must prepare it for the VMware Private AI Foundation with NVIDIA. This preparation involves installing the necessary ESXi VIBs to enable NVIDIA vGPU (C-Series) capabilities and GPU metrics at the host level, licensing the solution, and performing other essential infrastructure tasks as outlined in Table 13. For more information refer to [Preparing VMware Cloud Foundation for a Private AI workload deployment](#).

Table 13. VMware Private AI Foundation with NVIDIA prerequisites

Category	Component	Function	Number	Settings / Notes	Deployment Reference
NVIDIA	License Server	Manage and distribute licenses for NVIDIA software products including vGPU (C-Series)	1	<ul style="list-style-type: none"> <li>Can be deployed by Cloud License Service (CLS) or Delegated License Server (DLS)</li> </ul>	<a href="#">NVIDIA License Server Documentation</a>
vSphere Content Library	DLVM template	Used by VCF Automation when serving requests from self-service catalog	As needed	<ul style="list-style-type: none"> <li>A set of VM images is shipped and maintained by VMware</li> <li>Stored as templates in a VCF content library</li> </ul>	<a href="#">Create a Content Library with Deep Learning VM Images for VMware Private AI Foundation with NVIDIA</a>
Aria Automation	NVIDIA licensing portal API key and client token file	Facilitates interactions with the NVIDIA Licensing Portal to obtain licenses from CLS or DLS	1	<ul style="list-style-type: none"> <li>API Key with the following access types:                             <ul style="list-style-type: none"> <li>Licensing State</li> <li>Software Downloads</li> <li>Enterprise</li> </ul> </li> </ul>	<a href="#">Requirements for Deploying VMware Private AI Foundation with NVIDIA</a>
	NGC API KEY	Authenticate access and download NGC catalog items	1	<ul style="list-style-type: none"> <li>Only NGC Key with NVAIE entitlements can access the assets on NGC.</li> <li>It's recommended to use an NGC API key that is not tied to an individual user's account and is managed by your IT department</li> </ul>	

## Validation

Table 14 presents product demonstrations of VMware Private AI Foundation with NVIDIA, highlighting its capabilities and providing deployment references for validation purposes. It is encouraged you to review these deployment procedures and capabilities to gain a comprehensive understanding of the platform. Additionally, curated official documentation and blogs to assist you in your deployment efforts are provided. For more use cases for implementing VMware Private AI Foundation with NVIDIA, check these two blogs ([part 1](#) and [part 2](#)) and this [talk](#) about a real customer RAG journey.

Table 14. VMware Private AI Foundation with NVIDIA use cases

Use Cases	Deployment Reference
Deploy a DLVM	<ul style="list-style-type: none"> <li>• <a href="#">Deploy a DLVM by Using a Self-Service Catalog</a></li> <li>• Section "Improved Experience for the Data Scientist" in <a href="#">VMware Private AI Foundation with NVIDIA a Technical Overview</a> (Blog)</li> </ul>
Deploying a RAG chatbot in a DLVM ( <a href="#">Product Demo</a> )	<ul style="list-style-type: none"> <li>• <a href="#">Deploy a Deep Learning VM with a RAG Workload</a> (Official doc)</li> <li>• Section "A Chatbot Application that Uses Retrieval Augmented Generation (RAG)" in <a href="#">VMware Private AI Foundation with NVIDIA - a Technical Overview</a> (Blog)</li> <li>• <a href="#">Building Production-grade AI-driven Apps on VMware Private AI Foundation with NVIDIA NIM</a> (Blog)</li> </ul>
Deploy an AI Kubernetes cluster	<ul style="list-style-type: none"> <li>• <a href="#">Deploy a GPU-Accelerated TKG Cluster by Using a Self-Service Catalog Item in VMware Aria Automation</a> (Official doc)</li> </ul>
Deploying a RAG chatbot app in a Kubernetes cluster ( <a href="#">Product Demo</a> )	<ul style="list-style-type: none"> <li>• <a href="#">VMware Private AI Foundation with NVIDIA - a Technical Overview</a></li> </ul>
Model Governance ( <a href="#">Product Demo</a> )	<ul style="list-style-type: none"> <li>• <a href="#">Storing ML Models in VMware Private AI Foundation with NVIDIA</a> (Official doc)</li> <li>• <a href="#">Onboarding Llama3 to the Private AI Model Gallery</a> (Blog)</li> </ul>
Integrate DSM into Service Catalog ( <a href="#">Product Demo</a> )	<ul style="list-style-type: none"> <li>• <a href="#">VMware Data Services Manager Design for Private AI Ready Infrastructure</a> (Official Doc)</li> <li>• <a href="#">Private AI Automation Services Enhancements in VMware Cloud Foundation 5.2.1</a> (Blog)</li> <li>• <a href="#">Step-by-step: Deploy DSM using Aria Automation</a> (Blog)</li> </ul>
GPU Monitoring Dashboard ( <a href="#">Product Demo</a> )	<ul style="list-style-type: none"> <li>• <a href="#">Monitoring VMware Private AI Foundation with NVIDIA</a> (Official doc)</li> <li>• Section "Monitoring GPU Consumption and Availability" in <a href="#">VMware Private AI Foundation with NVIDIA - a Technical Overview</a> (Blog)</li> </ul>



## Performance

To check whether the virtual infrastructure meets your performance requirements or establish a baseline, this section explains how to run and validate the [GenAI-Perf](#) benchmark, which is a convenient command-line tool for measuring the throughput and latency of generative AI models served through an inference server. A performance comparison between a bare metal and virtualized configurations is provided. The results show that the difference between virtual and bare metal fall into a statistically negligible variance. The following instructions are adapted from this step-by-step documentation on [Using GenAI-Perf to Benchmark](#), and are provided as reference examples and do not guarantee performance

### Benchmarking GenAI-Perf

#### Deploy a DLVM

Use the service catalog in Aria Automation to deploy a DLVM with 4x H100-80c vGPU connected via NVSwitch and NVLink. The DLVM's settings, NIM, and Triton-inference-server container used in benchmarking are specified in Table 15.

Table 15. DLVM Settings and Components used in benchmarking

Component	Settings
DLVM	<ul style="list-style-type: none"> <li>• 24 vCPU, 320 GB memory, 256 GB disk</li> <li>• Nvidia:4@nvidia_h100xm-80c%NVLink</li> <li>• Enable UVM by setting pciPassthru[0~3].cfg.enable_uvm = 1</li> <li>• pciPassthru.64bitMMIOSizeGB = 1024 GB</li> </ul>
<a href="https://nvc.io/nim/meta/llama-3.1-70b-instruct">nvc.io/nim/meta/llama-3.1-70b-instruct</a>	<ul style="list-style-type: none"> <li>• 1.3.3</li> <li>• Profile: tensorrt_llm-h100-fp8-tp4-pp1-throughput</li> </ul>
tritonserver	24.10-py3-sdk

## Launch NIM in a DLVM

After deploying the DLVM, the scripts in Figure 12 to launch NIM are followed. To initiate a TensorRT-LLM backed NIM, set the two additional environment variables specified in Table 16. For other NIM configuration parameters, refer to the [NIM's Getting Started](#) and [Configuring a NIM](#) page.

Table 16. NIM Launch Parameters to launch TensorRTLLM-backed NIM

Parameter	Value	Consideration
NCCL_CUMEM_ENABLE	0	Disables NCCL's cuMem allocator to help avoid certain NCCL-related issues, e.g., increased memory overhead during CUDA graph captures.
shm-size	16 GB or higher	Increasing the shared memory size is beneficial for applications that require interprocess communication or large shared memory segments. For instance, Docker containers are allocated 64 MB of shared memory by default

Figure 12 launches a Docker NIM container for the `llama-3.1-70b-instruct` model with `tensorrt_llm-h100-fp8-tp4-pp1-throughput` profile. To use a different NIM profile, replace the `NIM_MODEL_PROFILE` with a different value in the `list-model-profiles` command. To launch a container for a different NIM, replace the value of `Repository` with the value of the other NGC `image list` command using the [NGC CLI tool](#) and change the value of `CONTAINER_NAME` to something appropriate.

Figure 12. TensorRT-LLM backed NIM launch scripts

```

# Export the NGC API key, then Docker login to NGC
export NGC_API_KEY=YOUR_KEY
echo "$NGC_API_KEY" | docker login nvcr.io --username '$oauthtoken' --password-stdin

# (Optional) List available NIMs. This step requires to install NGC CLI.
ngc registry image list --format_type csv nvcr.io/nim/*

# Choose a container name for bookkeeping
export CONTAINER_NAME=Llama-3.1-70b-instruct

# The container name from the previous ngc registry image list command
Repository=nim/meta/llama-3.1-70b-instruct
Tag=1.3.3

# Choose a LLM NIM Image from NGC
export IMG_NAME="nvcr.io/${Repository}:${Tag}"

# Choose a path on your system to cache the downloaded models
export LOCAL_NIM_CACHE=~/.cache/nim
mkdir -p "$LOCAL_NIM_CACHE"

# List the compatible profiles of the NIM
docker run --rm --runtime=nvidia -e NGC_API_KEY --gpus=all $IMG_NAME list-model-profiles

export PORT=8000
export NCCL_CUMEM_ENABLE=0
export NIM_MODEL_PROFILE=tensorrt_llm-h100-fp8-tp4-pp1-throughput

docker run -it --rm --name=$CONTAINER_NAME-$PORT \
  --runtime=nvidia \
  --gpus all \
  -e NGC_API_KEY \
  -v "$LOCAL_NIM_CACHE:/opt/nim/.cache" \
  -u $(id -u) \
  -p $PORT:8000 \
  -e NIM_MODEL_PROFILE \
  -e NCCL_CUMEM_ENABLE \
  --shm-size=16GB \
  $IMG_NAME

```

## Launch GenAI-Perf

Next, while the GPU-enabled NIM container from the previous section remains active and running, interactively launch a Triton inference server container without GPU support. Within this CPU-only container, the GenAI-Perf using the parameters outlined in Table 15 to conduct a warm-up load test on the NIM backend are initiated. For additional model input parameters, refer to [this link](#). Consequently, the GenAI-Perf tool in the CPU-only Triton container sends prompt requests to the GPU-enabled NIM container, enabling validation of its functionality.

Table 17. GenAI-perf launch parameters

Parameters	Value
Input Sequence Length	200
Output Sequence Length	50
Output Sequence Std	10
Concurrency	10

**Note:** With concurrency=N, GenAI-Perf maintains N active inference requests during profiling. For example, with a concurrency of 4, it sustains 4 simultaneous requests, issuing a new request as each one completes.

For understanding additional metrics and parameters to run GenAI-Perf, consult [Metrics](#) and [Parameters and Best Practices](#) page.

Figure 13. Setting Up GenAI-Perf and warm-up load test

```
# Launch a triton inference server container with only CPU
export RELEASE="24.10"
docker run -it --rm --runtime=nvidia \
  --net=host \
  -v $(pwd):/workspace/host \
  nvcr.io/nvidia/tritonserver:${RELEASE}-py3-sdk \
  /bin/bash

# Log in with your Huggingface credential for accessing llama-3 tokenizer
pip install huggingface_hub
huggingface-cli login

# Run GenAI-Perf within triton inference server container for warm up
export INPUT_SEQUENCE_LENGTH=200
export INPUT_SEQUENCE_STD=0
export OUTPUT_SEQUENCE_LENGTH=50
export CONCURRENCY=10

genai-perf profile \
  -m meta/llama-3.1-70b-instruct \
  --endpoint-type chat \
  --service-kind openai \
  --streaming \
  -u localhost:8000 \
  --synthetic-input-tokens-mean $INPUT_SEQUENCE_LENGTH \
  --synthetic-input-tokens-stddev $INPUT_SEQUENCE_STD \
  --concurrency $CONCURRENCY \
  --output-tokens-mean $OUTPUT_SEQUENCE_LENGTH \
  --extra-inputs max_tokens:$OUTPUT_SEQUENCE_LENGTH \
  --extra-inputs min_tokens:$OUTPUT_SEQUENCE_LENGTH \
  --extra-inputs ignore_eos:true \
  --tokenizer meta-llama/Meta-Llama-3.1-70B-Instruct \
  -- \
  -v \
  --max-threads=256
```

## Performance comparison of virtual vs. bare metal

Next, a sweep across varying concurrency levels—1, 2, 5, ..., up to 125 is performed—and results for larger workloads with a Input Sequence Length (ISL) and Output Sequence Length (OSL) pairs to represent a summarization chatbot (ISL=7000, OSL=1000).

Table 18 details the hardware configurations used for running the workload on both bare-metal and virtualized systems. The key distinction is that the virtualized setup uses virtualized NVIDIA H100 GPUs, labeled as H100-80c vGPU.

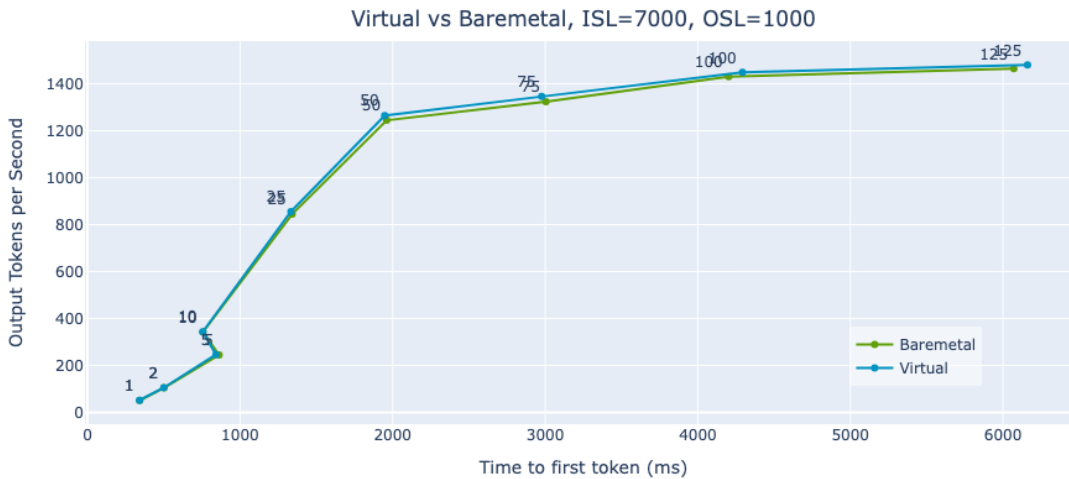
To ensure a fair comparison, NIM with Tensor Parallelism using 4 H100 GPUs is launched. In the bare-metal setup, only 4 out of 8 GPUs are utilized, allowing for a direct comparison to the virtual configuration, which uses 4 NVIDIA vGPUs (C-Series) attached to the VM. On GenAI perf, online mode is used by setting `--streaming` to benchmark online application. The throughput (output tokens per second) for summarization chatbot is then plotted against the Time to First Token (TTFT) in Figure 16.

**Note:** The performance results presented in the following two subsections reflect observed results under specific conditions. Identical performance cannot be guaranteed, as it may vary depending on your hardware, software configuration, or workload. These results are also subject to be improved with future releases or updates to the software stack.

Table 18. Bare metal vs. virtual server configurations for virtualized H100

Component	Bare metal	Virtual
Logical processors	208	24 allocated to the VM <b>(184 available for other VMs/workloads)</b>
GPUs	8x NVIDIA H100-SXM-80GB 4x NVSwitch, NVLink Gen4	4x NVIDIA H100-80c vGPU via NVLink (4@nvidia_h100xm-80c%NVLink)
Memory	2TB	256 GB <b>(1.74 TB for other VMs/workloads)</b>
Storage	8x 3.84TB NVMe SSD	256 GB VM Hard Disk on vSAN
OS	Ubuntu 22.04.3	Ubuntu 22.04.03 DLVM in vSphere 8.0.3
GPU Drivers	Data Center Driver 550.144.03	<ul style="list-style-type: none"> <li>NVIDIA vGPU (C-Series) Host and Guest v17.5</li> <li>550.144.02 (Host), 550.144.03 (Guest)</li> </ul>
CUDA	12.4	12.4
NIM	<ul style="list-style-type: none"> <li>llama-3.1-70b-instruct</li> <li>Tag: 1.3.3</li> <li>Profile: tensorrt_llm-h100-fp8-<b>tp4</b>-pp1-throughput</li> </ul>	<ul style="list-style-type: none"> <li>llama-3.1-70b-instruct</li> <li>Tag: 1.3.3</li> <li>Profile: tensorrt_llm-h100-fp8-<b>tp4</b>-pp1-throughput</li> </ul>
tritonserver	24.10-py3-sdk	24.10-py3-sdk
GenAI-Perf	0.0.11	0.0.11

Figure 14. Throughput and TTFT: virtual vs baremetal across concurrency levels (1~125).



**Plot description:**

- The x-axis represents latency (TTFT).
- The y-axis represents throughput (tokens per second).
- Each point with the same color on the plot corresponds to a measurement from the same underlying model (llama3-70b) and device configuration (type of GPUs in use).
- Points are connected by lines, with the concurrency level indicated (1, 2, ..., 125).

**Interpreting the graph:**

- The optimal points lie closest to the top-left corner of the plot.
- Higher points indicate higher throughput.
- Points further to the left indicate lower latency.

Figure 14 helps illustrate the tradeoff between latency and throughput in inference. A higher concurrency leads to better throughput and, consequently, improved GPU utilization, but it increases latency.

- **For online applications**, where minimizing latency is crucial, it's important to set a maximum acceptable TTFT. For example, if your chatbot requires a TTFT of no more than 1 second, you should select the highest concurrency level that doesn't exceed this threshold.
- **For offline applications**, where latency is less of a concern, a higher concurrency level can be used to maximize throughput.

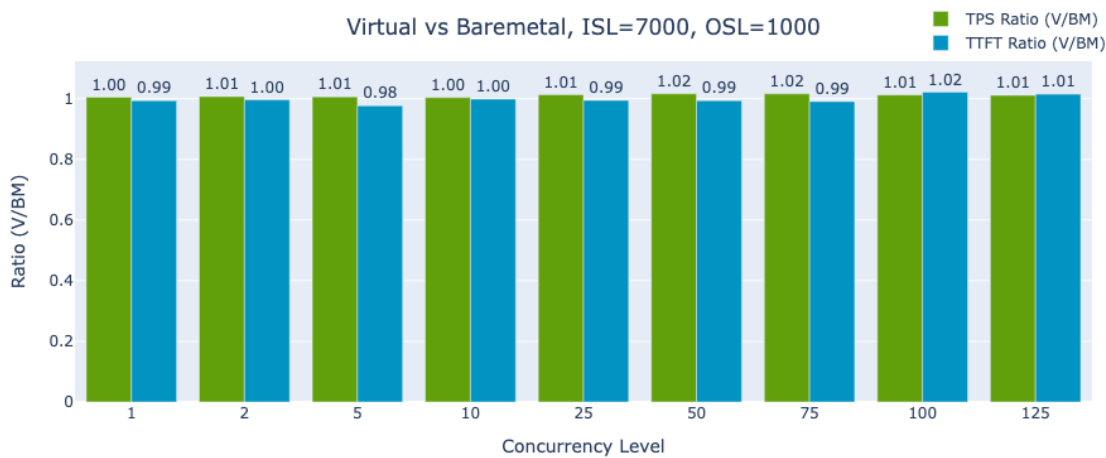
# VMware Private AI Foundation with NVIDIA on HGX Servers: Reference Design for Inference

Figure 15 compares throughput and TTFT ratios between virtual and bare metal setups across varying concurrency levels, evaluating summarization use case with distinct ISL and OSL pairs. For throughput, higher values indicate better performance with virtual GPUs compared to bare metal, and vGPUs deliver 1%~2% better throughput. For TTFT, lower values indicate lower latency, and we find that at certain concurrency levels, virtual GPUs exhibit 1%~2% lower latency, while at other concurrency levels, NVIDIA vGPU (C-Series) show up to 2% higher latency.

Additionally, only 24 out of the total 208 logical CPU cores were utilized for inference in both configurations, leaving the remaining 184 logical CPU cores available for other tasks in the data center. Similarly, we used only 256 GB of CPU memory for the inference workload, while reserving 1.74 TB for other applications, highlighting the benefits of resource isolation.

From this, we conclude that NVIDIA vGPU (C-Series) offer near-bare-metal performance for this workload, positioning them in the "Goldilocks Zone." This means they strike an ideal balance between delivering near-native performance and offering the advantages of easier data center management and cost savings.

Figure 15. Throughput and TTFT Ratio: virtual vs. bare metal across concurrency levels (1~125)

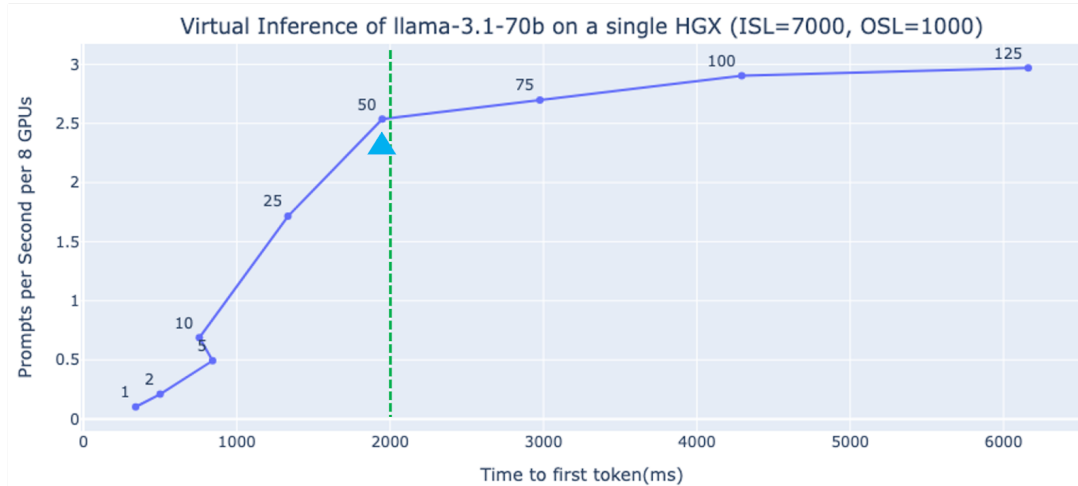




## Inference sizing guidance

Based on the GenAI performance benchmarks, you can estimate your system's capability for specific use cases. For example, a summarization chatbot with an input sequence length (ISL) of 7000 and an output sequence length (OSL) of 1000 requires a Time to First Token (TTFT) of under 2 seconds. Using a single HGX server with 8x H100 GPUs, you could deploy two DLVMs, each utilizing 4 NVIDIA vGPUs (C-Series). According to the benchmark data, this configuration achieves a peak performance of approximately 2.54 prompts per second on an 8x H100 setup, as shown in Figure 16.

Figure 16. Throughput (prompts per second per 8 GPUs for ISL=7k, OSL=1k) on a single HGX server



Assuming an 8-hour working day, this results in 73.2k requests per day ( $2.54 \times 28,800$ ). If each user sends 3 requests daily, this system could support approximately 24.4k daily active users ( $73.2k \div 3$ ). This workload would process around 512 million input tokens ( $7000 \times 73.2k$  requests) and 73 million output tokens ( $1000 \times 73.2k$  requests) per day. Similarly, for a different use case with (ISL=200, OSL=200) is also listed in Table 17. In both cases, the number of daily active users is served with only one HGX H100 server. To support more users, additional replicas across multiple HGX servers working in parallel would be required.

These estimates demonstrate how the benchmark results can be extrapolated to assess capacity for real-world applications, such as chatbots, based on your infrastructure and performance requirements. Additionally, the data can be used to conduct a Total Cost of Ownership (TCO) analysis—calculating costs per input/output token—by factoring in your total on-prem infrastructure expenses, including server costs, hosting costs, and software licensing.

## VMware Private AI Foundation with NVIDIA on HGX Servers: Reference Design for Inference

Table 19. Estimation of Daily Active User using 1x HGX server

Metric	7000 in, 1000 out	200 in, 200 out
Input Tokens	7000	200
Output Tokens	1000	200
TTFT Requirement	< 2s	< 500 ms
Peak Prompts/s per 8x H100 NVIDIA vGPU (C-Series)	2.54	48.1
Concurrency level	50	125
Requests (k) per Day (8h)	73.2	1382.4
Requests per Person	3	3
Daily Active Users (k)	24.4	460.8
Total Input Million Tokens/Day	512.1	276.5
Total Output Million Tokens/Day	73.2	276.5

If you haven't purchased GPU servers yet, we recommend reviewing VMware's [LLM Inference Sizing and Performance Guidance](#) blog. It provides hardware-based metrics for theoretical calculations, independent of benchmark results. Additionally, you can contact VMware's Professional Services team, who can offer tailored recommendations and support to ensure your AI deployment is efficient and cost-effective.

## Conclusion

The integration of VMware Private AI Foundation with NVIDIA) on NVIDIA-Certified HGX Systems offers a robust and scalable infrastructure to address the growing demands of AI inference workloads in enterprise environments.

By optimizing GPU utilization and offering a cloud-like interface for data scientists, this architecture empowers teams to maximize their resources while maintaining governance and compliance. The seamless integration of VMware Cloud Foundation and VMware Private AI Foundation with NVIDIA ensures that IT teams can efficiently manage infrastructure and enforce security policies while offering data scientists the autonomy to innovate and focus on AI model development.

The detailed deployment considerations, product validation, and benchmarking results presented in this paper offer organizations the best of both performance and governance in the private cloud environment.

## Additional information

- [VMware Private AI Foundation with NVIDIA webpage](#)
- [Three Reasons Customers Choose VMware Private AI from Broadcom - Tech Field Day](#)
- [VMware Private AI Foundation Capabilities and Features Update from Broadcom - Tech Field Day](#)
- [AI Model Security and Governance - Broadcom VMware Private AI Model Gallery Demo - Tech Field Day](#)
- [Real-World Customer Journey with VMware Private AI from Broadcom - Tech Field Day](#)
- [DeepSeek-R1 Now Live With NVIDIA NIM](#)

## About the authors

**Dr. Yuankun Fu** is a performance engineer at Broadcom focusing on optimizing AI and HPC performance.

**Agustin Malanco** is a solution architect at Broadcom focusing on AI platforms and solutions.

**Dr. Ramesh Radhakrishnan** is a distinguished engineer leading the AI Platforms and Solutions team at Broadcom.

## Acknowledgments

The authors thank Justin Murray, Vrushal Dongre, Shobhit Bhutani, and Julie Brodeur from Broadcom's VMware Cloud Foundation division and Joe Cullen from NVIDIA for reviewing and improving the paper.



