# Next Era of AI - AI Factories

Frank Lin, NVIDIA Senior Solutions Architect
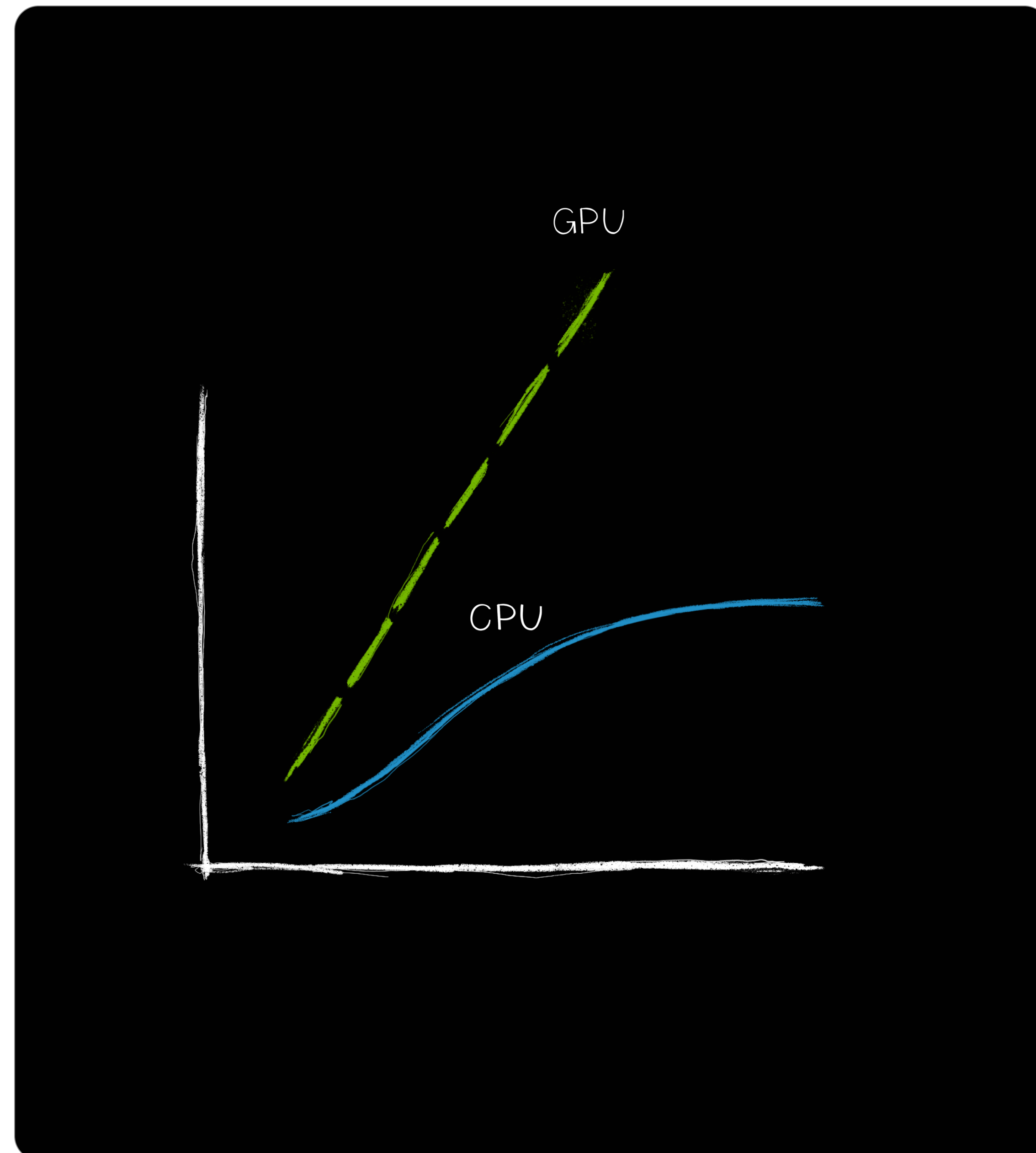
# Next Wave of AI



PERCEPTION AI
SPEECH RECOGNITION
DEEP RECSYS
MEDICAL IMAGING

2012 ALEXNET

GENERATIVE AI
DIGITAL MARKETING
CONTENT CREATION

AGENTIC AI
CODING ASSISTANT
CUSTOMER SERVICE
PATIENT CARE

PHYSICAL AI
CODING ASSISTANT
CUSTOMER SERVICE
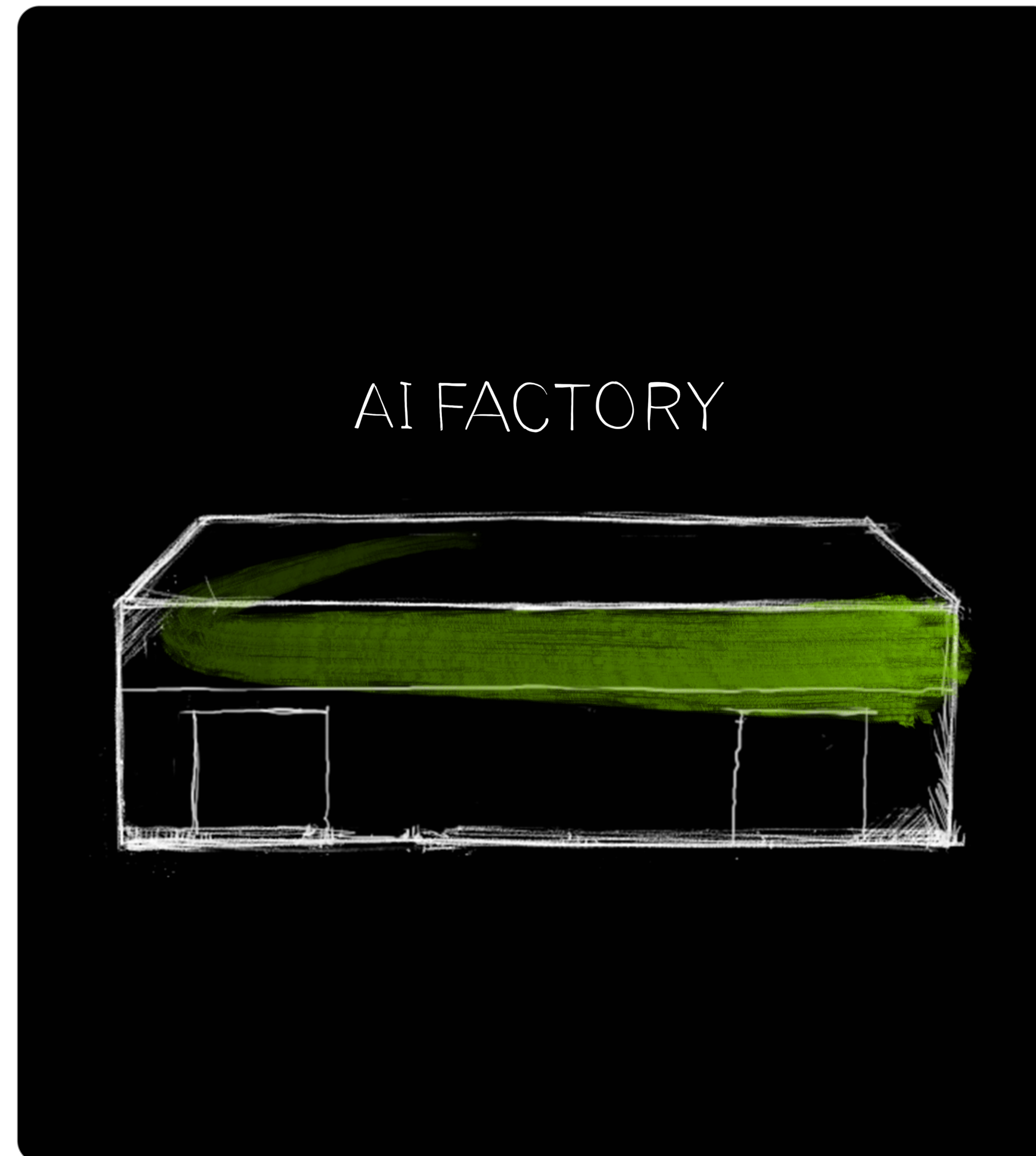PATIENT CARE

NVIDIA.

# Three Major Trends in Computing



Traditional → Accelerated



Data Centers → AI Factories



Generative AI → Physical AI

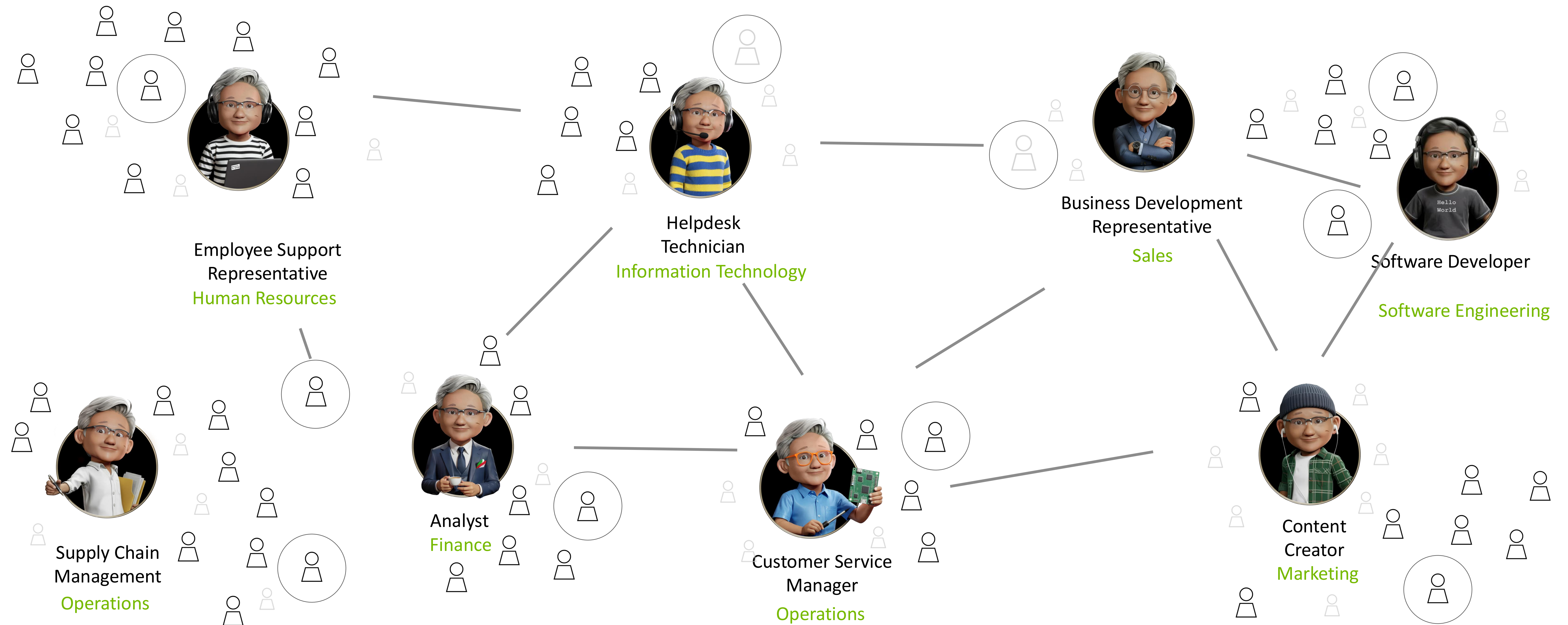# Trillion-Dollar Global IT Investment Shifting to AI Factories

**92%** of enterprises investing in AI

**50%** will use AI agents to achieve business value

**33%** find complexity top barrier for adoption
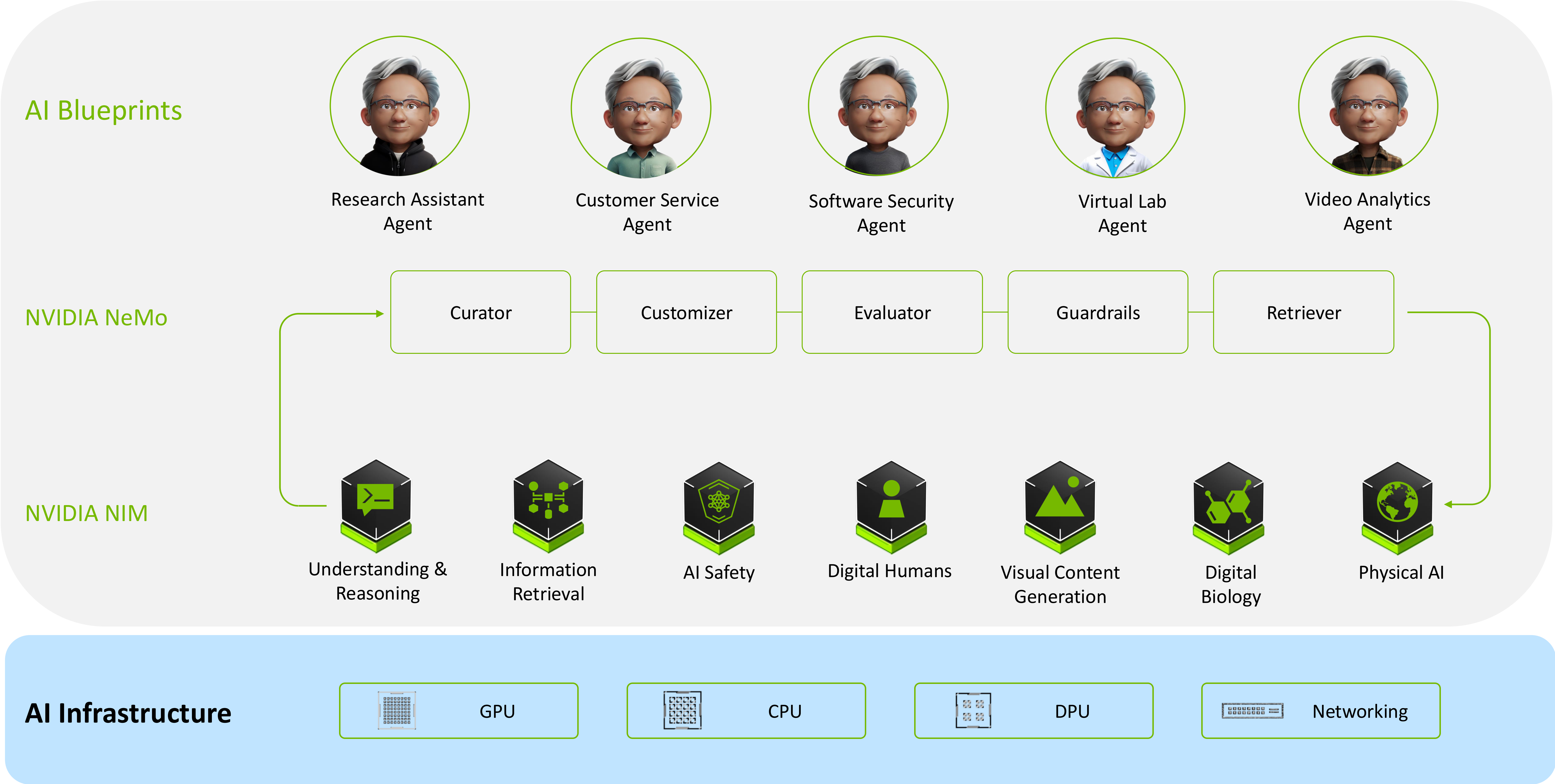
**1%** have mature AI deployments



**NVIDIA.**

# Agents Work Together to Solve Complex Problems



**Employee Support Representative**
Human Resources

**Helpdesk Technician**
Information Technology

**Business Development Representative**
Sales

**Software Developer**
Software Engineering

**Supply Chain Management**
Operations

**Analyst**
Finance

**Customer Service Manager**
Operations

**Content Creator**
Marketing

NVIDIA.

# NVIDIA Provides the Building Blocks for Agentic AI

**AI Blueprints**

Research Assistant Agent

Customer Service Agent

Software Security Agent

Virtual Lab Agent

Video Analytics Agent

**NVIDIA NeMo**

| Curator | Customizer | Evaluator | Guardrails | Retriever |

**NVIDIA NIM**

Understanding & Reasoning

Information Retrieval

AI Safety

Digital Humans

Visual Content Generation

Digital Biology

Physical AI

**AI Infrastructure**

GPU

CPU

DPU

Networking

NVIDIA

# AI-Q: AI Agent Interface to Enterprise Data Stores

## AI Data Platform Reference Design



Compute
RTX PRO 6000 Blackwell Server Edition
NVIDIA BlueField

Networking
NVIDIA Spectrum-X SN5600

Storage
NVIDIA BlueField

NVIDIA AI-Q

Critique

Plan

Prompt

Report

NVIDIA Llama
Nemotron Reason

Tool Use

Calculator

Web
Search

Semantic
Query

SQL
Query

Generate
Summary

# AI Factory Output Drives Revenue

## High throughput multiplied by high interactivity = total token output

Max Perf/W -> Max Revenue

Max Perf/W/TCO -> Max Gross Margin

Token Volume

Factory Tokens per second

50X

Hopper

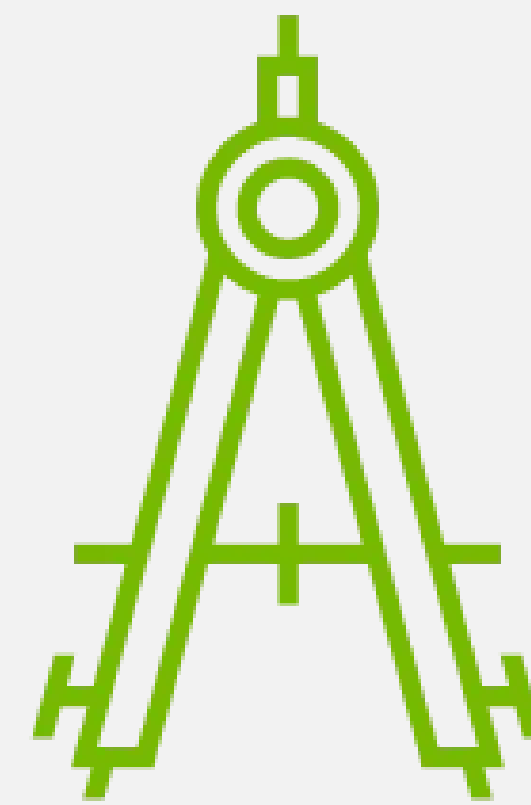Blackwell

TPS for 1 User          $ / M Token

NVIDIA.

# Today's AI Challenges Require AI Factories

## AI workloads require optimized full stack solutions

### Design Complexity

Spans project prioritization,
data acquisition,
and infrastructure

### Deployment and Cost

Infrastructure, security,
and customization

### Time to Value

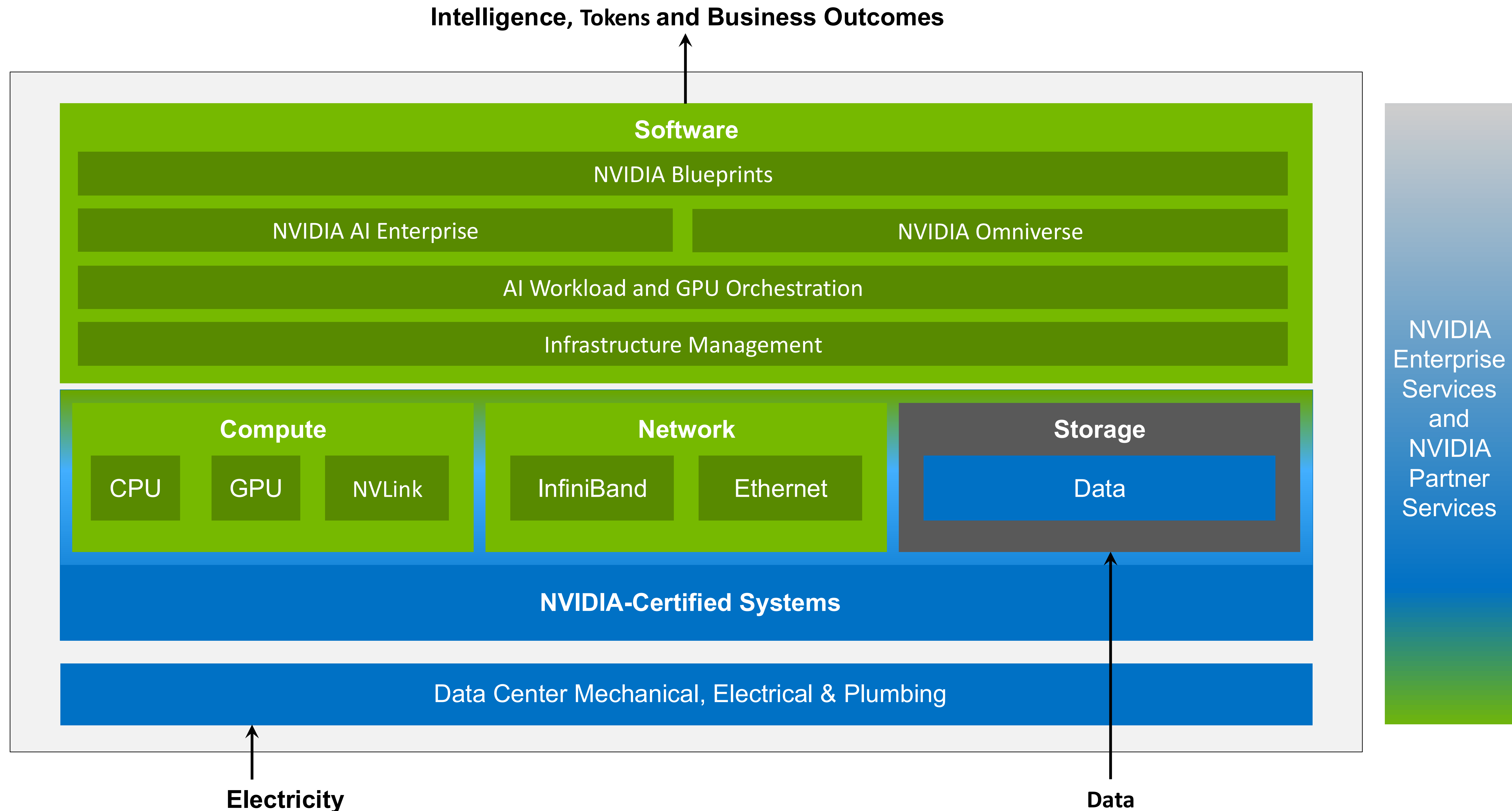Resource management,
time-to-first-train,
time-to-inference

NVIDIA.

# Enterprise Reference Architectures

# NVIDIA Provides a Full Stack for AI Factories

### Built on NVIDIA customer-validated data center reference architectures

**Intelligence, Tokens and Business Outcomes**

**Software**

NVIDIA Blueprints

| NVIDIA AI Enterprise | NVIDIA Omniverse |
|---|---|

AI Workload and GPU Orchestration

Infrastructure Management

| **Compute** | **Network** | **Storage** |
|---|---|---|
| CPU GPU NVLink | InfiniBand Ethernet | Data |

**NVIDIA-Certified Systems**

Data Center Mechanical, Electrical & Plumbing

NVIDIA Enterprise Services and NVIDIA Partner Services

**Electricity**

**Data**

# Every Enterprise Will Have an AI Factory

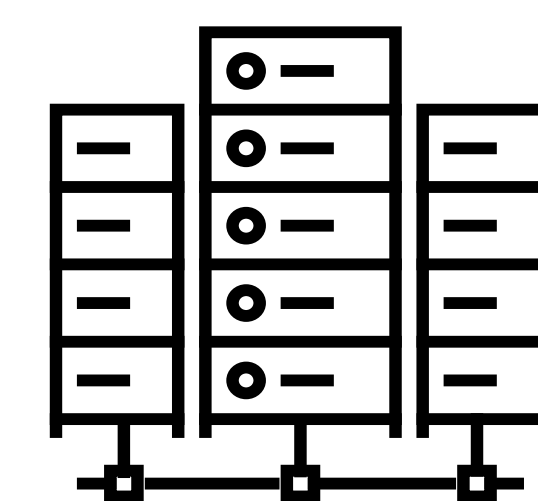| Enterprise Owned | Enterprise Rented |
|---|---|

### Enterprise AI Factory

Enterprise AI Factory Validated
Design

**NVIDIA Enterprise RAs**
NVIDIA Compute
NVIDIA Networking
Partner Ecosystem

### NCP AI Factory

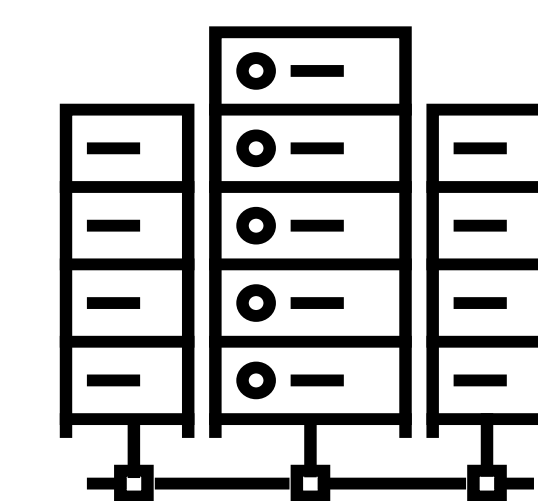NCP Software Stack
NVIDIA AI Enterprise

NVIDIA NCP RAs
NVIDIA Compute
NVIDIA Networking
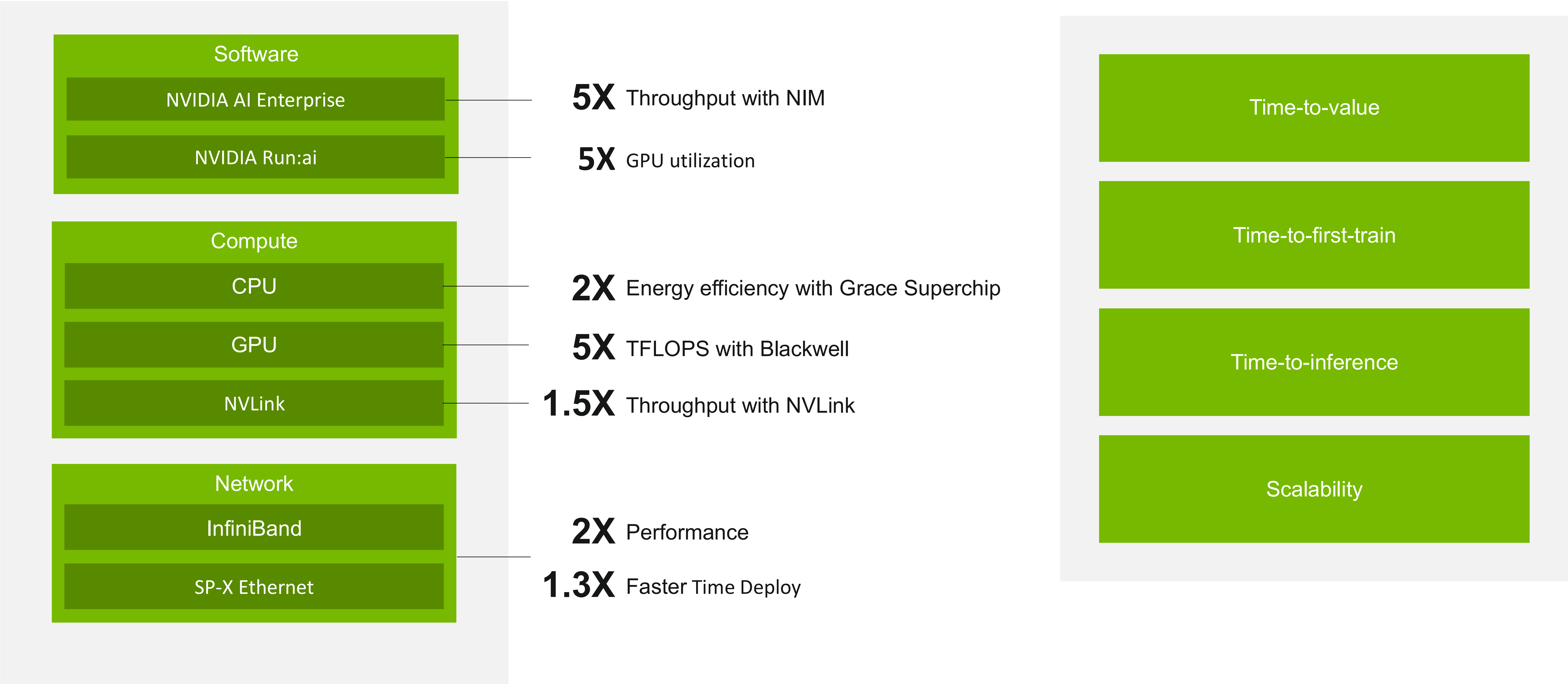Partner Ecosystem

### CSP AI Factory

CSP Software Stack
NVIDIA AI Enterprise

CSP Designed Infrastructure
NVIDIA Compute
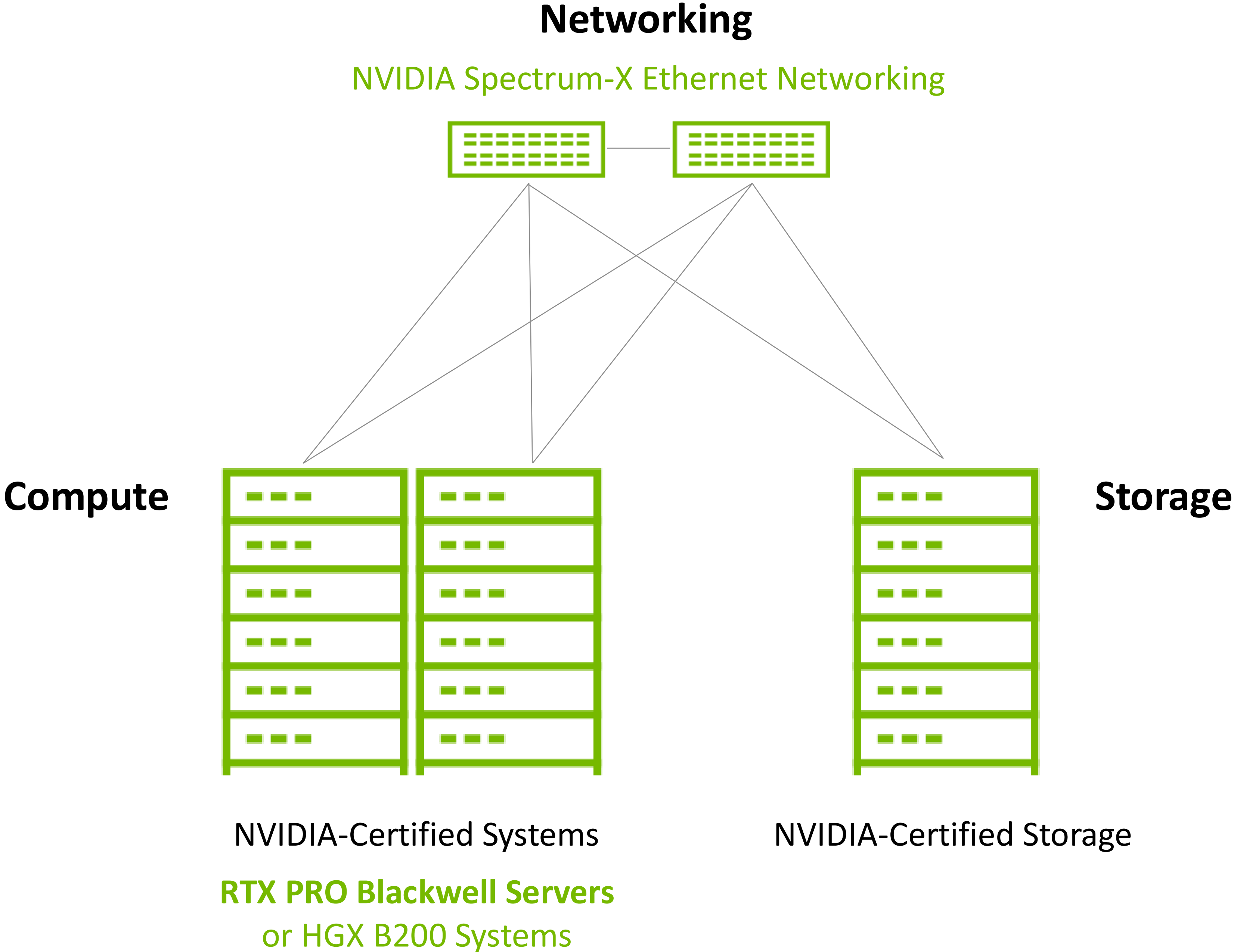NVIDIA Networking
Partner Ecosystem

# AI Factory Compounds Benefits Of NVIDIA Reference Architectures

**Software**

NVIDIA AI Enterprise ——— **5X** Throughput with NIM

NVIDIA Run:ai ——— **5X** GPU utilization

**Compute**

CPU ——— **2X** Energy efficiency with Grace Superchip

GPU ——— **5X** TFLOPS with Blackwell

NVLink ——— **1.5X** Throughput with NVLink

**Network**

InfiniBand ——— **2X** Performance

SP-X Ethernet ——— **1.3X** Faster Time Deploy

Time-to-value

Time-to-first-train

Time-to-inference

Scalability

Up to these speeds, may vary based on input/output and model

NVIDIA

# Building on NVIDIA Enterprise Reference Architectures

**Networking**

NVIDIA Spectrum-X Ethernet Networking

**Compute**

**Storage**

NVIDIA-Certified Systems

**RTX PRO Blackwell Servers**
or HGX B200 Systems

NVIDIA-Certified Storage

Time to value

Scalability

Manageability

Security

# Introducing: Enterprise Reference Architectures

**Comprehensive full-stack design recommendations for building high-performance, scalable data center infrastructure.**

- NVIDIA-Certified Systems
  - Optimized Scale-Up & Scale-Out Configurations
- High-Performance AI Networking
  - Spectrum-X
- AI Software Stack
  - NVIDIA AI Enterprise

**Deployment Guides for Multiple Workloads**

- LLM, RAG, NIM, and NIM Agent Blueprints
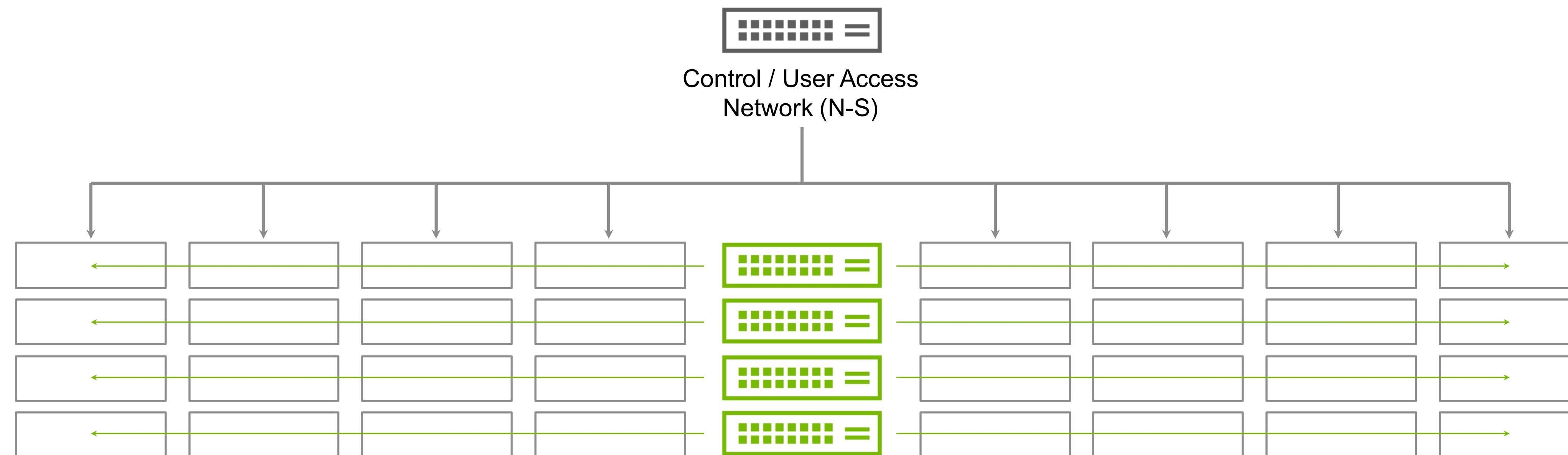
**Flexible Sizing for Expansion Needs**

- Multiple discrete design points for 32 - 512 GPUs
  - Optimized resource utilization



NVIDIA.

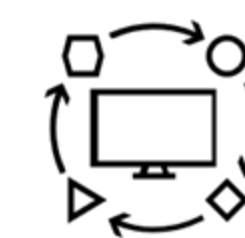# Only NVIDIA Networking Delivers the Fabric for AI Factories

Tightly coupled | Isolated | High-bandwidth | Low jitter | Predictable performance at any scale

Control / User Access
Network (N-S)

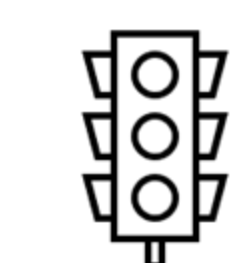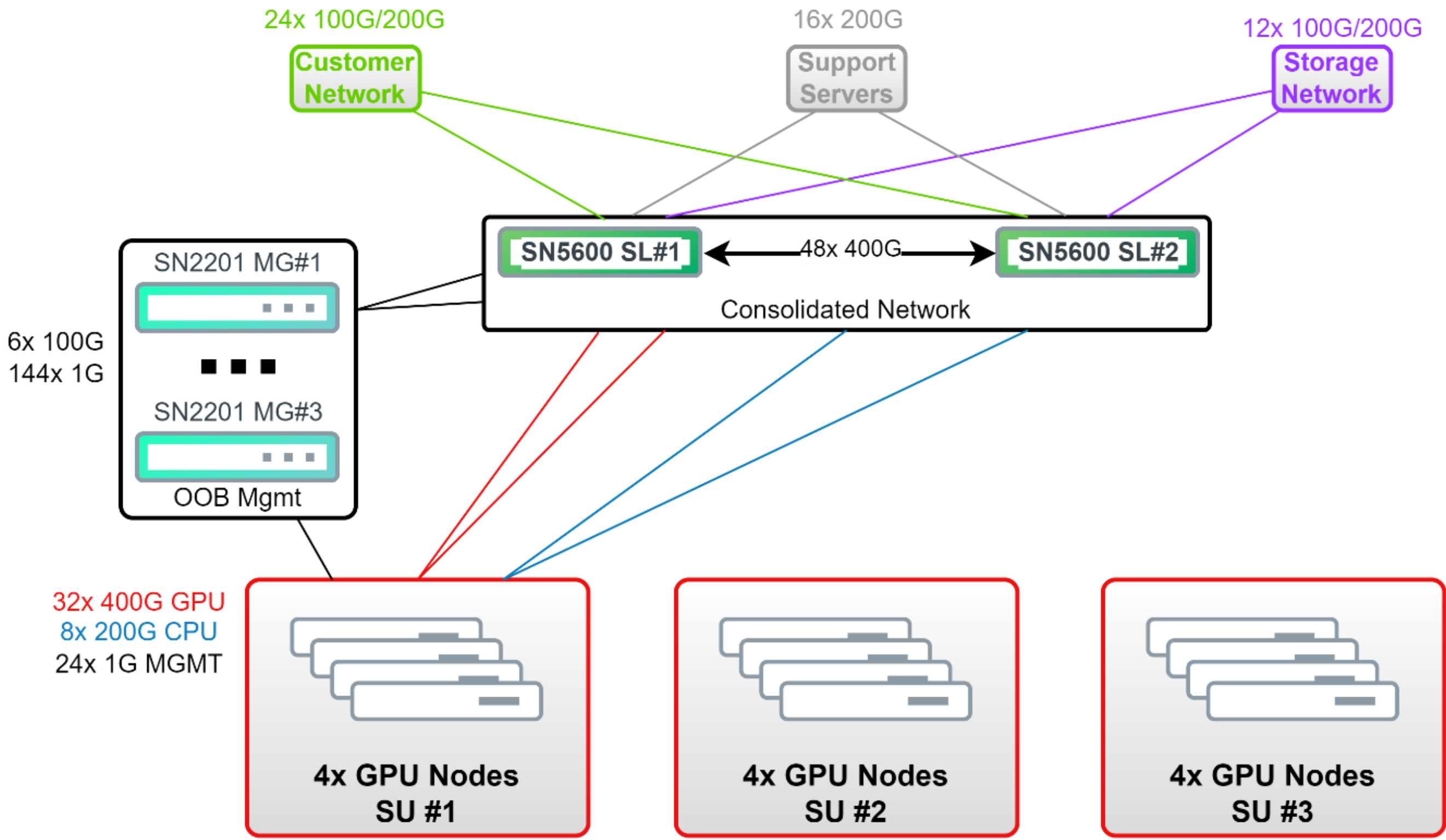| Control / User Access Network (N-S) | AI Fabric (E-W) |
|---|---|
| Loosely-coupled applications, no isolation required | Tightly-coupled processes, tenant isolation required |
| TCP (low bandwidth flows and utilization) | RDMA (high bandwidth flows and utilization) |
| High jitter tolerance | Low jitter tolerance |
| Heterogeneous traffic, statistical multi-pathing | Bursty network capacity, predictive performance |

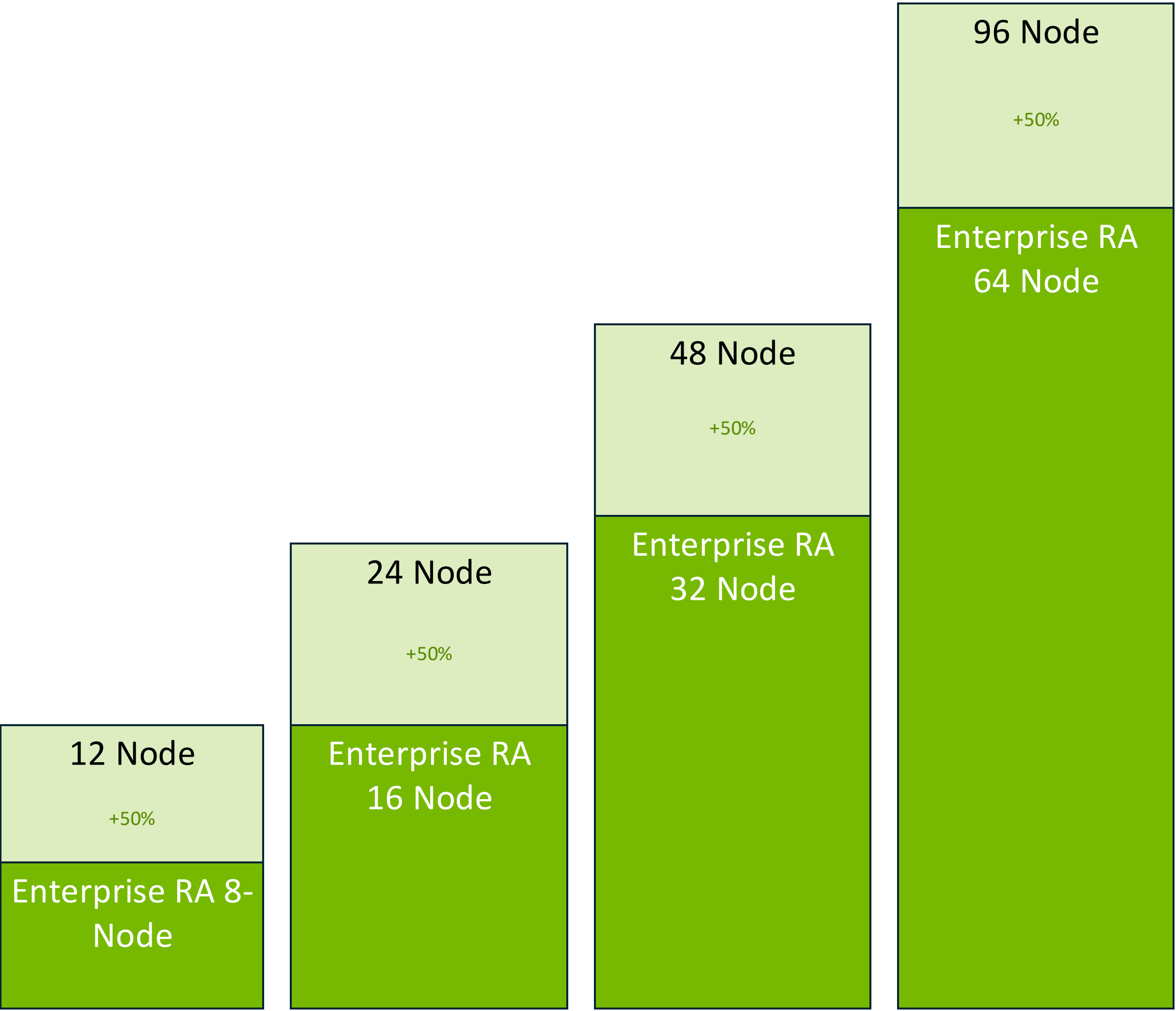NVIDIA.

# Extensible Designs Enable Efficient Scaling

## Enterprise RA Networking w Spectrum-X



**Rail-Optimized Topology**

24x 100G/200G
Customer Network

16x 200G
Support Servers

12x 100G/200G
Storage Network

SN2201 MG#1

SN5600 SL#1 ←→ 48x 400G ←→ SN5600 SL#2

Consolidated Network

6x 100G
144x 1G

SN2201 MG#3

OOB Mgmt

32x 400G GPU
8x 200G CPU
24x 1G MGMT

4x GPU Nodes
SU #1

4x GPU Nodes
SU #2

4x GPU Nodes
SU #3

Example : HGX H100 3SU Cluster, 96 GPUs
1SU = 4 Nodes

**Efficient scaling at multiple design points**

96 Node
+50%
Enterprise RA
64 Node

48 Node
+50%
Enterprise RA
32 Node

24 Node
+50%
Enterprise RA
16 Node

12 Node
+50%
Enterprise RA 8-Node

# Resources

## Enterprise Reference Architecture Announcement

| Resources |
| --- |

**Enterprise RA**
- Announcement Blog
- Web
- Enterprise RA Whitepaper

**NVIDIA-Certified**
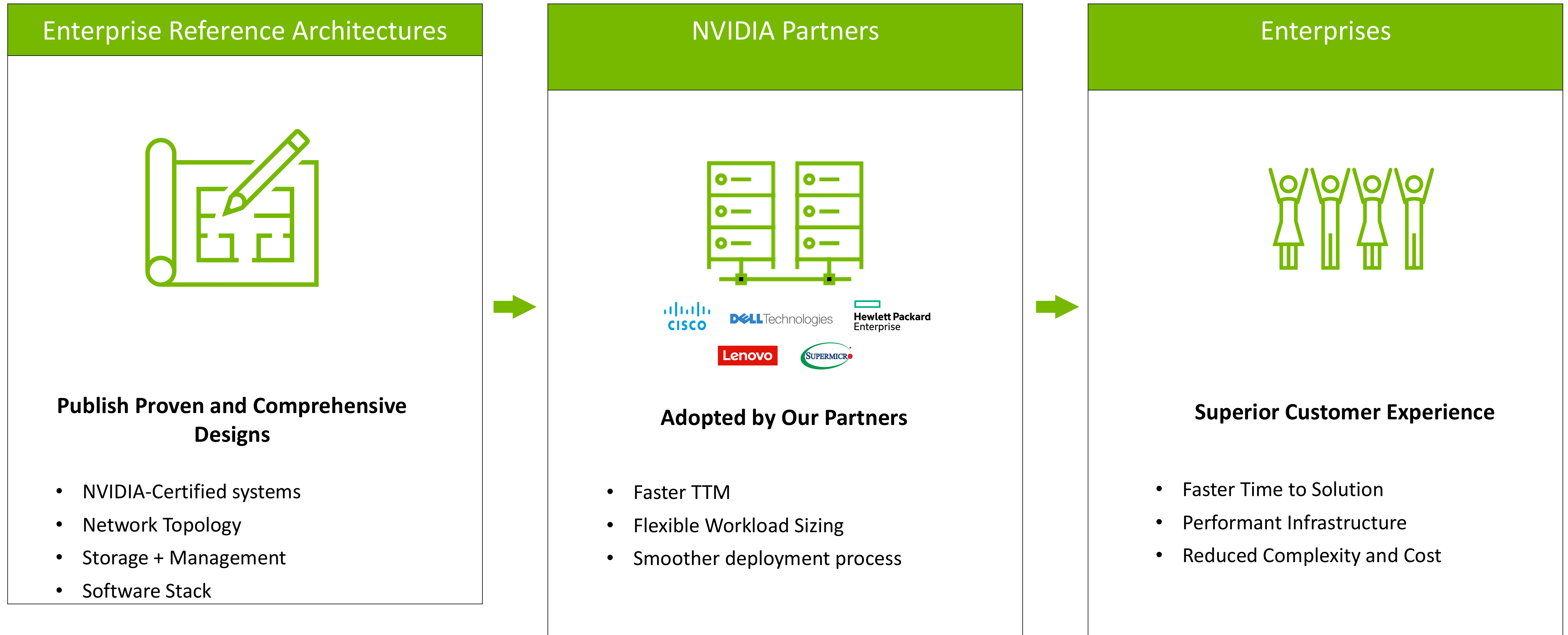- Whitepaper
- Data Sheet

# Helping Enterprises Build AI Factories

Enterprise Reference Architecture

## Enterprise Reference Architectures



**Publish Proven and Comprehensive Designs**

- NVIDIA-Certified systems
- Network Topology
- Storage + Management
- Software Stack

## NVIDIA Partners



CISCO    DELL Technologies    Hewlett Packard Enterprise

Lenovo    SUPERMICR⊙

**Adopted by Our Partners**

- Faster TTM
- Flexible Workload Sizing
- Smoother deployment process

## Enterprises



**Superior Customer Experience**

- Faster Time to Solution
- Performant Infrastructure
- Reduced Complexity and Cost

NVIDIA.

# Understanding DGX BasePOD Vs DGX SuperPOD

## NVIDIA DGX BasePOD

**Scalable, foundational architecture**

## NVIDIA DGX SuperPOD

**Physical twin of NVIDIA's infrastructure**

VS

| | |
|---|---|
| • Flexible reference architectures | • Turnkey data center **product** |
| • Powered by Base Command | • Powered by Base Command |
| • Validated against key **benchmarks** | • **Certified performance** for the most **complex workloads** |
| • Foundation for **partner branded offerings** | • **No customization, no partner re-branding** |

**I need:**
• Choice of flexible vs performance optimized designs
• Inclusion of non-SuperPOD certified storage

**I need:**
• A replica of NVIDIA infrastructure
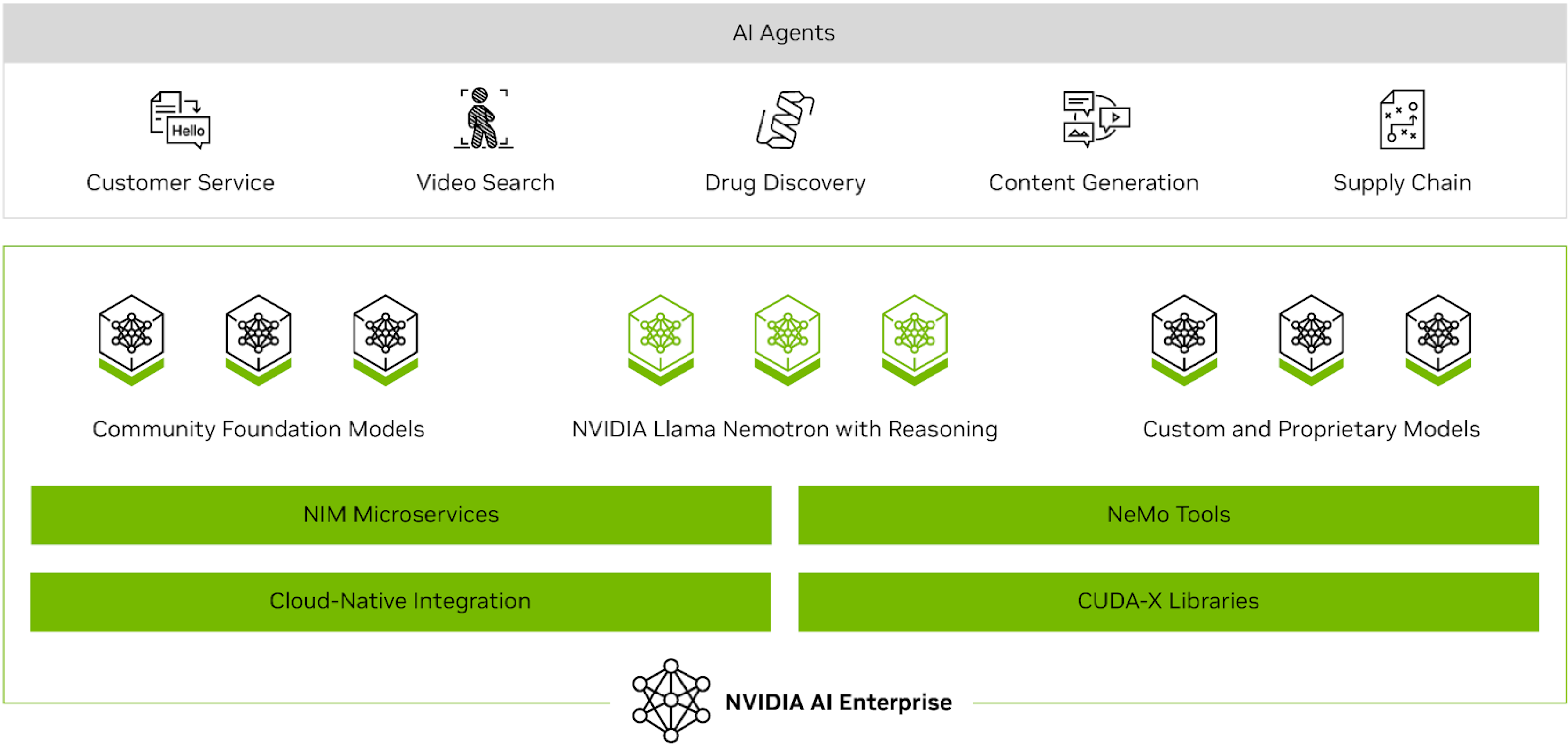• A turnkey deployment
• Full-stack support of the entire deployment

# Fastest Path to Production

# NVIDIA AI Enterprise

What's Inside?



**Comprehensive collection of preconfigured NIM microservices** for efficient inferencing of state-of-the-art foundation models for any use case.

**Powerful, ready-to-use NeMo training, evaluation, and guardrailing tools** and RAG building blocks to accelerate time to deployment.
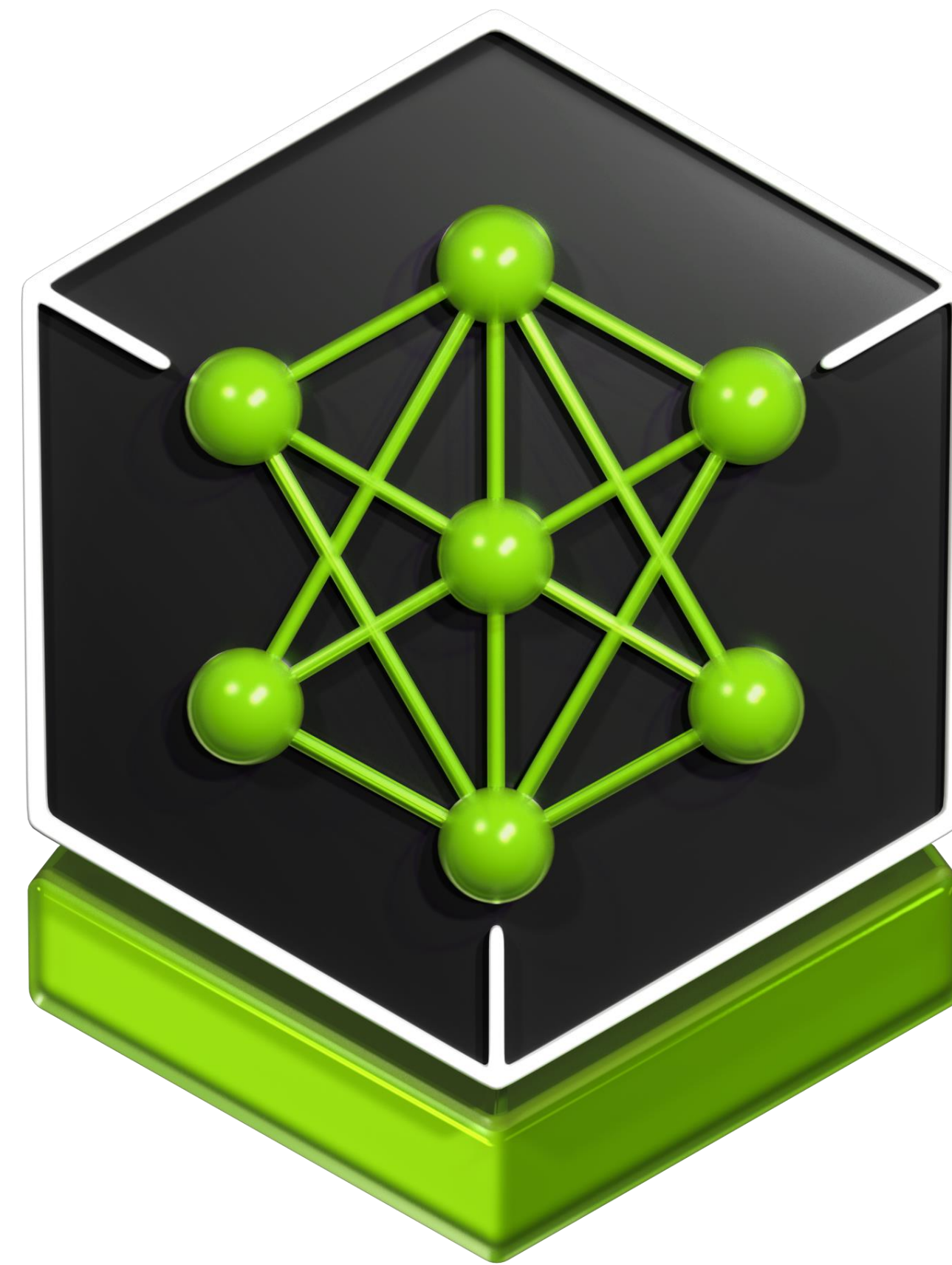
**Infrastructure software** to help manage AI clusters at scale, across the edge and data center, both bare-metal and virtualized.

## Onboard Your AI Agents!

# NVIDIA NIM: Optimized AI Models Run Up to 5X Faster

## Community models – partner models – NVIDIA models



NVIDIA INFERENCE MICROSERVICE

Pre-Trained AI Models
Packaged and Optimized to Run Across
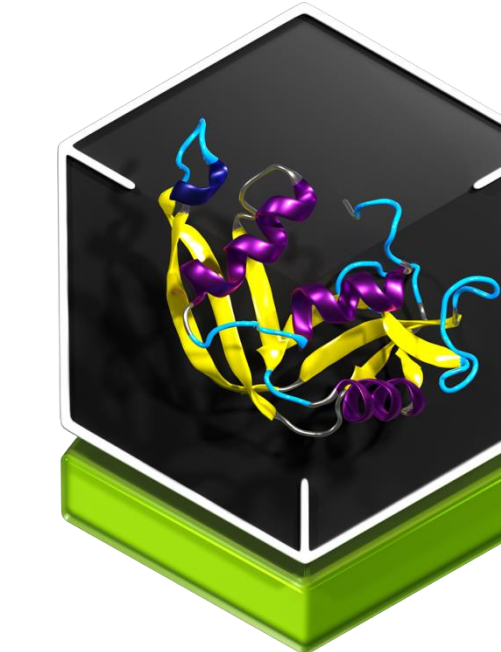CUDA Installed Base

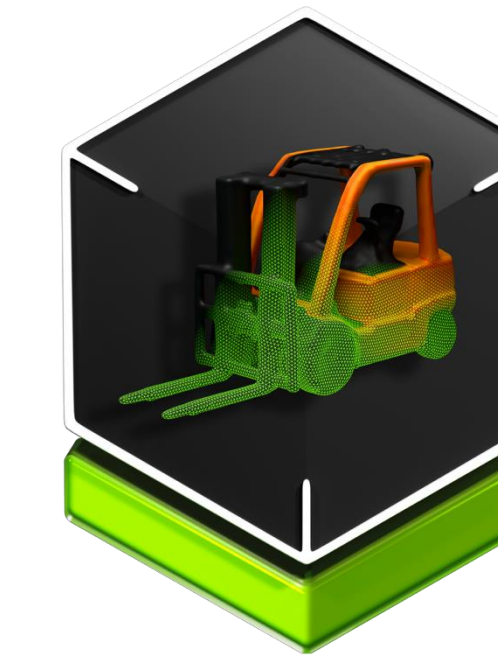Speech     Digital Human     Computer Vision     Biology     Simulation

Language     Regional Language     Vision Language     RAG

ADEPT    gettyimages    Google    Meta    MIT    MISTRAL AI_    NVIDIA    shutterstock    snowflake

# NVIDIA NIM Optimized Inference Microservices

## Rapidly deploy reliable building blocks for accelerated generative AI anywhere



**Portable** Run cloud-native microservices anywhere, maintaining security and control of data and apps

**Easy to Use** Move fast with the latest agentic AI building blocks for reasoning, retrieval, images and more, deployed in minutes with standard APIs

**Enterprise Supported** Gain confidence with stable APIs, quality assurance, continuous updates, security patching, and support

**Performance** Optimize accuracy, latency and throughput to meet requirements with lowest TCO

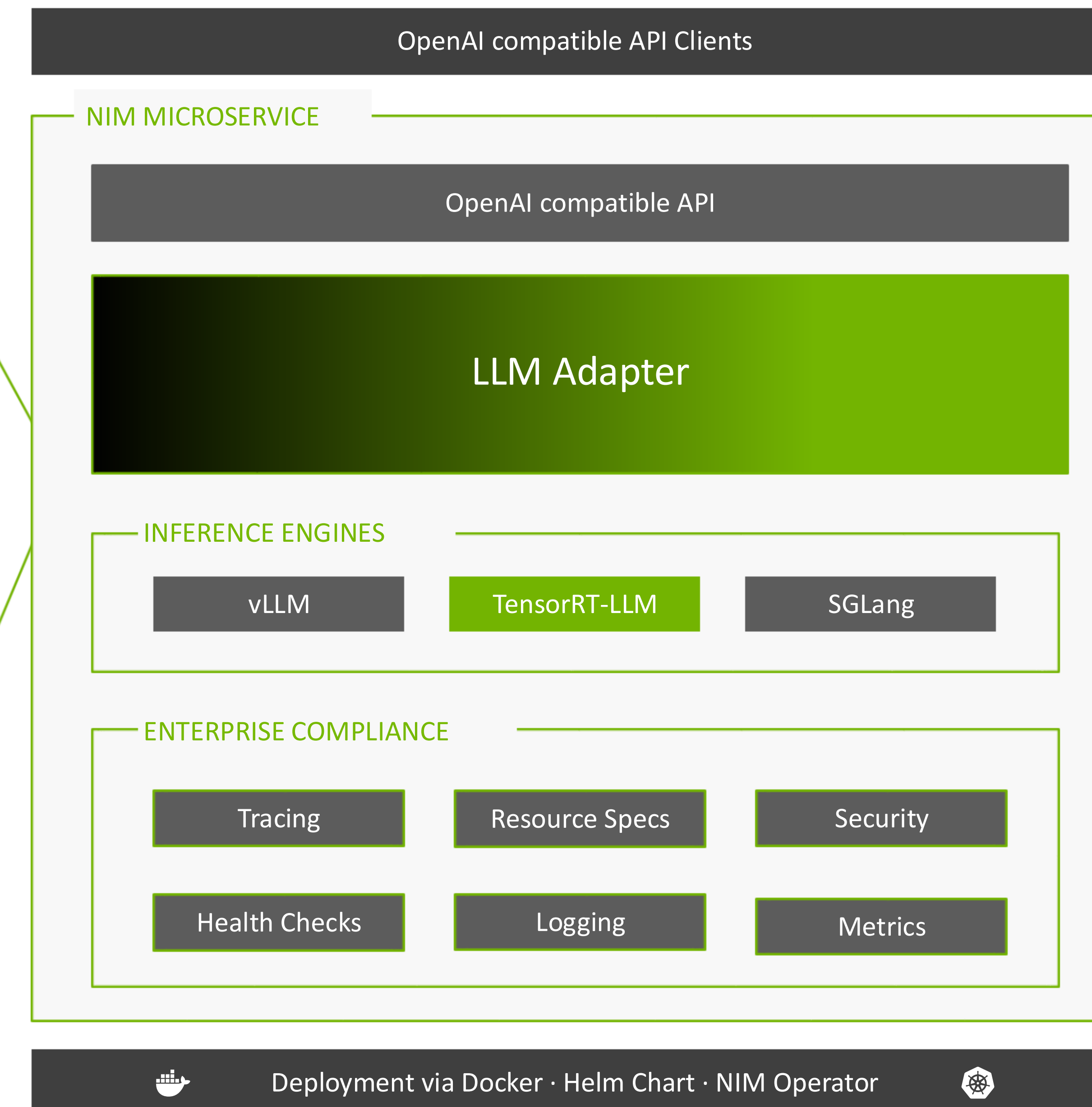# Enterprise-ready Inference for a World of LLMs

Single NIM container with multiple inference backends for rapid, reliable deployment of a broad range of LLMs

## 100K+

public and private LLMs

> `Docker run nim_container [your_LLM]`

vLLM     NVIDIA    SGL

---

**OpenAI compatible API Clients**

**NIM MICROSERVICE**

OpenAI compatible API

**LLM Adapter**

**INFERENCE ENGINES**

| vLLM | TensorRT-LLM | SGLang |

**ENTERPRISE COMPLIANCE**

| Tracing | Resource Specs | Security |
| Health Checks | Logging | Metrics |

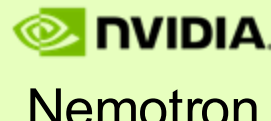Deployment via Docker · Helm Chart · NIM Operator

NVIDIA

# NVIDIA NIM is the Fastest Path to AI Inference

Reduces engineering resources required to deploy optimized, accelerated models

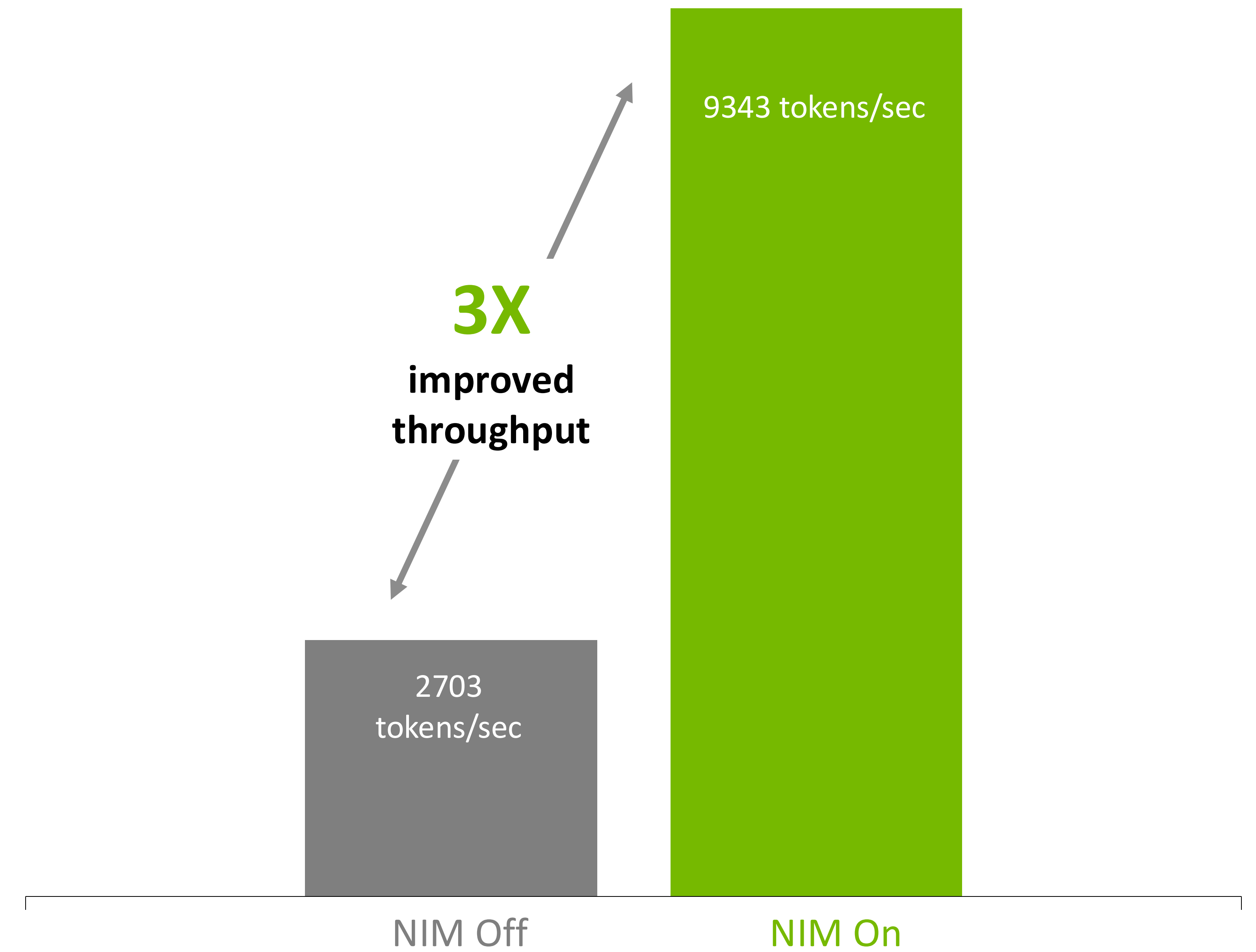| | NVIDIA NIM | Do It Yourself |
|---|---|---|
| Deployment Time | 5 minutes | 1 week + |
| API Standardization | Industry standard protocol<br>OpenAI for LLMs, Google Translate for Speech | Implement the API layer for each domain and model family according to industry standard specifications |
| Optimized Engines | Pre-built engines for NVIDIA and community models<br>MISTRAL AI_    Meta    starcoder    NVIDIA Nemotron | Build your own engine and manually customize for workload and hardware specific requirements |
| Pre and Post Processing Pipelines | Pre-built with optimized pipeline engines to handle pre/post processing (tokenization) | Implement custom logic |
| Model Server Deployment | Automated | Manual setup and configuration |
| Customization | LoRA is supported, more planned | Create custom logic |
| Container Validation | Extensive workload specific QA support matrix validation | No validation |
| Enterprise Support | Delivered with NVIDIA AI Enterprise<br>Security and CVE scanning/patching and tech support | Self supported |

NVIDIA

# Up to 5x Cost Savings

## Improved Efficiency Reduces Overall Cost of Solution

**Llama3-70B on 4xH100 SXM**

1346 tokens/sec

**5X**

**improved throughput**

250 tokens/sec

NIM Off

NIM On

**Llama-3-8B on 1xH100**

9343 tokens/sec

**3X**

**improved throughput**
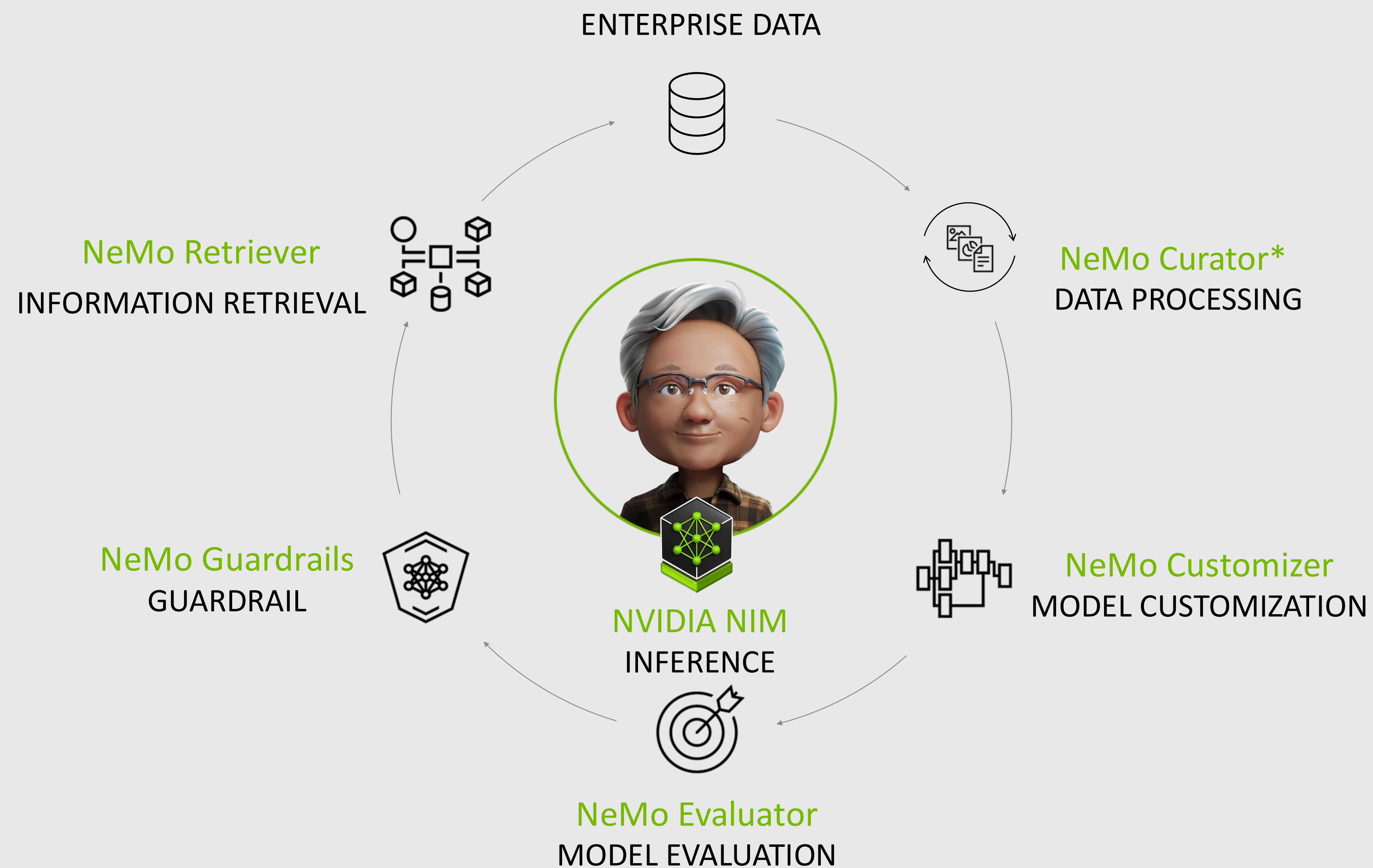
2703 tokens/sec

NIM Off

NIM On

Llama 3-70b-instruct, input token length: 7,000, output token length: 1,000. Concurrent client requests: 100. 4xH100 SXM NVLink. NIM Off: FP16, TTFT: ~120s, ITL: ~180ms. NIM On: FP8. TTFT: ~4.5s, ITL: ~70ms.

# NVIDIA NeMo Microservices

## Modular End-to-End Platform to easily build Data Flywheels

ENTERPRISE DATA

NeMo Curator*
DATA PROCESSING

NeMo Customizer
MODEL CUSTOMIZATION

NVIDIA NIM
INFERENCE

NeMo Evaluator
MODEL EVALUATION

NeMo Guardrails
GUARDRAIL

NeMo Retriever
INFORMATION RETRIEVAL

**Easy Setup**
Modular microservices, deployable with standard APIs

**Broad Ecosystem Support**
Integrated in popular open frameworks and AI software platforms

**Enterprise-grade**
Secure, stable, and supported software

**Run Anywhere**
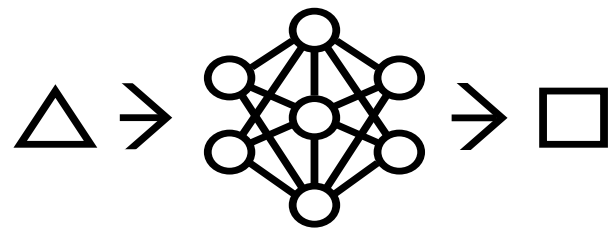Provides higher security, privacy, and flexibility

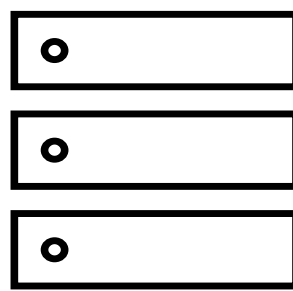**NVIDIA**

# NVIDIA Blueprints

Available on build.nvidia.com

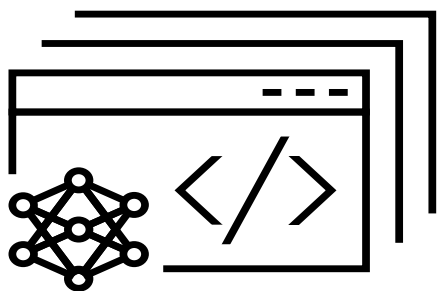NVIDIA NIM & microservices

Blueprints

Reference Application

Sample Data

Reference Code

Architecture

Customization Tools

Orchestration Tools

AI Agents

# NVIDIA Provides the Building Blocks for Agentic AI



**AI Blueprints**

Research Assistant Agent

Customer Service Agent

Software Security Agent

Virtual Lab Agent

Video Analytics Agent

**NVIDIA NeMo**

| Curator | Customizer | Evaluator | Guardrails | Retriever |

**NVIDIA NIM**

Understanding & Reasoning

Information Retrieval

AI Safety

Digital Humans

Visual Content Generation

Digital Biology

Physical AI

**AI Infrastructure**

GPU

CPU

DPU

Networking

NVIDIA.

# NVIDIA AI Enterprise Supported Software

For additional information, see NGC Catalog

| SDKs and Frameworks | | |
|---|---|---|
| Maxine | Riva | Modulus |
| TAO Toolkit | MONAI | CUDA |
| DeepStream | Clara | CUDA Toolkit |
| Metropolis | Clara Holoscan | ACE |
| NeMo | RAPIDS | Kubernetes Device Plugin |
| TensorRT | cuOpt | PyTorch Geometric |
| Triton Inference Server | Merlin | PyTorch |
| TensorFlow | Morpheus | DGL |

| NIM Microservices | | |
|---|---|---|
| Full list at build.nvidia.com/nim | | |
| **Infrastructure Software Collection** | | |
| NIM Operator | GPU Operator | Network Operator |
| NVIDIA Data Center Driver | NVIDIA vGPU for AI | NVIDIA DOCA Driver for Networking |
| Base Command Manager | | |
| **Extended Life Software Branches** | | |
| Production Branch (9-month) | | Long-Term Support Branch (3 year) |

# Summary

# Value Proposition

Enterprise Reference Architectures for NVIDIA Partners and Customers

## Optimized Performance

Comprehensive cluster design recommendations built upon tested and validated technologies to ensure peak computing performance for generative AI workloads; NVIDIA AI Enterprise, NIM, Training, Inference

## Reduced Complexity

Avoid design and planning pitfalls when deploying infrastructure with detailed guidance on server, cluster, and network configuration, minimizing setup errors and accelerating deployment timelines.
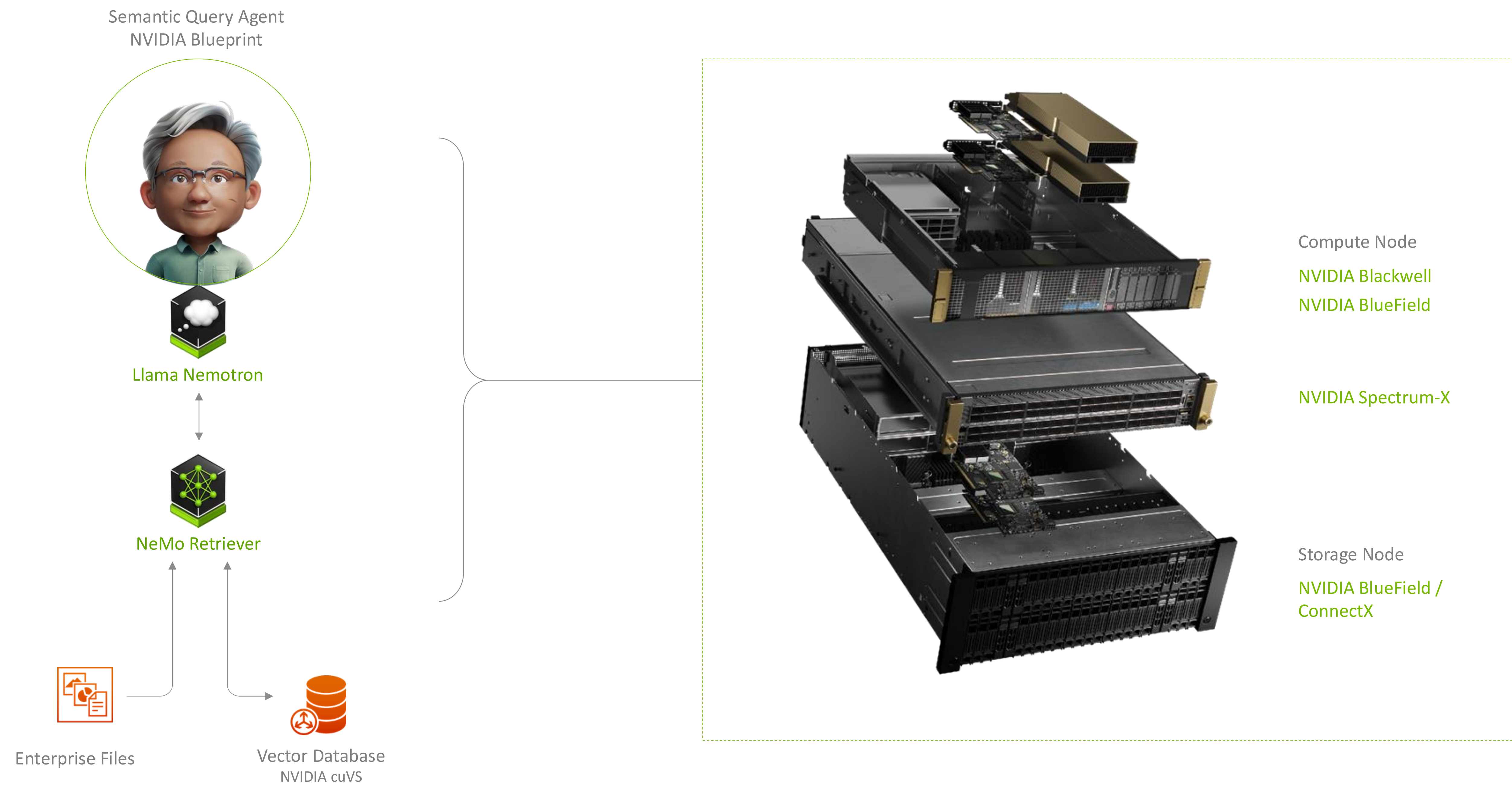
## Flexibility & Scale

Improve resource utilization and eliminate over-provisioning with discrete design points for compute, network, and storage infrastructure based upon deployment scale.

# NVIDIA AI Factory for Enterprise

## New class of infrastructure serving knowledge instead of data



Semantic Query Agent
NVIDIA Blueprint

Llama Nemotron

NeMo Retriever

Enterprise Files

Vector Database
NVIDIA cuVS

Compute Node

NVIDIA Blackwell
NVIDIA BlueField

NVIDIA Spectrum-X

Storage Node

NVIDIA BlueField /
ConnectX

PURESTORAGE®

NVIDIA.

# Enterprise AI Factory Stack

**Agentic AI**



Supply Chain    Customer Service    Marketing    Video Analytics

**Physical AI**



Robotics    Digital Twins

**HPC**



CUDA-X Libraries

---

**Security**

**AI Platform**

NVIDIA AI Enterprise

| NVIDIA NIM | NVIDIA NeMo |

**Physical AI**

NVIDIA Omniverse

NVIDIA Cosmos

**HPC**

NVIDIA CUDA

**Data Connectors**

**Repository**

NVIDIA AI Enterprise Application Software and Tools

**Observability**

**Orchestration and Management**

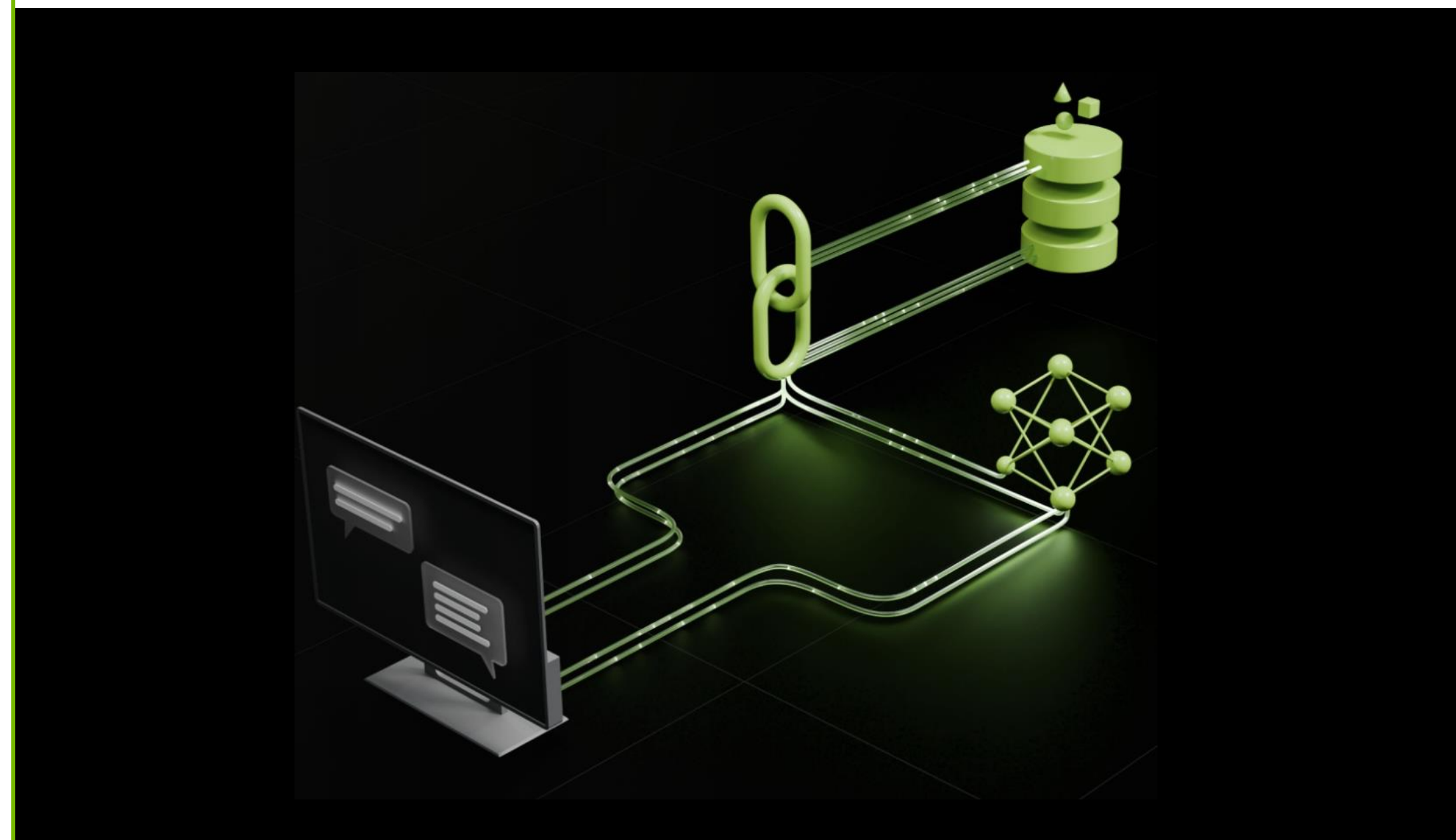| NVIDIA-Certified Storage | NVIDIA-Certified Systems | NVIDIA Networking |

NVIDIA.

# Getting Started With NVIDIA AI Enterprise

Experience for free on NVIDIA-hosted infrastructure or deploy on your own infrastructure

## NVIDIA API Catalog

- Try out sample applications from your web browser

- Free API access to experiment and prototype with NVIDIA-optimized microservices



## Developer Access

NVIDIA Developer Program members can

- Download any NIM containers

- Self-host for research and testing (up to 16 GPUs)



## 90-Day Evaluation

- Free evaluation licenses for POC

- Running on compatible on-prem or cloud accelerated infrastructure

- Access to exclusive features