# AI Factories – Transforming Data Centers with NVIDIA and HPE

Eric Kang 康勝閔, 資深解決方案架構協理, NVIDIA

October 16, 2025

# Trillion-Dollar Global IT Investment Shifting to AI Factories
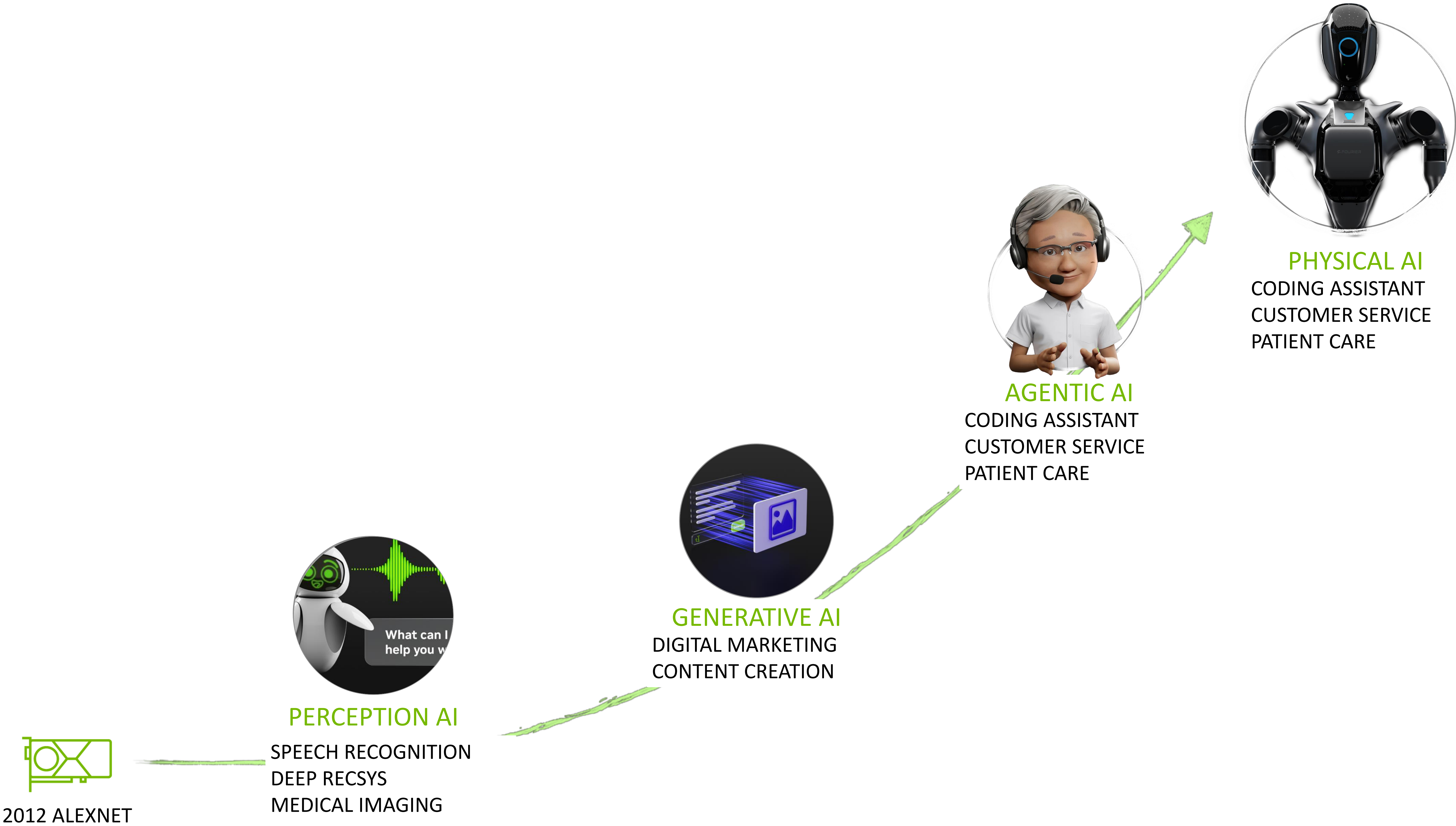
**92%**    of enterprises investing in AI

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

**33%**    find complexity top barrier for adoption

**1%**     have mature AI deployments
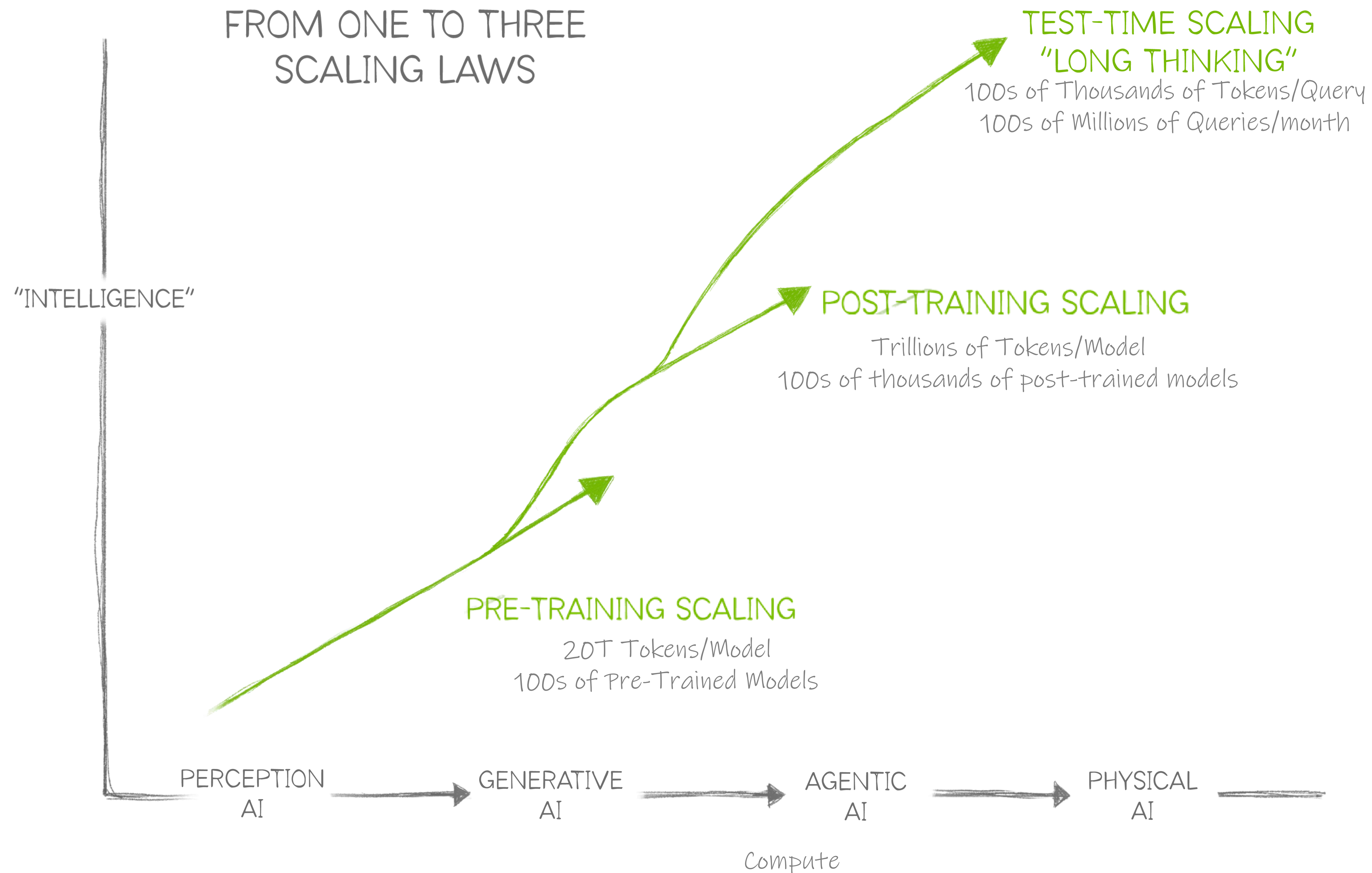
NVIDIA.

# Evolution of AI
## Agentic AI Enables More Powerful AI Applications



**PHYSICAL AI**
CODING ASSISTANT
CUSTOMER SERVICE
PATIENT CARE

**AGENTIC AI**
CODING ASSISTANT
CUSTOMER SERVICE
PATIENT CARE

**GENERATIVE AI**
DIGITAL MARKETING
CONTENT CREATION
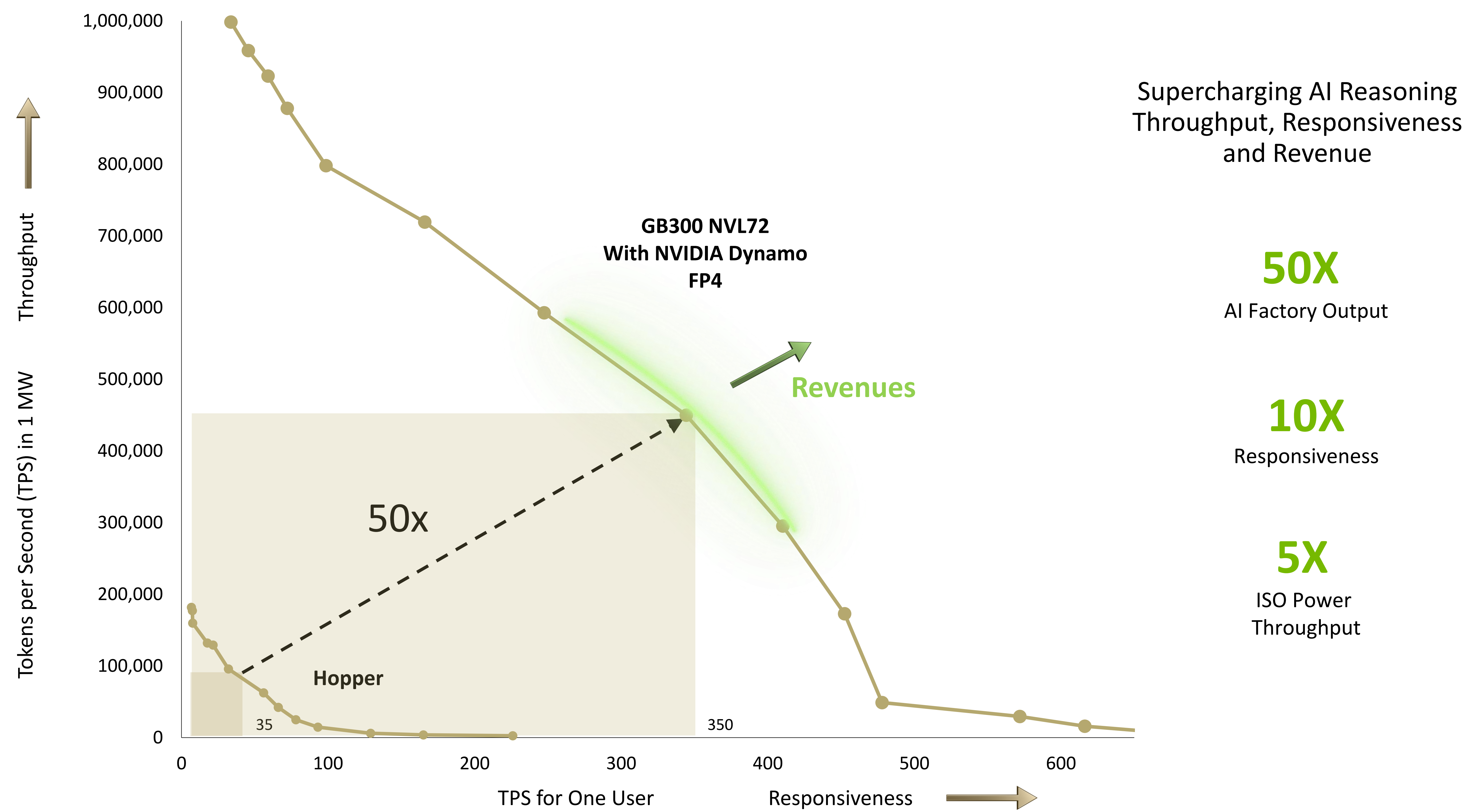
**PERCEPTION AI**
SPEECH RECOGNITION
DEEP RECSYS
MEDICAL IMAGING

2012 ALEXNET

# AI Scaling Laws Drive Exponential Demand for Compute

New "long thinking" required for agentic and physical AI

FROM ONE TO THREE
SCALING LAWS

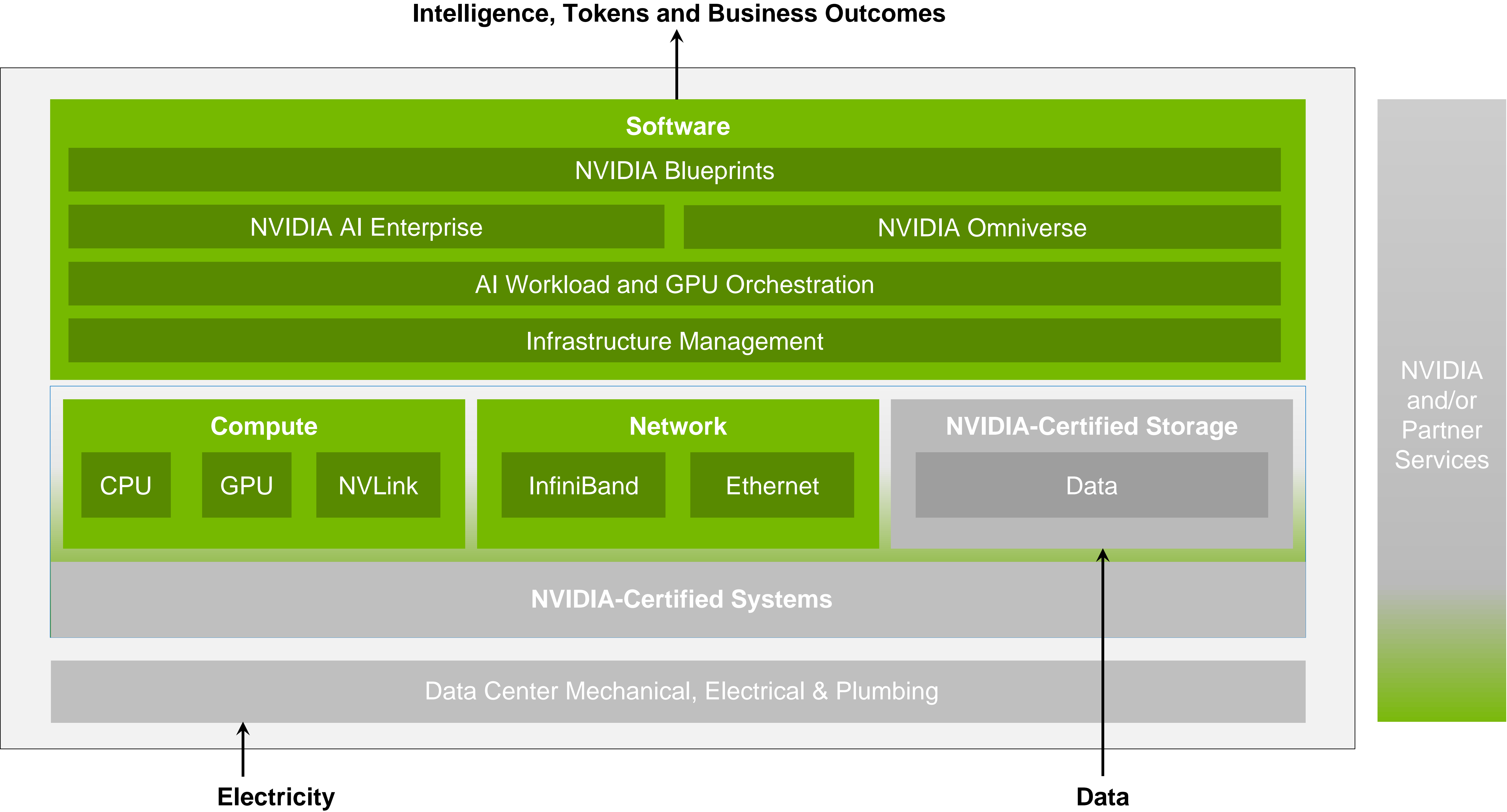TEST-TIME SCALING
"LONG THINKING"
100s of Thousands of Tokens/Query
100s of Millions of Queries/month

"INTELLIGENCE"

POST-TRAINING SCALING
Trillions of Tokens/Model
100s of thousands of post-trained models

PRE-TRAINING SCALING
20T Tokens/Model
100s of Pre-Trained Models

PERCEPTION
AI → GENERATIVE
AI → AGENTIC
AI → PHYSICAL
AI

Compute

NVIDIA.

# AI Factory Output Drives Revenue

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



Supercharging AI Reasoning Throughput, Responsiveness and Revenue

**50X**
AI Factory Output

**10X**
Responsiveness

**5X**
ISO Power Throughput

GB300 NVL72
With NVIDIA Dynamo
FP4

Revenues

50x

Hopper

Throughput

Tokens per Second (TPS) in 1 MW

TPS for One User

Responsiveness

# NVIDIA Provides a Full Stack for AI Factories

## Built on NVIDIA 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**

CPU · GPU · NVLink

**Network**

InfiniBand · Ethernet

**NVIDIA-Certified Storage**

Data

**NVIDIA-Certified Systems**

Data Center Mechanical, Electrical & Plumbing

NVIDIA and/or Partner Services

**Electricity**

**Data**

# Transforming Data Centers into AI Powerhouses

## Enterprises Need for Purpose-Built AI Factories for the Age of AI Reasoning

**Traditional Data Center**
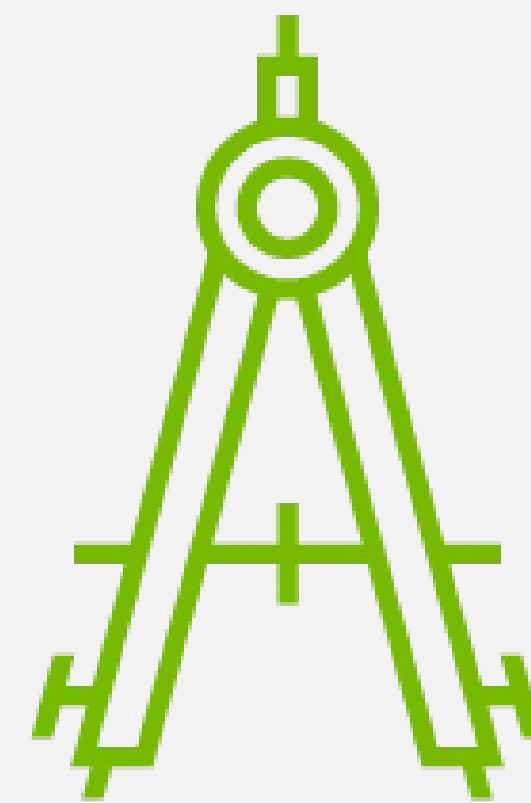Systems of Record / Systems of Engagement / Single Server Workloads

**Today's AI Factories**
AI Scale-Up and Scale-Out Workloads

| Traditional Data Center | Today's AI Factories |
|---|---|
| Cost Center | Profit Engine |
| Retrieval | Generative |
| X86 CPU | Accelerated |
| 20 kW/Rack | >130 kW/Rack |
| Air Cooled | Liquid and Air Cooled |
| 100k NICs | 400k Network Fabric |

NVIDIA

# Today's AI Challenges Require AI Factories

## AI Workloads Require Optimized Full Stack Solutions

**Design Complexity**

Spans project prioritization,
data acquisition, infrastructure,
and sizing

**Deployment and Cost**
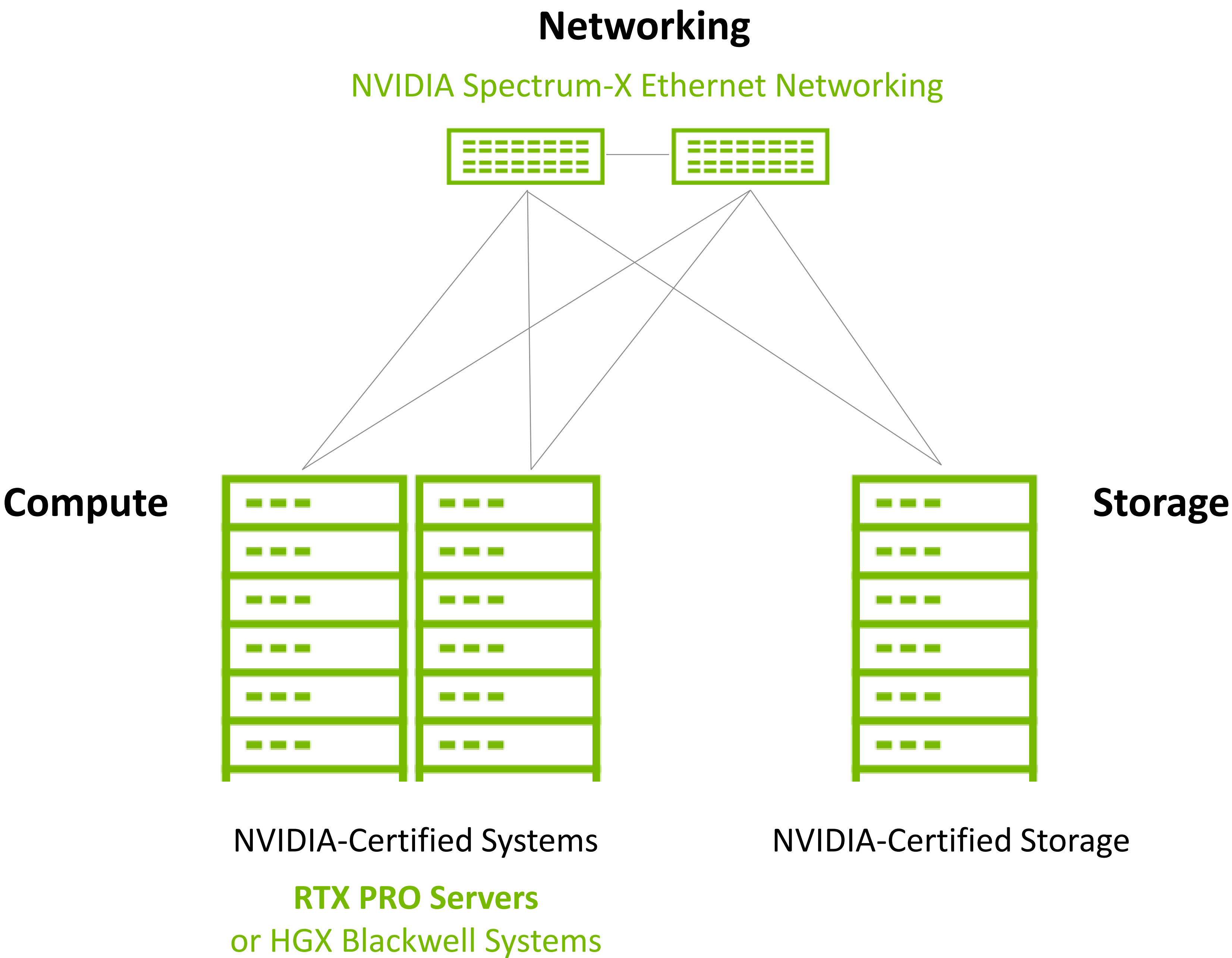
Infrastructure, security,
and customization

**Time to Value**

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

# NVIDIA Enterprise AI Factory Validated Design

Building on NVIDIA Enterprise Reference Architectures

**Networking**

NVIDIA Spectrum-X Ethernet Networking

**Compute**

**Storage**

NVIDIA-Certified Systems

NVIDIA-Certified Storage

**RTX PRO Servers**
or HGX Blackwell Systems

✓ **Time to value**     ✓ **Scalability**     ✓ **Manageability**     ✓ **Security**

NVIDIA.

# NVIDIA Accelerated Computing for Enterprise

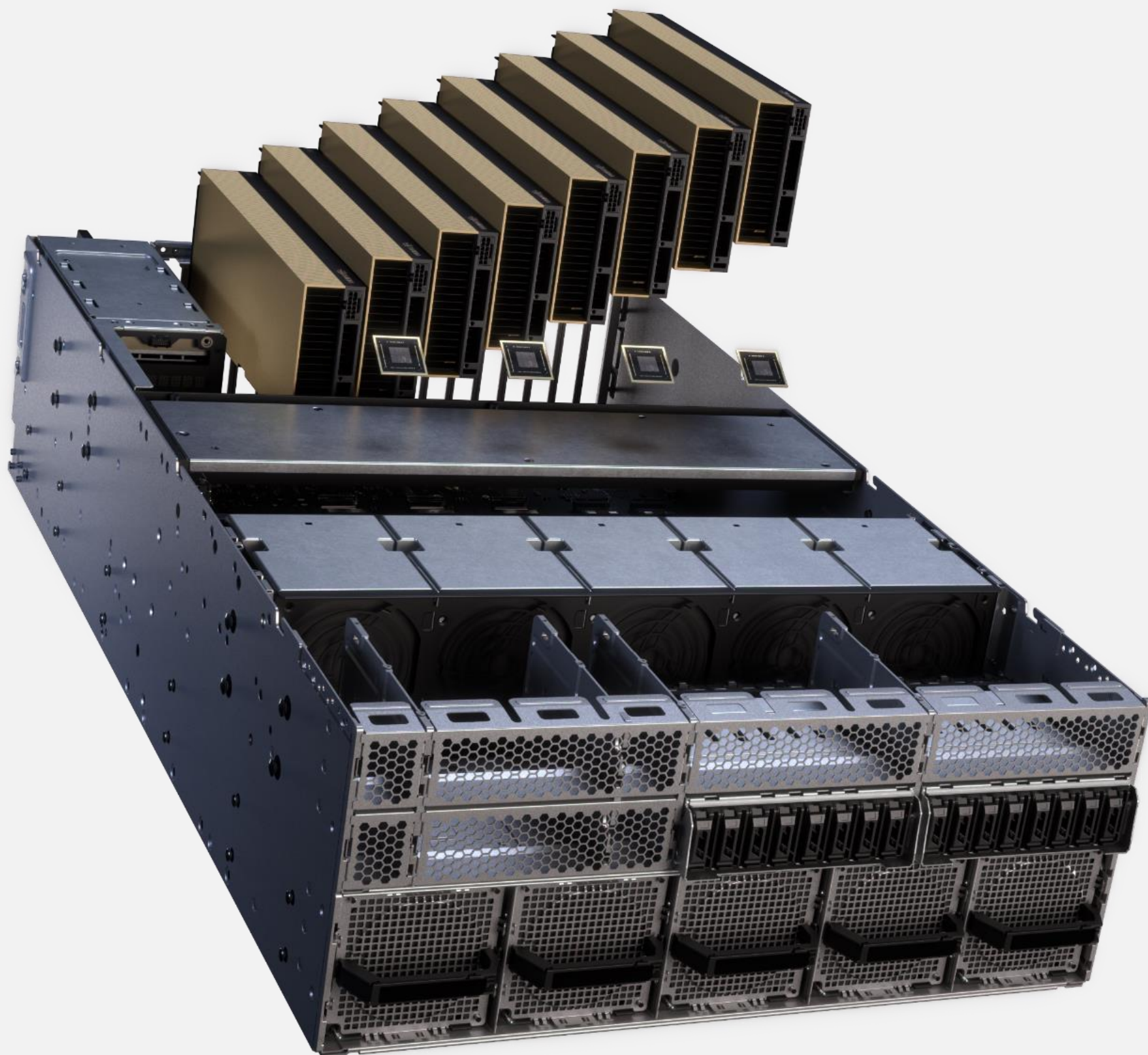## From Hopper to Blackwell, Enterprise AI Factory Building Blocks

LLM Inference and Enterprise AI
Available for Deployments Now

Flexibility and Performance
for Enterprise AI and Industrial AI

Best Air-Cooled Platform
for LLM Training and Inference



**H200 NVL**

**RTX PRO Server**

**HGX B200**

Hopper, incl 5 yr NVIDIA AI Enterprise license
8kW, Air-Cooled PCIe, x86, with NVL4 and FP64
Enterprise RA (2-8-5)

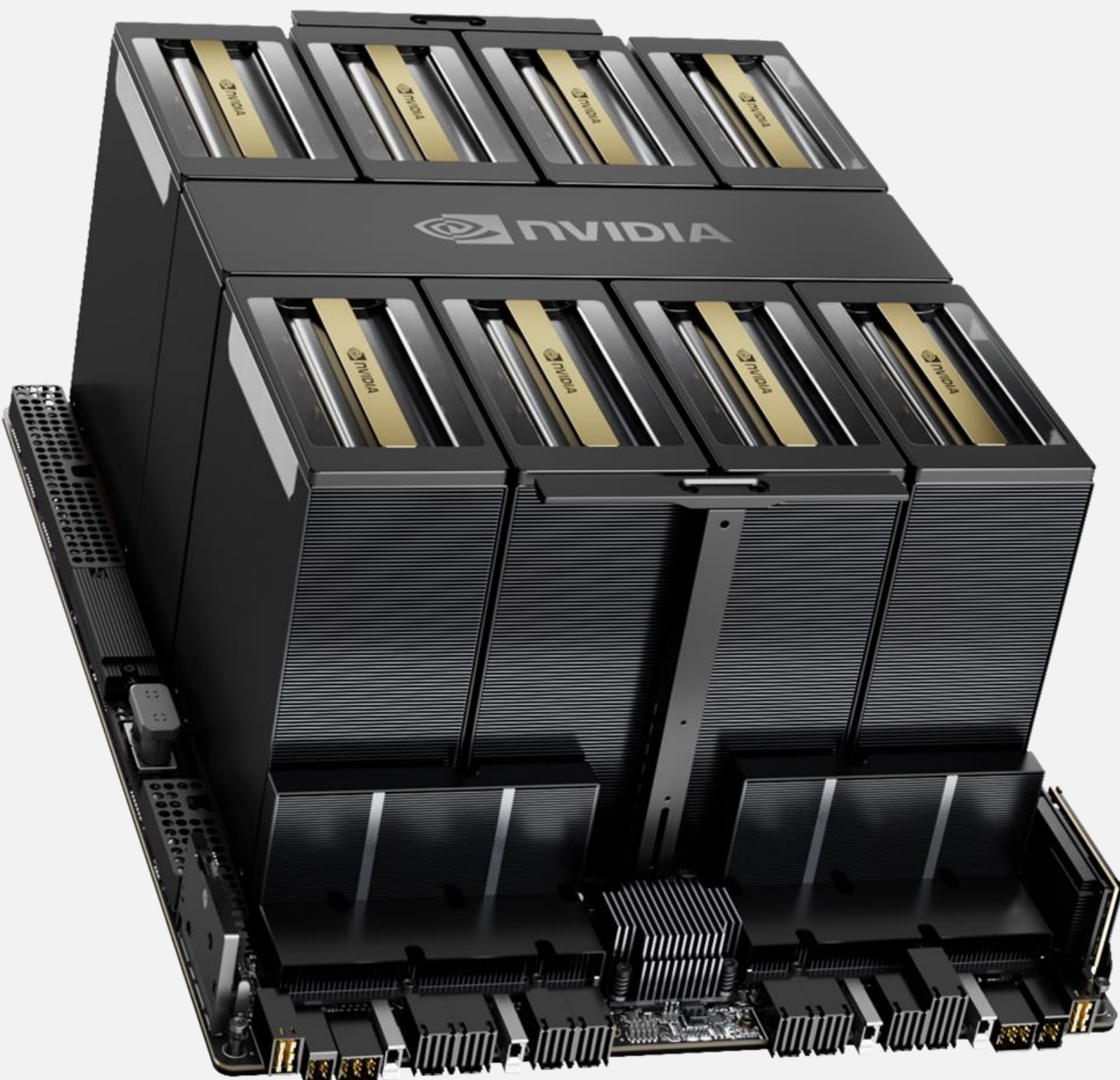Blackwell, Enterprise DC Compatibility, Multi-Workload
7kW, Air-Cooled PCIe, x86, RTX Graphics, vGPU
Enterprise RA (2-8-5)

Blackwell, Optimized Performance and TCO
14kW, Air-Cooled HGX, x86
Enterprise RA (2-8-9)

* 2-8-5 config (PCIe-Optimized) = 2 CPUs, 8 GPUs, 5 network adapters; 2-8-9 (HGX) = 2 CPUs, 8 GPUs, 9 network adapters

# RTX PRO 6000 Blackwell Server Edition

The Most Powerful Blackwell Data Center Platform for AI and Visual Computing

## Breakthrough Multimodal AI Inference

- 5th-Gen Tensor, 2nd-Gen Transformer Engine, FP4
- Full Media Pipeline: 4 NVENC/ NVDEC/ NVJPEG
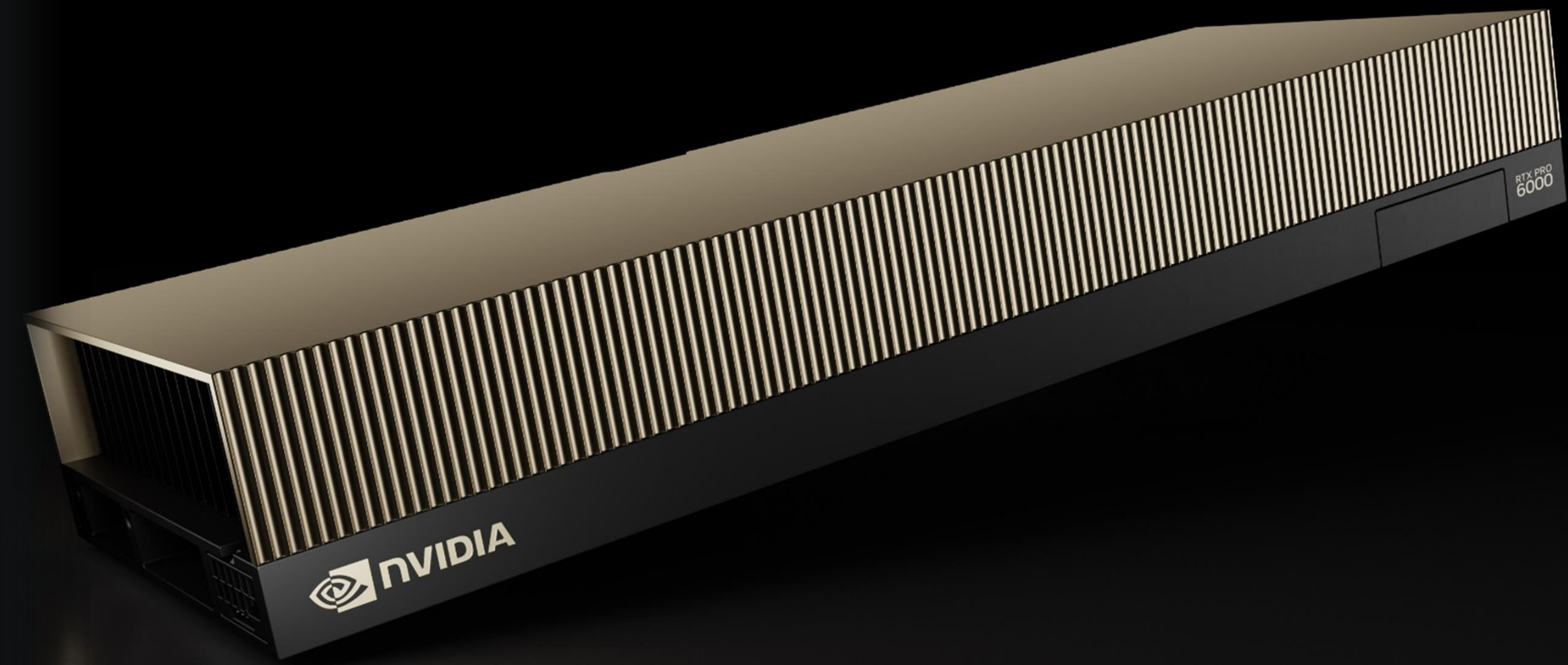
## Powerful Graphics and Visual Computing

- 4th-Gen RTX, Neural Shaders, DLSS 4

## Data Center Ready

- 96GB GDDR7,1.6 TB/s Memory BW, 128MB L2 Cache
- Multi-Instance GPU (MIG), TEE Confidential Compute

## Performance Specs

- ✓ 188 Ray Tracing Cores
- ✓ 752 Tensor Cores
- ✓ 24,064 Cuda Cores

- ✓ Peak FP4 AI Performance: 3.7 PFLOPS
- ✓ Peak RT Core Performance: 354.5 TFLOPS

Dual Slot, FHFL I Up to 600W

# Modern Enterprises Have Diverse Accelerated Workloads

## From Agentic AI and Physical AI to AI-Enabled Applications



AGENTIC AI        INDUSTRIAL & PHYSICAL AI        SCIENTIFIC COMPUTING, DATA ANALYTICS, & SIMULATION        VISUAL COMPUTING        ENTERPRISE APPLICATIONS
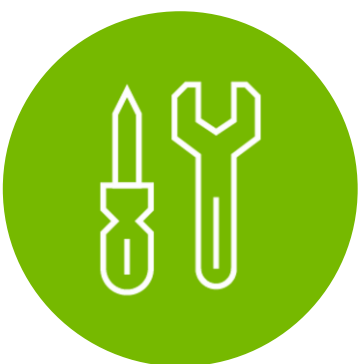
NVIDIA AI Enterprise

NVIDIA Omniverse

NVIDIA CUDA-X Microservices

HPE ProLiant 380a Gen12

### NVIDIA AI Computing by HPE
Co-developed solutions to simplify enterprise AI

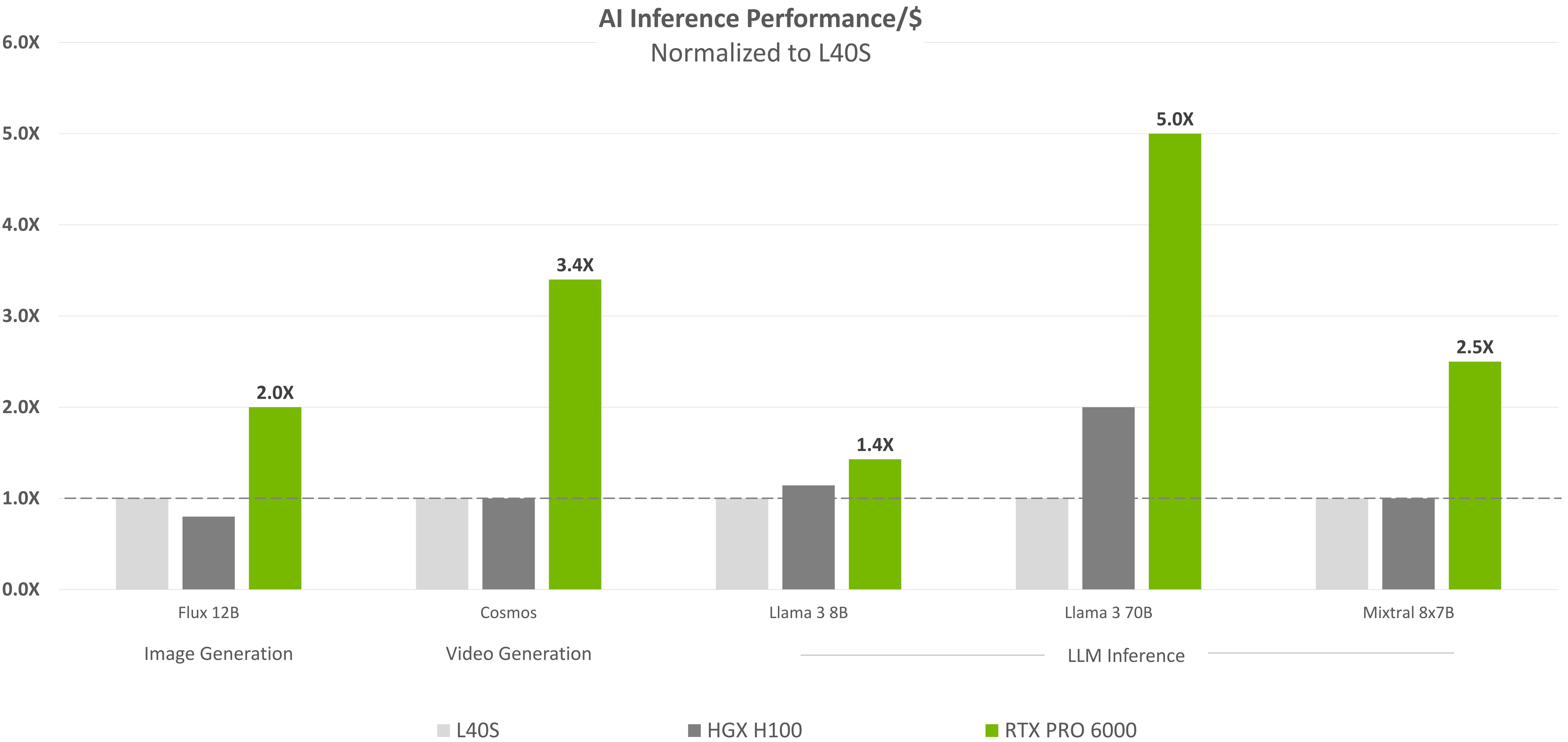| | | |
|---|---|---|
| Turnkey Private Cloud | PEOPLE | Inference, Tuning & Training |
| AI Services & Training | | Virtual Assistants |
| AI Optimized Systems | TECHNOLOGY | Process Automation |
| Enterprise-grade Control | ECONOMICS | Content & Product Creation |

NVIDIA RTX PRO 6000
Blackwell Server Edition

NVIDIA.

# Best Server for Enterprise AI
## RTX PRO Server Up to 5X Better Price Performance

**AI Inference Performance/$**
Normalized to L40S



Flux 12B — 2.0X
Cosmos — 3.4X
Llama 3 8B — 1.4X
Llama 3 70B — 5.0X
Mixtral 8x7B — 2.5X

Image Generation | Video Generation | LLM Inference

Legend: ■ L40S  ■ HGX H100  ■ RTX PRO 6000

Performance/$ = Performance / TCO for a Single Node (Server + Power Costs) for L40S and HGX H100 compared to RTX PRO 6000 Blackwell Server Edition
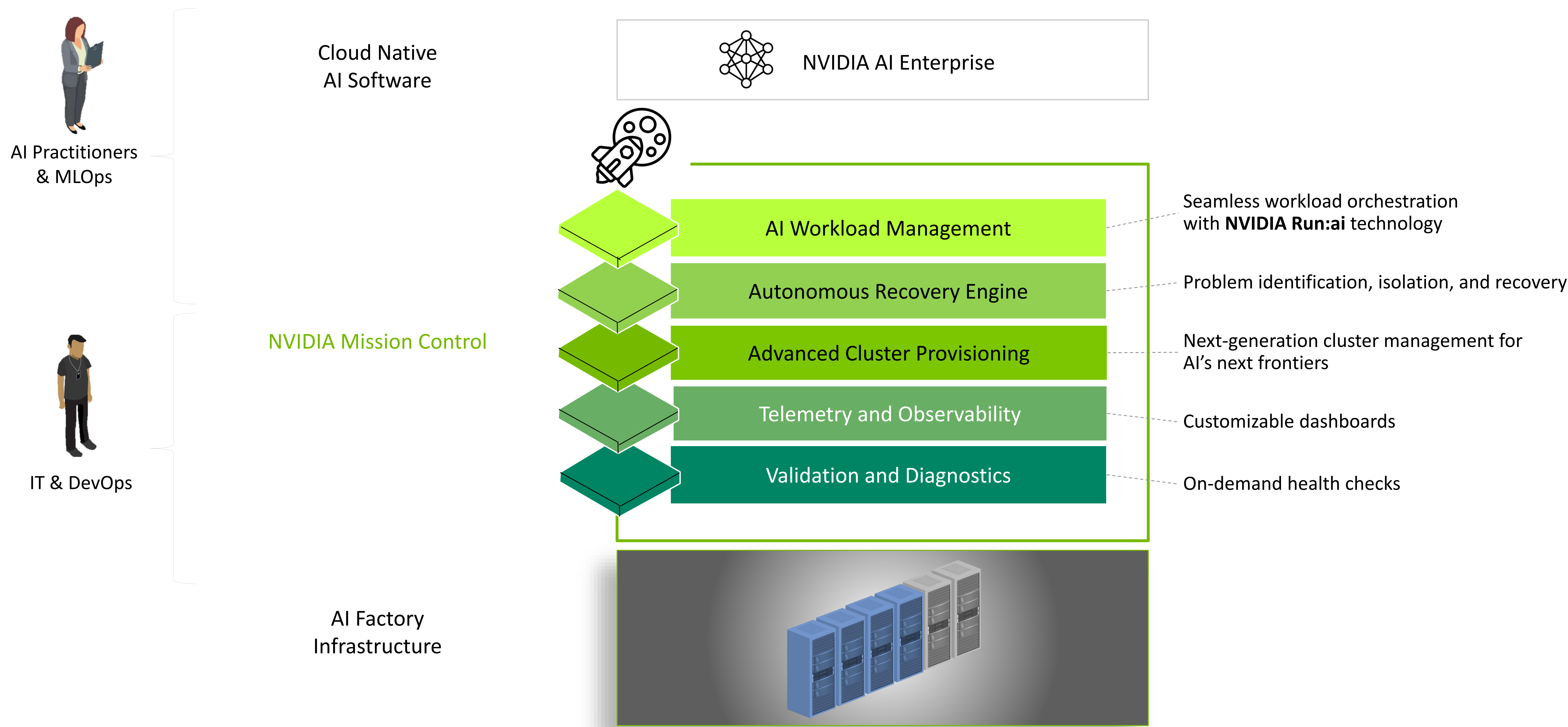
NVIDIA.

# Building Blocks For Enterprise AI Factories



## NVIDIA Certified

**Servers**     **Storage**

**High performance, scalable, and secure accelerated computing platform for AI & HPC**

PERFORMANCE          MANAGEABILITY

SECURITY          SCALABILITY

## NVIDIA Enterprise Reference Architectures

**Proven and Comprehensive Recommendations for Scaling Enterprise Deployment**

NVIDIA-Certified Servers | Network Topology Management | Software Stack

- Faster Time to Solution
- Performant & Scalable Infrastructure
- Reduced Complexity and Cost
- Supportable design patterns and scalable units

## Adopted by NVIDIA and HPE

**HPE | NVIDIA**

**Informing designs to speed TTM with an enhanced deployment for AI Factories of various sizes**

Superior Customer Experience

- Faster Time to Deploy
- Optimized Performance
- Reduced Complexity and Cost
- Faster Time to Resolution

# State of the Art Infrastructure Management

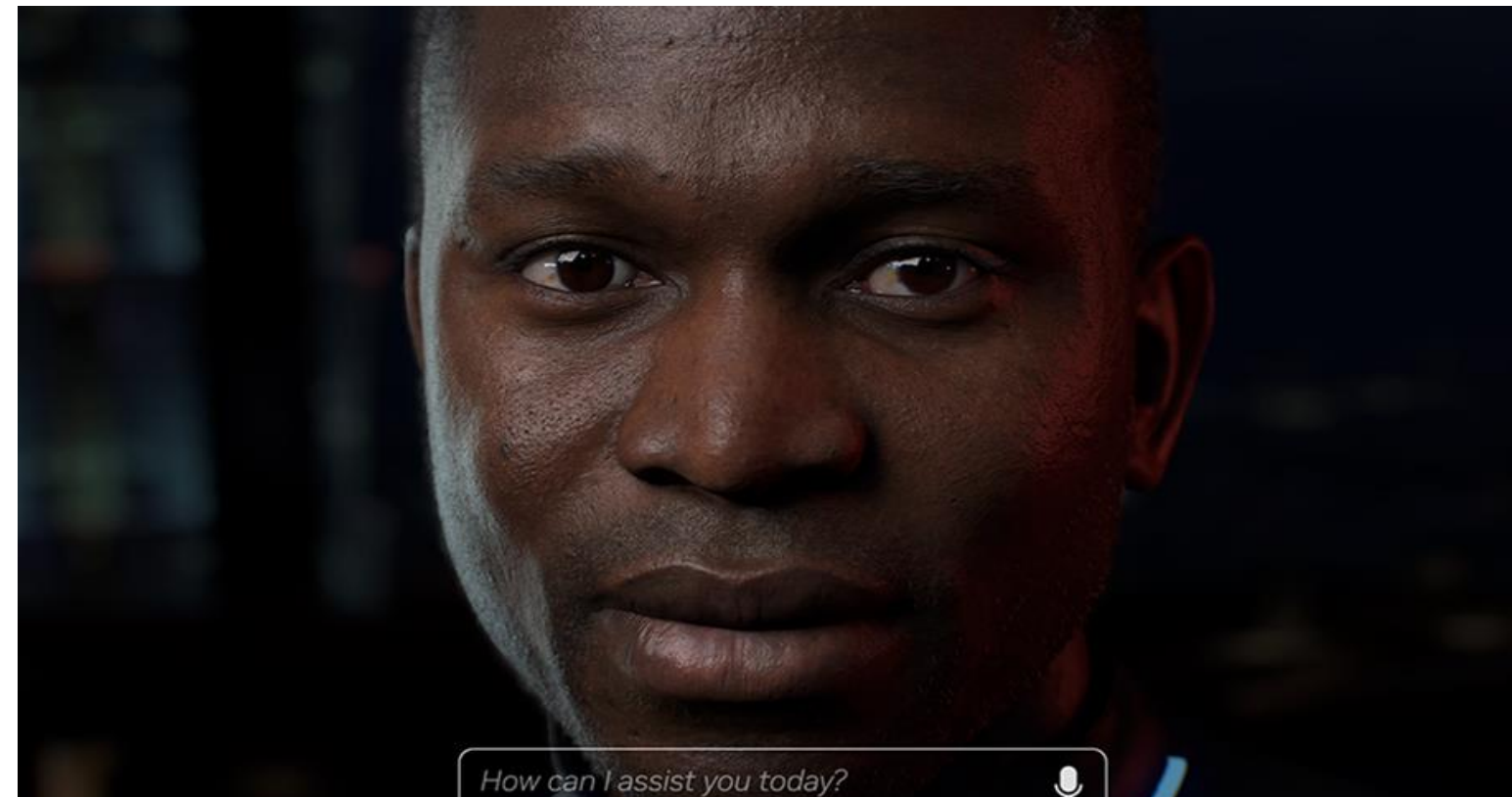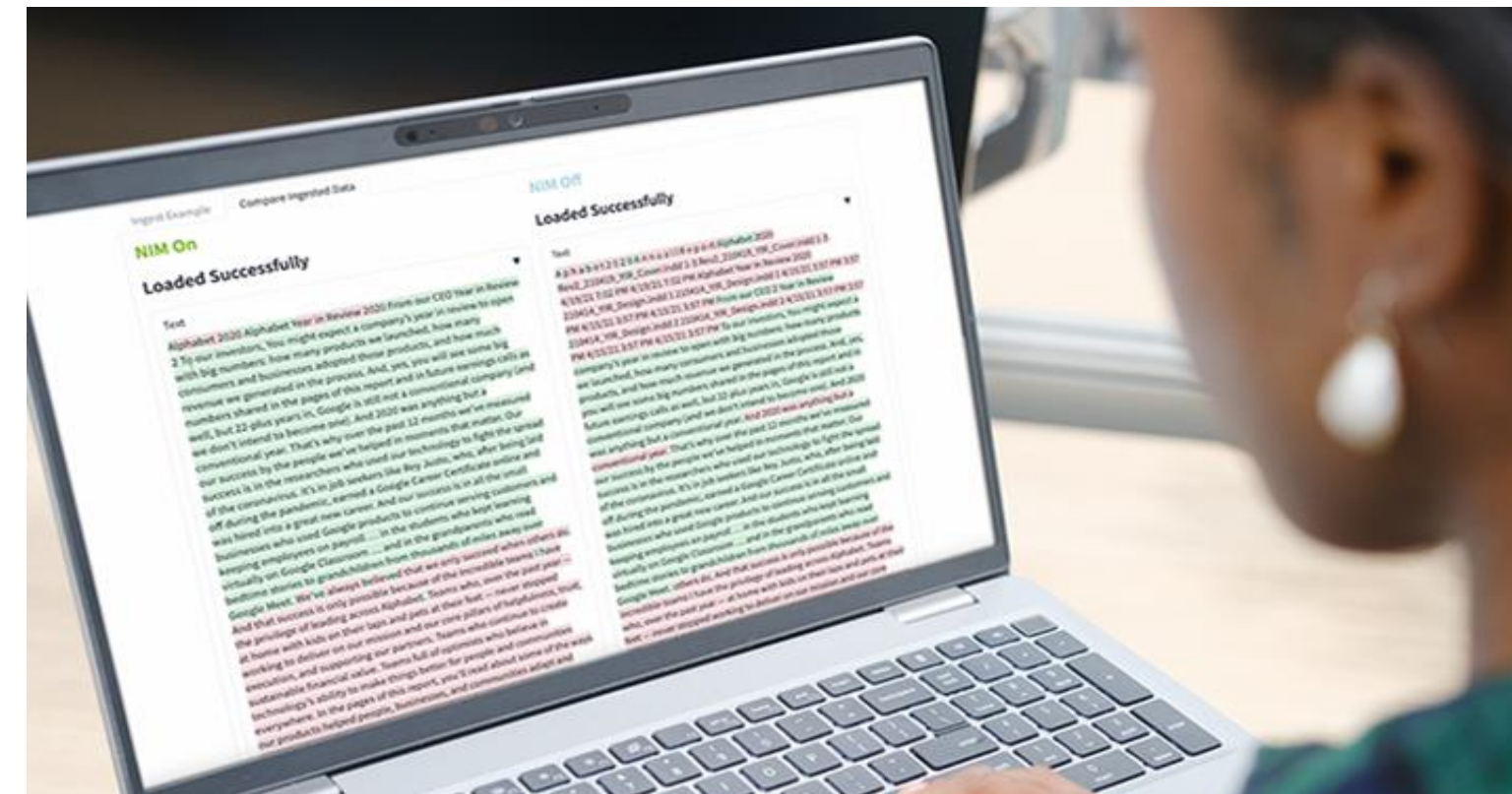## Empower model builders and enterprise IT with full stack intelligence

Cloud Native
AI Software

NVIDIA AI Enterprise

AI Practitioners
& MLOps

NVIDIA Mission Control

| | |
|---|---|
| AI Workload Management | Seamless workload orchestration with **NVIDIA Run:ai** technology |
| Autonomous Recovery Engine | Problem identification, isolation, and recovery |
| Advanced Cluster Provisioning | Next-generation cluster management for AI's next frontiers |
| Telemetry and Observability | Customizable dashboards |
| Validation and Diagnostics | On-demand health checks |

IT & DevOps

AI Factory
Infrastructure

# NVIDIA Blueprints

Available on https://build.nvidia.com/

### Digital Humans
### for Customer Service



NVIDIA AI Blueprint

### Multimodal PDF Data Extraction
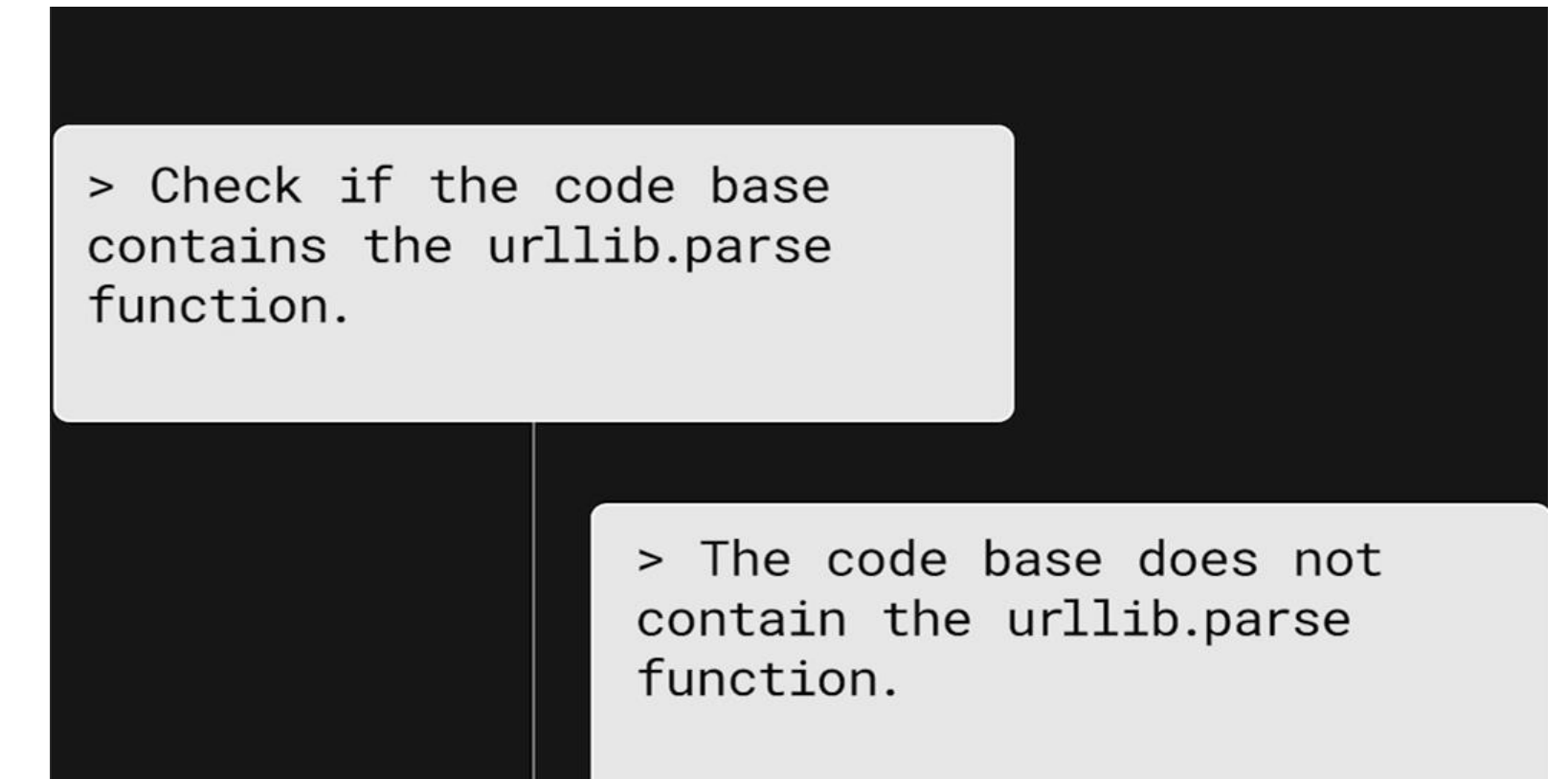### for Enterprise RAG



NVIDIA AI Blueprint

### Generative Virtual Screening
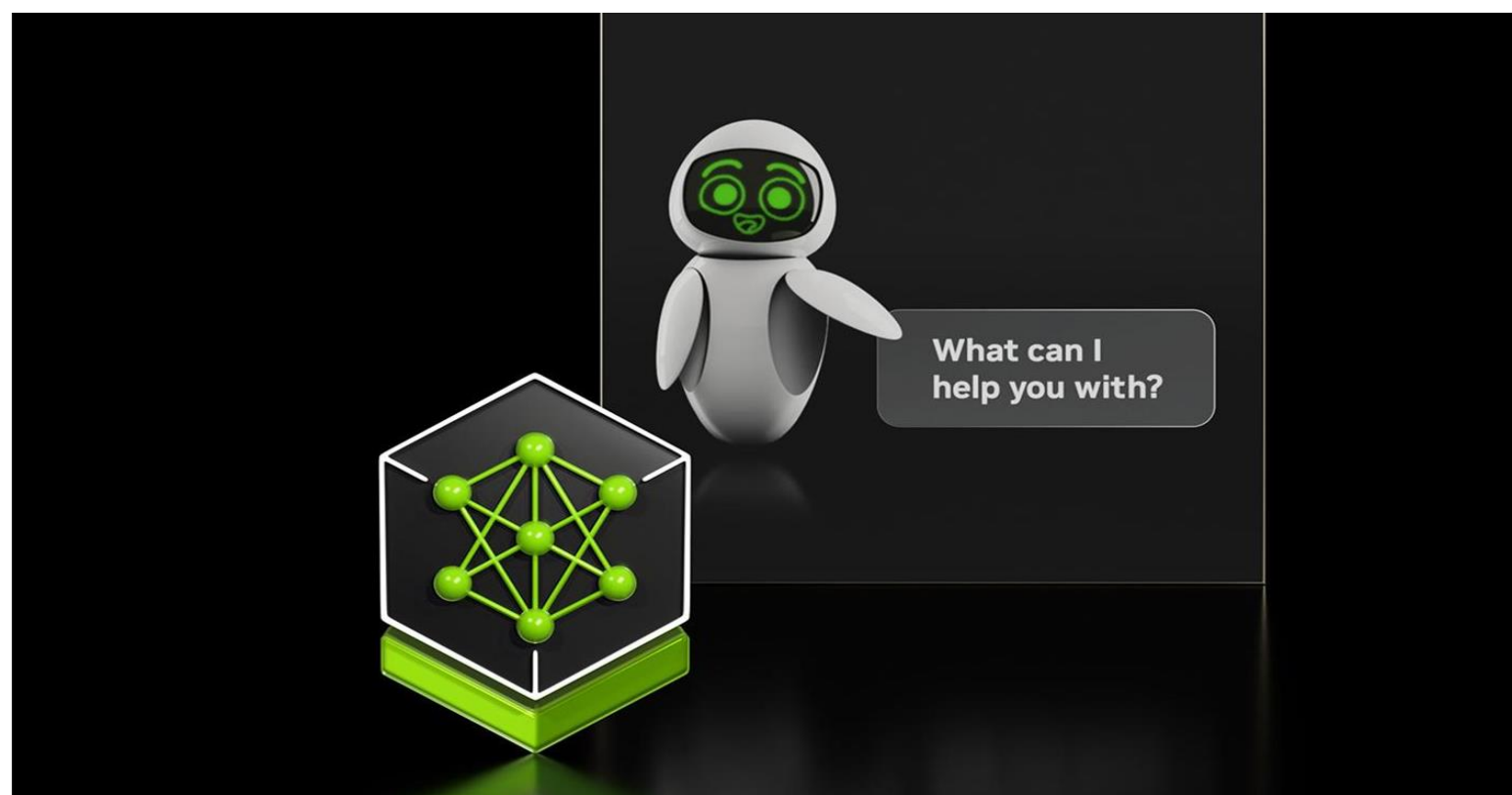### for Drug Discovery



NVIDIA BioNeMo Blueprint

### Vulnerability Analysis
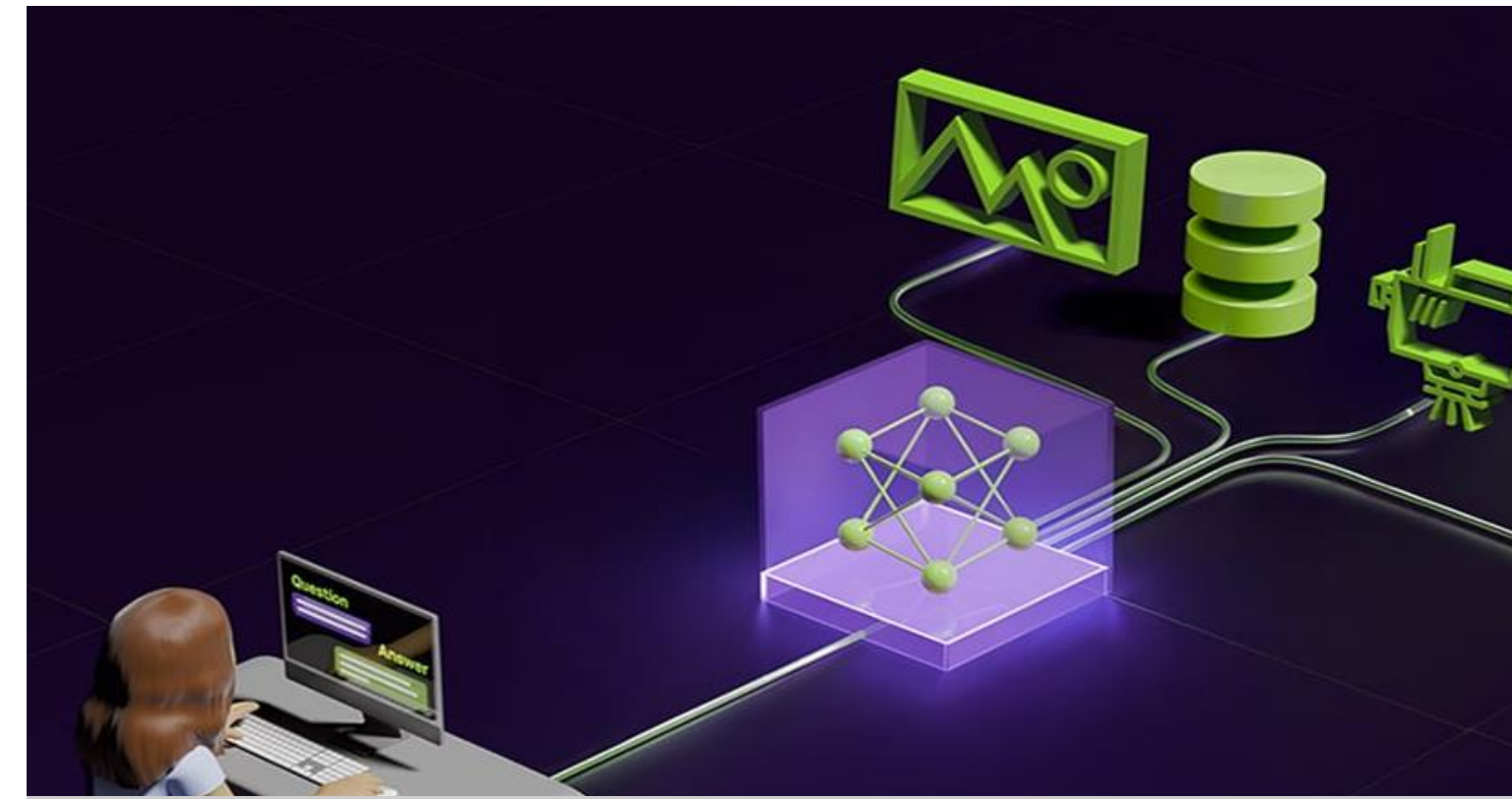### for Container Security



NVIDIA AI Blueprint

### AI Virtual Assistants
### for Customer Service



NVIDIA AI Blueprint

### Visual AI Agent
### for Video Search and Summarization



NVIDIA AI Blueprint

### 3D Conditioning for
### Precise Visual Generative AI



NVIDIA Omniverse Blueprint

### Build a Digital Twin for
### Interactive Fluid Simulation



NVIDIA Omniverse Blueprint
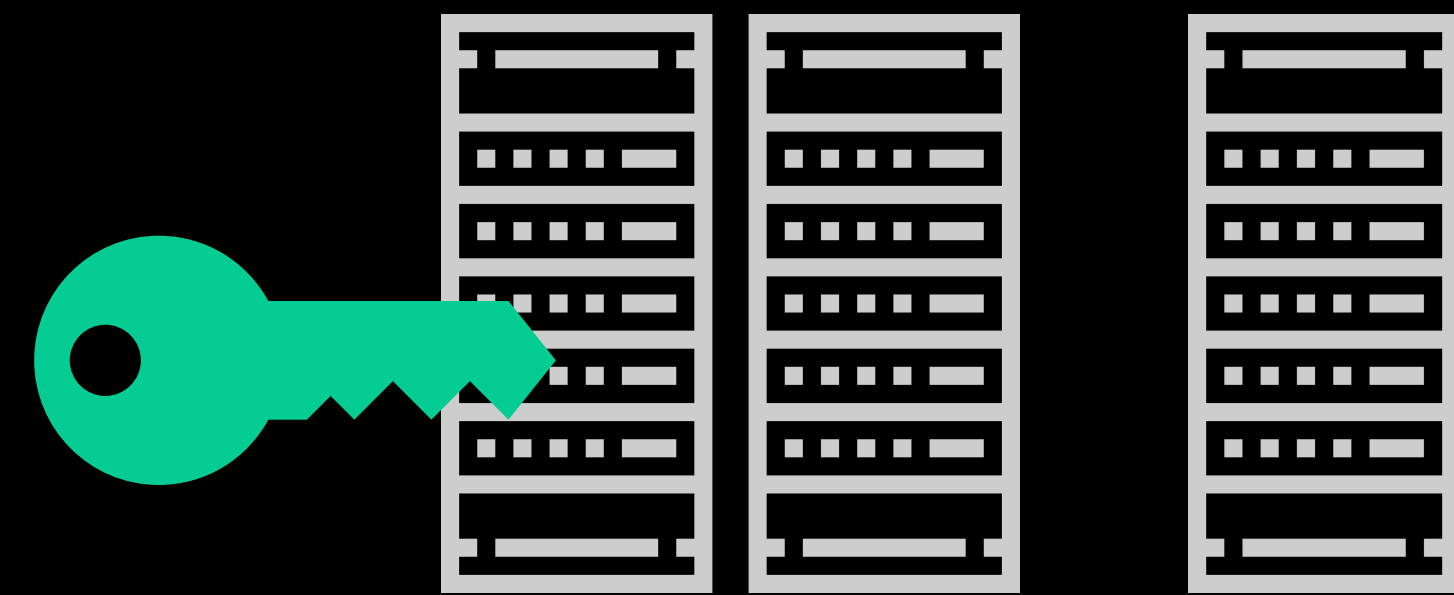
NVIDIA.

# Expanded AI factory portfolio from HPE
## for every AI ambition, across clouds, cores and countries
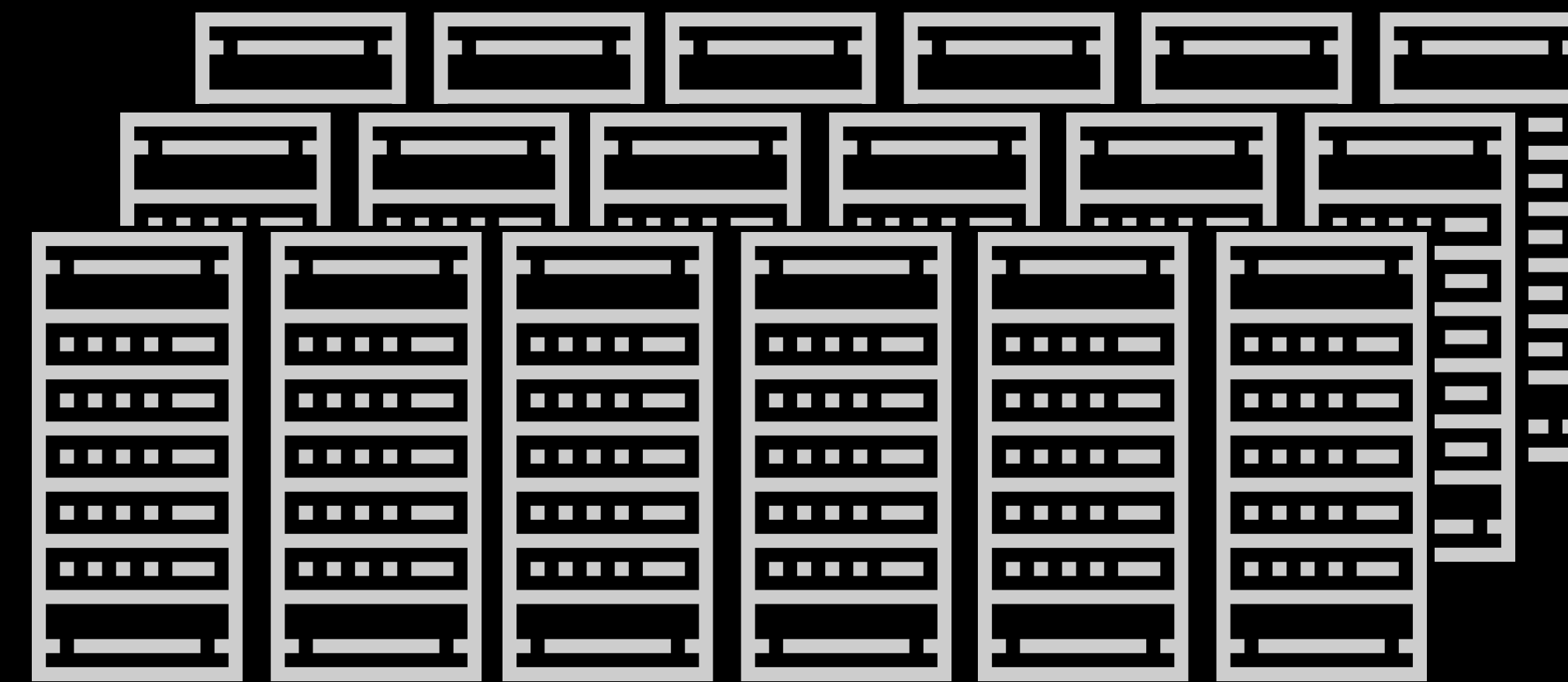
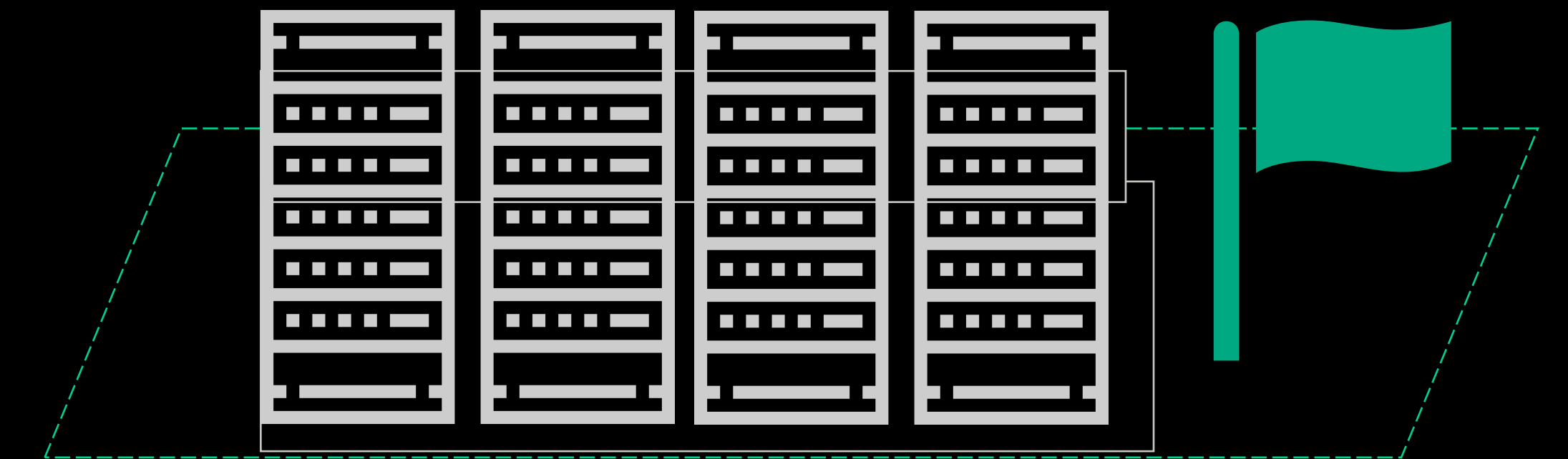| **Turnkey AI factory** Enterprises | **AI factory at scale** Model builders & SP's | **Sovereign AI factory** Governments, public sector |
|---|---|---|

Common control plane: HPE Morpheus and HPE OpsRamp

Turnkey, engineered systems                    Customized, validated solutions
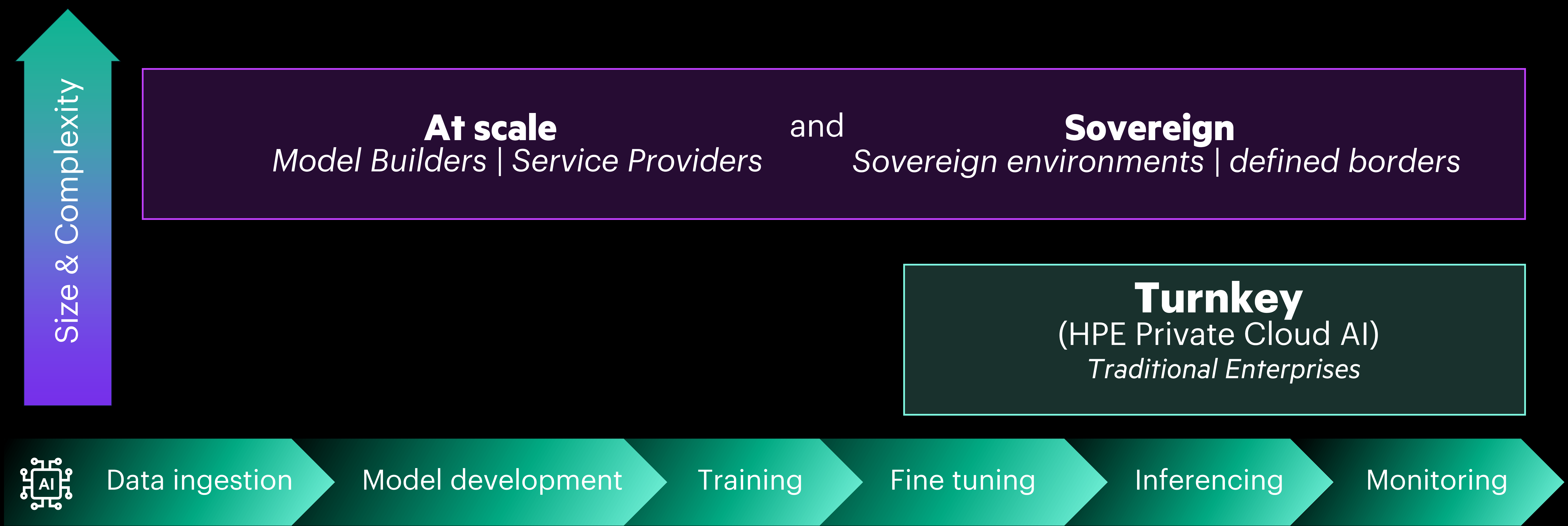
Infrastructure | Software | Services | Ecosystem | Sustainability

# The AI factory portfolio
## For every AI ambition

Size & Complexity

**At scale** and **Sovereign**
*Model Builders | Service Providers* *Sovereign environments | defined borders*

**Turnkey**
(HPE Private Cloud AI)
*Traditional Enterprises*

Data ingestion → Model development → Training → Fine tuning → Inferencing → Monitoring

# Thank You