# AI Acceleration with NVIDIA

Dr. Arts Yang, Sr. Product Architect DGX & Cloud

arts@nvidia.com
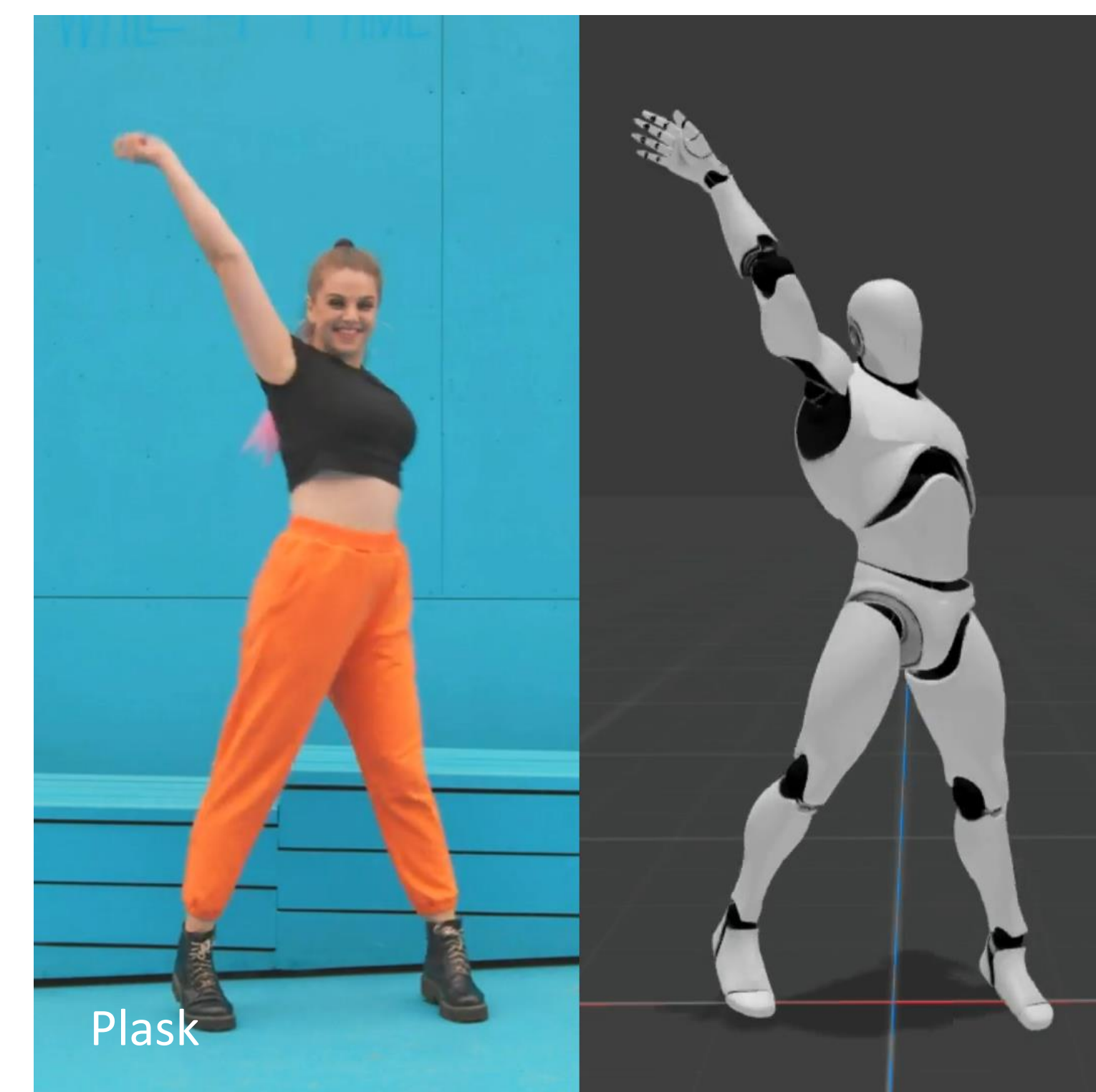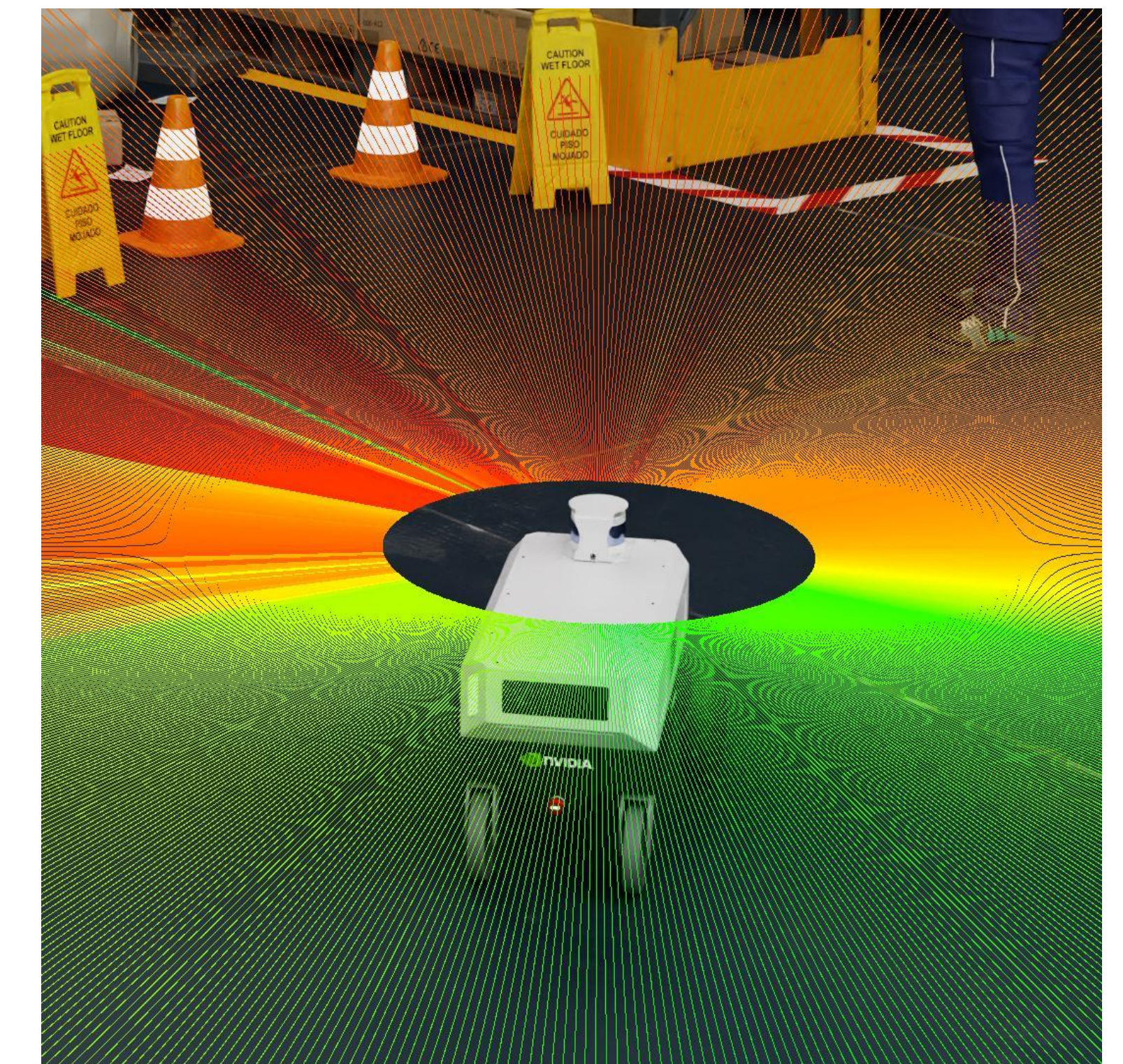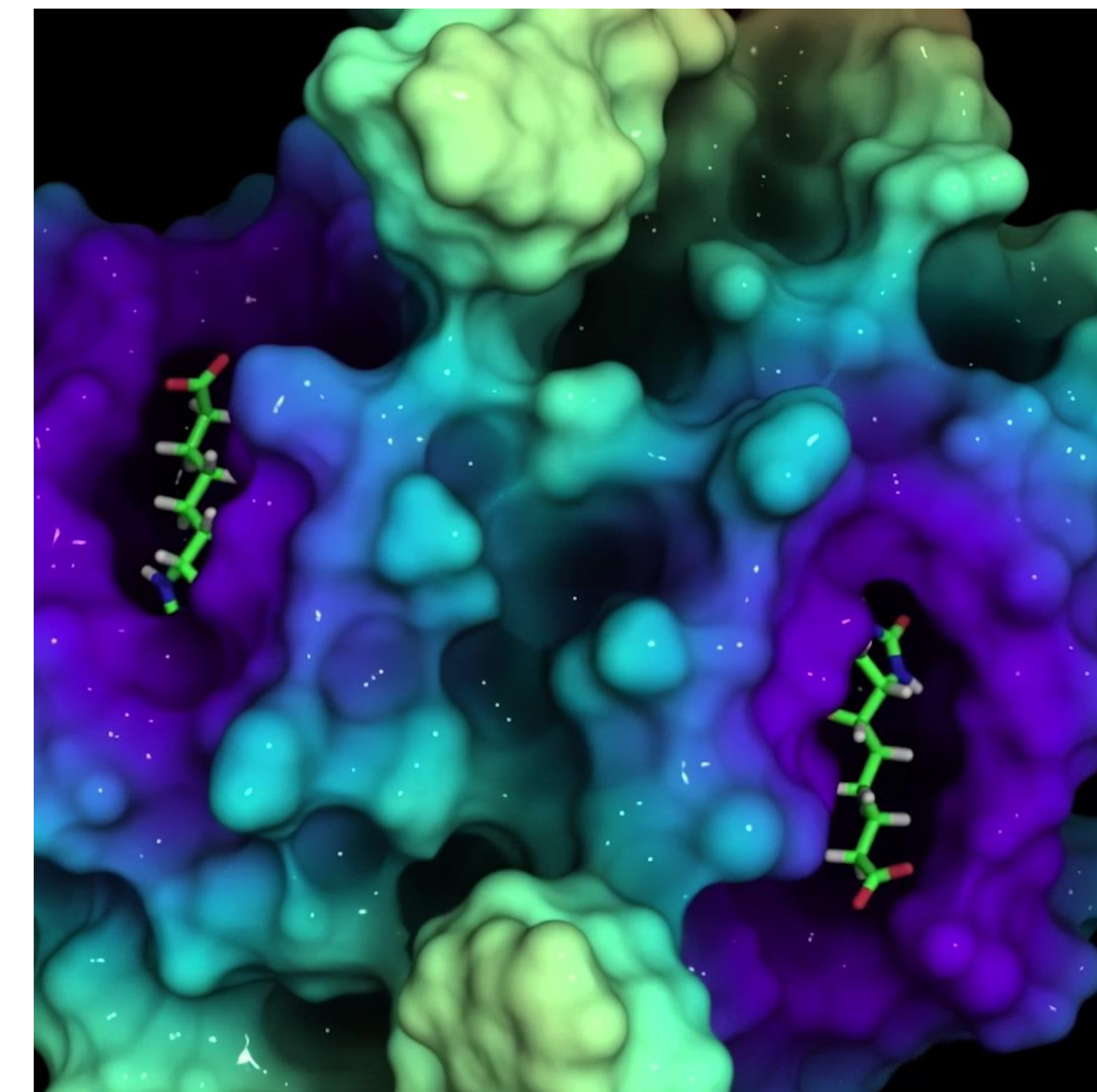
AI DAYS in Prague

25. Jan. 2024

# The In-Person GTC Experience Is Back

Come to GTC—the conference for the era of AI—to connect with a dream team of industry luminaries, developers, researchers, and business experts shaping what's next in AI and accelerated computing.

From the highly anticipated keynote by NVIDIA CEO Jensen Huang to over 600 inspiring sessions, 200+ exhibits, and tons of networking events, GTC delivers something for every technical level and interest area.

Be sure to save your spot for this transformative event. You can even take advantage of early-bird pricing when you register by February 7.

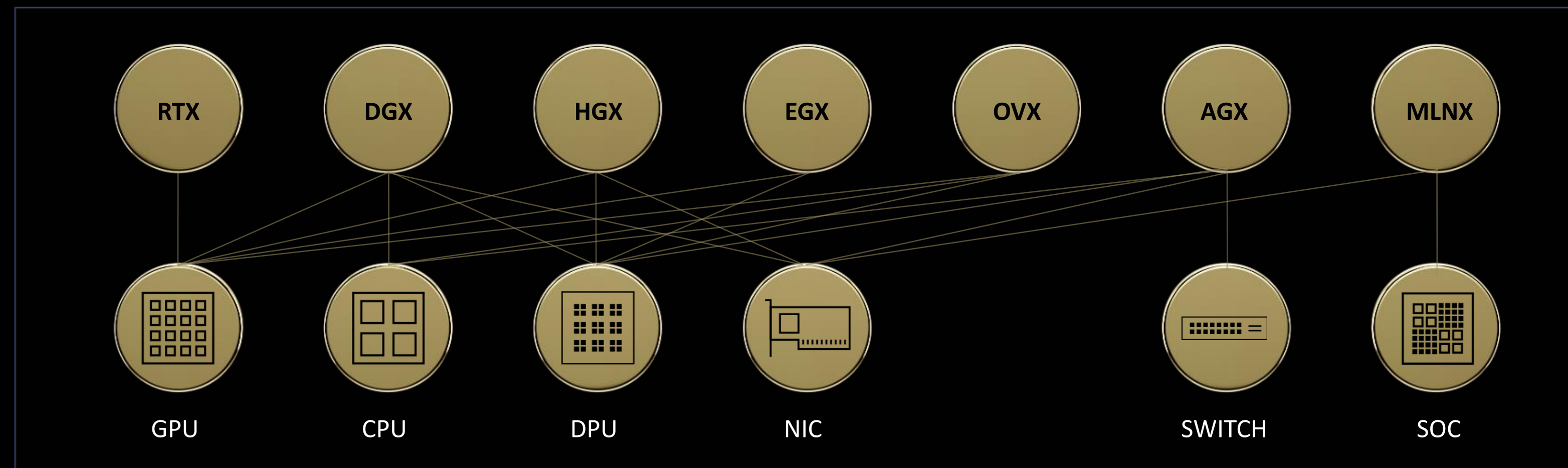**March 18-21, 2024 | www.nvidia.com/gtc**

APPLICATION FRAMEWORKS

MODULUS · MONAI · MAXINE · NEMO · AVATAR · DRIVE · ISAAC · METROPOLIS · HOLOSCAN

PLATFORM

NVIDIA AI · NVIDIA OMNIVERSE

ACCELERATION LIBRARIES

RTX
CUDA-X
CUDA

SYSTEM SOFTWARE

Magnum IO · DOCA · Base Command · Forge

HARDWARE

RTX · DGX · HGX · EGX · OVX · AGX · MLNX

GPU · CPU · DPU · NIC · SWITCH · SOC

# Generative AI Market and Use Cases

# What's the Hype!

## Generative AI is here

## Key Players

**Microsoft**
Reinvent search with Generative AI

**OpenAI** ChatGPT
110 million users in 2 months

**NVIDIA.**
Generative AI computing power & platform of the world

**Meta**
Business group focused on Generative AI

Enterprises & Startups
**1600+**

**Google**
Generative AI changing the main products

## Generative AI Possibilities

TEXT
AUDIO
IMAGE
3D
VIDEO
DNA
PROTEIN
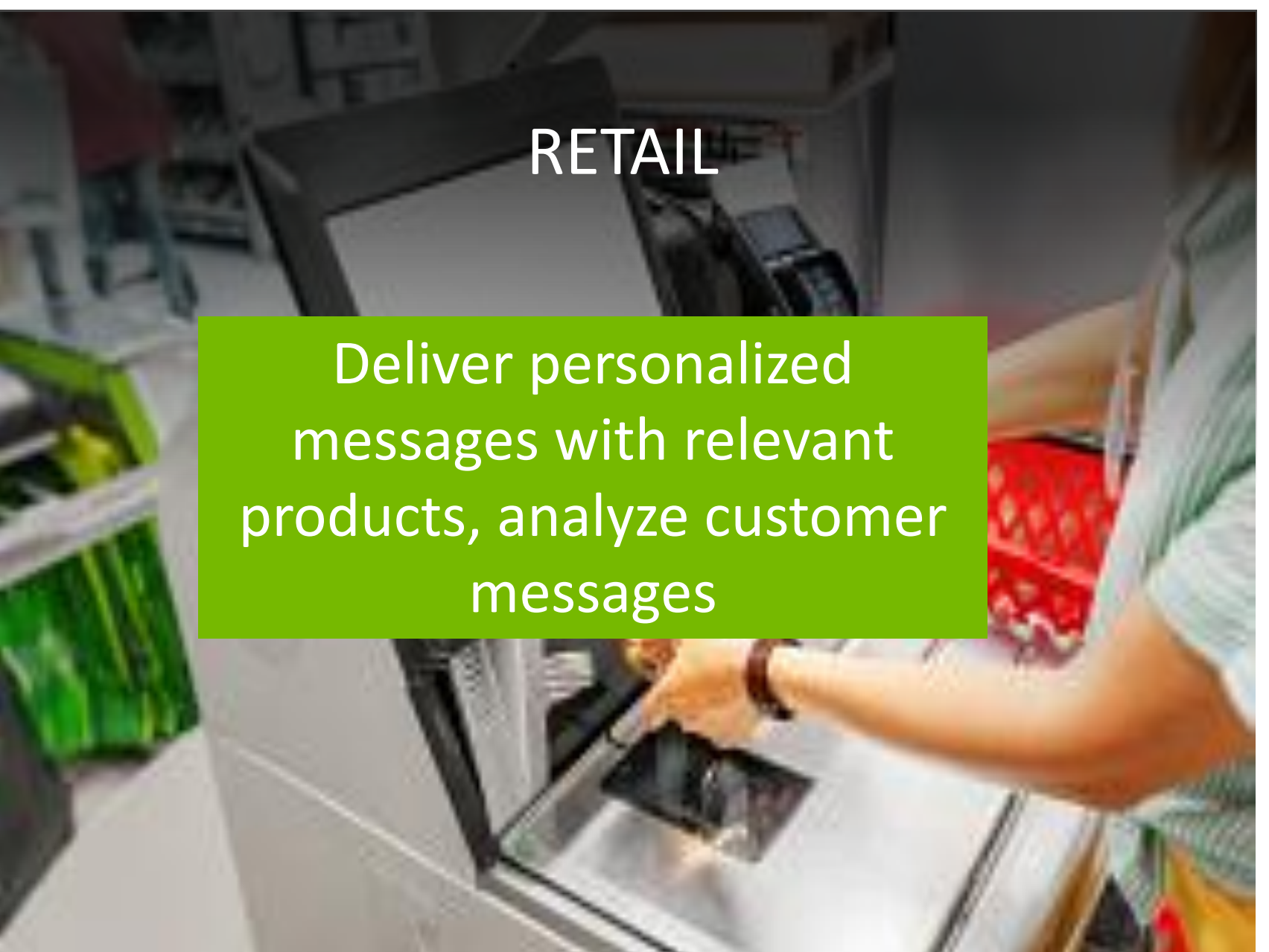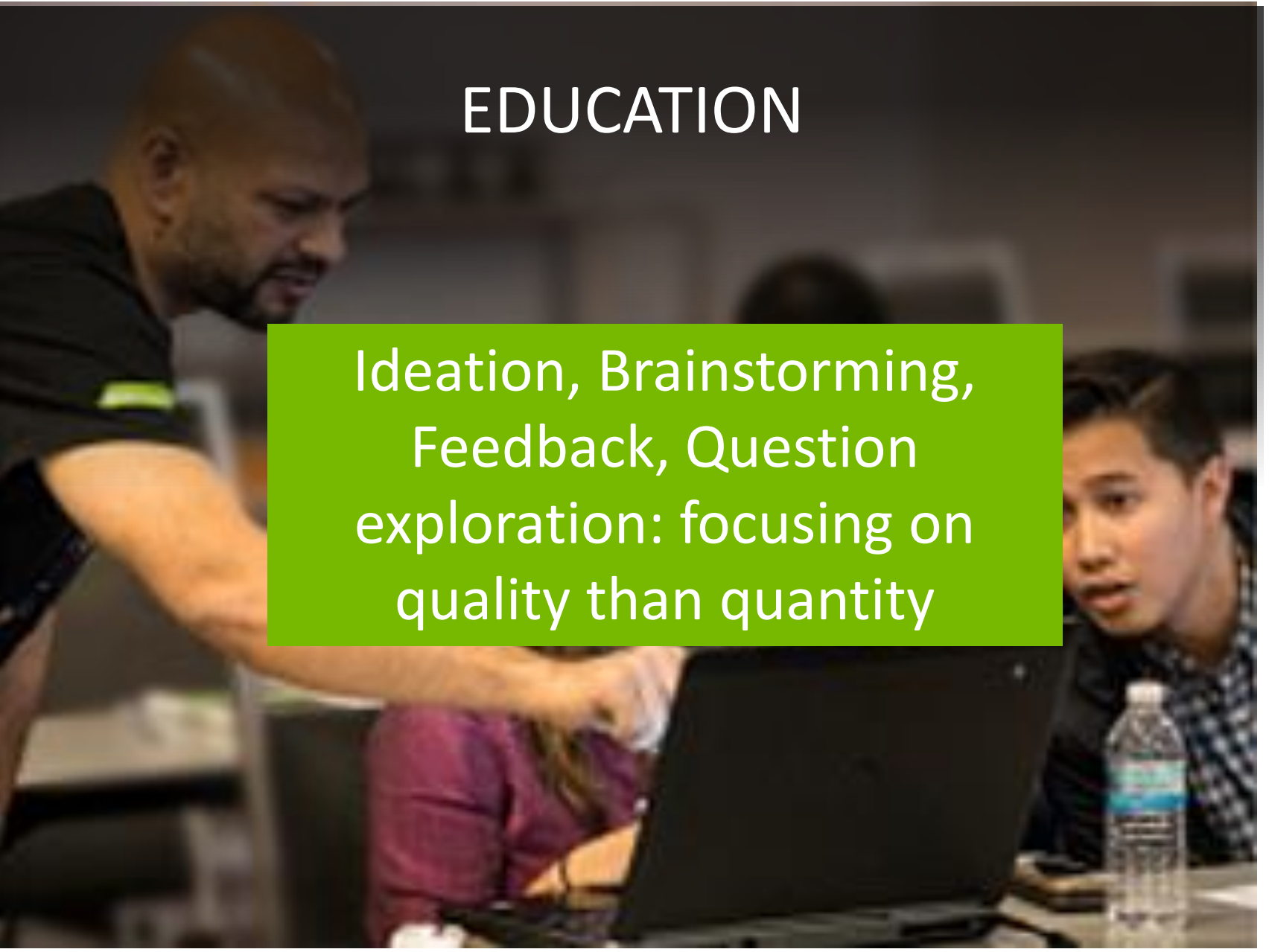MOLECULE
ANIMATION

ANIMATION
MOLECULE
PROTEIN
DNA
VIDEO
3D
IMAGE
AUDIO
TEXT

## Generative AI is Having a Moment!
Captivating World's Attention | Happening Fast | Solving Problems in a Way Never Done Before

**nVIDIA.**

# Benefits of Generative AI on Various Business

Generative AI is transforming every industries



**HEALTHCARE & PHARMA**

Drug Discovery, Predicting patient outcomes, & summarizing Doctor's notes

**EDUCATION**

Ideation, Brainstorming, Feedback, Question exploration: focusing on quality than quantity

**eCOMMERCE**

Virtual shopping assistant & tailored messaging

**MARKETING**

Copy generation, New age advertising, Headline generation, Personalized images and text

**RETAIL**

Deliver personalized messages with relevant products, analyze customer messages

**MEDIA & ENTERTAINMENT**

Improve efficiency of running new ideas

NVIDIA.

# Generative AI Risks

Challenges/Limitations of Generative AI

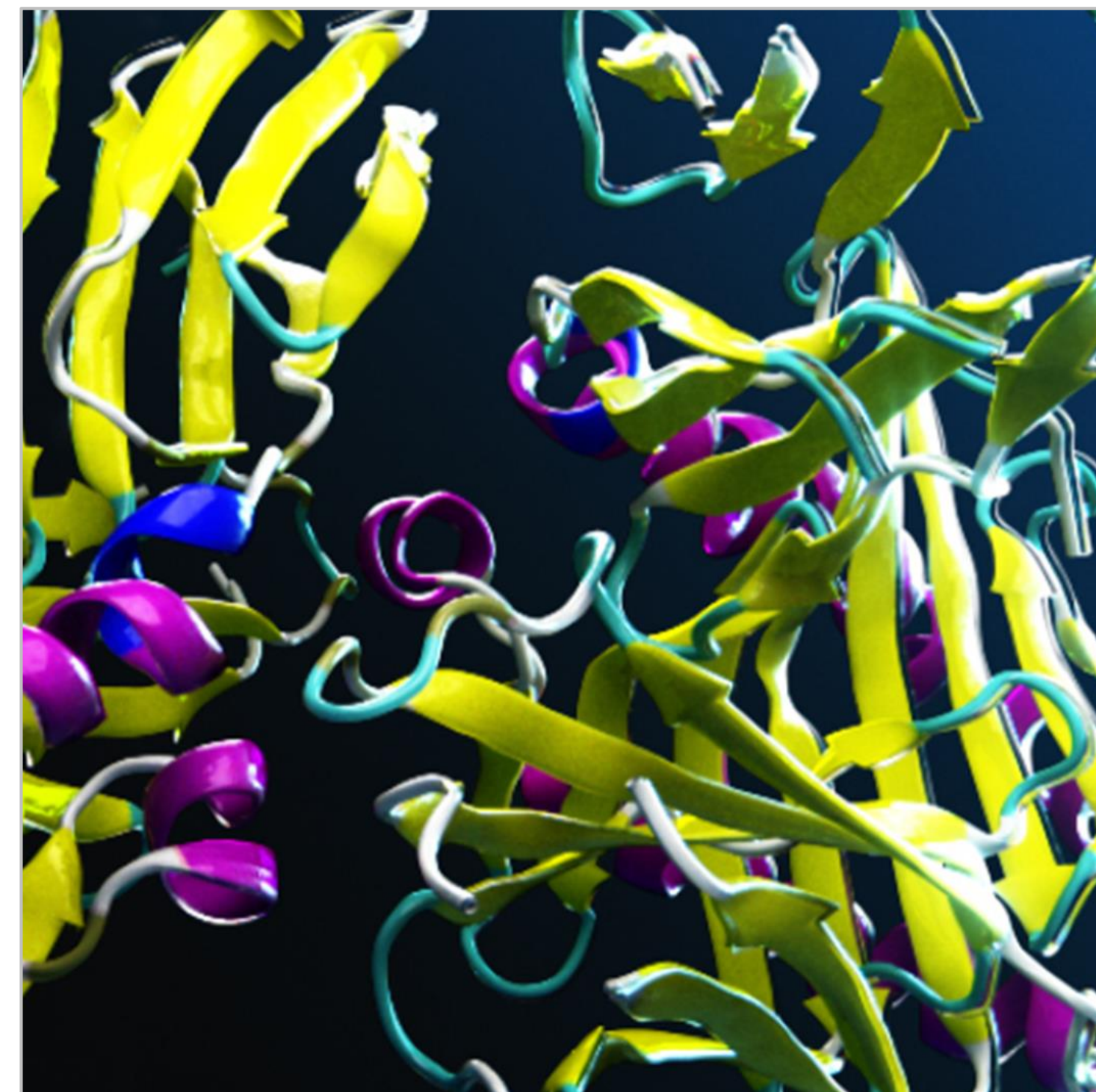| Data Privacy & Security | IP Rights & Copyright | Biases, Errors, & Limitations | Ethical Implication | Malevolent Activities |

# NVIDIA AI Foundations
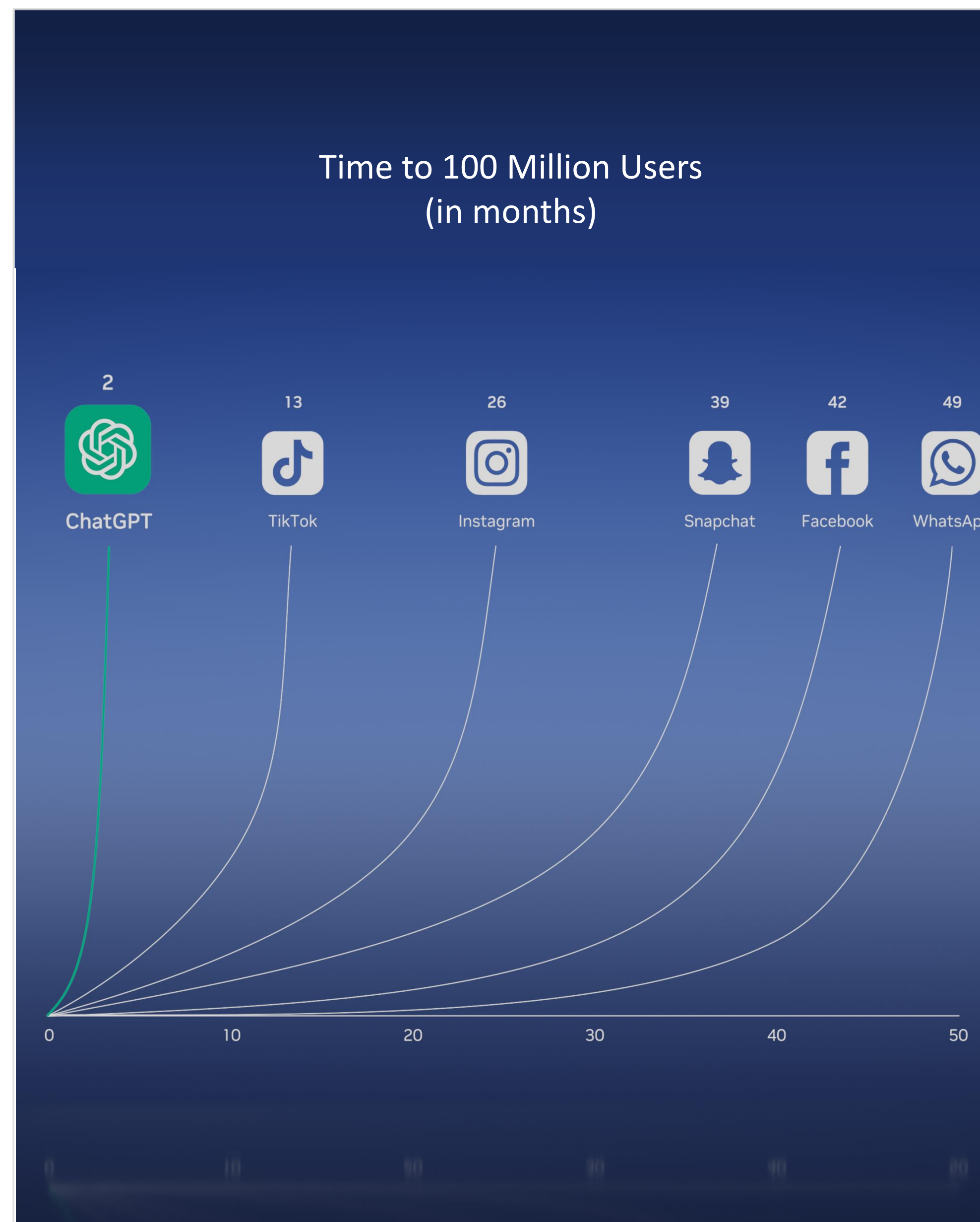


NEMO



BIONEMO



PICASSO

NVIDIA AI ENTERPRISE

NVIDIA DGX Cloud

# The iPhone Moment of AI is Here
## Every major application and workflow is going to include AI
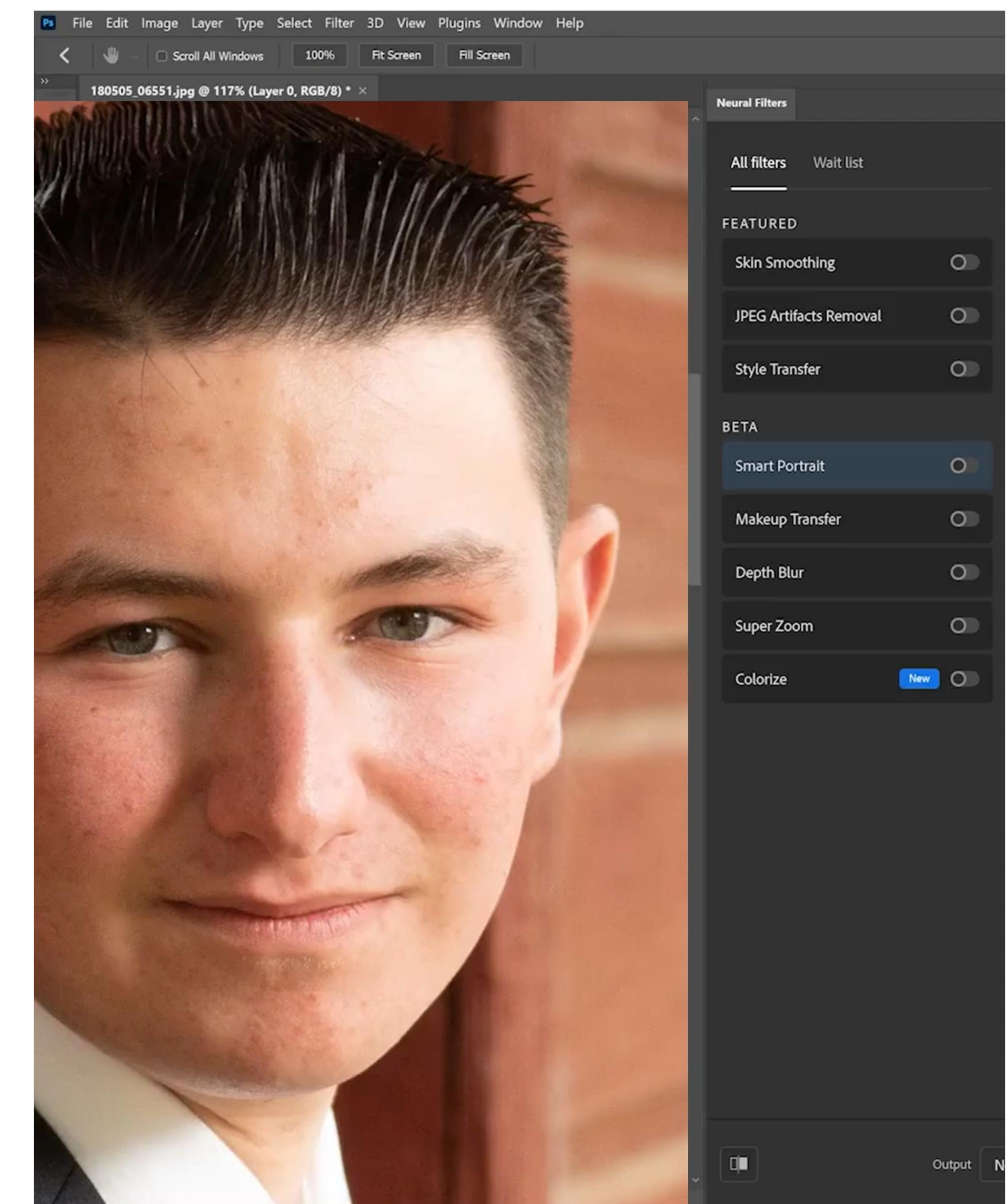


**CHATBOTS**
Fastest Growing Application Ever

**GENERATIVE ART**
Over 200M+ Users

**AI-AUGMENTED APPLICATIONS**
ISVs Accelerating AI Integration

# Generative AI Technology Preview

# AI Playground on NGC

## Experience Top Generative AI Models via Simple Web UI

Explore top NVIDIA & OSS AI models on NGC

Experience on DGX Cloud via AI Playground

One-click install to RTX PC

Fine tune with AI Workbench
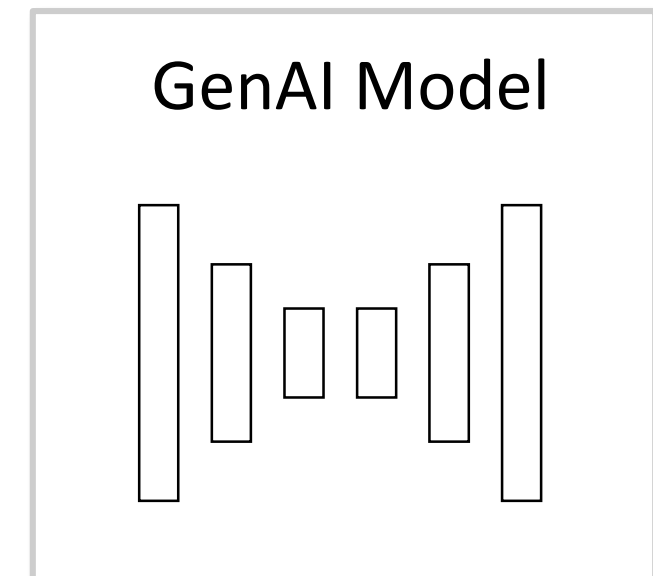


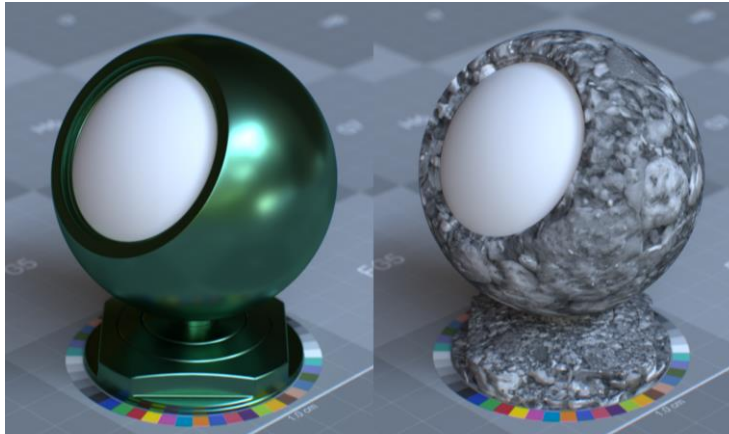**Browser-based UI**

| DGX Cloud | Windows RTX PC / WS |
|-----------|---------------------|

# AI Playground on NGC

## Demo

# Populating 3D Worlds With Generative AI

# Adobe Photoshop Example
## Adding Content with Generative Fill

# Generative AI Transforming Workflows Across Industries



Architecture

Product Design

Film / Video

3D FX / Game Dev

Marketing

Photography

# September Early Access for NVIDIA AI Workbench

Announcing at SIGGRAPH

- New development platform for data scientists & AI/ML engineers

- Brings Gen AI development to RTX GPUs on PCs & Linux

- Streamlines NVIDIA HW & SW for a broad audience

- Integrates with NGC, GitHub & Hugging Face

- One-click push of local work into DGX Cloud

- Easy hybrid deployment on local machines & clouds

- Free and available for self-service install from NGC

- EA starting in September

# NVIDIA AI Enterprise

## End to end AI software

**MLOps**

**AI Applications**

**NVIDIA AI Enterprise**

### Workload and Infra Management

**Model Deployment**
Triton Management Service

**Cloud Native Management and Orchestration**
GPU Operator/Network Operator

**Cluster Management**
Base Command
Manager Essentials

**Infra Acceleration Libraries**
Magnum IO, vGPU, CUDA

### AI Use Cases and Workflows

Hello
Nemo for Gen AI/LLM

Riva for Speech AI

Metropolis for Video Analytics

MONAI for Medical Imaging

Morpheus for Cybersecurity

CuOpt for Route Optimization

More

### AI Development

**Data Prep**
RAPIDS

**Model Training**
TAO, PyTorch/TensorFlow

**Deploy at Scale**
Triton Inference Server

**Simulate and Test**
TensorRT

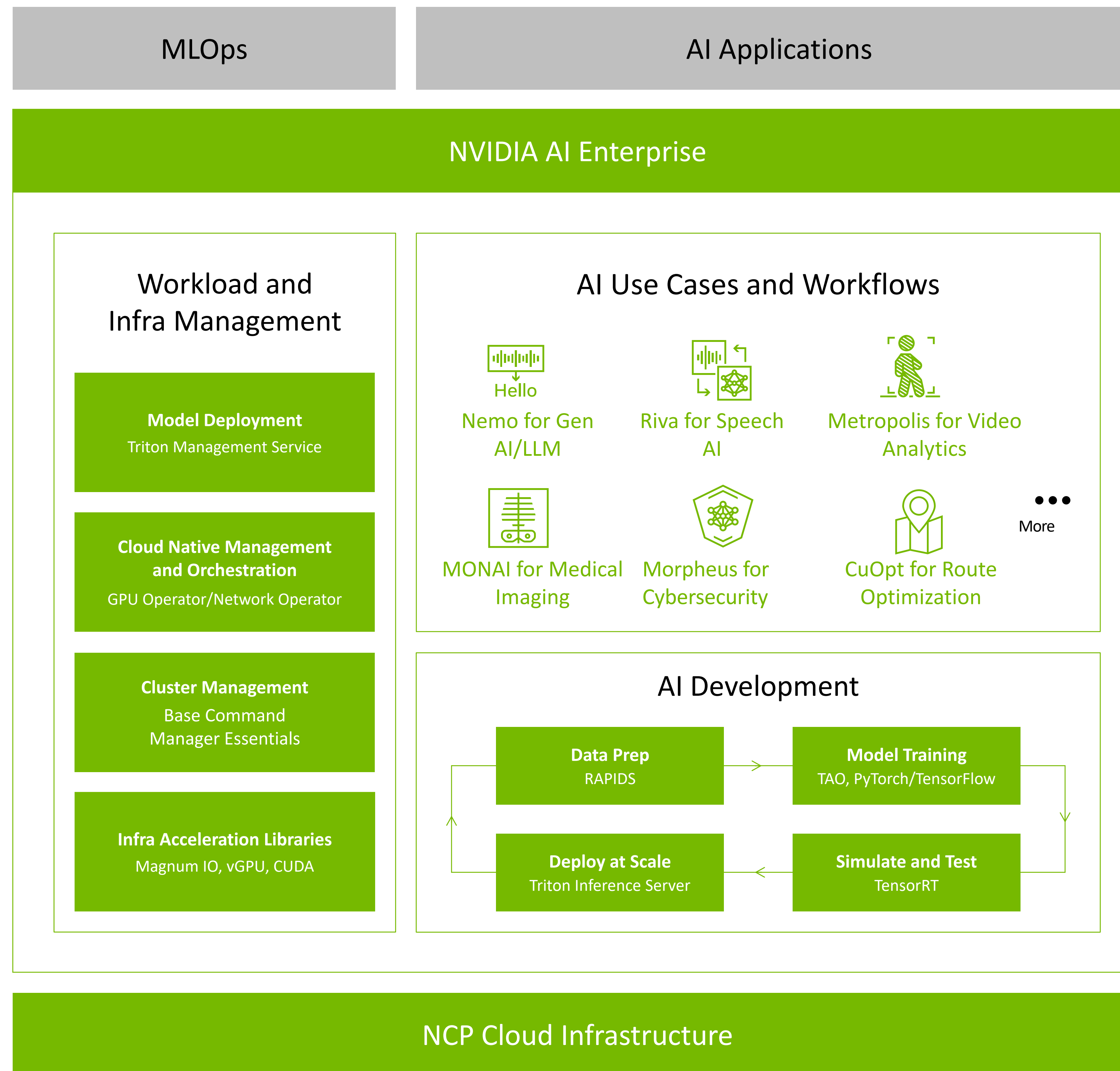**NCP Cloud Infrastructure**

- Cloud Native, Hybrid Optimized
- Enterprise Grade OSS Components
- Secure And Scalable
- Technical Standard Support 9 X 5, Premium 24x7

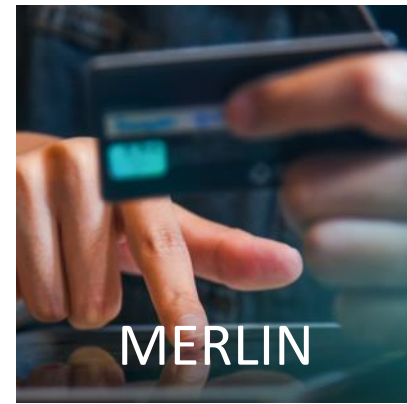# NVIDIA Software

# Productivity with Base Command

## Enterprises tools that drive the value of AI investment

**NVIDIA AI Enterprise**

AI and Data Science Tools / Frameworks

NVIDIA RAPIDS™ | NVIDIA TAO Toolkit | NVIDIA Tensor RT™ | AI Frameworks | NeMo Megatron

- **NVIDIA AI Enterprise** (included in DGX SW Bundle)

Ready-to-use, fully supported software that **speeds developer success**

## NVIDIA Base Command

### AI Workflow Management and MLOps

- **DGX-Ready Software Partners**

Get more models from **prototype to production**

### Job Scheduling & Orchestration

KUBERNETES | SLURM

- Kubernetes
- Slurm

Ensure hassle-free **execution of every developer's jobs**

### Cluster Management

PROVISIONING | MONITORING | CLUSTERING | MANAGING

- **Bright Cluster Manager**
- **Base Command Manager** (DGX SuperPOD)

Effortlessly **scale and manage** one node to thousands

### Network/Storage Acceleration Libraries & Management

NETWORK IO | STORAGE IO | IN-NETWORK COMPUTE | IO MANAGEMENT

- **Magnum IO, UFM, NetQ**

Accelerate end-to-end **infrastructure performance**

### Operating System
DGX OS Extensions for Linux Distributions

- **DGX OS**

**Maximize system uptime security** and reliability

Single Node → Multi-Node

DGX H100 SuperPOD
NVIDIA Reference Architectures

# Introducing Base Command Manager Essentials

Purpose-built for Enterprise AI Infrastructure Management

## Infrastructure Provisioning

- Maintain a secure, up-to-date, and reliable AI infrastructure

## Workload Management

- Easily provide data scientists with all the tools and resources they need

## Resource Monitoring

- Get detailed data for informed decision-making

Data Center

Cloud

NVIDIA.

# Powerful Image Management and Cluster Configuration

Greatly increases admin productivity and prevents configuration drift



Head Node
(HA deployment)

Group 1

Group 2

Group 3

**Provision systems by categories**

- Share common software stack

**Software stack synchronization**

- Occurs at each reboot

**Post-install configuration**

- Networking & Security
- User authentication
- GPU settings
- BMC & Power settings
- *and more*

# Selected NVIDIA GPUs

# NVIDIA Hopper

The Engine for the World's AI Infrastructure



World's Most Advanced Chip

Transformer Engine

2nd Gen MIG

Confidential Computing

4th Gen NVLink

DPX Instructions

H100 SXM

H100 NVL

H100 PCIE

Includes NVIDIA AI Enterprise

# NVIDIA H100

Unprecedented Performance, Scalability, and Security for Every Data Center

### Highest AI and HPC Performance

4PF FP8 (6X)| 2PF FP16 (3X)| 1PF TF32 (3X)| 60TF FP64 (3.4X)

3.35TB/s (1.5X), 80GB HBM3 memory

### Transformer Model Optimizations

6X faster on largest transformer models

### Highest Utilization Efficiency and Security

7 Fully isolated & secured instances, guaranteed QoS

2nd Gen MIG | Confidential Computing

### Fastest, Scalable Interconnect

900 GB/s GPU-2-GPU connectivity (1.5X)

128GB/s PCI Gen5

*FP8, FP16, TF32 performance include sparsity. X-factor compared to A100*

# Delivering the AI Center of Excellence for Enterprise

Best-of-breed infrastructure for AI development built on NVIDIA DGX

## NVIDIA DGX H100

The World's Proven Choice for Enterprise AI



8x NVIDIA H100 GPUs  |  32 PFLOPS FP8 (6X)  |  0.5 PFLOPS FP64 (3X)
640 GB HBM3 | 3.6 TB/s (1.5X) BISECTION B/W

4th Generation of the World's Most Successful Platform
Purpose-Built for Enterprise AI

## DGX SuperPOD WITH DGX H100



32 DGX H100  |  1 EFLOPS AI
QUANTUM-2 IB  |  20TB HBM3  |  70 TB/s BISECTION B/W (11X)

1 ExaFLOPS of AI Performance in 32 Nodes
Scale as Large as Needed in 32 Node Increments

# NVIDIA L40S

The highest performance universal GPU for AI, graphics, and video

| | | |
|---|---|---|
| **LLM Fine Tuning** | **LLM 1st Token Latency** | **Small LLM Training** |
| **8 hours** | **<30 ms** | **3 days** |
| Llama 2-70B 1 Billion Tokens[1] | Llama 2-13B Inference 225/20[2] | Llama 2-7B 100 Billion Tokens[3] |
| **Image Gen AI** | **LLM E2E Latency** | **Full Video Pipeline** |
| **>82** | **<0.5 s** | **184** |
| Images per Minute[4] | Llama 2-13B Inference 225/20[5] | AV1 Encode Streams[6] |

Dual-Slot l FHFL l 350W

Preliminary performance, subject to change
1. Fine-Tuning Llama 2-70B SFT, 1 Billion Tokens; 64x L40S (simulated).
2. Llama 2-13B, ISL=225, OSL=20, BS=1, FP8, 1x L40S. 1st Token Latency.
3. Llama 2-7B, 100 Billion Tokens; 64x L40S (simulated).
4. Image Generation, Stable Diffusion v2.1, 50 iterations, 512 x 512 resolution; 1x L40S.
5. Llama 2-13B, ISL=225, OSL=20, BS=1, FP8, 1x L40S. E2E Latency.
6. Concurrent Encoding Streams; 720p30; 1x L40S.

# L40S Generates > 82 Images/min for Image Gen AI Inference

## Incredible Performance Across Different Image Sizes and Resolutions



Absolute Performance

| | SD (512 x 512) | SD (1024 x 1024) | SDXL (1024 x 1024) |
|---|---|---|---|
| Images/min | 82 | 17 | 11 |

# NVIDIA Grace for Cloud, AI and HPC Infrastructure

## Grace CPU Superchip
### CPU Computing



CPU-based applications where absolute performance, energy efficiency, and data center density matter, such as scientific computing, data analytics, enterprise and hyperscale computing applications

## GH200 Grace Hopper Superchip
### Large Scale AI & HPC



Accelerated applications where CPU performance and system memory size and bandwidth are critical; tightly coupled CPU & GPU for flagship AI & HPC. Most versatile compute platform for scale out.

NVIDIA.

# NVIDIA GH200 Grace Hopper Superchip

Built for the New Era of Accelerated Computing and Generative AI

**Most versatile compute**
Best performance across CPU, GPU or memory intensive applications

**Easy to deploy and scale out**
1 CPU:1 GPU node simple to manage and schedule for for HPC, enterprise, and cloud

**Best Perf/TCO for diverse workloads**
Maximize data center utilization and power efficiency

**Continued Innovation**
Grace and Hopper-Next in 2024



900GB/s NVLink-C2C | 624GB High-Speed Memory
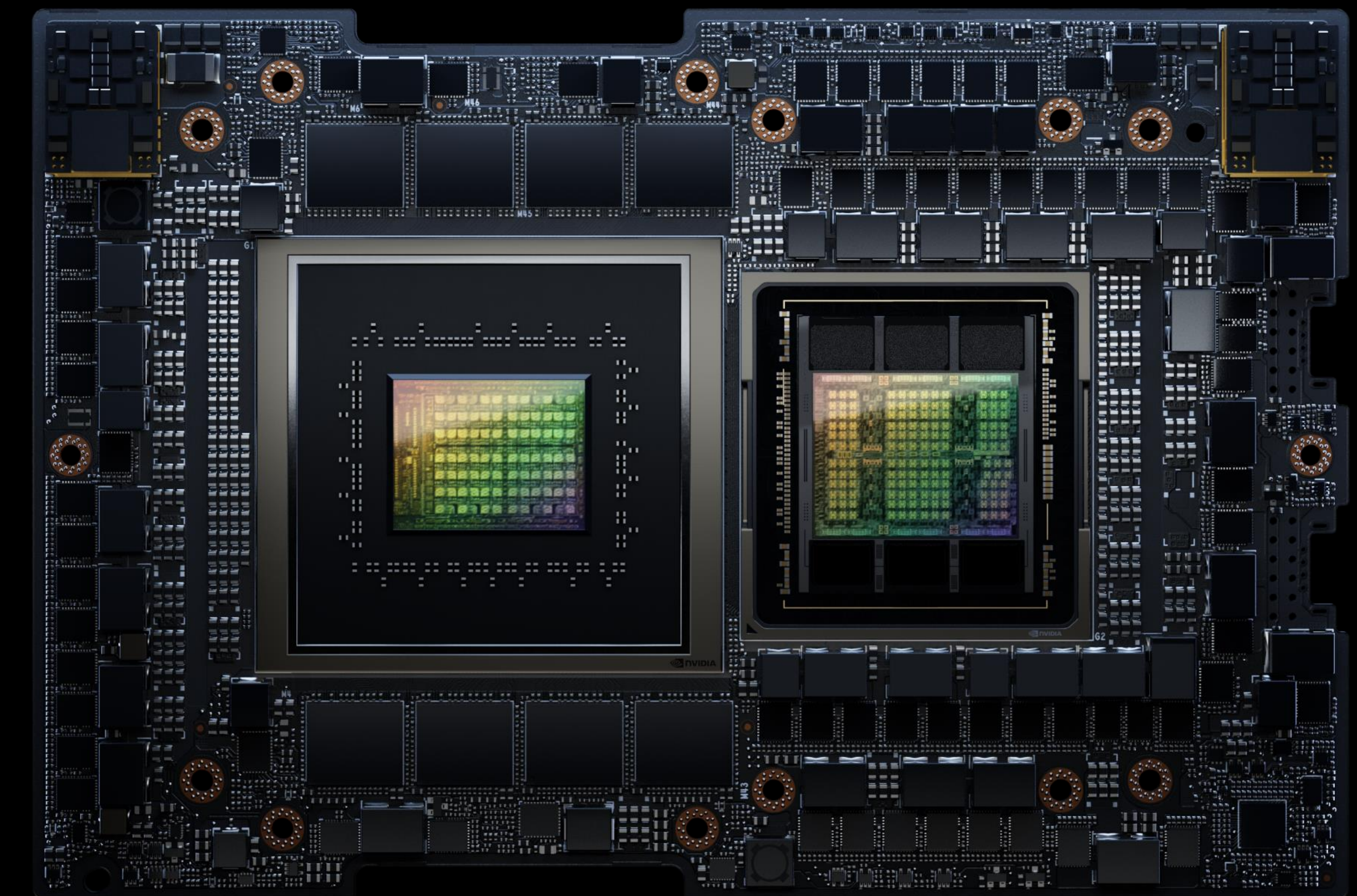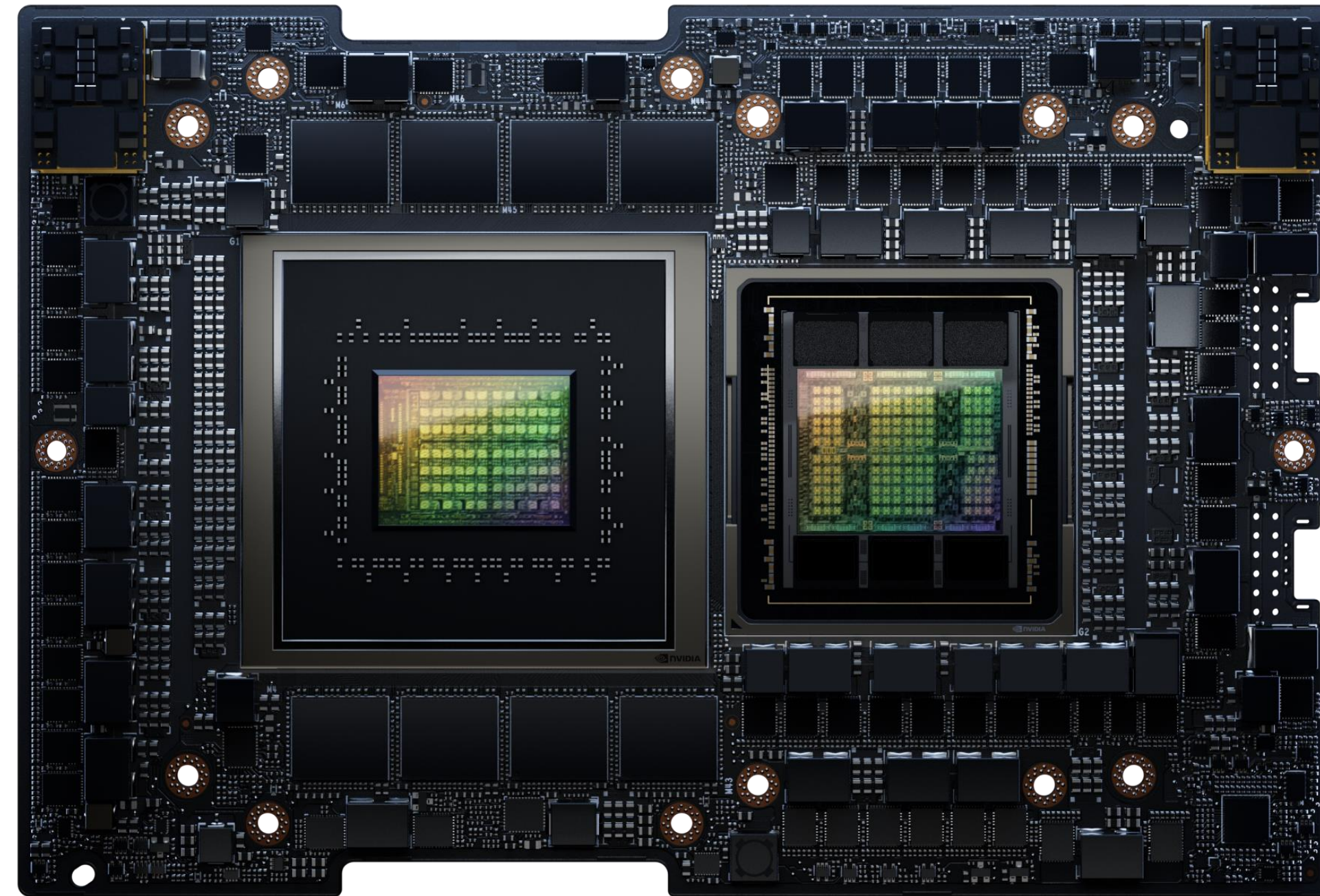4 PF AI Perf | 72 Arm Cores

# NVIDIA GH200 Grace Hopper Superchip
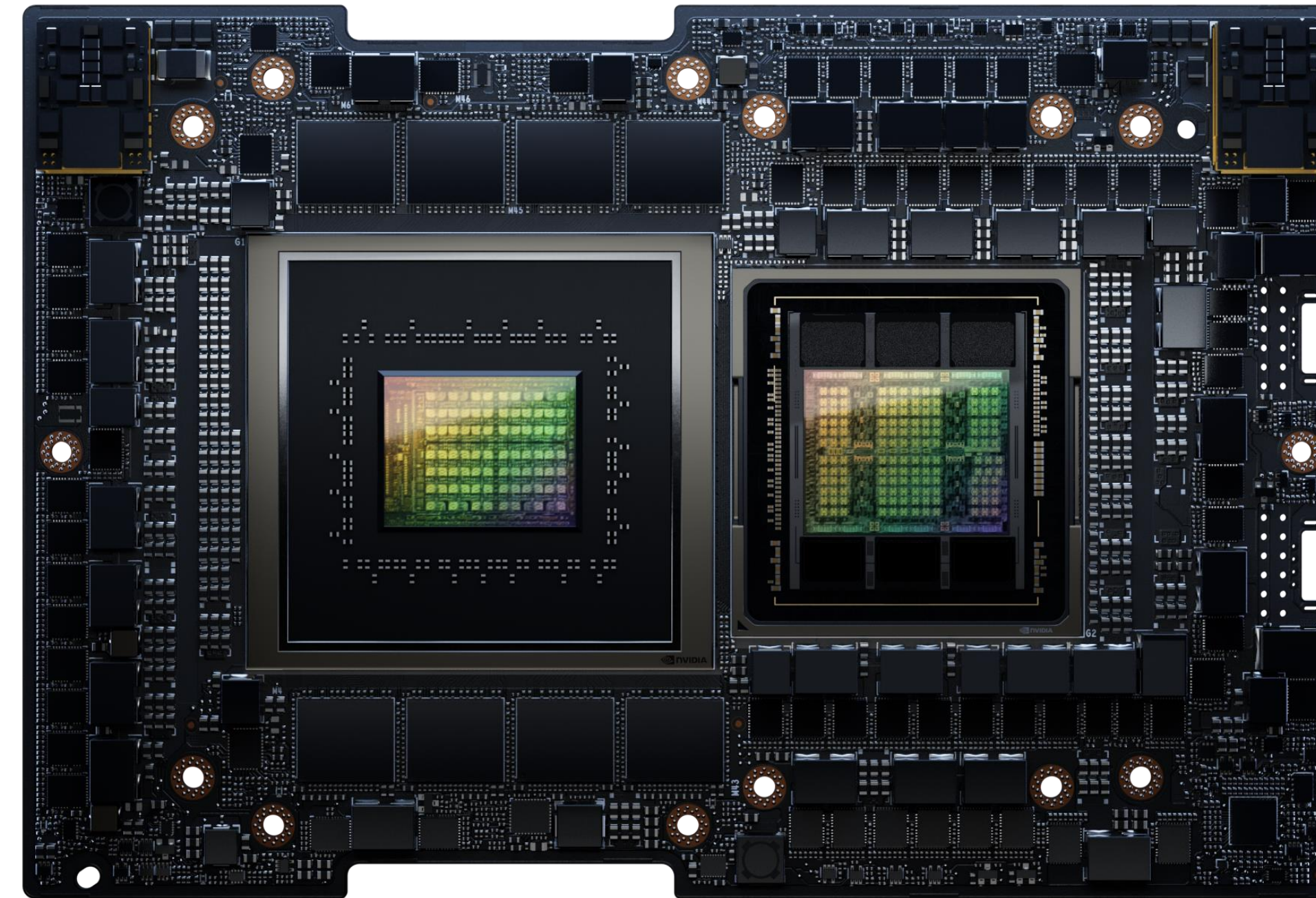
## Processor For The Era of Accelerated Computing And Generative AI



72 Core Grace CPU  |  4 PFLOPS Hopper GPU
96 GB HBM3  |  4 TB/s  |  900 GB/s NVLink-C2C

- 7X bandwidth to GPU vs PCIe Gen 5
- Combined 576 GB of fast memory
- 1.2x capacity and bandwidth vs H100
- Full NVIDIA Compute Stack

### GH200 with HBM3
Available for order



72 Core Grace CPU  |  4 PFLOPS Hopper GPU
144 GB HBM3e  |  5 TB/s  |  900 GB/s NVLink-C2C

- World's first HBM3e GPU
- Combined 624 GB of fast memory
- 1.7x capacity and 1.5x bandwidth vs H100
- Full NVIDIA Compute Stack

### GH200 with HBM3e
Available Late Q2 2024



144 Core Grace CPU  |  8 PFLOPS Hopper GPU
288 GB HBM3e  |  10 TB/s  |  900 GB/s NVLink-C2C

- Simple to deploy MGX-compatible design
- Combined 1.2 TB fast memory
- 3.5x capacity and 3x bandwidth vs H100
- Full NVIDIA Compute Stack

### NVLink Dual GH200 System
Available Q2 2024
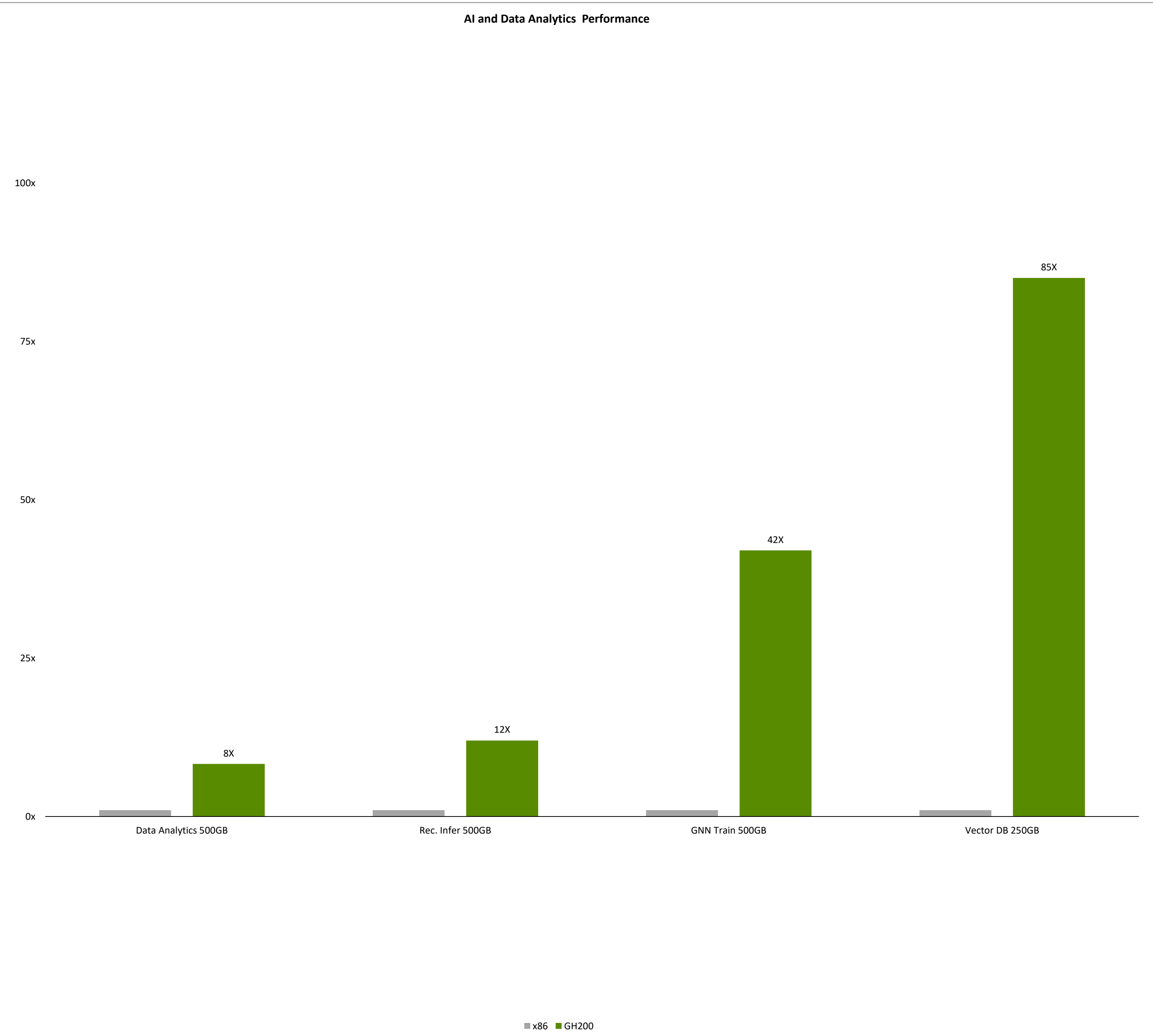
# GH200 Grace Hopper AI Inference Platform

## Versatile Scale Out with Unmatched Performance

### Memory Intensive

**AI and Data Analytics Performance**

- 100x
- 85X — Vector DB 250GB
- 75x
- 50x
- 42X — GNN Train 500GB
- 25x
- 12X — Rec. Infer 500GB
- 8X — Data Analytics 500GB
- 0x

Legend: ▦ x86  ■ GH200

### GPU Intensive

**LLM Infer Performance**

- 225X
- 203X — LLAMA 2 70B
- 200X
- 175X
- 150X
- 125X
- 100X
- 75X
- 50X
- 25X
- 1X
- 0X

Legend: ▦ x86  ■ GH200 144GB

### Use Cases

**LLM**
Conversational AI
Domain Knowledge



**Recommender Systems**
eCommerce
Personalized Content



**Vector Database**
Fraud Detection
Drug Discovery



**GNN**
Computer Vison
Recommenders

NVIDIA.

# NVIDIA Networking

# NIC
## Network Interface Card

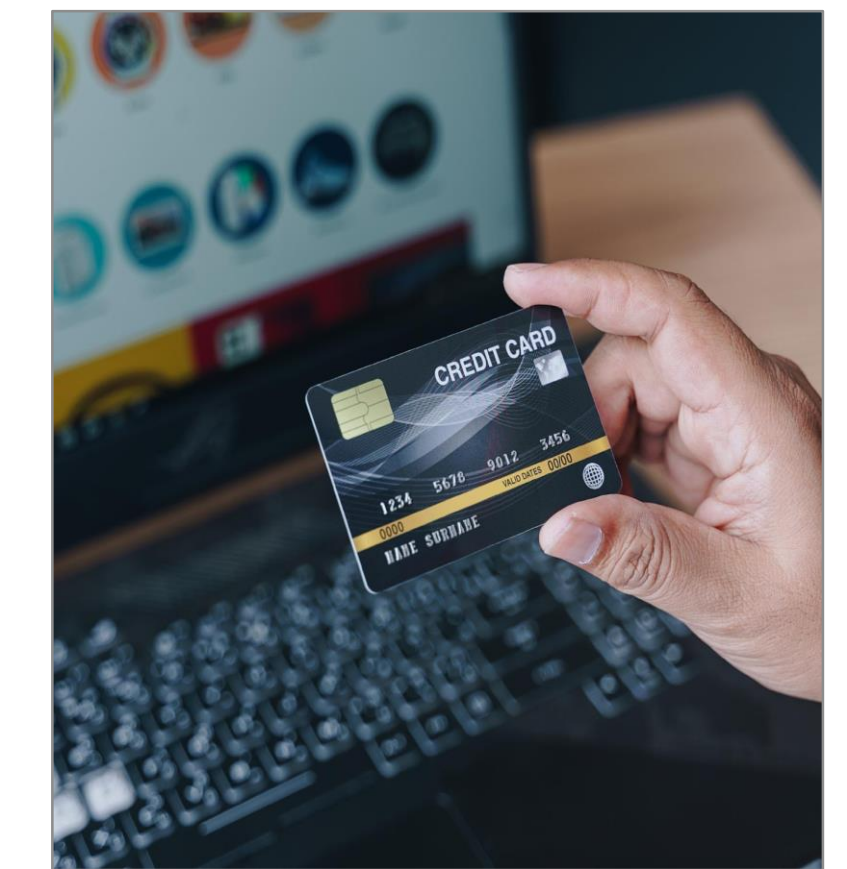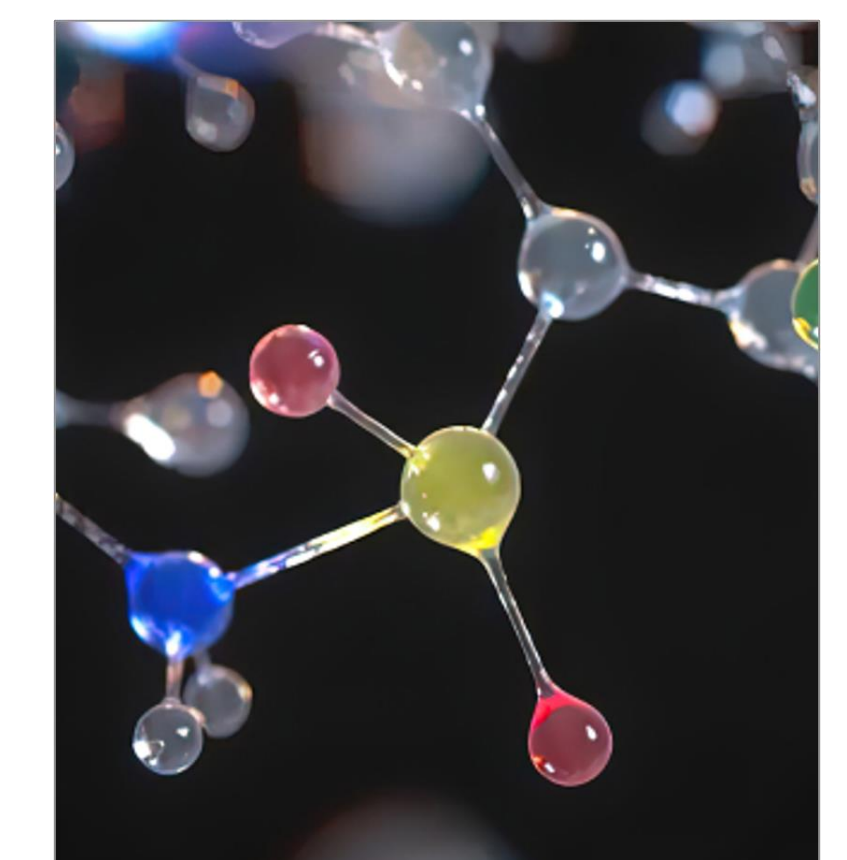A network interface controller (a.k.a. network adapter) is a computer HW component that connects a computer to a computer network.

Early network interface controllers were commonly implemented on expansion cards that plugged into a computer bus…

Modern network interface controllers offer advanced features such as interrupt and DMA interfaces to the host processors, support for multiple receive and transmit queues, …

*https://en.wikipedia.org/wiki/Network_interface_controller*



2019 ConnectX-6 200gE NIC



2020 ConnectX-7 400gIB/E NIC



2007 ConnectX-EN 10gE NIC – First Ethernet NIC



2001 InfiniBridge 10g IB HCA – Mellanox First
Host Channel Adapter (HCA)

# Recommended Compute Nodes

NVIDIA DGX Systems & Qualified and Certified HGX H100 Servers



## NVIDIA ConnectX-7 Networking Adapters

The ConnectX®-7 networking adapter is the latest ConnectX® line. It can provide 25/50/100/200/400Gb/s of throughput. NVIDIA recommended systems use the single-port ConnectX-7® to provide best performance in deployments with NDR. Specifications are available here.

# NVIDIA Networking Platforms

Accelerated Networking Solutions for the Era of AI



Quantum-2 InfiniBand

BlueField-3 Networking

Spectrum-X Ethernet

**Supercomputing Networking Platform**
AI Factories and Cloud-Native Supercomputing

**Infrastructure Compute Platform**
Offload, Accelerate, and Isolate Data Center Infrastructure

**Hyperscale Networking Platform**
Purpose-built Ethernet Networking for AI Clouds

# TCP and Sockets

RFC: 793

TRANSMISSION CONTROL PROTOCOL

DARPA INTERNET PROGRAM

PROTOCOL SPECIFICATION

September 1981

prepared for

Defense Advanced Research Projects Agency
Information Processing Techniques Office
1400 Wilson Boulevard
Arlington, Virginia  22209

by

Information Sciences Institute
University of Southern California
4676 Admiralty Way
Marina del Rey, California  90291

RFC 793 – Transmission Control Protocol (TCP)

```
SocketFD = socket(PF_INET, SOCK_STREAM, IPPROTO_TCP);
 if (SocketFD == -1) {
   perror("cannot create socket");
   exit(EXIT_FAILURE);
 }

 memset(&sa, 0, sizeof sa);

 sa.sin_family = AF_INET;
 sa.sin_port = htons(1100);
 res = inet_pton(AF_INET, "192.168.1.3", &sa.sin_addr);

 if (connect(SocketFD, (struct sockaddr *)&sa, sizeof sa) == -1) {
   perror("connect failed");
   close(SocketFD);
   exit(EXIT_FAILURE);
 }

 n = write(SocketFD,buffer,strlen(buffer));
 if (n < 0)
     error("ERROR writing to socket");

 n = read(SocketFD,buffer,MAX_BUFFER_SIZE);
 if (n < 0)
     error("ERROR reading from socket");

 shutdown(SocketFD, SHUT_RDWR);

 close(SocketFD);
```

BSD/POSIX Sockets
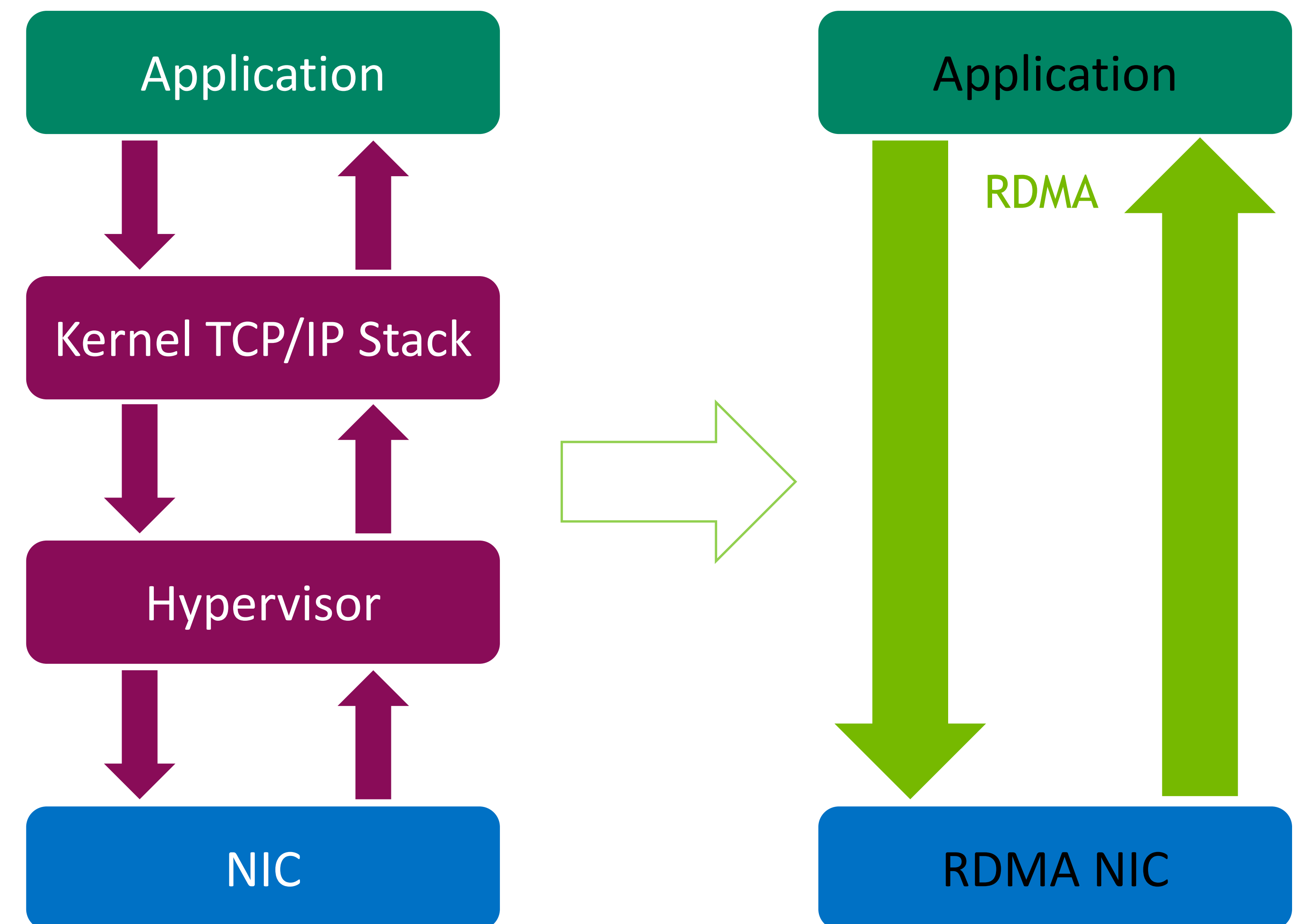
# RDMA – Remote DMA
## InfiniBand and RoCE (RDMA over Converged Ethernet)

- ▸ Messages send/receive, remote DMA and remote atomics

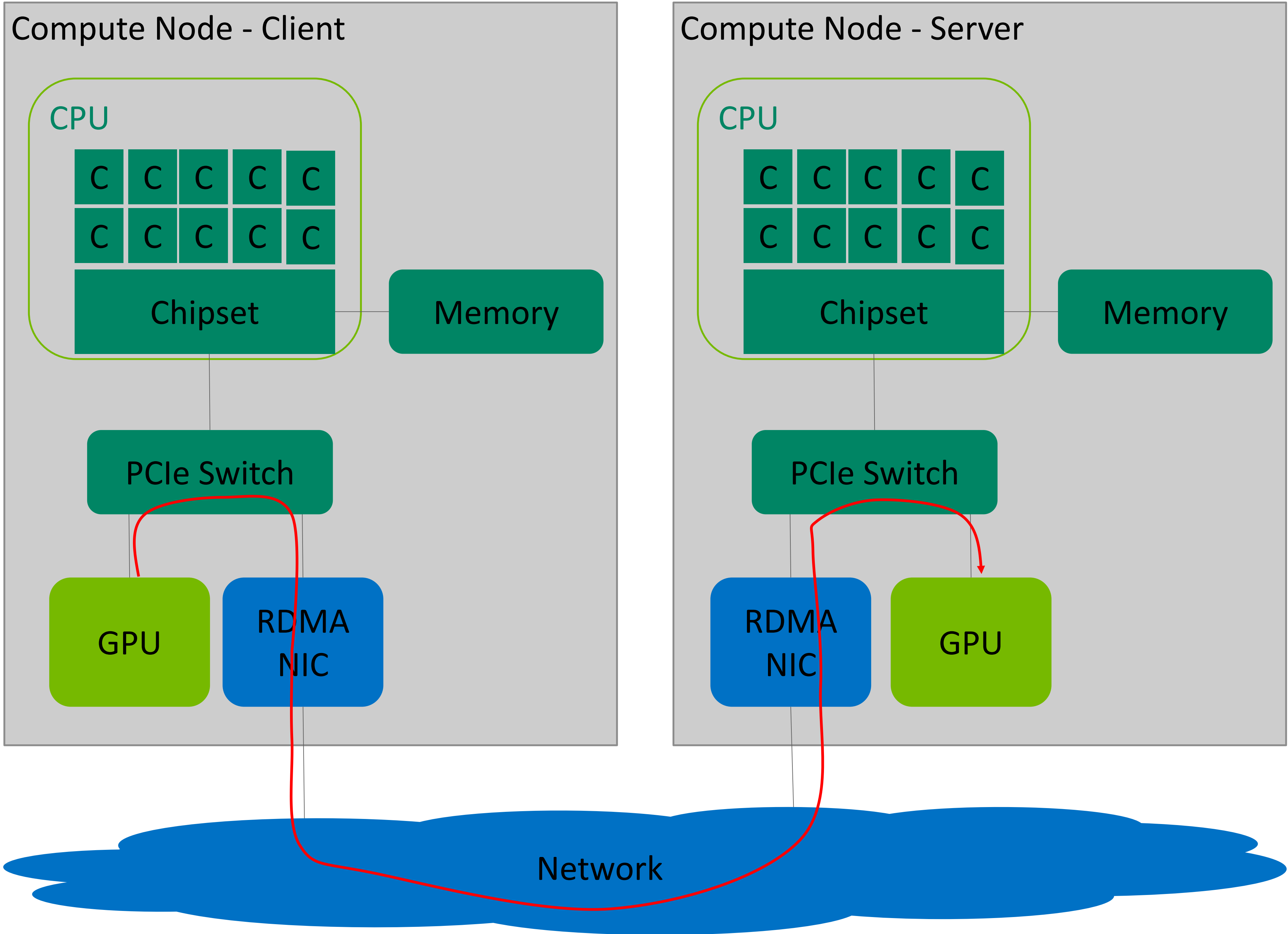- ▸ Hardware transport

- ▸ Kernel and hypervisor bypass

Advantages

- ▸ Lower latency – 10us → 700ns

- ▸ Higher message rate – 215M messages/s

- ▸ Lower CPU utilization – 0%



Application

Kernel TCP/IP Stack

Hypervisor

NIC

Application

RDMA

RDMA NIC

**NVIDIA.**

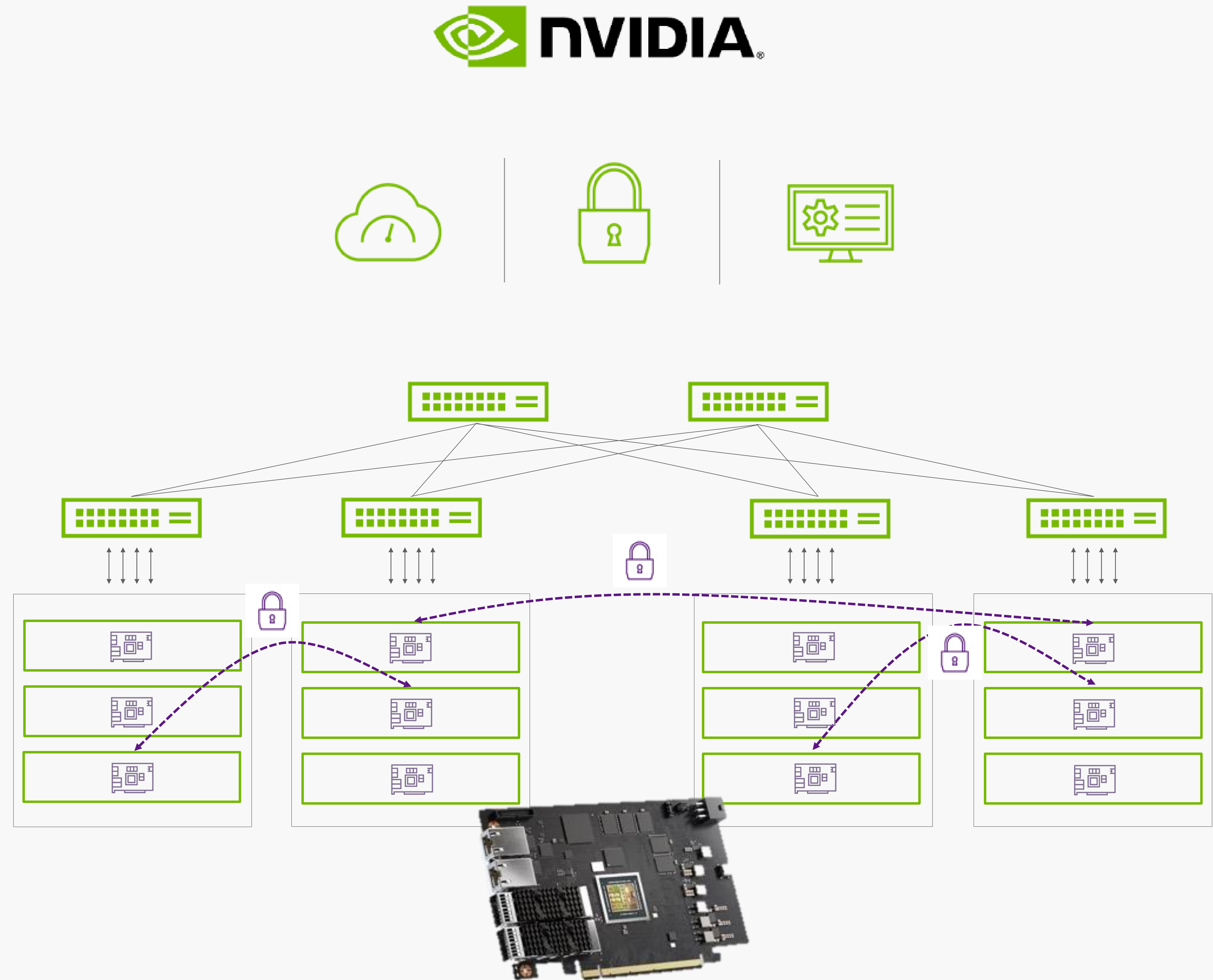# RDMA enables Peer to Peer
## GPU Direct

# DOCA Host-based networking

## Focus areas

- Controller-less VPC networking for BMaaS

- Accelerated Routing and EVPN on Host

- Simplifying underlay network fabric
  - End-to-end IP Fabric from the host
  - L3 ECMP to replace proprietary LAG/MLAG

- Advanced ToR switch features on DPU
  - Reduce risk with Whitebox deployments (SONiC, Switchdev)
  - Save on 3rd party switch feature licenses

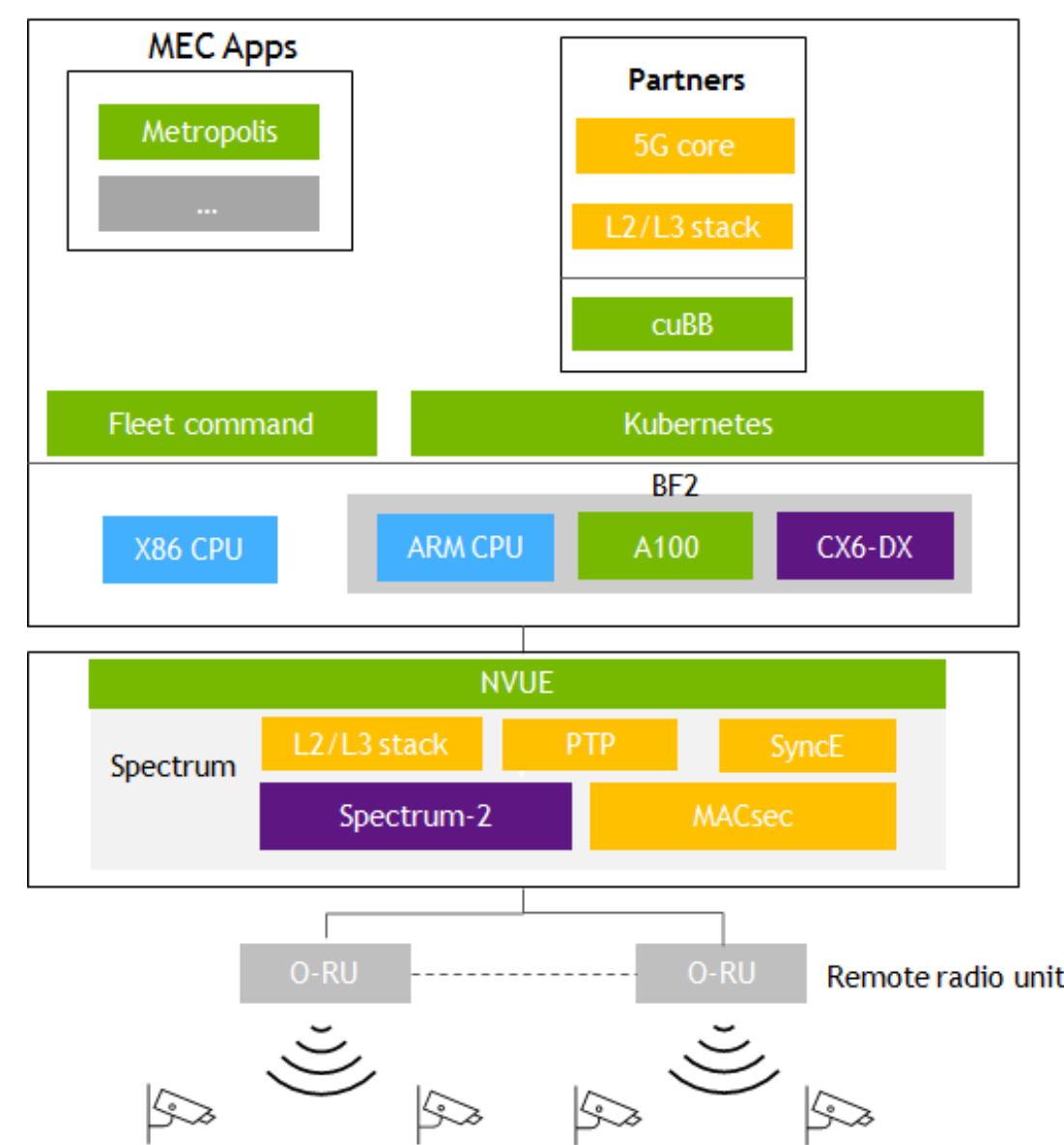- NVIDIA E2E value-add with Cumulus and HBN



NVIDIA Cumulus Apps on DPU
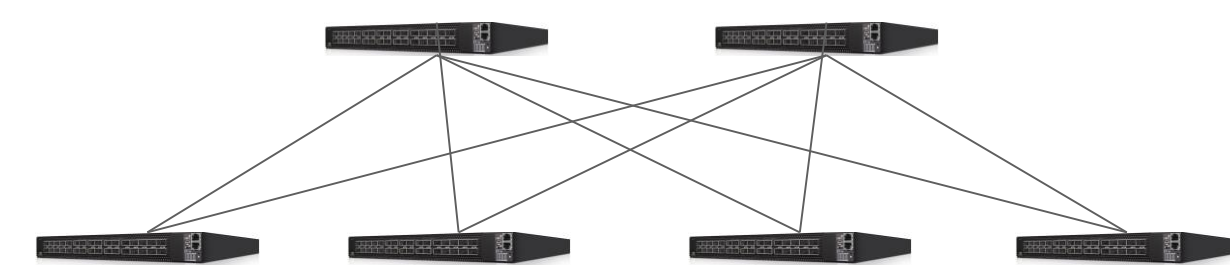
# NetQ in Full Stack NVIDIA Solutions

## Telemetry Data Collection, Network Monitoring and Troubleshooting



### Aerial (EGX POD)

### DGX/HGX

### OVX POD / Super POD

### Host Based Networking (HBN)

### AI Powered by Morpheus Pipeline

- Switch & DPU inventory
- Validations

- Metrics monitoring
- RoCE monitoring

- WJH and other events
- Trace & flow telemetry

- Network data collection
- APIs for integration

# Getting Started with NVIDIA Reference Architectures

# Reference Architecture Components
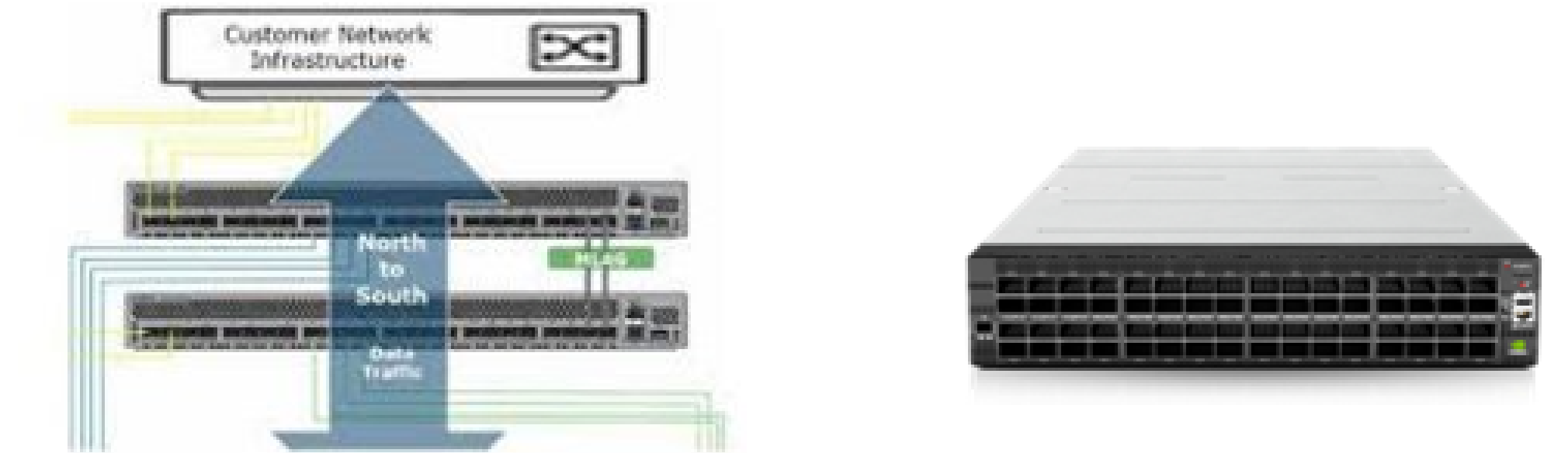
## Compute Node



- HGX H100 8 GPUs Compute node design
- Ensures performance within the compute node

## E-W Network



- 8-Rails-optimized, non-blocking fat tree IB topology
- Validated with 3 production generations of SuperPOD
- Highly scalable with large locality groups
- SharpV3 for multi tenancy
- Managed by UFM & BCM
- Easy to validate and operate

## Tenant Access Network: N-S



- Each HGX has 2 x 200GbE BF3 connections to Leaf switches
- BF3 enables a true, zero-trust BMaaS
- Infrastructure isolation from the compute workload
- Enhanced security
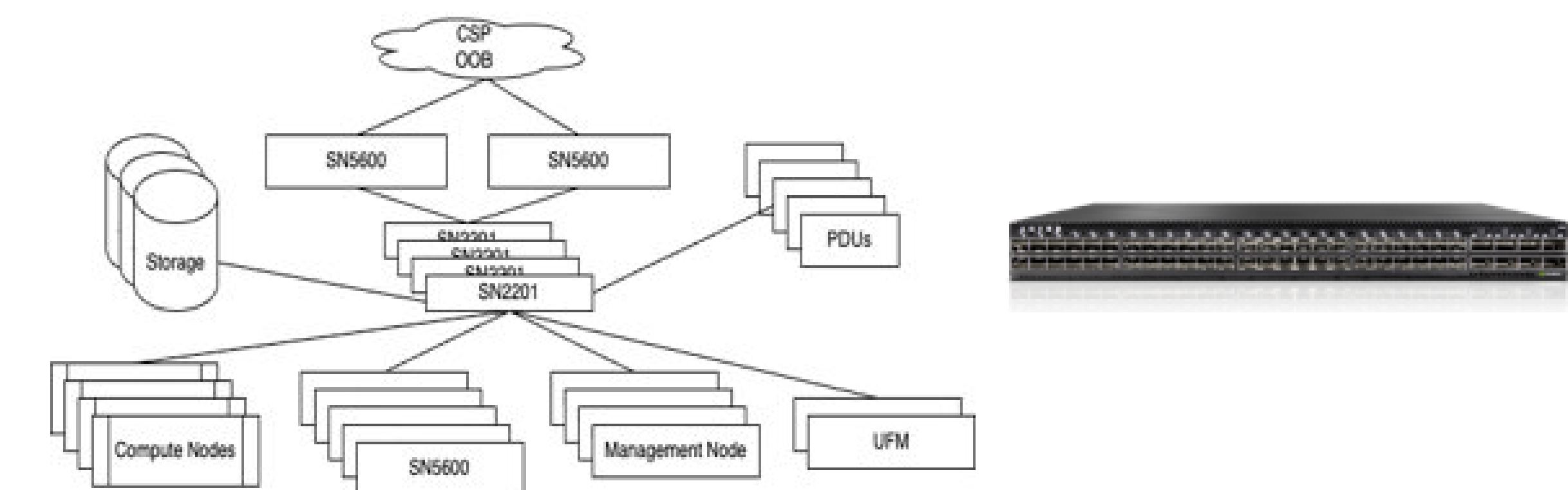- Storage offloaded and accelerated by BF3

## High Speed Storage (HSS)



- Clear minimum recommended bandwidth for storage
- Optimized for best performance of LLM Training
- Enabling a choice of storage solutions
- Per POD or centralized HSS placement

## Outer Ring Storage (Data Lake)



- Raw-data storage
- Minimal requirements for a Datalike storage
- Can be a unified solution with HSS
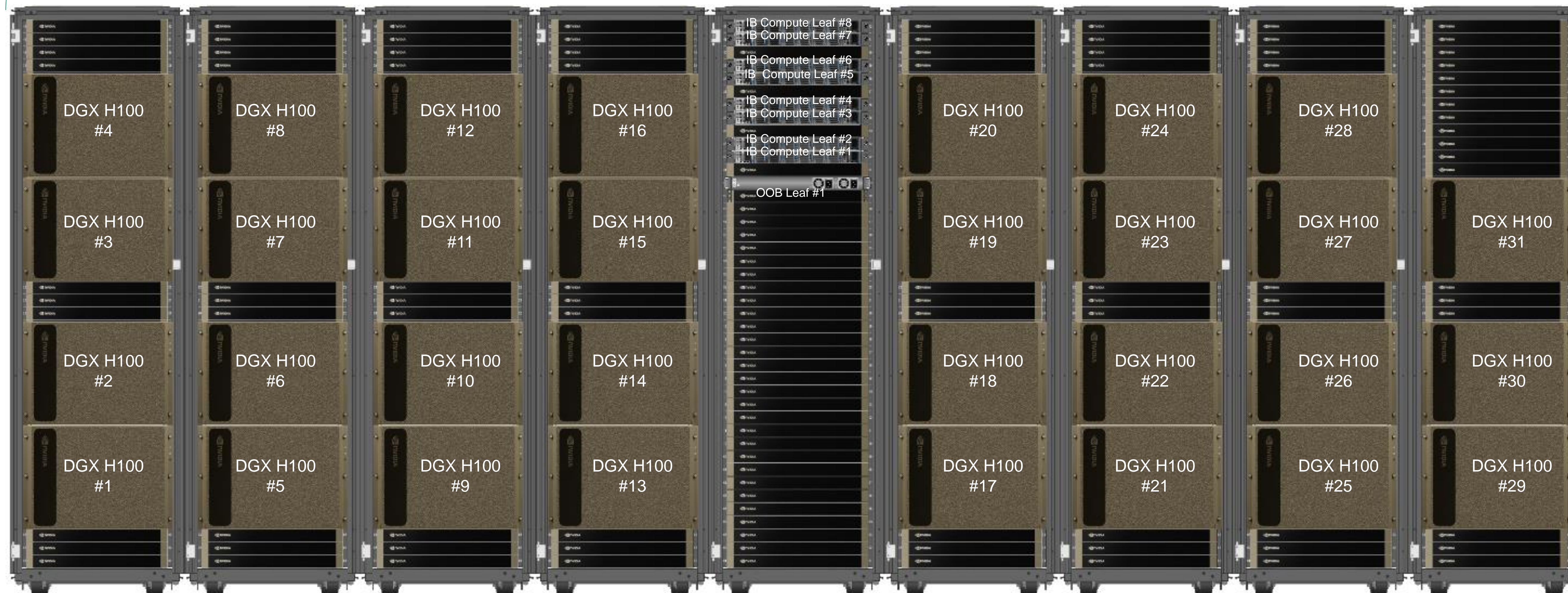
## Out Of Band Management



- 1GbE management network
- Providing monitoring and management of all DC devices
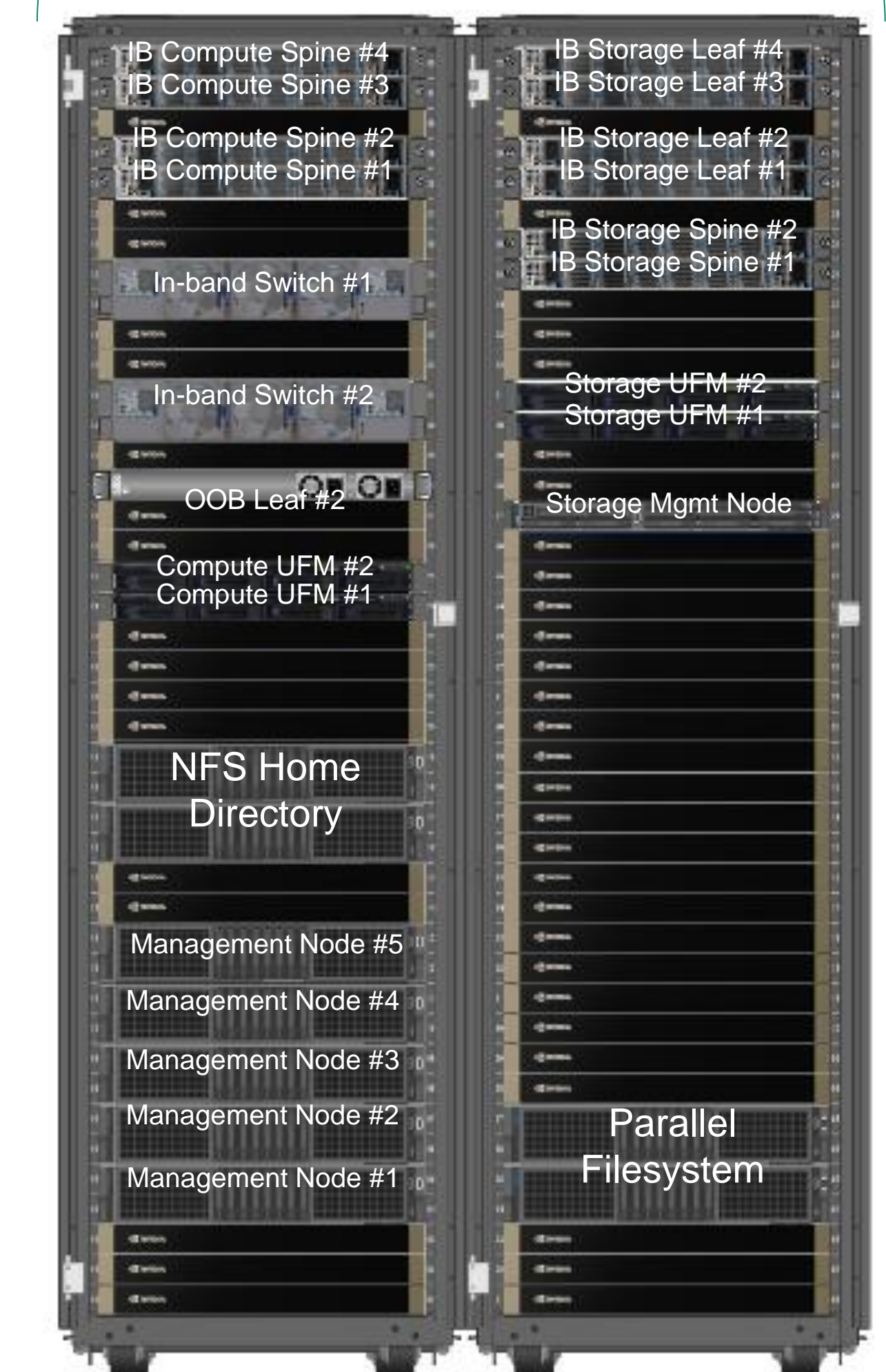- Enabling integration of different APIs between all devices

# Flagship for AI Training
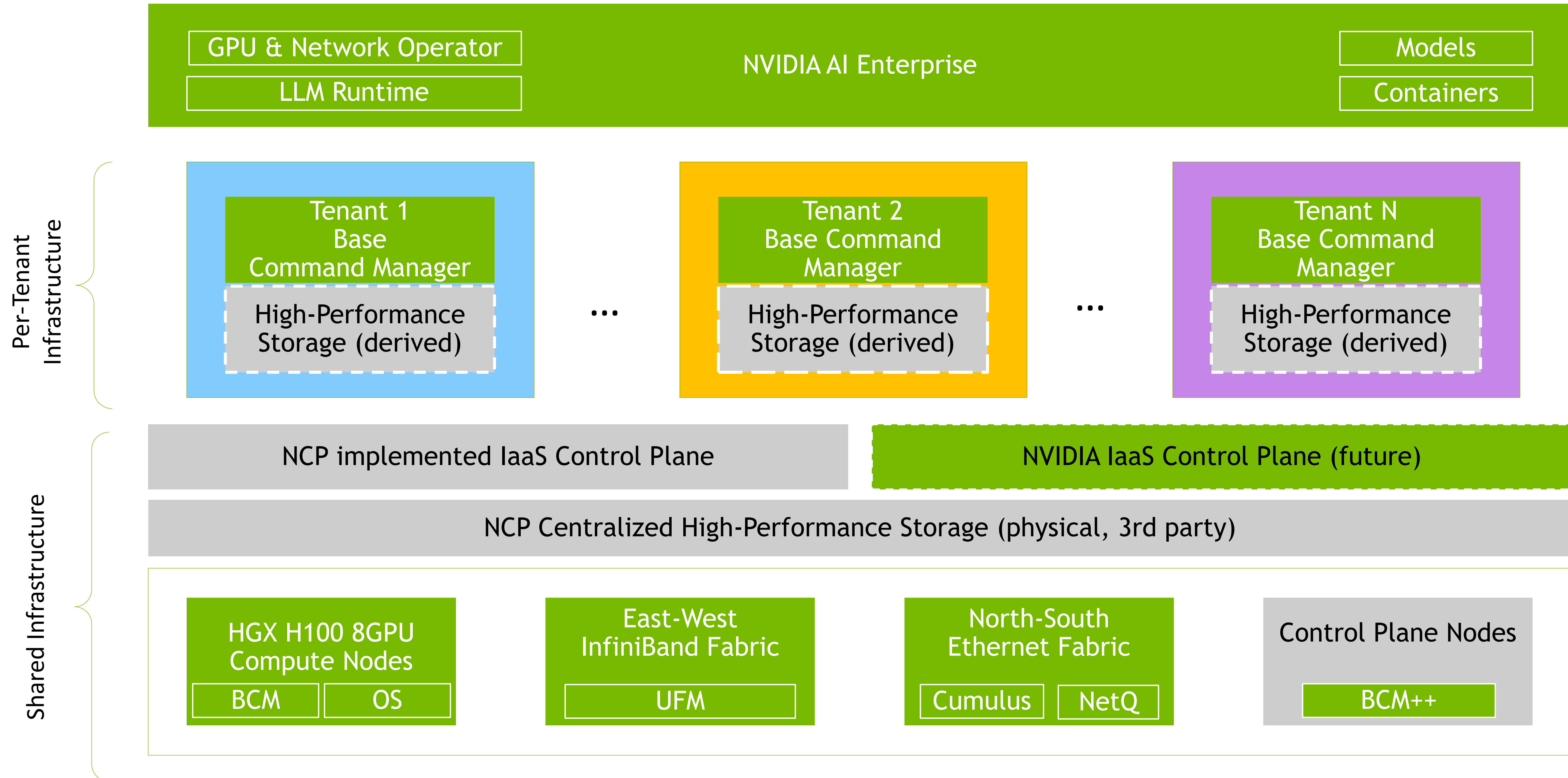
## NVIDIA DGX SuperPOD with H100
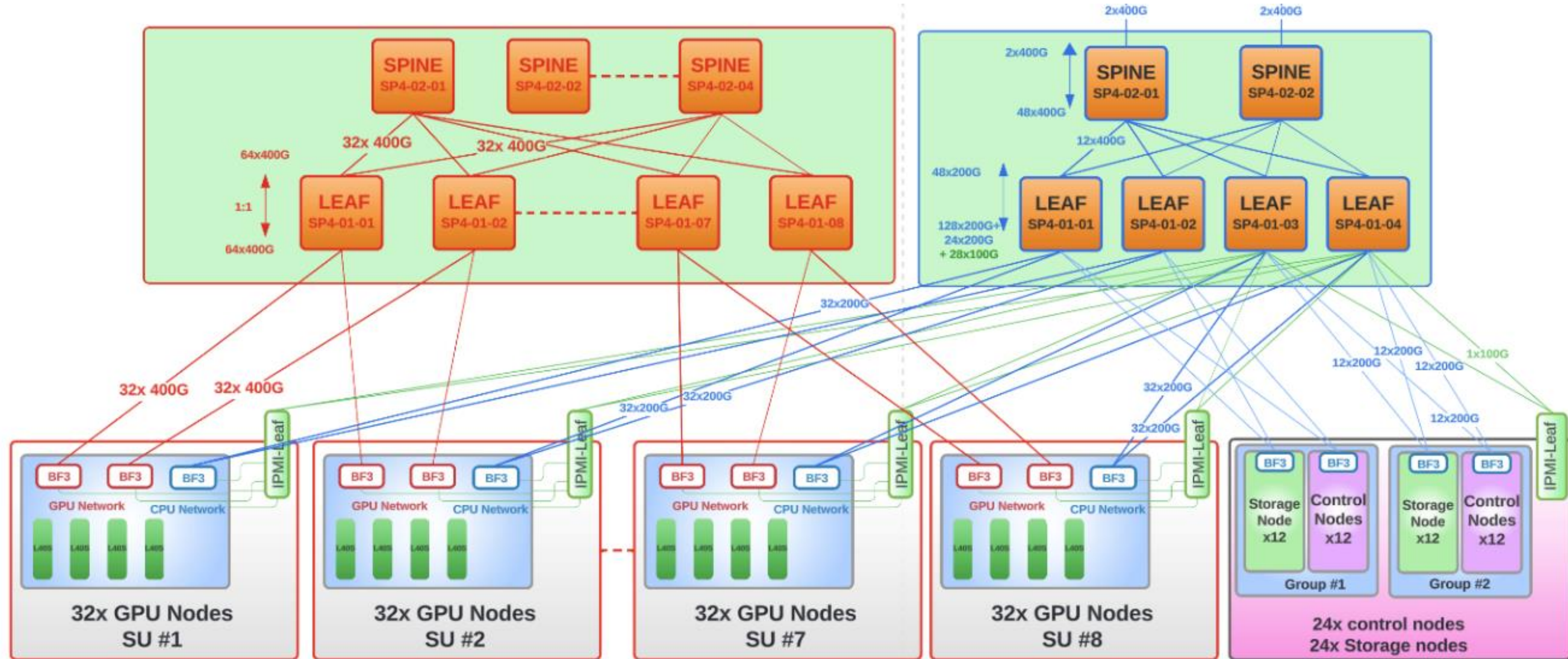


Scalable Unit (SU)

Utility Racks

# Taking the "SuperPOD" to Cloud

## NCP Reference Architecture

# The Versatile Answer for Inference and Graphics

NVIDIA L40s Reference Architecture

# New Option for Fine-Tuning and Inference
NVIDIA Reference Architecture with GraceHopper Superchip GH200