| | |
|---|---|
| **Kick-off Meeting (Nov 16) Virtual** | • 02:00PM - 02:05PM: Welcome and event overview (Jay, CK)<br>• 02:05PM – 02:10PM: NCHC opening (王順泰組長@NCHC )<br>• 02:10PM – 02:15PM: Hackathon team opening (Bharat)<br>• 02:15PM - 03:00PM: Round table self-introduction (Team & Mentor).<br>• 3 mins for each team lead<br>• 1 mins for two mentors per team<br>• 03:00PM - 03:05PM: 5 mins break<br>• 03:05PM - 03:15PM: Introduction to computing resources (Kuan-Ting)<br>• 03:15PM - 04:00PM: Introduction to Nsight Analysis Tools (Leo Chen)<br>• 04:00PM - 04:30PM: breakout rooms (Team & Mentor) |
| **Day 1 (November 23) Virtual** | • 02:00PM - 03:00PM: Scrum #1 (5 mins presentation per team) |
| **Day 2 (November 30) Virtual** | • 02:00PM - 03:00PM: Scrum #2 (5 mins presentation per team) |
| **Day 3 (Dec 07) In-Person** | • 10:00 AM - 10:30 PM: Welcome and event description<br>• 10:30 AM - 12:00 PM: Final presentation (12 mins presentation +3 minutes QA per team)<br>• 12:00 PM - 01:30 PM: Lunch time<br>• 01:30 PM - 03:00 PM: Final presentation (12 mins presentation +3 minutes QA per team)<br>• 03:00 PM - 04:00 PM: Wrap-up session |

# Team Name

NVIDIA.

# NVIDIA GH200 Grace Hopper Superchip

Built for the New Era of AI Supercomputing

**CPU to GPU Bandwidth**
## 900GB/s
NVLink-C2C

**GPU Memory Bandwidth**
## 4.9TB/s
HBM3e

**Energy Efficiency**
## 52X
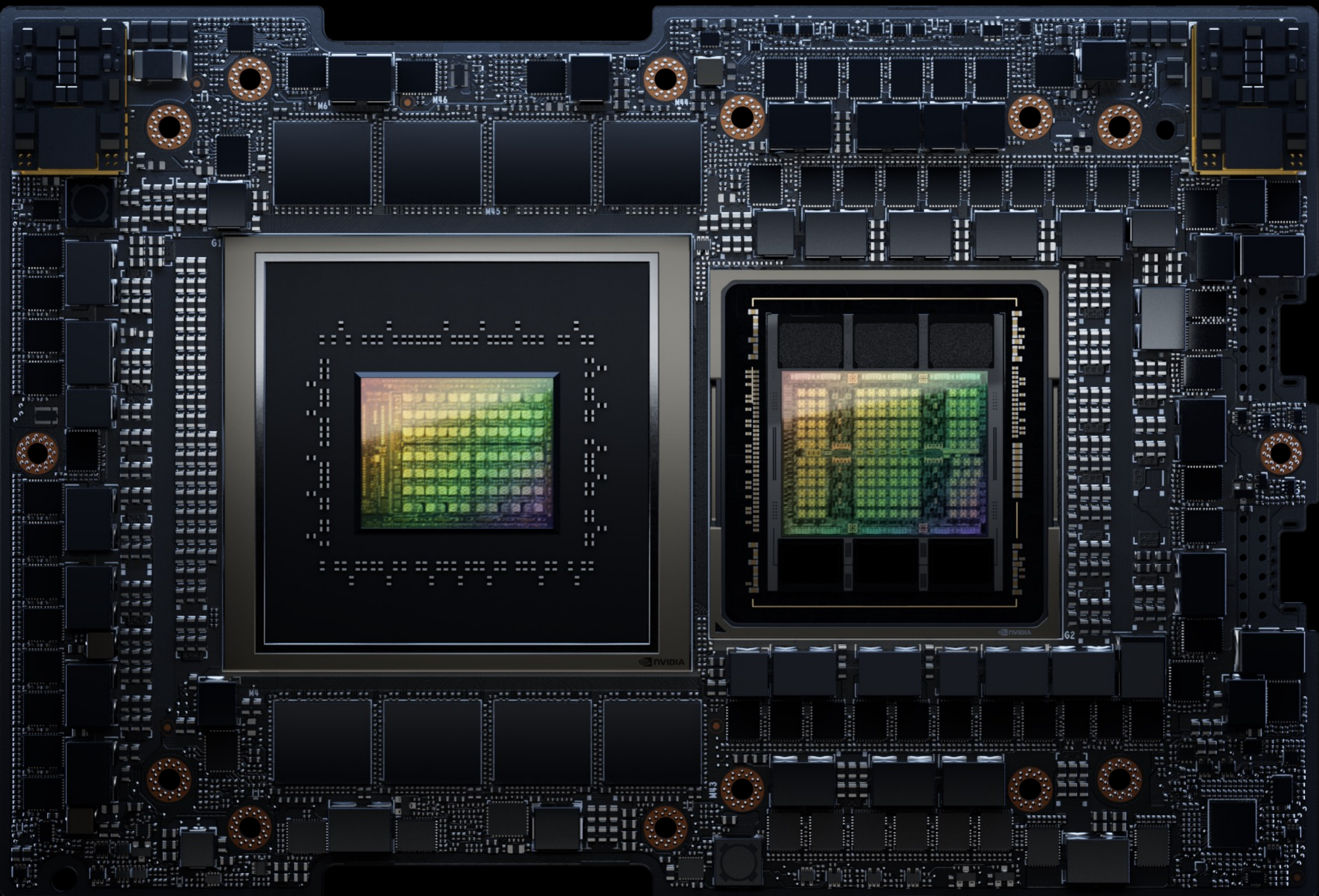MILC Efficiency vs 2S x86 CPUs

**QFT Quantum Simulation**
## 90X
Performance vs 2S x86 CPUs

**Llama 2 70B Inference**
## 100X
Performance vs 2S x86 CPUs

624GB High-Speed Memory | 4 PF AI Perf | 72 Arm Cores

**NVIDIA**

# Now in Full Production – NVIDIA GH200 Grace Hopper Superchip

## Built for the New Era of AI Supercomputing

**Most versatile compute**
Best performance across CPU, GPU or memory
intensive applications
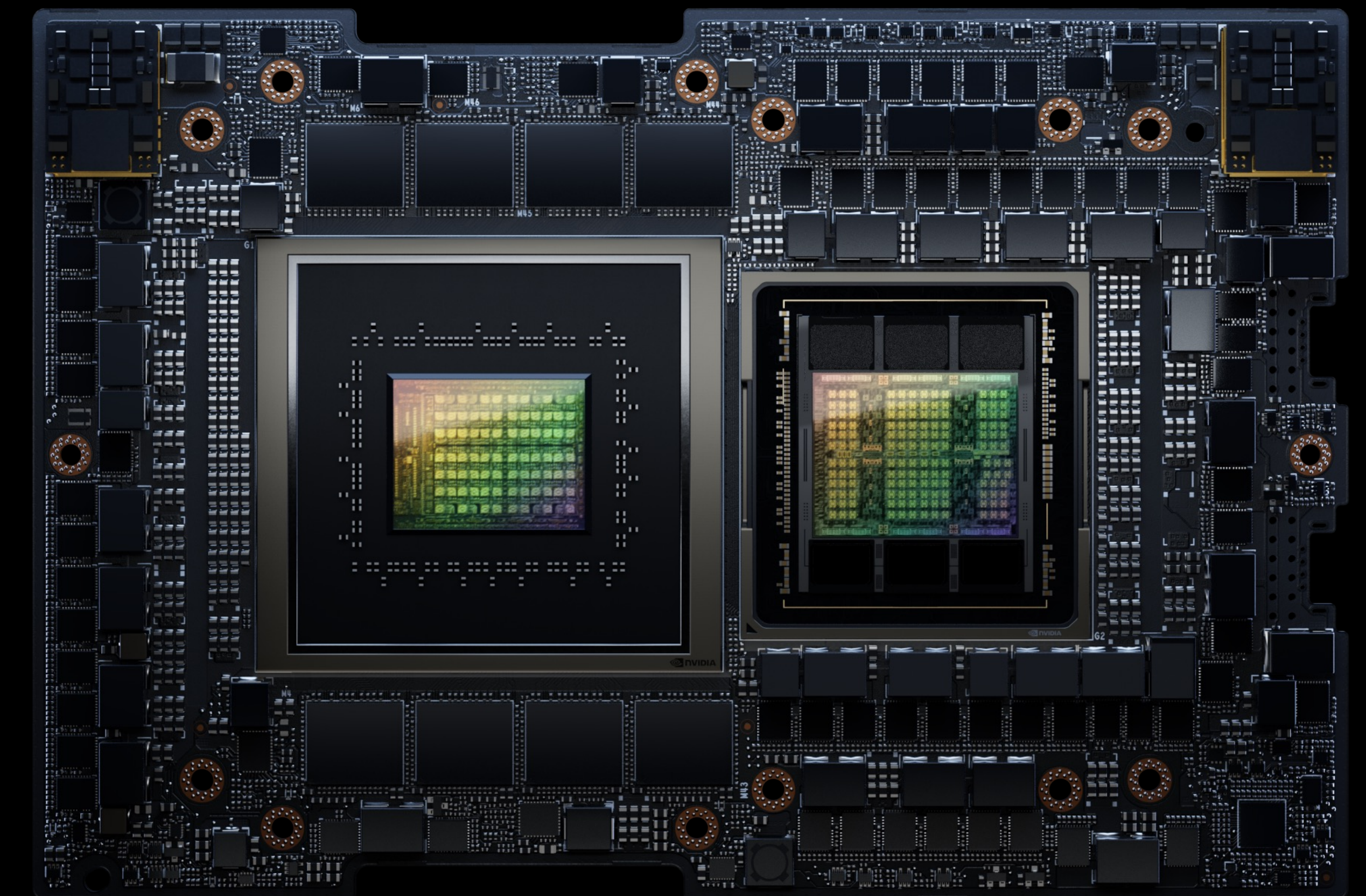
**Easy to deploy and scale out**
1 CPU:1 GPU node simple to manage and schedule
for for HPC, enterprise, and cloud

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

**Now Available in Launchpad**
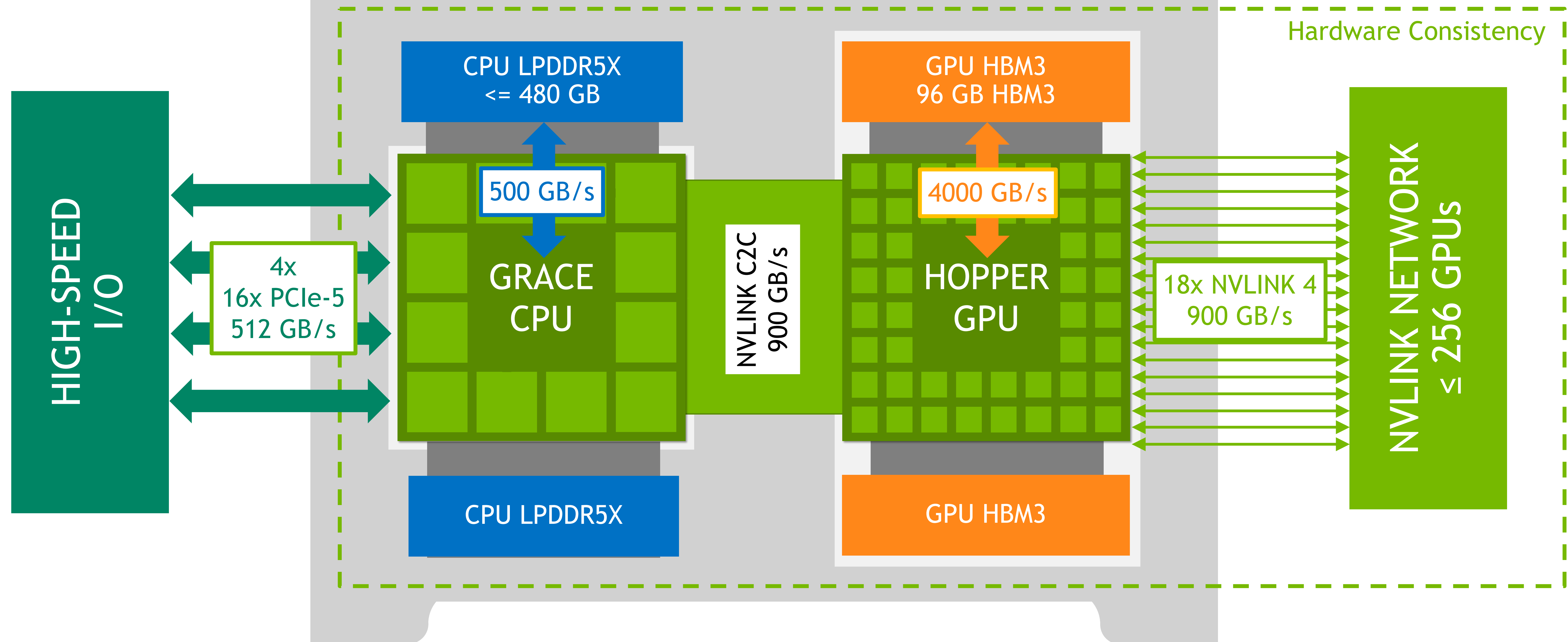Available Starting Early December; Sign-up Today

https://www.nvidia.com/en-us/launchpad/grace-hopper/

Shipping Globally with Early Access from CSPs

ASRock Rack   ASUS   DELL Technologies   CoreWeave   EVIDEN an atos business   Fii Foxconn Industrial Internet

GIGABYTE   Hewlett Packard Enterprise   Inventec   Lambda   Lenovo   NVIDIA   ORACLE CLOUD Infrastructure

PEGATRON   QCT   SUPERMICRO   VULTR   Wistron   wiwynn   zt Systems

# Grace Hopper Superchip
## GPU can access CPU memory at CPU memory speeds

**NVIDIA Grace Hopper Superchip**

Hardware Consistency

HIGH-SPEED I/O

CPU LPDDR5X <= 480 GB

GPU HBM3 96 GB HBM3

500 GB/s

4000 GB/s

4x 16x PCIe-5 512 GB/s

GRACE CPU

NVLINK C2C 900 GB/s

HOPPER GPU

18x NVLINK 4 900 GB/s

NVLINK NETWORK ≤ 256 GPUs

CPU LPDDR5X

GPU HBM3

https://resources.nvidia.com/en-us-grace-cpu/nvidia-grace-hopper

NVIDIA

# GH200 Grace Hopper HPC Platform

## Unified Memory and Cache Coherence for Next Gen HPC Performance

### Partially GPU Accelerated Apps

Big performance gains with no code changes
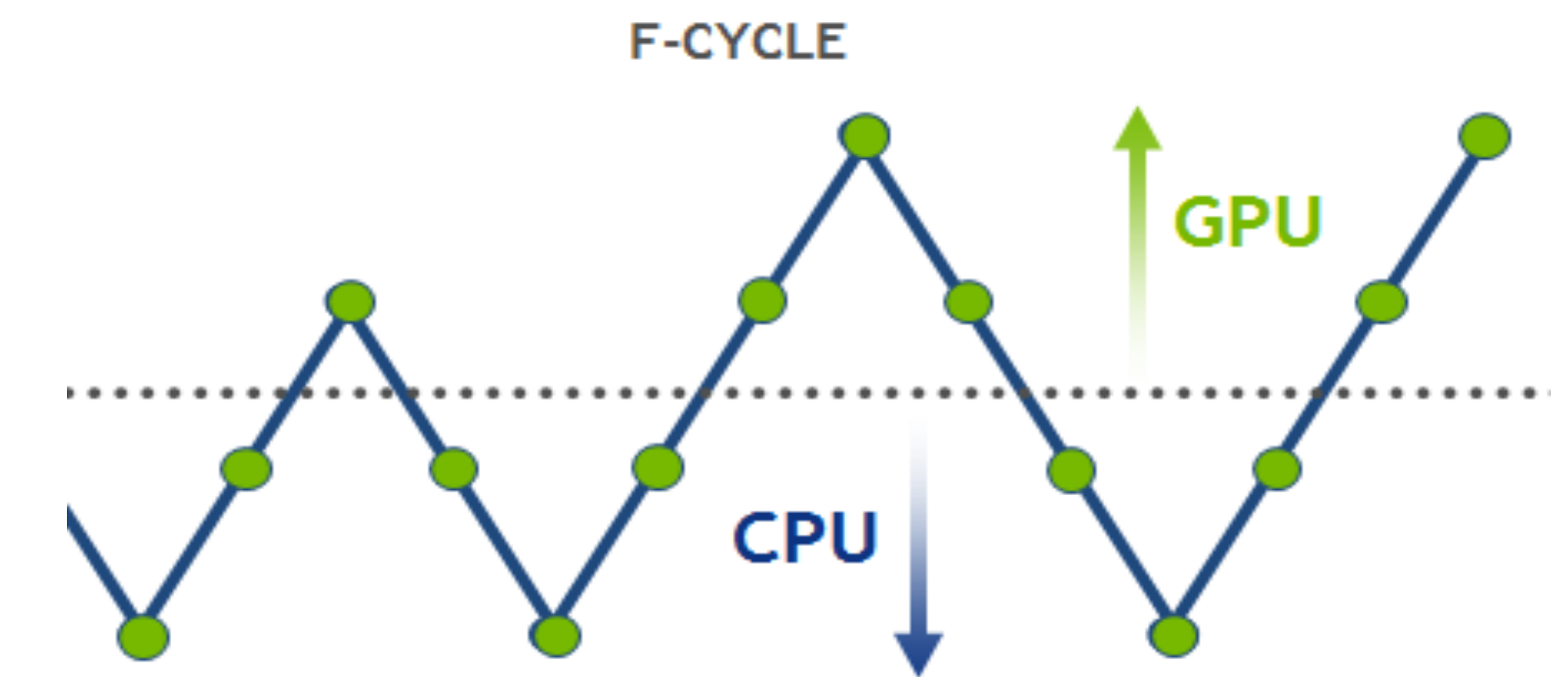


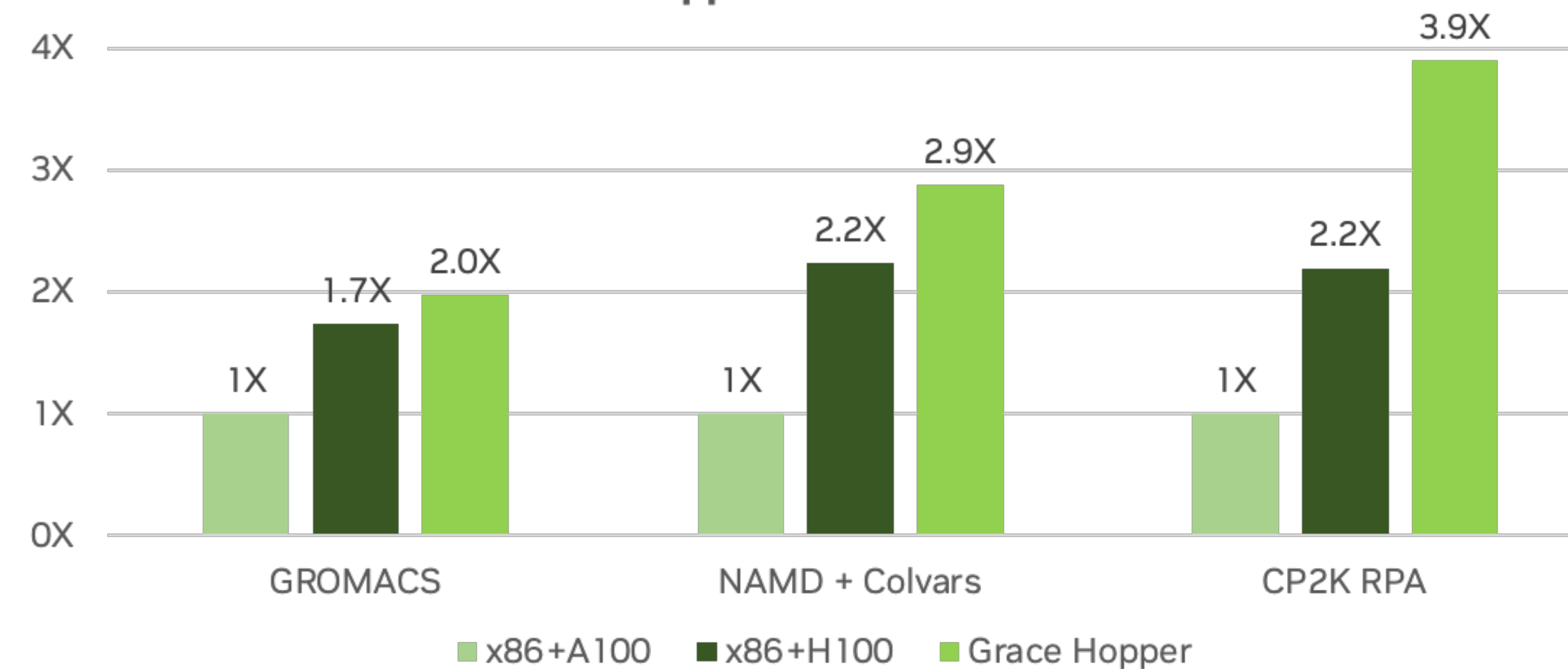### No More PCIe Bottleneck

NVLink-C2C is 7X PCIe BW



### CPU & GPU Cache Coherence

Incremental code changes yield big gains



**Grace Hopper HPC Performance**



**Fast Access Memory**

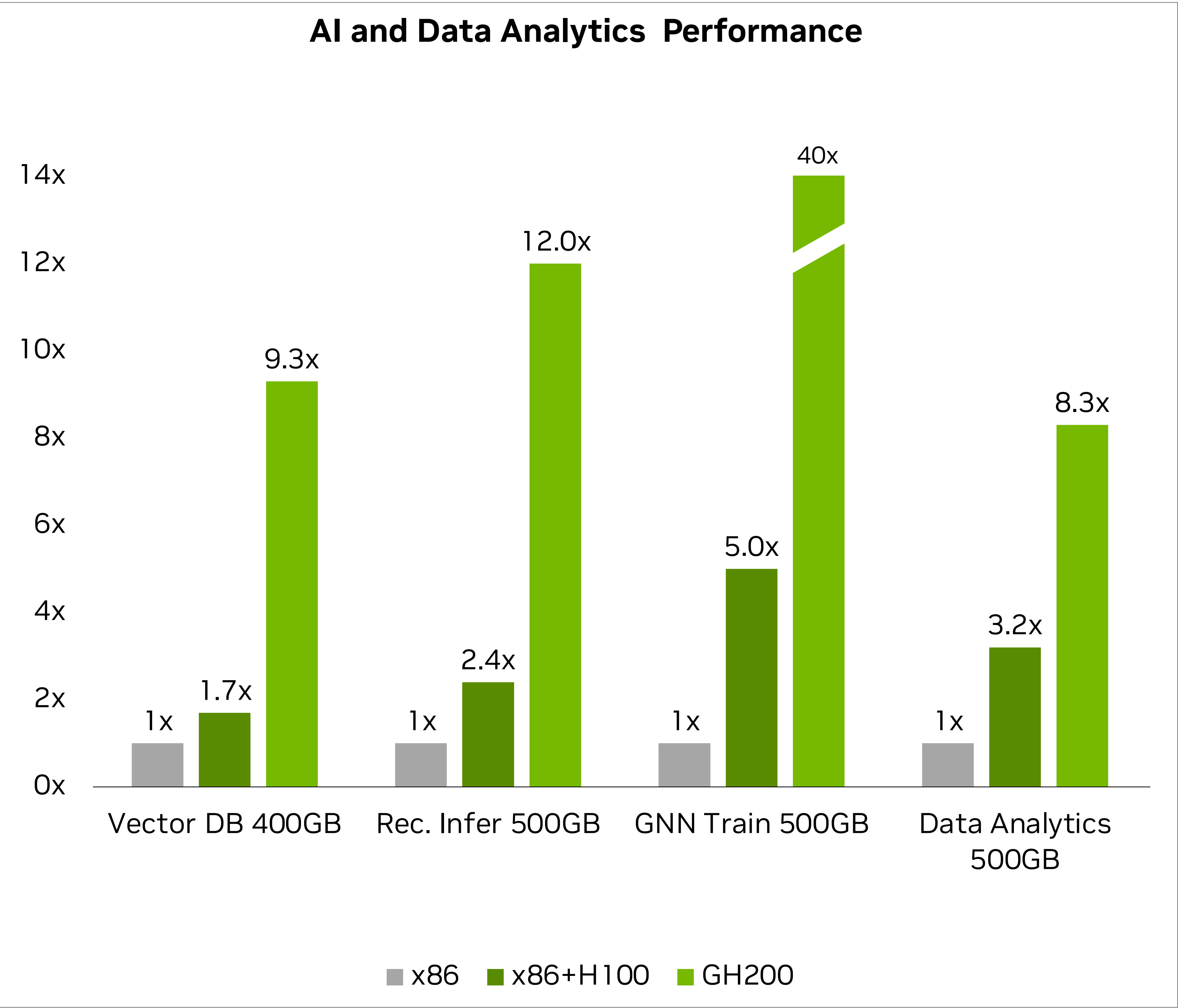## 624GB

**Memory Bandwidth**

## 5TB/s

NVIDIA

# GH200 Grace Hopper AI Inference Platform

## Versatile Scale Out with Unmatched Performance

## Memory Intensive

### AI and Data Analytics Performance



Legend: ■ x86 ■ x86+H100 ■ GH200

- Vector DB 400GB: 1x, 1.7x, 9.3x
- Rec. Infer 500GB: 1x, 2.4x, 12.0x
- GNN Train 500GB: 1x, 5.0x, 40x
- Data Analytics 500GB: 1x, 3.2x, 8.3x

## GPU Intensive

### LLM Infer Performance



Legend: ■ x86 ■ x86+H100 ■ GH200

- LLAMA 65B: 1x, 122x, 284x
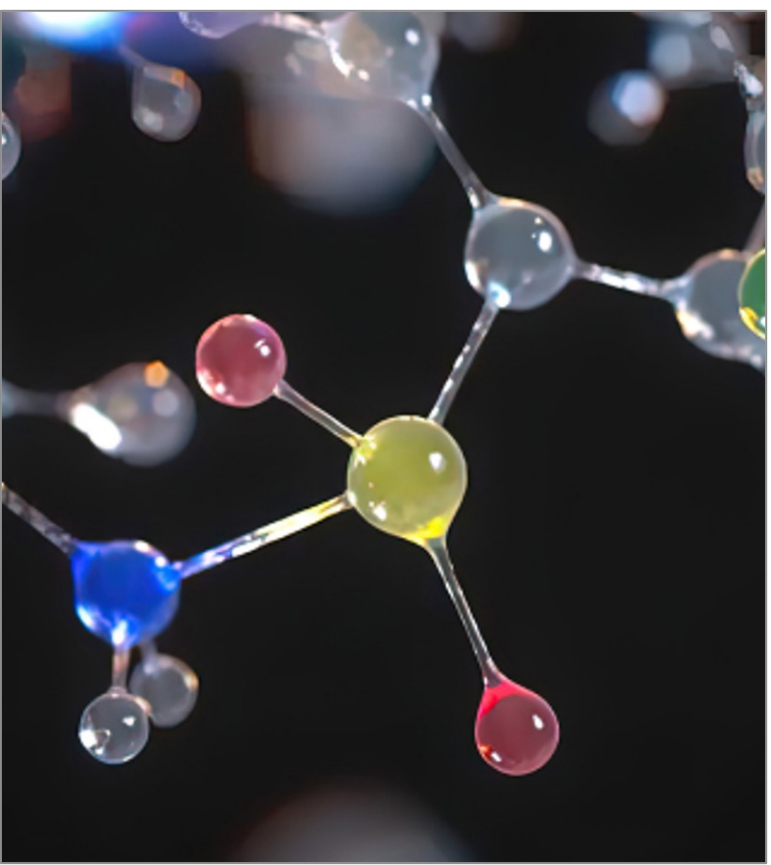
## Use Cases

**LLM**
Conversational AI
Domain Knowledge

**Recommender Systems**
eCommerce
Personalized Content

**Vector Database**
Fraud Detection
Drug Discovery

**GNN**
Computer Vison
Recommenders

NVIDIA.

# Simplifying GPU Programming for HPC with NVIDIA Grace Hopper Superchip

624GB High-Speed Memory | 4 PF AI Perf | 72 Arm Cores