**NVIDIA**

# GTC 2025 直送！AIコンピューティング最新情報
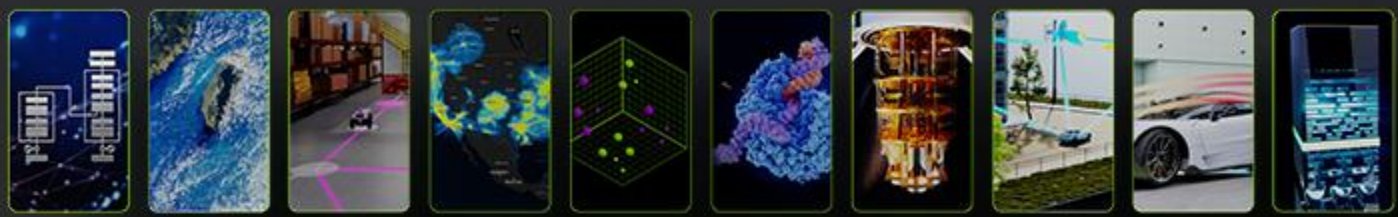
2025/4/9
佐々木 邦暢 (@_ksasaki)
エヌビディア合同会社

# NVIDIA AI & HPC プラットフォーム

NIM
CUDA-Accelerated
Agentic AI Libraries

Omniverse
CUDA-Accelerated
Physical AI Libraries

CUDA-X Libraries

CUDA • DOCA • NCCL
Cluster-Scale Software
System Software
Chip Software

Accelerated
Software Stack

GB200 NVL72 SuperPOD

Grace Blackwell
MGX Node

NVLink Switch

Quantum Switch

Spectrum-X Switch

Chips Purpose-Built for AI Supercomputing
GPU | CPU | DPU | NIC | NVLink Switch | IB Switch | Enet Switch

# Blackwell があらゆる場所に
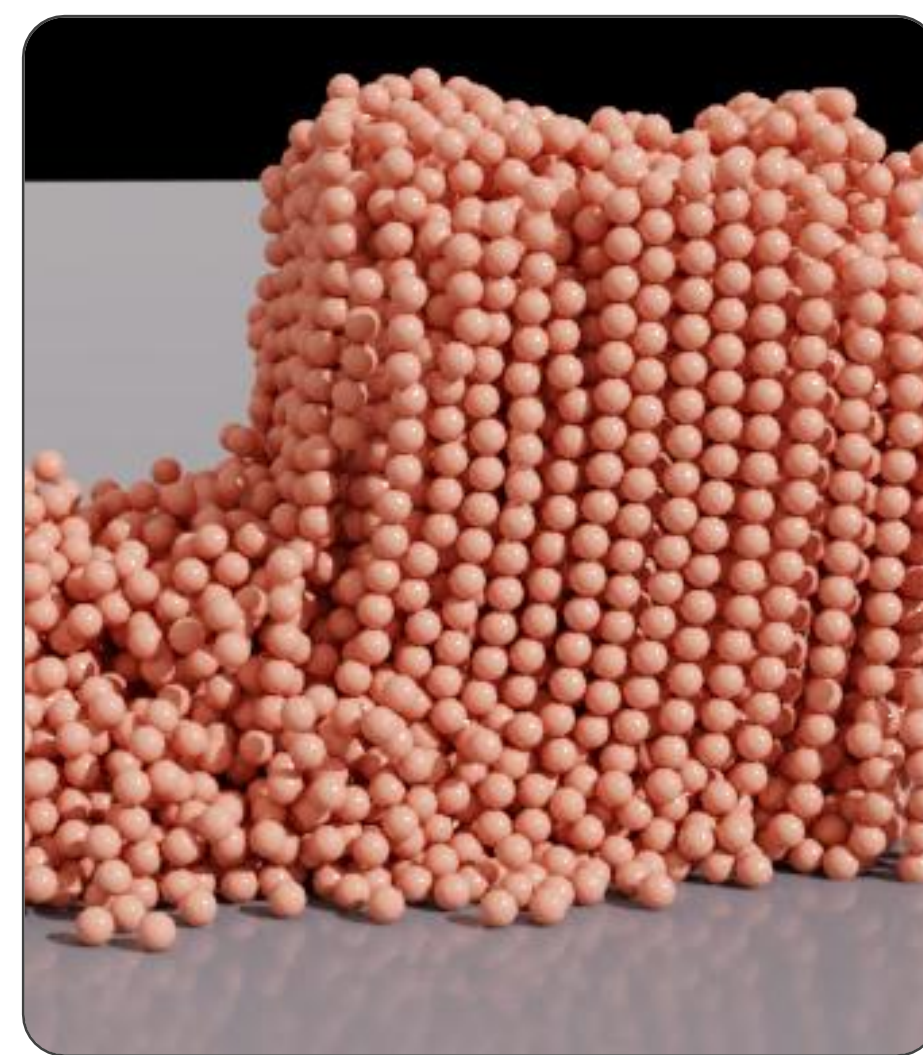


Fastest Ramping Product in NVIDIA History
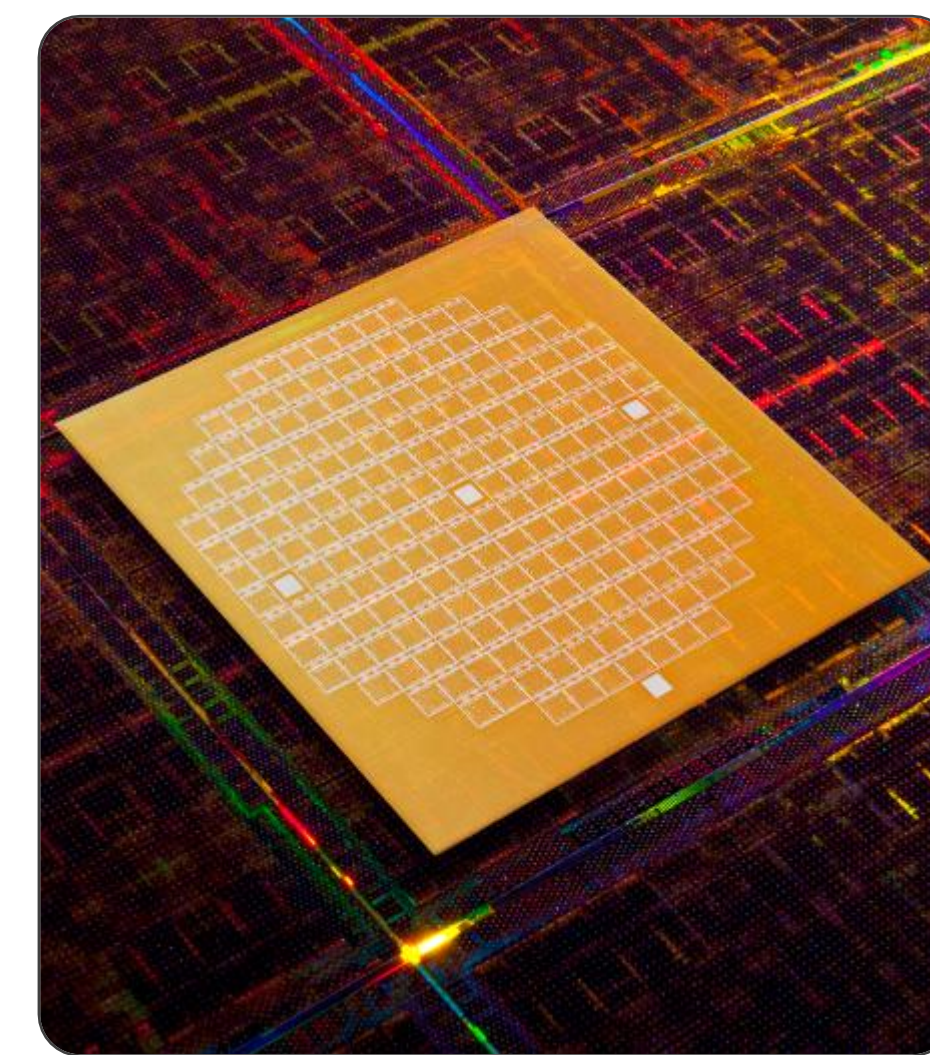
# CUDA-X が多様なアプリケーションを高速化
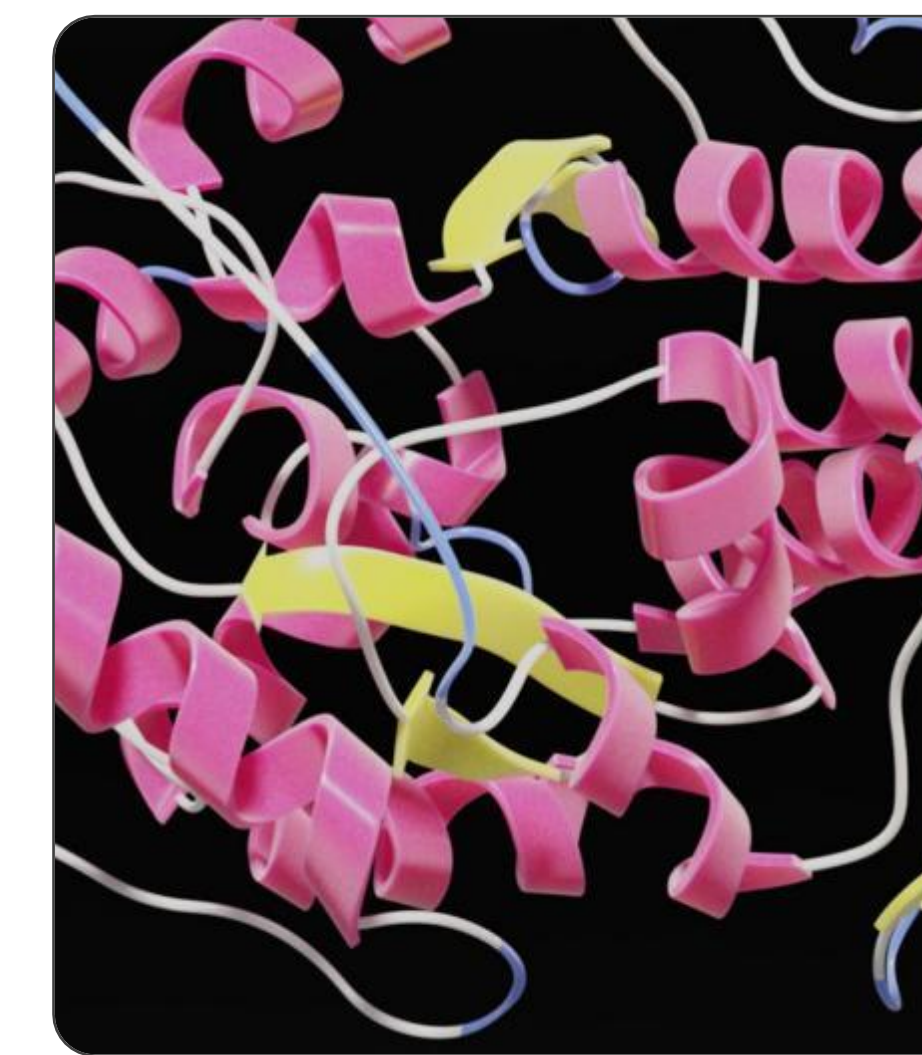


**cuDSS**
CAE



**PhysicsNeMo**
AI Physics



**Warp**
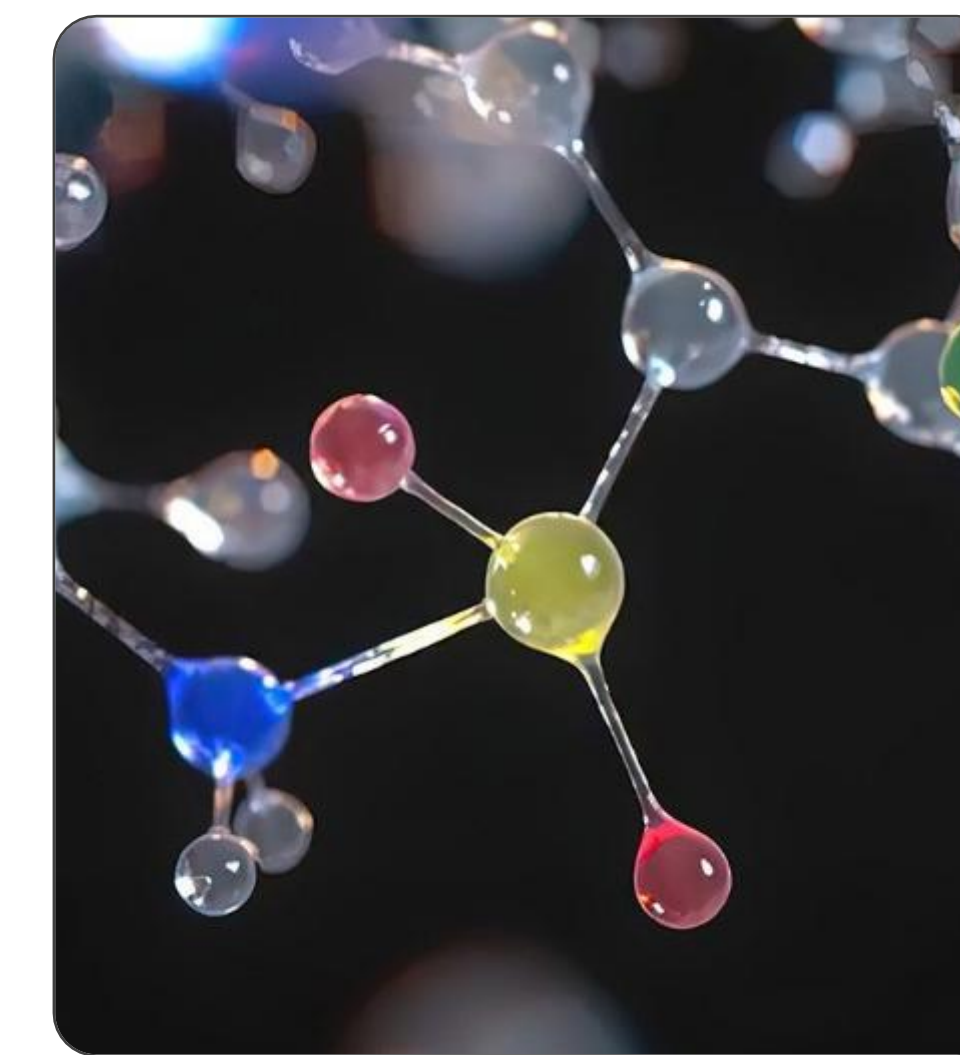Physical Simulation



**cuLitho**
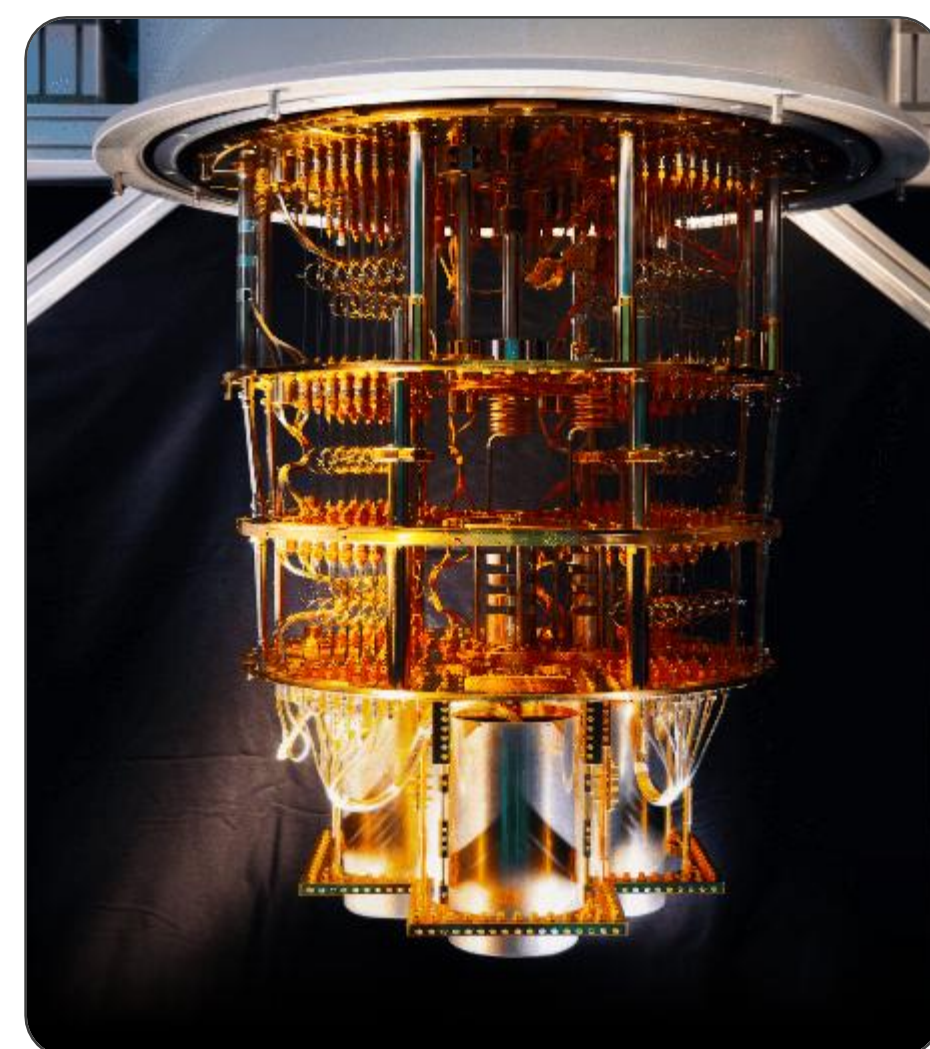Computational
Lithography



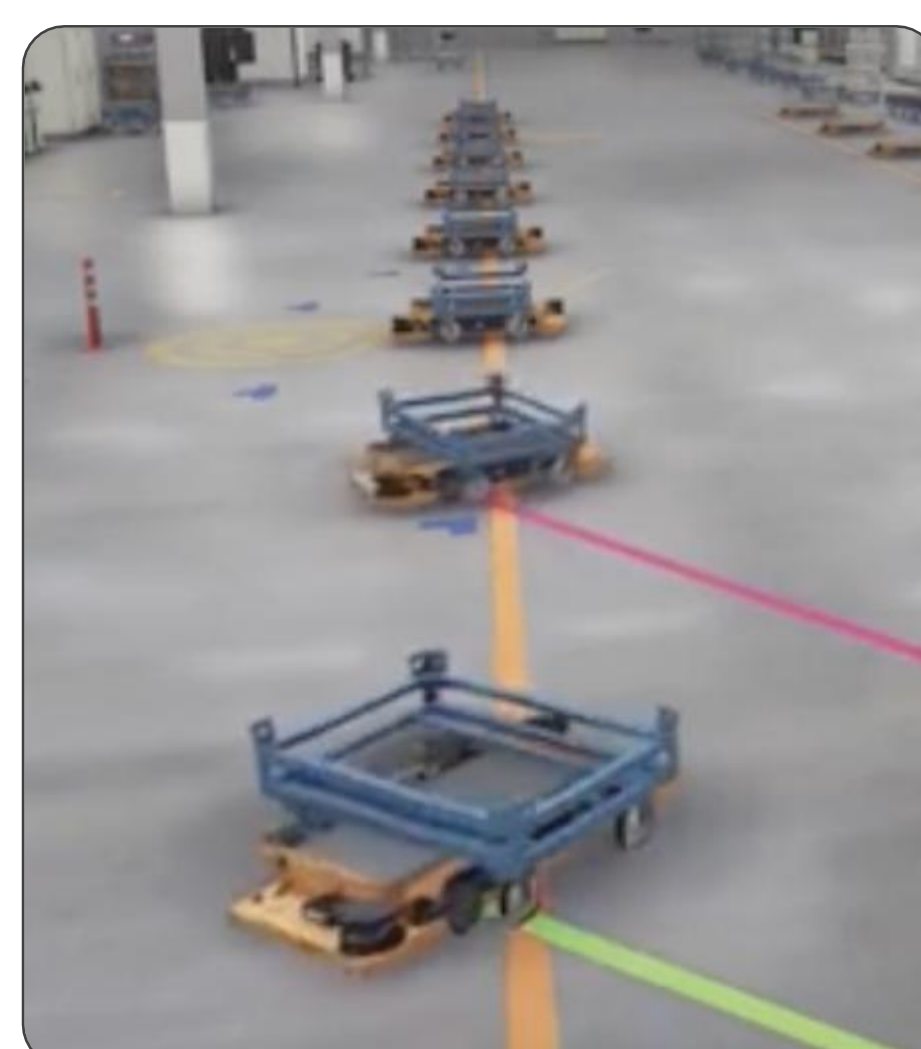**ALCHEMI**
AI Materials Science



**cuEquivariance**
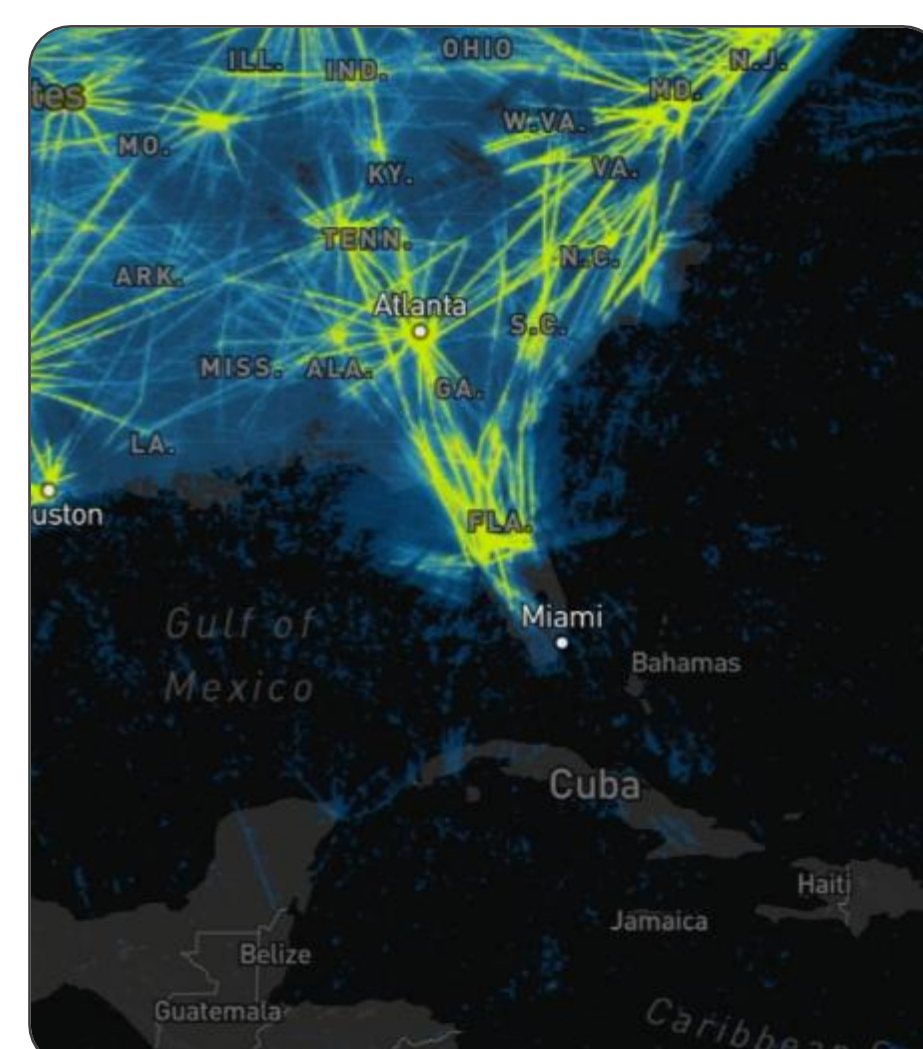Drug & Materials
Discovery



**Parabricks**
Gene Sequencing



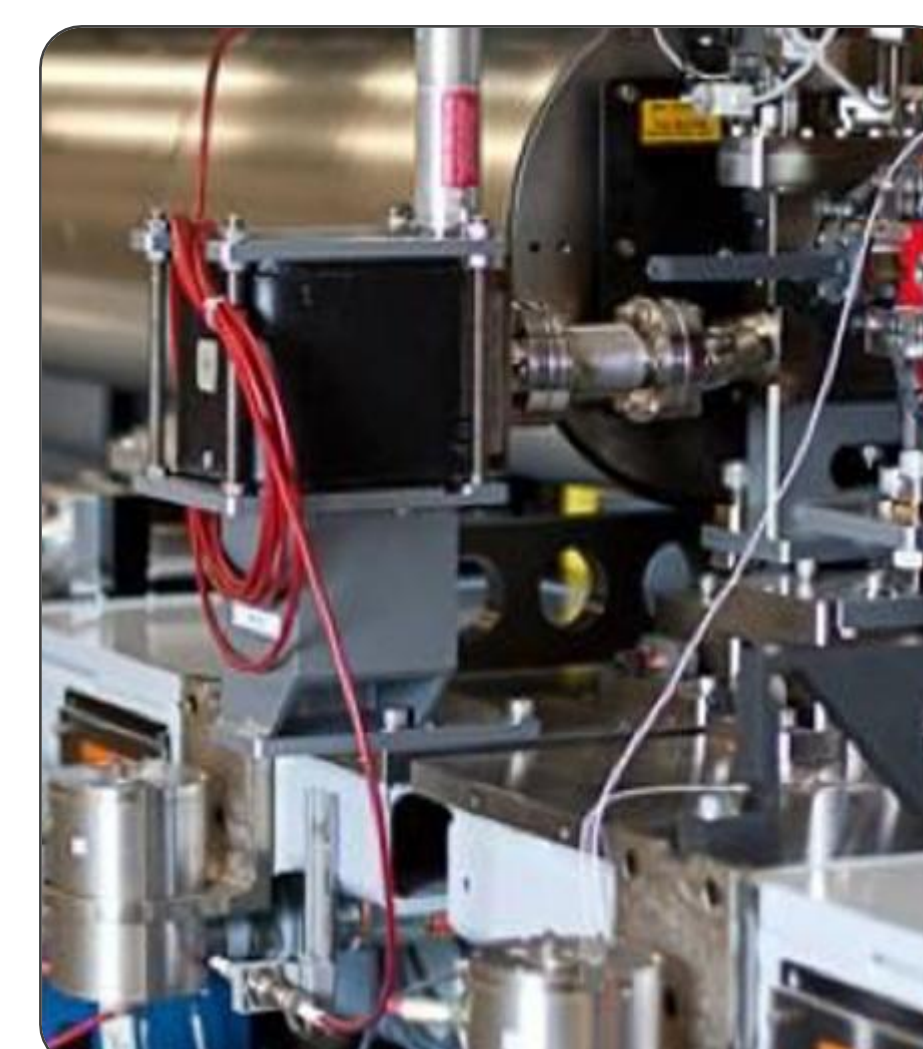**CUDA-Q**
Quantum Computing



**cuOpt**
Decision Optimization
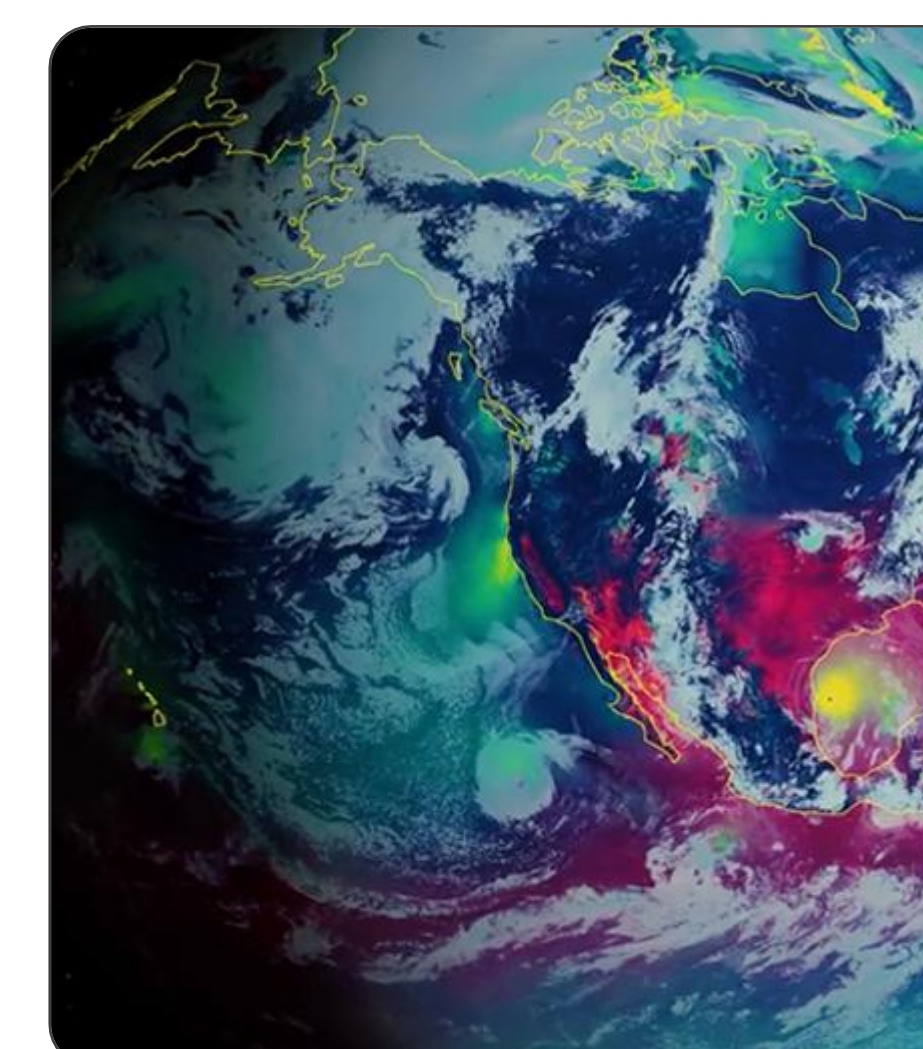


**cuDF**
Data Processing



**cuPyNumeric**
Numerical Computing



**Holoscan**
Edge HPC



**Earth-2**
Weather Analytics

# Blackwell で加速する計算工学のエコシステム



50X

40X

30X

20X

20X

SPICE　　Litho　　DEM　　FEA　　CFD

Boom Supersonic, Rescale, and NASA Fun3D

**SIEMENS**
Maserati with Siemens Simcenter STAR-CCM+

**cādence**®
NVIDIA data-center digital twins with Cadence Fidelity

**Ansys**
Volvo accelerating with Ansys Fluent

**ALTAIR**　**Ansys**　**AUTODESK**　**BeyondMath** POSSIBLE IS POSSIBLE　**cādence**®　**COMSOL**　**engys**® Open-source CFD for Enterprise　**HEXAGON**　**luminary**

**M-Star**　**NAVASTO** An AUTODESK Company　**NEURAL CONCEPT**　**nTop**　**PHYSICSX**　**rescale**　**SIEMENS**　**SIMSCALE**　**VORTEX**　**VOLCANO** Platforms

NVIDIA

# AI スケーリング則による演算需要の指数関数的な高まり

## "長考" する AI が推論をスケーリング

# リーズニング時代の **Post-Training** スケーリング

Teaching models to think requires a lot of compute

Update the Model

Training Question

**Training Model**
1 Trillion Parameters

100 Answer Candidates
10,000 Tokens Each

Verify Score

**100**
Topic Areas

**1 Million**
Questions Per topic

100 Topics x 1 Million Questions x 100 Answers x 10,000 Tokens

**100 Trillion Tokens**

7X More Than What Llama 3 Was Pre-Trained On

# リーズニングによるトークン数が急激に増大

Cost increasing with compute – requires full stack innovation to lower cost/token

## KNOWLEDGE
Pre-trained models simulate internet knowledge
100 tokens | $0.50 / M Token

NVIDIA とは？

NVIDIAは、グラフィックス・プロセッシング・ユニット（GPU）、高性能コンピューティング・ハードウェア、人工知能（AI）技術の設計・製造を専門とする多国籍テクノロジー企業です。1993年にジェンスン・フアン、クリス・マラコウスキー、カーティス・プライムによって設立されました。

NVIDIAはカリフォルニア州サンタクララに本社を置き、コンピューターグラフィックス、ゲーム、AIの分野におけるリーディングカンパニーの一つとなっています。

## REASONING
Reasoning models think through problems
~10K tokens → 100X cost

家族8人で円卓で夕食をとるのですが、どこに座ればいいのか迷っています。義理の両親は私の両親と仲が悪く、妹は弟の食事の手伝いをするために弟の隣に座らなければなりません。妻は母の隣に座りたがりません。それぞれどこに座ればいいのでしょうか？

8人家族が円卓に円満に座れるように、あらゆる制約を考慮しながら最適な配置は以下の通りです。

あなた 妻 お父さん お母さん 妹 弟 義母 義父

主な制約：
義理の両親と両親：妹、弟、そしてあなたで区切る。
妹と弟：助け合うために隣同士に座る。
妻とお母さん：隣同士に座らない（あなたとお父さんで区切る）。

## ACCELERATING AI FACTORY VALUE
>$1/M Token  |  >1T Param  |  >300 TPS

Need for
increasing tokens

Increasing revenue
customer experience

**Decrease cost**

# NVIDIA Dynamo

# NVIDIA Dynamo

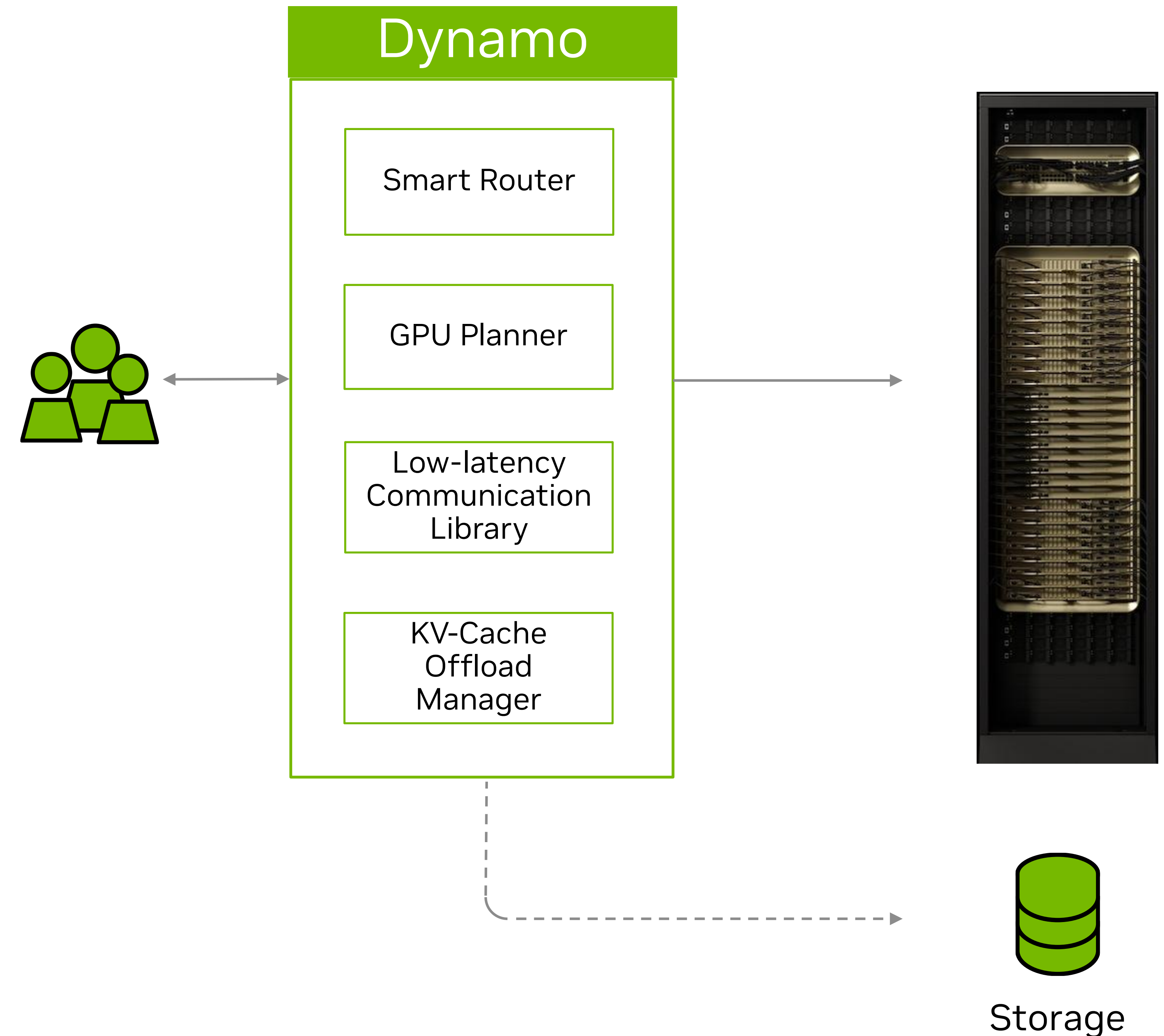## リーズニング AI を支える推論ソフトウェア

### 分散型および非集約型の推論サービス

**30X**

AI Factory
Throughput
& Revenue

Deepseek R1
Based models

**2.5X**
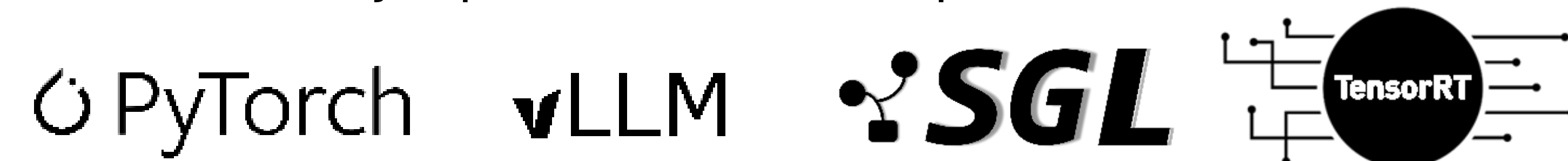
AI Factory
Throughput
& Revenue

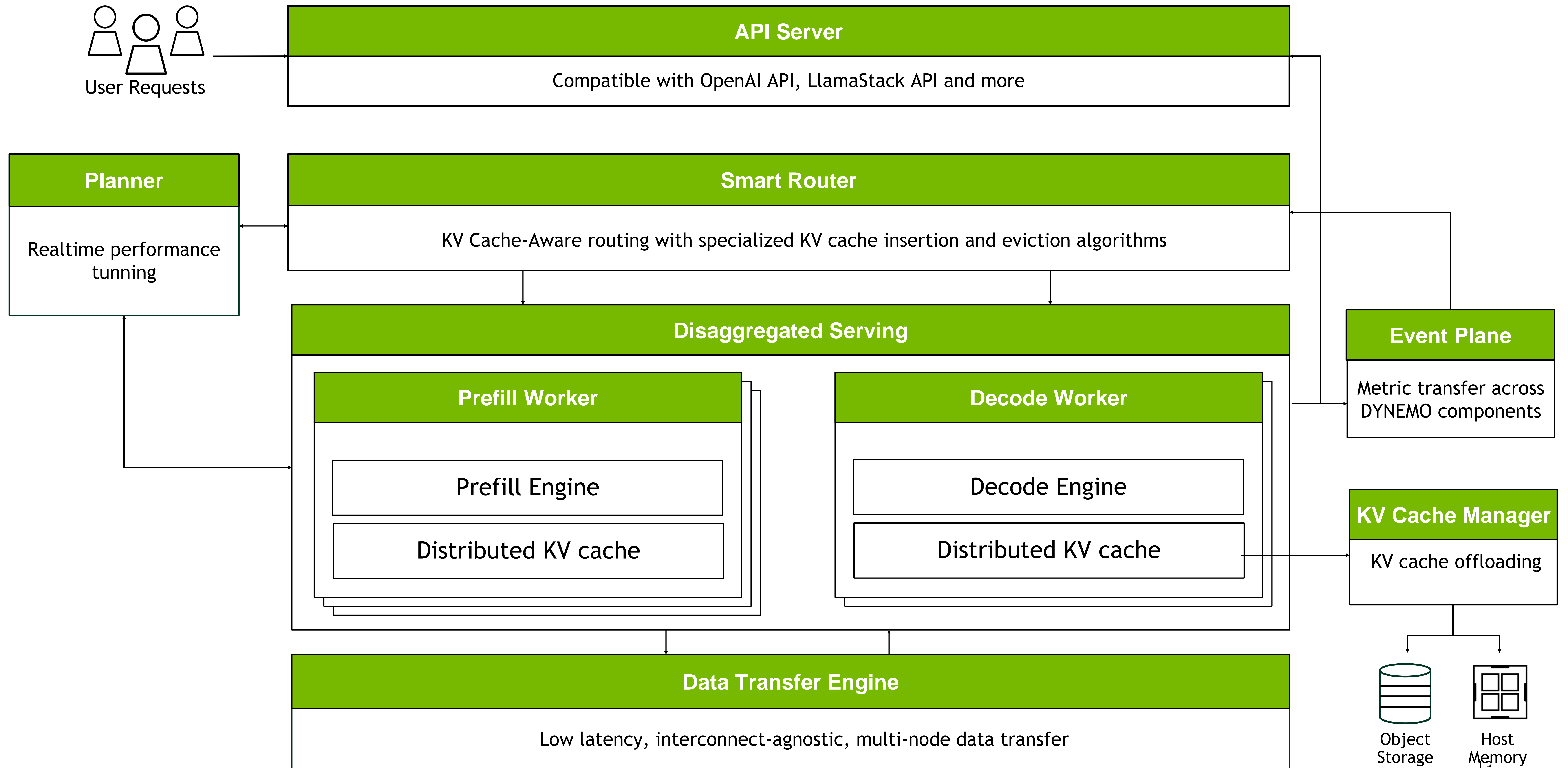Llama
Based Models

**1000+**

GPU Scale for
a single query

Fully Open Source and Open Backend

PyTorch    vLLM    SGL    TensorRT

**Dynamo**

Smart Router

GPU Planner

Low-latency
Communication
Library

KV-Cache
Offload
Manager

Storage

NVIDIA

# Architecture and Components

**User Requests**

**API Server**

Compatible with OpenAI API, LlamaStack API and more

**Planner**

Realtime performance tunning

**Smart Router**

KV Cache-Aware routing with specialized KV cache insertion and eviction algorithms

**Disaggregated Serving**

**Prefill Worker**

Prefill Engine

Distributed KV cache

**Decode Worker**

Decode Engine

Distributed KV cache

**Event Plane**

Metric transfer across DYNEMO components

**KV Cache Manager**

KV cache offloading

**Data Transfer Engine**

Low latency, interconnect-agnostic, multi-node data transfer
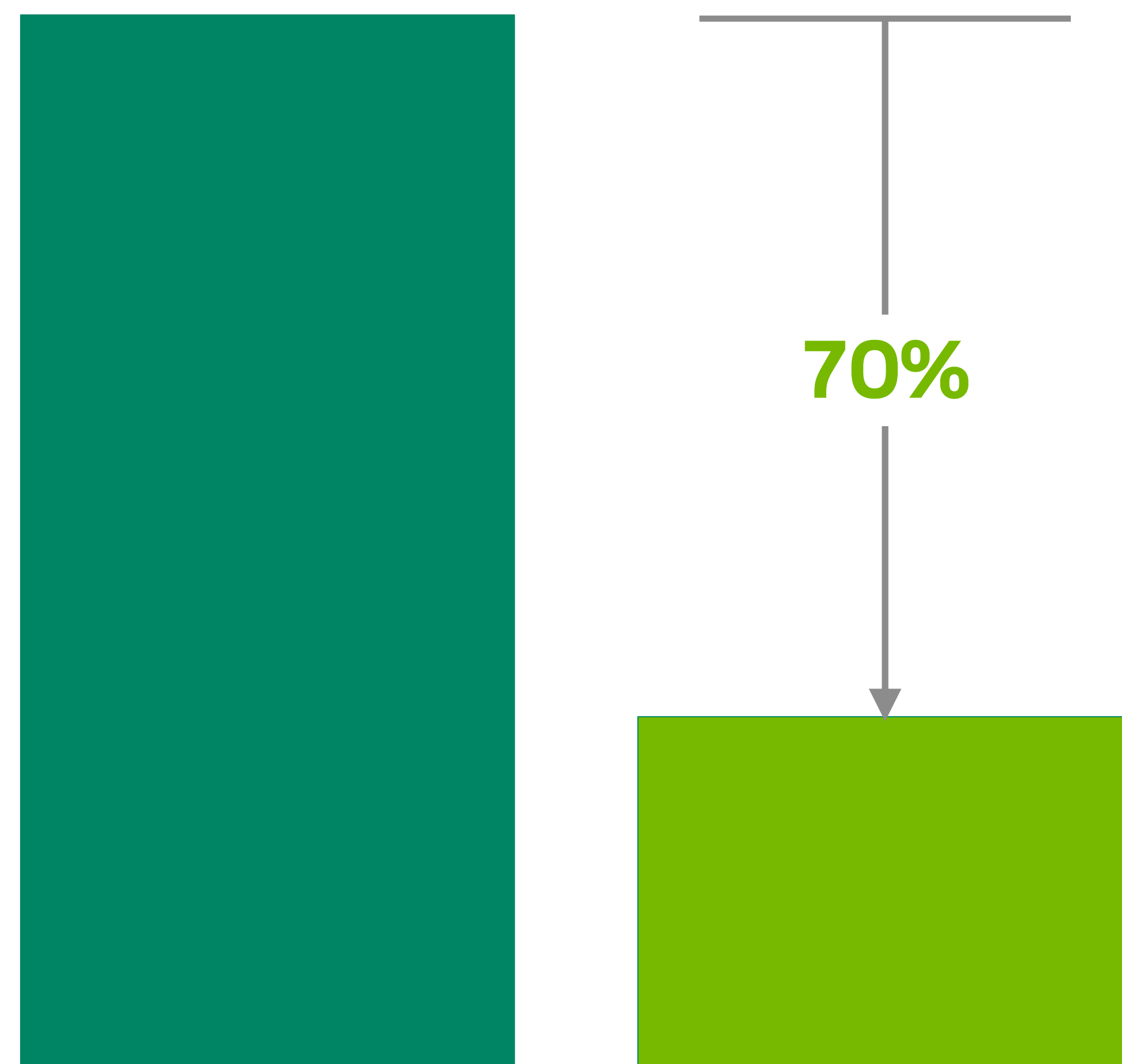
Object Storage

Host Memory

13

# NVIDIA Dynamo: Smart Router

## Reducing costly re-computation of KV cache

**DeepSeek-R1 Distill Llama 70B | NVIDIA HGX-H100**
(Lower is Better)



**Time to First Token**

70%

**Avg. Request Latency**

50%

■ NVIDIA Dynamo w/ Random Routing    ■ NVIDIA Dynamo w/ Smart Router

■ NVIDIA Dynamo w/ Random Routing    ■ NVIDIA Dynamo w/ Smart Router
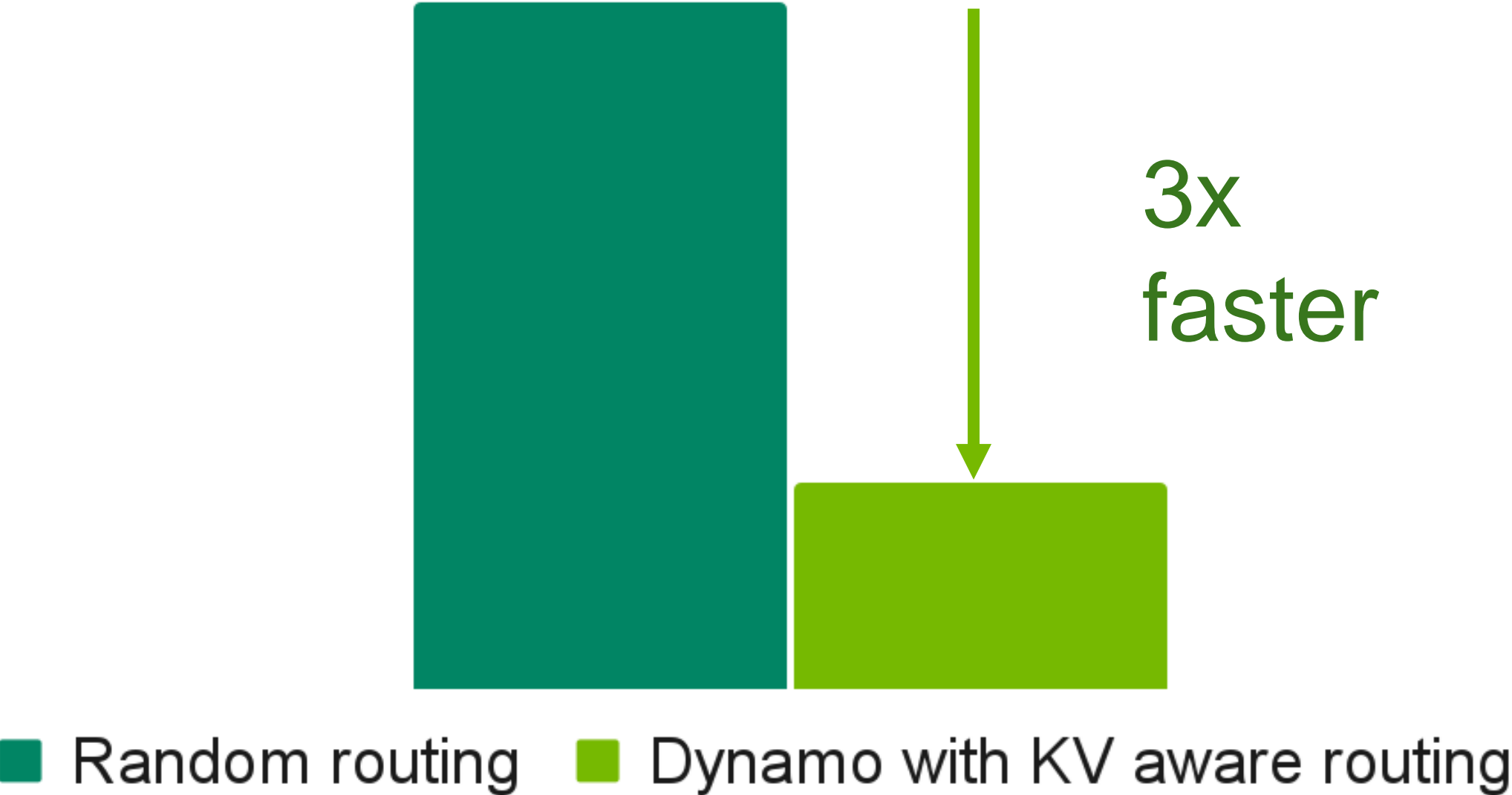
# KV Cache Aware Routing

## Significant boost in TTFT and End to End Latency with real data (100K requests with R1)



**Time To First Token**

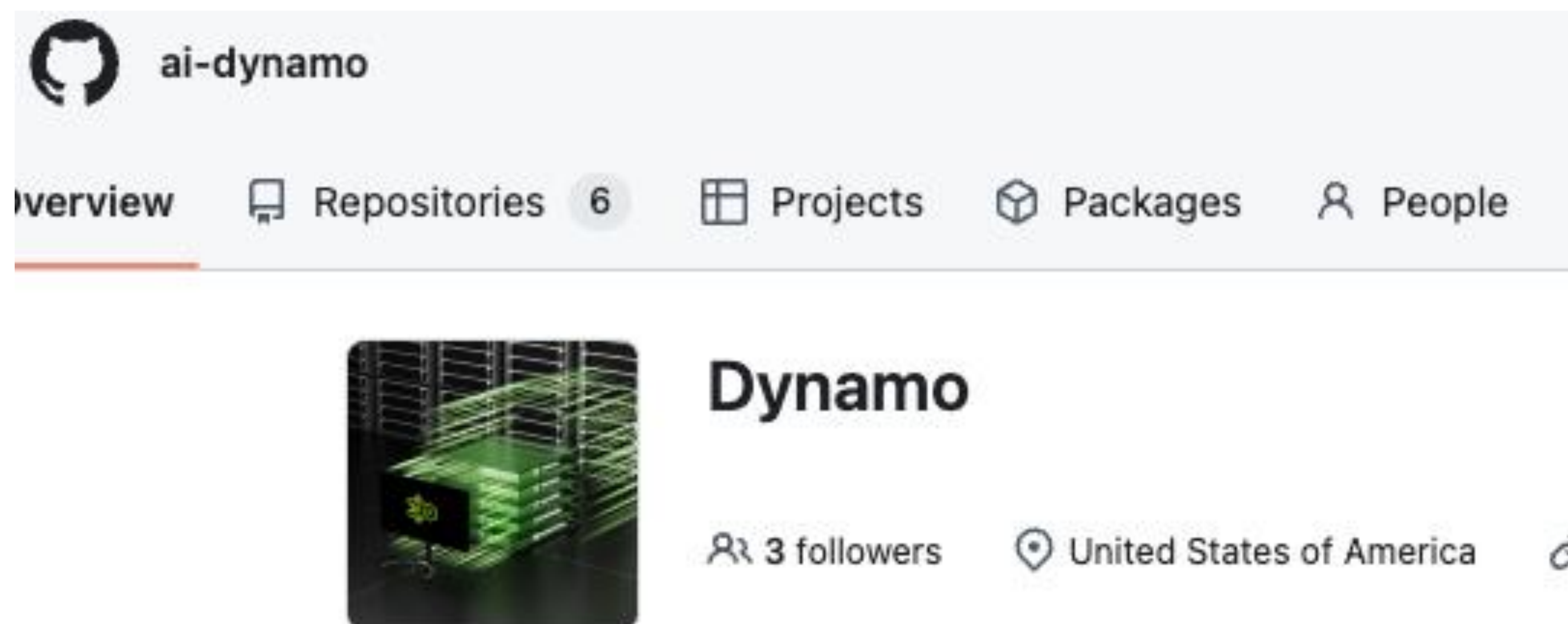■ Random routing  ■ Dynamo with KV aware routing

3x faster

**Avg request latency**

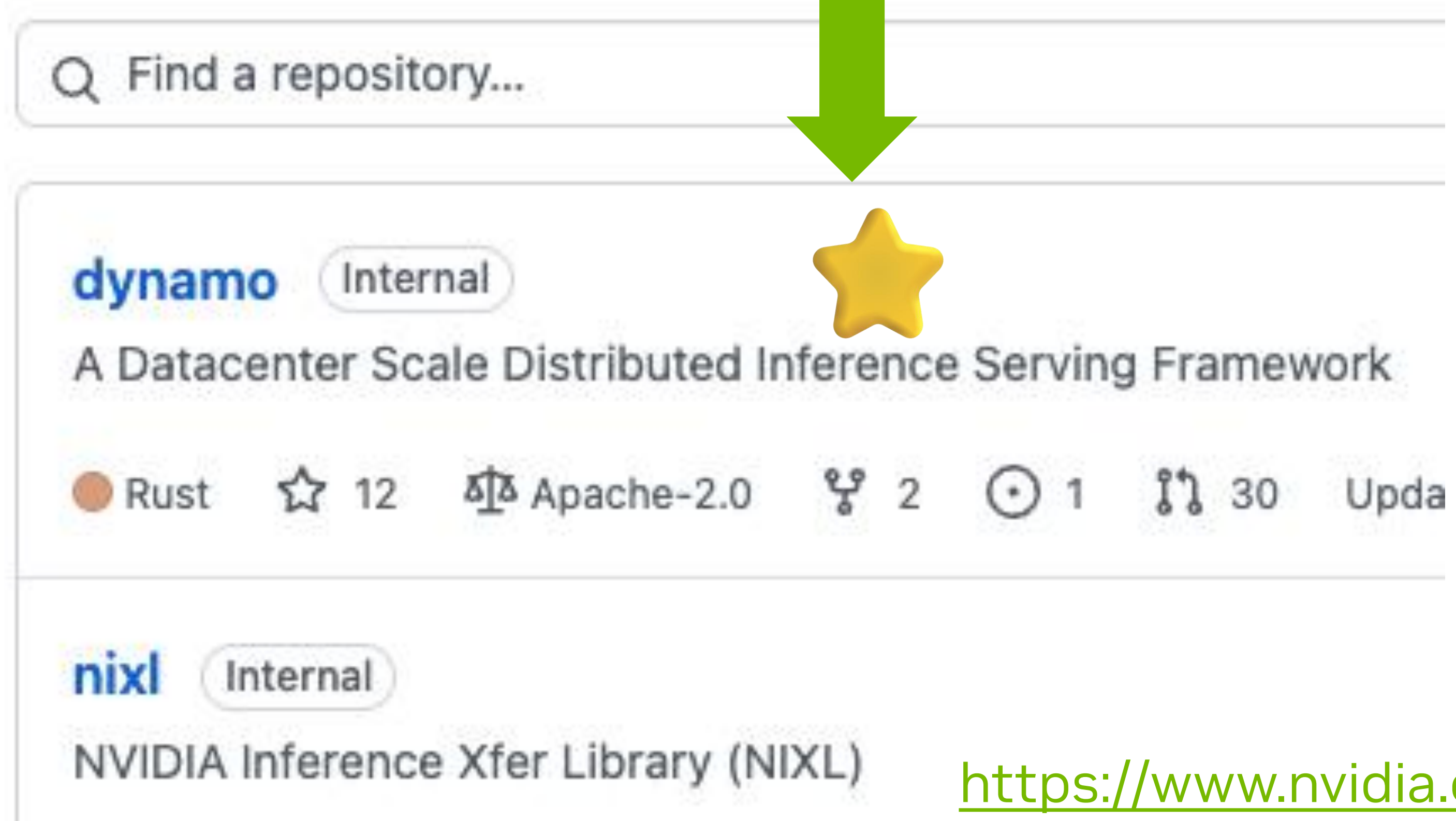■ Random routing  ■ Dynamo with KV aware routing

2x faster

**Tested with R1 Distilled Llama 70B over 2 nodes of 8 x H100s with vLLM 0.7.3**

15

# NVIDIA Dynamo as of Today (=GTC week)



**Github: nvidia.com/dynamo**

- Apache 2 license and public CI

- Pip wheels on PyPi

- Rust for perf and Python for extensibility

- Discord for developer community

- Dynamo CLI
  - dynamo run: Quick start with model, input and output
  - dynamo serve: Construct graph of workers and serve
  - (EA) dynamo build: Containerize
  - (EA) dynamo deploy: Deploy to K8

- Three backends: TRT-LLM, vLLM, & SGLang

- Disaggregated serving with TRT-LLM and vLLM

- KV aware routing with TRT-LLM and vLLM

- (EA) KV manager with vLLM

- NIXL for RDMA and TCP (fallback for AWS EFA)

https://www.nvidia.com/en-us/on-demand/session/gtc25-S73042/
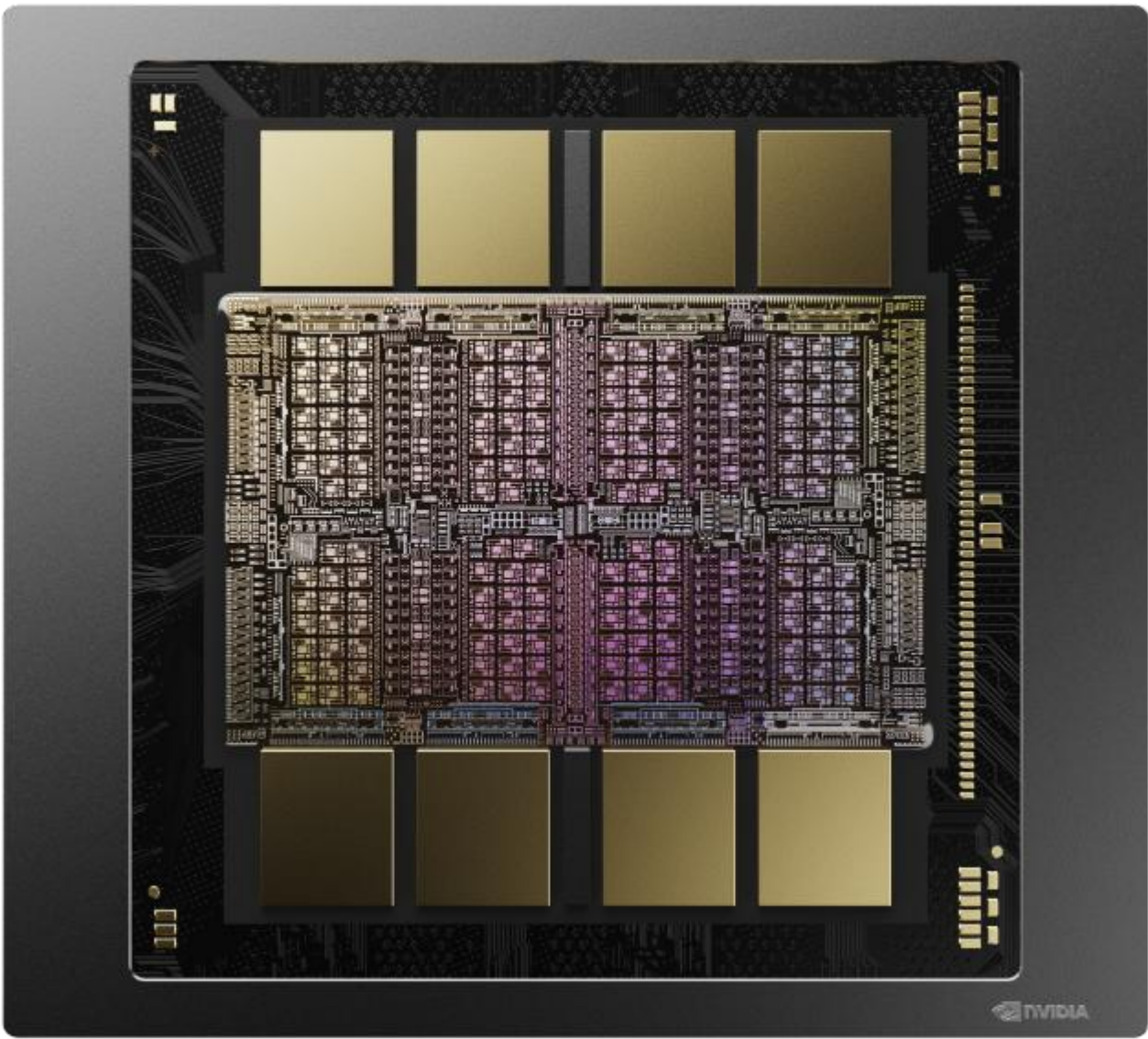
# Blackwell Ultra

# Announcing Blackwell Ultra

## Built for the Age of AI Reasoning

NVIDIA GB300 NVL72

BLACKWELL ULTRA GPU
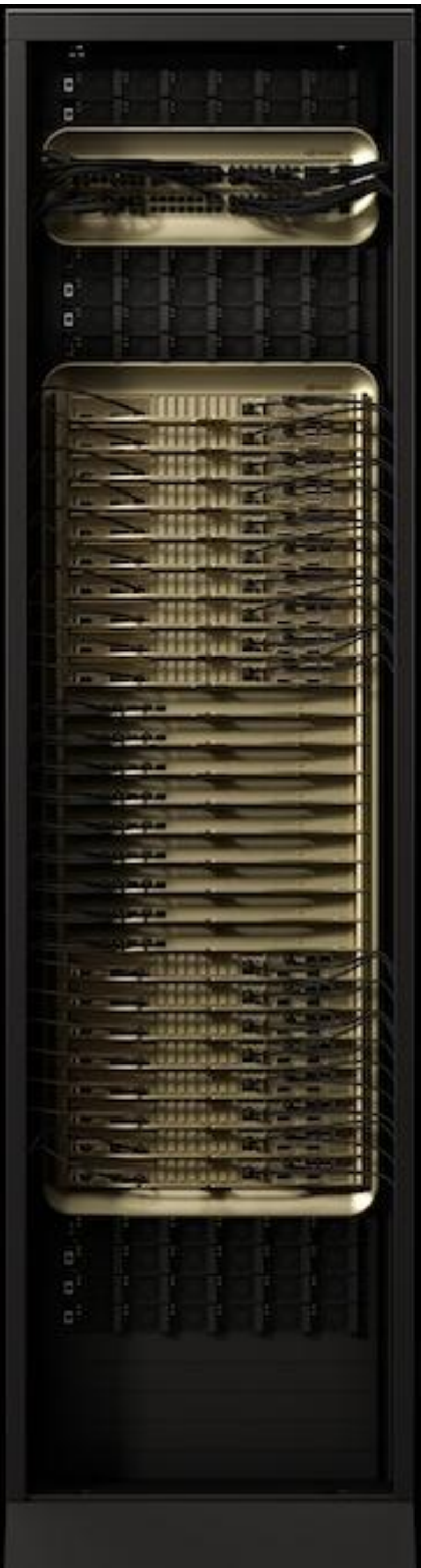
Blackwell 288GB GPU | 1.5x more FP4 Inference

**50X** AI Factory Output

Upgraded NVL72 Design for Improved Energy Efficiency

| | |
|---|---|
| **Upgraded FP4** | 15PF Dense |
| **HBM Memory** | Up to 288GB HBM3e |
| **Attention** | 2.5x Hopper |

| | |
|---|---|
| **FP4 (Dense)** | 1.1 ExaFLOPS |
| **HBM Memory** | 20 TB |
| **Fast Memory** | 40 TB |
| **Networking** | 14.4 TB/s |

NVIDIA

# NVIDIA GB300 NVL72

Built for the Age of AI Reasoning

Increasing AI Reasoning
Throughput and Responsiveness

## 50X

AI Factory Output

Interactive DeepSeek-R1 671B

| H100 | GB300 NVL72 |
|---|---|
| 35 TPS for one user | 350 TPS for one user |
| 1.5 Mins | 10 Seconds |

*Projected performance subject to change*

*Compared to HGX H100 | AI Reasoning on Deepseek R1*

# NVIDIA Blackwell Ultra AI Factory Output



Tokens per Second (TPS) in 1 MW — Throughput

GB300 NVL72
With NVIDIA Dynamo
FP4

50x

Hopper

35

350

TPS for One User

Responsiveness

DeepSeek R1 ISL = 32K, OSL = 8K, GB300 NVL72 with FP4 Dynamo disaggregation. H100 with FP8 In-flight batching. Projected performance subject to change.

NVIDIA

# DGX GB300

Accelerate Real-Time State of the Art Inference Models

---

- Based on the NVIDIA GB300 NVL72 rack architecture

- Provides the foundation for NVIDIA DGX SuperPOD with DGX GB300

- Powered by Grace Blackwell Ultra Superchips connected with fifth-generation NVIDIA NVLink

- 36 Grace CPUs and 72 Blackwell Ultra GPUs

- 1.4 exaFLOPS of AI performance and 38TB of fast memory

- Massive shared memory space to accelerate the most data-intensive workloads

# NVIDIA DGX GB200

Always-available enterprise infrastructure for mission-critical AI

- The building block of DGX SuperPOD with DGX GB200 systems

- Based on the NVIDIA GB200 NVL72

- Provides a fully-integrated, ready-to-scale infrastructure solution for generative AI

- Built with 36 GB200 Superchips and fifth-gen NVLink

- Connects 36 Grace CPUs and 72 Blackwell GPUs for compute intensive workloads

- 1.4 exaFLOPS of AI performance and 30TB of fast memory

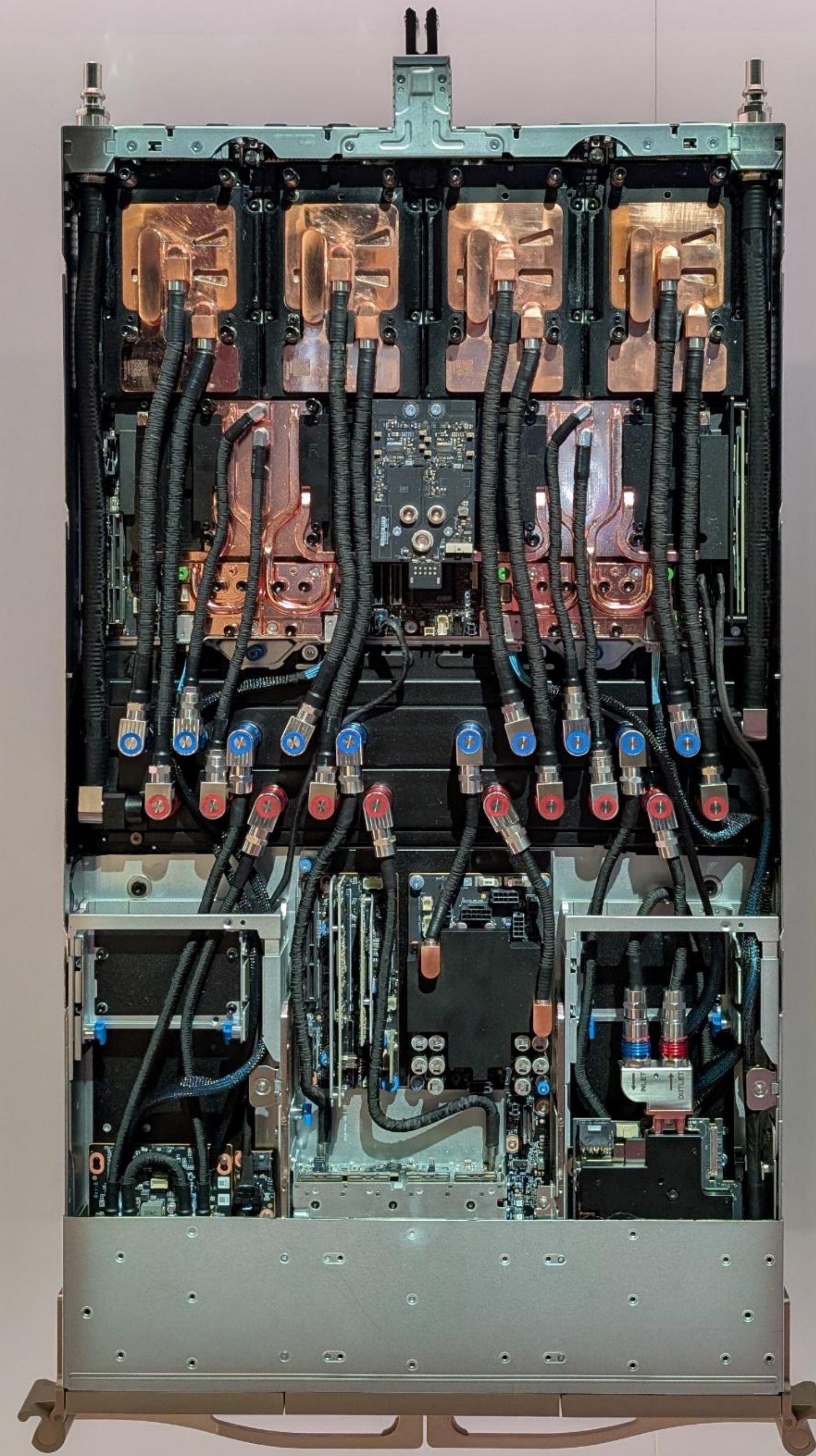- Handles the most complex generative AI workloads
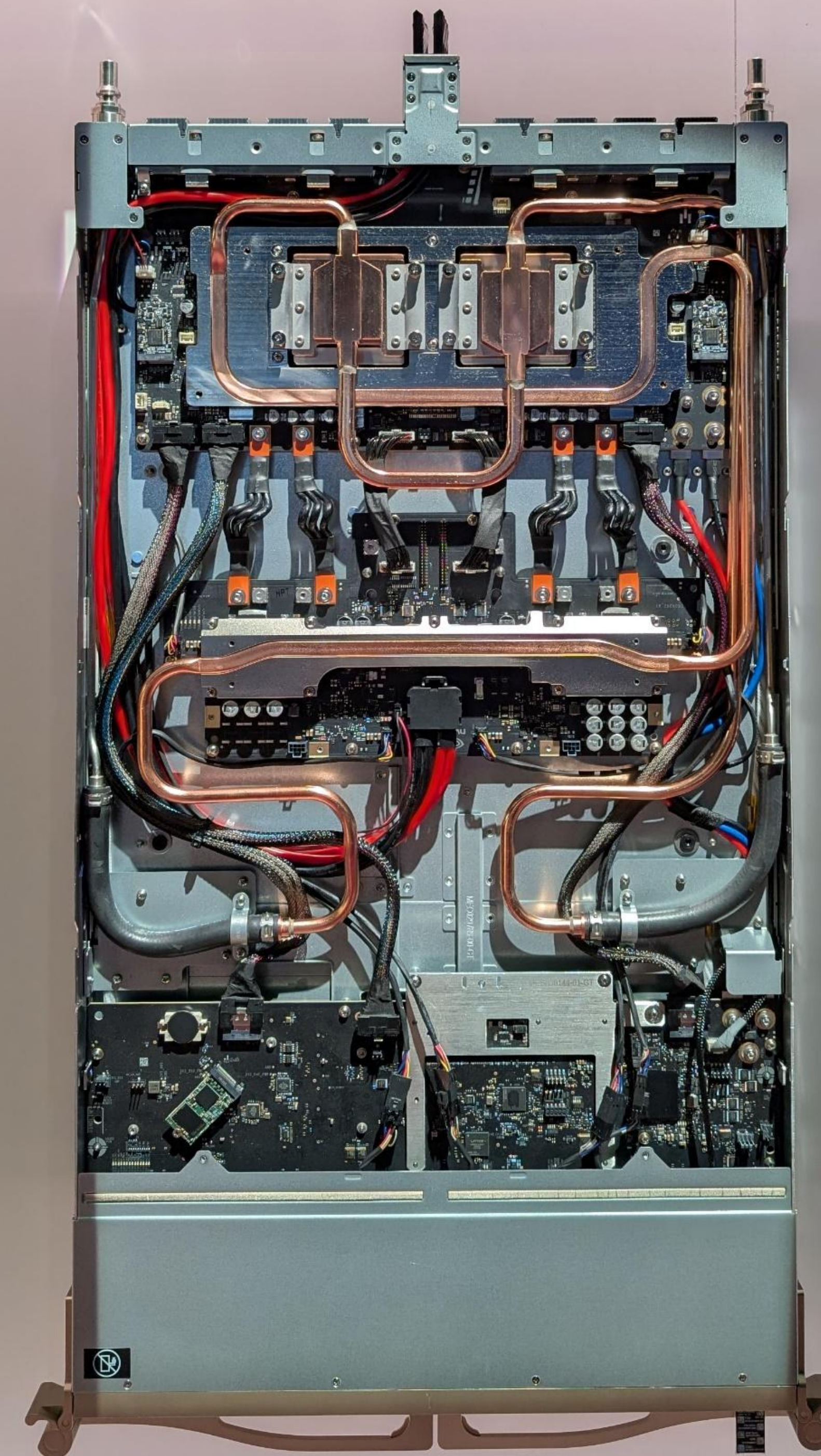


NVIDIA.

NVIDIA GB300
Compute Tray

NVIDIA GB300
NVLink Switch Tray
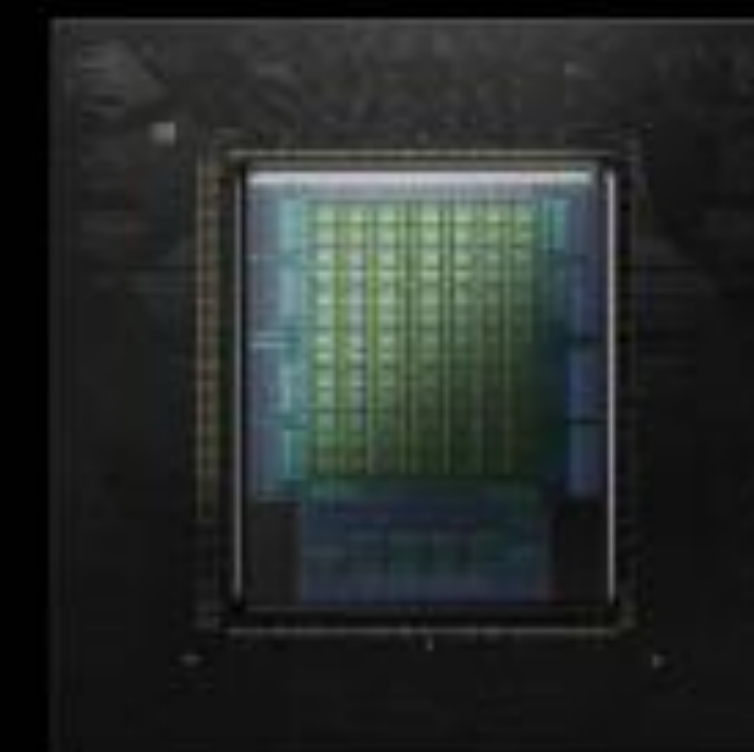
NVIDIA

# Vera Rubin NVL144

Second Half 2026

3.6 EF FP4 Inference
1.2 EF FP8 Training
**3.3X GB300 NVL72**

13 TB/s HBM4
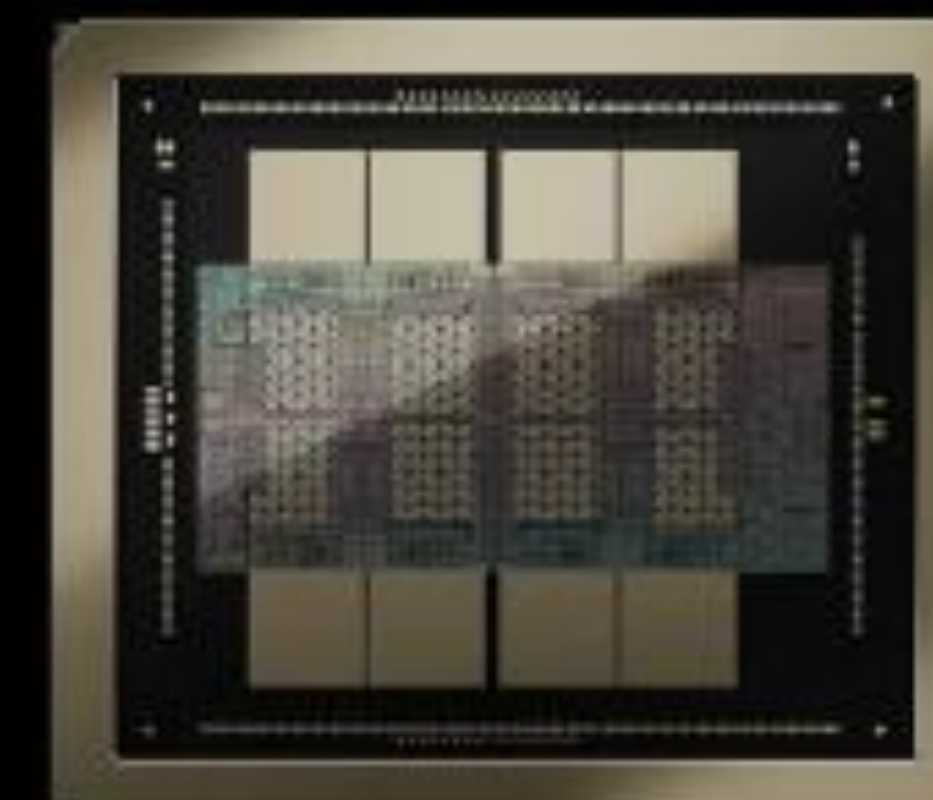75 TB Fast Memory
**1.6X**

260 TB/s NVLink6
**2X**

28.8 TB/s CX9
**2X**

*Oberon Rack*
*Liquid Cooled*

## Vera

88 Custom Arm Cores
176 Threads
1.8 TB/s NVLink-C2C

## Rubin

2 Reticle-Sized GPUs
50PF FP4 | 288GB HBM4

# Rubin Ultra NVL576

Second Half 2027
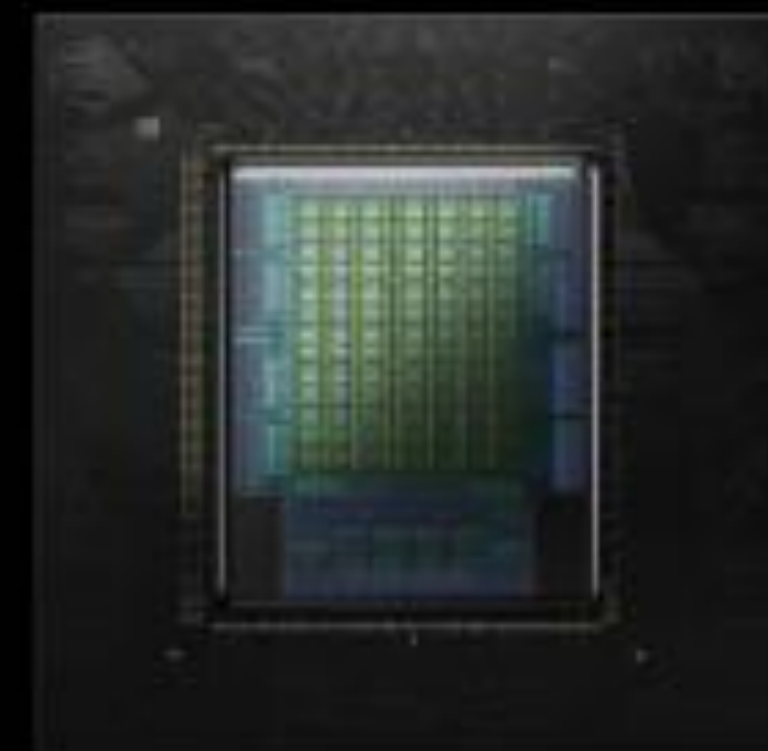
15 EF FP4 Inference
5 EF FP8 Training
14X GB300 NVL72
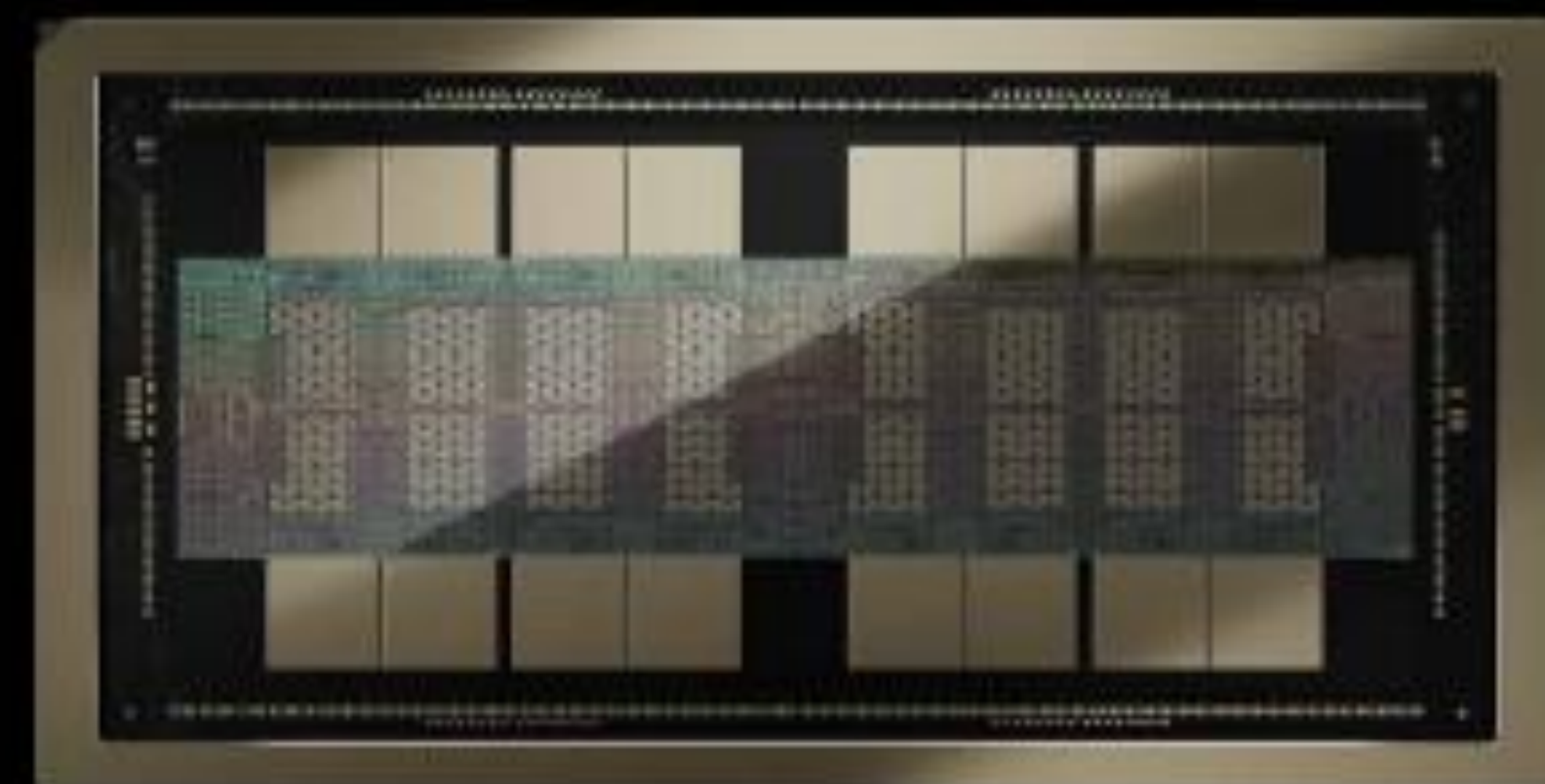
4.6 PB/s HBM4e
365 TB Fast Memory
8X

1.5 PBs NVLink7
12X

115.2 TB/s CX9
8X

*Kyber Rack*
*Liquid Cooled*

## Vera

88 Custom Arm Cores
176 Threads
1.8 TB/s NVLink-C2C

## Rubin Ultra

4 Reticle-Sized GPUs
100PF FP4 | 1TB HBM4e

# DGX B300

Accelerated Infrastructure for the Era of AI Reasoning

- Newest air-cooled DGX system with NVIDIA Blackwell Ultra GPUs

- All new system design seamlessly integrates into NVIDIA MGX or traditional enterprise racks

- **2.3TB of GPU memory**, enabling training and inference of complex models

- Equipped with NVIDIA ConnectX-8 high speed networking at **800Gb/s**

- Delivers **72 PFLOPS AI training** and **144 PFLOPS AI inference** performance

- Purpose-built platform for the era of AI reasoning, setting a new bar for LLM inference

**10U Chassis | ~14 kW system**
**Designed for the modern data center**

# NVIDIA DGX B200

The foundation of the modern AI data center

- Air-cooled DGX system with 8X NVIDIA Blackwell GPUs

- **1.4TB of GPU memory**, enabling training of large generative AI models , **64 TB/s Bandwidth**

- **1.8 TB/s** NVLink GPU-to-GPU Bandwidth

- Purpose-built, unified platform for every workload from training, to fine-tuning, to inference

- Delivers **3X AI training** and **15X AI inference** performance as previous generation (DGX H100)

10U Chassis | ~14.3 kW system
Deployable in today's data centers

# Blackwell Ecosystem

## Blackwell Ultra Coming Later 2025
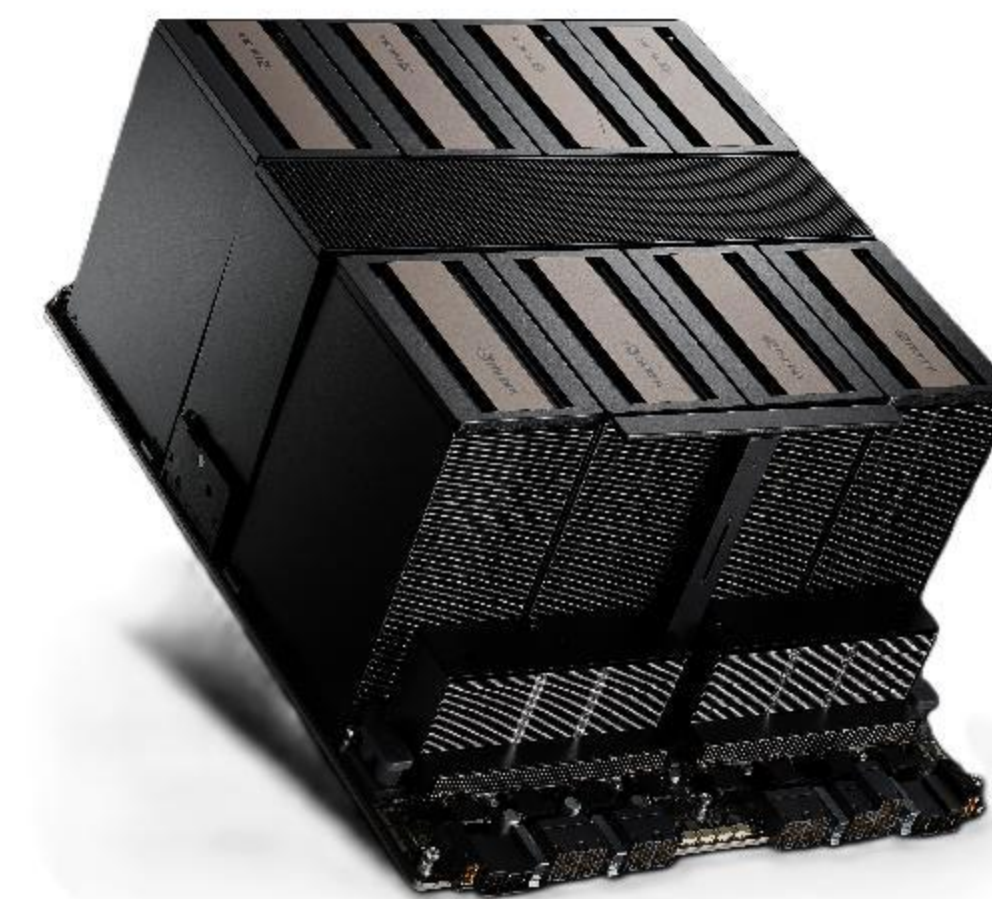


GB200 NVL72
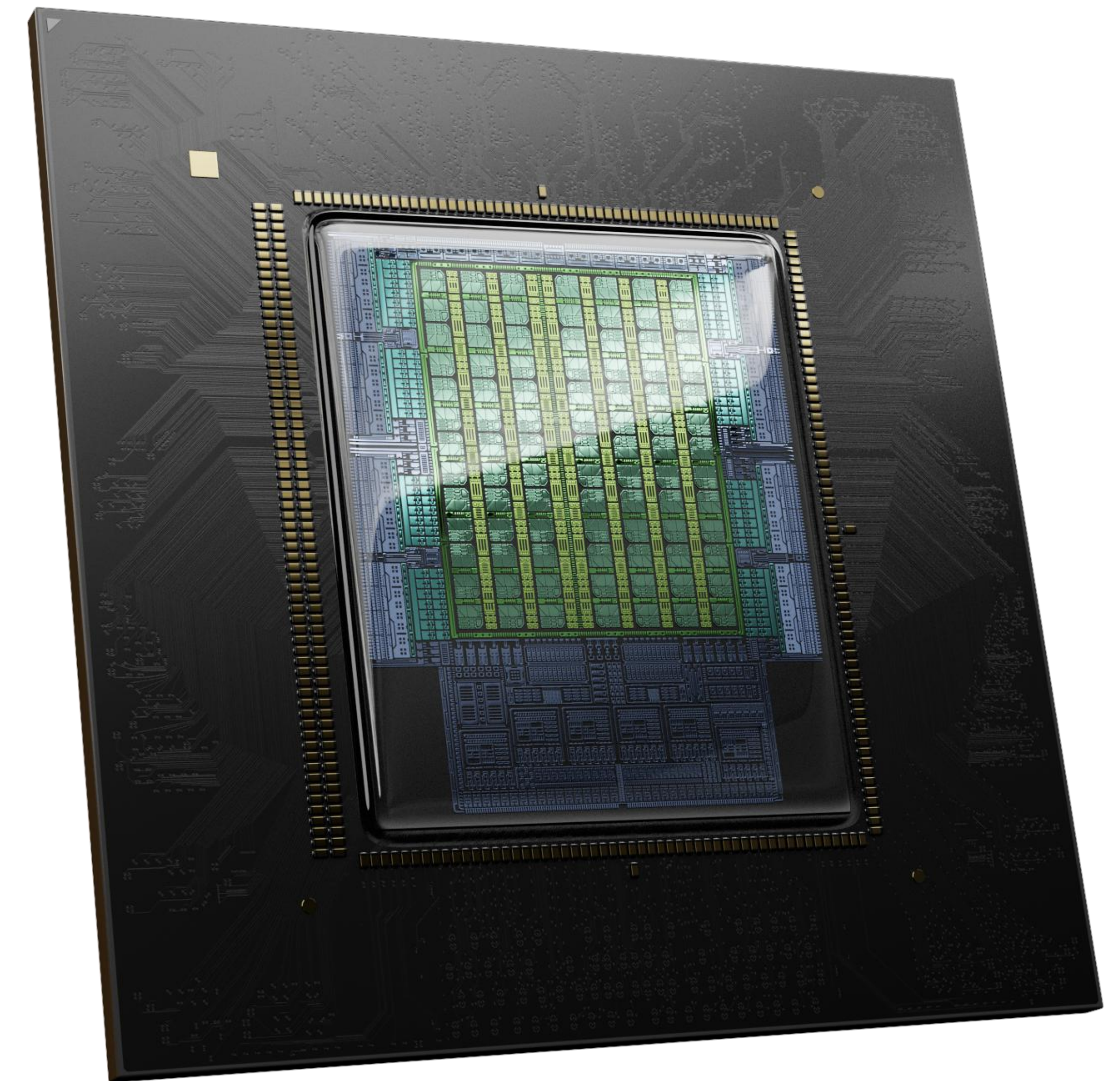
HGX B200

HGX B300 NVL16

GB300 NVL72

# Vera CPU / Networking

# Introducing NVIDIA Vera:
# A Next-Generation CPU
## AI Factories, compute and memory intensive CPU workloads

- **>2X CPU Compute Capability, 2.4x threads vs. Grace,**
  - 88 Cores with Spatial Multi-Threading

- **5X Memory BW per Watt**
  - Memory power per socket under 50W vs 280W MRDIMM

- **>3x Memory Capacity**
  - 1.5 TB of coherent LPDDR5X in Vera Rubin platforms
  - 2 TB of DDR5

- **>2x Bisection Bandwidth vs. x86**
  - Single NUMA design for optimal tuning out-of-box

- **7x Faster GPU connectivity vs. traditional CPU**
  - 1.8 TB/s NVLink-C2C CPU:GPU bandwidth vs PCIe Gen 6

*Projected performance subject to change*

# Spectrum-X Ethernet Accelerates the World's Largest AI Supercomputers



Scaling Compute to 400K GPU AI

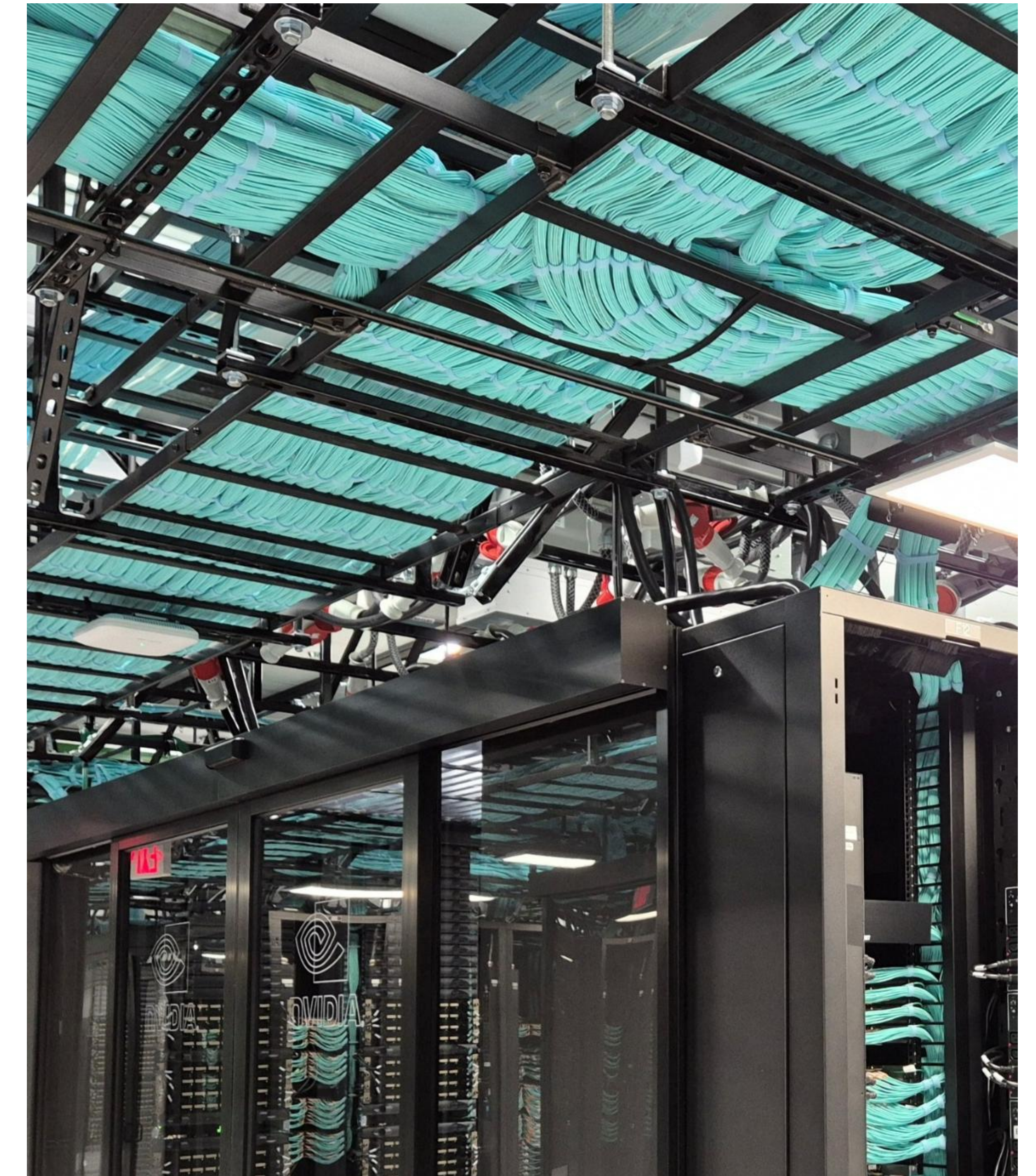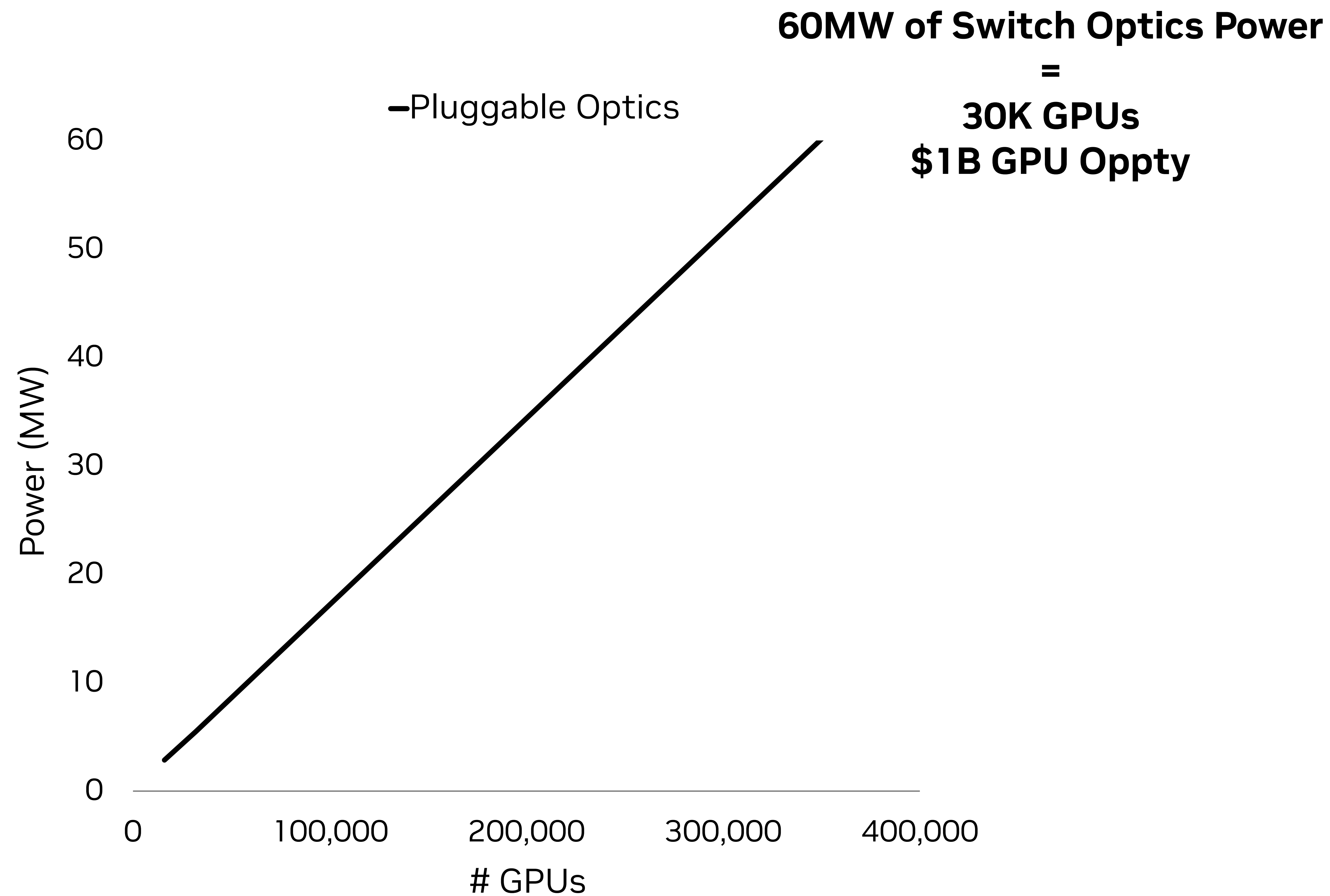| 100K | 400K | 2.4M | 40 MW |
|------|------|------|-------|
| Servers | GPUs | Optical Transceivers | Transceiver Power |

NVIDIA

# Power and Reliability Challenges of AI Scale-Out and Density

**60MW of Switch Optics Power**
**=**
**30K GPUs**
**$1B GPU Oppty**

**—Pluggable Optics**

Power (MW)

60

50

40

30

20

10

0

0          100,000      200,000      300,000      400,000
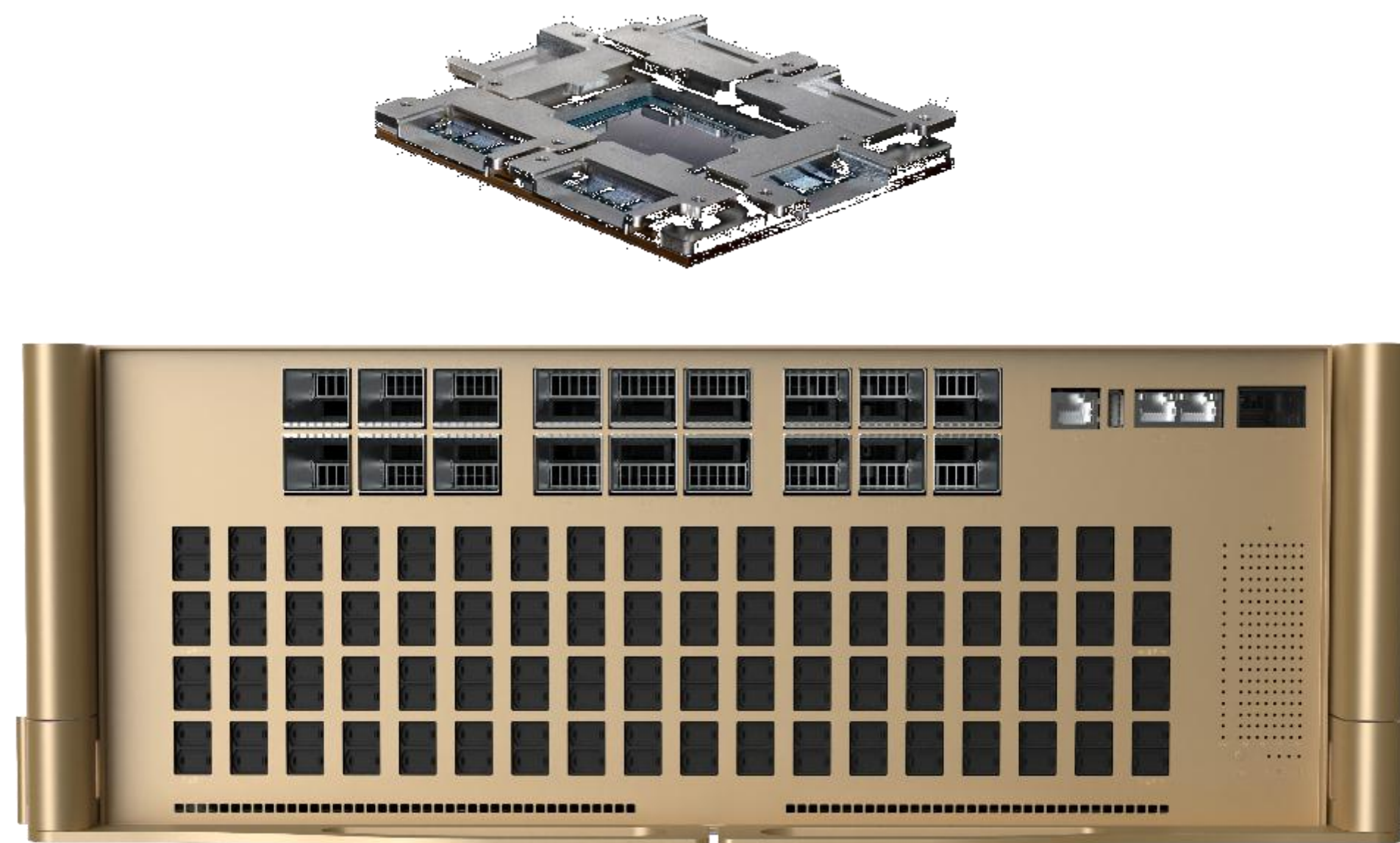
# GPUs

# Announcing NVIDIA Photonics Switch Systems

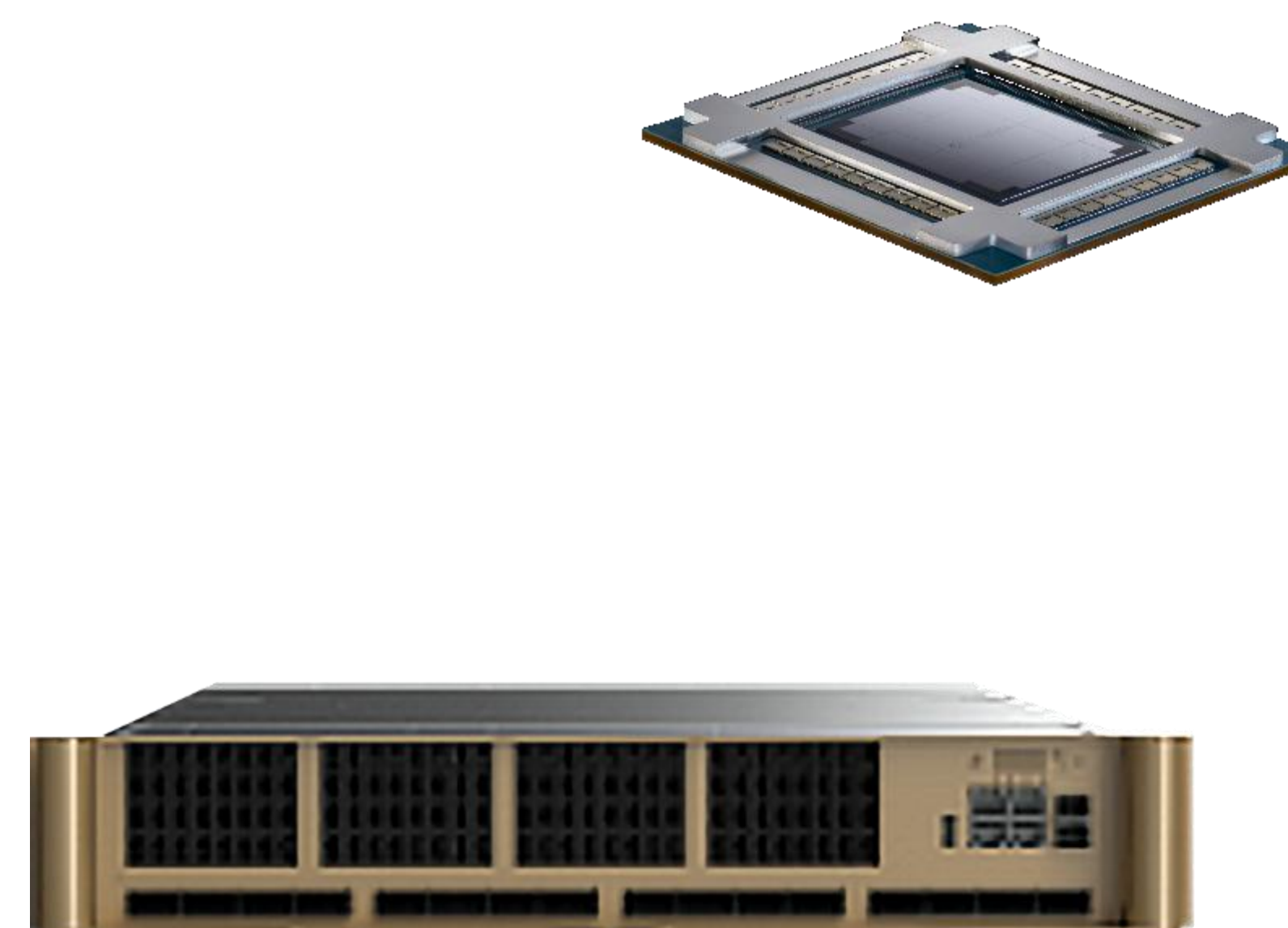Co-packaged optics networking switches to scale AI factories to millions of GPUs

## Quantum-X Photonics



Quantum 3450-LD

**115Tb/s**
144 ports of 800G
(576 ports x 200G)
Liquid cooled

## Spectrum-X Photonics



Spectrum SN6810

Spectrum SN6800

**102.4Tb/s**
128 ports of 800G
(512 x 200G)
Liquid cooled

**409.6Tb/s**
512 ports of 800G
(2048 x 200G)
Liquid cooled
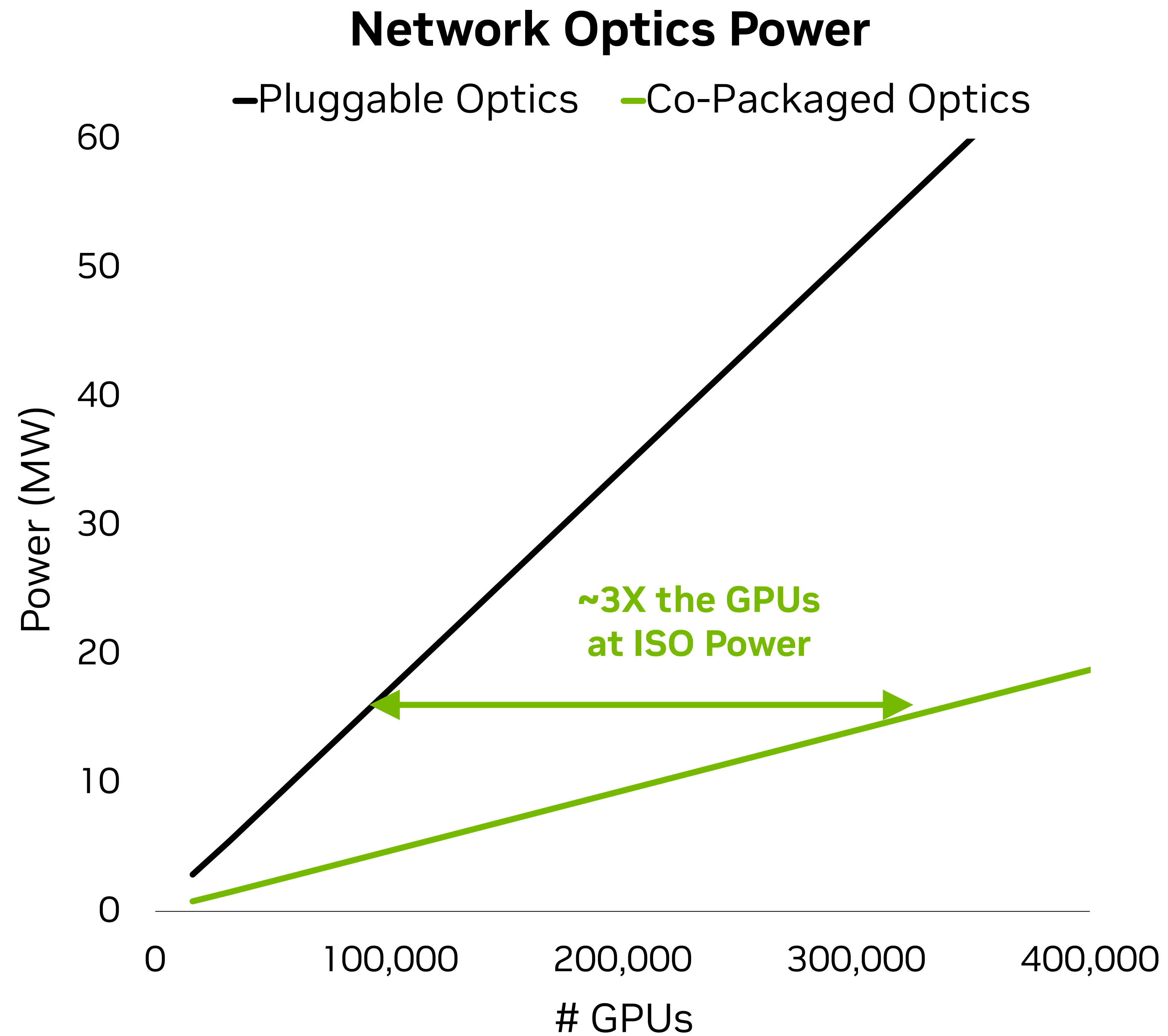
**3.5X**
Power efficiency

**10X**
Higher resiliency

**1.3X**
Faster time to deploy

*Projected performance subject to change*

NVIDIA

# NVIDIA Photonics Solves Power and Reliability Challenges of AI Scale-Out

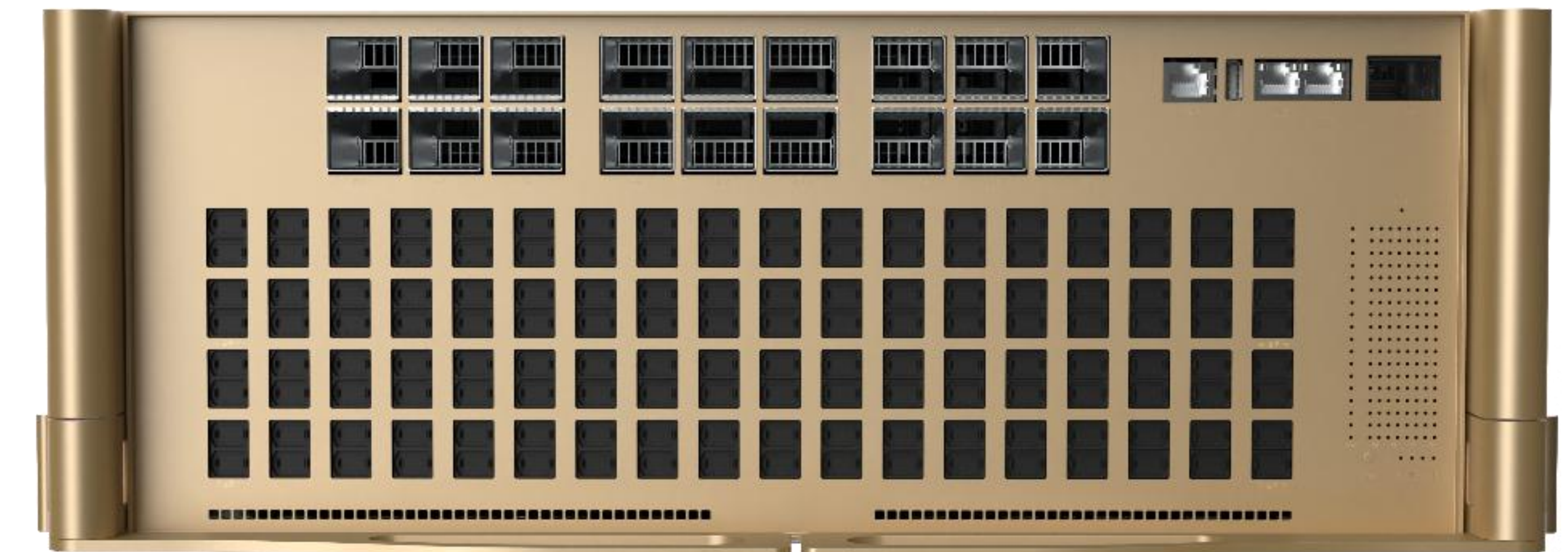Co-packaged silicon photonics networking switches to scale AI factories to millions of GPUs

## Network Optics Power

— Pluggable Optics — Co-Packaged Optics



**~3X the GPUs at ISO Power**

Power (MW) vs # GPUs

**72**
Transceivers Replaced

**432**
Fewer Lasers

**3.5X**
Power efficiency

**10X**
Higher resiliency

**1.3X**
Time to operation

# DGX Spark