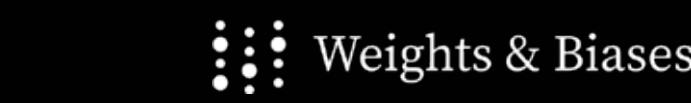
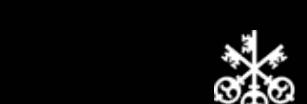
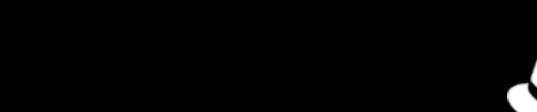
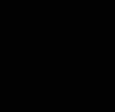
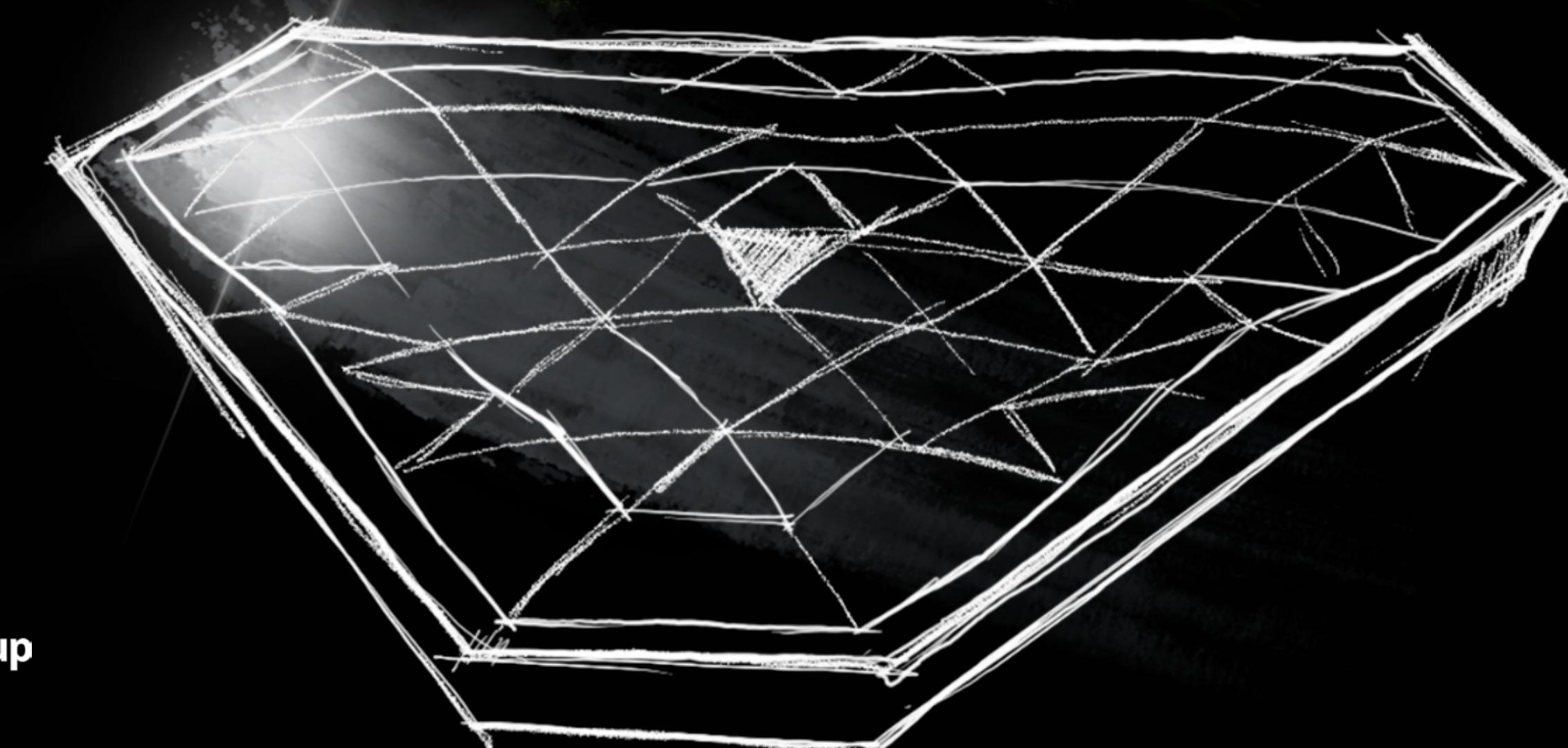


nVIDIA

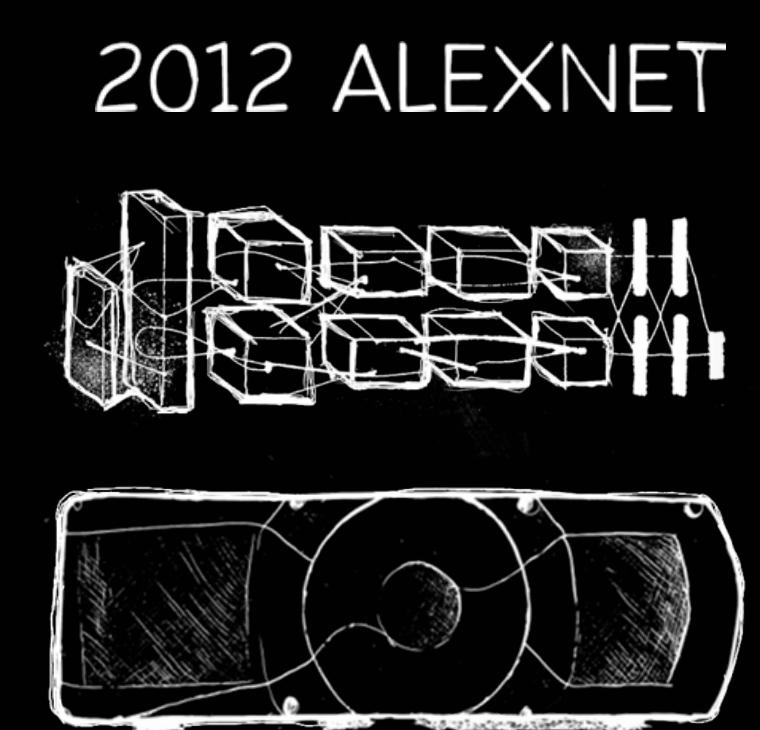


# GTC25









2012 ALEXNET

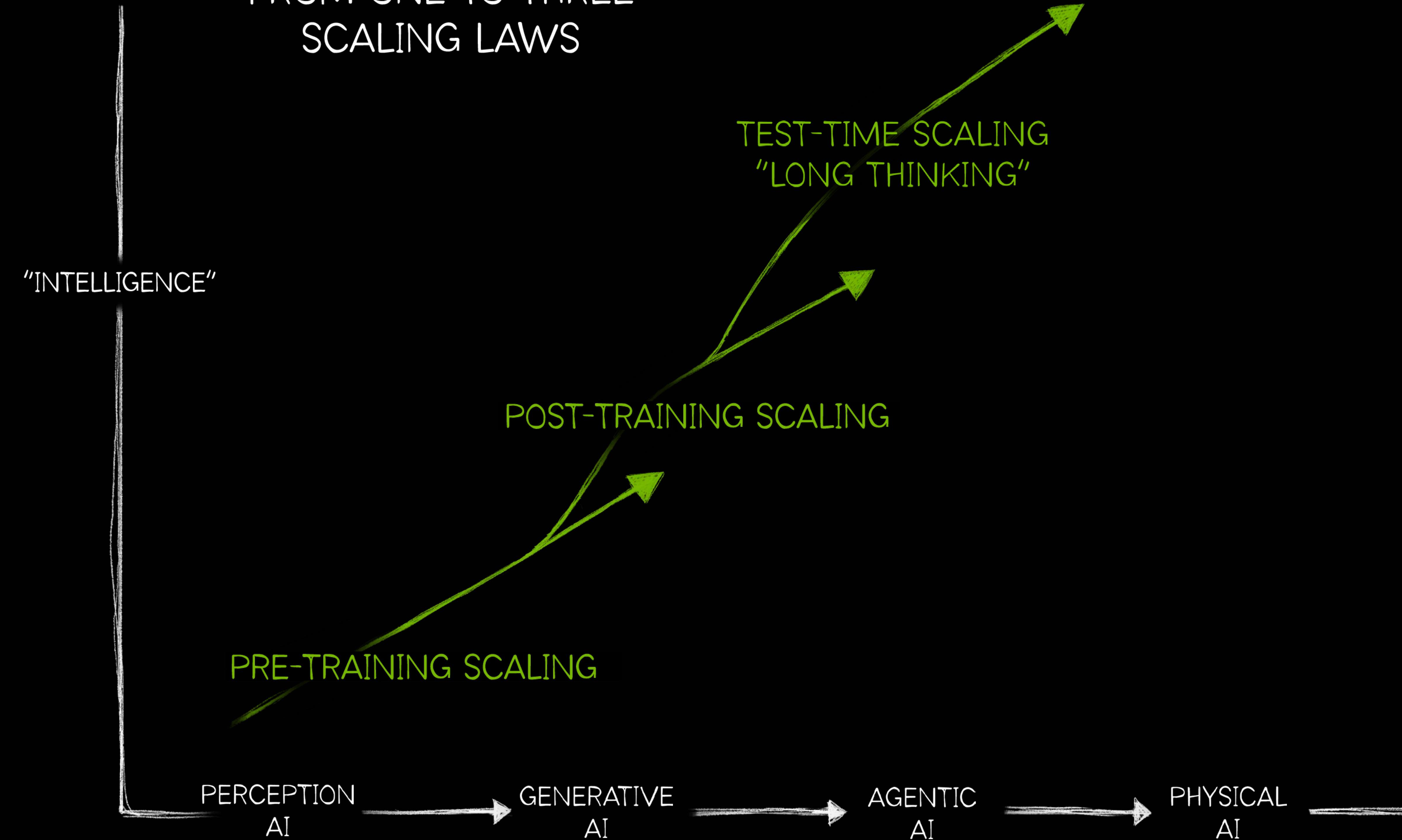
PERCEPTION AI  
SPEECH RECOGNITION  
DEEP RECSYS  
MEDICAL IMAGING

GENERATIVE AI  
DIGITAL MARKETING  
CONTENT CREATION

AGENTIC AI  
CODING ASSISTANT  
CUSTOMER SERVICE  
PATIENT CARE

PHYSICAL AI  
AUTONOMOUS VEHICLES  
GENERAL ROBOTICS

## FROM ONE TO THREE SCALING LAWS



## FROM ONE TO THREE SCALING LAWS

"INTELLIGENCE"

PRE-TRAINING SCALING

POST-TRAINING SCALING

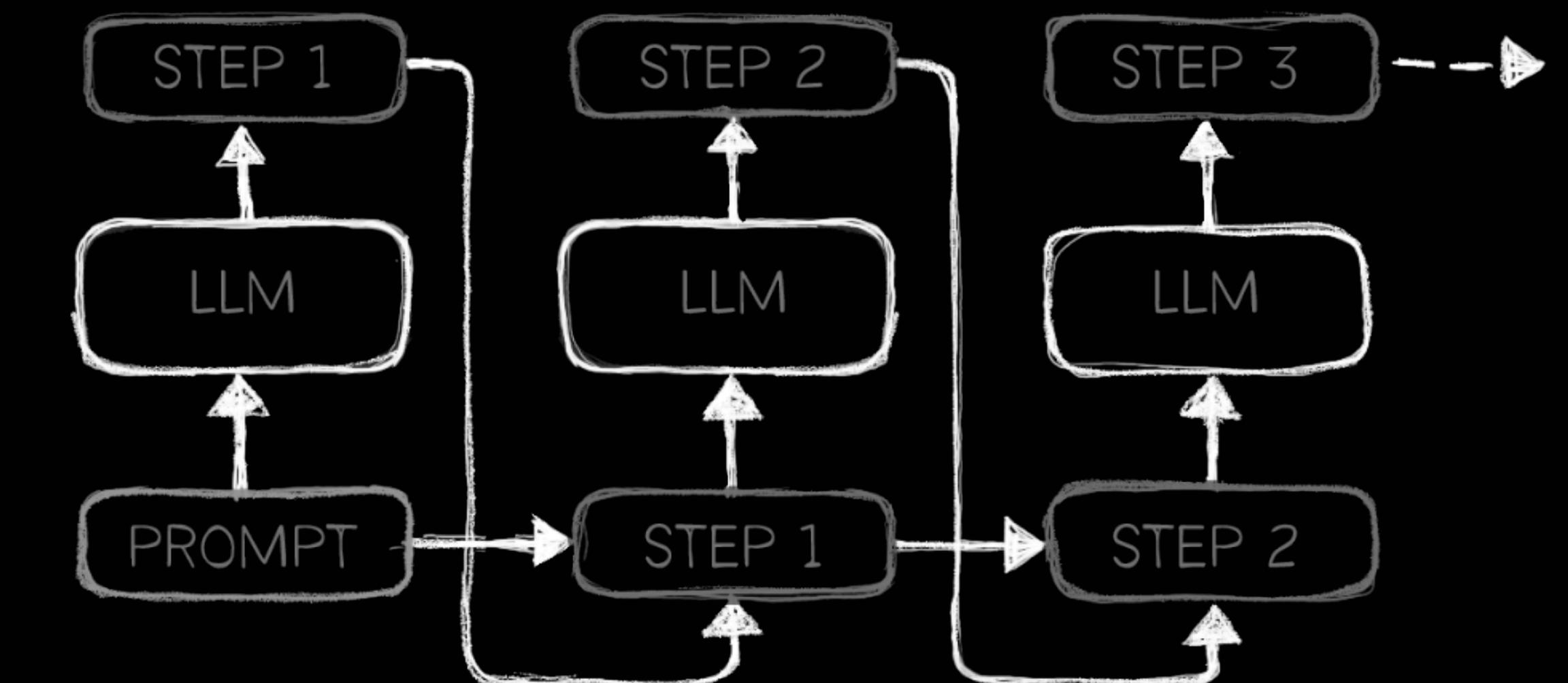
TEST-TIME SCALING  
"LONG THINKING"

PERCEPTION  
AI

GENERATIVE  
AI

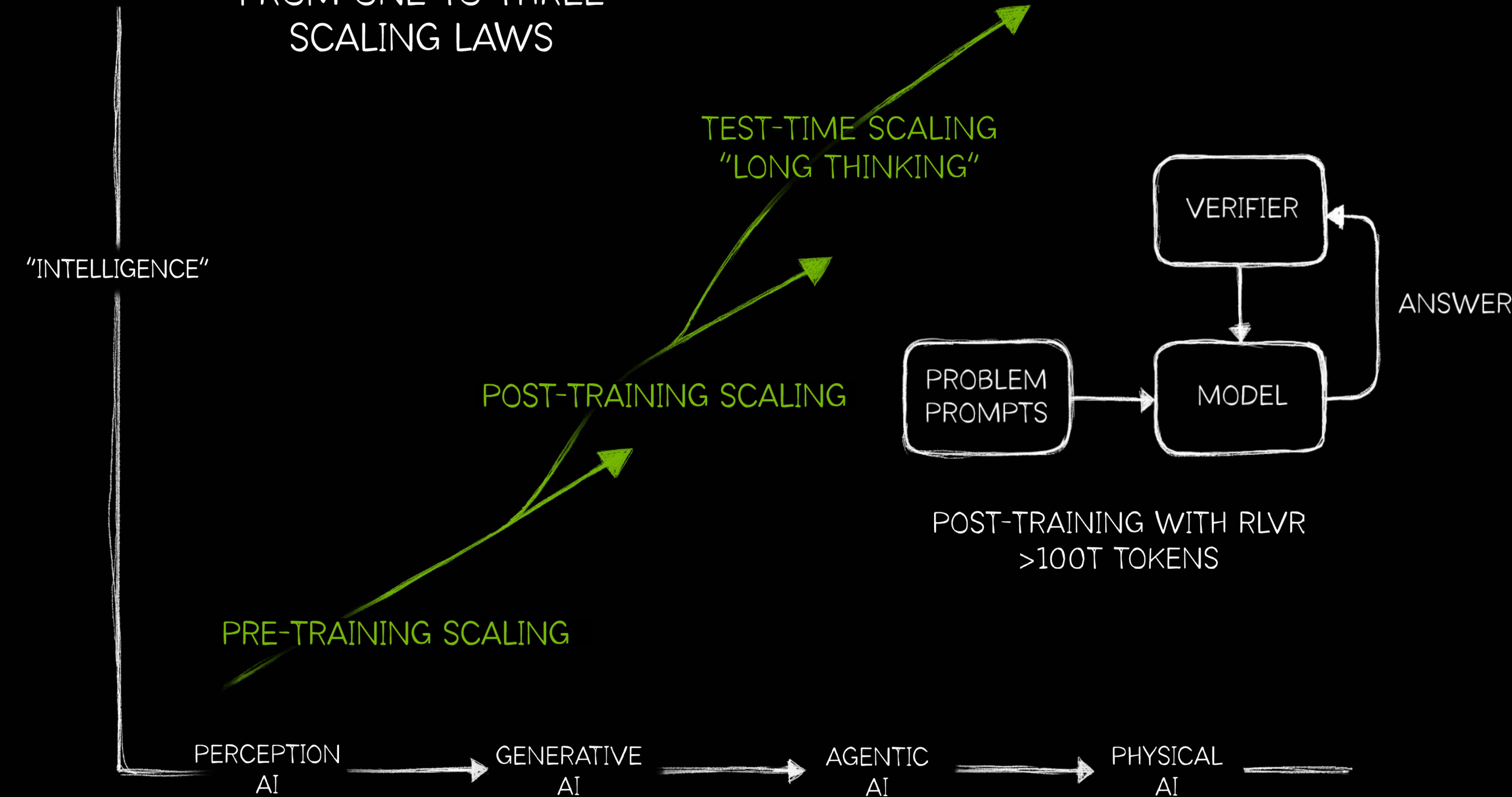
AGENTIC  
AI

PHYSICAL  
AI

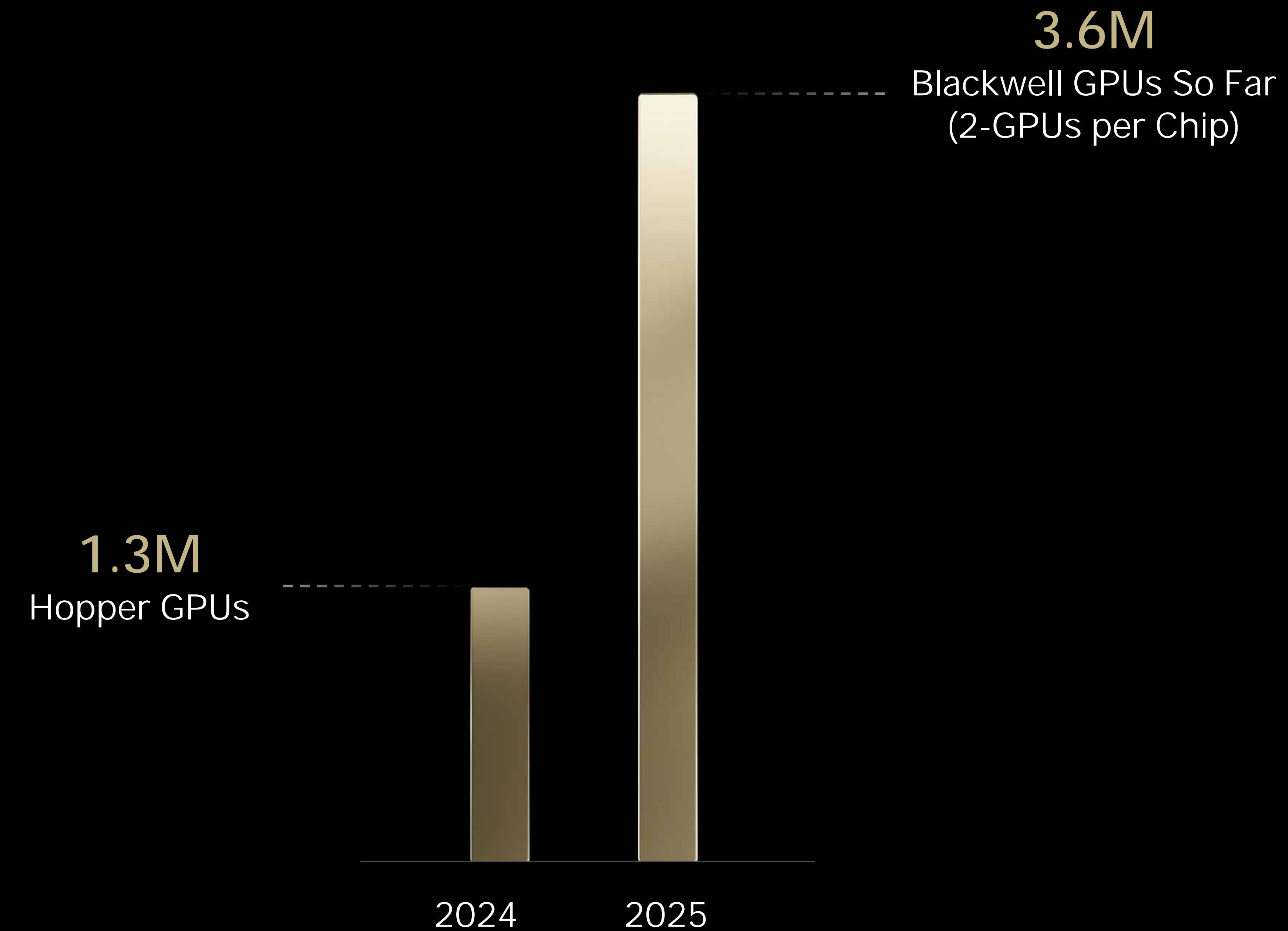


REASONING AI INference  
COMPUTE  
>100X ONE-SHOT

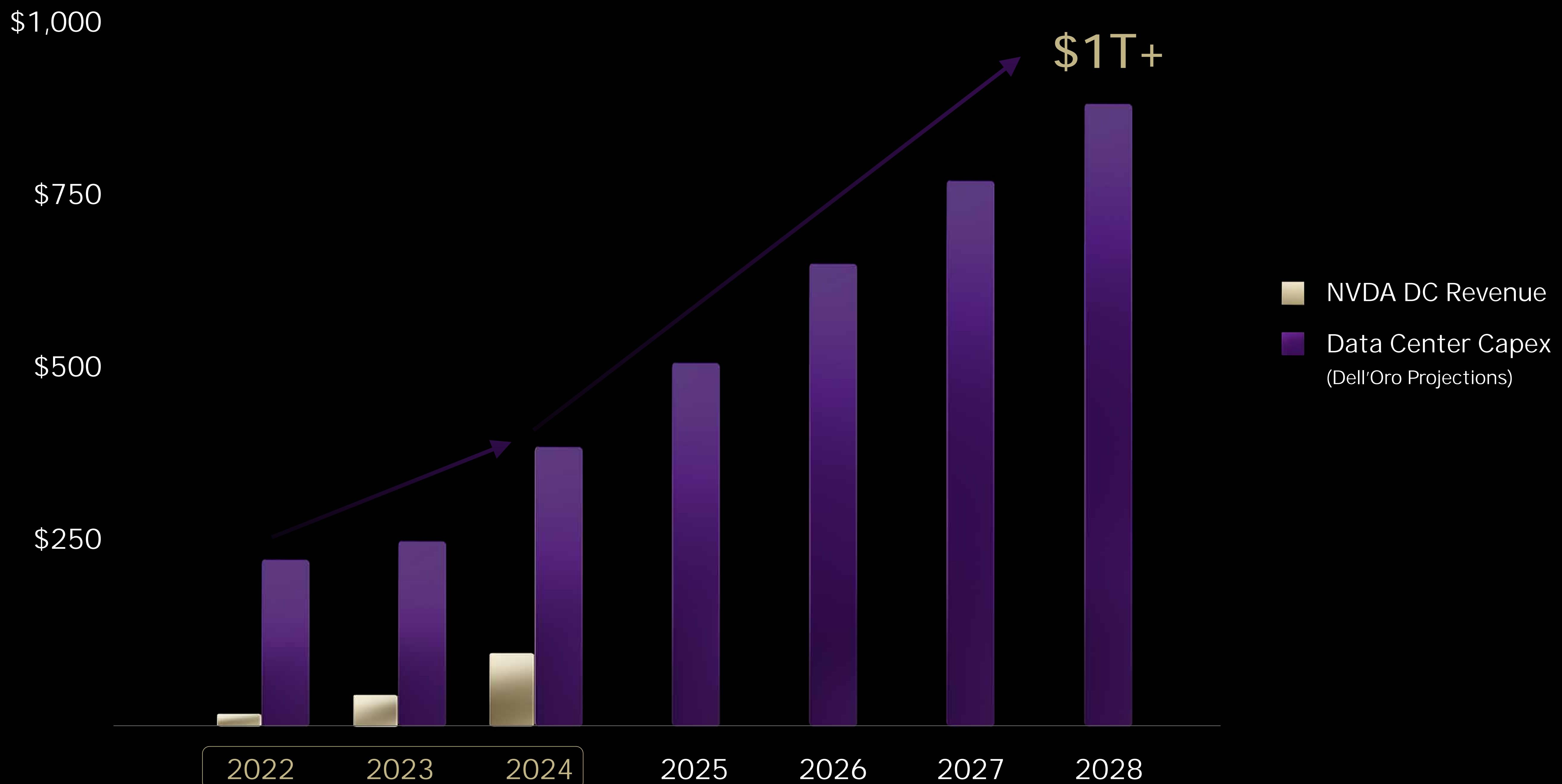
## FROM ONE TO THREE SCALING LAWS



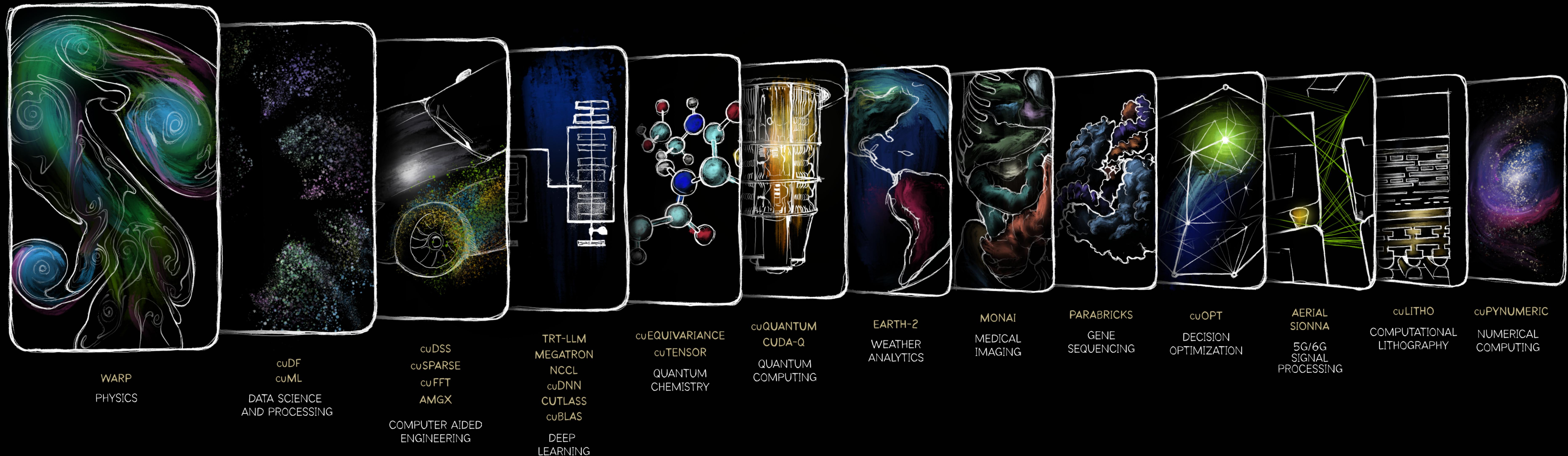
## Top 4 US CSPs



# Computing At Inflection Point

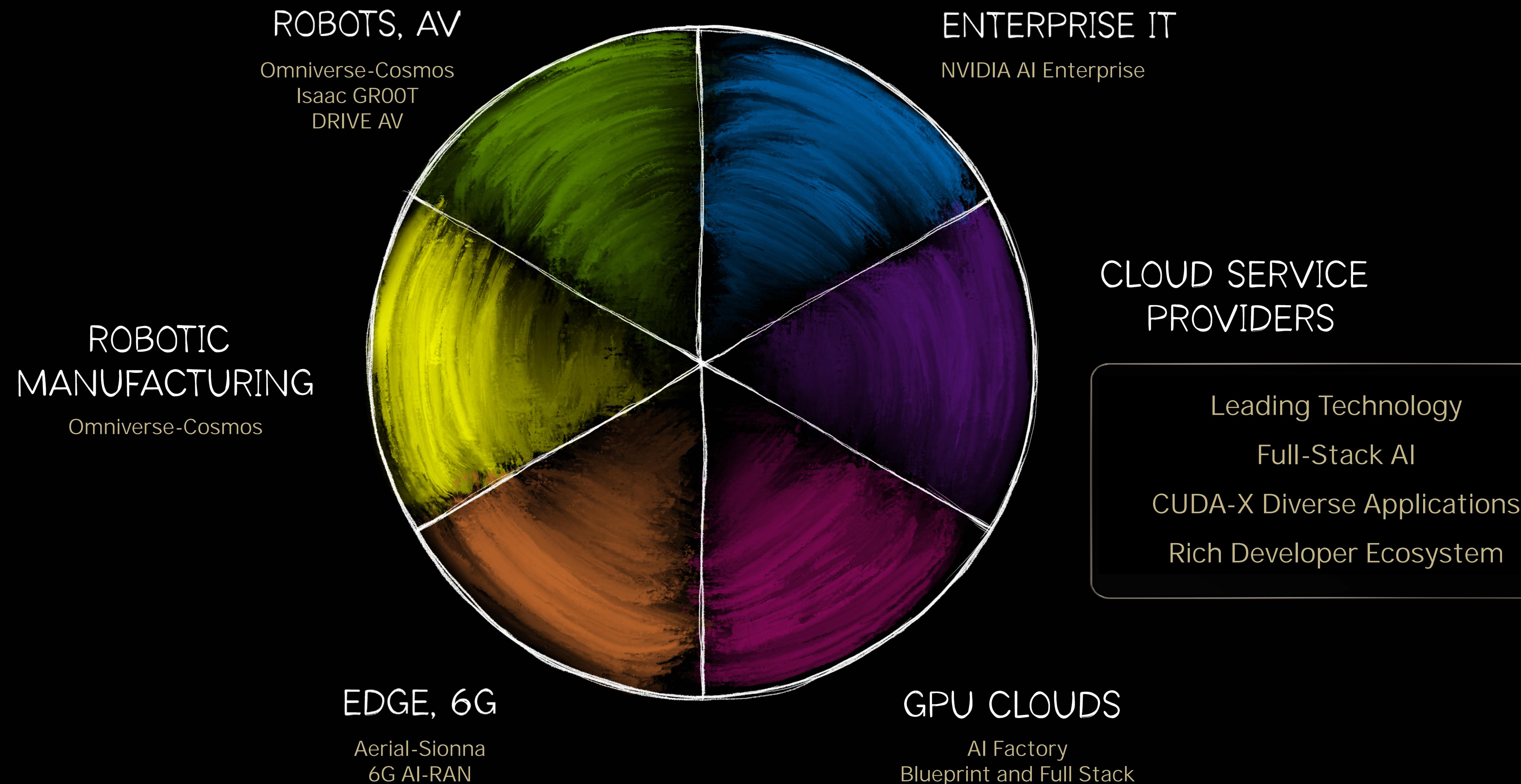


# CUDA-X FOR EVERY INDUSTRY





# AI FOR EVERY INDUSTRY







# Announcing NVIDIA Halos

## Chip-to-Deployment AV Safety System

EN

AI Conference March 17-21, 2025 | Keynote March 18 | Exhibits March 18-21 | Workshops March 18-20 | San Jose, CA & Virtual

Register Now | Log In

NVIDIA GTC

Keynote Session Catalog Agenda Attend Sponsors & Exhibitors More

### Autonomous Vehicle Safety Day

Safety is the top priority in autonomous vehicles. In this workshop, we'll showcase how NVIDIA is accelerating the development of safe and secure AVs.

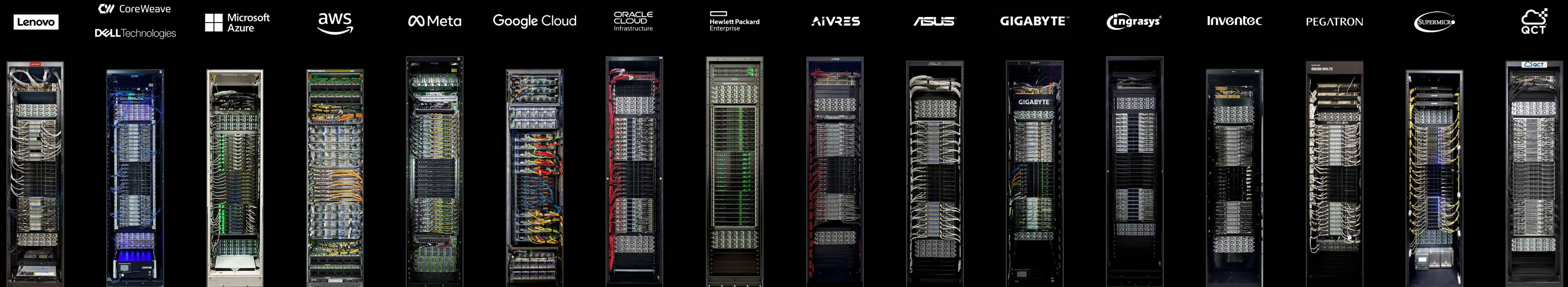
[View All Autonomous Vehicle Safety Sessions](#)

**Featured Sessions**

- Introduction to NVIDIA's Strategy for AV Safety [S74743]**  
Ed Schroming, Dr. Research Scientist, NVIDIA  
Wednesday Mar 21 | 10:00 AM - 11:00 AM PDT
- Guardrails for AV Safety Across the Product Life Cycle [S74744]**  
Marco Manelli, VP Safety, NVIDIA  
Wei Lin, VP of Automotive, NVIDIA  
Tuesday Mar 20 | 2:00 PM - 3:00 PM PDT
- Safety Regulation and Standardization in the Era of AI-Based AVs [S73722]**  
Marco Manelli, VP Safety, NVIDIA  
Tim Kuehne, Manager Regulatory Affairs, NVIDIA  
Wednesday Mar 21 | 1:00 PM - 2:00 PM PDT
- Navigating the High-Stakes Safety Challenges of Autonomous Driving [S73722]**  
Ricardo Manelli, VP Safety, NVIDIA  
Warren Markel, Vice President, Policy, ANSD National Auto Safety Council  
Melissa Wade, Sr. Director of Government Affairs, Autolbla  
Sneha Ray Singh, Head of Systems Engineering, Nuro  
Tuesday Mar 19 | 1:00 PM - 2:00 PM PDT
- Exploring NVIDIA's Strategy for AI-Defined Safety [CWE74730]**  
Jonas Nilsson, Director of Safety Engineering, Head of AI Safety, NVIDIA  
Ricardo Manelli, VP Safety, NVIDIA  
Agnieszka Skarzynski, Senior Software Engineer, NVIDIA  
Hui Li, Manager of Safety Engineering, NVIDIA  
Ed Schroming, Dr. Research Scientist, NVIDIA  
Tuesday Mar 20 | 1:00 PM - 2:00 PM PDT
- From Research to Production: Transforming AV Technology With AI [S72701]**  
Jonas Nilsson, Director of Safety Engineering, Head of AI Safety, NVIDIA  
Ricardo Manelli, VP Safety, NVIDIA  
Boris Isenbeck, Senior Research Scientist and Manager, Autonomous Vehicle Research, NVIDIA  
Xiaohou Wu, VP of Automotive, NVIDIA  
Tuesday Mar 19 | 1:00 PM - 2:00 PM PDT
- Accelerate the Future of AI-Defined Vehicles and Autonomous Driving [D740000]**  
Xiaohou Wu, VP of Automotive, NVIDIA  
Wednesday Mar 20 | 8:00 AM - 9:00 AM PDT
- Developing Next-Gen AVs with Physical AI-Powered World Foundation Models [D740001]**  
Norm Marks, Worldwide Field Ops, VP-Autonomous Driving, NVIDIA  
Natalia Dzhigadze, Director of Solutions Architecture and Engineering, Autonomous, NVIDIA  
Tuesday Mar 20 | 9:00 AM - 10:00 AM PDT
- Advancing AI Development With Sensor Simulation and Cosmos [D740002]**  
Gautham Shirodkar, Sr. Manager, Vehicle Metrics, NVIDIA  
Learn how the NVIDIA Orin-based Xavier chip for AI Simulation and Control is enabling AV developers to enhance their simulation and drive sensor fusion more efficiently.  
Tuesday Mar 20 | 10:00 AM - 11:00 AM PDT
- Advancing AI to Build the Most-Trusted Driver [S72644]**  
Diego Angelini, VP, Head of Research, Waymo  
Learn how some of Waymo's latest AI research, written in TensorFlow, is expanding the capabilities of Waymo's autonomous driving stack, and three of the leaders that help Waymo.  
Tuesday Mar 19 | 1:00 PM - 2:00 PM PDT
- AI's Next Frontier: Taking Autonomy from Digital to Physical Reality [S73038]**  
Robert Urman, Co-Founder and CEO, Wayve  
We'll explore the exciting potential of generative AI beyond the digital realm, driving into its real-world applications and the unique challenges innovators face in deploying safe, scalable, AI-powered autonomous vehicles.  
Tuesday Mar 19 | 1:00 PM - 2:00 PM PDT
- Crafting a New Automotive Experience With End-to-End Embodied AI [S72691]**  
Alex Kendall, Co-Founder and CEO, Wayve  
Knight Rider, The Batmobile, Iron Man, Star Wars... the list of automotive vehicles for sequels, but what does the road to autonomy look like? Alex Kendall, co-founder and CEO of



# Grace Blackwell in Full Production



CoreWeave

Lenovo

Microsoft Azure

DELL Technologies

aws

Meta

Google Cloud

ORACLE CLOUD Infrastructure

Hewlett Packard Enterprise

AiRES

ASUS

GIGABYTE™

Ingrasys®

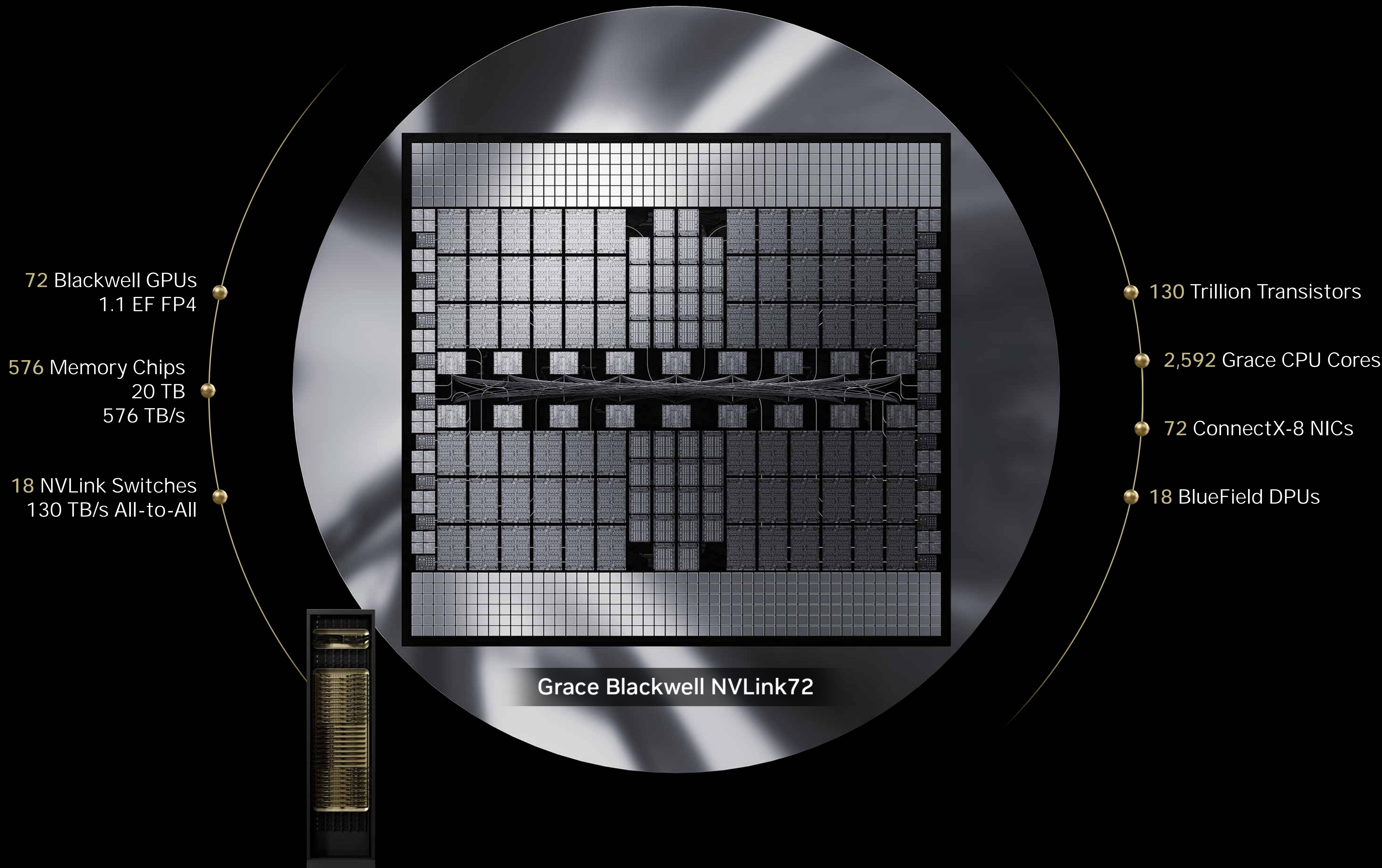
Inventec

PEGATRON

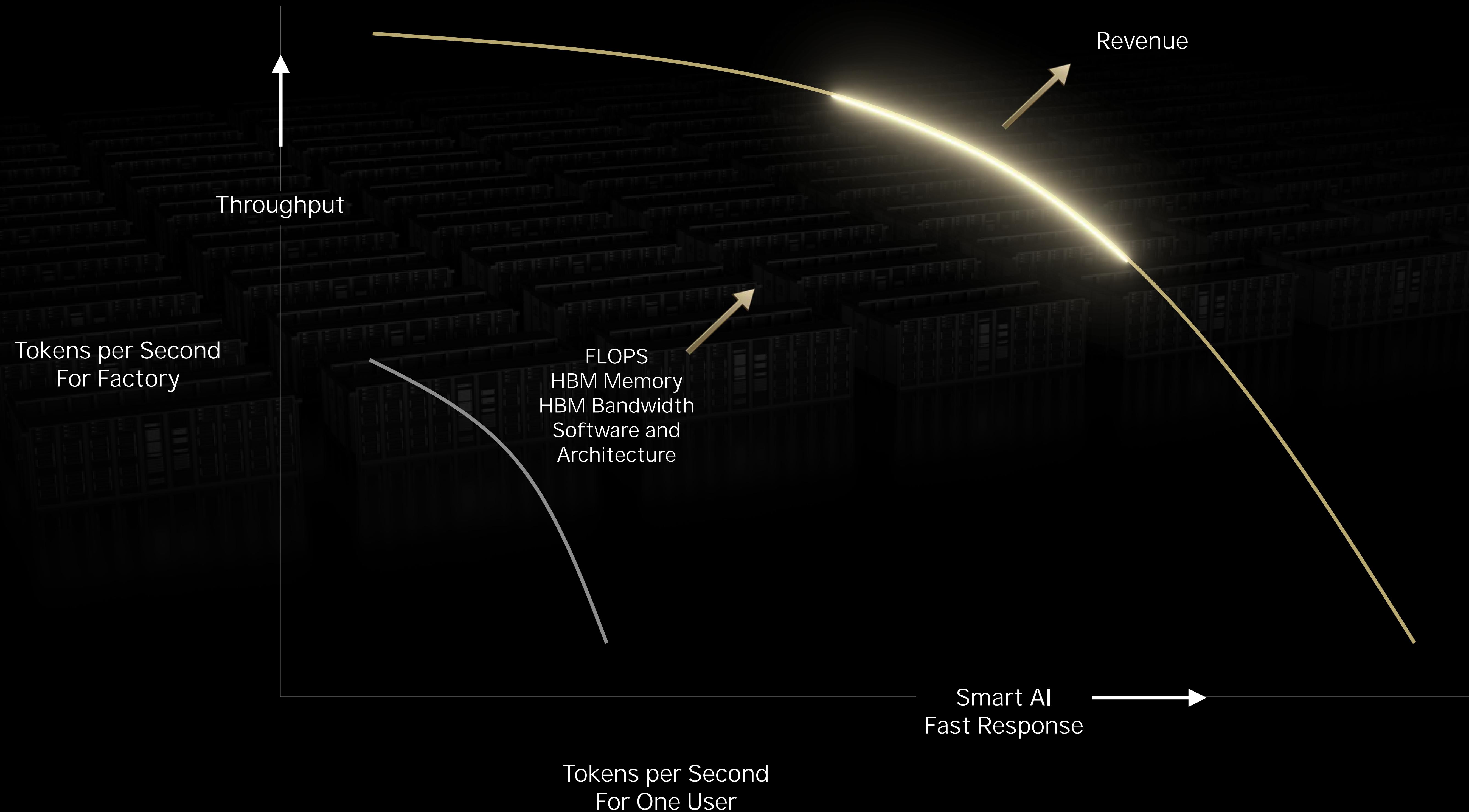
SUPERMICRO

QCT

# NVIDIA Blackwell System



# Inference At-Scale is Extreme Computing

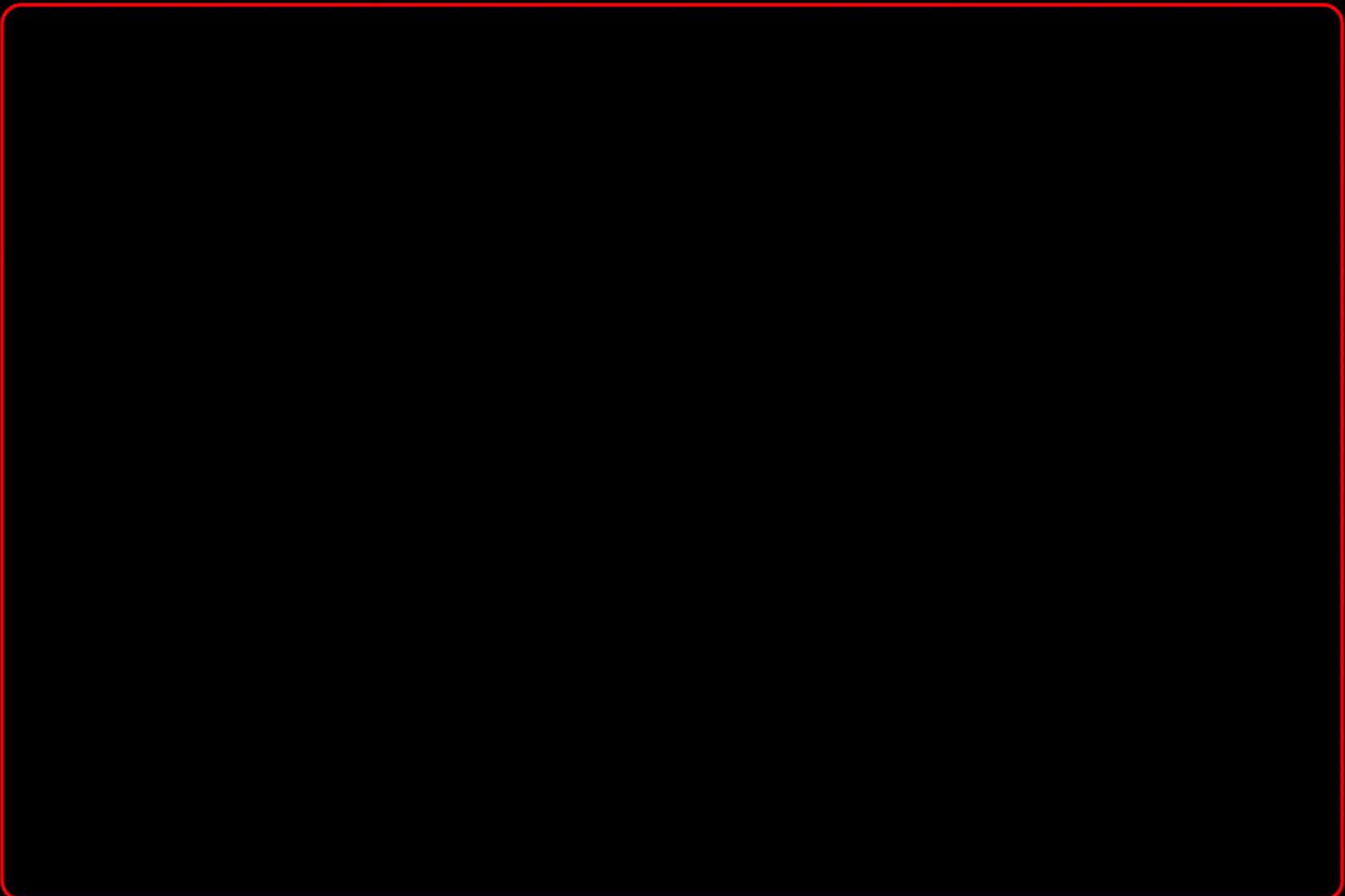


*I need to seat 7 people around a table at my wedding reception, but my parents and in-laws should not sit next to each other. Also, my wife insists we look better in pictures when she's on my left, but I need to sit next to my best man. How do I seat us on a round table? But then, what happens if we invited our pastor to sit with us?*

 Send

## Traditional LLM Model

Tokens: **439**



## Reasoning Model

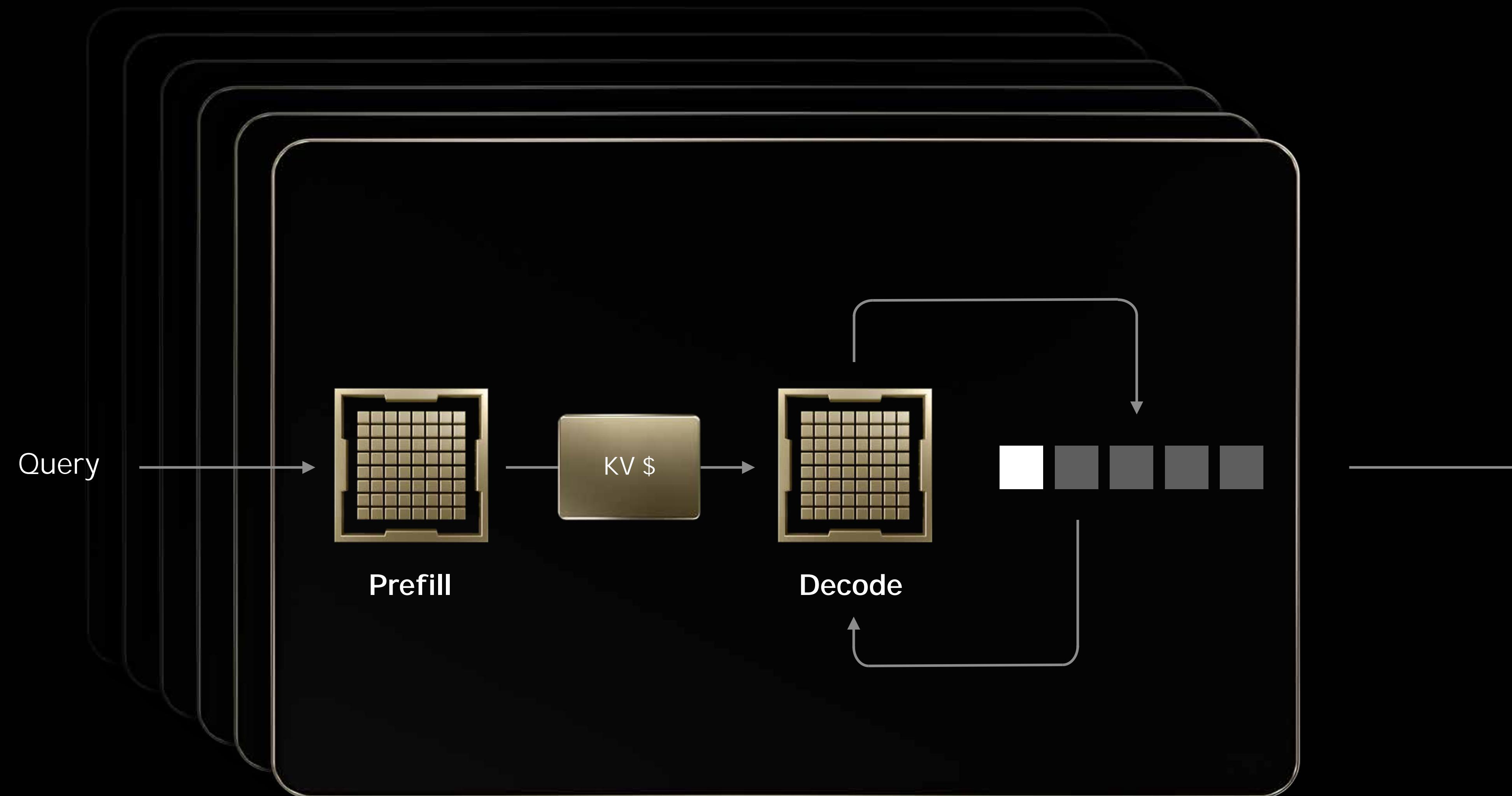
Tokens: **8,559**



**20x More Tokens  
150x More Compute**

# Announcing NVIDIA Dynamo

Distributed Inference Serving Library



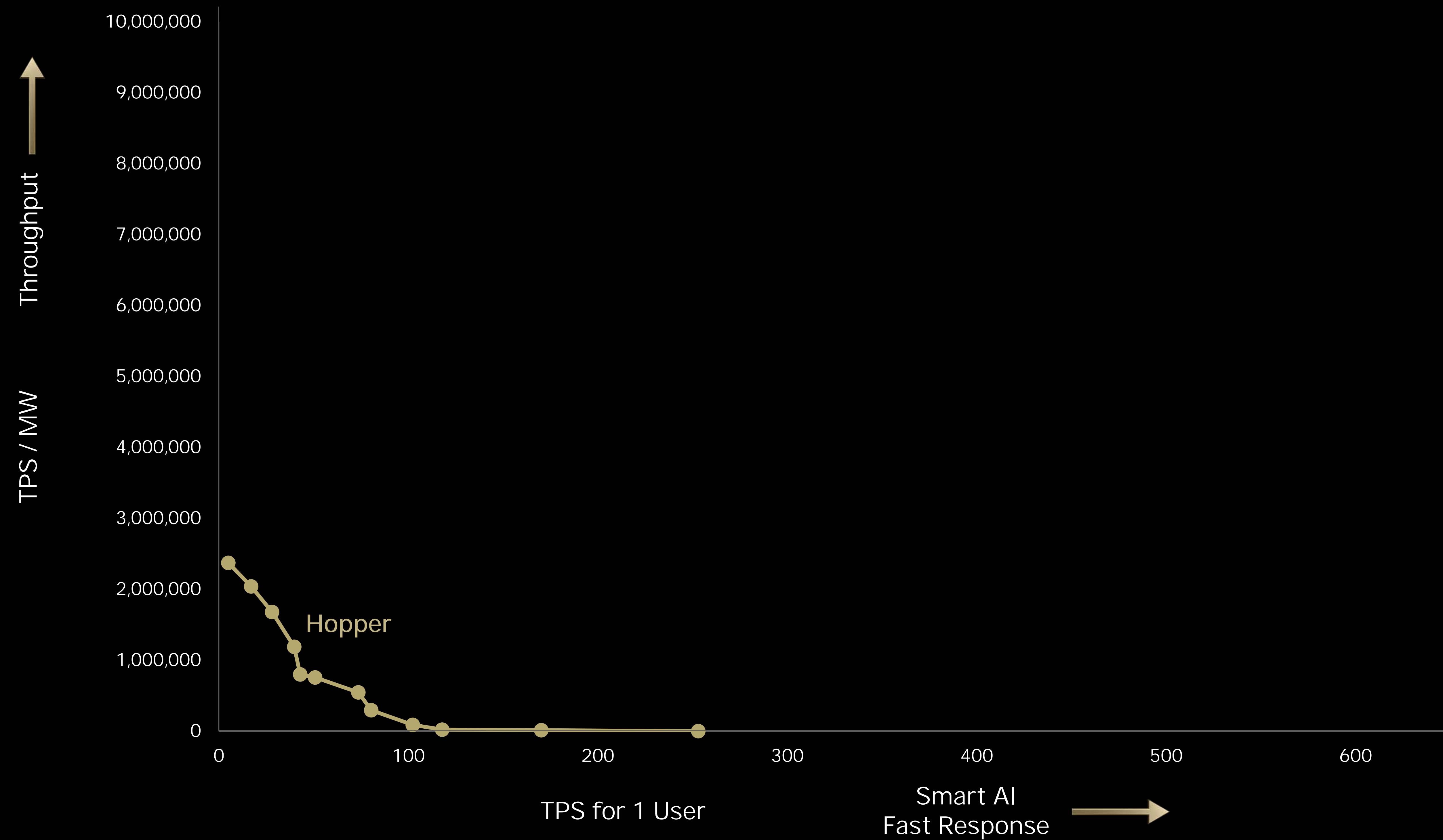
Disaggregated Inference

GPU Resource Allocation

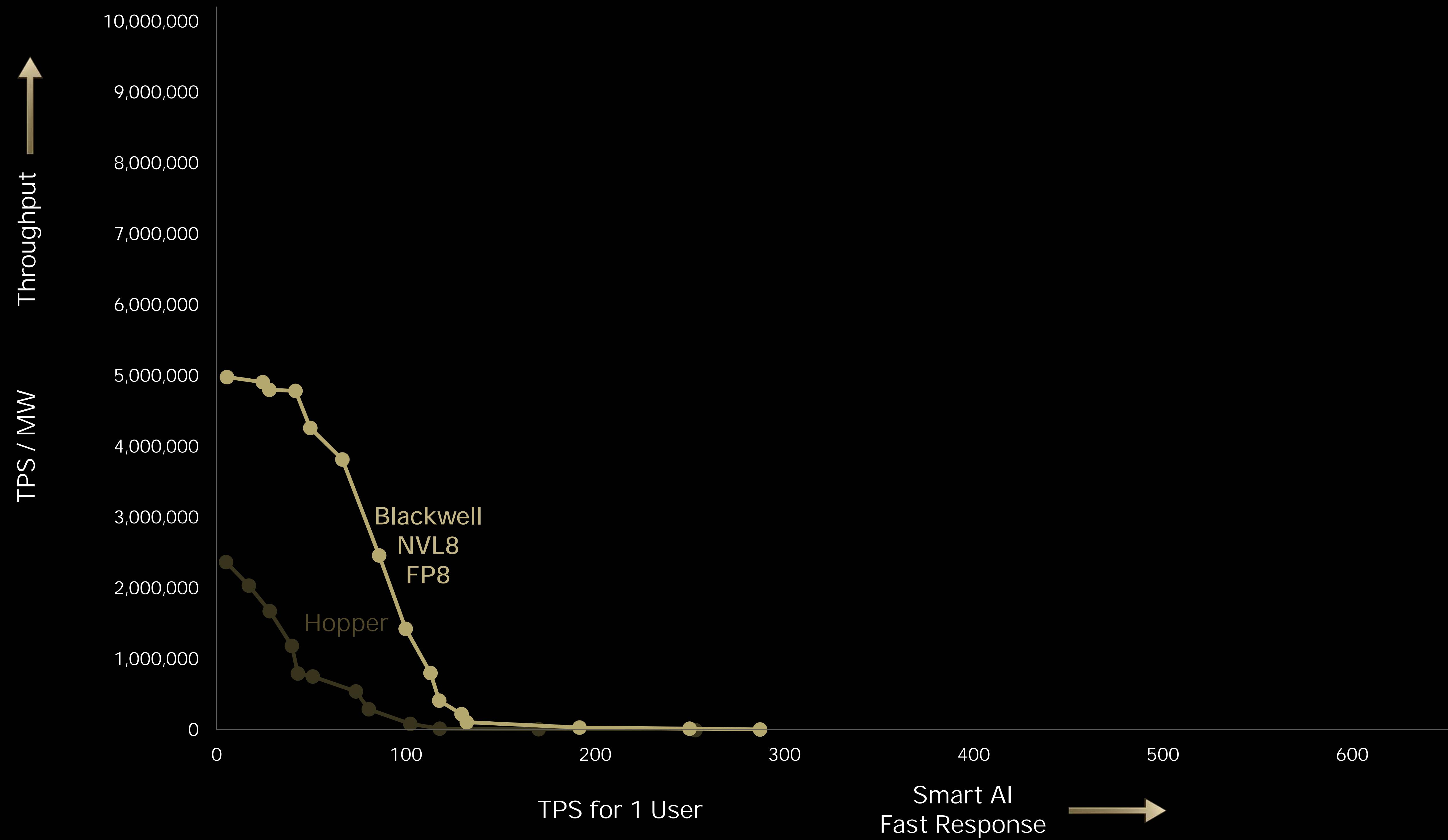
KV Cache Routing

Communication Library (NIXL)

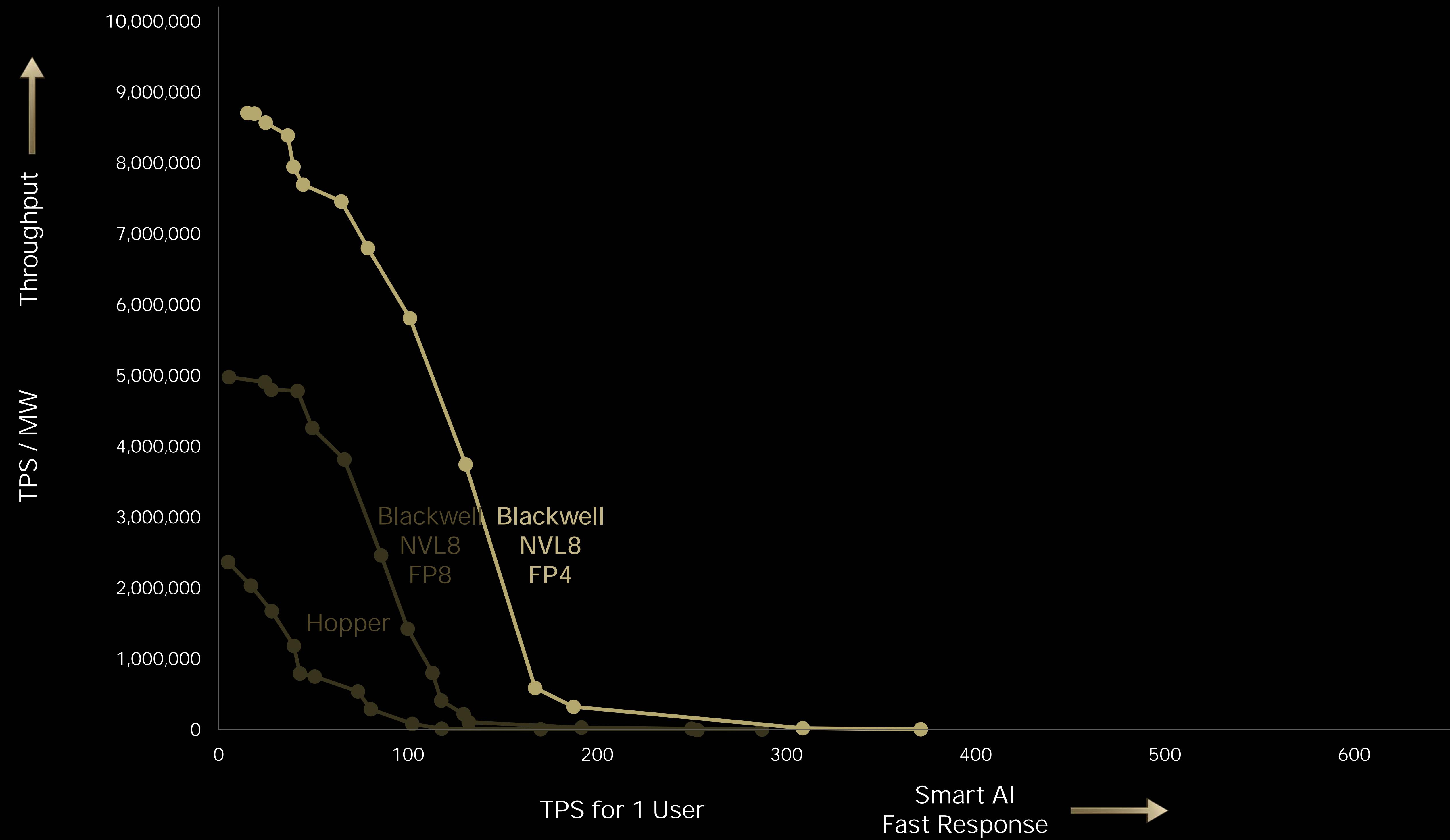
# Blackwell Giant Leap in Inference Performance



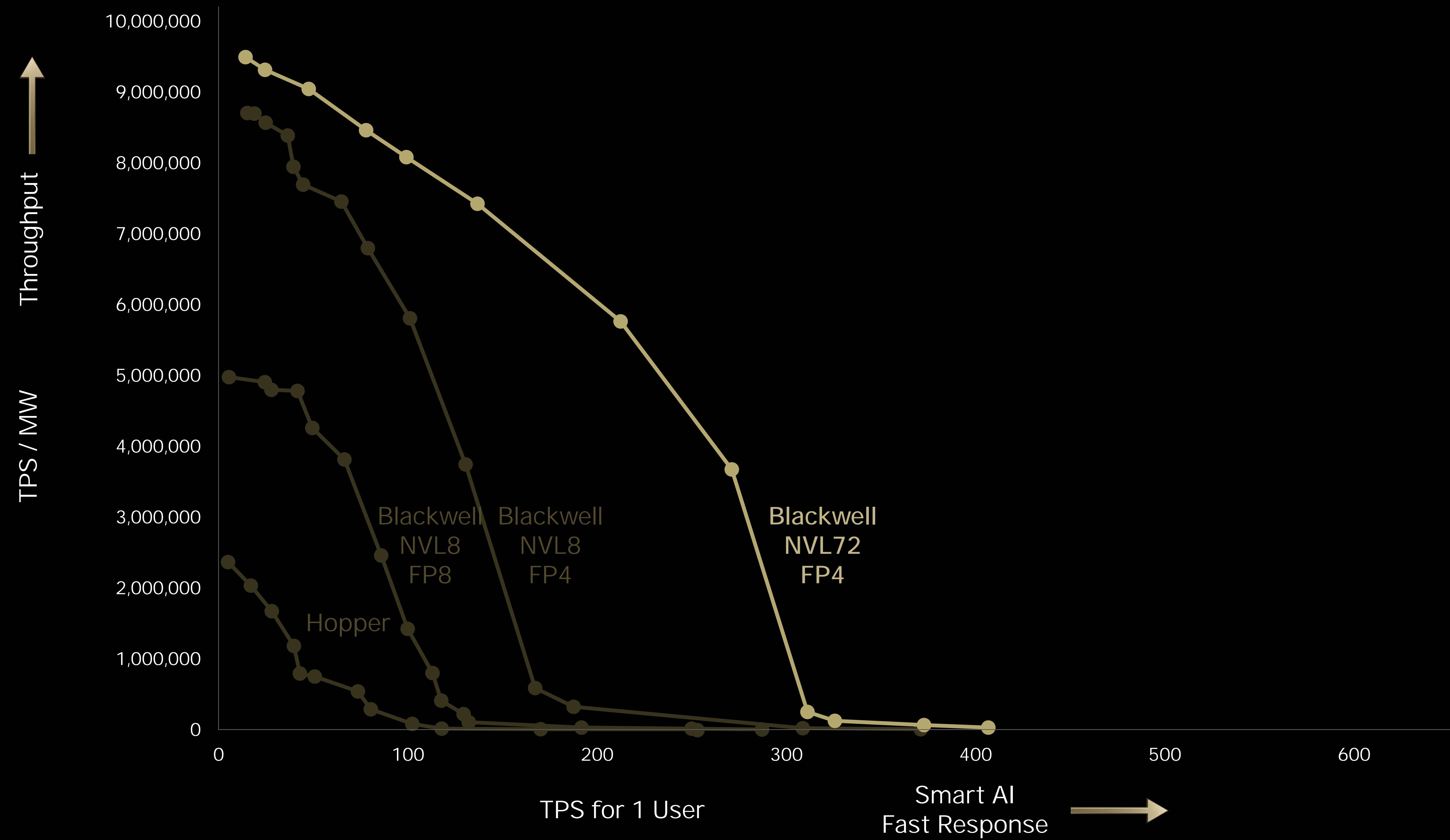
# Blackwell Giant Leap in Inference Performance



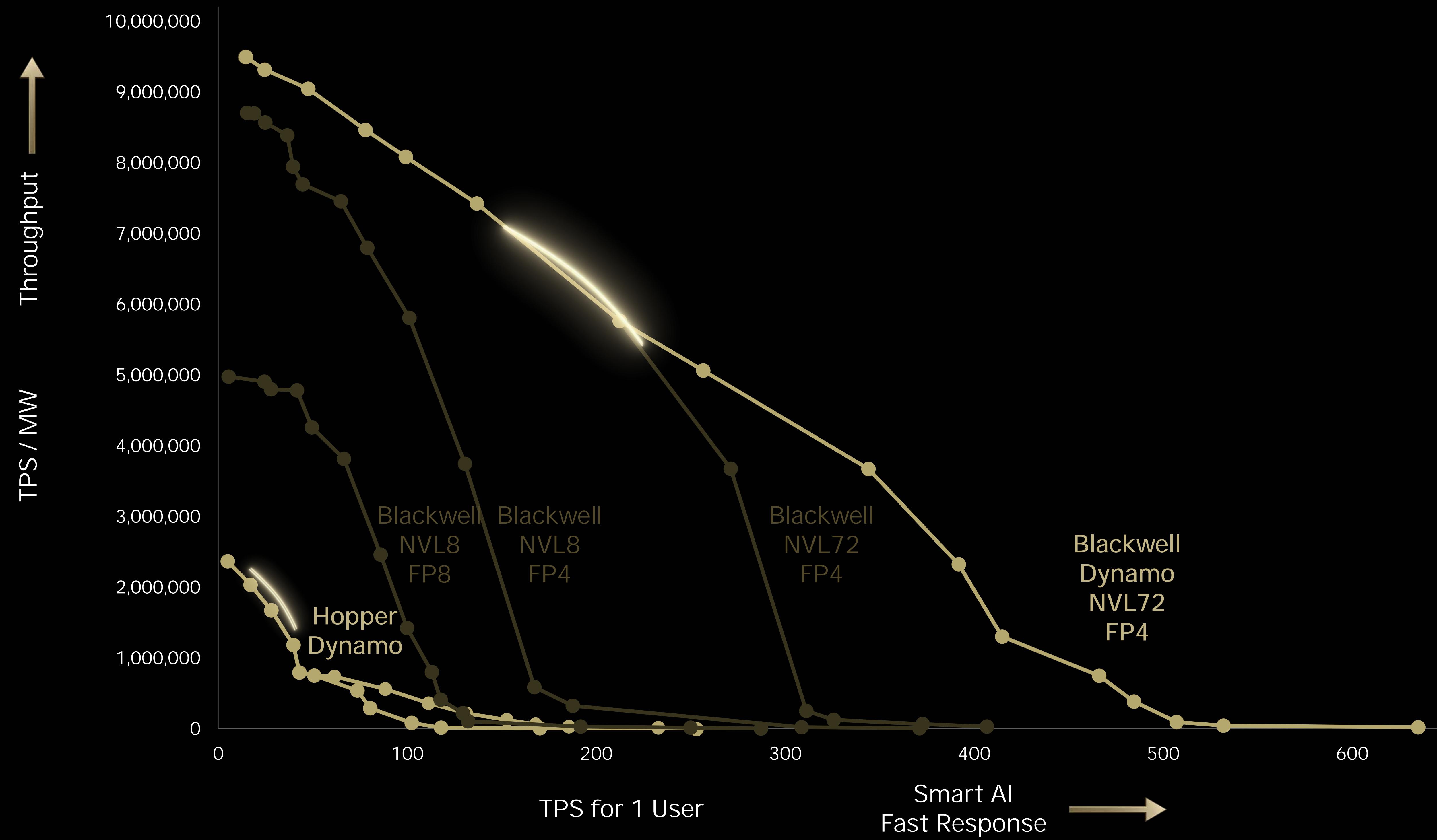
# Blackwell Giant Leap in Inference Performance



# Blackwell Giant Leap in Inference Performance

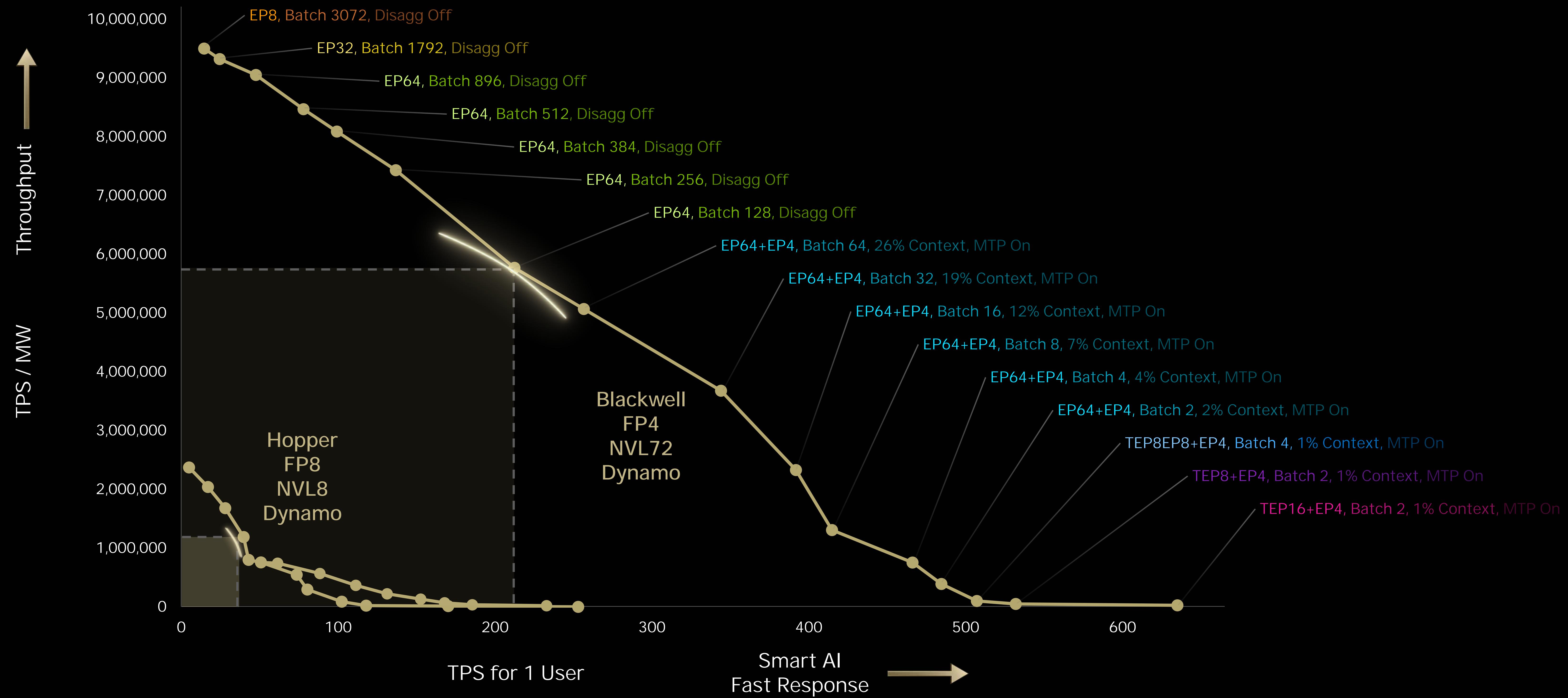


# Blackwell Giant Leap in Inference Performance



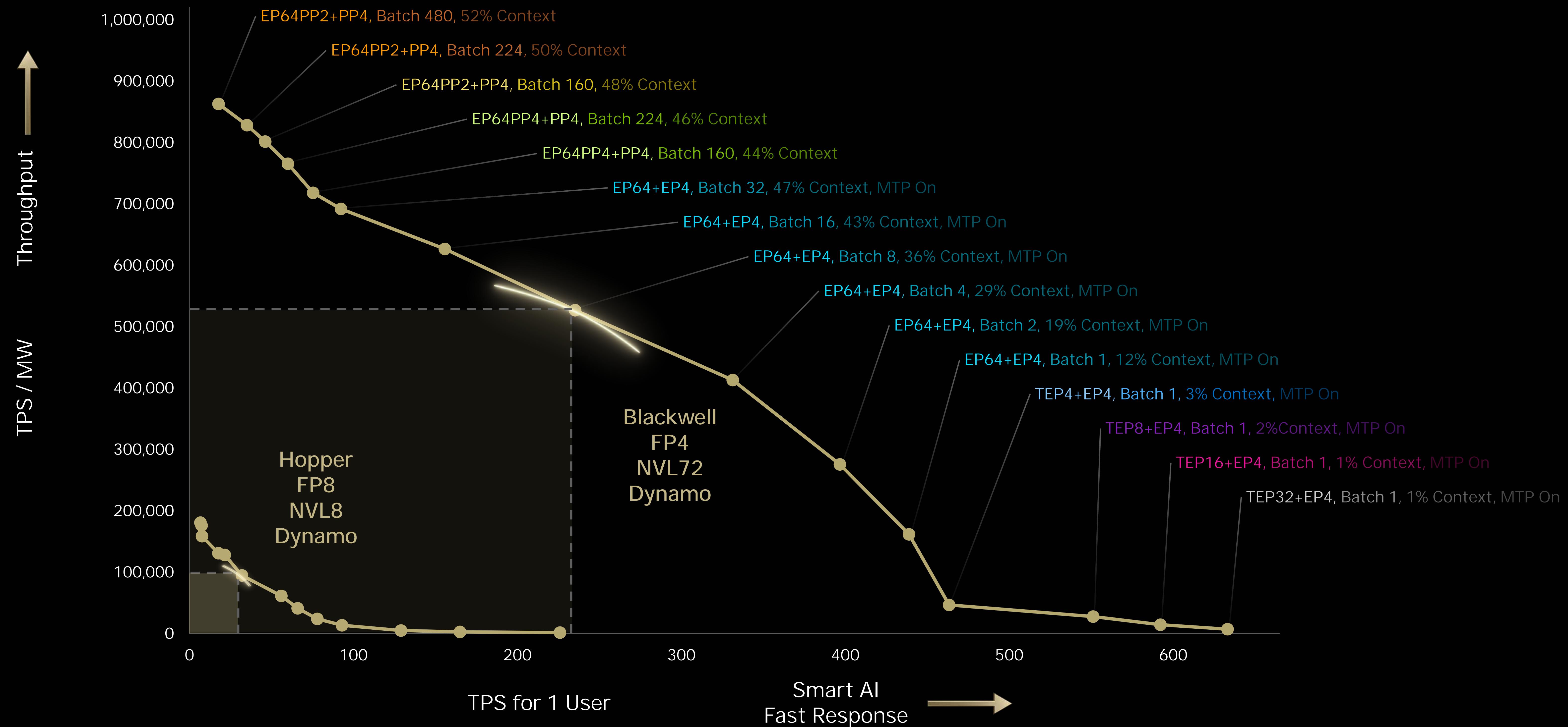
# Blackwell 25X Hopper

FP4, NVL72, Dynamo, and TRT-LLM Continuous Optimization  
1K ISL / 2K OSL



# Blackwell 40X Hopper

FP4, NVL72, Dynamo, and TRT-LLM Continuous Optimization  
32K ISL / 8K OSL





**100 MW AI Factory**

GPU Dies

**H100 NVL8**

45K

Racks

1,400

Data Center Productivity

300M

# Blackwell 40X Hopper Inference Performance

NVLink Flops Per Watt ~ AI Factory Output



100 MW AI Factory

GPU Dies

H100 NVL8

45K

GB200 NVL72

**85K**

Racks

1,400

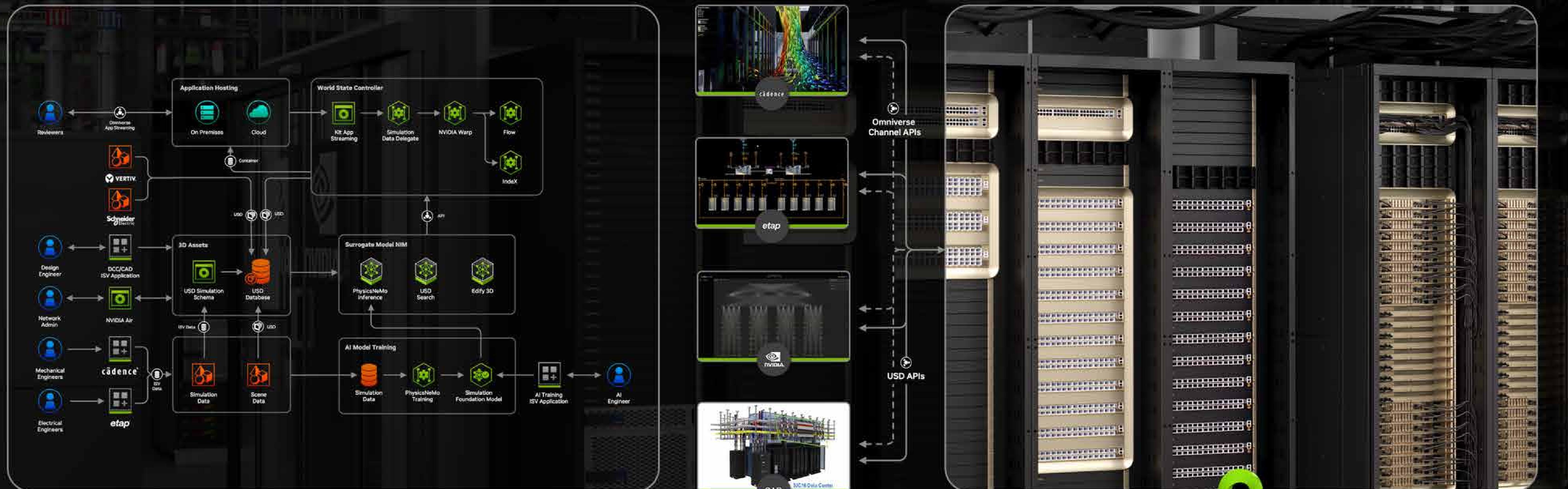
**600**

Data Center Productivity

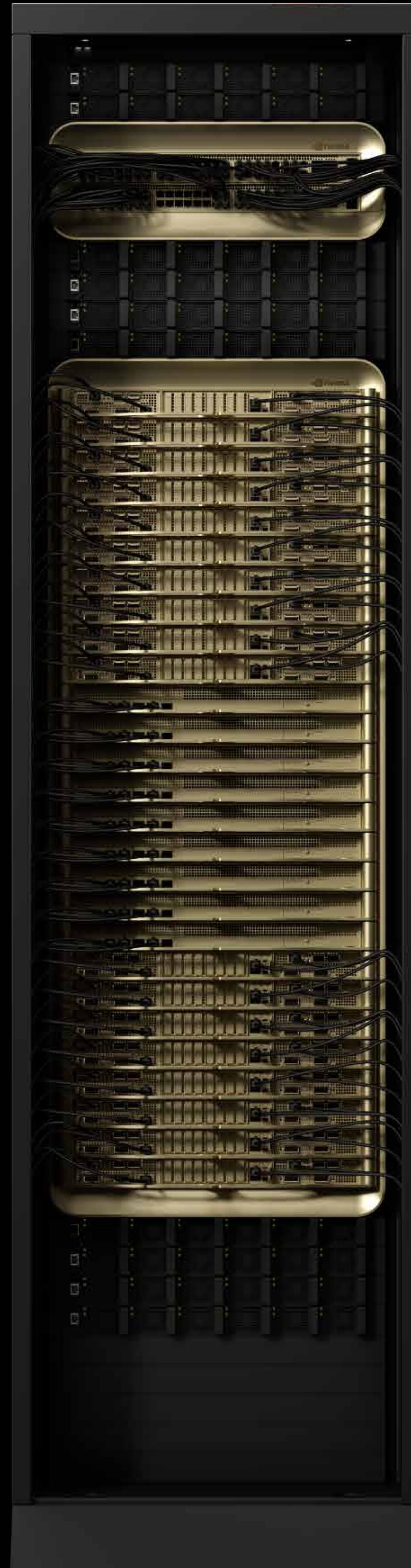
300M

**12,000M**

**40X More Token Revenue**



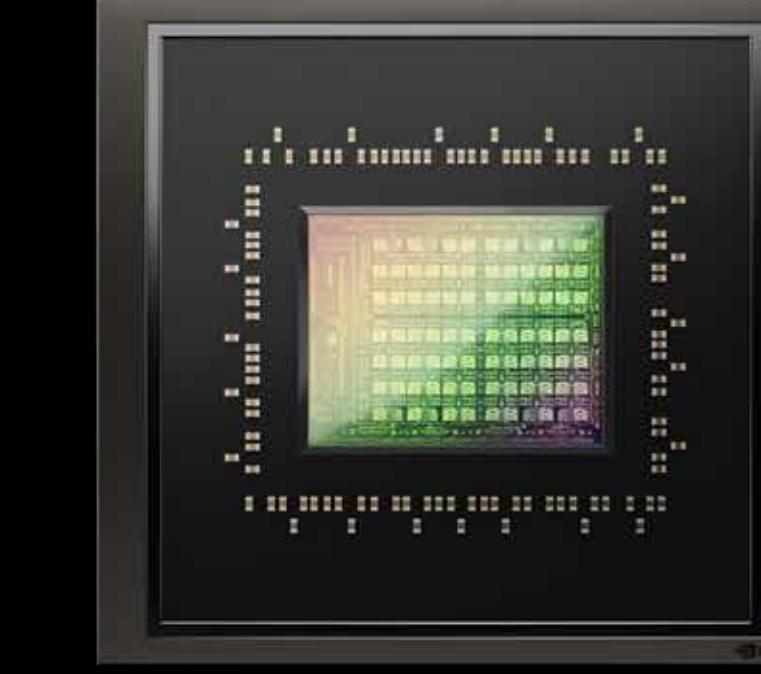




# Blackwell Ultra NVL72

Second Half 2025

Grace



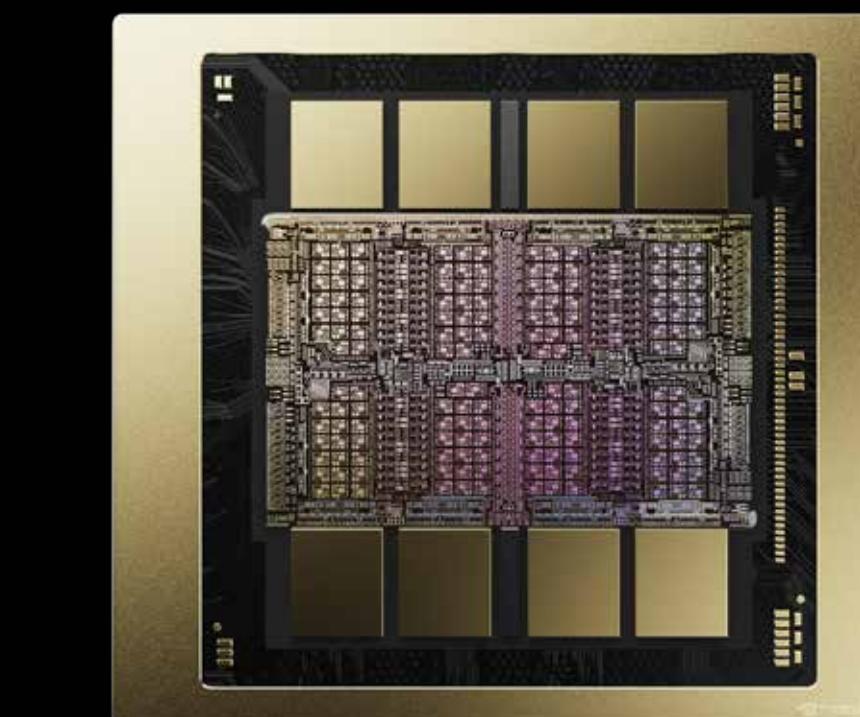
1.1 EF Dense FP4 Inference  
0.36 EF FP8 Training  
**1.5X NVL72**

New Attention Instructions  
**2X**

20 TB HBM | 40 TB Fast Memory  
**1.5X**

14.4 TB/s CX8  
**2X**

Blackwell Ultra



2 Reticle-Sized GPUs  
15PF Dense FP4 | 288GB HBM3e

*Oberon Rack  
Liquid Cooled*



# Vera Rubin NVL144

Second Half 2026

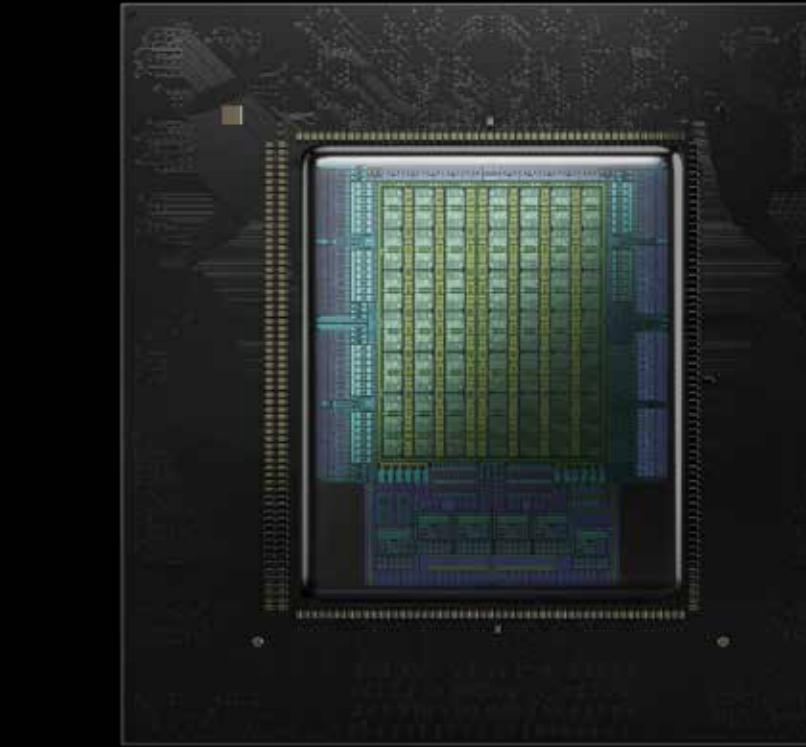
3.6 EF FP4 Inference  
1.2 EF FP8 Training  
**3.3X GB300 NVL72**

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

260 TB/s NVLink6  
**2X**

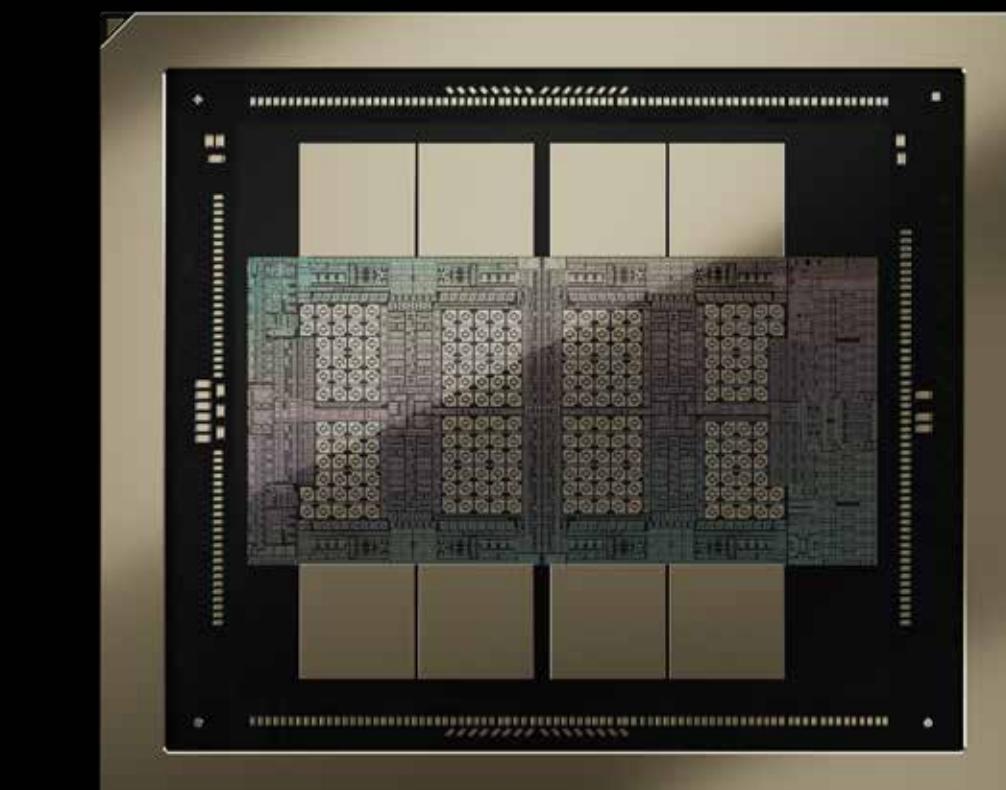
28.8 TB/s CX9  
**2X**

Vera



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

Rubin



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

*Oberon Rack  
Liquid Cooled*



*Kyber Rack*  
*Liquid Cooled*

# Rubin Ultra NVL576

Second Half 2027

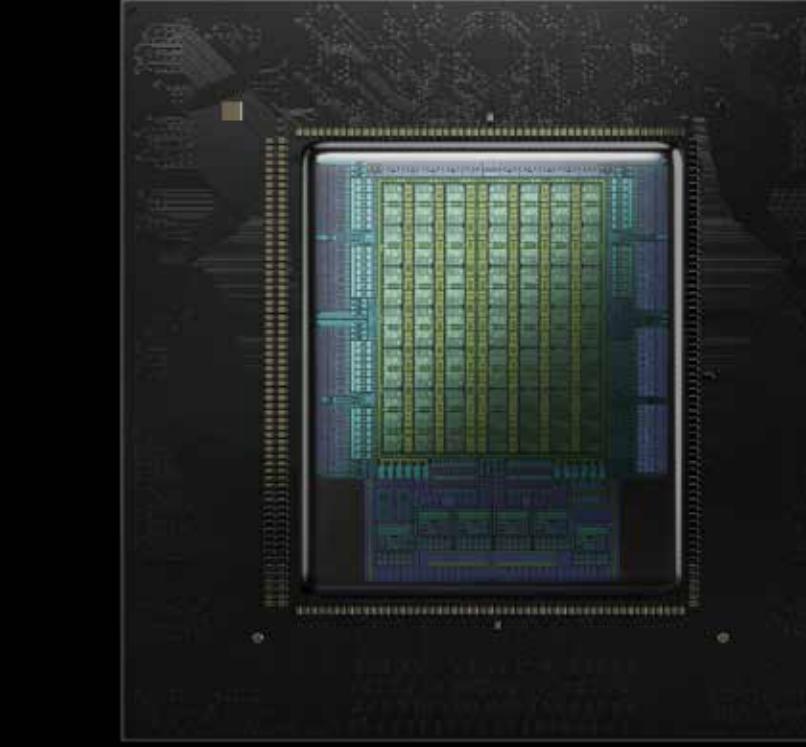
15 EF FP4 Inference  
5 EF FP8 Training  
**14X GB300 NVL72**

4.6 PB/s HBM4e  
365 TB Fast Memory  
**8X**

1.5 PBs NVLink7  
**12X**

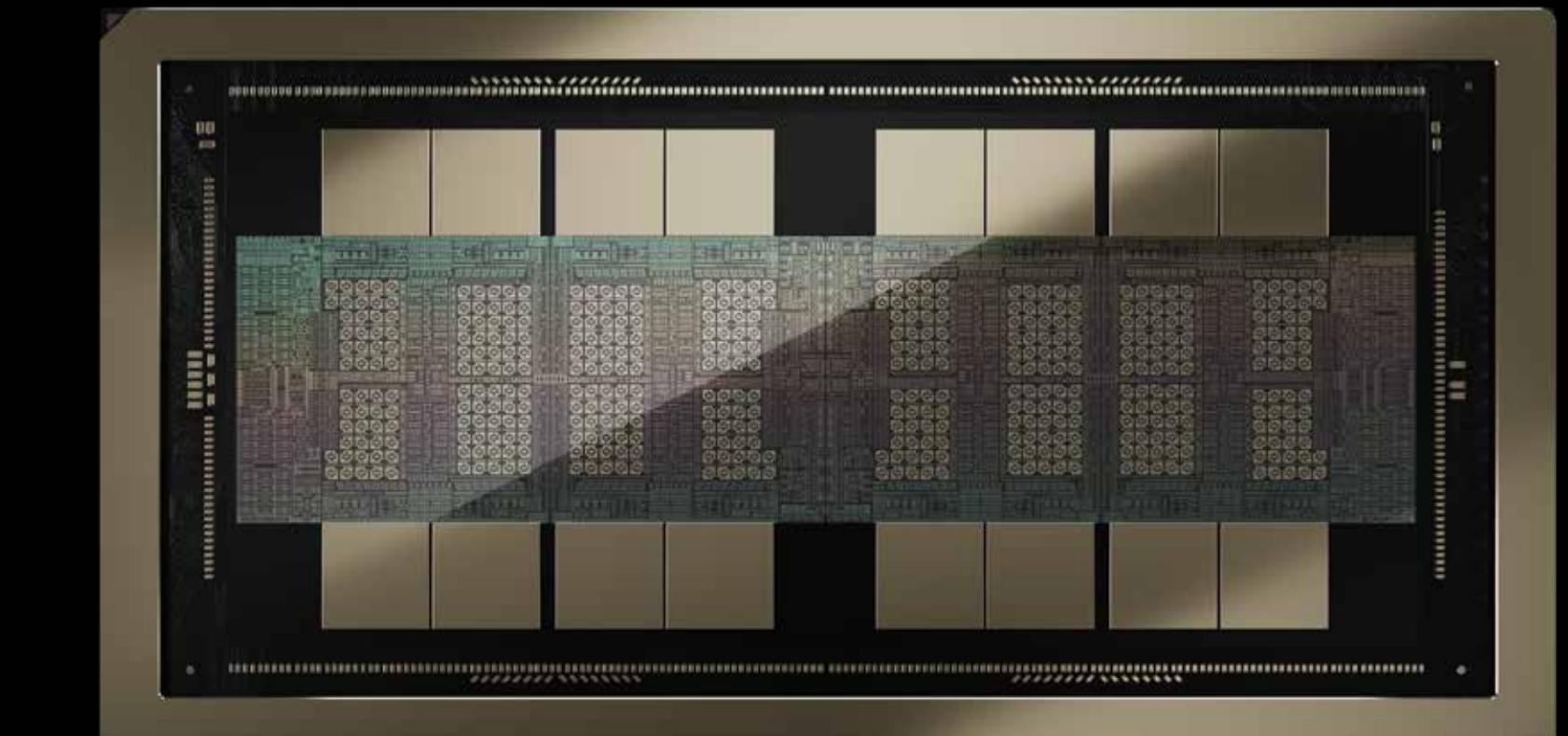
115.2 TB/s CX9  
**8X**

Vera



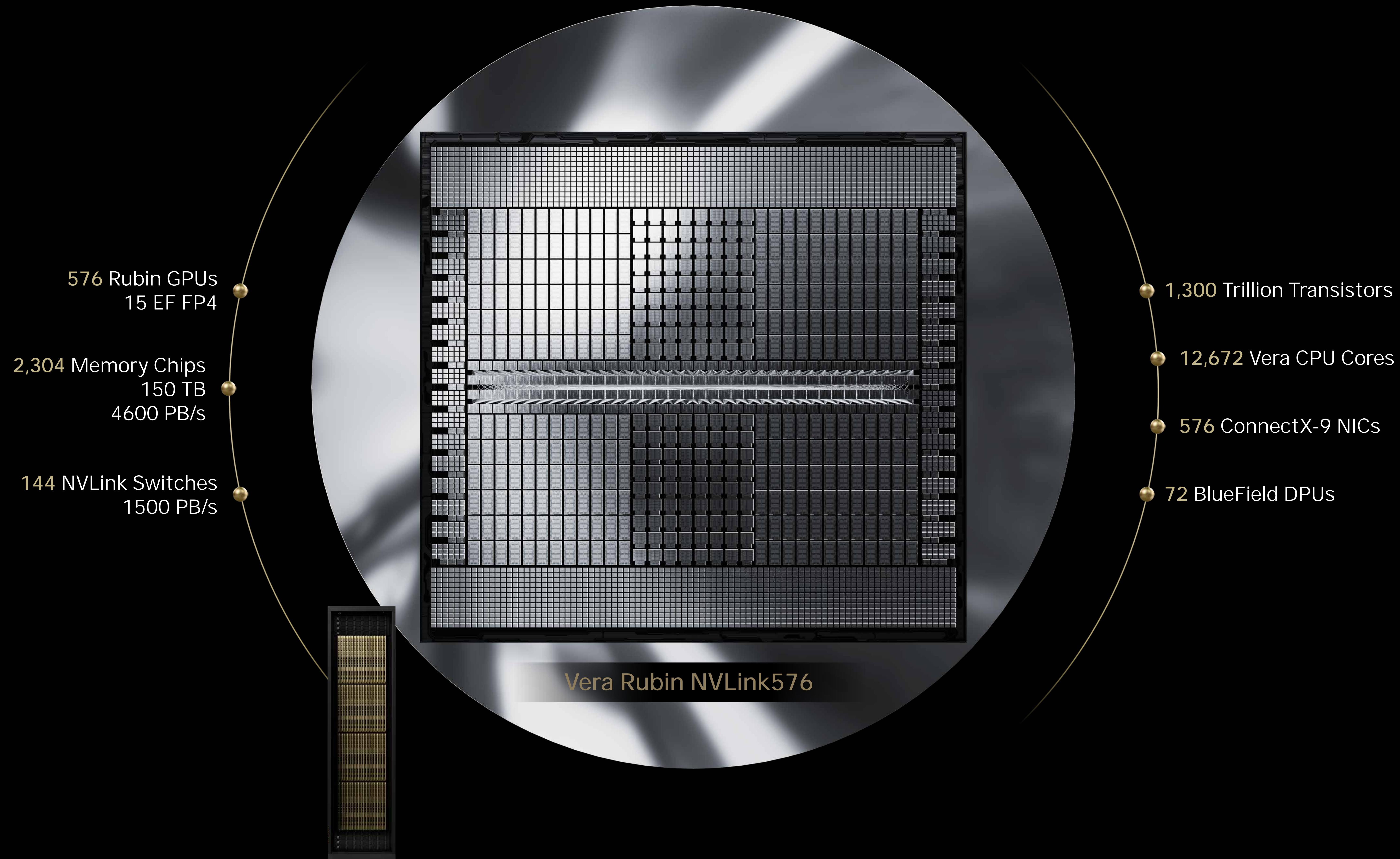
88 Custom Arm Cores  
176 Threads  
1.8 TB/s NVLink-C2C

**Rubin Ultra**



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

# NVIDIA Rubin System



# NVIDIA Blackwell System

72 Blackwell GPUs  
1.1 EF FP4

576 Memory Chips  
20 TB  
576 TB/s

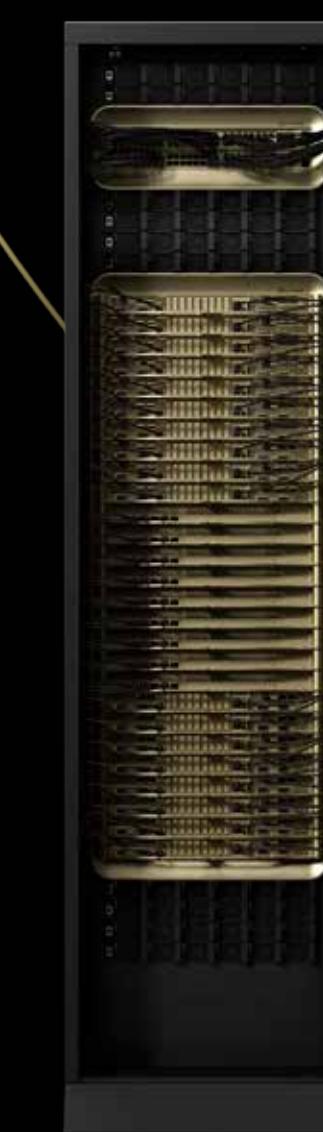
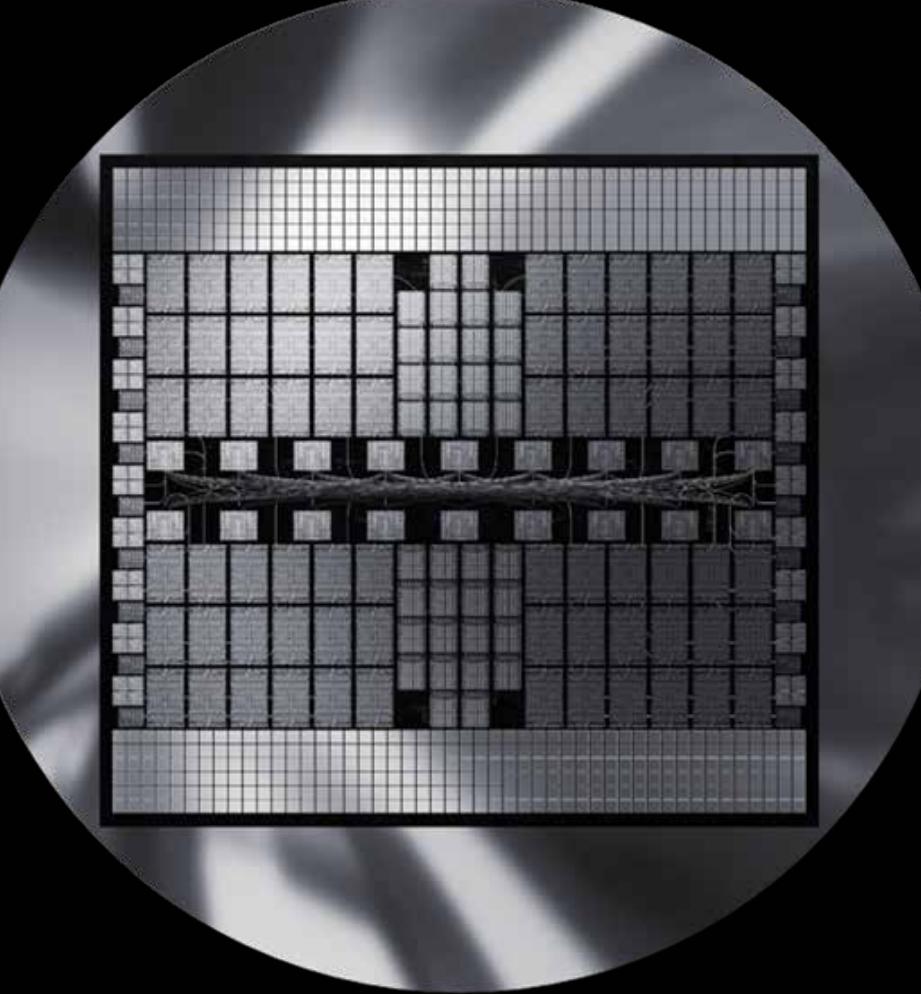
18 NVLink Switches  
130 TB/s

130 Trillion Transistors

2,592 Grace CPU Cores

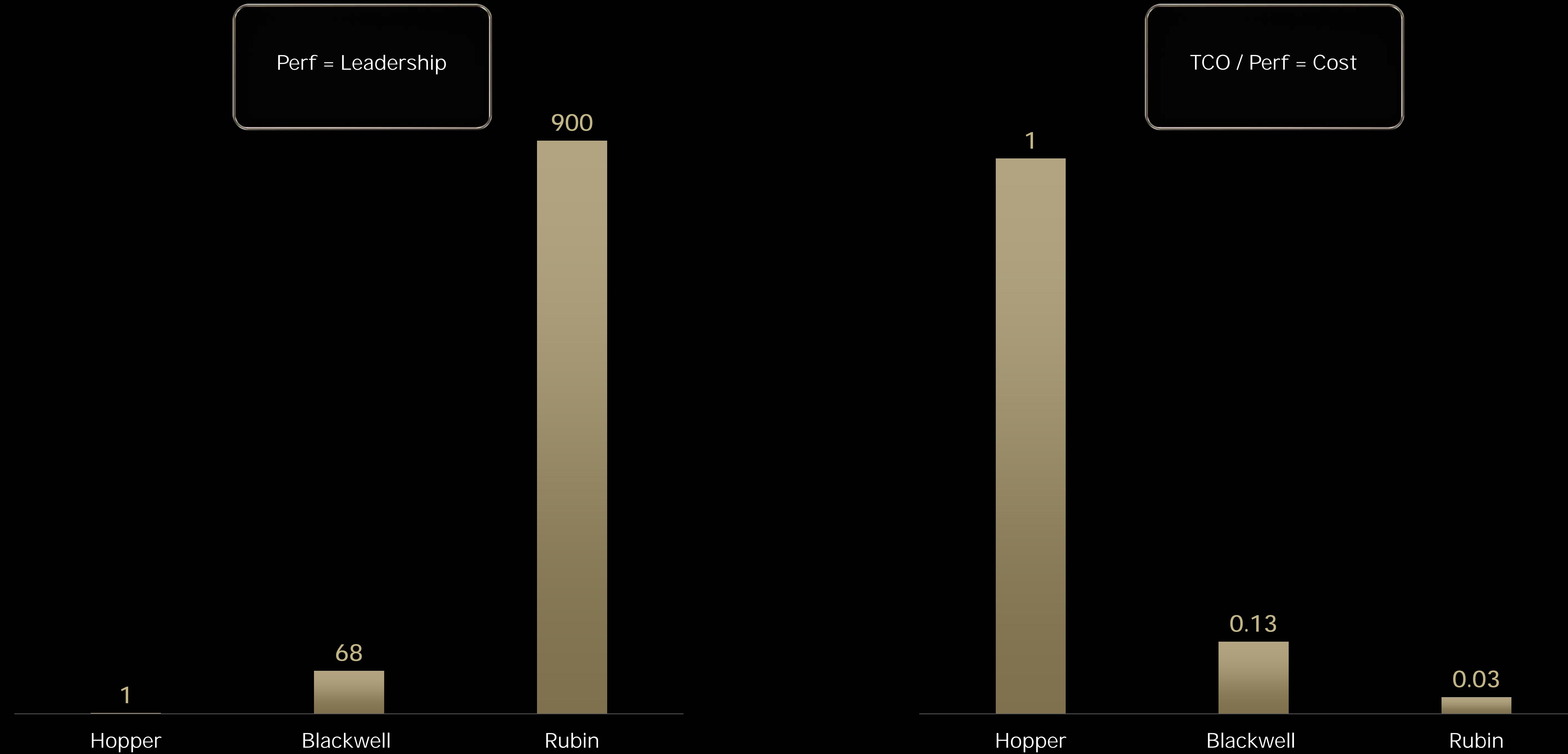
72 ConnectX-8 NICs

18 BlueField DPU



Grace Blackwell NVLink72

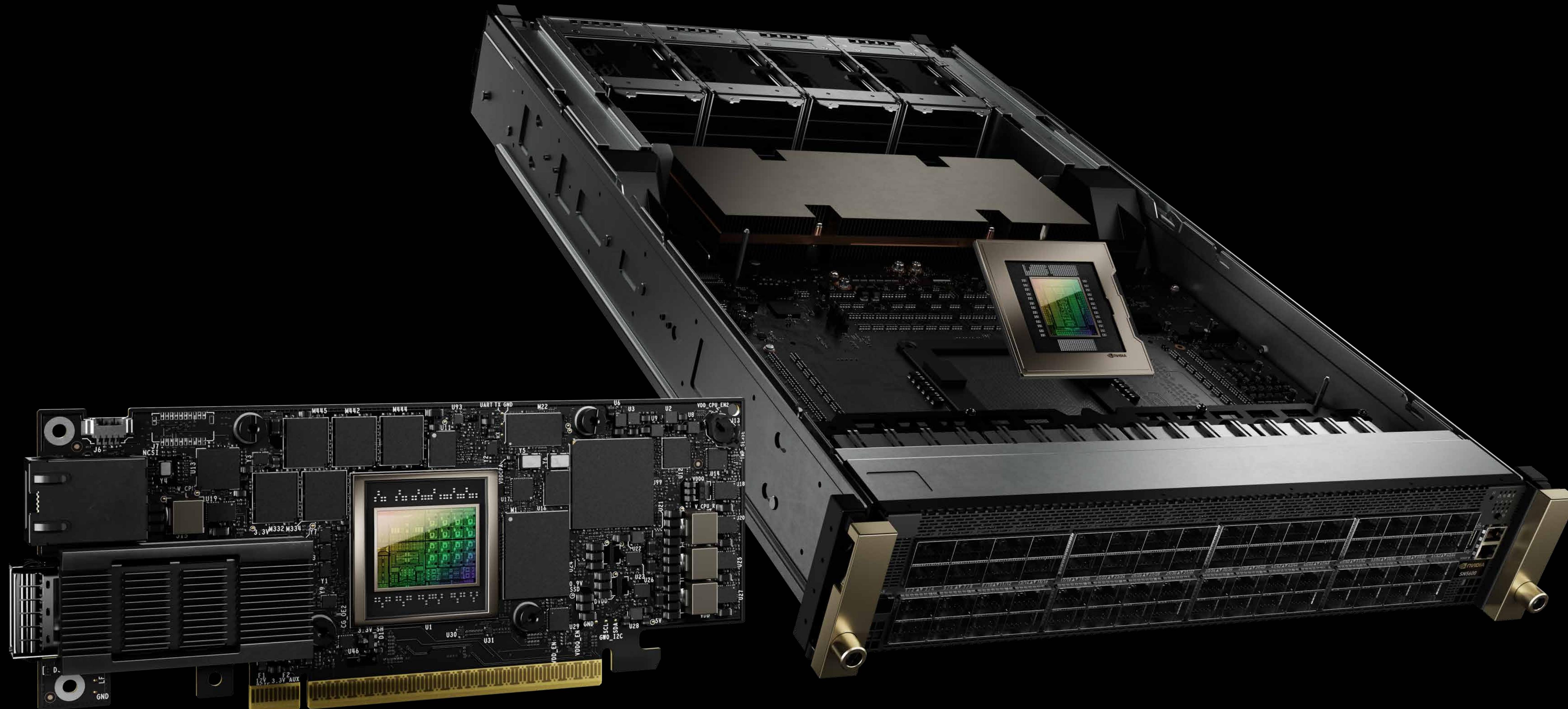
# Fast Roadmap Drives AI Factory Economics



\*Without sparsity

# Spectrum-X “Supercharged” Ethernet

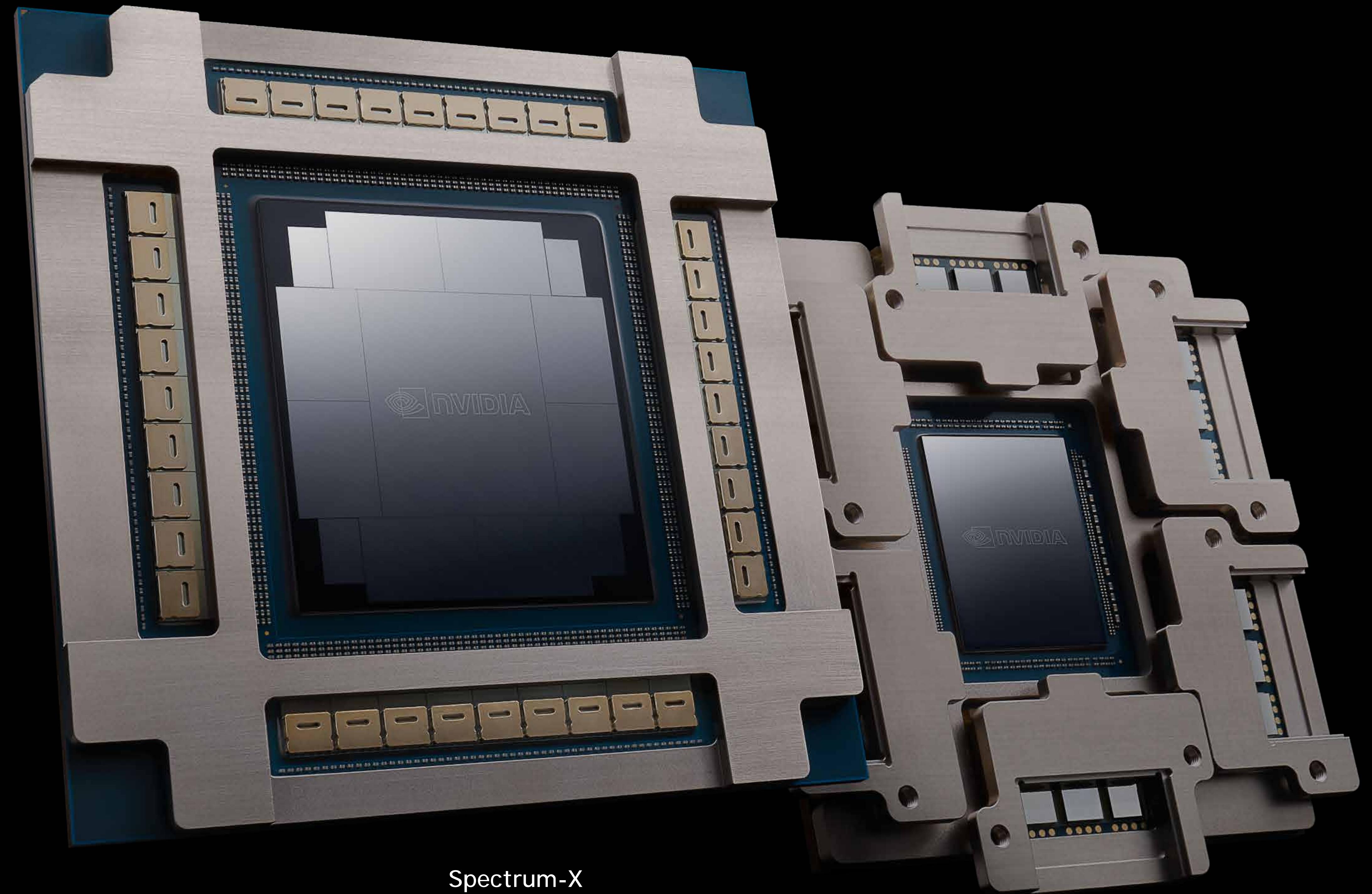
Powering the World’s Most Advanced AI Clouds and Factories



Spectrum-X SuperNIC

Spectrum-X Ethernet Switch

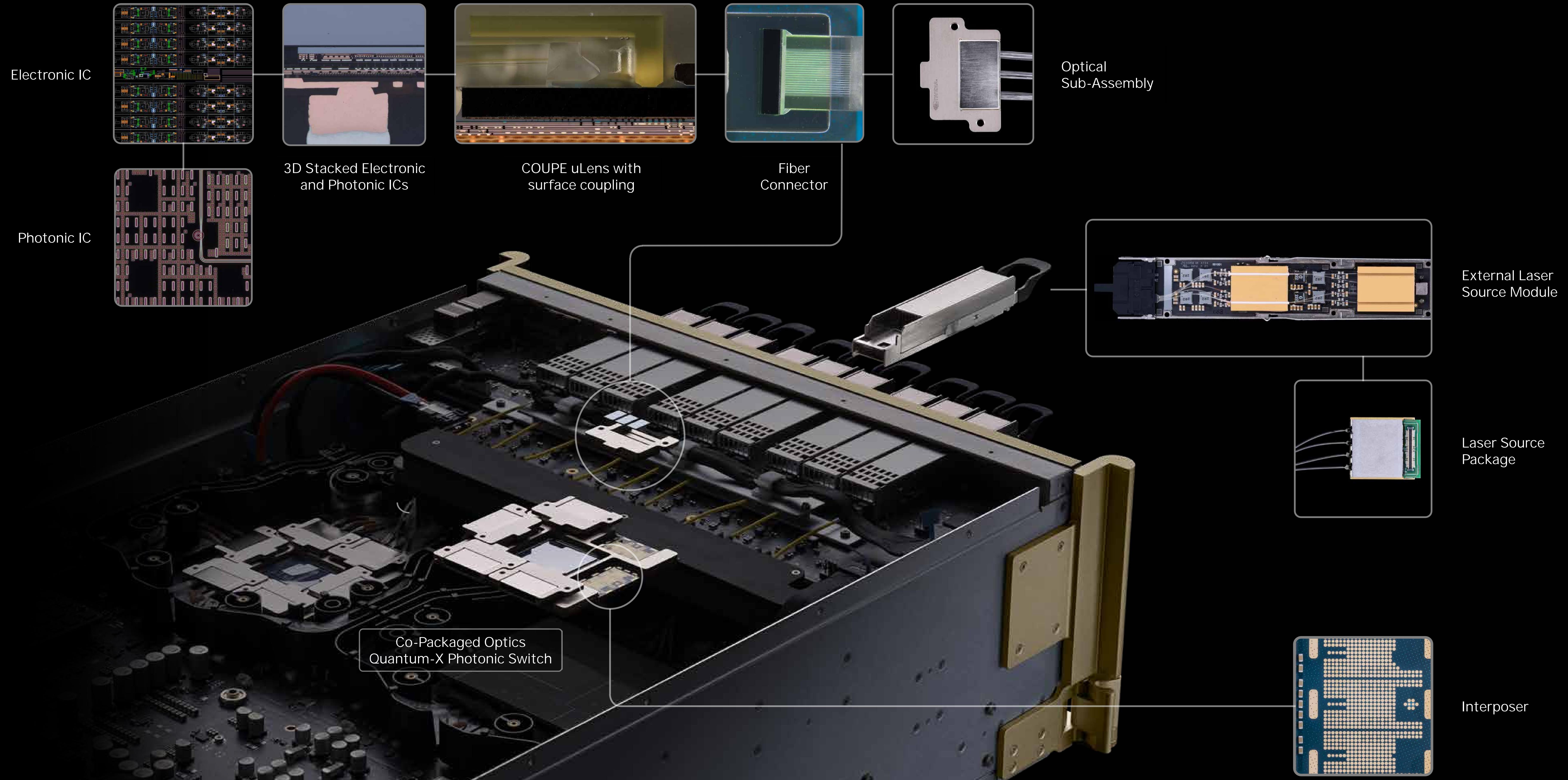


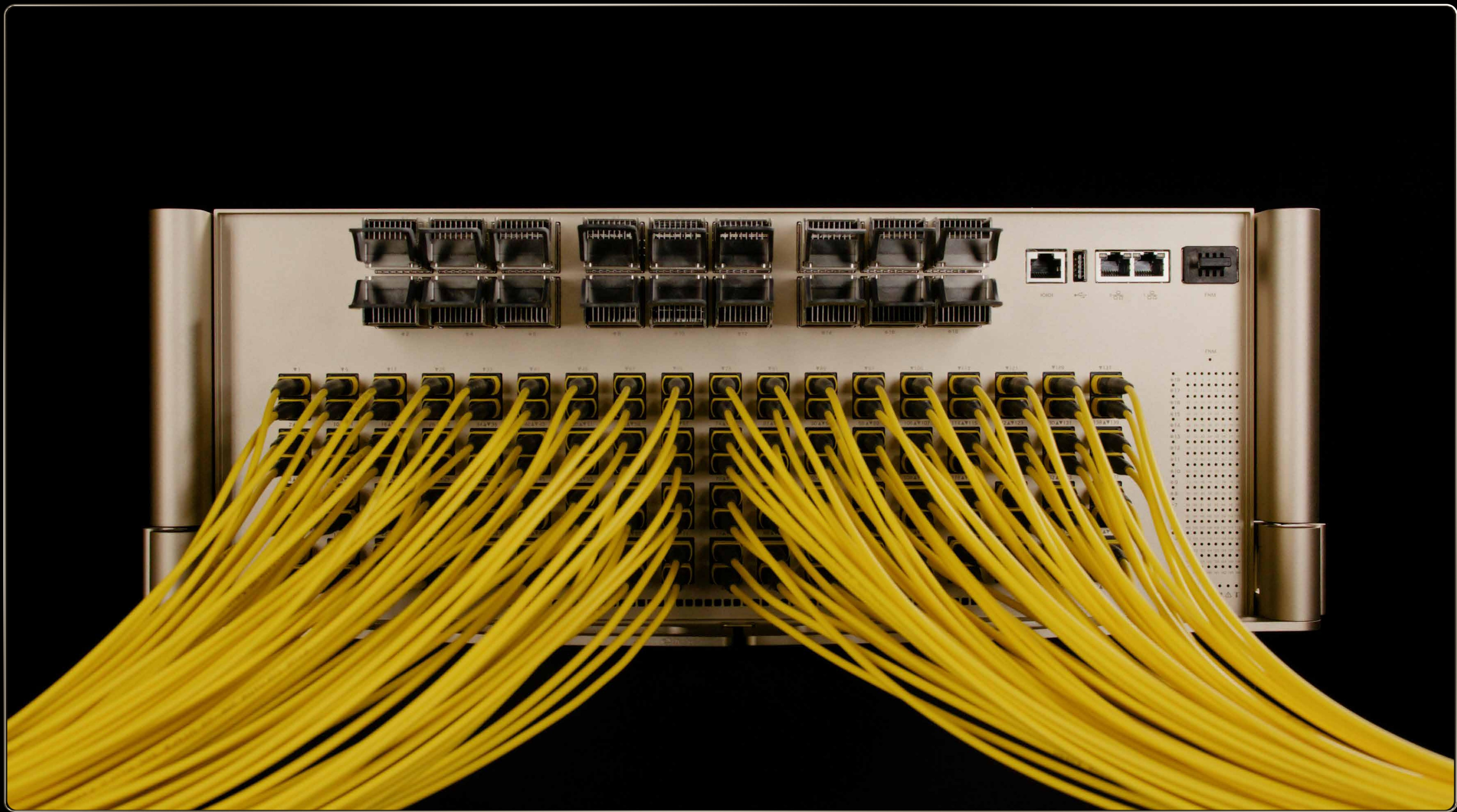


## NVIDIA Photonics

CPO Co-Invention With Ecosystem Partners

- 1<sup>st</sup> 1.6T Silicon Photonics CPO Chip - New Micro Ring Modulators (MRM)
- 1<sup>st</sup> 3D-Stacked Silicon Photonics Engine with TSMC Process
- High-Power, High-Efficiency Lasers
- Detachable Fiber Connectors
- 100's of Patents, Licensed to Partners

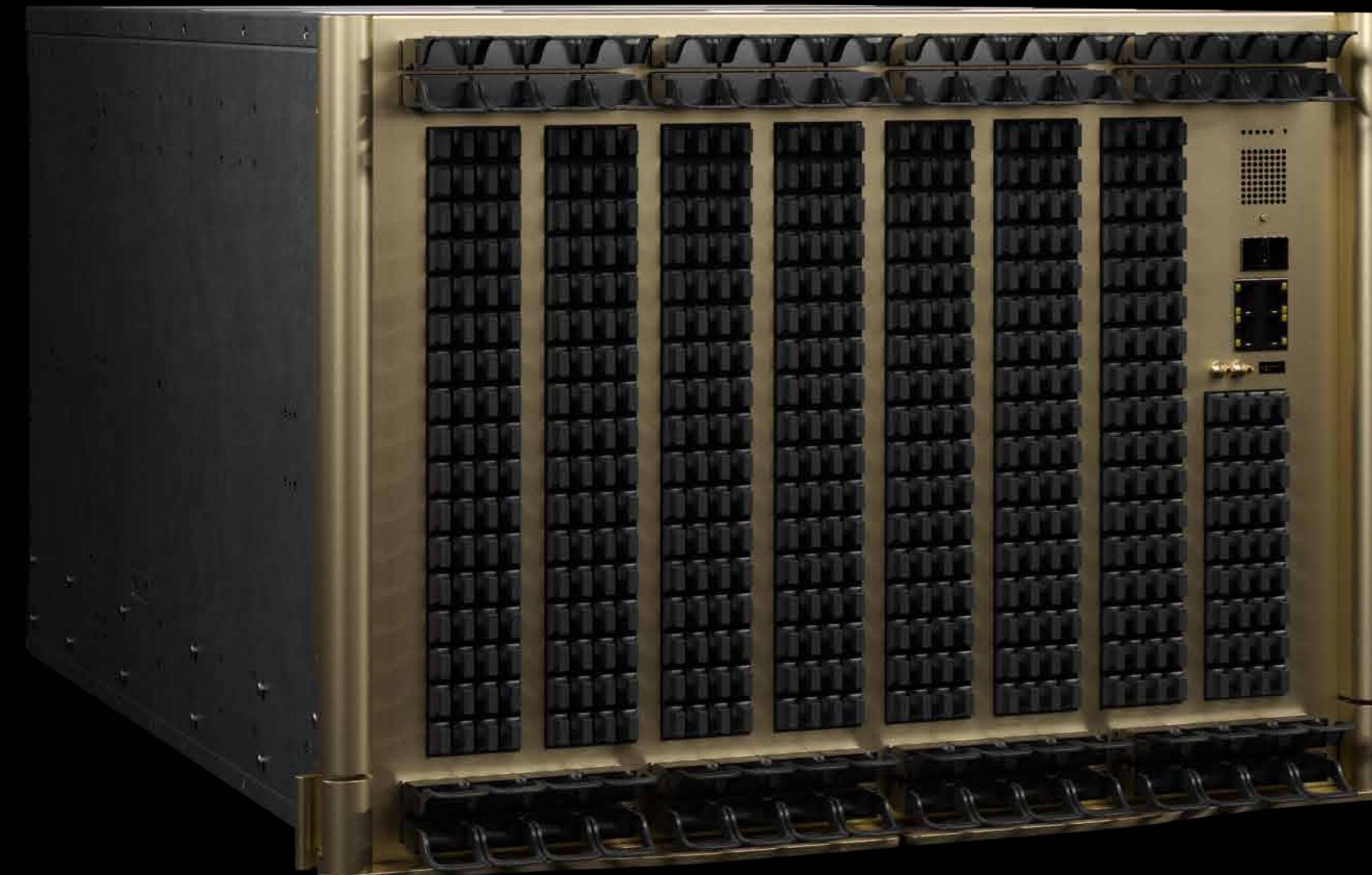




# Announcing NVIDIA Photonics Switch Systems

World's Most Advanced Co-Packaged Optical Switches

**Spectrum-X Photonics**  
2<sup>nd</sup> Half 2026



128 Ports of 800G | 512 x 200G

512 Ports of 800G | 2K x 200G

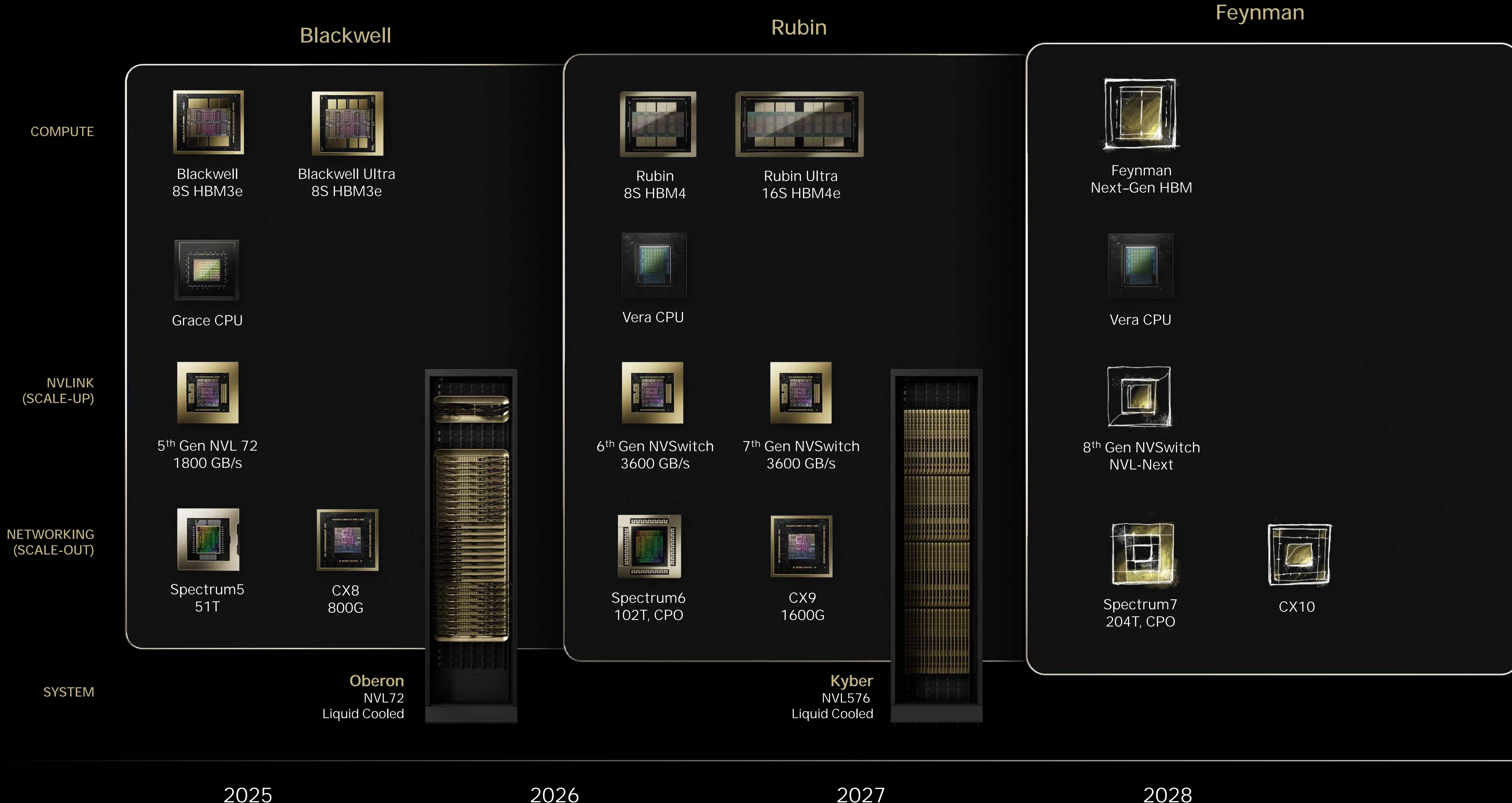
**Quantum-X Photonics**  
2<sup>nd</sup> Half 2025



144 Ports of 800G | 576 x 200G

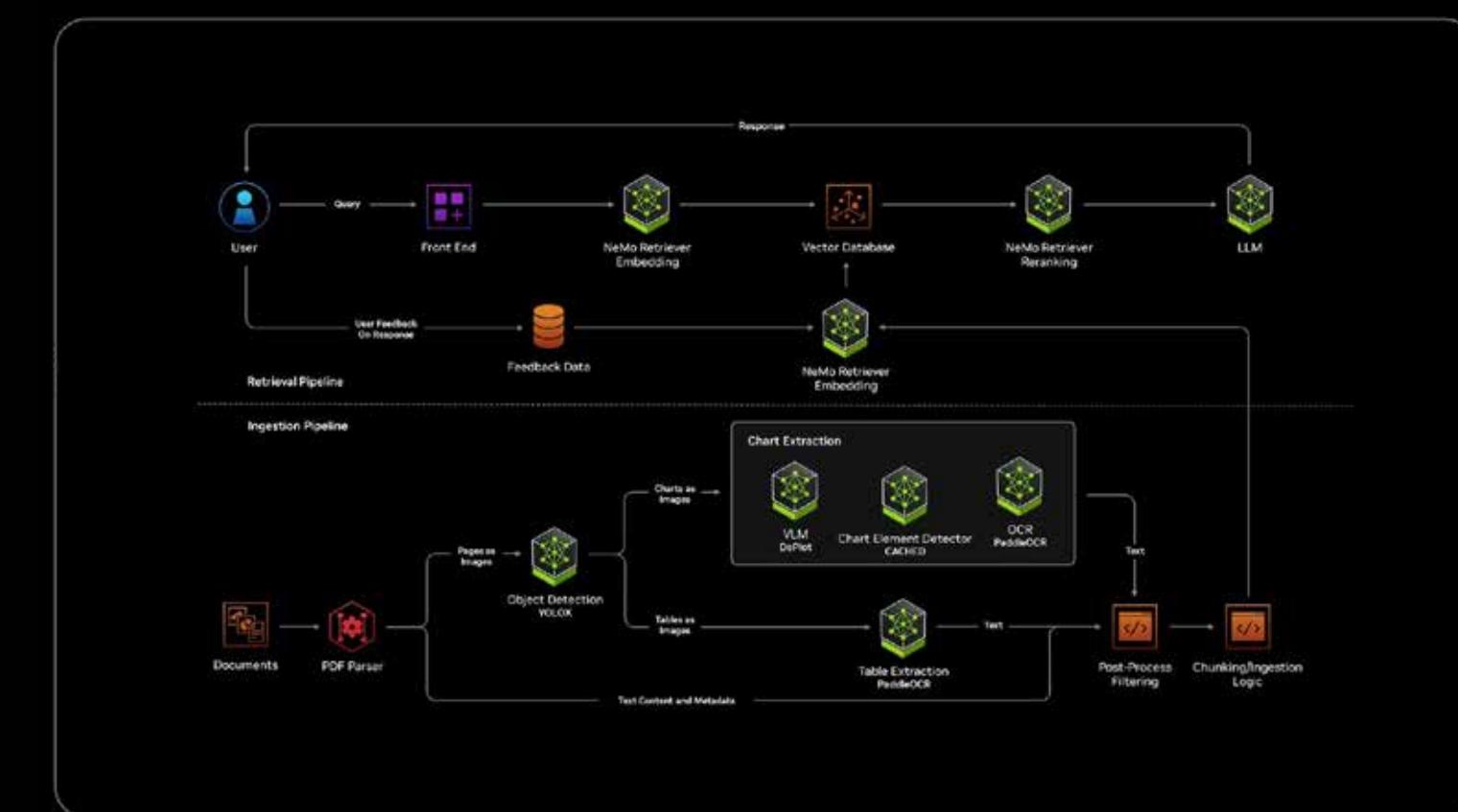
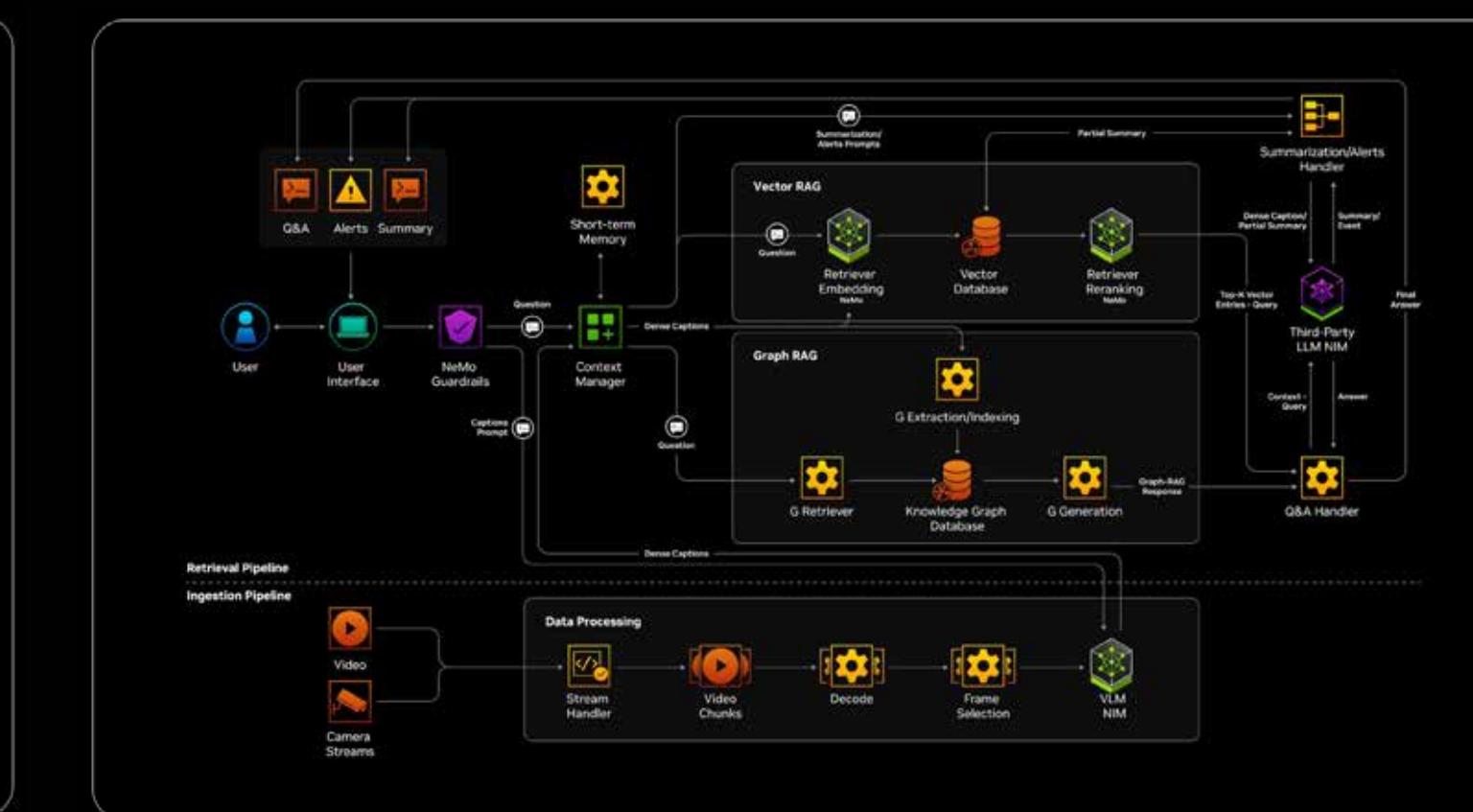
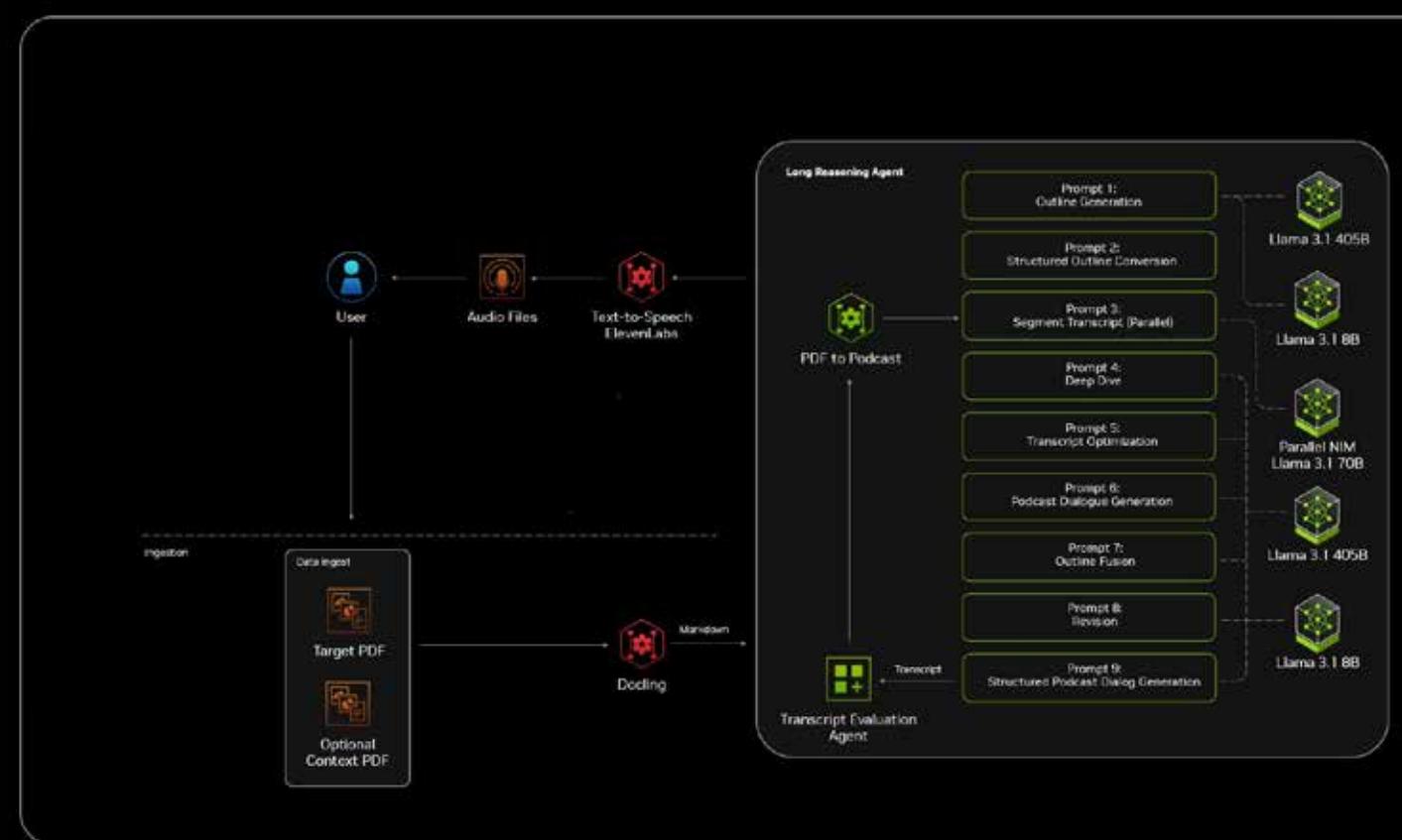
# NVIDIA Paves Road to Gigawatt AI Factories

One-Year Rhythm | Full-Stack | One Architecture | CUDA Everywhere

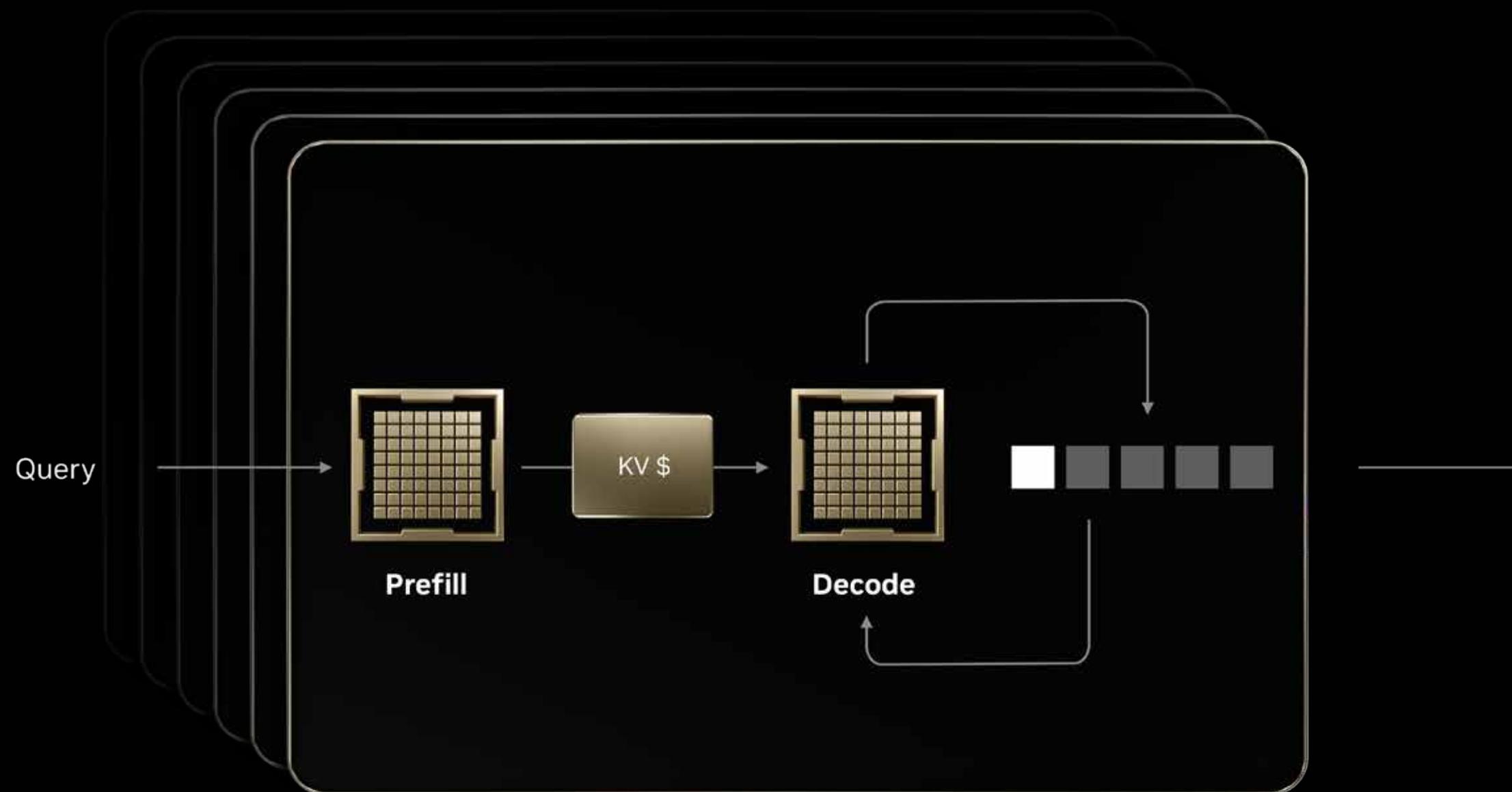


# Reinventing \$500B Enterprise IT For the Age of AI

AI



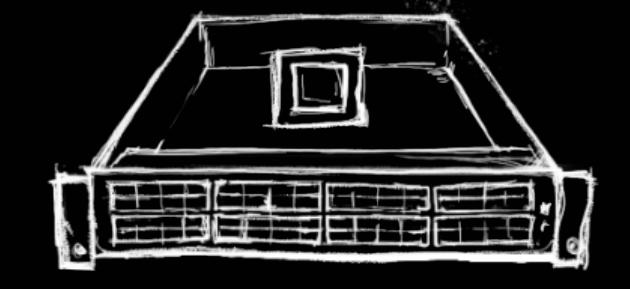
OS



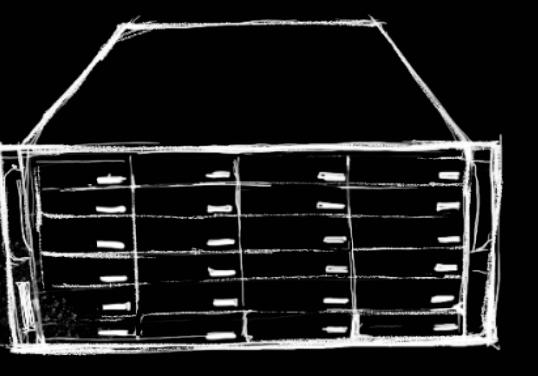
INFRASTRUCTURE



Compute

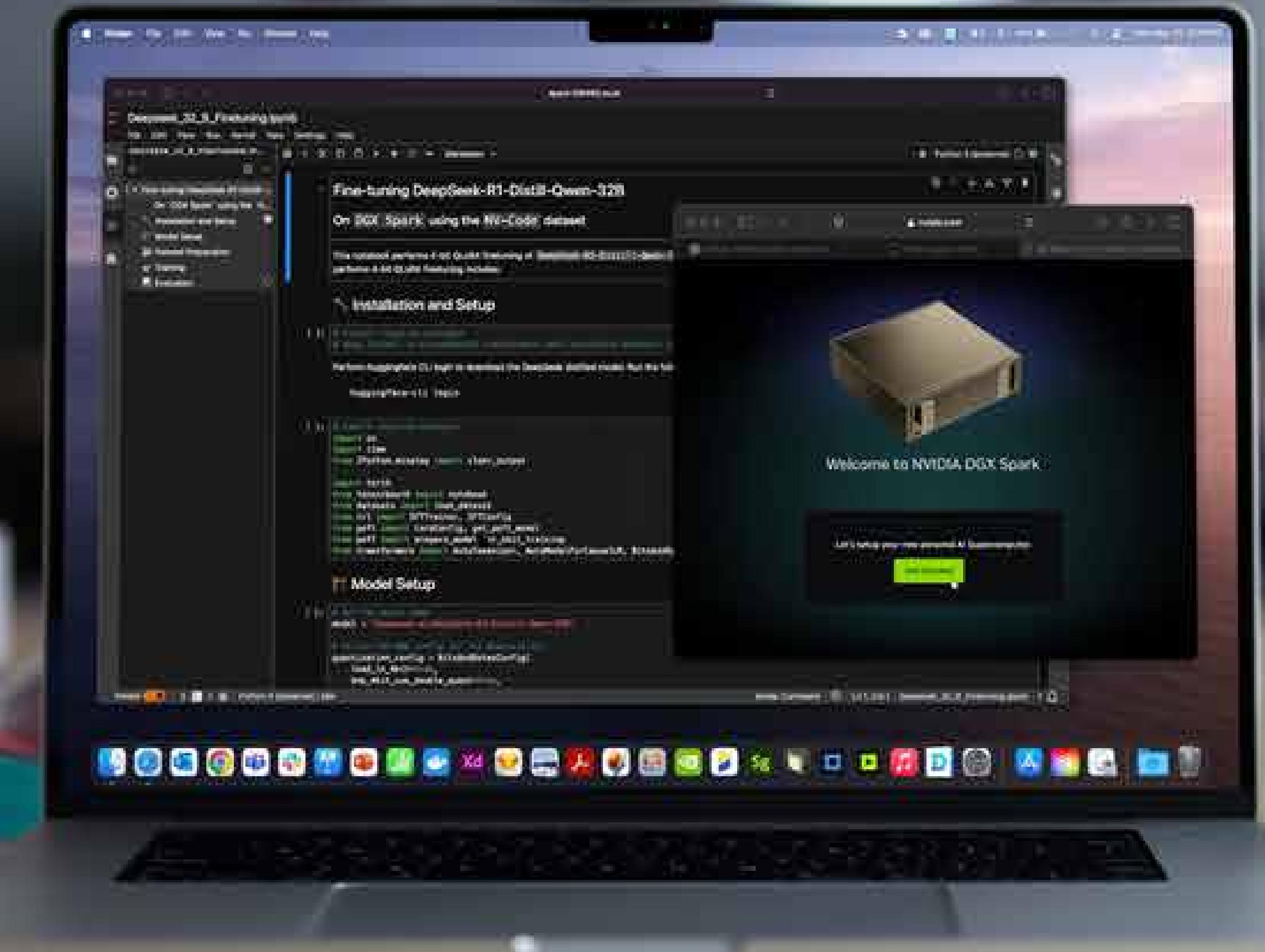


Networking



Storage







# DGX Station

The Ultimate Workstation  
for AI and Data Science

- GB300 Superchip
- 784 GB Unified System Memory
- 20,000 AI TFLOPS
- ConnectX-8 SuperNIC

ASUS

BOXX

DELL Technologies

hp

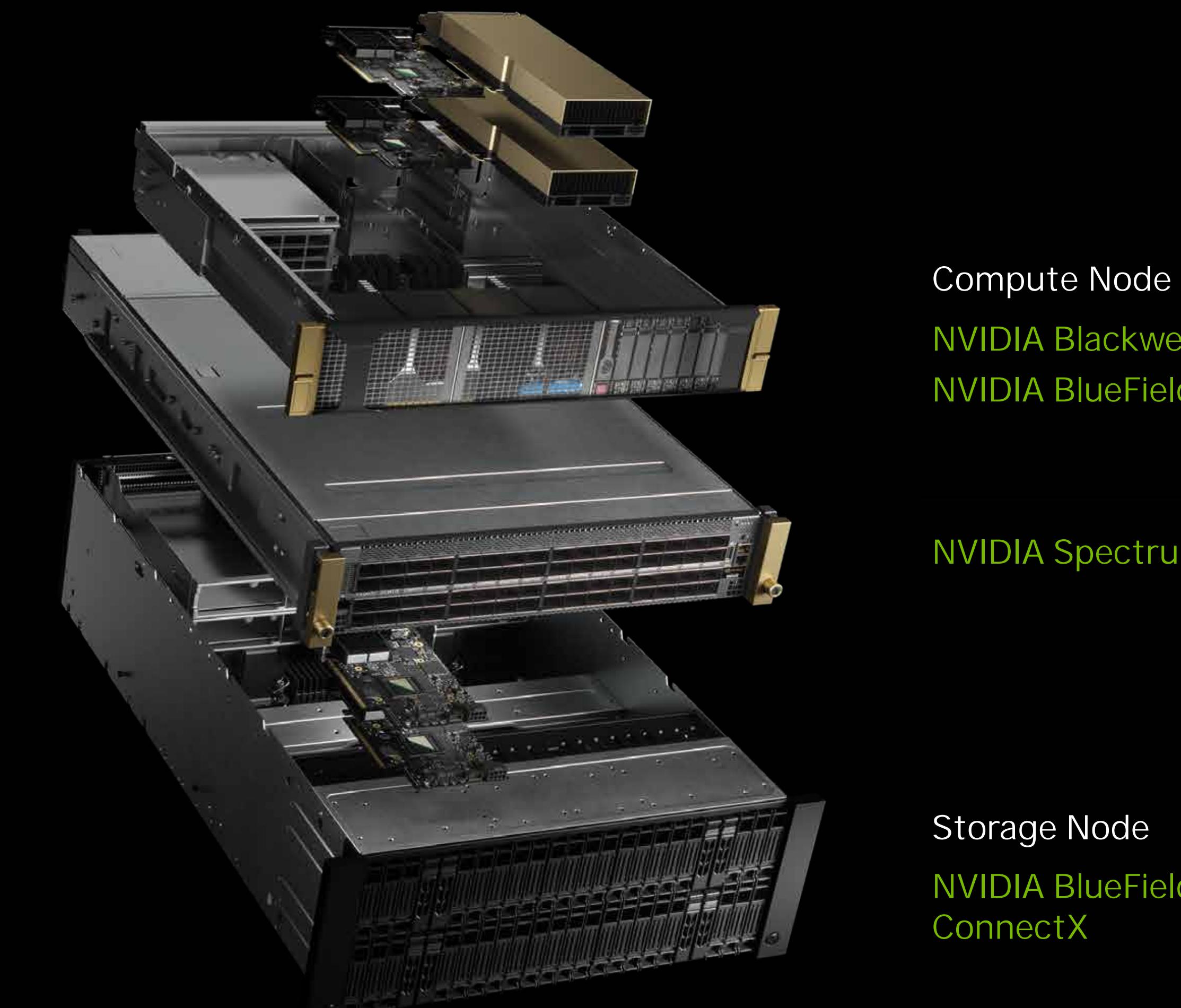
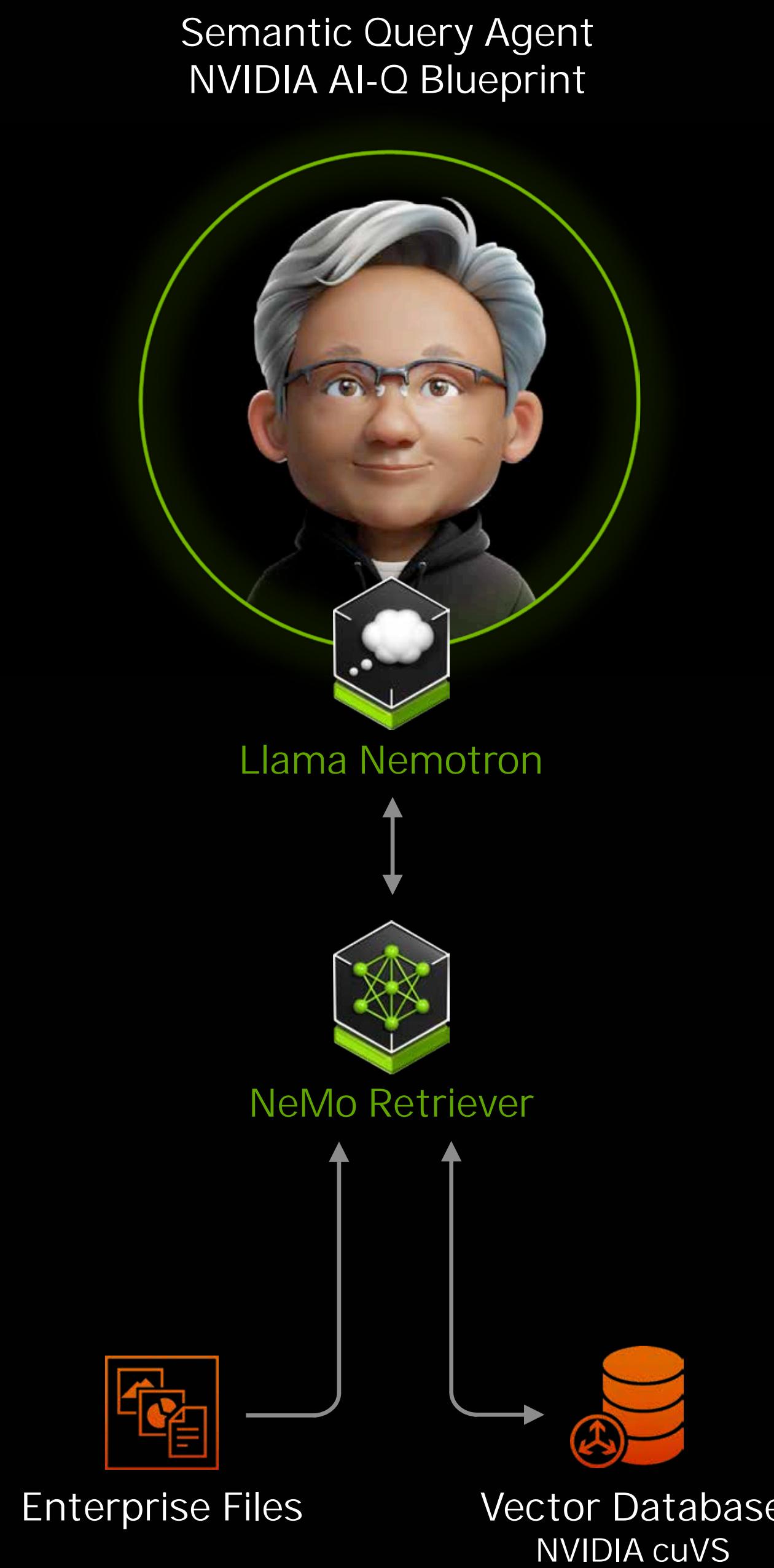
Lambda

SUPERMICRO

# NVIDIA AI Infrastructure for Enterprise Computing



# Announcing Storage Leaders Build AI Data Platforms for Enterprise AI





## INDUSTRY'S BROADEST END-TO-END NVIDIA AI ENTERPRISE INFRASTRUCTURE



AI PCs



Spark



Station



Compute



Networking



Storage



Blackwell Compute Racks

2,000+ Customers

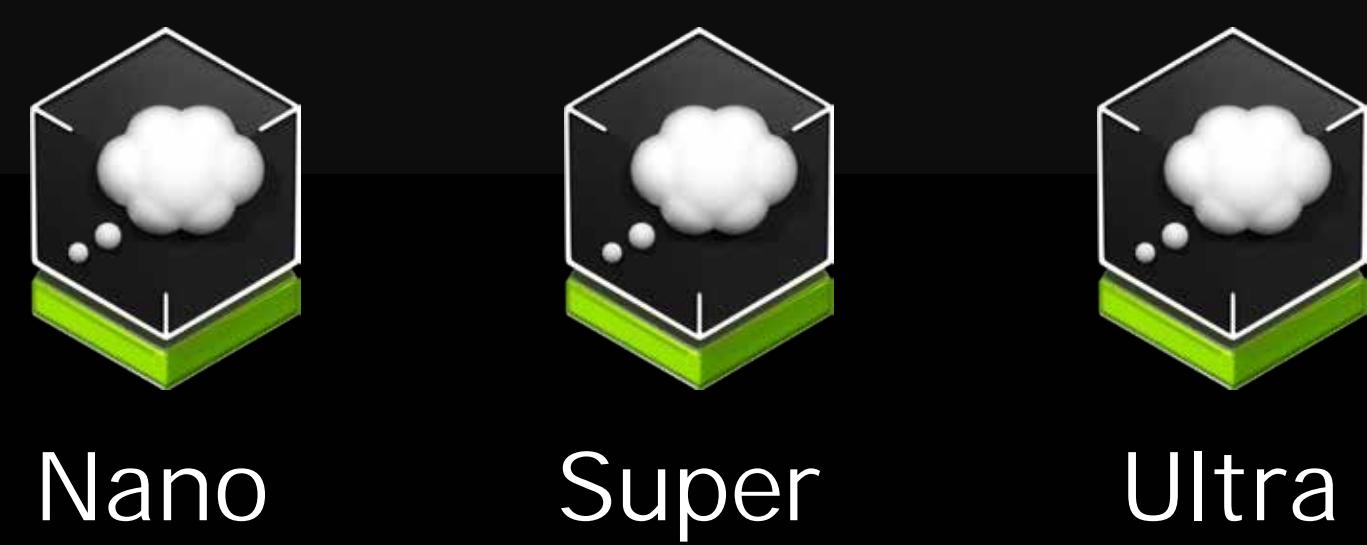
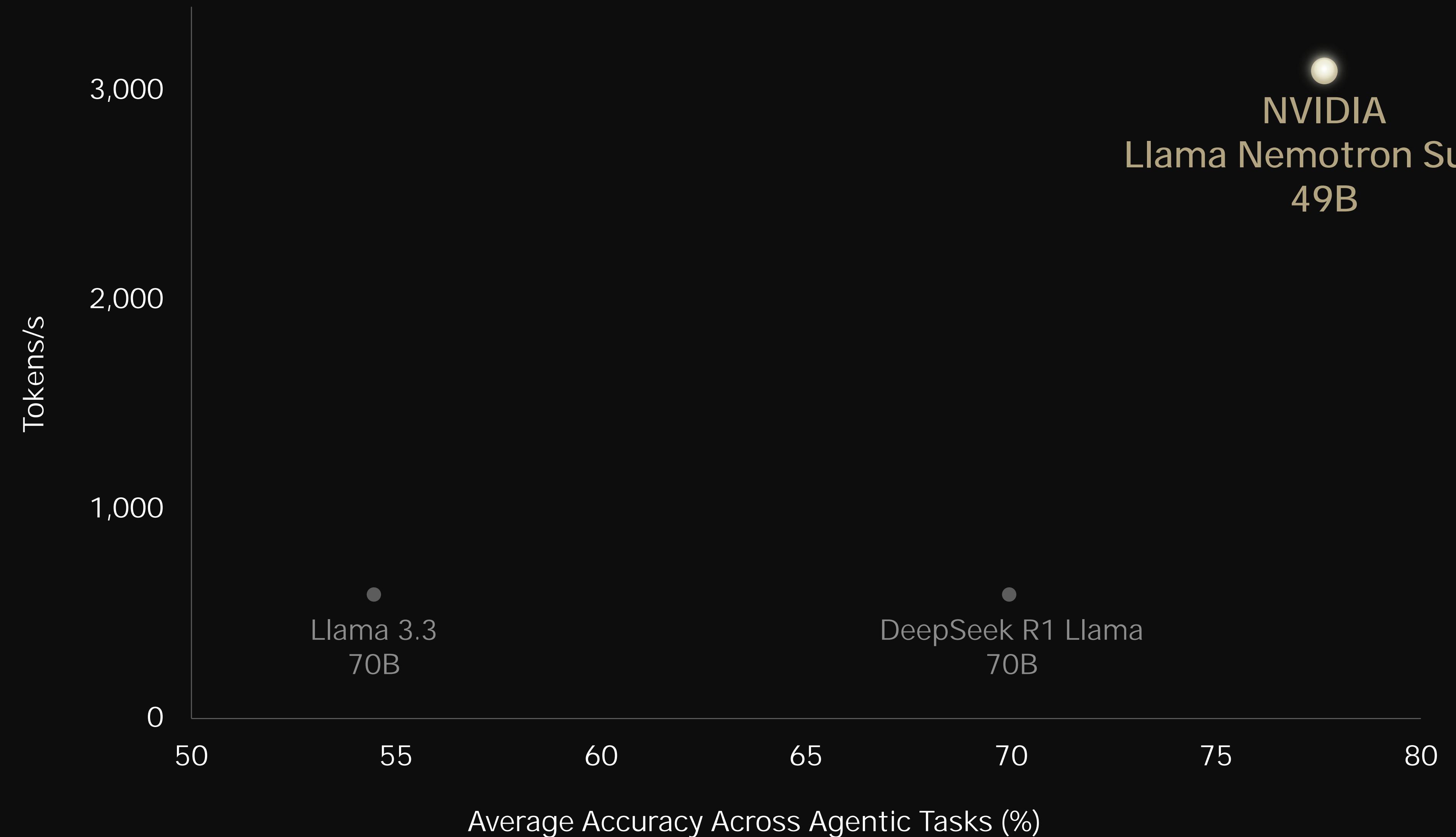


100+ New releases

49+ Exaflops

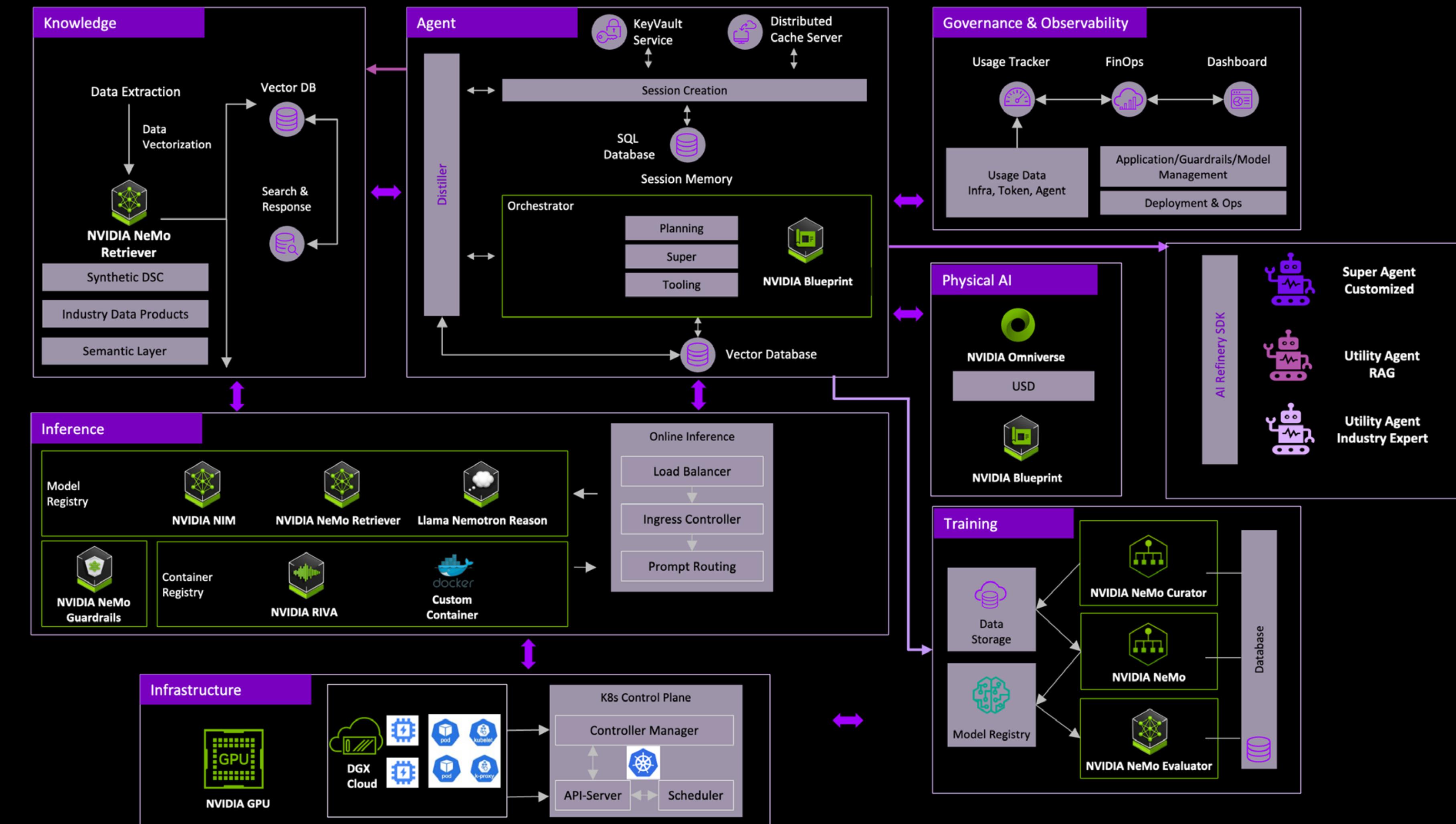
# Announcing NVIDIA Llama Nemotron Reasoning

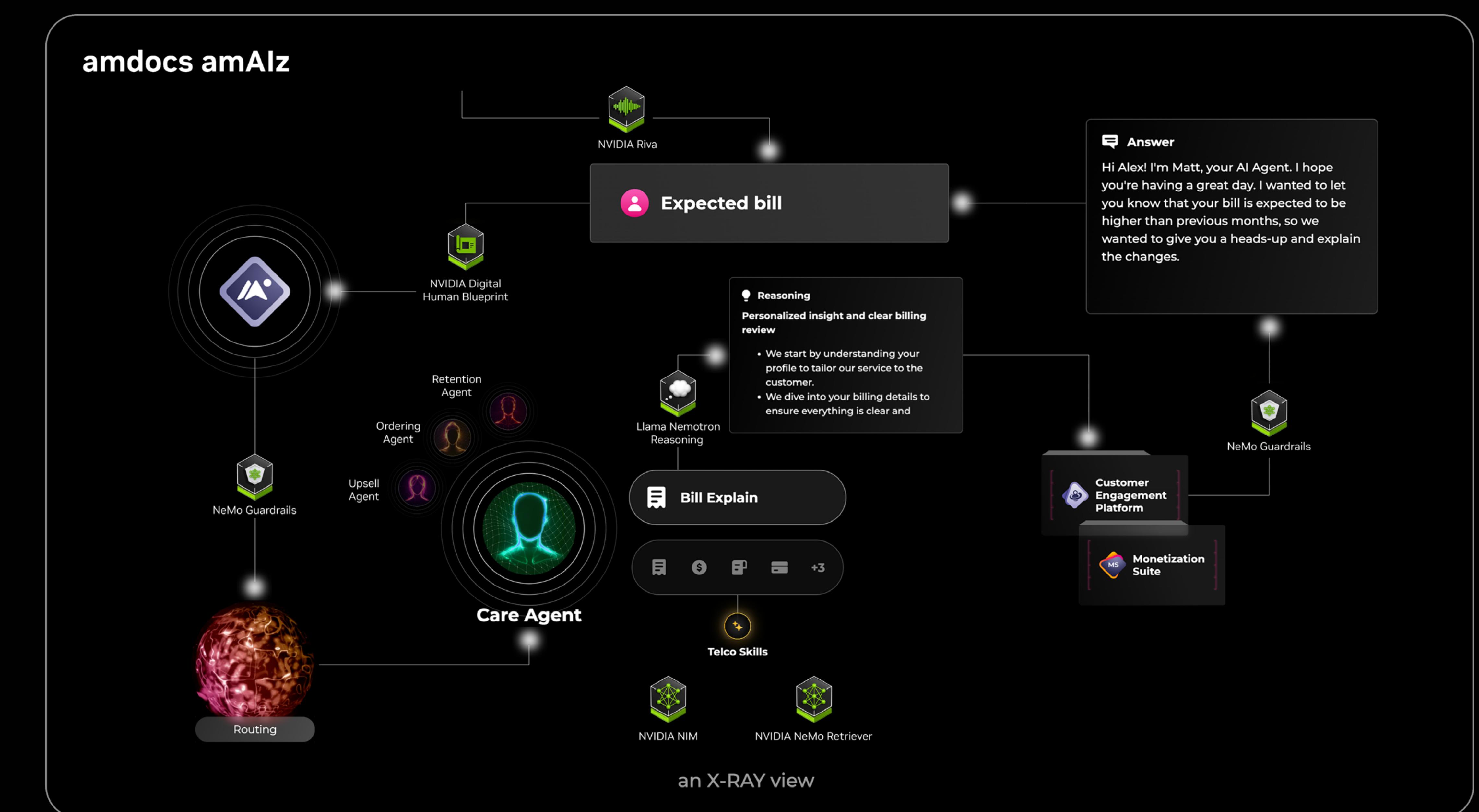
Distilled, Quantized, Aligned, and Optimized by NVIDIA



# accenture

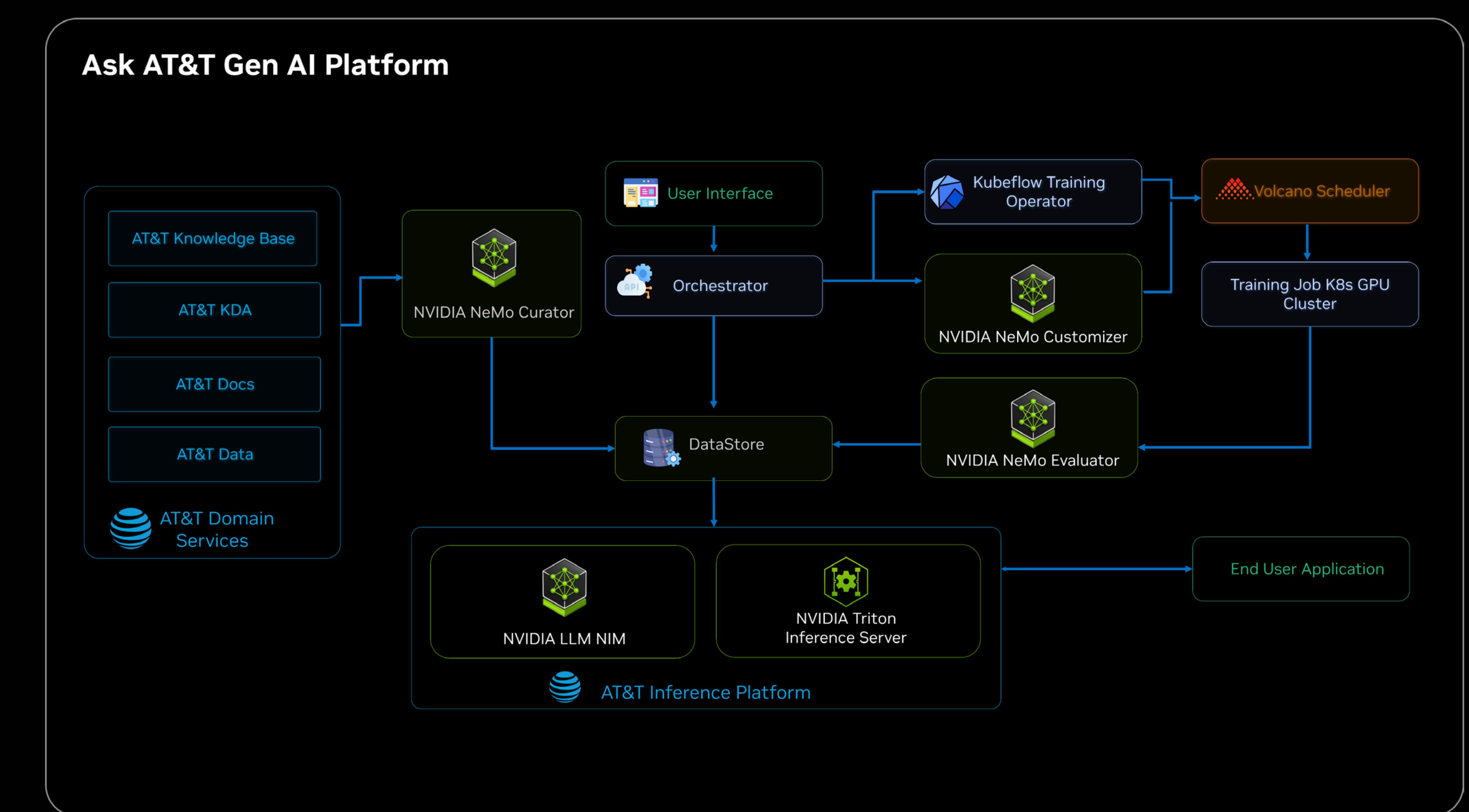
## Accenture AI Refinery



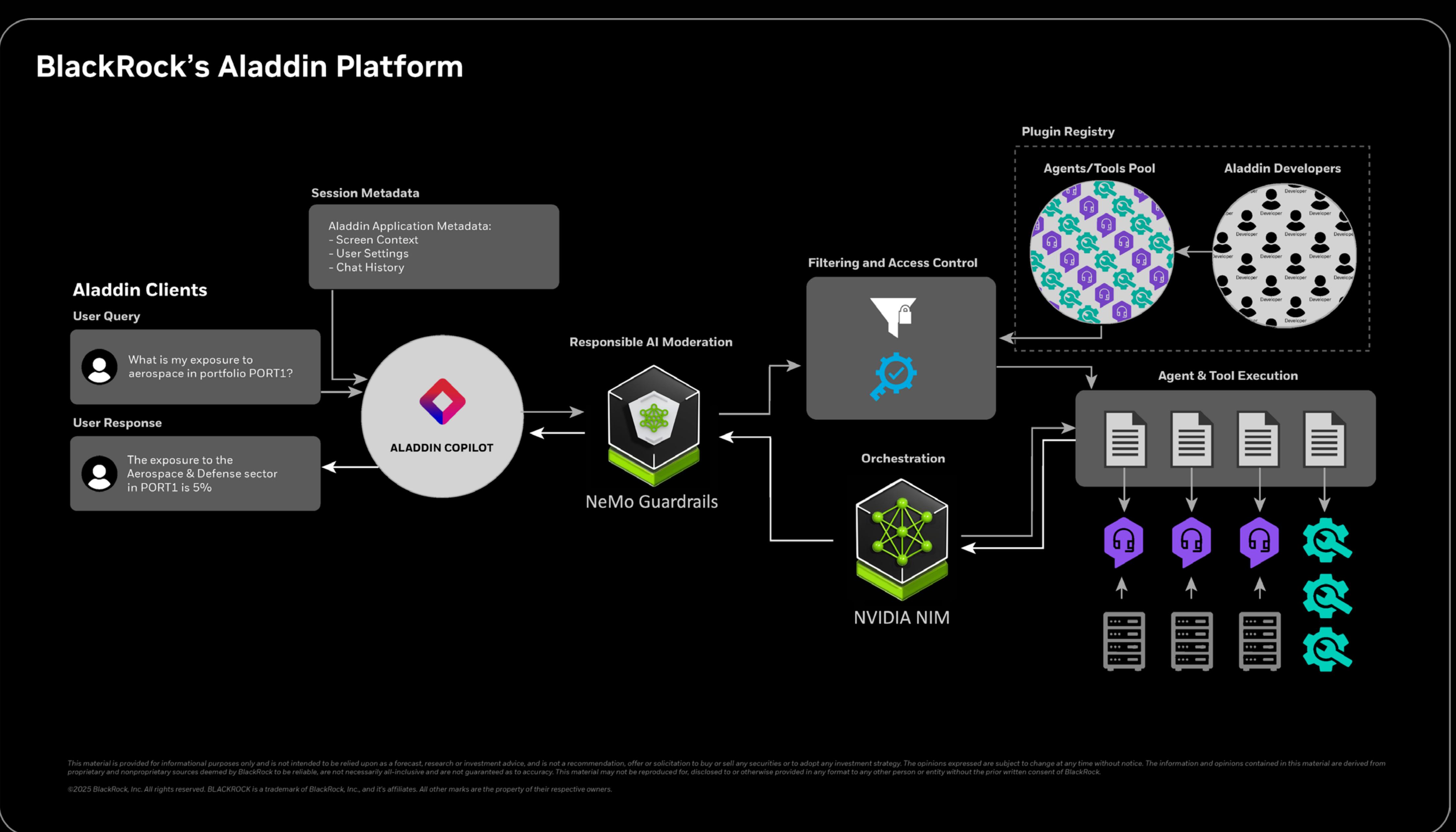




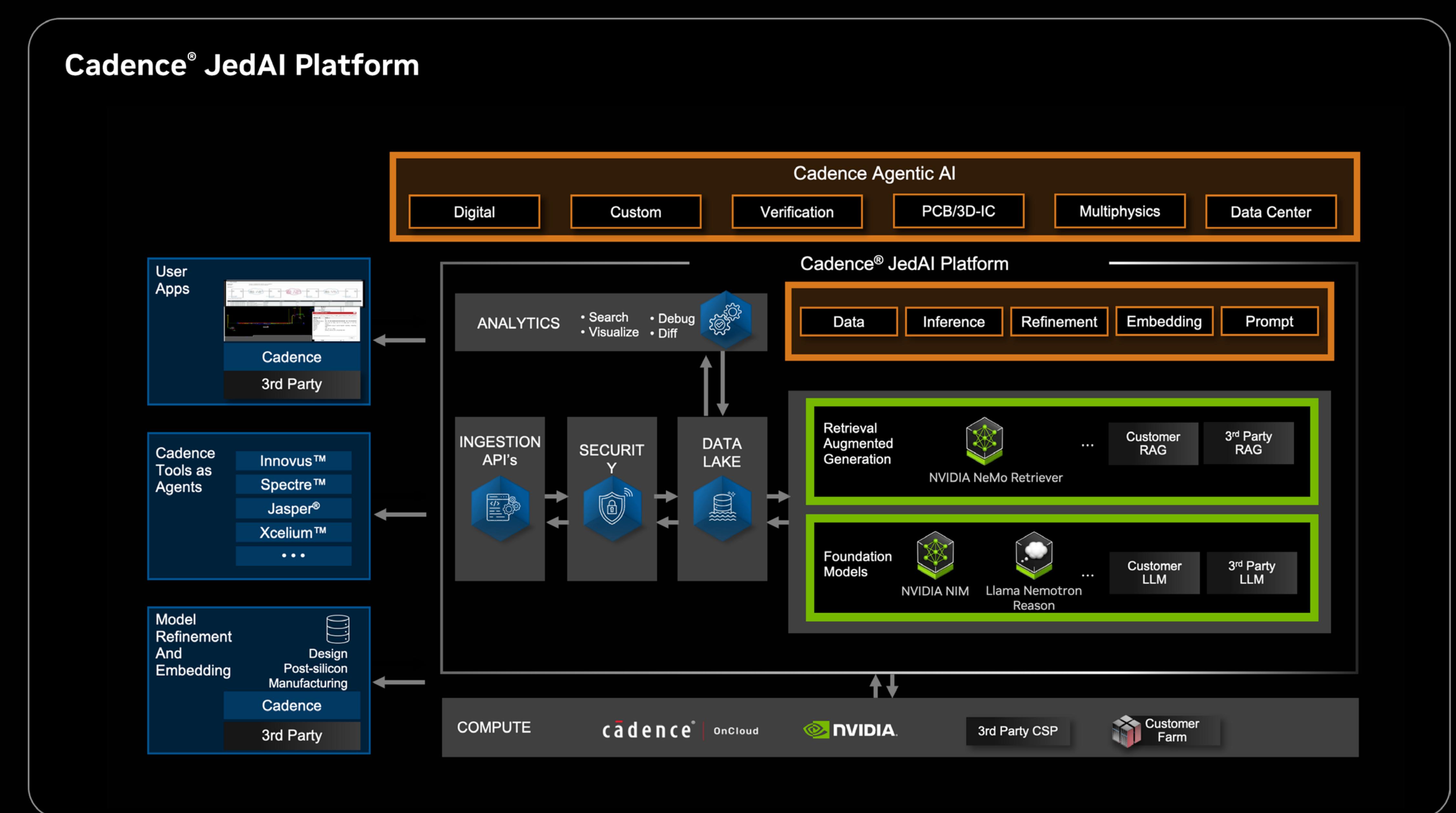
# AT&T



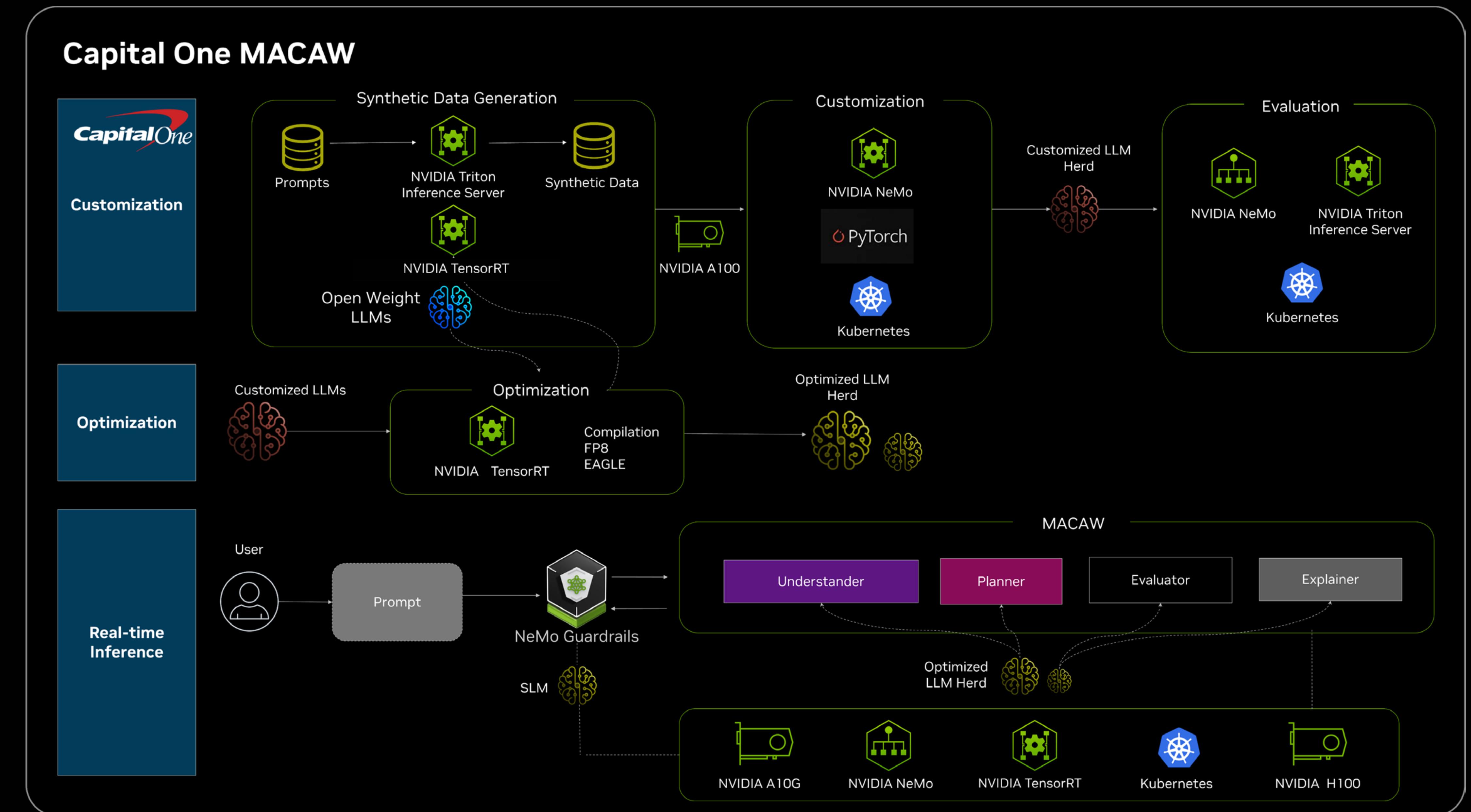
# BlackRock®



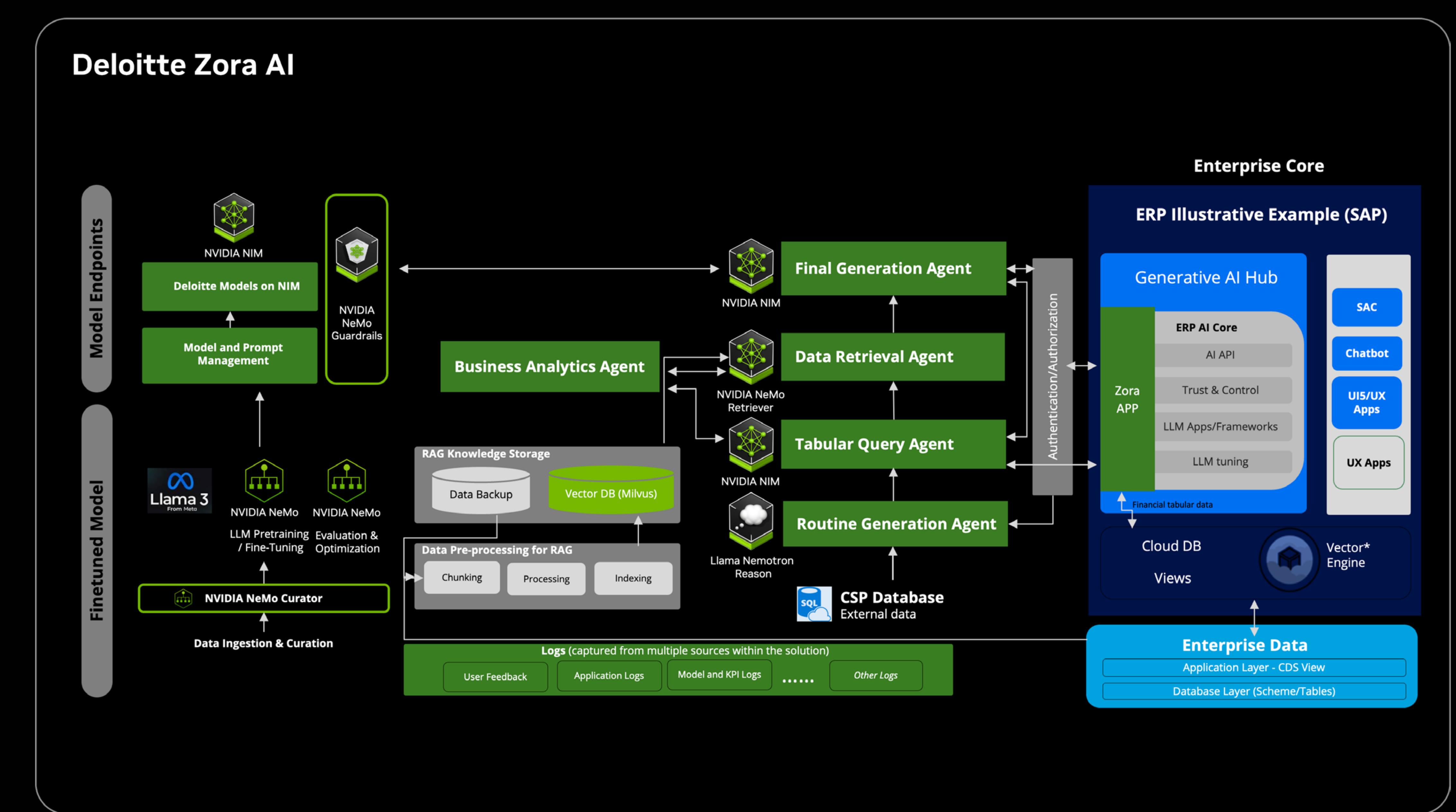
# cadence®

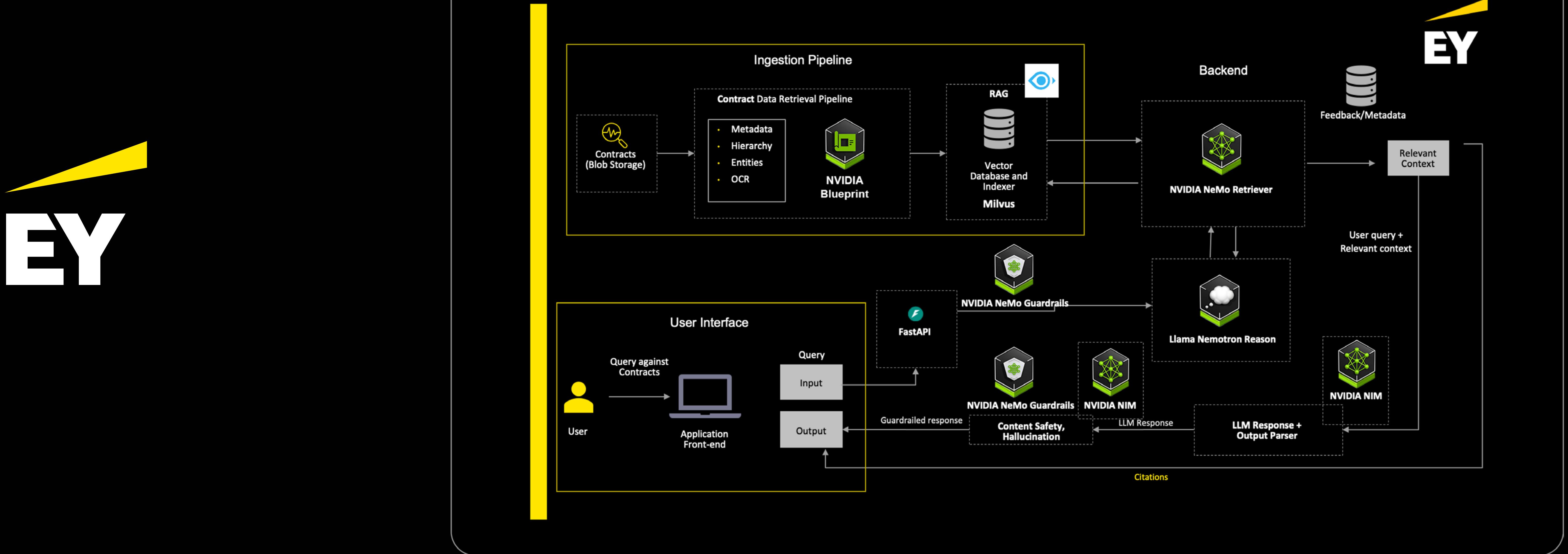


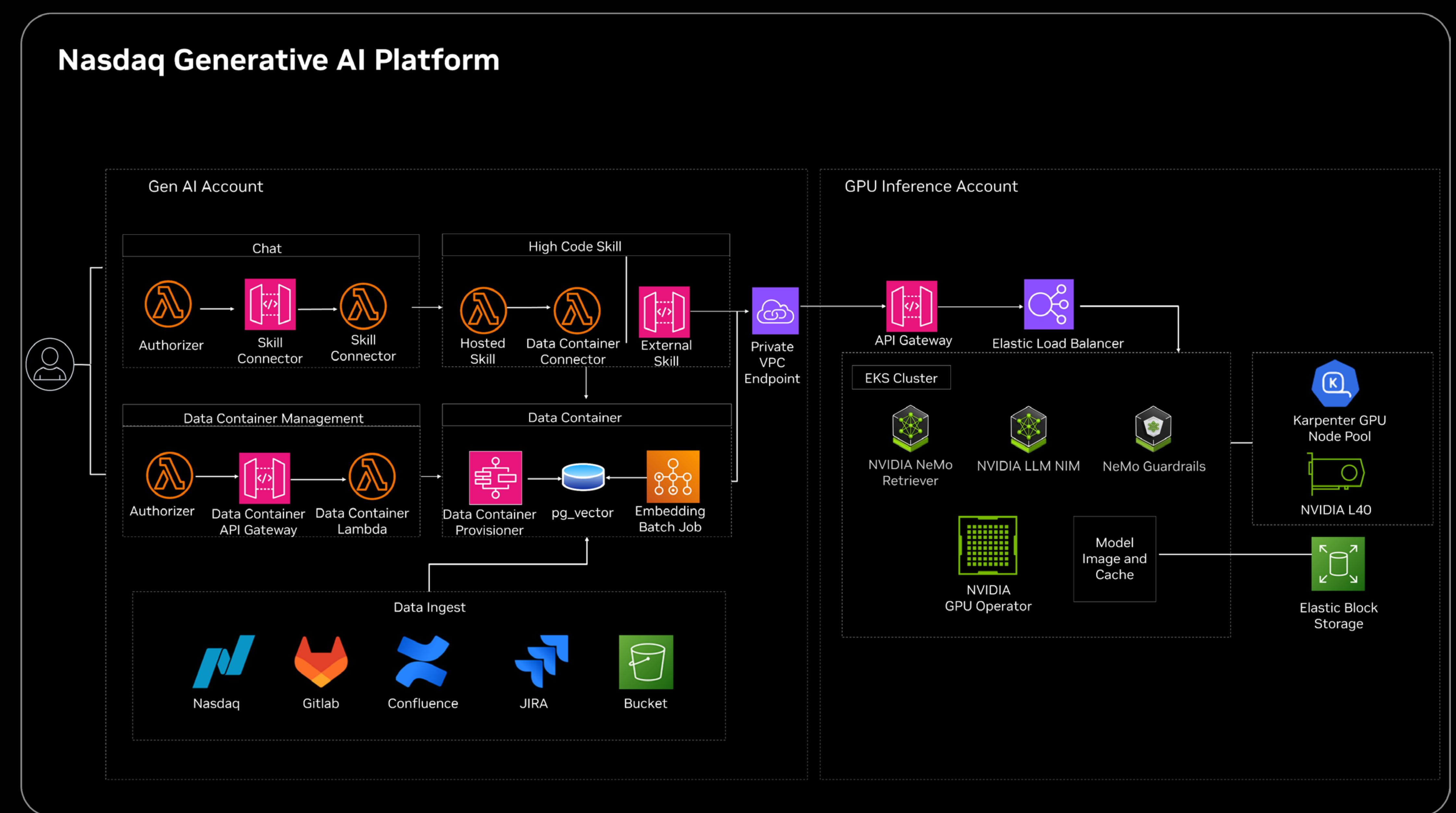
# CapitalOne



# Deloitte.

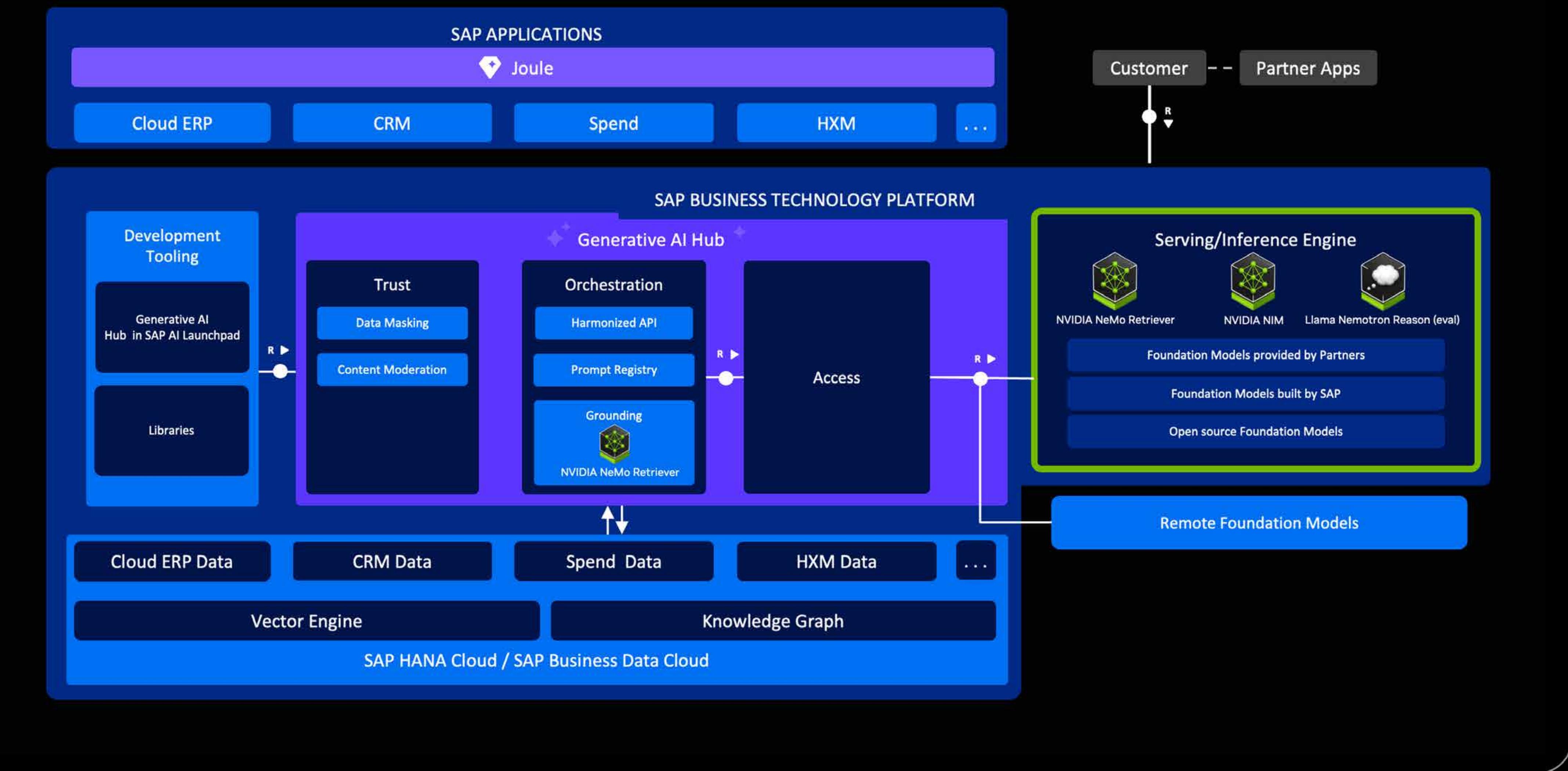




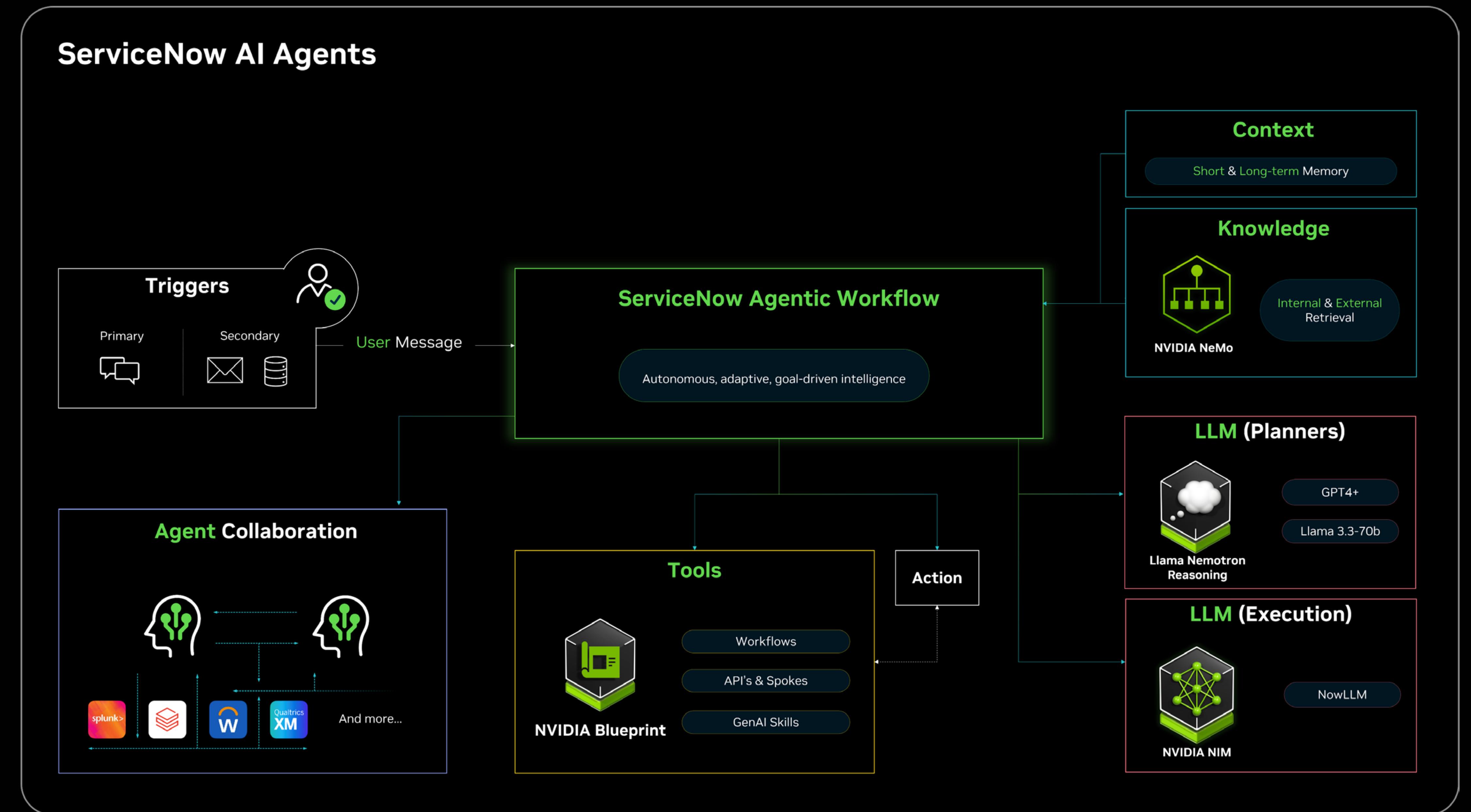




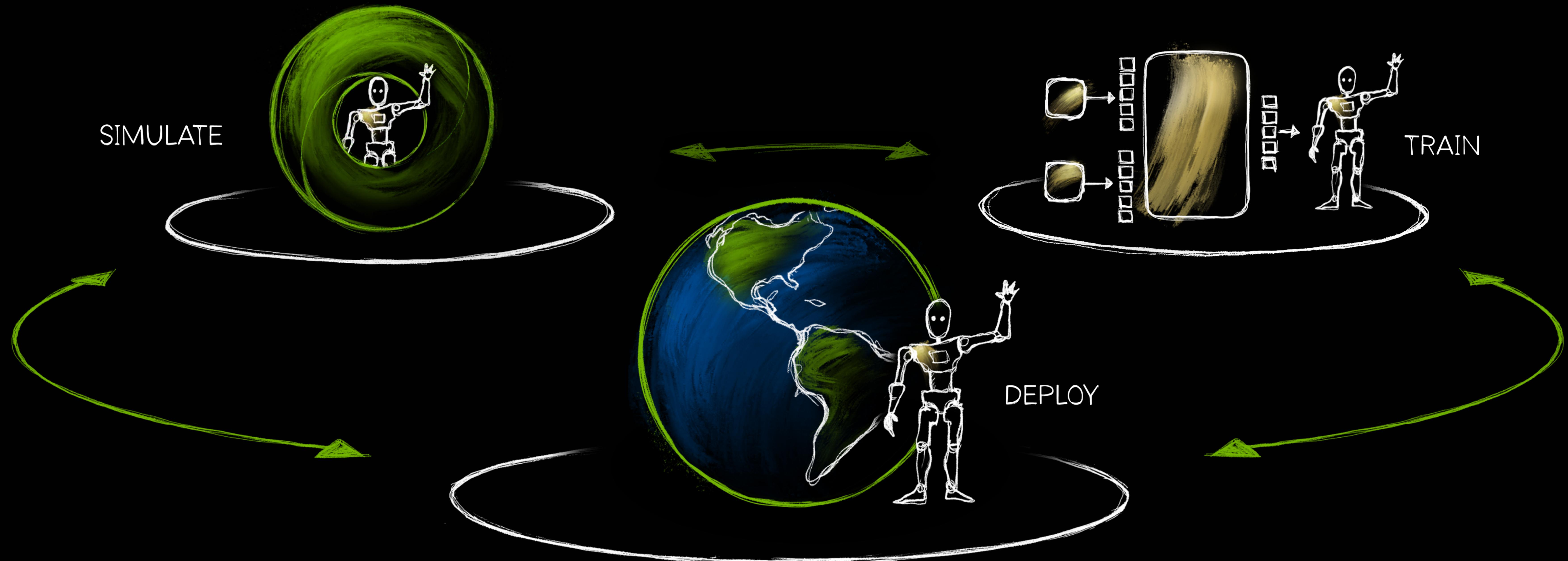
## SAP Business AI



# servicenow®

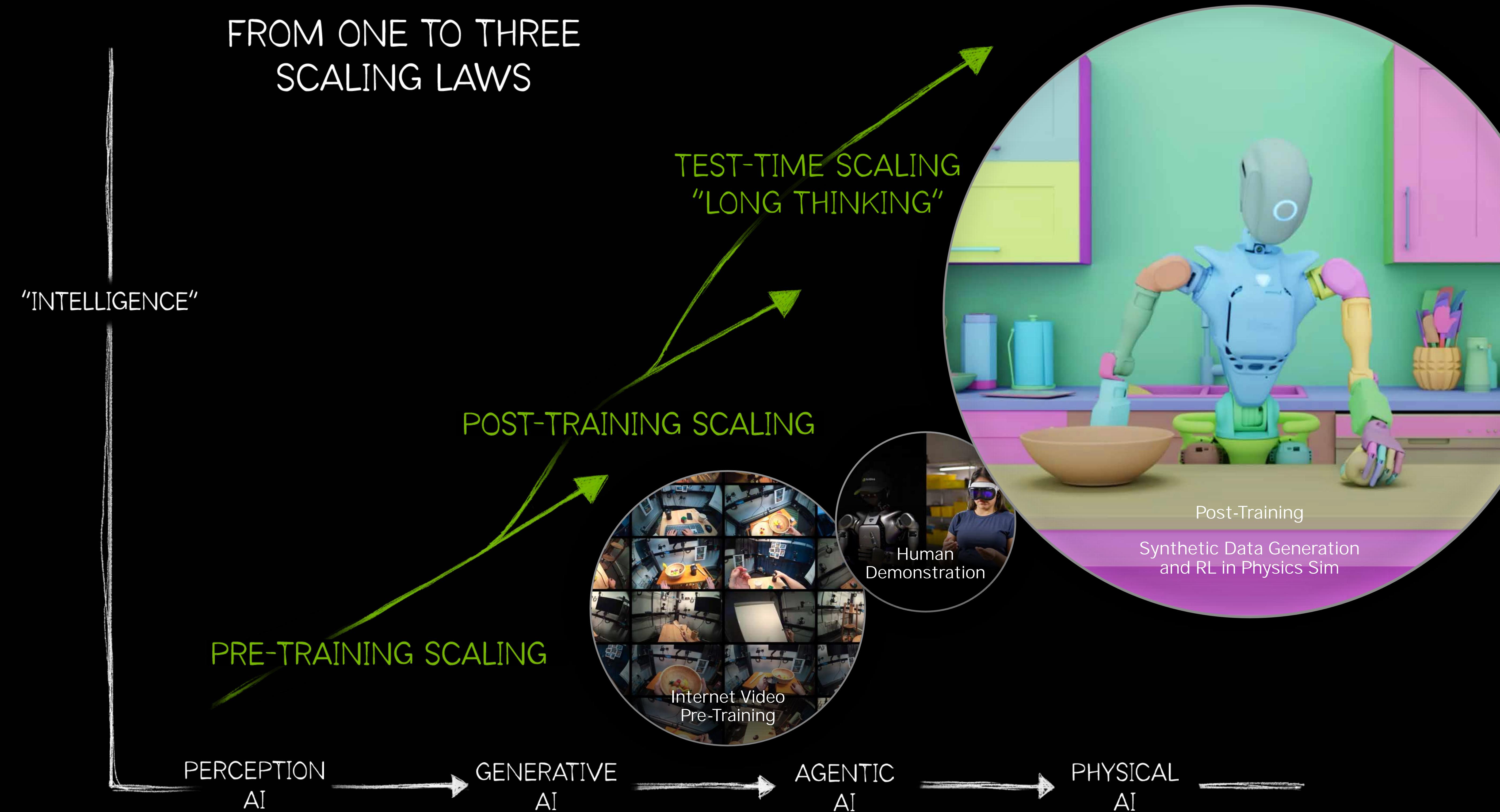


# ROBOTS – THE NEXT MULTI-TRILLION DOLLAR INDUSTRY



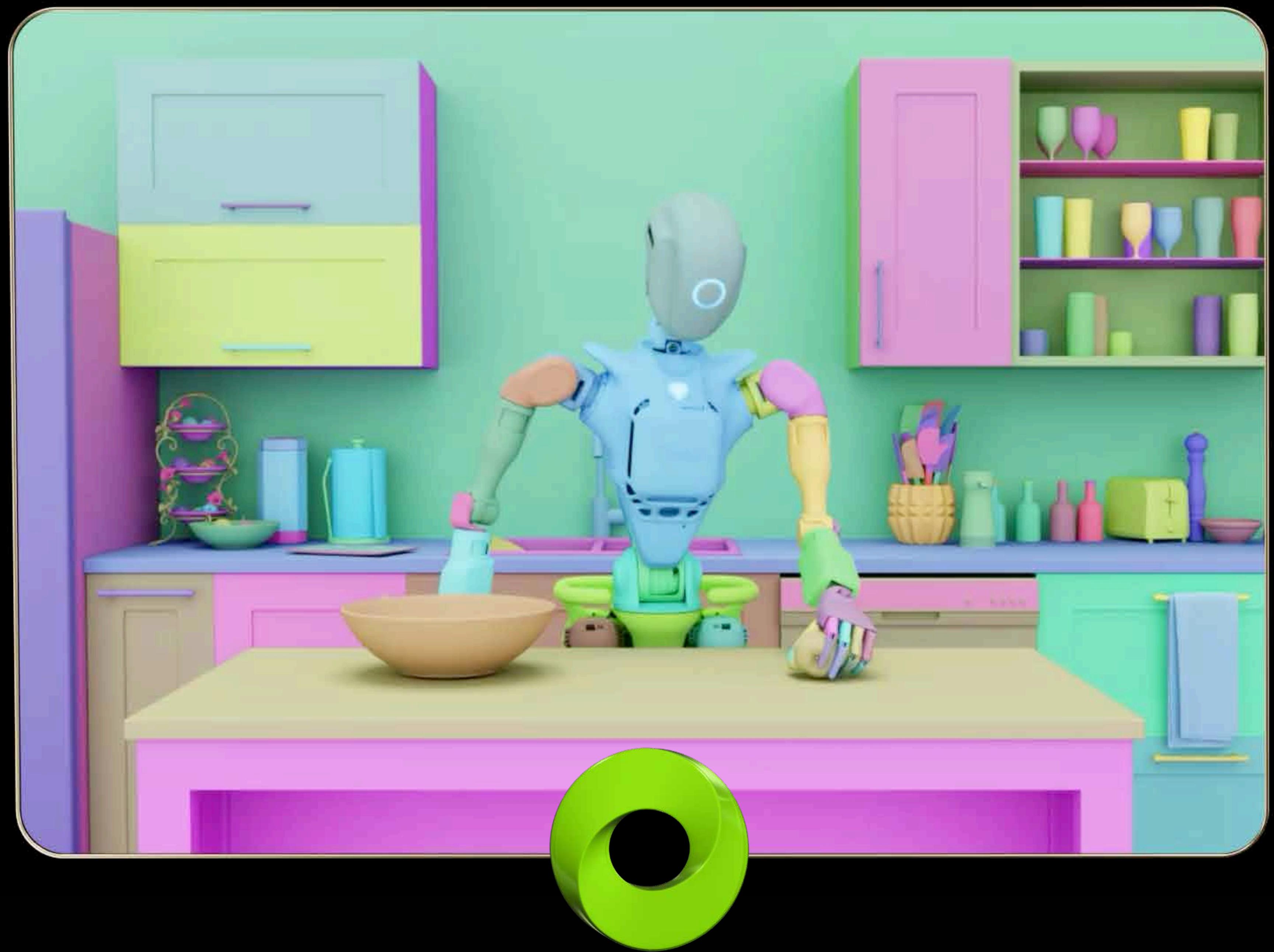


## FROM ONE TO THREE SCALING LAWS



# NVIDIA Omniverse With Cosmos

Physical AI Digital Twin Operating System



NVIDIA Omniverse

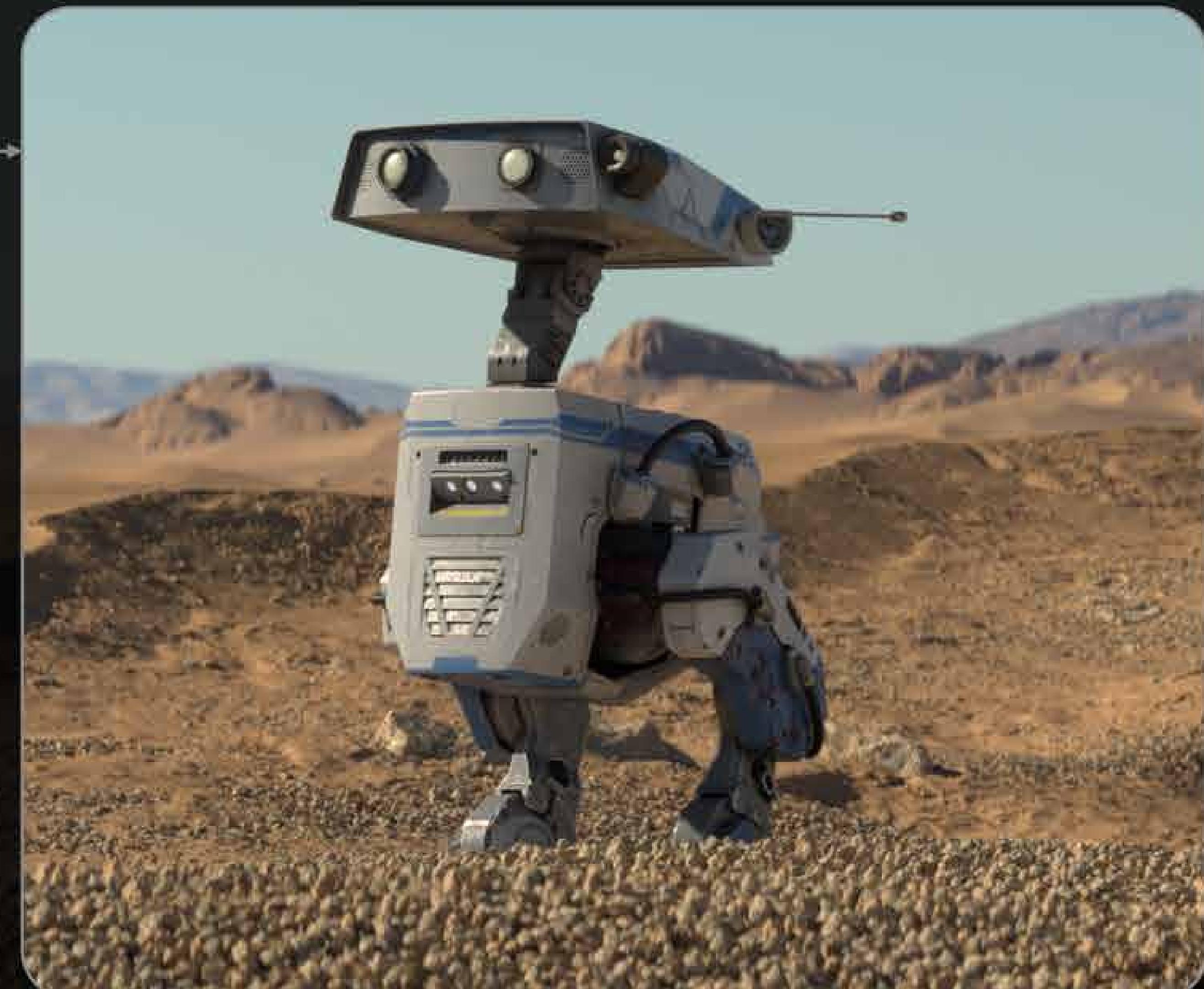
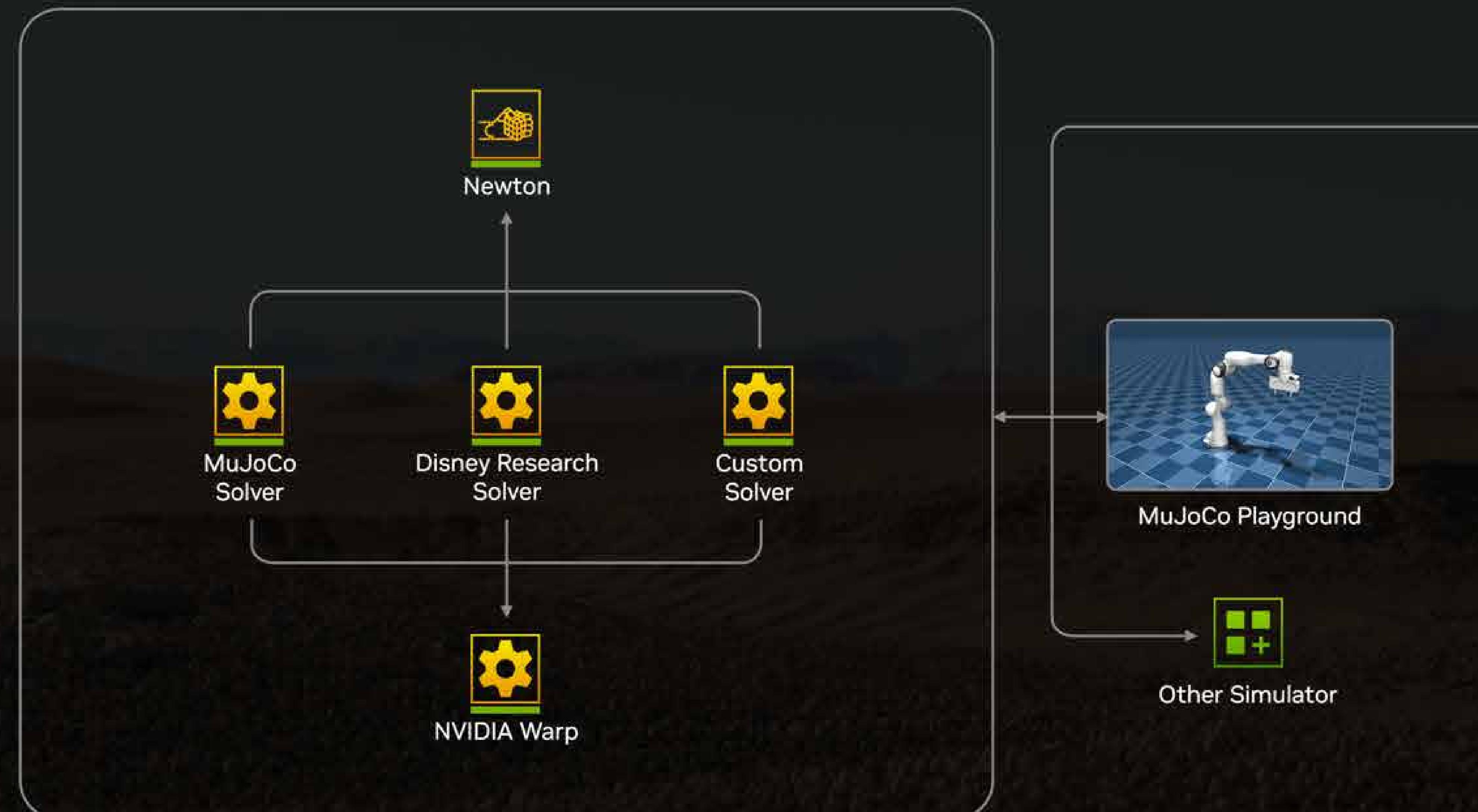


Cosmos

Disney Research

NVIDIA

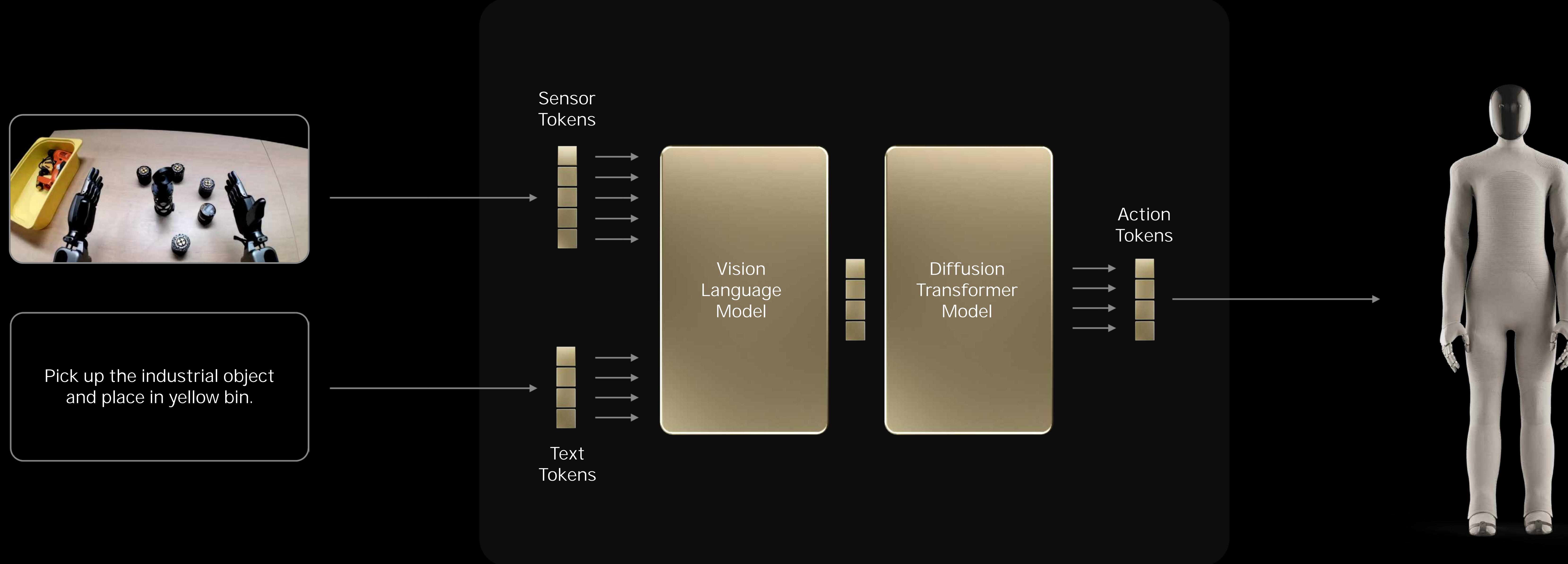
Google DeepMind



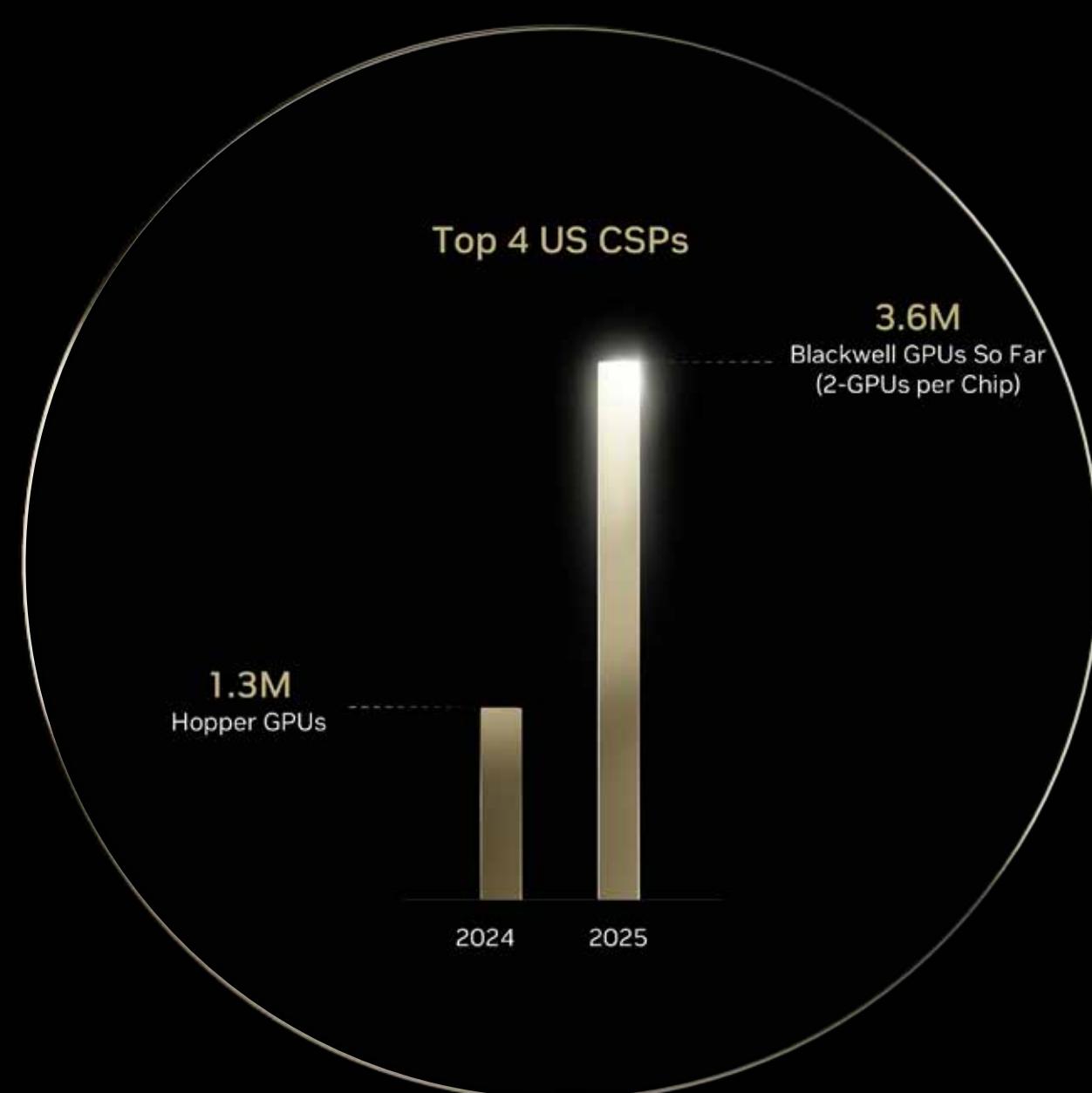
Isaac Lab



# Announcing NVIDIA Isaac GR0OT N1 Humanoid Foundation Model

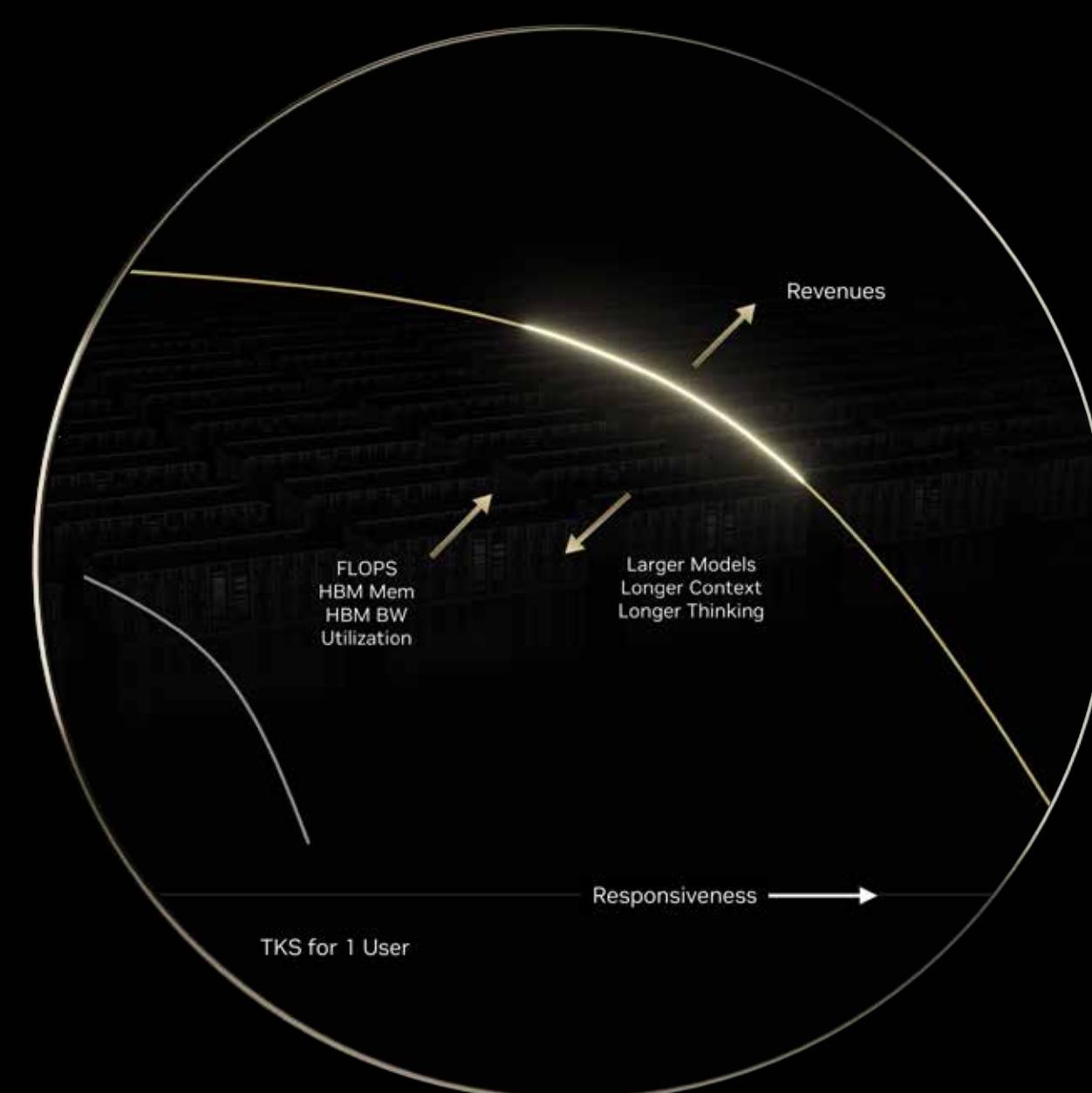


# NVIDIA GTC 2025



BLACKWELL IN FULL PRODUCTION

\$1T Computing Inflection Point



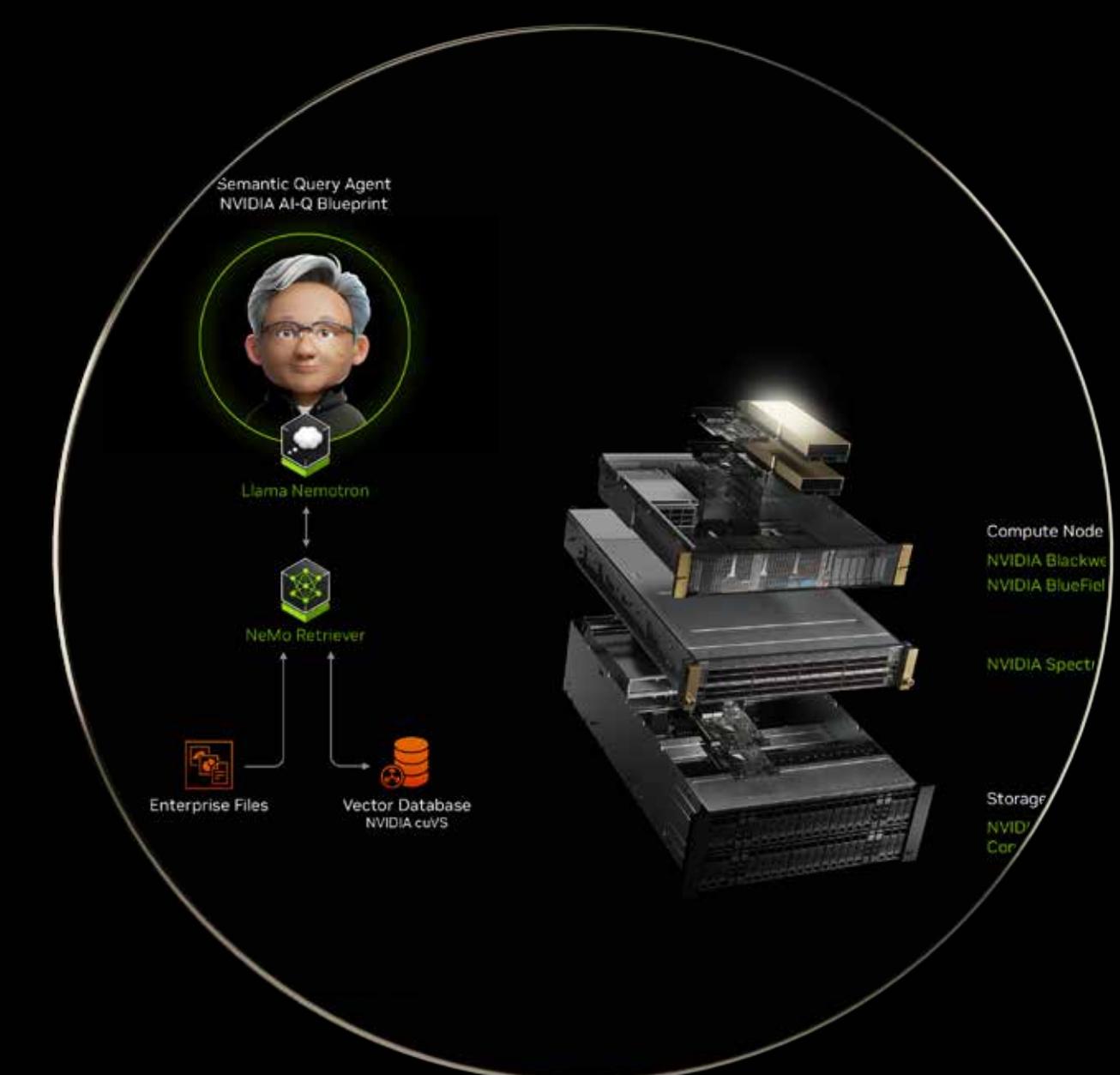
BLACKWELL NVL72 AND DYNAMO 40X FOR INFERENCE

Reasoning 100X One-Shot Blackwell 40X Hopper



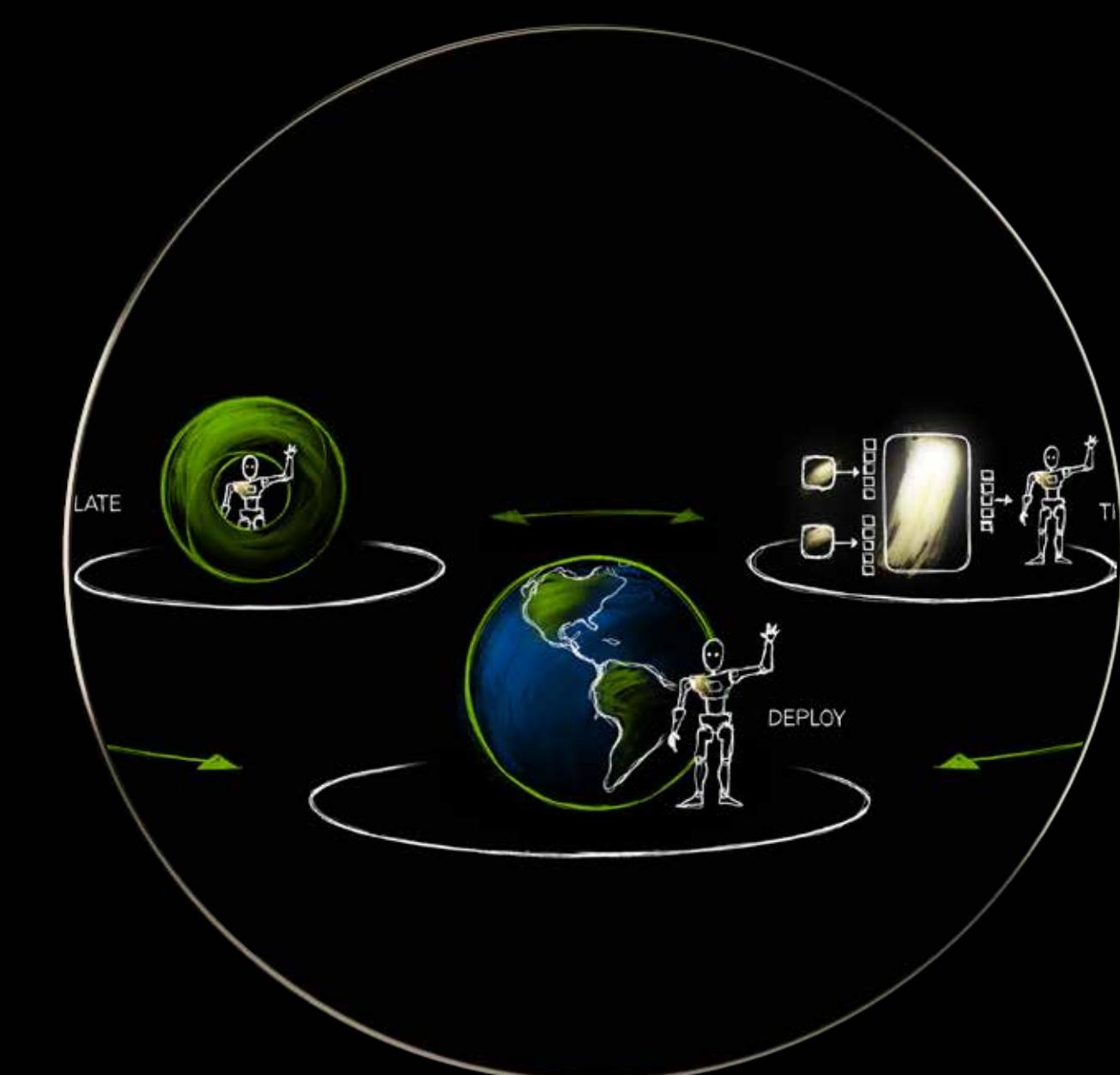
VERA RUBIN  
NVIDIA PHOTONICS  
AI INFRASTRUCTURE  
FOR AI CLOUDS

Annual Rhythm for the World to Build-Out AI Infrastructure



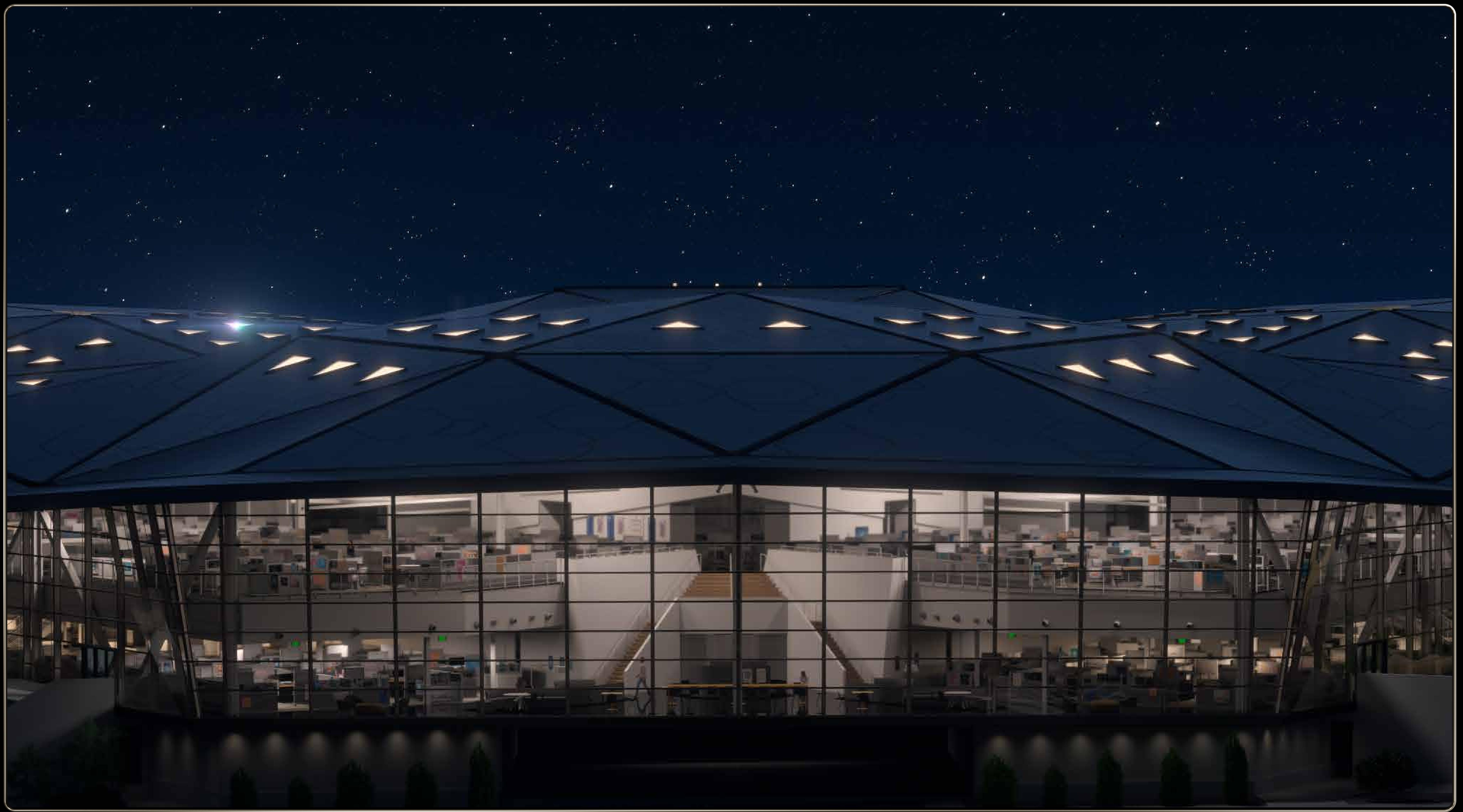
AI INFRASTRUCTURE FOR \$500B ENTERPRISE IT

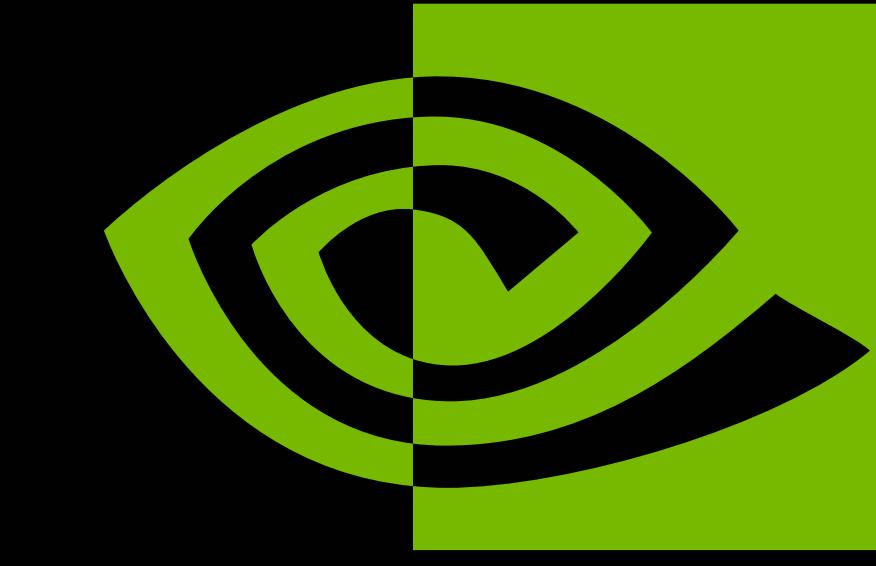
New Compute, Networking, Storage, Software



AI INFRASTRUCTURE FOR ROBOTS

Physical AI For \$50T Industrial and Robotics





nVIDIA