

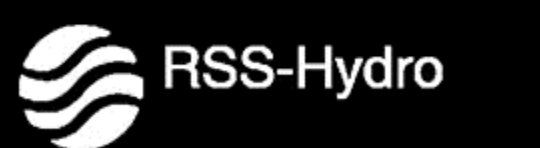
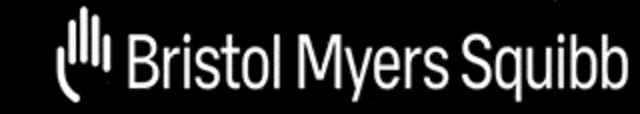


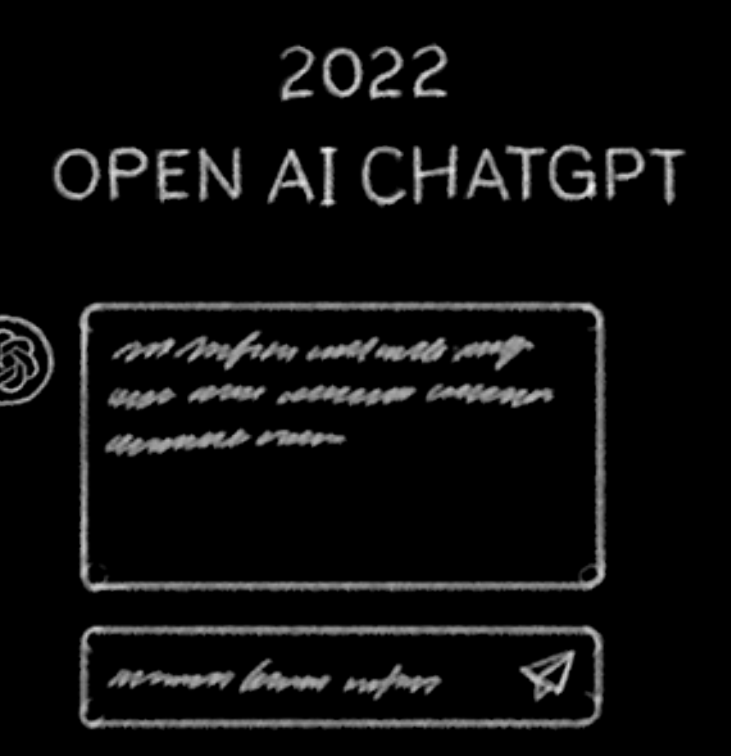
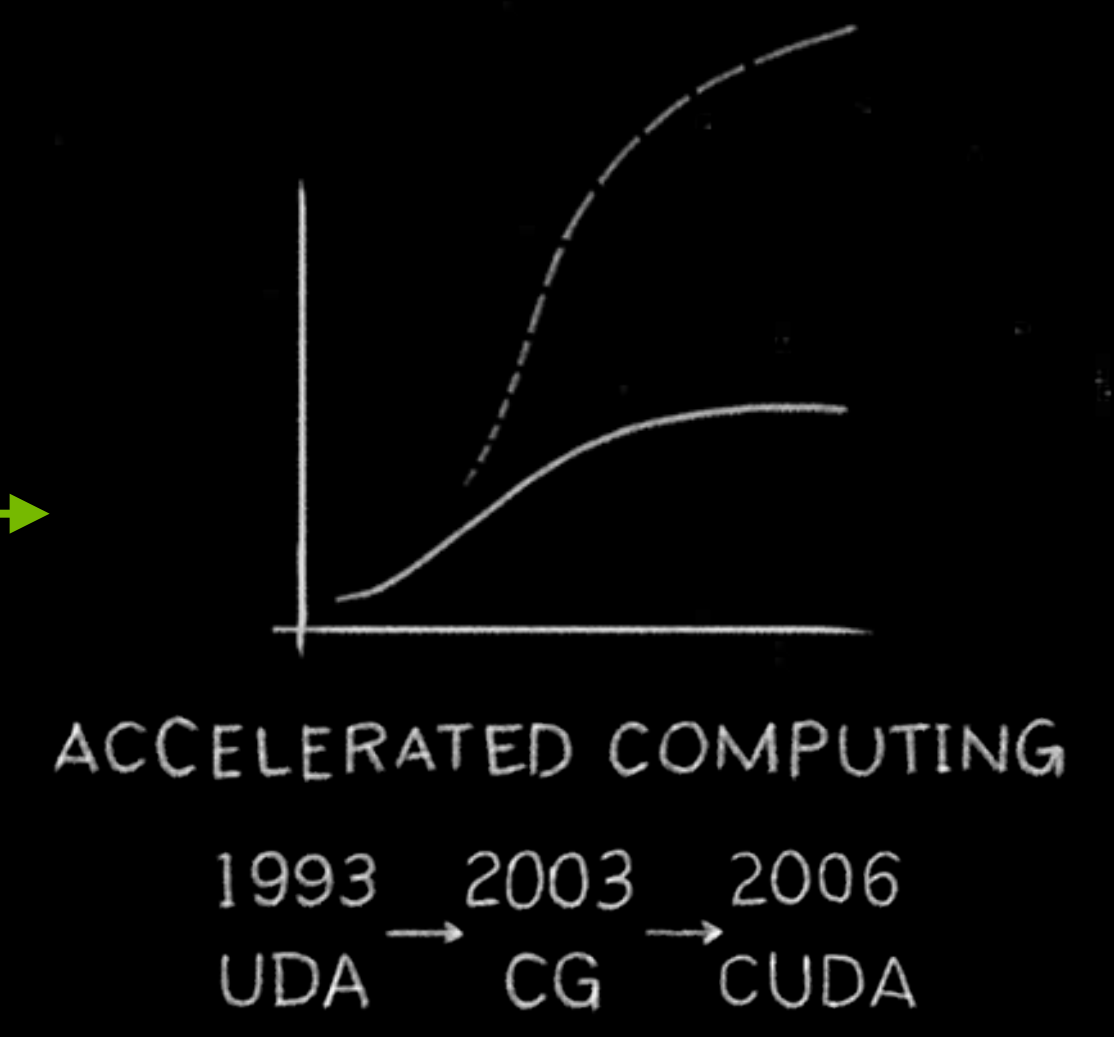
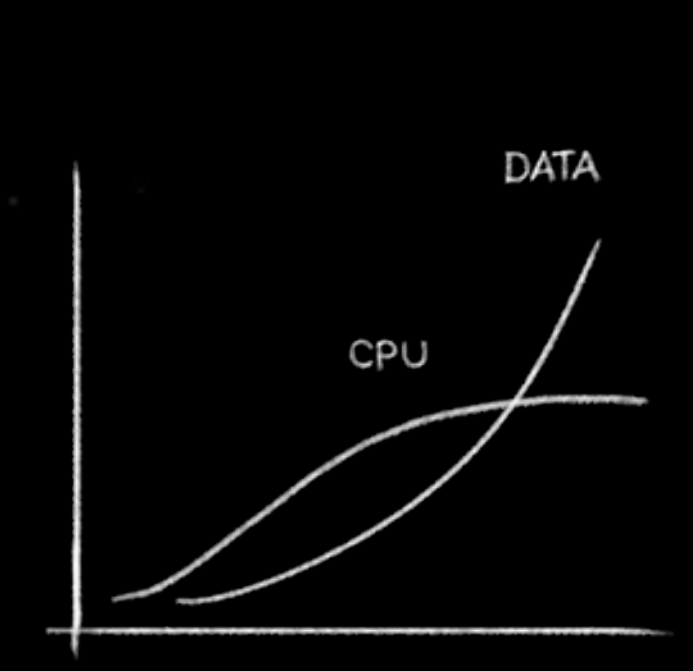


Except for the historical information contained herein, certain matters in this presentation are forward-looking statements. These forward-looking statements and any other forward-looking statements that go beyond historical facts that are made in this presentation are subject to risks and uncertainties that may cause actual results to differ materially. Important factors that could cause actual results to differ materially include: global economic conditions; NVIDIA's reliance on third parties to manufacture, assemble, package and test NVIDIA's products; the impact of technological development and competition; development of new products and technologies or enhancements to NVIDIA's existing product and technologies; market acceptance of NVIDIA's products or NVIDIA's partners' products; design, manufacturing or software defects; changes in consumer preferences and demands; changes in industry standards and interfaces; unexpected loss of performance of NVIDIA's products or technologies when integrated into systems and other factors.

NVIDIA has based these forward-looking statements largely on its current expectations and projections about future events and trends that it believes may affect its financial condition, results of operations, business strategy, short-term and long-term business operations and objectives, and financial needs. These forward-looking statements are subject to a number of risks and uncertainties, and you should not rely upon the forward-looking statements as predictions of future events. The future events and trends discussed in this presentation may not occur and actual results could differ materially and adversely from those anticipated or implied in the forward-looking statements. Although NVIDIA believes that the expectations reflected in the forward-looking statements are reasonable, the company cannot guarantee that future results, levels of activity, performance, achievements or events and circumstances reflected in the forward-looking statements will occur. Except as required by law, NVIDIA disclaims any obligation to update these forward-looking statements to reflect future events or circumstances. For a complete discussion of factors that could materially affect NVIDIA's financial results and operations, please refer to the reports NVIDIA files from time to time with the SEC, including NVIDIA's most recent Annual Report on Form 10-K, Quarterly Reports on Form 10-Q, and Current Reports on Form 8-K. Copies of reports NVIDIA files with the SEC are posted on NVIDIA's website and are available from NVIDIA without charge.

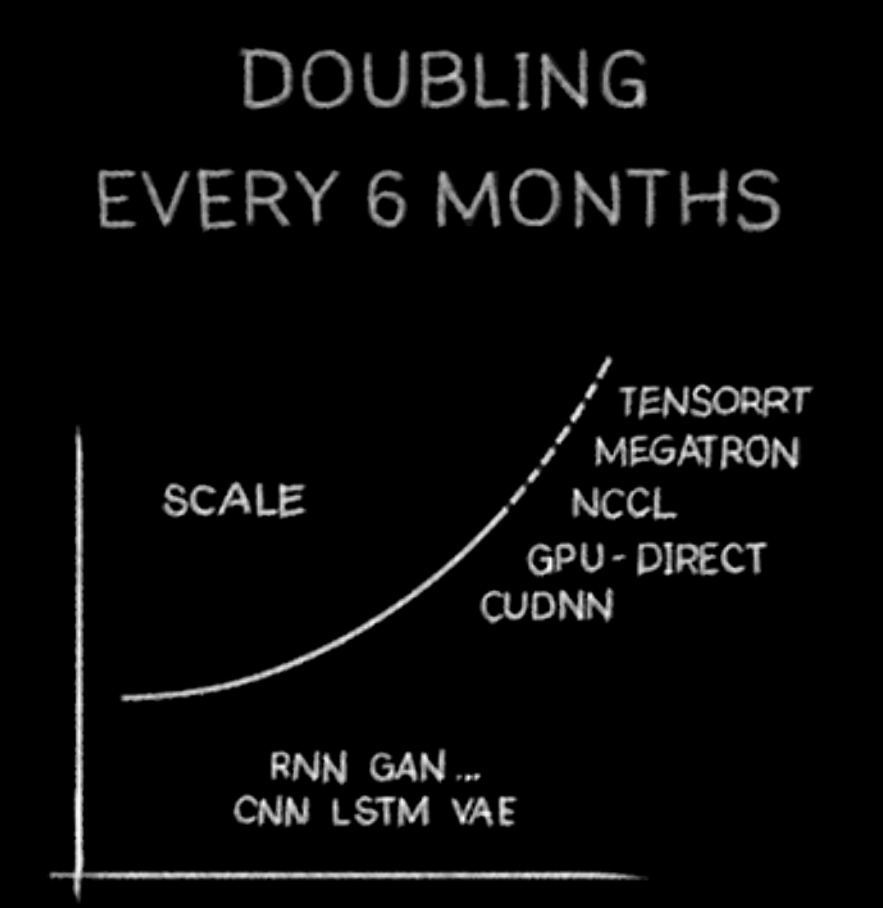
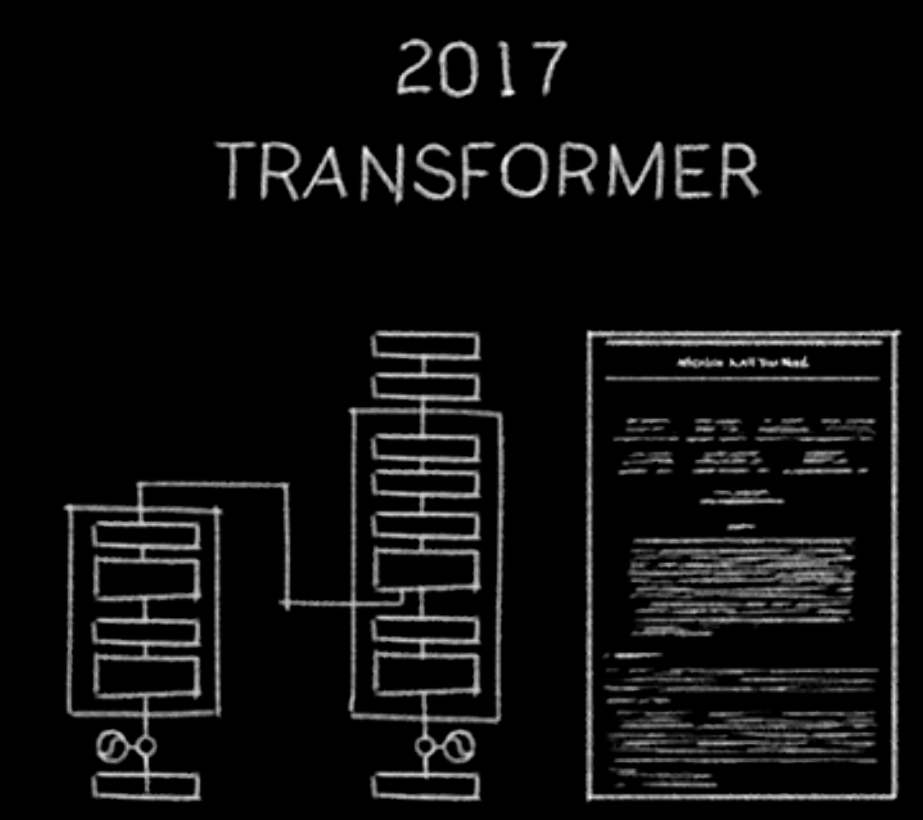






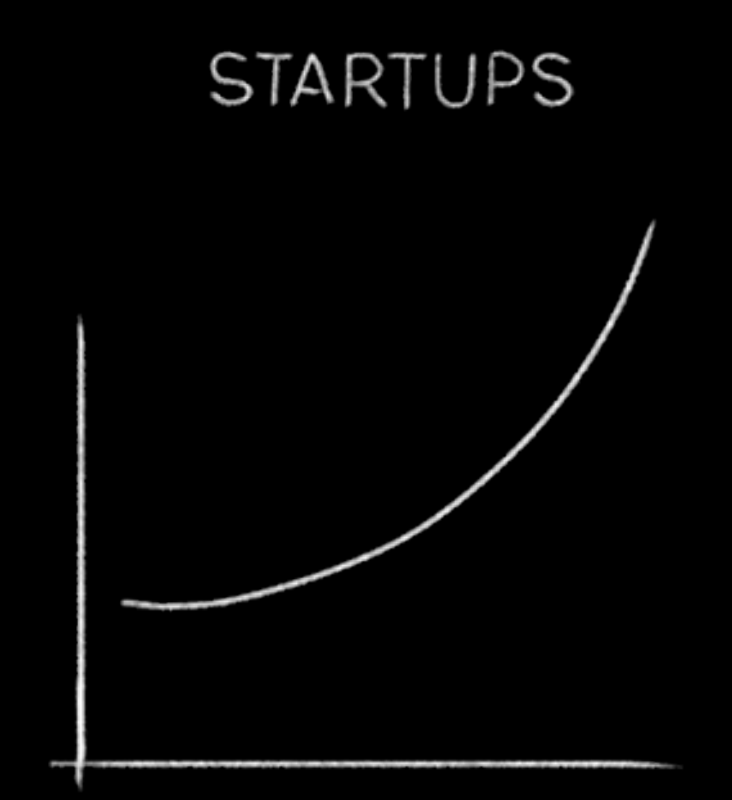
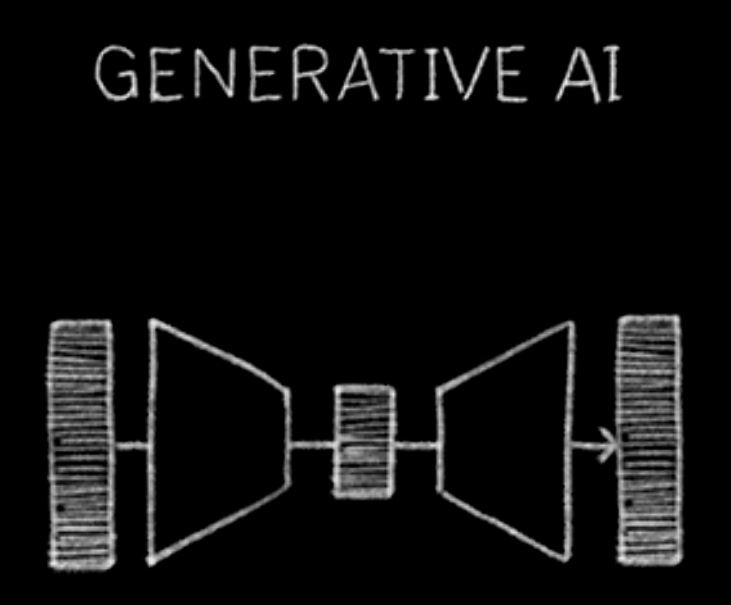
LEARN EVERYTHING

PROTEIN	LANGUAGE
SOUND	PHYSICS
3D	VIDEO
IMAGES	MANIPULATION
	GESTURE

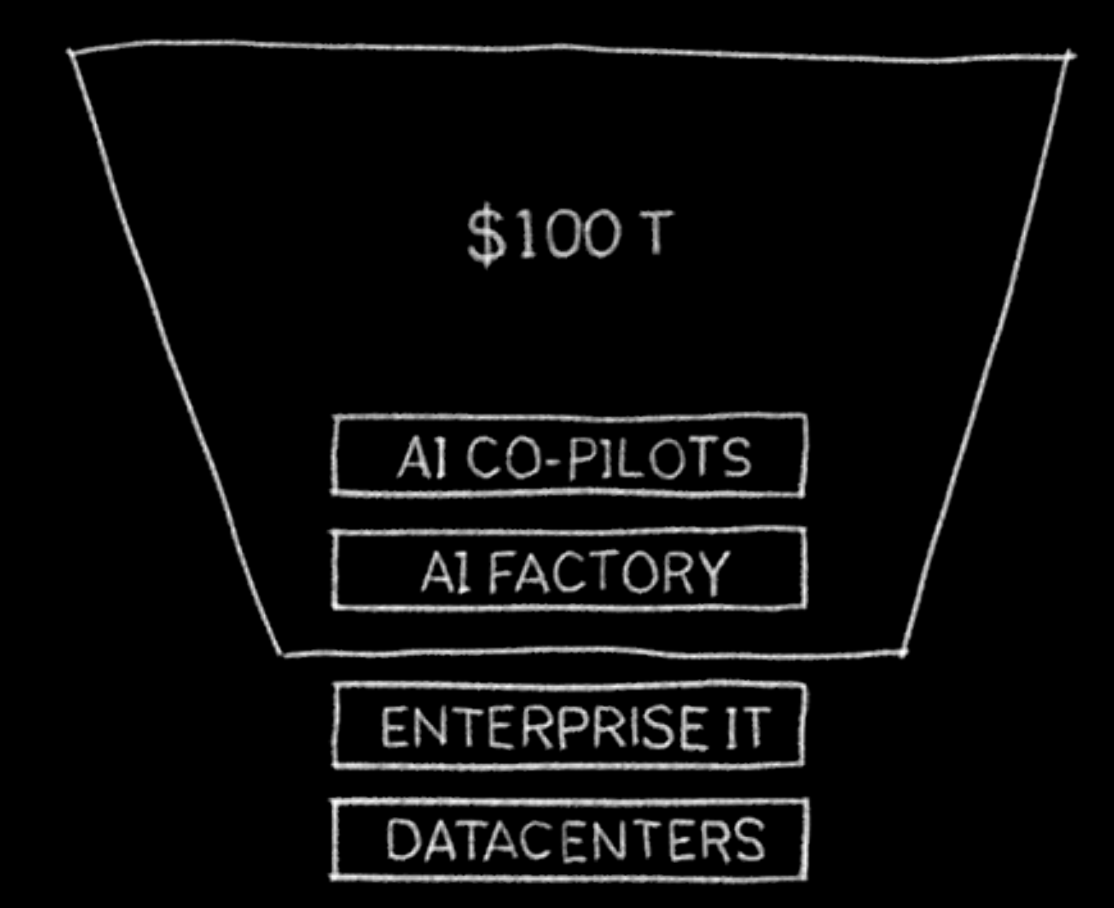


FINE TUNING
GUARDRAILING
ALIGNMENT
PROMPT ENGINEERING
VECTOR DB
COT & TOT

RAG
MULTI-MODAL
AGENTS



A NEW
INDUSTRIAL
REVOLUTION



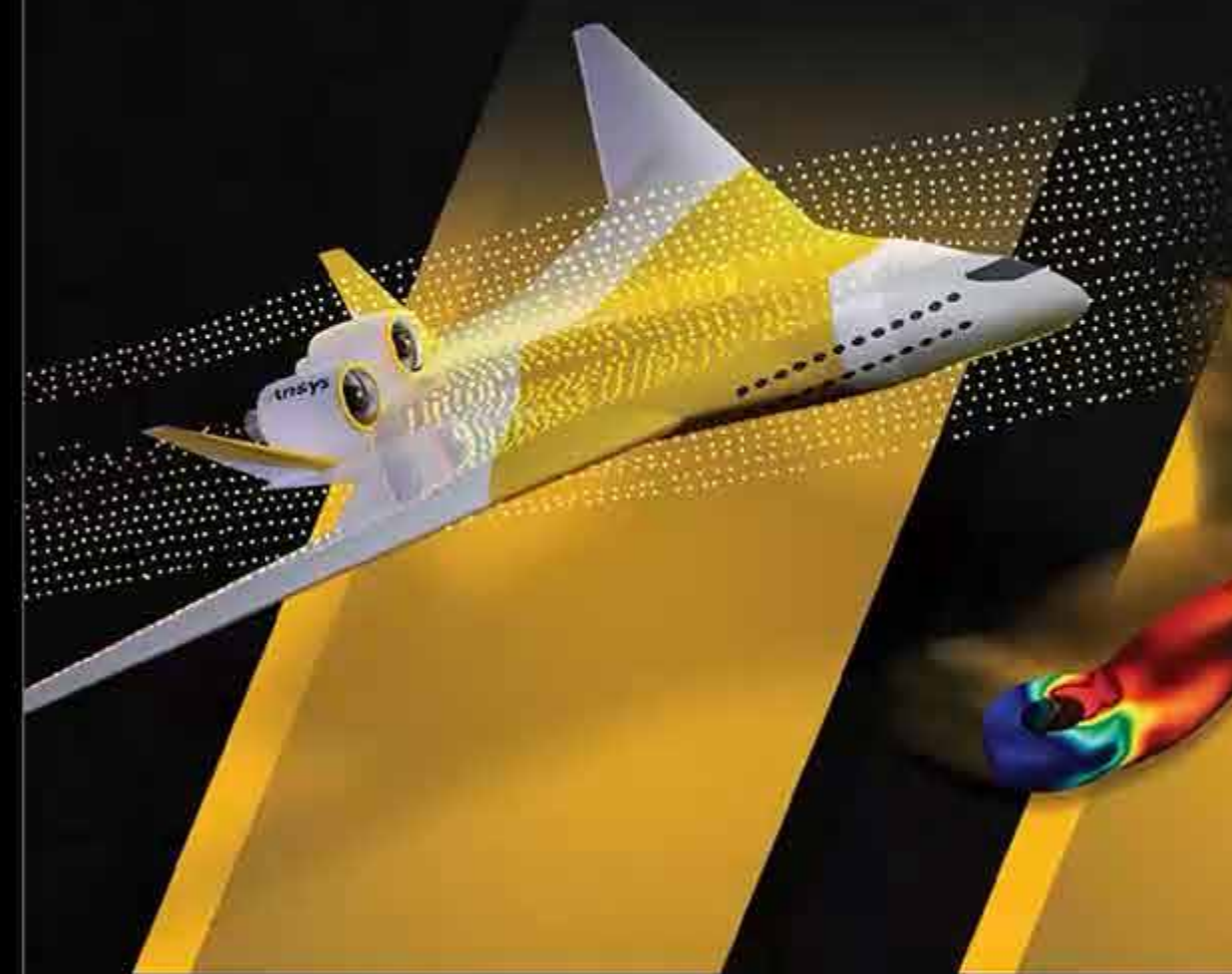




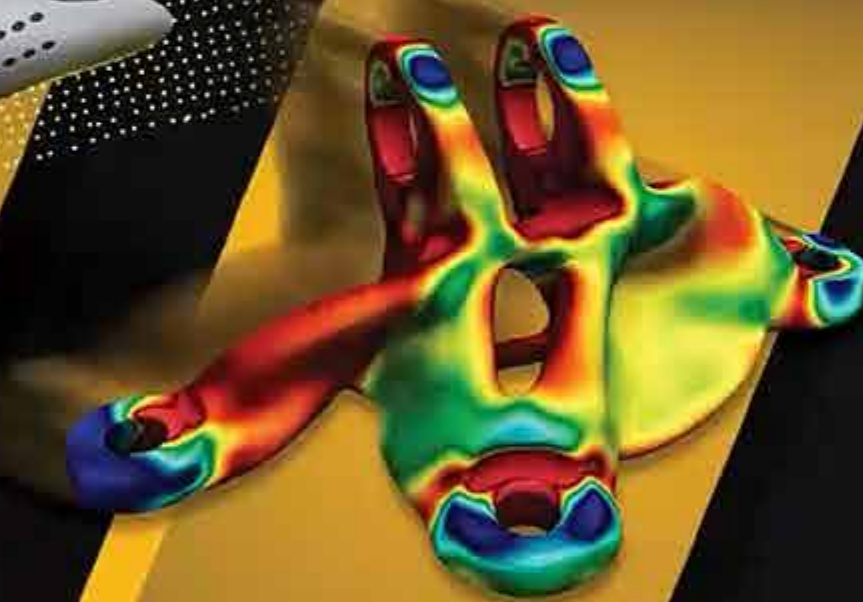
Ansys, Rescale



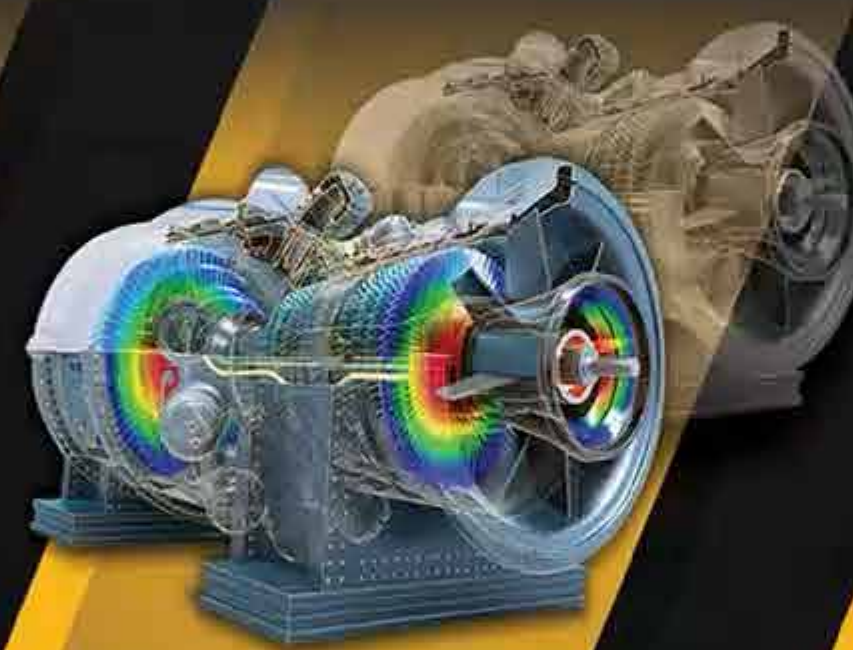
Building a New Era of CAE with Accelerated Computing & Generative AI



ANSYS SOLVERS
OPTIMIZED FOR
NVIDIA GPUS



NVIDIA AI ACCELERATED
ANSYS SIMULATION



PHYSICS-BASED
DIGITAL TWINS



NVIDIA POWERED
ANSYS LLMS



TRANSFORM 6G RESEARCH
WITH ANSYS PERCEIVE EM
INTEGRATION IN OMNIVERSE



ACCELERATED NVIDIA
DEVELOPMENT
WITH ANSYS

Ansys / Powering Innovation That Drives Human Advancement



SYNOPSYS: MISSION CRITICAL FOR NVIDIA SILICON SUCCESS

DECADES OF COLLABORATION ACROSS FULL EDA SUITE POWERS ACCELERATED COMPUTING

13X

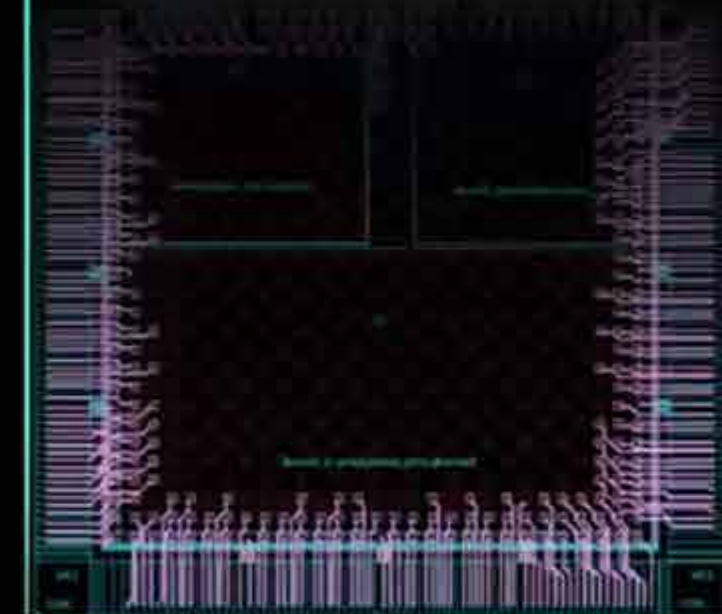
Verification Functional Verification



- Synopsys VCS
- NVIDIA L40
- NVIDIA Grace Hopper

10X

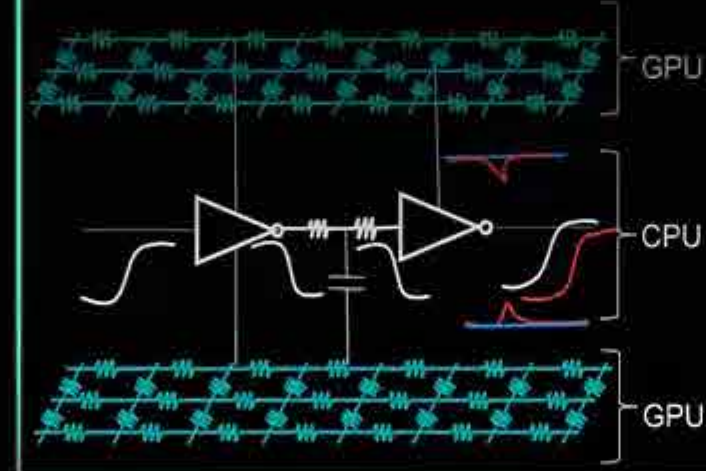
Design Place and Route



- Synopsys Fusion Compiler
- NVIDIA Grace Hopper

15X

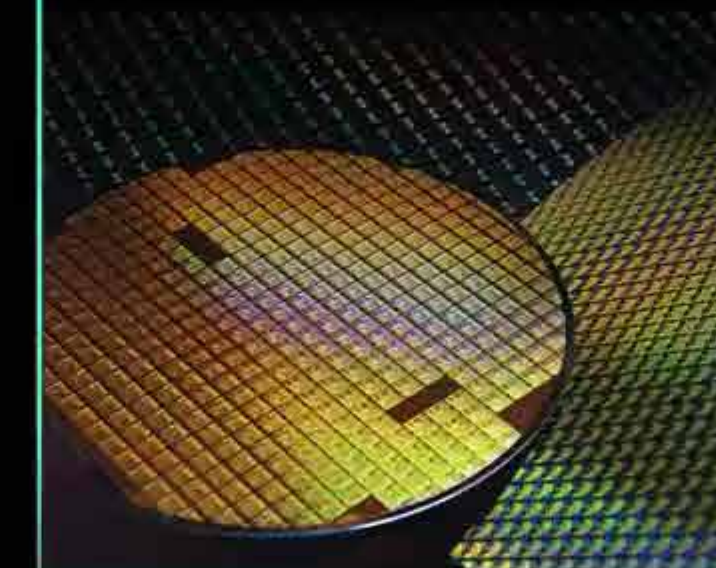
Simulation SPICE Simulation



- Synopsys PrimeSim
- NVIDIA Hopper
- NVIDIA Grace Hopper

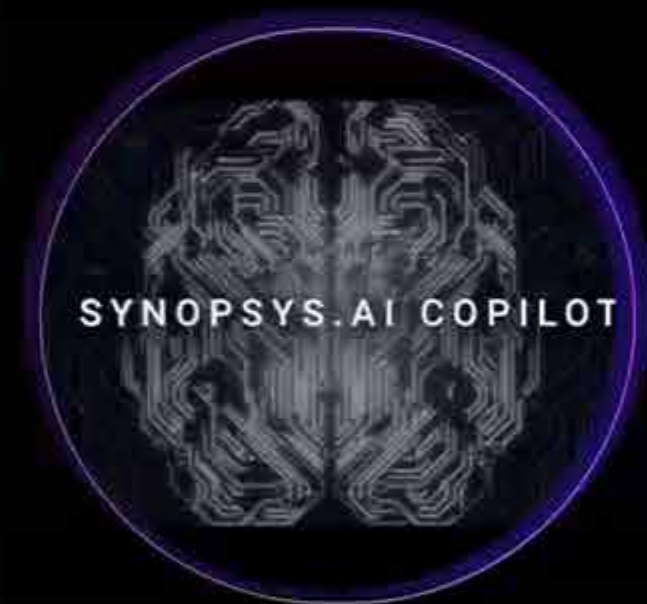
15X

Manufacturing Computational Lithography



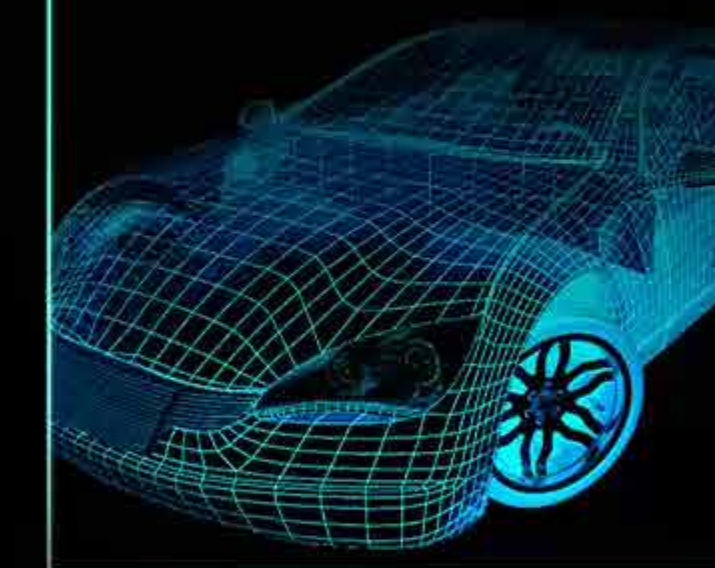
- Synopsys Proteus
- NVIDIA cuLitho
- NVIDIA Grace Hopper

Generative AI Industry's 1st LLM-Based GenAI EDA Solution



- Synopsys.ai
- NVIDIA NeMo & NIM
- NVIDIA DGX

Systems Software Testing & Validation of Automotive Software



- Synopsys Electronics
Digital Twin, vECU, TPT
- NVIDIA Omniverse

*Performance Speed-Up Based on Projected Results

cadence + NVIDIA

Deep collaboration on EDA, SDA, Digital Biology and AI

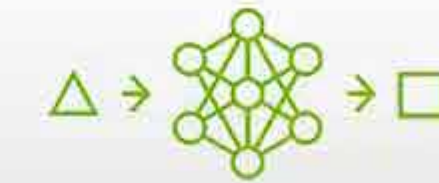
cadence.ai



NVIDIA NeMo



NVIDIA Modulus



NVIDIA BioNeMo

Cerebrus

Virtuoso Studio

Verisium

Optimality

Allegro X AI

Chip NeMo

EDA



Digital

Custom

Verification

SDA



3D-IC

PCB

Multiphysics

OMNIVERSE



Environmental &
Data Center Digital Twin

BIO



Molecular
Sciences



Cadence Palladium
& Protium



NVIDIA
HGX H100



Cadence
Millennium



NVIDIA
OVX L40



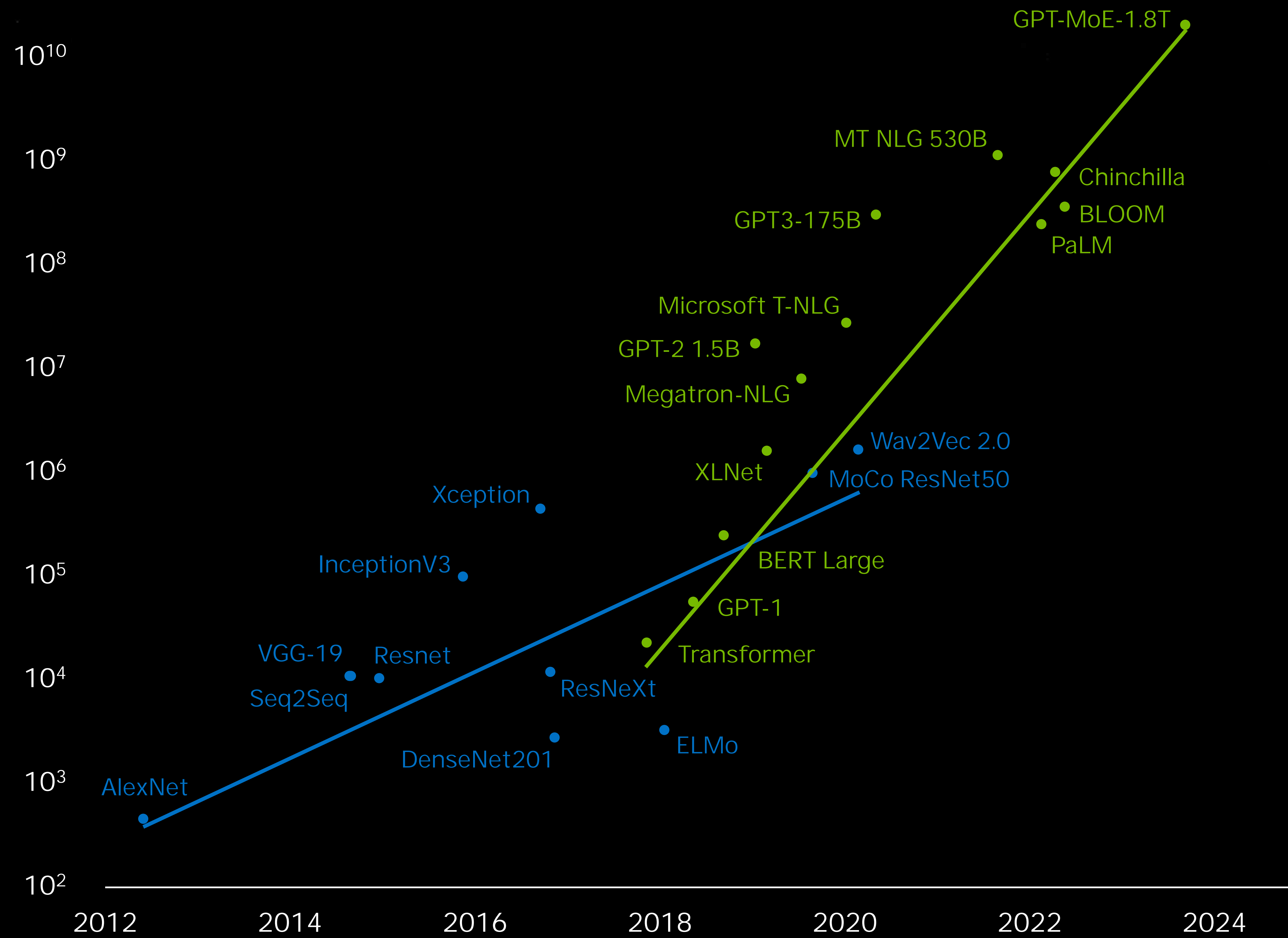
NVIDIA Grace Hopper &
Grace Blackwell

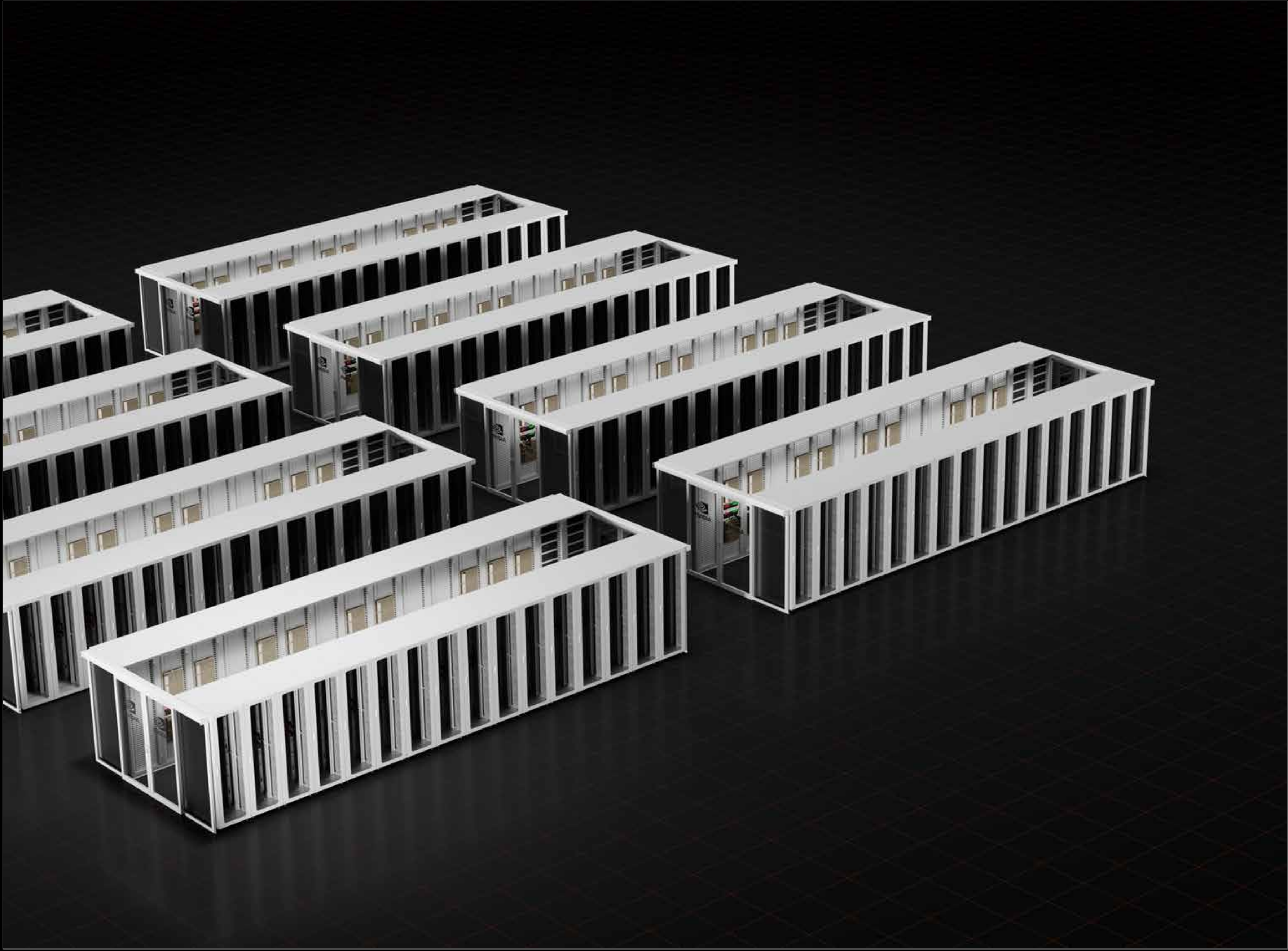
AI
GenAI Co-Pilots

Principled
Simulation +
Optimization

Accelerated
Compute

Training Compute PFLOPs





SELENE
2021

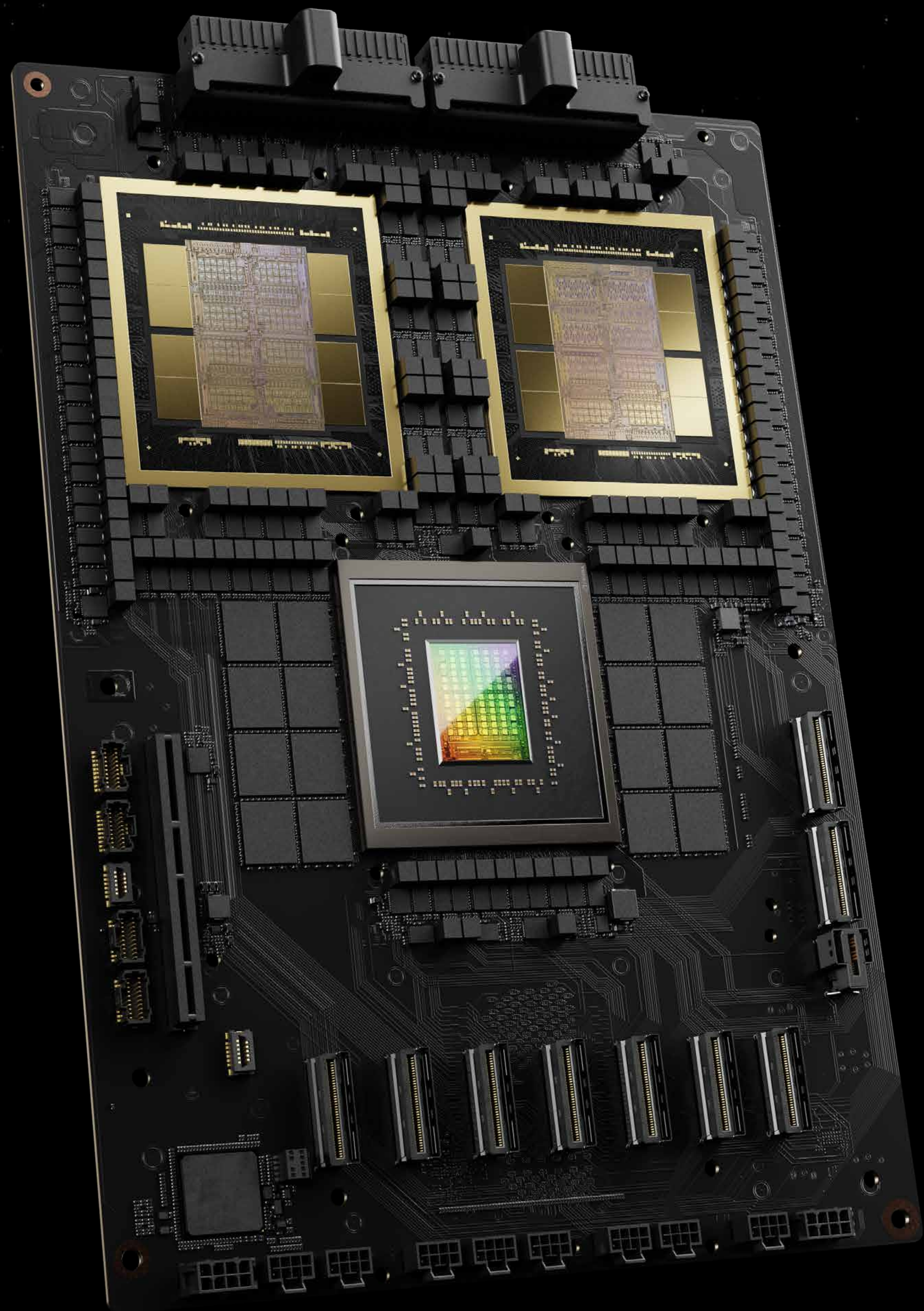
4,480 A100 GPUs
3 EF AI Compute
112 TB/s Interconnect BW



EOS
2023

10,752 H100 GPUs
43 EF AI Compute
1,100 TB/s Interconnect BW





ANNOUNCING NVIDIA BLACKWELL PLATFORM FOR TRILLION-PARAMETER SCALE GENERATIVE AI



AI SUPERCHIP
208B Transistors



2nd GEN TRANSFORMER ENGINE
FP4/FP6 Tensor Core



5th GENERATION NVLINK
Scales to 576 GPUs



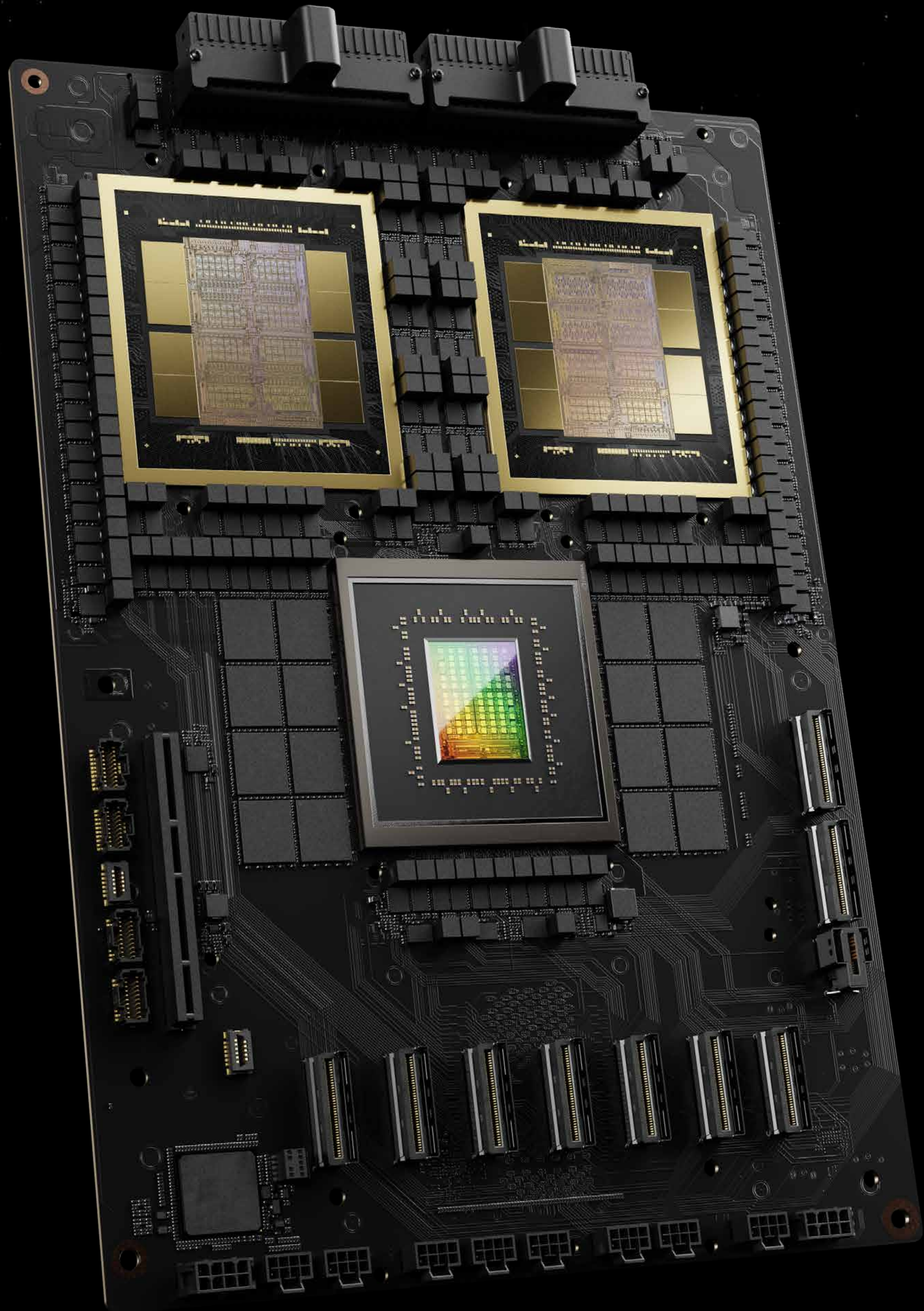
RAS ENGINE
100% In-System Self-Test



SECURE AI
Full Performance
Encryption & TEE

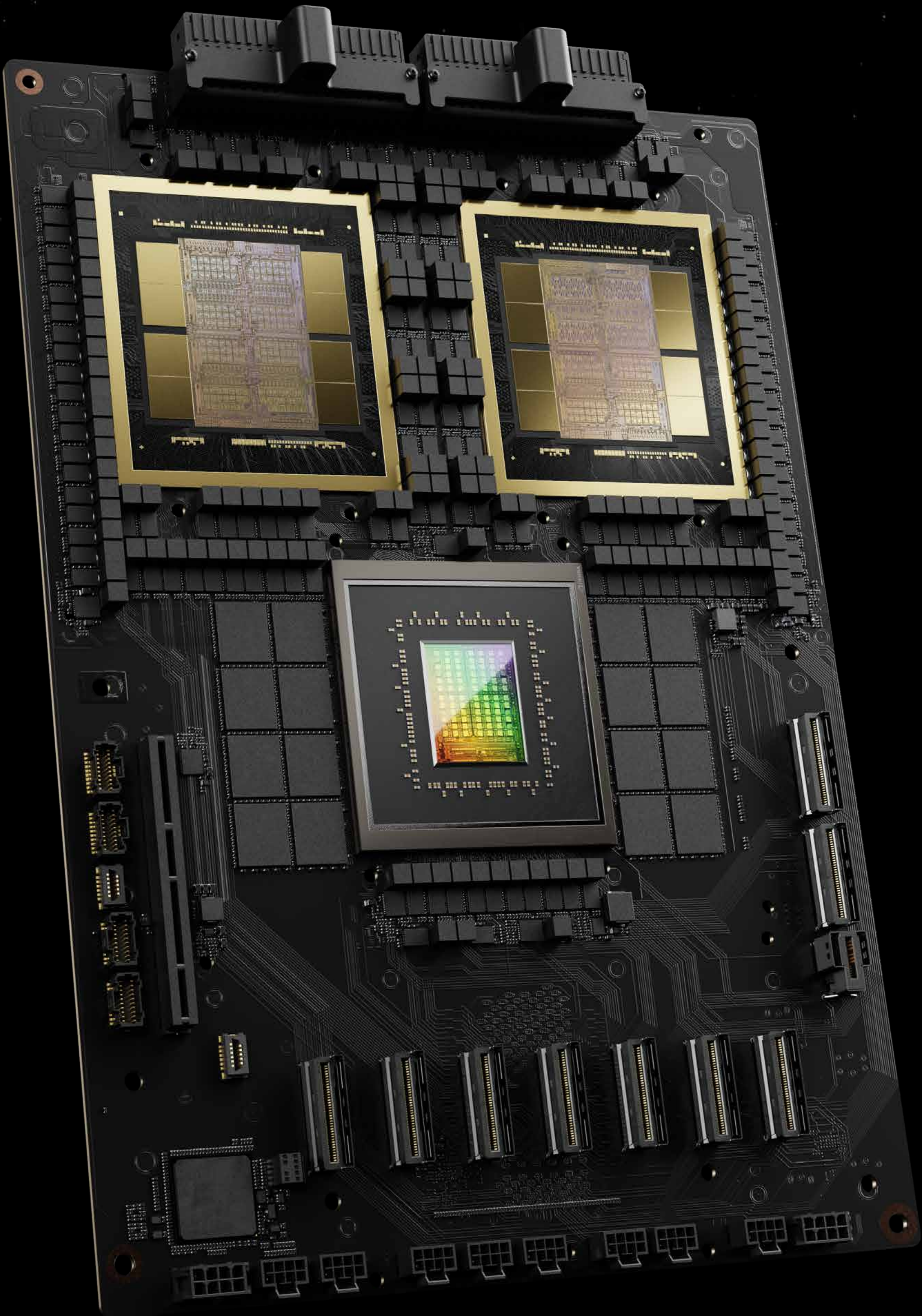


DECOMPRESSION ENGINE
800 GB/sec

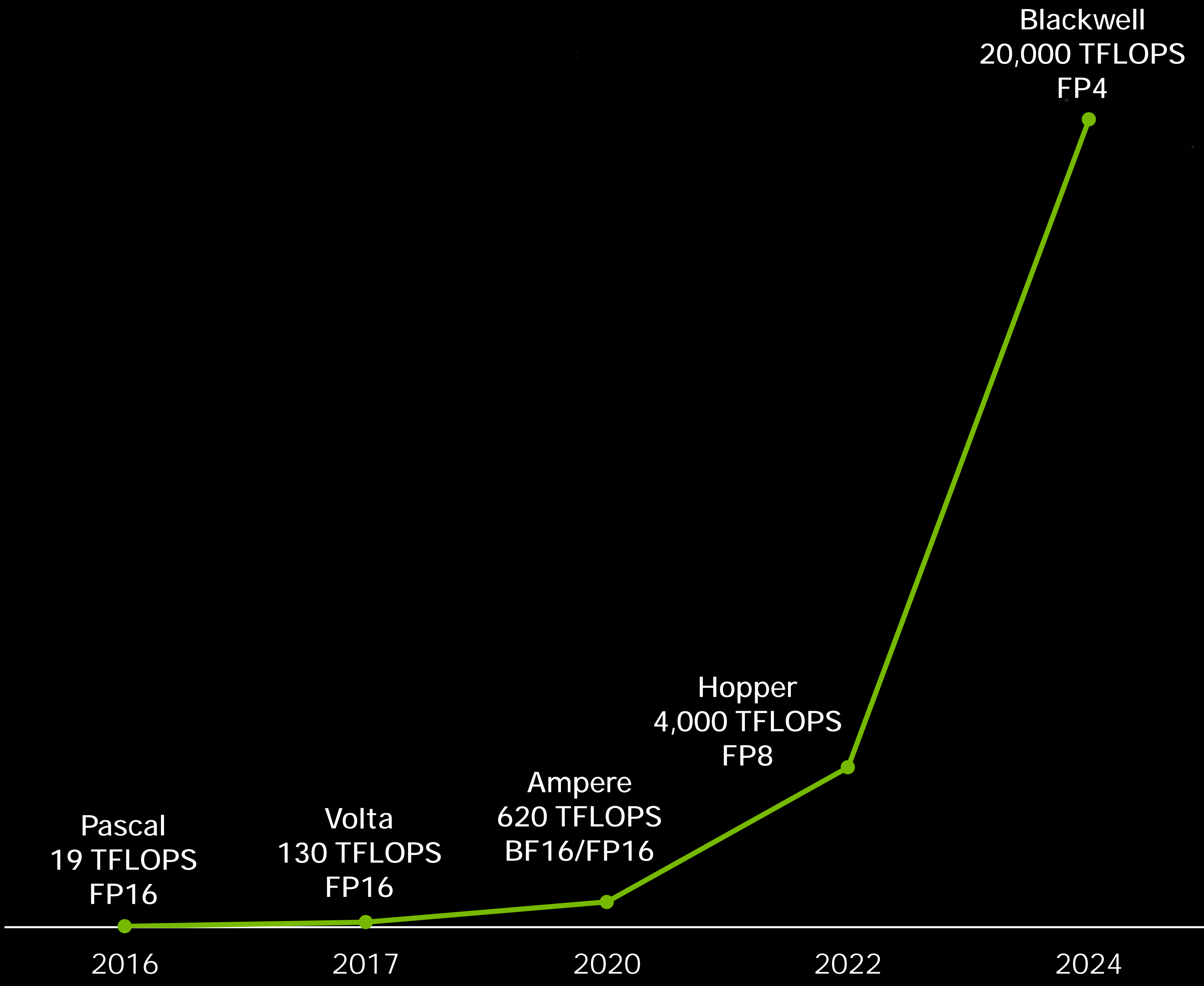


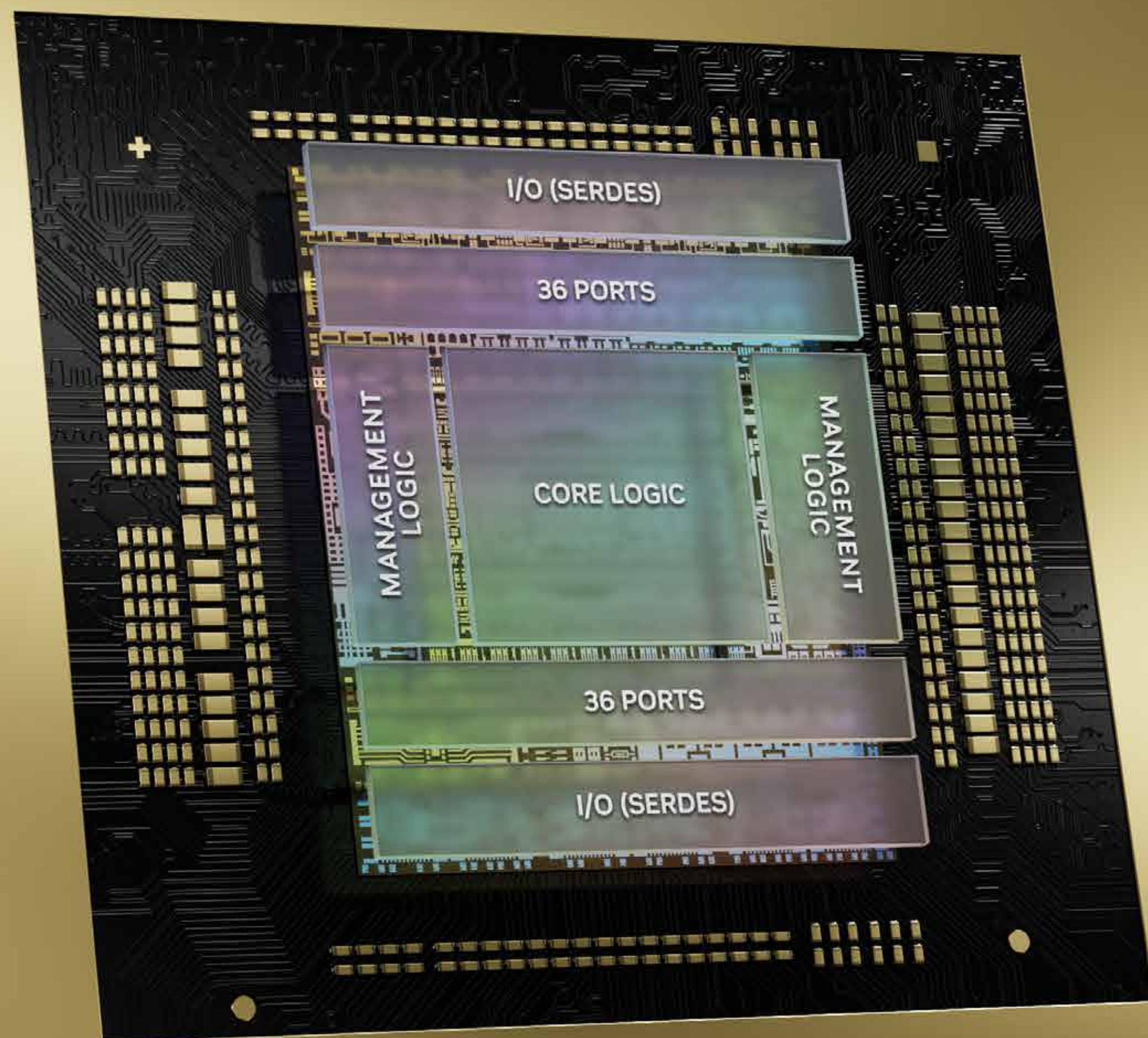
Blackwell GPU

FP8	20 PFLOPS	2.5X Hopper
NEW FP6	20 PFLOPS	2.5X
NEW FP4	40 PFLOPS	5X
HBM Model Size	740B param	6X
HBM Bandwidth	34T param/sec	5X
NVLINK All-Reduce with SHARP	7.2 TB/s	4X



1000X AI Compute in 8 Years





NVLink Switch Chip

50B Transistors in TSMC 4NP

72-Ports Dual 200 Gb/sec SerDes

4 NVLinks at 1.8TB/sec

7.2TB/sec Full-Duplex Bandwidth

SHARP In-Network Compute - 3.6 TFLOPS FP8



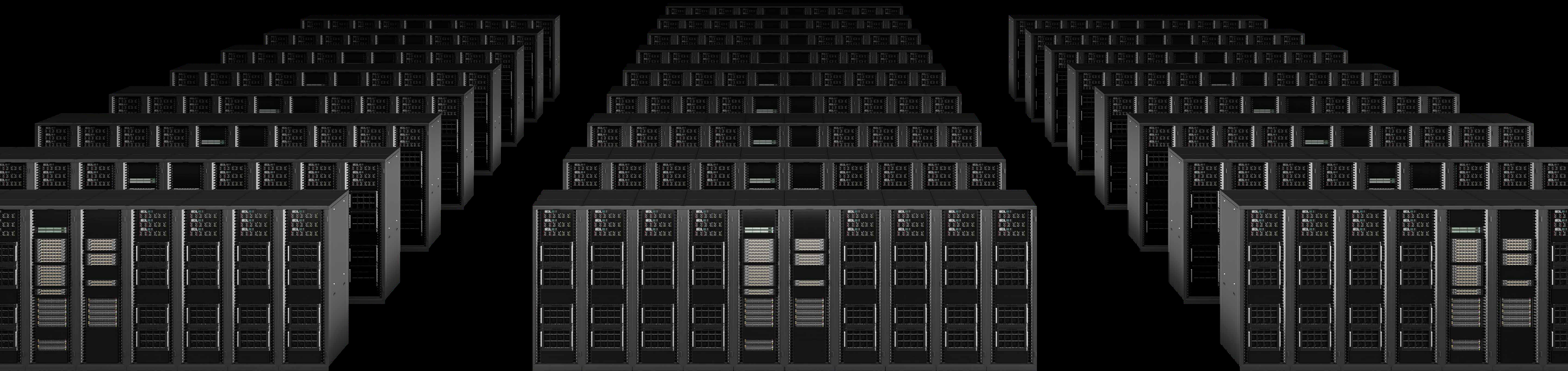
DGX GB200 NVL72
1 Giant GPU

Training FP8	720 PFLOPS	22X
Inference FP4	1.44 ExaFLOPS	45X
Multi-Node All-to-All	130 TB/sec	18X
Multi-Node All-Reduce	260 TB/sec	36X



Train GPT-MoE-1.8T in 90 Days

Hopper
8000 GPUs | 15MW



Train GPT-MoE-1.8T in 90 Days

Blackwell GB200 NVL72
2000 GPUs | 4MW

1/4th the Power



GPT-MoE 1.8T
Inference (seqlen=32k/1k, FTL=5s)

Throughput per GPU
Tokens per Second

160

140

120

100

80

60

40

20

0

0

10

20

30

40

50

Interactivity per User Tokens per Second

Multi-Dimensional Optimization:

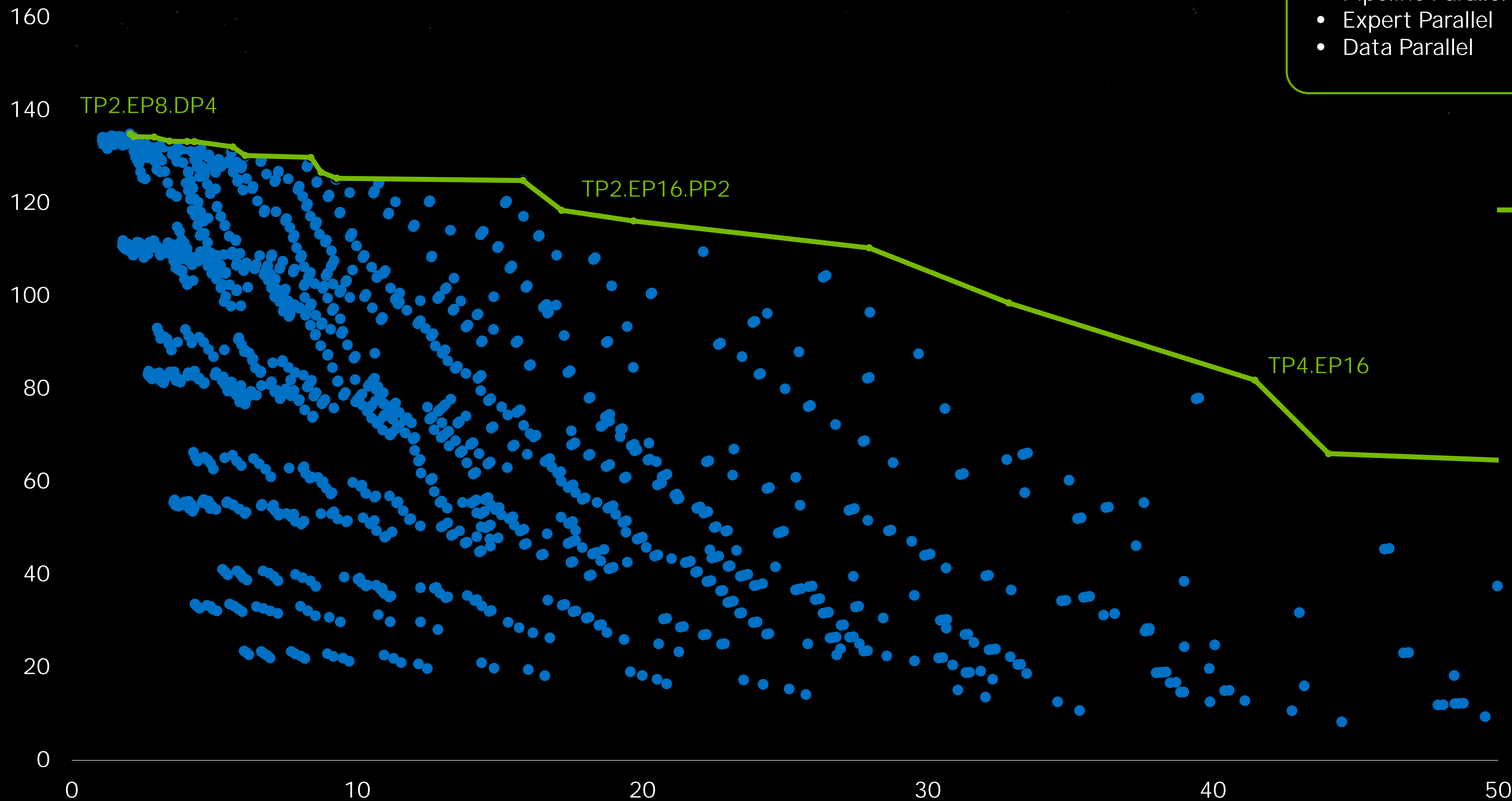
- Tensor Parallel
- Pipeline Parallel
- Expert Parallel
- Data Parallel

TP2.EP8.DP4

TP2.EP16.PP2

TP4.EP16

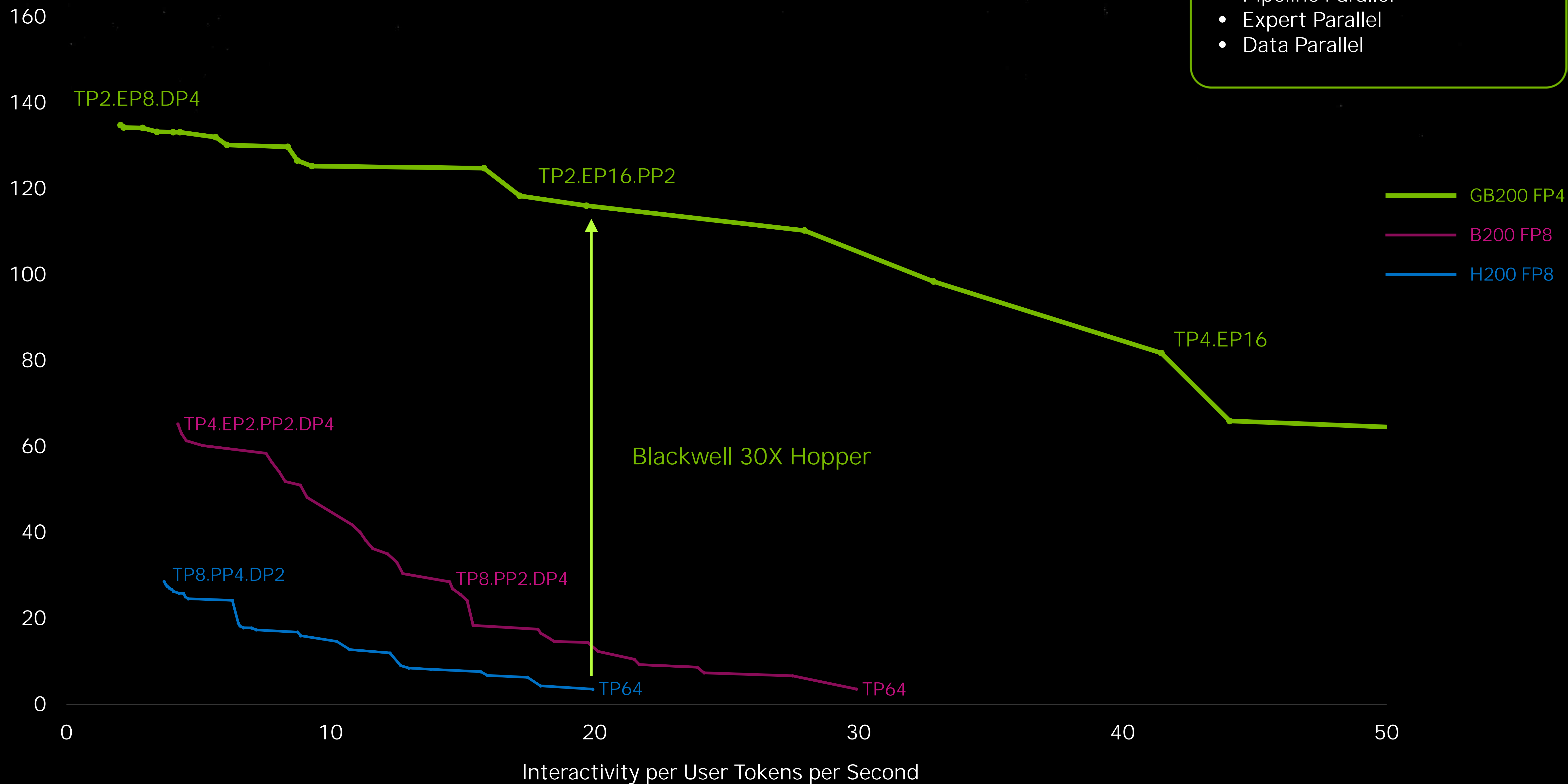
GB200 FP4



GPT-MoE 1.8T Inference (seqlen=32k/1k, FTL=5s)

Throughput per GPU
Tokens per Second

- Multi-Dimensional Optimization:
- Tensor Parallel
 - Pipeline Parallel
 - Expert Parallel
 - Data Parallel





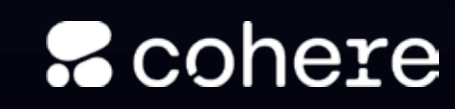
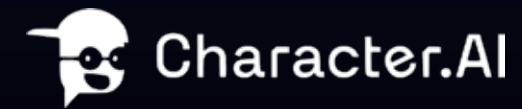
Google Cloud



ORACLE
CLOUD
Infrastructure

ADEPT

AI21labs



Inflection



together.ai



AIVRES



ASUS



DELL Technologies

EVIDEN



FUJITSU

GIGABYTE



IBM Cloud



Inventec



Lenovo



PEGATRON



SoftBank



wlstron





AWS and NVIDIA expand joint AI offerings

OVER A DECADE OF CO-INNOVATION IN THE CLOUD

AWS AI Security

NVIDIA BLACKWELL COMING TO AWS

AWS NITRO + KMS-ENCRYPTED EFA

NVIDIA DGX Cloud

GB200

Amazon EC2

B100



Project Ceiba

AWS-BUILT AI SUPERCOMPUTER FOR NVIDIA INTERNAL R&D

414 EXAFLOPS | 20,736 GB200 | GEN 4 AWS EFA

Optimizing price performance of FMs running on accelerated computing



Transforming industries with the power of AWS and NVIDIA

ADVANCED ROBOTICS SIMULATION



AI-DRIVEN DRUG DISCOVERY





Google Cloud

Data science & analytics

Dataproc

Dataflow

RAPIDS

AI platforms & frameworks

Vertex AI

GKE

DGX Cloud

JAX

XLA

NVIDIA AI Enterprise

TensorRT

Triton

NeMo

Accelerated computing

A3 (H100)

G2 (L4)

Blackwell

Google DeepMind



ORACLE



Sovereign
AI

Enterprise
Gen AI

Healthcare
AI

ORACLE
Database

ORACLE
Cloud Applications

ORACLE
CLOUD
Infrastructure

ORACLE CLOUD
Infrastructure Generative AI

ORACLE CLOUD
Infrastructure Supercluster

ORACLE
Distributed Cloud





Powering the most sophisticated AI innovators

ADEPT



Inflection



Largest NVIDIA InfiniBand connected supercomputer



Most deeply integrated cloud for NVIDIA DGX Cloud



Microsoft and NVIDIA Announce



NVIDIA Clara integration with Azure for healthcare



NVIDIA NIM available on Azure



New Azure VM series with NVIDIA H100 NVL



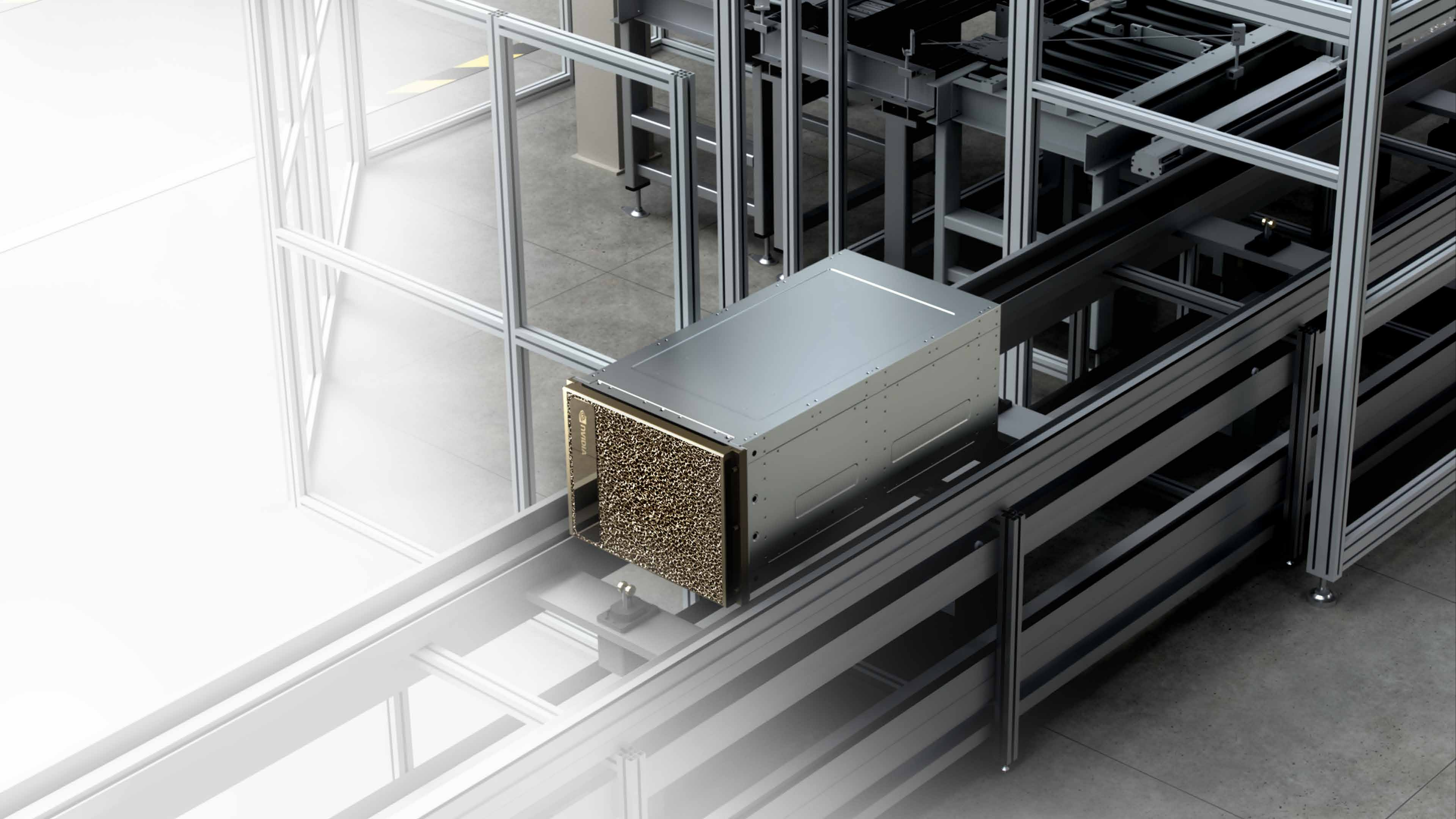
Microsoft Azure adopts NVIDIA GB200 for customers and AI services



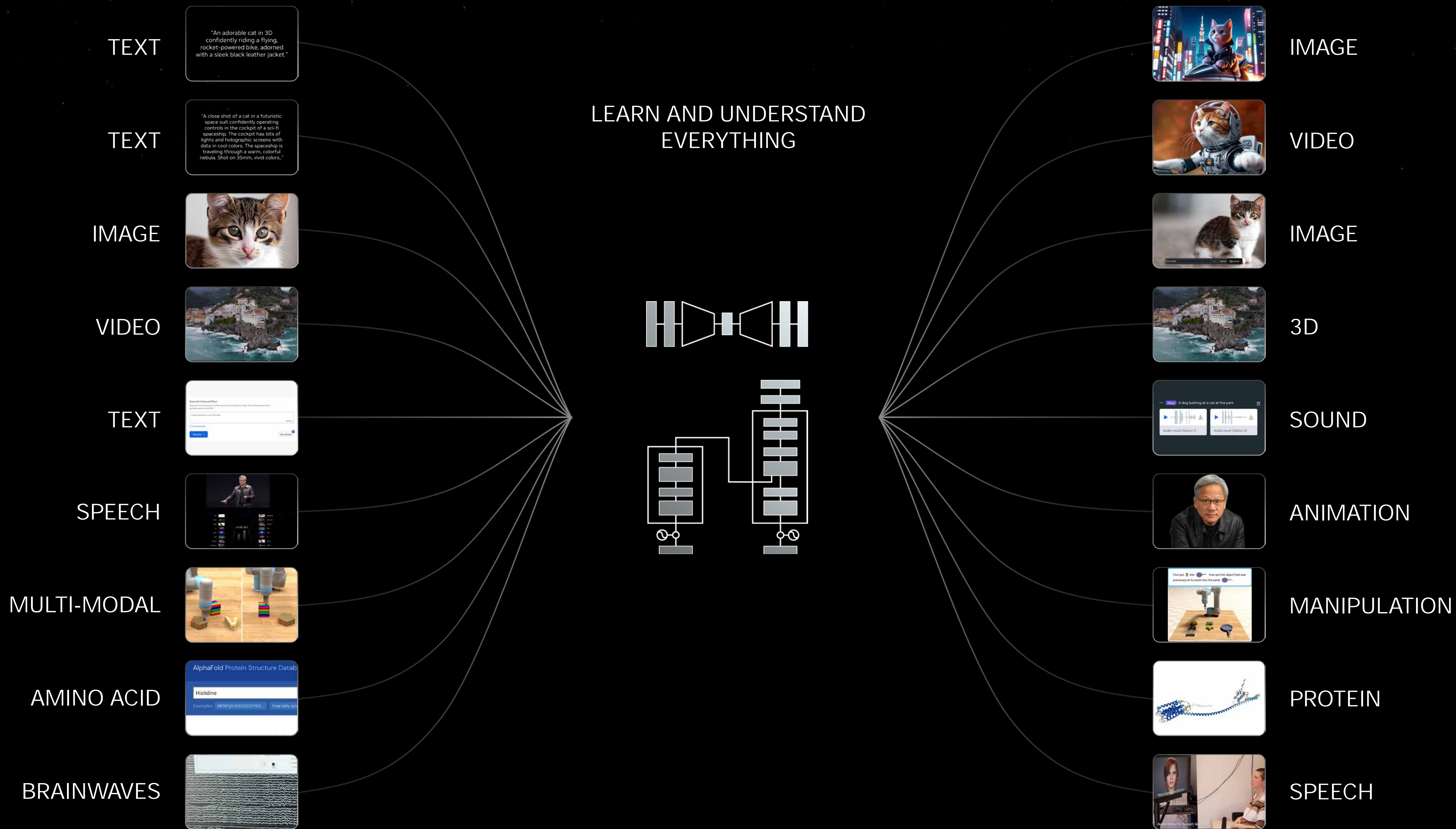
NVIDIA DGX Cloud Integration with Microsoft Fabric



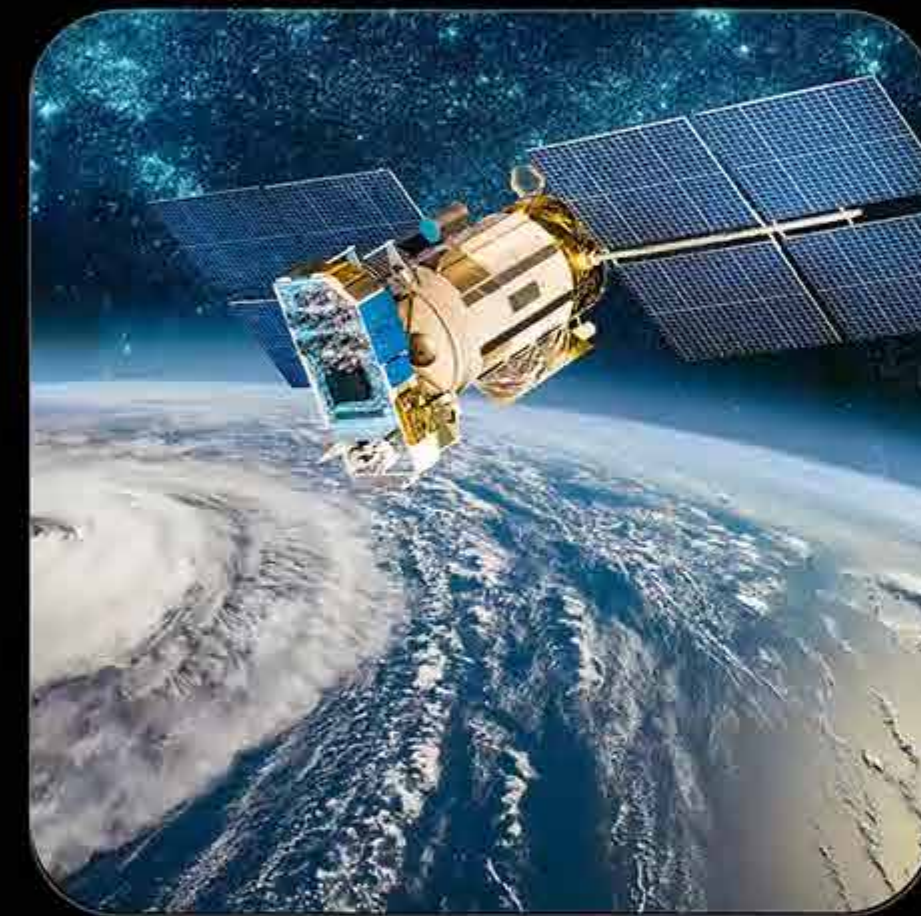
NVIDIA Omniverse APIs on Azure











TWCo PROCESSING 400 TB OF WEATHER DATA PER DAY

One of the largest collections of global weather data, producing 25B+ forecasts per day for 2.2B locations globally



TWCo GRAF FORECAST MODELLING

Next-generation AI forecast modeling



TWCo WEATHER ENGINE

Providing actionable insights to improve social resilience to extreme weather and make weather a competitive advantage for businesses

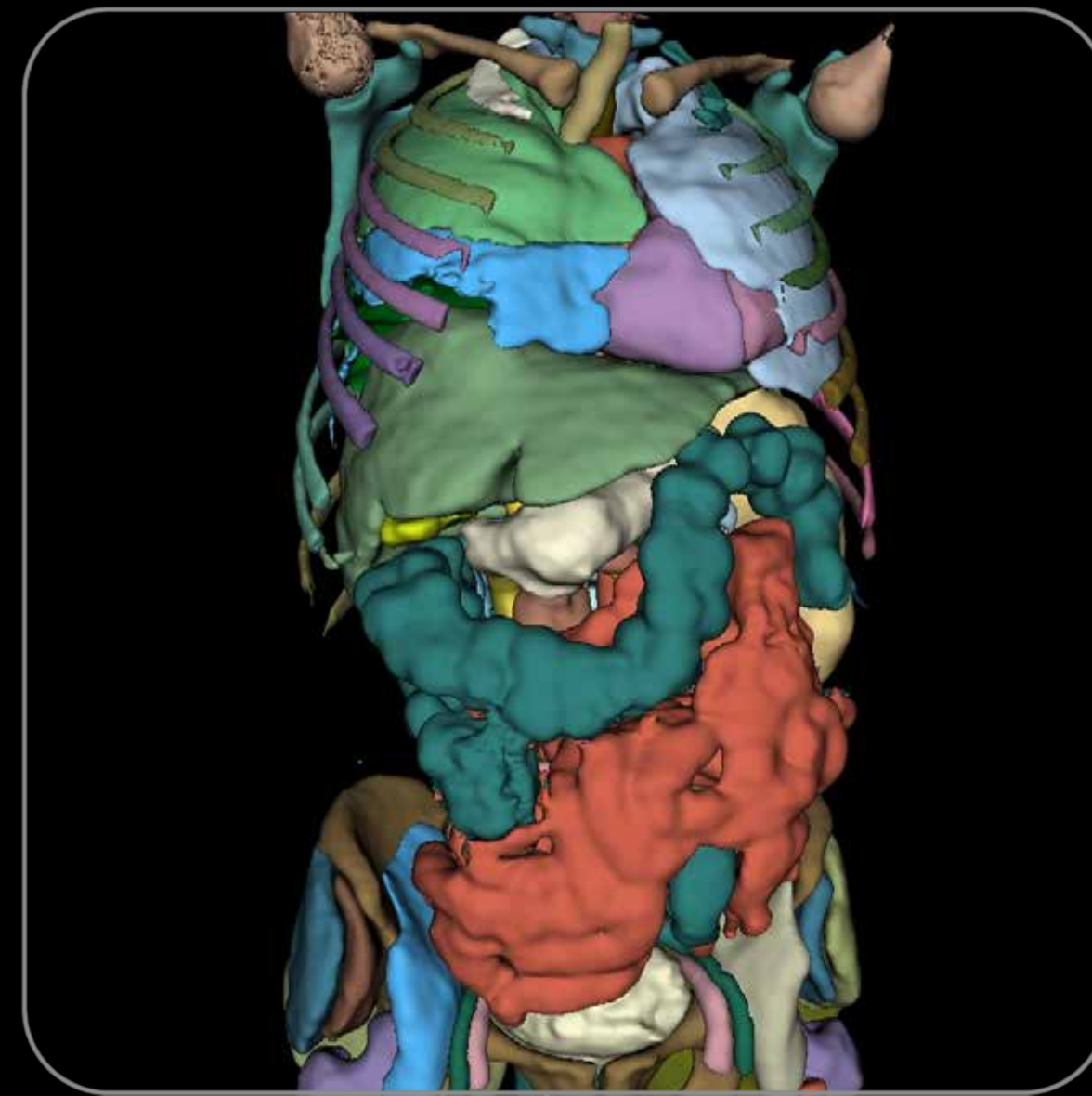


TWCo WEATHERVERSE

Advanced visualization for better decision making

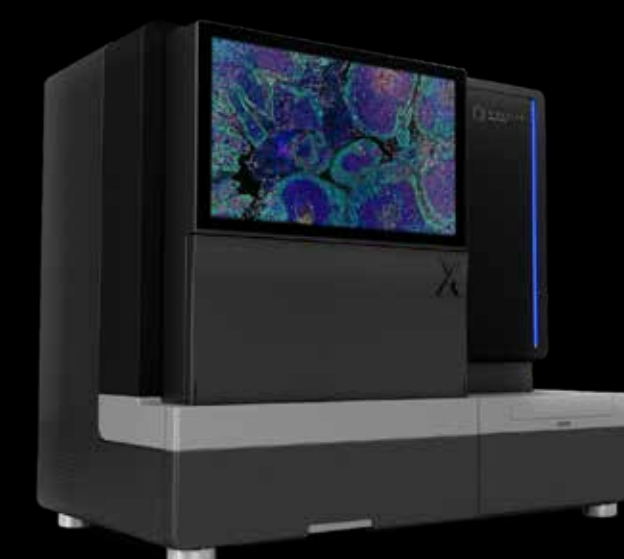
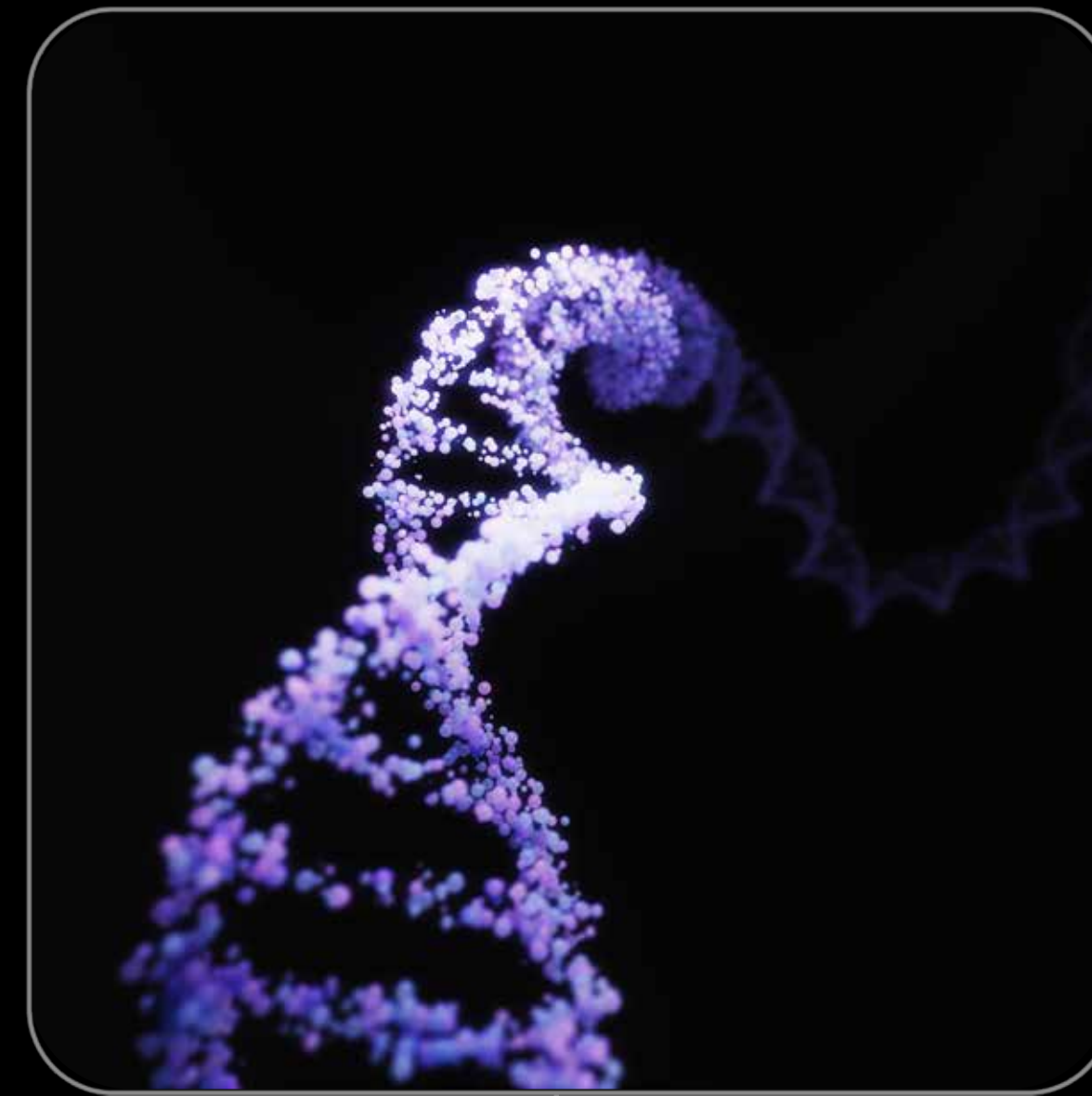
NVIDIA HEALTHCARE

IMAGING & ROBOTICS



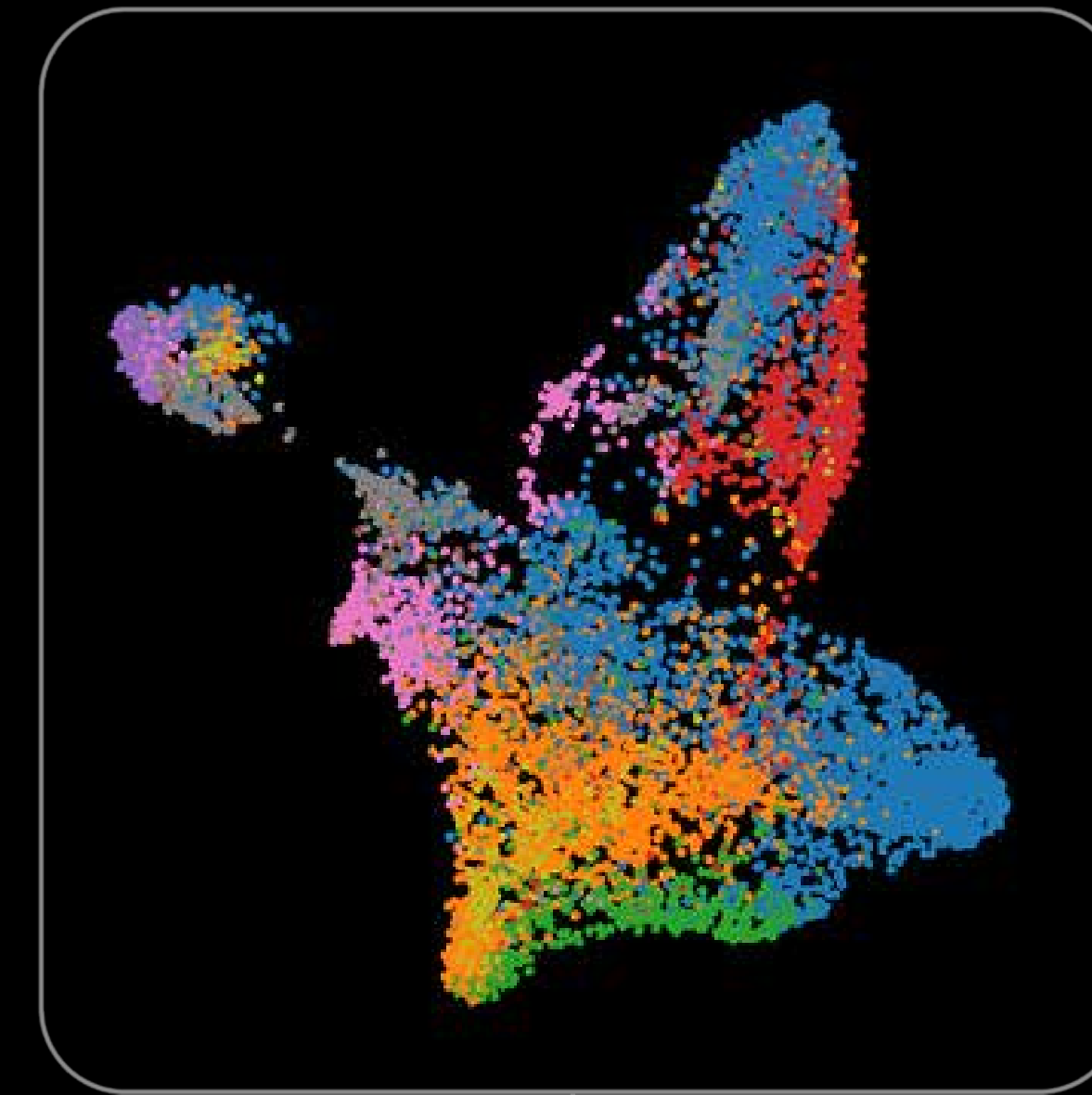
MONAI | Holoscan

GENOMICS



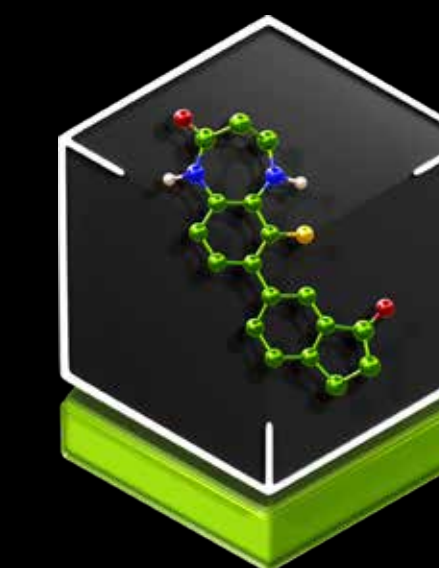
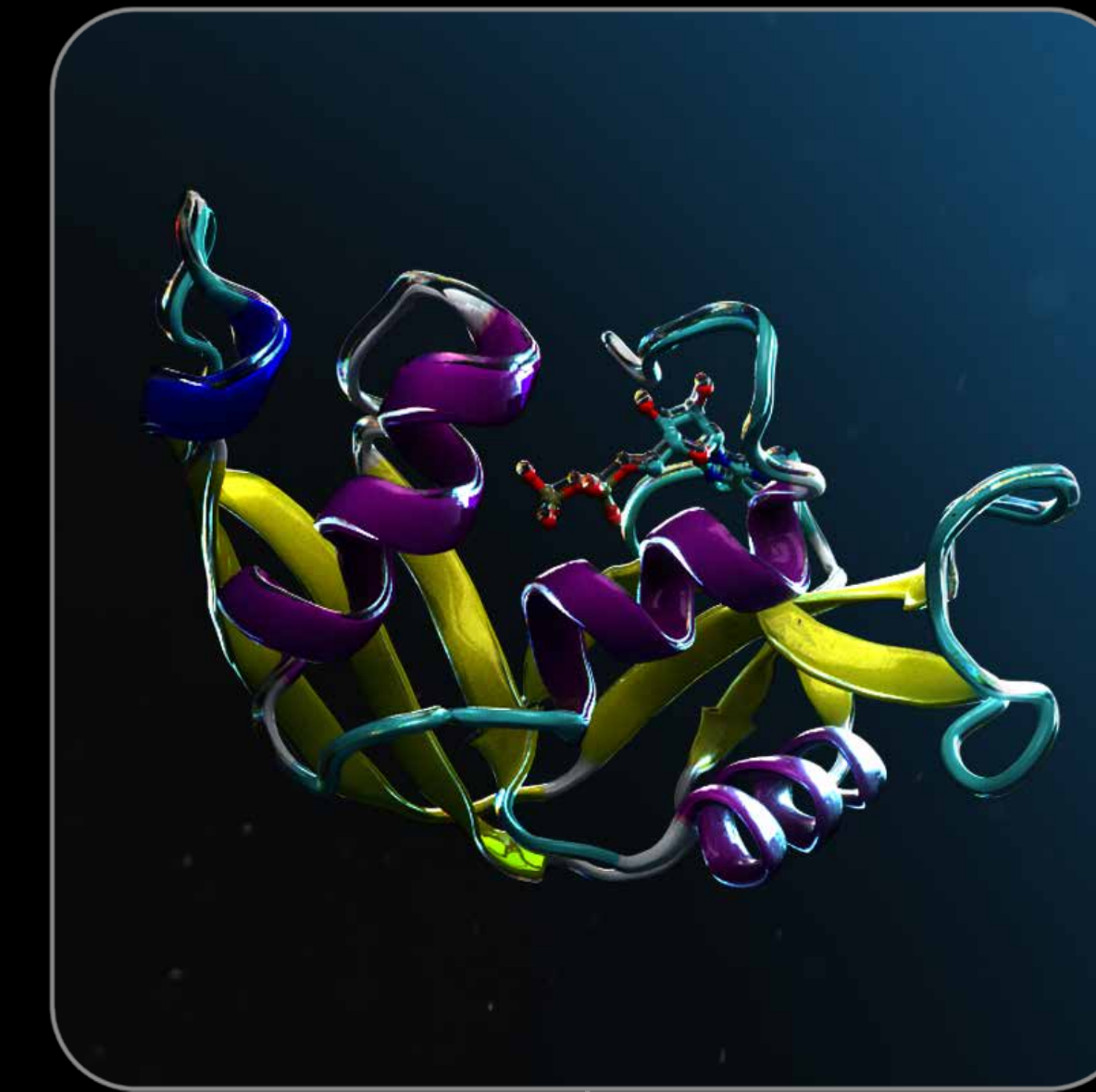
Parabricks

BIOINFORMATICS

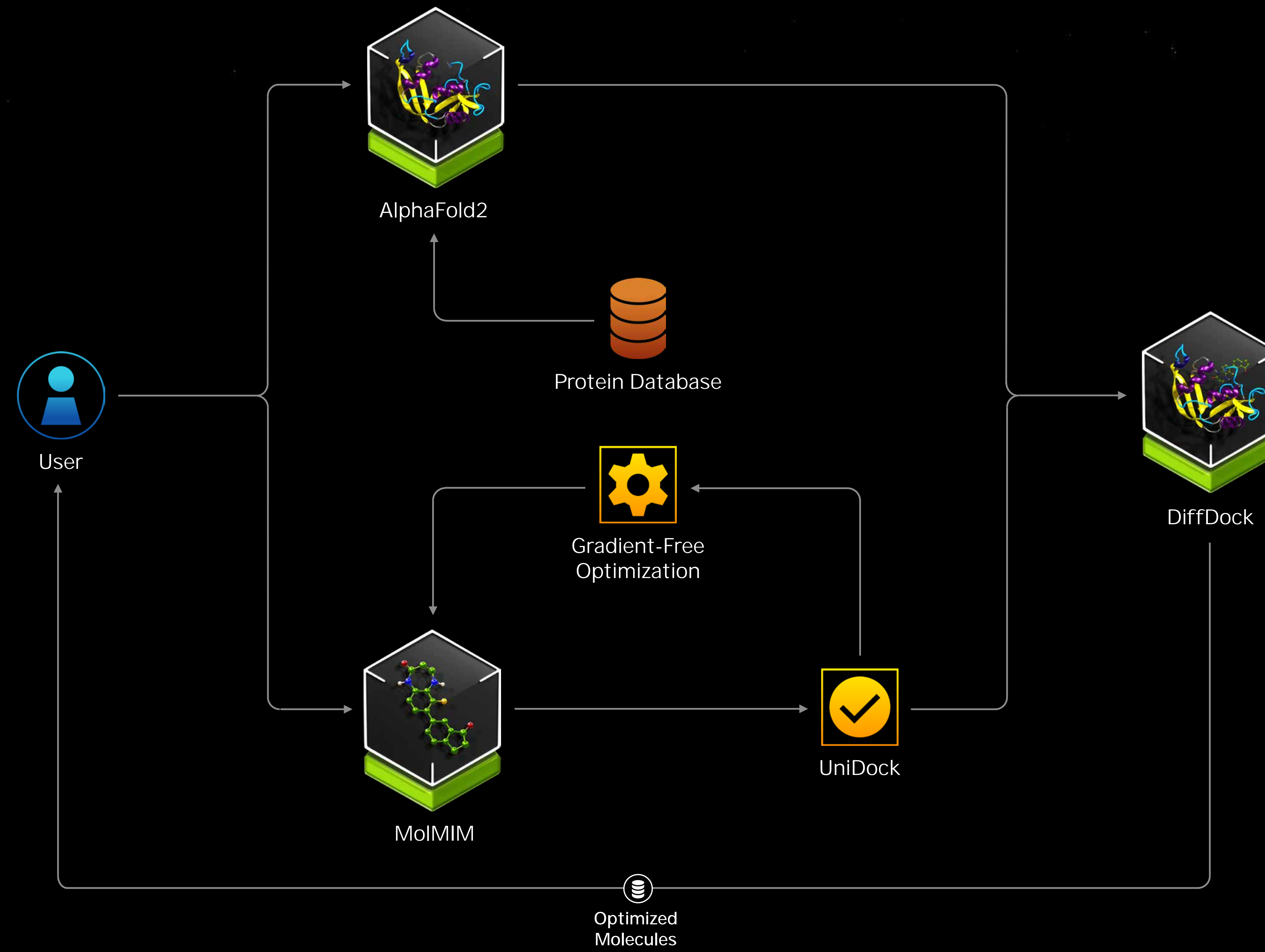


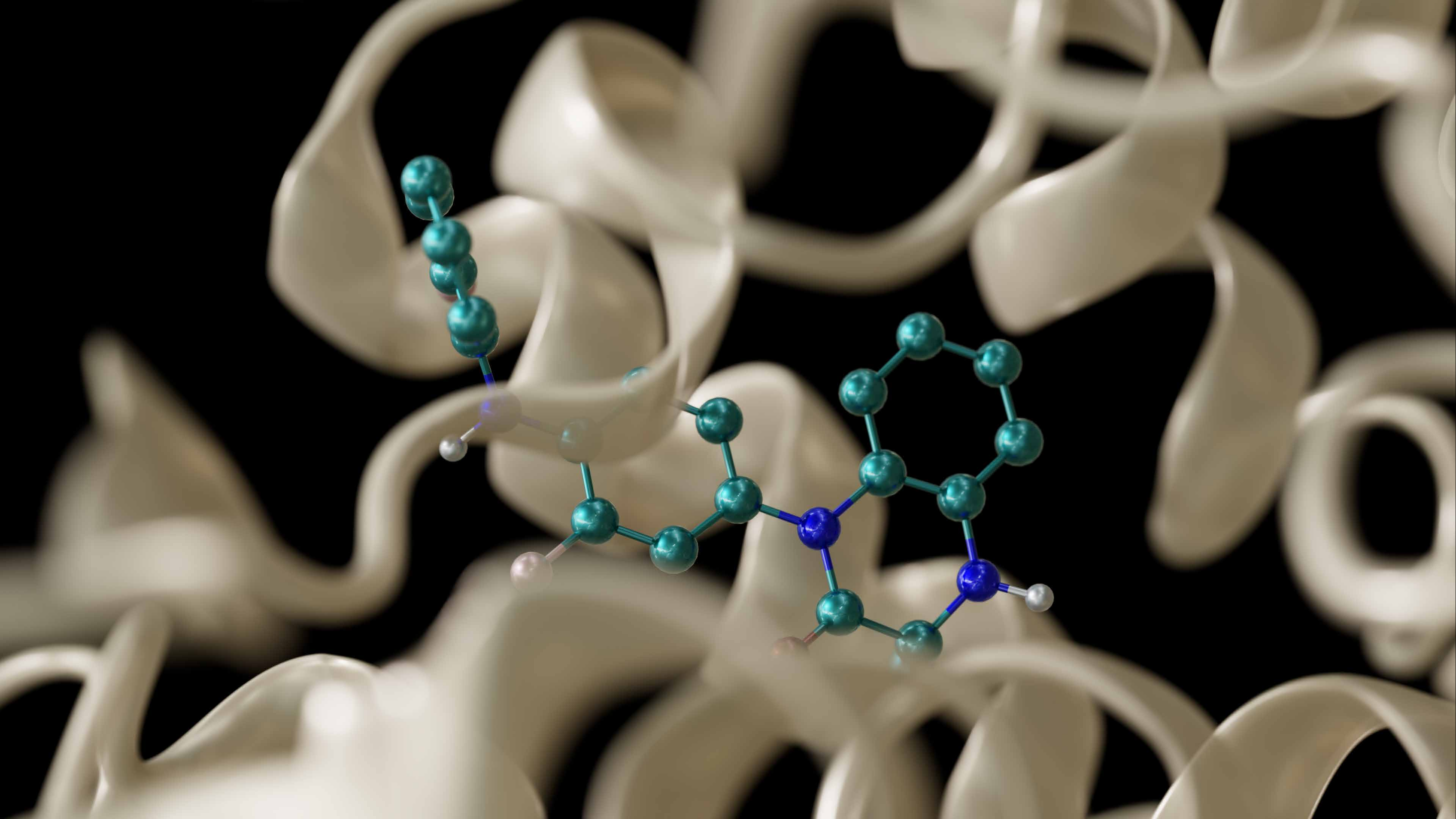
RAPIDS

DRUG DISCOVERY



BioNeMo





Google

Google DeepMind



ADEPT

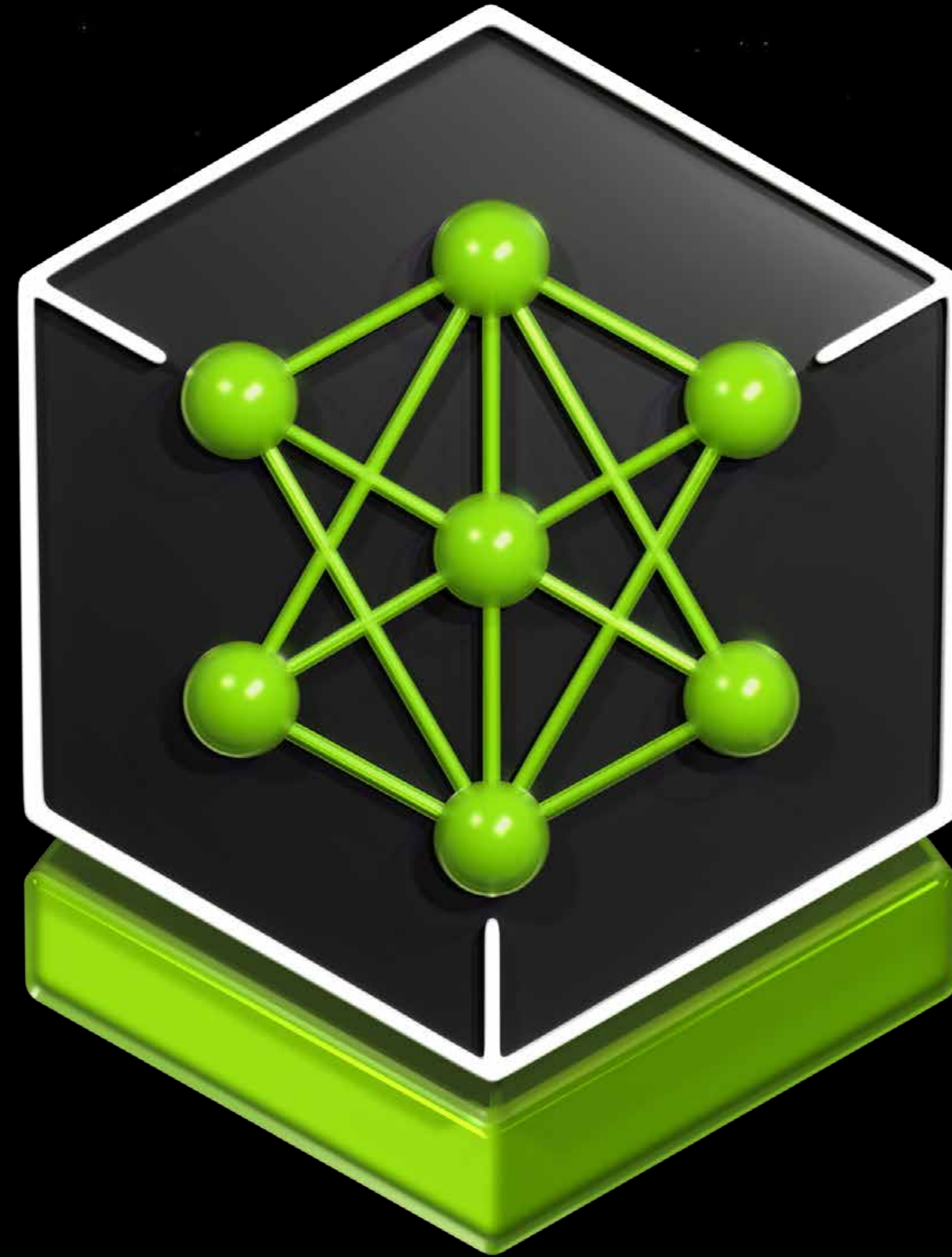
AI21 labs



Inflection

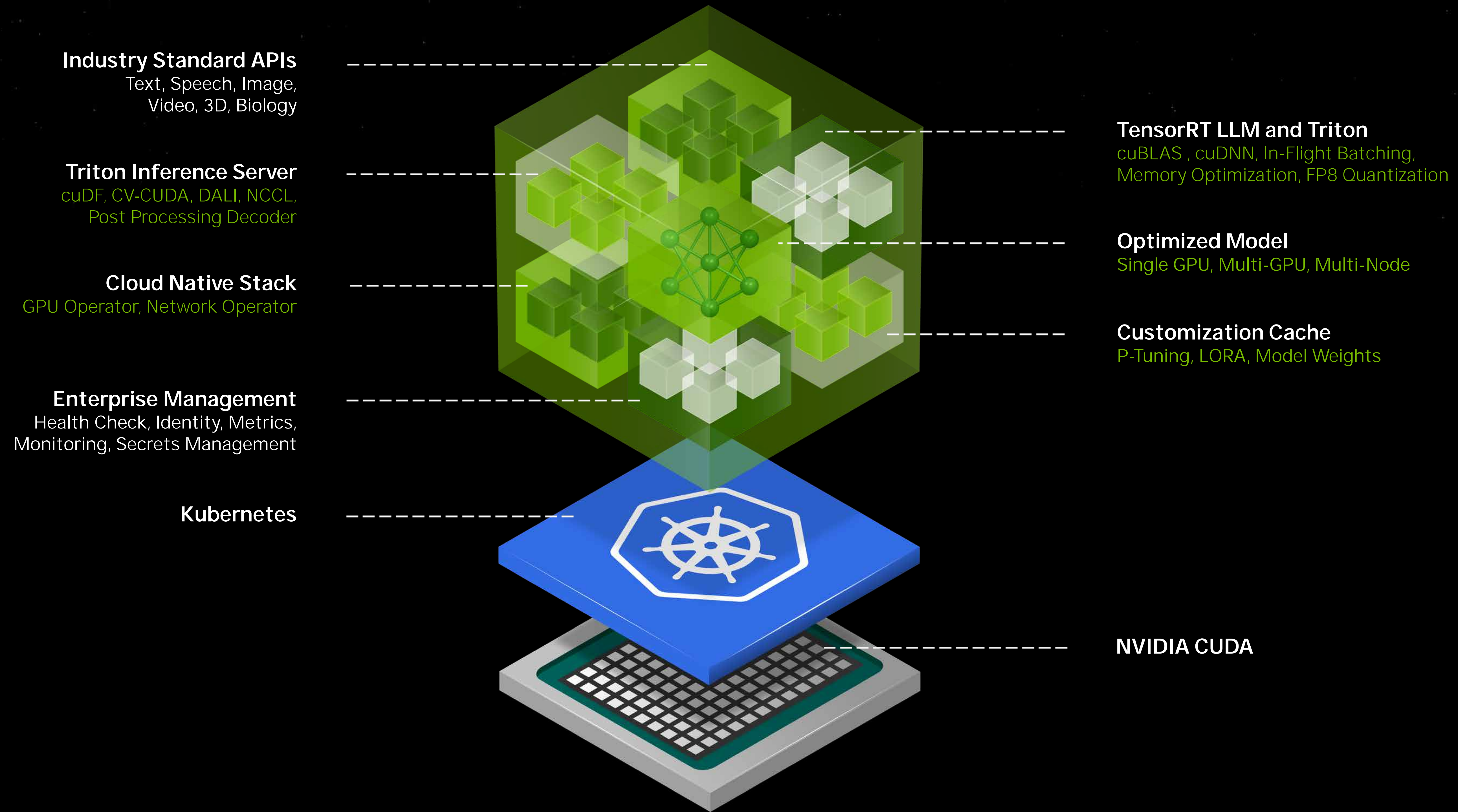
∞ Meta





NVIDIA INFERENCE MICROSERVICE

Pre-Trained AI Models
Packaged and Optimized to Run Across
CUDA Installed Base



100's of Millions of CUDA GPUs Installed Base

Discover

MODELS

Reasoning

Visual Design

Speech & Translation

Video

Biology

DATA PROCESSING

Embedding

Retrieval

OPTIMIZATION

Inference

Route Planning

INDUSTRIES

Gaming

Healthcare

Automotive

Industrial

Top Open Foundation Models

The leading open models built by the community, optimized and accelerated by NVIDIA's enterprise-ready inference runtime



stability-ai
stable-diffusion-xl
image generation text-to-image



google
gemma-7b
chat language generation



stability-ai
stable-video-diffusion
image-to-video video generation



cohere
aya-101
text-to-text multilingual

Open Full Page

Input

Try Python Node.js Shell

Input Prompt ⓘ

View Examples

A happy dog hanging out at the park

View Parameters

Reset Parameters

Run

Output

Preview 250K



Trending Now

The latest and most popular additions to the list



cohere
command-r
text-to-text language generation



gettyimages
edify-image
text-to-image image-to-image



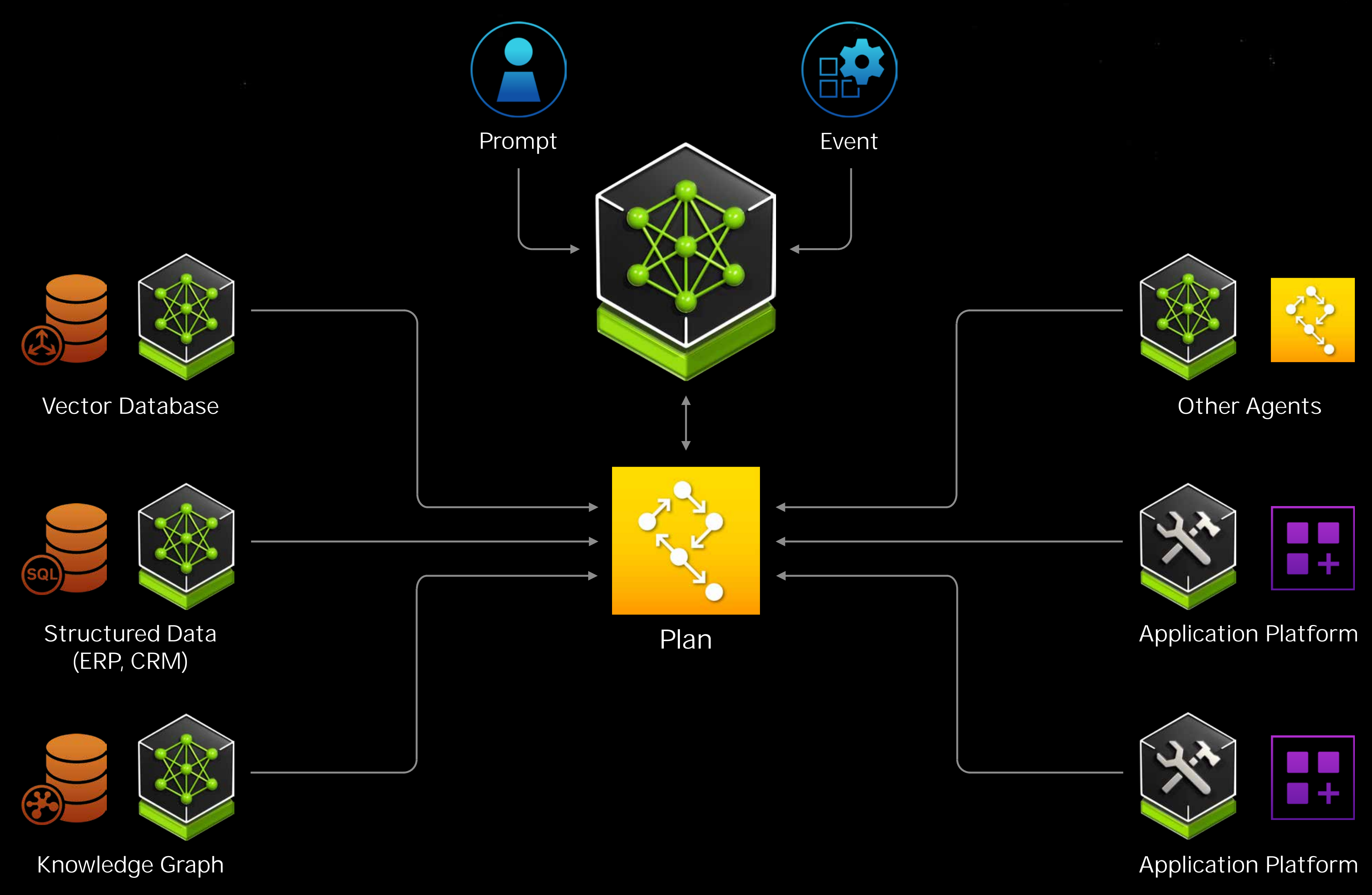
shutterstock
edify-3d
3D generation text-to-3D



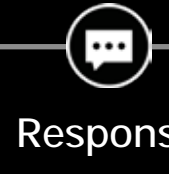
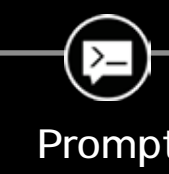
meta
codellama-70b
code generation text-to-code

Explore by Collection

Discover new use-cases and the right set of APIs to turbocharge your enterprise



What is CTL? Write an example CTL test.




NVIDIA ChipNeMo
Llama 2 70B
Before Fine-Tuning

CTL stands for "Combinational Timing Logic" and it's a type of digital circuit used in ASIC (Application-Specific Integrated Circuit) design. CTL is used to implement finite state machines, counters, and other sequential logic circuits.

Here is an example of a CTL Test for a simple finite machine:

1. Set up the initial state: Set EN to 0 and wait for the state machine to settle into state S0.
2. Test the transition from S0 to S1: Set EN to 1 and wait for the state machine to transition to state S1. Verify that the output signal (Q) is set to 1.
3. Test the transition from S1 to S0: Set EN to 0 and wait for the state machine to transition back to state S0. Verify that the output signal (Q) is set to 0.
4. Test the reset behavior: Set EN to 0 and wait for the state machine to settle into state S0. Then, assert the reset signal (RST) and verify that the output signal (Q) is set to 0.
5. Test the clock behavior: Set EN to 1 and wait for the state machine to transition to state S1. Then, pause the clock and verify that the output signal (Q) remains set to 1. Resume the clock and verify that the state machine transitions back to state S0 when the clock is restarted.

What is CTL? Write an example CTL test. 

Prompt

Response

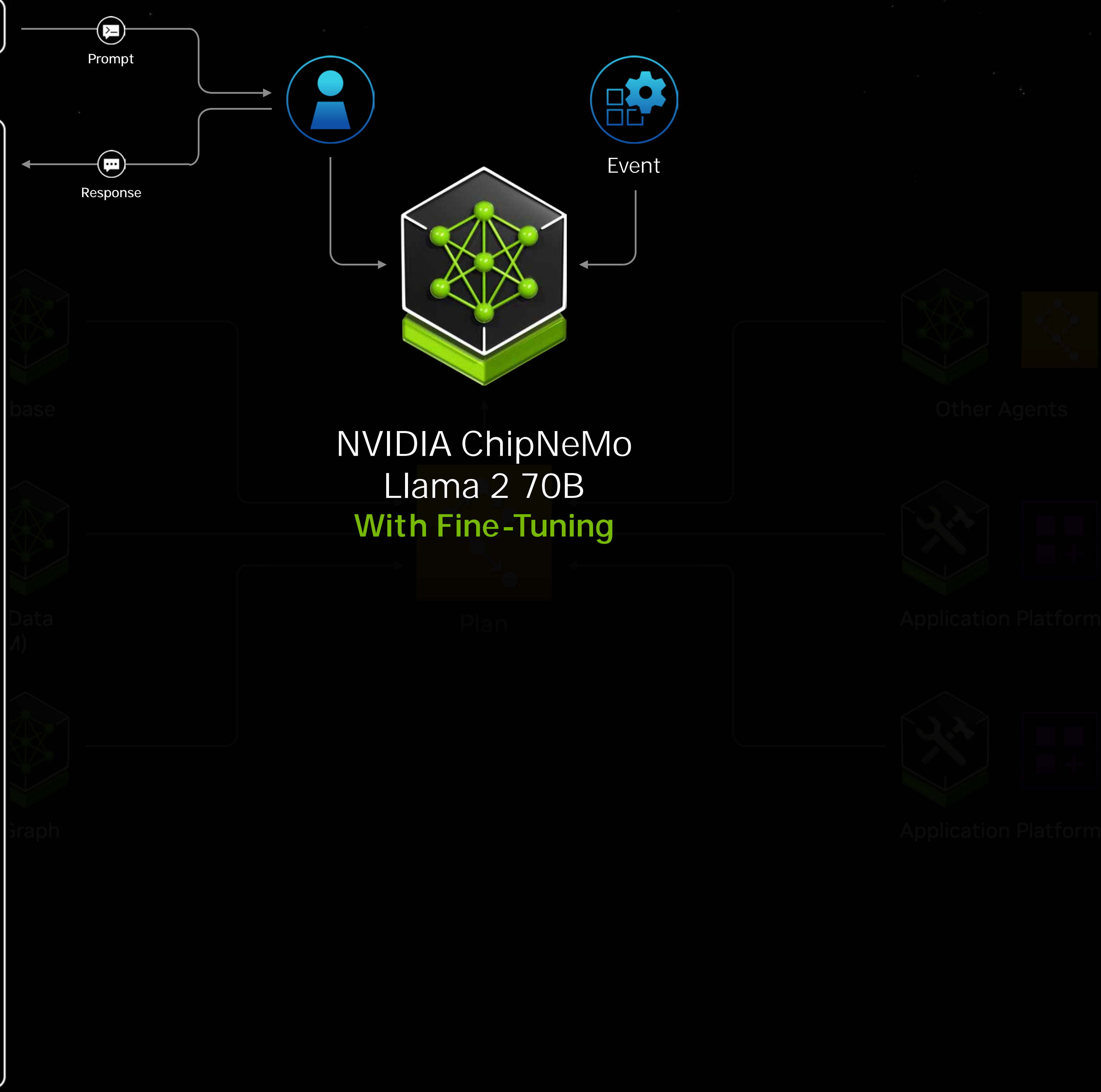
Event

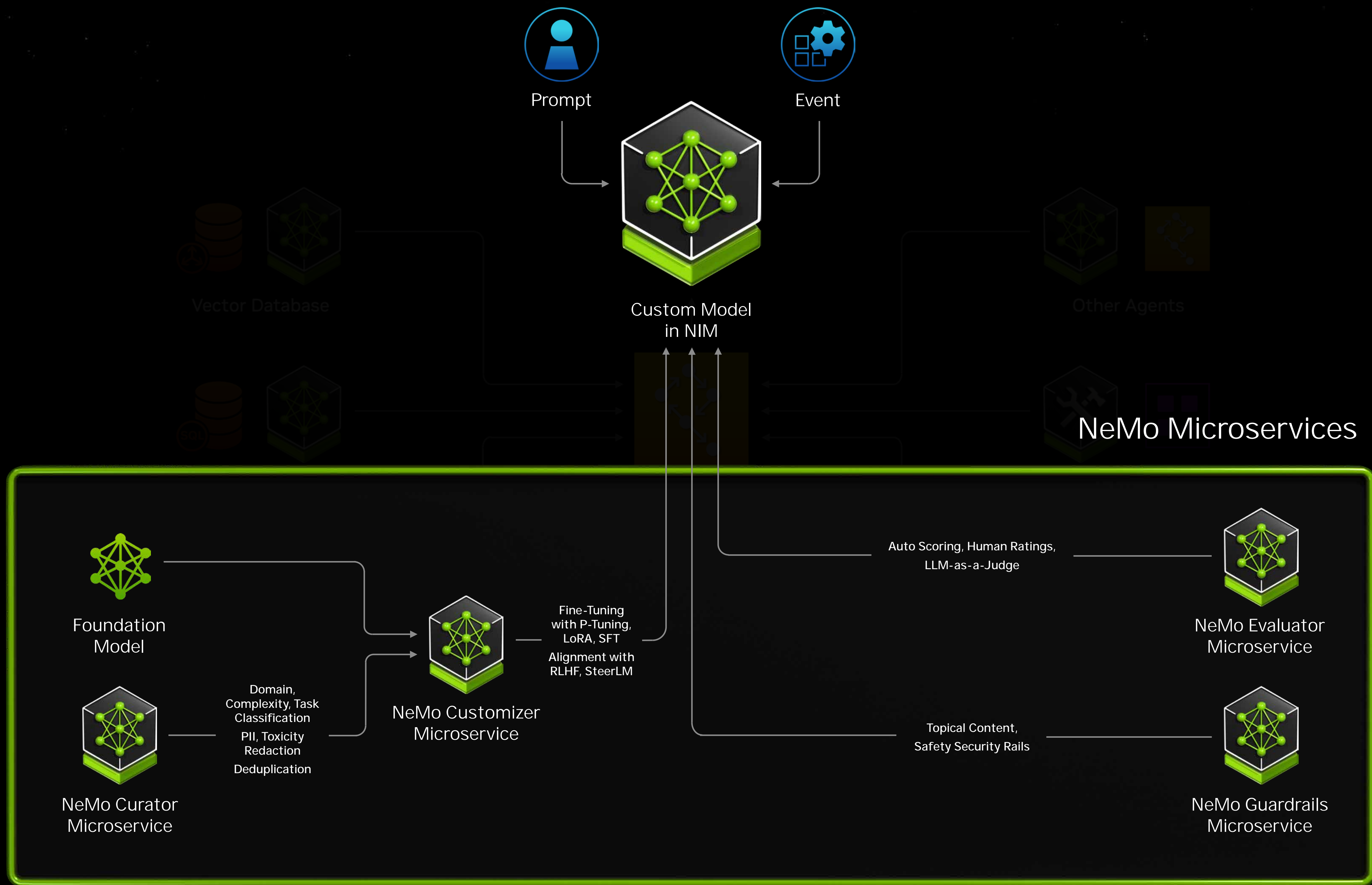
CTL stands for "Compute Trace Library". It is a trace generation library used to create functional and performance tests for the GPU architecture. CTL provides an interface to write tests in C++ that can be compiled into a trace that runs on various simulation models, such as Fmodel, RTL, and Emulation.

Here is an example of a simple CTL test that allocates a buffer, initializes it with some values, and then reads back the values to verify that they are correct:

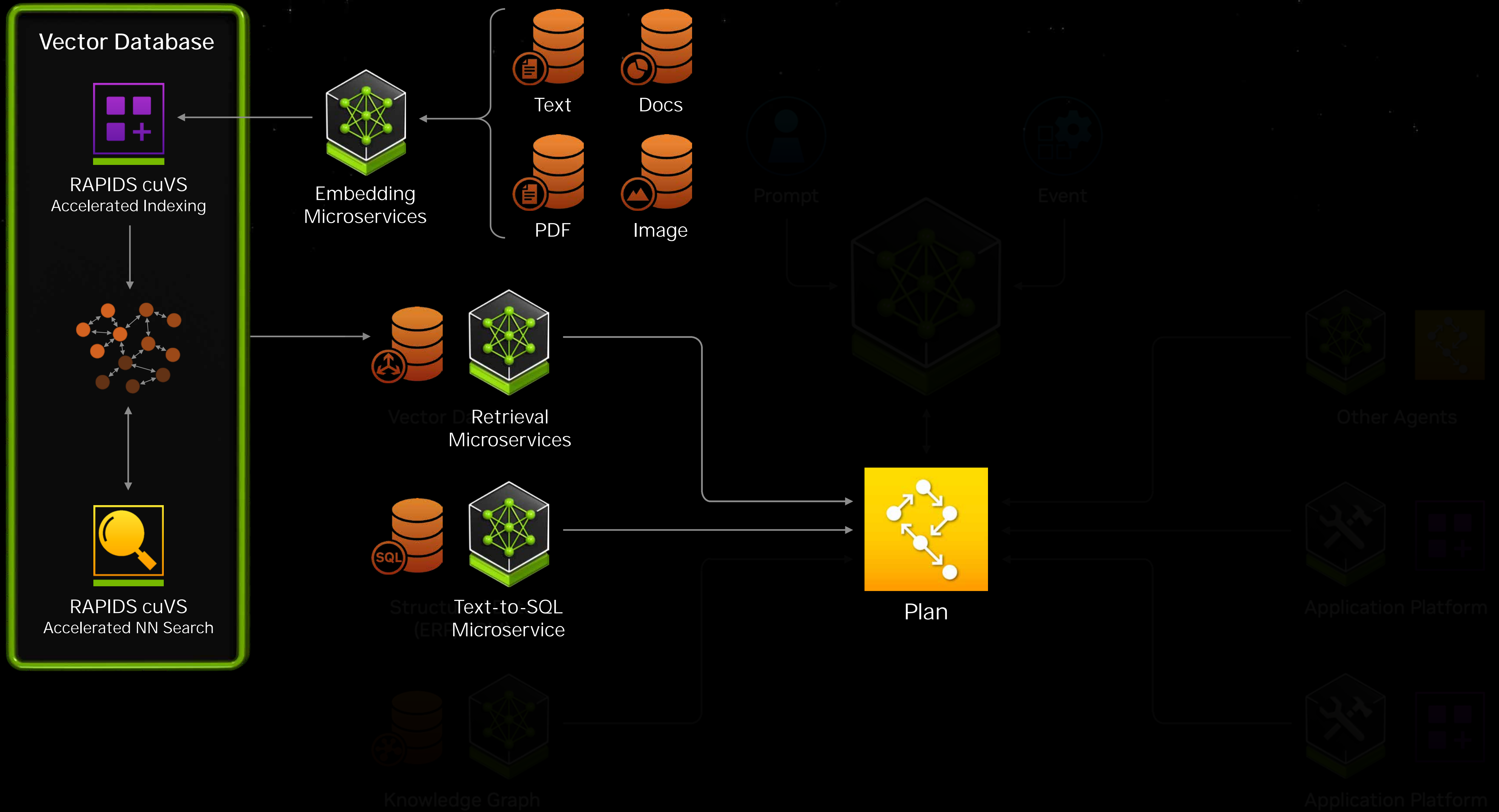
```
#include "ctl.h"
using namespace ctl;
int main(int argc, char* argv[]) {
    // Set up the test
    ctlAddTest("simple_buffer", "", "ga100");
    // Allocate a buffer and initialize it
    // with some values
    unsigned* buffer;
    ctlMalloc((void**)&buffer, 1024, "BUFFER",
    CTL_GLOBAL_MEMORY);
    ctlMemset(buffer, 0, 1024);
    ctlMemcpy(buffer, (void*)0x12345678, 4,
    ctlMemcpyHostToDevice);
    // Set up the task
    ctlCreateTask("TASK", "kernel");
    ctlSetQmdCtaRasterDimension("TASK", 1, 1,
    1);
    ctlSetQmdCtaThreadDimension("TASK", 1, 1,
    1);
    // Set up the kernel
    char program[] =
    "!!SPA7.0\n"
    ".THREAD_TYPE COMPUTE\n"
    ".option
```

NVIDIA ChipNeMo
Llama 2 70B
With Fine-Tuning

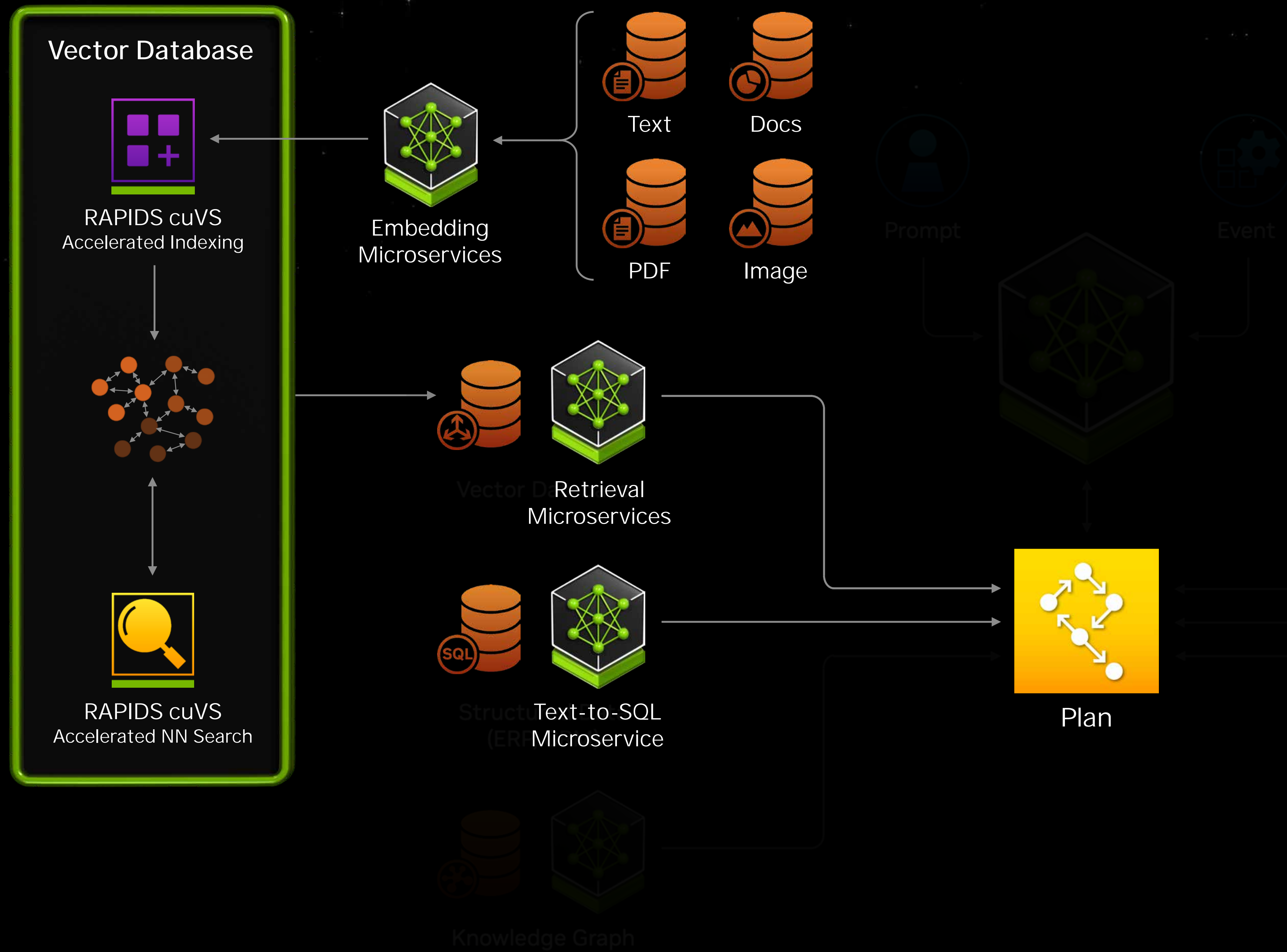




NeMo Retriever



NeMo Retriever



Chat with PDF

The screenshot shows a document with the following sections:

Setting New Records in MLPerf Inference v3.0 with Full-Stack Optimizations for AI

NVIDIA AI and H100 Tensor Core GPU deliver record results

NVIDIA H100 Tensor Core GPUs, which made their MLPerf Training debut just 5 months ago, set new per-accelerator performance records across all MLPerf Inference v3.0 workloads. Looking at the NVIDIA single node DGX H100 results this round, performance increased by up to 17% in just 6 months on the same hardware through software improvements alone. Compared to the NVIDIA A100 Tensor Core GPU submission in MLPerf Training v2.1, the latest H100 submission delivered up to 3.1x more performance per accelerator.

Beating SOTA Inference Performance on NVIDIA GPUs with GPUNet

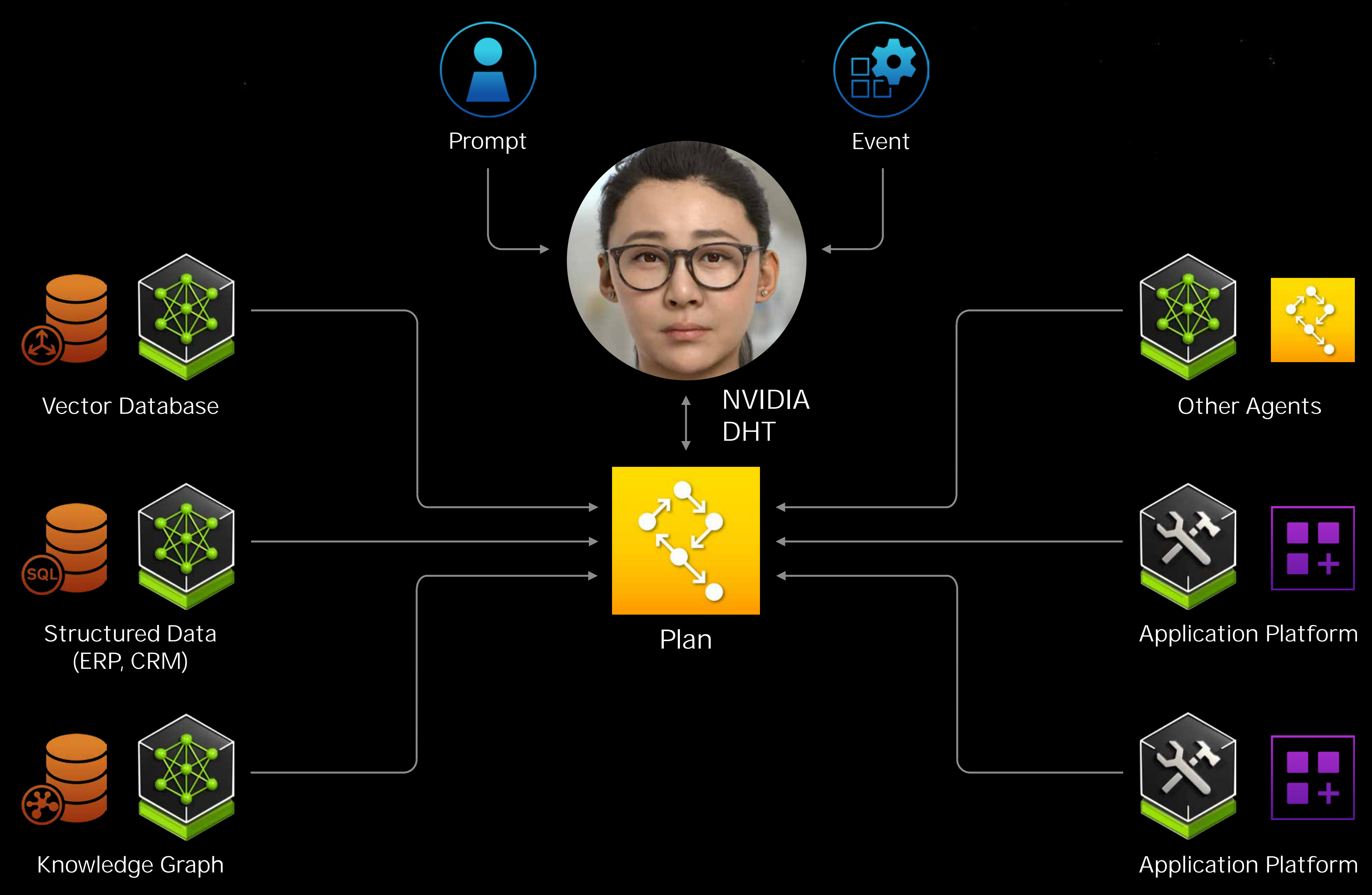
GPUNet is a class of convolutional neural networks designed to maximize the performance of NVIDIA GPUs using NVIDIA TensorRT.

GPUNet NAS design methodology

Efficient architecture search and deployment-ready models are the key goals of the NAS design methodology. This means little to no interaction with the domain experts and efficient use of cluster nodes for training potential architecture candidates. Most important is that the generated models are deployment-ready.

Crafted by AI

Benchmark	Max Scale Records (records)
Large language model (GPT-3)	109
Natural language processing (BERT)	6.13 (8 seconds)
Recommendation (DLRMv2)	1.81
Object detection, heavyweight (Mask R-CNN)	1.47
Object detection, lightweight (RetinaNet)	1.51





Relevant

Reliable

Responsible



Joule

A copilot that truly understands your business

Embedded AI capabilities

Cloud ERP

Supply Chain
Management

Human Capital
Management

Spend Management
& Business Network

Customer Relationship
Management

Business Technology
Platform

AI Foundation

on Business Technology Platform

SAP Business AI Ecosystem



NeMo Retriever

NIM

NeMo Guardrails

Accelerated by NVIDIA

Nemotron Models

RAPIDS

cuOpt



THE INTELLIGENT WORKFLOW COMPANY

Now Assist **Generative AI** Experiences powered by ServiceNow Platform

AI Research & Development

Enterprise Foundation Models

AI Trustworthiness

StarCoder

Engineering & Applied Science

Domain Specific AI Models

Enterprise Workflows

ServiceNow Global Cloud Services

AI Embedded Workflows

Privacy and Security

Grounded in Customer Data

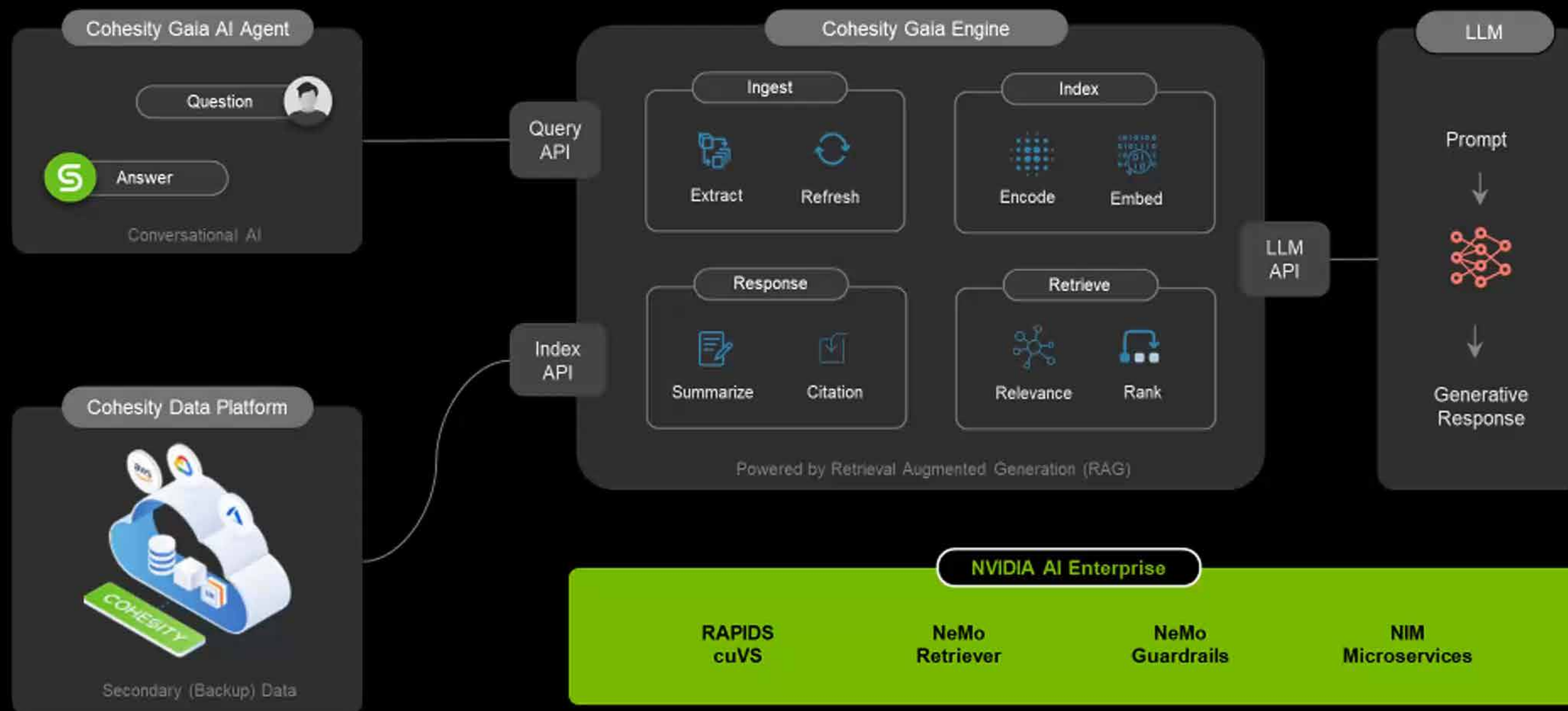
POWERED BY  NVIDIA SINCE 2018

DGX Cloud DGX SuperPOD NeMo TensorRT NIM Microservices NVIDIA AI Enterprise

COHESITY



Cohesity Gaia – First to Provide AI-Powered Business Insights from Secondary Data





DATA + AI – ALL ON NVIDIA

APPLICATIONS



Snowflake Copilot



Snowflake Document AI



**SNOWPARK
CONTAINER SERVICES**

Model Training, Fine Tuning



SNOWFLAKE CORTEX

Serverless AI, LLMs, Search



ENTERPRISE DATA

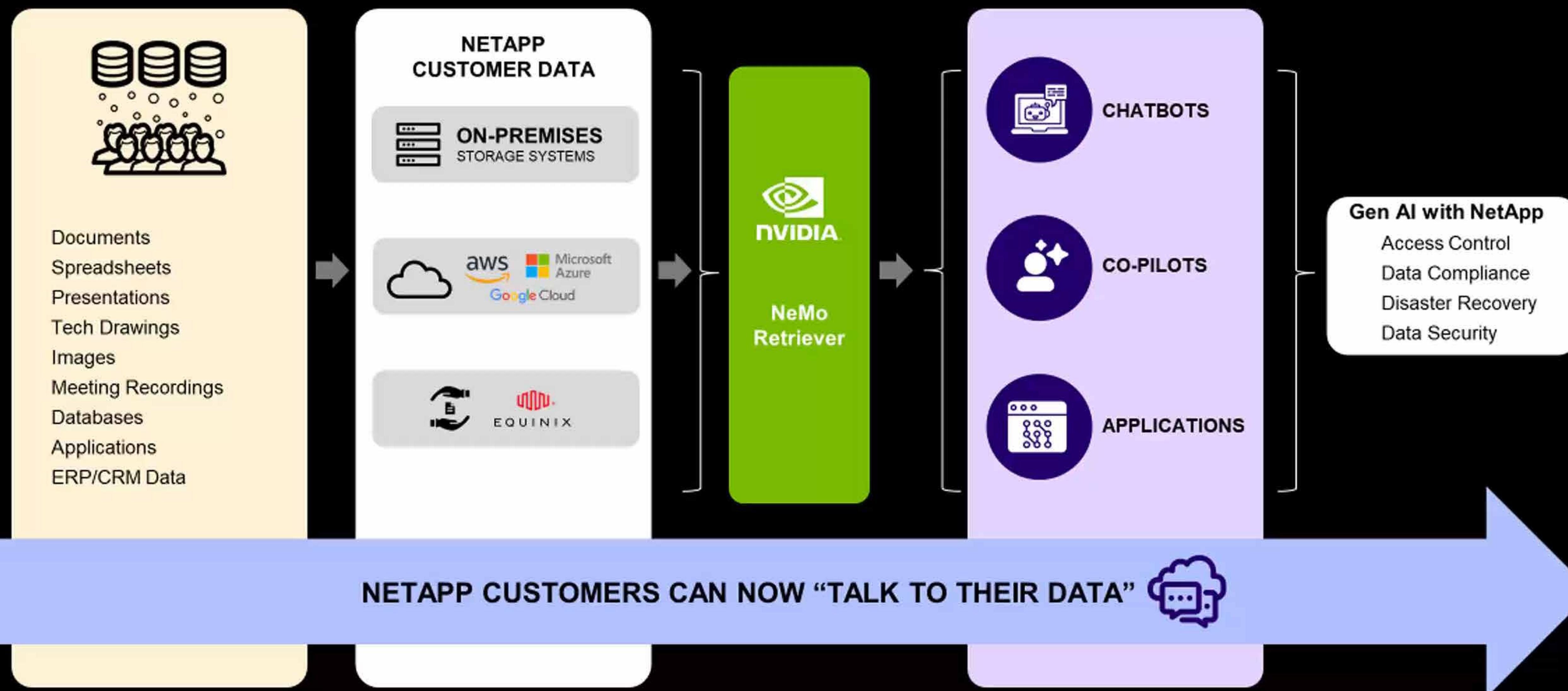
ACCELERATED BY NVIDIA

NeMo Retriever

Triton Inference Server

TensorRT

NetApp Unlocks Exabytes of Data for Secure, Private Gen AI



Dell AI Factory with NVIDIA

INDUSTRY'S FIRST END-TO-END ENTERPRISE AI SOLUTION

Model Lifecycle

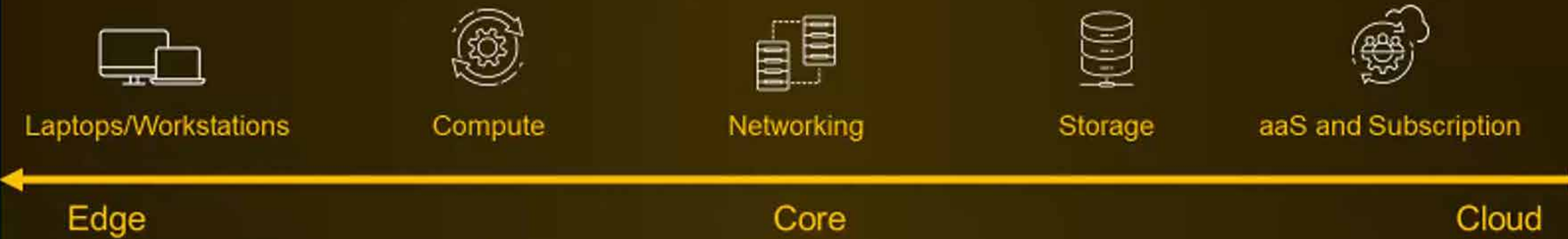


Dell Professional Services



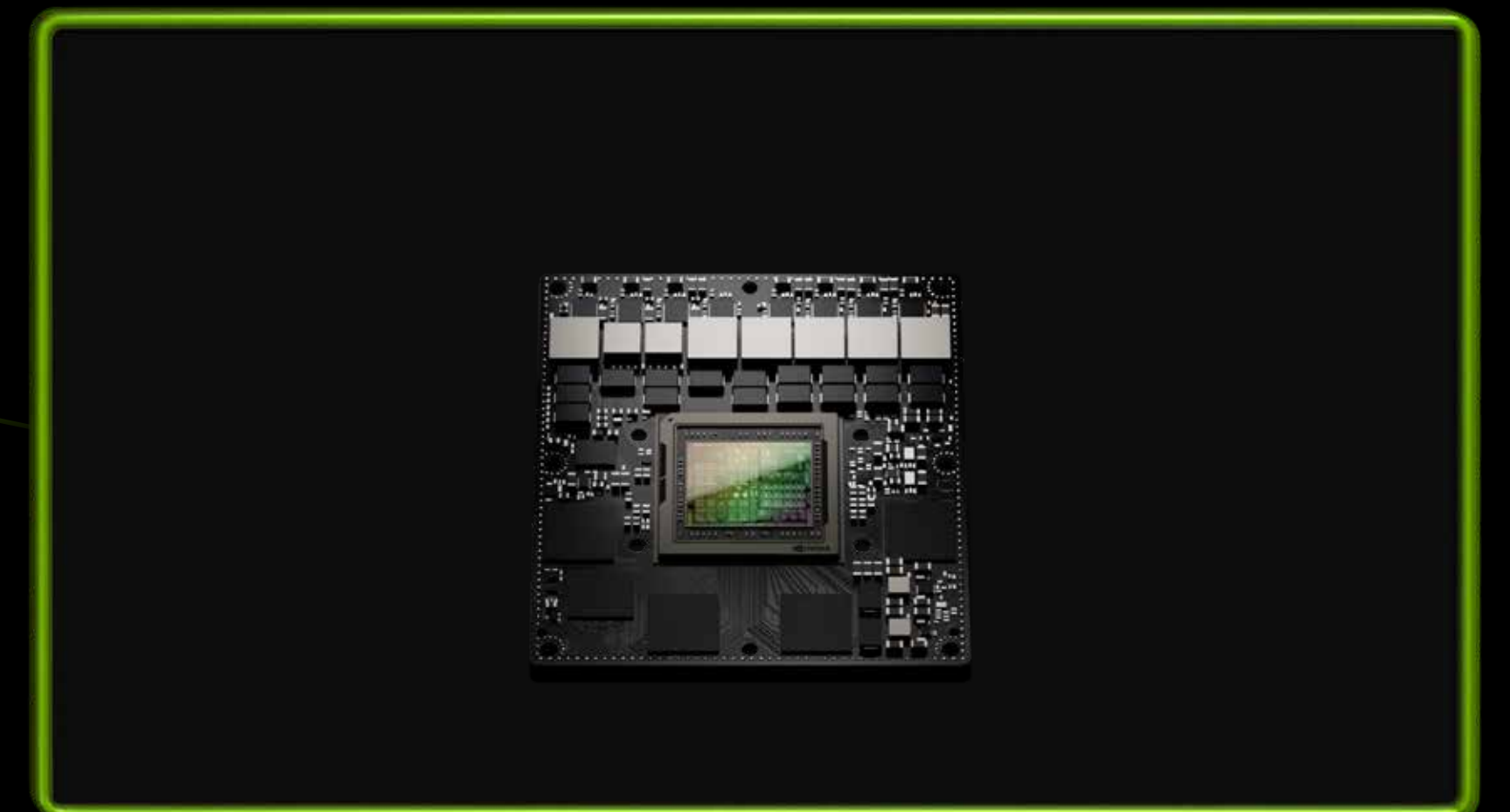
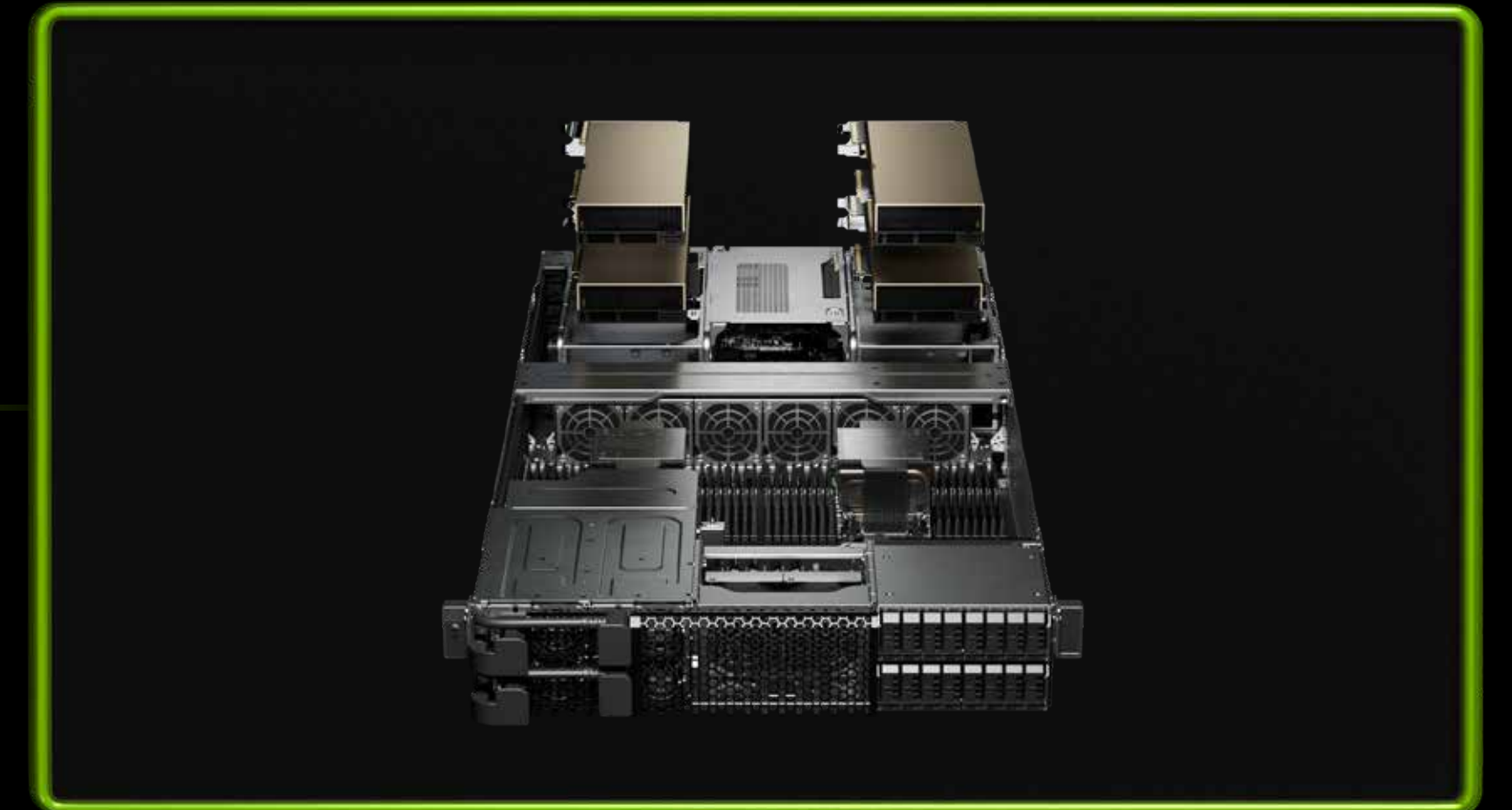
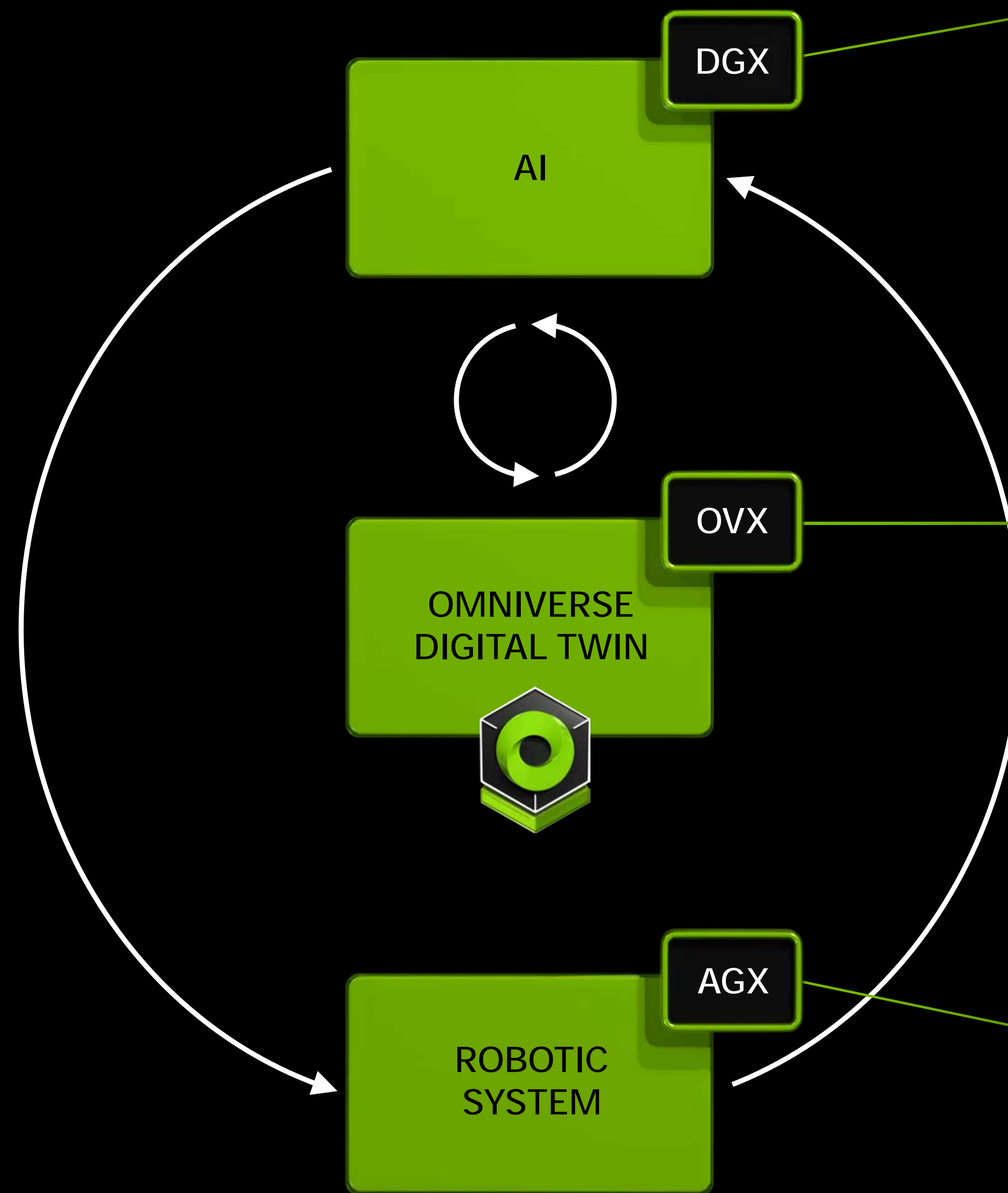
NVIDIA AI Enterprise Software

Dell and NVIDIA Infrastructure



Broadest AI solutions portfolio from desktop to data center to cloud, all in one place



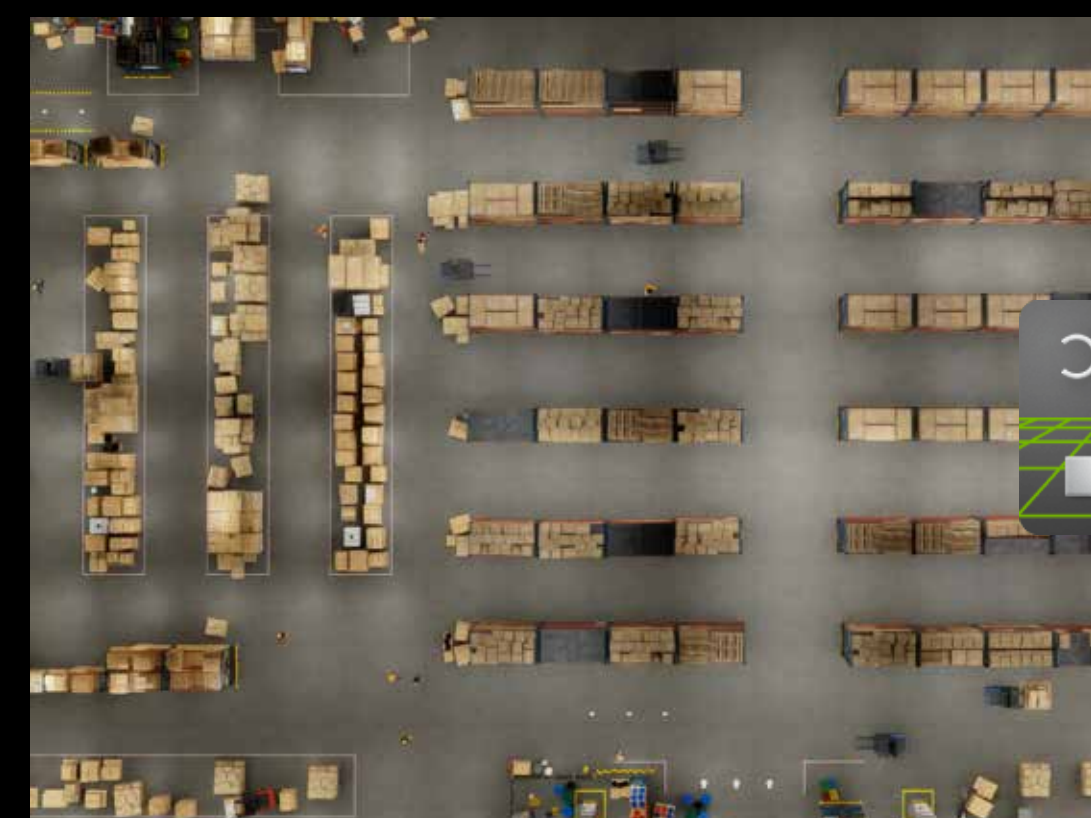


AMR Robot Digital Twins



Isaac Sim

Warehouse Digital Twin



Isaac Sim

Worker Digital Twins



Isaac Sim

USD APIs

Omniverse Channel APIs

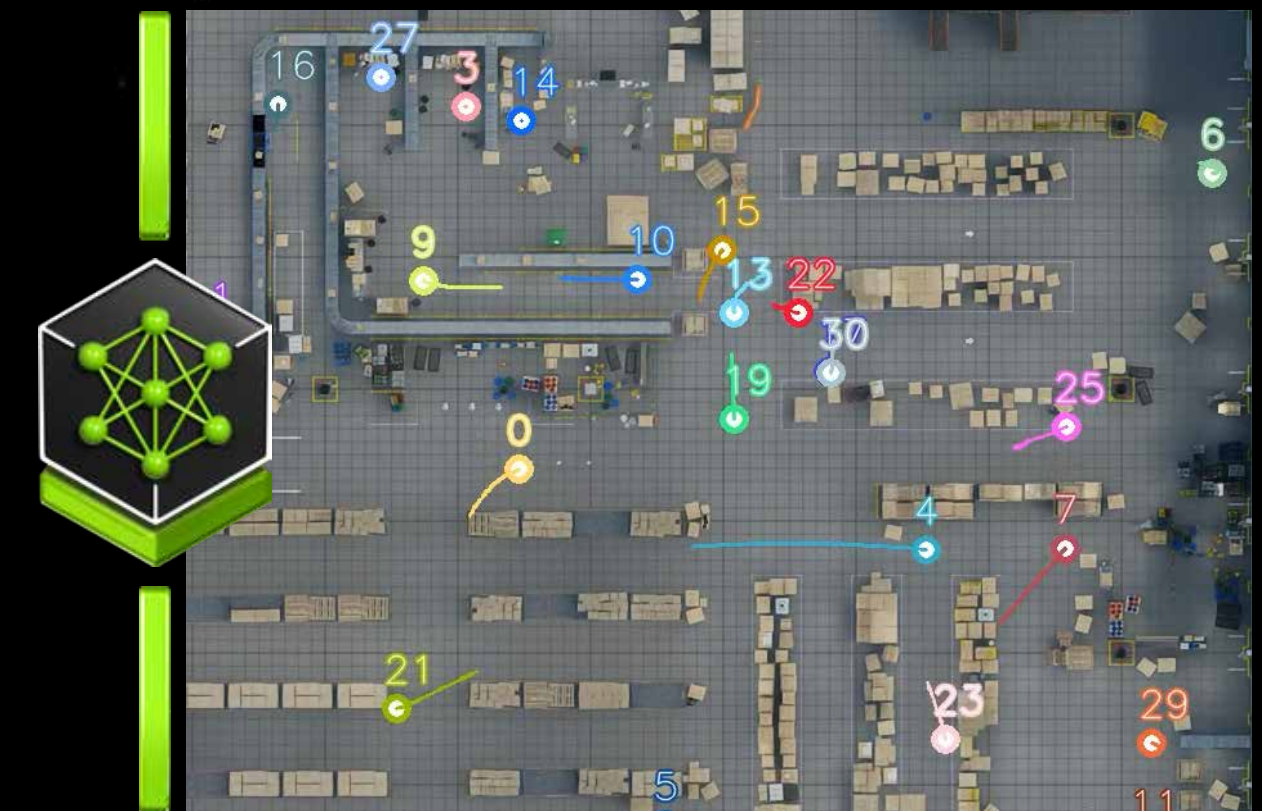


NVIDIA Omniverse Cloud



NVIDIA DGX Cloud

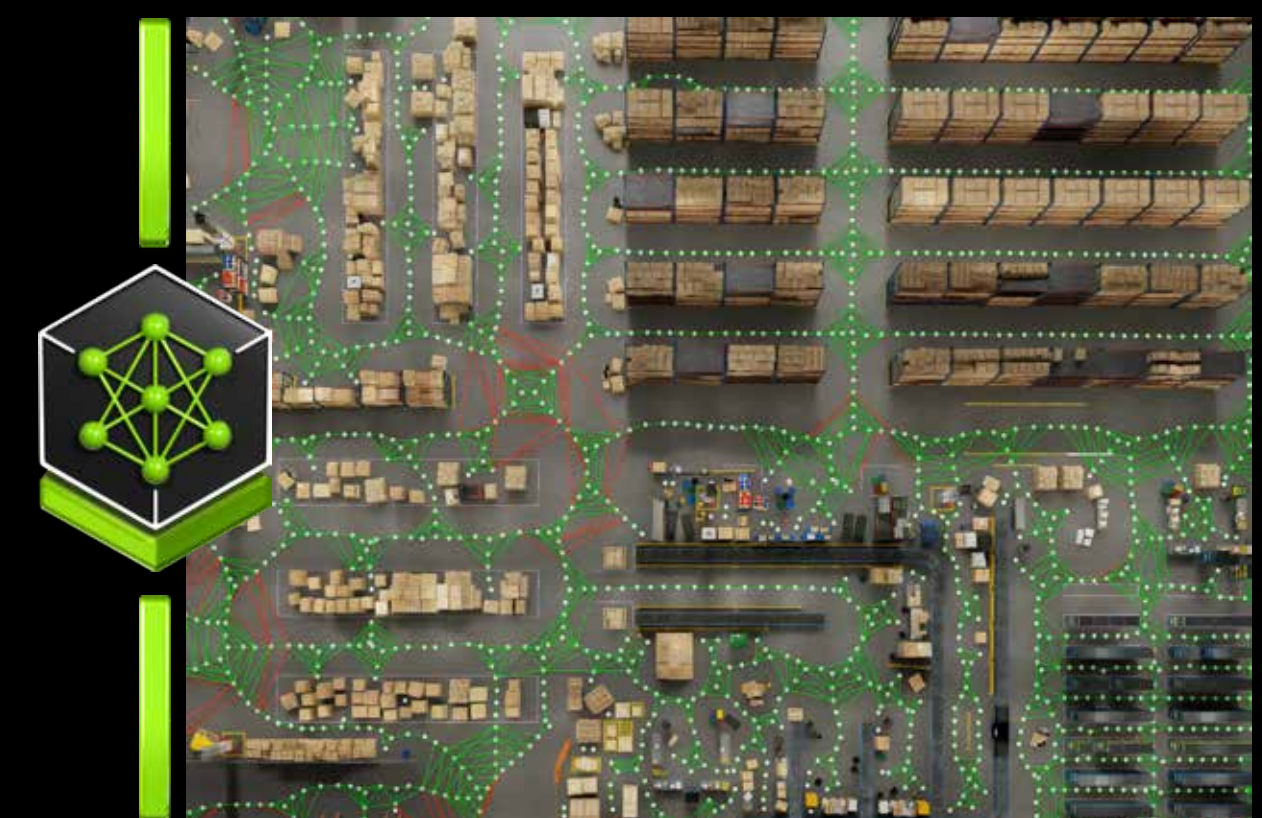
Sensor Fusion



Metropolis

Occupancy APIs

Route Planning

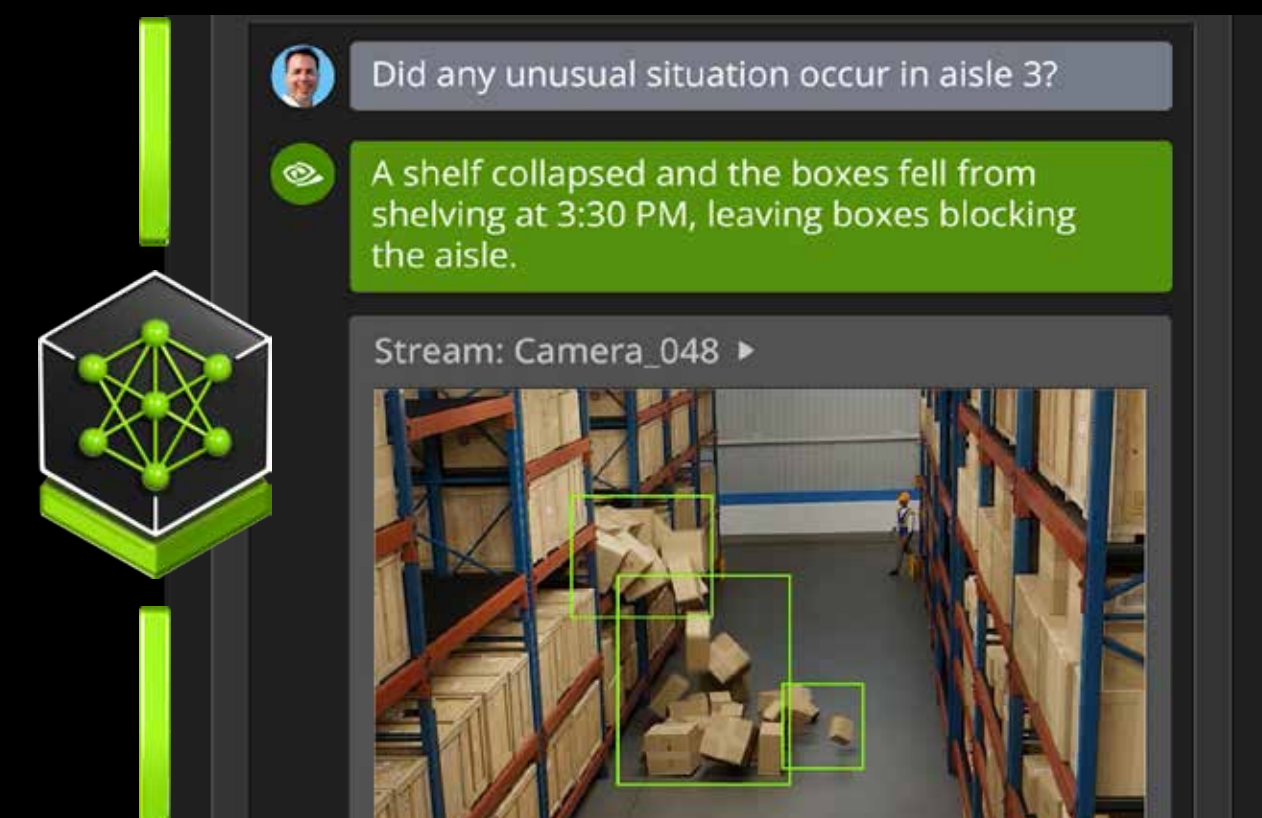


cuOpt

cuOpt APIs

VLM APIs

Operator UI



Metropolis



LESTAK-10

DIGITAL TWIN

FRAGILE

DIGITAL TWIN

DIGITAL TWIN

DIGITAL TWIN

DIGITAL TWIN

DIGITAL TWIN

DIGITAL TWIN

DIGITAL TWIN

DIGITAL TWIN

DIGITAL TWIN

DIGITAL TWIN

DIGITAL TWIN

DIGITAL TWIN

DIGITAL TWIN

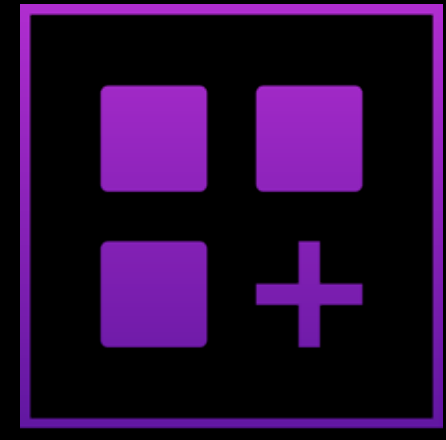
DIGITAL TWIN

DIGITAL TWIN

DIGITAL TWIN

DIGITAL TWIN

DIGITAL TWIN



Application

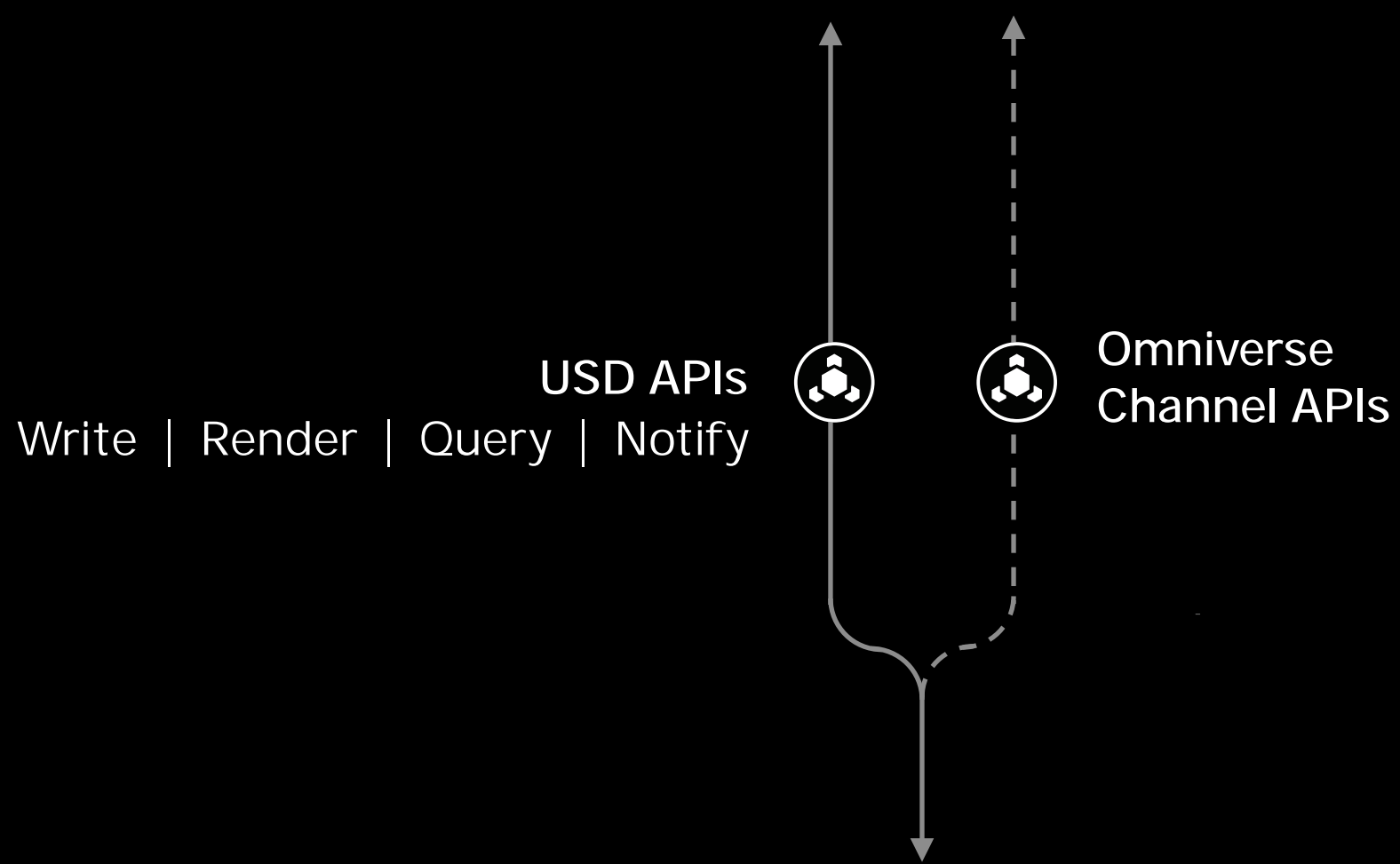
USD APIs
Write | Render | Query | Notify

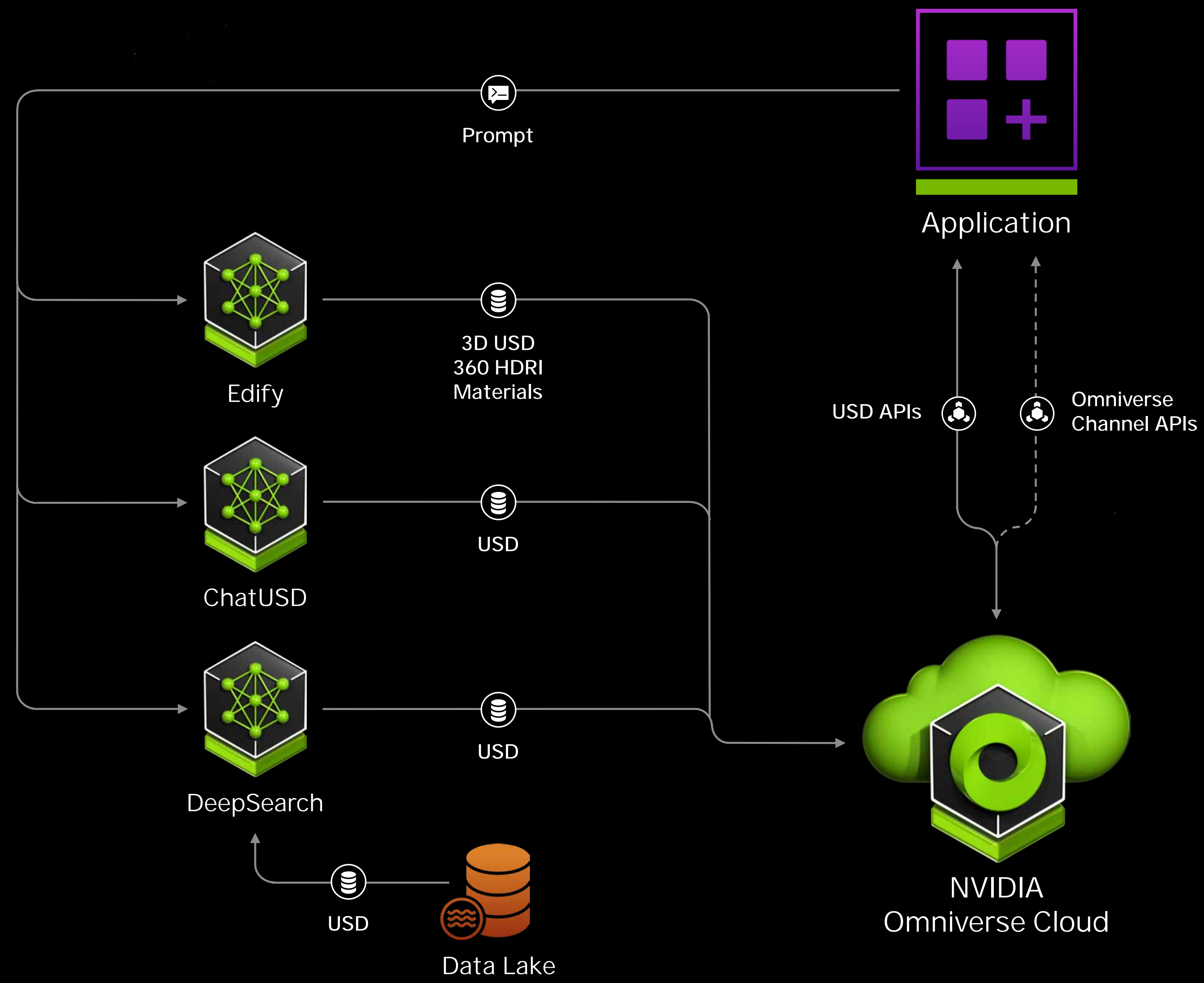


Omniverse
Channel APIs



NVIDIA
Omniverse Cloud





SIEMENS

nVIDIA


Building the Next Era of Industrial Digitalization Together

Immersive digital twins




Siemens Xcelerator NVIDIA Omniverse Cloud

Generative AI




Siemens Xcelerator NVIDIA AI Enterprise

Industrial Edge AI & Robotics

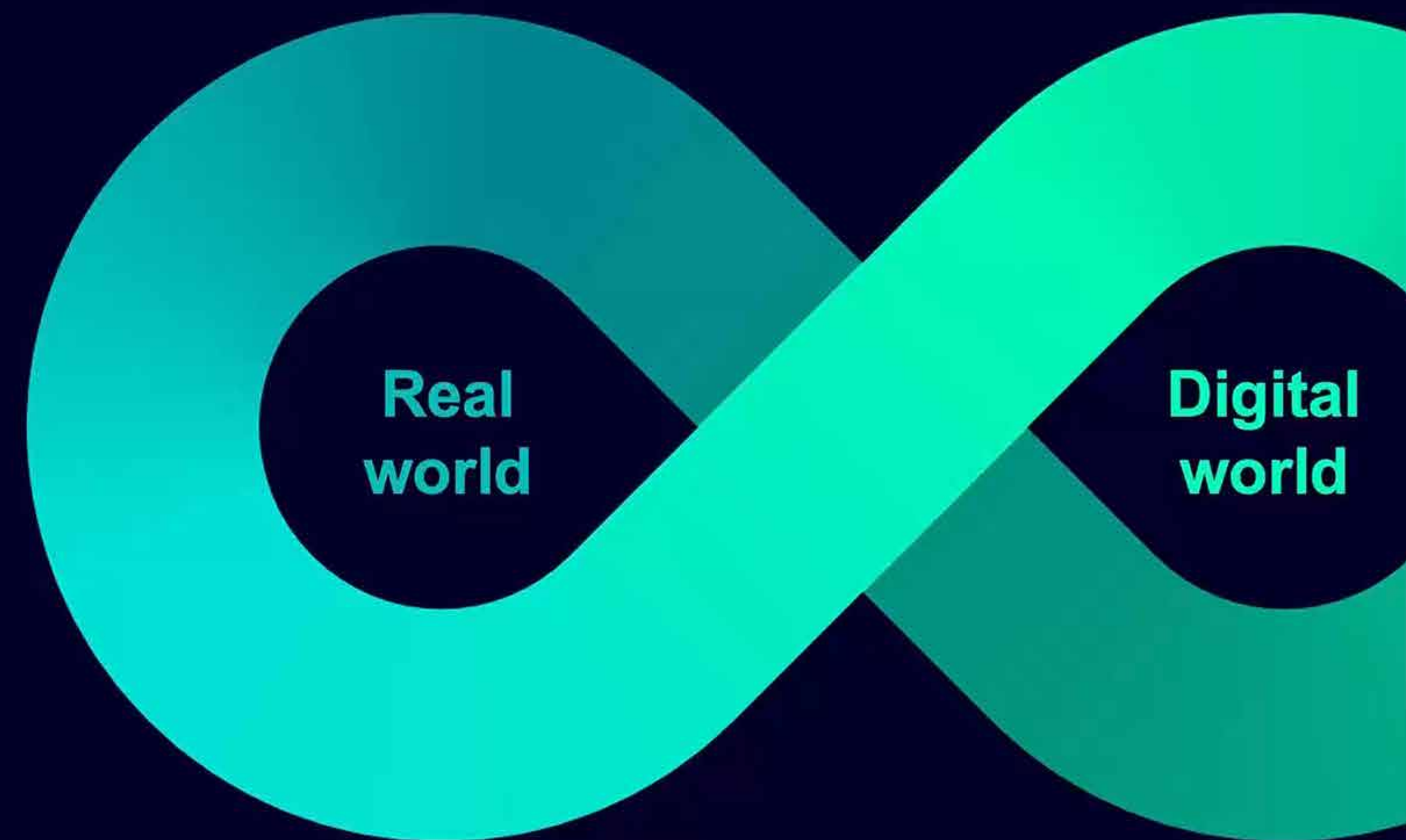


Siemens Xcelerator NVIDIA Robotics

Industrial AI-Physics Simulation

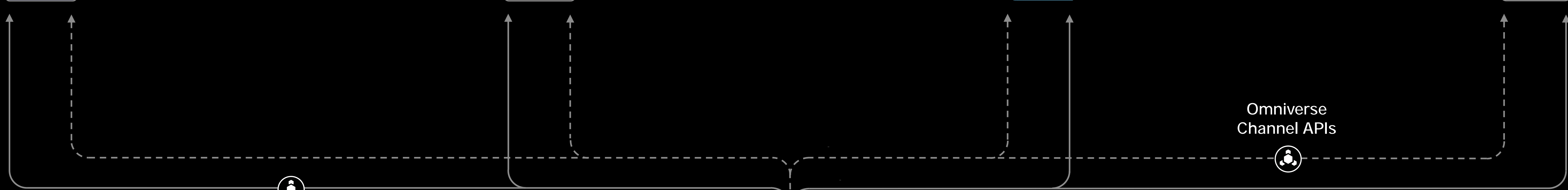
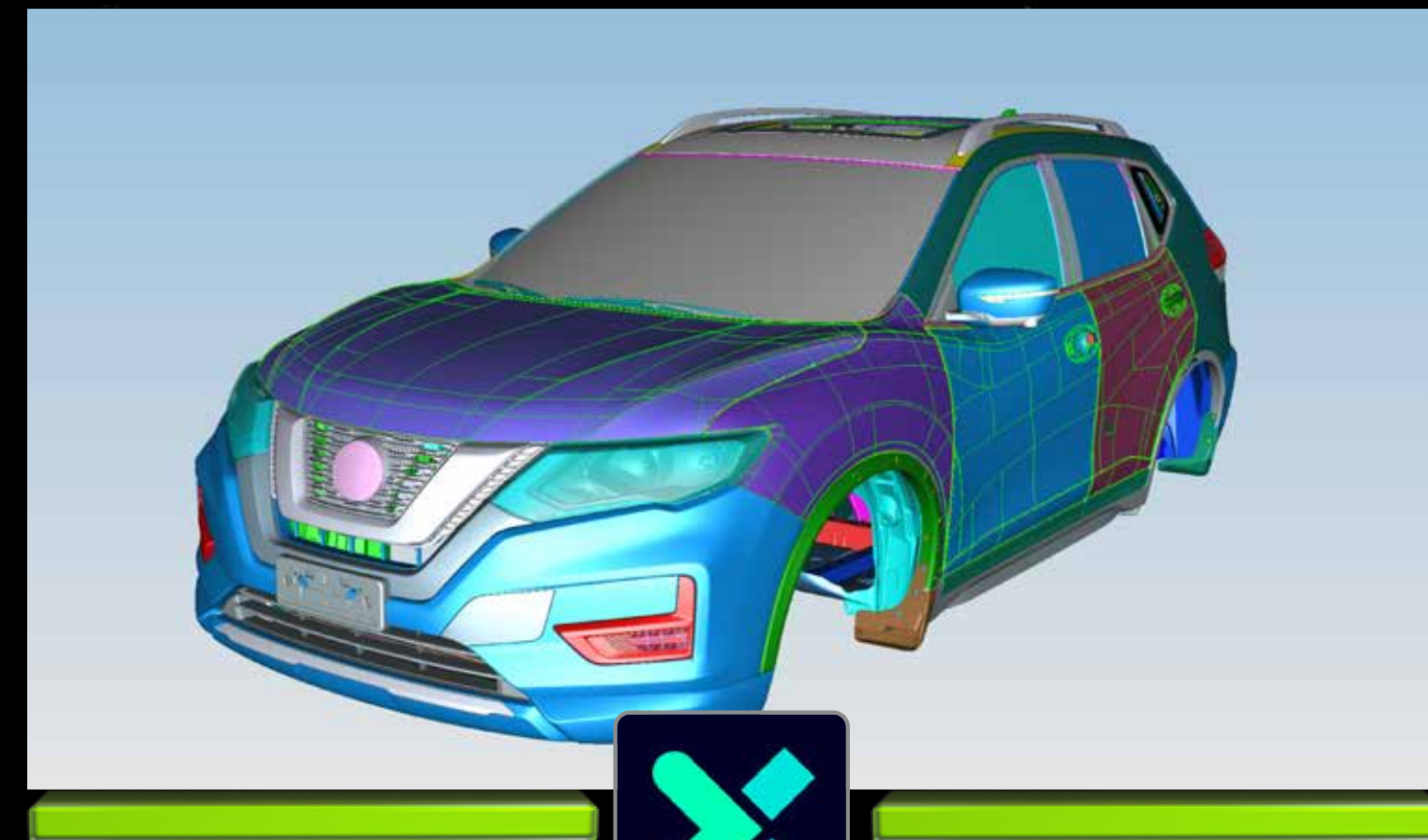


Siemens Xcelerator NVIDIA Accelerated Computing



SIEMENS





USD APIs

Omniverse Channel APIs



NVIDIA Omniverse Cloud



Team Hyundai

GENERAL

Dashboard

Workspaces

ORDERS

New Asset

My Orders

Cart

PRESETS

Templates

USER MANAGEMENT

Members



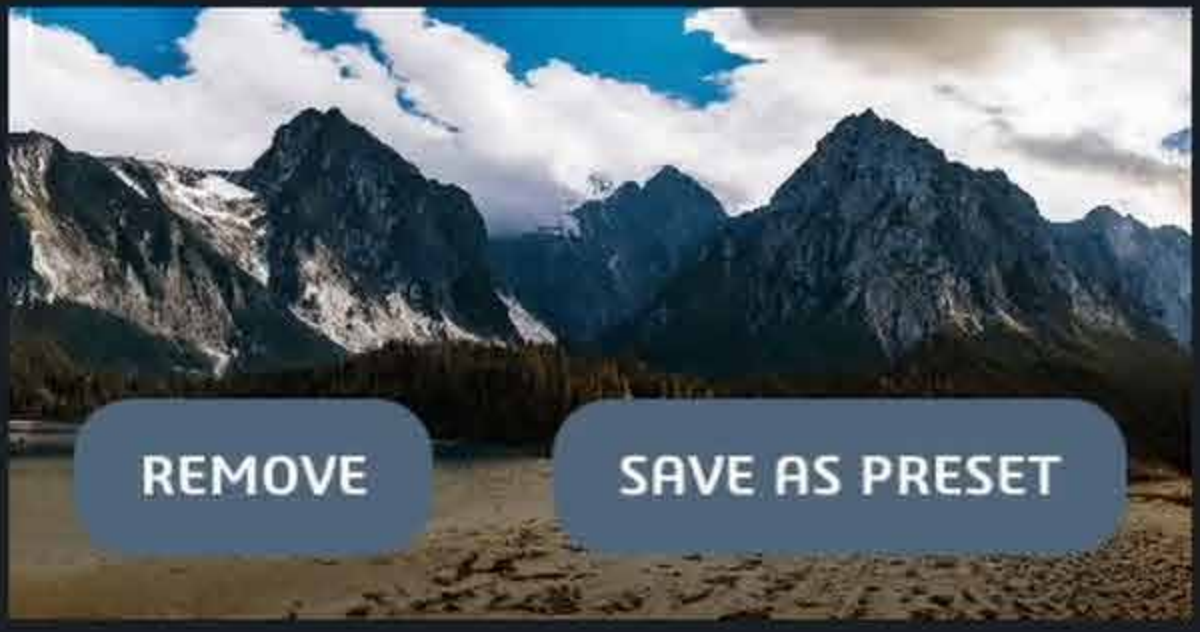
- 🏠
- 🗃️
- 🌊
- 👤
- 💡
- 📺

Environments

3D

360°

2D



Vertical Horizon Offset

Slider control for Vertical Horizon Offset, set to 0

Environment Rotation

Slider control for Environment Rotation, set to 0

Prompt: Breathtaking swiss mountain landscape with a large open space in a cinematic afternoon atmosphere.

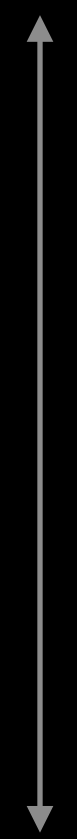
GENERATE

ADD TO TEMPLATE

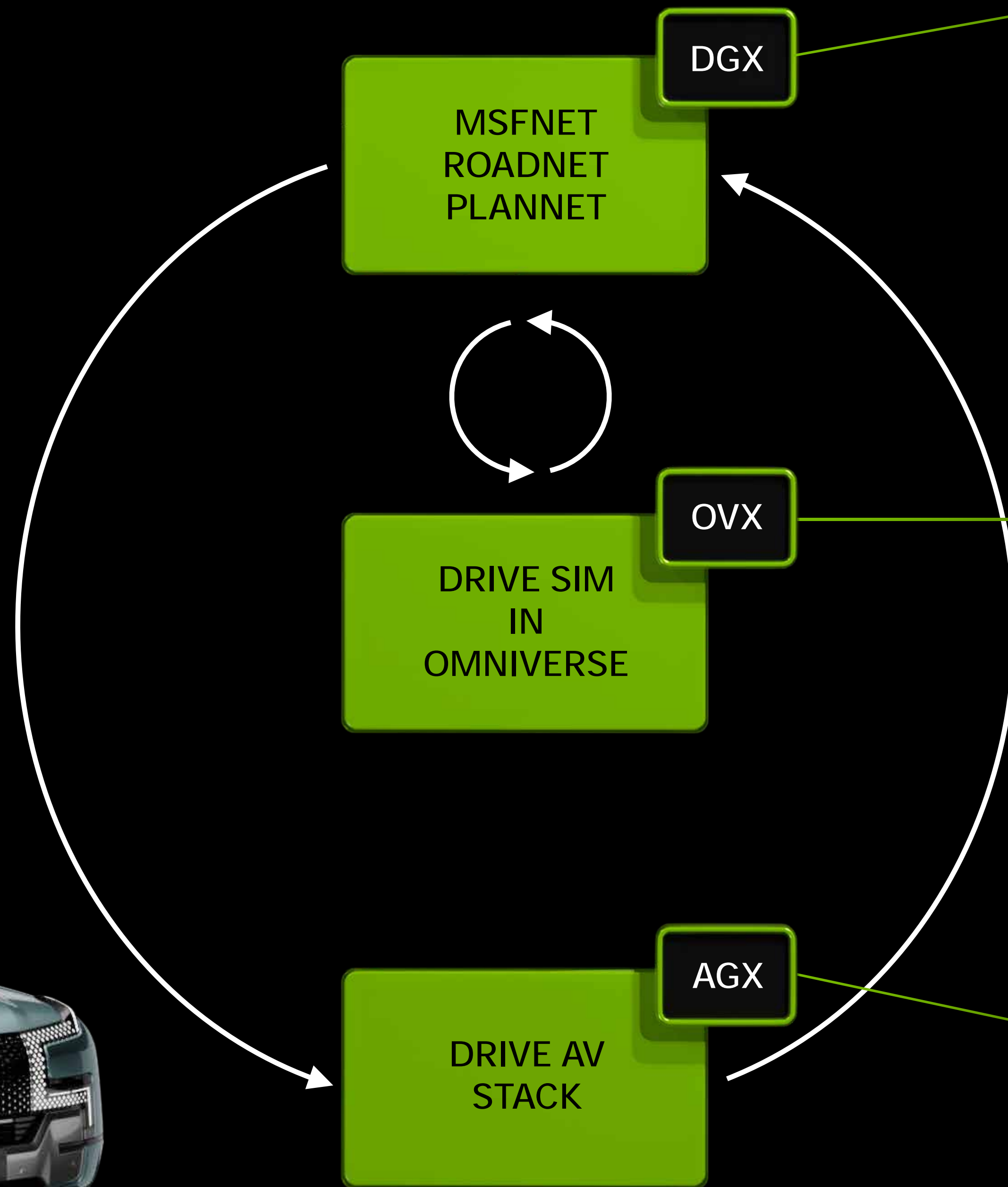
ADD TO CART



Apple Vision Pro



NVIDIA
Omniverse Cloud



DGX

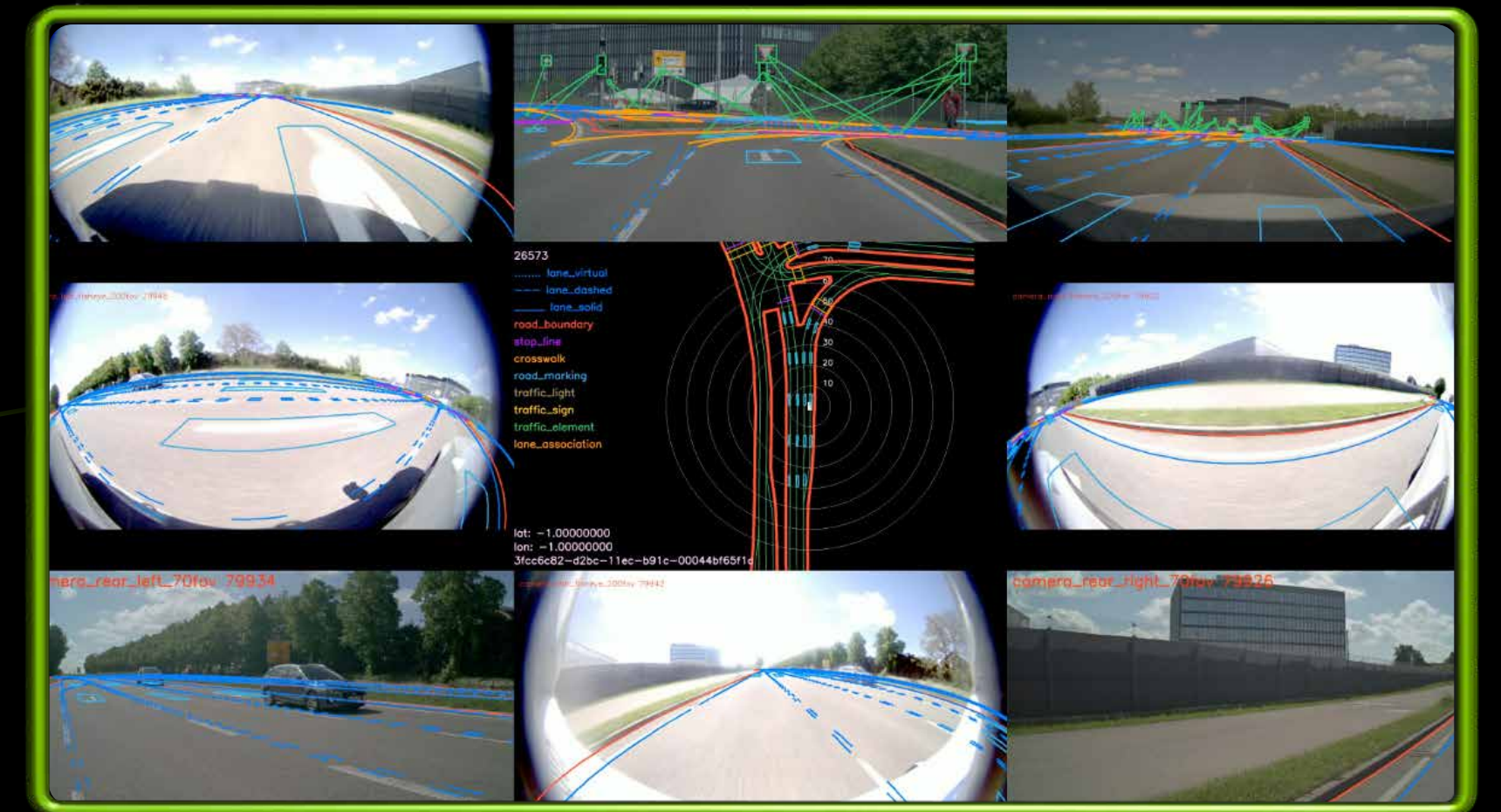
MSFNET
ROADNET
PLANNET

OVX

DRIVE SIM
IN
OMNIVERSE

AGX

DRIVE AV
STACK



DRIVE THOR
ASIL-D AV COMPUTER & STACK

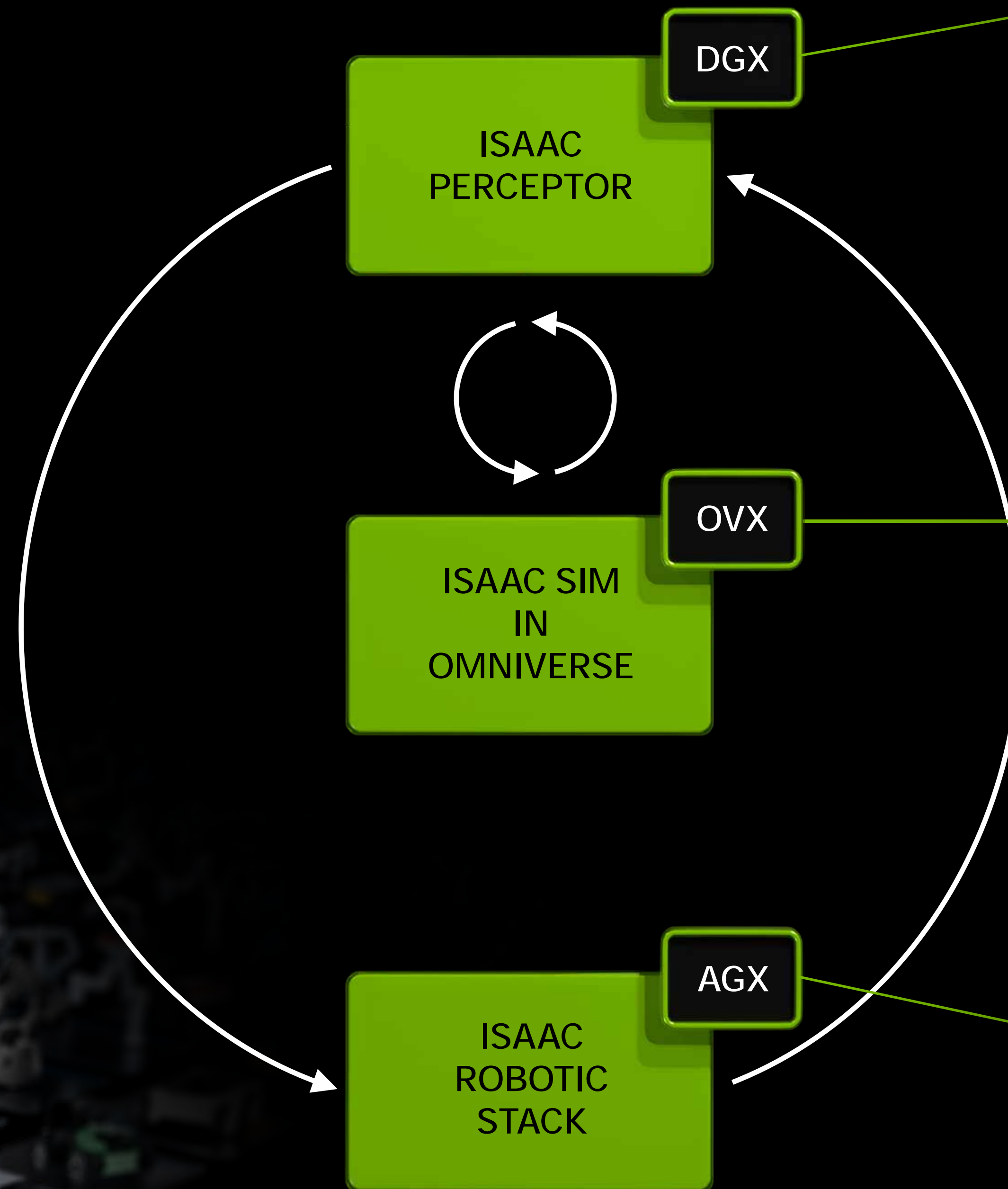
NVIDIA Robotics Platform Adoption

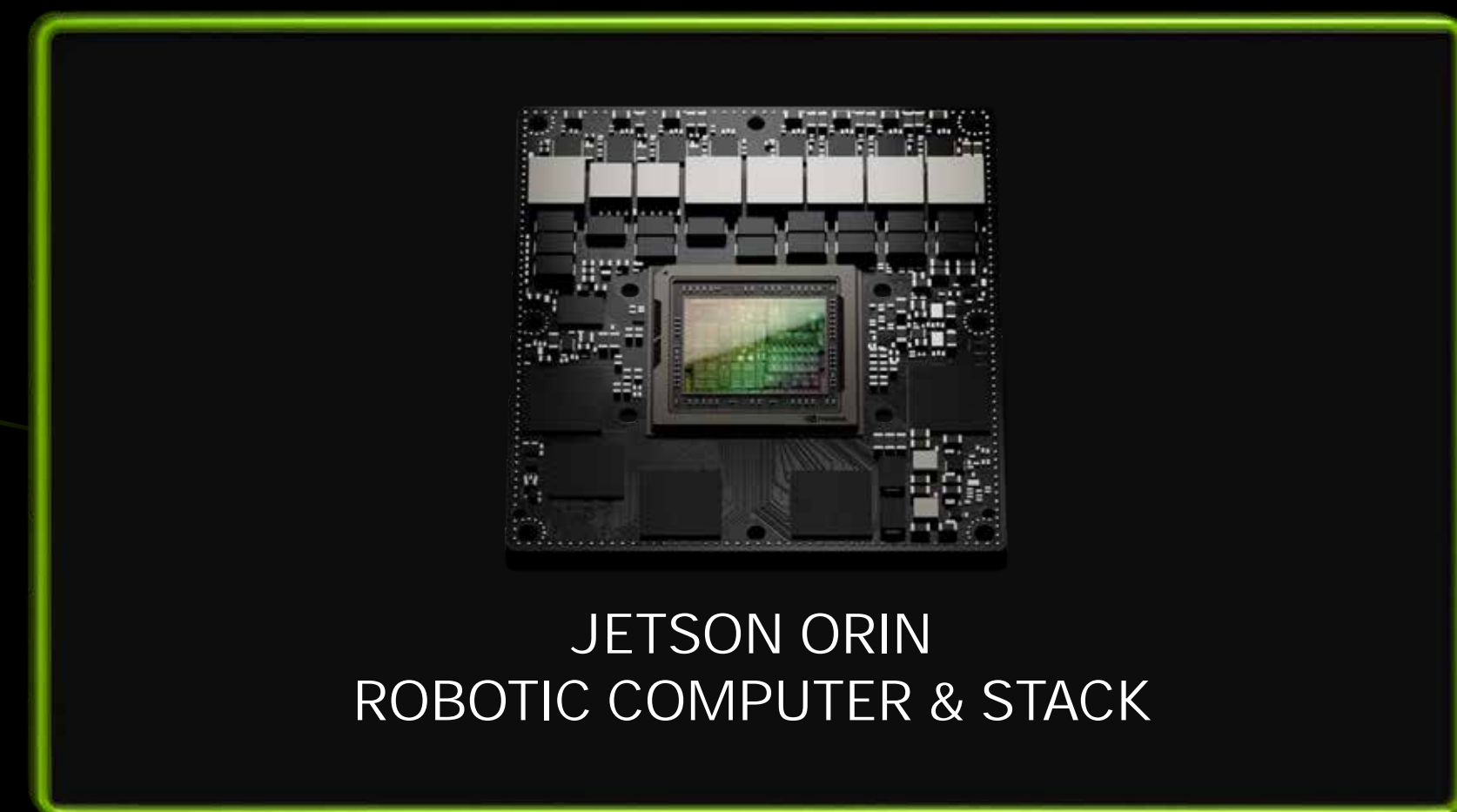
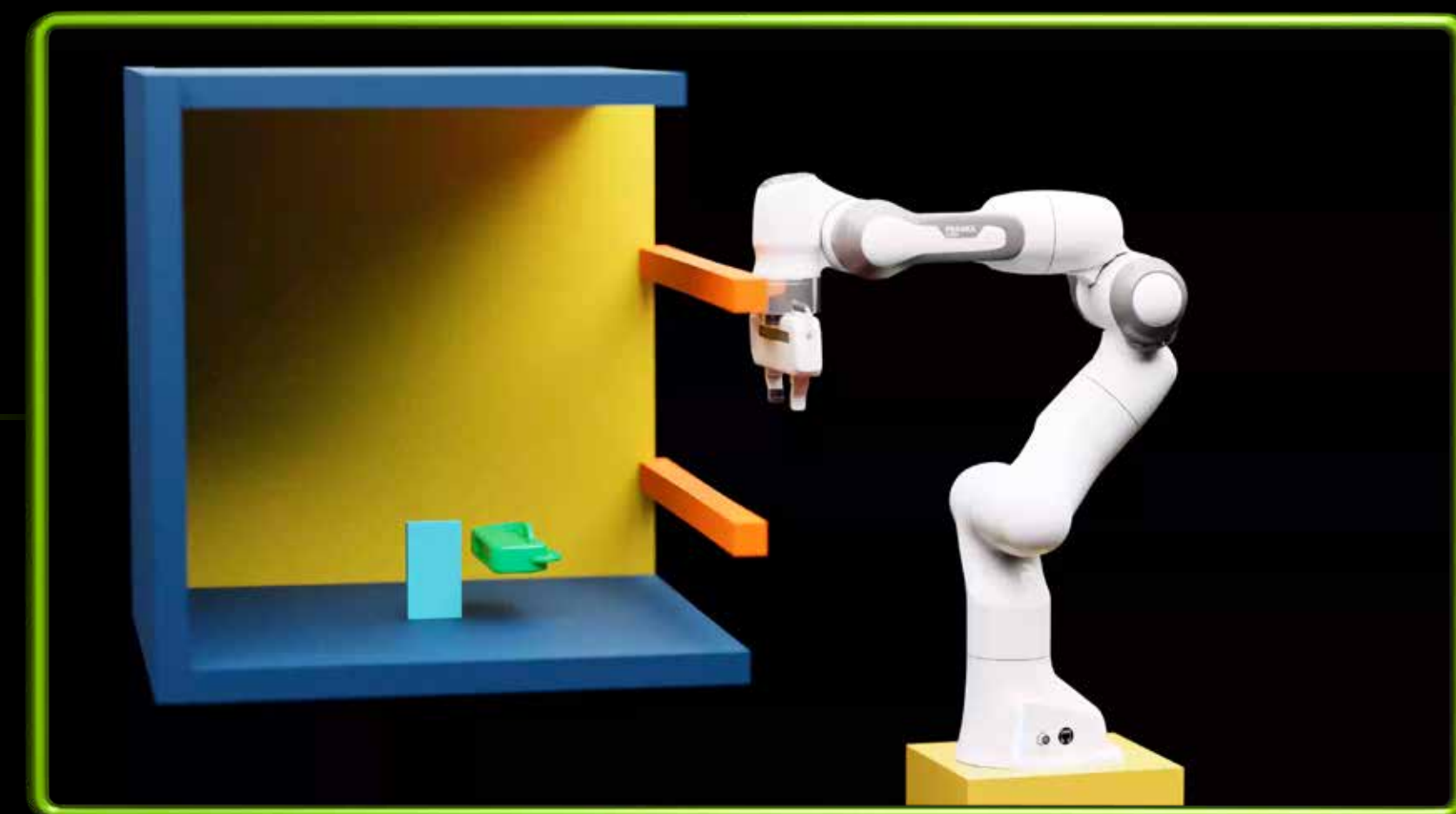
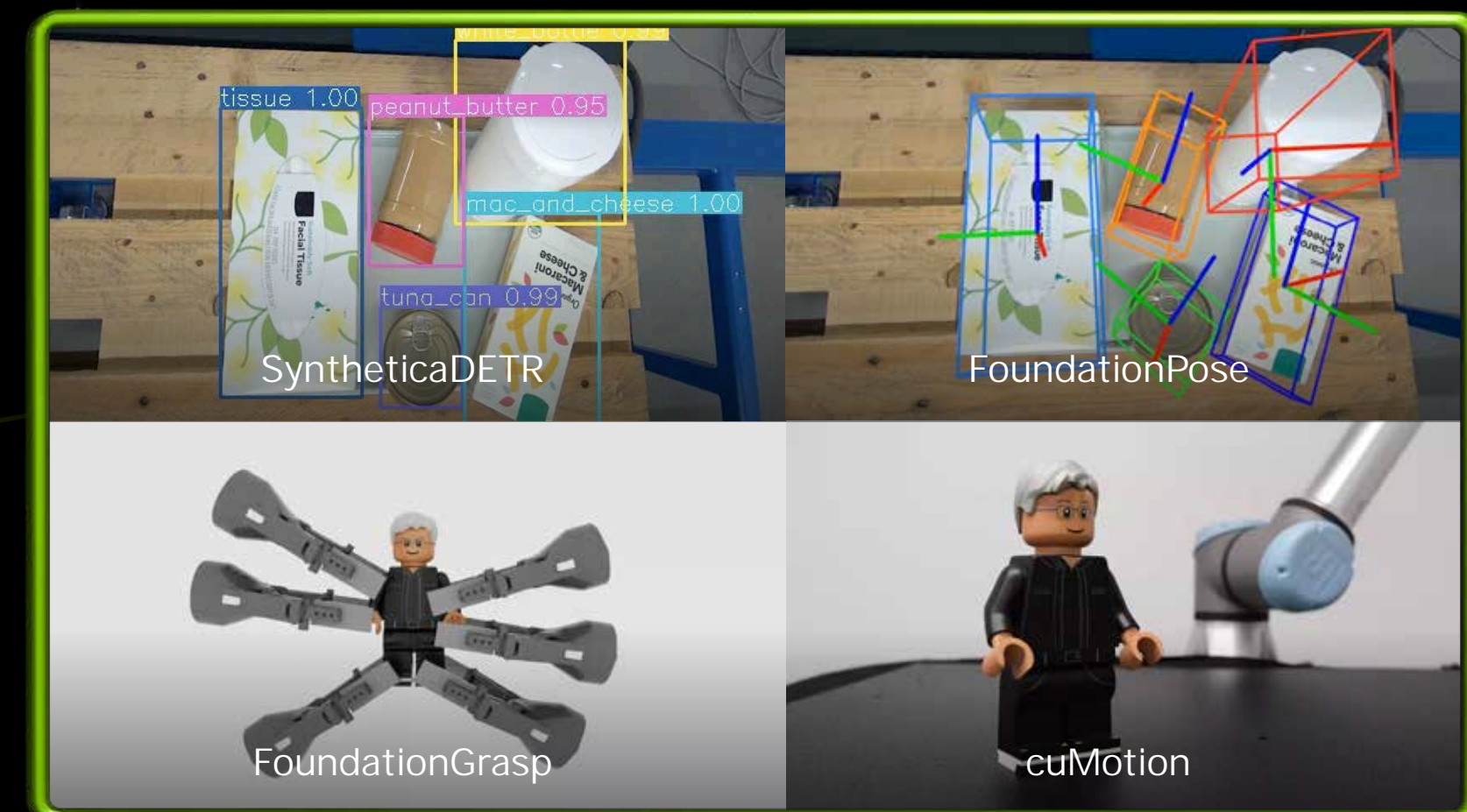
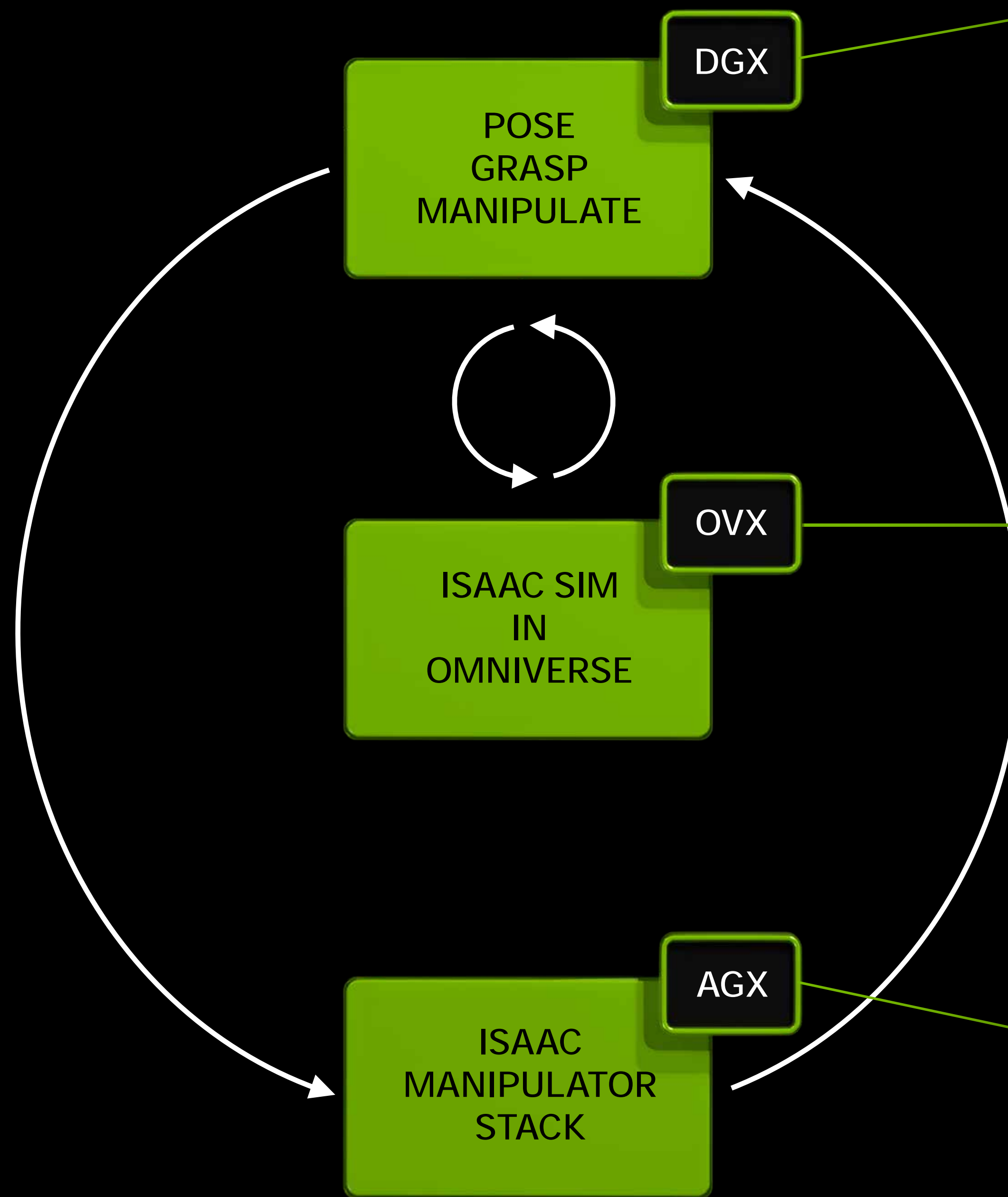
1.3M

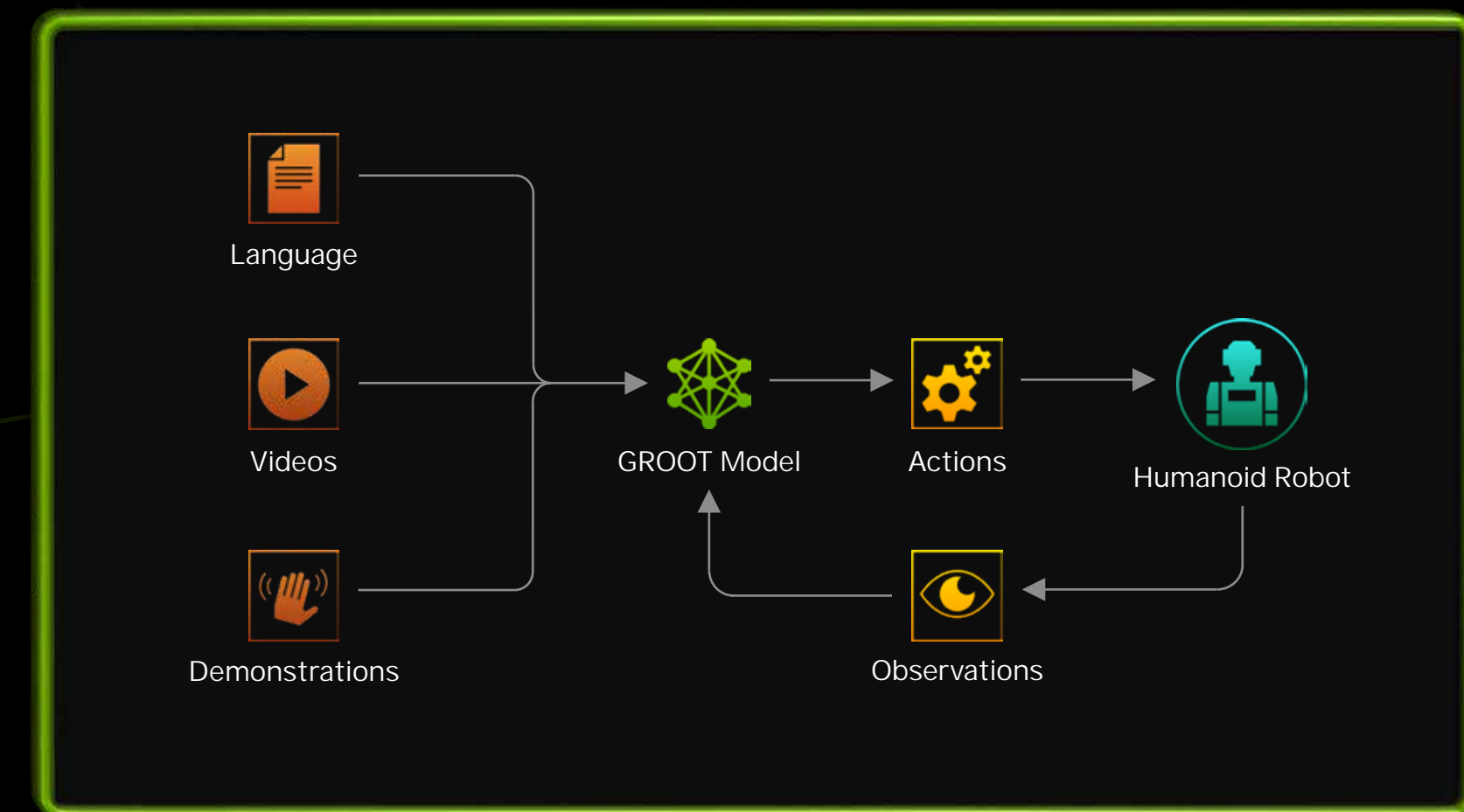
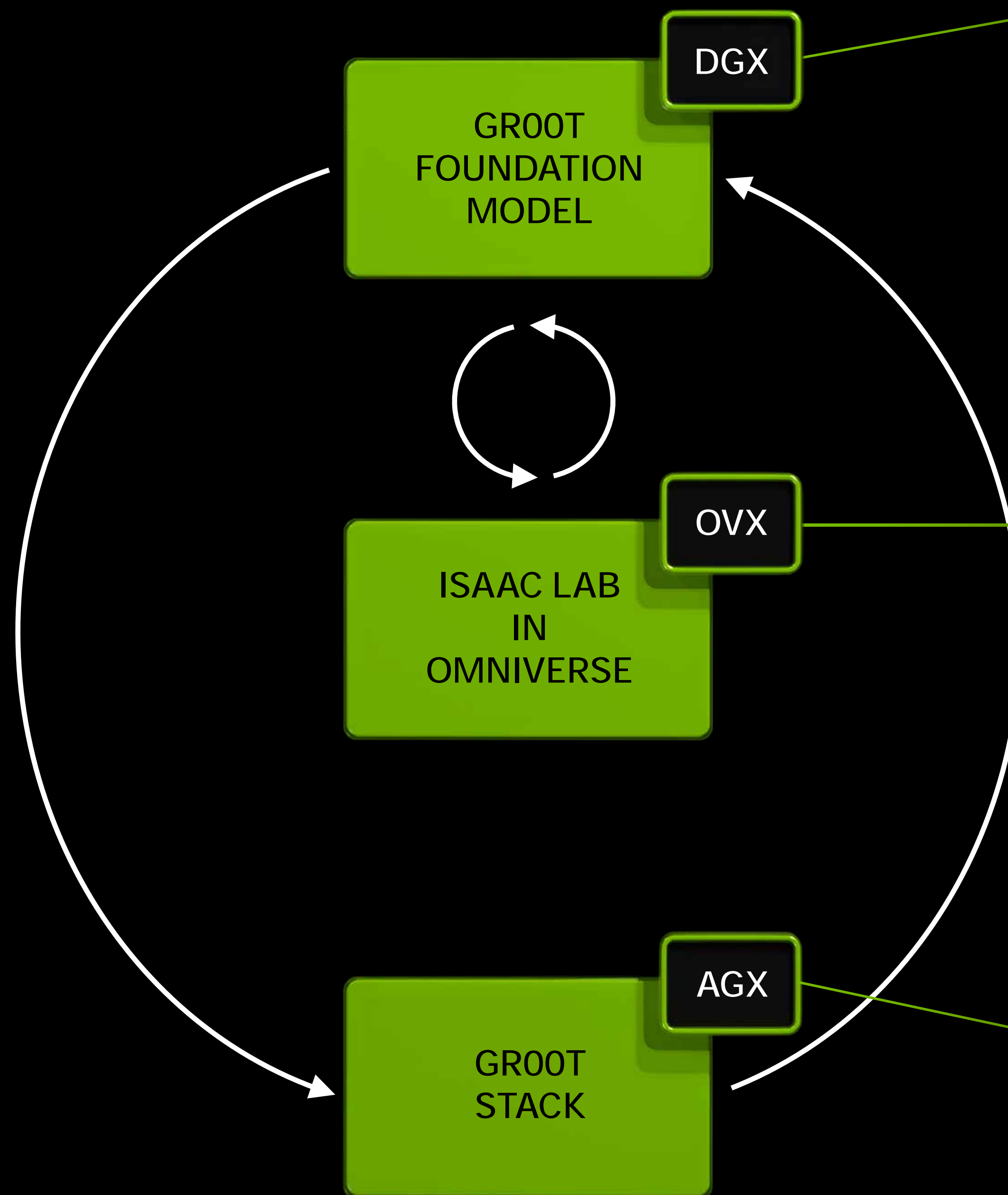
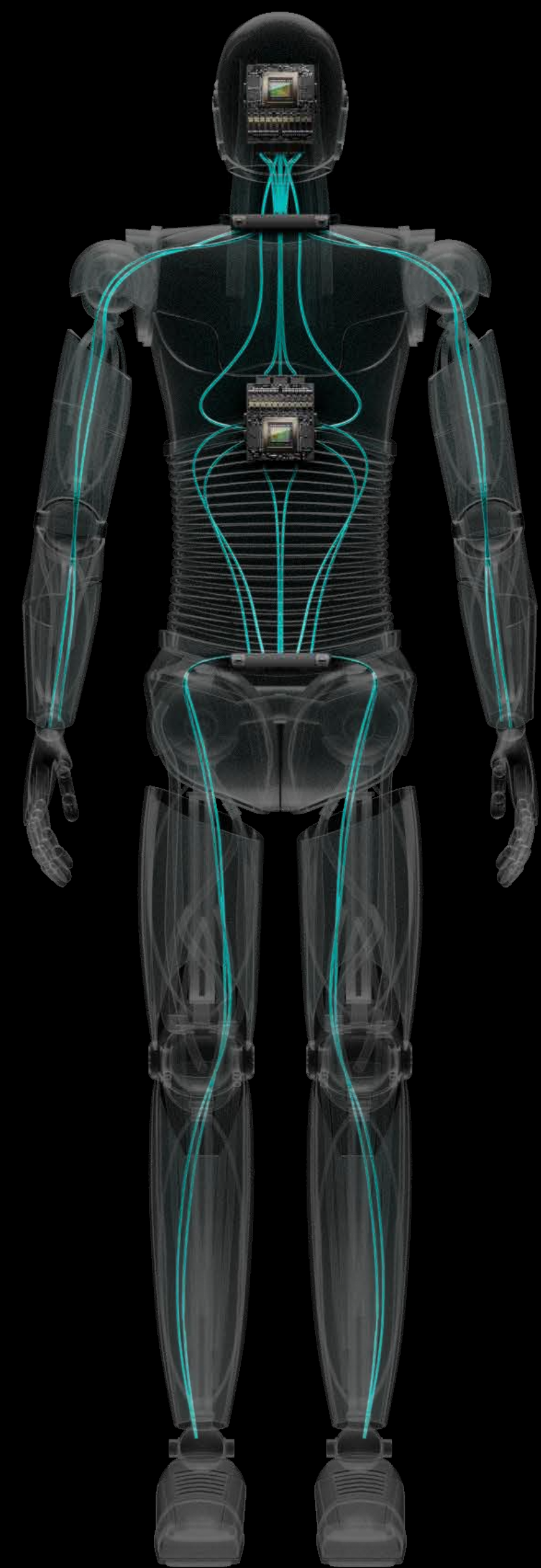
Robotics Developers

100,000 ROS Developers

6,000 Companies Developing on Orin



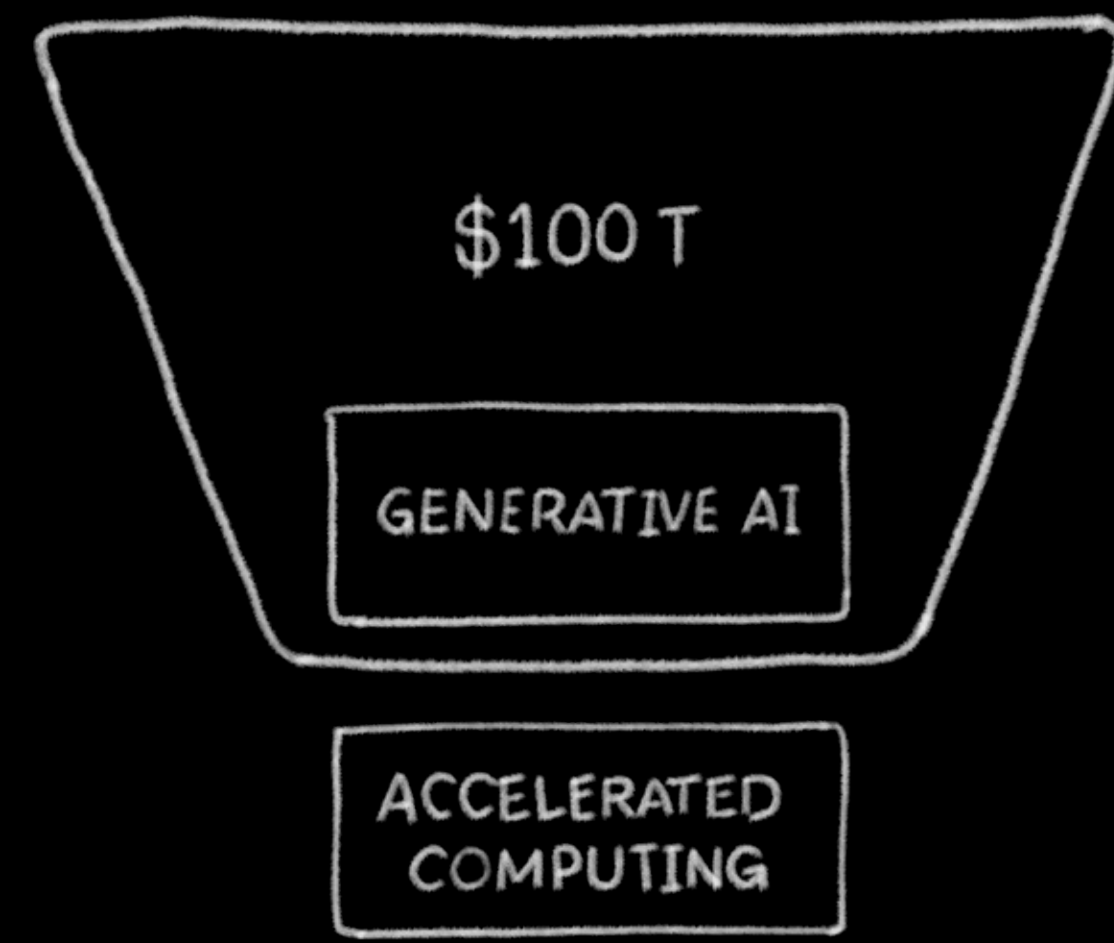




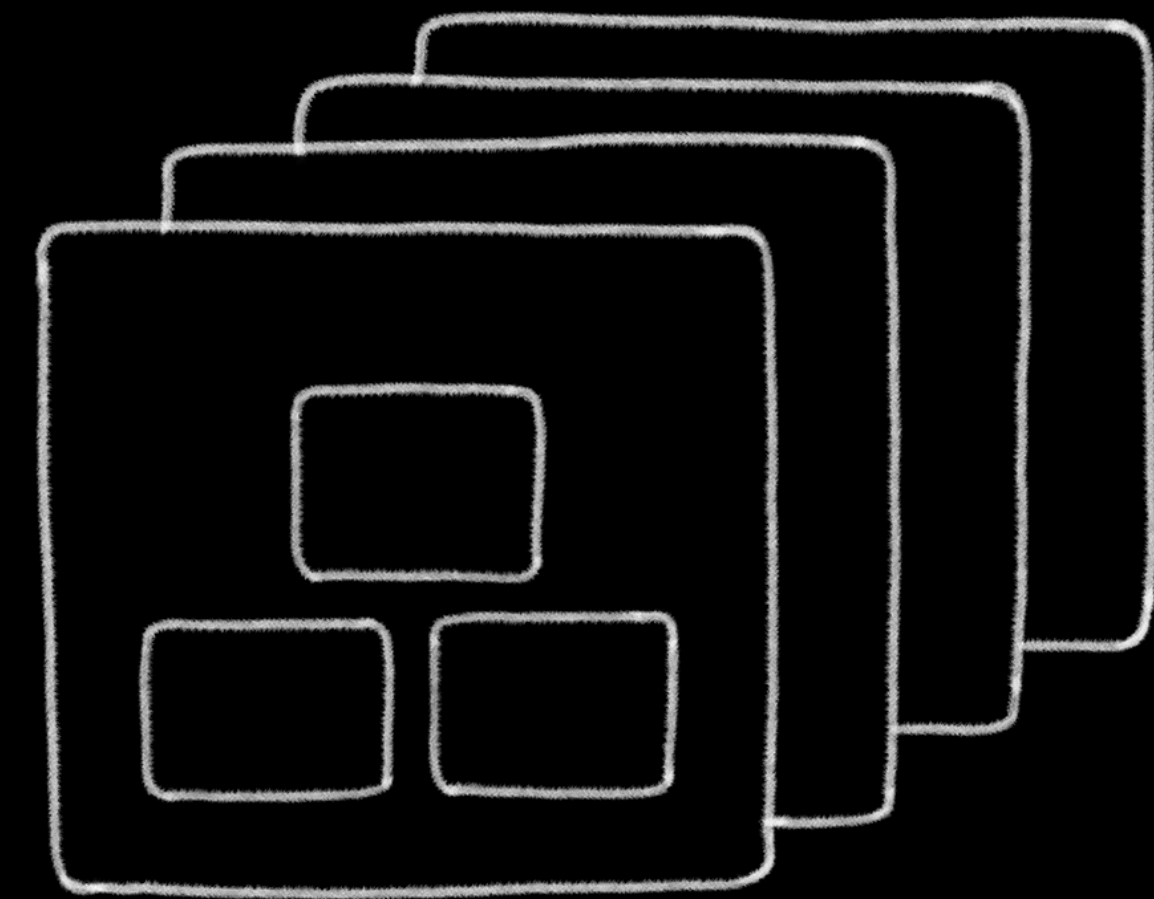


A 01

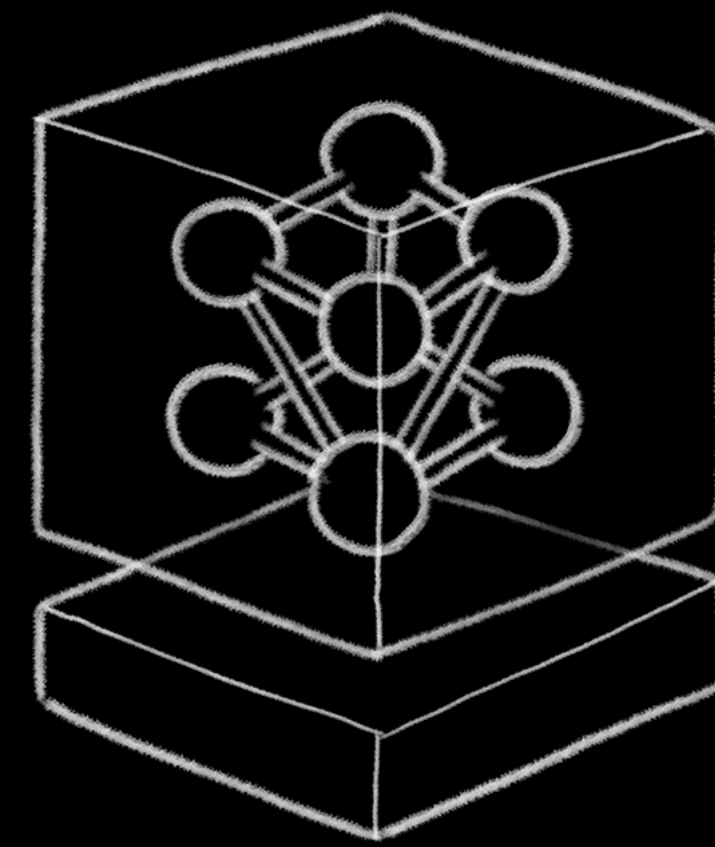
A NEW INDUSTRIAL REVOLUTION



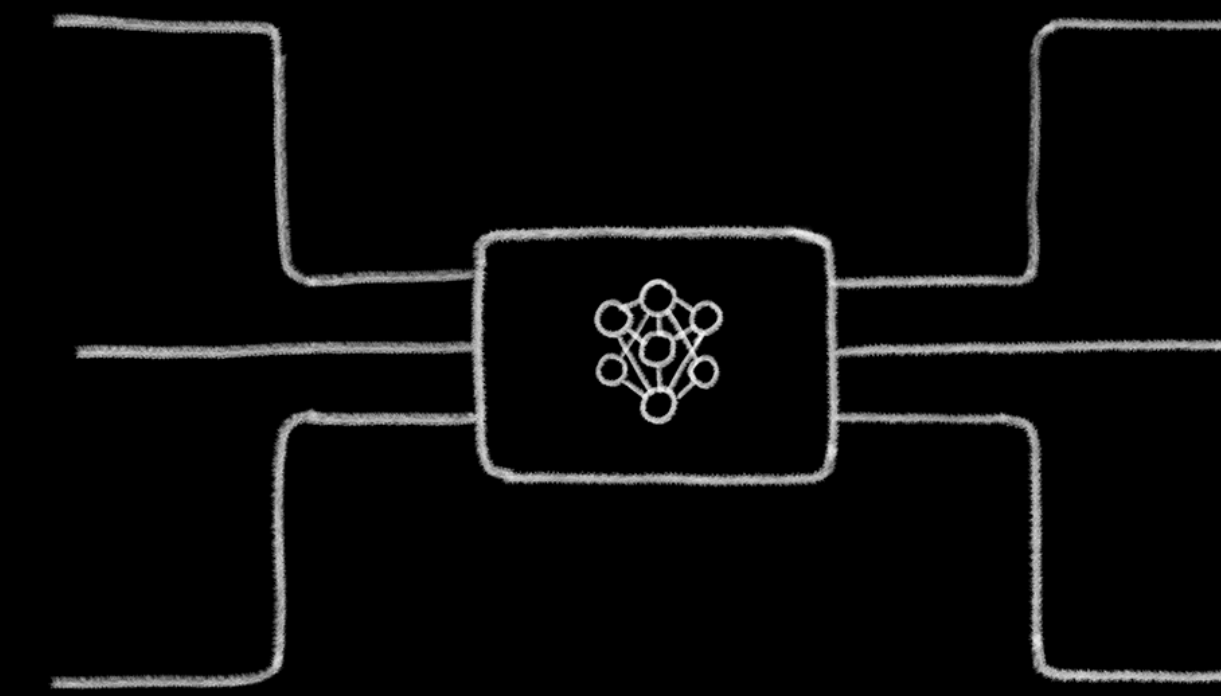
NEW INDUSTRY



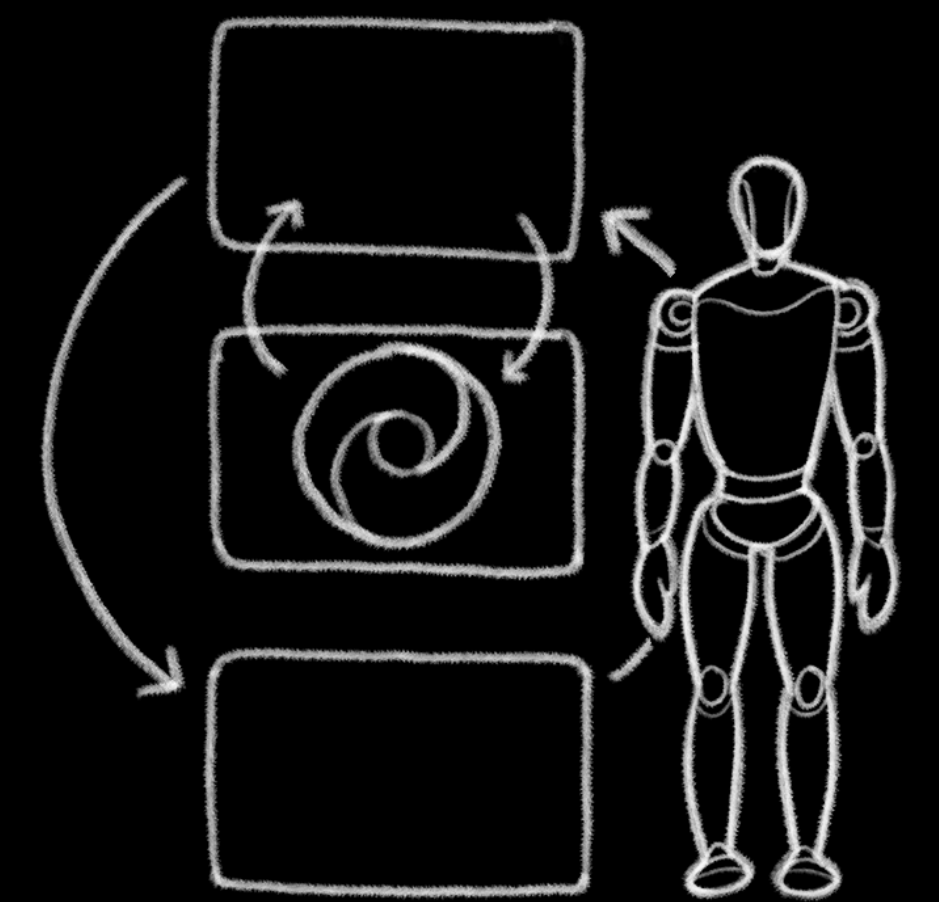
BLACKWELL
PLATFORM



NIMs

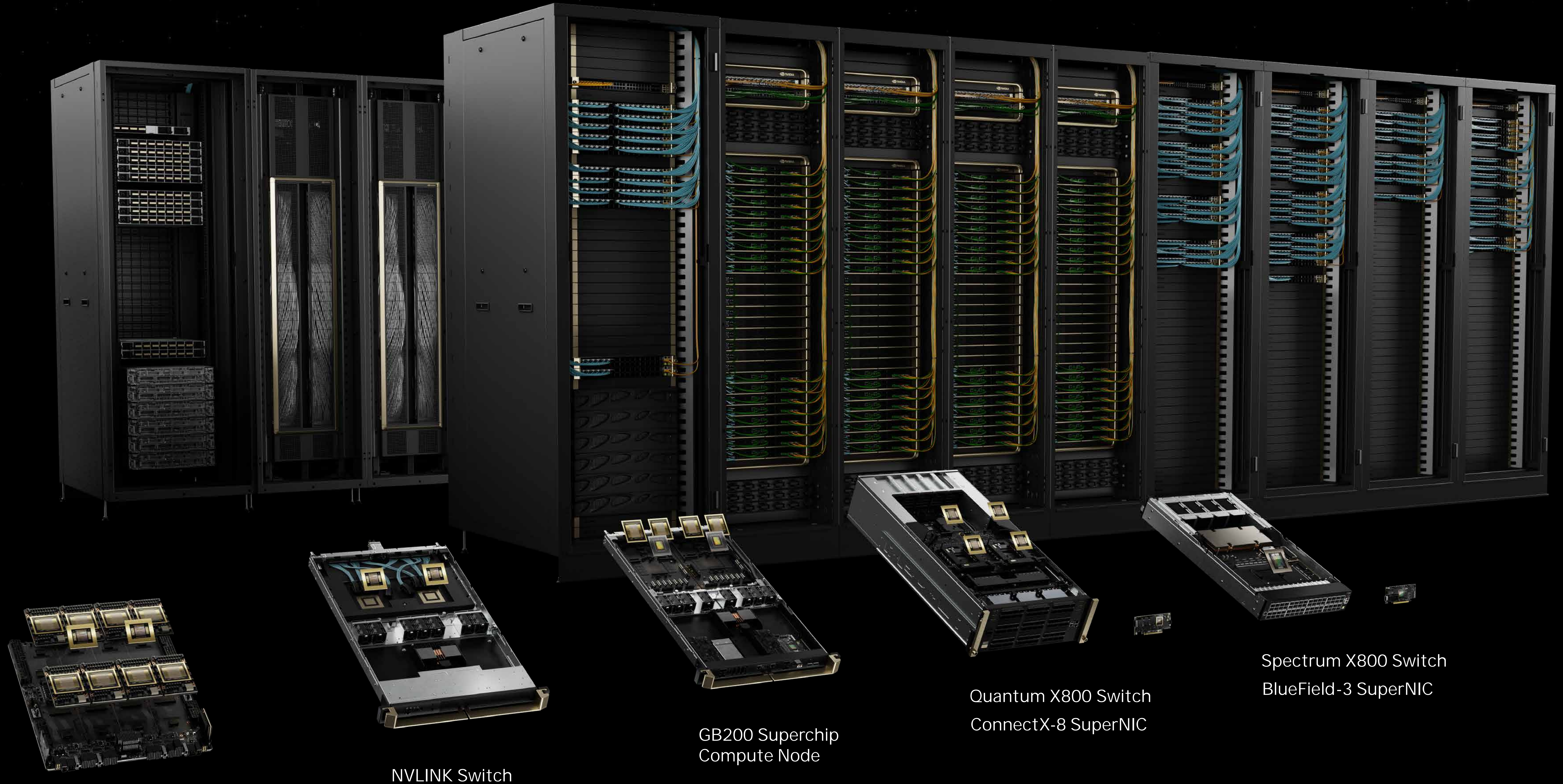


NEMO AND
NVIDIA AI FOUNDRY



OMNIVERSE AND
ISAAC ROBOTICS

NVIDIA Blackwell Platform



HGX B100

NVLINK Switch

GB200 Superchip
Compute Node

Quantum X800 Switch
ConnectX-8 SuperNIC

Spectrum X800 Switch
BlueField-3 SuperNIC

