# III FÓRUM DE INOVAÇÃO

# AIOT: AMPLIANDO AS FRONTEIRAS DA INTELIGÊNCIA ECONECTIVIDADE

Realização:

















#### III FÓRUM DE INOVAÇÃO

AIOT: AMPLIANDO AS FRONTEIRAS DA INTELIGÊNCIA E CONECTIVIDADE

# Plataformas NVIDIA para Desenvolvimento em HPC e IA

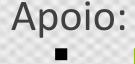
Fabio Alves de Oliveira – NVIDIA America Latina Higher Education and Research

Realização:















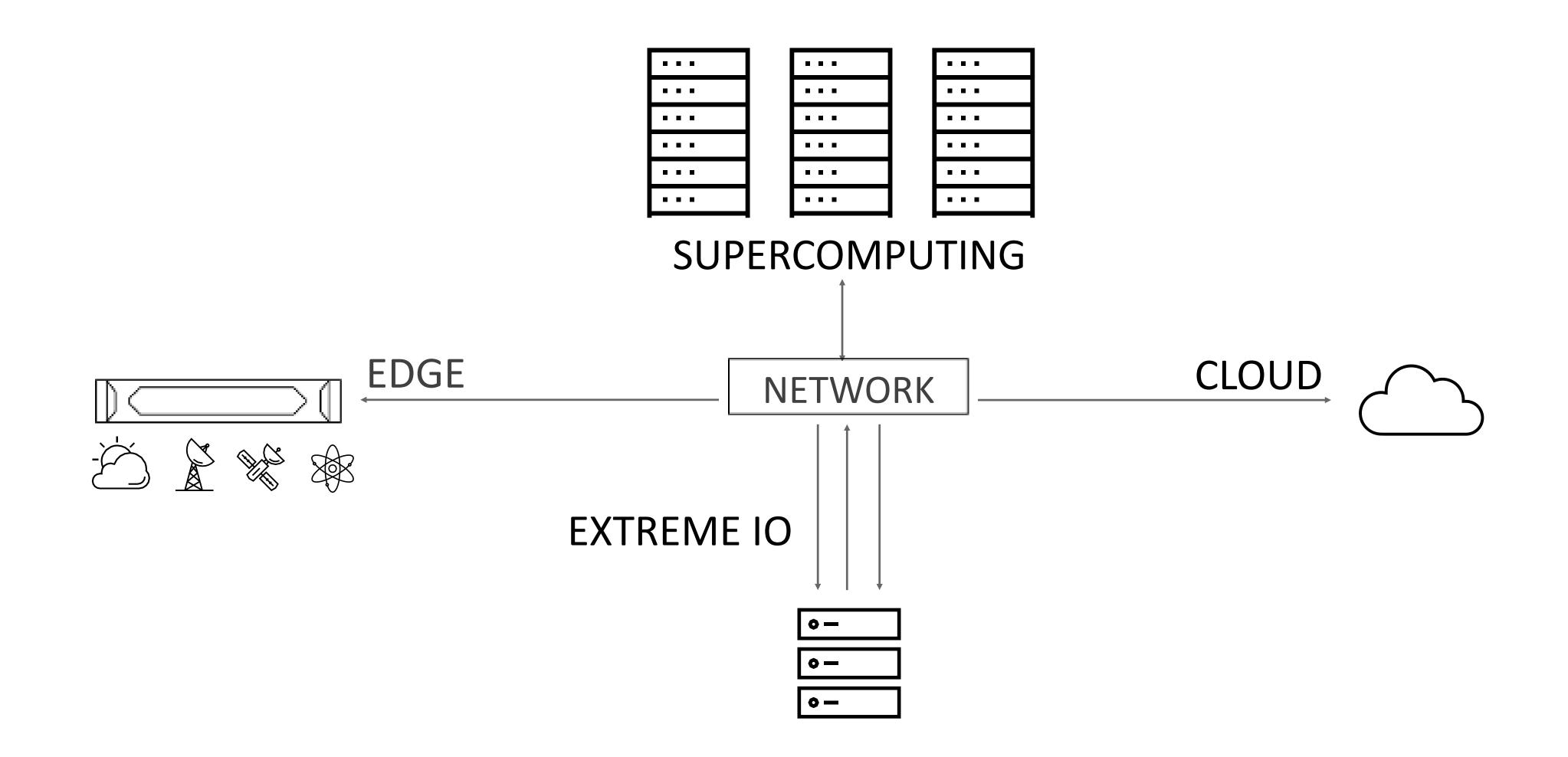


# Workloads of the Modern Supercomputer

EDGE SIM + AI SIMULATION DIGITAL TWIN QUANTUM COMPUTING

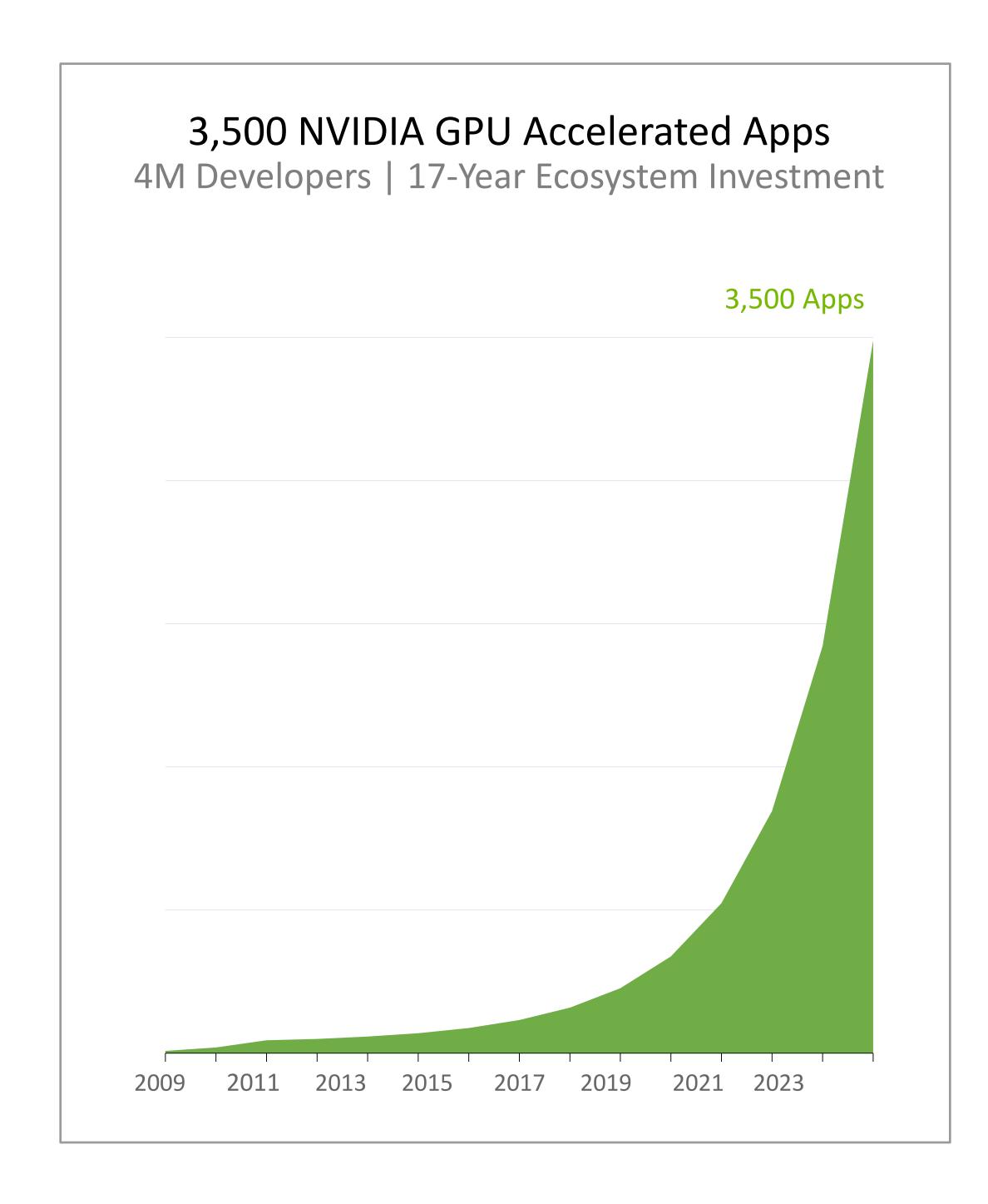
\*\*COMPUTING\*\*

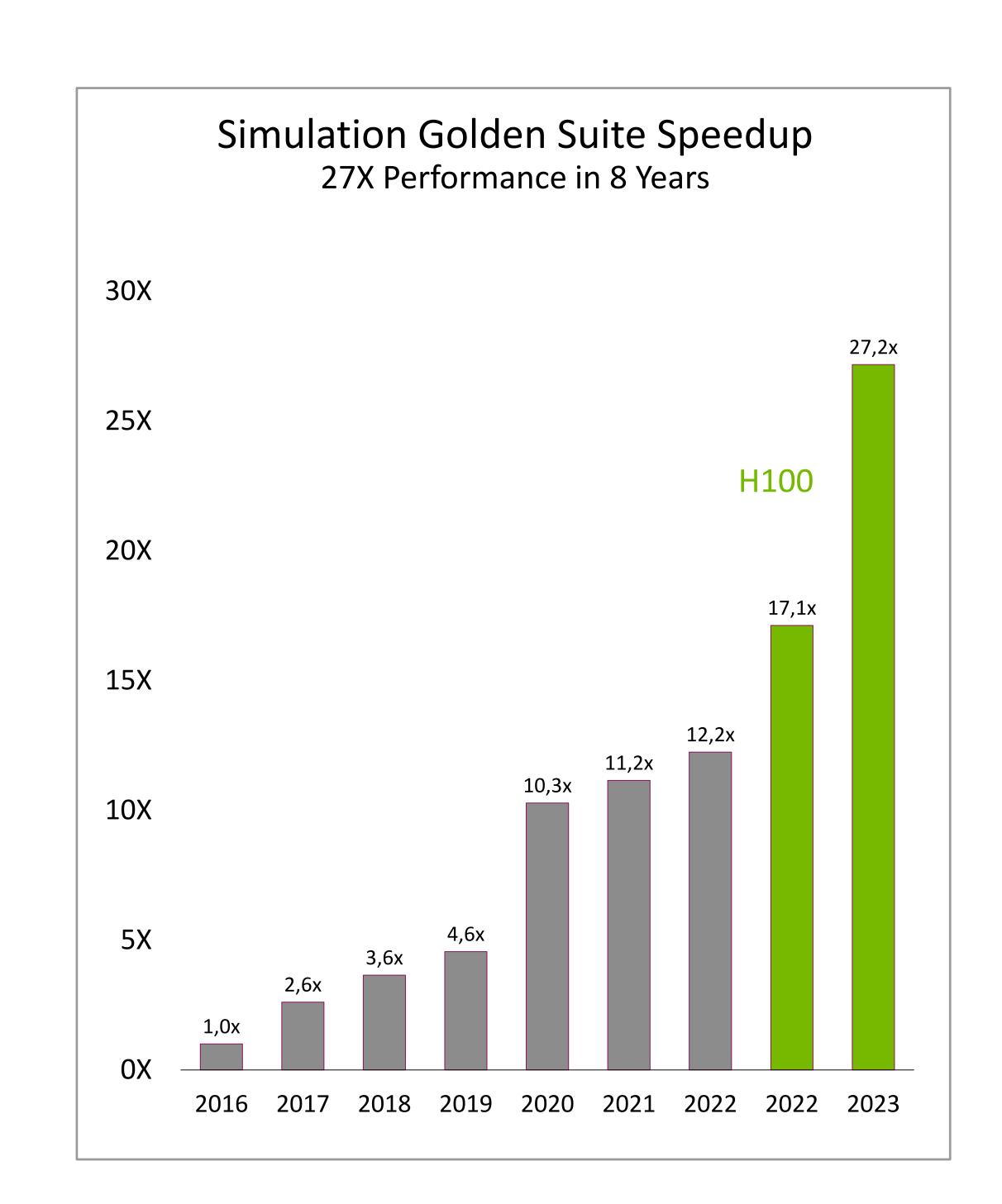
\*\*COMPUTI

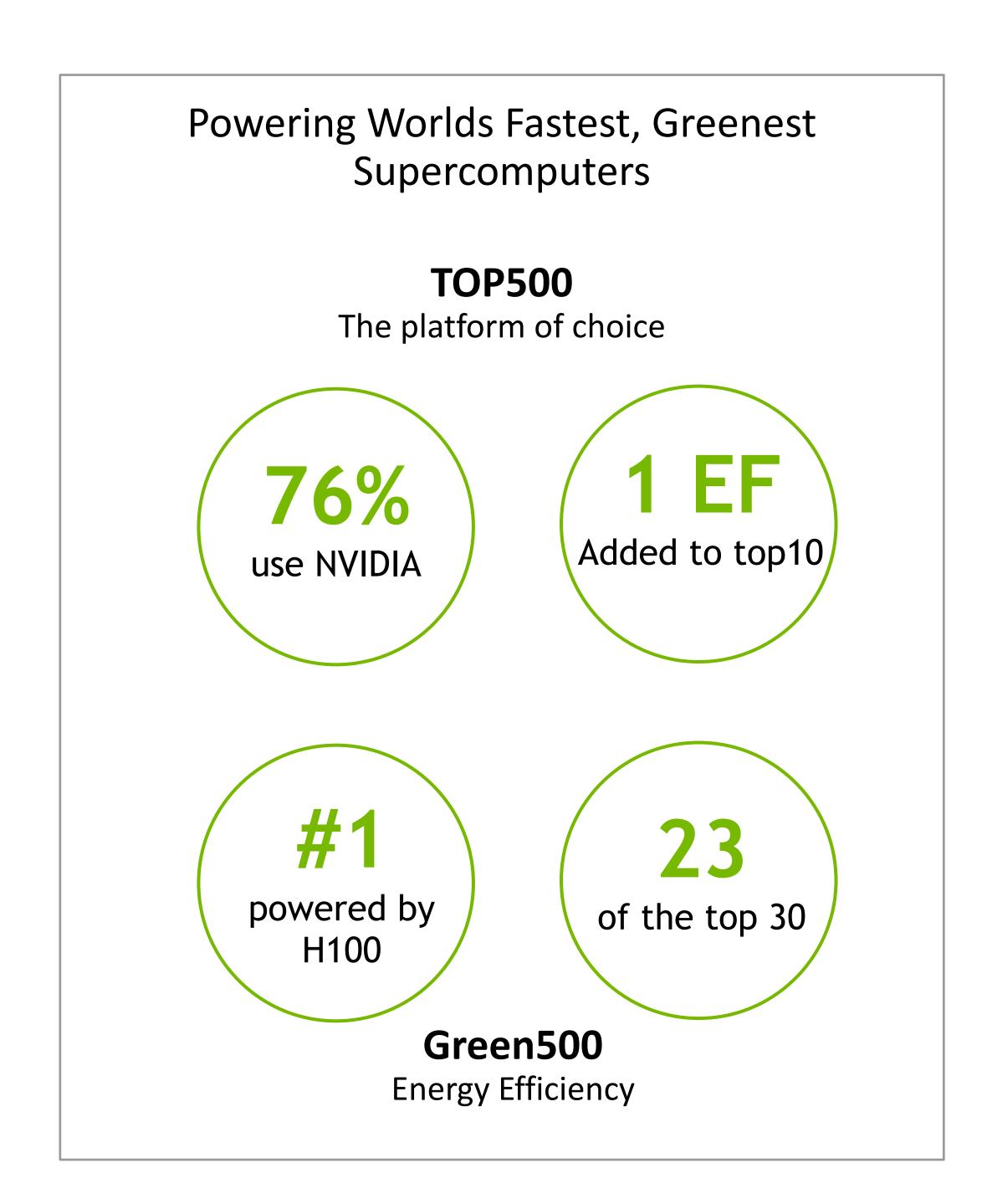


# World's Leading Scientific Computing Platform

Delivering Breakthrough Application Performance







# **HPC Reinvented with AI**

	Experiments Simulation Viz	Edge HPC + AI Simulation Digital Twin Computing    Computing
FEATURE	PRE-EXASCALE	EMERGING POST EXA-SCALE
USAGE	BATCH	INTERACTIVE & DISTRIBUTED
WORKLOAD	SINGLE SIMULATION/ENSEMBLES	SIMULATION/ENSEMBLES, AI TRAINING AND INFERENCE
EXPERIMENTS	OFFLINE DATA ANALYSIS FOR EXPERIMENTS	MIX OF REAL-TIME ANALYSIS, STEERING AND OFFLINE
DIGITAL TWINS	IN-SITU VISUALIZATION	INTERACTIVE COMBINATION OF SIMULATION AND OBSERVATIONAL DATA
QUANTUM COMPUTING	SIMULATION	PREPARING FOR A HYBRID MODEL
PROGRAMMING MODELS	FORTRAN, C++, MPI, OPENMP	STANDARD PARALLELISM SUPPORT IN FORTRAN, C++, MPI, OPENMP, OPENACC, PYTHON, JULIA, PYTORCH, JAX, TENSORFLOW
CLOUD	GRID	BURST CAPABILITIES, FASTER REFRESH CYCLE, ACCESS TO LATEST TECHNOLOGY AT SCALE

Quantum

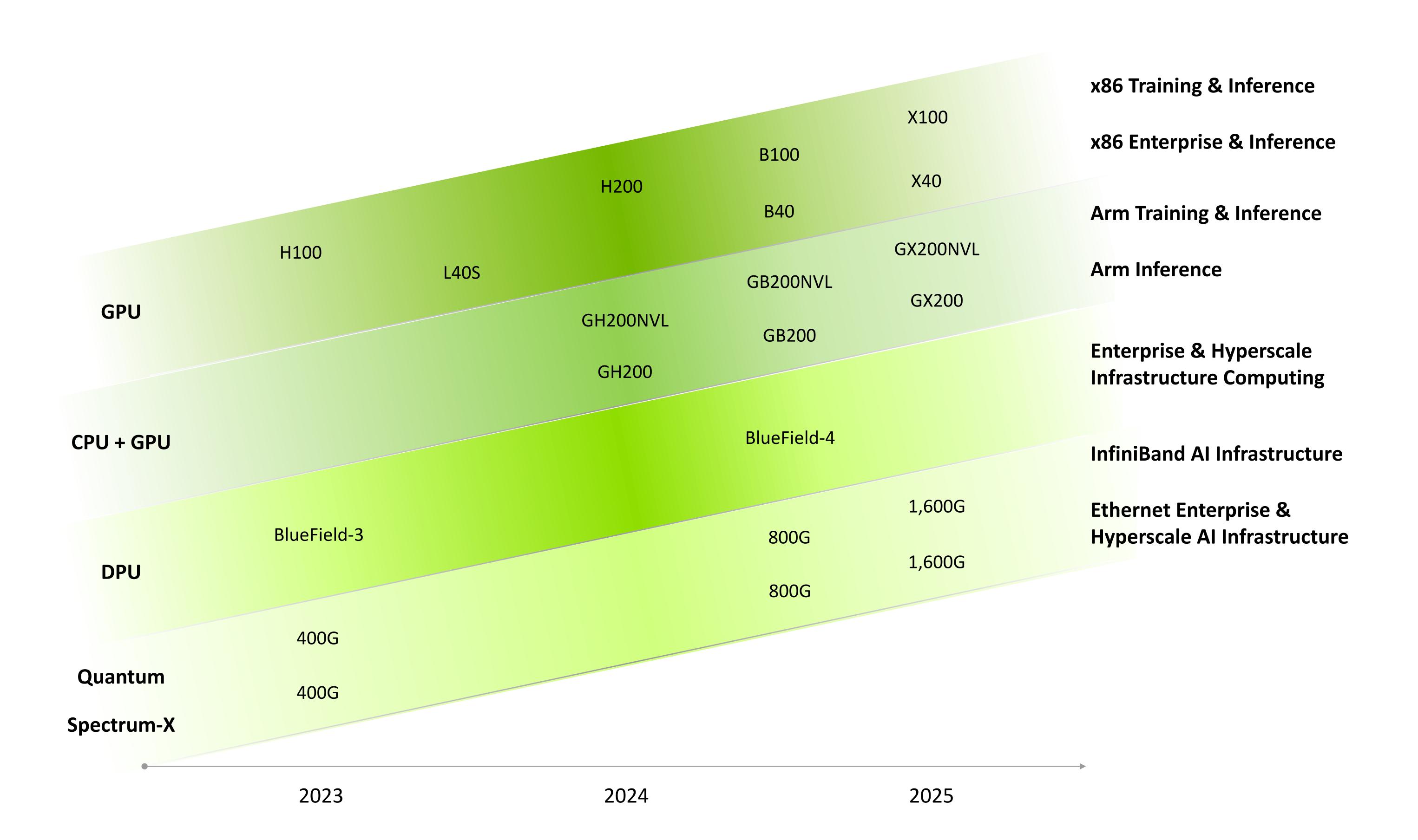


HOLOSCAN

METROPOLIS

# **NVIDIA AI - One Architecture | Train and Deploy Everywhere**

One-Year Rhythm





# **NVIDIA GH200 Grace Hopper Superchip**

Built for the New Era of Al Supercomputing

**CPU to GPU Bandwidth** 

900GB/s

NVLink-C2C

**GPU Memory Bandwidth** 

5TB/s

HBM3e

**Energy Efficiency** 

50X

MILC Efficiency vs 2S x86 CPUs

**QFT Quantum Simulation** 

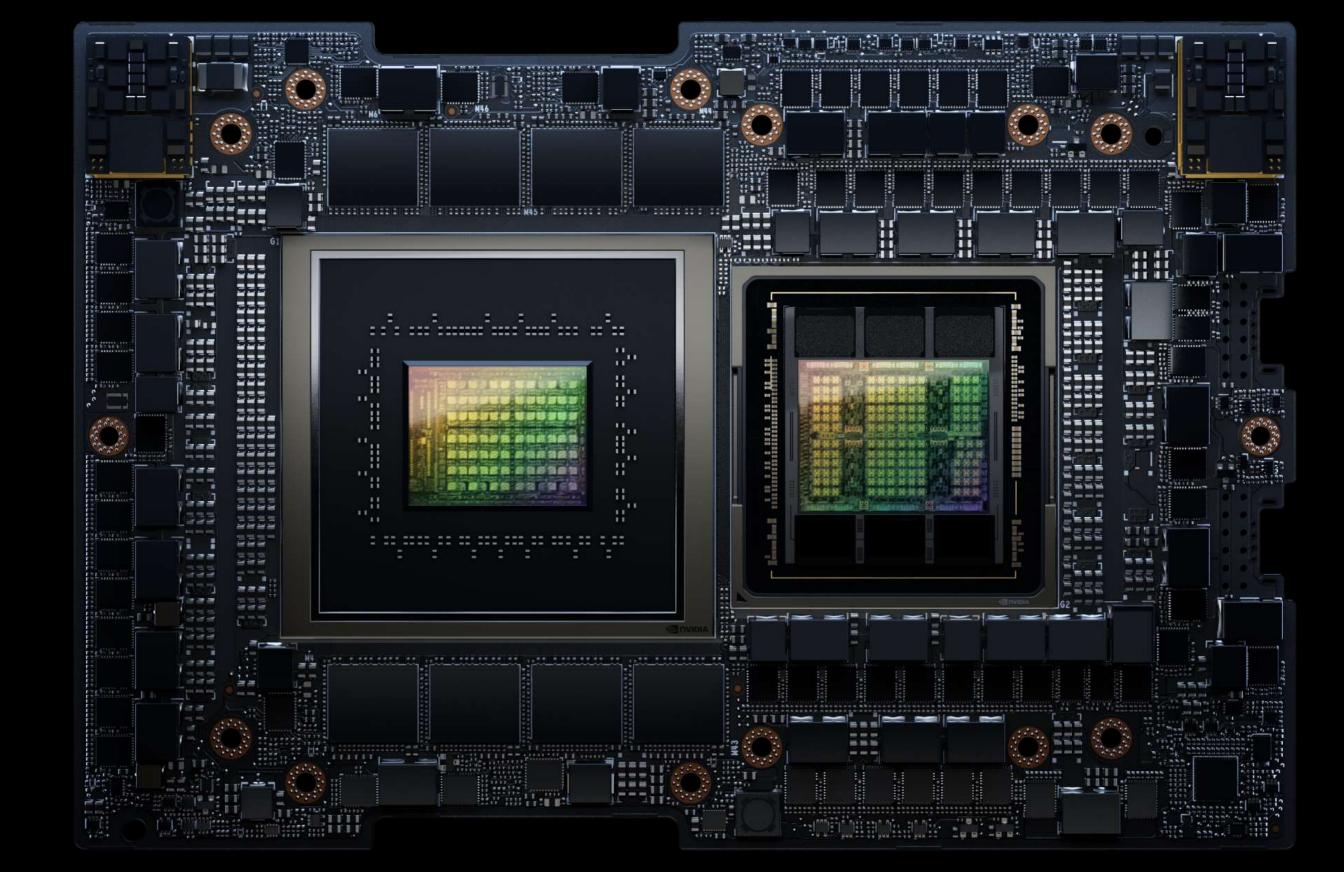
90X

Performance vs 2S x86 CPUs

**LLM Inference** 

200X

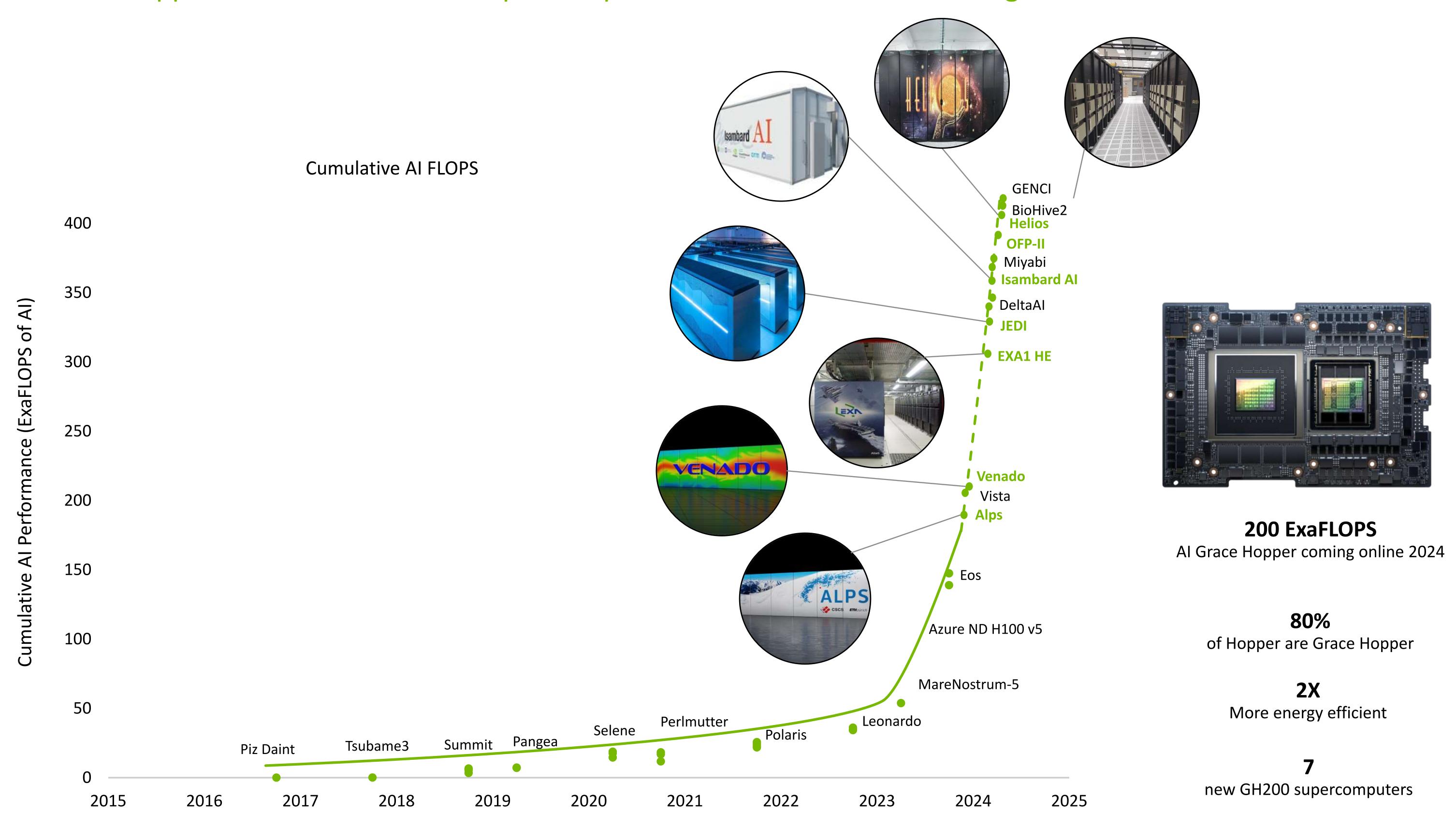
Performance vs H100 80GB



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

## Grace Hopper Powers Al Supercomputing Datacenters

Grace Hopper Will Deliver 200 Exaflops of AI performance for Groundbreaking Research





# **NVIDIA Grace CPU Superchip**

Breakthrough Performance and Efficiency for the Modern Data Center

**CPU + Memory Power** 

500W

**Grace CPU Superchip TDP** 

**Memory Bandwidth** 

1 TB/s

LPDDR5X

Green 500

7 Grace systems

Performance, Energy Efficiency with GH200

**Energy Efficiency** 

**2X** 

Performance vs x86 CPU

Graph Analytics: GAP BS Breadth First Search

Weather

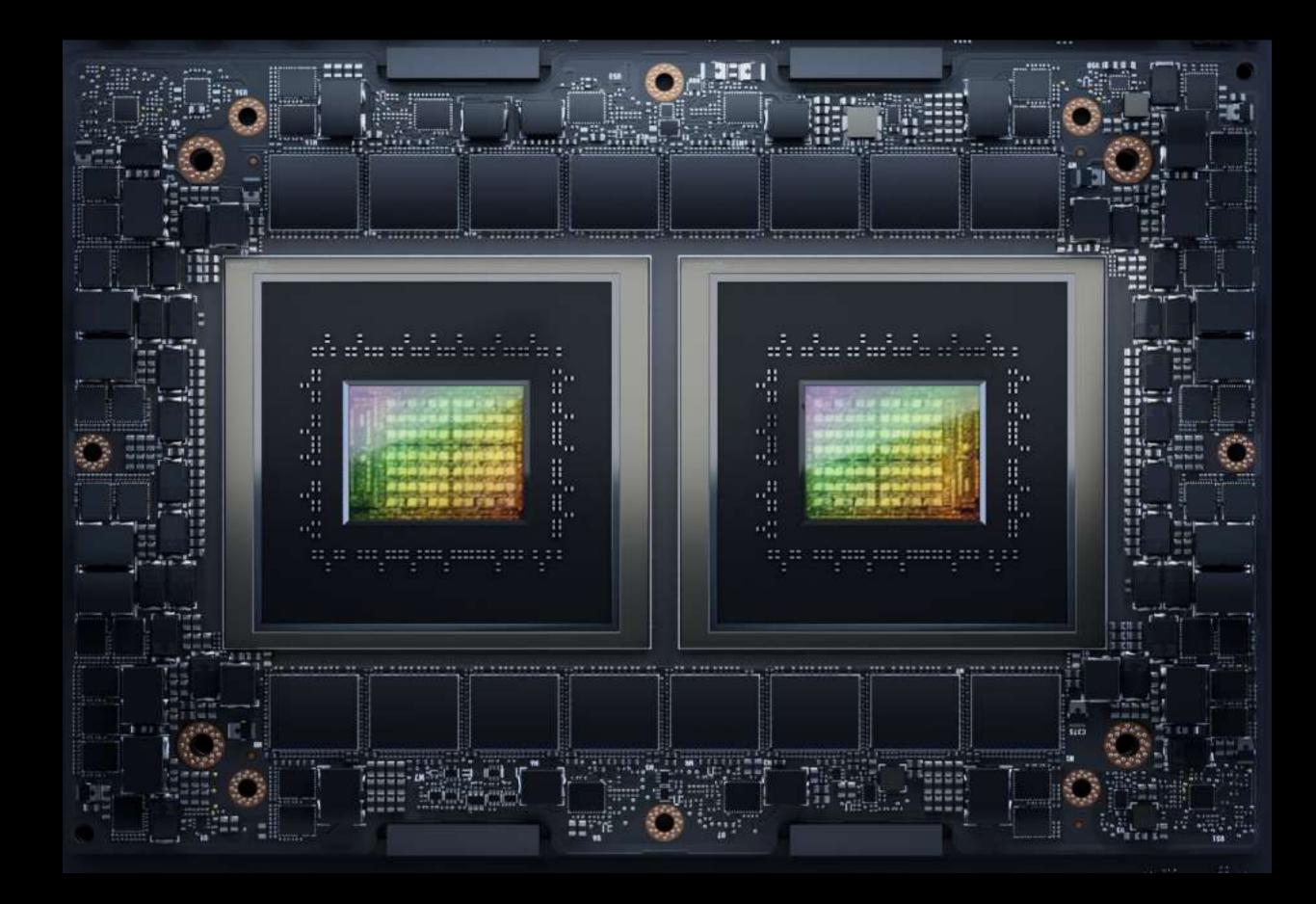
1.3X

Performance vs x86 CPU

**Graph Analytics** 

**2X** 

Performance vs x86 CPUs

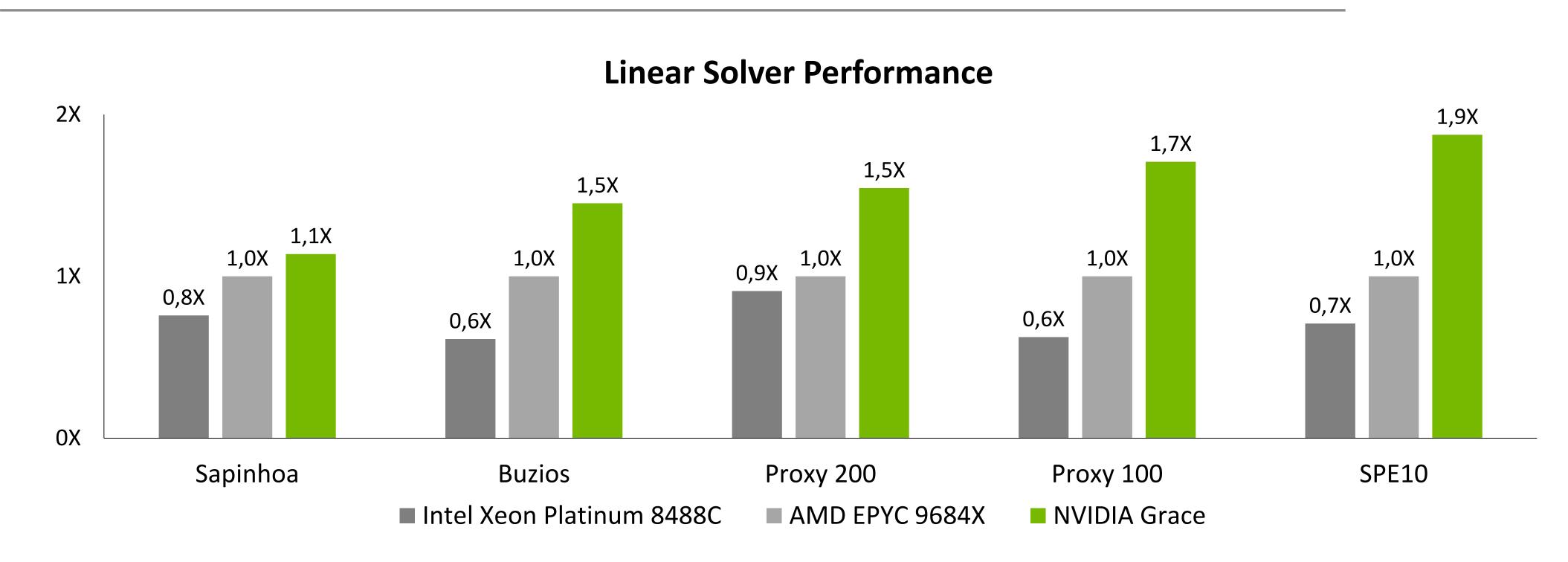


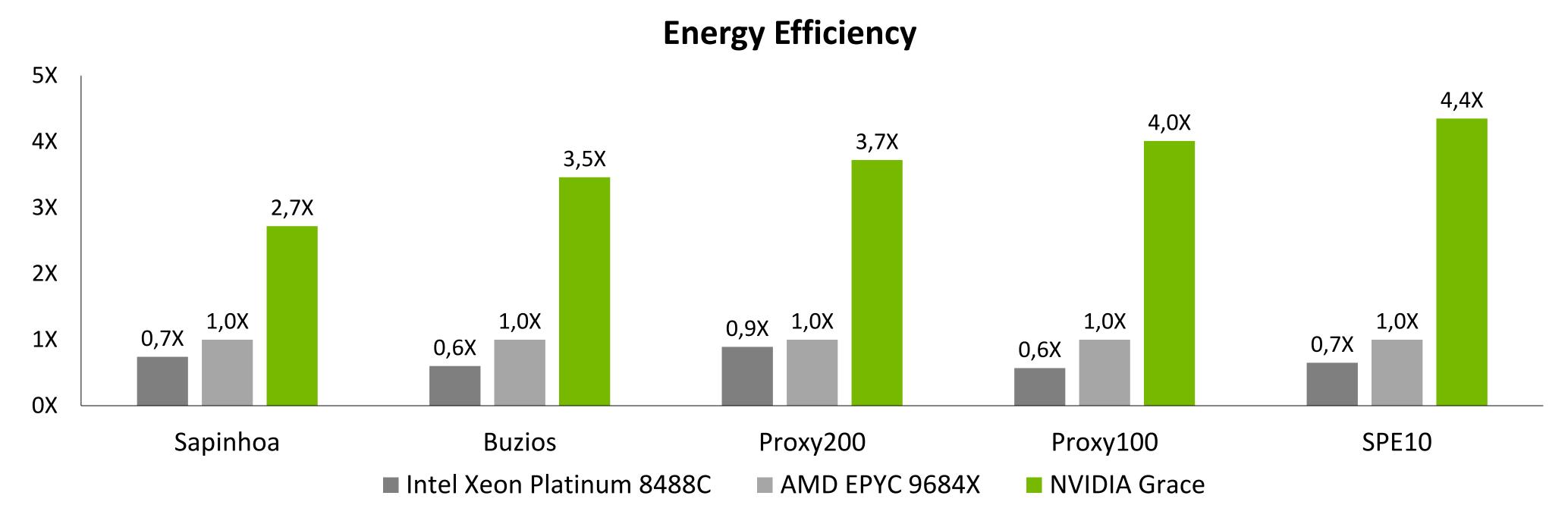
144 Arm Neoverse V2 Cores | 234MB L3 Cache 3.2 TB/s NVIDIA Scalable Coherency Fabric | 960GB LPDDR5X



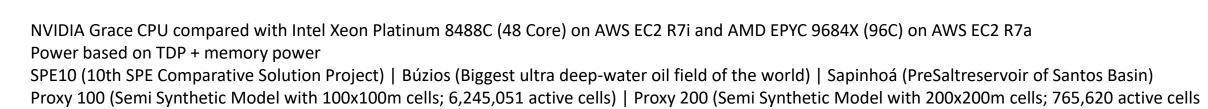
## Petrobas Reservoir Simulation

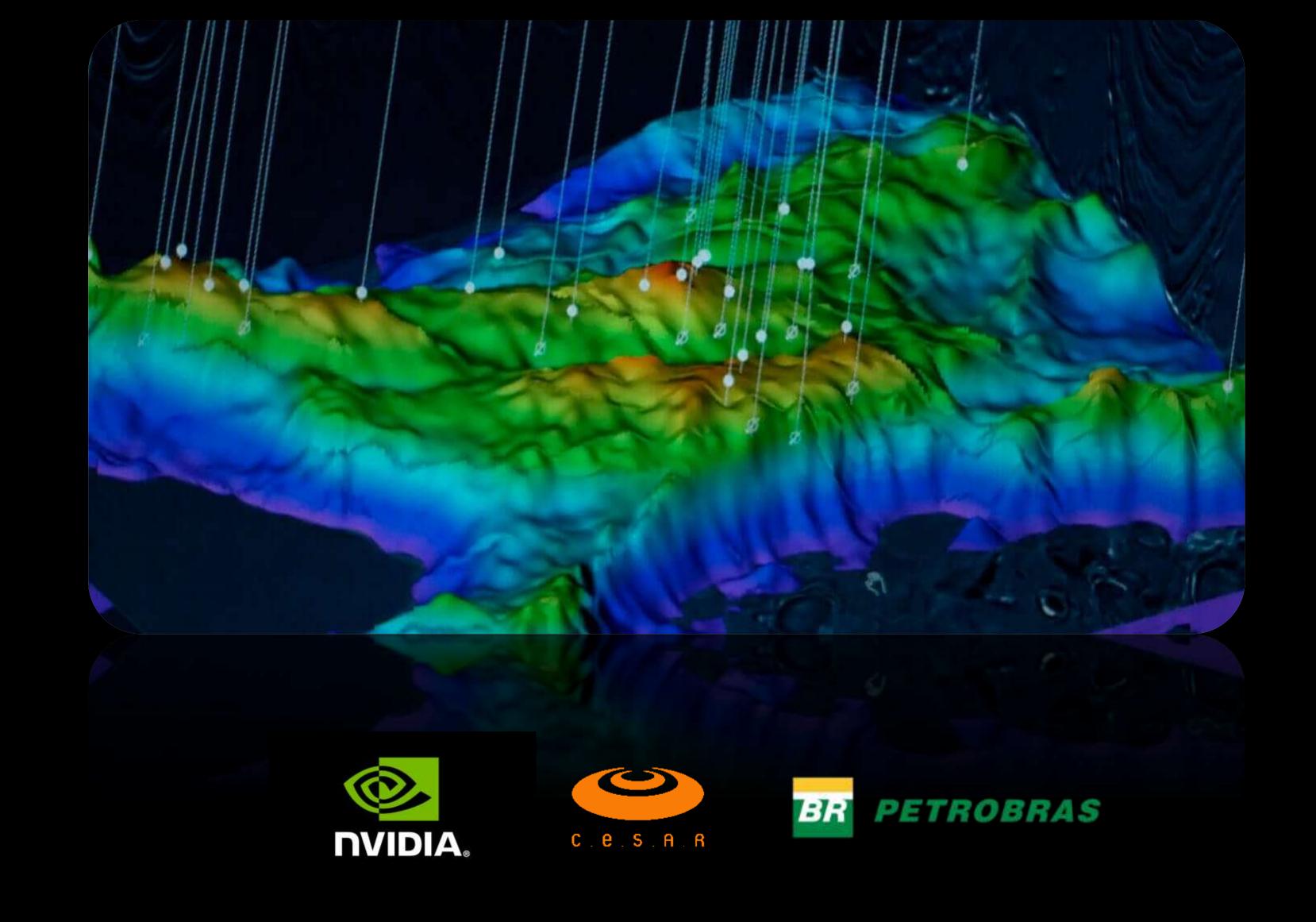
NVIDIA Grace CPU delivers over 4X energy efficiency







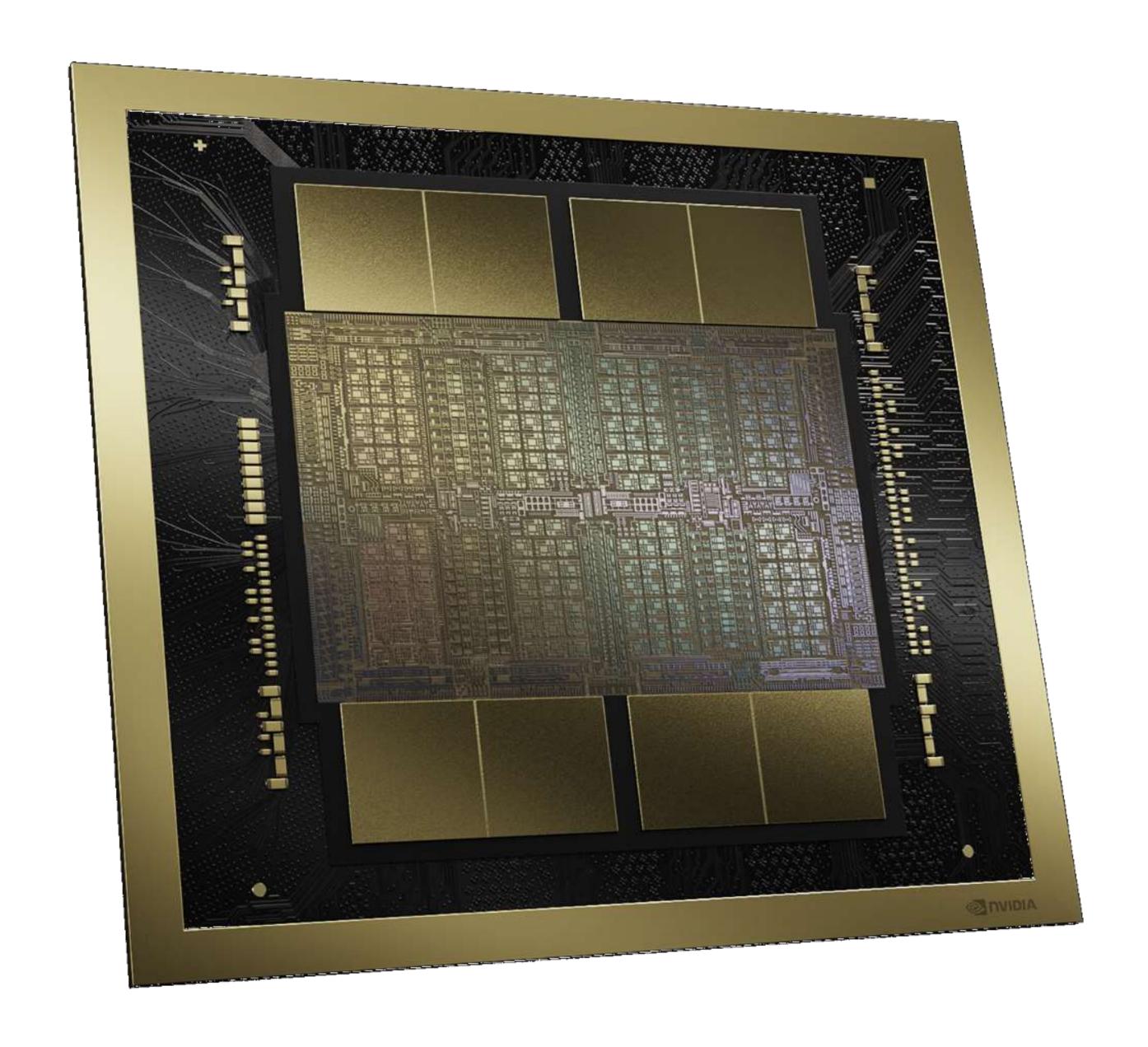






## **NVIDIA Blackwell**

The Engine of the New Industrial Revolution



Built to Democratize Trillion-Parameter Al

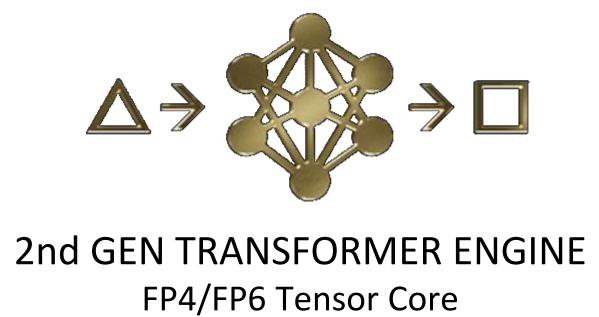
20 PetaFLOPS of Al performance on a single GPU

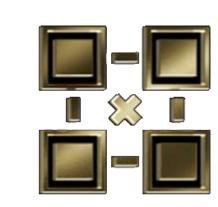
4X Training | 30X Inference | 25X Energy Efficiency & TCO

Expanding AI Datacenter Scale to beyond100K GPUs

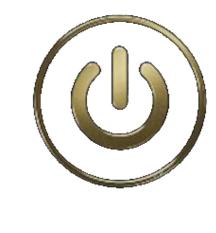


AI SUPERCHIP 208B Transistors



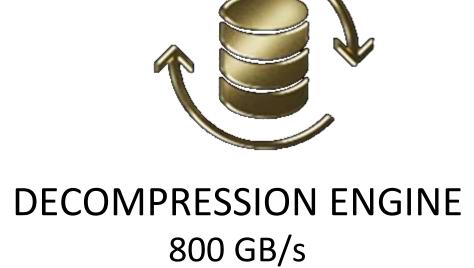


5<sup>th</sup> GENERATION NVLINK Scales to 576 GPUs



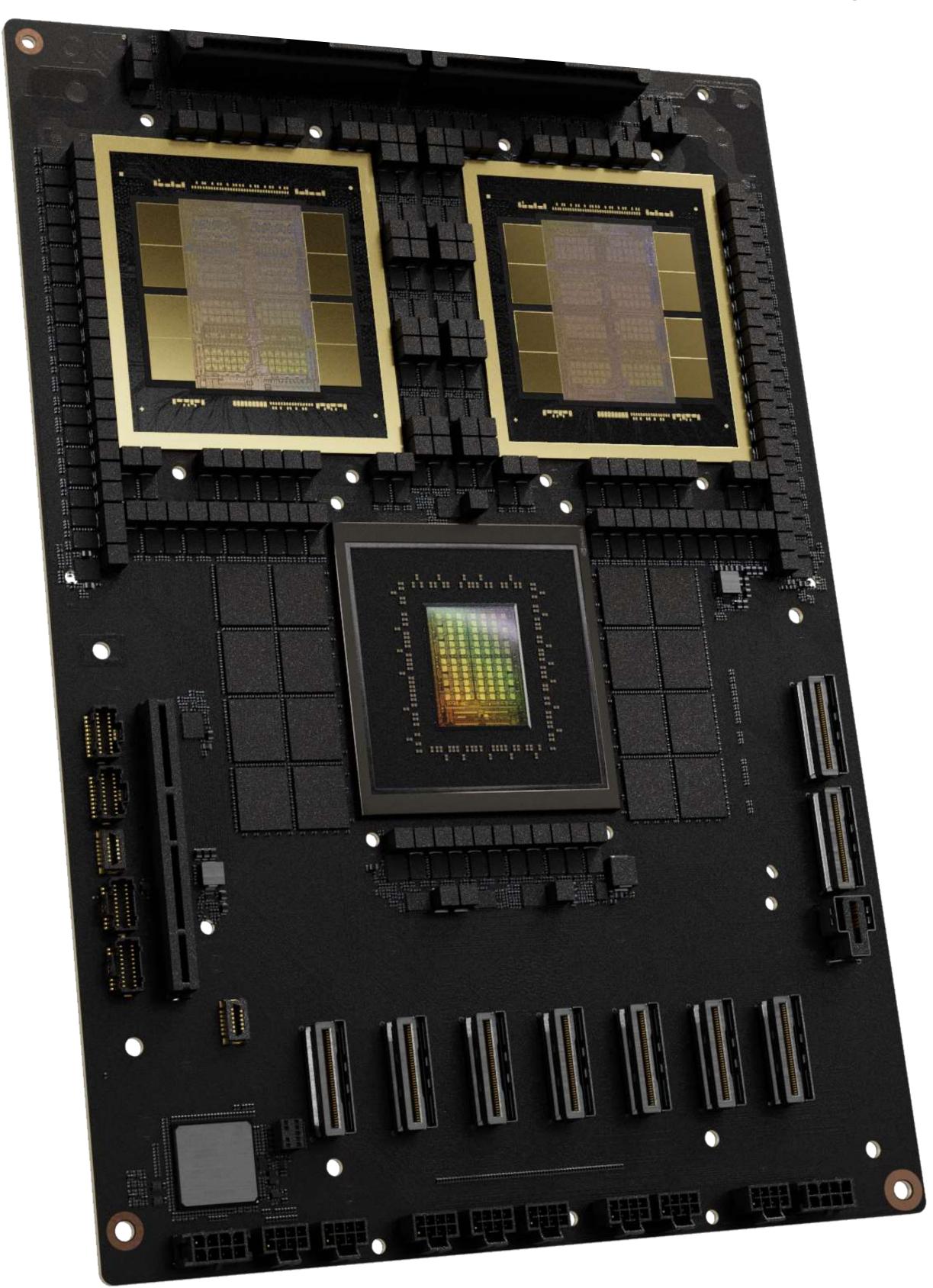
RAS ENGINE 100% In-System Self-Test





## GB200 SUPERCHIP

Optimized for Supercomputer-Scale Science



72 Grace CPU Arm cores

40 PetaFLOPS FP4 Al Inference

20 PetaFLOPS FP8 Al Training

16 TB/s of GPU memory bandwidth

864 GB Fast Memory

## DGX GB200

## Delivers New Unit of Compute



DGX GB200 - 36 GRACE CPUs

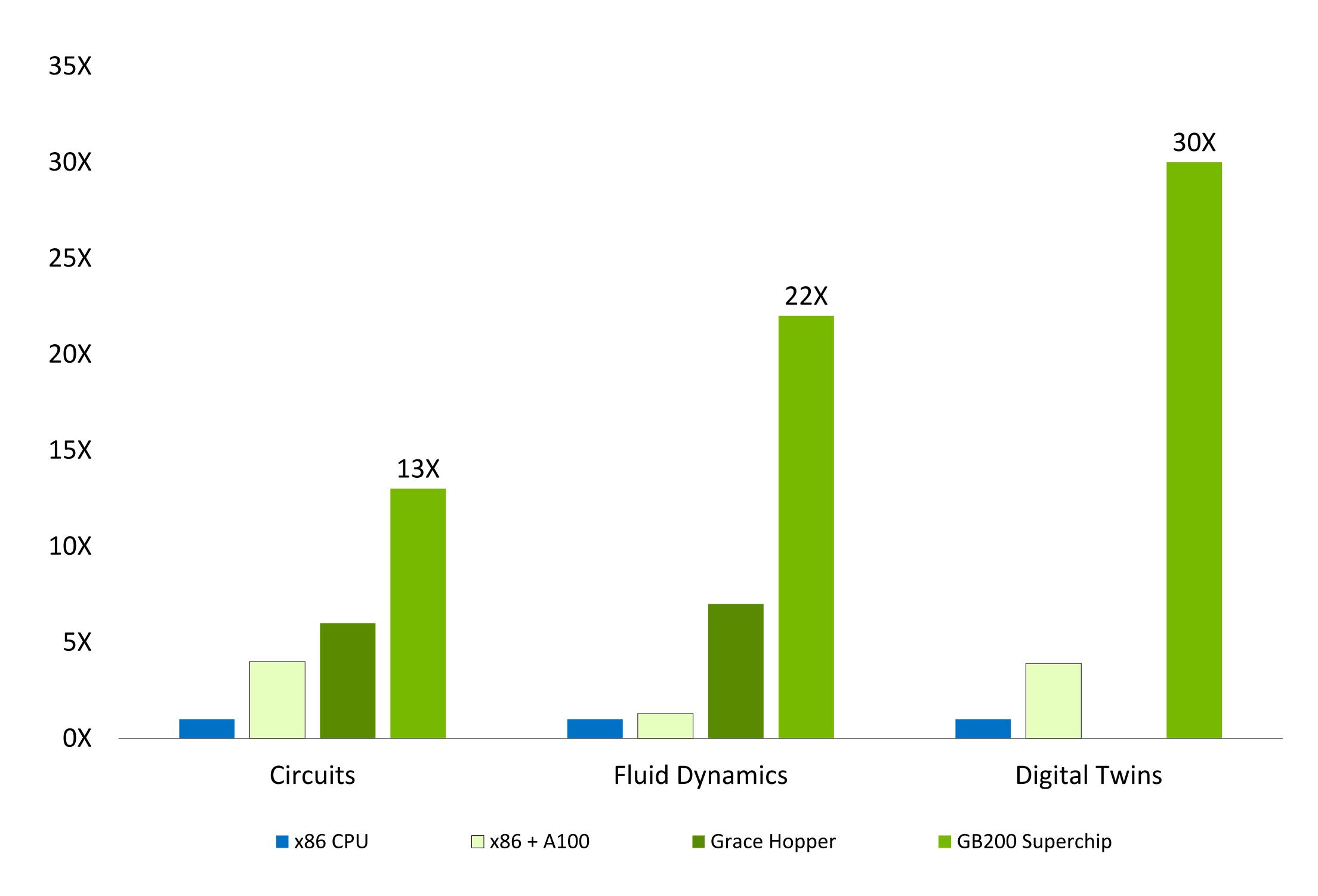
72 BLACKWELL GPUs

Fully Connected NVLink Switch Rack

Training
 Inference
 NVL Model Size
 Multi-Node All-to-All
 Multi-Node All-Reduce
 720 PFLOPs
 1,440 PFLOPs
 27T params
 130 TB/s
 260 TB/s

# Blackwell Pushes Boundaries of Engineering

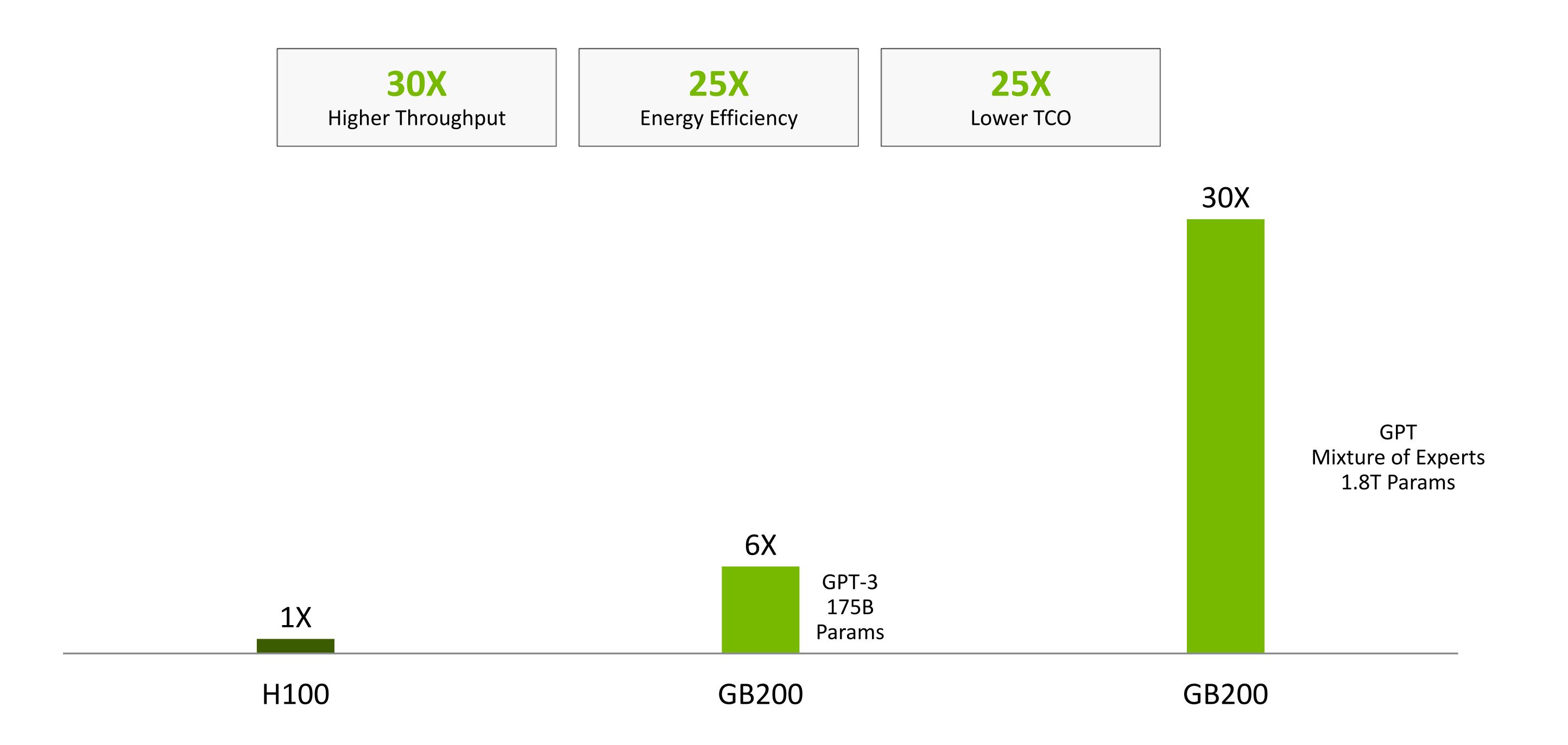
Supercharging Double Precision Performance for Product Design Simulations





# Blackwell — Driving the Era of Generative Al

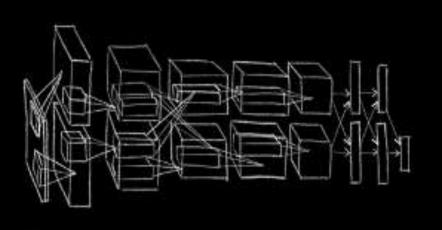
30X realtime inference, 25X improved energy efficiency

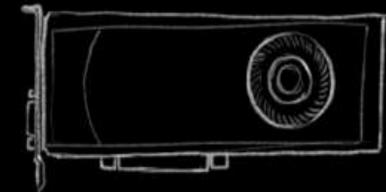






2012 ALEXNET "FIRST CONTACT"

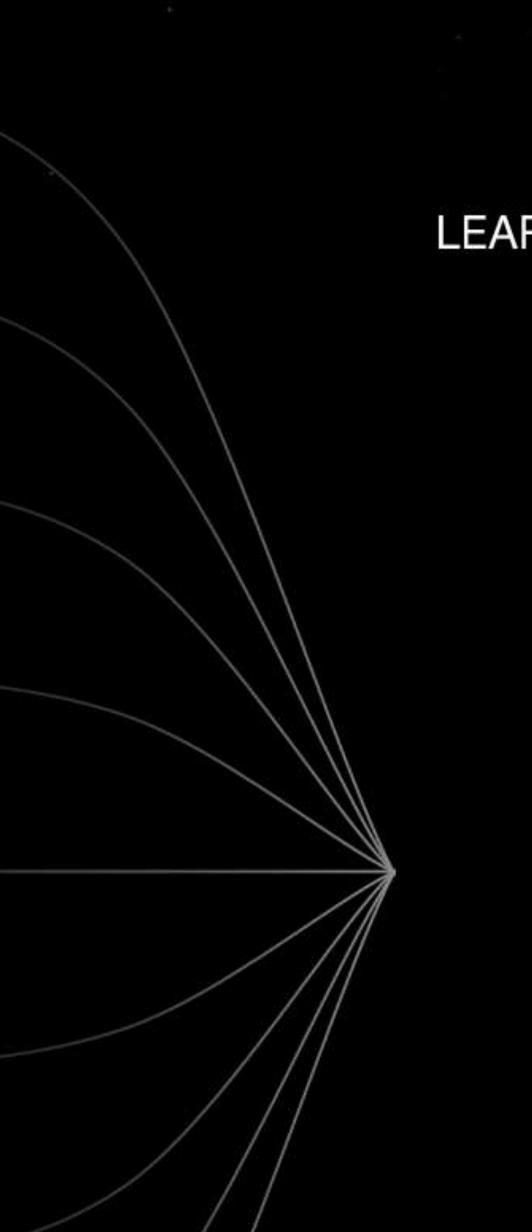




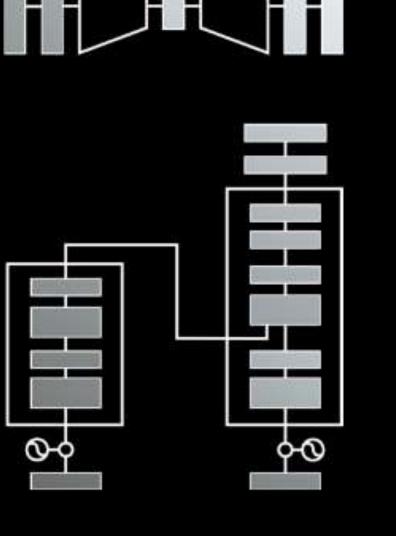
"CAT"

"An adorable cat in 3D confidently riding a flying, rocket-powered bike, adorned with a sleek black leather jacket." **TEXT** "A close shot of a cst in a futuristic space suit confidently operating controls in the cockpit of a sci-fi spaceship. The cockpit has lots of **TEXT IMAGE VIDEO TEXT** SPEECH **MULTI-MODAL AMINO ACID** 

**BRAINWAVES** 



## LEARN AND UNDERSTAND EVERYTHING





**IMAGE** 



**VIDEO** 



**IMAGE** 



3D



SOUND



ANIMATION



MANIPULATION



**PROTEIN** 



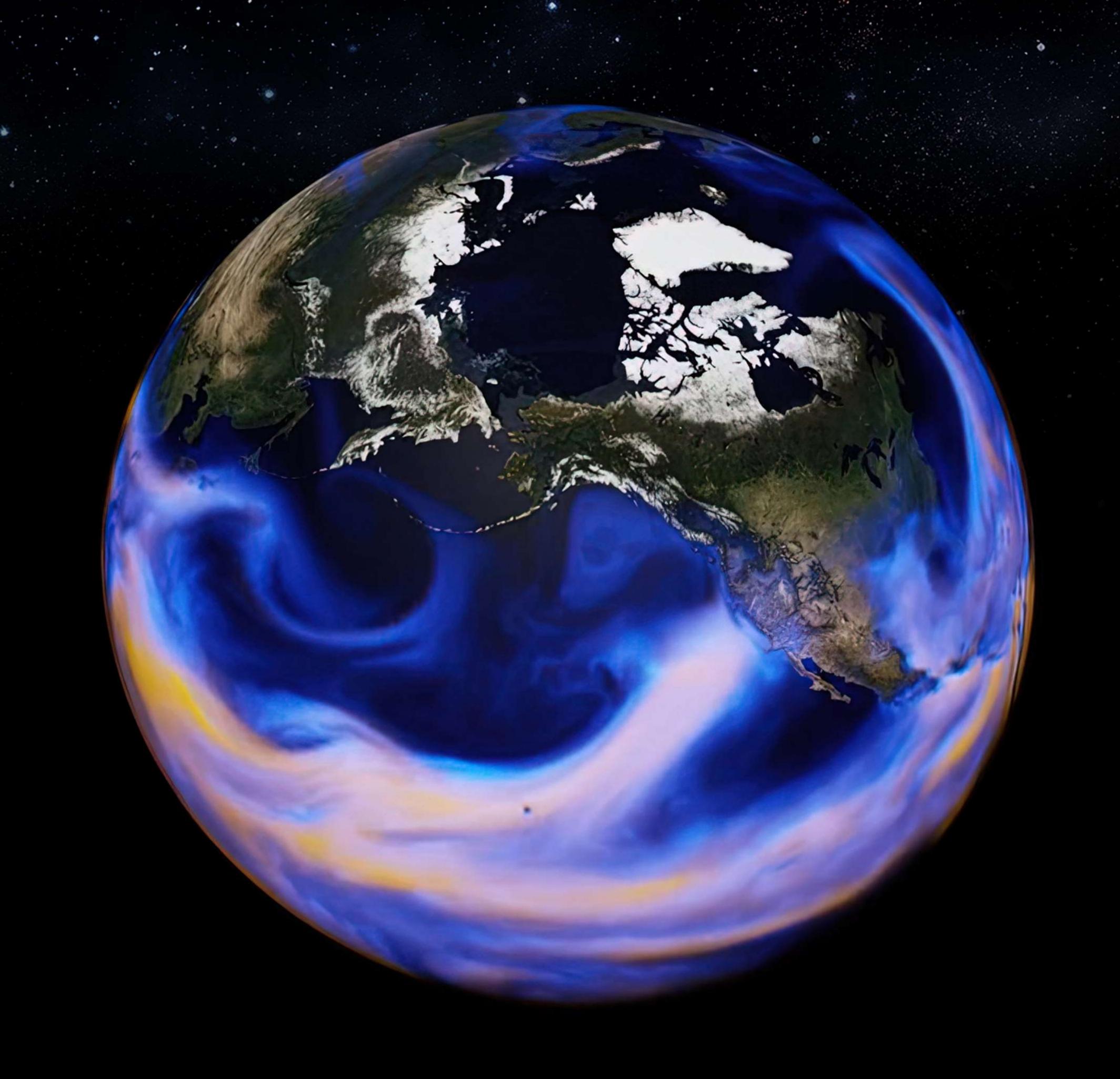
SPEECH

# Understand Planet Earth for building climate resilience

Digital Twins for Solving the world's biggest challenges, including climate change

Al for Disaster risk monitoring

Al for Modeling & Simulation



# Accelerating AI & HPC to Transform the Public Sector

Sovereign Foundation Models



Radar & Signal Processing



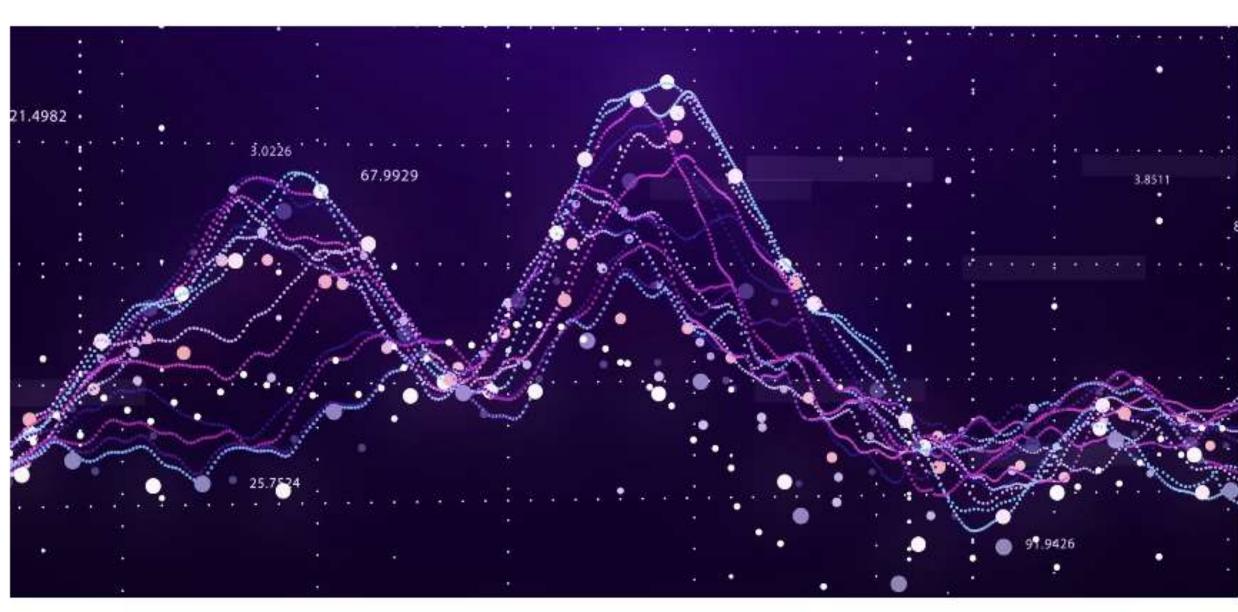
Geospatial Intelligence



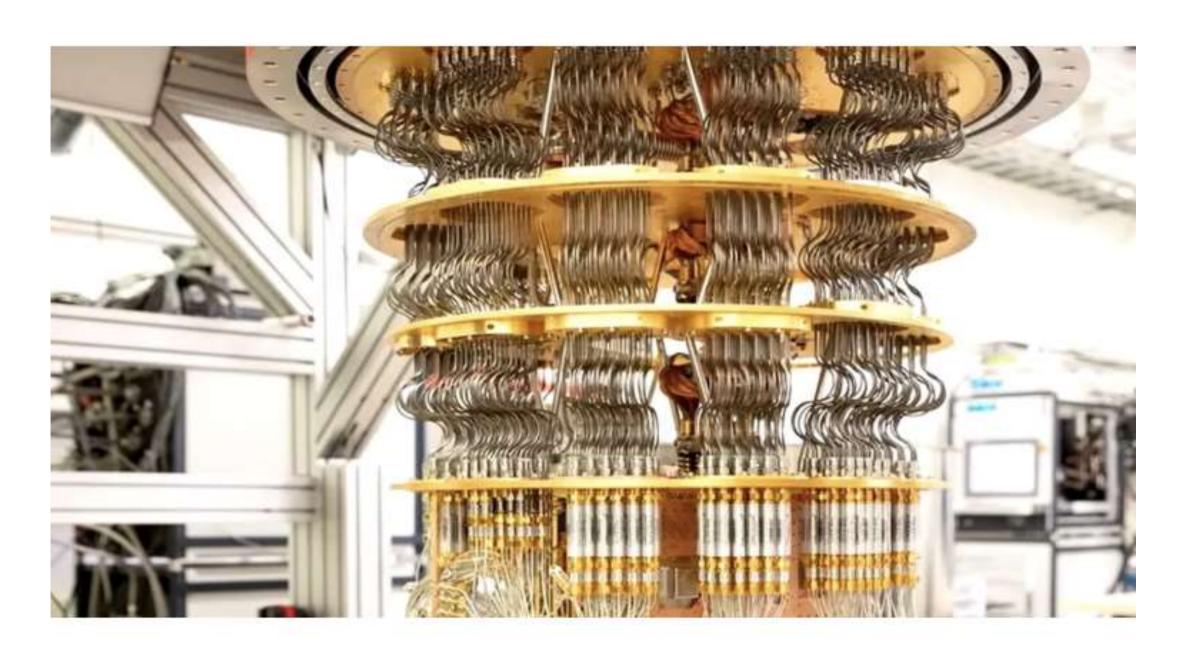
Digital Twins & Simulation



Data Analytics



Quantum Computing



# Generative Al Accelerates Policy Imperatives



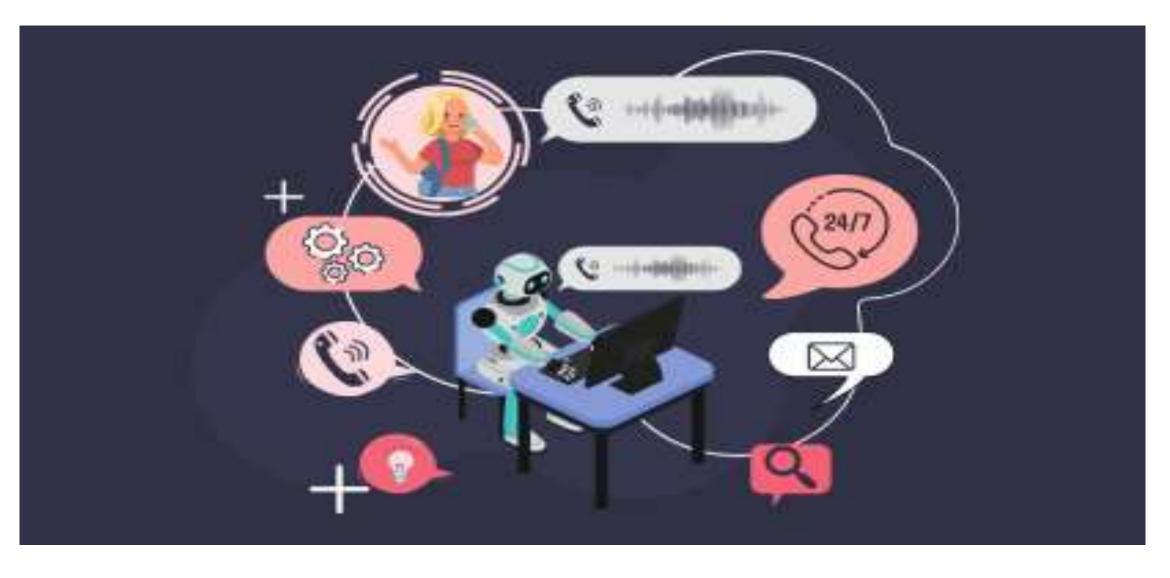
### **Protect National Security**

Anomaly Detection | Counter-Disinformation
Training and Simulation | Autonomous Systems and Robotics



Education and Skills Development

Adaptive personalized learning materials | Student Analytics Virtual Teaching Assistants | Education Policy Analysis



Personalized Citizen Services

Virtual Assistance for Social Services | Sentiment Analysis Language Translation Services | Citizen GPT



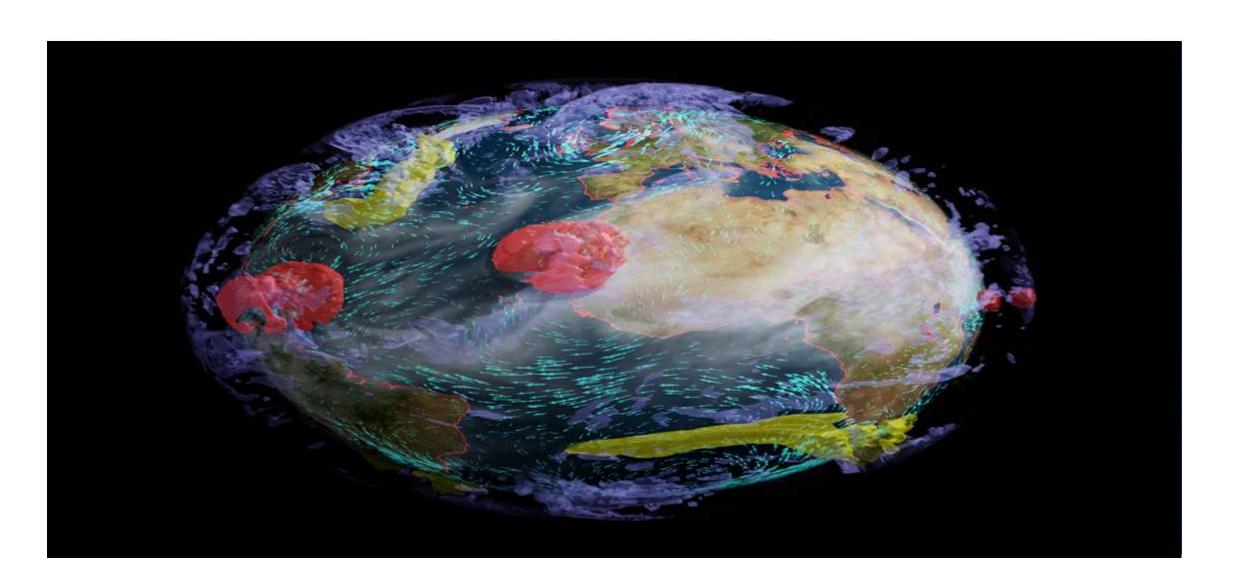
**Urban Planning and Smart Cities** 

Urban Design and Simulation | Public Space Optimization



### Regulatory Compliance and Enforcement

Automated Compliance Monitoring (i.e.., safety, environmental, Transportation, and vehicle emissions)

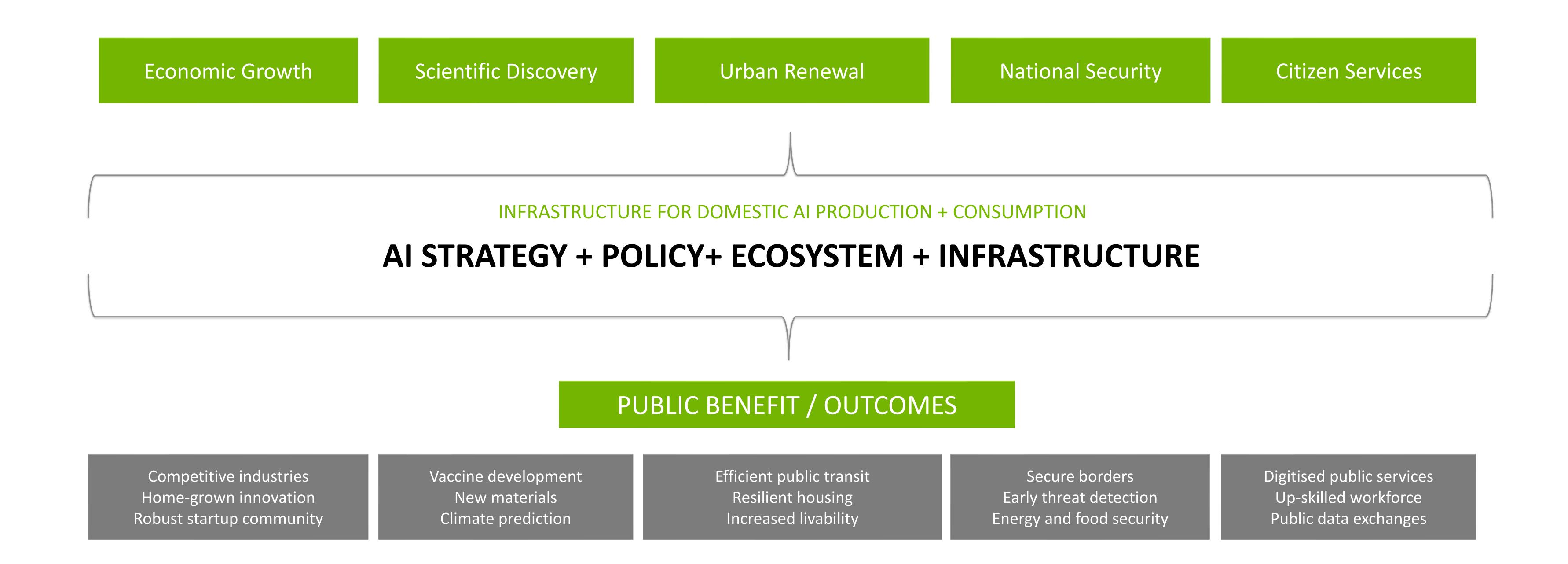


Climate Action and Policy

Environment GPT – Scenario Generation | Climate Policy Assessment and Adaptation Planning



# Building an Al Nation Digital Transformation in the Age of Al



# AI FACTORY IS THE NEW CRITICAL INFRASTRUCTURE

### INSTRUMENT FOR SCIENTIFIC DISCOVERY

- Human condition
- Fundamental research
- Industrial innovation

### **ENGINE FOR ECONOMIC GROWTH**

- Start-up ecosystem
- New jobs/sectors
- Retain talent; Workforce productivity

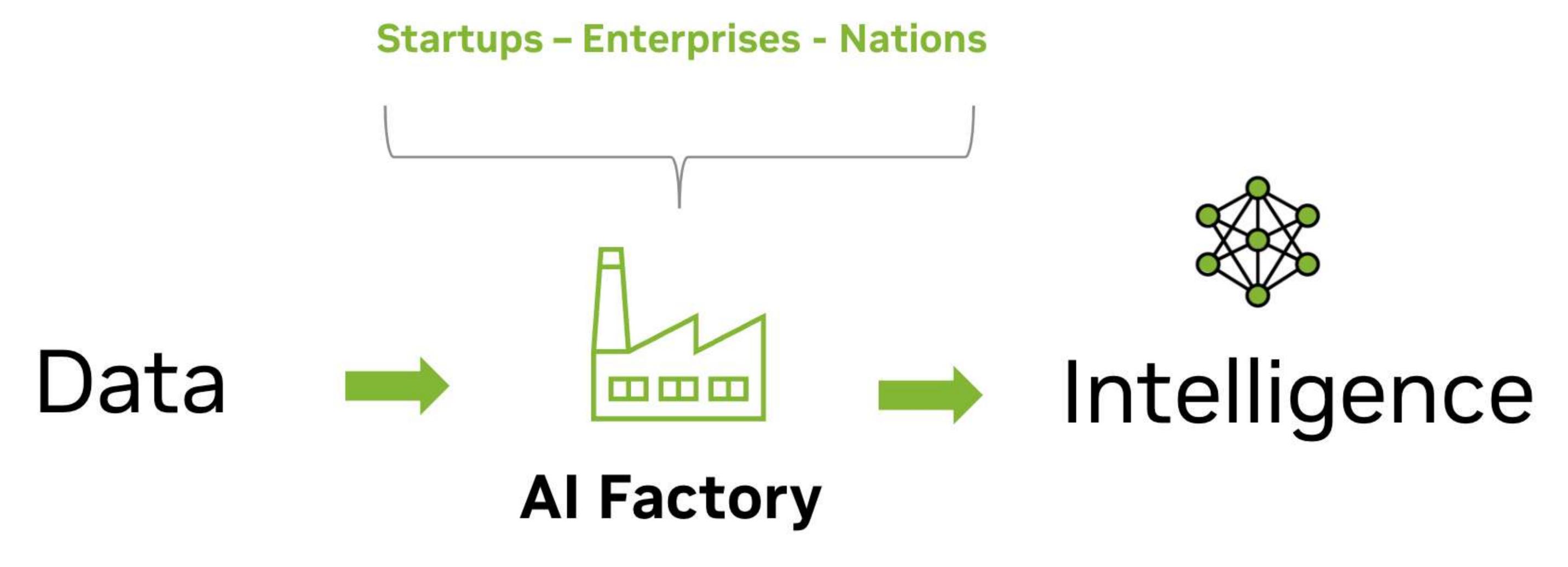
## PLATFORM FOR PUBLIC SECTOR INNOVATION

- Improve access to services
- Expand citizen services
- Defend national interests



# Al Factories

Manufacturing "intelligence"



Accelerated compute + network + storage designed for Al

Access to Al platform, tools and expertise

# Sovereign Al

## New Al-based Supercomputers to Foster Scientific Innovation



Isambard-AI | Isambard-3
United Kingdom
5,280 GH200 GPUs | Grace CPUs with over 55,000 Arm Neoverse V2
Cores



Gefion

Denmark

1,528 H100 GPUs | NVIDIA Quantum-2 InfiniBand | NVIDIA CUDA Quantum



Jean Zay
France

1,456 NVIDIA H100 GPUs | Liquid Cooled |
400Gb/s NVIDIA Quantum-2 InfiniBand

# Plano Brasileiro de Inteligência Artificial (PBIA)

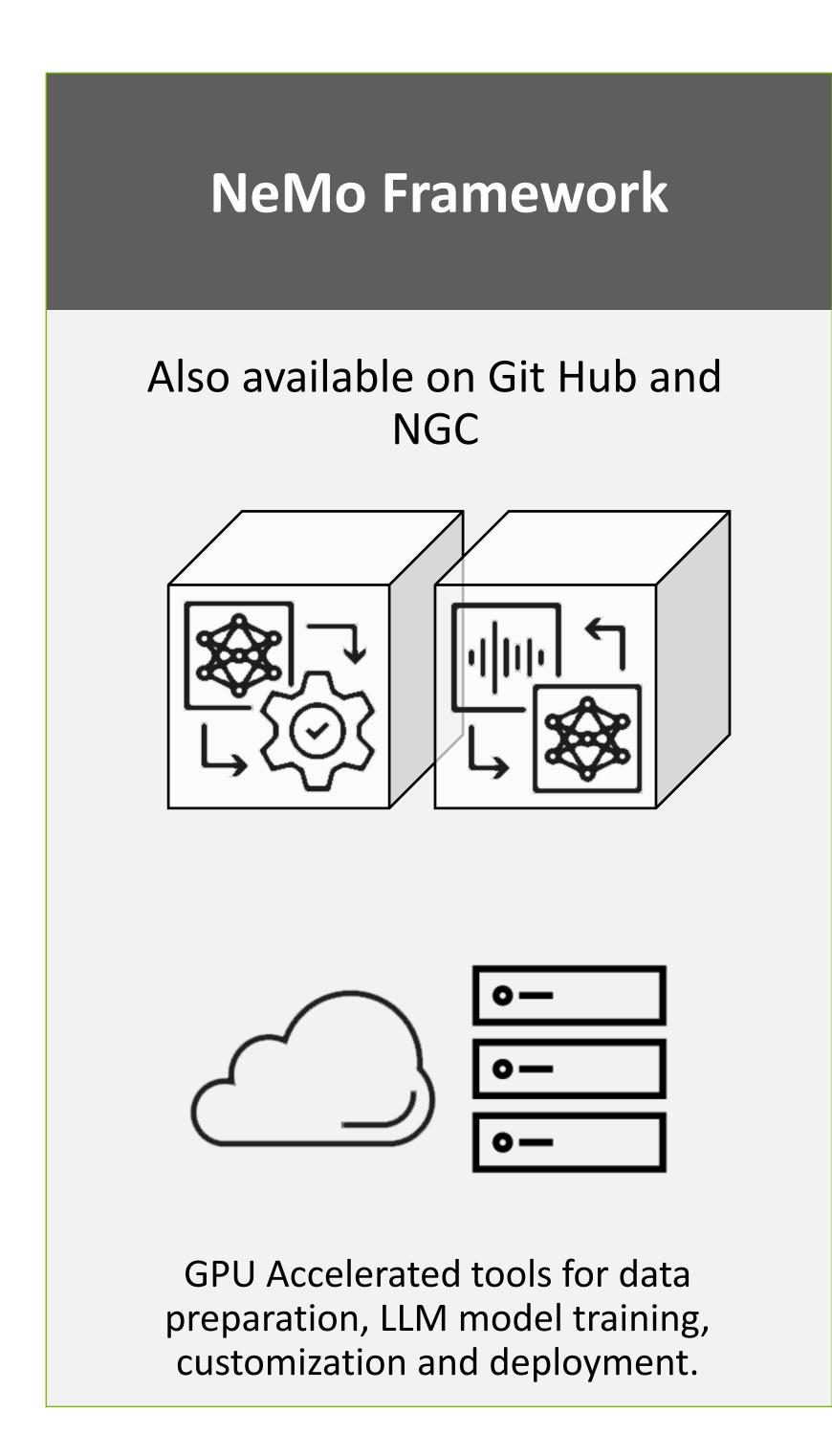
IA para o Bem de Todos

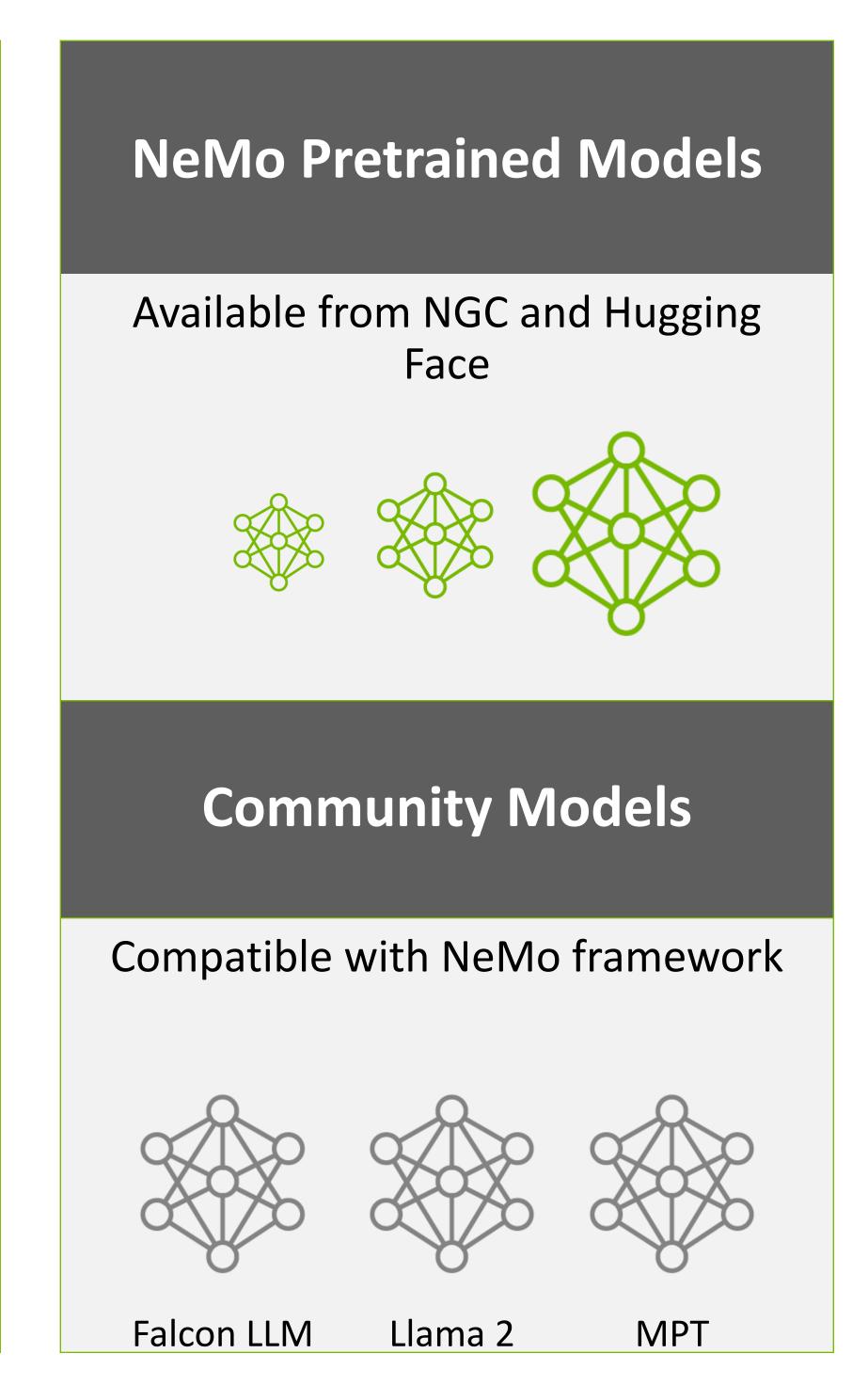


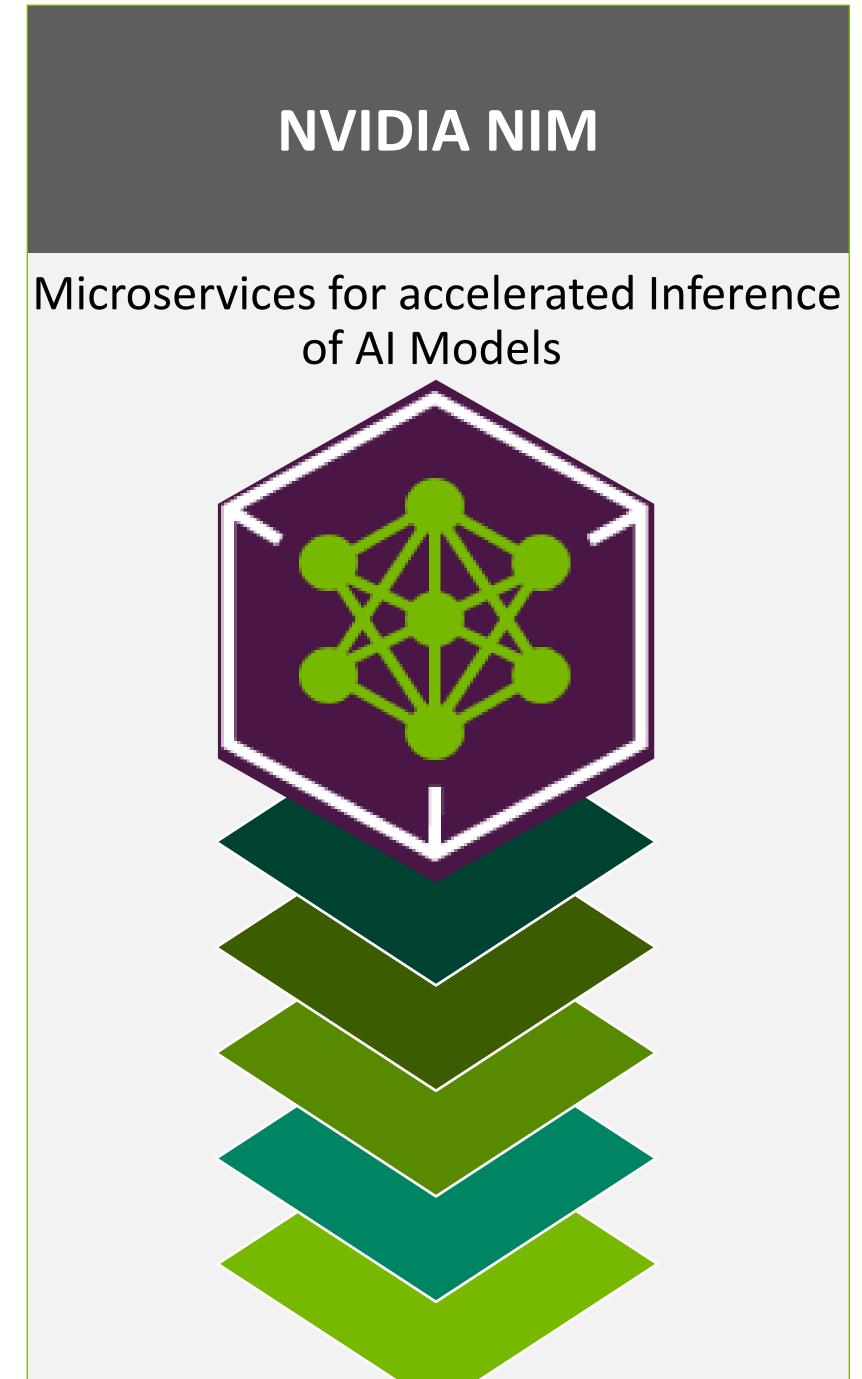


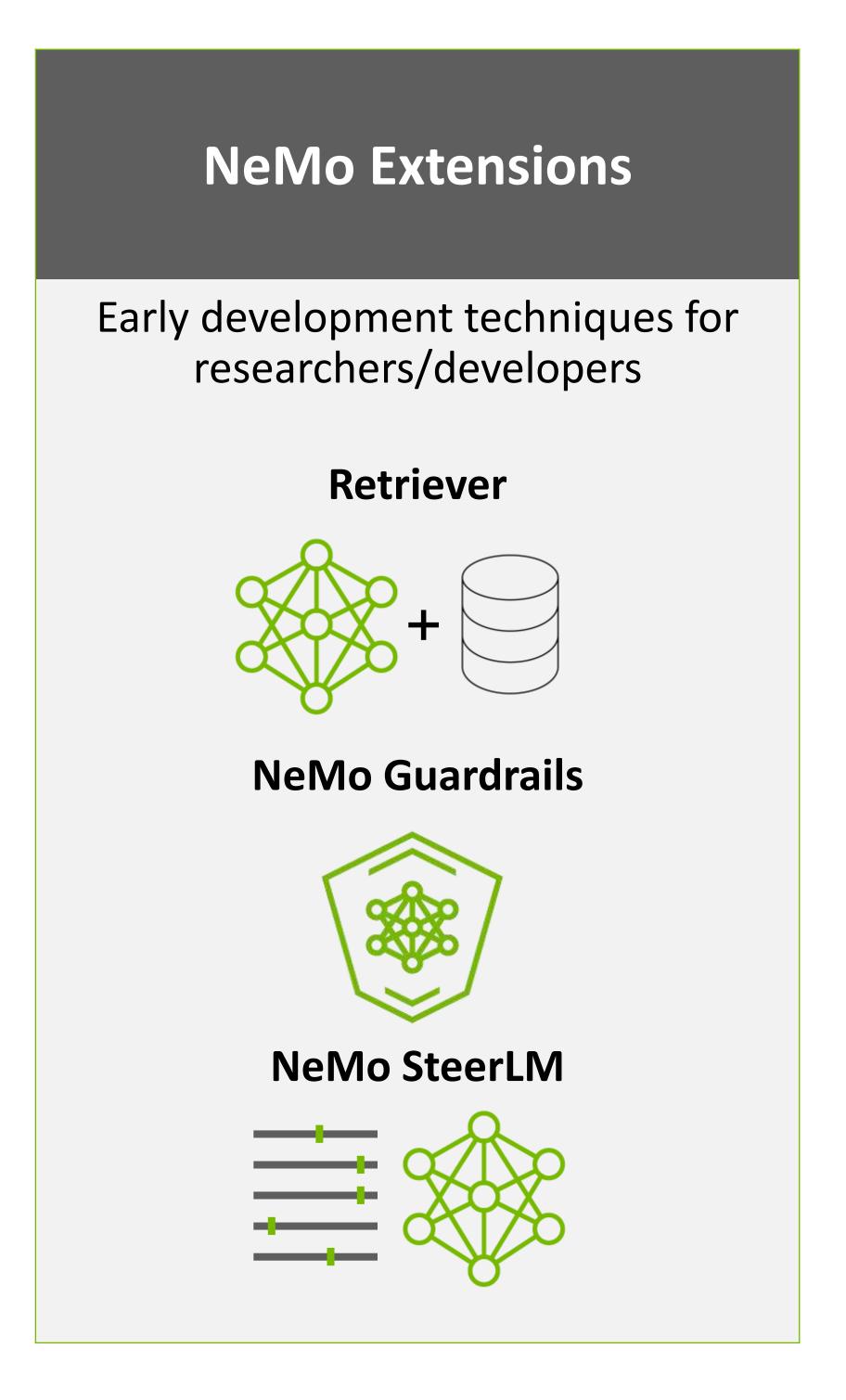
## **NVIDIA LLM Portfolio**

Multiple Products to Help Customers Build Custom LLMs





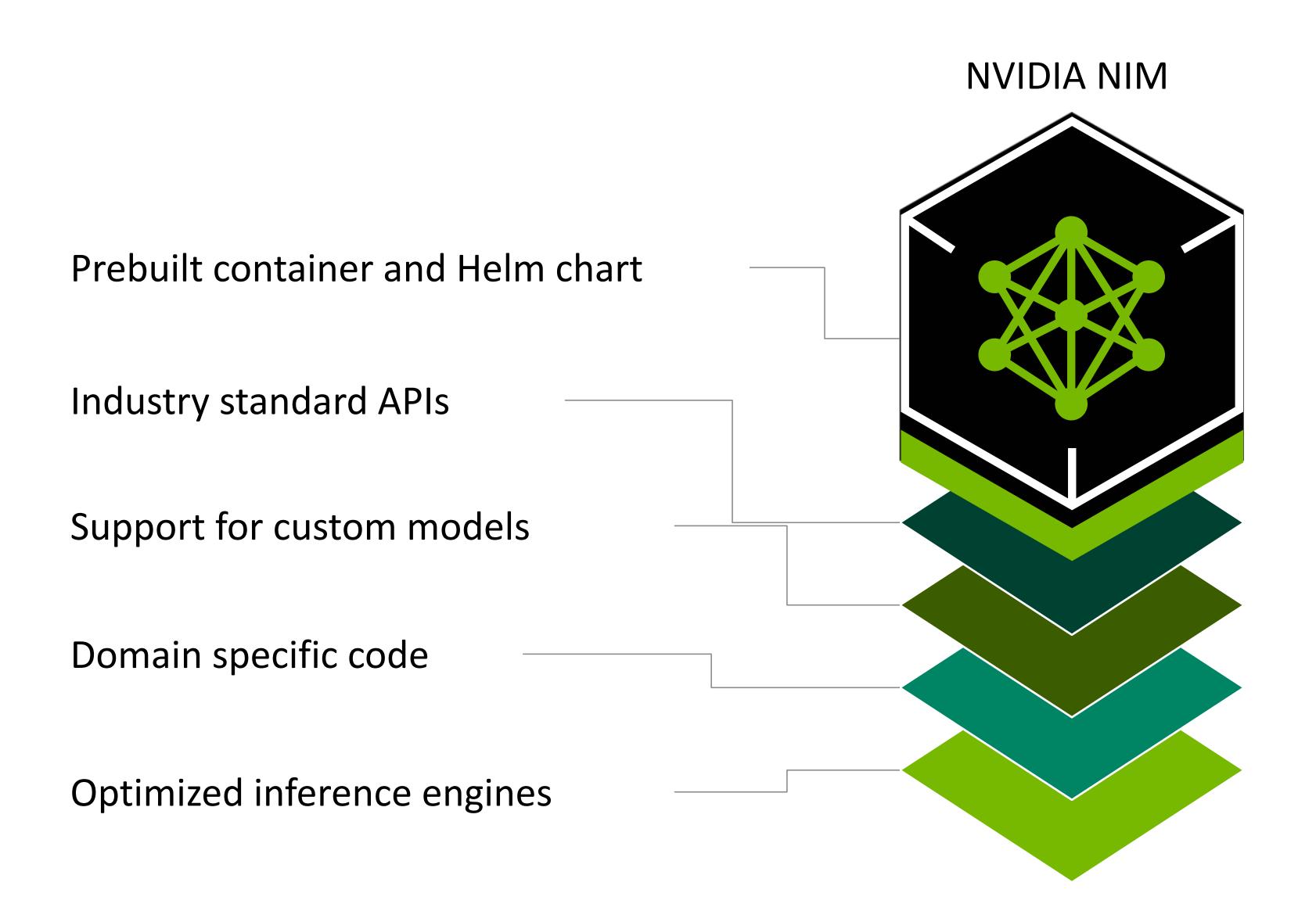




https://www.nvidia.com/en-us/ai/

## **NVIDIA NIM: Inference Microservices for Generative Al**

Accelerated runtime for generative Al



Deploy anywhere and maintain control of generative Al applications and data

Simplified development of AI application that can run in enterprise environments

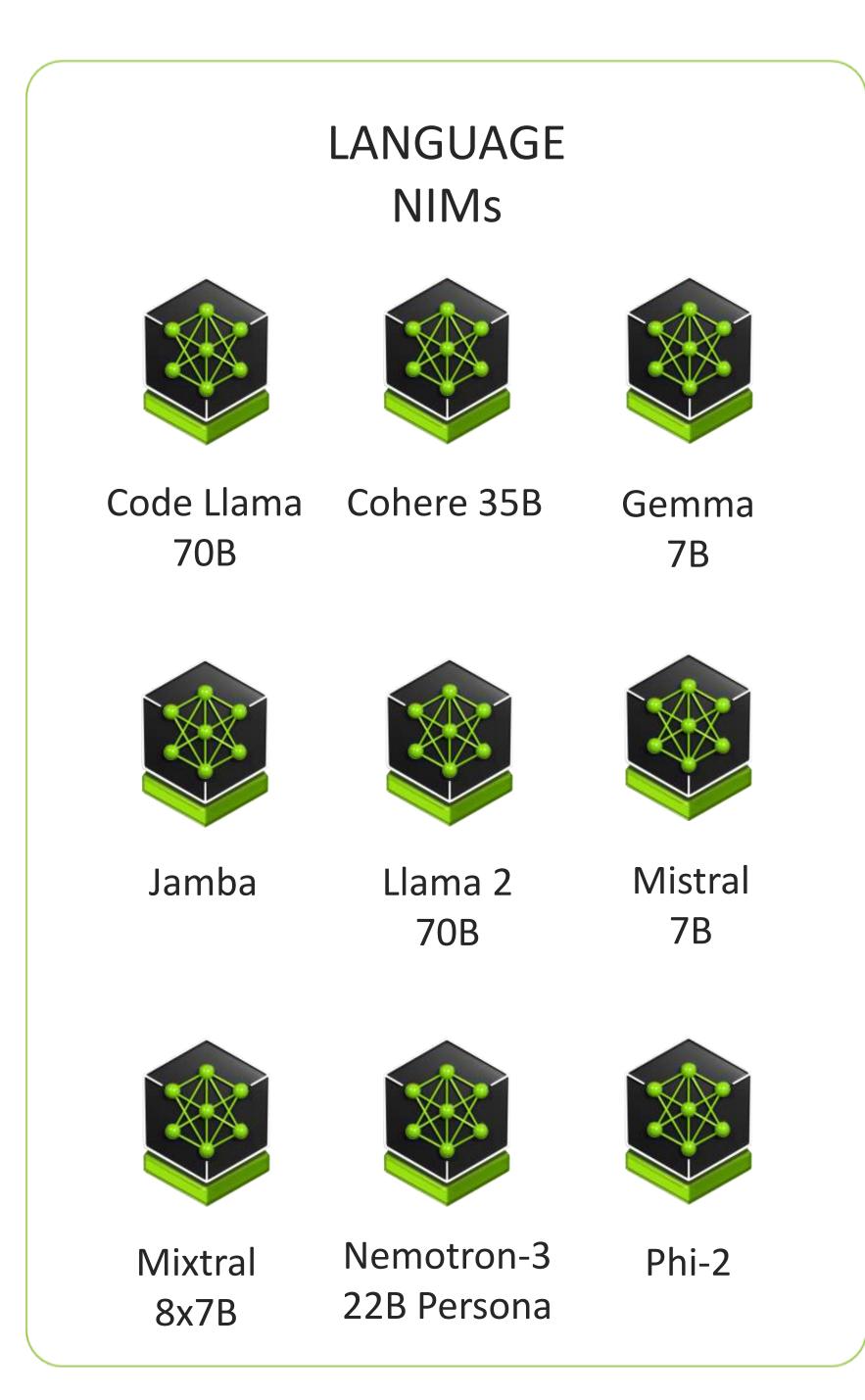
Day 0 support for all generative AI models providing choice across the ecosystem

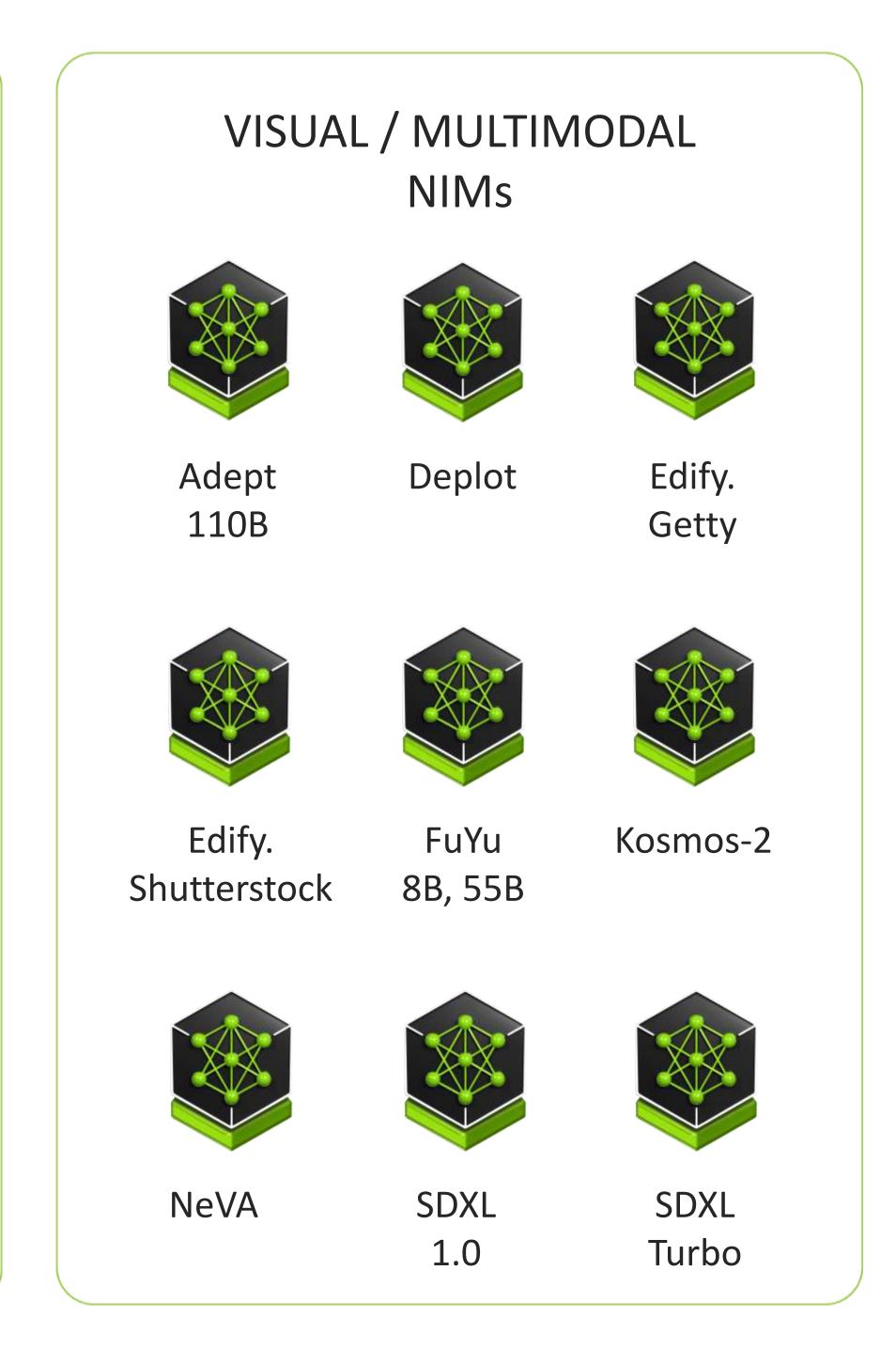
Improved TCO with best latency and throughput running on accelerated infrastructure

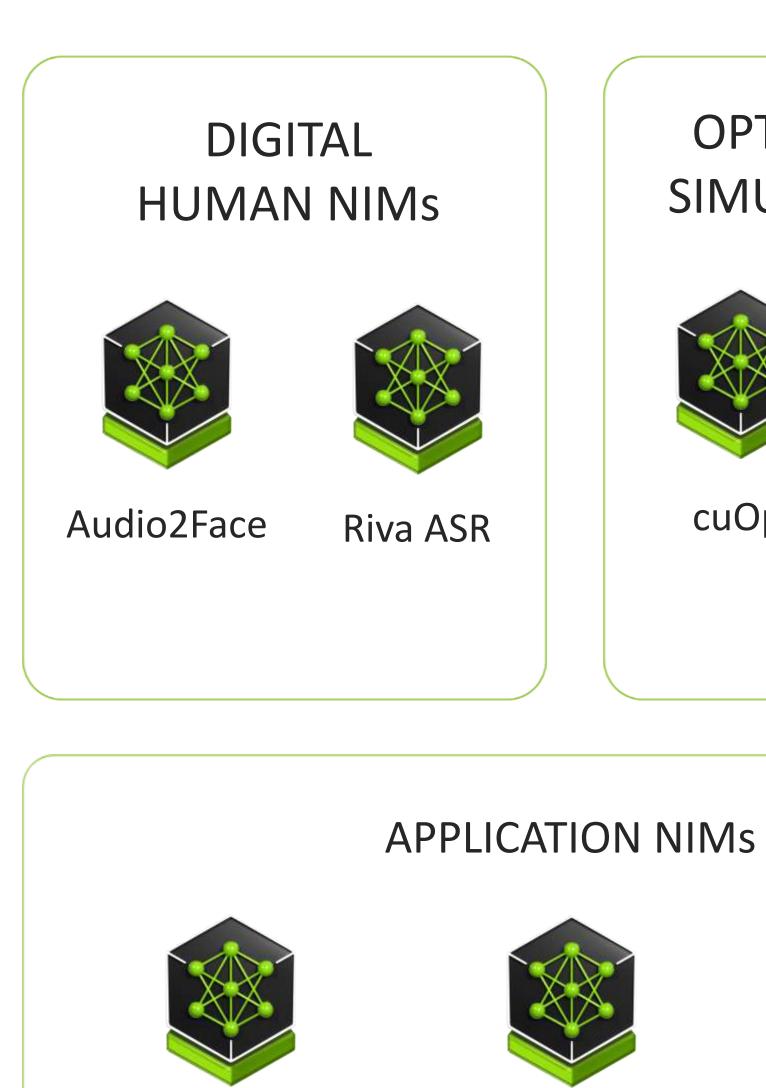
Best accuracy for enterprise by enabling tuning with proprietary data sources

Enterprise software with feature branches, validation and support

# **NVIDIA NIM for Every Domain**





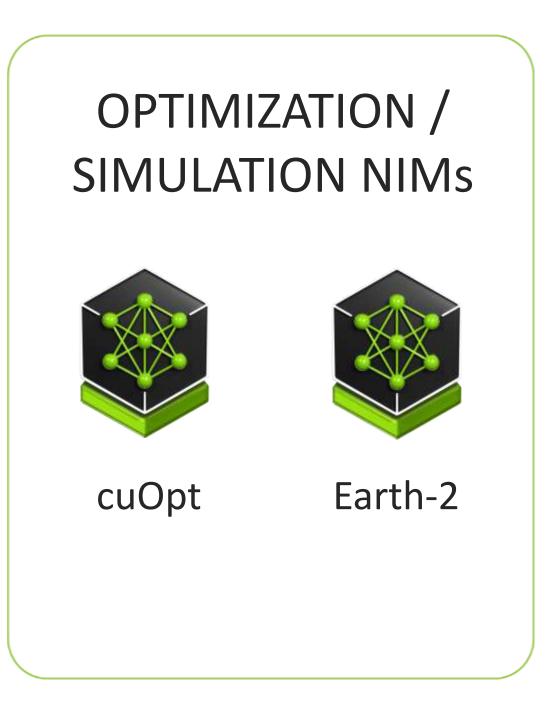


Llama

Guard

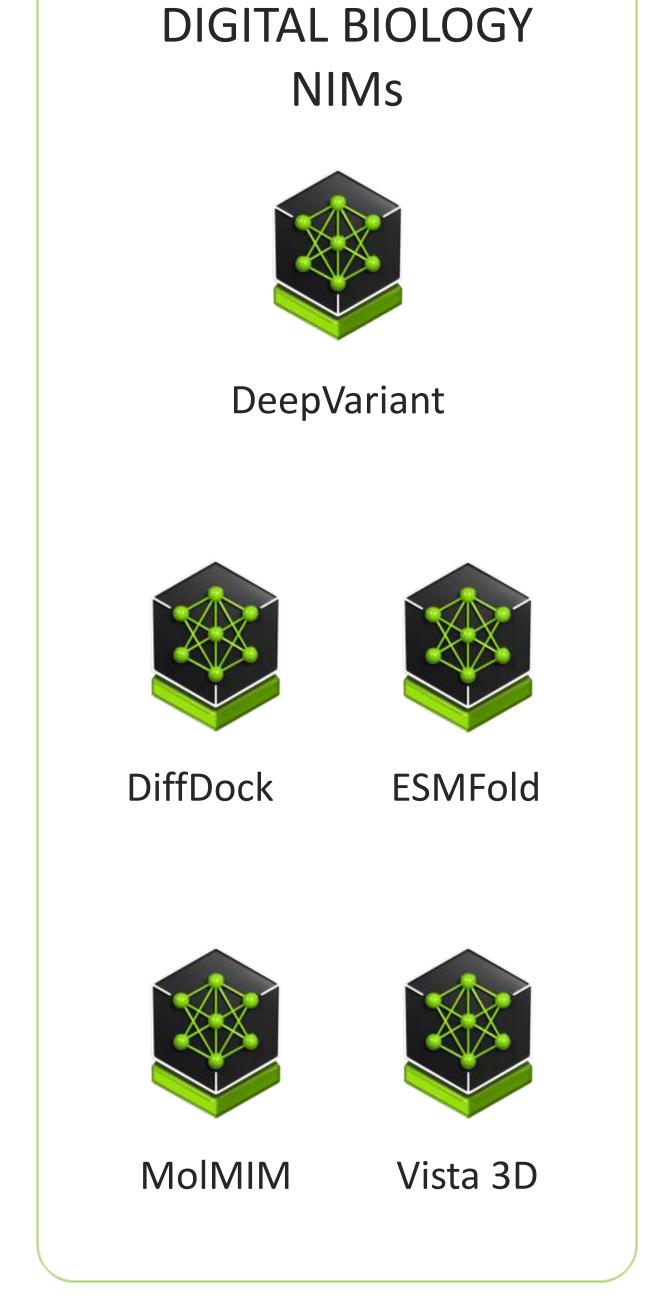
Retrieval

Embedding



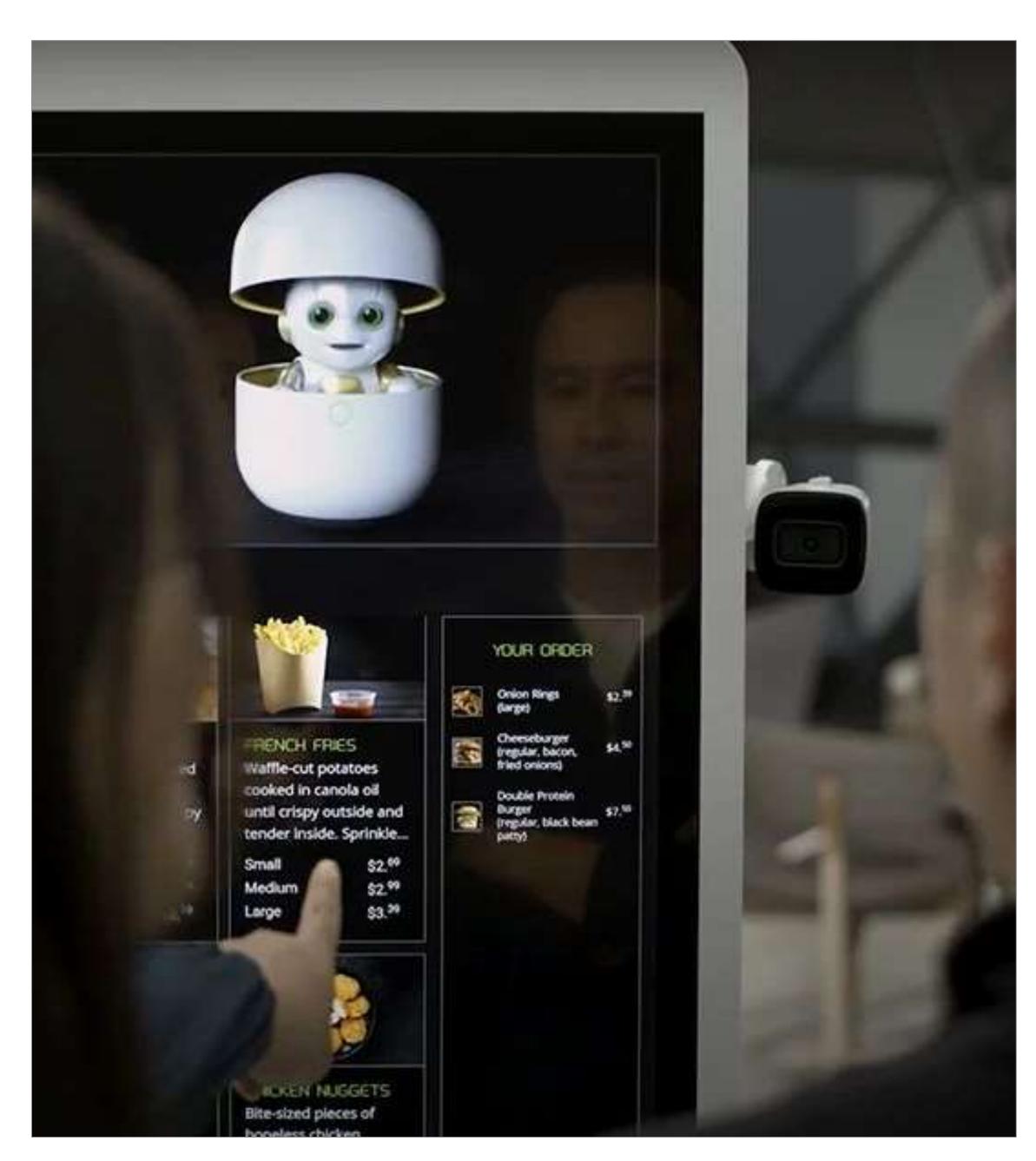
Retrieval

Reranking



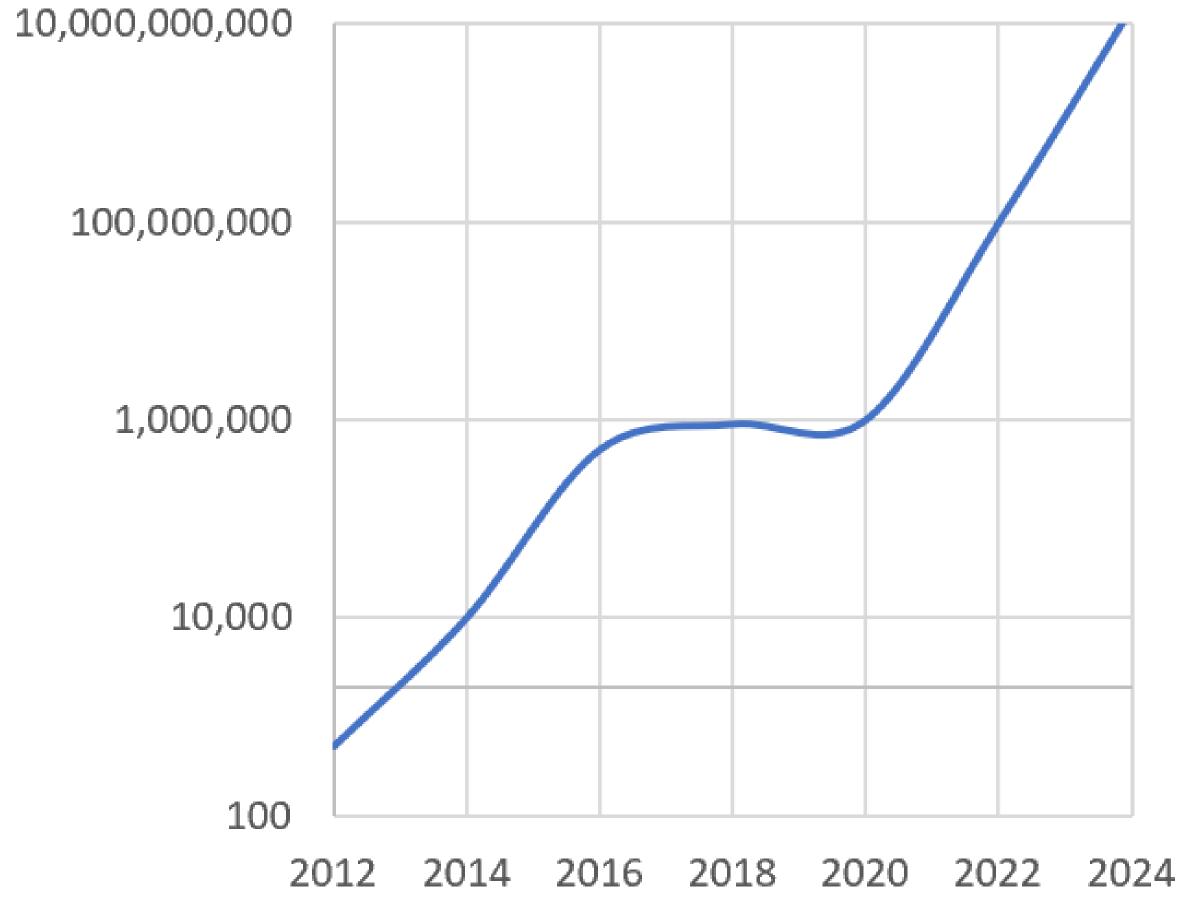


# **Exploding Energy Usage In Data Centers**

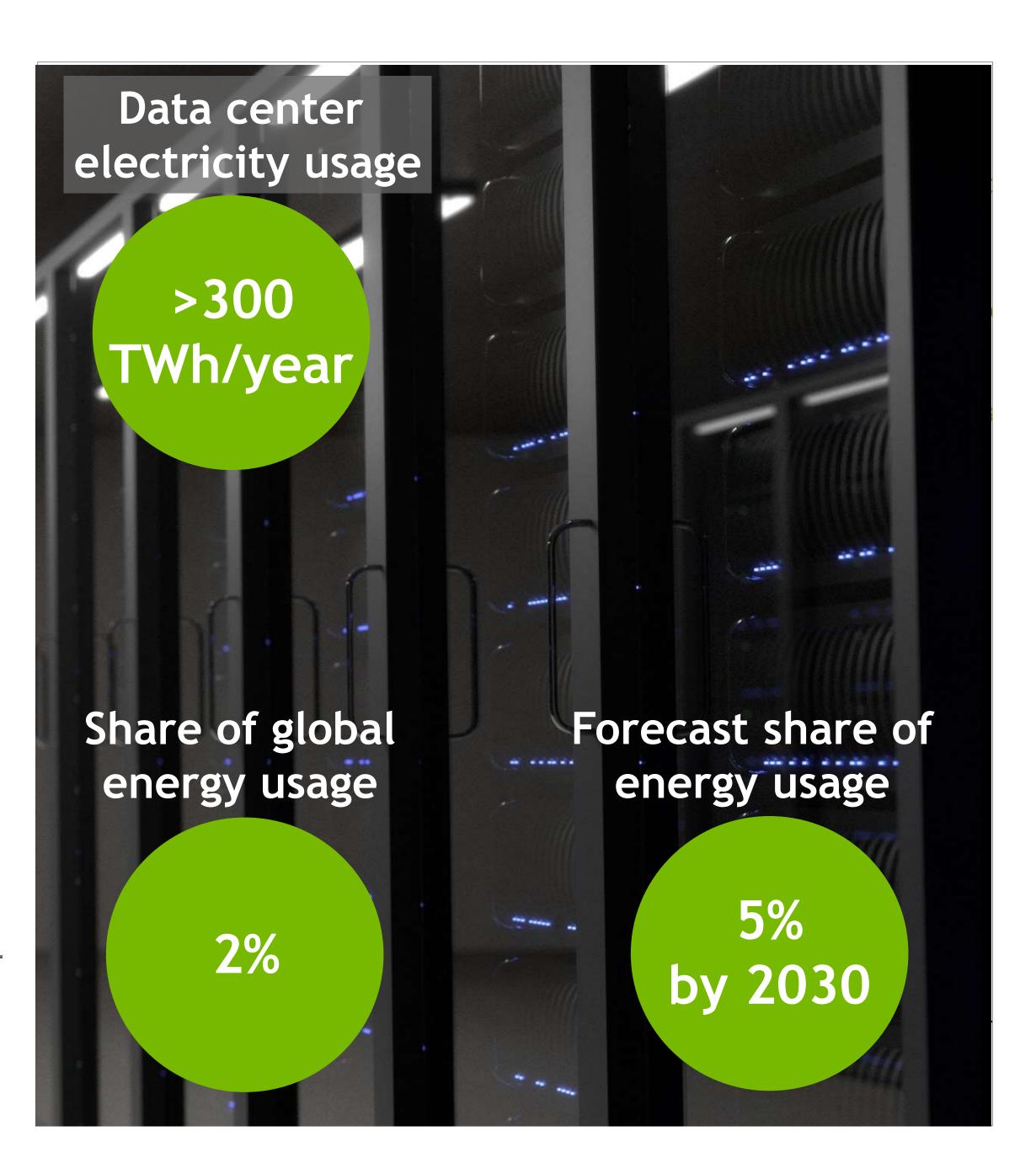


Massive Al Models Deliver New Use Cases
Transformer networks power Conversational Al and Metaverse

#### Transformer Models PetaFLOPS to Train



Model Sizes Demanding More Compute



Data Center Need to Become More Efficient



# **Accelerated Computing Is Sustainable Computing**

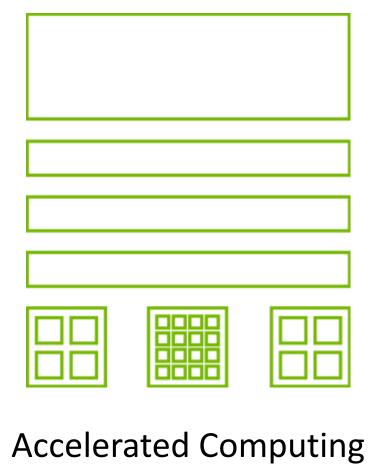
Full Stack, 3 Chips, Data Center Scale

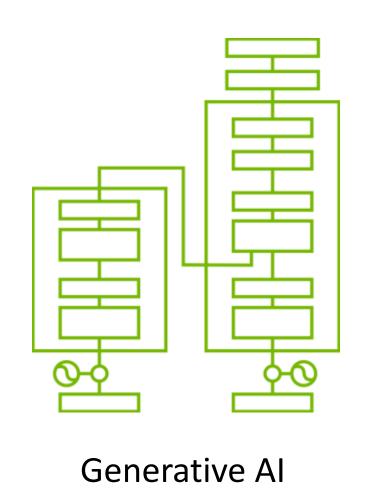
48 Million CUDA Downloads

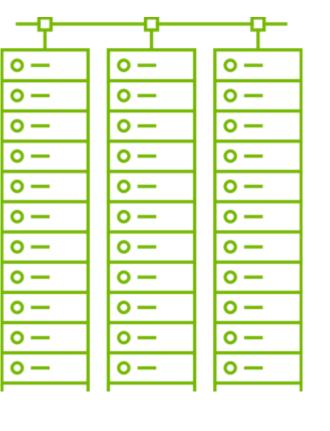
4.5 Million Developers

3,000 Applications

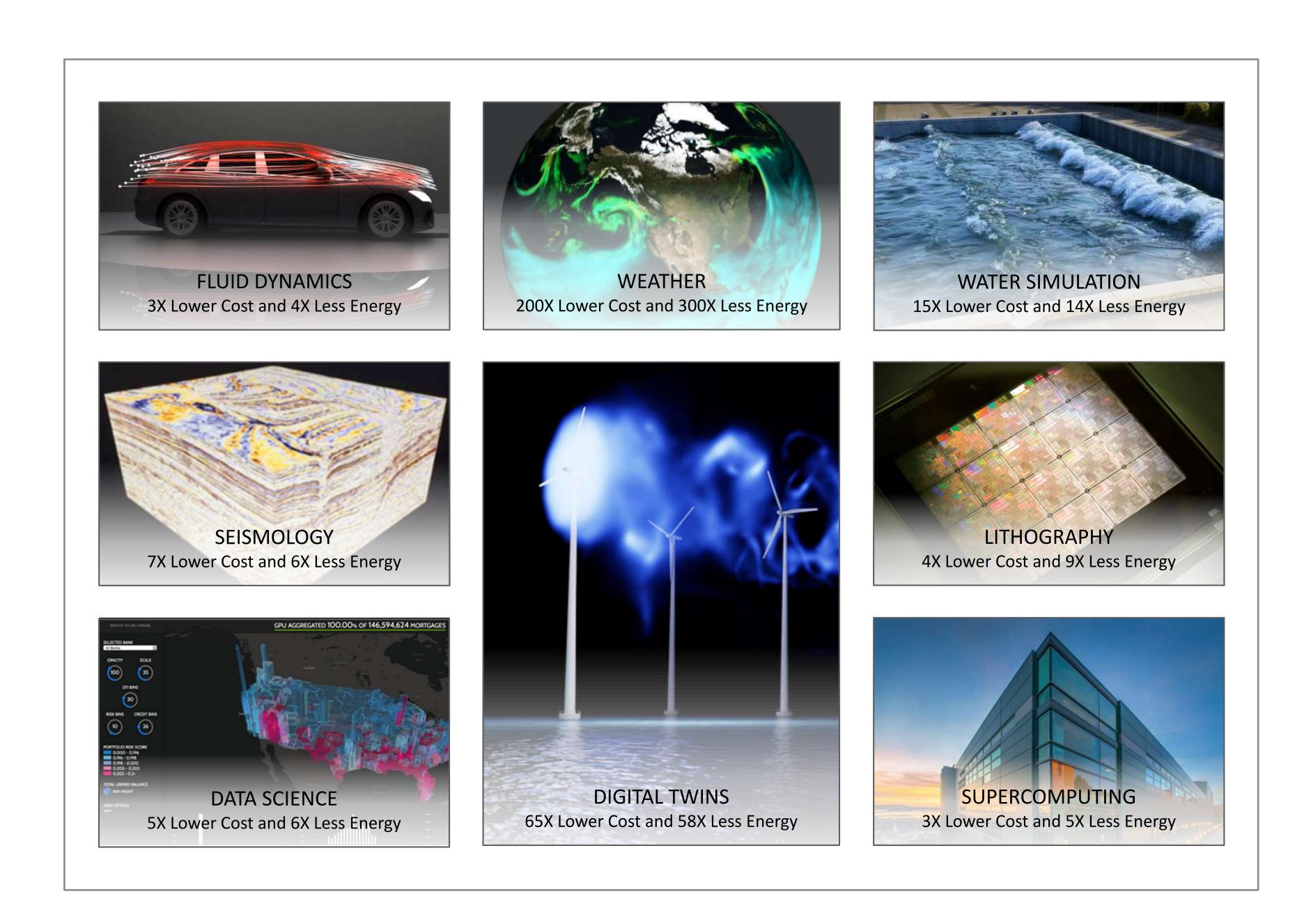
600 SDKs & Al Models







Data Center Scale

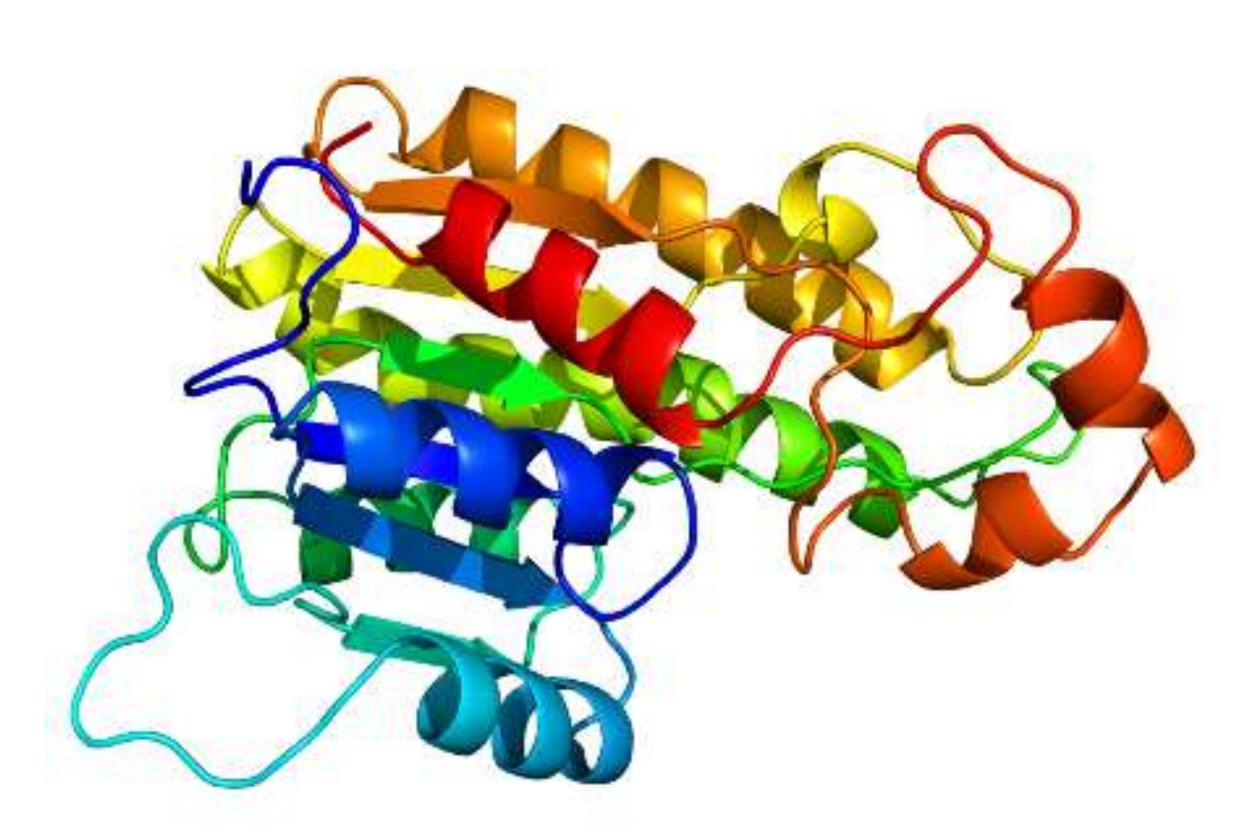


# HPC Powered AI Applications Accelerated by NVIDIA Platform

23X Improvement in Energy Efficiency of LLMs Adapted for Scientific Use Cases

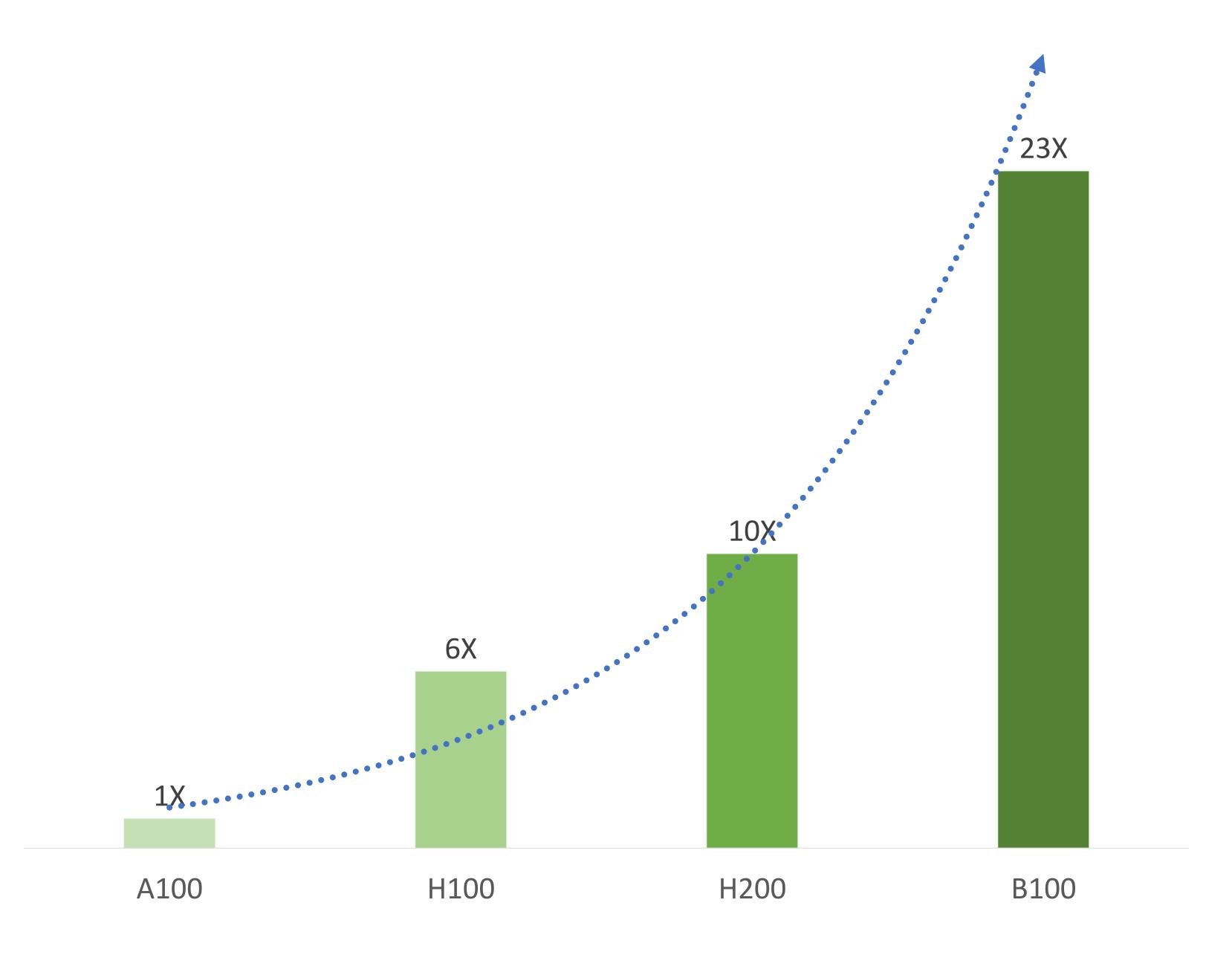
#### **Drug Discovery and Genomics**

Transformer based networks like AlphaFold and MegaMolBART enable protein sequencing for drug discovery



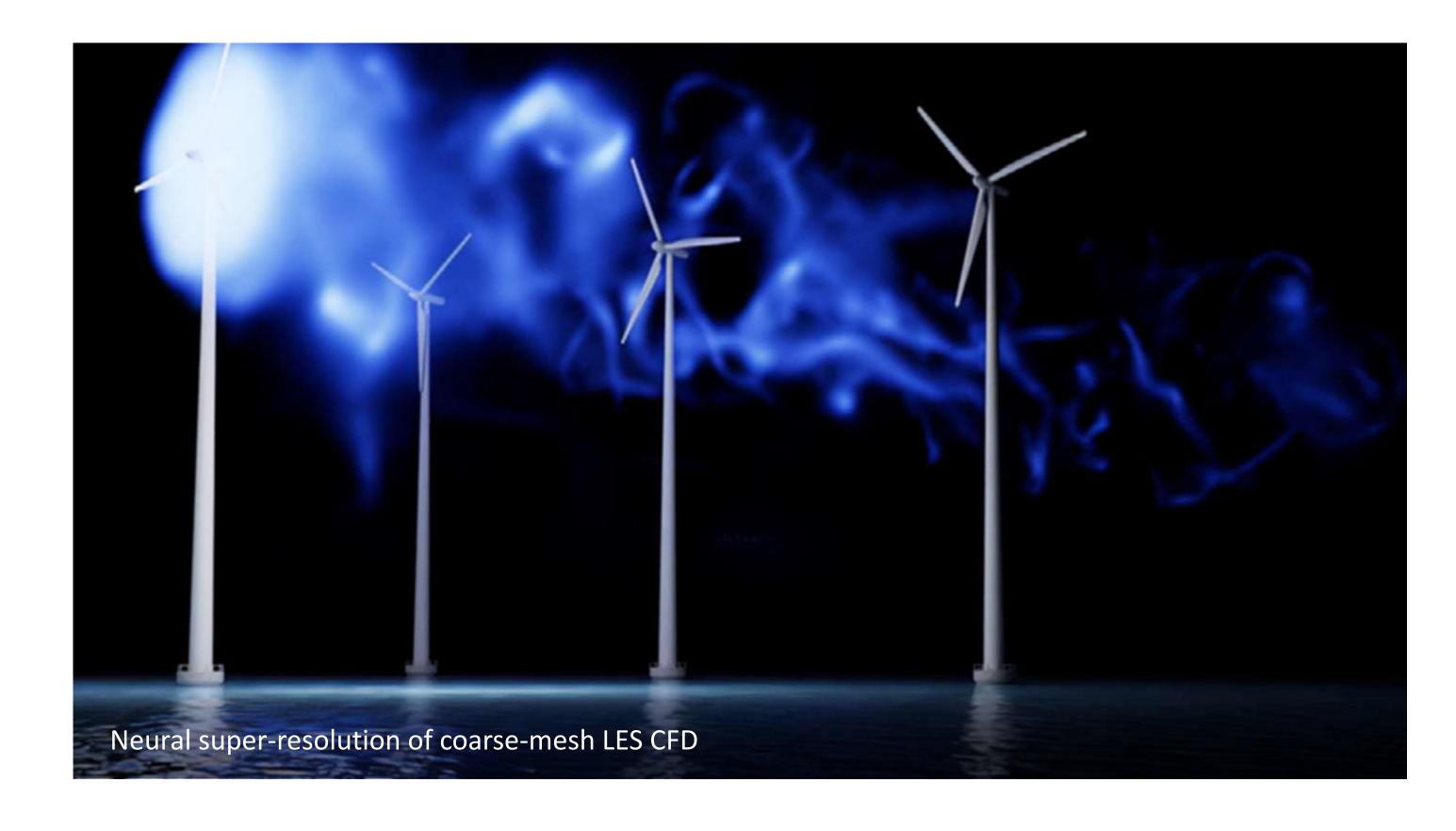
Al outperforms lab discovery techniques at Cost, speed, permutations

GPT-3 175B LLM Inference Energy Efficiency Performance
Higher is Better



#### Digital Twins With Physics-ML

Al-accelerated, physics-based, high-fidelity surrogates for interactive digital twins



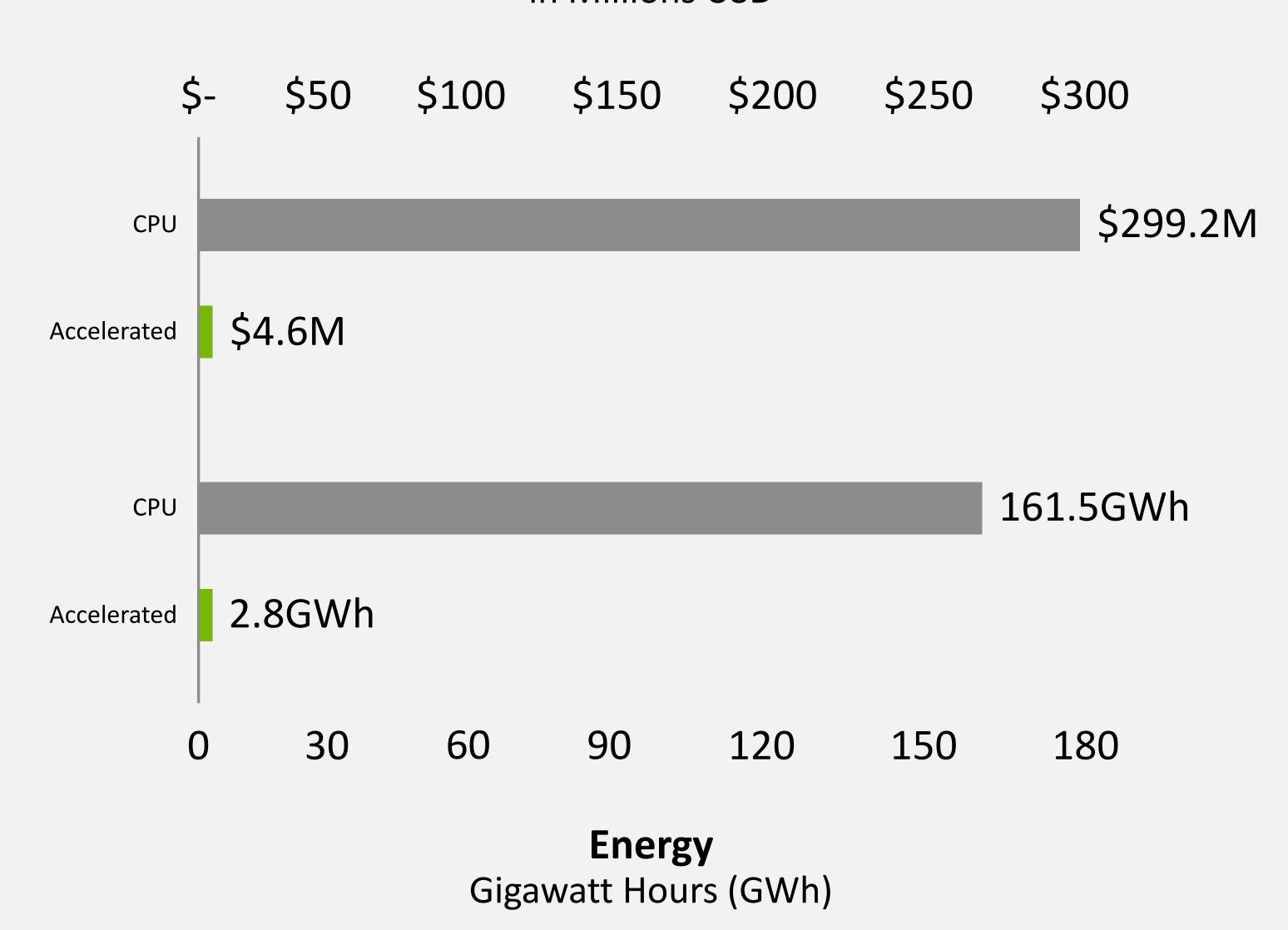
<sup>&</sup>quot;By leveraging NVIDIA Modulus, a GPU-based physics-ML framework, we want to accelerate complex, multi-physics simulations with AI-powered surrogates. The reduced computational time enables us to develop energy-efficient digital twins for a sustainable, reliable, and affordable energy ecosystem."



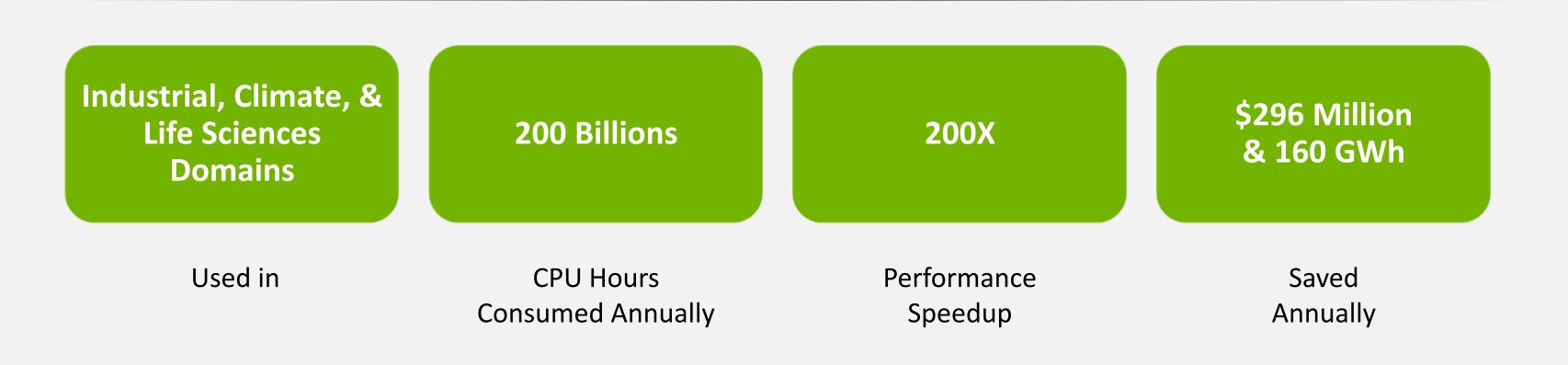
#### 65X Lower Cost and 58X Less Energy

#### **Acquisition Cost**

in Millions USD



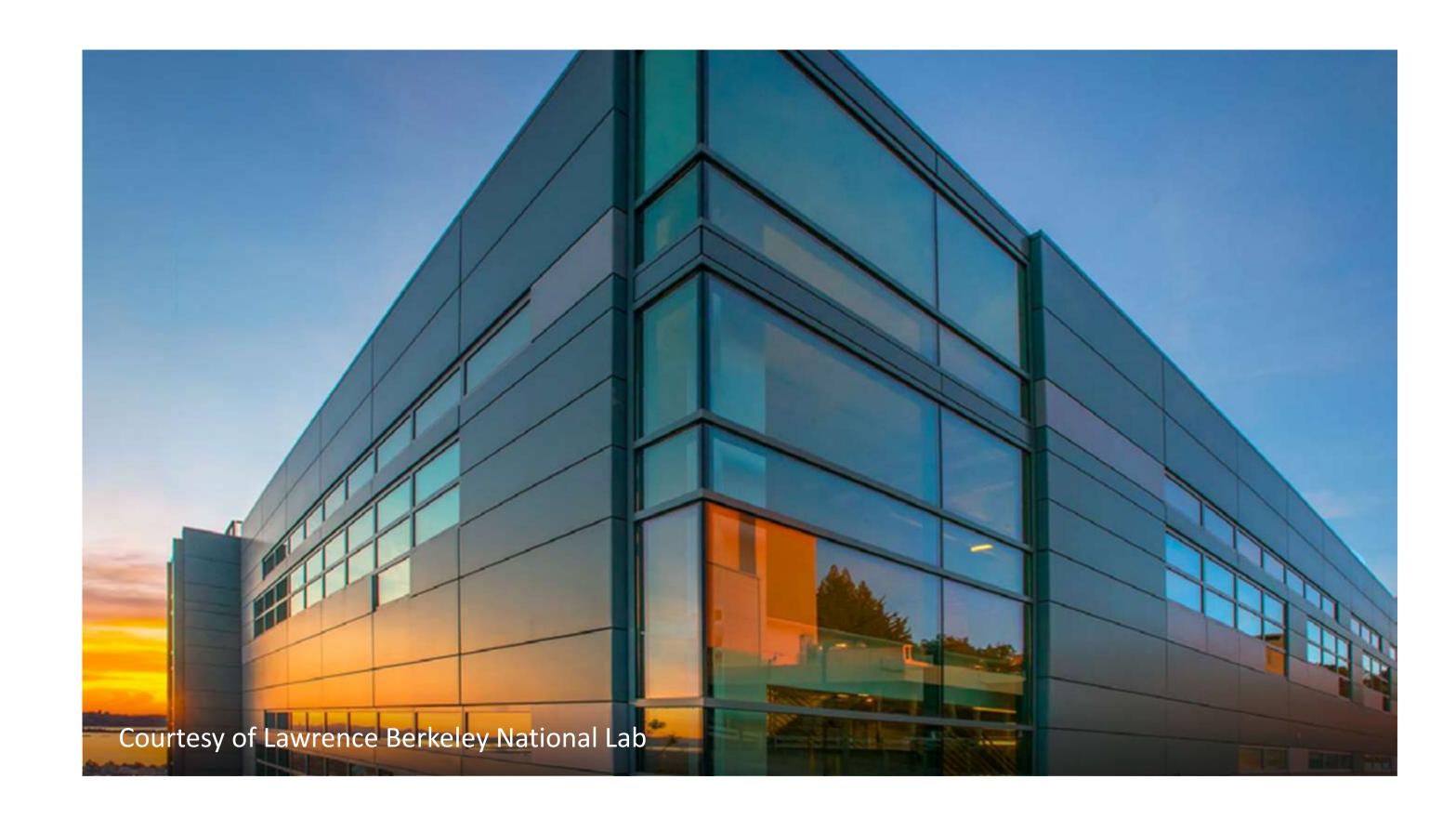
Based on measured performance of Modulus versus open-source CFD solvers like OpenFOAM. CPU: Intel Ice Lake 8380. GPU A100 80GB Tensor Core.



<sup>—</sup> Georg Rollmann, Head of Advanced Analytics and AI, Siemens Energy

# Driving Scientific Discovery

Application efficiency and performance for every supercomputing center



<sup>&</sup>quot;Researchers need energy-efficient, performant ways for their simulation and AI applications to support their work. At NERSC, we've adopted accelerated computing to support research across a broad variety of science applications."

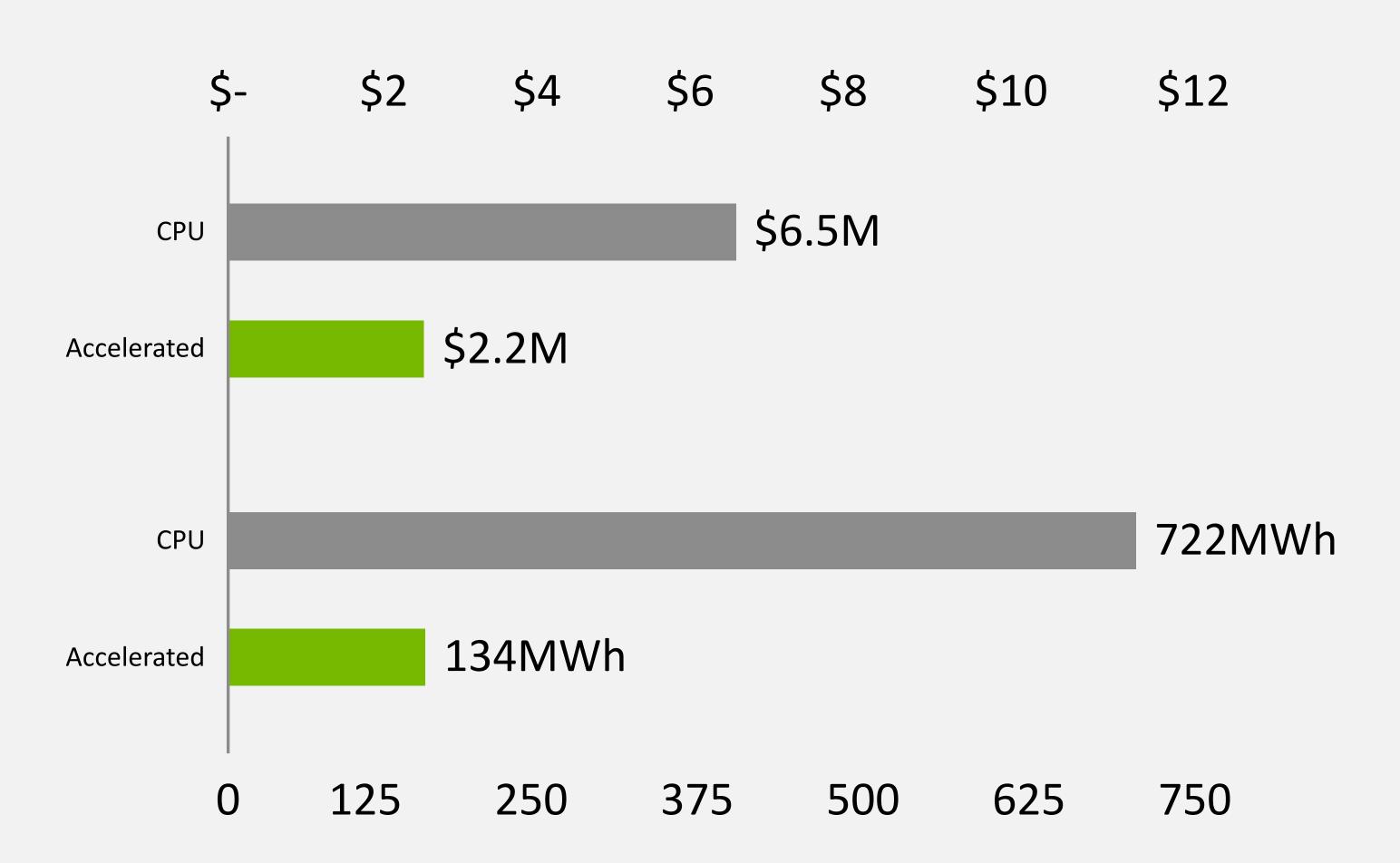
Nick Wright, Chief Architect, NERSC

#### **NVIDIA**

#### 3X Lower Cost and 5X Less Energy

#### **Acquisition Cost**

in Millions USD



# **Energy**Megawatt Hours (MWh)

Based on the geometric mean of application speedups and energy consumption vs. dual AMD 7763 / benchmark applications / BerkleyGW / DeepCAM / EXAALT / MILC and monthly prices on Microsoft Azure for instances standard\_NC96ADS\_A100\_v4 and standard\_D96ads\_v5

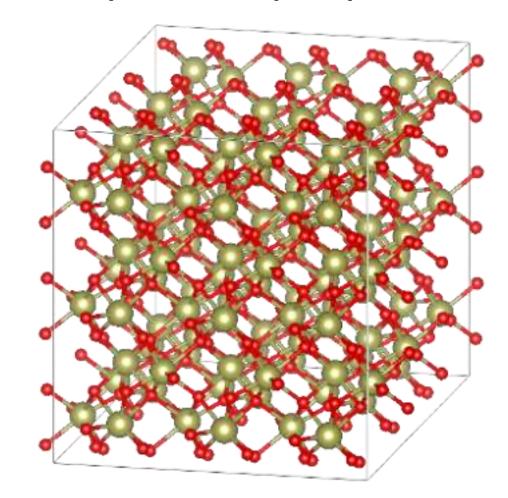


## Magnum 10 Multi-Node Scaling Performance

Magnum IO NVIDIA Common Collectives Library Delivers up to 2.3X Energy Efficiency

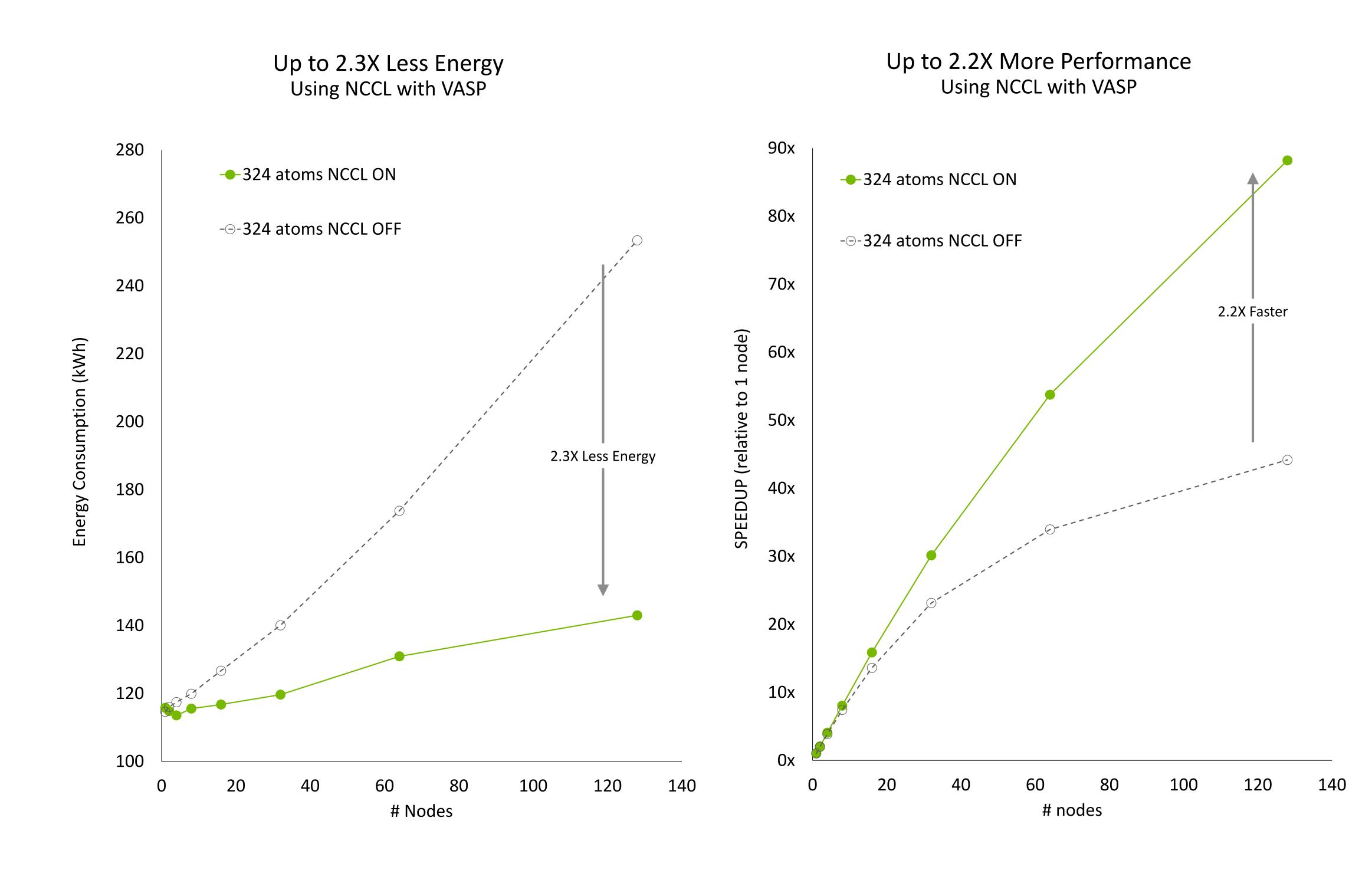


Performs quantum mechanical calculations for systems of atoms to predict properties



**Challenge:** Computation grows as n<sup>3</sup> with the number of atoms

**Solution:** Multi-node accelerated computing







# Learn. Innovate. Succeed.

Grow and validate your skills with training from the NVIDIA Deep Learning Institute.



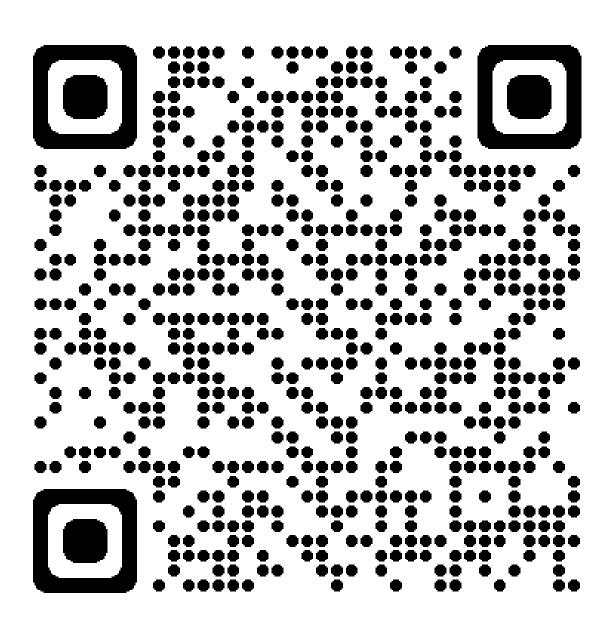
## Developers

http://developer.nvidia.com

"NVIDIA agora possui um ecossistema robusto e autossustentável de software, universidades, startups e parceiros que permitiram que ela se tornasse referência de um universo recém-criado." FORBES

Em menos de dois anos depois dobramos o número de desenvolvedores no nosso programa, isso mostra o grande interesse nas nossas plataformas.

Faça parte deste ecossistema:





# Build and Grow Your Generative Al Startup

**NVIDIA** Inception











#### **NVIDIA** Inception

https://www.nvidia.com/pt-br/startups/

O NVIDIA Inception é uma plataforma de aceleração para startups de AI, ciência de dados e HPC que oferece suporte, expertise e tecnologia essenciais para inserção no mercado.

"Participamos de uma série de briefings e eventos da NVIDIA em que fomos apresentados a pessoas que se tornaram nossos clientes. No GTC, você conhece outras empresas que estão desenvolvendo tecnologias complementares e investidores interessados na área que podem ajudar sua empresa." Adam Bry, Cofundador e CEO, Skydio







# Explore Breakthroughs in Al and Beyond With the Best of GTC





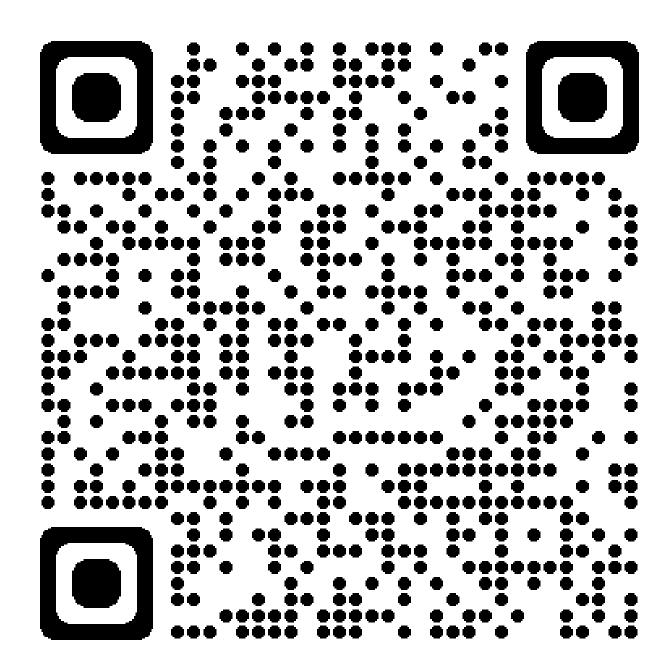


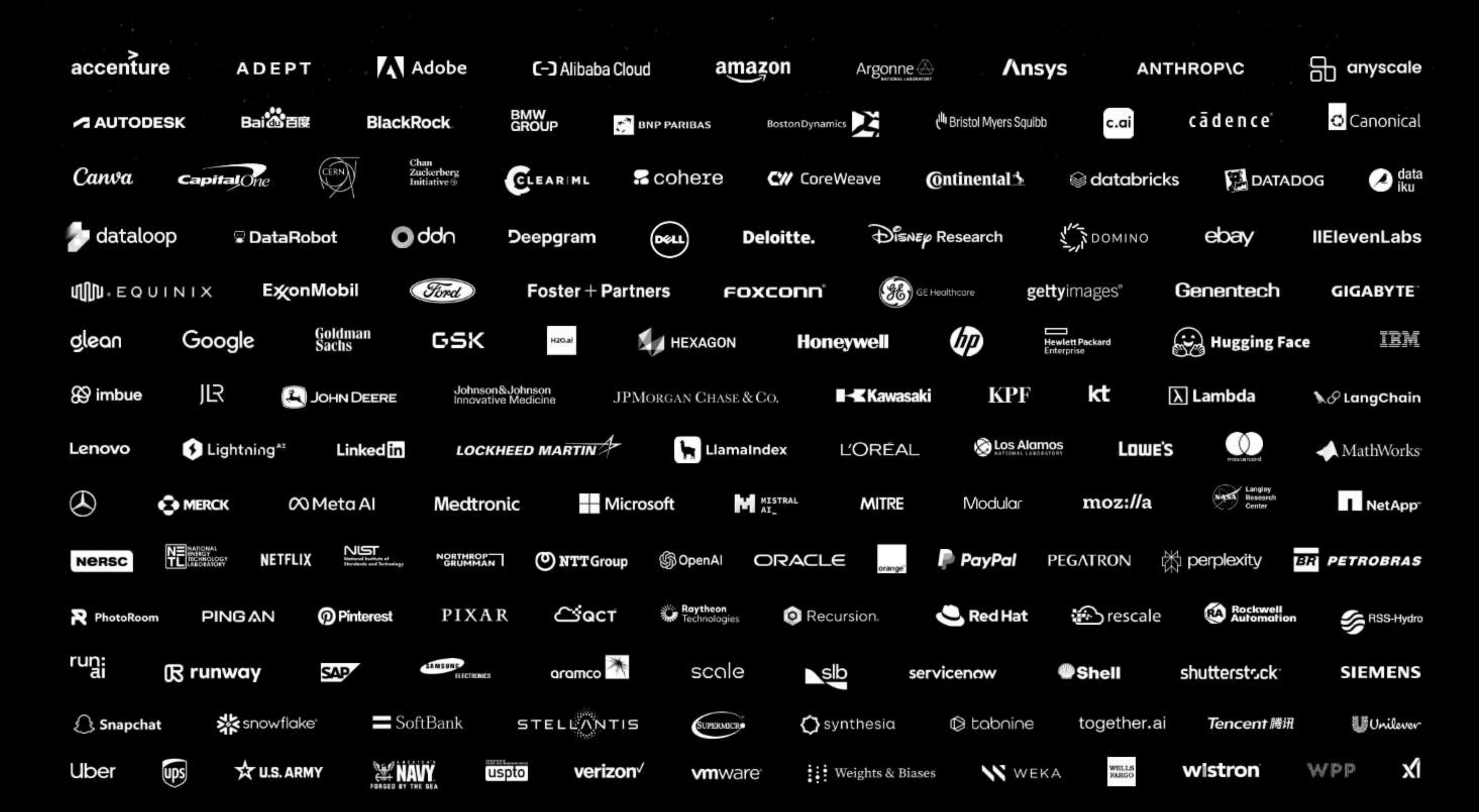
#### **NVIDIA GTC**

https://www.nvidia.com/gtc/

Com mais de 650 sessões gratuitas sobre o que há de mais moderno em em IA generativa, metaverso, grandes modelos de linguagem ampla (LLMs), robótica, computação em nuvem e muito mais.

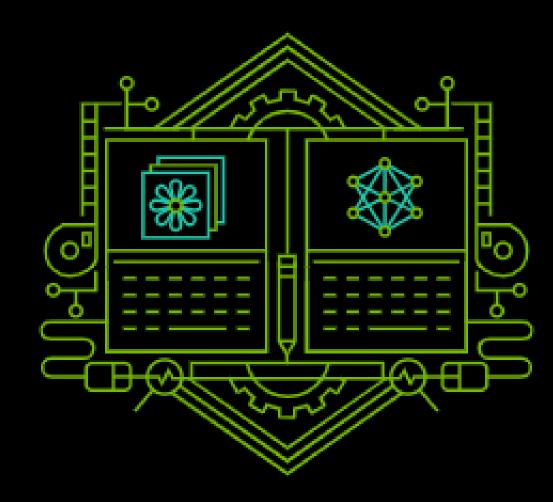
Os participantes podem obter insights e treinamento sobre as tecnologias avançadas que transformam as diversas indústrias.





NVIDIA DLI Teaching Kits

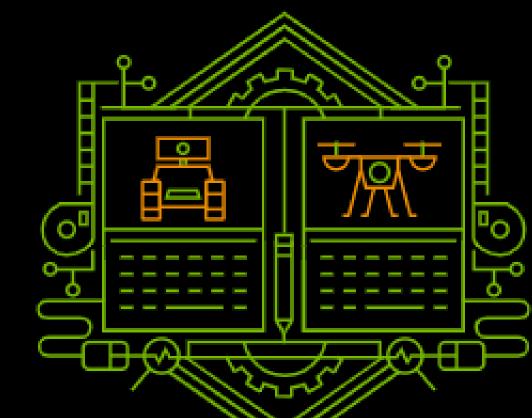
The Foundation for Next-Generation Innovators

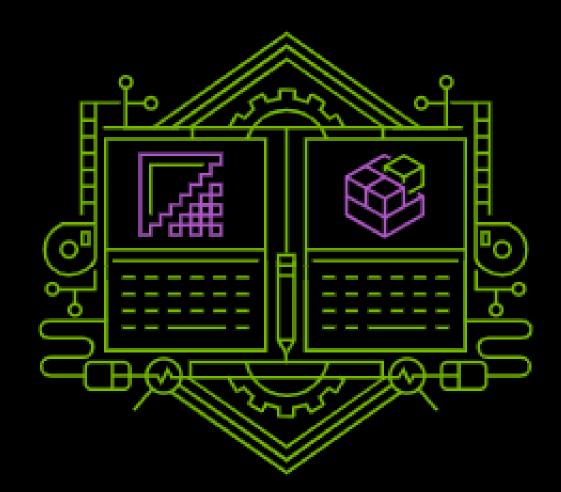


#### Deep Learning

Co-developed with Professor Yann LeCun and his team at New York University (NYU), this Teaching Kit leverages the latest computing frameworks and techniques to explore introductory and advanced deep learning topics, from image classification to generative adversarial networks (GANs) to natural language processing (NLP).

 Watch Video: Bringing AI to the Classroom: NVIDIA's Deep Learning Teaching Kit





#### Accelerated Computing

Co-developed with Professor Wen-Mei Hwu and his team from the University of Illinois (UIUC) and Professor Sunita Chandrasekaran and her team from the University of Delaware, this Teaching Kit covers introductory and advanced topics such as parallel programming APIs, programming tools and techniques, principles and patterns of paralgorithms, and processor architecture features and constraints.

 Watch Video: Accelerated Computing Teaching Kit for University **Educators: Introduction and Use Cases** (60 Minutes)



# DLI Teaching Kits

https://www.nvidia.com/en-us/training/teaching-kits/

- Removendo Barreiras para o Ensino de Novas Tecnologias
  - Economiza Tempo de Planejamento de Aulas e Laboratórios
  - Reduz Custos e Necessidades de Infraestrutura
  - Aborda Teoria Acadêmica e Fundamentos
  - Oferta de Curso Singular e Abrangente com Suporte
- Deep Learning
- Accelerated Computing
- Edge Al and Robotics
- Data Science
- Graphics and Omniverse
- Science and Engineering
- Generative AI (NEW)





# Obrigado!

#### Fabio Alves de Oliveira

NVIDIA – Gerente de Negócios fdeoliveira@nvidia.com

Realização:















