AAEON an ASUS assoc. co. · acer · ADLINK · ADVANTECH · aetina · AIMobile Industrial Mobile Systems · aira · 安宏生醫 AnHorn Medicines Innovation · New Therapeutics · APMIC · AppWorks

ASE · ASRock Rack · ASUS · Avalanche Computing · AVerMedia · AXIOMTEK · B!gGo · CHENBRO · ChoozMo 集仕多股份有限公司 · 中原大學 Chung Yuan Christian University

COLORFUL · COMPAL · COOLER MASTER · Coretronic · DeepMentor · DeepRad.AI · DELTA · DMKTZ · EDOM TECHNOLOGY CO., LTD. · EverFocus Your Safety. Our Focus

Footprint-AI Making machine learning for everyone · FORTUNE AI TECHNOLOGIES · FOXCONN HON HAI TECHNOLOGY GROUP · Garage+ · GIANT GROUP · GIGABYTE · GLIACLOUD · GMI

HOMEE.AI · 義守大學 I-SHOU UNIVERSITY · INFINITIES 數位無限軟體 · ingrasys · INNO3D · Inventec · InWin · KENMEC 廣運 · KYEC

Lanner · LEADTEK · LEDA TECHNOLOGY · 律果科技 Legal Tech Inc. · Linker Vision · LITEON 光寶科技 光電事業部 · Manli · MEDIATEK

MetAI · MiTAC · msi · MEI · National Cheng Kung University 國立成功大學 1931 · 國立臺灣科技大學 NATIONAL TAIWAN UNIVERSITY OF SCIENCE AND TECHNOLOGY

國立陽明交通大學 NATIONAL YANG MING CHIAO TUNG UNIVERSITY · 財團法人國家實驗研究院 NCHC 國家高速網路與計算中心 National Center for High-performance Computing · Neousys TECHNOLOGY · NEXCOM · NHRI · onyx Smart Healthcare an ASUS assoc. co.

PALIT · PEGATRON · PROFET AI · QCT · Quanta Computer · SHIH CHIEN UNIVERSITY · SOLOMON · 南臺科技大學 · SPIL 矽品精密 Siliconware

SPINGENCE · StarFab accelerator · Stream Teck · SUPERMICRO · TAMKANG UNIVERSITY · TM ROBOT · TRI innovation · thermaltake · TREND MICRO · tsmc

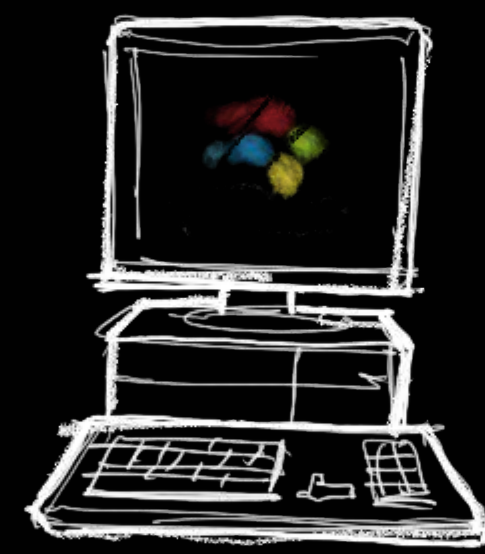東海大學 TUNGHAI UNIVERSITY · TWS TAIWAN WEB SERVICE · UMC · Unimicron 欣興電子 · Vecow · Wistron · wiwynn · YUAN · ZOTAC
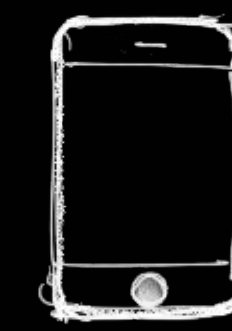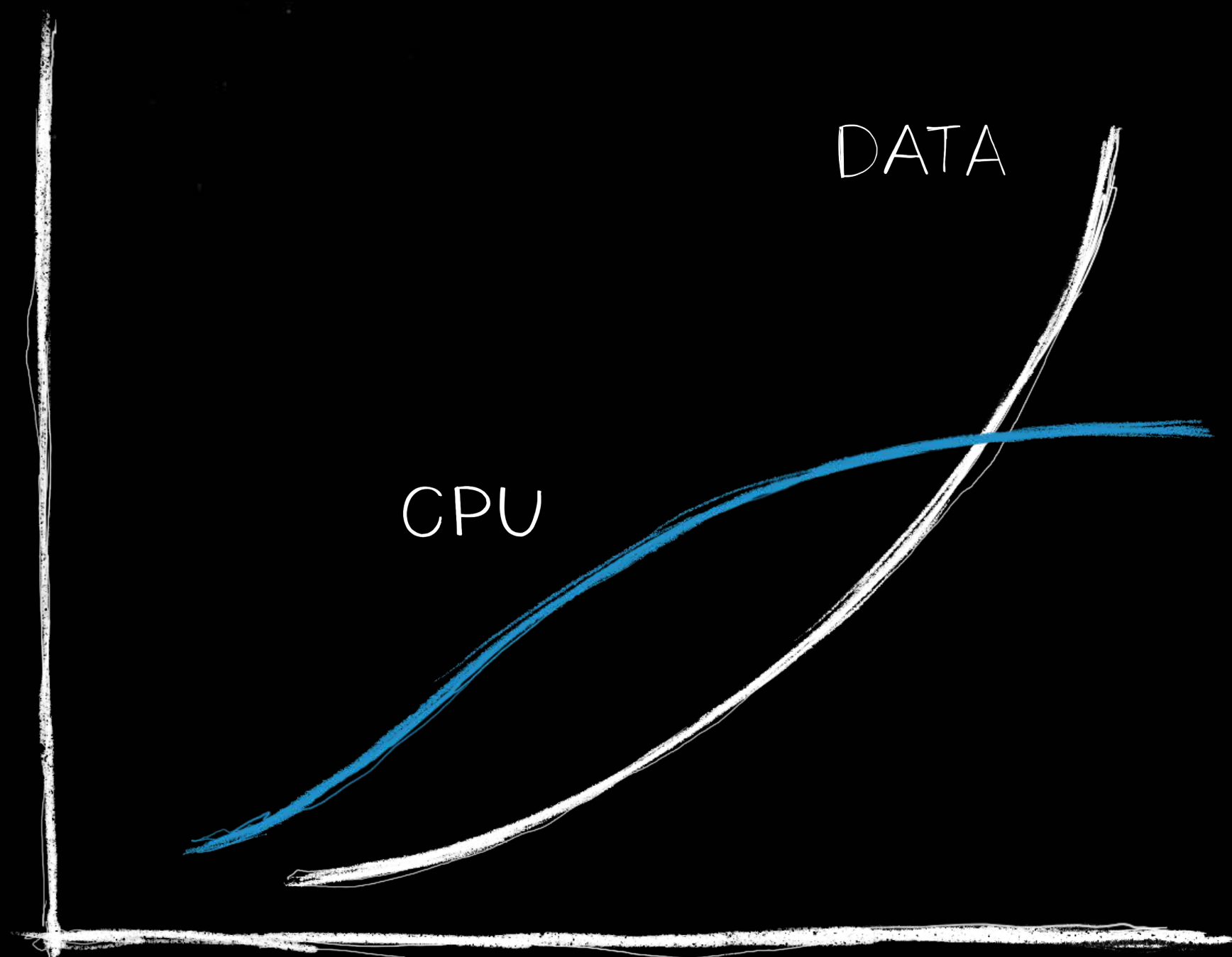
1964
IBM S/360

1995
WINDOWS 95
PENTIUM

2007
IPHONE
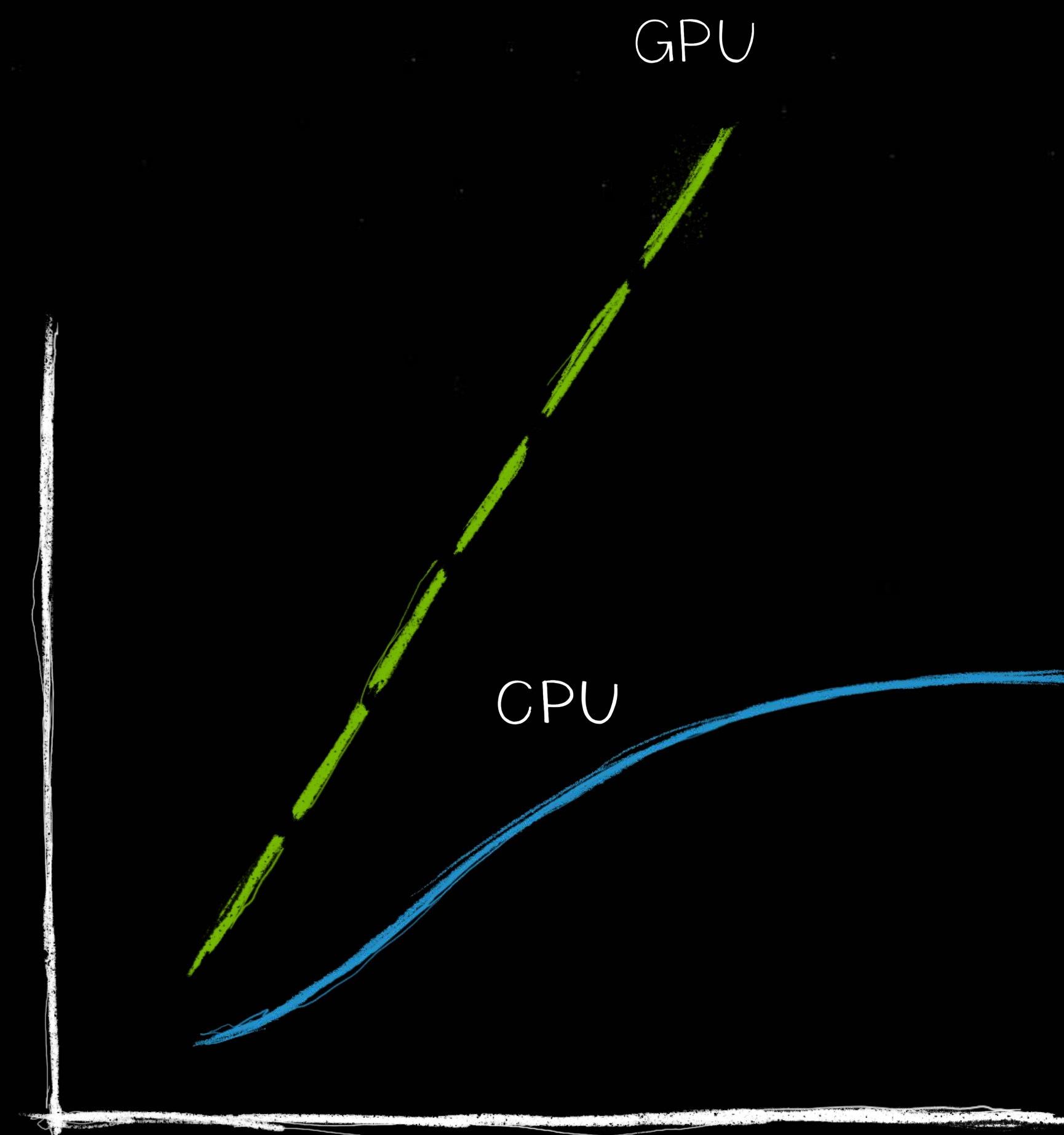
2007
INTERNET
MOBILE CLOUD

新的運算時代正在開始

A NEW COMPUTING AGE IS STARTING
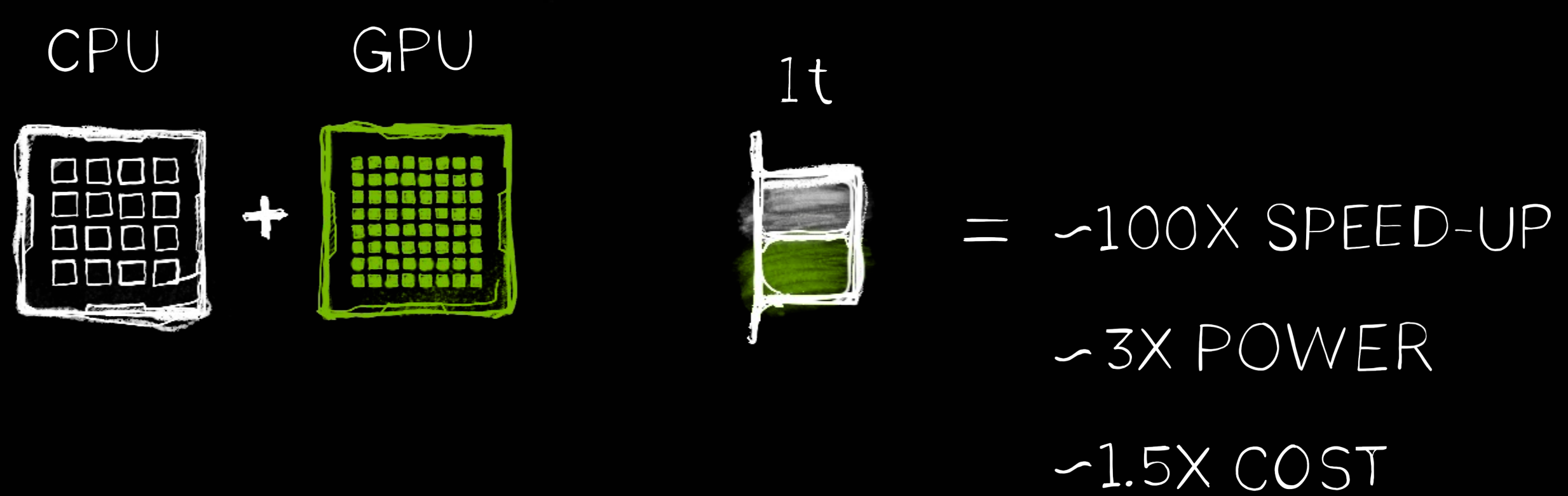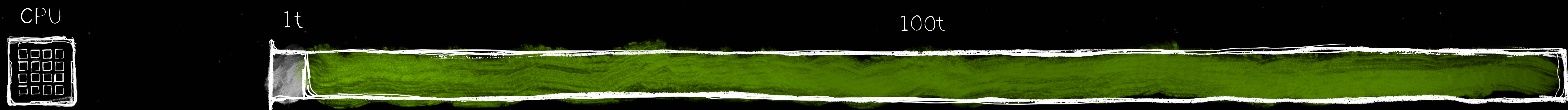
CPU SCALING SLOWS

...AND
COMPUTE DEMAND GROWS
EXPONENTIALLY

GPU-ACCELERATED COMPUTING

2006
CUDA

加速每個應用程式
ACCELERATE EVERY APPLICATION

CPU

1t                                                                                100t

CPU   +   GPU

1t

= ~100X SPEED-UP

~3X POWER

~1.5X COST

60X PERF / $     OR     98% SAVINGS
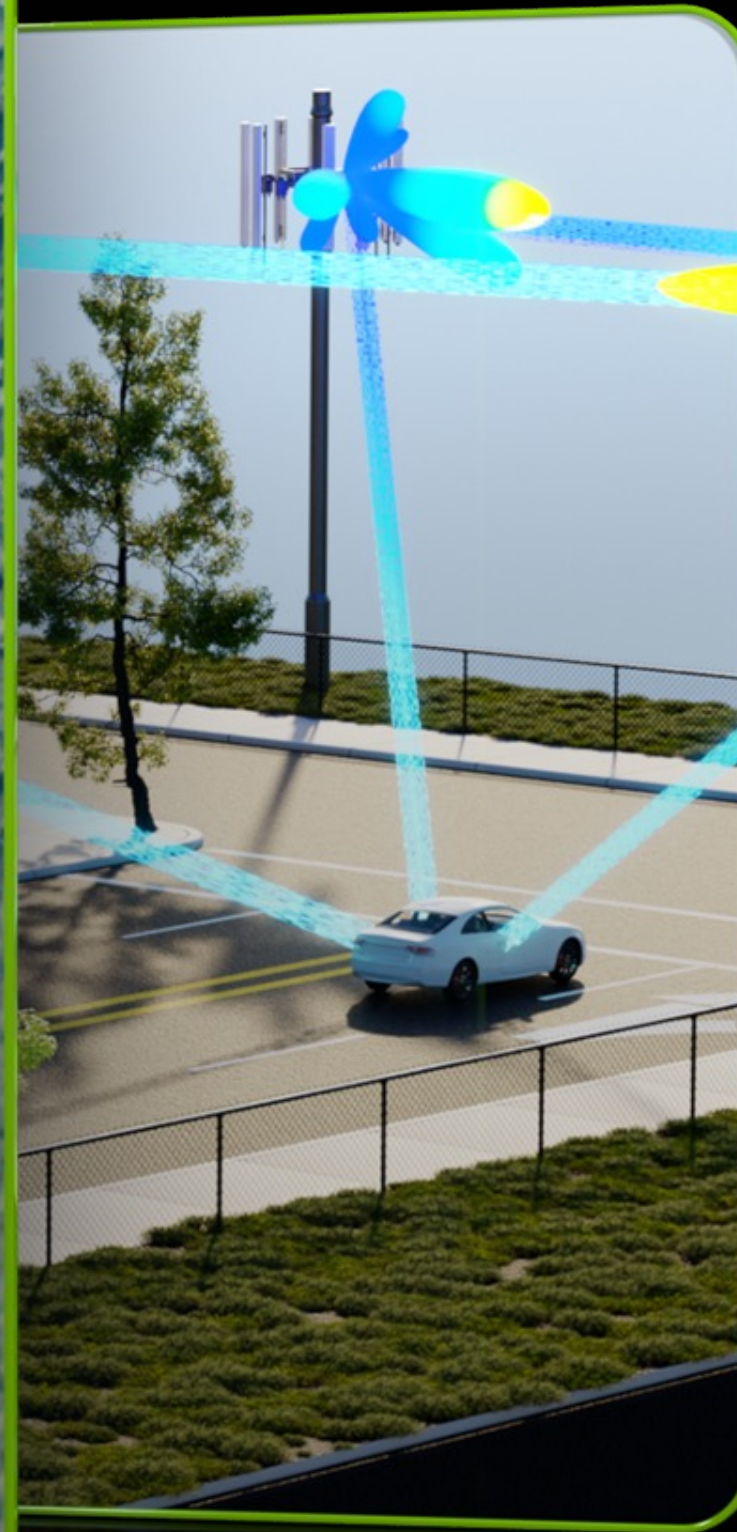
30X PERF / W     OR     97% SAVINGS

買越多 … 省越多
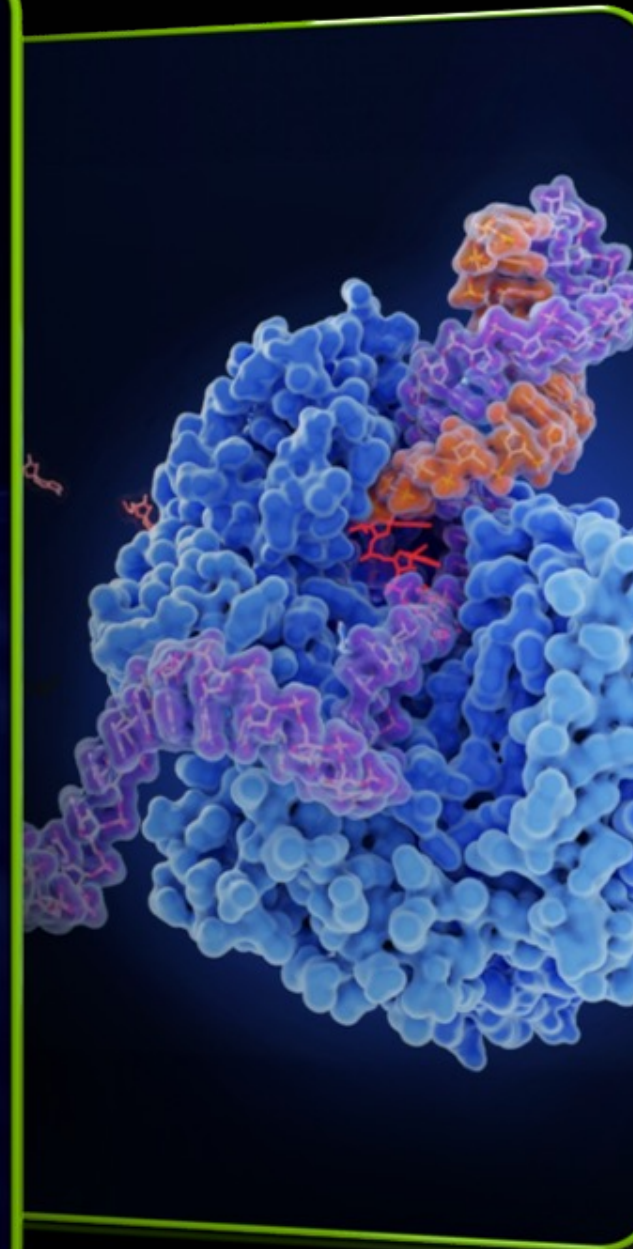
"THE MORE YOU BUY… THE MORE YOU SAVE"
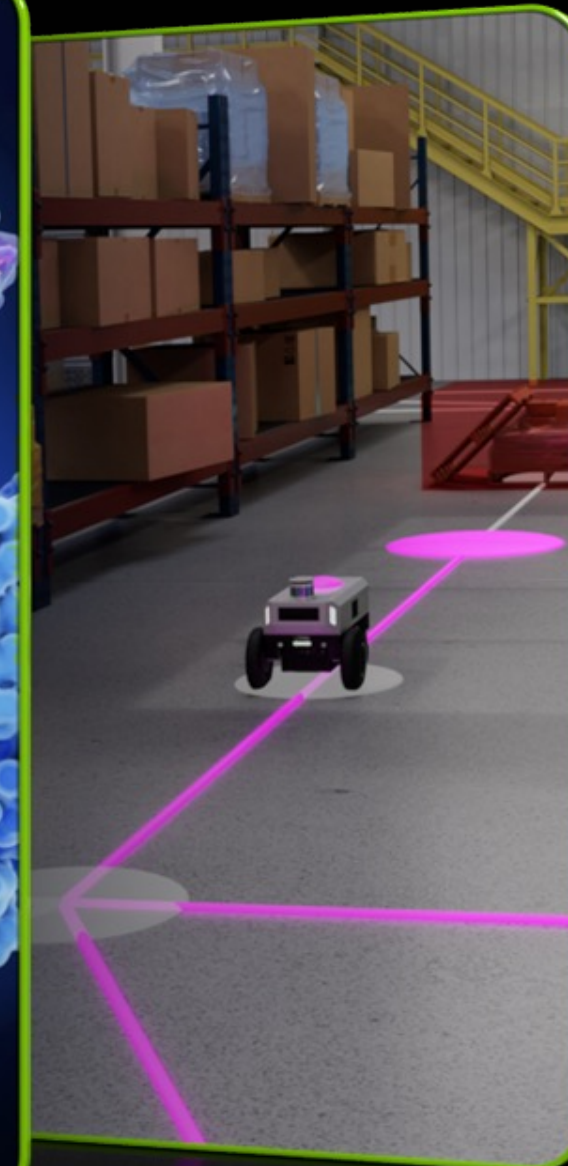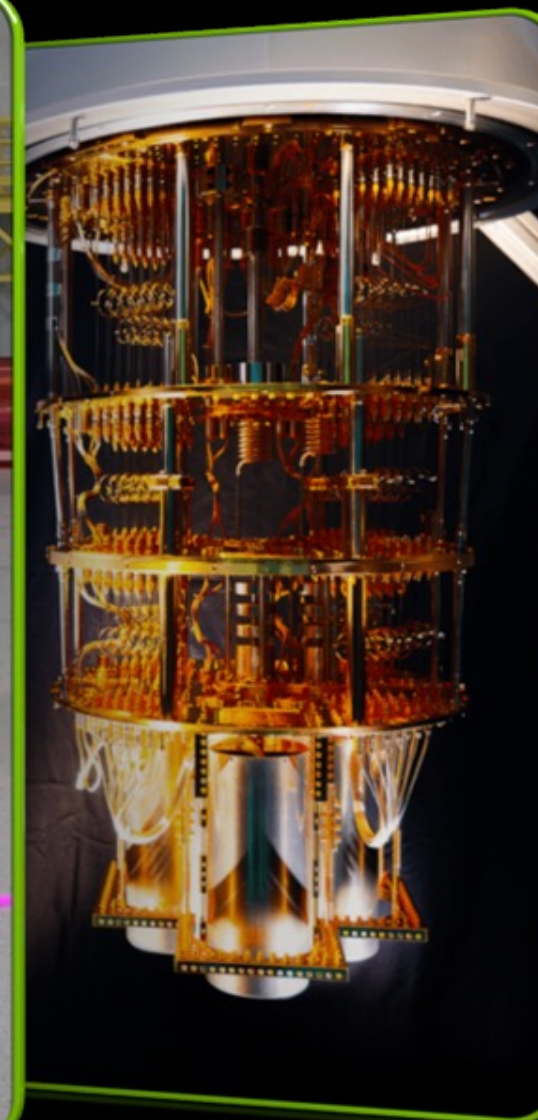
cuDNN
Deep Learning

Modulus
AI Physics

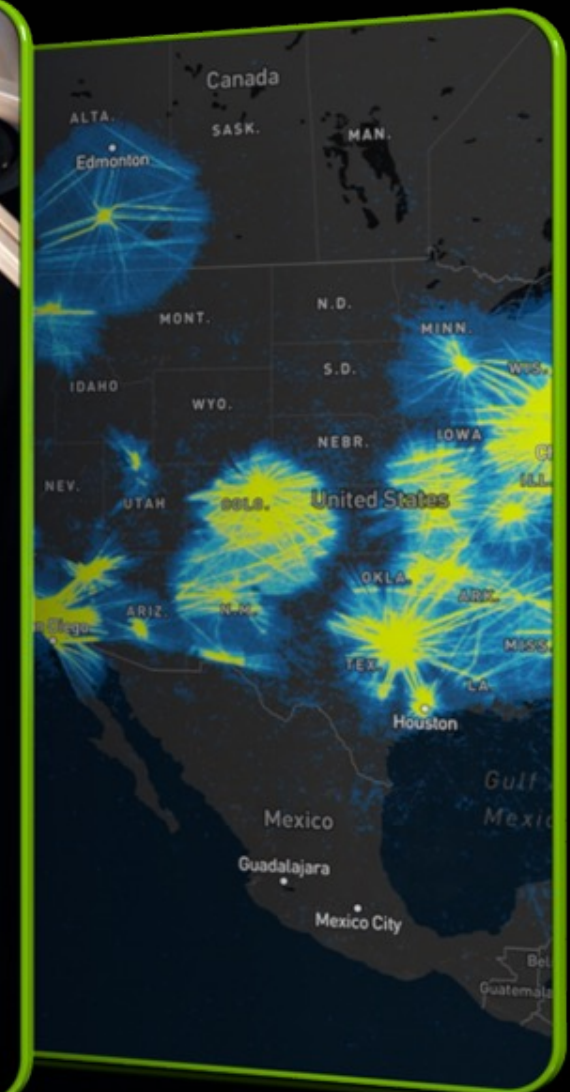Aerial RAN
AI Radio

cuLITHO
Computational Lithography

Parabricks
Gene Sequencing

cuOPT
Combinatorial Optimization

cuQUANTUM
QC Simulation

cuDF
Data Processing

NVIDIA CUDA 函式庫開拓新市場
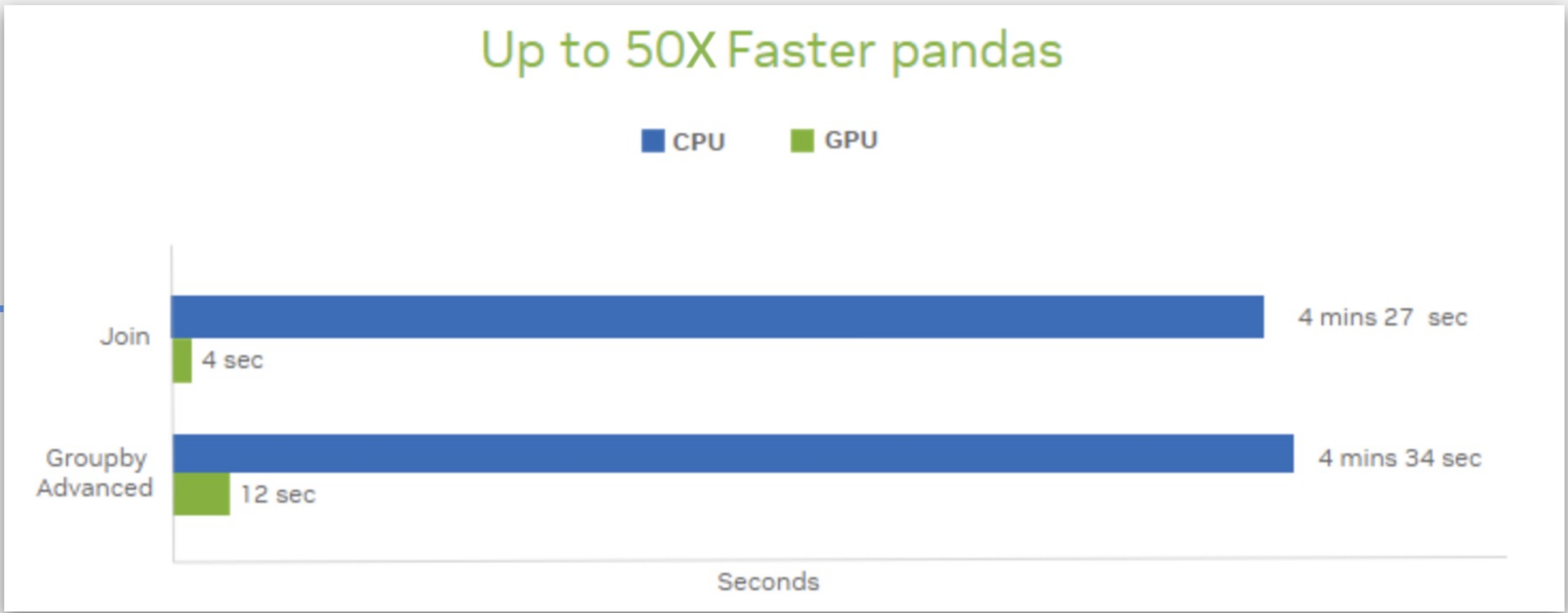
NVIDIA CUDA LIBRARIES OPEN NEW MARKETS

NVIDIA 宣佈推出 GOOGLE COLAB 中的 PANDAS-CUDF

NVIDIA ANNOUNCES PANDAS-CUDF IN GOOGLE COLAB

CUDA 實現良性循環

CUDA ACHIEVES VIRTUOUS CYCLE

2012
ALEXNET
"FIRST CONTACT"

DATA

2016
FIRST DGX DELIVERED

DATA

2017
TRANSFORMER

DATA

ChatGPT 3.5 ⌄                                                          JE

Email for plumber quote    Make me a personal webpage    Overcome procrastination    Design a fun coding game
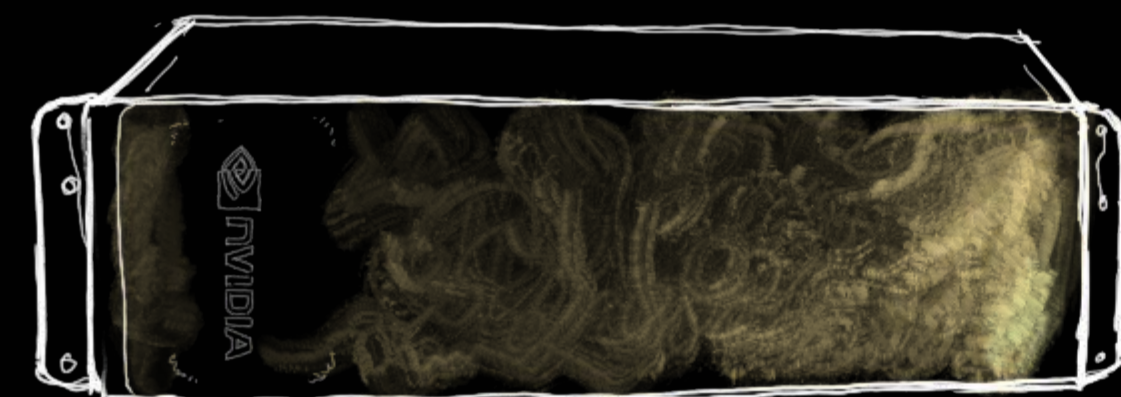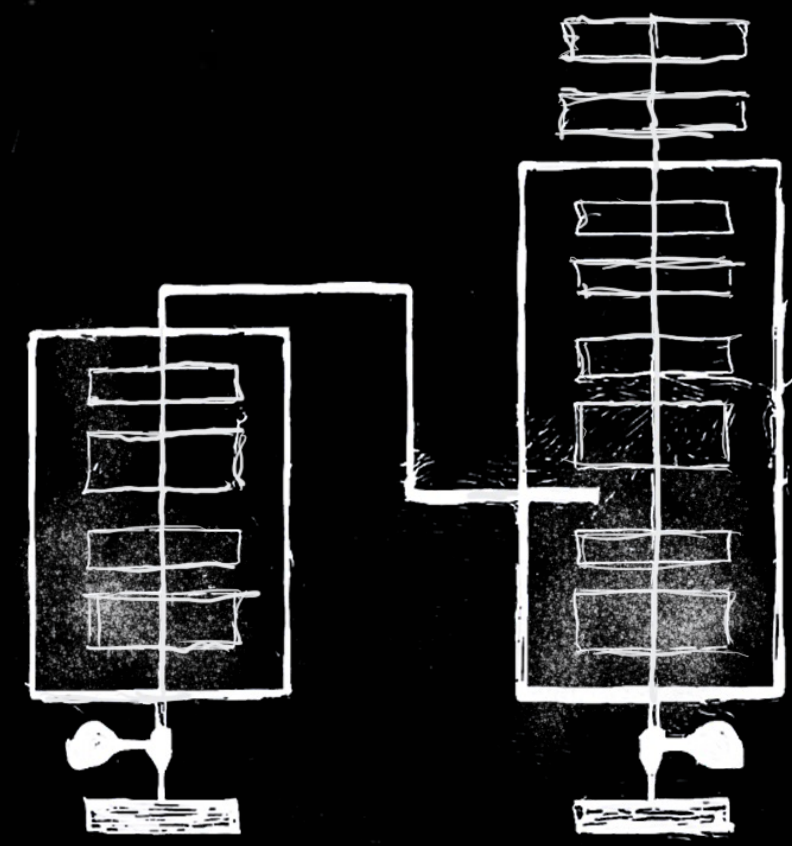
Message ChatGPT

ChatGPT can make mistakes. Check important info.

LLM
Large Language Model

2023
CHATGPT
"THE BIG BANG OF AI"

DATA

AI FACTORY

$100T

MANUFACTURING
TRANSPORTATION
HEALTHCARE
COMPUTING

新產業革命

"A NEW INDUSTRIAL REVOLUTION"

SOFTWARE FACTORY

AI FACTORY

TOOLS

RETRIEVAL

INSTRUCTIONS

**CPU**

SKILLS

GENERATIVE

LLMS

**GPU**

生成式人工智慧推動全棧重塑

GENERATIVE AI DRIVES FULL STACK REINVENTION

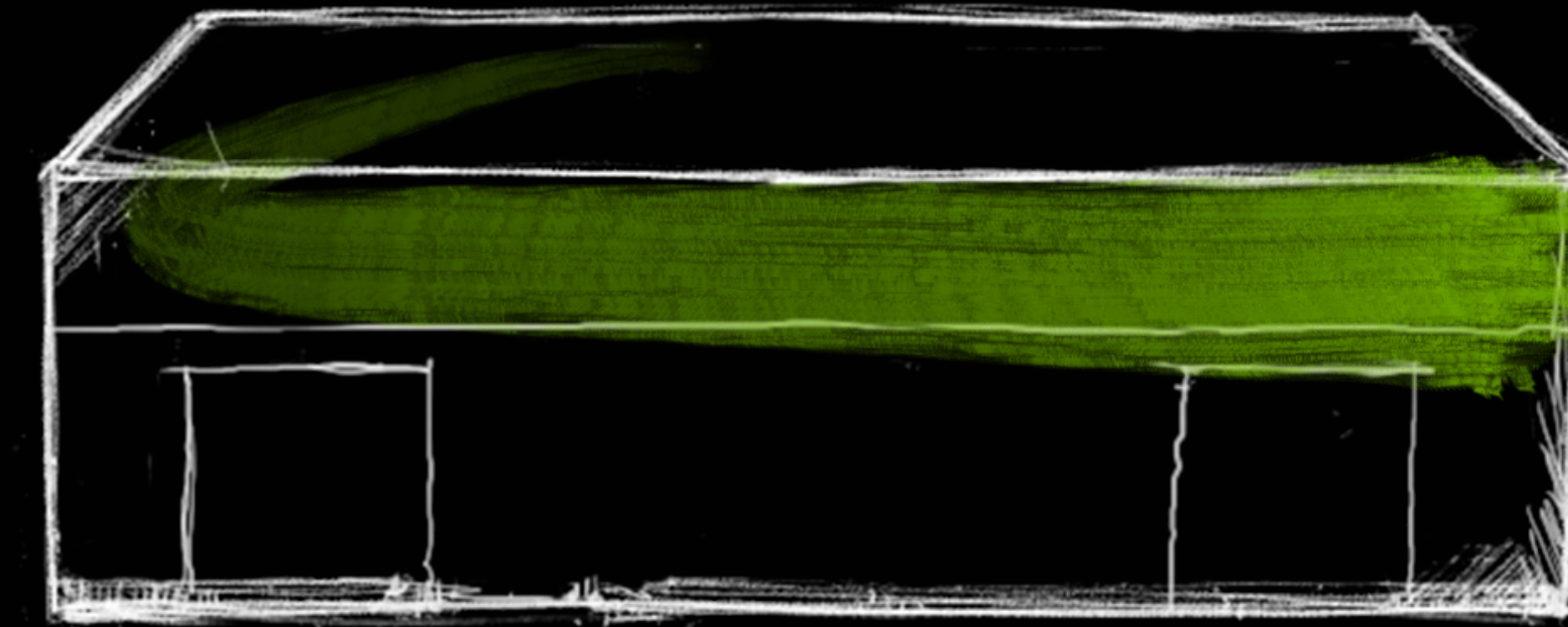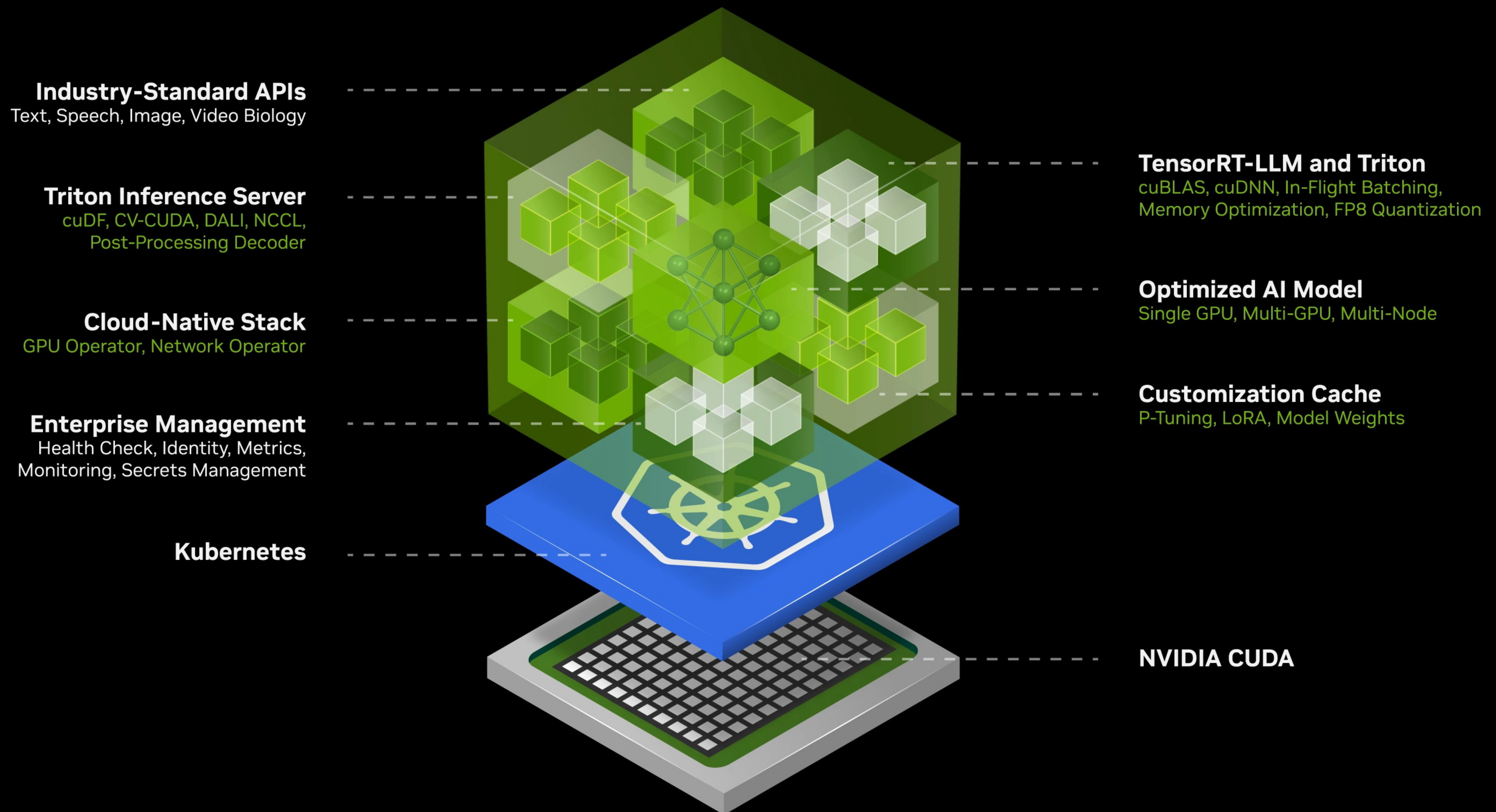**NVIDIA INFERENCE MICROSERVICE**

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

**Industry-Standard APIs**
Text, Speech, Image, Video Biology

**Triton Inference Server**
cuDF, CV-CUDA, DALI, NCCL,
Post-Processing Decoder

**Cloud-Native Stack**
GPU Operator, Network Operator

**Enterprise Management**
Health Check, Identity, Metrics,
Monitoring, Secrets Management

**Kubernetes**

**TensorRT-LLM and Triton**
cuBLAS, cuDNN, In-Flight Batching,
Memory Optimization, FP8 Quantization

**Optimized AI Model**
Single GPU, Multi-GPU, Multi-Node

**Customization Cache**
P-Tuning, LoRA, Model Weights

**NVIDIA CUDA**

**Installed Base of 100s of Millions of CUDA GPUs**

**Speech**  **Digital Human**  **Computer Vision**  **Biology**  **Simulation**

**Language**  **Regional Language**  **Vision Language**  **RAG**

AI.NVIDIA.COM

生成式人工智慧實現數位護士、客服、導師等

GENERATIVE AI ENABLES DIGITAL NURSE, CUSTOMER SERVICE AGENTS, TUTORS, ETC

ASUS TUF A14 / A16          ASUS Zephyrus G16          ASUS ProArt PX13 / P16          MSI Stealth A16 AI⁺

發表新款 RTX AI 電腦

現已超過 200 款 RTX AI 筆記型電腦 | 高達 700 AI TOPS | 7X 生成式人工智慧

ANNOUNCING NEW RTX AI PCs

Now Over 200 RTX AI Laptops | Up to 700 AI TOPS | 7X Generative AI

PHYSICAL AI
2X/ 3 MON

TRANSFORMER
2X/ 6 MON

ALEXNET
2X/YR

運算模型規模呈指數級成長

MODEL COMPUTE SCALE GROWS EXPONENTIALLY

AI SUPERCHIP
208B Transistors
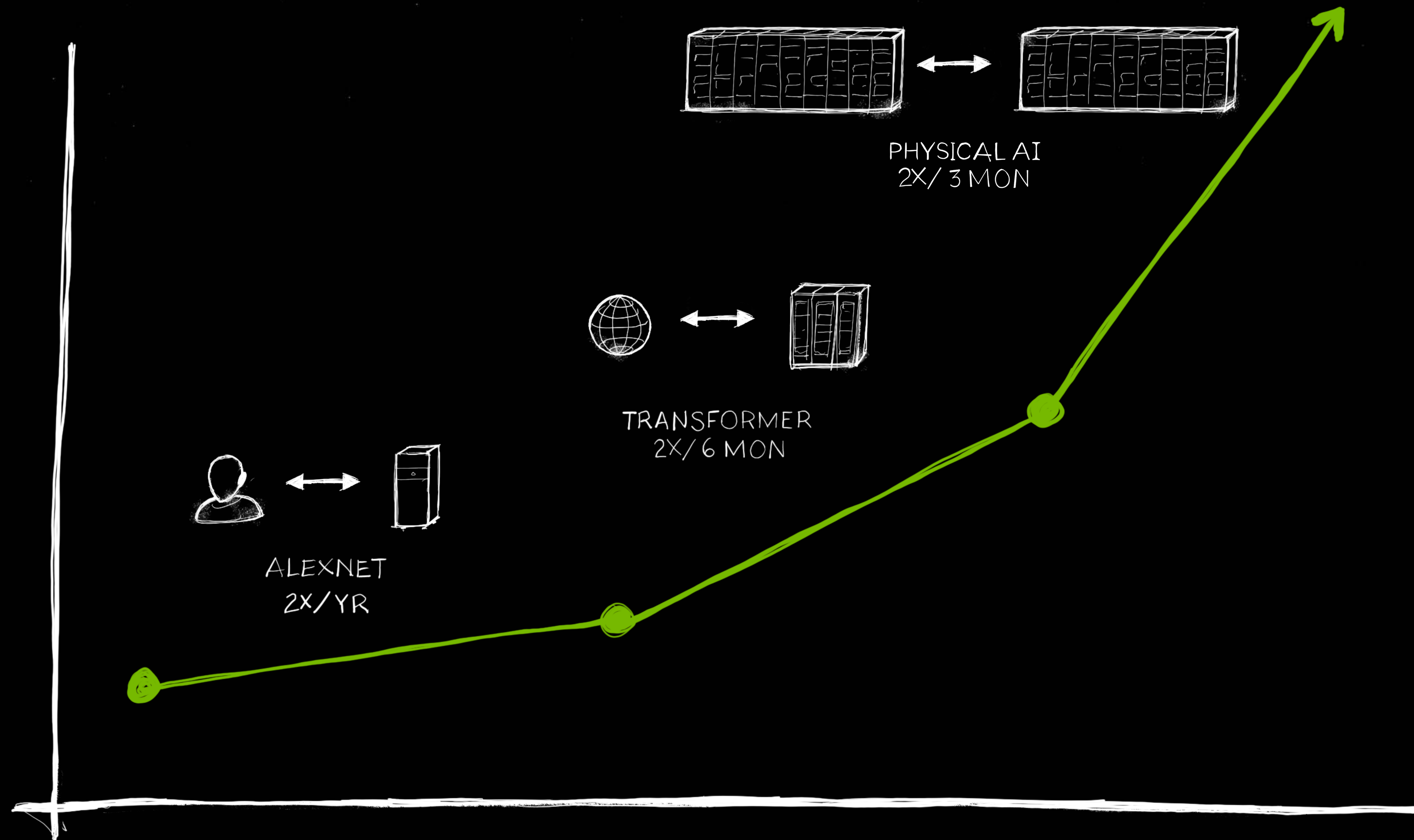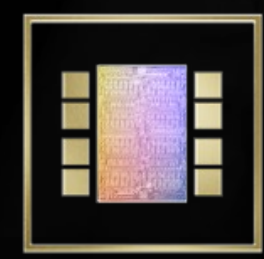
TRANSFORMER ENGINE
FP4/FP6 Tensor Core

SECURE AI
Full Performance
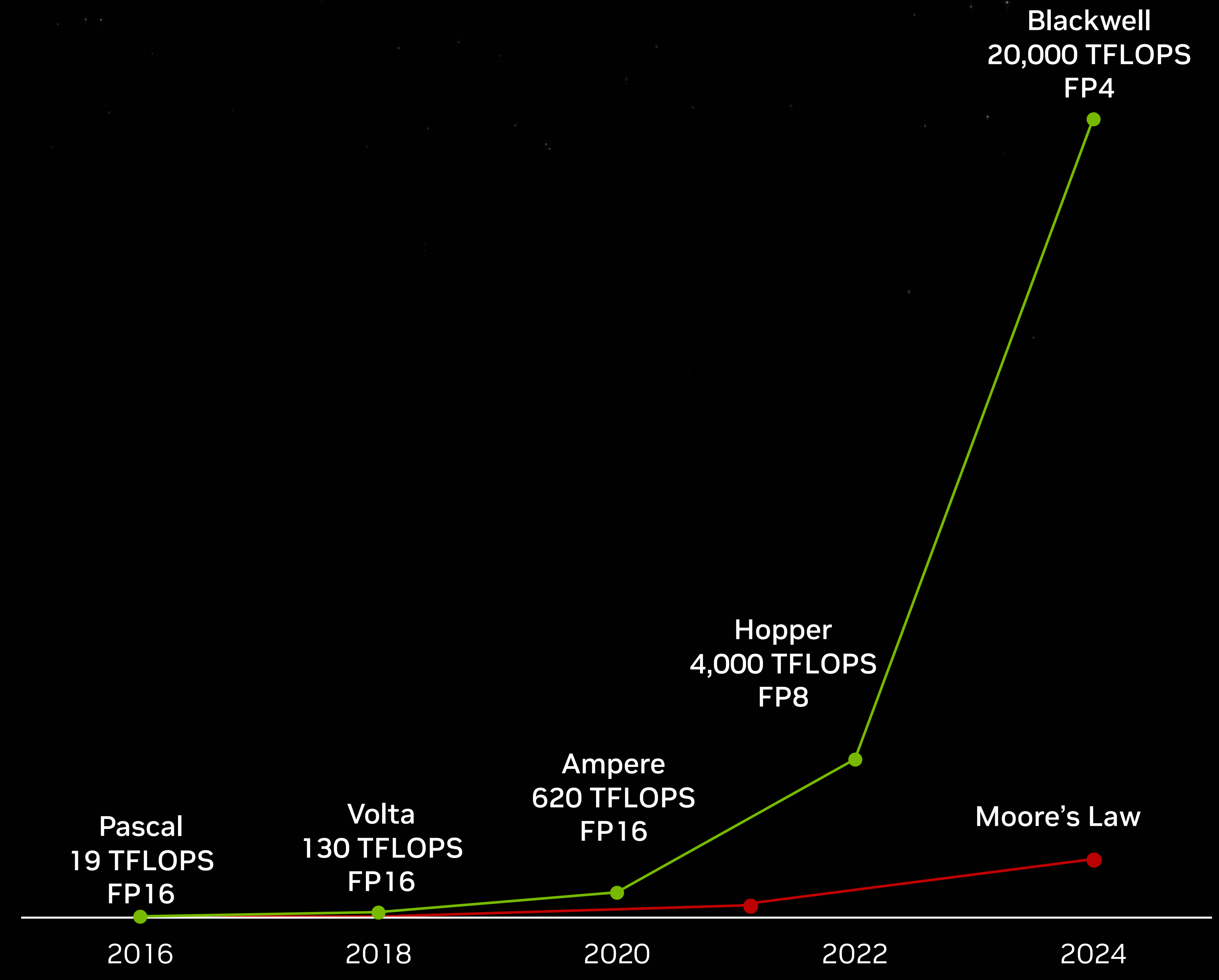Encryption and TEE

5th GENERATION NVLINK
Scales to 576 GPUs

RAS ENGINE
100% In-System Self-Test

DECOMPRESSION ENGINE
800 GB/sec

NVIDIA BLACKWELL 平台
兆級參數規模的生成式人工智慧

Blackwell
20,000 TFLOPS
FP4

Hopper
4,000 TFLOPS
FP8

Ampere
620 TFLOPS
FP16

Moore's Law

Pascal
19 TFLOPS
FP16

Volta
130 TFLOPS
FP16

2016    2018    2020    2022    2024

8 年內 1,000X 的人工智慧運算
1,000X AI COMPUTE IN 8 YEARS

Pascal
1,000+ GWh

Volta
140 GWh

Ampere
40 GWh

Hopper
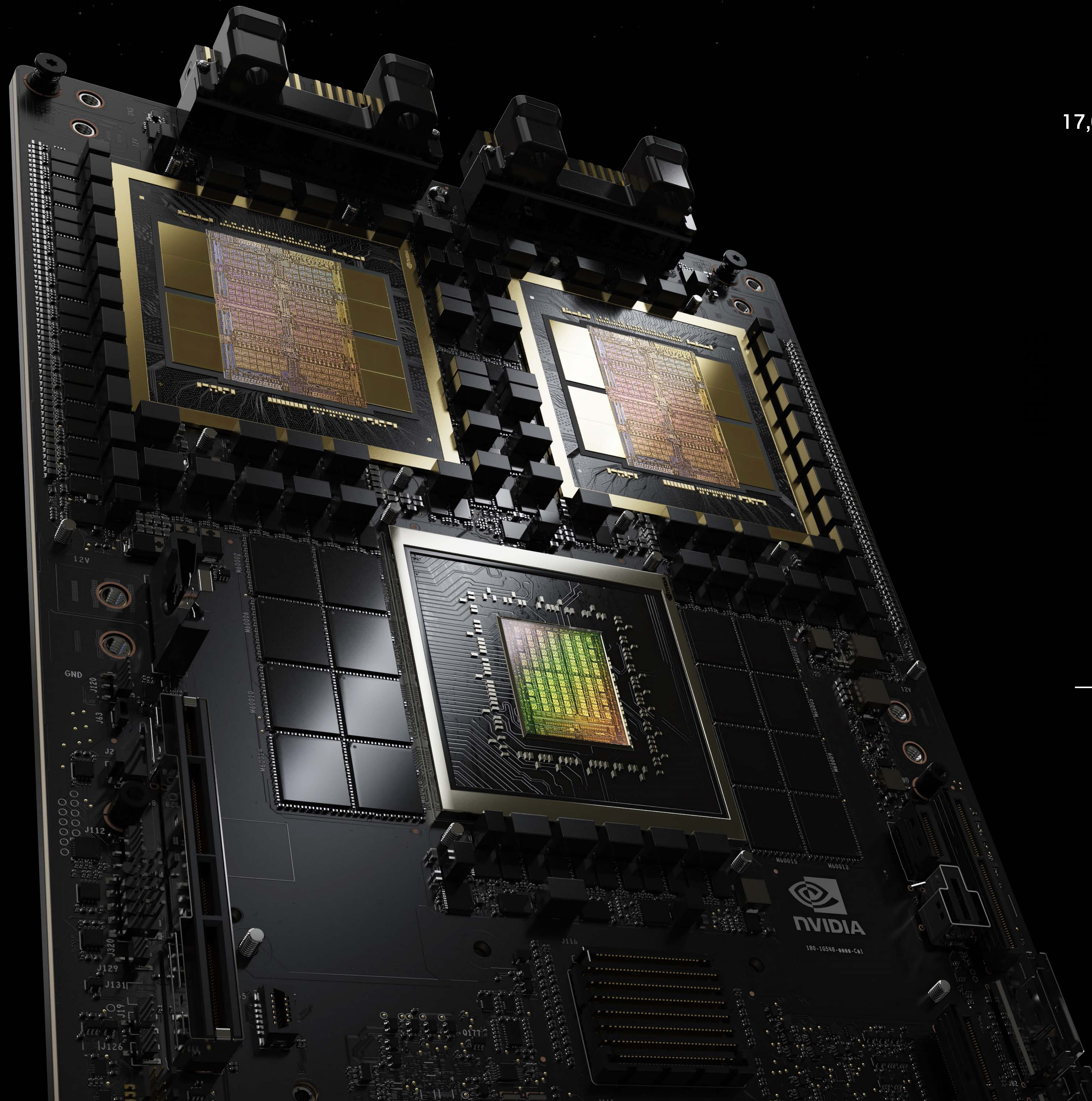13 GWh

Blackwell
3 GWh

2016 2018 2020 2022 2024

8 年內節能至 1/350
訓練 GPT4-1.8T

350X ENERGY REDUCTION IN 8 YEARS
TO TRAIN GPT4-1.8T

Pascal
17,000 Joules/Token

Volta
1,200 Joules/Token

Ampere
150 Joules/Token

Hopper
10 Joules/Token

Blackwell
0.4 Joules/Token

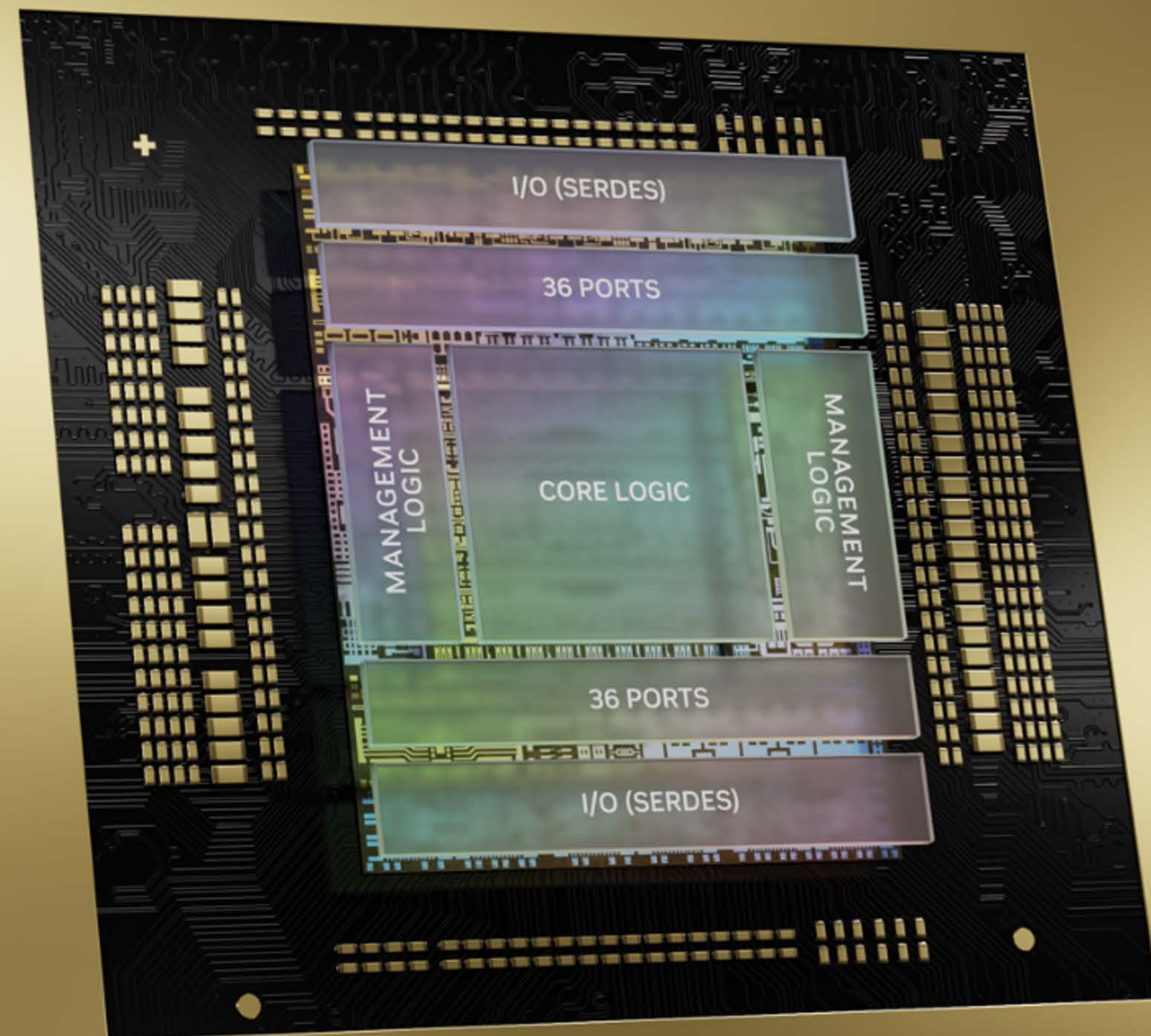2016  2018  2020  2022  2024

8 年內節能至 1/45,000
生成 GPT4-1.8T 詞元

45,000X ENERGY REDUCTION IN 8 YEARS
TO GENERATE GPT4-1.8T TOKENS

|  | DGX BLACKWELL | DGX HOPPER |  |
|---|---|---|---|
| NVLink Domain | 72 | 8 | 9X |
| NVLink BW (TB/s) | 130 | 7 | 18X |
| AI FLOPS (PF) | 1,440 | 32 | 45X |
| Power (kW) | 100 | 10 | 10X |

NVLink Switch
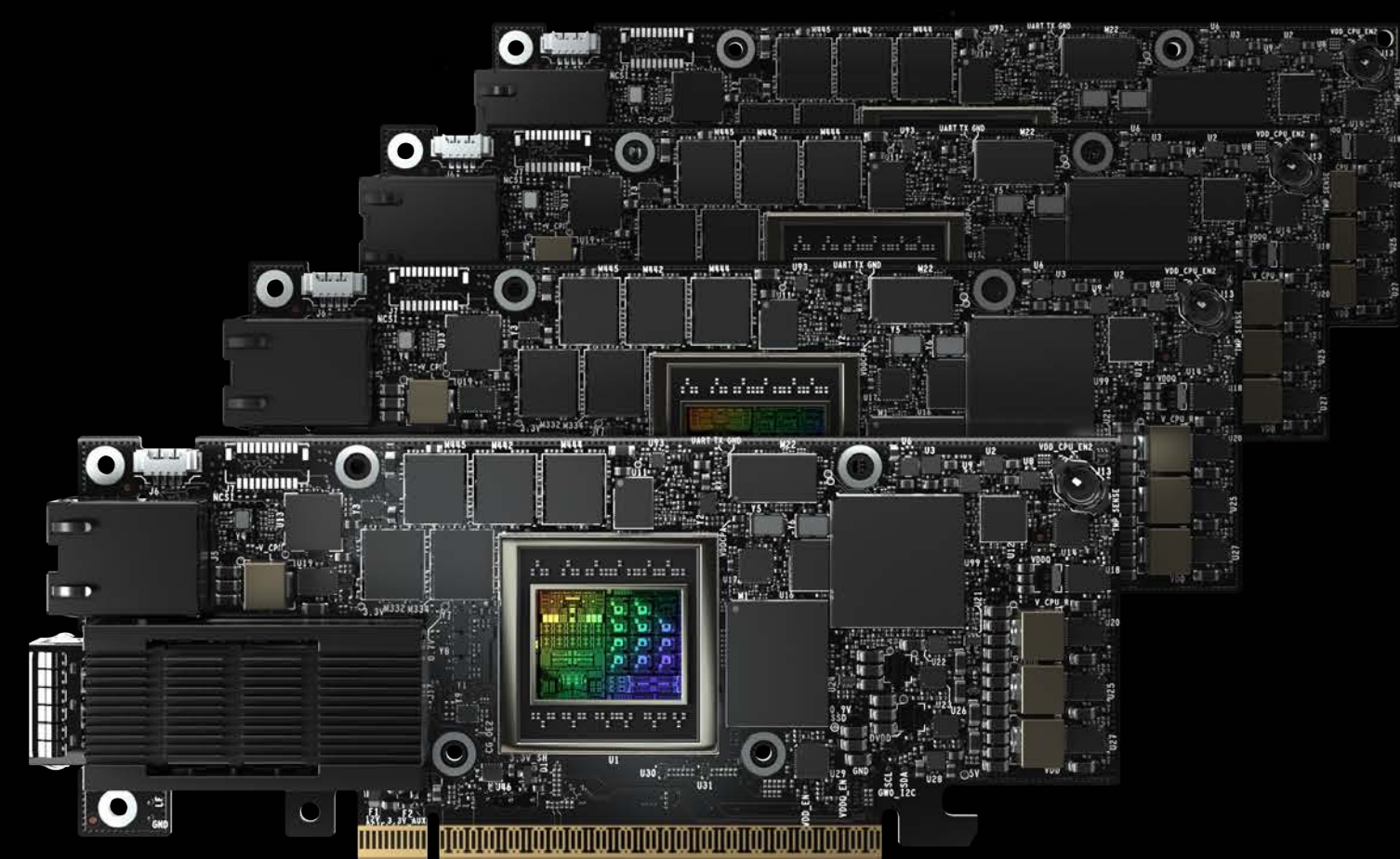Chip

50B Transistors in TSMC 4NP

72-Ports 400G SerDes
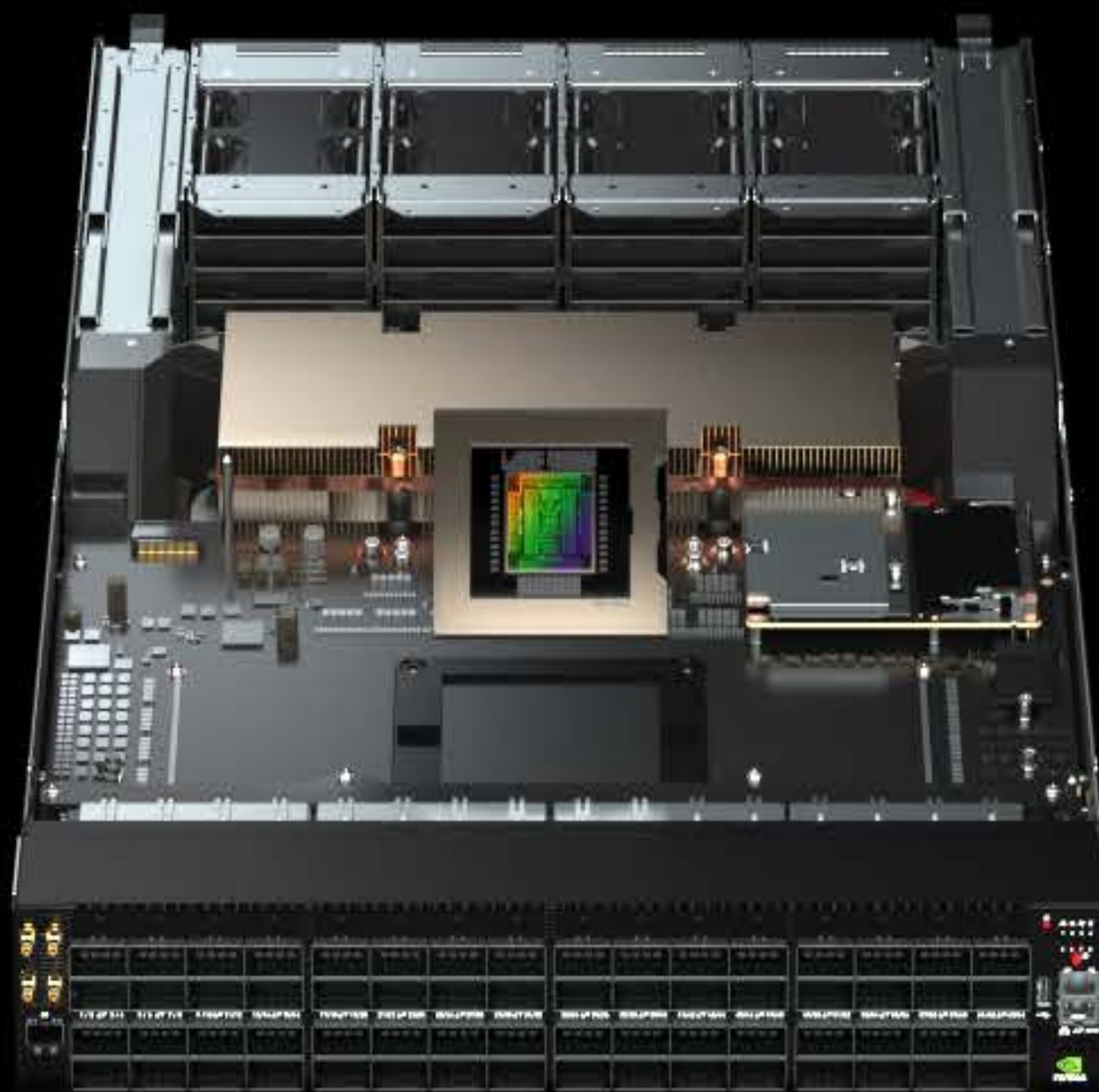
4 NVLinks at 1.8TB/sec

7.2TB/sec Full-Duplex Bandwidth

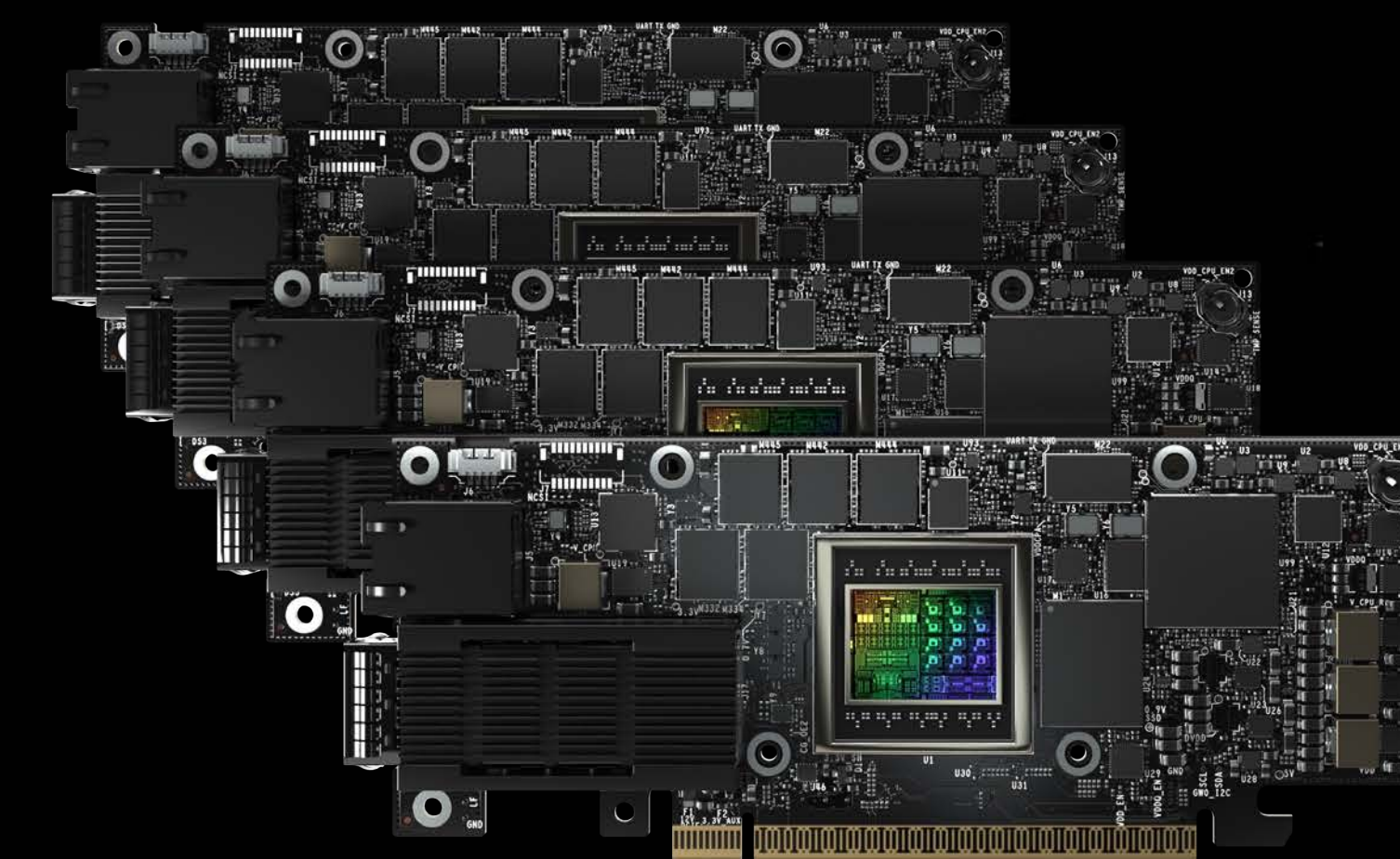SHARP In-Network Compute – 3.6 TFLOPS FP8

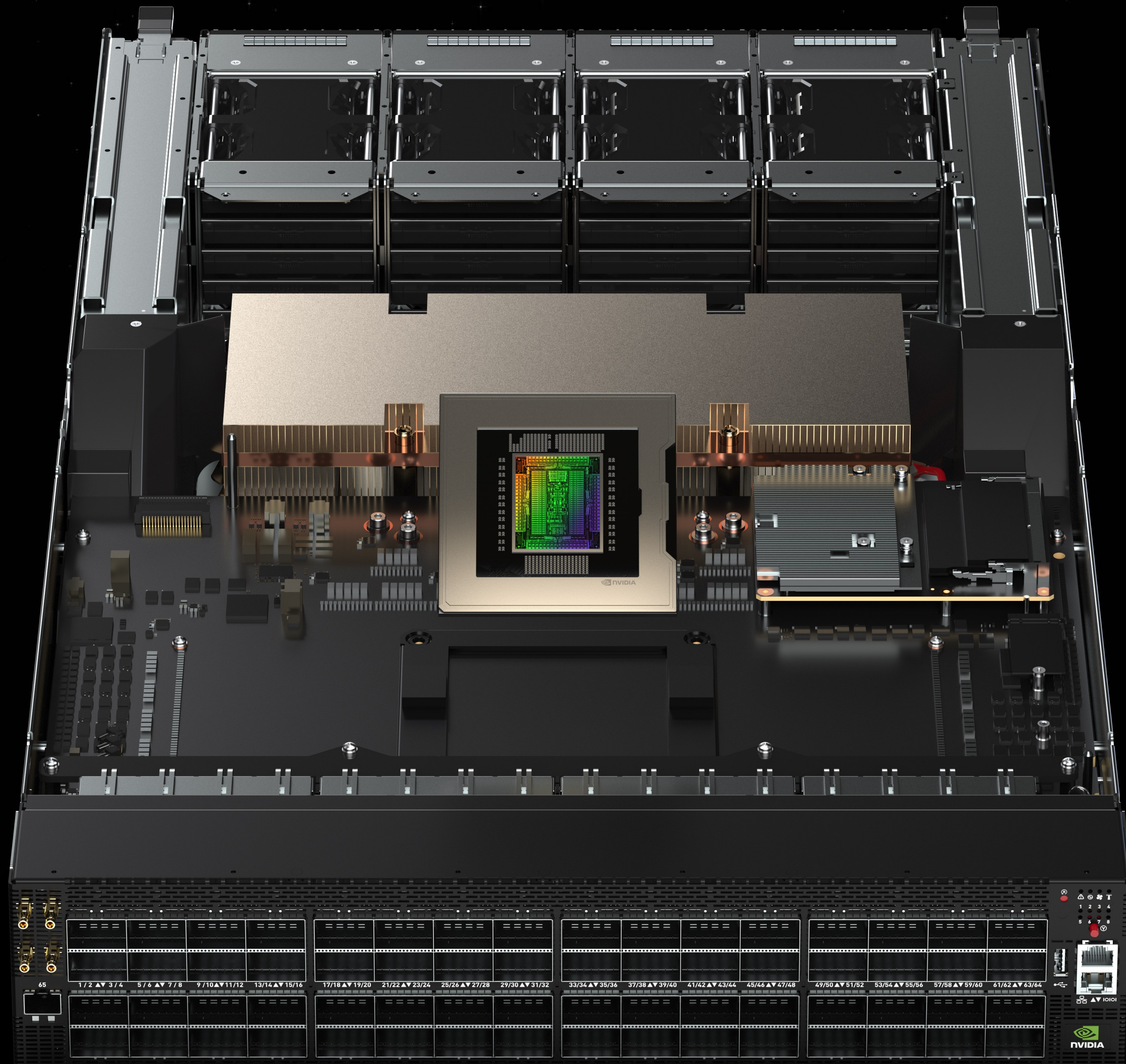BlueField-3
400G SuperNIC

Spectrum-X800
Ethernet Switch
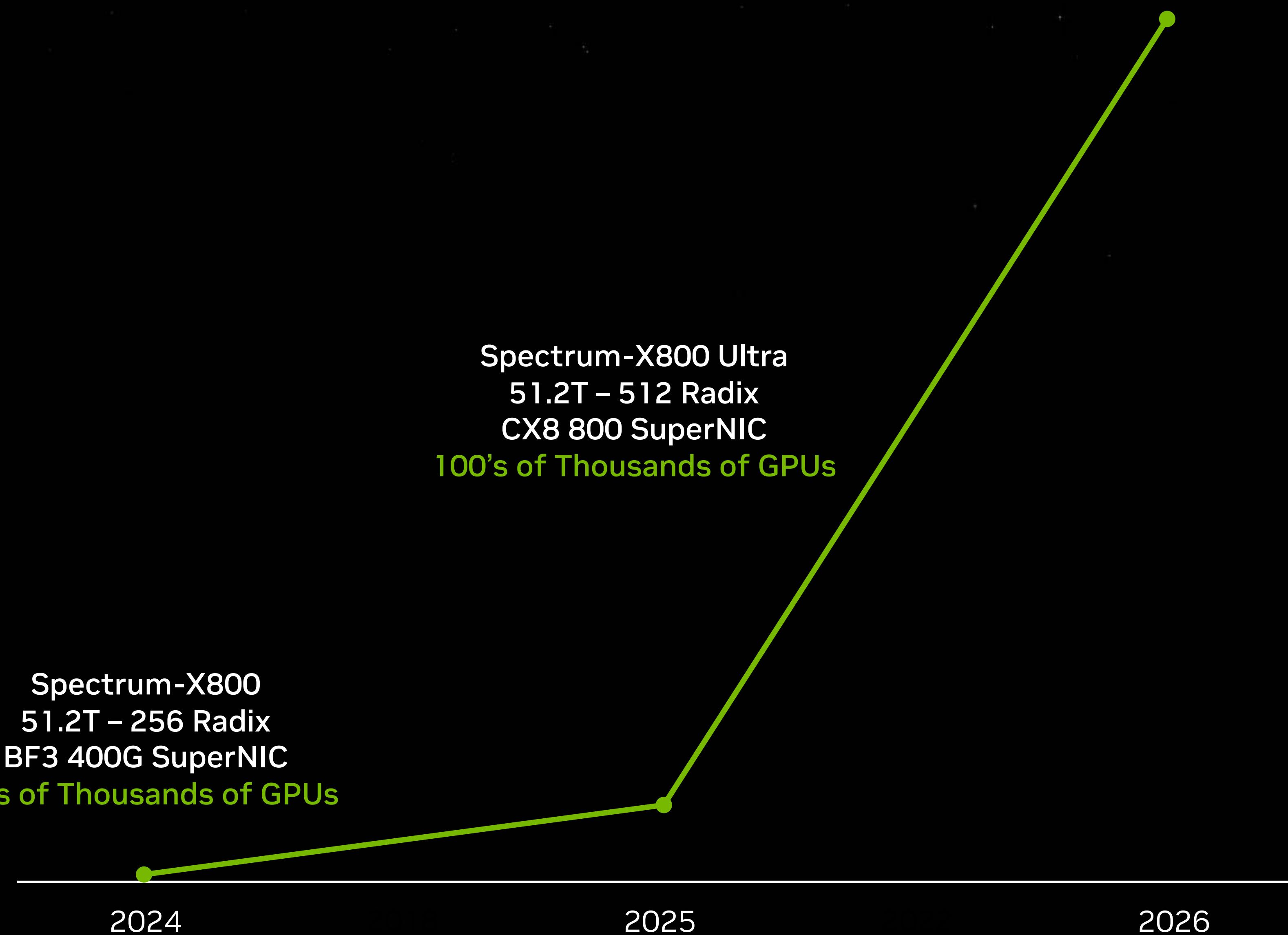
1.6X

Effective
Bandwidth

NVIDIA SPECTRUM-X
為生成式人工智慧推進乙太網路
Supercharging Ethernet for Generative AI

Spectrum-X1600
102.4T – 512 Radix
CX9 1600G SuperNIC
Millions of GPUs

Spectrum-X800 Ultra
51.2T – 512 Radix
CX8 800 SuperNIC
100's of Thousands of GPUs

Spectrum-X800
51.2T – 256 Radix
BF3 400G SuperNIC
10's of Thousands of GPUs

2024          2025          2026

年度 SPECTRUM-X 節奏
擴展到數百萬個 GPU

ANNUAL SPECTRUM-X RHYTHM
SCALING TO MILLIONS OF GPUS

aws
Google Cloud
Microsoft Azure
ORACLE CLOUD Infrastructure

ADEPT
AI21 labs
Character.AI
cohere
essential AI
Hugging Face
Inflection

Meta
MISTRAL AI_
OpenAI
perplexity
Recursion
Tesla
together.ai
X

AiVRES
APPLIED DIGITAL
ASRock Rack
ASUS
CISCO
CoreWeave
Crusoe
DELL Technologies

EVIDEN
FOXCONN HON HAI TECHNOLOGY GROUP
FUJITSU
GIGABYTE
Hewlett Packard Enterprise
IBM Cloud
indosat OOREDOO HUTCHISON
Inventec

Lambda
Lenovo
NEXGEN CLOUD
NORTHERN DATA GROUP
PEGATRON
QCT
Scaleway
Singtel

SoftBank
SUPERMICRO
wistron
wiwynn
YOTTA
YTL COMMUNICATIONS
zt Systems

Hopper Platform

GPU

Hopper GPU
6S HBM3

Hopper+ GPU
6S HBM3e

CPU

Grace CPU

NVLINK

NVLink Switch
900 GB/sec

NIC

CX7
SuperNIC

BF3
SuperNIC

SWITCH

Quantum-X400
InfiniBand Switch

資料中心規模 ‧ 一年節奏 ‧ 技術限制 ‧ 一個架構
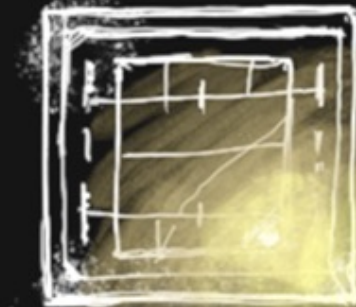DATACENTER SCALE ‧ ONE-YEAR RHYTHM ‧ TECHNOLOGY LIMITS ‧ ONE ARCHITECTURE
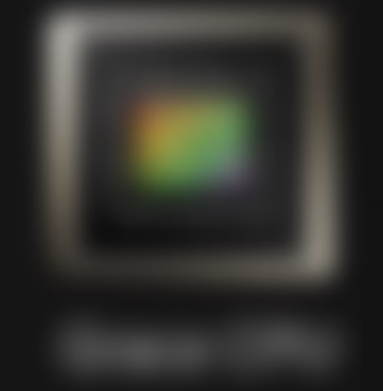
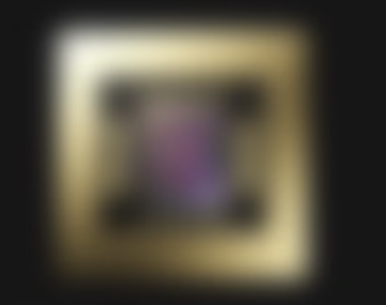**Hopper Platform**

**Blackwell Platform**
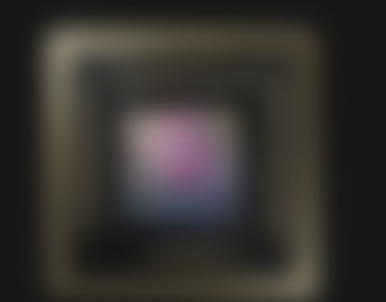
GPU

Blackwell GPU
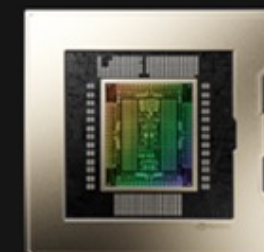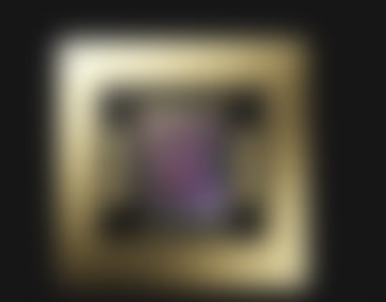8S HBM3e

Blackwell Ultra GPU
8S HBM3e 12H

CPU

NVLINK

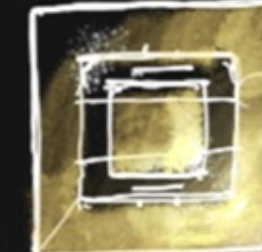NVLink 5 Switch
1800 GB/sec

NIC

CX8 SuperNIC

SWITCH

Spectrum-X800
Ethernet Switch

Quantum-X800
Switch

Spectrum Ultra X800
Ethernet Switch 512-Radix

資料中心規模 · 一年節奏 · 技術限制 · 一個架構

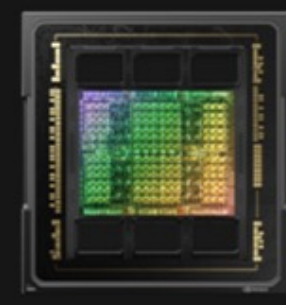DATACENTER SCALE · ONE-YEAR RHYTHM · TECHNOLOGY LIMITS · ONE ARCHITECTURE

GPU CPU NVLINK NIC SWITCH

Hopper Platform
Blackwell Platform
Rubin Platform

Rubin GPU
8S HBM4

Rubin Ultra GPU
12S HBM4

Vera CPU

NVLink 6 Switch
3600 GB/sec

CX9 SuperNIC
1600 Gb/sec

X1600
IB/Ethernet Switch

資料中心規模 ・ 一年節奏 ・ 技術限制 ・ 一個架構
DATACENTER SCALE ・ ONE-YEAR RHYTHM ・ TECHNOLOGY LIMITS ・ ONE ARCHITECTURE

Hopper Platform
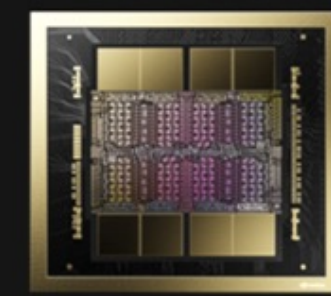
Blackwell Platform

Rubin Platform

GPU

CPU

NVLINK

NIC
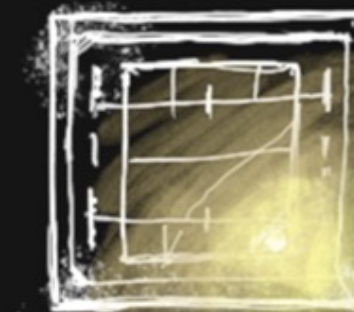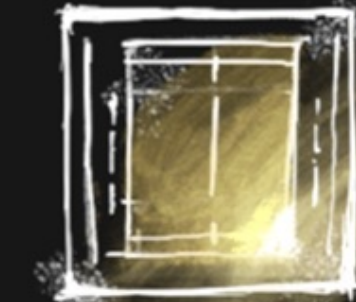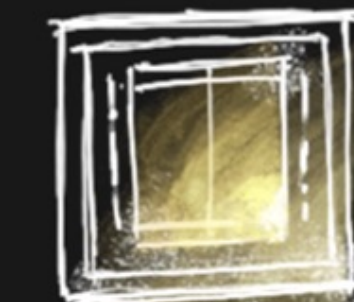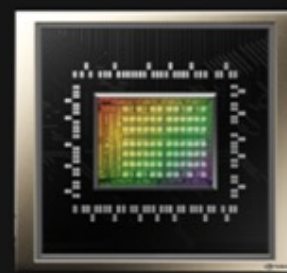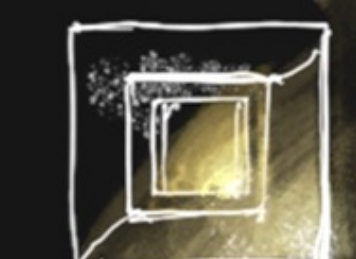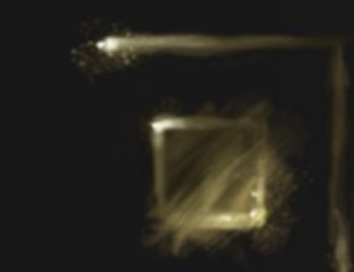
SWITCH

Hopper GPU
6S HBM3

Hopper+ GPU
6S HBM3e

Grace CPU

NVLink Switch
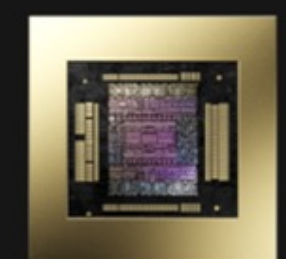900 GB/sec

CX7
SuperNIC

BF3
SuperNIC

Quantum-X400
Infiniband Switch

Blackwell GPU
8S HBM3e

Blackwell Ultra GPU
8S HBM3e 12H

NVLink 5 Switch
1800 GB/sec

CX8 SuperNIC

Spectrum-X800
Ethernet Switch

Quantum-X800
Switch

Spectrum Ultra X800
Ethernet Switch 512-Radix

Rubin GPU
8S HBM4

Rubin Ultra GPU
12S HBM4

Vera CPU

NVLink 6 Switch
3600 GB/sec

CX9 SuperNIC
1600 Gb/sec

X1600
IB/Ethernet Switch

2022          2023          2024          2025          2026          2027

資料中心規模 · 一年節奏 · 技術限制 · 一個架構

DATACENTER SCALE · ONE-YEAR RHYTHM · TECHNOLOGY LIMITS · ONE ARCHITECTURE
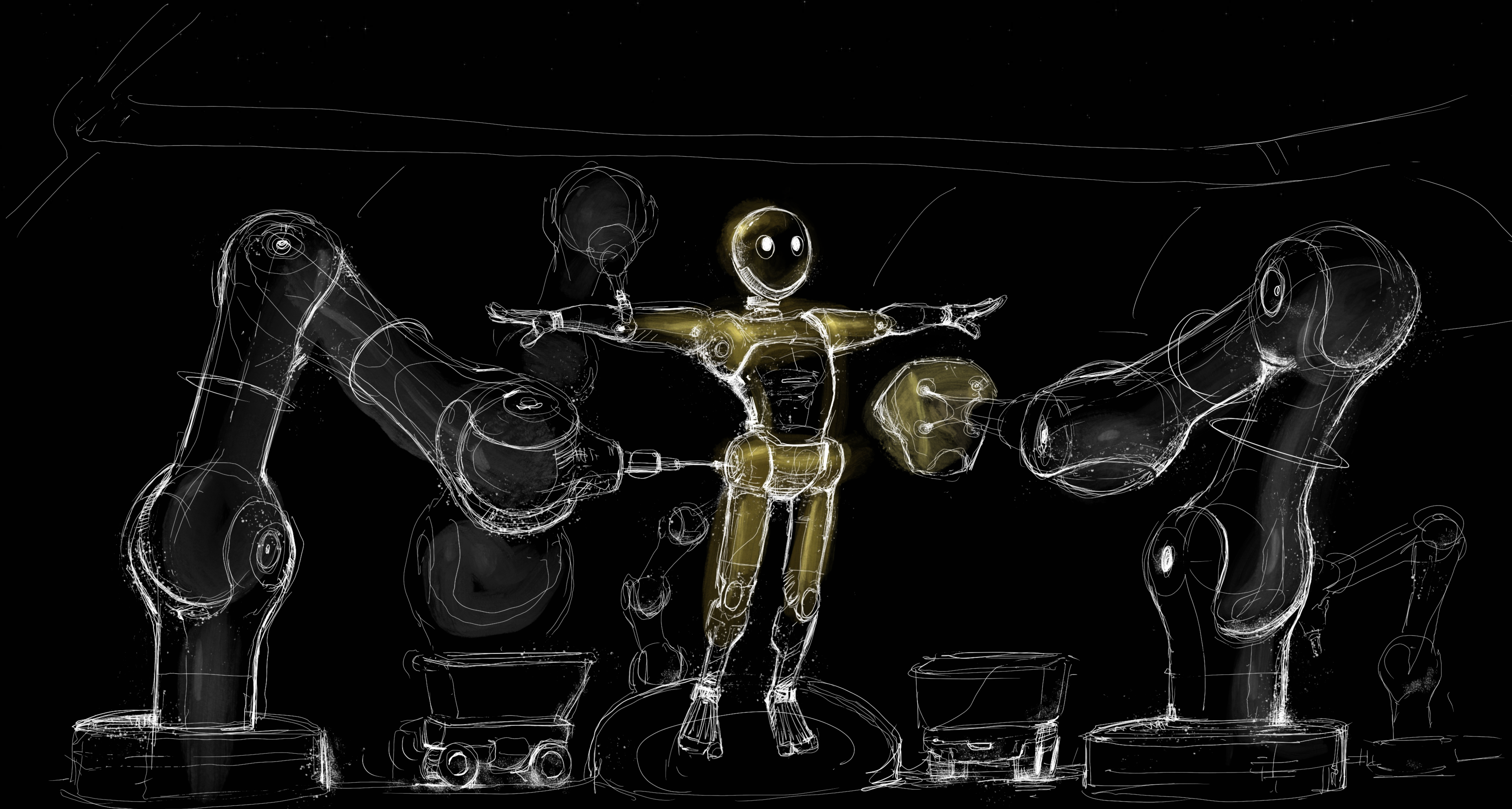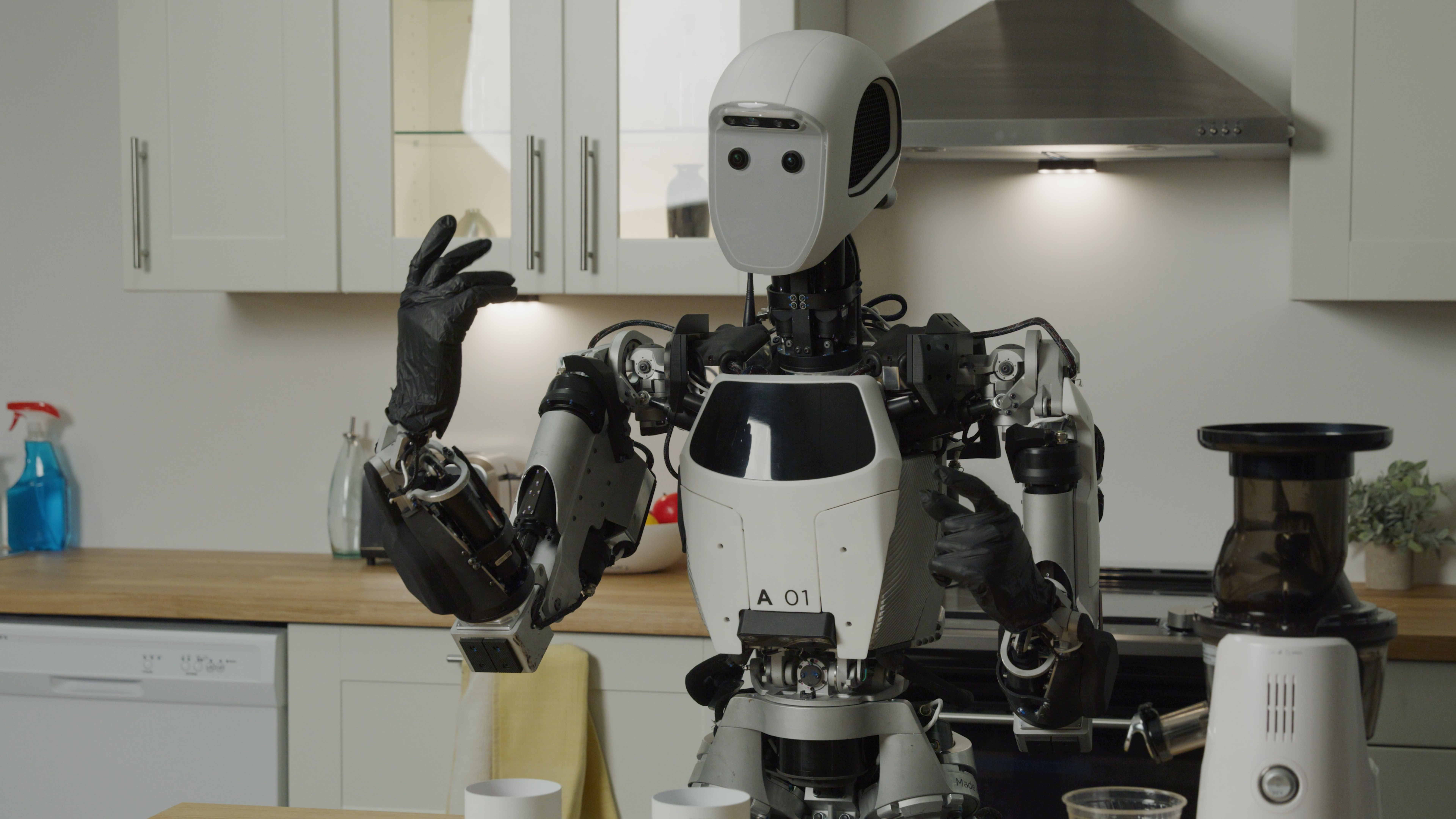
HGX B200

NVLink Switch

GB200 Superchip
Compute Node

Quantum-X800 Switch
ConnectX-8 SuperNIC

Spectrum-X800 Switch
BlueField-3 SuperNIC

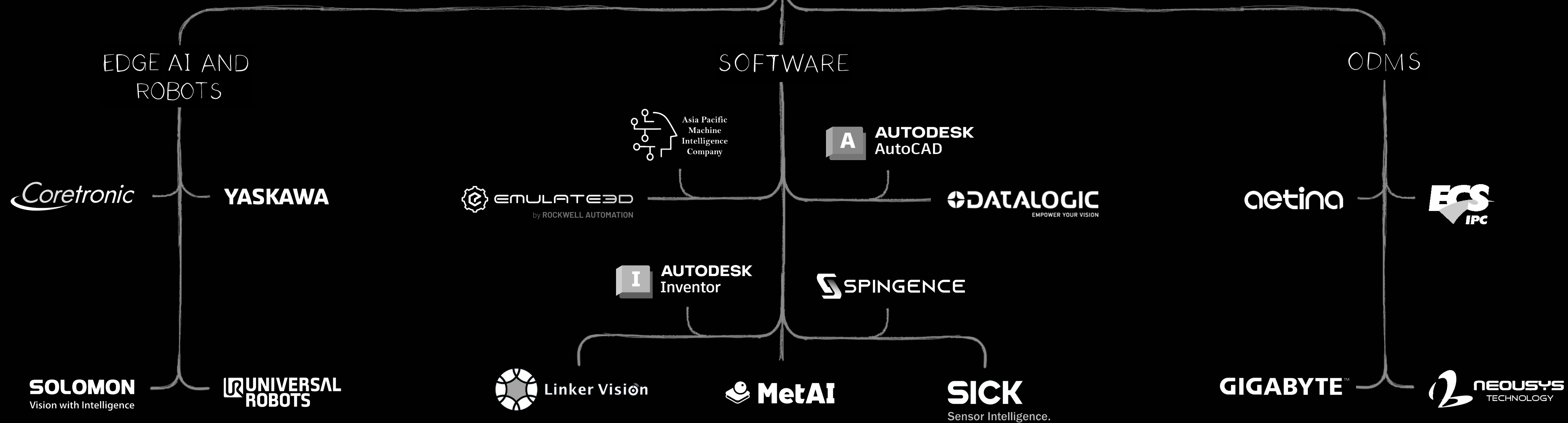NVIDIA BLACKWELL 平台

NVIDIA BLACKWELL PLATFORM

# WAREHOUSE ECOSYSTEM

**GIANT GROUP**

INTEGRATOR **KENMEC**

EDGE AI AND ROBOTS

SOFTWARE

ODMS

Asia Pacific Machine Intelligence Company

**AUTODESK** AutoCAD

*Coretronic* **YASKAWA**

**EMULATE3D** by ROCKWELL AUTOMATION

**DATALOGIC** EMPOWER YOUR VISION

aetina

**ECS** IPC

**AUTODESK** Inventor

**SPINGENCE**

**SOLOMON** Vision with Intelligence

**UNIVERSAL ROBOTS**

Linker Vision

**MetAI**

**SICK** Sensor Intelligence.

**GIGABYTE**

**neousys** TECHNOLOGY

ISAAC   DRIVE

HOLOSCAN   METROPOLIS

NVIDIA OMNIVERSE

# FACTORY ECOSYSTEM

**Foxconn Industrial Internet**

INTEGRATOR — **FOXCONN**®

## EDGE AI AND ROBOTS

ASM    EPSON®

FARobot®    JUSDA

kurtz ersa    TM ROBOT

TRI innovation    ViTrox

## SOFTWARE

Ansys    ASPEED

AUTODESK    cādence®

creo®    FlexSim An AUTODESK Company

SIEMENS    SketchUp    VISUAL COMPONENTS

## ODMS

Innoconn    Ingrasys®

SIEMENS

---

ISAAC    DRIVE    HOLOSCAN    METROPOLIS

**NVIDIA OMNIVERSE**

wistron®

FANUC Robot CRX-20iA/L

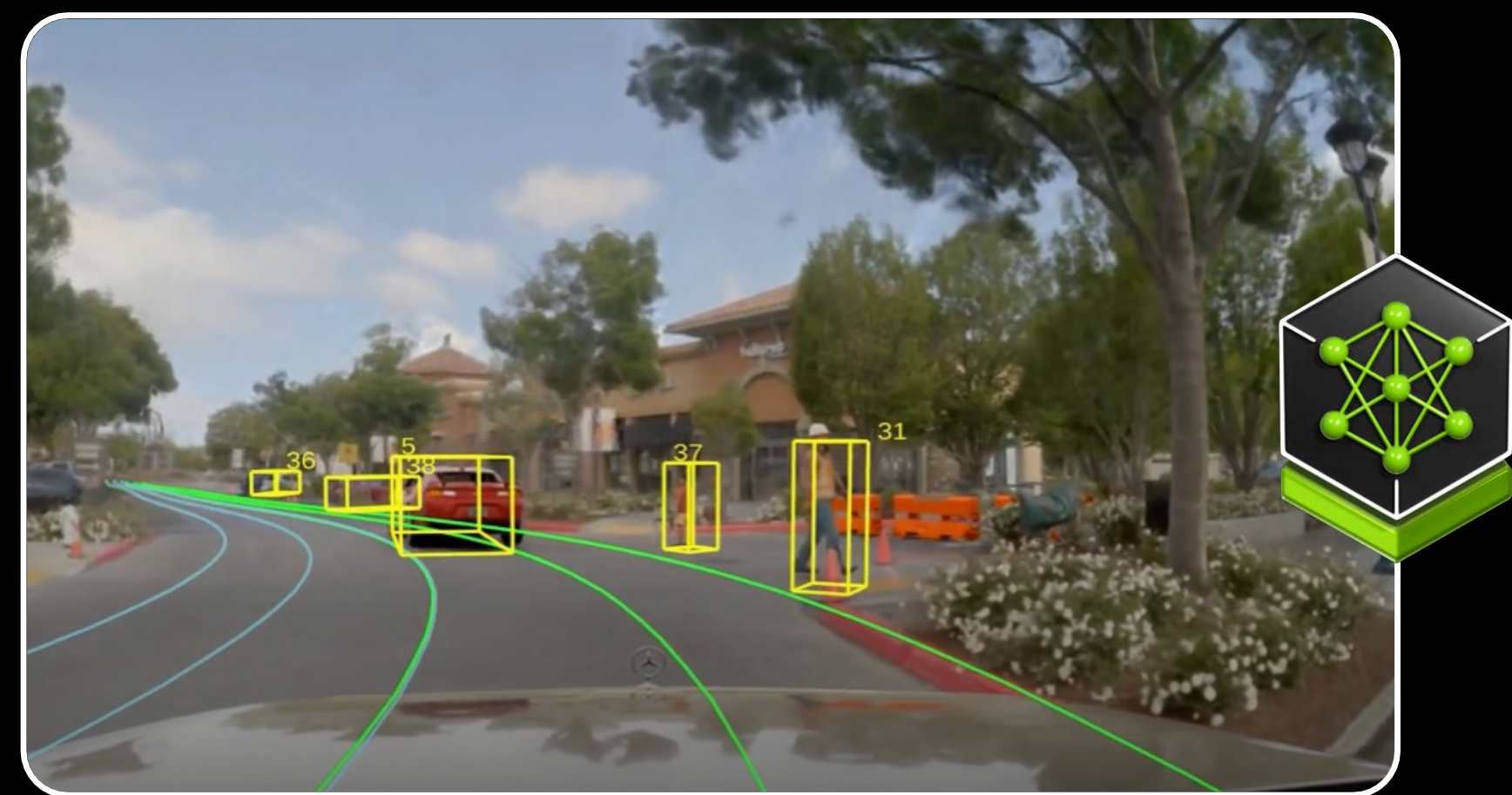NVIDIA Omniverse

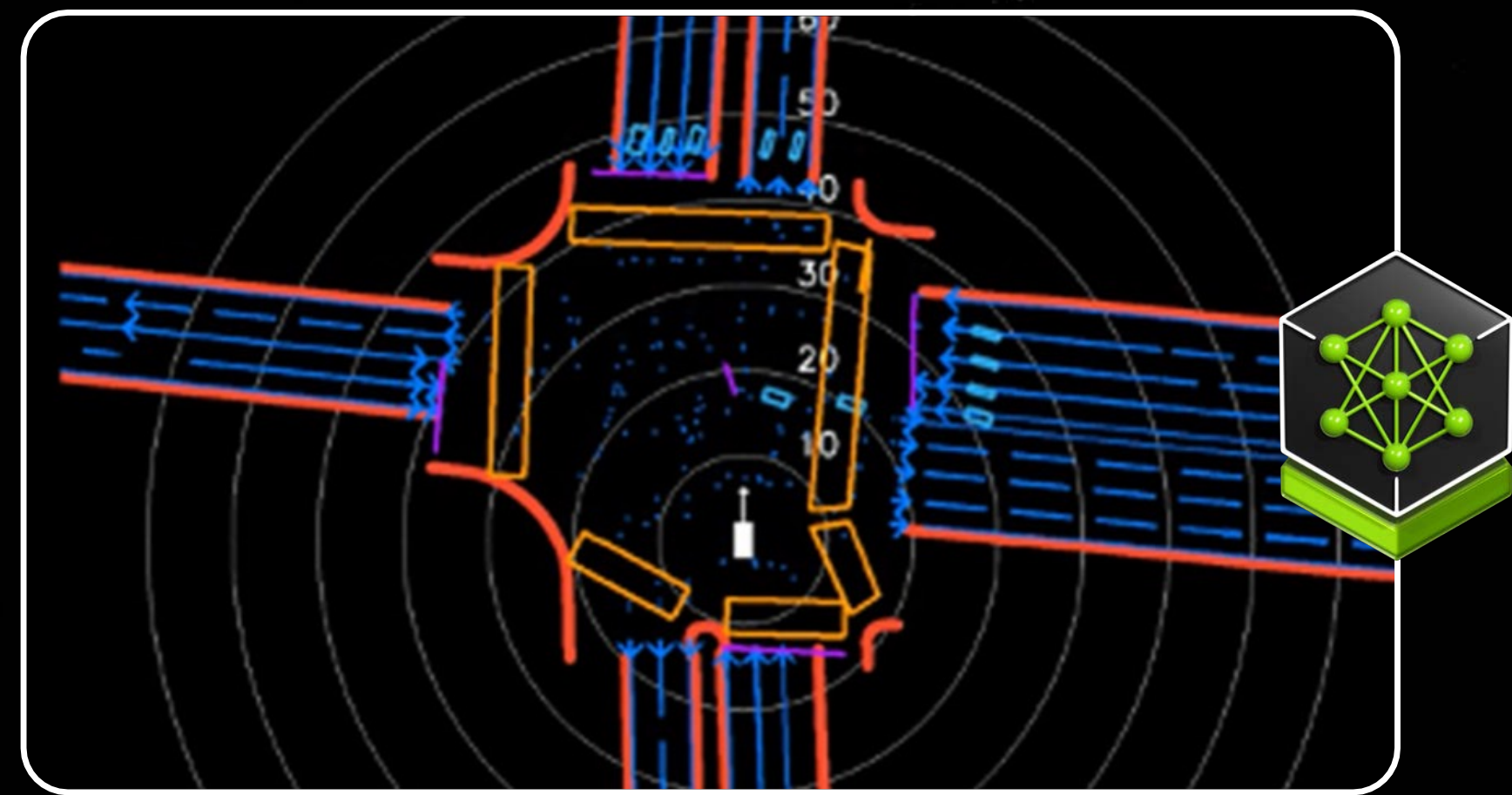NVIDIA Isaac Manipulator
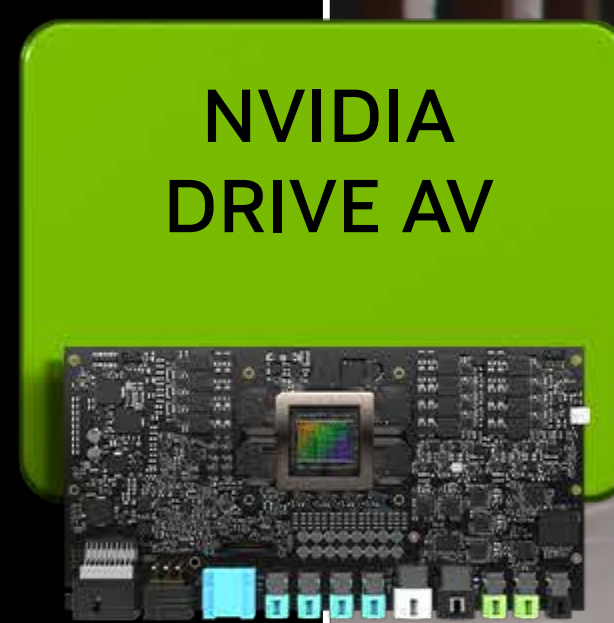
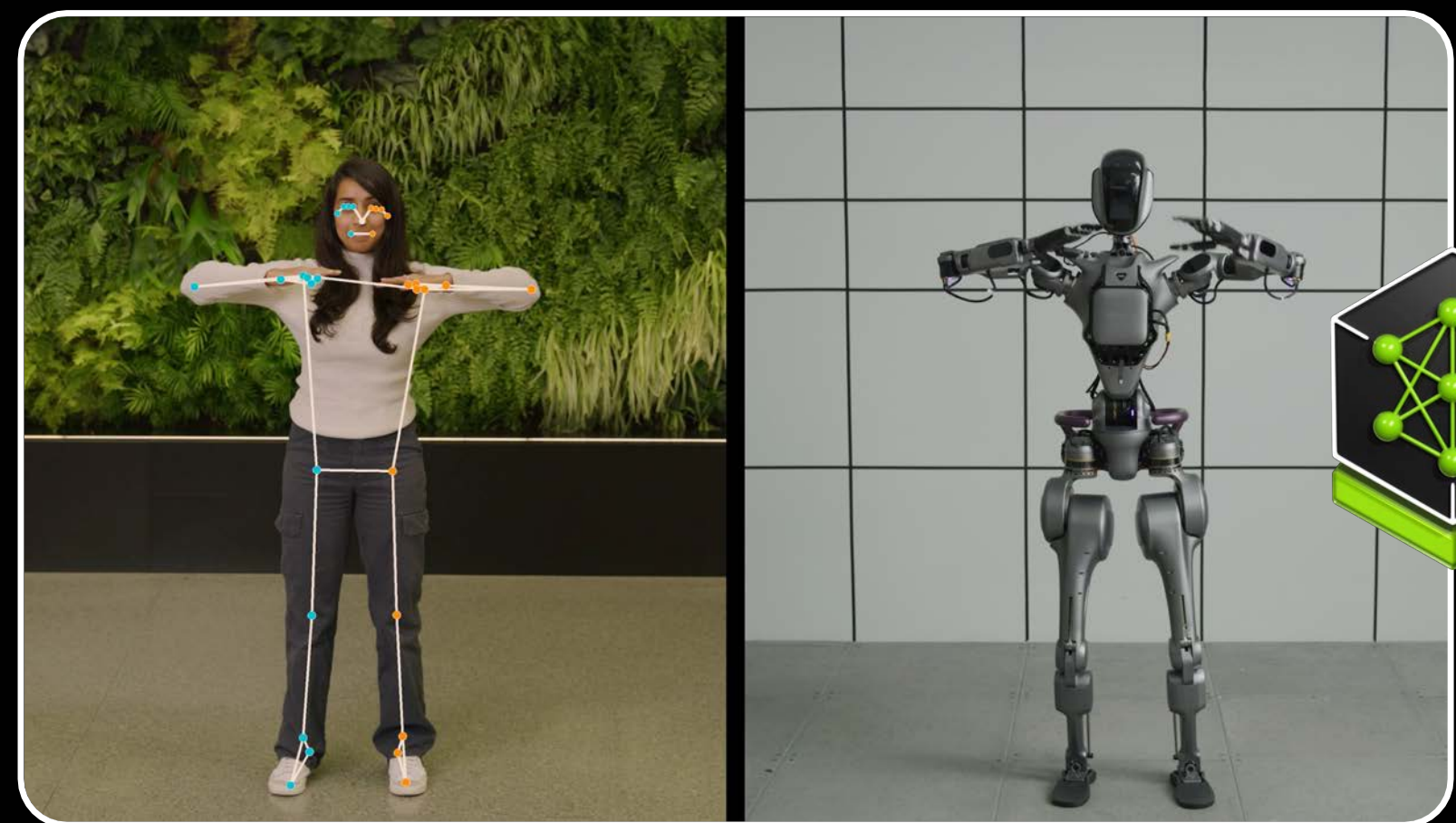NVIDIA AI

Isaac ROS 3.0 Available Now on Github
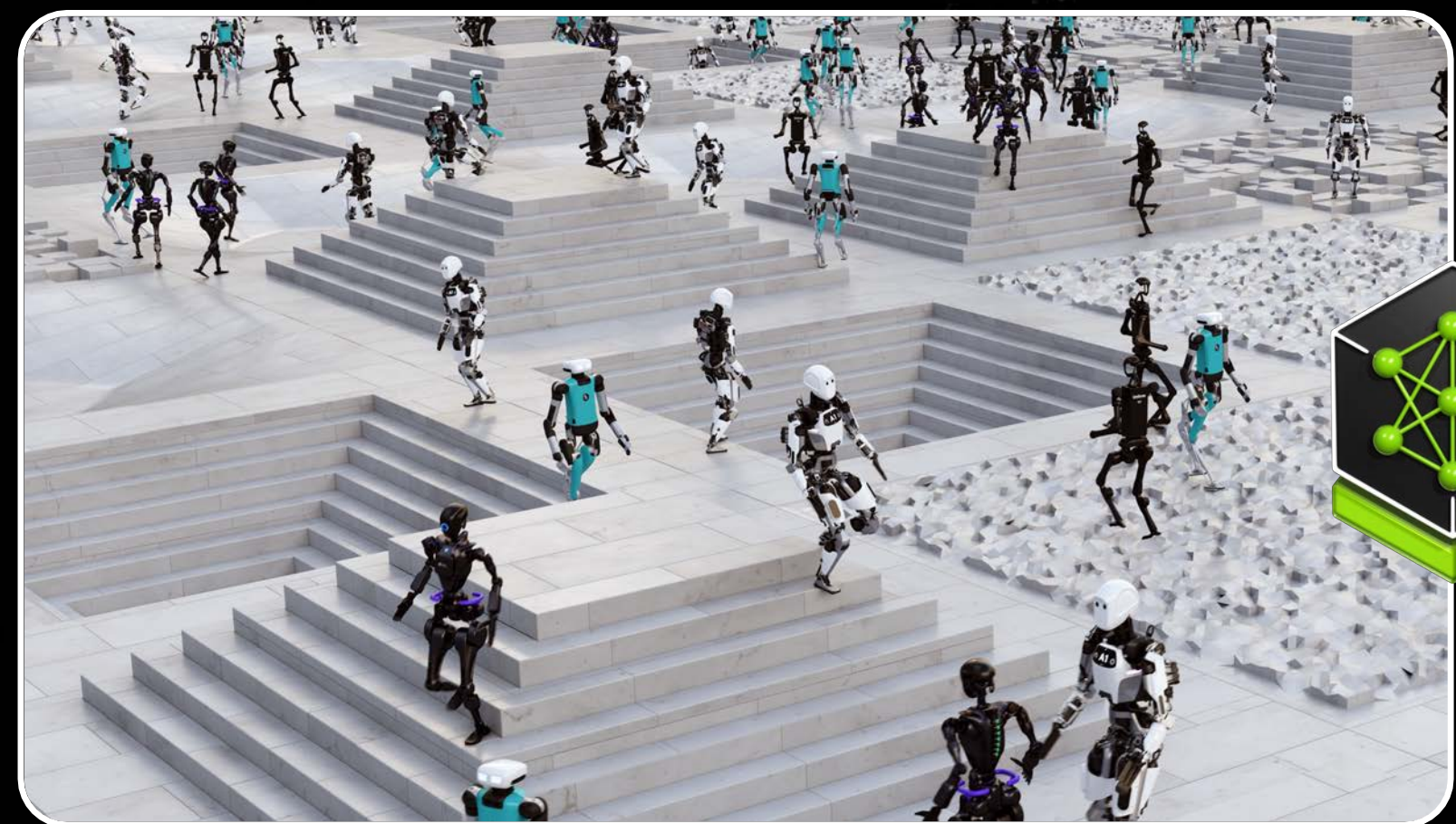
ArcBest

ArcBest Vaux
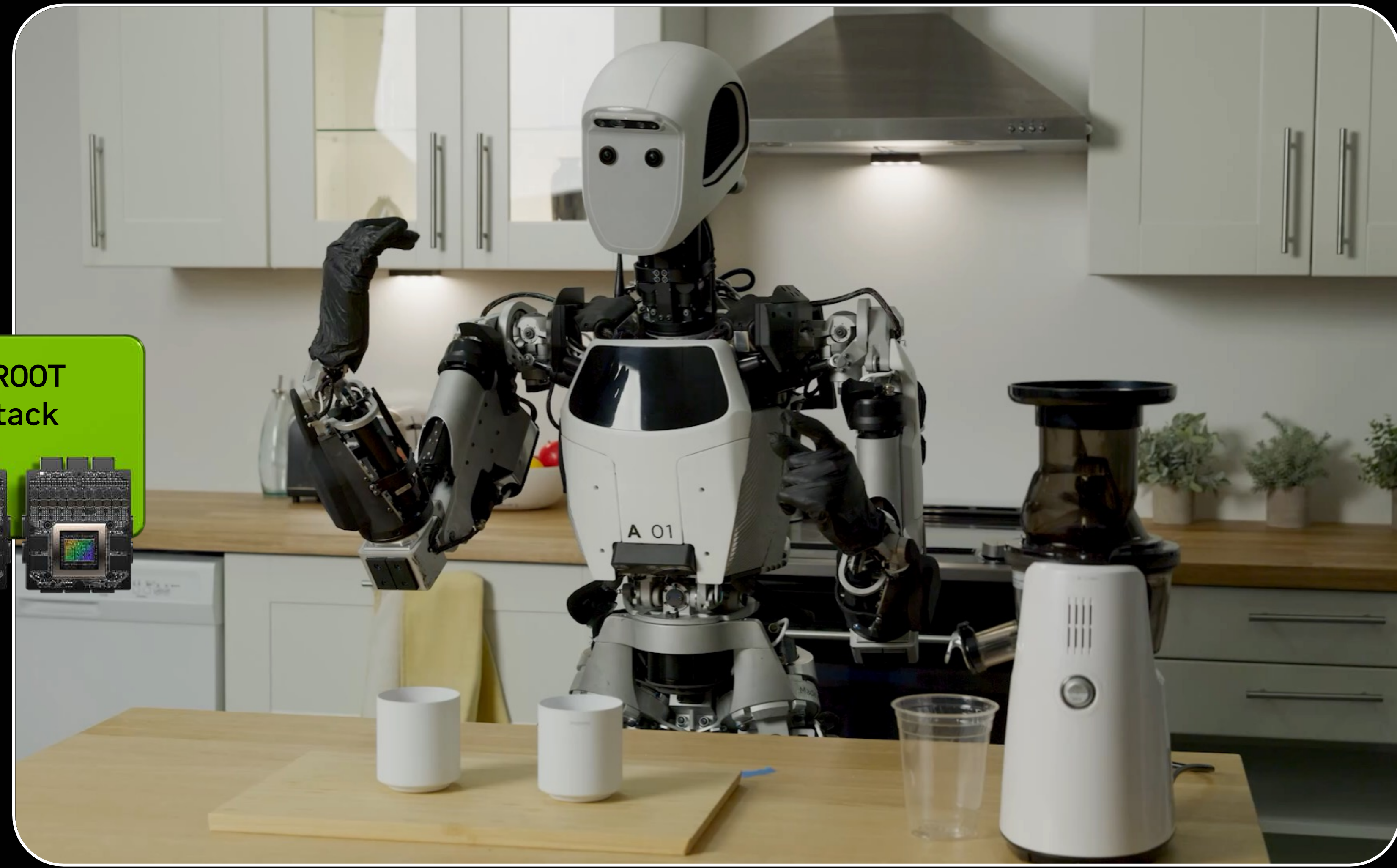
NVIDIA Isaac Perceptor

NVIDIA Omniverse

NVIDIA AI

NVIDIA Omniverse

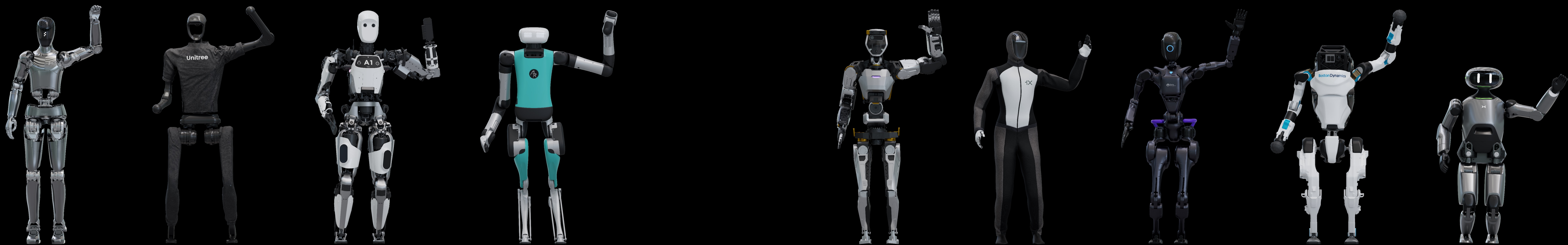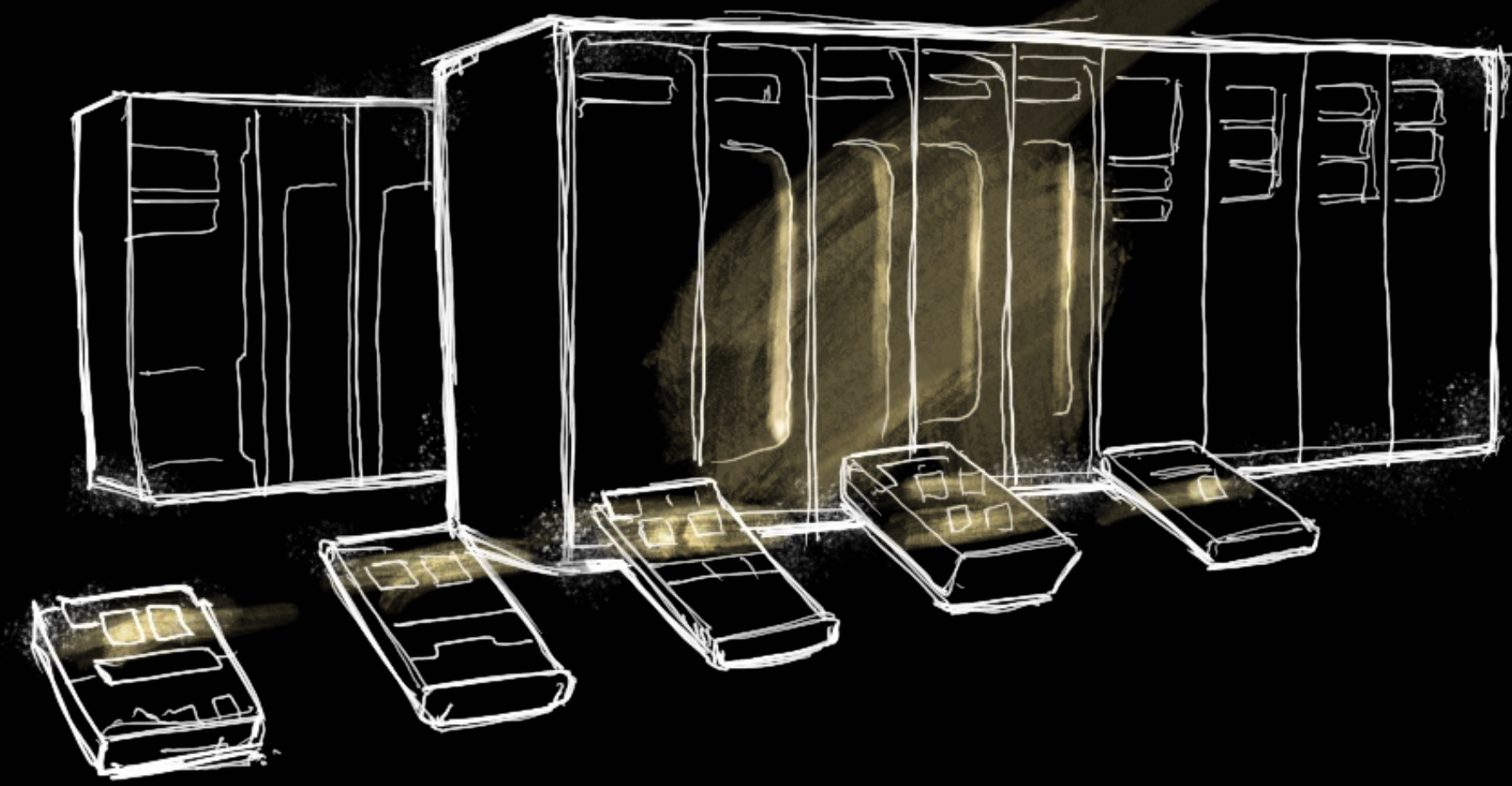NVIDIA AI

NVIDIA DRIVE AV

NVIDIA Omniverse

NVIDIA AI

GR00T
Stack

AI FACTORY

BLACKWELL

NIMS

OMNIVERSE/ ROBOTICS

新產業革命

"A NEW INDUSTRIAL REVOLUTION"