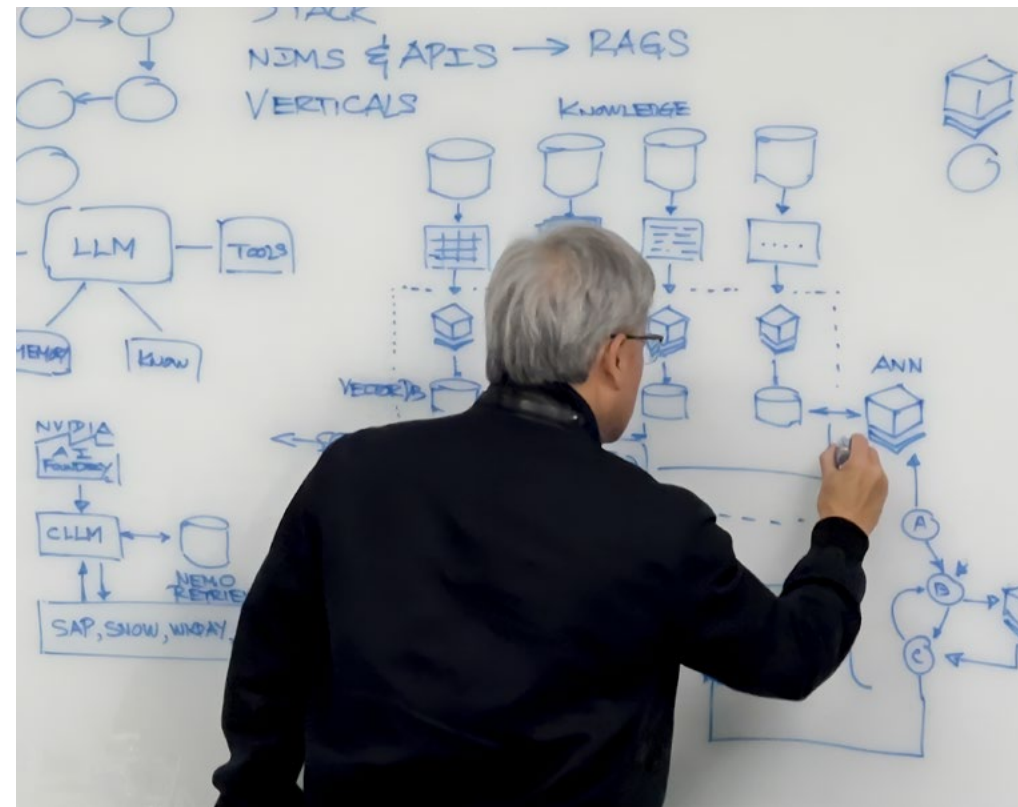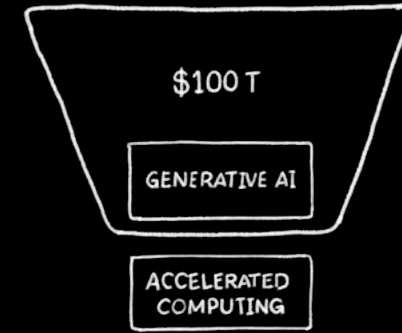**NVIDIA**

**GTC 2024**

Highlights

"GTC is a front-row seat to what's happening in AI."
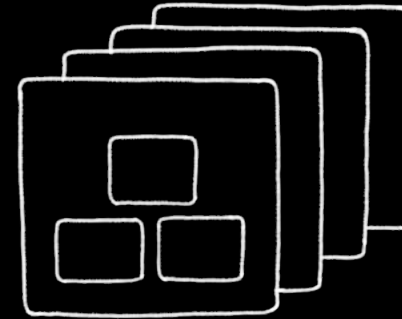
*Bloomberg*
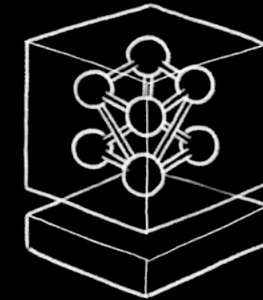
**A New Industrial Revolution**

$100 T

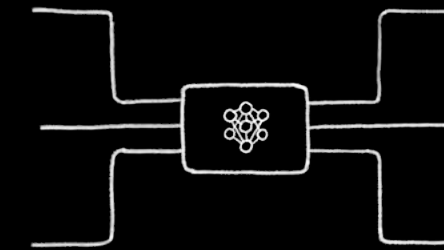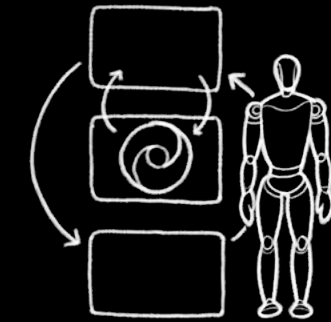GENERATIVE AI

ACCELERATED COMPUTING

NEW INDUSTRY

BLACKWELL PLATFORM
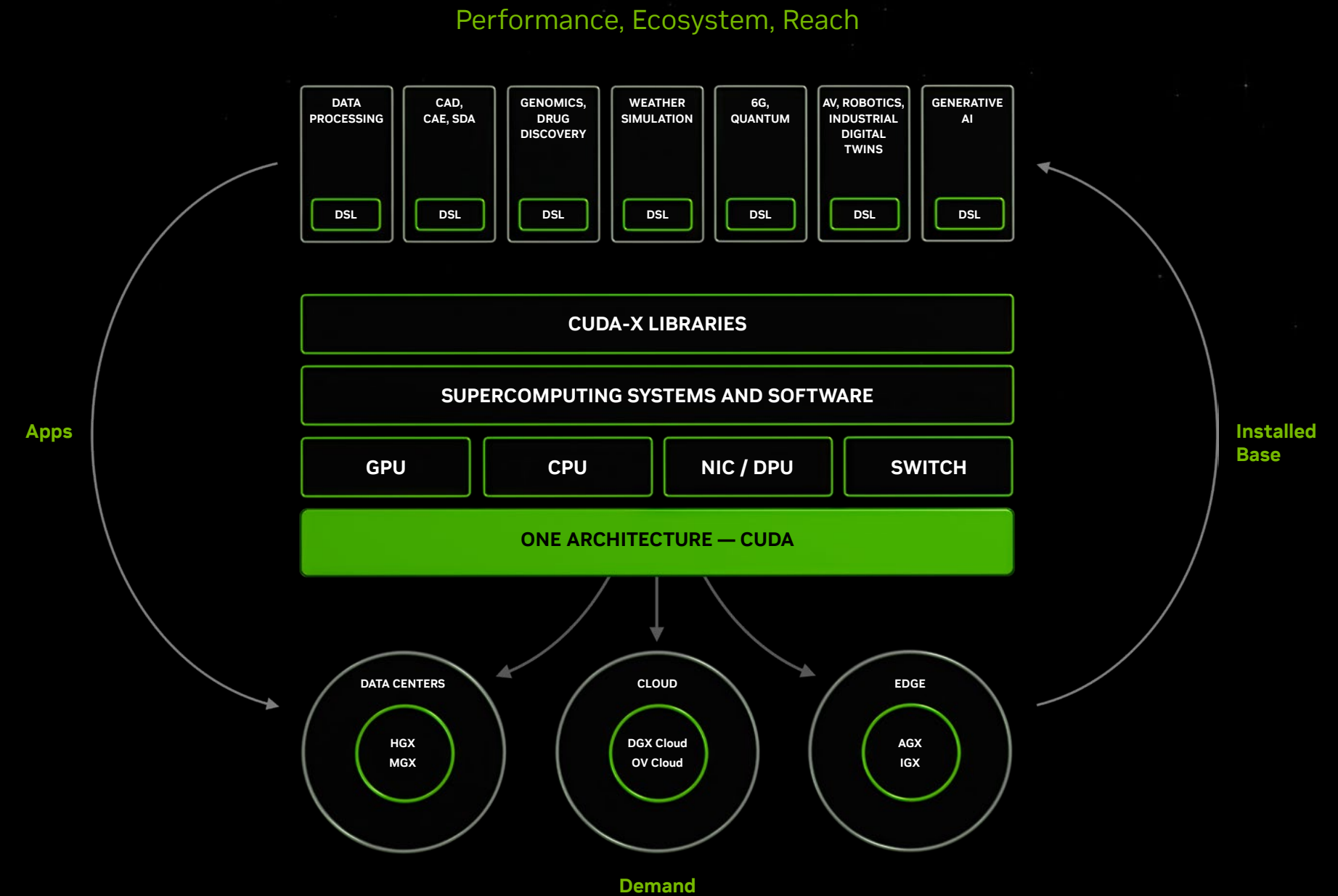
NIMs

NEMO AND NVIDIA AI FOUNDRY

OMNIVERSE AND ISAAC ROBOTICS

# "NVIDIA's moat is its software and ecosystem."

*The Edge Singapore*

The global NVIDIA ecosystem now approaches 5 million developers—capping a record-setting year for new developers. Forty thousand companies have worked with NVIDIA. There are now more than 3,300 GPU-accelerated applications. More than 1,600 generative AI companies are building on NVIDIA.

**Performance, Ecosystem, Reach**

| DATA PROCESSING | CAD, CAE, SDA | GENOMICS, DRUG DISCOVERY | WEATHER SIMULATION | 6G, QUANTUM | AV, ROBOTICS, INDUSTRIAL DIGITAL TWINS | GENERATIVE AI |
|---|---|---|---|---|---|---|
| DSL | DSL | DSL | DSL | DSL | DSL | DSL |

**CUDA-X LIBRARIES**

**SUPERCOMPUTING SYSTEMS AND SOFTWARE**

| GPU | CPU | NIC / DPU | SWITCH |
|---|---|---|---|

**ONE ARCHITECTURE — CUDA**

Apps

Installed Base

**DATA CENTERS**
HGX
MGX

**CLOUD**
DGX Cloud
OV Cloud

**EDGE**
AGX
IGX

Demand

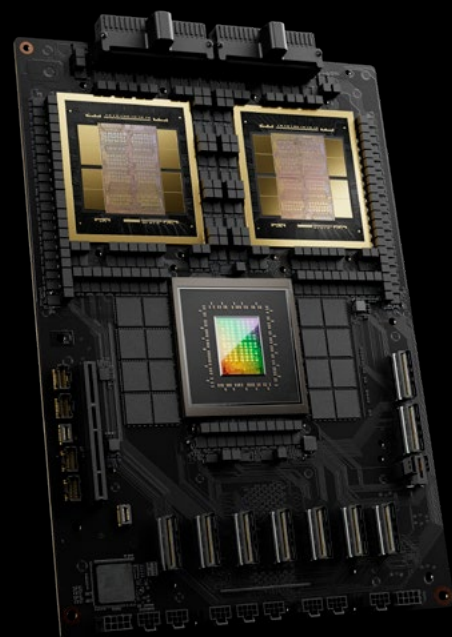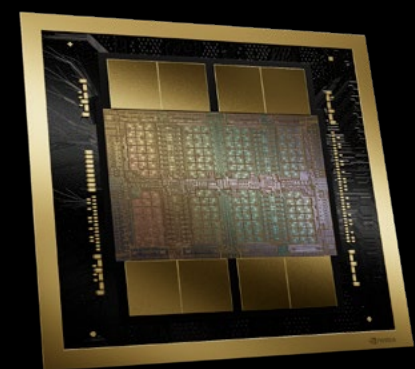"NVIDIA's new Blackwell chip is key to the next stage of AI."

*Bloomberg*

The NVIDIA Blackwell platform has arrived—enabling organizations everywhere to build and run real-time generative AI on trillion-parameter large language models at up to 25X less cost and energy consumption than its predecessor.

# "Meet Blackwell, the new GPU for the AI era."

*Engineering.com*

## NVIDIA Blackwell Platform Arrives to Power a New Era of Computing
March 18, 2024

GTC—Powering a new era of computing, NVIDIA today announced that the NVIDIA Blackwell platform has arrived — enabling organizations everywhere to build and run real-time generative AI on trillion-parameter large language models at up to 25x less cost and energy consumption than its predecessor.

The Blackwell GPU architecture features six transformative technologies for accelerated computing, which will help unlock breakthroughs in data processing, engineering simulation, electronic design automation, computer-aided drug design, quantum computing and generative AI — all emerging industry opportunities for NVIDIA.

"For three decades we've pursued accelerated computing, with the goal of enabling transformative breakthroughs like deep learning and AI," said **Jensen Huang, founder and CEO of NVIDIA.** "Generative AI is the defining technology of our time. Blackwell is the engine to power this new industrial revolution. Working with the most dynamic companies in the world, we will realize the promise of AI for every industry."

Among the many organizations expected to adopt Blackwell are Amazon Web Services, Dell Technologies, Google, Meta, Microsoft, OpenAI, Oracle, Tesla and xAI.

**Sundar Pichai, CEO of Alphabet and Google:** "Scaling services like Search and Gmail to billions of users has taught us a lot about managing compute infrastructure. As we enter the AI platform shift, we continue to invest deeply in infrastructure for our own products and services, and for our Cloud customers. We are fortunate to have a longstanding partnership with NVIDIA, and look forward to bringing the breakthrough capabilities of the Blackwell GPU to our Cloud customers and teams across Google, including Google DeepMind, to accelerate future discoveries."

**Andy Jassy, president and CEO of Amazon:** "Our deep collaboration with NVIDIA goes back more than 13 years, when we launched the world's first GPU cloud instance on AWS. Today we offer the widest range of GPU solutions available anywhere in the cloud, supporting the world's most technologically advanced accelerated workloads. It's why the new NVIDIA Blackwell GPU will run so well on AWS and the reason that NVIDIA chose AWS to co-develop Project Ceiba, combining NVIDIA's next-generation Grace Blackwell Superchips with the AWS Nitro System's advanced virtualization and ultra-fast Elastic Fabric Adapter networking, for NVIDIA's own AI research and development. Through this joint effort between AWS and NVIDIA engineers, we're continuing to innovate together to make AWS the best place for anyone to run NVIDIA GPUs in the cloud."

**Michael Dell, founder and CEO of Dell Technologies:** "Generative AI is critical to creating smarter, more reliable and efficient systems. Dell Technologies and NVIDIA are working together to shape the future of technology. With the launch of Blackwell, we will continue to deliver the next-generation of accelerated products and services to our customers, providing them with the tools they need to drive innovation across industries."

**Demis Hassabis, cofounder and CEO of Google DeepMind:** "The transformative potential of AI is incredible, and it will help us solve some of the world's most important scientific problems. Blackwell's breakthrough technological capabilities will provide the critical compute needed to help the world's brightest minds chart new scientific discoveries."

**Mark Zuckerberg, founder and CEO of Meta:** "AI already powers everything from our large language models to our content recommendations, ads, and safety systems, and it's only going to get more important in the future. We're looking forward to using NVIDIA's Blackwell to help train our open-source Llama models and build the next generation of Meta AI and consumer products."

**Satya Nadella, executive chairman and CEO of Microsoft:** "We are committed to offering our customers the most advanced infrastructure to power their AI workloads. By bringing the GB200 Grace Blackwell processor to our datacenters globally, we are building on our long-standing history of optimizing NVIDIA GPUs for our cloud, as we make the promise of AI real for organizations everywhere."

**Sam Altman, CEO of OpenAI:** "Blackwell offers massive performance leaps, and will accelerate our ability to deliver leading-edge models. We're excited to continue working with NVIDIA to enhance AI compute."

**Larry Ellison, chairman and CTO of Oracle:** "Oracle's close collaboration with NVIDIA will enable qualitative and quantitative breakthroughs in AI, machine learning and data analytics. In order for customers to uncover more actionable insights, an even more powerful engine like Blackwell is needed, which is purpose-built for accelerated computing and generative AI."

**Elon Musk, CEO of Tesla and xAI:** "There is currently nothing better than NVIDIA hardware for AI."

Named in honor of David Harold Blackwell — a mathematician who specialized in game theory and statistics, and the first Black scholar inducted into the National Academy of Sciences — the new architecture succeeds the NVIDIA Hopper™ architecture, launched two years ago.

## [...] to Fuel Accelerated Computing and Generative AI

[...] technologies, which together enable AI training and real-time LLM inference for models scaling up to 10 [...]

**[...]ful Chip** — Packed with 208 billion transistors, Blackwell-architecture GPUs are manufactured using a [...]MC process with two-reticle limit GPU dies connected by 10 TB/second chip-to-chip link into a single, unified [...]

**[...] Transformer Engine** — Fueled by new micro-tensor scaling support and NVIDIA's advanced dynamic range [...]ms integrated into NVIDIA TensorRT™-LLM and NeMo Megatron frameworks, Blackwell will support double [...]del sizes with new 4-bit floating point AI inference capabilities.

**[...]VLink** — To accelerate performance for multitrillion-parameter and mixture-of-experts AI models, the latest [...] NVLink® delivers groundbreaking 1.8TB/s bidirectional throughput per GPU, ensuring seamless high-speed [...]ng up to 576 GPUs for the most complex LLMs.

[...]well-powered GPUs include a dedicated engine for reliability, availability and serviceability. Additionally, the [...]adds capabilities at the chip level to utilize AI-based preventative maintenance to run diagnostics and forecast [...]maximizes system uptime and improves resiliency for massive-scale AI deployments to run uninterrupted for [...]s at a time and to reduce operating costs.

[...]d confidential computing capabilities protect AI models and customer data without compromising performance, [...]native interface encryption protocols, which are critical for privacy-sensitive industries like healthcare and [...]

**[...]ine** — A dedicated decompression engine supports the latest formats, accelerating database queries to deliver [...]nce in data analytics and data science. In the coming years, data processing, on which companies spend tens [...]nnually, will be increasingly GPU-accelerated.

[...] Blackwell Superchip connects two NVIDIA B200 Tensor Core GPUs to the NVIDIA Grace CPU over a 900GB/s [...]chip-to-chip interconnect.

[...]mance, GB200-powered systems can be connected with the NVIDIA Quantum-X800 InfiniBand and Spectrum [...]ance, also announced today, which deliver advanced networking at speeds up to 800Gb/s.

[...]onent of the NVIDIA GB200 NVL72, a multi-node, liquid-cooled, rack-scale system for the most [...]oads. It combines 36 Grace Blackwell Superchips, which include 72 Blackwell GPUs and 36 Grace CPUs [...]generation NVLink. Additionally, GB200 NVL72 includes NVIDIA BlueField®-3 data processing units to enable [...]on, composable storage, zero-trust security and GPU compute elasticity in hyperscale AI clouds. The GB200 [...]0x performance increase compared to the same number of NVIDIA H100 Tensor Core GPUs for LLM [...] reduces cost and energy consumption by up to 25x.

[...]ngle GPU with 1.4 exaflops of AI performance and 30TB of fast memory, and is a building block for the newest [...]

[...]00, a server board that links eight B200 GPUs through NVLink to support x86-based generative AI platforms. [...]working speeds up to 400Gb/s through the NVIDIA Quantum-2 InfiniBand and Spectrum-X Ethernet networking [...]
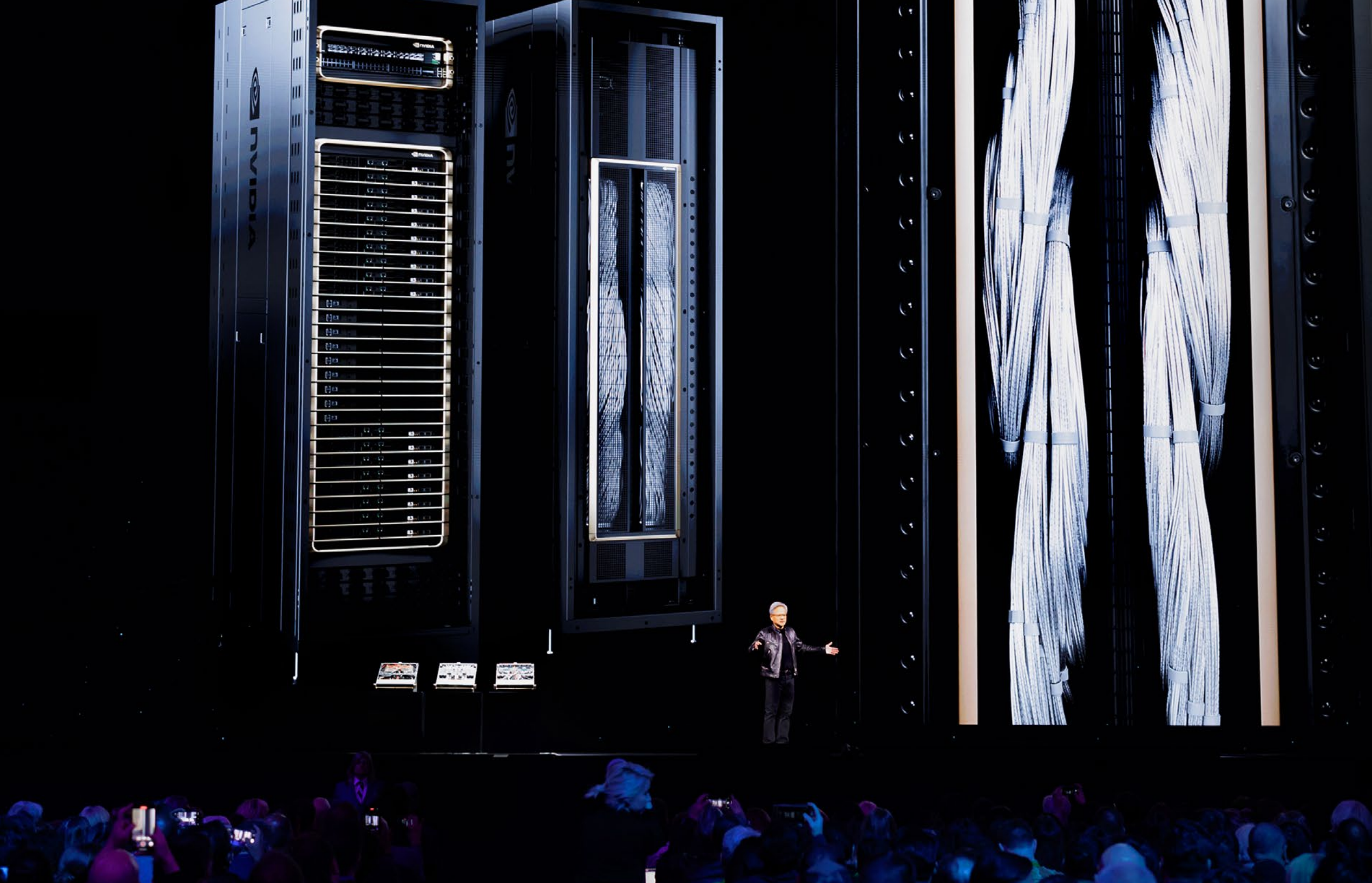
## [...] Blackwell Partners

[...] will be available from partners starting later this year.

[...]soft Azure and Oracle Cloud Infrastructure will be among the first cloud service providers to offer [...]ances, as will NVIDIA Cloud Partner program companies Applied Digital, CoreWeave, Crusoe, IBM Cloud and [...]uds will also provide Blackwell-based cloud services and infrastructure, including Indosat Ooredoo Hutchinson, [...]nde EU Sovereign Cloud, the Oracle US, UK and Australian Government Clouds, Scaleway, Singtel, Northern [...] Yotta Data Services' Shakti Cloud and YTL Power International.

[...]le on NVIDIA DGX™ Cloud, an AI platform co-engineered with leading cloud service providers that gives [...]dicated access to the infrastructure and software needed to build and deploy advanced generative AI models. [...]Oracle Cloud Infrastructure plan to host new NVIDIA Grace Blackwell-based instances later this year.

[...]ard Enterprise, Lenovo and Supermicro are expected to deliver a wide range of servers based on Blackwell [...] products, as are Aivres, ASRock Rack, ASUS, Eviden, Foxconn, GIGABYTE, Inventec, Pegatron, QCT, Wistron, Wiwynn and ZT Systems.

"NVIDIA moved the ball forward with the latest iteration of its speedy NVLink technology."
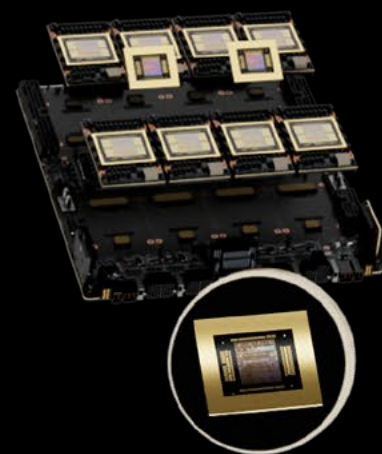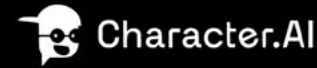
*HPCWire*

To accelerate performance for multitrillion-parameter and mixture-of-experts AI models, the latest iteration of NVIDIA NVLink delivers groundbreaking 1.8TB/s bidirectional throughout per GPU. This ensures seamless high-speed communication among up to 576 GPUs for the most complex LLMs.

"The new AI platform could be a game changer."

*Barron's*

Generative AI is the defining technology of our time. Blackwell is the engine to power this new industrial revolution. Working with the most dynamic companies in the world, we'll realize the promise of AI for every industry.

**HGX B100**

**NVLink Switch**

**GB200 Superchip Compute Node**

**Quantum X800 Switch ConnectX-8 SuperNIC**

**Spectrum X800 Switch BlueField-3 SuperNIC**

aws    Google Cloud    Microsoft Azure    ORACLE CLOUD Infrastructure

## "Widespread adoption anticipated."

*VentureBeat*

Among the many organizations expected to adopt Blackwell are Amazon Web Services, Dell Technologies, Google, Meta, Microsoft, OpenAI, Oracle, Tesla, and xAI.

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

Industry Standard APIs
Text, Speech, Image,
Video, 3D, Biology

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

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

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

Cloud Native Stack
GPU Operator, Network Operator

Customization Cache
P-Tuning, LORA, Model Weights

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

Kubernetes

NVIDIA CUDA

100's of Millions of CUDA GPUs Installed Base

"NVIDIA launches NIM to make it smoother to deploy AI models into production."

*TechCrunch*

NVIDIA Inference Microservices are a new way to package and deliver AI software. The curated selection of microservices adds a new layer to NVIDIA's full-stack computing platform. This layer connects the AI ecosystem of model developers, platform providers, and enterprises with a standardized path to run custom AI models.

"NVIDIA has virtually recreated the entire planet—and now it wants to use its digital twin to crack weather forecasting for good."

*TechRadar*

To help combat the $140 billion in economic losses due to extreme weather brought on by climate change, we announced the Earth-2 digital twin cloud platform for simulating and visualizing weather and climate at unprecedented scale.

"NVIDIA packages inference to deliver generative AI for healthcare."

*Electronics Weekly*

The new suite of NIM for healthcare offers advanced imaging, natural language and speech recognition, and digital biology generation, prediction, and simulation.

"NVIDIA is working to bring AI robots to life."

*Barron's*

We announced a collection of robotics pretrained models, libraries, and reference hardware. We also announced Project GR00T, a general-purpose foundation model for humanoid robots, designed to further our work driving breakthroughs in robotics and embodied AI.
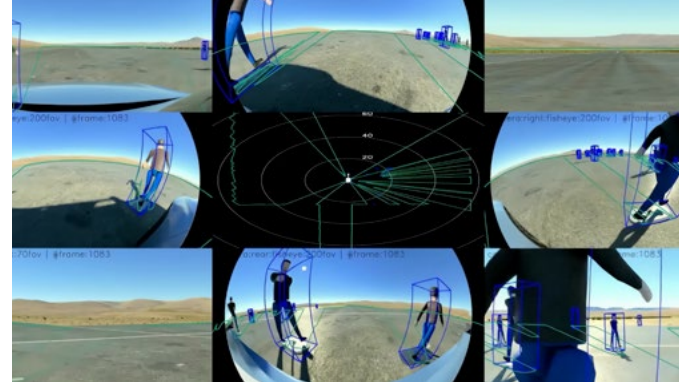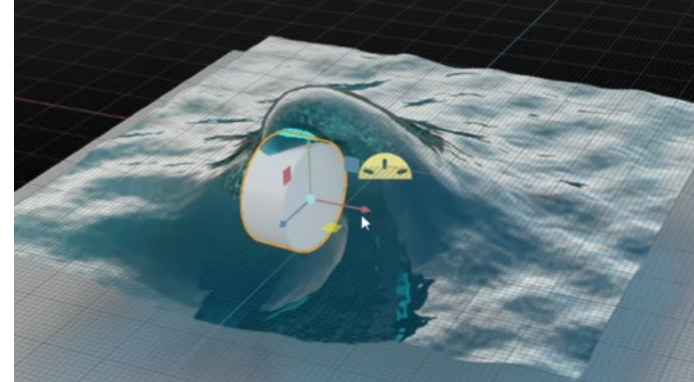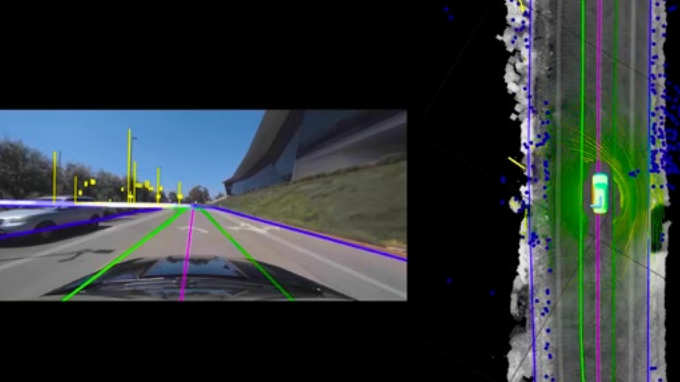
"NVIDIA Omniverse Cloud APIs will elevate digital twins for a new industrial revolution."

*VentureBeat*

NVIDIA Omniverse Cloud will be available as APIs, extending the reach of the world's leading platform for creating industrial digital twin applications and workflows across the entire ecosystem of software makers.

Siemens

# What Analysts Said...

"Move over
Taylor Swift..."

Bernstein

"The soothsayer of Santa Clara
did not disappoint."

Cantor

"[NVIDIA] sits on the
cusp of an entirely new wave
of demand."

UBS

"The leader in AI showcasing innovation
across the full stack
of accelerated computing."

Cowen

"No one else in the industry
can match this capability."

Melius

"[NVIDIA]'s platform expansion
is remarkable...period!"

Wells Fargo

# What Press Said...



"NVIDIA's CEO painted a vision of AI turbocharging computing power."

*The New York Times*

"NVIDIA expects to win outsized chunk of data center spending."

*Bloomberg*

"Event was a daunting reminder of the speed at which NVIDIA is moving."

*Financial Times*

"With Blackwell GPUs, AI gets cheaper and easier, competing with NVIDIA gets harder."

*The Next Platform*

"The NVIDIA frenzy over artificial intelligence has come to this: AI Woodstock."

*The Wall Street Journal*

"GTC can serve as a preview of where the entire field is going."

*Fast Company*

**nVIDIA GTC**

# Explore and Share

@nvidia #gtc24

"GTC exists to inspire the world on the art-of-the-possible."

**300,000** Registrations

**19,000** In-Person Attendees

**1,100** Sessions

**33M** Keynote Views

**29,000** Press Articles

**300** Partner Sponsorships

"GTC: Spearheading AI and accelerated computing innovation."

*siliconANGLE*

Transforming AI Panel: the authors of the seminal research paper—*Attention Is All You Need*—that introduced the transformer neural network architecture came together at GTC.

"NVIDIA's GTC conference, not surprisingly, was an absolute whirlwind."

*VentureBeat*

"GTC 2024 was the single most important event in the history of the technology industry."

*SiliconANGLE theCUBE*

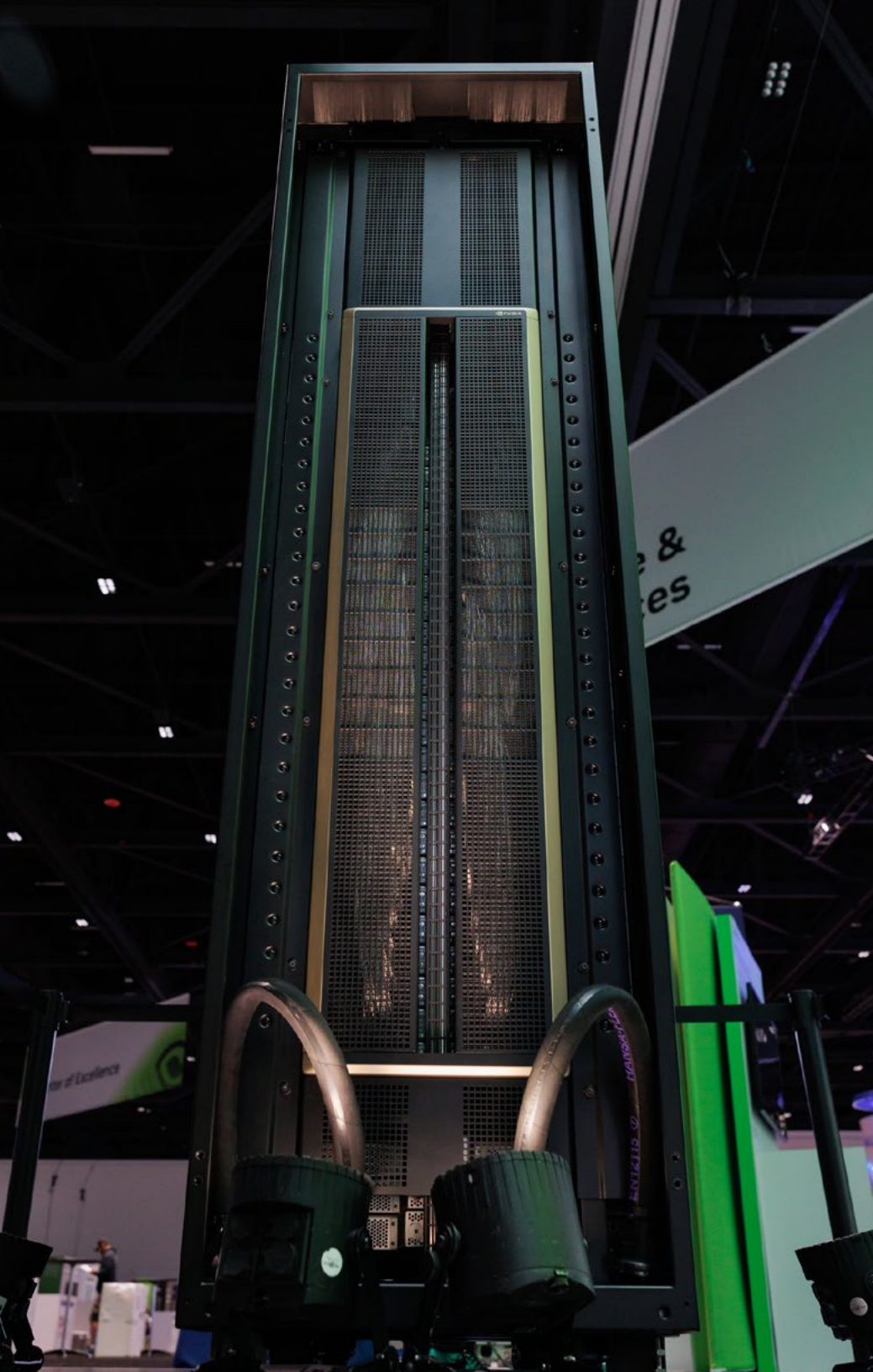"NVIDIA's conference captured the industry's attention."
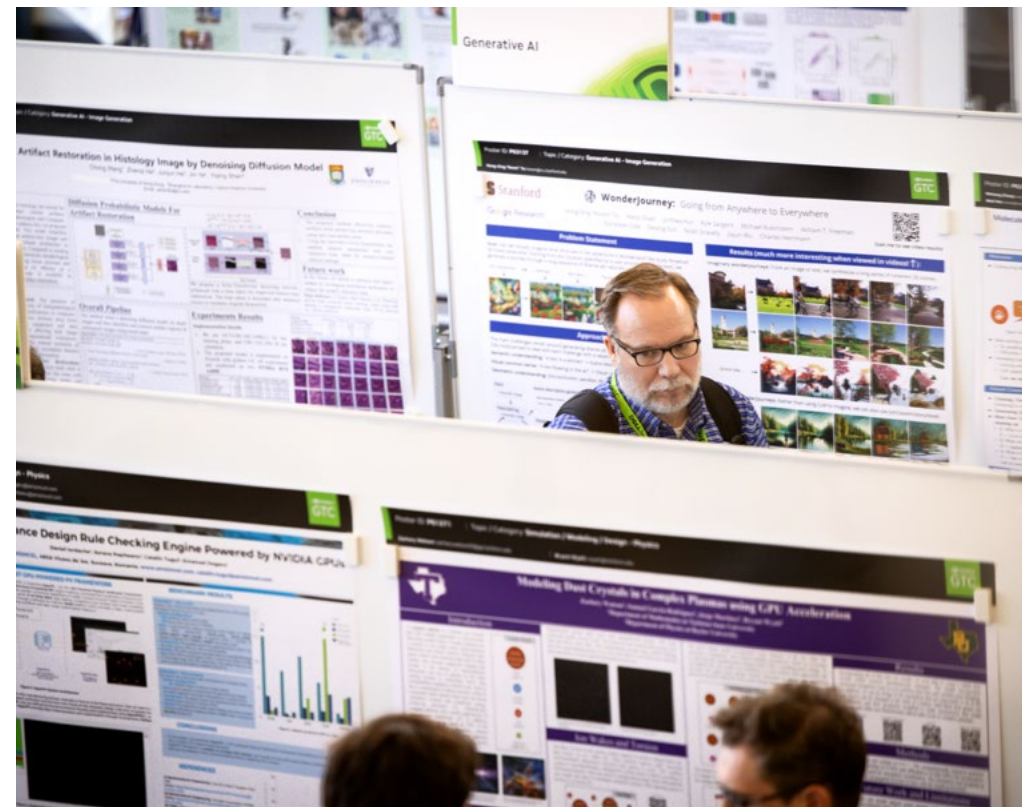
*Yahoo Finance*
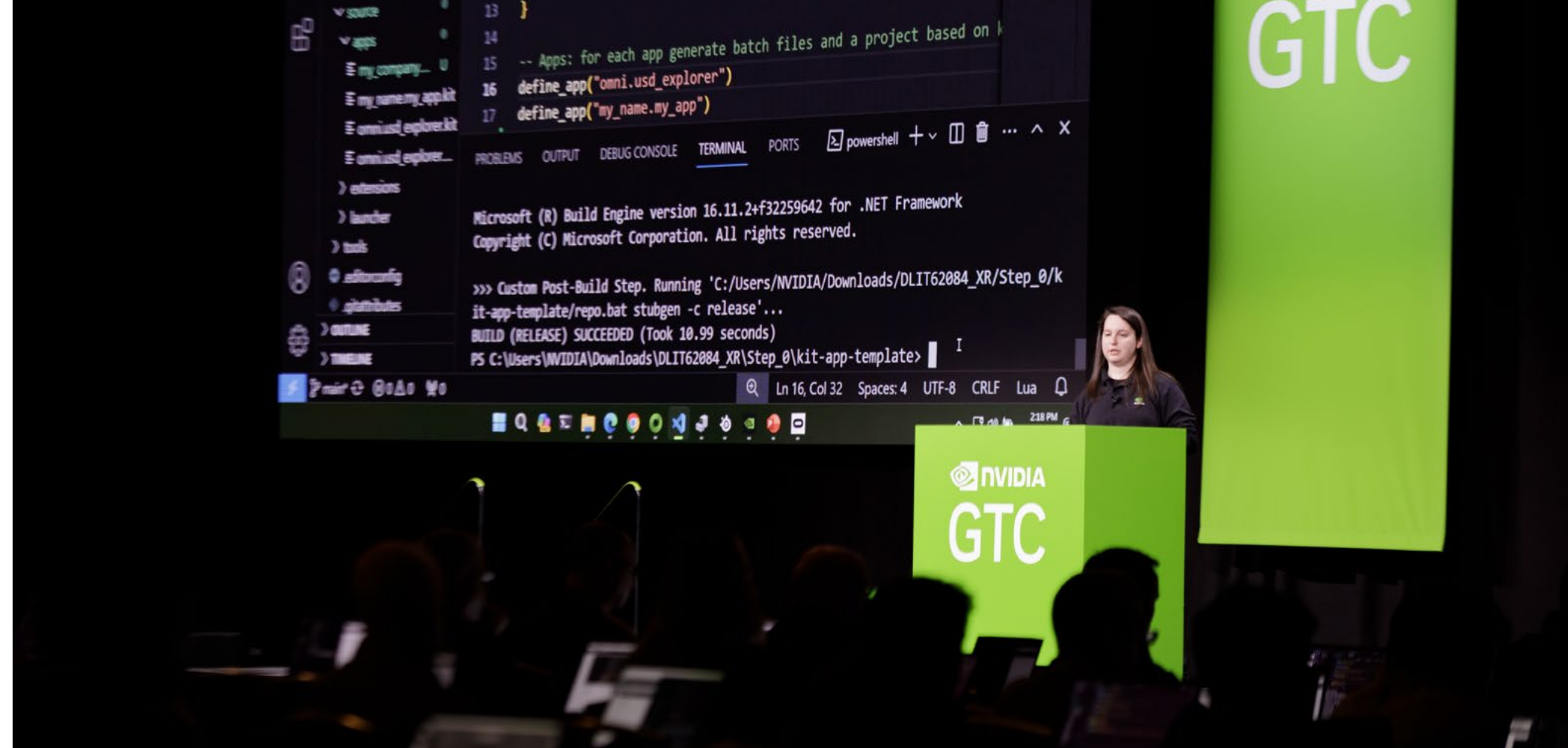
"The Generative AI Future Is Now."

*HPCWire*

Engagement with brilliant minds was evident, from 'Connect with Experts' sessions to watching industry leaders speak, and even casual encounters in the networking lounge.

"For three decades we've pursued accelerated computing, with the goal of enabling transformative breakthroughs like deep learning and AI. Generative AI is the defining technology of our time.

Blackwell is the engine to power this new industrial revolution. Working with the most dynamic companies in the world, we will realize the promise of AI for every industry."

Jensen Huang

"NVIDIA's AI boom is only getting started."