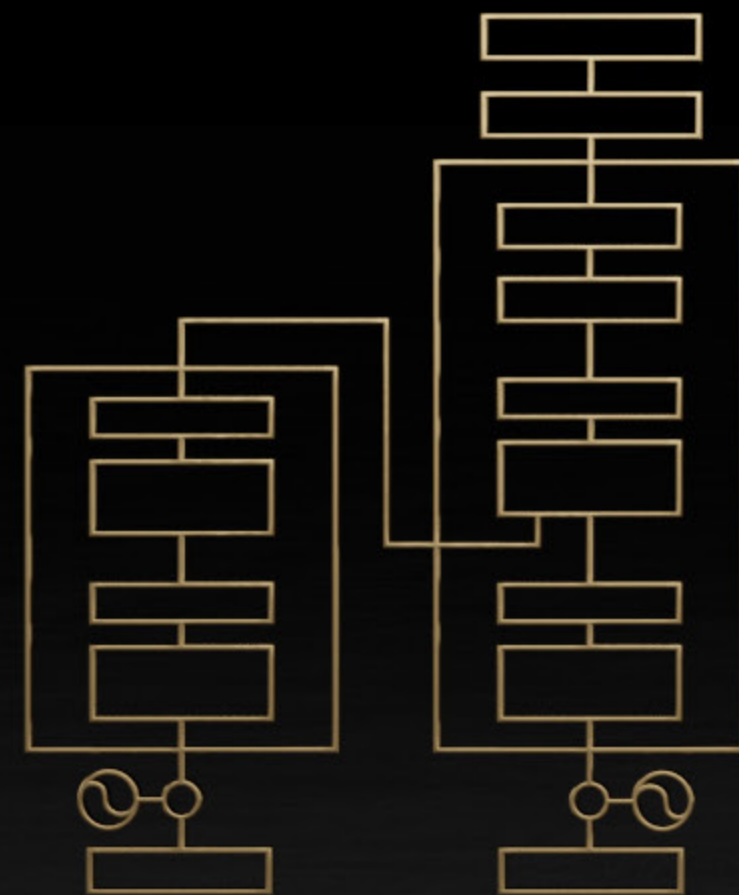Accelerated Computing
is the Path Forward

AI is Software that
Writes Software

Data Center is the
New Unit of Computing

AI-on-5G Kickstarts the
4th Industrial Revolution

Autonomous Systems in
Real and Virtual Worlds

# NEW NVIDIA TECHNOLOGIES

Omniverse
Isaac

Megatron
Drug Discovery
Quantum Computing
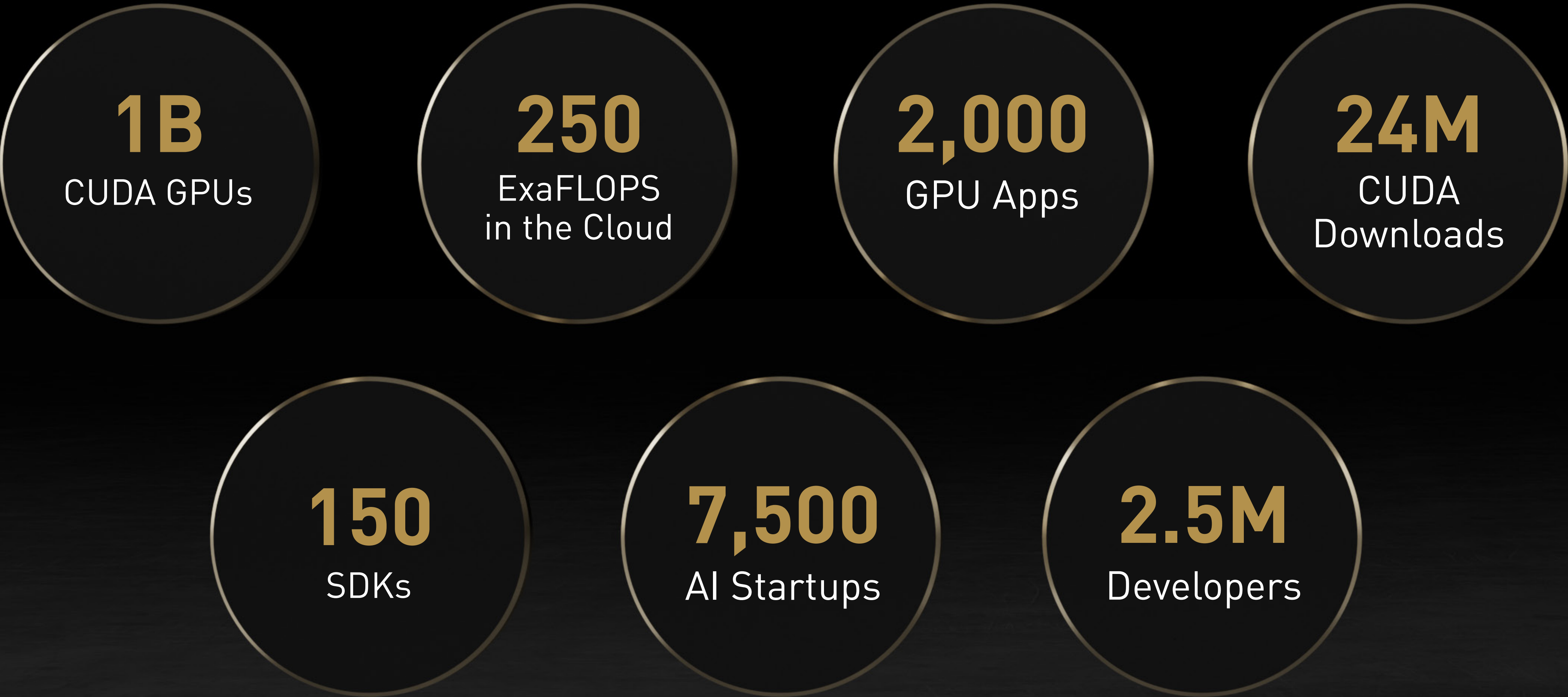
Jarvis
Merlin
Maxine
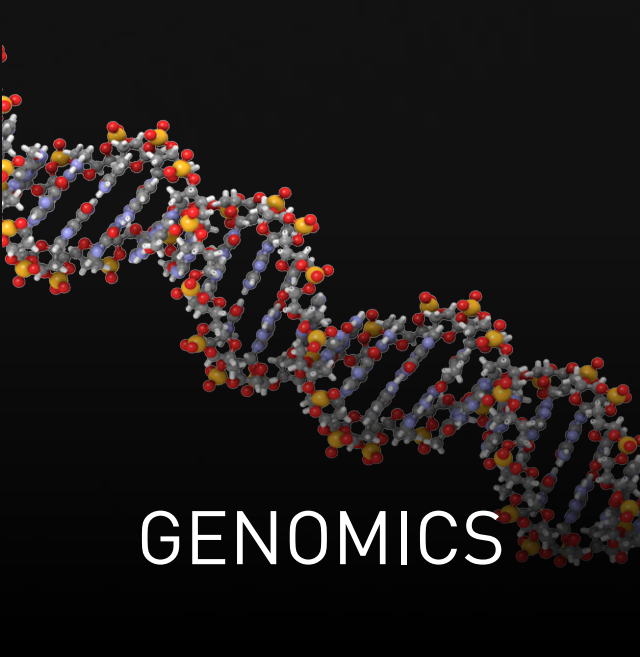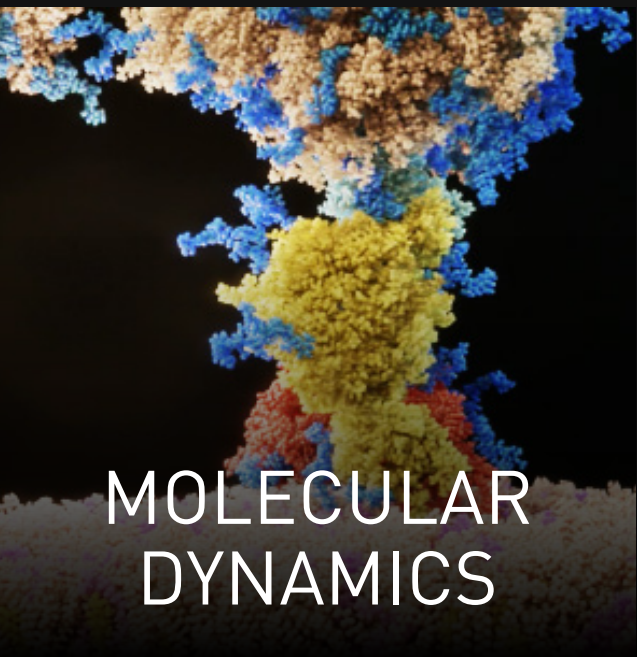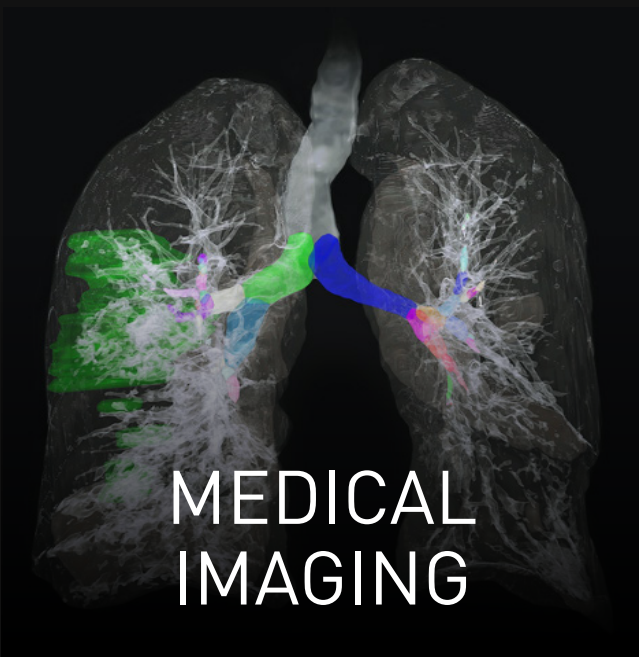Morpheus
NVIDIA AI

DRIVE

RTX

DGX
Grace
BlueField
DOCA

EGX
5G

Hyperion
Atlan
Orin

# NVIDIA IS A COMPUTING PLATFORM
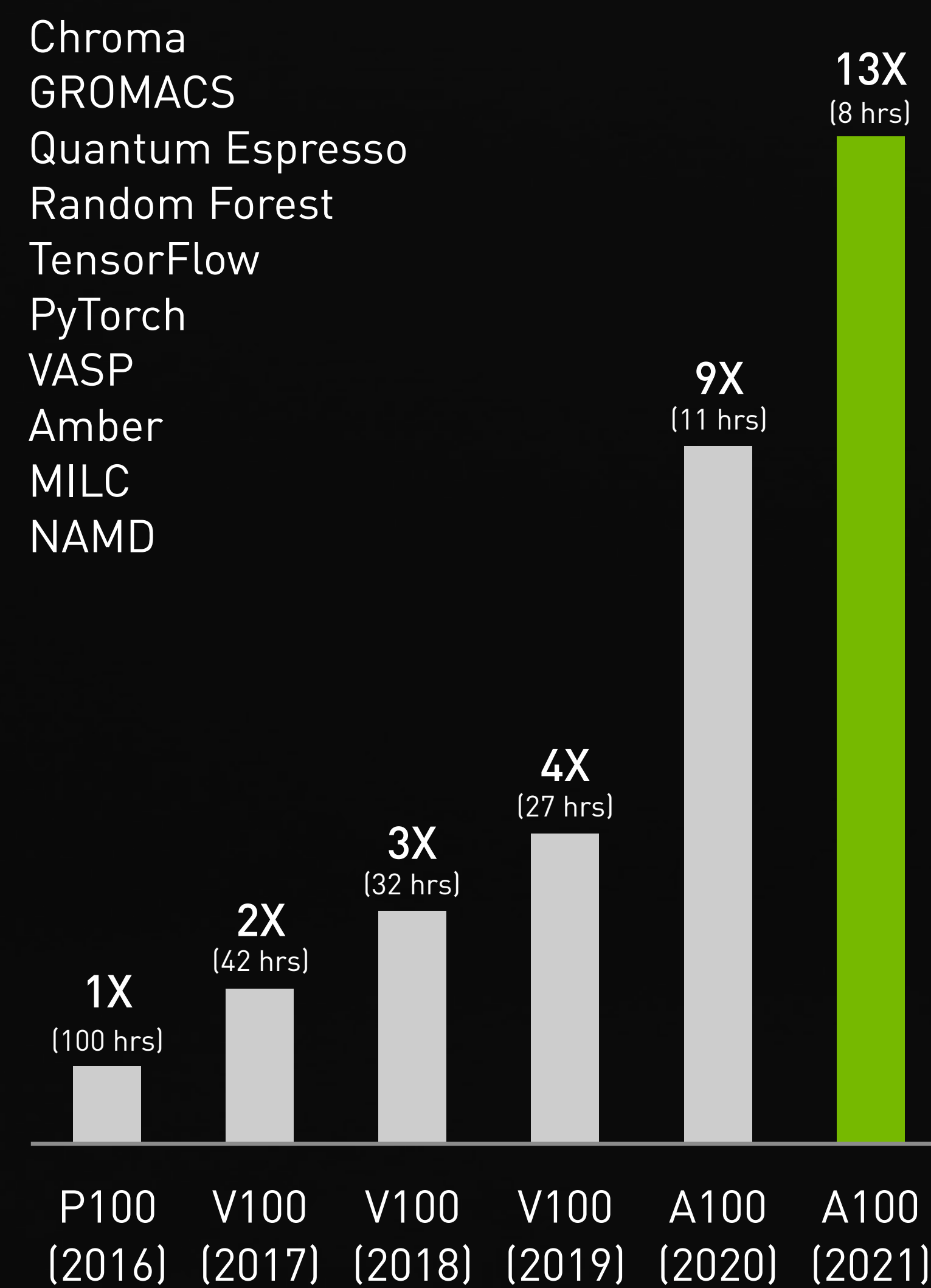
**1B** CUDA GPUs

**250** ExaFLOPS in the Cloud

**2,000** GPU Apps

**24M** CUDA Downloads

**150** SDKs

**7,500** AI Startups

**2.5M** Developers

## AN INSTRUMENT OF SCIENCE

ASTRONOMY & ASTROPHYSICS

CLIMATE & WEATHER

SCIENTIFIC VISUALIZATION

MEDICAL IMAGING

MOLECULAR DYNAMICS

GENOMICS

GAME DEVELOPMENT

RENDERING & RAY TRACING

AR / VR

Leeza SOHO, Beijing
By ZAHA HADID ARCHITECTS

# GIVING SCIENTISTS A TIME MACHINE

## FULL-STACK OPTIMIZATION
### 13x in 5 years

Chroma
GROMACS
Quantum Espresso
Random Forest
TensorFlow
PyTorch
VASP
Amber
MILC
NAMD

13X
(8 hrs)

9X
(11 hrs)

4X
(27 hrs)

3X
(32 hrs)

2X
(42 hrs)

1X
(100 hrs)

P100
(2016)

V100
(2017)

V100
(2018)

V100
(2019)

A100
(2020)

A100
(2021)

## MULTI-GPU MULTI-NODE SCALE
### NAMD Version 3.0

▲ 2.14 baseline  ●— 3.0a7 En

ns/day (higher is better)

120

80

40

0

0    2    4    6    8

GPUs

## LARGEST AI + MD SIMULATION
### 305 Million Atoms



## EXASCALE HPC AI
### HPL-AI and HPL Performance

HPL-AI

HPL

10000

1000

100

PFLOPS

Summit
Sierra
JUWELS
Fugaku
Leonardo
Perlmutter

2018    2019    2020    2021

# FOR THE DA VINCIS OF OUR TIME

## SCIENTISTS AT GTC

**Yoshua Bengio**
University of Montreal
Quebec AI Institute

**Yann LeCun**
Facebook
New York University

**Geoffrey Hinton**
University of Toronto
Google
Vector Institute

**Daphne Koller**
Insitro
Coursera
Stanford

**Jürgen Schmidhuber**
Dalle Molle Institute for
AI Research

**Raquel Urtasun**
University of Toronto

**Alvy Ray Smith**
Pixar
Altamira

**Abhay Parasnis**
Adobe

**Kim Libreri**
Epic Games

**Rommie Amaro**
University of
California, San Diego

**Soumith Chintala**
Facebook

**Rose Yu**
University of
California, San Diego

## TALKS AT GTC

| AI | 5G | IOT & EDGE |
|----|----|----|
| QUANTUM COMPUTING | SPEECH NLU RECOMMENDERS | SELF-DRIVING CARS |
| CYBERSECURITY | DIGITAL TWINS | ROBOTICS |

## LEADERS AT GTC

Adobe · Alibaba Cloud · arm · Audi · aws

Baidu 百度 · Baker Hughes · BAYER · BMW · EPIC GAMES · ERICSSON

FACEBOOK · Ford · Google · INDUSTRIAL LIGHT & MAGIC · Kroger

Microsoft · NASA · NTT · ORACLE Cloud Infrastructure · PayPal

Pinterest · Pfizer · PIXAR · Red Hat · SIEMENS Healthineers · Snapchat

Spotify · Tencent Cloud · TESCO · tsmc · Twitter

Verizon · vmware · VW · Walmart · WELLS FARGO

# NVIDIA OMNIVERSE

AI   Path-Tracing   USD   Materials   Physics

NVIDIA OMNIVERSE

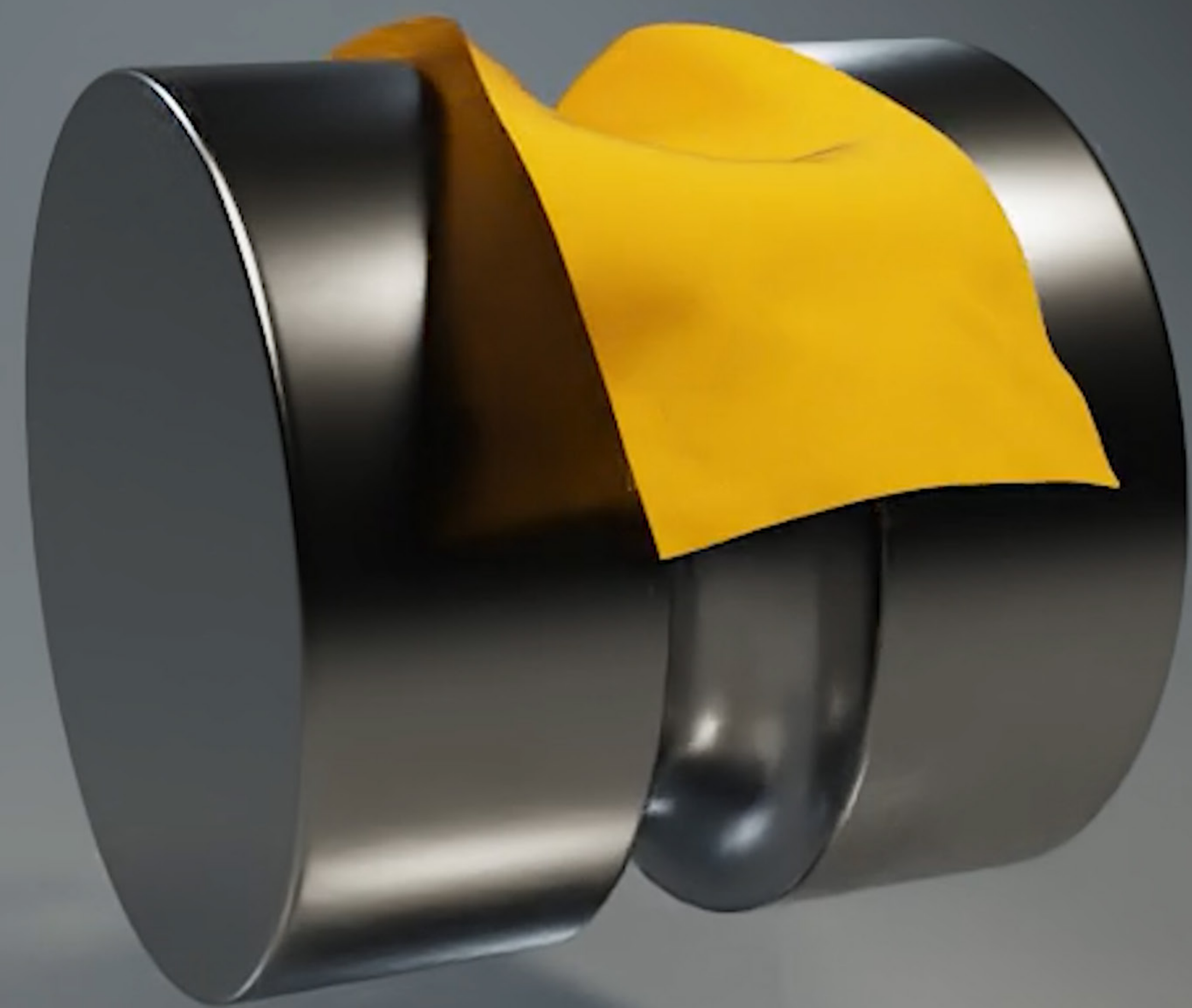# COLLABORATING AND SIMULATING IN OMNIVERSE



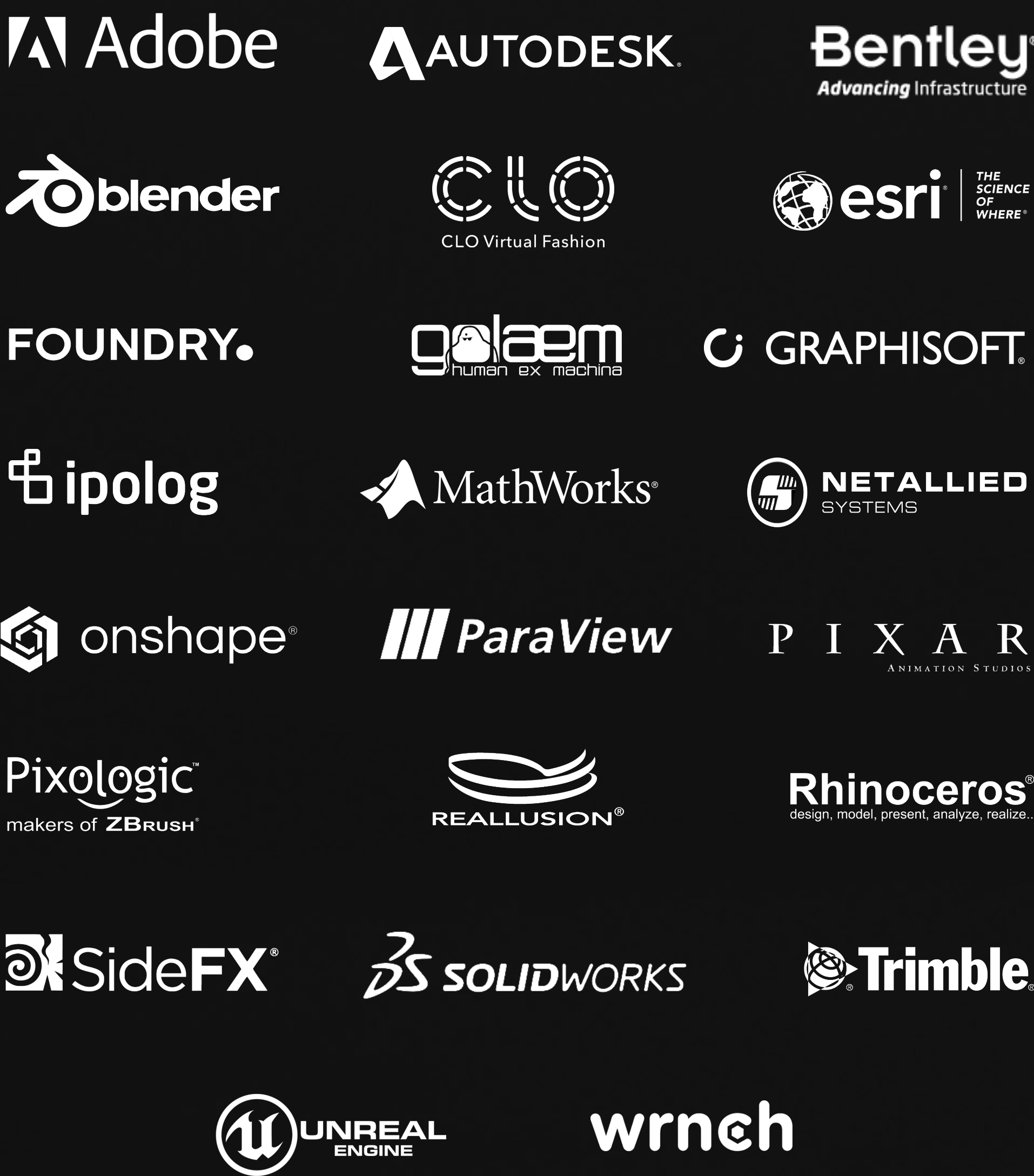10M DESIGNERS | 20M CREATORS | 1M SCIENTISTS | 2M DEVELOPERS | 40M ENGINEERS

**NVIDIA ISAAC DIGITAL TWIN IN OMNIVERSE**

# CONNECTING AND CREATING WITH OMNIVERSE

## SOFTWARE

Adobe · AUTODESK · Bentley Advancing Infrastructure
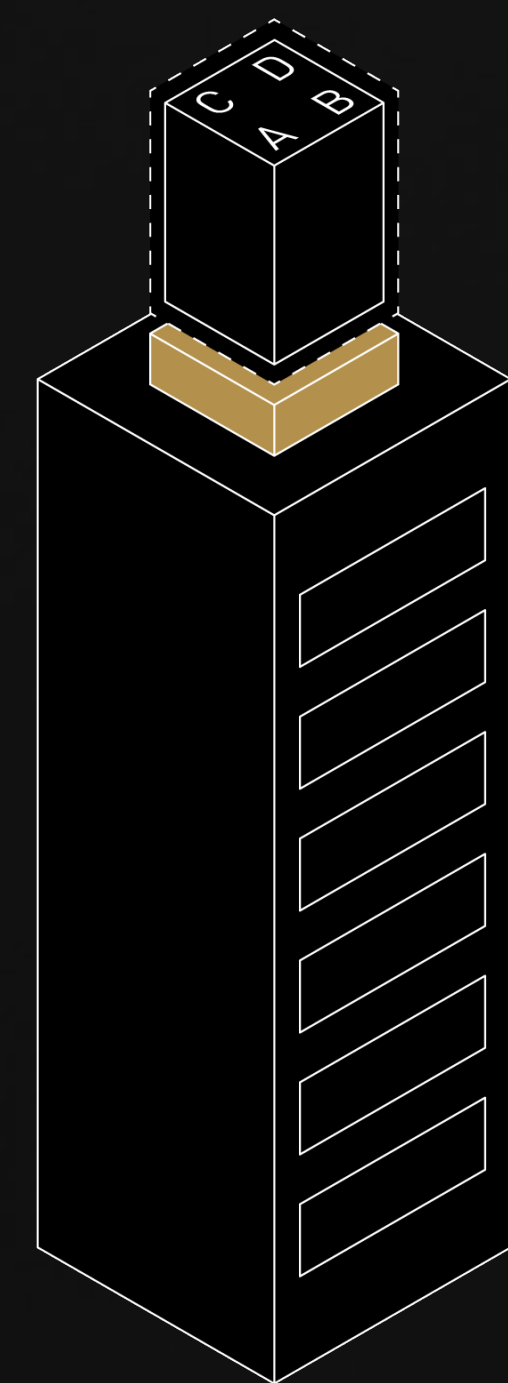
blender · CLO Virtual Fashion · esri THE SCIENCE OF WHERE

FOUNDRY · galaem human ex machina · GRAPHISOFT

ipolog · MathWorks · NETALLIED SYSTEMS

onshape · ParaView · PIXAR Animation Studios

Pixologic makers of ZBrush · REALLUSION · Rhinoceros design, model, present, analyze, realize...

SideFX · SOLIDWORKS · Trimble

UNREAL ENGINE · wrnch

## SYSTEMS

ASUS · BOXX

CISCO · DELL Technologies

HP · Lenovo

PNY · SUPERMICRO

## PIONEERS

BMW GROUP · ERICSSON · Foster + Partners

Framestore · IBI · Gensler

INDUSTRIAL LIGHT & MAGIC · KITESTRING · KPF

Lightcraft TECHNOLOGY · TANGENT · TURBOSQUID

W–B WOODS BAGOT · WPP · VOLVO

# DATA CENTER IS THE NEW UNIT OF COMPUTING



Monolithic

# DATA CENTER IS THE NEW UNIT OF COMPUTING



Monolithic

Software-Defined
Data Center

# DATA CENTER IS THE NEW UNIT OF COMPUTING



Monolithic

Software-Defined
Data Center

Disaggregated,
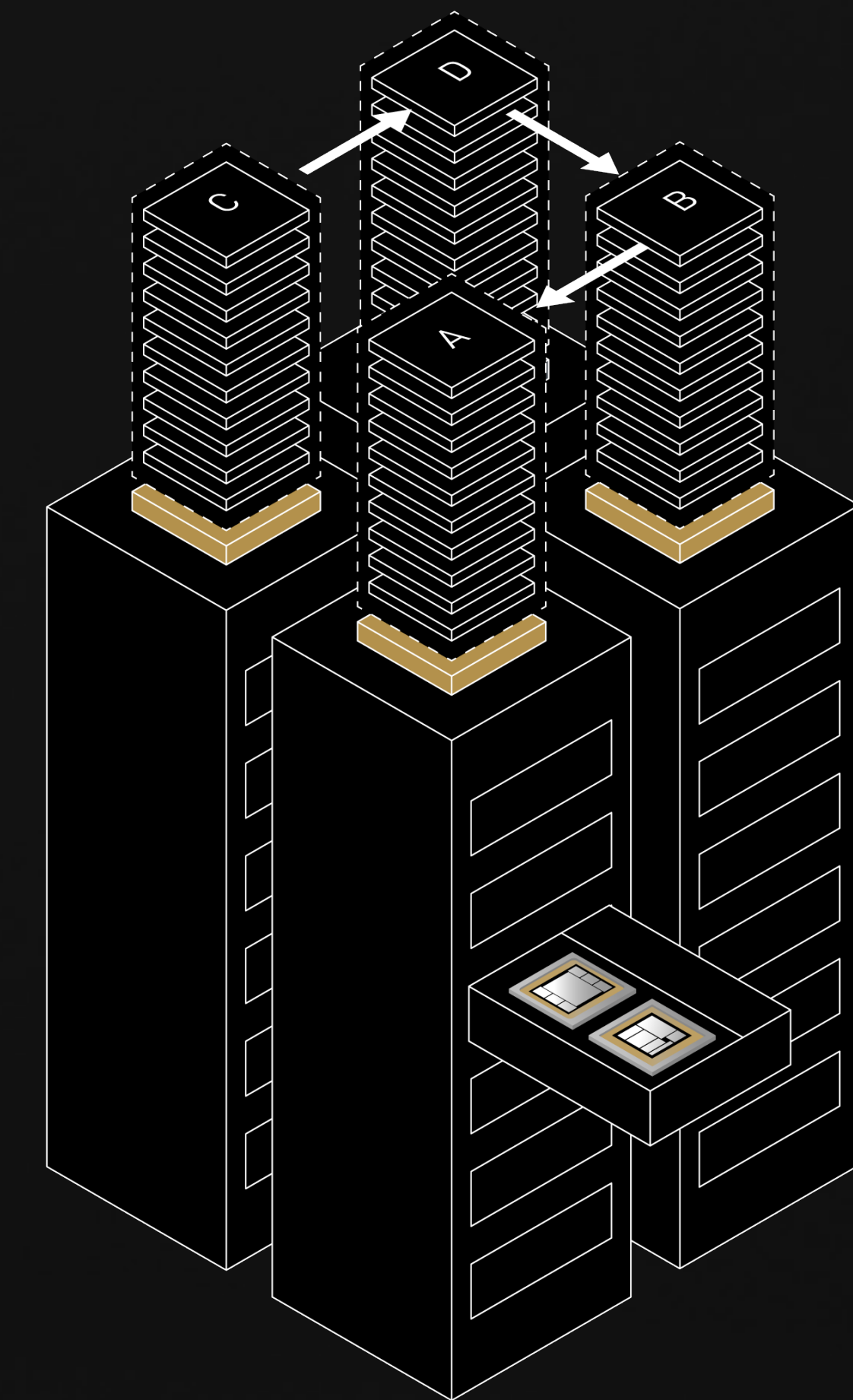Micro-Services, Scaled Out

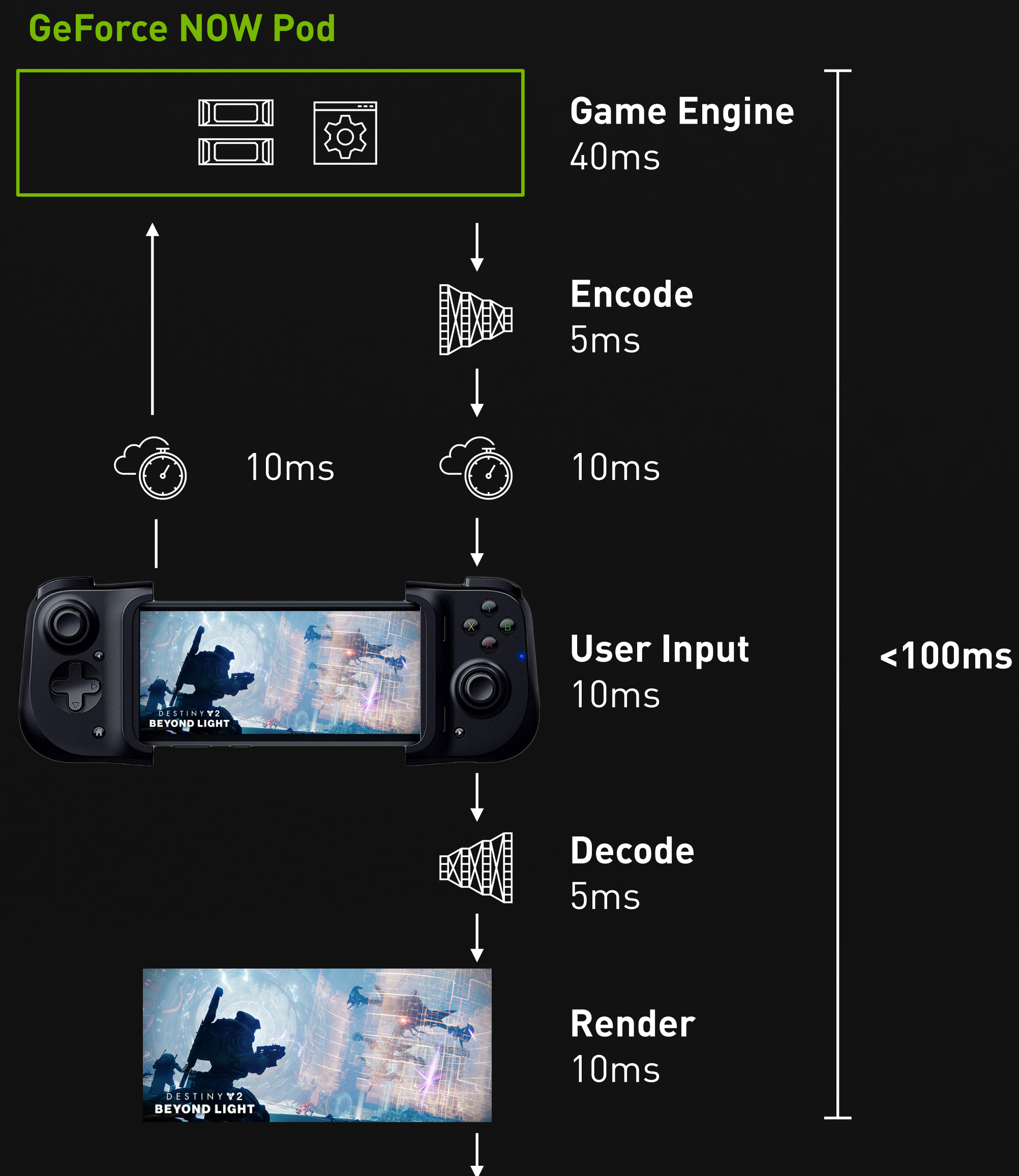# DATA CENTER IS THE NEW UNIT OF COMPUTING



Monolithic

Software-Defined
Data Center

Disaggregated,
Micro-Services, Scaled Out

GPU-Accelerated
Computing

# DATA CENTER IS THE NEW UNIT OF COMPUTING



Monolithic

Software-Defined
Data Center

Disaggregated,
Micro-Services, Scaled Out

GPU-Accelerated
Computing

DPU-Accelerated
Data Center Infrastructure

# BLUEFIELD SECURES AND ACCELERATES GEFORCE NOW CLOUD GAMING

Isolate and Secure Infrastructure | High Quality-of-Service | More Concurrent Users

**GeForce NOW Pod**

**Game Engine**
40ms

**Encode**
5ms

10ms    10ms

**User Input**
10ms

**Decode**
5ms

**Render**
10ms

<100ms

# BLUEFIELD SECURES AND ACCELERATES GEFORCE NOW CLOUD GAMING

## Isolate and Secure Infrastructure | High Quality-of-Service | More Concurrent Users



**GeForce NOW Pod**

Game Engine
40ms

Encode
5ms

10ms    10ms

User Input
10ms

<100ms

Decode
5ms

Render
10ms

TRADITIONAL CLOUD
GAMING SERVER

Game | Graphics Driver | Video

SDDC | Security | Telemetry

NAT | DDOS | Reverse Proxy

Ethernet NIC

# BLUEFIELD SECURES AND ACCELERATES GEFORCE NOW CLOUD GAMING

## Isolate and Secure Infrastructure | High Quality-of-Service | More Concurrent Users

# ANNOUNCING
# NVIDIA BLUEFIELD-3

## 400 Gbps Data Center Infra Processor

Offloads and Accelerates Data Center Infrastructure

Isolates Application from Control and Management Plane

Powerful CPU – 16x Arm A78 Cores

Process Networking, Storage, and Security at 400 Gbps

22 Billion Transistors

DATA PATH ACCELERATOR

CONNECTX-7

PCIe GEN 5.0

DDR5 MEMORY INTERFACE

ARM CORES

ACCELERATION ENGINES

# EXPONENTIAL GROWTH IN DATA CENTER INFRASTRUCTURE PROCESSING

Cloud-Native | Disaggregation | Micro-Services | AI | Zero-Trust Security



100X

**BlueField-4**
64 Billion Transistors
160 SPECint
1000 TOPS
800 Gbps

10X

**BlueField-3**
22 Billion Transistors
42 SPECint
1.5 TOPS
400 Gbps

1X

**BlueField-2**
7 Billion Transistors
9 SPECint
0.7 TOPS
200 Gbps

DOCA — ONE ARCHITECTURE

2020          2022          2024

SPECrate2017_int

# NVIDIA DGX
## The Computer of AI Researchers

**9** of 10
Top
Universities

**8** of 10
Top
Global Telcos

**7** of 10
Top
U.S. Hospitals

**7** of 10
Top
Consumer
Companies

**6** of 10
Top
U.S. Banks

**7** of 10
Top
Car Makers

**10** of 10
Top
Aerospace
Companies

**DGX SUPERPOD**
AI Data Center As-a-Product

**DGX A100**
AI Data Center
Building Block

**DGX STATION A100**
AI Data Center-in-a-Box

# ANNOUNCING NVIDIA DGX STATION 320G

Workgroup AI Supercomputer-in-a-Box

Plug-into-the-Wall Instant AI Infrastructure

2.5 petaFLOPS

320 GB at 8 TB/sec

7.68 TB NVMe

28 MIGs

1500W and < 37db

$149,000 or $9,000/Month Subscription

ANNOUNCING THE NEW DGX SUPERPOD

World's First Cloud-Native Supercomputer | Secured by NVIDIA BlueField | Multi-Tenant Bare-Metal Performance

A100 80GB

BLUEFIELD-2

BASE COMMAND

# NVIDIA CLARA DISCOVERY

**ANNOUNCING**

# RECURSION BUILDS PHARMA AI SUPERCOMPUTER WITH NVIDIA DGX SUPERPOD

## BioHive-1 Aims to Decode Biology and Industrialize Drug Discovery

Recursion OS Built on NVIDIA DGX SuperPOD
Generates, Analyzes, and Derives Insight from
Biological and Chemical Datasets

Generate up to 9 Million Images, or Approximately
80 Terabytes of Data, Across up to 1.5 Million
Experiments per Week

# ANNOUNCING NVIDIA CUQUANTUM
## Research the Computer of Tomorrow on the Most Powerful Computer Today

### GPU-ACCELERATED QUANTUM SIMULATIONS

$$|\dot\Psi\rangle = U|\Psi\rangle$$

cuQuantum

### STATE VECTOR SIMULATION
Scales to 10's of Qubits

10 Days

2 Hours

### TENSOR NETWORK SIMULATION
Scales to 1000's of Qubits

9 Years

4 Days

■ Dual-CPU  ■ DGX A100

"Using the Cotengra/Quimb packages, NVIDIA's new cuQuantum SDK, and the Selene supercomputer, we've generated a sample of the Sycamore quantum circuit at depth=20 in record time (less than 10 minutes). This sets the benchmark for quantum circuit simulation performance and will help advance the field of quantum computing by improving our ability to verify the behavior of quantum circuits."

—Johnnie Gray, Research Scientist, Caltech
Garnet Chan, Bren Professor of Chemistry, Caltech

Footnotes: State Vector- 1,000 circuits , 36 qubits depth =10, complex 64 | CPU: Qiskit (IBM) on Dual AMD EPYC 7742 | GPU: Qgate  (NVAITC) on DGX-A100
Tensor Network - 53 qubits, depth 20 | CPU: Estimated Quimb (Caltech) on Dual AMD EPYC 7742 | GPU: Quimb (Caltech) on DGX-A100

# DIVERSE DATA CENTER ARCHITECTURES

ENTERPRISE COMPUTING

HYPERSCALE

SCIENTIFIC COMPUTING

STORAGE SERVERS

ACCELERATED HPC

ACCELERATED HYPERSCALE

# DATA-COMPUTE DEMAND GROWING FASTER THAN SYSTEM BANDWIDTH

GPU Starved by CPU Memory and PCIE Bandwidth

| | | |
|---|---|---|
| GPU | 8,000 | GB/sec |
| CPU | 200 | GB/sec |
| PCIE Gen 4 | 16* | GB/sec |
| Mem-to-GPU | 64 | GB/sec |

DDR4

HBM2e

x86

GPU

GPU

GPU

GPU

\* Effective bandwidth between CPU and each A100 GPU in half of a DGX A100 today illustrated

# DATA-COMPUTE DEMAND GROWING FASTER THAN SYSTEM BANDWIDTH

GPU Starved by CPU Memory and PCIE Bandwidth

| | | |
|---|---|---|
| GPU | 8,000 | GB/sec |
| CPU | 200 | GB/sec |
| PCIE Gen 4 | 16* | GB/sec |
| Mem-to-GPU | 64 | GB/sec |

DDR4

HBM2e

x86

GPU

GPU

GPU

GPU

* Effective bandwidth between CPU and each A100 GPU in half of a DGX A100 today illustrated

# A NEW COMPUTING ARCHITECTURE FOR AI AND DATA SCIENCE

30X Increase System Memory to GPU

| | | | |
|---|---|---|---|
| GPU | 8,000 | GB/sec | |
| CPU | 500 | GB/sec | |
| NVLINK | 500 | GB/sec | |
| Mem-to-GPU | 2,000 | GB/sec | 30X |

LPDDR5x

HBM2e

| GRACE | GPU |
|---|---|
| GRACE | GPU |
| GRACE | GPU |
| GRACE | GPU |

**ANNOUNCING NVIDIA GRACE**

CPU Designed for Giant-Scale AI and HPC Accelerated Computing

**ANNOUNCING**
# THE WORLD'S FASTEST SUPERCOMPUTER FOR AI

20 Exaflops of AI

Powered by NVIDIA Grace CPU and
Next Generation NVIDIA GPU

HPC and AI for Scientific and Commercial Apps

Advance Weather, Climate, and Material Science

# SUPERCOMPUTING COMMUNITY EMBRACES ARM

"Alps will use NVIDIA's novel Grace CPU to converge AI technologies and classic supercomputing in one single powerful data center infrastructure."

**CSCS**
Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

-Thomas Schulthess
Director CSCS

"Thanks to NVIDIA's new Grace CPU, we'll be able to deliver advanced scientific research using high-fidelity 3D simulations and analytics with data sets that are larger than previously possible."

**Los Alamos**
NATIONAL LABORATORY

-Thom Mason
Director LANL

# 3 CHIPS. YEARLY LEAPS. ONE ARCHITECTURE.

# EXPANDING ARM ECOSYSTEM BEYOND MOBILE

| CLOUD | SCIENTIFIC COMPUTING | EDGE AND ENTERPRISE | PC |
|---|---|---|---|
| **NVIDIA AI**<br>**NVIDIA GeForce NOW** | **NVIDIA AI**<br>**NVIDIA HPC** | **NVIDIA AI**<br>**NVIDIA DeepStream** | **NVIDIA AI**<br>**NVIDIA RTX** |
| NVIDIA GPU | NVIDIA A100    NVIDIA BlueField-2 | NVIDIA T4 | GeForce RTX 30 Series |
| AWS Graviton | Ampere Altra | Marvell OCTEON | MediaTek MT819x |

# WAVES OF AI

AI Computing          Cloud          5G Industrial Edge          Robotics

# ANNOUNCING NVIDIA EGX ENTERPRISE PLATFORM

# ANNOUNCING NVIDIA AERIAL A100
## AI-on-5G

ERICSSON  FUJITSU  MAVENIR  Capgemini  Radisys

**NVIDIA AERIAL**

5G CORE

| IPSEC/TLS | DPI | ASAP2 Acceleration |

5G RAN

| 5T FOR 5G | GPU Direct | cuVNF | cuPHY |

...chnologies  FUJITSU  GIGABYTE  H3C

Lenovo  QCT  SUPERMICRO

NVIDIA
BlueField-2

NVIDIA
A100/A30

NVIDIA
A40/A10

**NVIDIA AERIAL
A100**

## GOOGLE CLOUD AND NVIDIA PARTNER TO DELIVER AI-ON-5G

Anthos-Enabled 5G Edge to Deliver Low-Latency,
Secure and Mission-Critical Edge AI Applications

Enables the Rapid Delivery of New Services and
Applications at the 5G Edge

Provides a Consistent Platform for Application
Deployments from Cloud to the Edge

Google Cloud

# EVERY USER AND WORKLOAD OF A DATA CENTER IS A SECURITY THREAT

# ANNOUNCING NVIDIA MORPHEUS

10.244.0.50 -> 10.244.0.14
SI: secret_keys

{
  "who": "32205660 pVycjWHH",
  "thing": "within there ImXxoVGN",
  "condition": "Nnpdtcvs 89234632",
  "key": "nearly appear n1kLSr8A 41142934 actually QZnPjQyd"

ANNOUNCING NVIDIA EGX ENTERPRISE PLATFORM

Enterprise-Ready Suite of AI and Data Science Software

DATA PREP → TRAIN → SIMULATE → ORCHESTRATE

NVIDIA AI ENTERPRISE

NVIDIA GPU CLOUD REGISTRY

| cuIO | cuDF | cuDNN | cuML | cuGraph | TAO | TensorRT | Triton |

NVIDIA CUDA, DOCA, vGPU, Magnum IO, **Aerial 5G**, **Morpheus**

VMware vSphere

NVIDIA CERTIFIED

Atos    DELL Technologies    FUJITSU    GIGABYTE    H3C
hp    inspur    Lenovo    QCT    SUPERMICRO

NVIDIA BlueField-2

NVIDIA A100/A30

NVIDIA A40/A10

NVIDIA Aerial A100

Convolutional Network (CNN)

Transformer

Reinforcement Learning

Autoencoders

Long Short-Term Memory (LSTM)

Generative Adversarial Network (GAN)

# NVIDIA'S AI

**DLSS**
AI Rendering

NATIVE 4K | DLSS 4K

**StyleGAN**
High-Res Image Generation

**GANcraft**
3D Scene Generation

**GANverse3D**
2D to 3D Render

**Sim2Real**
Quadruped Locomotion

**Face Vid2Vid**
Audio-Driven Character Animation

**BioMegatron**
Medical Language Model

**3DGT**
3D Ground Truth

**SimNet**
Physics-Informed AI Model

**OrbNet**
AI for Quantum Chemistry

# NGC PRE-TRAINED MODELS

## Production-Quality AI Models

Trained by Experts for Enterprise Deployment

Credentials to Find Models You Trust

Continuously Updated to be State-of-the-Art

Adapt with NVIDIA TAO and Orchestrate
with NVIDIA Fleet Command

Reference AI Code Samples to Ease
Application Development

---

**NVIDIA. NGC | CATALOG**

Download

### GAZE ESTIMATION

**Application**
Gaze detection for a person - point of regard (X, Y, Z) and gaze vector (theta and phi).

**Popularity**

**Domain**
Computer Vision

**Usage**
Unrestricted

**License**
TLT Licence

**Training Dataset**
Proprietary dataset with more than 220k images.

**Performance**
T4 (1698 FPS)

Xavier (704 FPS)

NX (510 FPS)

CV   Computer Vision   DL   Deep Learning
Gaze   TLT

Expand Credentials

### Model Overview
The model described in this card detects a person's eye gaze point of regard (X, Y, Z) and gaze vector (theta and phi). The eye gaze vector can also be derived from eye position and eye gaze points of regard.

### Model Architecture
GazeNet is a multi-input and multi-branch network. The four input for GazeNet consists: Face crop, left eye crop, right eye crop, and facegrid. Face, left eye, and right eye branch are based on AlexNet as feature extractors. The facegrid branch is based on fully connected layers. Please see the paper in the citations for an example of the model architecture.

### Training Algorithm
The training algorithm optimizes the network to minimize the root mean square error between predicted and ground truth point of regards.

### How to use this model
Primary use case for this model is to detect eye point of regard and gaze vector. The model can be used to detect eye gaze point of regard by using appropriate video or image decoding and pre-processing. In the TLT Computer Vision Inference Pipeline, gaze estimation network results are used to determine whether the subjects are looking at the camera. See the following image for an illustration of eye gaze estimation usage.

### Input
GazeNet is a multi-input network, which takes in face crop image, left eye crop image, right crop image, and facegrid.

Face Image which is gray scale. 224 x 224 x 1

### Output
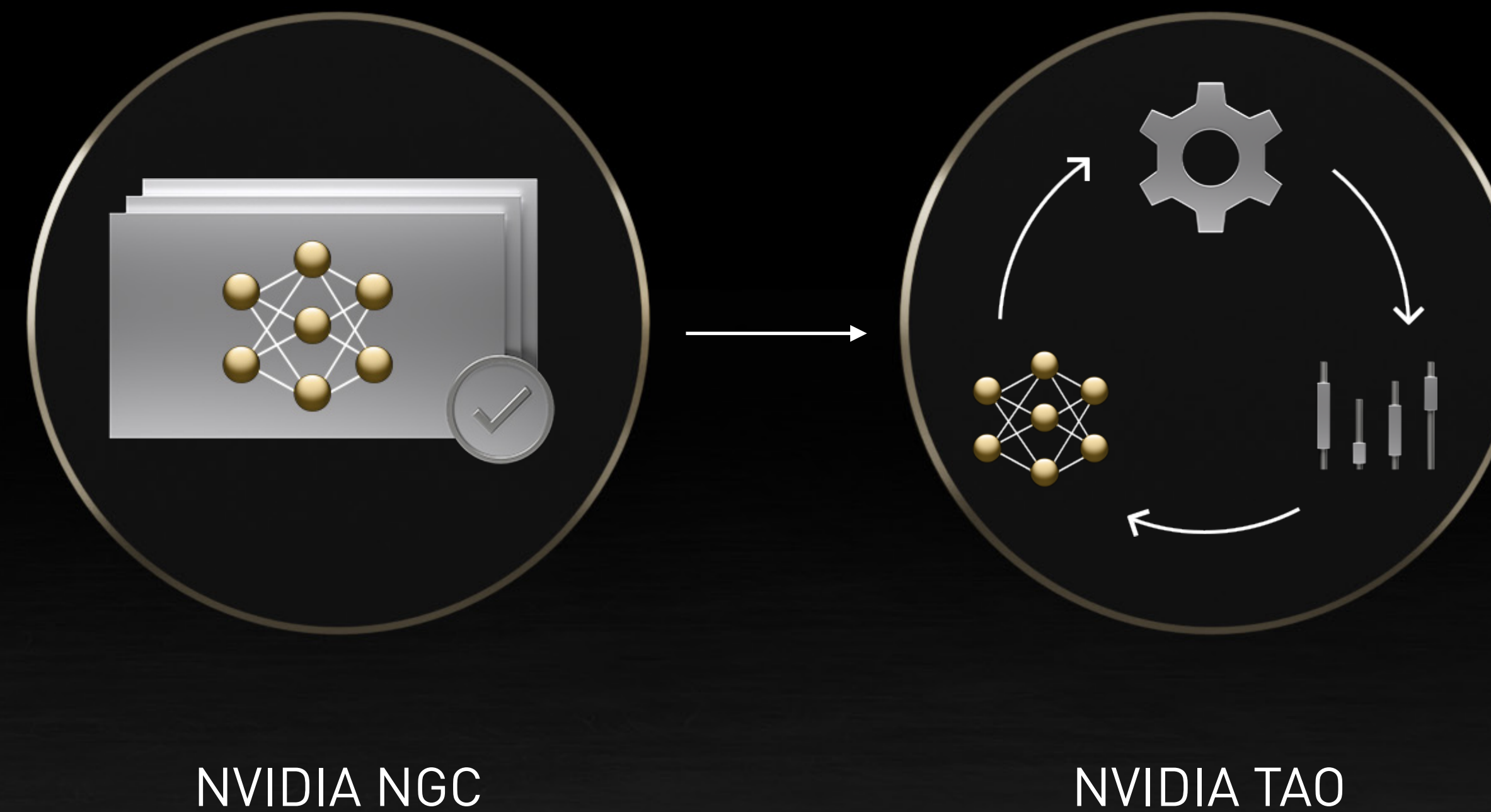3D point of regards (X, Y, Z) and gaze vector (theta and phi)

**ANNOUNCING**
# NVIDIA TAO FRAMEWORK
Train | Adapt | Optimize

Customize Pre-Trained Models for
Domain-Specific Applications

Federated Learning Enables Model Training
Collaboration while Protecting Data Privacy

Produce State-of-the-Art Models in Hours

NVIDIA NGC

NVIDIA TAO

**ANNOUNCING**
# NVIDIA FLEET COMMAND
## Securely Orchestrate AI Fleet at the Edge of the Network

Control and Manage Millions of AI-Powered Devices from Any Cloud
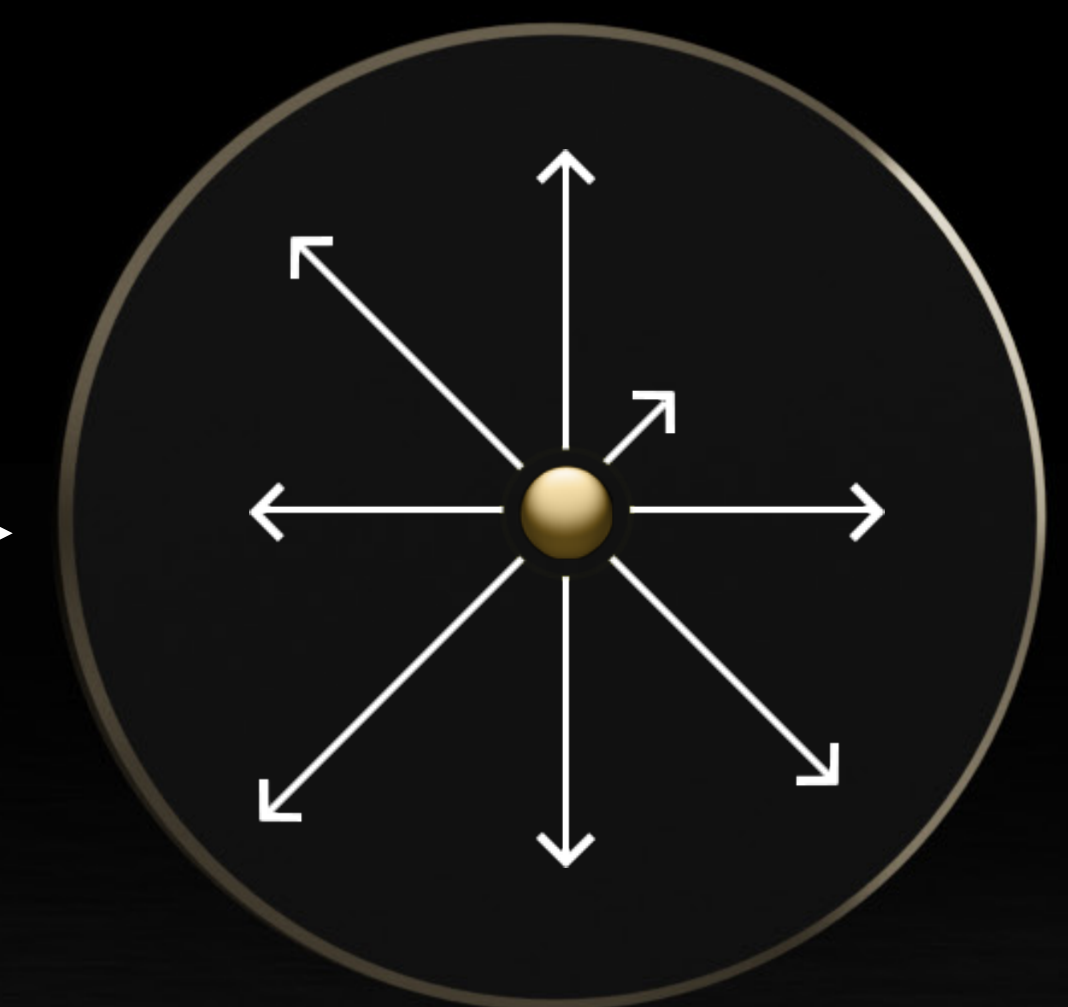
Secure from Boot, Attestation, Uplink and Downlink, to Confidential AI Enclave

Centrally Monitor Health and Remotely Fix Edge Systems

NVIDIA NGC          NVIDIA TAO          NVIDIA FLEET COMMAND

# NGC PRE-TRAINED MODELS
## Production-Quality AI Models

Trained by Experts for Enterprise Deployment

Credentials to Find Models You Trust

Continuously Updated to be State-of-the-Art

Adapt with NVIDIA TAO and Orchestrate
with NVIDIA Fleet Command

Reference AI Code Samples to Ease
Application Development

# NVIDIA JARVIS

## State-of-the-Art Conversational AI

GPU-Accelerated ASR, NLU, TTS

Interactive Performance

Customize with NVIDIA TAO

Orchestrate with NVIDIA Fleet Command

Scale-Out with NVIDIA Triton

Run in Every Cloud and at the Edge

## Settings

### ASR Model
None

### Translate
En -> Ja

[STOP]

## Jarvis transcription

We're also making tremendous progress in translation and now offer support for five languages.

You can see how fluid the translation is into Japanese.

And it's running in real time with under one hundred milliseconds latency for each sentence.

また、翻訳も飛躍的に進歩しており、現在5カ国語に対応しています。

日本語の翻訳の流動性がわかります。

そして、各文で 100 ミリ秒未満のレイテンシでリアルタイムで実行されています。

**NVIDIA MAXINE**

SOTA AI to Reinvent Virtual Collaboration

SDKs for Video, Audio, and Augmented Reality

AI-Face Codec 10x Lower Bandwidth vs H.264

Jarvis for Conversational AI

Deploy on Client, in Data Center and Every Cloud

Download Today: developer.nvidia.com/maxine

AI FACE CODEC

EYE CONTACT

MACHINE TRANSLATION

SPEECH RECOGNITION

that's so much better thanks to NVIDIA Maxine.

No es ideal. Ahora vamos a volver a activar todas las características de Maxine.

# ANNOUNCING NVIDIA TRITON INFERENCE SERVER

**MODELS**

Infinite AI Models
CV, ASR, NLU, TTS, RecSys

Multiple Frameworks
TensorFlow, Pytorch,
ONNX, TensorRT

Varying Services
Batch, Real-Time,
Streaming

**TRITON INFERENCE SERVER**

Any Model
Any Framework
Any Service Requirement
CPU or GPU

GPU

CPU

# THOUSANDS OF COMPANIES DOING COOL THINGS WITH NVIDIA AI

Analyze 225K
Network Events per Second

**BEST BUY**

Accurately Detect Diseases
in 145M Hearts per Year

GE Healthcare

Identify Trends in over 300B Pins
for Better Search Results

Tasteful Recommendations
from 600K Restaurants

Postmates

Personalized Playlists
for over 345M Listeners

Spotify®

Award-Winning Customer Care
Using Real-Time ASR

T·Mobile™

Real-Time Analytics
on 7B Packages per Year

UNITED STATES
POSTAL SERVICE

Intelligent Search with SOTA NLU
for 1.2B Users

WeChat

# NVIDIA DRIVE AV

**DATA COLLECTING, TESTING** — DRIVE Hyperion

**MAPPING, TRAINING** — DGX A100

**SIMULATION** — DRIVE Constellation

**AV DRIVING** — AV on DRIVE AGX

**AV WORLD MODEL** — IX on DRIVE AGX

# TOPS IS THE NEW HORSEPOWER

Orin

250

TOPS

Xavier

Parker

0

| 2018 | 2020 | 2022 |
|------|------|------|
| 1 TOPS | 30 TOPS | 254 TOPS |
| 7 SPECint | 17 SPECint | 25 SPECint |

SPECrate2017_int

# ANNOUNCING ORIN CENTRAL COMPUTER
## 1 Chip – 4 Domains



**CLUSTER**

FUSA Graphics

OS #1

**INFOTAINMENT**

Entertainment Graphics & Multimedia

OS #2

**PASSENGER INTERACTION & MONITORING**

FUSA CV, AR, & Conversational AI

OS #4

**AV WITH CONFIDENCE VIEW**

FUSA CV & AR

OS #3

ORIN VEHICLE COMPUTER

Centralized & Software Defined

Tightly Integrated with AV Applications

Containerized & Easy to Upgrade

SOTA Security and Functional Safety

**ANNOUNCING**
# HYPERION 8 AV PLATFORM
State-of-the-Art Advances for
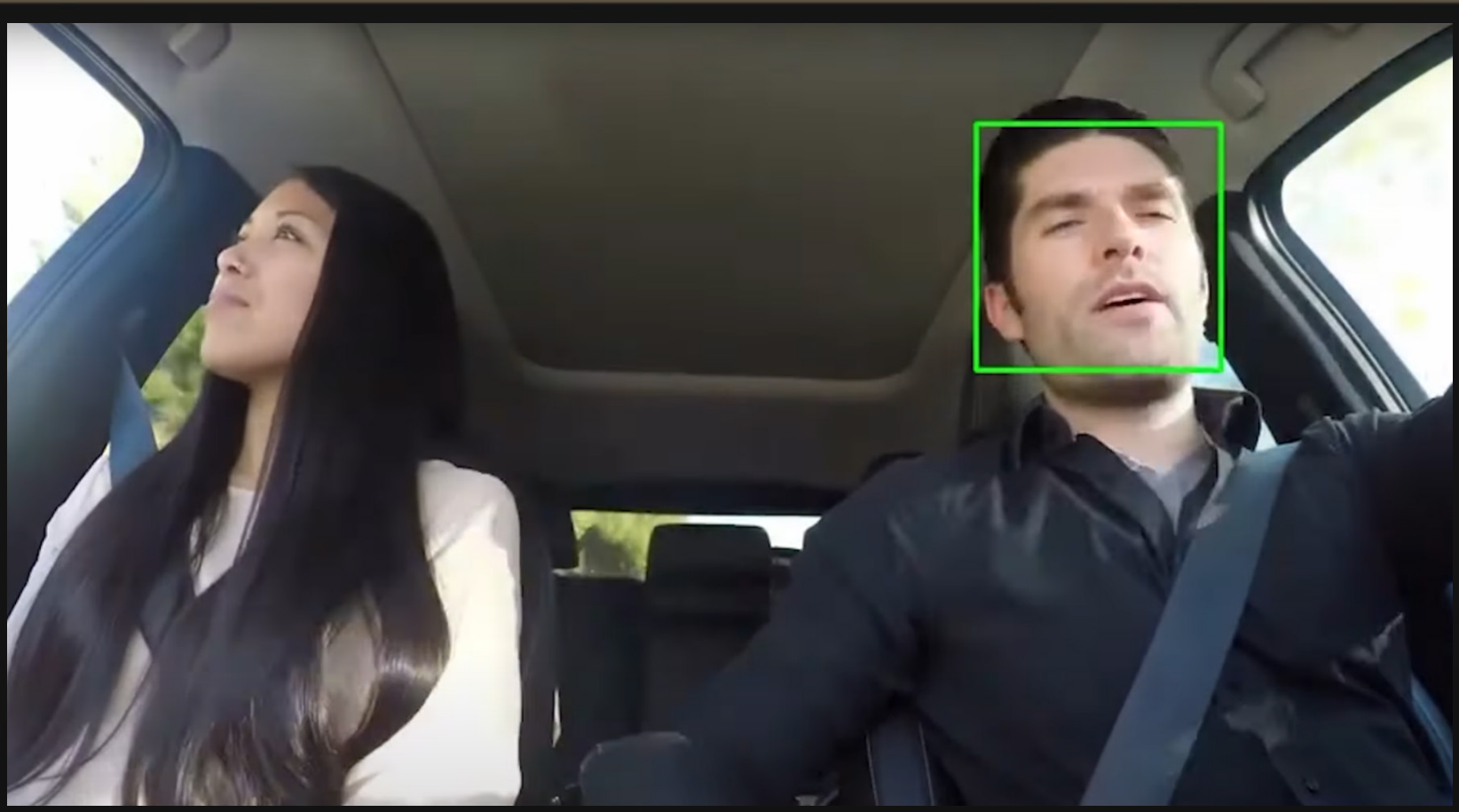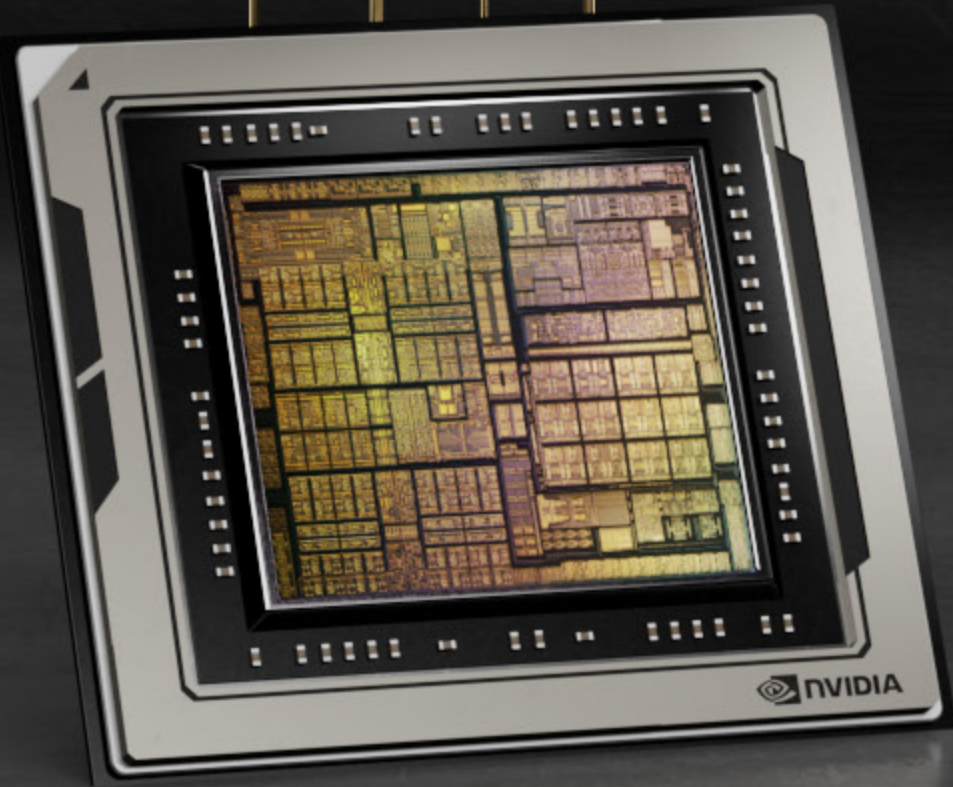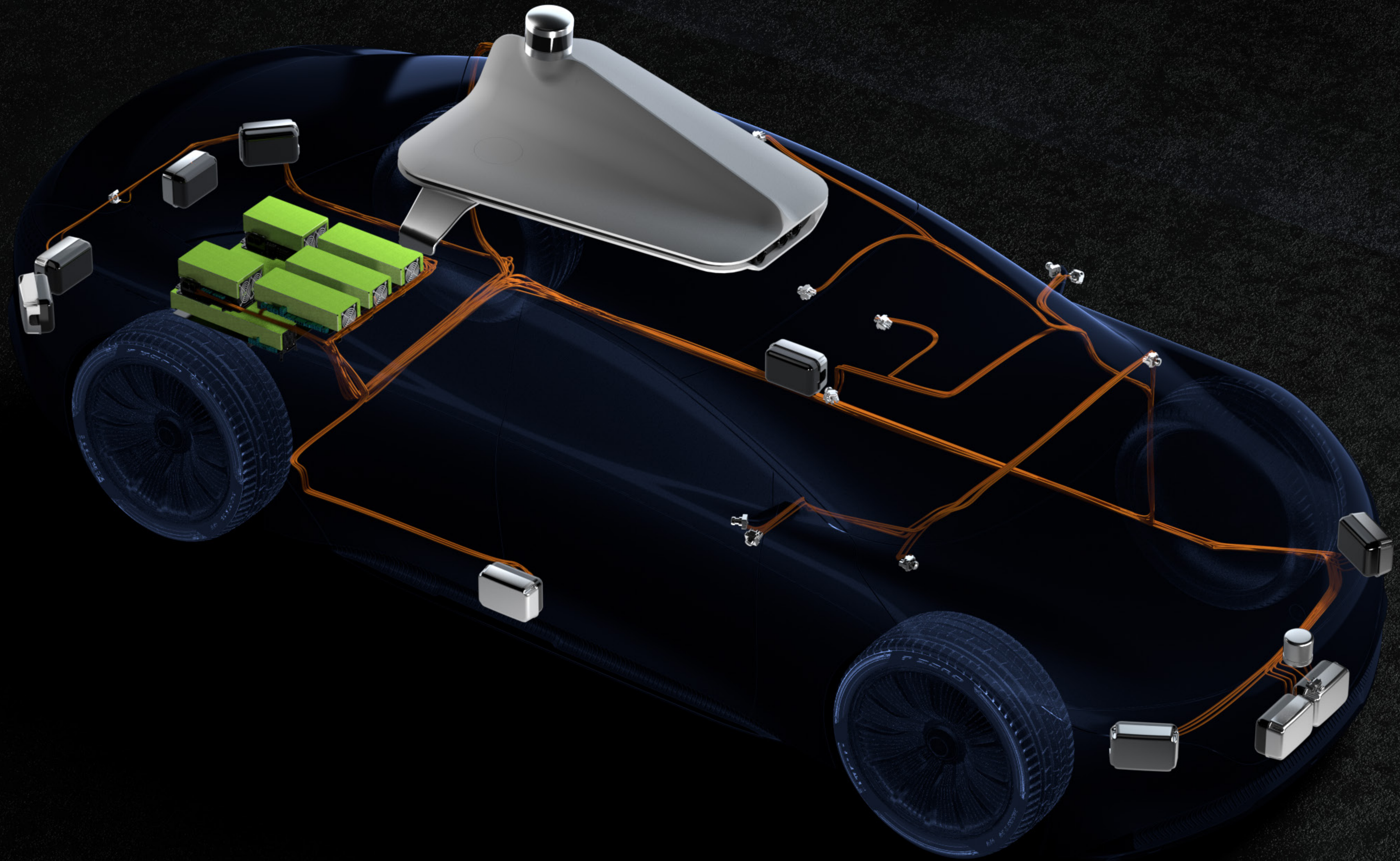Data Collection, Development and Testing

2x Orin AV Computer

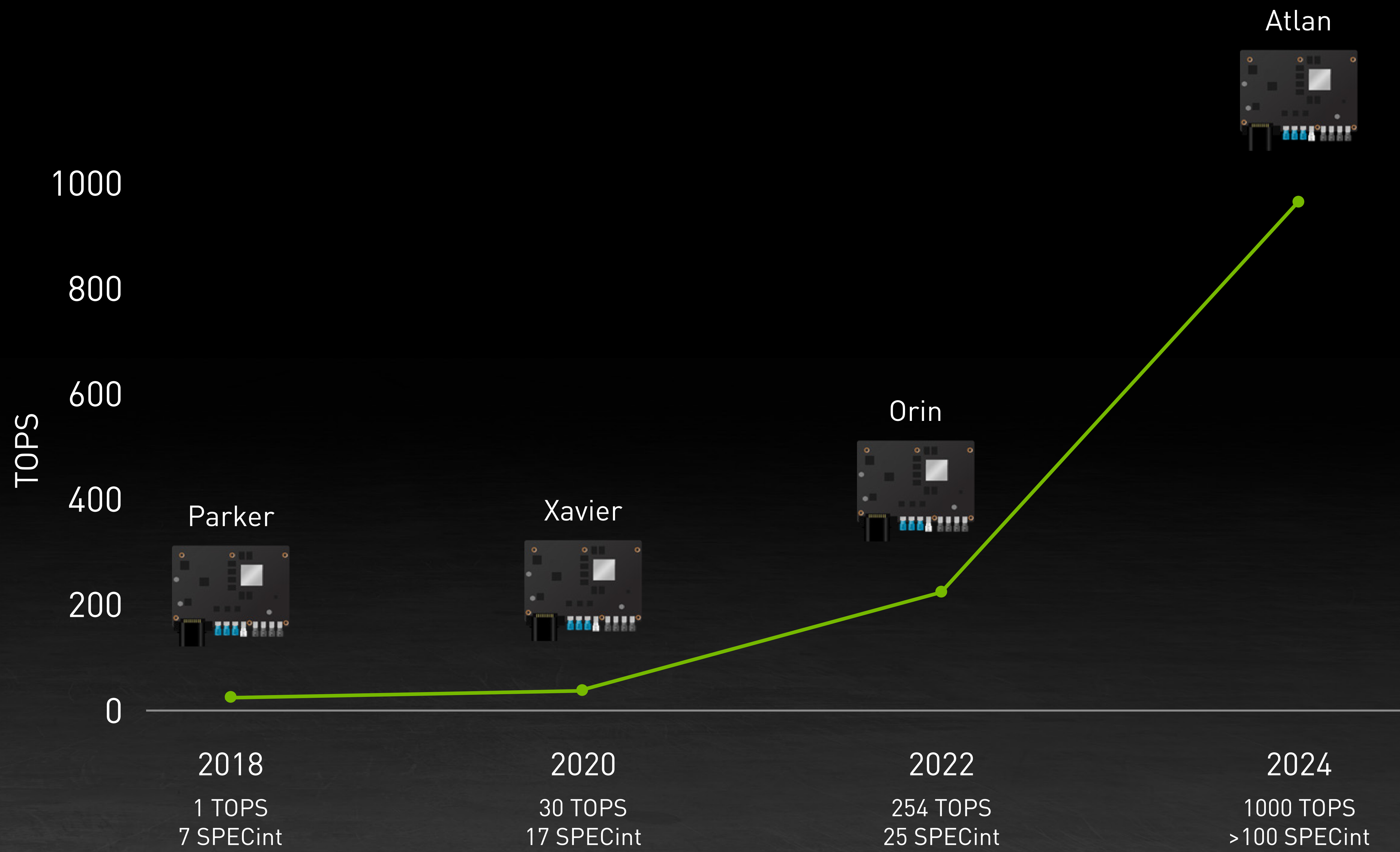1x Orin IX Computer

4x Orin + 4x MLNX 3D GT Data Recorder

Sensor Suite: 8 Cameras (8MP), 4 Fisheyes (3MP),
3 In-Cabin, 9 Radar, 2 Lidar

Source Access to AV & IX Software Repository

OTA Ready

**TOPS IS THE NEW HORSEPOWER**

TOPS

1000

800

600

400

200

0

| 2018 | 2020 | 2022 | 2024 |
|------|------|------|------|
| 1 TOPS | 30 TOPS | 254 TOPS | 1000 TOPS |
| 7 SPECint | 17 SPECint | 25 SPECint | >100 SPECint |

Parker

Xavier

Orin

Atlan

SPECrate2017_int

# A NEW BREED OF TECHNOLOGY COMPANIES

# THE WORLD'S BIG BRANDS HAVE GIANT OPPORTUNITIES

# TRANSFORMING 10 TRILLION MILES A YEAR INTO A SERVICE

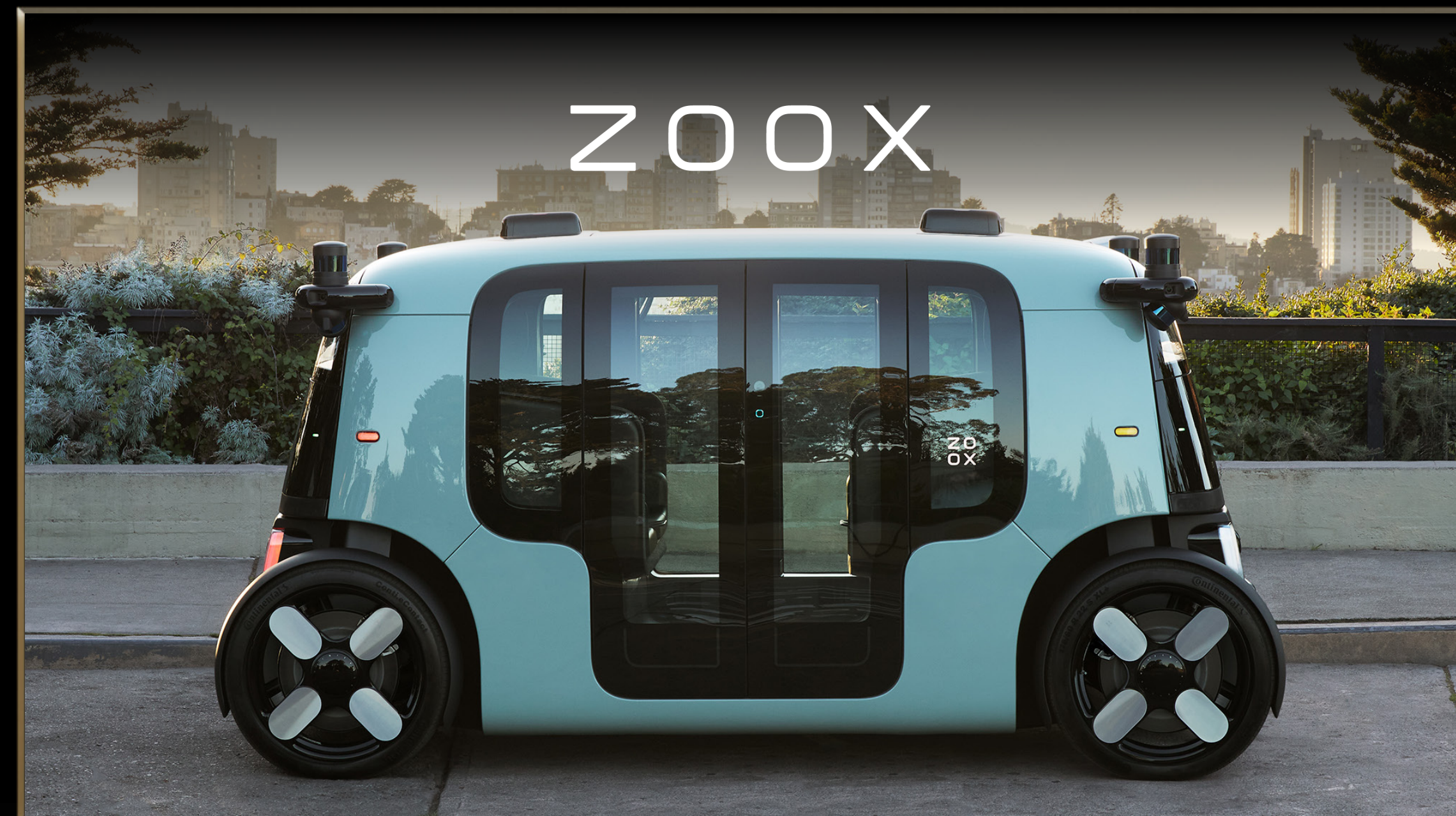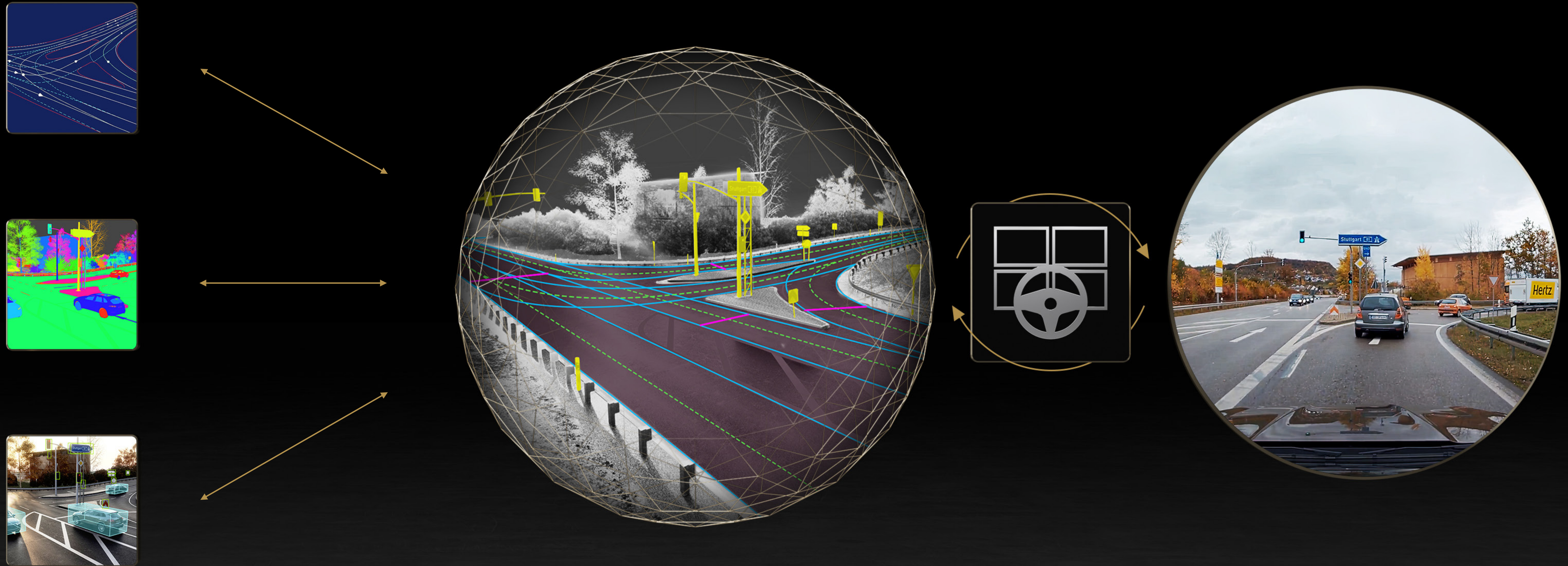# THE INTERNET OF ATOMS

NVIDIA DRIVE DIGITAL TWIN IN OMNIVERSE

**ANNOUNCING NVIDIA DRIVE SIM POWERED BY OMNIVERSE**

PHYSICALLY ACCURATE SENSORS    CLOUD-NATIVE WORKFLOW    SCALABLE ARCHITECTURE    OPEN | MODULAR | EXTENSIBLE
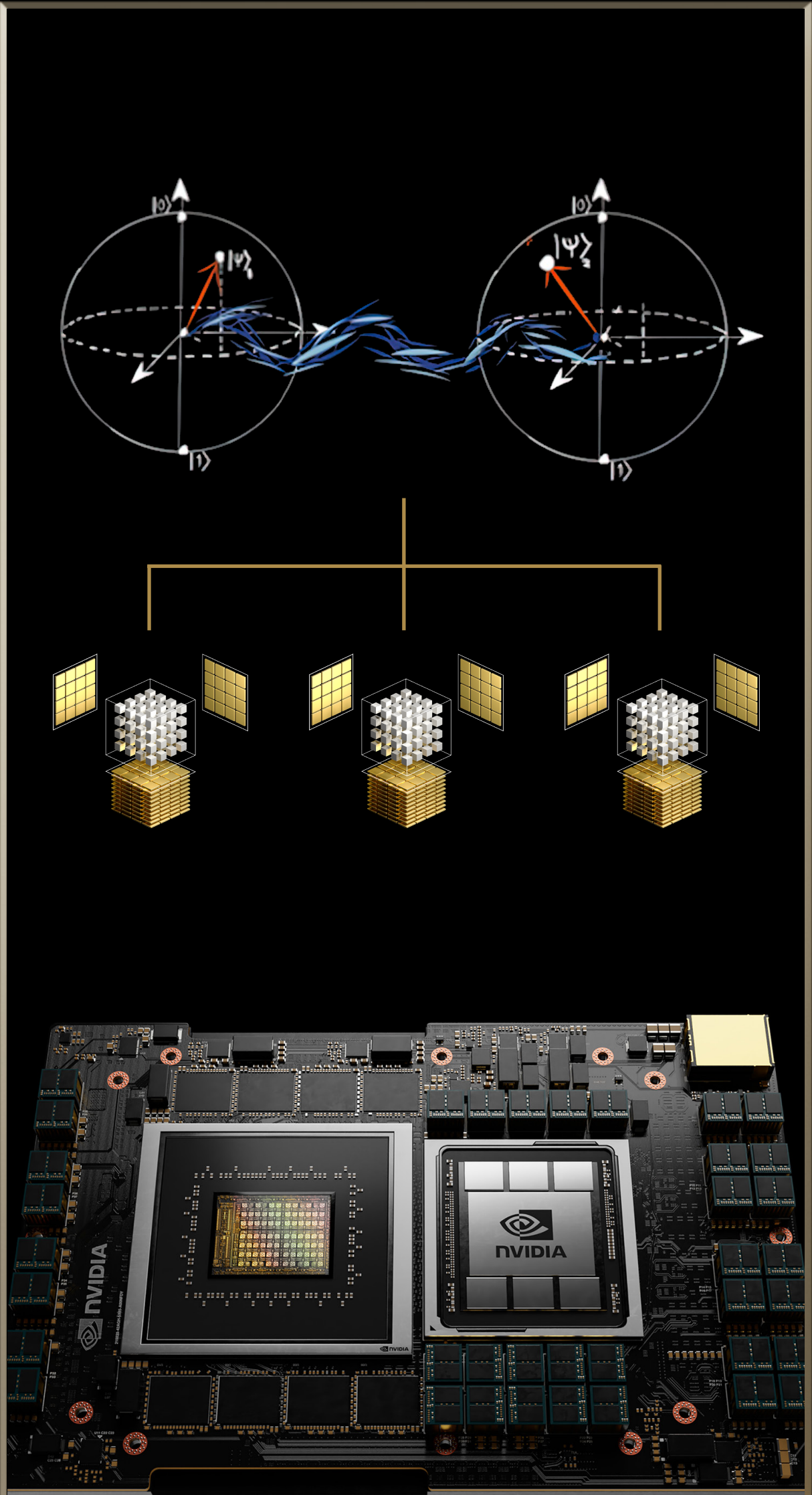
# NEW NVIDIA TECHNOLOGIES



**Omniverse**
**Isaac Sim**

**Megatron | Drug Discovery**
**Quantum Computing**
**DGX | Grace**
**BlueField-3 | DOCA 1.0**

**Jarvis | Merlin**
**Maxine | Morpheus**
**NVIDIA AI**
**EGX Aerial 5G**

**DRIVE Sim**
**Hyperion 8**
**Atlan**
**Orin**

PRODUCT OF ITALY

TRAIN | ADAPT | OPTIMIZE

Here we are

TAO
EXTRA VIRGIN
OLIVE OIL

Mado olive oil is
luxury you can drizzle.
Punchy and fresh, it's made
from 100% Italian olives,
all picked by hand
and pressed the same
day. It's an olive oil so
tasty it can spirit any dish
straight to gastronomic heaven.

Also try our olive oil crushed with fresh lemons.

Produced in Italy for Mado Ltd London UK