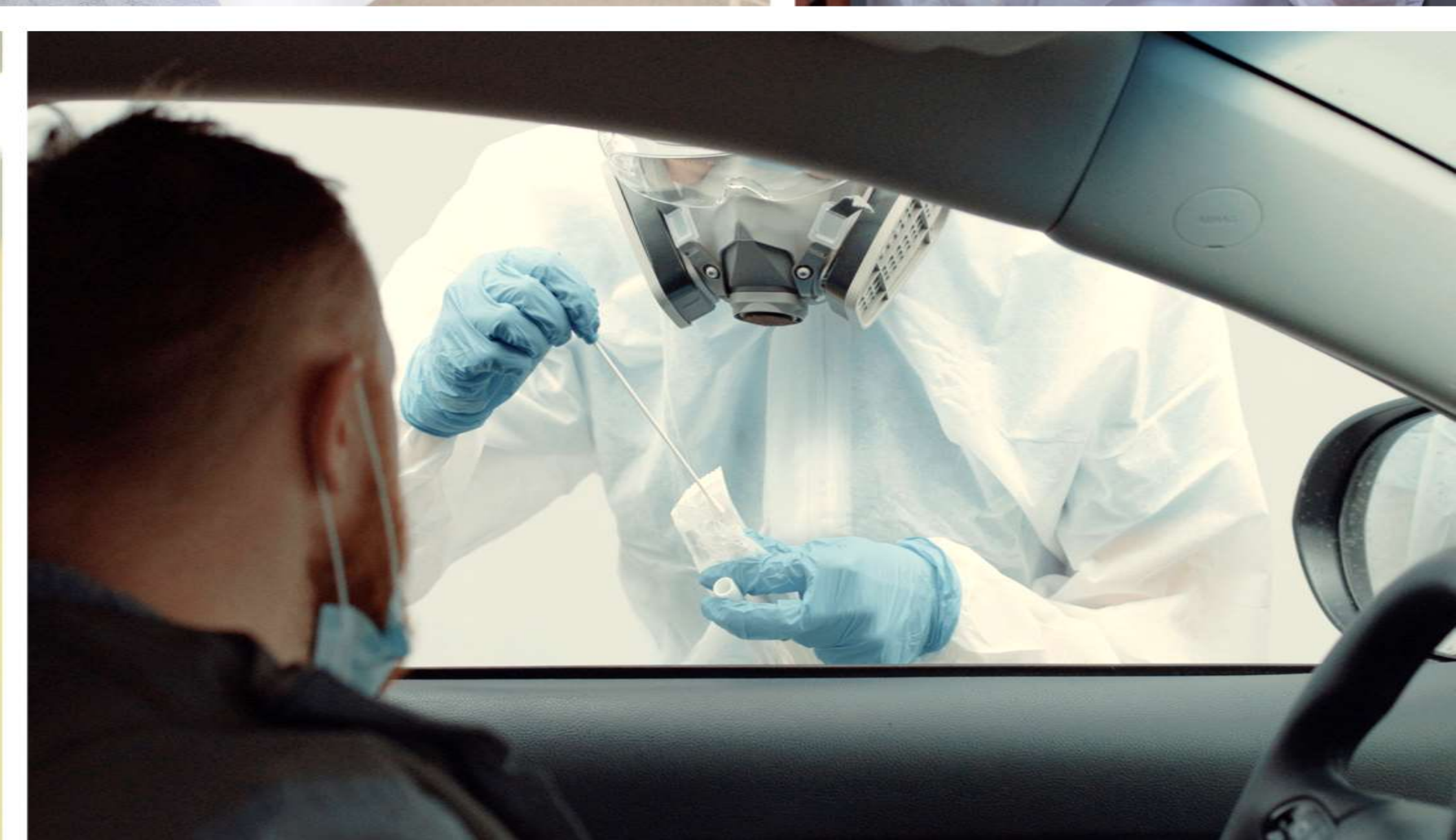
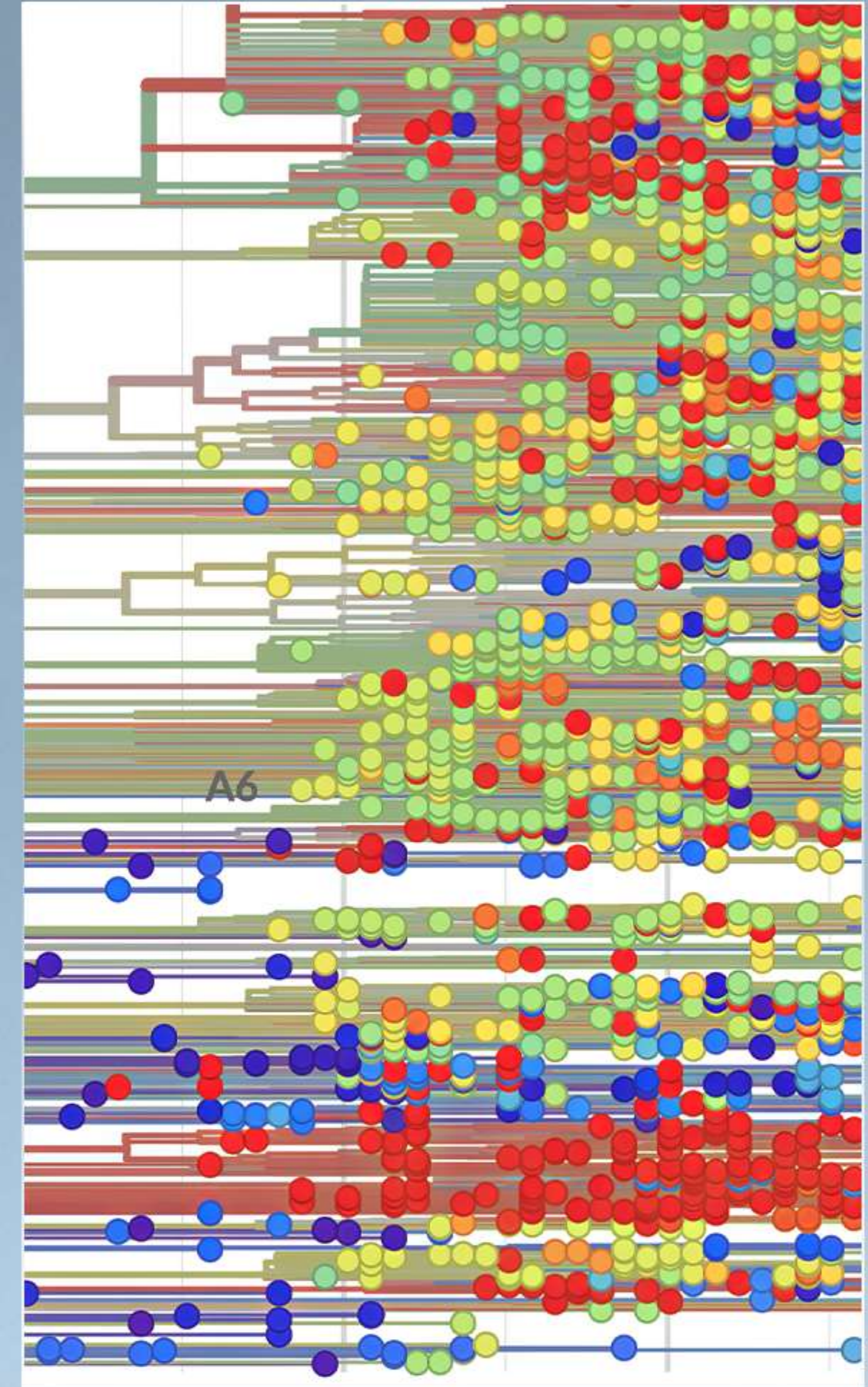




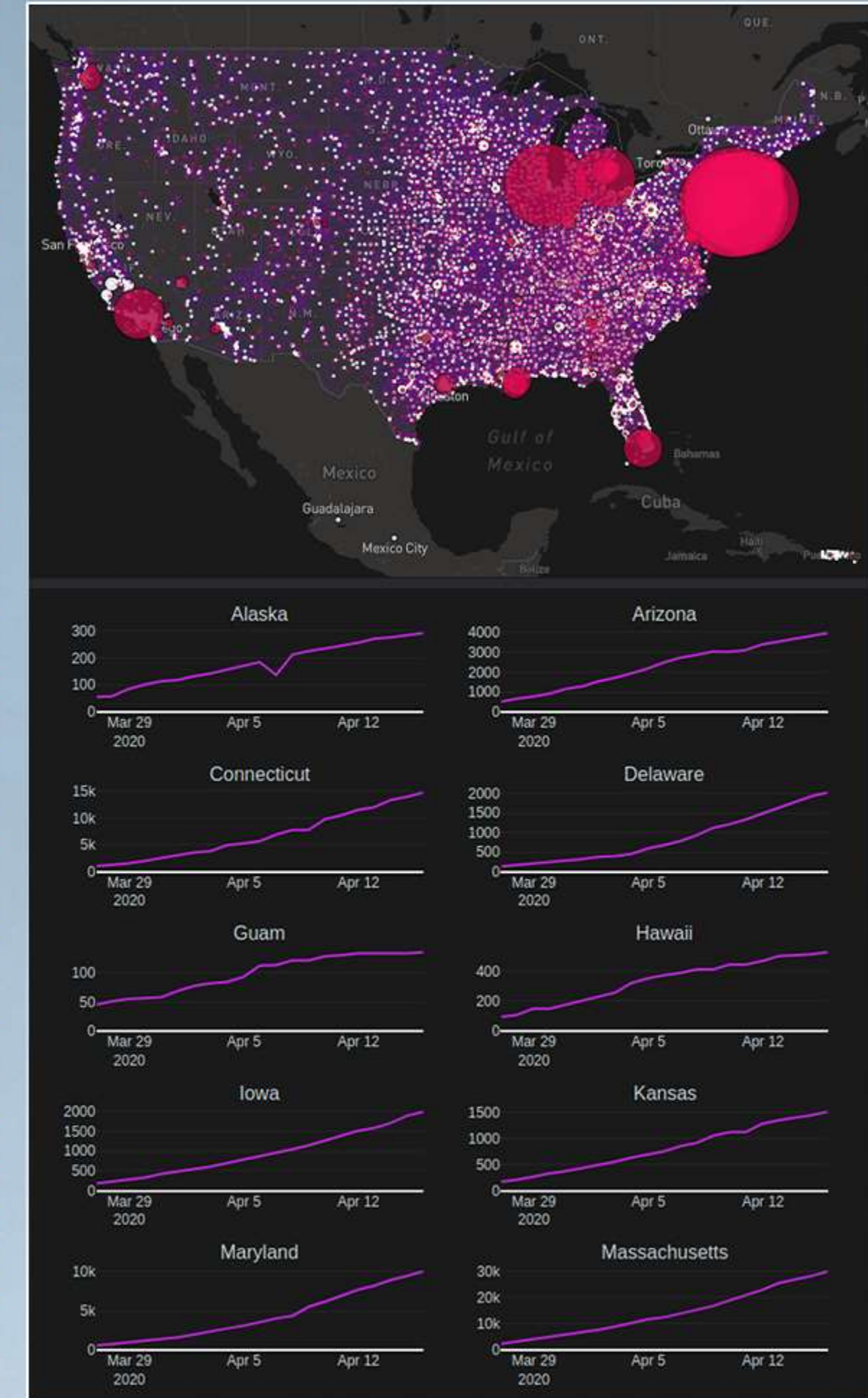
nVIDIA



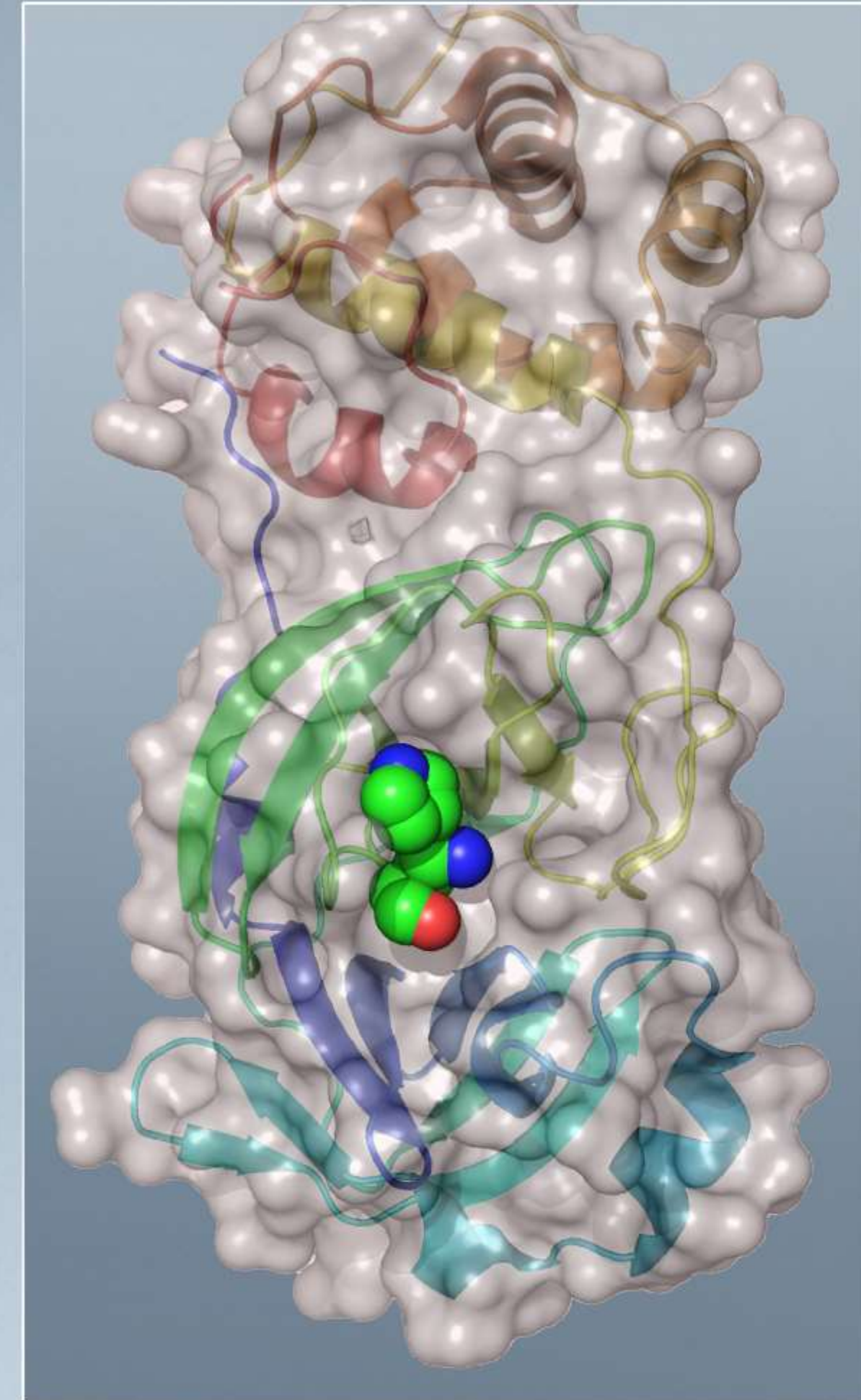
NVIDIA の COVID-19 との戦い



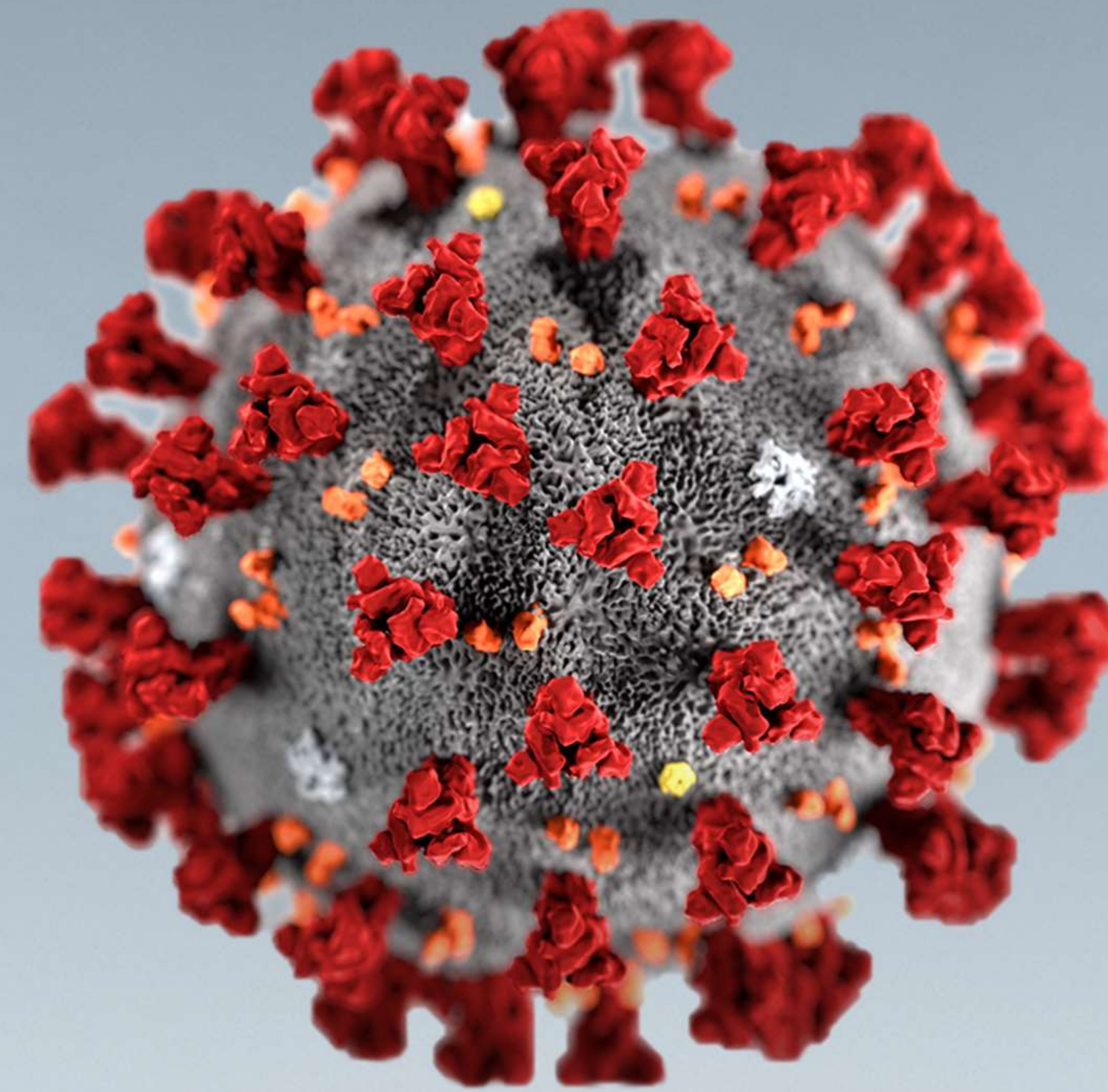
Oxford Nanopore
ウイルスのゲノムを
7 時間で解析



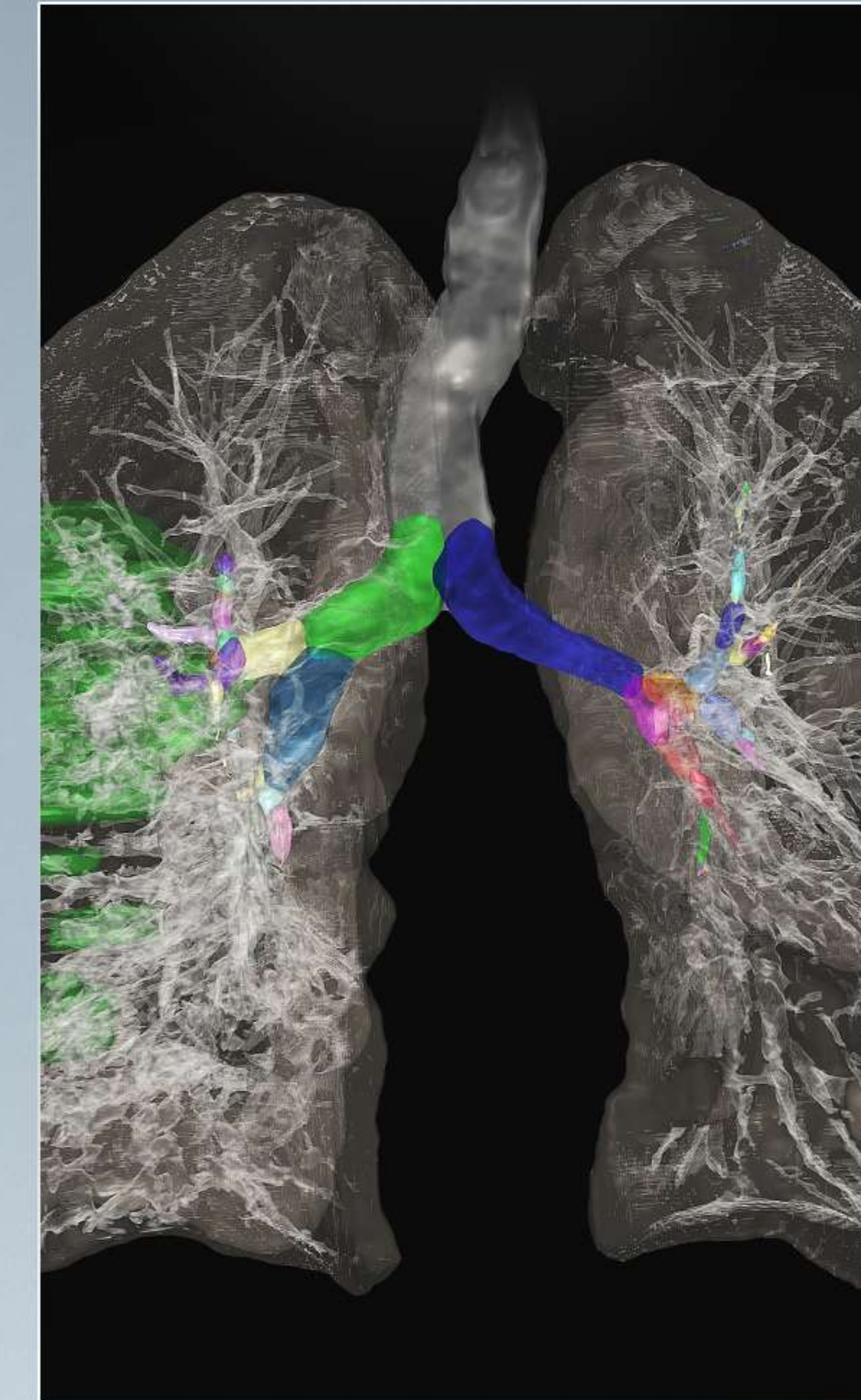
Plotly, NVIDIA
リアルタイムの
感染率分析



ORNL, Scripps
10 億の薬物化合物の
スクリーニングを
1 年から 1 日に短縮



Structura, NIH, UT Austin
CryoSPARC
初のウイルス スパイク
プロテインの 3 次元構造



NIH, NVIDIA
AI による COVID-19 の分類



Kiwibot
ロボットによる医療品の配達



Whiteboard Coordinator
AI による体温上昇の
スクリーニング システム

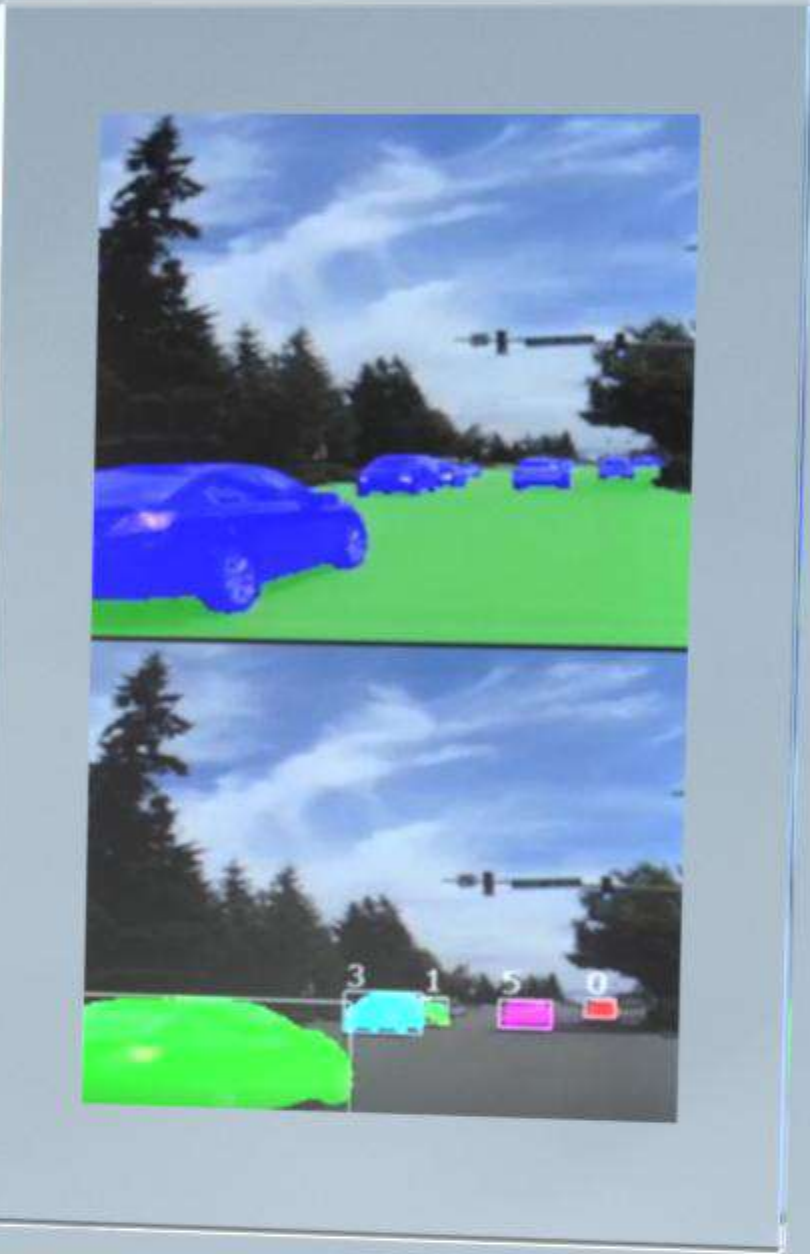
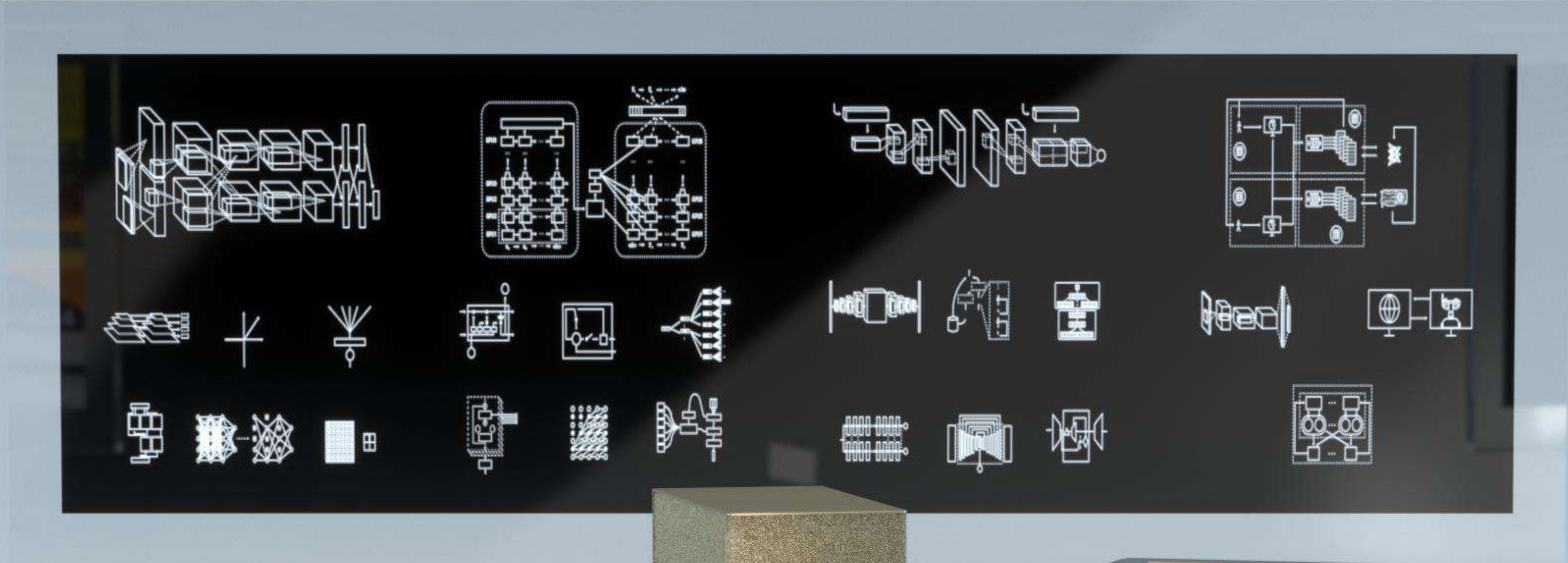
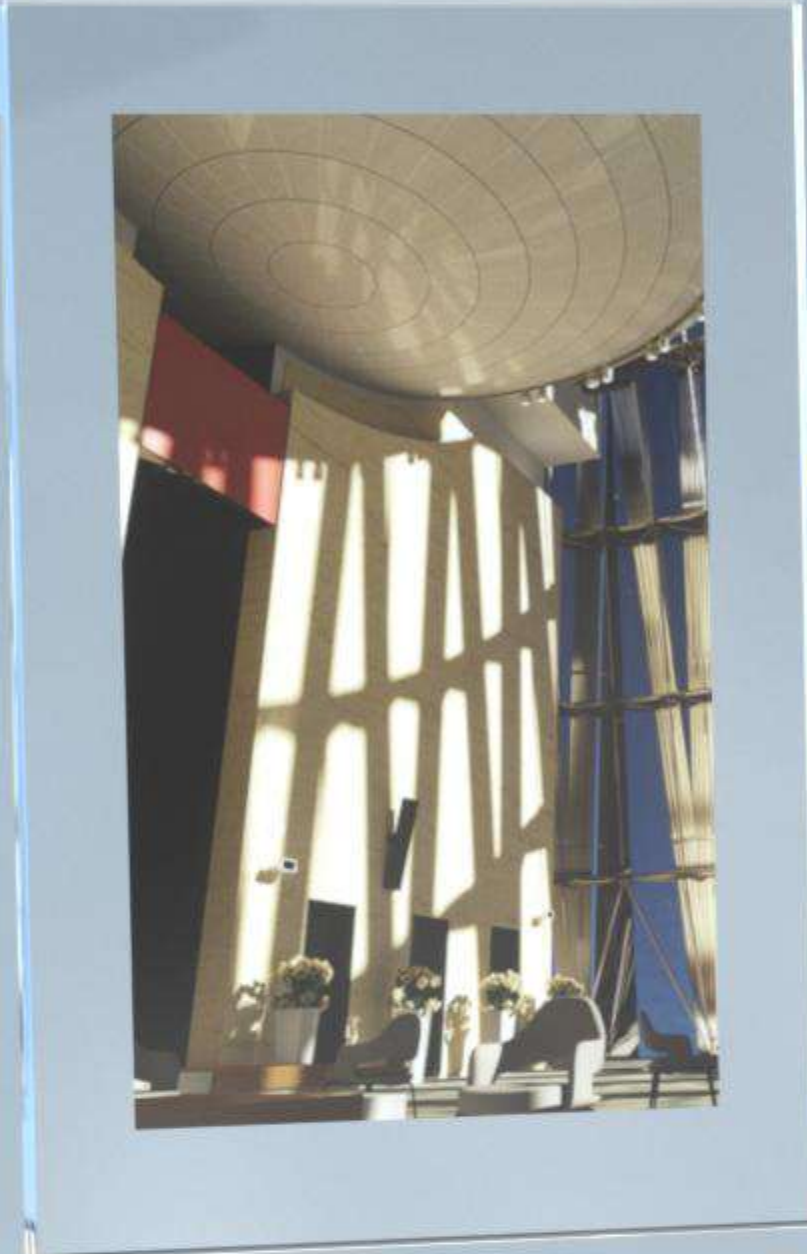
封じ込め

緩和

治療

追跡と監視

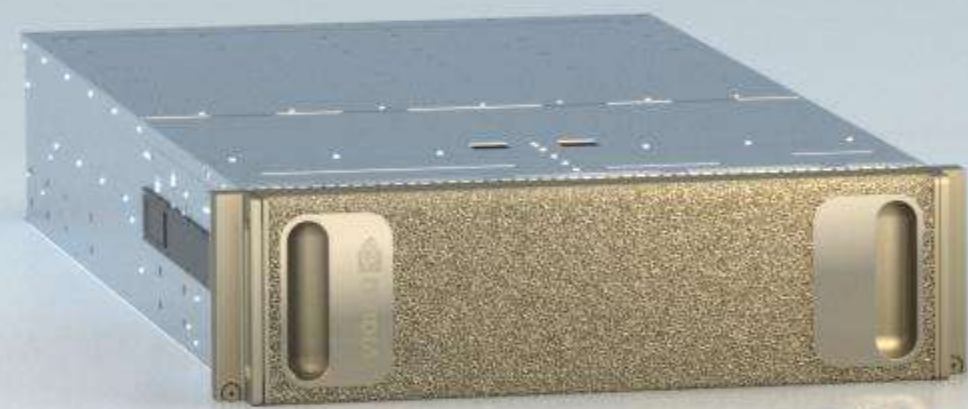
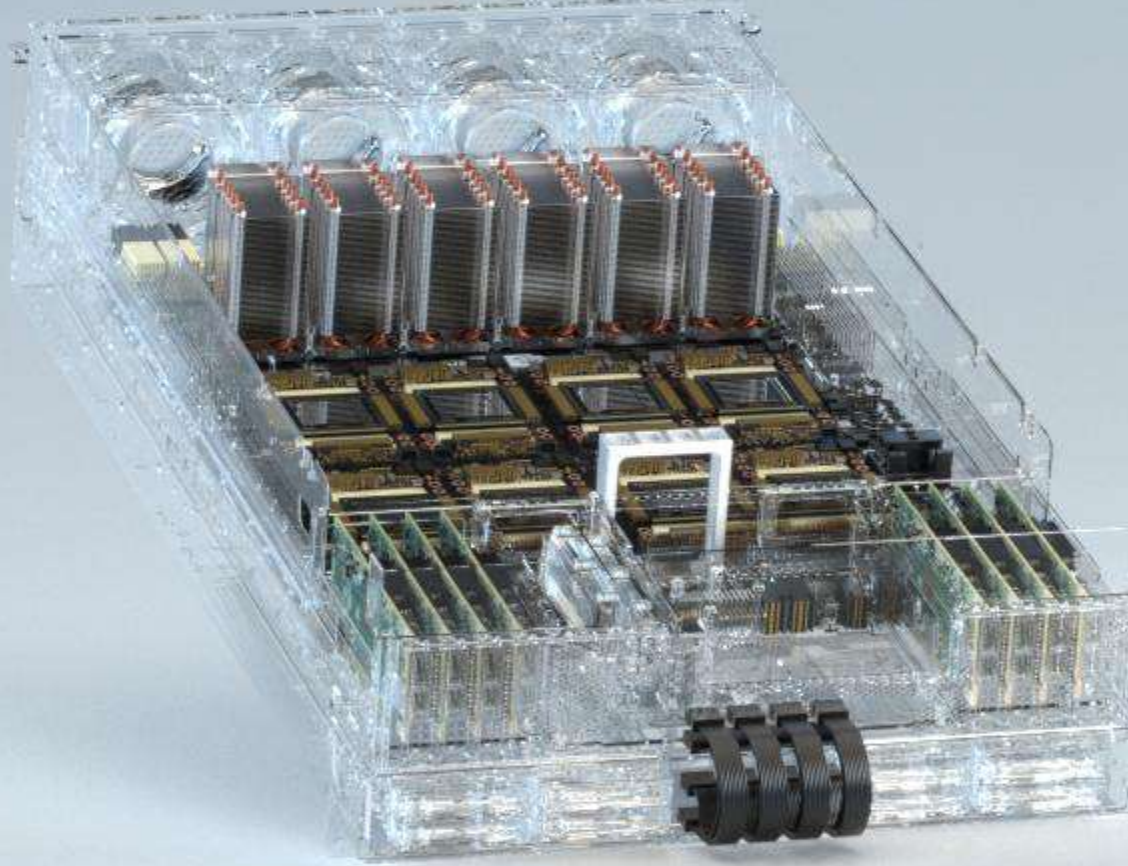
私たちの時代のダヴィンチのためのコンピューティング



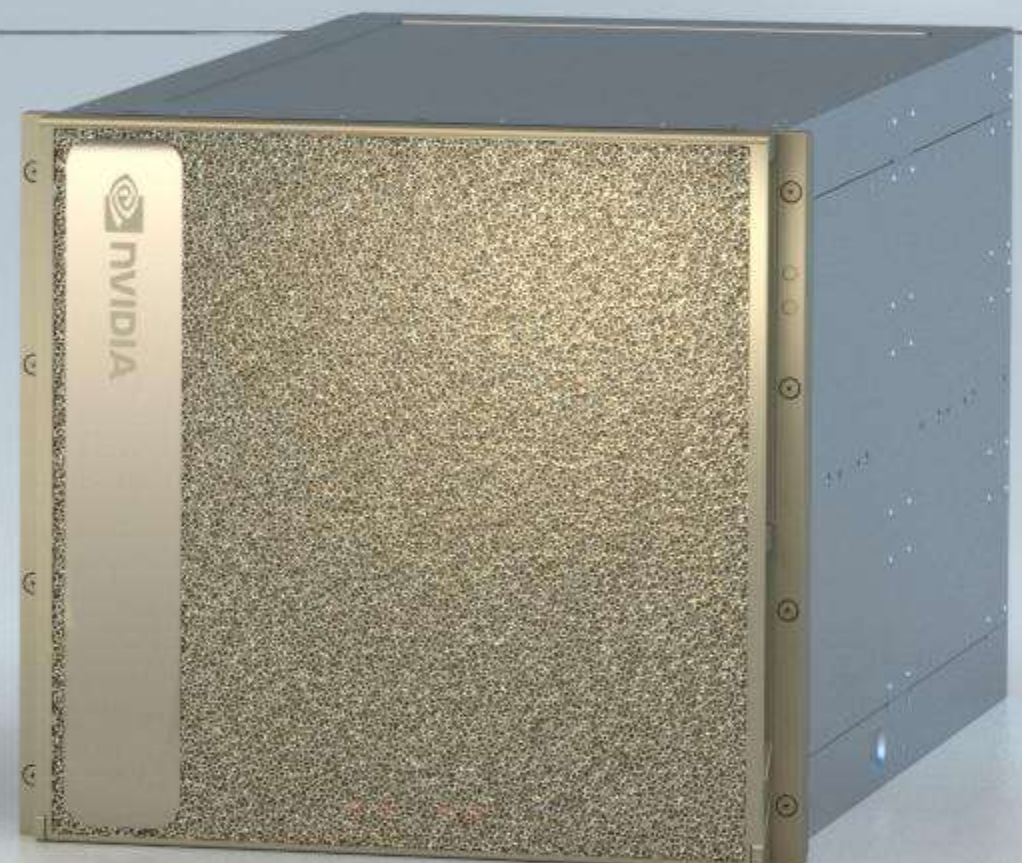
RTX



HGX



DGX



EGX

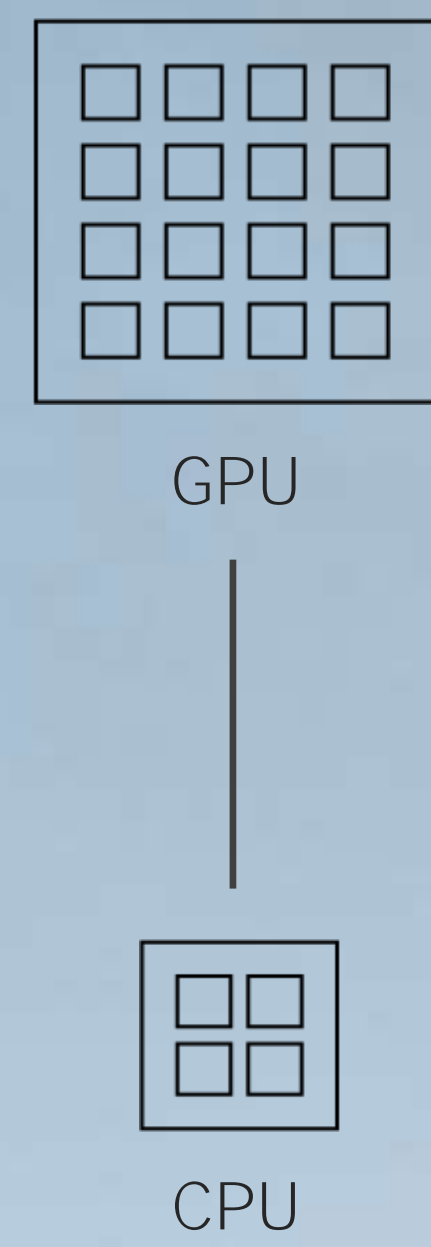


AGX

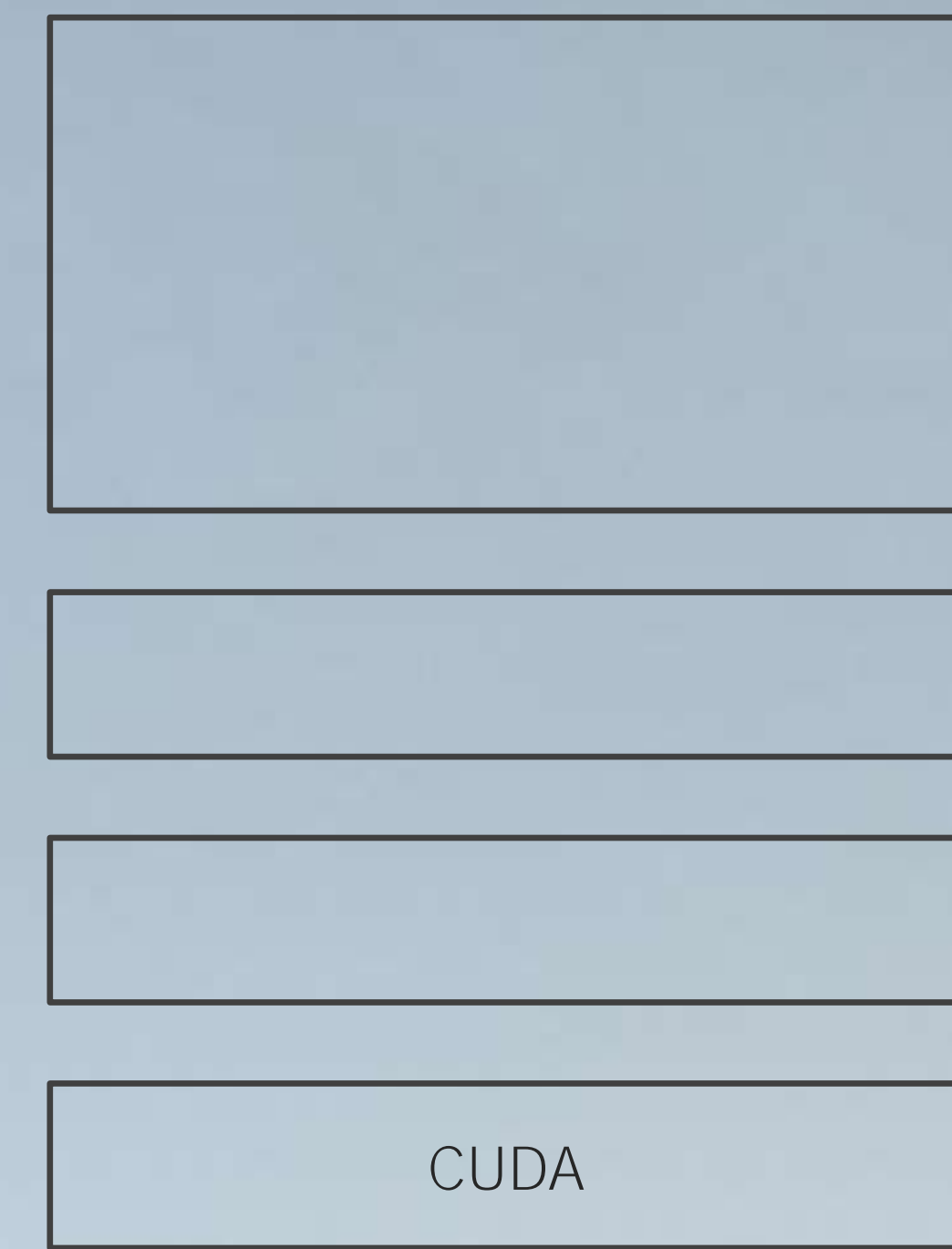




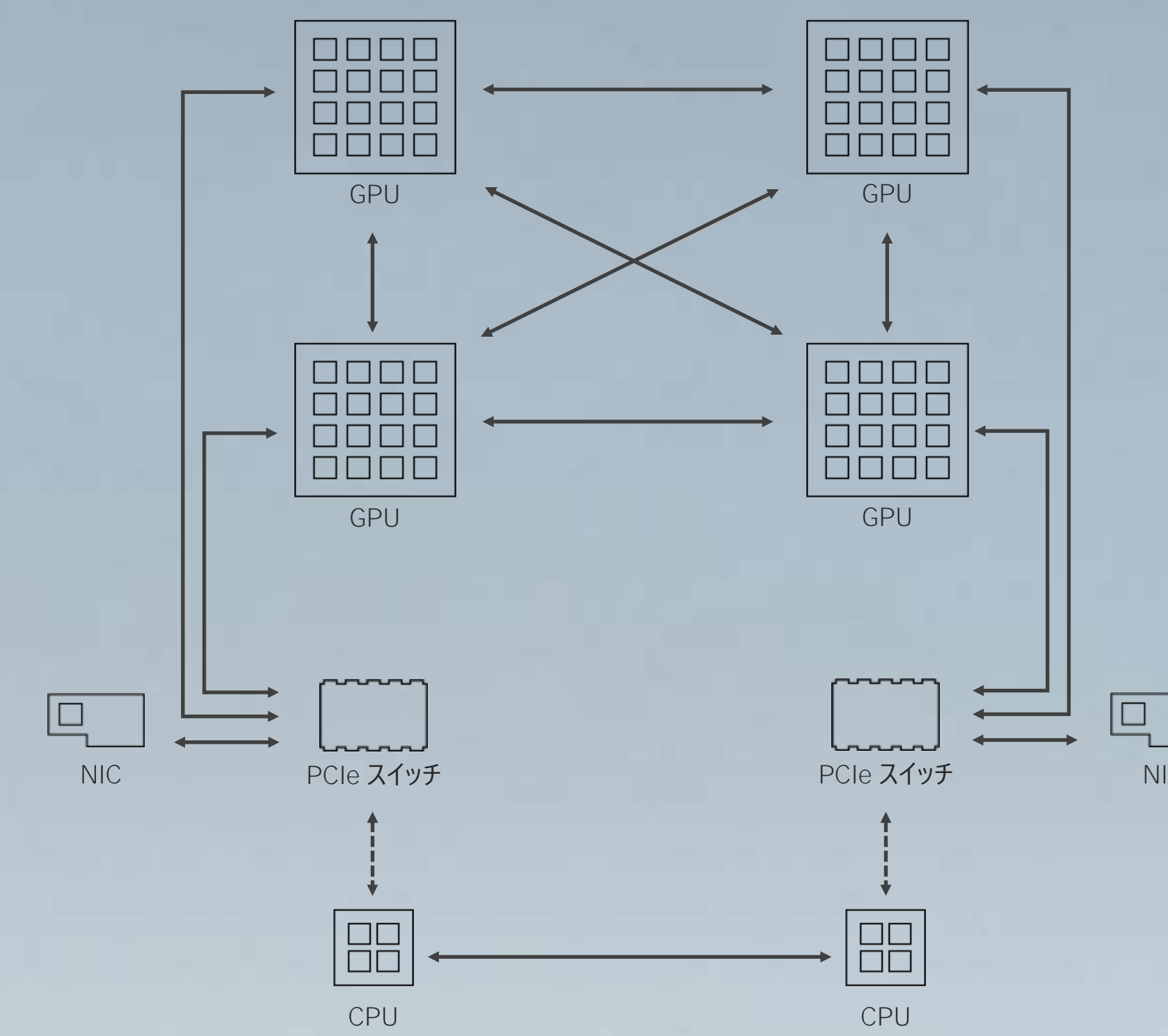
アクセラレーテッド コンピューティングの 25 年



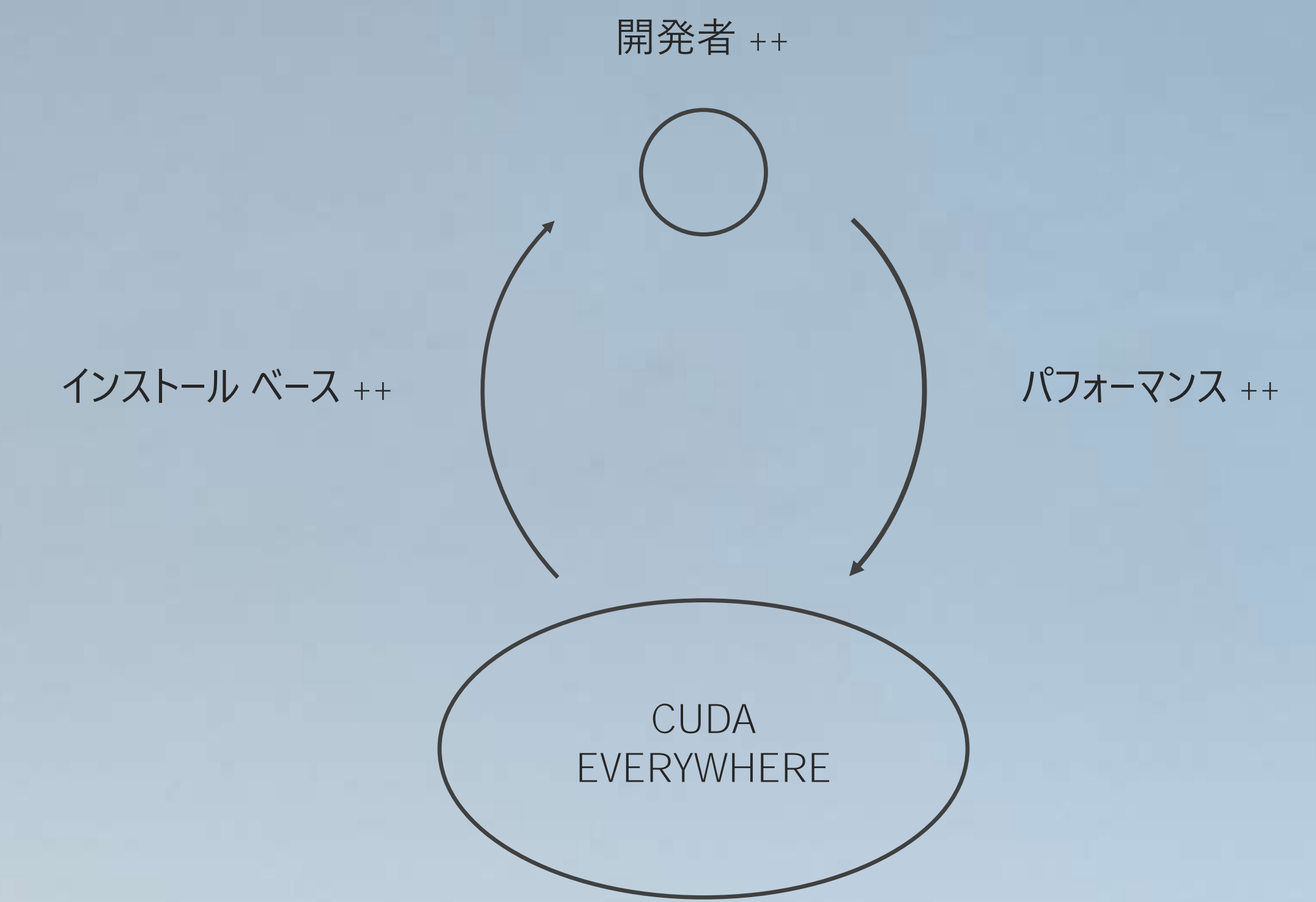
X-FACTOR スピードアップ



フル スタック

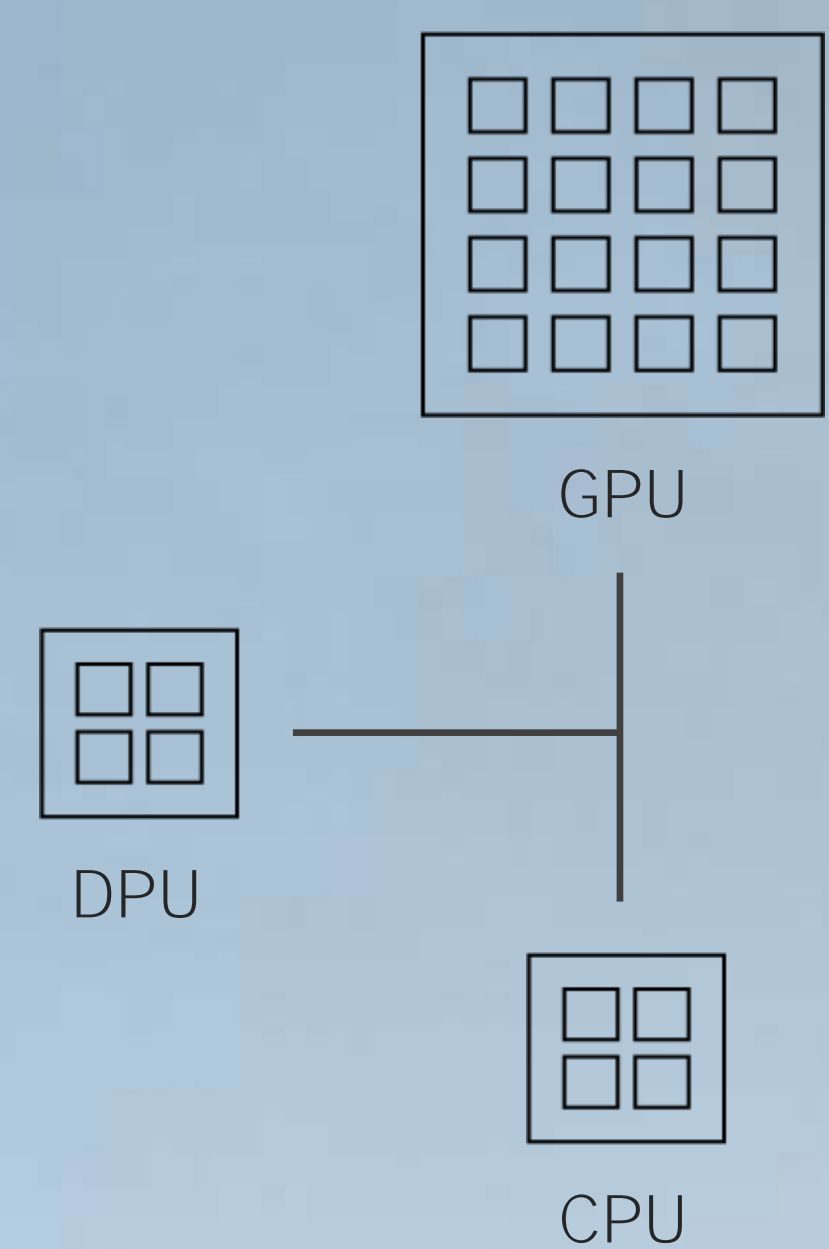


システム



1 つのアーキテクチャ

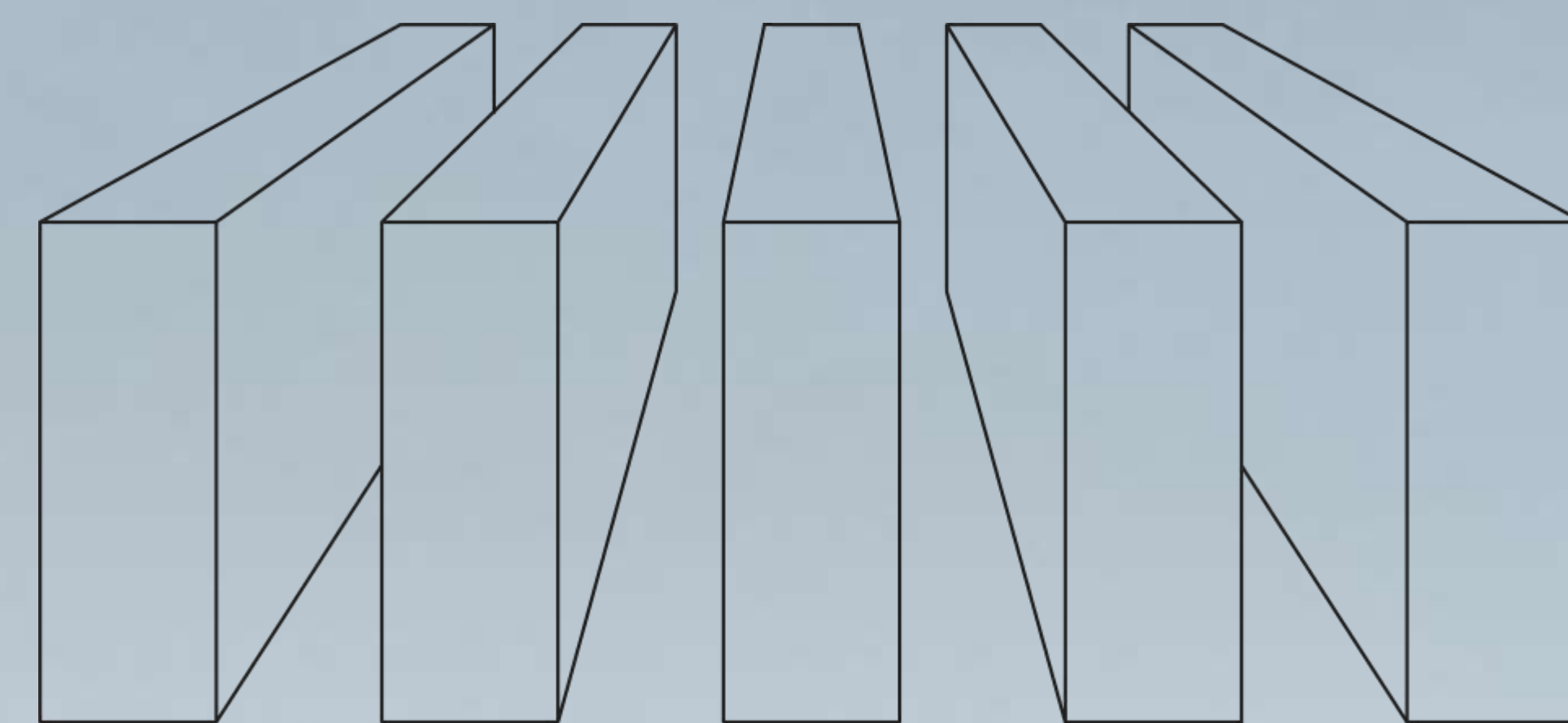
データセンター規模のアクセラレーテッド コンピューティング



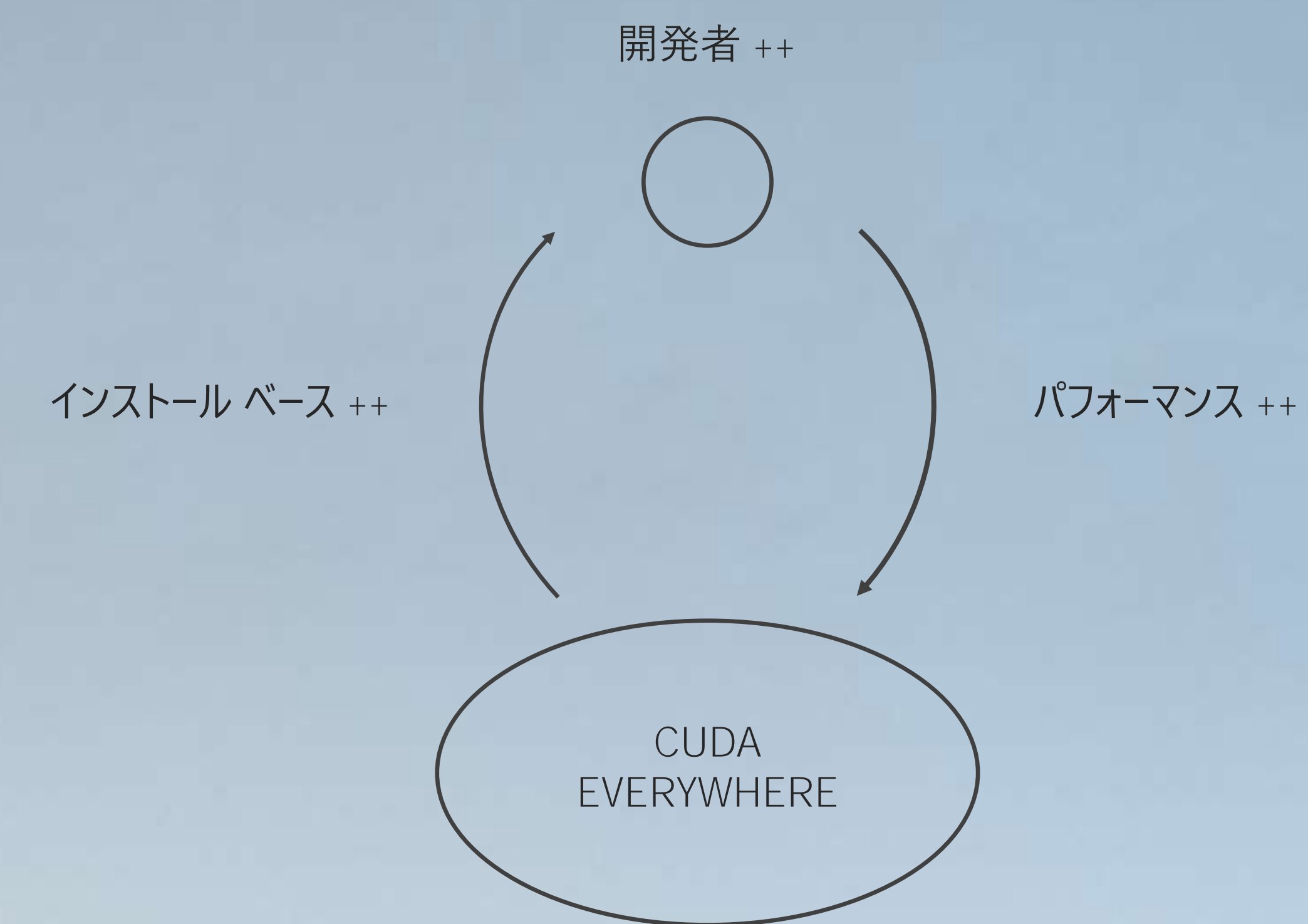
X-FACTOR スピードアップ



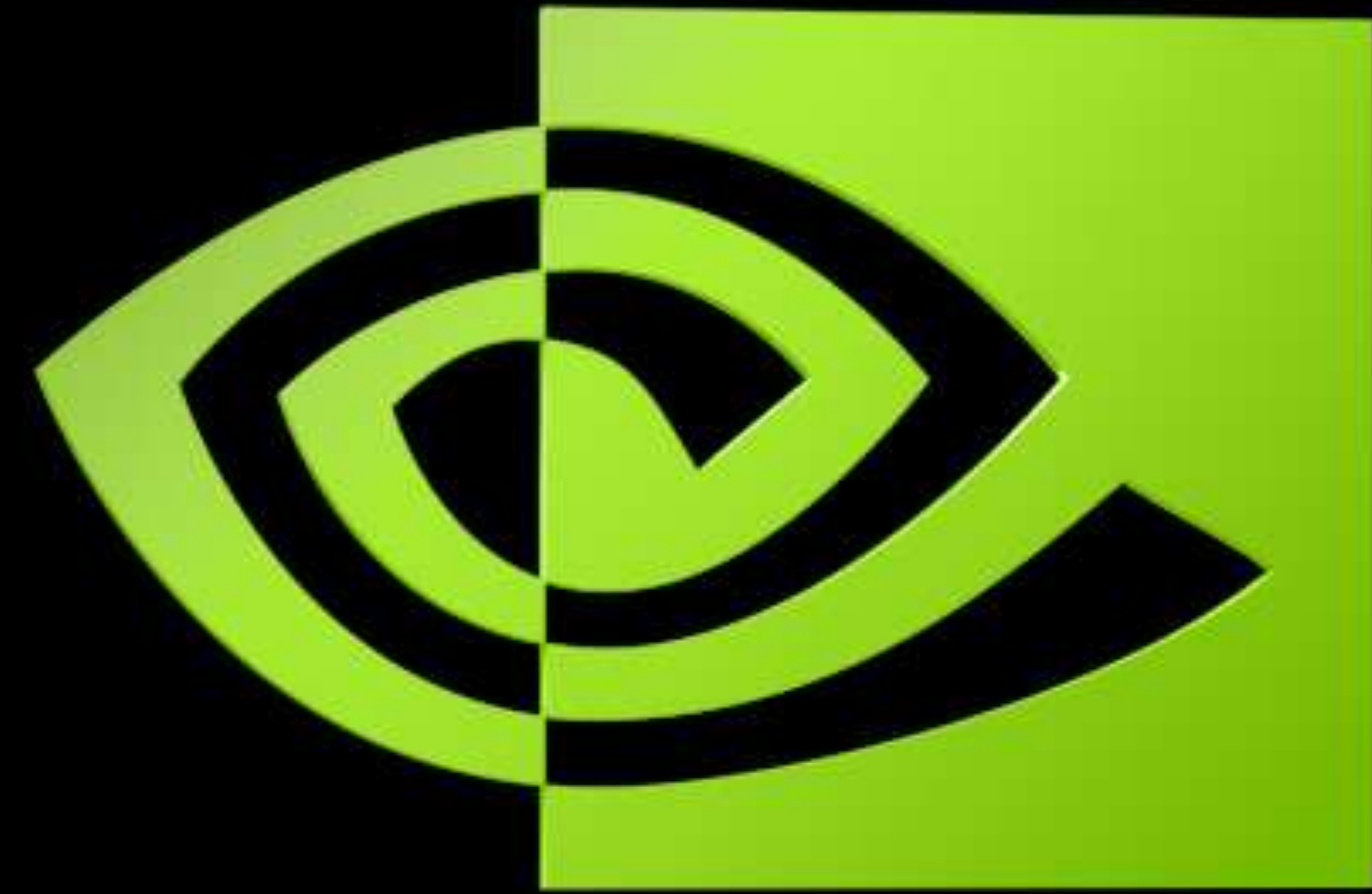
フル スタック



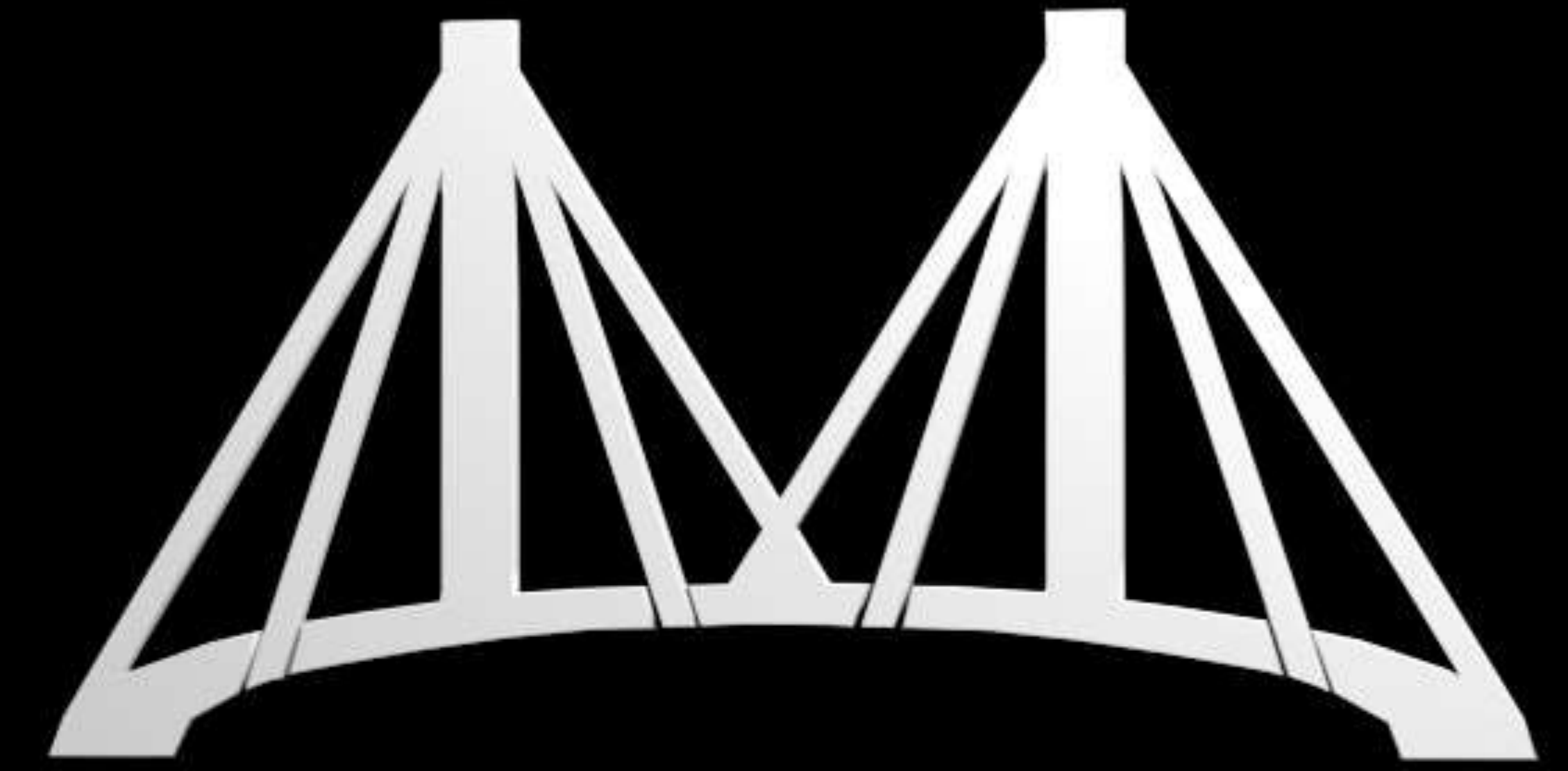
システム



1 つのアーキテクチャ



nVIDIA®



Mellanox®
TECHNOLOGIES

開発者にとって素晴らしい年

50
新しい SDK

1800 万
開発者



2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020 YTD

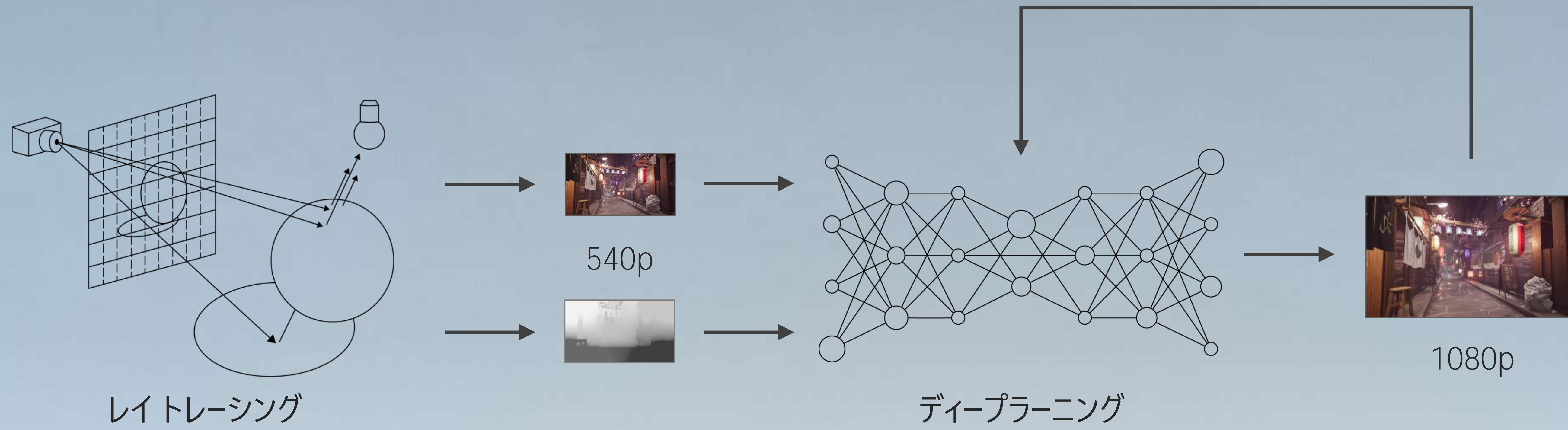
2018年8月13日、SIGGRAPHでNVIDIAはRTXを発表、
コンピュータグラフィックスの新時代の幕が上がる。

「30年以上前から約束されていたリアルタイムレイトレーシングが、
予想より10年早くやって来ました。」

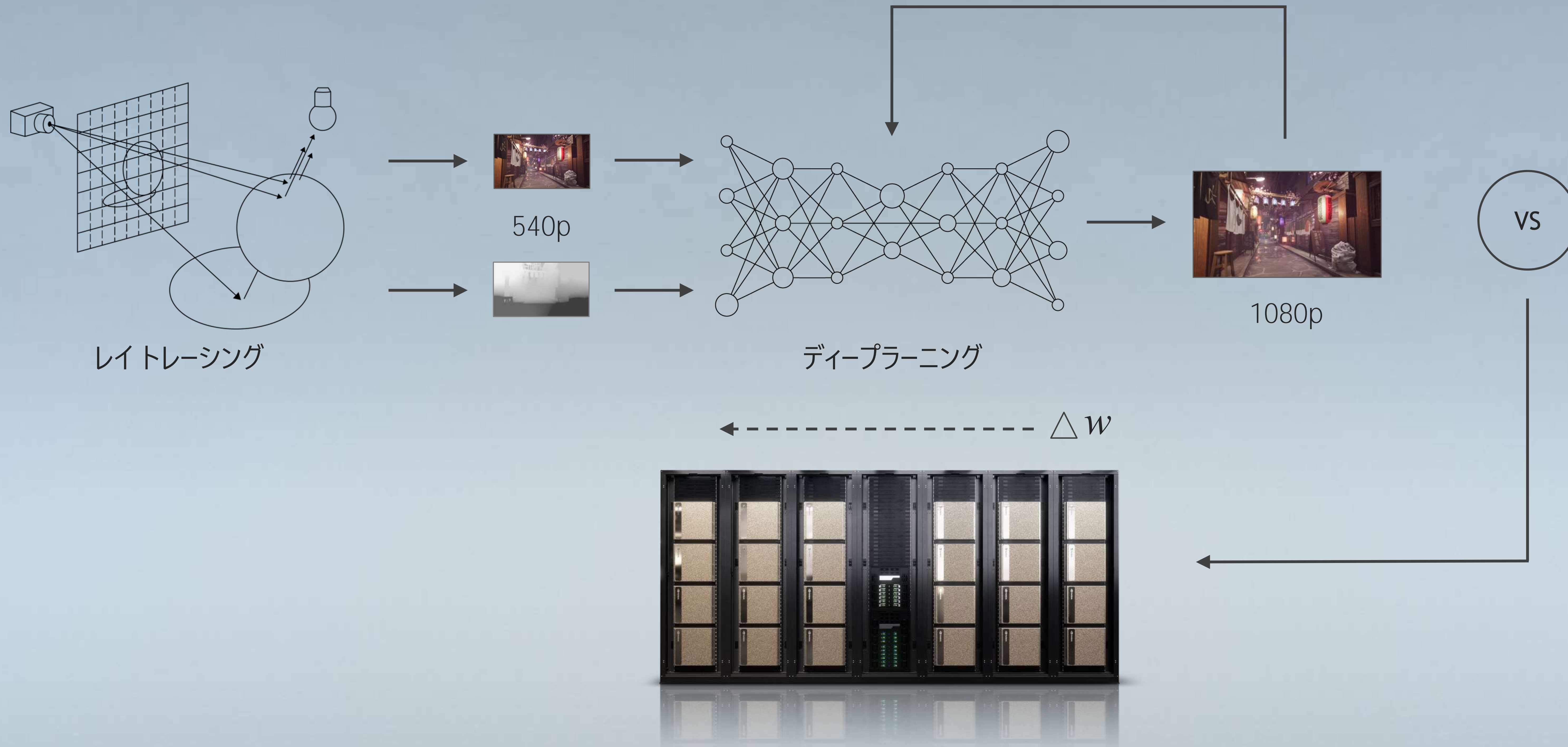
- Jon Peddie, JPR



NVIDIA RTX コンピュータグラフィックスの新時代 — レイトレーシングと AI

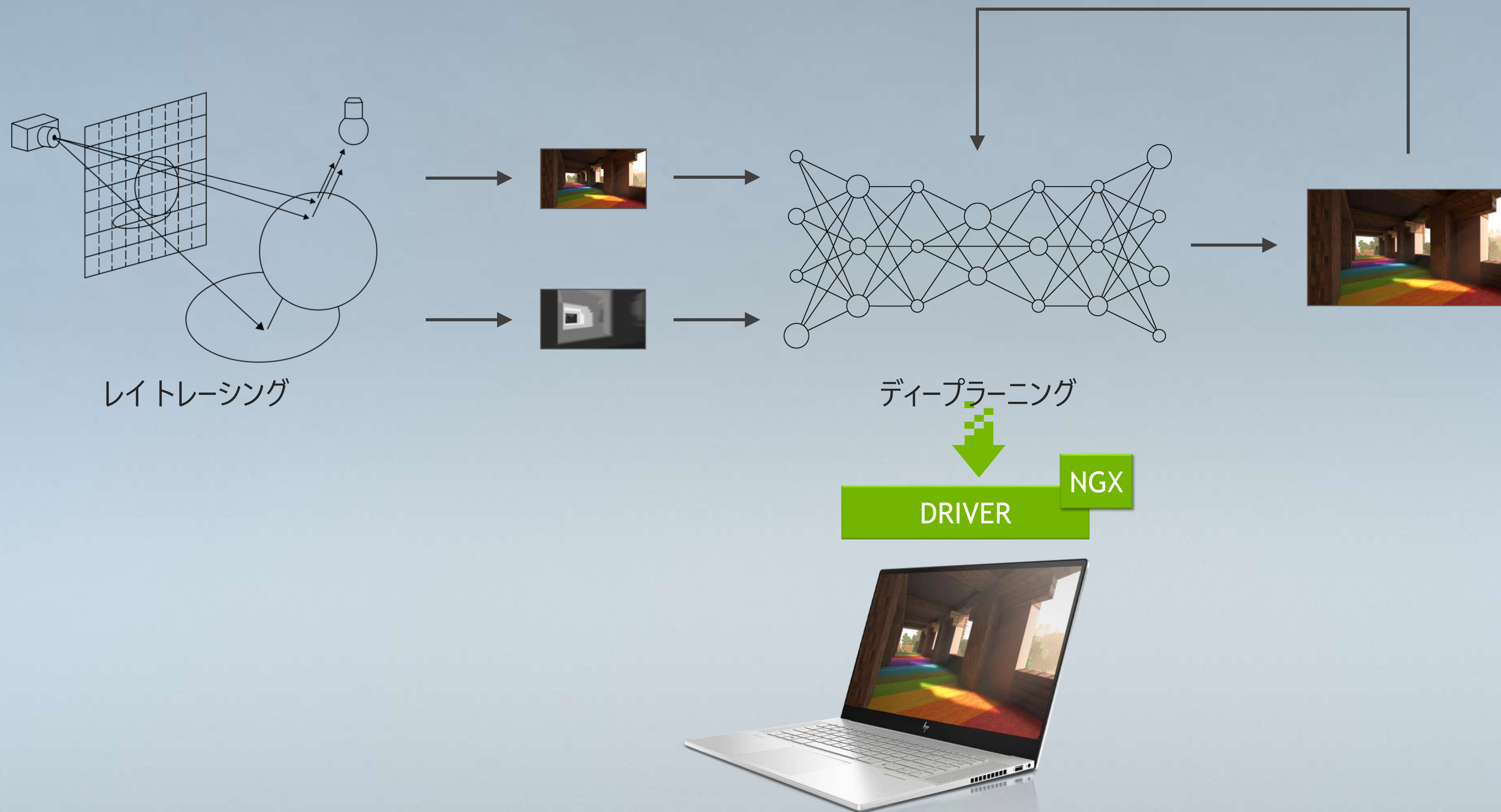


NVIDIA RTX コンピュータグラフィックスの新時代 — レイトレーシングと AI



スーパーコンピュータでレンダーされた 16K Ground Truth

NVIDIA RTX コンピュータグラフィックスの新時代 — レイトレーシングと AI





Ground
Truth 16K



Native
720p



DLSS 1.0
720p > 1080p



DLSS 2.0
720p > 1080p



Native
1080p



Ground
Truth 16K

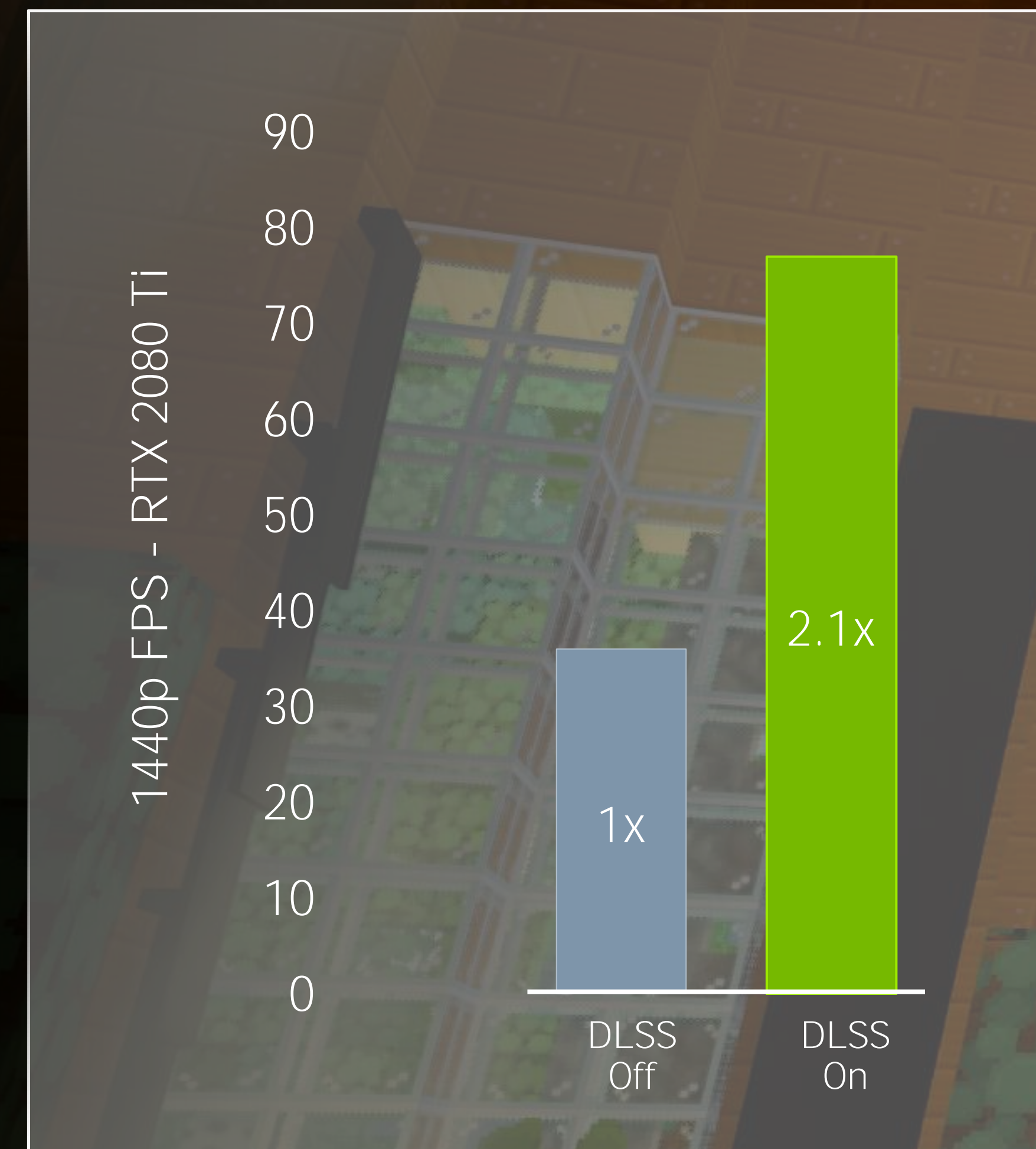


Native
540p



DLSS 2.0
540p > 1080p

MINECRAFT



「ゲームチェンジャー」

- Digital Foundry

「なんてゴージャスなんだ」

- PCWorld

「素晴らしいに他ならない」

- IGN

「見事だ」

- PC Gamer

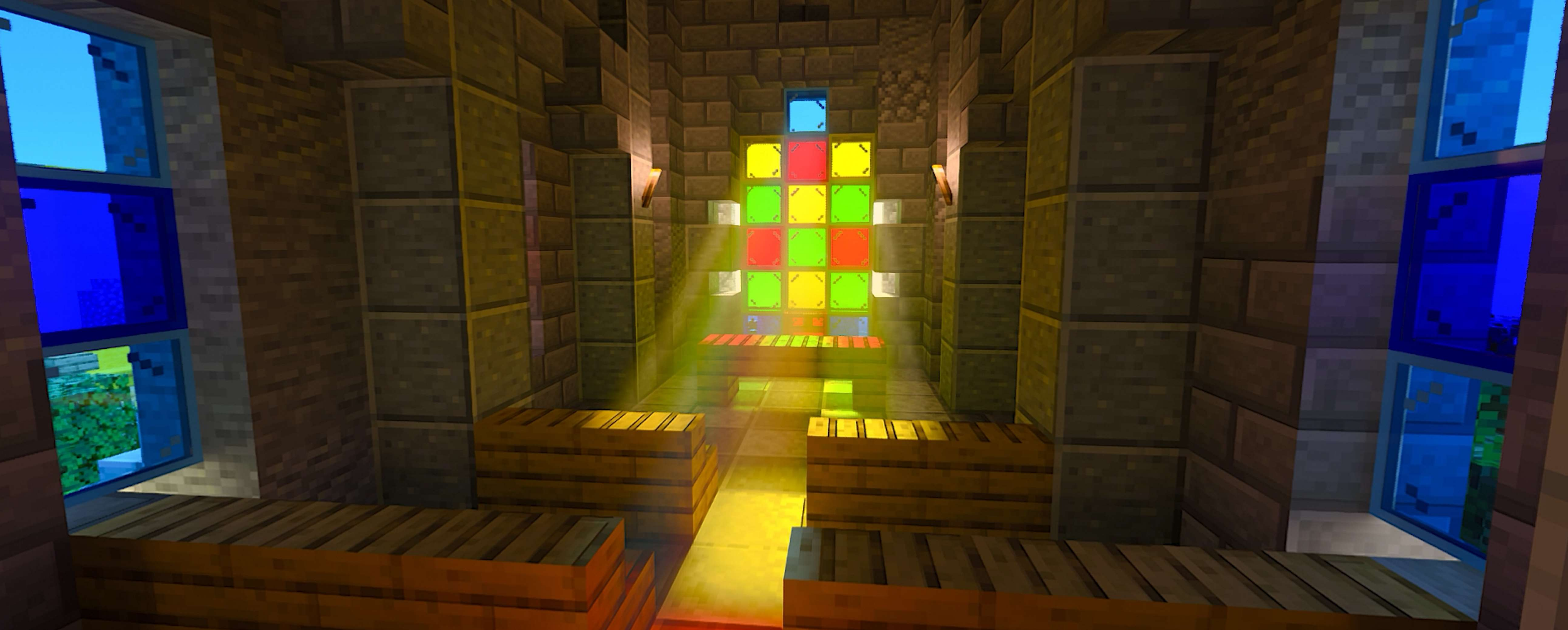
「驚いて開いた口が塞がらない」

- Trusted Reviews



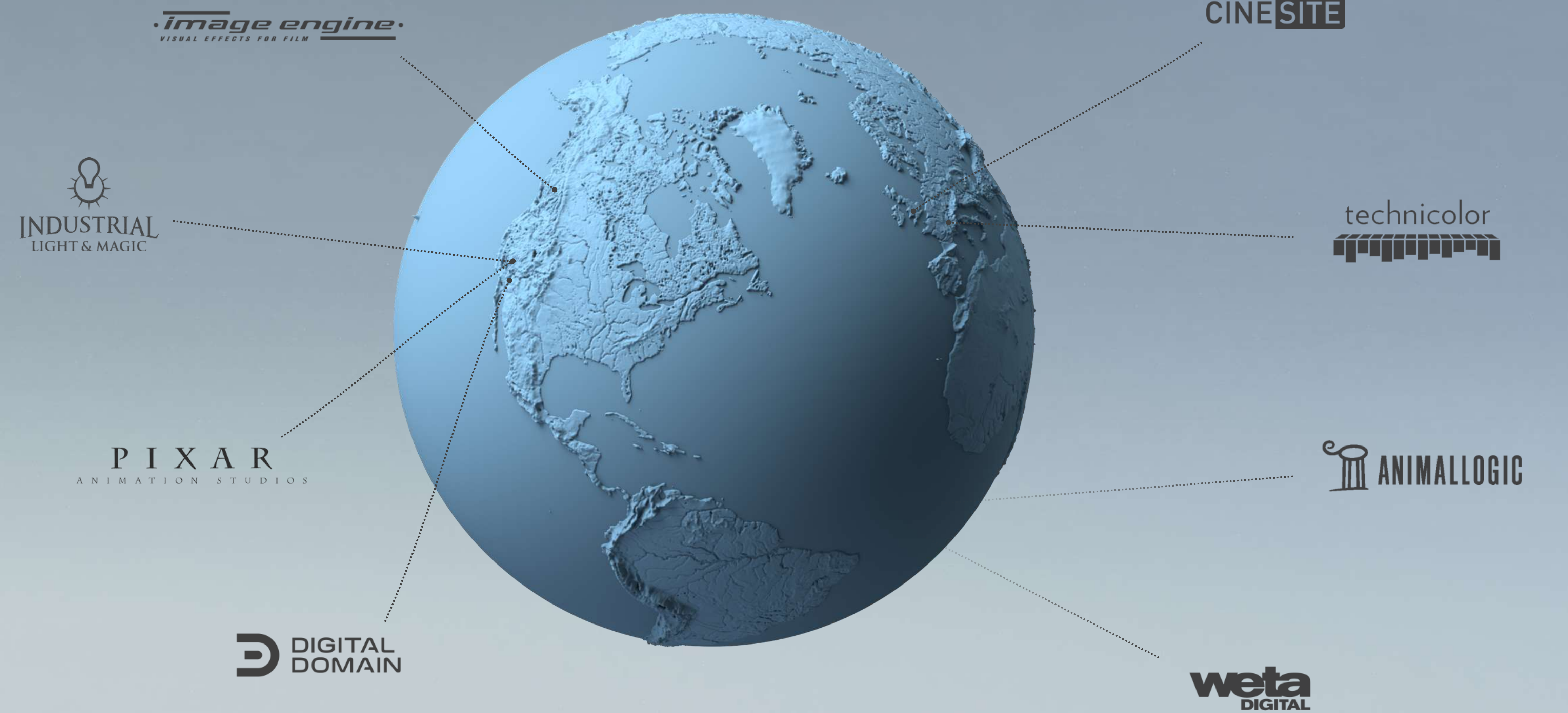
MINECRAFT

RTX
ON



3D は非常に複雑

さまざまなツールと巨大なデータ セット
多様な専門家による大規模なチーム
複数のロケーションとスタジオ
高価



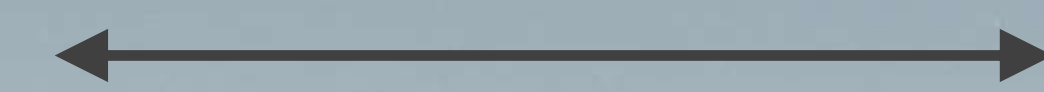
NVIDIA OMNIVERSE

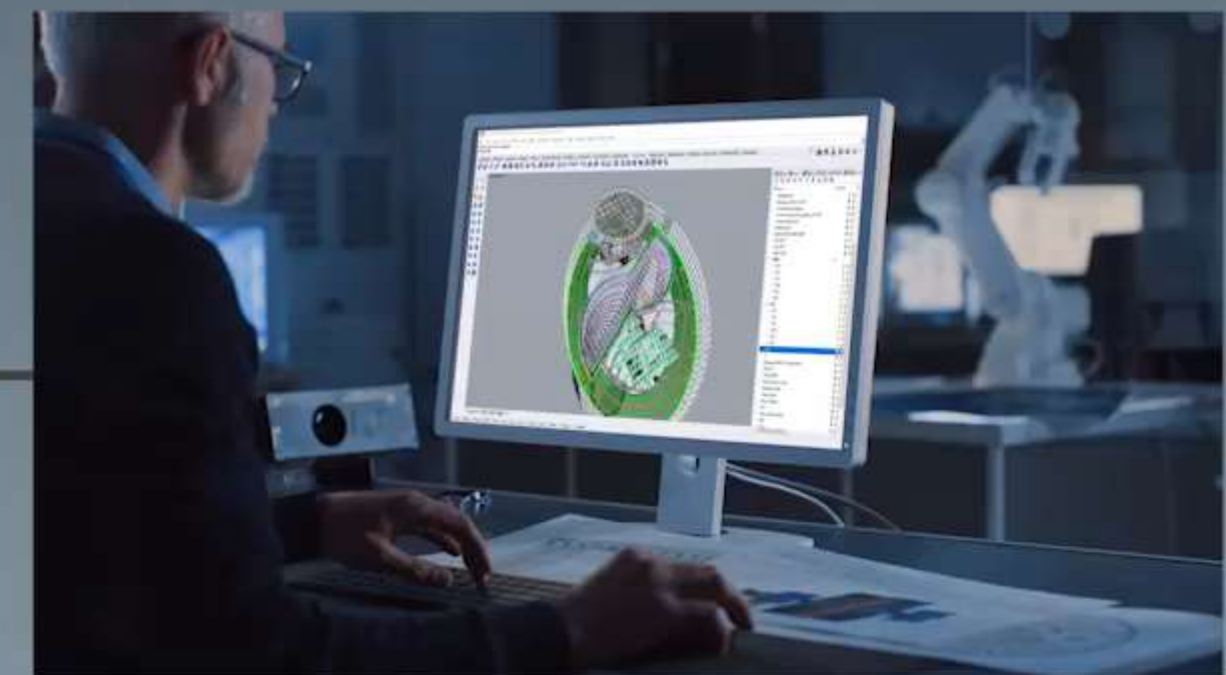
デザインワークフローのコラボレーションプラットフォーム

USD (Universal Scene Description) ベースに構築

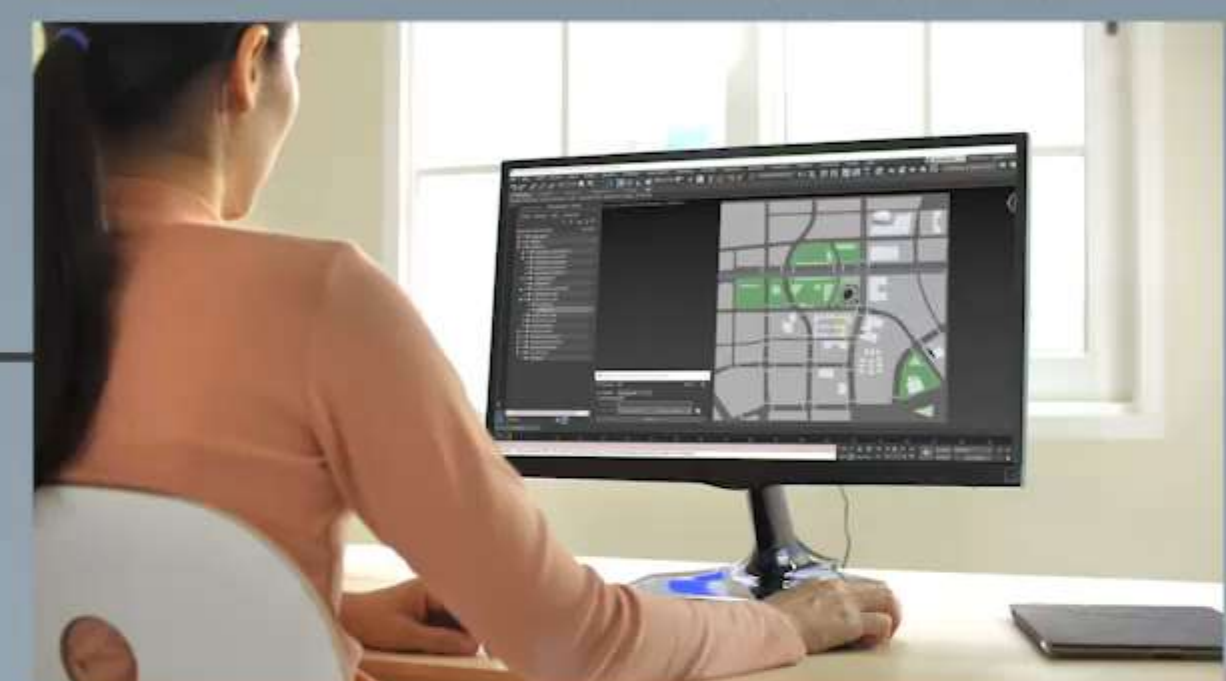
マテリアルと物理をインタラクティブレンダラーに搭載

PC と Linux に対応、ストリーミング端末には Mac、Android が対応

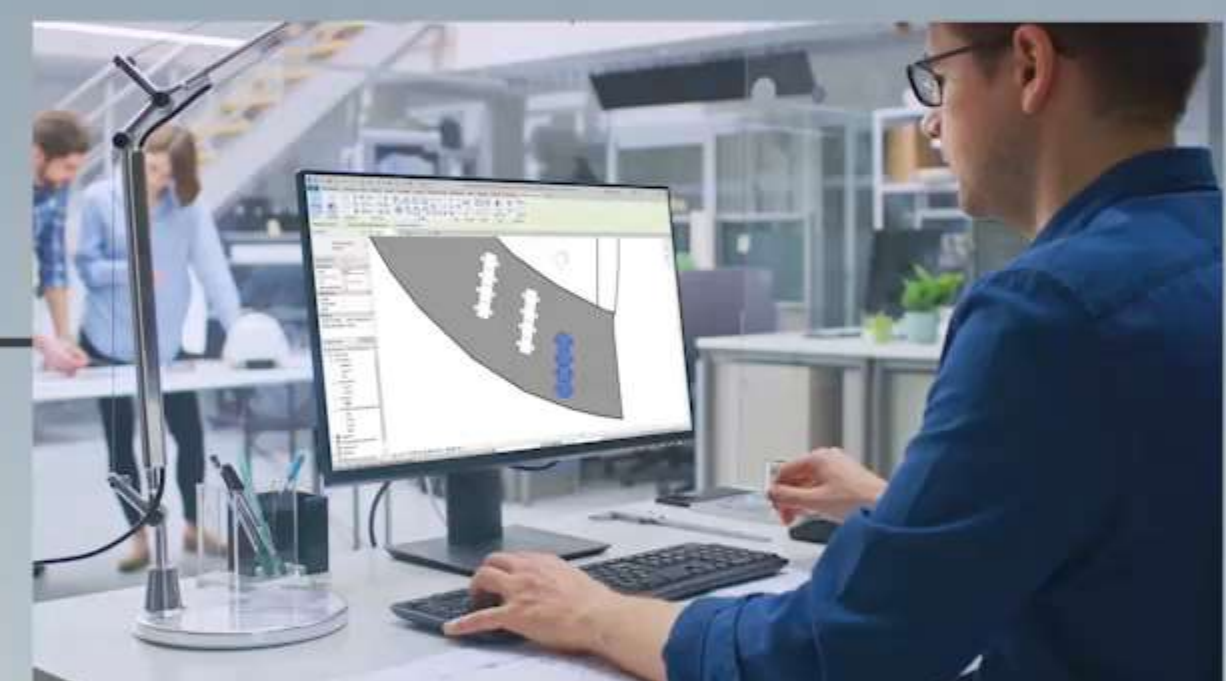




RHINO



MAX



REVIT



AR





HOOPKA
OIL COLOUR
L'ARTISTE
POUR LE
L'ARTISTE
27 ml (1.25 US Fl. Oz.)

発表： リモート コラボレーションに最適化された NVIDIA RTX サーバー

Omniverse によるデザイン ワークフロー コラボレーション
完全なレイトレーシングによるグローバル イルミネーションに対応したインタラクティブ プロダクション レンダリング
デザインとシミュレーションで検証された Quadro 仮想ワークステーション

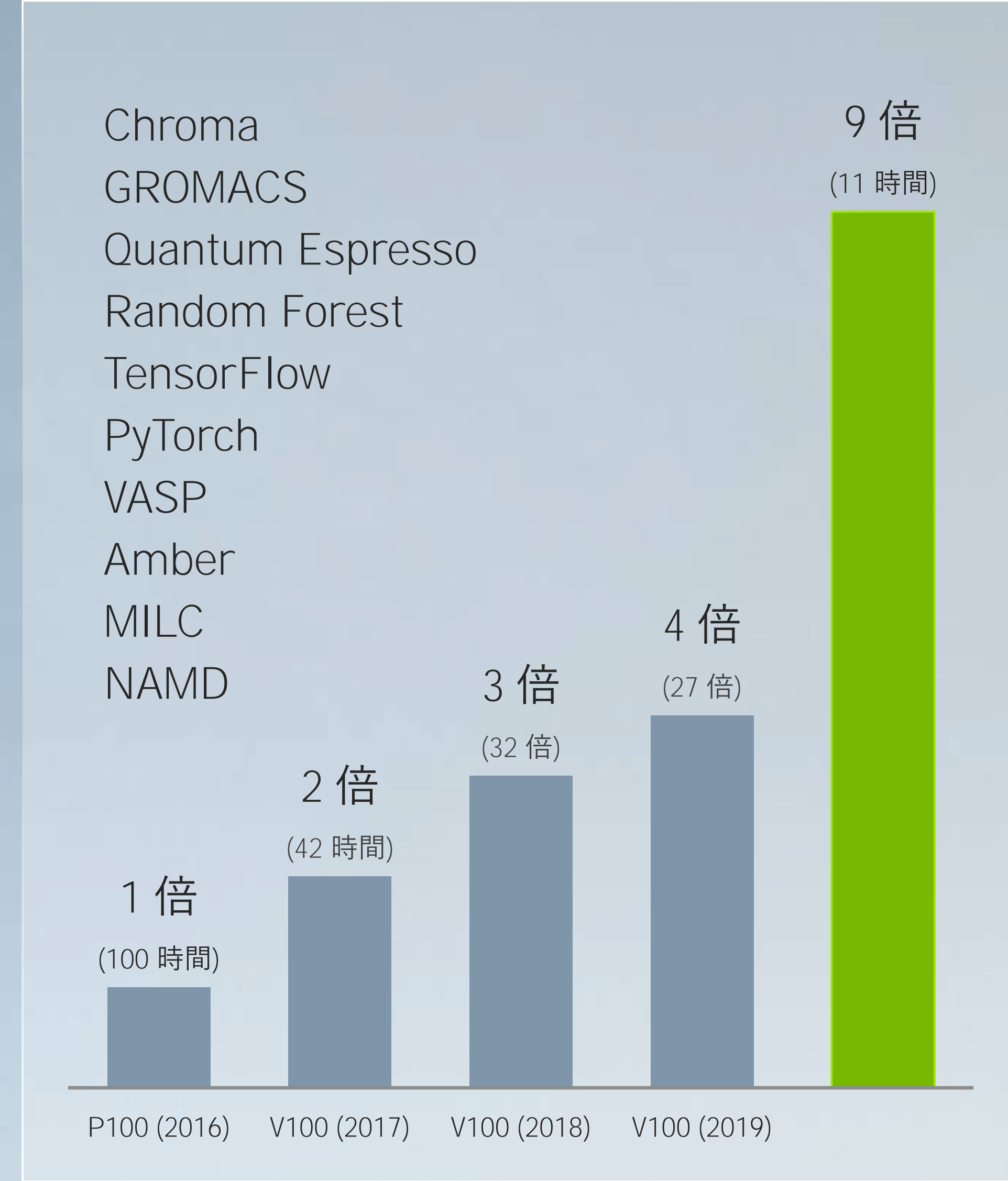
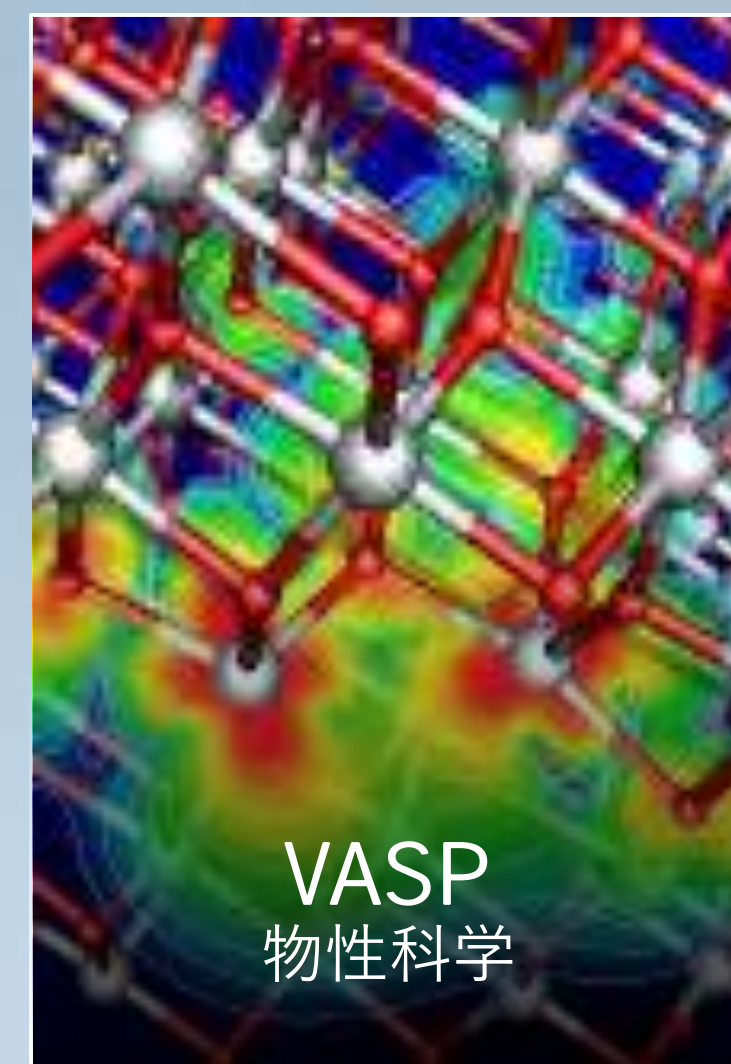
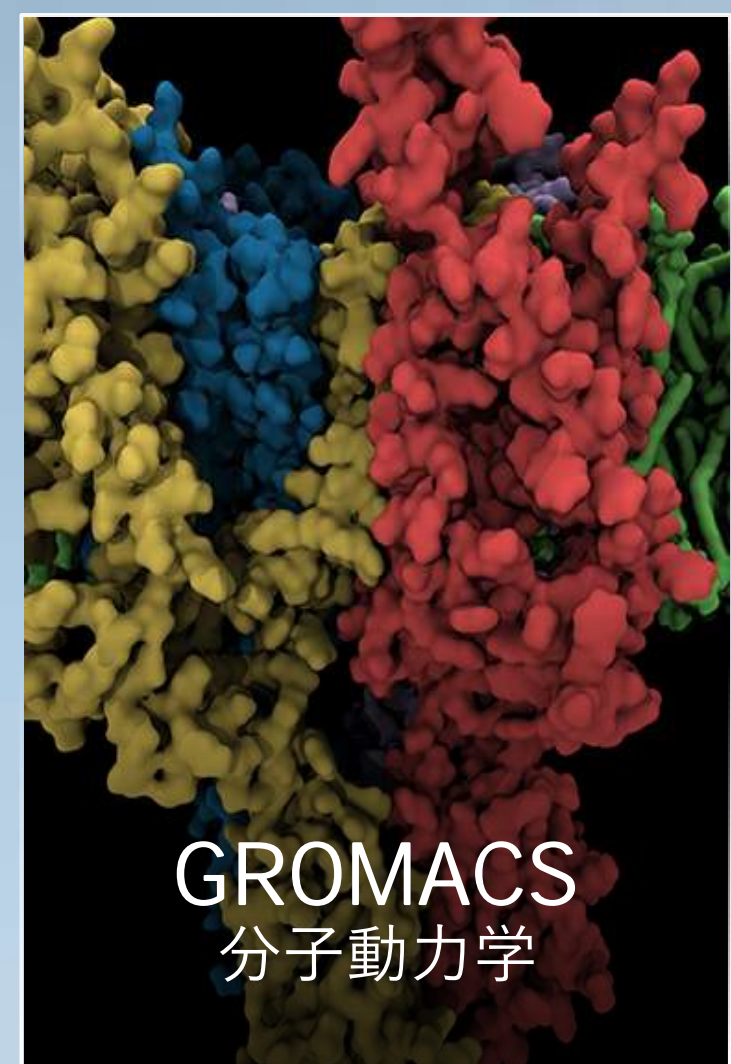
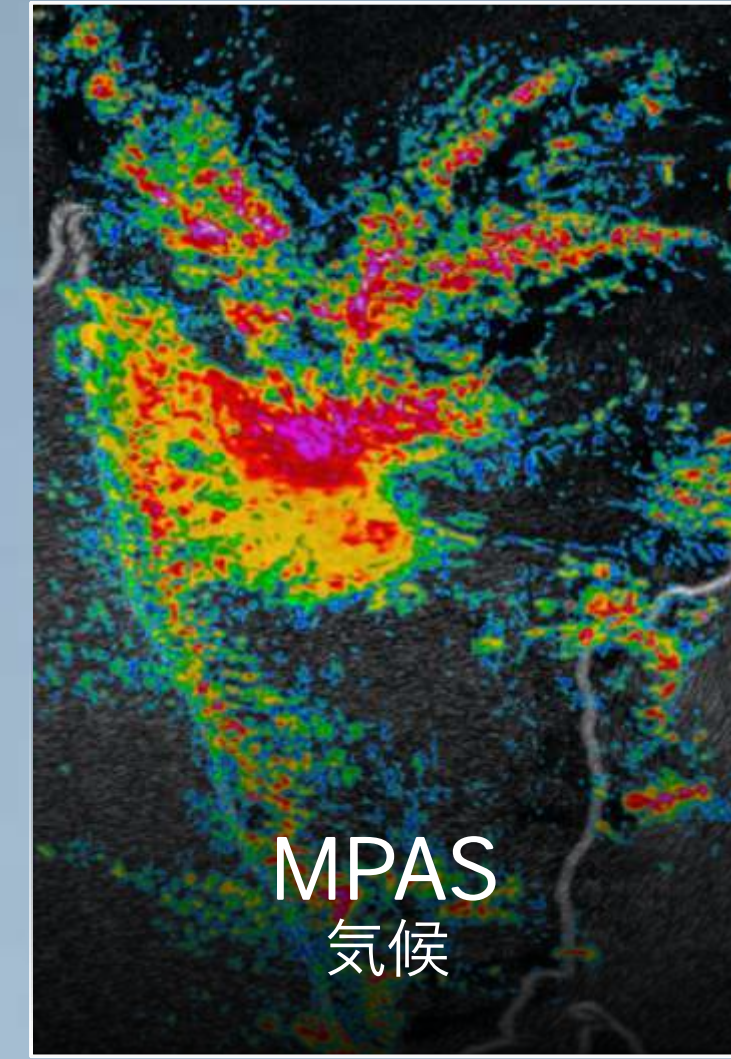
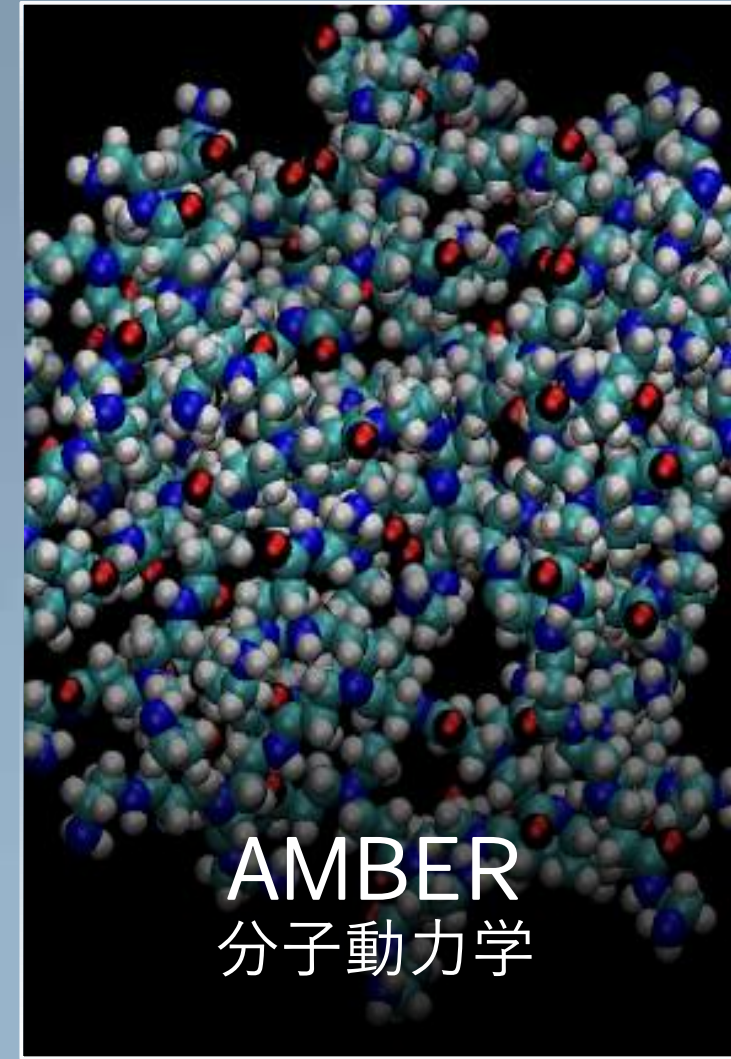
NVIDIA RTX サーバーは出荷を開始

BOX **DELL** Technologies

Hewlett Packard
Enterprise **SUPERMICR**



NVIDIA HPC

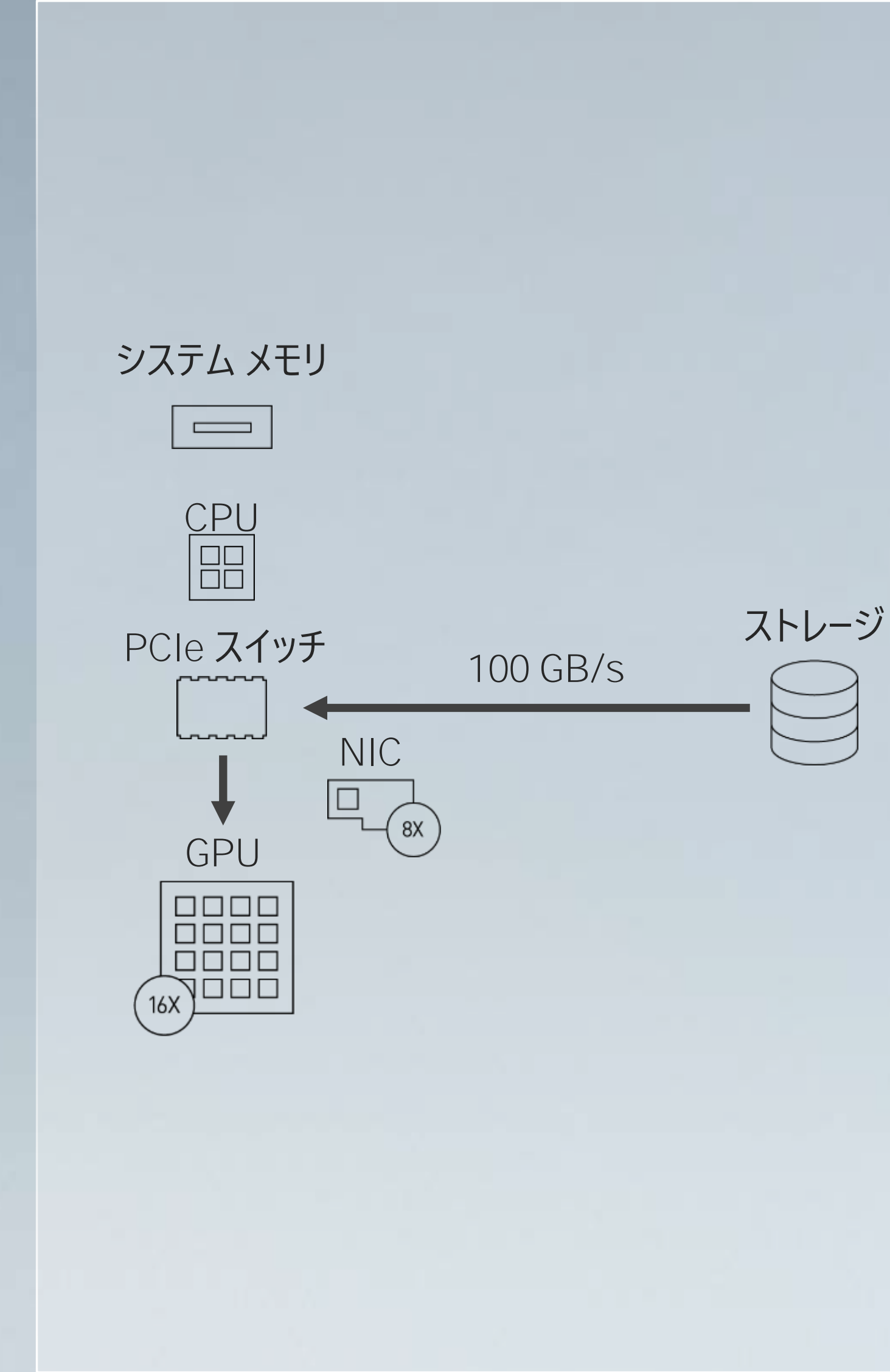


700 以上のアプリが
CUDA で高速化

9 倍のパフォーマンス
過去 4 年で



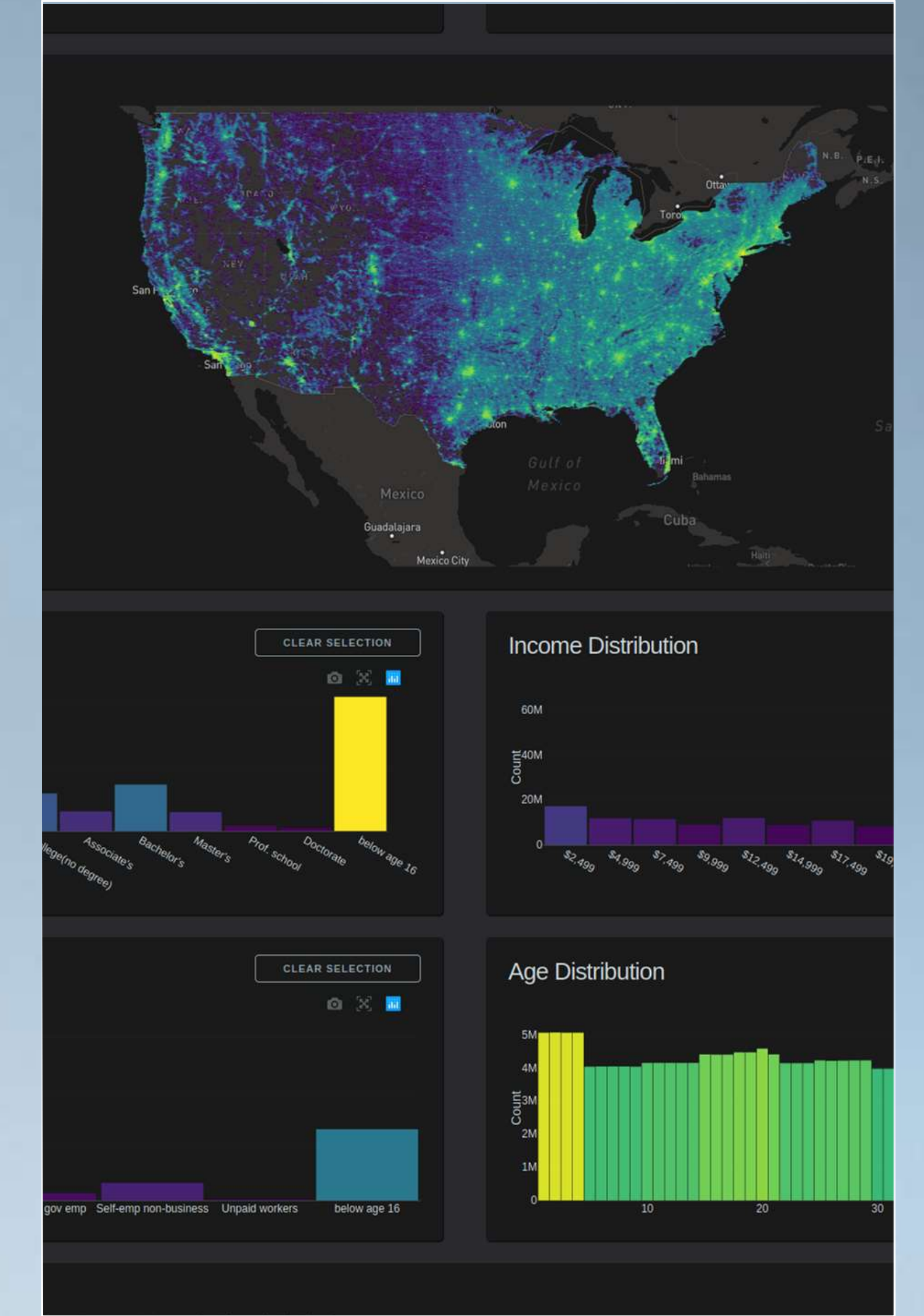
CUDA ON ARM



NVIDIA MAGNUM-IO
I/O アクセラレーション

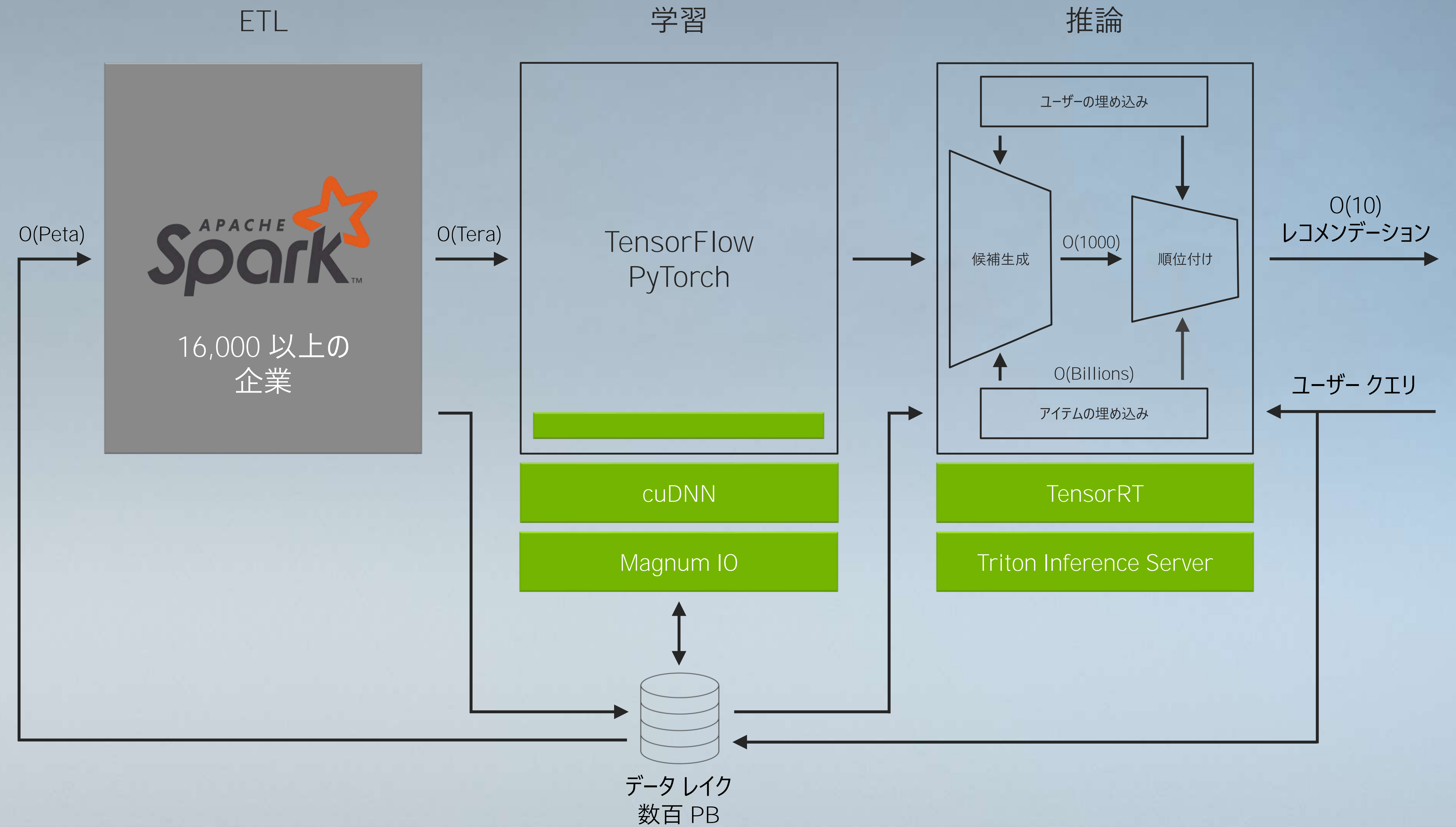


NVIDIA PARABRICKS
ゲノミクス

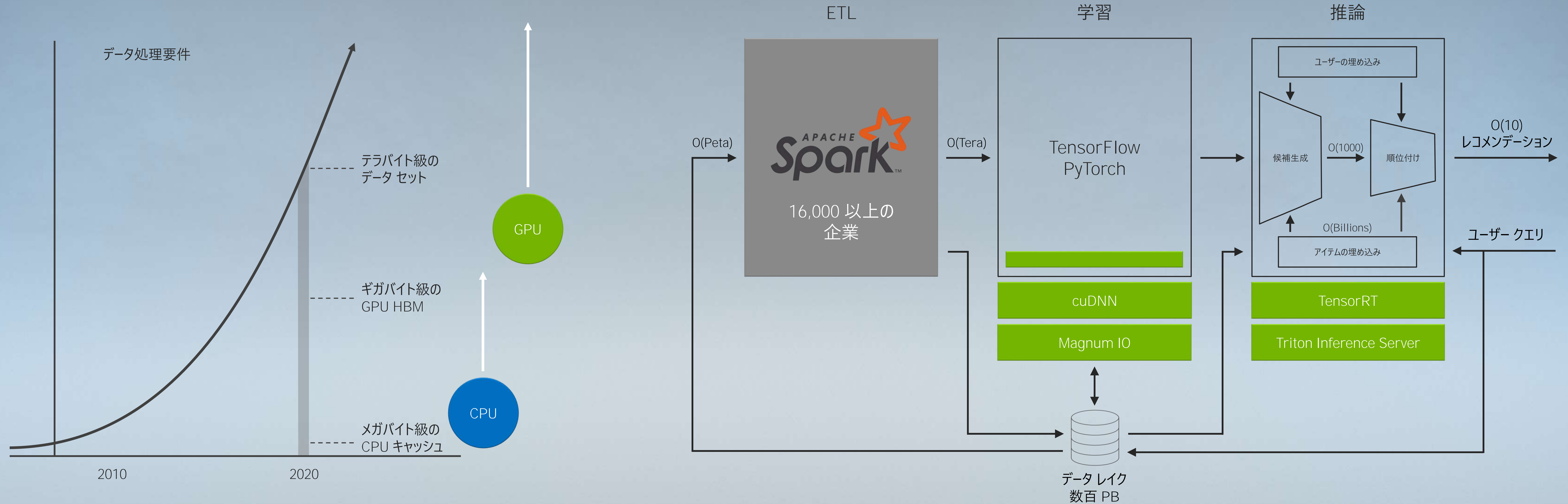


NVIDIA RAPIDS
データ分析

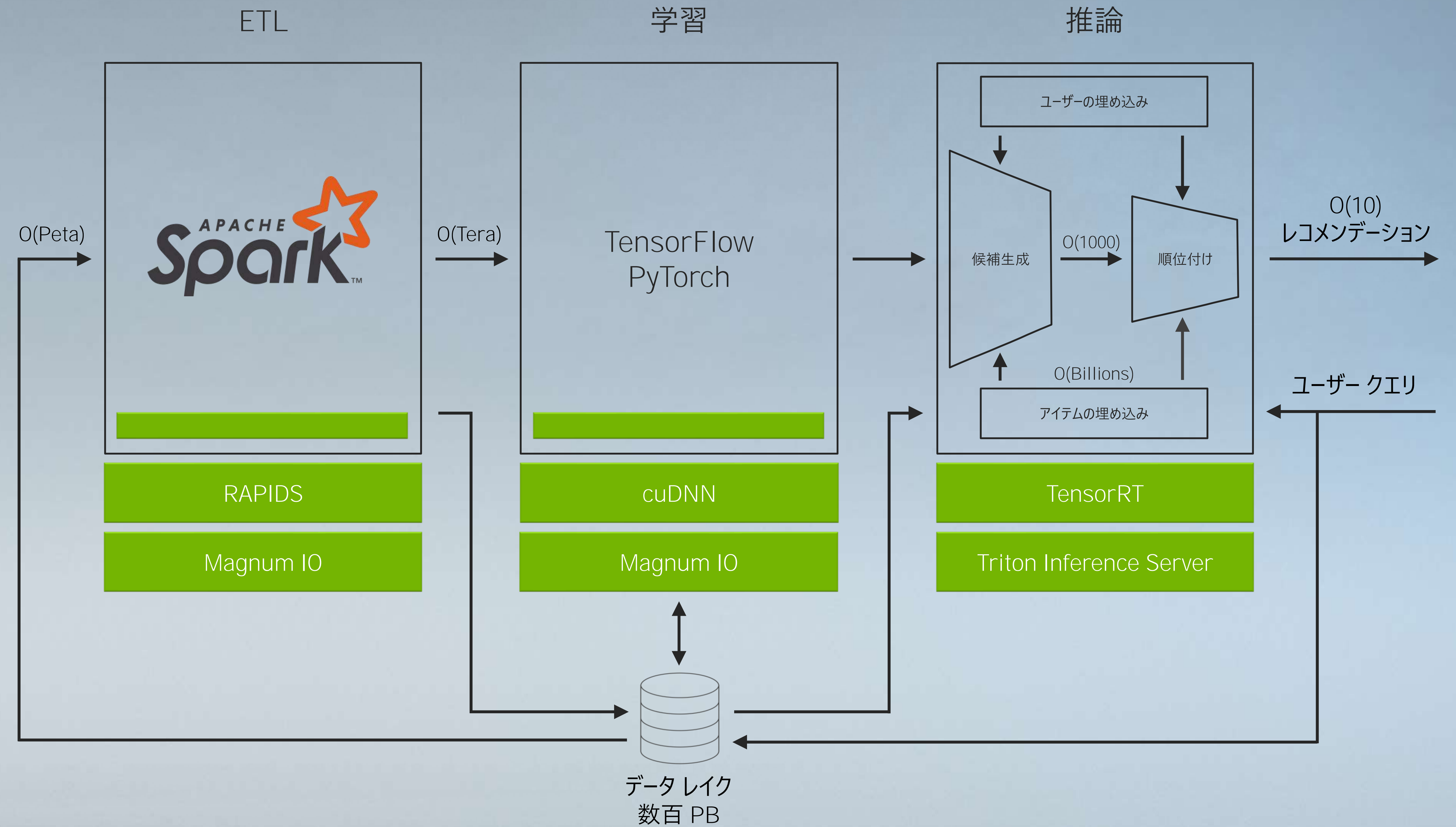
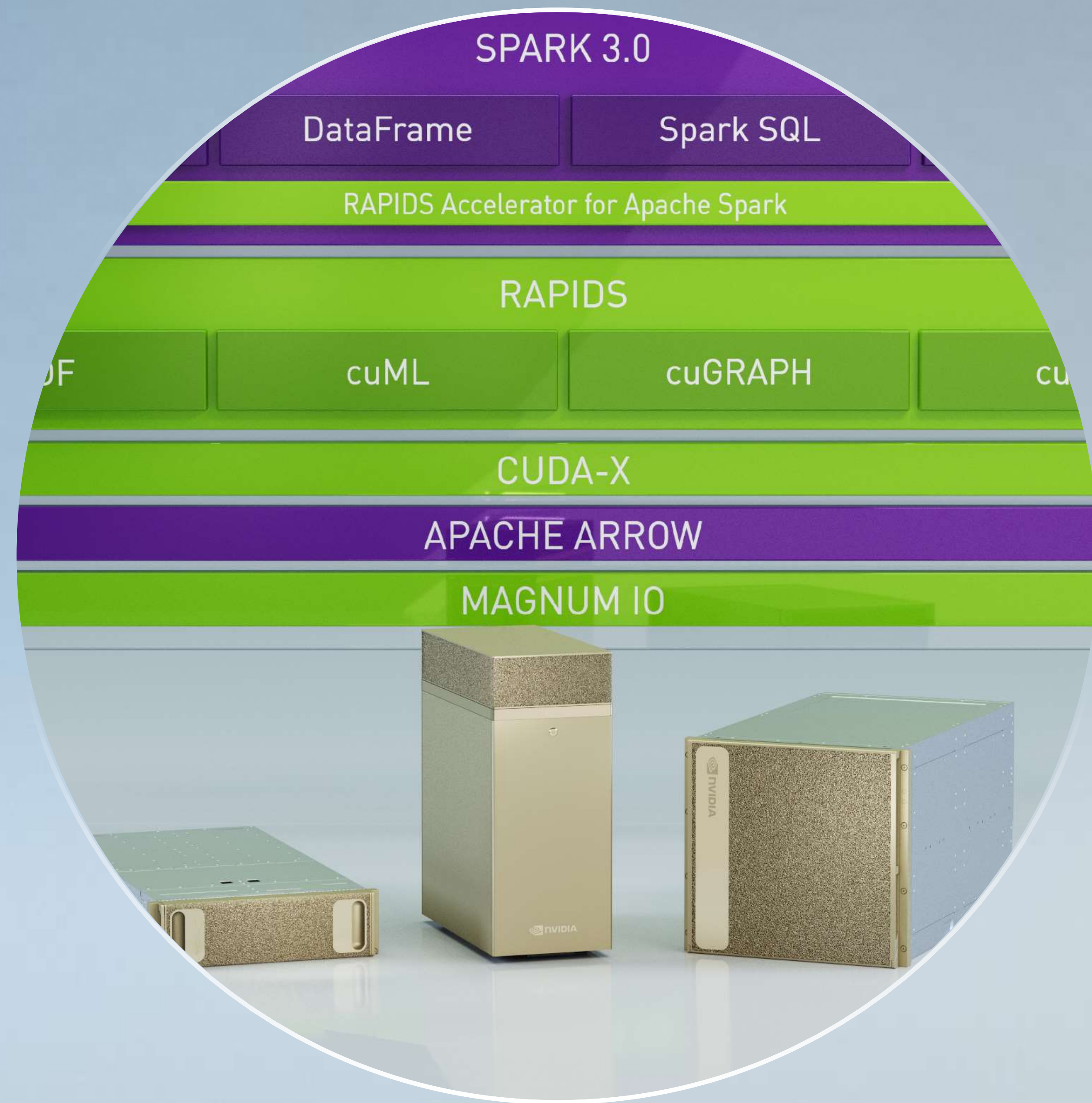
機械学習パイプラインは HPC の課題



機械学習のデータが飛躍的に増大



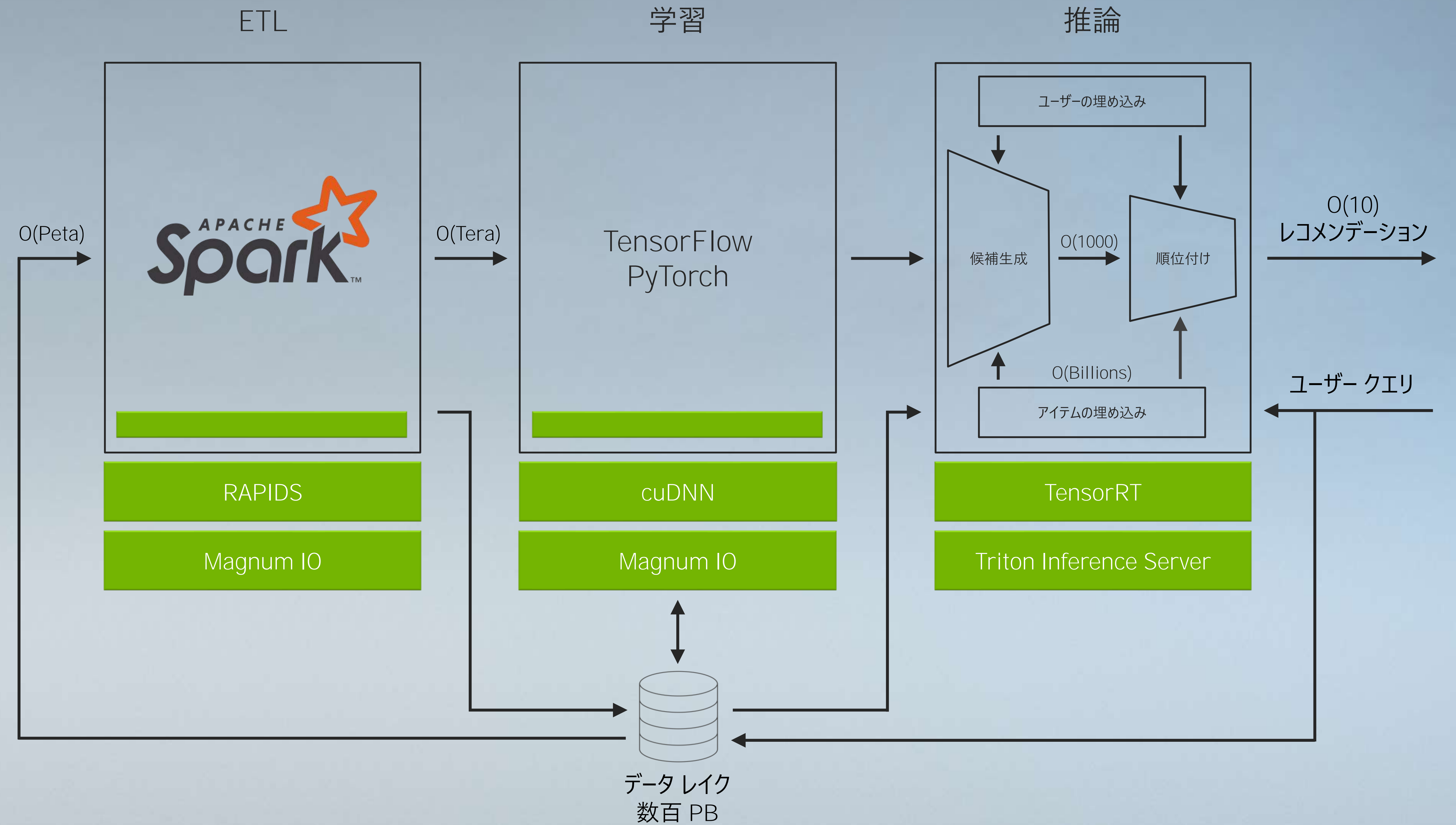
発表: NVIDIA が SPARK 3.0 を高速化



SPARK 3.0 は最先端技術の基に構築 RAPIDS が ETL ベンチマークを打ち破る



TPCx-BB @ SF 10K で 17 GB/秒のスループット
\$1M | 18.2U CPU システム | 2 ラック | 16 kW



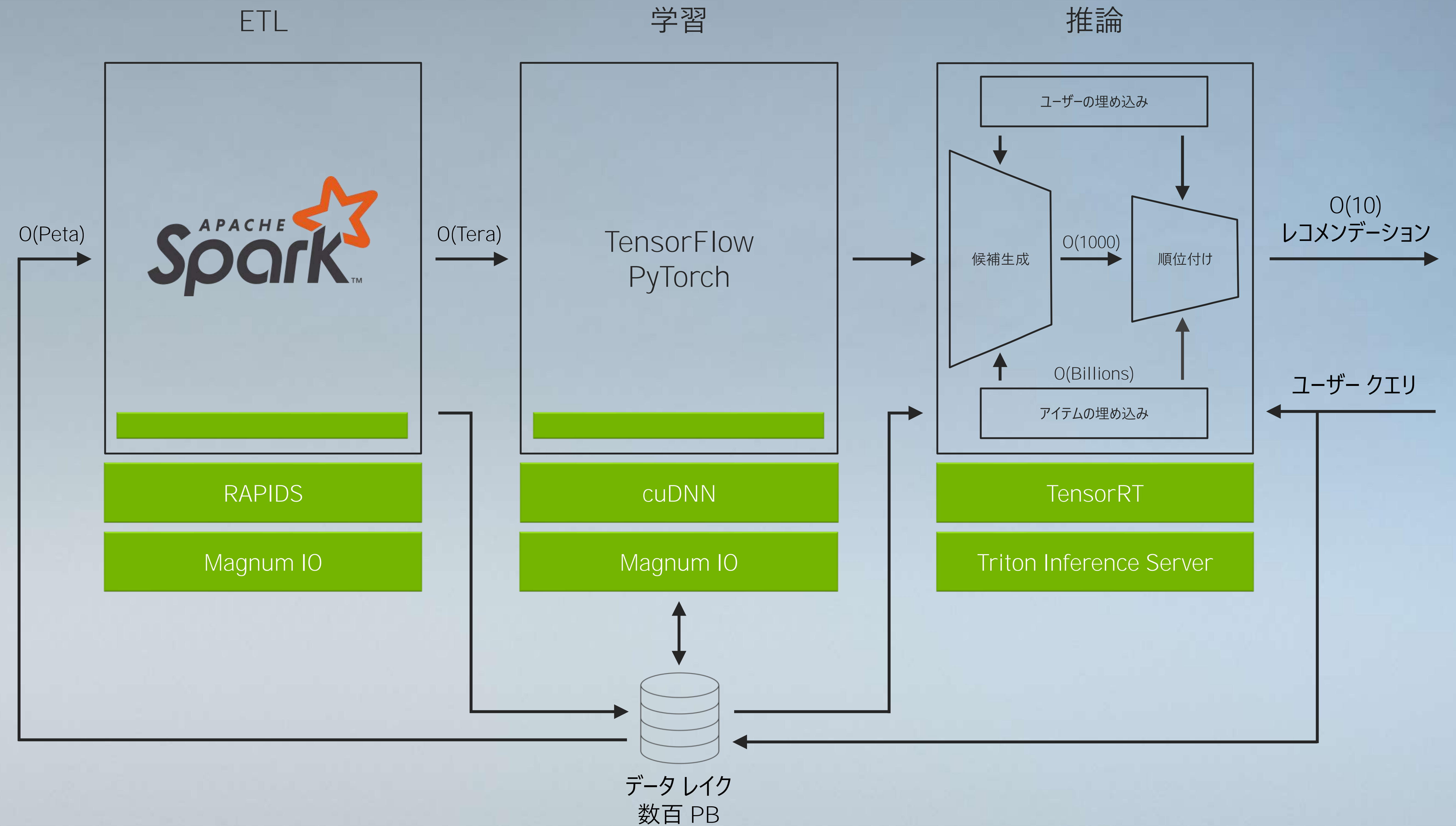
SPARK 3.0 は最先端技術の基に構築 RAPIDS が ETL ベンチマークを打ち破る



\$2M 163 GB/秒

RAPIDS を実装した TPCx-BB @ SF 10K で 163 GB/秒のスループット

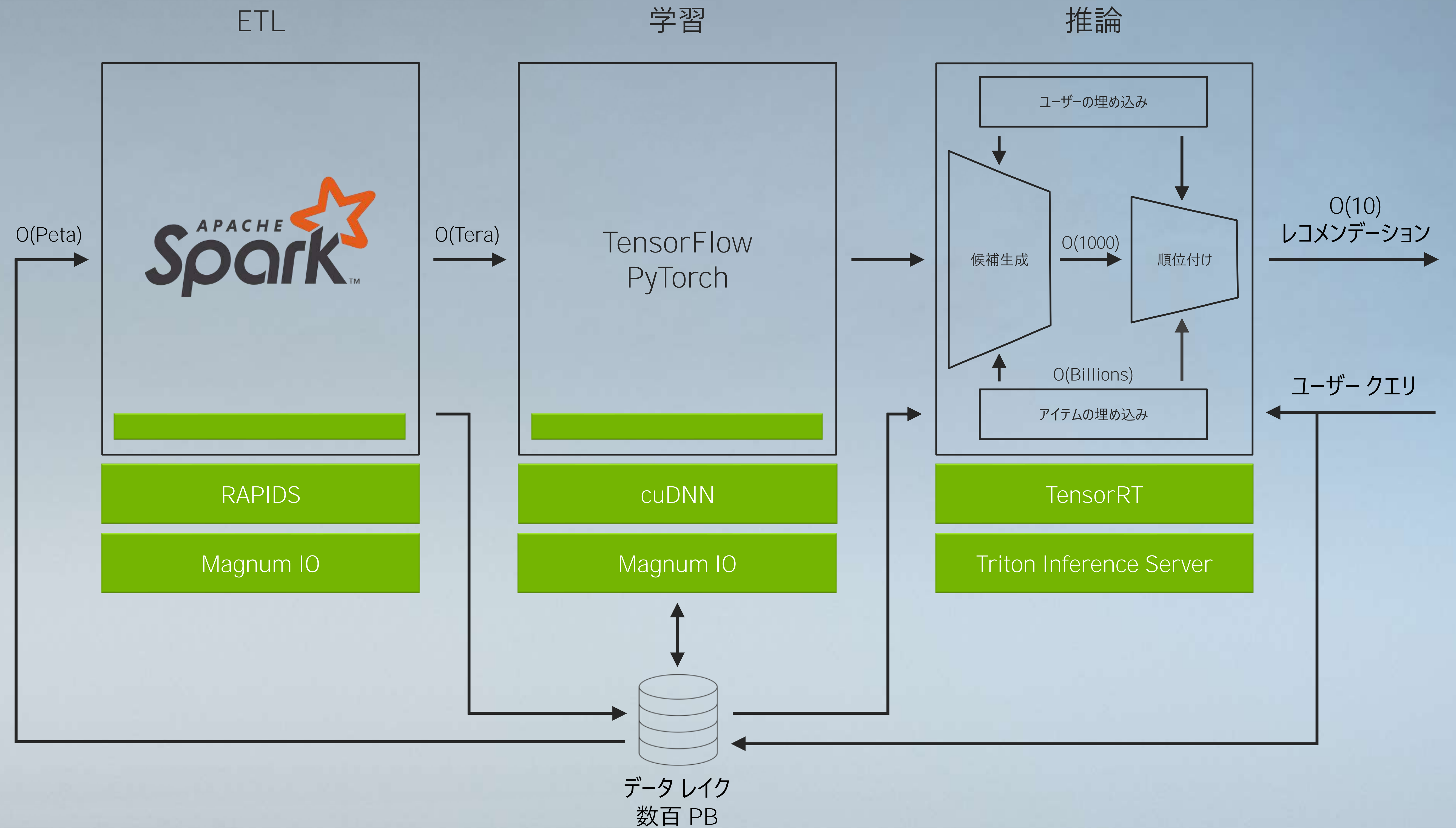
\$2M | 16 台の DGX-1 | 2 ラック | 56 kW



SPARK 3.0 は最先端技術の基に構築 RAPIDS が ETL ベンチマークを打ち破る



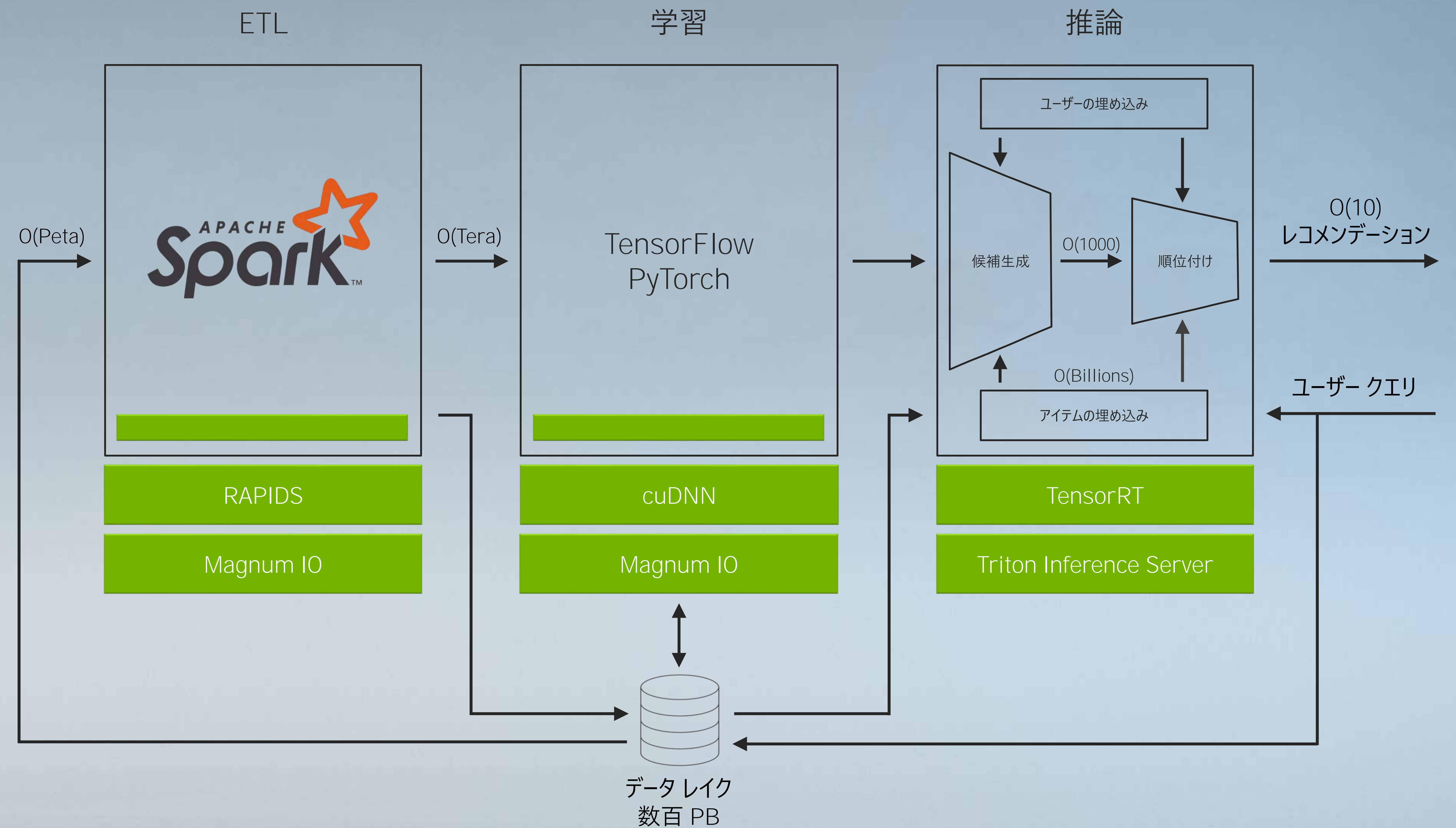
RAPIDS を実装した TPCx-BB @ SF 10K で 163 GB/秒のスループット
\$2M | 16 台の DGX-1 | 2 ラック | 56 kW



SPARK 3.0 は最先端技術の基に構築 RAPIDS が ETL ベンチマークを打ち破る



TPCx-BB @ SF 10K で 163 GB/秒相当のスループット
\$10M | 167 台の 2U CPU サーバー | 11 ラック | 140 kW

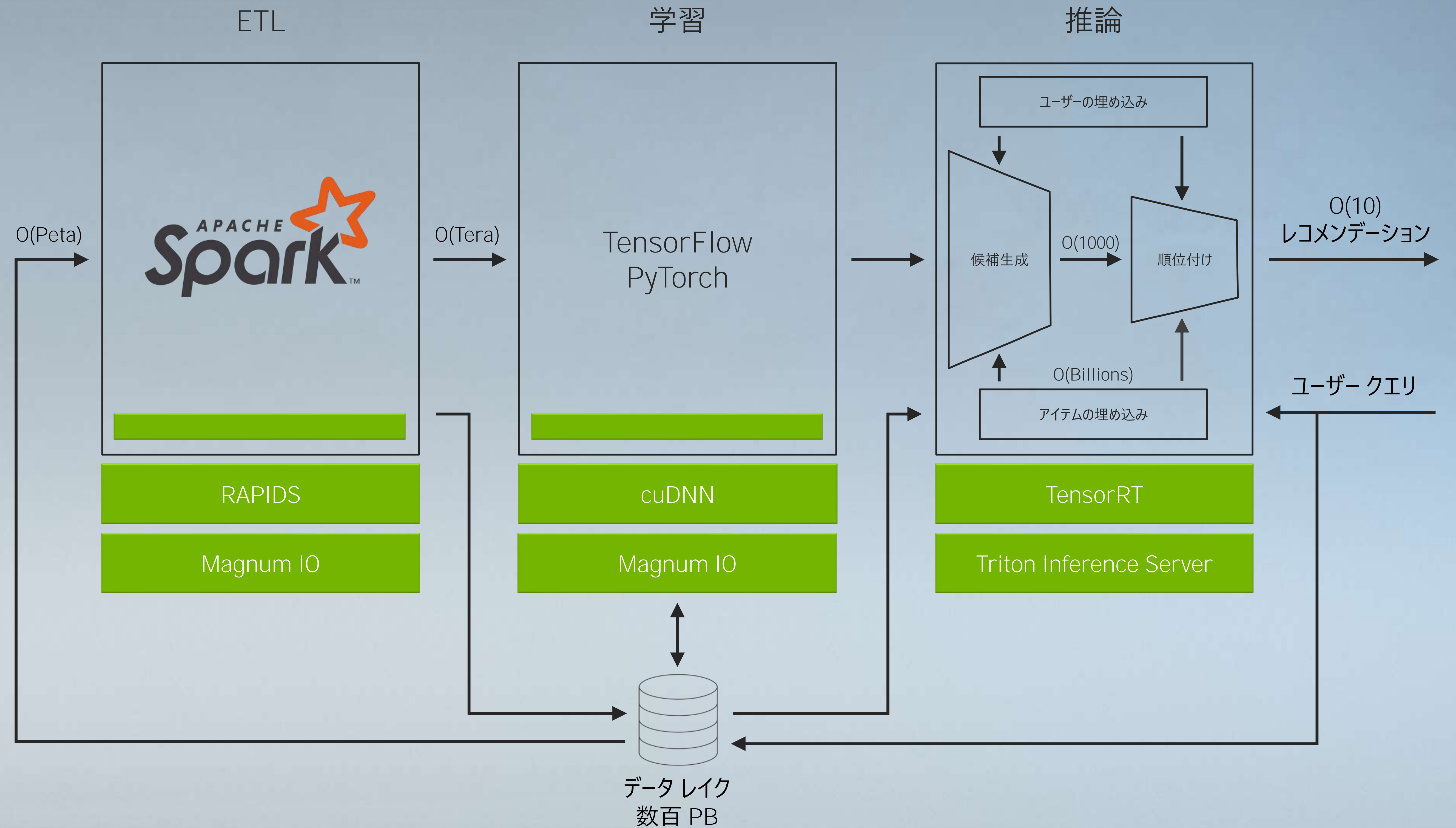


SPARK 3.0 は最先端技術の基に構築 RAPIDS が ETL ベンチマークを打ち破る



RAPIDS を実装した TPCx-BB @ SF 10K で 163 GB/秒のスループット

\$2M | 16 台の DGX-1 | 2 ラック | 56 kW



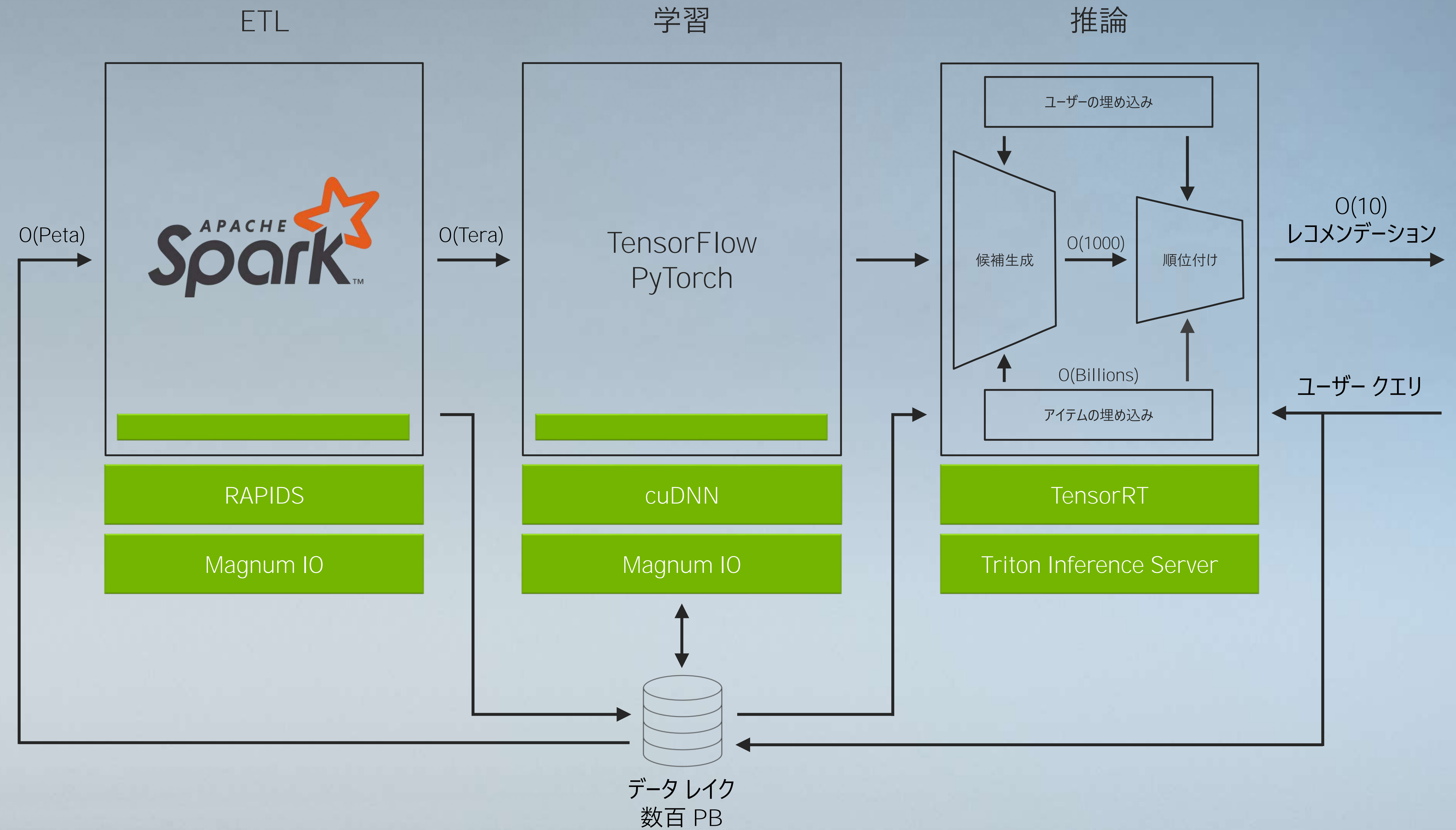
発表: DATABRICKS が NVIDIA で高速化



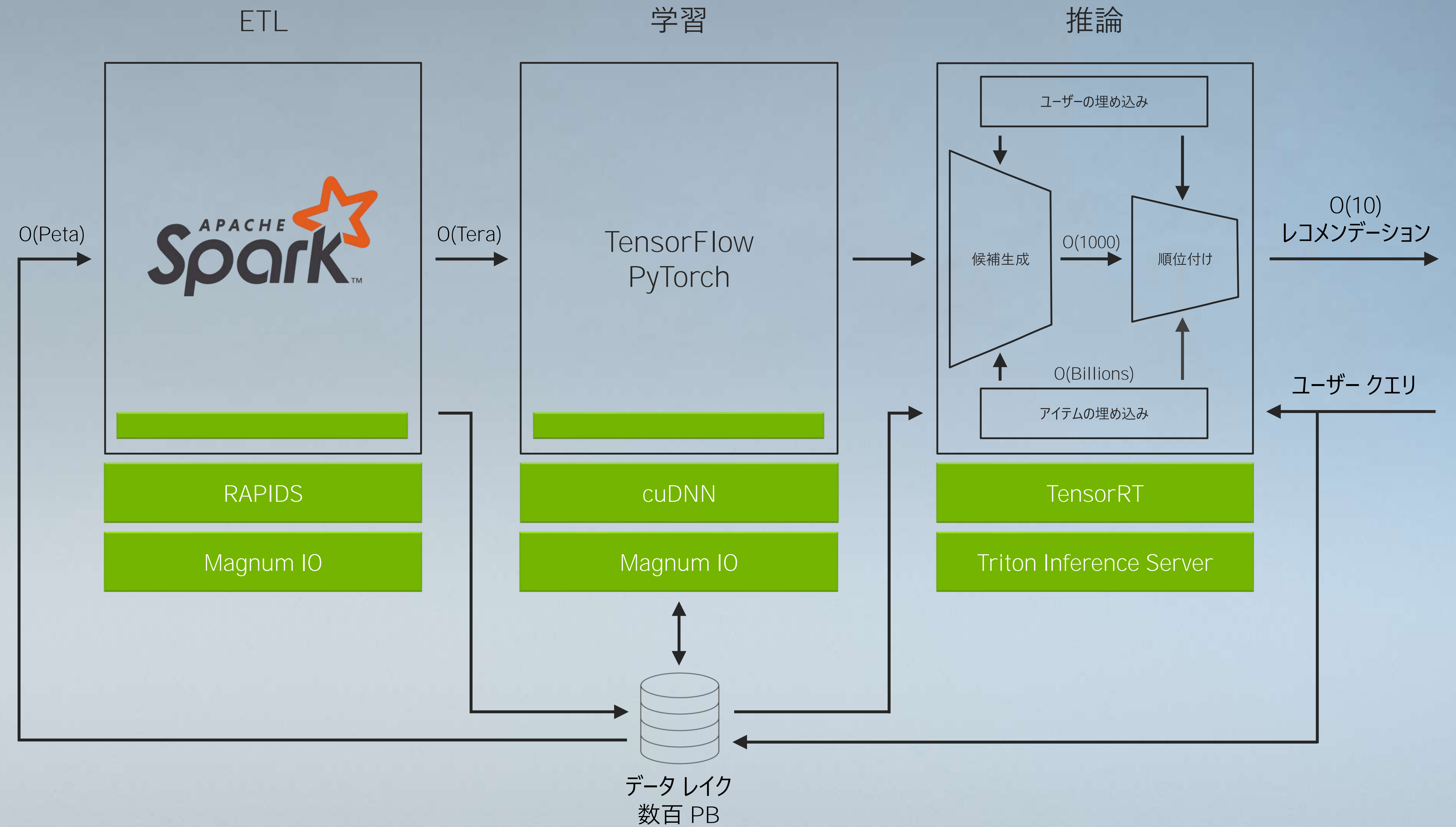
databricks

「これらの貢献は、Apache Spark 3.0 と Databricks を使用でデータパイプライン、モデルの学習とスコアリングを高速化し、より多くのブレイクスルーと洞察の獲得につながります。」

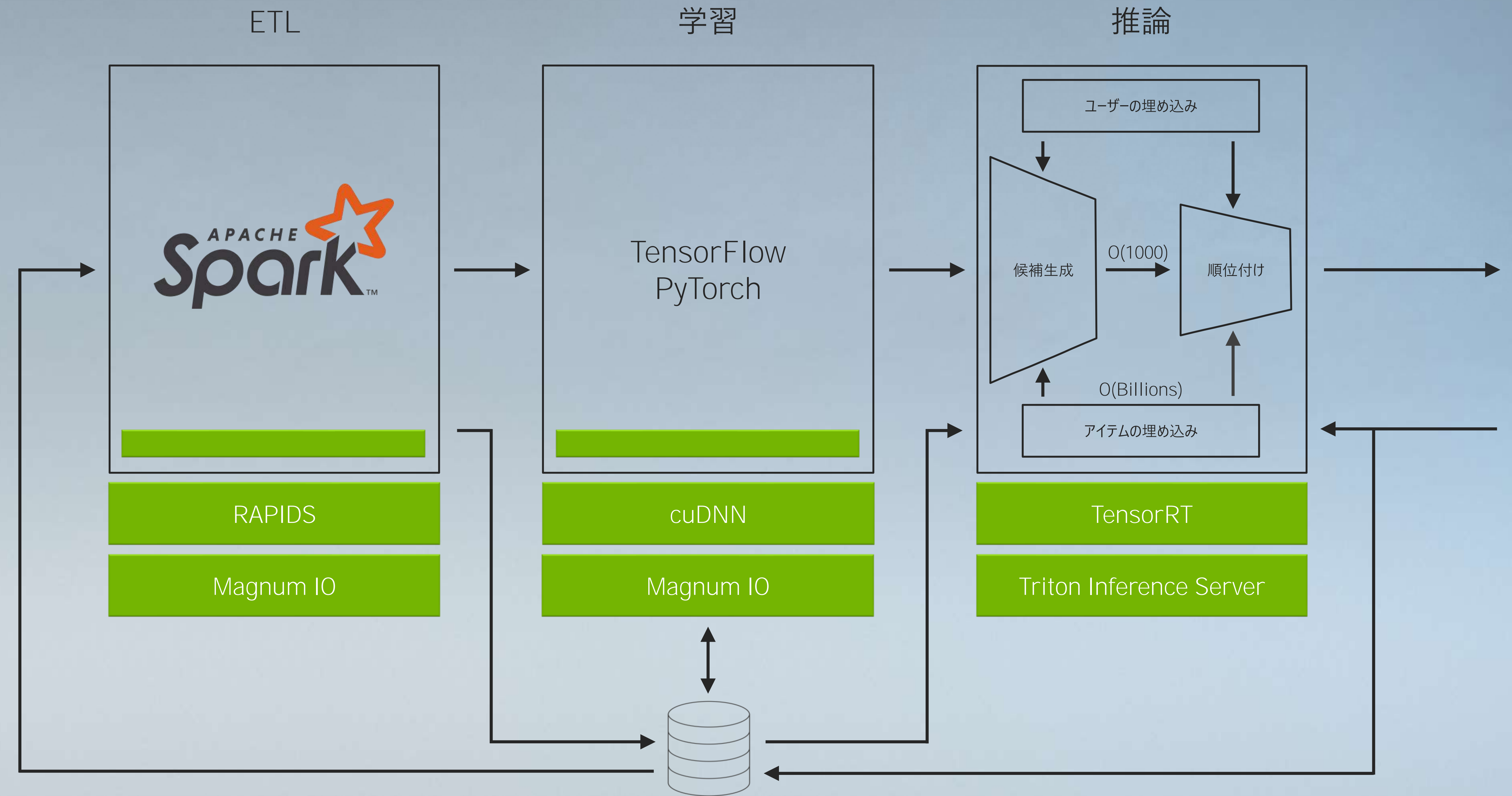
— Matei Zaharia, original creator of Apache Spark and chief technologist at Databricks



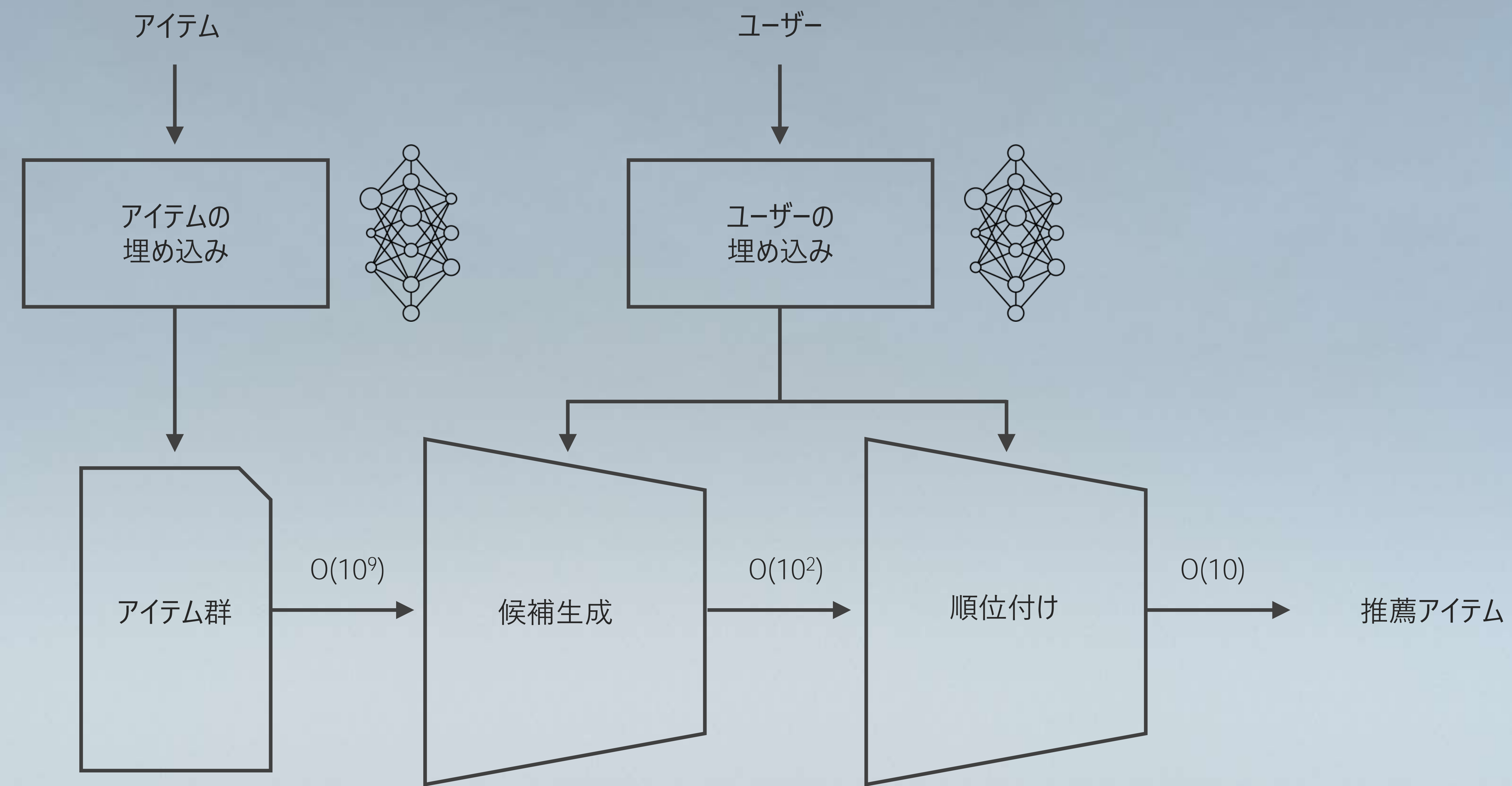
発表: クラウド分析プラットフォームが NVIDIA で高速化



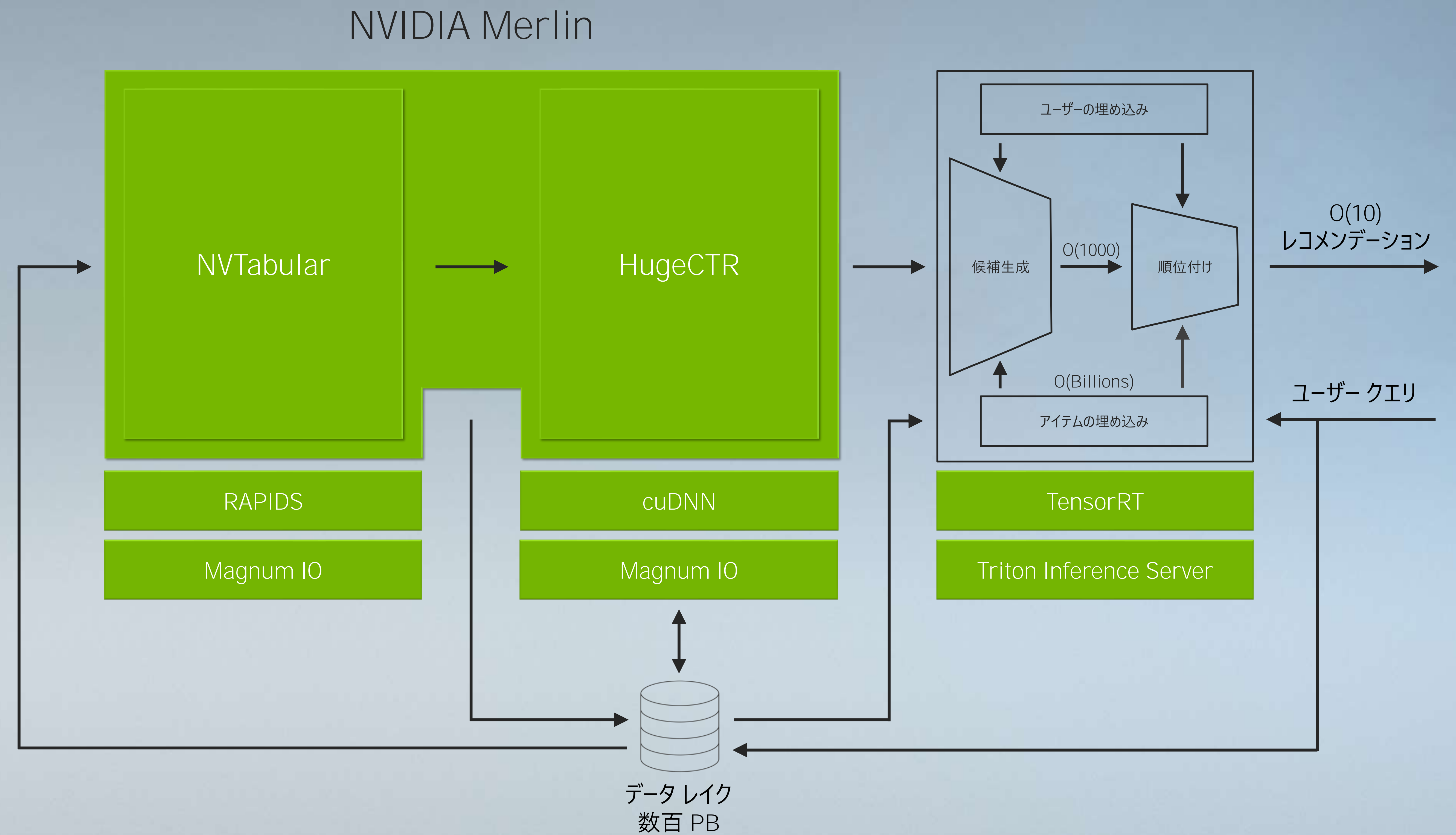
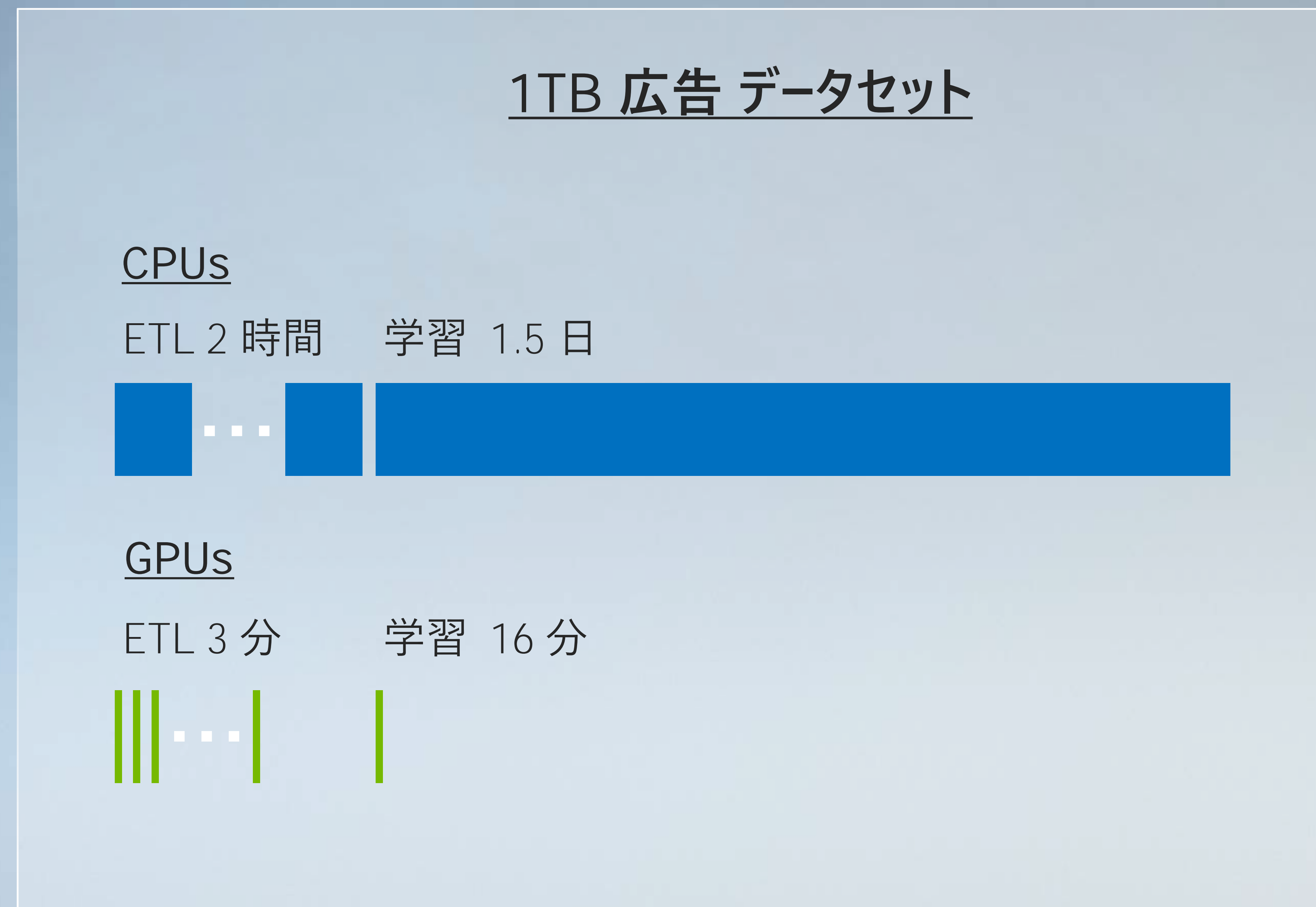
NVIDIA AI



レコメンダー システムはパーソナライズされたインターネットのエンジン



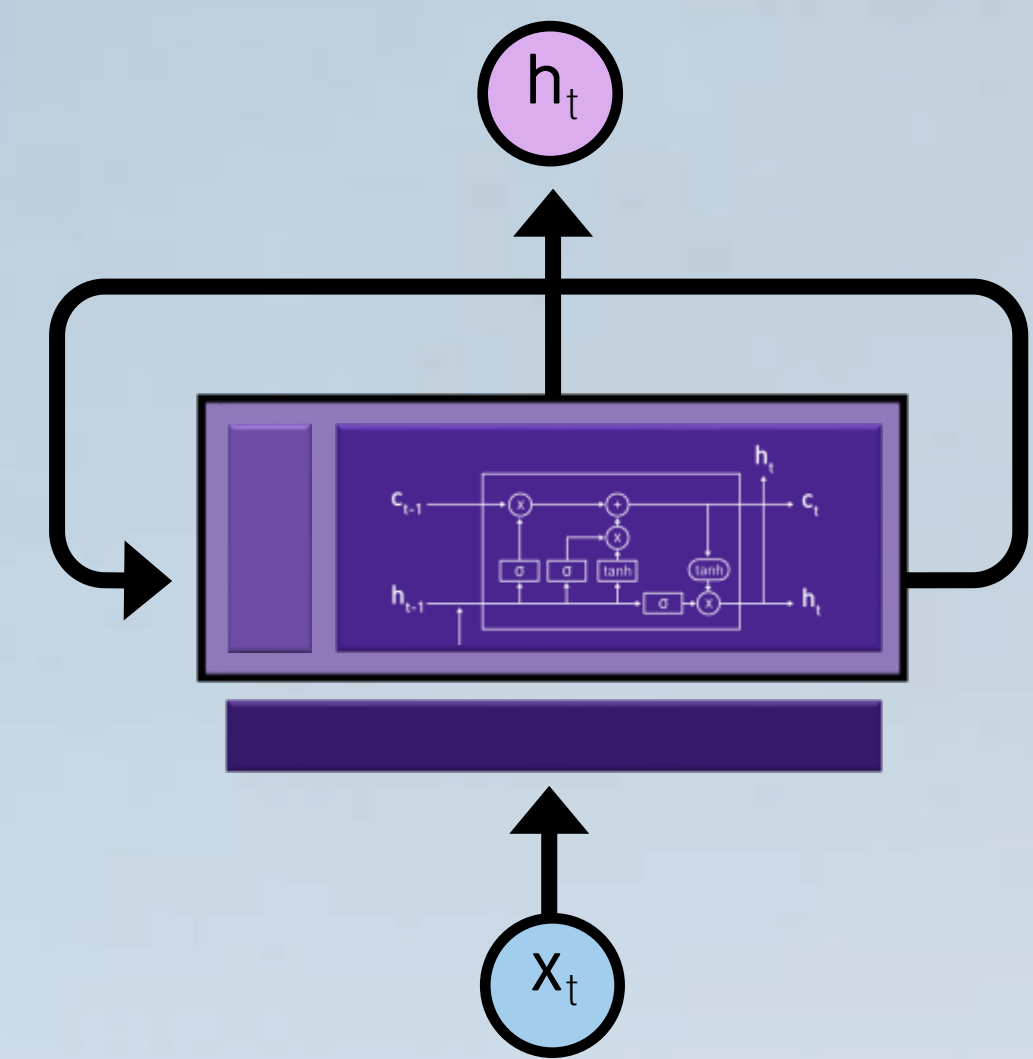
発表:
NVIDIA MERLIN — ディープレコメンダー アプリケーションフレームワーク



NVIDIA AI 推論

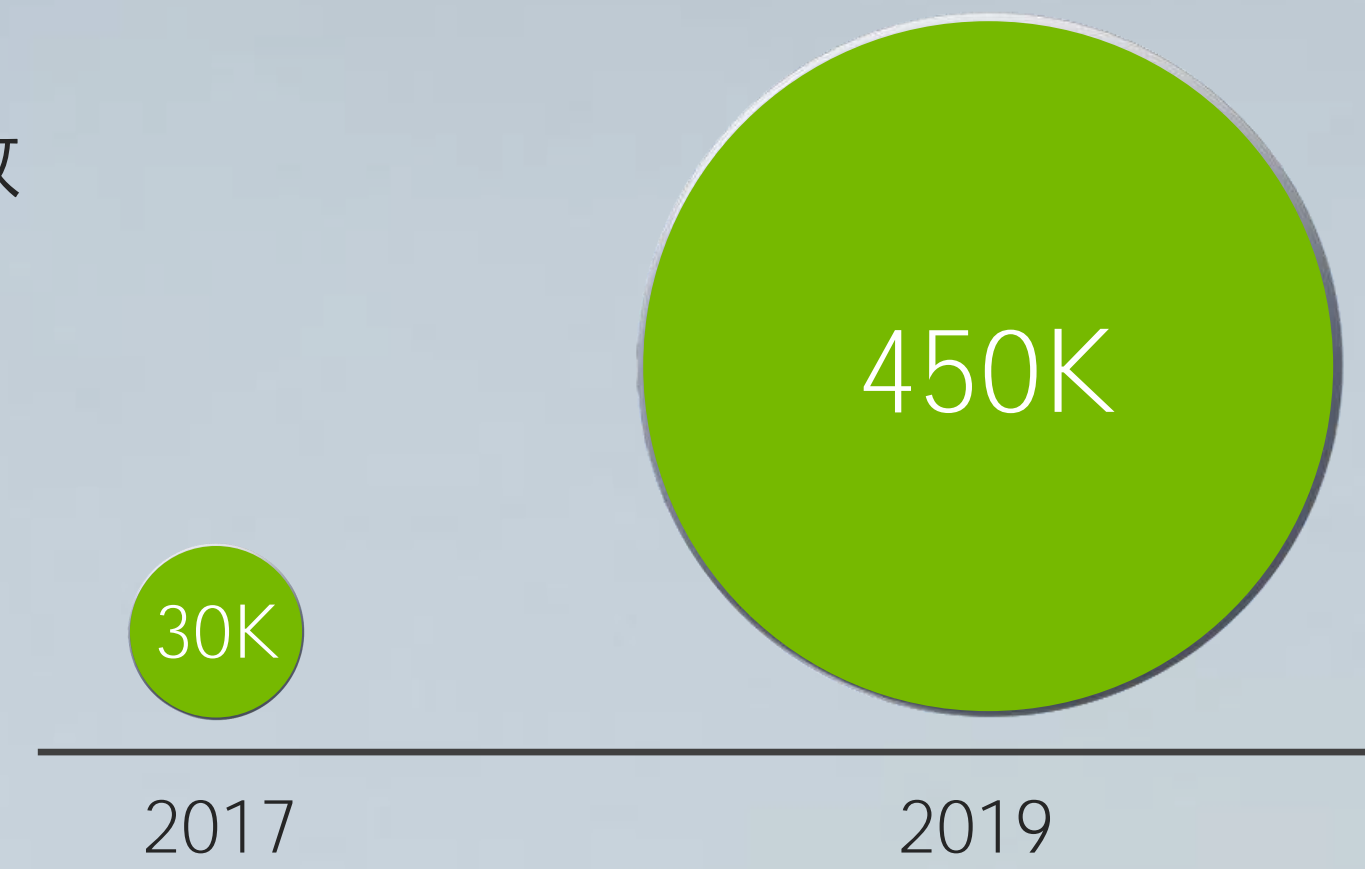
過去 2 年で TENSORRT のダウンロードが 15 倍

TENSORRT 7

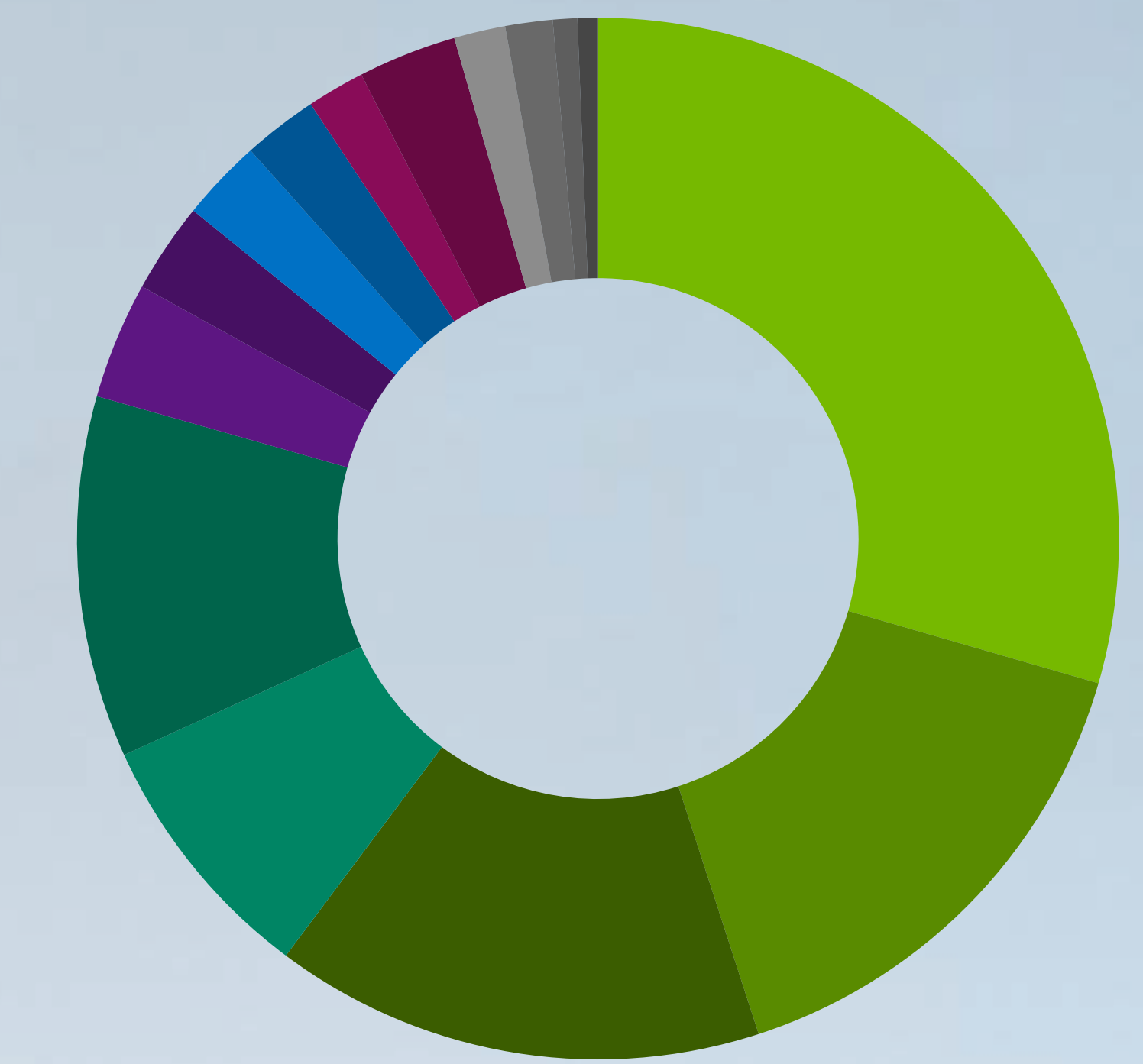
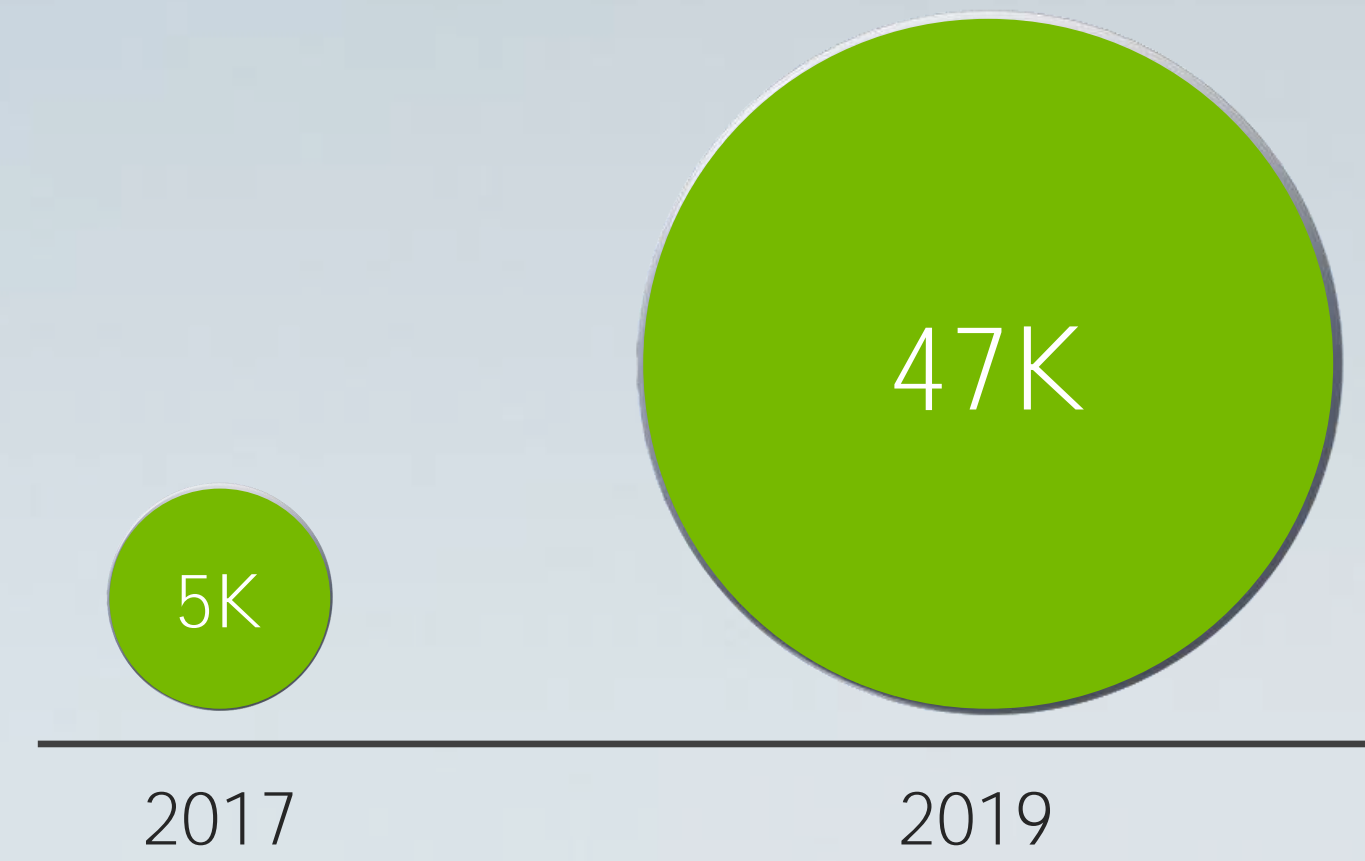


CNN、Transformer、RNN
1000 以上の最適化されたカーネル
自動混合精度 FP32、FP16、INT8

15 倍
ダウンロード数

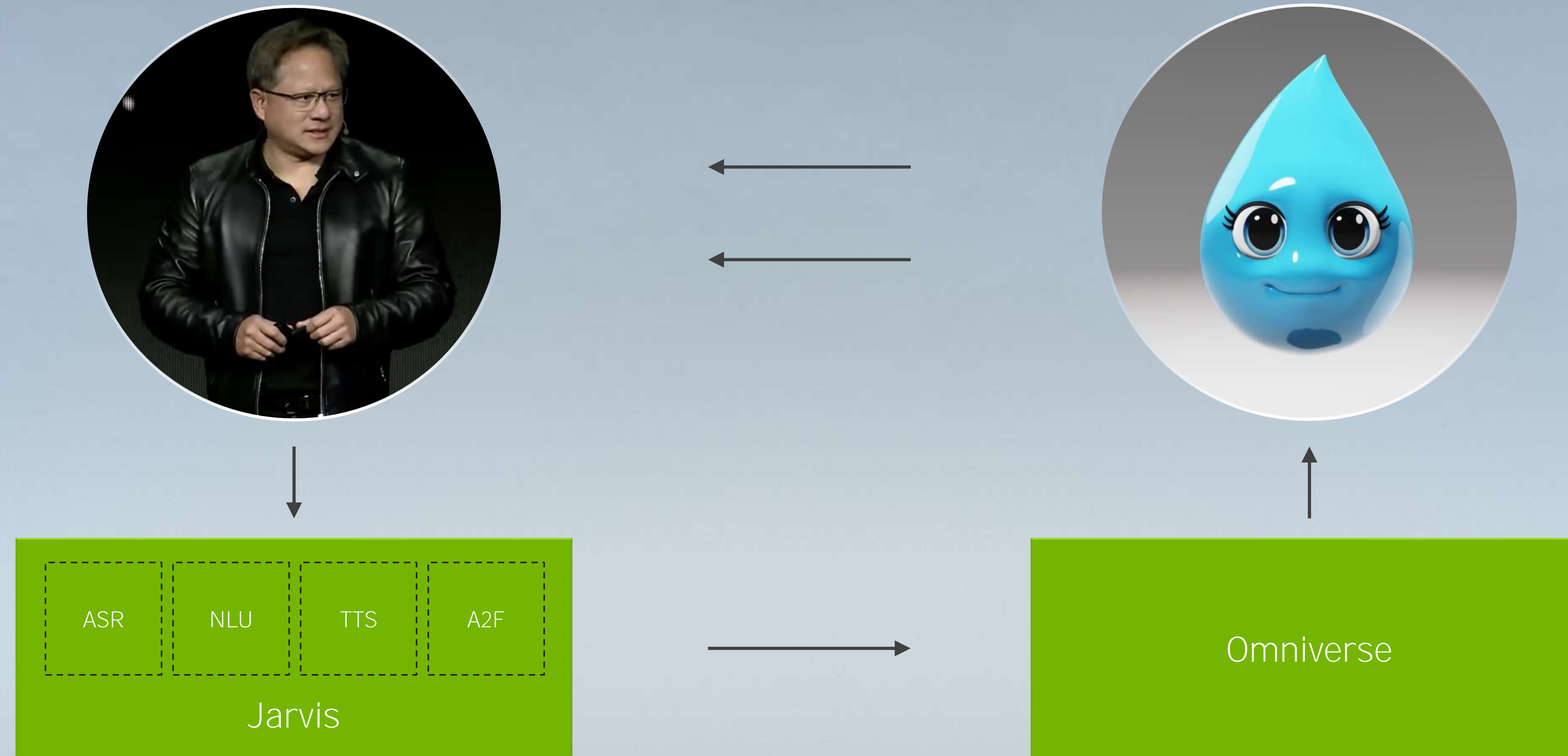


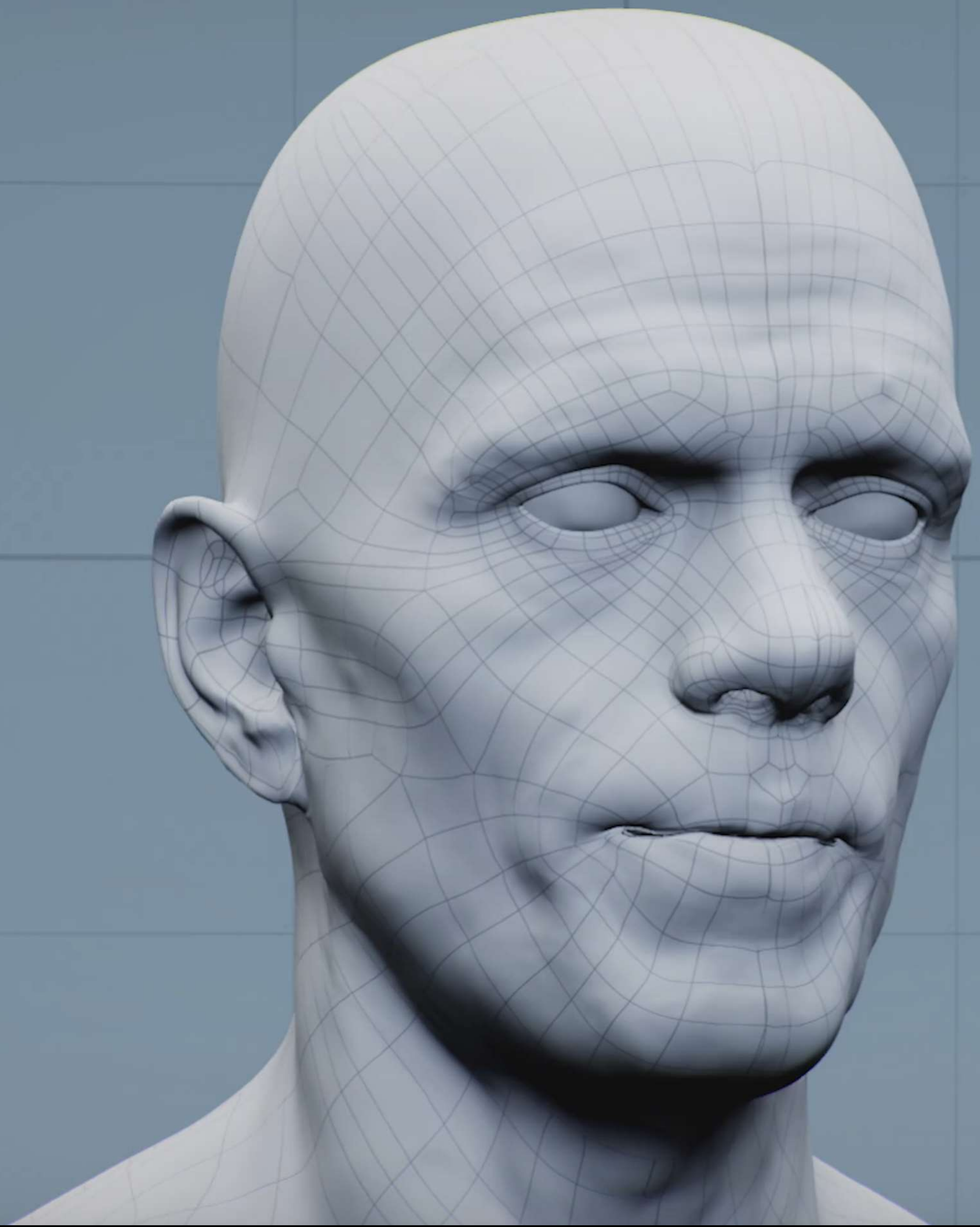
10 倍
開発者数



- Software
- IT Services
- Other
- Healthcare & Life Sciences
- Manufacturing
- Public Sector
- Consulting Services
- Research / Higher Ed
- Automotive
- Internet / Telecom
- Hardware / Semiconductor
- Cloud Services
- Financial Services
- Energy / Oil & Gas

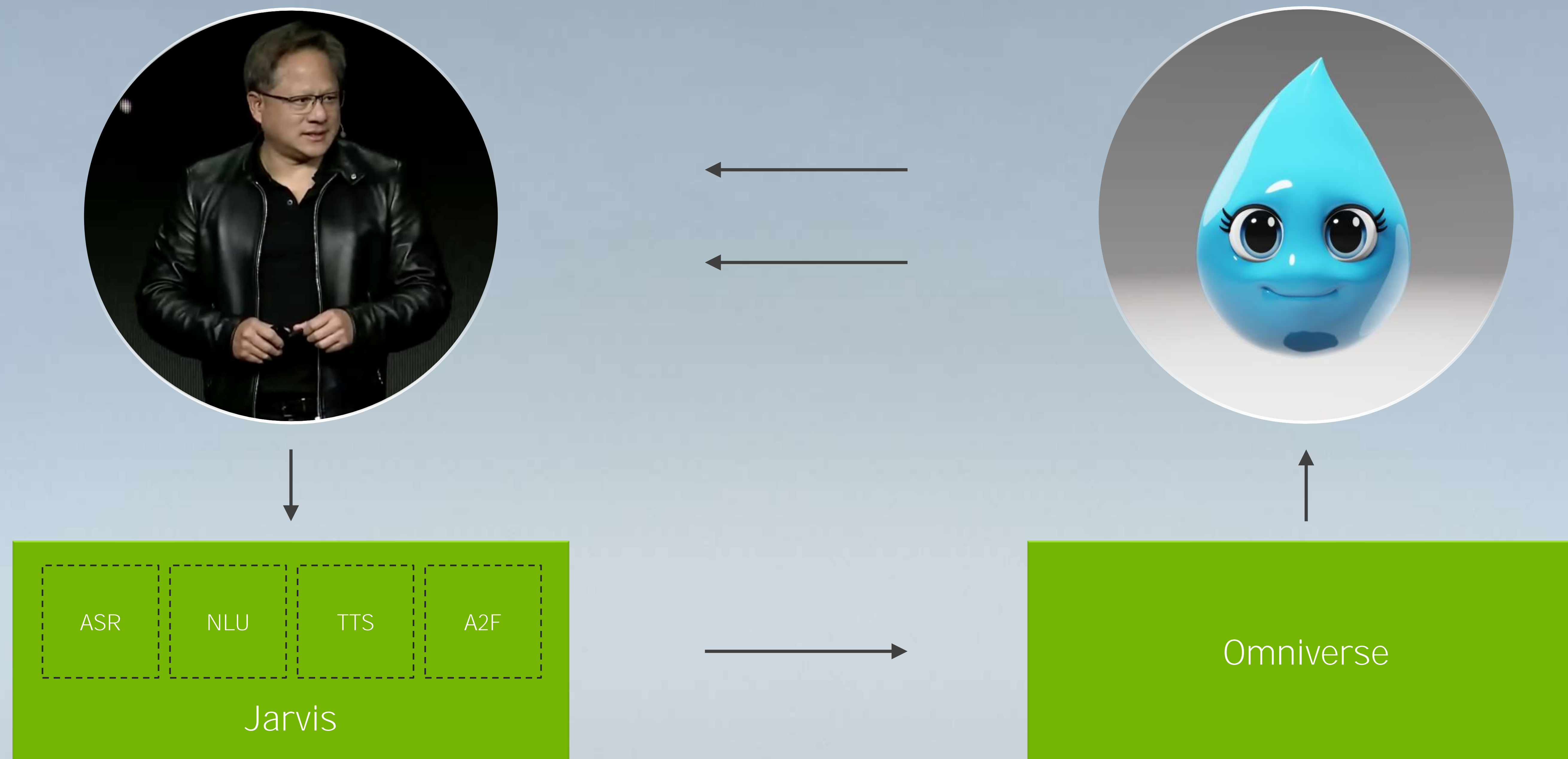
発表:
NVIDIA JARVIS —マルチモーダル 対話型 AI サービス フレームワーク





JARVIS デモ

MISTY — インタラクティブ 3D チャットボット





発表: NVIDIA JARVIS — マルチモーダル 対話型 AI サービス フレームワーク

NVIDIA JARVIS



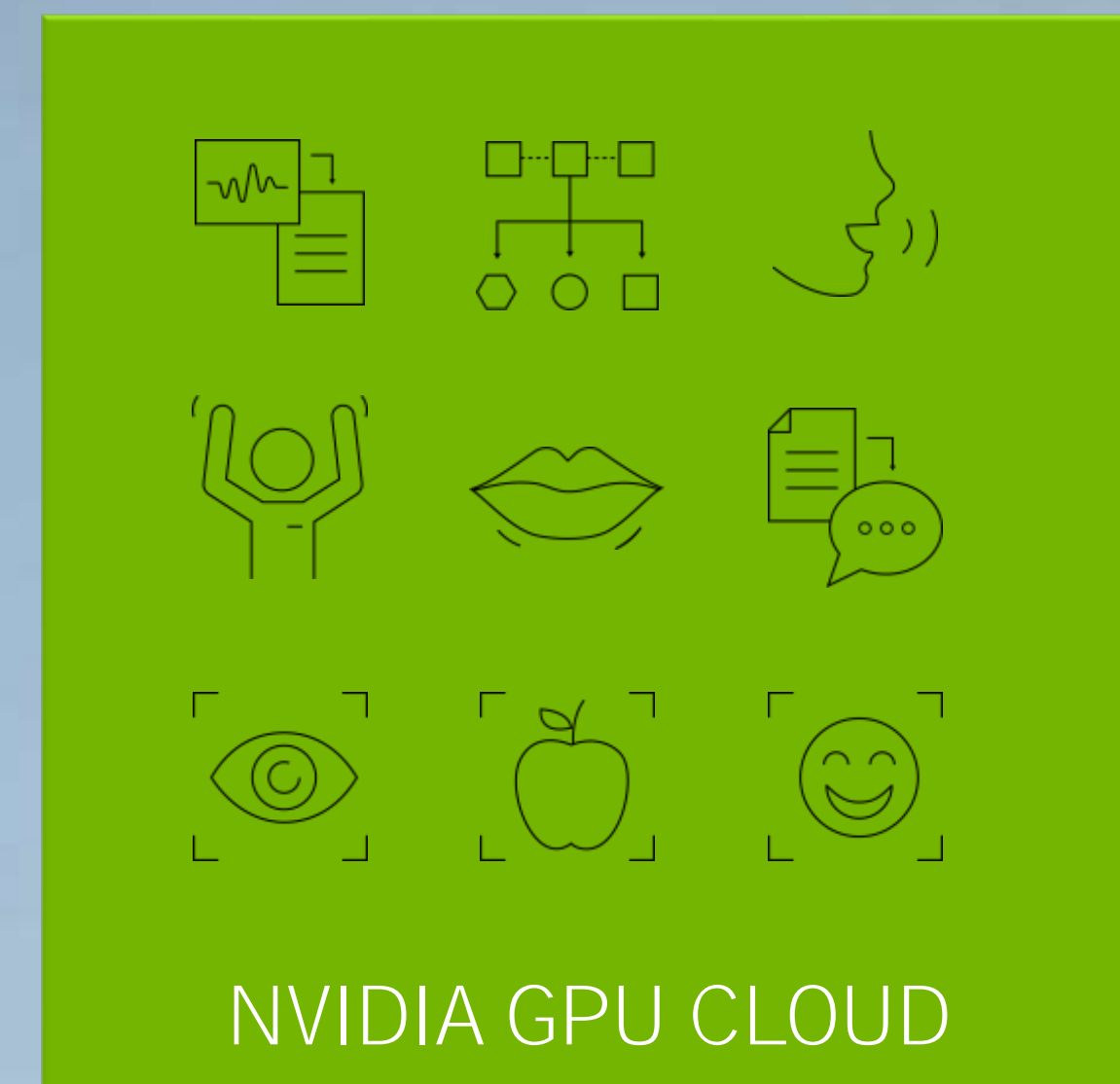
ビデオ
音声

マルチスピーカー
文字起こし

JESSICA: What will you have ready for Wednesday?
DOUGLAS: I expect to have early designs of the packaging.
JESSICA: Great.

発表: NVIDIA JARVIS —マルチモーダル 対話型 AI サービス フレームワーク

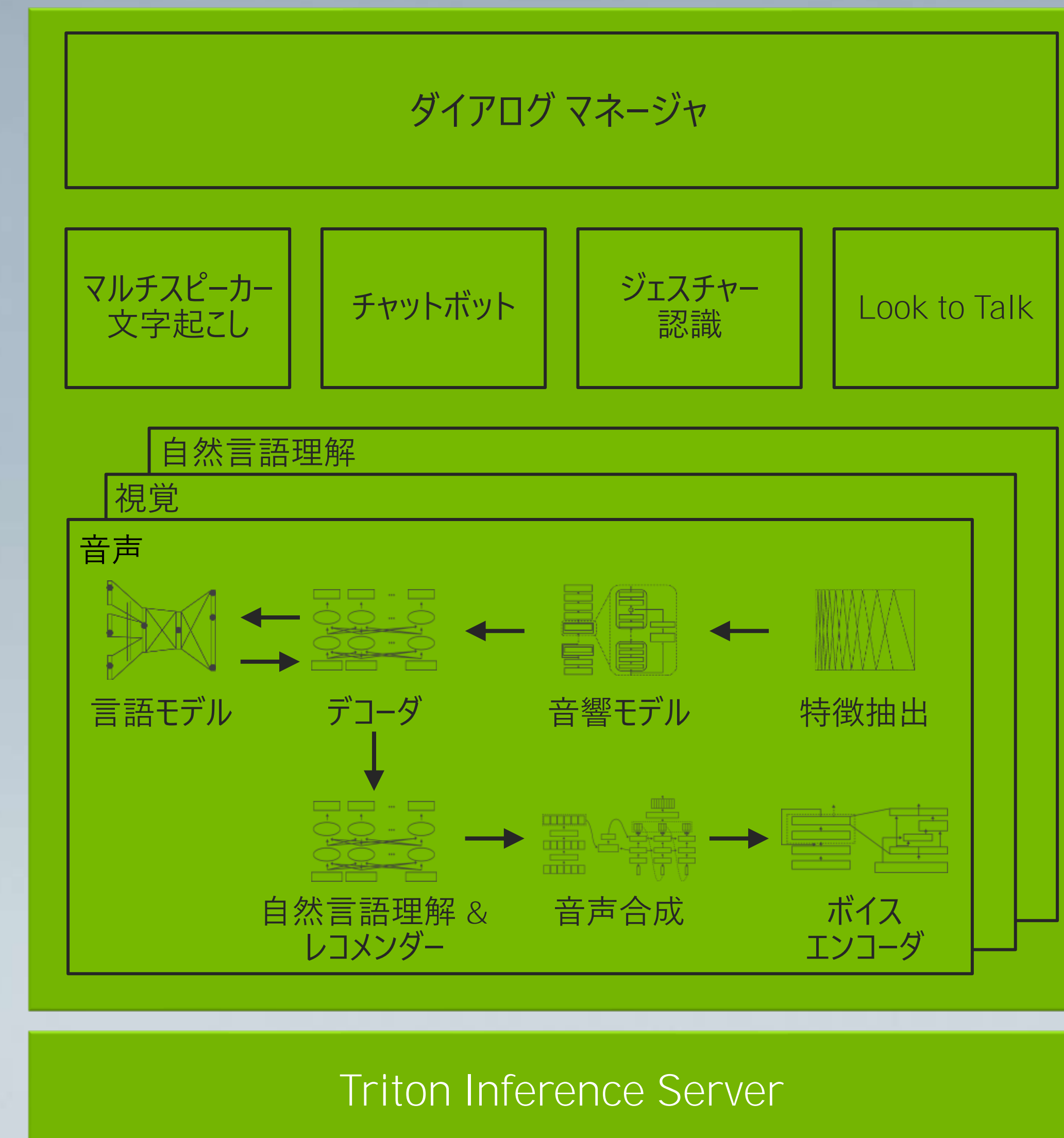
事前学習モデル



再学習



NVIDIA JARVIS



JESSICA: What will you have ready for Wednesday?

DOUGLAS: I expect to have early designs of the packaging.

JESSICA: Great.

早期アクセス プログラムに参加
developer.nvidia.com/nvidia-jarvis

対話型 AI が産業を変革中



ビデオ会議
字幕、翻訳、文字起こし
1日あたり2億回の会議



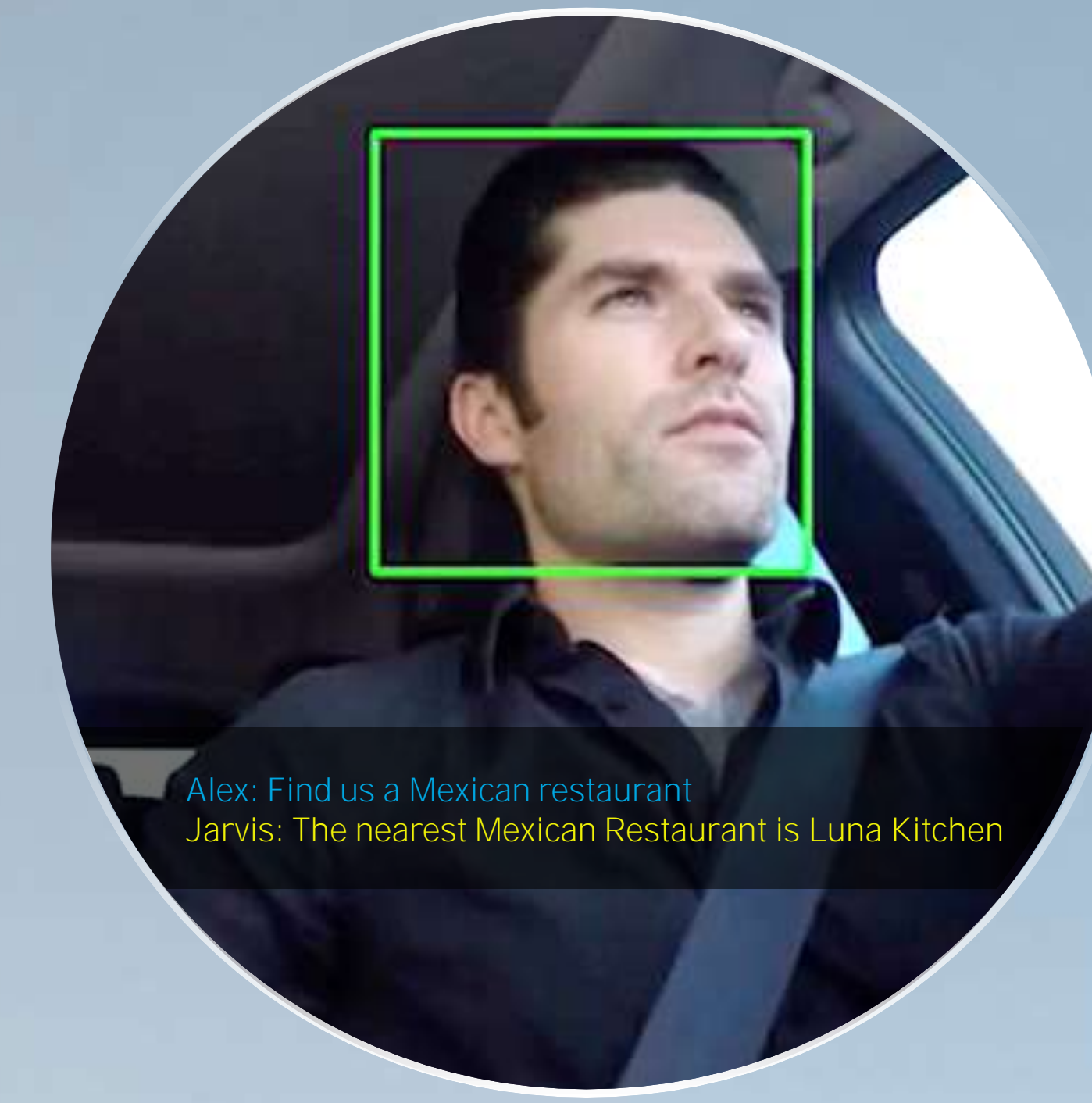
コールセンター
1日あたり5億コール



スマートスピーカー
年間1億5000台販売



リテールアシスタント
1200万の小売店



車内アシスタント
年間750万台の新車

KENSHO + S&P Global

Microsoft

NUANCE

Square

voca.ai



fast.ai

KALDI

ESPnet

spaCy



Alibaba Cloud

aws

Baidu 百度

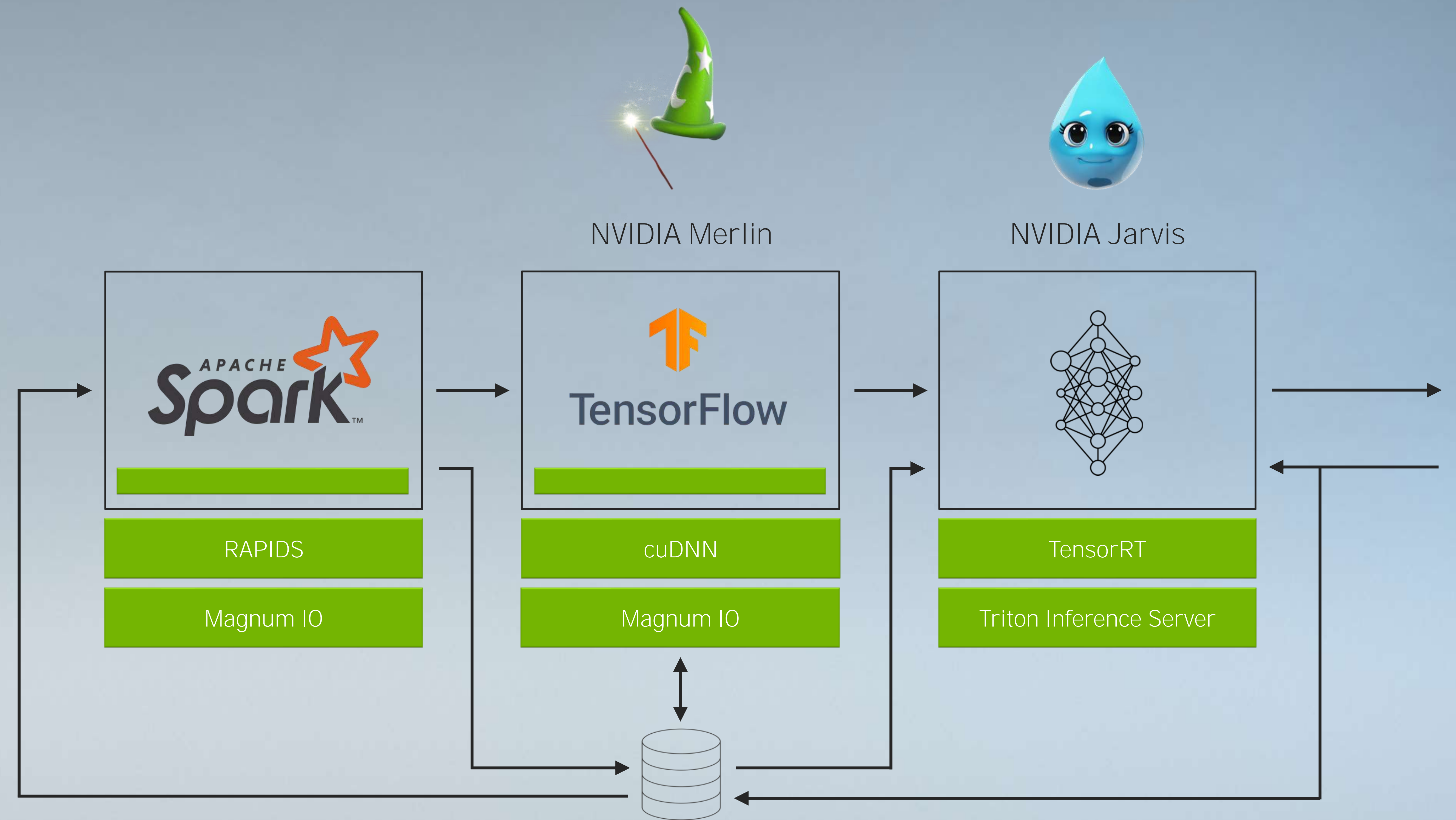
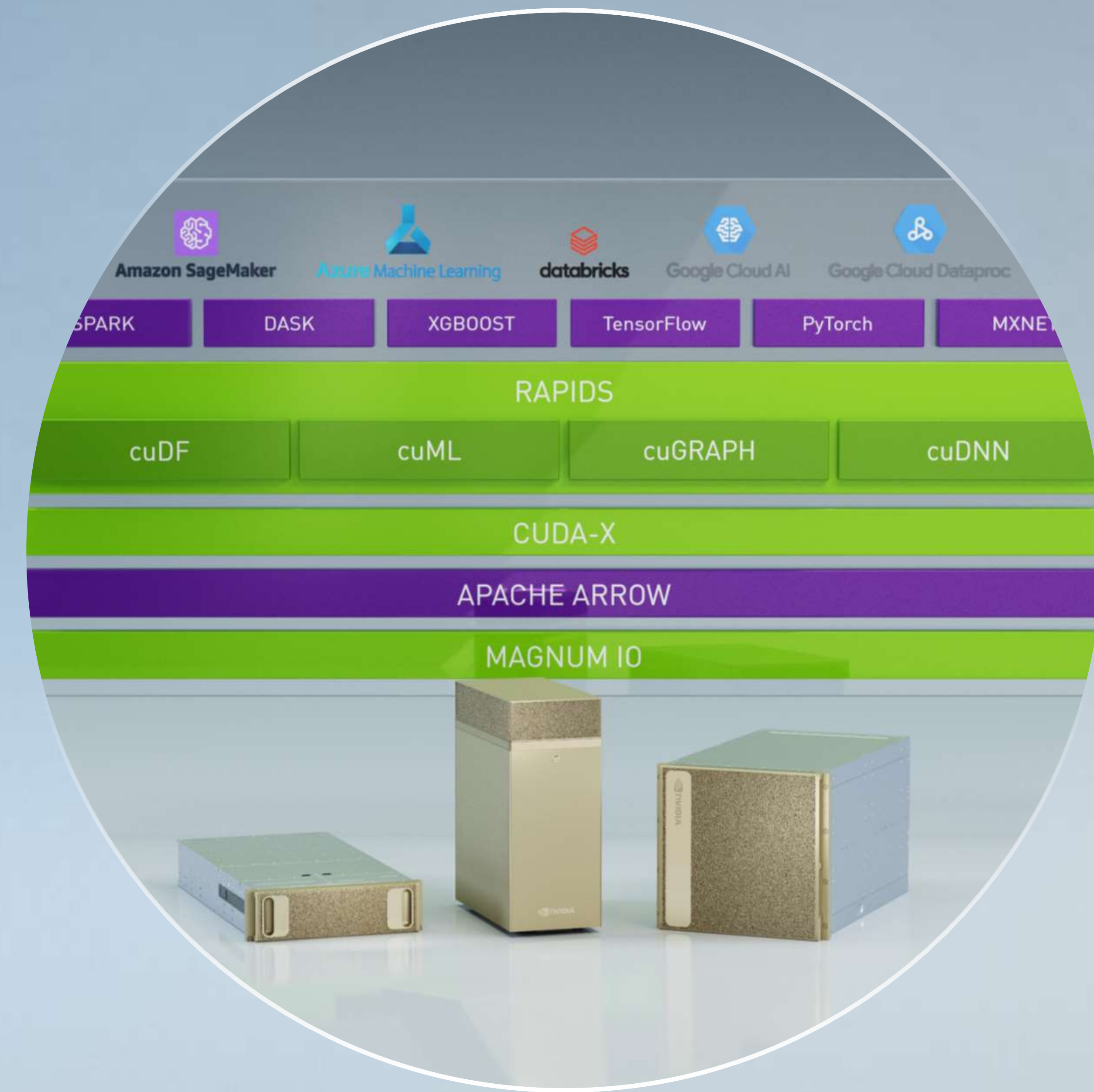
Google Cloud

Microsoft Azure

ORACLE
Cloud Infrastructure

Tencent
Cloud

NVIDIA AI

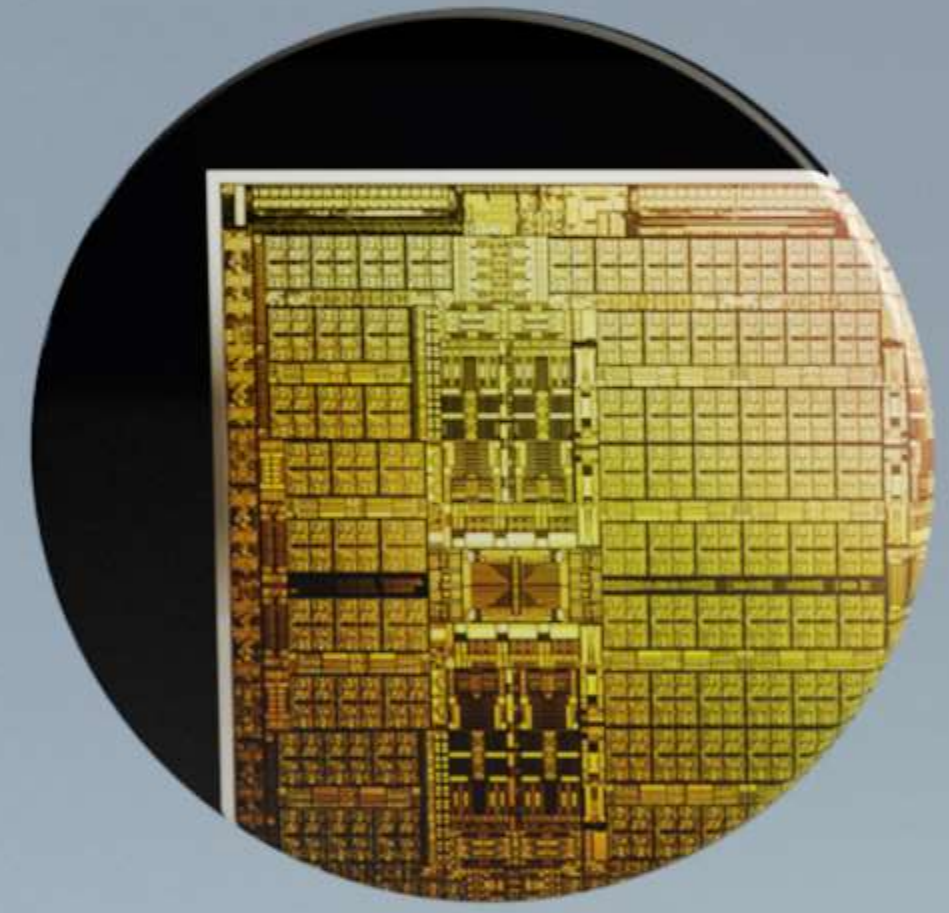


現代のクラウド データ センター

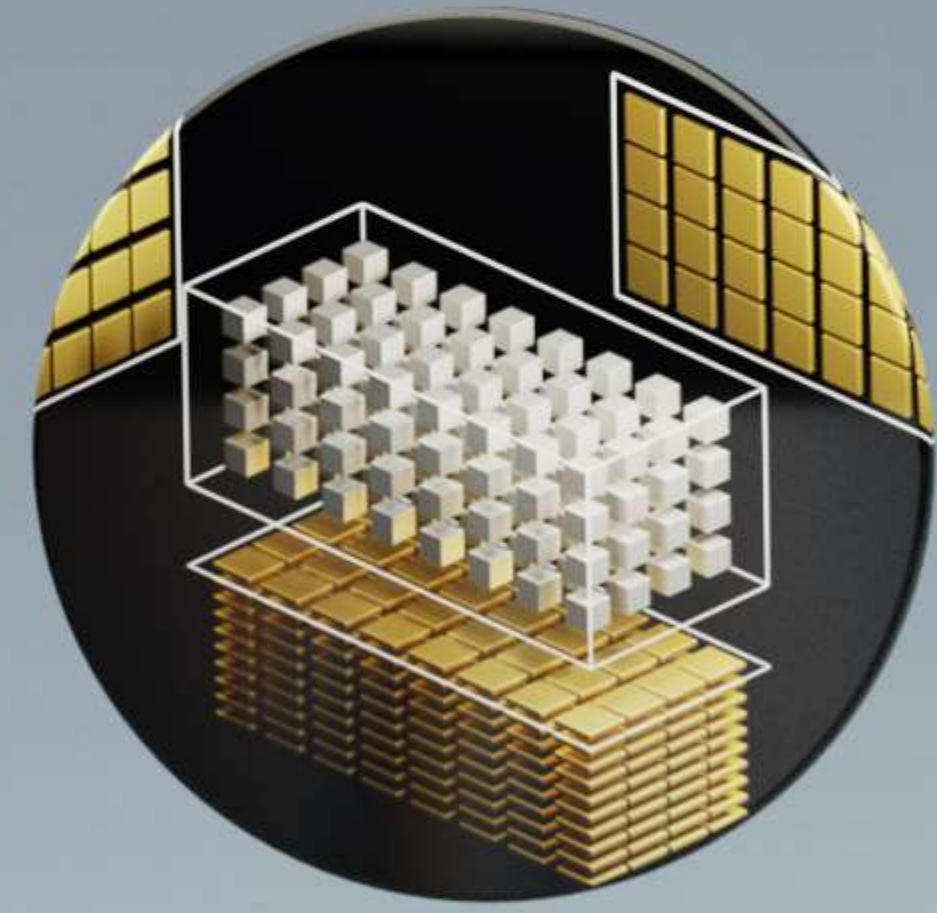
多様なアプリケーション | スケールアップとスケールアウトのワークロード | 飽くなき需要



発表: NVIDIA A100



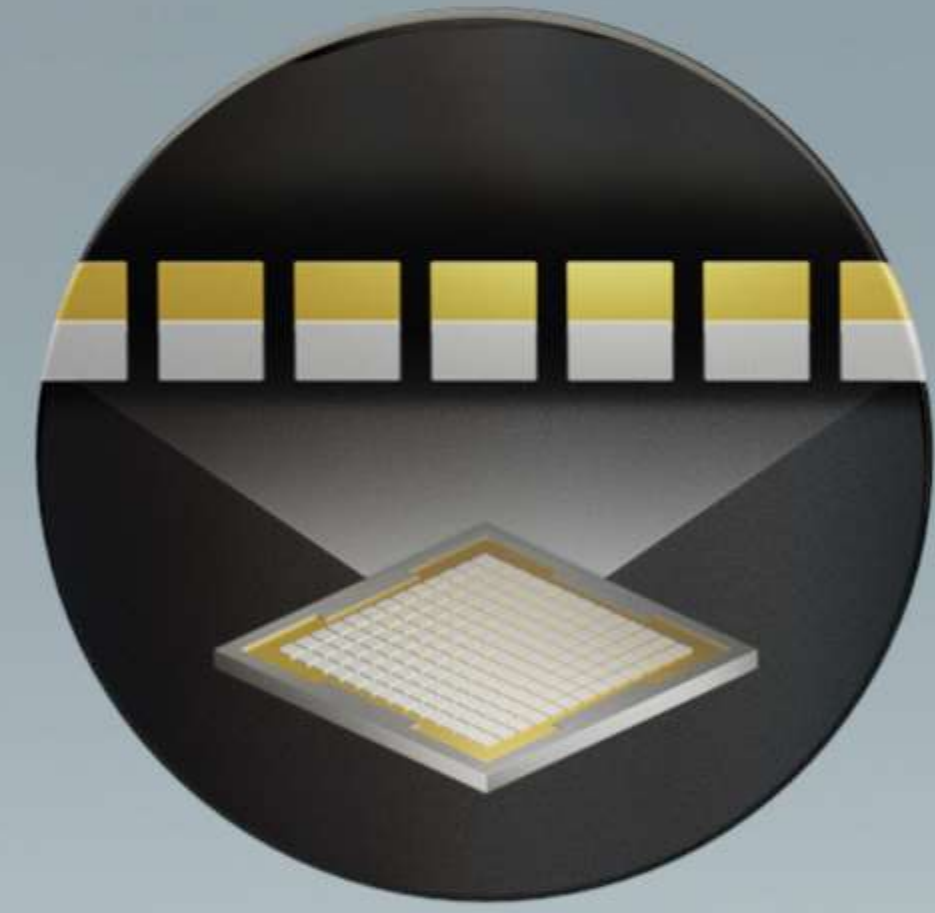
540 億トランジスタ



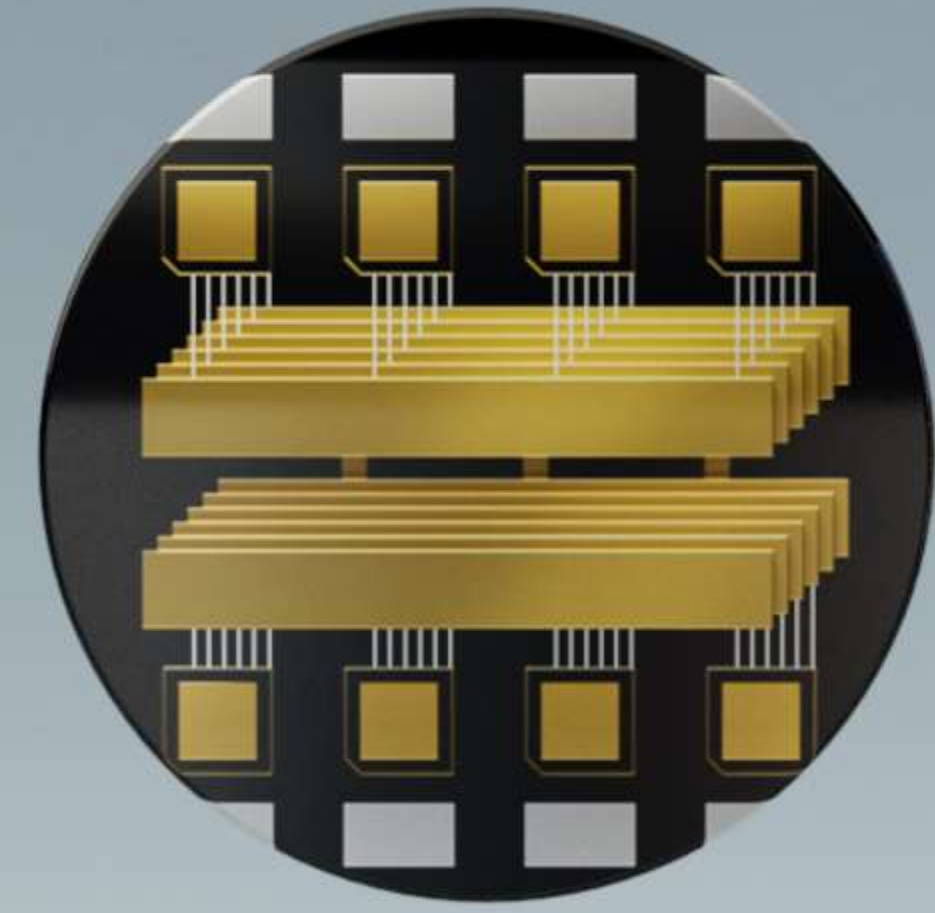
第 3 世代 TENSOR コア



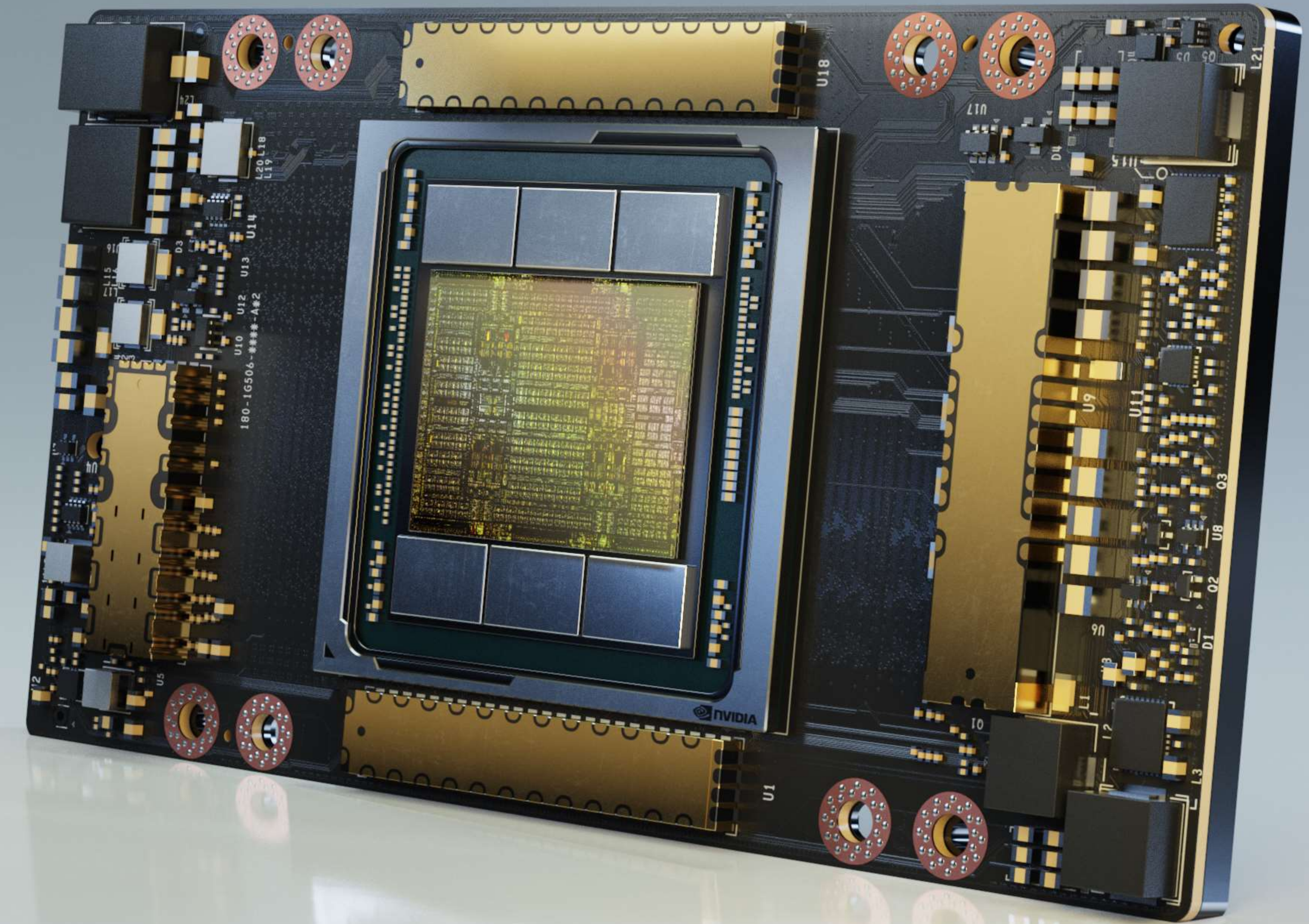
スパースティ アクセラレーション



MIG

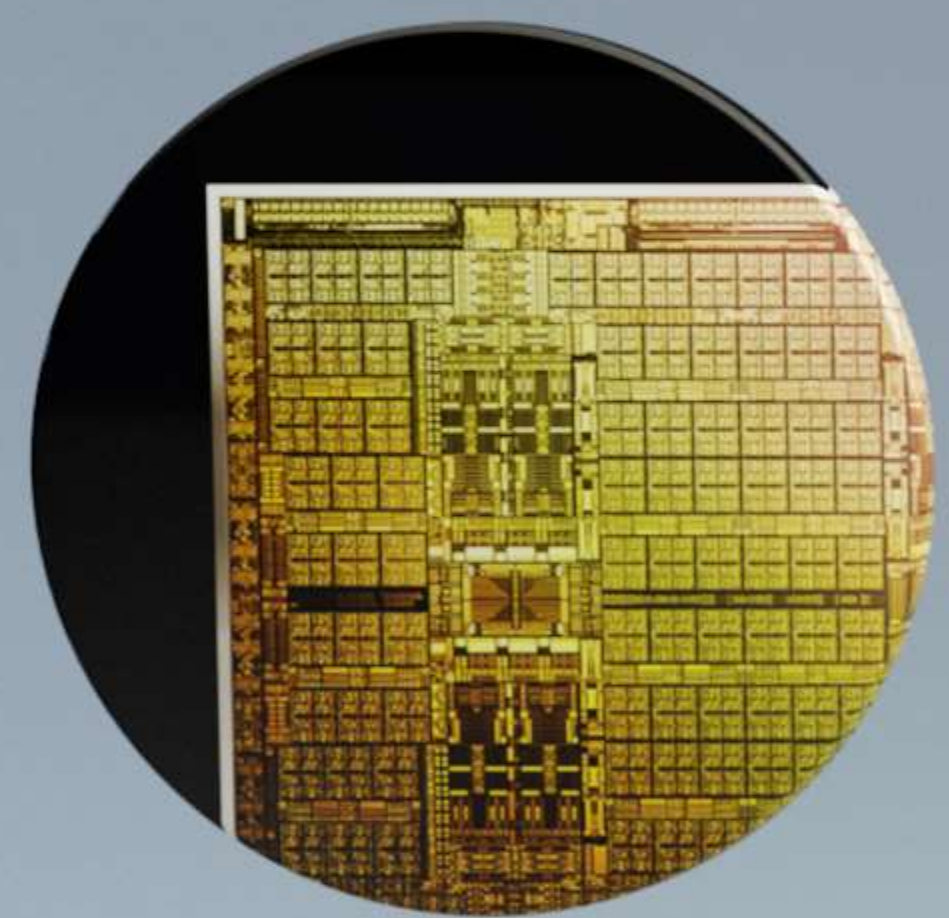


第 3 世代 NVLINK & NVSWITCH

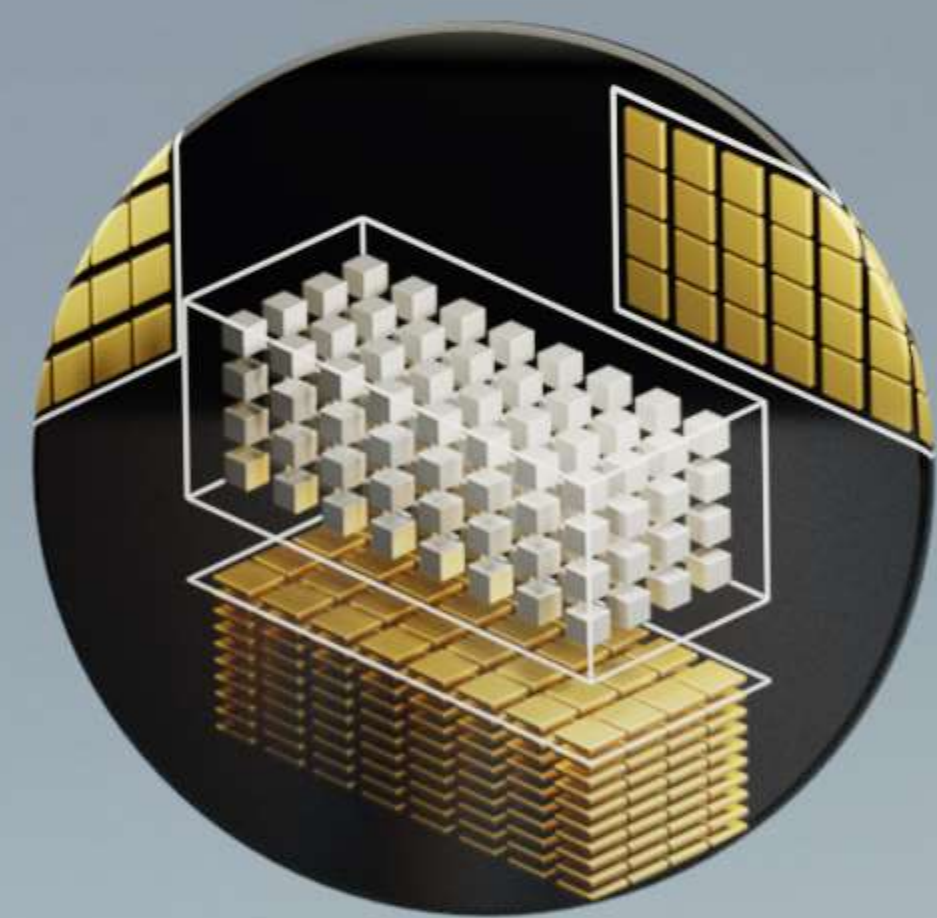


発表: NVIDIA A100

TSMC 7nm | HBM2 — 1.6 テラバイト / 秒 | 3D チップ スタック



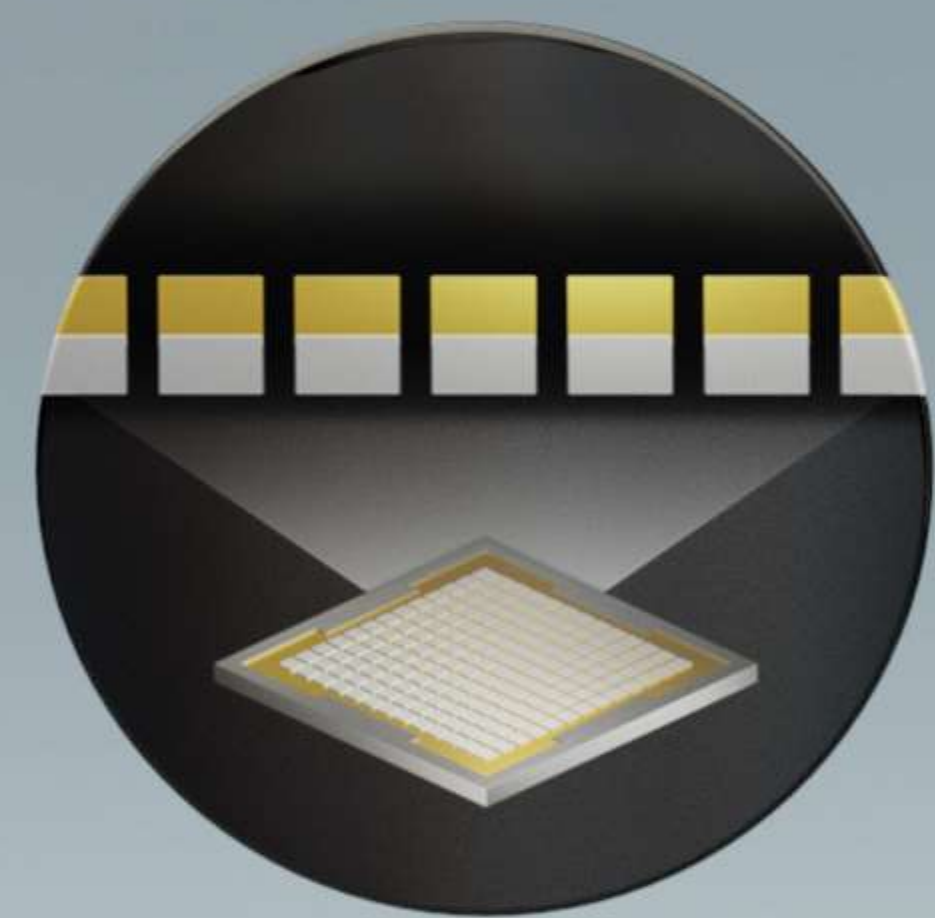
540 億トランジスタ



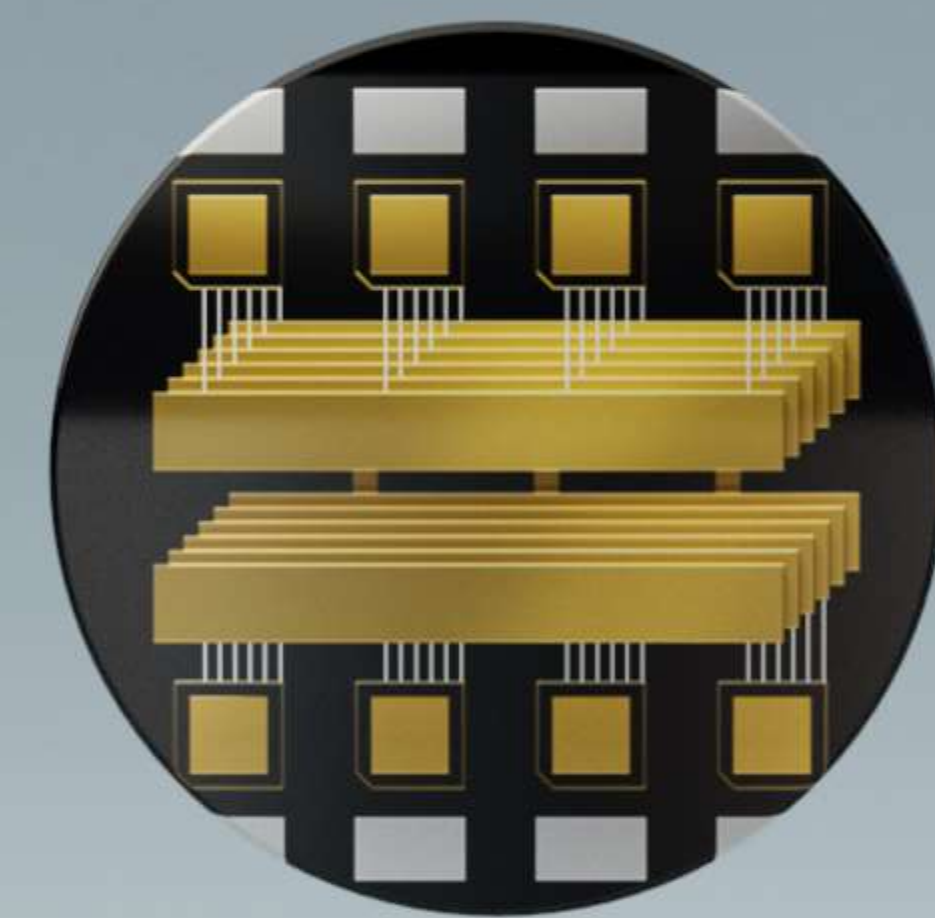
第 3 世代 TENSOR コア



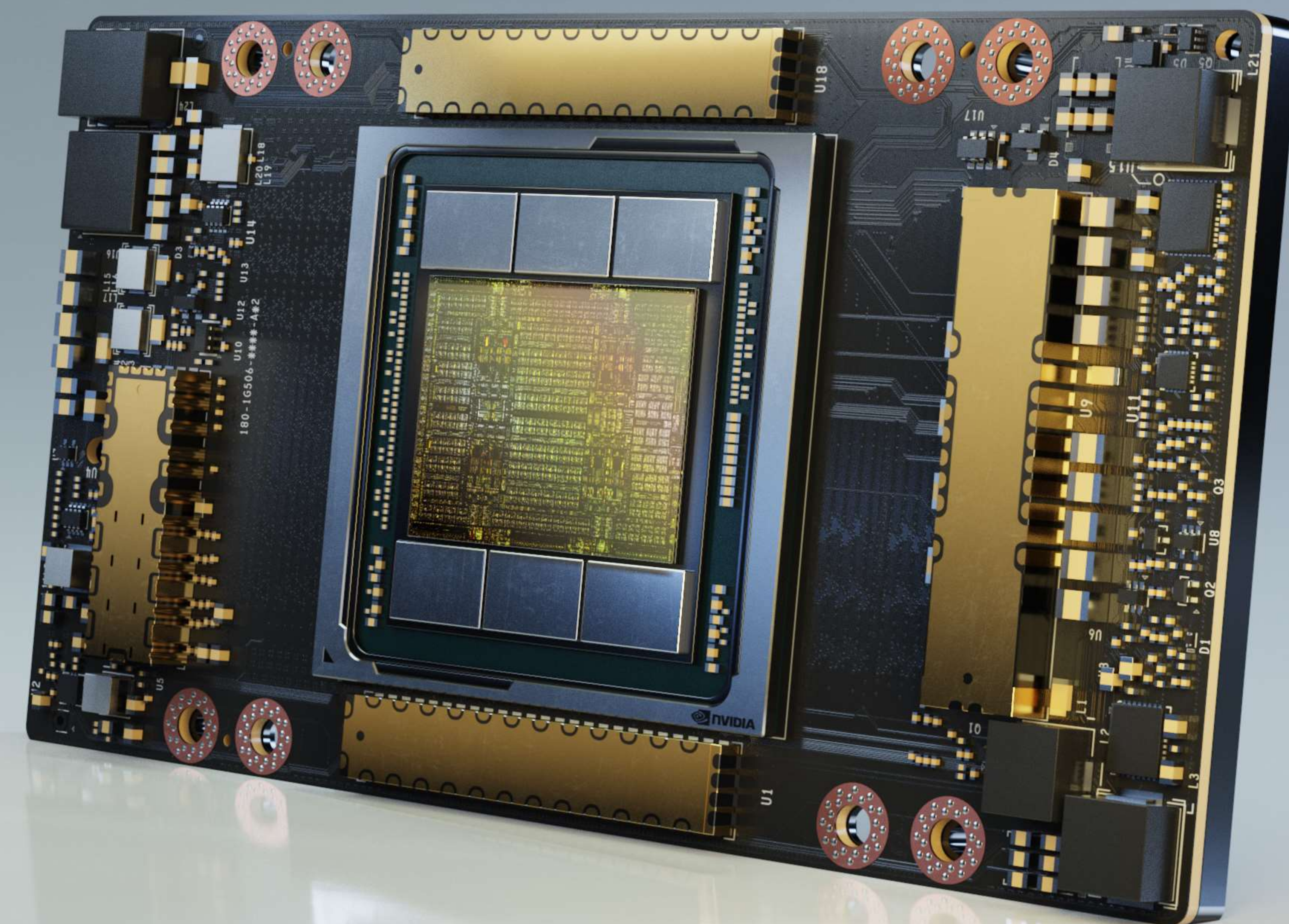
スパースティ アクセラレーション



MIG

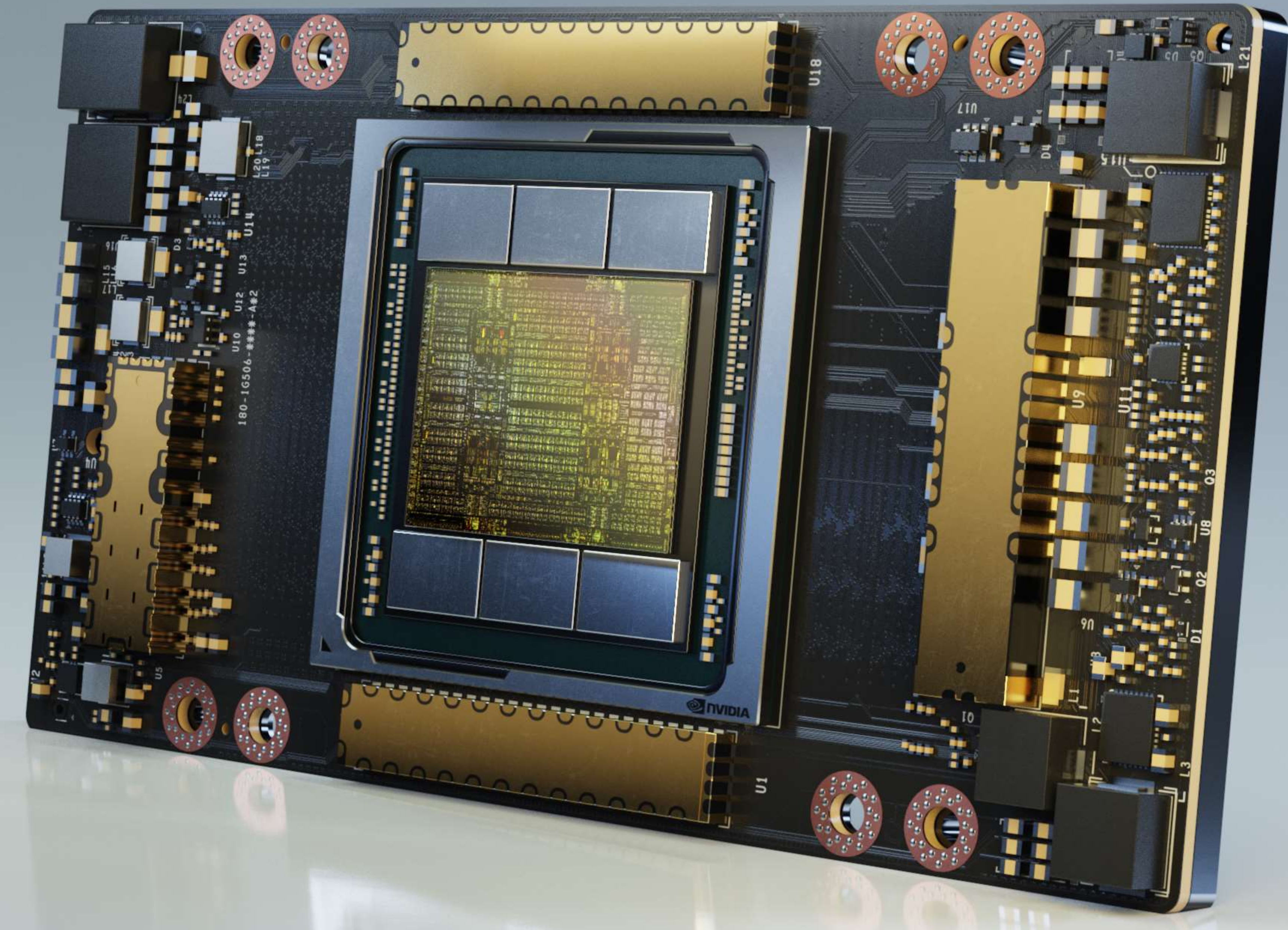
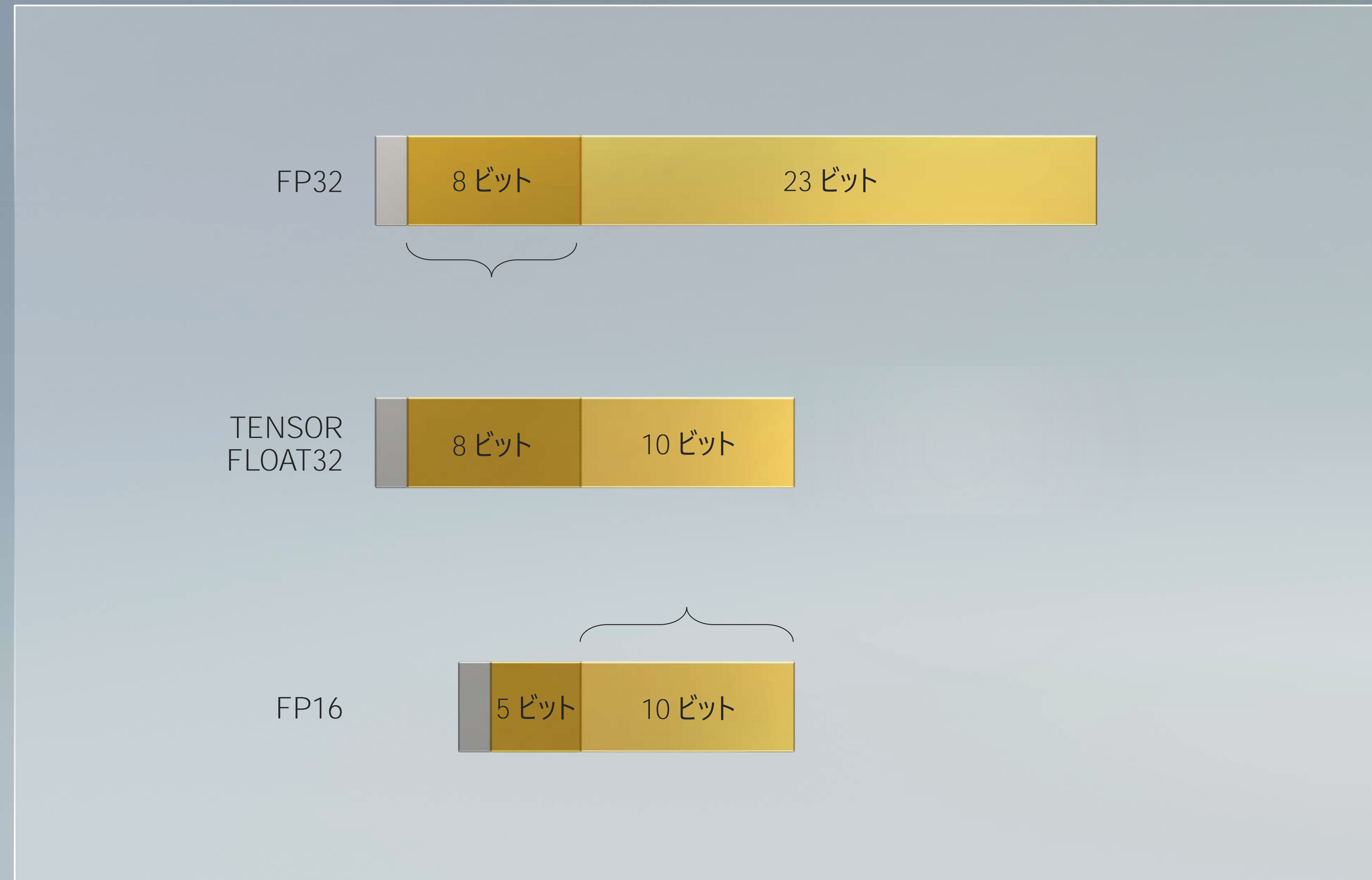


第 3 世代 NVLINK & NVSWITCH



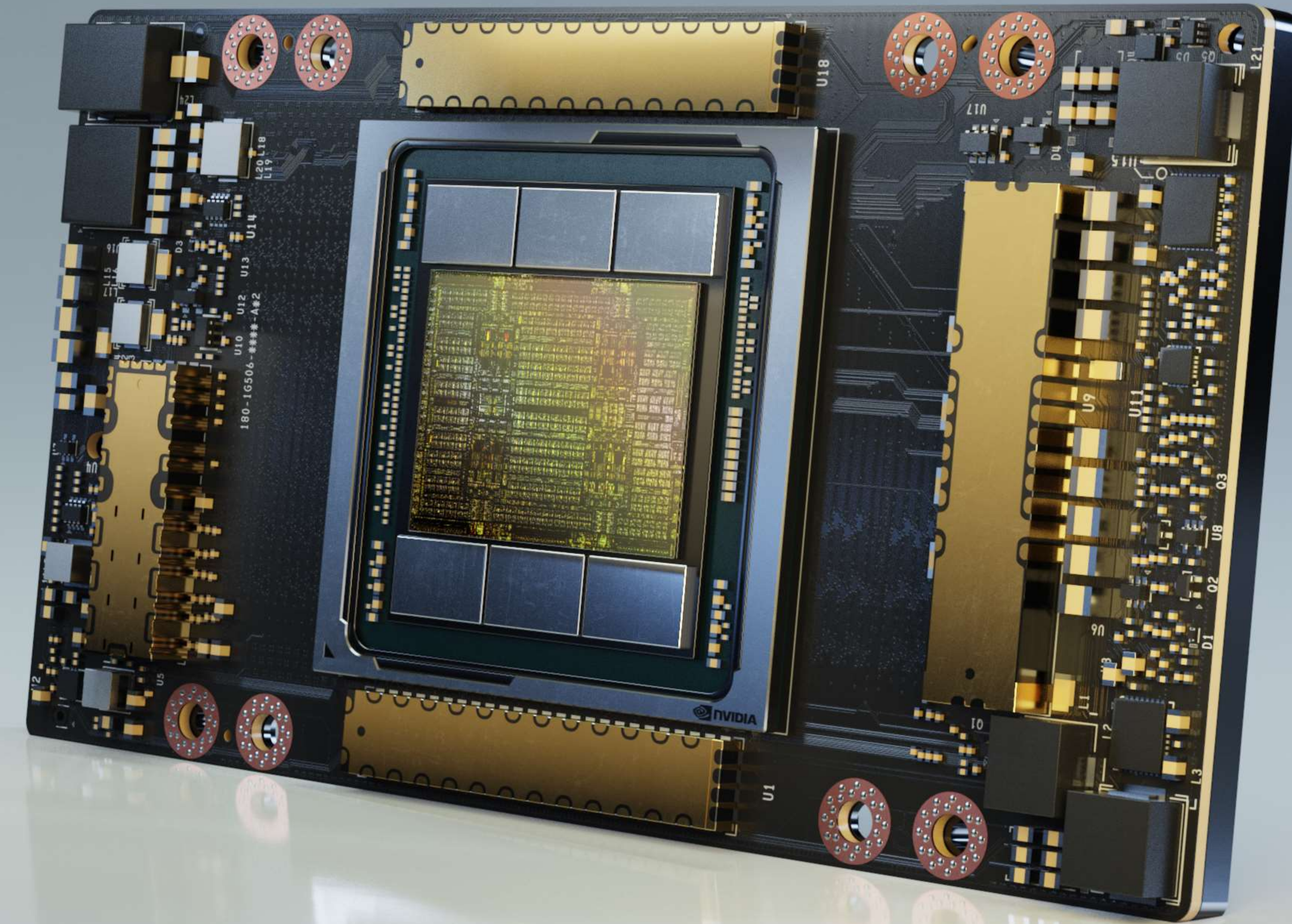
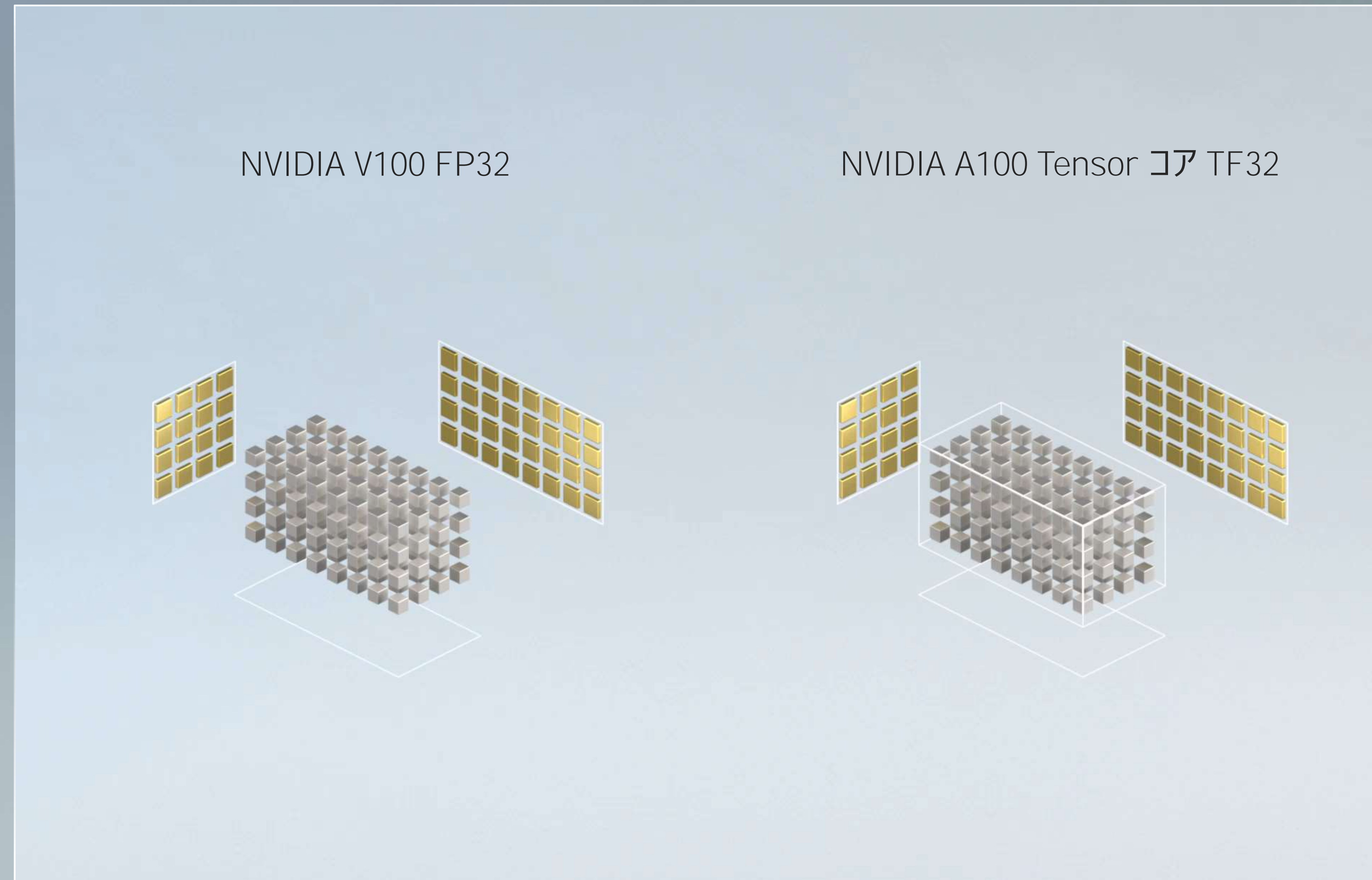
新しい TF32 TENSOR コア

FP32 の範囲と FP16 の精度 | FP32 の入力と FP32 でのアキュムレーション | コードの変更なしで学習を高速化



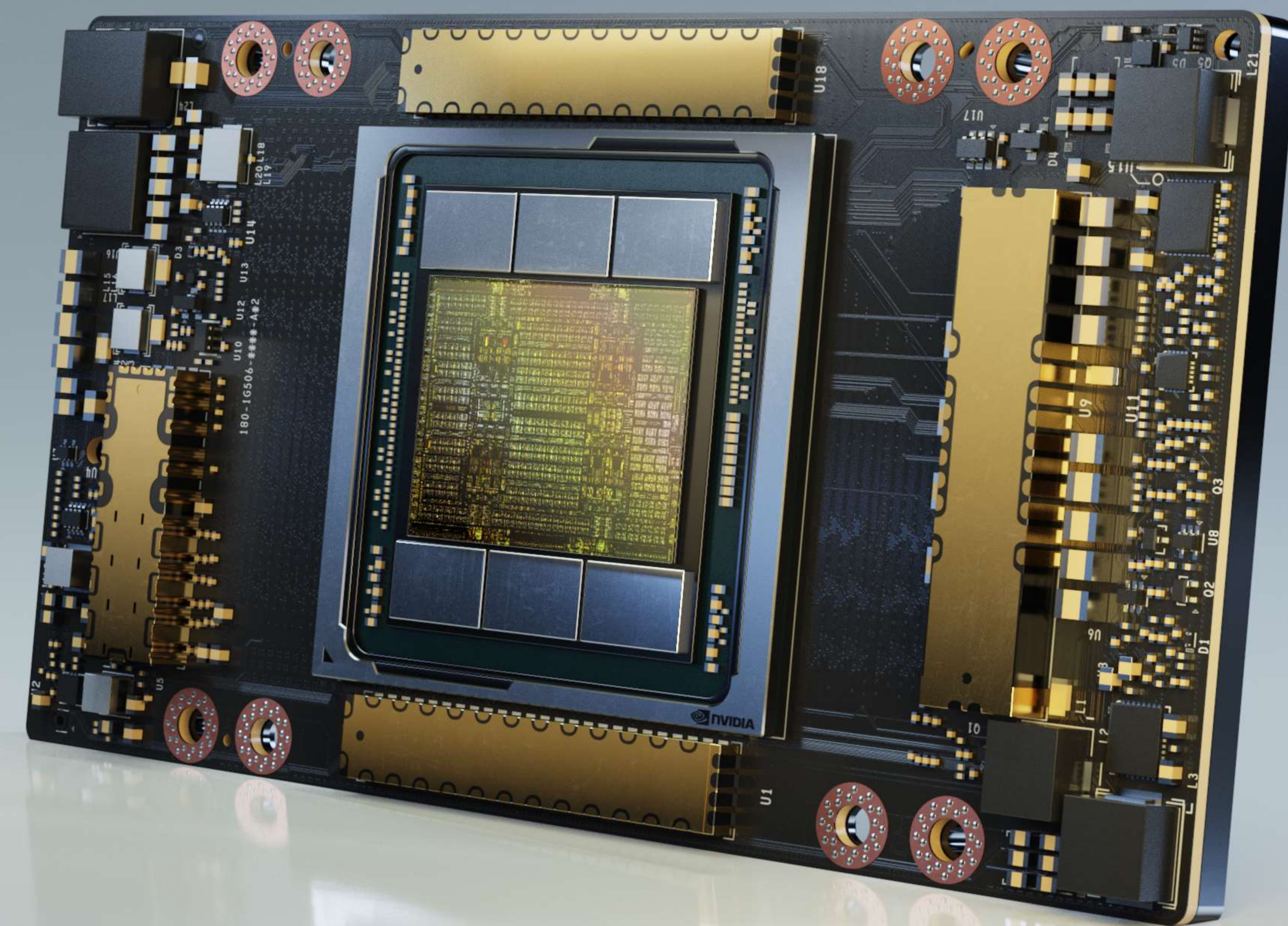
新しい TF32 TENSOR コア

FP32 の範囲と FP16 の精度 | FP32 の入力と FP32 でのアキュムレーション | コードの変更なしで学習を高速化

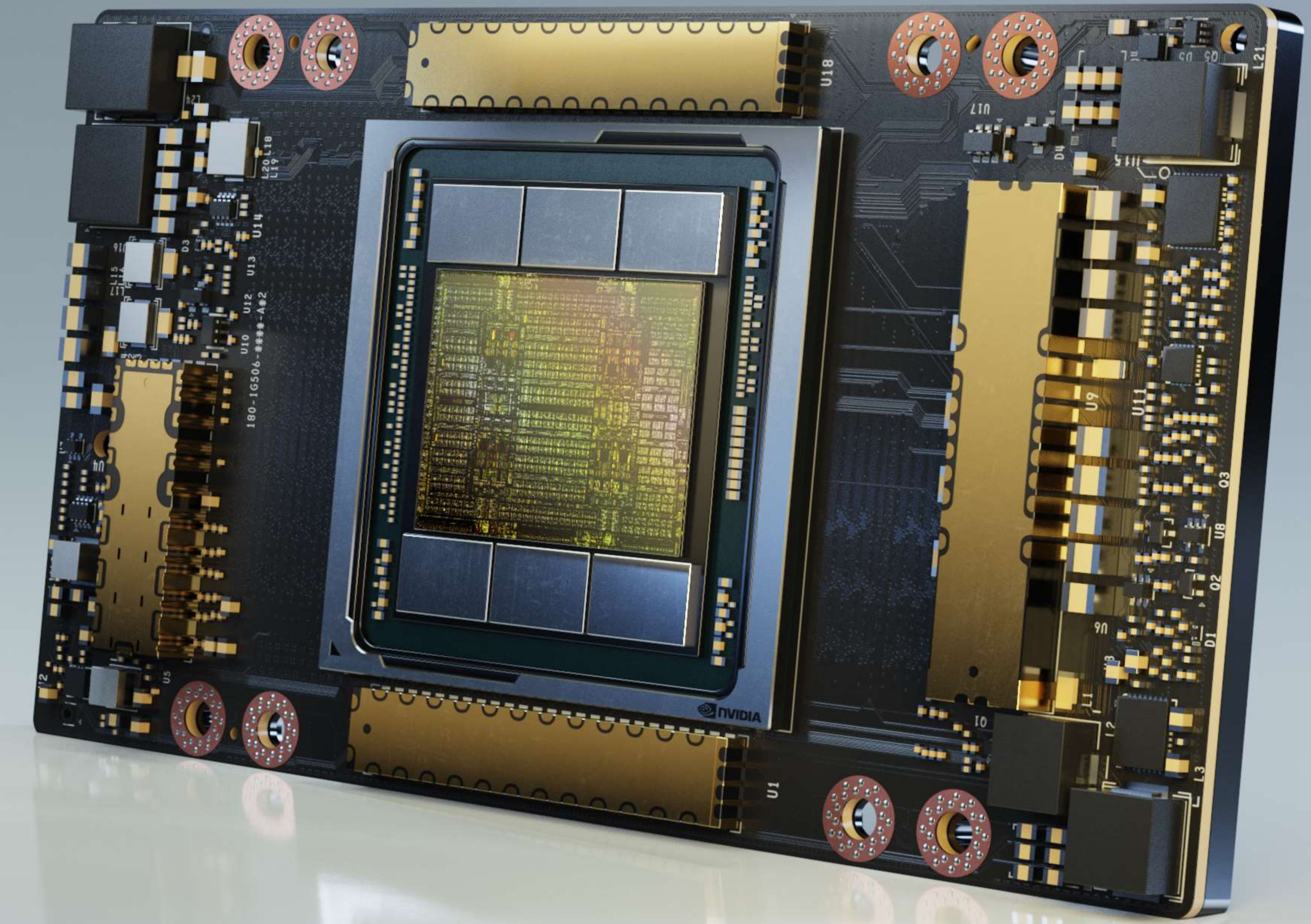
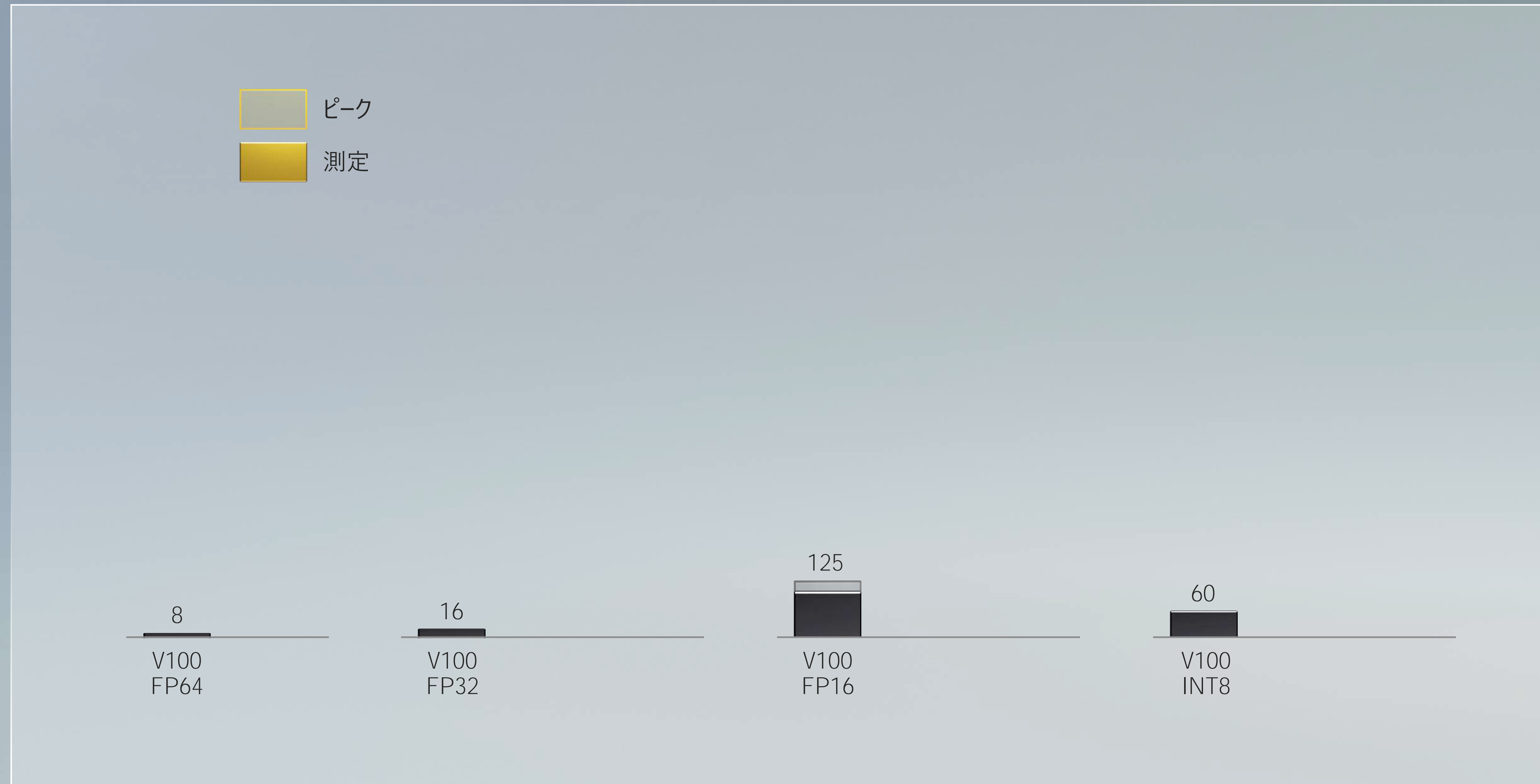


疎行列向けの新しい TENSOR コア アクセラレーション

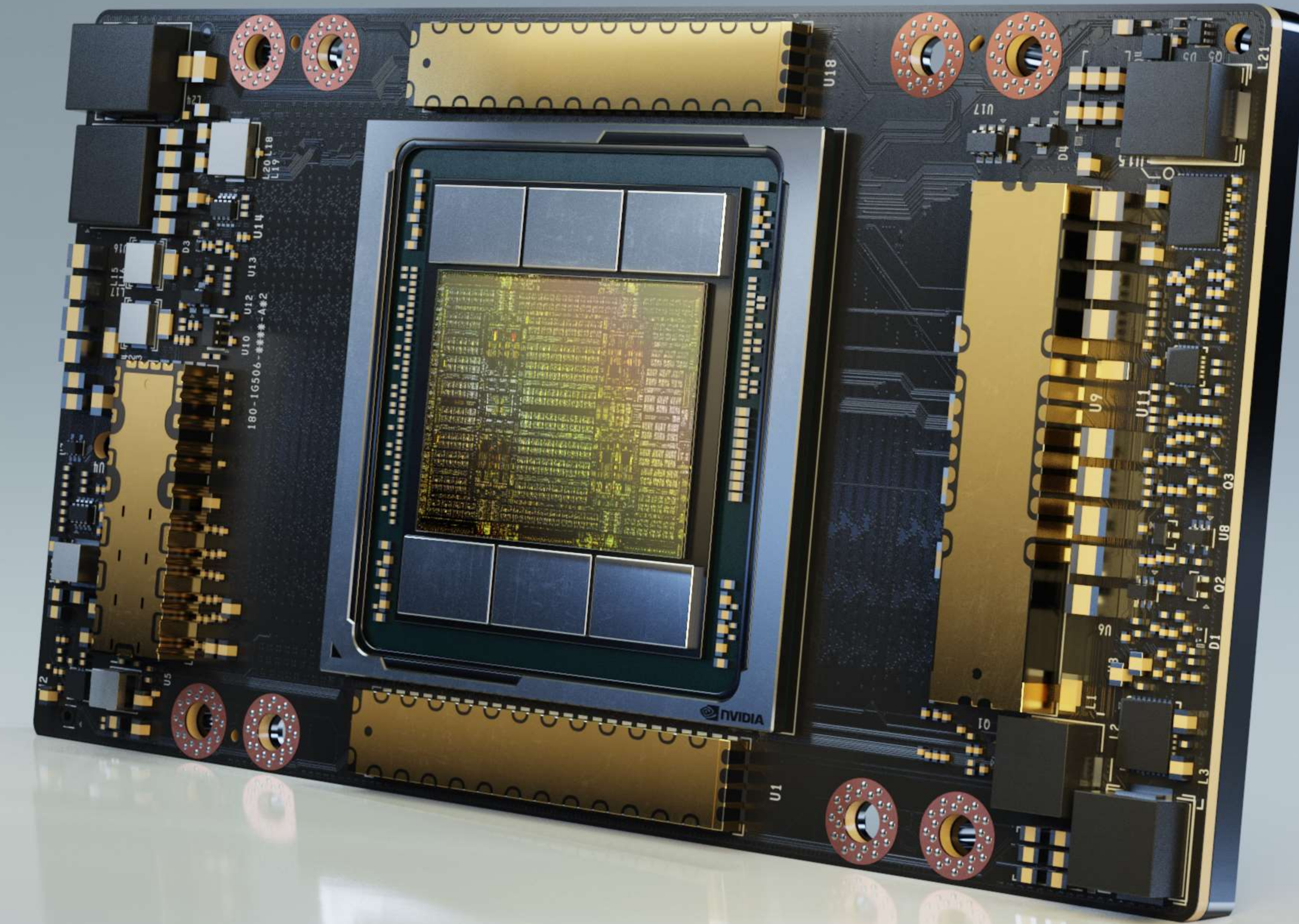
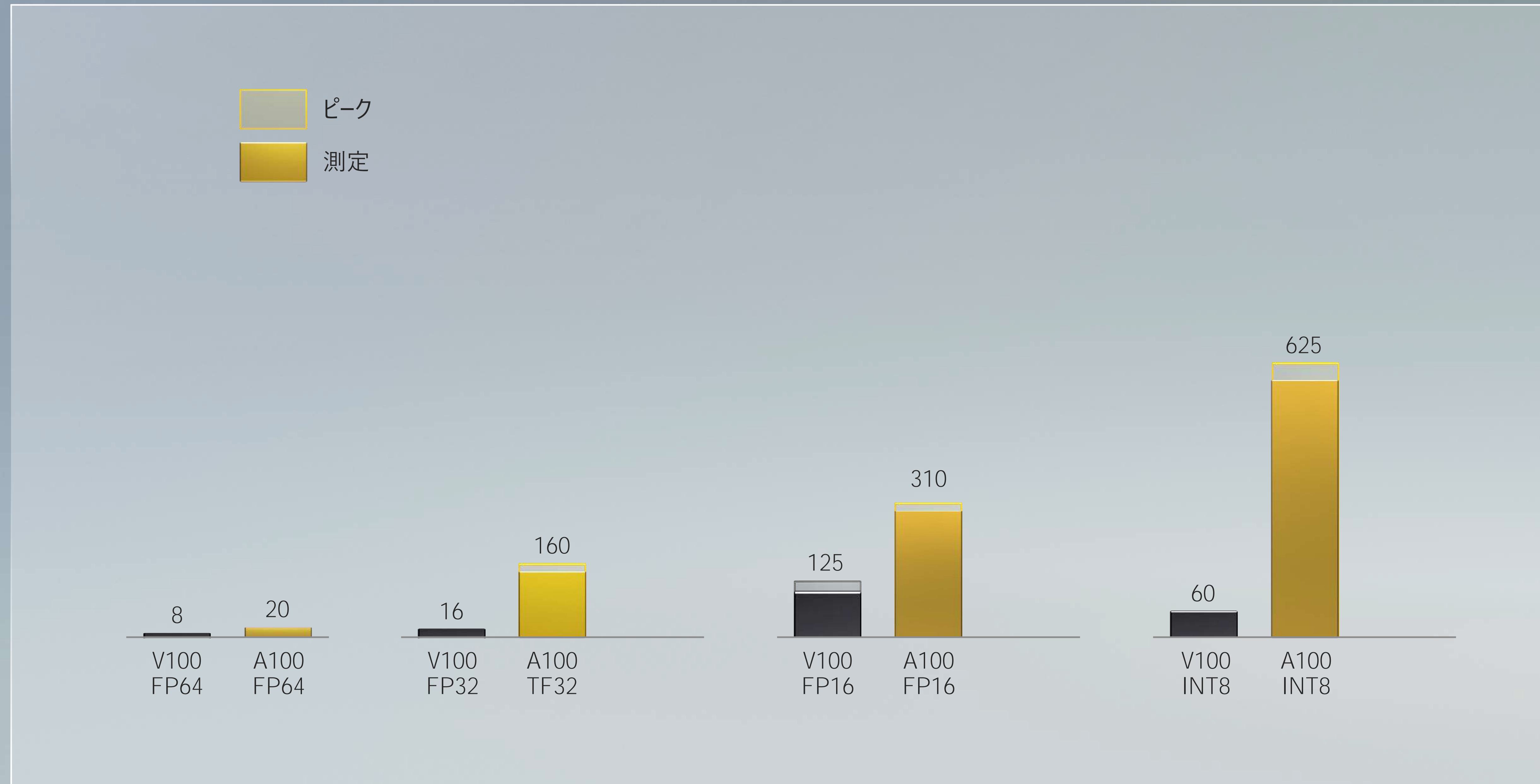
疎な AI Tensor 命令を最適化 | 2 倍高速な実行 | TF32、FP16、BFLOAT16、INT8、INT4 で対応



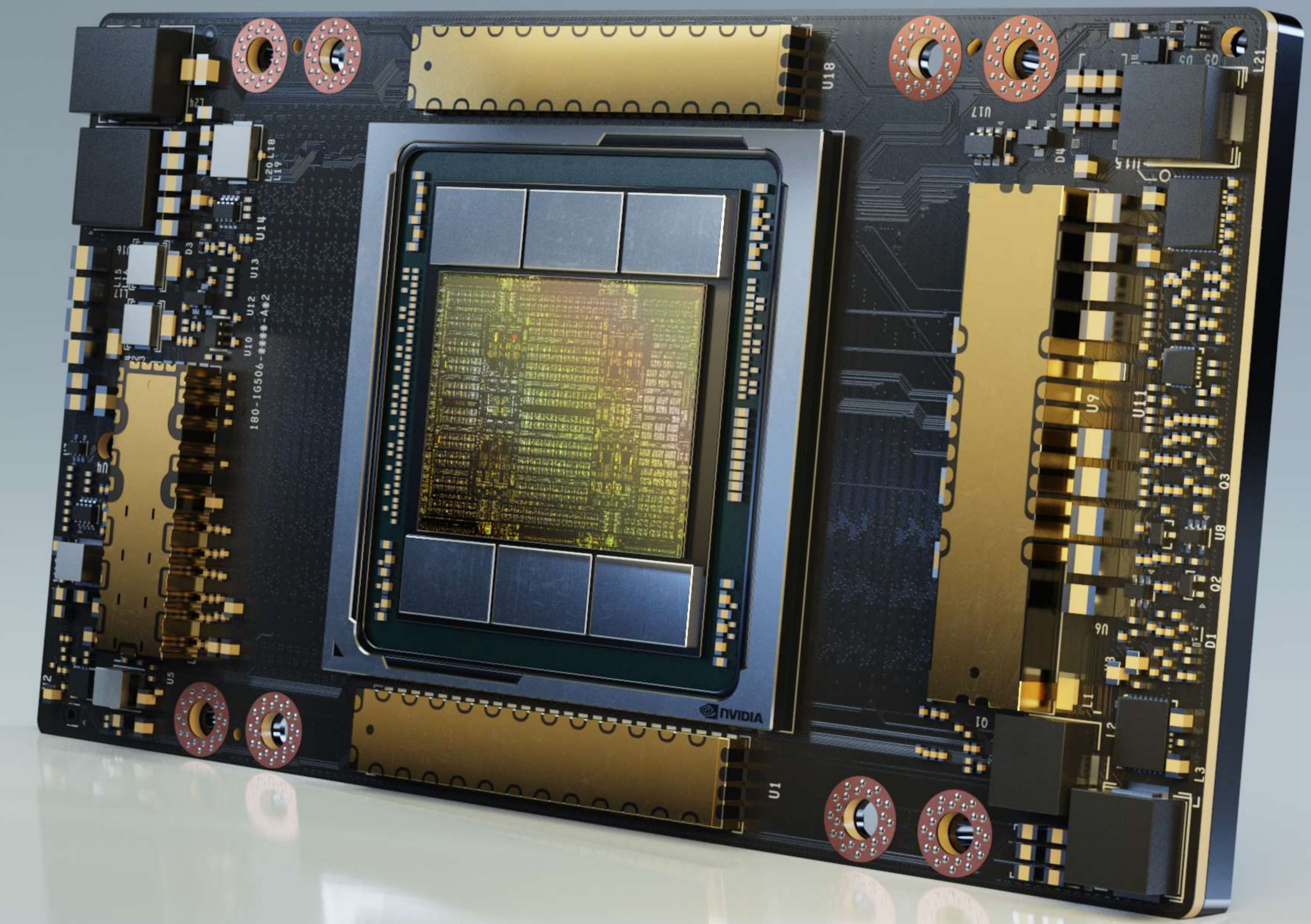
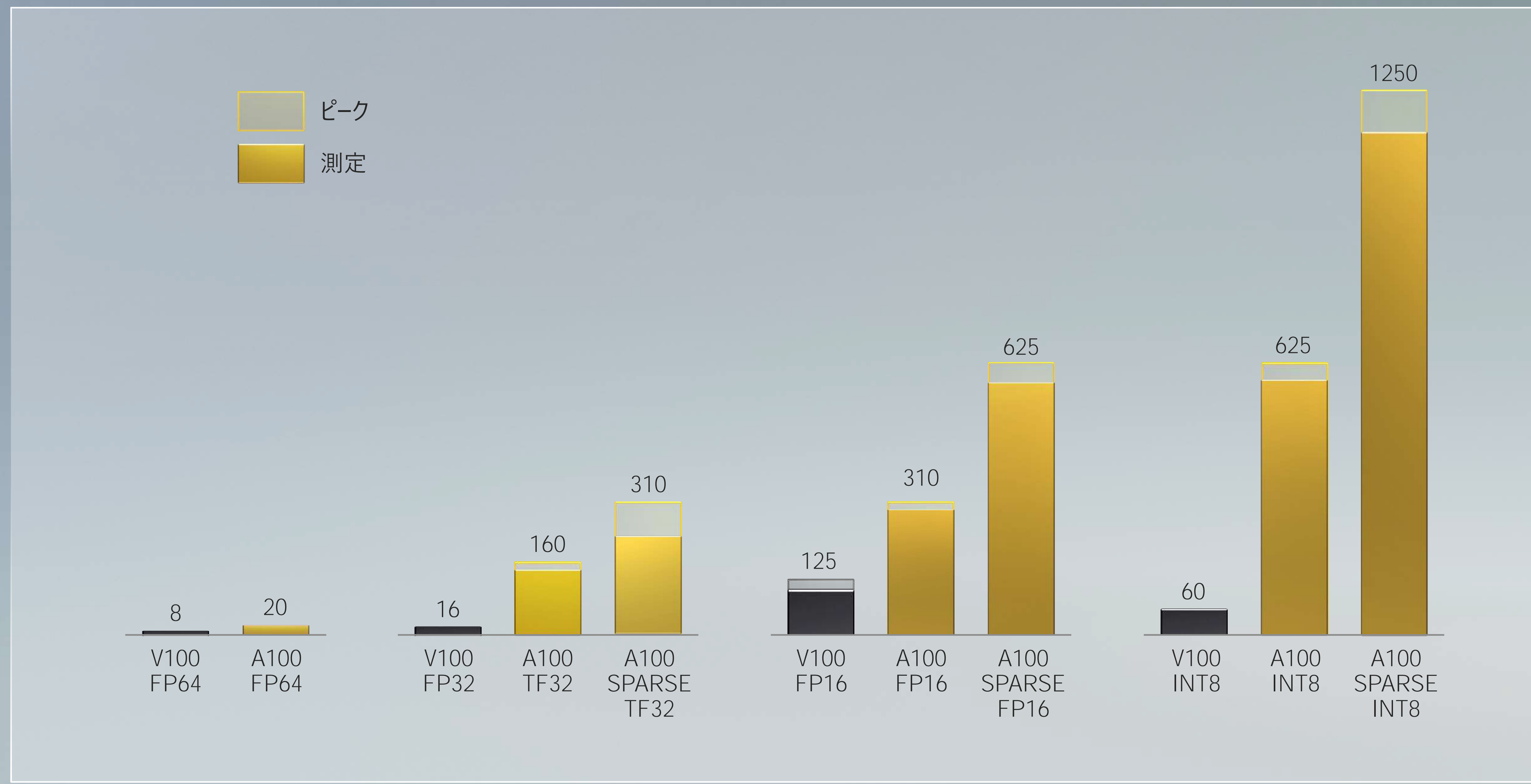
発表: NVIDIA A100 最大の世代間の飛躍



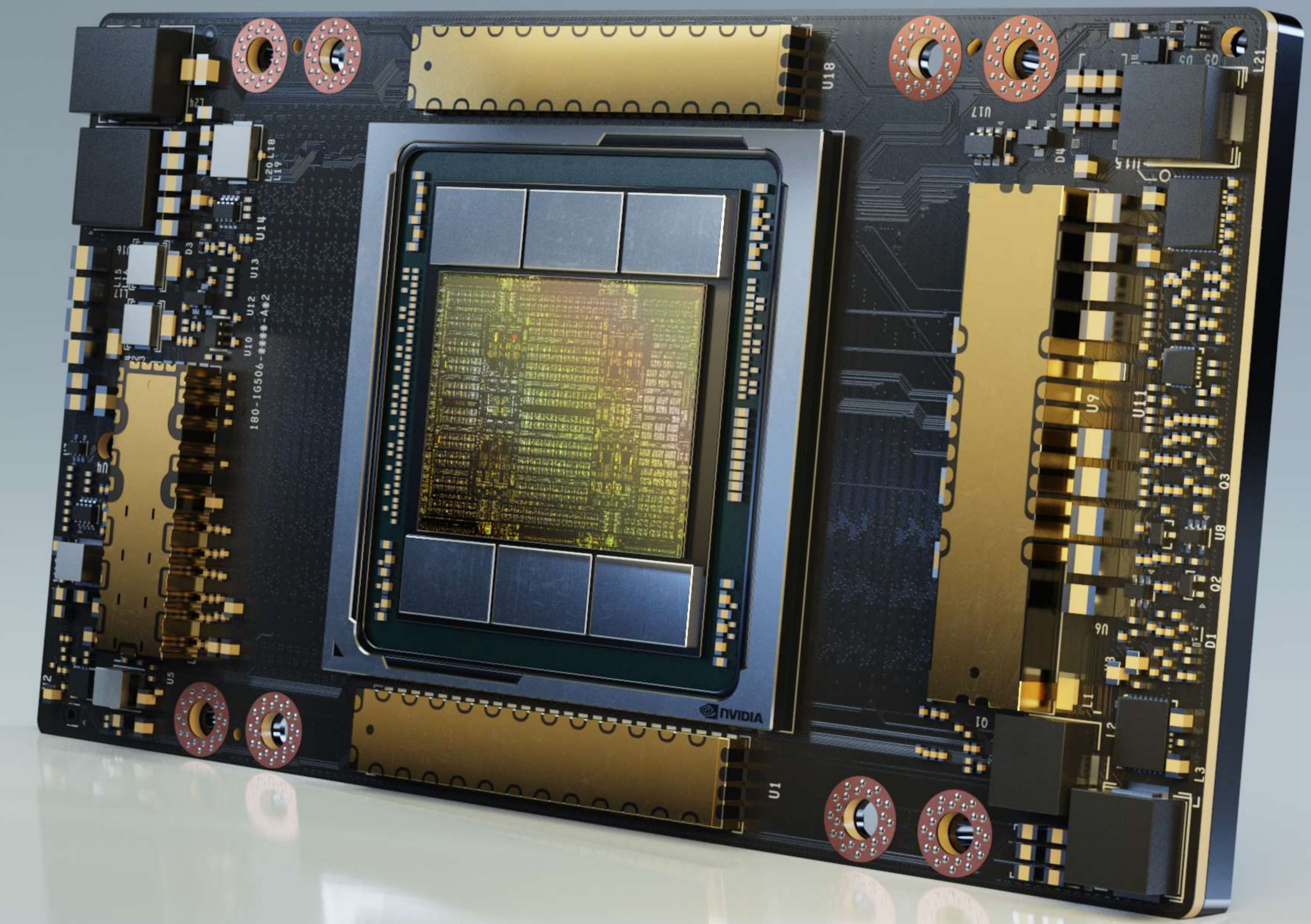
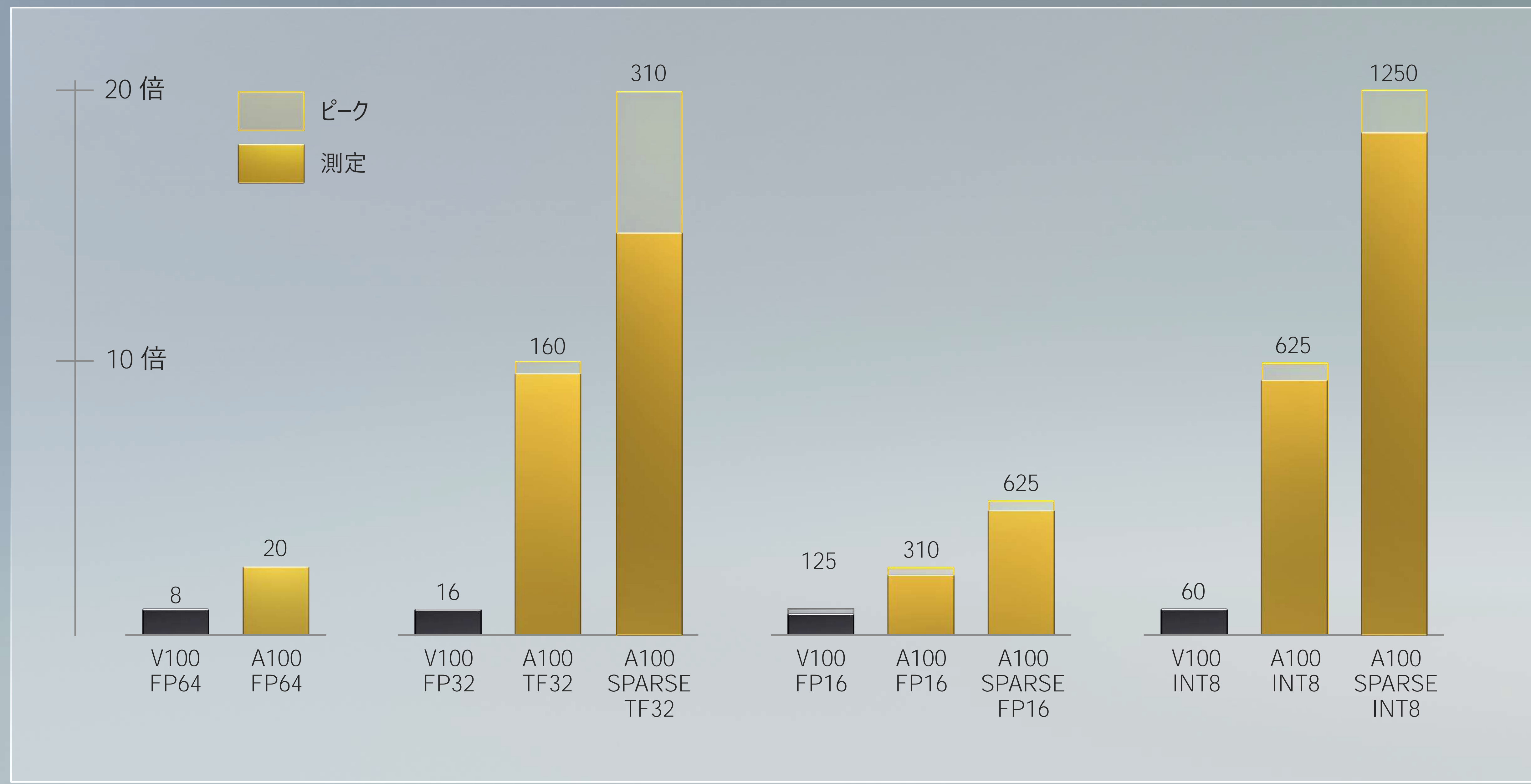
発表: NVIDIA A100 最大の世代間の飛躍



発表: NVIDIA A100 最大の世代間の飛躍

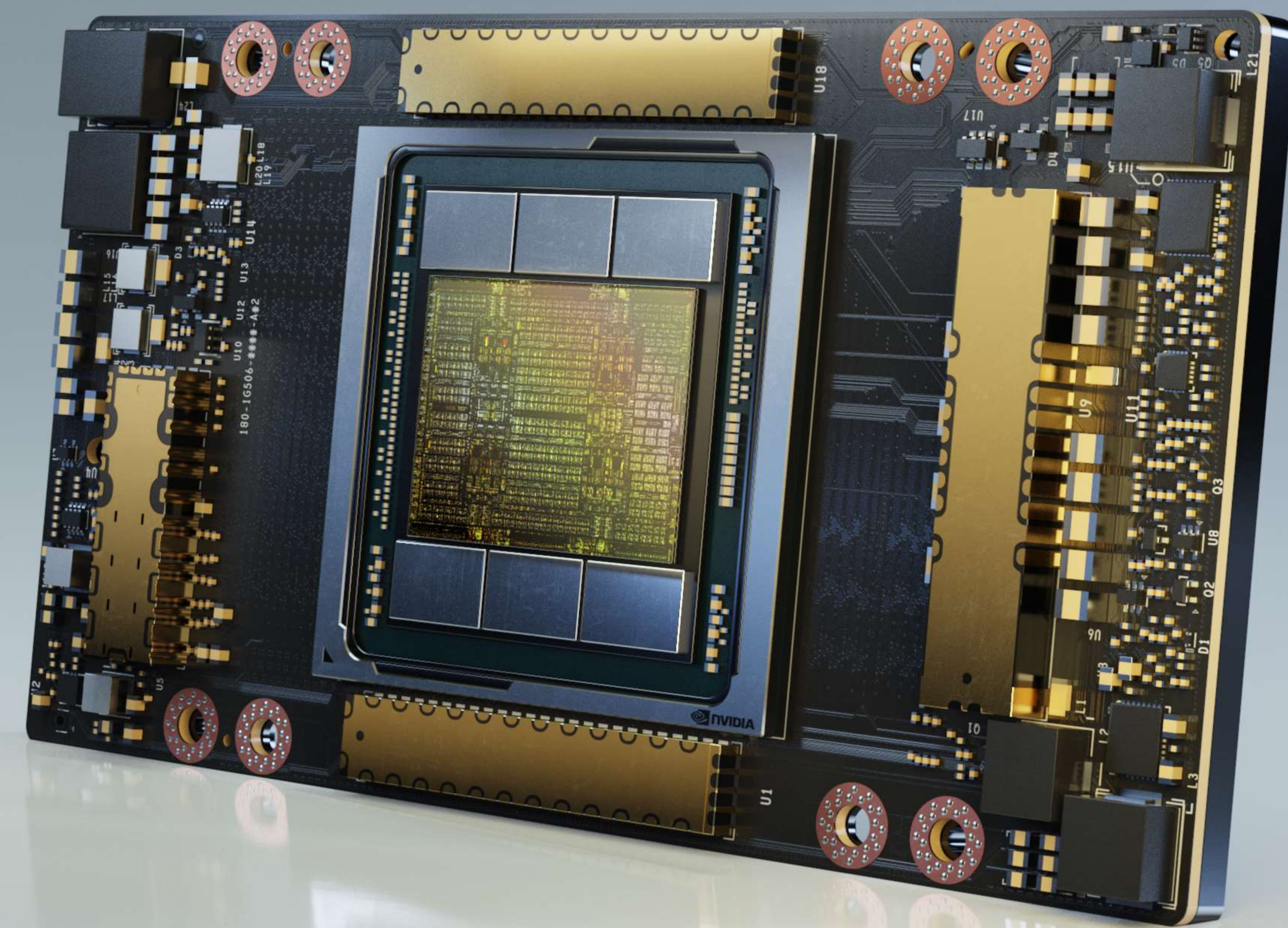
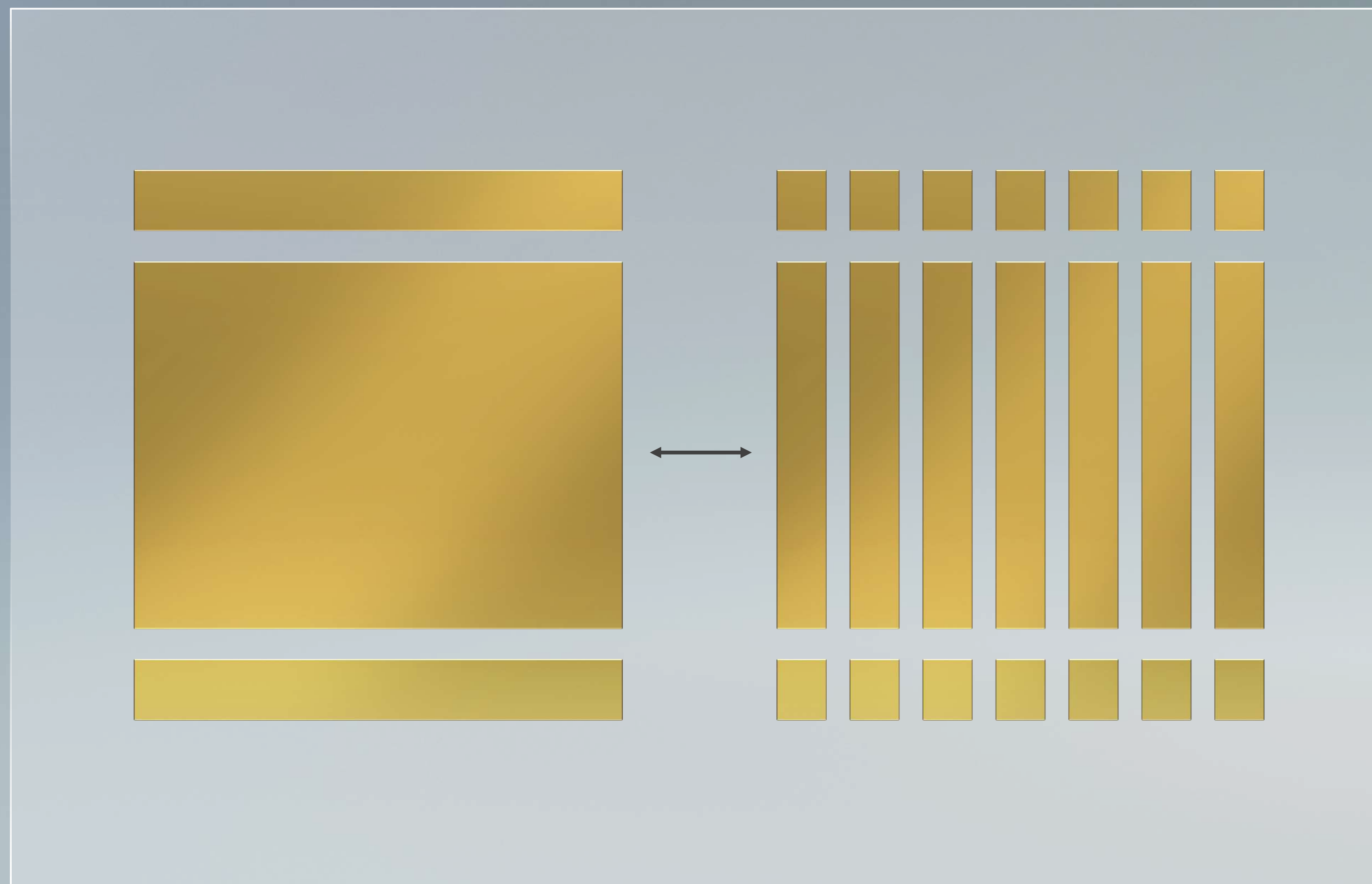


発表: NVIDIA A100 最大の世代間の飛躍

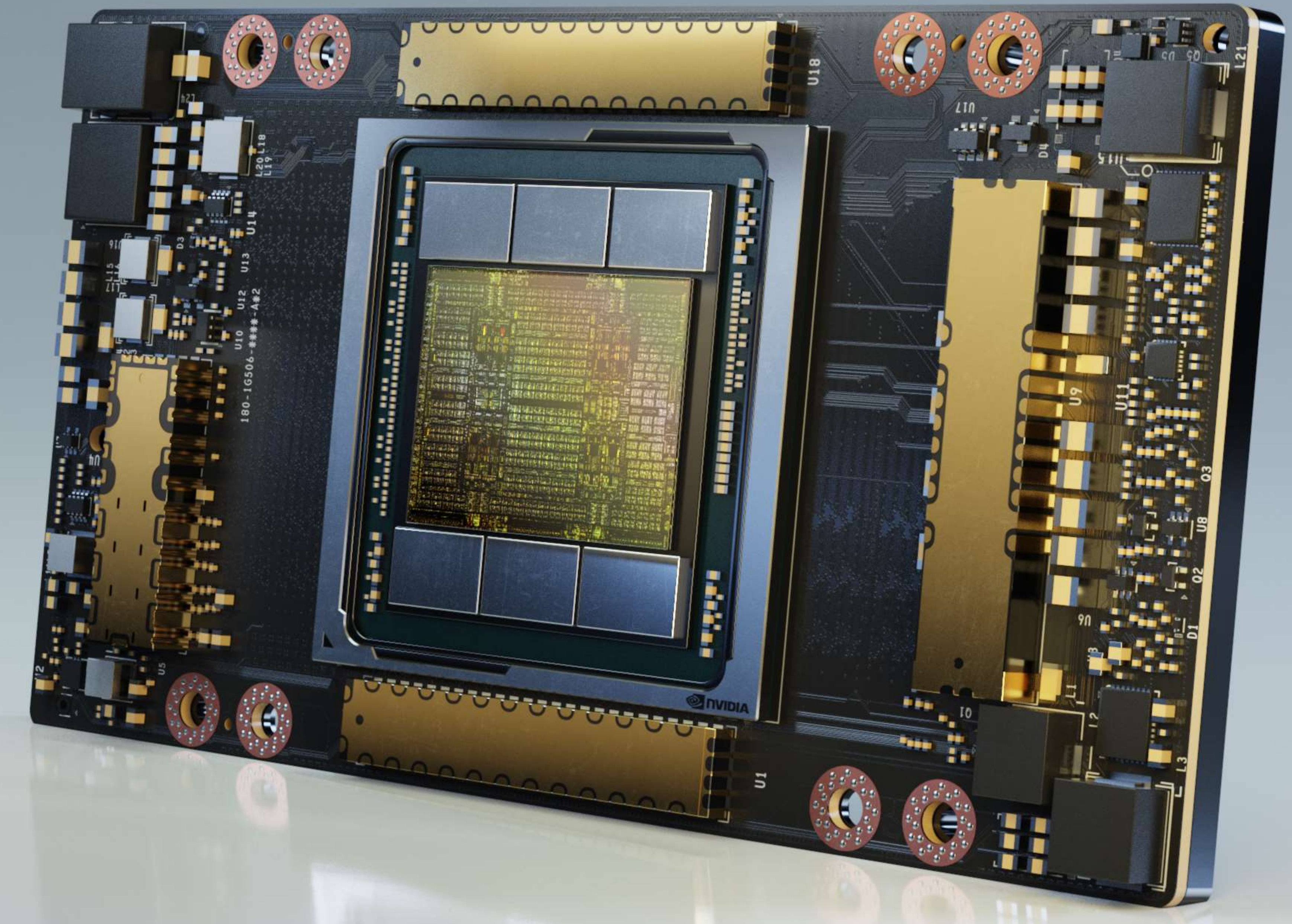


エラスティック GPU コンピューティングのための新しいマルチインスタンス GPU

GPU あたり複数のインスタンスを同時に動かし V100 の 7 倍高いスループット



学習と推論のアクセラレーションを統合



BERT 事前学習のスループット、Pytorch を使用、(2/3)Phase 1 と (1/3)Phase 2 を含む | Phase 1 Seq Len = 128、Phase 2 Seq Len = 512 V100: FP32 精度を使用した 8 基の V100 を搭載した DGX-1 サーバー A100: TF32 精度を使用した 8 基の A100 を搭載した DGX A100 サーバー | BERT Target 推論 | T4、V100: TRT 7.1、精度 = FP16、バッチサイズ = 256 | A100 MIG: リリース前の TRT、バッチサイズ = 94、精度 = INT8、疎性高速化を有効


Batch One Inference Performance on V100

Audio Classification and BERT Question and Answer



What is the native region of the bird I am hearing?

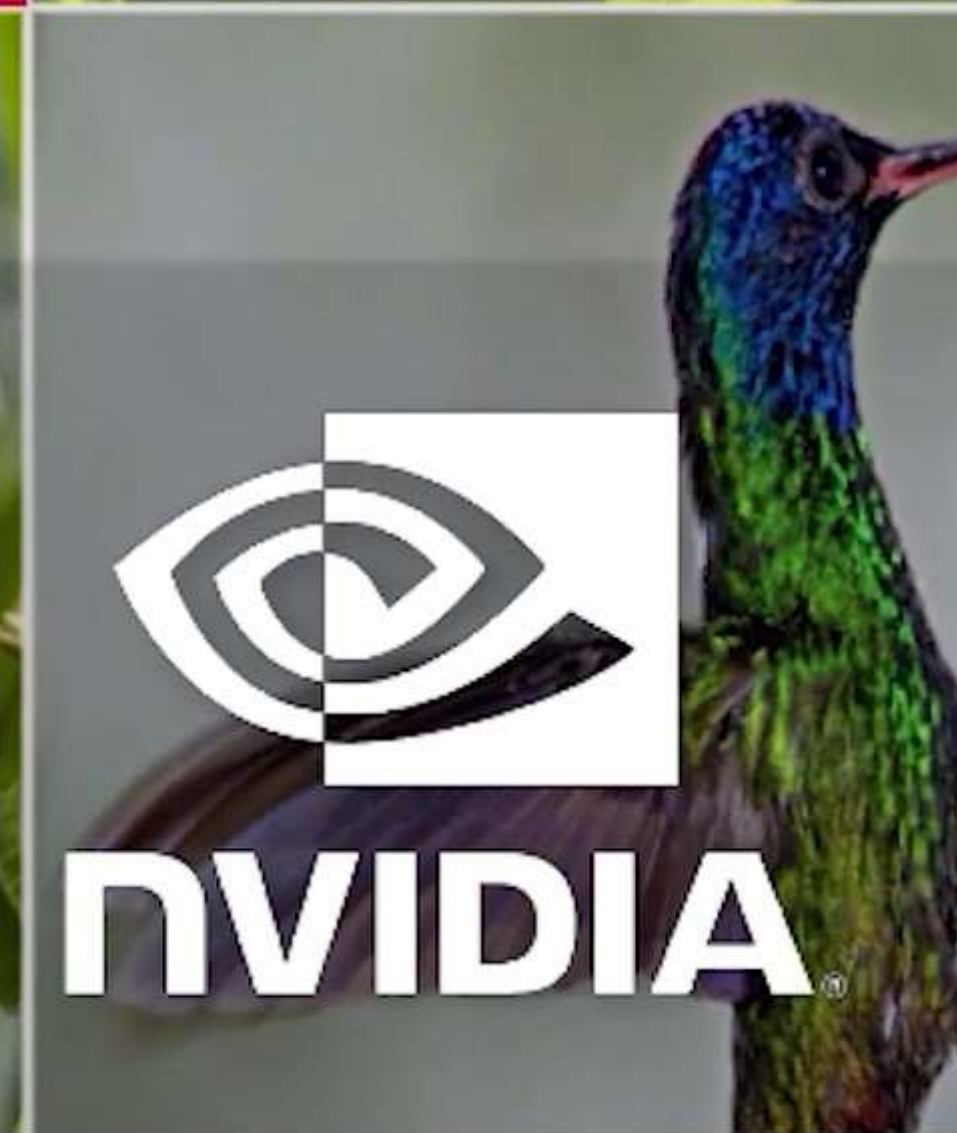
Bronze-winged Parrot



Audio
17.3 ms

NLP
3.4 ms

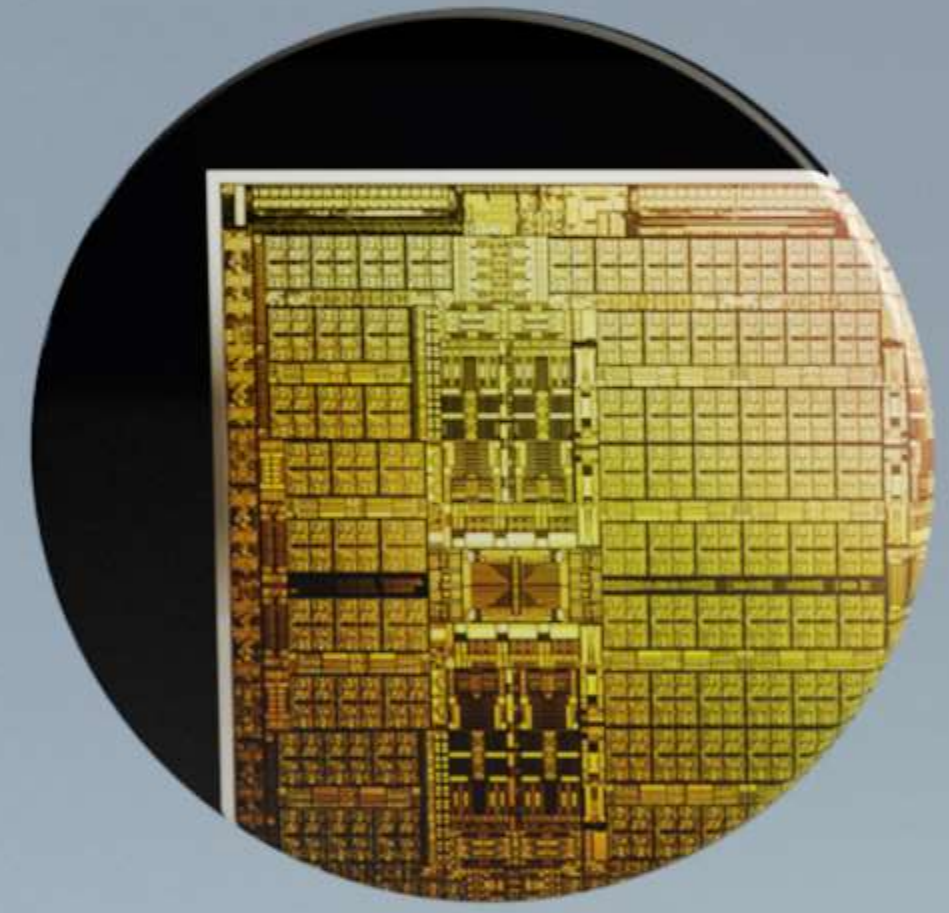
north america



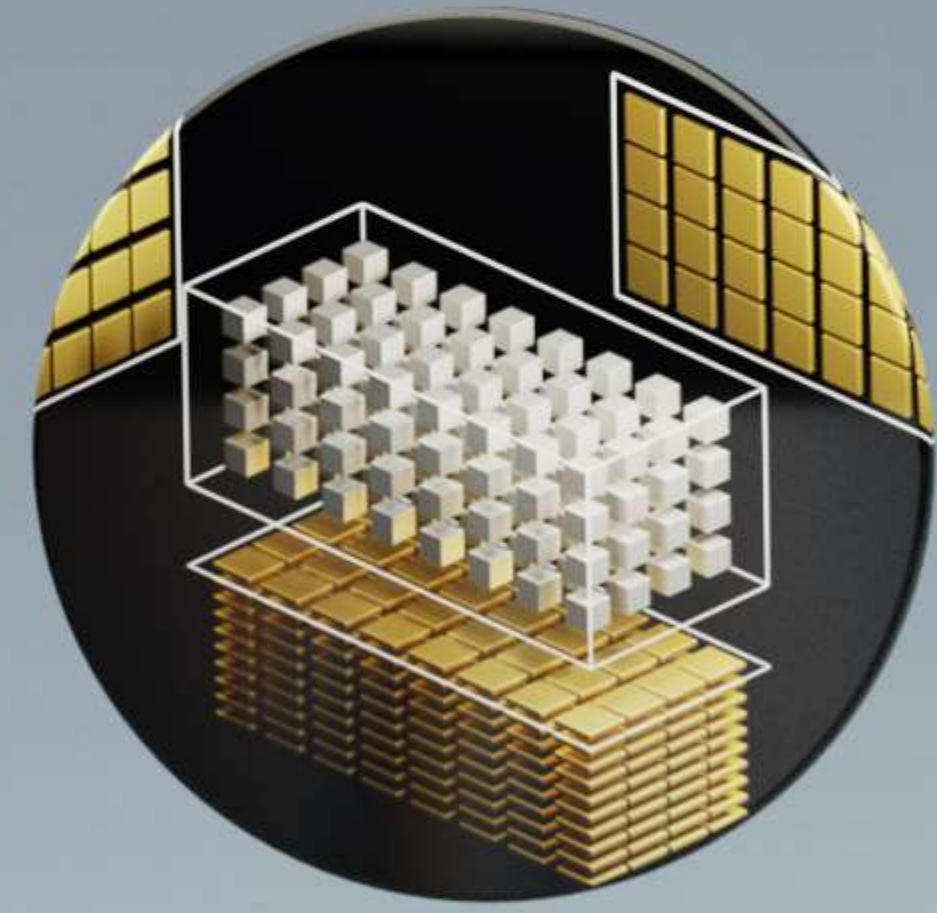
Queries Per Sec: 22.5



発表: NVIDIA A100



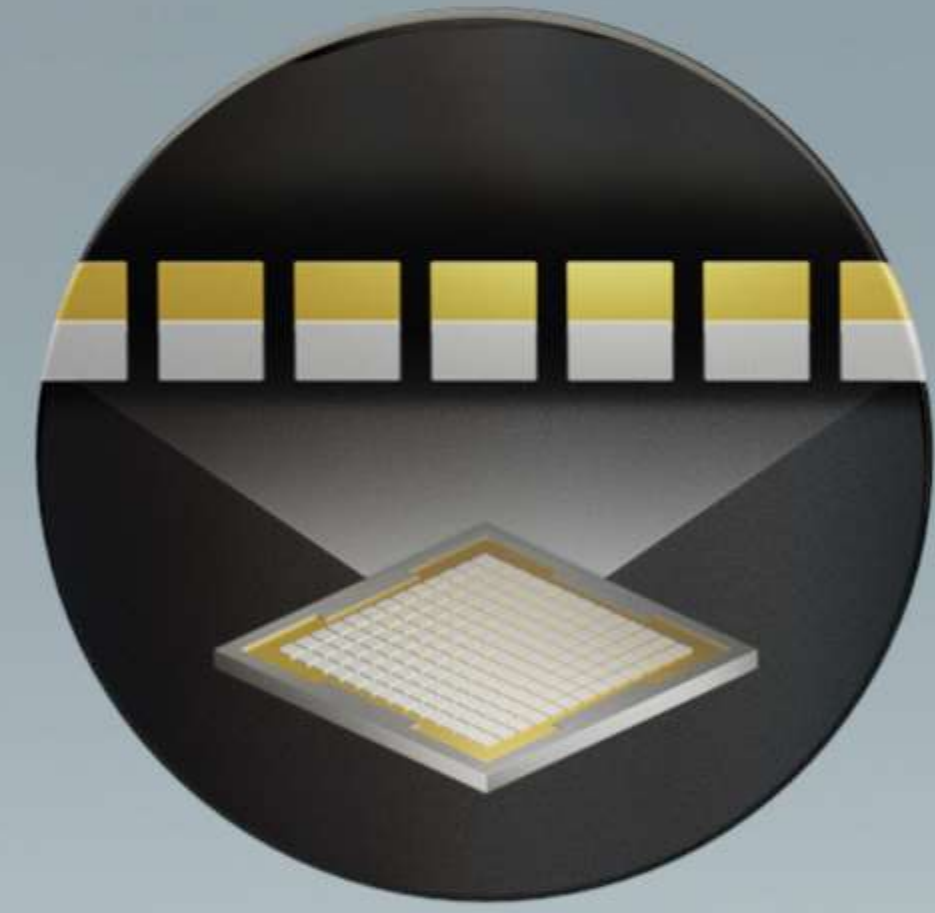
540 億トランジスタ



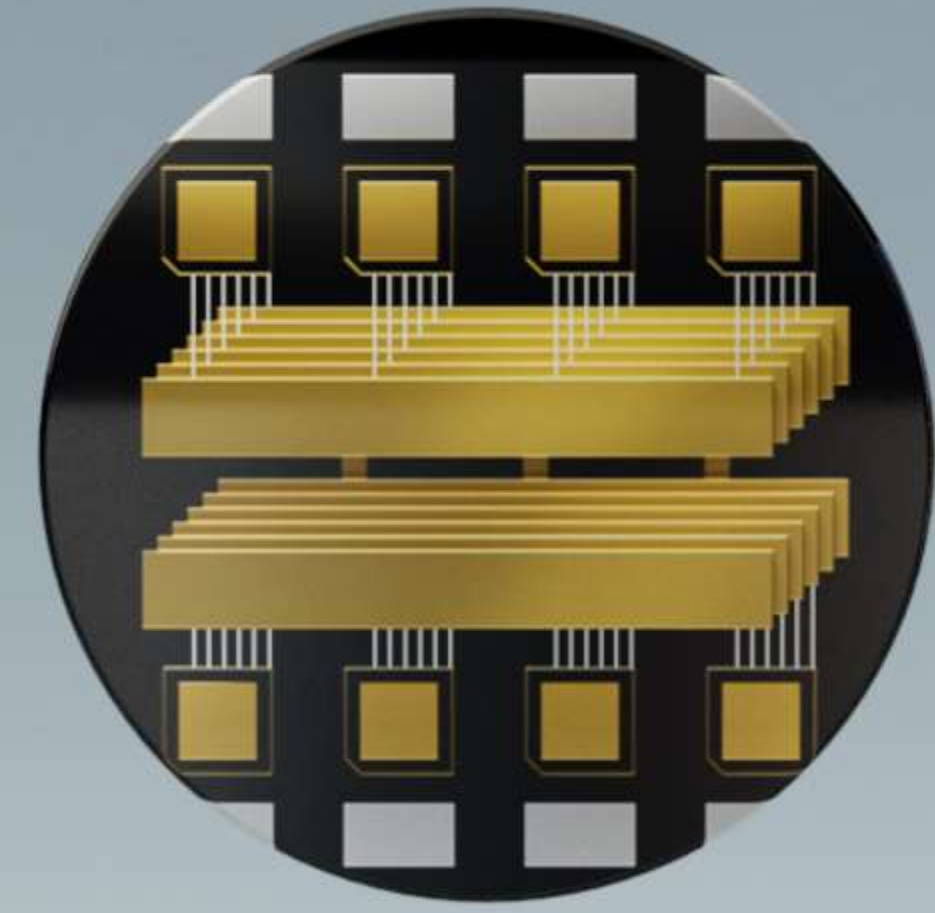
第 3 世代 TENSOR コア



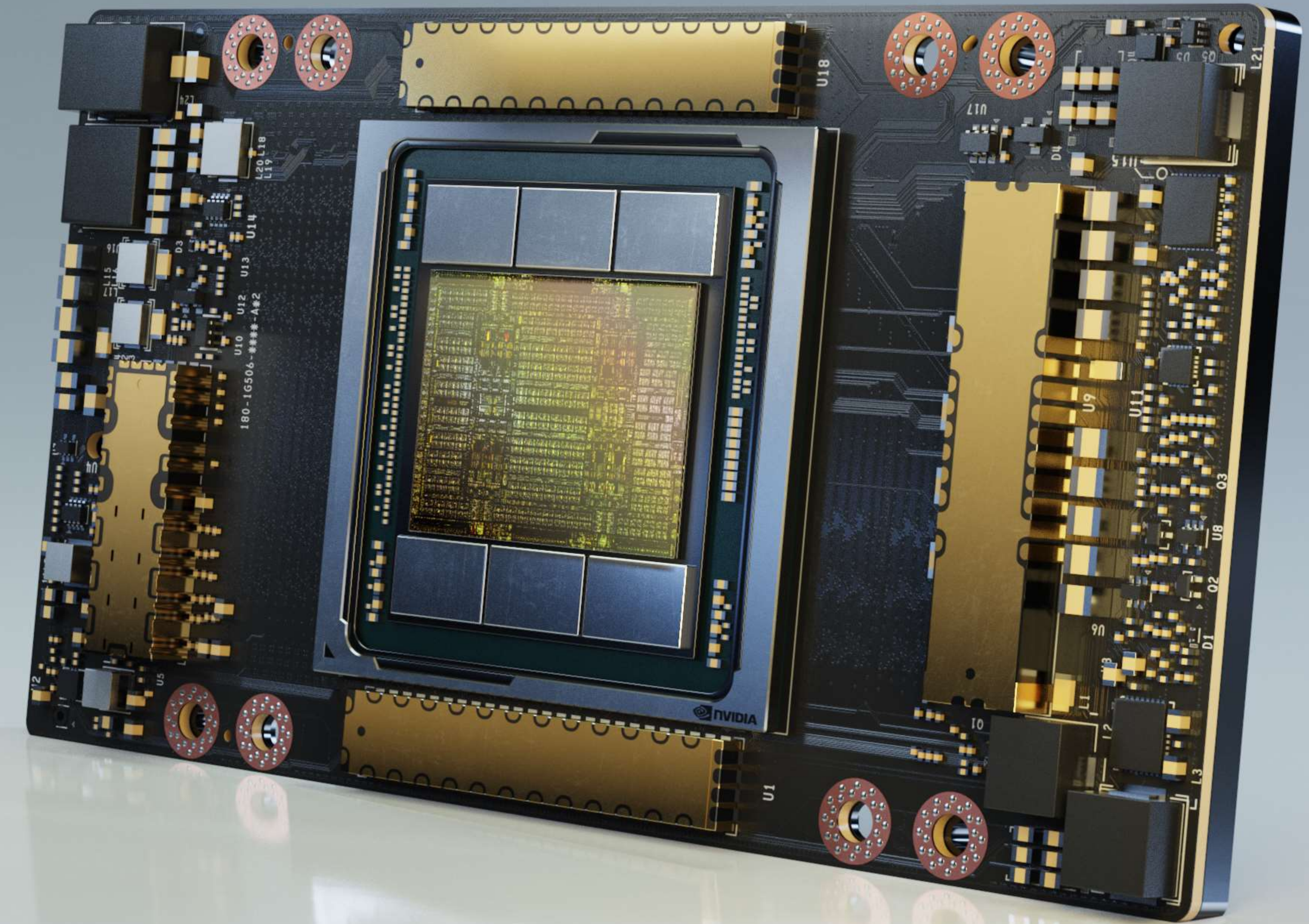
スパースティ アクセラレーション



MIG



第 3 世代 NVLINK & NVSWITCH



発表:
NVIDIA DGX A100
第3世代 統合 AI システム

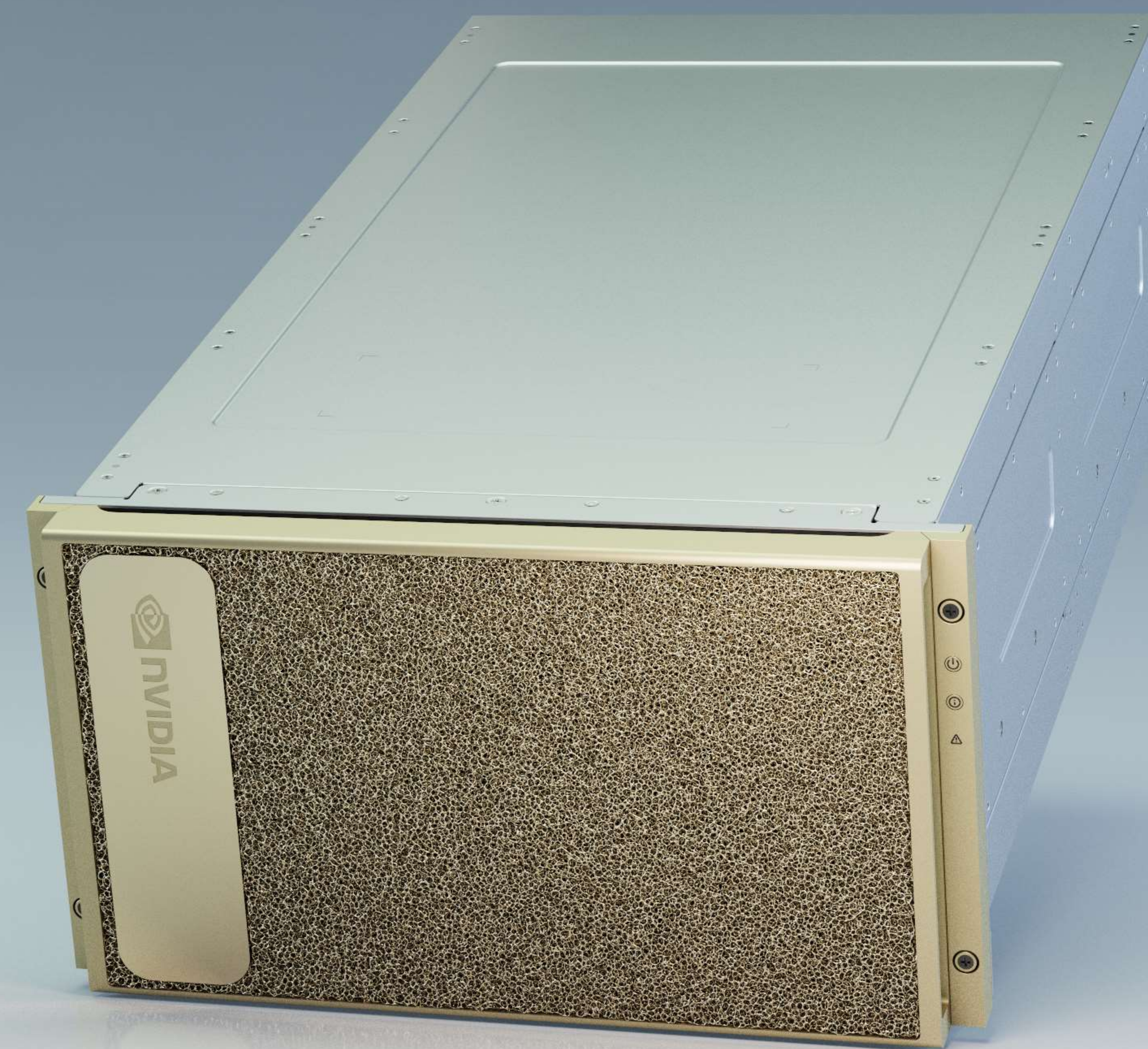
シングルノードで 5 PetaFLOPS の性能

エンドツーエンドのデータサイエンスと AI のための統合システム

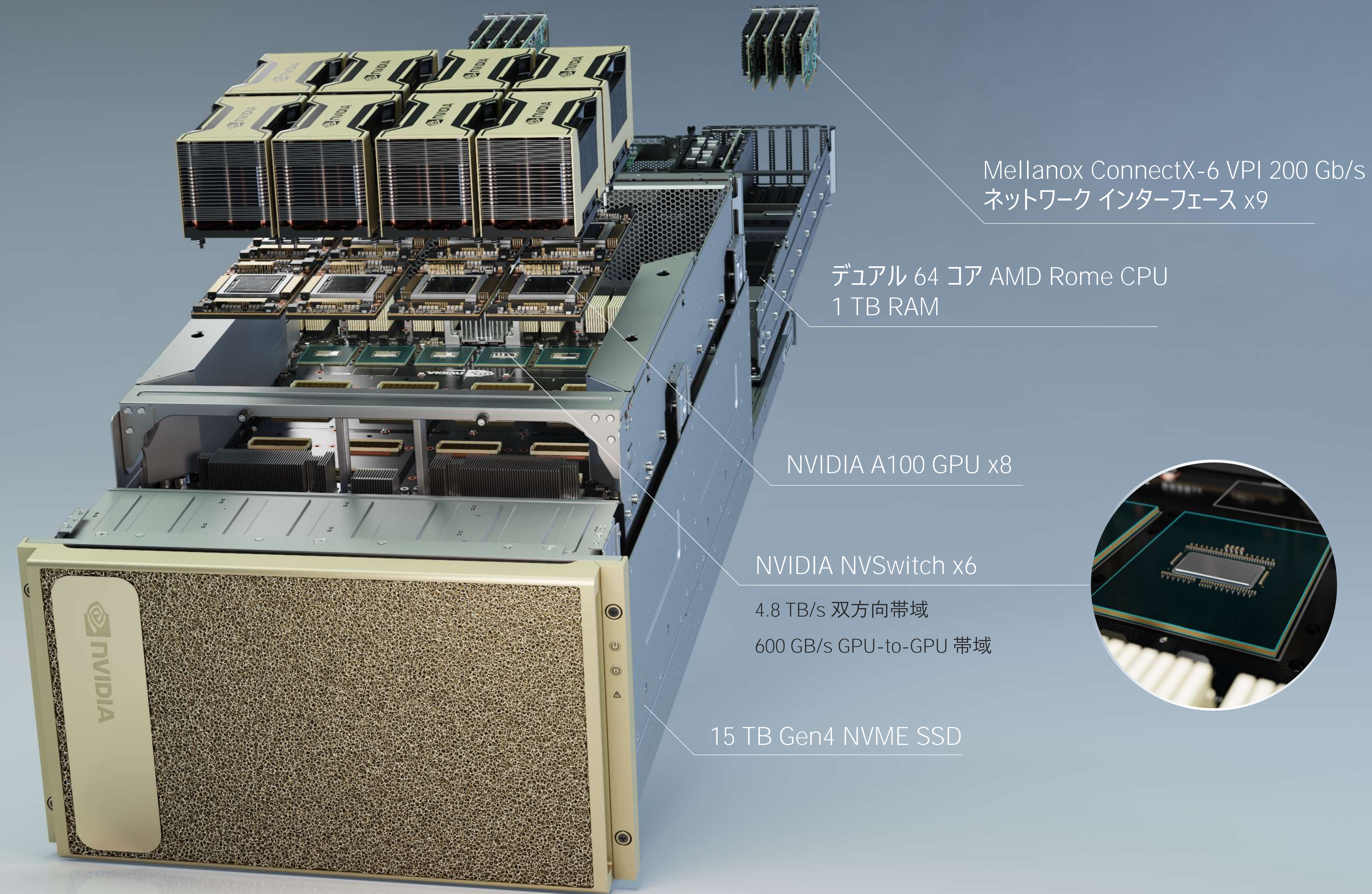
完全にアクセラレーションされたスタック — Spark 3.0、RAPIDS、TensorFlow、PyTorch、Triton

エラスティックなスケールアップまたはスケールアウト コンピューティング

Mellanox ネットワーキングによる高いスケーラビリティ



発表:
NVIDIA DGX A100
第3世代統合AIシステム
シングルノードで5 PetaFLOPS の性能



発表:
NVIDIA DGX A100
第3世代統合AIシステム
シングルノードで5 PetaFLOPS の性能

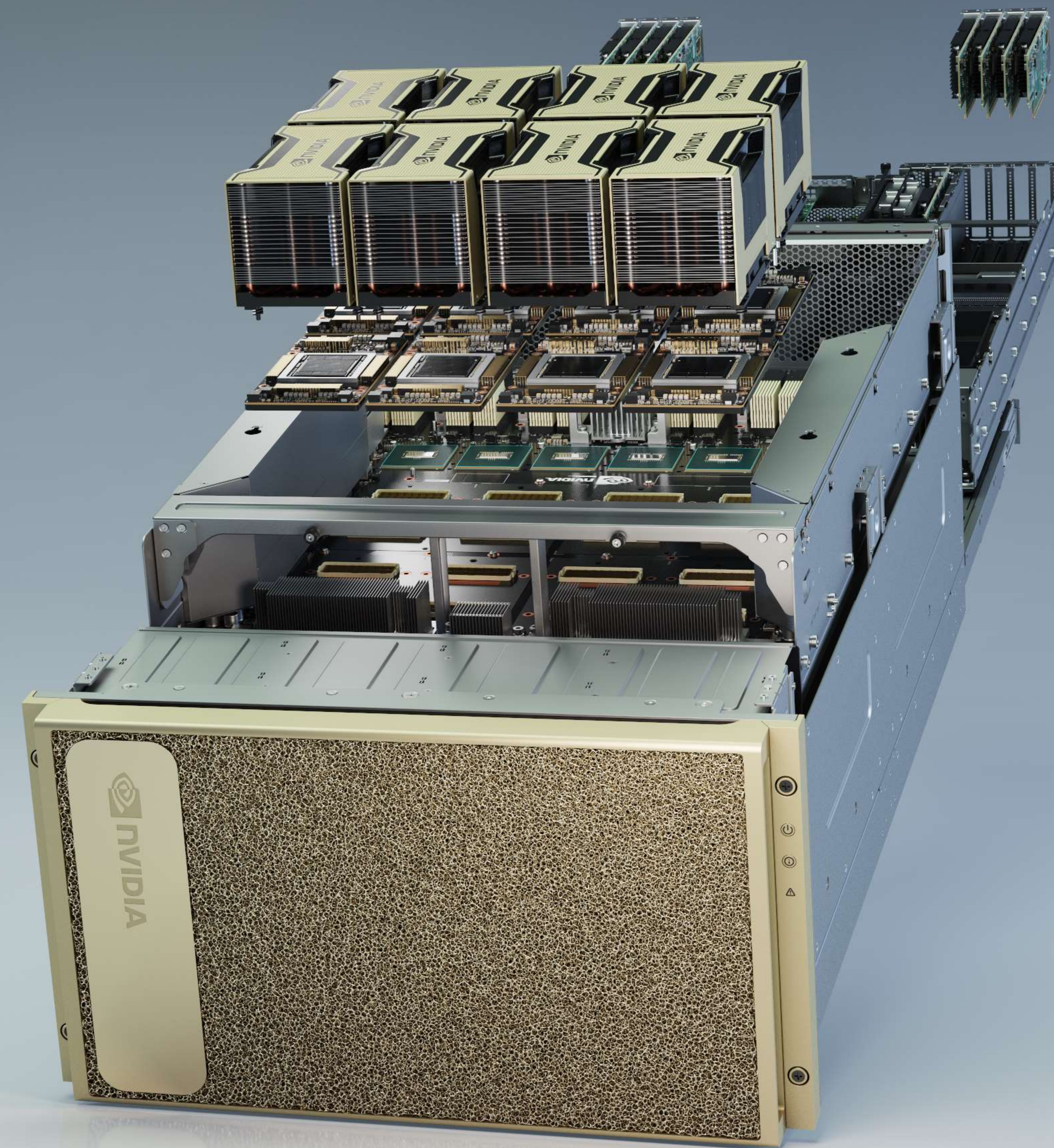
INT8 10 PetaOPS ピーク

FP16 5 PFLOPS ピーク

TF32 2.5 PFLOPS ピーク

FP64 156 TFLOPS ピーク

Tensor コア + スパーシティ アクセラレーション



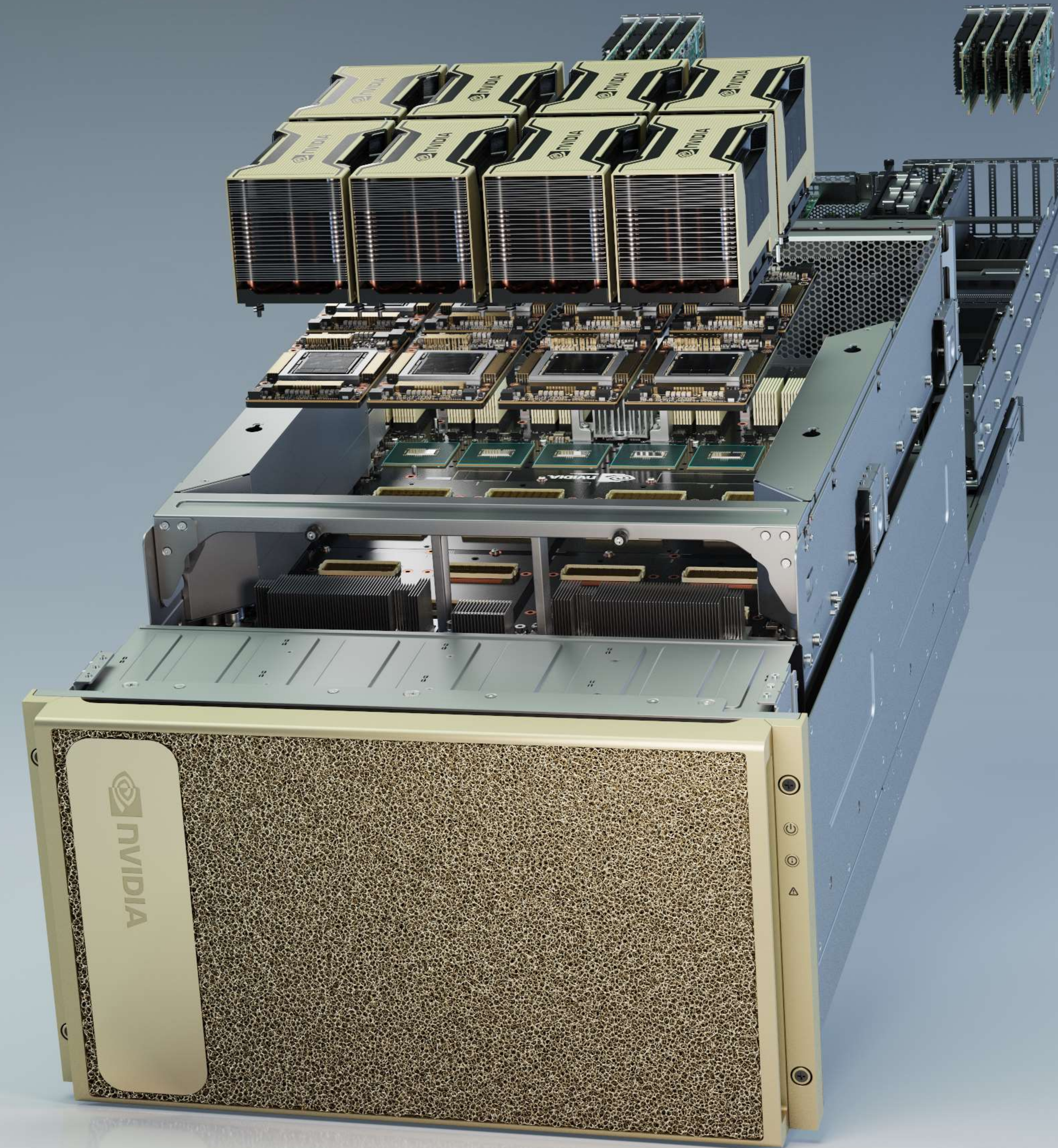
発表:
NVIDIA DGX A100
第3世代統合AIシステム
シングルノードで5 PetaFLOPS の性能

150 倍 AI 演算

40 倍 メモリ帯域幅

40 倍 IO 帯域幅

ハイエンド CPU サーバーとの比較



発表:
NVIDIA DGX A100
第3世代統合AIシステム
シングルノードで5 PetaFLOPS の性能

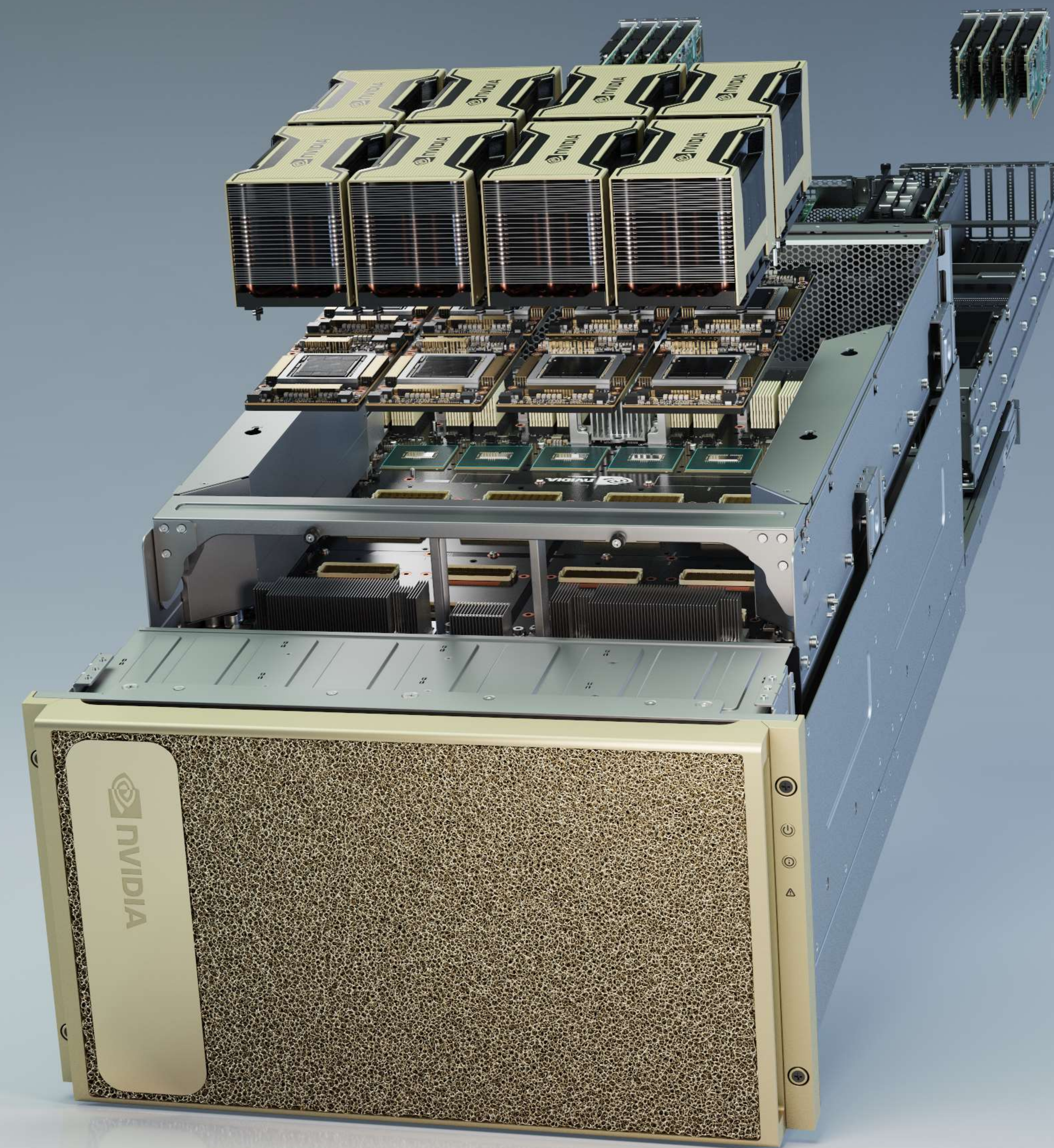
150 倍 AI 演算

40 倍 メモリ帯域幅

40 倍 IO 帯域幅

ハイエンド CPU サーバーとの比較

199,000 ドルで販売中

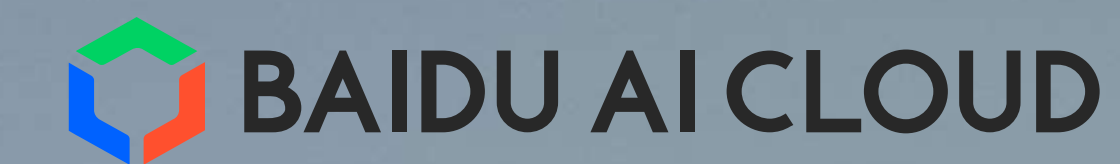


発表:

NVIDIA A100 ライトハウス カスタマー

業界のリーダーが選ぶ エラスティック データ センター アクセラレータ

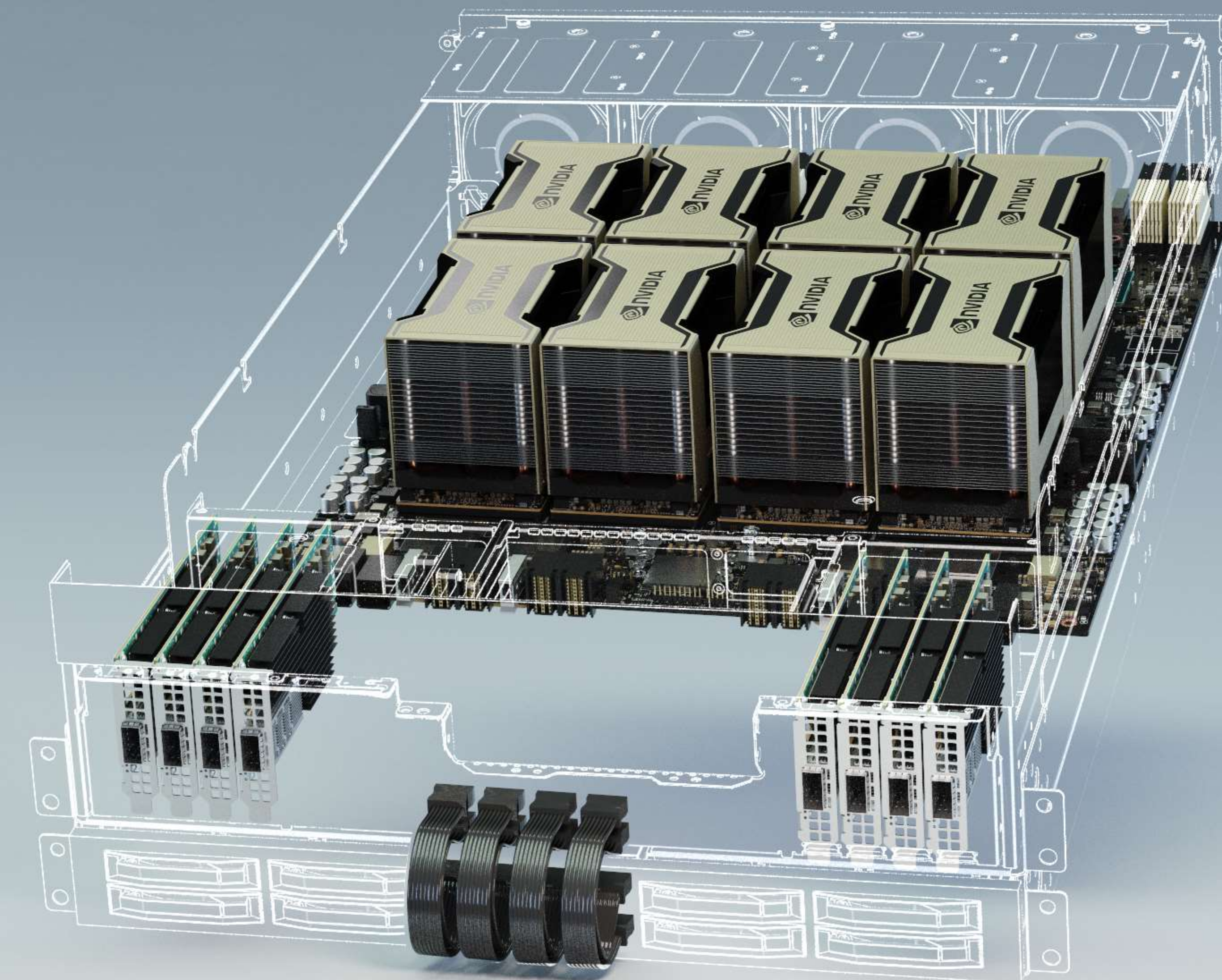
クラウド



Google Cloud

Tencent Cloud

システム



今日の AI データ センター

50 台の DGX-1 システム : AI 学習向け

600 台の CPU システム : AI 推論向け

1,100 万ドル

25 ラック

630 kW



\$11M 630 kW

DGX A100 AI

5 台の DGX A100 システム : AI 学習と推論向け

100 万ドル

1 ラック

28 kW

\$1M 28 kW

1/10

費用

1/20

消費電力



今日の AI データ センター

50 台の DGX-1 システム : AI 学習向け

600 台の CPU システム : AI 推論向け

1,100 万ドル

25 ラック

630 kW



DGX A100 AI

5 台の DGX A100 システム : AI 学習と推論向け

100 万ドル

1 ラック

28 kW



ページランクの事例

Common Crawl データセット
2.6TB グラフ - 1280 億エッジ

3,000 CPU サーバー - 105 ラック

520 億エッジ / 秒

A photograph of a server room with multiple rows of black server racks. The racks are filled with server hardware, and the room is lit by overhead fluorescent lights. The floor is a dark, polished surface. The text '520 億エッジ / 秒' is overlaid in white on the right side of the image.

ページランクの事例

Common Crawl データセット
2.6TB グラフ - 1280 億エッジ

外部 NVLINK で接続された 4 台の DGX A100

6880 億エッジ / 秒

13 倍

性能

1/75

費用



ページランクの事例

Common Crawl データセット
2.6TB グラフ - 1280 億エッジ

3,000 CPU サーバー - 105 ラック



ページランクの事例

Common Crawl データ セット
2.6TB グラフ - 1280 億エッジ

外部 NVLINK で接続された 4 台の DGX A100





発表:
NVIDIA DGX A100 SUPERPOD

140 台の DGX A100 システム (1,120 基の A100)
170 台の Mellanox Quantum 200G InfiniBand スイッチ
280 Tb/s ネットワーク ファブリック - 15km の光ケーブル
4 PB のオールフラッシュ ネットワーク ストレージ
700 PFLOPS の AI 性能
3 週間以内で構築



NVIDIA が SATURNV を拡張

拡張前

1,800 台の DGX システム

1.8 ExaFLOPS

4 セットの DGX SuperPOD を追加

560 台の DGX A100 = 2.8 ExaFLOPS

4.6 ExaFLOPS 総演算性能

発表:
NVIDIA DGX A100
第3世代統合AIシステム
シングルノードで5 PetaFLOPS の性能

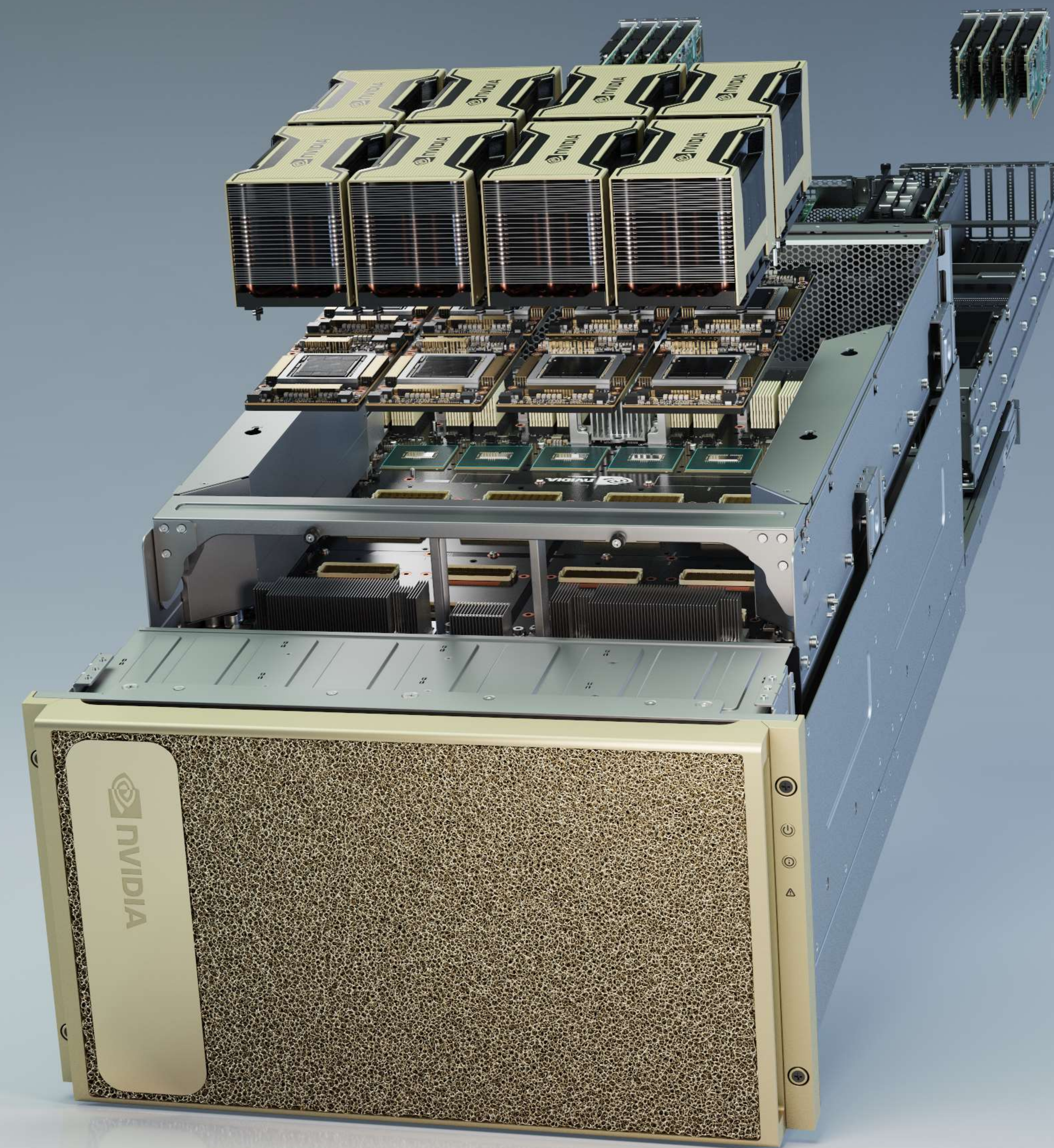
150 倍 AI 演算

40 倍 メモリ帯域幅

40 倍 IO 帯域幅

ハイエンド CPU サーバーとの比較

199,000 ドルで販売中



SMART EVERYTHING 革命

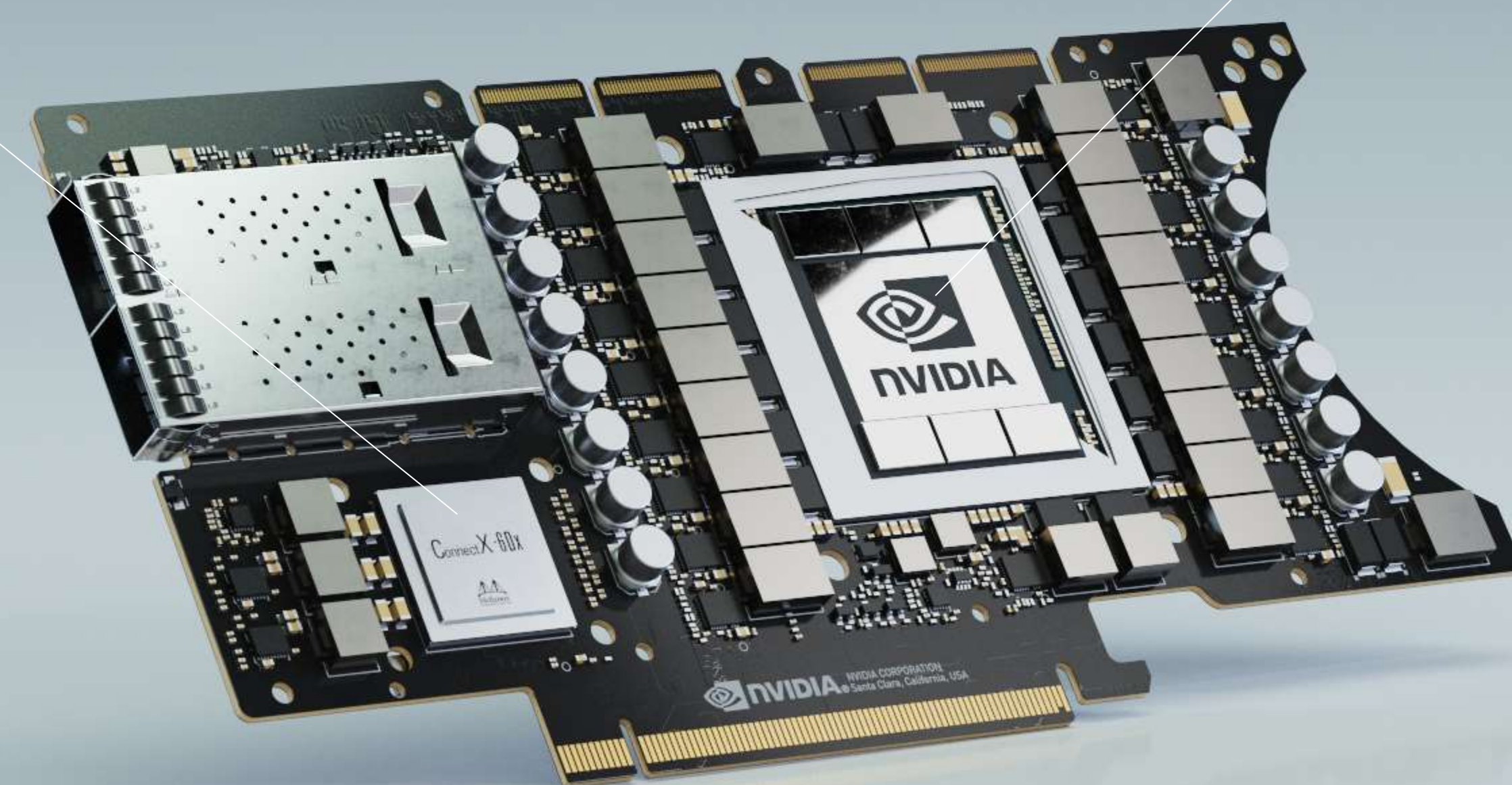


常時オン | 即座に感知-推論-行動 | 遠隔 | 数兆個

発表: MELLANOX CX6 DX 搭載 NVIDIA EGX A100

NVIDIA Mellanox ConnectX-6 DX

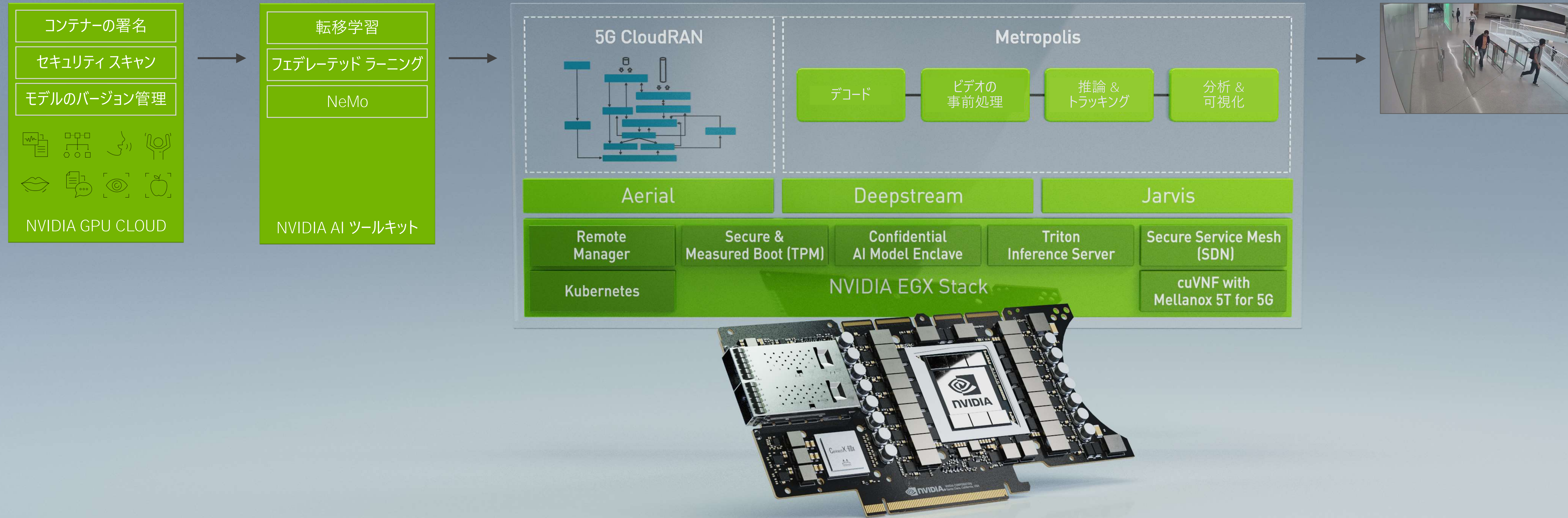
Dual 100 Gb/s Ethernet または InfiniBand
ラインスピード TLS/IPSec 暗号化エンジン
通信事業者向けタイムトリガ方式転送技術:
Time Triggered Transmission Tech for Telco
(5T for 5G)
ASAP² SR-IOV と VirtIO オフロード



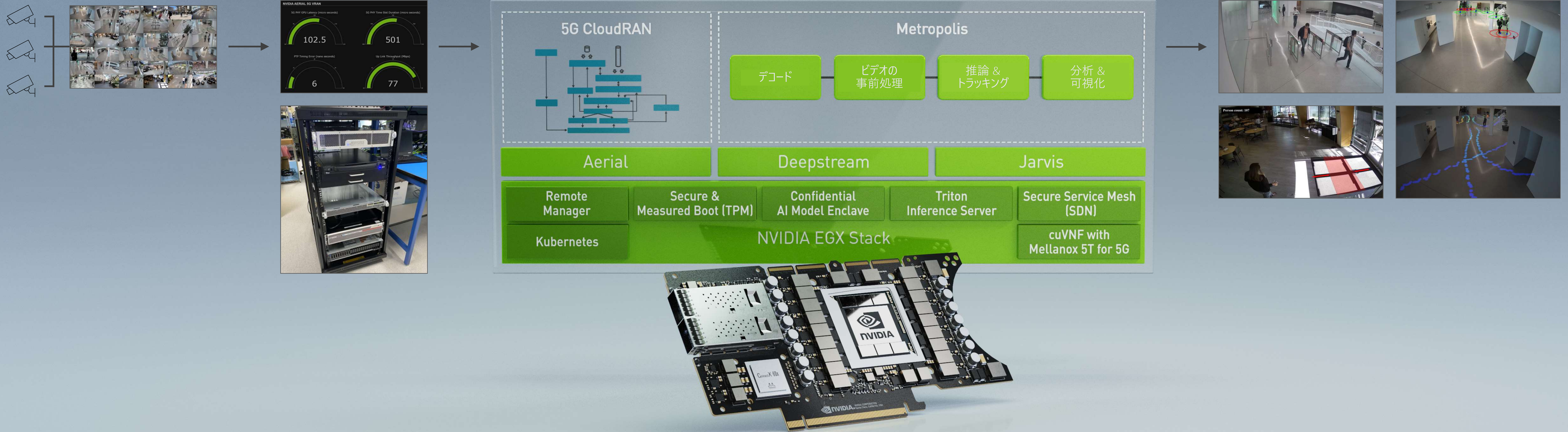
NVIDIA Ampere GPU

第3世代 Tensor コア
対話型 AI 向けの新しいセキュリティエンジン
セキュアで認証済みのブート

発表: NVIDIA EGX エッジ AI プラットフォーム



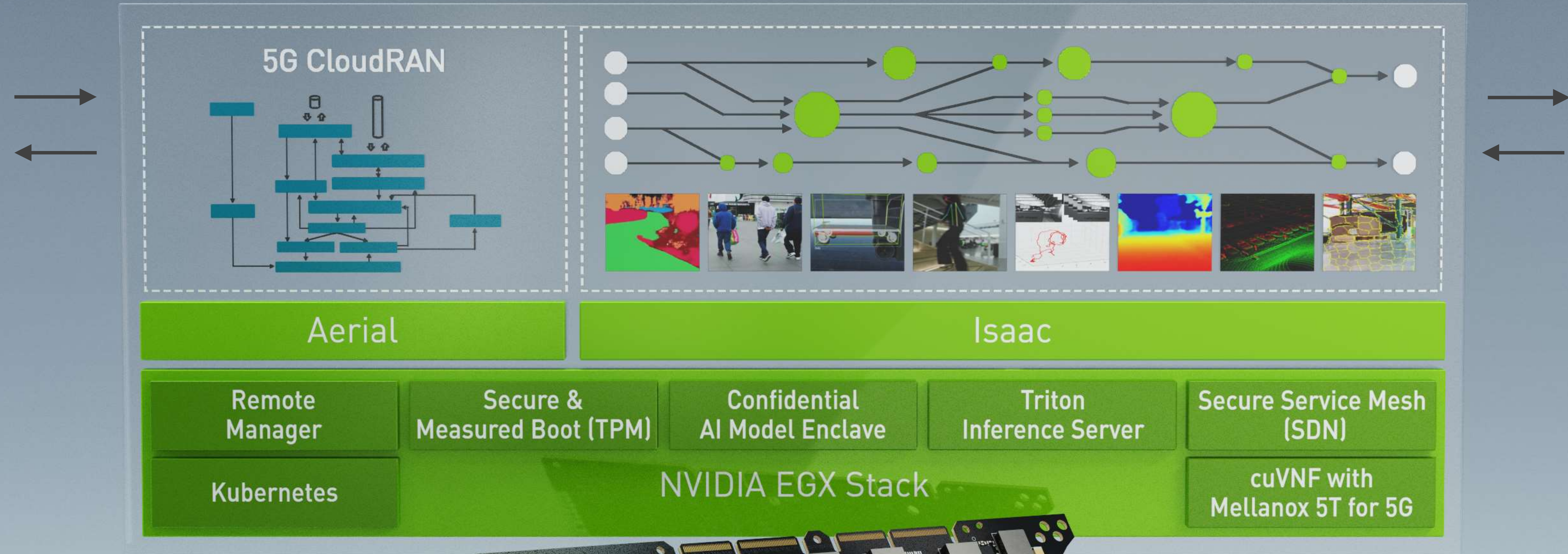
NVIDIA EGX で動作する METROPOLIS ビデオ AI と AERIAL 5G



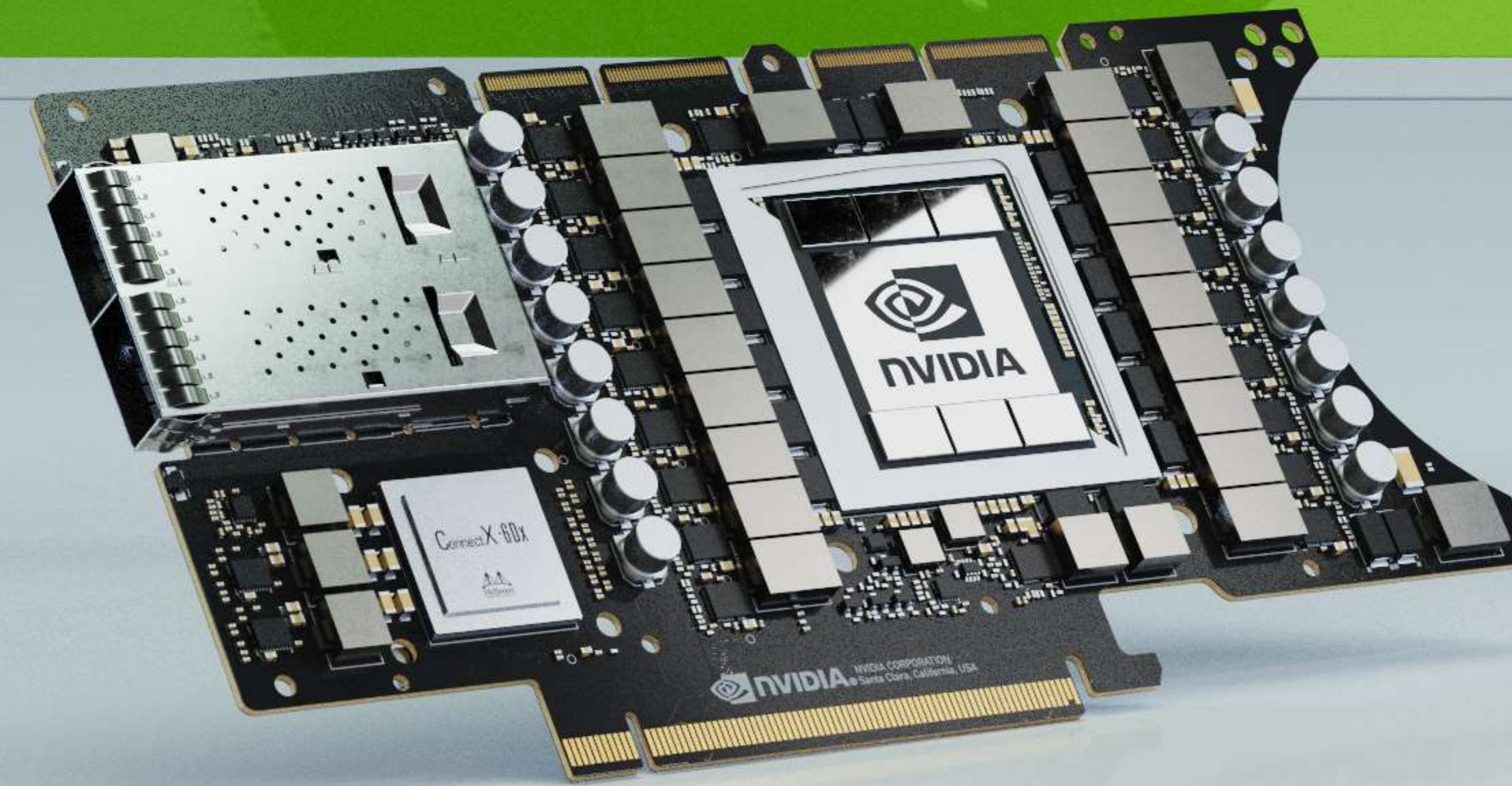
NVIDIA EGX で動作する ISAAC ロボティクス ファクトリと AERIAL 5G



実際の工場



バーチャル工場
デジタル ツイン







発表: BMW が NVIDIA ISAAC ロボティクスを採用

“The Power of Choice”

40 以上の BMW 車種、車あたり 100 オプション
99% の注文がカスタマイズ/ユニーク
2¹⁰⁰ 種類の可能な構成

“Raw Parts In, Parts Trays Out”

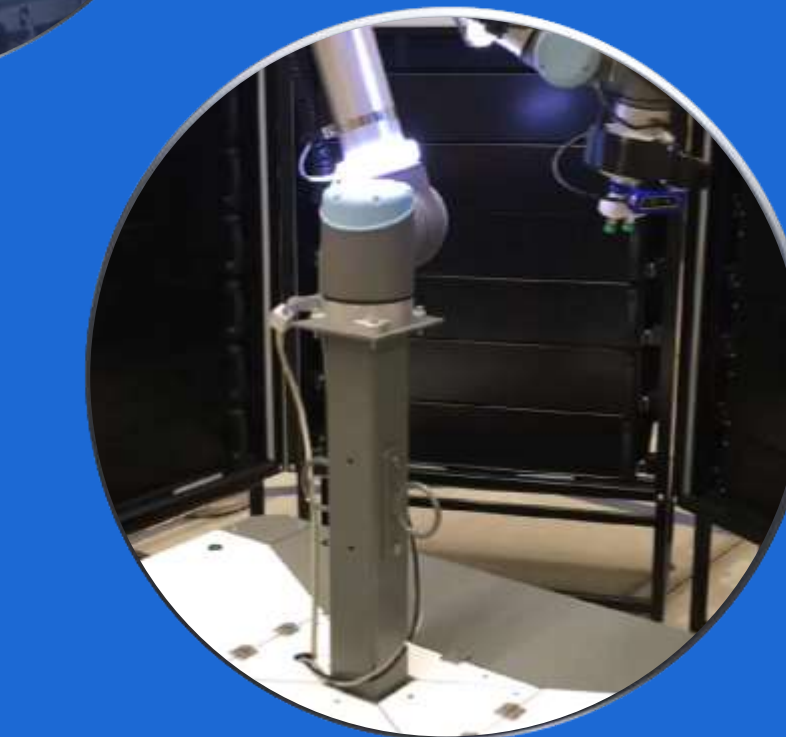
3,000 万のパーツが毎日到着
1,800 社のサプライヤーから、31 の工場
23 万のパーツナンバー

“Just in Time, Just in Sequence”

ラインあたり最大 10 台
56 秒ごとに新しい車を生産



SplitBot



PickBot



PlaceBot



Smart Transport Robot (STR)



SortBot

NVIDIA EGX エコシステム

5G & CloudRAN



セキュリティ & ネットワーキング



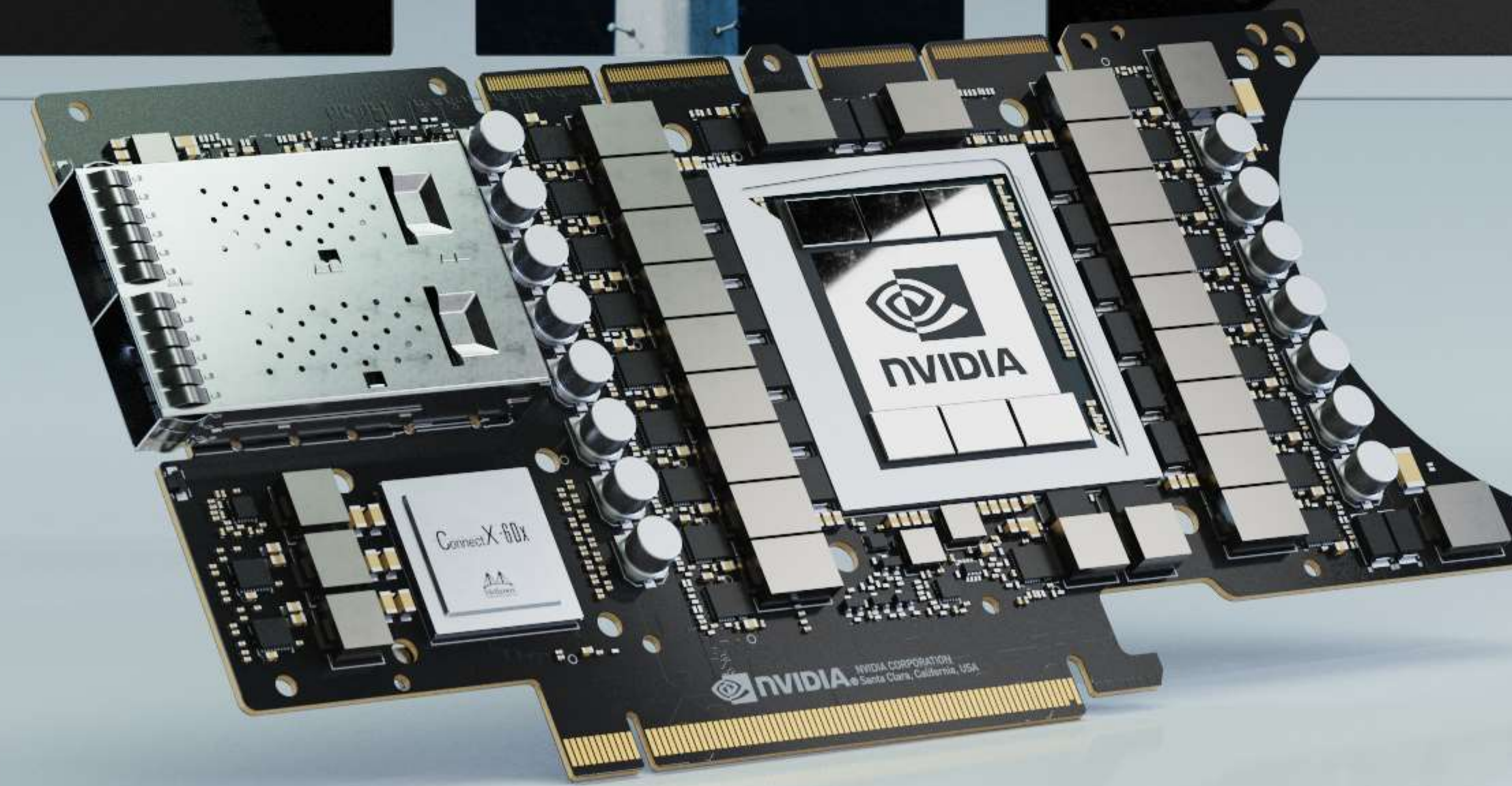
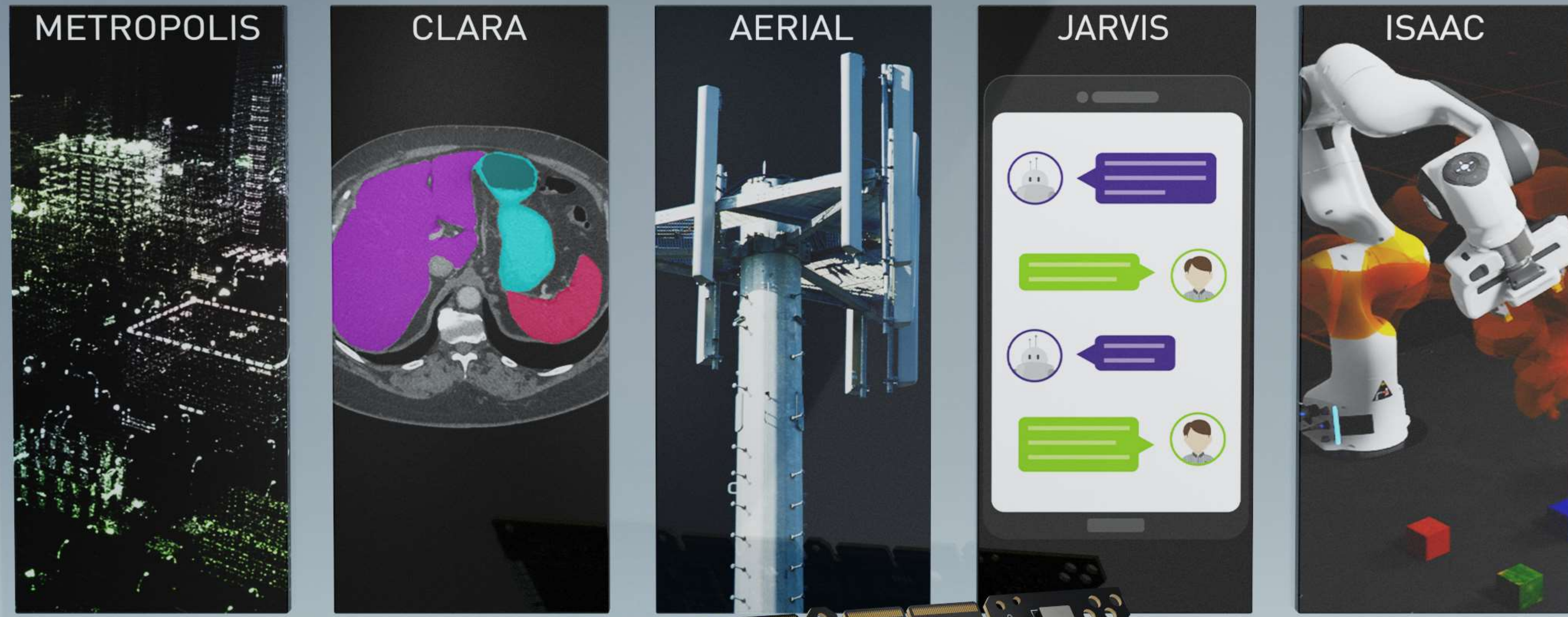
インフラストラクチャ



システム



クラウド



対話型 AI



インテリジェント ビデオ分析



ロボティクス



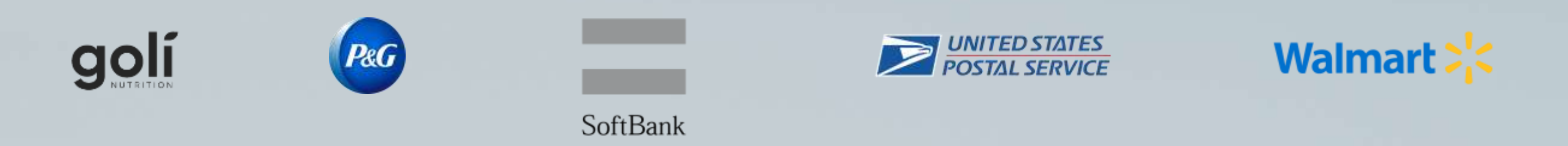
医療



工業



業界のリーダー

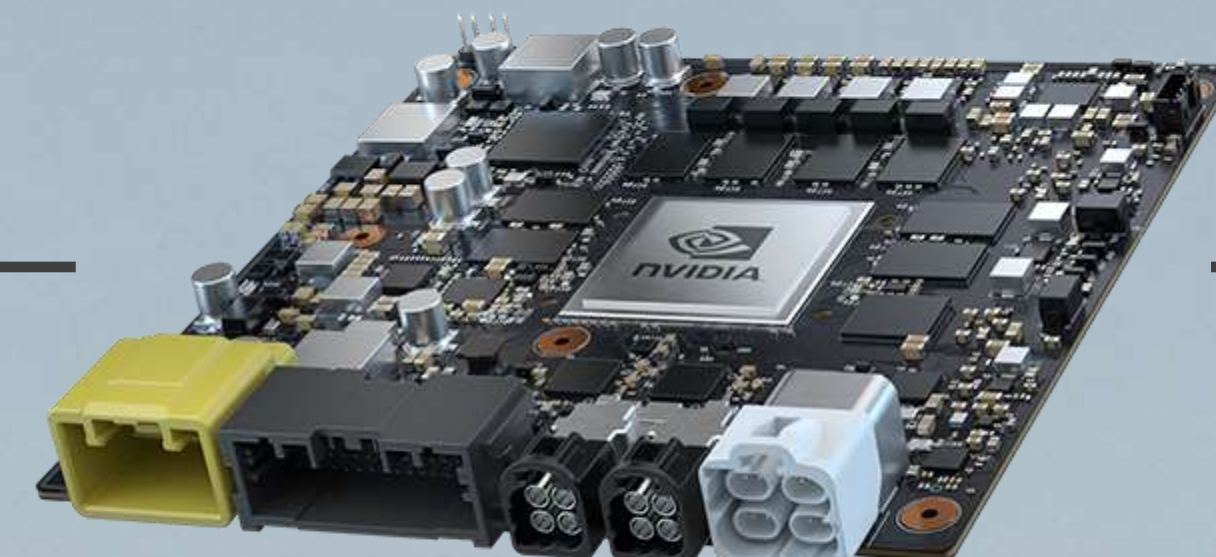


ORIN と AMPERE 搭載の NVIDIA DRIVE 5W から 2,000 TOPS まで — ひとつのプログラマブル アーキテクチャ

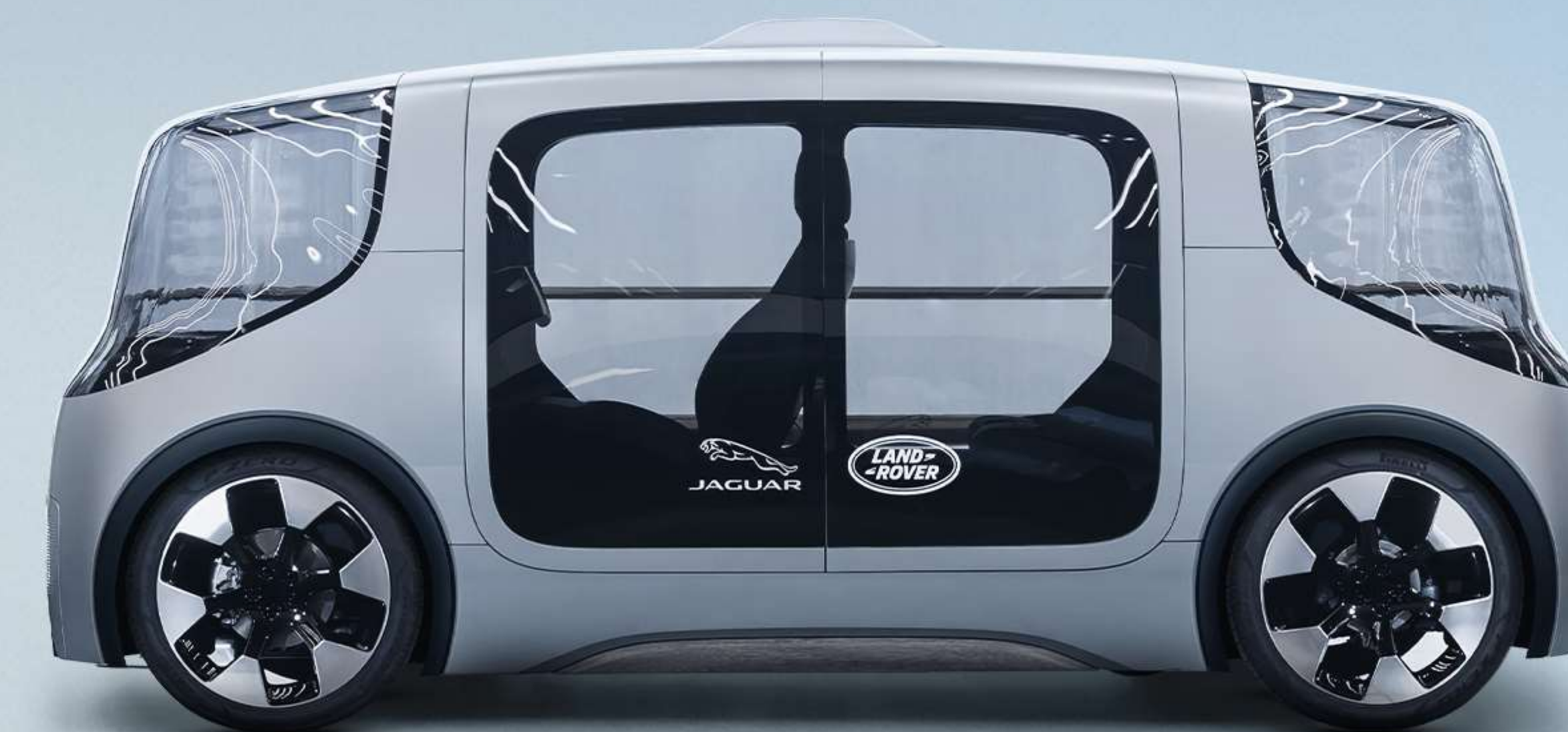
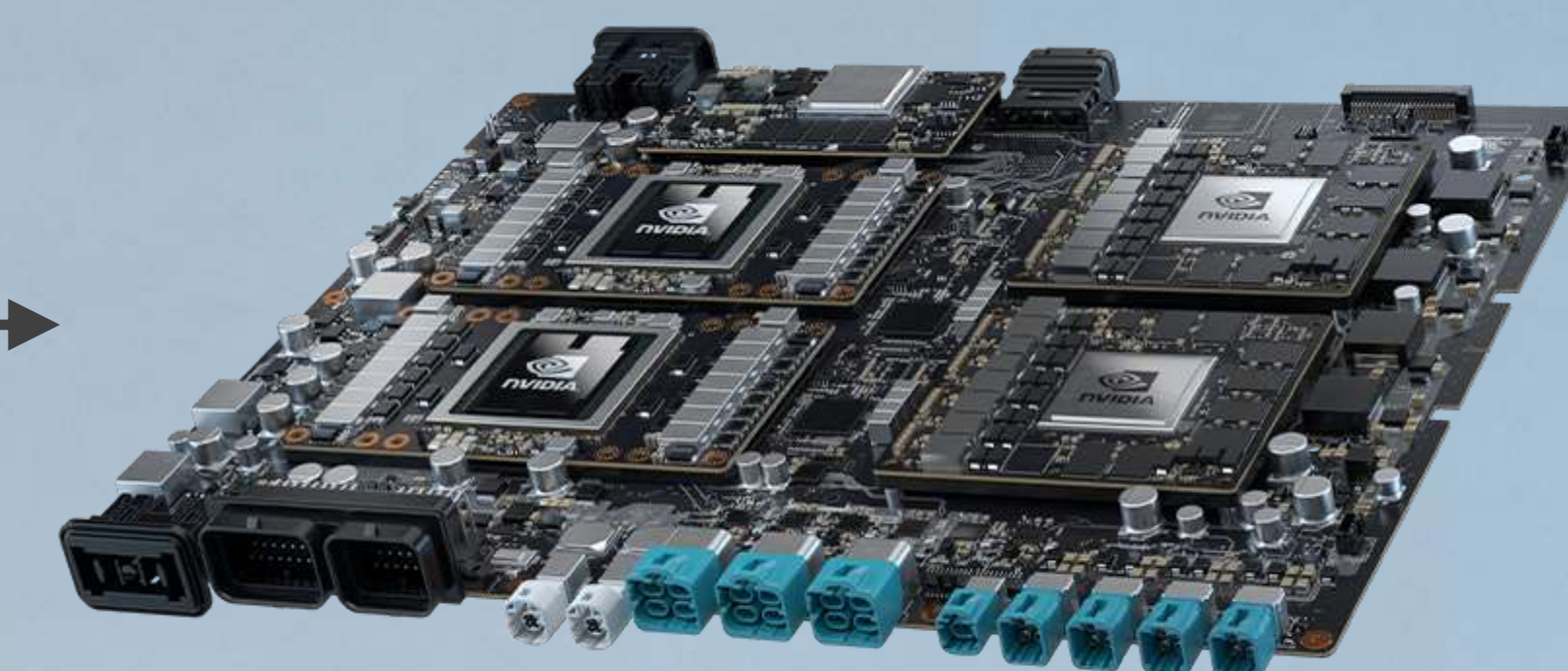
ADAS
Windshield NCAP
10 TOPS、5W



L2+
Autopilot
200 TOPS、45W

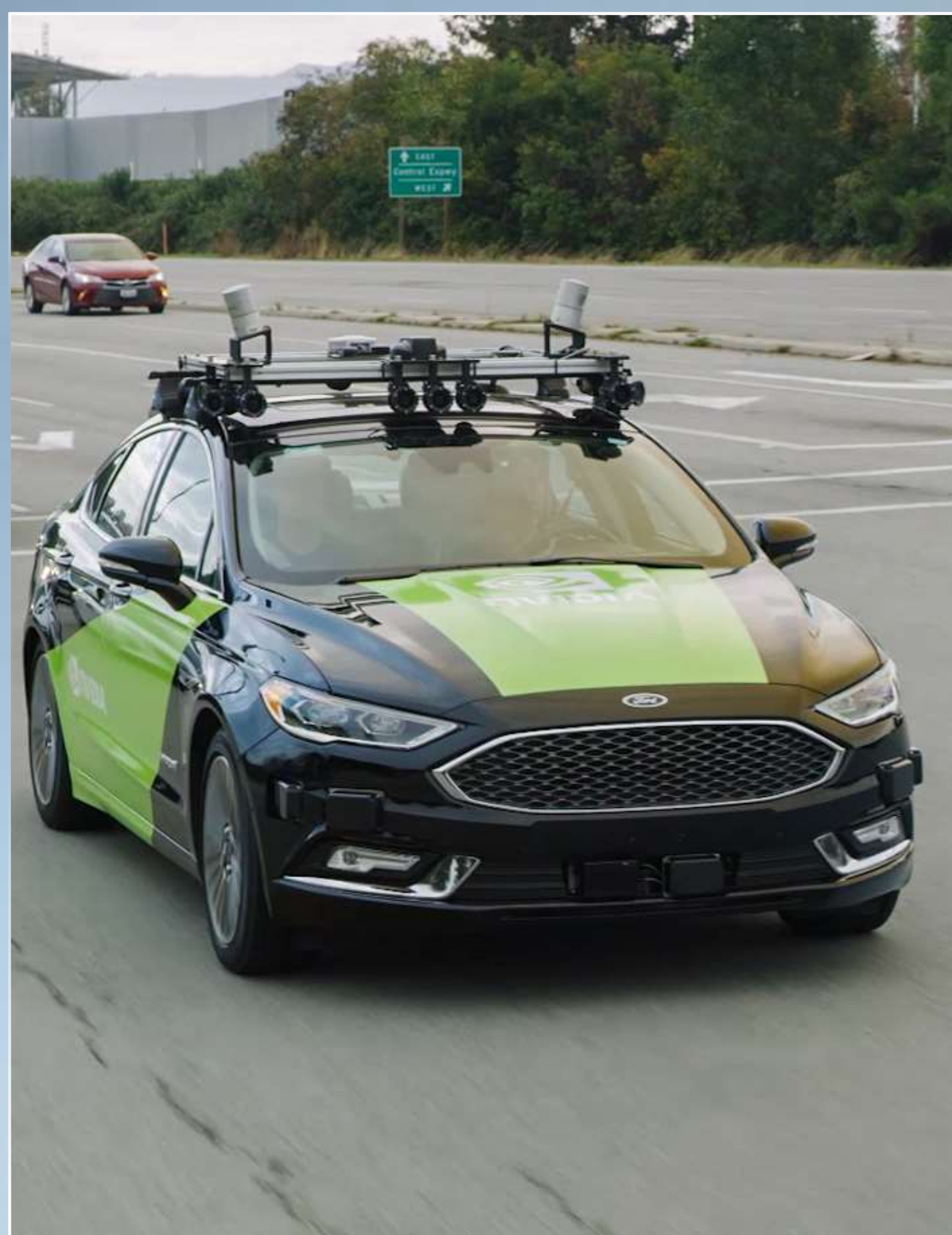


L5
Robotaxi
2,000 TOPS、800W

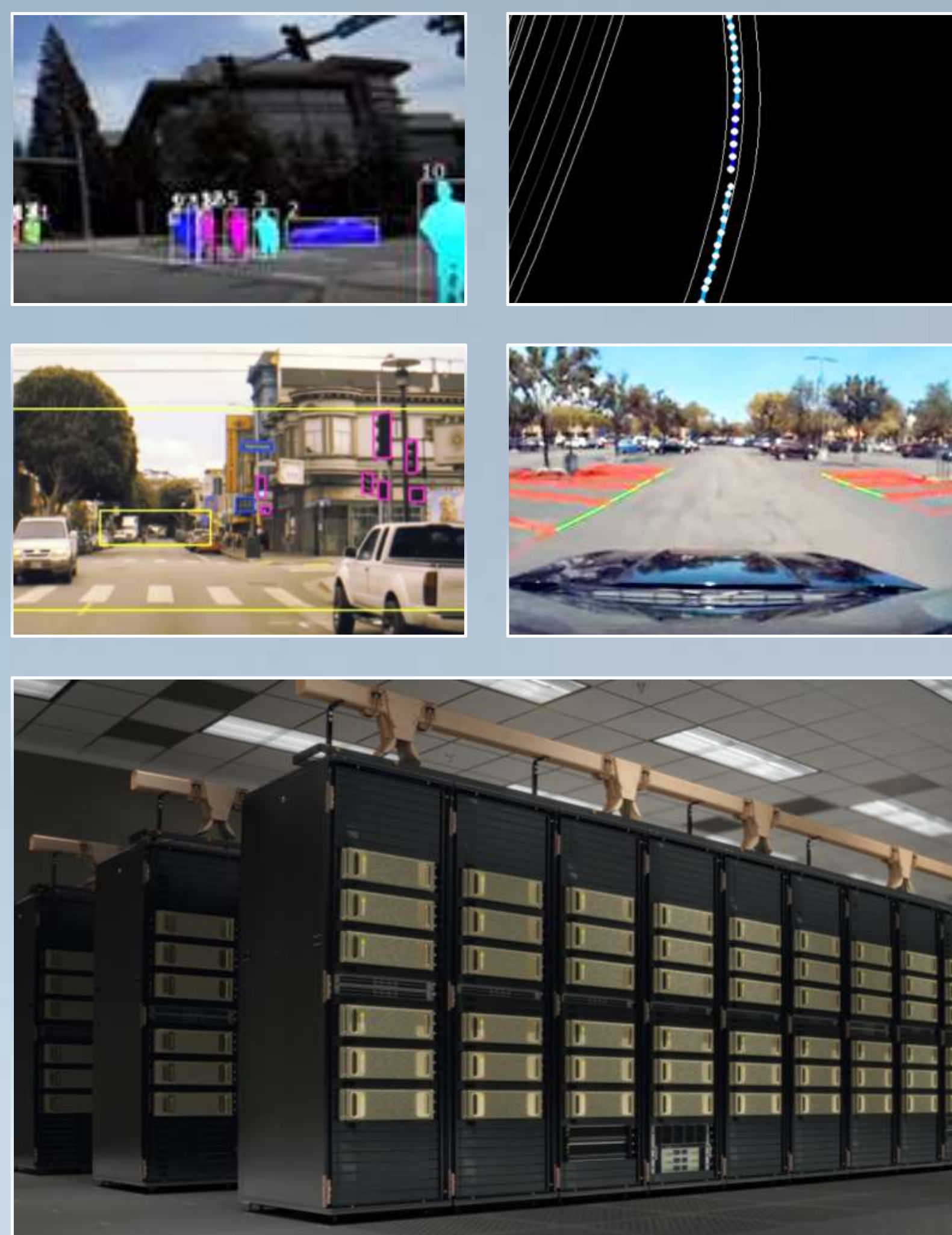


NVIDIA DRIVE — ソフトウェア定義 AV プラットフォーム

データ収集



モデルの学習



シミュレーション



DRIVE AV



DRIVE IX



DRIVE RC





NVIDIA DRIVE グローバル エコシステム

乗用車



トラック



サプライヤー



移動サービス



スタートアップ



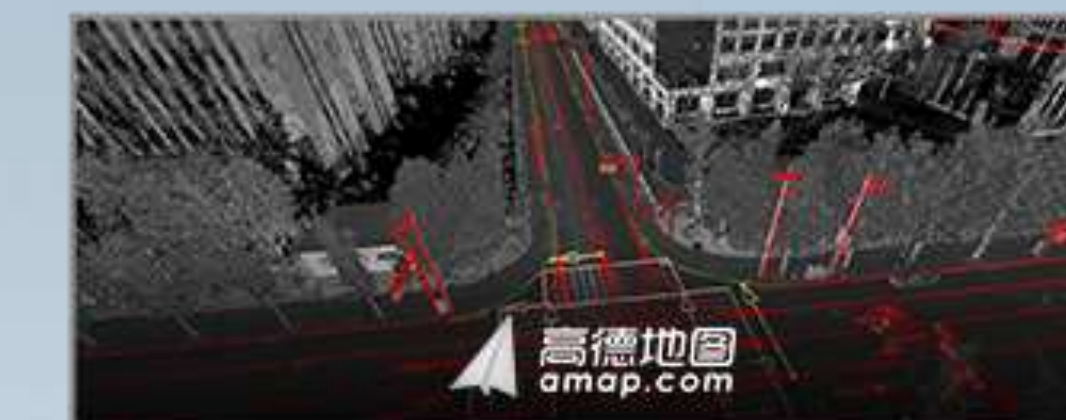
ソフトウェア



マッピング



シミュレーション



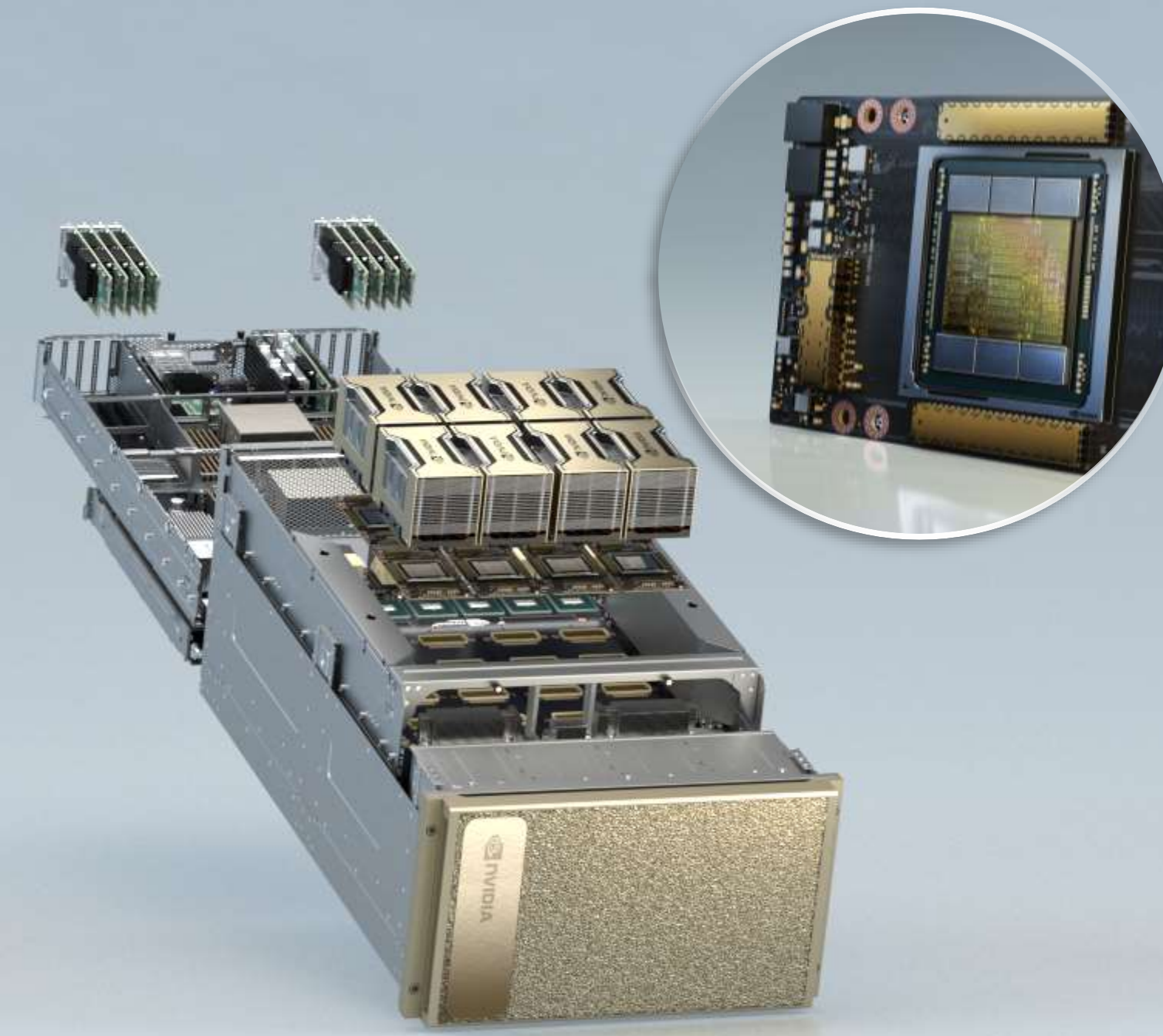
GTC 2020 発表



データセンタースケール
コンピューティング



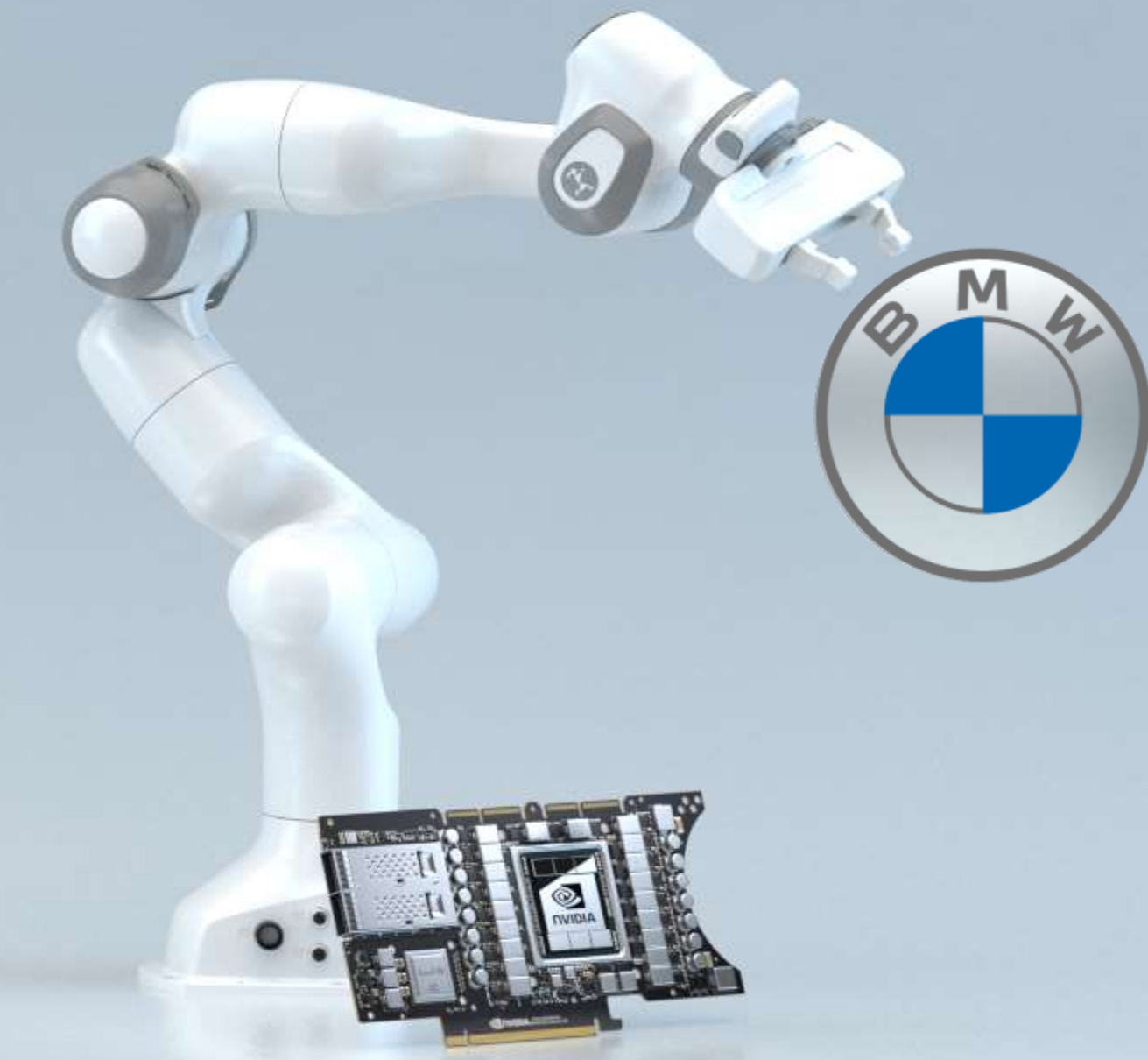
Omniverse RTX サーバー



A100 と DGX A100



NVIDIA AI



EGX と ISAAC



nVIDIA