



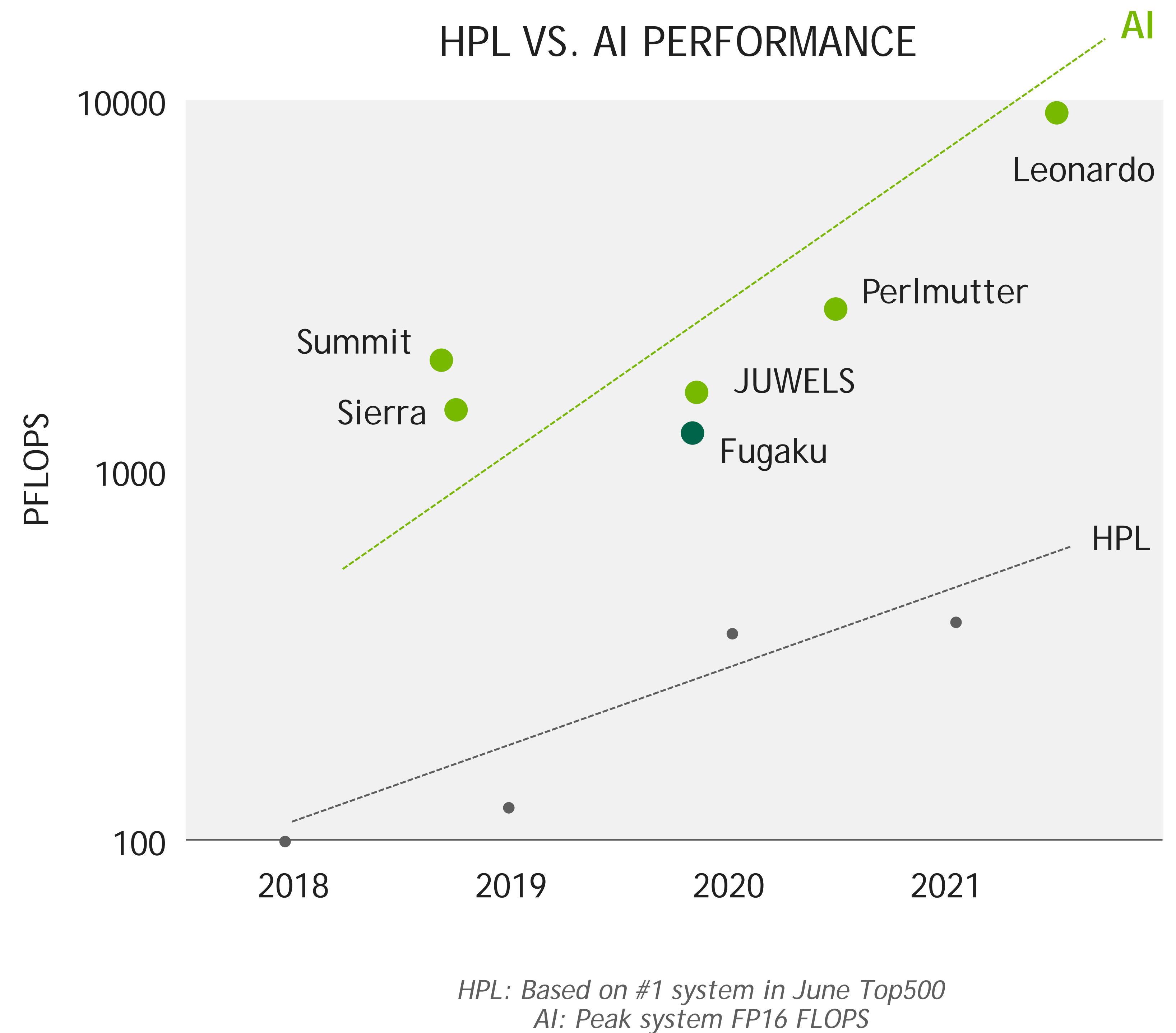
GPUから見たワークロードと高速化の歴史と展望

エンタープライズ事業本部 事業本部長 井崎 武士

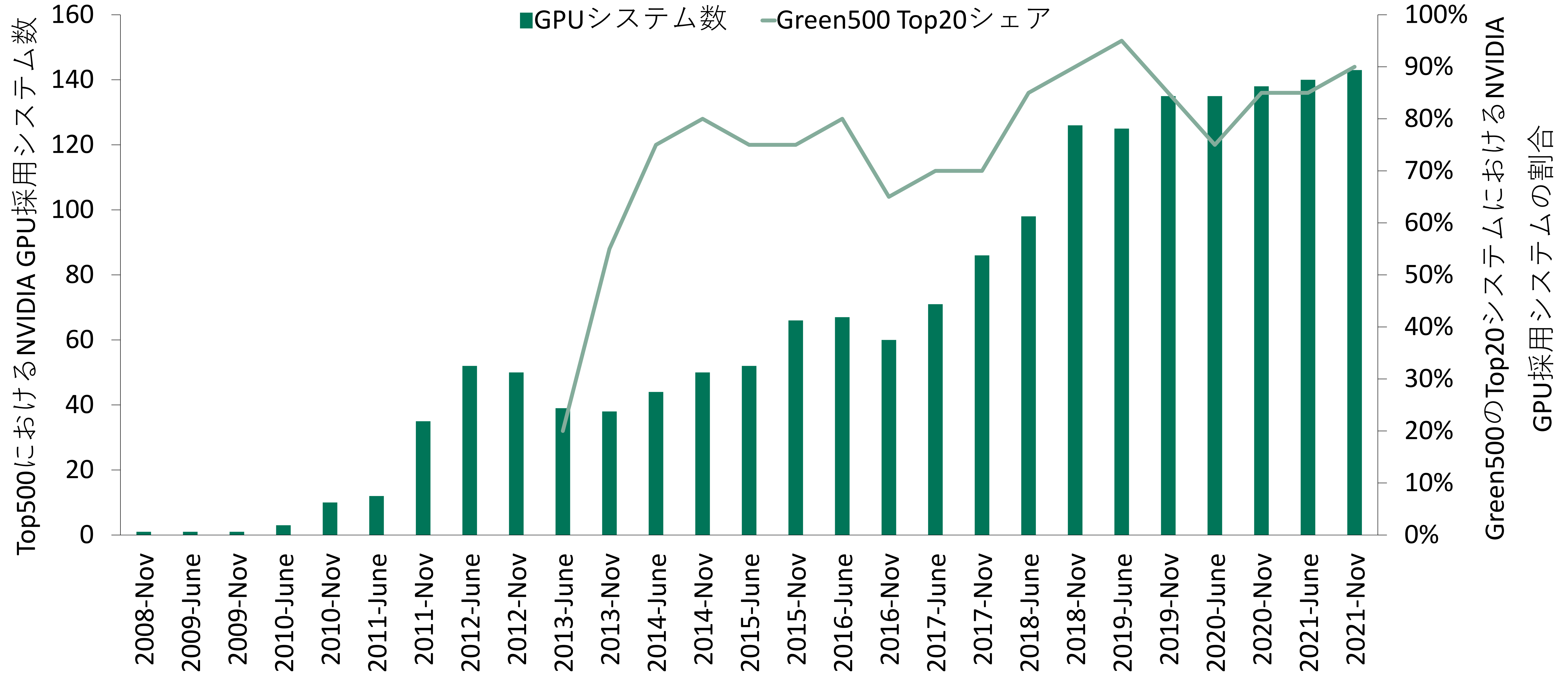
アプリケーションの多様化と支えるインフラの柔軟性



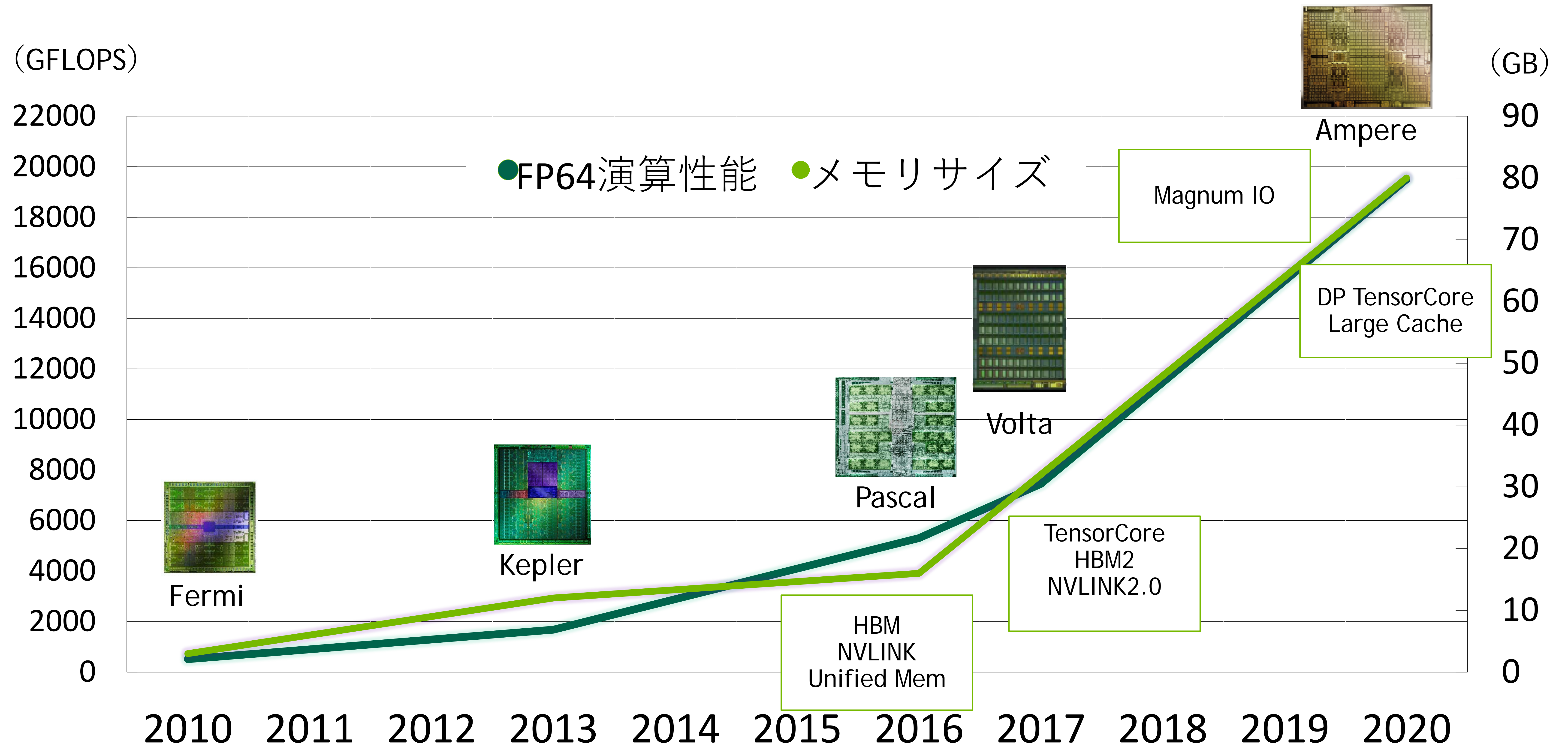
エクサスケールの時代
AI スーパーコンピューティング
AI が現代のHPCの新たな成長ドライバー



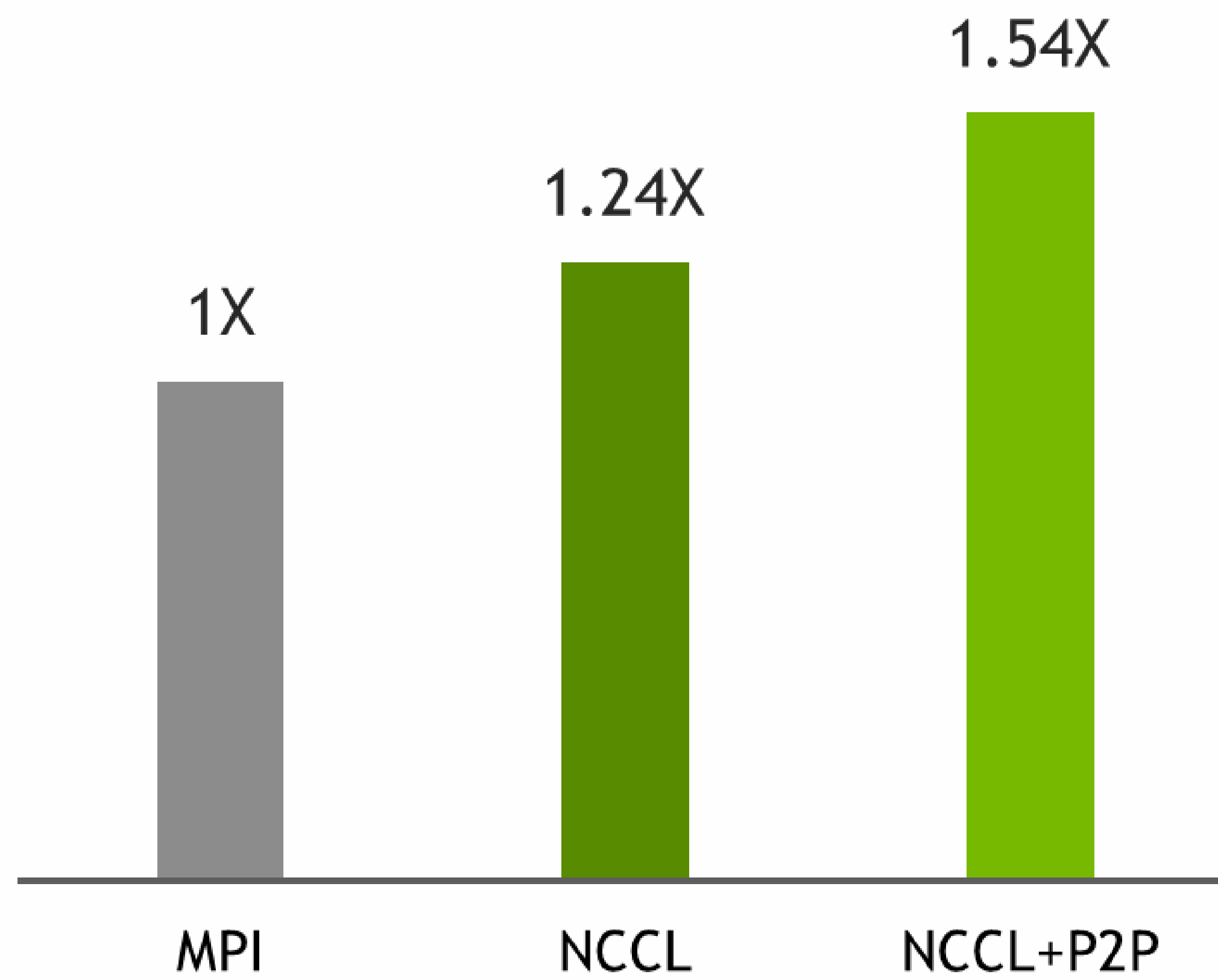
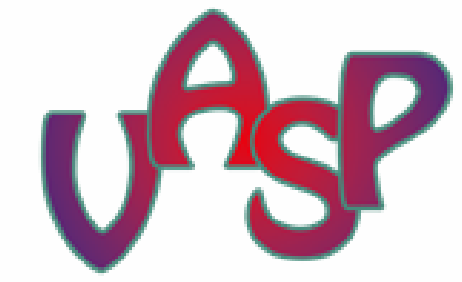
TOP500・GREEN500におけるNVIDIA GPU採用システム数の推移



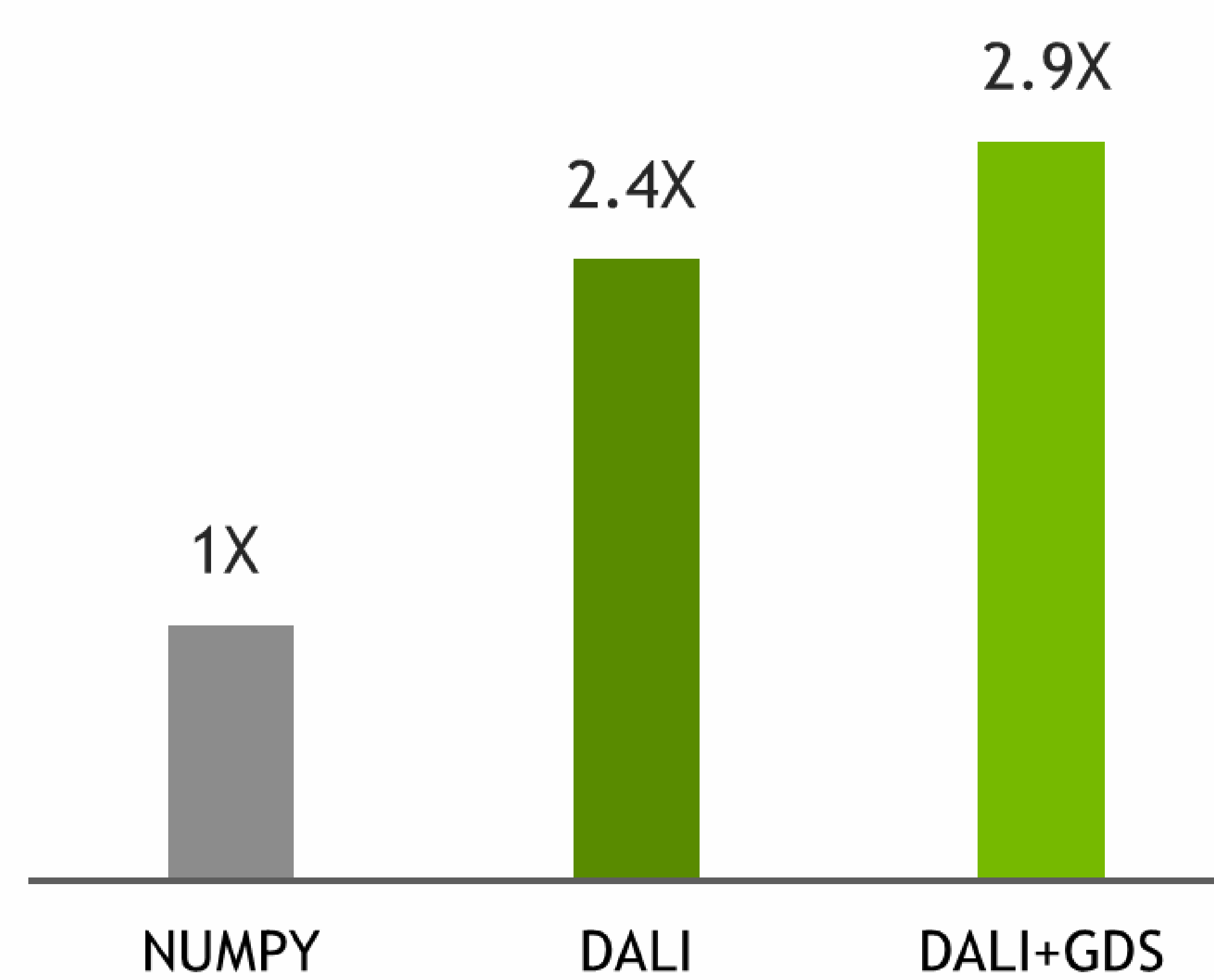
GPUハードウェア性能の推移



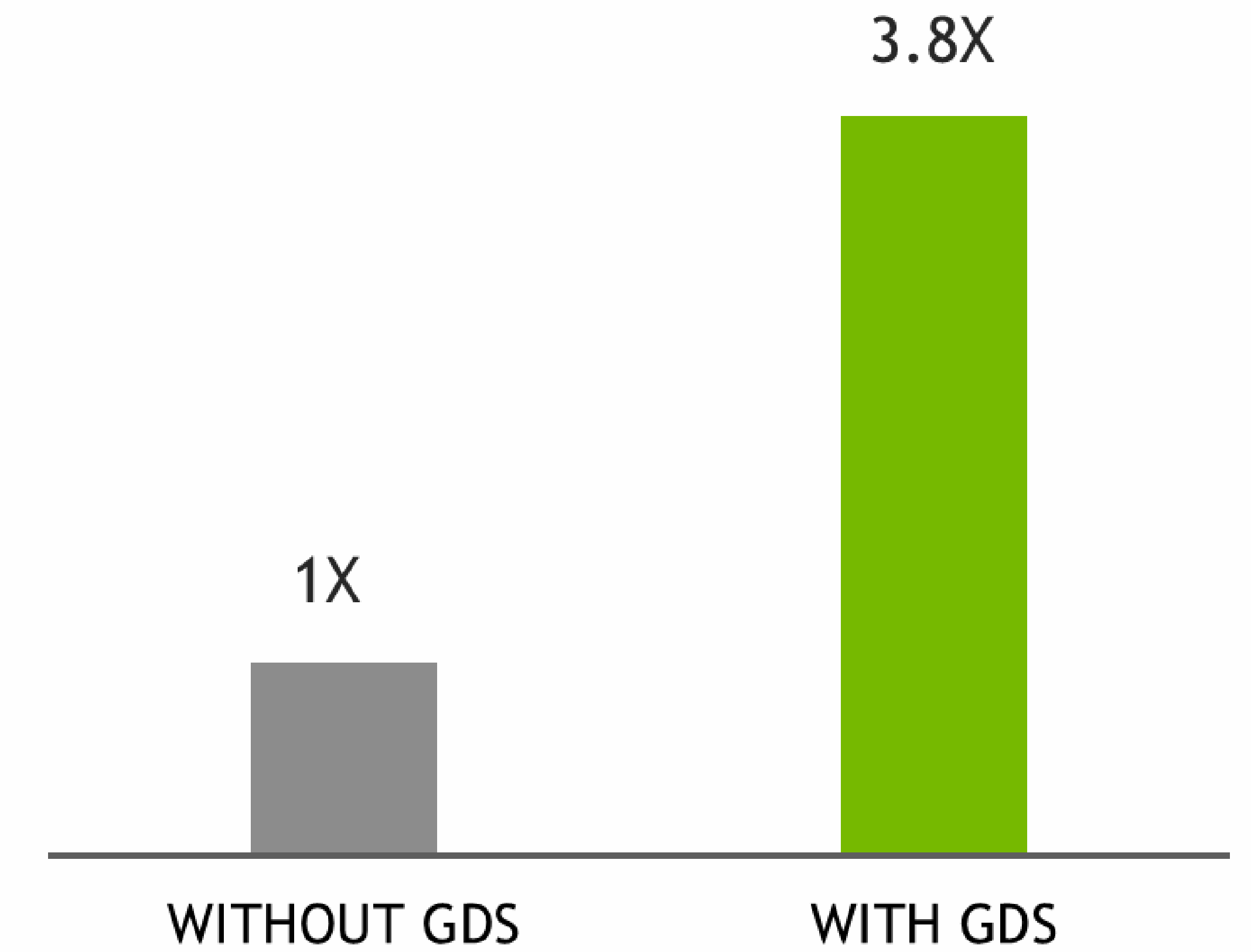
MAGNUM IO によるベネフィット



NCCL P2Pを用いてリンモレイヤーのシミュレーション性能を向上



気象シミュレーションにおける異常気象現象のセグメンテーション



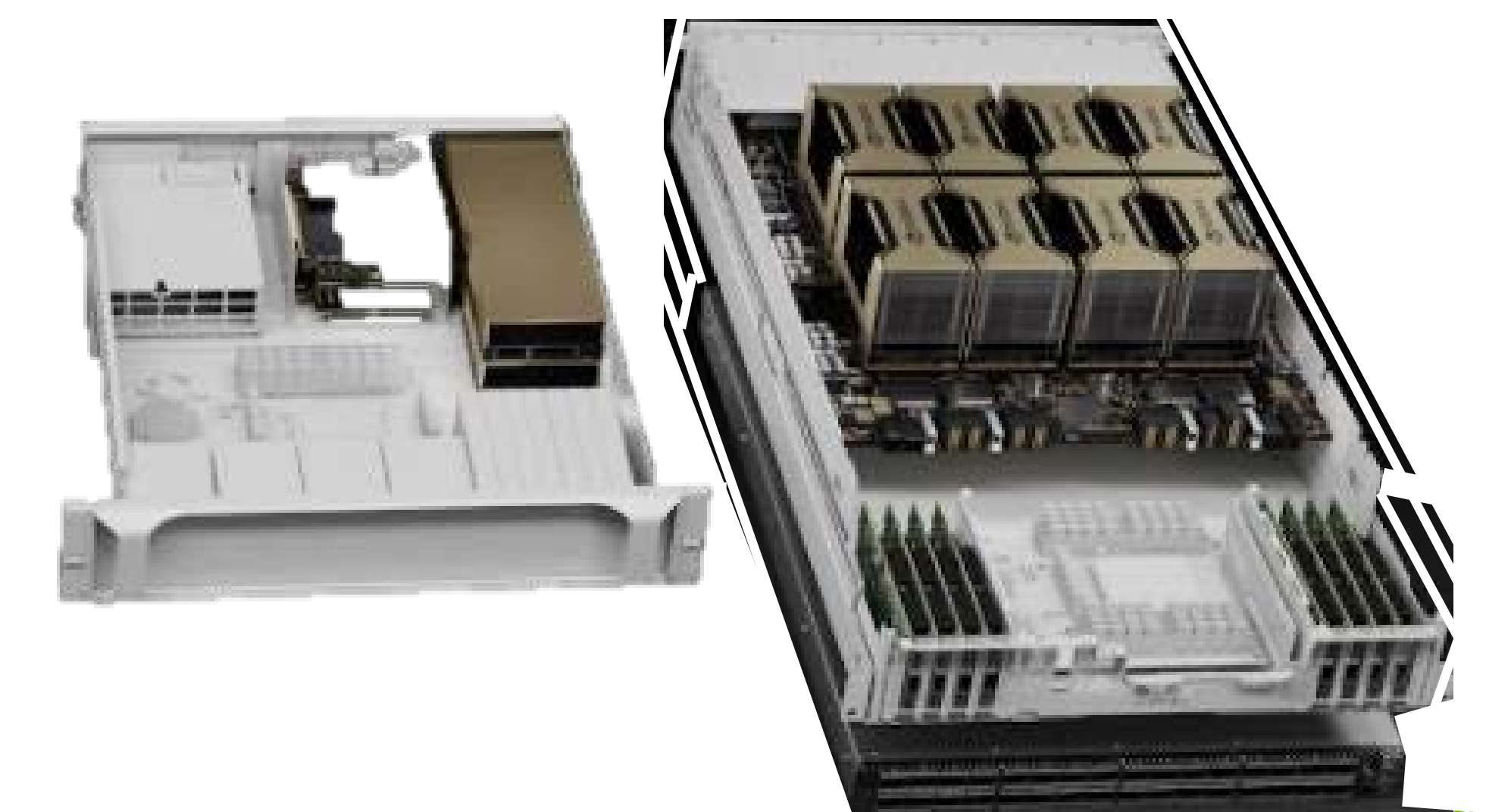
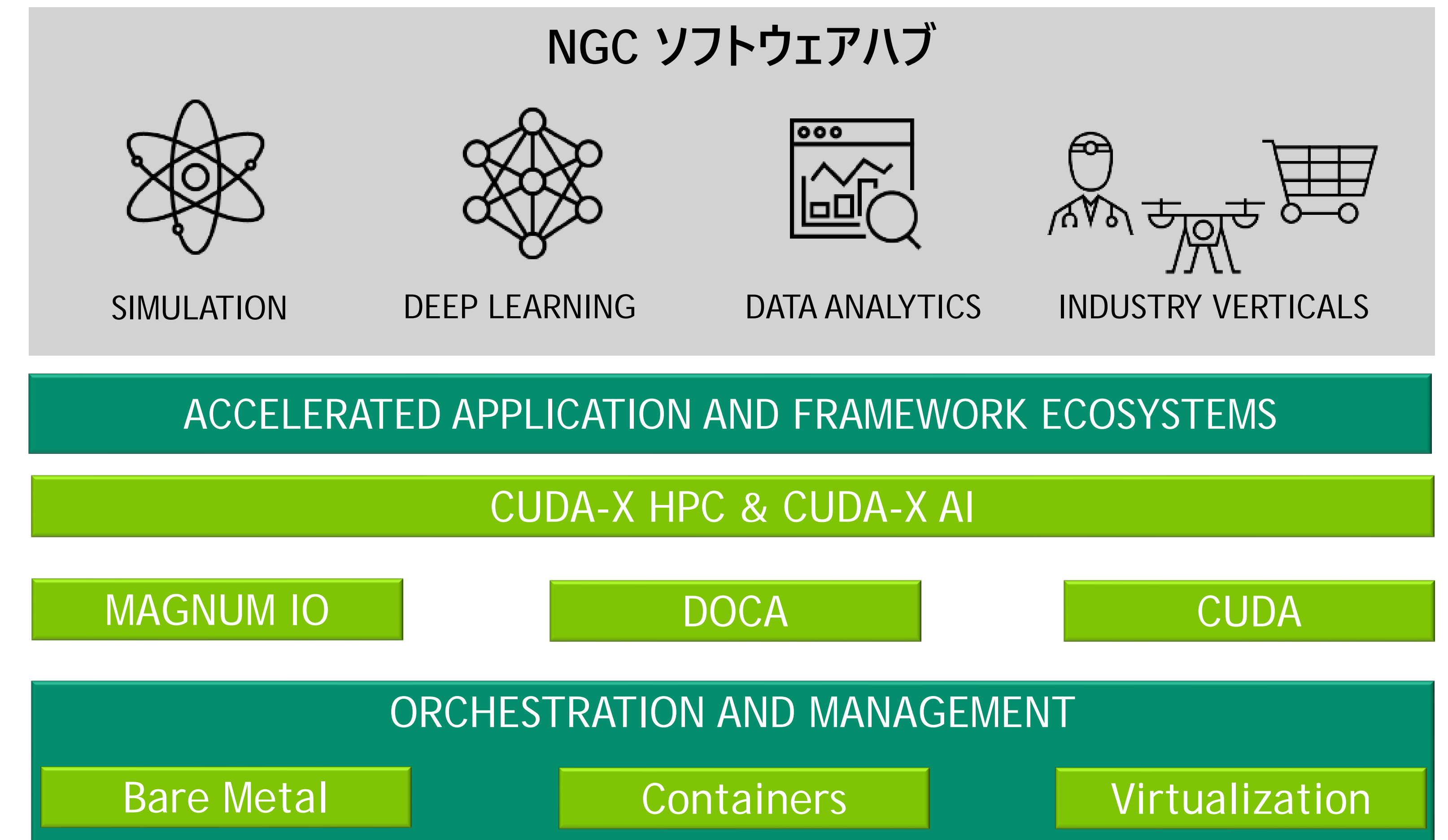
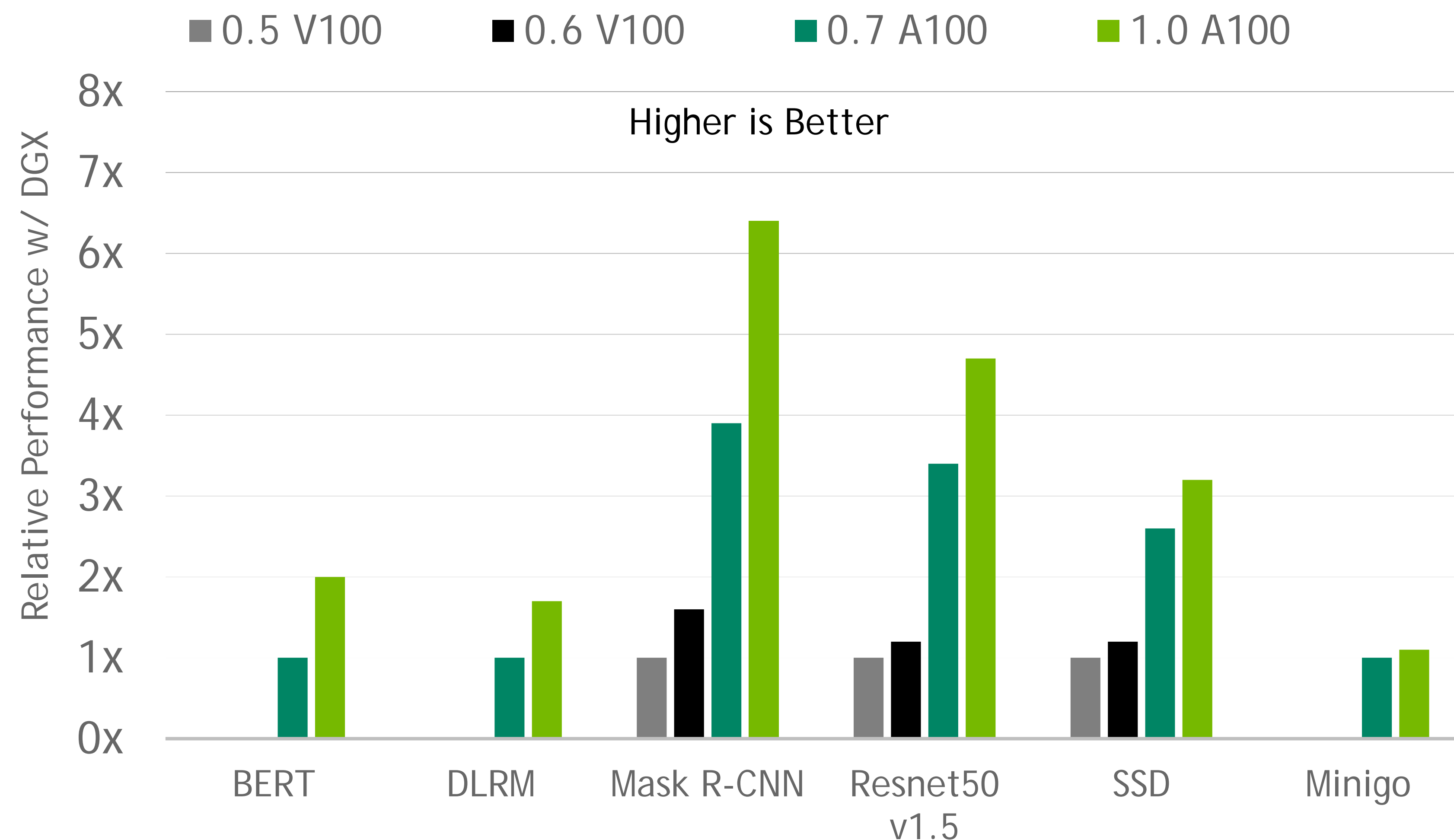
DGX A100における帯域最大のリモートファイルリード

NVIDIA データセンタープラットフォーム



2.5年でGPU当たり最大6.5倍の性能向上

NVIDIA AI はフルスタックの改善で継続的に向上



Results normalized for throughput due to higher accuracy requirements on latest round of MLPerf 1.0

MLPerf ID 0.5/0.6/0.7/1.0 comparison:

BERT: 0.7-19/1.0-1060 | DLRM: 0.7-17/1.0-1058 | Mask R-CNN: 0.5-13/0.6-9/0.7-19/1.0-1060 | ResNet50v1.5: 0.5-12/0.6-8/0.7-18/1.0-1059 |

SSD: 0.5-13/0.6-9/0.7-19/1.0-1059 | MiniGo: 0.7-20/1.0-1059

MLPerf name and logo are trademarks. See www.mlperf.org for more information.

NVIDIA プラットフォームの多様なプログラミング手法

CPU, GPU, Network

標準言語を使った高速化

ISO C++, ISO Fortran

```
std::transform(par, x, x+n, y, y,  
              [=](float x, float y){ return y +  
a*x; }  
);
```

```
do concurrent (i = 1:n)  
  y(i) = y(i) + a*x(i)  
enddo
```

```
import cunumeric as np  
...  
def saxpy(a, x, y):  
  y[:] += a*x
```

インクリメンタル ポータブル 最適化

OpenACC, OpenMP

```
#pragma acc data copy(x,y) {  
...  
std::transform(par, x, x+n, y, y,  
              [=](float x, float y){  
                return y + a*x;  
              });  
...  
}  
  
#pragma omp target data map(x,y) {  
...  
std::transform(par, x, x+n, y, y,  
              [=](float x, float y){  
                return y + a*x;  
              });  
...  
}
```

プラットフォーム特化

CUDA

```
__global__  
void saxpy(int n, float a,  
           float *x, float *y) {  
  int i = blockIdx.x*blockDim.x +  
          threadIdx.x;  
  if (i < n) y[i] += a*x[i];  
}  
  
int main(void) {  
  ...  
  cudaMemcpy(d_x, x, ...);  
  cudaMemcpy(d_y, y, ...);  
  
  saxpy<<<(N+255)/256,256>>>(...);  
  
  cudaMemcpy(y, d_y, ...);  
}
```

高速化ライブラリ

Core

Math

Communication

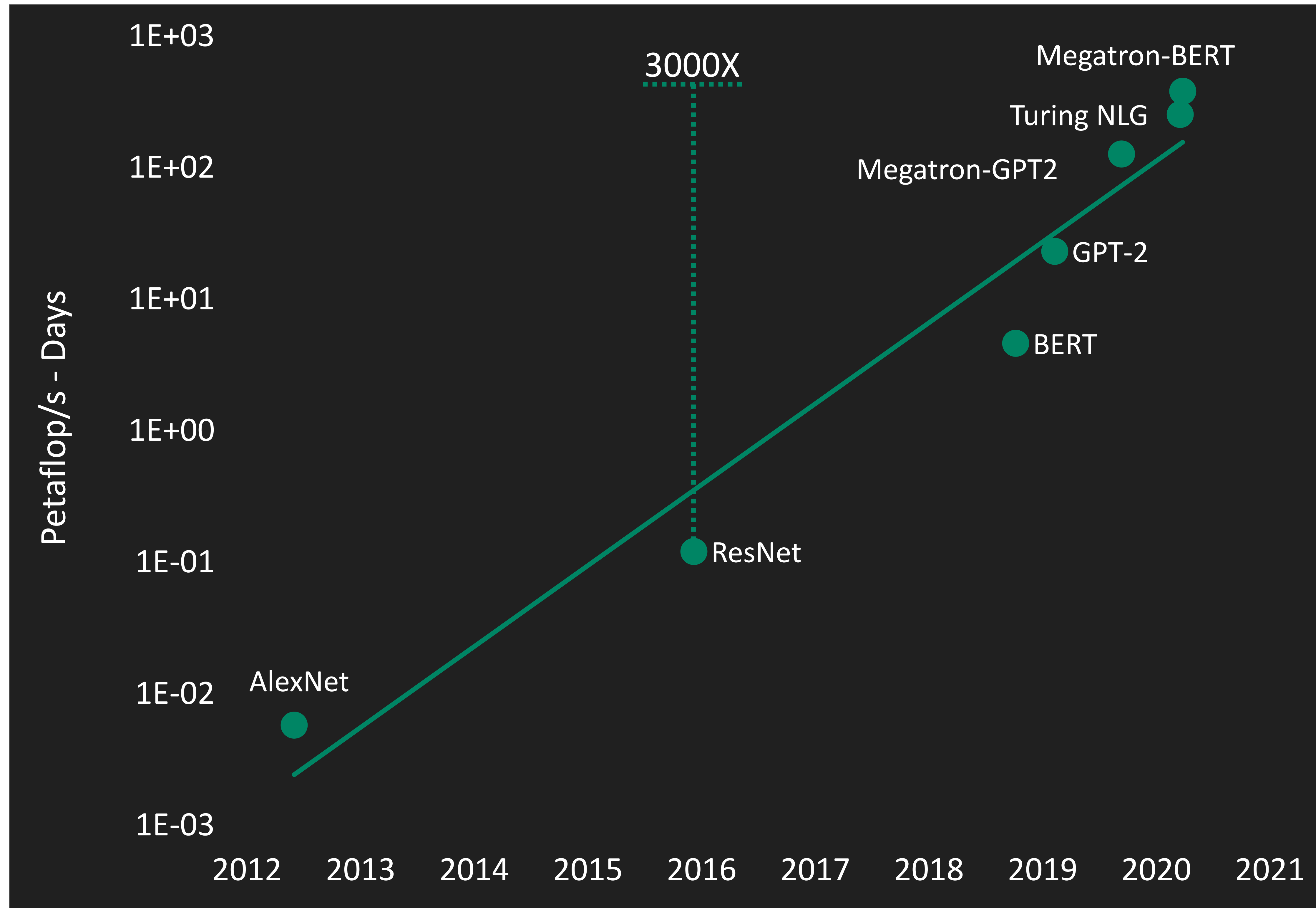
Data Analytics

AI

Quantum

学習に要求される計算性能の飛躍的拡大

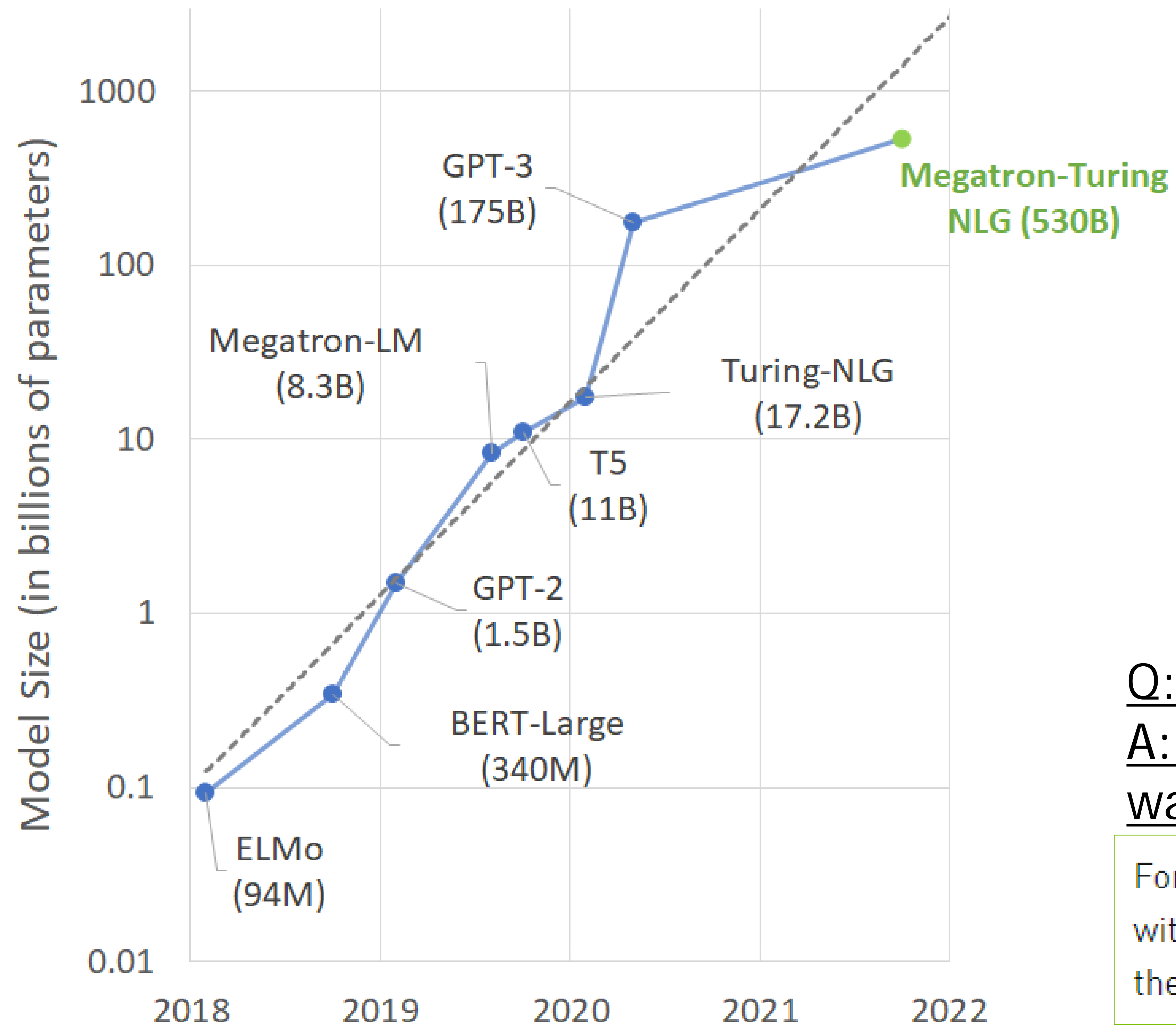
指数関数的な計算要求



Voltaの時代から大規模モデルの学習に必要な計算性能は3000倍になっている

深層分散学習の状況

モデルは引き続き巨大に



- 2021/10/11時点の最新モデルは 530B パラメータにまで拡大
 - 単純計算で 2,120 GB (in FP32)
 - モデルをすべてメモリにロードするだけで、8xA100 (80GB) サーバが、3 台強必要

Q: 実際どうやって扱っているの？

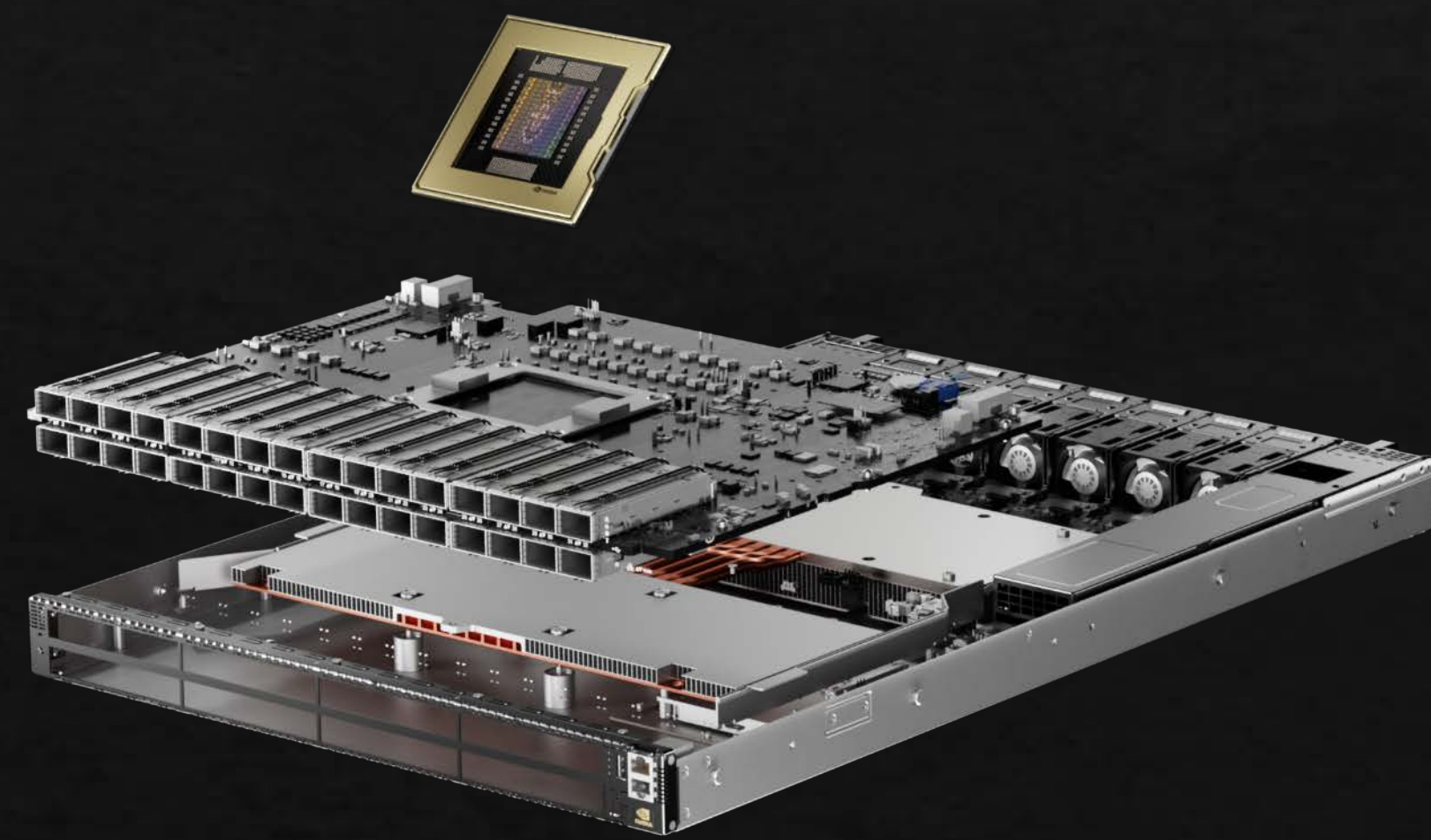
A: 280xA100 (=35 nodes) で 8-way tensor parallel & 35-way pipeline parallel のモデル並列 & クラスタ全体でデータ並列

For example, for the 530 billion model, each model replica spans 280 NVIDIA A100 GPUs, with 8-way tensor-slicing within a node and 35-way pipeline parallelism across nodes. We then use data parallelism from DeepSpeed to scale out further to thousands of GPUs.

NVIDIA QUANTUM-2

400G NDR InfiniBand

QUANTUM-2 SWITCH



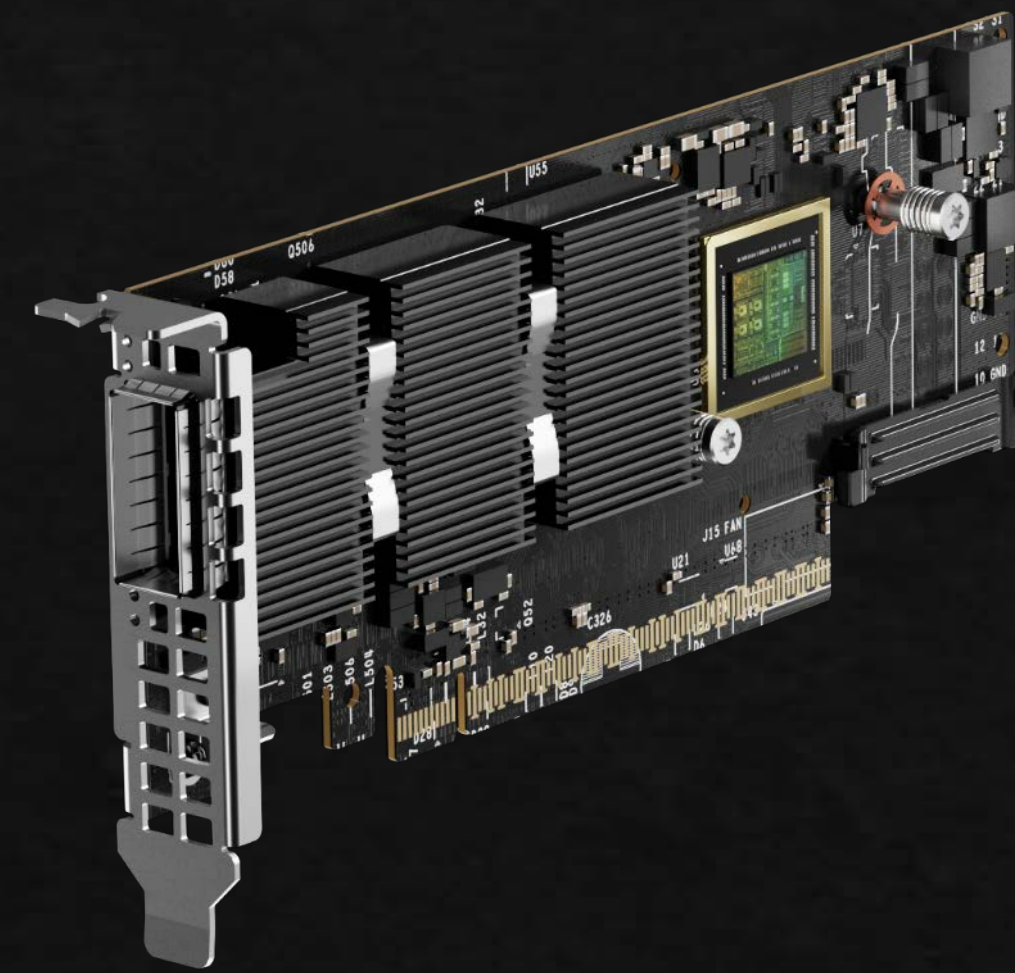
400 Gbps 64 ポート / 200 Gbps 128 ポート

32倍超の AI 高速化エンジン

3倍のスイッチングスループット

サンプル出荷中

CONNECTX-7 INFINIBAND

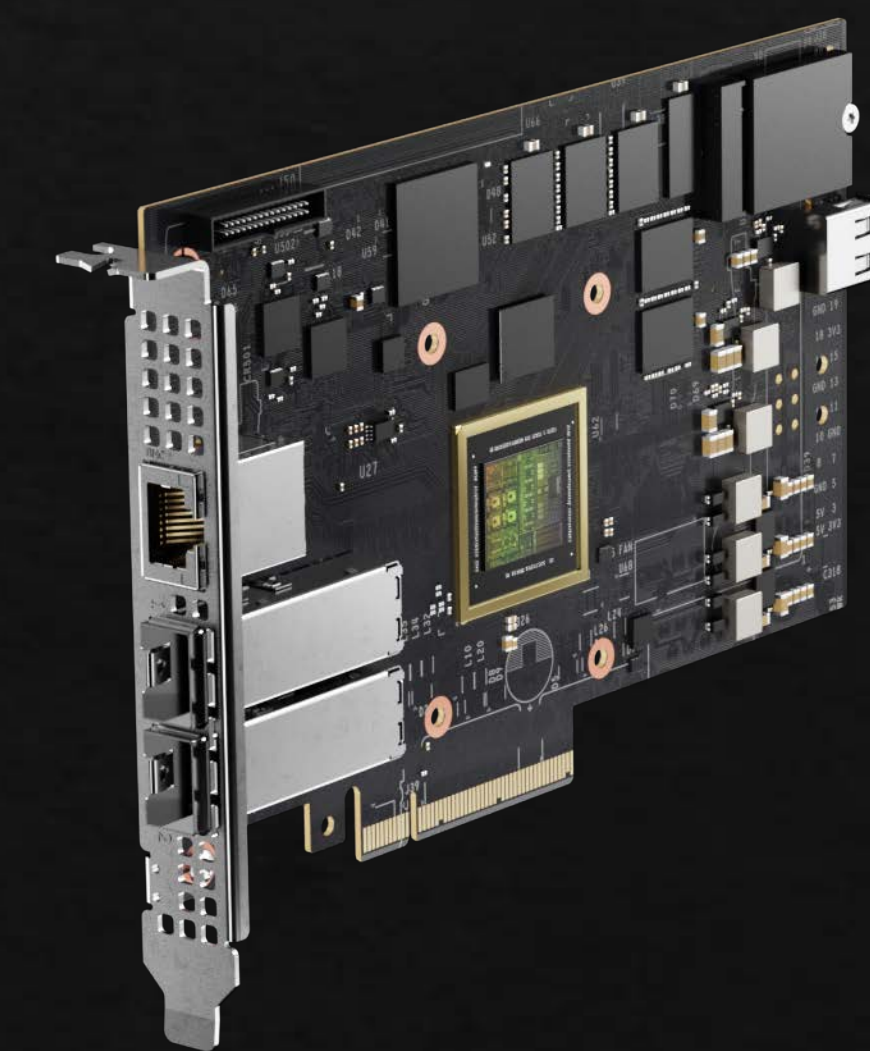


4倍のネットワーク内コンピューティング性能

2倍の GPUDirect スループット

2022年1月サンプル予定

BLUEFIELD-3 INFINIBAND



16 Arm 64 ビットコア | 400 Gbps 暗号化アクセラレータ

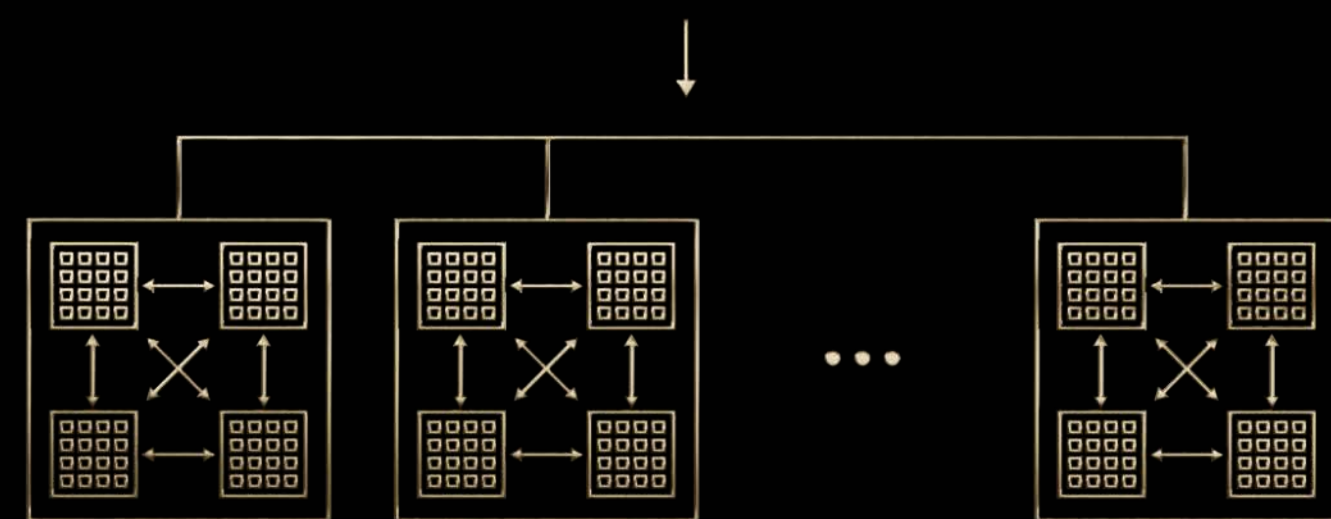
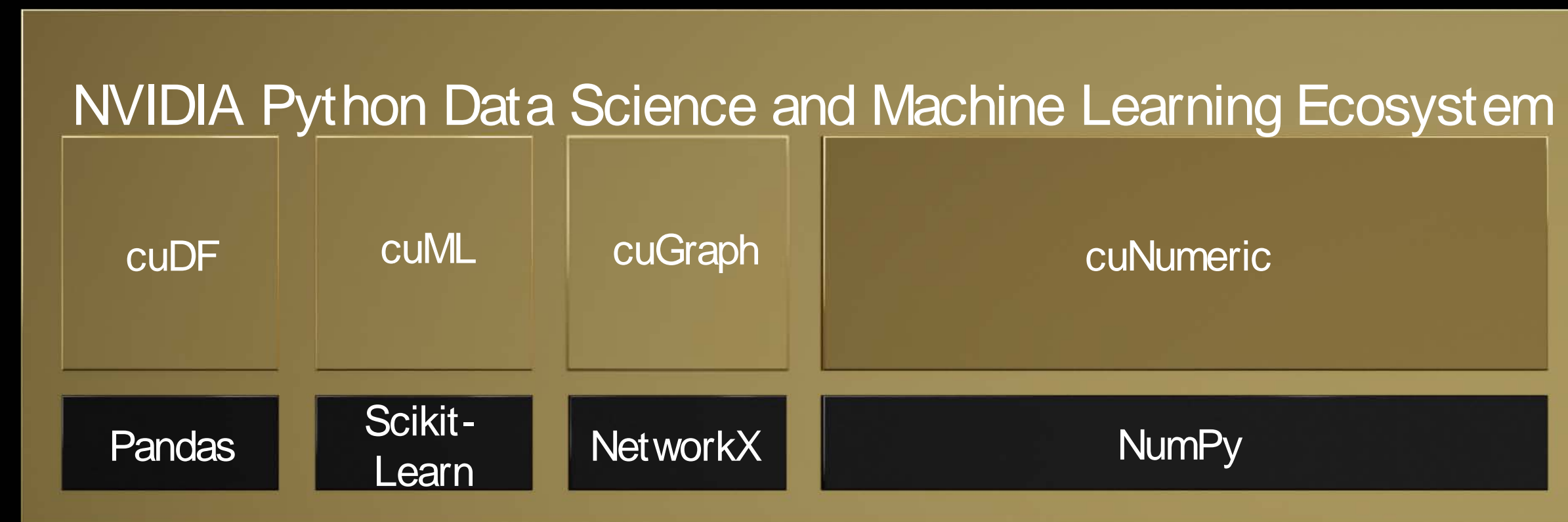
4倍のネットワーク内コンピューティング性能

2倍のGPUDirect スループット

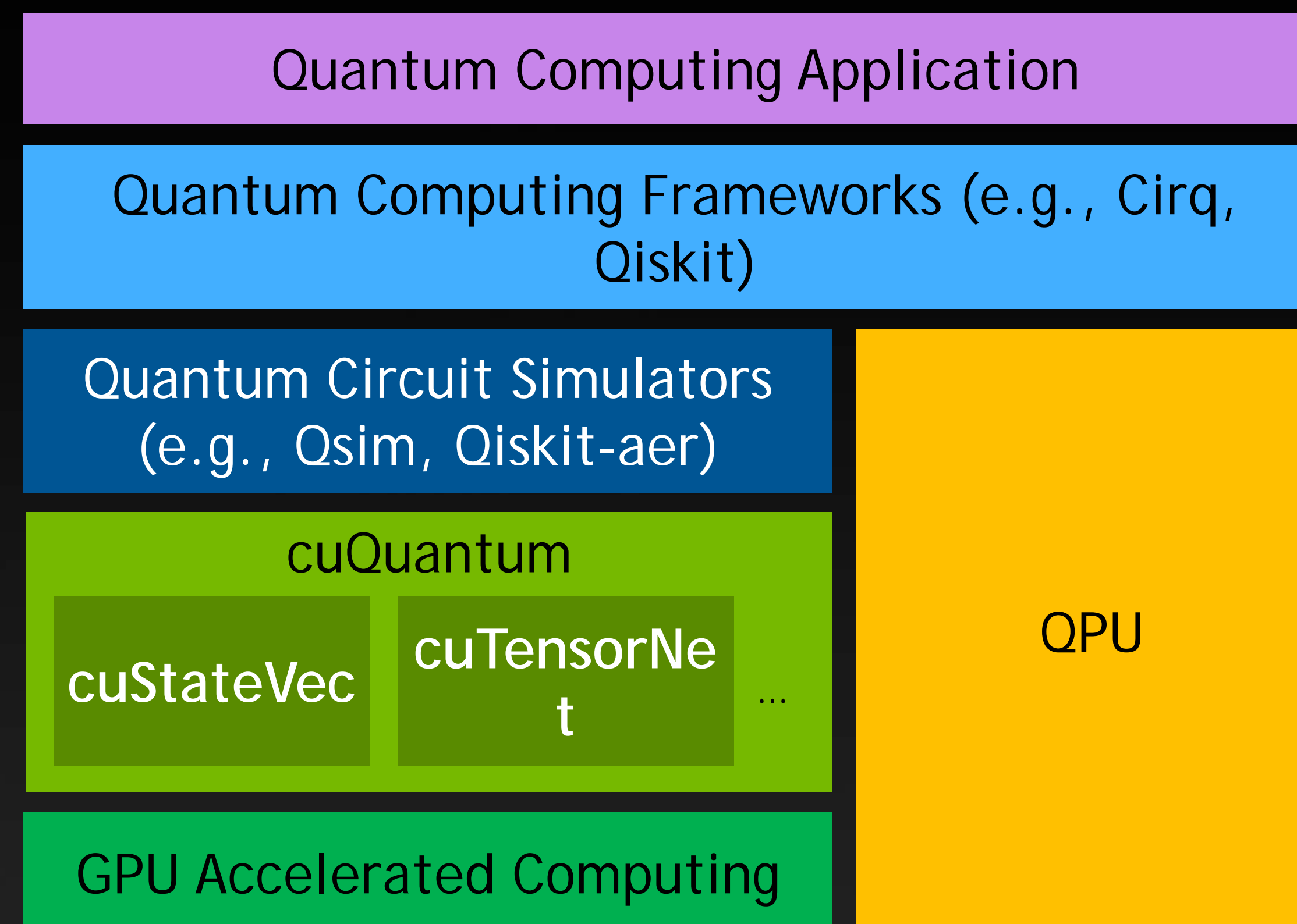
2022年5月サンプル予定

広がるコンピューティング需要

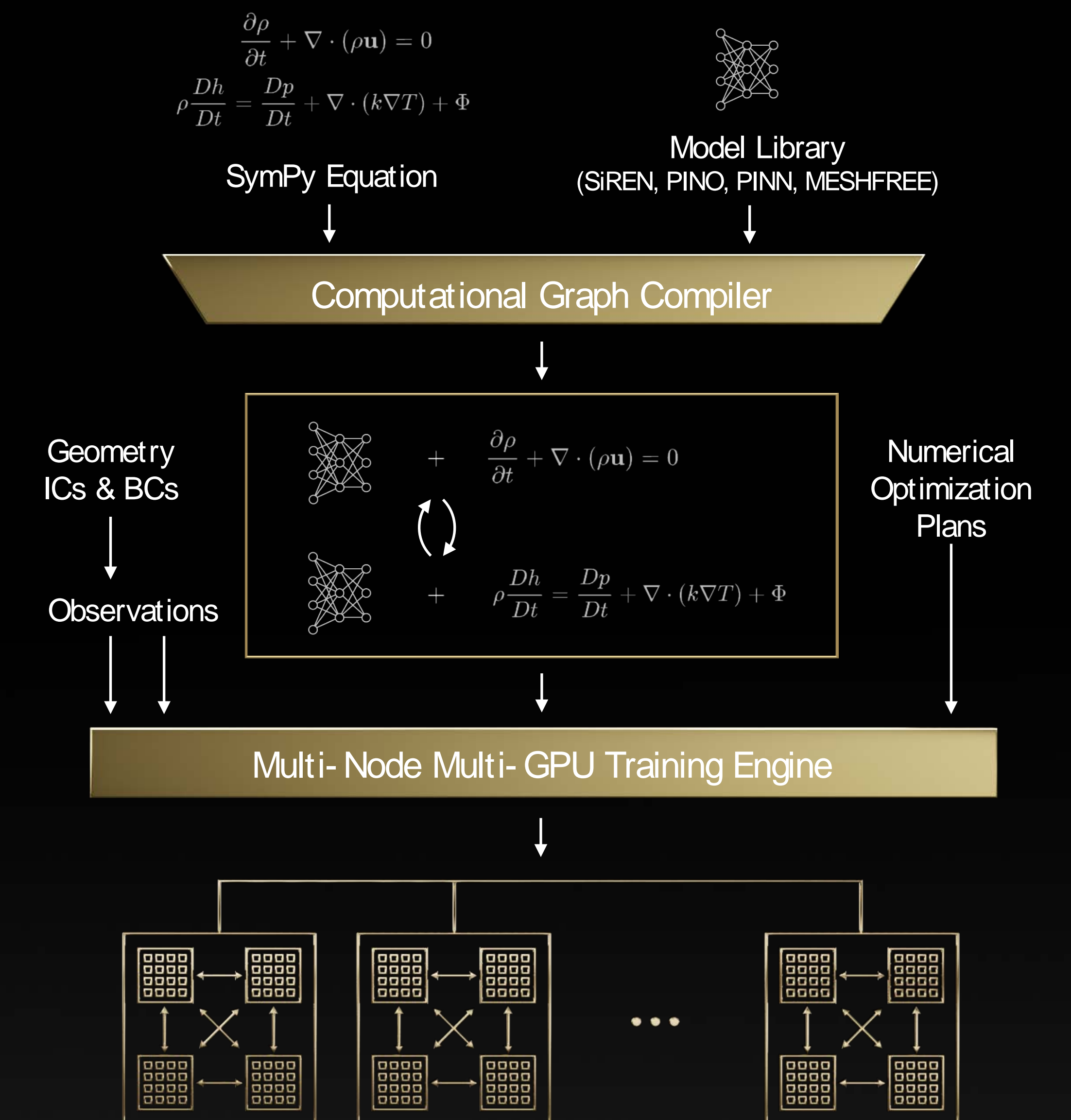
• cuNumeric



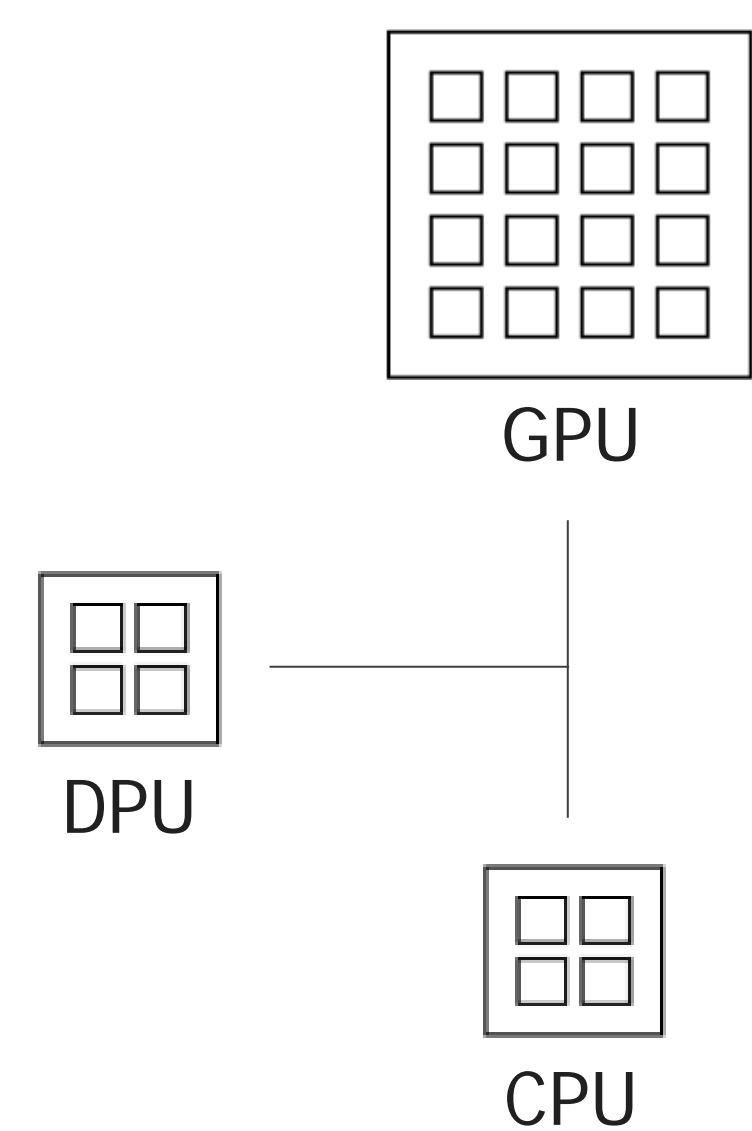
• cuQuantum



• Modulus



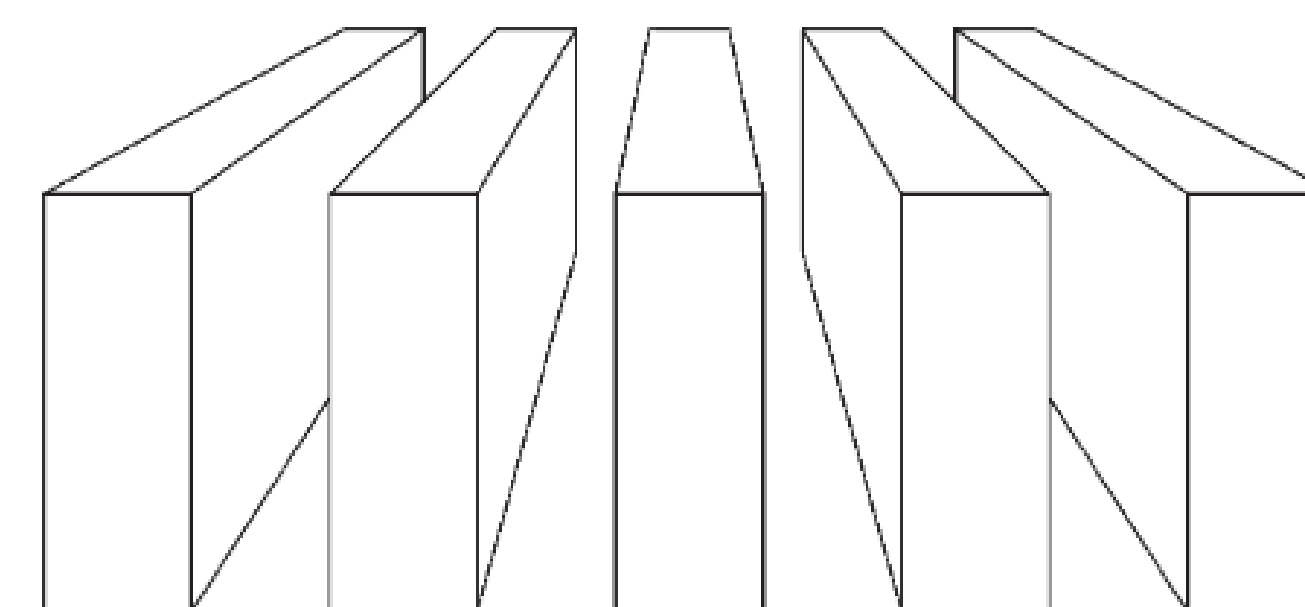
コンピューティング高速化の方針



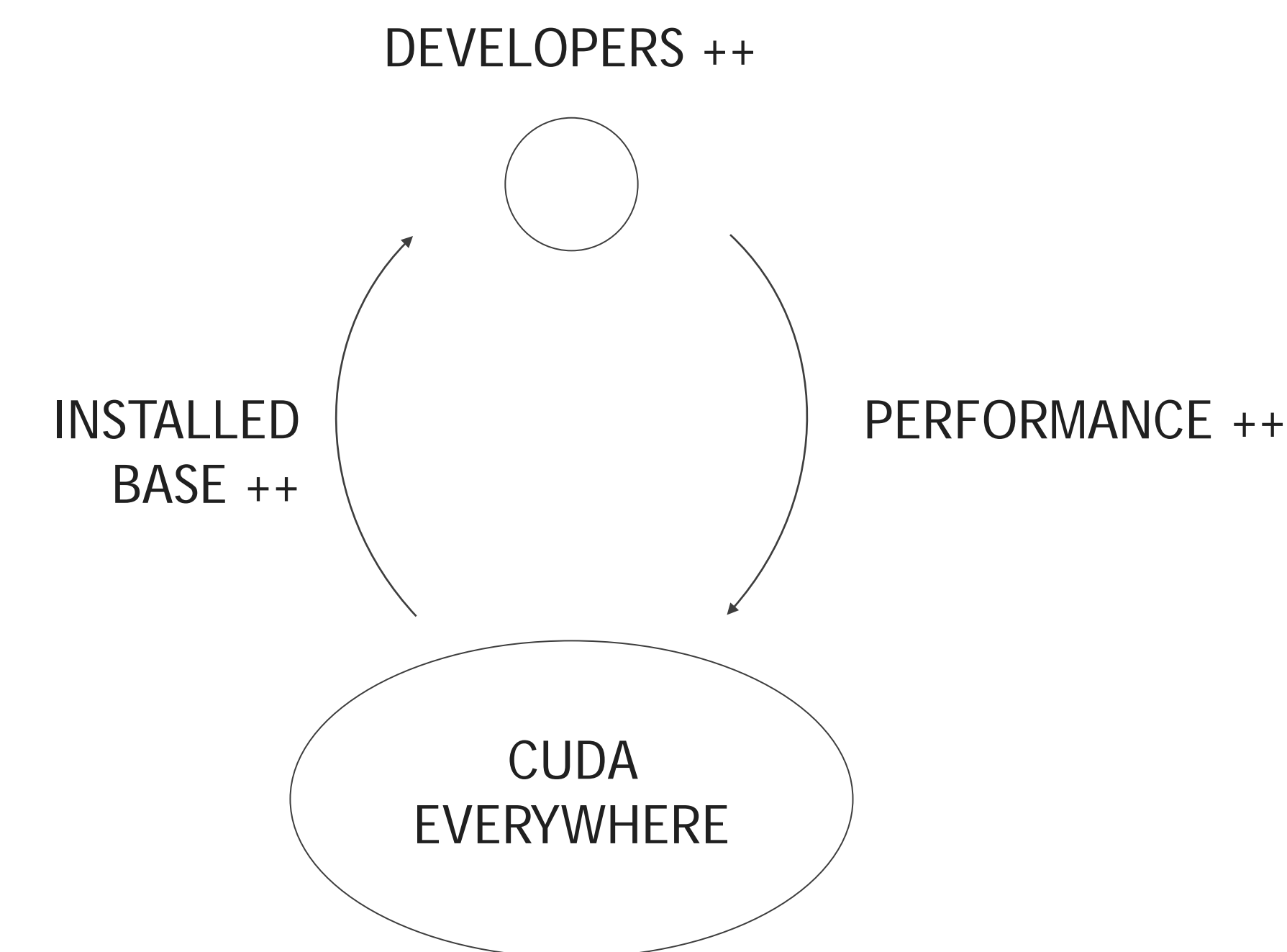
X-要素 高速化



フルスタック



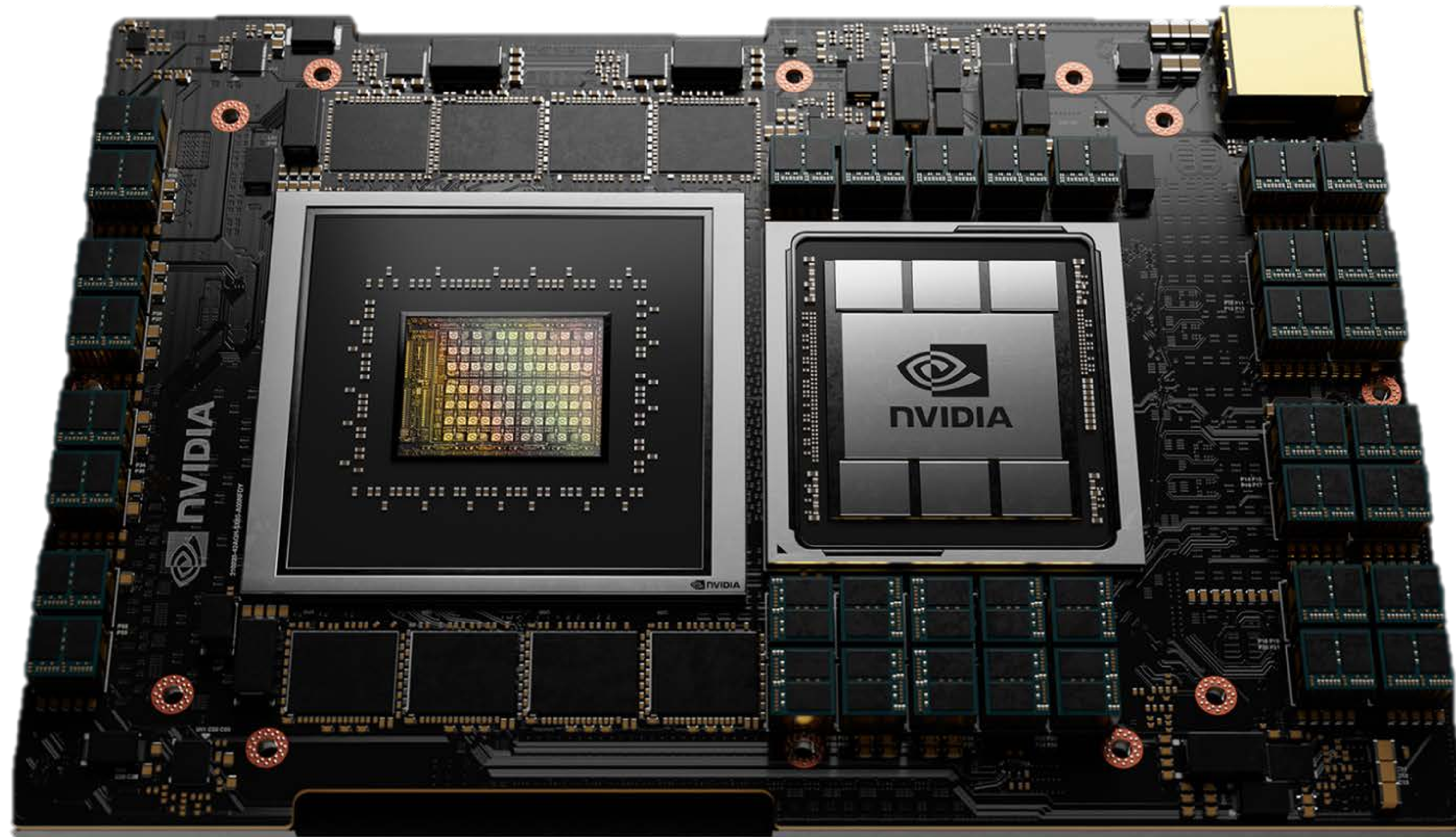
データセンタースケール



ワン・アーキテクチャー

NVIDIA GRACE

巨大スケールのAIおよびHPCアプリケーションに向け設計したブレークスルーCPU



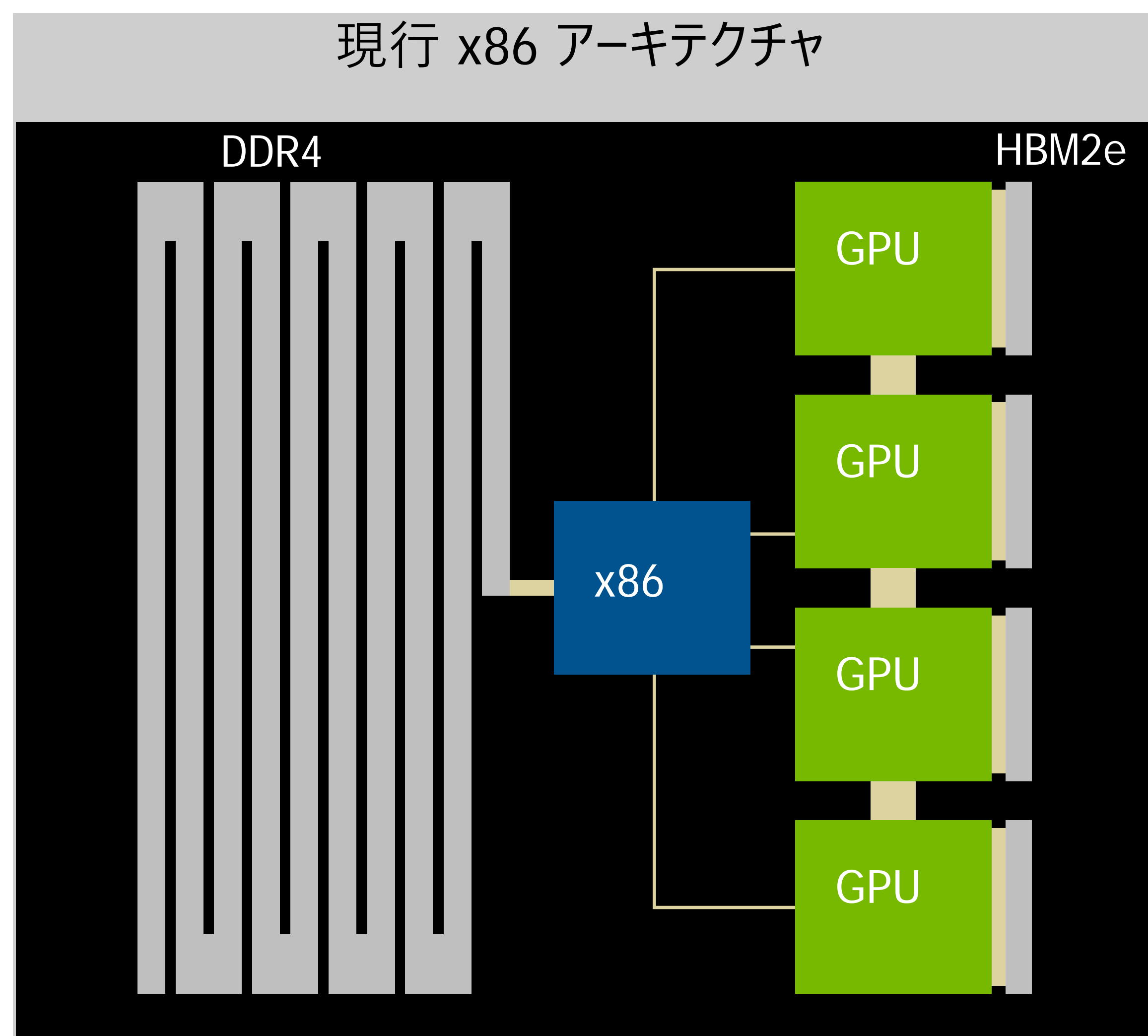
最高速の相互接続
>900 GB/s キャッシュコヒーレント NVLink
CPU-GPU間通信 (14倍)
>600GB/s CPU-CPU間 (2倍)

最高のメモリバンド幅
500GB/s LPDDR5x ECC付
2倍以上のバンド幅
10倍のエネルギー効率

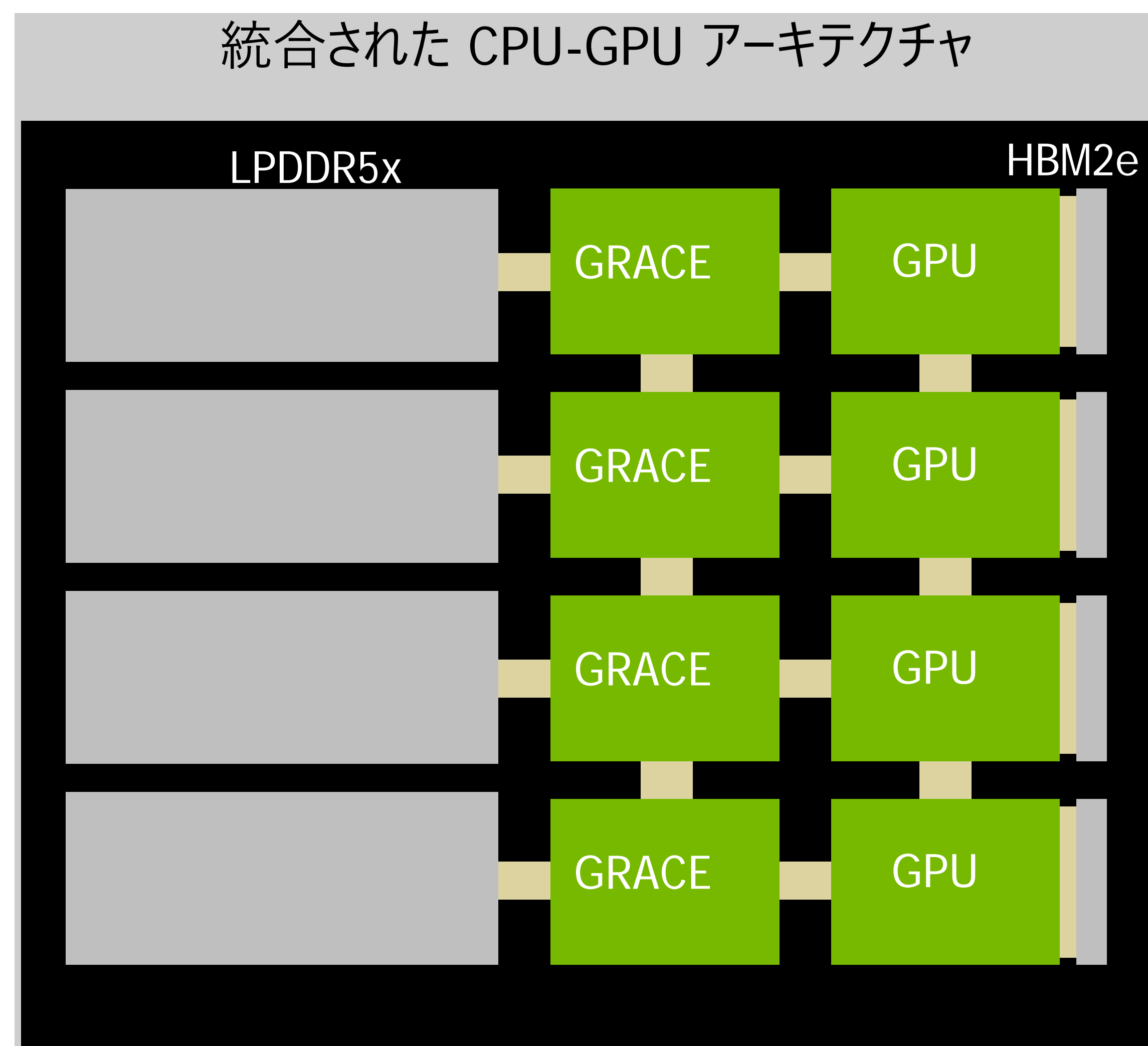
次世代 ARM NEOVERSE コア
>300 SPECrate2017_int_base est.
2023年出荷予定

テラバイトレベルに高速化されたコンピューティング

新たなワークロードに向けたアーキテクチャの革新



GPU	8,000	GB/sec
CPU	200	GB/sec
PCIe Gen4 (Effective Per GPU)	16	GB/sec
メモリ-GPU間	64	GB/sec
2TBの転送に30秒		



GPU	8,000	GB/sec
CPU	500	GB/sec
NVLink	500	GB/sec
メモリ-GPU間	2000	GB/sec
2TBの転送に1秒		

3 DAYS FROM 1 MONTH
1Tモデルの学習ファインチューン

0.5Tモデルでのリアルタイム推論
シングルノードでのインタラクティブなNLP推論

Bandwidth claims rounded to nearest hundred for illustration.
 Performance results based on projections on these configurations Grace : 8xGrace and 8xA100 with 4th Gen NVIDIA NVLink Connection between CPU and GPU and x86: DGX A100.
 Training: 1 Month of training is Fine-Tuning a 1T parameter model on a large custom data set on 64xGrace+64xA100 compared to 8xDGX A100 (16xX86+64xA100)
 Inference: 530B Parameter model on 8xGrace+8xA100 compared to DGXA100.

CPU/GPU/DPU のロードマップ

