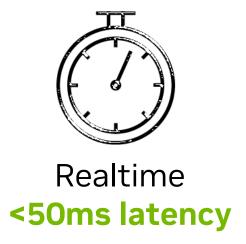


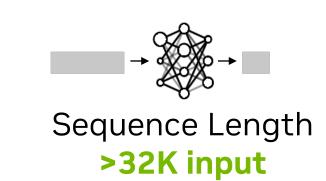
サステナブルなコンピューティングの ための取り組み

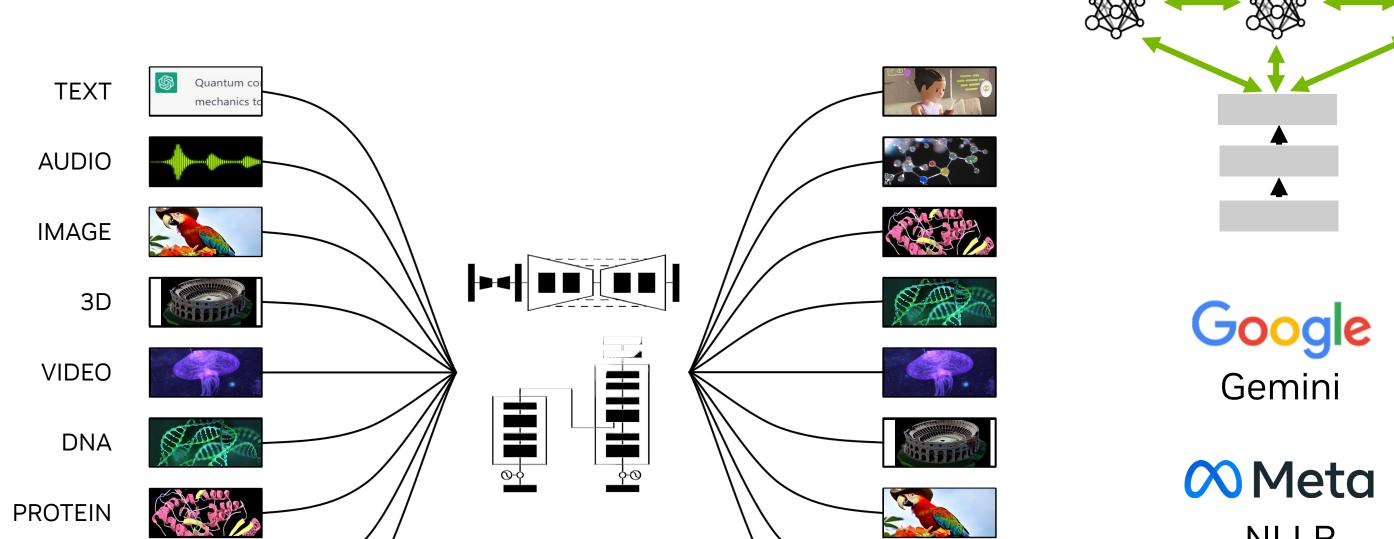
Dec 5,2024 – Hiroshi Aiko, Sr. Marketing Manager, NVIDIA.

生成AIの次の時代へ





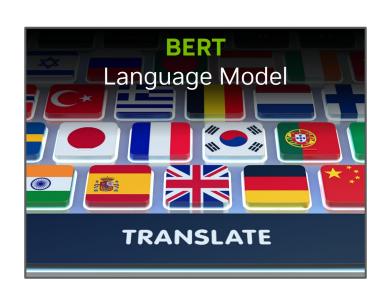






NLLB



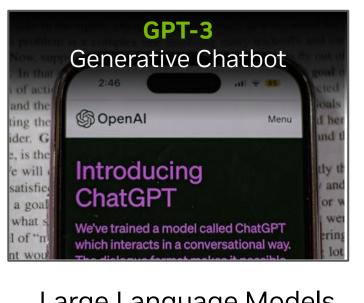


RESNET-50

Image Classification

Labeled Datasets





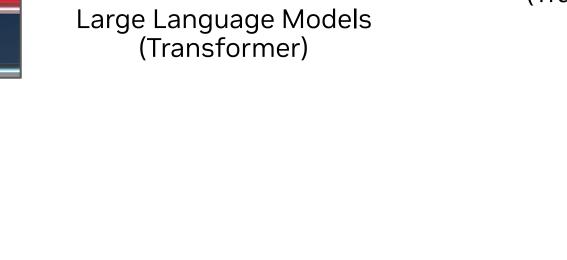
MOLECULE

ANIMATION

Large Language Models (Transformer)



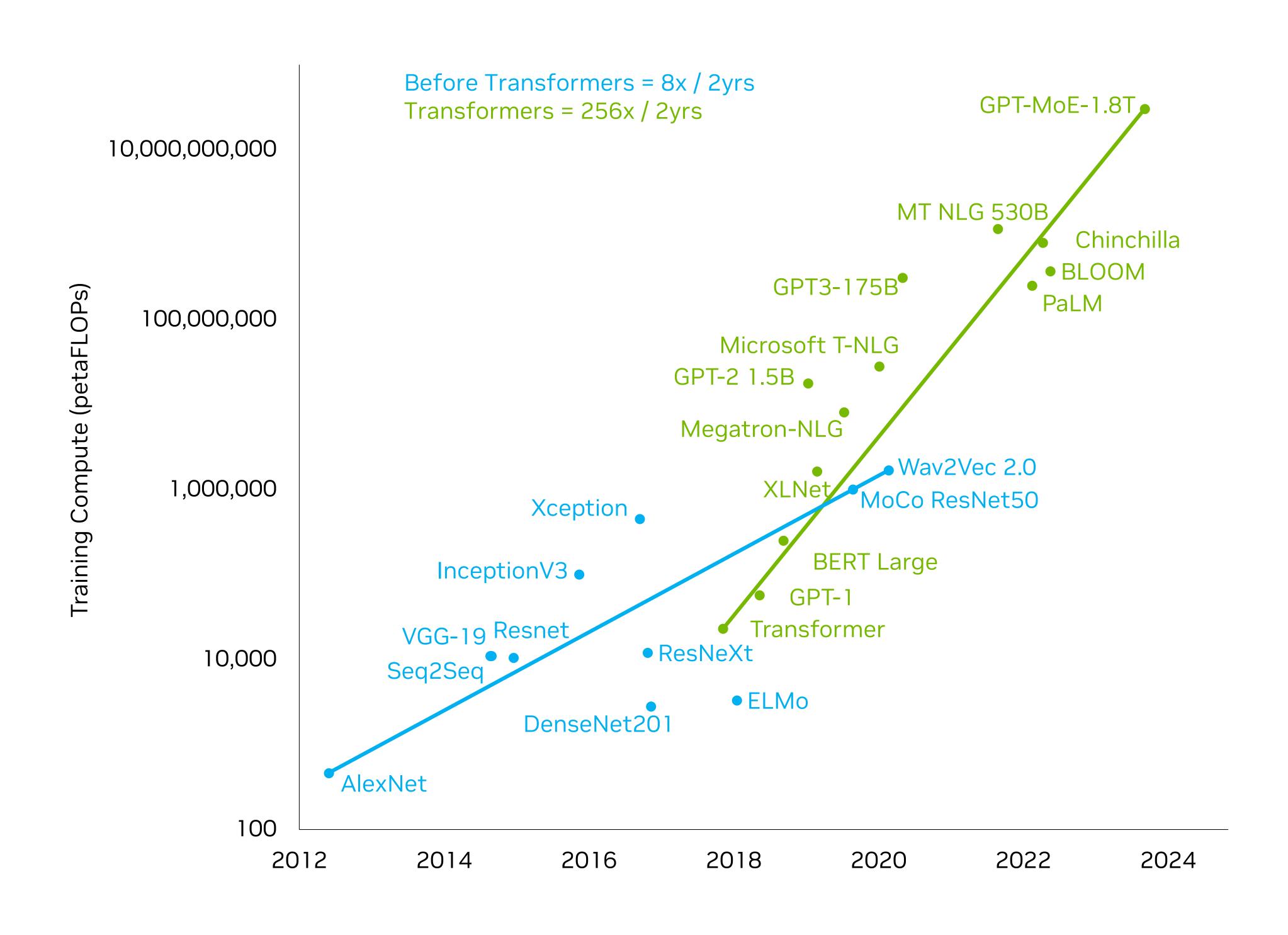
Unlabeled Datasets



Generative Al Multimodal Generative Al Mixture of Experts (MoE) Production GenAl Inference

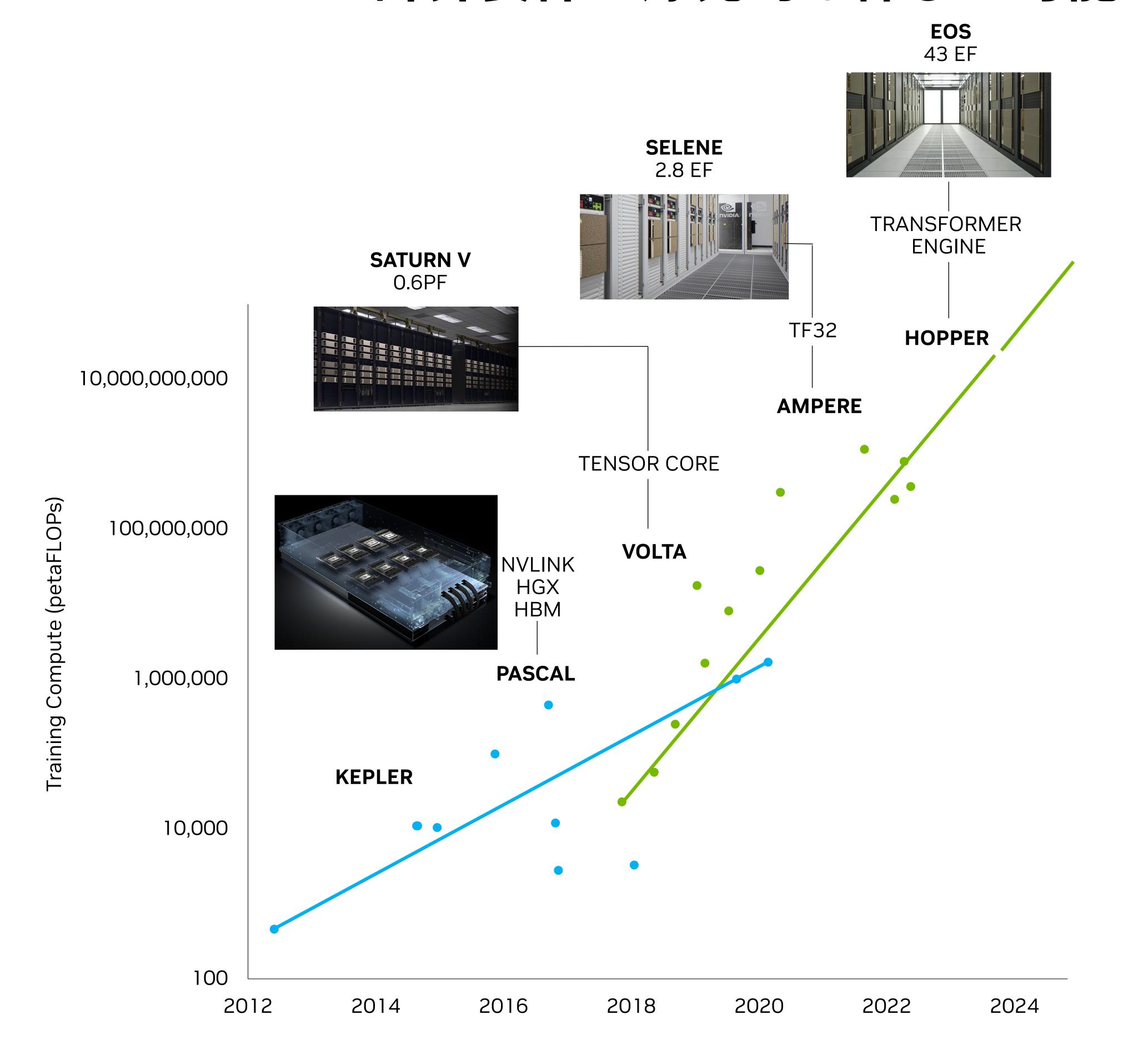
Quantum col

爆発的に増加するAIの計算要件





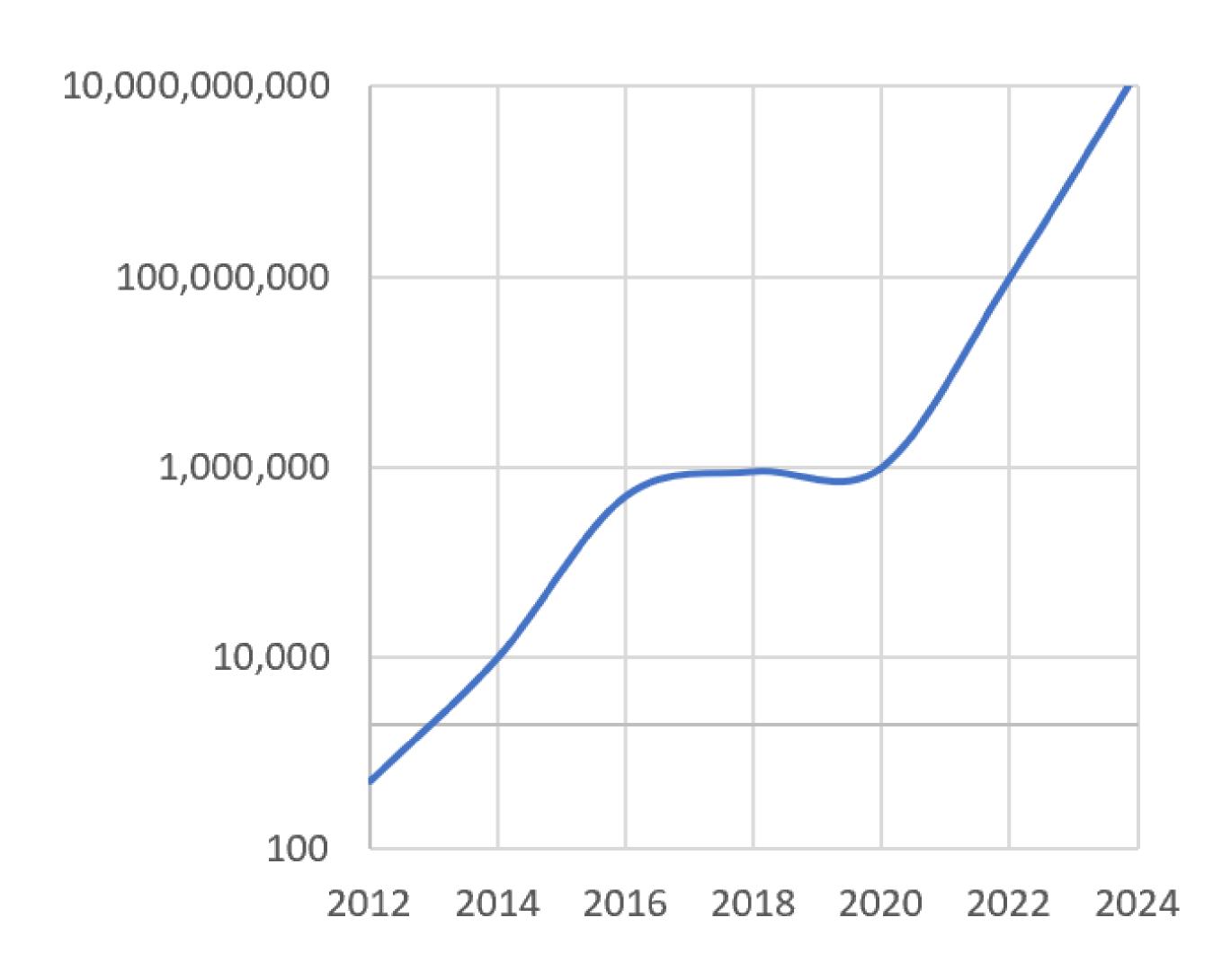
NVIDIA は AI 計算要件の爆発的な伸びを可能に





データセンターにおけるエネルギー使用量の爆発的増加

Transformer Models PetaFLOPS to Train



より多くのコンピューティングを要求する モデルサイズ



AIと HPC において GPU は CPU より大幅なエネルギー効率向上を実現

現在出荷中の NVIDIA GPU でデータセンターは今日からエネルギーを節約

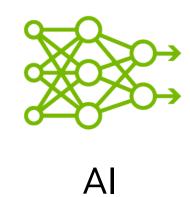


NVIDIA Blackwell GPUは、AI や HPC のワークロードにおいて、CPU の20 倍のエネルギー効率を発揮します。



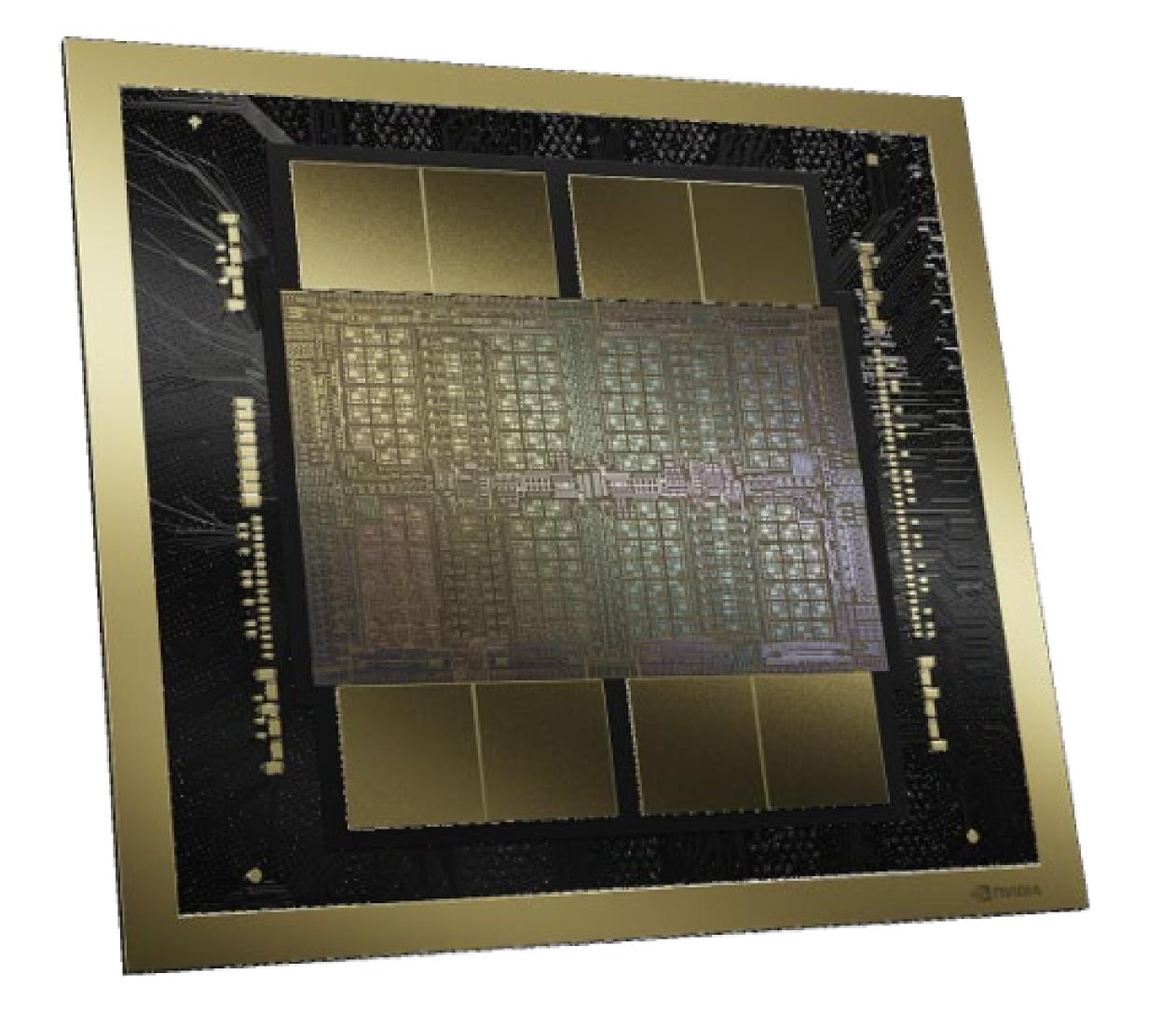
AI、HPC、データ分析のワークロードにおいて、CPU のみサーバーが、すべて GPU に移行した場合に節約できる**年間エネルギー量**

計算集約型用途













25M Metric Tons CO₂



に相当

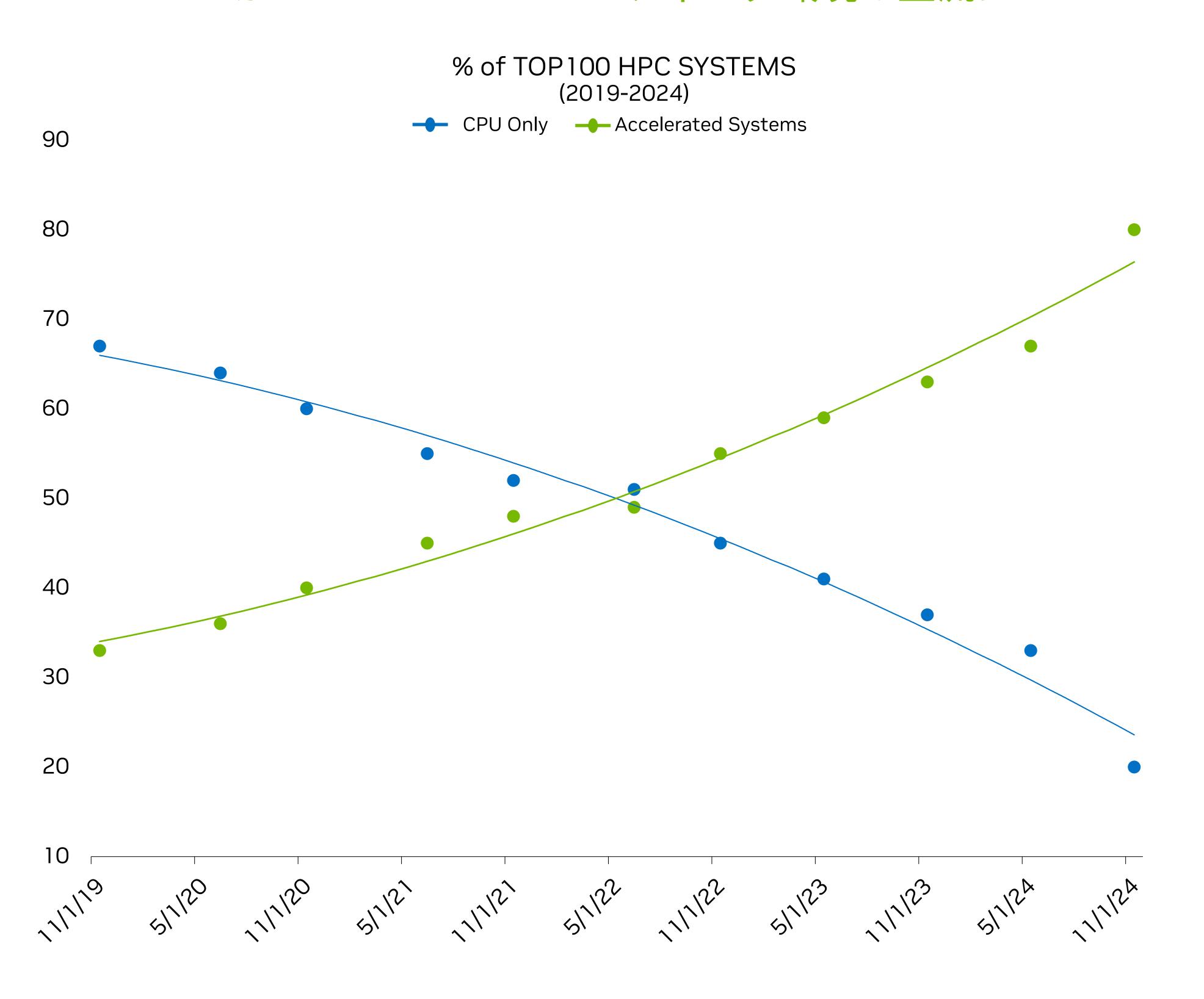
NVIDIA GB200 GPU の優れたエネルギー効率

AI、HPC、データ分析ワークロードの、CPU から GPU への移行による大幅なコスト削減



アクセラレーテッドコンピューティングが スーパーコンピューティングを変革

GPUがスーパーコンピューティング環境の主流に





Green500トップ10でのご採用

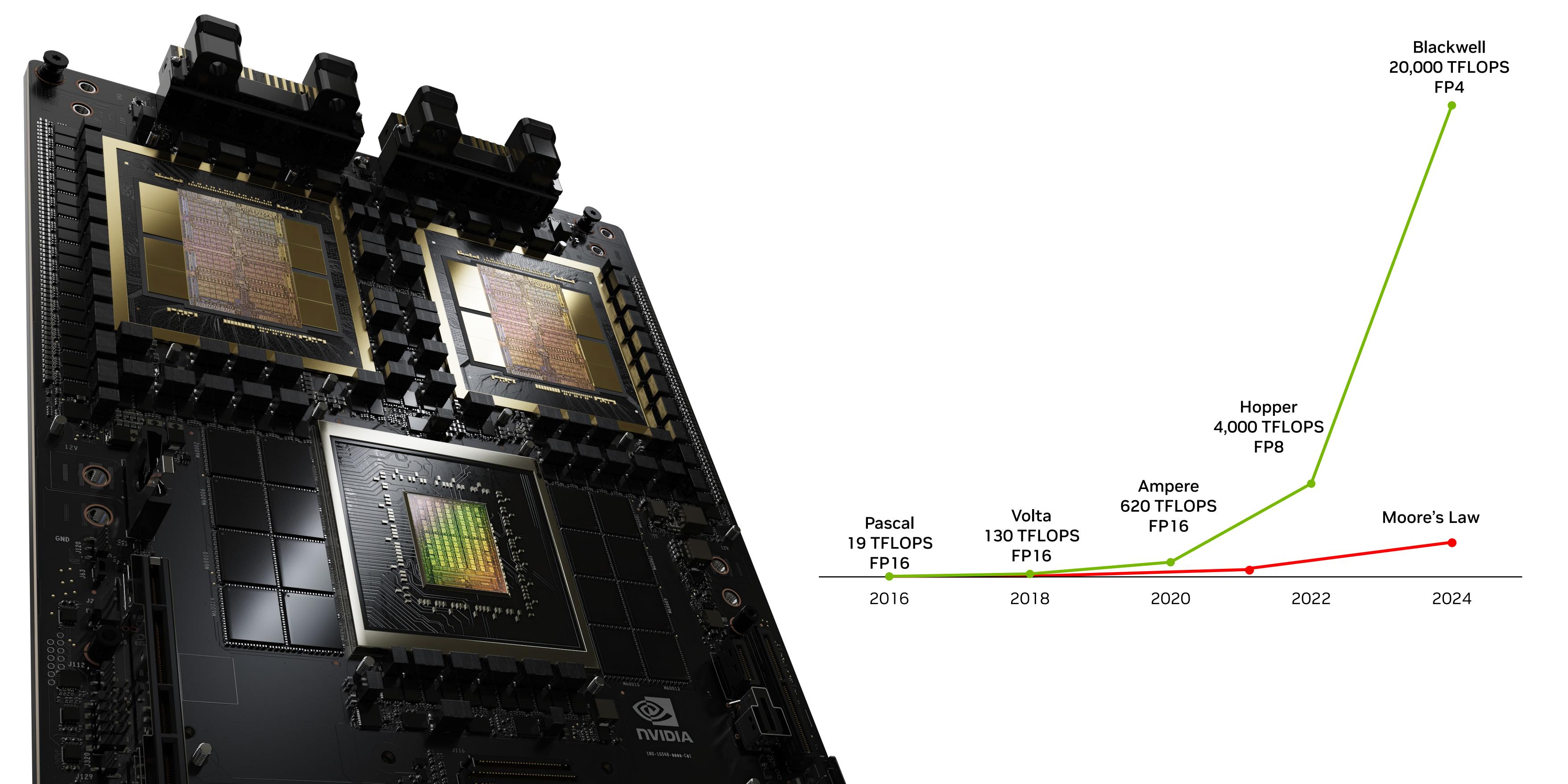
Green 500 上位 10 スパコン中 8 スパコンでのご採用 (2024年11月)

# Name	Computer
1JEDI	BullSequana XH3000, Grace Hopper Superchip 72C 3GHz, NVIDIA GH200 Superchip, Quad-Rail NVIDIA InfiniBand NDR200
2ROMEO-2025	BullSequana XH3000, Grace Hopper Superchip 72C 3GHz, NVIDIA GH200 Superchip , Quad-Rail NVIDIA InfiniBand NDR200, Red Hat Enterprise Linux
3Adastra 2	HPE Cray EX255a, AMD 4th Gen EPYC 24C 1.8GHz, AMD Instinct MI300A, Slingshot-11, RHEL
4Isambard-AI phase 1	HPE Cray EX254n, NVIDIA Grace 72C 3.1GHz, <u>NVIDIA GH200 Superchip</u> , Slingshot-11 Lenovo ThinkSystem SD665-N V3, AMD EPYC 9334 32C 2.7GHz, Nvidia H100 SXM5 94Gb , Infiniband
5 Capella	NDR200, AlmaLinux 9.4
JETI - JUPITER Exascale 6Transition Instrument	BullSequana XH3000, Grace Hopper Superchip 72C 3GHz, NVIDIA GH200 Superchip , Quad-Rail NVIDIA InfiniBand NDR200, RedHat Linux and Modular Operating System
7Helios GPU	HPE Cray EX254n, NVIDIA Grace 72C 3.1GHz, NVIDIA GH200 Superchip , Slingshot-11
8Henri	ThinkSystem SR670 V2, Intel Xeon Platinum 8362 32C 2.8GHz, NVIDIA H100 80GB PCIe, Infiniband HDR ThinkSystem SD665-N V3, AMD EPYC 9354 32C 3.25GHz, Nvidia H100 94Gb SXM5, Infiniband
9HoreKa-Teal	NDR200
10rzAdams	HPE Cray EX255a, AMD 4th Gen EPYC 24C 1.8GHz, AMD Instinct MI300A, Slingshot-11, TOSS



NVIDIA Blackwell が飛躍的な演算能力を実現

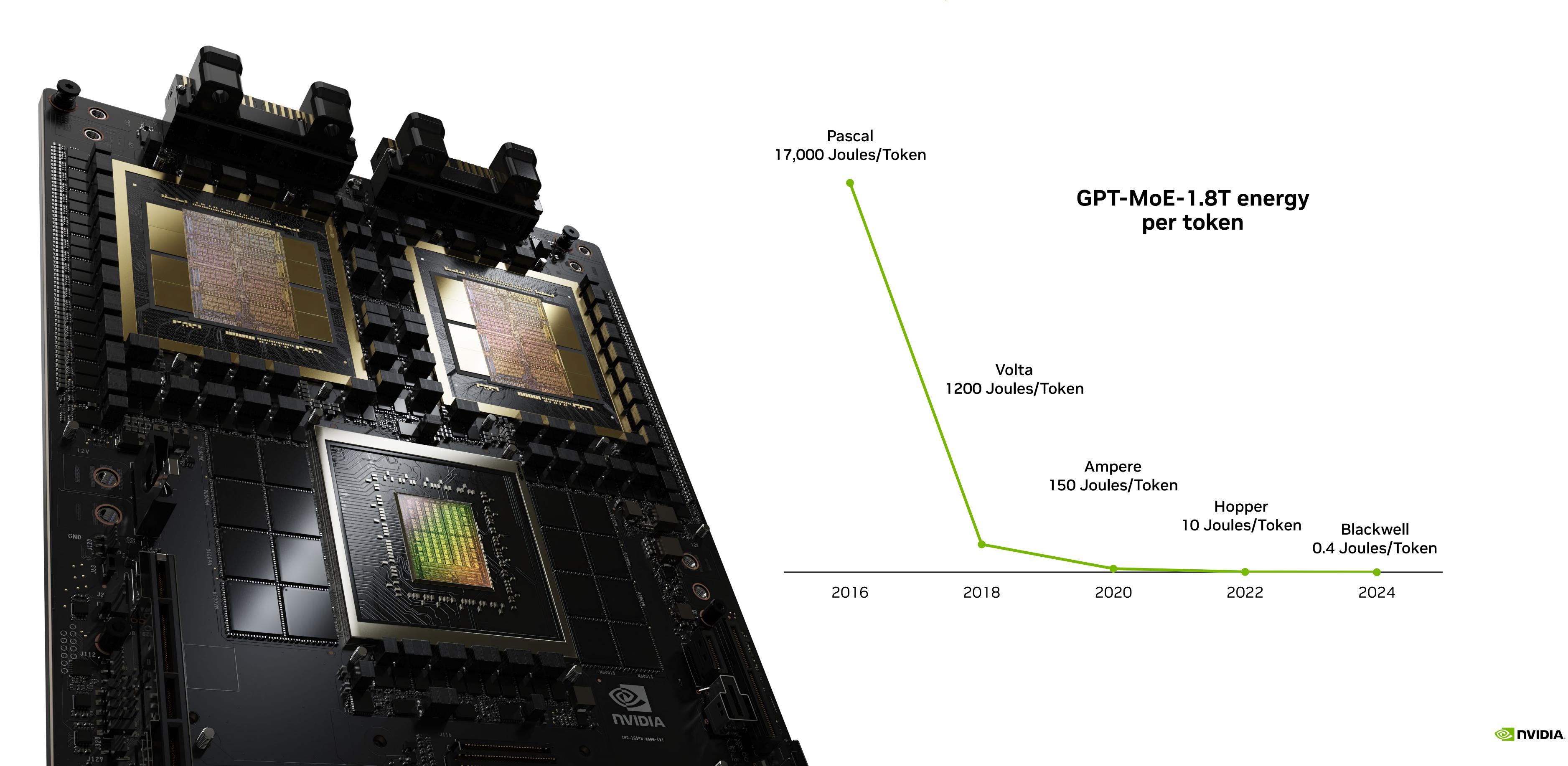
8年で1000倍の AI 演算能力





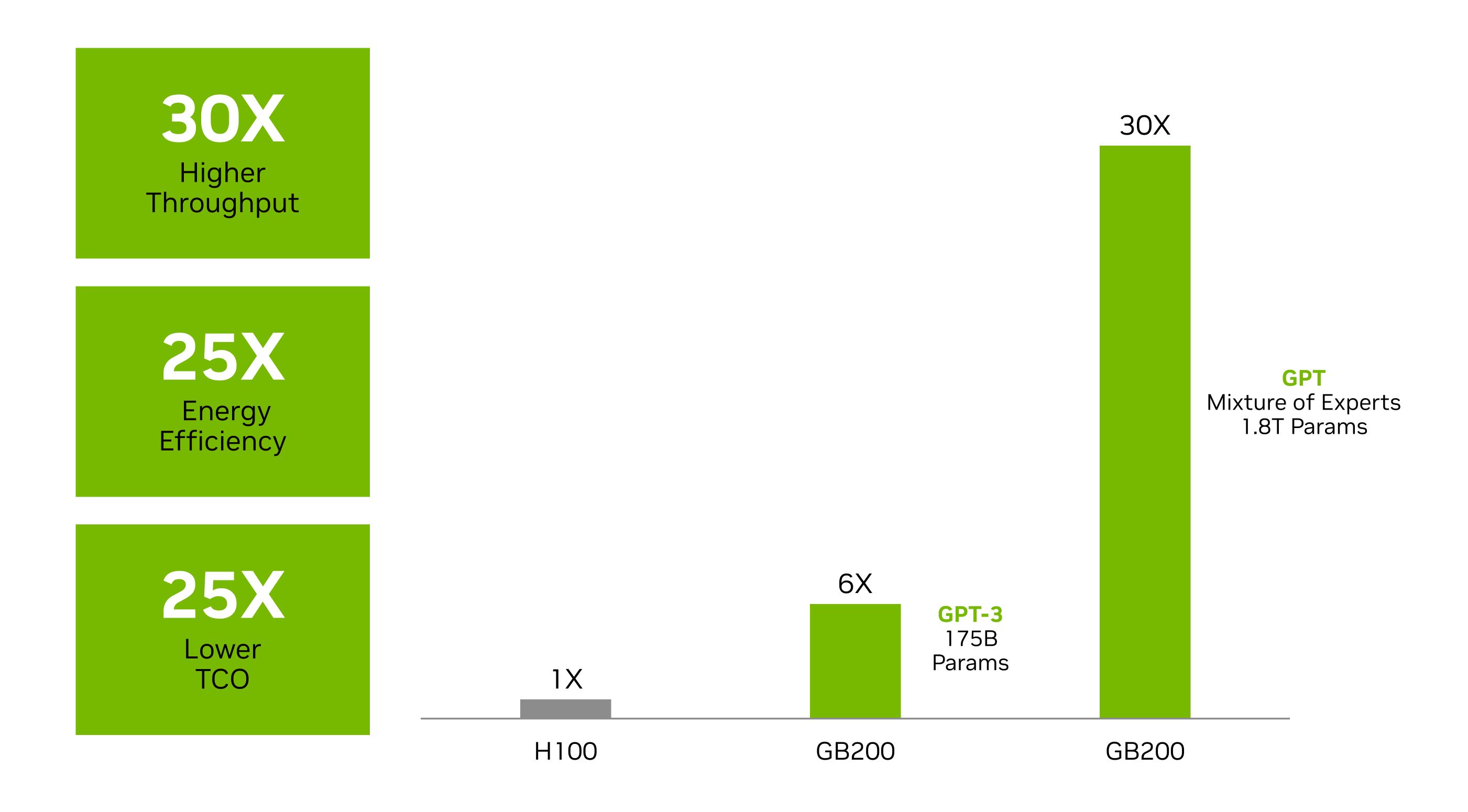
LLM推論のエネルギー効率は向上し続けている

トークンに必要なエネルギーは8年間で45,000分の1に減少



Blackwell は大規模モデルのパフォーマンスを向上

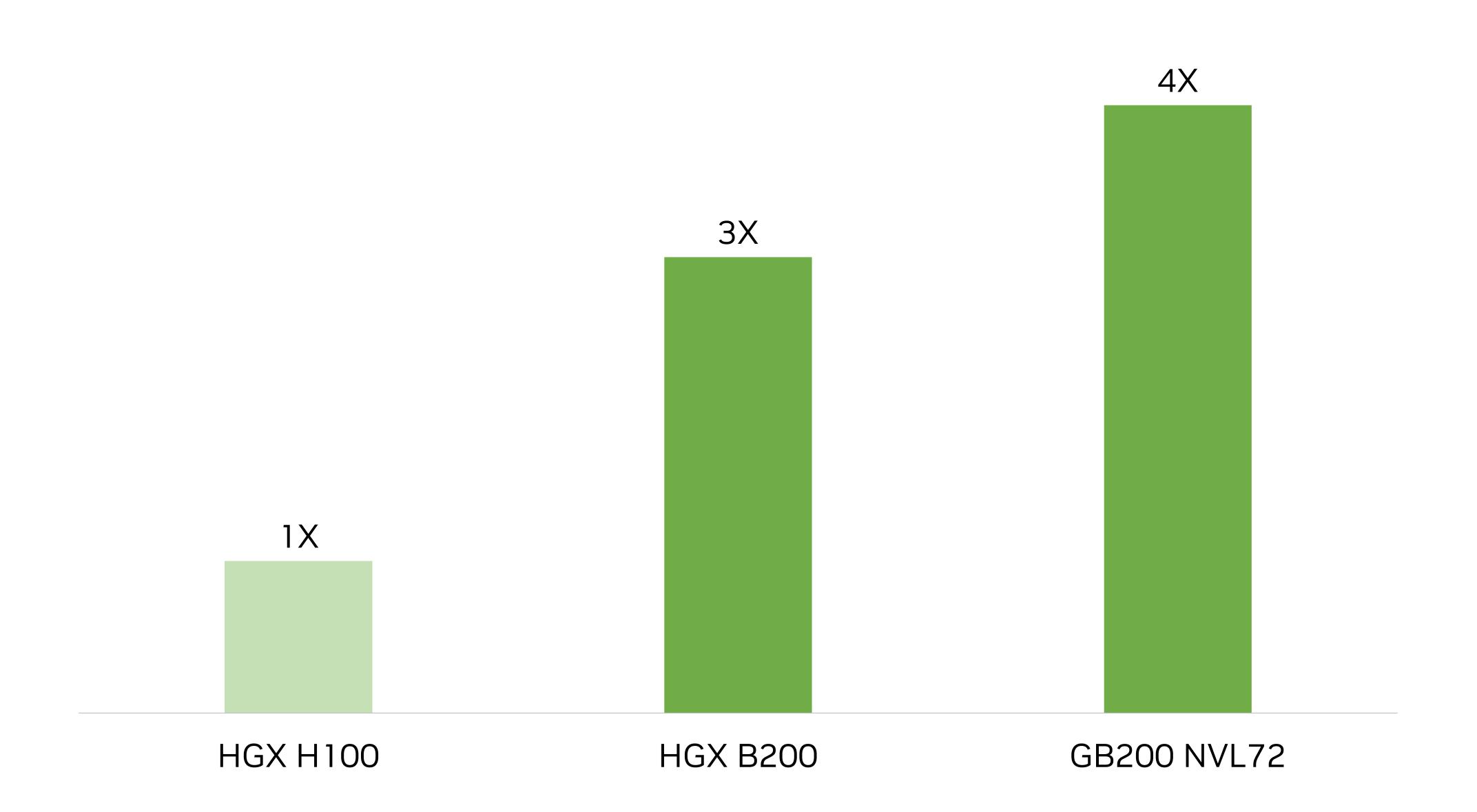
推論スループットとエネルギー効率は桁違いに



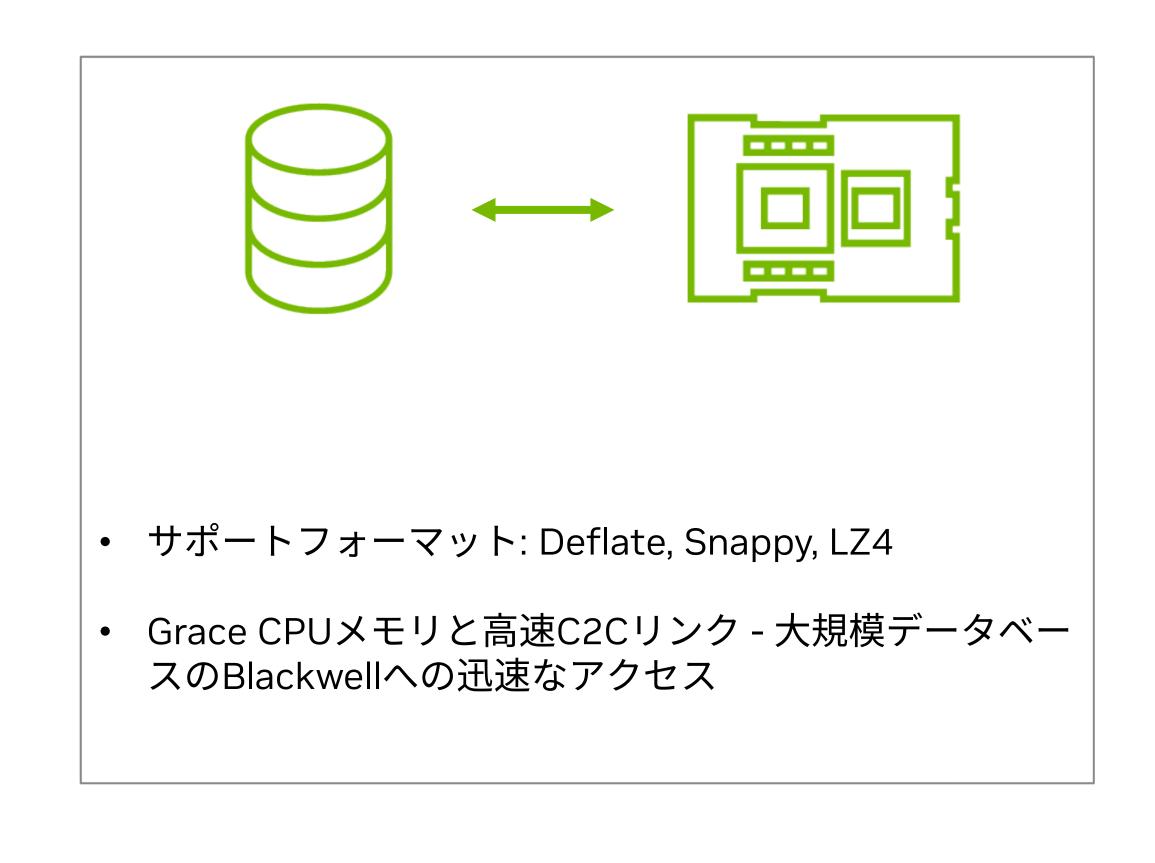


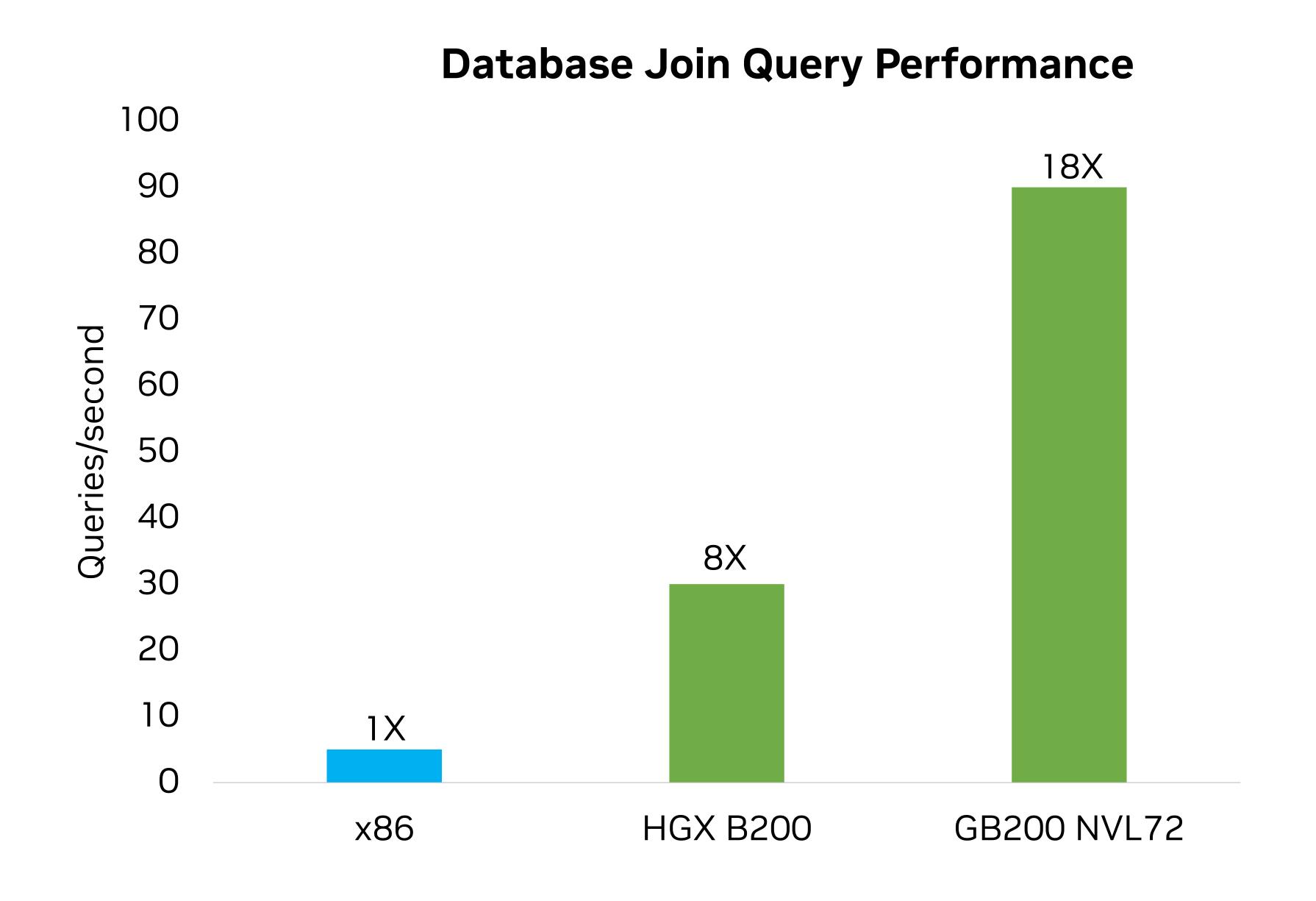
AIトレーニングのパフォーマンス向上

GPT-MoE-1.8T Model Training Speed-Up



Decompression Engineでデータ処理を高速化









DPU (Data Processing Unit) がデータセンターのエネルギー効率を高める

CPU から重要なインフラ処理をオフロードしてエネルギー消費を削減



インフラ オフロード









ストレージ セキュリティ ネットワーク

NVIDIA Bluefield DPU



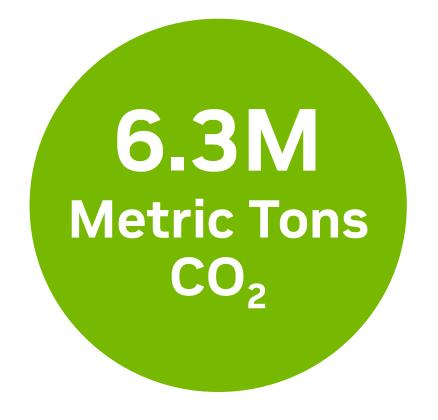


CPU による OVS インフラ処理を DPU に移行した場合のサーバ1台あたりのエネルギー削減量





すべての CPU インフラ運用を DPU に 移行した場合に節約できる年間エネル ギー量





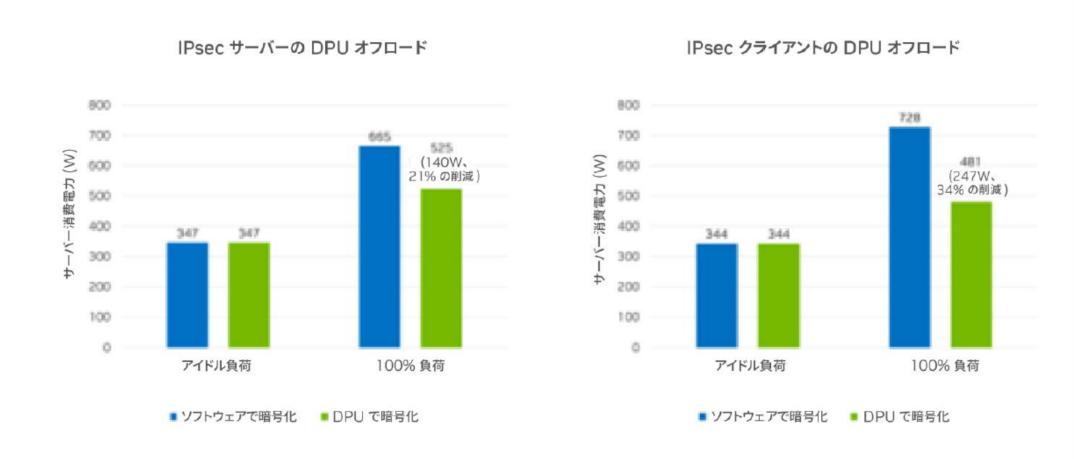
DPU による年間エネルギー節約額



サーバ1台あたりの DPU 節電効果

NVIDIA ホワイトペーパー「DPU による電力仕様の効率化」より

図 5 および 6: IPsec の暗号化、復号化の実証試験により、BlueField-2 ハードウェア へのオフロードが IPsec 全体をソフトウェアで処理する場合と比べて 顕著に省電力効果を示すことが判明。IPsec クライアントと IPsec サ ーバーの両方で顕著な省電力



予想どおり、サーバーの負荷が 0% (負荷なし = 暗号化トラフィックなし = オフロードするものなし) の場合に省電力効果は最小であり、負荷が 100% (大量の暗号化処理をオフロード) の場合に最大でした。つまり、DPU とサーバー効率化による最大の省電力効果は、各サーバーを可能な限り 100% に近い負荷で稼働させることで達成できます。この条件は、データセンターが一般的な戦略として打ち出す、必要なサーバーの台数を最小にし、ハードウェア アクセラレーションによるオフロードのメリットを最大にすることにも合致します。

表 3 IPsec の暗号化を CPU から BlueField DPU にオフロードすることで 得られた省電力効果⁶

100% 負荷時に IPsec を BlueField-2 にオフロード	サーバーあたりの電力使用量(削減)	サーバー 10,000 台を 3 年間運 用した場合の電力コスト (\$0.15/kWh で計算)
IPsec サーバー、ソフトウェアで暗 号化	665W	2,620 万ドル
IPsec サーバー、暗号化を DPU に オフロード	525W (140W、21% の削減)	2,070 万ドル (550 万ドルの削減)
IPsec クライアント、ソフトウェア で暗号化	728W	2,870 万ドル
IPsec クライアント、暗号化を DPU にオフロード	481W <u>(247W、34% の削減</u>)	1,900 万ドル (870 万ドルの削減)



持続可能なコンピューティングに向けて

- •近年のプロセッサ効率や性能の改善は性能要求に追いついていない
 - 結果として、データセンターやスパコンなどでの消費電力量は増加を続けている
- •非効率なコンピューティングによる無駄を削減
 - AI や HPC は、GPU をはじめとするアクセラレータの活用で何十倍ものエネルギー効率
 - AI 以外でも、暗号化や各種データ処理は専用のプロセッサ(DPU など)を活用

ハードウェアやソフトウェアを含むあらゆるレイヤーにおいて 最適化や効率化を図ることでエネルギー効率が向上し 持続可能なコンピューティングにつながる



