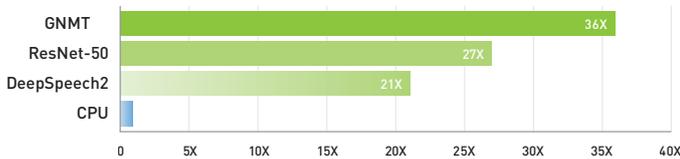


NVIDIA T4 TENSOR コア GPU

AIの大規模なトレーニングと推論を加速

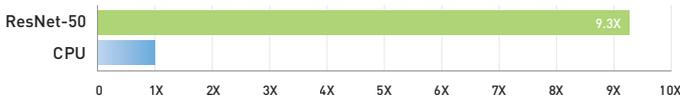
NVIDIA® T4 GPU は、スケールアウト環境のパフォーマンスを飛躍的に向上させる世界最高クラスのアクセラレーターであり、あらゆるサーバーを高速化できます。幅広い最新アプリケーションを加速できるよう、NVIDIA Turing™ Tensor コアを搭載し、革新的な多精度コンピューティングを実現しました。エネルギー効率に優れ、70ワットの小型 PCIe フォームファクターに収まったこの最先端の GPU は、スケールアウトサーバーに最適化され、新進の AI ワークロードのために設計されています。

推論性能



NVIDIA T4 GPU 1 基とデュアルソケット Xeon Gold 6140 CPU 搭載サーバーの比較

トレーニング性能



NVIDIA T4 GPU 2 基とデュアルソケット Xeon Gold 6140 CPU 搭載サーバーの比較



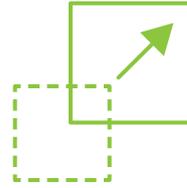
仕様

GPU アーキテクチャ	NVIDIA Turing
NVIDIA Turing Tensor コア	320
NVIDIA CUDA® コア	2,560
単精度演算性能	8.1 TFLOPS
混合精度演算性能 (FP16/FP32)	65 TFLOPS
INT8 精度演算性能	130 TOPS
INT4 精度演算性能	260 TOPS
GPU メモリ	16 GB GDDR6 300 GB/秒
ECC	対応
相互接続帯域幅	32 GB/秒
システム インターフェイス	x16 PCIe Gen3
フォームファクター	ロープロファイル PCIe
冷却方式	パッシブ
計算 API	CUDA、NVIDIA TensorRT™、ONNX

データセンター アクセラレーションをさらに強化する圧倒的パフォーマンス



小型のフォームファクターと70Wの電力設計によってT4はスケールアウトサーバー向けに最適化されており、CPUの50倍という驚異的なエネルギー効率が実現されるため、運用コストを劇的に削減できます。この2年間でNVIDIAの推論プラットフォームの効率性は10倍以上向上し、分散AIトレーニングおよび推論のための最もエネルギー効率に優れたソリューションであり続けています。



NVIDIA T4 データセンター GPUは、分散コンピューティング環境に理想的なユニバーサルアクセラレーターです。革新的な多精度コンピューティングにより、ディープラーニングや機械学習のトレーニングや推論、ビデオコード変換、仮想デスクトップを高速化します。あらゆる種類のAIフレームワークやネットワークがサポートされ、劇的にパフォーマンスと効率が向上されることで、大規模な環境を最大限に有効活用できるようになります。



Turing Tensor コア テクノロジーと多精度コンピューティングを採用したことで、FP32 から FP16、INT8、さらには INT4 の精度での演算に画期的なパフォーマンスが発揮され、さまざまなAIワークロードに対応できます。トレーニング性能はCPUと比べて最大9.3倍、推論性能は最大36倍に向上しています。

NVIDIA T4 の詳細については、www.nvidia.com/T4 をご覧ください。