

第5世代インテル® Xeon® スケーラブル・プロセッサとインテル® AIエンジンがAI処理全体のパフォーマンスを大幅に向上

65%

インテル® Xeon® プロセッサ上で稼働するデータセンターのAI推論の割合¹

最大

14倍向上

インテル® AMX BF16内蔵の第5世代インテル® Xeon® プロセッサ使用時のリアルタイム・オブジェクト検出(検知)推論パフォーマンス(SSD-ResNet34)を第3世代インテル® Xeon® プロセッサと比較した場合²

最大

9.9倍向上

インテル® AMX BF16内蔵第5世代インテル® Xeon® プロセッサ使用時のリアルタイム自然言語処理推論(BERT-large)パフォーマンスを第3世代インテル® Xeon® プロセッサと比較した場合。また、同じ比較条件でワット当たりの性能は最大7.7倍向上³

最大

8倍向上

第5世代インテル® Xeon® プロセッサ使用時のバッチのレコメンデーション・システムの推論パフォーマンス(DLRM)を、第3世代インテル® Xeon® プロセッサと比較した場合。また、同じ比較条件でワット当たりの性能は6.2倍向上⁴

AIのワークロードやユースケースは、データの前処理、昔からあるマシンラーニング(ML)をはじめ、自然言語処理や画像認識といったディープラーニングの用途まで広範囲に及びます。インテル® Xeon® スケーラブル・プロセッサを使用することで、このようなAI処理全体のパフォーマンスを大幅に向上させることが可能です。各プロセッサには、マシンラーニングやデータ分析、ディープラーニングなど、特定のAIワークロード向けに最適化されたアクセラレーターが内蔵されています。

あらゆる用途のAIに対応

AIは今日のビジネスに欠かせないものとなっており、影響範囲は多種多様かつ非常に重要なワークロードにおよびます。マシンラーニングやディープラーニングは、コア・エンタープライズ・アプリケーションから自動音声応答機能に至るまで、ビジネスの基本的な構成要素になりつつあります。AIを大規模に運用できるかどうかは、データの前処理からトレーニング、導入までの長期的な開発プロセス次第です。各ステップにはそれぞれ独自の開発ツールチェーン、フレームワーク、ワークロードがあり、そのすべてにおいて特有のボトルネックが発生し、それぞれ異なる要因でコンピューティング・リソースに負荷をかけます。インテル® Xeon® スケーラブル・プロセッサに内蔵されているアクセラレーターが、処理全体の実行環境を整え、AIパフォーマンスを全体的に向上させることが可能です。

要求が非常に厳しい新たなワークロードをサポートするための統合アクセラレーター、インテル® アクセラレーター・エンジン

第5世代インテル® Xeon® スケーラブル・プロセッサは、汎用コンピューティングに優れているだけでなく、日々重要性の高まるAIワークロードの多くをサポートする基盤としての役割を果たします。また、CPUのディープラーニング推論とトレーニングの高速化をはかるために設計されたAIアクセラレーター、インテル® アドバンスド・マトリクス・エクステンション(インテル® AMX)が内蔵されています。内蔵AIアクセラレーターのおかげで、ディスクリット・アクセラレーターにかかる追加コストや複雑性を回避できます。最新世代のインテル® Xeon® プロセッサは、パラメータ数が200億(20B)未満の大規模言語モデル(LLM)に適しており、ほとんどのお客様との合意サービス水準(SLA)を満たしています。⁵ インテル® AMXは、トランスファー・ラーニング(転移学習)やファインチューニングにも優れているため、ハードウェア追加の必要なしに、わずか4分でモデルをトレーニングできます。データセンターの推論の65%がインテル® Xeon® プロセッサで実行されています。そのためお客様は、GPUインフラストラクチャーに移行する際の複雑性に煩わされることなく、汎用AI向けの既存アーキテクチャーの恩恵を受けられます。

第5世代インテル® Xeon® スケーラブル・プロセッサとインテル® アクセラレーター・エンジンが導く未来のイノベーション

インテル® Xeon® プロセッサを使用しているなら、オンプレミス、クラウド、エッジなど実行環境を問わず、内蔵のインテル® アクセラレーター・エンジンで、ビジネスを新たな高みに導くことが可能です。データ保護の強化やインフラストラクチャーの活用など、さまざまな利点を得られます。インテル® アクセラレーター・エンジンは、仮想CPUと物理CPUの使用率を向上させ、コア・ソリューション・ライセンスのコストを削減します。すなわち、これらの内蔵アクセラレーターは、アプリケーションのパフォーマンス向上、コスト削減、プラットフォーム・レベルの効率向上を実現します。



採用事例: Intel® Xeon® スケーラブル・プロセッサによる現実世界のアクセラレーション

Tencent Cloud: Intel® Xeon® スケーラブル・プロセッサにより、リアルタイムの音声合成を実現。

[ストーリーを読む\(英語\)](#)

Gunpowder: 第4世代Intel® Xeon® スケーラブル・プロセッサ搭載 Google Cloud C3 インスタンスで、レンダリング・パフォーマンスの高速化を実現。

[ストーリーを読む](#)

Intel® アドバンスド・マトリクス・エクステンションによるディープラーニングの高速化

Intel® AMX は、第5世代Intel® Xeon® スケーラブル・プロセッサでディープラーニングのトレーニングや推論を行うためのIntelの最先端テクノロジーです。Intel® AMX は、自然言語処理、レコメンデーション・システム、画像認識のようなワークロードに最適で、AMX BF16 内蔵の第5世代Intel® Xeon® プロセッサは、第3世代と比較して、リアルタイム・オブジェクト分類推論パフォーマンスが最大7.2倍、ワット当たりの性能が5.3倍向上します。⁶

またIntel® AMX は、AI モデルのワークロードを向上させるだけでなく、多くのお客様がすでに運用しているプラットフォームのSLAを満たします。第5世代Intel® Xeon® スケーラブル・プロセッサは、5段階のターボレシオを追加し、ハイパフォーマンス・コンピューティングやAIなどのベクトル演算、マトリクス演算に親和性のあるワークロードに対してターボ周波数を向上させることが可能です。

Intel® AMX は、CPU コアのIntel® アドバンスド・ベクトル・エクステンション512 (Intel® AVX-512) と比較すると、スループット (Ops/Cycle) が高く、行列乗算のパフォーマンスが向上しています。⁷ そのため、ディープラーニングのトレーニングのワークロードがより迅速に完了するだけでなく、多くのお客様がすでにビジネスで運用しているプラットフォームのSLAを満たすことが可能です。

自然言語処理と生成AIをサポート

Intel® AMX 内蔵の第5世代Intel® Xeon® スケーラブル・プロセッサは、ハードウェアの追加を必要することなく、自然言語処理のパフォーマンスが大幅に向上します。Intelの各種ライブラリーはTensorFlowやPyTorch向けに最適化されているため、開発者は、内蔵AIアクセラレーターの利点をすぐに利用できます。開発者は、異なるハードウェア環境からコードを移植するという、時間とコストがかかる可能性のあるプロセスを、より簡単に行えるようになります。

Intel® AMX 内蔵の第5世代Intel® Xeon® スケーラブル・プロセッサは、ディープラーニングの推論とトレーニングを高速化することで、総保有コスト (TCO) の均衡を保ちつつSLAを満たします。リアルタイムのユーザー行動シグナルに加え、時間や場所のようなコンテキスト (状況) を考慮するディープラーニング・ベースのレコメンダー・システムで活用されています。

第5世代Intel® Xeon® プロセッサは、人間中心のコンテンツを模倣する生成AIモデルの実行も可能で、大規模言語モデルやText-to-image生成 (テキストから画像への生成) をサポートします。生成AIにより特化したタスクについては、専用のIntel® Gaudi® AIアクセラレーター、Intel® データセンターGPU、その他のハードウェア・コンポーネントを使用して、CPUの機能拡張が可能です。

Intel® AVX-512でMLを高速化

Intel® Xeon® プロセッサは、ウェブサイト向けSSL暗号のハッシュ化、巨大なデータベースの処理、そして製薬研究やチップ設計、F1エンジンのシミュレーションを可能にします。

数世代にわたり改良されてきたIntel® AVX-512は、Intel® Xeon® スケーラブル・プロセッサの各クロックサイクルにおいてより多くのオペレーションを可能にするだけでなく、並列処理アプリケーションのパフォーマンス向上も実現しています。Intel® AVX-512の命令セット・アーキテクチャー (ISA) には、AI、HPC、ネットワーキング、ストレージといった多様なワークロードのパフォーマンス強化のために構築された拡張機能を備えています。

新世代のプロセッサのターボ・パフォーマンスについては、ターボレシオが4段階から5段階に増えています。そのため、Intel® AMXやIntel® AVX-512を活用すると、特定のHPCやAIワークロードのターボ周波数が向上します。

少ないステップで処理を高速化

第5世代Intel® Xeon® スケーラブル・プロセッサに内蔵されているIntel® AVX-512は、スマートで美しい数学を多用して、通常のコンピューティング作業を凝縮・結合・融合し、ステップの数を減らします。基本要素の例として、5回のクロックサイクルを要する「 $3 \times 3 \times 3 \times 3 \times 3$ 」の演算をCPUに命令するとします。あるいはCPUが1回のサイクルで処理できるよう、 3^5 (3の5乗)の命令にすることもできます。Intel® AVX-512は、このようなロジックを採用し、AIで最も困難な作業を含む数百ものワークロードに特化した作業に応用しています。

1刻みよりも遥かに速い8刻みを採用

Intel® AVX-512の「512」とは、命令が1回のクロックサイクルにおけるCPU処理のビット数を増やすという第2の方法を意味しています。40年前の16ビットPCの登場は非常に印象的でしたが、その後すぐに32ビットマシンに取って代わりました。今日のスマートフォンは64ビットで動作します。ビット数は、CPUが1回のクロックサイクルごとにアドレス指定できるレジスタの数 (CPUがデータを保持するメモリスロット) を示しています。名前が示すように、Intel® AVX-512はレジスタ数が512ビットに拡張されています。アプリケーションがIntel® AVX-512を活用する際、レジスタの数を単に拡張するだけで、CPUが基本とする64ビットの速度よりも最大で8倍高速に動作させられることとなります。96まで数える際に、8、16、24...と数えるか、1、2、3...と数えるかの違いです。

より少ない電力でより強力なAIを実行するエンジン

Intel® AI エンジン内蔵のIntel® Xeon® スケーラブル・プロセッサは、必要なハードウェア・リソースが少なく済むため、より強力な電力効率に優れたソリューションでAIワークロードの実行が可能になります。

アクセラレーター・エンジンを内蔵したIntel® Xeon® スケーラブル・プロセッサは、TCOの削減や、要求の厳しいAIワークロードの投資収益率 (ROI) の向上など、ワークロードの結果改善に貢献しています。

Intel® Xeon® プロセッサにより、実質的に自動で動く高速AI

Intel® Xeon® スケーラブル・プロセッサのAIアクセラレーターは、CPUの命令セット・アーキテクチャー (ISA) に組み込まれています。これは、利用可能なソフトウェアが準備されていて、それらを活用できることを意味しています。Intelのソフトウェア・エンジニアは、オープンソースのAIツールチェーンを常に最適化し、成果物をコミュニティに還元しています。例えばTensorFlow 2.9には、最適化されたIntel® oneAPI ディープ・ニューラル・ネットワーク・ライブラリー (Intel® oneDNN) が最初から付属しています。最新版のTensorFlowをダウンロードすると、自動的にIntel® Optimization for TensorFlowが使用可能です。

データサイエンティストや開発者は、AI処理に使用されるアプリケーションのために、オープンソースで無料のIntel® ディストリビューションやライブラリー、開発環境をダウンロードして、Intel® Xeon® スケーラブル・プロセッサ用のISAに内蔵されているすべてのアクセラレーターを利用できます。データサイエンティストやAIの開発者は、ツールをコーディングし直し、Intel® AVX-512用に再コンパイルする必要がなくなります。なぜなら、それらはすでにIntel® AVX-512用に最適化されているからです。

今日の企業や組織は、より高いワークロード・パフォーマンスをインフラストラクチャーから引き出し、電力効率の向上とコストの削減を行う必要があります。Intel® Xeon® スケーラブル・プロセッサに統合された専用のIntel® AIエンジンは、ビジネスにとって最も重要なAIワークロードを最大限活用するのに役立ちます。

ビジネスにとって最も重要なAIワークロードのために、Intel® アクセラレーター・エンジン内蔵のIntel® Xeon® スケーラブル・プロセッサを通じて、どんなことを実現できるかについては、詳細情報も参照してください。

詳細情報

[Intel® Xeon® スケーラブル・プロセッサ上のAIとディープラーニング](#) ›

[Intel® アドバンスド・ベクトル・エクステンション512 \(Intel® AVX-512\)](#) ›

[Intel® AI アナリティクス・ツールキット \(英語\)](#) ›

[Intel® ハードウェア上およびソフトウェア上での開発 \(英語\)](#) ›

AI やマシンラーニングのために最適化されたIntel® Xeon® プロセッサを活用して、クラウドや自社のインフラストラクチャーでのAI ワークロードのアクセラレーションを始めてみませんか。

[詳細情報 \(英語\)](#) ›



- 2022年12月時点でのAI 推論ワークロードを実行するデータセンター・サーバーの、Intelによる世界のインストール・ベース市場モデリングに基づく。
- [intel.com/processorclaims/ \(英語\)](#) : 5th Gen Intel Xeon® Scalable processors (A21) を参照してください。実際の結果は異なる場合があります。
- [intel.com/processorclaims/ \(英語\)](#) : 5th Gen Intel Xeon® Scalable processors (A19) を参照してください。実際の結果は異なる場合があります。
- [intel.com/processorclaims/ \(英語\)](#) : 5th Gen Intel Xeon® Scalable processors (A20) を参照してください。実際の結果は異なる場合があります。
- 2023年12月時点でのIntel社内モデリングに基づく。
- [intel.com/processorclaims/ \(英語\)](#) : 5th Gen Intel Xeon® Scalable processors (A22) を参照してください。実際の結果は異なる場合があります。
- [https://edc.intel.com/content/www/us/en/products/performance/benchmarks/vision-2022/ \(英語\)](https://edc.intel.com/content/www/us/en/products/performance/benchmarks/vision-2022/) のSession Benchmark #41と#42を参照してください。実際の結果は異なる場合があります。

通知と免責事項

性能は、使用状況、構成、その他の要因によって異なります。詳細については、[Performance Indexのウェブサイト](#)を参照してください。

性能の測定結果は構成情報に記載された日付時点のテストに基づくものです。また、公開中のすべてのアップデートが適用されているとは限りません。構成の詳細については、補足資料を参照してください。絶対的なセキュリティを提供できる製品またはコンポーネントはありません。

実際のコストと結果は異なる場合があります。

ワークロードおよび構成については、第5世代Intel® Xeon® スケーラブル・プロセッサ [www.intel.com/processorclaims/ \(英語\)](#) を参照してください。実際の結果は異なる場合があります。

Intelのテクノロジーを使用するには、対応したハードウェア、ソフトウェア、またはサービスの有効化が必要となる場合があります。

Intelは、サードパーティーのデータについて管理や監査を行っていません。ほかの情報を参考にして、正確さを評価してください。

アクセラレーターの有無は、SKUによって異なります。製品に関する追加情報に関しては、[Intel製品の仕様情報ページ](#)にアクセスしてください。

Intel® アドバンスド・ベクトル・エクステンション (Intel® AVX) は、特定のプロセッサ演算に高スループットを提供します。処理能力特性の違いにより、AVX 命令を使用しても、a) 一部の演算が定格周波数より遅くなる場合があります。b) Intel® ターボ・ブースト・テクノロジー2.0による一部の演算では、任意のターボ周波数や最大ターボ周波数に達しない場合があります。性能は、ハードウェア、ソフトウェア、システム構成により異なります。詳細については、<https://www.intel.co.jp/content/www/jp/ja/architecture-and-technology/turbo-boost/intel-turbo-boost-technology.html> を参照してください。

Intelは人権を尊重し、人権侵害の発生を回避するように尽力しています。詳しくは[Intelの世界的人権主義](#)をご覧ください。Intelの製品とソフトウェアは、国際的に認められている人権を侵害しない、または侵害の原因とならないアプリケーションに使用することを目的としています。

© Intel Corporation. Intel, インテル, Intel ロゴ, その他のIntelの名称やロゴは、Intel Corporationまたはその子会社の商標です。その他の社名、製品名などは、一般に各社の表示、商標または登録商標です。