

## インテル® アドバンスト・マトリクス・エクステンション (インテル® AMX) による AI ワークロードの高速化

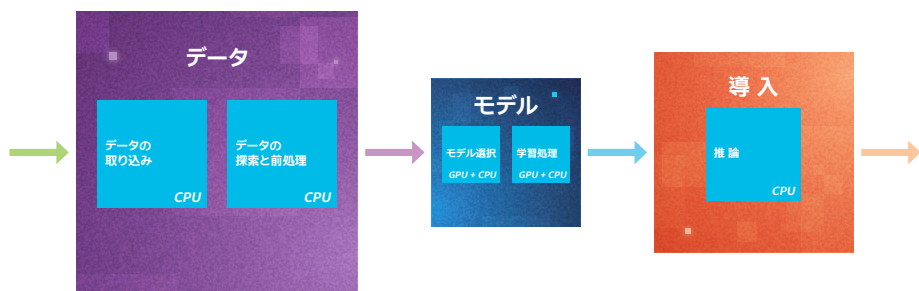
第4世代インテル® Xeon® スケーラブル・プロセッサとインテル® AMXによりAI機能が進化し、前世代と比較して推論と学習処理のパフォーマンスが3倍~10倍向上<sup>1</sup>



### AIパイプラインの最適化

企業はさまざまなシナリオに人工知能 (AI) を適用してメリットを得ることができます。書籍や映画のレコメンドサービスから、大規模なeコマースサイトを動かす小売デジタル・ソフトウェア、チャットボットや機械翻訳の自然言語処理 (NLP) まで、適用範囲は多岐にわたります。しかし、複雑な環境や膨大なデータセットを理解し、以前は解明できなかった問題を解決するという、AIを価値あるツールにしている特性に秘められた、ビジネスにさらなる変革をもたらす可能性は計り知れません。ある調査によると、2025年までには、新たにリリースされるエンタープライズ・アプリケーションの90%にAI機能が組み込まれるようになると見込まれています。<sup>2</sup>

### AIパイプライン



外側の3つのボックスは、AIパイプラインのステージを、内側の5つのボックスは、AIワークロードを表しています。ボックスの大きさは、AIパイプライン内でプロセッサ・アクティビティが担う相対的なレベルを示しています。

図1. AIパイプライン内のAIワークロードとプロセッサ・アクティビティ

AIパイプライン最適化のために導入すべきは、内蔵AIアクセラレーターの1つであるインテル® アドバンスト・マトリクス・エクステンション (インテル® AMX) を搭載した第4世代インテル® Xeon® スケーラブル・プロセッサです。インテル® AMXは、AIアプリケーションにおけるCPUの用途として最も知られている推論と、比較的機能の多い学習処理とのバランスを意図して設計されました (図1を参照)。<sup>3</sup> インテル® Xeon® スケーラブル・プロセッサは、データセンターでAI推論ワークロードを実行しているプロセッサ・ユニットの (インストール・ベースで) 70% を占めており、新たなAI導入にインテル® AMX搭載の第4世代インテル® Xeon® スケーラブル・プロセッサを選択することは、AIワークロードを高速化する効率的かつコスト効率の高いアプローチとなります。<sup>4</sup>

## 内蔵アクセラレーターの事例

インテル® ディープラーニング・ブースト (インテル® DL ブースト)搭載の第3世代インテル® Xeon® スケーラブル・プロセッサをAIに導入すれば、ITチームはすぐにも顧客のサービスレベル・アグリーメント (SLA) を満たすことが可能ですが、インテル® AMXを搭載した第4世代インテル® Xeon® スケーラブル・プロセッサを使用すると、状況はさらに一変します。

図2は、インテル® AMXによりPyTorchでリアルタイム推論を実行したパフォーマンスが前世代と比較して最大5.7倍～10倍に向上した結果を示しています。また図3は、インテル® AMXによりPyTorchで学習処理を実行したパフォーマンスが前世代と比較して最大3.5倍～10倍に向上した結果を示しています。<sup>5</sup> このパフォーマンスの結果から、インテル® AMXによって顧客満足度の向上を促進できることが分かります。すでに使い慣れているソリューションに、大幅なパフォーマンス向上を実現する内蔵アクセラレーターのインテル® AMXが統合されたことで、AIアプリケーションに最適なCPUの選択に迷うことはなくなります。

### インテル® AMX 搭載の第4世代インテル® Xeon® スケーラブル・プロセッサにより、リアルタイム推論のパフォーマンスが前世代と比較して最大5.7倍～10倍向上 (値が大きいほど高性能)

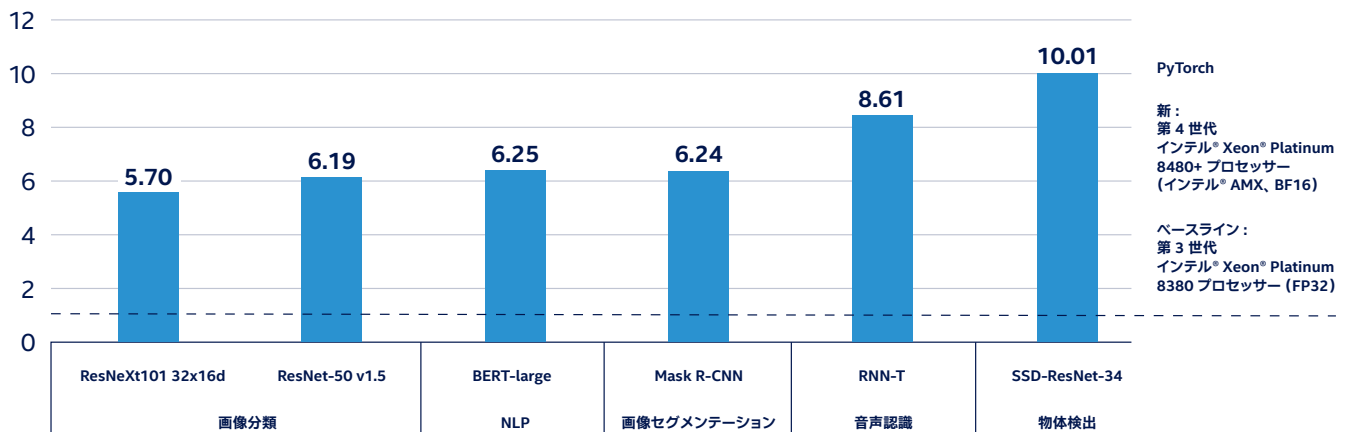


図2. インテル® AMX 搭載の第4世代インテル® Xeon® スケーラブル・プロセッサにより、PyTorchのリアルタイム推論パフォーマンスが向上<sup>5</sup>

### インテル® AMX 搭載の第4世代インテル® Xeon® スケーラブル・プロセッサにより、学習処理のパフォーマンスが前世代と比較して最大3.5倍～10倍に向上 (値が大きいほど高性能)

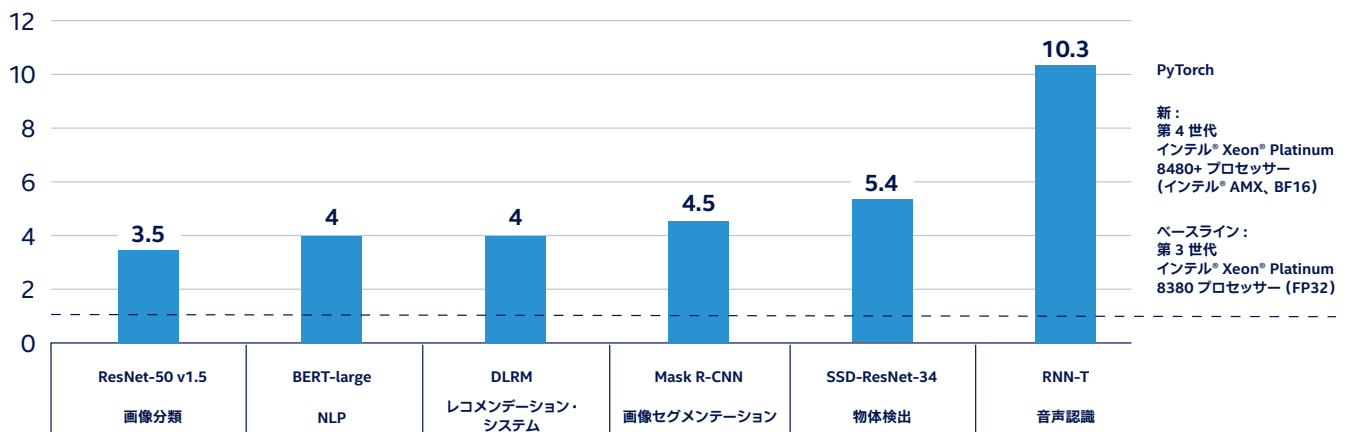


図3. インテル® AMX 搭載の第4世代インテル® Xeon® スケーラブル・プロセッサにより、PyTorchの学習処理パフォーマンスが向上<sup>5</sup>

図4は、初代のインテル® Xeon® スケーラブル・プロセッサから世代を追うごとに増加したコア数よりも高い比率で、インテル® AMXがパフォーマンス向上を実現していることを示しています。

## ムーアの法則とアクセラレーター

適切なワークロードに適切なコンピューティング・エンジンを割り当て

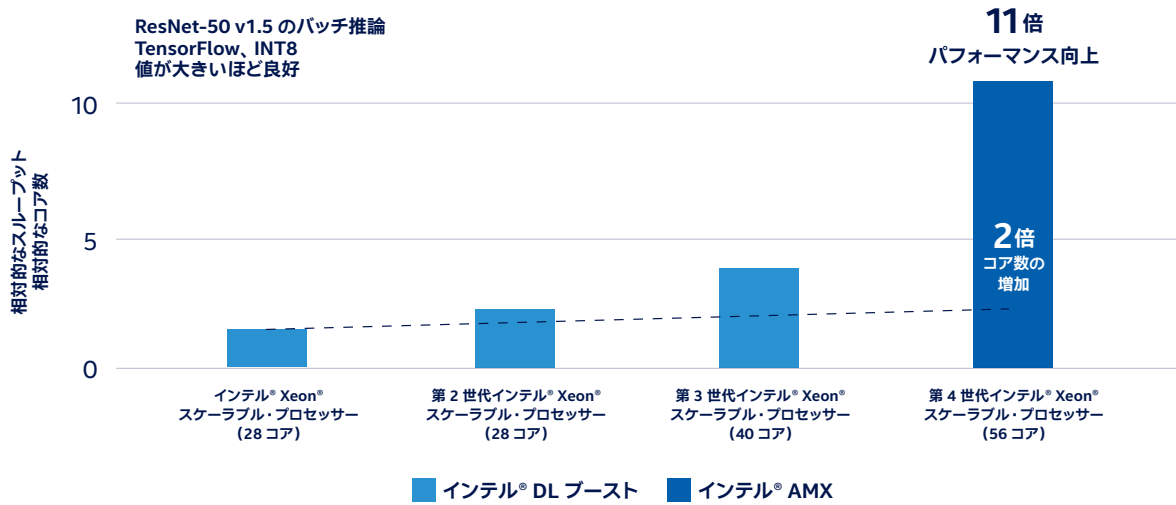


図4. 初代のインテル® Xeon® スケーラブル・プロセッサをベースラインとして、これまでの世代とは異なりインテル® AMXが実現するパフォーマンスの伸びは非線形<sup>6</sup>

### インテル® AMTとは

インテル® AMXは、第4世代インテル® Xeon® スケーラブル・プロセッサでディープラーニング (DL) 学習処理 / 推論ワークロードの最適化を可能にする内蔵アクセラレーターです。第4世代インテル® Xeon® スケーラブル・プロセッサはインテル® AMXを使用することで、汎用的なコンピューティングの最適化とAIワークロードの間の素早い切り替えが可能になります。例えば市街地での走行を得意とする車が、突然F1レースのパフォーマンスに切り替わる状況を想像してみてください。このようなタイプの柔軟性が、第4世代インテル® Xeon® スケーラブル・プロセッサには備わっています。開発者は、インテル® AMXの命令セットのメリットを活かしてAI機能をコーディングすることも、プロセッサの命令セット・アーキテクチャー (ISA) を使用してAI以外の機能をコーディングすることも可能です。インテルは、oneAPIのDLエンジンであるインテル® oneAPI ディープ・ニューラル・ネットワーク (oneDNN) ライブラリーを、AIアプリケーションに広く利用されているオープンソース・ツール (TensorFlow, PyTorch, PaddlePaddle, ONNX など) に統合しました。

### インテル® AMXアーキテクチャー

インテル® AMXのアーキテクチャーは、2つのコンポーネントで構成されています(図5を参照)。

- 1つ目のコンポーネントはタイルです。タイルは8つの2次元レジスターで構成されており、各タイルのサイズは1キロバイトです。タイルには大きなデータチャンクが格納されます。
- 2つ目のコンポーネントは、タイル行列乗算 (TMUL) です。TMULはタイルに接続されたアクセラレーター・エンジンであり、AIの行列乗算を実行します。



図5. 2Dレジスタファイル (タイル) とTMULで構成されるインテル® AMXアーキテクチャー

インテル® AMXでは、AIワークロードに必要な行列乗算で、INT8とBF16の2つのデータ型をサポートしています。

- INT8は、AIでよく使用される単精度浮動小数点形式であるFP32の精度が不要な場合に、推論に使用されるデータ型です。INT8データ型はFP32よりも精度が低い分、コンピューティング・サイクルあたりに処理できるINT8演算数が多くなります。
- BF16は、ほとんどの学習処理に十分な精度を提供するデータ型です。必要に応じて、推論の精度を高めることもできます。

この新しいタイル型アーキテクチャーにより、インテル® AMXは前世代と比較してパフォーマンスを大幅に向上させます。インテル® アドバンスト・ベクトル・エクステンション 512 ベクトル・ニューラル・ネットワーク・インストラクション (インテル® AVX-512 VNNI) を実行する第3世代インテル® Xeon® スケーラブル・プロセッサと比較して、インテル® AMXを実行する第4世代インテル® Xeon® スケーラブル・プロセッサでは、1サイクル当たり256回のINT8演算ではなく、2,048回のINT8演算を実行できます。また、図6に示すように、1サイクル当たり64回のFP32演算ではなく、1,024回のBF16演算を実行することも可能です。<sup>7</sup>

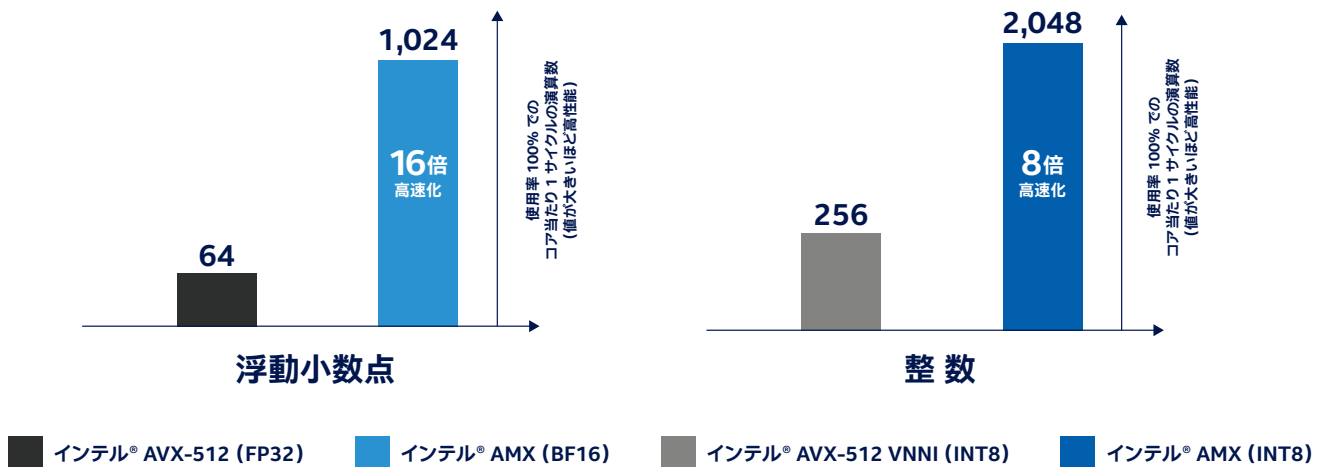


図6. インテル® AMXは、INT8/BF16データ型で、インテル® AVX-512VNNIよりも高いパフォーマンスを発揮<sup>7</sup>

## AIのユースケース

インテル® AMX 搭載の第4世代インテル® Xeon® スケーラブル・プロセッサは、幅広いDLユースケースに導入できます。



### レコメンドシステム

おすすめの映画や本の紹介、ターゲット広告の表示など、カスタマイズされたエンドユーザー体験を提供。リアルタイムのユーザー行動シグナルを反映し、時間や場所などのほぼリアルタイムのコンテキスト機能を備えた、DLベースのレコメンドシステムを構築します。



### 自然言語処理 (NLP)

言語推論やマシンラーニング (ML) を含む NLP アプリケーションは、2026年までに世界の市場規模が806億8,000万ドルに達すると予測されており<sup>8</sup>、企業がセンチメント分析、チャットボット、機械翻訳などのさまざまな機能を採用 / 拡張するために不可欠となっています。



### 小売業のeコマース・ソフトウェア・ソリューション

AIに最適化されたPyTorchやTensorFlowなどのフレームワークに加え、DL推論 / 学習処理を駆使し、トランザクション時間を最小限に抑え、ピーク需要に対応することで、収益の拡大と同時に、格別の顧客体験を提供します。

## インテル® AMXの利用開始

インテル® AMXによってパフォーマンスを向上させるために必要な作業はほとんどありません。これは、インテル® oneAPI ディープ・ニューラル・ネットワーク (oneDNN) でデフォルトのフレームワークが最適化されているためです。Windows/Linuxオペレーティング・システム、カーネルベースの仮想マシン (KVM)、広く利用されているハイパーバイザーは、インテル® AMX 命令セットを公開しています。TensorFlowやPyTorchといったオープンソースのフレームワークでは、INT8/BF16演算が自動的に最適化されます。インテル® ディストリビューションのOpenVINO™ ツールキットを活用することで、コーディングの知識がほとんどなくても、AI推論を自動化、最適化、調整し、実行できます。開発者が行う必要があるのは、インテル® ニューラル・コンプレッサーを使用して、学習処理モデルをINT8データ型に量子化することだけです。

## インテル® Xeon® スケーラブル・プロセッサによるAIの高速化

インテル® AMX を搭載した第4世代インテル® Xeon® スケーラブル・プロセッサに移行することで、まだ眠っているAIの可能性をビジネスで活用できます。すでにデータセンターで広く導入されているインテル® Xeon® スケーラブル・プロセッサの幅広い基盤の上に構築される、高速化された全く新しい行列乗算演算で、AI学習処理 / 推論の卓越したパフォーマンスを体験してください。

インテルのAIとインテル® AMXの詳細については、<https://www.intel.co.jp/ai/> をご覧ください。



- <https://www.intel.com/processorclaims/> (英語): 4th Generation Intel® Xeon® Scalable Processorsの[A16]および[A17]を参照。結果は異なる場合があります。
- 「Top Artificial Intelligence (AI) Predictions For 2020 From IDC and Forrester」(Forbes, 2019年11月)。 <https://www.forbes.com/sites/gilpress/2019/11/22/top-artificial-intelligence-ai-predictions-for-2020-from-idc-and-forrester/#4fef9821315a> (英語)
- 「With AMX, Intel Adds AI/ML Sparkle to Sapphire Rapids」(The Next Platform, 2021年8月)。 <https://www.nextplatform.com/2021/08/19/with-amx-intel-adds-ai-ml-sparkle-to-sapphire-rapids/> (英語)
- AI 推論ワークロードを実行するデータセンター・サーバーを対象に世界全体のインストール・ベースでインテルが実施した市場モデリングに基づく(2021年12月時点)。
- PyTorch モデルのパフォーマンス構成。PT-NLP BERT-large: 8480 : 1ノード、インテル® Xeon® Platinum 8480 プロセッサー x2 を搭載した製品出荷前のプラットフォーム、総メモリ容量 1,024GB (16スロット/64GB/DDR5-4800)、ucode 0x2b0000a1、インテル® ハイパースレッディング・テクノロジー(インテル® HT テクノロジー)有効、インテル® ターボ・ブースト・テクノロジー有効、CentOS Stream 8, 5.15.0, 1TB インテル® SSDSC2KW256G8 (PT)/Samsung SSD 860 EVO (TF) x1, BERT-large, Inf: SquAD1.1 (seq len=384), bs=1 (インスタンス当たり4コア), bs=n (インスタンス当たり1ソケット), bs: FP32=1,56, インテル® AMX BF16=1,16, インテル® AMX INT8=1,56, Trg: Wikipedia 2020/01/01 (seq len=512), bs: FP32=28, インテル® AMX BF16=56 (1インスタンス, 1ソケット), フレームワーク: <https://github.com/intel-innersource/frameworks.ai.pytorch.private-cpu/tree/d7607bdd983093396a70713344828a989b766a66/>, ModelZoo: <https://github.com/IntelAI/models/tree/spr-launch-public/>, PT: 1.13, IPEX: 1.13, oneDNN: v2.7 (2022年10月24日に実施したインテル社内テストで測定)。8380: 1ノード、インテル® Xeon® Platinum 8380 プロセッサー x2, 総メモリ容量 1,024GB (16スロット/64GB/DDR4-3200), ucode 0xd000375, インテル® HT テクノロジー有効、インテル® ターボ・ブースト・テクノロジー有効、Ubuntu 22.04 LTS, 5.15.0-27-generic, インテル® SSDSC2KG960G8 x1, BERT-large, Inf: SquAD1.1 (seq len=384), bs=1 (インスタンス当たり4コア), bs=n (インスタンス当たり1ソケット), bs: FP32=1,56, INT8=1,56, Trg: Wikipedia 2020/01/01 (seq len=512), bs: FP32=28, インテル® AMX BF16=56 (1インスタンス, 1ソケット), フレームワーク: <https://github.com/intel-innersource/frameworks.ai.pytorch.private-cpu/tree/d7607bdd983093396a70713344828a989b766a66/>, ModelZoo: <https://github.com/IntelAI/models/tree/spr-launch-public/>, PT: 1.13, IPEX: 1.13, oneDNN: v2.7 (2022年10月24日に実施したインテル社内テストで測定)。8380: 1ノード、インテル® Xeon® Platinum 8480 プロセッサー x2 を搭載した製品出荷前のプラットフォーム、総メモリ容量 1,024GB (16スロット/64GB/DDR5-4800)、ucode 0x2b0000a1、インテル® HT テクノロジー有効、インテル® ターボ・ブースト・テクノロジー有効、CentOS Stream 8, 5.15.0, 1TB インテル® SSDSC2KW256G8 (PT)/Samsung SSD 860 EVO (TF) x1, DLRM, 推論: bs=n (インスタンス当たり1ソケット), bs: FP32=128, インテル® AMX BF16=128, インテル® AMX INT8=128, 学習処理 bs: fp32/インテル® AMX BF16=32K (1インスタンス, 1ソケット), テラバイトの Criteo データセット, フレームワーク: <https://github.com/intel-innersource/frameworks.ai.pytorch.private-cpu/tree/d7607bdd983093396a70713344828a989b766a66/>, ModelZoo: <https://github.com/IntelAI/models/tree/spr-launch-public/>, PT: 1.13, IPEX: 1.13, oneDNN: v2.7 (2022年10月24日に実施したインテル社内テストで測定)。PT-ResNet-34: 8480 : 1ノード、インテル® Xeon® Platinum 8480 プロセッサー x2 を搭載した製品出荷前のプラットフォーム、総メモリ容量 1,024GB (16スロット/64GB/DDR5-4800)、ucode 0x2b0000a1、インテル® HT テクノロジー有効、インテル® ターボ・ブースト・テクノロジー有効、CentOS Stream 8, 5.15.0, 1TB インテル® SSDSC2KW256G8 (PT)/Samsung SSD 860 EVO (TF) x1, DLRM, 推論: bs=n (インスタンス当たり1ソケット), bs: FP32=128, インテル® AMX BF16=128, インテル® AMX INT8=128, 学習処理 bs: fp32/インテル® AMX BF16=32K (1インスタンス, 1ソケット), テラバイトの Criteo データセット, フレームワーク: <https://github.com/intel-innersource/frameworks.ai.pytorch.private-cpu/tree/d7607bdd983093396a70713344828a989b766a66/>, ModelZoo: <https://github.com/IntelAI/models/tree/spr-launch-public/>, PT: 1.13, IPEX: 1.13, oneDNN: v2.7 (2022年10月24日に実施したインテル社内テストで測定)。8380: 1ノード、インテル® Xeon® Platinum 8380 プロセッサー x2, 総メモリ容量 1,024GB (16スロット/64GB/DDR4-3200), ucode 0xd000375, インテル® HT テクノロジー有効、インテル® ターボ・ブースト・テクノロジー有効、Ubuntu 22.04 LTS, 5.15.0-27-generic, インテル® SSDSC2KG960G8 x1, DLRM, 推論: bs=n (インスタンス当たり1ソケット), bs: FP32=128, インテル® AMX BF16=128, インテル® AMX INT8=128, 学習処理 bs: fp32/インテル® AMX BF16=32K (1インスタンス, 1ソケット), テラバイトの Criteo データセット, フレームワーク: <https://github.com/intel-innersource/frameworks.ai.pytorch.private-cpu/tree/d7607bdd983093396a70713344828a989b766a66/>, ModelZoo: <https://github.com/IntelAI/models/tree/spr-launch-public/>, PT: 1.13, IPEX: 1.13, oneDNN: v2.7 (2022年10月24日に実施したインテル社内テストで測定)。PT-ResNet-50: 8480 : 1ノード、インテル® Xeon® Platinum 8480 プロセッサー x2 を搭載した製品出荷前のプラットフォーム、総メモリ容量 1,024GB (16スロット/64GB/DDR5-4800)、ucode 0x2b0000a1、インテル® HT テクノロジー有効、インテル® ターボ・ブースト・テクノロジー有効、CentOS Stream 8, 5.15.0, 1TB インテル® SSDSC2KW256G8 (PT)/Samsung SSD 860 EVO (TF) x1, ResNet-50 v1.5, 推論: bs=1 (インスタンス当たり4コア), bs=n (インスタンス当たり1ソケット), bs: FP32=1,64, インテル® AMX BF16=1,64, インテル® AMX INT8=1,116, 学習処理 bs: FP32/インテル® AMX BF16=224 (1インスタンス, 1ソケット), Coco 2017, フレームワーク: <https://github.com/intel-innersource/frameworks.ai.pytorch.private-cpu/tree/d7607bdd983093396a70713344828a989b766a66/>, ModelZoo: <https://github.com/IntelAI/models/tree/spr-launch-public/>, PT: 1.13, IPEX: 1.13, oneDNN: v2.7 (2022年10月24日に実施したインテル社内テストで測定)。PT-RNN-T: 8480 : 1ノード、インテル® Xeon® Platinum 8480 プロセッサー x2 を搭載した製品出荷前のプラットフォーム、総メモリ容量 1,024GB (16スロット/64GB/DDR5-4800)、ucode 0x2b0000a1、インテル® HT テクノロジー有効、インテル® ターボ・ブースト・テクノロジー有効、CentOS Stream 8, 5.15.0, 1TB インテル® SSDSC2KW256G8 (PT)/Samsung SSD 860 EVO (TF) x1, ResNet-50 v1.5, 推論: bs=1 (インスタンス当たり4コア), bs=n (インスタンス当たり1ソケット), bs: FP32=1,64, インテル® AMX BF16=1,64, インテル® AMX INT8=1,116, 学習処理 bs: FP32/インテル® AMX BF16=128 (1インスタンス, 1ソケット), ImageNet (224 x 224), フレームワーク: <https://github.com/intel-innersource/frameworks.ai.pytorch.private-cpu/tree/d7607bdd983093396a70713344828a989b766a66/>, ModelZoo: <https://github.com/IntelAI/models/tree/spr-launch-public/>, PT: 1.13, IPEX: 1.13, oneDNN: v2.7 (2022年10月24日に実施したインテル社内テストで測定)。PT-RNN-T: 8480 : 1ノード、インテル® Xeon® Platinum 8480 プロセッサー x2 を搭載した製品出荷前のプラットフォーム、総メモリ容量 1,024GB (16スロット/64GB/DDR5-4800)、ucode 0x2b0000a1、インテル® HT テクノロジー有効、インテル® ターボ・ブースト・テクノロジー有効、CentOS Stream 8, 5.15.0, 1TB インテル® SSDSC2KW256G8 (PT)/Samsung SSD 860 EVO (TF) x1, Resnet101 32x16d, 推論: bs=1 (インスタンス当たり4コア), bs=n (インスタンス当たり1ソケット), bs: FP32=1,64, インテル® AMX BF16=1,64, インテル® AMX INT8=1,116, ImageNet, フレームワーク: <https://github.com/intel-innersource/frameworks.ai.pytorch.private-cpu/tree/d7607bdd983093396a70713344828a989b766a66/>, ModelZoo: <https://github.com/IntelAI/models/tree/spr-launch-public/>, PT: 1.13, IPEX: 1.13, oneDNN: v2.7 (2022年10月24日に実施したインテル社内テストで測定)。8380: 1ノード、インテル® Xeon® Platinum 8380 プロセッサー x2, 総メモリ容量 1,024GB (16スロット/64GB/DDR4-3200), ucode 0xd000375, インテル® HT テクノロジー有効、インテル® ターボ・ブースト・テクノロジー有効、Ubuntu 22.04 LTS, 5.15.0-27-generic, インテル® SSDSC2KG960G8 x1, Resnet101 32x16d, 推論: bs=1 (インスタンス当たり4コア), bs=n (インスタンス当たり1ソケット), bs: FP32=1,64, INT8=1,116, ImageNet, フレームワーク: <https://github.com/intel-innersource/frameworks.ai.pytorch.private-cpu/tree/d7607bdd983093396a70713344828a989b766a66/>, ModelZoo: <https://github.com/IntelAI/models/tree/spr-launch-public/>, PT: 1.13, IPEX: 1.13, oneDNN: v2.7 (2022年10月24日に実施したインテル社内テストで測定)。8380: 1ノード、インテル® Xeon® Platinum 8380 プロセッサー x2, 総メモリ容量 1,024GB (16スロット/64GB/DDR4-3200), ucode 0xd000375, インテル® HT テクノロジー有効、インテル® ターボ・ブースト・テクノロジー有効、Ubuntu 22.04 LTS, 5.15.0-27-generic, インテル® SSDSC2KG960G8 x1, Resnet101 32x16d, 推論: bs=1 (インスタンス当たり4コア), bs=n (インスタンス当たり1ソケット), bs: FP32=1,64, INT8=1,116, ImageNet, フレームワーク: <https://github.com/intel-innersource/frameworks.ai.pytorch.private-cpu/tree/d7607bdd983093396a70713344828a989b766a66/>, ModelZoo: <https://github.com/IntelAI/models/tree/spr-launch-public/>, PT: 1.13, IPEX: 1.13, oneDNN: v2.7 (2022年10月24日に実施したインテル社内テストで測定)。PT-MaskRCNN: 8480 : 1ノード、インテル® Xeon® Platinum 8480 プロセッサー x2 を搭載した製品出荷前のプラットフォーム、総メモリ容量 1,024GB (16スロット/64GB/DDR5-4800)、ucode 0x2b0000a1、インテル® HT テクノロジー有効、インテル® ターボ・ブースト・テクノロジー有効、CentOS Stream 8, 5.15.0, 1TB インテル® SSDSC2KW256G8 (PT)/Samsung SSD 860 EVO (TF) x1, MaskRCNN, 推論: bs=1 (インスタンス当たり4コア), bs=n (インスタンス当たり1ソケット), bs: FP32=1,112, インテル® AMX BF16=1,112, 学習処理 bs: FP32/インテル® AMX BF16=112 (1インスタンス, 1ソケット), Coco 2017, フレームワーク: <https://github.com/intel-innersource/frameworks.ai.pytorch.private-cpu/tree/d7607bdd983093396a70713344828a989b766a66/>, ModelZoo: <https://github.com/IntelAI/models/tree/spr-launch-public/>, PT: 1.13, IPEX: 1.13, oneDNN: v2.7 (2022年10月24日に実施したインテル社内テストで測定)。8380: 1ノード、インテル® Xeon® Platinum 8380 プロセッサー x2, 総メモリ容量 1,024GB (16スロット/64GB/DDR4-3200), ucode 0xd000375, インテル® HT テクノロジー有効、インテル® ターボ・ブースト・テクノロジー有効、Ubuntu 22.04 LTS, 5.15.0-27-generic, インテル® SSDSC2KG960G8 x1, MaskRCNN, 推論: bs=1 (インスタンス当たり4コア), bs=n (インスタンス当たり1ソケット), bs: FP32=1,112, 学習処理 bs: FP32=112 (1インスタンス, 1ソケット), Coco 2017, フレームワーク: <https://github.com/intel-innersource/frameworks.ai.pytorch.private-cpu/tree/d7607bdd983093396a70713344828a989b766a66/>, ModelZoo: <https://github.com/IntelAI/models/tree/spr-launch-public/>, PT: 1.13, IPEX: 1.13, oneDNN: v2.7 (2022年10月24日に実施したインテル社内テストで測定)。推論: ResNet-50 v1.5: ImageNet (224 x 224), SSD ResNet-34: Coco 2017 (1200 x 1200), BERT-large: SquAD1.1 (seq len=384), Resnet101: ImageNet, Mask RCNN: COCO 2017, DLRM: テラバイトの Criteo データセット, RNN: LibriSpeech, 学習処理: ResNet-50 v1.5: ImageNet (224 x 224), SSD ResNet-34: COCO 2017, BERT-large: Wikipedia 2020/01/01 (seq len=512), DLRM: テラバイトの Criteo データセット, RNN: LibriSpeech, Mask RCNN: COCO 2017。
- INT8 測定のソフトウェア構成: TensorFlow ResNet-50 v1.5, 推論: BS=116 (INT8), ソケット当たり1インスタンス, oneDNN v2.7, インテルが最適化したTensorFlow 2.10, 2022年10月24日(第3世代および第4世代インテル® Xeon® スケールアップ・プロセッサー), 2022年7月19日(第2世代および初代のインテル® Xeon® スケールアップ・プロセッサー)に実施したインテル社内テストで測定。ハードウェア構成: 第4世代インテル® Xeon® スケールアップ・プロセッサーのハードウェア構成(測定済み): 2ソケットのインテル® Xeon® Platinum 8480 プロセッサー(56コア, 350W 熱設計電力(TDP))を搭載した製品出荷前のプラットフォーム、DDR5 総メモリ容量 1TB (8チャネル/64GB/4,800MHz), BKC 01 を使用、インテル® AMX (INT8, BF16) を使用、CentOS Stream 8, インテル® AMX カーネル (5.15), 測定値は異なる場合があります。第3世代インテル® Xeon® スケールアップ・プロセッサーのハードウェア構成(測定済み): 1ノード、インテル® Xeon® Platinum 8380 プロセッサー x2 (40コア/2.30GHz, 270W TDP), DDR4 総メモリ容量 1TB (8スロット/64GB/3,200MHz), ucode 0xd0002f2, インテル® HT テクノロジー有効、インテル® ターボ・ブースト・テクノロジー有効、Ubuntu 20.04.2 LTS (Focal Fossa), 5.4.0-73-generic, インテル® SSDSC2CW480A3 OS ドライブ x1。第2世代インテル® Xeon® スケールアップ・プロセッサーのハードウェア構成(測定済み): 1ノード、2ソケットのインテル® Xeon® Platinum 8280 プロセッサー(28コア), インテル® HT テクノロジー有効、インテル® ターボ・ブースト・テクノロジー有効、総メモリ容量 384GB (12スロット/32GB/2,933MHz), BIOS: SE5C620, 86B.02.01.0013.12152020065 (ucode: 0x500320a), CentOS Stream 8, 4.18.0-383.el8.x86\_64, インテル® Xeon® スケールアップ・プロセッサーのハードウェア構成(測定済み): 1ノード、2ソケットのインテル® Xeon® Platinum 8120 プロセッサー(28コア), インテル® HT テクノロジー有効、インテル® ターボ・ブースト・テクノロジー有効、総メモリ容量 384GB (12スロット/32GB/2,666MHz), BIOS: SE5C620.86B.0X.01.0117.021220182317 (ucode: 0x200b062), Ubuntu 20.04.2 LTS, 5.4.0-73-generic。
- CPU 利用率 100% を想定した、コア当たりの1サイクルの行列乗算演算 + 累積演算のピーク・アーキテクチャー能力に基づく。2021年8月時点。完全なワークロードと構成の詳細については <https://www.intel.com/PerformanceIndex/> (英語): Architecture Day 2021 を参照。結果は異なる場合があります。
- 世界の NLP 市場規模の洞察: 「Natural Language Processing (NLP) Market Size, Share & COVID-19 Impact Analysis, By Deployment (On-Premises, Cloud, Hybrid), By Enterprise Size (SMEs, and Large Enterprises), By Technology (Interactive Voice Response (IVR), Optical Character Recognition (OCR), Text Analytics, Speech Analytics, Classification and Categorization), By Industry Vertical (Healthcare, Retail, High Tech, and Telecom, BFSI) and Regional Forecast, 2021-2028」(Fortune Business Insights, 2021年6月)。 <https://www.fortunebusinessinsights.com/industry-reports/natural-language-processing-nlp-market-101933#> (英語)

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、<https://www.Intel.com/PerformanceIndex/>(英語)を参照してください。

性能の測定結果は、構成に示されている日付時点のテストに基づいています。また、現在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やコンポーネントはありません。

コストと結果は状況によって異なります。

インテル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。

インテルは、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。ほかの情報も参考にしてデータの正確さを評価してください。

Intel、インテル、Intelロゴ、その他のインテルの名称やロゴは、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。

## インテル株式会社

〒100-0005 東京都千代田区丸の内3-1-1

<http://www.intel.co.jp/>

©2023 Intel Corporation. 無断での引用、転載を禁じます。

2023年8月