

## コンフィデンシャル・コンピューティングで AI 推論を高速化

Fortanix は、インテル® アクセラレーター・エンジンのワークロードを最適化することで、第4世代インテル® Xeon® スケーラブル・プロセッサ・ファミリーを使用して、クラウドでの AI 推論の安全性と高速化を支援します。



コンフィデンシャル AI では、セキュリティとパフォーマンスの両方が重要です。第4世代インテル® Xeon® スケーラブル・プロセッサ・ファミリーは、AI 推論の安全性と高速化を支援する設計になっています。インテル® アクセラレーター・エンジンはインテル® Xeon® スケーラブル・プロセッサ・ファミリーに内蔵された専用アクセラレーターであり、今日急成長を遂げているワークロードの多くに、パフォーマンスと電力効率のメリットをもたらします。

インテル® ソフトウェア・ガード・エクステンションズ (インテル® SGX) とインテル® アドバンスト・マトリクス・エクステンション (インテル® AMX) の両方を使用する [Fortanix Runtime Encryption® \(RTE\) プラットフォーム](#) 上のワークロードでは、TensorFlow Resnet50 を実行する場合のパフォーマンスが最大 7.57 倍<sup>1</sup>、Bert-Large を実行する場合のパフォーマンスが最大 5.26 倍向上します。<sup>2</sup>

インテル® AMX などのアクセラレーターにより、AI 推論ワークロードは、インテル® SGX などのハードウェアベースのセキュリティと組み合わせられた場合でも、優れたパフォーマンスを発揮できます。

### コンフィデンシャル・コンピューティングがクラウドベース AI の 安全性を支援

Fortanix RTE は [インテル® SGX によるコンフィデンシャル・コンピューティング](#) を使用して、OS やそのほか実行中のプロセスにプレーンテキストのアプリケーション・コードを露出することなく、暗号化されたデータの汎用演算を可能にします。インフラストラクチャーが侵害された場合、また悪意のある組織内の人間がルートパスワードを保持している場合でも、アプリケーションは暗号により保護されます。

コンフィデンシャル・コンピューティングでは、ほかのソフトウェア、コラボレーター、クラウド・プロバイダーなどに機密データを露出することなく使用する方法で、インサイトの導出や AI モデルのトレーニングが可能です。これにより、「機密性が高い」「規制されている」などの理由で以前は分析やその他の目的のために有効化できなかったデータを使用した業務改革が可能になります。

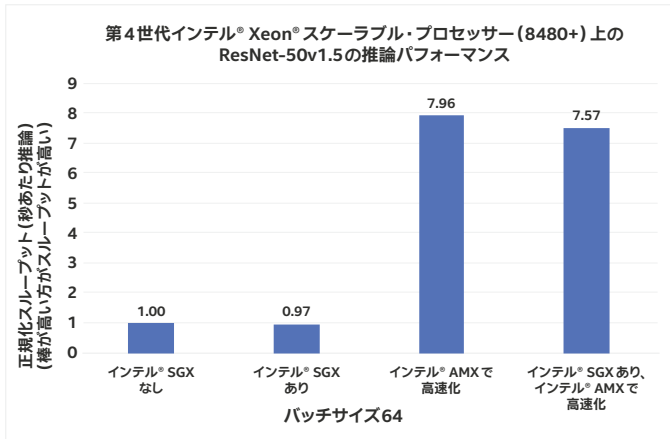


図1. Intel® AMX および Intel® SGX 搭載の第4世代Intel® Xeon® スケーラブル・プロセッサを使用した ResNet50 推論ワークロードのパフォーマンス<sup>1</sup>

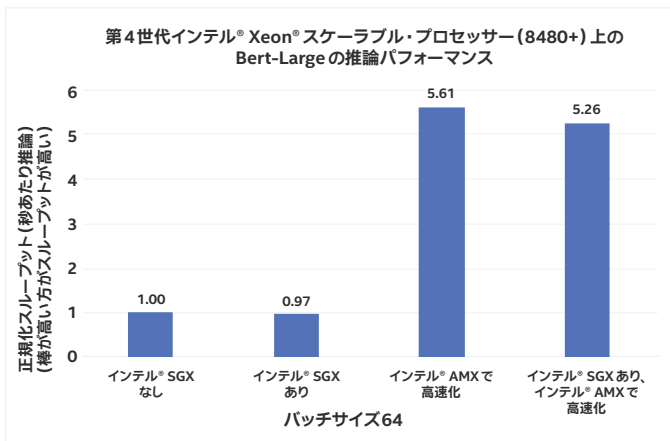


図2. Intel® AMX および Intel® SGX 搭載の第4世代Intel® Xeon® スケーラブル・プロセッサを使用した Bert-Large 推論ワークロードのパフォーマンス<sup>2</sup>

### Intel® AMX で AI 機能を高速化

Intel® AMX は、ディープラーニング・トレーニングと CPU 上の推論のパフォーマンスを向上させる新しい内蔵アクセラレーターであり、自然言語処理、推奨システム、イメージ認識などのワークロードに最適です。Intel は第4世代 Intel® Xeon® スケーラブル・プロセッサと Intel® AMX で AI 機能をさらに進化させています。推論やトレーニングのパフォーマンスも前世代の Intel® Xeon® スケーラブル・プロセッサに比べて向上しています。<sup>3</sup>

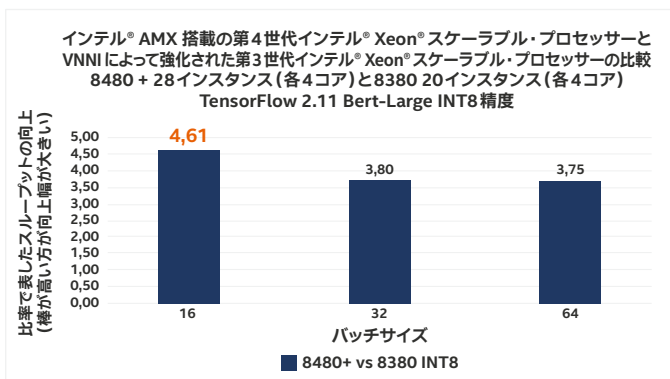


図3. 第3世代Intel® Xeon® スケーラブル・プロセッサと第4世代Intel® Xeon® スケーラブル・プロセッサのINT8精度におけるセキュアDL 推論ワークロードのパフォーマンスの比較<sup>3</sup>

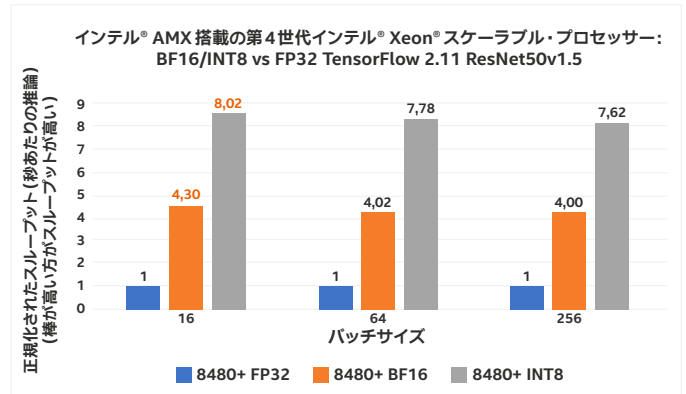


図4. 第4世代Intel® Xeon® スケーラブル・プロセッサ上の TensorFlow ResNet50 のさまざまなバッチサイズでのセキュアDL 推論ワークロードのパフォーマンス<sup>4</sup>

Intel® SGX と Intel® AMX を両方搭載した Fortanix RTE を使用した場合、第4世代 Intel® Xeon® スケーラブル・プロセッサ上の INT8 精度における Bert-Large 推論ワークロードのパフォーマンスは、第3世代 Intel® Xeon® スケーラブル・プロセッサ上と比較して最大 4.61 倍向上します。<sup>3</sup>

### 市場最多の内蔵アクセラレーターを手中に

第4世代 Intel® Xeon® スケーラブル・プロセッサ・ファミリーは、市販されている CPU で最も多くのアクセラレーターを内蔵しており、新たなワークロード、特に AI を活用するワークロードのパフォーマンスを高めるために役立ちます。

パフォーマンスの向上に加えて、第4世代 Intel® Xeon® スケーラブル・プロセッサ・ファミリーは、絶えず変化する脅威環境の中でデータを保護する高度なセキュリティ・テクノロジーを備え、ビジネス上のインサイトを導出する新たな機会をもたらします。Intel® SGX など、セキュリティ向けに設計されたハードウェア対応の機能を有効にした場合でも、Intel® AMX などの Intel® アクセラレーター・エンジンにより AI 推論ワークロードの卓越したパフォーマンスが実現します。

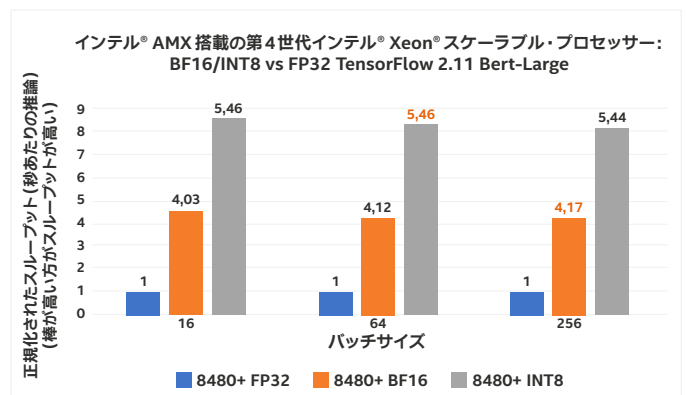


図5. 第4世代Intel® Xeon® スケーラブル・プロセッサ上の TensorFlow Bert-Large のさまざまなバッチサイズでのセキュアDL 推論ワークロードのパフォーマンス<sup>5</sup>

第4世代インテル® Xeon® スケーラブル・プロセッサであれば、パフォーマンスとセキュリティを両立できます。

詳細情報は、[intel.co.jp/xeonaccelerated/](https://intel.co.jp/xeonaccelerated/) および [fortanix.com/platform/runtime-encryption/](https://fortanix.com/platform/runtime-encryption/) (英語) を参照してください。



1. インテル® SGX およびインテル® AMX を搭載した第4世代インテル® Xeon® スケーラブル・プロセッサ上でTensorFlow ResNet50推論ワークロードを実行した場合、パフォーマンスが最大7.57倍向上。以下の構成情報を参照してください。
2. インテル® SGX およびインテル® AMX を搭載した第4世代インテル® Xeon® スケーラブル・プロセッサ上でBert-Large 推論ワークロードを実行した場合、パフォーマンスが最大5.26倍向上。以下の構成情報を参照してください。
3. INT8精度でBert-Large 推論ワークロードを実行した場合のパフォーマンスは、インテル® SGX およびインテル® AMX を搭載した第4世代インテル® Xeon® スケーラブル・プロセッサ上では前世代と比較して最大4.61倍向上。以下の構成情報を参照してください。
4. インテル® SGX およびインテル® AMX を搭載した第4世代インテル® Xeon® スケーラブル・プロセッサ上でTensorFlow ResNet50推論ワークロードを実行した場合、FP32との比較で、INT8精度でのパフォーマンスが最大8.02倍、BF16精度でのパフォーマンスが最大4.30倍向上。以下の構成情報を参照してください。
5. インテル® SGX およびインテル® AMX を搭載した第4世代インテル® Xeon® スケーラブル・プロセッサ上でBert-Large 推論ワークロードを実行した場合、FP32との比較で、INT8精度でのパフォーマンスが最大5.46倍、BF16精度でのパフォーマンスが最大4.17倍向上。以下の構成情報を参照してください。

#### 構成の詳細

TEST-1: 2022年11月21日にインテルが社内で実施したテストの結果。1ノード、インテル® Xeon® Platinum 8380 CPU @ 2.30GHz x2、40コア、HT 無効、ターボ有効、メモリー合計512GB (16x32GB DDR4 3200MT/s [run @3200MT/s]), BIOSバージョンSE5C6200.86B.0022.D64.2105220049、ucodeバージョン0xd000375、OSバージョンUbuntu 22.04.1 LTS、カーネルバージョン6.0.6-060006-generic、Fortanixの安全なエンクレーブ内のワークロード/ベンチマーク・デンプラーニング推論、フレームワークバージョンTensorFlow 2.11、モデル名 & バージョンResNet50v1.5/Bert-Large

TEST-2: 2022年11月21日にインテルが社内で実施したテストの結果。1ノード、インテル® Xeon® Platinum 8480+ CPU @ 2.0GHz x2、56コア、HT 無効、ターボ有効、メモリー合計512GB (16x32GB DDR5 4800MT/s [run @4800MT/s]), BIOSバージョン3A05、ucodeバージョン0x2b000070、OSバージョンUbuntu 22.04.1 LTS、カーネルバージョン6.0.6-060006-generic、Fortanixの安全なエンクレーブ内のワークロード/ベンチマーク・デンプラーニング推論、フレームワークバージョンTensorFlow 2.11、モデル名 & バージョンResNet50v1.5/Bert-Large

#### 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](#) のサイトをご覧ください。

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。

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。

© Intel Corporation. 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。