

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) によるマシンラーニング (ML) モデルの最適化

BERT (Bidirectional Encoder Representations from Transformers) モデルのスループットは、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を使用すると、前世代と比較してパフォーマンスが2倍から3倍向上します。^{1, 2}

このソリューション概要では、融合最適化という手法を導入するためのプロキシモデルとして、レイヤー数 12、隠れ層の次元 768、ヘッド数 12、シーケンス長 128 (トークンサイズ) の標準 BERT モデルを使用しています。

概要

BERT (Bidirectional Encoder Representations from Transformers) は、自然言語処理 (NLP) で広く使用されている ML モデルとその技術です。登場以来、BERT は膨大な数の NLP タスクのレコードを更新するために利用されてきました。また、実際のコアバウンド・アプリケーションでも優れたパフォーマンスを発揮します。

BERT は、検索、機械翻訳、マンマシン・インタラクション、その他の NLP タスクを対象に、多岐にわたるユーザーシナリオで採用されています。BERT のパフォーマンスが向上すればアプリケーションのユーザー・エクスペリエンスに直接影響するだけでなく、1秒当たりのクエリー数 (QPS) におけるスループット・レートが上昇するため、エンジニアは、モデルを最適化してパフォーマンスを向上させるさまざまな方法が検討されてきました。

例えばテンセントの StarLake Lab では、高度なクラウド・コンピューティング、人工知能 (AI)、セキュリティ、ストレージ、ネットワーク・テクノロジーを調査し、データセンターのパフォーマンスを向上させると同時に、データセンターの総保有コスト (TCO) を削減するソリューションを研究しています。テンセントの Machine Learning Platform Department (MLPD) は、AI プラットフォームの中核部門として、同社のインターネットやテクノロジー・ビジネス全体でイノベーションを推進する取り組みを常時行っています。MLPD は、コンピューター・ビジョン、音声認識、グラフ計算、NLP など、幅広い分野に及ぶ研究開発を行っており、開発されたソリューションは、ソーシャルメディア、パーソナライズド広告、ゲーミング AI、レコメンド機能、検索の主要なシナリオに広く使われています。これらのさまざまな技術領域のアプリケーションにおいて、BERT が重要な役割を果たしています。

インテルは、テンセントの MLPD や StarLake ラボラトリーと緊密に協力し、第4世代インテル® Xeon® スケーラブル・プロセッサの内蔵アクセラレーターであるインテル® AMX を活用した BERT 推論の最適化に取り組んでいます。インテル® AMX 内蔵の第4世代インテル® Xeon® スケーラブル・プロセッサを搭載したシステムで BERT モデルを実行したところ、BERT モデルでの [INT8] のスループットが2倍、[BF16] のスループットが3倍に向上する可能性があることが実証されました。^{1, 2} インテル® AMX とソフトウェアの最適化を組み合わせた統合ソリューションにより、テンセントは一貫したサービス体験を提供する機能や TCO を最適化する機能の進化を目指しています。

テンセントのソーシャル・アプリケーション最適化

テンセントのソーシャル・アプリケーションは、世界で10億人を超えるアクティブユーザーをつないでいます。テンセントのソーシャル・アプリケーションの中で最も人気のあるものは2011年にリリースされ、2018年には世界最大のスタンドアロン・モバイル・アプリになりました。このアプリにはメッセージ交換、ネット通話、一対多の動画配信、オンラインゲーム、ビデオ会議などの機能が揃い、さまざまな用途で使えることから「中国での生活はこれだけで完結する」と言われています。さらに、写真共有や動画共有、位置情報共有機能も備えています。

このアプリのユーザーの大部分は、メッセージや記事、ミニプログラム、ショート動画、音楽などの人気コンテンツを探す際に、アプリに搭載された検索エンジンを使用する傾向にあります。検索エンジンにとっての主な課題は、いかに大規模なクエリーを処理し、検索結果に素早く対応できるかです。このような課題のソリューションとして、インテル® AMXを利用した深い最適化によりTCOを削減すると同時に、この検索エンジンで使用していた既存の汎用インフラストラクチャーを活用し、アプリケーションの検索体験全体の底上げができると考えられています。

融合最適化

従来の融合最適化は、BERTベースのモデルの12層のレイヤーを1つの大規模なオペレーション(op、演算処理)に融合するというFP32のソリューションを通じて実現されていました。現在は、インテル® AMXによりさらに深い融合最適化に必要な機能が提供されています。

MatMulやBatchMatMulのオペレーションは、BERT側の時間を大量に消費する可能性があるため、MatMulとBatchMatMulを最適化することがパフォーマンス向上のポイントとなります。過去の実験から計算処理やメモリアクセスを削減すれば、パフォーマンスを最適化できます。モデルから不要なオペレーションを削除すれば、計算量も命令の数も減らせます。また複数のオペレーションを1つにまとめると、アクセスしたデータを次に使用するまでキャッシュ内に保持できるので、メモリアクセスを減らせます。

図1に示すように、テンセントのMLPDとインテルは、この考え方に基づいて特に時間のかかるプロセスを「Fused BERT op」という1つの大きなオペレーションに変更しました。

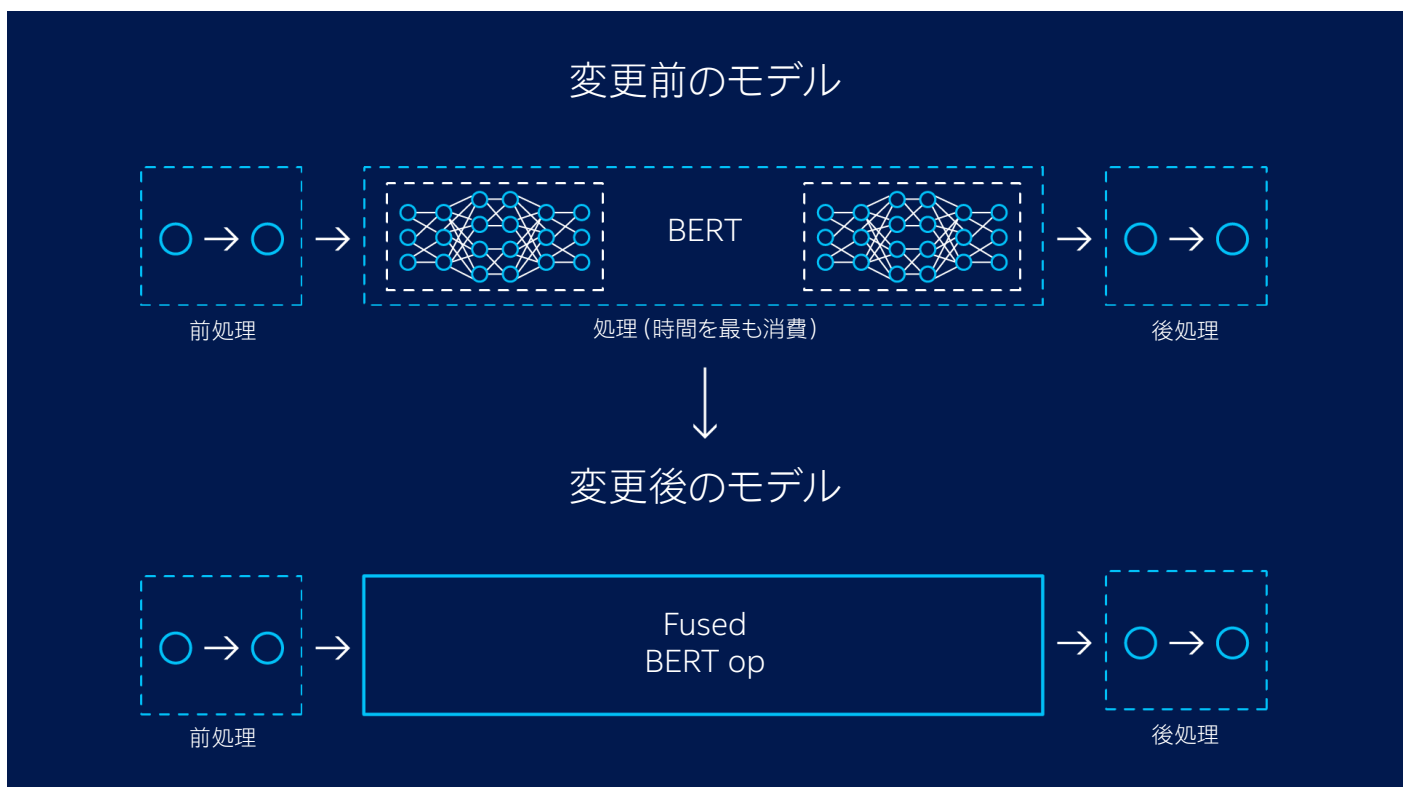


図1. 複数のオペレーションを1つの「Fused BERT op」にマージしてBERTを最適化

「Fused BERT op」は、次の方法で最適化を実現します。

- QKV MatMul (QKV: Query/クエリー、Key/キー、Value/値) の入力 が等しいことに着目し、これらのオペレーションの重みを1つの大きな重み行列にまとめ、1つの大きなQKV MatMulにマージしました。これらの重みが1つにマージされた後は、それぞれの重みと出力のメモリーが連続しなくなります。そのため、QKVの出力を次のオペレーションの入力として使用する場合は適切な「ストライド (間隔)」を設定する必要があります。この最適化フローを図2に示します。

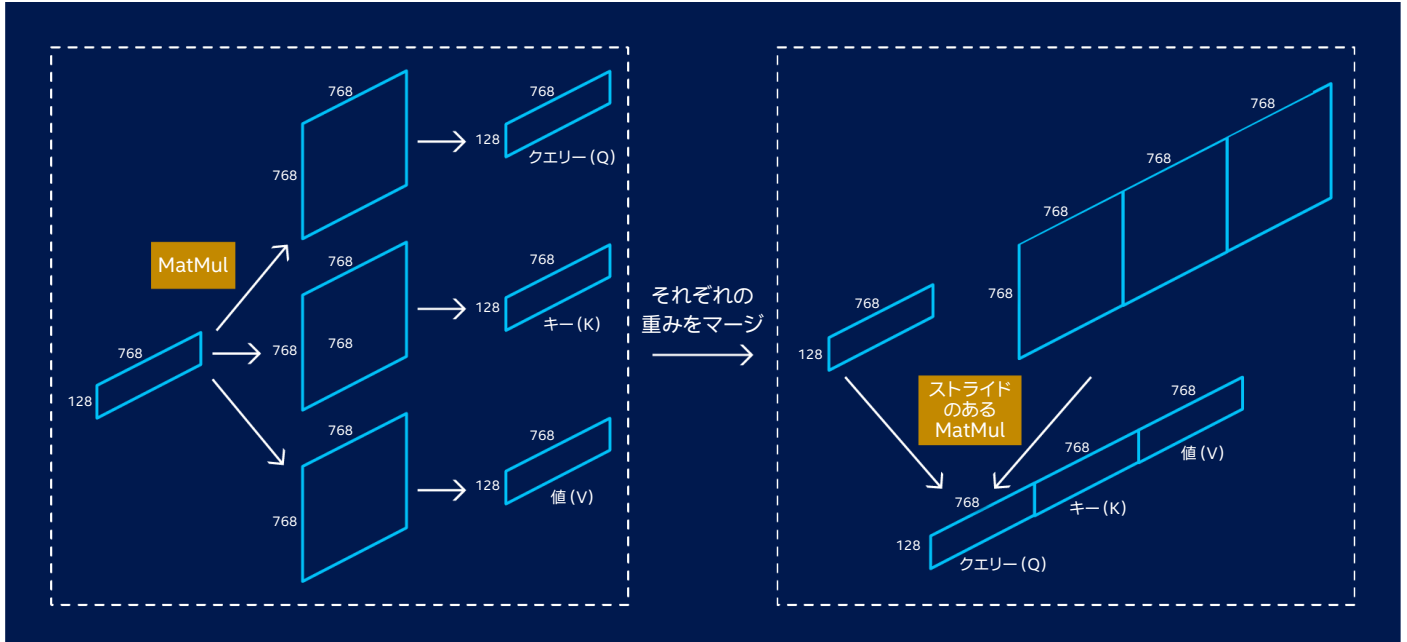


図2. QKV MatMul 最適化フローチャート

- インテル® oneAPI ディープ・ニューラル・ネットワーク・ライブラリー (インテル® oneDNN) がストライド (間隔) を持った BatchMatMul をサポートするので、BatchMatMul の前後の置き換えオペレーション (op) が削除できます。これにより、大量のメモリーアクセスと計算処理が節約できます。この最適化フローを図3に示します。

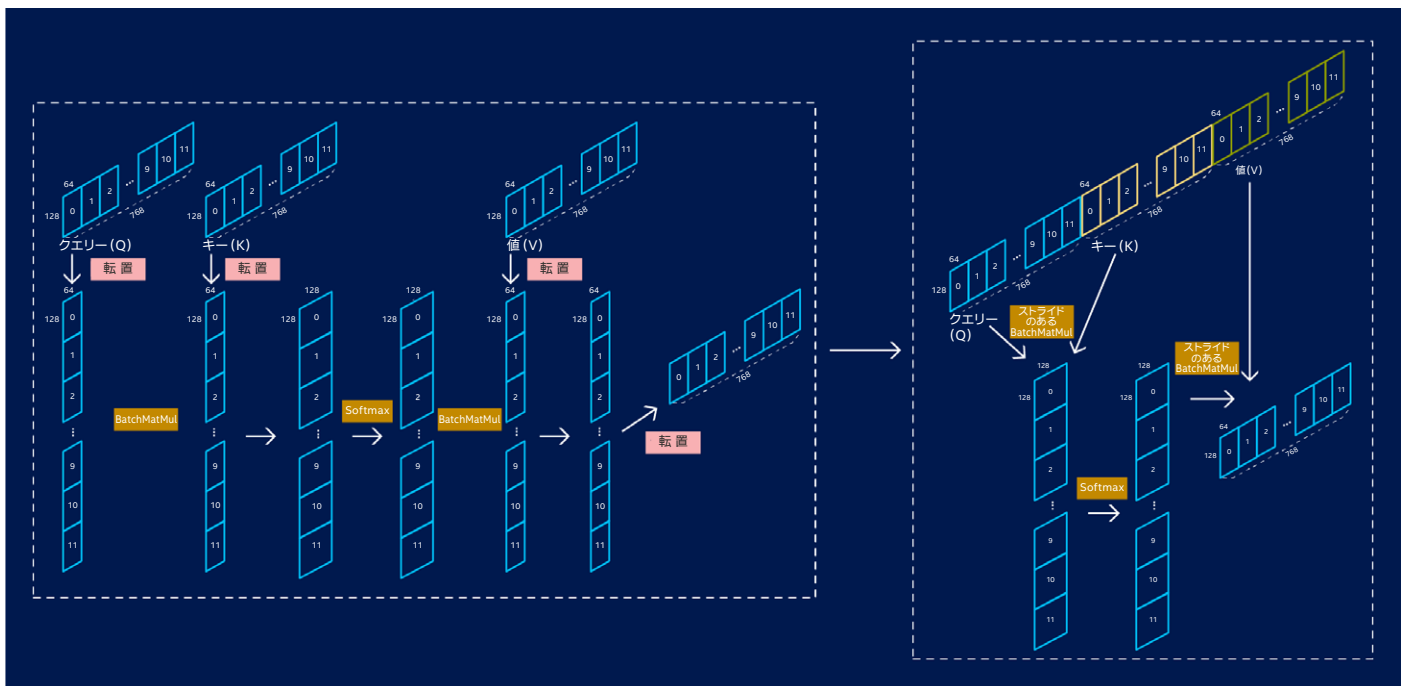


図3. BatchMatMul 最適化フローチャート

- インテル® oneDNN は、BiasAdd と特定の後処理 (OutputScale、Sum、Relu、Gelu、Tanh など) で MatMul をサポートするため、特定のオペレーションを MatMul のプリミティブに統合することができます。これらのオペレーション (op) を統合すると、キャッシュの使用効率が向上します。つまり、関連するデータを複数のタスクによって呼び出せるよう、キャッシュ内のデータをすぐに利用できる状態で保持できます。

Feature Dense 最適化

NLPのタスクでは、通常、フィーチャーやデータの長さが異なるため、バッチを形成するために大量のパディングを挿入する必要があります。その結果、不要な計算オーバーヘッドが大量に発生します。インテルはテンセントのMLPDと共同で、BERTモデル向けにFeature Dense最適化ソリューションを開発しました。この最適化により計算オーバーヘッドが削除され、タスクのパフォーマンスが大幅に向上します。また、バッチサイズが大きいほどパフォーマンスが向上します。図4では、BERTモデルの一部に絞って、この最適化がどのように機能するかの概要を示しています。

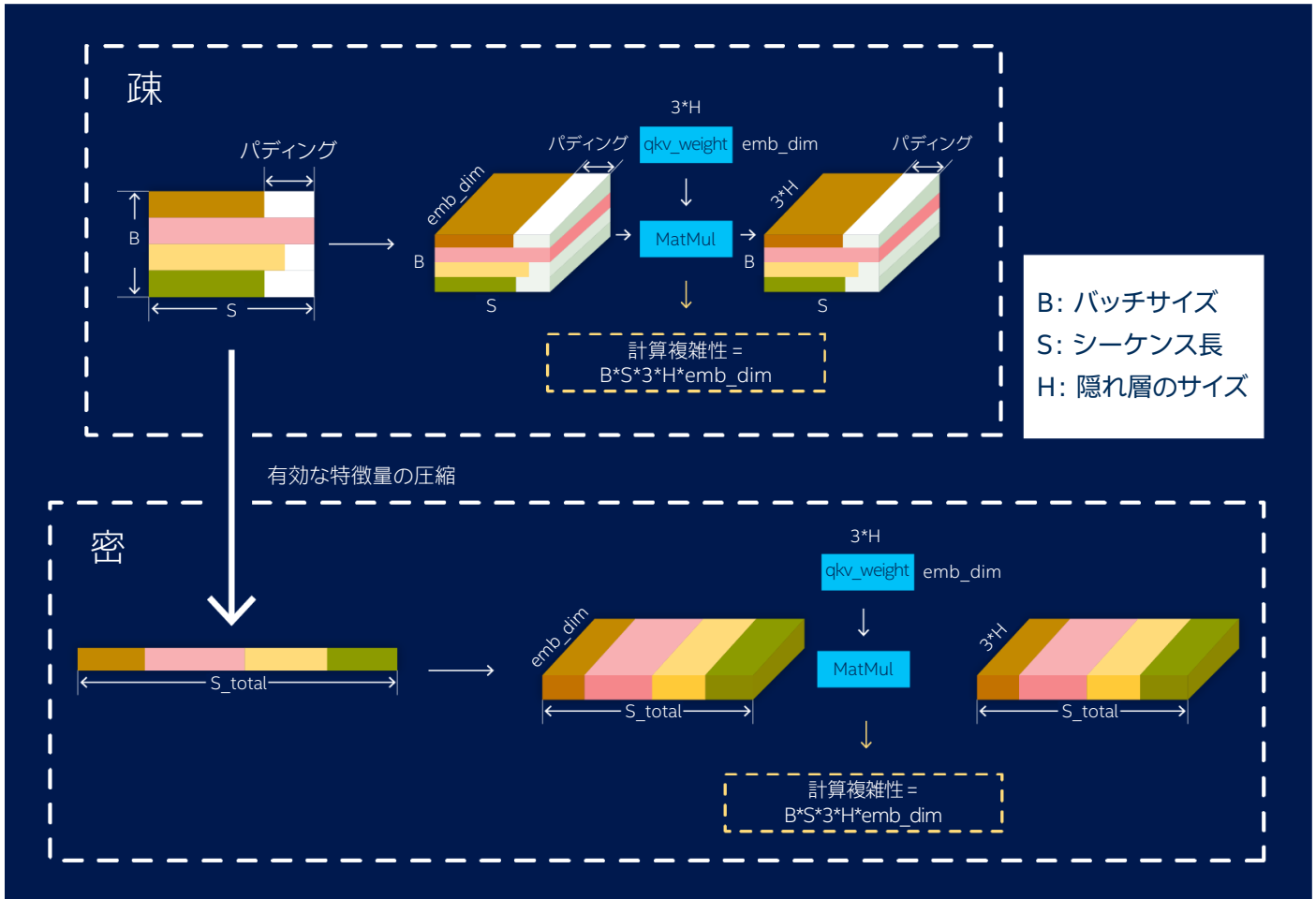


図4. Feature Dense 最適化ソリューション

図4に示すように、Feature Dense最適化ソリューションは、入力からパディングを排除し、各バッチのデータを1つずつつなげて、1次元のデータを形成します。埋め込み、QKV MatMul、その後のMatMulオペレーション (op: 図では省略) などの後続オペレーションは、圧縮されたデータを使うので、図4のように計算が大幅にシンプルになります。

BF16/INT8の最適化

上記のようにFP32の最適化によりBERTのパフォーマンスが大幅に向上しますが、さらに改善が可能です。メモリーのアクセスをさらに減らせば、パフォーマンスがもっと向上します。これは、計算中のデータサイズを小さくすることにより実現します。入力と重みを含むFP32のデータは、計算の前にBF16/INT8データに変換できます。

では、高いパフォーマンスを維持しながらBF16/INT8計算をサポートするのはどのようなプラットフォームでしょうか？ 第3世代インテル® Xeon® スケーラブル・プロセッサ (Cooper LakeおよびIce Lake) には、ベクトル乗算向けのVNNI (INT8) 命令をサポートする、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) が内蔵されています。またCooper LakeはBF16の数値フォーマットもサポートしています。このため、FP32最適化ソリューションを使ってBF16またはINT8を最適化できます。テストの結果では、BF16やINT8を最適化すると、FP32ソリューションと比較してパフォーマンスが大幅に向上することが確認されています。

インテル® AMXを内蔵した第4世代インテル® Xeon® スケーラブル・プロセッサ

BF16/INT8の最適化を行った上で、さらに最適化が可能でしょうか？可能です。インテル® AMXは、第4世代インテル® Xeon® スケーラブル・プロセッサに内蔵されたアクセラレーターです。インテル® AMXは、より大きな2次元メモリーイメージのサブアレイ(部分配列)で表される、2次元レジスター(タイル)のセット、そしてタイルのオペレーション(op)を可能にするアクセラレーターを備えた64ビット・プログラミング・パラダイムを提供します。最初の実装は「tile matrix multiply unit」を表すTMULユニットです。

図5は、インテル® AMXアーキテクチャーの概念図です。インテル®アーキテクチャー・ホストは、アルゴリズム、メモリーブロック、ループ・インデックス、ポインター演算を担っています。タイルの読み込みと格納、アクセラレーター・コマンドは、マルチサイクル実行ユニットのTMULに送られます。

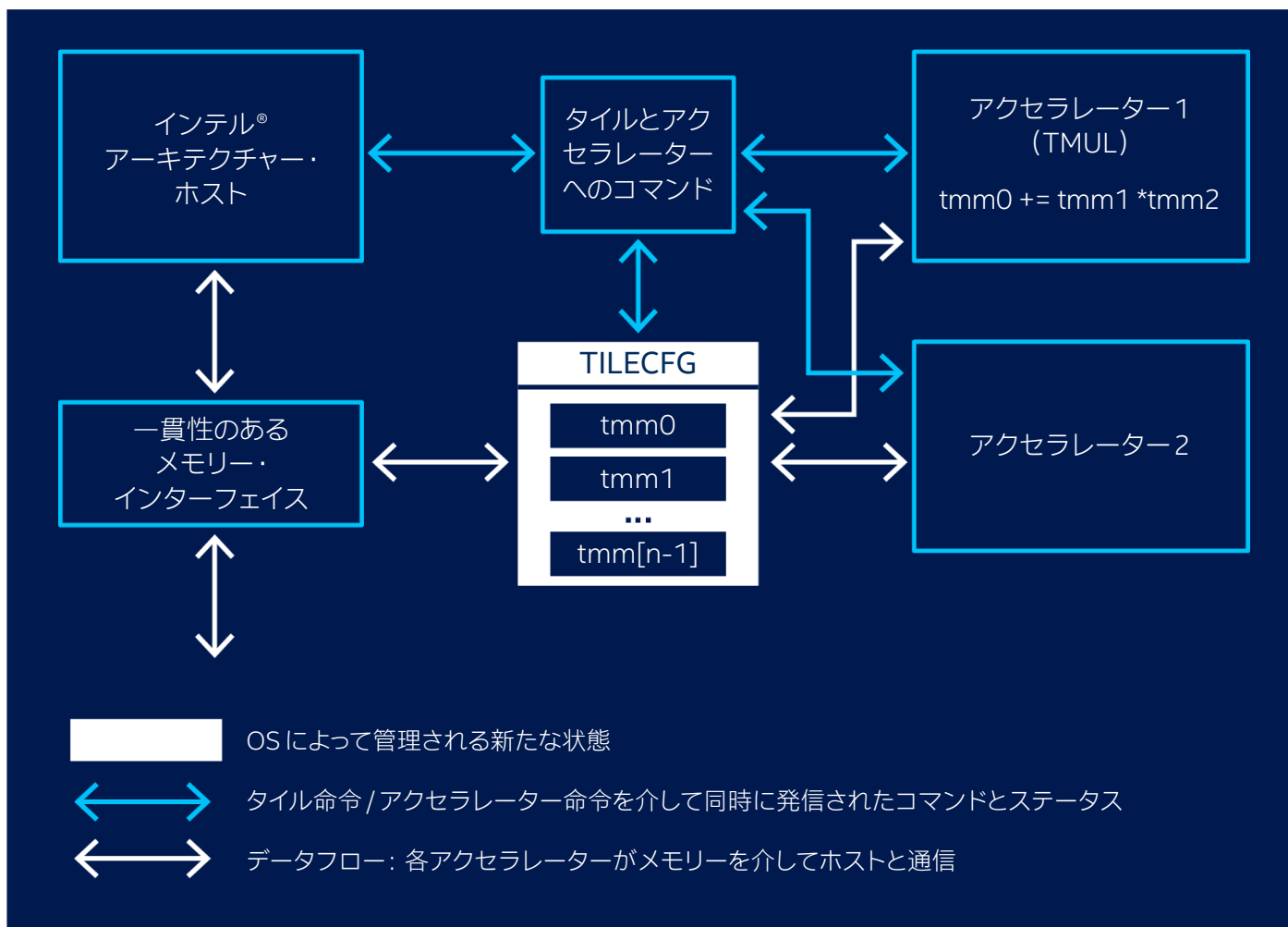


図5. インテル® AMXアーキテクチャー

行列積(行列乗算)を実行するには、非常に複雑な計算が必要です。行列乗算のオペレーションが必要になるたびに複雑なプログラムを書くのでは、コストパフォーマンスが低下します。

インテル® oneDNNを使えば、行列計算の実行をシンプルにできます。作業に必要なのは、一部の後処理でMatMulプリミティブを呼び出して、いくつかのパラメーター(m、n、k、ストライド、データアドレスなど)を渡すことです。以前のインテル®マス・カーネル・ライブラリー(インテル® MKL)と同じようにインテル® oneDNNは、タイル・レジスター・ファイルの設定、メモリーからのデータ読み込み、後処理の行列乗算の実行、結果をメモリーに格納する処理、タイル・レジスター・ファイルの解放など、残りの作業を完了させます。インテル® oneDNN経由でインテル® AMXを使う方法はプログラマーにとってわかりやすいので、プログラムをシンプルにできます。

インテル® AMXを内蔵した第4世代インテル® Xeon® スケーラブル・プロセッサで動作するBERTの全体的なフローチャートを図6に示します。

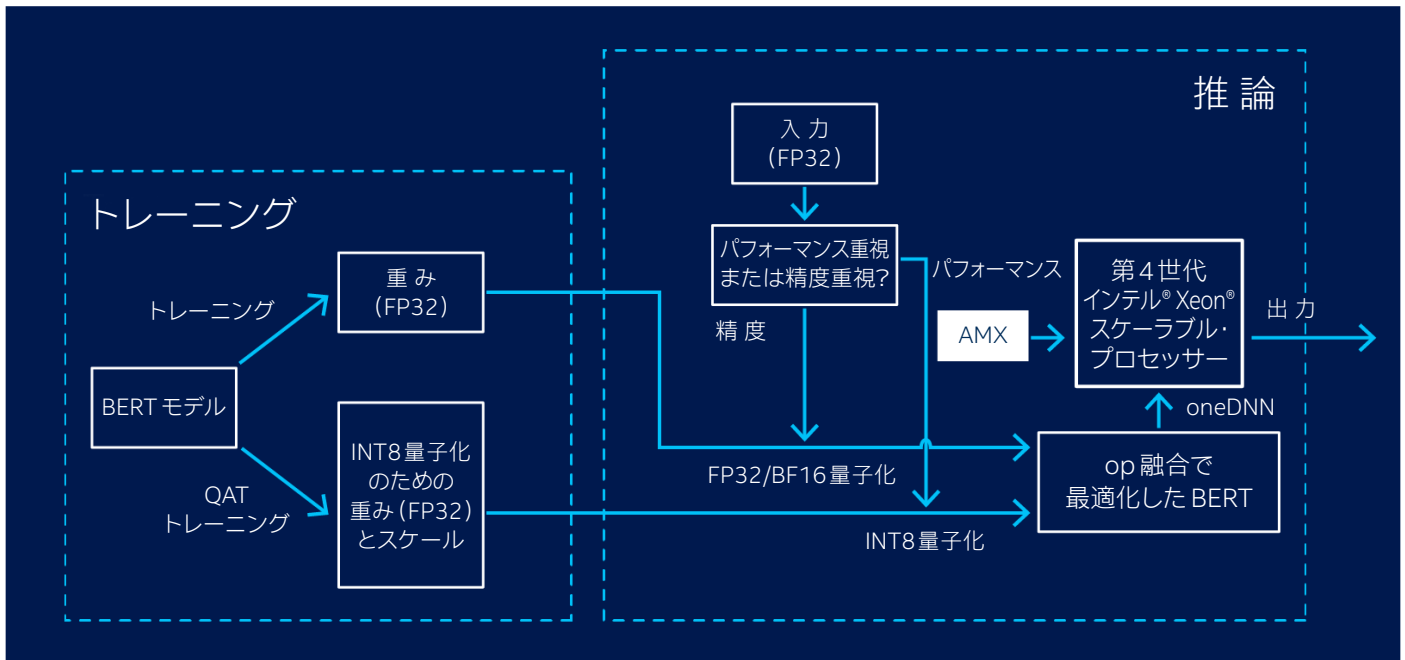


図6. 第4世代インテル® Xeon® スケーラブル・プロセッサ内蔵インテル® AMX使用時のBERT全体のフローチャート

パフォーマンスの比較

インテル® AMXによりBERTのパフォーマンスを大幅に向上できます。プロセッサの世代間でパフォーマンスの優位性を明らかにするには、さまざまなプラットフォームでパフォーマンスを比較する必要があります。インテルとの協業を通じて、テンセントのStarLake Labは、x86マイクロアーキテクチャに対する深い理解と、パフォーマンス・チューニングの経験を活用し、パフォーマンスの比較と最適化の作業に大きく貢献し、BERTテクノロジーを使用した際の第3世代インテル® Xeon® スケーラブル・プロセッサと第4世代インテル® Xeon® スケーラブル・プロセッサでのパフォーマンス検証で貴重なデータが得られました。

1つのソケットで複数の最適化インスタンスを実行し、各インスタンスのレイテンシーは一定のままでした。図7は、インテル® AMX内蔵の第4世代インテル® Xeon® スケーラブル・プロセッサを第3世代インテル® Xeon® スケーラブル・プロセッサと比較すると、INT8とBF16の双方において、それぞれ2.05倍および3.02倍とシステム・パフォーマンスが大幅に向上したことを示しています。^{1,2}

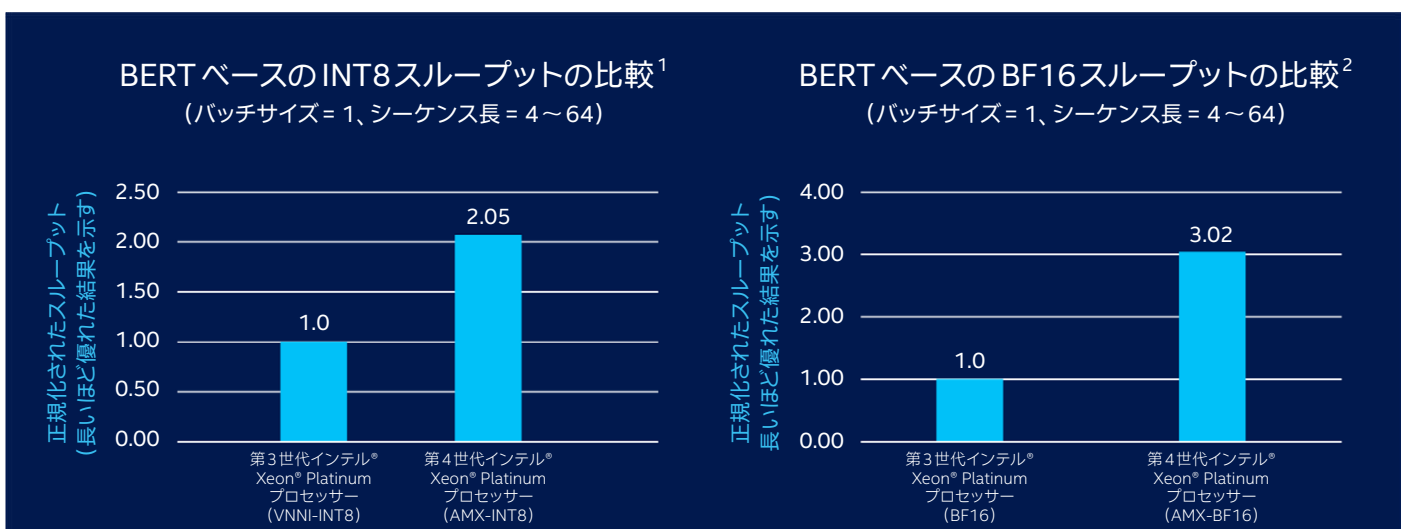


図7. インテル® AMX が内蔵された第4世代インテル® Xeon® スケーラブル・プロセッサと前世代でのBERTモデルのスループット比較

まとめ

インテル® AMX を内蔵した第4世代インテル® Xeon® スケーラブル・プロセッサは、BF16/INT8 TMUL コンピューティング・ユニットと、それに関連する命令により、行列乗算のパフォーマンスを大幅に向上させることが実証されました。インテルとテンセントは、インテル® AMX により、前世代と比較して BERT モデルのスループットが2倍から3倍向上することを証明しました。^{1,2} 現在、テンセントは最適化した BERT モデルを活用し、サービス体験の向上を実現しています。



1. BERT-base16 INT8スループットの比較

比較基準 (Baseline): Tencent TriRivers, BIOSバージョン1.08.00, OS: CentOS Linux release 8.5.2111, カーネル4.18.0-348.7.1.el8_5.x86_64, マイクロコード0xd000375, IRQ バランス: 有効, CPU: インテル® Xeon® Platinum プロセッサ, コアあたりのスレッド数2, ソケット数2, NUMA ノード数2, プリフェッチャー L2 HW, L2 Adj., DCU HW, DCU IP, ターボ有効, PPINs b59090a6f33f7966, b591426010edf86a, 電源およびパフォーマンスのポリシー: パフォーマンス, TDP 270ワット, 周波数ドライバー intel_pstate, 周波数ガバナンス performance, 周波数 (MHz) 2701, Max C-State 9, 搭載メモリー 960GB (15x64GB DDR4 3200MT/s [3200MT/s]), Huge Page サイズ2048kB, Transparent Huge Pages [always], 自動 NUMA バランシング: 有効, NIC: MT2892 Family [ConnectX-6 Dx] x2, デバイスx1, イーサネット・インターフェイスx1, ドライバースマリー: 111.8G INTEL SSDSCKHB12 x1, ワークロードおよびバージョン: INT8向け BERT 最適化, コンパイラー GCC 8.5, ライブラリー インテル® oneDNN-master-0721, テスト実施日: 2022年8月5日。

新規比較対象 (New): Intel Corporation ArcherCity, BIOSバージョンEGSDCRB1.SYS.0090.D03.2210040200, OS: CentOS Linux 8, カーネル5.16.0, マイクロコード0x2b0000c0, IRQ バランス: 有効, CPU: インテル® Xeon® Platinum プロセッサ, コアあたりのスレッド数2, ソケット数2, NUMA ノード数2, プリフェッチャー L2 HW, L2 Adj., DCU HW, DCU IP, ターボ有効, PPINs 31461920530bed98, 3143931fcf7f6036, 電源およびパフォーマンスのポリシー: パフォーマンス, TDP 350ワット, 周波数ドライバー intel_pstate, 周波数ガバナンス performance, 周波数 (MHz) 2494, Max C-State 9, 搭載メモリー 512GB (16x32GB<仕様外>4800MT/s [4800MT/s]), Huge Page サイズ2048kB, Transparent Huge Pages [always], 自動 NUMA バランシング: 有効, NIC: イーサネット・コントローラー I225-LM x1, QSFP 向けイーサネット・コントローラー E810-C x1, ドライバースマリー: 349.3G INTEL SSDPE21K375GA x1, 1.5T INTEL SSDPEDMD016T4 x1, 1.9T INTEL SSDPEKNW020T8 x1, ワークロードおよびバージョン: INT8向け BERT 最適化, コンパイラー GCC 8.5, ライブラリー インテル® oneDNN-master-0721, テスト実施日: 2022年10月19日。

2. BERT-base16 BF16スループットの比較

比較基準 (Baseline): Tencent QinghaiLake, BIOSバージョン1.02.00, OS: Red Hat Enterprise Linux 8.2 (Ootpa), カーネル4.18.0-193.el8.x86_64, マイクロコード0x7002502, IRQ バランス: 有効, CPU: インテル® Xeon® Platinum プロセッサ, コアあたりのスレッド数2, ソケット数4, NUMA ノード数4, プリフェッチャー L2 HW, L2 Adj., DCU HW, DCU IP, ターボ有効, PPINs 07ab90bc7f220116, 07be7bc039fedc8f, 07abcfbe92ff7f9b, 07aba8bfff01f278d, 電源およびパフォーマンスのポリシー: パフォーマンス, TDP 175ワット, 周波数ドライバー acpi-cpufreq, 周波数ガバナンス performance, 周波数 (MHz) 2695, Max C-State 9, 搭載メモリー 1536GB (24x64GB DDR4 3200MT/s [3200MT/s]), Huge Page サイズ2048kB, Transparent Huge Pages [always], 自動 NUMA バランシング: 有効, NIC: イーサネット・インターフェイスx1, ドライバースマリー: 447.1G SSSSTC ER2-GD480 x1, 894.3G Micron_5100_MTFD x1, ワークロードおよびバージョン: BF16向け BERT 最適化, コンパイラー GCC 8.3, ライブラリー インテル® oneDNN-master-0721, テスト実施日: 2022年8月8日。

新規比較対象 (New): Intel Corporation ArcherCity, BIOSバージョンEGSDCRB1.SYS.0090.D03.2210040200, OS: CentOS Linux 8, カーネル5.16.0, マイクロコード0x2b0000c0, IRQ バランス: 有効, CPU: インテル® Xeon® Platinum プロセッサ, コアあたりのスレッド数2, ソケット数2, NUMA ノード数2, プリフェッチャー L2 HW, L2 Adj., DCU HW, DCU IP, ターボ有効, PPINs 31461920530bed98, 3143931fcf7f6036, 電源およびパフォーマンスのポリシー: パフォーマンス, TDP 350ワット, 周波数ドライバー intel_pstate, 周波数ガバナンス performance, 周波数 (MHz) 2552, Max C-State 9, 搭載メモリー 512GB (16x32GB<仕様外>4800MT/s [4800MT/s]), Huge Page サイズ2048kB, Transparent Huge Pages [always], 自動 NUMA バランシング: 有効, NIC: イーサネット・コントローラー I225-LM x1, QSFP 向けイーサネット・コントローラー E810-C x1, ドライバースマリー: 349.3G INTEL SSDPE21K375GA x1, 1.5T INTEL SSDPEDMD016T4 x1, 1.9T INTEL SSDPEKNW020T8 x1, ワークロードおよびバージョン: BF16向け BERT 最適化, コンパイラー GCC 8.5, ライブラリー インテル® oneDNN-master-0721, テスト実施日: 2022年10月19日。

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、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。インテルは、明示されているか否かにかかわらず、いかなる保証もいたしません。ここにいる保証には、商品適格性、特定目的への適合性、および非侵害性の黙示の保証、ならびに履行の過程、取引の過程、または取引での使用から生じるあらゆる保証を含みますが、これらに限定されるわけではありません。

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。ほかの情報を参考にして、正確さを評価してください。

© 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。