# 時空間並列性を活用した Vision Transformer の FPGA 実装

## An FPGA Implementation of Vision Transformer Exploiting Spatial and Temporal Parallelism

〇松井友宏, ウィッデヤスーリヤハシタムトゥマラ, 田中大介, 張山昌論

Tomohiro Matsui, Hasitha Muthumala Waidyasooriya, Daisuke Tanaka, Masanori Hariyama

#### 東北大学

#### Tohoku University

キーワード: Vision Transformer (Vision Transformer), 再構成可能コンピューティング (reconfigurable computing), FPGA (Field Programmable Gate Array), 機械学習 (machine learning)

連絡先: 〒 980-8579 仙台市青葉区荒巻字青葉 6-3-09 東北大学大学院 情報科学研究科 張山・ウィッデヤスーリヤ研究室

松井友宏, Tel.: 022-795-7155 E-mail: matsui.tomohiro.t6@dc.tohoku.ac.jp

1. はじめに

Transformer アーキテクチャ<sup>1)</sup> は、機械学習分 野における画期的な進歩を示すものであり、そ の柔軟性と高い拡張性により、自然言語処理か らコンピュータビジョンに至るまで、幅広い応用 に活用されている。Vision Transformer (ViT) 2) は、その中でも特に強力な画像解析モデルと して登場した。ViT は自己注意 (self-attention) メカニズムを活用することで、視覚データ中の 長距離依存関係(long-range dependencies)を 効果的に捉え、畳み込みニューラルネットワーク (CNN) に代わる有力な手法として注目を集めて いる。Vision Transformer の計算で特徴的なの は、極めて高いデータ並列性 (data parallelism) である。大量のデータが外部メモリから繰り返 し読み込まれ、大規模な並列演算によって処理 され、その後再びメモリに書き戻される。この 一連の処理が多数の層で繰り返されるため、頻繁かつ膨大なメモリアクセスが発生する。GPUは高い外部メモリ帯域と多数の演算コアを備えており、このような処理に適しているが、一方で極端に多い外部メモリアクセスによって消費電力が大幅に増加するという問題がある。

エネルギー効率を改善する有望な方法の一つは、外部メモリアクセスの削減である。しかし、単にメモリアクセス回数を制限すると、並列計算に必要なデータ供給が不足し、演算性能のボトルネックを生じる可能性がある。したがって、計算の並列性を損なうことなく、消費電力を抑えるためには、より洗練されたアーキテクチャ的な工夫が必要である。この課題に対処するため、著者らは先行研究3)において、Multi-head attention機構における「時間的並列性(temporal parallelism)」を活用するアーキテクチャを提案した。本論文では、そのアプローチをVi-

sion Transformer 全体に拡張し、時間的並列性 と空間的並列性(spatial parallelism)の両方を 探究する。

本研究では、Vision Transformer の計算全体 を対象として、マクロブロックレベルからマイ クロ演算単位に至るまで、複数の粒度で空間的・ 時間的並列性を分析する。特に、中間結果を外 部メモリに書き戻すことなく再利用することで、 データ再利用効率を最大化する技術に焦点を当 てている。時間的並列性は主に、外部メモリア クセスを削減しつつ十分な並列処理量を維持す るために用いられる。一方で、空間的並列性は、 データフローの最適化やパイプラインのボトル ネック段階の処理時間短縮に活用される。この ように、メモリアクセスを積極的に削減しつつ 深くパイプライン化されたアーキテクチャを実 装するには、Field-Programmable Gate Array (FPGA) のような柔軟なハードウェアプラット フォームが最適である。FPGA はカスタム化さ れた省電力アクセラレータ設計を可能にする。 本研究では、提案アーキテクチャを Intel Agilex 7 FPGA 上に実装した。実験結果として、本設 計は高性能に最適化された20コアCPU実装を 上回る性能を示した。

### 2. 空間的および時間的並列性

#### 2.1 空間的および時間的並列性

Fig. 1 は、2 種類の基本的な並列処理手法、すなわち空間的並列性と時間的並列性を示している。Fig. 1(a) に示すように、空間的並列性では、複数の演算要素(Processing Elements: PEs)を用いて、大量のデータを同時並行的に処理する。それぞれの PE が異なるデータに対して並行して計算を行うことで、スループット(処理速度)を向上させることができる。多くの場合、同じハードウェアユニットを異なる種類の計算に時間的に再利用することもある。空間的並列性を効果的に活用するためには、システム

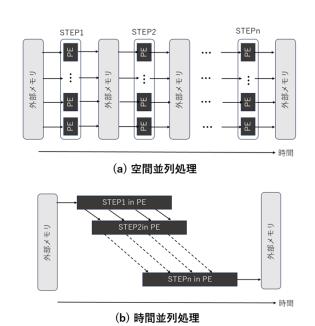


Fig. 1 空間的並列処理と時間的並列処理

が外部メモリへの並列アクセスを行う必要がある。しかし、これにより高いメモリアクセス帯域幅が要求される。同様に、中間結果も同時に外部メモリへ書き戻す必要がある。これは、内部メモリ資源が通常、全ての処理データを保持するには不十分であるためである。十分なメモリ帯域が確保できる場合、空間的並列性によって高い処理速度を実現することが可能である。しかし、外部メモリへの頻繁かつ大規模なアクセスは、非常に高い消費電力を引き起こす原因となる。特に、GPUを用いた処理ではこの問題が顕著である。

一方、時間的並列処理は、ある程度の並列性を維持しながらも、外部メモリアクセスを大幅に削減することに重点を置いている。Fig. 1(b)に示すように、異なる処理をパイプライン状に順番にスケジューリングし、各段階で得られた中間結果を外部メモリに書き戻さずに、次の段階で直接再利用する。これにより、演算処理が連続的かつ効率的に行われる。各ステップでは、外部メモリからアクセスまたはストリームされるデータ量はごく一部に限られ、そのデータが複数のPE間で再利用される。最終的な結果のみが、全ての計算パイプラインを通過した後で

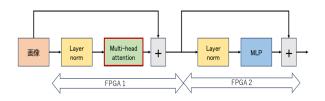


Fig. 2 Vision Transformer Encoder の処理 フロー

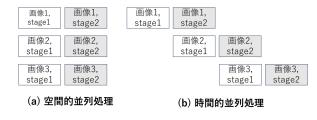


Fig. 3 マクロレベルの Vision Transformer 並 列処理

外部メモリに書き戻される。このアプローチに より、外部メモリアクセスを劇的に削減でき、 結果として消費電力を大幅に低減できる。しか し、こうした中間データの再利用を可能にする ための演算スケジューリングを正確に設計する ことは非常に難しい。特に、CPUや GPU のよ うな汎用プロセッサでは、スレッド間で頻繁な同 期が必要となり、その際に大きなオーバーヘッ ドが発生する。このオーバーヘッドが、せっか く削減したメモリアクセスによる性能向上効果 を相殺してしまうことがある。それに対して、 FPGA のようなハードウェアプラットフォーム は、時間的並列性の活用に非常に適している。 FPGA では、演算処理のスケジューリングを細 かく制御でき、同期のための複雑な制御を必要 とせずに中間データをシームレスに再利用する ことが可能である。この結果、外部メモリアク セスの大幅な削減が実現し、電力消費を大きく 低減できる。

### 2.2 Vision Transformer における並列 性の分析

ここでは、数百枚程度の画像からなるバッチ入 力を処理するケースを考える。Fig. 2は、Vision Transformer のマクロレベルでの計算フローを 示しており、大きく2つのステップに分けるこ とができる。第1ステップ:レイヤ正規化の後に Multi-Head Attention を実行する。第2ステッ プ:再びレイヤ正規化を行い、その後に多層パー セプトロン (Multi-Layer Perceptron: MLP) の 計算を行う。これらの2つのステップは、Fig. 3 に示すように、空間的並列処理または時間的並 列処理のいずれかで実行できる。Fig. 3(a) に示 すように、空間的並列性では、複数の埋め込み 画像を同時に処理し、それぞれのデータに並行 してアクセスする。この手法は、十分なメモリ帯 域幅が確保できる場合には、非常に高いスルー プットを実現できる。一方で、Fig. 3(b) に示す ように、時間的並列性では、2つのステップを異 なる埋め込み画像に対してパイプライン形式で 順次処理する。この方法では、中間データを再 利用することによって外部メモリアクセスを大 幅に削減できるが、得られる並列性の度合いは 限定的であり、大きな性能向上を実現するには 不十分な可能性がある。したがって、より高い 性能を達成するためには、各ステップの内部構 造をさらに詳細に分析し、より細かい粒度での 時間的並列処理を特定することが重要である。

#### 2.3 Multi-head-attention の並列処理

Fig. 4は、ステップ1の計算フローを示しており、このステップはレイヤ正規化、Multi-Head Attention、および加算の処理から構成されている。Multi-Head Attention は、複数の独立した「ヘッド(head)」から構成されており、それぞれのヘッドは入力データの異なる側面に焦点を当てて学習を行う。これにより、モデルは入力シーケンス内の異なる位置関係に同時に注目で

き、多様な関係性を効果的に捉えることができる。Multi-head-attention機構では、入力シーケンスはそれぞれのヘッドに対して、異なる線形変換を用いて複数のクエリ(Q)、キー(K)、およびバリュー(V)に射影される。各ヘッドは独立して、クエリとキーを比較してアテンションスコアを計算し、そのスコアをバリューに適用してコンテキストベクトルを生成する。

すべてのヘッドから得られたコンテキストベクトルは、その後連結され、線形変換を通して処理されることで、アテンションブロック全体の最終出力が得られる。以下の式は、自己注意機構における単一ヘッドの計算処理を表している。

Attention(
$$\mathbf{Q}, \mathbf{K}, \mathbf{V}$$
) = softmax( $\frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{d_k}}$ )  $\mathbf{V}$ 
(1

$$\mathbf{Q} = \mathbf{X}\mathbf{W}_Q, \quad \mathbf{K} = \mathbf{X}\mathbf{W}_K, \quad \mathbf{V} = \mathbf{X}\mathbf{W}_V$$
 (2)

ここで、X は位置エンコーディング (positional encoding)が付加された後の入力データ を表す。また、 $W_O$ ,  $W_K$  および  $W_V$  は、それ ぞれ クエリ (Query)、キー (Key)、および バリュー(Value) の射影に対応する重み行列 である。これらの演算は、バッチ内の各入力画 像に対して個別に行われ、すべてのアテンショ ンヘッドに対して繰り返し実行される。マルチ ヘッド・アテンション機構は、非常に高い空間 的並列性を備えている。すべてのアテンション ヘッドは同時に並列処理することが可能であり、 さらに各ヘッドの内部でも、バッチ内のすべて の入力データを並行して処理できる。加えて、 クエリ(Q)、キー(K)、バリュー(V)を計 算するための行列積演算も同時に実行可能であ る。それぞれの行列積演算自体も細かい粒度で 並列化が可能であるため、この段階は空間的並 列処理による高速化に非常に適している。実際、 複数のアテンションヘッドをまたいで複数の画 像を同時に処理することで、非常に高い並列性

を実現できる。しかしこの手法では、埋め込み 画像データと重み行列の両方を外部メモリから 同時にアクセスする必要があるため、メモリ帯 域幅の要求が非常に大きくなるという問題があ る。一方、時間的並列処理では、各ヘッドの処 理を順番に開始し、全体の5つの計算ブロック がパイプライン的に動作するようにスケジュー リングされる。この方法では、空間的並列処理 と比較して達成できる並列性の度合いは低いも のの、一度に1つのヘッド分のデータだけをメ モリから読み出せばよいため、外部メモリアク セスを大幅に削減できる。

総じて言えば、時間的並列性はスループットをある程度犠牲にする代わりに、より高い電力効率を実現する。一方で、空間的並列性は高い性能を実現できるが、より大きなメモリ帯域を必要とする。最適な設計は、利用可能なメモリ帯域幅と電力制約とのトレードオフによって決まる。

#### 2.4 MLP の並列処理

Fig. 5は、ステップ2におけるMLP計算の構造を示している。この計算は5つの計算ステージで構成されており、その中にはレイヤ正規化、行列積、およびGELU活性化関数が含まれている。Fig. 5に示す通り、5つのステージがパイプライン的に動作し、任意の時点でメモリから読み出す必要のあるデータ量を削減できる。この方法では、空間的並列処理と比較して並列度は低下するものの、一度に1枚の画像分のデータだけを扱えばよいため、外部メモリアクセスを大幅に削減できる。

## 3. Vision TransformerのFPGA アクセラレータ

第 II 章で説明したように、時間的並列性は、 データの再利用効率を高めるための重要な要素 であり、これが消費電力の削減に直結する。し

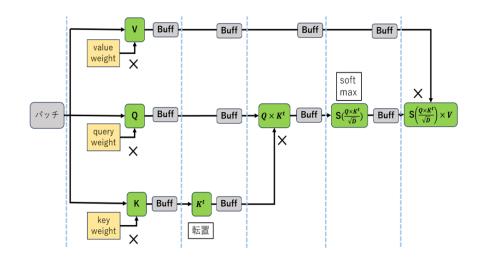


Fig. 4 Multi-head-attention の FPGA 処理の流れ

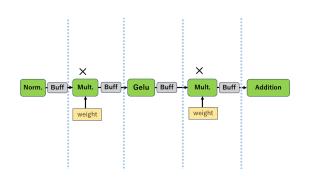


Fig. 5 MLP の FPGA 処理の流れ

たがって本研究では、可能な限り時間的並列性を優先的に利用する設計方針を採用している。一方で、空間的並列性は、マクロパイプライン内のボトルネックとなるステージの処理時間を短縮するために活用される。たとえば、マルチヘッド・アテンションにおける行列積は最も大きなボトルネックであるため、複数の行列積演算を同時に処理することで、空間的並列性を利用して処理時間を短縮している。

行列積の実装には、Fig. 6 に示すように 4) で用いられたシストリックアレイ構造を採用している。各行および列における PE の数は、処理速度や FPGA リソース制約など、設計要件に応じて調整可能である。すべての PE はパイプライン的に同時並行して動作するため、この構造では高いスループットを実現しつつ、アレイ全体でデータ再利用を最大化することができる。このアプローチは、計算効率を高めるだけでなく、

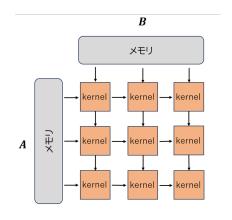


Fig. 6 MLP の FPGA 処理の流れ

外部メモリアクセスを削減して電力消費を抑えるという効果もある。同様の手法で、レイヤ正規化、行列転置、およびGELU活性化といった他のコンポーネントも、複数のマイクロステージに分割して実装している。このような分割設計により、中間データの効率的な再利用が可能となり、メモリアクセスの最小化とともに、性能および電力効率の向上が実現されている。

このように、時間的並列性と空間的並列性を 組み合わせたハイブリッド並列戦略により、低 消費電力を維持しながらさらなる高速化を実現 できる。このような大規模なマルチ FPGA アク セラレータシステムの設計と評価は、今後の課 題として残されている。しかしながら、本アー キテクチャは、高スループットかつ高エネルギー 効率な Vision Transformer 実装に向けた明確な 性能スケーリングの道筋を示している。

#### 3.1 評価

評価のために BittWare IA-840F FPGA ボードを使用した。このボードには Intel Agilex-7 AGF027 FPGA が搭載されている。FPGA カーネルのコンパイルには Intel OneAPI バージョン 2024.1 と Quartus Prime Pro 23.1 を使用した。比較対象として用いた CPU は、Intel Xeon Silver 4316(20 コア)である。CPU上の行列演算を高速化するために Intel MKL(Math Kernel Library)を使用した。Table 1にエンコーダのパラメータを示す。すべての計算は 単精度浮動小数点(FP32)で実行した。

Table 1 Vision Transformer 主要パラメータ

Image size	32
Number of channels	3
Patch size	8
MLP size	3072
Hidden size	768
Number of heads	12
Batch size	200

Table 2 FPGA Resource Utilization Summary

Resource	Utilization	Percentage
Registers	1,747,339	_
Logic	715,410	78%
DSP	3,563	42%
RAM blocks	3,465	26%
Memory	7.26 MB	22%
Clock frequency	292 MHz	_

マルチヘッドアテンション計算における FPGA リソース使用率を Table 2 にまとめる。主要な 制約要因はロジックブロックの使用量であり、 一方で浮動小数点演算に用いる DSP (Digital Signal Processor) の使用率は 42% に留まった。

また提案手法は、CPU 実装や従来 FPGA 実装<sup>3)</sup>と比較して高速の処理を実現した。この性

能向上は主に、本アーキテクチャで採用した高い並列処理度によるものである。現時点では、主なボトルネックは行列積にある。FPGAリソースの活用をさらに最適化することで、演算ユニット(PE)の数を増やし、さらなる性能向上が可能と考えられる。しかし、高い並列度を持つ完全な Vision Transformer 全体を単一 FPGA 上に実装するのは、現状ではリソース容量を超えている。そのため今後は、マルチ FPGA 構成による拡張を進め、並列性をより一層高めるとともに、Vision Transformer 全体のエンドツーエンド加速を実現する予定である。

### 4. まとめ

本研究では、Vision Transformer における時 間的並列性と空間的並列性の両方を分析し、外 部メモリアクセスと電力消費を削減するために 時間的並列性を優先した深いパイプライン構造 を提案した。提案アーキテクチャに基づき、マ ルチヘッドアテンション計算を FPGA 上に実 装し、その性能を評価した。実験結果から、提 案アクセラレータは並列マルチコア CPU 実装 よりも高い性能を示し、さらなる高速化の可能 性を有することが確認された。また、同一ハー ドウェアプラットフォーム上での従来 FPGA 実 装 $^{3)}$ よりも優れた性能を達成している。提案し たアーキテクチャは本質的にスケーラブルであ り、計算パイプラインを分割することで、複数 の FPGA に効率的に分散実行できる構造を備 えている。今後の研究では、この手法を Vision Transformer 全体に拡張し、マルチ FPGA 環境 で実装することを目指す。これにより、大規模 トランスフォーマモデルに対して、高スループッ トかつ省電力なアクセラレーションを実現する ことを目的とする。

## 参考文献

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017.
- 2) Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.
- 3) Hasitha Muthumala Waidyasooriya, Masanori Hariyama, and Daisuke Tanaka. FPGA-Based Deep-Pipelined Architecture for Vision Transformer's Multi-Head Attention. In 25th Workshop on Synthesis And System Integration of Mixed Information Technologies, pages 160–163, June 2024.
- 4) Hasitha Muthumala Waidyasooirya, Takuro Fukuda, and Masanori Hariyama. Scalable Architecture Targeting HBM-Based FPGAs for Complex Matrix Multiplication. In 23rd International Conference on Parallel and Distributed Computing, Applications and Technologies, December 2022.