深層学習コンパイラTVMの ベクトルマルチコア向けコード生成手法の検討

大西 文彬^{1,a)} 大髙 凌聖¹ 藤田 一輝¹ 末次 智貴¹ 川角 冬馬¹ 北村 俊明¹ 笠原 博徳¹ 木村 啓二¹

概要:自動運転車やスマートロボットなどの IoT デバイスのような組み込み機器においても,高度な判断 や制御をするために,深層学習による推論処理が広く利用されつつある.しかしながら,画像入力に対す る推論処理で多用される畳み込み演算は計算量が大きく,消費電力の増大とそれによる発熱量の増大を招 いてしまい,その結果組み込み機器に大容量のバッテリーが必要となったり,冷却装置のための大きなス ペースが必要となったりする懸念がある.これはスマートロボットなどにおいては行動自由度の低下や筐 体の大きさにも影響を与える.これに対して筆者等は,コンパイラ協調の OSCAR ベクトルマルチコアに よる,深層学習処理の高速かつ低消費電力実行の実現を目指している.本マルチコアでは,既存の多くの 学習モデルを利用すべく,そのコンパイルツールチェーンに深層学習コンパイラ TVM を取り入れる.本 稿では OSCAR 自動並列化コンパイラと TVM によるベクトルマルチコア用コンパイルツールチェーンの 有効性を示すべく,TVM によるベクトルマルチコア用コード生成手法を提案・実装する.さらに提案手 法を実装した TVM を OSCAR コンパイラおよび NEC を組み合わせ,ResNet の学習済み深層学習モデ ルを NEC のベクトルマルチコアである SX-Aurora TSUBASA 上で評価した.評価では上記 TVM の生 成コードを OSCAR コンパイラで並列化し,さらに NEC コンパイラによりベクトル化した.評価の結果, 提案手法実装前後の同コア実行時の実行時間を比較すると,1コア実行時に 14.0 倍,8コア時に 27.1 倍の 速度向上をそれぞれ得た.

Investigation of code generation techniques for vector multicore targeting using the deep learning compiler TVM

Fumiaki Onishi^{1,a)} Ryosei Otaka¹ Kazuki Fujita¹ Tomoki Suetsugu¹ Tohma Kawasumi¹ Toshiaki Kitamura¹ Hironori Kasahara¹ Keiji Kimura¹

1. はじめに

自動運転車やスマートロボット等,組み込み機器におい ても深層学習による画像認識や音声認識をはじめとする高 度な処理を要求するものが増えている.このような深層学 習の処理には,多数の配列の畳み込み演算などにより,大 きな計算コストがかかる.

このような処理を高速に行う方法として,多数の高速な 演算器や高バンド幅のメモリの導入が挙げられるが,処理

1 早稲田大学

の余裕度を確保するために対象となるアプリケーションの 要求性能を過剰に超過するハードウェア資源の導入は不要 な消費電力の増大やそれに伴う発熱を招き,大容量バッテ リーや大きな冷却装置の導入を必要とする. IoT デバイス を始めとする多くの組み込み機器では筐体の大きさや機 材設置の自由度の観点から,これらの大消費電力や高発熱 に伴うハードウェアフットプリントの増大は最小限とす べき項目である.すなわち,とりわけスマートロボット等 の高機能組み込み機器における高性能・高電力効率のコン ピュータシステムの導入が肝要である.

高性能・高電力効率のシステムを実現する方法として、

Waseda University

 $^{^{\}rm a)} \quad nueshitone@kasahara.cs.waseda.ac.jp$



図1 TVM と OSCAR を用いた深層学習モデルのコンパイルフロー

プログラムを低電圧かつ低動作周波数の複数コアを持つマ シン上で、電圧や動作周波数を適切に設定しながら並列に 実行することで、高動作周波数のマシン上で実行した場合 以上の性能を達成する方式が提案されている [1]. このよ うな低消費電力なコンピュータシステムを実現するために は、対象アプリケーションの挙動に応じて適切なコア数と 動作周波数及び電圧等を細かく制御する自動並列化コンパ イラとマルチコアアーキテクチャの協調が必要である.筆 者等は、各コアが、グローバル空間にマップされ分散共有 メモリとしても使用されるローカルメモリとベクトルアク セラレータを持つ OSCAR ベクトルマルチコアと、それに 対応するためのベクトル化や並列化、メモリ最適化、電力 最適化を行う OSCAR 自動並列化コンパイラを開発してい る [2]. ベクトルアーキテクチャでは、1つの命令で複数の データ要素を同時に処理できるため、OSCAR ベクトルマ ルチコアはデータ並列性の高い処理に対して特に優れた性 能を発揮する.

一方,深層学習のアプリケーション開発には Tensorflow[3] や PyTorch[4] などのフレームワークを用いること が多い. これらのフレームワークによる深層学習モデル を共通フォーマットである ONNX[5] を介して深層学習コ ンパイラ TVM[6] で処理することにより,OSCAR 自動並 列化コンパイラで並列化可能な C プログラムを生成でき る [2]. すなわち,フレームワークを用いた豊富な学習モデ ルを OSCAR 自動並列化コンパイラと OSCAR ベクトル マルチコアで利用可能となる. このような深層学習モデル のコンパイルフローを図 1 に示す.

このコンパイルフローでは、Tensorflow や PyTorch な どの様々なフレームワークを用いて学習された機械学習モ デルを、深層学習コンパイラ TVM を用いて逐次 C ソース コードに変換し、OSCAR 自動並列化コンパイラに入力す る. OSCAR 自動並列化コンパイラでは、ホスト CPU に 向けた並列化 C ソースコードと、ベクトルアクセラレータ に向けた C ソースコードを出力する.最後にこれらの C ソースコードを用いて実行バイナリを生成する.このよう に,TVM を用いることで機械学習モデルを OSCAR 自動 並列化コンパイラに入力することができる.

しかしながら,現在の TVM はベクトル化に適した再内 側ループが十分に長いコードを生成せず,その生成コード はベクトル化効率が低い.そこで本稿では,学習モデル中 のテンソルのレイアウトを変形し,その生成コードがベクト ル化に適した十分な長さの再内側ループとなるようなコー ド生成手法を提案する.さらに提案手法を実装した TVM を OSCAR 自動並列化コンパイラと組み合わせて NEC の ベクトルマルチコアである SX-Aurora TSUBASA[7]上で 評価した結果について報告する.本評価では図1に沿って, TVM が生成した学習モデルのコードを OSCAR 自動並列 化コンパイラの入力としマルチコア並列化を適用する.さ らに,並列化したコードを NEC コンパイラに入力しベク トル最適化を適用し実行バイナリを生成した.

以下,2節でOSCAR 自動並列化コンパイラについて,3 節で TVM の概要について説明し,4節でコード生成と高 速化のための実装について述べる.5節で性能評価の結果 について報告し,6節でまとめる.

2. OSCAR 自動並列化コンパイラ

OSCAR 自動並列化コンパイラは、C 言語あるいは Fortran で書かれた逐次プログラムをソースコードとし、 OpenMP と互換性のある OSCAR API[8] による並列化 コードを生成する Source to Source の自動並列化コンパイ ラである. OSCAR 自動並列化コンパイラでは、プログラ ムの基本ブロックやループ、関数呼び出しどうしを並列化 する粗粒度タスク並列化、ループイタレーションレベルで 並列化を行う中粒度並列化、基本ブロック内のステートメ ント間の並列化を行う近細粒度並列化の、3 種類の並列化 を階層的に行うマルチグレイン並列化を行う [9].

OSCAR コンパイラはフロントエンド、ミドルパス、バッ クエンドの3つのモジュールから構成される.まず、フロ ントエンドで入力の C 言語や Fortran で書かれたコード の構文解析や字句解析をし、その結果をコンパイラの中間 表現に変換する.次に、ミドルパスで、プログラムを粗粒 度タスク並列処理の単位である基本ブロック、繰り返しブ ロック、および関数呼び出しに分割しこれらをマクロタス ク(MT)とする.これら MT 間のコントロールフロー解析 とデータ依存解析の結果からマクロフローグラフ (MFG) を生成する.さらに MFG から MT 間のデータ依存とコン トロール依存を解析し、各 MT の最早実行可能条件、す なわち MT 間の並列性を解し、MT 間の並列性を図的に表 現するマクロタスクグラフ (MTG)を生成する.さらに、 キャッシュ最適化やメモリ最適化および各プロセッサへの コード割当をし,最適化された中間表現を生成する.最後 に,バックエンドでは,ミドルパスで生成された中間表現 を OSCAR API ディレクティブによって並列化された C 言語や Fortran のプログラムに変換する.

本稿においては,機械学習の学習済みモデルを TVM を 用いて C 言語の逐次ソースコードに変換した後,並列化 を行うために使用した. TVM と OSCAR を用いた機械学 習モデルのコンパイルフローは図 1 に示したとおりであ る.学習済みモデルを本稿で提案した拡張を施した TVM を用いて C ソースコードに変形し,OSCAR コンパイラを 通すことで並列化済み C ソースコードを出力し,NEC の SX-Aurora TSUBASA コンパイラで生成した実行バイナ リにより評価した.

3. TVM

TVM (Tensor Virtual Machine) は、オープンソースの 深層学習コンパイラおよび深層学習推論ランタイムである. 深層学習コンパイラは、学習済みの深層学習モデルを様々 なバックエンドに向けて最適化してコード生成する. TVM は複数の深層学習フレームワーク (PyTorch, TensorFlow, Keras など)をサポートしており、各フレームワークで構 築された学習モデルを取り込むことができる. また, TVM は複数のバックエンドを持つことで CPU, GPU, DSP な どの複数のターゲットハードウェア用の実行コードを生成 可能である. TVM では入力としたモデルを、グラフレベ ルとオペレーターレベルの最適化を行い、ターゲットハー ドウェア用バックエンドでハードウェアにあわせた効率的 なコードを生成する. グラフレベルの最適化では、複数の 演算をまとめて中間結果のメモリアクセスを削減するオ ペレータ融合や、静的に決定できる値を事前に計算して、 実行コストを節約する定数畳み込みなどの最適化を行う. オペレータレベルでの最適化では、ループ構造や並列パ ターンなどの実行詳細を決定するスケジュール変換などの 最適化を行う.本稿では、学習済みの機械学習モデルを、 OSCAR コンパイラに入力する C ソースコードへ変換する ために使用した.

TVM の全体の構成は図2のとおりである.

図2に示すように、TVM に入力されたモデルは、それ ぞれのモデルの形式に対応するフロントエンドによって、 Computational Graph と呼ばれるグラフレベルの中間表現 に変換される. Computational Graph の概要は図3の通り である.

Computational Graph は抽象構文木 (AST) の形で表さ れ、Python のコードや図3左下のようなテキストの形で 表される. TVM では、この中間表現に対して、畳み込み 演算と活性化関数の計算の融合などのグラフレベルでの 最適化を行う. その後、最適化済みのグラフを用いてオペ







⊠ 3 Computational Graph[11]



☑ 4 Tensor Expression[6]

レータレベルでの最適化を行ったうえで、コードを生成す る.オペレータレベルの最適化は Tensor Expression と呼 ばれる図4のような形式で表される中間表現を使って実行 される.

この中間表現は、テンソルの出力形状や演算ルールなど をもとに記述される.オペレータレベルの最適化には、多 重ループの複数のループの融合や、ループネストのソート などの最適化が含まれる.本稿では、この過程に変更を加 え、ベクトルアーキテクチャを持つ計算機上での機械学習 推論の高速化を図った.

4. 提案する TVM の拡張

本節では本稿で提案する、深層学習推論のベクトルアク



図 5 NCHW レイアウトの概要 [12]

セラレータによる高速化を可能とする TVM の拡張につい て述べる.

4.1 テンソルレイアウトの変換

一般にループを含むプログラムをベクトル化する場合, メモリアクセスが連続的であることや,ハードウェアのベ クトル演算の幅とループ内のデータの幅が一致していると いうことが重要である.メモリアクセスが連続的であるこ とでデータを連続的に読み込むことができ,メモリバンド 幅を最大限に活用できる.また,キャッシュラインの効率 的なプリフェッチやキャッシュヒットの増加にもつなが り,データの再利用性が高まるため,実行時間の短縮が期 待できる.このため,データへのアクセスが連続でかつベ クトル化するループ長を十分な長さになるようにテンソル のレイアウトを変更した.

TVM では、バックエンドが C 言語でコード生成する場合、テンソルを NCHWc のレイアウトで扱う. ここで N はバッチサイズ、C はチャネル数、H は高さ、W は幅で ある. NCHWc のレイアウトは、図 5 のような ONNX や PyTorch などで用いられる NCHW のレイアウトのチャネ ル内の複数の要素を1つの channel-wise としてさらに細か く並べる、図 6 のようなデータ構造となっている.

CNNにおける畳み込み演算では、出力チャネル(図5の C)方向にベクトル化できる[12].また、内部的な次元を 出力特徴マップに追加することで、連続したメモリアクセ スが可能となる[13].本稿では、ベクトル演算における同 時演算性能を大きくするために、ベクトル化される最内側 のループ長を十分な長さになるよう、channel-wiseの大き さを拡大した.channel-wiseの大きさを拡大することで、 図6のように、ベクトル化して同時計算できるテンソルの 数を増やし、推論処理の高速化を図った.評価に使用した NEC SX-Aurora TSUBASA のベクトル長が 256 であるこ とから、各演算の出力テンソルにおける channel-wise 方向 の大きさが最大 256 となるように変形した.さらに、それ



図 6 NCHWc レイアウトの概要

に伴いカーネルのレイアウト変形も行った.

4.2 出力コードの形状

実際にテンソルレイアウトの編集をする前の出力 C コードは図 7 のとおりであった.

このコードでは最内側のループ長が4となっている.次に、テンソルレイアウトの編集をした後の畳み込み層の出 力コードを図8に示す.

提案するテンソル変換の拡張を施した TVM による出力 コードでは,最内側のループなど,チャネル方向の内側の ループ長は 256 となっている.

5. 性能評価

5.1 評価方法

本節で提案するベクトルマルチコア用拡張を適用した TVM の性能評価について説明する.評価には、スキップ 接続を導入することで勾配消失問題を軽減し、深いネット ワークの構築を可能とした ResNet[14] より, ResNet50 の ImageNet1000 画像分類の学習済みのための ONNX 形式の モデルを用いた.本モデルを拡張前後の TVM を用いて C 言語の逐次ソースコードに変換し、それぞれ OSCAR 自動 並列化コンパイラを用いて並列化などの最適化を行った場 合と, OSCRA 自動並列化コンパイラを用いない場合と両 方のソースコードを評価環境 NEC SX-Aurora TSUBASA のネイティブコンパイラである NCC を用いてコンパイル し、ターゲット向けのバイナリコードを生成した. ター ゲットとなる NEC SX-Aurora の最大使用コア数は 8 コア であるため, OSCAR 自動並列化コンパイラを用いた評価 においては 1,2,4,8 コアに向けた並列化コードを生成し、そ の並列性を評価した.評価では画像一枚あたりの処理時間 に加えて、ベクトル長を測定した. OSCAR コンパイラを 用いた場合については、評価用画像を入力し出力コードを 実行した結果となる出力バイナリを、同様のモデルの実行 結果の出力と比較して、並列化後も等価なコードが実行さ

情報処理学会研究報告

IPSJ SIG Technical Report

```
int32_t ax0_ax1_fused_ax2_fused;
for (ax0_ax1_fused_ax2_fused = 0; ax0_ax1_fused_ax2_fused < 896; ++ax0_ax1_fused_ax2_fused) {</pre>
  static float conv2d_NCHWc_global[56];
      int32_t ic_outer;
      for (ic_outer = 0; ic_outer < 64; ++ic_outer) {</pre>
        ł
          int32_t kh;
          for (kh = 0; kh < 3; ++kh) {</pre>
            {
              int32_t kw;
              for (kw = 0; kw < 3; ++kw) {
                {
                  int32_t ic_inner;
                  for (ic_inner = 0; ic_inner < 4; ++ic_inner) {</pre>
                    {
                       int32_t ow_c_inner;
                       for (ow_c_inner = 0; ow_c_inner < 14; ++ow_c_inner) {</pre>
                         ſ
                           int32 t oc block c:
                           for (oc_block_c = 0; oc_block_c < 4; ++oc_block_c) {</pre>
                             conv2d_NCHWc_global[(((ow_c_inner * 4) + oc_block_c))] = (conv2d_NCHWc_global
                               [(((ow_c_inner * 4) + oc_block_c))] + (data_pad[(((((((ic_outer * 3364) + ((
                               ax0_ax1_fused_ax2_fused % 14) * 232)) + (kh * 116)) + (ow_c_inner * 8)) + (
                               kw * 4)) + ic_inner))] * arg_placeholder1[(((((((ax0_ax1_fused_ax2_fused /
                               14) * 9216) + (ic_outer * 144)) + (kh * 48)) + (kw * 16)) + (ic_inner * 4))
                               + oc_block_c))]));
                           }
                           . . .
```

図7 編集前の畳み込みループ

 表 1 評価対象		表 2	NEC SX-Aurora TS	SUBASA Vector Engine の構成
TVM	Ver.0.8.0[15]		Vector Core	8cores
モデル	ResNet50(ImageNet1000 画像分類)[16]		動作周波数(VE)	1.4GHz
入力画像	$224 \times 224 \times 3$ (cat.png) [17]		メモリ容量	24GB
			L1d Cache	32KiB / core

れていることを確認した.一方で NCC のみを用いて実行 した場合は,1コアでは正答となるものの,2コア以上の 並列化を行った場合には正答とならなかった.使用したモ デル,TVM,推論に用いた画像などは表1の通りである.

5.2 評価環境

ベクトルアクセラレータを搭載した NEC SX-Aurora TSUBASA 上で評価した. 評価に用いた NEC SX-Aurora TSUBASA の構成は表 2 の通りである.

尚,評価に用いたモデルは A100-1 で,Type 10C の Vector Element (VE)を1機搭載している.VE の動作周波 数は1.4GHz で,6MiB の LLC と 64 個のレジスタがある. また,最大のベクトル長は 256 要素である.評価に使用し たコンパイラ NCC のバージョンは 3.5.1 である.最適化 は O3 の最適化で,OSCAR コンパイラを使用しない場合 は-mparallel のオプションを付け並列化した一方,OSCAR コンパイラを使用して並列化した場合は OpenMP フラグ に対応するため,-fopenmp のオプションを付け,NCC の

Vector CoreScores動作周波数 (VE)1.4GHzメモリ容量24GBL1d Cache32KiB / coreL1i Cache32KiB / coreL2 Cache256KiB / coreL3 Cache (LLC)16MiBベクトルレジスタ数64 (プロセッサあたり)ベクトル長256 要素

並列化を無効化する-mno-parallel オプションをつけた.

5.3 テンソルレイアウト変更前の評価結果

まず,提案手法適用前の TVM を用いて ResNet50 の学 習済みモデルから得られたコードを OSCAR コンパイラを 用いて並列化して実行した.評価結果を表3に示す.実行 時間は1コアで約2秒あり,8コア時の速度向上率は1コ ア時に対して2.3 倍となった.

また,OSCAR コンパイラを用いずに NCC で並列化を 行った場合の評価結果は表4に示す.NCC で並列化を行っ た場合には、マルチコア実行時にも実行時間は高速化され なかった.

情報処理学会研究報告

IPSJ SIG Technical Report

int32_t oh;				
for (oh = 0; oh < 14; ++oh) {				
{				
<pre>int32_t ow_outer;</pre>				
<pre>for (ow_outer = 0; ow_outer < 14; ++ow_outer)</pre>				
{				
{				
int32_t kh;				
<pre>for (kh = 0; kh < 3; ++kh) {</pre>				
{				
int32_t kw;				
<pre>for (kw = 0; kw < 3; ++kw) {</pre>				
{				
<pre>int32_t ic_inner;</pre>				
<pre>for (ic_inner = 0; ic_inner < 256;</pre>				
++ic_inner) {				
{				
<pre>int32_t oc_block_c;</pre>				
<pre>for (oc_block_c = 0;</pre>				
oc_block_c < 256; ++				
<pre>oc_block_c) {</pre>				
<pre>conv2d_NCHWc_global[(</pre>				
oc_block_c)] = (
conv2d_NCHWc_global[(
oc_block_c)] + (data_pad				
[(((((oh * 14848) + (kh *				
7424)) + (ow_outer * 512)				
) + (kw * 256)) + ic_inner				
))] * arg_placeholder1				
[(((((kh * 196608) + (kw *				
65536)) + (ic_inner *				
256)) + oc_block_c))]));				
}				
••••				

図8 編集後の畳み込みループ

_	表 3	拡張前の TVM 出	力の C コードの実行時間
	コア数	実行時間 [ms]	速度向上率 (対1コア比)
	1	2,077.9	_
	2	$1,\!489.3$	1.4
	4	974.7	2.1
	8	901.4	2.3

表 4 拡張前の TVM 出力の C コードの実行時間 (NCC のみ)

コア数	実行時間 [ms]
1	2,181.3
2	2,181.3
4	2,181.2
8	2,181.2

5.4 テンソルレイアウト変更後の評価結果

次に,提案するテンソルレイアウトの変更の処理を実装 した TVM を用いて評価した.提案手法を実装した TVM から得られたテンソルレイアウト変更適用のコードを OS-CAR コンパイラを通さず,そのまま NCC を通して実行し た場合 (NCC のみ)と,同じコードを OSCAR コンパイラ



図 9 テンソルレイアウト変更後の実行時間

を用いて並列化したうえで NCC に通してコンパイルした 場合 (OSCAR+NCC) の評価結果は図 9 に示す. 図 9 よ り、テンソルレイアウトの変更を行う前の OSCAR コンパ イラを用いて並列化した場合の実行時間は1コアで2077.9 ミリ秒,8コアで901.4ミリ秒となっているのに対し、テ ンソルレイアウト変更適用のコードでは、OSCAR コンパ イラを用いて並列化した場合には1コアで151.8 ミリ秒, 8 コアで 42.0 ミリ秒, NCC のみを用いた場合でも1 コア で 367.5 秒,8 コアで 180.0 ミリ秒となっており,テンソ ルのレイアウト変更により、NCCのみの場合とOSCAR コンパイラを通した場合の両方で、シングルコア実行時と マルチコアでの実行時のいずれの場合も実行時間は短縮さ れた. テンソルのレイアウト変更前の OSCAR コンパイ ラを用いた場合の評価結果と比較すると、シングルコア時 は,NCCのみで5.7倍,OSCARコンパイラを用いた場合 は13.7 倍であり、さらに8コア時には21.5 倍とそれぞれ 性能向上が得られた.

次に並列化の状況を調べるため、シングルコア実行時に 対するマルチコア実行時の高速化率を調べた. この結果を 図 10 に示す. テンソルのレイアウト変更前の OSCAR コ ンパイラ使用時には1コア時に対する8コア時の高速化率 は 2.3 倍であり、テンソル変形後には NCC のみで 2.0 倍、 NCC+OSCAR で 3.6 倍となった. テンソルのレイアウト 変更前の並列化時の速度向上率は変更後の NCC のみのも のよりも大きくなっており、また、テンソルのレイアウト を変更したうえで OSCAR コンパイラを用いた場合は、並 列化時の速度向上率が他の 2 つよりも大きくなった.

推論の処理のための総実行命令数は図 11 に示す. テン ソルレイアウトを変更する前, NCC で並列化した場合(図 11 拡張前(NCC のみ))には 1 コアで 19.8 億命令, 8 コア で 89.2 億命令, OSCAR コンパイラを用いて並列化した場 合(図 11 拡張前(NCC+OSCAR))には 1 コアで 18.0 億 命令, 8 コアでは 42.4 億命令が実行されていたのに対し, テンソルのレイアウトを変更したコードで OSCAR コンパ

98.2

OSCAR+NCC



IPSJ SIG Technical Report





図 11 総実行命令数

イラを用いた場合(図 11 NCC+OSCAR)には,1コアで 2.4 億命令,8コアでも4.1 億命令となり,実行された命令 の数が大幅に減少していることがわかる.このように総実 行命令数が減少していることは,実行時間の削減の1つの 要因となっている.また,テンソルのレイアウト変更後の OSCAR有りの場合(図 11 NCC+OSCAR)と無しの場合 (図 11 NCCのみ)を比較すると,1コアの場合はOSCAR なしで2.3 億命令,OSCAR有りの場合は2.4 億命令なの に対し,8コアに並列化した場合は,OSCARなしの場合 は7.7 億命令,OSCAR有りの場合は4.1 億命令となって おり,OSCARコンパイラを用いて並列化した場合には命 令数の増加を抑えられていることがわかる.テンソルのレ イアウト変更によりベクトル命令数やベクトル長を大きく は大きくなり,これにより総実行命令数が削減できた.

テンソルのレイアウト変更前後のベクトル命令の割合を 12 に示す.変更する前は NCC のみ(図 12 拡張前(NCC のみ))で 77.0%, OSCAR コンパイラを用いた場合(図 12 拡張前(NCC+OSCAR))に 83.5%であったベクトル命令 の割合は,テンソル変形により, NCC のみ(図 12 NCC の み)で 97.4%, OSCAR コンパイラを併用した場合(図 12 NCC+OSCAR)には 98.2%となり,ベクトル命令の割合 が大きくなっていることがわかる.

図 12 全命令に占めるベクトル命令の割合

NCCOA

拡張前(NCCのみ) 拡張前(NCC+OSCAR)

83.5

また,平均のベクトル長を図 13 に示す. テンソルのレ イアウト変更前は NCC のみ (図 13 拡張前 (NCC のみ)) で 14.9, OSCAR コンパイラを用いた場合(図 12 拡張前 (NCC+OSCAR)) でに 15.3 だったベクトル長は, 変更後 には, NCC のみ (図 12 NCC のみ) で 166.8, OSCAR + NCC (図 12 NCC+OSCAR) で 160.2 となり、 テンソルの レイアウト変更によってベクトル長は 10 倍以上に大きく なっていることがわかる. さらに、ベクトル演算の総実行 時間は、レイアウト変更前に OSCAR コンパイラを用いる と 255.7 ミリ秒だったのに対し、変更後には 76.5 ミリ秒と なっており、ベクトルロード命令の実行時間も1701.5ミリ 秒から 59.5 ミリ秒に短縮されている.ベクトル長が大き くなることで、ベクトル命令の実行回数が減少しているこ とがわかる.NCC のみの場合と比較してもベクトルロー ド命令の実行時間は 255.7 ミリ秒から 59.5 ミリ秒に減少し ており、メモリへのアクセス回数が最適化により減ったこ とで実行時間の短縮に繋がった.

以上のように、ベクトル長が大きくなり、かつベクトル 命令の割合も大きくなったことで、実行時間の大幅な削減 に繋がった.

6. まとめ

100.0

90.0

80.0

70.0

60.0

50.0

30.0

20.0

10.0

▶ (%) (%)

↓ ______ 40.0 77.0

本稿では,深層学習コンパイラ TVM を用いたベクトル マルチコア向けコード生成手法を提案した.深層学習の学 習済みモデルを深層学習コンパイラ TVM を用いて C ソー スコードに変換し,OSCAR コンパイラを用いて並列化な どの最適化を行ったうえで実行バイナリを生成した.学習 モデルを C ソースコードに変換する過程で,ベクトル化の 効果を意識して,最内側のループ長が大きくなるようにテ ンソルレイアウトを変形した.この結果,テンソルレイア ウトの変形前の同コア比での推論処理の実行時間は,1コ アでは 13.5 倍,8コアでは 21.5 倍にそれぞれ高速化され



図 13 平均ベクトル長とベクトル命令実行時間

た.また,テンソルレイアウトを変形した C ソースコード を用いた場合,OSCAR コンパイラを通したときと NCC のみでコンパイルしたときとを比較すると,1コア実行時 と 8 コア実行時との速度比は 2.0 から 3.6 まで大きくなっ た.以上から,深層学習モデルを,TVM を用いて C ソー スコードに変換し,これを OSCAR コンパイラで解析する ことで,自動で並列化を行うことが確認できた.さらに, TVM が生成する C ソースコードのテンソルレイアウトが, 畳み込みの最内側のループ長が大きくなるように変換する ことで,ベクトルプロセッサの実行効率が向上し,より高 速な実行が可能となることを確認した.

謝辞 本研究の成果の一部は JST【ムーンショット型研 究開発事業】【JPMJMS2031】の支援を受けたものです.

参考文献

- Hirano, T., Yamamoto, H., Iizuka, S., Muto, K., Goto, T., Wake, T., Mikami, H., Takamura, M., Kimura, K. and Kasahara, H.: Evaluation of Automatic Power Reduction with OSCAR Compiler on Intel Haswell and ARM Cortex-A9 Multicores, Proceedings of the 27th International Workshop on Languages and Compilers for Parallel Computing (LCPC), LCPC (2014).
- [2] Kasahara, H., Kimura, K., Kitamura, T., Mikami, H., Morita, K., Fujita, K., Yamamoto, K. and Kawasumi, T.: OSCAR Parallelizing and Power Reducing Compiler and API for Heterogeneous Multicores: (Invited Paper), 2021 IEEE/ACM Programming Environments for Heterogeneous Computing (PEHC), pp. 10–19 (online), DOI: 10.1109/PEHC54839.2021.00007 (2021).
- [3] Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M. et al.: TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems, arXiv preprint arXiv:1603.04467 (2016).
- [4] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L. et al.: PyTorch: An Imperative Style, High-Performance Deep Learning Library, arXiv preprint arXiv:1912.01703 (2019).
- [5] Bai, J., L. F. Z. K. e. a.: ONNX: Open Neural Network Exchange, GitHub (online).

- [6] Chen, T., Moreau, T., Jiang, Z., Zheng, L., Yan, E., Cowan, M., Shen, H., Wang, L., Hu, Y., Ceze, L., Guestrin, C. and Krishnamurthy, A.: TVM: An Automated End-to-End Optimizing Compiler for Deep Learning, arXiv preprint arXiv:1802.04799 (2018).
- NEC Corporation: NEC SX-Aurora TSUBASA Architecture - NEC Vector Engine Processor, available from (https://www.nec.com/en/global/solutions/ hpc/sx/architecture.html) (accessed 2023-06-25).
- [8] Kimura, K., Mase, M., Mikami, H., Miyamoto, T., Shirako, J. and Kasahara, H.: OSCAR API for Real-Time Low-Power Multicores and Its Performance on Multicores and SMP Servers, *Languages and Compilers for Parallel Computing* (Gao, G. R., Pollock, L. L., Cavazos, J. and Li, X., eds.) (2010).
- [9] Kasahara, H., Honda, H., Mogi, A., Ogura, A., Fujiwara, K. and Narita, S.: A Multi-grain Parallelizing Compilation Scheme for OSCAR (Optimzally Scheduled Advanced Multiprocessor), *Proc. 4th Intl. Workshop on LCPC*, pp. 283–297 (1991).
- [10] TVM: Design and Architecture (online), available from (https://tvm.apache.org/docs/arch/index.html) (accessed 2023-06-25).
- TVM: Introduction to Relay IR (online), available from (https://tvm.apache.org/docs/arch/relay_intro. html) (accessed 2023-06-25).
- [12] 大髙凌聖,小池穂乃花,磯野立成,川角冬馬,北村俊明, 見神広紀,納富 昭,木村貞弘,木村啓二,笠原博徳: 各コアがローカルメモリを持つ組み込みベクトルマルチ コアでの畳み込み層演算の評価,2023-EMB-62(32), pp. 1-6 (2023).
- [13] Das, D., Avancha, S., Mudigere, D., Vaidynathan, K., Sridharan, S., Kalamkar, D., Kaul, B. and Dubey, P.: Distributed Deep Learning Using Synchronous Stochastic Gradient Descent, arXiv preprint arXiv:1602.06709 (2016).
- [14] He, K., Zhang, X., Ren, S. and Sun, J.: Deep Residual Learning for Image Recognition, CoRR, Vol. abs/1512.03 (2016).
- [15] APACHE: TVM: Open deep learning compiler stack for cpu, gpu and specialized accelerators, GitHub (online), available from (https://github.com/apache/ tvm/tree/v0.8) (accessed 2023-06-25).
- [16] ONNX: ONNX Models: ResNet model, GitHub (online), available from (https://github.com/onnx/models/ tree/main/vision/classification/resnet/model/ resnet50-v2-7.onnx) (accessed 2023-06-25).
- [17] DMLC: MXNet.js: Javascript Package for Deep Learning in Browser, GitHub (online), available from (https:// github.com/dmlc/mxnet.js/blob/master/data) (accessed 2023-06-25).

正誤表

下記の箇所に誤りがございました.お詫びして訂正いたします.

訂正箇所	誤	E
1ページ	1 コア実行時に 14.0 倍, 8 コア時に	1 コア実行時に 13.7 倍, 8 コア時に
14 行目	27.1 倍の速度向上をそれぞれ得た.	21.5 倍の速度向上をそれぞれ得た.