

GPU を用いた行列-行列積の実装と性能評価

榎田 匠[†] 田中 輝雄[†] 藤井 昭宏[†]
工学院大学[†]

1 はじめに

高い理論演算性能を持つ GPU (Graphics Processing Unit) で性能を最大限に発揮させるプログラムを作成するには、多数のスレッドによる超並列やメモリへのコアレスアクセスなど GPU プログラミングの技法を駆使する必要がある[1].

本研究では、行列-行列積の実装と性能評価を行った。演算効率を極限まで高めた数値計算ライブラリ cuBLAS[2]に対し、CUDA を用いた実装で高い性能を得ることを目標とした。性能を引き出すため、効率的な Shared Memory の利用しかたに着目した。Shared Memory 使用量を変更することで、問題サイズごとに SM (Streaming Multiprocessor) の稼働率を最大限に上げることを試みた。

2 GPU ハードウェア

本研究では、産業技術総合研究所の ABCI[3]に搭載の NVIDIA 社製 GPU Tesla V100[4]を用いた。GPU の概要図を図 1、Tesla V100 の仕様を表 1 に示す。GPU は複数の SM を搭載し SM には演算器 (CUDA Core) や Shared Memory などが搭載されている。Shared Memory と Global Memory は階層構造であり、Shared Memory のほうがより高速である。

Tesla V100 は SM 数 80、SM の倍精度演算器 (FP64 Core) 数 32 である。倍精度演算の理論演算性能は約 7.8TFLOPS である。

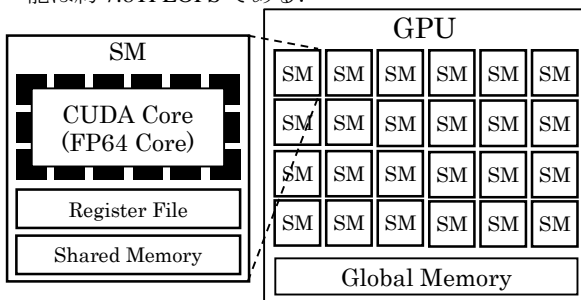


図 1 GPU の概要図

表 1 Tesla V100 の仕様

SMs	80
FP64 Cores / SM	32
GPU Boost Clock	1530 MHz
Peak FP64 TFLOPS	7.8
Memory Size	16GB
Shared Memory Size / SM	96 KB
Register File Size / SM	256 KB

Implementation and Performance Evaluation of Matrix-Matrix Multiplication using GPU

Takumi Sakakida[†], Teruo Tanaka[†] and Akihiro Fujii[†]
Kogakuin University[†]

表 2 CUDA の Thread 構造と GPU 構造

Thread 構造	GPU 構造
Thread	CUDA Core
Thread Block (Thread の集合)	SM (CUDA Core の集合)
Grid (Thread Block の集合)	GPU (SM の集合)

3 CUDA のマルチスレッディング

プログラム作成には NVIDIA 社の CUDA を使用した。表 2 に CUDA プログラミングの Thread 構造と GPU のハードウェア構造を示す。両者の構造が対になっていることがわかる。そのため、GPU 構造を考慮し Thread を生成すれば、より高い性能を得られる。そこで本研究では SM と Thread Block の関係について詳細に評価を行った。

4 実効 Thread Block 数

各 SM への処理の振り分けは、Thread Block 単位で行われる。SM の稼働上限は Thread 数 (2048)、Register 使用量 (256KB)、Shared Memory 使用量 (96KB) の 3 要素であり、これに収まる数だけ Thread Block が同時に稼働できる。3 要素すべてを半分使用する Thread Block を多数生成した場合には、各 SM は同時に 2 Thread Block が稼働可能となる。このように実際に同時稼働可能な Thread Block 数を実効 Thread Block 数と呼ぶ[5].

5 行列-行列積の実装

行列-行列積の実装には Nath らの実装[6]と数値計算ライブラリ cuBLAS の Thread 構成などを参考にした。図 2 に実装の概要図を示す。行列-行列積を $C=A \times B$ としたとき、 C を 128 行 64 列の部分行列に分割し、各 Thread Block がそれを計算する。

Thread Block は A を s 列、 B を s 行ずつ Shared Memory に格納する。そして 128 行 64 列の部分行列をさらに 8 行 8 列の部分行列に分割し、各 Thread が外積形式で計算する。

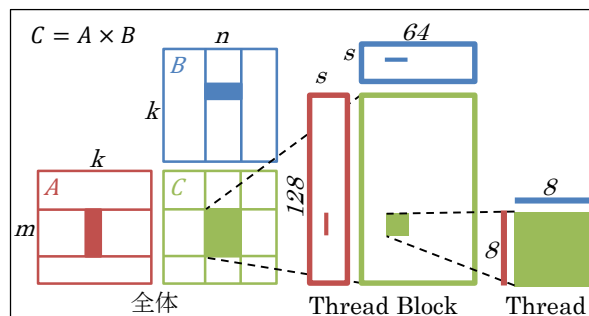


図 2 行列-行列積の実装

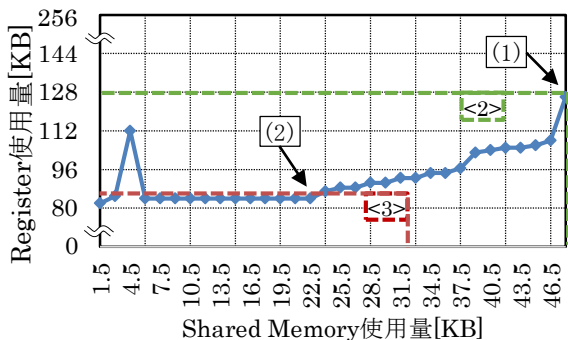


図3 Shared Memory と Register の使用量

表3 計測時の Thread Block 設定

	設定値		プロファイラから得た結果	
	Thread	Shared Memory	Register	実効 Thread Block 数
(1)	128	48KB	126KB	2
(2)	128	22.5KB	84KB	3

6 予備実験: Register 使用量の取得

実効 Thread Block 数を決める要素のうち、Register 使用量についてはプロファイラで取得した。図3に Shared Memory と Register の使用量を示す。また、実効 Thread Block 数の境界線を図3内に破線で示す。図3の領域<2>は Shared Memory と Register の使用量が4節に示した上限の半分以下のため実効 Thread Block 数は2である。同様に領域<3>は実効 Thread Block 数が3であり、本実装では実効 Thread Block 数を2と3に設定可能である。

7 評価実験

Thread Block 内の計算量は一定のまま Thread Block 数を1刻みで増加させて計測を行うため、行列サイズ M を128、 K を4096で固定し、 N を64から47360まで64ずつ増加させた。

表3に今回計測を行った2種類の Thread Block 設定について示す。図3の(1)と(2)の設定とした理由は、各実効 Thread Block 数の領域内で、最も Shared Memory、Register を使用するためである。

計測結果を図4に示す。横軸は生成されている Thread Block 数を表し、数が大きいほど問題の行列サイズが大きくなる。縦軸は実測性能を表す。グラフ(1)の時、実効 Thread Block 数は2であり、80SM 全てが2Thread Block 稼働可能な160Thread Block の倍数を頂点とする鋸波が確認できる。同様にグラフ(2)では240Thread Block の倍数を頂点とする鋸波で確認できる。この結果から、生成する Thread Block 数/SM 数の倍数で実効 Thread Block 数を指定すれば、高い性能を得られるといえる。今回のように限られた実効 Thread Block 数でしか稼働できない時も、それらのうち性能の良いほうを選択し実行することで性能低下を抑制できる。さらに、複数の実効 Thread Block 数での稼働を組み合わせることで高い性能を得られる。例として生成される Thread Block 数が560の時、実効 Thread Block 数2で320Thread Block を、実効 Thread Block 数3で240Thread Block を分割して実行を行った。結果を図4の菱形に示す。分割せず単一の

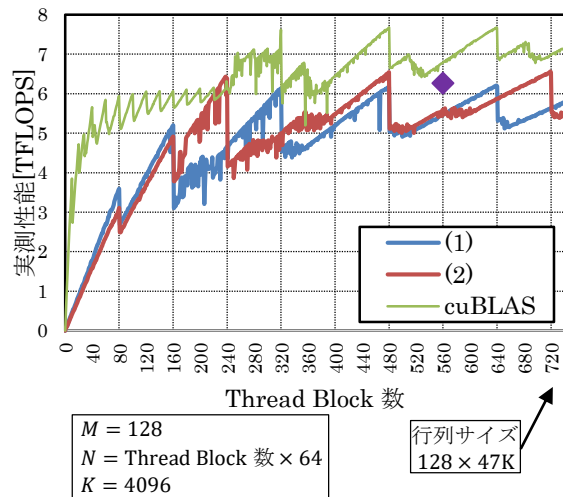


図4 実効 Thread Block 数による性能の差

実効 Thread Block 数で実行するよりも性能が向上した。また、cuBLASについても同じ行列サイズで計測を行った。240Thread Block 周辺ではほぼ同程度の性能を得られ、9割を超える性能を得られる Thread Block 数も確認できた。

8 おわりに

本研究では、GPU上で行列-行列積を実装し、実効 Thread Block 数を変更することによる性能の違いを評価した。生成される Thread Block に対して適切な実効 Thread Block 数を指定することでプログラムの性能を向上させることができた。その結果、特定の問題サイズではcuBLASに対して同程度の性能を得られた。

今回は実効 Thread Block 数2と3の場合のみ計測することができた。その他の実効 Thread Block 数でも計測と評価を行うのが今後の課題として挙げられる。

謝辞 本研究の一部はJSPS 科研費 JP18K19782, JP18K11340, JP15K15998 の助成を受けたものです。

参考文献

- [1] Hisa Ando, GPUを支える技術 技術評論社(2017)
- [2] NVIDIA: cuBLAS, <https://developer.nvidia.com/cublas> (2019/1/11 アクセス)
- [3] 産業技術総合研究所: ABCI, <https://abci.ai/> (2019/1/11 アクセス)
- [4] NVIDIA TESLA V100 GPU ARCHITECTURE, <http://images.nvidia.com/content/volta-architecture/pdf/volta-architecture-whitepaper.pdf> (2019/1/11 アクセス)
- [5] Best Practices Guide, <https://docs.nvidia.com/cuda/archive/9.2/cuda-c-best-practices-guide/> (2019/1/11 アクセス)
- [6] R. Nath, S. Tomov, and J. Dongarra, "An Improved Magma Gemm For Fermi Graphics Processing Units.", *The International Journal of High Performance Computing Applications*, 24(4), pp. 511–515 (2010)