

科学研究費助成事業（科学研究費補助金）研究成果報告書

平成 24 年 4 月 1 日現在

機関番号：12608

研究種目：若手研究（A）

研究期間：2010～2011

課題番号：22680002

研究課題名（和文） GPUによるFFT計算の自動チューニング手法の研究

研究課題名（英文） Auto-tuning FFT using GPU

研究代表者

額田 彰 (NUKADA AKIRA)

東京工業大学・学術国際情報センター・産学官連携研究員

研究者番号：40545688

研究成果の概要（和文）：NVIDIA社製のCUDA対応GPU向けの自動チューニングFFTライブラリであるNukadaFFTライブラリを開発した。その性能は多くの場合にNVIDIA社のCUFFTライブラリを上回る。また複数GPU版についても複数GPUを搭載するシングルノードと複数ノード版を実装し、さらなる速度向上を達成した。

研究成果の概要（英文）：We have developed the NukadaFFT library, which is an auto-tuning FFT library for NVIDIA CUDA GPUs. It outperforms NVIDIA's CUFFT library in many cases. We also implemented a multi-GPU version for both single-node with multi-GPUs and multi-node, and achieved further speed-up.

交付決定額

（金額単位：円）

	直接経費	間接経費	合計
2010年度	5,000,000	1,500,000	6,500,000
2011年度	4,000,000	1,200,000	5,200,000
年度			
年度			
年度			
総計	9,000,000	2,700,000	11,700,000

研究分野：総合領域

科研費の分科・細目：情報学・ソフトウェア

キーワード：ソフトウェア工学

1. 研究開始当初の背景

GPUの利用による各種数値計算の高速化は今日既に十分に普及してきたといえる。従来のアクセラレータは多数の演算コアを搭載したもので演算量の極めて多い計算において高速化を実現するものであったが、GPUはメモリバンド幅においても優れたアクセラレータであるためより広範囲な計算に対応することが可能である。高速フーリエ変換(FFT)は信号処理、画像解析、マルチメディア処理のような小規模なものからスーパーコンピュータを利用するような大規模なシミュレーションまで幅広く用いられる重要

な計算であるが、演算性能だけでなくメモリアクセス性能の要求も高いため従来のアクセラレータでは対応できなかった。

2. 研究の目的

高速なFFTライブラリは常に多くのアプリケーション開発者等から必要とされている。特にGPU製品はおよそ年に1回新しいアーキテクチャ世代の製品が発売されるという状況であるため、新アーキテクチャに即対応することが可能である自動チューニング機能を搭載するライブラリを開発する。

3. 研究の方法

GPUによるFFT計算の高速化は単純ではない。GPUのメモリアクセスは汎用CPUのメモリやベクトルプロセッサのそれとも異なる性質をもち、理論ピーク性能に近いメモリアクセス性能を得るためにはメモリアクセスの局所性や各スレッド間でアクセス箇所を調整することが不可欠であり、またキャッシュメモリの容量に限られるGPUではキャッシュを搭載するCPU用のアルゴリズムとベクトルプロセッサ向けのアルゴリズムのどちらでもない新たな計算手法を開発する必要がある。

またFFTの演算量を $O(N \log N)$ に削減するアルゴリズムはデータサイズに依存するため、ライブラリ化する場合にはデータサイズ毎に自動最適化を行う仕組みが不可欠であった。

また現在主流であるCUDAはNVIDIA社のGPUにのみ対応する。より広範なアーキテクチャに対応するOpenCLでの実装も検討する。

4. 研究成果

GPUにおける自動最適化のパラメータは大きく分けて以下の3点に絞られる。

(1) 入力サイズから基底への因数分解方法、及び各基底の計算順序

これはCPUでも考慮しなければならないパラメータであるが、GPUの場合には演算数だけでなくスレッド数やレジスタ数、ロード・ストア命令数などにも影響するため最適なパラメータはGPUアーキテクチャに大きく依存する。

(2) スレッド数

GPUは多数のスレッドが交代で命令を実行する仕組みになっており、メモリアクセス時にはデータが実際に利用されるまでは他のスレッドが命令実行を行うことでメモリアクセスレイテンシを隠ぺいすることが可能である。FFTの計算ではメモリアクセス効率が重要で、最大のデータ転送スループットを得られるように最適なスレッド数を選択する必要がある。

(3) シェアードメモリアクセスパターン

1組のFFTの計算をスレッドブロックが担当する場合、スレッド間で頻りにデータの交換が必要である。CUDAではスレッド間のデータ交換にはオンチップにある高速にアクセス可能なシェアードメモリを用いる。シェアードメモリは16個や32個などのバンクに分かれており、別々のバンクは同時にアクセス可能である。FFTにおけるスレッド間のデータ交換パターンに対してバンクコンフリクトが起らないように最適化するべく自動的にパディング挿入パターンを決定する。

これらの3つのパラメータを基本的には網羅的探索として全ての組み合わせで試行し、性能が最高となるパラメータを得る。

ライブラリとして実装する場合に、全ての入力サイズ及び各パラメータに対応するコードを予め生成しておくことは現実的ではない。そのため与えられた入力サイズに対して全パラメータに対応するコードを生成し、試行して最適なものを選択する方式を採用する他にない。全パラメータの試行にはかなり時間を要する。試行には①カーネル生成、②コンパイル、③ロード、④デバイスへのデータ転送、⑤カーネル実行、⑥ホストへのデータ転送、⑦結果のチェックの7つのタスクで構成される。特にホストCPU側で行われる②、④、⑥、⑦のタスクが時間がかかる。中でも②が一番時間がかかるため、C for CUDA言語ではなくよりローレベルなCUDA PTX言語でカーネルをメモリ上に生成してコンパイルすることで大幅な時間短縮を実現した。

OpenCL環境ではCUDAとほぼ同じ機能が提供されており、特にカーネル部分のコードは似ている。唯一の差異はFFTで利用する三角関数のテーブルの扱いで、CUDAではテクスチャメモリを利用していたがOpenCLではコンスタントメモリを使用した。AMD製RADEON HD 7970を用いてCUDA+NVIDIA製GPUを超える性能を実現することができた。一方、性能のポータビリティはOpenCLでは確保されていない。NVIDIA製GPUではCUDA版と比べるとOpenCL版の性能は現時点ではかなり劣る。今後ドライバの成熟を待てばこの差は縮まることが期待される。

さらに複数GPUへの対応も進めている。特にGPUのデバイスメモリ容量は限られるため、実アプリケーションでは容量を確保するために複数GPUを利用するケースも少なくない。GPU間の転送はPCI-ExpressインターフェイスやInfiniBandなどのノード間インターコネクタの性能に大きく依存し、実行時間の大部分をGPU間のall-to-all通信が占めることになる。CUDAではversion 4.0からこのGPU間通信を強化しており、PCI-Expressネットワークを介したGPU間の直接通信や、InfiniBand HCAとの協調動作などがサポートされた。これらの機能を活用することで、通信が全く必要ないシングルGPU実行時と比べても4GPU搭載システムで約2倍、64ノード64GPUのクラスタで最大13倍の性能向上を実現した。通信関連の各種自動最適化は今後の課題である。

シングル GPU 用の FFT ライブラリは NukadaFFT ライブラリという名前で既に一般公開しており、現時点（平成24年3月31日）でのダウンロード数は576である。

5. 主な発表論文等

（研究代表者、研究分担者及び連携研究者には下線）

〔雑誌論文〕（計3件）

- ① 遠藤 敏夫, 額田 彰, 松岡 聡. スーパーコンピュータ TSUBAME 2.0 における Linpack 性能1 ペタフロップス超の達成. 情報処理学会論文誌コンピューティングシステム, 査読有, Vol. 4, No.4 (ACS 35), pp.169-179, 2011年10月.
- ② 額田彰. CUDA による高速フーリエ変換, 応用数理, 査読有, 第20巻, 第2号, pp.37-43, 応用数理学会, 2010年6月.
- ③ 遠藤敏夫, 額田彰, 松岡聡. 異種アクセラレータを持つ TSUBAME スーパーコンピュータの Linpack 評価, 応用数理, 査読有, 第20巻, 第2号, pp.29-36, 応用数理学会 2010年6月.

〔学会発表〕（計18件）

- ① Akira Nukada, Yutaka Maruyama, Satoshi Matsuoka. “High Performance 3-D FFT using multiple CUDA GPUs”, In Proceedings of the Fifth Workshop on General Purpose Processing using Graphics Processing Units (GPGPU-5) in conjunction with ACM ASPLOS XVII, London, UK, ACM Press, March 3rd, 2012.
- ② 遠藤 敏夫, 松岡 聡, 額田 彰, 長坂 真路, 四津 匡康, “グリーンスパコン TSUBAME2.0 における電力危機対応運用, 情報処理学会研究報告, Vol. 2011-ARC-197/HPC-132, 札幌, 2011年11月28日.
- ③ Takashi Shimokawabe, Takayuki Aoki, Tomohiro Takaki, Akinori Yamanaka, Akira Nukada, Toshio Endo, Naoya Maruyama, and Satoshi Matsuoka, “Peta-scale Phase-Field Simulation for Dendritic Solidification on the TSUBAME 2.0 Supercomputer”, In Proc. of 2011 ACM/IEEE International Conference for High Performance, Networking, Storage, and Analysis (SC’11), Seattle, ACM Press, Nov. 15th, 2011. (Technical Paper and Gordon Bell Award finalist.)
- ④ Shuntaro Yamazaki, Akira Nukada, Masaaki Mochimaru, “Hamming Color Code for Dense and Robust One-shot 3D

Scanning”, In Proc. of the 2011 British Machine Vision Conference, Dundee, Scotland, Springer, Aug. 30th, 2011.

- ⑤ Akira Nukada, “Fast Fourier Transform for AMD GPUs”, AMD Fusion Developer Summit 2011, Bellevue, WA. Jun. 15th, 2011.
- ⑥ 遠藤 敏夫, 額田 彰, 松岡 聡. スーパーコンピュータ TSUBAME 2.0 における Linpack 性能1 ペタフロップス超の達成. 先進的計算基盤システムシンポジウム (SACIS2011), 秋葉原, 2011年5月27日.
- ⑦ Akira Nukada, Hiroyuki Takizawa, and Satoshi Matsuoka. “NWCR: A Transparent Checkpoint-Restart Library for NVIDIA CUDA”, In Proc. of 20th Heterogeneity in Computing Workshop (HCW 2011), in conjunction with IPDPS 2011, Anchorage, AK, USA, May 16th, 2011.
- ⑧ Leonardo Bautista Gomez, Akira Nukada, Naoya Maruyama, Franck Cappello and Satoshi Matsuoka. Low-overhead diskless checkpoint for hybrid computing systems, In Proc. of International Conference on High Performance Computing (HiPC 2010), Goa, India, Dec.20th, 2010.
- ⑨ Ali Cevahir, Cevdet Aykanat, Ata Turk, B. Barla Cambazoglu, Akira Nukada and Satoshi Matsuoka. Efficient PageRank on GPU Clusters, 情報処理学会研究報告, Vol. 2010-HPC-128, No.21, 札幌, 2010年12月17日.
- ⑩ 遠藤 敏夫, 額田 彰, 松岡 聡. ヘテロ型スーパーコンピュータ TSUBAME 2.0 の Linpack による性能評価, 情報処理学会研究報告, Vol. 2010-HPC-128, No.11, 札幌, 2010年12月16日.
- ⑪ 長坂仁, 丸山直也, 額田彰, 遠藤敏夫, 松岡聡. GPUにおけるモデルに基づいた電力効率の最適化, 情報処理学会研究報告, Vol. 2010-HPC-128, No. 2, 札幌, 2010年12月16日.
- ⑫ Takashi Shimokawabe, Takayuki Aoki, Chiashi Muroi, Junichi Ishida, Kohei Kawano, Toshio Endo, Akira Nukada, Naoya Maruyama and Satoshi Matsuoka, An 80-Fold Speedup, 15.0 TFlops, Full GPU Acceleration of Non-Hydrostatic Weather Model ASUCA Production Code, In Proc. of the 2010 ACM/IEEE conference on Supercomputing (SC’10), New Orleans, IEEE Press, Nov. 17th, 2010.

- ⑬ Akira Nukada and Satoshi Matsuoka. "NukadaFFT: An Auto-Tuning FFT Library for CUDA GPUs". NVIDIA GPU Technology Conference 2010, Research Summit Poster, San Jose, September 22nd, 2010.
- ⑭ Hitoshi Nagasaka, Naoya Maruyama, Akira Nukada, Toshio Endo and Satoshi Matsuoka. Statistical Power Modeling of GPU Kernels Using Performance Counters, Proceedings of the First International Green Computing Conference (IGCC' 10), Chicago, Aug. 17th, 2010.
- ⑮ Ali Cevahir, Akira Nukada, and Satoshi Matsuoka. High Performance Conjugate Gradient Solver on Multi-GPU Clusters Using Hypergraph Partitioning, Computer Science - Research and Development, Vol. 25, No. 1-2, Springer, (Proceedings of the 2010 International Supercomputing Conference (ISC' 10), Hamburg, Germany, May 31st, 2010).
- ⑯ Akira Nukada and Satoshi Matsuoka, "Fast Fourier Transform using CUDA GPUs", ETHZ - Tokyo Tech Workshop: Computing with GPUs, Cells, and Multicores, Zurich, Switzerland, May 11th, 2010.
- ⑰ Naoya Maruyama, Akira Nukada, and Satoshi Matsuoka. A High-Performance Fault-Tolerant Software Framework for Memory on Commodity GPUs, In Proceedings of 24th IEEE International Parallel and Distributed Processing Symposium (IPDPS 2010), Atlanta, April 20th, 2010.
- ⑱ Toshio Endo, Akira Nukada, Satoshi Matsuoka, and Naoya Maruyama. Linpack Evaluation on a Supercomputer with Heterogeneous Accelerators, In Proceedings of 24th IEEE International Parallel & Distributed Processing Symposium (IPDPS 2010), Atlanta, April 21st, 2010.

[図書] (計1件)

- ① Tamito Kajiyama, Akira Nukada, Reiji Suda, Hidehiko Hasegawa, and Akira Nishida. Toward Automatic Performance Tuning for Numerical Simulations in the SILC Matrix Computation Framework, Chapter 11 of "Software Automatic Tuning : From Concepts to the State-of-the-Art Results", pp. 175-192, Springer, Sep. 2010.

[その他]
ホームページ等

下記 URL にて本研究の成果の一部である NukadaFFT ライブラリソフトウェアを公開.
<http://matsu-www.is.titech.ac.jp/~nukada/nufft/>

6. 研究組織

(1) 研究代表者

額田 彰 (NUKADA AKIRA)
東京工業大学・学術国際情報センター・産学官連携研究員
研究者番号 : 40545688

(2) 研究分担者

()

研究者番号 :

(3) 連携研究者

()

研究者番号 :

