

科学研究費助成事業 研究成果報告書

平成 29 年 6 月 9 日現在

機関番号：13501

研究種目：基盤研究(C) (一般)

研究期間：2014～2016

課題番号：26400197

研究課題名(和文) CPU/GPU混在環境におけるタイルLU分解アルゴリズムの実行時自動チューニング

研究課題名(英文) Runtime autotuning of tile LU factorization for CPU/GPU hybrid environments

研究代表者

鈴木 智博 (SUZUKI, Tomohiro)

山梨大学・総合研究部・准教授

研究者番号：70235977

交付決定額(研究期間全体)：(直接経費) 3,400,000円

研究成果の概要(和文)：本研究は、並列性の高い行列分解アルゴリズムであるタイルアルゴリズムをCPUとGPUからなる計算環境上に効率的に実装し、大規模密行列のLU分解を高速化することを目的とする。そのために、適応的タイルサイズチューニングのためのデータ構造の検討、効率的なタスクスケジューリング手法、性能モデルの構築、新しいピボット選択戦略について検討を行った。この中で、タイルサイズチューニング、性能モデルの構築に関しては満足できる成果が得られた。また、スーパーコンピュータ上に行列分解のタイルアルゴリズムの実装を行った。

本研究に関して、14件の口頭発表(うち査読付き国際会議論文2件)を行った。

研究成果の概要(英文)：The purpose of this research is to implement the tile algorithms for matrix decomposition efficiently on a CPU/GPU computing environment. This implementation makes it possible to speed up the LU decomposition for a large-scale dense matrix. For this purpose, data structures for adaptive tile size tuning, efficient task scheduling method, construction of the performance model and new pivoting strategy were examined. Among them, satisfactory results were obtained regarding the tile size tuning and the performance model. We also implemented the tile algorithm for matrix decomposition on the GPU supercomputer TUBAME 2.5. About this research, 13 oral presentations and two international conference papers with peer review were given.

研究分野：高性能計算

キーワード：タイルアルゴリズム タイルサイズチューニング CPU-GPU実装 LU分解 ピボット選択

1. 研究開始当初の背景

科学技術計算の大規模化に伴って、行列分解の大規模化、高速化の要求が高くなっている。これらの要求に対して、CPU だけを使用するのではなく、GPU によって計算を加速する試みが多く見られる。マルチコア CPU や GPU のハードウェア性能を最大限に引き出すために、アルゴリズムには高い並列性が求められる。タイルアルゴリズム[参考文献(1)]は高い並列性を持った行列分解アルゴリズムであり、近年注目を集めている。タイルアルゴリズムに基づくライブラリ開発プロジェクトも公開されている [参考文献(2), (3)]。

行列分解のうち、連立一次方程式の解法に使用される LU 分解と、各種前処理や最小二乗法の解法などに使用される QR 分解は科学技術計算で多用される基本的かつ重要な計算である。我々はこれまで、マルチコア CPU 環境、スーパーコンピュータ、CPU/GPU 混在環境上に QR 分解のタイルアルゴリズムを実装し、特にタスクスケジューリング方法に着目して効率化を行ってきた。

タイルアルゴリズムは行列を小行列 (タイル) に分割し、タイル毎に処理を行う (図 1)。処理は「タイル分解」部と「タイル更新」部に大別できる。タイル更新部はデータ並列性が高く GPU での実行に向くが、タイル分解部は条件分岐や逐次処理が多く CPU での実行に向いている。CPU/GPU 混在環境では、行列全体でなく処理するタイルのみを GPU メモリに転送することで、GPU の小さなメモリ空間でも大規模行列を扱うことが可能となる。

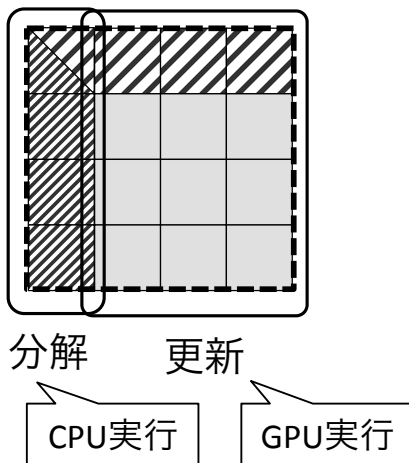


図 1 タイル LU 分解の概要

<参考文献>

- (1) A. Buttari, J. Langou, J. Kurzak and J. J. Dongarra, "A class of parallel tiled linear algebra algorithms for multicore architectures", *Parallel Computing*, 35, pp 38 - 53 (2009).
- (2) PLASMA Project, <http://icl.cs.utk.edu/plasma/>.

(3) FLAME Project

<http://www.cs.utexas.edu/~flame/web/>.

2. 研究の目的

LU 分解、QR 分解などの行列分解は科学技術計算で多用される基本的かつ重要な計算である。本研究の目的は、並列性の高い行列分解アルゴリズムであるタイルアルゴリズムを複数のマルチコア CPU と GPU からなる環境上に効率的に実装し、大規模密行列の LU 分解を高速化することである。

このために、実行時間に大きな影響を及ぼすパラメータであるタイルサイズを、新たに構築するタイルアルゴリズムの性能モデルに基づき実行時に適応的にチューニングする機構を実装する。さらに、計算精度に大きな影響を及ぼすピボット選択法について、タイルアルゴリズム向きの手法を新たに開発し、その有効性を示す。

3. 研究の方法

- (1) タイルサイズの変更やタイルの位置交換が容易に行えるような行列のデータ構造を設計する。
- (2) タイルサイズのチューニングのために、行列サイズ、CPU-GPU 間データ転送時間などさまざまなパラメータが実行時間に与える影響を評価し性能モデルを構築する。得られた性能モデルを使用して、プログラム各部の実行時間測定機構から得られる測定データからタイルサイズを決定し、これを実行中に適応的に変更する機構を実装し、複数の環境において性能評価を行う。
- (3) タイルピボット選択戦略として、並列性の高いトーナメント方式のピボット選択とタイル行入れ替え機構をアルゴリズムに組み込む。テスト行列のサイズ、条件数、タイルサイズを変化させて性能評価を行い、タイルピボット選択法とタイルサイズが誤差に与える影響について定量的な解析を行う。

4. 研究成果

平成 26 年度:

- (1) 行列のデータ構造の検討について、実行時タイルサイズ変更を想定した実装において、他の処理とタイルサイズ変更をオーバーラップさせることで、十分実用的な時間でタイルサイズ変更が行えることを確認した。[学会発表(10, 12)]
- (2) ピボット選択について、タイル内のみでピボット選択を行う方式であるインクリメンタルピボットリングではタイルサイズを大きくしないと十分な計算制度が得られないことが分かった。これは実行時間の意味で最適なタイルサイズ

ではないため、制度、速度の両面で効果的なピボット戦略が必要であることを確認した。[学会発表(10)]

- (3) その他、縦長行列向け行列分解アルゴリズムの高速化の検討[図 2、学会発表(14)]と、これらのスーパーコンピュータへの実装、CPU-GPU 混在環境における効果的なタスク分散方法の検討[学会発表(11, 13)]を行った。

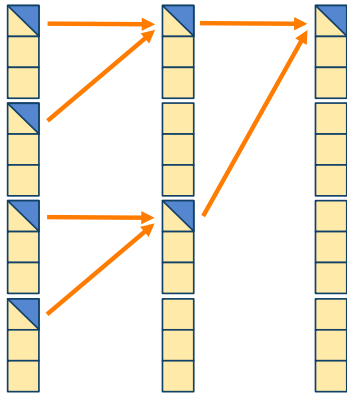


図 2 縦長行列では、縦方向の逐次性が性能低下要因となるため、並列リダクションを導入する

平成 27 年度:

- (4) 性能モデルの構築について、研究当初は性能モデルに基づくタイルサイズチューニングを計画していたが、平成 26 年度中に方針変更を行い、実行環境の並列計算資源に対して十分な量のタスクを供給することをチューニングの方針とし、新たなチューニング指標を導入した。これにより、共有メモリ環境において行列サイズ、並列計算資源に応じて、速度の意味で最適に近いタイルサイズを高速に設定可能となった。[学会発表(5, 8, 9)]
- (5) CPU/GPU 混在環境向けのチューニング機構について、CPU と GPU では速度の意味で最適なタイルサイズが異なるため、GPU 向きに比較的大きなタイルサイズを設定し、CPU での処理ではこれを細分する再帰的タイル分割アルゴリズムを実装した。[図 2、学会発表(6, 7)]
- (6) ピボット戦略に関して、これまで想定していたトーナメント方式の戦略は性能上の大きなボトルネックとなり、分散メモリ環境においては実用的な実装が得られなかった。ピボット戦略のようなリダクション演算を分散メモリ環境上で効率的に行うために、データ構造にリストを導入することで高速化できる可能性があることが分かった。

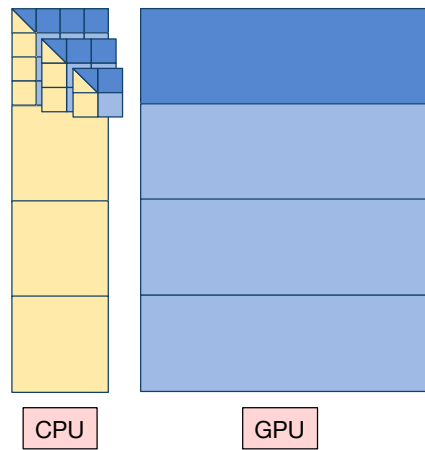


図 2 CPU で扱うタイルは小さく、GPU では大きくなるよう再帰的に細分されたタイル

平成 28 年度:

- (7) ピボット選択について満足な性能がえられないため再検討を行い、タイル列分解ではタイル分割せずに、逐次アルゴリズムよりも高速な再帰アルゴリズムを適用することとした。これにより、タイルアルゴリズム向けのインクリメンタルピボットリングよりも速度、精度の両面で有利な実装が得られた。これはタイル LU 分解と再帰 LU 分解のハイブリッド版という意味合いのアルゴリズムと言える。[図 3]

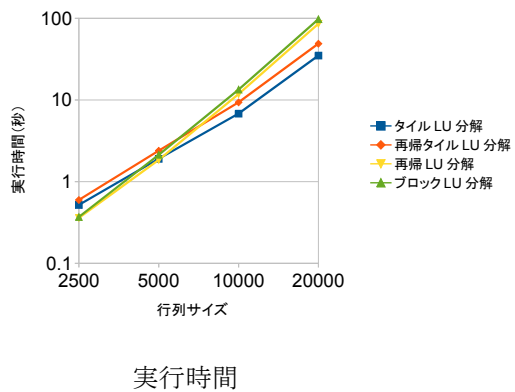
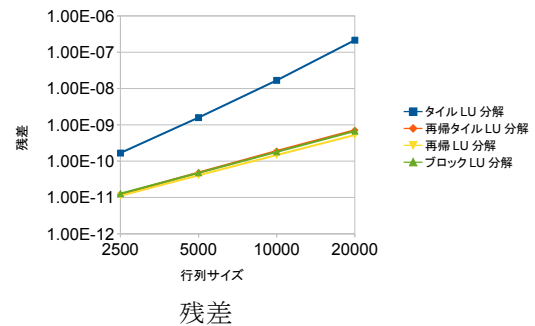


図 3 各種 LU 分解の精度と速度の比較 (2016 年度山梨大学工学部コンピュータ理工学科卒業論文より転載)

- (8) 適応的なタイルサイズ変更を行うチューニング機構は実現できなかった。しかし、平成 27 年度に導入した事前実験を必要とする高速なチューニング機構は、適応的なタイルサイズチューニング機構と比べて低いコストで十分な成果が得られるものと評価した。[学会発表 (1, 2, 3, 4)]

表 1 高速な枝刈り探索よりも更に高速にタイルアルゴリズムのチューニングが実施できていることを示す表。この表は、パラメータ候補の探索時間を表す。表中の Pruned search は従来法の中でも高速な枝刈り探索、Reduced は提案手法 (学会発表 (1) より転載)

	SSRFB	Pruned search		Reduced	
		pdgeqrf	DS	pdgeqrf	DS
Xeon	1.5	18.2	16.9	9.8	9.4
FX10	3.3	31.6	27.6	17.0	15.4
Opteron	2.0	17.5	16.7	9.5	8.9

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)

[学会発表] (計 14 件)

- (1) T. Suzuki, “Faster method for tuning the tile size for tile matrix decomposition”, Proceedings of IEEE 10th International Symposium on Embedded Multicore/Many-core SoCs (MCSoc-16), 査読有, pp. 329 – 336 (2016/09/22, リヨン (フランス)).
- (2) 高柳雅俊, 鈴木智博, “京を用いた CAQR アルゴリズムの性能評価”, 自動チューニング研究会第 15 回オープンアカデミックセッション, (2016/10/24, 山梨大学武田キャンパス (山梨県甲府市)).
- (3) 高柳雅俊, 鈴木智博, “縦長行列におけるタイル CAQR アルゴリズムの性能評価”, 情報処理学会第 148 回ハイパフォーマンスコンピューティング研究会, (2017/03/09, 大月ホテル和風館 (静岡県熱海市)).
- (4) 高柳雅俊, 鈴木智博, “クラスタ型ヘテロジニアス環境におけるタイル QR 分解”, 日本応用数学会 2016 年度年会, (2016/09/12, 北九州国際会議場 (福岡県北九州市)).
- (5) T. Suzuki, “Improvement internode communication for tile QR decomposition for multicore cluster method”, Proceedings of IEEE 29th International Parallel & Distributed Processing Symposium (IPDPS2015), 査読有, pp. 1214 – 1220, (2015/05/29, ハイデラバード (インド)).
- (6) 高柳雅俊, 鈴木智博, “CPU/GPU 混在環
- 境における再帰的タイル QR 分解の動的スケジューリング実装”, 日本応用数学会第 12 回研究部会連合発表会, (2016/03/04, 神戸学院大学ポートアイランドキャンパス (兵庫県神戸市)).
- (7) 高柳雅俊, 鈴木智博, “CPU/GPU 混在環境における再帰的タイル QR 分解”, 日本応用数学会 2015 年度年会, (2015/09/11, 金沢大学角間キャンパス (石川県金沢市)).
- (8) 高坂知弘, 鈴木智博, “タイル CAQR の MPI/OpenMP ハイブリッド並列化”, 日本応用数学会 2015 年度年会, (2015/09/11, 金沢大学角間キャンパス (石川県金沢市)).
- (9) 鈴木智博, “共有メモリ環境上でのタイル QR 分解のタイルサイズチューニング”, 情報処理学会第 151 回ハイパフォーマンスコンピューティング研究会, 情報処理学会研究報告 Vol. 2015-HPC-151, No. 21, pp. 1-7, (2015/10/01, 沖縄産業支援センター (沖縄県那覇市)).
- (10) T. Suzuki, “Implementation of Tile Algorithms for Matrix Decomposition on CPU/GPU Systems”, Conference on Advanced Topics and Auto Tuning in High-Performance Scientific Computing (ATAT2015), (2015/02/27, 台北 (台湾)).
- (11) 鈴木智博, “CPU/GPU 混在環境におけるタイル LU 分解アルゴリズムの実行時自動チューニング”, 第 6 回自動チューニング技術の現状と応用に関するシンポジウム (ATTA2014), (2014/12/25, 東京大学本郷キャンパス (東京都文京区)).
- (12) 鈴木智博, “クラスタシステム向けタイル QR 分解のタイルサイズチューニング”, 情報処理学会第 146 回ハイパフォーマンスコンピューティング研究会, 情報処理学会研究報告 Vol. 2014-HPC-146, No. 15, pp. 1-7, (2014/10/03, 沖縄産業支援センター (沖縄県那覇市)).
- (13) 高柳雅俊, 鈴木智博, “マルチ GPU 環境におけるタイル CAQR アルゴリズムの実装”, 日本応用数学会 2014 年度年会, (2014/09/05, 政策研究大学院大学 (東京都港区)).
- (14) 小嶋弘樹, 鈴木智博, “スーパーコンピュータを使用した列方向並列化による行列分解計算の高速化”, 日本応用数学会 2014 年度年会, (2014/09/05, 政策研究大学院大学 (東京都港区)).

[その他]

ホームページ等

- <https://site.google.com/site/tomohirosuzuki20140211/>
- http://erdb.yamanashi.ac.jp/rdb/A_Index.Main

6. 研究組織

(1) 研究代表者

鈴木智博 (SUZUKI, Tomohiro)
山梨大学・総合研究部・准教授
研究者番号：70235977