

受付 ID	17a50
分野	数値解析

## GPU クラスタにおける通信回避アルゴリズムによる

### 行列計算の実装と性能最適化

#### Implementation and performance optimization of communication-avoiding algorithms for matrix computation on GPU cluster

松本 和也

会津大学コンピュータ理工学部

#### 1. 研究目的

GPU クラスタにおいて利用 GPU 数を増やしたときには、その通信時間の増加が性能ボトルネックとなり高い演算性能を発揮するための障壁となる場合が少なくない。近年では増大する演算と通信の性能差に対応するために通信回避 (Communication-avoiding) アルゴリズムや通信削減アルゴリズムと呼ばれるアルゴリズムが注目を集めている。本プロジェクトでは、GPU クラスタの高い演算性能を有効に利用する Krylov 部分空間反復法の通信回避アルゴリズムの実装の実現を目指し研究開発を行った。

本プロジェクトは、HA-PACS/TCA において GPU を演算に利用した通信回避 GMRES 法 (CA-GMRES) の実装に関する研究開発を行うことを計画の中心とした。CA-GMRES の実装では Matrix Powers Kernel (複数回の行列ベクトル乗算)、Tall-skinny QR 分解 (行数が列数と比べて極めて大きい行列に対する QR 分解) を主演算とする各種の行列計算・ベクトル計算が必要である。通信回避アルゴリズムは必要な演算量が多少増大する代わりに必要な通信回数を減らすものであり、単純にアルゴリズムを実装したとしても性能が向上するとは限らず実装の最適化を行うことが求められる。そこで本プロジェクトでは、CA-GMRES の実装の性能を大きく決定する処理を明確にし、高性能な CA-GMRES 実装の実現を目指した。

#### 2. 研究成果の内容

HA-PACS/TCA において通信回避アルゴリズム CA-GMRES の実装を行い、その評価結果から CA-GMRES はよく知られた Krylov 部分空間法ソルバ実装 (GCR 法と GMRES の実装) の性能よりも高速に解を求められることを確認した。特に利用計算ノード数が多い場合において CA-GMRES の実装が高い性能を示すことを確認した。

本プロジェクトで対象とした核融合シミュレーションコード GT5D に現れる行列・ベクトルデータ向けの Matrix Powers Kernel (MPK) の実装については、行列ベクトル乗算を複

数回実行する実装に対する性能的な有利性は確認できなかった。MPKはノード間の通信回数は削減できるが演算量とノード間通信量は増加してしまうため、対象としたデータにおいては性能が向上しなかったと考えられる。

Tall-Skinny QR分解を行う通信回避型アルゴリズムである Cholesky QR法の有利性を調べるために、QR分解を行う他のアルゴリズム（古典 Gram-Schmidt法、修正 Gram-Schmidt法）を実装し性能比較を行った。比較の結果、Cholesky QR法はノード間通信回数を減らすとともに高性能な行列演算ルーチンを利用できるという特徴により、それらのQR分解の実装よりも高速であることが確認できた。

HA-PACS/TCAにおいてCA-GMRES実装の性能を測定し評価を行った結果については、論文にまとめて投稿することを予定している。

### 3. 学際共同利用として実施した意義

HA-PACS/TCAは国内で運用されている高性能なGPUクラスタ計算機であり、本プロジェクトは学際共同利用によりHA-PACS/TCAを利用することなしでは進めることが困難であった。特に計算ノードあたり4基のGPUが搭載されているHA-PACS/TCAにおいて実装したプログラムの性能を測定し評価できたことは意義があった。

### 4. 今後の展望

今後は他の研究機関で運用されているGPUを搭載した並列計算機システムにおいて研究を継続する予定である。現実装はGT5Dコードにおいて現れる行列・ベクトルデータ向けのものであり、今後は他のアプリケーションコードにおいても動作するように実装の修正を行うことを予定している。また、CA-GMRES実装の主演算であるMPKやTSQRの実装に関しては性能改善の余地があり、それらの実装を改良することも今後の課題である。

### 5. 成果発表

#### (1) 学術論文

#### (2) 学会発表

1. Yasuhiro Idomura, Takuya Ina, Akie Mayumi, Susumu Yamada, Kazuya Matsumoto, Yuuichi Asahi, and Toshiyuki Imamura. "Application of a communication-avoiding generalized minimal residual method to a gyrokinetic five dimensional Eulerian code on many core platforms," In Proceedings of the 8th Workshop on Latest Advances in Scalable Algorithms for Large-Scale Systems (ScalA '17), No. 7, Nov. 2017, Denver, USA, 2017.

#### (3) その他

使用計算機	使用計算機 に○	配分リソース※	
		当初配分	追加配分
HA-PACS/TCA	○	3360	0
COMA			
Oakforest-PACS			
※配分リソースについてはノード時間積をご記入ください。			