

科学研究費助成事業 研究成果報告書

平成 29 年 6 月 21 日現在

機関番号：32613

研究種目：基盤研究(C) (一般)

研究期間：2013～2016

課題番号：25330144

研究課題名(和文)ベクトル演算と通信量削減によるマルチコア向け4倍精度反復法アルゴリズムの開発

研究課題名(英文)Development of double-double precision iterative method for multi-core processor based on vector operation and communication avoidance

研究代表者

田中 輝雄 (Tanaka, Teruo)

工学院大学・情報学部(情報工学部)・教授

研究者番号：90622837

交付決定額(研究期間全体)：(直接経費) 3,300,000円

研究成果の概要(和文)：1. AVX/AVX2を用いた倍々精度演算ライブラリDD-AVXを開発。また、Cプログラム上のdouble型データをGMPによる任意多倍長型に変換するツールXev-GMPを開発した。Xev-GMPはMPIにも対応できる。
2. 通信量削減方式として、Chebyshev基底共役勾配法に対して、集団通信を削減するCBCGR法を提案し評価した。また、Matrix Power Kernel法において重複を排除する計算パターンを考案した。
3. 自動チューニングにおいて、2つの性能パラメタの同時推定方式を自動チューニング基盤pp-OpenATに実装。さらに3つ以上の性能パラメタに対して、反復線形探索方式を提案した。

研究成果の概要(英文)：1. We developed DD-AVX, a library of Double-Double (DD) precision matrix and/or vector operations accelerated by AVX and AVX2 SIMD instructions. We also developed Xev-GMP, a directive-based automatic code generation for a C code with multiple-precision floating-point-operation data from a C code with double precision data. The GMP code uses the GNU Multiple Precision Arithmetic Library. Xev-GMP can also support MPI libraries.
2. We proposed and evaluated CBCGR method which reduces MPI collective communication for the Chebyshev base conjugate gradient method in Massively parallel processing. We also devised a calculation pattern to remove the overlap in the Matrix Power Kernel method.
3. We implemented the simultaneous estimation method of two performance parameters on automatic tuning base pp-OpenAT. The enhanced version of pp-OpenAT is made public. Furthermore, we suggested a repetitive linear search method for the practical use of performance parameters more than three.

研究分野：高性能計算

キーワード：倍々精度計算 任意多倍長計算 AVX GMP 通信量削減 MPI 自動チューニング 性能パラメタ

1. 研究開始当初の背景

(1) 4倍精度演算の必要性と課題

近年、計算機パワーの拡大により、いままですべて解くことが難しかった「解」の精度を出しにくい、いわゆる悪条件問題に対する疎行列の反復解法による大規模シミュレーションのニーズが高まってきている。このような悪条件問題に対しては、その精度を確保するために、4倍精度あるいはそれ以上の高精度演算が必要になる。この課題の解決策として、反復解法ライブラリ Lis では、インテル社プロセッサ・アーキテクチャのベクトル命令 SSE2 を活用した 4 倍精度演算が実装されている。なお、ここでの 4 倍精度とは、2 つの倍精度型データを用いて倍精度演算の組み合わせで実現する、いわゆる Double-Double 型である。したがって、4 倍精度計算を実現するためには倍精度演算に比べて約 2 倍のデータ領域と約 20 倍の演算量（演算の構成比率により異なる）が必要となる。また、一般に測定データなど入力データは 4 倍精度で提供されることはないため、係数（疎）行列自体は倍精度で保持し、作業用ベクトルのみを 4 倍精度として扱うことになる。この場合、疎行列-ベクトル積では 4 倍精度と倍精度の混合演算が必要となる。

(2) ベクトル演算機構の利用と研究状況

一方、プロセッサ・アーキテクチャの動向に目を向けると、プロセッサコア内の演算能力向上として、Intel 社により、ベクトル機能として SSE2 から SSE4/AVX へ機能拡張が行われている。さらに来年度投入される Xeon Phi ではさらなる機能拡張が予定されており、プロセッサ内のコア数も現在の 10~12 コアから、Xeon Phi では 40~50 コアが実用化される。私たちの研究グループでは、これまで反復解法ライブラリ Lis で実装されている SSE2 利用コードに対し、AVX 利用コードを作成・実装し、その特性を分析してきた。その結果、アーキテクチャの拡張の効果により、AVX では並列度の拡張(2 並列 4 並列)以上の性能向上を可能とすることや、マルチコア環境では、通常、4 コア程度でメモリ(主記憶)からプロセッサ・コアへのデータ性能が限界となり、複数コアを用いた並列効果が得られなくなることを見出した。したがって、今後のマルチコア、あるいは、メニーコアと呼ばれる多数のコアを用いる環境において高い性能を追求するためには、さらなる工夫が必要となる。

(3) データ通信量の削減によるマルチコアの有効利用

マルチコア環境においてはメモリからコアへのデータ通信時間を極力抑えることが重要となる。そのような中で、通信量を極力減らした新しいアルゴリズムが提案されている。繰り返し行う同一の行列に対する行列ベクトル積をまとめることにより、行列デ

ータをキャッシュから外すことなく再利用することができる。本研究で扱う 4 倍精度計算は、通常の倍精度計算より演算比率が高いため、この通信量を削減するアルゴリズムを用いることにより、相乗効果で多くの数のコアまで有効に利用できる可能性が高い。

(4) 実行時自動チューニング機構の適用

さらに、コア数の増大に伴い、50 を超えるコア数を有するハイエンド構成から、コア数の少ない廉価版までプロセッサレベルにおいても多くの構成があり、さまざまな容量・性能のメモリ構成との組み合わせを考える必要がある。そのような状況下では、システムの性能を左右する性能パラメータを見出し、その性能パラメータを自動的にシステムごとに最適化する自動チューニングが必要となる。特に、主に対象とする疎行列の非零要素の構造はプログラムの実行時に決定されることが多く、いかにコストをかけずに実行時に自動チューニングを行うかが重要となる。現在、我々はこの実行時自動チューニングを低コストで行う研究を推進しており、フレームワークはすでに開発済みである。

2. 研究の目的

収束性の改善と高速性を両立させるため、高精度演算を用いた(疎)行列を対象とした大規模反復計算ライブラリの実現を目指している。その1ステップとして、マルチコア環境での4倍精度演算(倍々精度演算)の実用化を研究している。本研究では、第1に、進化しつつある Intel 社ベクトル演算機構の特性を明らかにし、そのプロセッサ・アーキテクチャを駆使して、CPU 性能の極限を追求した倍々精度演算処理の実現をはかる。第2に、コア数が増加した際のメモリ コア間のデータ転送の改善により、より多くのコア数で有効となるアルゴリズムを開発する。第3に、自動チューニング技術を用いて自動的に精度切り替えを実現し、ユーザの負担を最小限にするライブラリの実現を目指す。

3. 研究の方法

(1) AVX およびその後継拡張機能を用いた4倍精度演算のプロトタイプ作成

すでに、私たちの研究グループでは AVX を用いた 4 倍精度による基本的な四則演算機能を作成している。その結果、SSE2 から AVX へのアーキテクチャ拡張の効果により、AVX による並列度の拡張(2 並列 4 並列)以上の性能向上を可能とすることや、マルチコア環境において、通常では 4 コア程度でメモリからプロセッサコアへのデータ性能が限界となり、複数コアを用いた並列効果が得られなくなることを見出した。それらをもとに、

Xeon Phi 向け AVX 拡張アーキテクチャにも対応し、最適化をはかる。マルチコア環境での 4 倍精度演算の BLAS レベルでの特性評価、高速化を行い、コア数向上に対する性能

効率劣化を最小限に抑える方式を開発する。さらに、主要ターゲットである線形ソルバで利用される疎行列ベクトル積計算に焦点を当て（この処理は倍精度行列と4倍精度ベクトルとの演算となる）特性評価、高速化を行う。

上記については、演算器の利用効率向とともに、マルチコアにおけるメモリ上のデータの分散配置の工夫等により、メモリとコア間でのデータ通信の効率化を進め、コア数向上に対する限界値を引き上げ、さらなるプロセッサ・コアの有効活用をはかる。逆に、対象とする問題規模および許される計算時間を固定すれば、計算に必要なプロセッサコア数を削減することが可能となり、省電力に寄与することができる。

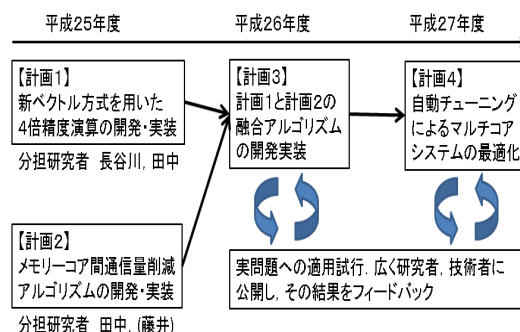
（2）通信量削減を目的とする新しいアルゴリズムの開発・実装

Demmel 氏らが提唱するメモリ コア間のデータ通信量を極力減らした新しいアルゴリズムでは、(疎)行列ベクトル積 Ax , A^2x , A^3x , A^4x ...を同時に計算することで、キャッシュ利用効率を上げる。そこでまず、(疎)行列ベクトル積（倍精度データ×倍精度データ）を高速に行うプログラムを実装する。行列Aの内容をこれら複数の行列ベクトル積で使いまわすことにより、効率的に計算をする。この演算を導入するためには、上記のような形の疎行列ベクトル積が複数出てくるように、数値解法を変形させる必要がある。すでに変形方法が研究されている解法を実装するとともに、その他の解法への拡張について研究を進める。次に、倍精度演算を対象とする検討を始める。計画(1)にて開発される基本的4倍精度演算機構をもとに、(疎)行列ベクトル積計算を4倍精度演算×倍精度計算へと切り替える。倍精度演算に比べて、4倍精度演算は、必要データ通信量は2倍、演算量は20倍であり、倍精度演算より多くのプロセッサ・コアを有効に利用できると考える。

（3）実行時自動チューニング機能の組み込み

反復処理中の主要計算である(疎)行列ベクトル積計算などに対し、実行時に行う自動チューニング機能の実装を行う。現在、私たちの研究グループは実行時自動チューニング機能付きライブラリ開発支援機構 ppOpenAT に、実行時自動チューニングにおけるチューニングの効率化機能の実装を進めている。さらに、実行するアプリケーションに対しては、たとえば、BICG 法等の代表的な解法において収束が悪くなると4倍精度演算に(自動的に)切り替える方法を含めて、4倍精度演算自体を動的に必要なに応じて選択できるような実行時自動チューニングも研究対象とする。

なお、これらの研究で開発されたプログラムは、試験的に公開していく。また、その知見・ノウハウについては適宜、論文としてまとめ、広く公開する。



4. 研究成果

（1）AVX およびその後継拡張機能を用いた4倍精度演算のプロトタイプ作成

倍々精度演算を Intel プロセッサのアクセラレータ機構 AVX/AVX2 を用いて高速化した DD-AVX を開発し、リリースした。演算処理の高速化により、倍々精度演算は倍精度演算より20倍演算量があるにも関わらず、データサイズの2倍まで、倍々精度演算の実行時間を倍精度演算の実行時間まで近づけることができた。また、さらなる高精度演算に対応するために、GMP (GNU Multi-Precision Library) を対象とし、double 要素を基本とするC言語をGMPによる任意多倍長型に自動変換するディレクティブ型の変換ツール Xev-GMP を開発した。また、開発した Xev-GMP は MPI の基本関数にも対応しており、評価実験により、データサイズに比例した通信時間で GMP の任意多倍長精度型の MPI 通信を可能とした。一般に GMP のプログラムは膨大な演算時間がかかるので、MPI の利用による大規模並列処理による GMP の適応拡大の可能性を広げた。

（2）通信量削減を目的とする新しいアルゴリズムの開発・実装

通信量削減方式については、高並列環境における共役勾配法 (CG 法) の集団通信の削減として、Communication-avoiding CG 法の一種としたチェビシェフ基底共役勾配法 (CBCG 法) の実装・評価を行なった。また、CBCG 法において、内積計算の集団通信の回数を減らした CBCGR 法を提案・評価した。また、反復処理において繰り返し行われる疎行列ベクトル積 (SpMV) をまとめて、疎行列のべき乗を最初に計算しておく Matrix Powers Kernel (MPK) において、重複演算を排除する計算パターンを考案し、その効果を示した。さらに、これら2つの手法を組み合わせた CBCGR-MPK 法を実装し、高並列計算環境での有効性を示した。

(3) 実行時自動チューニング機能の組み込み

自動チューニングに関しては、祖行列の非ゼロ要素パターンなど、実行時にしかわからない情報を取り込むために、実行時自動チューニングが重要となる。この実行時自動チューニングの実用化に向けた基本技術として、複数の性能パラメタの同時推定方式を検討した。2つの性能パラメタに対する同時推定に対して、2つの性能パラメタを用いて評価関数 d-Spline を拡張したモデル化した逐次追加型性能パラメタ推定法 (IPPE/d-Spline2) を自動チューニング基盤 pp-OpenAT に実装し公開した。さらに、性能パラメタがさらに増やした場合の同時推定を実用化のために、高速な1次元 d-Spline を反復して線形探索する方式を提案し、4性能パラメタの同時推定での効果を確認し、さらなる複数性能パラメタ同時推定の実現の目途を立てた。

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4件)

- [1] 熊谷 洋佑, 藤井 昭宏, 田中 輝雄, 深谷 猛, 須田 礼仁, 共役勾配法への種々の通信削減手法の適用と評価, 情報処理学会論文誌コンピューティングシステム(ACS), Vol.9, No.3, pp.1-13, 2016.8.
- [2] Teruo Tanaka, Ryo Otsuka, Akihiro Fujii, Takahiro Katagiri, Toshiyuki Imamura, Implementation of d-Spline-based incremental performance parameter estimation method with ppOpen-AT, Scientific Programming 22, pp.299-307, 2014.8.
- [3] 菱沼 利彰, 藤井 昭宏, 田中 輝雄, 長谷川 秀彦, AVX2 を用いた倍精度 BCRS 形式疎行列と倍々精度ベクトル積の高速化, 情報処理学会論文誌コンピューティングシステム(ACS), Vol.7, No.4, pp.25-33, 2014.7.
- [4] 坂本 真貴人, 藤井 昭宏, 田中 輝雄, Strassen のアルゴリズムを用いた行列積自動チューニングライブラリ, 電子情報通信学会論文誌 D, Vol. j97-D, No.3, pp.405-413, 2014.3.

[学会発表](計46件)

- [1] Masayoshi Mochizuki, Akihiro Fujii, Teruo Tanaka, Fast Multidimensional Performance Parameter Estimation with Multiple One-dimensional d-Spline Parameter Search, The Twelfth International Workshop on Automatic Performance Tuning(iWAPT), in 31st IEEE International Parallel & Distributed Processing Symposium

(IPDPS), 2017.5.(査読あり)

- [2] Tanaka Teruo, Masayoshi Mochizuki, Guqing Fan, Akihiro Fujii, Two topics about fitting function d-Spline for realization of practical AT, 2017 Conference on Advanced Topics and Auto Tuning in High-Performance Scientific Computing (ATAT in HPSC 2017), 2017.3. (査読なし)
- [3] Toshiaki Hishinuma, Takuma Sakakibara, Akihiro Fujii, Teruo Tanaka, Shoichi Hirasawa, Xev-GMP: Automatic code generation for GMP multiple-precision code from C code, 19th IEEE International Conference on Computational Science and Engineering (CSE 2016), 2016.8. (査読あり)
- [4] Toshiaki Hishinuma, Hidehiko Hasegawa, Teruo Tanaka, SIMD Parallel Sparse Matrix-Vector and Transposed-Matrix-Vector Multiplication in DD Precision, 12th International Meeting on High Performance Computing for Computational Science(VECPAR2016), 2016.6. (査読あり)
- [5] Teruo Tanaka, Enhancement of Functionality of ppOpen-AT with d-Spline based Incremental Performance Parameter Estimation, International Workshop on Software for Peta-scale Numerical Simulation (SPNS2015). (招待講演), 2015.12.
- [6] Riku Murata, Jun Irie, Akihiro Fujii, Teruo Tanaka, Takahiro Katagiri, Enhancement of Incremental Performance Parameter Estimation on ppOpen-AT", Special Session: Auto-Tuning for Multicore and GPU (ATMG-15) In Conjunction with the IEEE 9th International Symposium on Embedded Multicore/Many-core Systems-on-Chip (MCSoc-15), 2015.9. (査読あり)
- [7] Yosuke Kumagai, Akihiro Fujii, Teruo Tanaka, Yusuke Hirota, Takeshi Fukaya, Toshiyuki Imamura, Reiji Suda, Performance Analysis of the Chebyshev Basis Conjugate Gradient Method on the K Computer", 11th International Conference on Parallel Processing and Applied Mathematics (PPAM2015), 2015.9. (査読あり)
- [8] Toshiaki Hishinuma, Akihiro Fujii, Teruo Tanaka, H. Hasegawa. Fast computation of double precision sparse matrix in BCRS and DD vector product using AVX2, 11th International Meeting High Performance Computing for Computational Science (VECPAR2014),

2014.6.(査読あり、ポスター)

- [9] Toshiaki Hishinuma, Akihiro Fujii,
Teruo Tanaka, Hidehiko Hasegawa, "AVX
Acceleration of DD Arithmetic between
a Sparse Matrix and Vector", 10th
International Conference on Parallel
Processing and Applied
Mathematics(PPAM2013), Workshop on
Numerical Algorithms on Hybrid
Architectures, 2013.9. (査読あり)
- [10] Teruo Tanaka, Makito Sakamoto,
Akihiro Fujii, "BLAS3-level Matrix
Multiply Automatic Tuning Implemented
within Strassen", Conference on
Advanced Topics and Auto Tuning in
High Performance Scientific Computing
(@^2HPSC), 2013.3. (査読なし)

〔図書〕(計0件)

〔産業財産権〕

出願状況(計0件)

取得状況(計0件)

〔その他〕

ホームページ等

<http://hpcl.info.kogakuin.ac.jp/>

6. 研究組織

(1) 研究代表者

田中 輝雄 (TANAKA Teruo)

工学院大学・情報学部(情報工学部)・教授

研究者番号: 90622837

(2) 研究分担者

長谷川 秀彦 (HASEGAWA Hisehiko)

筑波大学・図書館情報メディア系・教授

研究者番号: 20164824

(3) 連携研究者

藤井 昭宏 (FUJII Akihiro)

工学院大学・情報学部(情報工学部)・准教授

研究者番号: 10383986