

階層メモリを有するスーパーコンピュータ のためのLU分解並列アルゴリズム

妹尾義樹、白戸幸正¹、西直樹

日本電気㈱C&Cシステム研究所、¹日本電気技術情報システム開発㈱

大規模LU分解を高速に実行する並列アルゴリズムを開発し、実験によりその有効性を確認した。LU分解の対象は非対称密行列である。実験は通産省大型プロジェクトで開発された、ローカルメモリ、共有メモリ、半導体拡張記憶、磁気ディスクの4階層メモリを有する並列型スーパーコンピュータ上で行い、メモリ階層間のデータ転送を最適化するとともに並列化を行った。32768×32768の密行列(データ量4GB)の分解を10時間40分で実行することができた。また16384×16384行列の分解では4台のプロセッサを用いて1台の2.8倍の高速化が達成できた。さらに磁気ディスクを除いた3階層の実行では1台の3.2倍の並列化加速率が得られ、アルゴリズムの有効性が実証できた。

An LU Decomposition Algorithm for Parallel Supercomputers

Yoshiki SEO, Yukimasa SHIROTO¹, Naoki NISHI

C&C Systems Research Laboratories, NEC Corporation

¹NEC Scientific Information System Development

4-1-1, Miyazaki, Miyamae-ku, Kawasaki, Kanagawa 213 JAPAN

A newly designed parallel algorithm for solving LU decomposition of huge dense matrices was developed for parallel vector supercomputers with a hierarchy of memory layers (i.e., local memories, shared memory, semiconductor extended storage and magnetic disk). Using four memory layers, an LU decomposition for a 32768x32768 dense matrix was calculated in 10 hours and 40 minutes on the HPP-LHS supercomputer system developed under the MITI (the Ministry of International Trade and Industry) Supercomputer Project. Required memory capacity for the gigantic matrix is 8GB, and the whole matrix data area was allocated to magnetic disk for this calculation. The execution speed with 4 processors was 2.8 times faster than that with 1 processor, even using a magnetic disk, and the algorithm was proved to be effective.

1. はじめに

ベクトル型のスーパーコンピュータはその超高速の演算能力により科学技術計算の分野では実用マシンとしての評価を定着させつつあるが、複雑な自然現象を解析するためにはさらに高速の演算能力と大量のメモリ容量が要求されている。

演算能力の向上に関してはプロセッサは並列化の方向に向かっており、NECのSX-3やCRAY-2/3、XMP、YMPはいずれもマルチプロセッサ構成となっている^{[1][2]}。また、メモリ容量に関しては半導体高速補助記憶を用いたメモリの階層化が有力な手段であり、ほぼすべての商用機が半導体の拡張記憶をサポートしている。

これらの階層メモリを有する並列型スーパーコンピュータ上で主記憶に入りきらない大量データを扱う大規模数値計算プログラムを高速に実行するためには、メモリ階層間のデータ転送の最適化と並列化が重要となる^[3]。

さらに、現在の並列型の商用機はローカルメモリを持たずにすべてのプロセッサが同一の共有メモリだけをアクセスする密結合方式が主流であるが、この方式では共有記憶へのアクセスの集中により結合台数に限界がある。さらに多数のプロセッサを結合しようとするとしても各プロセッサがローカルメモリを持たざるを得なくなる。この場合にはローカルメモリ共有メモリというメモリ階層が存在し、拡張記憶を用いて計算を行う際には3階層のメモリ階層について各メモリ階層間のデータ転送を最適化する必要がある。

本論文では、基本計算として最も重要なLU分解について、メモリ階層間データ転送の最適化と並列化を行う直接法のアルゴリズムを示す。過去においては2階層のメモリ階層を持つスーパーコンピュータ(シングルプロセッサ+拡張記憶)上でデータ転送を最適化するLU分解アルゴリズムが考案され、9600元の連立一次方程式が解かれている^[4]。また、メモリ階層のない密結合の並列計算機上でLU分解を解く並列アルゴリズムもいろいろと開発されている^{[5][6]}。しかしながら、ローカルメモリを有する並列型スーパーコンピュータ上でしかも3階層以上のメモリ階層を用いて大規模問題を高速に解くアルゴリズムは今回が初めてである。

本アルゴリズムは通産省スーパーコンピュータプロジェクトで開発されたHPP-LHSスーパーコンピュータシステムを用いて過去に例のない32768×32768という大規模なLU分解を実行し、その有効性が確認できた。この実機評価結果についても本稿で報告する。

2. HPP-LHSスーパーコンピュータ

図1にHPP-LHSスーパーコンピュータシステムの概要を示す。本システムは4台のマルチプロセッサであり、最大演算性能は10GFLOPSである(1PEは4GFLOPSで、残りの3PEは2GFLOPS)。各プロセッサは64MBのローカルメモリ(LSU)を有し、1GBの共有メモリ(CSU)に接続されている。LSU、CSUともにワードアドレスでアクセス可能であるが、CSUに対して高速にアクセスするためにはプログラムが明示的にブロック転送を指定する必要がある。

LHSは大容量高速記憶装置で、論理的には2の68乗バイトという大規模な論理空間をサポートする。アクセスの単位は1MBのブロックアドレスである。物理的には4GBの半導体高速記憶部と15GBの磁気ディスク(バッキングストア)の2階層から構成されるが、LHS制御プログラムの提供するワンレベルファイル機能^[7]により、プログラムはバッキングストアの容量の単一ファイル空間と考えてLHSを利用できる。つまり、アクセス(READ)すべきデータが高速記憶部に存在する場合にはハードウェアで実装されたアドレス変換機構により高速に(1GB/秒)アクセスが可能であり、高速記憶部に存在しない場合にはLHS制御プログラムにより自動的にバッキングストアから高速記憶部にデータが転送された後にアクセスされる。この場合にはLHSへのアクセスの実効性能は1.1MB/秒である。

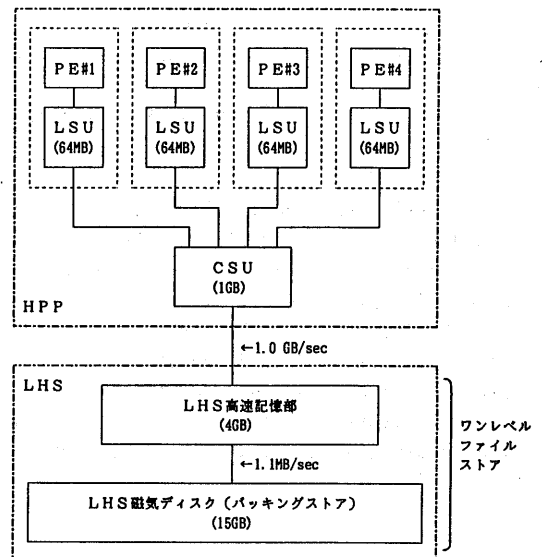


図1 LHS-HPPシステム構成

3. LU分解のデータ転送／並列化アルゴリズム

HPP-LHSシステムにおいてLU分解を行うことを考える。普通は演算器により近い階層のメモリ（以降メモリの階層が上位であるという）ほど高速のデータ転送が可能であるが容量は小さくなる。したがって下位のメモリ階層ほどデータ転送の影響が大きく、こちらから優先してデータ転送を最適化すべきである。各記憶階層間のデータ転送を最適化するために、再帰型データ分割方式を考案した。これは外積形のガウス消去法アルゴリズムをもとにしてデータ依存関係（定義参照の順番）が変化しないように実行の順序を入れ換えて問題を分割するものである。この分割により記憶階層間データ転送の最適化と同時に効率の良い並列化をも行うことができる。

3.1 前提条件

本論文で説明するアルゴリズムでは、LU分解の諸元とデータの格納およびLU分解を行う基本算法について以下のものを前提としている。

まずLU分解は $n \times n$ の2次元の係数行列に対して行うものとし、係数行列は非対称でかつ密行列であるとする。LHSにはすべての係数行列を収容できる記憶域が確保されており、LU分解終了後にはこのエリアに計算結果が格納されるものとする。この際には図2に示すようにFORTRANのデータ格納法と同じく2次元の係数行列がカラムワイズにページ分割されており（各ページが p 列で全体で q ページ）、ひとつ上位の階層の記憶システムとの間でブロック番号の指定によるブロック転送が可能であるとする。また任意の隣接する記憶階層間で直接にデータを転送することができ、他の記憶階層を経由する必要はないものとする。

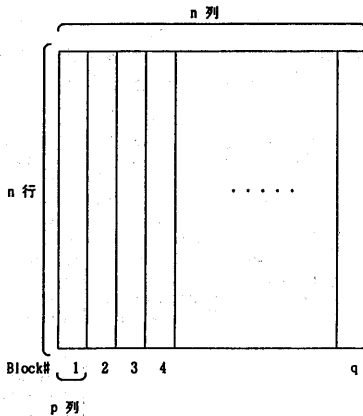


図2 最下位メモリ上の係数行列データ構造

LU分解の実行は基本的には直接法であるガウス消去法に基づいて行われるものとする。ガウス消去法は外積形でも内積形でもどちらでも用いることができるが、本稿では外積形をベースとして説明を行う。外積形のガウス消去法は図3に示すアルゴリズムであり、係数行列は a_{ij} で $n \times n$ のサイズであるとする。演算の進行する様子を図4に示す。k列に対する処理を

- ① a_{kk} の逆数を計算し、
 - ② a_{ki} ($i=k+1, \dots, n$)をこの逆数で割り、
 - ③ ②で計算した a_{ki} を用いて $k+1$ カラムから n カラムまでのデータを更新する。
- として、 k が1から n まで繰り返せば良い。

```

do k=1,n          --- (1)
  akk=1/akk
  do i=k+1,n      --- (2)
    aik=aik*akk
  end do
  do j=k+1,n      --- (3)
    do i=k+1,n
      aij=aij-aik*akj
    end do
  end do
end do
    
```

図3 ガウス消去法（外積形）

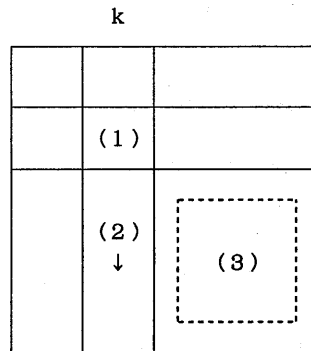


図4 データ参照パターン

(1), (2), (3)は図3のループ処理に対応

3.2 ブロック分割方式

まず、2階層の場合についてデータ転送の最適化について説明する。

上のガウス消去法の順序で演算を進めると、各ピボット列に対する操作の度にピボット列より右側のデータをすべてアクセスする必要があり、非常に多量のデータ転送が必要となる。

そこで処理のローカリティーを高めるために処理を分割する。すなわち、係数行列を複数のブロックにカラムワイズに分割し、全体の消去操作をブロック内消去とブロック間消去に分ける。各ブロックが p 列で q ブロックに分割された場合のアルゴリズムを図5に示す。(a)がブロックを単位として進行する全体処理を示し、(b)、(c)が対応するブロック間消去、ブロック内消去の処理内容である。上位階層のメモリに2ブロック分の作業エリアが確保できれば、図6のように処理を進めることができる。ブロック k の更新についてブロック k を作業エリア2にフェッチし、作業エリア1にブロック 1 からブロック $(k-1)$ を順次フェッチしてブロック間消去を行う。そしてブロック k のブロック内消去を行った後にブロック k を下位メモリに書き戻す。このブロックを単位とする消去の進行は内積形ガウス消去法と同等であり、更新ブロックに対して参照ブロックが走査されるために外積形の処理に比べて走査されるブロック(作業エリア1)の書き戻しが不要であるという利点がある。

内積形ガウス消去法に対して、ブロック分割して演算のローカリティーを高める基本概念は文献[4]で述べられており、この方法を適用して9600元の連立一次方程式が解かれている。

```

do k=1,q
  do j=1,k-1
    (b)
    ブロック j とブロック k の
    ブロック間消去
  end do
  (c)
  ブロック k のブロック内消去
end do
(a) ブロック分割アルゴリズム (全体構成)

do l=(j-1)*p+1,j*p
  do m=(k-1)*p+1,k*p
    do i=l+1,n
      aim=aim-aij*ami
    end do
  end do
end do
(b) ブロック j とブロック k のブロック間消去

do l=(k-1)*p+1,k*p
  all=l/all
  do j=l+1,k*p
    alj=alj*all
  end do
  do j=l+1,k*p
    do i=l+1,k*p
      aij=aij-all*ali
    end do
  end do
end do
(c) ブロック k のブロック内消去

```

図5 ブロック分割アルゴリズム

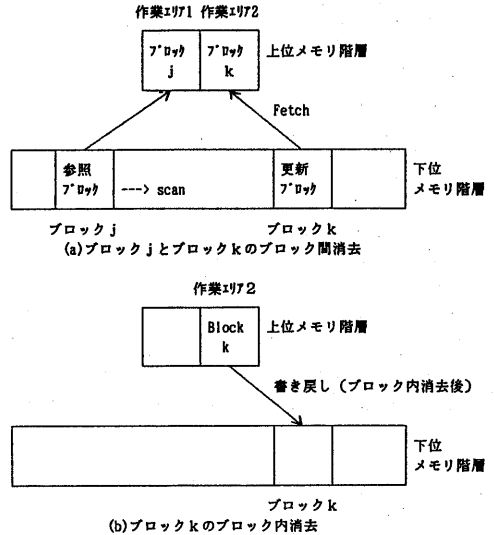


図6 メモリ階層間データ転送の様子

2階層のメモリ階層ではこの方法でデータ転送を最適化することができる。2階層メモリの上位メモリ(小容量で高速の方、たとえば主記憶)の容量が m であったとすれば $m/2$ をブロックサイズとして下位メモリ(たとえば拡張記憶)をブロック分割し、上の問題分割法を適用すればよい。

3.3 再帰的ブロック分割

3階層の場合にはこのブロック分割を再帰的に適用すれば下位のメモリ階層間から順にデータ転送を最適化できる。つまり係数行列に対してブロック分割の手法を適用することにより最下位とその一つ上位のメモリ階層間のデータ転送を最適化することができる。このときのブロック内消去については、ブロックをさらにサブブロックに分割することによってさらに上位のメモリとのデータ転送を最適化できる。ブロック間消去操作も簡単にサブブロックへの分割が可能である。これを模

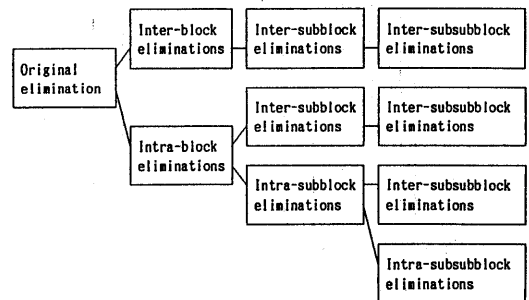


図7 ブロック分割

式的に表すと図7のようになる。4階層の場合にはサブブロックの処理をさらに細分化することにより対応が可能であり、この分割は原理的には分割後の最小のブロックが1列のブロックになるまで適用可能で、4階層以上の場合にも対応できる。

さらに、このブロック分割アルゴリズムでは係数行列の同じ列の要素はすべて同一のブロックに含まれるため、ピボットリングを簡単に行うことができる。

3. 4 再帰型データ分割方式を用いた並列化

前節で述べたブロック化を繰り返し用いてブロックの細分化を進めると処理の大半がブロック間消去に帰着される。たとえば $n \times n$ のサイズの係数行列に対してブロック化を繰り返し、あるレベルで分割されたサブブロックのサイズが n 行 p 列 ($n = p \times q$) であったとすると、処理すべきブロック内消去の演算量の合計は n が大きい場合には

$n^3 (2/3q^2 + 1/q - (q+1)/2q^2)$ となる。これは q (分割数) が大きくなり、 $1/q^2$ の項を無視すると $n^3/2q$ となる。ガウス消去法における全体の演算量は $2/3n^3$ であるから、仮に $q = 32$ とするとブロック内消去に要する演算量は全体のわずか 2.3% となる。したがって、主にブロック間消去を対象として並列化を行えば良いことになる。共有記憶のレベルで分割されたブロックに対するブロック間消去の並列化について説明する。

共有メモリ上に参照ブロックと更新ブロックがロードされ、このブロック間消去が並列化される様子を図8に示す。更新ブロック、参照ブロックをさらにサブブロックに分割し、各ローカルメモリに2つのサブブロック用の作業エリアが確保できるとする。まず、更新ブロックの先頭から4つの連続するサブブロックをローカルメモリにフェッチし、これに対して参照ブロックを先頭から順次同一サブブロックを各プロセッサの作業エリア1にブロードキャストする。更新サブブロックの参照サブブロックを用いた更新が終了すれば更新サブブロックを共有記憶に書き戻す。以上の処理を順次、更新ブロックの残りのサブブロックについて繰り返す。ここでも参照サブブロックの走査を内側ループとすることが、共有記憶への書き戻しデータ量を減らすために重要となる。

上でブロック内消去の演算量が少ないことを述べたが、さらにプロセッサ台数が増えた場合の並列化においてはブロック内消去の並列化も重要となる。本稿では省略するがブロック内消去についても並列化が可能であり、以下で示す実験につい

てはこの部分も並列化を行っている。

また、並列化はすべて係数行列の列を単位として行われるため、一列の内部の処理はベクトル化によりベクトル演算器での高速処理が可能である。

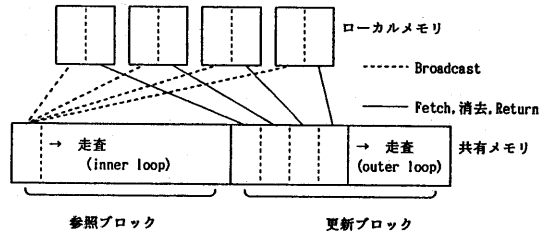


図8 ブロック間消去の並列化

4. 実験

4. 1 評価の前提条件

まず、LU分解の実測についての前提条件を示す。係数行列(初期データ)はプログラム中で生成しており、ディスク等から読み込んだ訳ではない。初期データをディスクにおいてスタートするとすれば、ディスクからLHS高速メモリ部へのデータ転送が余分にかかることになるが、ステージング制御によるプリロードを行えば、今回行った評価とはほぼ同等の時間で演算を実行できる(別ジョブとしてのプリロード処理に時間はかかる)。実験に用いた係数行列 A_{ij} は次式により生成した。

$$A_{ij} = \text{if } (i < j) \text{ then } i \text{ else } i+1$$

また、プログラムはスーパーコンピュータで開発された *Ph i l* 言語により記述した。LHSのアクセスはシステムで用意された関数で記述し、並列化は *Ph i l* 言語が提供する非同期サブルーチン呼び出しを用いて記述した。

4. 2 実測結果

LHS-HPPシステムを用いてLU分解を行なった場合の結果を表1に示す。16384×16384と32768×32768の係数行列に対して実測評価を行なった。32768×32768の実行においてはローカルメモリ各32MB、共有記憶768MB、LHS高速メモリ部4GB、LHSバッキングストア8GBを係数行列のためのメモリ領域として用いた。16384×16384の実行に際してはそれぞれの領域を4分の1として測定している。

表1により32768×32768の実行を10時間40分で実行できたことが分かる。このうちLHSのアクセスに要した時間は約2時間19分であり、全体の22%である。この内の大半がLHSのディスク-高速記憶部のデータ転送とLHSの初期化で占められている。LHSの初期化とはあるブロックへの最初のアクセス時に要するアドレス変

換テーブルへの登録処理などである。演算精度の評価は分解されたL*Uと分解前のAの行列式の差をAの行列式で割った値で行い、7桁以上の精度が得られている。

次に図8に16384×16384の実行における並列化加速率（対1PE性能向上比）を示す。実線は表1と同一条件で行った場合であり、4台で2.8倍の速度向上が得られている。データ転送を工夫すれば磁気ディスク装置を用いても4台のプロセッサを有効に利用できることが分かる。

破線はLHS高速メモリ部の容量を2GBとして、LHSバッキングストアにアクセスが行かない条件で測定を行ったものである。この場合はシリアル処理であるLHSアクセス時間が減少し、4台で3.2倍の速度向上が得られた。

5. まとめ

LHS-HPPシステムを用いた実測評価において、過去に例のない3万2千元余りの膨大な元数のLU分解を10時間40分で実行することができた。高速実行のポイントは4階層メモリの各階層間データ転送量の最適化と、並列実行にある。今回新たに開発したアルゴリズムはこれら2点のポイントを満たすものとして今後重要な役割を果

係数行列のサイズ		16384*16384	32768*32768
ELAPSE time (秒)		7586	38387
HPP内処理時間 (秒)		5505	30062
LHSアクセス (HIT時)	転送量 (GB)	69	274
	時間 (秒)	66	263
	速度 (MB/秒)	1040	1040
LHSアクセス (MISS HIT時)	転送量 (GB)	1.9	7.8
	時間 (秒)	1741	6965
	速度 (MB/秒)	1.1	1.1
LHS 初期化時間 (秒)		274	1097
Error Ratio for determinant		1.3×10^{-8}	5.2×10^{-8}

表1 実測結果(4PE実行) '90 1月 富士通(株)沼津工場にて実測
LHSでの(MISS HIT)はLHSバッキングストア-高速メモリ部間のデータ転送量、時間、データ転送速度を示す。

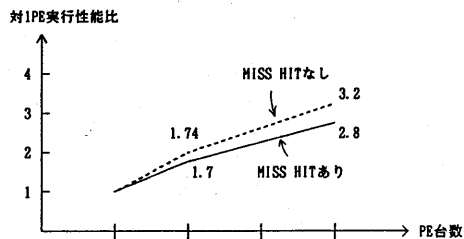


図9 16384×16384実行での並列化加速率

たすものとする。

謝辞

本研究は通商産業省工業技術院大型プロジェクトの一環として新エネルギー・産業技術開発機構(NEDO)から委託を受けて、実施したものである。

LU分解プログラムの開発において様々な助言を頂いた京都大学の津田教授と岡部助手に感謝致します。またLU分解プログラムのデバッグ、実機評価でお世話頂いた橋本氏を始めとする富士通(株)の方々には御礼申し上げます。

本研究の実施にあたりご指導頂いた、C&Cシステム研究所の石黒所長、大野主管研究員、コンピュータシステム研究部の小池部長、中崎課長に感謝します。また、本プロジェクト遂行にあたり様々なご協力を頂いた関係諸機関各位に深く御礼申し上げます。

参考文献

- [1] J. L. Larson, "Multitasking on the Cray X-MP-2 Multiprocessor", IEEE COMPUTER, Vol. 17, No. 7 (July 1984), pp. 62-69.
- [2] T. Watanabe and H. Matsumoto and P. D. Tannenbaum, "Hardware Technology and Architecture of the NEC SX-3/SX-X Supercomputer System", Proceedings of Supercomputing '89 (Reno, Nevada, November 1989), pp. 842-846.
- [3] T. Tsuda and Y. Seo, "Supercomputing External Multidimensional FFT", Journal of Information processing, Vol. 11, No. 2, 1988.
- [4] T. Tsuda and Y. Okabe, "Use of Semiconductor Extended Storage as Extended Main Storage for Large-Scale Supercomputing", Proceedings of the Second International Conference on Supercomputing Vol. 1 of three volumes, (Santa Clara, Calif., May 1987), pp. 176-183.
- [5] I. Bar-on, "A Practical Parallel Algorithm for Solving Band Symmetric Positive Definite Systems of Linear Equations", ACM Trans. Math. Softw., Vol. 13, No. 4 (December 1987), pp. 323-332.
- [6] V. Conrad and Y. Wallach, "Iterative Solution of Linear Equations on a Parallel Processor System", IEEE Trans. on Computers, Vol. C-26, No. 9 (September 1977), pp. 838-846.
- [7] 西ほか「大容量高速記憶装置(LHS)の基本方式」信学会CPSY90-6 1990 6月