

汎用目的マルチプロセッサ・システム OSCAR 上での スパース線形方程式求解の並列処理

笠原 博徳^{*} 中山 晴之^{*} 高根 栄二^{*} 橋本 親^{**}

^{*}早稲田大学 理工学部 電気工学科 ^{**}富士ファコム制御株式会社

本論文では、任意個のプロセッサ・エレメントから成るマルチプロセッサ・システム上で、スパース連立一次方程式を、最小時間で効率良く解く並列処理手法について述べる。本手法は、LU分解に基づくガウスの消去法やクラウト法などの、直接法の並列処理を目的とし、各プロセッサ・エレメントで実行されるマシン・コードの自動生成を行うところに特徴がある。すなわち、コンパイル時に筆者等が開発したスタティック・スケジューリング・アルゴリズム CP/MISF および DF/IHS を用いて、最適スタティック・スケジューリングを得ることにより全体の処理時間を最小にし、同時に諸々のオーバーヘッドを最小化するマシン・コードを生成する。本手法の有効性は、8086 と 8087 とを用いた実験用マルチプロセッサ・システム上で、さらに、筆者等が現在開発中であるマルチプロセッサ・システム OSCAR (Optimally Scheduled Advanced Multiprocessor) 上で検証される。

Parallel Processing for The Solution of Sparse Linear Equations on OSCAR (Optimally Scheduled Advanced Multiprocessor)

Hironori KASAHARA^{*}, Haruyuki NAKAYAMA^{*}, Eiji TAKANE^{*} and Shin HASIMOTO^{**}

^{*}Department of Electrical Engineering, Waseda University

3-4-1, Ohkubo, Shinjuku-ku, Tokyo, 160, Japan

^{**}FUJIFACOM Co. Ltd.

1-1, Fuji-cho, Hino-shi, Tokyo, 191 Japan

This paper describes an efficient parallel processing scheme for the solution of sparse linear equations on multiprocessor supercomputing system composed of arbitrary number of processor elements. The parallel processing scheme is aimed at the parallel processing for direct solution methods such as Gaussian elimination algorithm and Crout algorithm based on LU factorization. In this scheme, a variety of overheads are also minimized by using the static scheduling algorithms CP/MISF and DF/IHS developed by the authors to obtain the optimal schedule at the stage of computation. The effectiveness of the proposed scheme is demonstrated on an experimental multiprocessor system using Intel 8086 and 8087, and on OSCAR (Optimally Scheduled Advanced Multiprocessor), a prototype multiprocessor supercomputing system being developed by the authors to extract advantageous features of static scheduling to the maximum extent.

1. まえがき

電子回路動特性解析、あるいは電力システムにおける電力潮流計算や過渡安定度計算等の種々の分野において、スパース線形方程式の高速な求解が要求されている。この要求を満たすため、従来より、大規模な連立一次方程式の並列処理解法について、多くの試みがなされてきた[1]-[8][20]。これまで発表された論文で取り上げられている並列処理システムのアーキテクチャは、ベクトル・プロセッサ[1][2]、シストリック・アレイ[3]、データフロー・マシン[4]、およびマルチプロセッサ[5]-[8][20]に分類することができ、また、取り上げられている解法は、直接法[5]-[8]と繰り返し法[4][21]とに分けることができる。本論文では、MIMD型のマルチプロセッサ・システム上での直接法の並列処理を扱う。

従来このテーマに対する並列処理手法もいくつか提案されており、例えば、ガウスの消去法を並列処理する際のスケジューリング法について、係数行列がスパース行列である場合をWing等が[5][6]、密行列であるような場合をSrinivasが[7]それぞれ提案している。また、Arnold等は、LU分解後の前進および後退代入のみに適用範囲を限定した、ダイナミック・スケジューリングを用いた並列処理手法を提案している[8]。

この中で、Wing等は、1回の除算および更新演算を1つのノードで、また、演算間の先行制約を2つのノードを結ぶアークでそれぞれ表すことにより、ガウスの消去法で要求される計算負荷を無サイクル有向グラフ(DAG)によって表現できることを示し、さらに、もしすべての除算および更新演算が1ユニット・タイムで実行されるならば、Hのアルゴリズム[9]を用いることによって、最適スケジュールが得られる場合が多いことを示した。これは、線形方程式で解く過程を無サイクル有向グラフで表し、それにスケジューリング理論が適用できることを示した点では大変評価されるものではあるが、すべての除算および更新演算の実行時間が1ユニット・タイムであるという仮定は非現実的である。一般に、更新演算は乗算と減算とを1回ずつ含むが、1回の除算の方が長くかかる。

Srinivasの方法は、係数行列がスパース行列である場合に適用できない点で問題があるし、これもまたWing等の方法と同様に、すべての演算の処理時間が1ユニット・タイムであるという非現実的な仮定を用いている。

一方、Arnold等の方法では、前進および後退代入のみ扱われており、LU分解の過程は考慮されていない。しかも、彼らが提案しているのは、並列処理の過程における演算処理時間の変動を考慮した、実時間のダイナミック・スケジューリング法である。しかし、電子回路解析で線形方程式を解く場合等は、係数行列の構造は全く変えずに要素の値のみを変えて繰り返し解くのが普通であり、その並列処理で重要となるのは、むしろLU分解の過程である。また、筆者等は、代入文のみから成る計算を処理する場合には、ダイナミック・スケジューリングよりもスタティック・スケジューリングを用いた方が効率が良いことを既に示している[18]。すなわち、スタティック・スケジューリングでは、コンパイル時に最適スケジュールを決定し、さらに、レジスタ利用の最適化等が行なえるため、実行時のオーバーヘッドが小さい効率的な並列処理が可能となる。

以上の事を踏まえ、本論文では、線形方程式の処理全体を、タスク(プロセッサ・エレメントへの割当て単位)処理時間が等しいというような制約の無い、一般的なスケジューリング問題として扱う。この異なる処理時間と任意の先行制約(タスク間の実行順序関係の制約)を持つタスク集合を、任意個のプロセッサ・エレメントに割当てるスケジューリング問題は、多くの研究者により15年以上の長期に渡り研究されてきたのにもかかわらず、最適化アル

ゴリズムが開発されていなかった非常に難しい問題であり、強NP困難[11]であることが知られている。

しかし、筆者等は、この問題に対して最適解あるいは高精度の近似解を得ることができる2つの実用的な最適化スケジューリング・アルゴリズムCP/MISFおよびDF/HS[12]を開発することに成功した。そして、スカラ代入文の処理が大部分を占める常微分方程式の解析と、小行列・小ベクトル計算を扱うロボット・アームの制御計算とに対して、これらのスケジューリング・アルゴリズムを応用することにより、その有用性を実証した[10][13]。

以下では、筆者等が開発したこれらのスケジューリング・アルゴリズムを用いた、LU分解に基づくガウスの消去法およびクラウト法の並列処理手法について述べると共に、この手法の有効性を、実験用マルチプロセッサ・システム上とスタティック・スケジューリングの性能を最大限に引き出すことを可能とするマルチプロセッサ・システムOS CAR (Optimally Scheduled Advanced Multiprocessor) 上で検証を行った結果について述べる。

2. 線形方程式の解法

本章では、次の連立一次方程式をLU分解に基づくガウスの消去法およびクラウト法を用いて解く場合を考える。

$$AX=B \quad (1)$$

ただし、Aは一般的なスパース行列であり、帯行列やブロック対角行列である必要はない。

これらのアルゴリズムは、大きく分けて2つの部分から成っている。すなわち、LU分解の部分

$$LU=A \quad (2)$$

と、その後の、2つのアルゴリズムに共通な前進及び後退代入の部分

$$LY=B \quad (3)$$

$$UX=Y \quad (4)$$

とである。ここで、行列A、L、およびUの各要素を、それぞれ a_{ij} 、 l_{ij} 、および u_{ij} で表す。ただし、 $1 \leq i, j \leq N$ において、 $l_{ij}=0$ ($i < j$)、 $u_{ij}=0$ ($i > j$)であり、 $u_{ii}=1$ ($1 \leq i \leq N$)であるとする。このとき、ガウスの消去法およびクラウト法におけるLU分解の過程は、それぞれ以下のように書くことができる。

```
Procedure GE
/Triangulation by Gaussian elimination algorithm/
begin
  for k=1 step 1 until N do
  begin
    for all  $a_{kj} \neq 0, k+1 \leq j \leq N$ , do
       $u_{kj} = a_{kj}/a_{kk}$ ;
    for all  $a_{ik} \neq 0, k \leq i \leq N$ , do
       $l_{ik} = a_{ik}$ ;
    for all  $a_{ij} \cdot u_{kj} \neq 0, k+1 \leq i, j \leq N$ , do
       $a_{ij} = a_{ij} - a_{ik} \cdot u_{kj}$ ;
  end
end
```

```
Procedure CR
/Triangulation by Crout algorithm/
begin
  for k=1 step 1 until N do
  begin
    for all not [ $a_{km}=0$  or
       $l_{km} \cdot u_{mj}=0$  ( $1 \leq m \leq k-1$ )],
       $k \leq i \leq N$ , do
       $l_{ik} = a_{ik} - \sum_{m=1}^{k-1} l_{im} \cdot u_{mk}$ ;
  end
```

$$\text{for all not } [a_{kj}=0 \text{ or } l_{km} \cdot u_{mj}=0 \text{ (} 1 \leq m \leq k-1 \text{)}], k+1 \leq j \leq N, \text{ do}$$

```
begin
   $a_{kj} = a_{kj} - \sum_{m=1}^{k-1} l_{km} \cdot u_{mj}$ ;
end
```

$$u_{kj} = a_{kj}/l_{kk} \quad (9)$$

```
end
end
```

また、前進および後退代入の過程は、それぞれ下のよう
に書ける。

```

Procedure FS
/Forward Substitution/
begin
  for k=1 step 1 until N do
  begin
    if bk≠0 then
      yk=bk/lkk; (11)
    for all lik·yk≠0, k+1≤i≤N, do
      bi=bi-lik·yk; (12)
    end
  end
end
Procedure BS
/Backward Substitution/
begin
  for k=N step -1 until 1 do
  begin
    if yk≠0 then
      xk=yk; (13)
    for all uik·xk≠0, 1≤i≤k-1, do
      yi=yi-uik·xk; (14)
    end
  end
end

```

Wing[5]とHuang[6]とは、(5)式で表される除算 (di
vide operation) と (7) 式で表される更新演算 (updat
e operation) との処理時間が、共に1ユニット・タイム
であると仮定したが、この仮定は非現実的である。一般的
には除算の方が更新演算よりも処理時間を多く必要とする。

本論文では、そのような非現実的な仮定は設定せずに、
二つの演算の処理時間はそれぞれ任意であるとする。さら
に、(8)式や(9)式のように各式の処理時間が非零要
素の位置により変化するクラウト法の並列処理も取り扱う。
しかし、ここでは、フィル・インを最小化したり[15]並列
性を最大にする[6]リオーダーリングや、丸め誤差に関する
議論は行わない。すなわち本論文では、適当なりオーダ
リングがあらかじめなされていると仮定する。言い替える
と、本論文で述べる並列処理は、SPICEのような従来の
電子回路シミュレータのコード生成法[14]に基づいて生
成された、FORTRAN算術代入文の集合の並列処理と
見なすことができる。

3. 並列処理手法

本章では、任意のスパース線形方程式を、任意個のプロ
セッサ・エレメントから成るマルチプロセッサ・システム
上で最小時間で解く並列処理手法について述べる。この手
法は、筆者等が開発した最適スタティック・スケジューリ
ング・アルゴリズム[12]を用いる点を最大の特徴とし、タ
スク生成 (グラニュラリティの決定)、タスク・グラフ表
現、タスク・スケジューリング、およびマシン・コード生
成の4ステップから構成されている。以下、各ステップの
詳細を述べる。

3.1 タスク生成

線形方程式を並列処理するためには、処理すべきすべ
ての計算をタスク (プロセッサ・エレメントへの割当ての単
位) に分割しなければならない。この時、並列性をできる
限り引き出す一方、データ転送やタスク間の同期に関連し
たオーバーヘッドを低く抑えることができるように分割す
る必要がある。ガウスの消去法やクラウト法では、スカ
ラ演算レベル、代入レベル、行および列単位の演算レベル、
あるいはこれらの中間のレベルなどのタスクのグラニュ
ラリティが考えられる。一見、グラニュラリティがもつとも
細かいスカラ演算レベルの場合が、並列性を最大限に引き
出して、最小の並列処理時間が得られるように思えるが、
必ずしもそうではない。このような分割は、オーバーヘッ
ドの増大を招くからである。このため、各々のマルチプロ
セッサ・システムの処理能力、プロセッサ間のデータ転送
能力、同期のオーバーヘッド、あるいはタスク・スケジ
ューラの性能などを十分考慮して、タスクのグラニュラリ

ティを決定しなければならない。筆者等は、筆者等が開発し
たスタティック・スケジューリング・アルゴリズムを用いる際
のタスクのグラニュラリティ決定に関する指針を、既に提
案している[12][13][18]。この指針によれば、タスクの平
均処理時間が、プロセッサ間のデータ転送や同期のオー
バーヘッドの時間のタスク数平均の3、4倍以上であること
が望ましい。

本論文では、ガウスの消去法には代入レベルのグラ
ニュラリティを、クラウト法には代入レベルの他にスカ
ラ演算レベルのグラニュラリティを、それぞれ選択した (た
だし以下で提案する手法は、どのようなグラニュラリティ
にも対応し得る。)。たとえば、図1 (a) の線形方程式
に対して、ガウスの消去法の場合は図 (b)、クラウト法
の場合は図 (c) のようなタスクが生成される。

一般に、実システム上では、行列AとL、U、およびベ
クトルBとX、Yは、それぞれ同じ領域に割付けるのが普
通であるので、このタスク生成において、(6)式、(8)
式、(9)式、および(13)式が、a_{ik}やy_kの単なる
コピーとなる場合には、その処理はタスクとして扱ってい
ない。図1に示すように、代入レベルのグラニュラリ
ティでタスク生成したガウスの消去法の場合 (以下、GE)、
LU分解と前進および後退代入とを合わせて、20個のタ
スクが生成されている。クラウト法の場合、代入レベル
のグラニュラリティで生成した場合 (以下、CO) は19
個、演算レベルにより近いグラニュラリティの場合 (以下、
CS) は22個のタスクが生成されている。後者の場合、
(8)式及び(9)式の右辺の総和の計算は、図1 (c)
に示すように、いくつかの乗算と加算とに分割される。

このとき、総和の中の積の項が奇数個の場合、例えば、
 $l_{54} = a_{54} - (l_{51} \cdot u_{14} + l_{52} \cdot u_{24} + l_{53} \cdot u_{34})$
のような場合には、右辺の計算は、次のように最適にタ
スク分割される。

$$t_1 = l_{51} \cdot u_{14}, \quad t_2 = l_{52} \cdot u_{24}, \quad t_3 = l_{53} \cdot u_{34}$$

$$t_4 = a_{54} - t_1, \quad t_5 = t_2 + t_3, \quad l_{54} = t_4 - t_5$$

$$A = \begin{pmatrix} x & x \\ xx & F \\ xxx & x \\ & x \\ & & x \end{pmatrix} \quad B = \begin{pmatrix} x \\ x \\ x \\ x \\ x \end{pmatrix}$$

x: 非零要素
F: フィル・イン

(a) 線形方程式の例

【LU分解】

- ① u₁₄ = a₁₄ / a₁₁
- ② a₂₄ = -a₂₁ · u₁₄
- ③ a₃₄ = -a₃₁ · u₁₄
- ④ u₂₄ = a₂₄ / a₂₂
- ⑤ a₃₄ = a₃₄ - a₃₂ · u₂₄
- ⑥ u₃₄ = a₃₄ / a₃₃
- ⑦ a₅₄ = -a₅₃ · u₃₄

【前進消去】

- ⑧ y₁ = b₁ / l₁₁
- ⑨ b₂ = b₂ - l₂₁ · y₁
- ⑩ b₃ = b₃ - l₃₁ · y₁
- ⑪ y₂ = b₂ / l₂₂
- ⑫ b₃ = b₃ - l₃₂ · y₂
- ⑬ y₃ = b₃ / l₃₃
- ⑭ b₅ = b₅ - l₅₃ · y₃
- ⑮ y₄ = b₄ / l₄₄
- ⑯ b₅ = b₅ - l₅₄ · y₄
- ⑰ y₅ = b₅ / l₅₅

【後退代入】

- ⑱ y₃ = y₃ - u₃₄ · x₄
- ⑲ y₂ = y₂ - u₂₄ · x₄
- ⑳ y₁ = y₁ - u₁₄ · x₄

(b) ガウスの消去法による
タスク

【LU分解】

- ① u₁₄ = a₁₄ / l₁₁
- ② a₂₄ = -l₂₁ · u₁₄
- ③ u₂₄ = a₂₄ / l₂₂
- ④ a₃₄ = - (l₃₁ · u₁₄
+ l₃₂ · u₂₄)

⑤ u₃₄ = a₃₄ / l₃₃

⑥ l₅₄ = -l₅₃ · u₃₄

【前進消去】 (⑦~⑩)

(b) ⑨~⑩と同じ

【後退代入】 (⑪~⑬)

(b) ⑪~⑬と同じ

(c) クラウト法によるタスク

- (4.1) t₁ = l₃₁ · u₁₄
- (4.2) t₂ = l₃₂ · u₂₄
- (4.3) a₃₄ = - (t₁ + t₂)

(d) (c) のタスク④右辺の
積和の分解

図1: タスク生成の例

3.2 タスク・グラフ表現

前節のような方法で生成されたタスク間には、フロー依存や出力依存といったデータ依存関係[16]がある。線形方程式の求解においては、出力依存関係のあるタスク間には、常にフロー依存関係が存在するので、後者のみで先行制約を決定できる。この先行制約は、「タスク・グラフ」[5][10][12]と呼ばれる無サイクル有向グラフ $G(N, A)$ で表現できる。ここで、 N は n 個のノードの集合を、 A はアークの集合をそれぞれ表す。このグラフ中、各ノードは1つのタスクを表し、ノード間のアークはタスク間の先行制約を表す。このグラフは、入口ノードと出口ノードとを一つずつ持ち、すべてのノードにこの二つのノードから到達し得るものとする。例えば、タスク集合図1(b)のタスク・グラフは、図2の様になる。

図中、ノード中の数字はタスク番号 i を、その外の数字はそのタスクのプロセッサ・エレメント上での処理時間をそれぞれ示す。ノード N_i から N_j に引かれたアークは、タスク T_j が T_i に先行するという半順序制約、すなわち先行制約を表す。また、二重線はクリティカル・パスを表す。クリティカル・パス長 m_c とは、与えられたタスク・グラフを任意台数のプロセッサで並列処理する際の、オーバーヘッドをまったく無視した場合の最小実行時間である[12]。ここで、タスクの処理時間といっても浮動小数点演算の処理時間はその数値によって変化するので、一般的には確定的な値を得ることができないという問題があるが、この問題は、各処理の平均所用時間を用いることにより解決できる[12][18]。この平均処理時間を用いた時、次節で述べるスケジューリング・アルゴリズムを用いると、タスク集合の平均並列実行時間を最小化できる。

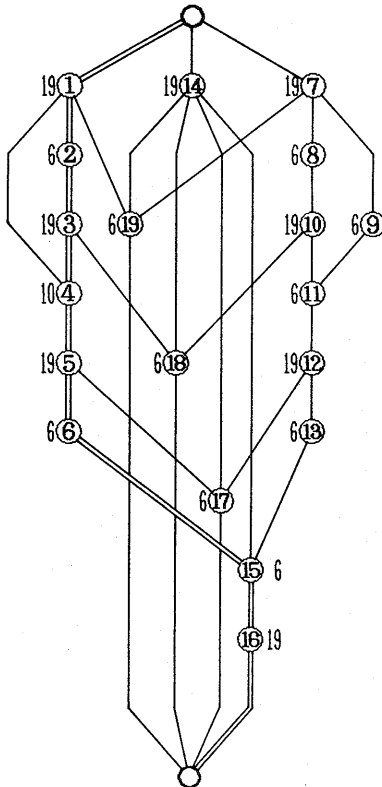


図2；タスクグラフ表現

3.3 スケジューリング・アルゴリズム

マルチプロセッサ上で、任意のタスク集合を効率良く並列処理するためには、各プロセッサへのタスクの割当てと一つのプロセッサに割当てられたタスク間の実行順序とを最適に決定しなくてはならない。最適な割当てと実行順序とを決定するこの問題は、並列処理時間あるいはスケジューリング長の最小化を目的関数とするマルチプロセッサ・スケジューリング問題として、長い間研究されてきたものである[10][12]。この問題を正確に定義すれば、能力の等しい m 台のプロセッサ、そこで計算すべき n 個のタスクから成るタスク集合 $T = (T_1, \dots, T_n)$ 、およびタスク間の先行制約が与えられたときに、その実行時間あるいはスケジューリング長を最小とする、ノンプリエンティブ・スケジューリング問題となる。しかし、この問題は、強NP困難問題として知られている非常に難しい問題であり、 $P = NP$ でなければ、擬多項式時間最適アルゴリズムばかりでなく、両完全多項式時間近似アルゴリズムさえも存在しない事が知られている。しかし、筆者等は、ヒューリスティック・アルゴリズム $CP/MISF$ と実用的な最適化アルゴリズム DF/IHS の開発に既に成功している[12]。前者は、その時間複雑度の低さから、短時間で近似解を得ることができるものであり、後者は、 $CP/MISF$ と深さ優先探索とを組み合わせたもので、最適解あるいは最適解との精度が保証されている近似解を得ることができるものである。スケジューリングを行なうコンピュータの処理性能が低い場合や、タスク数が多い場合等は、 $CP/MISF$ が適している。詳細は、文献[12]を参照して頂きたい。

3.4 マシン・コード生成

実マルチプロセッサ・システム上で効率良い並列処理を実現するためには、スケジューリングの結果を考慮した最適なマシン・コードを、そのシステムに合わせて生成しなければならない。スケジューリングの結果を見れば、1つのプロセッサ・エレメントで実行されるタスクおよびそのタスク間の実行順序、別のプロセッサに割当てられたタスクからデータを受け取るためのおよびその待ち時間、どのタスクとどのタスクとの間で同期を取る必要があるか等といった情報が得られる。これらの情報に従って、1つのプロセッサに割当てられたタスク用のコードを生成し、同期やデータ転送のためのコードを付け加えることによって、各プロセッサ用のマシン・コードを生成することができる。さらに、上述の情報を最大限に活用することにより、種々のオーバーヘッドを最小化し、コードを最適化することができる。例えば、タスクの割当てや実行順序の情報を活用すれば、プロセッサ内で行なわれるタスク間のデータ授受のためのレジスタの使用を最適化することができ、これはダイナミック・スケジューリングを用いるシステムでは不可能である。その結果、処理時間を大幅に減らすことができる。また、データの待ち時間の期待値が分かれば、データ待ちのタスクがそのデータの有無を確認する(データ・レバ同期)ために頻繁にバス・アクセスする結果生じるバスのデータ転送能力低下を防止することができる。さらに、同期を取るタスク、タスクの割当てや実行順序の情報を注意深く考慮すれば、同期のオーバーヘッドを最小化することも可能である。たとえば、図3に示すように、タスクA、BおよびCがプロセッサ1に、Dがプロセッサ2にそれぞれ割当てられているとする。このとき、タスクBはタスクAの完了を示すフラグをチェックする必要はない。この2つのタスクは、共に同じプロセッサに割当てられているからである。BC間も同様である。また、タスクBがDの終了をチェックしているのでCはDの終了をチェックする必要は無い。

各プロセッサ・エレメントに対して上述のように生成さ

れた最適マシン・コードは、各プロセッサ・エレメントのローカル・メモリにロードされ、非同期に実行される。以上の並列処理手法の4つのステップは、ホスト・コンピュータで自動的に実行される。

4. 実験用マルチプロセッサ上での性能評価

本章では、図4に示す実験用マルチプロセッサ・システム上でスパース線形方程式を並列処理した結果を示す。このシステムは、スタティック・スケジューリング・アルゴリズムを用いた並列処理手法の実用性と有用性を実証するために開発したものであり、このため非常に単純なハードウェア構成となっている。このように単純なシステムで我々が提案する手法を効率良く並列処理できたならば、この手法の有用性が証明されるばかりでなく、さらに高性能なマルチプロセッサ・システム上ではもっと効率良く並列処理できることが約束されることになる。図4に示すように、このシステムは、7つのプロセッサ・エレメント (PE)、1つの共有メモリ (CM)、単一バス、およびホスト・コンピュータ (HC) と PE 部間のデータ転送を制御するシリアル通信ユニット (SCU) から成る。各 PE は、16ビットマイクロプロセッサ、インテル8086 (クロック: 5MHz)、8086のコプロセッサで浮動小数点演算用の8087およびローカル・メモリ (16KB ROMと32KB RAM) から成る。各PE間のデータ転送は、CMを通してのみ可能である。タスク間の同期には、「バージョン・ナンバー」方式[13][18]が採用されている。本手法では、行列やベクトルの一つの要素が格納されている同一アドレスへの「書き込み」回数が、バージョン・ナンバーに相当する。ライタタスクは、転送すべきデータを共有メモリに書き込んだ後に、そのデータのバージョン・ナンバーをそのタスクにあらかじめ割り振られたナンバーに書き換え、一方リーダタスクは、目的とするデータのバージョン・ナンバーがそのタスクにあらかじめ割り振られたナンバーと同じかどうかをチェックする。バスは、16ビットのデータ幅を持ち、そのデータ転送能力は、1MW/sである。バス・アクセスは、各PEとSCUとに事前に割り振られたバス・アクセス優先順位に基づく分散バス・アービタ形式で制御されている。

このマルチプロセッサ・システム上で1つのスパース線形方程式を並列処理した結果を図5に示す。この例のスパース線形方程式の係数行列Aは、図6のように 20×20 で非零要素を10%含むものであり、GEでは64個、C0では65個、また、CSでは48個のタスクが、それぞれ生成される。

図5において、3本の点線は、上記の線形方程式を3通りの方法で解くときの、各々のタスクのスケジューリング長、すなわち、オーバーヘッドが全く無いとしたときの実行時間の理論値を示している。これに対して実線は、実マルチプロセッサ・システム上での実行時間の実測値を示す。縦軸には実行時間を、横軸にはPEの台数をそれぞれとっており、ここで、実行時間の理論値 (スケジューリング長) は、各PE台数の各アルゴリズムに対する最小実行時間を示している。この理論値は、プロセッサ台数の逆数にほぼ比例して減少し、プロセッサが6、7台のときにクリティカル・パス長 (GEとC0の894 [μ s]、CSの787 [μ s]) に一致する。

実線は、スケジューリング結果を基に生成したマシン・コードの実マルチプロセッサ・システム上での並列処理時間の実測値を表す。プロセッサ台数が1台のときの実行時間は、GE、C0およびCSの各アルゴリズムに対して、それぞれ5.14 [ms]、5.40 [ms] および4.82 [ms] であるのに対して、3台になると、それぞれ2.70 [ms] (1/1.90)、2.70 [ms] (1/1.9) および2.29 [ms] (1/2.10) に、7台では、そ

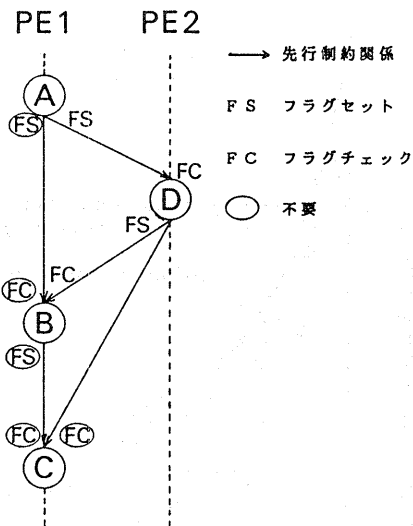


図3; タスク間同期オーバーヘッドの最小化

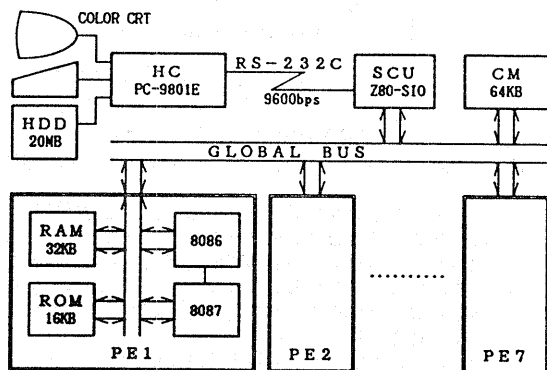


図4; 実験用マルチプロセッサ・システムの構成

れぞれ1.60 [ms] (1/3.21)、1.58 [ms] (1/3.42) および1.37 [ms] (1/3.52) にまで減少する。実測値とスケジューリング長との間の差は、プロセッサ間のデータ転送、タスク間の同期やバス・アクセス競合といった処理上の種々のオーバーヘッドを表している。一見、このオーバーヘッドは小さくないように見えるが、このマルチプロセッサ・システム上で共有メモリを介して2プロセッサ間で64ビットのデータを転送するのに72.6 [μ s] 要するのに対し、各アルゴリズムにおけるタスク平均処理時間が、それぞれ84 [μ s]、83 [μ s] および100 [μ s] と小さく、システムのデータ転送性能が低いことを考えると、この結果は、極めて満足できるものであり、この手法の有用性と実用性が確認されたと言えることができる。また、もしプロセッサ間でデータをロードキャストする能力があれば、ピボット要素 a_{ij} をすべてのプロセッサに同時に転送することができるし、他のプロセッサのローカル・メモリに直接アクセスできれば、2台のプロセッサ間でデータ転送するためにバスと共有メモリとをアクセスする回数を2回 (共有メモリへの1回の「書き込み」とそれから1回の「読み取り」) から1回 (他のプロセッサのローカル・メモリへの1回の「書き込み」) へ減らすことができるので、さらにオーバーヘッドを減らすことが可能となる。

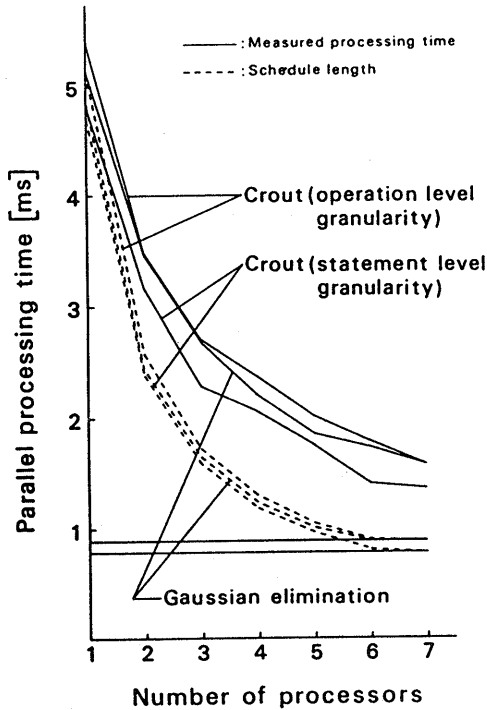


図5；実験用マルチプロセッサ・システム上での並列処理

ここで注意しなければならないのは、ここで採用した3つのアルゴリズムの中でタスクの大きさが最も大きいCSの場合が、最も効率良い並列処理を実現していることである。すなわち、タスクの粒度を小さくすればするほど並列実行時間が最小化するとは限らないということが示されている。

以上より、提案する並列処理手法を用いて実マルチプロセッサ・システム上で効率良く並列処理を行うことが可能であると結論できる。このように、長い間非実用的であると考えられてきたスタティック最適スケジューリングが実用的な手法であることを世界で初めて示したことは、非常に意義のあることと思われる。

5. OSCAR上での性能評価

本章では、現在筆者等が開発中であるプロトタイプ・マルチプロセッサ・システムOSCAR上での、本並列処理手法の性能評価について述べる。OSCARは、図7に示すように、プロセッサ・クラスタを複数持つ階層的なマルチプロセッサ・システムである。OSCAR開発の目的は、従来効率的な並列処理が難しいと考えられていた、スカラ代入文の多いFortranプログラムやスパース線形方程式求解等のアプリケーションの並列処理を、スタティック・スケジューリングとダイナミック・スケジューリングとの長所を最大限に引き出すことにより、効率良く行うことである。

現在1プロセッサ・クラスタのハードウェアがほぼ完成しており、このシステム上で本並列処理手法をインプリメントした。1つのプロセッサ・クラスタ(PC)は、16個のプロセッサ・エレメント(PE)、3個の共有メモリ、1個のローカル・コントロール・プロセッサ、および3本の共有バスから成る。各PEは、浮動小数点演算を含むすべての命令を1クロック(200ns)で実行するRIS

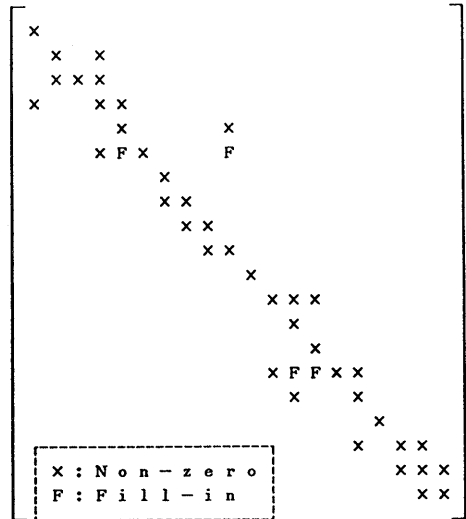


図6；線形方程式の係数行列の例

Cライクな32ビットカスタムメイド・プロセッサ、64個の汎用レジスタ、256KWのローカル・データ・メモリ、他のPEと通信するための2KW2ポート・メモリ、128KWの命令メモリ2バンク、およびDMAコントローラから成る。このDMAコントローラを用いれば、実行中に、共有メモリや他のPEの2ポート・メモリへの高速データ・ブロック転送や、共有メモリから1つの命令メモリ・バンクへの命令セットのダイナミックなローディングが可能となる。プロセッサ間のデータ転送のために、ブロードキャスト・モード、他のPEの2ポート・メモリへの直接データ転送モード、および共有メモリを介する間接データ転送モードの3種類のデータ転送モードが用意されている。それぞれのモードは、1ワード・データ転送にもブロック・データ転送にも使うことができる。各共有メモリに、3本のバスから同時にアクセスしても構わない。3本のバスのデータ転送能力は、総合すると最大60MW/sである。こうしたデータ転送能力は、前章で述べた実験用マルチプロセッサ・システムと比べて、データ転送オーバーヘッドを著しく軽減できるものと考えられる。さらに、OSCARの各PEはすべての命令を1クロックで実行できるため、従来問題となっていたタスク処理時間の推定が容易となり、厳密なスケジューリングを行えると共に、タスク間データ授受のためのレジスタ利用の最適化、データ転送モード及びタイミングの最適化によるデータ転送オーバーヘッドの最小化等、クロックレベルの厳密な最適化が可能となる。これらの最適化により生成されたマシンコードの実実行イメージを図8示す。図8中の記号をPE3の例に説明すると、まず[LD30][L25R][WAIT][→PE5]とあるが、これは、「タスク30のデータをメモリからレジスタにロードし」、「タスク25を実行後その出力データをレジスタに残し」、「バスの空くのを待った」後、「その出力データをPE5に1対1転送すること」を表わしている。ここでバスが空くのを待っているのは、PE1がブロードキャスト([BC]、バスを3本共占有する)を行なっているためである。次の[U2526R]は、「タスク25がレジスタに残したデータを使用してタスク26を実行後そのデータをレジスタに残す」ことを表わしている。その下の方の[WAIT][FC44][38R]では、「PE5からタスク44の出力デー

ーヘッドが通常時の2倍程度になっていると思われる。この通常状態での処理時間を上述のシミュレーション法により予測したのが図中の一点鎖線であり、7PEで1PEの場合の約1/3.25(40[μs])となり、さらに低オーバーヘッドの並列処理が可能となることがわかる。

5. 結び

本論文では、スタティック最適スケジューリングを用いたスパース線形方程式の並列処理手法を提案すると共にその有効性・実用性を汎用16ビット・マイクロプロセッサを用いた実験用マルチプロセッサ・システム上と現在開発中のプロトタイプ・マルチプロセッサ・システムOSCAR上で検証した。スパース線形方程式の効率良い並列処理は従来難しいとされてきたにもかかわらず、本論文では本手法を用いて、任意個のプロセッサから成る2種の異なるアーキテクチャのマルチプロセッサ・システム上で効率良く、すなわち最小実行時間で並列処理することが可能であることを示した。

さらに、この論文によって、長い間非実用的であると考えられてきたスタティック・スケジューリング・アルゴリズムが、実マルチプロセッサ・システム上で効率良い並列処理をするために非常に有用な手法であることが示された。このことは並列処理とスケジューリング理論との両方の研究にとって非常に意義のあることと考えられる。

また、OSCARは現在まだ動作確認テスト中でバス性能を本来の半分に落として実行しており、今後通常の動作モードで実行すれば、データ転送オーバーヘッドが更に小さくなり、より効率良い並列処理が実現できるものと予想される。

参考文献

- [1] D.A.Calahan and W.G.Ames, "Vector processors: models and applications," IEEE Trans. Circuits Syst., vol.CAS-26, pp.726-732, Sept. 1979.
- [2] F.Yamamoto and S.Takahashi, "Vectorized LU decomposition algorithms for large-scale circuit simulation," IEEE Trans. Computer-Aided Design, vol.CAD-4, pp.232-238, July 1985.
- [3] H.T.Kung and C.E.Leiserson, "Algorithm for VLSI processor arrays," in Introduction to VLSI Systems C.Mead and L.Conway, Eds. Reading, MA: Addison-Wesley, 1980, pp.271-292.
- [4] D.A.Reed and M.L.Patrick, "Iterative solution of large, sparse linear systems on a static data flow architecture: performance studies," IEEE Trans. Comput., vol.C-34, pp.874-880, Oct. 1985.
- [5] O.Win and J.W.Huang, "A computation model of parallel solution of linear equations," IEEE Trans. Comput., vol.C-29, pp.632-638, July 1980.
- [6] J.W.Huan and O.Wing, "Optimal parallel triangulation of sparse matrix," IEEE Trans. Circuit Syst., vol.CAS-26, pp.726-732, Sept. 1979.
- [7] M.A.Srinivas, "Optimal parallel scheduling of Gaussian elimination DAG's," IEEE Trans. Comput., vol.C-32, pp.1109-1117, Dec. 1983.
- [8] C.P.Arnold, M.I.Parr and M.B.Dewe, "An efficient parallel algorithm for the solution of large sparse linear matrix equations," IEEE Trans. Comput., vol.C-32, pp.265-273, Mar.1983.
- [9] T.C.Hu, "Parallel sequencing and assembly line problems," Operations Res. vol.9, pp.841-848, 1961
- [10] E.G.Coffman Jr.(ed.), Computer and Job-shop

Scheduling Theory. New York: Wiley, 1976.

- [11] M.R.Garey and D.S.Johnson, Computers and Intractability: A Guide to the Theory of NP-Completeness. San Francisco: Freeman, 1979.
- [12] H.Kasahara and S.Narita, "Practical multi-processor scheduling algorithms for efficient parallel processing," IEEE Trans. Comput., vol.C-33, pp.1023-1029, Nov. 1984.
- [13] H.Kasahara and S.Narita, "Parallel processing of robot-arm control computation on a multi-microprocessor system," IEEE J. Robotics and Automation, vol.RA-1, pp.104-113, June 1985.
- [14] F.G.Gustavson, W.Liniger and R.A.Willoughby, "Symbolic generation of an optimal Crout algorithm for sparse systems of linear equations," J.ACM, vol.17, pp.87-109, Jan. 1970.
- [15] R.D.Berry, "An optimal ordering of electronic circuit equations for sparse matrix solution," IEEE Trans. Circuit Theory, vol.CT-18, pp.40-50, Jan. 1971.
- [16] D.A.Padua, D.J.Kuck and D.H.Lawrie, "High-speed multiprocessor and compilation techniques," IEEE Trans. Comput., vol.C-29, pp.763-776, Sep. 1980.
- [17] E.B.Fernandez and B.Bussel, "Bounds on the number of processors and time for multiprocessor optimal schedules," IEEE Trans. Comput. vol.22, pp.745-751, Aug. 1973.
- [18] H.Kasahara and S.Narita, "An approach to supercomputing using multiprocessor scheduling algorithms," in Proc. IEEE 1st International Conf. on Supercomputing Systems, pp.139-148, Dec. 1985.
- [19] C.V.Ramamoorthy, K.M.Chandy and M.J.Gonzalez, Jr., "Optimal scheduling strategies in a multiprocessor system," IEEE Trans. Comput., vol.C-21, pp.137-146 Feb. 1972.
- [20] H.Kasahara and S.Narita, "Load distribution among real-time control computers connected via communication media", in Proc. IFAC 9th World Cong Pergamon Press: Oxford. 1984.
- [21] H.Amano, T.Boku, T.Kudo an H.Aiso, "(SM)²-11", Proc. Int. Sympo. on Computer Architecture, pp. 100-107, 1985.