

マルチプロセッサシステム MUGENのソフトウェア構成

Software architecture for the Multiprocessor System MUGEN

高木 康志⁺ 堀口 進⁺ 川添 良幸⁺⁺ 重井 芳治⁺
Yasushi Takaki Susumu Horiguchi Yoshiyuki Kawazoe Yoshiharu Shigei

⁺東北大学工学部情報工学科
Department of Information Engineering, Tohoku University

⁺⁺東北大学情報処理教育センター
E C I P, Tohoku University

あらまし バス結合型のマルチプロセッサシステムでは、バスの競合によりシステムのスループットが低下する。このバス競合を緩和する方法のひとつとして、プロセッサを複数のクラスタに分けるクラスタ方式が提案されている。本稿では、このクラスタ方式を採用した試作マルチプロセッサシステムMUGEN上にインプリメントしたシステムソフトウェアの構成および並列計算機の性能に関して検討している。特に、並列処理を記述できるC言語para-CおよびそのプリプロセッサP3Cと、試作システムのオペレーティングシステムの構成について議論する。

Abstract A clustered multiprocessor system consisted of 32 processing elements has been designed and built. To implement and develop application programs on the multiprocessor system, the system software is significantly important. In this paper, we show the system software implemented on the clustered multiprocessor system (MUGEN). The system software contains a monitor, para-C and a pre-processor para-C (P3C) which translates a serial program into a parallel program written in para-C. The system performance is experimentally evaluated by using the Livermore loops.

1. はじめに

高速な計算機システムを実現する方法として、多数のプロセッサ上で処理を並列に実行する並列計算機の研究が盛んに行なわれている。この並列計算機の実現方法のひとつに、安価で高性能になったマイクロプロセッサを多数用いて、マルチマイクロプロセッサシステムを構成することが考えられ、また実際に様々な構成の試作機が製作されている[1][2][3][4]。

マルチプロセッサシステムにおけるプロセッサの結合方式としては種々のものが提案されているが、広く使用されているもののひとつに共有メモリとバスを用いた方式がある。そのう

ち、一本のバスに複数のプロセッサを接続する単一バス方式は、ハードウェア量が少なく構成が容易であるが、プロセッサ数が増加するとバス競合によるオーバーヘッドが増大し、処理効率が著しく低下する。これを緩和するために、プロセッサをいくつかのクラスタ単位にまとめるクラスタ方式が考案されている[1][2][4]。この方式は、バスを階層化するため異なるクラスタ内に属するプロセッサ間の通信はオーバーヘッドが大きくなるが、比較的大規模なシステムに対しても、処理効率の向上に比べてハードウェアの付加量が少なくすむという利点を持つ[4]。

このような高速処理を目的とした並列計算機では、ハードウェアが複雑になり、種々の資源を有効に利用する必要があるために、従来の計算機以上にシステムソフトウェアの重要性が増している。一般に、処理を高速に実行するためには、並列計算機の結合方法などのハードウェアの特徴に合わせてオペレーティングシステムを構成しなければならない。更に、プログラムを効率良く記述するための並列処理言語も重要である。

本稿では、クラスタ方式バス結合型マルチプロセッサ MUGEN^{[4][5][6]} 上にインプリメントしたシステムソフトウェアの特徴ならびに並列処理計算機の性能に関して検討している。次章に、試作クラスタ方式バス結合型マルチプロセッサの構成を簡単に示す。第3章では、既にインプリメントされている並列処理言語 *para-C*^{[7][8][9]} ならびに従来のC言語で記述されたプログラムを *para-C* に変換するプリプロセッサ (PPC)^[12] を中心に、プログラミング環境について述べる。また、リバモアループを用いた並列計算機の性能評価を行い、システム効率の評価・検討を行なう。第4章では、試作システムにおいて必要とされるオペレーティングシステムの構成について議論している。

2. クラスタ方式バス結合型

マルチプロセッサ MUGEN

バス結合によるマルチプロセッサは様々な問題に対して柔軟に適應することができ、汎用性の高い並列計算機を構築できる。バスの結合形態としては、単一バス、クロスバ、複線バスなどが代表的である。それぞれの特徴をまとめると、まず単一バスシステムは、バス調停を行なうロジックが簡潔であり、ネットワークに関するハードウェア量が少なく構成が容易である。また、クロスバは最大でメモリと同数までのプロセッサが同時にメモリアクセス可能である。複線バスはシステム内に複数のバスを有しているので、バスの本数までのメモリアクセスを同時に行なうことが可能である^[10]。

バス結合マルチプロセッサでは、プロセッサ数が増加するにつれてバス競合の生じる確率が高くなり、システムのスループットが低下する。相互結合網のコストとマルチプロセッサのスループットの比を結合の価格性能比と定義す

ると、単一バスは非常によい価格性能比を得ることができ、それに対してクロスバは価格性能比の点からは問題がある。しかし、単一バスは、プロセッサ数が大きくなるに従ってメモリアクセス競合によりスループット自体が飽和するため、大規模なシステム構築には向かない。これに対してクロスバはメモリアクセスが多くなっても性能が比較的低下せず、スループットが増加する。従って、単一バスとクロスバの中間的な特性を持ち、メモリアクセス競合が小さいシステムを構築することが望まれる。

メモリアクセス競合を緩和するアプローチとして、系内のプロセッサをクラスタ分けする方法がある。このアプローチは価格性能比の点から好ましいと考えられる。ある程度プロセッサ台数が多く競合が生じているとき、価格性能比は単一バスに比較して非常に改善される。しかし、クラスタ分けした場合には、他クラスタの共有メモリアクセスする場合、いったん1レベル上の階層にあるマスタプロセッサ (MP) を経由するため、クラスタ内の共有メモリへのアクセスよりも時間がかかる。また、クラスタ外へのアクセスが多いときにはMPがネックとなってシステムの性能が低下する。そのため、これを避けるよう注意してプログラムを作成し

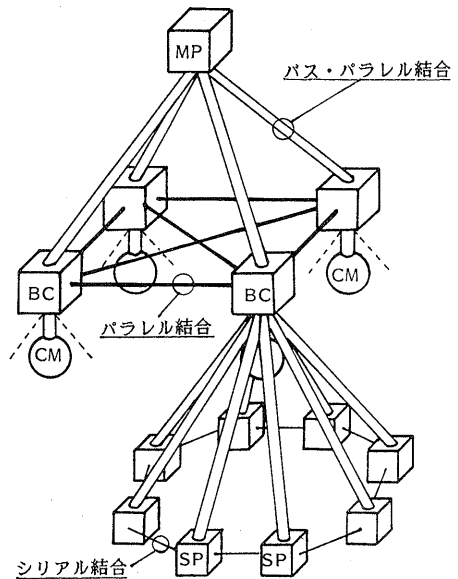


図1 システム構成概念図

MP (MASTER PROCESSOR)	Number CPU ROM RAM PIO SIO CTC	1 Unit Z80A-CPU 8 KB 56 KB Z80A-PIO×6 (12 Ports) Z80A-SIO (2 Ports) Z80A-CTC
BC (BUS CONTROLLER)	Number CPU ROM RAM PIO	4 Units Z80A-CPU 8 KB 8 KB Z80A-PIO×8 (16 Ports)
CM (CLUSTER COMMON MEMORY)	Number RAM	4 Units 32 KB
SP (SLAVE PROCESSOR)	Number CPU RAM PIO SIO CTC	32 Units Z80A-CPU 32 KB Z80A-PIO (2 Ports) Z80A-SIO (2 Ports) Z80A-CTC

表1 試作システムのハードウェア仕様

なければならない。

このクラスタ方式を用いたバス結合マルチマイクロプロセッサの試作システムとして、32台のスレーブプロセッサ (SP) を持つ試作システム MUGEN^[4] を取り上げる。システムの構成概念図を図1に示す。また、試作システムのハードウェア仕様を表1に示す。システム全体を管理するマスタプロセッサ (MP) の下に4つのクラスタがあり、それぞれのクラスタは、クラスタ内のバスを管理するバスコントローラ (BC)、プロセッサ間のデータ受け渡しに使用されるクラスタコモンメモリ (CM)、および並列処理を実行する8つのSPで構成されている。

各SPには連続した番号が付けられている。クラスタ0のSP0をスレーブナンバー0とし、以下、クラスタ3のSP7まで各々に0から31までの番号が割り当てられている。これによりMPはSP全体を容易に管理できる。

それぞれのBC間はパラレルポートで完全結合され、異なったクラスタ間でのデータ通信を補助している。また隣接するSP間はシリアルポートで接続され、これによってクラスタ内の通信の競合を緩和することができる。バスはアドレスバス16ビット、データバス8ビットから成る。各SPとBC間および各BCとMP間はパラレルポートを通して結ばれる。

MPは現在RS232Cのシリアル回線(9600bps)によりホストコンピュータと直接に接続され、システムの操作およびプログラムの開発などはホストコンピュータ上で行っている。

各プロセッサの機能と構成について述べる。まず、マスタプロセッサ (MP) は、プログラムを実行する際に、並列実行が可能なプログラムを32台のSPのローカルメモリ (LM) に転送して演算処理を行わせる。処理された結果は4つのCMから読み込む。CPUにはZ80A-CPUが用いられており、メモリはROMが8KB、RAMが56KBで構成されている。メモリの8000H番地からは9バンクに分けられ、MPのLMの他、各クラスタのCMおよびSPのLMを直接アクセスすることができる。

スレーブプロセッサ (SP) は、MPから送られたプログラムを実行してクラスタ内のCMに結果のデータを転送する。CPUにMPと同じZ80A-CPUを用いており、LMとして32KBのRAMを持っている。8000H番地以降は自クラスタのCMをアクセスする時に用いる。

バスコントローラ (BC) は、MPがSPのLMにプログラムをロードする際のSPの制御や、SPがCMをアクセスする際のバス競合の仲裁などのバスの管理を行う。さらに、MPが起動をかけたクラスタ内のSPが処理を終了し、必要なデータが揃ったことをMPに通報する処理を行なう。

3. MUGENのプログラム開発環境

MUGEN上での応用プログラム開発が効率よく行え、システムの機能を十分に発揮できるためには、高級言語が利用できることが必要不可欠である。MUGENでは、ホストコンピュータのCP/M-80上で動作するC言語のコンパイラを改良して、本システム上で実行可能なオブジェクトを生成できるコンパイラ (parser)^[8] を開発し、利用している。C言語は、浮動小数点演算などの数値処理が可能で、またアセンブラレベルの細かい記述もできる。アセンブラで記述したプログラムを新しい関数や手続きとしてライブラリに加えることも容易で、さらにライブラリのソースが公開されているため既存の関数を本システム上で実行可能にすることも容易に行える。

実際の応用プログラム開発および実行は、以下のように行なわれる。まず、ホストコンピ

ユーザ上でプログラムを作成する。そのプログラムを、ホスト上でコンパイルした後、MUGENシステム上に転送・実行する。

3.1. 並列処理言語 para-C

para-Cは、並列処理をプログラム上に陽に示す方法を取っており、そのためのキーワードとして cobegin, process, body および forall が追加されている。para-Cのキーワードを表2に、para-Cによるプログラム例を図2に示す。para-Cでは、並列処理の起動は、Cのメインプログラムにあたる関数mainによってのみ行なうことができる。cobegin 文と、それに続くprocess 文によって、並列処理が起動される。並列処理の単位はプロセスで、各プロセスはSP上で逐次的に処理される。並列処理される関数は各SPにプロセスとして静的に割り付けられ、process 文によって、対応するSPが起動される。body 文はその関数がプロセスの本体であることを示す。プロセスの実行結果として得られるデータは、各クラスタのCMを介して、MP上で実行されるメインプログラムに引き渡される。プロセス間の通信は、関数 send, receive を用いて1対1で行なわれる。

3.2. para-CのプリプロセッサP3C

para-Cでは、並列処理をプログラム上に陽に記述する方法を取る。この方法によ

```

main()
{
    int i;
    long n, result, answer, start[32], end[32];

    scanf("%ld", &n);
    for(i = 0; i < 32; ++i)
        start[i] = i*n/32 + 1;
        end[i] = (i+1)*n/32;
    cobegin
        ①データの転送、SPの起動
        forall(i = 0; i < 32; ++i)
            ②process(i; sum(start[i], end[i]));
    for(answer = 0, i = 0; i < 32; ++i) {
        ③hreceive(i, "%ld", &result);
        answer += result;
        ④データの受け取り
    }
}
body sum(lower, upper)
long lower, upper;
{
    int i;
    long total;

    for(total = 0, i = lower; i <= upper; ++i)
        total += i;
    ⑤hsend("%ld", total);
    ⑥結果の転送
}

```

図2 para-Cのプログラム例

cobegin	並列処理の起動を宣言
process	関数のSP割り付けおよびプロセスの起動
body	関数がSP上のプロセスであることを宣言
forall	プロセス起動の繰り返し
hsend	MPへデータを送信
hreceive	SPからのデータを受信

表2 para-Cのキーワード

て、コンパイラは問題の並列性を自動的に検出する必要がなくなり、また、並列処理の単位をプログラマが任意に決定できるので、実効的なスループットを上げることができる。その反面、並列処理を行なうためには、プログラマが並列処理を直接記述する必要があり、通常のC言語で記述されたプログラムをpara-Cによって並列に実行することはできない。この問題を解決するために、通常のC言語で記述されたプログラムをpara-C用のプログラムに変換するプリプロセッサ(P3C) [12]を作成した。これを用いることによって、通常のC言語で記述されたプログラムを、自動的にpara-C用プログラムに変換し、高速に並列実行することができる。

C言語で記述された応用プログラムは、図3の手順で処理される。まず、ホストコンピュー

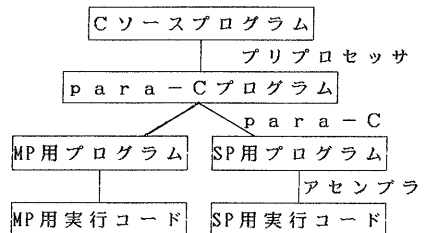


図3 プログラムの処理手順

タ上でプリプロセッサP3Cによってpara-C用のプログラムに変換され、次に、para-Cコンパイラによって、MUGENのプロセッサであるZ80のアセンブリ言語に変換される。この時、1つのソースプログラムから、MP用とSP用の2つのアセンブリプログラムが出力される。このそれぞれをアセンブラによって機械語プログラムに変換し、MUGEN上に転送して実行する。

para-Cの文法上の制限により、プリプロセッサで処理できるC言語のプログラムには、幾つかの制限が設けられる。構造体や共用体は使用できず、浮動小数点演算では、doubleの演算はfloat(32bits)に変換される。また、プロセスは逐次的に処理され、プロセスからプロセスを呼び出すことは出来ないので、再帰的な関数は、直接並列なプロセスへは変換されない。

3.3. リバモアループによる性能評価

C言語でプログラムを実行した際のシステムのスループットを測定するために、リバモアループによるベンチマークテストを行なった。

リバモアループ^[13]は、約20年前に科学計算コンピュータのFORTRAN処理能力を測定するためにリバモア研究所で作られたベンチマーク用プログラムである。プログラムは、一般的な大規模計算の中核コードを集めた14個のDOループで構成される。今日では計算機の処理能力を測定する一般的な検査プログラムとして、広く利用されている。各ループの内容は以下の通りである。

- (1) 流体
- (2) M.L.R.内積
- (3) 内積
- (4) 帯型連立一次方程式
- (5) 三角化消去(下三角形)
- (6) 三角化消去(上三角形)
- (7) 状態方程式
- (8) P.D.E.積分
- (9) 整数予測
- (10) 差分予測
- (11) 総和
- (12) 差分
- (13) 2次元粒子推進
- (14) 1次元粒子推進

このリバモアループをP3Cを用いて並列処理検出を行い、MUGENシステム上での処理時間を測定した。測定結果を表3に示す。表の左側の数値は、MP1台のみで実行した場合の処理時間、右側の数値は、SPを使用して測定した処理時間である。ここでの処理時間は大量のデータ転送時間も含んでいる。数値が欠けているループは、P3Cを用いてアルゴリズム上並列化が難しいなどの理由で並列処理の実行時間は測定できなかったものである。並列化による速度向上比は、2から17倍と大きなばらつきが見られる。これは、演算のアルゴリズムによって並列化できる度合いが違ふことと、SPに転送するデータ量が大きく違ふ事による。

次に、ループ1の場合の速度向上比の台数による変化を図4に示す。台数の増加に従って、速度比の増加の割合が小さくなっていくことが分る。

4. MUGEN用OSの構成

現在MUGENシステムは、各プロセッサのメモリ上に常駐するモニタによって制御されている^[11]。しかし、マルチプロセッサシステムにおいて実用的なMIMDのプログラムを動作させるためには、各プロセッサへのプロセスの割付、プロセス間の通信、それに伴うメモリの管理等の、より高度な機能をシステムで提供しなければならない。また、並列処理プログラムの開発ではデバッグが複雑になるため、プロ

No.	内容	演算時間 (sec.)	
		Z80A	MUGEN
1	流体	973	84
2	M.L.R.内積	944	136
3	内積	1150	140
4	帯型連立一次方程式	606	-
5	三角化消去(下)	1084	-
6	三角化消去(上)	1089	-
7	状態方程式	924	55
8	P.D.E.積分	796	371
9	整数予測	692	81
10	差分予測	348	189
11	総和	409	-
12	差分	535	84
13	2次元粒子推進	713	-
14	1次元粒子推進	994	319

表3 リバモアループの実行時間

プログラムの実行時に検査を行なう機能が重要となる。このため、システム上でのプログラムの開発・実行環境を整備する必要がある。試作システム MUGEN 用オペレーティングシステム (OS) は、このような要求に基づいて設計する必要がある。

MUGEN は、クラスタ化によってシステムが階層化されている。またプロセッサが互いに通信を行う際にはシリアルポートなどの複数のバスが存在し、また各バスごとにデータの授受に必要な時間が異なる。これらのハードウェアの特徴を十分に考慮して設計すべき MUGEN 用 OS に必要な条件をまとめる [14][15]。

- (1) システム全体の統一管理 (階層内管理、通信バス管理等)
- (2) 汎用的なプロセス管理
- (3) 整備された開発・実行環境
- (4) 簡便な入出力、ファイル管理

MUGEN においては、システムの大きさの点で、メモリを多量に消費するような大規模な構成は現実的でない。したがって、OS の基本的な機能に重点を置き、その中でシステムの特徴を生かした構成を考える。具体的には、クラスタ構成に適合したプロセスの管理や、複数のバスを有効に利用するプロセッサ間通信の実行などが挙げられる。これらの条件から検討した OS の構成を図 5 に示す。入出力・ファイル処理は、ターミナルプロセッサ (TP) を接続し、一括して管理する。

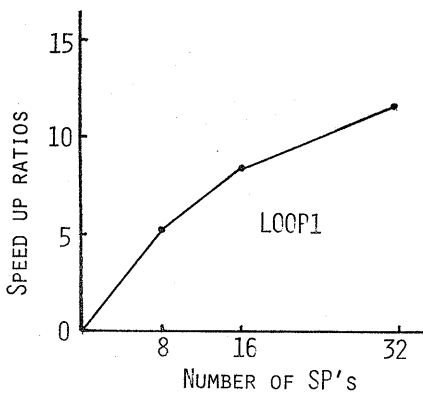


図 4 リバモアループの速度向上比

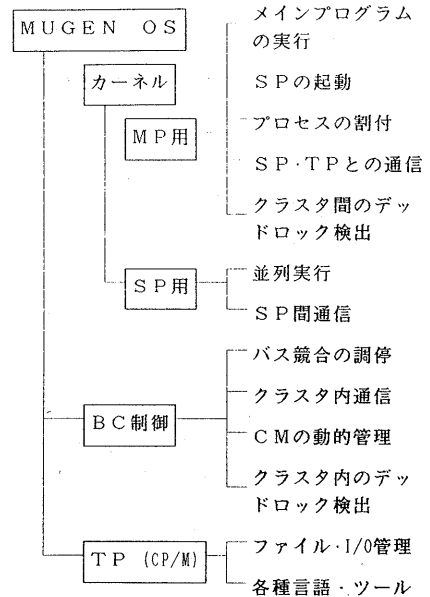


図 5 MUGEN 用 OS の構成

4. 1. BC によるクラスタ管理

システムのクラスタ構成を生かすため、クラスタ内の管理を BC によって統括的に行う。内容として、以下のようなことが挙げられる。

- 1) CM の動的管理
- 2) SP 間通信およびデッドロック管理
- 3) バス競合の調停
- 4) MP-SP 間の通信管理

まず CM の管理については、現在は起動前に静的に各 SP の使用する範囲を割当てている。これを OS によって一括して動的に管理する。各クラスタ 32 キロバイトの CM を、256 個のブロックに分割して管理する。MP および SP は CM を使用する場合、必要なブロック数を BC に知らせる。BC は必要な数の連続したブロックが確保できる場合は、プロセッサに先頭ブロックへのポインタを渡し、CM の使用を許可する。このように CM をブロックに分けて管理することにより、管理のテーブルの大きさを抑えることができる。さらに、アクセスのためのポインタも短くなり、その転送による通信時間を減らすことができる。

SP 間の通信方法としては、基本的に各クラスタの CM を用いる。しかし、バスの競合によ

る性能低下をできるだけ抑えるため、ポートで接続された各バスを有効に活用しなければならない。各バスのバンド幅が異なるため、通信するプロセッサやデータ量に応じて最適なバスをOSが自動的に選択する必要がある。他クラスタのSPと通信する場合は、MPが仲介を行なう。また、隣接するSPによってパイプラインを構成する場合、他クラスタ内SPも含めて16台以上をリング状に構成することもハードウェア的に可能であるが、この場合も、他クラスタ内SPと接続されるSPは例外処理としてMPが管理する。

4. 2. プロセス割付

MUGENのSPへのプロセス割付について考える。外部の事象によってプロセスが起動される場合に対応して、動的にプロセスをSPに割付ける機能が必要になる。MUGENのクラスタ型システムである特徴を生かすためには、相互関連の強いプロセスは同一のクラスタに割付ける必要がある。SPの起動はMPが直接行っており、したがってプロセスのSP割付・起動はMPが行なう。このため、動的に起動されたプロセスはMPが管理し、起動したSPに物理的に最も近いSPに割当てる。

4. 3. プロセッサ間通信

同一クラスタ内のSP間通信について考察する。この場合、CMによる通信のほかに、隣接するSP間のシリアル接続を利用することができる。しかし、表4に示すように、シリアルポートによる通信は転送するバイト数に比例して通信時間が増加する。さらに、隣接していないSPとの通信にシリアルポートを用いる場合は、SP間の間隔に比例した転送時間の増加に加えて、プロセッサの割り込み処理のオーバーヘッドも増加する。このため、このバスは数バイトの、極めて短い情報を転送する場合のみが実用的である。このシリアルポートによるバス

	1バイト	10バイト	使用バス
SP ↔ CM (競合なし)	135 μS	165 μS	バス
SP ↔ SP (SIO使用)	120	1200	ポート

表4 SP間の通信時間

は、通信の同期を取るためのメッセージ伝達のみを使う通信方法も考えられる。

5. まとめ

クラスタ方式を用いたバス結合マルチマイクロプロセッサの試作システムMUGENにおけるソフトウェアの構成について、システムソフトウェアを中心に検討した。今後は、OSの全体のインプリメントと、それに伴うシステム評価を行なう予定である。

《謝辞》日頃、御指導戴く東北大学奈良久教授ならびに中村維男助教授に深く感謝いたします。

参考文献

- [1] A.K.Jones and E.F.Gehringer: "The Cm^{*} Multiprocessor Project: Research Review", Dep.of Computer Science, Carnegie-Mellon Univ., CMU-CS-80-131 (1980).
- [2] 小原: "階層構造のMIMD型スーパーコンピュータ", 情報処理, vol.25, No.5, pp.480-490, 1984.
- [3] 白河他: "並列計算機PAX-128", 信学論(D), J67-D, No.8, pp.853-860 (昭59-8)
- [4] 中田他: "クラスタ方式を用いた共有メモリ型並列処理システムの試作", 情処第30回全大, 5B-3.
- [5] 中田他: "クラスタ方式を用いたメモリ結合型並列処理システムのバス制御方式の検討", 情処第29回全大, 5B-9.
- [6] 中田他: "クラスタ方式共有メモリ型並列処理システムの処理効率の測定", 昭和60年度信学総合全大, 1712.
- [7] 中津山他: "バス結合型マルチプロセッサシステムの高位言語", 昭60情処全大 3D-5.
- [8] 中津山他: "para-Cコンパイラと試作機の性能評価", 昭和61年度電気関係学会東北支部連合大会, 2J-19
- [9] 中津山他: "並列計算機MUGENにおける並列処理言語", 昭和61年度信学総合全大, 1481.
- [10] 中津山他: "並列計算機のソフトウェアと性能評価", 信学技報, CAS86-120 (昭61-11).
- [11] 高木他: "クラスタ式マルチプロセッサのシステムソフトウェア", 情処第33回全大, 4C-6.

[12] 高木他、"並列処理言語para-Cのプリプロセッサについて", 昭和62年度電気関係学会東北支部連合大会

[13] 唐木幸比古 : "Formula Translator の動向", Computer Today, vol7, No.2, pp35-42, 1984

[14] A.S.Tanenbaum and R.van Renesse : " Distributed Operating Systems ", ACM Computing Surveys, vol.17, No.4, pp.419-470, 1985

[15] G.R.Andrews and F.B.Schneider : " Concepts and Notations for Concurrent Programming", ACM Computing Surveys, vol. 15, No.1, pp.3-43, 1983