

## 並列 / 分散ソフトウェア開発支援環境の構築

野口善昭<sup>†</sup> 金丸弘樹<sup>†</sup>  
東由美<sup>†</sup> 岩根雅彦<sup>†</sup>

細粒度並列処理マルチプロセッサ MSBM と汎用マルチプロセッサ MDBM/FMM を制作して、細粒度静的コードスケジューリングアルゴリズムや動的スケジューラの開発などの細粒度並列処理の研究にスタンドアロンで使用してきた。研究環境を向上させるためだけでなく、並列処理と分散処理の効率のよい融合方法を研究するために、MSBM と MDBM/FMM を LAN に接続して、Linux 上でのクライアント / サーバモデルを採用した。マルチプロセッサで並列実行されるアプリケーションプログラムに対してコントラクタという概念を導入して抽象化することによってハードウェア構成と分離した。そして複数のコントラクタを一つのマルチプロセッサで並列実行可能とした。コントラクタは動作環境ヘッダ、コーディネータ、マネージャおよびワーカで構成される。コーディネータはサーバとアプリケーションプログラムとのインタフェース、マネージャはアプリケーションプログラムのなかの親プロセスであり、ワーカは子プロセスとしての計算主体である。マネージャおよびワーカは要素プロセッサをまたがって並列実行される。動作環境ヘッダはコントラクタの効率的な要素プロセッサへの割当てに使用する。特に実行ノードリストによりコントラクタのほかのノードコンピュータへの遠隔手続き呼び出しの自動転送を可能にした。簡単な例で基礎実験を行った。実験結果からこの方式の有効性が確認された。

## The Parallel and Distributed Applications Development Environment

YOSHIAKI NOGUCHI,<sup>†</sup> HIROKI KANAMARU,<sup>†</sup> YUMI HIGASHI<sup>†</sup>  
and MASAHICO IWANE<sup>†</sup>

The parallel computer MSBM for fine grain tasks and the general purpose parallel computer MDBM/FMM has been developed. We study fine grain size parallel processing such as development of fine grain size code scheduling algorithms, dynamic scheduling systems and so on. MSBM and MDBM/FMM has been used as a stand-alone parallel computer. We connect MSBM and MDBM/FMM to LAN, and we adopt client/server model on Linux operating system not only to improve an environment for research but also to study how to fuse parallel and distributed processing. For this purpose, we has been introduced the contractor that separates hardware configurations from an application program. It consists of an operation environment header, a coordinator, a manager, and workers. A coordinator is an interface between server and an application program. A manager is parent process and a worker is child process in an application program. The manager and the worker compute using some processors. An operation environment header is used for efficient assigning contractor to processor. Specially, server sends a remote procedure call to other node-computer automatically. We made a simple test. It shows that this method is effective.

### 1. はじめに

複数のプロセッサを密結合接続した並列計算機が開発されてきた。一方ではローカルエリアネットワーク LAN による分散処理が一般化している<sup>1)</sup>。密結合並列計算機をスタンドアロンで使用するかまたは LAN のノードコンピュータとして使用してほかのノードコン

ピュータから遠隔手続き呼び出しによってアプリケーションプログラムを実行している。また並列計算機をクラスタとして取り扱い高速ネットワークによってそれらを接続して超並列計算機を実現している<sup>2)</sup>。我々は細粒度マルチプロセッサ MSBM<sup>3)4)</sup> および汎用マルチプロセッサ MDBM/FMM<sup>5)</sup> を開発してきたが、これらは MSDOS のもとでスタンドアロンで細粒度静的コードスケジューリングアルゴリズムの開発<sup>6)</sup>、並列コンパイラの開発、細粒度動的スケジューラの開発<sup>7)</sup> に使用してきた。MSBM と MDBM/FMM を前者の立場で LAN に接続して、ソフトウェア開発支援環境を向上させると

<sup>†</sup>九州工業大学工学部電気工学科  
Department of Electrical Engineering, Faculty of Engineering, Kyushu Institute of Technology

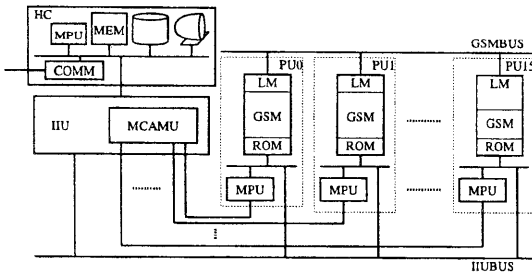


図1 MSBM ハードウェア構成

ともに、複数の並列計算機による細粒度並列処理と分散処理の効率的な融合処理の実現可能性を探る。本報告ではMSBMとMDBM/FMMのハードウェア構成を簡単に述べ、両マルチプロセッサに共通化したLinux上で構築されたスタンドアロン用ソフトウェア構成を述べる。そしてネットワーク構成と融合処理の実現方法を述べ、簡単な実験を行い、その結果について検討した。

## 2. 並列計算機のハードウェア構成

### 2.1 MSBM

MSBMは細粒度並列処理の多重プログラミング環境における静的バリア同期管理SBMとグループ共有メモリGSMの評価のために開発したマルチプロセッサである。一つのアプリケーションプログラムAPの実行に割り当てられる要素プロセッサPUをプロセッサグループ(単にグループ)、そのAPの中でバリア同期をとるPUの組をバリアグループとよぶ。MSBMではプログラムの実行開始前にバリア同期機構MCAMU<sup>8)</sup>に設定しプログラム終了後にその設定を解除するSBMをおこなった。

MSBMは図1に示すようにホストコンピュータHC、統合インタフェースIIUおよび16台のPUから構成される。HCは80386MPU、8MBメモリ、各種入出力装置からなるパーソナルコンピュータである。IIUは、DMAによるHCメモリとPUメモリ間のマルチキャストデータ転送機構、PUからHCへの割り込みの中継機構、SBMに適したMCAMUで構成される。PUは8088+8087MPU、256kBメモリ、グループ内PUメモリ間で共有メモリを実現するGSMインタフェース、PUからHCへ種々の処理を割り込みによって依頼するサービスコールSVCやHCとPU間のデータ転送に使用するIIUインタフェース、PUとMCAMU間のバリア信号で構成される。PUレジスタのうちグループ番号を保持するGPNO、局所メモリと共有メモリの境界を示すGSBASE、HCからPUへの小規模通信のためのRCTLPUはHCとPUの両方から設定できる。

### 2.2 MDBM/FMM

MDBM/FMMは細粒度並列処理の多重プログラミング環境における動的バリア同期管理DBMと細粒度、中

粒度、粗粒度並列処理におけるPU間通信の検証のために開発したマルチプロセッサである。DBMはバリアグループのMCAMUへの設定をAPの実行中にも可能としたものである。このためにPUからMCAMUへの設定時間は極力小さくしなければならない。汎用マルチプロセッサでは細粒度、中粒度、粗粒度並列処理を効率よくこなす必要があるので、共有メモリ構成から局所メモリ構成までPUメモリを段階的に選択できる再構成可能メモリRSMが望ましい。そこでPUメモリを局所メモリLM、グループ共有メモリGSM、システムで共有するシステム共有メモリSSMを任意の大きさに設定できるようにした。LMのみ必要で大量のプロセッサ間通信をおこなうAPに対してはPUメモリ間でDMA転送をおこなうILLIAC網を用意した。図2にMDBM/FMMの構成を示す。MDBM/FMMは80486MPUをもつ

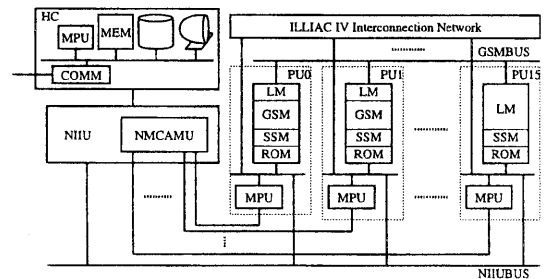


図2 MDBM/FMM ハードウェア構成

たHC、32-64台のPU(8088+8087)と新統合インタフェースNIUで構成される。MSBMと異なる箇所は高速バリア同期管理機構、RSM、NEWS転送である。NEWS転送はグループ内またはシステム内のすべてのPUのメモリ間でおこなわれる。隣接PUへ大量のデータ転送をおこなうMIMD型の粗粒度並列化APでは線形結合で十分であるのでILLIAC網はPU割り当ての容易さを確保するためである。

## 3. 並列計算機のソフトウェア構成

### 3.1 概要

MSBM, MDBM/FMMのソフトウェアはオペレーティングシステムLinux, PUなどの資源管理をおこなうエグゼクティブとアプリケーションプログラムであるコントラクタから構成され、コントラクタはコーディネータとマネージャおよびワーカから構成される。エグゼクティブはLinux上でつねに動作しており、キーボード入力からのコントラクタ実行の指示またはコントラクタの実行終了を待っている。PU待ちキューの中から、図3に示したコントラクタファイルにある動作環境ヘッダの必要PU台数とPU割り当て形態(離散的PU割り当てまたは列/行方向に連続したPU割り当て)をもとに、最小一致アルゴリズムによって、割り当て可能なコ

ントラクタにPUを割り当てる。同時にコーディネータをHCにローディングしてLinuxの一つのプロセスとして実行する。

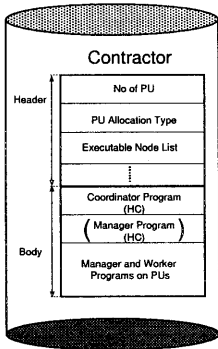


図3 コントラクタファイル

コーディネータはワーカをPUまたはHCにローディングして実行を指示するがこの手続きはコントラクタによって異なる。コーディネータがEXITを実行すればコントラクタは終了する。図4に示すように一つの並列計算機上で複数のコントラクタの実行が可能である。手間はかかるけれど各ユーザが都合のよいようにコントラクタを構成できる。これは細粒度並列処理用オペレーティングシステムなどの開発のために自由度を持たせたことによる。

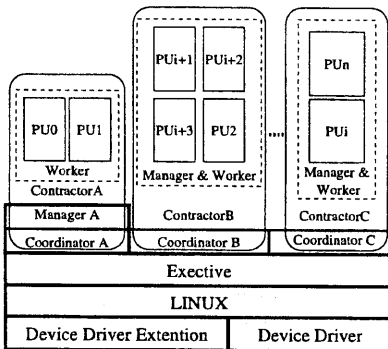


図4 ソフトウェア概要

### 3.2 デバイスドライバ

デバイスドライバはHCのLinux上に組み込まれ、IIUまたはNIU(以後IIU)を介してHCとPU間のプログラムやデータの転送、PU間同志の大量のデータ転送、HCからのバリア同期管理、HCからPUの制御、PUからHCへのSVC受付をおこなう。デバイスファイルによりPU指定、グループ指定、ブロードキャスト、MCAMU指定、SVCのPU指定をおこなう。これによりHCとPU間の情報転送

には通常のREAD, WRITE, LSEEKにより実行できる。またIIU内のDMACレジスタなどの設定やPU内のRCTLPUレジスタ、GSMBASEレジスタなどの設定にはIOCTLを用いる。IOCTLのIIU用特殊関数を表1に示す。またNEWS転送のためにIIU\_NEWS(from\_address,to\_address,size,direction), SVC受付のためにIIU\_SVC(PU\_no,count,data,current\_time)のシステムコールを用意した。

表1 IIU用IOCTLコマンド

Command	Operation
RESET	Send reset signal to IIU and All PUs
SET_GSBASE	Set GSM Base register
SET_GPNO	Set Group No register
SET_RCTLPU	Se Control PU register
NEWS_Transfer	Set DMAC and Start NEWS Transfer
SET_SVCGROUP	Set SVC Group No

### 3.3 モニタ

PUモニタはPUメモリの8KBのROM上におかれている。アイドル時にはRCTLPUを検査してその値により異なる適切な処理をおこなう。HCのIOCTLによってあるいはPU自身によってRCTLPUが書き換えられたとき、実行中の一連の処理の終了後にPUメールボックスをみてその処理をおこなう。表2に処理内容、図5にメールボックスのメモリマップを示す。

表2 処理内容

RCTLPU	Operation
FF	No Operation
FE	Start PU Program
FC	RESERVED
?	?
00	RESERVED

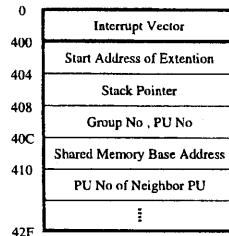


図5 PUメールボックス

PU上のマネージャおよびワーカの実行開始はHCからのRCTLPUの設定により、メールボックスに示された実行開始から実行させる。HCへのSVC要求や

MCAMの更新などはマネージャまたはワーカが直接出力命令を発行することによってモニタ処理は最小限にしている。動的スケジューリングを含んだオペレーティングシステムをPUに搭載したときは、マネージャおよびワーカの実行開始はこの方法とは当然異ったものとなる。

#### 4. ネットワーク構成

##### 4.1 分割利用と協調動作

図6に示すようにMSBMとMDBM/FMMは一つのノードコンピュータとしてWSやPCと同等に研究室LANで接続されている。またルータであるWSより学外にも接続されている。WSまたはPC上のクラ

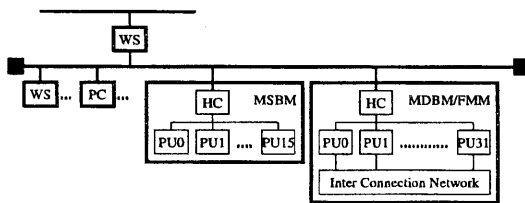


図6 ネットワーク構成

アントプロセスからMSBMまたはMDBM/FMMにあるコントラクタを実行するために、MSBMあるいはMDBM/FMMサーバプロセスにエグゼクティブを組み込んだことにより、サーバプロセスとクライアントプロセスが同一ノードコンピュータにあっても遠隔手続き呼び出しと局所手続き呼び出しは同じ取り扱いができる。サーバプロセスはクライアントプロセスから送付されたメッセージすなわちコントラクタファイル名を受け取り、ヘッダにある動作環境を調べて十分な空きPUが存在してすぐに実行可能であればこのノードコンピュータで実行する。実行が待たされるときは、コントラクタが実行できるノードコンピュータを示すコントラクタ実行ノードリストを調べて次のノードコンピュータにクライアントプロセスからのメッセージを転送する。このメッセージが実行ノードリストを一巡して最初のサーバプロセスに到着したときこのコントラクタをPU待ちキューに登録する。このようにユーザがノードコンピュータを直接指定しなくても空きノードコンピュータでコントラクタを実行できるようにした。またMSBM、MDBM/FMMでの複数のコントラクタの協調動作は通常のLinux上の協調プロセスとして実現した。図7に並列/分散ソフトウェア構成を示す。

##### 4.2 クライアントとサーバ

クライアントプロセスからMSBM(MDBM/FMM)サーバプロセスにコントラクタの実行を依頼するとき、well\_knownなメッセージプールWKNMPとtemporaryなメッセージプールTMPMPを使用してプロセ

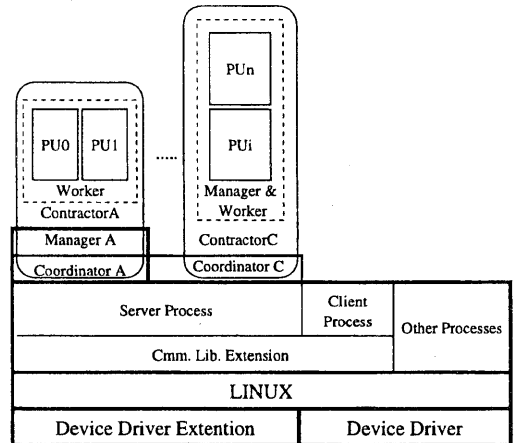


図7 並列/分散ソフトウェア構成

ス通信をおこなう。コントラクタ名とコントラクタがクライアントとデータの送受をおこなうTMPMP名をWKNMPに書き込んでサーバプロセスに送信する。サーバプロセスはコントラクタのコーディネータプロセスを生成してTMPMP名を渡す。コーディネータプロセスはTMPMPからデータを受け取るとともに、動作環境ヘッダをみてマネージャやワーカをHCやPU上に生成してコントラクタを実行する。コーディネータプロセスは実行結果をTMPMPに書き込みクライアントプロセスに送信する。クライアントプロセスでは実行結果を受け取って処理を続行する。図8にこの過程を示す。

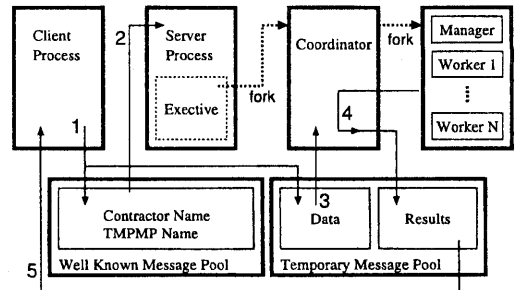


図8 遠隔手続き呼び出し

これらのメッセージの送受信のために、送信用として非ブロック型手続きであるsend(pool\_name,queue\_name,message,size), 受信用としてブロック型手続きであるrecv(pool\_name,queue\_name,message,size)をライブラリとして用意する。なおメッセージプールには複数のキューを持つことができる。このようにTMPMPを導入したのはMSBMとMDBM/FMM上で同一のコントラクタの実行可能を確保することとオーバーヘッドの削減を図るためである。

## 5. 実験

### 5.1 実験方法

コントラクトの自動転送による複数ノードコンピュータの分割利用と協調動作の実験を行う。クライアントプロセスからノードコンピュータへ送付された複数の並列実行可能なコントラクトを同時に処理し、PU待ち状態のコントラクトをほかのノードコンピュータへ送付することにより分割利用と協調動作の確認とその実行時間を測定する。

ノードコンピュータとしてMSBMとMDBM/FMMを用い、MSBMでは8台、MDBM/FMMでは16台のPUを使用する。クライアントプロセスを実行するマシンはSun SPARCstation 10(SS10)を使用する。

各コントラクトは並列実行可能な8個のワークを持ち8PUを使用する。PU割り当て形態は離散的PU割り当てを用いる。コーディネータはPUへのワークプログラム及びデータのローディング、MCAMへのバリアグループの登録、SVCグループの設定、グループの設定、実行開始、終了確認を行う。ワークは6553600回ループし1回のループごとにバリア同期をとるプログラムを実行する。ワークはSCVを用いて処理の終了をコーディネータへ連絡する。ワークプロセスの動作はマネージャを必要としないためマネージャは省略した。実行ノードリストについては、MDBM/FMMではMSBM、MSBMではMDBM/FMMを持つ。

クライアントプロセスは10個のコントラクトの実行要求メッセージをMDBM/FMMサーバに送り、すべてのコントラクトの実行終了を待つ。コントラクトの終了はクライアントプロセスが用意したメッセージプールTMPMPのキューにコーディネータからの終了メッセージが到着したことにより確認する。

PU上で動作するワークプログラムはアセンブリ言語で記述し、HC上で動作する通信ライブラリ、サーバ、コーディネータ及びSS10上で動作するクライアントプログラムはC++言語で記述した。ネットワークプロトコルはTCP/IPを用い、BSD socket インタフェースを用いて実装した。

動作の確認はサーバの動作記録を基に行った。実行時間の測定はサーバの動作記録のタイムスタンプとワークの実行時間の2つを測定した。ワークの実行時間はコーディネータがワークに実行開始要求を送ってから、実行終了を通知するSVCを受けとるまでの時間で測定した。

### 5.2 実験結果

サーバの動作記録をもとに得られた動作結果を図9に示す。 $C_0 \sim C_9$ はクライアントプロセスから駆動されたコントラクトを表す。また $C_n$ の添字はクライアントプロセスから要求を発行した順番を示す。 $G_0, G_1, G_2$ はグループを表わしそれぞれ8PUで構成される。

矢印はサーバによって自動的に送信されたコントラクトを表す。空欄の部分はアイドル状態であることを示し、斜線部分はコーディネータがワークプログラムの実行環境の準備を行っているオーバヘッドを示す。

10個のコントラクトのうち8個のコントラクトがMDBM/FMMを2つのグループに自動分割し同時に実行され、2個のコントラクトがMSBMに自動送付され処理が行われた。10個のコントラクトを自動分割、自動送付することなしに逐次実行する場合に比べ約2.5倍の速度向上が得られた。また、コーディネータによって得られた $C_0, C_3, C_6, C_9$ の実行時間と、サーバの実行記録のタイムスタンプによって得られた $T_{G_0}$ から、図9の斜線で示されるオーバヘッドを計算したところ、 $T_{G_0}$ の4.4%がオーバヘッドであることが得られた。

$C_5$ は5番目に駆動されたコントラクトであるがこの結果では最後に実行が行われている。この原因について検討する。 $C_0, C_1$ はそれぞれ $G_0, G_1$ で実行が開始される。残りの $C_2 \sim C_9$ は実行が待たされるので実行ノードリストよりMSBMのサーバに送られる。 $C_2$ はMSBMで実行されるが $C_3 \sim C_9$ は再びMDBM/FMMのサーバに送付される。この時点で $C_3 \sim C_9$ はMDBM/FMMサーバのPU待ちキューに登録される。 $C_0, C_1$ はほぼ同時刻に実行を終了しMDBM/FMMサーバは新たに $C_3, C_4$ の実行を開始する。 $C_5$ は資源不足のため実行が待たされるので、MSBMサーバに送られる。MSBMは $C_2$ を実行中であるため、 $C_5$ は再びMDBM/FMMサーバへ戻る。このときPU待ちキューの最後に登録されるため $C_5$ の実行は最後に行われる。

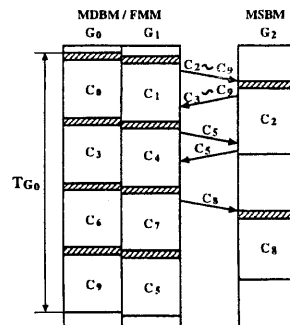


図9 コントラクトの動作履歴

## 6. むすび

複数の細粒度並列処理可能なマルチプロセッサをLANのノードコンピュータとして接続して細粒度並列処理と分散処理の効率的な融合を実現するソフトウェアシステムを構築した。マルチプロセッサ上で並列処理されるAPに対してコントラクトなる概念を導入して、コン

トラクタの構成要素を動作環境を示すヘッダとコーディネータ、マネージャ、ワーカと抽象化することによってハードウェア構成と分離することが可能となった。コントラクタと実行ノードリストによってマルチプロセッサの分割利用と有効利用が実現できた。

細粒度並列処理ノードコンピュータや細粒度並列処理マルチプロセッサオンチップのソフトウェア構築の基盤を与えたが、今後はこのシステムの各種プログラムによる評価およびこのシステム上で細粒度並列処理可能なマルチプロセッサのオペレーティングシステムの開発、細粒度並列コンパイラを開発する予定である。

### 参 考 文 献

- 1) 藤田, 篠田, 今泉(訳): 平行分散システム, トップラン(1996)
- 2) 富田真治: 超並列処理計算機 JUMP-1 のアーキテクチャ, 情報処理, Vol.36, No.6, pp528-537(1995)
- 3) 岩根, 本石, 野口, 米沢: 細粒度マルチプロセッサ MSBM, 情報処理論文誌, Vol.37, No.6, pp1196-1205(1996)
- 4) 岩根, 野口, 茶屋道, 本石, 重松: 細粒度並列計算機 MSBM の開発, 九州工大研究報告(工学), Vol67, pp.61-68(1995)
- 5) 岩根, 茶屋道, 本石, 立川, 浦崎: マルチマイクロプロセッサ MDBM/FMM の開発, 九州工大研究報告(工学), Vol.68, pp.49-56(1996)
- 6) 岩根, 濱田, 小島, 松田: 式の分割による並列化アルゴリズム ESH, SWOPP'96(秋田)(1996)
- 7) 松田, 武石, 野口, 岩根: 細粒度並列処理のためのハイブリッドスケジューリングシステム: スケジューラによるバリア同期管理手法, 情報処理研資 OS-70, pp.113-120, SWOPP'95(別府), (1995)
- 8) 岩根, 重松, 定家, 茶屋道, 金沢: FPGA によるバリア同期用機能メモリの開発, 九州工大研究報告(工学), Vol.66, pp.45-52(1994)