

汎用マルチマイクロプロセッサ FMM の構成と性能評価

野口善昭, 本石彰, 茶屋道宏貴, 岩根雅彦
九州工業大学工学部

少量のハードウェアで同期と通信のオーバーヘッドを削減するために、包含検索機能をもった拡張 CAM による半動的ハードウェアバリア同期機構と書き込み選択放送機能をもった拡張グループ共有メモリ構成を提案し、これらの機構の評価のための実験機として FMM を設計した。バリア同期機構からバリア識別情報を削除することによって生じるバリアの不一致を、拡張 CAM へのバリアパターンの半動的な登録削除方法を確立することによって解決した。並列化された Whetstone ベンチマークを使用して MSBM 計算機でシミュレーションによる評価をしたところ、6PU 以上で並列実行したときは逐次実行と比べて計算時間が 0.46 倍となった。

Architecture and Performance Evaluation of FMM: A General Purpose Multi-Microprocessors

Yoshiaki Noguchi, Akira Motoishi, Hirotaka Chayamichi, Masahiko Iwane
Faculty of Engineering, Kyushu Institute of Technology

In general purpose multi-microprocessor systems, decreasing the overhead of interprocessor communication and synchronization for running fine to coarse grain size parallel programs. FMM with the semi-dynamic barrier mechanism(SDB) and the extended group shared memory system(EGSM) to reduce the overhead is shown. The barrier pattern in the modified CAM can be updated dynamically at the beginning of running the procedure using the SDB. The EGSM consisting of the system shared memory, the group shared memory and the local memory has the ability of multicasting. The performance of FMM is evaluated using the paralleled Whetstone benchmark on MSBM.

1 はじめに

汎用マルチプロセッサでは細粒度から粗粒度までの並列性をもった多種多様な複数のプログラムが効率的に並列処理される必要がある。並列処理の効率を妨げる大きな要因はプロセッサ間の同期と通信によるオーバーヘッド¹⁾である。プロセッサ間同期の方法の一つにバリア同期があり、その機構はソフトウェアまたはハードウェアで実現されている。細粒度の並列性をもったプログラムではソフトウェアのバリア同期はバリア共有変数の操作による遅れが大きいのでハードウェアのバリア同期機構^{2, 3)}が提案されている。フェジーバリア⁴⁾はバリア同期識別情報とマスクによって、エラストティックバリア⁵⁾は同期カウンタとマスクによって任意のプロセッサ間のバリア同期を可能にしているが、配線量、回路量ともに増大する。SBM⁶⁾はバリア同期情報の代わりにバリア同期間で全順序関係を導入しバリアキューによって単一のバリア同期列の任意のプロセッサ間の同期をとる。DBM⁷⁾はバリア同期間に半順序関係を導入してバリアキューとCAM (Content Addressable Memory)によって任意のプロセッサ間の複数のバリア同期を同時に可能としている。SBMおよびDBMでは解の収束の程度によってループの繰り返し回数が異なるときには同期の管理が困難である。

密結合マルチプロセッサにおけるプロセッサ間通信は共有メモリによって行われるのでプロセッサと共有メモリを接続する相互接続網や共有メモリモジュールにおいて競合がおこる^{8, 9)}。集中型共有メモリにおける相互接続網としてハードウェア量の少ないバス結合が多く使用されている。プロセッサが数十台以上になるとキャッシュメモリによりバス競合を削減させているが、ハードウェア量が増加するだけでなくキャッシュコヒーレンス問題が生じる¹⁰⁾。各プロセッサが局所メモリと共有メモリをもつ部分共有マルチリードメモリ

では共有メモリの読み出し時の競合を減少できるが共有メモリをプロセッサに分散させるためにアドレス変換機構が必要である^{11, 12)}。疎結合マルチプロセッサにおける通信は入出力操作と同様に行われるので直接通信路のないプロセッサ間ではオーバーヘッドが生じるが、隣接したプロセッサへ一度に多量のデータ通信を行い、粗粒度の並列性をもつある種のプログラムに対しては格子結合網が有効である¹³⁾。MSBM¹⁴⁾がハードウェアバリア同期機構とグループ共有メモリ構成をもった細粒度並列計算機として開発されているが、効率よく静的コードスケジューリングを行うには制約が大きい¹⁵⁾。

細粒度の並列性をもった多種多様な複数のプログラムを効率よく処理するために、拡張CAMを用いたハードウェアによる半動的バリア同期機構、プログラムごとに独立した共有メモリ空間とマルチキャスト機能による部分分散共有メモリ構成、格子結合網による隣接プロセッサ間通信機能をもったマルチマイクロプロセッサを設計した。

本論文では拡張CAMを用いたハードウェアによる半動的バリア同期機構、FMMのハードウェア構成および並列化 Whetstone ベンチマークによる性能予測について述べる。

2 半動的バリア同期機構

2.1 前提

共有データのための部分共有メモリとバリア同期機構をもった多数のマイクロプロセッサPUとホストコンピュータHCで構成されたマルチマイクロプロセッサシステムにおいて、マスタスレーブ型オペレーティングシステムのもとで複数の応用プログラムが多重実行される環境で次の前提を考える。

並列処理される一つの応用プログラムはその各々がプロセデュアと呼ばれる一つの主プログラムと幾つかのサブプログラムにより構成される。さらにプロセデュアはマイクロプロセデュアと呼

ばれる逐次あるいは並列に実行される命令列で構成される。ユーザがオペレーティングシステムのコマンドにより実行することができる一つの応用プログラムをジョブと呼び、プロセデュア、マイクロプロセデュアに対するそれぞれの処理の実体をプロセス、マイクロプロセスと呼ぶ。ジョブを構成するプロセデュアは個別にコンパイラなどによって静的スケジューリングがなされて、バリア同期が設定され、並列化されている。リンカーおよびローダによりこれらのプロセデュアはジョブとして結合されジョブ単位でスケジューリングを行い、ジョブを固定のマイクロプロセデュアの集合と考えてジョブ実行前にマイクロプロセデュアを要素プロセッサへ割り当てる静的割り当て¹⁶⁾を行う。このとき並列に動作可能なマイクロプロセデュアは異なった要素プロセッサに割り当てる。図1において丸印はマイクロプロセデュア、点線はプロセデュア、一点破線はジョブである。そして実行中のジョブが終了したときPUを他のジョブに解放するものとする。なお一つのジョブに割り当てられた要素プロセッサの集合をグループと呼ぶ。

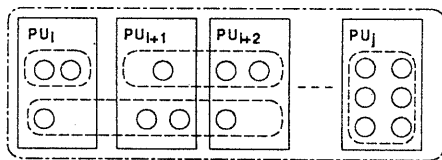


図1: プロセデュアのPU割り当て

2.2 拡張CAM¹⁷⁾による半動的同期機構

通常のバリア同期のほか、バリア同期識別情報の削除によって生じる不都合を補うダミーバリア (Dummy Barrier) 同期、および広がりのある領域内での同期を許すファジーバリア (Fuzzy Barrier) 同期を導入する。これらの同期を実現するために DmBar (Dummy Barrier) 命令、AdBar (Advanced Barrier) 命令、Bar (Bar-

rier) 命令を PU 出力命令として用意する。PU は DmBar 命令を実行するとバリア同期ユニット MCAMU にバリア同期要求を送り処理を続行する。PU は AdBar 命令を実行すると MCAMU にバリア同期要求を送り処理を続行するが、このとき同期が成立しておれば MCAMU は PU に同期成立を知らせる。PU が Bar 命令を実行したとき、すでに AdBar 命令が実行されて同期成立が知らされていれば PU は処理を続行する。まだ同期が成立していなければ同期成立まで処理を中断する。PU が Bar 命令に先だつて AdBar 命令を実行していなければバリア同期要求を MCAMU に送り処理を中断する。そして同期が成立すれば MCAMU は PU に同期成立を知らせて PU の処理を再開させる。すなわち Bar 命令単独で通常のバリア同期、DmBar 命令でダミーバリア同期、AdBar 命令と Bar 命令のペアでファジーバリア同期を実現する。

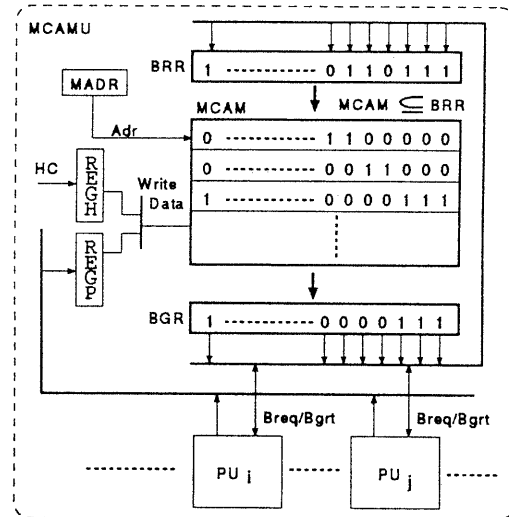


図2: MCAMUの構造

図2に示すようにバリア同期機構はMCAMUと個々のPU内のBarif (Barrier interface)で構成される。またMCAMUはシステム内PU台数に等しいビット長をもつたBRR (Barrier Re-

quest Register),BGR (barrier Grant Register) および MCAM(Modified Content Addressable Memory)で構成される。特に MCAMU には、エントリの有効、無効を示す v フラグが1ビット付加してある。PU が DmBar 命令, AdBar 命令, Bar 命令を実行したとき Barif は最初に Dm フラグを検査して、セットしていれば MCAMU からの Bgrt (Barrier Grant) 信号を受け取るまで Wait 信号により PU の実行を中断させる。Dm フラグがセットしていないときまたは Bgrt 信号の到着後、DmBar 命令では Dm フラグをセットして Breq (Barrier Request) 信号を MCAMU に送り PU の処理を続行する。AdBar 命令では Ad フラグをセットして Breq 信号を MCAMU に送り PU の処理を続行する。Bar 命令では Ad フラグと Gr フラグの両方がリセットしていれば Breq 信号を MCAMU に送ると共に Wait 信号により Bgrt 信号が到着するまで PU の実行を中断する。単に Gr フラグのみがリセットしていれば Bgrt 信号が到着するまで PU の実行を中断する。Gr フラグがセットしていれば Gr フラグをリセットし PU の処理を続行する。Barif から送付された Breq 信号は BRR の当該ビットをセットする。MCAMU は BRR を MCAM の探索データとして使用して包含関係が成立すれば MCAM を BGR に読み出す。すなわち、

$$\text{if } MCAM_{ij} \subseteq BRR_j \ (j = 0, 1, \dots, n-1)$$

$$\text{then } BGR_j := MCAM_{ij} \ (j = 0, 1, \dots, n-1)$$

但し、 i は MCAM 内の i 番目のエントリ、 j はエントリ i 内の j 番目のビット、 n はシステムの PU 総数を表す。また包含関係は次式で与えられる。

$$MCAM_{ij} \subseteq BRR_j = \neg MCAM_{ij} \vee BRR_j$$

BGR に読み出されたデータは Bgrt 信号として Barif に同期の成立を伝える。同時に BGR と BRR の排他的論理和 EXOR をとって BGR に格納する。Bgrt 信号を受け取った Barif は Dm フラグをリセットし、さらに Ad フラグがセットしていれば Gr フラグのセットと Ad フラグのリセット

を行ったのちに、PU が処理を中断していれば再開する。なお、検索データに包含されるすべてのパターンの論理和が出力される。MCAM への書き込みは HC からも PU からも可能である。

2.3 バリアパターンの管理

グループ内 PU 間でのバリア同期は Bar 命令、ダミーバリア同期は DmBar 命令およびファジーバリア同期は AdBar 命令と Bar 命令のペアで実現できるが、どのマイクロプロセスやサブプロセスと同期をとるかあらかじめ MCAM に設定する必要がある。このようなバリア同期をとる PU の組を示すビット列をバリアパターンと呼ぶ。マイクロプロセス間またはジョブ間のバリアパターンの MCAM への登録はジョブの実行開始時に HC が行うかあるいは PU が `init_bar` プリミティブを実行することによって行われる。MCAM からの削除はジョブの終了時に HC が行うか PU が `reset_bar` プリミティブを実行することにより行われる。

`init_bar(arg1,arg2,...,entry)` : $arg1, arg2, \dots$ で与えられたバリアパターンを MCAM 上に登録し、そのときのエントリを `entry` に返す。そして v フラグをセットする。

`reset_bar(entry)` : `entry` で示した MCAM の v フラグをリセットする。

図3においてノードはマイクロプロセス、矢印は実行順序、 \times は Bar 命令、 Δ は DmBar 命令を表す。マイクロプロセス 7(3) で生産されたデータをマイクロプロセス 11(12) で消費するときマイクロプロセス 2(8) の直後の DmBar 命令は本来不要であるがバリアパターンの MCAM への登録削除回数を削減してオーバーヘッドを少なくするために挿入した。マイクロプロセス 9 の直後の DmBar 命令も本来不要である。PU0 で `init_bar` が実行されるとバリアパターンが MCAM に登録される。バリア同期 B1, B2, B3 では

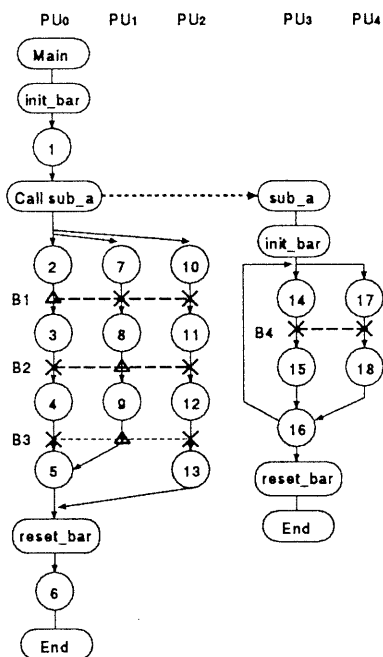


図 3: ジョブの実行順序

PU0, PU1, PU2 がバリア命令またはダミーバリア命令を実行するのでバリア同期が成立する。このバリアパターンはPU0が reset_bar を実行したとき無効となる。sub_a はループ中にバリアを必要とするサブルーチンとする。PU3, PU4 に対するバリアパターンはPU3の init_bar によってMCAMに登録され reset_bar によって抹消される。このように同じジョブでもバリア同期の設定は任意に独立に行うことができるので繰り返し回数が不定のループ中のバリア同期も問題なく処理できる。また、プロセッサ同志の静的スケジューリングはリンカーによって行われているので、通常メインルーチンはサブルーチンの終了を待たなくてよい。終了と待合わせる必要のあるときは粒度が大きいため共有変数による待合わせで十分である。

3 ハードウェア構成

3.1 システム構成

FMM は図4で示すようにHC (Host Computer), IIU (Integrated Interface Unit), および最大64台のPU (Processing Unit) から構成されている。

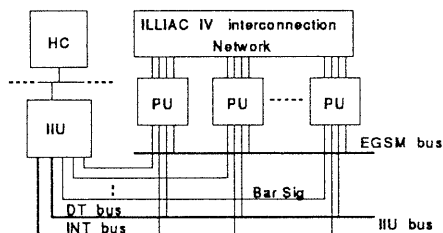


図 4: FMM ハードウェア構成

HC はパーソナルコンピュータ (MPU:i80386) を使用しユーザーインタフェースとプログラミング環境を提供するとともに、IIU, PU の制御を行う機能をもつ。

IIU はDTU (Data Transfer Unit), MCAMU (Modified CAM Unit), INTU (Interrupt Unit) からなり、HC-PU間のインタフェースを構成する。DTUはDMAによるHC-PU間及び隣接PU間のデータ転送機能を提供する。HC-PU間の転送はHCバスとDTバスを使用しPU単位、グループ単位またはすべてのPUに転送が可能である。隣接PU間はILLIAC網を使用し各PUは4方向の隣接PUとのデータ転送が可能である。MCAMUはバリア同期機構を使用するためのバリアプロセッサであり、任意のPUの集合で同期を取ることができる。MCAMUへのバリアパターンの登録、削除はHCからはHCバスを、PUからはINTバスを介して行われる。INTUはINTバスを介したPUからのサービス要求をその種類に応じて処理する。図5にパラメータデータのフォーマットを示す。31ビット目が0のときは割り込み要求データで1のときはバリア

パターン更新要求データである。Mod(Modifier)はIIUでの処理種別を指定する。Modがバリアパターン更新要求であれば、MCAMのAdrとPosで指定した位置にDataのバリアパターンを書き込む。

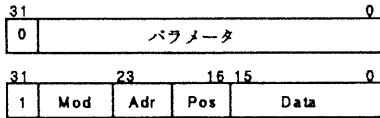


図 5: INT bus データ交換フォーマット

PUは8088MPU+8087FPU, 256KBメモリ, IIUインタフェース, EGSMインタフェース, NEWSインタフェースから構成される。そして、DMA転送やデータのグループ共有のためにPU番号とグループ番号の識別機能をPU内部にもつ。PUメモリは拡張グループ共有メモリEGSM構成をとり、LM(Local Memory), GSM(Group Shared Memory), SSM(System Shared Memory), からなる。GSMにはグループで共有するデータを格納することができる。

3.2 拡張グループ共有メモリEGSM

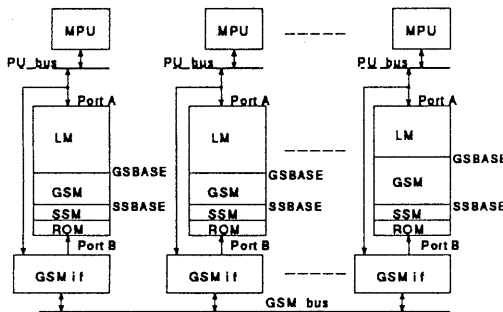


図 6: GSM 構成

拡張グループ共有メモリEGSM (Extended Group Shared Memory) 構成を図6に示す。EGSM構成は各PUのメモリの内容をグループごとに完全一致させるグループ共有メモリGSMおよび、システム内のすべてのPUでメモリの内

容を一致させるシステム共有メモリSSMを実現する。特にSSMにはバリア管理のようなシステムに共通な情報を置く。共有データの参照時は常に自PUのメモリから読み出せばよいため、不要なバスアクセスを排除でき、バス競合による性能低下を軽減できる。また、各PUでは共有メモリ領域をLMと同様にマッピングするため、通常のメモリ命令でアクセスできる。GSM/SSM内容のグループ内一致は、あるプロセッサがGSM/SSMに書き込みを行ったときGSMバスに所属グループの番号とSSMアクセスを表す信号を出力させ、これを識別情報とするマルチキャストを同一アドレスへ行うことで実現される。GSM/SSM領域への書き込みかLM領域の書き込みかの判断は各PU内のレジスタGSBASE/SSBASEとアドレスを比較することで行う。メモリデバイスとしてDPM(Dual Port Memory)を用い、GSBASEレジスタによりLMとGSM, SSBASEレジスタによりGSMとSSMに分ける。DPMのPortAをプロセッサ内部に、PortBをGSMバスに接続することで、GSM内容の一致のためのGSMバスを介したメモリ書き込みが、プロセッサ内部のメモリアccessを阻害することがなくなり性能低下を軽減できる。

GSM/SSMの使用の際に必要な相互排除もGSMバスの機能として提供されており、GSM転送用のグループロックとSSM転送用のシステムロックの機能がある。

4 MSBMによる性能予測

4.1 基本性能と並列化ベンチマーク

FMMの基本性能として、PUからのバリアパターンの更新、データ転送およびバリア同期の処理能力を設計データから算出した。バリアパターンの更新に10μsを要した。HCとPU間のデータ転送能力およびILLIAC網によるPUとPU間のデータ転送能力は500KB/s, EGSMバスによるPUとPU間のデータ転送能力は700KB/s

であった。バリア同期の処理時間として同期成立からPUの動作が再開するまでの時間とすると1.5 μ sであった。

細粒度静的スケジューリング法¹⁵⁾により Whetstone ベンチマークのモジュールを各々単独に並列化を行って MSBM で実行したところ、各モジュールの並列度およびモジュール2を1とした相対実行時間は次のとおりとなった。

module2: 並列度 3, 実行時間 1, 逐次時間 1.1
 module3: 並列度 4, 実行時間 6.6, 逐次時間 8.9
 module4: 並列度 1, 実行時間 3.3, 逐次時間 3.3
 module6: 並列度 5, 実行時間 9.2, 逐次時間 18.5
 module7: 並列度 4, 実行時間 16.5, 逐次時間 37.8
 module8: 並列度 2, 実行時間 54.2, 逐次時間 60.2
 module9: 並列度 3, 実行時間 20.4, 逐次時間 24.0
 module11: 並列度 3, 実行時間 16.7, 逐次時間 17.3

Whetstone ベンチマークの実行にあたり使用PU台数をあらかじめ与え、その台数以内で並列実行できるモジュールは並列実行するようにモジュールレベルのスケジューリングを行った。また並列化されたモジュールの実行開始直後にバリアパターンを MCAM に登録し、終了直前にそのバリアパターンを無効とした。

4.2 性能評価

図7に使用可能PU台数以内でモジュールの実行順序を入れ換えることなしにデータ依存性を考慮して、できる限りモジュールを並列実行させたときの相対実行時間の最良値を示す。縦軸はすべてのモジュールを逐次実行したときの実行時間を1としたときの相対実行時間、横軸は使用可能PU台数である。最大性能は6PUを使用して、次のスケジューリングを行ったときに得られた。このときの相対実行時間は0.46であった。

mod2-mod3-mod6-mod8

mod4-mod7-mod9-mod11

7PU以上でこのプログラムを実行しても6PUでの実行以上の性能向上はできないのはモジュール

8の並列実行時間がほかのモジュールの2倍以上であることに起因する。

MSBMのsingle-barrier MIMDモードで実行させたとき、5PUの使用で相対実行時間は0.75となった。一方、FMMでは5PUの使用で相対実行時間は0.6となった。これはモジュール8と9が並列実行でき、モジュール9の終了後に11を実行することによって短縮できた。

このプログラムではバリアパターンの登録、削除に要する実行時間に比べてモジュールの実行時間が非常に大きいので並列実行の効果が十分に得られた。これによりもう少し短い時間間隔でバリアパターンの登録、削除を行ってもかなりの性能向上が見込めると推測できる。

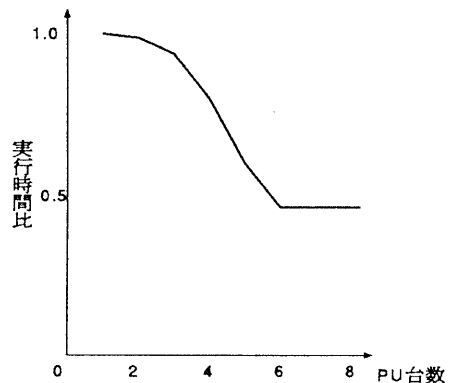


図7: Whetstoneの実行時間

5 むすび

細粒度の並列処理を効率よく行うための同期機構としてハードウェアによる半動的バリア同期機構とこの機構を効率よく使用するための拡張グループ共有メモリ構成を示した。そして、これらの機能を評価するためのテストベッドとして開発中の汎用並列計算機FMMのハードウェア構成と基本性能について述べた。Whetstoneベンチマークを用いてFMMの性能予測を行ったところ逐次実行に比べて相対実行時間は0.46、MSBMに比べ

て相対実行時間は0.61(=0.46/0.75)となり有効性が確認できた。

今後の課題として、FMMを製作すると共に半動的バリア同期機構に適した静的スケジューリングアルゴリズムを開発することである。

参考文献

- 1) 富田, 末吉: 並列処理マシン, オーム社(1989).
- 2) S.F.Lundstrom: Applications Considerations in the System Design of Highly Concurrent Multiprocessors, IEEE trans. on Computer, Vol.C-36, No.11, pp1292-1309(1987).
- 3) H.dietz et al: Static Synchronization Beyond VLIW, Supercomputing 89, pp416-425(1989).
- 4) R.Gupta: The Fuzzy Barrier: A Mechanism for High Speed Synchronization of Processors, 3rd Int. Conf. on ASPLOS, pp.54-63(1989).
- 5) 松本: Elastic Barrier: 一般化されたバリア型同期機構, 情処論, Vol.32, No.7.pp886-896(1991).
- 6) M.T.O'Keefe and H.G.Dietz: Hardware Barrier Synchronization: Static Barrier MIMD (SBM), 1990 Int. Conf. on Parallel Processing, pp.I5-I42(1990).
- 7) M.T.O'Keefe and H.G.Dietz: Hardware Barrier Synchronization: Dynamic Barrier MIMD (DBM), 1990 Int. Conf. on Parallel Processing, pp.I43-I46(1990).
- 8) D.J.Kuck, E.S.Davdson and D.H.Lawrie: Parallel Supercomputing Today and the Cedar Approach, 信学論, Vol.J71-D, No.8, pp.1361-1374(1988).
- 9) L.n.Bhuyan, Q.Yang and D.P.Agrawal: Performance of Multiprocessor Interconnection Networks, Computer, pp25-37, Feb.(1989).
- 10) M.Dubois and C.Scheurich: Synchronization, Coherence and Event Ordering in Multiprocessors, Computer, pp9-21, Feb.(1988).
- 11) 松田 他: マルチプロセッサシステム PARK 上での並列 Prolog 処理系の実現, 情処論, Vol.29, No.1, pp38-46(1988).
- 12) 高橋義造編: 並列処理機構, 丸善(1989).
- 13) T.Hoshino: Ainvitation to the world of PAX, Computer, pp68-79, May(1986).
- 14) 本石彰, 濱田智雄, 米沢敏夫, 岩根雅彦: 細粒度マルチマイクロプロセッサMSBMの構成と性能評価, SWoPP 琉球'94, (1994).
- 15) 濱田智雄, 野口善昭, 前川孝徳, 岩根雅彦: 細粒度コードスケジューリング手法の提案とその評価, SWoPP 琉球'94, (1994).
- 16) 渡辺: 並列処理概説, コロナ社(1991).
- 17) 岩根雅彦, 重松保弘, 定家健治, 茶屋道宏貴, 金沢高貴: FPGAによるバリア同期用機能メモリの開発, 九州工業大学研究報告(工学), No.66, pp45-52(Mar.1994).