

細粒度マルチプロセッサ MSBM

岩根雅彦[†] 本石 彰[†]
野口善昭[†] 米澤敏夫^{††}

細粒度並列処理を効果的に行うためには高速なプロセッサ間の同期と通信が必要である。多重プログラミング環境に対して、少量のハードウェアを持った同期と通信機構を提案する。1つは包含検索機能付き機能メモリによるバリア同期機構であり、1つはマルチキャスト機能を持ったグループ共有メモリ構成である。これらを検証するためのテストベッドとして細粒度並列計算機 MSBM を開発した。MSBM は1台のホストコンピュータ、グループ共有メモリを持つ16台のプロセッサ、機能メモリによるバリア同期機構を持つインタフェースから構成される。タスク内で同期をとるプロセッサ組(バリアグループ)をバリア同期機構にタスク実行開始前に登録し、実行中は動的に変更せずに終了後に削除する。これを静的バリア同期管理 SBM とよび、ただ1つのバリアグループをタスクに設定する SsBM とタスク内で複数のバリアグループを設定する SmBM を考える。SsBM での doacross ループによる評価では16プロセッサで8.39倍の速度向上を得た。さらに Whetstone ベンチマークで性能評価を行った。SsBM での速度向上比は個々のモジュールで1.03から2.29、Whetstone 全体では1.32であった。SmBM での速度向上比は Whetstone 全体で2.35である。2つのモジュールを多重実行したときの実行時間の増加は0.00から0.09であった。MSBM における同期機構および通信機構の有効性が確かめられた。

Multiple Static Barrier Management MIMD (MSBM) Multi-processors for Fine Grain Tasks

MASAHIKO IWANE,[†] AKIRA MOTOISHI,[†] YOSHIAKI NOGUCHI[†]
and TOSHIO YONEZAWA^{††}

The high-speed interprocessor communication and synchronization are required to achieve the efficient processing of fine grain parallelism. Two mechanisms with small amount of hardware to improve them are proposed under the multiprogramming environment. One is the hardware barrier synchronization (HBS) using the modified content addressable memory. The other is the group shared memory organization (GSM) with the ability of multicasting. The parallel computer for fine grain tasks, called MSBM, has been developed to evaluate these mechanisms as a test bed. MSBM consists of 1 hostcomputer, 16 processors with GSM and the interface unit with HBS. The processors synchronizing each other are registered into HBS before the execution of the program and are removed from HBS after the execution. This is called the static barrier management (SBM). The speedup ratio is 8.39 for doacross-loop program using 16 processors. The speedup ratio is 1.03 to 2.29 for each module of Whetstone program and 1.32 for whole Whetstone program under the single barrier registration per the program. The speedup ratio is 2.35 for whole Whetstone program under the multiple barrier registration. The increment of the execution time is 0.00 to 0.09 under the multiprogramming of 2 modules.

1. はじめに

コンパイラなどでプログラム全域からの潜在的な並列性の抽出およびスケジューリングを静的に行い、複

数のプロセッサで細粒度並列処理させることが注目されているが、細粒度並列処理を行えるマルチプロセッサは少ない¹⁾。

細粒度並列処理ではプロセッサ間のデータ通信と同期が頻繁に行われ、これがオーバーヘッドとなる。このオーバーヘッドを削減するために、精密なスケジューリングによってプロセッサ間の同期を極力排除する方法²⁾やハードウェアによるバリア同期機構³⁾が提案されている。後者には、バリア同期概念の拡張と

[†] 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

^{††} (株)東芝北九州工場

Kitakyushu Works, Toshiba Corporation

してのファジーバリア⁴⁾, エラスティックバリア⁵⁾, ウルティメイトバリア⁶⁾がある. またバリア同期をとりあうプロセッサ組 (バリアグループ) の指定方法として, 同期識別タグによる指定⁴⁾, マスクレジスタによる方法^{4),5)}, バリアキューを用いた SBM⁷⁾ および CAM (Content Addressable Memory) を用いた DBM⁸⁾ がある. ファジーバリアおよびエラスティックバリアでは各プロセッサに分散された同期用ハードウェアで同期がとられる. 一方 SBM および DBM のバリア同期ではバリアプロセッサ⁹⁾ と呼ばれる 1 つの専用の同期用ハードウェアで同期処理が行われる. このようなバリア同期機構は単一タスクシステムを主な対象としており, タスク間の実行順序が不定である多重プログラミングシステムについて十分な考慮がされていない.

プロセッサ間通信として単一バスによる共有メモリ構成が多く使用されている. 細粒度並列処理のようなプロセッサ間通信が頻繁に行われる環境では, バスでの競合が頻発するためにバスの多重化¹⁾ やキャッシュメモリ¹⁰⁾ によりバスでの競合を軽減させているが, ハードウェア量の増大やキャッシュコヒーレンス問題が生じる. また各プロセッサが局所メモリと共有メモリを持つ部分分散型共有メモリのひとつである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競合を軽減できるが共有メモリをプロセッサ分散させるためにならぬかのアドレス変換機構¹⁰⁾ が必要である.

そこで, 細粒度の並列性を持った複数のタスクを効率よく処理するために, MCAM (Modified CAM)¹¹⁾ によるバリア同期機構およびタスク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を有する細粒度並列計算機 MSBM (Multiple Static Barrier Management MIMD) を機能の評価および並列アルゴリズムの開発のためのテストベッドとして開発した.

本論文では, バリアグループ管理の方法について述べ, つづいて MCAM によるバリア同期機構, 部分分散共有メモリであるグループ共有メモリ構成, 細粒度並列計算機 MSBM のハードウェアシステム構成, さらに doacross プログラムによる基本性能の評価および Whetstone ベンチマークプログラムを用いた単一タスク実行ならびに多重タスク実行での評価について述べる.

2. 基本設計

2.1 静的バリア同期管理

並列処理される 1 つの応用プログラムはメインプログラムやサブプログラムようないくつかのプログラ

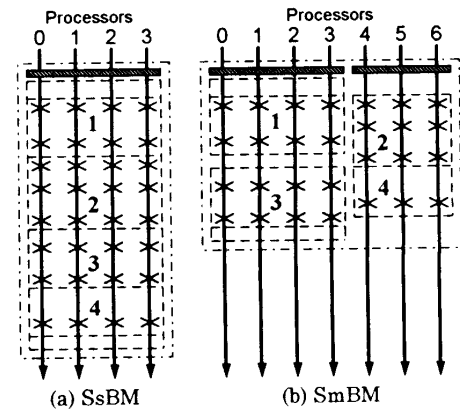


図 1 SBM
Fig. 1 SBM.

ムモジュールから構成されている. 1 つの応用プログラムの構成要素である個々のプログラムモジュールをプロセッサとよぶ. プロセッサごとに静的コードスケジューリング¹²⁾ を行ってバリア同期を挿入しておく. これらのプロセッサのリンクをとったのち, プロセッサを静的に割り当てて並列化されたタスクを形成する. すなわちオペレーティングシステムによりスケジューリングされて実行される 1 つの応用プログラムをタスクとよぶ. 1 つのタスクに静的に割り当てられたプロセッサの集合をプロセッサグループとよぶ.

タスクの実行開始前にバリア同期機構にこのタスクでバリア同期をとりあうプロセッサ組すなわちバリアグループを設定する. このタスクの実行中はバリアグループを動的に変更しないでタスク終了後にバリアグループを削除する. この方法を静的バリア同期管理 SBM* (Static Barrier Management) とよぶ.

SBM において 2 つの異なるバリア同期管理形態 SsBM (Static single Barrier Management) と SmBM (Static multiple Barrier Management) を考える. SsBM では 1 つのタスクにただ 1 つのバリアグループを設定して, タスクに割り当てられたすべてのプロセッサで同期をとる. SmBM では 1 つのタスクに複数のバリアグループを設定し, 独立かつ並列に複数のバリア同期をとる. 1 つのプロセッサ内では SsBM と同様にプロセッサに割り当てられたすべてのプロセッサでバリア同期をとるが, 並列実行可能なプロセッサにそれぞれ 1 つのバリアグループを設定して並列に実行する. 図 1 の (a) に SsBM, (b) に SmBM を示す. 図 1 の斜線での塗りつぶしはバリアグループのバリア同期機構への設定を, × はバリア同期を示

* 文献 7) の SBM とは異なる.

す。また一点鎖線はタスクを、破線はプロセッサを示す。SsBMではバリアグループは1つで、プロセッサは逐次実行される。SmBMでは2つのバリアグループの設定のもとにプロセッサ1, 3とプロセッサ2, 4が並列に実行される。

SBMではタスクまたはプロセッサごとに1つのバリアグループを設定するために本来必要としないバリア同期をダミーとして挿入しなければならないが、バリアグループの登録および削除がタスクごとに1回ですむ。またSBMではタスクごとに並列度が異なるので、使用していないプロセッサに別のタスクを割り当てて多重実行する。システムを多重プログラミング環境、つまり多重静的バリア同期管理 (Multiple SBM) とすることでシステム全体のスループットの向上をはかる。

2.2 バリア同期機構

2.2.1 単純なハードウェアバリア同期機構

ソフトウェアによるバリア同期ではバリア変数によって同期識別タグを持たせて任意のプロセッサ間のバリア同期をとる。これをそのままハードウェア化した場合にはハードウェア量が大きくなるので図2に示す同期識別タグを削除したバリア同期機構が一般的である^{3),6)}。単一タスクの実行では静的コードスケジューリングによってバリア同期の実行順序が一意に決定している。このようなバリア同期機構では、バリア同期をとらないプロセッサを同期をとる前にマスクすることによって任意のプロセッサ間でバリア同期をとることができる。これには各プロセッサ (PU) に1ビットのマスクが必要である。また n 本のバリア同期要求信号線とこれらの信号の論理積 (Wired AND) をとって各プロセッサに通知する n 本のバリア同期成立信号線が必要である。すなわち n 個のフリップフロップと $2n$ 本の配線が必要である。

多重プログラミング環境においてタスク間の実行順序は静的コードスケジューリングでは決定できない。したがってバリア同期の実行順序は一意には決まらない。正しくバリア同期をとるためには n 台のプロセッサからのバリア同期要求をそれぞれ独立にほかのプロセッサに伝達するための $n \times (n - 1)$ 本の信号線と、

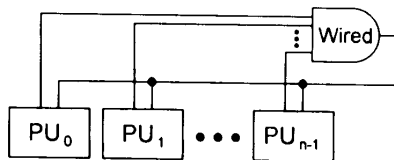


図2 バリア同期検出機構

Fig. 2 Barrier synchronization detection mechanism.

ほかのプロセッサからのバリア同期要求をマスクするための n ビットのマスクレジスタが各プロセッサに必要である。すなわち任意のプロセッサ間でバリア同期をとるために n^2 個のフリップフロップと $n \times (n - 1)$ 本の配線が必要である。

2.2.2 MCAMによるバリア同期機構

多重プログラミング環境ではタスクごとのバリア同期が互いに干渉して性能低下を起こさないように複数のバリアグループを並列に検索する機構が必要である。これを実現するためのMCAMによるバリア同期機構を図3に示す。このバリア同期機構はシステム内プロセッサ (PU) 台数に等しいビット長を持ったBRR (Barrier Request Register), BGR (Barrier Grant Register), MCAMで構成される。BRRおよびBGRの各ビットはそれぞれPUとBreq (Barrier request) 信号およびBgrt (Barrier grant) 信号で接続されている。またPUにはバリア同期用命令が用意されている。

多重静的バリア同期管理におけるこのバリア同期機構の基本動作はタスクの実行開始前にMCAMにバリアグループを書き込むことから始まる。次にPUが通常バリア同期命令であるBar命令を実行すると、Breq信号をバリア同期機構に送ったのちPUは実行を中断する。バリアグループのすべてのPUがBar命令を実行すると、バリア同期機構からこれらのPUにBgrt信号が送られてPUは実行を再開する。

バリア同期機構では送付されたBreq信号はBRRの当該ビットをセットする。BRRをMCAMの探索データとしてMCAMのすべてのエントリを同時に探索する。包含関係が成立すれば当該バリアグループがMCAMからBGRに読み出される。すなわち、

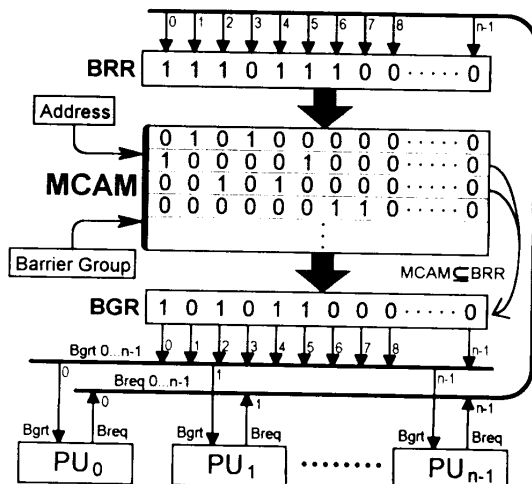


図3 MCAMによるバリア同期機構

Fig. 3 Barrier synchronization mechanism using MCAM.

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

ただし、 i は MCAM 内 i 番目のエントリ、 j は MCAM, BRR, BGR 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 $(MCAM_{ij} \subseteq BRR_j) \equiv (\neg MCAM_{ij} \vee BRR_j)$ である。なお、包含されるバリアグループが MCAM に複数存在する場合はそれらのバリアグループの論理和が MCAM から BGR に出力される。BGR に読み出されたバリアグループは $Bgrt$ 信号として PU に同期の成立を伝える。同時に BGR によって同期が成立した PU に対応する BRR のビットをリセットする。

マスクレジスタを持った単純なハードウェアバリア同期機構では同一 PU を同時には異なるバリアグループの構成要素にできない。同様の制約のもとで MCAM によるバリア同期機構で任意のバリアグループを設定するためには、基本的に $\lfloor n/2 \rfloor \times n$ ビットの MCAM 容量と $2n$ 本のバリア同期用信号が必要である。しかし包含論理は一致論理 $(MCAM_{ij} \wedge BRR_j) \vee (\neg MCAM_{ij} \wedge \neg BRR_j)$ よりも簡単であるので MCAM は通常の CAM に比べて少ないハードウェアで構成できる。

2.3 グループ共有メモリ

プロセッサ間通信は共有メモリによる方法が簡単である。しかし頻繁にプロセッサ間通信が発生する細粒度並列処理、さらにその多重プログラミング環境で、共有メモリのアクセス競合を低減させるために分散型共有メモリを考える。

共有メモリを個々のプロセッサに分散させて、プロセッサグループ（単にグループとよぶ）ごとに任意の大きさの共有メモリ空間をとることができるグループ共有メモリ（GSM）構成をとる。GSM は図 4 に示すような部分分散共有メモリである。各プロセッサのメモリ空間は境界レジスタ（SMBASE）により局所メモリ（LM）と GSM に分けられる。同一グループのプロセッサの SMBASE は同一の値を設定する。MPU が局所バスを介してプロセッサのメモリの全アドレス空間にアクセスできるように、局所バスには 2 ポートメモリ（DPM）の 1 つのポート（portA）が接続されている。DPM のもう 1 つのポート（portB）は GSM バスインタフェースに接続されている。MPU からメモリの全アドレス空間に読み出しおよび書き込みが可能であるが、書き込み時には GSM バスインタフェースの SMBASE とアドレスを比較し、GSM への書き込みならば GSM バスにグループ番号、アドレス、データがマルチキャストされる。受け取り側のプ

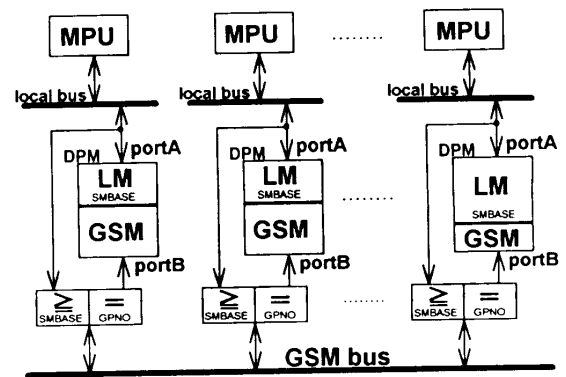


図 4 GSM 構成

Fig. 4 GSM organization.

ロセッサではグループ番号を比較し、一致すれば GSM バス上のデータを PortB から DPM に書き込む。タスクによっては GSM 空間を多く必要とすることがあるため、DPM の任意の容量を GSM として使用できるように SMBASE を任意に設定可能とし、不必要な GSM バスアクセスを排除してバストラフィックの削減をはかる。

また、GSM の相互排除処理のために lock/unlock 機構を用意する。lock はグループ内の他の PU の動作を停止させ、unlock は再開させる。あるプロセッサが lock 要求または unlock 要求を発行すると、その要求をプロセッサが属するグループの番号とともに GSM バスに乗せる。それ以外のプロセッサで GSM バスのグループ番号が自分のグループ番号と一致するものは、lock ならば処理を中断し、unlock ならば中断していた処理を再開させる。このように、タスクごとに独立して相互排除処理を行える。

3. ハードウェア構成

3.1 システム構成

並列計算機 MSBM は図 5 のようにホストコンピュータ（HC）、統合インタフェースユニット（IIU）、および 16 台のプロセッサ（PU）から構成されている。HC-IIU 間は HC バスで、IIU-PU 間は DT バス、INT バスおよび BAR 信号線で、PU 相互は GSM バスでそれぞれ結合されている。

HC は 80386MPU、8MB メモリ、各種入出力装置および IIU インタフェースから構成されたパーソナルコンピュータである。MS-DOS が稼動しておりユーザインタフェース、入出力管理などを行う。さらに IIU のための BIOS レベルの基本ソフトウェアが稼動しており、PU 管理やバリアグループ管理などを行うタスク管理ソフトの基礎となっている。タスク実行に際し

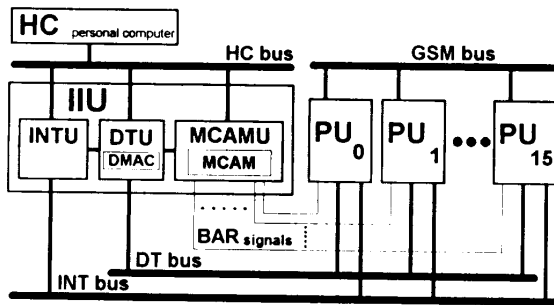


図5 MSBMのハードウェア構成
Fig. 5 Hardware organization of MSBM.

では必要数のPUでグループを形成し、コード転送するとともに同期機構にタスクで使用する全バリアグループを登録する。そして実行終了後にPU解放とともにバリアグループの同期機構からの登録解除を行う。

IIUはデータ転送ユニット(DTU)、MCAMユニット(MCAMU)および割り込みユニット(INTU)で構成される。DTUはDMAを用いてPU、グループまたは全PU単位でHCメモリとPUメモリ間のデータ転送を行う。MCAMUは一種のバリアプロセッサであり、16ビット×8のMCAMを持ちMCAMによるバリア同期機構を実現している。各PUとの接続は基本設計ではBreq信号とBgrt信号の2本であったが、BAR信号1本で実現し配線数を削減している。またバリアグループのMCAMへの登録をHCからのみとすることでハードウェアを簡素化している。INTUはPUからHCへの割り込み処理要求であるサービスコール(SVCC)の中継処理を行う。PUからSVCC要求を受け取るとPUからパラメータを受け取り、HCへ割り込みをかけパラメータを渡す。

3.2 プロセッサ

PUは図6に示すようにMPU、メモリ、各種内部レジスタおよび各種インタフェース(i/f)から構成され、割り当てられた逐次コードを各々独立したクロックで実行する。

MPUは8088MPU+8087NDPである。メモリは基本設計ではすべてDPMとしたがROM:32KB、RAM:208KBおよびDPM:16KBで構成し費用削減を行っている。ROMには初期化ルーチンやHCとの基本的なインタフェースとなる最低限のソフトウェアが格納されている。RAMおよびDPMの一部でLMを構成しコードや局所データを格納する。DPMのその他の空間はグループ内の共有データを格納するGSMを構成している。

内部レジスタとして、PU選択レジスタ、GSM境界レジスタ、PU番号レジスタ、グループ番号レジスタ、

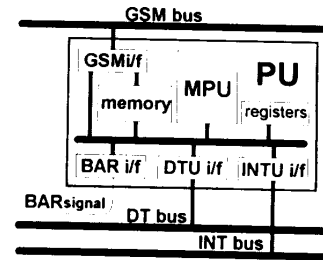


図6 PUの構成
Fig. 6 Organization of PU.

SVCCのパラメータ用レジスタおよび実行開始指示などのHCとの小規模通信用のレジスタなどが用意されている。DTU*i/f*はHC-PU間データ転送時のPUバスの制御を行う。PU、グループおよび全PUの転送単位の選択は、PU選択レジスタへのHCからの書き込みによって行われる。

INTU*i/f*はMPUがSVCC要求発行の出力命令を実行したときINTUへSVCCの中継処理要求を行い、承諾されるとINTUとパラメータ(16ビット×2)の転送を行う。PU間の要求調停はデジチェーン接続で行われる。

各PUのGSM*i/f*はGSMバスで接続されGSM構成とlock/unlock機構を実現する。GSM書き込み要求とlock/unlock要求はPU番号順に優先順位が設定された同じ調停機構で同等に調停される。

BAR*i/f*ではMPUの出力命令による通常のバリア同期Bar(Barrier)命令のほかにAdBar(Advanced Barrier)命令を加えた。そしてAdBar-Bar間をバリア領域とし、領域での同期を許すファジーバリア同期⁴⁾を実現している。先だつAdBar命令なしのBar命令実行では同期要求を行い実行を中断し同期成立まで待つ通常のバリア同期命令となる。AdBar命令を実行した場合、同期要求するが実行を続ける。AdBar命令実行済みでBar命令を実行したとき、同期が成立していれば実行を中断せず、成立していなければ実行を中断し同期が成立するまで待ち、ファジーバリア同期となる。

PUの各*i/f*、レジスタおよびIIUの各ユニットは、FPGA(Field Programmable Gate Array)を用いて製作されている。BAR*i/f*およびMCAMUなどのバリア同期関連のハードウェアの1PUあたりの割合は、8088MPU、8087NDPおよびメモリを除いたハードウェアに対し5%程度であり、8088MPU、8087NDPを含めたハードウェア全体に対しては無視できる規模である。図7にMSBMの全景を示す。

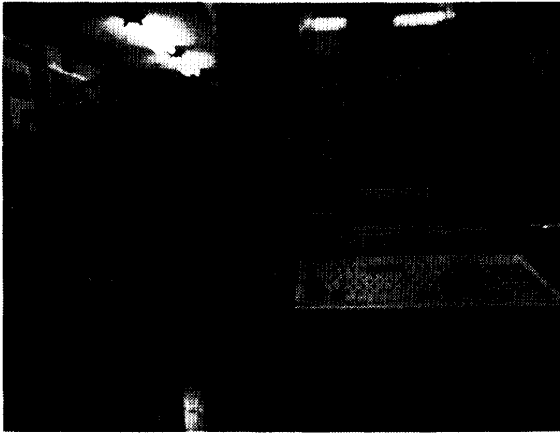


図7 MSBMの全景
Fig. 7 Picture of MSBM.

4. 性能評価

4.1 doacross ループによる評価

do ループの並列処理によってバリア同期の基本性能の評価を行う。図8(a)に示す簡単なプログラムをバリア同期を用いて並列化した細粒度 doacross ループプログラム (b) を実行する。逐次版の実行時間 T_s および各PU数による実行時間 T_p を測定し、逐次版に対する各PU数での速度向上比 T_s/T_p を算出した。その結果を図9に示す。

並列処理される部分の割合が大きいため、大きな速度向上が認められ、16PUで8.39倍の速度向上を得ている。ただしPU数が増加するほど速度向上の伸びがなくなっている。これはバリア同期によって毎回すべてのPUがほぼ同時に共有メモリに書き込みを行うため、競合による遅れがPU数に比例して増加してしまうためである。なおバリア同期を使用しないループディストリビューション¹⁾による並列処理では、16PUでの速度向上比は7.72であった。

4.2 Whetstone ベンチマークによる評価

4.2.1 Whetstone プログラムの並列化

Whetstone¹³⁾プログラムは異なる特徴を持つ分離性の高い複数のモジュールで構成されている。そこでまず異なる性質のタスクに対する性能を求める意味で各モジュールを1つのタスク(モジュールタスクとよぶ)と見なして並列化した。すなわちモジュールごとに Duplication-Scheduling Heuristic¹²⁾およびソフトウェアパイプライン手法¹⁴⁾に従い、手作業でできる限り並列化を行って SsBM のバリア同期を挿入した。モジュール8以外は単項および二項演算を基本とした変数レベルでの並列化を行い、モジュール8は変数レベルでの並列化ができないので機械語レベルでの並列

```

p = number of PU (n >> p)
k = 0(PU0),1(PU1),2(PU2)...
for i:=0 to n-1 do begin
  a[i] := b[i];
  c[i] := a[i+1];
end
(a) serial

for i:=k to n-1 step p do begin
  t := a[i+1];
  Barrier;
  a[i] := b[i];
  c[i] := t;
end
(b) parallel (doacross)

```

図8 doacross プログラム

Fig. 8 doacross programs.

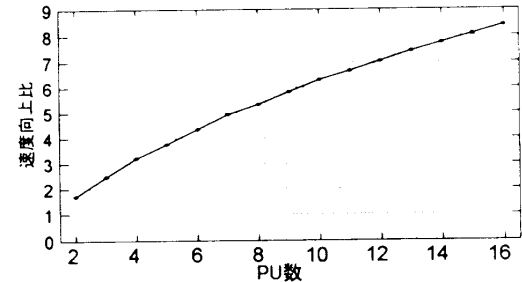


図9 doacross プログラムの速度向上比

Fig. 9 Speedup ratios of doacross programs.

化を行った。表1に評価に使用した各モジュールの並列化プログラムの特徴を示す。バリア同期間隔は並列化コードから得られる推定値であり、複数で範囲表記(～)でないものは、同期がそれらの間隔でその順序で繰り返されることを意味する。またSMWRはタスク中の全メモリアクセスに対するGSM書き込みアクセスの比率である。並列化率は並列化タスク中の並列化部分の比率を示し、限界値はアムダールの法則¹⁵⁾によって得られた速度向上比の限界値を表す。なお、MCAMによるバリア同期は要求から成立まで待ち合わせがない場合4クロックで行える。

さらにWhetstoneプログラム全体をSsBMおよびSmBMによって並列化した。SsBMでは上述の並列化された各モジュールを逐次実行する形となる。最大のPU数がモジュール6の5であるため、PU数が5に満たないモジュールはそれに合わせてバリア同期のみのコードを余剰PUで実行する。SmBMでは各モジュールを並列実行する形となる。各モジュール間の依存関係を維持し、並列実行可能なモジュールを並列実行させるが、可能な限りPU数は小さくする。図10にスケジューリング結果を示す。タスク全体のPU数は7であり、2つのバリアグループをMCAMに登録して実行する。並列実行するモジュールの呼び出しおよび終了待ち合わせはGSM上の共有変数で行う。

4.2.2 モジュールタスクにおける評価

各Whetstoneモジュールにおいて、MCAMによるバリア同期機構を用いたSsBMでの並列化タスク

表1 並列版モジュールの特徴

Table 1 Characteristics of parallel modules.

module	PU 数	バリア同期間隔 [マシンクロック]	SMWR[%]	並列化率 [%]	限界値
2	4	301	7.27	60.58	1.83
3	4	307	6.30	61.18	1.85
6	5	220~371	5.58	71.43	2.33
7	4	2478~2589	1.52	75.51	2.31
8	2	919	4.15	24.56	1.14
9	3	198, 32, 35	5.93	99.02	2.94
11	3	1237, 1748, 140	1.33	5.77	1.04

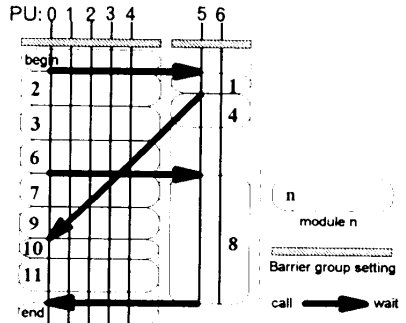


図10 SmBMでのモジュールのスケジューリング
Fig. 10 Modules scheduling under SmBM.

(HWB-P)とソフトウェアによって必要なプロセッサ間でバリア同期をとる並列化タスク (SWB-P) を、逐次タスク (S) と比較して評価を行った。SWB-Pではバリア同期を共有メモリ上のカウンタとビジーウェイトによって実現し、カウンタ更新時には相互排除処理 (lock/unlock) を行っている。各モジュールタスクのS, HWB-P, SWB-Pの実行時間 T_s , T_{hwbp} , T_{swbp} を測定し、Sに対するHWB-P, SWB-Pの速度向上比 T_s/T_{hwbp} , T_s/T_{swbp} を求めた。その結果を図11に示す。ただしバリア同期間隔を粒度の基準としモジュールタスクを左から粒度が大きい順に列挙している。

HWB-Pではモジュール7, 6の速度向上比は2.29, 2.02でその他のモジュールは逐次プログラムより少し良い程度である。一方, SWB-Pではモジュール7を除いて著しい性能低下が起きており, 最悪の速度向上比は0.21である。モジュール7は粗粒度の並列性を持ちGSMアクセスが少なく, 限界値も2.30である。HWB-Pでは限界値近くまで速度が向上しているが, SWB-Pではバリア同期によるオーバーヘッドのために速度向上比は1.65である。

モジュール11, 8は中粒度, モジュール3, 2, 6, 9は細粒度の並列性を持っている。これらのモジュールではソフトウェアによるバリア同期処理時間がバリア同期間の本来の処理時間と比較して大きいことが性能

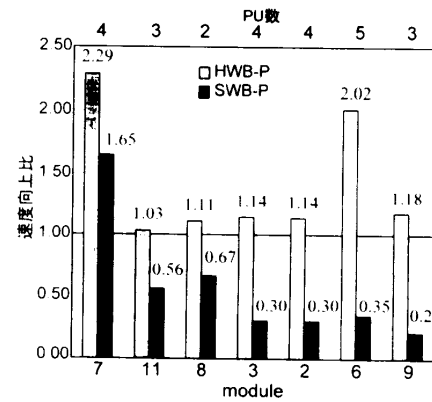


図11 逐次版Sに対するHWB-P, SWB-Pの速度向上比
Fig. 11 Speedup ratios of HWB-P, SWB-P to serial modules.

低下の原因になっている。すなわちSWB-Pでは粒度が小さいほどバリア同期によるオーバーヘッドが顕著に現れてくる。モジュール11, 8は並列化部分が少ないのでその限界値は1.04, 1.14と小さいが, HWB-Pでは限界値に近い速度向上比1.03, 1.11を得ている。モジュール3, 2, 6は並列化部分は60~70%であるが, 大きいSMWRのために通信によるオーバーヘッドが相対的に大きく, 速度向上比は限界値1.85, 1.83, 2.33よりも1.14, 1.14, 2.02と低くなっている。モジュール9は並列化部分が99%で限界値は2.94であるが, 最も細粒度であるうえにデータ依存が逆依存関係にある。すなわち逆依存解消のための一時変数の導入によるオーバーヘッド, 非常に短いバリア同期間隔のほかにすべてのPUが同時にGSMに書き込むことによる競合の発生のために速度向上比は1.18と小さい。このようにHWB-Pにおける速度向上比の低下は共有メモリ構成が主な原因である。そしてMCAMによるバリア同期機構を用いたSsBMは細粒度並列処理に有効であることが示された。なおモジュール8以外のモジュールもより細かい機械語の並列化を行えば, もう少し速度向上比の向上が見込まれる。

表2 WhetstoneのSsBM, SmBMによる速度向上比
Table 2 Speedup ratios of Whetstone under SsBM, SmBM.

	SsBM	SmBM
PU数	5	7
速度向上比	1.32	2.35

4.2.3 プログラム全体における評価

Whetstoneプログラムの各モジュールをプロセッサ(モジュールプロセッサとよぶ)とし、プログラム全体をタスクとする。SsBM および SmBM による並列化タスクを逐次タスク S と比較し評価する。S に対する SsBM, SmBM の速度向上比を求め、その結果を表 2 に示す。

SsBM では並列化されたプロセッサを逐次実行する。すなわち、5台のプロセッサでプロセッサグループおよびバリアグループを構成して各モジュールプロセッサを逐次実行している。SsBM における同期通信のオーバーヘッドを除いたときの理想的な速度向上比は 1.54 である。これはプログラム全体の実行時間の 40% 近くを占めるモジュール 8 の速度向上比が低いことが主な原因である。速度向上比の実測値は 1.32 であり、理想値より少し低い程度であった。

一方 SmBM では 7 台のプロセッサでプロセッサグループを構成し、5 台のプロセッサで 1 番目のバリアグループを、2 台のプロセッサで 2 番目のバリアグループを構成する。モジュール 2, 3, 6, 7, 9, 10, 11 が 1 番目のバリアグループで、モジュール 1, 4, 8 が 2 番目のバリアグループで実行される。モジュール 2, 3, 7, 9, 11 は並列化されたプロセッサである。またモジュール 1, 4, 10 は逐次プロセッサであるのでバリア同期は存在しない。1 番目のバリアグループで実行されるモジュールプロセッサと 2 番目のバリアグループで実行されるモジュールプロセッサが並列に実行される。各モジュールの Thwbp から予測されるモジュール実行のクリティカルパスは 2, 3, 6, 8 であり、逐次タスクに対する SmBM の予測速度向上比は 2.37 である。モジュール自身の並列化に加えてモジュール間の並列化を行っていることにより、大きく性能が向上している。速度向上比の実測値は 2.35 であり予測値と大きく違わない。予測値がバリア同期機構でのバリアグループ間の干渉が無いときの値であることを考えると、そういった干渉による性能低下が起きていないことがわかる。

このように SmBM は PU 数が増大するが速度向上比に優れ、SsBM は速度向上比は劣るが PU 数を抑えることができる。

表3 多重実行タスクの単独実行タスクに対する実行時間比
Table 3 Execution time ratio of multiple tasks to single task.

測定タスク (module)	同時実行タスク (module)						
	7	11	8	3	2	6	9
7	1.00	1.00	1.00	1.00	1.00	1.01	1.00
11	1.00	1.00	1.00	1.00	1.01	1.01	1.01
8	1.00	1.00	1.00	1.01	1.01	1.02	1.01
3	1.01	1.00	1.02	1.03	1.02	1.02	1.02
2	1.01	1.00	1.00	1.02	1.02	1.08	1.02
6	1.01	1.01	1.04	1.03	1.05	1.01	1.07
9	1.02	1.01	1.04	1.03	1.02	1.09	1.00

4.2.4 多重プログラミング環境での評価

多重プログラミング環境下でのタスクの多重実行時には、MCAM によるハードウェアバリア同期機構によりバリア同期の競合による性能低下は起こらないが、GSM 書き込み時のバス競合による性能低下が予想される。2つのモジュールを多重実行した場合の実行時間を単独実行時の実行時間と比較し、多重プログラミング環境の性能評価を行う。1つのモジュールタスクが GSM バス優先権の高い PU 群で繰り返し実行されている状態で各モジュールタスクを任意のタイミングで実行開始させたときの平均実行時間 T_m を測定し、実行時間比 $T_m/Thwbp$ を算出した。その結果を表 3 に示す。

各タスクの GSM バスアクセス間隔のばらつきから、競合の発生にランダム性があり平均的な性能低下は小さくなっており、最大でも 1.09 倍である。SMWR が小さなタスクの性能低下はわずかである。SMWR が大きなタスクどうしでの性能低下は少し大きなものがあるが、その平均値は 1.02 であり多重実行ではほとんど影響はないと見なせる。また GSM バスアクセス間隔がほぼ等しい同一モジュールタスクどうしの多重実行では大きな性能低下が起きていない。これは一度競合が発生すれば GSM バスアクセスにずれが生じ、以後、競合の発生が少なくなるためと考えられる。

以上から多重プログラミング環境における GSM 構成の有効性が確認できた。

5. む す び

バリア同期管理方式として SsBM, SmBM の管理形態からなる SBM を示した。さらに静的コードスケジューリングされたタスクの多重プログラミング環境を実現する同期通信機構として、MCAM によるハードウェアバリア同期機構と GSM 構成を提案した。そしてテストベッドとして開発した並列計算機 MSBM を使用して、細粒度 doacross ループプログラムによ

る評価を行った結果、16PUで8.39倍の速度向上を得た。さらに、Whetstoneベンチマークプログラムの、モジュールタスクにおける評価、SsBMとSmBMでのプログラム全体における評価、および多重プログラミング環境での評価を行った。その結果、モジュール単体ごとのSsBMでは1.03倍から2.29倍の速度向上がみられ、ソフトウェアバリア同期によって多くのモジュールで速度低下がみられた。また、Whetstone全体の並列処理では逐次処理の場合に対して、SsBMで1.32倍、SmBMでは2.35倍の速度向上がみられ、SmBMによる多重処理の有効性が示された。一方、2つの細粒度タスクの多重実行では、GSM構成によって最悪で単独実行時の1.09倍の実行時間に抑えられることが認められた。以上から、MCAMによるハードウェアバリア同期機構とGSM構成の多重プログラミング環境での有効性が確認できた。

今後の課題は、ファジーバリア同期の有効性の確認、タスク実行中にバリアグループを半動的にあるいは動的に処理することで並列処理効率を向上させることである。

謝辞 本研究を遂行するにあたり多大なご助力ならびにご支援をいただいた(株)東芝北九州工場の方々に深謝いたします。

参 考 文 献

- 1) 笠原博徳：並列処理技術，コロナ社(1991).
- 2) 笠原博徳，成田誠之助，橋本 親：OSCAR (Optimally Scheduled Advanced Multiprocessor) のアーキテクチャ，電子情報通信学会論文誌 D, Vol.J71-D, No.8, pp.1440-1445 (1988).
- 3) Lundstrom, S.F.: Applications Considerations in the System Design of Highly Concurrent Multiprocessors, *IEEE Trans. Comput.*, Vol.C-36, No.11, pp.1292-1309 (1987).
- 4) Gupta, R.: The Fuzzy Barrier: A Mechanism for High Speed Synchronization of Processors, *Proc. Third Int. Conf. on ASPLOS*, pp.54-63 (Apr. 1986).
- 5) 松本 尚：Elastic Barrier：一般化されたバリア型同期機構，情報処理学会論文誌，Vol.32, No.7, pp.886-896 (1991).
- 6) 高木浩光，有田隆也，曾和将容：重複可能なバリア型同期のためのスケジューリングアルゴリズムとその性能，電子情報通信学会研究会資料 CPSY91-15, pp.91-97 (1991).
- 7) O'Keefe, M.T. and Dietz, H.G.: Hardware Barrier Synchronization: Static Barrier MIMD (SBM), *1990 Int. Conf. on Parallel Processing*, Vol.1, pp.35-42 (1990).
- 8) O'Keefe, M.T. and Dietz, H.G.: Hardware Barrier Synchronization: Dynamic Barrier MIMD (DBM), *1990 Int. Conf. on Parallel Processing*, Vol.1, pp.43-46 (1990).
- 9) Dietz, H., Schwederski, T., O'Keefe, M. and Zaafrani, A.: Static Synchronization Beyond VLIW, *Supercomputing 89*, pp.416-425 (1989).
- 10) 高橋義造編：並列処理機構，丸善株式会社(1989).
- 11) 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, pp.45-52 (Mar. 1994).
- 12) El-Rewini, H., Lewis, T.G. and Ali, H.H.: *Task Scheduling in Parallel and Distributed Systems*, Prentice Hall (1994).
- 13) Curnow, H.J. and Witchmann, B.A.: A Synthetic Benchmark, *Comput. J.*, Vol.19, No.1, pp.43-49 (1976).
- 14) Monica, L.: Software Pipelining: An Effective Scheduling Technique for VLIW Machines, *SIGPLAN Conf. on Programming Language Design and Implementation*, ACM (June), Atlanta, GA, pp.318-328 (1988).
- 15) Hennessy, J.L. and Patterson, D.A. (富田眞治，村上和彰，新實治男訳)：コンピュータ・アーキテクチャー設計・現実・評価の定量的アプローチ—，日経 BP 社 (1992).

(平成6年12月14日受付)

(平成8年3月12日採録)

岩根 雅彦 (正会員)



1946年生。1968年京都大学工学部数理工学科卒業。同年(株)東芝に入社。1980年同社退職後トヨタ自動車(株)，豊田工業大学を経て，1988年より九州工業大学工学部教授。工学博士。計算機アーキテクチャの研究に従事。電子情報通信学会会員。

本石 彰 (学生会員)



1970年生。1994年九州工業大学工学部電気工学科卒業。1996年九州工業大学大学院工学研究科電気工学専攻博士前期課程修了。同年(株)ニコンに入社。並列計算機アーキテクチャの研究に従事。

**野口 善昭** (学生会員)

1971年生。1995年九州工業大学工学部電気工学科卒業。現在、同大学大学院工学研究科電気工学専攻博士前期課程在学中。並列計算機アーキテクチャの研究に従事。またオペレーティングシステムとプログラミング言語にも興味を持つ。

**米澤 敏夫**

1944年生。1968年金沢大学理学部物理学科卒業。同年(株)東芝に入社。1992年同社北九州工場生産技術部長。1994年より同社四日市工場工場長。1974年大内記念賞、1979年全国発明特別賞を受賞。半導体関連技術の開発、研究に従事。
