

階層型通信メモリバス方式(H-COM) を用いたプロセッサ間通信方式

加久間 勝、中川 秀敏、白石 一彦

(金沢工業大学)

プロセッサ数の多いマルチ・マイクロ・プロセッサ・システムを対象として、多数のプロセッサ間の情報交信を、効率的にかつ経済的に行える新しいプロセッサ間通信方式を提案している。第一は、対称個別通信メモリバス方式(SICOM)である。これは、各プロセッサ間をICMと呼ぶユニットで完全結合したものであり、処理のオーバーヘッドが小さく、高速な通信が期待できる。第二は、階層化した通信メモリバス方式(H-COM)である。これは、各プロセッサ間を疎結合にし、通信制御プロセッサ(CP)と主通信制御プロセッサ(MCP)を経由させることにより、ICMの数を大幅に削減し、経済化を図っている。

“INTER-PROCESSOR COMMUNICATION SYSTEM BY HIERARCHICAL
COMMUNICATION MEMORY PATH SYSTEM” (in Japanese)

by Masaru KAKUMA, Hidetoshi NAKAGAWA, Kazuhiko SHIRAISHI

(Kanazawa Institute of Technology, 7-1, Oigigaoka, Nonoichi-machi
Ishikawa-ken, 921, Japan)

This paper describes two new inter-processor communication systems for the multi-micro processor system. The one is called SICOM (Symmetrical Individual Communication Memory Path System). In this system, all processors are connected by ICM (Individual Communication Memory) each other. SICOM system is expected less communication delay time and less processing over-head. The other is called H-COM (Hierarchical Communication Memory Path System). This system consists of three stages that are peripheral processor (PE), communication processor (CP) and main communication processor (MCP). PE is connected to CP through ICM, and CP is connected to MCP through ICM. H-COM is less expensive system than SICOM.

電子計算機が商用化されて以来、単体としての電子計算機の処理能力の向上に、不断の努力がはらわれてきた。その発展に寄与したのは、主として電子部品技術の進歩によるものであって、計算機アーキテクチャの変革によるものではなかった。しかし、最近のLSI技術の進歩によって、これまでのシングルプロセッサの構成から、マイクロプロセッサ程度の処理装置を多数使用し、これらの処理装置を並行実行させる事により、電子計算機の処理能力を向上させる複合計算機システムが着目されるようになってきた。これらのシステムには、ノイマン型計算機をベースとした負分散を図った負分散型マルチマイクロプロセッサシステムや、処理装置にそれぞれ異なる機能を分担させて並行処理する機能分散型マルチプロセッサシステムがあり、さらに並行処理が容易なデータ駆動原理を基礎にした非ノイマン型のデータフローマシンなどがある。

これらのシステムは広義の意味でのマルチプロセッサシステムであり、処理装置相互間で情報を更新しながら処理を実行するから、処理装置相互間の情報交信（以下プロセッサ間通信方式と称す）の性能の是非が、システムの性能を支配する重要な要因になると考えられる。

今後は処理装置（以下プロセッサと称す）の数がますます増大するものと考えられるので、これら多数のプロセッサ間の通信を効率的にかつ経済的に行える新しいプロセッサ間通信方式の開発が必要になってくる。この場合、システム開発を容易に進めるために各プロセッサの内部処理とプロセッサ間の通信制御との独立性を高めること、

及び通信時のプロセッサでのオーバーヘッドを大幅に減少させることが重要である。

プロセッサでの独立性を高める方法としては、プロセッサ間通信に関する制御を、各プロセッサでの内部処理と全く分離し、専用のハードウェアで行う方法が考えられる。本論文で述べている対称個別通信メモリバスシステム（SICOM）と階層型通信メモリバスシステム（HCOM）はこの理念を基礎としたものである。

まずSICOM方式は、個別メモリユニット（ICM）と称する2者間のアクセス制御機構を持つ、一種の共有メモリを、各プロセッサ間に配置して相互通信を行う方式である。この方式では、各プロセッサ間にそれぞれ個別の通信バスが設置されるので、

- (1) バス系の輻輳による待ち合わせが発生しない。
- (2) データ転送制御手順が不要になる。

などの利点が生じ、大幅なオーバーヘッドの削減が可能になる。しかし、SICOM方式は各プロセッサ間に個別にICMユニットを配置するため、 n 台のプロセッサの場合、 nC_2 個すなわち $n(n-1)/2$ 個のICMユニットが n^2 に比例して増大し、経済的な問題が生じる。そこで、SICOM方式の構成を変え、 n 台のプロセッサに対して n 台程度のICMユニットで相互通信を可能にする方式として、HCOMを考案した。

SICOM方式は完全対称（密結合）の構成を採っているのに対し、HCOM方式はSICOM方式の利点を十分に生かし、2階層からなるモジュール結合（疎結合）の構成を採っている。これは、16台の周辺プロセ

セッサを I C M ユニット結合した、通信制御プロセッサ C P を 1 つのグループとし、更にこれを 16 個の I C M ユニットで結合した主通信制御プロセッサ M C P の 2 段構成となり、272 個の I C M ユニットを用いて最大 256 台のプロセッサ間通信を可能にしている。

2. S I C O M 方式

2. 1 S I C O M 方式による電子計算機システムの構成

S I C O M (Symmetrical Individual Communication Memory Path System) によるシステム構成を図 1 に示す。各周辺プロセッサ (P E) はプロセッサ間に配置された個別通信メモリユニット (I C M ユニット) によって完全結合の構成をとり、プロセッサ間の通信はこの I C M ユニットのアクセスすることにより行われる。実行プログラムはホストコンピュータ (H C) より、それぞれ H C と P E に対応する I C M ユニットに分割されて供給される。各 P E はその供給されたプログラムに基づいて各 P E とプロセッサ間通信を行いながら処理を進めていく、その実行結果は I C M ユニットの通して H C へ送られるようにシステムが構成されている。このように各プロセッサは内部処理を中断することなくプロセッサ間での通信を行うことができる。また、各プロセッサは I C M ユニットにより完全結合の構成を採っているためハードウェア構成を変えずにソフトウェアの変更のみで、アレイプロセッサやパイプラインプロセッサは勿論、データフローマシンとしても使用することが可能である。

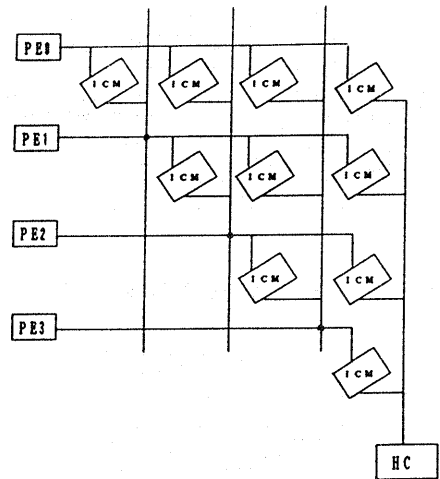


図 1 S I C O システム構成

2. 2 S I C O M の動作

I C M ユニット内のデータの流れを図 2 に、制御動作フローを図 3 に示す。各 P E から対応する通信メモリをリードするときは、相手の状態をかく乱することなく同時にかつ自由にアクセスできる。 P E から通信メモリにライトされたデータはメモリに格納されると同時にバッファに対してもそのデータとアドレスがセットされる。さらにその対応する P E の通信メモリへその P E がメモリフェッチしていない周期を利用して、バッファ内のアドレスに基づいて対応する通信メモリへデータを転送する。試作したシステムのサイクルスチールのタイミングチャートを図 4 に示す。

2. 3 ICMユニットの構成

図5にICMユニットの全体ブロック図を示す。ICMユニットは図1で示した様に2つのCPUの間で双方向通信を可能にするユニットでありどちらのCPUも送信側、または受信側になることがある、従ってAブロック、Bブロックが完全対称の構成となり、各PE間に配置され相互通信を可能にする。

メモリーコントロール回路は、CPUアクセスの時、及びサイクルスチールの時、両方の信号を得てメモリのチップセレクト信号及びライト信号を生成するためのものです。

リクエスト回路は、CPUからの書き込みが発生した場合、それを相手側のスチールコントロール回路に知らせる信号と相手側のバッファへの書き込みタイミングを発生する。

スチールコントロール回路では、相手側のリクエスト回路で発生した信号を受けて、スチールコントロール信号を発生すると共に2つのバッファへの切り替えを行う。

バスセレクト回路は、メモリに接続されているアドレスバス及びデータバスを、CPUからの書き込み時にはCPUに、メモリースチール時にはバッファからの出力へ切り替えを行う。

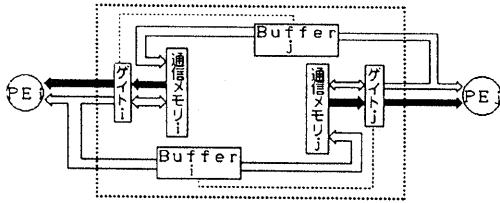


図2 データ転送過程概念図

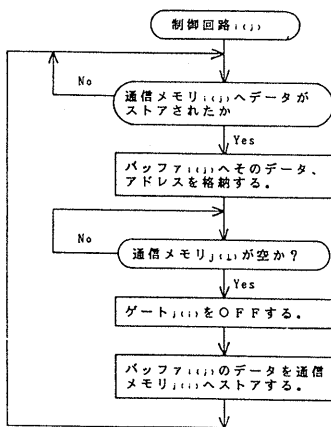


図3 制御動作フロー

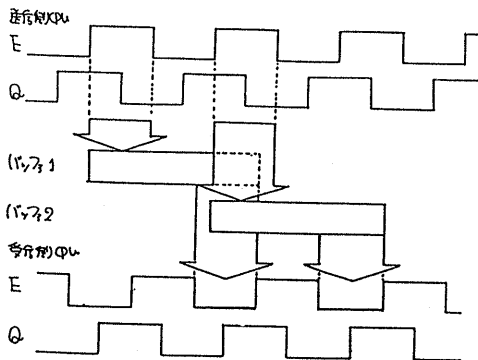


図4 サイクルスチールの
タイミングチャート

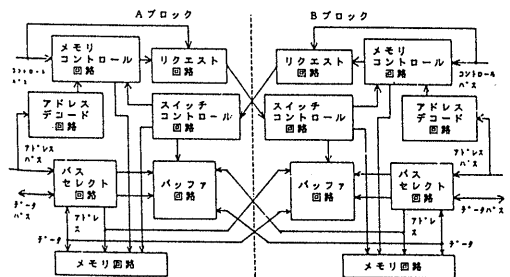


図5 ICMユニットのブロック図

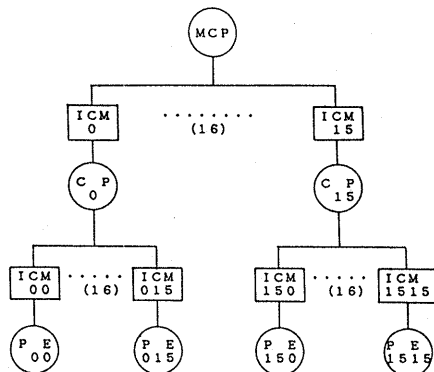
3. H-COM方式

3. 1 H-COMによる電子計算機システムの構成

H-COM (Hierarchical Communication Memory Path System) によるシステム構成を図6に示す。H-COMは16台の周辺プロセッサをICMユニットにより結合した通信制御プロセッサ(CP)を1つのグループとし、さらにこれを16個のICMユニットで結合した、主通信制御プロセッサ(MCP)の2段構成となる。従って272個のICMユニットを使い256台のプロセッサ間通信が可能となる。

1つのグループ内の通信はPE-ICM-CP-ICM-PEの順で通信が行われ、他グループとの通信にはPE-ICM-CP-ICM-MCP-ICM-CP-ICM-PEの手順で通信が行われる。

CPとMCPは通信制御のみを専用に行うプロセッサであり、今回の試作ではPEと同じマイクロプロセッサを用い、プログラム制御方式によって通信制御を行っている。データ転送遅延時間に関する性能の面では、CP及びMCPの性能が支配的になる。当初から予想したとうりであるが、通信制御をソフトウェアで行う今回の方式は、転送遅延時間が大きくなり、スループットも減少し、実用性の面では問題がある。しかし、この通信制御のアルゴリズムは簡単なものであるから、CP及びMCPを電子回路による布線論理制御方式で実現し、性能の大幅な改善を図ることは、比較的容易であると想定している。



CP: 通信制御プロセッサ
MCP: 主通信制御プロセッサ

図6 H-COMのシステム構成図

3. 2 通信制御方式

プロセッサ間通信の制御を行う方式として、リングバッファ方式とフラグセンス方式を検討した。

3. 2. 1 リングバッファ方式

(1) 構成の概要

リングバッファ方式は、ICMの通信メモリ上にリングバッファを構成しメッセージパケット単位でデータの転送を行うものである。図7にメッセージパケットの形式を示す。パケット形式はヘッダが3バイト、データ(DAT)が最大61バイトまでの可変長とする。ヘッダ情報は、受信PE番号(RN)、送信PE番号(SN)、メッセージデータ長(LEN)から成り、PE番号は上位4ビットがグループ番号、下位4ビットがグループ内の相対PE番号(0~15)を表現する。

ICM上でのリングバッファの構成を図8に示す。バッファは固定長のバッファを複数個組合せ、読み出し用ポインタ(r)と書き込み用ポインタ

(w) を付けリングバッファを構成する。リードポインタ (r) とライトポインタ (w) は、図の矢印の方向へ進むが互いに相手を追い越すことはな
 区、待ち行列制御を ICM 内で行うことができ、プロセッサの通信からの独立性が高くなり、非同期通信が可能と成る。なお、リードの時はポインタを 1 つ更新した場所より読み込み、ライトの時はそのポインタの示す場所に書き込む。実際に ICM 上では通信メモリを送信側と受信側に分割してそれぞれにリングバッファを構成する、又リングバッファを制御するリードポインタ、ライトポインタも同じく通信メモリ上に配置する。

(2) 動作の概要

CP は各 PE の通信バッファのポインタをルックイン方式で走査し通信要求が発生していたならば受信プロセッサへ転送する。その時、受信プロセッサがグループ内であるかグループ外であるか判断し、グループ内であればその宛先のプロセッサへパケットを転送する。グループ外であるなら MCP にパケットを転送する、以後この動作を繰り返す。MCP は各 CP からの送信バッファをルックイン方式で走査し通信要求が発生したならば宛先となるプロセッサのあるグループの CP へパケットを転送する、以後この動作を繰り返す。

PE_i から PE_j についての転送手順について述べる。この場合 PE_i と PE_j が同一グループ間での転送なのか、他グループに対しての転送かで、その通信処理は約 3 倍に増加する。

(I) PE_i は ICM_i にメッセージパケットを生成し、送信側ライトポインタを更新する。

(II) CP は ICM_i から通信要求を確認すると、宛先の ICM_j が転送可能か判断し、可能であれば転送を行い、ICM_i の送信側リードポインタと ICM_j の送信側ライトポインタを更新する。

(III) PE_j は受信要求を確認すると ICM_j にメッセージパケットを取り込み、ICM_j の受信側リードポインタを更新する。

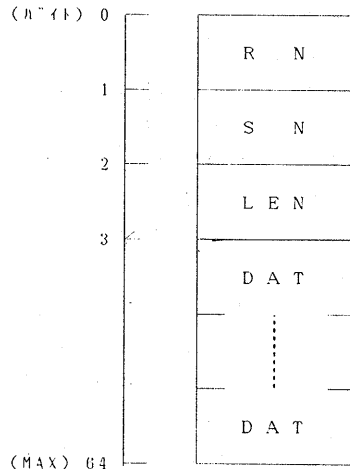


図 7 パケットの形式

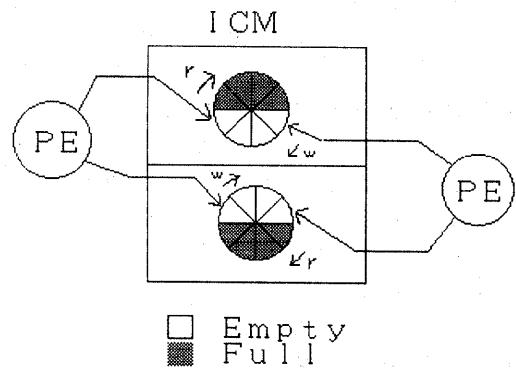


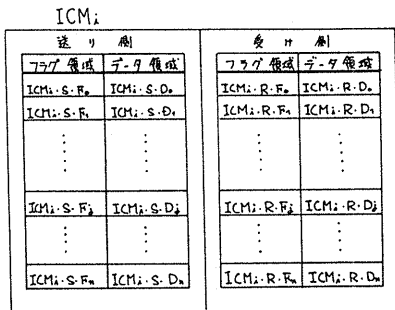
図 8 ICM 内のリングバッファの構成

3. 2. 2 フラグセンス方式

(1) 構成の概要

フラグセンス方式は、ICM内の通信メモリを仮想的に対応する周辺プロセッサ（PE）に対するバッファメモリとして固定的に分割し、転送データの有無をフラグによって表示する方式である。CPとPE間のICM内の通信メモリを図9のように分割する。

- (I) 送り側（PE → CP）と受け側（CP → PE）に2分割する。
- (II) それぞれをさらにフラグ領域とデータ領域に分割する。
- (III) 送り側のフラグ領域のアドレスは受信PEの装置番号と1対1に対応させる。
- (IV) 受け側のフラグ領域のアドレスは、送信PEの装置番号と1対1に対応させる。



ICM_i-S-F_j : ICM_i-Send-Flag_j
 ・ PE_j 宛 PE_i 宛 送信データに付随するフラグ領域

ICM_i-S-D_j : ICM_i-Send-Data_j
 ・ PE_j 宛 PE_i 宛 送信データのバuffer領域

ICM_i-R-F_j : ICM_i-Receive-Flag_j
 ・ PE_i の PE_j 宛 受信データに付随するフラグ領域

ICM_i-R-D_j : ICM_i-Receive-Data_j
 ・ PE_i の PE_j 宛 受信データのバuffer領域

図9 ICM内の分割

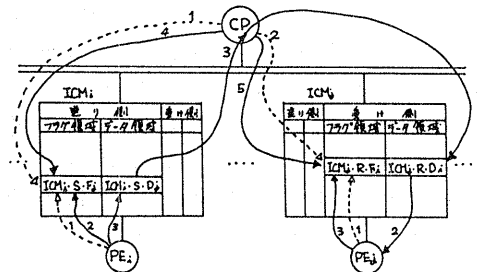
(2) 動作の概要

PE_i から PE_j へデータを転送する場合の通信の動作を図10に示す。

- (I) PE_i は PE_j と CP 間に配置された ICM_i の送り側のフラグ領域の j 番地をフェッチし、空いていればフラグを立てると

共に、データ領域の j 番地に転送するデータを格納する。

- (II) CP は常時 ICM の送り側のフラグ領域を走査しており、ICM_i の j 番地にフラグを検出すると、
 - (a) ICM_j の受け側のフラグ領域の i 番地のフラグを調べ空いているかどうか判断する。
 - (b) 空いていると、ICM_i の送り側のデータ領域の j 番地のデータを ICM_j の受け側のデータ領域の i 番地に転送する。
 - (c) ICM_j の受け側のフラグ領域の i 番地にフラグを立てると共に ICM_i の送り側のフラグ領域の j 番地のフラグをクリアする。
- (III) PE_j は常時 ICM_j の受け側のフラグ領域を走査しており、i 番地にフラグを検出すると、
 - (a) PE_i からの転送データと判断する。
 - (b) ICM_j のデータ領域の i 番地のデータを PE のメモリへ転送する。
 - (c) ICM_j のフラグ領域の j 番地のフラグをクリアする。



(番号はデータの流れる方向、2の数字は動作手順を示す)

←...: スキャン
 ←: リードアウト

図10 PE_i から PE_j に対する通信動作

4. むすび

高速で通信オーバーヘッドの小さいプロセッサ間通信方式であるS I C O M方式では、プロセッサ間をI C Mにより完全結合するために、I C Mの数がプロセッサ数の2乗に比例する経済的な欠点がある。これを改善するために、I C Mを階層構造にした新しいプロセッサ間通信方式であるH - C O M方式を提案し、その構成概要を説明した。この方式ではI C Mの数がほぼプロセッサ数に類似するので経済化の目的は達成しているが、C P及びM C Pでの制御をプログラム制御で行っているため、当初から予想したことではあるが転送遅延時間が増大する欠点がある。この問題を解決するためには、C P及びM C Pでの制御を電子素子を用いた布線論理化することが必要である。

S I C O M方式およびH - C O M方式について、I C Mを試作し、周辺プロセッサにM C 6 8 0 9を用いた小規模のシステムを構築し、機能の正常性を確認した。

参考文献

- (1) 中川、加久間：対称個別通信メモリバスシステム(S I C O M)によるプロセッサ間通信方式、信学会研究資S E 8 3 - 3 8 (1 9 8 3)
- (2) 加久間、中川：対称個別通信メモリバスシステム(S I C O M)による並行処理計算機の構成、信学会研究資E C 8 4 - 2 0 (1 9 8 4)