

可変構造型並列計算機の通信システム

3T-3

森 眞一郎 甲斐康司 村上和彰 福田 晃 末吉敏則 富田眞治
(九州大学)

1. はじめに

現在我々は、128 台の PE (Processing Element) を 128×128 の多重化クロスバー網 (MC-net : Multiplexed Crossbar Network) で相互接続し、PE間結合形態およびメモリ構造を可変としたマルチプロセッサシステム「可変構造型並列計算機」を開発中である。本システムでは PE 間通信をメッセージ交換により実現している。本稿では、このメッセージ交換を行なうための具体的な通信制御方式、および通信システムについて述べた後、その性能予測を行なう。

2. 通信制御方式

本システムでは、当初、MC-net 上での通信プロトコルとして、SEND 信号と ACK 信号のハンドシェイクによるクロック非同期転送方式を採用する予定であった。しかし、この方式では、クロスバー LSI および PE の MC-net インターフェースの動作周波数を 20MHz としても、MC-net の最大転送速度は 10MByte/sec にも満たず、クロスバー LSI の性能を十分引き出すことができない。これは、ハンドシェイク方式では、PE-クロスバー LSI 間を接続するケーブルの遅延時間 (往復で 100ns 以上の見積り) が、SEND 信号のサイクル時間を決定することによる。このような性能低下を防ぎ、クロスバー LSI を最大限に動作させる

ために、ハンドシェイク方式をやめてパイプライン制御によるクロック同期転送方式を採用することにした。この方式は、PE 間を結ぶ MC-net 上の通信路を遅延要素ごとに分割して、これを 1 つのパイプラインステージと見なし、ステージ間にあるラッチを同一クロックで制御する方式である。このときパイプライン・ピッチは各ステージの遅延のうち最大のもので決まる。

本 MC-net では図 1 に示すように、通信路を

- ① MCU (MS) - クロスバー LSI 入力ポート間
- ② クロスバー LSI の入出力ポート間
- ③ クロスバー LSI 出力ポート - MCU (MR) 間

の 3 つのステージに分割し、3 段パイプライン構成としている。また、MC-net は半二重回線であるためパイプライン上でも双方向のデータ転送を可能にしなければならない。そこで、本システムでは、ハンドシェイクによる方向制御シーケンスをデータ転送とオーバーラップして実行し、これを可能にしている。なお、クロック周波数は 20MHz とし、最大転送速度 20MByte/sec を予定している。

3. 通信システム

本システムの通信システムは図 1 に示すように、大きく 3 つに分割できる。以下、それらの機能について簡単に説明する。

3.1 ネットワーク・コントローラ (NETC)

通信システム全体を制御するユニットである。ホストからの指示に従ってネットワークを任意の構成に設定するよう、集中制御を行う。具体的には、通信機構を統括する 2 種類のシステムクロックの作成、MC-net のモード制御などの機能がある。

3.2 多重化クロスバー網 (MC-net)

8×8 のクロスバー LSI を 256 個用いて構成した 128×128 のクロスバー網である。MC-net は、各クロスバー LSI がもつ時分割多重化機能 (プリセットモード時) により、各種の PE 間結合形態をエミュレートすることができる。以下 MC-net の 3 つの動作モードについて述べる。

① デマンドモード : PE からの回線接続要求に対して回線の接続を動的に行なう。

② プリセットモード : プログラムの実行前に、PE 間の接続形態を時系列としてクロスバー LSI 内のスイッチ制御メモリに記憶させ、実行時には、このメモリの内容を順次読み出し、その内容に従って、クロスバー LSI 内および MC-net 全体の回線の接続を行う。

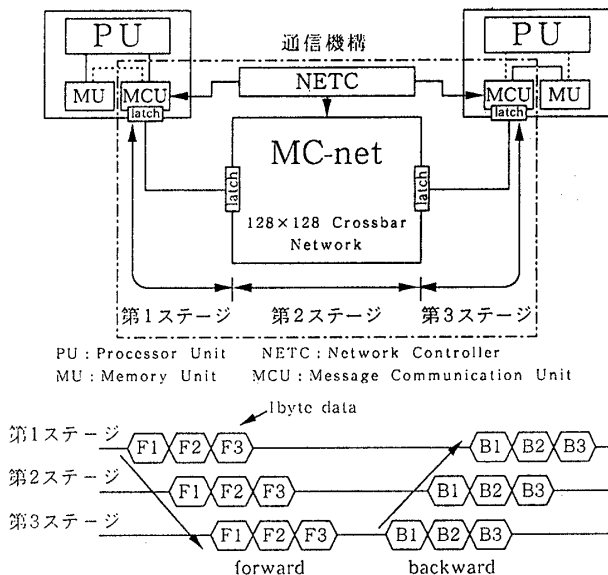


図 1 パイプライン制御によるクロック同期転送方式

③パーティション・モード：MC-netをデマンドモードで動作する領域とプリセットモードで動作する領域とに分割し、両モードが同時に混在することを許すモードである。

3.3 メッセージ通信ユニット (MCU)

MCUはPUとMC-netとの間のインターフェースの役割を果たしており、以下の3つの機能を提供している。

- ①OSのメッセージ交換機能のサポート：パケットの作成・分解、転送などの基本的な機能に加え、放送機能やマルチキャスト機能、さらにキューイング機能などをサポートする。
- ②MC-netとのインターフェース：NETCからの指示に従い、MC-netの動作モードに合ったデータ転送プロトコルでメッセージの転送を行なう。
- ③共有メモリウインドウ (SMW) へのアクセス：本システムは、分散メモリ構成を採るが、他PEのローカル・メモリ (リモート・メモリと呼ぶ) を直接アクセスすることを許している。MCUは、MC-netを用いてリモート・メモリ・アクセスを実現する。

4. 交換メッセージの種類

通信システムがサポートするメッセージは大きく2種類に分けられる。

- ①プロセス間メッセージ：通常のデータ送出メッセージの他に、他PEへの割込み、リセット、またはホールド、といった要求を伝えるためのメッセージが含まれる。
- ②SMWアクセス・メッセージ：メモリ共有型密結合マルチプロセッサにおいて必要なメッセージであり、SMWへのREAD/WRITEアクセス・メッセージ、Atomic LOAD-STORE (不可分) メッセージ、これらのアクセスに対するREPLYメッセージ、および2種類のキャッシュバージ・メッセージの5つのメッセージがある。

5. 通信システムの性能予測

以上述べた通信システムについて、プリセットモード時の性能予測を次の2つの項目について行なう。

5.1 プロセス間メッセージの転送速度

プロセス間メッセージ転送時のMC-netの実効転送速度を図

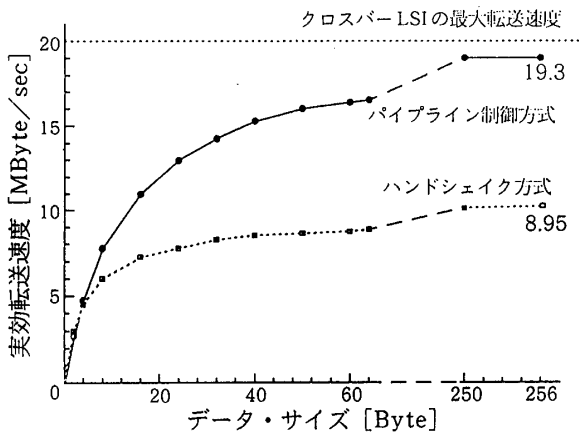


図2 プロセス間メッセージの実効転送速度

2に示す。ここでは、受信側MCUの受信バッファには十分な空き領域が存在すること、を仮定している。

MC-netの物理的な転送速度20MByte/secに対して256Byteデータ転送時の実効転送速度(最大実効転送速度)は19.0MByte/secとなる。よって、MC-netの利用効率は95%となり、最大実効転送速度はハンドシェイクによる非同期転送時の2.1倍に向上する。

5.2 SMWアクセス時のアクセス時間

図3はアクセスサイズとアクセス時間の関係を表わしている。ただし、メモリ側PEでの内部バス競合は起こらないことを仮定している。

図3よりSMWアクセス時のキャッシュ、ローカルメモリ、グローバルメモリへの4Byte READアクセス時間比は約1:3:24となることがわかる。本システムのキャッシュシステムはストアスルー方式を取っているため、実際には8Byteを超えるWRITEアクセスは生じない、従ってMC-netを経由したSMWアクセス時のアクセス時間はREADアクセス時1.31μs(1Byte時)~2.86μs(32Byte時)、WRITEアクセス時1.31μs(8Byte時)となる。

6. まとめ

以上、可変構造型並列計算機の通信システムについて述べた。パイプライン制御によるクロック同期転送方式を用いることにより2倍の通信速度が期待される。

参考文献

- 1) 村上ほか：“可変構造型並列計算機のシステム・アーキテクチャ”，情報処理学会「コンピュータアーキテクチャ」シンポジウム論文集，Vol.88，No.3，pp.165-174（1988年5月）
- 2) 森ほか：“可変構造型並列計算機のPE間メッセージ通信機能”，情報処理学会「並列処理」シンポジウム論文集（1989年2月）
- 3) 濱口ほか：“可変構造型並列計算機のプロセッサ・ユニット”，本大会論文集（1989年3月）
- 4) 蒲池ほか：“可変構造型並列計算機のメモリ・アーキテクチャ”，本大会論文集（1989年3月）

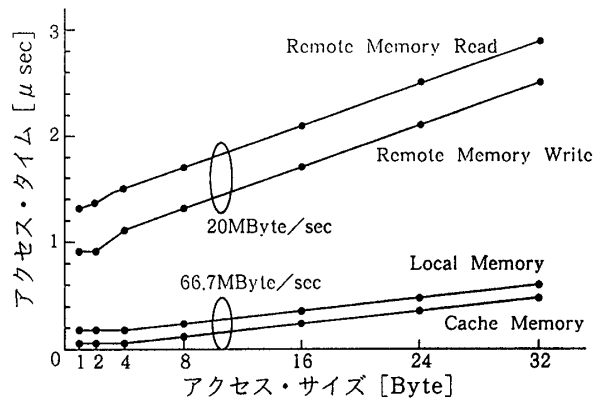


図3 共有メモリへのアクセス・タイム