

デュアル・バス型マルチプロセッサ・システム

4C-9

小原盛幹¹⁾、佐藤森芳²⁾、ワタリ・シゲル³⁾、永松礼夫³⁾、吉田紀彦⁴⁾、森下巖³⁾

1) 日本アイ・ビー・エム株式会社 サイエンス・インスティテュート
 2) 株式会社 東芝 3) 東京大学工学部 4) 九州大学工学部

1 はじめに

我々は共有バス型アーキテクチャのマルチプロセッサ・システムとして、デュアル・バス型の方式を提案した[1]。ここではこの方式で実際に試作したハードウェアとその性能について述べる[2]。

本システムは、並列に処理されるプロセス間で同期をとるための(メッセージと呼ばれる)小容量のデータを送るメッセージ・バスと、それ以外の大容量のデータを送るデータ・バスの2つのバスを共有する、共有バス型マルチプロセッサ・システムである(図1)。

実際に試作したシステムは、4台のコンピュータ・モジュール(CM)から構成されている。各CMはCPU(MC68000, 8MHz)、ローカル・メモリ、メッセージ・バスのコントローラなどからなっている。CPUとローカル・メモリには、既製品のVMEバス規格のボードを使用した。また、各CPUはプログラムとデータを自分のローカル・メモリに持っていて、独立に命令を実行する。

2 デュアル・バス方式

2.1 メッセージ・バス

各CMは送信用、受信用のメール・ボックスを持っている。CPUはメール・ボックスへの読み書きによって、他のCPUとの通信を行う。メッセージ・バスはこのメール・ボックス内のメッセージを転送する。

メッセージ・バスのコントローラは送信コントローラ(Sender)、受信コントローラ(Receiver)、Arbiterから構成されている。Senderは送るべきメッセージがあると、Arbiterに対して送信要求を出し、送信許可を受け取るとバスにメッセージを出力する。メッセージ転送の最初のサイクルでは宛先プロセスのIDが送信される。Receiverはバスを常に監視していて、自分宛のメッセージがあれば、自分の状態を返送する。Receiverが受信可能状態であれば、Senderはメッセージの続きを送信し、Receiverはメッセージを受信用メール・ボックスに取り込む。Receiverが受信不能状態であれば、Senderは送信をやめる。後者の場合、通常は再送信を行う。

メッセージ・バスはなるべく高速であることが望ましい。CPUのバス・サイクルとメッセージ・バスのサイクルの時間差を吸収するため、受信用、送信用にそれぞれメッセージ・バッファを持つことにした。これは前述のメール・ボックスに相当するが、メッセージ1つ分のスロットのみを持つ。2番目以降のメッセージはCPUのローカル・メモリ上に置かれる。

Arbiterはトークン・リング方式を採用した。バスの予約権を示すトークンがチェーンのリング上を回っている。トークンを受け取ったCMはバスの要求があるときにはバスの予約を行う。バス要求がないか、バスの使用を開始すると、トークンを次のCMに回す。ここでバスの予約を行っているのは、バス使用中に次の使用者を決定するためである。

実際に試作した回路では、トークンは62.5nsで次のCMへ伝わる。バス幅は1word(16bit)、バス上の転送時間は125ns/wordである。また、論理デバイスとして主にLS-TTLを使用した。ただし、一部高速性を要求されるところにS-TTL、FAST-TTLを使用している。ICの個数は1コントローラ(Sender、Receiver、Arbiter)当たり、約60個である。バッファにはアクセス・タイム35nsのスタティックRAMを使用している。

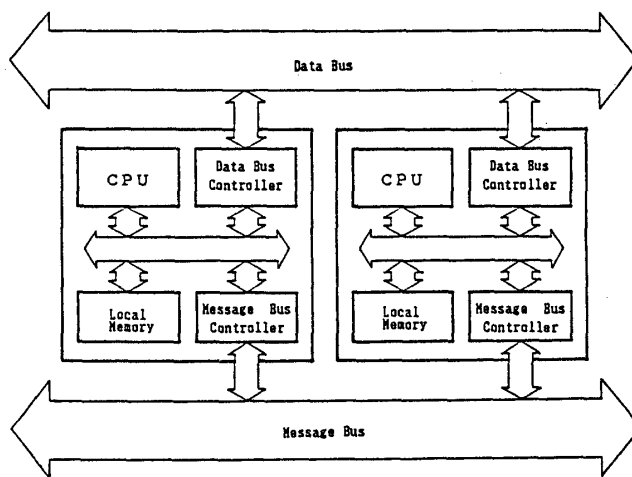


図1 デュアル・バス型マルチプロセッサ・システム

Design of Dual Bus Multiprocessor System

M. Ohara¹⁾, T. Satoh²⁾, S. Watari³⁾, L. Nagamatsu³⁾, N. Yoshida⁴⁾, I. Morishita³⁾

1) Science Institute, IBM Japan, Ltd.

2) TOSHIBA CORPORATION

3) University of Tokyo

4) Kyushu University

2.2 データ・バス

今回の試作システムで使用した既製品のCPUボードでは、ローカル・メモリをVMEバスからアクセスできない。このため、VMEバス上の共有メモリを経由してCM間のデータ転送を行っている。しかし、ローカル・メモリをVMEバスからアクセスでき、ローカル・メモリ間でDMA転送ができるようにすれば、VMEバス上より高速なデータ・バスを実現することができる。

3 メッセージ・バスの評価

メッセージは最初、CPUのローカル・メモリ上に組み立てられる。このため、CPU間のメッセージ通信は次の手順で行われる。

- 1) ローカル・メモリ上のメッセージをメッセージ・バスのバッファに転送する。
- 2) メッセージ・バスを経由して、相手のCPUのメッセージ・バスのバッファに転送する。
- 3) 相手のCPUが、バッファ上のメッセージを自分のローカル・メモリに読み込む。

2) は、メッセージ・バスのコントローラによって行われるが、1) と3) はソフトウェアによって行われる。1) ~3) を合計すると、CPU間の通信時間は3.88 μ s/wordであった。さらに、メッセージ・バスにはコンテンションが発生する。シミュレーションによって、コンテンションのある場合の通信時間を求めた。図2はメッセージの発生間隔やCMの個数を変化させたときの通信時間である。ここで、メッセージの発生間隔は、最大メッセージ発生間隔 (I_{max}) 以下の正値をとる一様乱数である。曲線Aは I_{max} を50 μ sとしてCMの台数を変化させたときの通信時間を示している。

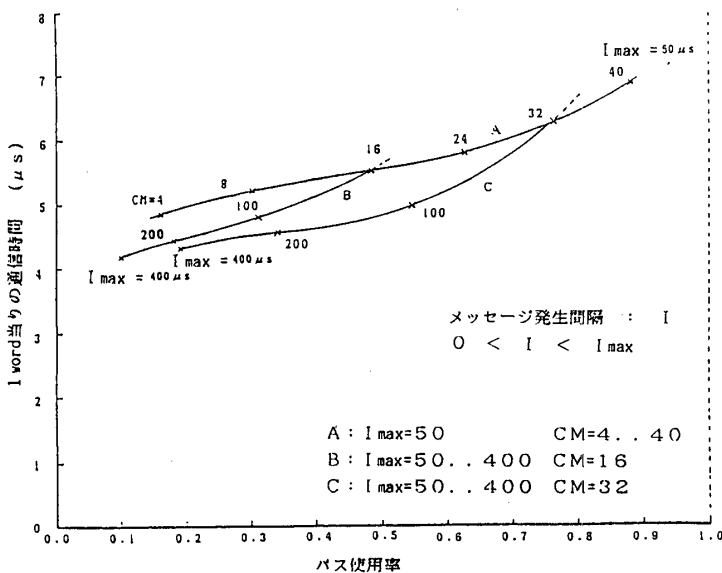


図2 シミュレーションによる通信時間

また曲線B、Cは、それぞれCMが16台、32台の場合に、 I_{max} を変化させたときの通信時間を示している。時間 I はCPUがOSやプロセスを実行する時間であるから、50 μ sは現在使用している68000としては非常に短い。また通常のTTLを使った場合、ファンアウトの制約から、いろいろ工夫を行ったとしても、1つのバスに接続できるCMの台数は高々32台である。このため、パラメータを変えても曲線Aと曲線Cには含まれた領域の中にある。最悪の場合でも通信時間がコンテンションのない場合の2倍以上になることはない。

4 システムとしての評価

本システムの上に、通信制御、プロセス管理などを行うOSのカーネルを実現してみた [3]。このOSでは、メッセージ通信にblocking send方式を採用し、メッセージの送信や受信のたびにプロセス切り換えを行っている。コンテンションがない場合、OSによるCPU間のメッセージ往復時間は845 μ s (メッセージ長=7 word) であった。

ここで、CPU間往復の転送時間は、54.3 μ s (3.88 μ s/word \times 7 word \times 2) であり、全通信時間の6%を占めるに過ぎない。プロセスの切り換えやその他のオーバーヘッドが大きいことがわかる。

5 おわりに

メッセージを転送するという意味では、十分な性能のメッセージ・バスを試作することができた。今後の課題は、通信にともなうOSのオーバーヘッドをなるべく少なくするようなハードウェアを提供することである。たとえば、現在バッファはメッセージを1つまでしか持つことができないが、全てのメッセージをコントローラのバッファ上に持ち、メッセージの管理を高速かつ容易にするなどの方法が考えられる。

参考文献

- [1] ワタリ、吉田、永松、森下：“メッセージ/データ通信 2重共有バス型マルチプロセッサシステム”、第32回情報処理学会全国大会論文集、5Q-6(1986)
- [2] 小原：“共有バス型マルチプロセッサシステムの研究”、東京大学工学部計数工学科卒業論文 (1986,3)
- [3] ワタリ：“Design of an Operating System Nucleus for Shared Bus Multiprocessor Systems”、東京大学工学部計数工学科修士論文 (1986,3)