

7L-2

階層構造並列計算機のための多数の結合回路を持つ
プロセッサノードの設計

廣田 勝久 高橋 義造
徳島大学工学部知能情報工学科
e-mail : hiro@n30.is.tokushima-u.ac.jp

1 はじめに

これまで、超並列計算機向けの相互結合網として、様々なものが研究されている。このうち多進木構造は、平均通信距離が短く、再帰的構造をもつためスケラビリティに優れているなど、超並列向きの特長がある。相互結合網として優れた特性を持つバイナリハイパーキューブの平均通信距離は、 n をプロセッサ台数とした場合、 $d(HC(2)) \sim 0.5 \log_2 n$

で表される。また、 k 進木の平均通信距離は、 $d(k) \sim 2 \log_k n$

で表わされる。これより k 進木の方が平均通信距離の短くなる k を求めると、 $k \geq 16$ となり、16進木以上では、平均通信距離について、バイナリハイパーキューブをしるぐ特性を持つことがわかる。

このように優れた特長を持った、図1のような多進木構造を持つ分散メモリ型並列計算機を構築するため、まず多数の接続ポートを持つプロセッサノード (PN) を設計した。以下の節では、マシンを構成するプロセッサノードに付いて述べ、さらにルータ、Network Interface Unit (NIU) の概要について説明する。

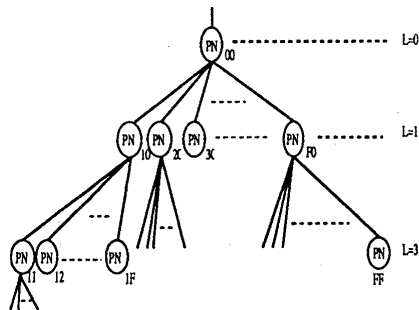


図1 16進木構造を持つ階層構造並列計算機

2 プロセッサノード

PNは、図2に示すように、MPU、ルータ、そしてMPU-ルータ間のインターフェイスを受け持つNIUから構成される。

Design of a processor node provided with many ports for hierarchical parallel computer.
Katsuhisa HIROTA and Yoshizo TAKAHASHI
Department of Information Science and Intelligent Systems, University of Tokushima.

MPUには、68Kをコアに持ち、周辺機能を1チップにまとめたTMP68301(12MHz)を用い、メモリは256KBのSRAMを実装している。

ノード間の通信は、専用のハードウェアであるルータとNIUが受け持つ。これによって通信は高速におこなわれ、またMPUの負荷も軽減される。このような構造のPNを複数組み合わせることによって、16進木構造の並列マシンを構築する。木構造のマシンでは、リーフ(葉)の部分にあたるPNはルーティングを行わないのでルータが不要である。16進木においては、リーフにあたるPNの割合が大きいので、ノード数に比較して、マシン全体におけるルータのハードウェア量が少なくて済む。

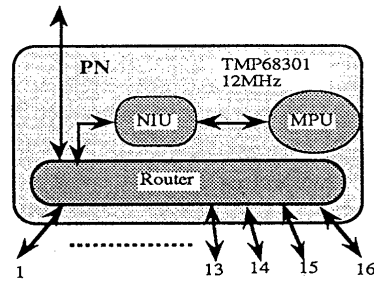


図2 プロセッサノード (PN)

3 ルータ

ルータは他のPNと接続する為のポートを多数持つ。16進木構造のマシンのルータには、上位につながる1本、下位への16本、及び自ノードに接続する1本の計18の双方向のポートがつながることになる。この様にルータには多数のノードがつながるので、高速な動作及び柔軟なフロー制御を実現することが要求される。

フロー制御の観点から、PN間の通信には、メッセージが通信回線を長時間占有する事のない様、バケット交換方式を採用し、バーチャルカットスルー (VCT) [1][2]方式を用いて、ネットワークのスループットを上げ、レイテンシを短縮するなど、ネットワークの高性能化を図る。

ルータのハードウェアは、データバス幅8ビット幅の、18ポート×18ポートクロスバスイッチを中心に構成し、20MB/secのデータ転送速度を実現する。

4 Network Interface Unit

NIUの構成を図3に示す。

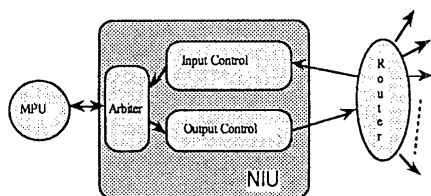


図3 Network Interface Unit

NIUはMPUから受け取った送信メッセージを、パケットに変換してネットワークに送り出すOutput制御部と、ネットワークから受け取ったパケットを、元のメッセージに復元し、MPUに受け渡す働きをするInput制御部、そしてこれらの各ユニットとMPU間のデータ転送の調停を行うアービタから構成される。各ユニットはLCAを用いて実現する。

Output制御部は、MPUからのコマンドによってDMAを起動し、メッセージを一旦NIU内のFIFOバッファに格納する。メッセージの最大長は4KWとしている。これ以上の長さのメッセージを送るときは何回かに分けて送るようにする。Output制御部はメッセージを受け取ると、それをパケット単位に分割してネットワークに送り出す。パケットサイズは67ワード固定長(ヘッダ3ワード、ボディ64ワード)である。(1ワード=16ビット)メッセージ送出の完了は、割り込みによってMPUに伝えられる。

Input制御部は128KWの受信バッファを持っている。ルータから自ノード宛に送られてきたパケットは、送信元PN、メッセージID別ごとに仕分けされ受信バッファ内に蓄えられて行く。受信バッファのデータ構造を図4に示す。

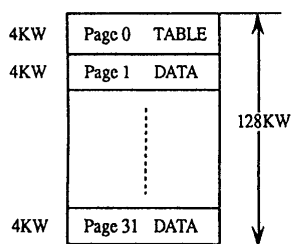


図4 受信バッファの構造

SRC ADRS(16)	
MSGID(16)	
COUNT(6)	LENGTH(6)
PAGE POINTER(16)	

図5 テーブルの1要素

受信バッファは4KWのページ単位で管理されている。このうち0ページは、現在保持しているメッセージの種類や状態などを記録しておくテーブル用の領域として用いる。テーブルの1要素を図4.3に示す。テーブル用のページを除いた残りの1ページから31ページまでを、メッセージの格納用に用いる。このため、一度に保持出来るメッセージの個数は31個以下に制限される。

メッセージを構成するパケットがすべて到着すると、Input制御部はこれを検知し、MPUに割り込みをかけてこれを知らせる。メッセージを構成するパケットがすべてそろったかどうかの判断は、テーブルのCOUNTとLENGTHの値を見て行う。COUNTは、今まで受け取ったパケットの個数で、LENGTHはメッセージを構成するパケットの個数である。この両者の値が等しくなれば、すべてのパケットが到着したと判断する。

並列計算機の通信網は、PN間の距離が分散システムなどに比べて短く、エラーレートも非常に低い。このため未到着のパケットがある場合には、タイムアウトなどによって、メッセージ全体を再送信するなどといった方法をとっても、さほどの効率低下にはつながらないと思われる。そのため、パケット再送などの、エラー訂正用の特別なハードウェア機構は設けていない。

メッセージが受信バッファ内に到着すると、MPUはDMAを起動してこれを受信バッファから抜き取る。

MPU側からは0ページが見える様になっており、到着メッセージのIDと、受信バッファ内での位置とサイズを知ることが出来る。MPUはこの情報を見てDMAのパラメータをセットする。DMA転送中は、パケットの受信は一時停止され、パケットはルータのバッファ内でブロックされる。

5. おわりに

現在プロトタイプを製作中であり、MPU部分の製作を終え、ルータとNIUの実装設計を行っている段階である。

参考文献

- [1]W.Dally. : Network and Processor Architecture for Message-Driven Computers ,in VLSI and PARALLEL COMPUTATION,ed. R.Suaya,G.Birtwistle. MORGAN KAUFMANN (1990)
- [2]S.Konstantinidou,L.Snyder:Chaos router: architecture and performance. IN Proc. Int. Symp. on Comput. Arch., pp.212-221,(1991)