

マルチPSI第2版のハードウェア構成

武田 保孝・中島 浩・益田 真直・浦原 和男
(三菱電機) (ICOT)

1 はじめに

マルチPSIシステムは、並列推論ソフトウェアの研究開発環境を提供することを目的とした並列計算機である。マルチPSIの開発は第1版と第2版の2段階に分けて行われており、既に完成した第1版は6台のPSI (Personal Sequential Inference Machine) を格子型ネットワークで接続したものである。現在開発中の第2版は第1版の評価結果を基に、プロセッサ数の増加及び接続ネットワークの改良による大規模化・高性能化を行なったものである。本稿では、接続ネットワークを中心としてマルチPSI第2版の特徴とハードウェア構成について述べる。

2 マルチPSI第2版の概要

マルチPSI第2版の全体構成を図1に示す。マルチPSI第2版は、フロントエンド・プロセッサとしてPSIを小形化・高性能化したPSI-IIを用いたバックエンド・マシンであり、第2版本体上では並列推論言語KL-1 (Kernel Language version 1) が実行される。

フロントエンド・プロセッサは、本体と接続するネットワークバス及びメンテナンスバスを介して入出力・デバッグ・保守などを行なう。本体は、64台のPE (Processor Element) を接続用ハードウェアを介して 8×8 の格子型ネットワークに接続したものであり、各PEのCPUはPSI-IIのCPUとはほぼ同じものである。PE間の通信を高速に行なうために、全PEが同期クロックで動作する。

マルチPSI第2版の特徴として、 P^3 (Processing Power Plane) を用いた負荷分散方式があげられる。 P^3 は処理能力が均一に分布するような底板平面であり、この平面を適宜分割しその各々をゴールに割り当てるこ

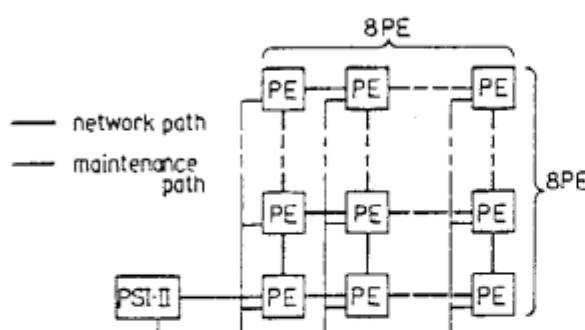
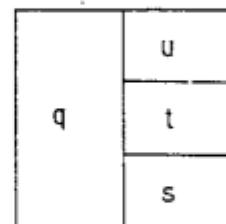


図1 マルチPSI第2版の全体構成

Hardware Configuration of the Multi-PSI/Version-2
Yasutaka Takeda, Hiroshi Nakashima, Kanae Masuda (Mitsubishi Electric Corp.)
Kazuo Taki (ICOT)



$p := q, r, s$
 $r := s, t, u$

図2 ゴールへの P^3 の割り付け

PE0	1	2
3	4	5
6	7	8

図3-1 静的負荷分散

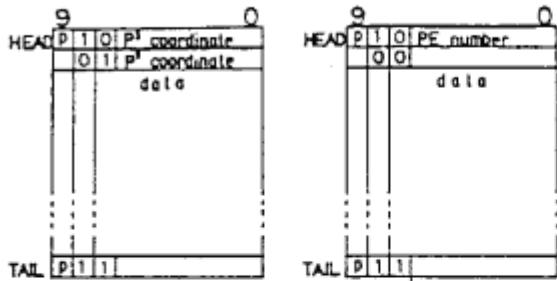
PE0	1	2
3	4	5
6	7	8

図3-2 動的再割り付け

により静的負荷分散を行なう。図2の例では、 p に割り付けられた P^3 領域を q と r で等分し、 r に割り付けられた領域をさらに3等分している。物理的なプロセッサとの対応付けは、 P^3 をプロセッサ数に分割し、各PEが自身に割り付けられた領域上のゴールを実行することにより行なう(図3-1)。計算過程においてPE間で負荷の不均衡が生じた場合、そのPE間の境界を動かすことにより受け持ち領域の面積を変えて負荷の均衡を図る(図3-2)。この方式は、局的に負荷の調整を行なえるため、大規模な並列計算機システムに適した方法と考えられる。

3 パケットのルーティング方式

マルチPSI第2版では、動的負荷分散を行なうために P^3 に於ける各PEの受け持ち領域は動的に変化する。そのため、個々のPEが全体の処理の分配状況を知ることは困難であり、かつ処理の局所性を考えると好ましくない。そこで、図4に示すように、PE間通信に用いるパケットを、その宛て先が P^3 の座標で示されたものと(a)、物理的なPE番号で指示されたもの(b)の2種類設け、ゴールの実行を要求するパケットは(a)を用いることとした。



(a) P3 座標表示 (b) PE番号表示
図4 パケットの構造

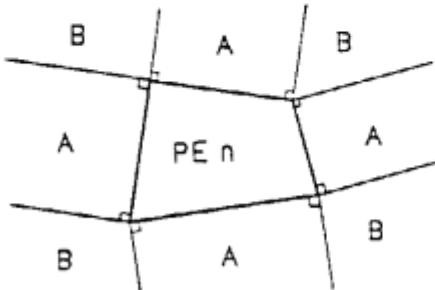


図5 パケットの転送方向決定法

このようなパケットの転送時には、パケットの通過経路上にあるPEも全体の状況を知らないため、パケットの転送方向決定はなんらかの局所的な情報で行なう必要がある。そこで、図5のように各PEが P^3 上の受け持ち領域の各頂点から辺に対する垂線を立て、パケットの目標座標が辺に垂直な区域Aに属するときはその辺を共有するPEにパケットを転送し、区域Bに属するときはその頂点を挟む両側の辺を共有する2つのPEのうちどちらか一方にパケットを転送することとした。このようにすれば転送経路上のPEの受け持ち領域とパケットの目標座標の距離は増加することなく、パケットは必ず目標座標に対応するPEに到着する。従って、各PEは P^3 における自身の受け持ち領域情報をだけから、パケットの正しい転送方向を求めることができる。

4 ハードウェア構成

図6に接続用ハードウェアの構成を示す。接続用ハードウェアはPEの内部バスに接続されており、その内部資源はCPUの特殊レジスタとしてアクセスできる。接続用ハードウェアは、網羅するPEの接続用ハードウェアへの4本のチャネル(Ch0~3)とCPUへのチャネル(Ch4)が5×4のスイッチで接続されている。パケットの転送は、同一チャネルへの複数の転送要求がないかぎり並列に行われる。各チャネルは、バリティを含めた10ビット単位で転送を行ない、1チャネルの1方向当たりの転送速度は5Mバイト/秒である。

Ch0~3には、 P^3 の座標またはPE番号に対応し

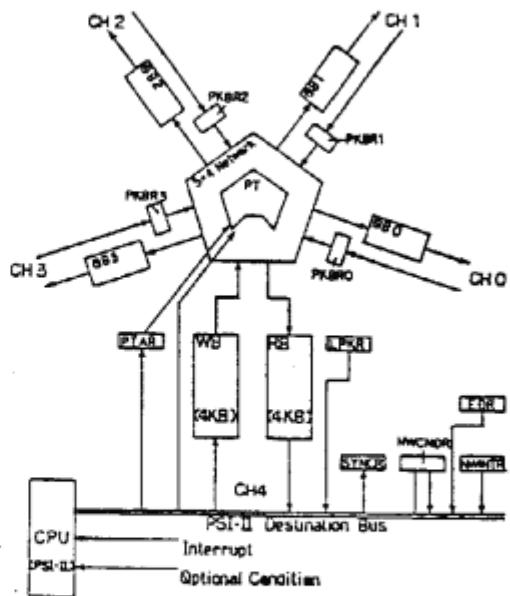


図6 接続用ハードウェアの構成

た転送方向を格納するPT(Path Table)がチャネル毎に設けられており、パケットの到着時にその宛て先を座標としてPTの内容を参照することにより転送方向が決定される。PTの内容は、負荷の調整が行われる度に更新されるが、更新の手間を軽減するために座標の下位ビットをマスクしてPTを引くことができるようになっている。即ち、PEの受け持ち領域から「遠い」座標点に関するテーブルを「粗く」することにより、更新すべきエントリを大幅に減らすことができる。転送方向の競合などで生じるネットワーク中の渋滞を緩和するために、Ch0~3の出力には48×10ビットの出力バッファが設けられている。

Ch4はPTを持たないため、CPUから送出されるパケットには転送方向を決めるために余分の2ビットが付加されている。またCPUへ渡すパケットには受けとったチャネル番号を示す2ビットが加えられている。Ch4は、入出力とも4K×12ビットのバッファを持ち、CPUから送出されるパケットはWBR(Write Buffer)中に全データが揃った時点でネットワークへの送出が開始され、CPUへのパケットもRBR(Read Buffer)内に全データが揃ってからCPUに対し到着が報告される。パケットの到着・送り出し及びネットワークの異常は、割込を用いてCPUに報告される。

5 おわりに

本報告では並列推論マシンである、マルチPSI第2版に付いてネットワーク接続用ハードウェアを中心とし、その特徴とハードウェア構成について述べた。マルチPSI第2版の特徴としては、 P^3 を用いた動的な負荷分散方式があげられる。