

# マルチPSIシステムとその接続方式

7B-1

木村康則 瀧和男 内田俊一

(財)新世代コンピュータ技術開発機構

## 1. はじめに

第5世代コンピュータの研究開発プロジェクトでは、並列推論マシンPIMの研究を主にアーキテクチャの観点から進めてきたが、その過程で、並列ソフトウェアに関する研究の重要性が明らかとなった[3]。マルチPSIシステムは、将来の大規模並列推論マシンの実現を目指して並列ソフトウェアの研究開発を進めていくために、PIMの完成に先駆けて研究開発ツールとして使用される疎結合の並列マシンであり、以下に示すスケジュールで研究開発を進めている[1]。

### (1) マルチPSI第1版(61年度)

要素プロセッサ(以下PEと略す)として前期に開発したPSIを6~8台格子状に接続したシステム。並列言語として並列論理型言語KLI(核言語第1版)を使用する[2]。処理系やネットワークハンドラなどはESPで書かれ、PE内の処理は、逐次型OSのSIMPOSDで進む。速度よりも開発効率を重視したシステムである。

### (2) マルチPSI第2版(62~63年度)

現在開発中のPSI-II[6]を16~64台接続したシステム。処理系は、KLIの機械語KLI-Bをファームウェアで実行することにより高速実行を目指しており、本格的な並列OS(PIMOS)を実装した実用的なシステムである。

### (3) PSI-II(61~62年度)

マルチPSI第2版のPEとして使うためにPSIを改良し、小型化、高速化したものである。また、単体でスタンドアロン型マシンとしても使えるように設計されている。

本稿では、マルチPSIシステムの研究項目などについて説明した後、現在開発を進めている第1版のPE間接続方式について報告する。

## 2. マルチPSIシステムの研究項目と役割

マルチPSIシステムで行おうとしている並列ソフトウェアの研究項目をまとめると以下になる[1]。

- (1) 論理型言語KLIとその処理系の開発
- (2) 並列OS(PIMOS)の基本部分の開発
- (3) PIMOSのうち大規模応用プログラムの効率的実行方式に関する研究

### 2.1 マルチPSI第1版の役割

第1版では、上記、(1)、(2)に重点を置いて研究開発を行う。すなわち、並列論理型言語KLIの分散環境での効率的な実現方式、KLIの上位言語となるべきシステム記述言語の開発や、並列OSのうち、メモリ管理、オブジェクトコードの分配、管理、入出力や割込みなどプロ

グラム実行に不可欠な部分の開発である。このとき、必要ならばファームウェアも含め、PSIシステムの変更も行う予定である。

### 2.2 マルチPSI第2版の役割

第2版では、第1版で得られた知見をもとに、大規模なプログラムを走らせるに足る高速な並列マシンを作成し、(3)の研究を行う。すなわち、プログラムが本来持つ並列性をいかにして引き出すか、局所性を考慮した、効率的な負荷分散方式はどうすればよいか、などについて検討する。これは、従来の逐次OSには無いが並列OSには必須である機能の検討と言うこともできる。また、中期PIM研究との関連においては、距離を意識したクラスタ間通信において、通信の局所性を保つ方法の検討と考えることができる[3]。

## 3. マルチPSI第1版の接続方式

### 3.1 概要

PSIは、元々スタンドアロン型のワークステーションとして開発したものであり、マルチプロセッサ用のハードウェア機能は持っていない。そこで、新たに接続用のハードウェア[5]を作成し、PSIを疎に結合した。

PE間の通信は、バケット交換によるメッセージ通信によって行われる。この処理の概要を図1に示す。

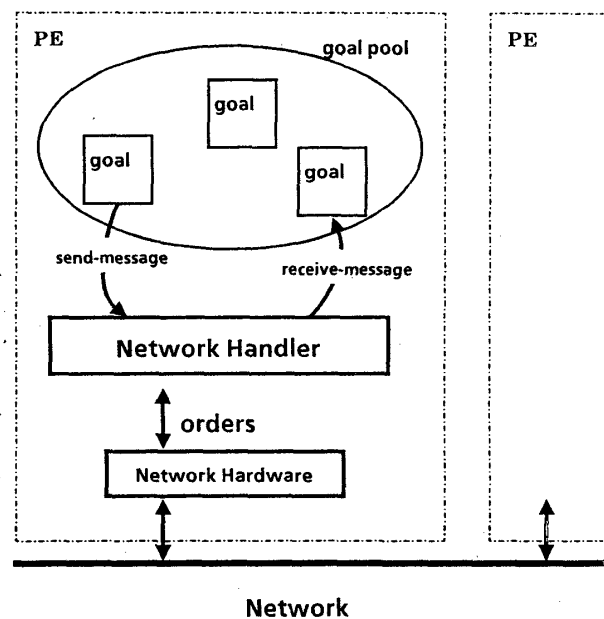


図1. 処理の概要

3.2 PE間通信処理

マルチPSI第1版システムにおいて或るPEの処理系から他のPEのデータをアクセスする場合の手順は以下の様になる。

- (1) KLI処理系において、外部参照タグを持ったデータが現われると、処理系内で管理しているテーブルを引き、参照先PEを求め、下位モジュールであるネットワークハンドラに対して参照要求をだす。
- (2) ネットワークハンドラでは、参照要求とデータから図2に示す様なパケットを構築し、バッファに格納する。具体的には、KLOで言うバイトストリングをパケット格納領域として用い、その先頭バイトに行先PE番号を、終了バイトに送出元PE番号を格納し、その間に変換されたデータをいれる。表1にパケットの通信オーダーを示す。
- (3) 次に、PE内の処理の隙間でためていたパケットをネットワークに送出する。具体的には、送出用に用意された組込述語に、パケットの格納されたバイトストリングと、その長さを与え、実行することにより行う。
- (4) ネットワーク上のパケットは、目的のPEまで自動的に転送される。
- (5) また、ネットワークよりパケットを読み込む時も組込述語により、バイトストリングに読み込まれたパケットをハンドラで処理系に意味のある形に再構成した後、処理系に渡す。

7	0
1	行き先PE番号
	パケット長(下位)
	パケット長(上位)
	PE間通信オーダー
	引数個数 N
	引数 1
	⋮
	引数 N
0	送り元PE番号

図2. パケット形式

3.3 PE間通信オーダー

パケットにうめ込む通信オーダーはKLI処理系と密接に関連しているものである。全部で14種用意したが、ここではその一部について簡単に説明する。

read, read\_answer は、他PEに対する値の問い合わせと返答に用いられ、read\_n\_level は、構造体などの深さnレベルまでの値の問い合わせに用いる。unify は、他PEに対する単一化要求である。throw\_goal は、負荷分散処理などの結果、他PEにゴールを投げるときに用いられる。第1版処理系では、ゴールを投げると、送出元PEのAND木には、“代理”が作られ、行き先PEには、この代理を指す“里親”メタコールが一般的には作られる。しかし、行き先PEに既にこの“里親”が存在するときには、二重作成をさげ、送出元PEに対して今作った“代理”を

消す要求を出さなければならない。このときに用いられるのがcancelである。このように、PE間にまたがるAND木を管理するためのオーダーを幾つか用意しており、kill, goal\_dead, goal\_terminates, goal\_failsなどがそれである(図3)。

表1. PE間通信オーダー

オーダー	引数
read	Var_id, Return_obj
read_n_level	Var_id, Return_obj
read_answer	Object_id, Value
unify	Var_id, Meta_call, Parent_record
throw_goal	Meta_call_id, Module_id, ...
goal_terminates	Parent_record, Reduction_count
goal_fails	Parent_record, Reduction_count
goal_dead	Parent_record, Reduction_count
kill	Parent_record_id
cancel	Parent_record_id
wake_up_goals	Var_id, Value

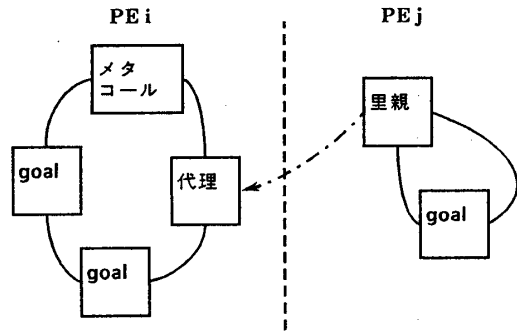


図3. 代理と里親ノード

3.4 ネットワーク制御用組込述語

ネットワーク制御用組込述語としては、送出用に、check\_write\_buffer, write\_buffer, 読取用に、check\_read\_buffer, read\_buffer, sense\_packet\_arrivalの5種を用意した。

4. おわりに

現在、ICOTにおいて、PSIを6台接続して第1版の試験、デバッグを行っているところである。今後は、第1版システム上でソフトウェアの開発を進めると共に、第2版システムに関する処理系及びネットワーク構成について検討を進める予定である。

<参考文献>

- [1] 瀧他: Multi-PSIシステムの概要, 第32回情報処理学会全国大会, 5Q-8, 1986-3
- [2] K. Ueda: GUARDED HORN CLAUSES, LPC '85, pp.225 1985-6, JAPAN
- [3] 後藤他: 並列推論マシンPIM-中期構想一, 本大会 3B-5
- [4] 宮崎他: Multi-PSIにおけるGHCの実行方式, LPC '86, pp.83 1986-6, JAPAN
- [5] 益田他: マルチPSIのネットワークハードウェア構成, 本大会 7B-2
- [6] 中島他: マルチPSI要素プロセッサPSI-IIのアーキテクチャ, 本大会 7B-3