

TR-0733

並列推論マシンPIMの
アーキテクチャ

中島 克人、益田 嘉直、中島 浩(三菱)
近藤 誠一、武田 保孝、村澤 靖
小森 隆三

January, 1992

© 1992, ICOT

ICOT

Mita Kokusai Bldg. 21F
4-28 Mita 1-Chome
Minato-ku Tokyo 108 Japan

(03)3456-3191 ~5
Telex ICOT J32964

Institute for New Generation Computer Technology

並列推論マシン PIM のアーキテクチャ

中島克人* 益田嘉直** 中島 浩*** 近藤誠一+ 武田保孝* 村澤 靖* 小森隆三*

1. まえがき

複雑で大規模な問題をより簡単にプログラムし、より高速に実行するという、ユーザーのコンピュータに対する要求はとどまるところを知らない。これに対して、コンピュータメーカーは半導体技術の目覚ましい進歩を背景にメモリの大容量化やプロセッサの高性能化という形で速度向上への期待にこたえてきた。しかし、單一プロセッサに関しては、今までのような性能向上をこれまでのペースで続けることが難しくなってきている。

バイブルайн、分岐予測などの従来からの大型計算機技術に加え、最近ではSuperscalar、VLIWなどのように複数の命令を同時に実行するための様々なアーキテクチャ上の改良が施されているが、それらによる速度向上には限界が見えてきている。半導体の集積度を十分に生かす道は、設計工数の点から見ても、もはやマルチプロセッサ化以外になくなりつつある。一方、ユーザーの立場からは、マルチプロセッサ化はプログラムを更に難しくする。それぞれのプロセッサを有效地に活用するための従来どおりのプログラミング技術に加え、複数のプロセッサを無駄なく働かせる、いわゆる負荷分散の技術が必要となるためである。

1982年から10年計画で始められた国家プロジェクト“第五世代コンピュータ”は上記の半導体技術の進歩を予測し、複雑な問題を高速に実行できる計算機の実現を、論理型言語と並列処理を組み付けることによって目指したものである。

三菱電機(株)はこのプロジェクトの推進母体である(財)新世代コンピュータ技術開発機構(ICOT)の委託を受け、この度プロジェクトの最終成果である並列推論マシンPIM⁽¹⁾の開発を完了した(図1)。

この論文では、大規模知識情報処理マシンを指向したPIMのアーキテクチャとその上に実装した高レベル並列言語KL1について述べる。

2. システムの特長

三菱電機(株)が、ICOTからの委託で開発したPIMは、PIM/mと称され、プロトタイプとしてやはり1987年に受託開発したマルチPSI⁽²⁾⁽³⁾と同様、二次元格子状に多数の要素プロセッサ(PE)を接続した疎結合型マルチプロセッサである。PIM/mではそれぞれのPEがローカルメモリを備えることから分散メモリ型とも称する。PIM/mの最大構成は

256PEである。

バス接続などで一つのメモリを複数PEで共有する共有メモリ型のマルチプロセッサでは、各PEにキャッシュメモリを配するとしても、メモリへのアクセス競合による性能低下のために、接続PE数はせいぜい8~16台程度が限界であると言われている。これに対してPIM/mのような分散メモリ型マルチプロセッサは、メモリへのアクセス競合が生じないことから、PEの台数に関する拡張性が高く、大規模並列マシンに適している。PIM/mで採用した二次元格子接続はPE間の最大距離がPE数の平方根に比例して増大するため、PE数が大きくなるにつれて任意のPE間の通信には配慮が必要となるが、実装が非常に容易であることから、1,000台規模程度までのシステムに適している。

分散メモリ型ではPE間の通信のオーバヘッドが大きいため、仕事の分散(負荷分散)が難しい。分散の単位(粒度といふ)を制御し、負荷分散とPE間通信オーバヘッドのバランス点を見いださなければならないからである。残念ながらこれをシステムで自動的に制御するような技術はまだ確立しておらず、ユーザーがそれぞれの問題に応じた負荷分散手法をソフトウェアで実装し、それをチューンアップしていくなくてはならない。そのためには、並列プログラムの記述を容易にし、負荷分散の技術開発をサポートするための強力なプログラミング言語と充実したプログラミング環境が必要である。

PIMで採用している並列論理型言語KL1(核言語第1版)は、①変数の値の待ち合わせによる同期機構によってバグの少ないコーディングができる、②宣言的な記述で小粒度

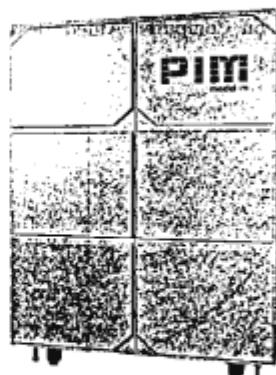


図1. PIM/mの筐体外観(32PE内蔵)
(最大構成は8筐体256PE)

74(682) *三菱電機情報電子研究所 ***同研究所(工博)
**同コンピュータ製作所 †(財)新世代コンピュータ技術開発機構

の並列計算を自然に記述でき、③仕事をどれだけまとめて実際の負荷分散の単位にするかは、プログラムのロジックを変更することなく後から指定できる、という特長をもつ。このため、ユーザープログラムのデバッグや性能改良、さらには新しい負荷分散手法の開発に大変適している。PIM/mではこのKL1を高速に実行するハードウェア及びファームウェアを備えている。

以下では、PIM/mのシステム構成、ハードウェア、言語とその実行方式などについて順に説明する。

3. ハードウェアアーキテクチャ

3.1 システム構成

最大構成はマルチ PSI では 64 PE であったが、PIM/m では 16×16 、すなわち 256 PE からなり、これらは 1 きょう(筐)体に 32 PE ずつ(図 1)、合計 8 筐体に収められる。PE

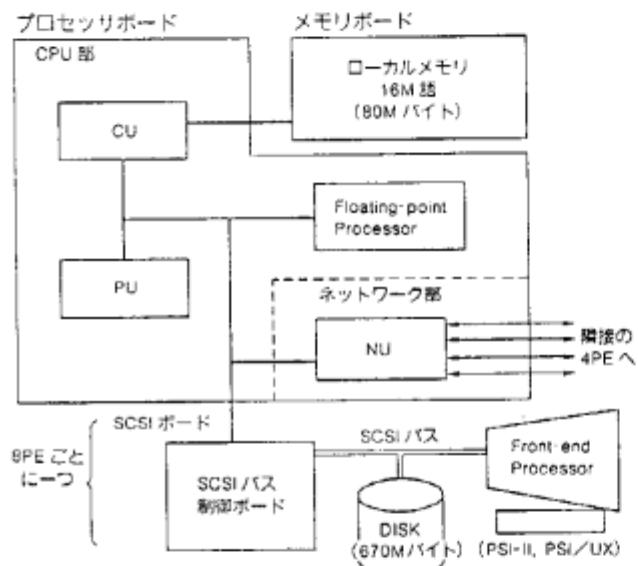


図 2. 要素プロセッサのブロック図

表 1. 要素プロセッサ(PE)諸元

	PIM/m	マルチ PSI
サイクル時間	1 枚 65ns	4 枚 200ns
使用 LSI	CMOS セルベース 0.8/ $1.0\mu m$ 3 種	CMOS ゲートアレー 8K/20K ゲート 9 種
マイクロ命令	64 ビット水平型	53 ビット水平型
バイオペラント	5 段	なし
キャッシュメモリ	データ 4 K 語 × 1 命令 1 K 語 × 1	データ、命令共通 4 K 語 × 1
ネットワーク	10 ビット / チャネル (ch) 3.84M バイト / 秒 · ch	10 ビット / チャネル (ch) 5M バイト / 秒 · ch
メモリボード	1 枚	4 枚
メモリ容量	80M バイト (16M ブロック)	80M バイト (16M ブロック)
使用素子	4 M ビット DRAM	1 M ビット DRAM

1 台当たりの性能はマルチ PSI の 2~3 倍となる。PE はメモリとプロセッサの 2 枚のボードから構成される(図 2)。メモリボードは 4 M ビット DRAM を使用することによって 80M バイトの容量を実現している。プロセッサボードは CPU 部とネットワーク部からなり、PU (Processing Unit), CU (Cache Unit), NU (Network Unit) の三つの CMOS セルベース LSI, 及び周辺回路から構成される。8 PE ごとに一つの PE には入出力用 SCSI バスの制御ボードを通して 670M バイトの内蔵磁気ディスクやフロントエンドプロセッサ (FEP) などが接続される。FEP には "MELCOM PSI-II" 若しくは "MELCOM PSI/UX" (4)(5) が用いられる。これらは大容量データベースや高度なマンマシンインターフェースを必要とする応用システムに活用される。

要素プロセッサの諸元をマルチ PSI のそれと比較して、表 1 に示す。

3.2 CPU 部⁽⁶⁾

CPU は論理型言語専用に開発された 2 種類の LSI (PU, CU) から構成され^(注 1)、5 段のバイオペラント、水平型マイクロプログラム制御、タグアーキテクチャ (8 ビットタグ + 32 ビットデータの 40 ビット / 語)、及びデータと命令にそれぞれのキャッシュメモリを配したハーバードアーキテクチャが特長となっている。KL1 言語の実行速度は後述の append プログラムで 530 K LIPS (Kilo Logical Inference Per Second) に達する。

5 段のバイオペラント(図 3)の各ステージの機能は以下のとおりである。

D : 命令キャッシュから読み出された機械語命令(後述)のデコードを行う。すなわち、機械語命令コードからオペラントのタイプ(レジスタ、即値、メモリアクセス)やオペラントに対する処理を決定する。
A : メモリアクセスオペラントのアドレス計算を行う。次命令若しくは分岐先命令のフェッチも行う。なお、分岐予測が失敗した場合や予測が不可能な場合には、E ステージの指示によって再フェッチが行われる。

R : A ステージで計算したアドレスに従って、必要であれば主記憶上のオペラントを読み出す。

S : E ステージでのマイクロ命令の読み出し、オペラントのセットアップ、及びデレファレンス機能がある。デレファレンスとは変数の参照ポインタの連鎖をたどることで、このステージには、R ステージで読み出したデータが更に他の変数セルへの参照ポインタである場合に、この連鎖の末端まで自動的にたどり続ける機能が備えられている。

(注 1) PU 及び CU は国からの成果移転を受け、当社の AI ワークステーション "MELCOM PSI/UX" の推論部にも使用されている。

E : 機械命令の主な処理を64ビット幅のマイクロプログラムで実行する。マイクロ命令は32K語の高速のマイクロ制御メモリに格納され、毎サイクル読み出されて実行される。図4はマイクロ命令フォーマットである。条件分岐、内部レジスタ間データ転送、レジスタ間演算、キャッシュアクセス、条件フラグ設定、内部カウンタ操作などが1命令で同時に実行できる。特に、タグ値などの判定から分岐完了までが2サイクルで実行されるため、KL1プログラムの実行のように動的な判定が頻繁な処理には大変有効である。

3.3 ネットワーク制御部

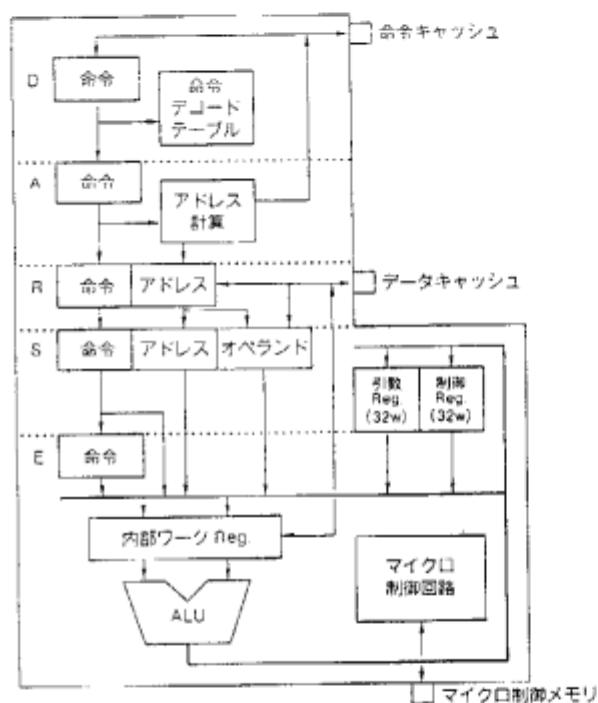


図3. 要素プロセッサの5段パイプライン

63	55 55	48 47	40 39	32				
*	CCR	DSTF	SC1F	SC2F	SD1F	SD2F	AS1F	AS2F
31	24 23	16 15	8 7	0				
ADF	ALF	JMPF	CNDF	*	EM1	EM2	EM3	

CCR<5> : キャッシュ/メモリ制御
DSTF<5> : 引数/制御レジスタ セット指定
SC1F<5> : 引数/制御レジスタ 選択指定(1)
SC2F<5> : 引数/制御レジスタ 選択指定(2)
SD1F<3> : 内部ワークレジスタ セット指定(1)
SD2F<2> : 内部ワークレジスタ セット指定(2)
AS1F<3> : ALU 左入力レジスタ指定
AS2F<2> : ALU 右入力レジスタ指定

ADF<3> : ALU 出力レジスタ指定
ALF<3> : ALU オペレーション制御
JMPF<3> : 分岐オペレーション制御
CNDF<4> : 分岐条件制御
EM1<6> : 単値/オペレーション修飾子(1)
EM2<6> : 単値/オペレーション修飾子(2)
EM3<4> : 単値/オペレーション修飾子(3)
*<計5> : 条件フラグ設定, カウンタ操作

図4. 要素プロセッサのマイクロ命令フォーマット

PE間を接続するネットワーク制御回路(NWC)はNUを中心に構成される。NWCには隣接4方向のPE、及び自CPUのための5組の双方向通信チャネルが接続されており、PE間メッセージの送受信及び転送を制御する(図5)。各チャネル幅は10ビットで、9本のデータ線、1本のbusy線(逆方向)からなる。各チャネルのバンド幅は3.84Mバイト/秒である。

図6はメッセージパケットの形式を表す。データ線の上位の2ビットにより、メッセージパケットの先頭バイト(以下“ヘッダ”という)と最後尾バイト(以下“テール”という)が判別される。NWCはヘッダにあるて先PE番号を読み取り、内部のルーティング用のテーブルを参照することによって転送すべき方向を決定する。

そして、その方向のチャネルが“busy”でなければそのチャネルにメッセージパケットを転送する。この時、そのチャネルはメッセージテールの転送完了まで保持される(“wormhole”ルーティングという)。さて先PE番号が自PEである場合はNWC内のリードバッファに取り込み、テールまでの取り込みが完了するとCPUに割込みを上げ、メッセージの解釈を要求する。CPUではKL1プログラムの処理の適当な区切りでこの割込みを処理する。CPUのメッセージ送信はNWCのライトバッファに一つの完全なメッセージを書き込み終わることにより、自動的に開始される。

このように、NWCはCPUから出入りするメッセージの操作を行うだけでなく、通り過ぎるメッセージの転送もCPUに負担をかけずに自動的に行うため、遠距離にあるPE間の通信を円滑に行うことができる。

4. 並列論理型言語 KL1

大規模な知識情報処理システムを実現するには、プログラムの生産性や保守性が非常に重要となる。並列論理型言語

KL1 (Kernel Language version 1)

(図7)は、並列でかつ複雑な問題を効率良くプログラミングできることをねらった高レベル並列言語で、PIMの核言語(ハードウェアとソフトウェアとのインタフェース言語)としてICOTで開発された。並列論理型言語 Flat GHC (Flat Guarded Horn Clauses)^[7]にシステム記述やスケジューリング及び負荷分散指定のための機能を付加したものである。

KL1は並列言語として次の特長をもっている。

(1) 同期は、未束縛変数への値の代入を待ち合わせる機能で実現される。

(2) 通信は、変数値の読み出し時に言語

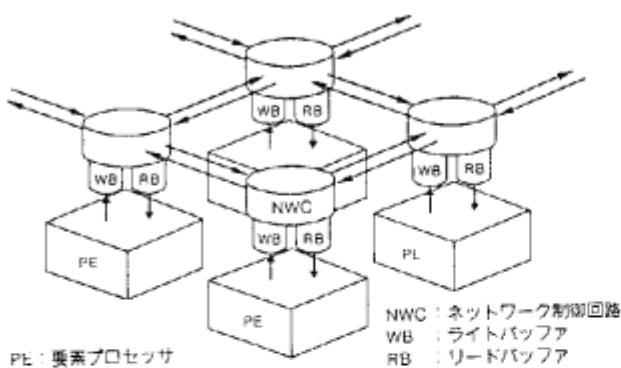


図5. ネットワーク制御回路(NWC)による
プロセッサ間結合

機能として自動的に行われる。

- (3) 論理変数の特長である单一入規則により、いったん定まった変数値は変化しないことから、他PEへの値のコピーが許される。したがって、2回目以降の変数値の読出しにはPE間通信が省略できる。
 - (4) 負荷分散は“プログラマ”と称するKL1述語呼出し時の附加情報で行われるため、プログラムの論理的な意味を変えずに負荷分散を指定・変更できる。
 - (5) メモリ領域の割付けや再利用も言語機能として自動的に行われる。これらにより、従来の手続き型言語に並列機構と同期機構を加えたようなものに比べ、同期のタイミングバグなどに悩まされることが非常に少なくなる。

5 KI-1 實行執行方式

KL 1 プログラムは Prolog での WAM⁽⁸⁾ に相当する抽象機械語命令にコンパイルされ⁽⁹⁾、PE のマイクロプログラムによって直接解釈実行される。図 8 はリストの結合を行う KL 1 プログラムの抽象機械語命令へのコンパイル例である。switch 系、wait 系及び read 系の命令が変数の待ち合わせと値のチェックを行い、get 系及び write 系の命令が主に未定義変数への値の代入を行う。put 系命令は変数値を引数レジスタに用意する。

他 PE への KL1 ゴールの送出にも専用機械語命令が用意されており、マイクロプログラムで解釈実行される。なお、送出先 PE での KL1 ゴールの実行に必要なプログラムコード、データなどは必要になった時点で自動的にコピーが送られ、不要になったものは廃棄される。そして、空いたメモリ領域は再利用される。ゴール実行のための資源管理やゴールの終了報告なども言語表層以下の言語処理マイクロプログラムによって PE 間通信を用いて行われる⁶⁰。

以上のように、メモリ管理やプロセス管理のような従来のオペレーティングシステムの基本機能を並列に制御する部分が、言語機能の一部としてマイクロプログラムによって効率良く行われ、その上位に位置するPIMのオペレーティング

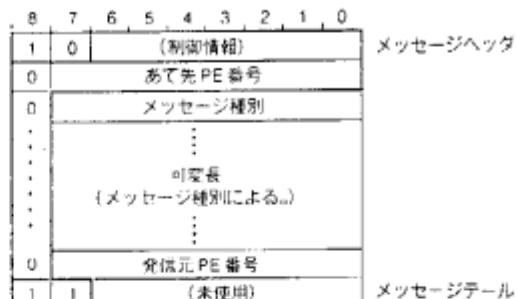


図6. メッセージパケットのフォーマット

(透視 coat の定義)

`goal(X, Y) :- X > 0 | goal1(X, Y), goal2(Y, Z). . . . クローズ1`
`goal(X, Y) :- X <= 0 | goal3(X, Y), goal4(Y, Z). . . . クローズ2`

ガード部 来テイニ部

並列論理型言語 KL1 によるプログラムは複数のクローズからなる述語定義の集まりからなる。 \mid (bar) の左側の条件(ガード部)が最初に成立したクローズが選ばれ、他のクローズの実行は放棄される。値が決まっていない変数の条件を調べようとして、クローズが選択できなかった場合はその述語呼び出しはサスペンドされる。選ばれたクローズのボディ一部ゴール(goal_1 と goal_2 、又は goal_3 と goal_4) は並列に実行され、ゴール間の共有変数(上の例では Y)によって通信や同期が行われる。ゴールにはユーザー定義(他の述語を呼び出す)、システム定義(算術演算を行うようなあらかじめシステムが提供している述語を呼び出す)、及びボディユニフィケーション("X=t(Y)"などのように表記し、変数への値の代入に用いられる)がある。

圖 7. 幫列論理型言語 L₁

システム(PIMOS^⑩という。)やユーザープログラムへの負担が軽くなるよう考慮されている。

6 たすけ

効率の良い並列実行のためにはどのようなタイプの問題に
対しても、以下のようなことが必要となる。

- (1) 並列度が高く、通信が1か所に集中しないように分散したアルゴリズム
 - (2) どのプロセッサもできるだけ遊ばせないようにする負荷分散
 - (3) (探索型の問題に対しては) 無駄になってしまふような自己計算を最低限で控えますスケジューリング

並列アルゴリズム、負荷分散などの問題のほかにも並列プログラミング言語の設計・改良、デバッグ環境の充実、故障に対する耐久性の増大など、並列処理技術の行く手には難しい課題が山積しており、並列ハードウェアの研究と足並みをそろえて一歩一歩進めていく必要がある。

現在 PIM/m 上では、KL1 言語処理マイクロプログラム

```

append([X | L1], L2, L) :-  

    true | L = [X | L0], append(L1, L2, L0).  

append([], L2, L) :-  

    true | L = L2.

```

(a) KL1 プログラム例

```

%%% 3つの引数が引数レジスタ R1～R3にセットされて、app//3が呼出される  

app//3 : try_me_else app//s  

    % 中断または失敗時の飛先を割り当てる  

    % APにセット  

    switch_on_error_list R1, app//n % R1が Listなら次、それ以外なら APに分岐  

    % ただし、未定義なら APに分岐  

    read_variable R4 % 入力 Listの CAR要素を R4に取り出す  

    read_variable R5 % 入力 Listの CDR要素を R5に取り出す  

    put_reused_structure R1, R2, 2 % 入力 Listの cell回収が可能なら再利用  

    % 回収不可なら新しいList cellを割り当てる  

    write_value R4 % 出力 Listの CARに R4を書き込む  

    write_variable R4 % 新しい変数 cellを割り当てる、それへのボ  

    % インタを出力 Listの CDRと R4に書き込む  

    get_bound_value R3, R1 % R3と R1(出力 List)をユニファイする  

    put_value R1, R5 % R1に R5をセット  

    put_value R3, R4 % R3に R4をセット  

    execute app//3 % app//3を再帰呼び出し  

app//n : wait_nil R1 % NILなら次、それ以外なら APに分岐  

    get_value R2, R3 % R2と R3をユニファイ  

    proceed % このクローズは終了し、他のコール  

    % を実行  

app//s : suspend app//3 % 中断か失敗かの判定とそれぞれの処理

```

(b) KL1B 命令列へのコンパイル例

図 8. リストを結合する KL1 プログラムのコンパイル例

の評価・改良のほか、ユーザー管理・ファイル管理・資源管理などを行いプログラム環境を提供する並列オペレーティングシステム PIMOS の機能拡張などを ICOT と協力して進めている。また、並列応用プログラムとしては、遺伝子情報処理の一環としてのたん(蟹)白質構造解析プログラム、判例を用いた法的推論システム、LSI-CAD システムの一環としての論理シミュレータ・セル配置プログラム・配線プログラム、囲碁システム、定理証明システムなどの開発が行われており、より実用的な並列処理システムに向けての研究が進められている。PIM//m 上でのこれらの成果が将来の並列計算機の進歩への大きな足掛かりとなることを期待したい。

最後に、PIM//m の研究開発に際して御指導をいただいた ICOT 研究部内田俊一郎長を始めとする関係各位に深く謝意を表す次第である。

参考文献

- Goto, A., Sato, M., Nakajima, K., Taki, K., Matsumoto, A. : Overview of the Parallel Inference Machine Architecture (PIM), Proc. of the International Conference on Fifth Generation Computer Systems, 208～

229 (1988)

- Taki, K. : The Parallel Software Research and Development Tool : Multi-PSI System, Programming of Future Generation Computers, Elsevier Science Publishers B.V. (North-Holland) (1988)
- 益田嘉直、中川智明、岩山洋明、武田保孝、中島浩、瀧和男：財新世代コンピュータ技術開発機構向けマルチ PSI システム、三菱電機技報, 62, No. 12, 1100～1105 (1988)
- 湯浅准介、上田尚純、松本明：三菱 AI ワークステーション《MELCOM PSI/UX シリーズ》の概要と特長、三菱電機技報, 65, No. 6, 594～601 (1991)
- 益田嘉直、田辺隆司、池田守宏、深沢雄、岩山洋明、中島浩：三菱 AI ワークステーション《MELCOM PSI/UX シリーズ》のハードウェアシステム、三菱電機技報, 65, No. 6, 602～607 (1991)
- Nakashima, H., Takeda, Y., Nakajima, K., Andou, H., Furutani, K. : A Pipelined Microprocessor for Logic Programming Languages, Proc. of the International Conference on Computer Design, 355～359 (1990)
- Ueda, K. : Guarded Horn Clauses : A Parallel Logic Programming Language with the Concept of a Guard, Technical Report TR-208, ICOT (1986)
- Warren, D. H. D. : An Abstract Prolog Instruction Set, Technical Note 309, Artificial Intelligence Center, SRI (1983)
- Kimura, Y., Chikayama, T. : An Abstract KL1 Machine and its Instruction Set, Proc. of the Symposium on Logic Programming, 468～477 (1987)
- Nakajima, K., Inamura, Y., Ichiyoshi, N., Rokusawa, K., Chikayama, T. : Distributed Implementation of KL1 on the Multi-PSI/V2, Proc. of the Sixth International Conference on Logic Programming, 436～451 (1989)
- Chikayama, T., Sato, H., Miyazaki, T. : Overview of the Parallel Inference Machine Operating System (PIMOS), Proc. of the International Conference on Fifth Generation Computer Systems, 230～251 (1988)