

ICOT Technical Memorandum: TM-0675

TM-0675

並列推論マシンPIM

市吉伸行

January, 1989

©1989, ICOT

ICOT

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

(03) 456-3191~5
Telex ICOT J32964

Institute for New Generation Computer Technology

並列推論マシン PIM

市吉 伸行

(財) 新世代コンピュータ技術開発機構 第4研究室

1 はじめに

第5世代プロジェクトは、1990年代の大規模知識処理システムのベースとなる第5世代コンピュータの開発を目指しているが、その特徴は、並列マシンと知識処理システムを結ぶ層として論理型の核言語を設定し、そこから下層および上層を設計して行くインサイド・アウトのアプローチを取っているところにある。このアプローチをとる大前提として、(1) 大規模知識処理システムの要求する計算能力に応えるためには並列実行が必須であること、(2) 複雑な知識処理システムを見通し良く設計するためには明確な意味論を持つ論理プログラミングが有効であること、(3) 論理プログラムは並列処理に向いていること、という認識がある。

本稿では、第5世代コンピュータシステムのエンジンとなる並列推論マシン(ハードウェア+言語処理系)とその上の並列オペレーティングシステムPIMOSについて解説する。初めに並列論理型言語KL1の設計思想および言語仕様を述べ、次にPIMのハードウェアとその上のKL1処理系の実現方式を概説し、PIMOSについては、その機能・特徴とOS保護の手法を紹介する。最後に、研究開発の現状と課題について触れる。

2 並列論理型言語 KL1

並列論理型言語KL1は、ストリームAND並列型言語GHC^[8]に制限を加えたFlat GHC(FGHC)をベースにして、OSを支援し、効率的並列実行を可能にするための実行管理、資源管理、優先順位および負荷分散指定の諸機能を追加したものである。

2.1 Flat GHC

Flat GHC(FGHC)はGHCに限られた組込み述語のみの出現をガード部に許すような制限を加えたも

のである。FGHCではガードの実行がネストせず、AND/OR木や変数環境の管理が不要なので、オーバヘッドの小さい処理系が実現できる。ガードにユーザ述語が書けないことにより表面的記述力は落ちるが、ガードにおけるテストをボディ・ゴールの組合せに変換できるので、本質的記述力は低下しない。

2.2 莊園

莊園とは、実行制御と資源管理を行なうためのメタプログラミング機能である。末尾再帰呼出しされるKL1のゴールが小粒度プロセスを表わすのに対し、莊園は1つのまとまった計算を定義し、莊園単位に実行の監視や制御が行なえるようになっている。莊園は組込み述語executeの呼出しによって生成される。

execute(Goal, Control, Report, Mask)

ここで、GoalはKL1データとして表現されたKL1ゴール、Controlは実行制御用のコントロールストリーム、Reportは実行監視用のレポートストリーム、Maskはこの莊園で捕捉する例外事象を指定するマスクである。生成された莊園の下でGoal(およびその子孫ゴール)が実行される。

莊園の下で新たにexecuteを呼出すことによりネストした莊園(子莊園と言う)を生成することもできる。

コントロールストリームにstart, stop, abortなどのコマンドを流すことにより、実行の開始(再開)、中断、放棄などができる¹。レポートストリームには、莊園内で起きた特定の事象を報告するメッセージが流れれる。莊園の終了(全てのゴールおよび子莊園が終了したこと)の報告、例外事象の発生の報告などである。

¹ KL1におけるストリームは、メッセージを要素とするリストとして表わされる。メッセージの送り手と受け手の間の同期はKL1の同期機構によって自然に実現される。

莊園は資源管理の単位でもある。資源とは計算量の目安であり、計算時間やメモリ消費量と大体の比例関係にある。莊園があらかじめ指定された量の資源を消費してしまうと実行中断状態になり、資源の不足がレポートストリームに報告される。莊園の実行管理プログラムは、必要に応じて、コントロールストリーム・コマンドにより消費できる資源の上限を増やすことができる。資源管理機構により、プログラムの暴走を防いだり、莊園単位のファアなスケジューリングが可能になっている。

KL1では、例外事象は、例外の種類によって決まるタグ値とマッチするマスク値を持っている一番内側の莊園のレポートストリームに次の形で報告される。

```
exception(Type, Goal, NewGoal)
```

ここで、Typeは例外の種類を表わすデータ、Goalは例外が発生したゴール、NewGoalは未束縛変数である。莊園の実行管理プログラムは、上記メッセージを受け取るとTypeおよびGoalを解析し、NewGoalを適切なゴールで具体化することによって実行を継続したり、場合によってはabortコマンドで単に実行を放棄することもできる。

例外事象には、ゴールの失敗(全ての節のガードが失敗した)、ボディ・ユニフィケーションの失敗、組込み述語例外(不正な人力引数、算術演算例外など)、raise述語の呼出しによるユーザ定義例外などがある。これらはユーザ定義例外を除けば全てFGHCプログラムにおける失敗である。

例外処理機能は、OSサービス提供やKL1の言語機能拡張に使うことができる。例えば、整数の3と構造型データ1+2のユニフィケーションの失敗時に再開ゴールとしてtrueを与えることで、(オーバヘッドはかなり大きいが)ユニフィケーションのセマンティクスを拡張できたりする。

2.3 優先度指定・負荷分散指定

優先順位および負荷分散は、ボディ・ゴールにプラグマと呼ばれる次のようなアノテーションを付けることによって指定する。

```
Goal@priority(Prio)
```

```
Goal@processor(Proc)
```

前者はGoalを優先度Prioで実行せよという指定であり、後者はGoalをプロセッサProcで実行せよという指

定である。プラグマはプログラムの論理的意味を変えない。

莊園毎に実行優先度の上下限が決まっており、実行優先度は、莊園の優先度幅の中の割合で指定したり(絶対指定)、親ゴールの優先度と上下限の間の割合で指定したり(相対指定)できる。物理的には 2^{12} 程度の細かい優先度レベルを提供するが、応用プログラムの実行効率を上げるためにこの細かい優先度指定が役立つことが示されている。また、ユーザ・プログラムの暴走のためにOSプロセスがなかなかスケジュールされないといった状況は、OSプロセスの優先度をユーザ・プログラムよりも高くすることで防げる。

MIMD型並列マシンの性能を十分に引き出すためには、実行負荷をよく分散しなければならない。しかし、単純な負荷の均等化はプロセッサ間通信オーバヘッドを大きくするし、KL1プログラムが対象とする非定型な問題では計算負荷の事前の予想が難しい、という困難がある。種々の負荷分散手法を試すために、現在のところ上記のように低レベルの指定法を提供している。

3 並列推論マシンのハードウェア

KL1を実行する並列マシンがParallel Inference Machine(PIM)[3]である。現在、ハードウェア詳細仕様の少しずつ異なるいくつかのPIMを設計開発中だが、それらはKL1-B[6]と呼ばれる共通の中間言語をサポートする。ここでは、その内のひとつであるPIM/pを紹介する。

PIM/pのアーキテクチャは図1に示すように階層的である。すなわち、下位のレベルとして8台の要素プロセッサがメモリを共有するクラスタがあり、上位のレベルとしてクラスタたちが高速ネットワーク(ハイペリュープ)で結合され、全体として当初128プロセッサの並列マシンが構成されるが、将来的には1000台程度の規模まで拡張可能となっている。

メモリ共有型アーキテクチャは共有バスがボトルネックになるため、要素プロセッサ10数台程度が限界と見られ、大規模並列マシンを実現できない。また、大規模ネットワークにおいては、配線が可能であるためには、何らかの局所性を持つ(つまり、プロセッサ間の距離に遠近がある)ような結合方式とならざるを得ない。これが、上位構造が局所性を持つネットワーク型になっている理由である。

一方、プロセッサ間の通信はメッセージ通信よりも

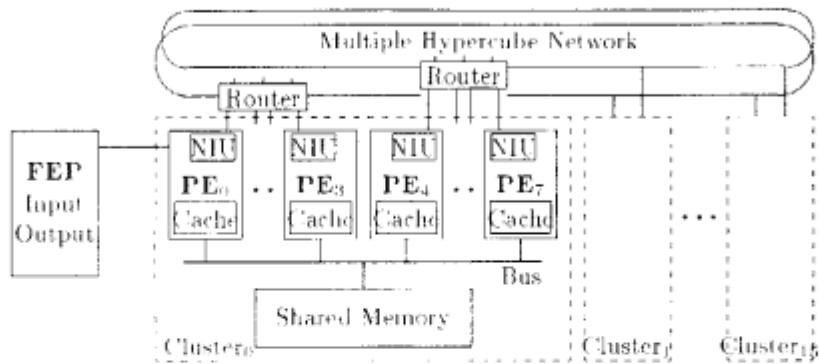


Figure 1: PIM/p のアーキテクチャ

共有メモリを介した通信の方がコストが安い。PIM の典型的なアプリケーションは、多くのプロセスが部分を受け持つて全体として 1 つの問題を解くような通信の多い型のものが予想されるので、台数性能比を上げるために通信コストはできるだけ抑えたい。下位レベルに密結合のクラスタを採用した理由は、相互間通信の多いプロセスを同一のクラスタに割当てることによって、通信オーバヘッドを低減できることにある。

プロセッサは並列記号処理向きに、タグ処理機能や後述の MRB の支援機能を持ったものを新規設計した。

クラスタ内では、分散したキャッシュ間の一貫性を保つこと、競合資源のロックを低いオーバヘッドで実現することが必要である。そのために、KL1 プログラムの実行特性に合わせたバス・プロトコルを設計した。例えば、プロデューサの書いたデータを別のプロセッサ上のコンシューマが読む場合にキャッシュ間データ移動が起るが、KL1 ではプロデューサがそのデータに再びアクセスしないので、移動元のキャッシュのデータを無効化することができる。また、共有されていないことの分かるデータはバス・サイクルを使わずにロックできる。

ネットワーク接続装置 (NIU) には、バイト / ワード変換とパケット・バッファリング機能があり、汎用入出力インターフェースも接続されている。

なお、クラスタ内の 6 プロセッサにフロント・エンド・プロセッサ (FEP) やディスク装置などの入出力機器が直結可能となっている。PIM 上のアプリケーションの多くは計算バウンドだと思われるが、自然言語処理などで、大きな (PIM の大きな実メモリにも入り切らないような) データベースを必要とする可能性があるので入出力バンド幅は大きめに取る予定である。

4 KL1 処理系

PIM の上で KL1 プログラムの実行を司るのが KL1 処理系である。KL1 处理系の開発において、次のような問題点があった。

- 分散したゴールの管理

プロセッサをまたがって実行される並列の実現、実行制御、終了検出。

- 分散したデータの管理

分散ユニファイケーションの実現、クラスタをまたがる参照ループ生成の防止。

- 効率のよいガーベージ・コレクタ

KL1 は、破壊的書き込みの許される言語やバックトラックによりメモリが自動的に開放される Prolog と比べてメモリ消費速度が速いので、トータルな実行性能を高く保つためには効率の良いガーベージ・コレクタ (GC) が必須。クラスタをまたがる GC も必要。

KL1 の基本実行モデルを記述した後、個々の問題点への解決策を紹介する。

4.1 基本実行モデル

我々の処理系は以下のようないくつかの実行モデルに従っている。AND 関係で結ばれたゴールの集まりであるゴール・ブルーとプロセッサの集合があり、各プロセッサはゴール・ブルーからゴールを取り出して、リダクションを試みる。プロセッサの数が増えると、同時に多数のゴールのリダクションが行なえるので、処理系のスループットが上がることになる。

これが基本実行メカニズムであるが、実際には、KL1ではゴール・プールは莊園に対応し、ゴールの失敗は例外の形でゴール・プールの外に伝えられる。

KL1プログラムは、プログラム・モジュール単位に WAM ライク²な中間言語 KL1-B にコンパイルされる。コンバイラはレジスタ・レベルの最適化やコミットできる節を高速に見つけるクローズ・インデキシングを行なっている。

4.2 莊園里親方式

クラスタをまたがった莊園においては、分散実行制御や分散終了検出が大きな問題となる。莊園には、実行可能状態や中断状態があり、もしその情報を集中化しているとゴールがスケジュールされる度にそのゴールの莊園を管理しているクラスタに状態を問い合わせねばならず、オーバヘッドが大きい。そこで我々の処理系では莊園の管理情報をクラスタ間に分散し、それを里親と呼んでいる[4]。

分散した莊園における実行終了の検出はいわゆる分散計算終了検出の問題であるが、重みつきレファレンス・カウントを応用した WTC 方式[7]により終了検出に必要なメッセージ量を低く抑えることに成功した。

4.3 分散データ管理と分散ユニフィケーション

分散処理系で全体を 1 つのアドレス空間とすると、クラスタ毎の独立した GC が難しくなるという難点がある。クラスタ内 GC によりデータのアドレスが変わると³、そのデータを参照している全てのポインタを更新しなければならないが、独立して動いている他のクラスタ内のポインタの更新は非常に難しい。

そこで KL1 処理系では、クラスタ内外でアドレス系を分け、各クラスタはクラスタ外からクラスタ内への参照を管理する輸出表と、クラスタ内からクラスタ外への参照を管理する輸入表という 2 つのアドレス変換テーブルを管理することとした。

データがクラスタ間に分散していると、ユニフィケーションにおいてクラスタ外のデータを読みだり、書き込んだりする場合が出てくるが、そのための read プロトコルと unify プロトコルを設けた。また、分散ユニフィケーションを無秩序に行なうと外部参照ボ

²WAM は、現在広く採用されている Prolog の抽象マシン語[9]である。

³データのアドレスを動かさないマーク・アンド・スワープ型の GC は処理時間やフッゲメンテーションの欠陥がある。

インタからなるループができる可能性があるので、クラスタ番号を比較して番号の大きいクラスタを指す外部参照は自クラスタの変数にバインドしてはいけない（その代わりに unify メッセージを出す）という規則を設けた⁴。

4.4 ガーベジ・コレクション

KL1 処理系における GC の重要性から、クラスタ内・クラスタ間の即時および一括型の GC が開発された。

MRB 方式[1]は、参照ポインタに Multiple Reference Bit (MRB) と呼ばれる多重参照情報を持たせて、MRB=OFF のポインタが單一参照であることを保証するようにしたものである。KL1 プログラムにおいては、大半のデータは單一参照されているが、MRB=OFF のポインタが消費されればそのデータの占めていたメモリ領域は回収できる。ポインタ側に MRB を持たせたために、回収時以外にはポインタの先にアクセスしなくて済むというメリットがある。これは、クラスタにおいてロックの回数を増やさずに済むことにつながる。MRB 情報は、ベクタ要素の更新やストリーム・マージの最適化およびデータ輸出入管理の簡略化にも利用している。

分散処理系では、クラスタ外から参照されているデータの回収が問題となるが、重みつき参照カウント原理を用いた WEC 方式[5]によって即時 GC を実現した。

MRB 方式と WEC 方式で回収できないデータがあるため、クラスタ内並列 GC[10] および一括型大域 GC を検討中である。

5 並列オペレーティングシステム PIMOS

並列オペレーティングシステム PIMOS[11, 2]は、並列推論マシン向けの並行ユーザ、マルチタスクの OS である。

オペレーティング・システムの目的を一言で言えば、計算機の持つ物理的な機能をユーザに使い易い形で提供することである。計算機の物理的機能は大雑把に言えば、計算能力と入出力機能だが、それをユーザに使い易い形で提供するとは、一つには計算機資源を明快なモデルの形に抽象化することであり、もう一つにはユーザが、故意または過失により、そのモデルに

⁴ 実際には間接的輸出の関係でもう少し複雑な規則となっている[8]。

違反した時に、計算機資源および違反していない部分を保護することである。並列オペレーティングシステム PIMOS は次のような基本方針で設計された。

- KL1 で記述された純粋な論理型 OS であること。

これには、OS を記述するレベルで既に計算機資源がユニフォームに抽象化されているというメリット(例えば、メモリ保護の問題が生じない)がある。KL1 は外部と相互作用するプロセスの記述においており、PIMOS では計算機の資源の多くを KL1 プロセスとしてユーザーに見せている。入出力デバイスはその典型である⁵。ユーザプログラムと OS はオブジェクトとメタの関係にあり、理想的には、OS はユーザプログラムを 1 レベル下のデータとして制御すべきだが、これは実行効率上の問題があり、現時点では実用的でない。そこで、前述のように KL1 では FGHC にメタプログラミング機能を組込んでいる。

- 集中型単一 OS であること。

PIMOS はプロセサ毎に独立して動作する OS の集合体ではなく、全体として一体であるシステムを並列実行可能な部分を並列実行する。非定型並列処理においては物理プロセッサを意識しない方が、自然にプログラムが書けるのである。ただし、PIM が上位層で疎結合マシンであることを考慮して、強い結びつきのあるプロセスたちを同じクラスタにマッピングするなどの措置は講じている(PIMOS の論理的構造への影響はない)。

PIMOS の具体的機能としては、タスクの生成・監視・制御、入出力デバイスの管理、シェル機能、などがある。ユーザ・プログラムはソフト例外を上げる(raise)ことにより、PIMOS への要求ストリームを獲得し、このストリームにメッセージを流すことによって必要な PIMOS 資源を得るようになっている。全ての PIMOS 資源はその資源特有のプロトコルを持ったストリームとしてユーザーに提供される。

メタプログラミング機能を「理想的に」実現していないために、PIMOS では OS をユーザから保護する必要がある。下の例を見よう。ユーザプログラムが N 文字を読むために getb というコマンドをファイルに送ったとする。

⁵これに対して Prolog において入出力機能がいかに Prolog のセマンティクスと合わない形で入っているかを思い起させたい。

```
..., Req = [getb(N,Str)|NewReq], ...
```

ここで、Req はファイルへの要求ストリームである。PIMOS のファイルハンドラ側には getb に対応して、次のような節がある。

```
file_handler([getb(N,Str)|NewReq],File) :-  
    true !,  
    readFromFile(File,N,StrRead,NewFile),  
    Str = StrRead,  
    file_handler(NewReq,NewFile).
```

もしここで、文字列が読み込まれる前にユーザ・プログラムが Str に例えば整数の 0 をユニファイしてしまったとすると、ファイルハンドラで読み込まれた文字列 StrRead と Str とのユニフィケーションが失敗してしまう。また、ユーザプログラムが文字数の N を具体化しないまま終了してしまうと、ファイルハンドラ側に終了しないゴルフができてしまう。

PIMOS では保護フィルタというインタフェース・ルーチンを導入することでこれを解決した。保護フィルタはユーザが PIMOS 資源へのストリームを獲得する度に、ユーザプログラムと PIMOS 資源の間に自動的に挿入されるフィルタ・プロセスであり、次の機能を持つ。

- ユーザプログラムからのメッセージで PIMOS 側で読む部分については、それが具体化されるのを待ち、メッセージとして正しい形をしていることを確認した上で、PIMOS 側に流す。
- ユーザプログラムからのメッセージで PIMOS 側で書く部分については、新たに未束縛変数に置き換えて PIMOS 側に渡し、PIMOS 側で具体化した後に元の変数とユニファイする。

保護フィルタの導入によって、ユニフィケーションの失敗やデッドロックはユーザタスク内で起きることになり、PIMOS は安全である。

また、ユーザタスクの強制終了時に PIMOS へのストリームを閉じるバルブ機構が用意されている。

6 現状と課題

PIM と PIMOS、およびそれらを結ぶ KL1 について概説してきた。現在、KL1 処理系がマルチ PSI⁶ に

⁶逐次型推論マシン PSI をメッシュ状のネットワークで結合した疎結合型並列マシン。マルチ PSI の要素プロセッサは PIM/p のク

実装されており、PIMOS暫定版が稼働している。また、1台の小型化 PSI 上でマルチ PSI をシミュレートする Pseudo マルチ PSI システムも開発されている。両システムとも徐々に ICOT 外の研究機関にリリースして並列プログラミング研究に役立てて頂く予定である。KL1 处理系が PIM/p に実装されるのは来年度半ば頃となろう。

現在の PIMOS では、入出力機能が全て FEP において実現されており、コンバイラなどの開発環境も PSI 上のクロスシステムに置かれているが、それらの機能を PIMOS 本体に移行していくつもりである。

本来 OS が管理すべき計算機資源の内、クラスタ毎のメモリ使用量などは、現在 KL1 の下に隠れてしまつており、きれいな形でそれらを言語レベルに持ち上げて管理の対象にすることは、論理 OS としての PIMOS の今後の課題のひとつであろう。

負荷分散の研究はこれから分野だが、研究の進展にともない将来的には、実験的な動的負荷分散機能を PIM システムに組込んで行きたい。

参考文献

- [1] T. Chikayama and Y. Kimura. Multiple reference management in Flat GHC. In *Proceedings of the Fourth International Conference on Logic Programming*, pages 276-293, 1987.
- [2] T. Chikayama, H. Sato, and T. Miyazaki. Overview of the parallel inference machine operating system (PIMOS). In *Proceedings of the International Conference on Fifth Generation Computer Systems 1988*, pages 230-251.
- [3] A. Goto, M. Sato, K. Nakajima, K. Taki, and A. Matsumoto. Overview of the parallel inference machine (PIM) architecture. In *Proceedings of the International Conference on Fifth Generation Computer Systems 1988*, pages 208-229.
- [4] N. Ichiyoshi, T. Miyazaki, and K. Taki. A distributed implementation of Flat GHC on the Multi-PSI. In *Proceedings of the Fourth International Conference on Logic Programming*, pages 257-275, 1987.
- [5] N. Ichiyoshi, K. Rokusawa, K. Nakajima, and Y. Inamura. A new external reference management and distributed unification for KL1. In *Proceedings of the International Conference on Fifth Generation Computer Systems 1988*, pages 904-913.
- [6] Y. Kimura and T. Chikayama. An abstract KL1 machine and its instruction set. In *Proceedings of the 1987 Symposium on Logic Programming*, pages 468-477, 1987.
- [7] K. Rokusawa, N. Ichiyoshi, T. Chikayama, and H. Nakashima. An efficient termination detection and abortion algorithm for distributed processing systems. In *Proceedings of the 1988 International Conference on Parallel Processing, Vol. 1 Architecture*, pages 18-22, 1988.
- [8] K. Ueda. Guarded Horn Clauses: A Parallel Logic Programming Language with the Concept of a Guard. Technical Report TR 208, ICOT, 1986.
- [9] D.H.D. Warren. An abstract Prolog instruction set. Technical Note 309, SRI International, 1983.
- [10] 佐藤正俊,後藤厚宏. KL1並列処理系の評価 - メモリ消費特性とGC. 並列処理シンポジウム JSPP'88 論文集.
- [11] 宮崎敏彦. 並列論理型言語 KL1 の実現方式と並列 OS の記述. 電子情報通信学会論文誌 D Vol. J71-D, No. 8, pages 1423-1432, 1988.

著者紹介

市吉 伸行

1979年東京大学理学部情報科学科卒業。1981年同大学院修士課程修了。1982年(株)三菱総合研究所入社。1984年より1年間米国バトル研究所にてAIを研修。1987年より(財)新世代コンピュータ技術開発機構へ出向。論理型言語処理系、並列プログラミングの研究開発に従事。情報処理学会、AAAI会員。

クラスターに相当すると考えてよい。FEP はシステム全体で最大 4 台まで接続可能となっている。