

並列推論マシンに於ける負荷分散方式の評価

An Evaluation of the Load Balancing Mechanism for a Parallel Inference Machine

山内 宗、 田中 英彦

Tsukasa YAMAUCHI* and Hidehiko TANAKA**

*日本電気株式会社
NEC Corp.

**東京大学工学部
University of Tokyo

あらまし 我々は、並列推論マシンで、Committed choice言語を実行する際の負荷分散方式についての評価・検討を行うために、ソフトウェア・シミュレーションを行った。シミュレーションでは、単一化プロセッサとメモリで構成される推論ユニットが2種の多段ネットワークによって多数並列に接続されたモデルを仮定した。本稿では、シミュレーション・モデルと処理方式、負荷分散方式について述べた後、シミュレーション結果について報告する。

Abstract There are many important problems in executing committed choice languages in a Parallel Inference Engine PIE. One of them is the load balancing problem. We made a simulation to evaluate more effective load balancing mechanism. In this paper, first we describe the simulation model, the execution mechanism and the load balancing mechanism, then we present the simulation results.

1. はじめに

現在我々は、高並列推論エンジンPIEの開発を進めている。PIEは、ゴール書換えモデルに基づき、論理型言語を高並列に実行する推論マシンである。現在迄のところ、PIEに於ける並列実行は、OR並列を主体に検討が進められてきた。従って、従来の処理方式のままでは、GHC、FLENG等のCommitted choice言語を効率良く実行することは、難しい。並列推論マシン上で、Committed choice言語を実行する際に問題となるのは、次の2点である。

①共有変数のアクセス競合と負荷分散

②実行制御とスケジューリング

ここでは、まず上記の①について、ソフトウェア・シミュレーションにより、評価、検討を行うことにする。

具体的には、

- ・負荷分散の方式
- ・相互結合網の構成

についての評価検討を行う。

シミュレーションの対象とする言語は、GHCとよく似たCommitted choice言語であるFLENGを用いることにする。

2. 論理型言語 FLENG

FLENGは、GHCと非常によく似た言語であり、次に挙げる点がGHCとの主な相違点である。

①ガード部が存在せず、定義節のヘッド・リテラルの単一化が成功すれば、すぐにコミットされる。

②ボディ部の各リテラルは、論理的AND関係にないので、どれか一つのゴールが失敗しても、他のゴールの実行には影響を与えず、全てのゴールは必ず実行される。ゴール間にAND関係を持たせたい場合には、プログラム中に明示的に書く必要がある。

また、GHCと同様にFLENGに於いてもヘッド・リテラルの単一化の際には、親ゴール中の変数を具体化することはできず、その場合は単一化はサスペンドする。

この様に、共有変数の扱いに関してはFLENGはGHCにかなり似ているが、実行制御に関する部分はかなり簡略化されている。この様にFLENGはGHCに比べて低レベルの言語であるが、ほとんどのGHCプログラムは容易にFLENGプログラムにコンパイルすることが可能であり、また、そのための処理系もすでに完成している。

3. PIEのシミュレーション・モデル

今回用いたPIEのシミュレーション・モデルを図1、2に示す。図1は全体構成を示し、推論ユニット(IU)が相互結合網によって並列に接続されて、PIEを構成している。図2は、各IUの内部構成を示し、一つの推論ユニットは、単一化の処理を行う単一化プロセッサ(UP)と定義節やゴール・フレーム(GF)を格納するメモリ・モジュール(MM)、共有変数等を格納する共有メモリ(SM)で構成されている。そして、UPは自分と対になっているMM、SMへはネットワークを介さずに高速にアクセスできる。

各IUは基本処理単位であるゴール・フレームをやり取りするための分配網(DN)と、共有変数の値、構造データ等をやり取りしたり、サスペンドしたゴール・フレームのアクティベ

イト等の制御を行うための共有メモリ網(SMN)の2つの相互結合網で結合されている。共有メモリへのアクセスとゴール分配の際の通信特性が大きく異なっていることが、二種のネットワークに分けた理由である。

SMNとDNは、どちらも4×4のクロスバー・スイッチング・ユニット(SU)を用いた多段ネットワーク(オメガ網)である。特に、DNは、自動負荷分散機能を持った、ネットワークとなっている。

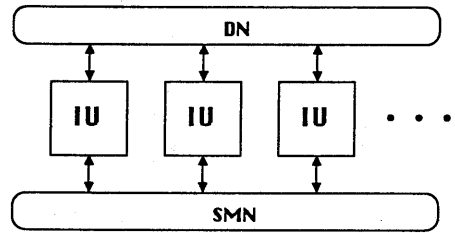
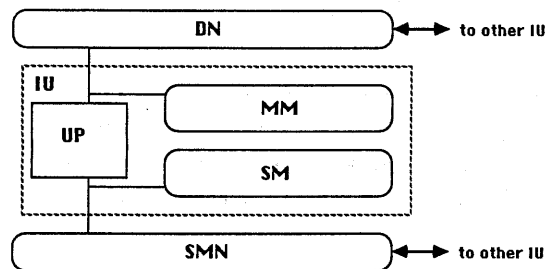


図1 システムの全体構成



IU: Inference Unit
UP: Unify Processor
MM: Memory Module
SM: Shared Memory
DN: Distribution Network
SMN: Shared Memory Network

図2 IUの内部構成

4. シミュレーションに於ける処理方式

4.1 ゴール・フレームの表現

PIEの基本処理単位はゴール・フレームであるが、今回のシミュレーションでは、そのゴール・フレームをリテラル単位で構成した。従って、コミット後に生成される、ボディ部の複数のリテラルは、それぞれ別のゴール・フレームとして別のIUに送られる。その際、SMNへのアクセスを減らすために、ローカルなSM中へのポインタは、デリファレンスした結果がGF中に書かれる。

4. 2 共有変数の配置と単一化

UPがGFの処理を行う際には、まず定義節中の全ての変数の領域をローカルなSM中にとり、undefに初期化する。そして、ヘッド・リテラルの単一化によって生じた結合は、その変数セルに書き込まれる。また、新GFを生成する(この処理をリダクションと呼ぶ)際には、変数セルのうち値がバインドされているものについてはその値を、undefのままのセルについてはそのセルへのポインタをGF中に書き込む。

4. 3 ヘッド・リテラル単一化時のサスペンド・アクティブイト

FLENGでは、次の2通りの場合にヘッド・ユニフィケーションのサスペンドが起こる。

- ①親ゴール中のundefの変数に値をバインドしようとした時
- ②定義節のヘッド・リテラル中のbind-to-non-variable annotationのある変数にundefの変数をバインドしようとした時

このような場合にサスペンドしたGFは、そのサスペンドの要因となった変数セルが具体化された時に実行可能となり、アクティブイトされる。GFのサスペンド、アクティブイト処理は次のように行われる。

- ①ヘッド・ユニフィケーション中に、サスペンドが起こるようなバインドをしようすると、MM中のGFのアドレスがサスペンドの要因となったundefセルの所に記憶される。
- ②GFの単一化において、1つの単一化も成功せず、1つ以上の単一化がサスペンドした場合は、GFはアクティブ・キューから外され、サスペンド・キューにつながる。
- ③undefセルが具体化されると、SMはサスペンドしていたGFをSMNを介してアクティブ・キューにつなぎ直す。

5. 負荷分散方式

5. 1 負荷分散の目的

GHCやFLENGの様なCommitted choice言語を並列推論マシン上で高速に効率良く実行させるためには、どの様な負荷分散方式を目指すべきであるか。並列マシンである以上、プロ

グラムの並列度を十分引き出せる様な負荷分散方式を目指すのはもちろんであるが、前述の様にCommitted choice言語は、共有変数のアクセス競合や、実行のサスペンド等の特徴が有り、これらが頻繁に生ずると処理速度に低下をもたらすのは明かであるので、本稿では以下の事項を目標として負荷分散方式を検討する。

- ①プログラムの並列度を十分に引き出す。
- ②共有変数へのアクセス競合を減らす。
- ③実行のサスペンドを減らす。

5. 2 自動負荷分散ネットワーク

前述のように、PIEにはゴール分配網と共有メモリ網の2種の相互結合網がある。特に、ゴール分配網は、自動負荷分散の機能をサポートしている。これは、ゴールの転送とは、逆の向きに行き先IUの負荷情報を伝送し、負荷の量が最小の行き先を選択して、そこにゴールを送りつけることにより、効率的な負荷の分配をするものである。

自動負荷分散ネットワークは、4入力4出力クロスバー・スイッチを基本とするスイッチング・ユニット(SU)を構成要素としたオメガ網であり、各SUは、次段のSUから送られてくる負荷情報を比較し、その中で最少の負荷情報を前段SUに伝送する。そして、負荷分散する時は、最少の負荷情報を伝えて来ている次段SUへの経路を設定する。

5. 3 負荷分散のパラメータ

負荷分散をするためには、まず、負荷量というものを何らかの形で評価する必要がある。PIEの基本処理単位はゴールであるので、負荷情報としてゴールの数というものが考えられる。但し、Committed choice言語においては、サスペンド状態であるゴールの数が非常に多く、これらのゴールの数を考慮に入れずにゴールの分配を行うと、それらのサスペンド状態のゴールが実行可能状態になった時にIUの負荷に大きな偏りを生じてしまう可能性がある。そこで、アクティブ(実行可能)なゴール数だけではなく、サスペンド状態のゴールの数も負荷情報の考慮に入れる必要がある。

負荷量の評価の次に考えなければならないの

は、リダクションによって新たに生成されたゴールを自分で実行するのか、あるいは、DNを通じて他のIUに送りつけるべきなのかを決定するための負荷量の閾値である。この閾値は、原則的には、GFを自分で処理するコストと他のIUに送りそのIUで処理してもらうコストの差で決まるのであるが、スケジューリングのことと考えると、それほど単純には決めることが出来ない。

負荷分散は全て自動的にするのか、あるいは、ある程度は、プログラマやコンパイラの指示も付け加えるべきなのかということは議論の分かれる点であるが、重要なことである。ここでは、プログラマやコンパイラが付け加えた指示があれば、それを尊重し、そのような指示がなくてもそれなりの性能を出すということをポリシーとする。

以上の点をまとめると、負荷分散のパラメータとしては、以下の3つが挙げられる。

- ① 負荷量を何で評価するか。
- ② GFを分配する時の閾値
- ③ プログラマやコンパイラによる負荷分散の指示(但し、具体的にどの様な指示を与えるかについては、後述する。)

6. シミュレーションの仮定

シミュレーションをする際に以下のことを仮定した。

- ① UPはPIEの試作UPとほぼ同程度の機能のハードウェアを仮定し、マイクロ命令レベルまでシミュレートする。
- ② 共有メモリ網は、4入力4出力クロスバーのスイッチング・ユニット(SU)を用いたオメガ網とし、1段あたりの接続時間は、2クロックとする。
- ③ SM、MMは、外部からのネットワークとローカルなUPから同時読みだし可能であり、書き込みを行う場合には、ロックがかけられる。
- ④ IU台数は16台、64台とする。

シミュレータはUNIX上のC言語で実装し、約11000行である。

シミュレーションに用いた例題は、以下のプログラムである。

- nreverse (リストの反転)

- qsort (リストのソート)
- primes (素数を求める)
- perm (リストの順番を入れ換えた結果を全て求める)

7. シミュレーション結果

7.1 ボディ部のゴール・リテラルの順番

GHCやFLENGの様なCommitted choice言語に於いては、ボディ部のゴール・リテラルは並列に実行されることを仮定しているの、本来は、その順番は処理に影響を与えないはずである。しかし、PIEは、新ゴールを生成するリダクションという処理を逐次的に行うので、その順番は処理に大きな影響を与える。前述の様にサスペンドを減らし、スムーズに実行が行われる様に負荷分散をするためには、データの流れる通りに実行されることが望ましいと考えられる。

図3にnreverseのFLENGプログラムを示す。このプログラムは、nrevの実行結果が返って来ない限りappendはすぐサスペンドしてしまうので、ボディ部の順番としては、nrev、appendの順にして、少しでも早くnrevの実行が行われるようにすることが望ましい。また、nreverseの場合は、最初はappendはサスペンドしてしまい、nrevだけが走り、それが終わって初めてサスペンドしていたappendが動き出す。そして、nrevが走っている間は、並列度が全く無く、その後で、appendが最高でリストの長さ分の並列度で動き始める。従って、理想的な形としては、図4の様に、各IUがパイプ・ライン状に動いてなるべく早くnrevの実行を終了させ、appendの実行に移ることが望ましい。

```
append([H|X],Y,Z):-append(X,Y,ZZ),unify(_,[H|ZZ],Z).
append([],X,Y):-unify(_,X,Y).

nreverse([H|X],R):-nreverse(X,XX),append(XX,[H],R).
nreverse([],R):-unify(_,R,[]).
```

図3 FLENGプログラムの例(nreverse)

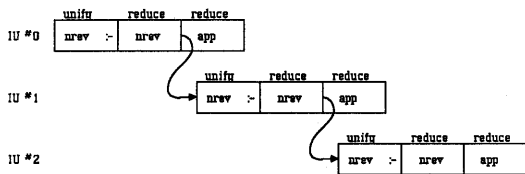


図4 nreverseの理想的実行

7. 2 負荷量の評価

負荷量の評価基準として、アクティブGFの数とサスペンドGFの数をを用いるということは前にも述べた。そして、1台のIUがアクティブGFを多く抱えているという状態は、プログラムの並列度を十分に出し切っていないということであり、あまり望ましくない。従って、アクティブGFの個数を重視する必要があると考えられるが、前述の様にサスペンドGFを無視するわけにも行かない。そこで、負荷量の評価基準として、

$$\text{負荷量} = \text{アクティブGF数} \times 256 + \text{サスペンドGF数}$$

という式を用いることにした。

7. 3 負荷分散の閾値

リダクションによって新しく生成されたGFを他のIUに送るか、それとも自分で処理するかという判断をするために負荷分散に閾値を設けるということを前述したが、これは、他のIUへGFを送ってそのIUで実行してもらうコストと自分で実行するコストの差で決まる。しかし、PIEのDNは、接続時間、スループット共に高速であり、通信コストは低い。よって、他のIUに実行してもらうのと自分で実行するのでは、あまりコストの差が無い。

従って、閾値はほとんどゼロに設定することが望ましいが、他のIUとのGFの差が一つしかない時にGFの無駄な転送キャッチ・ボールをしてしまうおそれがあるので、アクティブなGFの数の差が1以下の場合は他のIUにGFを送りつけないことを原則とする。

7. 4 プログラマやコンパイラによる負荷分散の指示

負荷分散の制御をあまり複雑にすると、そのためのオーバーヘッドが無視できないものとなるので、負荷分散を全く自動のみで行うことには、少し無理があると考えられる。そこで、ある程度は、プログラマやコンパイラがプログラムに書き加えた付加情報に従って、負荷分散を制御するということが考えられる。特に、急激な負荷の変動が起きた時（GFの数が急に増えるのではなく、サスペンドGFが急に次々とアクティブになる場合等）には、全自動の負荷分散方式では分散に乱れが生じる事がある。そのような場合に負荷分散をスムーズに行うために、付加情報は有効である。また、サスペンド、アクティブが頻繁に起きるプログラム（素数生成etc.）では、付加情報によってサスペンドの頻度を減らすことにより、実行をスムーズにすることが出来る。

7. 4. 1 負荷分散の付加情報

具体的には、付加情報として次の事項が考えられる。

- ① GFを他のIUへは送らず、自分のMMに残す。
- ② DNからの負荷情報と自分の負荷量を比較し閾値を見て、他のIUへ送るかどうかを判断する。（全自動の負荷分散方式と同じ）
- ③ 負荷量の比較、閾値との比較等はせず、GFをDNへ放出する。

7. 4. 2 付加情報を与える目安

実際のFLENGプログラムにどの様な基準で前述の3種の付加情報を割り当てるかについてであるが、これは以下の方針に従うことになる。

- ① ボディ部の各リテラルはデータの流れ通りに並んでいると仮定し、ジェネレータあるいはサスペンドを解除することが可能と考えられるシステム述語（unify等）等は、なるべく他のIUへ送ることが望ましい（③の指示を与える）
- ② ボディ部の最後尾のリテラルは、システム述語ならば自分のMMに残す（①の指示）がそうでない場合は、②の指示にする。その理由は、自分のMMに残すという指示は、かなりきびしいものであり、使い方を誤ると1台のIUへの負荷の集中が起こりかねず、新しいゴールを生

成ることが無い、システム述語に限った方が
良いと考えられたからである。また、unify述
語は、SMへのアクセスが多いのでその点を考
慮しても、やはりローカルMMに残すべきであ
る。

③GHCのガード部をFLENGに書き直した
場合は、そのためにシステム述語(比較等の)
が付加されることになるが、これは、真っ先に
実行してもらいたいものであるから、やはり、
③の指示を与える。

④リテラルの性質のよくわからないものに対し
ては、従来通りの指示(②の指示)を与える。

7. 4. 3 付加情報の効果

自動負荷分散ネットワークの効果、付加情報
の効果を比較するために、シミュレーション結
果を図5~7に示す。図5は、比較のためにラン
ダムに負荷分散を行った場合についてのシミ
ュレーション結果である。図6は自動負荷分散
ネットワークの機能のみを使って負荷分散を行
った場合のシミュレーション結果である。また
図7は、ユーザによる付加情報を考慮して負荷
分散を行った場合のシミュレーション結果であ
る。図5と図6を比較すると、図6の方が実行
速度が一樣に速く、Committed choice言語に於
いても自動負荷分散ネットワークが有効である
ことがわかる。また、図6と図7を比較すると、
nrev,primesは図7の方が実行が早く終了して
いる。これは、稼働プロセッサ数の時間変化を
見てもわかるように、平均並列度及び並列度の
立ち上がりの両方が向上したためであると考え
られる。特に、nrevの場合は、付加情報を与え
た事によって、appendの立ち上がりがスムーズ
になっている事がわかる。また、permに対して
は、付加情報の効果がほとんど無いが、これは、
付加情報が主にサスペンドを減らすことを目的
としているので、permの様にサスペンドをしな
いプログラムに対しては意味が無いためと考え
られる。

8. おわりに

以上のシミュレーション結果より、ユーザに
よる付加情報を用いて負荷分散を行う場合に
ついては以下のことがわかった。

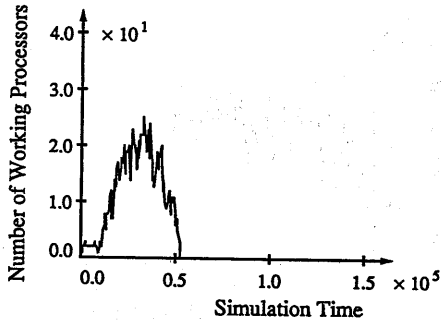
・付加情報を用いることにより、ゴールがサス
ペンドする頻度、時間等を減少させることが出
来ようになり、実行速度が速くなる。

・サスペンドを起こさないような例題において
は、あまり付加情報の効果は見られない。

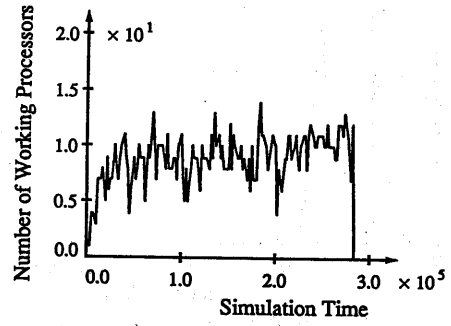
・Committed choice言語では、パイプライン並
列で各プロセスが動いていることが多いので、
そのパイプをなるべく乱さないような付加情報
が効果的である。

<参考文献>

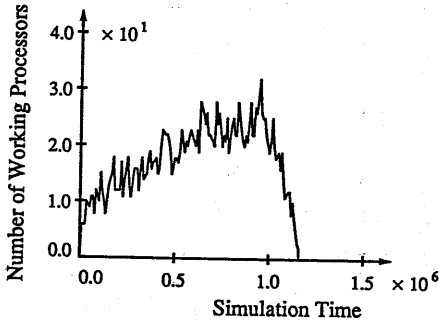
- [1] Moto-oka, T., Tanaka, H., et al,
"The Architecture of a Parallel Inference
Engine -PIE-", FGCS '84, ICOT.
- [2] 山内、小池、野田、田中、"高並列推論エ
ンジン実験環境PIEEE -自動負荷分散ネ
ットワーク-", 第34回情処全大、4P-3、1987
.
- [3] 坂井、小池、田中、元岡、"動的負荷分散
を行う相互結合網の構成"、情処論、Vol.27、
No.5、1986.
- [4] 垂井、田中、"並列推論マシンにおけるス
トリーム並列言語の実行方式の評価"、Proc.
of Logic Programming Conference '87、ICOT、
1987.
- [5] 山内、田中、"PIEにおけるストリーム
並列言語を指向したネットワーク構成法"、第
35回情処全大、3C-8、1987.
- [6] Nilsson, M., Tanaka, H., "FLENG
Prolog - The Language which turns
Supercomputers into Parallel Prolog
Machines", The Logic Programming
Conference '86, ICOT
- [7] 宮崎、瀧、"Multi-PsiにおけるFlat GHC
の実行方式"、Proc. of Logic Programming
Conference '86, ICOT, 1986.



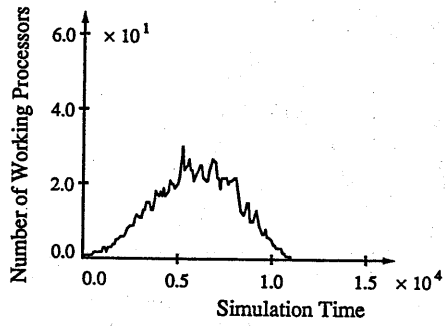
(nrev-60, 推論ユニット 64 台)



(qsort-256, 推論ユニット 16 台)

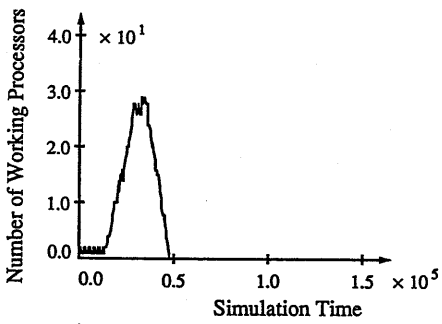


(primes-1000, 推論ユニット 64 台)

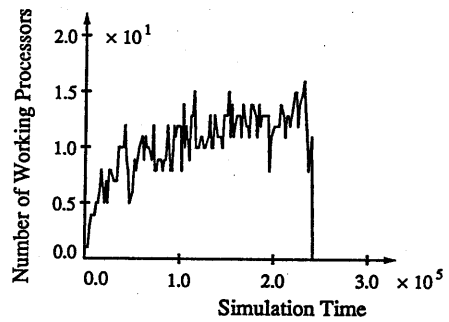


(perm-4, 推論ユニット 64 台)

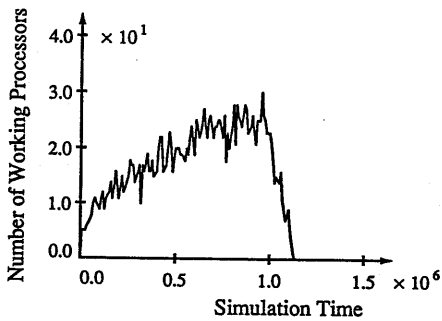
図 5 ランダム負荷分散の場合のシミュレーション結果



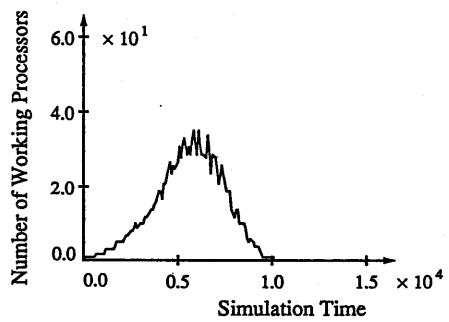
(nrev-60, 推論ユニット 64 台)



(qsort-256, 推論ユニット 16 台)



(primes-1000, 推論ユニット 64 台)



(perm-4, 推論ユニット 64 台)

図 6 全自動の負荷分散のみの場合のシミュレーション結果

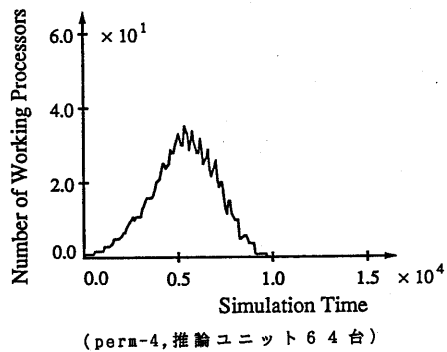
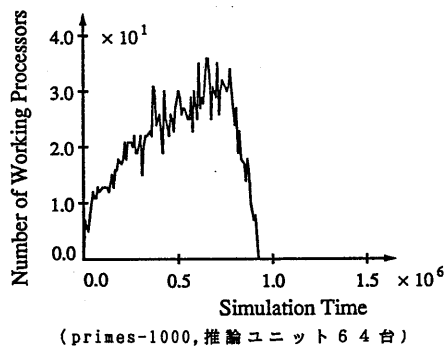
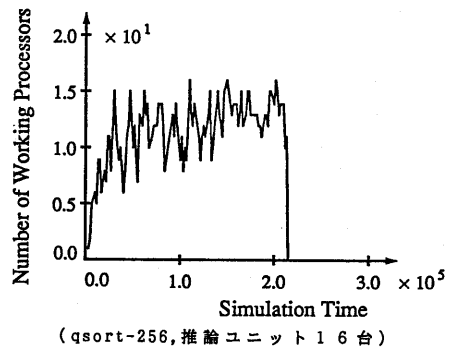
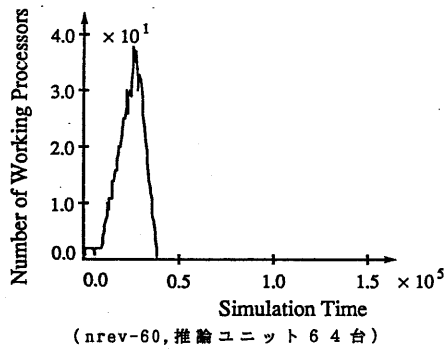


図7 付加情報を用いて負荷分散を行った場合