

# 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

垂井 俊明, 丸山 勉, 田中 英彦

( 東京大学 工学部 )

我々は、並列推論マシンでGHCを初めとするストリーム並列言語を実行する際に問題となる、共有変数へのアクセス等についての評価・検討を行なうために、ソフトウェアシミュレーションを行なった。シミュレーションでは、単一化プロセッサと共有メモリからなるユニットが多数台並列に接続されたモデルを仮定した。単一化プロセッサが他のユニットの共有メモリにアクセスする場合は、共有メモリネットワークを通じてアクセスが行なわれる。本稿では、そのシミュレーションモデルと処理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。

## A Preliminary Evaluation of the Execution Mechanism of Stream Parallel Languages

Toshiaki TARUI, Tsutomu MARUYAMA and Hidehiko TANAKA

Faculty of Engineering, University of Tokyo  
3-1 Hongo 7-choume, Bunkyo-ku, Tokyo, 113 Japan

There are two important problems to execute stream parallel languages in parallel inference machine. One is the control of goals that are suspended and then activated. The other is the access frequency to shared variables. We made a simulation to evaluate these two problems. In this paper, first we describe the control mechanism of goals and management of shared variables, then we present the simulation results of these two problems.

## 1. はじめに

我々は高並列推論エンジンPIE[1]の開発を進めている。PIEはゴール書き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における並列実行は、OR並列を主体に検討が進められて来た。従って現在のPIEのアーキテクチャそのままでは、GHCを始めとするストリーム並列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ストリーム並列言語を実行する際、問題になるのは次の2点である。

- ① 共有変数へのアクセスと負荷分散
- ② ゴール間の実行制御

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、検討を行なうことにする。具体的には、

- ・共有変数へのアクセスの頻度
- ・SM間ネットワークの構成
- ・負荷分散の方式

についての評価検討を行なう。

共有変数については、現在のPIEで用いられている集中型の構造メモリ[2]を利用することも可能であるが、アクセスの集中型の問題が生じる可能性が大きいため、集中型に代わり分散型の共有メモリ(SM)を導入するのが適当であると考えられる。SMと単一化プロセッサ(UP)が1つのユニットを構成し、各UPが同一ユニット内のSMに高速にアクセスすることを可能にすることにより、他のユニットのSMへのアクセスをおさえ、高速な処理を行なうことができる。

シミュレーションの対象とする言語については、ゴール間の実行制御が不要であるFLENGを用いることにする。これにより①のみに対する評価をより明確にすることができる。

## 2. 論理型言語FLENG

FLENG[3]はGHCと良く似た言語であり、以下の点がGHCとの主な相違点である。

- ① ガードが存在しない。従って定義節のヘッドリテラルの単一化が成功すると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。従って、全てのボディ部のリテラルは必ず実行され、あるゴールが失敗しても他のゴールには影響を与えない。論理的な AND関係が必要な場合は、プログラム中に明示的に書く必要がある。

また、FLENGにおいてもGHCと同様に、ヘッドリテラルの単一化を行なう場合親ゴールの変数を具体化することはできず、この単一化はサスペンドする。

以上述べたように、FLENGはGHCと共有変数に対するアクセス等については同一であるが、実行制御の部分は大幅に簡略化されている。また、FLENGはGHCより低レベルな言語であるが、ほとんどのGHCのプログラムはFLENGに容易にコンパイルすることができ、すでにコンパイラも完成している。

## 3. シミュレーションモデルの全体構成

図1に今回シミュレーションを行なうシステムの全体構成を示す。単一化等の処理を行なう単一化プロセッサ(UP)と共有メモリ(SM)の対が基本処理要素である推論ユニット(IU)を構成し、IUが多数台並列に接続されてシステム全体を構成している。

システムにはUP間を結び、基本処理単位であるゴールフレーム(GF)の分配等を行なう分配網(DN)、SM間を接続し共有変数の値や構造データ等のやり取りを行なったり、サスペンド状態にあるGFの起動を行なうための通信に用いられる構造メモリ網(SMN)の2種類のネットワークがある。DNについては、従来PIEで用いられてきた多段結合網[4]を用いる。各IU間での負荷分散は、DNが各IUの負荷に応じGFを自動的に分配することにより行なわれる。SMNについては現在のところどのような構成が最適であるか分からないため、とりあえず理想的な(接続時間一定、閉塞なし)ネットワークを仮定してシミュレーションを行なった。また同一IU内のUPとSMは密に結合されており、SMNを通さずに直接、高速にアクセスすることが可能である。これにより、負荷分散をうまく行なえ

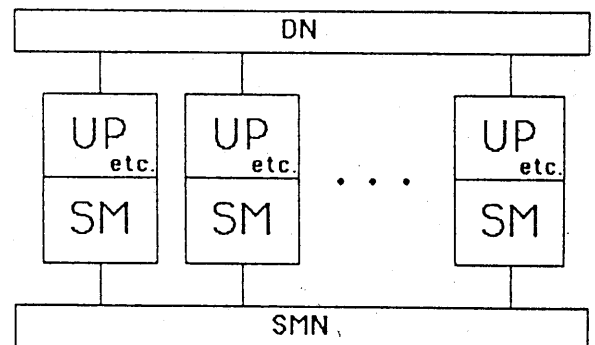


図1 システムの全体構成

はSMNへのアクセスの回数を減らすことが可能になる。

図2に推論ユニットの内部構成を示す。各部の機能は以下の通りである。

・DM (定義節メモリ)

プログラムを保持する。全てのDMは同じ定義節を持つ。

・UP (単一化プロセッサ)

単一化及びゴールの生成を行なう。UPからはDM、MM、Cache、及び同一IU内のSMの内容を直接アクセスすることができる。他のIUのSMをアクセスする場合は、SMNを通してアクセスを行なう。

・MM (メモリモジュール)

ゴールフレーム (GF) の格納、管理等を行なう

・SM (共有メモリ)

共有変数及び構造データを記憶する。

・Cache (キャッシュ)

他のIUのSMの内容のキャッシュである。ここで言うキャッシュとは、外部SMのアドレス (IU番号、IU内アドレス) によりSMの内容を検索することのできる連想メモリである。Cache中には他のSMの中の共有変数の値、構造データ、及びundefであるセルのアドレスが記憶される。

#### 4. シミュレーションにおける処理方式

##### 4.1 ゴールフレームの表現

今回の処理モデルでは、処理の基本単位となるゴールフレーム (GF) はリテラル毎に分割されている。従って単一化の結果、新たに (AND関係にある) 複数のリテラルが生成された場合、それらは別々のGFとして異なるMMに分配される。その際SMNへのアクセスを減らすため、GFを生成するUPと

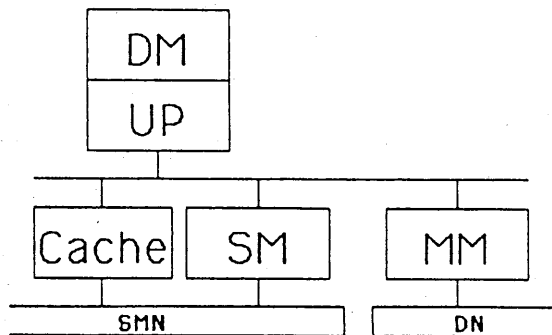


図2 IUの内部構成

同じ (ローカルな) SM中にある構造データは、GFの構造部としてGFのリテラル部と共にコピーされる。また、同様にローカルなSMへのポインタについても、デリファレンスした結果がGF中にコピーされる。

##### 4.2 共有変数の配置と単一化

UPにおいてGFの処理を行なうとき、定義節中の変数は全てローカルなSM中におかれ、undefに初期化される。そしてヘッドリテラルの単一化中におきた結合は、その変数セルに書き込まれる。定義中の構造データが結合された場合は、構造データはSM中にコピーされ、変数セルにはそこへのポインタが書き込まれる。新GFの生成時には、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、その値がGF中にコピーされ、undefのままのものについては、そのセルへのポインタがGF中に書き込まれる。

##### 4.3 単一化のサスペンドとアクティブイト

FLENGにおいて単一化がサスペンドするのは以下の2つの場合である。

- ・ヘッドリテラルの単一化中に、親ゴールのundefの変数を具体化しようとしたとき
- ・定義節のヘッドのbind-to-non-variable annotationのある変数にundefの変数を単一化しようとしたとき

サスペンドしたGFはundefの変数に値がバインドされると再び実行することができるようになる。GFのサスペンド及びアクティブイトの処理は以下のように行なわれる。

- ①ヘッドリテラルの単一化中に、サスペンドが起るようなバインドをしようとする (実際にはバインドをしようとしてundefのセルを読みに行った時点で)、MM中のGFのアドレス (IU番号、IU内アドレス) がSM中のサスペンドの原因となったundefの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、GFはMM中のアクティブなキューから外され、サスペンド状態になる。
- ③その後undefのセルに値がバインドされると、SMはサスペンドしていたGFにSMNを通してアクティブイトコマンドを送り、GFは再び実行される (それと同時に新しい値がCacheに書き込まれる)。

UPで実行中のGFに対してアクティブイトコマンドが来た場合は、UPは実行を中止し、Cache中の新しい値を使ってGFの処理をやりなおす。

#### 4.4 Cache における処理

各IUは他のIUのSMの内容のキャッシュを持つ。UPが他のSMのデータをアクセスする場合は、先ずCacheの内容を調べ、該当するエントリがある場合(ヒットした場合)はその値を用いて処理が続けられ、外部SMへのアクセスは行なわれない。Cacheがヒットしなかった場合にはSMNを通じて値の読み出しが行なわれ、読み出された値はCacheに登録される。また、他のSM中のundefのセルのアドレスについても、1つのGFの処理が終了するまでここに記憶される。

#### 4.5 システム述語 unify における処理

FLENGにおいて親ゴールの変数への値の書き込みは、システム述語 unifyにより行なわれる。その際に、もしUPが書き込もうとしたSM中の変数セルに既に他のUPにより値が代入されていた場合は(書こうとしたセルがundef以外だった場合)、UPは書き込まれていた値を読み出し、その値と書き込もうとした値との単一化を行なう。

#### 4.6 負荷分散の方式

IU間の負荷分散はDNが各IUの負荷(IU内のGFの数)に応じ、GFを分配することにより行なわれる。UPにより生成されたGFは、接続可能なIUのうちで一番アクティブなGFの数が少ないIU(アクティブなGFの数が同じ時は、一番サスペンドしたGFの数が少ないIU)に送られる。ただし、同一IU内のMMにアクティブなGFがない場合は、GFはDNを通さずにローカルなMMに戻される。

#### 5. シミュレータでの仮定

シミュレーションに当たっては以下のことを仮定した。

- ① UPは現在の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 ネットワークは一定時間でかならずつなげる閉塞のない理想的な回線交換網とする。
- ③ SMNをつなげるのに要する時間については、0~20(クロック)の間で変化させ、接続時間の実行時間への影響を調べる(それ以外の測定では5クロックとした時のデータを測定した)。
- ④ Cacheは連想メモリであり、必要なデータを高速に検索することができる。

⑤ Cacheの容量は十分大きい。

⑥ SM、MMはマルチバンクになっており、外部からのネットワークとローカルなUPから同時に読み出しを行なうことができる。ただし、書き込みを行なう場合にはロックがかけられる。

シミュレータはUNIX上のC言語で実装し、約10000行である。

シミュレーションに使用した例題は以下の4つである。

- ・nreverse120 (NREV) - 長さ120のリストの反転
- ・qsort64 (QSORT) - 長さ64のリストのソート
- ・primes600 (PRIMES) - 600までの素数を求める
- ・perm5 (PERM) - 長さ5のリストの順番を入れ替えた結果を全て求める(全探索)

#### 6. シミュレーション結果

##### 6.1 並列度

IUを1~64台まで変えた場合の台数効果を図3に示す。また、IU64台の時の各例題の稼動IU台数、及びGF数の変化を図4~図11に示す。また、表1にサスペンド、アクティブの回数を示す。

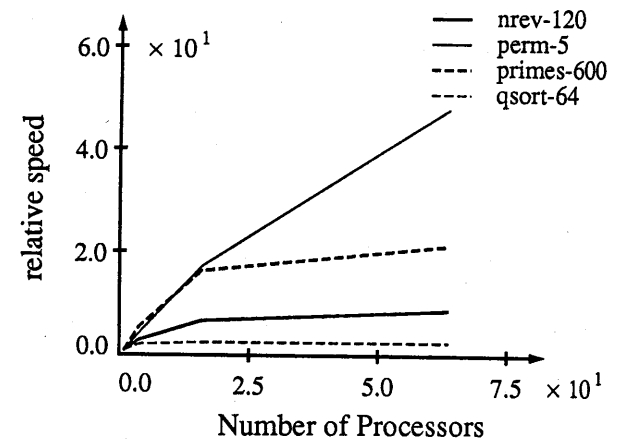


図3 台数効果

	nrev120	qsort64	primes600	perm5
num. of suspended goals	539	584	20246	0
num. of hook	701	1644	41811	301
num. of activate command	579	1378	34701	301
suspended in MM	539	584	20246	0
executed by UP	31	62	1285	23
no goals to activate	9	732	13170	278

表1 サスペンド、アクティブの回数

PERMはOR並列な全解探索問題なので、サスペンド等は生じておらず、ストリーム並列処理は行なわれていないため、非常に高い並列度が得られている。PRIMESは素数のフィルタがストリーム並列に動くため、比較的大きい並列性を持っており、suspendするGFも多い。NREVは前半のリストを分割する部分には全く並列性は見られないが、後半のappendによりリストを結合する部分は、多数のunifyが並列に動くため高い並列度を示す。またゴール数のグラフ(図8)からは、前半のリストの分割が終わった所で、appendが120個サスペンドしていることがわかる。QSORTも同様に最後のリストを連結する部分の並列度が高くなっているが、差分リストを使っているため並列度はそれほど高くない。

## 6.2 SMへのアクセス数

表2にSMへのアクセス数を示す。NREVや

PRIMES等では定義節の変数以外へのローカルSMへのアクセスが外部SMへのアクセスの倍程度あり、共有変数へのアクセスを考慮に入れた負荷分散を行っていないにもかかわらず、ある程度のローカリティがあることがわかる。PERMで外部SMへのアクセスが非常に少ないのは、PERMはもともとORパラレルな問題であり、解の収集の時以外には外部SMへのアクセスがほとんど行なわれないからである。

	nrev120	qsort64	primes600	perm5
total access to SM	131499	22086	337887	91465
access to external SM	29397	2647	84575	978
access to local SM	102102	19439	253312	90487
access to Def's variable	29161	4267	101346	13535

表2 SMへのアクセス回数

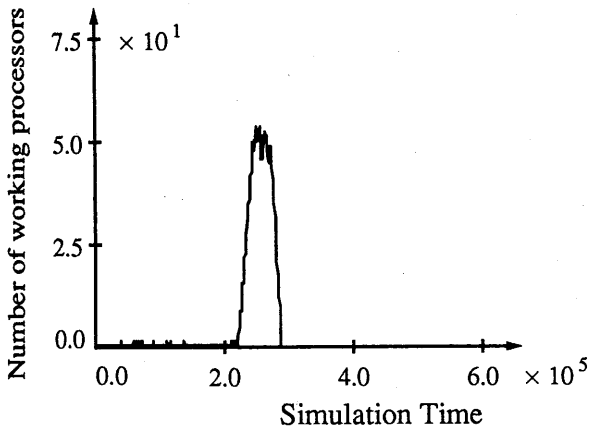


図4 稼働プロセッサ台数の変化(NREV)

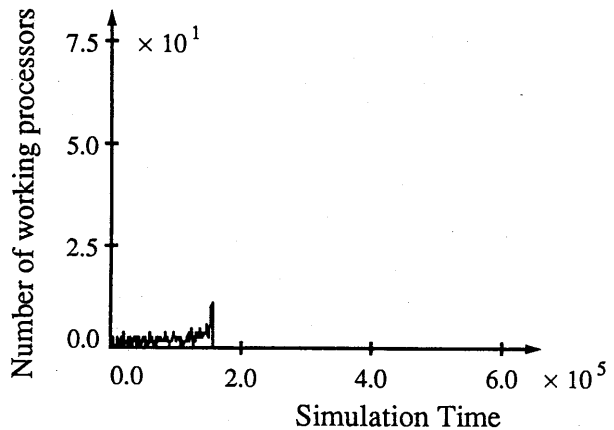


図5 稼働プロセッサ台数の変化(QSORT)

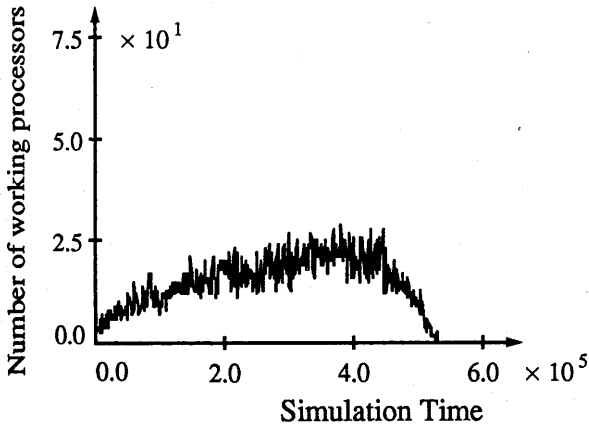


図6 稼働プロセッサ台数の変化(PRIMES)

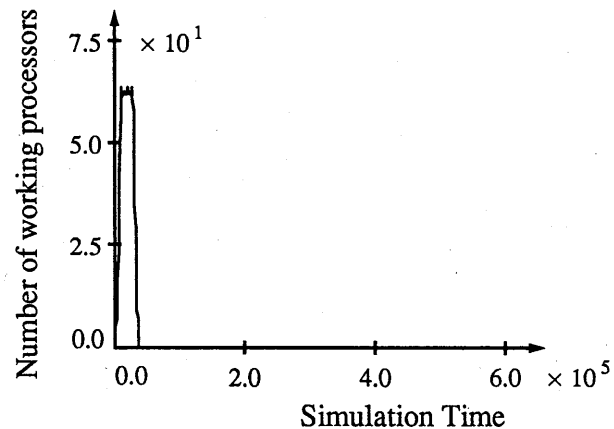


図7 稼働プロセッサ台数の変化(PERM)

### 6.3 キャッシュの効果

表 3に外部SMへの読みだしと Cacheのヒット数を示す。Cacheのヒット数は多くて2~3割であり、そのうち undefのヒットが半分程度を占めている。また注目すべき点として、以前に実行した別のGFがエントリしたデータがヒットすることはほとんどないことがわかる。これは、外部SMへのポインタは、一度値が読出されるとその値は新GFの生成時にGF中にコピーされるため、ほとんどのポインタは1回程度しか読まれないためである。従ってCacheがヒットするうちのほとんどは、ある入力ゴールに対して複数の定義節の単一化を行なっているときに生じるものである。従って、undef以外のCache内のデータも1つのGFの実行が終わった時点で消去してしまえば良いと考えられる。

	nrev120	qsort64	primes600	perm5
external SM read	22143	1989	62866	498
total cache hit	907	737	19700	18
undef	436	320	6430	0
except undef	471	417	13270	18
entried by other goals	8	2	44	0

表 3 キャッシュの効果

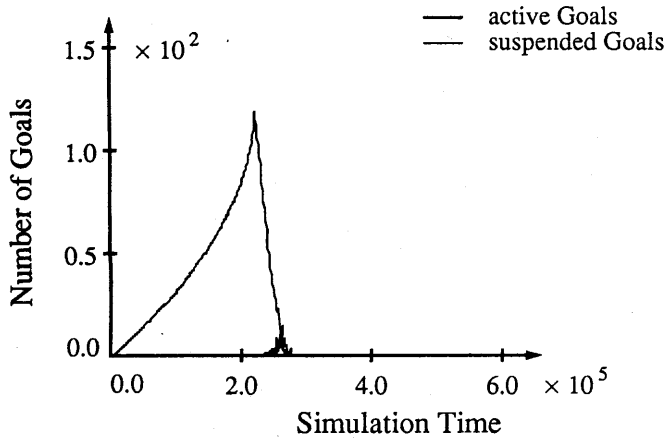


図 8 GFの数の変化(NREV)

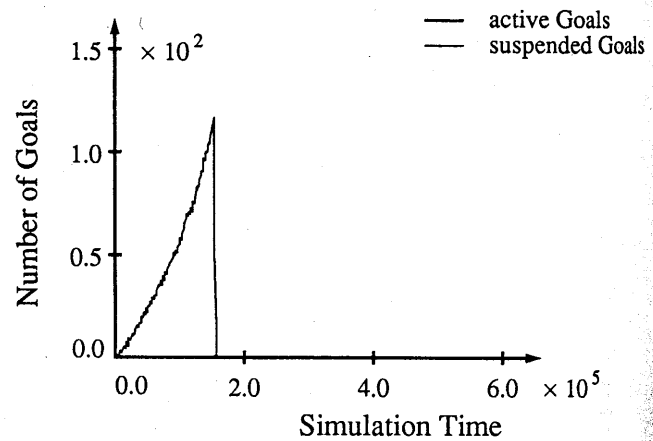


図 9 GFの数の変化(QSORT)

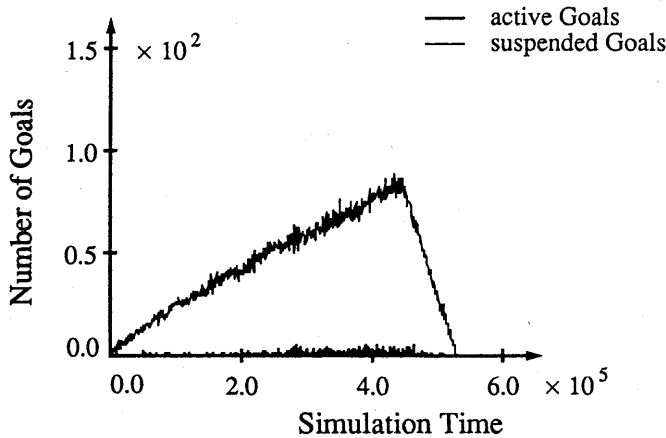


図 10 GFの数の変化(PRIMES)

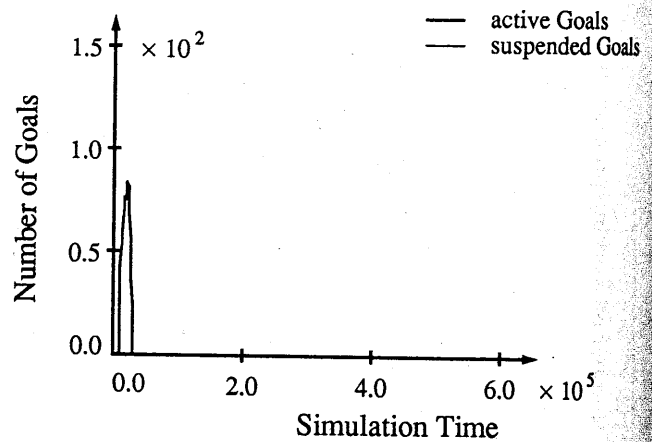


図 11 GFの数の変化(PERM)

#### 6. 4 SMNの接続時間の影響

図12にSMNをつなげるのに要する時間を0~20(クロック)と変化させた場合の実行時間の変化を示す。PERMやQSORTのような外部SMへのアクセスが比較的少ない例題では実行時間はほとんど変化しないが、PRIMESのように外部SMへのアクセスやサスペンド等が多い例題では速度にかなりの変化が見られる。従ってSMNの接続時間はストリーム並列計算機の性能にかなり影響を与えることになる。

#### 6. 5 SMNへのアクセスの頻度

ここではSMNへのアクセスが多いNREVとPRIMESの結果により検討を行なう。図13、図14にSMNへの同時アクセス数の最大値と平均値の時間変化を、図15~図18にSMNへの読み出し要求と書き込み要求の頻度を示す。NREVの後半のように、64台ほとんど全てのプロセッサが単一化を行なっている場合には、SMNを通じた外部SMへの読み出し要求が約2クロックに1回、書き込み要求が約6クロックに1回生じている(これは、1回の単一化に約50クロックかかっていることから理解できる値である)。そしてその時にはSMNに最大12個(平均でも6個)の同時アクセスがおきている。(1回当りの平均アクセス時間は、SMNを接続する時間が5クロックの時で約10クロックである。)従ってこの方式で64台のIUをつなげるとSMNへのトラフィックはかなり多くなることがわかる。

#### 7. おわりに

以上のように、シミュレーション結果よりSMNへのトラフィックはかなり大きく、SMNの構成によってシステム全体の性能が左右されることがわかっ

た。SMNの構成についてはバス、多段網等が考えられるが、シミュレーション結果より、SMNをバスで接続した場合は接続可能なIUの台数はせいぜい16台程度であり、システムを階層構造にする必要がある。また、SMNを多段網にした場合の性能についても、より詳細に検討してみる必要がある。今後は

- ・ SMNの構成の検討
  - ・ 共有変数の他のアクセス方法(マルチリード・ワンライト方式等)についての検討
  - ・ より効率的な負荷分散方式についての検討。
- 等を行なって行きたい。

#### <参考文献>

- [1] Moto-oka, T, Tanaka, H, et al, "The Architecture of a Parallel Inference Engine -PIE-", FGCS '84, ICOT.
- [2] 平田、田中、元岡: "PIEにおける構造メモリの構成について"、情報処理学会、アーキテクチャワークショップインジャパン '84、1984.
- [3] Nilsson, M, Tanaka, H, "FLENG Prorog- The Language which turns Supercomputers into Parallel Prolog Machines", The Logic Programming Conference '86, ICOT.
- [4] 坂井、田中、元岡: "動的負荷分散を行なう互結合網"、信学技報、1985年8月.
- [5] 垂井、丸山、田中: "PIEにおける並列論理型言語FLENGの実行方式"、情報処理学会第33回全国大会、5B-4、1986.

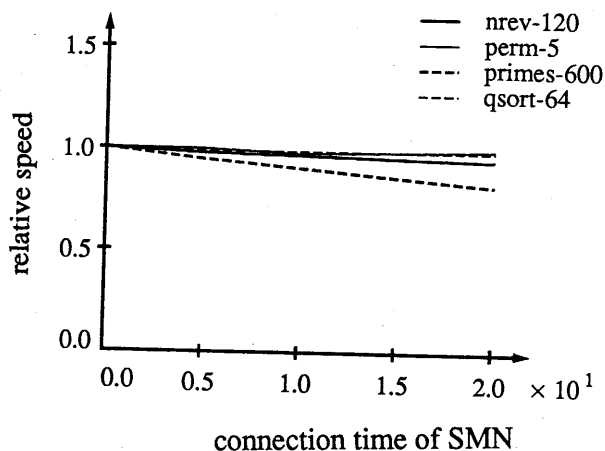


図12 SMNの接続時間の影響

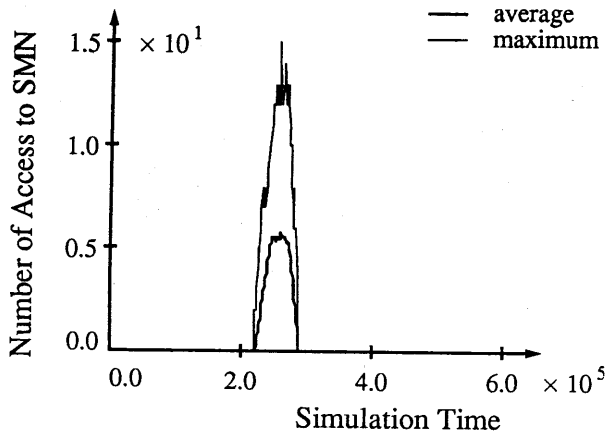


図13 SMNへの同時アクセス数(NREV)

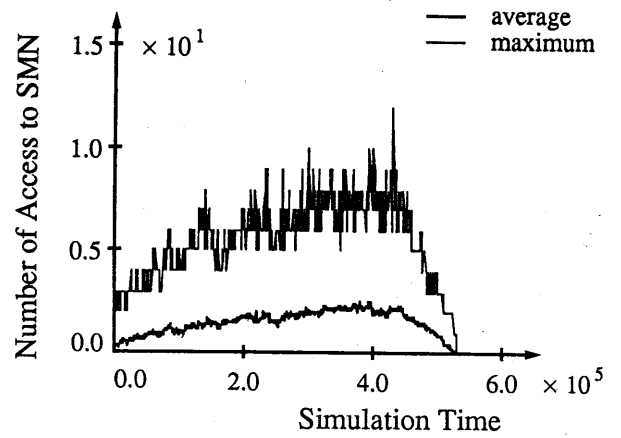


図14 SMNへの同時アクセス数(PRIMES)

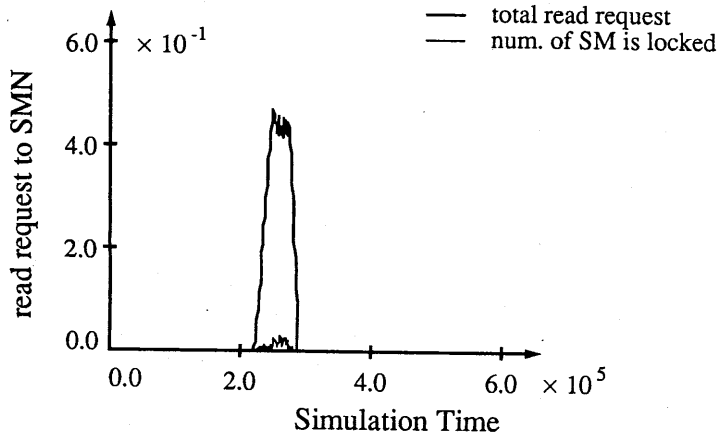


図15 SMNへの読み出し要求(NREV)

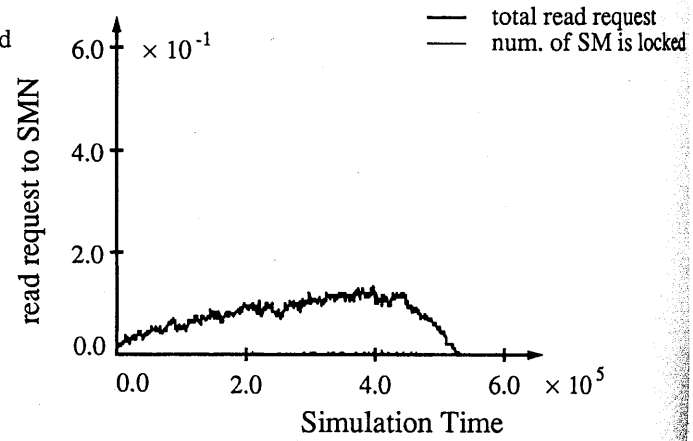


図16 SMNへの読み出し要求(PRIMES)

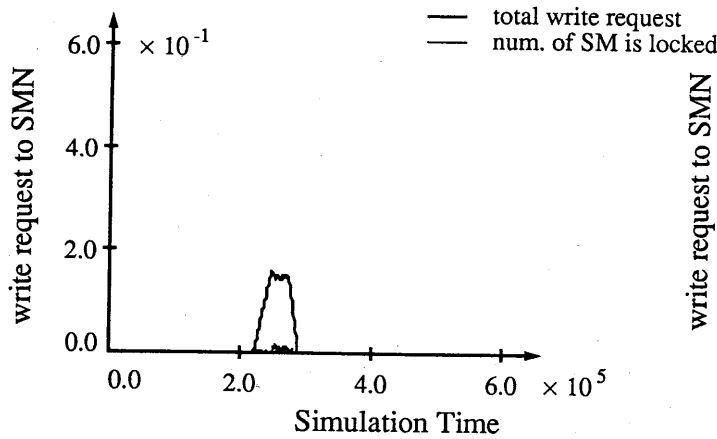


図17 SMNへの書き込み要求(NREV)

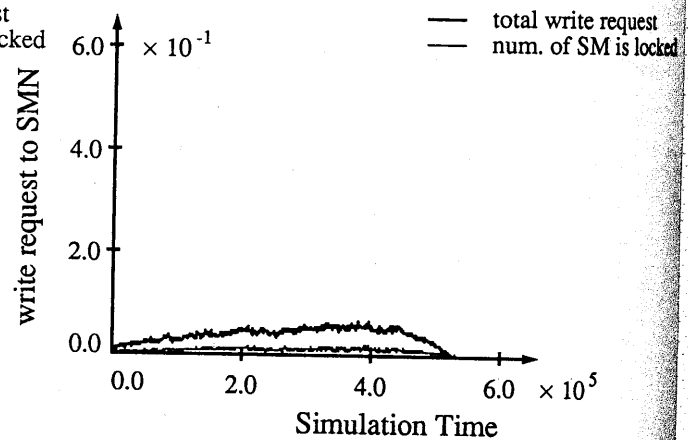


図18 SMNへの書き込み要求(PRIMES)