

## ASURAの解析モデル

城 和貴

(株)クボタ  
コンピュータ事業推進室

大阪市浪速区敷津東1丁目2番47号

### Abstract

分散共有メモリ型マルチプロセッサASURAに対するセミ・マルコフ過程を利用した解析モデルを提案する。提案されるモデルはコヒーレンス制御やネットワーク競合による待ち状態を容易に記述できるため、大規模で複雑なアーキテクチャを持つ並列計算機の性能評価に利用できる。モデル化の目的は、ASURAのキャッシュ・ヒット率及びデータのリード/ライト比が、キャッシュ・コヒーレンス制御を含むリクエストや、それに対する待ち状態に、どのような影響を与えるかを解明し、性能評価を与えることである。さらに、構築されたモデルを用いて、プロセッサ利用率、通常データ/コヒーレンス制御リクエスト及びそれらの待ち時間について実際の評価を行なった。

セミ・マルコフ過程、並列計算機、解析モデル、性能評価

## An Analytic Model for the ASURA

Kazuki JOE

Office of Computer Business, KUBOTA Corporation

1-2-47, Shikitsuhigashi, Naniwaku, Osaka, Japan

*E-mail: joe@kocb.astem.or.jp*

### Abstract

This paper presents a discrete time model of memory and network interference for the cache coherence mechanism in the ASURA, which is a caching cluster-based distributed shared memory multiprocessor computer. It differs from earlier models in its ability to model variable waiting time both for normal requests and cache coherent requests. Actually, the aim of our model is to analyze how the ratio of cache hit and read/write requests can affect the performance from the view point of normal/coherent requests and their waiting states. Furthermore, we evaluate the performance, IO request rates, processor utilization etc., of the ASURA using this analytic model.

SemiMarkov Processing, Parallel Computer,  
Analytic Model, Performance Evaluation

## 1 はじめに

プロセッサ・クラスタ方式の並列計算機は将来のスーパーコンピュータへの有力な一方式として、研究開発が推進されてきた。特に共有メモリを持つものは[8]並列コンパイラ開発の観点からも望ましいものである。さらにキャッシュ制御方式の研究[1]に伴い、キャッシュを持ったクラスタ方式の共有メモリ型並列計算機が精力的に研究されている[7][5]。いずれもネットワークの性能を補うものとしてキャッシュを有効利用しており、スケラブルなクラスタ型共有メモリ方式の並列計算機という新しい方向に向かっている。

このような大規模な並列計算機アーキテクチャの妥当性を確かめるのに、シミュレーションは有効な手法である反面、数万数十万規模のプロセッサの動きを検証するのに要するコストは新たな問題になるであろう。

一方、並列計算機の性能評価を行なうのに、確率モデルや確率過程を用いた解析モデルは低コストで有効な手法である。しかしながら、これまでに提案されてきた並列計算機システムの解析モデル[4][6][10][2][3][9]では、キャッシュを持ったアーキテクチャは対象にされておらず、また、リクエストに対するネットワーク競合等の待ち状態を記述しているモデルも少ない。

我々は既に大規模な並列処理計算機への第一歩として、クラスタ・ベースの階層型マルチプロセッサ・システムASURAを提案し[15]、ASURA全体を確率モデルで表した性能評価[13]、ASURAクラスタに着目した解析モデル[14]について報告している。

そこで本稿では、マルコフ連鎖の各状態で任意の時間滞在出来るセミ・マルコフ過程[11]を階層的に構成し、クラスタ型並列計算機に適用することに着目し、ASURA全体を表現した解析モデルを提案する。

この解析モデルの目的は、ネットワークで結合されたキャッシュを持つクラスタ型分散共有メモリ方式の並列計算機、ASURAに対する、定性的な解析とその評価である。

本稿では、ASURA全体のモデル化のための仮定を述べ、モデル化のための状態記述の定義を行ない、具体的な計算手順を示した後、性能評価を行なう。

## 2 準備

### 2.1 ASURAの定義

性能評価のモデル化のために、ASURAを次のように定義する。 $P$ 個のプロセッサ(プライベート・キャッシュ(PCH)を持つ)と1つの<sup>1</sup>ローカル・メモリ(LM)、ネットワーク・インターフェイス(NIF)がバス結合され、プロセッサ・クラスタ(PC)を形成する。ただし、NIFはグローバル・メモリ(GM)とグローバル・キャッシュ(GCH)、PC間の通信を制御するコントローラからなる。(PC間の通信はPC内のバスに影響を

与えないことに注意)  $M$ 個のPCはリング結合されてクラスタ・グループ(CG)を構成し、 $K$ 個のCGがクロスバ結合することによりASURA全体が定義される。従って、全体では $P \times M \times K$ 個のプロセッサとPCH、それぞれ $M \times K$ 個のローカル・メモリ・モジュール、グローバル・メモリ・モジュール、GCH、 $M \times K$ 本のバス、 $K$ 個のリング、1個のクロスバによって並列計算機システムを構成することになる。メモリ階層性を示すために、メモリ・ロケーションをローカル・メモリ(LM)、PC内グローバル・メモリ(PCGM)、CG内グローバル・メモリ(CGGM)、リモート・グローバル・メモリ(RGM)で区別する。なお、PCH、GCHのキャッシュ・コヒーレンス・プロトコルはそれぞれ、イリノイ型の変形[12]、Synapse[15]に従う。

### 2.2 セミ・マルコフ過程

セミ・マルコフ過程(以後SMPと呼ぶ)に関する詳細は文献[11]に譲るとして、ここではSMPについての簡単な説明を行う。

SMPとは $K$ 個の状態からなる確率過程である。状態 $i$ においては平均 $\eta_i$ の間滞在し、状態 $j$ に $p_{ij}$ の確率で遷移する<sup>2</sup>。もしSMPがエルゴード的な既約なエンベデッド・マルコフ連鎖(以後EMCと呼ぶ)を持つなら、状態 $i$ の定常分布 $P_i$ は次式で表される。

$$P_i = \frac{\pi_i \eta_i}{\sum_{j=1}^K \pi_j \eta_j} \quad (1)$$

ただし、 $\{\pi_i\}$ はEMCの定常分布である。本稿におけるモデルでは、後に述べるようにエルゴード的な成分が1つだけからなるEMCで表されているので、式1が利用できる。また、SMPの状態 $i$ から脱出する確率 $\lambda_i$ は次のようにして求められる。

$$\lambda_i = \frac{P_i}{\eta_i} = \frac{\pi_i}{\sum_{j=1}^K \pi_j \eta_j} \quad (2)$$

### 2.3 モデルの仮定

本稿ではASURAのキャッシュ・コヒーレンス制御を考慮に入れた定性的評価を行なうために、プロセッサおよびPCの状態遷移をマルコフ連鎖によってそれぞれ表し、これらをEMCと見なすことにより、SMPを定義し、任意のリクエストの待ち時間を考慮した解析モデルを構築する。モデルを簡略化するために、次のような仮定を導入する。

1. 全てのプロセッサからのリクエストが全て独立であるような並列プログラムが稼働中である。言い換えれば特殊な同期操作は本モデルでは記述出来ないことになる。
2. プロセッサとPCの状態は、計算(アイドル)状態、ネットワークを握ってリクエストを実行している状態、その状態への待ち状態。
3. プロセッサ内オン・チップ・キャッシュにヒットした状態は計算状態と解釈。
4. PCHにヒットした状態は、バスを使わないリクエストと解釈し、同時に計算も行なう。
5. 他のプロセッサの状態に影響を与えるようなリクエストは、その作用が自分自身に影響を与えないものとする。

<sup>1</sup>実際にはインターリーブされているが、バスとメモリのサイクル比が同じと考えた方がモデルが簡単であるためこのような仮定を用いる

<sup>2</sup>各状態の滞在時間がすべて1である時、そのSMPは通常のマルコフ連鎖と等価なことに注意された。

6. プログラムは安定稼働している状態、すなわち、起動時のページングやキャッシュ・ヒット率の悪さは考慮しない状態を仮定する。またキャッシュは既に使い切った状態で、それに伴うリプレースについては考慮する。
7. ファイルI/O等によるバスへのリクエストは考慮しない。

### 3 モデル化

#### 3.1 モデル化の概要

ASURAは階層的なネットワークおよび複雑なメモリ・システムを持つため、全体を一つのモデルで表すことは困難である。本稿ではモデルの簡略化のため、まず、PC内バスの影響を直接受けるプロセッサとPCHを一つのモデル（以後Pモデルと呼ぶ）として表し、PC間ネットワークに直接影響を受けるNIF（GCHを含む）を別のモデル（以後Gモデルと呼ぶ）として表す。

Pモデルは[14]で構築したASURAクラスタのモデルの自然な拡張となる。つまり、ASURAクラスタのモデルにPC間通信のための状態を追加することで、Pモデルが形成される。また、GモデルにおけるPC間通信のメモリ階層に応じたリクエスト率は、[13]によるものとする。

PモデルにおいてGCHおよび共有メモリに対してなされるリクエスト（通常のデータ・リクエストとMLI保証のためのGCH制御のリクエスト）はPモデルからGモデルへの出力として、また、GモデルにおいてMLI保証のためになされるPCHへの制御はGモデルからPモデルへの出力として定義される。

PモデルおよびGモデルを使ってASURA全体を表すためには、まずPモデルに対する適当な初期値（このときPC間通信はないものとする）を与え、Gモデルへの出力を得る。次に得られたGモデルへの入力と適当な初期値によって、GモデルからPモデルへの出力を得る。以後、PモデルおよびGモデルが安定状態に達するまでこの操作を繰り返すことにより、ASURA全体の定常状態を知ることができる。

#### 3.2 Pモデル

##### 3.2.1 状態の定義

Pモデルのモデル化を行なうために、状態定義を行なう。PモデルはASURAクラスタのモデルにクラスタ間通信のための状態を追加したものであるため、[14]に、グローバル・データに対するPCHのミス・ヒットに起因するリクエスト（この場合、GCHもしくは共有メモリが参照される。）およびその待ち状態、他のPCへのコヒーレント制御のためのリクエストおよびその待ち状態、の4状態を追加することとする。

$pCOM$  計算状態、もしくはプロセッサ内キャッシュのアクセス  
 $pRh$  PCHからの読み込み  
 $pWh$  PCHへの書き込み

$pWhIV$   $pWh$ により発生するインバリデーション、もしくはMLI保証のためGCHから要求されるインバリデーション  
 $\overline{pWhIV}$   $pWhIV$ への待ち状態  
 $pRc$  そのライン属性が *Dirty* ではないローカル・データのキャッシュ・ミスによる読み込み  
 $\overline{pRc}$   $pRc$ への待ち状態  
 $pRd$  そのライン属性が *Dirty* であるローカル・データのキャッシュ・ミスによる読み込み  
 $\overline{pRd}$   $pRd$ への待ち状態  
 $pRdWB$  *Dirty*  $pRd$ により発生するライトバック  
 $\overline{pRdWB}$   $pRdWB$ への待ち状態  
 $pWc$  そのライン属性が *Dirty* ではないローカル・データのキャッシュ・ミスによる書き込み  
 $\overline{pWc}$   $pWc$ への待ち状態  
 $\overline{pWcIV}$   $pWc$ により発生するインバリデーション  
 $pWcIV$   $pWcIV$ への待ち状態  
 $pWd$  そのライン属性が *Dirty* であるローカル・データのキャッシュ・ミスによる書き込み  
 $\overline{pWd}$   $pWd$ への待ち状態  
 $\overline{pWdWB}$   $pWd$ により発生するライトバック  
 $pWdWB$   $pWdWB$ への待ち状態  
 $pRP$  ラインのリプレースにより発生するライトバック、もしくはMLI保証のためグローバル・キャッシュから要求されるライトバック  
 $\overline{pRP}$   $pRP$ への待ち状態  
 $pRW$  グローバル・データのPCHミス・ヒットによる読み書き  
 $\overline{pRW}$   $pRW$ への待ち状態  
 $pCO$  他のPCへのコヒーレント制御命令  
 $\overline{pCO}$   $pCO$ への待ち状態

図1は定義された状態間の遷移を示す。図1において、 $H$ はPCHのヒット率、 $R$ はプロセッサからの通常のリクエストに対するリード・リクエストの割合（従って、 $1-R$ はライト・リクエストを意味することになる。プロセッサのオペレーションに対するリクエスト率ではないことに注意。）、 $D$ はリクエスト要求のあるデータを含むラインが *Dirty* である確率、 $C$ はリクエスト要求のあるデータを含むラインが自分以外の少なくとも一つのキャッシュで *Valid* である確率、 $W$ は（バスに対して）リクエストを出した時に待ち状態に陥る確率、 $X$ はキャッシュが既に一杯である確率を意味する。この時、プロセッサの各サイクル毎のバスに対するリクエスト率 $\phi$ が必要となる。ここで言うバスに対するリクエストとは、PCHにヒットしないリクエストを意味する。また、 $L$ は通常のデータ・リクエストのアドレスがローカル・メモリにある確率、 $V$ は他のPCにコヒーレント制御のためのリクエストを要求する確率を表す。図1からも明らかなように、このSMPをEMCと見るとこれはエルゴード的であるので、式1を使って定常分布を求めることができる。

EMCの遷移確率行列は図1のように表される。先に述べたEMCを用いてSMPを形成するには、各状態の平均滞在時間 $\{\eta_i\}$ が必要である。まずASURAクラスタの仕様から定まるものとしては、表1がある。ただし、各値はプロセッサのサイクル数を表す。状態COMでの滞在時間は、先に述べたように通常計算時間に加えて、プロセッサ内キャッシュに対するアクセス時間も含めたもの

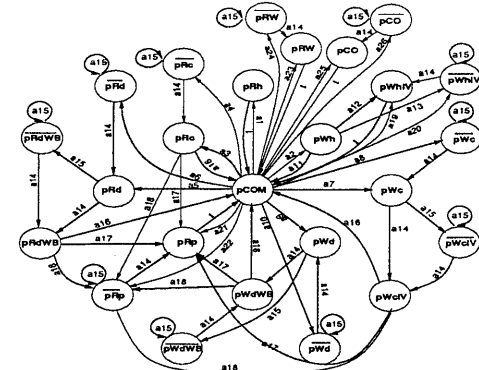


Figure 1 includes a table of transition rates and a matrix of transition probabilities. The table lists parameters like  $a_{11}$  through  $a_{24}$  with their corresponding mathematical expressions. The matrix below shows the transition probabilities between states.

	a11	a12	a13	a14	a15	a16	a17	a18	a19	a20	a21	a22	a23	a24
a11	...	...	...	...	...	...	...	...	...	...	...	...	...	...
a12	...	...	...	...	...	...	...	...	...	...	...	...	...	...
a13	...	...	...	...	...	...	...	...	...	...	...	...	...	...
a14	...	...	...	...	...	...	...	...	...	...	...	...	...	...
a15	...	...	...	...	...	...	...	...	...	...	...	...	...	...
a16	...	...	...	...	...	...	...	...	...	...	...	...	...	...
a17	...	...	...	...	...	...	...	...	...	...	...	...	...	...
a18	...	...	...	...	...	...	...	...	...	...	...	...	...	...
a19	...	...	...	...	...	...	...	...	...	...	...	...	...	...
a20	...	...	...	...	...	...	...	...	...	...	...	...	...	...
a21	...	...	...	...	...	...	...	...	...	...	...	...	...	...
a22	...	...	...	...	...	...	...	...	...	...	...	...	...	...
a23	...	...	...	...	...	...	...	...	...	...	...	...	...	...
a24	...	...	...	...	...	...	...	...	...	...	...	...	...	...

図 1: Pモデルの状態遷移図とそのEMCの遷移確率行列

とするので、20サイクル程度が妥当かと思われる。一方、待ち状態に関してはこのような明確な値は出てこない。なぜなら、通常のデータ・リクエストに加えてキャッシュ・コヒーレンス制御のためのリクエストも各待ち状態に影響を与えるからである。あるいは、PC外部へのリクエストもGモデルの状態によってその遅延時間が動的に変化する。従って、各待ち状態の滞在時間(WTとする)と外部へのリクエスト時間(GRW, CO)は適当な初期値を与え、以後バスへのアクセスを行なう各状態の確率をもとに繰り返して計算することによって収束させることとする<sup>3</sup>。

pRh	pWh	pWhIV	pRc	pRd	pRdWB
10	10	33	121	163	33
pWc	pWcIV	pWd	pWdWB	pRP	pCOM
121	33	163	33	33	20

表 1: 各状態の平均滞在時間

SMPの定常分布は、EMCの定常分布  $\{\pi_i\}$  と滞在時間  $\{\eta_i\}$  を使って式 1から求められる。EM

<sup>3</sup> Pモデルでのリクエストはプライオリティを仮定しない。すなわち、各待ち状態の平均滞在時間は同じと仮定する。

Cの定常分布  $\{\pi_i\}$  は

$$\begin{pmatrix}
 H R (1 - V) (1 - W) (1 - Y) \\
 (1 - W) \{ C H (1 - R) (1 - V) (1 - Y) + (1 - Z) \} \\
 W \{ C H (1 - R) (1 - V) (1 - Y) + (1 - Z) \} \\
 (1 - D) (1 - H) L R (1 - V) (1 - W) (1 - Y) \\
 (1 - D) (1 - H) L R (1 - V) (1 - W) (1 - Y) \\
 D (1 - H) L R (1 - V) (1 - W) (1 - Y) \\
 D (1 - H) L R (1 - V) (1 - W) (1 - Y) \\
 D (1 - H) L R (1 - V) (1 - W) (1 - Y) \\
 (1 - D) (1 - H) L (1 - R) (1 - V) (1 - W) (1 - Y) \\
 (1 - D) (1 - H) L (1 - R) (1 - V) (1 - W) (1 - Y) \\
 D (1 - H) L (1 - R) (1 - V) (1 - W) (1 - Y) \\
 D (1 - H) L (1 - R) (1 - V) (1 - W) (1 - Y) \\
 D (1 - H) L (1 - R) (1 - V) (1 - W) (1 - Y) \\
 (1 - W) \{ D (1 - H) L (1 - V) X (1 - Y) + Y Z \} \\
 W \{ D (1 - H) L (1 - V) X (1 - Y) + Y Z \} \\
 (1 - H) (1 - L) (1 - V) (1 - W) (1 - Y) \\
 V W (1 - Y)
 \end{pmatrix}$$

$$2 - w + (1 - V)(1 - V)CH(1 - R) - HW + (1 - H)L(1 - R + D(R + X)) \quad (3)$$

で求められる。また、各状態からの脱出確率は式 2で求められる。

### 3.2.2 Pモデル中の変数定義

このようにして定められたモデルを評価するには、モデルに対する適当なパラメータが必要である。Pモデルでは静的なパラメータとして、PCHのヒット率H、データ・リクエストに対するリードの割合R、ローカル・メモリに対するリクエストの割合L、全てのデータに対するリード/ライト・データの割合RW.ratio<sup>4</sup>を考慮することとする。また、動的なパラメータとしては、PC外部にコヒーレント制御命令を要求する確率V、他のPCからコヒーレント制御要求が来る確率Y、その時の要求がライトバックである確率Zがあるが、これらは適当な初期値を与え、以後Gモデルとの相互作用で収束させる。Pモデル内部だけで使う動的なパラメータとしては、プロセッサのリクエスト率 $\phi$ がある。また、PモデルでのSMPの状態iの定常分布をP<sub>i</sub>、脱出確率をλ<sub>i</sub>、平均滞在時間をη<sub>i</sub>と表すこととする。

まず、先に上げたパラメータとPC内のプロセッサPを使って、変数D, C, X<sup>5</sup>, W, V, WTを表す。これらは、[14]より、

$$D = (P - 1)H(1 - R)(1 - H(1 - R))^{P-2} \quad (4)$$

$$C = 1 - (1 - HRRW.ratio)^{P-1} \quad (5)$$

$$x = (1 - (HC(1 - R)\phi + (1 - H)(1 - D)(1 - R)\phi))^{P-1} \quad (6)$$

LMおよびGMへのリクエストを出す状態の集合をQ<sub>LM</sub>, Q<sub>GM</sub>とすると、

$$\begin{aligned}
 Q_{LM} &= \{pWcIV, pRc, pRd, pRdWB, pWc, pWcIV, \\
 &\quad pWd, pWdWB, pRp\} \\
 Q_{GM} &= \{pRW, pCO\}
 \end{aligned} \quad (7)$$

あるプロセッサの状態がi ∈ Q<sub>LM</sub> ∪ Q<sub>GM</sub>の時、次のサイクルでも引続き状態を変えない確率はP<sub>i</sub> - λ<sub>i</sub>であるから、バスが空いていない確率をBUSYと

<sup>4</sup>データをリード・オンリー、ライト・オンリー、リード/ライトと分けた場合、キャッシュ・コヒーレンス制御に関係するのは同期変数等のリード/ライト・データである。

<sup>5</sup>本モデルではキャッシュを使い切った状態を考えているが、キャッシュを使い切っていない状態はX = 0とすることにより容易に記述出来る。

すると、

$$BUSY = (P-1) \left\{ L \sum_{i \in Q_{LM}} (P_i - \lambda_i) + \frac{(1-L) \sum_{i \in Q_{LM}} (P_i - \lambda_i)}{MK} \right\} \quad (9)$$

他のプロセッサからのリクエストとの競合に負ける確率は、

$$1 - \frac{1 - (1 - \varphi)^P}{P\varphi} \quad (10)$$

よって、待ち状態に陥る確率  $W$  は

$$W = BUSY + (1 - BUSY) \left( 1 - \frac{1 - (1 - \varphi)^P}{P\varphi} \right) \quad (11)$$

サイクル毎のプロセッサのリクエスト率  $\varphi$  は次のようにして求める。ある状態から抜け出した時にリクエストを発行する状態は、 $pCOM$  および、待ち状態である。前者はキャッシュ・ミスの場合にリクエストが発行される。よって、式 2 を使って次のようにして求める。

$$\begin{aligned} \varphi = & (1-H)\lambda_{pCOM} + L(\lambda_{pWcIV} + \lambda_{pRc} \\ & + \lambda_{pRd} + \lambda_{pRdWB} + \lambda_{pWc} + \lambda_{pWcIV} \\ & + \lambda_{pWd} + \lambda_{pWdWB} + \lambda_{pRp}) \\ & + (1-L)(\lambda_{pRW} + \lambda_{pCd}) \quad (12) \end{aligned}$$

$WT$  は待ち状態にいる時、パスが次に空くまでの平均時間であるので、

$$WT = \sum_{i \in Q_{LMUQGM}} P_i \eta_i \quad (13)$$

によって求まる。

### 3.3 G モデル

P モデルではプロセッサの状態推移を PCH の動きに注目してモデル化を行なったが、G モデルではプロセッサの代わりに NIF の状態推移を GCH の動きに注目してモデル化を行なうこととする。従って、P モデルでの計算状態の代わりに、G モデルではリクエストがない状態、アイドル状態、を状態推移の中心に考える。

#### 3.3.1 状態の定義

G モデルのモデル化を行なうために、状態定義を行なう。以下は NIF の GCH、共有メモリ、P C 間ネットワークに対する各状態を表している。

$gIDLE$	アイドル状態
$gRh$	GCH からの読み込み
$gWh$	GCH の書き込み
$gWhIV$	$gWh$ により発生するインバリデーション
$gWhIV$	$gWhIV$ への待ち状態
$gRc$	そのライン属性が <i>Dirty</i> ではないグローバル・データの GCH ミス・ヒットによる読み込み
$\overline{gRc}$	$gRc$ への待ち状態
$gRd$	そのライン属性が <i>Dirty</i> であるグローバル・データの GCH ミス・ヒットによる読み込み
$\overline{gRd}$	$gRd$ への待ち状態
$gWc$	そのライン属性が <i>Dirty</i> ではないグローバル・データの GCH ミス・ヒットによる書き込み

$\overline{gWc}$	$gWc$ への待ち状態
$gWcIV$	$gWc$ により発生するインバリデーション、もしくは MLI 保証のため、PC から要求されるインバリデーション
$\overline{gWcIV}$	$gWcIV$ への待ち状態
$gWd$	そのライン属性が <i>Dirty</i> であるグローバル・データの GCH ミス・ヒットによる書き込み
$\overline{gWd}$	$gWd$ への待ち状態
$gWB$	ターティ・ラインへの読み書きにより発生するライトバック
$\overline{gWB}$	$gWB$ への待ち状態
$gRP$	ラインのリプレースにより発生するライトバック
$\overline{gRP}$	$gRP$ への待ち状態

図 2 は定義された状態間の遷移を示す。図 2 において、 $h$  は GCH のヒット率、 $r$  は通常のリクエストに対するリード・リクエストの割合、 $d$  はリクエスト要求のあるデータを含むラインが *Dirty* である確率、 $c$  はリクエストのあるデータを含むラインが自分以外の少なくとも一つのキャッシュで *Valid* である確率、 $w$  は (P C 間ネットワークに対して) リクエストを出した時に待ち状態に陥る確率、 $x$  はキャッシュが既に一杯である確率、 $u$  は P C からコピー制御命令が要求される確率を意味する。この時、各サイクル毎の NIF に対するリクエスト率  $\psi$  が必要となる。このリクエストは P モデルでの状態  $pRW$  から発行される。あるプロセッサから NIF に対して発行されるリクエストは、[13] で考察したように、リクエスト先が P CGM、CGGM、RGM であるので、それらの確率をそれぞれ  $SPCGM$ 、 $SCGGM$ 、 $SRGM$  とする。これら以外のリクエストは、LM に対してなされるため、次式が成り立つ。

$$L + SPCGM + SCGGM + SRGM = 1 \quad (14)$$

図 2 からも明らかなように、この SMP を EM C と見るとこれはエルゴード的であるので、式 1 を使って定常分布を求めることが出来る。また、その推移確率行列は図 2 のように表される。

G モデルでの平均滞在時間は P モデルと違って一意に定まらないものが多い。これは、G モデルの場合はリクエスト先が同一場所ではないため、アクセス時間が異なるからである。

まず、滞在時間が一意に定まるものとしては、ネットワークを使わないグローバル・キャッシュのアクセスがある。 $\eta_{gRh}$ 、 $\eta_{gWh}$  の値としては、ローカル・メモリのアクセス時間と同じ (121 サイクル) とする。次に、ネットワークを使うリクエストが発行された場合、平均滞在時間はメモリのアクセス時間を  $\Gamma$  サイクル、P C 間ネットワークの遅延時間を  $\Delta_i$  サイクル ( $i \in \{PCGM, CGGM, RGM\}$ ) とすると、

$$\frac{\sum_{i \in \{PCGM, CGGM, RGM\}} S_i \Delta_i}{1-L} + \Gamma \quad (15)$$

となる。例えば  $\Delta_{PCGM} = 5$ 、 $\Delta_{CGGM} = 10$ 、 $\Delta_{RGM} = 50$ 、メモリ・アクセスの場合、 $\Lambda = 256$ 、インバリデーションの場合、 $\Lambda = 20$ 、ライトバックの場合、 $\Lambda = 20$ 、などにする。これらのリクエスト状態に対する待ち時間は P モデル同様、適当な初期値を与えた  $w$  に対して繰り返し計算することにより収束させる。最後に、アイドル状態

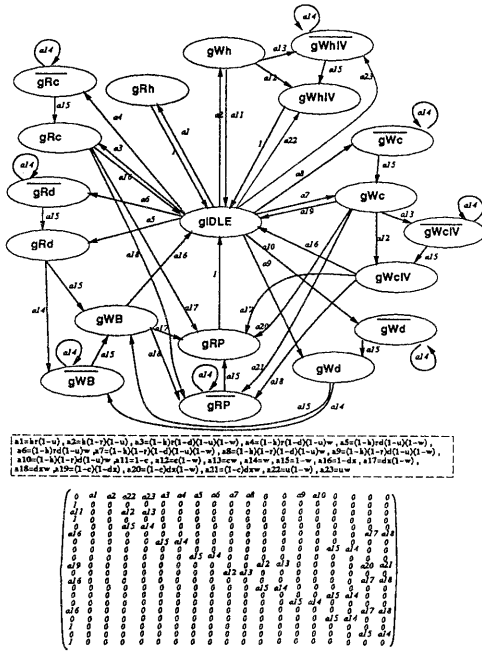


図2: Gモデルの状態遷移図とそのEMCの遷移確率行列

の平均滞在時間 $\eta_{gIDLE}$ は、Pモデルでの計算状態の平均滞在時間に対応するものであるが、これはPモデルの $P_{RW}$ によって変化する。

SMPの定常分布は、EMCの定常分布 $\{\pi_i\}$ と先に定義した滞在時間 $\{\eta_i\}$ を使って式1から求められる。EMCの定常分布 $\{\pi_i\}$ は次でもとめられる。また、各状態からの脱出確率は式2で求められる。

$$\left( \begin{array}{c} \lambda r(1-u) \\ c(h(1-r)(1-u) + u(1-w)) \\ (ch(1-r)(1-u) + u(1-w)) \\ (1-d)(1-h)r(1-u)(1-w) \\ (1-d)(1-h)r(1-u)(1-w) \\ d(1-h)r(1-u)(1-w) \\ c(1-d)(1-h)(1-r)(1-u)(1-w) \\ c(1-d)(1-h)(1-r)(1-u)(1-w) \\ d(1-h)(1-r)(1-u)(1-w) \\ d(1-h)(1-r)(1-u)(1-w) \\ d(1-h)(1-r)(1-u)(1-w) \\ d(1-h)(1-r)(1-u)(1-w) \end{array} \right) \frac{1}{2 + w + (1-u)(c(1-d+dh)(1-r) + d(1-h)(1+x) - h w)} \quad (16)$$

### 3.3.2 Gモデル中の変数定義

Pモデル同様、Gモデルでは静的なパラメータとして、GCHのヒット率 $h$ 、データ・リクエストに対するリードの割合 $r$ 、メモリ階層に対応するリクエストの割合 $SPCGM, SCGGM, SRGM$ 、を考慮することとする。動的なパラメータとしてはPモデル同様 $d, c, x, w, wt$ 、下位のPCから共有データへのリクエストが来る確率 $\psi$ がある。また、GモデルでのSMPの状態 $i$ の定常分布を $P_i$ 、脱出確率

を $\lambda_i$ 、平均滞在時間を $\eta_i$ と表すことにする。  
先に上げたパラメータとを使って、変数 $d, c, x, w, wt$ を表す。これはPモデルの時と同様にして、

$$d = (MK - 1)h(1-r)(1-h(1-r))^{MK-1} \quad (17)$$

$$c = 1 - (1 - hRRW\_ratio)^{MK-2} \quad (18)$$

$$x = (1 - (\psi h(1-r)(1 - (1-c)^{MK-1})))^{MK-1} \quad (19)$$

PC間ネットワークに対してリクエストを出した時、待ち状態に陥る確率はPモデルと同様、ネットワークの状態とアービトレーションによって定義される。まず、ネットワークにリクエストを出す状態の集合を $\Omega$ とする。

$$\Omega = \{gWcIV, gRc, gRd, gRdWB, gWc, gWcIV, gWd, gWdWB, gRp\} \quad (20)$$

あるプロセッサの状態が $i \in \Omega$ の時、リクエスト先がPCGMの場合は待ち状態に陥ることはない。CGGMの場合、リング・ネットワークが空いていない確率を $B_C$ とすると、

$$B_C = \frac{1}{M} \sum_{i \in \Omega} (P_i - \lambda_i) \quad (21)$$

ネットワーク内での競合に負ける確率 $Lost_C$ は、

$$1 - \frac{1 - (1 - \psi)^M}{M\psi} \quad (22)$$

リクエスト先がRGMの場合、クロスバ・ネットワークが空いていない確率を $B_R$ とすると、

$$B_R = \frac{K-1}{K} \sum_{i \in \Omega} (P_i - \lambda_i) \quad (23)$$

ネットワーク内での競合に負ける確率 $Lost_R$ は、

$$1 - \frac{1 - (1 - \frac{\psi}{K})^K}{K\psi} \quad (24)$$

で求まる。よって、待ち状態に陥る確率 $w$ は

$$w = \frac{SCGGM}{1-L} (B_C + (1-B_C)Lost_C) + \frac{SRGM}{1-L} (B_C + (1-B_C)B_R) + (1-B_C)(1-B_R)Lost_R \quad (25)$$

PCからのリクエスト率 $\psi$ は、次のようにして求める。まず、待ち状態を考えない場合の初期値としての $\psi$ は、 $(1-H)(1-L)\varphi$ の確率で各PCから発行されるリクエストであるから、

$$1 - \{1 - (1-H)(1-L)\varphi\}^P \quad (26)$$

となる。以後はPモデルと同様、待ち状態を考慮したリクエスト率 $\psi$ の修正を次のようにおこなう。

$$\bar{\Omega} = \{gWcIV, gRc, gRd, gRdWB, gWc, gWcIV, gWd, gWdWB, gRp\} \quad (27)$$

$$\psi = (1-h)\lambda_{gIDLE} + \sum_{j \in \bar{\Omega}} \lambda_j \quad (28)$$

$w$ は待ち状態にいる時、ネットワークが次に空くまでの平均時間であるので、

$$wt = \sum_{i \in \Omega} P_i \eta_i \quad (29)$$

によって求まる。

### 3.4 PモデルとGモデルのI/F

MLI属性の保証をするためのPモデルとGモデルとの入出力に関わる変数 $V, Y, Z, u$ および、

その相互作用に影響をうける滞在時間 $\eta_{pRW}$ ,  $\eta_{pCO}$ ,  $\eta_{gIDLE}$ の定義を行なう。

まず、PモデルからGモデルへの出力変数を定義する。

他のPCにコヒーレント制御のためのリクエストを要求する確率 $V$ は以下のようにして求める。PCHにおいて、ライト・ヒットもしくはライト・ミスが起こった場合、そのラインを含むGCHのラインがValidである時、他のGCHに対するインバリデーション要求が必要となる。よって、

$$V = (P_{pWcIV} + P_{pWhIV}) * c \quad (30)$$

Gモデルにおいて、下位のPCからコヒーレント制御命令を要求される確率 $u$ は、 $V$ を用いて、

$$u = 1 - (1 - V)^P \quad (31)$$

によって表される。

次に、GモデルからPモデルへの出力変数を定義する。Gモデルにおいてインバリデーションやライトバックが発生した場合、そのラインの一部をキャッシングしているPCHがあると、MLI属性を保証するためには、そのPCHにもコヒーレント制御命令を送らなければならない。

GCHでDirtyなライン、もしくはGCHでValidなラインでも、そのラインをさらにキャッシングしているPCHがDirtyであった場合、インバリデーション要求が起きると、下位のPCHにもインバリデーション要求を起こさなければならない。PCHでDirtyの場合は、インバリデーションの代わりに、ライトバック要求がなされる。ライトバック要求に対しても同様である。

GモデルからMLI属性保証のために、あるPCHにインバリデーションを流す確率は、 $\frac{1}{P} \{ (P_{gWhIV} + P_{gWcIV})(c + dC) + P_{gRP}(c + dC) \}$

$$\text{同じく、リプレース要求を出す確率は、} \quad \frac{1}{P} \{ (P_{gWhIV} + P_{gWcIV} + P_{gRP})dD \} \quad (33)$$

よって、他のPCからコヒーレント制御要求が来る確率 $Y$ は、式32と式33を足し合わせたものになる。また、その時の要求がライトバックである確率 $Z$ は、 $Y$ に対する式33の割合である。

最後に、残された滞在時間の定義を行なう。Pモデルでの状態、 $GRW$ は、Gモデルでの通常のデータ・アクセスに要する時間であるから、次のように定義できる。

$$\begin{aligned} GRW = & hR\eta_{gRh} + \\ & h(1-R)\{\eta_{gWh} + c(\eta_{gWhIV} + w\eta_{gWhIV})\} + \\ & hR(1-d)\{\eta_{gRc} + w\eta_{gRc} + dx(\eta_{gRP} + w\eta_{gRP})\} + \\ & hRd\{\eta_{gRa} + w\eta_{gRa} + \eta_{gWB} + w\eta_{gWB} + \\ & \quad dx(\eta_{gRP} + w\eta_{gRP})\} \\ & h(1-R)(1-d)\{\eta_{gWc} + w\eta_{gWc} + \\ & \quad c(\eta_{gWcIV} + w\eta_{gWcIV}) + dx(\eta_{gRP} + w\eta_{gRP})\} \\ & h(1-R)d\{\eta_{gWd} + w\eta_{gWd} + \eta_{gWB} + w\eta_{gWB} + \\ & \quad dx(\eta_{gRP} + w\eta_{gRP})\} \end{aligned} \quad (34)$$

同様に、 $CO$ は次式で求められる。

$$CO = \eta_{gWhIV} + w\eta_{gWhIV} \quad (35)$$

Gモデルでの $gIDLE$ の滞在時間は、リクエスト率からも求まるであろうが、ここでは次のようにして求めた。

$$\eta_{gIDLE} = \frac{P_{gIDLE} \sum \pi_i \eta_i}{(1 - P_{gIDLE}) \pi_{gIDLE}} \quad (36)$$

### 3.5 モデルの計算手順

このようにして定義したモデルを実際に計算する手順は以下の通りである。

1. パラメータとして、PCH、GCHのヒット率、リード・リクエストの割合を与える。
2. サイクル毎のプロセッサのリクエスト率とPモデルおよびGモデルでの平均待ち時間に適当な初期値を与える。
3. PモデルのEMCの定常分布 $\{\pi_i\}$ を求める
4. Pモデルの平均滞在時間 $\{\eta_i\}$ を求める
5. 各状態からの脱出確率 $\{\lambda_i\}$ を求める
6. 動的な変数 $(W, X, C, D)$ のアップデート。
7. 前回のリクエスト率と新しく求めたリクエスト率の差が適当な値よりも大きければ3に戻る
8. GモデルのEMCの定常分布 $\{\pi_i\}$ を求める
9. Gモデルの平均滞在時間 $\{\eta_i\}$ を求める
10. 各状態からの脱出確率 $\{\lambda_i\}$ を求める
11. 動的な変数 $(w, x, c, d)$ のアップデート。
12. GモデルからPモデルへの入力パラメータ $(Y, Z, V)$ の計算。
13. 前回のリクエスト率と新しく求めたリクエスト率の差が適当な値よりも大きければ8に戻る
14. 以上の操作を適当な回数繰り返す

## 4 評価

SMPを用いたASURAのモデルに対する評価対象として、プロセッサ利用率、通常及びコヒーレント制御のリクエスト時間の割合、それらの待ち時間の割合等を考える。これらの計算手順については文献[14]を参照されたい。

ASURAの評価に対するパラメータとして、 $r = 0.7, 0.7 < h < 1.0$ としたものの結果をグラフ3、4に示す。ただし、他のパラメータは、 $L = 0.4, S_{PCGM} = 0.3, S_{CGGM} = 0.2, S_{RGM} = 0.1, RW\_Ratio = 0.1$ である。また、PCHとGCHのヒット率、リードの割合は同じ値を使用した。

図3では、各プロセッサが、PCにおいてキャッシュ・ヒット率とどのような関係を持つかを示している。プロセッサの利用率からすると、キャッシュ・ヒット率が90%以上なければ、さほど性能が発揮されず、主にメモリ・リクエストとコヒーレンス制御に時間を費やされることがわかる。

図4では、各PCが、ASURA全体におけるキャッシュ・ヒット率とどのような関係を持つかを示している。PCは絶対数が多いため、キャッシュ・ヒット率が高くなってもコヒーレンス制御命令が多発するため、アイドル時間がほとんど増えないことがわかる。

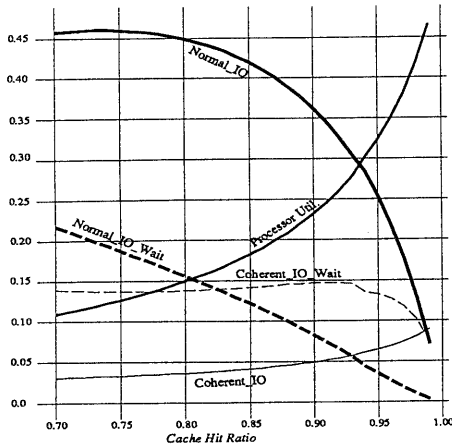


図 3: P-Model Performance under  $R = 0.7, 0.7 < H < 1.0$

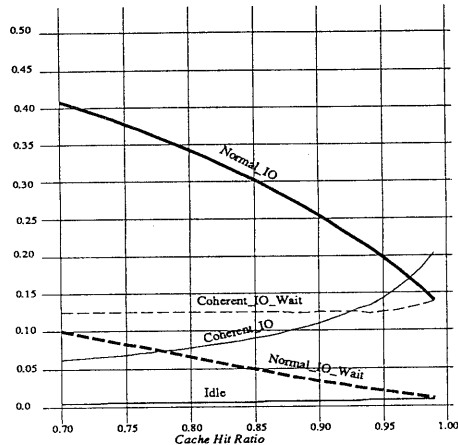


図 4: G-Model Performance under  $R = 0.7, 0.7 < H < 1.0$

## 5 結論

分散共有メモリ型並列計算機ASURAに対する解析モデルを提案し、それを実際に計算した結果、クラスタ内ではクラスタ内ネットワークの負荷を軽くするために高キャッシュ・ヒット率が、クラスタ間では coherence 制御命令を減らすことが、全体の性能向上に必要であることが判明した。

### 謝辞

日頃ご討論頂く京都大学工学部富田教授並びに同研究室の諸氏に感謝致します。また、研究の機会を与えて頂いた(株)クボタ山口部長並びに名古屋大学工学部阿草教授に感謝いたします。最後に、ここまでプロジェクトを推進してくれたのは

ASURAプロジェクト・チーム諸氏の努力の賜であり、諸氏に敬意を表します。

## 参考文献

- [1] A. Agarwal, Richard Simoni, John Hennessy, and Mark Horowitz. An evaluation of directory schemes for cache coherence. In *Proceedings of the International Symposium on Computer Architecture*, pages 280-289, 1988.
- [2] Santosh G. Abraham and Edward S. Davidson. A communication model for optimizing hierarchical multiprocessor systems. In *Proceedings of the 1986 International Conference on Parallel Processing*, pages 467-474, 1986.
- [3] Dharma P. Agrawal and Imadeldin O. Mahgoub. Performance analysis of cluster-based supersystems. In *Proceedings of the International Conference on Supercomputer System*, pages 593-602, 1985.
- [4] Laxmi N. Bhuyan. A combinatorial analysis of multibus multiprocessors. In *Proceedings of the 1984 International Conference on Parallel Processing*, pages 225-227, 1984.
- [5] David R. Cheriton and Hendrik A. Goosen. Paradigm: A highly scalable shared-memory multicomputer architecture. *IEEE Computer*, pages 33-46, 1991.
- [6] Chita R. Das and Laxmi N. Bhuyan. Bandwidth availability of multiple-bus multiprocessors. *IEEE Transactions on Computers*, 34(10):918-926, 1985.
- [7] Daniel E. Lenoski. The design and analysis of dash: A scalable directory-based multiprocessor. Technical Report CSL-TR-92-507, Stanford University, CSL, 1992.
- [8] David J. Kuck, Edward S. Davidson, Duncan H. Lawrie, and Ahmed H. Sameh. Parallel supercomputing today and the cedar approach. *Science*, 231(2):967-974, 1986.
- [9] Imadeldin O. Mahgoub and A. K. Elmagarmid. Performance analysis of a generalized class of m-level hierarchical multiprocessor systems. *IEEE Transactions on Parallel and Distributed Systems*, 3(2):129-138, 1992.
- [10] T. N. Mudge, J. P. Hayes, and D. C. Winsor. Multiple bus architectures. *IEEE Computer Magazine*, 20(6):42-48, 1987.
- [11] S.M. Ross. *Applied Probability Models with Optimization Applications*. Holden-Day, 1970.
- [12] 内藤 潤, 城和貴, 松野 宏昭, and 新田 博之. ASURA クラスタの性能評価. Technical Report 92-ARC-97-10, 情報処理学会計算機アーキテクチャ研究会, 1992.
- [13] 城和貴, 柳原 守, David Fraser, 田中 高士, 新田 博之, 森 眞一郎, 齋藤 秀樹, and 富田 眞治. 分散共有メモリ型マルチプロセッサ「ASURA」の階層性とその評価. Technical Report 92-ARC-95-1, 情報処理学会, 1992.
- [14] 城和貴 and 内藤 潤. セミ・マルコフ過程を用いた ASURA クラスタのモデル化. Technical Report 92-ARC-97-9, 情報処理学会, 1992.
- [15] 森 眞一郎, 齋藤 秀樹, 五島 正裕, 富田 眞治, 田中 高士, David Fraser, 城和貴, and 新田 博之. 分散共有メモリ型マルチプロセッサ「阿修羅」の概要. Technical Report 92-ARC-94-6, 情報処理学会, 1992.