

## セミ・マルコフ過程を用いたASURAクラスタのモデル化

城 和 貴      内 藤 潤

(株)クボタ  
コンピュータ事業推進室

### Abstract

分散共有メモリ型マルチプロセッサASURAの構成要素であるクラスタに着目し、その基本性能について解析する。モデル化の目的は、ASURAクラスタのキャッシュ・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。このモデル化にあたっては、待ち状態の記述が容易であるセミ・マルコフ過程を利用した。さらに、構築されたモデルを用いて、プロセッサ利用率、バンド幅、リクエスト及びその待ち時間について実際の評価を行なった。

## An Analytic Model for the Performance of the ASURA Cluster using a Semi-Markov Processing

Kazuki JOE      Jun NAITO

Office of Computer Business, KUBOTA Corporation  
ASTEM RI  
17 Chudoji, Minami-machi, Shimogyo-ku, Kyoto 600 Japan

*E-mail: {joe, naito}@kocb.astem.or.jp*

### Abstract

This paper presents a discrete time model of bus interference for the cache coherence mechanism in the ASURA cluster, which is a basic element of the cluster-based shared memory multiprocessor computer ASURA. 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. The model produces as output the processor utilization, bandwidth, requests and their waiting probabilities. This modeling capability is attained by using a Semi-Markov process.

## 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ティングへの有力な一方式として、商用 [11]、研究用 [4][2][3] 共に研究開発が推進されている。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 [13]、ASURA 全体を確率モデルで表した性能評価について報告している [12]。

[12] では、主に各プロセッサのキャッシュ・ヒット率とクラスタ間のネットワークのトラフィックについて報告しているが、ここではコヒーレンス制御を始めとするキャッシュ制御のためのリクエストは考慮していない。一般に密結合型並列計算機の解析モデルを考えた場合、キャッシュ制御に伴うリクエストを予測することは難しく、例えば時間的な状態記述を行えるモデル等を利用しなければ、モデル自体が複雑なものになってしまう。

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 [7] に着目し、ASURA 全体の詳細な解析モデル化の第一歩として、その構成要素である ASURA クラスタのモデル化を試みた。ASURA クラスタは、それ自身が密結合型の並列計算機であるため、キャッシュ・コヒーレンス制御のためのリクエストがネットワークに対してどのような影響を与えるかを解析することは、一般の密結合型並列計算機の解析としてはもとより、クラスタ・ベース型の大規模並列計算機の構成要素のデータとしても有用である。

本稿では ASURA について簡単に触れた後、セミ・マルコフ過程についての概略、本モデル構築の上での仮定を述べ、状態記述の定義を行い、具体的な計算手順を示した後、簡単な性能評価を行う。

## 2 ASURA の概要

ASURA は (株) クボタと京都大学とが共同で開発している実験システムである [13][9]。図 1 にそのアーキテクチャを示す。ASURA はメモリの階層性に重点を置いたプロセッサ・クラスタ方式の密結合型マルチプロセッサ・システムである。

プロセッシング・エレメント (PE) はオンチップの 1 次キャッシュと 4 MB の 2 次キャッシュを持つ R4000MC である。これらのローカル・キャッシュはスヌープ方式のキャッシュ・コヒーレンス制御によるイリノイ・プロトコル [1] を一部変更したものに従う。

8 個の PE と 256 MB のローカル・メモリ、ネットワーク・インタフェース (NIF) をバスによ

て結合し、1 つの ASURA クラスタを構成する。このローカル・メモリはクラスタ内の PE によって共有される。NIF は 16 MB のグローバル・メモリと 32 MB のグローバル・キャッシュから構成される。グローバル・メモリは ASURA システムの分散共有メモリとして全ての PE によって共有される。グローバル・キャッシュは ASURA クラスタ内の PE で共有され、フルマップ・ディレクトリ方式のキャッシュ・コヒーレンス制御によるシナプス・プロトコル [1] に従う。また、グローバル・キャッシュは上位のキャッシュとの間に MLI 特性を持つため [13] 1 次及び 2 次キャッシュはグローバル・キャッシュをキャッシング出来る。

複数の ASURA クラスタを階層型のネットワークによって結合することにより ASURA システムが形成される。現段階では、4 つの ASURA クラスタをレジスタ・インサージョン・リング [8] を一部変更した方式でリング結合しクラスタ・グループを形成する。さらに 32 クラスタ・グループをクロスバ結合することにより、最大 1024 CPU からなるシステムを構成することが出来る。

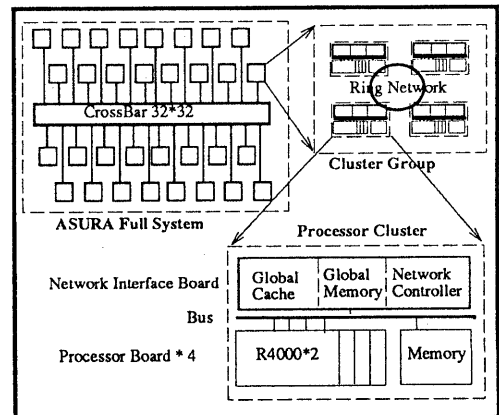


図 1: ASURA のアーキテクチャ

## 3 準備

### 3.1 セミ・マルコフ過程

並列計算機性能評価のモデル化にマルコフ連鎖を利用したものが報告されている [5][6]。マルコフ連鎖による並列計算機システムのモデル化は、モデルの簡潔さという利点がある反面、各状態の時間が一定であるという制約があるため、プロセッサからのリクエストに対する待ち状態を記述しにくいという欠点がある。これに対し、セミ・マル

コフ・モデルは各状態の時間が任意に設定出来るため、キャッシュ・コヒーレンス制御を含む複雑なリクエストが発生しうる共有メモリ型並列計算機システムのモデル化に適するものと思われる。

セミ・マルコフ過程（以後SMPと呼ぶ）に関する詳細は文献[7]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>1</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。

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

### 3.2 モデルの仮定

性能評価のモデル化のために、ASURAクラスタの構造を次のように定義する。N個のプロセッサと同数のキャッシュ・メモリ、1つのメモリ・モジュールがバス結合され、スヌープ方式のキャッシュ・コヒーレンス制御を行う。コヒーレンス・プロトコルはイリノイ型[1]を一部変更したものに従う。なお、ASURAクラスタの詳細なアーキテクチャについては文献[10]を参照されたい。

本稿ではASURAクラスタのキャッシュ・コヒーレンス制御を考慮に入れた性能を評価するために、プロセッサの状態遷移をSMPによって表し、これをエンベデッド・マルコフ連鎖と見なすことにより、セミ・マルコフ過程を定義し、メモリ・リクエストの待ち時間を考慮したモデルを形成する。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とす。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. プロセッサの状態は、計算を行なっている状態、バス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と考える。

<sup>1</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い。

3. プロセッサ内キャッシュにヒットした状態は計算状態と解釈する。
4. 2次キャッシュ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、その間に同時に計算も行なっているとす。
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るものと解釈する。ただし、インバリデーシンの様に、複数のプロセッサに影響を与えるものに関しては、便宜上、リクエストの確率を引き上げることで近似する。
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。

## 4 モデル化

### 4.1 状態の定義

SMPを用いたASURAクラスタのモデル化を行なうために、プロセッサの状態定義を行なう。プロセッサの状態を大別して

- 1) 計算
- 2) 通常のリード/ライト・リクエスト
- 3) 2)に伴う待ち状態
- 4) 2)に付随するキャッシュ制御リクエスト
- 5) 4)に伴う待ち状態

と考える。ASURAクラスタのプロセッサは2で示したように2種類のキャッシュを持つが、本モデルにおいてはチップ内キャッシュに対するリクエストは状態1)に含まれるものとする。従って、各状態においてのリード/ライトは、2次キャッシュもしくはバス(メモリ)に対するリクエストを意味する。以上のことを踏まえて本モデルでは次のような状態定義を行なった。

COM	計算状態、もしくはプロセッサ内キャッシュのアクセス
Rh	キャッシュからの読み込み
Wh	キャッシュへの書き込み
Wh I V	Whにより発生するインバリデーシオン
Wh I V	Wh I Vへの待ち状態
Rc	そのライン属性がダーティではないデータのキャッシュ・ミスによる読み込み
Rc	Rcへの待ち状態
Rd	そのライン属性がダーティであるデータのキャッシュ・ミスによる読み込み
Rd	Rdへの待ち状態

$\overline{RdWB}$	Rdにより発生するライトバック
$\overline{RdWB}$	$\overline{RdWB}$ への待ち状態
$\overline{Wc}$	そのライン属性がダーティではないデータのキャッシュ・ミスによる書き込み
$\overline{Wc}$	$\overline{Wc}$ への待ち状態
$\overline{WcIV}$	$\overline{Wc}$ により発生するインバリデーション
$\overline{WcIV}$	$\overline{WcIV}$ への待ち状態
$\overline{Wd}$	そのライン属性がダーティであるデータのキャッシュ・ミスによる書き込み
$\overline{Wd}$	$\overline{Wd}$ への待ち状態
$\overline{WdWB}$	$\overline{Wd}$ により発生するライトバック
$\overline{WdWB}$	$\overline{WdWB}$ への待ち状態
$\overline{Rp}$	ラインのリプレースにより発生するライトバック
$\overline{Rp}$	$\overline{Rp}$ への待ち状態

このように定められたSMPはASURAクラスタ内の1つのプロセッサの状態記述をしており、全体のモデル化はN個の等価なSMPで表す。このため、本来他のプロセッサに対してライトバックやインバリデーションに伴うライン・リプレースを要求する代わりに、要求したプロセッサ自身がその動作を行なうものとする。これは3.2の仮定に従うものである。

#### 4.2 状態遷移図

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(2次)キャッシュのヒット率、 $r$ 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ダーティである確率、 $c$ はリクエスト要求のあるデータを含むラインがクリーンである確率、 $w$ 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(この時、各サイクル毎のリクエスト率 $R$ が必要となる)、 $x$ はキャッシュが既に一杯である確率を意味する。図2からも明らかのように、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出来る。

状態COMから遷移するのは(プロセッサ内キャッシュ以外への)リクエストが発生した場合のみである。このリクエストが(2次)キャッシュ・ヒットするかしないか、キャッシュ・ミスした場合にそのラインがダーティか否か、バスが空いているかどうか、リクエストはリードかライトか、によってそれぞれ $\overline{Rh}$ 、 $\overline{Rc}$ 、 $\overline{Rd}$ 、 $\overline{Rc}$ 、 $\overline{Rd}$ 、 $\overline{Wh}$ 、 $\overline{Wc}$ 、 $\overline{Wd}$ 、 $\overline{Wc}$ 、 $\overline{Wd}$ に遷移する。 $\overline{Rh}$ 、 $\overline{WhIV}$ 、 $\overline{Rp}$ からはCOMにのみ遷移する。 $\overline{Wh}$ からはライト・ヒットしたラインがクリーンであった場合

はインバリデーションを流さなければならないので、バスの状態によって $\overline{WhIV}$ もしくは $\overline{WhIV}$ に、クリーンでない場合はCOMに遷移する。 $\overline{WhIV}$ からは直接COMに遷移する。なぜならインバリデーションを受けとるキャッシュでの当該ラインはクリーンであるため、バスへのリクエストは起こらないからである。 $\overline{Rc}$ からはキャッシュに空きがあったり、空きがなくてもリプレースされるラインがダーティでなければ、バスへのリクエストは起こらずCOMに遷移するが、キャッシュに空きがなく、かつリプレースされるラインがダーティである場合は、当該ラインのライトバックを行なうためにバスの状態によって、 $\overline{Rp}$ もしくは $\overline{Rp}$ に遷移する。 $\overline{Rd}$ からは当該ラインのライトバックを行なうために、バスの状態によって $\overline{RdWB}$ 、もしくは $\overline{RdWB}$ に遷移する。なお、先に述べたようにここでのライトバックは実際は他のキャッシュでの動作であるが、モデルを簡単にするために、ライトバックは自分のキャッシュで行なうこととする。 $\overline{RdWB}$ からは $\overline{Rc}$ のと同様、COM、もしくはバスの状態によって $\overline{Rp}$ か $\overline{Rp}$ に遷移する。 $\overline{Wc}$ からはインバリデーションを流さなければならないので<sup>2</sup>、バスの状態によって $\overline{WcIV}$ もしくは $\overline{WcIV}$ に遷移する。 $\overline{WcIV}$ からは $\overline{Rc}$ のと同様、COM、もしくはバスの状態によって $\overline{Rp}$ か $\overline{Rp}$ に遷移する。 $\overline{Wd}$ からは当該ラインのライトバックを行なうために、バスの状態によって $\overline{WdWB}$ もしくは $\overline{RdWB}$ に遷移する。ここでのライトバックも $\overline{RdWB}$ のと同様、自分のキャッシュが行なうことを仮定している。 $\overline{WdWB}$ からは $\overline{Rc}$ のと同様、COM、もしくはバスの状態によって $\overline{Rp}$ か $\overline{Rp}$ に遷移する。待ち状態 $\overline{WhIV}$ 、 $\overline{Rc}$ 、 $\overline{Rd}$ 、 $\overline{RdWB}$ 、 $\overline{Wc}$ 、 $\overline{WcIV}$ 、 $\overline{Wd}$ 、 $\overline{WdWB}$ 、 $\overline{Rp}$ はバスの状態に応じてそれぞれ $\overline{WhIV}$ 、 $\overline{Rc}$ 、 $\overline{Rd}$ 、 $\overline{RdWB}$ 、 $\overline{Wc}$ 、 $\overline{WcIV}$ 、 $\overline{Wd}$ 、 $\overline{WdWB}$ 、 $\overline{Rp}$ に遷移する。EMCの遷移確率行列は上記以外の遷移の確率を0として、図3のように表される。

#### 4.3 滞在時間

$\overline{Rh}$	$\overline{Wh}$	$\overline{WhIV}$	$\overline{Rc}$	$\overline{Rd}$	$\overline{RdWB}$
10	10	33	121	163	33
$\overline{Wc}$	$\overline{WcIV}$	$\overline{Wd}$	$\overline{WdWB}$	$\overline{Rp}$	
121	33	163	33	33	

表1: 各状態の平均滞在時間

先に述べたEMCを用いてSMPを形成するに

<sup>2</sup>イリノイ型と違ってクリーン・エクスクリューシブの状態がないのでダーティでないなら必ずインバリデーションを流すことになる。



あるプロセッサがあるラインをダーティで持っている確率は  $h(1-r)$  だから、自分以外のあるプロセッサがダーティで持ちそれ以外のプロセッサがダーティで持っていない確率  $d$  は次式で表される。

$$d = (N-1)h(1-r)(1-h(1-r))^{N-1} \quad (5)$$

同様に、自分以外の少なくとも一つのプロセッサがあるラインをリード・ヒットでキャッシュに持ちそのラインがリードかつライトの属性を持つ確率  $c$  は次式で表される。

$$c = 1 - (1-hrRW\_ratio)^{N-1} \quad (6)$$

キャッシュ・インしたときに既にキャッシュに空きがない状態は次のようにして求められる。キャッシュを使い切った時には、他からのインバリデーションが来ないとキャッシュに空きは出来ない。よって、特定のキャッシュが他からインバリデーションを受ける確率を  $inv$  とすると、キャッシュに空きがない確率は、

$$x = (1-inv)^{N-1} \quad (7)$$

となる。(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。) あるプロセッサがインバリデーションを流す確率は、自分がクリーンなラインをキャッシュに持ち、そのラインにライト・アクセスを行なうか、もしくは、自分の持っていないダーティではないラインに対してライトアクセスを行なう場合である。よって、

$$inv = hc(1-r)R + (1-h)(1-d)(1-r)R \quad (8)$$

によって  $inv$  は求まる。

バスに対してリクエストを出した時、待ち状態に陥る確率とは、そもそもバスが空いてなかった場合、バスは空いているが他のプロセッサからのリクエストとの競合に負けた場合とに分けて考えられる。まず、バスが空いていない確率は次のようにして求める。バスへのリクエストを出す状態の集合を  $S$  とする。

$$S = \{WcIV, Rc, Rd, RdWB, Wc, WcIV, Wd, WdWB, Rp\} \quad (9)$$

あるプロセッサの状態が  $i \in S$  の時、次のサイクルでも引続き状態を変えない確率は  $P_i - \lambda_i$  であるから、バスが空いていない確率を  $BUSY$  とすると、

$$BUSY = (N-1) \sum_{i \in S} (P_i - \lambda_i) \quad (10)$$

となる。ただし、 $P_i$  は SMP の状態  $i$  での定常分布である。次に、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、

$$1 - \frac{1 - (1-R)^N}{NR} \quad (11)$$

で求まる。なぜなら、あるサイクルで少なくとも一つのプロセッサがリクエストを出す確率は  $1 - (1-R)^N$ 、その時のリクエスト数は  $NR$  だからである。よって、待ち状態に陥る確率  $w$  は

$$w = BUSY + (1 - BUSY) \left(1 - \frac{1 - (1-R)^N}{NR}\right) \quad (12)$$

となる。

サイクル毎のプロセッサのリクエスト率  $R$  は、次のようにして求める。まず、状態が  $i \in S$  の時、 $\eta_i$  の時間の間に一回だけリクエストが起きる。さらに、状態  $S$  の待ち状態の集合  $\bar{S}$  を

$$\bar{S} = \{\overline{WcIV}, \overline{Rc}, \overline{Rd}, \overline{RdWB}, \overline{Wc}, \overline{WcIV}, \overline{Wd}, \overline{WdWB}, \overline{Rp}\} \quad (13)$$

とすると、状態が  $i \in \bar{S}$  の時はサイクル毎にリクエストを発行していると考えられる。(勿論、実際のインプリメンテーションとは異なる。) よって、サイクル毎のプロセッサのリクエスト率は、

$$R = \sum_{i \in S} \frac{P_i}{\eta_i} + \sum_{i \in \bar{S}} P_j \quad (14)$$

である。

$Wt$  は待ち状態にいる時、バスが次に空くまでの平均時間であるので、

$$Wt = \sum_{i \in S} P_i \eta_i \quad (15)$$

によって求まる。

なお、 $x$ 、 $w$ 、 $R$ 、 $Wt$  はモデルに与えるパラメータだけでは決まらず、後に述べるように、繰り返し計算することにより収束させて求める。

#### 4.6 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、キャッシュ・ヒット率  $h$ 、リード・リクエストの割合  $r$ 、リード/ライト・データの割合  $RW\_ratio$ <sup>5</sup> を定める。
2. サイクル毎のリクエスト率  $R$  と平均待ち時間  $Wt$  に適当な初期値 (例えば、 $R$  に  $1 - \frac{1}{N}$ 、 $Wt$  には  $\eta_i (i \in S)$  の平均) を与える。

<sup>5</sup>本稿ではこの割合を5%と見積もったが、本来は並列化プログラムに左右される値である。

3. EMCの定常分布  $\{\pi_i\}$  を求める (式 4)
4. 平均滞在時間  $\{\eta_i\}$  を求める (表 1と式 15)
5. 各状態からの脱出確率  $\{\lambda_i\}$  を求める (式 2)
6. R (式 14)、x (式 7)、w (式 12)、Wt (式 15) を求める
7. 前回のRと新しく求めたRの差が適当な値よりも大きければ3に戻る

## 5 評価

### 5.1 評価対象

SMPを用いたASURAクラスタのモデルに対する評価対象として、バンド幅、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。

バンド幅 (Bandwidth) の定義を、プロセッサのサイクルあたり、(セカンダリ) キャッシュ及びメモリに対するアクセプトされたリクエスト数とすると、

$$N \sum_{i \in \{Rh, Wh, Rc, Rd, Wc, Wd\}} \frac{P_i}{\eta_i} \quad (16)$$

プロセッサ利用率 (Processor Utilization) は、キャッシュにアクセス中にプロセッサは計算を行なうことができるという仮定 3.2 のもとでは、

$$\sum_{i \in \{COM, Rh, Wh\}} P_i \quad (17)$$

全体に対する通常のメモリ・リクエストに要する時間の割合 (Normal IO request) は、

$$\sum_{i \in \{Rc, Rd, Wc, Wd\}} P_i \quad (18)$$

その待ち時間の割合 (Normal IO wait) は、

$$\sum_{i \in \{Rc, Rd, Wc, Wd\}} P_i \quad (19)$$

同じく、コヒーレンス制御のリクエスト時間の割合 (Coherent IO request) は、

$$\sum_{i \in \{WhIV, RdWB, WcIV, WdWB, Rp\}} P_i \quad (20)$$

その待ち時間の割合 (Coherent IO wait) は、

$$\sum_{i \in \{WhIV, RdWB, WcIV, WdWB, Rp\}} P_i \quad (21)$$

で表される。

### 5.2 評価結果

ASURAクラスタの評価に対するパラメータとして、

1.  $r = 0.8, 0.7 < h < 1.0, RW\_ratio = 0.1$
2.  $h = 0.8, 0.7 < r < 1.0, RW\_ratio = 0.1$
3.  $h = 0.97, 0.7 < r < 1.0, RW\_ratio = 0.1$

としたものの結果をグラフ 4、5、6に示す。

図 4では、キャッシュ・ヒット率と性能の関係を示している。キャッシュ・ヒット率が95%以下の場合にはプロセッサ利用率が極端に悪く、97%程度まではコヒーレント制御リクエストに対する待ち時間が多いため、実際に高性能を発揮出来るのは、ヒット率97%以上である。

図 5では、リード・リクエスト率と性能の関係を示している。まず、全体の傾向を知るためキャッシュ・ヒット率を低目(80%)に設定して評価した結果、リクエストの待ち状態が8割前後にも達しており、リード・リクエスト率の影響を読みとりにくい。

図 6では、望ましいキャッシュ・ヒット率におけるリード・リクエスト率と性能の関係を示している。プロセッサ利用率はリード・リクエスト率が90%を越える付近から急速に向上しているのが分かる。さらに、コヒーレント制御に対する待ち時間は急減少している。

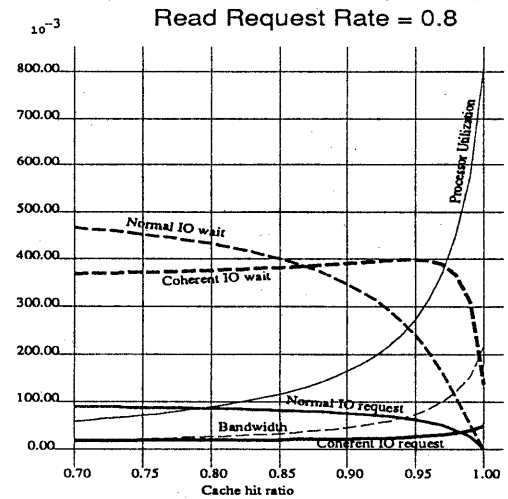


図 4:  $r = 0.8, 0.7 < h < 1.0$  の性能

## 6 結論

分散共有メモリ型並列計算機ASURAの構成要素であるクラスタのキャッシュ・コヒーレント制

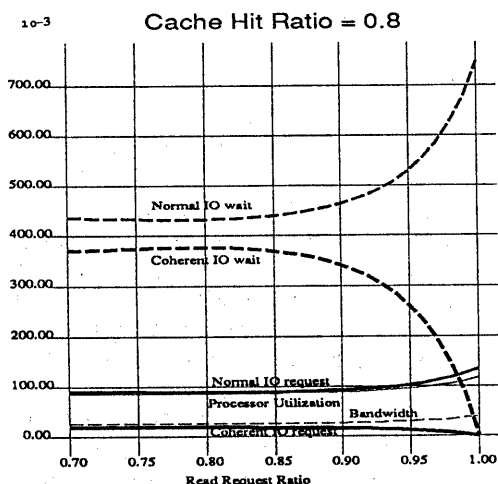


図 5:  $h = 0.8, 0.7 < r < 1.0$  での性能

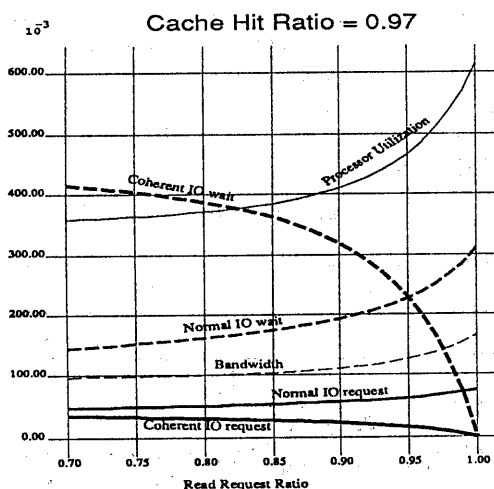


図 6:  $h = 0.97, 0.7 < r < 1.0$  での性能

御を考慮した基本性能を、セミ・マルコフ過程を用いることにより解析モデルで表し、キャッシュのヒット率及びリード/ライト命令の比率がクラスタ全体の性能にどのように影響を与えるかを示した。その結果、全体の性能はキャッシュのヒット率は勿論のこと、リード/ライト命令の比率にも左右されることが確認された。

本モデルはASURA全体の解析モデル構築のための基礎的な部分であり、キャッシュ・コヒーレント制御を考慮に入れた大規模なクラスターベースの並列計算機のモデル化を行なうことが今後の研究課題である。

## 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。

## 参考文献

- [1] James Archibald and Jean-Loup Baer. Cache coherence protocols: Evaluation using a multiprocessor simulation model. *ACM Transactions on Computer Systems*, 4(4):273-298, 1986.
- [2] Daniel Lenoski et al. The directory-based cache coherence protocol for the dash multiprocessor. In *Proceedings of the International Symposium on Computer Architecture*, pages 148-159, 1990.
- [3] David R. Cheriton et al. Paradigm: A highly scalable shared-memory multicomputer architecture. *IEEE Computer*, pages 33-46, 1991.
- [4] David J. Kuck, Edward S. Davidson, and Duncan H. Lawrie Ahmed H. Sameh. Parallel supercomputing today and the cedar approach. *電子情報通信学会論文誌*, J71-D(8):1361-1374, 1988.
- [5] Marco Ajmone Marsan and Mario Gerla. Markov models for multiple bus multiprocessor systems. *IEEE Transactions on Computers*, 31(3):239-248, 1982.
- [6] Ibrahim H. Onyuksel and Keki B. Irani. A markovian queueing network model for performance evaluation of bus-deficient multiprocessor systems. In *Proceedings of the 1983 International Conference on Parallel Processing*, pages 437-439, 1983.
- [7] S.M. Ross. *Applied Probability Models with Optimization Applications*. Holden-Day, 1970.
- [8] Andrew S. Tanenbaum. *COMPUTER NETWORKS*. Prentice Hall, 1981.
- [9] 齋藤 秀樹, 森 眞一郎, 城 和貴, David Fraser, 田中 高士, and 富田 眞治. イベント対応型キャッシュ・コヒーレンス制御方式とそのバリア同期への応用. Technical Report 92-ARC-95-2, 情報処理学会計算機アーキテクチャ研究会, 1992.
- [10] 内藤 潤, 城 和貴, 松野 宏昭, and 新田 博之. ASURA クラスターの性能評価. 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.
- [11] 日本アライアントコンピュータ. *The CAM-PUS/800 Supercomputer*. 技術資料, 1991.
- [12] 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.
- [13] 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.