

## 資源共有型並列計算機“砂丘”

荒川 修・橋本 正巳・井上 倫夫・小林 康浩  
鳥取大学工学部

本研究では、資源共有型マルチマイクロプロセッサシステム“砂丘”を開発している。本システムは、64台の汎用マイクロプロセッサを接続可能で、64Mバイトの容量の主メモリを持つ。並列処理プログラムの同期を容易に実現するために、システムコントローラ（割り込み制御回路）、排他制御回路などのハードウェアを特別に設けている。また、演算結果を分かりやすく表示するためのグラフィックディスプレイ、マンマシンインターフェースとシステムの制御を行うためのフロントエンドプロセッサを設けている。本報告では、フロントエンドプロセッサ、システムエリア及び同期通信方式について報告する。

## The Multiprocessor System "SAKYU" of the Shared-Memory Type

Osamu ARAKAWA · Masami HASHIMOTO · Michio INOUE · Yasuhiro KOBAYASHI  
Faculty of Engineering, Tottori University

Our multiprocessor system, called "SAKYU", consists of 64 microprocessors and has main memories (shared type) of 64MB a total. To achieve parallel processing regularly, various ideas are taken into this system, and are realized as hardware; e.g., System-Controller, Semaphore, Arbiter of memory contention. As the I/O equipment several graphic displays are furnished. Front-End-Processors are provided for practising good man-machine interface and for controlling the whole system easily.

Among them, the function of the Front-End-Processors, the role of System-Area devices and the procedure of the synchronized communication is described in this paper.

## 1. はじめに

筆者らは、実験室レベルで特定ユーザーが利用する数値シミュレーションマシンとして、並列計算機“砂丘”を開発している[8]~[10]。“砂丘”は、複数のマイクロプロセッサ(μP)を密に結合したマルチμPシステムである。一般に、この種のシステムでは、メモリを共有することによる各μPから共有メモリへのアクセスから競合する機会が増し、実質的に携わるμPの数が増えないことになる。μPの接続台数が、数十~百台程度であれば、アクセスバスを分散することによりアクセス競合を緩和ないし回避することが可能である[10]。共有メモリ方式は、データの授受が高速

かつ容易に行えることから、さまざまな処理アルゴリズムに対応できるような利便性のよいシステムの構築が可能である。

本報告では、まず現在開発中の並列計算機“砂丘”のシステム構成を述べ、次に同期通信方式について述べる。またシステムの稼働率とテストプログラムによる実測値について述べる。

## 2. システム構成

図1は、筆者らが製作を進めている並列計算機“砂丘”のシステム構成図である。本システムは、最大64台のプロセッサユニット(PU)が接続可能な密結合型マルチμPシステムである。その

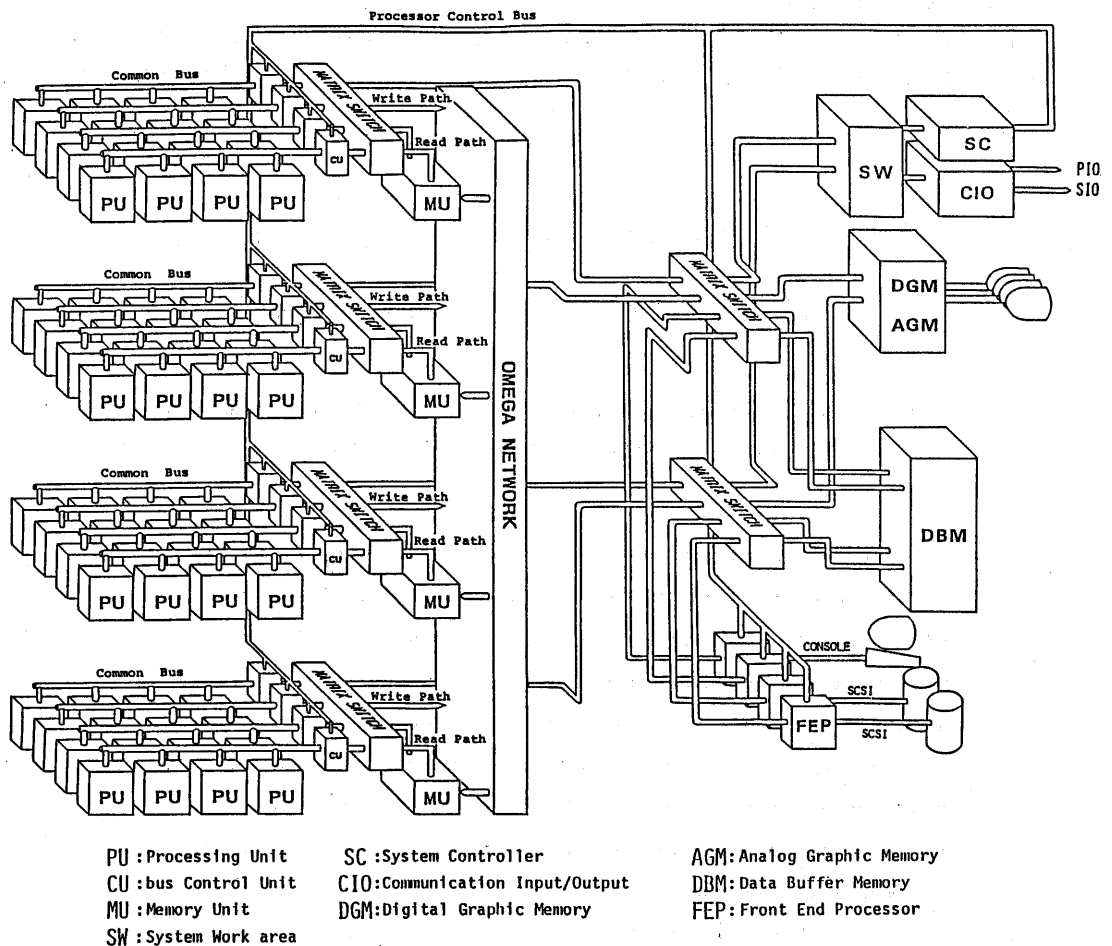


図1 並列計算機“砂丘”のシステム構成

構成はバスコントロールユニット（BCU）と最大4台のPUを1グループとし（単一バスで接続）、それらの4グループをまとめて1ブロックとする。システム全体で4ブロックをマトリックススイッチ（MTX）を介して主メモリ、オメガネットワーク、システムエリアに接続する。MTXは4×4のバススイッチである。4入力、すべてBCUに接続する。4出力は、1出力がシステムエリアのアクセスに、他の3出力が主メモリのアクセスに使用される。主メモリのアクセス（3出力）は、2出力が2インターリーブ構成のリードアクセスバスとして使用され、別の1出力がライトアクセスバスとしてオメガネットワークに接続される。

### 2. 1 主メモリ

主メモリは、合計64Mバイトの容量を持ち、演算データを格納する。通常のプログラムでは、PUの共有メモリアccessは、リードアクセスの方がライトアクセスより多く行われる。このため主メモリを4つのユニットに分割し、各PUブロック毎に置く。各メモリユニットのリードアクセスバスは2インターリーブでMTXの2出力と接続した。このため同じPUブロック内でのみリードアクセス可能である。またライトアクセスバスはオメガネットワークを通して全てのメモリユニットへ接続される。ライトアクセスバスにオメガネットワークを使用したことにより、主メモリは次の様な利用形態が考えられる。

#### (1) マルチリード・ワンライト方式

オメガネットワークの一斉放送モードを用いることにより4つの主メモリユニットの内容を同一に保持することができる。全PUが同じデータで一斉に処理をするときに、見かけ上のメモリ容量は4分の1になるが、同一PUブロックからのリードアクセスの競合のみを考えれば良く、アクセス競合を緩和できる。

#### (2) バイライン接続

あるPUブロックで所定の処理した結果を、隣のPUブロックのメモリユニットにライトアクセスする。そしてライトされたデータに新たな処理

をし、結果をさらに隣のメモリユニットにライトアクセスするということを反復する。

#### (3) 個別接続

それぞれのPUブロックのライトアクセスを、リードアクセスで使用するのと同じメモリユニットに対して実行する。つまり各ブロックが4つのメモリユニットをそれぞれ独立して使用する。

### 2. 2 フロントエンドプロセッサ（FEP）

PUが演算専用なのに対して、FEPはマンマシンインターフェースやシステムの制御の部分を受け持つ。FEPにも複数のプロセッサを用意してそれぞれの機能を分散させる。マンマシンインターフェースを受け持つFEPは、ユーザーからの要求の受け付けやプロセスの制御を行う。入出力装置を制御するFEPは、プリンタなどのキャラクタ型デバイスと、ハードディスクなどのブロック型デバイス向けに個別に用意する。

### 2. 3 システムエリア

2段目MTXの出力に接続されているメモリ領域をシステムエリアと呼ぶ。このエリアは大きく分けてグラフィックメモリ、システムワークエリアから成る。2段目MTXの入力には、PU側からは1段目MTXの出力と接続し、FEPもここに接続する。このエリアは全てのプロセッサからアクセス可能である。

#### (1) グラフィックメモリ

グラフィックディスプレイが膨大な計算結果を分かりやすく図示するために用意される。デジタルグラフィックは主に計算結果をグラフなど2次元的に図示するために使い、アナロググラフィックは主に色の違いや濃淡で3次元的なものを表すのに用いる。

#### (3) システムワークエリア

このエリアには、並列処理プログラムを効率よく実行するためのハードウェアとして、システムコントローラ、メールボックス、排他制御回路などを置く。また、他の機器との入出力を行うためのコミュニケーションI/Oが設けてある。

### 3. 同期・通信方式

マルチプロセッサシステムでは、複数のプロセス相互の協調をとりながら、まとまった仕事を進行させなければならない。そのためには、プロセスの実行を調整して正しい結果となるよう実行順序を保障しなければならない。それには割当てられたプロセスの同期と、互いの進行状況や中間結果を他のプロセスに通報することである。

本システムではこれら同期・通信を効率よく行うための専用回路として、システムコントローラ、排他制御回路、プロセス番号割当回路を設け同期・通信操作に伴うオーバーヘッドを小さくするようにした。

#### (1) システムコントローラ

システムコントローラとして、割り込み制御回路を設けており、全ての $\mu P$ 群に接続している。各 $\mu P$ は、この制御回路及びシステムワークエリア内のメールボックスを用いることにより、プロセス起動、 $\mu P$ 間の通信、同期等を行うことが可能になる。また、上記の操作を行う $\mu P$ は、制御回路内に設けた同期フラグを管理することにより、処理の終了を確認することができる。

#### (2) 排他制御回路

排他制御を実現するために、共有変数に対して行う不可分のリード・モディファイ・ライト動作をハードウェアで実現した。共有エリアの一部を排他制御用メモリとして確保し、このメモリに対してはリードとライト動作を不可分的に実行する。具体的には、 $\mu P$ のリード時に出力データのMSBの値に応じて、回路内の制御部で自動的にライトサイクルが実行され、排他制御が行われる。

#### (3) プロセス番号割当回路

多数のプロセスを限られた数の $\mu P$ で効率よく並列処理するには、各 $\mu P$ の処理するプロセスの番号管理が必要である。本システムではこの機能をハードウェア化している。共有変数を順序づけし、その番号をカウンターで管理し、 $\mu P$ からのリード動作が終わると自動的に指示値を増すようにした。これにより $\mu P$ は共有変数に対して不可分のリード・モディファイ・ライト動作を行わずに、リード動作だけで番号管理ができる。

### 4. 検討

資源共有型のマルチ $\mu P$ システムでは、原理的に共有バス上でアクセス競合が発生する。したがって、各 $\mu P$ の共有資源アクセス要求を能率よく制御し、競合による待ち時間の累積を如何に抑制するかがシステムを効率よく運用する上で重要な課題である。

本システムで扱う問題は、特に限定せずに科学技術計算全般を予定している。例えば、ある事象をモデル化し、偏微分方程式による数値シミュレーションを行ったり、一次元・二次元のFFTを行うような場合には、大きな配列データを主メモリ中に置くことになる。各 $\mu P$ はこの配列からデータを読み込み、演算を行った後に、配列のデータを更新することを繰り返す。このような計算を行う場合、各 $\mu P$ が分担する作業は、ほぼ均等に割りつけられて（アルゴリズム上では最適）いる。したがって、アクセス競合が頻繁に起こり易く、起これば全 $\mu P$ にアクセス待ちが生じ、処理速度の向上は望めない。アクセス競合の機会を減らすため、主メモリにはマルチリード・ワンライトメモリ方式、2インターリーブ方式、2ポート方式を採用している。

$\mu P$ が共有メモリをアクセスするとき、共有メモリ本来のアクセス時間、共有バスでの信号遅延時間及びアクセス競合調停のための待ち時間が必要となる。 $\mu P$ の持つ能力をどれくらい利用できたかを表すために、ここでは、稼働率を純粋に計算に専念できた時間の割合として定義し、検討を述べる。

#### 4. 1 並列動作時の各 $\mu P$ の稼働率

筆者らは、各 $\mu P$ の待ち時間の累積が最大となる（均等負荷で動作している）ときのシステムの性能を陽的に解析できることを示した[6]。そこでシステムのハードウェア特性（共有メモリのアクセス時間、サイクル時間等）及びソフトウェア特性（タスクの処理時間、共有メモリのアクセス回数）とを用いて、並列動作時の各 $\mu P$ の平均稼働率を次のように表すことができる。

$$P(n, Th) = \frac{100}{1 + \frac{t_{ac} + m_n \cdot t_s}{Th}} \quad (1)$$

ここで、

$$m_n = \begin{cases} 0 & (n \leq i_B \cdot n_0) \\ \frac{n}{i_B} - n_0 & (n > i_B \cdot n_0) \end{cases} \quad (2)$$

ただし、

$n$  : 同時に動作している  $\mu P$  の台数  
 $Th$  : タスクの平均シンクタイム [S]  
 ( $\mu P$  が共有メモリをアクセスする  
 平均時間間隔)

$t_s$  : 共有資源のサイクル時間 [S]

$t_{ac}$  : 共有資源をアクセスするとき各  $\mu P$  が  
 要する平均アクセス時間

$$t_{ac} = t_B + t_A + \frac{1}{2} t_s \quad (3)$$

$t_A$  : 共有資源のアクセス時間 [S]

$t_B$  : BCU, MTX等での遅延時間 [S]

$i_B$  : 共有バスのインターリーブ数

$n_0$  : 同一の共有バスでアクセス競合による  
 待ち時間の累積なしに動作できる  $\mu P$   
 の台数

$$n_0 = 1 + \frac{Th}{t_s} \quad (4)$$

#### 4. 2 マルチリード・ワンライトメモリ方式 による稼働率

マルチリード・ワンライトメモリ方式では、並列稼働が可能な  $\mu P$  の台数及び稼働率は、主メモリへのリードとライトのアクセス比に依存する。主メモリへのアクセス全体を1としたときのリードアクセスの割合を  $r$  とすれば、(4)式から、リードアクセスにおいて競合による待ち時間の累積なしに並列稼働が可能な  $\mu P$  の台数は、主メモリを4ユニット設け、それぞれのリードバスを2インターリーブとしたので、

$$n_R = 8 \left( 1 + \frac{Th/r}{t_s} \right) / r \quad (5)$$

と表せる。ライトアクセスにおける台数は、

$$n_W = \left( 1 + \frac{Th/(1-r)}{t_s} \right) / (1-r) \quad (6)$$

と表せる。(5)(6)式の関係を図2に示す。

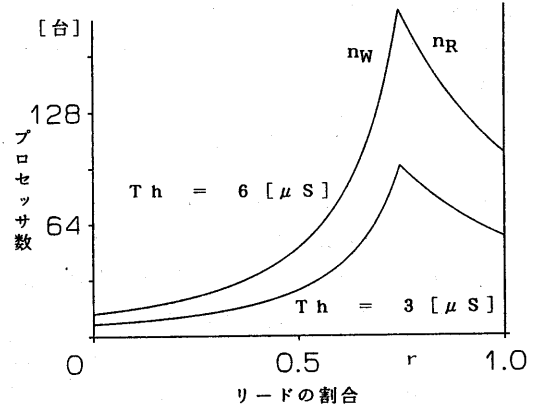


図2 メインメモリ利用時のリード・ライト比と競合を起こさずに稼働できる  $\mu P$  数

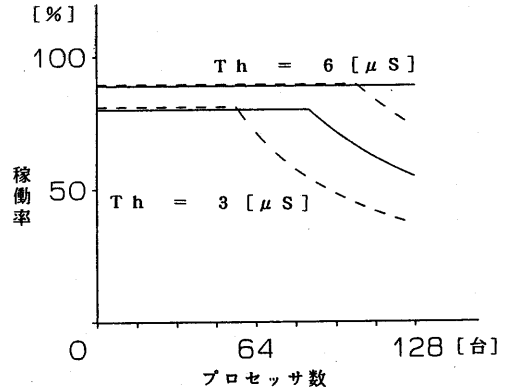


図3 メインメモリ使用時の  $\mu P$  数と稼働率

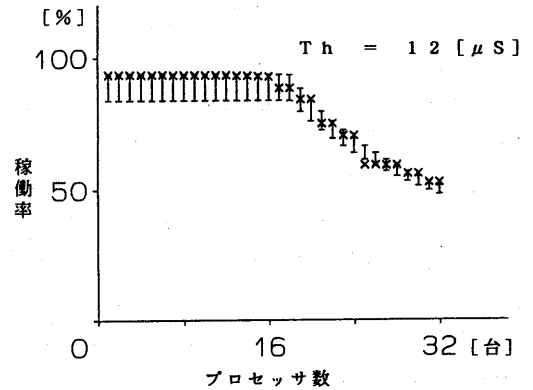


図4 並列計算機“砂丘”のシステムメモリにおける稼働率の実測値と理論値

この図から  $r = 0.75$  付近であるとき、最もアクセス競合が起こりにくいことがわかる。また稼働率を表す式は、

$$P(n) = \frac{100}{1 + \frac{t_{ac} + m \cdot t_a}{2h/r} + \frac{t_{ac}' + m' \cdot t_a}{2h/(1-r)}} \quad (7)$$

と表せる。(7)式において、 $r = 0.8$ のときの  $\mu P$ の台数と稼働率の関係を図3に示す。破線は共有メモリを8インターリーブ方式で構成した場合の稼働率である。この図より接続する  $\mu P$ の台数が同じ場合、マルチリード・ワンライトメモリ方式は8インターリーブ方式に比べてシンクタイムをより短くできることがわかる。

#### 4. 3 稼働率の実測値

32台の  $\mu P$ を16台のBCUに2台ずつ接続し、システムメモリを使用してシンクタイム12 [ $\mu S$ ]のときの本システムの稼働率を測定した。その結果を示したものが図4で×印が実測値を表している。I印は理論値である。I印の上端は各  $\mu P$ がシステムメモリをアクセスするときのアクセスタイムの最短の場合を示しており、下端は最長の場合を示している。この図より稼働率の実測値が理論値とよく一致していることがわかる。先の(1)式は  $\mu P$ の能力の利用率を表している。この式をもとにシステムのハードウェア特性の評価を行うことができる。具体的にはバスの多重度、高速メモリの評価等をシステムの基本設計時に行うことができる。

#### 5. おわりに

以上、資源共有型並列計算機“砂丘”の共有メモリの制御方式、機能及び  $\mu P$ 間の同期・通信方式を紹介した。数値シミュレーション等を行う場合には主メモリに対するアクセスが頻発することから、競合を抑制し主メモリを能率よく共同利用するために、マルチリード・ワンライトメモリ方式の採用が効果的であることを示した。また稼働率の実測値が理論値とよく一致することも示した。

今後の課題は本システムを使用して各種の具体的な問題に対する評価をすることである。

#### 参考文献

- [1] 白川 他: “並列計算機 PAX-128” 通信学論, Vol. J67-D, No. 8, pp. 853-860, Aug. (1984)
- [2] 出口 他: “コンピュータグラフィックシステム LINKS-1における画像生成の高速化手法” 情報処理論文誌, Vol. 25, No. 6, pp. 944-952, Nov. (1984)
- [3] Rodrigue. G.: “Parallel Computation” Academic Press (1982)
- [4] Paker. Y: “Multi-microprocessor” Academic Press (1983)
- [5] 井上・小林: “マルチマイクロプロセッサシステム  $\alpha-16$ のアーキテクチャ” 情報処理論文誌, Vol. 25, No. 4, pp. 632-639, July (1984)
- [6] 井上・小林: “ $\alpha-16$ マルチマイクロプロセッサシステムの性能評価” 情報処理論文誌, Vol. 25, No. 4, pp. 640-646, July (1984)
- [7] 井上・小林: “資源共有型マルチマイクロプロセッサシステムにおけるアクセス競合の調停について” 電子情報通信, 回路とシステム研究会資料, CAS84-206, pp. 9-16 (1984)
- [8] 山根 他: “並列計算機“砂丘”のハードウェアアーキテクチャ” 電子情報通信, コンピュータシステム研究会資料, CPSY87-3, pp. 15-22, June (1987)
- [9] 加納 他: “マルチマイクロプロセッサシステムの大容量共有メモリの一構成法” 情報処理学会 計算機アーキテクチャ研究会資料, CA-66-2, pp. 1-8, July (1987)
- [10] 井上 他: “マルチプロセッサシステム“砂丘”の共有メモリアーキテクチャについて” 情報処理学会 計算機アーキテクチャ研究会資料, CA-66-2, pp. 9-16, Nov. (1989)