

# ストレージバス結合並列計算機

A Parallel Computer with Storage Bus

有田 五次郎

Itsujiro Arita

兼田 考三

Genzo Kaneda

村上 健一郎

Ken-ichiro Murakami

九州大学 工学部 情報工学科

The Department of Computer Science and Communication Engineering, Kyushu University

## 1. はじめに

高多重並列計算機における大きな問題は、オペレーション間の同期比、データアクセスの競合である。最も一般的な並列処理機構であるMIMD型並列計算機は、この二つの問題のため高多重化は困難とされてきた。

ここで述べる並列計算機は、待ちなし並列プログラムの概念を用いた同期ハードウェアを持ち、ストレージバスの概念によってアクセス競合を減少させた、高多重化可能な並列計算機である。

## 2. ストレージバス

ストレージバスは並列度が増加した時にアクセスタイムは多少増加するが、アクセス競合は増加せず、しかもハードウェアの増加がCPUの数に比例するものである。

2.1 3次元8ノードストレージバス  
3次元8ノードのストレージバスの構成を図1に示す。

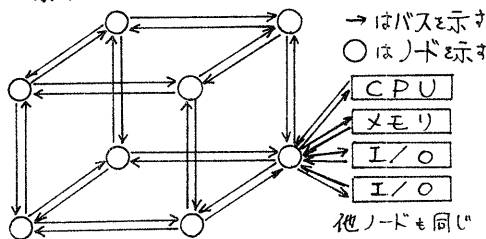


図1 3次元8ノードストレージバスの構成

ノードから出ている7個のバスすなわち4個のデバイス及び他の3個のノードへのバスは同一のバスである。従って、この並列計算機では同時に異なる、たゞ8個のデータ転送バスが存在する事に存す。

ノードの構成を図2に示す。ノードから外へ出るストレージバスは、他の3個のノードの  $BUF_x^i, BUF_y^j, BUF_z^k$  ( $i, j, k \neq x$ ) 計3ヶを経由して終わることとなる。

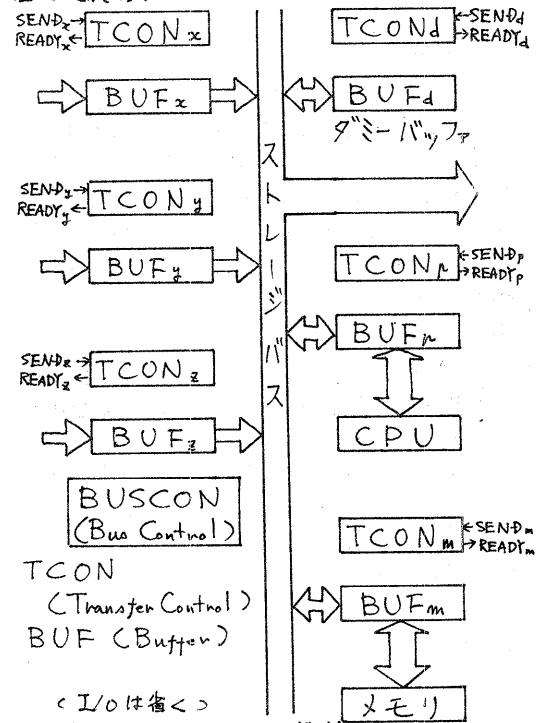


図2 ノードの構成

ストレージバス上のデータ及びコントロール信号を表1に示す。各ノードには図3に示す様に各ノードにシステムアドレス、各デバイスにデバイスアドレスが付与されている。

SSAはストレージバスを介してデータを転送する時の出発地ノードのシステムアドレスであり、DDAは目的地ノードのシステムアドレスである。

SDA, DDAはそれぞれ出発地、目的地のデバイスアドレスである。

PCNTはBUF<sub>n</sub>, BUF<sub>m</sub>からストレージバスにデータが出力される時に7リアセキ。各ノードのBUF<sub>x</sub>, y, zを通過する度に1加算される。この値が、最も離れたノード間の距離以上になるとTCONにおいて、そのストレージバス上のデータは抹消される。BUF<sub>d</sub>では、PCNTは何も加算されない。

SENDは前段のノードからストレージバス上にデータが出力されている事を示し、READYはBUFが空である事を前段のノードに知らせる、ストレージバスのコントロール信号である。

### データ

MA (Memory Address)	15 bit
M (Modifier)	16 bit
MDP (Memory Data Parity)	1 bit
R/W (Read or Write)	1 bit
TAS (Test and Set)	1 bit
SSA (Source System Address)	6 bit
SDA (Source Device Address)	2 bit
DSA (Destination System " )	6 bit
DDA (Destination Device " )	2 bit

### コントロール

SEND	1 bit
READY	1 bit

表1 ストレージバス上のデータ及びコントロール

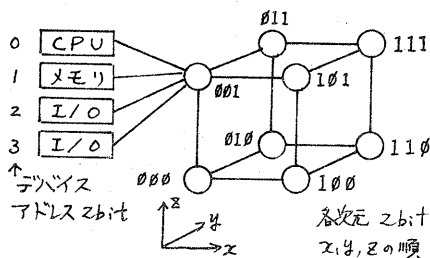


図3 システムアドレス-デバイスアドレス

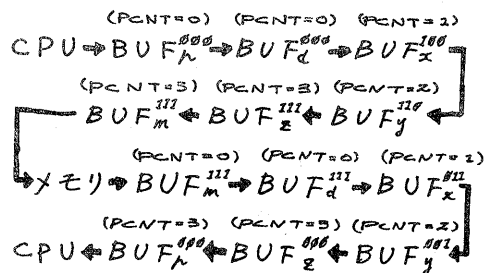
### ス.2 データ転送

このシステムにおけるデータ転送では、ノード間の距離が同じであれば、どの経路を通過しても構わない。(ルーチング) 図2, 図3を用いて実際のデータ転送の流れを示す。

1) CPUがプログラム(自ノードのメモリ)をフェッチする場合。



2) システムアドレス000のノードのCPUがシステムアドレス111のノードのメモリをアクセスする場合。



右上の数字はシステムアドレスを示す。

データ転送可能な次元の次段のノードがすべてREADYである時にはx, y, zの次元の順に優先度が低くなる。上記の2)は6通りある経路のうち1例である。

### ス.3 ダミーバッファ

1つの次元にノードが3個以上あるシステムの場合、非同期転送を行なうとデッドロックを起す。これを防ぐ為に、各次元に1個以上の余分なバッファを持つ必要がある。このシステムでは各ノードに1個のダミーバッファを設けている。

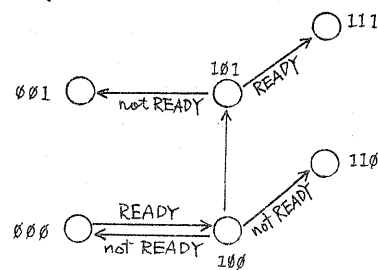


図4 ダミーバッファの効用

α, β, γ の方向の転送においてその行先がすべて not READY である時に、このダミーバッファにデータがストアされる。そして、not READY のどれかが READY になると、その次元・次段のノードへデータが転送される。

このダミーバッファには他に次にあげるものの効用がある。

1) CPU が他ノードのメモリをアクセスしようとする時はこのダミーバッファを介してデータを転送する。これにより CPU とストレージバスをたぐ TCONμ のハードウェア量が軽減する。

2) 図3から抜き出した図4の状態において次の様なデータ転送が行われた場合を考える。

- a) ノード100のCPUがノード001のメモリをアクセスの為にノード101へデータを送る。
- b) ノード000のCPUがノード111のメモリをアクセスする為にノード100へデータを送る。(これはαの優先度が最も高いのである。)

この時に、ノード100からノード101へ入ったデータが BUF<sub>101</sub>に入ると、ノード000からノード100へ入ったデータはノード100の BUF<sub>100</sub>を通ってノード111へ行く事ができ転送が速くなる。

すなわち、転送待ちのデータの一時退避場所として使用でき、他の転送可能データの妨害とならぬ。

また、これが有効に使えるのは、ノードの通過時間がCPUのサイクルタイム

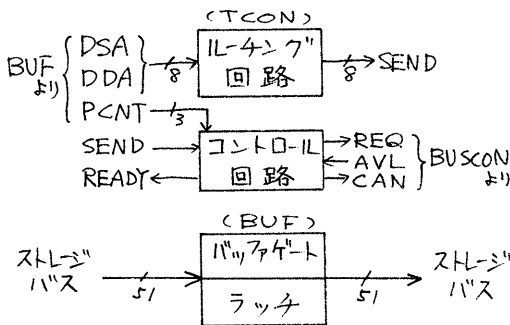


図5 TCON と BUF (α, β, γ, δ) ブロック図

メモリのアクセスタイムに比べて短いからである。もし、ノードの通過時間が同程度以上であれば、この効用は薄れてくる。

#### エ. 4 転送制御及ガルーチアップ回路

TCON はデータ転送の制御を行なう回路で前段のノードからのデータの受け取り、次に転送すべきノード又は CPU, メモリ, ダミーバッファの決定そしてデータの送り出しを行なう BUF はラッチとバッファゲートから成る。

図5に TCON と BUF (α, β, γ, δ) のブロック図, 図6に動作図を示す。

TCON は前段のノードからの SEND を受け付けた時に、果して転送可能な次元の次段のノードが READY であるかどうかを調べる前に、BUSCON に対して REQ を出してしまふ。そこで、もしすべての該当する READY がすべて not READY であり、どこにもデータが転送できない時、またたとえ READY であっても BUSCON の優先順位の為に他の同一ノード内の TCON にバス使用許可をとられてしまふ、すべて not READY になった時に一度受け取ったバス使用許可 AVL を返上しなければならぬ。その為に CAN という信号が出力される。これはいわゆる先行制御であり、それより優先順位の低い REQ が同時に出力されていた時のみ妨害となる。

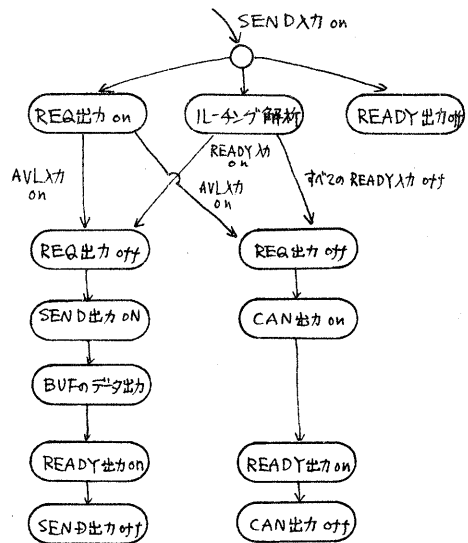


図6 TCON と BUF の動作図

$T_k$  ( $k=x, y, z, d, u, v, n, m$ ) を  $BUF_k$  へのデータ転送可能を示す信号とする。FSA をノード固有のシステムアドレスとする。デバイスアドレスを

- 0 = n (CPU)
- 1 = m (メモリ)
- 2 = u (I/O)
- 3 = v (I/O)

とする。

その時、ルーチング及び今まで述べたデータ転送を行なうためには、下記の論理式が成り立たなければならぬ。

$$T_x = DSA_0 \oplus FSA_0 \cdot DSA_1 \oplus FSA_1$$

$$T_y = DSA_2 \oplus FSA_2 \cdot DSA_3 \oplus FSA_3$$

$$T_z = DSA_4 \oplus FSA_4 \cdot DSA_5 \oplus FSA_5$$

$$T_d = T_x + T_y + T_z$$

$$T_n = \overline{T_x} \cdot \overline{T_y} \cdot \overline{T_z} \cdot DDA_0 \cdot DDA_1$$

$$T_m = \overline{T_x} \cdot \overline{T_y} \cdot \overline{T_z} \cdot DDA_0 \cdot DDA_1$$

$$T_u = \overline{T_x} \cdot \overline{T_y} \cdot \overline{T_z} \cdot DDA_0 \cdot DDA_1$$

$$T_v = \overline{T_x} \cdot \overline{T_y} \cdot \overline{T_z} \cdot DDA_0 \cdot DDA_1$$

$$SEND_x = \overline{T_x} \cdot READY_x$$

$$SEND_y = \overline{T_y} \cdot READY_y \cdot SEND_x$$

$$SEND_z = \overline{T_z} \cdot READY_z \cdot SEND_x \cdot SEND_y$$

$$SEND_d = \overline{T_d} \cdot READY_d \cdot SEND_x \cdot SEND_y \cdot SEND_z$$

$$SEND_n = \overline{T_n} \cdot READY_n$$

$$SEND_m = \overline{T_m} \cdot READY_m$$

$$SEND_u = \overline{T_u} \cdot READY_u$$

$$SEND_v = \overline{T_v} \cdot READY_v$$

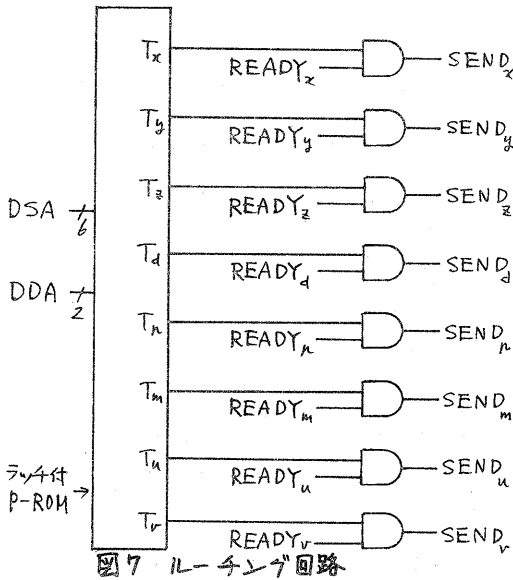


図7 ルーチング回路

これらの論理式を実現する為に図7の回路を用いる。ラッチ付P-ROMの一部を処理しているのがこの回路の特徴である。

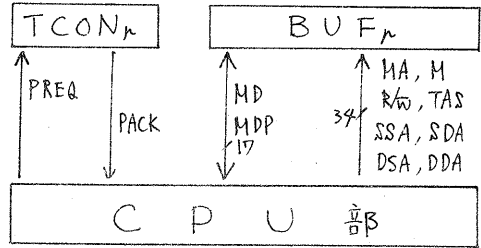


図8 TCON<sub>n</sub>及びBUF<sub>n</sub>とCPU部のインタフェイスブロック図

### 2.5 プロセッサインタフェイス回路

TCON<sub>n</sub>及びBUF<sub>n</sub>はCPU部とストレージバスとのインタフェイスを担持。TCON<sub>n</sub>及びBUF<sub>n</sub>とCPU部のインタフェイスを図8にブロック図で示す。

ここでM (Modifier) bitは0の時CPUが直ノードのメモリを、1の時CPUが他ノードのメモリをアクセスする事を示す。このM bitの使用は、CPU・メモリが他ノードへデータを転送する時にダミーバッファを使用する事、この事によって前述のルーチング回路がTCON<sub>n</sub>から省かれる。これはTCON<sub>m</sub>についても同様であり、ハードウェア量の軽減の為に行なわれている。

図9にTCON<sub>n</sub>, BUF<sub>n</sub>のブロック図を示す。

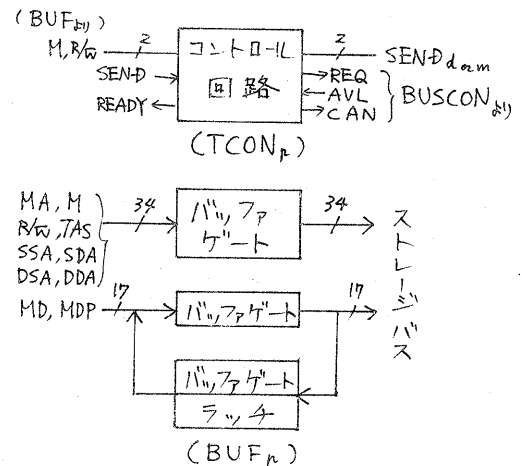


図9 TCON<sub>n</sub>, BUF<sub>n</sub>ブロック図

TCONの動作を図10に示す。TCON<sub>m</sub>はPREQにより起動される。そしてストレージバスを占有してデータを送り出した後、READYの状態ではSEND（メモリからの応答）待ちとなる。

ここで、メモリからの応答（リード、ライト共）が返ってきた時に、果たしてそのデータが本当に目的ノードからの応答なのかどうかを確認する必要がある。その為にはSSA, SDAとDSA, DDAと比較しなければならぬが、ここではTCON<sub>m</sub>のハードウェア量を軽減する為に省略する。メモリ側からSSA, SDAは送られてこない。

### 2.6 メモリインタフェース回路

TCON<sub>m</sub>及びBUF<sub>m</sub>はメモリ部とストレージバスのインタフェースを受け持つ。TCON<sub>m</sub>及びBUF<sub>m</sub>とメモリ部のインタフェースを図11にブロック図で示す。

メモリ部はメモリセルのデータを送り出すだけであるからDSA, DDA, SSA, SDAはメモリ部には入力されない。

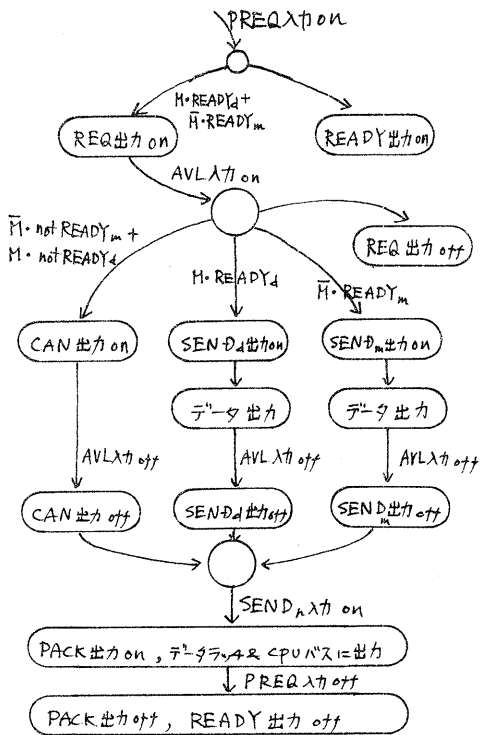


図10 TCON<sub>m</sub>の動作図

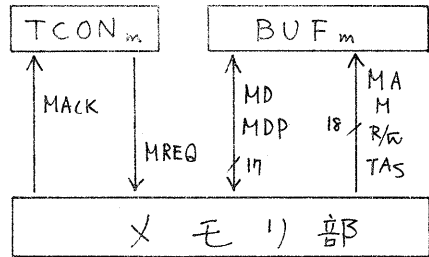
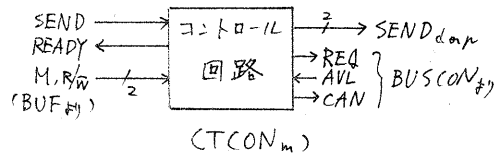


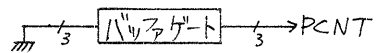
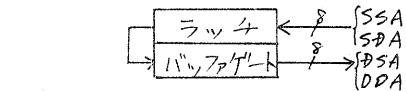
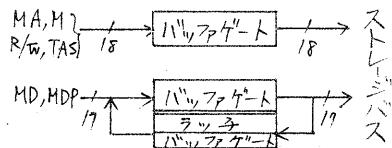
図11 TCON<sub>m</sub>及びBUF<sub>m</sub>とメモリ部のインタフェースブロック図

図12にTCON<sub>m</sub>とBUF<sub>m</sub>のブロック図を示す。TCON<sub>m</sub>と同じくルーチング回路が省略されている。

図13にTCON<sub>m</sub>の動作図を示す。TCON<sub>m</sub>はSEND<sub>m</sub>によって起動される。そして、MREQをメモリ部へ出力し、メモリのアクセスを開始させ、メモリ部からのMACK入力によってデータを得た後、ストレージバスを占有しCPUへデータを転送する。



(TCON<sub>m</sub>)



(BUF<sub>m</sub>)

図12 TCON<sub>m</sub>とBUF<sub>m</sub>のブロック図

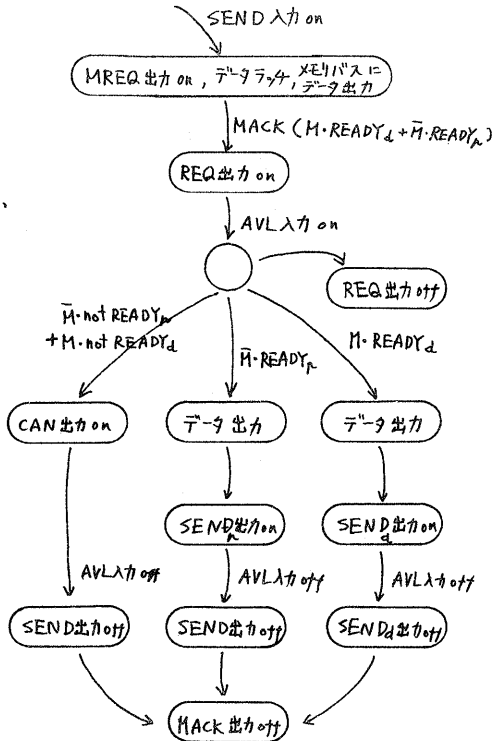


図13 TCON<sub>m</sub>の動作図

ス.7 バスコントロール回路

BUSCONはTCON<sub>k</sub> (k=α, β, γ, λ, μ, ν) から出されるストレージバス使用要求REQと優先度 (λ, μ, ν, α, β, γ, δ, ε, ζ, η, θ) の順に低くなる) を考慮して調停する。ブロック図を図14に示す。図15に動作図を示す。クロックを用いたサンプリング方式であるので、多少時間のロスが生じるがラッチを二段持っているので、引き続きREQがくる時には単位時間当たりの処理数は多くなる。いわゆるパイプライン処理と同様である。

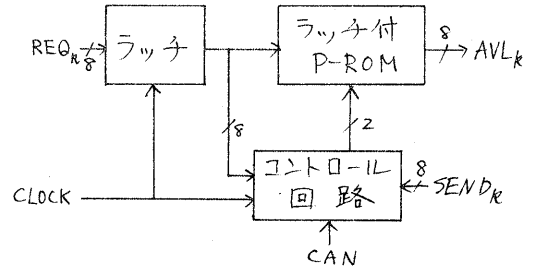


図14 BUSCONのブロック図

ス.8 転送時間

BUFからBUFへデータが転送されるのに要する時間は、構成素子にスタイプのTTLを使用した場合に素子のtyp.値で計算して50μm程度である。ただし、BUSCONがサンプリング方式なのでこれは平均値である。

従ってプログラムのアクセスの時のデータの流しはス.2の1)より、BUFからBUFへの転送が、BUF<sub>n</sub>→BUF<sub>m</sub>, BUF<sub>m</sub>→BUF<sub>n</sub>の2回であるから500nsメモリのアクセスタイムが延びる事になる。これはCPUのメモリサイクルの約25%増になる。この程度であれば、CPUとメモリの間に特別なバスを設けた場合のハードウェア量の増加と相殺できるとであろう。

他ノードのメモリをアクセスする場合、システムの次元数をm, 1次元のノード数を等しくαとすると、そのシステム内の最も離れたノード間の遅延時間Tは

$$T = 250(\alpha + m \cdot \alpha + \alpha) \quad (1)$$

で表わされる。ただしαは途中に通過したゲートバッファの個数である。

ここでシステムのノード数Nは

$$N = \alpha^m \quad (2)$$

で表わされる。

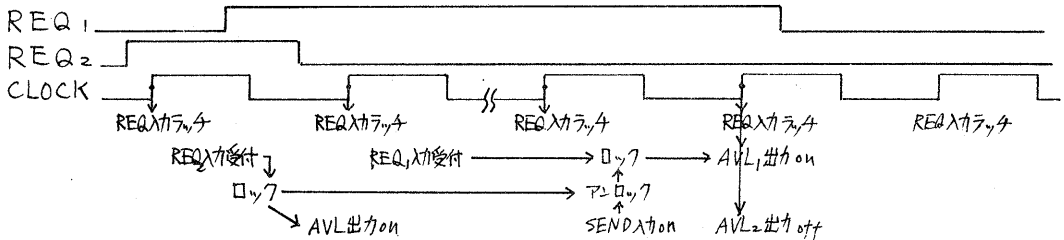


図15 BUSCONの動作図

(1)式と(2)式を比較してノード数が指数関数で少くなるのに、最大遅延時間Tは1次関数でしかかえないのが良くわかる。

実際、3次元(8)ノードの場合(1)式より

$$T = 250(2 + 3 \times 2) \quad (k=0) \\ = 2.0 \mu s$$

3次元(4)ノードの場合

$$T = 250(2 + 3 \times 4) \\ = 3.5 \mu s$$

4次元(8)ノードの場合

$$T = 250(2 + 4 \times 3) \\ = 3.5 \mu s$$

4次元(256)ノードの場合

$$T = 250(2 + 4 \times 4) \\ = 4.5 \mu s$$

にか過ぎない。これだけメモリのアクセス時間が伸びただけである。

### 3. CPU・メモリ

#### 3.1 待ちなし並列プログラム

並列プログラムはグラフ的に図16(a)のように表わされる。各ノードはオペレーション、各アークは制御の移動を表わし、同一列のノードは同一プロセッサで実行される。

このプログラムは適当な変換によって(b)のような木構造に変換できる。(b)はプロセッサが first come first serve で処理している時は(a)と同価であり、自然に同期がとれている。このようなプログラムを待ちなし並列プログラムと呼び制御が渡ると一度に実行されるオペレーションの集合をタスクと呼ぶ。

並列プログラムを実行するプロセッサに要求される命令は次の二つである。

- 1) 他プロセッサ上にタスクを発生させるための命令
- 2) 1つのタスクを終了後、次のタスクを実行するための命令

これらは比較的簡単なハードウェアで実現でき、従って同期のオーバーヘッドも少ない。この並列計算機では図16(b)の並列プログラムを実行させる事を目的としている。

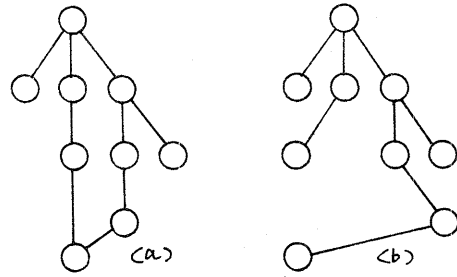


図16 並列プログラム

#### 3.2 CPU部・メモリ部の構成

今回の並列計算機製作においては、ハードウェア量の軽減に主眼を置いた。なぜなら、このハードウェア量の少くがシステムが実現するかどうかの鍵を握るからである。その為にマイクロプログラム方式の東芝製TOSBAC40-Lという5チップ構成のCPUを採用した。表3にその主な仕様を示す。又、TAS(Teapot Set)の弟子がチップに出ているメモリの実行制御に非常に都合が良い。

図17にCPU部、メモリ部のブロック図を示す。通常のCPU・メモリに比べて外部レジスタが付加されているに過ぎない。表3に外部レジスタの構成を示す。

EPIRはそのノードのメモリに登録されている実行可能なタスクのIDを示す。タスクは最高8個(8bit)まで可能である。

CCRはソフトウェアでコンソール制御を行なう為のもので、ステップ動作、コンソール割り込みが可能である。

FIFOSTATUSはタスク管理ハードウェアとして使われるFIFO(タスクキュー)のEMPTY, FULLを検出する為のレジスタで、リードゼロフ、ライトゼロフとなっている。

RSAは現在実行中のタスクを起動したノードのシステムアドレスを示す。

PIDは現在実行中のタスクのIDを示す。

NPSWは新しく実行するタスクのPSWである。

RMAは現在実行中のタスクの戻り先のメモリアドレスを示す。

以上の4つはFIFOから4回プルする事によって得られ、RSA, PID, RMAはCPU側のRSAR, PIDR, RMARに移される。

### 3.3 タスク管理ハードウェア

この並列計算機ではタスク管理の為に各CPUがタスクキュー(FIFO)を持っている。

1つのタスクは、RSA, PID, RMA, NPSWで表わされる。タスクキューには10個のタスクが登録できる。

このタスクキューを操作する為の拡張減減語命令(マイクロ命令による)は、PB(Parallel Branch), EXT(Exchange Task), CKTC(Check Task)の3つがあり、更にストレージバス使用の為にSDSA(Set Destination System Address)命令がある。これらにより並列プログラムの実行が可能となる。

PB命令は、他プロセッサ上にタスクを発生させる命令で、そのプロセッサのFIFOにタスクを書き込む。

EXT及びCKTC命令はFIFOからタスクを取り出し、そのタスクがEPIRに登録されているかどうかを調べ、登録されている場合は実行する。登録されていない場合は、それが共有タスク(I/OIL-チャンネル等)かどうかをNPSWをみて調べ、そうであれば実行する。そうでなければそのタスクはアボートされたタスクと見なされ実行されない。EXTは、FIFOがEMPTYの場合にタスクが新しく書き込まれるまで待つがCKTCは待たずに次の命令に進む。

図18に各命令のマイクロ命令による処理のフローチャートを示す。

演算論理	2進並列, 2の補数表示, 固定/浮動小数点表示
命令数	126種
アドレス空間	64 K Byte
アドレス指定	全アドレス直接指定, インテックス修飾
レジスタ	ゼネラルレジスタ 16bit x 16個 (内15個は、インテックスレジスタ兼用) 浮動小数点レジスタ 32bit x 8個 (主記憶上 X'0' ~ 'F')
制御方式	マイクロプログラム制御
語形式	命令語 16bit (RR, RF型), 32bit (RW, RX型) 固定小数点数 1, 8, 16, 32 bit/data 浮動小数点数 32 bit/data 制御語 (psw) 32 bit
割り込み	内部割り込み 6レベル 外部割り込み 基本外部割り込み 又は、オートマチックI/O

表2 TOSBAC 40-L仕様

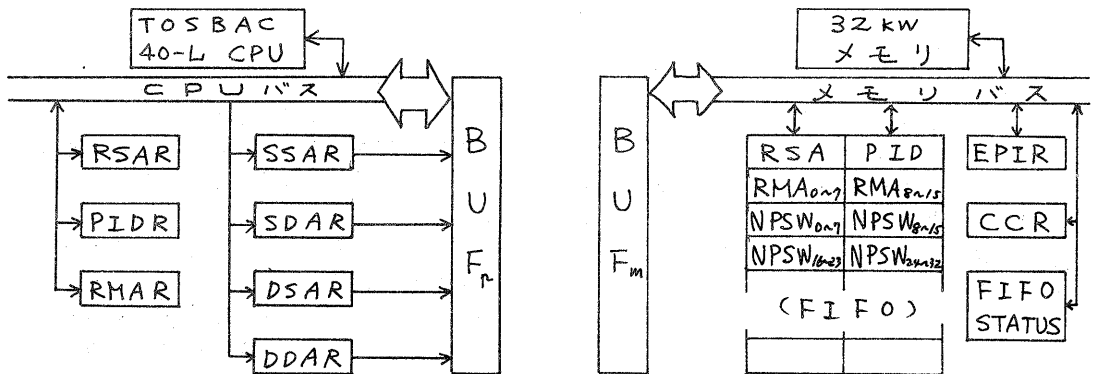


図17 CPU部・メモリ部



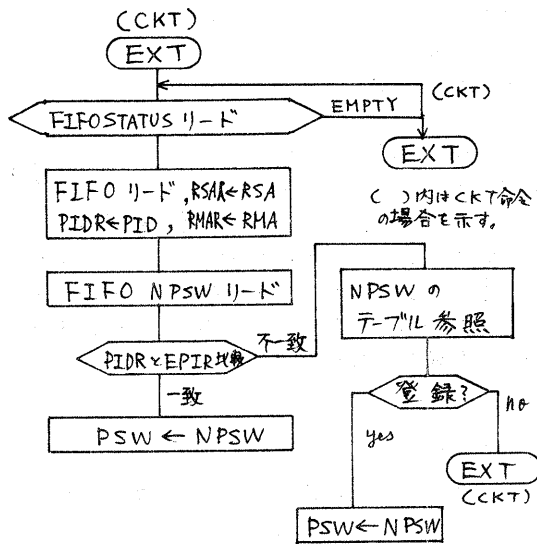


図 18 (a) EXT.CKT 命令フローチャート

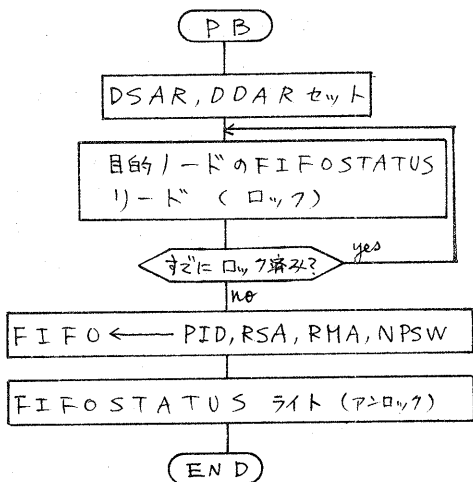


図 18 (b) PB 命令フローチャート

アドレス	bit	機能
Z0	8bit	EPIR (Execution Process ID Register)
Z1	3bit	CCR (Console Control Register)
BE	8bit	SSAR (6bit), SDAR (2bit)
BF	8bit	DSAR (bit), DDAR (2bit)
C0	5bit	FIFO STATUS
C2	8bit	(RSA or RMA <sub>0-7</sub> or NPSW <sub>0-7</sub> or NPSW <sub>16-23</sub> )
C3	8bit	(PID or RMA <sub>8-15</sub> or NPSW <sub>8-15</sub> or NPSW <sub>24-32</sub> )
C4	8bit	PIDR (Process ID Register)
C5	8bit	RSAR (Return System Address Register)
C6, C7	16bit	RMAR (Return Memory Address Register)

) FIFO

表 3 外部レジスタの構成

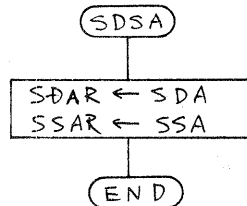


図 18 (c) SDSA 命令フローチャート

#### 4. 実装

表 4 に I/O ノード分について使用 IC の個数を示す。使われている IC は NAND ゲート, フリップフロップ, ラッチ付バッファゲート, バッファゲートがほとんどであり, 特殊な IC は ラッチ付 P-ROM ぐらいである。

図 19 にボードの形状, 図 20 にコネクタの接続図を示す。ボード 4 枚で I/O ノード (I/O 除く) を構成し, 50cm x 45cm x 10cm のラックに収納される。そしてこのラックに入る信号線は 50pin, 20pin フラットケーブルがそれぞれ 3 本出ていく信号線が 50pin, 20pin フラットケーブルでそれぞれ 4 本である。

TCON <sub>x, y, z</sub>	15ヶ
TCON <sub>d</sub>	14ヶ
TCON <sub>m</sub>	8ヶ
TCON <sub>p</sub>	7ヶ
BUF <sub>x, y, z, d</sub>	7ヶ
BUF <sub>m</sub>	10ヶ
BUF <sub>p</sub>	11ヶ
BUSCON	6ヶ
CPU 部	45ヶ
メモリ部	53ヶ

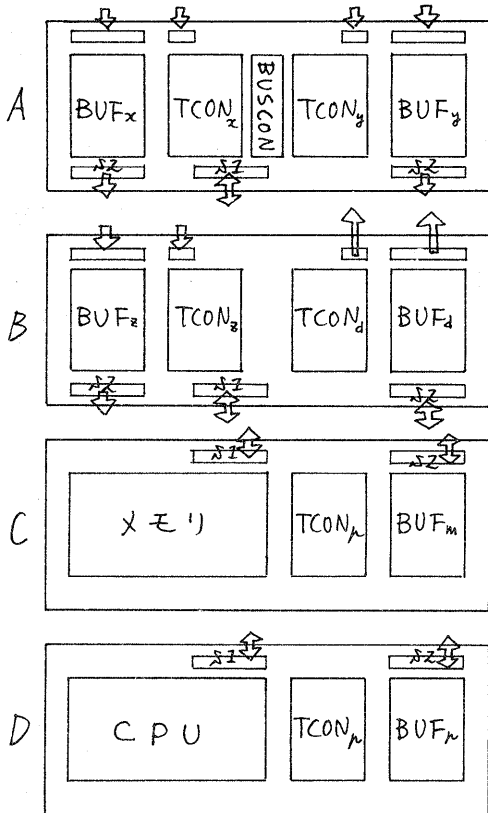
計 176ヶ

表 4 使用 IC の個数

### 5. おわりに

以上述べた並列計算機は現在プリントパターンを描き終り、エッチング、穴開けを外部に発注しているところである。

この並列計算機上のソフトウェアはまだほとんど未開発であり、ハードウェアが出来次第取り掛からなければならぬ。



□, □ ... 50 pin, 20 pin のフラットコネクタ

同じ記号のコネクタ同士接続する。

ボードの大きさは A, B, C, D 共 180mm × 410mm

図19 ボード図

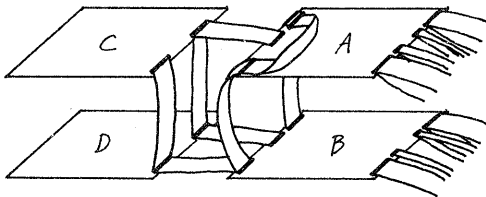


図20 コネクタ接続図

IPLの為に、I/Oとしてホストコンピュータをいづれかのノードに接続し、各ノードのメモリにプログラムをロードするつもりである。

この並列計算機では、シーケンスコントロール、アドレスデコーダ、ヒラキ付P-R O M を多く用いている。これがハードウェア量の軽減に大きく貢献している。特にCPU、メモリ部において著しい。しかし、P-R O M のアクセスタイムが50nsとSSI L S T T L に比べて遅いので、高速を目指した所へは応用できない。

### (参考文献)

1. 有田「並列処理システムの一構成法(I) ストレージバス」  
S54 電気四学会九州支部大会論文集
2. 有田, 兼田「並列処理システムの一構成法(II) 3次元ストレージバスの実現」  
S54 電気四学会九州支部大会論文集
3. 有田「MIMD型高多重複合計算機-MMC」(I) システムアーキテクチャ」  
S55 電子通信学会全国大会
4. 有田, 兼田, 村上「MIMD型高多重複合計算機MMC」(II) ストレージバス」  
S55 電子通信学会全国大会
5. 有田, 村上, 兼田「MIMD型高多重複合計算機MMC」(III) CPU, メモリ部」  
S55 電子通信学会全国大会
6. 村上「ストレージバス結合による複合計算機システム的设计」  
九州大学 工学部 昭和54年卒論
7. 「TOSBAC-40L アプリケーションガイド」LSI」東芝