

多重チャネルリングバスに於ける

ブロードキャスト伝送制御手順

——機能分散型計算機の内部接続方式——

喜連川 優、田中 英彦、元岡 達 (東大 工学部)

1. はじめに

近年、データベースの普及に伴い、ソフトウェアサポートによるDBMSの性能の限界が次第に明らかになり、データベース管理の諸機能をハードウェアでサポートする事を試みたデータベースマシンの研究が着手されている。データベースに於ける検索処理の中には、潜在的に高度の並列性を見出す事が出来、多くのメモリ、プロセッサモジュールを互いに結合したアーキテクチャが採られる場合が多い。更に、検索処理の基本は内容による読み出し(Content Addressing)であるが、連想メモリを用いて大規模なデータベースを格納するにはコストの面から実際的でなく、従って二次記憶のフルスキャンによって擬似的に連想機能を実現する事がなされてきた。我々も同様の考えから、より発展させた形のマシン、“可変構造多重処理データベースマシン”の設計、製作をすすめている^{1,2)}。図1に示される本マシンはメモリモジュールから複数のプロセッシングモジュールに対してデータをブロードキャストし、並列に処理する事により高速化を狙っている。データのフルスキャンによる連想性は、アーキテクチャ的にはデータブロードキャストにより実現され、これはメモリコンフリクトのない環境を生み出している為、駆動プロセッサ台数に略比した処理速度の向上が期待できる。しかしながら、多くのプロセッサに対しデータをブロードキャストする場合、一般に一台でもそのデータを読み誤ると当該オペレーションの正当性が保証されない場合がある。従って、この様なアーキテクチャを採用した場合、何らかの伝送制御手順によって、データの送受を正しく実現する事が不可欠となる。

ここではリングバス上でのブロードキャストを実現する際バス上のビットエラー及びバスインターフェイス部での同期障害に対して回復可能な伝送制御手順についての一手法(Broadcast Hand Shaking)及び逆に多くのモジュールから別の一つへ効率よくデータをあつめる為の手順(3 Phase Hand Shaking)について報告する。

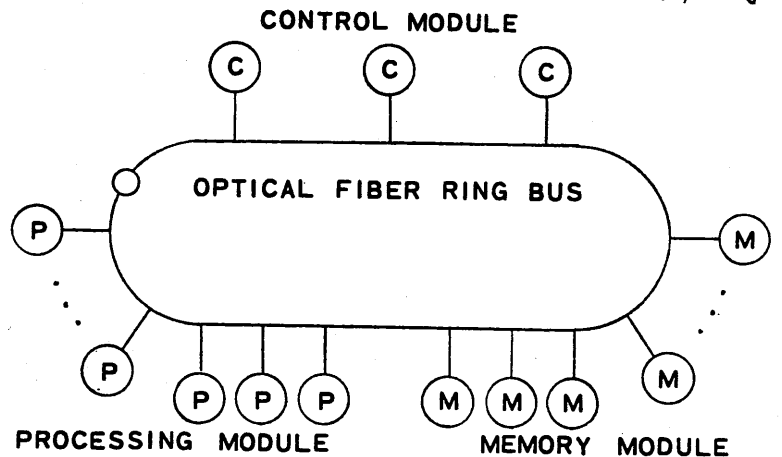


FIG. 1 ARCHITECTURE OF DATA MANIPULATION UNIT

2. リングバスデータ形式

リングバスは図2に示される如く固定長チャネル時分割多重方式を採るものとする。ここでプロセッサはメモリから送られたデータに対し処理を施すが、これはバスの転送レートではなくデータの出力レートによって定まる為、プロセッサはその処理結果を当該データチャネルよりいくつか後のチャネルに出力する。

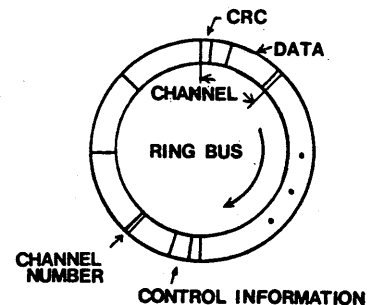


FIG. 2 RING BUS DATA FORMAT

3. Broadcast Hand Shaking

この伝送手順を図3に示す。一台のメモリから多数のプロセッサに対してデータをブロードキャストする際の伝送手順であるということによってこの様に呼んだ。

3-1 設計方針

1) 伝送制御手順の為の制御情報は全てコントロール部に入れ伝送制御の簡素化を図る。

チャンネルはデータ、コントロール两部分に分かれプロセッサは前者に対して処理を施し、後者に処理結果を書き込むが、ここに伝送制御用の情報を含め、データ部分には一切含まない。両者をまとめ、(データ、コントロール)を一つのユニットとして回復処理を行なう事によって、伝送手順の簡素化を図る事が出来る。

2) プロセッサを群として扱う。

個々のプロセッサを意識してデータを送出するのでは制御できる台数には限度がある。従って任意台数のプロセッサに対してデータ伝送出来る為には、対象を群としてとらえる必要があり、ここでは結合プロセッサのIDではなく総数に着目した。

3) 一台のプロセッサの脱落も許さない。

ブロードキャストは常に全プロセッサが正しいデータを受信した事を確認しつつ行なう。従って一台でも読み損じた場合には全プロセッサに対して同一データを再送し、当該プロセッサが回復するまで他を行たせる。

3-2 動作概説

▶ メモリの動作

- 1) データパートが来るとデータを出力しコントロールパートが来るのを待つ。
- 2) コントロールパートが来ると、1)で送ったデータも全てのプロセッサが正しく受信したかどうかチェックし、正常なら、ACKを送り、次のデータを用意し、1)へ。異常なら、NACKを送り、リカバリモードに入り、3)へ。
- 3) NACKの場合、次のデータパートで前々回のデータを再送し、プロセッサの回復を待つ。
- 4) プロセッサが回復すると、1)へ。

▶ プロセッサの動作

- 1) データパートが来るとデータを受信し、正しく受信された場合は処理を施し、コントロールパートが来るのを待つ。
- 2) コントロールパートが来て、ACKの場合にはデータを受信した事を書き込むと共に処理結果を出力し、1)へ。
- 3) NACKの場合には、前回及び前々回のデータに対するオペレーションとバックアウトし、1)へ。

以上は正常なプロセッサ及びメモリの振舞であり、障害回復の仕方は省いている。

3-3 総数チェックによる障害検出

メモリは自分に結合している多くのプロセッサとその総数が等しいかどうかという事によって動作監視する。即ち、メモリはコントロールパートの一つのフィールドをカウンタとして用い、結合されたプロセッサの総数に等しい値を書き込んで送出する。プロセッサは1だけデクリメントして値を返す事になると、カウンタはプロセッサ群を通り、再び戻って来た時には0になっている筈である。0によって、プロセッサ群正常とみなす事にする。(総数チェック) このカウンタ長は $\log_2 N$ (N : 結合プロセッサ総数)

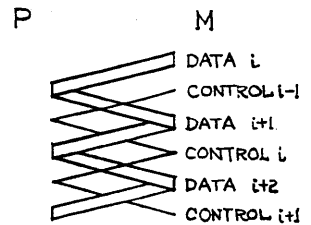


Fig.3. Broadcast Hand Shaking

3-4 障害回復

総数チェックにより障害が検出された時の回復手順について考える。プロセッサは当該データに対する処理に一定の時間が 必要な為、その処理結果とデータの入っているチャネルよりも幾つか後のチャネルに書込む。その為、メモリにとってはデータパートが戻ってきた時には、次のデータを送らすべきか、或は再送すべきかはまた不明であり、プロセッサ群の動作が正常である事を期待して次のデータと出力する。その後コントロールパートが来、制御情報が得られ、その段階で始めて前回に送出したデータが正常に受信されたかどうかを判明する。この遅延が制御を複雑にしている。勿論、一つのデータ毎に確認すれば簡単な制御で済むが、データ伝送率は半分になってしまう。障害を検出するとメモリはNACKを送り再送モードであることを示し、次のデータパートで前回のデータを送出する。この再送データに対する結果も同様に遅れて到着するので、メモリは回復対象データとそれに続くデータの二つを組にして回復するまで送出を繰返す。この際、再送データの前はNACKと次はACKと交互に送る。障害プロセッサ自身の回復動作については、障害を次の二つに分けて考える事にする。

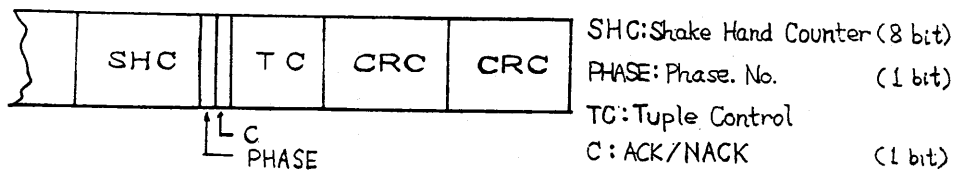
a. 一台のプロセッサが障害を起こした場合

b. 二台以上のプロセッサが引き続いて障害を起こした場合

aの場合、障害を起こしたプロセッサは回復した時点で、メモリ及び他のプロセッサ群が自分の回復を行っている筈であり、NACKの次のデータへリンクすればよい。aでは回復対象が一つに確定している為、制御は容易であるが、実際にはbの場合がある。複数台のプロセッサが全く同時に障害を起こした場合aに帰着されるが連続してプロセッサが障害を起こした場合が問題となる。プロセッサ P_0 が障害を起こし当該チャネルを読みとばし、プロセッサ P_1 はそれは読んだが次を読みとばしたとする。その時、全システムの回復対象は P_0 の読みとばしたデータであり、 P_1 のそれではない。 P_1 はシステムが自分の回復を待っていると考えてはならない。即ち障害を起こしたプロセッサにとって自分より一つ前に他のプロセッサが障害を起こしている可能性がある為、回復対象がいずれであるか判定出来ぬはならない。そこで、ここではPhase Bit (1bit)を導入し、データに対し、0,1と交互に付ける事とした。この1ビットにより、回復対象を決定出来る。例えば $\{(1\text{NACK}), D_0, (0\text{ACK}), D_1\}$ なる系列が送られてくると、プロセッサの最後に受信したフェイズが1の場合には、フェイズ1のNACKが来る事により自分が回復対象である事がわかり、その次のフェイズ0のデータより処理を再開する。0の場合には自分より一つ前に障害を起こしたプロセッサがあり、それに対するNACKを読みとばした事になる。従ってフェイズ0の処理をバックアウトし、処理を再開する。他のプロセッサの障害回復中に障害を起こしたプロセッサの回復も同様にして行なえる。回復対象に対する制御情報を読みとばすと、回復後もやはりこのデータの回復を行っている筈であり、実質上回復対象の決定は不要となる。

以上、障害を起こしたプロセッサが回復出来る為にはそれだけの情報と与える必要があり、リングバスに於けるブロードキャストではカウンターを用いACK/NACKの他にもう1bit Phase Bitを付加する事でこれが可能である事が示された。

3-5 コントロールパートフォーマット



3-6 状態遷移図

説明は略するが、上述の伝送手順を状態遷移図に表わすと、メモリは図6、プロセッサは図7、図8の如くなる。状態変数の記述は表1参照。(Phase 0のみ)

4. 3-Phase Hand Shaking

多くのプロセッサ群からメモリに対して効率よくデータと出力する為の回復可能な手続として、以下の方式を採用した。図4に示される如く3つのフェイズを互いに重なった形になっている為、3 Phase Hand Shaking と呼んだ。

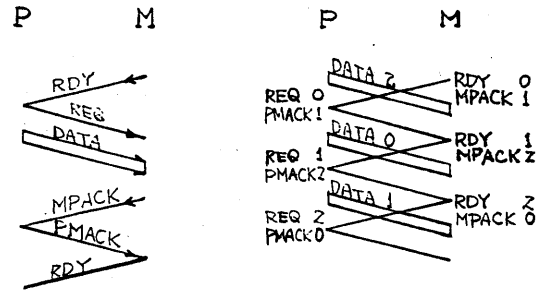


Fig. 4. 3-Phase Hand Shaking

4-1 設計方針

この伝送制御手順を設計する上で、次の点と考慮した。

1) データパートとコントロールパートの分離

Broadcast Hand Shaking との一貫性を保つため

2) データ伝送効率の向上

データ伝送は3つの Phase がインターリーブした形になっており、無駄はない。更に一つの Phase がプロセッサリカバリに使われていても他の二つの Phase は有効である様、効率向上を計っている。

4-2 動作概説

▶ プロセッサの動作

- 1) コントロールパートで RDY になっているものを探す。
- 2) REQ を出力するとともに、自分の ID を書き込む。
- 3) そのコントロールパートに対するデータパートが来ると、データと出力する。
- 4) メモリがデータを正しく受けとったかどうか、MPACK/MPNACK を待つ。
- 5) MPACK の場合はメモリから返事と受け取ったという事をメモリへ PMACK として送る

6) メモリが次の入力に移ったのを見て、1)へ戻り新たな出力に取りかかる。

5) MPNACK の場合は、1)へ戻り再出力する。

▶ メモリの動作

- 1) RDY を送る。
- 2) RDY が REQ として戻ってくると、データチャンネルを読み、特定パターンと出かしておく。
- 3) データが正常に読み取れると MPACK を、エラーが検出されると MPNACK を、プロセッサへ送る。
- 4) ACK を送った場合、プロセッサからの ACK を待ち、戻ってくると 1)へ。NACK の場合は、そのまま 1)へ。

以上はプロセッサ、メモリが障害を起こさず正常にデータをやりとりする場合である。

4-3 障害回復

ここでは、障害(同期ずれ)回復と考へに入れた場合、先の手順が最小限必要であることを示す。障害の発生する時点は種々考へられるが、論理的に図5の如く、3つ(T_0 , T_1 , T_2)に分ける事が出来る。プロセッサの動作に着目して考えると、

T_0 : RDYをみつけREQを出したが、データを出す前に障害と起こした。

T_1 : データは出力したが、メモリからのMPACKを受け取る前に障害と起こした。

T_2 : メモリからのMPACKを受け取り、PMACKを送ったが、それがメモリに届いたかどうか確認する前に障害と起こした。

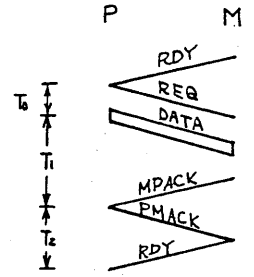


Fig. 5.3-Phase Hand Shaking
に於ける 障害分類

1) T_0 の場合の回復

データパートには制御情報を入れない事にした為、プロセッサはまずコントロールパートでRDYを探し、対応するデータチャネルを確保する。プロセッサがRDY→REQとしたにも拘わらず、データは出力せずに障害と起こした時の検出は、メモリがデータパートに常時特定パターンを出力してあり、それ以外のデータであった時はじめて読みとる事によって実現出来る。メモリはNACKを送り、プロセッサは再出力すればよい。

2) T_1 の場合の回復

プロセッサの動作の5)(PMACK)は、一見冗長に思われるが T_1 からの障害回復を行なおうとすると不可欠となる。例えば、PMACKを略し、MPACKをやり捨てにした場合と考えると、プロセッサがデータを送出した後、メモリ側で正常に受け取られたかどうか知るべきでない。従って回復可能である為には、もしデータをメモリが正常に受信した場合、プロセッサがMPACKを受け取ったという証拠にメモリはPMACKを受けとらねばならないという事になる。従って、 T_1 でプロセッサが障害と起こし、メモリがMPACKを送ったにも拘わらずPMACKが戻ってこない場合、メモリはプロセッサが回復するのを待つ。この「待つ」という動作によって不確定時間障害を起こすプロセッサのリカバリが可能になっている。MPNACKの場合も同様にプロセッサの回復を待つ事が出来るがここではNACKの時はやり捨てにした。メモリはMPACKを出力したにも拘わらずプロセッサよりPMACKが戻ってこないとダミーサイクルを繰返し、当該Phaseに於て回復中のプロセッサからのPMACKを待つ。PMACKが戻ってくると、次データの入力に取りかかる。プロセッサは回復すると自分が前に出力したチャネルに対して自分の回復を待つてくれているかどうか調べる。もし待つてくれている場合にはデータが正しく転送されているのでPMACKを送り次の仕事に取りかかる。待つてくれない場合には正しく転送されなかったことになり、RDYを探して再出力する必要がある。

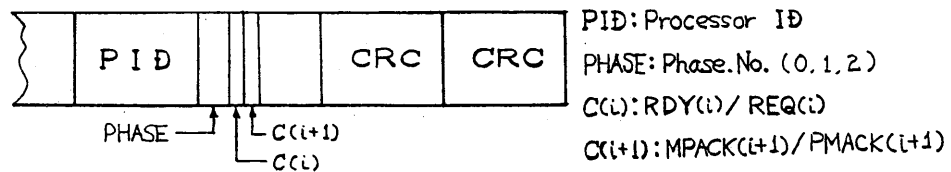
次にプロセッサの識別について考へる。メモリはある一台のプロセッサとデータのやり取りを行ない、その回復を行なう必要がある。この為には個々のプロセッサを意識せざるを得ない。つまり引き続き二台のプロセッサが同一Phaseにリンクし両方共障害と起こしたが、後者のみデータが正常に受け取られた場合、二つの

プロセッサにとってメモリが一体誰の回復を待っているのか識別出来る必要がある。データを正常に受信しなかった場合も回復を行なうならば一つのチャネルの一つのPhaseに結合し得るプロセッサは一台だけであるから、識別する必要はない。しかし、一方でデータの出カ効率を上げる事が望まれ、NACKをやり捨てにする事はこれに等与するが、逆に一つのPhaseに複数の回復中のプロセッサがリニアする可能性が生ずる為、それ等を区別する必要がある。ここではBroadcast Hand Shakingの際の総数チップ用のフィールドをプロセッサIDの為に使い、NACKをやり捨てたして、データを正常に受け取った時のみ回復の対象とする事にした。回復中のメモリはRDYの代わりに、始めからREQにしてコントロールパートを出力する。その際、本来ならばプロセッサ自身がPIDを書き込むのであるが、あらかじめ、回復対象であるプロセッサのIDを書いておく。その後のデータチャネルは読みとばし、MPACKを送り、PMACKが戻って来るのを待つ。プロセッサは回復後、当該Phaseにリンクし、自分のIDが書かれているかどうかと調べる。

3) T₂の場合の回復

プロセッサはメモリよりMPACKを受け取り、PMACKを返送したが、その後障害と起こした為、PMACKがメモリに正しく届いたかどうか不明になる。この場合のメモリの回復方法はT₁の場合に等しい。つまり、PMACKがエラーを起こすと、メモリにとってはプロセッサからデータは正常に受け取ったが、MPACKに対する応答であるPMACKが戻って来ないという状態に変わりはない。プロセッサの回復手法はT₁のそれと異なる。即ち、プロセッサは既にメモリがデータを受け取った事実は知っており、PMACKがメモリに届いたかどうかという事のみが問題となる。T₁の場合にはデータ自体が届いたかどうか問題であった。プロセッサはT₂と同様に当該Phaseにリンクし、そこに自分のIDが書かれてあると、PMACKが正しく届かなかった事を示し、それを再送する。書かれていないと、正しく届いた事を示し、次のデータ出力に取りかかる。

4-4 コントロールチャネルフォーマット



4-5 状態遷移図

説明は略するが、上述の伝送手順を状態遷移図に表わすと、メモリは図9、プロセッサは図10の如くなる。

- 但し・Phaseの0~2は各々対等であるから、ここではPhase 0の部分だけ表わすことにする。
- ・実線は、正常動作中での遷移を表わし、破線は障害が生じた場合の回復動作に伴う遷移を表わす。
 - ・状態、変数の記述は表2参照。(Phase 0のみ)

5. おわりに

多数のプロセッサにデータをブロードキャストし並列処理を試みる場合、データの正しい送受と保証する為、何らかの伝送手続きが必要となる。ここでは、障害の検出だけでなく、回復可能な制御手順について検討した。これによって、障害時点から再び処理を開始することが可能となり、検出と再試行による手法と比べて性能の向上が期待できる。又制御用データが大きなオーバーヘッドとならない様、最小限の情報量によって実現することを念頭に置き、全体として僅かなビット数で押えることができた。

参考文献

- 1) 喜連川、赤松、田中、元岡 “可変構造多重処理データベースマシンの構成”
S.55 情報処理学会全国大会 2F-5
- 2) “ ” “ ” 2F-6
- 3) 「高速リングバスシステムハードウェア仕様書」
パターン情報処理システム技術研究組合 S.52

Memory

- D°: Phase 0 のデータ待ち (正常)
- C°: Phase 1 のコントロール待ち (”)
- D₁°: 回復対象データ (Phase 1) の再送待ち
- C₁°: 回復状態 NACK 確認

Processor

- D°: Phase 0 のデータ待ち (正常)
- C°: Phase 0 のコントロール待ち (”)
- C₀°: データ入力異常終了後コントロール待ち
- D₁°: Phase 1 のデータ待ち (回復中)
- C₁°: NACK 待ち (”)
- D₀°: Phase 0 のデータ再送待ち (”)
- C₁°: 再送データ正常入力後コントロール待ち (”)
- C₁°: 再送データ異常入力後コントロール待ち (”)
- R°: D°, C°, C₁° での障害回復後の状態
- R₁°: その他での障害回復後の状態

Memory

- C0: 総数チェック正常
- PH: Phase 0
- ACK: ACK 帰着
- CRC: データエラー
- D: データパケット到着
- C: コントロールパケット到着

Processor

- ACK: Memory より ACK
- DN: チャネル No. - 致
- PH: Phase 0
- CRC: データエラー
- D: データパケット到着
- C: コントロールパケット到着

表1 Broadcast Hand Shaking に於ける状態の記述

- Memory
- C₁: REQ 待ち
 - D₁: 入力データ待ち
 - D₀: REQなしデータ待ち
 - C₄: データ正常入力後MPACK送待待ち
 - C₃: データ異常入力後MPNACK送待待ち
 - C₆: MPACK後PMACK待ち
 - C₅: MPNACK確認待ち
 - C₂: PMNACK受信後コントロール待ち
 - D₂: タミデータ待ち
 - D₃: REQ異常入力後データ待ち
 - C₀: D₃での障害後の状態
- Processor
- C₀: RDY待ち
 - D₀: データ出力待ち
 - C₁: Phase 1のコントロール待ち
 - C₂: MPACK待ち
 - C₃: PMACK到着確認
 - R₁: C₁, C₂での障害後の状態
 - R₂: C₃での障害後の状態
 - I: アイドル状態

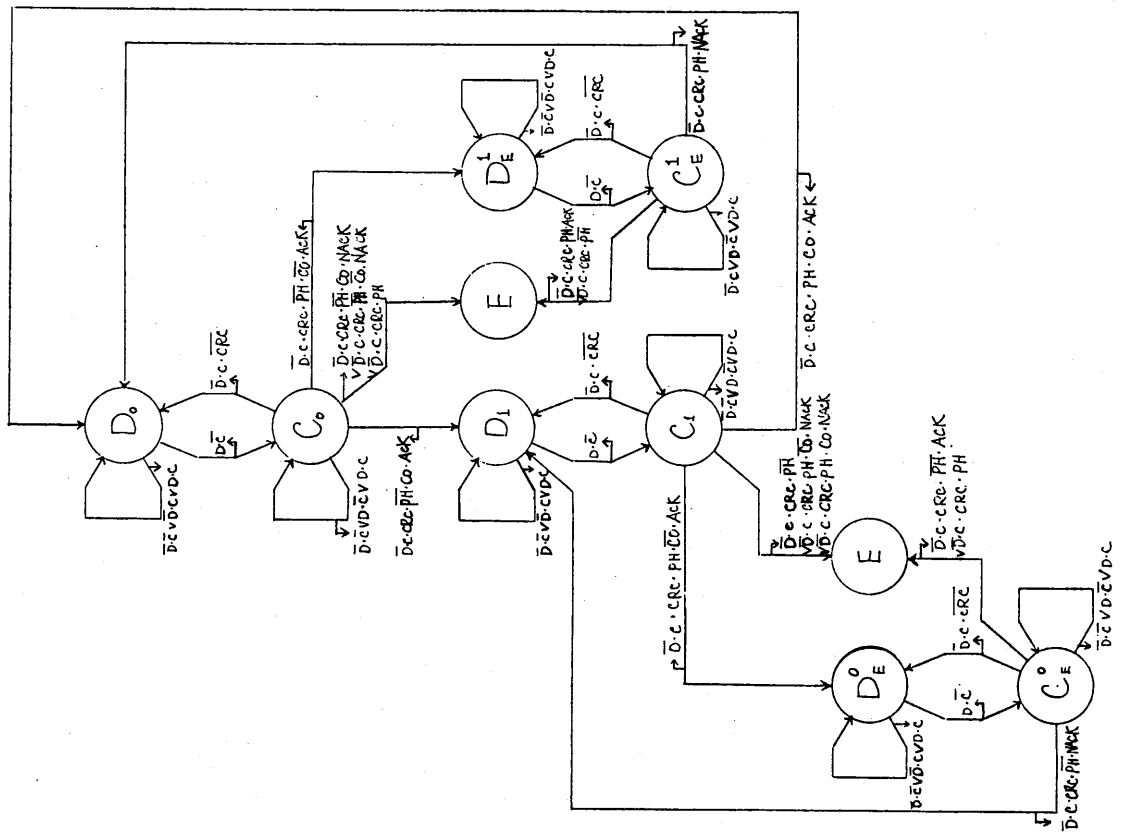
- Memory
- $\overline{RDY}(=REQ)$: Processorからの出力要求
 - PMACK: ProcessorよりACK
 - PH(0~2): Memoryの出力したPhaseNo.
 - SP: データが特定パターン
 - WAIT: データ入力禁止
 - CRE, D, C, DN: 表1と同じ

- Processor
- RDY: 当該チャネルが空
 - MPACK: MemoryよりACK
 - SID: IDフィールドに自分のIDが書かれている。
 - DMD: データ出力要求
 - PH(0~2), CRE, D, C, DN: Memoryと同じ

表2 3-Phase Hand Shaking に於ける状態の記述

図6 Broadcast Hand Shaking
Xレジスタル状態遷移図

図6



Broadcast Hand Shaking
 プロセス間モジュール状態遷移図II

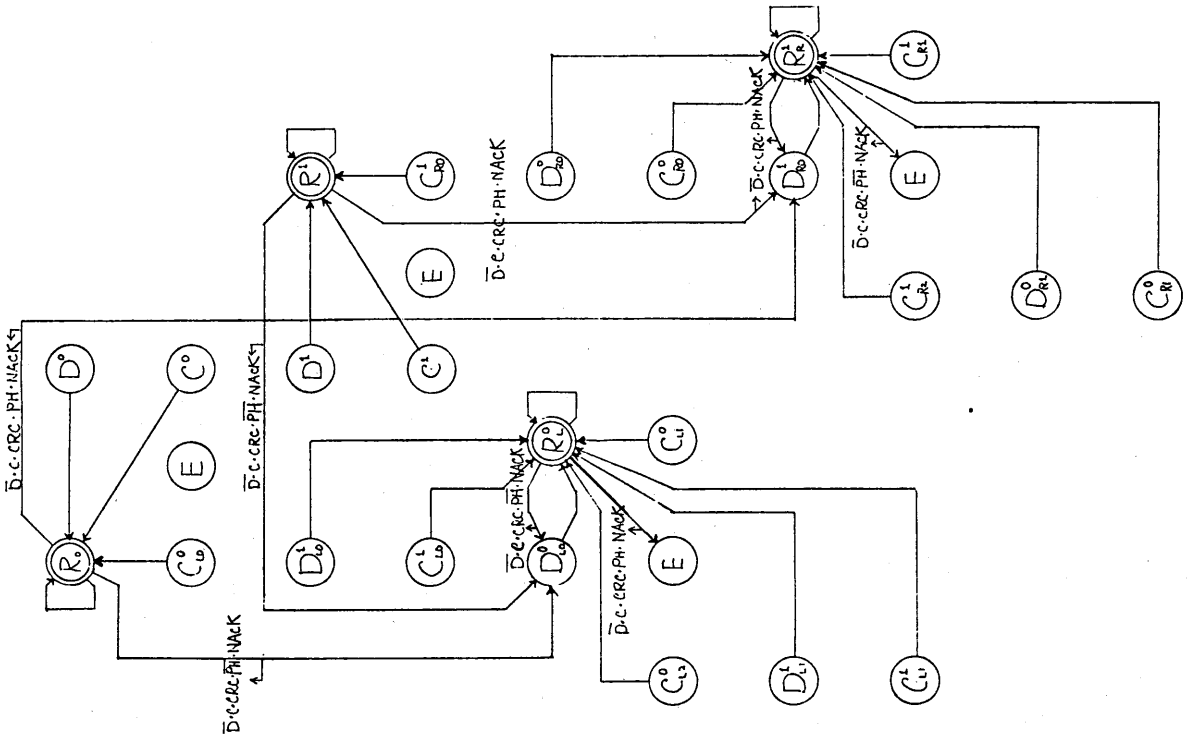


図7 Broadcast Hand Shaking
 プロセス間モジュール状態遷移図I

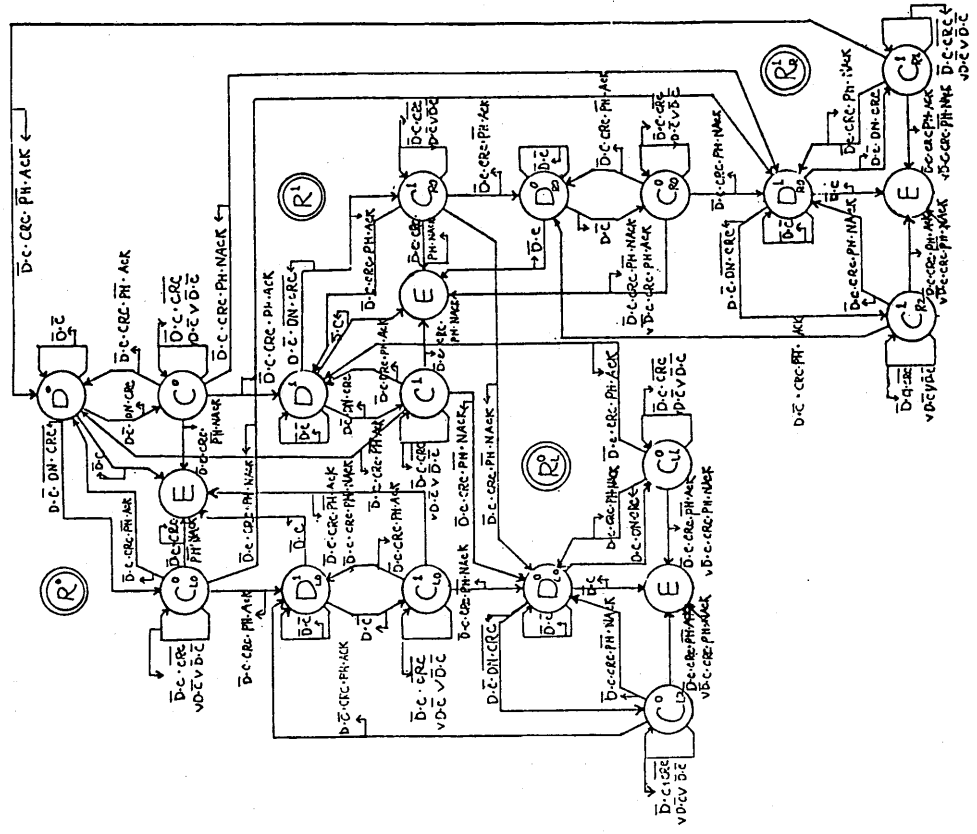


図 10 3-Phase Hand Shaking
7セグメントモジュール状態遷移図

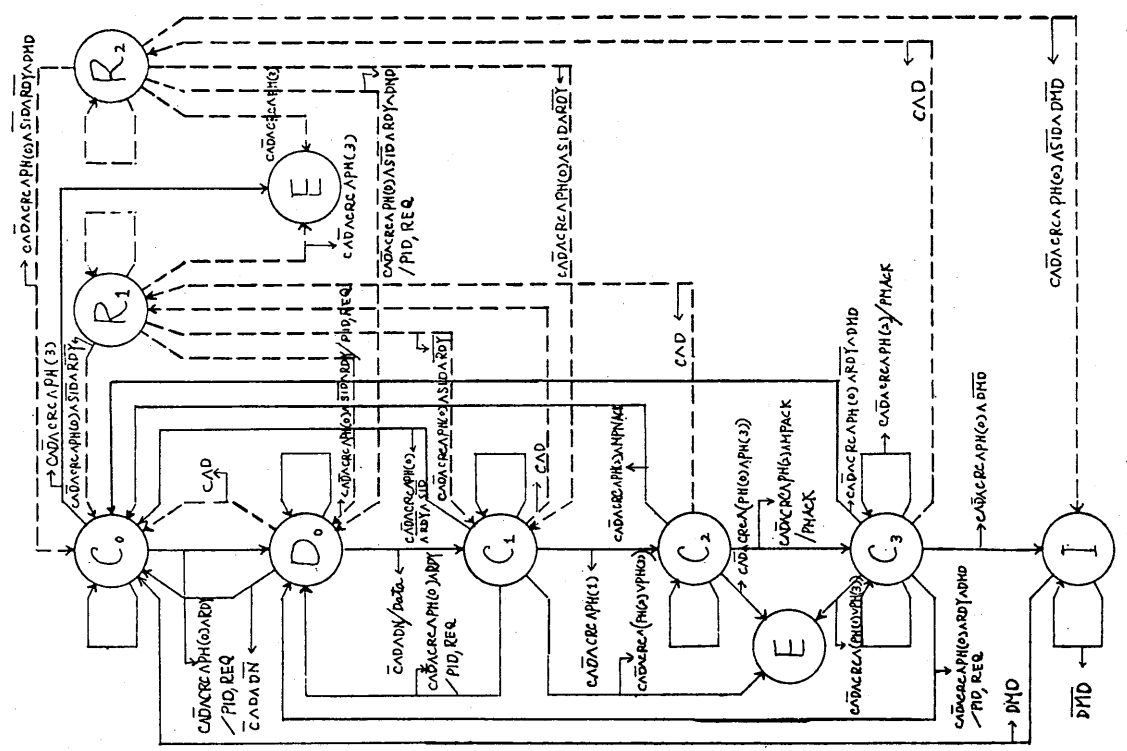


図 9 3-Phase Hand Shaking
Xモリモジュール状態遷移図

