

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4700773号
(P4700773)

(45) 発行日 平成23年6月15日 (2011. 6. 15)

(24) 登録日 平成23年3月11日 (2011. 3. 11)

(51) Int. Cl.

F I

G O 6 F 15/177 (2006. 01)

G O 6 F 15/177 6 8 2 F

G O 6 F 15/167 (2006. 01)

G O 6 F 15/177 6 8 2 J

G O 6 F 15/177 6 7 2 F

G O 6 F 15/167 B

請求項の数 25 外国語出願 (全 73 頁)

(21) 出願番号 特願平10-340924
 (22) 出願日 平成10年10月26日 (1998. 10. 26)
 (65) 公開番号 特開平11-282820
 (43) 公開日 平成11年10月15日 (1999. 10. 15)
 審査請求日 平成17年10月5日 (2005. 10. 5)
 審判番号 不服2008-19108 (P2008-19108/J1)
 審判請求日 平成20年7月28日 (2008. 7. 28)
 (31) 優先権主張番号 08/957298
 (32) 優先日 平成9年10月24日 (1997. 10. 24)
 (33) 優先権主張国 米国 (US)

(73) 特許権者 398038580
 ヒューレット・パッカード・カンパニー
 HEWLETT-PACKARD COM
 PANY
 アメリカ合衆国カリフォルニア州パロアル
 ト ハノーバー・ストリート 3000
 (74) 代理人 100059959
 弁理士 中村 稔
 (74) 代理人 100067013
 弁理士 大塚 文昭
 (74) 代理人 100084009
 弁理士 小川 信夫

最終頁に続く

(54) 【発明の名称】 スイッチをベースとするマルチプロセッサシステムに使用するための順序サポート機構

(57) 【特許請求の範囲】

【請求項 1】

複数の接続されたプロセッサノード及び共用メモリを有するコンピュータシステムであって、当該ノードのそれぞれが少なくとも1つのプロセッサと共用メモリの一部とを含んでおり、少なくとも一つのプロセッサは複数の、共用メモリにアクセスするためのメモリ要求コマンドを発生し、コンピュータシステムは、複数の通信チャンネルを有し、当該通信チャンネルは、ノードを相互接続し、コンピュータシステムを通して前記少なくとも一つのプロセッサによって発生されたメモリ要求コマンド及び一つあるいはそれより多いメモリ応答コマンドへの論理的に独立しているパスを与え、それに対する応答として発生されたデータ・メッセージを戻すものである、コンピュータシステムにおいて、

前記共用メモリへのメモリ要求コマンドに応じて発生された一つあるいはそれより多いメモリ応答コマンドを順序付ける、各ノードに配置されたシリアル化ポイントと、

第1の順序付けされた通信チャンネルであって、当該第1の順序付けされた通信チャンネルを通じて、前記順序付けされたメモリ応答コマンドが発送される、第1の順序付けされた通信チャンネルと、

第2の順序付けされていない通信チャンネルであって、当該第2の順序付けされていない通信チャンネルを通じて、戻りデータ・メッセージが発送される、第2の順序付けされていない通信チャンネルと、

各ノードに配置されたトランザクション追跡テーブルであって、当該トランザクション追跡テーブルは、保留中の要求に関する共用メモリのアドレス、コマンド、及びコマンド

10

20

IDを格納して、共用メモリからのデータに対して保留中の要求を識別し、順序付けされたメモリ応答コマンドと順序付けされていない戻りデータ・メッセージの間の相対的な順序を示すものであるトランザクション追跡テーブルと、を備え、

トランザクション追跡テーブルに格納されたコマンド及びコマンドIDから、トランザクション追跡テーブルに格納された保留中の要求と同じアドレスに対する後続の要求が、遅延されるべきか、無視されるべきかが決定されることを特徴とするコンピュータシステム。

【請求項2】

データに対する各保留中の要求に対して、前記トランザクション追跡テーブルに、順序付けされていない戻りデータ・メッセージが、関連する、順序付けされたメモリ応答コマンドを有することを示すための指示手段を更に備えている請求項1に記載のコンピュータシステム。

10

【請求項3】

指示手段は、トランザクション追跡テーブルへの少なくとも一つの順序付けされたチャンネル上で発生されたマーカー・メモリ応答コマンドを更に含む請求項2に記載のコンピュータシステム。

【請求項4】

トランザクション追跡テーブルは、複数のエントリを更に備え、

各エントリは、複数のノードの他の1つのメモリ位置のアドレスを記憶するためのものであり、

20

そして各エントリは、複数の状態ビットであって、その夫々が、関連する要求の対応する状態を指示するために選択的に設定され得る、複数の状態ビットを更に含む請求項3に記載のコンピュータシステム。

【請求項5】

トランザクション追跡テーブルは、順序付けされたチャンネル上のマーカー・コマンドがノードへ戻されたかどうかを指示する第1状態のビットセットを含む請求項4に記載のコンピュータシステム。

【請求項6】

トランザクション追跡テーブルは、

第2状態のビットセットであって、順序付けされていない戻りデータ・メッセージがノードへ戻されたかどうかを指示するために選択的に設定され得る、第2状態のビットセットと、

30

状態ビットセットの中の第1と第2ビットの両方を有するトランザクション追跡テーブルからエントリを取除く手段、

とを更に含む請求項5に記載のコンピュータシステム。

【請求項7】

トランザクション追跡テーブルに記憶されたアドレスへ発生されるメモリ要求コマンドであって、順序付けされたチャンネル上のマーカー・コマンドの受信を指示するためにそのアドレスに対応するトランザクション追跡テーブルの第1ビットがセットされる前に受け取られたメモリ要求コマンドを無視するための手段を更に備えた請求項6に記載のコンピュータシステム。

40

【請求項8】

無視されたメモリ要求コマンドは無効要求である請求項7に記載のコンピュータシステム。

【請求項9】

アドレスへ発生されたメモリ要求コマンドを無視する手段は、その要求を発生したプロセッサがそのアドレスをトランザクション追跡テーブルに入力させたプロセッサに対応する場合だけその要求を無視する請求項7に記載のコンピュータシステム。

【請求項10】

トランザクション追跡テーブルに記憶されたアドレスへ発生されるメモリ要求コマンド

50

を、順序付けされたチャンネル上のメモリ応答コマンドが受け取られるまで遅延する手段を更に備え、そのメモリ要求コマンドは、そのアドレスに対応するトランザクション追跡テーブルの第1ビットがセットされる前に受け取られている請求項6に記載のコンピュータシステム。

【請求項11】

アドレスに関連した所望のバージョンのデータがノードへ戻されるまでメモリ要求コマンドを更に遅延させる請求項10に記載のコンピュータシステム。

【請求項12】

アドレスをトランザクション追跡テーブルへ入れるようにした複数のプロセッサの一つへそのアドレスに関連したデータの所望のバージョンが戻されるまでメモリ要求コマンドを更に遅延させる請求項10に記載のコンピュータシステム。

【請求項13】

マルチプロセッサコンピュータシステムの共通のアドレスへ発生される複数のメモリ要求コマンド間の順序を維持する方法であって、前記のマルチプロセッサコンピュータシステムは、スイッチを経て接続された複数のマルチプロセッサノードを備え、マルチプロセッサノードの各々は、少なくとも2つのプロセッサと共用メモリの一部とを含んでおり、

前記マルチプロセッサコンピュータシステムは、複数の通信チャンネルを有し、該通信チャンネルは、前記マルチプロセッサノードを前記スイッチを経て相互接続し、前記マルチプロセッサコンピュータシステムを通して前記プロセッサの少なくとも1つによって発生されたメモリ要求コマンド、1つあるいはそれより多いメモリ応答コマンド、及び前記メモリ応答コマンドに対する応答として発生された戻りデータ・メッセージへの論理的に独立しているパスを与えるものであり、

前記複数の通信チャンネルが、

第1の順序付けされた通信チャンネルであって、該第1の順序付けされた通信チャンネルを通じて、前記順序付けされたメモリ応答コマンドが發送される、第1の順序付けされた通信チャンネルと、

第2の順序付けされていない通信チャンネルであって、該第2の順序付けされていない通信チャンネルを通じて、戻りデータ・メッセージが發送される、第2の順序付けされていない通信チャンネルと、から成り、

前記マルチプロセッサコンピュータシステムは、前記共用メモリへのメモリ要求コマンドに応じて発生された1つあるいはそれより多いメモリ応答コマンドを順序付ける、各マルチプロセッサノードに配置されたシリアル化ポイントを有する、方法において、

リモートマルチプロセッサノードの共用メモリの一部におけるそれぞれのアドレスに対してメモリ応答コマンドと順序付けられていない戻りデータ・メッセージの間の相対的な順序を識別するために、保留中の要求に関する共用メモリのアドレス、コマンド、及びコマンドIDを格納するための、各マルチプロセッサノードに配置されたトランザクション追跡テーブルにおいて、前記のマルチプロセッサノードの各々から前記のスイッチへ送られるメモリ要求コマンドのアドレスリストを維持し、

トランザクション追跡テーブルを通じて、要求の相対的な順序と関連するリモートなメモリ要求コマンドを追跡し、

トランザクション追跡テーブルの中で、メモリ要求コマンドの相対的な順序に対するリモートなメモリ要求コマンドのシリアル位置を識別する、ことを特徴とする方法。

【請求項14】

請求項13に記載の方法であって、

メモリ要求コマンドのアドレスのリストは、少なくとも一つの順序付けされたチャンネルにおける参照情報の順序を識別するものであり、

他のチャンネル上のトランザクションに対して、少なくとも1つの順序付けされたチャンネルにおける参照情報の識別された順序を再編成するステップを更に含む方法。

【請求項 15】

共用メモリの部分へのリモートな参照情報毎に、順序付けされた参照情報に対応するメモリの共用部分に関連したマルチプロセッサノードのアドレスリストへ、その順序付けされた参照情報が対応するリモートな参照情報を有することを指示する段階を更に含む請求項 14 に記載の方法。

【請求項 16】

指示段階が更に、少なくとも 1 つの順序付けされたチャンネル上でアドレスリストへコマンドを発生する段階を含む請求項 15 に記載の方法。

【請求項 17】

アドレスリストは複数のエントリを備え、各エントリは、別のマルチプロセッサノードのメモリ位置をアドレスする参照情報のアドレスを記憶するためのものであり、そして関連する要求の状態を指示するための複数の状態ビットを含んでいる請求項 16 に記載の方法。

10

【請求項 18】

状態ビットは、順序付けされたチャンネルのコマンドがマルチプロセッサノードに返送されたかどうかを示す第 1 ビットセットを更に含んでいる請求項 17 に記載の方法。

【請求項 19】

状態ビットは、リモートな参照情報がマルチプロセッサノードに返送されたかどうかを指示するための第 2 ビットセットを更に含み、そして前記の方法は、状態ビットのセットの第 1 ビットと第 2 ビットの両方を有するトランザクション追跡テーブルからエントリを除去する段階を更に含んでいる請求項 18 に記載の方法。

20

【請求項 20】

アドレスリストに記憶されたアドレスへ発生される要求であって、順序付けされたチャンネルにおけるコマンドの受信を指示するためにそのアドレスに対応するアドレスリストの第 1 ビットがセットされる前に受け取られた当該要求を無視する段階を更に含む請求項 19 に記載の方法。

【請求項 21】

無視される要求は無効の要求である請求項 20 に記載の方法。

【請求項 22】

アドレスへ発生された要求を無視する段階は、その要求を発生したプロセッサが、そのアドレスをアドレスリストに入力させたプロセッサに対応する場合だけ要求を無視する請求項 21 に記載の方法。

30

【請求項 23】

アドレスリストに記憶されたアドレスへ発生される参照情報を、順序付けされたチャンネルのコマンドが受け取られるまで遅延する段階を更に含み、当該参照情報は、そのアドレスに対応するアドレスリストの第 1 ビットがセットされる前に受け取られている請求項 22 に記載の方法。

【請求項 24】

参照情報は、アドレスに関連したデータの所望のバージョンがマルチプロセッサノードに返送されるまで更に遅延される請求項 23 に記載の方法。

40

【請求項 25】

参照情報は、アドレスに関連したデータの所望のバージョンが、そのアドレスをアドレスリストに入力させた複数のプロセッサの 1 つに返送されるまで更に遅延される請求項 24 に記載の方法。

【発明の詳細な説明】**【0001】****【発明の属する技術分野】**

本発明は、一般に、コンピュータアーキテクチャの分野に係り、より詳細には、分散型共用メモリマルチプロセッサシステムに係る。

【0002】

50

【従来の技術】

この分野で良く知られているように、対称型のマルチプロセッサコンピュータは、高性能のアプリケーション処理を行うことができる。通常の対称型マルチプロセッサコンピュータシステムは、バスによって互いに接続された多数のプロセッサを備えている。対称型マルチプロセッサシステムの1つの特徴は、メモリ空間が全てのプロセッサ間で共用されることである。1つ以上のオペレーティングシステムがメモリに記憶され、種々のプロセッサ間でのプロセッサ又はスレッドの分散を制御する。

異なるプロセッサ又はスレッドが多数の異なるプロセスを同時に実行できるようにすることにより、所与のアプリケーションの実行速度を著しく高めることができる。理論的に、システムの性能は、マルチプロセッサシステムにおけるプロセッサの台数を増加するだけで改善することができる。実際には、ある飽和点を越えてプロセッサを追加し続けると、単に通信ボトルネックが増えるだけとなり、従って、全システム性能を制限することになる。

10

【0003】

例えば、図1Aには、共通の相互接続バスを経て互いに接続された8個のプロセッサを含む典型的な公知のマルチプロセッサシステム2が示されている。動作中に、各プロセッサ3a-3hは、共用相互接続バス5を経て互いに他のプロセッサ及び共用メモリ4と通信する。図1Aの対称型マルチプロセッサ構成は、今日までに構築されたマルチプロセッサについて充分である。しかしながら、より高速のマイクロプロセッサの出現に伴い、通常の共用相互接続バスは、接続されたマイクロプロセッサの潜在的な全性能を十分に働かせることができない。プロセッサとメモリとの間の唯一の通信リンクは、共用バスであるから、バスはプロセッサからの要求で急速に飽和状態となり、各プロセッサがシステムバスへのアクセスを得よう試みるときに遅延が増大する。それ故、プロセッサは、高い速度で動作することができるが、性能に関する制限ファクタは、システムバスの使用可能な帯域である。

20

【0004】

通信帯域巾は、SMPシステムの性能において重要なファクタである。帯域巾は、SMPシステムにおけるノードの対又はサブセットの間で均一ではないから、業界では、SMPシステムの通信帯域巾を決定するために「二等分帯域巾」測定を使用している。二等分帯域巾は、次のように決定される。システムを等しい計算能力（等しいプロセッサ数）の2つの部分に区分化する全ての考えられる方法が確かめられている。各区分に対し、2つの区分間に維持し得る帯域巾が決定される。全ての維持し得る帯域巾の最小値は、相互接続の二等分帯域巾である。2つの区分間の最小帯域巾は、最悪の通信パターンが存在するときにマルチプロセッサシステムにより維持できる通信帯域巾を指示する。従って、大きな二等分帯域巾が望まれる。

30

【0005】

公知技術では、バス飽和の問題を克服するために、多数の相互接続アーキテクチャ即ち「トポロジ」が使用されている。これらのトポロジは、メッシュ、トーラス（円環体）、ハイパーキューブ（超立体）及び拡張ハイパーキューブを含む。

【0006】

40

【発明が解決しようとする課題】

例えば、メッシュ相互接続は、図1Bにシステム7として示されている。メッシュネットワークの主な利点は、簡単で且つ配線が容易なことである。各ノードは、少数の他の隣接ノードに接続される。しかしながら、メッシュ相互接続は、3つの重大な欠点を有する。第1に、メッセージは、それらの行先に到達するために平均的に多数のノードを横断しなければならない、その結果、通信の待ち時間が長くなる。第2に、二等分帯域巾は、他のトポロジに対するものであるから、メッシュトポロジの場合に十分に計測しない。最後に、各メッセージはメッシュ内の異なる経路を進行するので、SMPシステム内には自然の順序付けポイントが存在せず、それ故、メッシュトポロジの実施を必要とするキャッシュコヒレンスプロトコルがしばしば非常に複雑なものとなる。

50

【 0 0 0 7 】

トラス、ハイパーキューブ及び拡張ハイパーキューブトポロジは、全て、ノードが種々の複雑な構成、例えば円環体構成又は立体構成で相互接続されたトポロジである。トラス、ハイパーキューブ及び拡張ハイパーキューブの相互接続は、メッシュ相互接続よりも複雑であるが、その待ち時間及び帯域巾は、メッシュ相互接続よりも優れている。しかしながら、メッシュ相互接続と同様に、トラス、ハイパーキューブ及び拡張ハイパーキューブトポロジは、自然の順序付けポイントを与えず、従って、これらのシステムの各々に対して複雑なキャッシュコヒレンスプロトコルを実施しなければならない。

共用メモリのマルチプロセッサシステムでは、プロセッサは、通常、将来アクセスされる見込みが高いと決定されたデータを記憶するために専用キャッシュを使用している。プロセッサは、それらの専用キャッシュからデータを読み取りそしてメモリへ書き戻すことなく専用キャッシュにおいてデータを更新するので、各プロセッサの専用キャッシュが一貫して即ちコヒレントに保持されるよう確保するための機構が必要となる。SMPシステムのデータのコヒレンス性を確保するのに使用される機構は、キャッシュコヒレンスプロトコルと称される。

10

【 0 0 0 8 】

物理的な相互接続部のトポロジ、帯域巾及び待ち時間に加えて、キャッシュコヒレンスプロトコルの効率も、システム性能の重要なファクタである。キャッシュコヒレンスプロトコルは、待ち時間、ボトルネック、非効率性又は複雑さを多数の仕方で導入する。

ロード及び記憶動作の待ち時間は、設計のプロトコルによって直接影響されることがしばしばある。例えば、あるプロトコルでは、全ての無効化メッセージがそれらのターゲットプロセッサへ送られそして確認メッセージがその元のプロセッサへ完全に返送されるまで記憶動作が完了したとみなされない。従って、記憶の待ち時間は、無効化がその行先へ送られるのを元のプロセッサが待機しなくてよいプロトコルよりも相当に長いものとなる。更に、確認は、システム帯域巾の相当の部分を消費する。

20

【 0 0 0 9 】

ボトルネックは、コントローラの高い占有度によりしばしば生じる。「占有度」とは、コントローラが要求を受け取った後に使用できなくなる時間の長さを示す用語である。あるプロトコルでは、直接的なコントローラは、メモリ位置に対応する要求を受け取ると、その前のコマンドに対応するある確認がディレクトリに到着するまで同じメモリ位置への他の要求に対して使用できなくなる。コントローラは、平均より高いレートで競合する要求を受け取る場合に、ボトルネックとなる。

又、キャッシュコヒレンスプロトコルの設計は、ハードウェアの複雑さにも影響する。例えば、あるプロトコルは、停滞及び公正さの問題を招き、これらは、付加的な機構で対処される。その結果、ハードウェアの複雑さが増大する。

30

【 0 0 1 0 】

そこで、オペレーションの待ち時間を最小にし、広い通信帯域巾を与え、コントローラの占有度を低くし、そして多数のプロセッサへと拡張することのできる対称的なマルチプロセッサシステムを提供することが要望される。

【 0 0 1 1 】

【課題を解決するための手段】

本発明は、少なくとも1つのプロセッサ及び共用メモリの一部分を含む多数のマルチプロセッサノードがスイッチを経て互いに接続された対称的なマルチプロセッサシステムに効果的に使用される。マルチプロセッサノードの各々にはトランザクション追跡(トラッキング)テーブル(TTT)が維持される。TTTは、ノードをスイッチに接続するノードのグローバルポートに存在してもよいし、或いはマルチプロセッサノードの少なくとも1つのプロセッサの各々に存在してもよい。

TTTは、マルチプロセッサノードから発生され及びそれにより受け取られる要求の順序を決定しそしてそれを強制するのに使用される。本発明の1つの特徴によれば、TTTは、マルチプロセッサノードへ返送される要求の順序を次のように決定するのに使用される

40

50

。各要求は、多数のトランザクションに細分化され、各トランザクションは、異なる仮想チャンネルを経て搬送される。少なくとも1つのチャンネルが順序付けされるが、他のチャンネルの返送データは、ばらばらの順序で受け取ることができる。コヒレンス性を維持するために、共通のアドレスへ発生されるトランザクションが順序正しく取り扱われるのが望ましい。本発明の1つの特徴によれば、マーカーパケットが、順序付けされたチャンネルを経てT T Tへ発生され、アドレスに関連したデータが依然処理されていることを指示する。このような構成では、T T Tは、マーカーパケットに続いて受け取られた上記順序付けされたチャンネルの他の要求を無視するか、又はデータが受け取られるまで遅延するように確保する。

【0012】

従って、本発明の1つの特徴によれば、各々少なくとも1つのプロセッサ及び共用メモリの一部を含む複数の接続されたマルチプロセッサノードを有するコンピュータシステムは、上記複数のマルチプロセッサノードの各々における複数のプロセッサに関連した追跡機構であって、上記複数のマルチプロセッサノードの1つにおける少なくとも1つのプロセッサの1つにより発生された共用メモリのリモート部分のアドレスへの要求の位置を、上記複数の接続されたマルチプロセッサノードにおける少なくとも1つのプロセッサにより上記アドレスへ発生された複数の他の要求に対して識別するための追跡機構を備えている。

本発明の更に別の特徴によれば、マルチプロセッサコンピュータシステムの共通のアドレスへ発生される複数の要求間の順序を維持するための方法が提供される。マルチプロセッサコンピュータシステムは、スイッチを経て接続された複数のマルチプロセッサノードを備え、各マルチプロセッサノードは、少なくとも1つのプロセッサ及び共用メモリの一部を含む。上記方法は、マルチプロセッサノードの各々からスイッチへ送られる要求のアドレスリストを維持して、リモートマルチプロセッサノードの共用メモリの一部における各アドレスに対して要求の相対的な順序を識別し、アドレスに関連した要求が満足されるまでアドレスをリストに維持するという段階を含む。

【0013】

【発明の実施の形態】

本発明の上記及び他の特徴は、添付図面を参照した以下の詳細な説明から明らかとなる。

本発明の1つの実施形態によれば、ハイアラキー式の対称的マルチプロセッサ(SMP)システムは、高性能スイッチを経て互いに接続された多数のSMPノードを備えている。従って、SMPノードの各々は、SMPシステムにおいてビルディングブロックとして働く。以下、1つのSMPノードビルディングブロックの要素及び動作を最初に説明し、その後、SMPシステムの動作を説明し、それに続いて、大規模のSMPシステムにおいてメモリのコヒレンス性を維持するために使用されるキャッシュコヒレンスプロトコルを説明する。

【0014】

SMPノードビルディングブロック

図2を参照すれば、マルチプロセッサノード10は、4つのプロセッサモジュール12a、12b、12c及び12dを備えている。各プロセッサモジュールは、中央処理ユニット(CPU)を備えている。好ましい実施形態では、デジタル・イクイップメント社で製造されたAlpha(登録商標)21264プロセッサチップが使用されるが、以下に述べるコヒレンスプロトコルをサポートすることのできるものであれば、他の形式のプロセッサチップも使用できる。

マルチプロセッサノード10は、多数のメモリモジュール13a-13dを含むメモリ13を備えている。このメモリは、32ギガバイトの記憶容量を備え、4つのメモリモジュールの各々が8ギガバイトを記憶する。各メモリモジュールは、多数のメモリブロックに分割され、各ブロックは、例えば、64バイトのデータを含む。データは、一般に、メモリからブロックで検索される。

【 0 0 1 5 】

更に、マルチプロセッサノード 1 0 は、接続された I / O バス 1 4 a を経て外部装置（図示せず）とマルチプロセッサノード 1 0 との間で行われるデータ転送を制御するための I / O プロセッサ（I O P）モジュール 1 4 を備えている。本発明の 1 つの実施形態では、I / O バスは、周辺コンピュータ相互接続（P C I）プロトコルに基づいて動作する。I O P 1 4 は、I O P キャッシュ 1 4 c 及び I O P タグ記憶装置 1 4 b を含む。I O P キャッシュ 1 4 c は、P C I バス 1 4 a を経て外部装置へ転送されるメモリ 1 3 からのデータのための一時的な記憶装置である。I O P タグ記憶装置 1 4 b は、外部装置とプロセッサとメモリとの間に移動されるデータに対するコヒレンス情報を記憶するための 6 4 エントリのタグ記憶装置である。

10

【 0 0 1 6 】

マルチプロセッサノードのメモリ 1 3 に記憶されたデータのコヒレンス性は、デューブリケートタグ記憶装置（D R A G）2 0 によって維持される。D T A G 2 0 は、全てのプロセッサ 1 2 a - 1 2 d により共用され、そして 4 つのバンクに分割される。各バンクは、関連するプロセッサにより使用されるデータに対応する状態情報を専用に記憶する。

D T A G、メモリ及び I O P は、A R B バス 1 7 と称する論理バスに接続される。プロセッサにより発生されるメモリブロック要求は、ローカルスイッチ 1 5 を経て A R B バス 1 7 にルート指定される。D T A G 2 0 及び I O P 1 4 は、プロセッサ及び I O P のキャッシュにおけるブロックの状態をルックアップし、そしてメモリブロックに対しそれらの状態を原子的に更新する。A R B バス 1 7 は、全てのメモリ参照に対してシリアル化ポイントとして働く。メモリ要求が A R B バスに現れる順序は、プロセッサが要求の結果を認知する順序である。

20

【 0 0 1 7 】

プロセッサモジュール 1 2 a - 1 2 d、メモリモジュール 1 3 a - 1 3 d 及び I O P モジュール 1 4 は、ローカルの 9 ポートスイッチ 1 5 を経て互いに接続される。インターフェイスモジュール 1 2 a - 1 2 d、1 3 a - 1 3 d 及び 1 4 の各々は、同数の両方向性クロック送信データリンク 1 6 a - 1 6 i によりローカルスイッチに接続される。1 つの実施形態では、データリンクの各々は、1 5 0 M H z のレートで動作するシステムクロックの各縁で 6 4 ビットのデータ及び 8 ビットのエラー修正コード（E C C）を送信する。従って、データリンク 1 6 a - 1 6 i の各々のデータ帯域巾は、2 . 4 ギガバイト / s である。

30

ローカルスイッチ 1 5 は、クオドスイッチアドレス制御チップ（Q S A チップ）1 8 及びクオドスイッチデータスライスチップ（Q S D チップ）1 9 を備えている。Q S A チップ 1 8 は、プロセッサモジュール I O P とメモリとの間のアドレス経路を制御するためのアービター（Q S A R B）1 1 を備えている。更に、Q S A チップ 1 8 は、以下に述べるようにローカルスイッチ 1 5 を通るデータの流れを制御するために Q S D チップ 1 9 を制御する。Q S D チップ 1 9 は、プロセッサモジュールと、メモリモジュールと、I O P との間の全てのデータ経路に対するスイッチ相互接続を与える。図 2 には示されていないが、以下に述べるように、マルチプロセッサノード 1 0 がグローバルポートを経て他のマルチプロセッサノードに接続された場合には、Q S D 及び Q S A がグローバルポートに対するスイッチ相互接続部を付加的に形成する。各プロセッサは、メモリデバイス 1 3 a - 1 3 d、他のプロセッサ 1 2 a - 1 2 d、I O P 1 4 のような使用可能なリソースの 1 つからデータを要求することもできるし、或いは他のマルチプロセッサノードのリソースからグローバルポートを経てデータを要求することもできる。従って、ローカルスイッチ 1 5 は、2 . 4 ギガバイトの広いバス帯域巾を維持しながら、種々のリソースから同時入力を受け入れることができねばならない。

40

【 0 0 1 8 】

ローカルスイッチは、多数の同時トランザクションを取り扱うことができる。各トランザクションは、通常、多数のリソース（メモリバンクや、データ経路や、待ち行列のような）を使用するので、ローカルスイッチの制御機能は非常に複雑になる。例えば、あるトラ

50

ンザクションは、そのトランザクションの段階0でメモリバンクを使用でき、段階1でメモリバンクからプロセッサポートへのデータ経路を使用でき、そして段階2でプロセッサポートからプロセッサへのデータ経路を使用できることを必要とする。ローカルスイッチアービター（Q S A 1 8のQ S A A R B 1 1）は、あるトランザクションが開始されると、各段階でトランザクションにより必要とされるリソースが必要に応じて使用できるように要求を裁定する。

【0019】

より重要なことに、アービターは、特定の要求が、他の要求の進行中に長時間にわたり（潜在的に不定に）裁定に負けることのないよう確保することにより、全ての要求及びプロセッサがリソースに対して公平なアクセスを得るように保証する。例えば、3つのリソースA、B及びCを要求するトランザクションTについて考える。このトランザクションTは、トランザクションの適当な段階に3つのリソース全部が使用できるよう保証されるまで裁定に勝てない。リソースが使用可能であることのみに基づいてアービターがその判断を行う場合には、トランザクションTは、A、B又はCの1つを使用する（他のリソースD、E等と共に）他のトランザクションが裁定に勝ち続ける間は、長時間にわたって成功しないことが考えられる。

【0020】

各々が多数のリソースを使用して完了するような非常に多数の同時要求を伴うスイッチにおいて公平な裁定を保証するのは、計算上複雑である上に、高速データ経路において遅延を増加し勝ちである。ここに示す装置においては、Q S A A R B 1 1が、特定のトランザクションをスケジュールする前に、1つのリソース（メモリバンク）のみについて裁定を行う。プロセッサに通じる待ち行列である第2のリソースは、Q S A A R B 1 1により第1のリソースについて裁定を行うときに、それが使用可能であるかどうかについてチェックする必要がない。というのは、Q S Dのアーキテクチャーがそのデータ経路を保証しそして待ち行列に通じる待ち行列スロットが常に使用できるからである。リソースに対する公平な裁定は、Q S A A R B 1 1に著しい複雑さを伴うことなく与えられる。

【0021】

本発明の1つの実施形態によれば、Q S Dは、対応する行先に通じるバッファに対してアップフロント裁定を必要とせず全てのリソース（プロセッサ、メモリ、I O P及びグローバルポート）からの入力を同時に受け取ることができる。次いで、全てのデータリソースは、データ経路又はスイッチにおける待ち行列スロットへのアクセスを裁定する必要なく、スイッチへデータを独立して送信することができる。というのは、Q S Dは、全てのリソースからのデータを実質的に同時に受信することのできる多数の同時挿入バッファを備えているからである。同時挿入バッファの2つの実施形態を以下に説明する。

【0022】

同時挿入バッファスイッチ

上記のように、マルチプロセッサノードにおけるプロセッサ12a - 12d、I O P 14及びメモリデバイス13a - 13dの各々は、マルチプロセッサノードのプロセッサ及びI O Pからの要求を取り扱うためのリソースとして働く。データは、各リソース要素と、要求を発する要素との間でパケットの形態で転送される。各パケットは、512ビットのデータと、64ビットのE C Cとを含む。上記したように、各データリンクは、64ビットのデータ及び8ビットのE C Cを150MHzクロックの各縁において搬送する。従って、Q S Dの外部には、パケット当たり8個のデータ転送サイクルがある。しかしながら、Q S Dの内部では、クロックの1つの縁においてのみデータが収集される。従って、Q S Dの内部のロジックの各クロックサイクルに対し、潜在的に128ビットのデータがデータリンクから受け取られる。各パケットは、512ビットのデータ及び64ビットのE C Cを含むので、Q S Dの内部では、各パケットごとに4つのデータ転送サイクルがあり、各Q S Dクロックサイクルに、128ビットのデータ及び16ビットのE C Cがプロセッサ、I O P又はメモリデバイスからQ S Dへ転送される。

【0023】

10

20

30

40

50

図3を参照すれば、QSD19は、5つの同時挿入バッファ(SIB)25a-25eを含むように詳細に示されている。各SIBは、要求側要素、即ちプロセッサ12a-12d又はIOPの1つに専用である。各SIBは、それに関連した要求側要素と、ノード内の他のリソース要素、即ちプロセッサ12a-12d、メモリ13a-13d、IOP14及び好ましくはグローバルポートとの間でパケットを転送するためのデータ経路を制御する。グローバルポートは、他のマルチプロセッサノードへの相互接続部として働き、以下に詳細に説明する。SIBは、スイッチへのアクセスに対し要求側要素間の裁定を必要とせず、スイッチに接続されたいずれのリソースからでも要求側要素によりパケットを同時に受信できるようにする。

【0024】

既に述べたように、QSAARB11は、スイッチ19への制御を与えるように接続される。QSAARB11には、メインアービター27が含まれる。このメインアービター27は、リソース(IOP、プロセッサ12a-12d及びメモリ13a-13d)とスイッチ19との間のデータの移動を管理する。プロセッサ12a-12d及びIOP14の各々は、ライン28a-28e上のリソースの1つへアクセスするための要求を発生し、これらは、メインアービター27に送られる。次いで、メインアービターは、各リソースが要求を受信できるときにこれらの要求をその関連リソースへ送る。リソースが要求を受け取るときに、スイッチ19の裁定は必要とされない。というのは、SIBの各々は、全ての入力からの入力を実質的に同時に即ち同じデータサイクル内に受け取ることができるからである。又、QSAARB11には、多数の個々のアービター23a-23dも含まれる。これらアービター23a-23dの各々は、プロセッサ12a-12dの関連する1つと、それに対応するSIB25b-25eとの間のデータ路を管理するのに使用される。IOP14とSIB25aとの間のデータ路を管理するために、IOP14には同様のアービター(図示せず)が含まれる。各プロセッサは、その関連SIBからデータを受け取ることができるので、その関連アービターは、接続されたデータ路にデータを送信する。

【0025】

従って、スイッチ19内の同時挿入バッファの使用により、要求側要素とリソースとの間の裁定経路は、2つの別々の区分に分割される。即ち、接続されたりソースからデータを受け取るために要求を発生しているプロセッサが使用できるかどうかに関わりなくプロセッサからの要求に応答してメインアービター27がリソースを裁定するところの第1裁定区分と、プロセッサに関連したアービターがスイッチからのデータを送信するためにプロセッサへのアクセスを裁定するところの第2裁定区分である。このような構成では、裁定が分離されるために、接続されたりソース各々への公平なアクセスが与えられるように保証することができる。

図4Aを参照すれば、SIB25aの1つの実施形態が詳細に示されており、これは、ライン36aを経て8個の接続されたマルチプレクサ34a-34hにMUX選択信号<31:0>を与えるように接続された入力アービター36を備え、MUX選択信号の4つが8個のマルチプレクサの各々に送られて、各マルチプレクサの9個の入力の1つが選択される。SIB25a-25dは全て同様の構造にされ、従って、その1つについてのみ詳細に説明する。上記したように、潜在的に10個のリソースがSIBに接続される。10個のリソースの1つは、SIBから出力を受信する要求側デバイスであり、一方、他の9個のリソースは、SIBに入力を与える。それ故、マルチプレクサ34a-34hの各々は、SIBに接続された9個のリソースから入力を受け取る。接続されたプロセッサの3つからの入力は、ラインPx、Py及びPzを経て受け取られる。第4のプロセッサ(SIBがIOPデバイスに関連するとき)又はIOPデバイス(SIBが1つのプロセッサに関連するとき)からの別の入力はラインPW/IOPを経て受け取られる。メモリバンク13a-13dからの入力は、各々、ラインmem0、mem1、mem2及びmem3を経て受け取られ、そしてグローバルポートからの入力は、グローバルポートラインを経て受け取られる。

10

20

30

40

50

【 0 0 2 6 】

マルチプレクサ 3 4 a - 3 4 h の各々からの各出力は、バッファ 3 2 の 8 個のバンクの 1 つに接続される。各バンクは 8 個のエントリを有し、各エントリは、1 2 8 ビットのデータ及び 1 6 ビットの E C C を記憶する。従って、S I B により受信されるデータの各パケットは、バッファ 3 2 の同じ行において 4 つの異なるバンクに書き込まれる。以下に述べるように、入力アービター 3 6 は、データを記憶するのに使用できるバッファのバンクを指示するための状態ビットを維持する。従って、1 つ以上のリソースから 1 2 8 ビットのパケットデータが受け取られる各サイクルごとに、入力アービター 3 6 は、バンクの使用状態に基づいて関連バンク 3 2 a - 3 2 h へパケットデータのサイクルを送信するために、各マルチプレクサ 3 4 a - 3 4 h における考えられる 9 個のリソース入力の 1 つを選択する。又、入力アービターは、ライン 3 6 b を経てマルチプレクサ 3 0 へバイパスデータも与える。入力アービターの状態ビットが、全てのバンク 3 2 a - 3 2 h が空であることを指示するときには、9 個のリソース入力の 1 つが入力アービター 3 6 を経て関連する要求側要素へ直接バイパスされる。

10

【 0 0 2 7 】

バンク 3 2 a - 3 2 h の各々は、マルチプレクサ 3 0 に接続される。マルチプレクサ 3 0 は、出力アービター 3 8 により制御される。S I B 2 5 a に関連する要求側要素が S I B からデータを受け取る準備ができ、そしてパケットの一部分が S I B のエントリに書き込まれると、出力アービターは、バンク 3 2 a - 3 2 h から要求側要素に 8 個のエントリの 1 つを供給する。或いは又、出力アービターは、いずれのバンクも転送保留データをもたずそして入力アービターからライン 3 6 b を経てデータが得られる場合には、ライン 3 6 b を経て要求側要素にバイパスデータを供給する。

20

動作中に、パケットデータの第 1 の 1 2 8 ビットが S I B に受け取られたときに、8 個のバンクの 1 つが、パケットデータの第 1 の 1 2 8 ビットを記憶するために選択される。本発明の 1 つの実施形態によれば、パケットデータの 1 2 8 ビットが受け取られる次の 3 サイクルの各々の間に、手前の書き込みを実行するのに使用したバンクに隣接するバンクが、パケットデータの次の 1 2 8 ビットを書き込むのに選択される。例えば、バンク 3 2 a が、ソース m e m 0 からパケットデータの第 1 サイクルを書き込むのに使用できるバンクとして選択された場合には、パケットデータの第 2 サイクルはバンク 3 2 b に書き込まれ、第 3 サイクルはバンク 3 2 c に書き込まれ、そして第 4 サイクルはバンク 3 2 d に書き込まれる。従って、パケットデータのその後のサイクルを書き込むためにどのバンクを使用すべきかの選択は、入力アービターにより選択されたバンクでスタートしそして各次々のパケット書き込みに対して隣接バンクに続くようにして回転ベースで実行される。その結果、受け取られたパケットがバッファ 3 2 の共通の行における 4 つのバンクにわたって分散される。

30

【 0 0 2 8 】

8 個のバンクが設けられ、そして本発明の 1 つの実施形態では、いずれの要求側要素においても保留となり得るリソース読み取りの最大数は 8 であるから、各書き込みサイクルの間に各リソースに対して少なくとも 1 つのバンクを使用できることが確保される。それ故、所与の瞬間に、全部で 8 個の保留の読み取り応答がスイッチによって受け取られた場合に、バンク 3 2 a - 3 2 h の各々を使用して、第 1 のパケットデータ書き込みサイクルを受け入れることができ、バンクの選択は、次の 3 つの書き込みサイクルについて回転される。

40

本発明の 1 つの実施形態では、S I B の各バッファは、先入れ先出し (F I F O) プロトコルのもとで動作する。パケットの 2 つの部分が同時に受け取られるので、それらに対しスイッチへ「読み込まれる」順序が選択される。リソースに対して裁定を行う要求側要素のロジックは、S I B と通信せず、そしてリソースに対して裁定するための他の要求側要素とも通信しないので、標準的なルールに従ってデータの完全性を確保する。例えば、リソースに固定の優先順位番号が指定される場合には、「低い番号の入力リソースからのデータが、常に、高い番号の入力リソースからのデータの前にスイッチに書き込まれる」と

50

というようなルールに従う。

【 0 0 2 9 】

上記のように、図 4 A に示す S I B の実施形態では、8 個のバンクの使用について説明した。というのは、要求側要素が所与の瞬間にもつことのできる保留メモリ要求の数が 8 に対応するからである。しかしながら、設計上の制約により、それより少数のバンクを設けることが必要な場合には、インターリーブ又は同様の技術を使用して多数のデータチャンクを共通のバンクの異なる位置に同時に書き込みできるように当業者によって容易に設計を変更することができよう。それ故、本発明は、図 4 A に示す特定の実施形態に限定されるものではない。

上記のように、動作中に、入力アービターは、リソースからデータを書き込むための適当なバンクを選択するためにバンクにおける入力の利用性に関する状態情報を維持する。S I B への入力を制御するための入力アービター 3 6 の実施形態が図 4 B に示されている。上記では 9 個の入力リソースについて述べたが、図 4 B には、明瞭化のために、2 つのリソース入力のための書き込みを制御するロジックが示されている。入力パケットデータがライン 3 5 を経て受け取られるときに、「入力 1」のような指示信号がラッチチェーン 4 0 に送られ、このラッチチェーンは、4 個のラッチ、フリップ・フロップ又は同様の状態装置を含む。ラッチチェーン 4 0 は、カウンタ機構として使用される。この例の目的として、4 つの次々のデータ転送サイクルにパケットデータが受け取られると仮定する。4 つのデータ転送サイクルの間に、入力信号はラッチチェーンを経て伝播する。ラッチチェーンには、オアゲート 4 6 が接続される。入力値がラッチチェーン 4 0 を経て伝播するときに、オアゲート 4 6 の出力がアサートされる。

【 0 0 3 0 】

オアゲート 4 6 の出力は、シフトレジスタ 4 8 へのシフト信号を与える。シフトレジスタは、S I B の各バンクについて 1 つずつ、8 個のビット位置を含む。シフトレジスタ 4 8 は、入力信号サンプルを最初に受信する際に、バンク選択ロジック 4 4 からのビットベクトルがロードされる。バンク選択ロジック 4 4 から受け取られたビットベクトルは、1 ビットがセットされるだけであり、ベクトル内のビットの相対的な位置が、パケットデータの書き込みを開始すべきバンクを指示する。

従って、バンク選択ロジック 4 4 は、パケットデータの第 1 サイクルの書き込み行先を制御する。バンク選択ロジック 4 4 は、利用性ベクトル 4 2 を入力として受け取り、利用性ベクトルにおけるビットの相対的な位置が、書き込みデータを受け取ることのできない関連バッファを指示する。

【 0 0 3 1 】

バンク選択ロジックがシフトレジスタ 4 8 へビットを与えるときに、シフトレジスタ 4 8 の値がデマルチプレクサ 4 9 に送られる。又、デマルチプレクサ 4 9 は、入力 1 ソースが接続されるところのマルチプレクサ 3 4 a - 3 4 h の入力の数値表示も入力として受け取る。例えば、デマルチプレクサ 4 9 は、「1」のマルチプレクサ選択値を用いてマルチプレクサ 3 4 a を経て入力 1 リソースデータが送られることを指示する「1」入力値を受け取る。選択されたバンクを指示するシフトレジスタ内のビットの位置に基づいて、値「1」が M U X 選択 < 3 1 : 0 > 信号 3 6 a の適当な位置へ伝播される。各入力ソースに対する各デマルチプレクサは、全ての M U X 選択信号を駆動し、それらの出力は、これらの信号がマルチプレクサ 3 4 a - 3 4 h を駆動する前にオアされる。

【 0 0 3 2 】

バンクエントリの書き込みの後に、シフトレジスタの内容がオアゲート 5 0 によりオアされ、利用性バンクベクトル 4 2 として記憶される。これは、次のサイクルの間に、どのバンクが到来する書き込みに対して使用できるかをバンク選択ロジック 4 4 により決定するために使用される。

ライン 4 6 a のシフト信号がアサートされる各サイクルに、シフトレジスタ 4 8 のビットが右へシフトされる。ビットが右へシフトするときには、M U X 選択信号 < 3 1 : 0 > の選択値も右へシフトされ、次の書き込み動作中に入力ソースを次の隣接バンクへ供給する

ようにさせる。

従って、ローカルQSDスイッチ内のSIBを使用することにより、多数の同時に受け取られた入力がある行先である要求側要素へ到達するように確保できる簡単且つ効率的なスイッチング機構が設けられる。このような構成では、リソースへのアクセスに対してソースがいったん裁定されると、ソースにより実行されねばならない全ての裁定が完了する。ソースは、リソースが常にスイッチバッファ32へのアクセスを得ることができるという事実依存する。ソースアービターが互いに独立して動作してリソースを管理することにより、最小限の複雑さで公平な裁定を確保する機構が設けられる。更に、SIBは、要求側要素の最大数の保留中読み取りに対してデータを記憶できるので、たとえ全てのリソースからデータが同時に受け取られても、バッファ32に対するリソースを裁定する必要はなく、リソースロジックの全体的な複雑さが低減される。

10

【0033】

図5には、図3に示すようにプロセッサ又はIOPデバイス(キャッシュを含む任意の要求側デバイス)へ接続することのできる同時挿入バッファ(SIB)61の第2の実施形態が示されている。SIB61は、9個のマルチプレクサ60a-60iを含み、そのうちの8個は、8個のバッファ62a-62hの各々に接続される。第9マルチプレクサ60iは、以下に述べるようにバイパス経路を与えるのに使用される。マルチプレクサ60a-60iの各々は、接続されたメモリデバイスmem0-mem3からの4つの入力、グローバルポートからの1つの入力、接続されたプロセッサからラインPx、Py及びPzを経て送られる3つの入力、そしてIOP(SIBに関連したデバイスがプロセッサの場合)又は別のプロセッサ(SIBに関連したデバイスがIOPの場合)からラインPW/IOPを経て送られる1つの入力を含む9つの入力を受け取る。

20

【0034】

バッファ62a-62hの各々は、4つの128ビットエントリを含む。従って、各入力バッファは、SIBにおいて次々のサイクル中に4つの128ビット部分で受け取られた1つの512ビット情報パケットを記憶する。各バッファには、4対1のマルチプレクサ64a-64hが各々接続される。これらのマルチプレクサ64a-64hは、関連バッファの4つの入力のうちの1つを選択して、マルチプレクサ66を経てSIBの出力へ供給するのに使用される。

図4Aについて上述したように、本発明の1つの実施形態では、各要求側要素がいかなる所与の瞬間にも異なるリソースに対してせいぜい8個の保留中読み取り参照を有するだけであるから、8個のバッファが含まれる。従って、図5には8個のバッファが示されているが、本発明はこれに限定されるものではない。むしろ、選択されるバッファの数は、関連するプロセッサ又はIOPデバイスのバッファ特性に依存する。

30

【0035】

動作中に、接続されたりソースの各々から入力を受け取られるときに、入力アービター67は、各マルチプレクサにおける入力ラインの1つを選択し、データの packets を空きバッファへ供給する。所与のリソースからのパケット書き込みの時間中に同じバッファが選択され、パケットの全ての部分が単一のバッファに維持される。パケットの少なくとも1つの部分がバッファに書き込まれると、それがマルチプレクサ66に送られ、関連する要求側要素の準備ができたときにその要求側要素へ供給される。或いは又、いずれのバッファにもパケットデータが存在しない場合には、マルチプレクサ60iを経、マルチプレクサ66を経てパケットデータを出力へ直接的に供給することによりバイパス経路を選択することができる。

40

【0036】

8個のバッファが設けられるので、SIBデバイス61は、接続されたりソースの各々から実質的に同時に(即ち、同じデータサイクルに)データを受け取ることができる。QSDにSIBを使用することにより、前記の実施形態の場合のように、SIBへのアクセスに対し要求側要素の間に裁定は必要とされない。その結果、リソースがローカルスイッチを使用する準備ができたときにローカルスイッチの利用性が保証される。更に、本来的に

50

公平な裁定機構が設けられる。というのは、スイッチに対する裁定の結果としてリソースへの要求が他のリソースへの他の要求により阻止されないからである。従って、裁定の複雑さを最小限に抑えながら最大のバス帯域巾を維持することのできる公平で且つ比較的簡単な構造体を与えられる。

【 0 0 3 7 】

従って、同時挿入バッファを使用して広いバス帯域巾をサポートするローカルスイッチを実施することにより処理リソースを最適に使用するマルチプロセッサノード 1 0 が提供される。更に、A R B バス 1 3 において参照の順序がシリアル化されるので、マルチプロセッサ 1 0 のメモリのコヒレンス性を容易に維持する中央順序付けポイントが設けられる。ローカルスイッチに接続されるプロセッサモジュールの数を増加することにより処理能力を高める可能性が存在するので、図 2 の 4 プロセッサノードローカルスイッチ構成体は、待ち時間の短いそしてコストの安い高性能のシステムを提供する。

10

【 0 0 3 8 】

大型の対称的マルチプロセッサシステム

モノリシックマルチプロセッサノードに含むことのできるプロセッサの数は、2つのファクタにより制限される。第 1 に、ローカルスイッチを経て互いに接続できるプロセッサの数は、ローカルスイッチを構成するチップにおいて使用できるピンの数により制限される。第 2 に、単一のモノリシックスイッチによりサポートされるデータ帯域巾が制限される。従って、接続されるプロセッサの数をある点を越えて増加すると、何ら性能利得が得られないことになる。

20

本発明の 1 つの実施形態によれば、ハイアラキースイッチを経て複数のマルチプロセッサノードを相互接続することにより大型の対称的なマルチプロセッサシステムを形成することができる。例えば、ハイアラキースイッチを経て 8 個のマルチプロセッサノードを接続して、32 個のプロセッサモジュール、8 個の I O P デバイス及び 256 ギガバイトのメモリを含む対称的なマルチプロセッサ (S M P) システムが形成される。説明上、ここでは、少なくとも 2 つのマルチプロセッサノードを含む S M P を大型 S M P と称する。以下に詳細に述べるように、S M P ノードにローカルスイッチを用いて少数のプロセッサを接続し、そしてハイアラキースイッチを用いて多数のノードを大型の S M P へと接続することにより、拡張可能な高性能システムを実現することができる。

30

【 0 0 3 9 】

マルチプロセッサノードをハイアラキースイッチ式ノードへと接続するために、マルチプロセッサは、グローバルなポートインターフェイスを含むように拡張される。例えば、図 6 には、変更されたマルチプロセッサノード 1 0 0 が示されている。図 2 のマルチプロセッサノードと同様に、ローカルスイッチ 1 1 0 は、4 つのプロセッサモジュール、4 つのメモリモジュール及び I O P モジュールを接続する。図 2 及び 6 の同様の要素は、同じ参照番号を有する。マルチプロセッサノード 1 0 0 のローカルスイッチ 1 1 0 は、図 2 のポート 1 6 a - 1 6 i と同様に構成された 9 個のポート 1 1 6 a - 1 1 6 i を含む 1 0 ポートスイッチである。付加的なポート 1 1 6 j は、グローバルリンク 1 3 2 を経てグローバルポート 1 2 0 へ至る全二重のクロック供給データリンクを形成する。

40

【 0 0 4 0 】

グローバルポートは、マルチプロセッサノードをハイアラキースイッチに接続し、大型の S M P を実現する。例えば、図 7 A を参照すれば、本発明の 1 つの実施形態において、8 × 8 のハイアラキースイッチ 1 5 5 を経て互いに接続された 8 個のノード 1 0 0 a - 1 0 0 h を含む大型の S M P システム 1 5 0 が示されている。これらノード 1 0 0 a - 1 0 0 h の各々は、図 6 に示すノード 1 0 0 と実質的に同一である。

ノード 1 0 0 a - 1 0 0 h の各々は、全二重クロック供給データリンク 1 7 0 a - 1 7 0 h の各々によりハイアラキースイッチ 1 5 5 に接続される。1 つの実施形態において、データリンク 1 7 0 a - 1 7 0 h は、150 M H z のクロック速度で動作され、従って、スイッチ 1 5 5 との間でデータをやり取りするための 2 . 4 ギガバイト / 秒のデータ帯域巾をサポートする。これは、最大 38 . 4 ギガバイト / 秒の生の相互接続データ帯域巾、

50

及び 19.2 ギガバイト / 秒の二等分データ帯域巾をスイッチに与える。

【 0 0 4 1 】

大型の S M P システムは、マルチプロセッサノード 1 0 0 a - 1 0 0 h の各々が全システムメモリのアドレス可能な部分を含むか又は物理的メモリの分割部分を含むような分散型共用メモリシステムである。本発明の 1 つの実施形態では、全システムメモリに 2^{43} 個の物理的アドレス位置が存在する。S M P マルチプロセッサシステム 1 0 0 の 1 つの実施形態は、「大フォーマット」及び「小フォーマット」と称する 2 つのアドレスフォーマットをサポートする。大フォーマットは、各ノードのプロセッサが動作するところの 4 3 ビットの物理的アドレスを、マルチプロセッサシステムに使用するための 4 3 ビットの物理的アドレスに直接マップする。大フォーマットアドレスを使用すると、物理的メモリアドレスのビット < 3 8 : 3 6 > をノード識別番号として使用することができる。アドレスビット 3 8 : 3 6 は、メモリスペースアドレスのホームノードを直接デコードし、一方、アドレスビット 3 8 : 3 6 の逆数は、I / O スペースアドレスのホームノードをデコードし、ここで「ホーム」とは、メモリスペース又は I / O スペースに関連したメモリ及び I / O デバイスが存在するところの物理的マルチプロセッサノードを指す。

10

【 0 0 4 2 】

小フォーマットのアドレスモードは、マルチプロセッサシステムに 4 つ以下のノードが存在することを仮定するものである。小フォーマットは、各ノードのプロセッサが 3 6 ビットの物理的にアドレスされたシステムで動作できるようにする。小フォーマットにおいて、物理的アドレスのビット 3 4 : 3 3 は、データ又は I / O デバイスのホームノード番号を識別する。

20

しかしながら、たとえ C P U が 3 6 ビットの物理的アドレスを用いて動作しても、マルチプロセッサシステムは、データ位置を特定するのに 4 3 ビットの物理的アドレスを一貫して使用し、物理的アドレスのビット 3 7 : 3 6 がデータ又は I / O デバイスのホームノード番号を識別する。従って、C P U により発生された小フォーマットアドレスと、データライン 1 3 a - 1 3 h を経てハイアラークスイッチ 1 5 5 へ送信されるものとの間で何らの変換が実行される。

【 0 0 4 3 】

マルチプロセッサシステム 1 5 0 のここに示す構成は、3 2 個のプロセッサ間に広帯域巾のキャッシュコヒレントな共用メモリを与えることができる。本発明の 1 つの実施形態による大型 S M P の別の実施形態が図 7 B に示されており、ここでは、2 つのマルチプロセッサノード 1 0 0 a 及び 1 0 0 b がハイアラークスイッチを使用せずに互いに接続される。むしろ、2 つのマルチプロセッサノードは、それらのグローバルポート出力を互いに接続することにより直接接続される。

30

図 7 B の 2 ノード実施形態が使用されるか、図 7 A のマルチノード実施形態が使用されるかに拘わりなく、大きなアドレススペース及び処理能力をもつマルチプロセッサシステムが得られる。

【 0 0 4 4 】

両実施形態において、システムメモリアドレススペース及び I / O アドレススペースは、全てのノード 1 0 0 a - 1 0 0 h 間にセグメントで物理的に分配される。システムの各ノードは、メモリスペースの物理的アドレスの上位 3 ビットを使用してアクセスされるメインメモリの一部分を含む。従って、各メモリ又は I / O アドレスは、1 つのノードのみにおける 1 つの唯一のメモリ位置又は I / O デバイスへとマップされる。従って、上位 3 つのアドレスビットは、メモリ又は I / O アドレスがマップされるノードである「ホーム」ノードを識別するためのノード番号を与える。各マルチプロセッサノードは、それらのホームノード又は他のマルチプロセッサノードに記憶された共用メモリの部分をアクセスすることができる。ホームノードがプロセッサ自身のノードであるところの共用メモリブロックにプロセッサがアクセス（ロード又は記憶）するときには、参照は、「ローカル」メモリ参照と称される。ホームノードがプロセッサ自身のノード以外のノードであるようなブロックを参照する場合には、参照は、「リモート」又は「グローバル」メモリ参照と

40

50

称する。ローカルメモリアクセスの待ち時間は、リモートメモリアクセスの待ち時間と異なるので、SMPシステムは、非均一メモリアクセス(NUMA)アーキテクチャを有すると言える。更に、システムはコヒレントなキャッシュを備えているので、システムは、キャッシュコヒレントなNUMAアーキテクチャと呼ばれる。

【0045】

ここに示すキャッシュコヒレントなNUMAアーキテクチャは、高い性能と低い複雑さに寄与する多数の特徴を含む。設計上の1つの特徴は、メッセージ間の順序の固執及び利用である。メッセージがある順序特性に基づいてシステムに流れるよう保証することにより、オペレーションの待ち時間を著しく短縮することができる。例えば、記憶オペレーションは、記憶が完了したとみなされる前に無効メッセージがそれらの最終的な行先プロセッサに供給されることを必要とせず、むしろ、無効メッセージが行先プロセッサへと通じるある順序付けされた待ち行列に入れられるや否や記憶が完了したとみなされる。

更に、ある順序が維持されるよう保証することにより、設計上、確認又は完了メッセージの必要性が排除される。メッセージは、それらがある待ち行列に入れられた順序でそれらの行先に到達するように保証される。従って、メッセージがその行先に到達したときに確認を返送する必要性が排除される。これは、システムの帯域巾を改善する。

【0046】

更に、事象順序及びメッセージ順序は、「ホットポテト」オペレーションを行うのに使用される。ある待ち行列に順序を利用することにより、ディレクトリ又はDTAGコントローラのようなコントローラは、単一ビジットにおいて要求をリタイアすることができる。他の要求との競合により要求を否定的に確認しそして再トライする必要はない。「ホットポテト」オペレーションの結果として、公平さ及び欠乏の問題が解消される。

設計に使用される第2の特徴は、仮想チャンネルである。仮想チャンネルとは、メッセージを「チャンネル」へと分類する構成であって、チャンネルは物理的なリソースを共用する(従って、「仮想」である)が、各チャンネルは、他のものとは独立して流れ制御される。仮想チャンネルは、システムのメッセージ間で流れに依存しそしてリソースに依存するサイクルを排除することにより、キャッシュコヒレンスプロトコルにおける停滞を排除するのに使用される。これは、選択されたメッセージを否定的に確認しそしてそれに対応するコマンドを再トライすることにより停滞を検出しそして停滞状態を解消する機構を用いた公知のNUMAマルチプロセッサにおけるキャッシュコヒレンスプロトコルとは対照的である。

【0047】

チャンネルの使用について以下に簡単に説明するが、詳細な説明は後で行う。上述したように、メッセージは、「チャンネル」と称する論理的なデータ路を用いて大型SMP内をルート指定される。本発明の1つの実施形態には、以下のチャンネルが含まれる。即ち、要求側プロセッサから、トランザクションのアドレスに対応するホームノードのARBバスへトランザクションを搬送するためのQ0チャンネルと、ホームARBバスから1つ以上のプロセッサ及びIOPへトランザクションを搬送するためのQ1チャンネルと、所有者プロセッサから要求側プロセッサへデータ記入トランザクションを搬送するためのQ2チャンネルとである。変更されたデータを書き込むためにプロセッサからメモリへビクティム(Victim)トランザクションを搬送するためにQ0Vicチャンネルを設けることもできる。更に、Q0Vicチャンネルは、ビクティムトランザクションの背後に保持しなければならないQ0トランザクションを搬送するのに使用できる。最後に、プロセッサからIOPへI/Oスペーストランザクションを搬送するためにQIOチャンネルが設けられる。

【0048】

チャンネルは、以下に示すようなハイアラキーを構成する。

(最低) QIO > Q0Vic > Q0 > Q1 > Q2 (最高)

以下に述べるように、停滞を回避するために、いずれのチャンネルのメッセージも、下位チャンネルのメッセージによって決して阻止されてはならない。順序付け特性及び仮想チ

10

20

30

40

50

チャンネルを形成しそして使用する機構の設計及び実施に関する詳細は、後で述べる。
従って、図 7 A 及び 7 B に示すように、大型 S M P は、図 2 の S M P ノードを任意の数だけ互いに接続することにより形成することができる。図 7 A 及び 7 B に示すような大型 S M P システムのオペレーションは、以下に 3 つの部分について説明する。第 1 に、大型 S M P に含まれるハードウェア要素について説明する。次いで、S M P のプロセッサ間にコヒレントなデータ共用を与えるキャッシュコヒレンスプロトコルについて説明する。更に、ハイアラキースイッチの仮想チャンネルのために設けられたサポート機構を含む仮想チャンネルの実施及び使用について説明する。

【 0 0 4 9 】

大型 S M P のハードウェア要素

マルチプロセッサノードの各々には、チャンネルを用いてコヒレントなデータ共用を実施するための多数の要素が設けられる。図 6 に戻ると、これらの要素は、ディレクトリ 1 4 0 と、D T A G 2 0 と、I O P タグ 1 4 b と、グローバルポート 1 2 0 と、ディレクトリ 1 4 0 とを備えている。更に、シリアル化ポイントのハイアラキーは、キャッシュコヒレンスプロトコルを容易にするために参照の順序を維持できるようにする。これら要素の各々について、以下に詳細に述べる。

グローバルポート

グローバルポート 1 2 0 は、マルチプロセッサノード 1 0 0 を、ハイアラキースイッチリンク 1 7 0 を経て 1 つ以上の同様に構成されたマルチプロセッサノードに直接接続できるようにする。各ノード 1 0 0 は対称的なマルチプロセッサシステムとして動作するので、システムにより多くのノードが追加されるにつれて、使用可能なアドレススペース及び処理能力が増加される。

【 0 0 5 0 】

図 8 は、グローバルポート 1 2 0 の拡張ブロック図である。グローバルポートは、トランザクション追跡テーブル (T T T) 1 2 2 と、ビクティムキャッシュ 1 2 4 と、マルチプロセッサノードからハイアラキースイッチへ送られるパケットを記憶するためのパケット待ち行列 1 2 7、1 2 2、1 2 3 及び 1 2 5 と、ハイアラキースイッチから受け取られるパケットを記憶するためのパケット待ち行列 1 2 1 とを備えている。グローバルポート 1 2 0 は、A R B バス 1 3 0 と、ローカルスイッチの 2 つの専用ポート即ち G P リンク入力 1 3 2 b 及び G P リンク出力 1 3 2 a とを経たノードの他のロジック (特に Q S A チップ) と通信する。

T T T は、マルチプロセッサノードにおいて保留中のトランザクション、即ちノードからグローバルポートを経て発生されて、他のマルチプロセッサノード又はハイアラキースイッチからの応答を待機しているトランザクションを追跡する。グローバルポートにコマンドが送られるたびに、T T T にエントリが形成される。対応する応答がノードに受け取られたときに、T T T エントリがクリアされる。T T T は、2 つの部分、即ち Q 0 T T T 及び Q 1 T T T で構成され、Q 0 及び Q 1 は、上記のように Q 0 及び Q 1 チャンネルを進むパケットを指す。エントリーが T T T にいかに割り当てられるか及びそれがいつリタイアされるかについては、以下に詳細に述べる。

【 0 0 5 1 】

又、グローバルポート 1 2 0 は、ビクティムキャッシュ 1 2 4 を含む。ビクティムキャッシュ 1 2 4 は、マルチプロセッサノードの各プロセッサから受け取られて別のマルチプロセッサノードのメモリに向けられるビクティム化データを記憶する。ビクティム化データとは、プロセッサのキャッシュ位置に記憶されてそのプロセッサにより変更されたデータである。変更データを記憶するキャッシュ位置に記憶する必要のある新たなデータがプロセッサに受け取られると、変更データは、ビクティム化されると言われ、ビクティムデータと称される。

ビクティムキャッシュ 1 2 4 は、プロセッサからリモートマルチプロセッサノードのメモリへ向けられたビクティムデータからのビクティムデータの一時的な記憶装置である。グローバルポートを経て別のノードへビクティムデータを送信するための機会があるときに

10

20

30

40

50

は、マルチプレクサ 167 は、ピクティムキャッシュ 124 からバス 170 の出力部分にデータを供給するように切り換えられる。グローバルポートにピクティムキャッシュを設けることにより、個々のプロセッサがグローバルシステムのメモリ書き込み待ち時間を待機せずに、プロセッサが各々のピクティムデータバッファを空にすることができる。むしろ、ピクティム書き込みは、使用できるデータサイクルがあるときに書き込みが実行されるようにグローバルポートにより制御される。ピクティムキャッシュからデータを解放する適切さに関連した幾つかの制御の問題があるが、これらは以下に説明する。

【0052】

D T A G 及び I O P タグ

D T A G 及び I O P タグは、小型の S M P システムにも含まれるが、これについては以下に詳細に述べる。D T A G 20 は、マルチプロセッサノードのプロセッサのキャッシュに記憶されたデータブロック各々に対する状態情報を記憶する。同様に、I O タグ 14 a は、I O P に記憶された各データブロックに対する状態情報を記憶する。ディレクトリは、どのマルチプロセッサノードがデータのコピーを記憶するかを識別するおおよその情報を与えるが、D T A G 及び I O タグは、マルチプロセッサノード内のどのプロセッサがデータのコピーを記憶するかに関する正確な指示を与えるのに使用される。それ故、D T A G 及び I O タグは、参照情報がマルチプロセッサノードに到達したときに、そのノードのどのプロセッサがターゲットとなるべきかを決定するのに使用される。

【0053】

図 6 に示すように、D T A G 20 及び I O P タグ 14 b は、Q S A チップ 18 に接続されたメモリ領域を参照するアドレスを監視するために A R B バス 130 に接続される。D T A G は、4 つのプロセッサ 12 a - 12 d に対応する 4 つのセグメントに分割される。各プロセッサは、メモリ 13 からのデータのサブセットを一時的に記憶するためのキャッシュ（図示せず）を備えている。各プロセッサのキャッシュに記憶されたメモリのブロックの上位アドレスビット（タグ）を記憶するためのタグ記憶装置が各キャッシュに関連される。D T A G 20 の各セグメントは、関連プロセッサのキャッシュタグの状態を指示するデータを維持する。処理ユニットの外部の D T A G 20 にタグのコピーを記憶することにより、システムは、A R B バスを経て受け取ったコマンドをフィルタし、そしてプロセッサのキャッシュのデータに関連した調査（読み取り）及び無効化コマンドのみを各プロセッサに供給することができる。I O P タグ 14 a は、I O P キャッシュ 14 c に記憶されたデータブロック各々の上位アドレスビットを記憶する。I O P タグ記憶装置は、プロセッサ 12 a - 12 d の各々に維持されたタグ記憶装置と同様である。

【0054】

D T A G 20 及び I O P タグ 14 a の各エントリは、多数の状態ビットを含む。D T A G 状態ビットは、次の 4 つの状態、即ち I n v a l i d（無効）、C l e a n（クリーン）、D i r t y _ _ N o t _ _ P r o b e d、及び D i r t y _ _ P r o b e d のうちの 1 つを指示する。I O P タグのエントリの状態ビットは、次の 2 つの状態、即ち V a l i d（有効）及び D i r t y（ダーティ）のうちの 1 つを指示する。「有効」ビットは、関連キャッシュの対応エントリに記憶されたデータが、メモリに記憶されたデータと一致することを指示する。「ダーティ」ビットは、関連キャッシュの対応エントリに記憶されたデータが関連プロセッサによって変更されそしてメモリに記憶されたデータに一致しないことを指示する。

【0055】

D T A G 20 及び I O P タグ 14 b は、マイクロプロセッサノード 100 の A R B バスにコマンドが現れるたびにアクセスされる。「無効」の状態がプロセッサ 1 の D T A G アクセスに応答して返送される場合には、ノードのプロセッサ 1 は、メモリアドレスに関連したデータの有効コピーを記憶しない。「有効」の状態が I O P タグ 14 a へのアクセスから返送される場合には、I O P キャッシュ 14 c がデータの有効コピーを記憶する。「クリーン」状態がプロセッサ 1 に対する D T A G アクセスに応答して返送される場合には、これは、プロセッサ 1 がメモリアドレスに対応するデータの無変更コピーを有するが、そ

10

20

30

40

50

のデータを読み取るための他のプロセッサによる試みがなされていないことを指示する。
Dirty__Not__Probedの状態がDTAGに 응답して返送される場合には、これは、プロセッサ1がメモリアドレスに対応するデータの変更コピーを有し、そしてプロセッサが最後にデータを変更して以来、少なくとも1つのプロセッサがデータを読み取る試みをしていることを指示する。

【0056】

ディレクトリオペレーション

一般に、ディレクトリは、関連マルチプロセッサノード（ホームノード）におけるメモリの各ブロックの所有権情報を与えるのに使用され、メモリのブロックは、一般に、メモリとSMPシステムのプロセッサとの間に転送される最小量のデータである。例えば、本発明の1つの実施形態において、ブロックは、パケットのサイズと同様であり、即ち512ビット（64バイト）のデータである。更に、ディレクトリは、どのマルチプロセッサノードがメモリデータのブロックのコピーを記憶するかを指示する。従って、読み取り型のコマンドの場合に、ディレクトリは、データの最新バージョンの位置を識別する。ピクティム型のコマンドの場合には、データの変更ブロックがメモリに書き戻される場合に、ディレクトリは、データの変更ブロックが現在のものであってメモリに書き込まねばならないかどうか決定するために検討される。それ故、ディレクトリは、参照情報がリモートマルチプロセッサノードのプロセッサにより発生されたものであるかローカルマルチプロセッサノードのプロセッサにより発生されたものであるかに拘わりなく、関連するマルチプロセッサノードのメモリブロックへの参照に対する第1アクセスポイントである。

【0057】

ディレクトリは、対応するノード100においてメモリ13の各64バイトのデータブロック（以下、キャッシュラインとも称する）に対して1つの14ビットエントリを記憶する。メモリ13と同様に、ディレクトリは、メモリアドレスがノードNに存在する場合に、対応するディレクトリエントリもノードNに存在するように、システムのノードにわたって物理的に分配される。

図9を参照すれば、ディレクトリエントリ140aの1つの実施形態は、所有者IDフィールド142及びノード存在フィールド144を含むように示されている。所有者IDフィールドは、各64バイトブロックに対する6ビットの所有者情報を含む。所有者IDは、ブロックの現在所有者を特定し、現在所有者は、システムにおける32個のプロセッサの1つ、又はシステムにおける8個のI/Oプロセッサの1つ、又はメモリのいずれかである。8ビットのノード存在情報は、システムの8個のノードのどれがキャッシュラインの現在バージョンを獲得したか指示する。ノード存在ビットは、同じノードにおける4つのプロセッサの累積状態を1ビットで表わすおおよそのベクトルである。共用データの場合には、2つ以上のノードが、情報を記憶する少なくとも1つのプロセッサを有する場合に、2つ以上のノード存在ビットがセットされる。

【0058】

時々、状態情報のある断片がDTAG又はディレクトリから得られる。このような場合、DTAGからの状態情報を使用するのが好ましい。というのは、これは非常に高速で検索されるからである。例えば、メモリアドレスの所有者プロセッサがそのアドレスに対しホームノードに配置される場合には、所有者IDを供給するのにDTAGが使用される。性能上の理由でDTAGによりサービスされない情報又は参照については、ディレクトリ140は、全てのコヒレンス性判断の焦点であり、従って、多数の機能を実行する。ディレクトリは、メモリデータブロックの所有者を識別する。所有者は、プロセッサ又はメモリのいずれかである。ディレクトリからの所有者情報は、データブロックの最新バージョンのソースを決定するために読み取り型コマンド（例えば、読み取り、読み取り-変更）により使用される。又、所有者情報は、以下に詳細に述べるようにピクティム化データをメモリに書き戻さねばならないかどうか決定するのににも使用される。

【0059】

全ての読み取り型コマンドに対して、データの所有者を識別するのに加えて、ディレクト

10

20

30

40

50

リは、プロセッサからの「クリーン - ダーティ (Clean-to-Dirty)」及び「シェアド - ダーティ (Shared-to-Dirty)」コマンドを分析するのににも使用される。「クリーン - ダーティ」コマンドは、プロセッサがそのキャッシュにおいて現在「クリーン」状態にあるキャッシュラインを変更するよう希望するときにプロセッサにより発生される。「シェアド - ダーティ」コマンドは、「ダーティ - シェアド」状態にあるキャッシュラインを変更するよう希望するときに発生される。これらのコマンドは、ホーム A R B バスに送られ、そこで、ディレクトリは、プロセッサがキャッシュラインの最新バージョンを有するかどうか決定する。もしそうであれば、コマンドは成功となり、プロセッサは、キャッシュラインを変更することが許される。さもなくば、コマンドは失敗となり、プロセッサは、最初に、キャッシュラインの最新バージョンを獲得しなければならない。これらの記憶型オペレーションは、ディレクトリのノード存在情報を使用して、成功又は失敗を決定する。

10

【 0 0 6 0 】

上記のように、ディレクトリの存在ビットは、記憶型コマンドが発生されたときに各データブロックのコピーでマルチプロセッサノードを識別する。記憶コマンドは、キャッシュラインの内容が更新されようとしていることを指示する。関連するディレクトリエントリの存在ビット 1 4 4 を検討することにより、記憶コマンドがディレクトリ 1 4 0 に受け取られたときに、存在ビットを有するノードを用いて、これらのマルチプロセッサノードをそのノードにおけるキャッシュラインのコピーで識別し、従って、各ノードにおけるキャッシュラインを無効化できるようにする。

従って、ディレクトリ及び D T A G は、ローカルマルチプロセッサのメモリにおける各データブロック及びローカルプロセッサのキャッシュに記憶された各データブロックに対する状態情報を与えるように協働する。ホームノードのディレクトリは、キャッシュブロックのコピーの状態に関するおおよその情報を供給する。次いで、無効化コマンドがディレクトリにより識別されたノードへと進み、そこで、D T A G がアクセスされて、コピー情報を更に改善する。従って、これらノードにおける D T A G は、各ノードのどのプロセッサがそれらのキャッシュにラインのコピーを記憶するか指示する。

20

【 0 0 6 1 】

T T T :

T T T は、マルチプロセッサノードからの保留中のトランザクション、即ち別のマルチプロセッサノード又はハイアラキースイッチからの応答を待機している参照を追跡するのに使用される。保留中トランザクションに関する情報は、関連メモリアドレスへのその後のコマンドを処理する際にキャッシュコヒレンスプロトコルにより使用される。

30

図 1 0 を参照すれば、T T T 1 2 2 の 1 つの実施形態は、アドレスフィールド 1 5 2 と、コマンドフィールド 1 5 4 と、コマンド I D フィールド 1 5 6 と、ビット 1 5 8 a - 1 5 8 c を含む多数の状態ビット 1 5 8 とを含むように示されている。アドレスフィールド 1 5 2 は、現在進行中であるトランザクションに対するキャッシュラインのアドレスを記憶し、一方、コマンドフィールドは、現在進行中であるトランザクションに対するキャッシュラインに関連したコマンドを記憶する。コマンド I D フィールド 1 5 6 は、コマンドフィールドに記憶されたコマンドを開始したプロセッサのプロセッサ番号を記憶する。状態ビット 1 5 8 は、コマンドが進行中であるときにコマンドの状態を表わす。或いは又、状態ビット 1 5 8 は、進行中であるコマンドの種々の特性をあらわすように使用されてもよい。

40

【 0 0 6 2 】

例えば、「記入」状態ビット 1 5 8 a は、読み取り型コマンドに回答して「記入」データ応答が受け取られたときに更新される。「シャドー」状態ビット 1 5 8 b は、グローバルポートを経て発生されたコマンドが「シャドー」型コマンド（以下に詳細に述べる）である場合にセットされる。A C K 状態ビット 1 5 8 c は、確認型応答を期待しているメッセージが応答を受信した場合にセットされる。応答が到着した場合に、このビットはクリアされる。T T T に含むことのできる全ての状態ビットが示されているのではないことに注意されたい。むしろ、以下の説明に関連のある状態ビットが含まれている。更に、メモリ

50

のコヒレンス性を維持するために必要と考えられれば、他の状態ビットを設けてもよく、従って、本発明は、ＴＴＴにおける特定のビット指定に限定されるものではないことが明らかであろう。

【００６３】

従って、ディレクトリ、ＤＴＡＧ、ＩＯＰタグ及びＴＴＴの各々は、ＳＭＰシステムにおけるキャッシュラインのコヒレンス性（以下、キャッシュコヒレンス性と称する）を維持するのに使用される。これら要素の各々は、ハイアラークスイッチ１５５に接続されたマルチプロセッサノード間にコヒレント通信を与えるためにグローバルポートとインターフェイスする。

【００６４】

シリアル化ポイント：

上記要素に加えて、各マルチプロセッサノードにシリアル化ポイントを設定することによりデータ共用コヒレンス性が維持される。本発明の１つの実施形態において、各マルチプロセッサノードにおけるシリアル化ポイントは、ＡＲＢバス１３０である。全てのＱ０参照は、ローカルプロセッサにより発生されたものであるかリモートプロセッサにより発生されたものであるかに拘わりなく、ＱＳＡによりＡＲＢバス１３０を経てディレクトリ１４０及びＤＴＡＧ２０へ供給される。参照がディレクトリ及び／又はＤＴＡＧをアクセスすると、それにより得られるＱ１チャンネルコマンドが厳密な順序でＡＲＢバスに出力され、ここで、順序は参照のシリアル化順序である。マルチプロセッサノードの各々にシリアル化ポイントを設定することにより、ＳＭＰにおいて実施されるデータ共用コヒレンスプロトコルが相当に簡単化される。

【００６５】

マルチプロセッサノードの各々にシリアル化ポイントを設定するのに加えて、ハイアラークスイッチ１５５は、ＳＭＰシステムに第２のシリアル化ポイントを与える。以下に詳細に述べるように、ハイアラークスイッチは、第１のシリアル化ポイントに導入されたコヒレンス性が大型のＳＭＰシステムに維持されるよう確保するある順序付けルールに適合する。 グローバルポート／ハイアラークスイッチインターフェイス：

図１１は、８個の入力ポート１５５ｉ０－１５５ｉ７及び８個の出力ポート１５５ｏ０－１５５ｏ７を含むハイアラークスイッチ１５５のブロック図である。ハイアラークスイッチ１５５の入力ポート１５５ｉ０－１５５ｉ７は、接続されたマルチプロセッサノード各々のグローバルポートからパケットを受け取る。ハイアラークスイッチの出力ポート１５５ｏ０－１５５ｏ７は、接続されたマルチプロセッサノード各々のグローバルポートへパケットを供給する。

【００６６】

本発明の１つの実施形態において、受信したパケットをバッファするためのバッファ１６０ａ－１６０ｈが各入力ポートに関連される。図１１の実施形態は、各入力に１つのバッファを示しているが、いかなる数の入力ポート間にバッファが共用されてもよい。各パケットは、５つのチャンネルのいずれか１つと関連される。本発明の１つの実施形態では、以下に述べるように、各入力バッファ１６０ａ－１６０ｈの部分が、あるチャンネルのパケットを専用に記憶するようにされる。従って、グローバルポートからハイアラークスイッチ１５５への流れ制御は、チャンネルベースで実行される。チャンネルベースでスイッチへのデータの流れを制御しそして入力バッファの部分を選択されたチャンネルに専用とすることにより、スイッチは、ＳＭＰシステムにおけるマルチプロセッサノード間で停滞のない通信を行う。

【００６７】

停滞のない通信を与えるのに加えて、ハイアラークスイッチ１５５は、更に、メモリのコヒレンス性を確保するためにＳＭＰシステムの順序付け制約をサポートするように設計される。順序付け制約は、スイッチ１５５から関連マルチプロセッサノードのグローバルポートへ送出されるパケットの順序を制御することにより課せられる。いずれかの入力バッファ１６０ａ－１６０ｈからのパケットは、マルチプレクサ１８２ａ－１８２ｈを経て

10

20

30

40

50

いずれかの出力ポートへ送られる。更に、以下に述べるように、スイッチ 155 は、パケットをマルチキャストリングすることができる。従って、1つの入力バッファからのパケットは、いかなる数の出力ポートに送ることもできる。グローバル出力ポートに順序を強制することにより、マルチプロセッサノード各々に得られるシリアル化順序を維持して、完全にコヒレントなデータ共用機構を SMP システムに形成することができる。

【0068】

ハイアラキースイッチにおける停滞の回避

上述したように、図 7 A の 8 個のノードの各々は、ハイアラキースイッチにデータを供給し、全てのノードがデータを同時に供給することもある。パケットは、異なる仮想チャンネルに供給される多数の異なるチャンネル形式 (Q0、Q0Vic、Q1、Q2 及び QIO) に分割され、ここで、仮想チャンネルとは、本質的に、他のチャンネルとの共通の相互接続部を共用するがその相互接続部のいずれかの端において独立してバッファされる特定形式のパケットに専用のデータ経路である。各ノードのグローバルポートとハイアラキースイッチとの間には 1 つのデータ経路しかないので、異なる仮想チャンネルからの全てのパケットは、1 つのデータ経路を使用してハイアラキースイッチに書き込まれる。

【0069】

8 個のノード 100a - 100h の各々は、ハイアラキースイッチへデータを送信することができるので、全てのメッセージがスイッチにより受信されて、スイッチから適当な順序で供給されるよう適切に確保するために、ある形式の制御が必要となる。更に、本発明の 1 つの目的は、対称的なマルチプロセッサシステムに停滞 (デッドロック) が生じないよう保証するために上位順序のパケット形式が下位順序のパケット形式により阻止されないよう確保することである。本発明の 1 つの実施形態では、最高順序から最低順序までのパケットの順序は、Q2、Q1、Q0、Q0Vic 及び QIO である。

本発明の 1 つの特徴によれば、スイッチの入力ポートに到着するパケットの流れ制御を行うための機構であって、上記の停滞回避ルールが常に満足されるよう確保する機構が提供される。更に、スイッチにおいて使用できるバッファは最適に利用されねばならず、そして最大の帯域巾が維持されねばならない。

【0070】

本発明の 1 つの実施形態によれば、ハイアラキースイッチへのデータの書き込みを制御するための制御装置は、パケットの各形式に対し、ハイアラキースイッチのバッファに専用スロットを設けることにより実施される。又、バッファは、任意の形式のパケットを記憶するのに使用できる多数の一般的なスロットも含んでいる。ハイアラキースイッチに専用のバッファスロットを設けることにより、上位順序のパケット形式が常にスイッチを通る経路を使用できるよう保証することによって停滞を回避することができる。更に、使用できる一般的スロット及び専用スロットの数を監視し、そしてバッファに記憶されるパケットの異なる形式の数を監視することにより、ハイアラキースイッチのバッファが容量に達したときにノードがバッファに書き込みするのを防止するような簡単な流れ制御機構を実施することができる。

【0071】

図 12 A には、多数のソースノードによる共通の行先バッファへの書き込みを制御するのに使用するための制御ロジックの一例が示されている。図 12 A のブロック図には、2 つの異なるノードのグローバルポート 120a 及び 120b が一例として示されている。

図 12 A において、ノード 100a 及び 100b のグローバルポート各々 120a 及び 120b の部分は、ハイアラキースイッチ 155 へ転送するために Q0 / Q0Vic、Q1、Q2 及び一般形式のパケット (Q0、Q0Vic、Q1、Q2 又は QIO パケットのいずれか) を各々記憶するためのエントリ 135a - 135b を含むバッファ 135 を備えて詳細に示されている。バッファ 135 にはマルチプレクサ 167a が接続され、GP アービター 134 からの選択信号を使用してリンクを経てハイアラキースイッチへ送るためにパケット形式の 1 つを選択する。

【 0 0 7 2 】

更に、各グローバルポートは、専用のカウントレジスタ 1 3 6 を備えている。この専用のカウントレジスタは、パケットの各 Q 0 / Q 0 V i c、Q 1 及び Q 2 チャンネル形式に対して、ハイアラキー 스위ッチ 1 5 5 において現在保留となっているチャンネル形式のパケットの数のカウントを記憶する。このカウントは、各チャンネル形式のパケットがハイアラキー 스위ッチへ転送されるときに増加され、そしてパケットがハイアラキー 스위ッチから転送されるときに減少される。

本発明の 1 つの実施形態において、ハイアラキー 스위ッチ 1 5 5 は、8 個の入力ソースの各々に 1 つのバッファを備えている。図 1 2 A には、2 つのグローバルポート 1 2 0 a 及び 1 2 0 b に対応する 2 つのバッファ 1 6 0 a 及び 1 6 0 b のみが示されている。本発明の 1 つの実施形態では、バッファ 1 6 0 a 及び 1 6 0 b の各々に少なくとも (m - 1) × n 個の専用スロットがあり、但し、m は、バッファに専用エントリを有する仮想チャンネル形式の数に対応し、そして n は、バッファを共用するノードの数に対応する。図 1 2 A の実施形態において、各バッファは、8 個のエントリを有する。エントリのうちの 5 つは、一般的エントリであり、グローバルポート 1 3 5 から送られたパケットの形式を記憶することができる。残りの 3 つのエントリの各々は、特定形式のパケットを専用記憶し、即ち 1 つのエントリは、Q 0 / Q 0 V i c パケットを専用記憶し、1 つのエントリは、Q 1 形式パケットを専用記憶し、そして 1 つのエントリは、Q 2 形式パケットを専用記憶する。

【 0 0 7 3 】

専用エントリがバッファ 1 6 0 a 及び 1 6 0 b の固定位置に存在するものとして示されているが、実際には、バッファのいずれの位置も専用のバッファ位置であり、即ちエントリの位置に拘わりなく、パケットの各特定形式ごとにバッファには常に 1 つの専用エントリがある。

ハイアラキー 스위ッチは、更に、各バッファ 1 6 0 a 及び 1 6 0 b に対し、専用カウンタ 1 6 2 a 及び 1 6 2 b と、フラグレジスタ 1 6 3 a 及び 1 6 3 b とを含む。図 1 2 A の実施形態において、専用カウンタ 1 6 2 a は、4 つのエントリを有し、その 3 つは、バッファ 1 6 0 a に現在記憶されている Q 0 / Q 0 V i c、Q 1 及び Q 2 パケットの数を記憶するためのもので、そして 1 つは、バッファに使用される一般的エントリの数のカウントを記憶するためのものである。フラグレジスタは、3 つのビットを含み、各ビットは、パケットの Q 0 / Q 0 V i c、Q 1 及び Q 2 形式の 1 つに対応し、そして関連する専用カウンタがゼロであるかどうか（即ち、その形式のパケットの専用エントリが使用されたかどうか）を指示する。従って、フラグレジスタの値は、その形式の少なくとも 1 つのパケットがバッファに記憶されたことを指示する 1 であるか、又はその形式のパケットがバッファに記憶されないことを指示する 0 である。

【 0 0 7 4 】

更に、ハイアラキー 스위ッチ 1 5 5 は、各バッファ 1 6 0 a 及び 1 6 0 b に対し、トランシットカウンタ 1 6 4 a 及び 1 6 4 b を各々含む。トランシットカウンタは、各ソースに対して、所与のデータサイクル中にトランシット状態であるいずれかの形式の保留中パケットの数を維持する。

所与のデータサイクル中にトランシット状態にあるパケットの数は、ハイアラキー 스위ッチとグローバルポートとの間の流れ制御待ち時間に直接関係している。流れ制御信号は、ハイアラキー 스위ッチからグローバルポートへ送られて、ハイアラキー 스위ッチへのデータの送信を停止するようにグローバルポートに通知する。流れ制御待ち時間 (L) は、ハイアラキー 스위ッチによる流れ制御信号のアサートと、グローバルポートによるデータ送信の停止との間に生じるデータ転送サイクルの数として測定される。

【 0 0 7 5 】

又、ハイアラキー 스위ッチは、各バッファ 1 6 8 a 及び 1 6 8 b の書き込みを制御するための書き込み制御ロジック 1 6 6 a 及び 1 6 6 b も備えている。この書き込み制御ロジックは、ライン 1 6 8 a に「流れ制御」信号をそしてライン 1 6 8 b に「確認 (A C K)

10

20

30

40

50

「信号 < 3 : 0 > をアサートすることにより関連バッファへのデータの流れを制御する。「流れ制御」及び A C K 信号は、各データ転送サイクルに送信される。上記のように、「流れ制御」信号は、接続されたグローバルポートによるパケットデータの送信を停止するのに使用される。ライン 1 6 8 b の A C K 信号 < 3 : 0 > は、パケットの専用形式の各々に対して 1 ビットを含み、そして接続されたグローバルポートに、その形式のパケットが関連バッファから解放されたことを通知するのに使用される。従って、A C K 信号は、グローバルカウンタにより、専用カウンタ 1 3 6 の値を増加するのに使用される。

【 0 0 7 6 】

書き込み制御ロジックは、バッファの使用可能な全ての一般的エントリが、ハイアラキー 스위ッチへのトランシット状態にある考えられる全てのパケットを受け入れるのに充分でないと決定されたときに、流れ制御をアサートする。使用可能な一般的スロットの数は、次の式 I により決定することができる。

式 I :

$$\text{Generic_count} = (\text{バッファサイズ}) - (\text{バッファに使用される一般的エントリの数}) - (\text{非アサートフラグの数})$$

使用可能な一般的エントリの数が決定されると、式 I I が真である場合に、流れ制御信号がアサートされる。

式 I I :

$$\text{Generic_Count} = (\text{トランシットカウンタ}) * (\text{バッファを使用するノードの数})$$

従って、書き込み制御ロジック 1 6 6 は、使用中の一般的及び専用のスロットの数、トランシットカウンタ及び全バッファサイズを監視し、「流れ制御」信号をいつアサートすべきかを決定する。

【 0 0 7 7 】

「流れ制御」信号をアサートしても、ソースノードのグローバルポートによる全ての送信は停止されない。グローバルポートは、専用パケット形式に対応する専用スロットがハイアラキー 스위ッチのバッファに使用できる場合に、専用パケットデータをハイアラキー 스위ッチに常に転送する。従って、専用カウンタにおけるいずれかの専用カウンタの値がゼロに等しい場合には、グローバルポートは、常に、対応する専用パケット形式のパケットデータを転送することができる。従って、バッファに専用エントリを設けることにより、ハイアラキー 스위ッチを通る 1 つの形式のパケットの進行が、そのスイッチを通る他のパケットの進行によって左右されないように効果的に保証される。

バッファ 1 6 0 a 及び 1 6 0 b に専用及び一般的なスロットを使用することにより、各パケット形式ごとに最小数のスロットを指定するだけでよい。トランシット状態のパケットの数を追跡することにより、流れ制御を微細な粒度で行うことができる。バッファの利用性及びバスの帯域巾の両方が最大にされる。例えば、X の一般的スロットしか使用できないときには、流れ制御が 1 サイクル放棄され、そして次のサイクルに再アサートされる。その結果、X までのメッセージを時間周期内に受け取ることができる。

【 0 0 7 8 】

図 1 2 B は、ハイアラキー 스위ッチへデータを供給するためにグローバルポートにより使用される方法を示すフローチャートである。このプロセスは、1 つの形式のパケットについて説明するが、他の形式のパケットにも容易に拡張できる。ステップ 1 6 9 では、ハイアラキー 스위ッチ 1 5 5 へ供給すべきパケットがバッファ 1 3 5 a - 1 3 5 d の 1 つに存在するかどうか G S アービター 1 3 4 において決定される。パケットがある場合には、ステップ 1 7 1 において、「流れ制御」信号の状態がアービター 1 3 4 により評価される。「流れ制御」信号がアサートされる場合には、ステップ 1 7 2 において、ハイアラキー 스위ッチにより送られるべきパケットの特定形式に対する専用カウンタを検査して、それがゼロに等しいかどうか決定される。専用カウンタがゼロに等しくない場合には、その形式のパケットに対するバッファ内の専用エントリが既に使用中であり、プロセスはステップ 1 7 0 へ戻り、そのパケット形式の専用カウンタがゼロに等しくなるまで又は流れ制御信号がデアサートされるまで、ステップ 1 6 9、1 7 1 及び 1 7 2 間をループする。

ステップ 172 において専用カウン트가ゼロに等しいと決定された場合には、ステップ 173 において、GP アービター 134 は、適当な選択信号をマルチプレクサ 167 へアサートし、所望のパケットをハイアラキースイッチ 155 へ送信する。ステップ 174 において、パケットの選択された形式に対応する専用カウン트가グローバルポートの専用カウントレジスタ 134 及びハイアラキースイッチ 155 の専用カウントレジスタ 162a において増加され、そしてフラグレジスタ 163a の関連フラグがアサートされる。

【0079】

上記のように、フラグレジスタ 163a は、一般的カウンタ及びトランシットカウンタと共に使用されて、次のデータサイクルに対する「流れ制御」信号の状態を決定する。図 13 には、ハイアラキースイッチによる「流れ制御」信号のアサートを制御するためのプロセスの一実施形態が示されている。ステップ 175 において、フラグレジスタ 163a が検査されて、ゼロに等しい専用カウンタエントリの数が計数される。上記のように、ゼロの数は、「流れ制御」がアサートされた後であってもバッファに接続された各ノードにより送られる潜在的な専用パケットの数を指示する。従って、図 11 の例においていずれのノードについても専用スロットが全く使用されない場合には、フラグレジスタの全てのエントリがゼロに等しくなり、従って、専用パケットのために指定されねばならないバッファ位置が 3 つあることを指示する。

【0080】

フラグレジスタ 163a の値が検査された後、ステップ 176 において、使用可能な全一般的スロットが上記式 I を用いて決定される。次いで、ステップ 177 において、各ノードのトランシットカウン트가決定される。上述したように、トランシットカウンタは、所与のデータサイクル中にグローバルポートとハイアラキースイッチとの間でトランシット状態にあるメッセージの数を示す。最悪の場合のトランシットカウンタは、流れ制御の待ち時間 L にバッファ N を使用するノードの数を乗じたものに等しい。しかしながら、本発明の 1 つの実施形態によれば、トランシットカウンタの決定には、「流れ制御」信号が手前のサイクル中にアサートされたかどうかを考慮される。上記のように、「流れ制御」信号が手前のサイクルにアサートされた場合には、グローバルポートとハイアラキースイッチとの間でトランシット状態となるパケットはない。例えば、手前の J 個の周期中に「流れ制御」がゼロであった場合には、J x N 個までのメッセージがトランシット状態となる。しかしながら、J - 1 個の手前のデータサイクル中に「流れ制御」信号がゼロであった場合には、(J - 1) x N 個のメッセージのみがトランシット状態となる。

【0081】

従って、本発明の 1 つの実施形態では、ソース（グローバルポート）と行先（ハイアラキースイッチ）との間の全待ち時間を検査すると共に、手前のデータサイクルにおけるソースと行先との間の相互作用を検査することにより、トランシット状態のパケットの数がインテリジェントに決定される。各ノードに対するトランシットカウン트가決定された後に、ステップ 178 において、上記の式 II を用いて保留中の専用パケット及びトランシット状態のパケットを受け入れるに十分な使用可能な一般的エントリがバッファにあるかどうかの判断がなされる。使用可能な一般的パケットの全数が、トランシット状態にあるパケットの数にバッファを共用するノードの数を乗じた値より少ない場合には、ステップ 178 において、「流れ制御」信号がグローバルポート 120a にアサートされ、ハイアラキースイッチ 155 へのデータの供給が阻止される。しかしながら、全カウンタが、潜在的に受け取られるパケットの数をバッファ 160a で受け入れできることを指示する場合には、「流れ制御」信号がアサートされず、プロセスは、次のデータサイクルのためにステップ 175 へ復帰する。

【0082】

従って、トランシット状態にあるメッセージの数と、流れ制御信号がアサートされた手前のサイクルの数とを追跡することにより、流れ制御は、グローバルポートをハイアラキースイッチに接続するデータリンクの利用性が最大となるよう確保するように微同調される。図 11 ないし 13 に示すバッファ書き込み制御ロジック及び方法は、ノードからハ

10

20

30

40

50

イアラークースイッチへのデータの送信に関して説明したが、本発明は、このような構成に限定されるものではないことに注意されたい。むしろ、本発明の1つの実施形態は、共通の受信器に信号供給する多数のソースがありそして停滞を回避する必要があるいかなる環境にも使用できる。

【0083】

チャンネル順序付け制約をサポートするハイアラークースイッチの機構：

ハイアラークースイッチからのデータの読み取りは、本質的に、パケットの順序と、パケット間のデータ依存性との両方が維持されるように入力バッファから多数の出力ソースへデータを供給することを含む。上述したように、パケットは種々のチャンネルに供給される。異なるチャンネルにおいてパケットに関連するのは、ある順序付け制約即ち依存性である。本発明の1つの実施形態では、1つの順序付け制約は、Q1チャンネルの全てのパケットが順序正しく維持されることである。別のパケット順序付け依存性は、優先順位の高いチャンネルを進行するパケットが、優先順位の低いチャンネルを進行するパケットによって阻止されてはならないことであり、チャンネルの優先順位は、最も高いものから最も低いものへ、Q2、Q1、Q0、Q0vic及びQIOである。順序の維持は、以下に述べる種々の技術を用いてSMP全体にわたり達成される。ハイアラークースイッチにおいては、データ依存性及びQ1チャンネル順序付けを満足するよう確保するために3つの基本的なガイドラインに従う。これらのガイドラインは、次の通りである。

【0084】

ガイドライン1：所与のハイアラークースイッチ入力ポートに受け取られた多数のQ1パケットが共通の出力ポートをターゲットとする場合には、Q1パケットは、それらが入力ポートに現れたのと同じ順序で出力ポートに現れる。

ガイドライン2：ハイアラークースイッチにおいて多数の入力ポートからのQ1パケットが共通の出力ポートへマルチキャストされるときには、Q1パケットは、それらがターゲットとする全ての出力ポートに同じ順序で現れる。

ガイドライン3：ハイアラークースイッチの多数の入力ポートからのQ1パケットの順序付けリストが多数の出力ポートをターゲットとするときには、Q1パケットは、全ての到来するQ1パケットの単一の共通の順序付けに合致するように出力ポートに現れる。各出力ポートは、共通の順序付けリストにおける幾つかの又は全てのパケットを送信することができる。

【0085】

コヒレンス性の目的で全体的なシステム順序を維持するのに加えて、スイッチから出力されるパケットを、アドレス及びデータバスの性能が完全に実現されるように順序付けすることも望まれる。例えば、図14は、HSリンク170のアドレス及びデータバス構造の利用を示すタイミング図である。

HSリンク170は、2対の単方向性アドレス及びデータバスによりマルチプロセッサノード100の各々に接続される。データバスは、512ビットのデータパケットを搬送し、そしてアドレスバスは、80ビットのアドレスパケットを搬送する。データパケットの送信は、アドレスパケットの送信の2倍のサイクル数を必要とする。書き込みコマンドのようなあるコマンドは、アドレス及びデータパケットの両方を含む。例えば、図14において、アドレスパケット179aは、データパケット179dに対応する。各コマンドがアドレス及びデータパケットの両方を含む場合には、アドレスバスの1つおきのアドレススロットがアイドル状態となる。しかしながら、読み取りコマンドのような多数のコマンドは、アドレスパケットしか含まず、データパケットを転送するためのデータバスのスロットを必要としない。従って、全体的なシステム性能を向上するためには、データ部分及びアドレス部分の両方が「バック」され、即ちHSリンクのアドレス及びデータ部分の各考えられるタイムスロットにアドレス及びデータが存在するような順序でバスから送出すべきパケットを選択するスイッチを有するのが好ましい。アドレス及びデータがHSリンクにおいて「バック」されるときには、HSリンクが最適に利用される。

【0086】

多数の入力ポートを経て多数のソースからデータを同時に受け取りそして多数の出力ポートを経て多数の行先ヘデータを供給できる一方、データ依存性を満足し、システム順序を維持し、そしてデータ転送レートを最大にすることのできるハイアラキースイッチを実施するための種々の実施形態が提供される。これらの種々の実施形態を、図15ないし18を参照して説明する。

図15には、上記順序付け制約を実施することのできるスイッチ181の1つの実施形態が示されている。図11について述べたように、スイッチ155は、複数のバッファ160a - 160hを含む。入力バッファの各々は、1書き込みポート/8読み取りポートバッファであり、8個の各入力の1つからパケットを受け取るように接続される。又、スイッチは、8個の出力ポートも含むが、1つの出力ポート、即ち出力ポート<0>のみに対するロジックが示されている。残りの出力ポートに対するロジックも同様であり、明瞭化のために、ここでは詳細に述べない。

【0087】

本発明の1つの実施形態では、各バッファの各エントリは、バッファのエントリに記憶されるパケットのチャンネルを識別するチャンネルフィールド185を含む。更に、各エントリは、一連のリンクインデックス186を含む。各リンクインデックスは、入力バッファ160a - 160hのエントリの1つに対するインデックスである。これらのリンクインデックスは、パケット順序付け制約に基づきバッファ160aから同じチャンネルを経て次々のパケットをアクセスするためのリンクリストアドレス構造体を形成するのに使用される。3つのリンクインデックスL1、L2及びL3が各エントリと共に含まれ、各リンクインデックスは、3つまでの順序付けリストの1つにおけるエントリの位置を識別する。

又、各エントリは、依存性フラグ189も含む。依存性フラグは、チャンネル間の依存性をマークするのに使用される。依存性フラグF1は、対応するエントリのパケットがQ1、QIO又はQ0Vicチャンネルを進行するパケットである場合にセットされる。依存性フラグF2は、対応するエントリのパケットがQ0又はQ0Vicチャンネルを進行するパケットである場合にセットされる。依存性フラグは、パケットの処理順序を次のように維持する上で助けとなる。

【0088】

概念的に、受け取ったパケットは、Q2チャンネル待ち行列、合成Q1/QIO/Q0Vicチャンネル待ち行列、合成Q0/Q0Vicチャンネル待ち行列、Q0Vicチャンネル待ち行列及びQIO待ち行列を含む5つの順序付けされた待ち行列に分割される。従って、パケットは、2つ以上の待ち行列に含まれる。ヘッドポインタは、各待ち行列ごとに1つのポインタ187a - 187eを含む。ヘッドポインタは、その待ち行列に対応するバッファにおける次のパケットを識別するバッファ160a - 160hのインデックスを与えるのに使用される。従って、ヘッドポインタ187は、Q2ヘッドポインタ187a、Q1/QIO/Q0Vicヘッドポインタ187b、Q0/Q0Vicヘッドポインタ187c、Q0Vicヘッドポインタ187d及びQIOヘッドポインタ187eを含む。パケットが入力バッファに最初書き込まれるときには、それが1つ以上の順序付けされた待ち行列に入れられる。1つ以上の順序付けされた待ち行列に入れられるときには、1つ以上の依存性フラグ189がアサートされる。チャンネルの形式及び依存性フラグが検査されて、チャンネル依存性を満足するように出力すべきバッファの適当なエントリが選択される。

【0089】

8個の入力バッファ160a - 160h各々の各エントリは、マルチプレクサ182へ送られる。マルチプレクサ182は、マネージャー180からの選択信号に応答して入力バッファの1つからパケットの1つを選択する。マネージャー180は、入力バッファ160a - 160hの64個の考えられる読み取りポートからのエントリを関連出力ポートの出力として選択する。マネージャー180は、全体的なシステム順序及びチャンネル依存性が満足されるようにパケットを選択する。

入力バッファ 160a - 160h の 1 つにパケットが受け取られるときには、エントリのチャンネルフィールドにチャンネル形式が書き込まれ、そしてそのエントリの関連フラグがフラグフィールド 189 においてアサートされる。上述したように、入力バッファの各エントリごとに、3 つのリンクインデックスがあり、その各々は、3 つの順序付けされた待ち行列の 1 つに対応する。本発明の 1 つの実施形態では、パケットを 3 つの異なる出力ポートにマルチキャストするために多数のリンクインデックスが使用される。マルチキャストされるべきパケットが入力バッファに記憶されるときには、それが 2 つ以上のリンクされたリストに入れられ、リンクされたリストの各々は、異なる出力ポートに対応する。その結果、異なる出力ポートに関連する出力マネージャーは、各々、異なるリンクリストインデックスを用いて同じ入力バッファエントリにアクセスすることができる。

10

【0090】

上述したように、リンクインデックス値は、バッファ 160a - 160h において対応する形式の次のパケットをアドレスするためのバッファインデックス値である。従って、リンクインデックス値は、対応する形式のその後のパケットがバッファに書き込まれるまで書き込まれない。その後のパケットがバッファに書き込まれるときには、その後のパケットのアドレスが手前のパケットのリンクインデックスに書き込まれ、これにより、そのチャンネル形式の次のパケットのインデックスを与える。各エントリは、3 つの考えられるリンクインデックスフィールドを含むので、手前のエントリにアドレスを書き込むのに加えて、2 ビットフィールド（図示せず）がアドレスと共に記憶され、順序付けリストを構成するために 3 つのリンクインデックスの適当な 1 つをエントリで識別できるようにする。

20

【0091】

マネージャー 180 は、出力ポートへ供給するためにバッファ 160a - 160h のパケットの 1 つを次のように選択する。上述したように、ヘッドポインタ 187a - 187e は、各待ち行列の最上部に対応するバッファインデックスを記憶する。所与のチャンネルに対するパケットを処理するとき、マネージャーは、対応するヘッドポインタにより指示されたエントリを選択する。1 つ以上のフラグ 189 がセットされ、そして高い優先順位のチャンネルに関連した待ち行列のパケットが処理されていない場合には、パケットは、その待ち行列内のより優先順位の高い全ての手前のパケットが処理されるまで処理されない。

30

例えば、出力マネージャーが Q0 形式のパケットを処理する場合に、Q1 / QIO / Q0 Vic 及び Q0 / Q0 Vic ヘッドポインタで指示されたエントリを検査する。パケットが Q0 チャンネルパケットであるが、Q1 パケットの処理がまだ完了していない場合には、エントリは処理されない。パケットの処理は、チャンネル Q1 又は Q0 パケットが既に処理されたことを指示する処理フラグ（図示せず）を各フラグ F1 及び F2 と共に与えることにより指示される。高い優先順位のチャンネルを有する待ち行列における全てのパケットの処理が行われると（処理フラグにより指示される）、そのエントリに関連したパケットは自由に処理される。

【0092】

40

あるエントリが処理のために選択されると、マネージャーは、そのエントリが存在する待ち行列に関連したヘッドポインタをバッファインデックスとして選択する。バッファインデックスはマルチプレクサ 182 へ送られ、そしてバッファエントリが出力ポートへ送られる。リンクインデックスはヘッドポインタへ返送され、そしてヘッドリストポインタがその待ち行列の次のパケットのバッファインデックスで更新される。従って、図 15 のスイッチ実施形態は、リンクリストデータ構造体、順序付けされた待ち行列及びフラグを用いて、出力ポートへパケットを与え、全体的なシステム順序が維持されるようにする。更に、多数のリンクインデックスを含むリンクリストデータ構造体は、マルチキャストパケット順序付けルールに固執しながらパケットをマルチキャストするための簡単な機構を形成する。

50

【 0 0 9 3 】

従って、図 1 5 の実施形態は、フラグ及び順序付けされた待ち行列を使用して、チャンネルの順序が維持されるようにする。図 1 6 には、所定の順序依存性に基づいて出力データを与えることのできるスイッチの第 2 の実施形態が示されている。図 1 6 の実施形態では、スイッチの各出力ポートに対してバッファ 2 0 0 が設けられる。バッファ 2 0 0 は、入力パケット受信経路 2 0 1 を経てバッファ 1 6 0 a - 1 6 0 h (図 1 1) の各々から入力を受け取るように接続され、入力バッファからのパケットは、パケットの行先に基づいて出力ポートの適当なバッファへ送られる。本発明の 1 つの実施形態では、バッファは、コラプス(collapsing) F I F Oとして実施されるが、当業者に知られた他のバッファアーキテクチャ x を使用することもできる。

10

【 0 0 9 4 】

バッファ 2 0 0 は、スイッチから送出されるべき種々のパケットを記憶するように示されている。バッファ 2 0 0 は、ここでは、5 つの異なるチャンネル Q 0、Q 1、Q 2、Q 3 及び Q 4 を経て送信されるパケットを記憶する。チャンネル Q 0 - Q 4 は、上記のチャンネル Q 0、Q 1、Q 2、Q 0 V i c 及び Q I O と同様ではない。むしろ、これらは、単にスイッチの出力動作を示すためにのみ使用される。従って、パケット Q 0 - Q 4 は、異なるチャンネルにおける一般的パケットを表わし、チャンネルの依存性は、図 1 6 A の流れ図において矢印に基づいて定められる。図 1 6 A において、あるチャンネルから別のチャンネルへ向けられた矢印は、第 1 チャンネルのパケットが出力ポートへ送られず、一方、第 1 チャンネルのパケットの前に受け取られた第 2 チャンネルのパケットは、スイッチによる処理が保留中であることを指示する。例えば、図 1 6 A において、チャンネル Q 0 のパケットは、チャンネル Q 3 のパケットの処理に依存するように示されており、従って、チャンネル Q 0 のパケットは、チャンネル Q 3 のパケットを「プッシュ」したと言える。図 1 6 A の流れ図に示された付加的な依存性は、チャンネル Q 1 のパケットがチャンネル Q 2 及び Q 3 のパケットをプッシュしたことを指示する。この場合も、図 1 6 A の流れ図で表わされた依存性は、既に述べた Q 0、Q 1、Q 2、Q 0 V i c 及び Q I O チャンネルの依存性を表わすものではないことに注意されたい。以下に述べるように、Q 0、Q 1、Q 2、Q 0 V i c 及び Q I O チャンネルにおけるパケットの依存性は複雑であり、従って、バッファ 2 0 0 の動作を容易に説明するために、一般的パケット及び依存性が与えられる。

20

30

【 0 0 9 5 】

上述したように、入力パケットは、スイッチの入力バッファ 1 6 0 a - 1 6 0 h の各々に正しい順序で受け取られ、そしてそのパケットにより指示された行先に基づいて、バッファ 2 0 0 のような出力バッファに正しい順序で供給される。各出力バッファの各パケットエントリ、例えば、エントリ 2 0 0 a は、パケットの送信及び受信ノードを指示するソース及び行先フィールドと、パケットが送信されるチャンネルを指示するチャンネルフィールドと、一連のビット 2 0 6 a - 2 0 6 e とを備えている。一連のビット 2 0 6 a - 2 0 6 e は、ハイアラキースイッチを経てパケットを供給する各チャンネルごとに 1 ビットを含む。例えば、図 1 6 の実施形態では、一連のビットは、チャンネル Q 0、Q 1、Q 2、Q 3 及び Q 4 の各々について 1 ビットを含む。

40

【 0 0 9 6 】

出力ポートに対して入力パケット受信経路に接続された書き込み制御ロジック 2 0 5 は、受信パケットのチャンネルに基づくと共に、図 1 6 A の流れ依存性図に示されたチャンネル間の依存性に基づいて一連のビットの各々の設定を制御する。又、以下に詳細に述べるように、書き込み制御ロジックは、静的又は動的に依存性を確認することによりビットを更新することができる。依存性を静的に確認するときには、チャンネルに対して定められた依存性が、バッファ内の他のパケットに拘わりなく適用される。依存性を動的に確認するときには、チャンネルの依存性が、バッファ 2 0 0 内の他のパケットのチャンネル及びアドレス行先を考慮して適用される。

一連のビットの各々には、対応するサーチエンジン 2 0 8 a - 2 0 8 e が接続される。各

50

サーチエンジンは、ビットの関連列をサーチして、列セットの対応ビットを有するバッファ200のエントリを選択する。選択されたエントリは、各列（又はチャンネル）ごとに、一連の信号S4 - S0により出力バッファマネージャー202へ指示される。チャンネル間の既知のデータ依存性に関連してサーチエンジンの各々により受信された選択信号を用いて、出力バッファマネージャーは、グローバルポート出力に供給するために出力バッファ200からのパケットの1つを選択する。

【0097】

動作中に、入力パケット受信経路201を経てパケットが受信されるときに、パケットのチャンネルは、書き込み制御ロジック205により評価され、そして一連のビット206a - 206eのうちの、そのチャンネルに対応するビットがアサートされる。図16において、パケットの形式を指示するためにセットされたビットは、「丸内のX印」で示され、そしてこれはチャンネル識別子フラグと称する。従って、図16では、パケット1がQ3形式のパケットである。図15の実施形態によれば、エントリのチャンネルを指示するビットをアサートするのに加えて、そのチャンネルのパケットがブッシュするところの各チャンネルに対してビットが付加的にアサートされる。これらビットの各々は、依存性フラグと称され、図16に「X」で示されている。それ故、Q0チャンネルパケットであるパケット2の場合に、Q3チャンネルパケットに関連したビットが付加的にアサートされる。というのは、図16Aの流れ図に示されるように、Q0パケットがQ3パケットをブッシュするからである。

【0098】

パケットがバッファ200に記憶され、そしてそれらの関連する一連のビット206a - 206eがアサートされるときには、ビットの各列に関連したサーチエンジン208a - 208eの各々が、ビットセットを有するバッファ内の第1エントリを選択する。それ故、サーチエンジン208aの選択値は、パケット2を指し、サーチエンジン208bの選択値は、パケット3を指し、等々となる。

S0 - S4信号は、マネージャー202に送られる。マネージャー202は、サーチエンジンによる選択信号のアサートに応答するのに加えて、システムに存在する依存性に応答して、パケットの1つを選択する。例えば、本発明の1つの実施形態によれば、チャンネルQ0にあるパケット2のようなパケットは、チャンネルQ0のサーチエンジン（208a）及びチャンネルQ3のサーチエンジン（208d）の両方が同じパケットを選択しない限り、スイッチから送出されない。従って、多数のフラグが所与のパケットに対してセットされたときに、マネージャー202は、セットされたフラグに対応するサーチエンジンの両方がその所与のパケットを選択しない限り、出力に対してそのパケットを選択しない。

【0099】

本発明の別の実施形態によれば、サーチエンジンが、その依存性フラグがセットされたためにエントリを選択した場合に、サーチエンジンは、依存性フラグをクリアしそしてバッファを下方に進んで、依存性フラグ又は認識フラグがセットされた次のエントリを選択することができる。このような構成では、サーチエンジンが他のチャンネルによりストールされて処理を保留にすることがないので、パケットの処理が改善される。

依存性を識別するために多数のフラグをアサートする作用は、パケットがスイッチを経て伝播するときにパケットの全体的なシステム順序を維持する上で助けとなる。例えば、図16において、Q0パケットとQ3パケットとの間の関係は、Q0チャンネルパケットが実行の前に各手前のQ3チャンネルパケットをブッシュすることである。従って、Q3チャンネルパケットの後に受け取られたQ0チャンネルパケットは、Q3パケットの前に実行されてはならない。パケット1は、パケット2のQ0チャンネルパケットの前に受け取られるQ3チャンネルパケットである。パケット2に対してビット206dをセットすることにより、パケット2のQ0パケットがパケット1のQ3パケットの前に出力ポートに発生されないよう確保することができる。というのは、マネージャー208は、S3及びS0の両方がパケット2への信号を選択するまでQ0パケットを選択しないからである。

S 3 値は、パケット 1 が処理されるまでパケット 2 を指さない。その結果、所与のチャンネルのパケットによりプッシュされた各パケットごとにビットをアサートすることにより、所与のチャンネルによりプッシュされたパケットが処理されるまでチャンネルが効果的に阻止される。その結果、全体的なシステム順序が維持される。

【 0 1 0 0 】

上記のように、図 1 6 のバッファ制御ロジックは、静的又は動的な依存性を確認するように動作される。静的な依存性とは、図 1 6 A の流れ図で示されたような依存性である。動的な依存性は、バッファ内の 2 つのパケット間に静的な依存性が実際に存在するかどうかを決定するためにバッファの内容を評価することにより確認される。静的な依存性は、メモリデータが SMP においてコヒレンス性を失わないよう確保する順序付けルールを形成するの
10
に使用される。しかしながら、データのコヒレンス性は、パケットがメモリデータの同じブロックをアクセスする場合にしか影響されない。それ故、動的な依存性は、バッファに既にあるパケットの行先アドレスを検査することによってバッファの内容を微細な粒度で検査して、異なるチャンネルの 2 つのパケット間に依存性が実際に存在するかどうか決定する。

【 0 1 0 1 】

バッファ 2 0 0 内のパケット間の依存性を動的に確認する 1 つの効果は、バッファ内のパケットを処理するのに必要な時間を短縮することである。例えば、上記のパケット 1 及びパケット 2 の動作を使用すると、Q 0 パケット 2 及び Q 3 パケット 1 が同じアドレスにマップしない場合には、何ら問題なく、Q 0 パケットを Q 3 パケットの前に処理することが
20
できる。手前の Q 3 パケットの処理を待機する際に受ける遅延時間が排除され、これにより、SMP システムの全体的な性能が改善される。

例えば、図 1 7 は、依存性を動的に確認することによるプロセスへのパケットの選択動作を示すフローチャートである。ステップ 2 2 0 において、パケットがバッファ 2 0 0 に受け取られる。ステップ 2 2 2 において、パケットのチャンネルに対するビットが書き込み制御ロジック 2 0 5 により一連のビット 2 0 6 においてセットされる。ステップ 2 2 4 において、バッファ 2 0 0 に記憶された手前のパケットが検査されて、パケットがプッシュする
30
ところのチャンネルのパケットがメモリの同じブロックにあるかどうか決定される。それらがメモリの同じブロックにある場合には、ステップ 2 2 6 において、パケットがプッシュするところのチャンネルにあり且つ同じメモリブロックに存在するパケットに対応するビットがアサートされる。従って、パケット 2 に対して図 1 6 の例を使用すると、パケット形式 Q 3 に対するビットは、パケット 1 がパケット 2 と同じメモリブロックをアクセスする場合だけアサートされる。従って、依存性を動的に確認することにより、全体的なシステム性能を向上しながらメモリコヒレンス性を維持することができる。

【 0 1 0 2 】

図 1 8 には、全体的なシステム順序を維持しながら、多数の入力ソースから受け取ったデータを多数の出力ソースへ出力する方法の別の実施形態が示されている。図 1 8 の実施形態は、図 1 6 の場合と同様の要素を含むように示されている。しかしながら、図 1 8 の書き込み制御ロジック 2 0 9 は、パケットの依存性を異なるやり方で分析することにより一連のビット 2 0 6 a - 2 0 6 e の各々を更新する。図 1 6 の場合のように、パケットが関連
40
チャンネルのものであることを指示するために、一連のビットの 1 つが各パケットごとにセットされる。しかしながら、チャンネルがプッシュするところのチャンネルの全てのパケットに対して付加的なビットをセットするのではなく、そのチャンネルのパケットをプッシュするところのチャンネルのパケットに対してビットがセットされる。

【 0 1 0 3 】

従って、図 1 8 の実施形態は、チャンネル識別フラグをセットするのに加えて、そのパケットによりマスク又は阻止された全てのチャンネルに対して付加的なビットがセットされる。例えば、図 1 8 の例において、パケット 1 は、Q 3 チャンネルパケットである。Q 3 チャンネルのパケットは、図 1 8 A の依存性流れ図に示すように Q 3 パケットが実行される
50
まで、Q 1 及び Q 0 パケットの実行を阻止する。従って、ビット 2 0 6 d、2 0 6 b 及

び206aがパケット1に対してセットされる。しかしながら、パケット2は、他のパケットの実行を阻止しないQ0パケットである。その結果、ビット206bのみがパケット2に対してセットされる。

従って、図18のスイッチ実施形態は、依存性を静的に確認することによりシステム順序を維持しながら出力ポートヘデータを供給する別の方法を提供する。図18のバッファ実施形態は、依存性を動的に確認するようには使用できないことに注意されたい。というのは、そのようにするには、データがバッファ200に書き込まれる前にデータのアドレスを知る必要があるからである。しかしながら、ここに述べる静的及び動的な方法は、全て、パケット間の依存性を満足するよう確保するために使用できる。

【0104】

従って、多数の入力ポートを経て多数のソースからデータを同時に受け取りそして多数の出力ポートを経て多数の行先ヘデータを供給できる一方、データ依存性を満足し、システム順序を維持し、そしてデータ転送レートを最大にすることのできるスイッチの3つの実施形態が説明された。1つの実施形態では、フラグを記憶する多数の待ち行列の使用により順序付け依存性が達成されそして依存性を識別するように待ち行列が選択されるリンクリストバッファ機構が説明された。第2および第3の実施形態では、スイッチの入力バッファからデータを正しい順序で受け取る出力バッファが、ある形式のパケットを阻止するのに使用される一連のビットを備えていて、データ依存性及びコヒレンス性制約を満足するように確保する。全ての実施形態において、潜在的な依存性の競合をマークするためにセットされるフラグを含む順序付け待ち行列の使用により、順序付け依存性が追跡される。フラグの順序付けリストを用いて依存性を識別することにより、バスの利用性を最大にししながら順序を維持し且つコヒレンス性を確保するためにマネージャーにより実行されるオペレーションの複雑さが簡単化される。

【0105】

キャッシュコヒレンスプロトコル

本発明の1つの実施形態におけるキャッシュコヒレンスプロトコルは、書き込み無効化所有権をベースとするプロトコルである。「書き込み無効化」とは、プロセッサがキャッシュラインを変更するときに、他のプロセッサキャッシュにおける効力のないコピーを無効化することを意味し、新たな値でそれらを更新するのではない。このプロトコルは、システム内のメモリであるかプロセッサ又はIOPの1つであるかに拘わりなくキャッシュラインに対する識別可能な所有者が常に存在するので、「所有権プロトコル」と称される。キャッシュラインの所有者は、必要なときにキャッシュラインの最新の値を供給する責任がある。プロセッサ/IOPは、キャッシュラインを「独占的に」又は「共用して」所有することができる。プロセッサがキャッシュラインの独占的所有権を有する場合には、システムに通知せずにそれを更新することができる。さもなくば、システムに通知し、そして他のプロセッサ/IOPキャッシュのコピーを潜在的に無効化しなければならない。

【0106】

キャッシュコヒレンスプロトコルの詳細な説明に入る前に、ハイアラキーネットワークに使用される全通信手順について最初に説明する。

図7Aについて述べたように、大型のSMPシステム150は、スイッチ155を経て互いに接続された多数のノードを含む。各ノードにおける各プロセッサは、メモリのデータにアクセスするコマンドを発生する。これらのコマンドは、ソースノード内で完全に処理することもできるし、又はアドレス及び要求の形式に基づいてシステムの他のノードへ送信することもできる。

アドレススペースは、メモリスペース及びIOスペースに分けられる。プロセッサ及びIOPは、専用キャッシュを使用して、メモリスペースアドレスのみに対するデータを記憶し、そしてIOスペースデータは、専用キャッシュには記憶されない。従って、キャッシュコヒレンスプロトコルは、メモリスペースコマンドのみに関連している。

【0107】

キャッシュコヒレンスプロトコルの重要な要素は、ロード及び記憶動作をシリアル化する

解決策である。キャッシュコヒレンスプロトコルは、各メモリアドレスXへの全てのロード及び記憶に順序を課さねばならない。この順序は、Xへの全ての「記憶」が順序付けされるものであり、即ち、第1記憶、第2記憶、第3記憶、等々とならねばならない。第i番目の記憶は、(I - 1)番目の記憶により決定されたようにキャッシュラインを更新する。更に、各ロードには最新の記憶が関連され、そこからロードはキャッシュラインの値を得る。この順序をここでは「ロード - 記憶シリアル化順序」と称する。

ここに述べるプロトコルの特性は、アドレスXに対するホームARBバスが、Xへの全てのロード及び記憶に対する「シリアル化ポイント」であることである。即ち、Xへの要求がXのホームARBバスに到着する順序は、対応するロード及び記憶がシリアル化される順序である。大型のSMPシステムに対するほとんどの公知のプロトコルは、この特性を有しておらず、従って、効率が悪く、複雑である。

10

【0108】

図2に示す小型のSMPノードシステムには、1つのARBバスがある。このバスは、小型SMPにおける全てのメモリロード及び記憶に対するシリアル化ポイントである。ARBバスに接続されたDTAGは、小型SMPのプロトコルに必要とされる全ての状態を捕獲する。大型のSMPシステムでは、ホームARBバスのDIRがプロトコルに対するおおよその状態を捕獲し、TTT及びDTAGは、より微細なレベルにおける状態情報を捕獲する。要求RがホームARBバスに到着すると、DIR、DTAG及びTTT状態が検査され、他のプロセッサへの調査コマンド及び/又はソースプロセッサへの応答コマンドを発生することができる。更に、DIR、DTAG及びTTTの状態が要求Rの「シリアル化」を反映するように自動的に更新される。従って、要求アドレスがRのアドレスに等しく且つ要求Rの後にホームARBに到着する要求Qは、ロード - 記憶シリアル化順序においてRの後に現れる。

20

【0109】

その結果、ホームARBバスは、メモリアドレスへの全ての要求に対し「シリアル化ポイント」と定義される。各メモリアドレスXに対し、対応する要求(RdMod又はCTD)がホームARBバスに到着する順序で記憶が見掛け上実行される。アドレスXへのロードは、ホームARBにおいて最後にシリアル化された記憶Xに対応するXのバージョンを得る。以下に述べるキャッシュコヒレンスプロトコルの前書きにおいて、「システム」という用語は、プロセッサ及びIOPを除く大型SMPの全ての要素を指す。プロセッサ及びシステムは、「コマンドパケット」又は単に「コマンド」を送信することにより互いに対話する。コマンドは、要求、調査及び応答の3つの形式に分類される。

30

プロセッサによりシステムに発生されるコマンド及びシステムによりプロセッサに発生されるコマンドは、所与のプロセッサのメモリシステムインターフェイスに基づく。SMPの動作を説明する目的上、デジタル・イクイップメント・コーポレーションからのAlpha(登録商標)システムインターフェイスの定義に基づいて発生される要求及びコマンドについて説明するが、他の形式のプロセッサも使用できることを理解されたい。

【0110】

要求は、ロード又は記憶動作を実行する結果として、データのコピーを得なければならないときにプロセッサにより発生されるコマンドである。又、要求は、システムからのデータの断片に対して独占的な所有権を得るのにも使用される。要求は、読み取りコマンド、読み取り/変更(RdMod)コマンド、ダーティへの変更コマンド、ピクティムコマンド、及びエビクト(Evict)コマンド(データのキャッシュラインが各キャッシュから除去される場合)を含む。

40

調査(Probe)コマンドは、データ及び/又はキャッシュタグ状態更新を要求する1つ以上のプロセッサへシステムにより発生されるコマンドである。調査コマンドは、送信読み取り(Forwarded Read)(FRd)コマンド、送信読み取り変更(Forwarded Read Modify)(FRdMod)コマンド、及び無効化コマンドを含む。プロセッサPがシステムへの要求を発生するときには、システムは、1つ以上の調査コマンドを他のプロセッサへ発生しなければならない。Pがキャッシュラインのコピーを要求する(読み取り要求で)場合には

50

、システムは、所有者プロセッサ（もしあれば）へ調査コマンドを送信する。Pがキャッシュラインの独占的所有権を要求する（CTD要求で）場合には、システムは、キャッシュラインのコピーをもつ1つ以上のプロセッサへ無効化調査コマンドを送信する。Pがキャッシュラインのコピー及びキャッシュラインの独占的所有権の両方を要求する（RdMod要求で）場合には、システムは、データのキャッシュラインのダーティコピーを現在記憶しているプロセッサにFRdコマンドを送信する。FRdコマンドに回答して、キャッシュラインのダーティコピーがシステムに返送される。又、送信読み取り変更（FRdMod）コマンドも、キャッシュラインのダーティコピーを記憶しているプロセッサにシステムにより発生される。FRdModに回答して、ダーティキャッシュラインがシステムに返送され、そしてキャッシュに記憶されたダーティコピーが無効化される。キャッシュラインを別のプロセッサにより更新すべきときには、キャッシュラインのコピーをキャッシュに記憶しているプロセッサに、システムにより無効化コマンドが発生される。

10

【0111】

応答は、プロセッサにより要求されたデータ又は要求に対応する確認を搬送するシステムからプロセッサ/IOPへのコマンドである。読み取り及びRdModコマンドの場合に、応答は、各々要求されたデータを搬送するFill又はFillModコマンドである。CTDコマンドの場合に、応答は、CTDの成功又は失敗を指示するCTD成功又はCTD失敗コマンドである。ピクティムコマンドの場合には、応答がピクティム・リリースコマンドである。

図19には、要求と要求との間の関係、及び個々のプロセッサにおける対応キャッシュラインの状態を説明するためのテーブルが示されている。又、図19は、キャッシュラインの要求及び状態の各々に対して得られる調査形式のコマンドも示している。カラム300及び300aは、プロセッサにより発生される要求を示し、カラム305及び305aは、システムの他のプロセッサにおけるキャッシュの状態を示し、そしてカラム320及び320aは、システムにより発生される調査コマンドを示す。

20

【0112】

図19のテーブルは、プロセッサAと称するプロセッサがシステムに要求を発生することを仮定している。プロセッサAのコマンドは、次いで、プロセッサBと称する1つ以上の他のプロセッサと相互作用する。プロセッサAによりアドレスされるキャッシュラインが、DTAG及び/又はディレクトリ情報を用いて決定されたプロセッサBのキャッシュに記憶される場合には、プロセッサBのキャッシュ状態が、プロセッサBへ調査コマンドを発生する必要があるかどうか及びどんな形式の調査コマンドを発生すべきかを決定する。以下、コヒレンスプロトコル及び機構について詳細に述べる。コマンドパケットがとる経路、各コマンド形式に対する状態情報のソース、及びそれにより生じるアクションが含まれる。全てのコマンドは、プロセッサ又はIOPから発生され、IOPの発生プロセッサは「ソースプロセッサ」と称する。要求に含まれるアドレスは、「要求アドレス」と称する。アドレスの「ホームノード」は、そのアドレススペースが要求アドレスをマップするところのノードである。要求は、ソースプロセッサが要求アドレスのホームノードである場合は「ローカル」と称し、さもなくば、「グローバル」要求と称する。ホームノードのARBバスは、「ホームARBバス」と称する。「ホームディレクトリ」は、要求アドレスに対応するディレクトリである。従って、ホームディレクトリ及びメモリは、要求アドレスに対するホームARBバスに接続される。

30

40

【0113】

プロセッサ又はIOPから発せられるメモリ要求は、まず、ホームARBバスにルート指定される。この要求は、それがローカルである場合にはローカルスイッチを経てルート指定され、それがグローバルである場合にはハイアラキースイッチを経て送られる。後者の場合には、ローカルスイッチ及びGPリンクを横断してGPに達し、次いで、HSリンクを経てハイアラキースイッチへ至り、次いで、GP及びホームノードのローカルスイッチを経てホームARBバスへ至る。

グローバル要求は、ソースノードのARBバスに最初に現れず、むしろ、GPリンクを経

50

てHSに直接ルート指定されることに注意されたい。公知のプロトコルでは、グローバル要求は、それが別のノードへ送出される前にソースノードの状態をアクセスする。本発明は、グローバル要求をHSへ直接発生することによりグローバル要求の平均待ち時間を短縮する。

【0114】

図20A - 20Jは、多数の基本的なメモリランザクションを例示するフローチャートである。

ローカル読み取り：

図20Aにおいて、ソースプロセッサ320からホームARBバスへ要求が送られる。ディレクトリ322は、どのプロセッサがメモリブロックを所有するか決定する。ローカルメモリ323が所有者である場合には、「短い記入」コマンドがホームARBバスからソースプロセッサ320へ発生される。

10

【0115】

グローバル読み取り：

図20Bにおいて、ノード325のプロセッサ320が、「ホーム」がノード326にあるメモリのキャッシュラインへ読み取りを発生すると仮定する。(グローバル)読み取りコマンドは、ライン327で示された経路を経てスイッチ324を通り「ホーム」ARBバス及びディレクトリ321へルート指定される。ノード326のメモリ330がキャッシュラインの所有者である場合には、「短い記入応答」を発生するノード326によりノード326からノード325へデータが返送される。

20

キャッシュラインが別のプロセッサ/IO Pにより現在所有されている場合には、要求されたキャッシュラインを得るために異なるステップが取られる。図20Cを参照すれば、プロセッサ320が、「ホーム」がノード326にあるメモリのキャッシュラインへ読み取りを発生する場合には、読み取りは、再び、経路327を経てホームARBバス及びディレクトリ321へルート指定される。ディレクトリ321のエントリは、上述したように、メモリの各キャッシュラインに対し、所有者情報を含む14ビットの状態情報を備えている。所有者情報は、この場合に、所有者をノード328におけるプロセッサ342として識別する。

ノード328が要求されたキャッシュラインを所有するというディレクトリの指示に応答して、2つの事象が生じる。第1に、「ホーム」ノードであるノード326は、ライン329で示すように、所有者プロセッサ342へ「送信読み取り」調査を発生する。同時に、ホームノード326は、ライン331で示すように、プロセッサ320へ「記入マーカー」応答を送信する。「記入マーカー」応答の役割は、以下で説明する。

30

【0116】

「送信読み取り」に応答して、プロセッサ342は、「記入」コマンドをプロセッサ320へ発生し、「記入」コマンドは、当該キャッシュラインを含む。「読み取り」要求に対するこの形式の応答は、データ返送に対して一連の3つのコマンドを必要とするので、「長い記入」と称される。従って、「読み取り」ランザクションは、メモリからの応答である「短い記入」と、所有者プロセッサからの応答である「長い記入」の2つの形式に分割することができる。

40

ローカルRdMod：

図20Dを参照すれば、ローカル読み取り変更ランザクションは、(1)キャッシュラインの現在バージョンのコピーを得ている全てのプロセッサに無効化調査が送られ、そして(2)FRMod及びFillModsが、Frds及びFillsに代わって所有者に送られる点を除くと、ローカル読み取りランザクションと同様に働くことが明らかである。図20D図において、ホームノードのディレクトリは、ローカルプロセッサ又はメモリがブロックを所有することを示す。ホームARBバスにおいて、ディレクトリ322は、ブロックの現在バージョンを得ている全ての外部ノードを識別する。無効化コマンドは、HS324へ送られ、全ての当該ノードはマルチキャストベクトルで識別される。HSは、ベクトルで識別された全てのノードへ無効化メッセージをマルチキャストする。

50

無効化メッセージは、各ノードにおいて A R B バスへ進み、そこで、D T A G は、それらを更にフィルタし、キャッシュラインの現在バージョンを有すると識別されたプロセッサ又は I O P のみへ無効化調査を送信する。

【 0 1 1 7 】

グローバル R d M o d :

図 2 0 E を参照すれば、読み取り変更トランザクションは、図 2 0 A 及び 2 0 B について述べた読み取りトランザクションと同様に作用することが明らかである。読み取り変更 (R d M o d) コマンドは、先ず、プロセッサ 3 2 0 からキャッシュラインのホーム A R B 及びホームディレクトリ 3 2 1 ヘルツ指定される。ホームノードであるノード 3 2 6 のメモリがキャッシュラインを記憶する場合には、要求されたデータを含む「短い記入変更」コマンドがノード 3 2 6 からプロセッサ 3 2 0 へ送られる。ディレクトリ 3 2 1 は、このトランザクションの結果として更新される。

「読み取り変更」コマンドは、プロセッサ 3 2 0 が、キャッシュラインの内容を変更できるようにキャッシュラインの独占的所有権を要求することを指示する。それ故、「短い記入変更」コマンドに加えて、ノード 3 2 6 は、キャッシュラインの現在バージョンのコピーを得ている他の全てのプロセッサに「無効化」コマンドを発生する。D I R は、1 つ以上のプロセッサがキャッシュラインの現在バージョンのコピーを得ているノードを識別する。D I R の存在ビットは、この情報を含む。D T A G は、キャッシュラインのコピーを得ている全てのホームノードプロセッサを識別する。各々の D I R 存在ビットがセットされた全てのノードに「無効化」が送信される。「無効化」を受信する各ノードにおいて、D T A G をアクセスして、どのプロセッサがキャッシュラインのコピーを現在記憶するかを決定する。「無効化」は、これらのプロセッサのみに送られる。I O P タグは、I O P がコピーを有するかどうか決定するのに使用され、もしそうであれば、I O P は「無効化」調査も受け取る。

【 0 1 1 8 】

要求を発しているプロセッサ以外のプロセッサが所有者である場合には、ホームノードは、「記入変更マーカー」、「送信読み取り変更」及びゼロ以上の「無効化」を 1 つのコマンドとして発生する。スイッチにおいて、コマンドは、全ての行先ノードにマルチキャストされる。各行先ノードにおいて、コマンドは、その要素に分離され、各ノードのグローバルポートは、各ノードにおいてどんなアクションをとるべきかを決定する。上記の例では、「送信読み取り変更」がプロセッサ 3 4 2 により処理され、そして「記入変更マーカー」がプロセッサ 3 2 0 により処理される。更に、D T A G エントリに基づき、ホームノード、「記入変更マーカー」を受け取るノード、及び「送信変更」を受け取るノードにおいて「無効化」が実行される。「送信読み取り変更」に応答して、ダーティデータが「長い記入変更」コマンドを経てプロセッサ 3 4 2 からプロセッサ 3 2 0 へ送られる。

【 0 1 1 9 】

従って、「読み取り変更」コマンドは、2 つ又は 3 つのノード接続即ち「ホップ」を実行することができる。本発明の 1 つの実施形態では、読み取り型コマンド (「読み取り」及び「読み取り変更」) のみが 3 つのホップを生じ、但し、第 3 のホップは「記入」型コマンド (「記入」又は「記入変更」) である。しかしながら、本発明は、以下に述べる追加コマンドを仮想チャンネル待ち行列に適当に割り当てることにより 3 つ以上のホップを必要とする他のトランザクションも含むように容易に変更できる。

C T D :

図 2 0 G 及び 2 0 H には、クリーン - ダーティ (C T D) 及び無効化 - ダーティ (I T D) の基本的な流れが示されている。図 2 0 G では、クリーン - ダーティは、ホームノードにおいてプロセッサ 3 2 0 からディレクトリ 3 2 1 へ発生される。プロセッサ 3 2 0 が更新を希望するところのクリーンキャッシュラインが現在のものであるか効力を失ったものであるかに基づいて、「確認」コマンド (A C K) 又は「非確認」 (N A C K) コマンドのいずれかがプロセッサ 3 2 0 へ返送される。対応的に、C T D は成功又は失敗と言える。更に、C T D が成功の場合にデータのキャッシュラインのコピーをもつものとしてディ

レクトリ 3 2 1 の存在ビットにより指示された全てのノードに「無効化」が送られる。

【 0 1 2 0 】

図 2 0 H に示すように、I T D コマンドは、C T D と実質的に同様に働く。しかしながら、I T D は決して失敗とならない。A C K が常にプロセッサ 3 2 0 に送られ、そしてデータのキャッシュラインのコピーを記憶するシステムの他のノードには「無効化」が送られる。

ローカル及びグローバル書き込みビクティム：

上記のように、書き込みビクティムコマンドは、ダーティデータをプロセッサのキャッシュから適当なホームメモリへ返送する。図 2 0 I 及び 2 0 J を参照すれば、書き込みビクティムの流れは、「ホーム」メモリが書き込みビクティムを発生するプロセッサと同じノードであるかどうかに基づいて若干異なることが明らかである。図 2 0 I に示すように、「ホーム」ノードがプロセッサのノードである場合には、プロセッサ 3 2 0 が書き込みビクティムを発生し、そしてデータは、同じノードのメモリへ直接送られる。

【 0 1 2 1 】

しかしながら、図 2 0 J に示すように、ビクティムデータがプロセッサとは異なるホームにある場合には、データが 2 つの段階で転送される。第 1 に、ビクティムキャッシュラインがプロセッサ 3 2 0 のキャッシュ（又はビクティムバッファ）から送出され、そしてプロセッサノードのグローバルポートにおけるビクティムキャッシュ（図 6 の要素 1 2 4 ）に記憶される。ビクティムキャッシュは、「ビクティムリリース」信号でプロセッサに応答し、プロセッサがそのビクティムバッファエントリを再使用できることを指示する。次いで、スイッチに使用可能な帯域巾が存在するときには、ビクティムデータは、「書き込みビクティム」コマンドによりビクティムキャッシュからホームプロセッサのメモリへ送られる。

【 0 1 2 2 】

ソースプロセッサ P によりホームメモリに送られたビクティムデータは、それがメモリに到達するときまでに効力を失うことがあることに注意されたい。このような場合に、ビクティムは、「失敗」と言われ、ホームメモリは更新されない。このようなケースは、P がキャッシュラインの所有権を獲得するときと、P のビクティムがホームディレクトリに到着するときとの間のインターバルに別のプロセッサがキャッシュラインの所有権を獲得したときに生じる。このような場合には、P のビクティムがホーム A R B に到達する前に、キャッシュラインに対する「無効化」又は「F r d M o d」調査をプロセッサ P に送信しなければならない。

ビクティムデータをメモリに書き込まねばならないかどうか決定するために、「書き込みビクティム」コマンドがホーム A R B バスに現れるときに、要求されたアドレスに対するディレクトリエントリがルックアップされる。ソースプロセッサが依然としてキャッシュラインの所有者であることをディレクトリが指示する場合には、ビクティムが成功となり、メモリを更新する。さもなくば、失敗となり、メモリは更新しない。いずれにせよ、ディレクトリ 3 2 1 においてビクティムに対して判断がなされると、「ビクティム A C K」コマンドがノード 3 2 5 のグローバルポートに返送され、ビクティムキャッシュは関連エントリをクリアすることができる。

【 0 1 2 3 】

この設計の 1 つの実施形態では、D T A G を使用して、「書き込みビクティム」コマンドがローカルである場合に「書き込みビクティム」コマンドの成功又は失敗を判断する。この特定の例（ローカル「書き込みビクティム」要求の例）では、D T A G 及び D I R の両方が「書き込みビクティム」要求の成功又は失敗を決定するに必要な情報を与えることができる。D T A G は、単に D T A G をベースとする機構が小型の S M P ノードハードウェアに対して既に設けられているという理由で D I R に代わって使用される。

キャッシュコヒレンスプロトコルの上記説明では、最も一般的な動作及びコマンド形式について述べた。これら機構は、以下に詳細に説明する。

上記のように、本発明の 1 つの実施形態では、2 つ以上の関連メッセージパケットを効率

10

20

30

40

50

化のために１つに結合することができる。結合されたパケットは、次いで、ＨＳ又はノードのＡＲＢバスにおいてその成分に分割することができる。例えば、ＨＳへのＦｒｄＭｏｄメッセージは、所有者プロセッサをもつノードへのＦｒｄＭｏｄメッセージと、キャッシュラインのコピーをもつノードへの「無効化」メッセージと、ソースノードへのＦｉｌｌＭａｒｋｅｒＭｏｄメッセージとに分割される。所有者プロセッサノードへのＦｒｄＭｏｄは、ノードのＡＲＢバスにおいて、所有者プロセッサへのＦｒｄＭｏｄメッセージと、ノードの他のプロセッサへのゼロ以上の「無効化」メッセージとに更に分割される。

【０１２４】

ビクティムコヒレンス性を維持するための遅延書き込みバッファ動作：

図２０Ｉ及び２０Ｊについて上述したように、ホームメモリに送られるビクティムデータは、「書き込みビクティム」がホームＡＲＢに到達する前に受け取られるキャッシュラインに対し「無効化」又はＦｒｄＭｏｄ調査が介在する結果としてそれが到着するときまでに効力を失うことがある。

ビクティムデータをメモリに書き込まねばならないかどうかを決定する１つの方法は、各書き込みビクティムコマンドに対してディレクトリエントリをルックアップすることである。ビクティム書き込みコマンドを発生するプロセッサがダーティ所有者であることをディレクトリが指示する場合には、ビクティムを進めることが許されねばならない。さもなければ、失敗となってしまふ。この方法が望ましい理由は、プロセッサとシリアル化ポイントとの間のビクティム書き込みコマンドを、シリアル化ポイントとプロセッサとの間の調査コマンドと一致させるための複雑な比較論理構造体の必要性が回避されるからである。

【０１２５】

この解決策は、データコヒレンス性の維持を簡単化するが、メモリ帯域巾が減少するという形態の性能欠陥を生じさせる。この構成によれば、システムがビクティム書き込みコマンドを実行するたびに、先ず、ディレクトリ状態をアクセスし、次いで、その状態を評価し、そして最終的に、その状態に基づいて、ビクティムデータのＤＲＡＭ書き込みを実行しなければならない。メモリ及びディレクトリは原子的にアクセスされるので、公知の設計方法に基づいてシステムがシステムが設計された場合に、全ビクティム書き込みサイクルは、ディレクトリルックアップ時間と、状態評価時間と、ＤＲＡＭ書き込み時間との和に等しくなる。このようなシステムは、全ビクティムサイクルがＤＲＡＭ書き込みのみで構成されるシステムに対して甚だしい性能上の不利益をこうむる。

【０１２６】

本発明の１つの実施形態は、メモリの各バンクに遅延書き込みバッファを設けることにより、このメモリバンク利用低下問題を克服する。ビクティム書き込みがメモリシステムへ発生されるたびに、メモリシステムは、次の機能を並列に実行することにより応答する。即ち、ビクティム書き込みデータをターゲットメモリバンクの遅延書き込みバッファに記憶しそしてそのブロックを「非書き込み可能」又は「無効」と表示し、ビクティム書き込みに関連したディレクトリ状態をアクセスし、そして現在ビクティム書き込みに代わって、「書き込み可能」又は「有効」と表示された既にバッファされたビクティム書き込みのＤＲＡＭ書き込みを実行する。ディレクトリアクセスが完了したときに、ビクティム書き込みに関連したディレクトリ状態が、ビクティム書き込みが成功したことを示す場合には、ビクティムが存在する遅延書き込みバッファが「書き込み可能」又は「有効」状態へと移行する。遅延書き込みバッファにおけるデータブロックの「書き込み可能」又は「有効」状態は、バッファのデータが、ＤＲＡＭメモリに記憶されたバージョンよりも最新のキャッシュラインのバージョンであることを指示する。バッファが「書き込み可能」又は「有効」と表示された場合には、そのデータが、メモリシステムへのビクティム書き込みのその後に発生によりＤＲＡＭへ書き込まれる。

【０１２７】

既に発生されたビクティム書き込みのＤＲＡＭ書き込みと並列にディレクトリルックアップを実行することにより、この実施形態は、全ビクティムサイクル時間を単一のＤＲＡＭ書き込み時間に減少する。この実施形態は、「書き込み可能」な又は「有効」なデータブ

ロックを多数のサイクルにわたり遅延書き込みバッファに保持し、そのサイクル中にバッファされたブロックへのその後の参照をメモリへ発生することができるので、遅延書き込みバッファは、連想アドレスレジスタを備えている。ビクティム書き込みブロックのアドレスは、その関連データが遅延書き込みバッファに記憶されるのと同時に連想アドレスレジスタに記憶される。その後の参照がメモリシステムへ発生されるときには、メモリシステムは、アドレスレジスタに対するアドレス一致により遅延書き込みバッファにおけるアドレスブロックを識別する。これは、メモリシステムが、D R A Mメモリの効力を失ったデータに代わってバッファからの最新のデータで遅延書き込みバッファのブロックへの全ての参照にサービスすることを意味する。

【 0 1 2 8 】

ビクティムデータの遅延書き込みバッファ動作を与える上記技術は、D T A G状態を直接含まずにD T A G状態を使用してデータブロックの有効性を決定するスヌーピーバスをベースとするシステムにも使用できる。

図 2 1 を参照すれば、遅延書き込み動作を与えるメモリ制御システムの 1 つの実施形態は、ディレクトリ 1 4 0 からライン 1 4 0 a を経て O w n e r _ _ M a t c h 信号を受け取るように接続されたメモリコントローラ 3 3 2 を含むように示されている。更に、メモリコントローラ 3 3 2 は、ディレクトリに入力されるコマンドを追跡するために Q S A R B 1 1 (ディレクトリ 1 4 0 にも信号供給する) からも入力を受け取る。

メモリコントローラ 3 3 2 は、遅延書き込みバッファ 3 3 6 を含む。遅延書き込みバッファ 3 3 6 の各エントリは、データ部分 3 3 6 a と、フラグ部分 3 3 6 b と、アドレス部分 3 3 6 c とを含む。本発明の 1 つの実施形態において、設計上の複雑さを最小限にするために、遅延書き込みバッファは、1 つのアドレス、データ及びフラグエントリのみを保持するが、本発明は、このような構成に限定されるものではない。

【 0 1 2 9 】

遅延書き込みバッファは、次のように動作する。動作中に、コマンド、アドレス及びデータが A R B _ _ B U S 1 3 0 を経て受け取られると、それらはディレクトリ 1 4 0 及びメモリコントローラ 3 3 2 へ送られる。メモリコントローラ 3 3 2 は、コマンド、アドレス及びデータを書き込みバッファ 3 3 6 に 1 トランザクション周期中 (ここでは 1 8 クロックサイクル中) 記憶する。トランザクション周期中に、ディレクトリ 1 4 0 がアクセスされ、そしてアクセスの結果が O W N E R _ _ M A T C H ライン 1 4 0 a にアサートされる。O W N E R _ _ M A T C H ラインは、メモリの更新を求めるプロセッサのプロセッサ I D が実際にデータのキャッシュラインの所有者であることをディレクトリエントリが指示する場合にアサートされる。O W N E R _ _ M A T C H 信号は、遅延書き込みバッファエントリ 3 3 6 のフラグ 3 3 6 b をセットするのに使用される。次に続くトランザクション周期中に、メモリバスが使用できそしてフラグ 3 3 6 b がアサートされた場合には、メモリ 3 3 4 に記憶データが書き込まれる。本発明の 1 つの実施形態では、書き込み動作のみがバッファされ、到来する読み取り動作は、遅延なくメモリバスをアクセスすることが許される。遅延書き込みバッファに記憶されたビクティムデータへのその後の読み取り動作は、遅延書き込みバッファからサービスされる。

【 0 1 3 0 】

図 2 2 は、遅延書き込み動作のタイミング図である。時間 T 0 に、読み取り 0 動作が A R B _ _ B U S に受け取られる。この読み取り動作は、D R A M 3 3 4 をアクセスするためにメモリへ直ちに伝播される。時間 T 1 に、書き込み 1 動作が A R B _ _ B U S に受け取られる。この T 1 サイクル中に、ディレクトリ 1 4 0 がアクセスされ、そして T 1 サイクルの終わりに、書き込み 1 アドレスの一致を示す O W N E R _ _ M A T C H 信号がアサートされる。その結果、遅延書き込みバッファエントリのフラグ 3 3 6 b がセットされる。時間 T 2 に、読み取り 2 動作が受け取られ、書き込み 1 動作の前にメモリへ送られる。時間 T 3 の間に、書き込み 1 動作に対応するフラグがアサートされた場合に、次の書き込み 3 動作が遅延書き込みバッファに受け取られると、書き込み 1 動作が D R A M 3 により処理するためにメモリに送られる。

【0131】

ローカルメモリの読み取りについては、遅延書き込みバッファのフラグビットをセットするのにDTAGも使用できることに注意されたい。ローカルメモリからのキャッシュラインの1つをローカルノードにおけるプロセッサのキャッシュの1つに記憶することができる。プロセッサの1つがキャッシュラインをビクティム化しそしてキャッシュラインが遅延書き込みバッファに書き込まれたときに、そのキャッシュラインのDTAGエントリを検査して、キャッシュラインがプロセッサの1つに常駐したかどうか決定することができる。キャッシュラインがプロセッサの1つに常駐した場合には、DTAGエントリの有効ビットを検査して、プロセッサがビクティム化するコピーが有効であることを確保する。DTAGにヒットがありそしてキャッシュラインが有効であった場合には、DTAGが遅延書き込みバッファのフラグをセットし、キャッシュラインをローカルメモリに書き込みさせる。これは、簡単なスヌーピーバスをベースとする（即ちディレクトリのない）システムがこの同じ簡単なアルゴリズムを適用できるようにする。

10

【0132】

従って、図21のメモリ制御ロジックは、読み取り動作を読み取りサイクルにおいて直ちに実行することができそして書き込み動作を各書き込みサイクルに実行できるようにする（たとえ遅延書き込みであっても）。その結果、ディレクトリのアクセスにより遅延をこうむることなくデータの定常流がDRAMに送られ、そしてコヒレンス性を維持しながら性能が高められる。遅延書き込みバッファ技術は、ビクティム書き込み動作に関連して説明したが、メモリ性能を改善するためにコヒレンス状態が集中され且つ一定保持されるようないかなるシステムにも使用することができる。

20

【0133】

仮想チャンネル：

従って、キャッシュコヒレンスプロトコルを実施するために、プロセッサと、ディレクトリと、メモリと、DTAGとの間に多数のメモリ参照が送信されることが明らかである。更に、各メモリ参照は、多数のトランザクション即ちホップをノード間に備え、メモリ参照のためのメッセージは、参照全体が完了する前に転送される。メッセージ間の依存性が参照を不定に阻止する場合には、マルチプロセッサシステムが停滞（デッドロック）状態となる。上記で簡単に述べたように、本発明の1つの実施形態は、仮想チャンネル流れ制御を使用することにより、ノード間のトラフィックをマネージしそして停滞を生じることなくデータコヒレンス性を維持する。仮想チャンネルは、相互接続ネットワークに停滞のないルートを形成するために最初に導入された。本発明の1つの実施形態によれば、仮想チャンネルは、更に、共用メモリコンピュータシステムのためのキャッシュコヒレンスプロトコルにおけるリソース停滞を防止するのににも使用できる。

30

【0134】

公知の関連するキャッシュコヒレンスプロトコルでは、2つの形式の解決策が使用されている。少数のプロセッサと少数の同時保留中要求とを有するシステムの場合には、実行中の任意の点に生じ得る考えられる最大数の応答を受け入れるに足る大きさの待ち行列及びバッファが設けられている。十分な待ち行列及びバッファスペースを設けることにより、メッセージが進行のために別のメッセージに決して影響されないよう保証している。多数の保留中要求を伴う大型のシステムでは、考えられる最大数の応答を受け入れるに足る大きさのバッファ及び待ち行列を設けることは実際的ではない。従って、停滞検出及び分析機構に接続された2チャンネル相互接続を使用して問題が解決される。第1に、相互接続部（プロセッサ及びメモリのようなシステム要素間にメッセージを移動するのに使用される論理的経路）は、2つのチャンネル、即ち要求チャンネル（又は下位チャンネル）と、応答チャンネル（又は上位チャンネル）とを使用する。これらのチャンネルは、一般に、物理的なものであり、即ち個別のバッファ及び待ち行列を使用する。第2に、潜在的な停滞を検出するために発見的手法が一般的に実施される。例えば、コントローラは、待ち行列がいっぱいでありそして待ち行列からある時間中にメッセージが出力されないときに潜在的な停滞を通知する。第3に、選択されたメッセージが否定的に確認されて、リソ

40

50

ースを解放し、他のメッセージを進行できるようにする停滞分析機構が実施される。否定的な確認メッセージは、それに対応するコマンドをリタイアさせる。

【0135】

上記の大型システムの解決策は、公平さ／欠乏の問題及び性能不利益の問題を含む2つの主たる問題を有している。あるメッセージが否定的に確認されるので、あるコマンドが長時間完了しない（潜在的に不定である）ことが考えられる。コマンドが所与の時間周期内に完了するよう保証されない場合には、そのコマンドを発生するリソースは、システムデータへの公平なアクセスを得ることができない。更に、リソースがシステムデータへの公平なアクセスを得ることができないために、データに対して欠乏状態となり、潜在的にシステムの停滞を生じさせる。更に、あるメッセージが否定的に確認され、従って、それらの行先に到達しないので、無効化メッセージのようなプロトコルメッセージは、それらが行先に首尾良く到達することを指示するための確認を発生しなければならない。更に、コントローラは、対応するコマンドが完了したとみなし得る前に全ての確認が受け取られるまで待機しなければならない。この非決定論的結果は、キャッシュコヒレンスプロトコルの全性能を低減するようなメッセージオーバーヘッド及び余計な待ち時間を生じさせる。

【0136】

本発明の1つの実施形態によれば、停滞回避に対する系統的及び決定論的解決策を採用したキャッシュコヒレンスプロトコルが使用される。潜在的な停滞を検出しそして矯正動作を行うのではなく、停滞が設計により排除される。従って、停滞検出及び分析機構の必要性がなくなる。第2に、メッセージは、停滞回避のための否定的確認ではなくなるので、「無効化」のようなプロトコルメッセージに対する確認が不要となり、それ故、帯域巾及び待ち時間が改善される。

仮想チャンネルの使用を説明する目的で、幾つかの有用な用語について最初に説明する。

依存性：メッセージM2が進行しない限りメッセージM1が進行できない場合に、メッセージM1はメッセージM2に「依存」と定義する。更に、依存性は、移行的であるとも定義する。本発明のキャッシュコヒレンスプロトコルを実施する場合に、リソース依存性及び流れ依存性の少なくとも2種類の依存性がある。M2が待ち行列スロットのようなリソースを解放するまでM1が進行できない場合に、M1はM2に「リソース依存」と定義する。M2が進行するまでM1が進行しないことをキャッシュコヒレンスプロトコルが必要とする場合には、M1はM2に「流れ依存」と定義する。例えば、キャッシュコヒレンスプロトコルは、ディレクトリがある状態に達するまでM1が阻止状態であり、そしてディレクトリの状態を所望の値にセットするのがM2であることを要求する。従って、M1からM2へのリソース又は流れ依存性のチェーンが存在する場合に、M1はM2に依存すると定義する。

【0137】

依存性サイクル：M1の進行がM2の進行に依存し；M2の進行がM3の進行に依存し；Mk-1の進行がMkの進行に依存し；そして最終的に、Mkの進行がM1の進行に依存するときに、1組のメッセージM1、MK（2）の間に「依存性サイクル」が存在すると定義する。メッセージのあるサブセットが依存性サイクルを形成するときにメッセージのシステムは停滞状態になる。M1はMkに依存し、Mkは次いでM1に依存するので、サイクル内のどのメッセージも進行することができない。

ここに開示する方法及び装置は、仮想チャンネルを使用して、キャッシュコヒレンスプロトコルにおける停滞を決定論的に回避する。キャッシュコヒレンスプロトコルの設計において必要とされるハードウェア機構及び従うべき1組のルールについて説明する。

【0138】

1つの実施形態において、キャッシュコヒレンスプロトコルは、全てのメモリ動作がせいぜい3段階で完了すると定める。各段階において、システムの要素間に1つ以上のメッセージが転送される。それ故、各段階は、「ホップ」と称される。ホップは、0、1及び2と番号付けされる。ホップ0では、プロセッサ又はIOPロセッサからの要求がホームディレクトリへ送られる。ホップ1では、ホームディレクトリにより発生されたメッセージ

が1つ以上のプロセッサ又はIOPプロセッサへ送られる。ホップ2では、メッセージが所有者プロセッサからソースプロセッサへ送られる。これらホップは、図23に示されている。

キャッシュコヒレンスプロトコルの顕著な特性は、全ての動作が所定数のホップ内に完了することである。ここに示す実施形態では、所定数が3であるが、本発明は、選択される数が比較的小さく且つ一貫したものである限り、特定のホップ数に限定されるものではない。この特性は、停滞を検出しそして停滞を解消するためのメッセージを失敗して再トライする機構を伴わずに、全てのメッセージをそれらの行先にルート指定できることを保証するための鍵である。

【0139】

上記のように、ここに示す実施形態では、最大ホップ数が3である。従って、システムは、各々Q0、Q1及びQ2と示された3つのチャンネルを備えている。これらのチャンネルは、システム相互接続部を通る論理的に独立したデータ経路である。これらのチャンネルは、物理的なものでもよいし、仮想のもの（或いは一部分物理的で且つ一部分仮想）でもよい。物理的なものであるときには、各チャンネルは、システム全体にわたり個別の待ち行列及びバッファリソースを有する。仮想のものであるときには、チャンネルは、待ち行列及びバッファリソースを共用し、以下に述べる制約及びルールを受ける。

3つのチャンネルは、ハイレアキーを構成し、Q0は最下位であり、Q1はその次であり、そしてQ2は最上位のチャンネルである。システムにおける停滞回避のための重要なルールは、チャンネルQiのメッセージが、Qiより下位のチャンネルのメッセージに決して依存しないことである。

【0140】

更に、本発明の1つの実施形態において、IOPシステムからの応答メッセージと、IOPシステムからのメモリスペースコマンドとの間の流れ依存性サイクルを排除するためにQIOチャンネルが追加される。

最後に、本発明の1つの実施形態では、ビクティムメッセージと、ビクティムメッセージが発生されるがビクティムメッセージが保留中である間に発生されるその後の依存性メッセージとに対して、Q0Vicチャンネルが使用される。

図20a - 20hに関連して上述したように、スイッチへ発生される所与のコマンドメッセージは、一連の多数の個別トランザクションを発生する。本発明の1つの実施形態において、所与のコマンドパケットに対する各個別のトランザクションは、チャンネルに割り当てられる。チャンネルは、本質的に、所与のコマンドパケットの完了段階及び依存性を定義する順序付けされた構造体を形成する。

【0141】

例えば、図23は、図20A - 20Jについて述べた動作の個別トランザクションにチャンネルを割り当てるところを示すフローチャートである。個別トランザクションは、次の用語で識別される。即ち、参照により生じる一連のトランザクションにおける第1トランザクションは、Q0又はQ0Vicトランザクションと称し、一連のトランザクションにおける第2トランザクションは、Q1トランザクションと称し、そして一連のトランザクションにおける第3トランザクションは、Q2トランザクションと称する。

Q0又はQ0Vicチャンネルは、まだディレクトリを訪れていないプロセッサ及びIOPからの初期コマンドを搬送する。従って、Q0/Q0Vicパケットの行先は、常に、ディレクトリである。Q0Vicチャンネルは、「書き込みビクティム」コマンドに対して特に指定され、一方、Q0チャンネルは、プロセッサ又はIOPにより開始された他の全ての形式のコマンドを搬送する。

【0142】

ステップ380で発生されるコマンドは、データを得るか又は状態を更新しようと求める。状態は、常に、データのアドレスに対応するホームディレクトリで得ることができる。ステップ382において、ホームディレクトリがアクセスされ、そして使用可能なキャッシュラインがホームメモリにより所有される（ディレクトリに対して）か、別のプロセッ

10

20

30

40

50

サにより所有されるかが決定される。いずれの場合にも、応答はQ 1チャンネルを経て発生される。ステップ382において、状態又はデータが第2ノードに得られると決定された場合には、ステップ384において、Q 1チャンネルの応答が第1ノードへ返送される。Q 1形式のトランザクションは、ShortFill、ShortFillMod、VicAck、CTD-ACK/NACK当を含む。

【0143】

ステップ382において、ホームノードがデータを所有せず、データがダーティであって別のプロセッサにより所有されると決定された場合には、ステップ386において、「送信読み取り」又は「送信読み取り変更」のQ 1形式のトランザクションがQ 1チャンネルを経てリモートノードへ発生される。

10

ダーティへと状態変化したデータを他のノードが共用することを指示するホームノードの状態チェックに応答するか、又は「読み取り変更」に応答する場合には、ステップ388において、無効化Q 1形式トランザクションがシステムの他の当該ノードに送られる。

従って、Q 1チャンネルは、第2の「ホップ」におけるパケットを搬送するためのものであり、第1のホップはディレクトリに対するものである。第2の「ホップ」の行先は、常にプロセッサであり、プロセッサは、元のコマンドを開始したノードにあるか、又はシステム内の別のリモートノードにある。

【0144】

Q 2チャンネルは、「長い記入」又は「長い記入変更」トランザクションのいずれかを搬送する。Q 2チャンネルは、第3の「ホップ」による第3ノードからのデータを、元のコマンドを開始したノードへ返送する。

20

Q 0 / Q 0 Vic、Q 1及びQ 2形式のコマンドへのコマンドの割り当ては、SMPシステムにおいて停滞のないメッセージ送信を確保するために次のように使用できる。図23のフローチャートは、4つの仮想チャンネル間の対話を示すが、本発明の1つの実施形態では、キャッシュコヒレンスを維持する目的で5つの仮想チャンネルを使用することができる。その追加チャンネルは、QIOチャンネルである。一般に、QIOチャンネルは、制御状態レジスタ(CSR)アクセスを含むIOアドレススペースへ全ての読み取り及び書き込みを搬送する。

【0145】

以下のテーブルIIは、チャンネル経路へのコマンドマッピングを例示するリストである。

30

テーブル I I :

Q10	CPUへの全I Oスペース要求	RdByteIO, RdWordIO, WrByteIO, WrWordIO	10
Q0	CPU又はI O Pからの全メモリ スペース要求	Rd, RdMod, Fetch, CTD, ITD, Vic, RdVic, RdModVic	
QOVic	データを転送するCPU又はI O Pからの全メモリスペース要求	WrVic, Full Cache Line Write, QV_Rd, QV_RdMod, QV_Fetch	
Q1	全送信コマンド	FRd, FrdMod, Ffetch	20
	全シャドーコマンド	SFRd, SFRdMod, SFetch, Sinval, Ssnap	
	短い記入	SFill, SfillMod	
	マーカー記入の全性質	FM, FMMod, Pseudo-FM, Pseudo-DMMod, FrdMod with FM	30
	その他	CTD-ACK, CTD-NACK, ITD-ACK, Vic-ACK, VicRel	
	I Oスペース応答	IOFillMarker, IOWriteACK	
	関連Consig	Invi-Ack, LoopComSig	40
Q2	長い記入	Fill, FillMod	
	I Oスペース記入	IOFill	

スイッチをベースとするシステムにおける仮想チャンネルの1つの実施形態は、各チャンネルに対して物理的に個別の待ち行列、バッファ又は経路を使用することを含む。或いは又、待ち行列、バッファ又はデータ経路は、チャンネル間で共用されてもよく、従って、真の「仮想」であってもよい。本発明の1つの実施形態では、これら技術の組み合わせを使用して、ハードウェアの最適な使用がなされる。

【0146】

図24には、2つ以上の仮想チャンネル間で単一バッファをいかに共用するかが示されている。バッファ400は、多数の「スロット」を含むように示されている。各スロットは、1つのチャンネルのみにより専用使用される。例えば、スロット402は、Q2型コマンドに専用の多数のバッファエントリを含み、スロット404は、Q1型コマンドに専

用の多数のバッファエントリを含み、等々となる。

残りのスロット 4 1 0 は、いずれのチャンネルについても、メッセージにより使用することができ、それ故、「共用」又は「一般的」スロットと称される。各チャンネルについてビジー信号が与えられる。ビジー信号は、バッファがそれ以上のメッセージを記憶できず、それ故、そのバッファに何も送信してはならないことを指示する。

【 0 1 4 7 】

所与のチャンネルに対する所与のリソースにおいてビジー信号がアサートされるときと、そのリソースにコマンドを発生するデバイスがビジー信号に応答して発生を停止するときとの間には待ち時間周期がある。この待ち時間の間に、1 つ以上のコマンドパケットがリ

10

ソースへ発生されることが考えられ、それ故、リソースは、コマンドが脱落しないように設計されねばならない。

それ故、受信器がビジー流れ制御信号をアサートした後にも、M 個のメッセージを受け入れることができねばならず、但し、M は、次の式 III で定められる。

式 III :

$$M = (\text{フレームクロックでの流れ制御待ち時間}) / (\text{フレームクロックでのパケット長さ})$$

「M」の値は、ここでは、チャンネル当たりに得られる専用スロットの数を定義する。

【 0 1 4 8 】

図 2 5 には、各チャンネルごとに個別のリソースを使用して仮想チャンネルが実施される例が示されている。2 つのノード 4 2 0 及び 4 2 4 の部分は、ハイアラキースイッチ (HS) 4 2 2 を経て互いに接続されて示されている。

20

グローバルポート 4 2 0 は、バス 4 2 1 a を経てスイッチ 4 2 2 から入力データを受け取り、そしてバス 4 2 1 b を経てスイッチ 4 2 2 にデータを送信するように接続される。同様に、グローバルポート 4 2 4 は、バス 4 2 3 a を経てスイッチ 4 2 2 にデータを送信し、そしてバス 4 2 3 b を経てスイッチ 4 2 2 からデータを受け取るように接続される。

データバス 4 2 1 a、4 2 1 b、4 2 3 a 及び 4 2 3 b の各々は、全ての形式のチャンネルコマンドを送信又は受信する。待ち行列機構 4 2 5 のような待ち行列機構は、各リソースの各入力及び出力端子に設けられる。この待ち行列機構は、多数の個々に制御されるバッファ 4 2 5 a - 4 2 5 e を備え、各バッファは、1 つの形式のチャンネルコマンドのみを専用に記憶する。バッファ 4 2 5 a は、Q 0 チャンネルコマンドのみを記憶し、バッファ 4 2 5 b は、Q 0 V i c チャンネルコマンドのみを記憶し、等々となる。

30

【 0 1 4 9 】

コマンドパケットが各リソースインターフェイスに受け取られるときに、コマンドの形式がパースされ、そしてパケットは、適当なバッファへ送られる。コマンドパケットがノードの適当なプロセッサ又は I O P へ送られる準備ができると、それらが適当なバッファから選択され、そして A R B バス及び Q S A (図 6) を経て送られる。各チャンネルごとに 1 つずつ、5 つのサーチエンジンがあり、各チャンネルに対して次のメッセージを探索する。

上記機構においては、各チャンネルが独立して流れ制御され、そしてシステム全体にわたりハイアラキーの最下位チャンネル以外の各チャンネルにスロットが指定される。これは、チャンネルがリソース依存性により下位チャンネルによって決して阻止されないことを保証する。上位チャンネルメッセージの移動は、下位チャンネルメッセージによるリソースの占有により阻止されない。

40

【 0 1 5 0 】

仮想チャンネル間で物理的バッファを共用する上記機構は、簡単なものである。より精巧な機構については、ハイアラキースイッチに関して最初に述べた。

仮想チャンネル：裁定及びコヒレンスプロトコル設計のルール

コヒレンスプロトコルにおいて停滞のないメッセージ送信を保証するためにはハードウェア機構のみでは不充分である。というのは、問題のリソース依存性の部分しか対処しないからである。全てのリソース及び流れ依存性サイクルを排除するために、多数の付加的な

50

裁定及びコヒレンスプロトコル設計ルールが適用される。

第1に、メッセージの進行は、下位チャンネルメッセージの進行に依存してはならず、この場合に、Q2は上位チャンネルであり、そしてQ0は下位チャンネルである。アービターは、各チャンネルの流れ制御を互いに独立して維持しなければならない。例えば、ピジーの流れ制御信号がQ1に対してアサートされるが、Q2に対してはアサートされない場合には、アービターは、Q2メッセージを進行させねばならない。保留中のコマンドパケットに対してリソースをサーチするのに使用される全てのサーチエンジンは、同じ特性をサポートしなければならない。

【0151】

第2に、2つ以上のチャンネル間に共用されるいかなるリソースも、下位のチャンネルが阻止された場合に上位のチャンネルが進行できるようにするために、上位のチャンネルの各々に対してある専用のスロットを含まねばならない。

第3に、全てのチャンネルコマンドは、一貫して作用しなければならない。Q0コマンドの終了点は、常に、ディレクトリである。Q1コマンド及びQ2コマンドの終了点は、常に、プロセッサである。終了点において、トランザクションを継続するために、それらを上位チャンネルへ移動しなければならない。例えば、Q0メッセージがディレクトリに到達したときには、Q0メッセージを発生することができず、Q1又はQ2メッセージを発生しなければならない。それ故、メッセージは、下位チャンネルメッセージへと分岐又は変換することはできない。

【0152】

他の点において分岐するトランザクションの場合には、同じか又は上位のチャンネルのメッセージしか形成できない。例えば、「送信読み取り変更」(Q1メッセージ)がハイアラキースイッチにおいて「送信読み取り変更」、「無効化」及び「記入変更マーカー」を形成するときには、これら全てのメッセージがQ1メッセージとなる。

従って、バスをベースとするシステム又はスイッチをベースとするシステムのいずれかに仮想チャンネルを設ける装置及び方法が提供される。仮想チャンネル及び上記の順序付け制約を使用することにより、参照は、ディレクトリによっていったんサービスされると完了することが保証される。その結果、NACK(1つのプロセッサが別のプロセッサにプロセスが完了しないことを指示する)及びリタイアを必要とする公知の複雑なプロトコルは排除される。

【0153】

5つまでの独立したチャンネルを伴う実施形態を示したが、本発明の1つの実施形態は、所与の数のチャンネルに限定されず又は対称的なマルチプロセッサシステムに限定されないことを理解されたい。むしろ、選択されるチャンネルの数は、各チャンネルに固有の制御及びハードウェアオーバーヘッドが与えられると、コヒレントな通信をサポートするのに必要な数でなければならない。従って、仮想チャンネル制御方法及び装置は、マルチプロセッサシステムにおいて高性能の、停滞のない通信を行えるようにする。

コヒレンス性を維持するためのディレクトリの動作：

以上に、基本的な通信構成を説明し、そしてSMPのノード間に通信が自由に流れるようにするための基本的な制御構造体が提供された。しかしながら、コヒレンス性のための鍵は、自由に流れるコマンドがシステム内の各プロセッサにより正しい順序で「取り扱われる」ように確保することである。SMPシステム内の全てのコマンドに対しシリアル化ポイントを与える機構は、各ノードにおけるディレクトリである。

【0154】

上述したように、全てのQ0形式コマンドは、先ず、関連メモリアドレスのホームディレクトリにアクセスする。いずれのコマンドに対してもホームディレクトリが最初にアクセスされるよう確保することにより各コマンドを共通のソースから正しい順序で検討することができる。

本発明の1つの実施形態では、シリアル化順序は、アドレスXに対するディレクトリからの裁定に勝った後にXに対するQ0コマンドがARBバスに現れるという順序である。「

ロード」形式のコマンドは、それに対応する読み取りコマンドがホームディレクトリにアクセスしたときに順序付けされる。「記憶」形式のコマンドは、それに対応する「読み取り変更」コマンドがディレクトリにアクセスするか又はそれに対応する「クリーン - ダーティ」コマンドがディレクトリにアクセスして A R B バスに現れるときに順序付けされる。

【 0 1 5 5 】

例えば、10個のコマンドの以下のシーケンスが種々のプロセッサ (P #) により共通のホームディレクトリへ発生されると仮定する。但し、 X_i は、キャッシュライン X の一部分である。

テーブル IV :

1	P 1 : 記憶 X_1 (1)
2	P 2 : ロード X_1
3	P 3 : ロード X_1
4	P 5 : ロード X_1
5	P 1 : 記憶 X_2 (2)
6	P 2 : 記憶 X_1 (3)
7	P 4 : ロード X_1
8	P 5 : ロード X_2
9	P 6 : ロード X_1
10	P 2 : 記憶 X_1 (4)

キャッシュラインのバージョンは、各記憶動作の結果として更新される。従って、コマンド 1 はバージョン 1 を形成し、コマンド 5 はバージョン 2 を形成し、コマンド 6 はバージョン 3 を形成し、そしてコマンド 10 はバージョン 4 を形成する。

【 0 1 5 6 】

シリアル化順序は、ディレクトリに到達する事象の各シーケンスがキャッシュライン X の正しいバージョンを得るように確保する。例えば、コマンド 2 ないし 4 は、バージョン 1 を得なければならない。プロセッサ P 1 のコマンド 5 が記憶を行うときには、全てのバージョン 1 キャッシュライン (プロセッサ P 2、P 3 及び P 5 における) に「無効化」を送信しなければならない。同様に、プロセッサ P 2 のコマンド 6 がバージョン 3 データで X を更新するときには、プロセッサ P 1 のバージョン 2 データを無効化しなければならない。

プロセッサ P 4、P 6 及び P 7 は、バージョン 3 データを得るが、これは、プロセッサ P 8 のバージョン 4 データの記憶により後で無効化される。共通のアドレスキャッシュライン X に対する多数のロード及び記憶動作は、システムにおいていかなる所与の時間にも進行し得ることを述べれば充分であろう。システムは、ロード及び記憶がディレクトリによりシリアル化順序で処理されるようにこれらのコマンドを処理する。

【 0 1 5 7 】

システムがシリアル化順序を維持しそして付随的にデータのコヒレンス性を維持するのを助けるために多数の技術が使用される。これらの技術は、Q 1 チャンネルコマンドの厳密な順序付け、C T D 明瞭化、「シャドウコマンド」、「マーカー記入」及び「遅延ビクティム書き込みバッファ動作」を含む。各技術について、以下に詳細に説明する。

【 0 1 5 8 】

Q 1 チャンネル順序付け :

コヒレンス性を維持するのに使用される第 1 の方法は、Q 1 チャンネル上を進行する全てのメッセージ、即ちディレクトリから送られる全てのメッセージが、先入れ先出し順序で進むように確保することである。即ち、ディレクトリから別のプロセッサ又は I O P へ送られる Q 1 型メッセージは、コマンドがディレクトリにおいてシリアル化された順序に基づいて送られる。

例えば、図 26 のサブシステムの例では、ノード 430 における第 1 プロセッサ P 1 (431) がキャッシュライン X をそのキャッシュ「ダーティ」に記憶すると仮定する。ノ

10

20

30

40

50

ド 4 3 2 におけるプロセッサ P 1 6 (4 3 3) は、Q 0 チャンネルに「X 読み取り (Read X)」を発生し、これは、ノード 4 3 6 における X のホームディレクトリ 4 3 7 へ送られる。又、ノード 4 3 2 におけるプロセッサ P 1 7 は、Q 0 チャンネルに「無効 - ダーティ」コマンドを発生し、これも、ノード 4 3 6 における X のホームディレクトリ 4 3 7 へ送られる。「X 読み取り」の受信に回答して、ディレクトリエントリに基づき、「送信 X 読み取り (Forwarded Read X)」が Q 1 チャンネルを経てプロセッサ P 1 (4 3 1) へ送られる。I T D の受信に回答して、ディレクトリエントリの状態に基づき、「無効化」がハイアラキースイッチ 4 3 5 へ送られ、これは、Q 1 チャンネルを経てプロセッサ P 1 及びプロセッサ P 1 6 へ「無効化」を送る。従って、同じ時点で、「X 無効化」及び「X 読み取り供給」が Q 1 チャンネルコマンドとして P 1 へ送られる。

10

【 0 1 5 9 】

Q 1 チャンネルのコマンドが順序ずれて実行することが許された場合には、「読み取り」の前に「無効化」が生じることがある。その結果、「読み取り」のための記入データがプロセッサ P 1 6 に送られないことになり、それ移行の動作の結果が予想し得ないものとなる。

しかしながら、チャンネル Q 1 のコマンドを正しい順序で保つことにより「読み取り」は「無効化」を受け取る前に処理され、コヒレンス性が維持される。

【 0 1 6 0 】

本発明の 1 つの実施形態では、チャンネル Q 1 についてのみ F I F O 順序が維持され、F I F O 順序とは、同じメモリアドレスに対応する全てのメッセージが F I F O 順序に留まることを意味する。しかしながら、本発明は、Q 1 チャンネルに対する順序を維持することのみに限定されるものではなく、チャンネルのいかなる組み合わせに対する順序の維持も含むように拡張することができる。

20

上記の順序付け手順を実施する 1 つの方法は、Q S A チップ (図 6) の Q S A R B 1 1 により実行される。Q S A R B は、全ての Q 0 トランザクションをノードのホームメモリスペースに対してシリアル化する。その結果、Q 1 パケットのシリアル流が発生されて、ノードのローカルプロセッサと、グローバルポート及びハイアラキースイッチを経てノードから離れたプロセッサとの両方に向けられる。

第 1 の順序付けルールを次に説明する。所与の Q S A R B により発生される全ての Q 1 パケットは、シリアルな順序で発生される。所与の Q S A R B からの幾つかの又は全ての Q 1 パケットがターゲットとする全てのプロセッサは、これらの Q 1 パケットを、それらが Q S A R B により発生された順序で見る。

30

【 0 1 6 1 】

このルールをサポートするために、Q S A チップは、ノード内の接続されたプロセッサとやり取りされる全ての Q 1 パケットに順序を維持する。グローバルポートのロジックは、ハイアラキースイッチと Q S A チップとの間に転送される全てのパケットに F I F O 順序を維持する。更に、ハイアラキースイッチは、所与の入力から所与の出力へ送られる全ての Q 1 パケットにも順序を維持する。

このルールは、1 つの Q S A R B からの Q 1 パケットと、別のノードの Q S A R B からの Q 1 パケットとの間に特定の順序を命令するものではないことに注意されたい。他のノードから受け取られた Q 1 パケットは、ハイアラキースイッチを経てホームノードにより発生された Q 1 パケットと次のようにシリアル化される。リモートノードのプロセッサをターゲットとする全ての Q 1 パケットは、リモートノードの Q S A R B により処理される。これらの Q 1 パケットは、ハイアラキースイッチによりリモートノードで発生された Q 1 パケットとシリアル化される。所与の Q S A R B からの Q 1 パケットの全ての受信者は、Q 1 パケットを、それらが Q S A R B においてシリアル化されたのと同じ順序で見なければならない。

40

【 0 1 6 2 】

図 2 7 は、多数の Q 0 及び Q 1 コマンドの順序付けが上記の順序付けガイドラインに基づいて S M P を通して処理されるところを示すブロック図である。ノード 4 4 0 のプロセッ

50

サ P x はコマンド Q 0 a を発生し、プロセッサ P y はコマンド Q 0 b を発生し、そしてプロセッサ P z はコマンド Q 0 c を発生すると仮定する。同じ時間中に、Q S A R B 4 4 1 は、プロセッサ P r 及び P q からの Q 1 メッセージをグローバルポート 4 4 3 から受け取る。

これらのメッセージは、次のように順序付けされる。Q S A R B 4 4 1 は、Q 0 a、Q 0 b 及び Q 0 c を処理して、Q 1 a、Q 1 b 及び Q 1 c 応答を発生する。これらの発生された Q 1 コマンドは、到来する Q 1 コマンドと合成されて、コマンドの順序付けされた流れを F I F O 4 4 2 へ供給し、ローカルプロセッサへと送る。F I F O コマンドの順序は、Q S A R B により処理されたコマンドの順序を反映する。

Q 1 a、Q 1 b 及び Q 1 c コマンドは、グローバルポート 4 4 3 へ送られ、リモートノードへ送信される。グローバルポートの出力バッファ 4 4 4 は、これらのコマンドを、それらが Q S A R B により処理されたのと同じ順序で記憶する。この順序は、図 1 4 - 1 9 について上述した方法を用いてメッセージがリモート C P U 4 5 4 へ送られるときにハイアラキースイッチ 4 4 6 により維持される。

【 0 1 6 3 】

図 2 7 A は、ハイアラキースイッチにおいて従う別の順序付けガイドラインを示す。上述したように、ハイアラキースイッチは、ハイアラキースイッチの所与の入力ポートに現れてハイアラキースイッチの共通の出力ポートをターゲットとする多数のパケットが、それらが入力ポートに現れたのと同じ順序で出力ポートに現れるよう確保することにより、順序を維持する。

【 0 1 6 4 】

図 2 7 B を参照すれば、上述したように、ハイアラキースイッチは、入力メッセージをマルチキャストする役目も果たし、即ち受け取った 1 つの Q 1 パケットを 2 つ以上の行先ノードに送信するという役目も果たす。スイッチによりマルチキャストされるパケットの一例は、無効化パケットである。ハイアラキースイッチの異なるポートから入力された多数のパケットが共通の出力ポートにマルチキャストされるときには、Q 1 パケットは、全ての出力ポートにおいて同じ順序で現れねばならない。例えば、パケット 1 及びパケット 2 の両方がハイアラキースイッチ 4 6 0 に受け取られる場合に、2 つのメッセージをプロセッサ 4 6 4 及び 4 6 6 にマルチキャストする 1 つの許された方法は、上記のように、メッセージ 2 がメッセージ 1 の前に両プロセッサに到着するようにすることである。別の許された方法は、メッセージ 1 のパケットがメッセージ 2 のパケットの前に両プロセッサに到着するようにすることである。しかしながら、2 つのプロセッサは、2 つのパケットを異なる順序で受け取ってはならない。

【 0 1 6 5 】

ハイアラキースイッチが従わねばならない別の順序付けルールは、多数の入力ポートからの Q 1 パケットの順序付けされたリストが共通の出力ポートをターゲットとするときに、Q 1 パケットが、全ての到来する Q 1 パケットの 1 つの共通の順序付けに合致する仕方

で出力ポートに現れるように確保することである。
例えば、図 2 7 C において、入力ポート 4 6 1 には、パケット 2 がパケット 4 の前に受け取られる。同様に、入力ポート 4 6 2 には、パケット 1 がパケット 3 の前に受け取られる。停滞を防止するには、これら命令の全体的な順序を遵守しなければならない。出力パケットを与える 1 つの許された方法は、パケット 3 を最初にノード 4 6 4 に送信し、そしてパケット 1 を最初にノード 4 6 6 に送信することである。この送信が図 2 7 C に示されている。別の許された出力は、パケット 2 及び 4 を受信者のプロセッサにより最初に受け取ることである。しかしながら、1 つのプロセッサがパケット 3 を最初に受け取りそして別のプロセッサがパケット 4 を最初に受け取る場合には、プロセッサがそれらの元のシーケンスの他のパケットの受信を待機してストールするので停滞が生じ得る。

それ故、Q 1 チャンネルにおいて順序が維持されるよう確保するルールが設けられる。本発明の 1 つの実施形態では、性能の理由で、Q 0 及び Q 2 チャンネルパケットを順序ずれて処理するのが望ましい。データの一貫性を確保するために、多数のコヒレンス性機構

10

20

30

40

50

が以下に述べるように設けられる。

【 0 1 6 6 】

ダーティへの変更の明瞭化

上述したように、Q 1 形式のコマンドのみが、ディレクトリに定義されたシリアル化順序で維持される。本発明の 1 つの実施形態では、Q 0 及び Q 2 コマンドは順序付けされない。従って、受け取られる Q 0 及び Q 2 コマンドの相対的なタイミングの結果としてディレクトリにコヒレンス性の問題が生じないように予防策がとられる。

発生する 1 つのコヒレンス性の問題は、ディレクトリエントリの構造によるものである。図 9 に示すように、各ディレクトリエントリは、所有権フィールドと、各ノードに対して 1 つの存在ビットを含む。存在ビットは、関連ノードの 4 つのプロセッサの 1 つにデータが存在することを示すおおよそのベクトルである。4 つのプロセッサのいずれかが動作すると、存在ビットがセットされる。従って、ノードのどのプロセッサが存在ビットをセットしたかに関してある種の曖昧さが生じる。この曖昧さは、ある場合にコヒレンス性の問題を引き起こす。

【 0 1 6 7 】

例えば、図 2 8 A 及び 2 8 B は、2 つのノード 4 7 0 及び 4 7 2 のブロック図である。ノード 4 7 0 [グローバルシステムのノード I D 3] は、プロセッサ P 1 2、P 1 3、P 1 4 及び P 1 5 を備え、一方、ノード 4 7 2 [グローバルシステムのノード I D 7] は、ノード P 2 8、P 2 9、P 3 0 及び P 3 1 を含む。

【 0 1 6 8 】

時間 T 0 - T 3 の種々の一連の周期における所与のキャッシュライン X のディレクトリエントリの状態は、図 2 8 B においてディレクトリ状態テーブル 4 5 5 に示されている。この例では、キャッシュライン X のホームノードは、ノード 4 7 0 又は 4 7 2 以外のノードである。

時間 T 0 において、キャッシュライン X の所有者は、所有者 I D 8 0 で示すようにメモリである。更に、時間 T 0 において、ノード I D 7 のプロセッサ 3 0 は、キャッシュライン X のクリーンなコピーを記憶する。

時間 T 1 において、プロセッサ 1 4 は、「記憶」コマンドを送信し、これは、「読み取りブロック変更 X」に変換され、そしてキャッシュライン X のホームディレクトリへ送られる。メモリが所有者であるから、プロセッサ P 1 4 は、メモリからデータを得ることができ、そしてキャッシュラインの所有者となる。キャッシュライン X の古いバージョンを無効化するためにノード 7 に無効化が送信され、そしてノード 7 の存在ビットがクリアされる。更に、プロセッサ P 1 4 は、そのノード存在ビット 4 5 6 (ビット 3) をセットする。キャッシュライン X は、変更及び記憶のためにホームメモリからプロセッサ P 1 4 へ送られる。

【 0 1 6 9 】

時間 T 2 に、プロセッサ 3 1 のような別のプロセッサが、キャッシュライン X の「読み取り」を発生する。この「読み取り」は、プロセッサ P 1 4 から「記入」を経てデータを得る。従って、時間 T 2 に、ディレクトリは、ノード I D 3 (プロセッサ P 1 4) 及びノード I D 7 (プロセッサ P 3 1) の両方が、ノード存在ビット 4 5 8 及び 4 5 6 で示すように、キャッシュライン X のコピーを記憶することを指示する。

時間 T 3 に、プロセッサ 3 0 により C T D が発生される場合には、システムの異なるプロセッサから見たキャッシュライン X の状態は、次の理由でインコヒレントとなる。C T D がディレクトリに到達すると、X のディレクトリエントリを読み取り、そしてそのノード、即ちノード I D 7 の存在ビット 4 5 8 が既にオンであるかどうか決定する。その結果、プロセッサ 3 0 は、次いで、C T D 要求において成功したと仮定する。プロセッサ 3 0 は、キャッシュライン X のプロセッサ 1 4 のコピーを無効化し、そしてディレクトリの所有者フィールドを更新する。この動作は、予想し得ない結果を招くことがある。というのは、プロセッサ P 1 4 がプロセッサ P 3 0 よりも最新のデータバージョンを記憶するからである。

10

20

30

40

50

【0170】

1つの問題は、プロセッサ30がプロセッサ14により形成されたキャッシュラインの古いバージョンをまだ記憶しており、そしてプロセッサ14がデータの最新のバージョンを無効化するように通知したことである。このような状態は、SMPシステムで重大なコヒレンスの問題を生じさせる。

上記問題を解消するのに使用できる幾つかの方法がある。その1つの方法は、システムの各プロセッサごとに1ビットを与えるようにディレクトリエントリの存在ビット拡張することである。従って、分解能がノードレベルからプロセッサレベルへ変更される。しかしながら、この解決策は、不都合なことに、ディレクトリのサイズを増大する。

【0171】

本発明の1つの実施形態は、同じアドレスへの保留中参照がそのノードに対してトランシット状態にあるときにCTDコマンドを低速化することにより上記曖昧さの問題を防止するより簡単な方法を提供する。同じアドレスに対して保留中の要求がある場合には、その以前の要求がリタイアするまでCTDが保持される。所与のノードのトランザクション追跡テーブル(TTT)(図10)を使用して、そのノードに対する保留中のグローバル参照を監視する。更に、CTDがTTTに受け取られた後に受け取った要求は、失敗となる。

図10を参照して述べたように、TTTは、完全に連想式の多機能制御構造体である。TTTは、2つの一般的なタスクを実行する。これは、その関連ノードにより発生された全てのリモート参照のアドレスを記憶する。従って、TTTは、そのトランザクションが完了したとみなされるまで、ノードにより発生された各リモートアクセスに対して1つの情報エントリを記憶する。更に、TTTは、ローカルアドレスの要求に応答して、過渡的なコヒレンス状態に関してコヒレンス情報を与える。従って、TTTは、アクセスがトランシット状態にある間にその状態を追跡するためのテーブルである。

他の処理システムは、いかなる瞬間にも所与のキャッシュラインへの1つの参照をトランシット状態にすることができる。トランシット状態にあるキャッシュラインへのその後の参照は、トランシット状態の参照が完了するまで阻止される。

【0172】

これに対し、ディレクトリにおけるコマンドのシリアル化と、チャンネル順序付けルールとにより、本発明のSMPは、同じキャッシュラインへの多数の参照を所与の瞬間に進行させることができる。その結果、SMPの全性能が改善される。

TTT522は、QSAチップ535のロジックにより、グローバルポートに発生されたトランザクションの状態を決定するのに使用される。グローバルポートへ応答を発生する前に、QSAは、まず、TTTにアクセスして、同じキャッシュラインへのどんな参照が保留中であるかを決定する。参照は、最後に受け取ったトランザクションに응答してTTTからリタイアしていない場合には保留中である。

【0173】

参照がTTTからいかにリタイアするかは、コマンドフィールド584に示された参照の形式に依存する。例えば、TTTに記憶するためにグローバルポートへ送られる「X読み取り」参照は、「ここに記入」588a及び「マーカーをここに記入」588bの両方の状態ビットを受け取ることを必要とする。(「マーカーの記入」は、いかに詳細に述べる。)CTD又はITDのような状態型の参照の場合に、TTTにおいてACK/NACKビット588cをセットすれば、そのエントリをリタイアするのに充分である。

【0174】

図29は、TTTを使用して曖昧なディレクトリエントリを排除するところを示すフローチャートである。ステップ500において、キャッシュラインXは、そのホームノードのメモリに記憶され、そしてノード7のプロセッサ30は、データのコピーを記憶する。ステップ502において、「ReadModX」がプロセッサ14により発生される。その結果、無効化がノード7に送られる。ステップ504において、プロセッサP31は、「ReadX」を発生し、これは、ノード7のTTTのエントリを次の状態で形成する。

10

20

30

40

50

【 0 1 7 5 】

アドレス	コマンドID	状態	Fill	Fmark	Shadow	ACK/NACK
X	Read 31		-	-	-	

ステップ506において、プロセッサP30は、CTDXを発生する。QSAチップは、CTD命令のアドレスを検査し、それがリモートCTDであることを決定し、そしてTTTへのGPリンクを経てグローバルポートへ送信する。TTTの内容は、以下に示す通りである。

10

アドレス	コマンドID	状態	Fill	Fmark	Shadow	ACK/NACK
X	Read 30		-	-	-	
X	Read 31		-	-	-	

図6について述べたように、グローバルポートは、TTTからの情報を使用して、どのコマンドをハイアラキースイッチから送出することが許されたかを決定する。本発明の1つの実施形態では、保留中の「読み取り」がトランシット状態にあるとTTTが決定した場合に、グローバルポートは、「読み取り」結果が返送されるまでCTDをスイッチへ送ることが防止される。

20

【 0 1 7 6 】

図29のフローチャートに示す例では、アドレスXへの保留中の読み取り要求は、TTTにより識別される。その結果、ステップ508において、CTDは、「読み取り」がもはや保留中ではなくなるまで、オフに保たれる。

「読み取り」は、「記入」及び「マーカー記入」の両方がノード7に返送されるまで保留となる。この時間中に、ステップ502においてReadModにより発生された無効化がノード7に到達し、各ノードのDTAGSを更新する。Xの無効化がTTTに到達すると、TTTは、TTTに保持されたCTDを失敗と表示し、これは直ちに解除される。ステップ510において、CTDが依然TTTにある場合には、グローバルポートを経て送信される。

30

従って、TTTを使用して、CTDコマンドを適当にオフに保持し又は失敗状態とすることにより、ディレクトリの存在ビットの曖昧さにより生じるコヒレンスの問題を排除することができる。

【 0 1 7 7 】

マーカーの記入：

プロセッサに対するほとんどの応答は、Q1チャンネルにおけるものであり、従って、上記のルールによれば、正しい順序が維持される。しかしながら、Q2チャンネルで受け取られたメッセージは、この順序制約を受けない。Q2型のメッセージは、「記入」及び「記入変更」を含む。

40

Q2型メッセージの到着は、ディレクトリにおいて明らかなように、シリアル化順序を表わさないで、返送データに潜在的な曖昧さが生じる。例えば、「無効化」がQ1を進行し、そして「記入変更」がQ2を進行するので、コヒレンスを維持するためにどの動作が順序において最初に生じるべきかを決定する何らかの方法が必要ではない。

【 0 1 7 8 】

例えば、図30を参照すれば、2つのノード520及び532が示されている。説明上必要なノードの部分しか示されていない。プロセッサP2(524)及びプロセッサP4(534)がキャッシュラインXのコピーを記憶すると仮定する。キャッシュラインXのホームノードは、ノード532である。

以下の説明において、次のパケットにより使用されるチャンネルは、異なる線を用いて指

50

示される。Q 0 コマンドは、単一線矢印で指示され、Q 1 コマンドは、二重線矢印で指示され、そしてQ 2 コマンドは、破線矢印で指示される。

プロセッサP 4 がキャッシュラインXの独占的所有権を得るためにC T D Xを発生すると仮定する。これに回答して、ディレクトリ存在ビット及びD T A G (図示せず)により、ディレクトリ5 4 2 は、ノード5 2 0 へ無効化を発生する。この無効化は、Q 1 チャンネルを経てノード5 2 0 のD T A G Sを更新し、そしてコピーを有する全てのプロセッサ(ここではプロセッサP 2)に無効化調査を送信する。

【0 1 7 9】

次いで、プロセッサP 1 は、Xのホームディレクトリ5 4 2 へR e a d M o d Xを発生する。上記のように、Xは、現在プロセッサP 4 により所有され、それ故、コヒレンスプロ
10
コルによれば、F o r w a r d e d R e a d M o d XがプロセッサP 4 へ送られる。プロセッサP 4 は、それに回答して、Q 2 チャンネルを経てプロセッサP 1 へF i l l M o d
dを発生する。

Q 2 チャンネルの通信は、Q 1 の通信とシリアル化されないので、C T D Xからの「無効化」がノード5 2 0 に到達する前にQ 2 のF i l l M o dがプロセッサP 1 に到達する可能性が存在する。その結果、P 1 のキャッシュには有効データが書きこまれるが、そのすぐ後で、D T A G SがノードにおけるXのコピーを無効化するようにセットされ、そしてP 2 及びP 1 に「無効化」が送られる。しかしながら、「無効化」は、P 2 のバージョンのみに対応し、P 1 におけるバージョンには対応しない。ここで、システムは、インコ
20
ヒレントな状態となる。ディレクトリ5 4 4 は、P 1 を所有者として記録するが、P 1 はまだ無効化されている。

【0 1 8 0】

本発明の1つの実施形態は、各ノードのグローバルポートに「マーカー記入」及びトランザクション追跡テーブル(図1 0)を使用することによりこの問題を克服する。

「マーカー記入(Fill Marker)」又は「マーカー記入変更(Fill Marker Mod)」は、ホームノードのメモリに現在記憶されていないデータに対する「読み取り」又は「読み取り変更」要求に回答して発生されるパケットである。即ち、「マーカー記入」又は「マーカー記入変更」は、「送信読み取り(Forwarded Read)」又は「送信読み取り変更(Forwarded Read Mod)」と同時に発生される。従って、「マーカー記入」及び「マーカー記入変更」は、
30
Q 1 チャンネルコマンドである。「送信読み取り」又は「送信読み取り変更」コマンドは、キャッシュラインを記憶するプロセッサに送られるが、「マーカー記入」又は「マーカー記入変更」の行先は、元の「読み取り」又は「読み取り変更」を供給したプロセッサである。

【0 1 8 1】

「マーカー記入」は、発生元プロセッサが、ディレクトリに生じるシリアル化順序を決定できるようにする。図3 1を参照すれば、「マーカー記入」の適用は、上記問題を次のように矯正する。前記したように、プロセッサ5 3 4 がXのC T DをXのホームディレクトリに発生し、その結果、「無効化」5 5 0 がQ 1 チャンネルを経てノード5 2 0 へ送られると仮定する。プロセッサP 1 (5 2 2)がR e a d M o d Xをリモートディレクトリに発生すると、その要求に対してT T Tエントリが発生される。この要求に対するT T
40
Tテーブルエントリの例が図3 2 に示されている。T T Tテーブルエントリは、「ここに記入」及び「マーカーをここに記入」状態ビットを含むことに注意されたい。これらビットの各々は、ノード5 2 0 のグローバルポートに各パケットが受け取られるのに応答してセットされる。T T Tエントリは、「記入」及び「マーカー記入」の両方が返送されるまでクリアされない。

【0 1 8 2】

図3 1に戻ると、上述したように、プロセッサ5 2 2 からのR e a d M o d Xは、プロセッサ5 3 4 へのF R d M o d Xを生じる。同時に、チャンネルQ 1 を経て、F i l l
M a r k e r M o d X 5 5 2 がプロセッサP 1 に返送される。「無効化」及びF i l l
1 M o d M a r k e rの両方が同じQ 1 チャンネルに送られる。チャンネルQ 2 のF i
50

11 Mod 554は、「無効化」の前にノード520に到着すると仮定する。グローバル参照の「タグ複製」状態は、Fill Mod又はFill Mod Markerの返送に応答して更新される。従って、Fill Modは、Xの所有権をプロセッサP1として表わすようにXのDTAG状態を更新させる。

【0183】

「無効化」550が、ノード520に到達する次の命令であると仮定する。TTTは、「送信読み取り」命令の状態を決定するためにアクセスされる。この点において、TTTエントリは、「ここに記入」ビットをセットするが、「マーカーをここに記入」ビットはセットされない。従って、TTTは、無効化及びリモート読み取り動作の相対的なタイミングに関する指示を与える。Q1コマンドのシリアル化のために、無効化は、プロセッサ522からのRdModXよりも早い時間にディレクトリ542に発生されたと推測でき、従って、Fill Modが新しいバージョンであり、プロセッサ522のデータコピーには無効化が適用されない。その結果、プロセッサP1のDTAGエントリは、無効化されない。

【0184】

上記実施形態は、TTTをグローバルポートに存在するものとして示したが、別の実施形態によれば、各ノードの各プロセッサは、ディレクトリへの要求を監視することにより共通のアドレスへのリモート要求の状態を追跡することができる。従って、「マーカー記入」は、単にTTTへ送られるのではなく、ディレクトリにより関連プロセッサへ送られる。

従って、TTTは、2つの目的を果たすことが明らかである。マルチプロセッサノードから送出されたコマンドの形式を監視することにより、TTTは、同じアドレスへの他のコマンドが完了するまで、あるコマンド(CTDのような)の送信を禁止することができる。更に、要求がQ2チャンネル(「マーカー記入」のような)へ移行したときにTTTに指示する表示機構を設けることにより、TTTを用いて、異なるチャンネルに返送されるコマンド(即ち、Q2記入及びQ1コマンド)間の相対的なタイミング指示を与えることができ、従って、メモリを崩壊することのあるコマンドがプロセッサへ送られるのを防止することができる。

【0185】

シャドーコマンド

上記説明から明らかなように、ローカルアクセスは、通常、リモートアクセスよりも相当に早い。従って、性能に関しては、ローカル及びリモートアクセスの両方がSMPシステムにおいて同時に生じることが許される。

しかしながら、ローカルアクセスの発生によりリモートアクセスに対して停滞の問題を生じさせる幾つの場合がある。例えば、図33Aを参照すれば、1つのプロセッサ562がキャッシュラインXにRdXを発生すると仮定する。キャッシュラインXのホームノードは、ノード560である。ノード560のディレクトリは、プロセッサ582がキャッシュラインを現在所有することを指示する。従って、Forwarded Read Xが582に送られる。

その後、ノード560のプロセッサ564がCTD Xを発生すると仮定する。上記のように、キャッシュラインXは、ノード560に対してローカルであり、CTDが成功すると、「無効化」をプロセッサP1に(及び図示のようにプロセッサP5にも)送る。

【0186】

図33Bを簡単に参照すれば、参考としてここに取り上げる本発明と同日に出願されたバンドレン氏等の「分散型データ依存性ストール機構(Distributed Data Dependency Stall Mechanism)」と題する特許出願に開示されたように、プロセッサP1のような各プロセッサは、同じキャッシュ位置に対する保留中の読み取りがある場合にキャッシュへの調査をストールするためのロジックを備えている。上記の例が与えられると、Read Xの作用は、ミスアドレスファイル(MAF)574にアドレスXを記憶することである。MAFの内容は、到来する調査に対して比較され、そして到来する調査とMAFのアドレス

10

20

30

40

50

間に一致があるときに、調査待ち行列がストールされる。

【 0 1 8 7 】

「記入」データがプロセッサ 5 8 2 から返送されるときに調査待ち行列が解除される。しかしながら、同じ形式のトランザクション（即ち、P 5 がリモート R d Y を実行し、次いで、P 6 が C T D Y を発生する）がノード 5 8 0 に生じる場合に、プロセッサ P 5 の調査待ち行列がストールされ、R e a d Y 要求が満足されるのを保留する。

P 2 により発生された「無効化」の後に P 5 から F o r w a r d e d R e a d Y が送られる状態で P 1 調査待ち行列がストールされるのと同時に、P 6 により発生された「無効化」の後にプロセッサ P 1 から F o r w a r d e d R e a d X が送られる状態で P 5 の調査待ち行列がストールされた場合には、停滞が生じる。

10

【 0 1 8 8 】

この停滞問題を防止するための多数の解決策が存在する。第 1 に、全ての参照をリモートとすることができ、即ち全ての参照を（ホームノードからの参照も）、それらがホームノードに送られる前にスイッチに送ることができる。全ての参照がリモートにされた場合には、上述した中央の順序付けルールに基づき、停滞状態は生じない。第 2 の解決策は、キャッシュラインへのいずれかの参照がリモートから送られたときに所与のキャッシュラインへの全ての参照をストールすることである。しかしながら、この解決策は、これまでのローカル動作の性能に著しく影響し、それ故、好ましいものではない。

【 0 1 8 9 】

本発明の 1 つの実施形態は、コマンドシャドー作用の使用によるローカル及びリモート参照の混合により課せられる潜在的な停滞を克服する。キャッシュライン X へのローカル参照がリモートプロセッサへ送られると、そのキャッシュラインへのその後の全ての参照がハイアラークスイッチへリモートから送られ、キャッシュラインのローカル参照及びその後の全ての参照が完了するまで、中央で順序付けされる。従って、まだシャドー状態であるキャッシュラインへの以前の参照は、キャッシュラインへの現在の参照もシャドー状態にする。図 3 4 及び 3 5 を参照して、上記例をシャドーコマンドの使用と共に説明する。図 3 5 は、T T T の内容を例示している。第 1 プロセッサ P 1 は、R d X をアービターに発生する。上述したように、プロセッサ P 5 への F R d X を生じ、これは T T T に記録される。その後、プロセッサ P 2 は、C T D X を A R B に発生する。A R B は、T T T を検査し、リモートプロセッサへ送られる保留中のローカル読み取りがあると決定し、そしてグローバルポートからプロセッサ P 5 へ I n v a l X を送る。又、この動作を表わすエントリも T T T に形成され、そのシャドービットがセットされる。

20

30

【 0 1 9 0 】

同時に、ノード 5 8 0 において、同様の一連のトランザクションが生じる。プロセッサ P 5 は、R d Y を発生し、これはノード 5 6 0 に送られると共に、P 5 アドレスをエントリに含ませることにより T T T に記録される。プロセッサ P 6 は、その後、C T D Y を発生する。ノード 5 8 0 のアービターは、C T D アドレスを T T T 内の保留中読み取りに対して一致させ、そして C T D Y をグローバルポートにわたり「シャドー」処理する。その C T D Y に対して T T T にエントリが形成され、このエントリは、T T T においてそのシャドービットをセットし、C T D Y が、Y への要求の適切な順序付けを確保するためにリモート送信されたローカル参照であることを指示する。

40

上述したように、両ノードにおいて調査シーケンスで「無効化」の後に F R d があるときに問題が生じる。「無効化」は、ここでは中央で順序付けされるので、両無効化を両方の「送信読み取り」の前にそれらの調査待ち行列へ送信できないことにはならない。というのは、それらは、共通点即ちハイアラークスイッチにおいてシリアル化されるからである。従って、図 3 6 を参照すれば、コマンドの入力シーケンスは、ハイアラークスイッチ 5 6 8 へ入力されるように示されている。許容し得る出力シリアル化順序は、順序 a - f として識別される。上記の Q 1 チャンネル順序付けルールによれば、ハイアラークスイッチへのパケット入力 of シリアル化順序がスイッチ出力に維持されることに注意されたい。それ故、上記の場合には、F R d は、行先ノードへ送られるときにその関連する「無

50

効化」に先行する。

【0191】

ノードの1つは、調査待ち行列に「無効化」を受け取り、その後、「送信読み取り」を受け取る。例えば、シリアル化順序を用いて、プロセッサP5の調査待ち行列は、Invalid Yによりストールされ、そしてFrd Xがストールされて、記入を保留する。しかしながら、この例では、Frd Yは、Invalid Xの後ではなく、従って、P5調査待ち行列を阻止しないように「記入」データを与えることができる。

リモート参照のためにデータが返送されるときには、その参照に対応するTTTエントリがドロップされる。元の参照をシャドー処理した他の参照がTTTに存在することがある。これらコマンドがハイアラキースイッチから受け取られるときには、シャドー処理されたコマンドの各々に対するTTTエントリもドロップされる。最終的に、リモートアクセス及びシャドーアクセスが全て完了し、そしてTTTがもはやキャッシュラインへマップするエントリを含まなくなると、そのキャッシュラインへのその後のローカル参照をシャドー処理する必要がなくなる。

【0192】

従って、シャドーコマンドの使用により、ローカル及びリモートコマンドの共存から生じるリソース依存性の停滞を、ハードウェアの複雑さを著しく増加せずに排除することができる。上記の例は、「送信読み取り」及びCTDの使用を含むが、シャドーコマンド方法は、他の形式の命令及びマルチプロセッサにも等しく適用できることに注意されたい。一般に、ローカルアドレスXへの参照が存在し、そしてローカルアドレスXへの以前のメッセージがリモートプロセッサ(TTTにより指示された)へ送られるか、又はXへの以前の参照がまだシャドー処理されるときには、Xへの現在の参照もシャドー処理される。更に、この方法は、上記の単なるマルチプロセッサ/スイッチハイアラキーよりも多数のハイアラキーレベルを含む他の形式のアーキテクチャーにも使用できる。例えば、上記方法は、多数のハイアラキーレベルを含み、コマンドがキャッシュラインへの以前の保留中参照のハイアラキーレベルに基づいて適当なハイアラキーレベルに送られるコンピュータシステムにも使用できる。

【0193】

従って、大型のSMPコンピュータシステムに使用するためのアーキテクチャ及びコヒレンスプロトコルについて説明した。SMPシステムのアーキテクチャは、多数のマルチプロセッサノードをスイッチに接続して最適な性能で動作することのできるハイアラキースイッチ構造体を備えている。各マルチプロセッサノード内には、マルチプロセッサノードの全てのプロセッサを最高の性能で動作できるようにする同時バッファシステムが設けられる。メモリはノード間で共用され、マルチプロセッサノードの各々にメモリの一部分が常駐する。

マルチプロセッサノードの各々は、メモリコヒレンス性を維持するための多数の要素、即ちビクティムキャッシュ、ディレクトリ及びトランザクション追跡テーブルを含む。ビクティムキャッシュは、リモートのマルチプロセッサノードに記憶されたメモリを行先とするビクティムデータを選択的に更新することができ、これにより、メモリの全性能が改善される。ディレクトリに関連して使用されて、メモリに書きこまれるべきビクティムを識別する遅延書き込みバッファを各メモリに含ませることによりメモリ性能が更に改善される。

【0194】

各ノードのディレクトリの出力に接続されたARBバスは、SMPを経て転送される全てのメッセージに対して中央の順序付けポイントとなる。本発明の1つの実施形態によれば、メッセージは、多数のトランザクションを含み、各トランザクションは、メッセージの処理段階に基づいて多数の異なる仮想チャンネルに指定される。従って、仮想チャンネルの使用は、システム順序を維持する簡単な方法を与えることによりデータのコヒレンス性を維持する上で助けとなる。仮想チャンネル及びディレクトリ構造体を使用すると、従来停滞を生じるキャッシュコヒレンス性の問題を回避することができる。

以上、本発明の好ましい実施形態を説明したが、その概念を組み込んだ他の実施形態も使用できることが当業者に明らかであろう。それ故、本発明は、上記の実施形態に限定されるものではなく、特許請求の範囲のみによって限定されるものとする。

【図面の簡単な説明】

【図 1 A】公知の対称的なマルチプロセッサコンピュータシステムのブロック図である。

【図 1 B】公知の対称的なマルチプロセッサコンピュータシステムのブロック図である。

【図 2】スイッチを備えた本発明によるマルチプロセッサコンピュータノードの 1 つの一実施形態を示すブロック図である。

【図 3】多数の同時挿入バッファを備えた図 1 のスイッチのデータ経路を示すブロック図である。

10

【図 4 A】図 3 に示す同時挿入バッファの 1 つの一実施形態を示すブロック図である。

【図 4 B】図 4 に示す同時挿入バッファの 1 つを制御するためのロジックの一実施形態を示すブロック図である。

【図 5】図 3 に示す同時挿入バッファの 1 つの別の実施形態を示すブロック図である。

【図 6】同様のノードの大きなネットワークへ接続するように拡張された図 2 のマルチプロセッサコンピュータノードのブロック図である。

【図 7 A】図 6 のマルチプロセッサノードと同様の多数のノードを使用して実施された SMP システムの一実施形態を示す図である。

【図 7 B】図 6 のマルチプロセッサノードと同様の多数のノードを使用して実施された SMP システムの別の実施形態を示す図である。

20

【図 8】図 6 のグローバルポートのブロック図である。

【図 9】図 6 のマルチプロセッサノードのディレクトリにおけるエントリーを示す図である。

【図 10】図 8 のグローバルポートに使用するためのトランザクション追跡テーブル (TTT) を示す図である。

【図 11】図 7 A において多数のノードを接続するためのハイアラキー式スイッチを示すブロック図である。

【図 12 A】停滞を排除するハイアラキー式スイッチ用の相互接続ロジックの一実施形態を示すブロック図である。

【図 12 B】図 12 A の相互接続ロジックの動作を示すフローチャートである。

30

【図 13】マルチプロセッサノードの 1 つから送信されるデータを停止する流れ制御を与えるために図 12 A の相互接続ロジックに使用される方法を示すフローチャートである。

【図 14】ハイアラキー式スイッチに対してバスを経て行われるアドレス及びデータパケットの転送を示すタイミングである。

【図 15】ハイアラキー式スイッチにおいて順序を維持するためのバッファロジックの一実施形態を示すブロック図である。

【図 16】ハイアラキー式スイッチに対して順序を維持するためのバッファロジックの別の実施形態を示すブロック図である。

【図 16 A】チャンネルの依存性を矢印で示す図である。

【図 17】図 16 のバッファロジックを動作する 1 つの方法を示すフローチャートである。

40

【図 18】ハイアラキー式スイッチにおいて順序を維持するためのバッファロジックの別の実施形態を示すブロック図である。

【図 18 A】チャンネルの依存性を矢印で示す図である。

【図 19】図 7 A 又は 7 B の SMP に使用するためのプロセッサ命令 - ネットワーク命令の変換を示すテーブルである。

【図 20 A】図 7 A 又は 7 B の SMP においてノード間にパケットを転送するための多数の通信流を示す図である。

【図 20 B】図 7 A 又は 7 B の SMP においてノード間にパケットを転送するための多数の通信流を示す図である。

50

【図 20C】図 7A 又は 7B の SPM においてノード間にパケットを転送するための多数の通信流を示す図である。

【図 20D】図 7A 又は 7B の SPM においてノード間にパケットを転送するための多数の通信流を示す図である。

【図 20E】図 7A 又は 7B の SPM においてノード間にパケットを転送するための多数の通信流を示す図である。

【図 20F】図 7A 又は 7B の SPM においてノード間にパケットを転送するための多数の通信流を示す図である。

【図 20G】図 7A 又は 7B の SPM においてノード間にパケットを転送するための多数の通信流を示す図である。

10

【図 20H】図 7A 又は 7B の SPM においてノード間にパケットを転送するための多数の通信流を示す図である。

【図 20I】図 7A 又は 7B の SPM においてノード間にパケットを転送するための多数の通信流を示す図である。

【図 20J】図 7A 又は 7B の SPM においてノード間にパケットを転送するための多数の通信流を示す図である。

【図 21】図 2 又は 6 のマルチプロセッサシステムに使用するためのメモリモジュールのレイアウトを示すブロック図である。。

【図 22】遅延書き込み動作のために図 21 のメモリモジュールにより使用される制御ロジックを示すタイミング図である。

20

【図 23】本発明の 1 つの実施形態においてキャッシュコヒレンス性を維持するためにチャンネルに対してマップされる個別のトランザクションの使用を示すフローチャートである。

【図 24】図 7A 又は 7B の SMP において仮想チャンネルを取り扱うための共用待ち行列構造体の一実施形態を示すブロック図である。

【図 25】図 7A 又は 7B の SMP のノード及びハイアラーキーチャンネルにおける個々のチャンネルバッファの一実施形態を示すブロック図である。

【図 26】仮想チャンネル間にある程度の順序が維持されない場合に生じる問題を説明するためのブロック図である。

【図 27A】図 7A 又は 7B の SMP においてコヒレントな通信を与えるための Q1 チャンネルにおける流れ及び順序付けの制約を示すブロック図である。

30

【図 27B】図 7A 又は 7B の SMP においてコヒレントな通信を与えるための Q1 チャンネルにおける流れ及び順序付けの制約を示すブロック図である。

【図 27C】図 7A 又は 7B の SMP においてコヒレントな通信を与えるための Q1 チャンネルにおける流れ及び順序付けの制約を示すブロック図である。

【図 28A】図 7A 及び 7B の SMP のディレクトリエントリーにおおよそのベクトル存在ビットがあるために生じる曖昧さの問題を説明するブロック図である。

【図 28B】図 7A 及び 7B の SMP のディレクトリエントリーにおおよそのベクトル存在ビットがあるために生じる曖昧さの問題を説明するブロック図である。

【図 29】図 28 に示す問題の結果として生じるデータの曖昧さを防止するために使用される方法を示すブロック図である。

40

【図 30】異なるチャンネルのパケットが順序づれて受け取られるために生じるコヒレンス性の問題を示すブロック図である。

【図 31】図 29 に示すコヒレンス性の問題を防止するための記入マーカーの使用を示すブロック図である。

【図 32】図 31 について述べた流れ間の命令の状態を表わす TTT のエントリを示す図である。

【図 33A】SMP システムにおけるダーティへの変更コマンドの作用を示すブロック図である。

【図 33B】SMP システムにおけるダーティへの変更コマンドの作用を示すブロック図

50

である。

【図 3 4】図 3 3 について述べた問題を矯正するためのシャドーコマンドの使用を示すブロック図である。

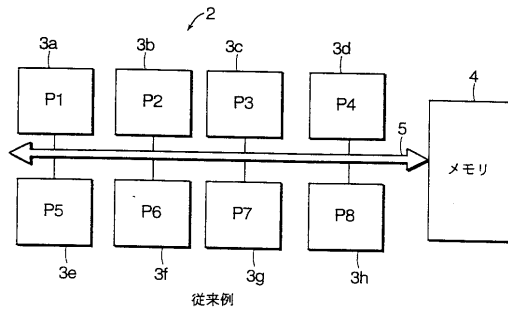
【図 3 5】図 3 4 について述べたフロー間の命令の状態を反映する T T T のエントリを示す図である。

【図 3 6】図 3 5 に示す例における許容し得る逐次順序付けを示すフローチャートである。

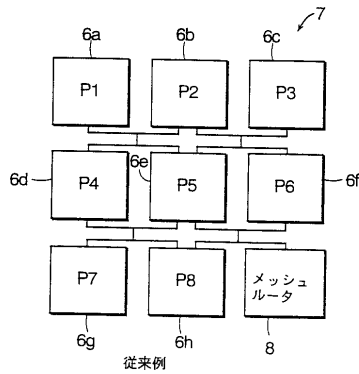
【符号の説明】

1 0	マルチプロセッサノード	
1 1	アービター (Q S A A R B)	10
1 2 a、1 2 b、1 2 c、1 2 d	プロセッサモジュール	
1 3	メモリ	
1 3 a - 1 3 d	メモリモジュール	
1 4	I / O プロセッサ (I O P) モジュール	
1 4 a	I / O バス	
1 4 b	I O P タグ記憶装置	
1 4 c	I O P キャッシュ	
1 5	ローカルスイッチ	
1 6 a - 1 6 i	データリンク	
1 7	A R B バス	20
1 8	Q S A チップ	
1 9	Q S D チップ	
2 0	デュープリケートタグ記憶装置 (D T A G)	
2 5 a - 2 5 e	同時挿入バッファ (S I B)	
2 7	メインアービター	
3 2	バッファ	
3 4 a - 3 4 h	マルチプレクサ	
3 6	入力アービター	
3 8	出力アービター	
1 2 2	トランザクション追跡テーブル	30
1 2 4	ビクティムキャッシュ	
1 4 0	ディレクトリ	

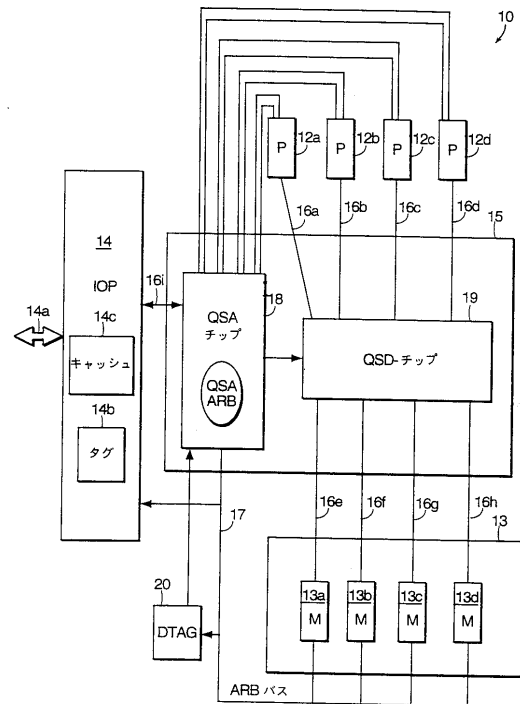
【 図 1 A 】



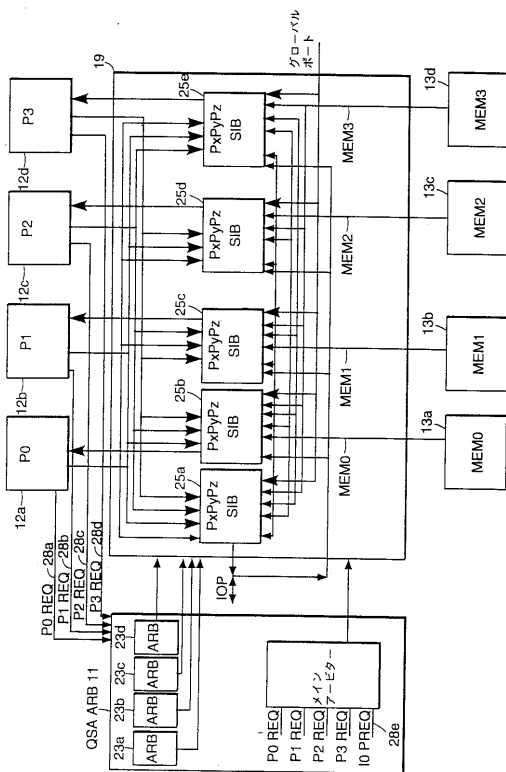
【 図 1 B 】



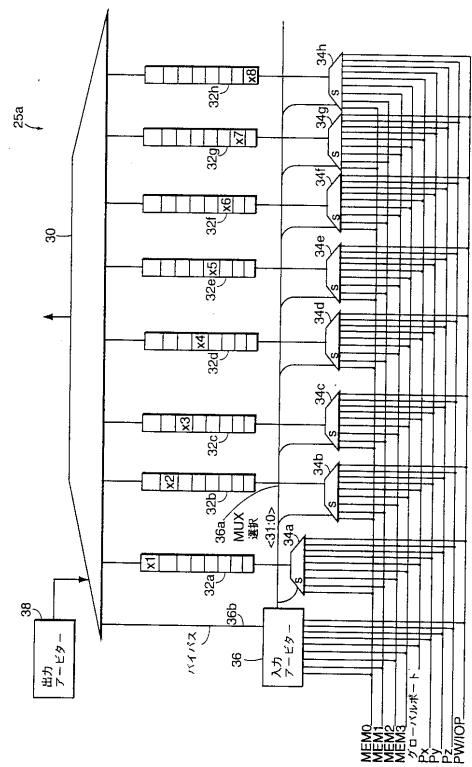
【 図 2 】



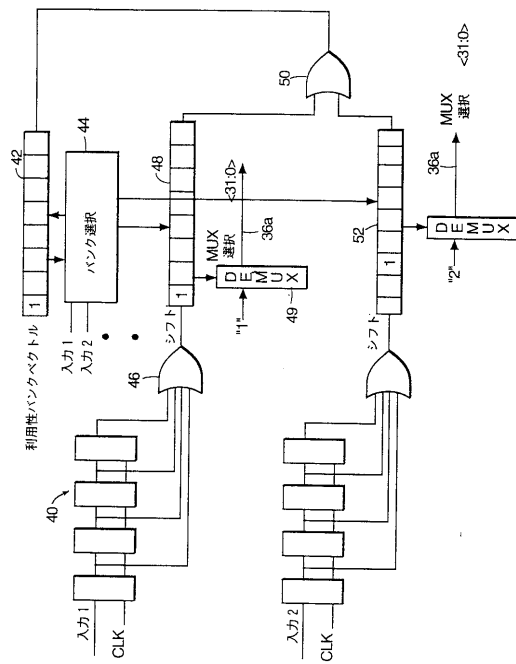
【 図 3 】



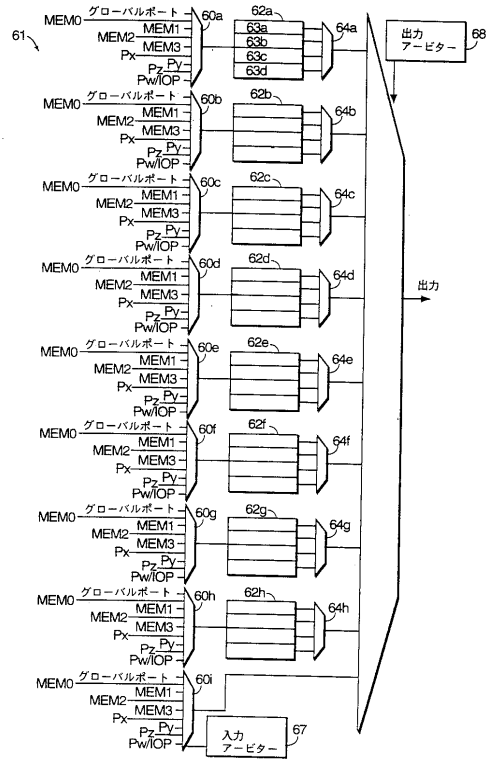
【 図 4 A 】



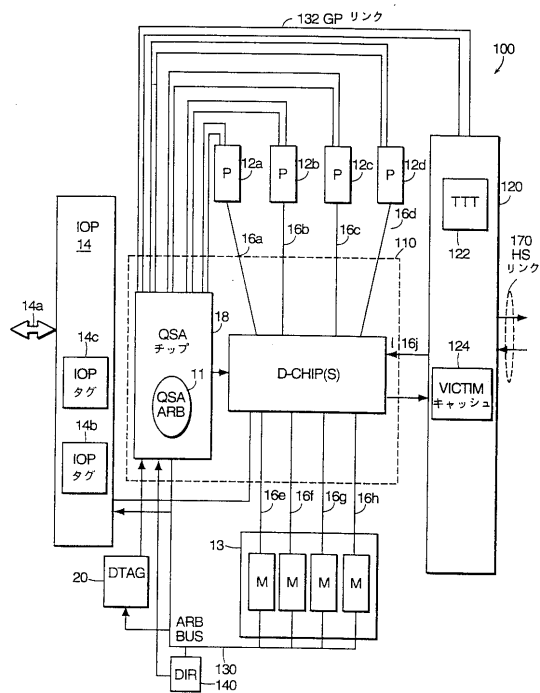
【 図 4 B 】



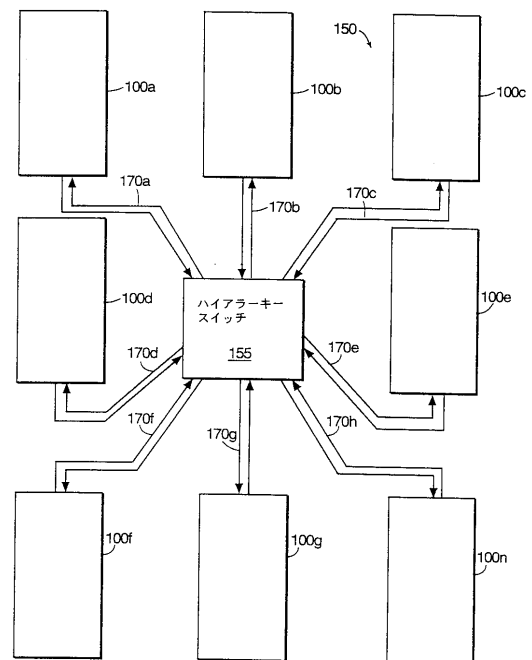
【 図 5 】



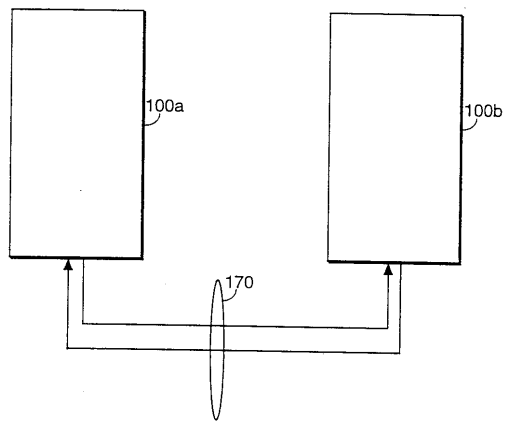
【 図 6 】



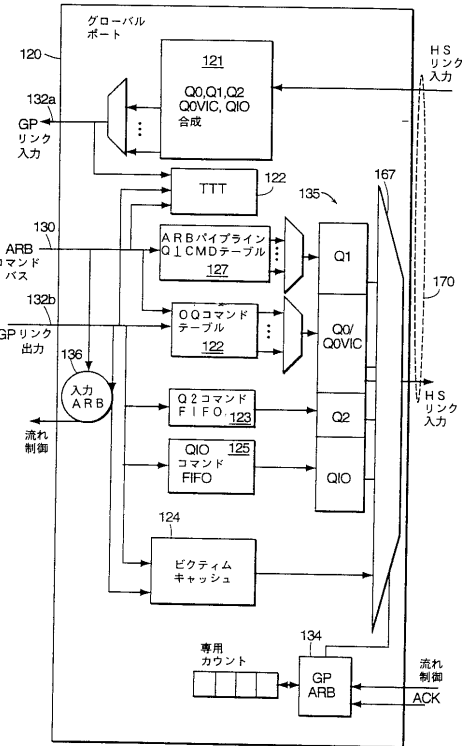
【 図 7 A 】



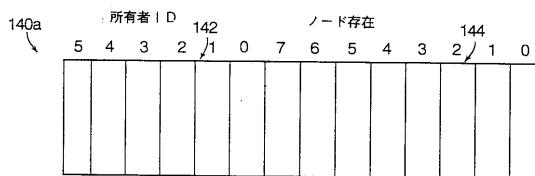
【 図 7 B 】



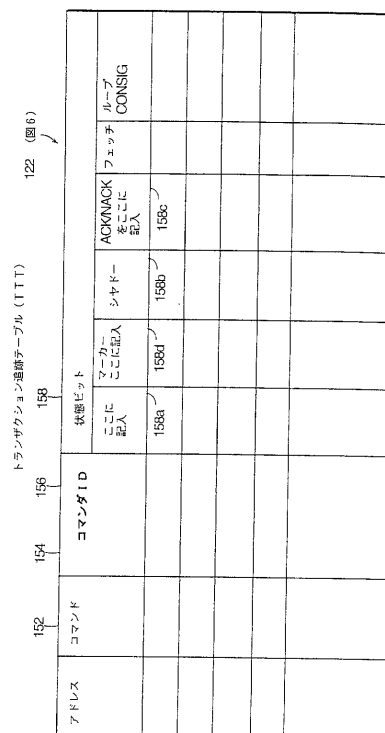
【 図 8 】



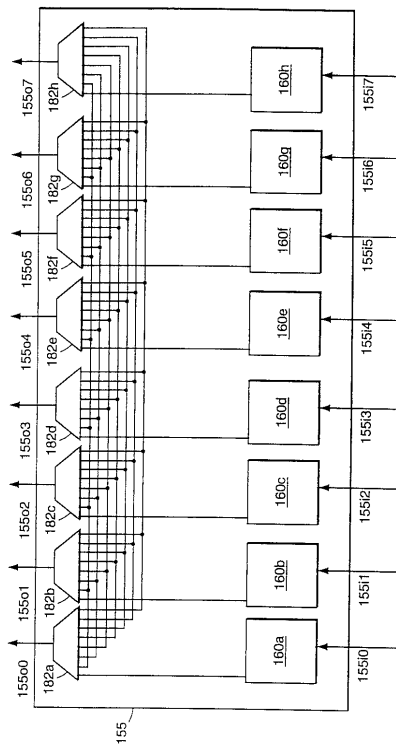
【 図 9 】



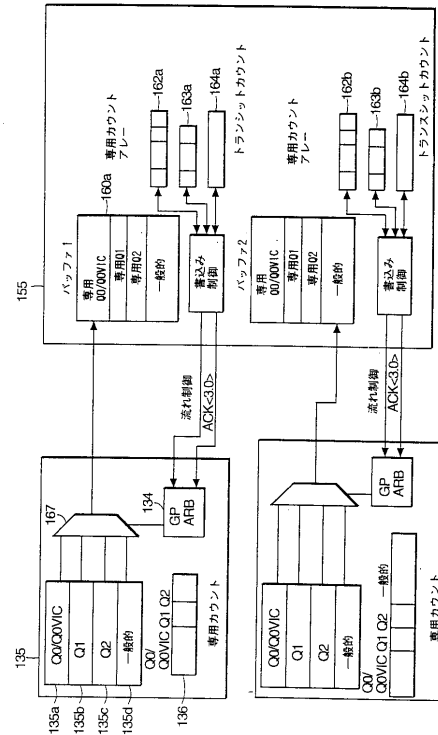
【 図 1 0 】



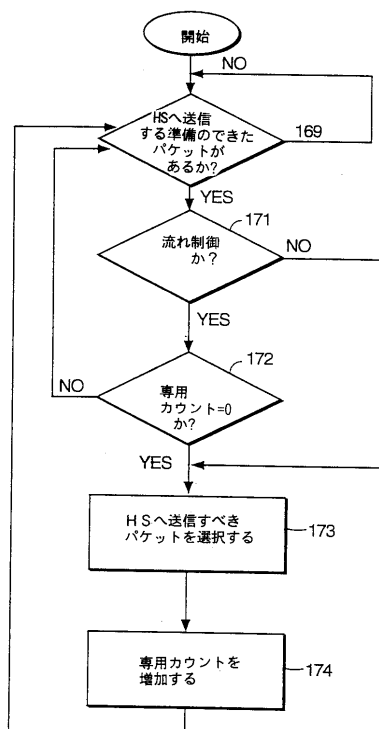
【図 1 1】



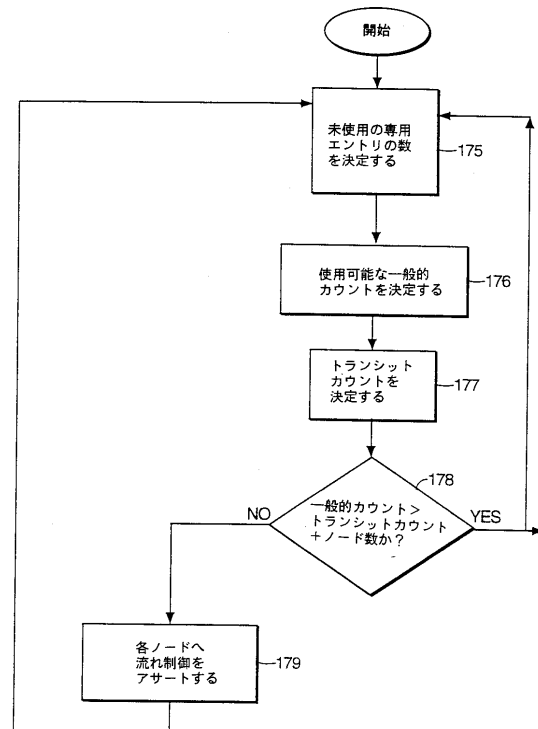
【図 1 2 A】



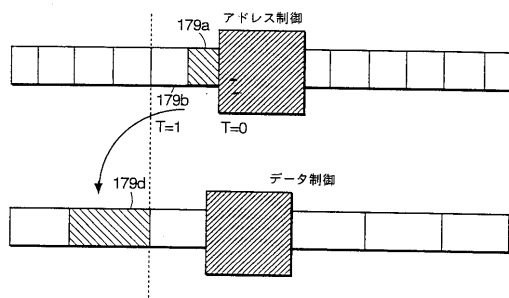
【図 1 2 B】



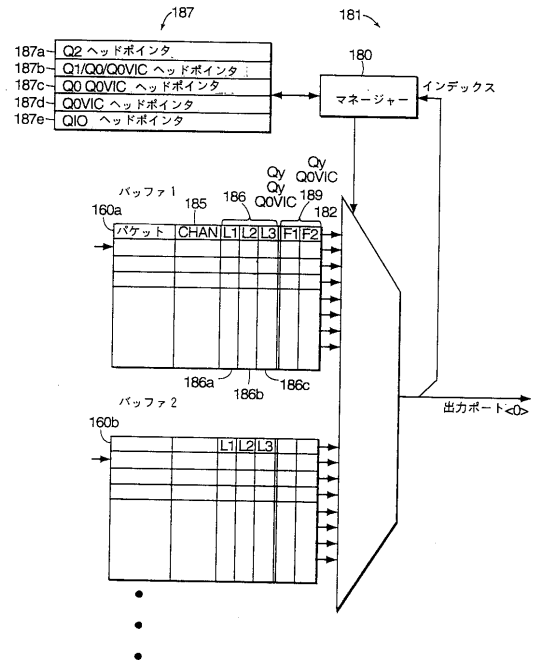
【図 1 3】



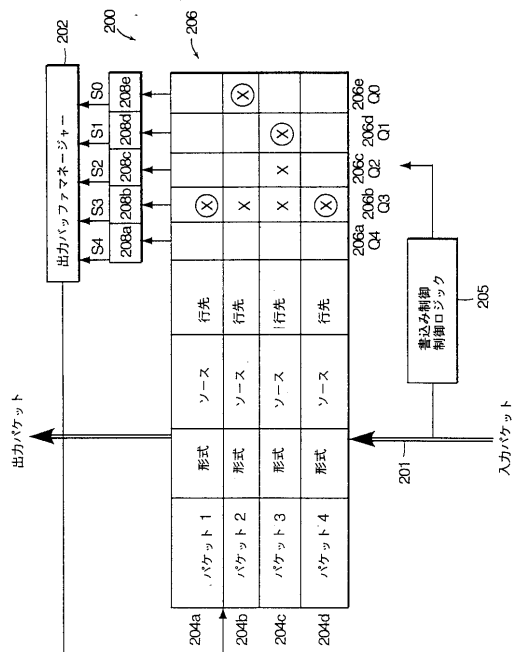
【図 14】



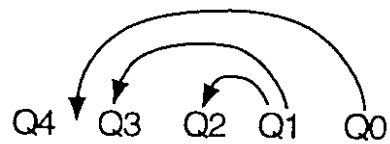
【図 15】



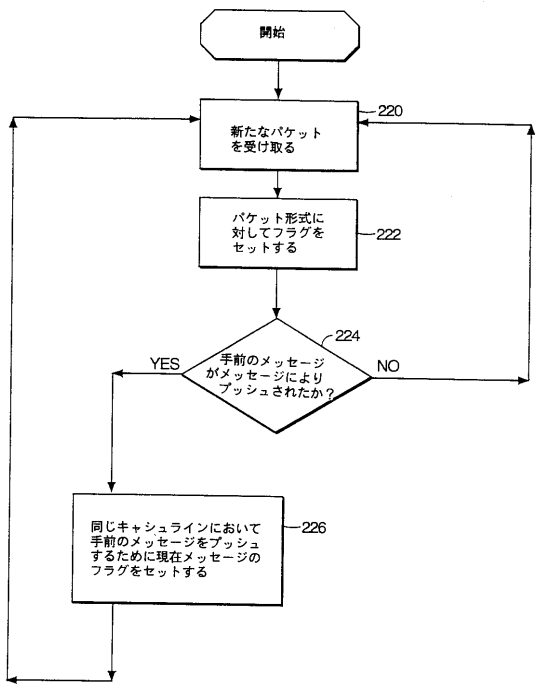
【図 16】



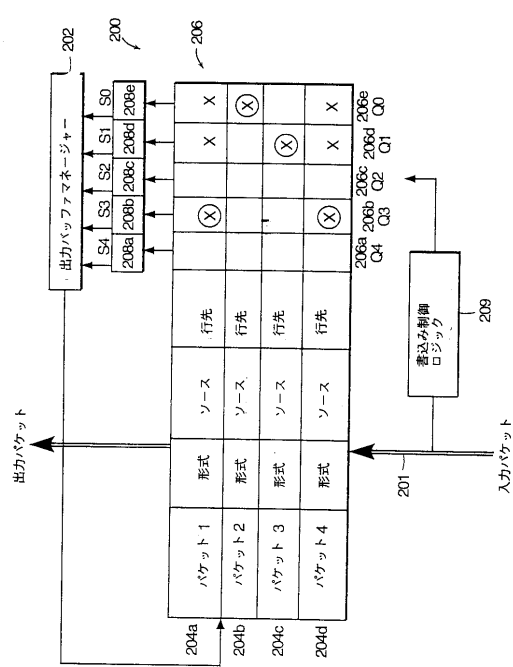
【図 16 A】



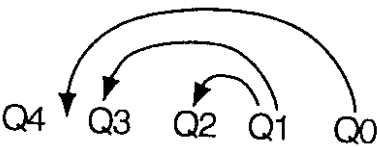
【図 17】



【図 18】



【図 18 A】

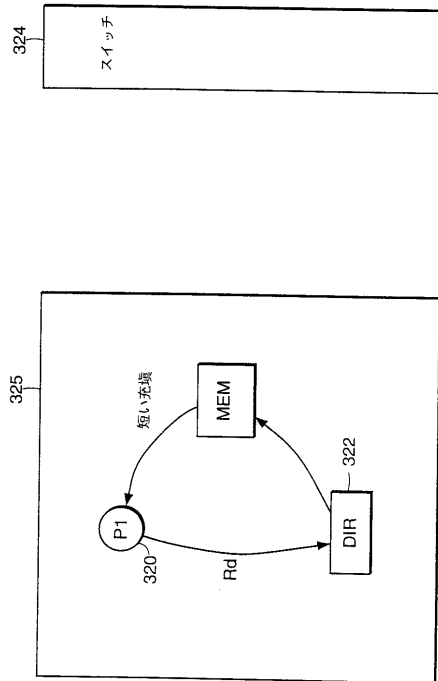


【図 19】

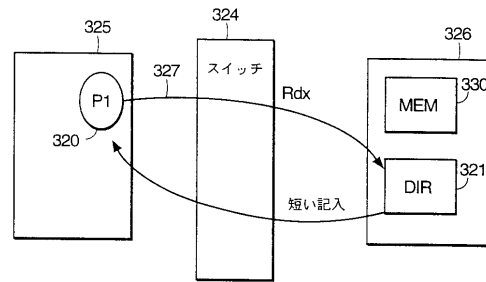
22/43

300	305	320	300a	305a	320a
ブロック読み取り	無効		クリーン	無効	
	クリーン		ダーティ	クリーン	無効化
	ダーティ	送信		ダーティ	N/A
	ダーティ共用	送信読み取り		ダーティ共用	無効化
ブロック読み取り変更	無効		共用ダーティ	無効	
	クリーン	無効化	クリーン	クリーン	無効化
	ダーティ	送信読み取り変更	ダーティ	ダーティ	N/A
	ダーティ共用	送信読み取り変更	ダーティ共用	ダーティ共用	N/A
フェッチ	無効		ダーティへのSTC変更	クリーン	無効化
	クリーン			ダーティ	N/A
	ダーティ	送信読み取り		ダーティ共用	無効化
	ダーティ共用	送信読み取り		ダーティ共用	無効化
ブロック読み取りビクティム	無効		無効ダーティ	無効	
	クリーン		クリーン	クリーン	無効化
	ダーティ	送信読み取り	ダーティ	ダーティ	無効化
	ダーティ共用	送信読み取り	ダーティ共用	ダーティ共用	無効化
ブロック読み取り変更ビクティム	無効		全ブロック書き込み	無効	
	クリーン	無効化	クリーン	クリーン	無効化
	ダーティ	送信読み取り変更	ダーティ	ダーティ	無効化
	ダーティ共用	送信読み取り変更	ダーティ共用	ダーティ共用	無効化
ブロックフェッチ	無効				
	クリーン				
	ダーティ	送信読み取り			
	ダーティ共用	送信読み取り			
ビクティムクリーン	任意の状態				
ビクティム	任意の状態				

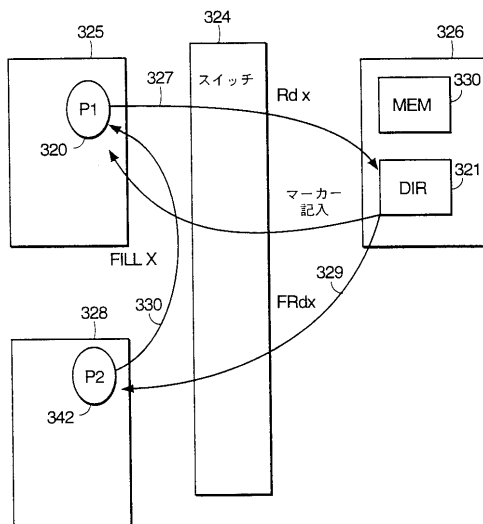
【図 20 A】



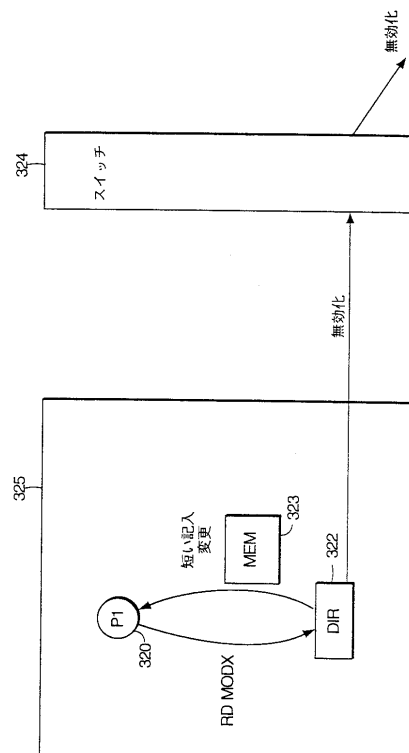
【図 20 B】



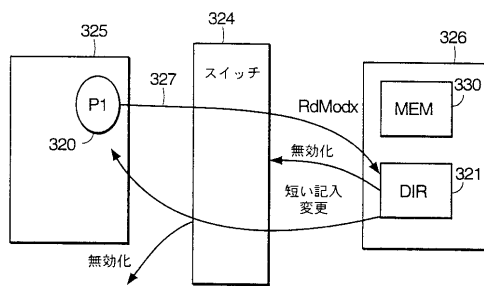
【図 20 C】



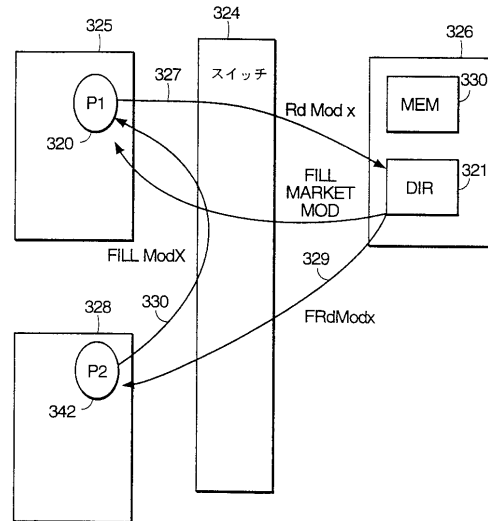
【図 20 D】



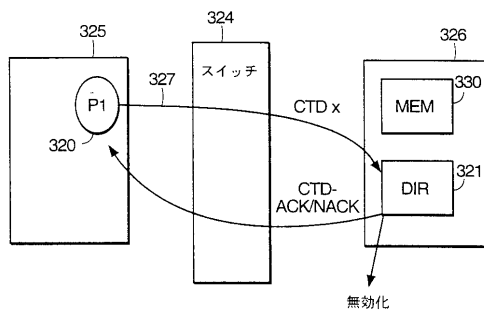
【図 20 E】



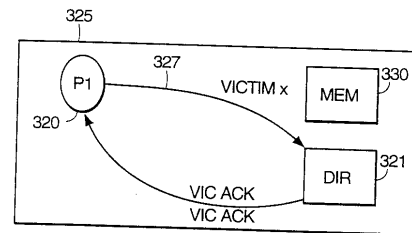
【図 20 F】



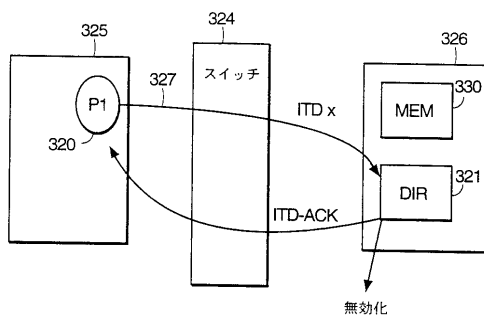
【図 20 G】



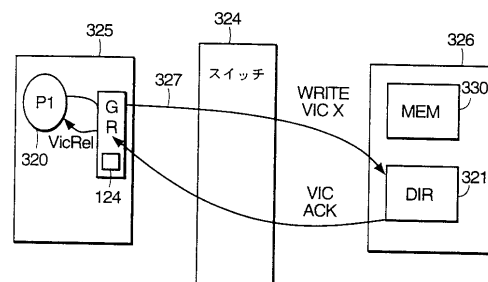
【図 20 I】



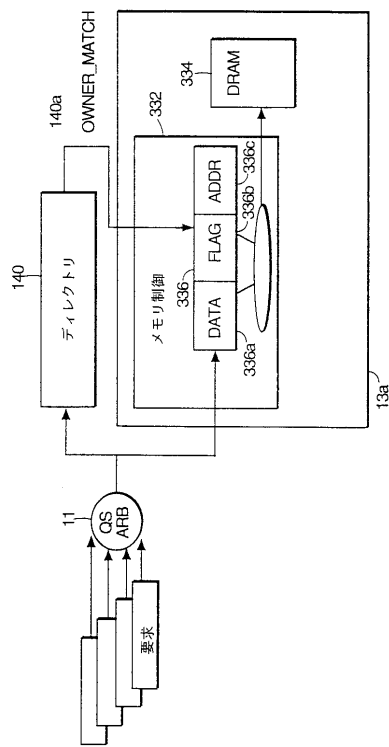
【図 20 H】



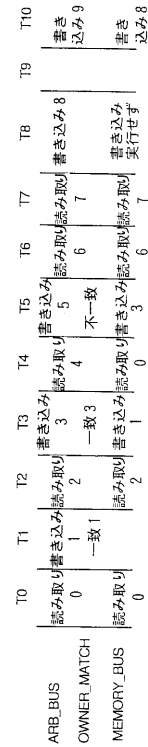
【図 20 J】



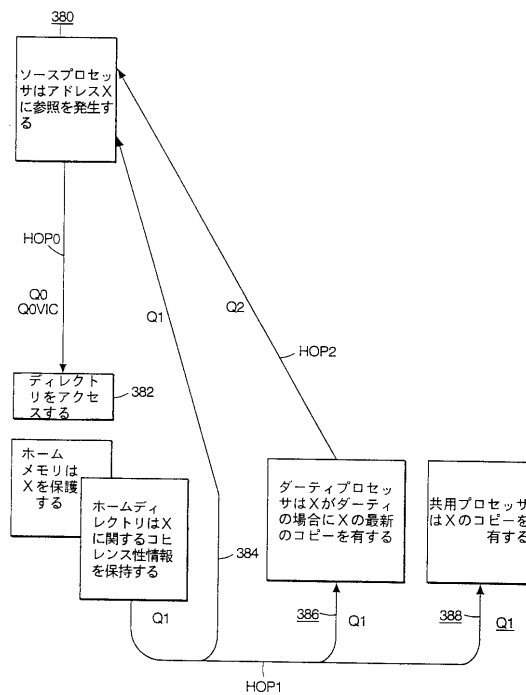
【図 2 1】



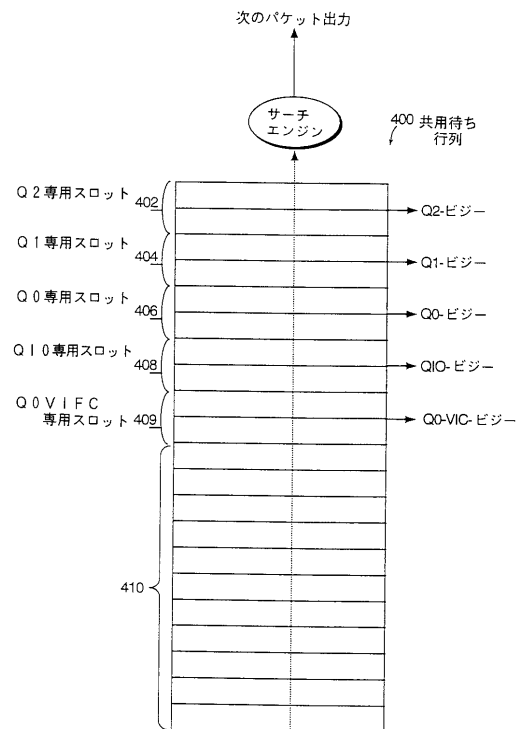
【図 2 2】



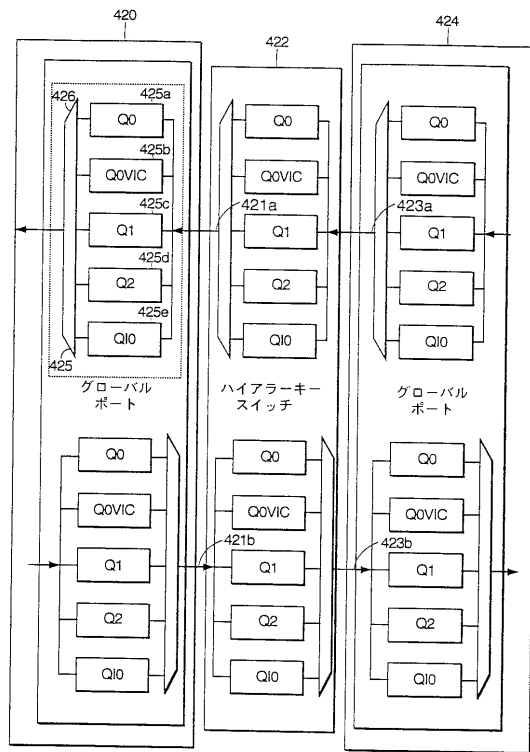
【図 2 3】



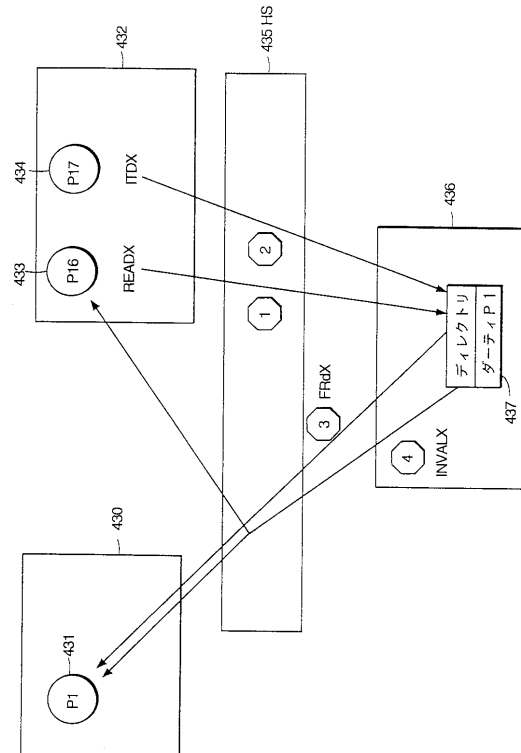
【図 2 4】



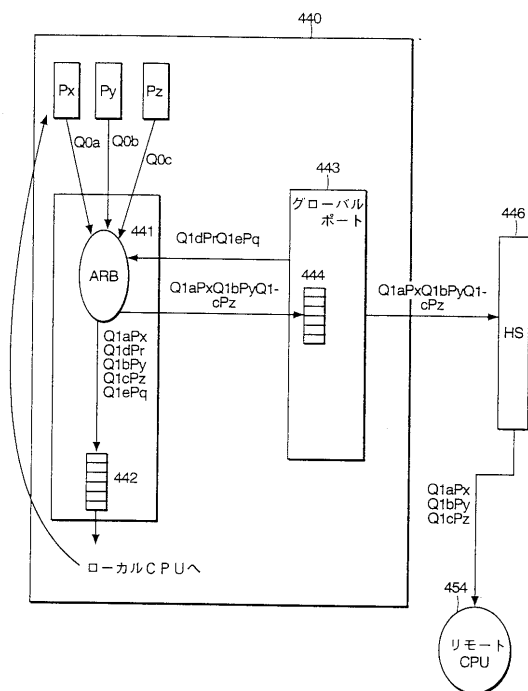
【図 25】



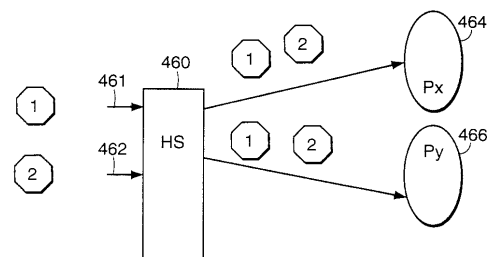
【図 26】



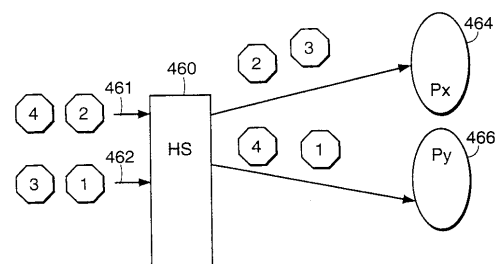
【図 27 A】



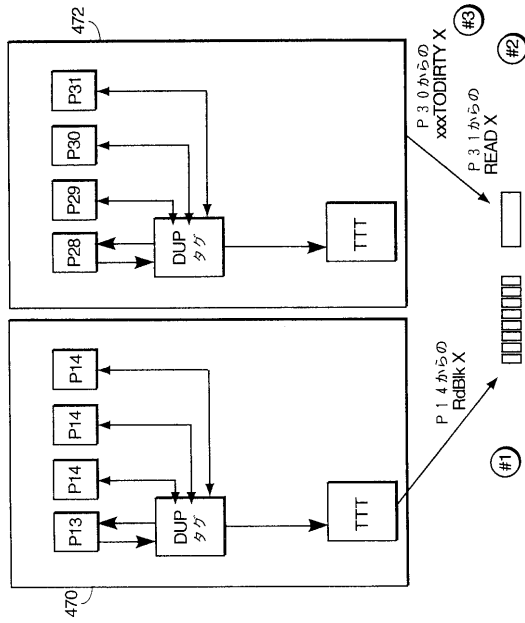
【図 27 B】



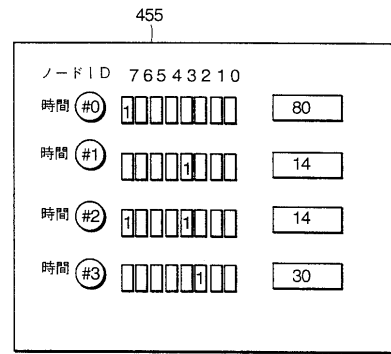
【図 27 C】



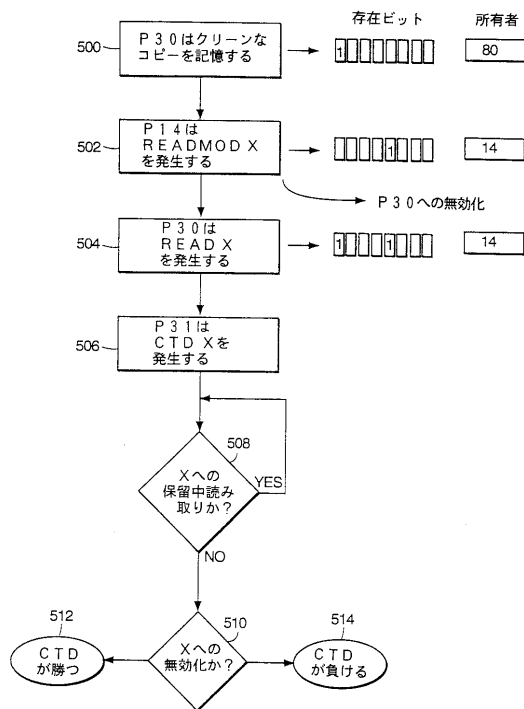
【図 28 A】



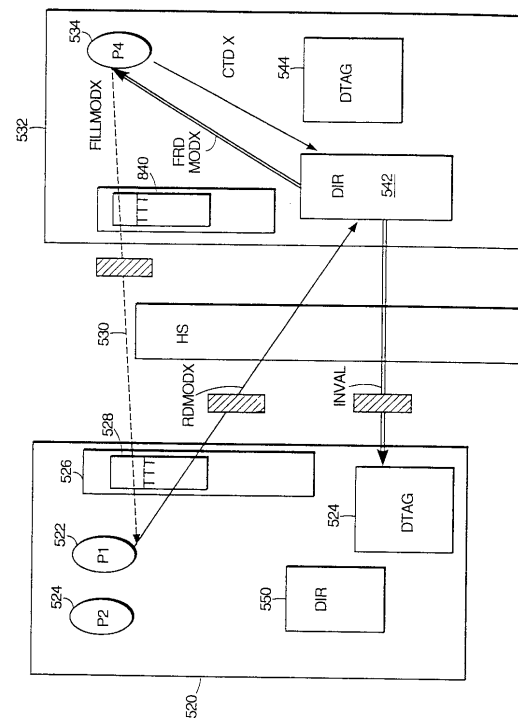
【図 28 B】



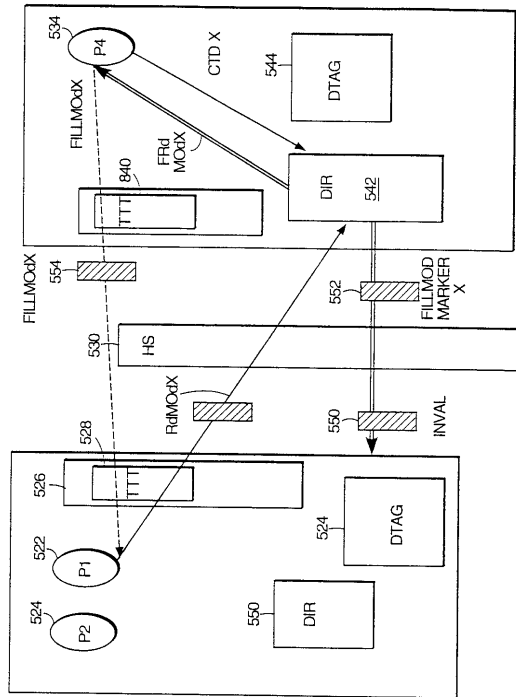
【図 29】



【図 30】



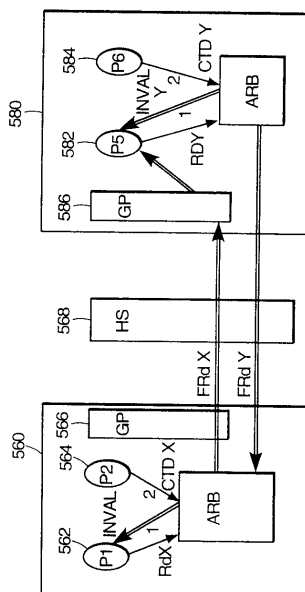
【図 3 1】



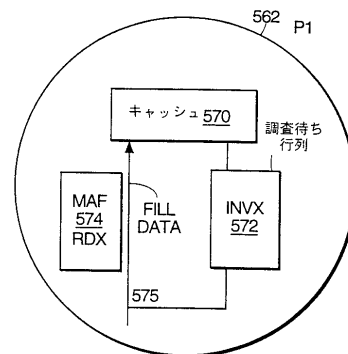
【図 3 2】

アドレス	コマンド	コマンドID	状態ビット	ACK/NACK ここに記入	フェッチループ CONSIST
X	RdmMod	P1	マーカーを ここに記入	シャドー	

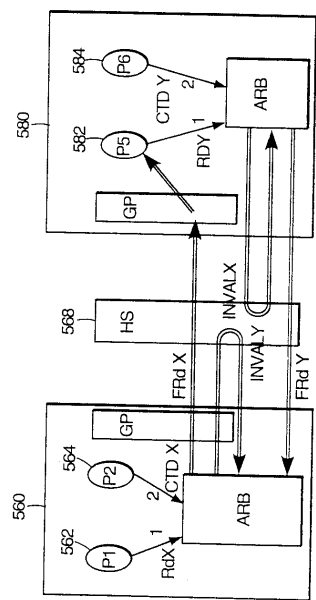
【図 3 3 A】



【図 3 3 B】



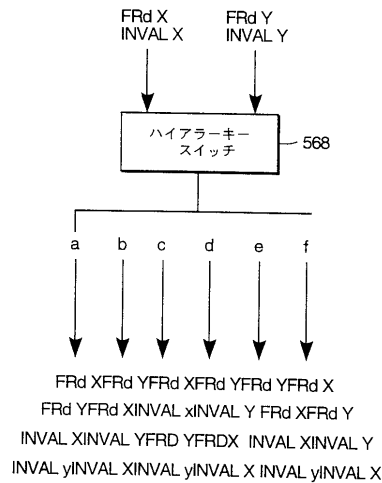
【図 3 4】



【図 3 5】

アドレス	コマンド	コマンドID	状態バス			ACK/NAK をここに記入	フェッチ	リレー CONSIG
			ここに 記入	マーカーを ここに記入	シャド			
X	FRD	P1			X			
X	INVAL	P2						

【図 3 6】



フロントページの続き

- (72)発明者 スティーヴン アール ヴァンドーレン
アメリカ合衆国 マサチューセッツ州 01532 ノースボロー デイヴィス ストリート 2
24
- (72)発明者 シモン シー ステイーリイ
アメリカ合衆国 ニューハンプシャー州 03051 ハドソン アンナ ルイス ドライヴ 8
- (72)発明者 マドハミトラ シャルマ
アメリカ合衆国 マサチューセッツ州 01545 シュローズバリー コモンズ ドライヴ 5
5 - 46
- (72)発明者 ディヴィッド エム フェンウィック
アメリカ合衆国 マサチューセッツ州 01545 アクトン ブラウン ベア クロッシング
297

合議体

審判長 吉岡 浩

審判官 清木 泰

審判官 田中 秀人

- (56)参考文献 特開平7 - 311751 (JP, A)
特開平4 - 328653 (JP, A)
特開平8 - 320827 (JP, A)