

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3884073号  
(P3884073)

(45) 発行日 平成19年2月21日(2007.2.21)

(24) 登録日 平成18年11月24日(2006.11.24)

(51) Int. Cl.	F I
<b>G06F 13/10 (2006.01)</b>	G06F 13/10 330A
<b>G06F 13/38 (2006.01)</b>	G06F 13/38 310A

請求項の数 1 (全 24 頁)

(21) 出願番号	特願平9-503399	(73) 特許権者	インテル・コーポレーション
(86) (22) 出願日	平成8年6月17日(1996.6.17)		アメリカ合衆国 95052 カリフォルニア州・サンタ クララ・ミッション カレッジ プーレバード・2200
(65) 公表番号	特表平11-514113	(74) 代理人	弁理士 山川 政樹
(43) 公表日	平成11年11月30日(1999.11.30)		
(86) 国際出願番号	PCT/US1996/010466	(74) 代理人	弁理士 黒川 弘朗
(87) 国際公開番号	W01997/000533		
(87) 国際公開日	平成9年1月3日(1997.1.3)	(74) 代理人	弁理士 紺野 正幸
審査請求日	平成15年5月21日(2003.5.21)		
(31) 優先権主張番号	08/490,651	(74) 代理人	弁理士 西山 修
(32) 優先日	平成7年6月15日(1995.6.15)		
(33) 優先権主張国	米国(US)	(74) 代理人	弁理士 鈴木 二郎

最終頁に続く

(54) 【発明の名称】 複数プロセッサ・システム内のプロセッサ間でメッセージを送信する装置

(57) 【特許請求の範囲】

【請求項 1】

ローカル・バスを介してローカル・プロセッサとメモリとに結合され、かつ第2のバスを介してバス・エージェントに結合されたメッセージ通信装置において：

a) インバウンド・フリー待ち行列であって、前記ローカル・プロセッサによりこの待ち行列の頭部ポインタによって示された該待ち行列の頭部から書き込まれ、かつ前記バス・エージェントによって使用されるために該待ち行列の尾部ポインタによって示された該待ち行列の尾部から読み出される空のメッセージ・バッファのハンドル情報であってバス・エージェントからローカル・プロセッサにメッセージを渡すために、該バス・エージェントによって使用され、かつ前記ローカル・プロセッサによって解放されるハンドル情報を記憶するインバウンド・フリー待ち行列と；

b) 前記インバウンド・フリー待ち行列内のハンドル情報を利用するためにこのインバウンド・フリー待ち行列に結合されたインバウンド・フリー回路手段であって、前記インバウンド・フリー待ち行列の頭部ポインタと尾部ポインタとに基づき、前記インバウンド・フリー待ち行列が「空」であるか否かを判定する判定手段と、前記インバウンド・フリー待ち行列の尾部から読み出したハンドル情報を保持するインバウンド・フリー・レジスタと

を備え、前記インバウンド・フリー回路手段によるハンドル情報に対する操作が、

(1) 前記インバウンド・フリー待ち行列が空でない場合には、該インバウンド・フリー待ち行列の尾部からハンドル情報をプリフェッチして前記インバウンド・フリー・レジスタ

10

20

タに格納し、

(2) 前記インバウンド・フリー待ち行列が空の場合には、該インバウンド・フリー・レジスタに該インバウンド・フリー待ち行列が空であることを示す「空」指示情報をロードし、かつ

(3) 一回のバス・トランザクションで、前記バス・エージェントが前記インバウンド・フリー・レジスタ内の情報を読み取ることを可能にするインバウンド・フリー回路手段と ;

c) インバウンド・ポスト待ち行列であって、前記バス・エージェントから読み出された後でこの待ち行列の頭部ポインタによって示された該待ち行列の頭部から書き込まれ、かつ前記ローカル・プロセッサによって使用されるために該待ち行列の尾部ポインタによって示された該待ち行列の尾部から読み出される、前記バス・エージェントからローカル・プロセッサに渡されるメッセージのハンドル情報であって前記インバウンド・フリー待ち行列から先に読み出されたハンドル情報を記憶するインバウンド・ポスト待ち行列と ;

d) 前記インバウンド・ポスト待ち行列内のハンドル情報を利用するためにこのインバウンド・ポスト待ち行列に結合されたインバウンド・ポスト回路手段であって、前記インバウンド・ポスト待ち行列の頭部ポインタと尾部ポインタとに基づき、前記インバウンド・ポスト待ち行列が「フル」であるか否かを判定する判定手段と、前記バス・エージェントからの書き込むべきメッセージのハンドル情報を保持するインバウンド・ポスト・レジスタと

を備え、前記インバウンド・ポスト回路手段によるハンドル情報に対する操作が、

(1) 前記インバウンド・ポスト待ち行列がフルである場合には、再試行信号を前記バス・エージェントに返し、

(2) 前記インバウンド・ポスト待ち行列がフルでない場合には、前記インバウンド・ポスト・レジスタに保持されたハンドル情報を前記インバウンド・ポスト待ち行列の頭部に書き込み、さらに

(3) 一回のバス・トランザクションで、前記バス・エージェントが前記インバウンド・ポスト・レジスタにハンドル情報を書き込むことを可能にするインバウンド・ポスト回路手段と

を備えたことを特徴とするメッセージ通信装置。

#### 【発明の詳細な説明】

##### 発明の背景

##### 1. 発明の分野

本発明は、複数プロセッサ・システム内の分野に関する。さらに詳細には、本発明は、複数プロセッサ・システム内のプロセッサ間でメッセージを送信する方法および装置に関する。

##### 2. 従来技術の説明

メッセージは、単にオペレーショナル・パラメータおよびデータを伝達するデータ構造である。メッセージは、1つまたは複数のプラットフォーム上で実行されている1つまたは複数のプロセス(すなわち、アプリケーション)によって生成される。プラットフォームは、メモリに関連するプロセッサまたはプロセッサのクラスタ、ローカル・メモリ・バス、およびメモリ入出力バスを含んでいる。プラットフォーム内のこれらの要素は、動作環境を構成する。

さらに、プラットフォームは、オペレーティング・システムの単一のインスタンスを実行する。言い換えれば、コンピュータ・システムは、単一のオペレーティング・システムが複数のプロセッサをサポートする分配処理システムである。メッセージは、ある特定のプラットフォーム上のプロセスの1つによって生成された後、処理のために他のプロセッサ・プラットフォームに送られる。

メッセージは、メモリ内に常駐し、命令および追加の情報のデータ・ブロックに対する他のポインタを含んでいる制御ブロックにポインタによって示される。例えば、制御ブロックは、特定の周辺装置(例えば、ハード・ディスク・ドライブ)を指定し、装置の指定さ

10

20

30

40

50

れたセクタからデータを読み取るよう要求する。

メッセージ送信は、プロセッサが「密に」結合された（すなわち、プロセッサが単一のキャッシュを共用する）対称複数プロセッサ・システム（SMP）内、およびプロセッサが共通のバス構造によって「粗に」結合された非対称複数プロセッサ・システム内のプロセッサ間で使用される。

メッセージを第1のプラットフォーム内のあるプロセッサから第2のプラットフォーム内の第2のプロセッサに送るとき、そのメッセージを向けられたプロセッサが自身の資源が自由なときにメッセージを処理することができるようにメッセージを待ち行列化しなければならない。

メッセージを待ち行列化する従来技術の方法は、主としてソフトウェア技法を使用して実施される。これらの方法は、共用された待ち行列構造への複数のアトミック・プロセッサ・アクセスを必要とする。例えば、単一のプロセッサ上で動作している複数のプロセスは、プロセッサによって共用されたメモリ内に配置されたメッセージの1つの待ち行列を共用する。プロセスの1つに対してアトミック・アクセスを行うために、オペレーティング・システムは、待ち行列へのアクセスを要求しているプロセスに、待ち行列に対するそのプロセス独占権（例えば、アトミック・アクセス）を与えるセマフォを付与する。セマフォは、単に共用されたデータ構造（すなわち、オペレーティング・システム・コンテキストの一部）に対するプロセス独占アクセス権を与えるオペレーティング・システム変数である。その場合、プロセスは、待ち行列に対してメッセージを追加したり、削除することができる。特定のプロセスは、セマフォを制御するとき、その待ち行列に対してアクセスを要求している他のプロセスをロックアウトする。他のプロセスは、共用された構造が使用できるようになるまで、第1のプロセスがセマフォを解放するのを待たなければならない。

複数プロセッサ・システムでは、複数のプロセッサがセマフォへのアクセス権を同時に得ようと試みることができる。したがって、同期のためにバス・ロック（すなわち、アトミック・アクセス）が必要となる。1つのプロセッサがバスをロックしている間、他のプロセッサは、第1のプロセッサがそのバスをロック解除するまでメモリ内の同じ共用された構造（すなわち、メモリ・ブロック）にアクセスすることはできない。セマフォはシステム・メモリ内にあるので、他のプロセッサは、セマフォを求めて争っていてもロックアウトされる。したがって、サスペンドすることができるソフトウェア・モジュール（すなわち、マルチタスク・オペレーティング・システム）内では、バス・ロックは使用できない。代わりに、これらのアプリケーション内でセマフォを獲得し、解放するとき、オペレーティング・システム・カーネルへのコールが必要となる。

上述の動作は、セマフォを待っている間またはバス・アクセスを待っている間、各プロセスがアイドルである時間のために非常に非効率的である。さらに、オペレーティング・システム・カーネルへの上述のコールは、高価なコンテキスト切換えを行う。

コンテキストは、単にアプリケーション（すなわち、アプリケーション・コードおよびデータ）に充てられるメモリ領域である。アプリケーション・コンテキストは、フラグ、変数、現在プロセスの状態を含んでいる。セマフォはアプリケーション・コンテキストと異なるコンテキスト（すなわち、オペレーティング・システム・コンテキスト）内のオペレーティング・システム変数であるので、システム資源はアプリケーション・コンテキストを切り換える必要がある。例えば、コンテキスト切換え中に、データ・ポインタを変更し、ポインタをスタック上に押し上げ、プロセス制御パラメータも変更する。

バス・ロック能力を有しない従来技術のコンピュータ・システムは、高度に複雑なアルゴリズムを使用して、プロセッサ間の同期を実施する。これらのシステムでは、性能がさらに低下する。

したがって、セマフォを使用せずに待ち行列への直接アクセスを効率的に可能にする方法および装置が必要である。

#### 発明の概要

複数プロセッサ・システム内のプロセッサ間でメッセージを送信する方法および装置。本

10

20

30

40

50

発明の方法および装置は、非対称複数プロセッサ・システム内のプロセッサ間でメッセージの通信を可能にする。非対称複数プロセッサ・システムは、単にプロセッサが様々なオペレーティング・システムを同時に実行しているシステムである。例えば、アプリケーション・プラットフォーム上のアプリケーション・プロセッサは、Windows NT<sup>TM</sup>など標準のアプリケーション・オペレーティング・システム・ソフトウェア上で動作している。しかしながら、入出力プラットフォーム上のプロセッサは、入出力動作（例えば、実時間オペレーティング・システム：RTOS）に適した特定のオペレーティング・システムを動作させている。特に、本発明は、1つまたは複数のプロセッサ・プラットフォーム上で実行されている1つまたは複数のプロセスからローカル・プロセッサを含んでいるプラットフォームへのメッセージを待ち行列化する迅速かつ直接的な機構を提供する。

10

本発明は、他のプラットフォームにメッセージ・バッファを割り当てるインバウンド・フリー待ち行列、および入出力プラットフォームの外部のプロセッサおよびバス・エージェントからメッセージをポストするインバウンド・ワーク待ち行列を提供する。さらに、本発明は、ローカル・プロセッサ（すなわち、入出力プラットフォーム用のプロセッサ）からのメッセージを、他のプラットフォーム上のプロセッサがこれらのメッセージを検索することができるように他のプロセッサ・プラットフォーム（すなわち、ホスト・プロセッサ）にポストするアウトバウンド・ワーク待ち行列を提供する。本発明はまた、ホスト・プロセッサがそれに対してメッセージ・バッファを解放するアウトバウンド・フリー待ち行列を提供する。この待ち行列は、ホスト・プロセッサがメッセージを処理した後、メッセージ・バッファをローカル・プロセッサに解放する。

20

本発明はホスト・プラットフォームと入出力プラットフォームとの間の非常に速くかつ効率的なハードウェア待ち行列インタフェースを提供するメッセージ・ユニットを使用してこれらの待ち行列を管理する。本発明は、単一のPCIバス・トランザクション・サイクル中に自由なメッセージ・バッファ、すなわち「エンプティ」インジケータを準備（すなわち、メッセージ・ユニット内のレジスタの読取り）することができる。さらに、本発明は、単一のPCIバス・トランザクション中にメッセージまたは「フル」インジケータのポストまたは検索（すなわち、メッセージ・ユニット内のレジスタの書込み）を可能にする。

本発明は、待ち行列をハードウェア・インタフェースを使用して管理するので、従来技術のソフトウェア待ち行列管理に勝るいくつかの利点を提供する。第1に、本発明はプロセスがフル待ち行列またはエンプティ待ち行列に対して待ち行列操作を実施しようとしたときデッドロックまたはロックアップを回避する。本発明のメッセージ・ユニットは、エンプティ・リストまたは待ち行列からフェッチしようとする試みが検出されたときエンプティ指示を迅速に戻す。同様に、本発明は、フル待ち行列にポストしようとする試みが検出されたとき特定の待ち行列がフルであるという指示を迅速に戻す。本発明は、最小のハードウェア資源を使用して効率的に実施することができる。

30

さらに、本発明は単一のバス・トランザクション内で待ち行列アクセスを実行するので、同期（すなわち、セマフォの獲得および解放）の必要がなく、またシステムの性能が大幅に改善される。待ち行列アクセスは、単に要素を待ち行列に追加すること、または要素を待ち行列から削除することである。待ち行列アクセスは、次の要素を見つけ出し、その要素を変更し、次の待ち行列アクセスのために次の要素を示す待ち行列記述子を修正する特定のタスクを含んでいる。これらのタスクは、本発明によって自動的に実施される。これらのタスクが完了する時間中、待ち行列は、他のプロセスが同じメッセージ・バッファを獲得したり、または他のメッセージを上書きしないようにロックされる。本発明は、単一のPCIバス・トランザクションが本来アトミックであることを利用するために1つのバス・トランザクション内で待ち行列アクセス（すなわち、トランザクションを実行しているバス・エージェントによる独占アクセス）を提供する。さらに、本発明は、読取り信号および再試行信号を介して自動的に同期を処理する。

40

さらに、本発明ではセマフォが不要であるので、システム資源をタイアップするコンテキスト切換えが不要である。メッセージ・ユニット内のレジスタに対する単一の読取りまた

50

は書込みだけで、特定の待ち行列にアクセスすることができ、また1つのバス・トランザクション内で読取りまたは書込みを行うことができるので、セマフォは不要である。  
本発明について、添付の図面において例を挙げて非限定的に説明する。図面中、同じ参照番号は同じ要素を示す。

#### 【図面の簡単な説明】

第1図は、本発明を実施する非対称複数プロセッサ・コンピュータ・システムのブロック図である。

第2図は、本発明を含む入出力プラットフォームを示す図である。

第3図は、本発明の一実施形態を示す図である。

第4図は、本発明の円形待ち行列を示す図である。

10

第5図は、本発明の円形待ち行列動作をさらに示す図である。

第6A図は、本発明のインバウンド・フリー状態機械を示す図である。第6B図は、インバウンド・フリー状態機械の状態図である。

第7A図は、本発明のインバウンド・ポスト状態機械を示す図である。第7B図は、インバウンド・ポスト状態機械の状態図である。

第8A図は、本発明のアウトバウンド検索状態機械を示す図である。第8B図は、アウトバウンド検索状態機械の状態図である。

第9A図は、本発明のアウトバウンド解放状態機械を示す図である。第9B図は、アウトバウンド解放状態機械の状態図である。

#### 好ましい実施形態の詳細な説明

20

第1図は、本発明を実施する複数プロセッサ・コンピュータ・システムのブロック図である。マルチプロセッサ・システム100は、ホスト・プロセッサ102を含んでいる。ホスト・プロセッサ102は、複数のプロセッサ（すなわち、密に結合されたプロセッサのクラスタ）を含んでいる。ホスト・プロセッサ102は、ホスト・バス103を介してホスト・メモリ104に結合される。メモリ・バス103はまた、ホスト・プロセッサ102およびメモリ104をホスト・チップ・セット105に結合する。ホスト・チップ・セット105は、メモリ・コントローラ、キャッシュ・コントローラ、およびメモリ・バス103と入出力（I/O）バス106（例えばPCIバス）とのインタフェースを提供するブリッジを含んでいる。

ホスト・チップ・セット105は、当技術分野において周知である。例えば、ホスト・プロセッサ102はIntel社製のPentium<sup>TM</sup>プロセッサである場合、適切なホスト・チップ・セット105は、これもIntel社製のTrident<sup>TM</sup>である。同様に、P6<sup>TM</sup>プロセッサが使用される場合、適切なホスト・チップ・セット105は、これもIntel社製のOrion<sup>TM</sup>チップ・セットである。ホスト・プロセッサ102、メモリ・バス103、ホスト・メモリ104、ホスト・チップ・セット105は、このマルチプロセッサ・システム100内ではホスト・プラットフォームと呼ばれる。

30

複数プロセッサ・システム100はさらに、第1のPCIバス106に結合された入出力プラットフォーム108を含んでいる。さらに、入出力プラットフォーム108は、第1のPCIバス106のアドレス空間と、入出力プラットフォーム108内に含まれているプロセッサのアドレス空間とのインタフェースを提供する。入出力プラットフォーム108はさらに、第1のPCIバス106を第2のPCIバス（図示せず）に結合するブリッジを含んでいる。

40

入出力プラットフォーム108は、ホスト・プロセッサ用の入出力サポート、および第1のPCIバス106および第2のPCIバスに結合されたデバイス（図示せず）を提供する。第2図は、本発明を含んでいる（前に第1図において要素108と呼んでいた）入出力プラットフォーム200を詳細に示す。入出力プラットフォーム200は、メモリ・コントローラ205からローカル・バス204を介してローカル・メモリ206に結合されたローカル・プロセッサ202を含んでいる。ローカル・プロセッサ202は、Intel 80960 JFプロセッサである。

アドレス・トランザクション・ユニット（ATU）218は、ローカル・バス204と（

50

前に第1図において要素106と呼んでいた)第1のPCIバス208とに結合される。アドレス・トランザクション・ユニット(ATU)218は、PCIバス208のアドレス空間内のアドレスをプロセッサ202アドレス空間内のアドレスに変換し、その逆も行う。したがって、PCIアドレス空間内のアドレスを有するPCIバス208に対するトランザクションは、メモリ・コントローラ205がローカル・メモリ206内の正確な位置またはMU210内の正確なレジスタ212にアクセスすることができるようにローカル・バス204アドレス空間に変換されなければならない。

ATU218は、ローカル・バス・トランザクションをPCIバス・トランザクションに変換するアウトバウンド・モジュール、PCIバス・トランザクションをローカル・バス・トランザクションに変換するインバウンド・モジュール、およびこのアドレス・トランザクションを管理する制御状態機械を含んでいる。本発明では、ATU218は、特定のPCIバス・トランザクションがMU210内のレジスタ212の1つにアクセスすることを検出するアドレス・デコードと考えられる。ATU218は、トランザクションがMU210内のレジスタ212の1つへのアクセスであることを検出した後、以下で説明するMU210内の制御状態機械214を開始するためにデータ・バス221を介して信号を送る。制御状態機械214は、MU210がトランザクションを受け取る準備ができていないことをATU218に通知するため、または要求側プロセスに再試行を通知するようATU218に指示するために、複数の信号をデータ・バス221を介してATU218に送る。

ローカル・バス・アービタ240は、ローカル・バス・マスタ(すなわち、MU210、ATU218のインバウンド・モジュール、およびローカル・プロセッサ202)のいずれかにローカル・バス204の制御権を付与する。アービトレーション回路240は、当技術分野において周知である。

メモリ・コントローラ205は、データ・バス224および225を介してローカル・メモリ206にアクセスするために備えられる。ローカル・バス204は、単一のデータ・バスとして示されているが、ローカル・バス204は、アドレス部分およびデータ部分から構成される。

バス・エージェント201は、ホスト・プロセッサまたは他の入出力プラットフォームである。さらに、バス・エージェント201は、第1図のホスト・メモリ104、ホスト・プロセッサ102、ホスト・チップ・セット105、メモリ・バス103を含んでいる。言い換えれば、バス・エージェント201は、それ自体サブシステムまたは任意のインテリジェント・バス・エージェントである。

メッセージ・ユニット(MU)210は、ローカル・バス204とATU218とに結合される。MU210は、本発明の教示を実施し、複数のレジスタ212および複数の状態機械214を含んでいる。これらのレジスタ212および状態機械214について、第3図に関して詳細に説明する。

第3図は、MU210内で実施された本発明を示す。MU210は、制御バス350を介してATU218に結合された複数の状態機械214を含んでいる。MU210はまた、複数のプリフェッチおよび一時レジスタ332を含んでいる。これらのレジスタ332は、データ・バス336を介してATU218に結合される。プリフェッチおよび一時レジスタ332はまた、データ・バス352を介して制御状態機械214によって制御される。レジスタ332はまた、ローカル・メモリ206にアクセスするためにデータ・バス334を介してローカル・バス204に結合される。

この実施形態では、MU210は、4つの円形待ち行列を使用したメッセージ送信方式を含んでいる。この実施形態では、4つのプリフェッチおよび一時レジスタ332がある。ホスト・プロセッサが円形待ち行列にデータを書き込むことを可能にする2つのレジスタが備えられる。ホスト・プロセッサが円形待ち行列の1つからデータを読み取ることを可能にする2つのレジスタが備えられる。

MU210はまた、データ・バス342を介して制御状態機械214に結合された複数の待ち行列ポインタ・レジスタ340を含んでいる。これらのレジスタ340は、待ち行列

10

20

30

40

50

207の頭部ポインタおよび尾部ポインタを記憶する。これらの待ち行列について、第4図および第5図に関して詳細に説明する。

#### 円形待ち行列

MU210は、4つの円形待ち行列207へのアクセスをバス・エージェント201に提供する。2つのインバウンド待ち行列および2つのアウトバウンド待ち行列がある。「インバウンド」および「アウトバウンド」は、活動メッセージの流れの方向を示す。「インバウンド」メッセージは、ローカル・プロセッサ202が処理すべきバス・エージェント201によってポストされた新しいメッセージであるか、またはバス・エージェント201が使用するために使用できるエンプティまたはフリーのメッセージ・バッファである。「アウトバウンド」メッセージは、ホスト・プロセッサ201が処理すべきローカル・プロセッサ202によってポストされた新しいメッセージであるか、またはローカル・プロセッサ202が使用するために使用できるフリーのメッセージ・バッファである。

一実施形態では、ホスト・プロセッサ/バス・エージェント201とローカル・プロセッサ202との間でメッセージを送るために使用される4つの円形待ち行列がある。インバウンド・メッセージを処理するために使用される2つのインバウンド待ち行列があり、アウトバウンド・メッセージを処理するために使用される2つのアウトバウンド待ち行列がある。インバウンド待ち行列の1つは、フリー待ち行列に指定され、インバウンド・フリー・メッセージ・ハンドルを含んでいる。メッセージ・ハンドルは、メッセージ・バッファの論理アドレスまたは物理アドレスである。他のインバウンド待ち行列は、ポスト待ち行列または作業待ち行列に指定され、インバウンド・ポスト・メッセージ・ハンドルを含んでいる。同様に、アウトバウンド待ち行列の1つはフリー待ち行列に指定され、他のアウトバウンド待ち行列はポスト待ち行列に指定される。

2つのアウトバウンド待ち行列は、ローカル・プロセッサ202がアウトバウンド・メッセージをポスト待ち行列内にポストし、外部ホスト・プロセッサ201からアウトバウンド・フリー待ち行列内に戻ったフリー・メッセージを受信することを可能にする。2つのインバウンド待ち行列は、バス・エージェント201がインバウンド・フリー待ち行列からフリー・メッセージ・バッファを獲得し、その後ローカル・プロセッサ202によって処理するためにそのバッファをインバウンド・フリー待ち行列にポストすることを可能にする。

円形待ち行列207のデータ記憶は、ローカル・メモリ206によって実施される。この特定の実施形態では、待ち行列内の各エントリは、メッセージ・ハンドルである32ビット・データ値である。さらに、待ち行列の読取りまたは書込みは、1つの待ち行列エントリに正確にアクセスすることができる。

各円形待ち行列は、頭部ポインタおよび尾部ポインタを有する。待ち行列への書込みは、待ち行列の頭部から行われ、読取りは、尾部から行われる。頭部ポインタおよび尾部ポインタは、ローカル・プロセッサ202上で動作しているソフトウェアまたはメッセージ・ユニット210によって増分される。頭部ポインタおよび尾部ポインタがどのようにしてローカル・プロセッサ202およびMU210によって増分されるかに関する詳細について、以下で説明する。

頭部ポインタおよび尾部ポインタは、それぞれの円形待ち行列内にオフセットされ、0から円形待ち行列サイズ-1に及ぶ(すなわち、ポインタを0からラベル付けし始める)。ポインタは、各待ち行列アクセスの後で増分される。頭部ポインタならびに尾部ポインタは、円形待ち行列サイズ(すなわち、待ち行列の終り)に達したときに0に丸められる。メッセージ・ユニット210は、ある条件のもとで、ローカル・プロセッサ202に対して割込みを発生するか、またはPCIバス割込み(すなわち、外部プロセッサに対する割込み)を発生する。一般に、ポスト待ち行列が書き込まれたとき、メッセージがポストされたターゲット・プロセッサを通知する割込みが発生する。

一実施形態では、各円形待ち行列のサイズは16Kバイト(4096ハンドル)から256Kバイト(65536ハンドル)までである。さらに、この実施形態では、4つのすべての待ち行列は、同じサイズであり、隣接している。したがって、円形待ち行列によって

10

20

30

40

50

必要とされるローカル・メモリの総量は、64 Kバイトから1 Mバイトまでである。これらの待ち行列は、ローカル・メモリ206内に常駐し、待ち行列の頭部ポインタおよび尾部ポインタは、MU210内のレジスタ内に常駐する。待ち行列サイズは、メッセージ・ユニット構成レジスタ(MUCR)内の待ち行列サイズ・フィールドによって決定される。MUCRの可能な1つのフォーマットを表1に示す。また、この実施形態では、4つのすべての待ち行列に対して1つのベース・アドレスがある。各待ち行列の開始アドレスは、待ち行列ベース・アドレスおよび待ち行列サイズ・フィールドに基づいている。ベース・アドレスは、これもMU210内に常駐する待ち行列ベース・アドレス・レジスタ(QBAR)内に記憶される。QBARの可能な1つのフォーマットを表2に示す。第6図から第9図に示す実施形態は、各待ち行列ごとに別個のベース・アドレスを含んでいる。

表1

MU構成レジスタ-MUCR			
ビット	デフォルト	読取り／書込み	説明
31:05	00000000H	読取り専用	予約
04:00	000000 <sub>2</sub>	読取り／書込み	待ち行列サイズ-このフィールドは、各円形待ち行列のサイズを決定する。4つのすべての待ち行列は同じサイズである。

表2

待ち行列ベース・アドレス・レジスタ-QBAR			
ビット	デフォルト	読取り／書込み	説明
31:20	000H	読取り／書込み	待ち行列ベース・アドレス-円形待ち行列のローカル・メモリ・アドレス。
19:00	000000H	読取り専用	予約

第4図は、本発明の4つの円形待ち行列を示す。ローカル・メモリ206内に常駐する2つのアウトバウンド待ち行列410および420と、2つのインバウンド待ち行列430および440がある。

ローカル・プロセッサ202は、アウトバウンド・ポスト待ち行列420の頭部に書込みを行うことによってアウトバウンド・メッセージ422をポストする。ホスト・プロセッサ201は、アウトバウンド・ポスト待ち行列420の尾部の読取りを行うことによってアウトバウンド・ポスト待ち行列420からポスト・メッセージを取り出す。

ホスト・プロセッサ201は、アウトバウンド・フリー待ち行列410の頭部に書込みを行うことによってアウトバウンド・メッセージ・バッファ412を解放する。ローカル・プロセッサ202は、アウトバウンド・フリー待ち行列410の尾部からフリー・メッセージ・バッファ414を読み取る。

ホスト・プロセッサまたはバス・エージェント201は、インバウンド・ポスト待ち行列430の頭部に書込みを行うことによってインバウンド・ポスト待ち行列430にインバウンド・メッセージ432をポストする。ローカル・プロセッサ202は、インバウンド・ポスト待ち行列430の尾部からこれらのポスト・メッセージを読み取る。ホスト・プ

10

20

30

40

50



ロセッサがインバウンド・ポスト待ち行列 4 3 0 に書込みを行ったとき、ローカル・プロセッサ 2 0 2 に対して割込み 4 3 6 が発生する。

ローカル・プロセッサ 2 0 2 によってアウトバウンド・ポスト待ち行列 4 2 0 に対してメッセージがポストされたとき、ホスト・プロセッサ 2 0 1 に対して割込み 4 2 6 が発生する。ここでは、P C I バス仕様改訂 2 . 0 によって指定されている割込みを使用する。

ローカル・プロセッサ 2 0 2 は、インバウンド・フリー待ち行列 4 4 0 の頭部に書込みを行うことによってこの待ち行列 4 4 0 にフリー・メッセージ・バッファ 4 4 2 を戻す。ホスト・プロセッサ/バス・エージェント 2 0 1 は、データ・バス 4 4 4 を介してインバウンド・フリー待ち行列 4 4 0 の尾部から読取りを行うことによってフリー・メッセージ・バッファを獲得する。

10

第 5 図は、アウトバウンド・フリー待ち行列 5 1 0、アウトバウンド・ポスト待ち行列 5 2 0、インバウンド・ポスト待ち行列 5 3 0、インバウンド・フリー待ち行列 5 4 0 を示す。

#### アウトバウンド・フリー待ち行列

アウトバウンド・フリー待ち行列 ( O F Q ) 5 1 0 は、ローカル・プロセッサ 2 0 2 が使用すべきバス・エージェント 2 0 1 によってそこに配置された ( すなわち、解放された ) エンプティ・メッセージ用のハンドルを保持する。ホスト・プロセッサ 2 0 1 は、アウトバウンド・ポート 5 1 6 内のレジスタに書込みを行うことによって O F Q 5 1 0 に対してメッセージ・バッファを解放する。O F Q 5 1 0 は、ローカル・プロセッサ 2 0 2 によって待ち行列尾部から読み取られ、ホスト・プロセッサ 2 0 1 によって待ち行列頭部に書き込まれる。頭部ポインタ ( O F H P ) 5 1 2 は、メッセージ・ユニット 2 1 0 によって維持される。アウトバウンド・フリー待ち行列尾部ポインタ ( O F T P ) 5 1 4 は、ローカル・プロセッサ 2 0 2 上で動作しているソフトウェアによって維持される。

20

アウトバウンド待ち行列ポート 5 1 6 にアクセスする P C I 書込みトランザクションの場合、M U 2 1 0 は、メッセージ・ハンドル ( すなわち、フリー・メッセージ・バッファに対するアドレス ) を、アウトバウンド・フリー頭部ポインタ・レジスタ ( O F H P R ) 9 2 6 内に記憶された頭部ポインタ ( O F H P ) 5 1 2 によって示されたローカル・メモリ 2 0 6 内の位置に書き込む。ローカル・メモリ・アドレスは、待ち行列ベース・アドレス・レジスタ + 3 \* 待ち行列サイズ + アウトバウンド・フリー頭部ポインタ・レジスタ ( O F H P R ) 9 2 6 である。O F H P R の可能な 1 つのフォーマットを表 3 に示す。

30

アウトバウンド待ち行列ポート 5 1 6 に書き込まれたデータがローカル・メモリ 2 0 6 に書き込まれたとき、M U 2 1 0 は O F H P 5 1 2 を増分する。

データがローカル・メモリ 2 0 6 に書き込まれるまで P C I 書込みトランザクションが M U 2 1 0 によって受け取られ、かつ O F H P 5 1 2 が増分されたときから、インバウンド待ち行列ポート 5 1 6 にアクセスしようと試みる P C I トランザクションは、待ち状態を挿入することによって遅延される。待ち状態を挿入している間に P C I 待ち時間違反が発生した場合、外部 P C I エージェント 2 0 1 に再試行が通知される。

ローカル・プロセッサ 2 0 2 は、アウトバウンド・フリー待ち行列尾部ポインタ ( O F T P ) 5 1 4 によって示されたローカル・メモリ位置を読み取ることによって O F Q 5 1 0 からメッセージ・バッファ・ハンドルを取り出す。ローカル・メモリ・アドレスは、待ち行列ベース・アドレス・レジスタ + 3 \* 待ち行列サイズ + アウトバウンド・フリー尾部ポインタ・レジスタ ( O F T P R ) 9 3 8 である。O F T P R の可能な 1 つのフォーマットを表 4 に示す。次いで、ローカル・プロセッサ 2 0 2 は、( 第 9 A 図に示される ) アウトバウンド・フリー尾部ポインタ・レジスタ ( O F T P R ) 9 3 8 内の O F T P 5 1 4 を増分する。

40

表 3

アウトバウンド・フリー頭部ポインタ・レジスタ－OFHPR			
ビット	デフォルト	アクセス	説明
31 : 19	0000H	読取り専用	予約
18 : 02	0000H	読取り／ 書込み	アウトバウンド・フリー頭部ポインタ・レジスタ－アウトバウンド・ポスト待ち行列に対する頭部ポインタのローカル・メモリ・オフセット
01 : 00	00 <sub>2</sub>	読取り専用	予約

10

表 4

アウトバウンド・フリー尾部ポインタ・レジスタ－OPTPR			
ビット	デフォルト	アクセス	説明
31 : 19	0000H	読取り専用	予約
18 : 02	0000H	読取り／ 書込み	アウトバウンド・フリー尾部ポインタ・レジスタ－アウトバウンド・フリー待ち行列に対する尾部ポインタのローカル・メモリ・オフセット
01 : 00	00 <sub>2</sub>	読取り専用	予約

20

#### アウトバウンド・ポスト待ち行列

アウトバウンド・ポスト待ち行列（OPQ）520は、ホスト・プロセッサ201が取り出し、処理するために、ローカル・プロセッサ202によってそこに配置されたポスト・メッセージのハンドルを格納している。ホスト・プロセッサ201は、アウトバウンド待ち行列ポート516内のレジスタの読取りを行うことによってOPQ520からメッセージを取り出す。ローカル・プロセッサ202は、待ち行列頭部に書込みを行うことによってOPQ520にメッセージを追加する。頭部ポインタ（OPHP）522は、ローカル・プロセッサ202によって維持される。尾部ポインタ（OPTP）524は、メッセージ・ユニット210によって維持される。

30

アウトバウンド待ち行列ポート516にアクセスするPCI読取りトランザクションの場合、MU210は、OPTP524によって示されたローカル・メモリ位置においてデータをプリフェッチする。ローカル・メモリ・アドレスは、待ち行列ベース・アドレス・レジスタ+2\*待ち行列サイズ+アウトバウンド・ポスト尾部ポインタ・レジスタ（OPTPR）826（第8A図に示す）である。OPQ520がエンプティでない（すなわち、頭部ポインタ522と尾部ポインタ524が等しくない）場合、メッセージ・ハンドルが要求側プロセッサ201に供給される。OPQ520がエンプティである（すなわち、頭部ポインタ522と尾部ポインタ524が等しい）場合、-1の値（FFFF.FFFFH）が要求側プロセッサ201に供給される。OPQ520待ち行列がエンプティでなく、かつMU210が尾部におけるデータのプリフェッチに成功した場合、MU210は、OPTPR826内の尾部ポインタ（OPTR）524を増分する。

40

上述のように、プリフェッチ機構は、頭部ポインタ522と尾部ポインタ524が等しい（すなわち、OPQ520がエンプティである）場合、-1の値（FFFF.FFFFH）を（以下で第8A図に関して説明する）プリフェッチ・レジスタ806内にロードする。メッセージがOPQ520に追加され、それがエンプティでなくなったときにORR8

50

06を更新するために、MU210内のプリフェッチ機構は、ORR806が(FFFF.FFFFH)を含んでいる場合、自動的にプリフェッチを開始し、アウトバウンド・ポスト頭部ポインタ・レジスタ(OPHPR)422がローカル・プロセッサ202によって書き込まれる。OPHPRの可能な1つのフォーマットを表5に示す。ローカル・プロセッサ202は、ローカル・プロセッサ202がOPQ520にメッセージを追加したときOPHPR422を更新する。

プリフェッチは、外部バス・エージェント201から見てアトミックに見えなければならない。プリフェッチが開始されたとき、アウトバウンド待ち行列ポート516内の(以下で第8A図に関して説明する)アウトバウンド検索レジスタ806にアクセスしようとするPCITランザクションは、プリフェッチが完了するまで待ち状態を挿入することによって遅延される。待ち状態を挿入している間にバス待ち時間違反が発生した場合、外部バス・エージェント201に再試行信号が通知される。

OPHP522がOPTP524に等しくない場合、ホスト・プロセッサ201に対してPCI割込みが発生する。OPHP522がOPTP524に等しい場合、割込みは発生しない。アウトバウンド・ドーベル・レジスタ内のアウトバウンド・ポスト待ち行列割込みビットは、OPHPR838の値とOPTPR828の値との比較の状態を示す。頭部ポインタ522と尾部ポインタ524が等しい場合、割込みはクリアされる。これは、ホスト・プロセッサ201がOPQ520をエンプティにするために十分な待ち行列エントリを読み取る場合に起こる。割込みは、ソフトウェアによって制御されるアウトバウンド・ドーベル・マスク・レジスタによってマスクされる。

ローカル・プロセッサ202は、頭部ポインタ(OPHP)522によって示されたローカル・メモリ位置にデータを書き込むことによってメッセージをOPQ520内に入れる。ローカル・メモリ・アドレスは、待ち行列ベース・アドレス・レジスタ+アウトバウンド・ポスト頭部ポインタ・レジスタ838である。OPTPRの可能な1つのフォーマットを表6に示す。次いで、ローカル・プロセッサ202は、アウトバウンド・ポスト頭部ポインタ・レジスタ838内のOPHP522を増分する。

表5

アウトバウンド・ポスト頭部ポインタ・レジスタ－OPHPR			
ビット	デフォルト	アクセス	説明
31:19	0000H	読取り専用	予約
18:02	0000H	読取り／ 書込み	アウトバウンド・ポスト頭部ポインタ・レジスタ－アウトバウンド・ポスト待ち行列に対する頭部ポインタのローカル・メモリ・オフセット
01:00	00 <sub>2</sub>	読取り専用	予約

表 6

アウトバウンド・ポスト尾部ポインタ・レジスタ - O P T P R			
ビット	デフォルト	アクセス	説明
31 : 19	0000H	読取り専用	予約
18 : 02	0000H	読取り／ 書込み	アウトバウンド・ポスト尾部ポインタ・レジスタ・アウトバウンド・ポスト待ち行列に対する尾部ポインタのローカル・メモリ・オフセット
01 : 00	00 <sub>2</sub>	読取り専用	予約

10

#### インバウンド・ポスト待ち行列

インバウンド・ポスト待ち行列 (IPQ) 530 は、ローカル・プロセッサ 202 が処理するために、バス・エージェント 201 によってそこに配置されたポスト・メッセージのハンドルを保持する。ホスト・プロセッサ 201 またはバス・エージェントは、インバウンド待ち行列ポート 536 内のレジスタに書込みを行うことによって IPQ 530 にメッセージをポストする。IPQ 530 は、ローカル・プロセッサ 202 によって待ち行列尾部から読み取られ、外部バス・エージェント 201 によって待ち行列頭部に書き込まれる。尾部ポインタ (IPTP) 534 は、ローカル・プロセッサ 202 上で動作しているソフトウェアによって維持される。頭部ポインタ (IPHP) 532 は、MU210 によって維持される。

20

インバウンド待ち行列ポート (IQP) 536 にアクセスする PCI 書込みトランザクションの場合、MU210 は、(第 7A 図に示される) インバウンド・ポスト頭部ポインタ・レジスタ (IPHPR) 724 内に記憶された IPHP 532 によって示されたローカル・メモリ位置にデータを書き込む。ローカル・メモリ・アドレスは、待ち行列ベース・アドレス・レジスタ + 待ち行列サイズ + インバウンド・ポスト頭部ポインタ・レジスタ (IPHPR) 724 である。IPHPR の可能な 1 つのフォーマットを表 7 に示す。IPTPR の可能な 1 つのフォーマットを表 8 に示す。

30

インバウンド待ち行列ポート 536 に書き込まれたデータがローカル・メモリ 206 に書き込まれたとき、MU210 は IPHPR 724 を増分する。データがローカル・メモリ 206 に書き込まれ、IPHPR 724 が増分されたとき、MU210 は、ローカル・プロセッサ 202 に対して割込みを発生する。この割込みは、インバウンド・ドーベル・レジスタのインバウンド・ポスト待ち行列割込みビットを設定することによって記録される。割込みは、ソフトウェアによって制御されるインバウンド・ドーベル・マスク・レジスタによってマスクされる。

表 7

インバウンド・ポスト頭部ポインタ・レジスタ - I P H P R			
ビット	デフォルト	アクセス	説明
31 : 19	0000H	読取り専用	予約
18 : 02	0000H	読取り／ 書込み	インバウンド・ポスト頭部ポインタ・レジスタ・インバウンド・ポスト待ち行列に対する頭部ポインタのローカル・メモリ・オフセット
01 : 00	00 <sub>2</sub>	読取り専用	予約

40

50

表 8

インバウンド・ポスト尾部ポインタ・レジスタ - I P T P R			
ビット	デフォルト	アクセス	説明
31 : 19	0000H	読取り専用	予約
18 : 02	0000H	読取り／ 書込み	インバウンド・ポスト尾部ポインタ・レジスタ・インバウンド・ポスト待ち行列に対する尾部ポインタのローカル・メモリ・オフセット
01 : 00	00 <sub>2</sub>	読取り専用	予約

10

#### インバウンド・フリー待ち行列

インバウンド・フリー待ち行列 540 は、バス・エージェント 201 が使用するために、ローカル・プロセッサ 202 によってそこに配置されたエンプティ・メッセージ・バッファのハンドルを保持する。ホスト・プロセッサ 201 は、インバウンド待ち行列ポート 536 内のレジスタの読取りを行うことによって I F Q 540 からメッセージ・バッファを割り振られる。I F Q 540 は、外部バス・エージェント 201 によって待ち行列尾部から読み取られ、ローカル・プロセッサ 202 によって待ち行列頭部に書き込まれる。頭部ポインタ (I P H P) 542 は、ローカル・プロセッサ 202 上で動作しているソフトウェアによって維持される。尾部ポインタ (I F T P) 544 は、M U 210 によって維持される。

20

インバウンド待ち行列ポート (I Q P) 536 にアクセスする P C I 読取りトランザクションの場合、M U 210 は、I F T P 544 によって示されたローカル・メモリ位置においてデータをプリフェッチする。ローカル・メモリ・アドレスは、尾部ポインタを記憶する待ち行列ベース・アドレス・レジスタ + インバウンド・フリー尾部ポインタ・レジスタ (I F T P R) 626 である。I F T P R の可能な 1 つのフォーマットを表 10 に示す。I F Q 540 がエンプティでない (すなわち、頭部ポインタと尾部ポインタが等しくない) 場合、ホスト・プロセッサまたはバス・エージェントによる次のアクセスのために I F T P 544 によって示されたデータが供給される。I F Q 540 がエンプティである (すなわち、頭部ポインタと尾部ポインタが等しい) 場合、- 1 の値 (F F F F . F F F F H) が要求側ホスト・プロセッサまたはバス・エージェントに供給される。I F Q 540 がエンプティでなく、かつ M U 210 が I F T P 544 によって示されたデータをプリフェッチした場合、M U 210 は、インバウンド・フリー尾部ポインタ・レジスタ (I F T P R) 626 内のポインタの値を増分する (第 6 A 図に示す)。

30

P C I 読取りアクセスに対する待ち時間を短縮するために、M U 210 は、I F Q 540 へのアクセスを予測するプリフェッチ機構を実施する。M U 210 は、I F Q 540 の尾部からデータをプリフェッチし、それを内部プリフェッチ・レジスタ内にロードする。P C I 読取りアクセスが行われたとき、データは、プリフェッチ・レジスタから直接読み取られる。

40

プリフェッチ機構は、頭部ポインタと尾部ポインタが等しい (すなわち、I F Q 540 がエンプティである) 場合、- 1 の値 (F F F F . F F F F H) をプリフェッチ・レジスタ 806 内にロードする。メッセージが I F Q 540 に追加され、それがエンプティでなくなったときにプリフェッチ・レジスタを更新するために、プリフェッチ機構は、プリフェッチ・レジスタが F F F F . F F F F H を含んでいる場合、自動的にプリフェッチを開始し、インバウンド・フリー頭部ポインタ・レジスタ (I F H P R) 638 が書き込まれる。I F H P R の可能な 1 つのフォーマットを表 9 に示す。ローカル・プロセッサ 202 上で動作しているソフトウェアは、I F H P 542 が I F Q 540 にメッセージを追加したときに I F H P 542 を更新する。

50

プリフェッチは、外部バス・エージェント 201 から見てアトミックに見えなければならない。プリフェッチが開始されたとき、インバウンド待ち行列ポート 536 内のインバウンド・フリー・レジスタにアクセスしようと試みる P C I トランザクションは、プリフェッチが完了するまで待ち状態を挿入することによって遅延される。待ち状態を挿入している間に P C I 待ち時間違反が発生した場合、M U 210 によって外部バス・エージェント 201 に再試行が通知される。

ローカル・プロセッサ 202 は、頭部ポインタ ( I F H P ) 542 によって示されたローカル・メモリ位置にデータを書き込むことによってメッセージを I F Q 540 内に入れる。ローカル・メモリ・アドレスは、待ち行列ベース・アドレス・レジスタ + インバウンド・フリー頭部ポインタ・レジスタ ( I F H P R ) 638 である。次いで、ローカル・プロセッサ 202 上で動作しているソフトウェアは、 I F H P R 638 を増分する。

インバウンド・フリー表

表 9

インバウンド・フリー頭部ポインタ・レジスタ - I F H P R			
ビット	デフォルト	アクセス	説明
31 : 19	0000H	読取り専用	予約
18 : 02	0000H	読取り／書込み	インバウンド・フリー頭部ポインタ・レジスタ・インバウンド・フリー待ち行列に対する頭部ポインタのローカル・メモリ・オフセット
01 : 00	00 <sub>2</sub>	読取り専用	予約

表 10

インバウンド・フリー尾部ポインタ・レジスタ - I F T P R			
ビット	デフォルト	アクセス	説明
31 : 19	0000H	読取り専用	予約
18 : 02	0000H	読取り／書込み	インバウンド・フリー尾部ポインタ・レジスタ・インバウンド・フリー待ち行列に対する尾部ポインタのローカル・メモリ・オフセット
01 : 00	00 <sub>2</sub>	読取り専用	予約

第 6 A 図は、本発明がどのようにして P C I バス上のバス・エージェントにフリー・メッセージ・バッファを割り振るのかを示す。M U 210 内で実現される。データは、ローカル・メモリ 206 内に配置されたインバウンド・フリー待ち行列 ( I F Q ) 540 からローカル・データ・バスを介してインバウンド・フリー・レジスタ ( I F R ) 606 まで移動する。このコンテキストでは、データは、特にメッセージ・バッファのアドレス ( すなわち、メッセージ・ハンドル ) を示す。その後、データは、インバウンド・フリー・レジスタ 606 からデータバス 608 を介して A T U 218 まで移動し、その後データバス 610 を介して P C I バス 208 上のバス・エージェントまで移動する。

M U 210 は、いくつかの制御信号を発生し、受信するフリー・メッセージ・バッファを割り振るインバウンド・フリー状態機械 612 を含んでいる。インバウンド・フリー状態機械 612 の状態図について、第 6 B 図に関して詳細に説明する。

I F Q 602 にメッセージ・バッファを要求するために、バス・エージェントは、 P C I

バス 208 およびデータバス 610 を介して読取りトランザクションを ATU 218 に送る。インバウンド・フリー・レジスタ 606 のアドレスを指定する読取りトランザクションは、ATU 218 によって検出される。ATU 218 が、バス・エージェントがインバウンド・フリー・レジスタ 606 を読み取ることを望んでいることを検出した後、ATU は、IFR\_\_Ready 信号 614 の状態をテストする。IFR\_\_Ready 信号 614 がアサートされた場合、ATU は、バス 608 を介して IFR 606 内のデータを ATU 218 に供給する PCI トランザクションを完了し、状態機械 612 に対して Read\_\_Inbound\_\_Free 信号 616 を発生する。

IFR\_\_Ready 信号 614 がアサート解除された場合（すなわち、状態機械 612 が準備できていない場合）、ATU 218 は、待ち状態を挿入し、IFR\_\_Ready 信号 614 がアサートされるまで Read\_\_IFR 信号 616 を送らない。IFR\_\_Ready 信号 614 は、IFR 606 内に状態データがある（すなわち、状態機械 612 がまだ IFR 606 へのデータのプリフェッチを完了していない）場合、アサート解除される。状態機械 612 は、Read\_\_IFR 信号 616 を受信した後、ローカル・バス・アービタ 240 に Memory\_\_Read\_\_Request 信号 618 を送り、IFR\_\_Ready 信号 614 をアサート解除する。許可信号 632 に基づいて、MU 210 は、単に IFQ 602 の適切な尾部アドレスをローカル・アドレス・バス 630 にアサートする。次いで、データがローカル・データ・バス 604 を介してローカル・メモリ 206 から IFR 606 に転送される（すなわち、IFQ 602 の尾部において値を読み取る）。MU 210 は、IFQ 602 の適切な尾部アドレスを計算する加算器 624 を含んでいる。加算器 624 は、インバウンド・フリー尾部ポインタ・レジスタ（IFTPR）626 およびインバウンド・フリーベースレジスタ（IFBR）628 の内容の合計を発生する。

IFQ 540 の尾部ポインタによって示されたデータがローカル・データ・バス 604 上にきた後、状態機械 612 は、ローカル・データ・バス 604 上のデータを IFQ 606 内にラッチするためにラッチ信号 634 を送り、増分信号 644 を IFTPR 626 に送る。したがって、次の使用できるメッセージ・バッファのプリフェッチが行われた。

MU 210 はまた、インバウンド・フリー頭部ポインタ・レジスタ（IFHPR）638 内の値とインバウンド・フリー尾部ポインタ・レジスタ（IFTPR）626 内の値とを比較するコンパレータ 636 を含んでいる。これら 2 つの値が等しい場合、コンパレータ 636 は、エンプティ信号 640 を発生する（すなわち、待ち行列内にフリー・メッセージ・バッファがない）。このエンプティ信号 640 は、状態機械 612 に送られ、状態機械 612 にプリセット信号をアサートさせる。プリセット信号 642 により、IFR 606 の内容がエンプティ指示に対して予約された（すなわち、有効なバッファアドレスでない）所定の値に設定される。バス・エージェントは、IFR 606 を読み取る場合、IFR 606 内に記憶されたプリフェッチされたデータか、または IFQ 602 がエンプティであることを示すプリセット値に直ちにアクセスする

第 6 B 図は、インバウンド・フリー状態機械 612 の状態図を示す。状態機械 612 は、エンプティ状態 650、プリフェッチ状態 652、ブライム状態 656 の 3 つの状態を有する。状態機械 612 は、エンプティ信号 654 がアサート解除されるまでエンプティ状態 650 にある。not\_\_Empty 信号は、状態機械 612 をエンプティ状態 650 からプリフェッチ状態 652 に遷移させ、状態機械 612 は、Memory\_\_Read\_\_Request 信号 618 を発生し、IFR\_\_Ready 信号 614 をアサート解除する。状態機械 612 は、許可信号 632 に基づいてプリフェッチ状態 652 からブライム状態 656 に遷移する。許可信号 632 を受信すると、状態機械 612 は、Latch\_\_IFR 信号 634、Increment\_\_IFTPR 信号 644 を出力し、IFR\_\_Ready 信号 614 をアサートする。状態機械 612 は、Read\_\_IFR 信号 616 が受信されたとき、ブライム状態 656 からプリフェッチ状態 652 に遷移し、エンプティ信号 654 はアサートされない。この遷移はまた、Memory\_\_Read\_\_Request 信号 618 を発生し、IFR\_\_Ready 信号 614 をアサートする。

状態機械 612 は、Read\_\_IFR 信号 616 が受信されたとき、ブライム状態 656

10

20

30

40

50

からエンプティ状態 650 に遷移し、エンプティ信号 640 はアサートされる。この遷移は、プリセット信号 642 を発生させる。

第 7 A 図は、本発明がどのようにしてバス・エージェントによって生成されたメッセージを、ローカル・メモリ 206 内に配置されたインバウンド・ポスト待ち行列 (IPQ) 530 内にポストするかを示す。

バス・エージェントがインバウンド・ポスト・レジスタ (IPR) 706 に書込みを行いたい場合、データは、PCIバス 208 からデータバス 702 を介して ATU 218 まで進み、次いでデータバス 704 を介して IPR 706 に進む。データは、IPR 706 内にラッチされた後、ローカル・データ・バス 604 を介してローカル・メモリ 206 内の IPQ 530 内に転送される。

10

ATU 218 は、IPR\_\_Ready 信号 716 の状態をテストする。IPR\_\_Ready 信号 716 がアサートされた場合、ATU 218 は、データを IPR 706 に供給し、状態機械 712 に対して Write\_\_IPR 信号 718 を発生することによって PCI トランザクションを完了する。

IPR\_\_Ready 信号 716 がアサートされない場合、ATU 218 は、待ち状態を挿入し、IPR\_\_Ready 信号 716 がアサートされたとき PCI トランザクションを完了する。要求側プロセスは、バスの制御権を保持し、PCI 待ち時間ルールが違反されなければ PCI トランザクションは完了する。

ATU 218 はまた、IPR\_\_Retry 信号 714 の状態をテストする。IPR\_\_Retry 信号 714 がアサートされた場合、PCI トランザクションは完了せず、要求側プロセスは、再試行を通知され、バスを解放し後で再度試行する。

20

MU 210 のインバウンド・ポスト状態機械 712 について、第 7 B 図に示される状態図によって詳細に説明する。状態機械 712 は、アイドル状態 750、ポスト状態 752、フル状態 754 の 3 つの状態を有する。状態機械 712 は、Write\_\_Inbound\_\_Post 信号 718 が ATU 218 によってアサートされたときアイドル状態 750 からポスト状態 752 に遷移する。Write\_\_Inbound\_\_Post 信号が状態機械 712 によって受け取られたとき、状態機械 712 は、Memory\_\_Write\_\_Request 信号 720 を発生し、IPR\_\_Ready 信号 716 をアサート解除する。状態機械 712 は、状態機械 712 がローカル・バス・アービタ 240 から許可信号 728 を受信したときポスト状態 752 からアイドル状態 750 に戻る。許可信号 728 を受信し、IPR データ 604 をメモリに書き込んだとき、状態機械 712 は、インバウンド・ポスト頭部ポインタ・レジスタ (IPHPR) 724 に対して増分信号 740 を発生し、また IPR\_\_Ready 信号 716 をアサートする。

30

状態機械 712 は、コンパレータ 734 からフル信号 738 を受信したときアイドル状態 750 からフル状態 754 に遷移する。フル信号 738 は、インバウンド・ポスト尾部ポインタ・レジスタ (IPTPR) 730 およびインバウンド・ポスト頭部ポインタ・レジスタ (IPHPR) 724 の内容がインバウンド・ポスト待ち行列 (IPQ) 530 がフルであることを示す場合、コンパレータ 734 によって生成される。フル信号 738 を受信したとき、状態機械 712 は、IPR\_\_Retry 信号 714 を ATU 218 にアサートする。

40

状態機械 712 は、フル信号 756 がフル信号 756 (すなわち、not\_\_Full) にアサート解除されたときフル状態 754 からアイドル状態 750 に遷移する。not\_\_Full 信号を受信したとき、状態機械 712 は、IPR\_\_Retry 信号 714 をアサート解除する。

コンパレータ 734 はまた、入出力プロセッサに対してローカル割込みを発生する割込み発生論理 (図示せず) に対して not\_\_Empty 信号 736 を発生する。not\_\_Empty 信号 736 を受信したときにローカル割込みを発生する論理は、当技術分野において周知である。この論理はまた、割込みレジスタを含んでおり、またソフトウェアによって制御され、割込みを選択的にマスクするマスクレジスタを含んでいる。

増分信号 740 は、IPHPR 724 に送られ、インバウンド・ポスト頭部ポインタを増

50



分する。加算器 7 2 2 は、I P H P R 7 2 4 の値 7 2 5 および I P B R 7 2 6 の値 7 2 7 を使用して、新しいインバウンド・ポスト頭部ポインタ 7 2 3 を計算する。このアドレス 7 2 3 は、ローカル・バス（すなわち、ローカル・アドレス・バス 6 3 0）を介してローカル・メモリにアクセスするためにメモリ・コントローラ 2 0 5 に送られる。

前に説明したように、M U 2 1 0 は、ローカル・アドレス・バス 6 3 0 上のアドレス 7 2 3 をアサートし、I P R 7 0 6 内で I P Q 5 3 0 の頭部にラッチされたデータ（すなわち、メッセージ・バッファのアドレス）の転送を可能にする。

第 8 A 図は、アウトバウンド検索状態機械 8 1 2、および本発明がどのようにしてホスト・プロセッサまたはバス・エージェントがアウトバウンド・ポスト待ち行列 5 2 0（O P Q）からポスト・メッセージを取り出すことを可能にするかを示す。ホスト・プロセッサまたはバス・エージェントがポスト・メッセージ・ハンドルを取り出したとき、データ（すなわち、メッセージ・バッファのアドレス）は、ローカル・データ・バス 6 0 4 を介してローカル・メモリ 2 0 6 内の O P Q 5 2 0 からアウトバウンド検索レジスタ（O R R）8 0 6 に進む。次いで、データは、O R R 8 0 6 からデータバス 8 0 8 を介して A T U 2 1 8 のアウトバウンド部分に送られる。次いで、データは、データバス 8 1 0 および P C I バス 2 0 8 を介してそれぞれホスト・プロセッサまたはバス・エージェントに送られる。状態機械 8 1 2 は、O R R 8 0 6 内の状態データを示すために O R R \_ R e a d y 8 1 4 をアサート解除する。O R R \_ R e a d y 信号 8 1 4 がアサート解除されたとき、A T U 2 1 8 は、O R R \_ R e a d y 信号 8 1 4 がアサートされるまで待ち状態を挿入する。これは、O R R 8 0 6 が有効なデータであることを示す。

M U 2 1 0 のアウトバウンド検索状態機械 8 1 2 について、第 8 B 図に示される状態図によって詳細に説明する。アウトバウンド検索状態機械 8 1 2 は、エンプティ状態 8 5 0、プリフェッチ状態 8 5 2、プライム状態 8 5 6 の 3 つの状態を有する。アウトバウンド検索状態機械 8 1 2 は、エンプティ信号 8 4 0 がアサート解除されたときエンプティ状態 8 5 0 からプリフェッチ状態 8 5 2 に遷移する。それに応答して、アウトバウンド検索状態機械 8 1 2 は、M e m o r y \_ R e a d \_ R e q u e s t 8 1 8 をローカル・バスアービトレーションユニット 2 4 0 に対してアサートし、許可信号 8 3 2 を待っている間、O R R \_ R e a d y 信号 8 1 4 をアサート解除する。許可信号 8 3 2 を待っている間、加算器 8 2 4 は、次のメッセージのアドレス（すなわち、尾部ポインタ）を計算し、このアドレスをローカル・アドレス・バス 6 3 0 上に配置する。

状態機械 8 1 2 は、許可信号 8 3 2 に基づいてプリフェッチ状態 8 5 2 からプライム状態 8 5 6 に遷移する。メモリ・コントローラ 2 0 5 は、アドレス 8 2 5 を使用し、O P Q 5 2 0 から適切なメッセージ・ハンドルを読み取る。このメッセージ・ハンドル（すなわち、ポインタ）は、ローカル・アドレス・バス 6 0 4 上に配置され、O R R 8 0 6 に転送される。次いで、状態機械 8 1 2 は、O P Q 5 2 0 からデータを O R R 8 0 6 内にラッチするラッチ \_ O R R 8 3 4 を発生し、また O P T P R 8 2 6 内に記憶された O P Q 5 2 0 の尾部ポインタを増分する I n c r e m e n t \_ O F T P R 信号 8 4 4 を発生する。このプリフェッチが完了し、新しいデータが O R R 8 0 6 内にラッチされた後、状態機械 8 1 2 は、P C I バス 2 0 8 からの他のトランザクションを完了する準備ができていることを A T U 2 1 8 に通知する前に O R R \_ R e a d y 信号 8 1 4 をアサートする。

状態機械 8 1 2 は、R e a d \_ O R P 信号 8 1 6 が発生したとき、プライム状態 8 5 6 からプリフェッチ状態 8 5 2 に遷移し、エンプティ信号 8 4 0 はアサートされない。それに応答して、状態機械 8 1 2 は、M e m o r y \_ R e a d \_ R e q u e s t 信号 8 1 8 をローカル・バス・アービタ 2 4 0 に対してアサートし、A T U 2 1 8 に対して O R R \_ R e a d y 信号 8 1 4 をアサート解除する。その結果、後のトランザクションは、プリフェッチが完了するまで O R R 8 0 6 の内容を読み取らない。

状態機械 8 1 2 は、エンプティ信号 8 4 0 がアサートされたときにアサートされた R e a d \_ O R P 信号を検知したときにプライム状態 8 5 6 からエンプティ状態 8 5 0 に遷移する。それに応答して、状態機械 8 1 2 は、プリセット信号 8 4 2 をアサートする。プリセット信号 8 4 2 により、O P Q 5 2 0 に読取りを要求しているトランザクションに O P Q

10

20

30

40

50

520がエンプティであることが通知されるようにORR806の内容がエンプティ指示に対して予約された値に設定される。

コンパレータ836がOPHPR838とOPTPR826の内容を比較し、各値が等しい場合、エンプティ信号840はアサートされる。not\_Empty\_OPQ520(すなわち、not\_Empty)は、ホスト・プロセッサ201が処理するために保留されているメッセージがあることを示す。本発明は、PCIバス仕様リリース2.0に指定されている割込み線を介してホスト・プロセッサ201に対して割込みを発生する論理(図示せず)を含んでいる。

第9A図および第9B図は、アウトバウンド解放状態機械912を示す。ホスト・プロセッサ201は、メッセージを処理した後、データ・バス904を通してPCIバス208を介してフリー・メッセージ・バッファポインタをATU218に戻し、アウトバウンド解放レジスタ(ORLSR)906内にラッチされる。次いで、フリー・メッセージ・バッファ・ハンドルは、アウトバウンド解放レジスタORLSR906からローカル・データ・バス604を介してアウトバウンド・フリー待ち行列(OFQ)510に送られる。フリー・メッセージ・バッファを解放するために、ホスト・プロセッサ201は、単に1つのバストランザクションサイクル中にそのフリー・メッセージ・バッファのアドレスをORLSR906に書き込む。

ATU218は、ORLSR\_Retry信号916およびORLSR\_Retry信号914の状態をテストする。ORLSR\_Retry信号914がアサート解除された場合、PCIトランザクション(すなわち、ORLSR906への書き込み)は完了しない。要求側プロセスは再試行を通知され、要求側プロセスは、バスの制御権を解放し、後で再度試行する。ORLSR\_Retry信号916がアサート解除された場合、ATU218は、ORLSR\_Retry信号916がアサートされるまで待ち状態を挿入する。ORLSR\_Retry信号916がアサートされたとき、ATU218は、状態機械912に対してWrite\_ORLSR信号918を発生し、データをORLSR906内にラッチする。

第9B図は、アウトバウンド解放状態機械912の状態図を示す。状態機械912は、フル状態954、アイドル状態950、ポスト状態952の3つの状態を有する。状態機械912は、フル信号940がコンパレータ936によってアサートされたときアイドル状態950からフル状態954に遷移する。このフル信号940に応答して、状態機械912は、ATU218に対してORLSR\_Retry信号914をアサートする。ORLSR\_Retry信号914が発生したとき、ORLSR906に対する書き込みトランザクションを示すプロセスは、後で再度試行するよう通知される。

状態機械912は、フル信号940がアサート解除されたときフル状態954からアイドル状態950に遷移する。アウトバウンド・フリー待ち行列OFQ510がフルでない場合、状態機械912は、ORLSR\_Retry信号914をアサート解除する(すなわち、OFQ510内に追加のフリー・メッセージ・ハンドル用の空きがある)。

状態機械912は、ATU218からWrite\_ORLSR信号918を受信したときアイドル状態950からポスト状態952に遷移する。Write\_ORLSR信号918はまた、フリー・メッセージ・ハンドルをORLSR906内にラッチする役目をする。Write\_ORLSR信号918がアサートされたことに応答して、状態機械912は、ローカル・バス・アービタ240に対してMemory\_Write\_Request信号918をアサートする。アービタからの許可信号932を待つ。次のフリー・メッセージ・ハンドルを書き込むべきOFQ510内の次の位置を計算する。状態機械912はまた、後のトランザクションがORLSR906内にラッチされているデータを上書きするのを防ぐためにORLSR\_Retry信号916をアサート解除する。

状態機械912は、ローカル・バス・アービタ240から許可信号932を受信したときポスト状態952からアイドル状態950に遷移する。それに応答して、アウトバウンド解放状態機械912は、Increment\_OFHPR信号944を介してOFHPR926内の頭部ポインタを増分する。状態機械912はまた、すでにORLSR906の

10

20

30

40

50

内容を記憶していること、および次のフリー・メッセージ・ハンドルを記憶すべきOFQアドレスを計算したことをATU218に示すORLSR\_\_作業可能信号916をアサートし、ORLSR906への次の書込みの準備ができる。

要約すると、ホスト・プロセッサは、そのハンドルをORLSR906に書き込むことによってOFQ510に対してフリー・メッセージ・バッファを解放する。OFQ510がフルである場合、要求側プロセスは、後で再試行するよう通知される。OFQ510がフルでない場合、フリー・メッセージ・バッファのハンドルはORLSR906内にラッチされる。次いで、状態機械912は、ローカル・バスへのアクセス権を得るためにローカル・バス・アービタ240からの許可信号932を待つ。ローカル・バスの制御権が付与された後、状態機械912は、前に計算した頭部ポインタ/アドレスによって示された位置においてORLSR906内にラッチされたデータをOFQ510に転送する。

10

以上、遠隔プロセスがセマフォの使用またはバスのロックなしにメッセージ・バッファを割り振り、次いでメッセージ・バッファを作業待ち行列にポストすることを可能にする方法および装置について説明した。

さらに、作業待ち行列からメッセージを取り出し、メッセージがホスト・プロセッサによって処理された後でメッセージをフリー待ち行列に解放する方法および装置について説明した。

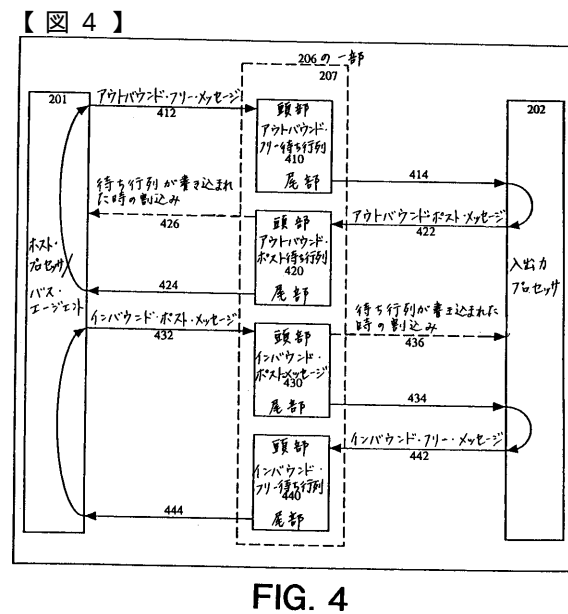
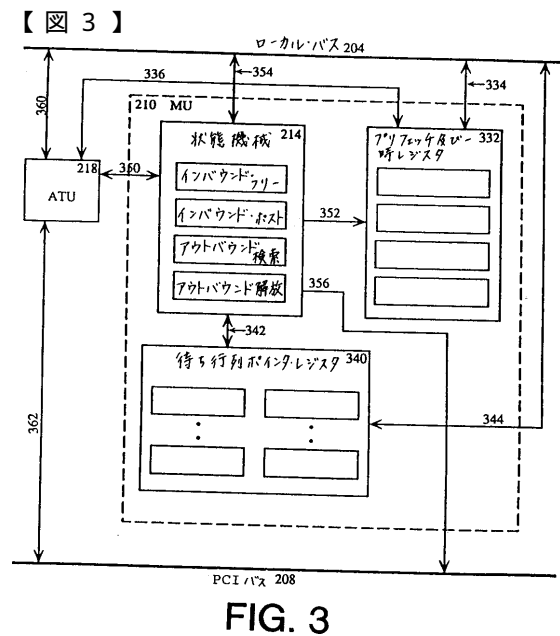
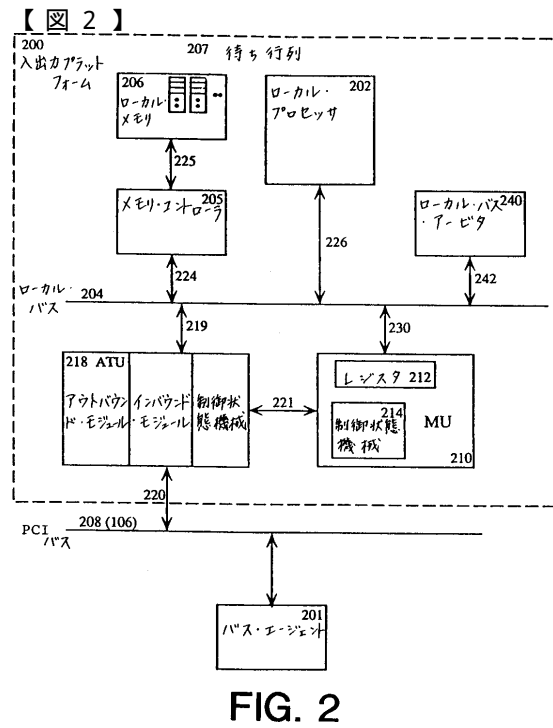
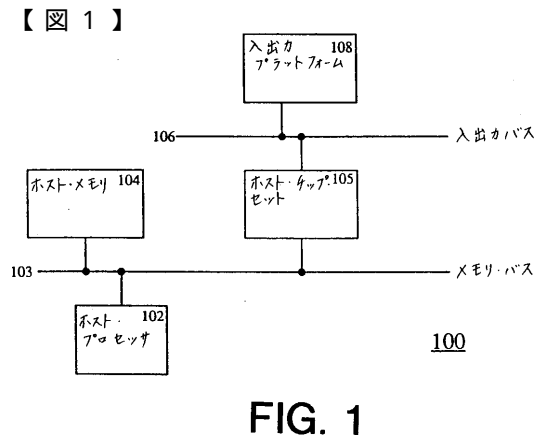
本発明はまた、スケーラビリティ、フレキシビリティ、および他のプラットフォームとの互換性を提供する。例えば、上述のように、インバウンドメッセージ待ち行列を含むすべてのプラットフォームは、プロセッサ間メッセージを容易に送ることができる。インバウンドメッセージ待ち行列を実施しない他のプラットフォームとの互換性については、プラットフォームのハードウェアを修正することなくアウトバウンド・メッセージ待ち行列がそのプラットフォームに同等の機能を供給する。さらに、本発明は、他のプロセッサがコンピュータシステム内に存在することを明確に知ることなく、他のプラットフォームが1つのプラットフォームのインバウンド待ち行列を同時に使用できるという抽象を可能にする。

20

したがって、本発明は、非対称マルチプロセッサ・システム内のプロセッサに対するハードウェア修正を必要とすることなくプロセッサ間で非常に効率的な形でメッセージを直接送る方法および装置を提供する。

以上、本発明について、その特定の例示的な実施形態に関して説明した。しかしながら、下記の請求の範囲に記載されている本発明のより広い精神および範囲から逸脱することなく本発明に様々な修正および変更を加えることができることが明らかであろう。したがって、明細書および図面は、限定的なものではなく、例示的なものと考えられたい。

30



【図 5】

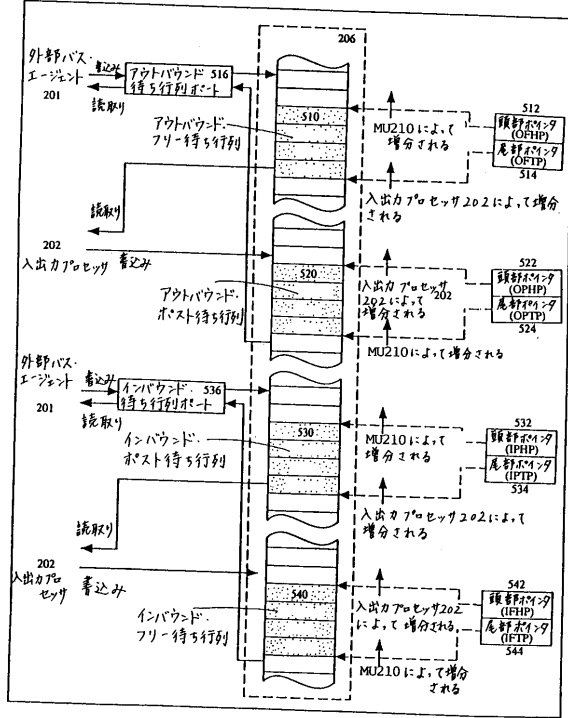


FIG. 5

【図 6 A】

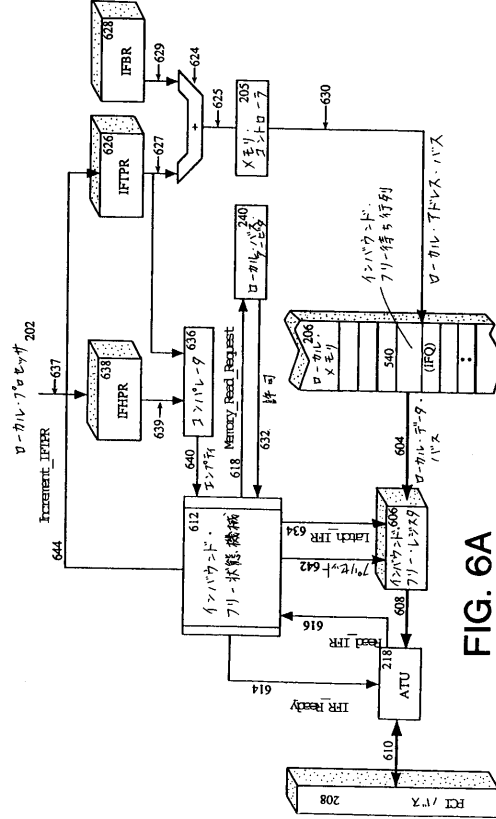


FIG. 6A

【図 6 B】

インバウンド・フリー状態機械の状態図

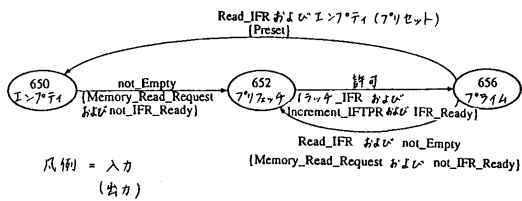


FIG. 6B

【図 7 A】

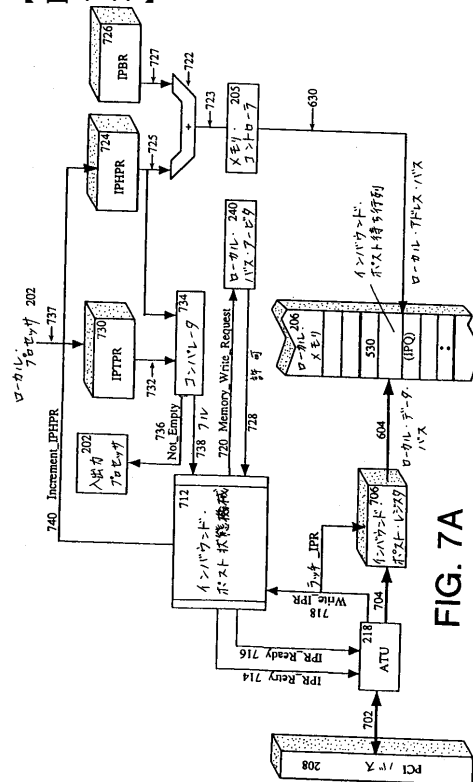
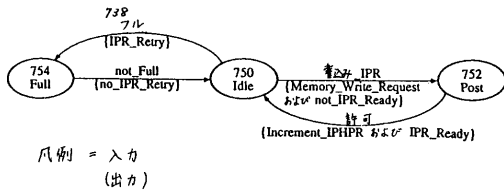


FIG. 7A

【 図 7 B 】

インバウンド・ポスト状態機械の状態図



**FIG. 7B**

【 図 8 A 】

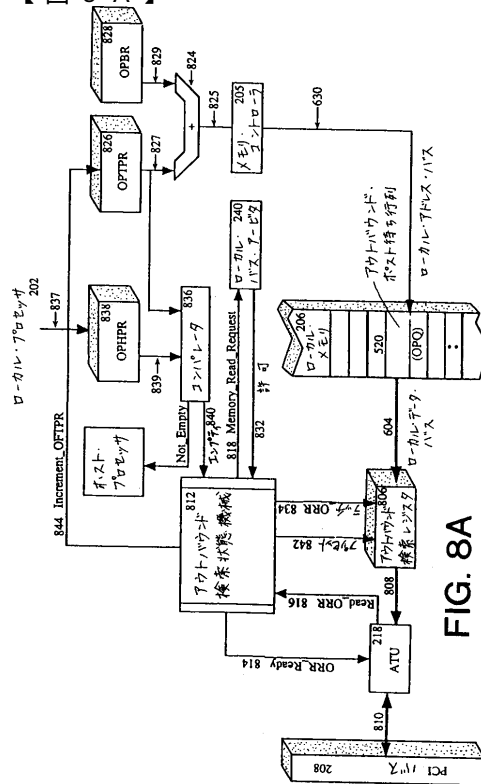
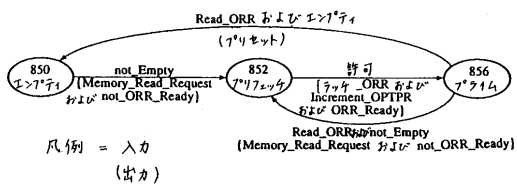


FIG. 8A

【 図 8 B 】

アウトバウンド検索状態機械の状態図



**FIG. 8B**

【 図 9 A 】

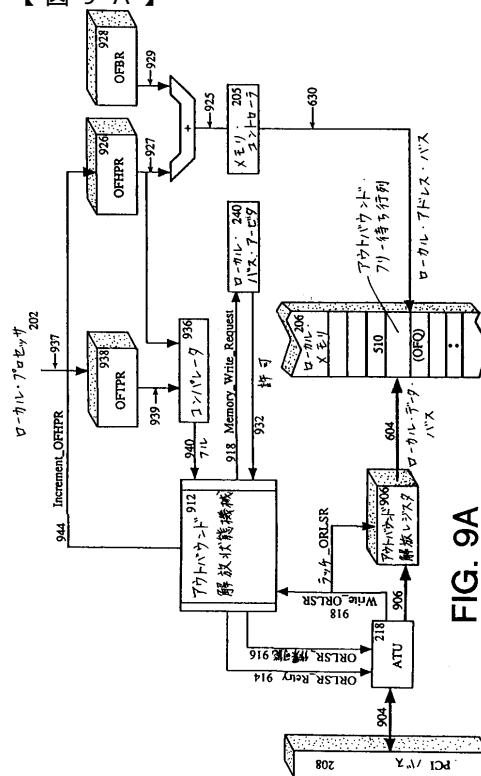


FIG. 9A

## 【図 9 B】

アウトバウンド解放状態機械の状態図

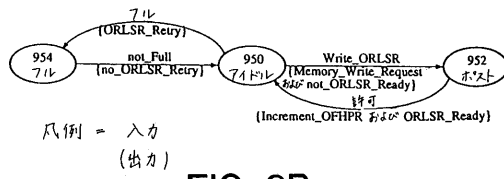


FIG. 9B

---

フロントページの続き

(72)発明者 デイビス, バリー

アメリカ合衆国・ 8 5 2 2 4 ・アリゾナ州・チャンドラー・ウエスト サラゴサ ストリート・ 7  
4 0

(72)発明者 フュトラル, ウィリアム・ティ

アメリカ合衆国・ 9 7 2 2 9 ・オレゴン州・ポートランド・ノースウエスト エルク ラン ドラ  
イブ・ 1 7 7 1 5

(72)発明者 ガーバス, エリオット

アメリカ合衆国・ 8 5 2 5 7 ・アリゾナ州・スコッツデイル・ノース ハイデン ロード・ 2 7 0  
0 ・ 3 1 0 6 番

審査官 鳥居 稔

(56)参考文献 特開平 0 5 - 2 8 2 1 6 6 ( J P , A )

特開平 0 2 - 1 5 8 8 5 8 ( J P , A )

特開平 0 5 - 2 6 8 2 9 1 ( J P , A )

特表平 1 1 - 5 1 3 1 5 0 ( J P , A )

(58)調査した分野(Int.Cl. , D B 名)

G06F 13/10

G06F 13/38

G06F 15/16 - 177