

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2004-54949

(P2004-54949A)

(43) 公開日 平成16年2月19日(2004.2.19)

(51) Int. Cl.⁷

G06F 13/12
G06F 12/08
G06F 12/10
G06F 15/17

F I

G06F 13/12 330G
G06F 12/08 531B
G06F 12/08 555
G06F 12/10 501Z
G06F 12/10 555

テーマコード(参考)

5B005
5B014
5B045

審査請求 未請求 請求項の数 10 O L (全 12 頁) 最終頁に続く

(21) 出願番号 特願2003-273509(P2003-273509)
(22) 出願日 平成15年7月11日(2003.7.11)
(31) 優先権主張番号 10/200247
(32) 優先日 平成14年7月23日(2002.7.23)
(33) 優先権主張国 米国(US)

(71) 出願人 503003854
ヒューレット・パカード デベロップメント カンパニー エル. ピー.
アメリカ合衆国 テキサス州 77070
ヒューストン 20555 ステイト
ハイウェイ 249
(74) 代理人 110000039
特許業務法人アイ・ピー・エス
(72) 発明者 デベンドラ・ダス・シャルマ
アメリカ合衆国・カリフォルニア州・サン
タクララ・アカシアコート2043
Fターム(参考) 5B005 KK13 MM01
5B014 FB04 FB05
5B045 BB54 DD12 EE03 EE08 KK06

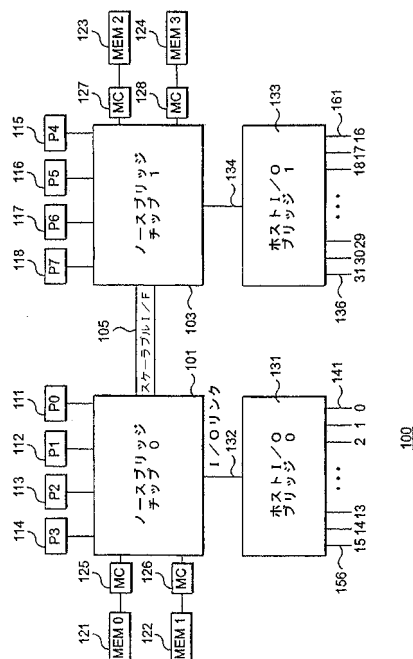
(54) 【発明の名称】 単一入出力ハブ下の複数ハードウェアパーティション

(57) 【要約】

【課題】 コンピュータ資源を細粒分割する方法および機構を提供する。

【解決手段】 複数のI/Oパスすなわちローブは、それぞれのローブがコンピュータシステム内の異なるパーティションに潜在的に属するように分割される。一実施形態では、到来キューおよび送出キューは、1つのローブに対するトランザクションが別のローブに干渉しないように設計される。各共通キューは、仮想キュー/バッファの集まりとして扱われる。コンポーネントは、イーサネット(登録商標)カード等の入出力(I/O)デバイスや他の入出力(I/O)デバイスである。キャッシュまたはTLBを有するコヒーレントI/Oキャッシュは、それぞれが異なるパーティションに潜在的に属する複数のローブで共有することが可能である。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

コンピュータシステム(100)においてハードウェアを分割する装置であって、前記コンピュータシステムは複数の入出力(I/O)パス(141~156、161~176)を備え、該I/Oパスはそれぞれホスト入出力(I/O)ブリッジ(131、133)に結合され、少なくとも1つのパーティションが前記複数のI/Oパスを包含し、装置は

アウトバウンドキュー(201、203、211~218)階層であって、

単一のI/Oパスに割り当てられる最下位階層の最下位アウトバウンドキュー(211~218)、および

1つまたは複数の上位レベルの上位アウトバウンドキュー(203)であって、前記1つまたは複数の上位レベルのアウトバウンドキュー中の上位アウトバウンドキューは、1つまたは複数の仮想キューセクションを含み、各最下位アウトバウンドキューは、前記1つまたは複数の仮想キューのうちの指定された仮想キューに対するデータの送受信を行う、1つまたは複数の上位レベルの上位アウトバウンドキュー(203)、を含むアウトバウンドキュー(201、203、211~218)階層と、

パーティション識別(ID)を特定のI/Oパス上で発信されたランザクションに割り当てるI/Oパスコントローラであって、前記パーティションIDは、前記特定のI/Oパスが割り当てられたパーティションに対応するI/Oパスコントローラとを備える装置。

【請求項 2】

前記ホストI/Oブリッジにキャッシュ(235)をさらに備え、該キャッシュはパーティションIDセクション(243)を含む請求項1記載の装置。

【請求項 3】

前記ホストI/Oブリッジは、キャッシュコントローラ(233)を含み、

該キャッシュコントローラは、大域共有メモリに送られない要求について読み出し/書き込み要求のアドレス情報とパーティションIDとを比較する

請求項2記載の装置。

【請求項 4】

前記キャッシュコントローラ(233)は、大域共有メモリに送られる要求のパーティションIDおよびアドレスを推測する

請求項3記載の装置。

【請求項 5】

変換ルックアサイドバッファ(TLB)(231)をさらに備え、該TLBはパーティションIDセクションを含み、特定のI/Oパスの前記パーティションIDが記録される請求項1記載の装置。

【請求項 6】

前記ホストI/Oブリッジは複数のレジスタ(237)を含み、該複数のレジスタのうちの1つまたは複数は、パーティションのオペレーティングシステムによってアクセスされ、他のパーティションのコンポーネントは、前記複数のレジスタのうちの前記1つまたは複数へのアクセスが拒絶される

請求項1記載の装置。

【請求項 7】

コンピュータシステムにおいてハードウェアを分割する方法であって、

ダイレクトメモリアクセス(DMA)の場合、読み出しまたは書き込み動作であるDMA要求を発した(initiated)コンポーネントのパーティション識別(ID)に対応するパーティションIDを生成すること(310)と、

リコール要求の場合、前記コンポーネントの前記アドレスおよび前記パーティションIDを、要求されたラインのアドレスおよびパーティションIDと比較することと、

10

20

30

40

50

割り込みトランザクションおよび非コヒーレントメモリ読み出し/書き込みトランザクションの場合、前記発信コンポーネントの前記パーティションIDを前記トランザクションに添付することとを含む方法。

【請求項8】

前記コンピュータシステムにおける各入出力(I/O)処理パスを特定することと、前記特定されたI/Oパスをそれぞれ別個のパーティションに分割することとをさらに含む請求項7記載の方法。

【請求項9】

前記コンピュータシステムにおける各入出力(I/O)処理パスを特定することと、前記特定されたI/Oパスのうちの2つ以上のI/Oパスを別個のパーティションに分割することとをさらに含む請求項7記載の方法。

10

【請求項10】

カスケード階層になった出力キューを提供することをさらに含み、第1レベルの出力キューは2つ以上の仮想キューを含み、該仮想キューはそれぞれ第2レベルの出力キューに対応し、該第2レベルの出力キューにおいてストールすると、前記関連する第1レベルの仮想キューのみがストールする請求項7記載の方法。

【発明の詳細な説明】

20

【技術分野】

【0001】

本技術分野は、多重かつ/または冗長サブシステムおよびコンポーネントを有するコンピュータシステムである。

【背景技術】

【0002】

現在のマルチプロセッサコンピュータシステムは、コンピュータシステムのハードウェア資源およびソフトウェア資源が各種パーティションに分けられるように通常分割される。

たとえば、マルチプロセッサコンピュータシステムでは、プロセッサは2つ以上のプロセッサグループに分けることができる。

30

たとえば、メモリバス等の資源もプロセッサパーティションに含めることができる。

【発明の開示】

【発明が解決しようとする課題】

【0003】

キャッシュまたはI/Oトランザクションルックアサイドバッファ(TLB)を有するコヒーレントな入出力(I/O)ハブチップは、複数のI/Oパスで共有することが可能である。

現在のソリューションでは、同じキャッシュの下にあるすべてのロープが1つのパーティションに属するものと仮定している。

40

現在のシステムではまた、I/Oハブ全体を1つのパーティションに配置している。

これでは、多すぎる資源(たとえば、I/Oカードおよび接続)が1つのパーティションにロックされることになる。

【課題を解決するための手段】

【0004】

本明細書において、コンピュータ資源を細粒分割する方法および機構を記載する。

特に、複数のI/Oパスすなわちロープは、それぞれのロープがコンピュータシステム内の異なるパーティションに潜在的に属するように分割される。

一実施形態では、到来キューおよび送出キューは、1つのロープに対するトランザクションが別のロープに干渉しないように設計される。

50

【0005】

各共通キューは、仮想キュー/バッファの集まりとして扱われる。

コンポーネントは、イーサネット(登録商標)カード等の入出力(I/O)デバイスや他の入出力(I/O)デバイスである。

しかし、本方法および本機構は、I/Oデバイス以外のコンピュータコンポーネントによる使用にも適合することができる。

【0006】

キャッシュまたはTLBを有するコヒーレントI/Oキャッシュは、それぞれが異なるパーティションに潜在的に属する複数のロープで共有することが可能である。

これまでのソリューションでは、キャッシュの下にあるすべてのロープが同じパーティションに属するものと仮定していた。 10

一実施形態では、複数のパーティションが1つのキャッシュ下に共存するため、パーティションの粒度が大きくなる。

【発明の効果】

【0007】

本発明によれば、コンピュータ資源を細粒分割する方法および機構が提供される。

【発明を実施するための最良の形態】

【0008】

詳細な説明では、同様の符号が同様の要素を指す添付図面を参照する。

本明細書において、個々の処理パスすなわちロープのレベルにおいてコンピュータハードウェアコンポーネントすなわち要素を分割することができる方法および機構を記載する 20

実施形態では、到来キューおよび送出キューは、1つのロープに対するトランザクションが別のロープに干渉しないように設計される。

各共通キューは、仮想キュー/バッファの集まりとして扱われる。

【0009】

コンピュータコンポーネントは、イーサネット(登録商標)カード等の入出力(I/O)デバイスや他の入出力(I/O)デバイスを含む。

しかし、本方法および本機構は、I/Oデバイス以外のコンピュータコンポーネントによる使用にも適合することができる。 30

【0010】

図1は、1つのI/Oハブの下で複数のハードウェアパーティションを使用して細粒分割を提供する例示的なコンピュータアーキテクチャ100のブロック図である。

アーキテクチャ100は、スケーラブルインタフェース105に共に結合されたノースブリッジチップ101および103を備える。

【0011】

ノースブリッジチップ101には、プロセッサ111~114と、メモリバス121および122と、が結合される。

ノースブリッジチップ103には、プロセッサ115~118と、メモリバス123および124と、が結合される。 40

メモリバス121~124は、メモリコントローラ125~128の制御下で動作することができる。

【0012】

アーキテクチャ100内の各ノースブリッジは、I/Oリンクすなわちシステムバスインタフェースを通してホストI/Oブリッジに接続することができる。

図1に示すように、ノースブリッジ101は、I/Oリンク132を通してホストI/Oブリッジ131に接続し、ノースブリッジチップ103は、I/Oリンク134を通してホストI/Oブリッジ133に接続する。

【0013】

ホストI/Oブリッジ131には、16本のロープ141~156が接続される。 50

ホストI/Oブリッジ133には、別の16本のロープ161~176が接続される。

図1に示すアーキテクチャ100のコンピュータコンポーネントすなわち要素の構成は、例としてのみのものであり、アーキテクチャ100は、図示するコンポーネントのそれぞれをより多く、またはより少なく備えてもよい。

【0014】

図2は、図1に示すホストI/Oブリッジ131等のホストI/Oブリッジのブロック図である。

ホストI/Oブリッジ131は、I/Oリンクすなわちシステムバスインタフェース132を有して示される。

ホストI/Oブリッジ131は、I/Oユニット137および139等の1つまたは複数のI/Oユニットを備えることができる。 10

【0015】

I/Oユニット137および139は、キャッシュ、変換ルックアサイドバッファ(TLB)、キュー、レジスタ、および制御要素等、他の様々なコンピュータコンポーネントを制御するかまたは備えることができる。

ホストI/Oブリッジ131に結合された16本のロープは、図示のようにI/Oユニット137および139の間で分割することができる。

個々のロープ141~156は、他のコンピュータコンポーネントに接続しうる。

【0016】

かかるロープはそれぞれ、ロープコントローラ(図示せず)の制御下で動作することができる。 20

図2に示す例では、ロープ141(ロープ0)は下位バスチップ181に接続し、下位バスチップ181はPCIバス182に接続する。

ロープ147および148(ロープ7および8)は下位バスチップ183に接続し、下位バスチップ183はPCI-Xバス184に接続する。

図2に示す他のロープは、他のコンピュータコンポーネントに接続することができる。

【0017】

図1および図2に示す各種コンピュータコンポーネントは、分割可能である。

分割は、ソフトウェアおよびハードウェアにおいて実現することができる。

図3は、かかる1つの分割方式を示したものである。 30

図示の例では、コンピュータコンポーネントをハードウェアにおいて2つのパーティションに分割することができ、パーティション0がプロセッサ111~114、メモリバス121および122、ならびにホストI/Oブリッジ131を含む。

【0018】

パーティション1は、プロセッサ115~118、メモリバス123および124、ならびにホストI/Oブリッジ133を含むことができる。

かかる分割は、粗粒分割と呼ぶことのできる例であり、パーティション0および1のそれぞれの内の各種コンポーネントは、パーティションの境界を越えてコンポーネントの共有が発生しないように厳密に分離されている。

【0019】 40

さらに、分割は、ホストI/Oブリッジレベル等の粗いレベルにおいて定義され、ホストI/Oブリッジに接続される個々のロープは、分割方式に明示的に包含されない。

言い換えれば、特定のホストI/Oブリッジに接続されたすべてのロープは、そのホストI/Oブリッジと同じパーティション内に含まれる。

同様に、すべてのプロセッサ(および複数のコアを有するプロセッサ内のコア)は、たとえば、ノースブリッジチップに結びついたパーティションと同じパーティションに包含される。

【0020】

代替の分割方式(図示せず)、すなわち細粒分割では、いくつかのコンピュータコンポーネントの共有が許される。 50

たとえば、最小計算粒度がプロセッサ（またはプロセッサ内のコア）レベルであり、メモリがページ単位で割り当てられ、I/O資源をロープの粒度まで下げて割り当てることができる。

この細粒分割は、コンピュータパフォーマンスおよびコスト削減に関して利点をもたらすが、個々のパーティションが互いに干渉しないように保証する機構および方法を必要とする。

【0021】

図4は、キュー中のデータの順方向進行を保証することによりパーティション干渉を回避する機構を示す。

10
ホストI/Oブリッジ131は、システムバスインタフェース132におけるアウトバウンド共通キュー201（キュー0）を有して示される。

キュー201は、アウトバウンドキュー203および204に供給することができる。

アウトバウンドキュー203（キュー1）は、共通キューとして動作し、アウトバウンドキュー211～218（キュー2～9）に供給することができる。

【0022】

アウトバウンドキュー211～218はそれぞれ図示のように特定のロープ（すなわち、ロープ149～156）に割り当てられている。

ロープ149～156のうちの本、たとえばロープ156に、アウトバウンドキュー218が空にならないという問題が生じた場合、上流すなわち共通キュー203およびその他のロープ149～155が影響を受け、故障することがある。

20

【0023】

さらに、他のパーティションが悪影響を受ける場合がある。

故障したロープ156からの逆圧が他のコンピュータコンポーネントに影響を与えないように、共通キュー203および201をそれぞれ独立したキューのグループとして扱うことができる。

たとえば、共通キュー203（キュー211～218に共通）を、下流キュー211～218それぞれに1つずつ、8個の独立したキューのグループとして扱うことができる。

【0024】

一実施形態では、キュー203は、キュー211～218それぞれにアドレス範囲を割り当てて、内部分割することができる。

30

共通キューの分割に対する代替として、タイムアウト機構220をキュー211～218に結合することができ、タイムアウト機構220は、関連するロープが故障した場合、指定された時間後にキュー211～218のいずれかの内容を破棄する。

そして、本明細書に述べる実施形態では、共通キューが情報およびデータを流し、関連するコンポーネントの考えられる付随的な故障による停滞を回避することができる。

【0025】

図5は、パーティション干渉を回避する他の機構を示す。

図5では、ホストI/Oブリッジ131のI/Oユニット137が、パーティションTLB231、キャッシュコントローラ233、キャッシュ235、パーティションレジスタ237、およびキュー239を備えて示される。

40

【0026】

キャッシュ235は、トランザクションターゲットを特定するセクション、および特定のトランザクションに関連する他の情報を提供するセクションを含む。

図示のように、キャッシュ235は、Cデータセクション241、Cタグセクション243、およびCスタットセクション245を含むことができる。

【0027】

Cデータセクション241は、特定のキャッシュラインに関連するデータを含む。

Cタグセクション243は、特定のキャッシュラインについてのアドレス情報を記録する。

アドレス情報は、パーティションID、特定のロープ、およびターゲットアドレスを含

50

むことができる。

C スタットセクション 2 4 5 は、特定のキャッシュラインに固有の情報を記録する。

【0028】

たとえば、C スタットセクション 2 4 5 は、特定のキャッシュラインが無効であることを示すことができる。

T L B 2 3 1、キャッシュコントローラ 2 3 3、キャッシュ 2 3 5、パーティションレジスタ 2 3 7、およびキュー 2 3 9 は連絡して、図 1 に示すアーキテクチャ 1 0 0 により表されるコンピュータシステムの動作中にパーティションが干渉しないように保証する。

上述したように、ホスト I / O ブリッジ 1 3 1 (本例では) がパーティション 0 に割り当てられる。

【0029】

ホスト I / O ブリッジ 1 3 1 のすべてのコンポーネントもまた、パーティション 0 に割り当てることができる。

ロープ 1 4 9 ~ 1 5 6 もパーティション 0 に割り当ててもよく (粗粒分割)、またはロープ 1 4 9 ~ 1 5 6 を他のパーティションに割り当ててもよい。

一実施形態において、ロープ 1 4 9 ~ 1 5 6 をそれぞれ固有のパーティションに割り当てることができる。

かかる各パーティションは、たとえば 1 ビットまたは 2 ビットのデータフィールドであり得るパーティション ID によって識別することができる。

【0030】

キャッシュ 2 3 5 は、トランザクションが各パーティション内の正しいコンポーネントに確実に送信されるように、パーティションのアドレスと併せてパーティション ID を使用することができる。

たとえば、直接メモリアクセス (DMA) 読み出し (または書き込み) 要求がロープに到着すると、ロープコントローラが、そのロープに関連するパーティション ID を生成する。

【0031】

次いで、アドレスおよびパーティション ID がキャッシュコントローラ 2 3 3 に送られる。

キャッシュコントローラ 2 3 3 は、パーティション ID およびアドレスを使用して、要求されたラインがキャッシュ 2 3 5 内に存在するかどうかを判断する。

【0032】

キャッシュミス (すなわち、要求されたラインが存在しない) がある場合、キャッシュコントローラ 2 3 3 は、フェッチ要求を適切なメモリコントローラに (たとえば、メモリバス 1 2 1 または 1 2 2 に (図 1 参照)) 送る。

フェッチ要求は、キャッシュラインアドレス、ならびに読み出し (書き込み) 要求を生成したロープのパーティション ID を含む。

【0033】

リコール要求とは、あるコンピュータコンポーネントがキャッシュラインの制御を有し、別のコンポーネントが同じラインを要求する場合に、キャッシュコントローラ 2 3 3 によって開始されるトランザクションである。

リコール要求の場合、キャッシュコントローラ 2 3 3 は、キャッシュ 2 3 5 内のラインへのアドレスと併せて要求しているコンポーネントのパーティション ID を比較する。

【0034】

この規則への唯一の例外は、大域共有メモリ (GSM)、すなわち複数のパーティションによって共有されるメモリであり得る。

一実施形態では、キャッシュコントローラ 2 3 3 は、<位置 ID、アドレス> が GSM であることを推測し、リコールの場合、そのパーティション ID を GSM の他のパーティション ID で置換する。

代替の実施形態では、メモリコントローラは複数のリコールを生成し、かかる各リコー

10

20

30

40

50

ルが G S M に属するすべてのパーティションのパーティション I D を有する。

【 0 0 3 5 】

割り込み、または非コヒーレントメモリ読み出し（書き込み）の場合、同じパーティション I D が添付される。

プロセッサ I / O 書き込みまたは読み出しの場合、チェックが行われてプロセッサが、書き込みを行うロープのパーティションに属するかどうかを調べる。

【 0 0 3 6 】

キャッシュフラッシュまたは I / O T L B ページの場合、信頼のおけるプロセッサあるいはそのパーティションのプロセッサがラインフラッシュ（または T L B エントリのページ）を許される。

10

したがって、たとえば、プロセッサ 1 1 1 ~ 1 1 4 は、ホスト I / O ブリッジ 1 3 1（図 1 参照）と同様にパーティション 0 に割り当てられていることから（本例では）、プロセッサ 1 1 1 ~ 1 1 4（図 1 参照）のみが I / O T L B 2 3 1 内のエントリのページを許される。

これにより、別のパーティション（たとえば、パーティション 1）がパーティション 0（またはコンピュータアーキテクチャ 1 0 0 の任意の他のパーティション）に属するキャッシュラインをフラッシュ（または T L B エントリをページ）しないよう保証される。

【 0 0 3 7 】

図 6 は、図 1 の処理パス（ロープ）のうちの 1 本に接続されたコンポーネントによって開始された直接メモリアクセス（D A M）読み出し動作 3 0 0 を示すフローチャートである。

20

動作 3 0 0 はブロック 3 0 5 において開始する。

ブロック 3 1 0 において、ロープ 1 5 6 に関連するコンポーネントが D M A 読み出し要求を生成する。

【 0 0 3 8 】

ブロック 3 1 5 において、ロープ 1 5 6 のロープコントローラが読み出し要求を受信する。

ブロック 3 2 0 において、ロープコントローラは、ロープ 1 5 6 に関連するパーティションのパーティション識別を生成する。

次いで、ブロック 3 2 5 において、ロープコントローラは、アドレスおよびパーティション I D をキャッシュコントローラ 2 3 3 に送る。

30

【 0 0 3 9 】

ブロック 3 3 0 において、キャッシュコントローラ 2 3 3 が、キャッシュラインがキャッシュ 2 3 5 に存在するかどうかを判断する。

キャッシュコントローラ 2 3 3 は、ラインがキャッシュ 2 3 5 に存在すると判断する場合、ブロック 3 3 5 において、要求されたラインをキャッシュ 2 3 5 から検索する。

【 0 0 4 0 】

キャッシュコントローラ 2 3 3 は、キャッシュミスが発生したと判断する場合、ブロック 3 4 0 において、フェッチ要求を適切なメモリコントローラに送る。

フェッチ要求は、キャッシュラインアドレスおよびパーティション I D を含む。ブロック 3 3 5 あるいは 3 4 0 のいずれか一方の後に、読み出し動作 3 0 0 が終了する。

40

【産業上の利用可能性】

【 0 0 4 1 】

本発明は、多重かつ／または冗長サブシステムおよびコンポーネントを有するコンピュータシステムに利用可能である。

【図面の簡単な説明】

【 0 0 4 2 】

【図 1】多重／冗長コンポーネントを用いたコンピュータアーキテクチャの図である。

【図 2】図 1 のコンピュータアーキテクチャにおいて使用されるホスト I / O ブリッジの図である。

50

【図3】図1のアーキテクチャで使用することが可能な分割方式を示す図である。

【図4】あるパーティションが別のパーティションによる干渉を受けないようにするために使用されるキューシステムの図である。

【図5】図2のホストI/Oブリッジのコンポーネントの図である。

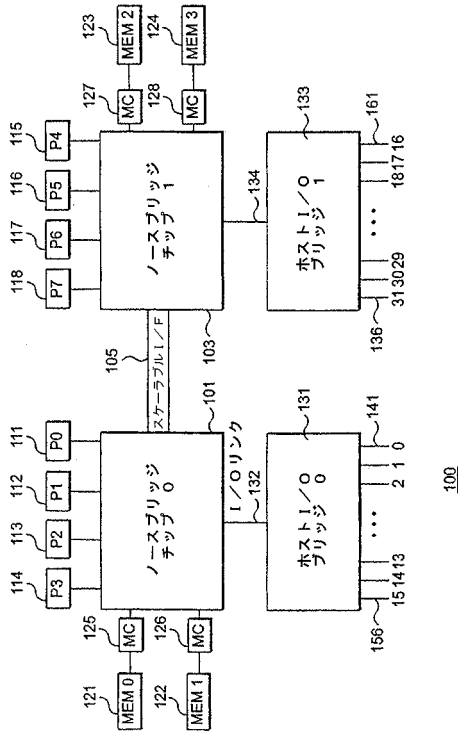
【図6】図1のコンピュータアーキテクチャの動作を示すフローチャートである。

【符号の説明】

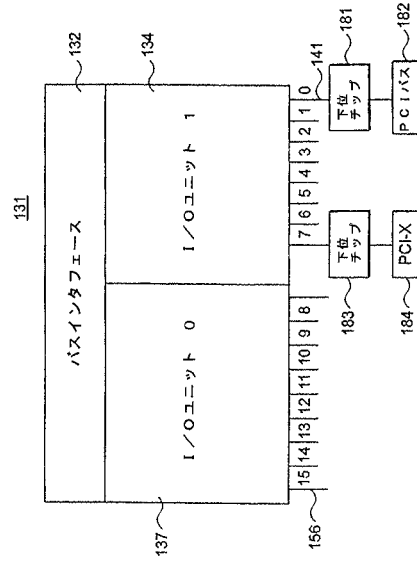
【0043】

101, 102 . . . ノースブリッジチップ、	
105 . . . スケーラブルインタフェース、	
111 ~ 118 . . . プロセッサ、	10
121 ~ 124 . . . メモリバス、	
125 ~ 128 . . . メモリコントローラ、	
131, 133 . . . ホストI/Oブリッジ、	
132, 134 . . . I/Oリンク、	
137, 139 . . . I/Oユニット、	
141 ~ 156, 161 ~ 176 . . . ロープ、	
181 . . . 下位バスチップ、	
182 . . . PCIバス、	
183 . . . 下位バスチップ、	
184 . . . PCI-Xバス、	20
201 . . . 共通キュー	
203, 204, 211 ~ 218 . . . アウトバウンドキュー	
220 . . . タイムアウト機構	
231 . . . パーティションTLB	
233 . . . キャッシュコントローラ、	
235 . . . キャッシュ、	
237 . . . パーティションレジスタ、	
239 . . . キュー、	
241 . . . Cデータセクション、	
243 . . . Cタグセクション、	30
245 . . . Cスタットセクション、	

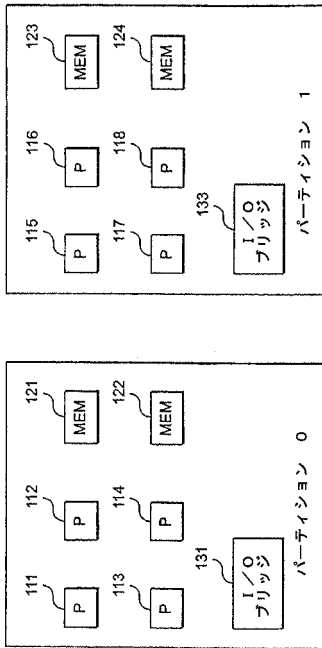
【 図 1 】



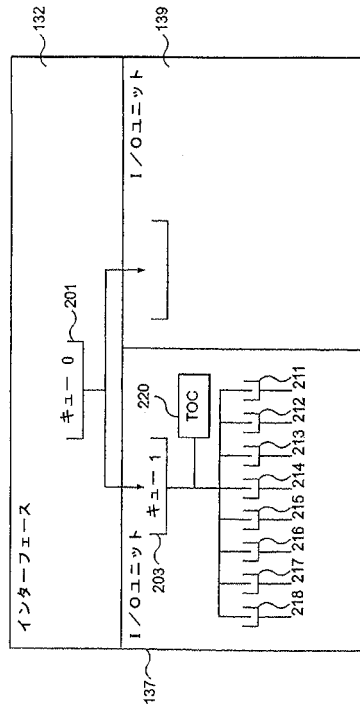
【 図 2 】



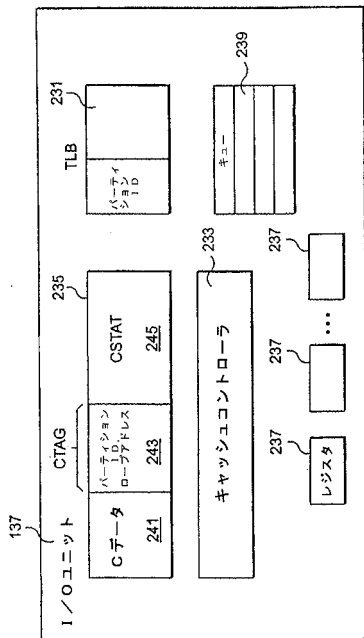
【 図 3 】



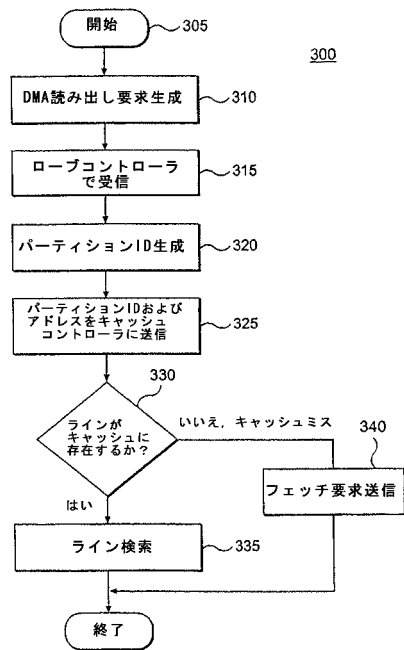
【 図 4 】



【図5】



【図6】



フロントページの続き

(51) Int.Cl.⁷

F I

テーマコード(参考)

G 0 6 F 15/177 6 7 0 C