

(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2004-30578
(P2004-30578A)

(43) 公開日 平成16年1月29日(2004.1.29)

(51) Int.Cl. ⁷	F I	テーマコード (参考)
GO6F 13/14	GO6F 13/14 33OE	5B014
GO6F 11/20	GO6F 11/20 31OA	5B034
GO6F 13/00	GO6F 13/00 301P	5B083

審査請求 未請求 請求項の数 10 O L (全 10 頁)

(21) 出願番号 特願2003-63486 (P2003-63486)	(71) 出願人 398038580
(22) 出願日 平成15年3月10日 (2003.3.10)	ヒューレット・パッカード・カンパニー
(31) 優先権主張番号 10/092603	HEWLETT-PACKARD COMPANY
(32) 優先日 平成14年3月8日 (2002.3.8)	アメリカ合衆国カリフォルニア州パロアルト
(33) 優先権主張国 米国 (US)	ハノーバー・ストリート 3000
(特許庁注：以下のものは登録商標) イーサネット	(74) 代理人 100099623 弁理士 奥山 尚一
	(74) 代理人 100096769 弁理士 有原 幸一
	(74) 代理人 100107319 弁理士 松島 鉄男

最終頁に続く

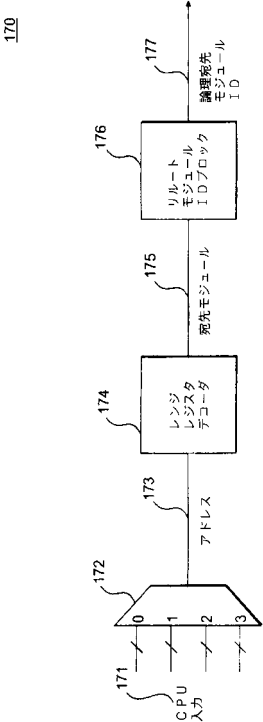
(54) 【発明の名称】 仮想入出力の相互接続メカニズム

(57) 【要約】

【課題】故障したコンポーネントにアドレス指定されたトランザクションを、代替のコンポーネントが要求できるようにする。

【解決手段】入力を多重化して、処理ユニット関連トランザクションに関連するアドレス173を生成するマルチプレクサ172と、前記アドレスを受信し、該アドレスに関連するトランザクションを受信するモジュールの宛先アドレス175を提供するレンジレジスタデコーダ174と、前記宛先アドレスを受信するリルートモジュールIDブロック176とを備えるアドレスデコードブロックと、前記コンピュータシステム中の1以上のオリジナルモジュールのアドレスを提供するオリジナルモジュールID182と、前記コンピュータシステム中の代替モジュールの論理的な宛先モジュールIDを提供するリマッピングモジュールID183とを備えるリルートモジュールIDブロックとを含む、仮想入出力の相互接続メカニズムを提供する。

【選択図】 図4



【特許請求の範囲】

【請求項 1】

1 以上のブリッジユニットによって連結された複数の入出力装置および複数の処理ユニットを備えるコンピュータシステムに使用するための仮想入出力の相互接続メカニズムであって、該相互接続メカニズムはアドレスデコードブロックとリルートモジュール ID ブロックとを含んでおり、
該デコードブロックが、
入力を多重化して、処理ユニット関連トランザクションに関連するアドレスを生成するマルチプレクサと、
前記アドレスを受信し、該アドレスに関連するトランザクションを受信するモジュールの宛先アドレスを提供するレンジレジスタデコーダと、
前記宛先アドレスを受信するリルートモジュール ID ブロックと
を含み、
前記リルートモジュール ID ブロックが、
前記コンピュータシステム中の 1 以上のオリジナルモジュールのアドレスを提供するオリジナルモジュール ID と、
前記コンピュータシステム中の代替モジュールの論理的な宛先モジュール ID を提供するリマッピングモジュール ID と
を含んでおり、該代替モジュールが前記コンピュータシステム中のオリジナルモジュールの機能を代替するものである相互接続メカニズム。

10

20

【請求項 2】

前記リルートモジュール ID ブロックは、有効ビット指示をさらに含み、該有効ビット指示が、前記オリジナルモジュールから前記代替モジュールへのトランザクションが有効な場合を示すものである請求項 1 に記載のメカニズム。

【請求項 3】

前記マルチプレクサが受信する前記入力は、前記複数の入出力装置からの入力である請求項 1 に記載のメカニズム。

【請求項 4】

前記アドレスが、入出力装置アドレスである請求項 1 に記載のメカニズム。

【請求項 5】

前記代替モジュールが、前記オリジナルモジュールにアドレス指定されたトランザクションを要求するプログラムを含むものである請求項 1 に記載のメカニズム。

30

【請求項 6】

前記オリジナルモジュールの状態が、前記代替モジュールにコピーされるものである請求項 1 に記載のメカニズム。

【請求項 7】

前記代替モジュールが、前記コンピュータシステムの非アクティブコンポーネントである請求項 1 に記載のメカニズム。

【請求項 8】

前記代替モジュールが、前記コンピュータシステムのアクティブコンポーネントである請求項 1 に記載のメカニズム。

40

【請求項 9】

コンピュータシステム中の故障したオリジナルモジュールから前記コンピュータシステム中の代替モジュールに仮想パスに沿ってトランザクションをリルーティングする方法であって、前記トランザクションが前記故障したオリジナルモジュールに対して開始されるものであり、
前記代替モジュールを前記故障したオリジナルモジュールの代替として識別するリマッピングモジュール ID を格納するステップと、
前記故障したオリジナルモジュールのトランザクションを受信するステップであって、前記トランザクションが前記故障したオリジナルモジュールのアドレスを含むものであるス

50

テップと、
前記アドレスを抽出するステップと、
該アドレスをデコードして、前記トランザクションを受信する前記故障したオリジナルモジュールのIDを提供するステップと、
該故障したオリジナルモジュールのIDを前記リマッピングモジュールIDと比較するステップと、
該リマッピングモジュールIDに基づいて前記トランザクションを前記代替モジュールにリルーティングするステップと
を含むものである方法。

【請求項10】

10

前記故障したオリジナルモジュールの状態を前記代替モジュールにコピーするステップと、前記故障したオリジナルモジュールにアドレス指定されたトランザクションを要求するように前記代替モジュールを再プログラムするステップとをさらに含む請求項9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、冗長サブシステムおよびコンポーネントを含むコンピュータシステムの技術分野に関する。

【0002】

20

【従来の技術】

現在のマルチプロセッサコンピュータシステムには、通常、主装置が故障した場合に使用することのできる1以上の冗長あるいは予備の装置が備えられている。たとえば、コンピュータシステムに2枚のイーサネットカードを装備し、第1のイーサネットカードが故障すると、ゼロもしくは最小のコンピュータダウンタイムで第2の（または予備の）カードを使用することができる。適切な冗長性を提供するために、現在のコンピュータシステムは、コンピュータシステムが分割される複数のパーティションについて予備装置を備えることができる。したがって、3つのパーティションを有するコンピュータシステムは、3つのパーティションそれぞれに1つの主装置と1つの予備装置とを含みうる。主装置および予備装置というこのコンフィギュレーションは、コンピュータシステムのコストを上げるとともに、コンピュータシステムレイアウトに対してさらなるスペースの制約を課す。

30

【0003】

【発明が解決しようとする課題】

【0004】

本発明の課題は、仮想ハードウェアパスを生成して、故障したコンピュータシステムコンポーネントにアドレス指定されたトランザクションを、代替コンピュータシステムコンポーネントが要求できるようにすることが可能な方法およびメカニズムを提供することである。

【課題を解決するための手段】

一実施形態では、コンポーネントは、イーサネットカードや他の入出力装置等の入出力装置（以下、「I/O」とよぶ）である。しかし、本方法および本メカニズムは、I/O装置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。

40

【0005】

オリジナルコンポーネント（または元のコンポーネント）と代替コンポーネントとは、同じタイプであることが好ましい。代替コンポーネントは、他のコンピュータシステム機能に現在使用されていてもよい（つまり、代替コンポーネントは、コンピュータシステムにおいてアクティブである）。あるいは、代替コンポーネントは、たとえば、インストールされた予備のように非アクティブであってもよい。

【0006】

一実施形態において、ハードウェアを使用して、故障中のまたは故障したコンポーネント

50

の間のバスを代替コンポーネント間のバスと同一に見えるようにする。故障したコンポーネント間の同じ物理バスは維持されるが、仮想バスが代替コンポーネントに対して確立される。そして、ソフトウェアを使用して、故障したコンポーネント間のアクティビティを保留し、故障したコンポーネントの状態を代替コンポーネントに再構築し、代替コンポーネントにおいて動作を再開する。そして、故障したコンポーネントに関するすべてのトランザクションまたはアクティビティが代替コンポーネントに行くことになる。この転送を確実にするために、レンジレジスタセットを使用してアドレス変換マッピングが呼び出される。プロセッサが、あるコンポーネントに行くアドレスを生成すると、そのアドレスをレンジレジスタと照らし合わせてチェックし、トランザクションをどのコンポーネントにルーティングすべきかを決定する。コンポーネントの故障によりトランザクションをリルーティングする必要がある場合には、レンジレジスタによって指定すべきリルートIDアドレス (reroute ID address) がマップテーブルに示される。

【0007】

特に、(故障した)オリジナルコンポーネントおよび代替コンポーネントの識別情報は、リルートモジュールIDブロック内に格納することができ、識別情報は、オリジナルコンポーネントが故障した場合には、適切な代替コンポーネントがリルートモジュールIDブロックを参照することによって識別することができるように、たとえばマップテーブルの使用などによって関連付けることができる。代替コンポーネントは、故障したコンポーネントにアドレス指定されたトランザクションを要求するとともに、故障したコンポーネントの状態を代替コンポーネントにコピーするために使用されるプログラムを含む。

【0008】

一実施形態では、1以上のブリッジユニットによって連結された複数のI/O装置および複数の処理ユニットを備えるコンピュータシステムで使用するための仮想入出力(I/O)相互接続メカニズムであって、アドレスデコードブロックを備え、アドレスデコードブロックは、入力を多重化して、プロセッサユニット関連トランザクションに関連するアドレスを生成するマルチプレクサと、アドレスを受信し、アドレスに関連するトランザクションを受信するモジュールの宛先アドレスを提供するレンジレジスタデコードと、宛先アドレスを受信するリルートモジュールIDブロックとを備えるメカニズムが提供される。リルートモジュールIDブロックは、コンピュータシステム中の1以上のオリジナルモジュールのアドレスを提供するオリジナルモジュールIDと、コンピュータシステム中の代替モジュールの論理的な宛先モジュールIDを提供するリマッピングモジュールID (remapping module ID) とを含み、代替モジュールはコンピュータシステム中のオリジナルモジュール(または元のモジュール)の機能を代替する。

【0009】

一実施形態では、コンピュータシステム中の故障したコンポーネントを動作コンポーネントに代替させる方法は、故障したコンポーネントを検出するステップと、故障したコンポーネントと同じタイプのコンポーネントが存在するかどうかを判定するステップとを含む。代替コンポーネントが存在する場合には、本方法は、故障したコンポーネントの間の (going to or coming from) 直接メモリアクセス等のすべてのアクティビティを保留するステップと、故障したコンポーネントの状態を代替コンポーネントにコピーするステップと、故障したコンポーネントをコンフィギュレーションを解除するステップと、故障したコンポーネントのハードウェアバスを代替コンポーネントにリマッピングするようにリルートモジュールIDを更新するステップと、代替コンポーネントのコンフィギュレーションレジスタを更新するステップと、故障したコンポーネントへの直接メモリアクセス等のアクティビティを再開するステップとを含む。代替コンポーネントが存在しない場合には、本方法はエラーハンドラを呼び出す。

【0010】

詳細な説明では、同じ符号が同じ要素を指す添付図面を参照する。

【0011】

【発明の実施の形態】

10

20

30

40

50

最近のコンピュータシステムは、互いに代替としての役割を果たすことができるいくつかの同様のコンポーネントを含みうる。たとえば、コンピュータシステムは、タイプ A のコンポーネント 4 つと、タイプ B のコンポーネント 4 つとを含みうる。4 つのタイプ A コンポーネントはすべて、ルーチンコンピュータシステム動作中に使用されることができる。すなわち、タイプ A コンポーネントに“予備 (spare)”はない。タイプ B コンポーネントの場合には、ルーチンコンピュータ動作中に 3 つが使用されることができ、4 番目のタイプ B コンポーネントはインストールされた予備でありうる。4 つのタイプ A コンポーネントの 1 つが故障した場合には、残りのタイプ A コンポーネントのうちの 1 つ以上を、故障したタイプ A コンポーネントの代替として利用することができる。タイプ B コンポーネントの 1 つが故障した場合には、インストールした予備のタイプ B コンポーネントを代替として利用することができる。

10

【0012】

図 1 は、4 つのタイプ A コンポーネントおよび 4 つのタイプ B コンポーネントを備えるコンピュータシステム 10 を示す。タイプ A およびタイプ B の各コンポーネントは、インタフェース接続 20 によってコンピュータシステム 10 の他のコンポーネント (図示せず) に連結される。4 つのタイプ A コンポーネントおよびコンポーネント 18 以外のすべてのタイプ B コンポーネントは、コンピュータシステム 10 の通常動作中に使用される。たとえば、コンポーネント 11 が故障した場合には、コンポーネント 13 (または、コンポーネント 15 あるいはコンポーネント 17) をコンポーネント 11 の代替とすることができる。コンポーネント 12 が故障した場合には、コンポーネント 18 をコンポーネント 12 の代替とすることができる。代替としてまたは加えて、コンポーネント 14 またはコンポーネント 16 を故障したコンポーネント 12 の代替とすることもできる。

20

【0013】

あるコンポーネントを別のコンポーネントで代替するには、故障したコンポーネントからのハードウェアパスを定義し、代替コンポーネントに対するハードウェアパスを、故障したコンポーネントのハードウェアパスと同一に見えるようにすることができる。そして、故障したコンポーネントを対象とした任意のランザクションは、代替コンポーネントに向けられることになる。したがって、コンポーネント 11 が故障し、コンポーネント 13 が代替として指定される場合には、コンポーネント 13 へのパス 23、20 は、コンポーネント 11 へのパス 21、20 と同一に見えるようにされる。以下、この概念を仮想化と呼ぶことにする。

30

【0014】

タイプ A または B コンポーネントのうちの 1 つの故障は、たとえば、失敗した直接メモリアクセス (direct memory access: 以下、「DMA」とよぶ) 試行中に検出することができる。コンピュータシステム 10 のハードウェア故障検出システム (図示せず) は、DMA の失敗を検出し、あるコンポーネントの別のコンポーネントでの代替化を完了するアルゴリズムを呼び出すことができる。故障したコンポーネントを代替する他に、コンポーネントの代替および仮想化は、あるタイプのコンポーネントをすべて取り外して、検査し、必要であれば修復し、交換するか、または、そのタイプの新しいコンポーネントと単に交換する定期的な予防のためのメンテナンス等の他の理由によっても行うことができる。

40

【0015】

図 2 は、仮想化が用いられるコンピュータシステムのより詳細な例である。コンピュータシステム 100 は、8 つの中央演算処理装置 (CPU) 101 ~ 108 を含む。各 CPU 101 ~ 108 は、図示のようにノースブリッジ 121 または 122 に連結される。ノースブリッジ 121 および 122 は、スケラブルインタフェース 120 によって接続される。ノースブリッジ 121 および 122 には、メモリ 124 およびメモリ 125 にも連結される。最後に、ノースブリッジ 121 にはサウスブリッジ 130 ~ 137 が連結され、ノースブリッジ 122 にはサウスブリッジ 140 ~ 147 が連結される。サウスブリッジ 140、144、130、132、および 136 には、イーサネットカード 154、15

50

5、151、152、および153がそれぞれ連結される。

【0016】

図2に示す各種ハードウェアコンポーネントを、いくつかの方式の1つに従ってパーティション化する(partitioned)ことができる。コンピュータシステム中のハードウェアコンポーネントのパーティション化は、コンピュータシステムパフォーマンスの最適化するための既知の技法である。例として、図3は1つの可能なパーティション方式を示す。パーティション0(160)は、CPU101、103、105と、いくつかのメモリ124と、イーサネットカード151と、他の入出力(I/O)カードおよび他のコンポーネント等の他のハードウェアコンポーネント(図示せず)とを含む。パーティション1(161)は、CPU102、106と、いくつかのメモリ124と、イーサネットカード152、154と、他のI/Oカードを含む他のハードウェアコンポーネント(図示せず)とを含む。パーティション2(162)は、CPU104、107、108と、いくつかのメモリ124と、イーサネットカード153と、他のI/Oカードを含む他のハードウェアコンポーネント(図示せず)とを含む。イーサネットカード155は、特定のパーティションいずれにも割り当てられない。

10

【0017】

次に図2および図3の双方を参照して、仮想化の実施(方法および装置)について詳細に述べる。具体的には、I/Oカードの仮想化(より具体的にはイーサネットカードの仮想化)を参照して説明する。しかし、コンピュータシステム100の他のハードウェアコンポーネントを、仮想化を用いてあるコンポーネントを同様の別のコンポーネントに代替することもできる。特定の例では、イーサネットカード152が故障する。故障したイーサネットカード152の機能を代替するため、イーサネットカード154からのハードウェアバスを故障したイーサネットカード152に対するハードウェアバスと同一に見えるようにすることによって、イーサネットカード154で代替することができる。すなわち、イーサネットカード154が、コンピュータシステム100の他のコンポーネントからはノースブリッジ121およびサウスブリッジ133に連結しているように見えるように、“仮想化”される。これは、イーサネットカード152に行くあらゆるトランザクションがイーサネットカード154にルーティングされることを意味する。加えて、イーサネットカード152に割り当てられたアドレス範囲は、イーサネットカード154によって要求されることになる。したがって、CPUがイーサネットカード152へのアドレスを生成する場合には、ノースブリッジ121および122はイーサネットカード154を、イーサネットカード152ではなく宛先として代替する。ピアツーピアトランザクション(peer-to-peer transaction)をイーサネットカード152にルーティングする必要がある場合には、ノースブリッジ121および122はピアツーピア転送(peer-to-peer transfer)をイーサネットカード154にルーティングする。加えて、イーサネットカード154は、以前はイーサネットカード152に割り当てられていたアドレス範囲を要求するようにプログラムされる。最後に、後述するように、イーサネットカード152の状態がイーサネットカード154にコピーされる。

20

30

【0018】

図4は、CPUがI/Oアクセスおよびイーサネットカード151~155へのハードウェアバスの仮想化を行うことができるように、ノースブリッジ121および122に組み込むことのできるアドレスデコードブロック170を示す。171において、CPU101~104がノースブリッジ121への入力を提供し、この入力がマルチプレクサ172において多重化されてアドレス173を提供する。そして、アドレス173がレンジレジスタデコード174に提供される。デコード174の出力は、(たとえば、ノースブリッジやサウスブリッジの)宛先175を含む。宛先175はリルートモジュールIDブロック176に提供され、リルートモジュールIDブロック176は論理宛先ID177を提供する。

40

【0019】

50

図5は、リルートモジュールIDブロック176を詳細に示す。ブロック176は、有効ビット列181と、オリジナルモジュールIDセクション182と、リマッピングモジュールIDセクション183とを含む。また、制御ブロック184も図示されている。オリジナルモジュールIDセクション182は、イーサネットカード151～155のうちの1以上のものについての識別情報を含む。この情報は、元々機能しているイーサネットカードを識別する。リマッピングモジュールIDセクション183は、元々機能しているイーサネットカードが故障した場合（または、代替が必要な他のアクションの場合）の代替イーサネットカードを識別する情報を含む。（たとえば、ビットが1に設定されている場合には）有効ビット列181は、オリジナルの故障したイーサネットカードから代替イーサネットカードへの変換（translation）が有効な場合を示す。

10

【0020】

リルートモジュールIDブロック176は、いくつかのエントリを含みうる。エントリ数は、いくつかの相互接続が代替を同時に受信することができるかに影響する（dictate）。たとえば、リルートモジュールIDブロック176が8個のエントリを含む場合には、最高で8個の代替または再配向を同時に行うことができる。各エントリは、エントリ（または、変換）が有効なことを示す有効ビットと、オリジナルモジュールIDと、代替モジュールIDとを含む。

【0021】

図6は、I/O仮想化プロセス200を示すフローチャートである。図6において、プロセス200は、図2に示すイーサネットカードの仮想化、特にイーサネットカード154で代替可能なイーサネットカード152の故障に関連する。プロセス200はブロック205において開始される。ブロック210において、管理ソフトウェアが、イーサネットカード152と同じタイプの予備イーサネットカードが存在し、利用可能であるかどうかを判定する。予備イーサネットカードが利用可能ではない場合には、プロセス200はブロック215に移り、エラーハンドラが呼び出される。ブロック210において、予備イーサネットカードが利用可能な場合には、イーサネットカード152が故障した状態であり、プロセス200はブロック220に移る。図示の例では、イーサネットカード154が存在し、故障したイーサネットカード152の代替として利用可能である。ブロック220において、管理ソフトウェアはDMAを中断する。次に、ブロック225において、イーサネットカード152の状態がイーサネットカード154にコピーされる。そして、ブロック230において、管理ソフトウェアはイーサネットカード152のコンフィギュレーションを解除する（deconfigure）。ブロック235において、管理ソフトウェアは、イーサネットカード152へのトランザクションを生成可能なコンピュータシステム100全体を通じてリルートモジュールIDブロックを更新する。更新は、有効ビット181を0から1に設定することと、イーサネットカード152のオリジナルモジュールIDをイーサネットカード152のサウスブリッジ側に設定することと、イーサネットカード152のリマッピングモジュールIDをイーサネットカード154のサウスブリッジ側に設定することを含む。

20

30

【0022】

ブロック240において、ノースブリッジ122およびサウスブリッジ144におけるコンフィギュレーションレジスタは、イーサネットカード154が元々イーサネットカード152に割り当てられていたアドレス範囲を要求するように更新される。ブロック245において、管理ソフトウェアは、イーサネットカード152へのDMAを再開する。ブロック250において、プロセス200が終了する。

40

【0023】

故障したイーサネットカード152を修復し、コンピュータシステム100に復帰させることができる。この場合には、復帰したイーサネットカード152は、故障したイーサネットの代替となりうる予備のイーサネットカードとしての役割を果たすことができる。

【0024】

上述した例示的な実施形態は、カード（または、モジュール）レベルでの代替またはパス

50

の仮想化に言及している。しかし、この代替は、コンピュータシステムにおいてカードレベルよりも低いレベルまたは高いレベルで行うことができる。

【図面の簡単な説明】

【図 1】冗長コンポーネントを採用するコンピュータシステムを示す概略図である。

【図 2】I/Oカードが故障した場合に、I/O相互接続の仮想化を使用して冗長性を提供するマルチプロセッサコンピュータシステムを示す概略図である。

【図 3】図 2 のシステムと併せて使用することが可能なパーティション方式を示す概略図である。

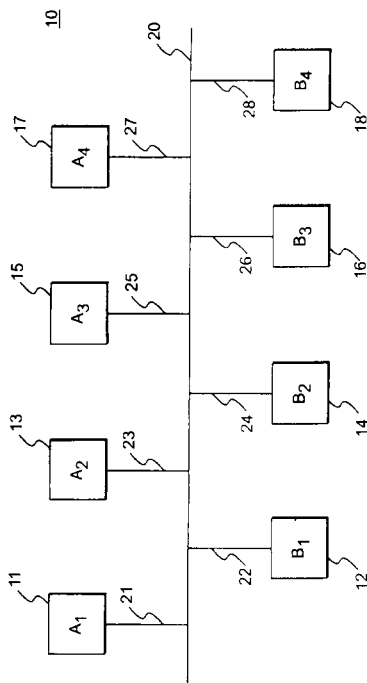
【図 4】図 2 のシステムと併せて使用するアドレスデコードブロックを示す概略図である。

【図 5】図 2 のシステムと併せて使用するリルートモジュールブロックを示す概略図である。

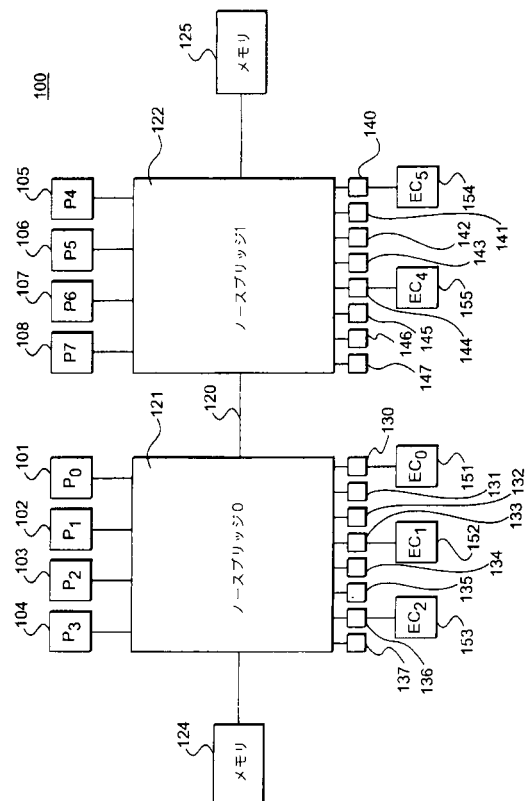
【図 6】ハードウェアパス仮想化方法を示すフローチャートである。

10

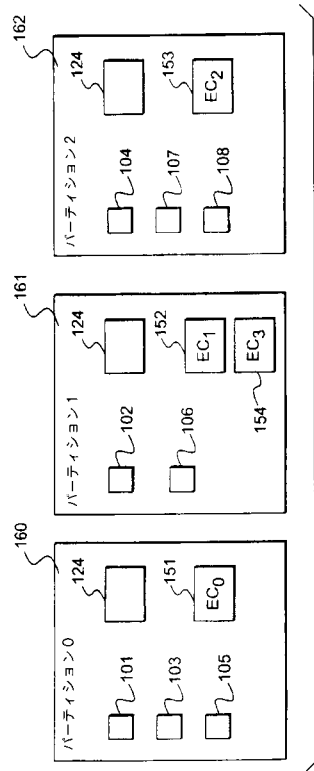
【図 1】



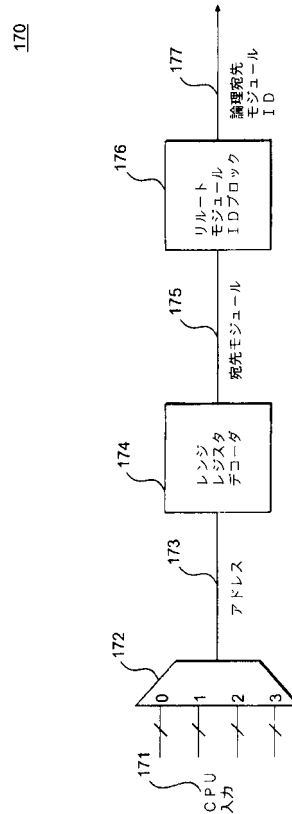
【図 2】



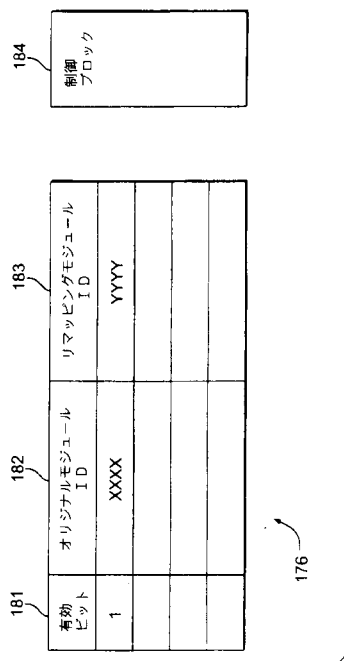
【図 3】



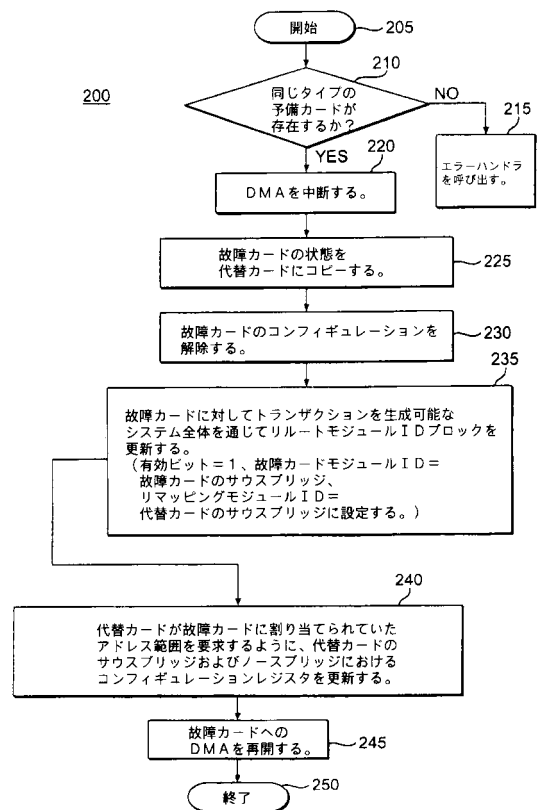
【図 4】



【図 5】



【図 6】



フロントページの続き

(72)発明者 ディベンドラ・ダス・シャルマ

アメリカ合衆国カリフォルニア州 9 5 0 5 0 , サンタ・クララ , アカシア・コート 2 0 4 3

(72)発明者 アシシ・グプタ

アメリカ合衆国カリフォルニア州 9 5 1 2 9 , サン・ノゼ , オラ・ストリート 5 6 3 7

F ターム(参考) 5B014 HA09 HB14 HC02

5B034 BB02 BB11 CC05

5B083 BB03 CC04 CD11 EE07 GG04