

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6134720号
(P6134720)

(45) 発行日 平成29年5月24日(2017.5.24)

(24) 登録日 平成29年4月28日(2017.4.28)

(51) Int.Cl.

F I

G 0 6 F 13/00 (2006.01)

G 0 6 F 13/00 3 0 1 J

請求項の数 15 (全 15 頁)

(21) 出願番号 特願2014-533980 (P2014-533980)
 (86) (22) 出願日 平成24年9月28日 (2012.9.28)
 (65) 公表番号 特表2014-532236 (P2014-532236A)
 (43) 公表日 平成26年12月4日 (2014.12.4)
 (86) 国際出願番号 PCT/GB2012/052413
 (87) 国際公開番号 W02013/050745
 (87) 国際公開日 平成25年4月11日 (2013.4.11)
 審査請求日 平成27年7月14日 (2015.7.14)
 (31) 優先権主張番号 1117230.1
 (32) 優先日 平成23年10月5日 (2011.10.5)
 (33) 優先権主張国 英国 (GB)

(73) 特許権者 512313285
 マイクロン テクノロジー, インコーポレ
 イティド
 アメリカ合衆国, アイダホ 83716-
 9632, ボイジー, サウス フェデラル
 ウェイ 8000
 (74) 代理人 100083806
 弁理士 三好 秀和
 (74) 代理人 100095500
 弁理士 伊藤 正和
 (74) 代理人 100111235
 弁理士 原 裕子

最終頁に続く

(54) 【発明の名称】 接続方法

(57) 【特許請求の範囲】

【請求項 1】

第1のデバイスと第2のデバイス間の通信の障害を管理する方法であって、

上記第1のデバイスと上記第2のデバイスの中間の第3のデバイスにおいて、上記第1のデバイスから上記第2のデバイスへ供給される第1のデータを傍受し、上記第1のデータは上記第1のデータが属するトランザクションを示すデータフィールドを含み、

上記第3のデバイスにおいて、上記第1のデータに対する応答が上記第1のデバイスによって受信されないという上記第1のデバイスと第2のデバイス間の通信の障害を検出し、

上記第1のデータに対する応答が上記第2のデバイスから上記第1のデバイスによって受信されないことに起因する上記第1のデバイスのサービスの中断を防止するために、上記第2のデバイスが該第1のデバイスによって使えないことを示す第2のデータが、上記第3のデバイスにおいて、生成され、上記第3のデバイスから上記第1のデバイスに送信されることによって上記第2のデバイスからの連続した通信をエミュレートし、上記第2のデータは上記第2のデータの送信元を示す一部を含み、上記第2のデータの上記一部は上記第2のデータが上記第2のデバイスに由来するように見えるように上記第2のデバイスを示す方法。

【請求項 2】

上記エミュレーションは、上記第3のデバイスにおいて、上記第1のデータの上記データフィールドによって示された上記トランザクションに少なくとも基づいて複数の格納さ

10

20

れたエラーメッセージから一つのエラーメッセージを選択し、上記選択されたエラーメッセージは上記第2のデータに含まれる請求項1に記載の方法。

【請求項3】

上記第2のデータは、該第2のデータが上記第1のデータに対する応答であることを示すデータを含む請求項1記載の方法。

【請求項4】

上記第1及び第2のデバイスの少なくとも1つは、サーバである請求項1乃至3のいずれか1項記載の方法。

【請求項5】

上記第1及び第2のデバイスの少なくとも1つは、I/Oデバイスである請求項1乃至4のいずれか1項記載の方法。

【請求項6】

上記第1のデバイスは、PCI Express接続によって上記第2のデバイスに接続される請求項1乃至5のいずれか1項記載の方法。

【請求項7】

上記第2のデータは、上記第2のデバイスが連絡できないことを示す請求項1乃至6のいずれか1項記載の方法。

【請求項8】

上記第2のデータは、上記第2のデバイスが故障したことを示す請求項1乃至7のいずれか1項記載の方法。

【請求項9】

上記第2のデータは、該第2のデータが破損していることを示す請求項1乃至8のいずれか1項記載の方法。

【請求項10】

上記第1のデバイスと第2のデバイス間の通信の障害の検出では、上記第1のデバイスと上記第2のデバイスを接続しているケーブルが切断されたことを検出することを含む請求項1乃至9のいずれか1項記載の方法。

【請求項11】

上記第1のデバイスと第2のデバイス間の通信の障害の検出では、

上記第3のデバイスから上記第2のデバイスに、第3のデータを送信し、

上記第2のデバイスから、上記第3のデータに対する応答が所定の時間内に受信されないときに、当該第1のデバイスと第2のデバイス間の通信の障害を決定することを含む請求項1乃至10のいずれか1項記載の方法。

【請求項12】

コンピュータが読み取り可能なプログラムを有するコンピュータが使えるストレージデバイスを含むコンピュータプログラム製品であって、上記コンピュータが読み取り可能プログラムは、第1のデバイスと第2のデバイスの間で実行されると、上記コンピュータに、

上記第1のデバイスから上記第2のデバイスへの第1のデータを傍受させ、上記第1のデータは上記第1のデータが属するトランザクションを示すデータフィールドを含み、

上記第1のデータに対する応答が上記第1のデバイスによって受信されないという上記第1のデバイスと上記第2のデバイス間の通信の障害を検出させ、

上記第1のデータに対する応答が上記第2のデバイスから上記第1のデバイスによって受信されないことに起因する上記第1のデバイスのサービスの中断を防止するために、第2のデータを生成して上記第1のデバイスに送信することによって上記第2のデバイスからの連続した通信をエミュレートさせ、上記第2のデータは上記第2のデバイスが上記第1のデバイスによって使えないことを示し、上記第2のデータは上記第2のデータの送信元を示す一部を含み、上記第2のデータの上記一部は上記第2のデータが上記第2のデバイスに由来するように見えるように上記第2のデバイスを示すようにさせるコンピュータプログラム製品。

【請求項13】

10

20

30

40

50

コンピュータが読取り可能なプログラムを格納したメモリデバイスを含むコンピュータプログラム製品であって、上記コンピュータが読取り可能プログラムは、第１のデバイスと第２のデバイスの間で実行されると、上記コンピュータに、

上記第１のデバイスから上記第２のデバイスへの第１のデータを傍受させ、上記第１のデータは上記第１のデータが属するトランザクションを示すデータフィールドを含み、

上記第１のデータに対する応答が上記第１のデバイスによって受信されないという上記第１のデバイスと上記第２のデバイス間の通信の障害を検出させ、

上記第１のデータに対する応答が上記第２のデバイスから上記第１のデバイスによって受信されないことに起因する上記第１のデバイスのサービスの中断を防止するために、第２のデータを生成して上記第１のデバイスに送信することによって上記第２のデバイスからの連続した通信をエミュレートさせ、上記第２のデータは上記第２のデバイスが上記第１のデバイスによって使えないことを示し、上記第２のデータは上記第２のデータの送信元を示す一部を含み、上記第２のデータの上記一部は上記第２のデータが上記第２のデバイスに由来するように見えるように上記第２のデバイスを示すようにさせるコンピュータプログラム製品。

【請求項１４】

第１のデバイスと第２のデバイス間の通信の障害を管理するコンピュータ装置であって、

プロセッサで読取り可能な命令を記憶するメモリと、

上記メモリに記憶された命令を読み出して、実行するプロセッサとを含み、

上記プロセッサで読取り可能な命令は、上記コンピュータを制御するように構成され、

上記第１のデバイスから上記第２のデバイスへの第１のデータを傍受させ、上記第１のデータは上記第１のデータが属するトランザクションを示すデータフィールドを含み、

上記第１のデータに対する応答が上記第１のデバイスによって受信されないという上記第１のデバイスと上記第２のデバイス間の通信の障害を検出させ、

上記第１のデータに対する応答が上記第２のデバイスから上記第１のデバイスによって受信されないことに起因する上記第１のデバイスのサービスの中断を防止するために、第２のデータを生成して上記第１のデバイスに送信することによって上記第２のデバイスからの連続した通信をエミュレートさせ、上記第２のデータは上記第２のデバイスが上記第１のデバイスによって使えないことを示し、上記第２のデータは上記第２のデータの送信元を示す一部を含み、上記第２のデータの上記一部は上記第２のデータが上記第２のデバイスに由来するように見えるように上記第２のデバイスを示すようにさせるコンピュータ装置。

【請求項１５】

自己修復ケーブルアダプタであって、接続された第１のデバイスから接続された第２のデバイスへの第１のデータを傍受し、上記第１のデータに対する応答が上記接続された第１のデバイスによって受信されないという上記接続された第１のデバイスと上記接続された第２のデバイス間の通信の障害を検出し、上記第１のデータに対する応答が上記接続された第２のデバイスから上記接続された第１のデバイスによって受信されないことに起因する上記接続された第１のデバイスのサービスの中断を防止するように、第２のデータを生成して上記接続された第１のデバイスに送信することによって上記接続された第２のデバイスからの連続した通信をエミュレートし、上記第２のデータは上記第２のデータの送信元を示す一部を含み、上記第２のデータの上記一部は上記第２のデータが上記第２のデバイスに由来するように見えるように上記接続された第２のデバイスを示し、上記第２のデータは上記接続された第２のデバイスが上記第１のデバイスによって使えないことを示し、上記第１のデータは上記第１のデータが属するトランザクションを示すデータフィールドを含む自己修復ケーブルアダプタを含む装置。

【発明の詳細な説明】

【技術分野】

【０００１】

本発明は、第1のデバイスと第2のデバイス間の通信の障害を管理する方法に関する。

【背景技術】

【0002】

コンピュータシステムでは、多くの場合、デバイス間でデータを送信することが必要であり、例えば、多くの場合、処理デバイスを複数の入出力デバイスに接続することが必要である。適切なデータ通信は、デバイスがリンクを介して互いにデータパケットを送信できるように、デバイスを接続することによって達成され、そして、リンクは有線リンク又は無線リンクであってもよい。

【0003】

被接続デバイスのシステムにとって、それらの被接続デバイスの1つの予期しない故障又は不稼働率を管理できることが、ますます重要になっている。例えば、システムは、ケーブルが間違っ

10

【0004】

例えば、サーバ(すなわち「インボックス」)内にI/Oアダプタを永久接続するために、PCI Expressだけが既に用いられている。I/Oデバイスをサーバ内に永久的にインストールするこのような構成の場合、I/Oデバイスの故障は、サーバ全体の故障として現れ、そのサーバ自体は、引き続き動作する必要があるとは考えられない。遠隔被接続PCI Express I/Oデバイス(すなわちサーバの外で接続されるI/Oデバイス)に対する関心が高まり、被接続デバイスの予期しない活線取外しをサーバが段階的にサポートする要求が、システムの障害許容力に対して重要になることを意味している。

20

【0005】

データパケットの送受信は、多くの場合、トランザクションの観点から説明される。トランザクションは、デバイス間で送信される1つ以上のデータパケットを含む。PCI Expressは、スプリットトランザクションモデルを実行し、そこでは、送信元デバイスは、要求データパケットを宛先デバイスに送信し、そして、応答における宛先デバイスからの完了データパケットを待つ。一般的に、オペレーティングシステムは、失敗したPCI Express トランザクションを段階的に取り扱うように構成されていない。例えば、サーバが要求データパケットを被接続デバイスに送信し、予想に反して、その要求に対する応答において完了データパケットを受信しない場合、サーバのオペレーティングシステムは、クラッシュする虞がある。PCI Expressに基づく現在の被接続システム自体は、PCI Expressに接続された送信元デバイスが予想に反して利用できなくなったとき、クラッシュする虞がある。

30

【0006】

PCI Expressの標準的な実装は、PCI Expressサブシステム、すなわち被接続PCI Expressデバイスの故障を取り扱う十分な手段を提供していない。これは、結果として、満足できるレベルの障害許容力を有する遠隔又は共用のI/Oシステムの開発を難しくしている。

40

【発明の概要】

【発明が解決しようとする課題】

【0007】

本発明の実施の形態の目的は、上で概説した問題の1つ以上を取り除く又は軽減することである。

【課題を解決するための手段】

【0008】

50

本発明の第1の特徴に基づいて、第1のデバイスと第2のデバイス間の通信の障害を管理する方法を提供し、この方法は、第1のデバイスと第2のデバイスの中間の第3のデバイスにおいて、第1のデバイスと第2のデバイスの通信の障害を検出し、第3のデバイスから第1のデバイスに、第2のデバイスが第1のデバイスによって使えないことを示す第1のデータを伝送する。

【0009】

このように、第1のデバイスは、第1のデバイスが故障する原因となる、第2のデバイスによる予期しない通信の損失にさらされることはない。第1のデバイスと第2のデバイス間の通信の損失を検出することによって、第3のデバイスは、適切なデータを第1のデバイスに送信して、第1のデバイスの故障を防止することができる。例えば、第2のデバイスが状態を変えたことを第1のデバイスに通知することによって、第1のデータは、通信の障害を第1のデバイス上で動作しているソフトウェア（OS、ドライバ、アプリケーション）から隠すことができる。例えば、第3のデバイスは、第2のデバイスの連続した存在を、第1のデバイスが「利用可能であるが、使えない」として解釈する状態にエミュレートすることができる。第1のデバイスは、（例えば第2のデバイスに代わるもので通信することによって）、サービスが中断することなく、「利用可能であるが、使えない」状態を段階的に取り扱うように構成することができる。

【0010】

第3のデバイスは、第1のデバイスによって定義される故障ユニット内であってもよく、一方、第2のデバイスは、第1のデバイスによって定義される故障ユニットの外であってもよい。

【0011】

第2のデバイスが使えないことを示す際に、第1のデータは、一時的かもしれない外部の問題のために、第2のデバイスが局所的に機能するが、その指定サービスを実行できないことを示すことができる。

【0012】

方法は、更に、第3のデバイスにおいて、第1のデバイスから第2のデバイスへの第2のデータを傍受してもよい。

【0013】

第1のデバイスと第2のデバイス間の通信が第3のデバイスを通るように、第3のデバイスは、第1のデバイスと第2のデバイスの間にあってよい。例えば、第3のデバイスは、ケーブルアダプタの形をとってもよい。あるいは、第1のデバイスから第2のデバイスに送信されるデータは、第1のデバイスと第2のデバイス間の更なるデバイスによって、第3のデバイスに転送されてもよい。

【0014】

第1のデバイスと第2のデバイス間の通信の障害の検出では、第2のデータに対する応答が第1のデバイスによって受信されないことを検出することを含んでもよい。第1のデバイスと第2のデバイス間の通信の障害の検出では、第1のデバイスと第2のデバイスを接続しているケーブルが切断されたことを検出することを含んでもよい。例えば、ケーブルの切断の検出では、以前に存在していた「ケーブル検出」信号がないことを検出することを含んでもよい。あるいは、第1のデバイスと第2のデバイス間の接続が無線接続の場合、通信の障害の検出では、無線接続の中断又は干渉を検出することを含んでもよい。

【0015】

第1のデータの送信元を示すことを目的とした第1のデータの一部は、第2のデバイスを示すことができる。第1のデータは、第1のデータが第2のデータに対する応答であることを示すデータを含むことができる。例えば、第2のデータは、第2のデータが属するトランザクションを示すデータフィールドを含むことができ、第1のデータは、それが同じトランザクションに属することを示すデータを含むことができる。

【0016】

第1及び第2のデバイスの1つは、サーバであってもよく、及び/又は第1及び第2の

デバイスの１つは、第１又は第２のデバイスの遠隔被接続リソースであってもよい。例えば、第１及び第２のデバイスの１つは、Ｉ／Ｏデバイス又は他の遠隔リソースであってもよい。

【００１７】

第１のデバイスは、PCI Express接続によって第２のデバイスに接続されていてもよい。

【００１８】

第１のデータは、第２のデバイスが連絡できないことを示してもよく、第１のデータは、第２のデバイスが故障したことを示してもよく、及び／又は第１のデータは、第１のデータが破損していることを示してもよい。例えば、第１のデータは、第２のデバイスが故障し、及び／又は第１のデータが破損したことを示す値が「１」を有する状態ビットを含むことができる。

10

【００１９】

第１のデバイスと第２のデバイス間の通信の障害の検出では、第３のデバイスから第２のデバイスに、第３のデータを伝送し、第２のデバイスから、第３のデータに対する応答が所定の時間内に受信されないときに、第１のデバイスと第２のデバイス間の通信のその障害を決定することを含んでもよい。

【００２０】

本発明の第２の特徴に基づいて、第１のデバイスと第２のデバイス間の通信の障害を管理する装置を提供し、この装置は、第１のデバイスと第２のデバイスの中間の第３のデバイスにおいて、第１のデバイスと第２のデバイス間の通信の障害を検出する手段と、第３のデバイスから第１のデバイスに、第２のデバイスが第１のデバイスによって使えないことを示す第１のデータを伝送する手段とを含む。

20

【００２１】

本発明の第３の特徴に基づいて、第１のデバイスと第２のデバイス間を接続する自己修復ケーブルアダプタを提供し、自己修復ケーブルアダプタは、被接続第１のデバイスと被接続第２のデバイス間の通信の障害を検出する検出器と、被接続第１のデバイスに、被接続第２のデバイスが被接続第１のデバイスによって使えないことを示す第１のデータを伝送する送信機とを含む。

【００２２】

発明の１つの特徴に関連して説明した多くの特徴が発明の他の特徴に関連して適用できることは言うまでもない。

30

【００２３】

本発明の特徴が適切なハードウェア及び／又はソフトウェアを含むあらゆる便利な方法で実現できることは言うまでもない。例えば、発明を実施するように配置されたスイッチングデバイスは、適切なハードウェアコンポーネントを用いて形成することができる。あるいは、プログラマブルデバイスは、発明の実施の形態を実現するようにプログラムすることができる。したがって、また、発明は、発明の特徴を実現する適切なコンピュータプログラムを提供する。このようなコンピュータプログラムは、有形のキャリア媒体（例えばハードディスク、ＣＤ－ＲＯＭ等）と、無形のキャリア媒体（例えば通信信号）とを含む適切なキャリア媒体で伝達することができる。

40

【図面の簡単な説明】

【００２４】

つぎに、本発明の実施の形態を、図面を参照し、一例として説明する。

【００２５】

【図１】Ｉ／Ｏデバイスがサーバのコンポーネントとして設けられた従来のＩ／Ｏ構成のブロック図である

【００２６】

【図１a】Ｉ／Ｏデバイスがサーバのコンポーネントとして設けられた従来のＩ／Ｏ構成のブロック図である

50

【 0 0 2 7 】

【図 2】 I / O デバイスがサーバに遠隔で接続された従来の I / O 構成のブロック図である。

【 0 0 2 8 】

【図 3】 複数の I / O デバイスがサーバに遠隔で接続された従来の I / O 構成のブロック図である。

【 0 0 2 9 】

【図 4】 I / O デバイスが本発明の実施の形態に基づいてサーバに接続された I / O 構成のブロック図である。

【 0 0 3 0 】

【図 5】 データパケットヘッダの略図である。

【 0 0 3 1 】

【図 6】 サーバが本発明の実施の形態に基づいて I / O デバイスに接続された I / O 構成のブロック図である。

【発明を実施するための形態】

【 0 0 3 2 】

図 1 に示すように、サーバ 1 は、CPU / PCI Express ルートコンプレックス (CPU / RC) 2 と、ネットワークインタフェースコントローラ (NIC) 3 とを含む。CPU / RC 2 は、PCI Express チップ間接続 4 によって、NIC 3 に接続されている。NIC 3 は、イーサネット (登録商標) 接続 6 (例えば、ケーブル接続又は無線接続であってもよい) によって、ネットワーク 5 に接続されている。NIC 3 及び適切なソフトウェアを使用することによって、サーバ 1 は、サーバ 1 のユーザに、ネットワーク 5 に対するアクセスを提供する。サーバ 1 の本発明に関係しない他の詳細は、説明を明確にするために省略されていることは言うまでもなく、それらの詳細は、当業者にとって容易に明らかである。

【 0 0 3 3 】

一般論として、図 1 のサーバ 1 は、単一の故障ユニットであると考えられる。サーバ 1 の内部コンポーネント (例えば NIC 3 又は接続 4) の故障は、サーバ 1 全体の故障であると考えられ、サーバ 1 が最早その意図された機能 (すなわちネットワーク 5 に対するアクセスを容易にすること) を提供できないと考えられる。NIC 3 自体が故障し、又は NIC 3 と CPU / RC 2 間の通信に障害が生じると、サーバ 1 が機能し続ける必要性はなくなり、そして、サーバ 1 は、機能しなくなる。

【 0 0 3 4 】

内部コンポーネントの故障とは対照的に、外部コンポーネント (すなわちサーバ 1 の外部の単一の故障ユニット) の故障、例えばイーサネット接続 6 又は下流のスイッチ (図示せず) は、サーバ 1 の故障であると考えられず、サーバ 1 自体は、動作し続ける必要がある。一般論として、交換するデバイスは、簡単に接続する (あるいは、例えば間違っただけで切断した場合に再接続する) ことができることを考えると、サーバ 1 に外的に接続されたコンポーネントの故障は、サーバ 1 の故障とすべきではない。したがって、サーバ 1 は、外部リソースが現在利用できないことを正確に通知する状況を段階的に取り扱うように構成されている。例えば、CPU / RC 2 と NIC 3 間のその通信が維持されているならば、イーサネット接続 6 の故障により、NIC 3 は、(適切なエラーメッセージを送信することによって) イーサネット接続 6 が利用できないことを CPU / RC 2 に通知する。CPU / RC 2 は、このメッセージを受信すると、サーバ 1 全体の故障を起こすことなく、イーサネット接続 6 上の伝送を用いる NIC 3 へのデータパケットの送信を止めることができる。

【 0 0 3 5 】

CPU / RC は、NIC からのエラーメッセージを受信すると、利用可能なバックアップリソースを利用することによって、その指定サービスを提供し続けることができる。図 1 a に示すように、サーバ 1 a は、PCI Express チップ間接続 4 a によって NIC 3 a に

10

20

30

40

50

、PCI Expressチップ間接続 4 b によってNIC 3 b に接続されたCPU/R C 2 a を備える。NIC 3 a、3 b のそれぞれは、それぞれのイーサネット接続 6 a、6 b によってネットワーク 5 に接続されている。NIC 3 a、3 b のそれぞれ自体は、ネットワーク 5 に対するアクセスを容易にすることができる。イーサネット接続 6 b が故障すると、例えば、対応するNIC 3 b は、イーサネット接続 6 b が利用できないことを示すエラーメッセージをCPU/R C 2 a に送信する。CPU/R C 2 a がこのようなエラーメッセージを段階的に取り扱うように構成されているとき、CPU/R C 2 a は、NIC 3 a 及びイーサネット接続 6 a を用いて、その指定機能を提供し続けることができる。

【0036】

図 2 の構成において、NIC 7 は、サーバ 8 に対して外的（すなわち、サーバ 8 によって定義される単一の故障ユニットの外）に收容されている。図 2 に示すように、サーバ 8 は、PCI Express接続 1 1 によってCPU/R C 1 0 に接続されたPCI Expressケーブルアダプタ 9 を含む。PCI Expressケーブルアダプタ 9 は、CPU/R C 1 0 をケーブル 1 4 によって、遠隔 I/O 機器 1 3 内に收容されたPCI Expressケーブルアダプタ 1 2 に接続する。また、遠隔 I/O 機器 1 3 は、NIC 7 を收容し、NIC 7 は、PCI Express接続 1 5 によってPCI Expressケーブルアダプタ 1 2 に接続されている。PCI Expressケーブルアダプタ 9、1 2 は、ケーブル 1 4 でPCI Express信号を運ぶのに必要なケーブル駆動及び信号調整機能を備えるだけである。すなわち、ケーブルアダプタ 9、1 2 は、サーバ 8 上で動作するソフトウェアのような論理機能を全く備えないが、図 2 のシステムは、論理的には図 1 のシステムと同じである。

【0037】

ケーブル 1 4 が故障又はNIC 7 自体が故障すると、CPU/R C 1 0 とNIC 7 間の通信は切断される。CPU/R C 1 0 自体は、NIC 7 によって処理されていないトランザクションに対する応答を受信しない。上述したように、未処理のトランザクションに対する応答の欠如は、結果的に、サーバ 8 全体の故障につながる。図 1 の構成におけるNIC 3（又はNIC 3 とCPU/R C 2 間の接続 4）の故障がサーバ 1 全体の故障となることは、通常、許容できるが、図 2 の構成におけるケーブル 1 4 及びNIC 7 の遠隔の本質は、これらのコンポーネントを簡単に交換又は再接続できることを意味する。図 2 の構成自体において、NIC 7 又はケーブル 1 4 のどちらかの故障をサーバ 8 の故障の原因としなければならないことは、望ましくない。

【0038】

さらに、図 3 の構成に示すように、サーバから離してI/Oリソースを配置することにより、複数のI/Oリソースを容易に配置でき、1つの被接続I/Oリソースが故障した場合に、次のリソースがバックアップとして機能することができる。図 3 に示すように、サーバ 2 0 のCPU/R C 1 9 は、ケーブルアダプタ 2 1 によって、ケーブル 2 3 により遠隔 I/O 機器 2 2 に、ケーブル 2 5 により遠隔 I/O 機器 2 4 に接続されている。遠隔 I/O 機器 2 2 は、ケーブル 2 3 に接続されたケーブルアダプタ 2 6 と、接続 2 8 によってケーブルアダプタ 2 6 に接続されたNIC 2 7 とを含む。遠隔 I/O 機器 2 4 は、ケーブル 2 5 に接続されたケーブルアダプタ 2 9 と、接続 3 1 によってケーブルアダプタ 2 9 に接続されたNIC 3 0 とを含む。各NIC 2 7、3 0 は、それぞれの接続 3 2、3 3 によってネットワーク 5 に接続されている。I/O 機器 2 2、2 4 の1つが故障し、NIC 2 7、3 0 の1つによってその通信が最早できない場合、あるいは接続 2 3、2 5 の1つが故障すると、サーバ 2 0 は、別のI/O機器 2 2、2 4 を用いて、その指定サービス（すなわちネットワーク 5 に対する接続）を提供し続けることができることは言うまでもない。しかしながら、I/O機器 2 2、2 4（又は接続 2 3、2 5）のどちらかが故障すると、CPU/R C 1 9 は、NIC 2 7、3 0 のうちの切断された1つに送信したデータパケットに対する応答を受信せず、結果的に、サーバ 2 0 の望ましくない故障につながる。

【0039】

図 4 は、本発明の実施の形態に基づいて変更された図 2 の一般的配置を示す。特に、図 4 の構成では、サーバによって定義される故障ユニット内のケーブルアダプタは、自己修

10

20

30

40

50

復ケーブルアダプタに置き換えられており、以下、その動作を更に詳細に説明する。しかしながら、一般論として、自己修復ケーブルアダプタは、致命的エラー（例えばP C I ケーブルの故障）を致命的でないエラーに変換する。

【 0 0 4 0 】

図 4 において、サーバ 3 5 は、PCI Express接続 3 8 によって自己修復ケーブルアダプタ 3 7 に接続されたC P U / R C 3 6 を含む。自己修復ケーブルアダプタ 3 7 は、C P U / R C 3 6 をケーブル 4 0 によって、遠隔 I / O 機器 3 9 内に収容されたPCI Expressケーブルアダプタ 3 8 に接続する。遠隔 I / O 機器 3 9 は、N I C 4 1 を収容し、N I C 4 1 は、PCI Express接続 4 2 によってPCI Expressケーブルアダプタ 3 8 に、イーサネット接続 4 3 によってネットワーク 5 に接続されている。

10

【 0 0 4 1 】

より詳細には、自己修復ケーブルアダプタ 3 7 は、自己修復ケーブルアダプタ 3 7 の下流のコンポーネント（例えばケーブル 4 0、N I C 4 1 等）の故障を監視するのに必要な論理及びハードウェアを備えている。例えば、自己修復ケーブルアダプタ 3 7 は、（例えば、供給される「存在検出」信号又は他のあらゆる適切な方法によって）ケーブルが抜かれたときを検出し、遠隔 I / O 機器 3 9 からのN I C 4 1 に関する問題を示すメッセージを監視するように構成されている。さらに、自己修復ケーブルアダプタ 3 7 は、サーバ 3 5 によって開始されたトランザクションに関する情報を判定して、記録するために、データパケットを調べ、そして、適切な応答を待つように構成されている。自己修復ケーブルアダプタ 3 7 は、オンボードタイマを用いて、その終了が故障の発生となる所定の時間待つことができる。例えば、時間の閾値は、応答データパケットを受信しないことの結果として、サーバ 3 5 がクラッシュすることになる時間よりも前に終了するように、設定することができる。また、自己修復ケーブルアダプタ 3 7 は、C P U / R C 3 6 とは関係なく、データパケットを生成して、N I C 4 1 に伝送し、そして、適切な応答を待つように構成されている（それによって、応答が所定の時間内にないことが、N I C 4 1 の故障となる）。自己修復ケーブルアダプタ 3 7 の下流のコンポーネントの故障が、複数の異なる原因のうちの 1 つ以上を有する可能性がある場合、当業者にとって容易に明らかなように、自己修復ケーブルアダプタ 3 7 が、このようなイベントを検出するあらゆる適切な手段を実装できることは言うまでもない。

20

【 0 0 4 2 】

自己修復ケーブルアダプタ 3 7 は、下流のコンポーネントの故障を検出すると、適切な致命的でないエラーメッセージをC P U / R C 3 6 に供給するように、N I C 4 1 をエミュレートする。致命的でないエラーメッセージは、N I C 4 1 が利用可能であるが、それを用いることができない状態にあることを示す（例えば、イーサネットケーブル 4 3 が抜かれた状態をエミュレートする）。

30

【 0 0 4 3 】

N I C 4 1 をエミュレートするために、自己修復ケーブルアダプタ 3 7 は、N I C 4 1 によって生成されたように見えるデータパケットを生成することができる。より詳細には、N I C 4 1 は、複数の独立したデバイス機能、すなわちPCI Expressプロトコルによってサポートされている最大 8 つの機能を有することができる。すなわち、N I C 4 1 は、C P U / R C 3 6 に対して、最大 8 つの別々のデバイスであるように見えることができる。N I C 4 1 の各デバイス機能は、一意的にそれらの機能を識別する対応した識別子を有する。N I C 4 1 の特定のデバイス機能から送信されたデータパケットは、トランザクション識別子を有し、トランザクション識別子は、データパケットを送信しているデバイス機能の識別子に対応したリクエスト識別子を含んでいる。

40

【 0 0 4 4 】

PCI Expressプロトコルによって用いられるデータパケットヘッダのフォーマットを、図 5 を参照して説明する。リクエスト識別子 4 5 は、データパケットを始めたデバイス機能を識別し、0 から 1 5 までのインデックスが付けられた 1 6 ビットを含む。リクエスト識別子 4 5 は、上位 8 ビットを占めるバス番号フィールド 4 6 と、中央の 5 ビットを占め

50

るデバイス番号フィールド47と、下位3ビットを占める機能番号フィールド48とを含むことが分かる。PCI Expressプロトコルを用いるとき、バス番号46、デバイス番号47及び機能番号48の組合せは、特定のデバイスによって提供される機能を一意的に識別する。

【0045】

図5に示すパケットヘッダは、更に8ビットからなるタグフィールド49を含む。上述したように、トランザクションは、要求データパケットと、1つ以上の対応する完了データパケットとで構成されている。各要求データパケットは、タグフィールド49に記憶されている値に関連する。各対応する完了データパケットは、タグフィールド49に記憶されている値と同じ値を有し、したがって、完了データパケットを適切な要求データパケットと関連させる。一意的なタグ値が、宛先デバイスから1つ以上の完了データパケットを要求する全ての未処理の要求に割り当てられる。タグフィールド49は、8ビットを有する場合、 2^8 個の考えられるタグ値を表すことができる。

10

【0046】

機能番号フィールド48は、要求を送信するデバイスに関連した機能の機能番号が設けられる。デバイスが有する機能が8つより少ない場合、機能番号フィールド48には未使用ビットがあってもよい。したがって、デバイスの機能を表すのに機能番号フィールド48の十分なビットだけを用い、タグフィールド49と論理的に組み合わせる仮想機能番号として、機能番号フィールド48の幾つかの未使用ビットを用いることが、理解される。1つの機能だけが設けられている場合、機能番号フィールド48の全てのビットは、タグフィールド49と論理的に組み合わせて、最大 2^{11} 個の未処理要求のサポートを行うことができる。

20

【0047】

自己修復ケーブルアダプタ37は、受信データパケットのトランザクション識別子を判定するために、そのデータパケット、特にタグフィールド及びリクエストIDを調べることができる。このデータは、自己修復ケーブルアダプタが備えるオンボードメモリに記憶することができる。そして、自己修復ケーブルアダプタ37は（供給するエラーメッセージの種類がトランザクションIDを必要とする場合）、CPU/R36が、メッセージがCPU/R36とNIC41間の未処理のトランザクションに関するものと信じるように、正しいトランザクション識別子を与えることができるデータパケットを生成する。同様に、自己修復ケーブルアダプタ37は、NIC41のリクエストIDを有するデータパケットを生成することができ、それによって、それらのデータパケットは、NIC41によって生成されたように見える。

30

【0048】

適切なエラーメッセージは、自己修復ケーブルアダプタ37によって登録されたトランザクションに関する情報に基づいて選択される。例えば、CPU/R36がレジスタ読出し要求データパケットをNIC41に送信した場合、下流のコンポーネントが故障して、CPU/R36とNIC41間の通信が切断されると、自己修復ケーブルアダプタ37は、エラーステータスビットが「1」の値に設定されたデータパケットを供給することができる。サーバ1上で動作しているソフトウェアは、データパケットの受信を、エラーが発生したことを示すように状態ビットが設定されていると解釈するが、エラーは、サーバ35の故障の原因とはならない。更なる例として、トランザクションがネットワーク5上で伝送されているデータに関する場合、自己修復ケーブルアダプタ37は、イーサネット接続43が切断されたことを示すメッセージを、CPU/R36に送信することができる。各エラーメッセージは、自己修復ケーブルアダプタ37のオンボードメモリにハードコードされていてもよい。そして、自己修復ケーブルアダプタ37上で動作している論理は、完了データパケットを待っているトランザクションを考慮して、格納されているメッセージのうちのどれが適しているかを決定することができる。

40

【0049】

さらに、CPU/R36に供給するエラーメッセージの選択は、自己修復ケーブルア

50

アダプタ 37 がエミュレートしている遠隔被接続デバイスの詳細によっても異なってもよい。例えば、遠隔被接続デバイスがストレージリソースの場合、イーサネット接続が切断されたことを示すエラーメッセージは、適切でないことになる。エラーメッセージ自体は、そのデバイスが局所的には機能するが、その指定サービスを実行することができないことを示す（すなわち、致命的エラーを致命的でないエラーに置き換える）ように、遠隔被接続デバイスの種類に合わせることができる。

【 0 0 5 0 】

本発明の上述した実施の形態は、2つの前提に基づいており、第1に、サーバ35に遠隔で接続されたあらゆるデバイスは、サーバ35が使えないデバイスと考えるが、サーバ35がクラッシュする原因とはならない状態を有し、そして、第2に、使えない状態は、I/Oデバイスの完全な動作状態と比較して、簡単である。

10

【 0 0 5 1 】

さらに、あるいは、被接続サーバ（すなわちサーバへの接続）が故障したときに、遠隔被接続I/Oデバイスが故障することを防止することが、望ましい。これは、限定されないが、特に、遠隔I/Oデバイスが複数のサーバに接続し、複数のサーバによって共有されるシステム（例えば、マルチルートI/O仮想化を実装したシステム）に適している。1つのサーバの故障が、遠隔I/Oデバイスが故障する原因になってはならず、したがって、そのサービスを残りのサーバ（複数のサーバ）に提供することができないことは言うまでもない。

【 0 0 5 2 】

20

図6は、本発明の実施の形態に基づいて変更された図2の一般的配置を示す。特に、図6の構成では、遠隔I/O機器内のケーブルアダプタは、自己修復ケーブルアダプタに置き換えられている。

【 0 0 5 3 】

図6において、サーバ50は、PCI Express接続53によってケーブルアダプタ52に接続されたCPU/RCS1を含む。ケーブルアダプタ52は、CPU/RCS1をケーブル56によって、遠隔I/O機器55内に収容された自己修復ケーブルアダプタ54に接続する。また、遠隔I/O機器55は、NIC57を収容し、NIC57は、PCI Express接続58によって自己修復ケーブルアダプタ54に、イーサネット接続59によってネットワーク5に接続されている。

30

【 0 0 5 4 】

自己修復ケーブルアダプタ54は、自己修復ケーブルアダプタ54の上流のコンポーネントの故障を検出するように構成されている。自己修復ケーブルアダプタ54は、上流のコンポーネントの故障を検出すると、サーバ50の存在をエミュレートするように構成されている。例えば、自己修復ケーブルアダプタ54は、NIC57によって発行された未処理のメモリ読出し/書込み要求に対する応答において、直ちに、致命的でないエラー状態を有する完了データパケットを発行するように構成することができる。

【 0 0 5 5 】

自己修復ケーブルアダプタ37に関して上述したように、NIC57に供給されるエラーメッセージは、完了データパケットが必要とされるトランザクションによって異なる。

40

【 0 0 5 6 】

サーバ及び遠隔リソースの構成が、自己修復ケーブルアダプタ37と、自己修復ケーブルアダプタ54との両方を含むことができることは言うまでもない。すなわち、ある構成において、自己修復ケーブルアダプタは、1つ以上のサーバによって定義される1つ以上のそれぞれの故障ユニット内と、1つ以上のそれぞれの遠隔被接続リソースによって定義される1つ以上の故障ユニット内との両方に設けることができる。

【 0 0 5 7 】

上述の具体的な実施の形態では、サーバに遠隔で接続されたPCI Expressイーサネットネットワークインタフェースカード（NIC）を含むシステムに関して説明している。しかしながら、本発明は、より一般的に応用できることは言うまでもない。実際、本発明は

50

、PCI Express以外の内部接続を利用し、NIC以外の遠隔デバイス（例えば、ファイバチャンネルホストバスアダプタ、ストレージコントローラ等）を有するシステム内で用いることができる。

【0058】

さらに、上述の説明は、サーバと、そのサーバの遠隔被接続リソースとを含むシステムに関するものであるが、本発明は、第1のサーバと第2のサーバ間の通信に用いることができることは言うまでもない。

【0059】

さらに、本発明は、被接続デバイスの他の構成、例えばマルチサーバI/O仮想化システムに適用することができる。実際、多くの実施の形態において、本発明を利用したシステムが、複数の独立した遠隔デバイスに対する接続をそれぞれ有する複数のサーバを含むことは言うまでもない。このような構成により、1つの遠隔デバイスに関連した故障が生じた場合、各サーバは、他の遠隔被接続デバイスを用いることによって、有効な動作を続けることができる。複数のサーバ又は複数の遠隔リソースに設けられている場合、自己修復ケーブルアダプタ37、54は、（必要に応じて、故障したサーバ又は遠隔リソースのアクティブトランザクションを考慮して）適切な致命的でないエラーメッセージを選択するために、どのサーバ又はどの遠隔リソースが故障したかの記録を維持する。

【0060】

さらに、上述では、自己修復ケーブルアダプタを、本発明の機能と、標準のケーブルアダプタの機能とを含む単一のデバイスとして説明したが、本発明は専用デバイスに実装することができることは言うまでもない。例えば、図2において、本発明を実装したデバイスは、CPU/RIC10とケーブルアダプタ9間、又はケーブルアダプタ9とケーブル14間に配置することができる。

【0061】

自己修復ケーブルアダプタ37、54は、フィールドプログラマブルゲートアレイ（FPGA）又は特定用途向け集積回路（ASIC）で実現することができる。しかしながら、自己修復ケーブルアダプタ37、54は、あらゆる適切な手段を用いて実現してもよいことは言うまでもない。

【0062】

上述では、サーバとI/Oデバイス間でデータパケットを伝送する発明の実施の形態を説明した。用語サーバは、広範囲のものを意図するものであり、あらゆるコンピュータ装置を網羅することを意図することは言うまでもない。

【0063】

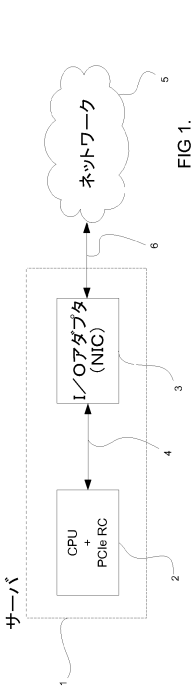
当業者にとって容易に明らかなように、本発明は、請求項の要旨を逸脱することなく、ここに教示する内容から変更及び応用することができる。

10

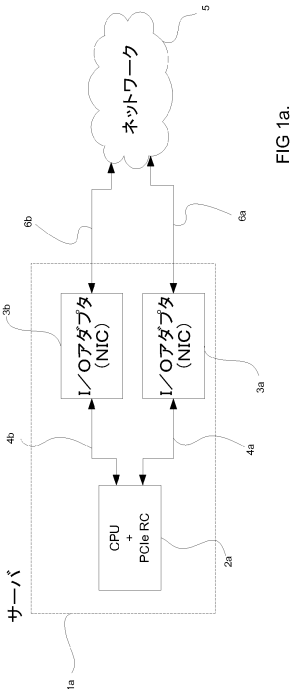
20

30

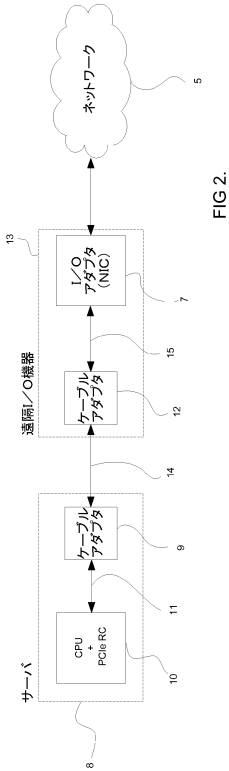
【図 1】



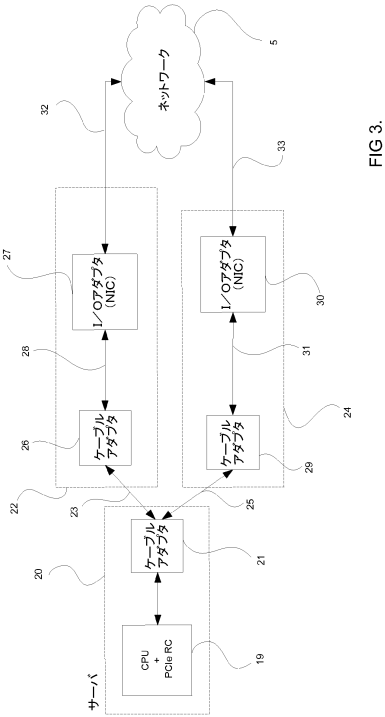
【図 1 a】



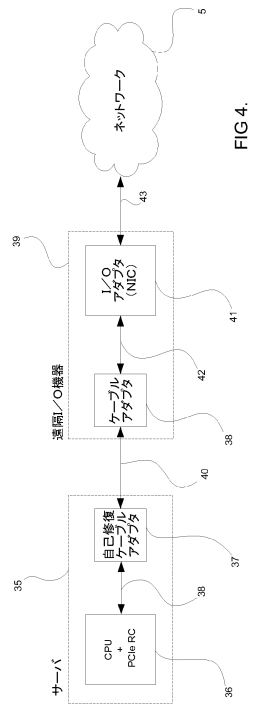
【図 2】



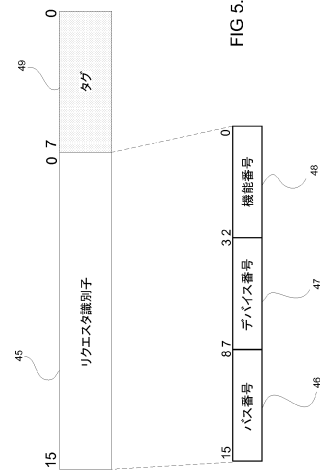
【図 3】



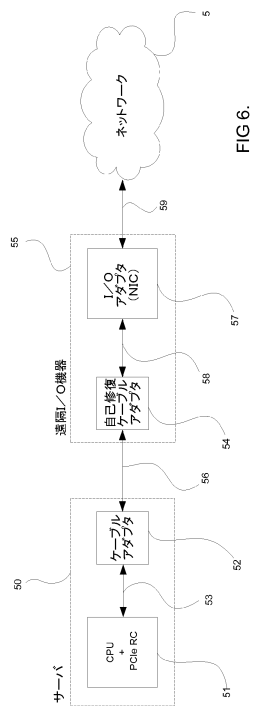
【 図 4 】



【圖 5】



【 図 6 】



フロントページの続き

(72)発明者 パイカルスキー、 マレク
イギリス国 エスケイ 1 1 8 ユーイー チェシャー マックルズフィールド ゴースワース ロ
ード 2 0

審査官 田中 啓介

(56)参考文献 特開 2 0 0 8 - 1 6 5 2 6 9 (J P , A)
特開 2 0 1 0 - 2 3 8 1 5 0 (J P , A)
特開 2 0 0 7 - 0 9 4 7 0 6 (J P , A)
特開 2 0 0 6 - 1 7 2 2 1 8 (J P , A)
特開 2 0 0 5 - 1 7 2 2 1 8 (J P , A)

(58)調査した分野(Int.Cl. , D B 名)
G 0 6 F 3 / 0 0、3 / 1 8
1 3 / 0 0、1 3 / 2 0 - 1 3 / 4 2