

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5111236号  
(P5111236)

(45) 発行日 平成25年1月9日(2013.1.9)

(24) 登録日 平成24年10月19日(2012.10.19)

(51) Int.Cl. F 1  
**G 0 6 F 1 3 / 3 6 ( 2 0 0 6 . 0 1 )** G 0 6 F 1 3 / 3 6 3 1 0 D

請求項の数 4 (全 9 頁)

<p>(21) 出願番号 特願2008-134509 (P2008-134509)                  (22) 出願日 平成20年5月22日 (2008.5.22)                  (65) 公開番号 特開2009-282773 (P2009-282773A)                  (43) 公開日 平成21年12月3日 (2009.12.3)                  審査請求日 平成22年6月8日 (2010.6.8)</p>	<p>(73) 特許権者 000005108                  株式会社日立製作所                  東京都千代田区丸の内一丁目6番6号                  (74) 代理人 110000442                  特許業務法人 武和国際特許事務所                  (72) 発明者 関 辰一郎                  神奈川県秦野市堀山下1番地 株式会社                  日立製作所 エンタープライズサーバ事業                  部内                  (72) 発明者 小国 哲                  神奈川県秦野市堀山下1番地 株式会社                  日立製作所 エンタープライズサーバ事業                  部内                  審査官 坂東 博司</p>
--	--

最終頁に続く

(54) 【発明の名称】 データ転送システム

(57) 【特許請求の範囲】

【請求項1】

PCI Express リンクを用いたデータ転送システムであって、  
 ルートコンプレックスと、PCI Express デバイスに接続されたアドインカードのそれぞれが接続される少なくとも1つのスロットと、パケットスイッチと、経路を選択する経路選択手段とを備え、

前記経路選択手段は、前記ルートコンプレックスと前記スロットとの間のデータ転送用の経路について、前記パケットスイッチを介する経路と、前記パケットスイッチを経由しない経路とを、前記アドインカードの性能に応じて選択することを特徴とするデータ転送システム。

【請求項2】

前記パケットスイッチと前記経路選択手段とがFPGA上に実装されたことを特徴とする請求項1記載のデータ転送システム。

【請求項3】

前記パケットスイッチと前記経路選択手段とがASIC上に実装されたことを特徴とする請求項1記載のデータ転送システム。

【請求項4】

前記経路選択手段は、外部から与えられるセレクト信号により、前記経路の選択を行うことを特徴とする請求項1記載のデータ転送システム。

【発明の詳細な説明】

## 【技術分野】

## 【0001】

本発明は、データ転送システムに係り、特に、用途に最適な経路を設定してデータ転送を制御するデータ転送システムに関する。

## 【背景技術】

## 【0002】

近年、情報処理装置に採用されるI/Oインターフェースは、性能要件の向上に伴い、PCI (TM)、PCI-X (TM)等のパラレルバスから、PCI Express (TM)等の高速シリアルインターフェースに置き換わりつつある。

## 【0003】

PCI Express は、従来のパラレルバスとは異なり、デバイス間を1対1で接続することができるため、ホストが備えるポート数以上のスロット数を実装したい場合等にしばしばPCI Express スイッチが使用される。

## 【0004】

PCI Express スイッチは、スロット数の拡張だけではなく、複数のホストによるI/O資源の共有、I/Oの冗長構成、I/Oの仮想化等の幅広い応用が可能である。なお、この種のPCI Express スイッチの応用に関する各種の技術が、例えば、非特許文献1等に記載されて知られている。

【非特許文献1】“PLX PCI Express Presentation”, October, 26 2007, PLX Technology Inc

## 【発明の開示】

## 【発明が解決しようとする課題】

## 【0005】

前述したように、PCI Express スイッチは、様々な応用が期待されるが、一方で、PCI Express カードやそれに関連するソフトウェアが、PCI Express スイッチを経由する使用方法に対応していない場合や、PCI Express スイッチを経由する際に増加するレイテンシが性能上無視できない場合等、PCI Express スイッチを経由したくない場合がある。

## 【0006】

本発明の目的は、前述したような点に鑑み、PCI Express スイッチを含むシステムにおいて、用途に最適な経路を設定することを可能にしたデータ転送システムを提供すること

## 【課題を解決するための手段】

## 【0007】

本発明によれば前記目的は、PCI Express リンクを用いたデータ転送システムであって、ルートコンプレックスと、PCI Express デバイスに接続されたアドインカードのそれぞれが接続される少なくとも1つのスロットと、パケットスイッチと、経路を選択する経路選択手段とを備え、前記経路選択手段は、前記ルートコンプレックスと前記スロットとの間のデータ転送用の経路について、前記パケットスイッチを介する経路と、前記パケットスイッチを経由しない経路とを、前記アドインカードの性能に応じて選択することにより達成される。

## 【発明の効果】

## 【0008】

本発明によれば、ユーザは、用途に応じた最適なデータ転送経路を選択して使用することが可能となる。

## 【発明を実施するための最良の形態】

## 【0009】

以下、本発明によるデータ転送システムの実施形態を図面により詳細に説明する。

## 【0010】

図1は本発明の第1の実施形態によるデータ転送システムの構成例を示すブロック図である。図1に示す本発明の第1の実施形態によるデータ転送システムは、3ポートのPCI

10

20

30

40

50

Express スイッチ（以下、単に、スイッチという）を、メインボード上に搭載して構成した例である。なお、PCI Express スイッチは、一般に、マルチプレクサ等のスイッチとの区別のために、パケットスイッチと呼ばれている。

#### 【0011】

そして、図1(a)に示す構成例は、メインボード100におけるルートコンプレックス110からスロット140a及びスロット140bまでの経路が、スイッチ130を経由するように設定が施された構成例を示し、図1(b)に示す構成例は、メインボード100におけるルートコンプレックス110からスロット140a及びスロット140bまでの経路が、スイッチ130を迂回するように設定が施された構成例を示している。また、図1(a)、図1(b)に示す例は、スイッチ130を経由する経路を作成するか否かを制御するために、ルートコンプレックス110とスロット140a及びスロット140bとの間に、マルチプレクサ120a、スイッチ130、マルチプレクサ120b、120cが設けられて図示のように構成されている。

10

#### 【0012】

前述において、ルートコンプレックス110は、PCI規格における最上位の親となるコントローラであり、一般には、図示しないホストコンピュータ等の情報処理装置に接続される。また、スロット140a、140bには、図1(a)の場合、アドインカード200a、200bが経路150j、150iを介して接続され、図1(b)の場合、アドインカード200c、200dが経路150j、150iを介して接続されている。これらのアドインカードには、ネットワーク、記憶装置、I/O等を含む図示しないPCI Express デバイスが接続され、また、これらのアドインカードは、接続されるPCI Express デバイスに対するコントローラの機能を有する。なお、ルートコンプレックス及びアドインカードの機能等については、後述する本発明の他の実施形態においても前述と同様である。

20

#### 【0013】

図1(a)に示す例において、ルートコンプレックス110のリンクは、16レーン幅の1ポートに設定され、経路150aを経てマルチプレクサ120aに接続される。また、マルチプレクサ120aは、ルートコンプレックス110とスイッチ130との接続を16レーン幅で経路150a、150bを経てリンクするように設定する。マルチプレクサ120bは、スイッチ130とアドインカード200aとの接続をスロット140aを介して、経路150f、150h、150jを経て16レーン幅でリンクするように設定する。マルチプレクサ120cは、スイッチ130とアドインカード200bとの接続をスロット140bを介して、経路150e、150g、150iを経て16レーン幅でリンクするように設定する。各マルチプレクサの設定は、共通のセレクト信号160を用いて行われる。このセレクト信号160は、図示しないボード管理用のコントローラを介して、ホストコンピュータから指示されるものであっても、また、図示しないサービスプロセッサを用いて利用者等により指示されるものであってもよく、後述する本発明の第2及び第4の実施形態の場合も同様である。

30

#### 【0014】

図1(b)に示す例において、ルートコンプレックス110のリンクは、8レーン幅の2ポートに分割され、経路150k、150lを経てマルチプレクサ120aに接続される。また、マルチプレクサ120aは、ルートコンプレックス110とスロット140aとを、マルチプレクサ120a及びマルチプレクサ120bを介して8レーン幅で経路150l、150d、150hを経てスイッチ130を介することなく接続するように直接リンクさせ、ルートコンプレックス110とスロット140bとを、マルチプレクサ120a及びマルチプレクサ120cを介して8レーン幅で経路150k、150c、150gを経てスイッチ130を介することなく接続するように直接リンクさせている。

40

#### 【0015】

前述で説明した図1(a)、図1(b)に示す本発明の第1の実施形態によれば、ルートコンプレックス110に接続された図示しないホストコンピュータは、用途に応じて、

50

スイッチ 130 を介したデータ転送経路を経由して、アドインカードに接続されたPCI Express デバイスを使用する経路と、スイッチ 130 を介することのないデータ転送経路を経由して、アドインカードに接続されたPCI Express デバイスを使用する経路とを使い分けることができる。

【0016】

図2は図1に示すPCI Express スイッチ 130 の内部構成を示すブロック図である。

【0017】

図2に示すように、PCI Express スイッチ 130 は、図1(a)、図1(b)に示しているルートコンプレックス 110 とリンクするポート 131 a と、スロット 140 a、140 b とリンクし、アドインカードに接続されるPCI Express デバイスとリンクするポート 131 b、131 c と、これらのポート 131 a ~ 131 c 相互間の接続を制御するスイッチ論理 135 とを備えている。また、図2に示すスイッチ 130 は、前述のポート 131 a ~ 131 c のそれぞれと、スイッチ論理 135 との間には、物理層 132 a ~ 132 c、データリンク層 133 a ~ 133 c、トランザクション層 134 a ~ 134 c が設けられている。

【0018】

前述において、物理層 132 a ~ 132 c は、転送データのシリアル化・非シリアル化、8ビットデータを10ビットデータの時間幅で転送するようにPCI規格による10b/8b変換等の処理を行う。データリンク層 133 a ~ 133 c は、主にリンクの管理とエラー検出・訂正との処理を行う。トランザクション層 134 a ~ 134 c は、トランザクションレイヤパケット(TLP)の分解・生成の処理を行うと共に、相手側のデバイスとのパケット交換のフローコントロールの処理を担う。スイッチ論理 135 は、各ポートのトランザクション層間のパケットをルーティングする。そして、スイッチ論理 135 は、ポート 131 a とポート 131 b との間、ポート 131 a とポート 131 c との間のデータ転送に限らず、ポート 131 b とポート 131 c との間のいわゆるpeer-to-peer転送も行うことが可能である。

【0019】

前述したように構成されるスイッチ 130 は、一般に、このスイッチによるパケットのルーティングを行う際に、前述した各階層の処理に起因するレイテンシが増加することが知られている。

【0020】

図3はスイッチ 130 のソフトウェアから見たトポロジを示す図である。

【0021】

PCI Express は、ソフトウェアの互換性を維持するため、PCI互換のコンフィグレーションメカニズムが採用される。各ポート 131 a ~ 131 c は、仮想PCIバス 137 a ~ 137 c を介して仮想PCI to PCIブリッジ 136 a ~ 136 c に接続される。仮想PCI to PCIブリッジ相互間の通信は、仮想PCIバス 138 上で行われる。コンフィグレーションソフトウェアは、各仮想PCIバス 137 a ~ 137 c、138 のそれぞれが異なるバス番号となるように、各仮想PCI to PCIブリッジのコンフィグレーション空間レジスタを設定する。

【0022】

図1(a)に示して説明した構成例では、ルートコンプレックス 110 とスイッチ 130 との間、及び、2つのアドインカード 200 a、200 b とスイッチ 130 との間が、全て16レーン幅でリンクしている。ルートコンプレックスと両アドインカードとの間の合計スループットは、ルートコンプレックス 110 とスイッチ 130 との間の16レーン幅に制限されるものの、アドインカード 200 a とルートコンプレックス 110 との間、アドインカード 200 b とルートコンプレックス 110 との間の通信が重ならない場合、それぞれ16レーン幅のスループットを享受することができる。また、図1(a)に示して説明した構成例の場合の最大スループットは、図1(b)に示して説明した構成例の場合における8レーン幅と比べて、アドインカード同士がpeer-to-peer転送をサポートする

10

20

30

40

50

場合のアドインカード200aとアドインカード200bとがスイッチ130を介して通信を行う16レーン幅となる。

【0023】

図1(b)に示して説明した構成例では、ルートコンプレックス110とアドインカード200c、200dが直接リンクするため、図1(a)に示した構成例の場合のように、スイッチ130を経由する際のレイテンシを増加させることがない。このため、アドインカードの性能がレイテンシに左右されやすい場合には、図1(a)に示す構成例のものを使用する場合に比較して、図1(b)に示す構成例のものを使用の方が性能面で有利となる。また、アドインカードを扱うソフトウェアがスイッチ130による仮想PCI-to-PCIブリッジを介した使用をサポートしていないような場合には、システムの構成を図1(b)に示す構成例とすることによりアドインカードを使用することが可能となる。

10

【0024】

前述した本発明の第1の実施形態によるデータ転送システムは、3ポートのPCI Express スイッチを用いるとしたが、この実施形態は、2ポートのPCI Express スイッチを用いても、あるいは、3ポート以上のPCI Express スイッチを用いてもよい。

【0025】

図4は本発明の第2の実施形態によるデータ転送システムの構成例を示すブロック図である。図4に示す本発明の第2の実施形態は、I/O仮想化機構を含む2ポートのスイッチをライザーカードに搭載した場合の構成例である。

【0026】

一般に、仮想化環境での性能を向上させる手段として、I/O仮想化機構が知られている。そして、I/O仮想化機構に対応していない既存のハードウェアを最小限の変更で対応させる方法として、メインボード上に搭載するライザーカードの追加あるいは変更という形態がある。本発明の第2の実施形態は、この場合の構成例である。

20

【0027】

そして、図4(a)に示す構成例は、メインボード300上に備えられたルートコンプレックス310とスロット320とがリンクするように経路330aにより接続され、ライザーカード400上に備えられたマルチプレクサ410aからスロット430までの経路が、スイッチ420を経由するように設定が施された構成例を示し、図4(b)に示す構成例は、メインボード300上に備えられたルートコンプレックス310とスロット320とがリンクするように経路330aにより接続され、ライザーカード400上に備えられたマルチプレクサ410aからスロット430までの経路が、スイッチ420を迂回するように設定が施された構成例を示している。また、図4(a)、図4(b)に示す例は、スイッチ420を経由する経路を作成するか否かを制御するために、メインボード300に搭載されるライザーカード400上にマルチプレクサ410a、スイッチ420、マルチプレクサ410b、スロット430が設けられて図示のように構成されている。

30

【0028】

図4(a)に示す構成例は、メインボード300と、ライザーカード400と、アドインカード200eとから構成される。ライザーカード400上のマルチプレクサ410a、410bは、スイッチ420を経由してメインボード300上のスロット320とライザーカード400上のスロット430とを接続するように経路330b、440a、440c、440dを設定する。また、図4(b)に示す構成例は、メインボード300と、ライザーカード400と、アドインカード200fとから構成される。ライザーカード400上のマルチプレクサ410a、410bは、スイッチ420を迂回してメインボード300上のスロット320とライザーカード400上のスロット430とを直接接続するように経路330b、440b、440dを設定する。

40

【0029】

本発明の第2の実施形態は、前述のようにしてアドインカード200e、200fにスイッチ420を経由する経路、あるいは、迂回する経路を提供することにより、I/O仮想化機構を使用する場合に図4(a)に示す構成とし、不要な場合やアドインカードやソ

50

フトウェアの制約により使用不可である場合に図4(b)に示す構成とすることにより、従来のシステムと互換性を保ちつつI/O仮想化機構を導入できるという効果を得ることができる。

【0030】

前述した本発明の第2の実施形態によるデータ転送システムは、2ポートのPCI Express スイッチを用いるとしたが、この実施形態は、3ポート以上のPCI Express スイッチを用いてもよい。

【0031】

図5は本発明の第3の実施形態によるデータ転送システムの構成例を示すブロック図である。図5に示す本発明の第3の実施形態は、図4により説明した本発明の第2の実施形態と同等の機能をFPGAにより実現した構成例である。

10

【0032】

FPGAは、論理を読み込んで論理を再構成することが可能なLSIであり、起動される毎に、その都度その機能を任意に変更可能なものである。本発明の第3の実施形態は、ライザーカード500上に、図4に示したスイッチ420と同等のスイッチ511を含むFPGA510を搭載して構成される。そして、この本発明の第3の実施形態は、FPGA510の初期化時に、読み込むROMの変更により、図5(a)に示すように、スイッチ511を経由する経路と、図5(b)に示すように、スイッチ511を迂回する経路との切り替えが可能である。

【0033】

20

前述したように構成される本発明の第3の実施形態においても、図4により説明した本発明の第2の実施形態と同等な効果を得ることができる。また、この本発明の第3の実施形態においても、スイッチ511として、3ポート以上のPCI Express スイッチを用いることができる。

【0034】

図6は本発明の第4の実施形態によるデータ転送システムの構成例を示すブロック図である。図6に示す本発明の第4の実施形態は、図5により説明した本発明の第3の実施形態と同等の機能をASICにより実現した構成例である。

【0035】

本発明の第4の実施形態は、ライザーカード600上にASIC610を搭載し、このASIC610上に、第2の実施形態の場合と同様に、マルチプレクサ611a、611b、及び、スイッチ612を実装して構成した例であり、セレクト信号614により、スイッチ612を経由する経路613a、613b、613d、613eと、スイッチ612を迂回する経路613c、613eとの選択を行うことが可能である。また、この本発明の第4の実施形態においても、スイッチ612として、3ポート以上のPCI Express スイッチを用いることができる。

30

【0036】

前述した本発明の各実施形態によれば、PCI Express を採用したシステムにおいて、用途に応じてPCI Express スイッチを経由する経路と、PCI Express スイッチを経由しない経路とを選択することができるため、単一のシステムで両方の経路が有するメリットを享受することができる。

40

【図面の簡単な説明】

【0037】

【図1】本発明の第1の実施形態によるデータ転送システムの構成例を示すブロック図である。

【図2】PCI Express スイッチの内部構成を示すブロック図である。

【図3】PCI Express スイッチのソフトウェアから見たトポロジを示す図である。

【図4】本発明の第2の実施形態によるデータ転送システムの構成例を示すブロック図である。

【図5】本発明の第3の実施形態によるデータ転送システムの構成例を示すブロック図で

50

ある。

【図6】本発明の第4の実施形態によるデータ転送システムの構成例を示すブロック図である。

【符号の説明】

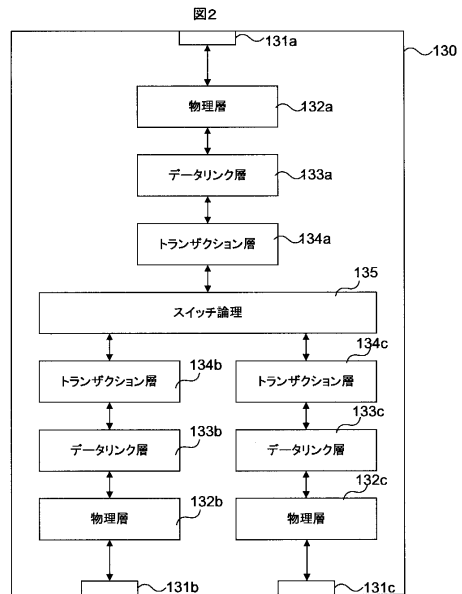
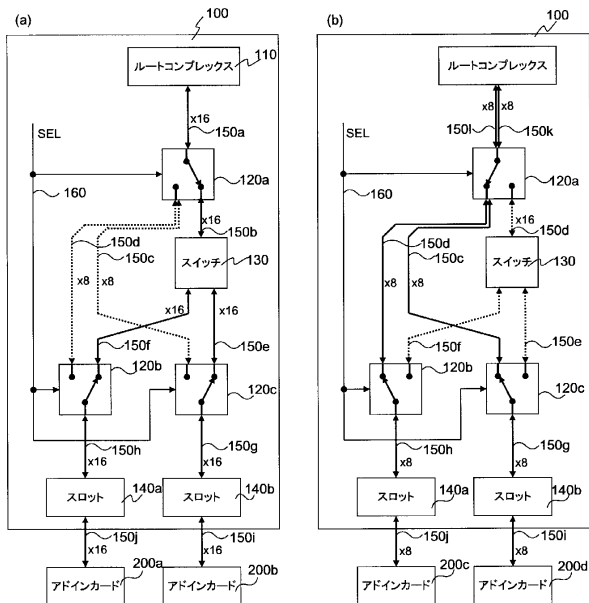
【0038】

- 100、300   メインボード
- 110、301、310   ルートコンプレックス
- 120a ~ 120c、410a、410b、611a、611b   マルチプレキサ
- 130、420、511、612   PCI Express スイッチ
- 140a、140b、320、430、520、620   スロット
- 200a ~ 200f   アドインカード
- 400、500、600   ライザカード
- 510   FPGA
- 610   ASIC

【図1】

【図2】

図1





---

フロントページの続き

- (56)参考文献 特開2005 - 166028 (JP, A)  
特開2004 - 135106 (JP, A)  
特開2003 - 69619 (JP, A)  
特開2003 - 69637 (JP, A)  
特開2004 - 56590 (JP, A)  
特開2001 - 320420 (JP, A)

- (58)調査した分野(Int.Cl., DB名)  
G06F 13/36