

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号  
特許第4909384号  
(P4909384)

(45) 発行日 平成24年4月4日 (2012.4.4)

(24) 登録日 平成24年1月20日 (2012.1.20)

(51) Int.Cl.	F I
H O 4 L 29/10 (2006.01)	H O 4 L 13/00 3 O 9 Z
G O 6 F 13/42 (2006.01)	G O 6 F 13/42 3 5 O C
G O 6 F 13/38 (2006.01)	G O 6 F 13/38 3 2 O A
H O 4 L 12/56 (2006.01)	H O 4 L 12/56 3 O O Z

請求項の数 19 外国語出願 (全 15 頁)

(21) 出願番号	特願2009-165061 (P2009-165061)	(73) 特許権者	591003943 インテル・コーポレーション 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 ブレーバード・2200
(22) 出願日	平成21年7月13日 (2009.7.13)	(74) 代理人	110000877 龍華国際特許業務法人
(65) 公開番号	特開2010-88106 (P2010-88106A)	(72) 発明者	ハリマン、デービット ジェイ. 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 ブレーバード・2200 インテル ・コーポレーション内
(43) 公開日	平成22年4月15日 (2010.4.15)		
審査請求日	平成21年7月14日 (2009.7.14)		
(31) 優先権主張番号	12/218,410		
(32) 優先日	平成20年7月15日 (2008.7.15)		
(33) 優先権主張国	米国 (US)		
		審査官	森谷 哲朗
			最終頁に続く

(54) 【発明の名称】 プロトコルスタックのタイミングの管理

(57) 【特許請求の範囲】

【請求項 1】

第 1 のプロトコルにしたがってデータを処理する第 1 のプロトコルスタックであって、  
前記第 1 のプロトコルスタックをトンネリングインターコネクにインターフェースする  
インターフェースロジックを有する前記第 1 のプロトコルスタックと、

前記第 1 のプロトコルスタックをリンクに結合する前記トンネリングインターコネク  
と

を備え、

前記インターフェースロジックは、前記トンネリングインターコネクによるトンネリ  
ングによって発生するタイミング遅延に少なくとも部分的に基づいて前記第 1 のプロトコ  
ルスタックの少なくとも 1 つのタイマを制御する

装置。

【請求項 2】

前記インターフェースロジックは、前記タイミング遅延に対応付けられているタイミン  
グ遅延情報を、前記第 1 のプロトコルスタックの少なくとも 1 つのスタックロジックの、  
タイミングについての制約を示すタイミング要件に対応付ける

請求項 1 に記載の装置。

【請求項 3】

前記インターフェースロジックは、前記タイミング遅延情報、及び前記タイミング遅延  
情報に対応付けられている前記タイミング要件に少なくとも部分的に基づいて、前記第 1

のプロトコルスタックの前記少なくとも1つのタイマを制御することによって、前記第1のプロトコルスタックのタイミングを変更するか否かを判断する

請求項2に記載の装置。

【請求項4】

前記インターフェースロジックは、前記タイミング遅延情報を前記タイミング要件に動的に対応付けして、前記第1のプロトコルスタックは、前記トンネリングインターコネクットの共通物理層または別の物理層に動的に結合される

請求項2に記載の装置。

【請求項5】

前記インターフェースロジックは、前記第1のプロトコルスタックの第1のスタックロジックが、前記第1のプロトコルの、タイミングについての制約を示すリンクタイミング要件を満たすべく、前記第1のスタックロジックに第1のクロック信号を供給する第1のクロックを、予め定められた時間にわたってディセーブルする

請求項1に記載の装置。

【請求項6】

前記トンネリングインターコネク트는、前記トンネリングインターコネク트의プロトコルによって、前記第1のプロトコルのパケットを前記リンクにトンネリングする

請求項5に記載の装置。

【請求項7】

前記リンクは、前記第1のプロトコルスタックと第2のプロトコルスタックとの間で共有する統合型インターコネクであり、前記第1のプロトコルスタックは、P C I e ( P e r i p h e r a l C o m p o n e n t I n t e r c o n n e c t E x p r e s s (登録商標))スタックである

請求項1から6のいずれか1項に記載の装置。

【請求項8】

前記トンネリングインターコネク트는、前記第1のプロトコルスタックに第1のスロットおよび第2のスロットを割り当てて、前記第2のプロトコルスタックに第3のスロットを割り当てる

請求項7に記載の装置。

【請求項9】

前記リンクに結合されており、前記トンネリングインターコネクによりトンネリングされたパケットを受信する受信機

をさらに備え、

前記受信機は、前記リンクに結合されているインターフェースロジックによって、割り当てられた前記第1のスロットおよび前記第2のスロットを利用する

請求項8に記載の装置。

【請求項10】

トンネリングインターコネク트에結合されている第1のプロトコルスタックのインターフェースロジックにおいて通信を受信する段階と、

前記通信の通信種類に基づいてテーブルにアクセスして、前記通信種類に対応付けられている、前記トンネリングインターコネクによるトンネリングによって発生するタイミング遅延を示すタイミング遅延情報を取得する段階と、

前記タイミング遅延情報が示す遅延に対応するべく、前記通信の通信種類について前記第1のプロトコルスタックの少なくとも1つのスタックロジックのタイミングを制御するべきか否かを決定する段階と、

前記決定する段階において変更するべきと決定した場合に、前記第1のプロトコルスタックの前記少なくとも1つのスタックロジックのタイミングを調整して、前記遅延に対応する段階と、

調整された前記タイミングを用いて前記第1のプロトコルスタックにおいて前記通信を処理する段階と

10

20

30

40

50

を備える方法。

【請求項 1 1】

前記テーブルは、不揮発性メモリに格納されており、前記トンネリングインターコネクと前記第 1 のプロトコルスタックとの間のマッピングを含む第 1 の部分と、前記トンネリングインターコネクと前記トンネリングインターコネクに結合されている第 2 のプロトコルスタックとの間のマッピングを含む第 2 の部分とを有する

請求項 1 0 に記載の方法。

【請求項 1 2】

前記第 1 のプロトコルスタックと前記第 2 のプロトコルスタックとの間で前記トンネリングインターコネクを共有する段階と、

前記第 1 のプロトコルスタックがあるスロットにおいて通信対象の情報を持たない場合に、前記第 2 のプロトコルスタックに前記第 1 のプロトコルスタックの前記スロットを与える段階と

をさらに備える、請求項 1 1 に記載の方法。

【請求項 1 3】

前記少なくとも 1 つのスタックロジックの前記タイミングは、前記少なくとも 1 つのスタックロジックに結合されているクロックをオフにすることによって調整される

請求項 1 0 から 1 2 のいずれか 1 項に記載の方法。

【請求項 1 4】

前記遅延に少なくとも部分的に基づいて、第 2 のスタックロジックのクロックを遅延させる段階

をさらに備える、請求項 1 3 に記載の方法。

【請求項 1 5】

前記クロックを遅延して、予め定められた時間が経過するまで、前記トンネリングインターコネクによって時間要件が満たされない場合にトリガされる、受信機からの肯定応答を受領していない旨を指し示すエラー信号が発行されないようにする段階

をさらに備える、請求項 1 4 に記載の方法。

【請求項 1 6】

リンクに結合されている物理層、および、前記物理層に結合されているプロトコルスタックを有する送信機と、

前記リンクを介して前記送信機に結合されており、第 1 のプロトコルにしたがってデータを処理する第 1 のプロトコルスタックを有する受信機と、

前記受信機に結合されているダイナミックランダムアクセスメモリ ( D R A M ) とを備え、

前記第 1 のプロトコルスタックは、トンネリング物理層を介して、前記リンクに前記第 1 のプロトコルスタックをインターフェースする第 1 のインターフェースロジックを含み、前記第 1 のインターフェースロジックは、前記トンネリング物理層によるトンネリングによって発生するタイミング遅延に少なくとも部分的に基づいて、前記第 1 のプロトコルスタックの少なくとも 1 つの第 1 のスタックロジックのタイミングを変更する

システム。

【請求項 1 7】

前記受信機はさらに、第 2 のプロトコルにしたがってデータを処理する第 2 のプロトコルスタックを有し、前記第 2 のプロトコルスタックは、前記タイミング遅延に少なくとも部分的に基づいて前記第 2 のプロトコルスタックの少なくとも 1 つの第 2 のスタックロジックのタイミングを変更する第 2 のインターフェースロジックを含む

請求項 1 6 に記載のシステム。

【請求項 1 8】

前記トンネリング物理層は、前記第 1 のプロトコルスタックまたは前記第 2 のプロトコルスタックを選択して前記送信機からのパケットを受信させるコントローラを含む

請求項 1 7 に記載のシステム。

10

20

30

40

50

## 【請求項 19】

前記第1のインターフェースロジックは、前記トンネリング物理層を介して前記送信機から受信したパケットの通信種類に基づいてテーブルにアクセスして、前記通信種類に対応付けられている、タイミング遅延情報を取得し、前記タイミング遅延情報が示す遅延に対応するべく、前記通信種類について前記第1のプロトコルスタックの前記少なくとも第1のスタックロジックのタイミングを変更するべきか否かを決定する

請求項16から18のいずれか1項に記載のシステム。

## 【発明の詳細な説明】

## 【背景技術】

## 【0001】

コンピュータプラットフォームは通常、さまざまなインターコネクトを用いて互いに結合された半導体素子を多数備えている。このようなインターコネクトまたはリンクは、プロトコルが異なっていることが多く、リンク上で実行される通信は、リンクが異なると、実行速度が異なり、準拠するプロトコルも異なる。一部のシステムでは、入出力(I/O)プロトコルの通信は、別のインターコネクトにおいてトンネリングされ得る。トンネリングは一般的に、第1のプロトコルに準拠した通信を受け取り、第2のプロトコルにしたがって動作するインターコネクトを介して当該通信を提供することを含む。例えば、第1のプロトコルのパケットに第2のプロトコルのヘッダを適用して、当該パケットをインターコネクトを介して送信することによって、第1のプロトコルのパケットがトンネリングされる。多くの場合、このようなプロトコルトンネリングは非常に上位のレベルで実行され、これら2つのプロトコルが有するソフトウェア抽象化が同じである一方、プロトコル間で共有されているハードウェアはない。このように、ソフトウェア互換性、性能、製品化までの時間において、上述したようなトンネリングによる利点は非常に小さい。

## 【図面の簡単な説明】

## 【0002】

【図1】本発明の一実施形態に係る、共有物理層を介したリンクへのプロトコルスタックの接続を示すブロック図である。

## 【0003】

【図2】本発明の別の実施形態に係る、共有物理層に結合されている複数の通信スタックを備えるシステムを示すブロック図である。

## 【0004】

【図3】本発明の一実施形態に係る方法を説明するためのフローチャートである。

## 【0005】

【図4】本発明の別の実施形態に係る、プロトコルスタックのインターフェースを動作させる方法を説明するためのフローチャートである。

## 【0006】

【図5】本発明の一実施形態に係るシステムを示すブロック図である。

## 【発明を実施するための形態】

## 【0007】

さまざまな実施形態によると、1以上の既存のI/Oプロトコルが、比較的低位のレベルで、別のインターコネクトにおいてトンネリングされ得る。本明細書では、「別のインターコネクト」を「トンネリングインターコネクト」と呼ぶことにする。一実施形態によると、このようなインターコネクトの一例として統合型(converged)IO(CIO)が挙げられる。CIOは、PCI Express(登録商標)仕様書の基本仕様書バージョン2.0(2007年1月17日発行)に従った(以降では、PCIe(登録商標)仕様書と呼ぶ)、PCIe(Peripheral Component Interconnect Express)プロトコルの通信、または、同様のプロトコルの通信をトンネリングするべく用いることができる。CIOの場合、PCIeハードウェアスタックの多くは直接実装されるので、ソフトウェア互換性、性能、および製品化までの時間において、利点が得られる。すなわち、低位のレベルのトンネリングにおいて、トンネ

リングされるプロトコルスタックは大半が実装されている。これとは対照的に、高位のレベルのトンネリングの場合は、ソフトウェアアーキテクチャは保存されるが、トンネリングされるプロトコルからのパケット、符号化、または有線用プロトコルメカニズムを必ずしも用いるわけではない。このように低位のトンネリングによって、P C I e プロトコルスタックのパケットを、例えば、トンネリングパケットにC I Oヘッダを適合させることによって、C I Oインターコネクタを介してトンネリングすることができる。このように送信されたトンネリングパケットを受信機が受信すると、受信機のC I Oプロトコルスタックはヘッダを復号化して、受信機の対応P C I e プロトコルスタックにP C I e パケットを渡すことができる。

【0008】

10

しかし、統合型インターコネクタに対するこのようなアプローチによって、抽象化という、より高位で行われるプロトコルのトンネリングとは対照的に、低位トンネリングに起因する問題が生じる。つまり、プロトコルのタイミングについて制約が潜在的に課されていることが多い。このような制約は、インターコネクタプロトコルのトンネリングされていない元々のインスタンス化では満足させられていたことが自明であるが、インターコネクタプロトコルをトンネリングする場合には、トンネリングに用いられるインターコネクタによる遅延が発生するために、管理するのがより困難になり得る。このような遅延は、トンネリングインターコネクタ自身によって生じる場合もあり、または、その他のトンネリングされたプロトコルによるトラフィックによって生じる場合もある。

【0009】

20

本発明の実施形態は、トンネリングインターコネクタを介してトンネリングを実行する場合に、トンネリングされるプロトコルの明示的および暗示的なタイマを管理するメカニズムを提供する。本明細書で説明する実施形態ではC I Oを介してトンネリングされるP C I e プロトコルを一例として用いるが、本発明の範囲はこれに限定されるものではなく、同様の原理は、その他のトンネリングインターコネクタ、および、有線インターコネクタおよび無線インターコネクタを始めとする、トンネリングに用いられるその他のインターコネクタにも応用され得ると理解されたい。

【0010】

インターコネクタに関する、明示的および暗示的な、タイミング要件は、大きく分けて2つのカテゴリーに分類され得る。本明細書では、リンクタイミング要件および掛け時計タイミング要件と呼ぶ。リンクタイミング要件は、リンクプロトコル等の低位レベルに対応付けられており、リンク動作を円滑化して、検証に関する問題を最低限に抑えることを目的として通常設けられている。掛け時計タイミング要件は、高位レベルにおいて観察可能なイベント、例えば、オペレーティングシステム(O S)およびアプリケーションソフトウェアに対応付けられている。リンクタイミング要件は、プロトコルトンネリングによって生じる遅延に直接影響を受ける可能性があり、本発明の実施形態によって問題解決を図る要件である。通常、リンクタイミング要件は、約10マイクロ秒( $\mu s$ )未満のオーダーであり得るが、掛け時計タイミング要件は、約10マイクロ秒( $\mu s$ )よりも大きい。掛け時計タイミング要件は、プロトコルのトンネリングによって生じる比較的短時間(例えば、マイクロ秒)の遅延に影響されない程度に十分長い時間値(例えば、ミリ秒( $ms$ ))に対応付けられていることが多いので、プロトコルトンネリングには基本的に影響されない。また、このような要件は、特定のインターコネクタプロトコルを伝送するべく用いられるハードウェアメカニズムに関わらず(本来のもの、または、トンネリング用)等しく望ましい、アプリケーションソフトウェアにおいてユーザに見える動作停止( $stall$ )の発生を防ぐというような、特性に対応付けられている。

30

40

【0011】

以下に記載する表1は、P C I e に対応付けられている多数のタイミング要件を一覧にしたものであり、各要件と本開示との関係を示す。「説明」部分の記載は、P C I E x p r e s s (登録商標)仕様書から引用していることに留意されたい。

【表 1】

説明	種類	注
肯定応答/非肯定応答(Ack/Nak) 送信および再生タイマ	リンク	絶対的な要件- トンネリングインターコネクトに よって時間要件が満たされない場合、 (スプリアス)エラーがトリガされる
リンク状態ゼロスタンバイ(L0s) 呼び出しポリシー： 「ポートは送信レーン、所定の アイドル条件(以下参照)が7 $\mu$ s以下 の期間にわたって満たされた場合 には、L0s状態に遷移させなければ ならない	リンク	リンクが使用されていない場合には、 リンク電力管理をトリガする。 これは、実際にどれが利用されたか ではなく、トンネリングインター コネクト割り当てに従って時間を カウントする場合の一例である。
リンク状態1(L1)エントリ交渉- 「上流の構成要素が、アイドル状態を 4シンボル時間以下として、 このDLLPを繰り返し送る」	リンク	暗示的なタイミング要件： この場合は、PCIeスタックに対して トラフィックがどのように見えるかを 管理して、挿入された遅延をマスクする
フロー制御の更新	リンク	ガイドラインであって要件ではない
PCI電力管理(PM)およびアクティブ 状態電力管理(ASPM)： 「L1を終了すると、下流の構成要素が、 フロー制御更新データリンク層 パケット、データリンク層パケット (DLLP)を、L1終了から1 $\mu$ s内に 開始された全てのイネーブル された仮想チャネル(VC)および フロー制御(FC)型に対して、 送信することが推奨される	リンク	ガイドラインであって要件ではない
L0s/L1終了レイテンシ	掛け 時計	これらのタイミングパラメータは、 PCIeを介したトラフィックに対する 影響を決定できるように設けられた ものであり、基本的なPCIeリンク動作 ではない
電力管理イベント(PME)- 「100ms(+50%/-5%)後に、要求元 エージェントのPME_Statusビットが クリアされていない場合、 PMEサービスタイムアウト メカニズムが切れて、PME要求元 エージェントをトリガして、 一時的に失われたPM/PME メッセージを再送させる」	掛け 時計	PMEが完全に失われてしまわないように 設けられている要件である； スプリアストリガを最低限に抑える ように特定の時間を選択すると同時に、 PMEが比較的タイミング良く処理される ように十分に短く設定される
出された要求に対する受諾の制限 (10 $\mu$ s)	掛け 時計	ファブリックの混雑によって生じる プラットフォームの観察可能な遅延を 制限するためのもの
フロー制御の最低更新頻度(30 $\mu$ s) および更新FCPタイマ-200 $\mu$ s	掛け 時計	フロー制御パケットの損失によって 生じる停止を制限するためのもの

表 1 は、一部を例示することを目的とするものであって、PCIeにおけるタイミング  
関連の要件を全て完全に網羅したものではないことに留意されたい。

10

20

30

40

50

## 【 0 0 1 2 】

リンクタイミング要件は、P C I e スタック自身によって「測定」されるので、P C I e スタックの時間に対する概念が変わると、このようなタイミング要件を感知する方法も変化し得る。このような構成を実現するべく、時間変更を実現するためのメカニズムが、スタックタイミングをいつ、そして、どのように変更するかを決定するためのハードウェア、ソフトウェアまたはファームウェアと共に、提供されるとしてよい。

## 【 0 0 1 3 】

P C I e スタックの時間の概念を変更するメカニズムは、複数の方法で実現することができる。例えば、P C I e スタックロジックにおけるさまざまな要素に対するクロックをゲーティングまたはオフに制御して、対象となるロジックの時間を止めることによって、実現することができる。尚、この方法によれば、利用されていないP C I e スタックロジックの電力消費を低減するという効果もさらに得られる。ほかの実施形態によると、時間をカウントすべき場合を指し示す明示的な制御信号を、P C I e スタックロジックに追加することができる。スタック全体を1つの単位として制御することは一般的に十分ではなく、スタックのサブ構成要素は、プロトコルメカニズムが異なれば、異なるロジックブロックに対する影響も異なるので、半独立制御が実行され得る。同様に、リンクタイミング要件を遵守するべく、全ての通信についてタイミング変更が必要となるわけではない。一実施形態によると、制御ロジックを用いてP C I e スタックの時間の概念をいつおよびどのように調整するかを決定するとしてよく、当該ロジックはP C I e スタックの一部であってよい。

## 【 0 0 1 4 】

図1は、P C I e スタック（およびその他のトンネリングされたプロトコル）が共有トンネリングリンクにインターフェースされる様子を示すブロック図である。一実施形態によると、共有トンネリングリンクは、C I O リンクであってよい。図1に示すように、システム10は、第1のスタック20aと、第2のスタック20b（プロトコルスタック20と総称する）とを備える。一実施形態によると、第1のプロトコルスタック20aは、P C I e スタックであってよく、第2のプロトコルスタック20bは、ユニバーサルシリアルバス（U S B）、ディスプレイインターコネクト、またはその他の同様のプロトコルスタックであってよい。図示の便宜上、P C I e プロトコルスタックの詳細のみを示す。具体的には、プロトコルスタック20aは、トランザクション層22と、データリンク層24と、インターフェースまたはバスケット層26とを有する。インターフェースまたはバスケット層26は、P C I e プロトコルとトンネリングプロトコルとの間のインターフェースとして機能する。このようなインターフェースロジックの動作の詳細は後述する。

## 【 0 0 1 5 】

図1に示すように、統合型I O層は、第1および第2のプロトコルスタック20と、リンク70との間に結合されるとしてよい。リンク70は、一実施形態によると、光リンク、電気リンク、またはその他の同様のリンクであってよい。図1に示すように、C I O プロトコルスタックは、C I O プロトコルトランスポート層30と、物理層の論理ブロック40と、物理層の電気ブロック50と、物理層の光ブロック60とを有するとしてよい。このように、ブロック40から60は、物理層と通信している複数のプロトコルによって共有されて、これらの複数のプロトコルの情報をリンク70によってトンネリングする共通物理層として機能する。

## 【 0 0 1 6 】

図2は、共有物理層に結合されている複数の通信スタックを備えるシステムを示す図である。具体的には、図2では、P C I e 送信（T x）および受信（R x）スタック20aに加えて、複数のその他の送信および受信スタック20b - 20dが設けられるとしてよい。同図に示すように、一对のマルチプレクサ35aおよび35b（マルチプレクサ35と総称する）が、これらのスタックと共有物理層40 - 60との間に、結合されているとしてよい。マルチプレクサ35は、プロトコルトランスポート層制御30の制御下で、動作させられるとしてよい。図2に示すように、C I O プロトコルトランスポート（P T）

層 3 0 は、P C I e およびその他のプロトコルをトンネリングするための、マルチプレクシングメカニズム（マルチプレクサ 3 5 a および 3 5 b に基づく）および制御メカニズムを実装する。P T 層制御 3 0 は、送信機の抽象化と、送信機とは別個である受信機の操作を実装する。このような構造はこの説明の残りの部分でも用いるが、本発明の実施形態は、例えば、同時に送信機および受信機を抽象化することによって、または、単一の双方向接続を用いることによって、送信機および受信機を別々に制御するその他の種類のインターコネク트에応用され得る。

#### 【 0 0 1 7 】

インターコネクットのタイミング制御を実装する方法は、実施形態によって異なるとしてよい。例えば、一部の实装によると、動的遅延バインディング（dynamic late binding）が生じるとしてよく、その結果、そのようなインターフェースロジックが、結合されるべきトンネリングインターコネク트를動的に決定でき、トンネリングインターコネク트에対応するべくプロトコルのタイミング要件を動的に制御する。ほかの実施形態によると、設計者は、システム開発時において、1 以上のプロトコルスタックによって使用されるべきトンネリングインターコネク트를決定するとしてよく、このため、当該トンネリングインターコネクによって影響され得るリンクタイミング要件が、システム設計時において決定され得る。このため、トンネリングインターコネクとプロトコルスタックとの間で、例えばインターフェースロジックにおいて、プロトコルスタックのタイミングを制御するロジックが組み込まれ得る。プロトコルスタックのタイミングの制御は、例えば、トンネリングインターコネクによって発生したさらなる遅延に対応するべくプロトコルスタックのタイミングの概念を変更することによって、実現する。

#### 【 0 0 1 8 】

図 3 は、リンクタイミング要件を処理するための上述した方法の実装例、つまり、プロトコルスタックが共通物理層または別の物理層に動的に結合され得るように、インターフェースロジックによって実装され得る動的遅延バインディングの実装例を示す図である。具体的に説明すると、図 3 は、例えば、プロトコルスタック（所与のプロトコルの標準スタックであってよい）と、さまざまなプロトコルのパケットをトンネリングすることができる統合型インターコネクのような共通物理層との間の通信のためのプロトコルスタックのインターフェースロジックで実装され得る方法 1 0 0 を説明するためのフローチャートである。図 3 に示すように、方法 1 0 0 は、トンネリングインターコネクのタイミング遅延情報を取得することによって開始されるとしてよい（ブロック 1 1 0）。この情報を取得する方法としては、さまざまな方法が実装され得る。例えば、一実施形態によると、共有物理層が、インターフェースロジックに対して所定の遅延情報の一覧を提供するとしてよい。これに代えて、インターフェースロジックは、共通物理層との間で発生しているパケット通信を分析して、タイミング遅延情報を決定するとしてもよい。より一般的には、一部の实施形態では所定の方法でタイミング情報を取得して、その他の実装でこのタイミング情報を動的に算出するとしてよい。それぞれについていくつかバリエーションが考えられ、例えば、人間が予め決定する場合と機械が予め決定する場合とがあってよく、または、算出される場合には、一度チェックするとしてもよいし、もしくは、定期的にチェックを繰り返すとしてもよい。尚、このような情報については例としてさまざまなものが考えられ、通信の性質および関連するロジックの実体に応じて、通信の種類によって発生する遅延も異なる。

#### 【 0 0 1 9 】

図 3 に戻って、ブロック 1 2 0 に進む。ブロック 1 2 0 では、タイミング遅延情報が、第 1 のプロトコルスタックのタイミング要件に対してマッピングされるとしてよい。一例を挙げると、プロトコルスタックのタイミング要件は、上記の表 1 において記載したように、リンク層通信によって異なるとしてよい。そして、ひし形 1 3 0 に進み、マッピングに基づいて、第 1 のプロトコルスタックのタイミングの概念を変更する必要があるか否かを決定するとしてよい。すなわち、共通物理層内に存在し得るレイテンシのために、当該プロトコルスタックの所定のロジックに対応付けられている 1 以上のタイマを、例えば、



加速、減速、ディセーブル等することによって、制御することができる。このようなタイミングの概念の変更が必要ない場合には、ブロック 135 に進み、プロトコルスタックの標準タイミングを用いてデータを送信および / または受信するとしてよい。

#### 【0020】

図3を参照しつつさらに説明すると、タイミングの概念を変更すべきであると決定される場合、ブロック 140 に進み、少なくとも1つのスタックロジックのタイミングを制御して、第1のプロトコルスタックのタイミングを変更するとしてよい。上述したように、このようなタイミングの変更は、タイマを制御すること、所定の間隔をカウントするように（またはしないように）ロジックを制御すること等に基づいて実現され得る。このようなタイミングの制御が実行された後、変更後のプロトコルスタックタイミングを用いて、  
10 所望のデータを送信 / 受信するとしてよい（ブロック 150）。さらに図3に示すように、続いて、通信、つまり、所定のトランザクションが完了したか否かを判断するとしてよい（ひし形 160）。完了していれば、当該方法を終了するとしてよい。完了していない場合は、ブロック 140 および 150 を繰り返すべく戻る。図3に図示する実施形態における具体的な実装に基づいて説明したが、本発明の範囲は、これに限定されない。

#### 【0021】

例えば、その他の実装例においては、所定のプロトコルスタックは、遅延が既知である既知のトンネリングインターコネクトを介してトンネリングされるように、システム設計が固定されるとしてもよい。したがって、システム設計時には、当該トンネリングインターコネクトに関して発生する可能性のあるいずれの遅延についても対応するべく、必要に  
20 応じて、さまざまなプロトコルトランザクションのタイミングの制御を実行するロジックが実装され得る。上記の表1は、リンク層タイミング要件の例を示している。

#### 【0022】

図4は、本発明の別の実施形態に係るプロトコルスタックのインターフェースを動作させる方法を示すフローチャートである。図4に示すように、方法200は、静的設計パラメータに基づいて、必要に応じて、プロトコルスタックのタイミングの概念を変更することができるインターフェースロジックによって実装され得る。図4に示すように、方法200は、トンネリングインターコネクトに対する通信またはトンネリングインターコネクトからの通信を受信することによって開始され得る（ブロック205）。このため、当該通信は、プロトコルスタックのインターフェースロジックにおいて、出力方向または入力  
30 方向で、受信される。インターフェースロジックではさまざまな種類の通信が処理されるとしてよい。例えば、肯定的応答（ACK）等のさまざまなプロトコルパケット、電力管理、フロー制御等に用いられる制御パケット、データパケットの送信および受信であってよい。

#### 【0023】

パケットの種類に基づいて、所定の通信の種類が、変更後タイミングに影響されるか否かを、インターフェースロジックにおいて決定するとしてよい（ひし形210）。例えば、インターフェースロジックは、所定の表を含むとしてもよいし、所定の表に対応付けられているとしてもよい（不揮発性メモリに存在するとしてよい）。この表は、トランザクションの種類と、その種類の通信についてプロトコルスタックのタイミングの概念を変更  
40 すべきか否かとを特定しており、さらに、適用可能な遅延の指定、タイミングを変更するべくインターフェースロジックが適用すべき制御方策の種類を示す命令またはその他の識別子を記載する。尚、この表は複数の部分に分割されているとしてよく、それぞれの部分は、所定のスタックに対応付けられており、それぞれの部分が、あるスタック - トンネリングインターコネクト関係専用のマッピングを提供している。

#### 【0024】

さらに図4を参照しつつ説明すると、変更が必要ない場合、標準のプロトコルスタックタイミングを用いて通信を処理するとしてよく、標準のプロトコルスタックタイミングを用いてデータを送信 / 受信してよい（ブロック220）。タイミングの概念を変更すべきであると決定される場合、ブロック230に進み、少なくとも1つのスタックロジックの  
50

タイミングを制御して、タイミングを変更するとしてよい。そして、変更後のプロトコルスタックタイミングを用いて所望のデータを送信／受信するとしてよい（ブロック240）。さらに図4に示すように、通信、つまり、所定のトランザクションが完了したか否かを決定するとしてよい（ひし形260）。完了している場合、当該方法を終了するとしてよい。完了していない場合、ブロック230および240を繰り返すべく戻る。リンクタイミング要件を処理する静的制御はこのように実現され得る。

#### 【0025】

上記の図3および図4に示すように、所定の通信の種類についてタイミング制御を変更することができる一方、その他の通信の種類については、変更なしで、通常のプロトコルスタックタイミングに従って処理が進められるとしてよい。以下の説明では、リンクタイ

10

#### 【0026】

一実施形態によると、PT層制御30は、PCIEに割り当てられる送信機「スロット」を提供し得る。尚、この送信機「スロット」は、送信すべきPCIEトラフィックが存在しない場合、その他の種類のトラフィックにも利用可能である。このため、第1のプロトコルスタックに割り当てられたスロットは、第1のプロトコルスタックが何も送信しない場合には、別のプロトコルスタックによって用いられるとしてもよい。同様に、受信機においても、PCIEトラフィックが受信されるはずだが、ほかの構成要素が送信すべきPCIEトラフィックを有していないか、または、別の種類のより優先度が高いトラフィ

20

#### 【0027】

「PCIE期間」という考え方をPCIEスタックに正確に伝える上で、受信時刻および送信時刻はある程度別個に考えることができる。表1に記載した、L0呼び出しポリシーおよび「L1が終了すると・・・」という要件のような、一部の状況によると、時刻が測定されるのは、1つの観点のみからである（このような場合は、送信機の観点である）。

#### 【0028】

しかし、Ack/Nakプロトコルの場合は、受信機および送信機の観点が両方とも考慮される必要がある。物理PCIEポートに基づいて送信パイプラインを伝播する際に特定のレイテンシを仮定して決定された、トランザクション層パケット(TLP)が送信された、PCIE送信機の観点での時間は、CIO送信パイプラインの遅延が異なっている場合には、不正確である場合がある。ほかの構成要素（つまり、受信機）は、自身のPCIEスタックが共有リンク上においてPCIE時間を割り当てられた場合にのみ応答が可能である。この割り当ては、自身の（受信機の）時間の概念を調整する必要がある旨を意味するものとして、受信機には感知される。PCIEスタックは、送信パイプラインの遅延を50ナノ秒(ns)と予測しているが、CIOリンクの送信パイプライン遅延は70nsだと仮定する。この場合、差分に対処するべく、20nsだけ、送信機の時間の概念を（この遅延を認識しているかどうかに応じて決まるプロトコルの側面について）停止または調整する必要がある。このようにすることで、送信機は、受信機（共有物理層によって遅延され得る）からのACK信号について適切な長さの時間だけ待機して、不適切にエラー信号が発行されないようにする。

30

40

#### 【0029】

受信機の場合は、ほかの構成要素の送信機がPCIEに対して割り当てた時間（利用した時間ではない）を考慮しなければならない。一部のケースでは、受信機に直接通知されるが、その他のケースでは、ほかの構成要素の受信機が時間の概念をどれだけ進めるべきかを指し示すメッセージ等のトンネリングプロトコルメカニズムが、トンネリングされたプロトコルのそれぞれについて、設けられるとしてよい。例えば、PCIE送信機に2つの100nsスロットが割り当てられるが、送信すべきPCIEトラフィックが少ないた

50

めに、そのうち1つのみが送信機によって利用される場合、受信機は200nsを考慮しなければならない。このように、ほかの構成要素が送信に利用可能なスロットを利用せずに、タイミングルールを守らない場合、ルールを守らなかったことは受信機において認識可能である。これは、(割り当てられたものではなく)使用された送信スロットのみを考慮する場合には、この限りでない。

#### 【0030】

尚、所定のプロトコルについては多岐にわたる最適化が可能であるとしてよい。例えば、既知の帯域幅トラフィックは、実際に与えられているリンク抽象化には関係なく、カウンタメカニズムを用いて考慮するとしてよい。プロトコルの受信割り当ておよび送信割り当てが等しいことが保証されている場合、片方のみ(例えば、送信機)が、他方(受信機)の時間の概念が一致すべきであると理解していると考えてよい。

10

#### 【0031】

先述したように、本発明の実施形態は、CIOまたはPCI等の具体的な内容に左右されるものでは一切なく、ディスプレイ、USB、ネットワーク等のトンネリングされるその他のプロトコルに適用され得る。本発明の実施形態はさらに、その他のトンネリングプロトコル/環境、例えば、有線または無線USBインターコネクトを介したトンネリングPCIEにも適用される。

#### 【0032】

本発明の一実施形態に応じてトンネリングを実行することによって、より汎用度が高いハードウェアから成る一般的なハードウェア群によって、より多くの個別のIOアプリケーションで要件を満たすことができる。例えば、プラットフォームは12個のUSBポート、8個のPCIEポート、およびさまざまな特定用途向けのポート(例えば、ディスプレイ)を備えるとしてよい。トンネリングによって、これらのポートは、例えば16個の複合型ポートから成るポート群にまとめることができ、各ポートは、従前のポートのいずれ(1または複数)としても利用することができる。

20

#### 【0033】

本発明の実施形態は、多くの異なる種類のシステムで実装することができる。図5は、シリアルリンクであるトンネリングインターコネクトによってコントローラハブに結合された複数のデバイスを備える、本発明の一実施形態に係るシステムを示すブロック図である。システム300は、コントローラハブ315に結合されている、プロセッサ305と、システムメモリ310とを備える。プロセッサ305は、マイクロプロセッサ、ホストプロセッサ、埋め込み型プロセッサ、コプロセッサ等のプロセッサを含む任意の処理素子を含む。プロセッサ305は、フロントサイドバス(FSB)306を介して、コントローラハブ315に結合されている。一実施形態によると、FSB306は、シリアル方式のポイントツーポイント(PtP)インターコネクトである。

30

#### 【0034】

システムメモリ310は、ランダムアクセスメモリ(RAM)、不揮発性(NV)メモリ等の、システム300が備えるデバイスがアクセス可能なメモリなど、任意のメモリデバイスを含む。システムメモリ310は、メモリインターフェース316を介して、コントローラハブ315に結合される。

40

#### 【0035】

一実施形態によると、コントローラハブ315は、PCIEインターコネクト階層構造において、ルートハブまたはルートコントローラである。コントローラハブ315の例を挙げると、チップセット、メモリコントローラハブ(MCH)、ノースブリッジ、入出力コントローラハブ(ICH)、サウスブリッジ、およびルートコントローラ/ハブがある。尚、コントローラハブ315は、シリアルリンク319を介して、スイッチ/ブリッジ320に結合される。入出力モジュール317および321は、インターフェース/ポート317および321と呼ばれることもあるが、コントローラハブ315とスイッチ320との間の通信を実現するべく層状プロトコルスタックを含む/実装する。一実施形態によると、スイッチ320には複数のデバイスを結合することができる。

50

## 【 0 0 3 6 】

スイッチ 3 2 0 は、デバイス 3 2 5 からのパケット / メッセージを、上流、つまり、階層構造を上方向に、コントローラハブ 3 1 5 に対してルーティングすると共に、下流、つまり、階層構造を下方向に、コントローラハブ 3 1 5 から離れるように、デバイス 3 2 5 へとパケット / メッセージをルーティングする。I O モジュール 3 2 2 および 3 2 6 は、スイッチ 3 2 0 とデバイス 3 2 5 との間の通信を実行するべく、層状プロトコルスタックを実装する。一実施形態によると、I O モジュール 3 2 6 は、複数のプロトコルスタック、つまり、スタック 3 2 7 および 3 2 8 のパケットをトンネリングするためのトンネリング物理層であってよい。デバイス 3 2 5 は、電子システムに結合される、任意の内部または外部のデバイスまたは構成要素を含む。例えば、I O デバイス、ネットワークインターフェースコントローラ ( N I C )、拡張カード、オーディオプロセッサ、ネットワークプロセッサ、ハードドライブ、ストレージデバイス、モニタ、プリンタ、マウス、キーボード、ルータ、ポータブルストレージデバイス、ファイヤワイヤデバイス、ユニバーサルシリアルバス ( U S B ) デバイス、スキャナー、およびその他の入出力デバイスが挙げられる。

10

## 【 0 0 3 7 】

コントローラハブ 3 1 5 にはさらに、シリアルリンク 3 3 2 を介して、グラフィクスアクセラレータ 3 3 0 が接続されている。一実施形態によると、グラフィクスアクセラレータ 3 3 0 は、M C H に結合され、M C H は I C H に結合されている。スイッチ 3 2 0 は、そして I O デバイス 3 2 5 は、I C H に結合されている。I O モジュール 3 3 1 および 3 1 8 もまた、グラフィクスアクセラレータ 3 3 0 とコントローラハブ 3 1 5 との間において通信を実行するべく、層状プロトコルスタックを実装する。

20

## 【 0 0 3 8 】

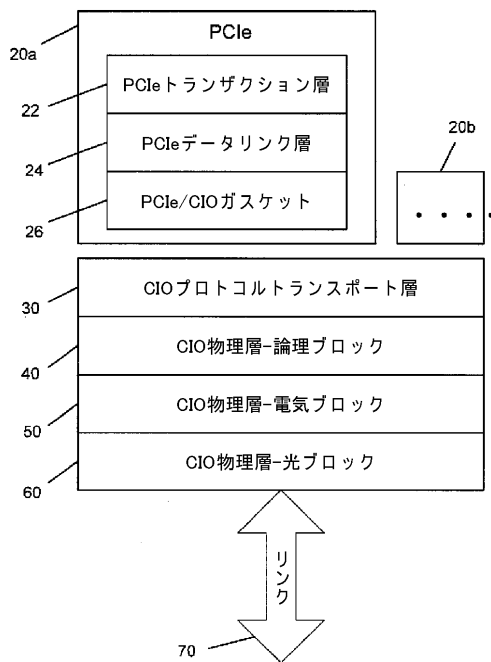
本発明の実施形態は、符号で実装されるとしてよく、命令を格納しており、当該命令を実行するようにシステムをプログラミングするべく用いられる格納媒体に格納されるとしてよい。このような格納媒体は、これらに限定されるものではないが、フロッピーディスク ( 登録商標 )、光ディスク、コンパクトディスクリードオンリーメモリ ( C D - R O M )、書き換え可能コンパクトディスク ( C D - R W )、および光磁気ディスク等の任意の種類のディスク、リードオンリーメモリ ( R O M )、ダイナミックランダムアクセスメモリ ( D R A M ) およびスタティックランダムアクセスメモリ ( S R A M ) のようなランダムアクセスメモリ ( R A M )、消去可能なプログラム可能リードオンリーメモリ ( E P R O M ( 登録商標 ) )、フラッシュメモリ、電氣的に消去可能なプログラム可能リードオンリーメモリ ( E E P R O M ( 登録商標 ) )、磁気カードあるいは光カード等の半導体デバイス、または、電子的に命令を格納するのに適しているその他の任意の種類の媒体を含むとしてよい。

30

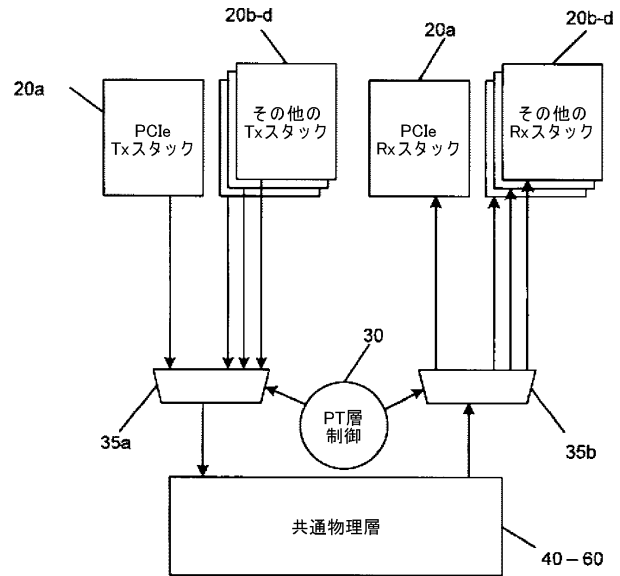
## 【 0 0 3 9 】

限られた数の実施形態に基づいて本発明を説明してきたが、当業者であれば、多くの変形および変更に想到するであろう。本願の請求項はこのような変形および変更を全て、本発明の真の精神および範囲内にあるものとして、含むものである。

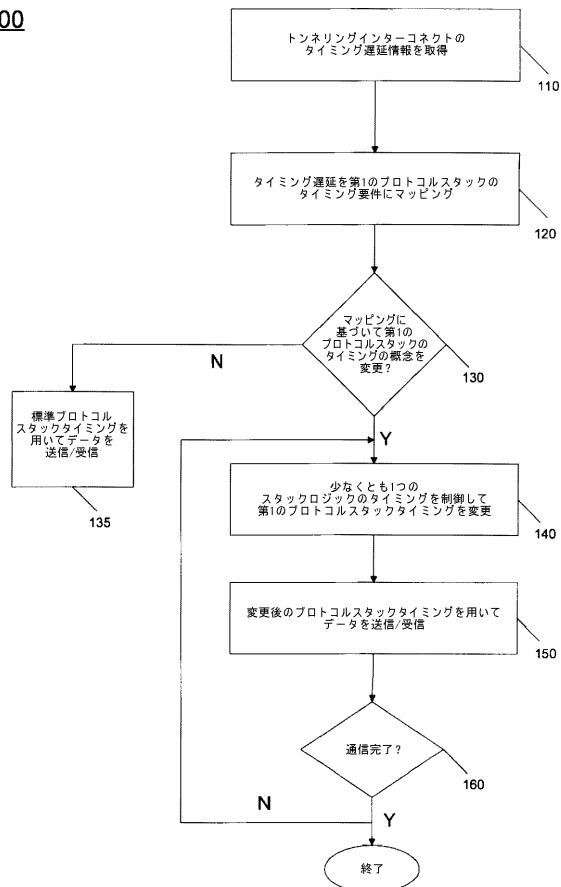
【図 1】  
10



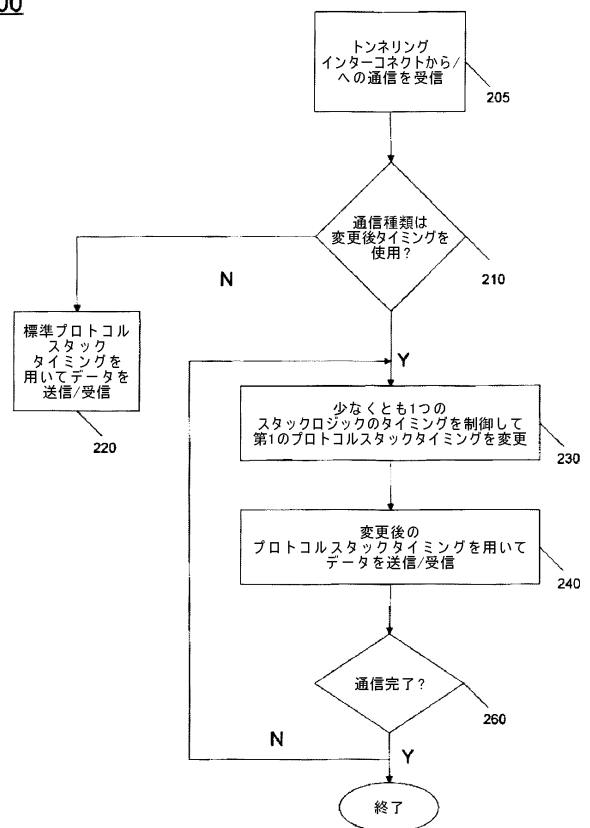
【図 2】



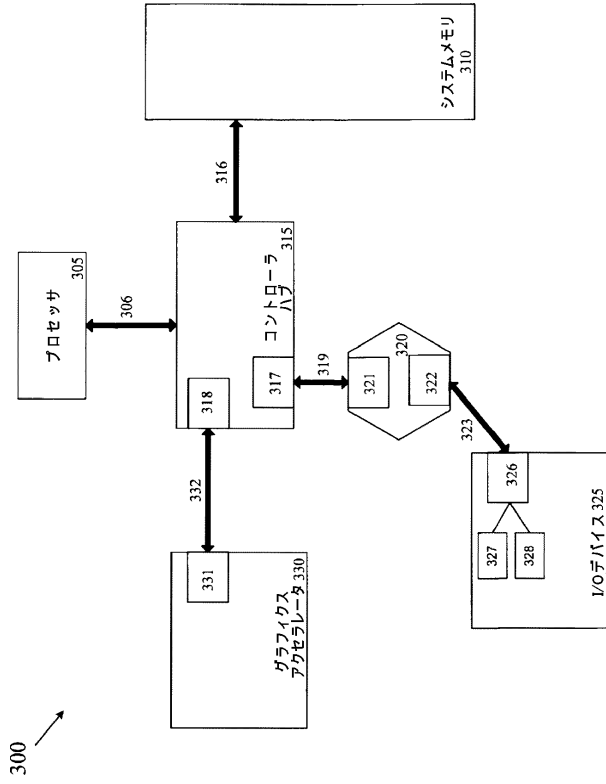
【図 3】  
100



【図 4】  
200



【図5】



---

フロントページの続き

(56)参考文献 国際公開第2005/117352(WO, A1)

特開2002-135264(JP, A)

特開2003-273953(JP, A)

特開2000-253071(JP, A)

特表2010-500807(JP, A)

(58)調査した分野(Int.Cl., DB名)

H04L 29/10

G06F 13/38

G06F 13/42

H04L 12/56