

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4391819号  
(P4391819)

(45) 発行日 平成21年12月24日(2009.12.24)

(24) 登録日 平成21年10月16日(2009.10.16)

(51) Int.Cl.

F I

G 0 6 F 13/36 (2006.01)

G 0 6 F 13/36 3 1 0 E

請求項の数 20 (全 19 頁)

(21) 出願番号 特願2003-536899 (P2003-536899)  
 (86) (22) 出願日 平成14年8月9日(2002.8.9)  
 (65) 公表番号 特表2005-505855 (P2005-505855A)  
 (43) 公表日 平成17年2月24日(2005.2.24)  
 (86) 国際出願番号 PCT/US2002/025278  
 (87) 国際公開番号 W02003/034239  
 (87) 国際公開日 平成15年4月24日(2003.4.24)  
 審査請求日 平成17年8月9日(2005.8.9)  
 (31) 優先権主張番号 09/978,349  
 (32) 優先日 平成13年10月15日(2001.10.15)  
 (33) 優先権主張国 米国 (US)

前置審査

(73) 特許権者 591016172  
 アドバンスト・マイクロ・ディバイズ・  
 インコーポレイテッド  
 ADVANCED MICRO DEVI  
 CES INCORPORATED  
 アメリカ合衆国、94088-3453  
 カリフォルニア州、サニペール、ピー・  
 オウ・ボックス・3453、ワン・エイ・  
 エム・ディ・プレイス、メイル・ストップ  
 ・68 (番地なし)  
 (74) 代理人 100099324  
 弁理士 鈴木 正剛  
 (74) 代理人 100111615  
 弁理士 佐野 良太

最終頁に続く

(54) 【発明の名称】 コンピュータ・システムの入出力ノード

(57) 【特許請求の範囲】

【請求項 1】

コンピュータシステムの入出力ノードであって、

第1ノードからの第1コマンドを第1通信経路を通じて受信するための第1受信ユニットと、

前記第1コマンドに対応する第1対応コマンドを第2通信経路を通じて第2ノードに送信するように結合された第1送信ユニットと、

前記第2ノードからの第2コマンドを第3通信経路を通じて受信するための第2受信ユニットと、

前記第2コマンドに対応する第2対応コマンドを第4通信経路を通じて前記第1ノードに送信するように結合された第2送信ユニットと、

を有し、前記第1受信ユニット、第2受信ユニット、第1送信ユニット及び第2送信ユニットのそれぞれは、コマンドを格納するための1つ以上のバッファを備え、

前記第1受信ユニット及び前記第2受信ユニットから選択されたコマンドを受信するよう結合されるとともに前記選択されたコマンドに対応するコマンドを周辺バスに送信するためのブリッジユニットと、を有し、前記ブリッジユニットは、更に、前記選択されたコマンドを対応するコマンドに変換するよう構成されたものであり、

前記第1受信ユニットと前記第1送信ユニットとの間に形成された第1通信リンクと、前記第2受信ユニットと前記第2送信ユニットとの間に形成された第2通信リンクと、を有し、前記第1通信リンク及び第2通信リンクは、パケットを伝送するための通信トンネ

10

20

ルを形成するものである、入出力ノード。

【請求項 2】

前記第 1 通信経路から前記第 2 通信経路及び前記周辺バスへのコマンド伝達と、前記第 3 通信経路から前記第 4 通信経路及び前記周辺バスへのコマンド伝達と、を制御するように結合された制御ユニットを更に有し、

前記制御ユニットは、トランザクションの種類、送信装置および宛先のバッファの空きの有無、トランザクションが転送トラフィックか注入トラフィックかであるか、のいずれかに基づいて、当該制御ユニットのバッファにあるトランザクション間の調停を行う、

請求項 1 記載の入出力ノード。

【請求項 3】

前記制御ユニットは、更に、前記周辺バスから前記第 2 通信経路及び前記第 4 通信経路へのコマンド伝達を制御するためのものである、

請求項 2 記載の入出力ノード。

【請求項 4】

前記ブリッジユニットは、更に、前記周辺バスから受信したコマンドに対応するコマンドを前記第 1 の送信ユニット及び前記第 2 送信ユニットに選択的に与える、

請求項 3 記載の入出力ノード。

【請求項 5】

前記制御ユニットは、更に、前記第 1 受信ユニット、前記第 2 受信ユニット及び前記ブリッジユニットから受信した複数の制御コマンドに基づいて前記コマンドの伝達を選択的に制御するためのものである、

請求項 4 記載の入出力ノード。

【請求項 6】

前記制御コマンドのそれぞれは、前記第 1 受信ユニット、第 2 受信ユニット及び前記ブリッジユニットによって受信した対応するコマンドの一部を含む、

請求項 5 記載の入出力ノード。

【請求項 7】

前記制御ユニットは、更に、前記制御コマンドを制御コマンド・バスを介して受信するためのものである、請求項 6 記載の入出力ノード。

【請求項 8】

前記周辺バスは、P C I バス ( P C I : Peripheral Component Interconnect ) である、

請求項 1 記載の入出力ノード。

【請求項 9】

前記周辺バスは、グラフィックス・バスである、

請求項 1 記載の入出力ノード。

【請求項 10】

前記第 1 通信経路及び第 3 通信経路及び前記第 2 通信経路及び第 4 通信経路は、HyperTransport ( 登録商標 ) リンクである、

請求項 1 記載の入出力ノード。

【請求項 11】

一つ以上のプロセッサと、

相互に接続されるとともに前記一つ以上のプロセッサの所定の一つにチェーン式に接続された 1 つ以上の入出力ノードと、を有し、各入出力ノードは当該入出力ノードに接続された他の入出力ノードと互いにパケットを送受信するようにされており、前記入出力ノードは、

前記 1 つ以上のプロセッサの所定の一つからの第 1 コマンドを第 1 通信経路を通じて受信するための第 1 受信ユニットと、

前記第 1 コマンドに対応する第 1 対応コマンドを前記 1 つ以上の入出力ノードの前記チェーンにおける次のノードへと第 2 通信経路を通じて送信するように結合された第 1 送信

10

20

30

40

50

ユニットと、

前記 1 つ以上の入出力ノードの前記チェーンにおける次のノードから第 3 通信経路を通じて第 2 コマンドを受信するための第 2 受信ユニットと、

を有し、前記第 1 受信ユニット、第 2 受信ユニット、第 1 送信ユニット及び第 2 送信ユニットのそれぞれは、コマンドを格納するための 1 つ以上のバッファを有し、

前記第 2 コマンドに対応する第 2 対応コマンドを前記 1 つ以上のプロセッサの所定の 1 つへと第 4 通信経路を通じて送信するように結合された第 2 送信ユニット、

前記第 1 受信ユニット及び前記第 2 受信ユニットから選択されたコマンドを受信するように結合されるとともに前記選択されたコマンドに対応するコマンドを周辺バスに送信するためのブリッジユニットと、を有し、前記ブリッジユニットは、更に、前記選択されたコマンドに対応するコマンドに変換するよう構成されたものであり、

10

前記第 1 受信ユニットと前記第 1 送信ユニットとの間に形成された第 1 通信リンクと、前記第 2 受信ユニットと前記第 2 送信ユニットとの間に形成された第 2 通信リンクと、を有し、前記第 1 通信リンク及び第 2 通信リンクは、パケットを伝送するための通信トンネルを形成するものである、コンピュータシステム。

【請求項 1 2】

前記入出力ノードは、前記第 1 通信経路から前記第 2 通信経路及び前記周辺バスへのコマンド伝達と、前記第 3 通信経路から前記第 4 通信経路及び前記周辺バスへのコマンド伝達と、を制御するように結合された制御ユニットを更に有し、

前記制御ユニットは、トランザクションの種類、送信装置および宛先のバッファの空きの有無、トランザクションが転送トラフィックか注入トラフィックかであるかに基づいて、当該制御ユニットのバッファにあるトランザクション間の調停を行う、

20

請求項 1 1 記載のコンピュータシステム。

【請求項 1 3】

前記制御ユニットは、更に、前記周辺バスから前記第 2 通信経路及び前記第 4 通信経路へのコマンド伝達を制御するためのものである、

請求項 1 2 記載のコンピュータシステム。

【請求項 1 4】

前記ブリッジユニットは、更に、前記周辺バスから受信したコマンドに対応するコマンドを前記第 1 の送信ユニット及び前記第 2 送信ユニットに選択的に与える、

30

請求項 1 3 記載のコンピュータシステム。

【請求項 1 5】

前記制御ユニットは、更に、前記第 1 受信ユニット、前記第 2 受信ユニット及び前記ブリッジユニットから受信した複数の制御コマンドに基づいて前記コマンドの伝達を選択的に制御するためのものである、

請求項 1 4 記載のコンピュータシステム。

【請求項 1 6】

前記制御コマンドのそれぞれは、前記第 1 受信ユニット、第 2 受信ユニット及び前記ブリッジユニットによって受信した対応するコマンドの一部を含む、

請求項 1 5 記載のコンピュータシステム。

40

【請求項 1 7】

前記制御ユニットは、更に、前記制御コマンドを制御コマンド・バスを介して受信するためのものである、

請求項 1 6 記載のコンピュータシステム。

【請求項 1 8】

前記周辺バスは、P C I バス ( P C I : Peripheral Component Interconnect ) である、

、

請求項 1 1 記載のコンピュータシステム。

【請求項 1 9】

前記周辺バスは、グラフィックス・バスである、

50

請求項 1 1 記載のコンピュータシステム。

【請求項 2 0】

前記第 1 通信経路及び第 3 通信経路及び前記第 2 通信経路及び第 4 通信経路は、HyperTransport (登録商標) リンクである、

請求項 1 1 記載のコンピュータシステム。

【発明の詳細な説明】

【技術分野】

【0 0 0 1】

本発明は、コンピュータ・システムの入出力 (I/O) に関し、より詳細には、I/O ノードにおけるトランザクション処理に関する。

【背景技術】

【0 0 0 2】

WO - A - 9 9 / 1 6 1 9 5 には、ファイバーチャネルデータリンクに接続された NL - ポート、NL - ポートとトランスミットプロトコルエンジンとの間に結合されたトランスミットフレーム FIFO バッファ、及び前記 NL - ポートと分離した受信プロトコルエンジンとの間に結合された分離受信フレーム FIFO バッファを備えたファイバーチャネル通信システムが開示されている。ホストメモリーインタフェースによって、ホストコンピュータが、トランスミットプロトコルエンジンと受信プロトコルエンジンとを備えた全二重通信プロセッサに結合される。受信プロトコルエンジンによって、上述のホストメモリーインターフェースを通じてホストメモリーへとフレームが送られる。

「ファイバーチャネル調停ループ」(1997, Solution Technology, Boulder Creek, CA, US, RW Kembel entitled "In-Depth Fibre Channel Arbitrated Loop") の 2 7 0 頁第 5 段落 ~ 2 7 5 頁第 2 段落には、デュアルループコンフィグレーションをポートバイパス回路とともに用いることが開示されている。デュアルループコンフィグレーションによって、冗長パスが各デバイスに対してループ状態で提供され、これにより、上述のループのうち一つが故障しても、他のループもしくは冗長ループを用いてデバイスのアクセスを可能としている。

典型的なコンピュータ・システムでは、1 個以上のプロセッサが、1 つ以上のバスを介して入出力 (I/O) デバイスと通信し得る。I/O デバイスは、その I/O デバイスに接続された周辺バスと、プロセッサに接続された共有バスとの間の情報の転送を管理する I/O ブリッジを介してプロセッサに結合されてもよい。さらに、I/O ブリッジは、システム・メモリと I/O デバイスとの間、あるいはシステム・メモリとプロセッサとの間の情報の転送を管理することもある。

【発明の開示】

【発明が解決しようとする課題】

【0 0 0 3】

しかし、バス・システムの多くは幾つもの難点を抱えている。例えば、バスに複数のデバイスを接続すると、バスに信号を送るデバイスの静電容量が比較的増大することがある。さらに、共有バスに複数の接続ポイントが存在すると、高い信号周波数において信号の反射が発生し、信号の整合性が失われる。このため、通常はバスの信号周波数を比較的長く保ち、信号の整合性を許容可能なレベルに維持している。信号周波数を比較的低くすると、信号帯域幅が減少し、バスに接続されたデバイスの性能が低下する。

【0 0 0 4】

共有バス・システムの別の欠点として、多数のデバイスを接続できず拡張性に欠けることがある。共有バスの利用可能な帯域幅はほぼ固定されている (また、デバイスをさらに接続すると、バスの信号周波数の低下を招く場合は、利用可能な帯域幅が減少することがある)。バスに (直接的もしくは間接的に) 接続されているデバイスが要求する帯域幅が、バスの利用可能な帯域幅を超えると、デバイスがバスにアクセスしようとしたときに、デバイスのストールが頻繁に発生し、その共有バスを含むコンピュータ・システムの全体的な性能が大幅に低下する可能性がある。I/O デバイスによって使用される共有バスの

例に、周辺機器インタフェース(peripheral component interconnect: P C I)バスがある。

【0005】

I/Oブリッジング・デバイスの多くが、何らかのバッファリング・メカニズムを使用して、P C Iバスから最終転送先バスへの転送を待っているトランザクションをバッファしている。しかし、バッファリングが、P C Iバス上でストールを発生させることもある。ストールは、一連のトランザクションがキューにバッファされて、転送先バスへの転送を待っているときに転送先バスでストールが発生してトランザクションが転送できなくなるという場合に発生する。この場合、これらの転送待ちのトランザクションを完了させるトランザクションがキューに到着し、このトランザクションは他のトランザクションの後

10

【0006】

一部のコンピュータ・システムでは、共有バスの難点の一部を解消するため、デバイス間すなわちノード間でパケット・ベースの通信を行っている。この種のシステムでは、ノードは、情報のパケットを交換することによって、相互に通信することができる。一般に、「ノード」とは、相互接続においてトランザクションに参加することが可能なデバイスを指す。例えば、この相互接続はパケット・ベースであってもよく、ノードはパケットを送受信するように構成されてもよい。一般に、「パケット」とは、パケットを送信する発信元すなわち「ソース」ノードと、このパケットを受信する宛先すなわち「ターゲット」ノードとの2つのノード間で行われる通信のことである。パケットがターゲット・ノードに到着すると、ターゲット・ノードは、パケットが伝達する情報を受け取って、その情報を内部的に処理する。ソース・ノードとターゲット・ノードとの間の通信経路に存在するノードが、ソース・ノードからターゲット・ノードにパケットを中継すなわち転送することもある。

20

【0007】

さらに、パケット・ベースの通信とバス・ベースの通信とを併用しているシステムもある。例えば、あるシステムが、P C IバスとA G Pなどのグラフィック・バスとを接続してもよい。P C Iバスはパケット・バス・インタフェースに接続されており、パケット・バス・インタフェースは、P C Iバス・トランザクションをパケット・トランザクションに変換して、パケット・バスで伝送されるようにする。同様に、グラフィックス・バスが、A G Pトランザクションをパケット・トランザクションに変換するA G Pインタフェースに接続されていてもよい。各インタフェースは、プロセッサのうちの1つに、または場合によっては別の周辺機器に接続しているホスト・ブリッジと通信し得る。

30

【0008】

P C Iデバイスがトランザクションを開始すると、パケット・ベースのトランザクションは、P C Iローカル・バス仕様に規定されているものと同じ順序付けルールによって制約され得る。これは、P C Iバス宛のパケット・トランザクションにも当てはまることがある。パケット・バス・インタフェースで発生し得るトランザクションのストールによって、パケット・バス・インタフェースでデッドロックが発生する可能性があるため、このような順序付けルールがパケット・ベースのトランザクションにも適用されている。このデッドロックが、逆にパケット・バスの側に、更なるストールを発生させる可能性がある。さらに、A G Pトランザクションが、データが正しく配達されるように、トランザクションの順序付けに関する一連のルールに従うこともある。

40

【0009】

I/Oノードの構成によっては、あるノードから別のノードへ、トランザクションが、

50

ホスト・ブリッジに向かう方向かホスト・ブリッジから出る方向に転送され得る。あるいは、特定のノードにおいてトランザクションがパケット・トラフィックに注入され得る。いずれの場合でも、トランザクションが通信経路を送信される際に、トランザクションを制御することができるI/Oノードのアーキテクチャが望ましい。

【課題を解決するための手段】

【0010】

コンピュータ・システムの入出力ノードにおけるタグ付けおよび調停メカニズムの種々の実施形態が開示される。一実施形態においては、コンピュータ・システムの入出力ノードであって第1通信経路からの第1コマンドを受信するための第1受信ユニットと、前記第1コマンドに対応する第1対応コマンドを第2通信経路に送信するように結合された第1送信ユニット(140)と、を有する。

10

【0011】

この入出力ノードは、第3通信経路から第2コマンドを受信するための第2受信ユニットと、前記第2コマンドに対応する第2対応コマンドを第4通信経路に送信するように結合された第2送信ユニットと、を有する。

【0012】

ある実施形態において、上述の通信経路は、ポイント・ツー・ポイント通信リンク、例えばハイパートランスポート(HyperTransport<sup>TM</sup>)としてもよい。更に、入出力ノードは、前記第1受信ユニット及び前記第2受信ユニットから選択されたコマンドを受信するよう結合されるとともに前記選択されたコマンドに対応するコマンドを周辺バス(152)に送信するためのブリッジユニット(150)を備えてもよい。

20

【0013】

他の実施形態において、上記入出力ノードは、前記第1通信経路から前記第1通信経路から前記第2通信経路及び前記周辺バスへのコマンド伝達と、前記第3通信経路から前記第4通信経路及び前記周辺バスへのコマンド伝達と、を制御するように結合された制御ユニットを更に有する。

【0014】

また、この制御ユニットは、前記周辺バスから前記第2通信経路及び前記第4通信経路へのコマンド伝達を制御するためのものとしてもよい。更に、この制御ユニットは、前記第1受信ユニット、前記第2受信ユニット及び前記ブリッジユニットから受信した複数の制御コマンドに基づいて前記コマンドの伝達を選択的に制御するためのものとしてもよい。

30

【発明を実施するための最良の形態】

【0015】

本発明は、様々に変形されたり他の形態とすることができるが、そのうちの特定の形態を、例として図中に図示され、かつ本明細書に詳細に記載する。しかし、図面および詳細な説明は、開示の実施形態に本発明を限定することを意図するものではなく、本発明が添付の特許請求の範囲によって規定される本発明の趣旨ならびに範囲に含まれる全ての変形例、均等物および代替例を含むことが意図にあることが理解されるべきである。

【0016】

40

図1を参照すると、コンピュータ・システムの一実施形態のブロック図が示される。このコンピュータ・システムは、各々がコヒーレント・パケット・バス15に相互接続されたプロセッサ10A~10Dを備える。コヒーレント・パケット・バス15の各々の部分は、各プロセッサ10A~Dとの間でポイント・ツー・ポイント・リンクを形成し得る。4個のプロセッサがポイント・ツー・ポイント・リンクを使用しているように図示されているが、プロセッサは4個でなくてもよいほか、他のタイプのバスをプロセッサ間の相互接続に使用してもよいという点が留意される。また、このコンピュータ・システムは、符号20, 30, 40によって示される3つのI/Oノードも備え、3つのI/Oノードは各々I/Oパケット・バス50B, 50Cによって鎖状(チェーン方式)に接続されている。I/Oパケット・バス50Aは、ホスト・ノード/プロセッサ10AとI/Oノード

50

20との間に結合されている。プロセッサ10Aは、I/Oパケット・バス50Aと通信を行うためのホスト・ブリッジを有することができるホスト・ノードとして示されている。プロセッサ10B~Dが、他のI/Oパケット・バス(図示なし)と通信を行うホスト・ブリッジを備えていてもよい。I/Oパケット・バス50A~Cによって形成される通信リンクは、ポイント・ツー・ポイント・リンクと呼ばれることもある。I/Oノード20は、一対の周辺バス25A, 25Bと接続されている。I/Oノード30はグラフィック・バス35と接続されており、I/Oノード40は別の周辺バス45と接続されている。

#### 【0017】

プロセッサ10A~10Dの各々の例に、Athlon(登録商標)マイクロプロセッサなどのx86マイクロプロセッサがある。さらに、I/Oパケット・バス50A~50Cなどのパケット・バスの一例として、非コヒーレントなHyperTransport(登録商標)がある。周辺バス25A, 25Bおよび周辺バス45の例として、周辺機器相互接続(PCI)バスなどの一般的な周辺バスがある。グラフィックス・バス35の例に、例えばアクセラレイティッド・グラフィックス・ポート(AGP)がある。しかし、上記以外のマイクロプロセッサおよび周辺バスを使用してもよいということが理解される。

#### 【0018】

3つのI/Oノードがホスト・プロセッサ10Aに接続されているが、ノードの数は3でなくてもよいほか、別のトポロジによってノード同士が接続されていてもよいということが留意される。図1には、理解しやすいようにチェーン方式のトポロジが示されている。

#### 【0019】

図中の実施形態においては、プロセッサ10Aのホスト・ブリッジが、I/Oノード20, 30または40などの下りノードから、上りパケット・トランザクションを受信し得る。あるいは、プロセッサ10Aのホスト・ブリッジが、例えば周辺バス25Aに接続され得る周辺機器(図示なし)などのデバイスに向けて、下り方向にパケットを送信することもある。

#### 【0020】

動作時に、I/Oノード20, 40は、PCIバス・トランザクションを、I/Oストリーム内を伝わる上りパケット・トランザクションに変換し得るほか、下りパケット・トランザクションをPCIバス・トランザクションに変換し得る。プロセッサ10Aのホスト・ブリッジ以外のノードから送信されるパケットは、全てプロセッサ10Aのホスト・ブリッジに向かって上り方向に向かった後に、他のノードに転送され得る。プロセッサ10Aのホスト・ブリッジから送信されるパケットは、全てI/Oノード20, 30または40などの他のノードに下り方向に流れる。本明細書中において、「上り」とは、プロセッサ10Aのホスト・ブリッジに向かうパケット・トラフィックの流れを指し、「下り」とは、プロセッサ10Aのホスト・ブリッジから離れる方向に向かうパケット・トラフィックの流れを指す。各I/Oストリームは、ユニットIDという識別子によって識別され得る。ユニットIDはパケット・ヘッダに含まれることもあれば、1つ以上のパケットに含まれる所定数のビットのことでもあるということが考察される。本明細書中において、「I/Oストリーム」とは、同じユニットIDを有し、このため同じノードから送信された全てのパケット・トランザクションを指す。

#### 【0021】

説明のため、周辺バス45にある周辺機器が、周辺バス25の周辺機器宛のトランザクションを開始したとする。トランザクションはまず、一意的なユニットIDを有する1つ以上のパケットに変換されて、次に上り方向に送られ得る。ここで、パケットの各々には、そのパケットを識別する固有の情報が符号化されている。例えば、ユニットIDが、パケット・ヘッダ内に符号化されていてもよい。さらに、トランザクションの種類も、パケット・ヘッダ内に符号化されていてもよい。各パケットは、発信元のノードを識別するユニットIDを割り当てられ得る。I/Oノード20が、周辺バス25にある周辺機器に下

10

20

30

40

50

りからパケットを転送しない可能性があるため、パケットが、プロセッサ10Aのホスト・ブリッジに向けて上り方向に送信され得る。プロセッサ10Aのホスト・ブリッジは、プロセッサ10Aのホスト・ブリッジのユニットIDを付けてこのパケットを下り方向に送信し、I/Oノード20は、このパケットを認識して、周辺バス25にある周辺機器のためにこれを要求する。次に、I/Oノード20は、パケットを周辺バス・トランザクションに変換して、このトランザクションを周辺バス25にある周辺機器に送信し得る。

#### 【0022】

パケット・トランザクションが上り方向または下り方向に伝わる際に、パケットは、1つ以上のI/Oノードを通過し得る。時として、このような通過点はトンネルと呼ばれることがあり、通過するI/Oノードはトンネル・デバイスと呼ばれることがある。上り方向から下り方向に、あるいは下り方向から上り方向に送信されるパケットは「転送(forwarded)」トラフィックと呼ばれることがある。さらに、特定のI/Oノードから発せられ、上りトラフィックに挿入されるパケット・トラフィックは、「注入(injected)」トラフィックと呼ばれることがある。

#### 【0023】

下記に詳述するように、I/Oノードは、I/Oノードに接続され得る種々のバスの順序付けルールを守るために、パケットのバッファリングのほかにトランザクションのリオーダーを行ってもよい。I/Oノードは、転送トラフィックと注入トラフィックとの両方による、トンネルに入るパケットおよびトンネルから出るパケットの流れを制御する制御ロジックも備え得る。

#### 【0024】

図2を参照すると、I/Oノードの一実施形態のブロック図が示される。このI/Oノードは、図1のI/Oノード20、30または40の代表例であり、以下、説明を簡単にするためにI/Oノード20と呼ぶ。図2のI/Oノード20は、コマンド・バス111を介して送信装置140に結合され、コマンド・バス112を介して周辺インタフェース150に結合されたトランザクション受信装置110を備え得る。また、I/Oノード20は、コマンド・バス121を介して送信装置130に結合され、コマンド・バス122を介して周辺インタフェース150に結合されたトランザクション受信装置120も備え得る。周辺インタフェース150は、コマンド・バス151を介して送信装置130、140に結合されているほか、周辺バス152にも結合されている。さらに、I/Oノード20は、制御コマンド・バス101を介して各受信装置および各送信装置に結合されているほか、周辺インタフェースに結合されたトランザクション制御ユニット100も備える。本明細書中において、コマンド・バスには、コマンド、制御およびデータ用の信号が含まれる。このため、それぞれのコマンド・バスを介してトランザクションまたはコマンドが送信されるといった場合、コマンド・ビットおよびデータ・ビットが含まれるということを目指す。

#### 【0025】

図中の実施形態においては、受信装置110と送信装置140とは、I/Oトンネルの一通信経路を形成しており、受信装置120と送信装置130とは、I/Oトンネルの別の通信経路を形成している。この通信経路はそれぞれ一方方向的であるため、いずれかの経路が上り経路か下り経路となる。このように、周辺インタフェース150からの注入トラフィックが、送信装置130、140のいずれかに提供される。

#### 【0026】

受信装置110、120の各々は、パケット・トランザクションを受け取って、受信バッファに入れる(図示なし)。各トランザクションが受信されると、受信したコマンドに含まれる情報の一部を含む制御コマンドが生成される。制御コマンドには、例えば、発信元ノードのユニットID、宛先の情報、データ・カウントおよびトランザクションの種類が含まれ得る。なお、制御コマンドは、他の情報を含んでいることもあれば、ここに挙げた情報の幾つかを含まないこともある。制御コマンドはトランザクション制御ユニット100に送信される。

10

20

30

40

50



## 【 0 0 2 7 】

周辺インタフェース 1 5 0 が周辺バス 1 5 2 からトランザクションを受信すると、周辺インタフェース 1 5 0 も、上記した制御コマンドと類似した情報を含む制御コマンドを生成し得る。また、周辺インタフェース 1 5 0 は、1 つ以上のバッファ内にトランザクションを記憶して、トランザクション制御ユニット 1 0 0 に制御コマンドを送信し得る。

## 【 0 0 2 8 】

トランザクション制御ユニット 1 0 0 は、受け取った各制御コマンドを、受信順に 1 つ以上のバッファ構造に記憶し得る。トランザクション制御ユニット 1 0 0 は、トランザクション制御ユニット 1 0 0 が自身のバッファに記憶している制御コマンドに基づいて、ソース・バッファ（すなわち受信装置および周辺インタフェースの少なくともいずれか）内で待機中の、対応するコマンドを送信する順序を決定し得る。図 3 ~ 図 6 に関して下記に詳述するように、トランザクション制御ユニット 1 0 0 は、トランザクションの種類、送信装置および宛先のバッファの空きの有無、トランザクションが転送トラフィックが注入トラフィックかであるかなどに基づいて、自身のバッファにあるトランザクション間を調停し得る。このように、トランザクション制御ユニット 1 0 0 は、I / O ノードのトンネルを通過するトランザクション全体の流れを調整する機能を果たすこともできる。

## 【 0 0 2 9 】

トランザクション制御ユニット 1 0 0 が処理すべきトランザクション間を調停したら、トランザクション制御ユニット 1 0 0 は、各々のソース・デバイスに対して、転送を待っているトランザクションを宛先デバイスに送信するよう指示する。例えば、トランザクション制御ユニット 1 0 0 は、自身のバッファから、受信装置 1 1 0 から送信装置 1 4 0 に転送しようとしているトランザクションを表す制御コマンドを選択する。トランザクション制御ユニット 1 0 0 は、コマンド・バス 1 1 1 を介して送信装置 1 4 0 にトランザクションを送信するように、受信装置 1 1 0 に対して通知する。次に、送信装置 1 4 0 は、トランザクションをチェーンの次のノードに送信する。次のノードは、上りまたは下りにある別の I / O ノードであることもあれば、図 1 のホスト・プロセッサ 1 0 A などのホスト・ノードのこともある。さらに、トランザクション制御ユニット 1 0 0 と送信装置 1 4 0 とは、受信バッファに空きが存在するかどうかをもう 1 つのノードに知らせるロジック（図示なし）を備え得る。

## 【 0 0 3 0 】

図 3 を参照すると、トランザクション制御ユニットの一実施形態のブロック図が示される。簡潔を期すと共に説明を簡単にするために、図 2 と対応する回路部品には同じ符号が付されている。トランザクション制御ユニット 1 0 0 は、各々 1 6 0 , 1 7 0 , 1 8 0 の符号が付された 3 つのスケジューラを有する。スケジューラ 1 6 0 , 1 7 0 , 1 8 0 の各々は、一対の仮想チャネル・コマンド・バッファと調停およびバッファ管理ユニットとを有する。スケジューラ 1 6 0 の仮想チャネル・コマンド・バッファは、V . C . F I F O 1 6 6 , 1 6 7 と示されており、調停およびバッファ管理ユニットは符号 1 6 8 によって示される。同様に、スケジューラ 1 7 0 の仮想チャネル・コマンド・バッファは、V . C . F I F O 1 7 6 , 1 7 7 と示されており、調停およびバッファ管理ユニットは符号 1 7 8 によって示され、スケジューラ 1 8 0 の仮想チャネル・コマンド・バッファは V . C . F I F O 1 8 6 , 1 8 7 として示されており、調停およびバッファ管理ユニットは符号 1 8 8 によって示される。

## 【 0 0 3 1 】

一般に、「仮想チャネル」とは、様々な処理ノード間でパケットを伝達する通信経路のことである。各仮想チャネルは、他の仮想チャネルの資源に依存しない（すなわち、ある仮想チャネル内で流れるパケットは、通常は物理的な伝送の点で、別の仮想チャネルの有無による影響を受けない）。パケットは、パケット・タイプに応じて仮想チャネルに割り当てられる。同じ仮想チャネル内のパケットは、相互の伝送が物理的に競合することがある（すなわち、同じ仮想チャネル内のパケット間での資源の競合が発生することがある）が、基本的には、別の仮想チャネルのパケットとの伝送は物理的に競合しない。

## 【 0 0 3 2 】

パケットによっては、他のパケットと論理的に競合する場合もある（すなわち、プロトコル、一貫性などの理由により、あるパケットが別のパケットと論理的に競合することがある）。論理的な理由やプロトコル上の理由により、1 番目のパケットが、2 番目のパケットよりも先に宛先ノードに到着しなければならない場合に、2 番目のパケットが1 番目のパケットの伝送を（競合する資源を占有することによって）物理的にブロックしてしまうと、コンピュータ・システムがデッドロックに陥る可能性がある。1 番目のパケットと2 番目のパケットとに別の仮想チャネルを割り当てると共に、異なる仮想チャネル内のパケットが相互の伝送をブロックできないような伝送媒体をコンピュータ・システムに実装することによって、デッドロックの発生しない運用を実現することができる。特に、異なる仮想チャネルのパケットが同じ物理リンクを介して伝送されるという点に留意されたい。しかし、伝送前に受信バッファを利用することが可能であるため、仮想チャネル同士がこの共有資源を使用している場合であっても、仮想チャネルが相互にブロックし合うことはない。

10

## 【 0 0 3 3 】

一観点から見ると、種々のパケット・タイプの各々（例えば、種々のコマンド符号化の各々）が独自の仮想チャネルを割り当てられてもよく、このため、一実施形態においては、各仮想チャネルに別個のバッファが割り当てられる。仮想チャネル毎に別のバッファが使用され得るため、ある仮想チャネルからのパケットが、別の仮想チャネルからのパケットと物理的に競合し得ない（これは、この種のパケットが、他のバッファに入れられるためである）。

20

## 【 0 0 3 4 】

各スケジューラは、特定の宛先と2つのソースとに対応している。図中の実施形態においては、スケジューラ160は、図2の送信装置130を宛先に、受信装置120および周辺インタフェース/ブリッジ150をソースに有するトランザクションを制御する。同様に、図3のスケジューラ170は、図2の送信装置140が宛先であり、受信装置110およびブリッジ150が送信元であるトランザクションを制御する。最後に、図3のスケジューラ180は、図2のブリッジ150が宛先であり、受信装置110および受信装置120が送信元であるトランザクションを制御する。図3において、各仮想チャネル・コマンド・バッファは、それぞれの受信装置またはブリッジから、その受信装置またはブリッジが受け取ったトランザクションに対応する制御コマンドを受け取る。制御コマンドには、どのスケジューラに制御コマンドを送信するかを指定する宛先ビットが含まれていてもよい。通常、制御コマンドには、宛先ビットが1つだけ設定されている。しかし、トランザクションがブロードキャスト・メッセージの場合には、複数のスケジューラが制御コマンドを受け取ることができるように、複数の宛先ビットが設定されていてもよい。

30

## 【 0 0 3 5 】

議論を簡略化するために、スケジューラ160についてのみ詳述する。制御コマンドが受信されてV・C・FIFO166または167に格納されると、トランザクションの種類に応じてそれぞれのFIFO部に格納される。V・C・FIFO166と167とは同一であるため、V・C・FIFO166のみについて詳述する。V・C・FIFO166は、ポスト済み、未ポストおよび応答の3種類のトランザクションに対応する3つのFIFO部を有する。制御コマンドは、受信順にそれぞれのFIFOに格納される。しかし、トランザクションの種類によっては、元のコマンドを生成した種々のバスまたはデバイスのタイプに関連する順序付けルールを保持するために、トランザクションが受信順に処理されないこともある。

40

## 【 0 0 3 6 】

図4～図6に関して下記に詳述するように、調停およびバッファ管理ロジック168は、1 番目、2 番目、それ以降に処理されるトランザクションを、V・C・FIFO166またはV・C・FIFO167のトランザクション間で調停するように構成され得る。例えば、応答コマンドより前にV・C・FIFO166に到着したポスト済み命令は、順序

50

付けルールにより、この応答コマンドより後に処理する必要がある。さらに、調停およびバッファ管理ロジック 168 は、公平ルールの組と次の I/O ノードまたはホスト・ブリッジの受信バッファの空きの有無とに基づいて、どの V.C.FIFO のトランザクションを処理すべきか調停し得る。宛先が図 2 のブリッジ 150 である場合、その調停ルールは、上記の調停ルールとは異なっているもよい。

なお、上述の例は、トランザクション制御ユニットの一実装例を示す物である。追加機能を実行することができる他の実装例を用いることも勿論可能である。

#### 【0037】

図 4 を参照すると、スケジューラの一実施形態のブロック図が示される。簡潔を期すと共に説明を簡単にするために、図 3 の回路部品と対応するものには同一の符号を付している。トランザクション・スケジューラ 400 は、調停および公平ロジック 450 に結合された仮想チャネル FIFO バッファ 410 を備える。トランザクション・スケジューラ 400 は、これも調停および公平ロジック 450 に結合された仮想チャネル FIFO バッファ 420 も備える。調停および公平ロジック 450 は FIFO バッファ 460 に接続されており、FIFO バッファ 460 はバッファ管理ロジック 470 に接続されている。バッファ管理ロジックの出力は、出力レジスタ 480 によってラッチされる。

#### 【0038】

図 3 に関して上記したように、仮想チャネル FIFO バッファ 410, 420 の各々は、例えば、図 2 の受信装置 110 またはブリッジ 150 などのそれぞれの送信元から制御コマンドを受信し得る。制御コマンドは、その制御コマンドが表すトランザクションの種類によって、仮想チャネル FIFO バッファ 410, 420 に格納され得る。具体的に言えば、制御コマンドは、ポスト済みコマンド、未ポスト・コマンド、応答コマンドのいずれかを表し得、このためポスト済み、未ポストまたは応答の各キューを表し得る。

#### 【0039】

図中の実施形態においては、調停および公平ロジック 450 は、調停ユニット 430, 440 および公平ユニット 445 を備える。調停ユニット 430 は、仮想チャネル FIFO バッファ 410 に記憶されている制御コマンドを 1 つ選択するように構成され得る。下記に詳述するように、この選択処理には、所定の調停アルゴリズムによって 1 つを「ウィナー（勝者）」として選択することが含まれ得る。同様に、調停ユニット 440 は、調停ユニット 430 と同様のアルゴリズムを使用して、仮想チャネル FIFO バッファ 420 に記憶されている制御コマンドを 1 つ選択するように構成され得る。次に、公平ユニット 445 は、調停ユニット 430, 440 によってウィナーとして選択されたトランザクションのうちの 1 つを選択し得る。公平ユニット 445 は、トランザクションが転送トランザクションか、注入トランザクションかに基づいて、公平アルゴリズムを使用し得る。また、調停ユニット 430, 440 は、スケジューラの出力先に応じて、次の I/O ノードあるいは図 1 に示したプロセッサ 10A のホスト・ブリッジの受信バッファなど、対応するトランザクションの宛先のバッファにあるバッファ空間を監視するロジック（図示なし）を有しているもよい。

#### 【0040】

図中の実施形態においては、制御コマンドがスケジューラ 400 内を伝わるのに 3 クロック・サイクルの遅延が生じ得る。仮想チャネル FIFO バッファ 410 および 420 をそれぞれ迂回するためのバイパス 415 およびバイパス 425 が図示されている。スケジューラ 400 がソースから制御コマンドを受け取り、ある仮想チャネル FIFO バッファの各キューが空である場合、仮想チャネル FIFO バッファをバイパスすることによって、1 クロック・サイクル節約することが可能となる。例えば、未ポスト制御コマンドが、現在空である仮想チャネル FIFO バッファ 410 に格納される場合を考える。トランザクションの宛先のバッファに使用可能なバッファ空間が存在することが調停ユニット 430 によって示される場合、調停ユニット 430 にあるロジックによって、未ポスト制御コマンドが仮想チャネル FIFO バッファ 410 を迂回して、直接 FIFO バッファ 460 に格納され得る。さらに、上記したように、公平ユニット 445 は、公平アルゴリズムに

応じてバイパスを許可し得る。このように、上記の例においては、1クロック・サイクル分だけ遅延を短縮することができる。本実施形態においては3クロック・サイクルの遅延が発生していたが、遅延が3クロック・サイクルよりも短い実施形態や長い実施形態も可能であるということが考察される。さらに、バイパス415, 425によって達成される実際の遅延の短縮が、本例よりも大きい場合もあれば小さい場合もある。

#### 【0041】

図4のFIFOバッファ460は、ウィナーの制御コマンドを受信し得る。図中の実施形態においては、FIFOバッファ460は2階層のバッファであるが、他の実施形態においては、FIFOバッファ460の階層数が多い場合もあれば少ない場合もあるということが考察される。

10

#### 【0042】

バッファ管理ロジック470は、図2の送信装置130または140、あるいはブリッジ150にあるバッファ空間を監視するように構成される。FIFOバッファ460にトランザクションが記憶されると、バッファ管理ロジック470は、次のバッファの空きの有無を確認して、バッファ空間に空きが生じるまで制御コマンドを保持するか、あるいは出力レジスタ480に出力する。制御コマンドが出力レジスタ480によってラッチされると、それぞれのトランザクションの送信元は、図2の送信装置130または140、あるいはブリッジ150に制御コマンドに対応するトランザクションを送信可能であることを通知される。

#### 【0043】

20

図5を参照すると、タグ付けロジックを備えたスケジューラの一実施形態のブロック図が示される。スケジューラ500は、仮想チャネルFIFOバッファ505に結合されたタグ付けロジック510を備える。仮想チャネルFIFOバッファ505は、ポスト済み、未ポストおよび応答の3種類のトランザクションに対応する3つのキューを有する。タグ比較/調停ロジック・ユニット520は、仮想チャネルFIFOバッファ505に結合されている。このほか、仮想チャネルFIFOバッファ505の分解図も示されている。この分解図には、未ポスト・キューと応答キューの各々は、対応するタグを1つ有することが示される。しかし、ポスト済みキューは、未ポストキューと応答キューとに対応する2つのタグを有する。

#### 【0044】

30

タグ付けロジック510は、各制御コマンドが仮想チャネルFIFOバッファ505に記憶される前に、その制御コマンドにタグを割り当て得る。制御コマンドには、図2の受信装置110またはブリッジ150などのソース・ユニットから制御コマンドを受信した順に、タグが割り当てられてもよい。タグは、制御コマンドの最後に追加されてもよい。

#### 【0045】

制御コマンドが仮想チャネルFIFOバッファ505の一番上の場所に達すると、タグ比較/調停ロジック・ユニット520は、3つの仮想チャネル間を調停して、ウィナーの制御コマンドを選択するように構成され得る。ウィナーは、順序付けルールのセットに基づいたアルゴリズムを使用して選択され、この順序付けルールのセットは、I/Oノードに接続されている周辺バスが従っている順序付けルールに対応したものであってもよい。一実施形態においては、この順序付けルールは、PCIの順序付けルールに対応したものであってもよい。別の実施形態においては、この順序付けルールは、AGPの順序付けルールに対応したものであってもよい。

40

#### 【0046】

図中の実施形態においては、仮想チャネルFIFOバッファ505には最大で16の場所を含むことができるため、タグは4ビットであり得る。しかし、他の実施形態においては、仮想チャネルFIFOバッファ505に含まれる場所の数が16以外であってもよく、このためタグのビット数が適宜変わってもよいという点が留意される。タグ付けロジック510は、仮想チャネルFIFOバッファ505に格納されている未ポスト制御コマンドおよび応答制御コマンドについて、現在のタグを監視するカウンタ・ロジック(図示な

50

し)を備え得る。このタグは、ポスト済み制御コマンドを受信したとき、現在のポスト済み制御コマンドより先で、かつ以前のポスト済み制御コマンドより後に、ブロック可能な未ポストの制御コマンドまたは応答制御コマンドを少なくとも1つ受信している場合に、未ポスト・カウンタまたは応答カウンタをそれぞれインクリメントするアルゴリズムに従って割り当てられ得る。本明細書中において、未ポスト制御コマンドまたは応答制御コマンドがブロック可能であるとは、仮想チャネルFIFOバッファ505内で、未ポスト制御コマンドまたは応答制御コマンドが、ポスト済み制御コマンドを追い越すことができることを示す特別なビットが、各制御コマンド内でセットされている場合を指す。一実施形態においては、この特別なビットをPassPWビットと呼ぶ。

#### 【0047】

タグ付け・アルゴリズムの使用を示すために、表1に、制御コマンドを受信する順序と、ポスト済み、未ポストおよび応答の3つのキューに入るタグを示す。1番目の列は、9個の制御コマンドを受信する順序を示す。2番目の列は、受信コマンドの種類を示す。3番目の列は、未ポスト・コマンドおよび応答コマンドに割り当てられたタグを、4番目および5番目の列は、ポスト済み制御コマンドによってインクリメントされ得る後の未ポスト・コマンドおよび応答コマンドのカウンタ値を示す。ポスト済み制御コマンドには2つのタグが割り当てられるため、未ポスト・カウンタおよび応答カウンタの各々の現在のカウンタ値に示される両方のタグが割り当てられる。仮想チャネルFIFOバッファ505の分解図は、表1のコマンドが格納される様子を示す。

【表1】

受信順	制御コマンド	PassPWビット	タグ値	未ポスト・カウンタ	応答カウンタ
1	ポスト済み1	0		0	0
2	応答1	0	0	0	0
3	ポスト済み2	0		0	1
4	未ポスト1	0	0	0	1
5	応答2	0	1	0	1
6	未ポスト2	0	0	0	1
7	ポスト済み3	0		1	2
8	応答3	0	2	1	2
9	未ポスト3	0	1	1	2

#### 【0048】

表2に、PassPWビットにより、表1に示した受信制御コマンドのタグがどのように変わるかを示す。未ポスト制御コマンドまたは応答制御コマンドにPassPWビットがセットされている場合、次のポスト済み制御コマンドが到着しても、それぞれのカウンタがインクリメントされない。例えば、表2において、応答1にPassPWビットがセットされているため、ポスト済み2の制御コマンドが受信されたときに、応答カウンタがインクリメントされない。しかし、未ポスト2の制御コマンドにPassPWビットがセットされており、ポスト済み3の制御コマンドによって、未ポスト・カウンタと応答カウンタとがインクリメントされる。これは、最後のポスト済み制御コマンドより後で、かつ現在のポスト済み制御コマンドより前に、PassPWビットがクリアされた状態で未ポスト1制御コマンドが受信され、このためカウンタをインクリメントするタグ付けルールを満たしているためである。PassPWビットがセットされると、未ポスト制御コマンドまたは応答制御コマンドがポスト済み制御コマンドを追い越し可能であることを示すと記載したが、このロジックが逆転している実施形態も可能であるということが考察される。

。

【表 2】

受信順	制御コマンド	P a s s PWビット	タグ値	未ポスト・ カウンタ	応答 カウンタ
1	ポスト済み1	0		0	0
2	応答1	1	0	0	0
3	ポスト済み2	0		0	0
4	未ポスト1	0	0	0	0
5	応答2	0	0	0	0
6	未ポスト2	1	0	0	0
7	ポスト済み3	0		1	1
8	応答3	0	1	1	1
9	未ポスト3	0	1	1	1

10

## 【0049】

図5に戻ると、タグ比較および調停ロジック・ユニット520は、調停中に、各制御コマンドの後に追加されたタグを比較することによって、仮想チャネルFIFOバッファ505からウィナーの制御コマンドを選択するように構成される。さらに、タグ比較および調停ロジック・ユニット520は、タグを比較する前に、各仮想チャネルについて、次のI/Oノードのバッファ空間に空きがあるかどうかを判定してもよい。ブロックされている仮想チャネルがあれば、そのチャネルは、このサイクルの調停の対象とならない。3つのチャネル全てが調停の対象となる場合は、未ポスト・チャネルと応答チャネルとの間でラウンド・ロビン方式で調停が行われ、ポスト済みチャネルと未ポスト・チャネルとの間の比較、およびポスト済みチャネルと応答チャネルとの間の比較が個別に行われる。最後に、タグ比較および調停ロジック・ユニット520は、公平アルゴリズムを使用してウィナーを決定し得る。この公平アルゴリズムについては、図6を参照して下記に詳述する。

20

## 【0050】

次に図6を参照すると、欠乏回避ロジックを備えたトランザクション・スケジューラの一実施形態のブロック図が示される。トランザクション・スケジューラ600は、調停回路650に結合された仮想チャネルFIFOバッファ610を備える。トランザクション・スケジューラ600は、これも調停回路650に結合された仮想チャネルFIFOバッファ620も備える。調停回路650はFIFOバッファ670に接続されており、FIFOバッファ670はバッファ管理ロジック680に接続されている。バッファ管理ロジックの出力は、出力レジスタ690によってラッチされる。

30

## 【0051】

上記の図3，図4と同様に、仮想チャネルFIFOバッファ610，620の各々は、例えば、図2の受信装置110またはブリッジ150などの入力元から制御コマンドをそれぞれ受信し得る。制御コマンドは、その制御コマンドが表すトランザクションの種類に従って、仮想チャネルFIFOバッファ610，620に格納され得る。例えば、制御コマンドは、ポスト済みコマンド、未ポスト・コマンド、応答コマンドのいずれかを表し得る、このためポスト済み、未ポストまたは応答の各キューを表し得る。

## 【0052】

40

図中の実施形態においては、調停回路650は、調停ユニット630，640および公平回路645を備える。調停ユニット630，640は、調停サイクルの間に、仮想チャネルFIFOバッファ610，620のそれぞれに記憶されている制御コマンドを選択するように構成され得る。さらに、公平回路645は、調停ユニット630か640かのいずれが、ウィナーのトランザクションを選択するかを決定し得る選択条件を提供してもよい。図7に関して下記に詳述するように、公平回路645は公平アルゴリズムを使用して、帯域幅を振り分ける調停の優先順位を確立し得る。公平アルゴリズムは、所定数の調停サイクルの間、あるトランザクションがブロックされていることを判定する欠乏回避ロジックを使用してもよい。

## 【0053】

50

図7を参照すると、図6に示した公平回路の一実施形態のブロック図が示される。公平回路645は、32個の3ビット・カウンタ0～31を備えた公平ユニット700を有し、3ビット・カウンタ0～31は8ビット・カウンタ705に結合されている。ラッチ710がカウンタ705に結合されている。挿入レート・ロジック715がラッチ710に結合されている。また、公平回路645は、3つの仮想チャネル・カウンタ755～757を備えた欠乏ユニット750も有し、仮想チャネル・カウンタ755～757は、欠乏閾値レジスタ760およびトランザクション選択ユニット775に結合されており、トランザクション選択ユニット775は公平ユニット700および欠乏ユニット750に結合されている。

#### 【0054】

一実施形態においては、トランザクションがトンネルを通過して転送される度に、そのトランザクションを送信したI/Oノードに対応する3ビット・カウンタ0～31がインクリメントされ得る。また、3ビット・カウンタ0～31のいずれかがインクリメントされる度に、カウンタ705もインクリメントされ得る。3ビット・カウンタ0～31のいずれかが1つがオーバーフローすると、カウンタ705の値がラッチ710によって取得される。取得した値は、ある時点での下りノードのトランザクション要求の発生レート(rate)を表し得る。挿入レート・ロジック715は、取得した値を使用して、そのノードの許容可能な挿入レートを計算することができる。

#### 【0055】

トランザクション選択ユニット775によって、図6の調停ユニット630または640のポインタが、調停サイクル中に考慮されている仮想チャネル・バッファをポイントするようになる。転送用の仮想チャネル・バッファの仮想チャネルにも、挿入用のバッファの同じ仮想チャネルにもトランザクションが存在する場合があります。図7において、トランザクション選択ユニット775は、公平ユニット700によって確立された優先順位に従って、この2つの仮想チャネル・バッファを交互に選択し得る。しかし、トランザクションの宛先のバッファに空きが存在しないなどの要因によって、ある仮想チャネルがブロックされている場合には、調停を続行している間、調停ロジックは、ブロックされているチャネルを避け、次の仮想チャネルを選択し得る。ブロックされているチャネルが利用可能になったときに、転送用のチャネルのみに空きが存在し、公平アルゴリズムによって転送用のチャネルは挿入用のチャネルよりも優先されるため、転送用のチャネルからのトランザクションが送信され得る。以前はブロックされていたチャネルが次に利用可能になると、現在の調停サイクルにおいて優先されないため、このチャネルは再度スキップされ得る。この状態はしばらく続き、このために挿入用の仮想チャネルが「欠乏」し得る。別の実施形態においては、ノードの構成によっては、同じような状況下で、注入チャネルのために転送チャネルが欠乏し得るということが考察される。

#### 【0056】

特定のチャネルが欠乏しないように、欠乏ユニット750は、トランザクションがブロックされた回数を監視し得る。トランザクション選択ユニット775が、調停の対象となっているがブロックされているトランザクションがどれであることを判定する度に、トランザクション選択ユニット775は、対応する仮想チャネル・カウンタ755～757をインクリメントする。欠乏閾値レジスタ760は、調停中のチャネルのスキップ回数の最大数に対応する値を保持している。欠乏閾値レジスタ760に記憶される値は、特定の時点における計算によって求めた要求の発生レートによって動的に変更され得る。仮想チャネル・カウンタ755～757の値のいずれかが、欠乏閾値レジスタに記憶されている値に達すると、次の調停サイクルで、ブロックされているトランザクションがトランザクション選択ユニット775によって選択されるように、対応する仮想チャネルの優先順位が変更され得る。このように、優先順位を動的に変更することによって、特定のチャネルの欠乏を回避することができる。上記のロジック構成は実装の一例に過ぎないという点が留意される。別法による実施形態においては、カウンタの数が異なるか、カウンタのビット数が異なるロジック構成を使用して、上記した機能を実現してもよいということが考察され

10

20

30

40

50

る。

【 0 0 5 7 】

また、上記の実施形態は、スケジューラの一実装であることも留意される。別法による実施形態においては、別の機能および異なる機能の少なくともいずれかを実行する別の実装を有していてもよいということが考察される。

【 0 0 5 8 】

上記の開示が完全に理解されれば、数多くの変形例および変更例が当業者にとって自明となるであろう。添付の特許請求の範囲はこのような変更例および変更例を全て包含するものと解釈されることが意図される。

【 産業上の利用可能性 】

10

【 0 0 5 9 】

本発明は、一般にコンピュータ・システムの入出力（I/O）に適用可能であり、より詳細には、I/Oノードにおけるトランザクション処理に適用可能である。

【 図面の簡単な説明 】

【 0 0 6 0 】

【 図 1 】 コンピュータ・システムの一実施形態を示すブロック図である。

【 図 2 】 I/Oノードの一実施形態の一実施形態を示すブロック図である。

【 図 3 】 トランザクション制御ユニットの一実施形態を示すブロック図である。

【 図 4 】 スケジューラの一実施形態を示すブロック図である。

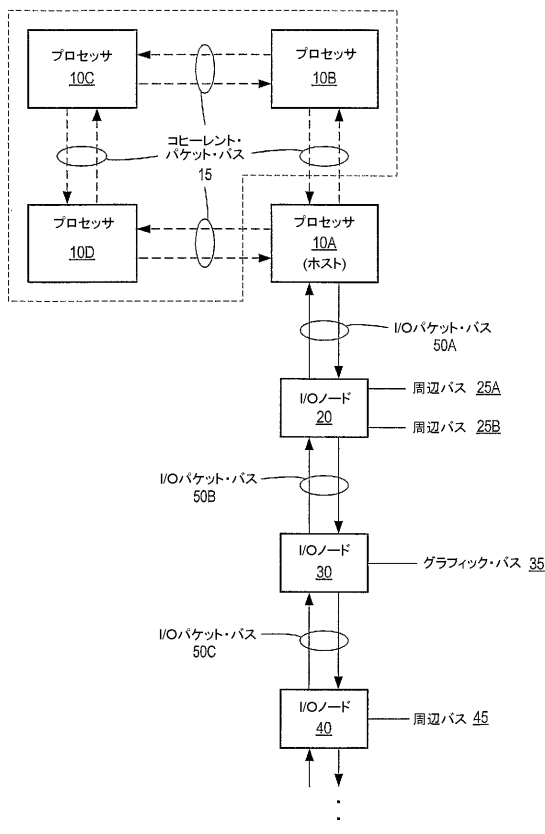
【 図 5 】 タグ付けロジックを備えたスケジューラの一実施形態を示すブロック図である。

20

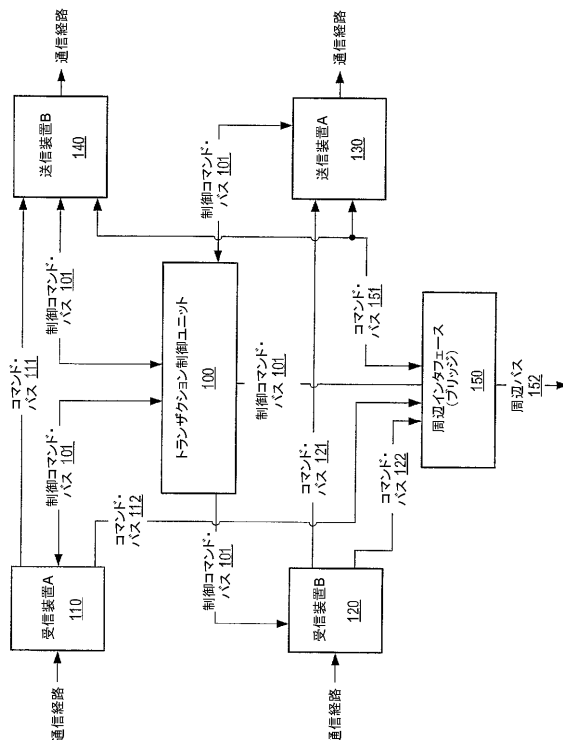
【 図 6 】 欠乏回避ロジックを備えたトランザクション・スケジューラの一実施形態を示すブロック図である。

【 図 7 】 公平回路の一実施形態を示すブロック図である。

【 図 1 】

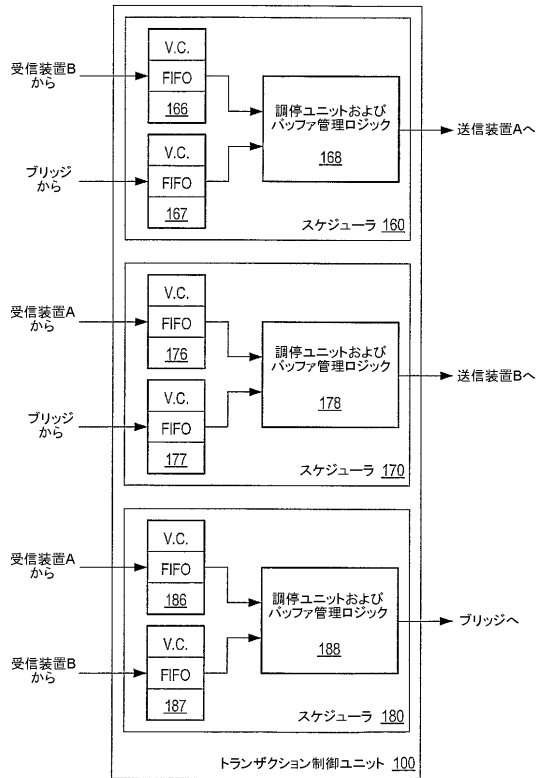


【 図 2 】

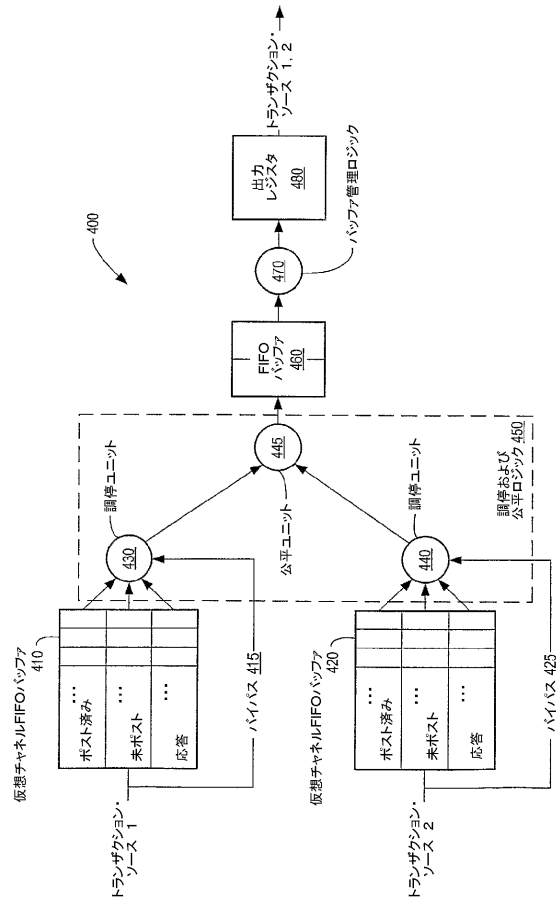




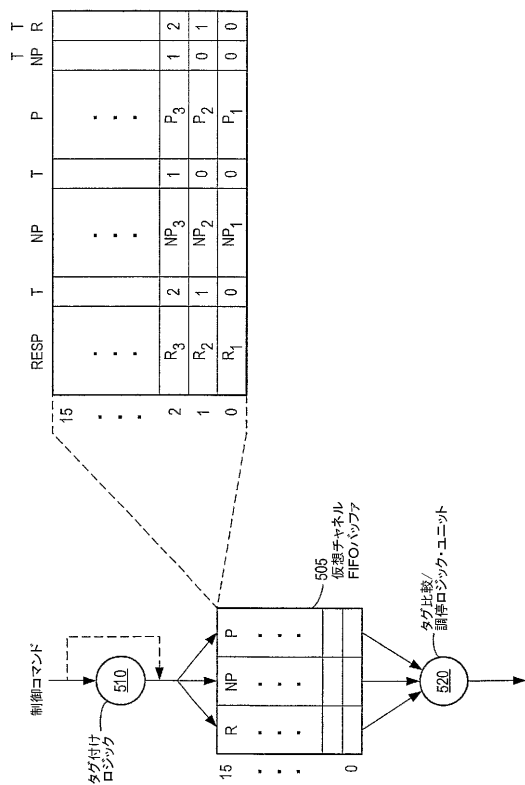
【図 3】



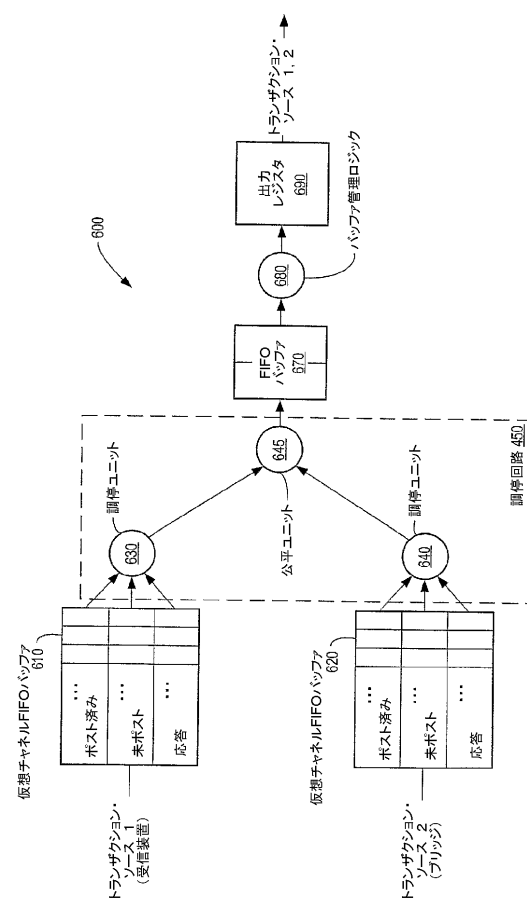
【図 4】



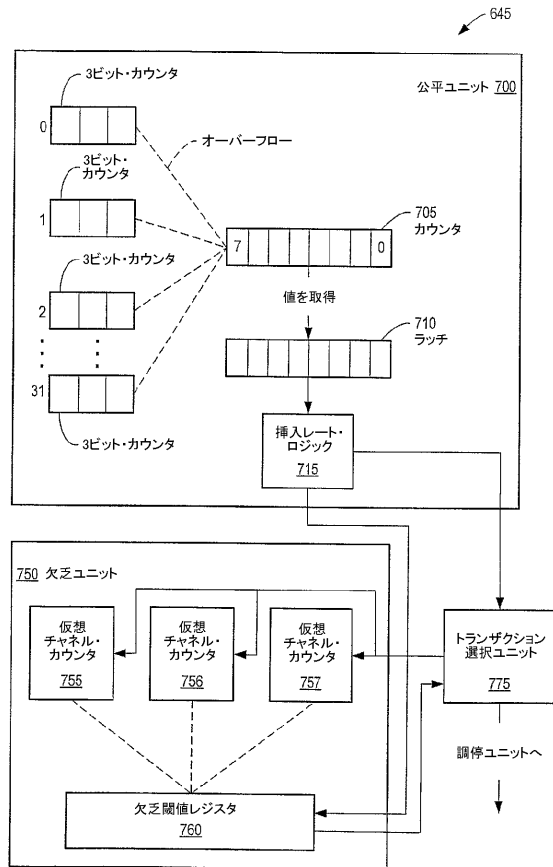
【図 5】



【図 6】



【図 7】



---

フロントページの続き

(74)代理人 100108604

弁理士 村松 義人

(72)発明者 スティーブン シー . エニス

アメリカ合衆国、テキサス州 78704、オースティン、エアロール ウェイ 2100 - ビー

(72)発明者 ラリー ディー . ヒューイット

アメリカ合衆国、テキサス州 78746、オースティン、ベンド オリバー ドライブ 6103

審査官 木村 貴俊

(56)参考文献 特開平11-143847(JP, A)

特開平08-297630(JP, A)

(58)調査した分野(Int.Cl., DB名)

G06F 13/20-13/378