

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4078446号  
(P4078446)

(45) 発行日 平成20年4月23日 (2008. 4. 23)

(24) 登録日 平成20年2月15日 (2008. 2. 15)

(51) Int. Cl.

H04L 12/56 (2006.01)

F I

H04L 12/56

F

請求項の数 21 (全 30 頁)

(21) 出願番号 特願平10-535787  
 (86) (22) 出願日 平成10年1月30日 (1998. 1. 30)  
 (65) 公表番号 特表2002-514366 (P2002-514366A)  
 (43) 公表日 平成14年5月14日 (2002. 5. 14)  
 (86) 国際出願番号 PCT/US1998/001985  
 (87) 国際公開番号 WO1998/036530  
 (87) 国際公開日 平成10年8月20日 (1998. 8. 20)  
 審査請求日 平成17年1月6日 (2005. 1. 6)  
 (31) 優先権主張番号 60/038, 025  
 (32) 優先日 平成9年2月14日 (1997. 2. 14)  
 (33) 優先権主張国 米国 (US)  
 (31) 優先権主張番号 08/993, 880  
 (32) 優先日 平成9年12月18日 (1997. 12. 18)  
 (33) 優先権主張国 米国 (US)

(73) 特許権者

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  
 アメリカ合衆国、 94088-3453  
 カリフォルニア州、 サニibel、  
 ビー・オー・ボックス・3453、 ワン  
 ・エイ・エム・ディ・プレイス、 メール  
 ・ストップ・68 (番地なし)

(74) 代理人

弁理士 深見 久郎

(74) 代理人

弁理士 森田 俊雄

(74) 代理人

弁理士 伊藤 英彦

最終頁に続く

(54) 【発明の名称】 探索可能なキャッシュ領域を備えるマルチコピーキュー構造

(57) 【特許請求の範囲】

【請求項 1】

ネットワークスイッチから送信されるフレームのコピーのカウントを維持するための構成であって、

エントリを受取ってキューに入れるよう構成されるマルチコピーキューを含み、各エントリはフレーム識別子およびコピー数を有し、各コピー数は正または負の値であり、正の値のコピー数は送信されるべきフレームのコピーの総数を表わし、負の値のコピー数は既に行なわれたフレームのコピーの1つについての送信を表わし、さらに、

正のコピー数を有するマルチコピーキューからのエントリを受取り、ストアするよう構成される探索可能な第1のメモリと、

マルチコピーキューを出るエントリを検査し、マルチコピーキューからのエントリが負のコピー数を有するときには、負のコピー数を有するエントリと同じフレーム識別子を有するエントリを求めて第1のメモリを探索し、第1のメモリ内のエントリのコピー数が1よりも大きければそのフレーム識別子を有する第1のメモリ内のエントリのコピー数をデクリメントし、コピー数が1と等しければ第1のメモリ内のそのエントリをデリートするよう構成されるバッファマネージャとを含み、

フレーム識別子は、フレームがストアされている第2のメモリ内の場所を指すフレームポインタである、構成。

【請求項 2】

第1のメモリはキャッシュメモリである、請求項1に記載の構成。

**【請求項 3】**

バッファ内にフレームをストアするよう構成される第 2 のメモリをさらに含む、請求項 1 に記載の構成。

**【請求項 4】**

第 2 のメモリ内のバッファは、そのバッファ内にストアされているフレームのためのコピー数がストア可能であるコピー数領域を備えたバッファヘッダを有する、請求項 3 に記載の構成。

**【請求項 5】**

バッファマネージャは、マルチコピーキューから出る正のコピー数を有するエントリを第 1 のメモリへと書込み、正のコピー数を有するエントリの書込の前に第 1 のメモリがフルであれば、第 1 のメモリから最古のエントリを除去し、そのエントリのコピー数を、バッファマネージャによって除去されたエントリ内のフレームポインタが指すフレームをストアするバッファのバッファヘッダのコピー数領域へと書込むようさらに構成される、請求項 4 に記載の構成。

**【請求項 6】**

バッファマネージャは、バッファマネージャによる第 1 のメモリの探索が負のコピー数を有するエントリと同じ識別子を有するエントリを第 1 のメモリ内でつきとめないとき、その負のコピー数を含むエントリ内のフレームポインタが指すフレームをストアしているバッファのコピー数領域からコピー数を取出すようさらに構成される、請求項 5 に記載の構成。

**【請求項 7】**

メモリ構成であって、  
エントリを受取ってキューに入れるよう構成されるマルチコピーキューを含み、各エントリはフレームポインタおよびコピー数を有し、各コピー数は正または負の値であり、正の値のコピー数は送信されるべきフレームのコピーの総数を表わし、負の値のコピー数は既に行なわれたフレームのコピーの 1 つについての送信を表わし、さらに、  
エントリをストアするよう構成される第 1 のメモリを含み、各エントリは第 2 のメモリ内にストアされているデータを指すデータポインタと、関連のデータ項目とを有し、さらに、

データポインタが指す第 2 のメモリ内の場所でデータをストアするよう構成される第 2 のメモリを含み、第 2 のメモリ内の場所はデータとともに関連のデータ項目をストアするよう構成され、したがって、第 2 のメモリは関連のデータ項目をストアするための第 1 のメモリの拡張部であり、

第 1 のメモリは探索可能なキャッシュであり、データはデータのフレームであり、データポインタはフレームポインタであり、関連のデータ項目は装置から送信されるべきフレームのコピーの数を示すコピー数である、メモリ構成。

**【請求項 8】**

エントリを検査し、エントリが負のコピー数を含むときに、負のコピー数を有するエントリと同じフレームポインタを有するエントリを求めてキャッシュを探索し、キャッシュ内のエントリのコピー数が 1 よりも大きければ、そのフレームポインタを有するキャッシュ内のエントリのコピー数をデクリメントし、コピー数が 1 と等しければキャッシュ内のエントリをデリートするよう構成されるバッファマネージャをさらに含む、請求項 7 に記載のメモリ構成。

**【請求項 9】**

バッファマネージャは、正のコピー数のエントリをキャッシュへと書込み、正のコピー数のエントリの書込の前にキャッシュがフルであれば、キャッシュから最古のエントリを除去し、そのエントリのコピー数を、バッファマネージャによってキャッシュから除去されたエントリ内のフレームポインタが指すフレームをストアする第 2 のメモリ内の場所のコピー数領域へと書込むようさらに構成される、請求項 8 に記載のメモリ構成。

**【請求項 10】**

バッファマネージャは、バッファマネージャによるキャッシュの探索が負のコピー数を有するエントリと同じ識別子を有するエントリをキャッシュ内につきとめないとき、負のコピー数を含むエントリ内のフレームポインタが指すフレームをストアする場所のコピー数領域からコピー数を取出すようさらに構成される、請求項 9 に記載のメモリ構成。

【請求項 1 1】

装置から送信されるべきデータ項目のコピーの数のカウントを維持するための構成であって、

エントリを受取ってキューに入れるよう構成されるマルチコピーキューを含み、各エントリはフレームポインタおよびコピー数を有し、各コピー数は正または負の値であり、正の値のコピー数は送信されるべきフレームのコピーの総数を表わし、負の値のコピー数は既

10

に行なわれたフレームのコピーの 1 つについての送信を表わし、さらに、エントリがストアされる第 1 のメモリを含み、第 1 のメモリ内にストアされる各エントリは複数のコピーが装置から送信されるべき異なるデータ項目に対応し、第 1 のメモリはキャッシュメモリであり、各エントリは、データ項目がストアされている第 2 のメモリ内の場所を指すポインタと、そのデータ項目のまだ送信されていないコピーの数を示すコピー数とを有し、さらに、

第 1 のメモリ内のエントリと同じポインタを有するデータ項目のコピーが装置から送信されるときに、第 1 のメモリ内にストアされているエントリの 1 つの中のコピー数を調節するよう構成されるマネージャ装置を含む、構成。

【請求項 1 2】

20

マネージャ装置は、装置内のエントリを受取り検査するようさらに構成され、前記装置内のエントリは、第 1 のメモリ内にストアされるべきエントリと、前記ポインタとポインタが指すデータ項目についての 1 つのコピーが装置から送信されたことを示す負のコピー数とを含むエントリとを含み、マネージャ装置はさらに、マネージャ装置が負のコピー数と第 1 のメモリ内にストアされているそのエントリのポインタに一致するポインタとを持つエントリを受取ると、第 1 のメモリ内にストアされているエントリの 1 つについてのコピー数を調節するよう構成される、請求項 1 1 に記載の構成。

【請求項 1 3】

装置はネットワークスイッチであり、マネージャ装置はマルチポートスイッチのバッファマネージャであり、第 1 のメモリはキャッシュメモリであり、データ項目はデータのフレームであり、ポインタはフレームポインタである、請求項 1 2 に記載の構成。

30

【請求項 1 4】

パケット交換網のためのマルチポートスイッチ構成であって、スイッチは単一のフレームの複数のコピーを送信するよう構成され、スイッチ構成は、

エントリを受取ってキューに入れるよう構成されるマルチコピーキューを含み、各エントリはフレームポインタおよびコピー数を有し、各コピー数は正または負の値であり、正の値のコピー数は送信されるべきフレームのコピーの総数を表わし、負の値のコピー数は既に行なわれたフレームのコピーの 1 つについての送信を表わし、さらに、

正のコピー数を有するマルチコピーキューからのエントリを受取ってストアするよう構成される探索可能な第 1 のメモリと、

40

マルチコピーキューを出るエントリを検査し、マルチコピーキューからのエントリが負のコピー数を含むときに、負のコピー数を有するエントリと同じフレームポインタを有するエントリを求めて第 1 のメモリを探索し、第 1 のメモリ内のエントリのコピー数が 1 よりも大きければ、そのフレームポインタを有する第 1 のメモリ内のエントリのコピー数をデクリメントし、コピー数が 1 と等しければ第 1 のメモリ内のそのエントリをデリートするよう構成されるバッファマネージャとを含み、

フレームポインタはフレームがストアされる第 2 のメモリ内の場所を指す、スイッチ構成

【請求項 1 5】

バッファ内のフレームをストアするよう構成される第 2 のメモリをさらに含む、請求項 1

50

4に記載のスイッチ構成。

【請求項16】

第2のメモリ内のバッファは、そのバッファ内にストアされているフレームのためのコピー数がストア可能であるコピー数領域を含むバッファヘッダを有する、請求項15に記載のスイッチ構成。

【請求項17】

バッファマネージャは、マルチコピーキューから出る正のコピー数のエントリを第1のメモリへと書込み、正のコピー数のエントリの書込の前に第1のメモリがフルであれば、第1のメモリから最古のエントリを除去し、そのエントリのコピー数を、バッファマネージャによって除去されたエントリ内のフレームポインタが指すフレームをストアするバッファのバッファヘッダのコピー数領域へと書込むようさらに構成される、請求項16に記載のスイッチ構成。

【請求項18】

バッファマネージャは、バッファマネージャによる第1のメモリの探索が負のコピー数を有するエントリと同じ識別子を有するエントリを第1のメモリ内につきとめないとき、負のコピー数を含むエントリ内のフレームポインタが指すフレームをストアするバッファのコピー数領域からコピー数を取出すようさらに構成される、請求項17に記載のスイッチ構成。

【請求項19】

ネットワークスイッチからのフレームの送信の数のカウントを維持する方法であって、探索可能な第1のメモリにエントリをロードするステップを含み、各エントリは、フレーム識別子と、ネットワークスイッチからまだ送信されていないフレームのコピーの数を示すコピー数とを含み、第1のメモリはキャッシュメモリであり、さらに、

第2のメモリにフレームをストアするステップと、

フレームが送信されるときに、送信されたフレームと同じ識別子を有するエントリを求めて第1のメモリを探索するステップと、

第1のメモリの探索が、送信されたフレームと同じフレーム識別子のエントリをつきとめると、第1のメモリ内のそのエントリのコピー数をデクリメントするステップと、

マルチコピーキューにエントリをロードするステップをさらに含み、各エントリはフレームポインタおよびコピー数を有し、各コピー数は正または負の値であり、正の値のコピー数は送信されるべきフレームのコピーの総数を表わし、負の値のコピー数は既に行なわれたフレームのコピーの1つについての送信を表わし、第1のメモリにエントリをロードするステップは、正のコピー数を有するマルチコピーキューから出るエントリのみをロードするステップを含む、方法。

【請求項20】

正のコピー数のエントリの書込の前に第1のメモリがエントリでフルであるかどうかを判断し、第1のメモリがフルであれば、第1のメモリから最古のエントリを除去し、そのエントリのコピー数を、除去されたエントリ内のフレームポインタが指すフレームをストアする第2のメモリ内の場所のコピー数領域へと書込むステップをさらに含む、請求項19に記載の方法。

【請求項21】

第1のメモリの探索が負のコピー数を有するエントリと同じポインタを有するエントリを第1のメモリ内につきとめないとき、負のコピー数を含むエントリ内のフレームポインタが指すフレームをストアする場所のコピー数領域からコピー数を取出すステップをさらに含む、請求項20に記載の方法。

【発明の詳細な説明】

発明の分野

本発明は通信分野に関連し、特に、データの複数のコピーを送信し、送信されたコピーの数のカウントを維持するための方法および装置に関する。

背景技術

EP-A-0 622 922は、通信システム内の共用メモリに接続された複数のユーザにデータまたはメッセージをマルチキャストする公知の方法および装置を開示している。説明されたシステムにおいては、メモリは複数のデータバッファとして組織化され、各々が独自の直接制御ブロックおよび間接制御ブロックを含む制御ブロックによって制御される。チェーン状に連結されたデータバッファから成るメッセージは、対応の直接制御ブロックをチェーン状に連結することによって形成され、関連の間接制御ブロックを用いることによって、各ユーザごとにデータを書き換えずに複数のユーザにマルチキャストされ得る。各直接制御ブロックは複製数のカウントを可能にするフィールドを有する。Lee, Tsern-Huei他による「共用バッファメモリスイッチにおけるマルチキャスト」( Multicasting in a shared buffer memory switch )、TENCON '93、北京、1993年10月19-21日、第1巻、IEEE第209-212頁は、マルチキャスト能力を備える公知の共用バッファメモリスイッチを開示している。

10

Schultz, Kenneth J他「CAMベース単一チップ共用バッファATMスイッチ」( CAM-based single chip shared buffer ATM switch )、Supercomm-ICC、ニューオーリンズ、1994年5月1-5日、第2巻、IEEE第1190-1195頁もまた、共用RAMを有する公知の共用バッファスイッチアーキテクチャを開示している。

パケット交換網(たとえば、イーサネット網)のような多くのネットワークシステムにおいて、単一のエンドステーションが同じデータをネットワーク上の複数の他のエンドステーションに送信することがしばしば望ましい。たとえば、従来の電子メールシステムでは、ユーザは同じ電子メールメッセージをその電子メールシステムに接続された異なる4人のユーザに送ることを望むかもしれない。

20

パケット交換網では、スイッチがそのポートを通してエンドステーションからデータのフレームを受信すると転送決定を行なう。フレームが多数のエンドステーションに送信されるべきであれば、スイッチはフレームを正しいポートに転送するための転送決定を行わなければならない。このようなシステムにおける問題の1つはフレームの正しい数のコピーが送信されることを確実にすることである。この問題は、情報がそれを通して送信され得る多数のポートをスイッチが有する場合に特に達成が困難である。たとえば、フレームのコピーがスイッチの3つの異なるポートから同時に送信されるのであれば、送信されたコピーの数の正確なカウントを維持することは問題である。

また、送信されたコピー数の正確なカウントを維持して、そのデータを一時的に保持するために用いられている資源が他のデータをストアするためにいつ再使用できるかがスイッチにわかるようにすることが重要である。この情報の提供が時宜を得たものでなければ、データをストアするための資源に不必要な輻輳およびコンテンションが起こり得る。

30

パケット交換網のネットワークスイッチのような装置から送信された、フレームのようなデータ項目のコピー数のカウントを維持する構成および方法が必要である。

#### 発明の概要

本発明は、添付の請求項1に記載の特徴によって特徴づけられる、ネットワークスイッチから送信されるフレームのコピーのカウントを維持するための構成を提供する。本発明は、添付の請求項9に記載の特徴によって特徴づけられる、装置から送信されるべきデータ項目のコピー数のカウントを維持するための構成をさらに提供する。さらに、本発明は、添付の請求項12に記載のステップによって特徴づけられる、ネットワークスイッチからのフレームの送信数のカウントを維持する方法を提供する。

40

上述および他の必要は、ネットワークスイッチから送信されるフレームのコピーのカウントを維持するための構成であって、マルチコピーキュー、探索可能なメモリおよびバッファマネージャを含む構成を提供する本発明の実施例によって満たされる。マルチコピーキューはエントリを受取り、キューに入れる。各エントリはフレーム識別子およびコピー数を有する。各コピー数は正または負の値であり、正の値のコピー数は送信されるべきフレームのコピーの総数を表わし、負の値のコピー数は既に行なわれたフレームのコピーの1つの送信を表わす。メモリは、正のコピー数を有するマルチコピーキューからのエントリを受取り、ストアするよう構成される。バッファマネージャはマルチコピーキューを出た

50

後のエントリを検査する。マルチコピーキューからのエントリが負のコピー数を含む（そのコピーが送信されたことを示す）とき、バッファマネージャは負のコピー数を有するエントリと同じフレーム識別子を有するエントリを求めて第1のメモリを探索する。バッファマネージャは、メモリ内のエントリのコピー数が1よりも大きければメモリ内のエントリのコピー数をデクリメントし、コピー数が1と等しければメモリ内のエントリをデリートする。

本発明のこの実施例の利点は、コピーが異なるポートから同時に送信され得るマルチポートスイッチ環境においても正確なカウント数の送信が行なわれることである。カウントはキャッシュメモリのようなメモリ内に各複数コピーフレームごとに維持される。このフレームのコピーが送信されるたびに、バッファマネージャはそのフレームのコピー数（すなわち、送信数のカウント）を変更する。フレームコピーの送信を示すエントリをキューに入れるためのキューを用いることによって、コピーの各送信を認識し、それに従ってメモリ内のカウントを調節する機会がバッファマネージャに与えられる。これは、全送信が同時に起こっても認識されるように、キューがバッファマネージャにコピーの送信を順次的に通知するためである。

上述の必要は、エントリをストアするよう構成される第1のメモリを含み、各エントリが第2のメモリ内にストアされるデータを指すデータポインタと関連のデータ項目とを有するメモリ構成を提供する本発明の別の実施例によって満たされる。第2のメモリが、第2のメモリ内のデータポインタが指す場所にデータをストアする。第2のメモリ内の場所は、そのデータとともに関連のデータをストアするよう構成され、したがって、第2のメモリは関連のデータ項目をストアするための第1のメモリの拡張部である。

上述の実施例によって与えられる利点の1つは、第2のメモリが必要に応じて第1のメモリの直接アドレス指定可能な拡張部となることができる点である。データと関連し、通常は第1のメモリ内にストアされるデータ項目がデータ自体とともに第2のメモリ内にストア可能である。したがって、この構成は必要な場合にのみデータ項目をストアするためのより大容量のメモリを与え、したがって、第1のメモリは小さく、比較的高速に保たれ得る。

本発明の別の実施例は、装置から送信されるべきデータ項目のコピー数のカウントを維持するための構成を提供し、エントリがストアされる第1のメモリとマネージャ装置とを含む。第1のメモリ内にストアされる各エントリは、複数のコピーが装置から送信されるべき異なるデータ項目に対応する。各エントリは、第2のメモリ内のデータ項目がストアされている場所を指すポインタと、データのいくつかのコピーがまだ送信されていないかを示すコピー数とを有する。マネージャ装置は、第1のメモリ内のエントリと同じポインタを有するデータ項目のコピーが装置から送信されるときに、第1のメモリ内にストアされているエントリの1つのコピー数を調節する。

本発明の別の実施例では、パケット交換網のためのマルチポートスイッチ構成が設けられ、スイッチは1つのフレームの複数のコピーを送信するよう構成される。スイッチ構成はエントリを受取り、キューに入れるよう構成されるマルチコピーキューを含み、各エントリはフレームポインタおよびコピー数を有する。各コピー数は正または負の値である。正の値のコピー数は送信されるべきフレームのコピーの総数を表わし、負の値のコピー数は既に行なわれたフレームのコピーの1つの送信を表わす。探索可能な第1のメモリが正のコピー数を有するマルチコピーキューからのエントリを受取り、ストアする。バッファマネージャがマルチコピーキューを出るエントリを検査する。マルチコピーキューからのエントリが負のコピー数を含むとき、バッファマネージャは負のコピー数を有するエントリと同じフレームポインタを有するエントリを求めて第1のメモリを探索する。バッファマネージャは、第1のメモリ内のエントリのコピー数が1よりも大きければそのフレームポインタを有する第1のメモリ内のエントリのコピー数をデクリメントする。コピー数が1と等しければ、バッファマネージャは第1のメモリ内のそのエントリをデリートする。

本発明のスイッチ構成は、送信の指示をキューに入れ、送信の指示がキューを出るときにメモリ内にストアされているコピー数を変更することによって、フレームコピーの送信数

10

20

30

40

50

のカウンタを維持する。負のコピー数をキューに入れ、それらを一度に1つバッファマネージャに与えて処理することによって、複数のポートを通しての同時のコピー送信がある場合でもフレームの送信されたコピーの正確なカウンタが維持され得る。

上述の必要は、ネットワークスイッチからのフレームの送信数のカウンタを維持する方法を提供する本発明によって満たされる。この方法は探索可能な第1のメモリにエントリをロードするステップを含み、各エントリは、フレーム識別子と、ネットワークスイッチからまだ送信されていないフレームのコピー数を示すコピー数とを含む。フレームが送信されるとき、送信されたフレームと同じフレーム識別子を有するフレームを求めて第1のメモリが探索される。第1のメモリの探索が送信されたフレームと同じフレーム識別子を有するエントリをつきとめると、第1のメモリ内のそのエントリのコピー数はデクリメントされる。

10

本発明の上述および他の特徴、局面および利点は、添付の図面を参照すると本発明の以下のより詳細な説明からより明らかとなるであろう。

#### 【図面の簡単な説明】

図1は、この発明の実施例に従って構成されたパケット交換システムのブロック図である。

図2は、この発明の実施例に従って構成され、図1のパケット交換システムに用いられるマルチポートスイッチのブロック図である。

図3は、この発明の実施例に従って構成された、図2のマルチポートスイッチのスイッチサブシステムの概略図である。

20

図4は、この発明の実施例に従って構成された、図3のスイッチサブシステムの単一の出力キューを示すブロック図である。

図5は、この発明の実施例による第1のタイプの出力キューを詳細に示す図である。

図6は、この発明の実施例による第2のタイプの出力キューを詳細に示す図である。

図7は、この発明の実施例に従って構成された、外部メモリのオーバフロー領域を詳細に示す図である。

図8は、この発明に採用されるリンクトリストデータ構造のブロック図である。

図9は、この発明の実施例によるフレームバッファヘッダフォーマットを概略的に示す図である。

図10は、この発明の実施例に従って構成された、図4のスイッチサブシステムのマルチコピー、リクレーンおよびフリーバッファプール領域を詳細に示す図である。

30

図11は、この発明の実施例に従って構成されたフリーバッファプール構造のブロック図である。

図12は、この発明の実施例に従って構成されたマルチコピーキューのブロック図である。

図13は、この発明の実施例に従って構成されたマルチコピーキャッシュの概略図である。

図14は、この発明の実施例に従って構成された、スイッチサブシステムのバッファマネージャのキュー部およびポートベクタFIFOのブロック図である。

#### 例示的な実施例の詳細な説明

40

イーサネット(IEEE 802.3)網などのパケット交換ネットワークにおけるスイッチを例に挙げてこの発明を説明する。しかしながら、以下に詳細に説明するように、この発明は他のパケット交換システムおよび一般的な他のタイプのシステムにも適用可能であることが明らかとなるであろう。

図1は、この発明が有利に採用され得る例示的なシステムのブロック図である。例示的なシステム10はイーサネット網などのパケット交換ネットワークである。パケット交換ネットワークは、ネットワークステーション間でのデータパケットの通信を可能にする統合マルチポートスイッチ(IMS)12を含む。ネットワークはたとえば10Mbpsのネットワークデータレートでデータの授受を行なう24個の毎秒10メガビットの速度(Mbps)のネットワークステーション14と、100Mbpsのネットワーク速度

50

でデータパケットの授受を行なう2つの100Mbpsネットワークステーション22といった、種々の構成を有するネットワークステーションを含み得る。したがって、スイッチ12はネットワークステーション14または22から受けたデータパケットをイーサネットプロトコルに基づく適切な宛先に選択的に転送する。

開示される実施例によると、100Mbpsネットワークステーション14は媒体17を介して、かつ半二重イーサネットプロトコルに従って、スイッチ12に対してデータパケットの授受を行なう。イーサネットプロトコルISO/IEC 8802-3 (ANSI/IEEE Std. 802.3, 1993 Ed.) は、すべてのステーション14が等しくネットワークチャネルにアクセスできるようにする半二重媒体アクセス機構を規定する。半二重環境のトラヒックは媒体17と区別されたりまたはそれより優先されることはない。各ステーション14はむしろ、媒体上のトラヒックを認識するために搬送波感知多重アクセス/衝突検出(CSMA/CD)を用いるイーサネットインタフェースカードを含む。媒体上の受信搬送波がデアサートされたことを感知することによりネットワークトラヒックの不在が検出される。送信するデータを有するステーション14はすべて、パケット間ギャップ期間(IPG)として公知である、媒体上の受信搬送波がデアサートされた後、予め定められた時間だけ待機することにより、チャネルにアクセスしようとする。複数のステーション14がネットワーク上に送信するデータを有する場合、ステーションの各々が、媒体上の受信搬送波の、デアサートが感知されたことに応答してIPG期間の後に送信を行なおうとするため、衝突が生じる。したがって、送信ステーションは、別のステーションが同時にデータを送信することにより衝突が生じていないかを判断するために媒体を監視する。衝突が検出されれば、両方のステーションが停止し、ランダムな期間だけ待機し、再度送信を試みる。

100Mbpsネットワークステーション22は好ましくは、提案されているフロー制御によるイーサネット規格IEEE 802.3x全二重-草案(0.3)に従う全二重モードで動作する。全二重環境は各100Mbpsネットワークステーション22とスイッチ12との間に双方向ポイントツーポイント通信リンクを設け、スイッチ12およびそれぞれのステーション22は衝突することなくデータパケットの送受信を同時に行なうことができる。100Mbpsネットワークステーション22の各々は、100ベース-TX、100ベース-T4または100ベース-FXタイプの100Mbps物理(PHY)装置20を介してネットワーク媒体17に結合される。スイッチ12は、物理装置20への接続をもたらす媒体独立インタフェース(MII)24を含む。100Mbpsネットワーク22は他のネットワークへの接続のためのサーバまたはルータとして実現され得る。

図1に示されるように、ネットワーク10は、スイッチ12と100Mbpsステーション14との間で送信されたデータパケットの時分割多重化および時分割非多重化を行なう一連のスイッチトランシーバ26を含む。磁気変成器モジュール19は媒体17上の信号の波形を維持する。スイッチ12は、時分割多重化プロトコルを用いて単一のシリアルノンリターンツーゼロ(NRZ)インタフェース23を介して各スイッチトランシーバ16に対するデータパケットの送受信を行なうトランシーバインタフェース18を含む。スイッチトランシーバ16はシリアルNRZインタフェース23からパケットを受信し、受信されたパケットを非多重化し、ネットワーク媒体17を介して適切なエンドステーション14にそのパケットを出力する。開示される実施例によると、各スイッチトランシーバ16は独立した4つの100Mbpsツイストペアポートを有し、スイッチ12が必要とするPINの数が4分の1に減少するようにするシリアルNRZインタフェースを介する4:1多重化を用いる。

スイッチ12は、意思決定エンジン、切換エンジン、バッファメモリインタフェース、構成/制御/状態レジスタ、管理カウンタ、ならびにネットワークステーション14および12のためのイーサネットポート間でデータパケットの経路制御を行なうためのMAC(媒体アクセス制御)プロトコルインタフェースを含む。スイッチ12はまた、インテリジェントな切換決定を行ない、後に説明するように、外部の管理エンティティに管理情報ベ

10

20

30

40

50



ース ( M I B ) オブジェクトの形式で統計的なネットワーク情報を与えるための優れた機能を有する。スイッチ 1 2 はさらに、スイッチ 1 2 のチップサイズを最小にするためにパケットデータの外部ストアおよびスイッチ論理を可能にするインタフェースを含む。たとえば、スイッチ 1 2 は、受信したフレームデータ、メモリ構造および M I B カウンタ情報をストアするための外部メモリ 3 6 へのアクセスをもたらし同期型ダイナミック R A M ( S D R A M ) インタフェース 3 4 を含む。メモリ 3 6 は 2 M b または 4 M b のメモリサイズを有する 8 0、1 0 0 または 1 2 0 M H z 同期型 D R A M であってもよい。

スイッチ 1 2 はさらに、外部管理エンティティが管理 M A C インタフェース 3 2 によってスイッチ 1 2 の全体的な動作を制御できるようにする、管理ポート 3 0 を含む。スイッチ 1 2 は、P C I ホストおよびブリッジ 2 8 を介して管理エンティティがアクセスできるようにする P C I インタフェース 2 6 をさらに含む。これに代えて、P C I ホストおよびブリッジ 2 8 が複数のスイッチデバイス 1 2 に対する拡張バスとしての役割を果たしてもよい。

スイッチ 1 2 は、1 つのソースから少なくとも 1 つの宛先ステーションに受信データパケットを選択的に送信する内部意思決定エンジン ( 図 2 ) を含む。内部意思決定エンジンには外部ルールチェッカが代用されてもよい。スイッチ 1 2 は外部ルールチェッカインタフェース ( E R C I ) 4 0 を含み、これは内部意思決定エンジンの代わりにフレーム転送決定を行なうために外部ルールチェッカ 4 2 が用いられるようにする。したがって、フレーム転送決定は、内部切換エンジンまたは外部ルールチェッカ 4 2 のいずれかによって行なわれ得る。

スイッチ 1 2 は、ポートごとのステータスをクロックに合わせて出力し L E D 外部論理 4 6 を駆動する、L E D インタフェース 4 4 をさらに含む。L E D 外部論理 4 6 は人間が読取ることができる L E D ディスプレイエレメント 4 8 を駆動する。発振器 3 8 はスイッチ 1 2 のシステム機能に 4 0 M H z のクロック入力を与える。

図 2 は、図 1 の統合マルチポートスイッチ ( I M S ) 1 2 のブロック図である。スイッチ 1 2 はそれぞれの 1 0 M \ b p s ネットワークステーション 1 4 間で半二重のデータパケットの送受信を行なうための 2 4 個の 1 0 M \ b p s 媒体アクセス制御 ( M A C ) ポート 5 0 ( ポート 1 から 2 4 ) と、それぞれの 1 0 0 M \ b p s ネットワークステーション間で全二重のデータパケットの送受信を行なうための 2 つの 1 0 0 M \ b p s M A C ポート 5 3 ( ポート 2 5 および 2 6 ) とを含む。上述のとおり、管理インタフェース 3 0 もまた M A C 層プロトコル ( ポート 0 ) に従って動作する。M A C ポート 5 0、5 3 および 3 0 の各々は、受信先入れ先出し ( F I F O ) バッファ 5 2 と送信 F I F O 5 4 とを有する。ネットワークステーションからのデータパケットは対応の M A C ポートで受信され、対応の受信 F I F O 5 2 にストアされる。受信されたデータパケットは対応の受信 F I F O 5 2 から外部メモリインタフェース 3 4 に出力されて、外部メモリ 3 6 にストアされる。受信されたパケットのヘッダもまた、内部ルールチェッカ 5 8 または外部ルールチェッカインタフェース 4 0 のいずれかである、意思決定エンジンに転送され、いずれの M A C ポートからデータパケットが出力されるかを決定する。具体的には、パケットヘッダは、スイッチ 1 2 が内部ルールチェッカ 5 8 または外部ルールチェッカ 4 2 を用いて動作するよう構成されているか否かに依存して、内部ルールチェッカ 5 8 または外部ルールチェッカインタフェース 4 0 に送られる。内部ルールチェッカ 5 8 および外部ルールチェッカ 4 2 は、所与のデータパケットに関する宛先 M A C ポートを決定するための意思決定論理を提供する。したがって、意思決定エンジンは、単一ポート、マルチプルポートまたは全ポート ( すなわちブロードキャスト ) のいずれかに所与のデータパケットを出力し得る。たとえば、各データパケットにはソースおよび宛先アドレスを有するヘッダが含まれ、意思決定エンジンは宛先アドレスに基づいて適切な出力 M A C ポートを特定する。これに代えて、宛先アドレスは、適切な意思決定エンジンが複数のネットワークステーションに対応するものと特定するバーチャルアドレスに対応してもよい。これに代えて、受信されたデータパケットは、( 1 0 0 M \ b p s ステーション 2 2 のうちの 1 つのルータを介する ) 別のネットワークまたは所定のグループのステーションを特定する I E E E 8 0 2 . 1 d

10

20

30

40

50

プロトコルに準拠するVLAN(バーチャルLAN)タグ付フレームを含んでもよい。したがって、内部ルールチェッカ58または外部ルールチェッカ42のいずれかがインタフェース40を介して、バッファメモリ36に一時的にストアされたフレームが単一のMACポートまたは複数のMACポートに出力されるべきかを決定する。

外部ルールチェッカ42を使用することにより、容量の増加、およびフレームが外部メモリに完全にバッファされる前にフレーム転送決定を可能にし、かつスイッチ12がフレームを受信する順からは独立した順で決定が行なわれるようにする、決定キューのうちランダムな順序付け、といった利点がもたらされる。

意思決定エンジン(すなわち内部ルールチェッカ58または外部ルールチェッカ42)は、データパケットを受信すべき各MACポートを特定するポートベクタの形式で転送決定をスイッチサブシステム56に出力する。ルールチェッカからのポートベクタは、外部メモリ36にデータパケットをストアするアドレス場所と、データパケットを受信して送信するためのMACポート(たとえばMACポート0から26)の識別子を含む。スイッチサブシステム56はポートベクタに特定されたデータパケットを外部メモリインタフェース34を介して外部メモリ36から取出し、取出されたデータパケットを特定されたポートの適切な送信FIFO54に与える。

付加的なインタフェースにより管理および制御情報が与えられる。たとえば、管理データインタフェース59は、MII管理仕様(IEEE 802.3u)に従うスイッチトランシーバ16および100M/bps物理装置20と制御およびステータス情報をスイッチ12が交換できるようにする。たとえば、管理データインタフェース59は、双方向管理データIO(MDIO)信号経路に時間基準を与える管理データクロック(MDC)を出力する。

PCIインタフェース26は、PCIホストプロセッサ28によって内部IMSステータスおよび構成レジスタ60にアクセスし、かつ外部メモリ36にアクセスするための、32ビットPCI改訂2.1に適合したスレーブインタフェースである。PCIインタフェース26は複数のスイッチデバイスのための拡張バスとしての役割も果たし得る。管理ポート30は標準7ワイヤ反転シリアルGPSインタフェースを介して外部MACエンジンにインタフェースされ、標準MAC層プロトコルによりホストコントローラがスイッチ12にアクセスできるようにする。

図3は、この発明の例示的な実施例に従う、図2のスイッチサブシステム56をより詳細に説明する図である。図2に示されるマルチポートスイッチ12の他のエレメントは、スイッチサブシステム56とこれらの他のエレメントとの接続を示すために図3に再度示される。スイッチサブシステム56はフレームの受信および転送を行なうためのコアスイッチングエンジンを含む。スイッチングエンジンを実現するために用いられる主な機能ブロックは、ポートベクタFIFO70と、バッファマネージャ72と、複数のポート出力キュー74と、管理ポート出力キュー75と、拡張バスポート出力キュー77と、フリーバッファプール104と、マルチコピーキュー90と、マルチコピーキャッシュ96と、リクレーンキュー98とを含む。これらの機能ブロックの動作および構成は後により詳細に説明するが、まず、個々のエレメントに関する後の説明に関連性を持たせるために、図3のスイッチサブシステム56の全体像を簡単に説明する。

ポートからマルチポートスイッチ12に入るフレームには基本的に2つのタイプがある。すなわち、単一コピーフレームとマルチコピーフレームとである。単一コピーフレームは、マルチポートスイッチ12によって他の1つのポートにのみ送られることとなる、ポートで受信されたフレームである。これとは対照的に、マルチコピーフレームは、1つのポートで受信され、1つより多い数のポートに送信されるフレームである。図3では、各ポートは別個のMAC50によって表わされ、それ自体の受信FIFO52および送信FIFO54を有する。

単一コピーまたはマルチコピーであるフレームは内部MACエンジン50によって受信される。フレームパケットがポートで受信されると、それは受信FIFO52に置かれる。各フレームはヘッダを有し、これは、内部ルールチェッカ58または外部ルールチェッカ

10

20

30

40

50

42のいずれかのルールチェッカに与えられる。ルールチェッカ42または58は、ヘッダの情報に基づいて、フレームパケットがどこから送り出されるかを決定し、すなわちいずれのポートを介してフレームパケットが送信されるかを決定する。

ルールチェッカ42または58が転送決定を行なうのと同時に、バッファマネージャ72はフリーバッファプール104からフリーバッファポイントを得る。このフリーバッファポイントは、受信FIFO52によってフレームがストアされることとなる外部メモリ36の場所である。バッファマネージャ72によってフリーバッファポイントがフリーバッファプール104から得られると、フリーバッファポイントによってポイントされるバッファはもはやフリーであるとは考えられない。フレームデータは、直接メモリアクセス(DMA)トランザクションでデータバス80を介して受信FIFO52から外部メモリ36に転送される。フレームはフリーバッファプール104から得られたフリーバッファポイントがポイントする場所にストアされるが、後に説明するように、フレームをストアするために多くの他のバッファが用いられてもよい。

ヘッダデータの他に、ルールチェッカ42または58はバッファマネージャ72からのフリーバッファポイントも受信する。このフリーバッファポイントはここではフレームポイントと呼ばれる。なぜなら、フレームがストアされる外部メモリ36でのメモリ場所をポイントするからである。ルールチェッカ42または58は、転送決定を行ないかつ「ポートベクタ」の形式で転送命令を発生するためにヘッダ情報を用いる。図示される例示的な実施例では、ポートベクタは、フレームが転送されるべき各出力ポートに対してセットされたビットを備えた28ビットベクタである。この全体像での例では、受信されたフレームは単一フレームであると想定する。したがって、ルールチェッカ42または58によって生成されたポートベクタには1つのビットしかセットされない。ポートベクタにセットされたビットはポートのうち特定のなものに対応する。

ルールチェッカ42または58はポートベクタFIFO70にポートベクタおよびフレームポイント(ならびに制御操作コードおよびVLANインデックス)を置く。ポートベクタはポートベクタFIFO70によって検査され、ポートベクタに関連したフレームポイントがどの特定の出力キュー74に入力されるべきかを決定する。ポートベクタFIFO70は適切な出力キュー74の一番上にフレームポイントを置く。これによりフレームの送信がキューとして維持される。

ある時点で、フレームポイントは出力キュー74を通過した後に出力キュー74の一番下まで到達する。バッファマネージャ72はそれが出力キュー74の一番下まで到達したときにフレームポイントを取り、フレームポイント書込バス86を介して正しいポートの適切な送信FIFO54にそのフレームポイントを送る。これによりフレームの送信がスケジュールされる。フレームポイントによってポイントされた外部メモリ36での場所からDMAトランザクションにおいて読出されたフレームデータは、適切な送信FIFO54に置かれ後に送信される。

マルチコピー送信は、ポートベクタが、フレームがそれらから送信されることとなるマルチプルポートを示す、セットされた複数のビットを有する点を除いて、単一コピー送信と同様である。フレームポイントは適切な出力キュー74の各々に置かれ、対応の送信FIFO54から送信される。

バッファマネージャ72は特殊な制御キューを用い、すなわち、フリーバッファプール104と、マルチコピーキュー90と、リクレームキュー98と、マルチコピーキャッシュ96とを用いて、受信フレームをストアするためにバッファを割当て、フレームがその指定された出力ポートに送信されると再度使用できるようバッファを取出すプロセスを管理する。後により詳細に説明するが、バッファマネージャ72はまた、出力キュー74ならびに制御キュー104、90および98のために外部メモリ36に「オーバフロー」領域を維持する。

この動作上の全体像を背景として、以下にスイッチサブシステム56の個々のセクションおよびさまざまな局面をより詳細に説明する。これらの局面のうち最初に説明するものは、この発明のさまざまな出力キュー74の構造である。10Mb/sポートおよび100

10

20

30

40

50

Mb/s 出力ポートに指定される出力キュー 74 の他に、管理ポート 30 のために出力キュー 75 が設けられ、拡張ポート 26 のために出力キュー 77 が設けられる。これらの出力キュー 75 および 77 は出力キュー 74 と同じ外部構成を有するが、後に説明するように、異なった内部構成を有する。

図 4 は、この発明の実施例に従う出力キュー 74 の外部構成を示すブロック図である。図 4 から明らかなように、この発明の出力キュー 74 は 3 部構成である。性能を最も高くするためには、チップ上のキュー構造のすべてを保持することが好ましいが（マルチポートスイッチ 12 を参照）、チップの占有面積に関する費用は非常に高い。これにより、チップが多数のエントリの切換を行ない、それらをキューとして維持する必要があるときにはジレンマが生じる。この発明は、チップ上に高性能な小容量セクションを含み、チップ外

10

にオーバフロー領域を含む、単一の出力キューを与えることによりこのジレンマを解消する。オーバフロー領域は、チップ上の領域よりも比較的性能が低いにも関わらず、所要の大容量のキューとしてキューが役割を果たすようにする。

図 4 の実施例に従うこの発明の単一論理出力キュー 74 は 3 つの物理セクションを有する。これらには、出力キュー書込側 76 と、出力キュー読出側 78 と、外部メモリ 36 にある出力キューオーバフロー領域（全体が 110 として示される）とが含まれる。出力キュー 74 のすべてに関する外部メモリ 36 へのアクセスは、前述のとおり外部メモリインタフェース 34 を介するものである。この発明は、現在の外部メモリのバースト的な性質を利用し、（フレームポインタなどの）データが、チップ 12 を外部メモリ 36 に接続するバス 84 を介してバースト状にチップの内外からオーバフローキュー領域 110 に送られるようにする。

20

出力キュー書込側 76 および出力キュー読出側 78 はチップ 12 上にある。書込側 76 および読出側 78 は小さくて値段の高い資源であると考えられる。これとは対照的に、出力キュー 74 の第 3 の部分を形成するオーバフロー領域 110 は大きくて比較的安価である。書込側 76 および読出側 78 により高い性能がもたらされ、オーバフロー領域を通る経路によっては低性能で大容量の経路がもたらされる。

動作時に、出力キュー書込側 76 はエントリを受信する。この発明に従うマルチポートスイッチ 12 の例示的な実施例では、エントリは、フレームの最初の 256 バイトがストアされる外部メモリの第 1 のバッファをポイントするフレームポインタである。しかしながら当業者には、出力キューの構成 74 はエントリとしてのフレームポインタに制限されず、マルチポートスイッチおよび他の技術の両方において、他のタイプのエントリをキューとして維持することに広く適用可能であることが明らかであろう。

30

エントリが出力キュー書込側 76 内を完全に移動し、その一番下まで到達すると、出力キュー 74 に関連した制御論理はエントリをどう処理するか決定する。出力キュー読出側 78 にスペースがあれば、出力キュー 74 のオーバフロー領域 110 は空いており、1 つまたはそれ以上のエントリが出力キュー書込側 76 から出力キュー読出側 78 に直接送られる。書込側 76 から読出側 78 に直接エントリを送ることはすべてチップ 12 上で行なわれるため、エントリは低レイテンシーで素早く完全に送られる。

出力キュー読出側 78 がいっぱいであり、出力キュー書込側 76 に少なくとも 1 バーストサイズの量のデータ（たとえばエントリの 16 バイト分）があれば、データはその出力キュー 74 のオーバフロー領域 110 にバースト状に書込まれる。出力キュー読出側 78 がいっぱいであり、かつ出力キュー書込側 76 にはまだ 1 バーストサイズの量のデータがないときは、エントリは出力キュー書込側に留まりさらに処理は行なわれない。最終的には、出力キュー読出側 78 は空になり、出力キュー読出側 78 に 1 バーストサイズの量のデータを収容する十分なスペースが生まれ、かつオーバフロー領域 110 にデータがあるときがくると、オーバフロー領域 110 から出力キュー読出側 78 に 1 バーストのデータが与えられる。

40

出力キュー構成において、読出側 78 は伝統的なキューとほぼ同様に作用する。なぜなら、エントリが 1 つずつ取出されるのはこの部分からであるからである。出力キュー書込側 76 は主に、データをバーストに組立てて外部メモリ 36 に書込むための回収機能を果た

50

す。したがって、この発明は単一の事象（エントリを出力キュー 74 に置くこと）をバースト事象に変える。書込側 76 は、蓄積されたデータが必要に応じて外部メモリ 36 のオーバフロー領域 110 にバーストされるようにする。比較的稀な場合にのみ必要となる機能に高価なチップ資源を提供するのではなく、輻輳時にオーバフロー領域 110 が安価なストレージを提供する。この発明はチップ外のオーバフロー領域 110 を利用するが、この領域 110 のアクセスは、1 度に多くのバイトの情報をバーストすることにより効率よく行なわれる。これは、単一のエントリがキューに対して書込まれたり読出されたりする従来のキュー構造とは対照的である。

動作時に、出力キュー 74 に到達するエントリが多ければ、これらのエントリはオーバフロー領域 110 に置かれ、チップ上のキュー 78 のオーバフローを回避するようにする。したがって、この発明のキュー構造を用いるとフレームの廃棄が大幅に防止される。また、オーバフロー領域 110 のためのメモリの合計量は、外部メモリ 36 のサイズを変更することにより容易に変更可能である。さらに、個々の特定のオーバフロー領域 110 のサイズは、出力キュー 74 の性能に影響を及ぼすことなくキューのサイズをカスタマイズするためにプログラム可能である。

典型的に、キューは、先入れ先出し構成を有する順序づけられた構成である。しかしながら、リクレームキュー 98 およびフリーバッファプール 104 などのいくつかのタイプのキューでは、エントリの順序は問題ではない。書込側 100 から読出側 102 にデータを直接送信することが可能であれば、この発明はそのキューに関するオーバフロー領域を迂回して情報がこの経路に直接送信されるようにする。これは、情報が順番によって影響を受けない限り、関連のオーバフロー領域に情報がある場合でも可能である。たとえば、バッファの再要求は順番によって影響を受けない。なぜなら、バッファがフレームにストアされる必要がなくなった後に、最終的にバッファがフリーバッファプール 104 のフリーリストに戻される順番は、いかなるものでも許容されるからである。したがって、データが順番によって影響を受けない場合に外部メモリ 36 のリクレームキュー 98 のオーバフロー領域 110 への書込についての帯域幅が生じるのを回避するために、読出側 102 にさらなるエントリのためのスペースがあるものと想定して、書込側 100 から読出側 102 に情報が直接送られる。リクレームキュー 98 は順番によって影響を受けないデータをキューとして維持するタイプのキューの一例である。しかしながら、順番によって影響を受けない他の多くのタイプのデータが種々の適用例で可能であるため、この発明のこの特徴は、他のタイプのデータをキューとして維持するキューにおいて有用性を見出す。

図 1 および図 2 に示されるこの発明の例示的な実施例のマルチポートスイッチには 28 個の出力キュー（各々が出力ポートと関連する）があり、すなわち、10 Mb/s ユーザポートに関するものが 24 個、100 Mb/s サバポートに関するものが 2 つ、管理ポートに関するものが 1 つ、そして拡張バスポートに関するものが 1 つある。出力キュー 74、75 および 77 は、フレームポインタが送信のためにキューとして維持されるときにそれらに一時的なストレージを提供する。キュー作業は、転送ポートベクタに示されるさまざまな出力キュー 74、75 および 77 に対してポートベクタ FIF070 がフレームポインタを書込むという形態をとる。

この発明のある好ましい実施例では、さまざまな出力キュー 74、75 および 77 は以下のフィールドのうちいくつかまたはすべてを含む。すなわち、単一コピービットと、フレームポインタと、制御操作コードまたは制御信号と、VLAN（バーチャルローカルエリアネットワーク）インデックスとである。単一コピービットは 1 つの出力ポートにのみ転送されることとなるフレームを示す。フレームポインタは外部メモリ 36 のフレームをポイントする。制御操作コードはフレームに関する特定のな情報（すなわち新たに得たフレームなど）を識別する。制御信号は制御操作コードからの情報を用いて、送信前にポートによってフレームがいかにかに処理されるかを示す。VLAN インデックスは、外部へのフレームに（必要であれば）挿入されるべき VLAN タグに対する基準を与える。しかしながら、この発明は種々のタイプのフィールドを有する他の出力キューにも適用可能であるため、これらのフィールドは例としてのみのものである。

10

20

30

40

50

第1のタイプの出力キュー74、すなわち10Mb/sポート出力キューの例示的な実施例の内部構成が図5に示される。10Mb/s出力キュー74は10Mb/sポートに転送されることとなるフレームのエントリを保持する。これらのキューの出力キュー書込側76は32個のエントリを保持し、出力キュー読出側78は図示される例示的な実施例において16個のエントリを保持するが、考えられる他のサイズのものもこの発明の範囲内である。10Mb/s出力キュー74は単一コピービットとフレームポインタ(14ビット)とを含む。この発明のマルチポートスイッチの例示的な実施例では、10Mb/sポートにはVLANTAGがないためVLANインデックスは必要ない。

第2のタイプの出力キュー74、すなわち100Mb/sポート出力キューの例示的な実施例の内部構成が図6に示される。100Mb/sポート出力キューは100Mb/sポートに転送されることとなるフレームのエントリを保持する。出力キュー書込側76はこのタイプの出力キューに64個のエントリを保持し、出力キュー読出側は16個のエントリを保持する。各エントリはVLANインデックスと、部分的な制御操作コード(ビット4-0)と、単一コピービットと、フレームポインタとを含む。

外部メモリ36の例示的なマップが図7に示される。外部メモリ36の全体の容量はたとえば4Mbであるが、種々の実施例において他の容量のメモリが採用されてもよい。この発明に従ってオーバフロー領域に外部メモリ36を使用することにより、外部メモリを変更するだけで出力キューのサイズを増減することができる。これは、キューとして維持する容量全体がチップの製造時に設定される、キュー構成がすべてチップ上にあるシステムよりも有利である。

スイッチ12のストア要件を満たすために、外部メモリ36の例示的な実施例は下記の領域にスペースを割当てる。すなわち、フリーバッファプールオーバフロー120と、リクレーンキューオーバフロー122と、マルチコピーキューオーバフロー124と、管理ポート出力キューオーバフロー126と、10Mb/sおよび100Mb/s宛先ポートの各々のための個々の出力キューオーバフロー128と、拡張バスポート出力キューオーバフロー130と、MIBカウンタ132と、グローバルフレームバッファプール134とである。

メモリ領域全体のBASEアドレスはチップ上のレジスタ60の中のメモリベースアドレスレジスタ内でプログラム可能である。外部メモリマップ内の各領域のBASEアドレスはレジスタセット内でプログラム可能である。領域長レジスタは不要である。所与の領域の長さは、マッピング内のその領域のBASEアドレスから次の領域のBASEアドレスまでの領域に等しい。

個々のオーバフロー領域の長さ(したがって容量)がプログラム可能であるため、各キューの容量全体がプログラム可能である。この発明のこの特徴により、必要に応じて容量の増大した特定の出力キューを提供するようにスイッチをカスタマイズすることが可能になる。

したがって、チップ12上の制御キューに適合しない後続のオーバフロー領域ストアエントリは外部メモリ36に置かれる。フリーバッファプールオーバフロー領域120はアドレスポインタをグローバルフレームバッファプール134中の未使用のバッファにストアする。リクレーンキューオーバフロー領域122は、必要でなくなったリンクトリストチェーンにフレームポインタをストアする。マルチコピーキューオーバフロー領域124は(キューとして維持されたフレームポインタについては)コピーナンバー「1」を、かつ(うまく送信されたフレームについては)コピーナンバー「-1」を付してフレームポインタをストアする。

後続のオーバフロー領域は、チップ上に入らない出力キューのエントリをストアする。管理ポート出力キューオーバフロー領域126は管理ポートへの送信を待機するフレームポインタをストアする。出力キューオーバフロー領域128は適切な10Mb/sまたは100Mb/sポートへの送信を待機するフレームポインタをストアする。拡張バスポート出力キューオーバフロー領域130は拡張バスポートへの送信を待機するフレームポインタをストアする。

10

20

30

40

50

M I B カウンタ領域 1 3 2 は、スイッチ 1 2 によって周期的に更新されるポートごとの統計をすべて含む。スイッチ 1 2 は M I B 統計をストアするための 8 ビットおよび 1 6 ビットカウンタをチップ上に維持する。スイッチ 1 2 は M I B データの損失を防止するために要求される周波数で、外部メモリ 3 6 の 3 2 ビットまたは 6 4 ビットの M I B カウンタを更新する。

グローバルフレームバッファプール 1 3 4 は、受信されたフレームデータをストアするリンクトリストのバッファを含む。任意の時点で、これらリンクトリストは有効フレームデータと無効になったバッファとを含み、無効になったこれらのバッファは、バッファマネージャ 7 2 によってフリーバッファプール 1 0 4 に戻されるか、または P C I ホストプロセッサ 2 8 の所有となる。

10

次に図 8 を参照して、いずれかの M A C ポートまたは P C I バスから受信されたフレームデータは、この発明の例示的な実施例におけるリンクトリストデータ構成のフォーマットで外部メモリ 3 6 にストアされる。リンクトリストを生成するために用いられるバッファ 1 4 0 の長さは 2 5 6 バイトであるが、発明の種々の実施例では他の長さのバッファ長が採用されてもよい。これらのバッファ 1 4 0 の各々へのアドレスポインタはスイッチ 1 2 内のフリーバッファプール 1 0 4 によってストアされる。

スイッチ 1 2 のポートのうち 1 つにフレームが受信されると、バッファマネージャ 7 2 はフリーバッファプール 1 0 4 からアドレスポインタを要求し、バッファ 1 4 0 をリンクしてフレームをストアするようにする。フレームをストアする外部メモリ 3 6 の第 1 のバッファに対するアドレスポインタが、そのフレームに対するフレームポインタになる。フレームポインタは、送信されることとなるフレームをキューとして維持するためのスイッチサブシステム 5 6 において用いられる。

20

バッファ 1 4 0 は、メモリの次のバッファの場所を示す各バッファヘッダ 1 4 2 のアドレスポインタによって互いに繋がれる。バッファヘッダ 1 4 2 はまた、バッファ 1 4 0 に含まれるフレームデータに関する他の情報を含む。図 9 a の例示的なバッファヘッダフォーマットに示されるように、先頭のバッファのヘッダは 1 2 バイトである。図 9 b に示されるように、後の各バッファのヘッダは 4 バイトである。外部メモリバーストは、2 バンク × 1 6 バイトの長さであるため、各バッファの実際のフレームストア容量は  $256\text{B} - 16\text{B} = 240\text{B}$  である。

図 9 a および図 9 b に示されるように、先頭および後のバッファヘッダフォーマットは下記のフィールドを含む。

30

バッファフォーマットビット：どのバッファフォーマットが使用中であることを示す。1 は 1 2 バイトの長さの先頭バッファフォーマットを示す。0 は 4 バイトである後のバッファのフォーマットを示す。バッファを繋ぐ際に残りのバッファの各々に関して用いられる。

E ビット（フレームマーカの最後）：フレームに関する最後のバッファであることを示す。E ビットがセットされていれば、チェーンにはこれ以上バッファはない。

C ビット（C R C エラー検出）：C R C エラーが受信機によって検出されたことを示す。C ビットが検出されると、送信機能は反転された C R C を意図的に送信する。

L ビット（整列エラー）：フレーム整列エラーが（C R C エラーとともに）受信フレームに検出されたことを示す。

40

O ビット（受信 F I F O オーバフロー）：受信 F I F O がオーバフローし、バッファのデータが有効でないかもしれないことを示す。

バッファ長さ：バッファヘッダの後の最初のバイトから始まる、バッファのデータフィールドにおいて有効なバイトの合計数。この長さにはオフセットバイト値は含まれるべきではない。

次のバッファポインタ：次のバッファに対するポインタ。次のバッファポインタは E ビットがセットされているときには有効でない。

オフセットバイト数：バッファのフレームデータセクションにおいてフレームの最初のバイトが始まる場所を示す。0 のオフセットは、データがバッファヘッダ 1 4 2 の後の最初のビットで始まることを意味する。0 のオフセットは、データがバッファの 1 6 番目のバ

50

イトに後続するバイトで始まることを示す。オフセットが0でない値の場合、フレームデータは16B + バッファの始まりからのオフセットの後に始まる。送信機能はオフセットバイトフィールドに示されるバイト数だけ飛び越す。

Pビット(ポートタイプ) : 入来する受信フレームのポートタイプを示す。0は10Mb/sポートを示し、1は100Mb/sポートを示す。このビットは、フレームが完全に受信されて外部メモリ36にバッファされる前に、フレームを拡張バスに転送するようスイッチ12をプログラミングする際に、タイムスタンプフィールドに関連してホスト28によって用いられる。

Tビット : 受信されたフレームのタイプを示す。タグ付またはタグ付でない場合がある。1はタグ付のフレームであることを示し、VLAN識別子フィールドは受信VLAN IDを含む。0はタグ付でないフレームを示し、VLAN IDは有効でない。

受信ポート番号 : フレームが受信されたポート番号を示す。

VLAN識別子 : 「タグ付」ポートから受信されたVLAN ID。フレームがタグ付でないポートから受信される場合、このフィールドは無効である。

Rビット(CRC再計算) : CRCを除去し送信機能において再計算する必要があることを示す。スイッチ12はタグ付フレームが受信されるとこのビットをセットする。さらに、ホスト28がフレームの内容を修正した場合、ホスト28はこのビットをセットしなければならない。スイッチ12がフレームを送信すると、スイッチ12はこのビットを検査して、既存のCRCを送信するか、CRCを除去してCRCを再計算するかを判断する。

Aビット(CRC追加) : フレームデータの最後にCRCがないことを示す。ホストはメモリに(CRCなしの)フレームを生成し、このビットをセットすることができる。スイッチ12はフレームの送信時にCRCを発生して追加する。Aビットがセットされている場合、フレームの長さにはCRCは含まれるべきでない。

Fビット(フォーマットビット) : フレーム長/タイムスタンプフィールドを特定する。0はフィールドが入来フレームのタイムスタンプであることを示す。1はフィールドが受信フレームのフレーム長であることを示す。

フレーム長/タイムスタンプ : Fビットに依存する。Fビットがクリアされていると、このフィールドは受信フレームの最初からのタイムスタンプを表わす。タイムスタンプは1μsの分解能を有する。Fビットがセットされている場合には、CRCおよび受信されたVLANタグの全てを含む受信フレームの長さの合計が示される。フレームが受信されると、スイッチ12は(タイマレジスタからの)タイムスタンプでこのフィールドをマークする。フレームが完全に受信される前に拡張バスフレームを転送するようホスト28によってスイッチ12がプログラミングされている場合、フレームデータを過度に読出すことなく外部メモリ36から取出すことができるデータを測定するために(受信ポートの速度とともに)タイムスタンプを用いることができる。フレーム全体が受信されると、スイッチ12はフレーム長をこのフィールドに書込みFビットをセットする。

コピー数 : ポートベクタFIFO70によって送信されるようにうまくキューとして維持されたコピーの数を示すために用いられる。このフィールドは、バッファマネージャ72が、新しいエントリのためのマルチコピーキャッシュ96にスペースを設ける必要がある場合に、フレームポインタのコピー数をストアするために用いられる。

図10は図3のスイッチサブシステム56のいくつかの要素を示す詳細図である。これらの要素は、フレーム記憶のためのバッファを与えるため、かつ、バッファがフレーム記憶のためにもはや必要とされなくなるとこれらのバッファを再要求し、再び使用可能にするために用いられる。上述のように、各出力キュー74、75(出力キュー77を除く)はフレームポインタをバッファマネージャ72に渡し、バッファマネージャ72はフレームポインタが指すフレームの送信をスケジューリングする。バッファマネージャ72は、1)スイッチ12の内部バスを管理し、2)出力キュー74への/からのフレームポインタのキュー入れ/出しを容易にし、3)バッファの位置を決め、フリーバッファプール104に戻すために制御キュー90、98を管理し、4)外部メモリ36を出入りするデータの流れを制御し、5)MIBおよびオーバーフロー領域を含むメモリ構造を維持するという機

10

20

30

40

50



能を制御する。バッファマネージャ 72 は全アクセスを外部メモリ 36 に割当ててするためのスケジューラ機能を含む。これらのアクセスには、1) 受信されたフレームデータを記憶バッファ 140 に書込み、2) 送信のために記憶バッファ 140 からフレームデータを読み出し、3) 出力キュー 74 および制御キュー 90、98 のためのオーバーフロー領域の各々にフレームポインタを維持し(すなわち、書込み、読出す)、4) MIB カウンタを更新することが含まれる。

バッファマネージャ 72 が所与のフレームポインタを適切な全出力キュー 74、75 にコピーした後、ポートベクタ FIFO 70 がコピーの数(「コピー数」)を計算し、フレームポインタおよびコピー数をマルチコピーキュー 90 の書込側に入れる。コピー数は、フレームが転送されるべきでないことを示す「0」、単一コピー送信を示す「1」、または

10

マルチコピー送信を示す「> 1」であり得る。これらの 3 つの場合を以下に説明する。コピー数が「0」であり、フレームポインタがセットされたビットのないヌル転送ポートベクタを有することが意味されているとき、ポートベクタ FIFO 70 はフレームポインタをリクレームキュー 98 の書込側 100 に直接渡す。バッファマネージャ 72 がリクレームキュー 98 を処理するときは、以下に述べるようにバッファマネージャ 72 がバッファのリンクリストチェーンを解体し、各「フリー」バッファごとのアドレスポインタをフリーバッファプール 104 の書込側 106 に戻す。

コピー数が「1」の単一コピー送信のとき、ポートベクタ FIFO 70 はフレームポインタ、制御信号/制御操作コードおよび VLAN インデックスを適切なポートの出力キュー 74 にコピーする。ポートベクタ FIFO 70 は出力キュー 74 内の単一コピービットをセットして(図 5 および図 6 参照)、これが単一の送信であることを示す。バッファマネージャ 72 はそのポートの出力キュー 74 からフレームポインタおよび単一コピービットを読み出すと、上述のように送信をスケジューリングする。バッファマネージャ 72 は、フレームがストアされている外部メモリ 36 において最初のバッファの位置を決めるためにフレームポインタを用いる。バッファマネージャ 72 はこの最初のバッファからバッファヘッダを読み出し、最初のバッファからデータを捕捉し、このデータを適切な MAC 送信 FIFO 54 に入れる。フレームが複数バッファにおよぶ場合を想定すると、そのフレームのためのチェーン内の全バッファを見つけ、送信するためのアドレスを、後続バッファへのリングがバッファマネージャ 72 に与える。データが送信のために FIFO 54 に一旦置かれると、バッファは不使用となり、フリーバッファプール 104 に戻され、結果として別のフレームデータをストアするために再割当される。

20

30

コピー数が 1 よりも大きいとき、ポートベクタ FIFO 70 はフレームポインタ、VLAN インデックスおよび制御信号/制御操作コードを適切な出力キュー 74 の各々にコピーする(キュー 74 に言及する際には、キュー 75、77 もまた言及されている)。ポートベクタ FIFO 70 は出力キュー 74 内の適切なフレームポインタのための単一コピービットをクリアし、コピー数が「> 1」であるフレームポインタをマルチコピーキュー 90 の書込側 92 に入れる。

バッファマネージャ 72 がフレームポインタおよびクリアされた単一コピービットを出力キュー 74 の 1 つから読出すたびに、バッファマネージャ 72 はフレームの送信をスケジューリングするが、コピー数「1」のフレームポインタを有するエントリがあるかどうかマルチコピーキャッシュ 96 をも調べる。コピー数「1」のフレームポインタがマルチコピーキャッシュ 96 に見つければ、バッファマネージャ 72 は、フレームの単一コピー送信の場合と同様に、送信のためにフレームをスケジューリングし、送信の間にバッファを再要求する。しかしながら、フレームポインタがマルチコピーキャッシュ 96 にないか、マルチコピーキャッシュ 96 におけるフレームポインタのコピー数が 1 よりも大きければ、バッファマネージャ 72 はフレームを送信するがバッファを再要求しない。送信を成功させた後、バッファマネージャ 72 はフレームポインタのコピーをコピー数「- 1」とともにマルチコピーキュー 90 の書込側 92 に入れる。

40

マルチコピーフレームが送信されるたびに、バッファマネージャ 72 はマルチコピーキャッシュ 96 内にコピー数「1」のフレームポインタを見つけられなかったならば、フレ

50

ムポインタのコピーをマルチコピーキュー 90 に入れる。したがって、いかなる所与の時間でも、マルチコピーキュー 90 はコピー数が「1」よりも大きいフレームポインタ、および/または、各々コピー数が「-1」である、同じフレームポインタのいくつかのコピーを含むことができる。

バッファマネージャ 72 は不使用となったバッファを再要求するためにマルチコピーキュー 90 およびマルチコピーキャッシュ 96 を絶えず処理する。バッファマネージャ 72 はマルチコピーキュー 90 を処理し、コピー数「>1」のフレームポインタを読出すと、この新しいエントリ（フレームポインタおよびコピー数）をマルチコピーキャッシュ 96 に入れようと試みる。マルチコピーキャッシュ 96 がフルであれば、バッファマネージャ 72 はその新しいフレームポインタのためにスペースを設ける。バッファマネージャ 72 は「より古い」マルチコピーキャッシュエントリを読出し、外部メモリ 36 内のそのバッファヘッダ内のこのエントリに対するコピー数を更新し、このエントリをマルチコピーキャッシュ 96 からクリアする。マルチコピーキャッシュ 96 内に使用可能な空きができると、バッファマネージャ 72 はマルチコピーキュー 90 からの新しいエントリをマルチコピーキャッシュ 96 に入れることができる。

バッファマネージャ 72 がマルチコピーキュー 90 を処理し、コピー数「-1」のフレームポインタを読出すと、それはマルチコピーキャッシュ 96 を探索して、デクリメントまたはデリートするためにコピー数「1」の対応するフレームポインタアドレスを探す。バッファマネージャ 72 がフレームポインタの一致を見つければ、それは 1) コピー数が「>1」であればマルチコピーキャッシュのフレームポインタをデクリメントするし、または 2) コピー数が「1」であればマルチコピーキャッシュのフレームポインタ/コピー数エントリをデリートし、フレームポインタをリクレームキュー 98 に入れる。

一致するフレームポインタが見つからなければ、バッファマネージャ 72 はコピー数を求めて外部メモリ 36（図 9 参照）におけるフレームポインタのバッファヘッダを探索する。メモリ内のコピー数が「1」であれば、バッファマネージャ 72 はフレームポインタをリクレームキュー 98 に入れる。メモリ内のコピー数が「>1」であれば、バッファマネージャ 72 はこのコピー数のフレームポインタをマルチコピーキャッシュ 96 に入れ、そのコピー数をデクリメントする。

バッファマネージャ 72 は、フレームポインタを読出してから、リンクトリストチェーンをたどり、バッファをフリーバッファプール 104 に戻すことによって、リクレームキュー 98 を絶えず処理する。この作用は、ヌルポートベクタを有し、ポートベクタ FIF0 70 によってリクレームキューに入れられていたフレームか、マルチコピー転送ベクタを有し、全コピーの送信を完了したフレームかのためのバッファを戻すのみである。単一コピーフレームにリンクされたバッファは、上述のようにそのフレームが送信されるときにフリーバッファプール 104 へと直接戻される。

出力キュー 74 と外部メモリ 36 内のそのオーバーフロー領域 110 とがフルであるために、ポートベクタ FIF0 70 が単一コピー転送ベクタのためのフレームポインタを出力キュー 74 に入れることができなければ、そのフレームは廃棄される。フレームポインタはリクレームキュー 98 に戻され、フレームの廃棄がスイッチの管理資源によって記録される。1 つ以上の出力キュー 74 と外部メモリ 36 内のそれらのオーバーフロー領域 110 とがフルであるために、ポートベクタ FIF0 70 がマルチコピー転送ベクタのための 1 つ以上のフレームポインタを入れることができなければ、そのフレームは使用可能なスペースのある出力キューへと転送されるのみであり、マルチコピーキュー 90 に入れられたコピー数はうまく入れられたフレームポインタを反映するのみである。フレームポインタが入れられなかったことは、フレームポインタがキューに入れられなかった各ポートごとにスイッチ管理資源によって記録される。全出力キュー 74 と外部メモリ 36 内のそれらのオーバーフロー領域 110 とがフルであるためにポートベクタ FIF0 70 がマルチコピー転送ベクタのためのどのフレームポインタも入れることができなければ、そのフレームポインタはリクレームキュー 98 に渡され、スイッチ管理資源にはそれに従い通知される。

10

20

30

40

50

マルチコピーキュー 90 は、フレームをストアするために用いられる全バッファ（すなわち、アドレスポインタ）がフリーバッファプール 104 に戻され得る前に、特定のマルチコピーフレームの何回の送信が完了されなければならないかを追跡するためにバッファマネージャ 72 が用いる高優先順位キューである。この出力キューの書込側 92 および読出側 94 はそれぞれ 64 エントリおよび 16 エントリを保持する。マルチコピーキュー 90 はマルチコピーキャッシュ 96 に入力を与え、マルチコピーキャッシュ 96 はいつバッファを再要求するかを決定するためにバッファマネージャ 72 によって用いられる。マルチコピーキューの内部構造を図 12 に示す。

出力キュー 74 にうまく入れることができたフレームポインタの数に基づいて、ポートベクタ FIF070 はフレームのフレームポインタコピーと「> 1」であるコピー数とをマルチコピーキュー 90 に入れる。特定のポートの出力キュー 74 がフルであれば、ポートベクタ FIF070 はフレームポインタのコピーをその出力キュー 74 に入れることができず、したがって、これをコピー数を決定する際の成功した事象として含めることはできない。

バッファマネージャ 72 が出力キューフレームポインタを読出し、単一コピービットが「0」である（すなわち、マルチコピー）ことを見つけるたびに、それは、これが最後の送信であることを示すコピー数「1」のフレームポインタを求めてマルチコピーキャッシュを調べる。この一致が見つからなければ、各バッファの内容が送信された後に不使用になったバッファをフリーバッファプール 104 に与えることによって、バッファマネージャ 72 は単一コピー送信の場合と同様にフレームを送信し、バッファを再要求する。一致が見つかれば、バッファマネージャ 72 はマルチコピーフレームを送信し、コピー数「- 1」のフレームポインタのコピーをマルチコピーキュー 90 に入れる。拡張バス出力キュー 75 または管理ポート出力キュー 77 へとキューに入れられたフレームのためのマルチコピーフレームポインタの（PCI インターフェイス 26 を介しての）使用をホストが終了すると、ホストはコピー数「- 1」のフレームポインタのコピーをフレームポインタレジスタを介してマルチコピーキューへと書込む。このレジスタは図 2 におけるレジスタ 60 のブロックに示されるレジスタの 1 つである。

出力キュー 74 と同様に、マルチコピーキュー 90 も入力経路および出力経路を備えて構成される。入力経路または書込側により、ポートベクタ FIF070 およびバッファマネージャはフレームポインタ / コピー数をマルチコピーキュー 90 に入れることができる。出力経路または読出側により、マルチコピーキュー 90 はフレームポインタ / コピー数をマルチコピーキャッシュ 96 に入れることができる。マルチコピーキューオーバーフロー 124 と呼ばれる、フレームポインタ / コピー数のためのさらなるストレージが外部メモリ 36 に設けられる。

フレームポインタ / コピー数が空のマルチコピーキュー 90 に書込まれると、それらは読出側 94 がフルになるまで書込側 92 から読出側 94 へと移動する。マルチコピーキュー 90 の書込側 92 に書込まれるさらなるフレームポインタ / コピー数は外部メモリ 36 内のマルチコピーキューオーバーフロー領域 124 に入れられる。一旦マルチコピーキュー 90 の読出側 94 とそのオーバーフロー領域 124 とがフルになれば、マルチコピーキューに入れられるさらなるフレームポインタ / コピー数が書込側 92 を満たし始める。

マルチコピーキュー 90 を通過するフレームポインタの順序は、マルチコピーキューの読出側 94 のスペースがクリアされると、フレームポインタ / コピー数がマルチコピーキューオーバーフロー領域 124 からマルチコピーキューの読出側 94 へと移動し、マルチコピーキューの書込側 92 からマルチコピーキューオーバーフロー領域 124 へと移動するようにして維持される。

マルチコピーキャッシュ 96 はマルチコピーキュー 90 と同様であるが、フレームポインタ / コピー数をスキャンするための探索可能な領域を設ける。マルチコピーキャッシュ 96 は 256 までのエントリを保持する。バッファマネージャ 72 はマルチコピーキュー 90 からフレームポインタを読出し、コピー数が「> 1」または「- 1」のいずれであるかによって、フレームポインタをマルチコピーキャッシュ 96 に入れるかそれ进行处理するか

10

20

30

40

50

する。

さらに、バッファマネージャ 7 2 が出力キュー 7 4 の読出側からフレームポインタを読出すごとに、バッファマネージャ 7 2 は送信をスケジュールする。単一コピービットが「0」である（マルチコピーフレームを意味する）ならば、バッファマネージャ 7 2 は、このフレームの最後の送信であることを示すコピー数「1」のフレームポインタを求めてマルチコピーキャッシュ 9 6 をスキャンする。一致があれば、バッファマネージャ 7 2 はフレーム送信の間にエントリを除去し、バッファをフリーバッファプールに戻す。一致がなければ、バッファマネージャは送信の終了時にコピー数「-1」のフレームポインタをマルチコピーキュー 9 0 に入れる。

バッファマネージャ 7 2 は周期的に、フレームポインタ / コピー数を読出し、それをマルチコピーキャッシュ 9 6 に入れるか処理することによってマルチコピーキュー 9 0 を処理する。これはフレーム送信から独立して行なわれる。バッファマネージャがコピー数「>1」のフレームポインタを読出すか、コピー数「-1」のフレームポインタを読出すかによって2つの場合が引き続いて生じる。

1) バッファマネージャ 7 2 がマルチコピーキュー 9 0 からコピー数「>1」のフレームポインタを読出す。マルチコピーキャッシュ 9 6 に空きがあれば、それは新しいエントリを書込む。マルチコピーキャッシュ 9 6 がフルであれば、バッファマネージャ 7 2 はキャッシュ 9 6 内のスペースをクリアしなければならない。これが行われるのは、マルチコピーキャッシュ 9 6 からより古いフレームポインタ / コピー数の1つを読出し、外部メモリ 3 6 内のフレームポインタのバッファヘッダをマルチコピーキャッシュ 9 6 内のコピー数で更新し、このキャッシュエントリをデリートすることによってである。一旦スペースが生じると、新しいフレームポインタ / コピー数がマルチコピーキャッシュ 9 6 に書込まれる。

2) バッファマネージャ 7 2 がマルチコピーキャッシュ 9 0 からコピー数「-1」のフレームポインタを読出す。バッファマネージャ 7 2 はコピー数「1」の一致するフレームポインタを求めてマルチコピーキャッシュ 9 6 を探索する。バッファマネージャ 7 2 がマルチコピーキャッシュ 9 6 内でフレームポインタの一致を見つけられるかどうかによって2つの場合が続く。

a) バッファマネージャ 7 2 がフレームポインタの一致を見つける。マルチコピーキャッシュ 9 6 のエントリのコピー数が「1」であれば、バッファマネージャ 7 2 はマルチコピーキャッシュエントリをデリートし、フレームポインタをリクレームキュー 9 8 に入れる。キャッシュエントリのコピー数が「>1」であれば、バッファマネージャ 7 2 はコピー数を「1」だけデクリメントする。

b) バッファマネージャ 7 2 がマルチコピーキャッシュ 9 6 内でフレームポインタの一致を見つけられない。これは、一致するフレームポインタが外部メモリ 3 6 内のフレームのリンクトリストチェーンのバッファヘッダに既に移動されていることを意味する。バッファマネージャ 7 2 はバッファヘッダに行って、コピー数を読出さなければならない。（メモリ内の）この値が「1」であれば、フレームはもはや必要ではなく、バッファマネージャ 7 2 はフレームポインタをリクレームキュー 9 8 に入れる。（メモリ内の）この値が「>1」であれば、バッファマネージャ 7 2 は（外部メモリ 3 6 内にあった）フレームポインタ / コピー数のコピーをマルチコピーキャッシュ 9 6 に入れ、コピー数を「1」だけデクリメントする。マルチコピーキャッシュ 9 6 がフルであれば、バッファマネージャはより古いフレームポインタ / コピー数の1つを外部メモリ 3 6 に移動させることによってスペースをクリアする。

リクレームキュー 9 8 はもはや必要とされないリンクトリストチェーンを指すフレームポインタを保持する。バッファマネージャ 7 2 は、マルチコピーキャッシュを処理してフレームポインタのコピー数が「1」である（すなわち、フレームの最後の送信がうまく終わった）ことを見出すと、フレームポインタのリクレームキューに書込む。さらに、ポートベクタ F I F O 7 0 は、1) フレームポインタのポートベクタがヌルであるか、2) 転送ベクタの全出力キューがフルであったのでフレームポインタがキューに入れられることがで

10

20

30

40

50

きなかったという条件下で、フレームポインタをリクレームキュー 98 に書込む。最後に、ホストは、拡張バス出力キュー 77 または管理ポート出力キュー 75 に対してキューに入れられた単一コピーフレームの使用を終えると、(フレームポインタレジスタを用いて)フレームポインタをリクレームキュー 98 に書込む。

バッファマネージャ 72 はリクレームキューのエントリを処理するとき、フレームポインタのリンクトリストチェーンをたどり、各バッファをフリーバッファプール 104 に戻す。リクレームキュー構造の内部構造は図示されないが、本発明の例示的实施例においてはフレームポインタ(14ビット)のみを含む。リクレームキューの書込側 100 は 64 エントリを保持し、リクレームキューの書込側 102 は 16 エントリを保持する。

出力キュー 74 と同様に、リクレームキュー 98 は入力経路および出力経路を備えて構成される。入力経路または書込側 100 によってバッファマネージャ 72 はフレームポインタをリクレームキュー 98 に入れることができる。出力経路または読出側 102 によってバッファマネージャ 72 はフレームポインタを読出し、関連の全バッファをフリーバッファプール 104 に戻すことができる。フレームポインタのためのさらなるストレージは外部メモリ 36 内に設けられるリクレームキューオーバーフロー領域 122 内に設けられる。

フレームポインタが空のリクレームキュー 98 に書込まれると、これらは読出側 102 がフルになるまで書込側 100 から読出側 102 へと移動する。リクレームキュー 98 の書込側 100 に書込まれるさらなるフレームポインタは外部メモリ 36 内のリクレームキューオーバーフロー領域 122 に入れられる。一旦リクレームキュー 98 の読出側 102 およびオーバーフロー領域 122 がフルになると、リクレームキュー 98 に入れられるさらなるフレームポインタが書込側 100 を満たし始める。

図 11 はフリーバッファプール 104 の内部構造の例示的实施例を示す。フリーバッファプール 104 は、外部メモリ 36 内の全フリーバッファ 140 を指すアドレスポインタを含んだ FIFO である。フレームが受信されると、バッファマネージャ 72 は入来するデータをストアするためにフリーバッファプール 104 から使用可能なアドレスポインタを捕捉する。バッファマネージャ 72 はまたフリーバッファプール 104 からのアドレスポインタを(要求される場合)ホストプロセッサ 28 に割当てて、ホストは、直接入力/出力スペースにおけるレジスタ 60 の中のフリーバッファプールレジスタを読出すか書込むことによってアドレスポインタを要求するかそれらをフリーバッファプール 104 に戻すことができる。フリーバッファプール 104 の書込側 106 および読出側 108 は本発明の例示的实施例においては各々 64 エントリを保持する。

フリーバッファプール 104 は(出力キュー 74 と同様に)入力経路および出力経路を備えて構成される。入力経路または書込側 106 により、バッファマネージャ 72 またはホスト 28 はアドレスポインタをフリーバッファプール 104 へと入れることができる。フリーバッファプール 104 の出力経路または読出側 108 により、バッファマネージャ 72 はアドレスポインタをホスト 28 に与え、またはプール 104 からアドレスポインタを引出して受信フレームデータをストアすることができる。使用可能なアドレスポインタのさらなるストレージ、フリーバッファプールのオーバーフロー領域 120 は上述のように外部メモリ 36 内に設けられる。

スイッチ 12 が起動すると、フリーバッファプールは読出側 108 からアドレスポインタを発生する。フレームが入来するときにフリーバッファプール 104 内のフリーリストが読出される。書込側 106 にトラフィック要求を扱うのに十分なバッファポインタがなければ、オーバーフロー領域 120 がより多くのバッファポインタを得るためにアクセスされる。

本発明のある実施例は、スイッチ 12 が開始されるとバッファポインタを与える有利な配置および方法を提供する。スイッチ 12 が初めに電源投入されるとき、外部メモリ 36 内のオーバーフロー領域 120 がバッファポインタを含むことは必要とされない。代わりに、バッファポインタはオンザフライで発生される。スイッチ 12 は電源投入されるとバッファポインタを発生し、それをオーバーフロー領域 120 に入れることができるが、この

10

20

30

40

50

ようなポインタは16,000個または32,000個存在することがあり、これによってスイッチ12の電源投入手順が遅くなるであろう。本発明は、電源投入時に全バッファがフリーであり、これらのバッファのアイデンティティが既知であるという事実を利用する。したがって、バッファポインタは電源投入後に必要とされるときに図10に示されるようにカウンタ180を用いて発生される。

フリーリストカウンタ発生器180がマルチプレクサ182の入力に接続される。フリーバッファプール104のフリーリストが開始時に空であるので、フリーリストカウンタ180はバッファポインタを発生する。一旦フリーリストが最高カウントに達すると、それはこれ以上バッファポインタを発生しない。

フレームパケットがスイッチ12において受信されると、フレームパケットは固定長バッファへと分解する。典型的にフレームはさまざまなサイズである。バッファは256バイトのサイズであり、バッファのデータ部分は240バイトである。バッファ内容の送信後、バッファポインタがリクレームキュー98に入れられるか、または、バッファチェーンをたどることができるならばフリーバッファプール104のフリーリストに直接入れられる。スイッチ12の動作の間、フリーバッファプール104に戻されるどのアドレスポインタも書込側106から読出側108へと移動する。読出側108がフルとなれば、さらなるアドレスポインタはオーバーフロー領域120に渡される。一旦読出側108およびオーバーフロー領域120がフルとなると、フリーバッファプール104に入れられるさらなるアドレスポインタがプール104の書込側106を再び満たし始める。

図13は本発明の実施例に従うマルチコピーキャッシュ96の内部配列の概略図である。上で簡単に述べたように、マルチコピーキャッシュ96へのエントリの時間順が維持される。本発明では、このように時間順が維持されるのは先行技術におけるようなタイムスタンプによってではなく、メモリ内の物理的順序によってである。本発明のマルチコピーキャッシュ96はまた有効性ビットの使用を避け、代わりに後述するように有効性を符号化する。

図13を参照すると、マルチコピーキャッシュ96は4ウェイセットアソシアティブメモリとして構成される。マルチコピーキャッシュ96へのエントリは上述のようにフレームポインタとそのコピー数とを含む。フレームポインタの最下位6ビットが、エントリがストアされるセットアソシアティブキャッシュ96内の行を決定する。本発明の図示される実施例では、キャッシュ96には64行が存在するが、キャッシュサイズが大きくなれば他の行数も制限されない。

セットアソシアティブキャッシュ96は4列に分割され、その各々が並行して探索される。バッファマネージャ72がエントリをキャッシュ96へとストアするとき、エントリは常に、第1の列の、フレームポインタの最下位6ビットによって示される行の最上位(51:39)ビットに入る。この行は読出され、全エントリが13ビット分右にシフトされ、行は再び書込まれる。実際にキャッシュ96に書込まれるエントリはフレームポインタの上位8ビットを含み、それはアドレスタグとフレームポインタに関連した5ビットコピー数を形成する。エントリがキャッシュ96から読出されると、フレームポインタはキャッシュ96の行数を指すビットおよびアドレスタグで再形成される。

行がフルであり、その行への新たなエントリが書込まれれば、キャッシュ96内の最も古いエントリがキャッシュ96から除去される。バッファヘッダ142に関して上述したように、除去されるフレームポインタに関連したコピー数は除去されるフレームポインタが指す外部メモリ内のフレームのバッファヘッダ142に書込まれる。したがって、外部メモリ36にストアされるフレーム(すなわち、バッファ140)はコピー数をストアするためのマルチコピーキャッシュ96のためのオーバーフロー領域となる。

本発明の有利な特徴の1つはセットアソシアティブキャッシュ96に別個の有効ビットが存在しないことである。コピー数が00000であるとき、エントリがもはや有効でないことをバッファマネージャ72はわかっており、エントリをキャッシュ96から除去する。これによってキャッシュ構成が簡素化される。本発明のキャッシュ96の別の利点は非常に高速な探索が行なわれ得ることである。これは、バッファマネージャ72がマルチコ

10

20

30

40

50

ピーキュー 90 を出たフレームポインタによって既に定められている単一の行を検査しさえすればよいためである。その行内の 4 つのエントリが並行して検査され、探索速度をさらに高める。4 ウェイセットアソシアティブメモリとして説明しているが、これは例にすぎず、メモリは本発明の範疇から逸脱せずに n ウェイセットアソシアティブ方式となり得る。

上の説明から、本発明がキャッシュにおけるエントリの行ごとの物理的位置決めによってキャッシュエントリの時間順（エージ）を維持すると理解されるべきである。すなわち、キャッシュ内のエントリの物理的位置がエントリの相対的エージを示す。エントリはメモリにおけるエントリの物理的再順序付けによってエージングされる。

本発明のある実施例はポートごとにスイッチ 12 によって切換えられるフレームのレイテンシをカスタマイズする。図 14 を参照すると、ポートベクタ F I F O 70 が受信ポートのプログラムされたスイッチモードを検査して、いつフレームポインタおよび関連の情報を送信ポートの適切な出力キュー 74 へと入れるかを決定する。第 1 のモード（低レイテンシモード）では、ポートベクタ F I F O 70 はいつフレームポインタを出力キュー 74 に入れるかに対して制限を与えない。第 2 のモード（中間レイテンシモード）では、ポートベクタ F I F O 70 はフレームの 64 バイトが受信されて初めてフレームポインタを出力キュー 74 に入れる。第 3 のモード（高レイテンシモード）では、ポートベクタ F I F O 70 はフレームが完全に受信されて初めてフレームポインタを出力キュー 70 に入れる。

いつポートベクタ F I F O 70 がフレームポインタを出力キュー 74 へと移動するかのタイミングを変えるいくつかの特殊な場合があり、それらは、1) 第 1 または第 2 のモードの 10 Mb/s ポートから 100 Mb/s ポートへのフレーム転送と、2) 管理ポート 30 へのフレーム転送と、3) 拡張バスポートへのフレーム転送とを含む。場合 1) では、10 Mb/s ポートから 100 Mb/s ポートへの速度不一致によって転送モードが強制的に第 3 の高レイテンシモードとされる。場合 2) では、管理ポートへと移動する全フレームが第 3 のモードのフレームである。場合 3) では、拡張バスポートへのどのフレーム転送も拡張バスポート 26 のスイッチモードを用いる。マルチコピーポートベクタが特殊な場合のポートの 1 つを含む場合、ポートベクタ全体に対するフレームポインタのキュー入れはポートベクタ内で表わされる最長レイテンシスイッチモードのそれになる。たとえば、フレームが第 1 または第 2 のモードのポートによって受信され、そのマルチコピー転送ポートベクタが管理ポート 30 を含めば、スイッチモードは第 3 のモードである。この場合、フレームが完全に受信されて初めてフレームポインタのコピーが全出力キュー 74 に入れられる。

スイッチモードをここでより詳細に説明する。入力（すなわち、受信）ポートに当てはまるスイッチモードが転送レイテンシ（一旦スイッチ 12 がフレームを受信し始めるとどの程度後にスイッチ 12 がフレームを転送するか）と出力ポートへのフラグメント/エラー伝搬を低減する能力とを決定する。第 2 の中間レイテンシモードは各ポートに対するデフォルトであるが、スイッチモードはレジスタ 60 ではポートごとにプログラム可能である。

これら 3 つのモデルのすべてにおいて、内部 MAC ポートの受信 F I F O 52 で受信されるフレームデータはできるだけ早く外部メモリ 52 内のバッファ 140 に転送される。ほぼ同時に、ルールチェッカ 42 または 58 が宛先アドレスおよびソースアドレス、受信ポート数、フレームポインタ、ならびにいくつかの付加的情報を受信し、適切なルックアップを行なう。一旦ルックアップが完了すると、ルールチェッカ 42 または 58 はフレームポインタおよび転送ポートベクタをポートベクタ F I F O 70 に戻す。

ポートベクタ F I F O はポートベクタ内で識別される出力ポートのための出力キュー 74 の書込側 76 にフレームポインタを入れる。受信ポートのスイッチモードは、ポートベクタ F I F O 70 がポートベクタ（およびフレームポインタ）を受取るときから、それがフレームポインタを出力キュー 74 に入れるときまでの間のレイテンシを規定する。これは以下の 3 つのモードに対して説明される。一旦フレームポインタが出力キュー 74 の読出

10

20

30

40

50

側 7 8 に移動すると、バッファマネージャ 7 2 はフレームポインタを讀出し、送信をスケジュールする。バッファマネージャはフレームポインタによって特定されるアドレスからフレームデータを移動させ始める。一旦 M A C ポートの送信 F I F O 5 4 がその開始点に設定されると（そして、データ送信のために媒体が使用可能であると想定すると）、フレーム送信が始まる。

第 1 のモードは最低のレイテンシを与えるように設計される。フレームはライン・レート速度で受信され、転送される。この第 1 のモードにおいてはネットワークエラーに対する保護がなく、これは、フレームがフラグメント（すなわち、 $< 64$  バイトの長さ）であるか C R C エラーを含むかが判断され得る前にフレームが送信のためにキューに入れられるためである。第 1 のモードにおいて、フレーム受信は出力ポートでのフレーム送信が始まるまでに完了していないかもしれない。受信フレームが短すぎる場合または無効な C R C で終る場合、受信 M A C は外部メモリ 3 6 内のバッファヘッダ 1 4 2 に印を付けてこれらの条件を示す。送信 M A C は、後に短すぎるものか無効な C R C で終るフレームの送信が始まれば M A C が無効な C R C を発生することを保証する。送信 M A C がフレーム送信を始めておらず、バッファヘッダ 1 4 2 が短すぎるものか無効な C R C で終るフレームを示している場合、バッファマネージャ 7 2 はフレームを出力ポートへと転送しない。

第 2 のモードはフレームを転送するための低レイテンシとあるネットワークエラーに対する保護とを与える。フレームは 64 バイト以上が受信された後に受信され、転送される。これによってスイッチ 1 2 がフレームのフラグメントをフィルタ処理する（すなわち、転送しない）ことが可能となるが、これは 64 バイトよりも大きい C R C エラーフレームを完全にはフィルタ処理しない。

第 2 のモードにおいては、受信 M A C で 64 バイトのしきい値を達成したフレームのフレームポインタは適切な出力キュー 7 4 に入れられる。最小の 64 バイトのしきい値を達成できないフレームはデリートされ、それらのフレームポインタは出力キュー 7 4 に入れられない。64 バイト以上の受信フレームが無効な C R C で終れば、受信 M A C は外部メモリ 3 6 内のバッファヘッダ 1 4 2 に印を付けてこの条件を示す。後に無効な C R C で終る 64 バイト以上のフレームの送信が開始されるときには、送信 M A C は不良な C R C で送信を終了する。送信 M A C がフレーム送信を開始しておらず、バッファヘッダ 1 4 2 が無効な C R C で終るフレーム（64 ビット以上）であることを示している場合、バッファマネージャはフレームポインタを（単一コピー転送のための）リクレーンキュー 9 8 または（マルチコピー転送のための）マルチコピーキュー 9 6 へと出力ポート 7 4 への転送なしに戻す。

第 3 のモードは 3 つのモードの中で最高レベルのネットワークエラー保護を与えるがより高い転送レイテンシを有するストアアンドフォワードモードである。フレームは、スイッチ 1 2 がそれらを出力ポートに転送する前に完全に受信される。このモードでは、スイッチ 1 2 は転送の前に全てのフラグメントおよび C R C エラーフレームをふるい分ける。第 3 のモードにおいて、一旦有効フレームが受信側でうまく完了すると（すなわち、有効な C R C を持ち、64 バイト以上であると）、フレームポインタが適切な出力キュー 7 4 に入れられる。受信エラー（無効 C R C、短すぎるもの（ $> 64$  バイト）等）で終るフレームはデリートされ、それらのフレームポインタは出力キュー 7 4 に入れられない。

ポートベクタ F I F O 7 0 は、受信ポートの選択されたモードと受信されたデータ量とに依存してポートベクタを出力キュー 7 4 に入れる決定を行なう。上述の実施例では、3 つのしきい値があるが他の実施例では異なる数のしきい値が存在する。例示的实施例では、これらのしきい値は 1 )  $n < 64$  バイトであるような  $n$  バイト（たとえば 6 バイト）の受信、2 ) 64 バイトの受信、および 3 ) 全フレームの受信である。

本発明はしきい値に基づいてフレームを出力キュー 7 4 へと転送する。ポートベクタ F I F O 7 0 は、受信されるデータタイプの量とポートがプログラムされたモードとに基づいて送信シーケンスを再び順序付ける。例示的实施例は受信されたデータの量に基づいて転送の決定を行なうが、本発明の他の実施例では、受信されるデータタイプのような他の要因に基づいて転送の決定が行われる。

10

20

30

40

50



本発明の転送方式を実施するにあたって、バッファマネージャ 72 はフレームポインタを受信ポートと関連付ける、キャッシュメモリ (CAM) 161 内のテーブル 160 を維持する。ポートベクタ F I F O 70 が新しいポートベクタおよびフレームポインタをルールチェッカ 42 または 58 から受信するたびに、それは関連付けを行なって受信ポートがフレーム受信を終えたかどうかを判断し、終えていなければどれほどのフレームが既に受信されているかを判断する。ポートベクタ F I F O 70 が受信ポートのアイデンティティに関する情報をルールチェッカ 42 または 58 から受信することはない。ポートベクタが受取る唯一のポートの何らかの識別を与える情報はフレームポインタである。

ポートベクタ F I F O 70 はフレームポインタでアドレステーブル 160 に問合せをする。フレームがなお受信されていればアドレステーブルは受信ポートを戻し、またはアドレステーブル 160 はフレームポインタを見つけることができないときはフレームが既に受信されたことを意味する。一旦フレームが完全に受信されると、フレームポインタがアドレステーブル 160 から移動される。これは、第 3 のしきい値 (フレーム完了) が満たされたことを意味する。したがって、フレームポインタは直ちに出力キュー 74 に入れられ得る。

アドレステーブル 160 が受信ポートを戻せば、ポートベクタ F I F O 70 がフレームポインタおよび関連の情報を保持領域 162 に入れ、その受信ポートからの 2 信号を監視し始める。これらの 2 信号は 3 つの事象のうちの 1 つを示す。第 1 の事象はポートが n バイトを受信するときに示される。その時点で、そのポートが第 1 のモードにあれば、ポートベクタ F I F O 70 がフレームポインタを適切な出力キュー 74 に送ることによってその処理を開始する。受信ポートが第 1 のモードになれば、ポートベクタ F I F O 70 は第 2 の事象の発生を示す信号が受信されるまで待機する。このポートが第 2 のモードにあれば、ポートベクタ F I F O 70 はフレームポインタを保持領域 162 から解放し、適切な出力キュー 74 に入れる。最後に、受信ポートが第 3 のモードにあれば、ポートベクタ F I F O 70 はフレームが完全であることを示すフラグの受信を待つ。各受信ポート (図 14 の参照番号 164) がこのフラグを維持し、この情報をポートベクタ F I F O 70 に提供する。フレームポインタに関連付けられたポートの決定はポートベクタ F I F O 70 次第である。ポートベクタ F I F O 70 は各ポートのモードを識別する情報を維持する。要約すると、フレームポインタが受信されると、ポートベクタ F I F O 70 は最初にバッファマネージャ 72 のアドレステーブル 160 に問合せをして受信ポートを決定し、その受信ポートのためのモードを決定し、受信ポートからのフラグを監視し、モードおよびフラグに従ってフレームポインタを解放する。

本発明が詳細に説明され、図示されたが、これは図示および例示のためのものにすぎず、限定するものとは受取られるべきでなく、本発明の精神および範疇が請求の範囲によってのみ限定されることが明らかに理解される。

10

20

30

【図 1】

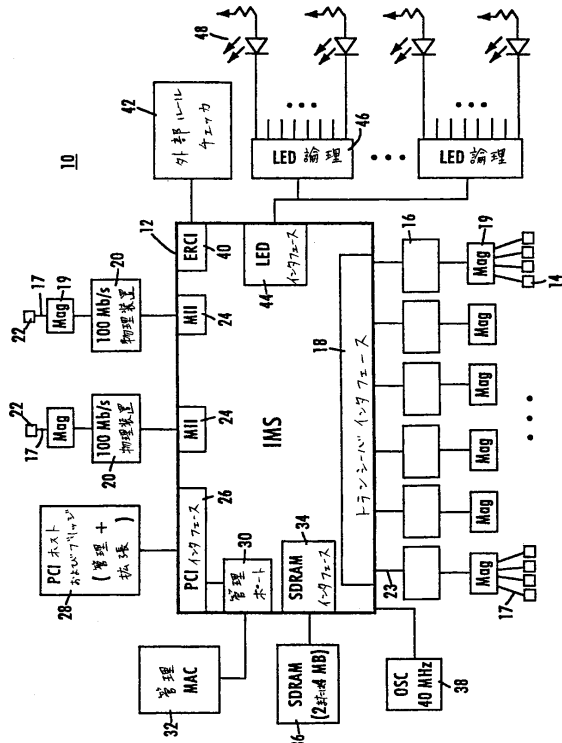
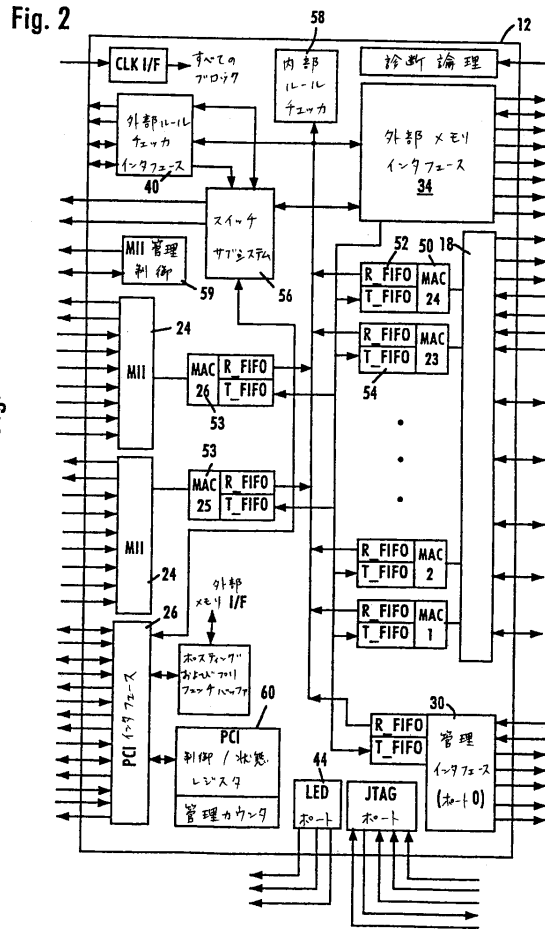


Fig. 1

【図 2】



【図 3】

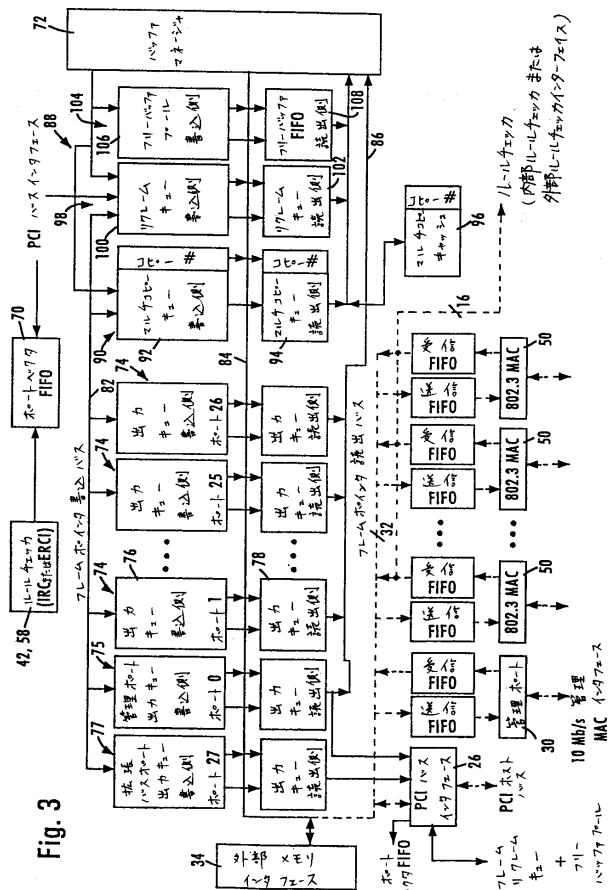


Fig. 3

【図 4】

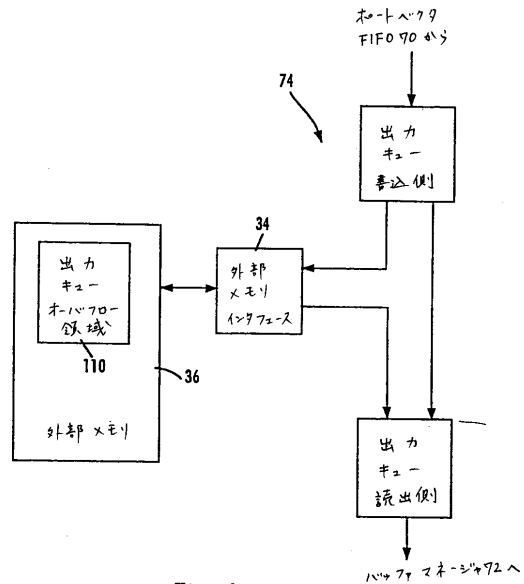
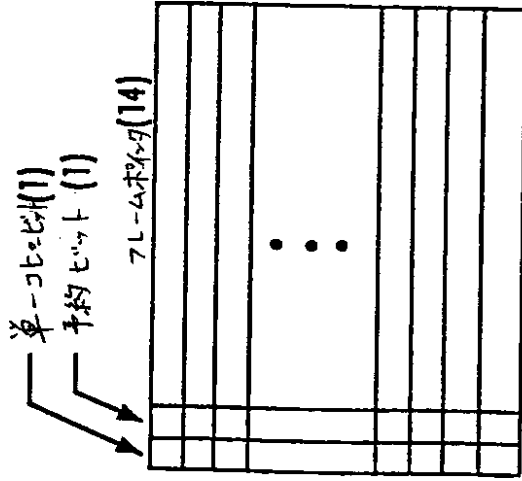


Fig. 4

【図 5】



【図 6】

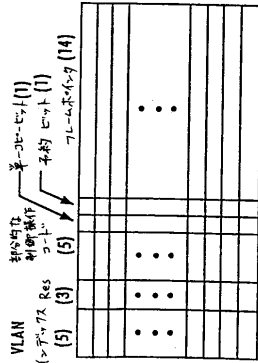


Fig. 6

【図 8】

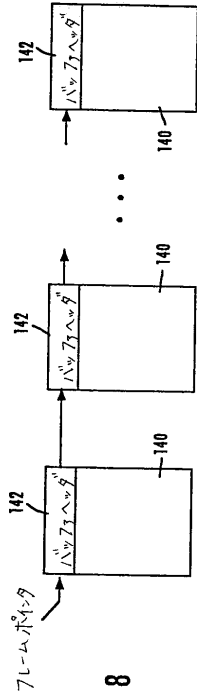


Fig. 8

【図 7】

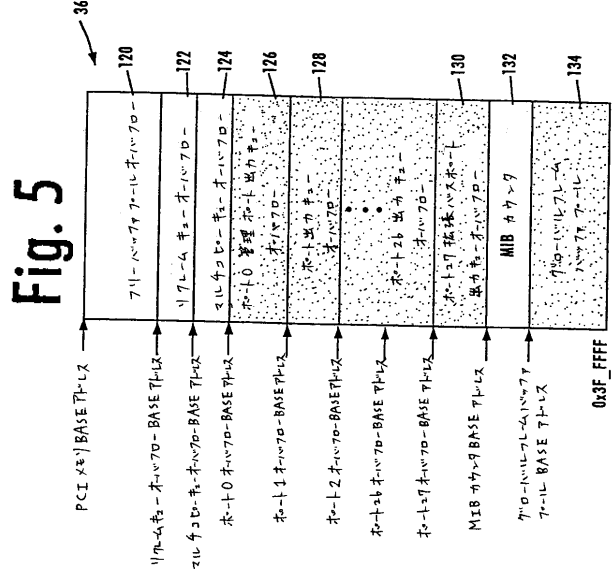


Fig. 5

Fig. 7

【図 9 a】

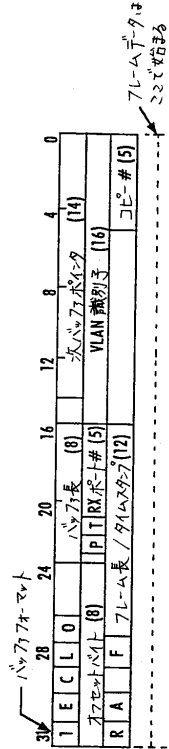


Fig. 9a

【図9b】

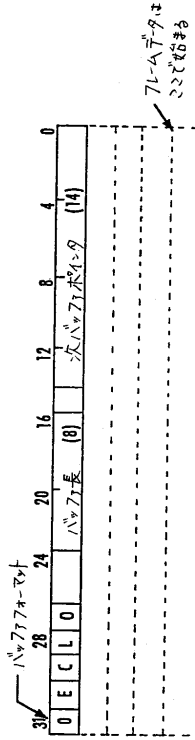


Fig. 9b

【図11】



Fig. 11

【図10】

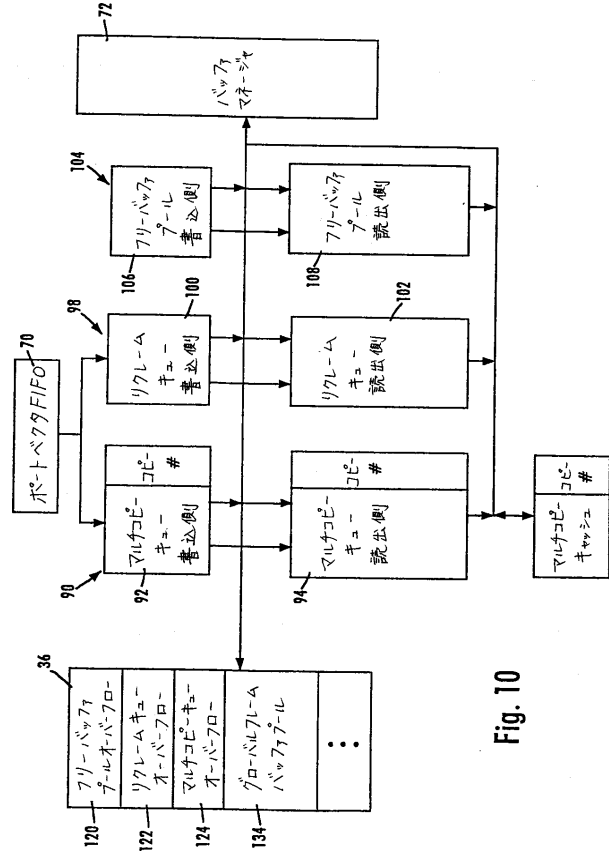


Fig. 10

【図12】

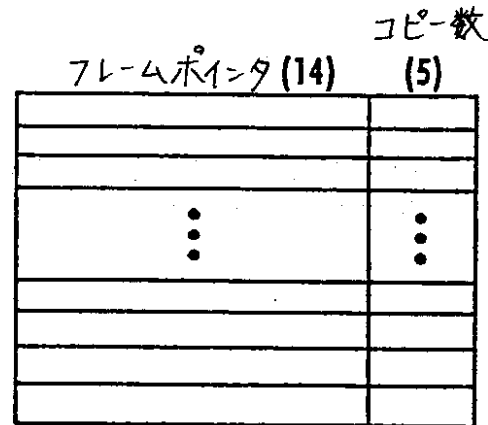


Fig. 12

【図 13】

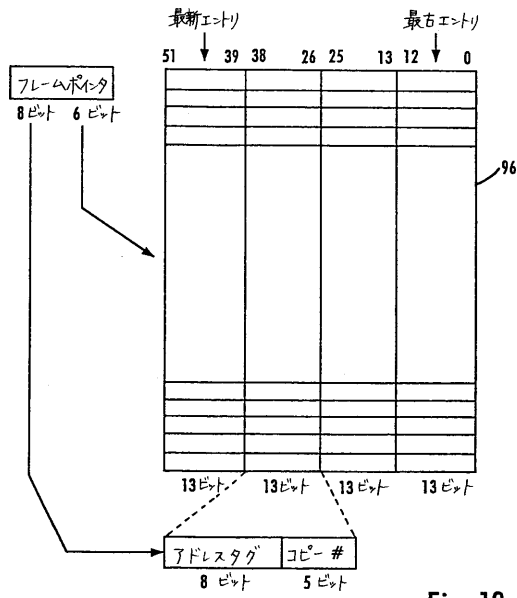


Fig. 13

【図 14】

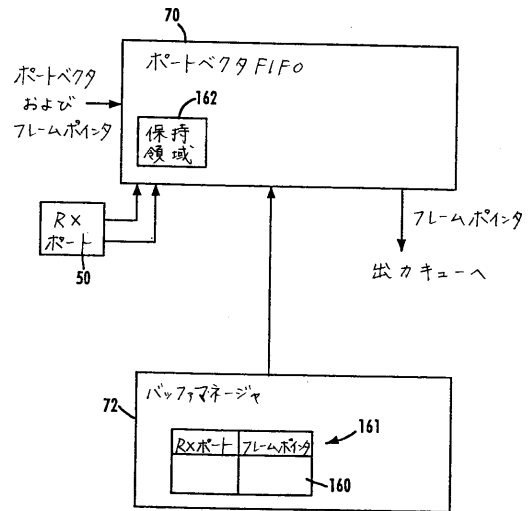


Fig. 14

---

フロントページの続き

## (74)代理人

弁理士 堀井 豊

## (72)発明者 エリムリ, バハディール

アメリカ合衆国、9 4 0 4 0 カリフォルニア州、マウンテン・ビュー、カリフォルニア・ストリート、2 1 0 1、ナンバー・1 0 9

## (72)発明者 ルナルデュー, トーマス・ジェイ

アメリカ合衆国、9 5 1 1 7 カリフォルニア州、サン・ノゼ、ブラックフォード・アベニュー、3 7 0 1

審査官 吉田 隆之

## (56)参考文献 特開平 8 - 8 9 0 6 ( J P , A )

特開平 4 - 1 7 5 0 3 4 ( J P , A )

特開平 6 - 3 3 8 8 9 9 ( J P , A )

特開平 6 - 3 3 4 6 5 2 ( J P , A )

## (58)調査した分野(Int.Cl. , D B 名)

H04L 12/56