

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4166285号

(P4166285)

(45) 発行日 平成20年10月15日(2008.10.15)

(24) 登録日 平成20年8月8日(2008.8.8)

(51) Int.Cl.

F I

H04L 12/56 (2006.01)

H04L 12/56

F

請求項の数 8 (全 28 頁)

(21) 出願番号	特願平10-535788	(73) 特許権者	591016172
(86) (22) 出願日	平成10年1月30日(1998.1.30)		アドバンスト・マイクロ・ディバイズ・
(65) 公表番号	特表2001-511979(P2001-511979A)		インコーポレイテッド
(43) 公表日	平成13年8月14日(2001.8.14)		ADVANCED MICRO DEVI
(86) 国際出願番号	PCT/US1998/001988		CES INCORPORATED
(87) 国際公開番号	W01998/036536		アメリカ合衆国、94088-3453
(87) 国際公開日	平成10年8月20日(1998.8.20)		カリフォルニア州、サニibel、ピー・
審査請求日	平成16年12月20日(2004.12.20)		オウ・ボックス・3453、ワン・エイ・
(31) 優先権主張番号	60/038,025		エム・ディ・プレイス、メイル・ストップ
(32) 優先日	平成9年2月14日(1997.2.14)		・68(番地なし)
(33) 優先権主張国	米国(US)	(74) 代理人	100064746
			弁理士 深見 久郎
		(74) 代理人	100085132
			弁理士 森田 俊雄

最終頁に続く

(54) 【発明の名称】 バッファを再要求するための方法および装置

(57) 【特許請求の範囲】

【請求項1】

パケット交換ネットワーク内においてネットワークスイッチのためにフレームをストアするバッファを再要求するための構成であって、

前記ネットワークスイッチからもはや送信される必要がないフレームがストアされた第1のバッファのバッファメモリにおける場所をポイントするフレームポインタをキューとして維持するように構成されたリクレーンキューと、

前記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数のカウント値をストアするよう構成された探索可能な第1のメモリと、

バッファメモリにおけるバッファの場所をポイントするバッファポインタをキューとして維持するように構成されたフリーバッファプールと、

フレームのフレームポインタを受信してそのフレームのコピーを送信し、このコピーの送信が、送信されるそのフレームの最後のコピーであるかどうかを判断するよう前記第1のメモリを探索し、そうであれば、前記ネットワークスイッチから前記そのフレームの最後のコピーを送信し、そのバッファの内容を送信した後に、前記フレームをストアする各バッファの前記バッファポインタを前記フリーバッファプールに送るよう構成されたバッファマネージャと、

複数のエントリをストアするよう構成されたマルチコピーキューとを含み、各エントリはフレームポインタおよび関連のコピー数を含み、正のコピー数は、前記関連のフレームポインタによってポイントされた前記フレームの送信されるコピーの数を示し、負のコピー

10

20

ー数は、前記関連のフレームポインタによってポイントされた前記フレームのコピーの送信を表わしており、

前記バッファマネージャはさらに、

前記フレームのコピーの送信が前記ネットワークスイッチによって送信されるそのフレームの最後のコピーではないと前記第1のメモリの探索により判断したときに、負のコピー数を有する前記エントリのうちの1つを前記マルチコピーキューにロードし、

前記マルチコピーキューを出るエントリのコピー数を検査し、正のコピー数を有する前記エントリを前記第1のメモリにストアし、負のコピー数を有するエントリについてはそのフレームポインタに等しいフレームポインタを有するエントリを前記第1のメモリ内において探索し、

10

前記第1のメモリを探索することにより、前記第1のメモリのエントリが、負のコピー数を有する前記エントリの前記フレームポインタに等しいフレームポインタを有することが示されたときに、前記第1のメモリの前記エントリの前記コピー数が1よりも大きい場合には、前記第1のメモリのエントリの前記コピー数をデクリメントするように構成される、構成。

【請求項2】

前記バッファマネージャは、前記第1のメモリを探索することにより、前記第1のメモリの前記エントリが、負のコピー数を有する前記エントリの前記フレームポインタに等しいフレームポインタを有することが示されたときに、前記第1のメモリの前記エントリの前記コピー数が1に等しい場合には、前記第1のメモリの前記エントリを削除し、前記第1

20

【請求項3】

前記バッファマネージャは、前記リクレーンキューから出るフレームポインタを得て、前記フレームポインタによってポイントされたフレームをストアする各バッファの前記バッファポインタを前記フリーバッファプールに送るようにさらに構成される、請求項2に記載の構成。

【請求項4】

単一のフレームをストアするバッファがリンクリストを形成するようにともに繋がれ、前記フレームの前記フレームポインタは前記リンクリストを形成する前記バッファのうちの第1のものの前記バッファポインタに等しく、前記リンクリストの各バッファは、そのフレームに関する前記リンクリスト内の次のバッファをポイントするバッファポインタを含み、前記リンクリスト内の最後のバッファは、前記次のバッファへの前記バッファポインタの代わりに、前記リンクリストの最後を示す表示を含む、請求項3に記載の構成。

30

【請求項5】

前記バッファマネージャは、前記リンクリスト内の前記バッファを検査し、前記フレームの前記リンクリスト内の前記バッファの前記バッファポインタを抽出し、前記フリーバッファプールにそれらを送るようにさらに構成される、請求項4に記載の構成。

【請求項6】

40

前記バッファマネージャによって受信された、送信のためのフレームの各フレームポインタは、前記フレームのコピーが1つ送信されるかまたは前記フレームのコピーが複数送信されるかを示すビットを有し、前記バッファマネージャは、前記ビットを検査し、前記フレームのコピーが複数送信されることを前記ビットが示す場合にのみ前記第1のメモリの探索を行ない、前記フレームの1つのコピーが送信されることを前記ビットが示す場合には、前記フレームの前記1つのコピーを送信し、そのバッファの内容を送信した後に、前記フレームをストアする各バッファの前記バッファポインタを前記フリーバッファプールに送るようにさらに構成される、請求項5に記載の構成。

【請求項7】

前記フリーバッファプールに結合され、データをストアするために使用できるようになっ

50

た前記バッファのバッファポインタを前記フリーバッファプールに戻すように構成された再要求回路と、

データのストアに使用できるフリーバッファを提供するよう、前記構成の起動時にバッファポインタを発生するカウンタと、

前記フリーバッファプールの出力に結合された第1の入力と、前記カウンタの出力に結合された第2の入力と、前記第1の入力と前記第2の入力との間で選択を行なう制御入力とを備えるマルチプレクサとをさらに含む、請求項1に記載の構成。

【請求項8】

前記構成の起動後に、データのストアにバッファが必要とされると、前記カウンタが最大カウント値に達するまで、前記マルチプレクサの前記第2の入力を選択し、その後、ストアのためにさらなるバッファが必要とされると、前記第1の入力を選択するための手段をさらに含む、請求項7に記載の構成。

【発明の詳細な説明】

発明の分野

この発明はデータストアの分野に関し、特に、データをストアするために用いられるバッファを再要求して、種々のデータをストアするためにバッファを再使用することができるようにするための方法および装置に関する。

背景技術

パケット交換ネットワーク（たとえばイーサネットネットワーク）などの受信および転送を行なう多くのシステムにおいては、データを転送する前に一時的にデータをストアするためにメモリバッファが用いられる。たとえば、パケット交換ネットワークでは、データのフレームはネットワークスイッチで受信され、このネットワークスイッチは適切なポートにフレームを転送し、フレームはデータがスイッチで受信されている間に個々のバッファにロードされる。フレームのサイズは通常、個々のバッファの容量よりも大きいため、単一のフレームのデータをストアするために複数のバッファが用いられる。

受信される都度データをバッファすることにより、種々の特徴の中でもとりわけ、頑健なエラー検査がフレームに対して行なわれるようになり、送信および受信ポート間でのレートを整合させることも可能になる。データのフレームがスイッチに到達すると、現在使用可能である（すなわち「フリーな」）バッファだけがデータのストアに用いられるため、データが送信される前にデータが上書きされることはない。データが送信済みか、またはさもなければ最早必要でなくなった場合、データをストアするバッファは新しいフレームのデータのストアのために再使用することができる。また、入来するフレームのデータをストアするためにバッファを確保しておくことは可能であるが、エラー検査により、または受信されたデータの量が不十分であることにより、フレームが破棄されてしまう。この状況においても、バッファは再使用のために戻されるべきである。

効率を高め、可能な限り入来フレームのストアに使用できるバッファを確保するために、再使用のためにバッファを戻すことは、それらが最早必要でなくなった後にできるだけ迅速に行なわれるべきである。フレーム全体が送信された後に複数のバッファが単一のフレームをストアできるようにするシステムでは、スイッチはフレームの送信を認識し、そのチェーンをストアしている個々のバッファを再要求するプロセスを開始する。同じフレームのコピーが複数送信される場合、スイッチは、フレームをストアしたバッファを再要求するプロセスを開始する前にそのフレームのすべてのコピーが送信されるまで待機する。したがって、フレームのコピーが1つ送信されようと、フレームのコピーが複数送信されようと、同じ再要求手順が用いられるであろう。

メモリにデータをストアするためにバッファを用いるシステムにおける別の考慮点は、これらのバッファをストアするためのポインタを発生することである。メモリ中のバッファがある場所をポイントするバッファポインタはシステムの起動時に発生され、バッファプールにロードされ、データをストアするためにバッファが用いられる際にそこから取出される。使用後にバッファを再要求すると、バッファポインタはバッファプールに戻される。システムが多数のバッファを使用する場合、システムの起動時に最初にバッファポイン

10

20

30

40

50

タを発生し、バッファプールにそれらをロードすることは比較的時間のかかるプロセスである。

EP-A-0 622 922には、通信システムにおいてデータをマルチキャストする際にバッファを解放するための典型的な公知の機構が開示されている。その図7Aを参照して、最後のメッセージの送信後に、直接制御ブロックはそのメッセージを含むバッファをポイントする。その後、間接制御ブロックが解放され、直接制御ブロックから複製の数が読出され、これは、送信されたものが最後のメッセージであったことを示す1に等しいであろう。複製の数は0までデクリメントされる。その後、複製の数は0に等しいと判断され、この判断に応答して初めて直接制御ブロックが解放される。

発明の概要

メモリのバッファからデータを送信する際に素早い態様でバッファを再要求し、他のデータをストアするためにバッファを再使用できるようにするための構成および方法が必要である。

この必要性および他の必要性は、パケット交換ネットワークにおけるネットワークスイッチに対し、フレームをストアするバッファを再要求するための構成を提供するこの発明によって満たされる。この構成は、フレームがストアされる初めのバッファのバッファメモリ内の場所をポイントするフレームポインタをキューとして維持するよう構成されたリクレーンキューを含む。探索可能な第1のメモリは、ネットワークスイッチから送信されるべきフレームのコピー数のカウント値をストアする。フリーバッファプールは、バッファメモリにおけるバッファの場所をポイントするバッファポインタをキューとして維持するよう構成される。バッファマネージャは、フレームのコピーを送信するためにそのフレームのフレームポインタを受信し、このコピーの送信が、送信されるべきフレームの最後のコピーであるかどうかを判断するために第1のメモリを探索する。この送信が最後のコピーに関するものであれば、バッファマネージャはそのフレームの最後のコピーがネットワークスイッチから送信されるようにし、フレームをストアする各バッファのバッファポインタを、そのバッファの内容が送信された後にフリーバッファプールに送る。

この発明の利点の1つは、送信されることとなる、次に送信されるフレームのコピーが、送信されるべき最後のコピーであるかと判断することが可能である場合、フレームをストアするバッファが、各バッファの内容が送信された後に再要求されることである。これにより、フレーム全体が送信されるまで待機し、そのフレームのコピーがこれ以上送信されないと判定して、そのときにはじめてフレームをストアしたバッファの再要求を開始する場合よりも、より素早く再使用のためにバッファが戻されるようになる。

前述の必要性および他の必要性は、ネットワークスイッチからのフレームの送信の後にそのフレームをストアしているバッファを再要求するための方法を提供するこの発明の別の実施例によって満たされる。発明の方法では、送信されるフレームのコピーが、送信されるべきフレームの最後のコピーであるかどうかを判断しようとする。送信されるフレームのコピーが送信されるべきフレームの最後のコピーであると判断された場合、フレームをストアしているバッファは、ネットワークスイッチから各バッファの内容が送信された後に再要求される。しかしながら、送信されるフレームのコピーが送信されるべきフレームの最後のコピーでないと判断された場合、フレームのコピーの送信を示す表示が、ネットワークスイッチからフレームが送信された後にキューとして維持される。また、送信されるフレームのコピーが送信されるべきフレームの最後のコピーであるかどうかを判断することができない場合にも、フレームのコピーの送信を示す表示が、ネットワークスイッチからフレームが送信された後にキューとして維持される。

この発明のこの実施例の1つの利点は、フレームの送信状態、すなわちそれが最後のコピーであるかどうかを判断する試みがなされることである。判断可能である場合、各バッファの内容が送信された直後にバッファを再要求することにより、再要求プロセスの速度を高めることができる。これが可能なのは、バッファのデータがそのフレームの別のコピーを送るために必要ではないことがわかっているためである。しかしながら、それが最後のコピーでないとき、またはそれが最後のコピーかどうか判断できない場合、送信を示す表

10

20

30

40

50

示がキューとして維持される。この表示は、送信されたコピー数のカウント値を維持するために用いられ得る。したがって、状況によっては、この発明は再要求プロセスの速度を高める。

発明のいくつかの実施例においては、フレームの単一のコピーのみが送られるか、または複数のコピーが送られるかが判断される。このような実施例では、単一のコピーの送信時にはバッファを再要求するプロセスの速度を高めることが可能である。なぜなら、バッファの内容はそのフレームの他のコピーの送信には必要でないことがわかっているからである。したがって、フレームの単一のコピーのバッファの再要求は、(複数のコピーフレームの最後のコピーであることがわかっている場合の)複数のコピーフレームの最後のコピーと同じ態様で行なうことができる。

10

また、起動時に時間のかからない、システムの起動時にバッファポインタを発生するための構成および方法が必要である。

この必要性および他の必要性は、メモリ内にデータをストアするためのバッファが使用できるようにするための構成を提供するこの発明によって満たされる。この構成は、複数のバッファポインタをストアするよう構成されたフリーバッファプールを含む。各バッファポインタは、関連のフリーバッファ、すなわちデータをストアするために現在使用できるバッファがある、メモリ中の別々の場所をポイントする。データをストアするために使用できるようになったバッファのバッファポインタをフリーバッファプールに戻す再要求回路が提供される。カウンタはこの構成の起動時にバッファポインタを発生して、データのストアに使用できるフリーバッファを提供する。マルチプレクサは、フリーバッファプールの出力に結合された第1の入力と、カウンタの出力に結合された第2の入力と、第1の入力と第2の入力との間で選択を行なう制御入力とを有する。

20

この発明のこの実施例の利点は、システムの起動時に、バッファポインタが使用できるようになる前に発生されてフリーバッファプールにロードされるのではなく、必要に応じて発生される前に発生されてフリーバッファプールにロードされるのではなく、必要に応じて発生されることである。32,000個のバッファポインタなどといった、発生され得るバッファポインタの数が非常に多い場合、データをストアするためのバッファの使用を開始する前に、これらのポインタのすべての発生およびロードを待機することは時間の無駄である。この発明では、バッファポインタのすべてが発生されるのを待機することなく、カウンタによって発生されたバッファポインタがポイントするバッファをデータのストアのために直ちに用いることができる。最終的に、スイッチの動作時にバッファポインタのすべてが発生されると、カウンタは最早必要でなくなり、マルチプレクサはバッファポインタのすべてがフリーバッファプールから取出されるように切換わる。

30

この発明の前述および他の特徴、局面および利点は、添付の図面に関連して読むとこの発明の以下の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。

40

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。

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

50

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ブル領域を詳細に示す図である。

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ブ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。

10

例示的な実施例の詳細な説明

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 10 M \ bps 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 (M \ bps) のネットワークステーション 14 と、100 M \ bps のネットワーク速度でデータパケットの授受を行なう 2 つの 100 M \ bps 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。

20

開示される実施例によると、10 M \ bps 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。

30

40

100 M \ bps 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 100 M \ bps 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。100 M \ bps ネットワークステーション 22 の各々は、100 ベース

50

- T X、1 0 0 ベース - T 4 または 1 0 0 ベース - F X タイプの 1 0 0 M \ b p s 物理 (P H Y) 装置 2 0 を介してネットワーク媒体 1 7 に結合される。スイッチ 1 2 は、物理装置 2 0 への接続をもたらす媒体独立インタフェース (M I I) 2 4 を含む。1 0 0 M \ b p s ネットワーク 2 2 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 1 0 は、スイッチ 1 2 と 1 0 M \ b p s ステーション 1 4 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 1 6 を含む。磁気変成器モジュール 1 9 は媒体 1 7 上の信号の波形を維持する。スイッチ 1 2 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 (N R Z) インタフェース 2 3 を介して各スイッチトランシーバ 1 6 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 1 8 を含む。スイッチトランシーバ 1 6 はシリアル N R Z インタフェース 2 3 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 1 7 を介して適切なエンドステーション 1 4 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 1 6 は独立した 4 つの 1 0 M \ b p s ツイストペアポートを有し、スイッチ 1 2 が必要とする P I N の数が 4 分の 1 に減少するようにするシリアル N R Z インタフェースを介する 4 : 1 多重化を用いる。

スイッチ 1 2 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 1 4 および 1 2 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 M A C (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 1 2 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す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

10

20

30

40

50

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100M\bps MACポート53(ポート25および26)とを含む。上述のとおり、管理インタフェース30もまたMAC層プロトコル(ポート0)に従って動作する。MACポート50、53および30の各々は、受信先入れ先出し(FIFO)バッファ52と送信FIFO54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MACポートで受信され、対応の受信FIFO52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FIFO52から外部メモリインタフェース34に出力されて、外部メモリ36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58または外部ルールチェッカインタフェース40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MAC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12が内部ルールチェッカ58または外部ルールチェッカ42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58または外部ルールチェッカインタフェース40に送られる。内部ルールチェッカ58および外部ルールチェッカ42は、所与のデータパケットに関する宛先MAC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MAC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、(100M\bpsステーション22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IEEE 802.1d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)の識別子とを含む。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GPSI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

10

20

30

40

50

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、リクレームキュー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VLAIN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

10

20

30

40

50

送信がキューとして維持される。

ある時点で、フレームポインタは出力キュー 7 4 を通過した後に出力キュー 7 4 の一番下まで到達する。バッファマネージャ 7 2 はそれが出力キュー 7 4 の一番下まで到達したときにフレームポインタを取り、フレームポインタ書込バス 8 6 を介して正しいポートの適切な送信 F I F O 5 4 にそのフレームポインタを送る。これによりフレームの送信がスケジュールされる。フレームポインタによってポイントされた外部メモリ 3 6 での場所から D M A 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 F I F O 5 4 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。フレームポインタは適切な出力キュー 7 4 の各々に置かれ、対応の送信 F I F O 5 4 から送信される。

10

バッファマネージャ 7 2 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 1 0 4 と、マルチコピーキュー 9 0 と、リクレームキュー 9 8 と、マルチコピーキャッシュ 9 6 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 7 2 はまた、出力キュー 7 4 ならびに制御キュー 1 0 4、9 0 および 9 8 のために外部メモリ 3 6 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 5 6 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 7 4 の構造である。1 0 M b / s ポートおよび 1 0 0 M b / s 出力ポートに指定される出力キュー 7 4 の他に、管理ポート 3 0 のために出力キュー 7 5 が設けられ、拡張ポート 2 6 のために出力キュー 7 7 が設けられる。これらの出力キュー 7 5 および 7 7 は出力キュー 7 4 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。

20

図 4 は、この発明の実施例に従う出力キュー 7 4 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 7 4 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 1 2 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。

30

図 4 の実施例に従うこの発明の単一論理出力キュー 7 4 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 7 6 と、出力キュー読出側 7 8 と、外部メモリ 3 6 にある出力キューオーバフロー領域（全体が 1 1 0 として示される）とが含まれる。出力キュー 7 4 のすべてに関する外部メモリ 3 6 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 3 4 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使用し、（フレームポインタなどの）データが、チップ 1 2 を外部メモリ 3 6 に接続するバス 8 4 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 1 1 0 に送られるようにする。

40

出力キュー書込側 7 6 および出力キュー読出側 7 8 はチップ 1 2 上にある。書込側 7 6 および読出側 7 8 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 7 4 の第 3 の部分を形成するオーバフロー領域 1 1 0 は大きくて比較的安価である。書込側 7 6 および読出側 7 8 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 7 6 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 1 2 の例示的な実施例では、エントリは、フレームの最初の 2 5 6 バイトがストア

50

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。

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。

10

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。

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を使用するが、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。

20

30

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

40

50

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（各々が出力ポートに関連する）があり、すなわち、10Mb/sユーザポートに関するものが24個、100Mb/sサー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FIFO70がフレームポインタを書込むという形態をとる。

10

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VLANTAG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。

20

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。

30

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。

40

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て。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

50

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。

MIBカウンタ領域132は、スイッチ12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12はMIB統計をストアするための8ビットおよび16ビットカウンタをチップ上に維持する。スイッチ12はMIB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36の32ビットまたは64ビットのMIB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134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72によってフリーバッファプール104に戻されるか、またはPCIホストプロセッサ28の所有となる。

次に図8を参照して、いずれかのMACポートまたはPCI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36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140の長さは256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140の各々へのアドレスポインタはスイッチ12内のフリーバッファプール104によってストアされる。

スイッチ12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72はフリーバッファプール104からアドレスポインタを要求し、バッファ140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36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56において用いられる。

バッファ140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142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142はまた、バッファ140に含まれるフレームデータに関する他の情報を含む。図9a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12バイトである。図9b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×16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256B - 1

10

20

30

40

50

6 B = 2 4 0 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。

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12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（CRCエラー検出）：CRC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CRC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（CRCエラーとともに）受信フレームに検出されたことを示す。

Oビット（受信FIFOオーバーフロー）：受信FIFOがオーバーフローし、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142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16番目のバ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1

10

20

30

40

50

μs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VLANTAG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ルする。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FIFO70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FIFO70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ト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FIFO70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FIFO70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ルする。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FIFO54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

10

20

30

40

50

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。

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、VLANN 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ルするが、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 に見つかれば、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ルし、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。

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。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。

10

20

30

40

50

バッファマネージャ 72 は、フレームポインタを読み出してから、リンク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 F I F O 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 F I F O 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 F I F O 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 F I F O 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。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 F I F O 70 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 F I F O 70 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み出し、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。一致が見つければ、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（P C I 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

10

20

30

40

50

成される。入力経路または書込側により、ポートベクタ F I F O 7 0 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 9 0 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 9 0 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 9 6 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 1 2 4 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 3 6 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 9 0 に書込まれると、それらは読出側 9 4 がフルになるまで書込側 9 2 から読出側 9 4 へと移動する。マルチコピーキュー 9 0 の書込側 9 2 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 3 6 内のマルチコピーキューオーバーフロー領域 1 2 4 に入れられる。一旦マルチコピーキュー 9 0 の読出側 9 4 とそのオーバーフロー領域 1 2 4 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 9 2 を満たし始める。

マルチコピーキュー 9 0 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 9 4 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 1 2 4 からマルチコピーキューの読出側 9 4 へと移動し、マルチコピーキューの書込側 9 2 からマルチコピーキューオーバーフロー領域 1 2 4 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 9 6 はマルチコピーキュー 9 0 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 9 6 は 2 5 6 までのエントリを保持する。バッファマネージャ 7 2 はマルチコピーキュー 9 0 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 9 6 に入れるかそれを処理するかする。

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。

10

20

30

40

50

a) 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る。マルチコピーキャッシュ 96 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 72 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 98 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 72 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 72 がマルチコピーキャッシュ 96 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 36 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 72 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 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) この値が「> 1」であれば、バッファマネージャ 72 は(外部メモリ 36 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 96 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 96 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 36 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 98 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 72 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 70 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 F I F O 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

10

20

30

40

50

カ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トがなければ、オーバーフロー領域120がより多くのバッファポイント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ト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がバッファポイントを含むことは必要とされない。代わりに、バッファポイント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トが発生し、それをオーバーフロー領域120に入れることができるが、このようなポイントは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使用する。したがって、バッファポイント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はバッファポイントが発生する。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トが発生しない。

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、バッファポイント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トとそのコピー数とを含む。フレームポイントの最下位6ビットが、エントリがス

10

20

30

40

50

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 FIF070 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 FIF070 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 FIF070 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 FIF070 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 FIF070 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 10Mb/s ポートから 100Mb/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、10Mb/s ポートから 100Mb/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

10

。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FIFO52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FIFO70に戻す。

ポートベクタFIFO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FIFO70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の読出側78に移動すると、バッファマネージャ72はフレームポインタを読出し、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MACポートの送信FIFO54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。

20

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CRC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CRCで終る場合、受信MACは外部メモリ36内のバッファヘッダ142に印を付けてこれらの条件を示す。送信MACは、後に短すぎるものか無効なCRCで終るフレームの送信が始まればMACが無効なCRCを発生することを保証する。送信MACがフレーム送信を始めておらず、バッファヘッダ142が短すぎるものか無効なCRCで終るフレームを示している場合、バッファマネージャ72はフレームを出力ポートへと転送しない。

30

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12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CRCエラーフレームを完全にはフィルタ処理しない。

40

第2のモードにおいては、受信MAC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74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74に入れられない。64バイト以上の受信フレームが無効なCRCで終れば、受信MACは外部メモリ36内のバッファヘッダ142に印を付けてこの条件を示す。後に無効なCRC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MACは不良なCRCで送信を終了する。送信MACがフレーム送信を開始しておらず、バッファヘッダ142が無

50

効なCRC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)リクレームキュー98または(マルチコピー転送のための)マルチコピーキュー96へと出力ポート74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12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12は転送の前に全てのフラグメントおよびCRC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CRC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74に

10

入れられる。受信エラー(無効CRC、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74に入れられない。ポートベクタFIFO70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74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74へと転送する。ポートベクタFIFO70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。

20

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FIFO70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FIFO70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。

30

ポートベクタFIFO70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FIFO70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を示される。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FIFO70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FIFO70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FIFO70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FIFO70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FIFO70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FIFO70次第である。ポートベクタFIFO70は各ポートのモードを識別する情報を維持する。要

40

50

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FIFO70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。

バッファマネージャ72はこの発明のスイッチ12においてさまざまな機能を果たす。バッファマネージャ72はこれらの機能を実現するための制御論理および/またはソフトウェアを含み、バッファマネージャ72の機能に関する上記の説明があれば当業者なら容易に実現できる。

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、限定するものとは理解されるべきでなく、本発明の精神および範疇が請求の範囲によってのみ規定されることが明らかに理解される。

10

【図1】

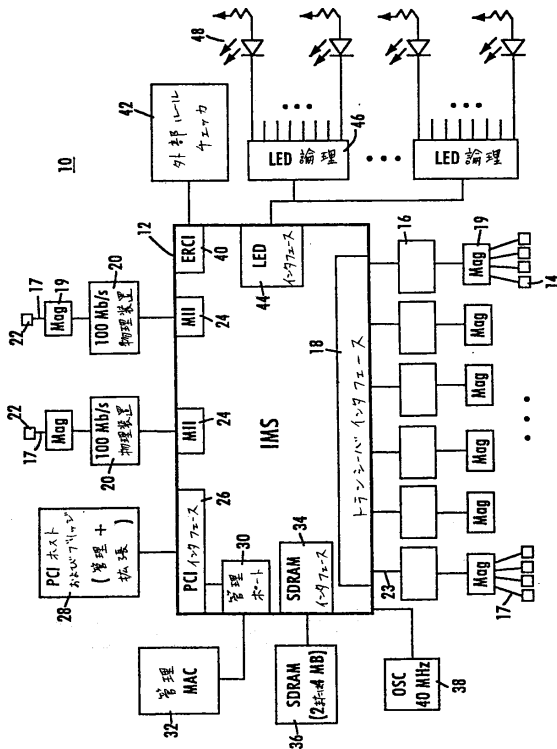
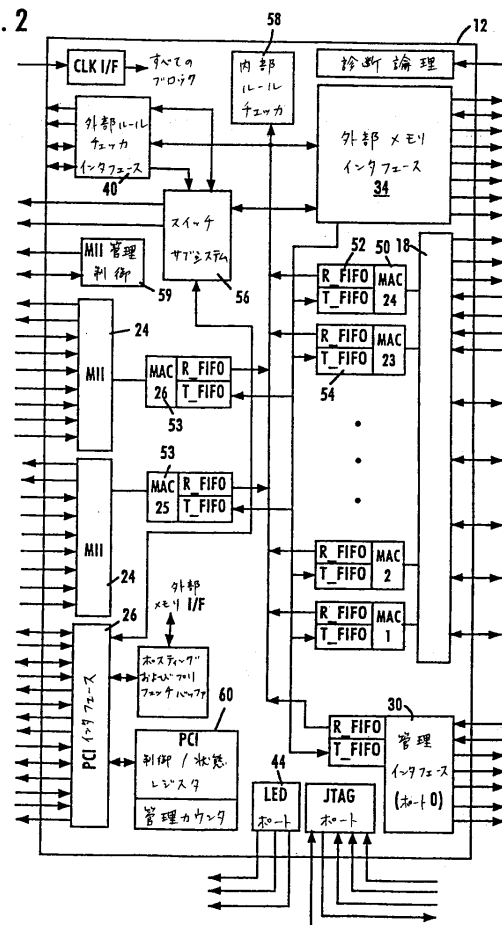


Fig. 1

【図2】

Fig. 2



【 図 3 】

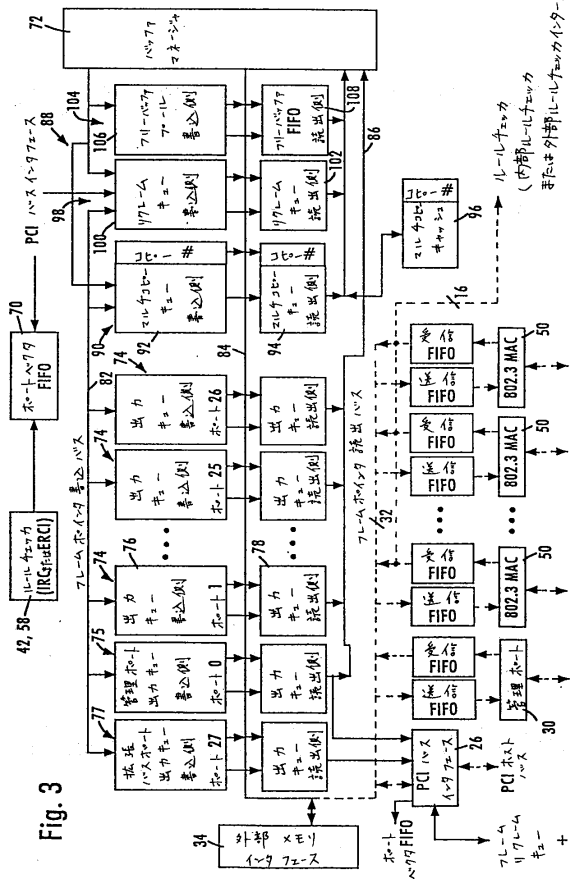


Fig. 3

【 図 4 】

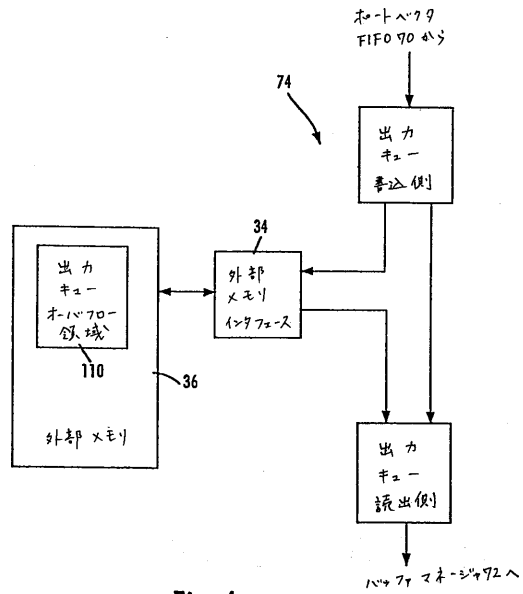


Fig. 4

【 図 5 】

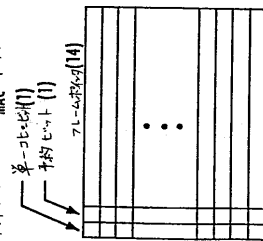


Fig. 5

【 図 6 】

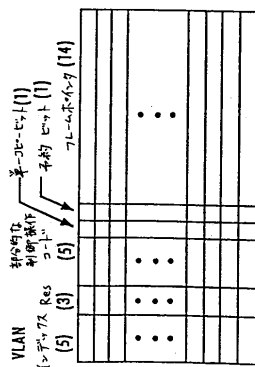


Fig. 6

【圖 7】

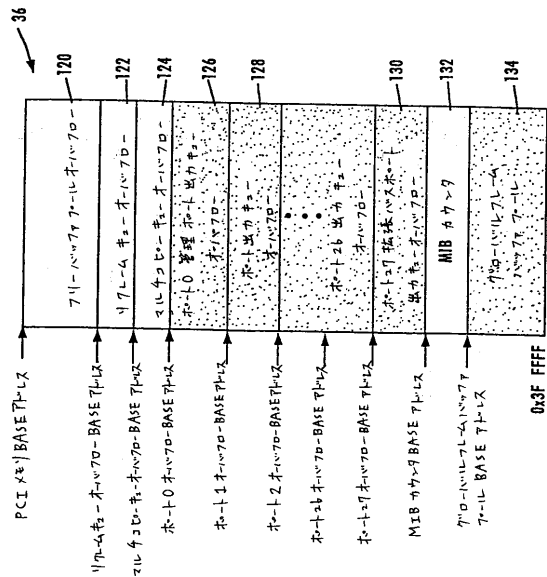


Fig. 7

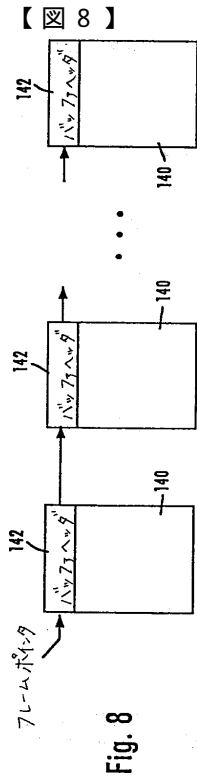


Fig. 8

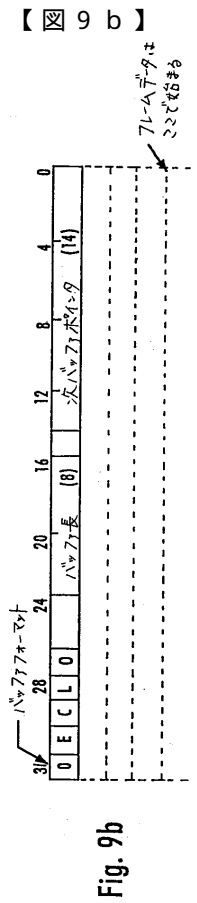


Fig. 9b

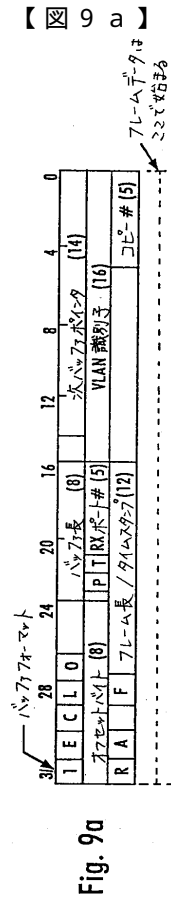


Fig. 9a

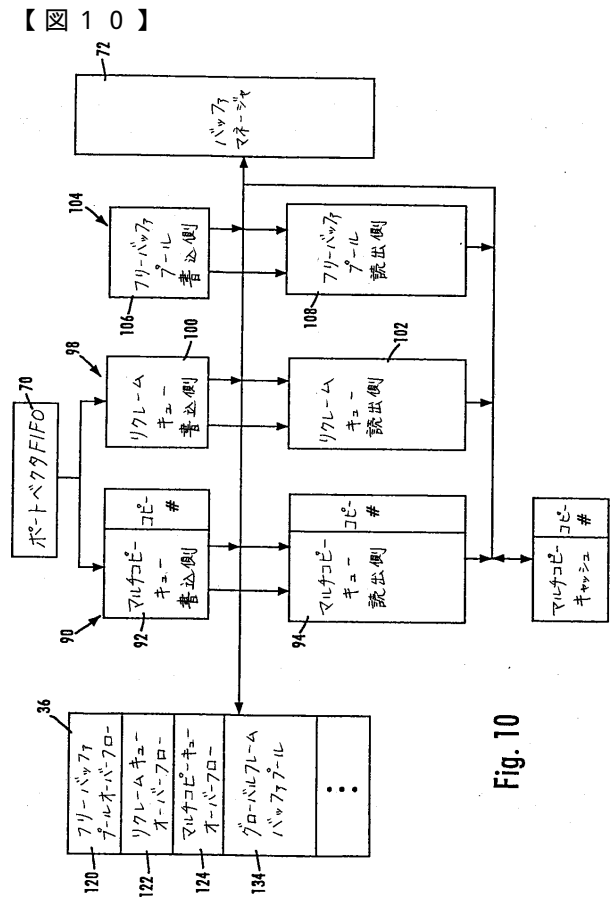


Fig. 10

【図 11】

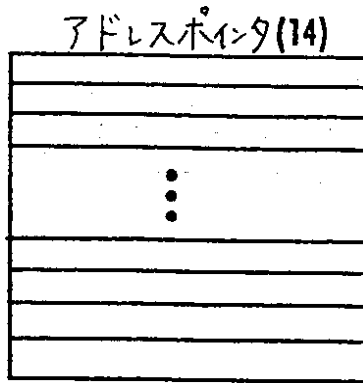


Fig. 11

【図 12】

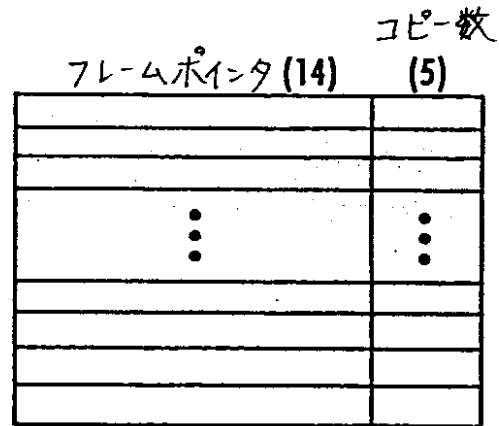


Fig. 12

【図 13】

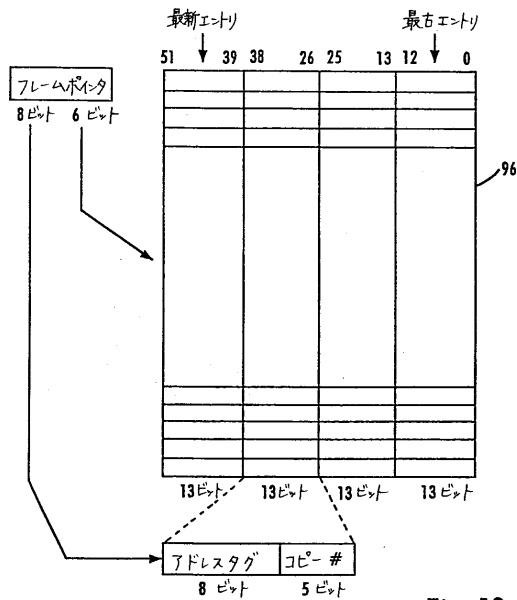


Fig. 13

【図 14】

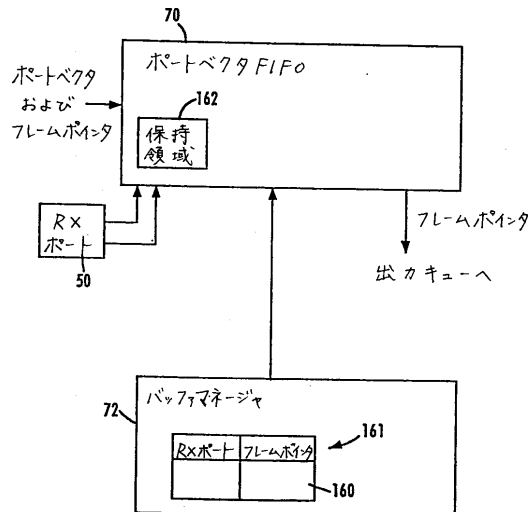


Fig. 14

フロントページの続き

(74)代理人 100096781

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・ジェファーソン

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

(72)発明者 エグバート, チャンダン

アメリカ合衆国、9 5 1 3 2 カリフォルニア州、サン・ノゼ、ブルームズバリー・ウェイ、3 6 3 2

審査官 吉田 隆之

(56)参考文献 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 8 - 8 9 0 6 (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