

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5593517号  
(P5593517)

(45) 発行日 平成26年9月24日 (2014. 9. 24)

(24) 登録日 平成26年8月15日 (2014. 8. 15)

(51) Int. Cl.

F I

H O 4 L 12/28 (2006. 01)

H O 4 L 12/28 2 O O D

H O 4 L 12/825 (2013. 01)

H O 4 L 12/825

請求項の数 12 (全 27 頁)

(21) 出願番号 特願2011-156939 (P2011-156939)  
 (22) 出願日 平成23年7月15日 (2011. 7. 15)  
 (65) 公開番号 特開2013-26680 (P2013-26680A)  
 (43) 公開日 平成25年2月4日 (2013. 2. 4)  
 審査請求日 平成25年11月5日 (2013. 11. 5)

(73) 特許権者 000005108  
 株式会社日立製作所  
 東京都千代田区丸の内一丁目6番6号  
 (74) 代理人 100114236  
 弁理士 藤井 正弘  
 (74) 代理人 100075513  
 弁理士 後藤 政喜  
 (74) 代理人 100120260  
 弁理士 飯田 雅昭  
 (72) 発明者 早川 仁  
 東京都国分寺市東恋ヶ窪一丁目280番地  
 株式会社日立製作所 中央研究所内  
 (72) 発明者 對馬 雄次  
 東京都国分寺市東恋ヶ窪一丁目280番地  
 株式会社日立製作所 中央研究所内  
 最終頁に続く

(54) 【発明の名称】 ネットワーク装置及び送信フレームの制御方法

(57) 【特許請求の範囲】

【請求項 1】

ノードとのリンクを流れる複数のトラフィックのフレームの送受信を行い、送信フレームの送信を行う送信部と受信フレームの受信を行う受信部とを備えるポートと、

前記ポートを介して受信した受信フレームが一時停止の指令のときには、前記ポートからの送信フレームの送信を所定時間まで一時的に停止し、

前記複数のトラフィックから選択した少なくとも1以上のトラフィックの送信帯域を制限し、

前記所定時間を経過した後、前記一時停止の指令を含む受信フレームを受信する間隔が拡大したときは、前記選択したトラフィックの送信帯域の制限を強化する制御部と、  
 を備えたことを特徴とするネットワーク装置。

10

【請求項 2】

請求項 1 に記載のネットワーク装置であって、

前記制御部は、

所定時間を経過した後、前記一時停止の指令を含む受信フレームを受信する間隔が拡大しないときは、前記選択したトラフィックの送信帯域の制限を緩和することを特徴とするネットワーク装置。

【請求項 3】

請求項 1 に記載のネットワーク装置であって、

前記制御部は、

20

予め設定した順序に従って前記送信帯域を制限する少なくとも 1 以上のトラフィックを選択することを特徴とするネットワーク装置。

【請求項 4】

請求項 1 に記載のネットワーク装置であって、

前記一時停止の指令は、P A U S E フレームであることを特徴とするネットワーク装置

。

【請求項 5】

請求項 1 に記載のネットワーク装置であって、

前記制御部は、

前記一時停止の指令を受信する間隔が拡大した場合には、当該選択したトラフィックを、輻輳を引き起こすトラフィックとして推定することを特徴とするネットワーク装置。

10

【請求項 6】

プロセッサと、メモリと、ノードとのリンクを流れる複数のトラフィックのフレームの送受信を行い、送信フレームの送信を行う送信部と受信フレームの受信を行う受信部とを備えるポートと、を備えたネットワーク装置で送信するフレームの輻輳を抑制する送信フレームの制御方法であって、

前記ポートを介して受信した前記受信フレームが一時停止の指令のときには、前記ポートからの前記送信フレームの送信を所定時間まで一時的に停止する第 1 のステップと、

前記複数のトラフィックから選択した少なくとも 1 以上のトラフィックの送信帯域を制限する第 2 のステップと、

20

前記所定時間を経過した後、前記一時停止の指令を含む受信フレームを受信する間隔が拡大したときは、前記選択したトラフィックの送信帯域の制限を強化する第 3 のステップと、

を含むことを特徴とする送信フレームの制御方法。

【請求項 7】

請求項 6 に記載の送信フレームの制御方法であって、

前記第 3 のステップは、

所定時間を経過した後、前記一時停止の指令を含む受信フレームを受信する間隔が拡大しないときは、前記選択したトラフィックの送信帯域の制限を緩和することを特徴とする送信フレームの制御方法。

30

【請求項 8】

請求項 6 に記載の送信フレームの制御方法であって、

前記第 2 のステップは、

予め設定した順序に従って前記送信帯域を制限する少なくとも 1 以上のトラフィックを選択することを特徴とする送信フレームの制御方法。

【請求項 9】

請求項 6 に記載の送信フレームの制御方法であって、

前記一時停止の指令は、P A U S E フレームであることを特徴とする送信フレームの制御方法。

【請求項 10】

40

請求項 6 に記載の送信フレームの制御方法であって、

前記第 3 のステップは、

前記一時停止の指令を受信する間隔が拡大した場合には、当該選択したトラフィックを、輻輳を引き起こすトラフィックとして推定することを特徴とする送信フレームの制御方法。

【請求項 11】

ノードとのリンクを流れる複数のトラフィックのフレームの送受信を行い、送信フレームの送信を行う送信部と受信フレームの受信を行う受信部とを備えるポートと、

前記ポートを介して受信した受信フレームが一時停止の指令のときには、前記ポートからの送信フレームの送信を所定時間まで一時的に停止し、

50

前記複数のトラフィックから選択した少なくとも1以上のトラフィックの送信帯域を制限し、

前記所定時間を経過した後、前記一時停止の指令を含む受信フレームを受信する間隔が拡大したときは、前記選択したトラフィック以外のトラフィックの送信を促進する制御部と、

を備えたことを特徴とするネットワーク装置。

【請求項12】

ノードとのリンクを流れる複数のトラフィックのフレームの送受信を行い、送信フレームの送信を行う送信部と受信フレームの受信を行う受信部とを備えるポートと、

前記ポートを介して受信した受信フレームが一時停止の指令のときには、前記ポートからの送信フレームの送信を所定時間まで一時的に停止し、

前記複数のトラフィックから選択した少なくとも1以上のトラフィックの送信帯域を制限し、

前記所定時間を経過した後、前記一時停止の指令を含む受信フレームを受信する間隔が拡大しないときは、前記選択したトラフィックの送信帯域の制限を緩和する制御部と、  
を備えたことを特徴とするネットワーク装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、コンピュータネットワークでフレームを転送する装置及び方法に関し、特に、トラフィックの輻輳を制御する技術の改良に関する。

【背景技術】

【0002】

コンピュータネットワークを実現するイーサネット（登録商標）は、データをフレームという単位に分割して送受信し、中継装置等の通信系路上でのフレーム廃棄を許容する標準規格となっている。フレーム廃棄によるデータ損失を許容しないアプリケーションに対しては、拡張規格としてIEEE 802.3xにフロー制御用のPAUSEフレームが規定されている。フロー制御は送信側と受信側を結ぶリンクに対して行われる。リンクの受信側で、輻輳等の理由でフレームを正しく受信できないと想定される場合、リンクの受信側から送信側にPAUSEフレームを送信する。PAUSEフレームを受信したリンクの送信側ではフレーム送信を一時停止する。これらリンクの送信側及び受信側での処理により、リンクのフロー制御が行われ、リンクの受信側でフレームを受信できないことによるフレーム廃棄を防ぐ。

【0003】

一方、サーバの高密度化や、Local Area Network（LAN）に用いられるイーサネットの高速化を背景として、これまでLANとは別に構築されていたストレージアクセスのためのStorage Area Network（SAN）を、イーサネットを用いて構築してLANに統合し、ネットワークインタフェースやケーブル、ネットワーク機器を削減して運用の効率化をはかるLAN/SAN統合の動きが加速している。また、サーバ仮想化の進展により、1台の物理サーバ上で複数の仮想マシンを稼働させ、1つのLANを複数のネットワークで共用する環境が一般化している。これらの環境では物理的に1本のイーサネットケーブル上に性質や要求の異なる複数のトラフィックを流す必要がある。これらのトラフィックを論理的に区別する手法としてIEEE 802.1Qに規定されるVirtual LAN（VLAN）が広く普及している。VLANでは各フレームにVLANタグを付加し、VLANタグに含まれる優先度およびVLAN IDでトラフィックを区別する。VLANを用いることで1本の物理的なリンクを論理的に分割して用いることが可能となる。

【0004】

LAN/SAN統合におけるストレージアクセスのトラフィックはフレーム廃棄を許容しない。しかしながら、PAUSEによるフロー制御は物理的なリンク単位で行うため、VLANにより論理的に分割されたリンクのいずれかで輻輳が発生してPAUSEフレー

10

20

30

40

50

ムが送信されて送信が一時停止されると、同じ物理的なリンクを共有する他の論理的なリンクを流れるトラフィックのフレーム送信も一時停止してしまう。すなわち、複数のトラフィックで1つのリンクを共有するため、一部トラフィックでの輻輳が他のトラフィックの通信を妨げるという問題が生じる。

【0005】

これに対する従来技術として、論理的なリンク毎に独立したP A U S Eフレームを規定する手法が知られている（例えば、特許文献1および非特許文献1）。

【0006】

特許文献1には「複数の入出力ポートから入力する伝送フレームの輻輳を検出する輻輳検出手段と、前記輻輳検出手段が輻輳を検出した入出力ポートにP A U S Eフレームを送出するP A U S Eフレーム送出手段とを備えたネットワークスイッチを有する輻輳制御システムにおいて、前記伝送フレームの代わりに仮想L A N情報を付加した拡張伝送フレームと、前記P A U S Eフレームの代わりに前記P A U S Eフレームに仮想L A N情報を付加した拡張P A U S Eフレームとを使用し、前記P A U S Eフレーム送出手段の代わりに前記拡張P A U S Eフレームを前記入出力ポートに送出する拡張P A U S Eフレーム送出手段を設けたことを特徴とする。」と記載されている。さらに「特に、前記仮想L A N情報をI E E E 8 0 2 . 3 x規格のV L A NタグのV L A N - I D値とすることを特徴とする。或いは、前記仮想L A N情報をI E E E 8 0 2 . 3 x規格のV L A Nタグのプライオリティ値とすることを特徴とする。」と記載されている。

【0007】

また、非特許文献1では、I E E E 8 0 2 . 1 Q b bとしてV L A Nタグのプライオリティ毎のP A U S E指示を含む拡張したP A U S Eフレームおよび当該P A U S Eフレームを用いた論理的なリンク毎のフロー制御を規定している。

【先行技術文献】

【特許文献】

【0008】

【特許文献1】特開2007-174152号公報

【非特許文献】

【0009】

【非特許文献1】IEEE P802.1Qbb/D2.3, "IEEE Draft Standard for Local and Metropolitan Area Networks - Virtual Bridged Local Area Networks - Amendment: Priority-based Flow Control", Publication No.5570060, Issue date Sept. 9 2010

【発明の概要】

【発明が解決しようとする課題】

【0010】

上記特許文献1および非特許文献1の処理では、リンクの送信側と受信側の双方がそれぞれ対応する必要があるため、既存の設備を全面的に置き換えまたは改造することが必要となる。

【0011】

また、上記特許文献1では、I E E Eの標準規格として規定されているP A U S EおよびV L A Nを組み合わせたものであるが、これらを組み合わせた処理について、上記標準には規定されていないため、既存の装置にそのまま適用することはできず、また、上記特許文献1の技術を利用した装置も広範に利用可能な状況ではない。

【0012】

一方、上記非特許文献1は現在策定中の標準規格であり、将来的に対応する装置が増えると考えられるが、非特許文献1の技術で扱うことが可能な論理的なリンクの数は8までという制限がある。これは仮想マシンの増加やトラフィックの優先度制御を考慮すると不十分であり、8を超えた数のトラフィックを流す場合は、1つの論理的なリンク内を複数のトラフィックが流れる状況が生じ、先に挙げた一部トラフィックでの輻輳が他のトラフィックの通信を妨げるという問題は解決できない。

## 【 0 0 1 3 】

したがって、上記従来技術では、標準規格に準拠し、技術導入時に必要な機器の変更が最小限で、フレーム廃棄を許容しない複数のトラフィックを扱い、かつ8を超える多数のトラフィックに対応可能な輻輳制御を行うことができないという課題があった。

## 【課題を解決するための手段】

## 【 0 0 1 4 】

本発明は、ノードとのリンクを流れる複数のトラフィックのフレームの送受信を行い、送信フレームの送信を行う送信部と受信フレームの受信を行う受信部とを備えるポートと、前記ポートを介して受信した受信フレームが一時停止の指令のときには、前記ポートからの送信フレームの送信を所定時間まで一時的に停止し、前記複数のトラフィックから選択した少なくとも1以上のトラフィックの送信帯域を制限し、前記所定時間を経過した後、前記一時停止の指令を含む受信フレームを受信する間隔が拡大したときは、前記選択したトラフィックの送信帯域の制限を強化する制御部と、を備える。

10

## 【発明の効果】

## 【 0 0 1 5 】

したがって、本発明によれば、既存のネットワーク機器をそのまま利用しながらフレームの廃棄を許容しない複数のトラフィックの輻輳制御が可能となる。

## 【図面の簡単な説明】

## 【 0 0 1 6 】

【図1】本発明の実施形態を示し、ネットワークシステムの構成の一例を示すブロック図である。

20

【図2】本発明の実施形態を示し、フレーム転送装置のネットワークシステムの構成の一例を示すブロック図である。

【図3】本発明の実施形態を示し、転送ポート管理表の一例を示す図である。

【図4】本発明の実施形態を示し、トラフィック設定表の一例を示す図である。

【図5】本発明の実施形態を示し、キュー設定表の一例を示す図である。

【図6】本発明の実施形態を示し、トラフィック履歴表の一例を示す図である。

【図7】本発明の実施形態を示し、帯域割当表の一例を示す図である。

【図8】本発明の実施形態を示し、輻輳トラフィック推定部における制限帯域計算処理の一例を示すフローチャートである。

30

【図9】本発明の実施形態を示し、輻輳トラフィック推定部における制限帯域補正処理の一例を示すフローチャートである。

【図10A】本発明の実施形態を示し、送信フレーム選択部における送信帯域計算処理の一例を示すフローチャートで、前半部である。

【図10B】本発明の実施形態を示し、送信フレーム選択部における送信帯域計算処理の一例を示すフローチャートで、後半部である。

## 【発明を実施するための形態】

## 【 0 0 1 7 】

以下、本発明の一実施形態を添付図面に基づいて説明する。

## 【 0 0 1 8 】

40

図1は、本発明の実施形態を示し、ネットワークシステムの構成の一例を示すブロック図である。

## 【 0 0 1 9 】

<システムおよびハードウェアの構成>

図1は本実施形態におけるネットワークシステムの構成を説明するブロック図である。図1において、ネットワークシステムはフレーム転送装置1の入出力ポート112-1にケーブル115-1を介して接続されたノード117-1、および、入出力ポート112-nにケーブル115-nを介して接続されたノード117-mから構成される。なお、ノード117-1~117-mは物理計算機やネットワーク装置あるいはストレージ装置等で構成される。また、フレーム転送装置1は、n個の入出力ポート112-1~112

50

- nを有する。以下では、入出力ポートの総称を112とし、ケーブルの総称を115とし、ノードの総称を117とする。

【0020】

本実施形態におけるネットワークシステムはノード117-1とノード117-mの間でケーブル115-1およびケーブル115-nとフレーム転送装置1を介したイーサネット（登録商標）による通信を実現する。ノード117-1～ノード117-mはイーサネットに準拠した入出力ポートを有し、P A U S Eフレームによるフロー制御をサポートしたネットワークノードであり、例えばネットワークインタフェースカードを装着した計算機や、ルータ装置、スイッチ装置等のネットワーク機器等である。ケーブル115はイーサネットによる通信に適合するケーブルである。フレーム転送装置1はイーサネットフレームを転送するブリッジの機能を持つネットワーク装置である。

10

【0021】

イーサネットについてはI E E E 802.3に規定されており、ブリッジ機能はI E E E 802.1DおよびI E E E 802.1Qに規定されている。本実施形態においては、フレーム転送装置1は、ブリッジ機能に加えて、イーサネットに関して、P A U S Eフレームによるフロー制御、特にI E E E 802.3bd及びI E E E 802.1Qbbに規定されるP F C (Priority-based Flow Control)を対象として、ノード117-1～ノード117-mの間で交換されるイーサネットフレームに関して、M A C (Media Access Control) アドレスやV L A N I D (Virtual Local Area Network Identifier)、T C P (Transmission Control Protocol) ポート番号等で区別される論理的なトラフィックに対して輻輳制御機能を提供する。ただし、本発明はイーサネットに限定されるものではなく、例えばF i b r e C h a n n e lにおけるB u f f e r - t o - B u f f e rクレジット等、受信側の状態について送信側に通知するフロー制御の機能を持ったプロトコルであれば、本実施形態で説明する輻輳制御方法を適用可能である。なお、本実施形態において、ノードは117-1～117-mのm台の場合で説明を行うが、ノードの個数は任意であり、例えば、より多くのノードを用意してケーブルによりフレーム転送装置1とそれぞれ接続し、より多くのノード間でフレーム転送装置1を介した通信をできるようにしてもよい。

20

【0022】

以下、他の図や図1の他の部分においても構成要素の数が限定されないことを同様の点線によって表現するものとする。また、本実施形態ではフレーム転送装置1およびノード117-1および117-mは物理的にそれぞれ異なる装置として実施した場合を説明するが、本発明の実施形態はこれに限定されるものではない。例えば、仮想マシンを稼働させる仮想化ソフトウェアを用いて物理的に同一装置内で複数の仮想マシンにより実施してもよい。仮想化ソフトウェアによる仮想化の実施ではケーブル115-1～ケーブル115-mは共有メモリやメモリ間コピー等に置き換えられる。仮想化ソフトウェアを用いた仮想マシンによる実施により、柔軟な接続関係の変更や、ハードウェアの利用効率を高めること等が可能となる。

30

【0023】

続いて、図1を用いてフレーム転送装置1内部のハードウェア構成を説明する。フレーム転送装置1は、演算処理装置101、記憶装置102、入出力装置103、ネットワークインタフェース104-1～ネットワークインタフェース104-nより構成され、共有バス106を通して互いに接続されている。本実施形態において、フレーム転送装置1はネットワークインタフェース104-1～104-nがそれぞれケーブル115-1およびケーブル115-nに接続され、ノード117-1～ノード117-mとイーサネットフレームを交換して、ノード117-1とノード117-mの間の通信を中継する。

40

【0024】

演算処理装置101はCentral Processing Unit (C P U) に代表される演算処理装置であり、記憶装置102上に展開されたプログラム（ソフトウェア）を実行する。なお、演算処理装置101を物理的または論理的に複数のC P Uにより構成してよい。記憶装置

50

１０２は演算処理装置１０１によって実行されるプログラムならびに当該プログラムが実行されることによって作成されるデータまたは読み込まれるデータおよび設定等を格納する。

#### 【００２５】

記憶装置１０２は例えばメモリ、ハードディスク、もしくは光ディスク等のデータやプログラムを保持可能な記憶媒体を含む装置（遠隔地に設置され、ネットワークインタフェース１０４ - １もしくは１０４ -  $n$ 、入出力装置１０３等を通して通信される装置を含む）、またはこれらを組み合わせたもので構成される。記憶装置１０２に格納されるプログラムには、全体制御プログラム１０７（図１においては、プログラムを「ＰＧＭ」と記述、以下同様）、ポート制御プログラム１０８ - １～１０８ -  $n$ が含まれる。また、記憶装置に格納されるデータの一例としてはキュー部１１０、転送ポート管理表３００、トラフィック設定表４００、キュー設定表５００、トラフィック履歴表６００、帯域割当表７００が含まれる。これらの詳細については後述するが、以下の説明においてデータの格納場所が明示されていないデータの保存・保持・記憶・記録・格納等の書き込みや参照・取り出し等の読み出しは記憶装置１０２に対して行うものとする。

10

#### 【００２６】

上記記憶装置１０２のプログラムはフレーム転送装置１を制御するための機能に関して論理的に区別したものであり、例えばすべてを１つのプログラムとして構成して処理させてもよいし、細分化して複数のスレッドやプロセス等に分離して実行したりしてもよい。なお、記憶装置１０２にはフレーム転送装置１を制御するためのプログラムにより作成された他のデータが格納されてもよい。本実施形態では記憶装置１０２にフレーム転送装置１を制御するためのプログラムが展開され、演算処理装置１０１により実行される。

20

#### 【００２７】

入出力装置１０３は、フレーム転送装置１に対して情報の入出力を行う装置である。入出力装置１０３は、例えば、スイッチ、キーボード、マウス、マイクロホン、ビデオカメラ、ディスプレイ、またはスピーカ等の機器、またはフレーム転送装置１にはこれらの機器に接続可能なインタフェースを備えて、フレーム転送装置１にこれらの機器を接続して機能させる形態をとってもよい。また、入出力装置１０３には信号ケーブルまたは電波及び赤外線等の無線を介して行うシリアル通信などの通信によるものも含まれる。

#### 【００２８】

フレーム転送装置１は、入出力装置１０３によって、フレーム転送装置１のユーザまたは管理者からの指示を受けたり、結果を出力したりすることが可能となる。ネットワークインタフェース１０４ - １および１０４ -  $n$ は、イーサネットフレームの送信および受信を行う入出力ポート１１２ - １～１１２ -  $n$ を備え、演算処理装置１０１の指示により、ポート制御プログラム１０８ - １および１０８ -  $n$ によってそれぞれ制御され、他の機器とイーサネットフレームによる通信を行う。以下、これらの入出力ポート１１２と当該入出力ポート１１２を制御するポート制御プログラム１０８をまとめてポートと呼ぶ。

30

#### 【００２９】

また、共有バス１０６はフレーム転送装置１の各構成要素間の通信を行うためのものであるが、本発明は共有バスに限定されるものではない。各構成要素間で必要な通信を行えるようになっていれば、共有バスを用いる以外の方法によって接続されてもよい。例えば、各構成要素間を直接接続することによって、要素間の接続を最適化し、処理に必要な消費電力を減らしたり、処理効率を上げたりすることが可能となる。

40

#### 【００３０】

<システムおよびハードウェアの変形例>

なお、本発明において、フレーム転送装置１の構成要素の個数は図１に限定されるものではなく、例えば入出力装置１０３を２つ用意して冗長化したり、機能分担したりしてもよい。また、本発明は、記憶装置１０２に格納したプログラムに格納したプログラムによって実現される機能のうち、一部または全部の機能をハードウェアとして実装してもよい。ハードウェアとして実装することにより、例えば処理の高速化や低消費電力化が可能と

50

なる。また、仮想マシンのように、当該ハードウェアが備える機能をプログラムとして実現してもよい。ハードウェアが備える機能をプログラムとして実現することによって、例えば、設置スペースを削減したり、管理を単純化したりできる。また、プログラムによって実現される機能の構成は後述する構成に限るものではなく、複数の機能が統合される構成、または1つの機能が複数に分割される構成でもよい。また、各機能によって実行される処理の順序についても、後述する順序に限るものではなく、処理の依存関係が許すならば、並列で同時に実行したり、順序を入れ替えて実行したりしてもよい。例えば、並列実行することによって処理時間を短縮したり、順序を入れ替えることによって待ち時間を減らしたりすることが可能となる。

#### 【0031】

##### <ソフトウェアの構成>

図2はフレーム転送装置1におけるソフトウェア構成を説明するブロック図である。なお、図において、四角はフレーム転送装置1で実行されるプログラムによって実現される処理ブロックを表す。また、処理ブロック間の矢印は、矢印の方向に情報または指示・指令が伝えられることを示す。処理ブロック間で伝えられる情報または指示・指令の内容は以下で説明する。

#### 【0032】

フレーム転送装置1におけるソフトウェアは全体制御プログラム107、ポート制御プログラム108-1~108-n(以下、制御プログラム108-1~108-nのうちのいずれかを表す場合は「ポート制御プログラム108」と記述する。複数ある他のコンポーネントについても同様)、キュー110、入出力ポート112-1~112-nの送受信機能から構成される。なお、入出力ポート112-1~112-nの送受信機能はハードウェアで構成してもよい。フレーム転送装置1全体は、入出力ポート112-1にてイーサネットフレームを受信して、キュー110に格納し、その内容をポート制御プログラム108-1で解析し、解析結果から全体制御プログラム107にて当該イーサネットフレームを転送すべきポートを決定し、キュー110に格納している当該イーサネットフレームを転送すべきと決定されたポートの入出力ポート112-nから出力することによりイーサネットブリッジとして機能する。

#### 【0033】

全体制御プログラム107は設定受付部207、輻輳検出部208、PAUSE送信ポート・時間決定部209、トラフィック設定部210および格納キュー選択部211、トラフィック履歴管理部212より構成され、複数のポートにまたがる可能性のある処理およびフレーム転送装置1全体の制御を行う。

#### 【0034】

ポート制御プログラム108-1~108-nは、割当帯域管理部213、輻輳トラフィック推定部(輻輳推定部)214、送信フレーム選択部215、送信停止タイマー(送信停止部)216、PAUSEフレーム作成部217、キュー制御部218、入出力制御部219およびフレーム解析部220より構成され、当該ポートに関連する入出力ポート112-1の制御およびキュー110の制御を行う。キュー110は論理的に分割されたキューである論理キュー221-1~221-iにより構成され、入出力ポート112-1~112-nから入出力されるイーサネットフレームを一時的に保持する。

#### 【0035】

なお、論理キュー221-1~221-iは、1つのイーサネットフレームについて、イーサネットフレームの内容や入出力されるポート等に基づいて、後述する格納キュー選択部211が分類したトラフィックに対応する。もちろん、論理キュー211の数は図2に示すような3に限定されるものではなく、例えば論理キューの数を1としてすべてのトラフィックのイーサネットフレームを格納してもよいし、論理キューの数を3より大きくしてより細かくトラフィックの種類毎などに分類して格納してもよい。また、論理キュー221-1~221-iの総称を論理キュー221とする。

#### 【0036】

10

20

30

40

50



また、格納キュー選択部 2 1 1 が行うトラフィックの分類は、1つのイーサネットフレームが複数のトラフィックに同時に分類される処理でもよく、その場合、1つのイーサネットフレームが複数の論理キュー 2 2 1 に格納されることとなる。この場合、記憶装置 1 0 2 上の表現としては、複数の論理キュー 2 2 1 に同一のイーサネットフレームをコピーして記憶装置 1 0 2 の制御を簡素化してもよい。あるいは、複数の論理キュー 2 2 1 で同一の記憶装置 1 0 2 を参照するようにしてメモリを節約してもよい。

#### 【0037】

1つのイーサネットフレームが複数のトラフィックに同時に分類されることで、例えばポートに接続されたノードで複数の仮想マシンが稼動する場合や、フレーム転送装置 1 のポートにブリッジが接続され、当該ブリッジに複数のノードが接続された場合にマルチキャストフレームの転送を簡素化できる。

10

#### 【0038】

入出力ポート 1 1 2 - 1 ~ 1 1 2 - n はフレーム送信部 2 2 4 およびフレーム受信部 2 2 5 から構成され、当該ポートに接続されたノード 1 1 7 とイーサネットフレームを送受信する。以下の説明では、当該ポートにおけるノードとの接続をリンクと呼び、リンクを通してイーサネットフレームが当該ノードと送受信できる状態にあるものとして説明する。

#### 【0039】

また、本実施形態では、トラフィックは V L A N I D や優先度情報、送信元や宛先等の通信経路等により区別される、ネットワークを流れる情報を示す。そして、イーサネットフレームまたはフレームは、トラフィックに流すデータの単位を表す。

20

#### 【0040】

<ソフトウェアの構成：事前準備>

続いて、図 2 の各構成要素間で送受信される支持・指令および情報について、フレーム転送装置 1 によるイーサネットフレームの転送を例に説明する。まず、イーサネットフレームの転送に先立って、設定受付部 2 0 7 は入出力装置 1 0 3 またはネットワークインタフェース 1 0 4 - 1 ~ 1 0 4 - n を介してフレーム転送装置 1 のユーザまたは管理者から入力された情報にしたがって、記憶装置 1 0 2 に各種の設定を格納する。

#### 【0041】

具体的には、設定受付部 2 0 7 へ入力された情報にしたがって、トラフィック設定部 2 1 0 からトラフィック設定表 4 0 0 と、キュー設定表 5 0 0 の内容を設定し、設定受付部 2 0 7 へ入力された情報とトラフィック設定部 2 1 0 のキュー設定表 5 0 0 の内容にしたがって、割当帯域管理部 2 1 3 を通して帯域割当管理表 7 0 0 の内容を格納する。

30

#### 【0042】

これらはフレーム転送装置 1 の製造時に予め記憶装置 1 0 2 に格納しておいてもよい。なお、トラフィック設定表 4 0 0、キュー設定表 5 0 0、帯域割当管理表 7 0 0 を含め、以下の説明において用いる表の具体的な内容については後ほど説明する。キュー制御部 2 1 8 はキュー設定表 5 0 0 に従ってキュー 1 1 0 に論理キュー 2 2 1 の作成・削除等の初期設定を指示し、論理キュー 2 2 1 を使用できるように準備する。

#### 【0043】

<ソフトウェアの構成：受信時の処理>

以下、まず、フレーム転送装置 1 で行われるイーサネットフレーム受信時の処理を説明する。イーサネットフレームはフレーム受信部 2 2 5 で受信されると、イーサネットフレームの情報をフレーム解析部 2 2 0 が解析し、フレーム解析部 2 2 0 が P A U S E フレームでない場合は転送すべきフレームとして格納キュー選択部 2 1 1 に解析結果を通知する。

40

#### 【0044】

フレーム解析部 2 2 0 が格納キュー選択部 2 1 1 に通知する解析結果は、送信先 M A C アドレス、V L A N I D、P r i o r i t y 等の、格納キュー選択部 2 1 1 がトラフィックとして分類するために必要な情報である。これらの情報は、送信するイーサネットフ

50

フレームに関連づけられた情報である。

【 0 0 4 5 】

格納キュー選択部 2 1 1 では、フレーム解析部 2 2 0 より通知された解析結果とトラフィック設定部 2 1 0 のトラフィック設定表 4 0 0 を用いて当該イーサネットフレームのトラフィック ID 4 0 5 ( 図 4 参照 ) を決定し、当該トラフィック ID 4 0 5 と当該イーサネットフレームが入力されたポートを表すポート ID 3 0 3 ( 図 3 参照 ) およびトラフィック設定部 2 1 0 のキュー設定表 5 0 0 ( 図 5 参照 ) から当該イーサネットフレームを格納する論理キューの ID 5 0 3 を決定する。そして、格納キュー選択部 2 1 1 がフレーム解析部 2 2 0 より通知された解析結果と転送ポート管理表 3 0 0 を用いて当該イーサネットフレームを転送するポート ID 3 0 3 を決定する。

10

【 0 0 4 6 】

格納キュー選択部 2 1 1 は決定されたポートのキュー制御部 2 1 8 と入出力制御部 2 1 9 に対してフレーム受信部 2 2 5 の当該イーサネットフレームを当該論理キューの ID 5 0 3 に格納するよう指示する。当該指示に従って、入出力制御部 2 1 9 は当該イーサネットフレームを当該論理キューに伝え、キュー制御部 2 1 8 は当該イーサネットフレームを当該論理キュー 2 2 1 に格納する。

【 0 0 4 7 】

また、格納キュー選択部 2 1 1 はトラフィック履歴管理部 2 1 2 に、当該イーサネットフレームのトラフィック ID とポート ID とフレームサイズを通知する。トラフィック履歴管理部 2 1 2 では格納キュー選択部 2 1 1 から通知された情報を現時刻とともにトラフィック履歴表 6 0 0 ( 図 6 参照 ) に保存する。トラフィック履歴表 6 0 0 への保存時に古い時刻のエントリを削除する等によりトラフィック履歴表 6 0 0 が使用するメモリサイズを節約してもよい。

20

【 0 0 4 8 】

以上により、リンクで受信されたイーサネットフレームが転送すべきフレームの場合、当該イーサネットフレームが適切な論理キューに格納される。このとき、格納キュー選択部 2 1 1 は、送信先 MAC アドレス、VLAN ID、Priority 等のフレームの種類に応じてトラフィックを分類し、分類結果に応じた論理キュー 2 2 1 へ当該フレームを格納する。

【 0 0 4 9 】

< ソフトウェアの構成 : 送信時の処理 >

続いて、各ポートの論理キュー 2 2 1 - 1 ~ 2 2 1 - i からイーサネットフレームを送信する処理について述べる。割当帯域管理部 2 1 3 は帯域割当管理表 7 0 0 ( 図 7 参照 ) の内容を輻輳トラフィック推定部 2 1 4 および送信フレーム選択部 2 1 5 に通知する。

30

【 0 0 5 0 】

送信フレーム選択部 2 1 5 は、割当帯域管理部 2 1 3 から通知された帯域割当管理表 7 0 0 の内容および輻輳トラフィック推定部 2 1 4 が通知する制限帯域と、キュー制御部 2 1 8 から通知された各論理キュー 2 2 1 の使用量と、送信停止タイマー 2 1 6 のタイマー値にしたがって、論理キュー ( 例えば 2 2 1 - 1 ) のイーサネットフレームを送信する。送信フレーム選択部 2 1 5 による送信すべきイーサネットフレームの選択方法については送信帯域計算処理として図 1 0 A、図 1 0 B で後述する。

40

【 0 0 5 1 】

送信フレーム選択部 2 1 5 による論理キュー 2 2 1 - 1 等からのイーサネットフレームの送信は、具体的には、格納キュー選択部 2 1 1 がキュー制御部 2 1 8 と入出力制御部 2 1 9 に対して論理キュー 2 2 1 - 1 のイーサネットフレームを読み出してリンクに送信するよう指示する。当該指示に従って、キュー制御部 2 1 8 は当該イーサネットフレームを当該論理キューから取り出してフレーム送信部 2 2 4 に送るよう論理キュー 2 2 1 - 1 に指示し、入出力制御部 2 1 9 はフレーム送信部 2 2 4 に論理キュー 2 2 1 - 1 から送られた当該イーサネットフレームを入出力ポート 1 1 2 - 1 に接続されたリンクに対して送信するように指示する。

50

## 【 0 0 5 2 】

また、入出力制御部 2 1 9 は送信したイーサネットフレームの情報を、具体的には属する論理キュー 2 2 1 - 1 ~ 2 2 1 - i のキュー ID と当該イーサネットフレームのフレームサイズを輻輳トラフィック推定部 2 1 4 に通知する。輻輳トラフィック推定部 2 1 4 は、フレーム解析部 2 2 0 から通知された P r i o r i t y 毎の P A U S E フレーム受信状況、割当帯域管理部 2 1 3 から通知された帯域割当管理表の内容、キュー制御部 2 1 8 から通知された論理キュー 2 2 1 - 1 ~ 2 2 1 - i の使用量および入出力制御部 2 1 9 から通知された送信するイーサネットフレームの情報と、送信停止タイマー 2 1 6 を参照して得た P r i o r i t y 毎の変数である P A U S E 状態に基づいて、各論理キュー 2 2 1 からイーサネットフレームの制限帯域を計算して送信フレーム選択部 2 1 5 に通知する。

10

## 【 0 0 5 3 】

輻輳トラフィック推定部 2 1 4 における制限帯域の計算方法については制限帯域計算処理として図 8 で後述する。以上により、論理キューからイーサネットフレームがリンクに送信される。

## 【 0 0 5 4 】

なお、イーサネットフレームの転送（ブリッジ）については、I E E E 8 0 2 . 1 D および I E E E 8 0 2 . 1 Q として標準化されており、これらの技術を適宜用いることができる。

## 【 0 0 5 5 】

< ソフトウェアの構成： P A U S E 受信 >

20

続いて、図 2 における P A U S E フレームの扱いについて述べる。P A U S E フレームを受信した場合および送信する場合のフレーム転送装置 1 における処理について説明する。まず、P A U S E フレームを受信した時のフレーム転送装置 1 が実行する処理について説明する。イーサネットフレーム受信時は、フレーム解析部 2 2 0 の解析により、当該イーサネットフレームが P A U S E フレームであることが検出される。

## 【 0 0 5 6 】

P A U S E フレームは転送すべきフレームではないのでフレーム解析部 2 2 0 は格納キュー選択部 2 1 1 に通知せず、P A U S E フレームの内容にしたがって各 P r i o r i t y で一時停止すべき時間を計算し、送信停止タイマー 2 1 6 に指示する。そして、フレーム解析部 2 2 0 は輻輳トラフィック推定部 2 1 4 に対して P r i o r i t y 毎に P A U S E フレームの受信状況を通知する。

30

## 【 0 0 5 7 】

具体的には、フレーム解析部 2 2 0 は輻輳トラフィック推定部 2 1 4 に P A U S E フレームにより指示された時間を通知する。送信停止タイマー 2 1 6 は格納キュー選択部 2 1 1 より指示された P r i o r i t y 毎の時間にタイマー値をセットする。送信停止タイマー 2 1 6 は一定時間毎に設定された各タイマー値を 0 になるまで減じる。送信停止タイマー 2 1 6 は、P A U S E フレームにより指示された時間が P A U S E 状態の解除を意味する 0 の場合は 0 をセットする。送信フレーム選択部 2 1 5 は送信停止タイマー 2 1 6 を参照し、特定の P r i o r i t y の送信タイマーが 0 より大きい場合は当該 P r i o r i t y に属するイーサネットフレームを含む論理キュー 2 2 1 からの送信を停止する。

40

## 【 0 0 5 8 】

以上により、P A U S E フレームに指定された時間だけ、当該リンクに対するフレーム転送装置 1 からのフレーム送信が入出力ポート 1 1 2 毎に一時停止され、P A U S E フレームによる輻輳制御が実現される。

## 【 0 0 5 9 】

< ソフトウェアの構成： P A U S E 送信 >

次に、フレーム転送装置 1 が P A U S E フレームをリンクへ送信する場合の処理について説明する。P A U S E フレームの送信は当該リンクが輻輳状態となり、ノード 1 1 7 から送信されたフレームをフレーム転送装置 1 で正しく受信できない恐れがある場合に行う。本実施形態においては、フレーム転送装置 1 のポート制御プログラム 1 0 8 が受信する

50

フレームを格納する論理キュー 221 の使用量が停止閾値（例えば、確保したキュー容量の 8 割）を超えた場合に、ノード 117 から送信されたフレームを受信できない恐れがあると判定して、当該論理キュー 221 に格納されるフレームを送信しているリンクに送信を停止させる P A U S E フレームを送信する。

【0060】

フレーム転送装置 1 は、当該論理キュー 221 の使用量が再開閾値（例えば、確保したキュー容量の 5 割）を下回った場合に受信できない恐れが消滅したと判定して、当該リンクへの送信を再開させる P A U S E 解除フレームを接続されているリンクに送信する。

【0061】

図 2 において、キュー制御部 218 は論理キュー（例えば、221 - 1）の使用量を輻輳検出部 208 に通知する。輻輳検出部 208 はキュー制御部 218 から通知された論理キューの使用量が停止閾値を越えたもしくは再開閾値を下回ったことを検知すると、P A U S E 送信ポート・時間決定部 209 に通知する。

【0062】

P A U S E 送信ポート・時間決定部 209 は、トラフィック履歴管理部 212 のトラフィック履歴表 600 とトラフィック設定部 210 のトラフィック設定表 400、キュー設定表 500 を参照して、当該論理キュー 221 に格納されるイーサネットフレームが入力されている（1 または複数の）ポートを検索する。P A U S E 送信ポート・時間決定部 209 は検索したポートの P A U S E フレーム作成部 217 に、当該論理キュー 221 の P r i o r i t y と、当該 P r i o r i t y での送信を一時停止させる時間を通知する。なお、当該論理キュー 221 に複数の P r i o r i t y に属するフレームを格納するなど、当該論理キューの P r i o r i t y が一意に定まらない場合は、複数の P r i o r i t y に関して同様の検索を P A U S E 送信ポート・時間決定部 209 が行う。また、前記検索を行う代わりにすべてのポートや当該論理キューの所属する V L A N に参加しているポートを対象として、検索に必要な時間を節約してもよい。

【0063】

送信を一時停止させる時間については、停止閾値を越えた場合は最大値である 65535（単位は 512 ビットを送信する時間）とし、再開閾値を下回った場合は 0 とする。本発明では、停止閾値を越えた場合の送信を一時停止させる時間が前記の 65535 に限るものではなく、例えば論理キュー 221 の空き容量や当該論理キューからの送出量、演算処理装置 101 の使用量等に基づいてさらに細かく時間設定をして送信再開 P A U S E フレーム（または P A U S E 解除フレーム）の送出を抑制するなどしてもよい。

【0064】

P A U S E フレーム作成部 217 は、P A U S E 送信ポート・時間決定部 209 からの指示を、当該 P r i o r i t y の時間として埋め込んだ P A U S E フレームを作成し、フレーム送信部 224 に通知する。

【0065】

P A U S E フレーム作成部 217 では、P A U S E 送信ポート・時間決定部 209 から指示された時間を P r i o r i t y 毎に記憶し、送信停止タイマー 216 と同様に 0 になるまで一定時間毎に前記指示され時間から 1 を減じ、各 P r i o r i t y での送信を一時停止させる時間の現在の値を参照可能にする。また、P A U S E フレーム作成部 217 では、P A U S E 送信ポート・時間決定部 209 から指示された P r i o r i t y 以外の P r i o r i t y に関しても一時停止させる時間を埋め込んだ P A U S E フレームを作成できるようにする。

【0066】

入出力ポート 112 のフレーム送信部 224 では、ポート制御プログラム 108 の P A U S E フレーム作成部 217 から通知された P A U S E フレームを、接続されたリンクに対して送出する。なお、P A U S E フレームについて直接フレーム送信部 224 による送出は行わず、いったん P A U S E フレーム用に用意した論理キュー 221 に格納し、送信フレーム選択部 215 で他の論理キュー 221 より優先してフレームを送出するという方

10

20

30

40

50

法にし、送出されるフレームの管理を一元化してもよい。

【0067】

以上により、当該リンクでのフレーム転送装置1によるフレームの受信が一時停止され、P A U S Eフレームによる輻輳制御が実現される。

【0068】

なお、P A U S Eフレームのフォーマットや一時停止時間等については、I E E E 802.3xおよびI E E E 802.3bd、I E E E 802.1Qbbとして標準化されており、本実施形態による説明があれば当業者によるP A U S E機能の実施は困難ではない。

【0069】

<ソフトウェアの構成のバリエーション>

図2においては、受信フレームをそのまま論理キュー221に格納する形での実施を説明したが、本発明はこれに限るものではなく、例えば入力バッファ、出力バッファのようにバッファを記憶装置102に設置してイーサネットフレームを格納し、格納キュー選択部211による論理キューの選択や送信フレーム選択部215による送信フレーム選択の処理を平準化してもよい。また、フレーム解析部220においてP A U S Eフレームでない場合は転送すべきフレームとしたが、例えばP A U S Eフレーム以外のフレームであっても転送する必要のないL L D Pフレームの取り扱いや、V L A Nタグの付加または除去等、一般にイーサネットスイッチやL2スイッチと呼ばれるネットワーク機器で行われる拡張を施してもよい。また、本実施形態ではトラフィックの分類に、送信元M A Cアドレス、送信先M A Cアドレス、V L A N I D、P r i o r i t yを用いたが、本発明におけるトラフィックの分類はこれに限定されるものではなく、例えば、入力されたポートや、イーサネットフレームに含まれる情報であるE t h e r T y p e、I P ( I n t e r n e t P r o t o c o l ) アドレス、T C Pのポート番号等を加え、各情報単独または複数の情報を組み合わせてトラフィックを決定する形で用いてもよい。たとえばT C Pのポート番号やI Pアドレスの情報をを用いることで、特定アプリケーションのトラフィックを区別して輻輳制御を行うことができるようになる。

【0070】

<表の説明>

続いて、図2の説明に用いた表である、転送ポート管理表300、トラフィック設定表400、キュー設定表500、トラフィック履歴表600、帯域割当管理表700の具体的な内容について例を用いて説明する。なお、これらの表は記憶装置102に保存される。

【0071】

図3に転送ポート管理表の例を示す。転送ポート管理表300は、M A Cアドレス301、V L A N I D302、ポートI D303を含むエントリで構成される。M A Cアドレス301はイーサネットにおけるエンティティを識別するためのアドレスである。V L A N I D302は、I E E E 802.1Qにて規定される、イーサネットフレームにおけるV L A Nの識別子で、異なるV L A N I Dを持つフレームを別のV L A Nとして扱うために用いられるものである。ポートI D303は複数あるポート（入出力ポートおよび対応するポート制御プログラムを表す）のうちの1つを識別するための識別子である。

【0072】

格納キュー選択部211では転送ポート管理表300を用いて、フレーム解析部220から通知された解析結果のうちの送信先M A CアドレスがM A Cアドレス301で、同じくV L A N I DがV L A N I D302と両方とも一致するエントリを検索し、該当するエントリのポートI D303を対象として決定する。格納キュー選択部211では転送ポート管理表300に前記一致するエントリが複数あった場合は複数のポートを対象として決定する。一方、前記一致するエントリが見つからなかった場合はV L A N I Dが一致する全ポートを対象とする。

【0073】

10

20

30

40

50

転送ポート管理表 300 は IEEE 802.1Q における Filtering Database に対応する。転送ポート管理表 300 はフレームの情報と照合して当該フレームを転送すべきポートを決定する機能を持てば、必ずしも図 3 に挙げた情報に基づかなくてもかまわない。例えば、入力ポートを識別するポート ID と転送先のポート ID 303 の対応表といった構成も考えられ、VLAN ID を省いた構成 (IEEE 802.1D に規定されるブリッジ機能) としてもよい。また、転送ポート管理表 300 は入力されたフレームの送信元 MAC アドレスと VLAN ID およびポート ID の組み合わせをエントリとして追加し、一定時間使用しなかったエントリを削除する MAC 学習により保守する。なお、本発明においてはこれに限るものではなく、例えば、固定的なエントリを作成し、MAC 学習の処理を省いてもよい。

10

#### 【0074】

図 4 にトラフィック設定表の例を示す。トラフィック設定表 400 は、送信元アドレス 401、送信先アドレス 402、VLAN ID 403、Priority 404、トラフィック ID 405 を含むエントリで構成される。送信元アドレス 401 は送信元の MAC アドレス、送信先アドレス 402 は送信先の MAC アドレスを格納する。VLAN ID 403 および Priority 404 は IEEE 802.1Q で規定される VLAN ID および PCP の値をそれぞれ表す。トラフィック ID 405 は論理的に区別して扱うトラフィックの識別子を格納する。

#### 【0075】

格納キュー選択部 211 ではトラフィック設定表 400 を用いて、フレーム解析部 220 から通知された解析結果のうちの対象フレームの送信先 MAC アドレスと送信先アドレス 402 を、送信元 MAC アドレスと送信元アドレス 401 を、VLAN ID と VLAN ID 403 を、Priority と Priority 404 をそれぞれ照合し、すべて一致したエントリのトラフィック ID 405 を取得し、転送ポート管理表 300 およびキュー設定表 500 と併せて対象フレームの格納すべき論理キューを決定する。

20

#### 【0076】

したがって、トラフィック設定表 400 ではトラフィックの分類に用いる情報と、一致した場合のトラフィック ID を対応付ける役割を担う。そのため、先に例示したようにトラフィックの分類に用いる情報は例えば Ether Type などでもよく、本発明ではトラフィック設定表 400 に含まれる情報は図 4 に限定されるものではない。一致するエントリが見つからず対象トラフィックが確定しない場合は、あらかじめ定められたトラフィック ID (例えば 0) を用いる。

30

#### 【0077】

図 5 にキュー設定表 500 の例を示す。キュー設定表 500 はトラフィック ID 501、ポート ID 502、キュー ID 503 を含むエントリで構成される。トラフィック ID 501 は図 4 に示したトラフィック設定表 400 のトラフィック ID 405 と同様、論理的に区別して扱うトラフィックの識別子を格納する。ポート ID 502 は図 3 の転送ポート管理表 300 のポート ID 303 と同様にポートを区別するための識別子を格納する。キュー ID 503 はキュー 110 における論理キュー 221-1 ~ 221-i を区別するための識別子を表す。格納キュー選択部 211 では転送ポート管理表 300 およびトラフィック設定表 400 の結果をキュー設定表 500 と照合し、トラフィック ID 501 とトラフィック ID 405、および、ポート ID 502 とポート ID 303 の双方が一致するエントリを検索し、対象フレームを格納する論理キューを検索されたエントリのキュー ID 503 から決定する。

40

#### 【0078】

格納キュー選択部 211 ではフレーム解析部 220 から通知された解析結果と、転送ポート管理表 300、トラフィック設定表 400 およびキュー設定表 500 を用いて、対象フレームを格納する論理キュー 221 を決定する。そのため、これら 3 種の表は必ずしも別に保存される必要はなく、フレーム解析部 220 から通知された解析結果から対象フレームを格納する論理キュー 221 を決定できればよい。例えば、転送ポート管理表 300

50

、トラフィック設定表 4 0 0 およびキュー設定表 5 0 0 を 1 つの表としてまとめることでデータ構造を簡素化することができる。

【 0 0 7 9 】

図 6 にトラフィック履歴表の例を示す。トラフィック履歴表 6 0 0 は、トラフィック ID 6 0 1、ポート ID 6 0 2、時刻 6 0 3、フレームサイズ 6 0 4 を含むエントリで構成される。トラフィック ID 6 0 1 は図 4、図 5 のトラフィック ID 4 0 5、トラフィック ID 5 0 1 と同様、論理的に区別して扱うトラフィックの識別子を格納する。ポート ID 6 0 2 は図 3、図 5 のポート ID 3 0 3 およびポート ID 5 0 2 と同様、ポートを区別するための識別子を格納する。時刻 6 0 3 は当該エントリが作成された時刻を表す値で、例えばある時点からの経過秒数を用いることができるがこれに限定されるものではない。例えば、一定周期でカウントアップされるカウンタを用いることもでき、この場合、ハードウェア実装の簡素化が期待できる。フレームサイズ 6 0 4 は対象とするイーサネットフレームのフレームサイズを表す。フレームサイズ 6 0 4 として格納する情報はフレームサイズに限定するものではなく、大まかなフレーム送信量を把握できればよい。例えば、フレームサイズ 6 0 4 として 1 を格納して処理を単純化してもよい。

10

【 0 0 8 0 】

P A U S E 送信ポート・時間決定部 2 0 9 では、トラフィック履歴表 6 0 0 を用いて対象トラフィックが入力されているポートを決定する。そのため、トラフィック履歴表 6 0 0 は対象トラフィックの入力ポートを調査可能であればよく、図 6 に示す形式に限定されるものではなく、例えば、各トラフィック ID と入力があったポート番号といった形式でもよい。形式を変更することでトラフィック履歴表 6 0 0 に必要なメモリ量を削減したりアクセス時間を削減したりできる。またエントリは追加のみでなく、例えば時刻 6 0 3 を現時刻と比較して一定時間経過したもの、フレームサイズ 6 0 4 が一定値より小さいもの等を削除したり参照時に除外したりしてもよい。このような対象エントリの限定により、例えば論理キュー 2 2 1 の輻輳に与える影響の大きなポートからのフレーム入力を抑制して他の通信に関しては制限せずに通すといった柔軟な輻輳制御がイーサネットフレーム受信時に可能となる。

20

【 0 0 8 1 】

図 7 に帯域割当管理表の例を示す。帯域割当管理表 7 0 0 はキュー ID 7 0 1 と割当帯域 7 0 2 を含むエントリで構成される。キュー ID 7 0 1 は図 5 のキュー ID 5 0 3 と同様、キュー 1 1 0 における論理キュー 2 2 1 を区別するための識別子を格納する。割当帯域 7 0 2 は当該論理キュー 2 2 1 からリンクに対して送信する保証帯域を格納する。割当帯域 7 0 2 は絶対的な送信帯域のほか、送信帯域の比のような形で指定してもよい。帯域割当管理表 7 0 0 を用いることで各論理キュー 2 2 1 からリンクに対して送信する際に保証すべき帯域を参照することが可能となる。

30

【 0 0 8 2 】

以上に説明した表のデータ構造は必ずしも表である必要はなく、例えばハッシュや 2 分木等でもよい。ハッシュを用いることで一致するエントリの検索時間短縮が期待できる。また、それぞれの照合は完全一致に限るものではなく、例えば部分一致等、与えられた情報に基づいて対象とするエントリを決定できればよい。例えば部分一致を用いることにより、表に必要なメモリ容量を削減することができる。

40

【 0 0 8 3 】

続いて、輻輳トラフィック推定部 2 1 4 における制限帯域計算処理および送信フレーム選択部 2 1 5 における送信帯域計算処理についてフローチャートを用いて説明する。なお以下の説明において、変数に格納された値は当該変数名により表すこととする。

【 0 0 8 4 】

< 輻輳トラフィック推定部の処理 >

図 8 は輻輳トラフィック推定部 2 1 4 で行う制限帯域計算処理の一例を説明するフローチャートである。輻輳トラフィック推定部 2 1 4 は、論理キュー 2 2 1 から読み出されてリンクへ送信されるイーサネットフレームについて、キュー ID によりトラフィックを区

50

別したうえで各トラフィックの送信帯域を計算する。そして、輻輳を引き起こすトラフィックを推定し、当該トラフィックが輻輳の要因と推定される場合には当該トラフィックの送信帯域を減らす。

【0085】

このため、輻輳トラフィック推定部214は、フレーム解析部220から通知されたPAUSEフレーム受信状況と、割当帯域管理部213から通知された帯域割当管理表の内容と、キュー制御部218から通知された論理キューの使用量と、入出力制御部219から通知された送信したイーサネットフレームの情報と、送信停止タイマー216を参照して得たPriority毎のPAUSE状態(タイマー値が0なら非PAUSE状態、それ以外ならPAUSE状態)を用いて、各トラフィックで送信可能な制限帯域を制限帯域計算処理により計算して送信フレーム選択部215に通知する。

10

【0086】

制限帯域計算処理は一定時間毎(例えば0.01秒毎)に繰り返して行う。なお、本発明では制限帯域計算処理実行の契機は一定周期に限るものではなく、例えば、PAUSEフレームの受信を制限帯域計算処理実行の契機としてもよいし、時間間隔をトラフィックの流量等に応じて変更してもよい。なお、以下の説明は、PFCの1つのPriorityに対する処理を対象としており、当該PriorityのPAUSEフレームにより当該Priorityに属する全トラフィックの送信が停止されるものとして説明する。これはPFCではない一般のPAUSEで当該リンクを流れる全トラフィックの送信が停止されるのと同様である。なお、PFCにおける複数のPriorityに対応するには、制限帯域計算処理をPriority毎に当該Priorityに属するトラフィックを対象として行う。以下の説明中で「トラフィック」は対象Priorityに属するトラフィックを、「キュー」は当該トラフィックのイーサネットフレームが格納されている論理キュー221をそれぞれ表す。また「保証帯域」は帯域割当管理表の内容である各トラフィックの割当帯域702を表す変数である。「累積送信量」は入出力制御部219から通知された送信済みのイーサネットフレームの情報を用いてトラフィック毎にこれまで転送したフレームサイズの合計値を保存する変数である。すなわち、各トラフィックの累積送信量は入出力制御部219から通知されるたびに加算していく。また変数「制限帯域」、「前回の制限帯域」の初期値は変数「保証帯域」の値とする。

20

【0087】

輻輳トラフィック推定部214は、制限帯域計算処理が開始状態S801から開始されると、ステップS802に進む。

30

【0088】

ステップS802では輻輳トラフィック推定部214が送信停止タイマー216を参照して取得したPAUSE状態であるか否かを判定し、真(PAUSE中)であればステップS803に進み、偽ならステップS804に進む。ステップS803では輻輳トラフィック推定部214が変数「PAUSE状態」を参照し、前回PAUSE状態ではなく今回PAUSE状態となった、すなわちPAUSEが開始されたか否かを判定し、真ならステップS805に進み、偽であればステップS818に進む。

【0089】

ステップS805では、輻輳トラフィック推定部214が変数「PAUSE間隔」を変数「前回のPAUSE間隔」に保存し、変数「PAUSE開始時刻」を変数「前回のPAUSE時刻」に保存し、変数「PAUSE開始時刻」に現在の時刻を代入し、変数「PAUSE間隔」にPAUSE開始時刻から前回のPAUSE時刻を減じた時間を代入し、ステップS806に進む。

40

【0090】

ステップS806では、輻輳トラフィック推定部214が各トラフィックについて変数「累積転送量」の値を変数「前回の累積転送量」に代入し、変数「累積送信量」の値を変数「累積転送量」に代入し、ステップS807に進む。ステップS807では、輻輳トラフィック推定部214が「前回のPAUSE間隔」および「PAUSE間隔」が有意であ

50



ることを判定し、真ならステップS 8 0 8に進み、偽ならステップS 8 1 7に進む。なお、ステップS 8 0 7における「有意」は、値が、変数が未定義値でなく、かつ、一定値（例えば0.05秒）より大きいことを表す。これは、P A U S E 間隔が初期状態であったり、誤差といえるほど小さかったりした場合に後述する制限帯域補正処理を行わないようにするためである。

【0091】

ステップS 8 0 8では、輻輳トラフィック推定部214が変数「帯域制限トラフィック数」と「帯域制限強化トラフィック数」にそれぞれ0を代入し、ステップS 8 0 9に進む。ステップS 8 0 9では、各トラフィックについて、ステップS 8 1 0からステップS 8 1 1までを繰り返す。

10

【0092】

ステップS 8 1 0では、輻輳トラフィック推定部214が後述する図9の制限帯域補正処理を実行して、ステップS 8 1 1に進む。ステップS 8 1 1では、輻輳トラフィック推定部214がステップS 8 0 9からの繰り返し処理を各トラフィックについて終了したら、ステップS 8 1 2に進む。

【0093】

ステップS 8 1 2では輻輳トラフィック推定部214が変数「帯域制限トラフィック数」が正であるか否かを判定し、真ならステップS 8 1 3に進み、偽ならステップS 8 1 7に進む。ステップS 8 1 3では輻輳トラフィック推定部214が変数「帯域制限強化トラフィック数」が0であるか否かを判定し、真ならステップS 8 1 4に進み、偽ならステップS 8 1 7に進む。

20

【0094】

ステップS 8 1 4では、輻輳トラフィック推定部214がトラフィックを選択し、帯域制限を実施し、ステップS 8 1 7に進む。このステップS 8 1 4におけるトラフィック選択の例として、輻輳トラフィック推定部214が、変数「トラフィック番号」（初期値は0）を参照し、対応するトラフィックの変数「前回の送信量」が正なら当該トラフィックを選択し、そうでなければトラフィック番号に1を加えて次のトラフィックを調べる。そして、輻輳トラフィック推定部214は、トラフィック番号がトラフィックの数以上となった場合は0に戻す。本発明においてステップS 8 1 4におけるトラフィック選択は前述の方法に限定されるものではなく、例えばランダムにトラフィックを選択するようにして処理を単純化してもよい。また、複数のトラフィックを順番に選択してもよい。さらに、例えば送信履歴に基づいて、直近にイーサネットフレームを多く送信しているトラフィックから優先して選択するようにしてもよい。また、ステップS 8 1 4における帯域制限は、選択したトラフィックの変数「新しい制限帯域」に、例えば、

30

$(2 \times \text{保証帯域} + \text{実効帯域}) / 3$

を代入する。各トラフィックの実効帯域は後述の制限帯域補正処理で計算される。なお、ステップS 8 1 4における「新しい制限帯域」の計算式は、制限帯域を保証帯域から実効帯域に近づけるためのものであり、本発明はこれに限るものではない。例えば、より保証帯域に近づけるようにして、帯域制限を緩やかにしてもよい。

【0095】

40

次に、ステップS 8 0 2の判定でP A U S E中でない場合のステップS 8 0 4では、輻輳トラフィック推定部214が、現在の時刻と変数「P A U S E 開始時刻」の差が一定時間（例えば変数「P A U S E 間隔」の2倍）より大きいか否かを判定し、真ならステップS 8 1 5に進み、偽ならステップS 8 1 8に進む。

【0096】

ステップS 8 1 5では、輻輳トラフィック推定部214が帯域制限を緩和して、ステップS 8 1 6に進む。ここで、ステップS 8 1 5における帯域制限の緩和は、各トラフィックについて、「新しい制限帯域」として、例えば、

$(\text{前回の制限帯域} + \text{保証帯域}) / 2$

を代入する。なお、ステップS 8 1 5における「新しい制限帯域」の計算式は、制限帯域

50

を保証帯域に近づけるものであり、本発明はこれに限定されるものではない。例えば、より保証帯域に近づけるようにして、帯域制限の緩和幅を大きくしてもよい。

【0097】

次に、ステップS816では、輻輳トラフィック推定部214が変数「PAUSE間隔」を大きくしてステップS817に進む。ここで、ステップS816はステップS815が実行される頻度を減らすために行うもので、帯域制限の緩和を緩やかにする効果がある。変数「PAUSE間隔」を大きくする例としては、PAUSE間隔に

(PAUSE間隔×2)

を代入するなどが考えられるが、本発明は必ずしもこれに限定されるものではなく、例えば定数値を足すなどでもよい。また、ステップS816を省略して処理を簡素化してもよい。

10

【0098】

次に、ステップS817では、輻輳トラフィック推定部214が各トラフィックについて、変数「前回の制限帯域」に変数「制限帯域」の値を代入し、変数「制限帯域」に「新しい制限帯域」の値を代入し、ステップS818に進む。ステップS818では輻輳トラフィック推定部214が送信停止タイマー216を参照して得たPAUSE状態を変数「PAUSE状態」に保存し、終了状態819に進む。終了状態819では制限帯域計算処理を終了する。

【0099】

以上の処理によって、輻輳トラフィック推定部214は、論理キュー221から読み出されてリンクへ送信されるイーサネットフレームについて、キューIDに基づいてトラフィックを区別したうえで各トラフィックの送信帯域を計算する。そして、輻輳を引き起こすトラフィックを推定し、当該推定されたトラフィックの送信帯域を減らして制限することが可能となる。そして、PAUSEを解除する際には、帯域制限を徐々に緩和することで保証帯域に戻すことができる。

20

【0100】

図9は輻輳トラフィック推定部214で行う制限帯域計算処理のステップS810で実行される制限帯域補正処理の一例を示すフローチャートである。

【0101】

制限帯域補正処理では輻輳トラフィック推定部214が、1つのトラフィックに関して制限帯域の値を補正する。したがって、制限帯域補正処理における変数はトラフィック毎に値を保持する。

30

【0102】

輻輳トラフィック推定部214は制限帯域補正処理が開始状態S901から開始されると、ステップS902に進む。ステップS902では変数「新しい制限帯域」に変数「保証帯域」の値を代入してステップS903に進む。ステップS903では、輻輳トラフィック推定部214が、変数「前回の送信量」に変数「送信量」の値を代入し、変数「送信量」に変数「累積転送量」から変数「前回の累積転送量」を減じたものを代入し、ステップS904に進む。

【0103】

ステップS904では、輻輳トラフィック推定部214は、変数「実効帯域」に(送信量/PAUSE間隔)を代入し、ステップS905に進む。ステップS905では、輻輳トラフィック推定部214が、「保証帯域」より「制限帯域」が小さい、すなわち帯域制限されているか否かを判定し、真ならステップS906に進み、偽なら終了状態913に進む。

40

【0104】

ステップS906では、輻輳トラフィック推定部214が、変数「前回の制限帯域」より「制限帯域」が小さい、すなわち帯域制限を前回強化したか否かを判定し、真ならステップS907に進み、偽ならステップS911に進む。

【0105】

50

ステップS907では、輻輳トラフィック推定部214が、「前回のPAUSE間隔」より「PAUSE間隔」が大きい、すなわちPAUSE間隔が拡大しているか否かを判定し、真ならステップS908に進み、偽ならステップS909に進む。ステップS908では、輻輳トラフィック推定部214が帯域制限を強化してステップS910に進む。ステップS908における帯域制限の強化は、例えば、「新しい制限帯域」に、

( 前回の制限帯域 + 2 × 実効帯域 ) / 3

を代入することで行う。ただし、この計算式は、制限帯域を小さくするためのものであり、本発明ではこれに限らない。例えば、「新しい制限帯域」に実効帯域を代入することとして、帯域制限の強化の幅を大きくしてもよい。

【0106】

ステップS910では輻輳トラフィック推定部214が、変数「帯域制限強化トラフィック数」を1増やし、ステップS912に進む。一方、ステップS909では、輻輳トラフィック推定部214は帯域制限を緩和してステップS911に進む。ステップS909における帯域制限の緩和は、例えば「新しい制限帯域」に、

( 2 × 前回の制限帯域 - 実効帯域 )

を代入することで行う。ただし、この計算式は、制限帯域を大きくするためのものであり、本発明ではこれに限らない。例えば、「新しい制限帯域」に保証帯域を代入することとしてもよい。

【0107】

ステップS911では輻輳トラフィック推定部214が、帯域制限を緩和してステップS912に進む。ここで、ステップS911における帯域制限の緩和は、例えば「新しい制限帯域」に、

( 制限帯域 × 1.01 )

を代入することで行う。ただし、ただし、この計算式は、制限帯域を大きくするためのものであり、本発明ではこれに限らない。例えば、「新しい制限帯域」に保証帯域を代入することとしてもよい。また、ステップ909を省略し、ステップ911にて制限帯域を大きくする処理を行ってステップ数を減らしてもよい。ステップS912では輻輳トラフィック推定部214が変数「帯域制限トラフィック数」を1増やし、終了状態913に進む。終了状態913では制限帯域補正処理を終了する。

【0108】

なお、図8に示した制限帯域計算処理において、各トラフィックの制限帯域が保証帯域より大きい場合は制限帯域に保証帯域を代入し、また制限帯域が実効帯域より小さい場合は制限帯域に実効帯域を代入することで、制限帯域の範囲を限定する。

【0109】

なお、本発明では、制限帯域の下限は実効帯域に限定されるものではなく、例えば制限帯域を0として対象トラフィックの送信を停止してもよい。また、帯域制限緩和の計算時のみ保証帯域を大きくしたり、帯域制限の計算時のみ実効帯域を小さくしたりして制限帯域の変化を大きくしてもよい。

【0110】

以上に説明した制限帯域計算処理により、輻輳トラフィック推定部214はトラフィック毎の制限帯域を計算して送信フレーム選択部215に通知し、フレーム転送装置1はトラフィック毎の送信帯域を制限する。

【0111】

制限帯域計算処理の要点は以下のようなものである。PAUSE開始時は輻輳が発生したことを意味するので、図8に示した制限帯域計算処理で帯域制限が強化されたトラフィックがなければ(ステップS813が真)、適当なトラフィックを選択して帯域制限を行う(ステップS814)。そして、図9に示した制限帯域補正処理では前記帯域制限の結果、PAUSE間隔が拡大していれば(ステップS907が真)、当該帯域制限が輻輳緩和に効果があったとみなして当該帯域制限を強化し(ステップS908)、拡大していなければ当該帯域制限の効果がないとみなして当該帯域制限を緩和する(ステップS909

10

20

30

40

50

）。

【 0 1 1 2 】

なお、P A U S Eにより一時停止している時間も含めた実効的な送信帯域（ステップS 9 0 4の実効帯域）は、当該トラフィックで使用できた帯域の実績とみなせるので、実効帯域を下限として帯域制限を行ってもよい。また、P A U S Eがしばらく来ない場合（ステップS 8 0 4が真）や帯域制限の強さが変化していない場合（ステップS 9 0 6が偽）、帯域制限をしていたトラフィックで利用可能な帯域が変化した可能性も考慮して輻輳状態が緩和されたとみなして帯域制限を緩和してもよい（ステップS 8 1 5およびステップS 9 1 1）。

【 0 1 1 3 】

すなわち、本発明では、P A U S Eでフロー制御されるリンクについて、P A U S Eフレーム受信側（すなわち制御対象のフローの送信側）において、P A U S Eの間隔の変化から輻輳状態の変化（間隔拡大は輻輳緩和、間隔縮小は輻輳激化）を検出する。そして、当該リンクに送信するトラフィックのうち、当該トラフィックの送信帯域制限（送信の抑制）により輻輳緩和（P A U S E間隔の拡大）が検出される場合に、当該トラフィックが輻輳の原因であると推定する。また、P A U S Eフレーム受信時に、適当なトラフィックを選択して送信を抑制してみて輻輳が緩和されるか否かを試し、輻輳が緩和されれば送信抑制を強化する。逆に輻輳が緩和されなければ送信抑制を緩和して他のトラフィックの送信を抑制してみるという試行錯誤により輻輳の原因と推定されるトラフィックを探索する。そして輻輳の原因であると推定されたトラフィックの送信を抑制することで、当該リンクのP A U S E発生を抑制して通信速度を改善する。本発明では制限帯域計算処理が輻輳制御の対象となるトラフィック数の制限要因を含まないため、必要に応じて対象トラフィック数を増減することができる。また、輻輳の原因と推定されるトラフィックの探索では、2分木探索にならって複数のトラフィックを選択して送信抑制して輻輳の状態変化を確認し、輻輳緩和が確認された複数のトラフィックの中でさらに一部トラフィックを選択して送信抑制を行うような、段階的な探索を行ってもよい。本発明の実施は、2分木探索のほか、リスト探索、木探索、グラフ探索など、一般に探索アルゴリズム、検索アルゴリズムと呼ばれる考え方をを用いても行うことができる。なお、本発明における前述した送信の抑止は、一部トラフィック以外のトラフィックの送信促進により、相対的に当該一部トラフィックの送信を抑制してもよい。また本発明は送信の抑制による輻輳の抑制に限るものではなく、一部トラフィックの帯域変更による輻輳状態の変化に基づいて、目的とする輻輳状態、例えば、あらかじめ定めたP A U S Eフレームの受信間隔となるように各トラフィックの送信帯域を変更してもよい。

【 0 1 1 4 】

< 送信フレーム選択部の処理 >

図10A、図10Bは送信フレーム選択部215における送信帯域計算処理を説明するフローチャートである。送信フレーム選択部215では、論理キュー221から読み出されてリンクへ送信されるイーサネットフレームについて、キューIDによりトラフィックを区別したうえで各トラフィックを送信する量を計算する。このため、送信フレーム選択部215は、割当帯域管理部213から通知された帯域割当管理表700の内容および輻輳トラフィック推定部214から通知された各トラフィックの制限帯域、キュー制御部218から通知された各論理キューの使用量、送信停止タイマー216を参照して得たPriority毎のP A U S E状態（タイマー値が0なら非P A U S E状態、それ以外ならP A U S E状態）を用いて、各トラフィックの送信帯域を送信帯域計算処理により計算する。そして、送信フレーム選択部215は、計算された前記送信帯域内でフレームが送信されるようにキュー制御部218および入出力制御部219を制御する。

【 0 1 1 5 】

前記送信帯域計算処理は一定時間毎（例えば0.01秒毎）に繰り返して行う。なお、本発明では送信帯域計算処理実行の契機は一定周期に限るものではなく、例えば、P A U S Eフレーム受信を送信帯域計算処理実行の契機としてもよいし、時間間隔をトラフィッ

10

20

30

40

50

クの流量等に応じて変更してもよい。以下、送信帯域計算処理を説明するが、処理の中で帯域とサイズ（量）を相互変換する必要がある場合、帯域は送信帯域計算処理を実行する時間間隔でサイズを割ることで求められ、サイズは帯域に送信帯域計算処理をかけることで求められる。以下の説明では、これらの変換方法への言及は明示的には行わないが、必要に応じて行うものとする。

【0116】

なお、サイズの単位の例としてはバイトまたはビット、帯域の単位の例としてはbps（ビット毎秒）とし、1バイトは8ビットとするが、本発明で用いる単位はこれに限らず、例えばオクテット、バイト毎秒等を用いてもよい。また、以下の説明中で「トラフィック」は対象Priorityに属するトラフィックを、「キュー」は当該トラフィックのイーサネットフレームが格納されている論理キュー221をそれぞれ表す。なお、送信帯域計算処理において、変数「制限帯域」は輻輳トラフィック推定部214から通知された各トラフィックの制限帯域を保持するものとし、変数「割当帯域」は割当帯域管理部213から通知された帯域割当管理表に基づいて各トラフィックに保証する帯域を保持するものとし、変数「キューの使用量」はキュー制御部218から通知された各論理キュー221の使用量を各トラフィックについて保持するものとし、各トラフィックについての変数はそれぞれのトラフィックで異なる値を保持できるものとする。

【0117】

送信フレーム選択部215は送信帯域計算処理が開始状態S1001から開始されると、ステップS1002に進む。ステップS1002では、送信フレーム選択部215が各トラフィックの変数「送信帯域」に変数「制限帯域」をそれぞれ代入し、ステップS1003に進む。ステップS1003では、変数「合計割当帯域」に全トラフィックの「割当帯域」の合計値を代入し、ステップS1004に進む。ステップS1004では、ステップS1005からステップS1030までの処理を繰り返す。

【0118】

ステップS1005では、各トラフィックについて、ステップS1006からステップS1012までの処理を繰り返す。ステップS1006では、送信フレーム選択部215が、変数「キューの使用量」を、前回送信帯域計算処理を実行してからの経過時間で割った値を変数「最大送信帯域」に代入し、ステップS1007に進む。

【0119】

ステップS1007では、送信フレーム選択部215が、変数「要求帯域」に変数「最大送信帯域」の値を代入し、ステップS1008に進む。ステップS1008では、送信フレーム選択部215が、当該トラフィックがPAUSE状態であるか否かを判定し、真の場合はステップS1009に進み、偽の場合はステップS1010に進む。ステップS1009では、送信フレーム選択部215が、変数「要求帯域」に0を代入しステップS1010に進む。

【0120】

ステップS1010では、送信フレーム選択部215は、変数「制限帯域」が変数「保証帯域」より小さく（すなわち帯域制限されており）変数「送信帯域」が変数「最大送信帯域」より小さいか否かを判定し、真ならステップS1011に進み、偽ならステップS1012に進む。ステップS1011では、変数「要求帯域」に変数「送信帯域」を設定し、ステップS1012に進む。

【0121】

ステップS1012では、送信フレーム選択部215がステップS1005からの繰り返し処理が終了したら、ステップS1013に進む。

【0122】

ステップS1013では、送信フレーム選択部215が変数「最大送信帯域」が0より大きいトラフィックが1トラフィックのみで、かつ当該トラフィックの変数「制限帯域」が変数「保証帯域」より小さい（すなわち帯域制限されている）か否かを判定し、真ならステップS1014に進み、偽ならステップS1015に進む。ステップS1014では

、送信フレーム選択部 215 が各トラフィックについて、変数「要求帯域」に変数「最大送信帯域」を代入し、ステップ S 1015 に進む。ステップ S 1015 では、変数「合計余剰帯域」、変数「合計不足帯域」、変数「合計割当済帯域」にそれぞれ 0 を代入し、ステップ S 1016 に進む。

#### 【0123】

ステップ S 1016 では、送信フレーム選択部 215 が各トラフィックについて、ステップ S 1017 からステップ S 1022 までの処理を繰り返す。ステップ S 1017 では、変数「分割帯域」に変数「送信帯域」の値を設定し、ステップ S 1018 に進む。ステップ S 1018 では、送信フレーム選択部 215 が変数「要求帯域」が変数「分割帯域」より小さいか否かを判定し、真ならステップ S 1019 に進み、偽ならステップ S 1020 に進む。ステップ S 1019 では、変数「合計余剰帯域」に（分割帯域 - 要求帯域）を加え、ステップ S 1021 に進む。

10

#### 【0124】

ステップ S 1020 では、変数「合計不足帯域」に（要求帯域 - 分割帯域）を加え、ステップ S 1021 に進む。ステップ S 1021 では、変数「合計割当済帯域」に分割帯域の値を加え、ステップ S 1022 に進む。ステップ S 1022 では、送信フレーム選択部 215 がステップ S 1016 からの繰り返し処理が終了したら、図 10B のステップ S 1023 に進む。図 10B のステップ S 1023 では、変数「合計余剰帯域」に（合計割当帯域 - 合計割当済帯域）

を加え、ステップ S 1024 に進む。ステップ S 1024 では、送信フレーム選択部 215 が、変数「合計不足帯域」および変数「合計余剰帯域」が十分小さいか否かを判定し、真なら終了状態 1031 に進み、偽ならステップ S 1025 に進む。なお、ステップ S 1024 における「十分小さい」とは例えば 1 bps 以下等、分配するに足らない帯域や制御可能な最小帯域等を表す。本発明では、これに限ることなく、例えばステップ S 1024 における判定を若干緩め（大き目）にして、繰り返し回数を削減してもよい。

20

#### 【0125】

ステップ S 1025 では、送信フレーム選択部 215 が、各トラフィックについて、ステップ S 1026 からステップ S 1029 までを繰り返す。ステップ S 1026 では、送信フレーム選択部 215 が、要求帯域が分割帯域より小さいか否かを判定し、真ならステップ S 1027 に進み、偽ならステップ S 1028 に進む。ステップ S 1027 では、変数「送信帯域」に要求帯域を代入し、ステップ S 1029 に進む。

30

#### 【0126】

ステップ S 1028 では送信フレーム選択部 215 が、変数「送信帯域」に（合計余剰帯域 × （要求帯域 - 分割帯域） / 合計不足帯域）を代入し、ステップ S 1029 に進む。なお、ステップ S 1028 における計算式は、合計余剰帯域の分配例を表しており、本発明はこれに限るものではない。例えば、1 つのトラフィックに合計余剰帯域をすべて与え、他のトラフィックへの分配を 0 とし、処理を簡素化してもよい。

#### 【0127】

ステップ S 1029 では、送信フレーム選択部 215 が、ステップ S 1025 からの繰り返し処理が終了したら、ステップ S 1030 に進む。ステップ S 1030 では、ステップ S 1004 からの繰り返し処理が終了したら、終了状態 1031 に進む。終了状態 1031 では送信帯域計算処理を終了する。

40

#### 【0128】

送信帯域計算処理の要点は以下の通りである。送信フレーム選択部 215 は、各トラフィックについて、送信帯域に制限帯域を代入して（ステップ S 1002）、以下を繰り返すことで送信帯域を調整する。送信フレーム選択部 215 は、要求帯域を PAUSE 状態であれば 0 とし（ステップ S 1009）、そうでなければ論理キュー 221 のフレーム送信に必要な帯域（最大送信帯域）とする（ステップ S 1007）。ただし、帯域が制限されており、要求帯域が送信帯域より大きい場合は、要求帯域は送信帯域とする（ステップ

50

S 1 0 1 1)。ただし、送信しようとしている（すなわちキューにフレームのある）トラフィックが、帯域制限された１トラフィックのみの場合は（ステップ S 1 0 1 3）、要求帯域はキューのフレーム送信に必要な帯域（最大送信帯域）とする（ステップ S 1 0 1 4）。そして、送信帯域と要求帯域を比べて、不足分及び余剰分を計算し（ステップ S 1 0 1 6）、不足のトラフィックに余剰分を再分配する（１０２５）。不足分及び余剰分が予め設定した一定値以下となって十分小さくなったら、余剰分の再分配を終了する（１０２４）。以上を繰り返す（ステップ S 1 0 3 0）。

#### 【 0 1 2 9 】

すなわち、本発明では、他のトラフィックがあり、かつ帯域制限されている場合には制限帯域を超えないように送信帯域を設定し、それ以外の場合は最大送信帯域を上限帯域にして、保証帯域を各トラフィックに保証しつつ、各トラフィックで使用していない余剰帯域を帯域が不足しているトラフィックに再配分して、送信帯域を決定する。本発明では送信帯域計算処理が送信対象となるトラフィック数の制限要因を含まないため、必要に応じて対象トラフィック数を増減することができる。

#### 【 0 1 3 0 】

##### < 変形例 >

なお、図２の説明で P F C における処理を前提として説明したが、本発明は P F C を前提としない実施形態も可能である。その場合も図２の構成や処理は同一で、P F C において１つの p r i o r i t y を用いると限定し、フレームの形式は I E E E 8 0 2 . 3 x に規定される P A U S E フレームとする。この場合、本発明は I E E E 8 0 2 . 3 x に規定される P A U S E をサポートしたイーサネットに適用可能となる。

#### 【 0 1 3 1 】

また、本発明は P A U S E の頻度（間隔）に基づいて送信する際のフレーム選択を変更することを含むため、図２で説明したようなフレーム転送装置でなくとも、１つ以上のネットワークインタフェースを持っていれば、実施可能である。このような実施形態を用いることで、計算機用の N I C （ネットワークインタフェースカード）や仮想マシンの仮想ネットワークインタフェースに対して本発明を実施することが可能となる。また、フレーム転送装置においてもポート単位で本発明を実施することも可能である。

#### 【 0 1 3 2 】

##### < まとめ >

以上のように、本発明のフレーム転送装置１によれば、リンクの送信先となるノード１１７から P A U S E フレームを受信した場合、P A U S E フレームで指定された時間までフレームの送信を一時停止し、送信先のノード１１７で輻輳の原因となっている論理トラフィックを推定する。そして、フレーム転送装置１は、当該ノード１１７への送信再開時には、輻輳の原因として推定された論理トラフィックの送信を抑制し、他の論理トラフィックを優先して送信する。これにより、送信先のノード１１７が P A U S E フレームを再度送信するのを防ぐことで当該ノード１１７の論理トラフィックが他の論理トラフィックに与える影響を抑制することができる。これらはすべて P A U S E によるフロー制御に基づいているため、確実な通信を保證できる。

#### 【 0 1 3 3 】

そして、本発明のフレーム転送装置１では、前記従来例のようにノード１１７は既存の構成のままでよく、新たな機能を付加する必要がないため導入コストを抑制できる。さらに、本発明のフレーム転送装置１では、前記従来例の非特許文献１のように、扱うことが可能な論理トラフィックの数が８以下に制限されることはなく、任意の数の論理トラフィックを制御することが可能となる。これにより、マルチコアプロセッサ等を利用した仮想化システムなどでは、仮想マシンの数がネットワークインタフェースの数に制約を受けることがなくなって、可用性を向上させることが可能となる。

#### 【 産業上の利用可能性 】

#### 【 0 1 3 4 】

以上のように、本発明は論理トラフィックを転送する装置や計算機及びネットワークシ

10

20

30

40

50

ステムに適用することができる。

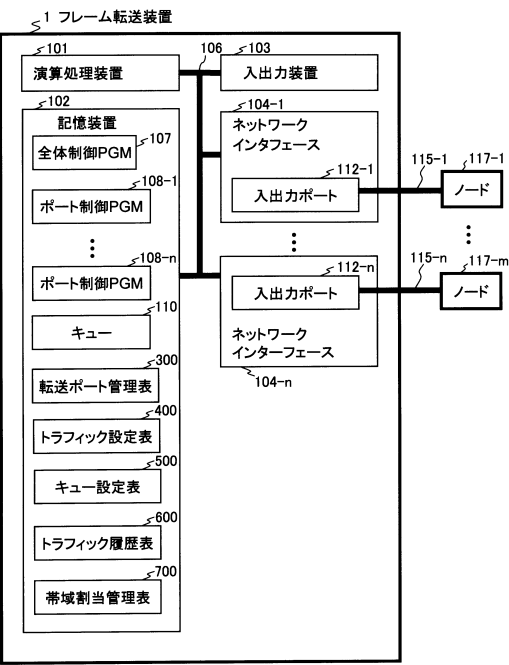
【符号の説明】

【0135】

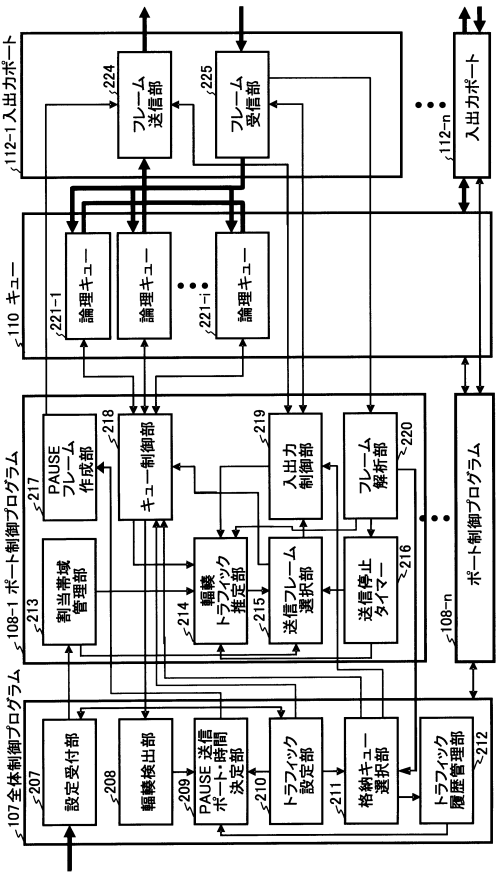
1	フレーム転送装置	
101	演算処理装置	
102	記憶装置	
103	入出力装置	
104 - 1 ~ 104 - n	ネットワークインタフェース	
107	全体制御プログラム	
108 - 1 ~ 108 - n	ポート制御プログラム	10
110	キュー	
112 - 1 ~ 112 - n	入出力ポート	
117 - 1 ~ 117 - m	ノード	
207	設定受付部	
208	輻輳検出部	
209	PAUSE送信ポート・時間決定部	
210	トラフィック設定部	
211	格納キュー選択部	
212	トラフィック履歴管理部	
213	割当帯域管理部	20
214	輻輳トラフィック推定部	
215	送信フレーム選択部	
216	送信停止タイマー	
217	PAUSEフレーム作成部	
218	キュー制御部	
219	入出力制御部	
220	フレーム解析部	
221 - 1 ~ 221 - i	論理キュー	
224	フレーム送信部	
225	フレーム受信部	30



【図 1】



【図 2】



【図 3】

300 転送ポート管理表

301 MACアドレス	302 VLAN ID	303 ポートID
MAC1	100	2
⋮	⋮	⋮

【図 6】

600 トラフィック履歴表

601 トラフィックID	602 ポートID	603 時刻	604 フレームサイズ
1	MAC 2	12:34:56	1200
⋮	⋮	⋮	⋮

【図 4】

400 トラフィック設定表

401 送信元アドレス	402 送信先アドレス	403 VLAN ID	404 Priority	405 トラフィックID
MAC 1	MAC 2	100	5	1
⋮	⋮	⋮	⋮	⋮

【図 7】

700 帯域割当管理表

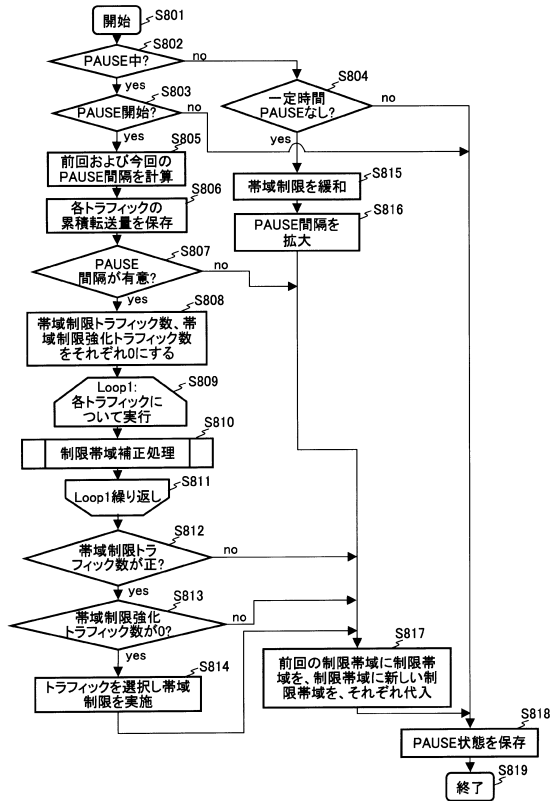
701 キューID	702 割当帯域
1	100Mbps
⋮	⋮

【図 5】

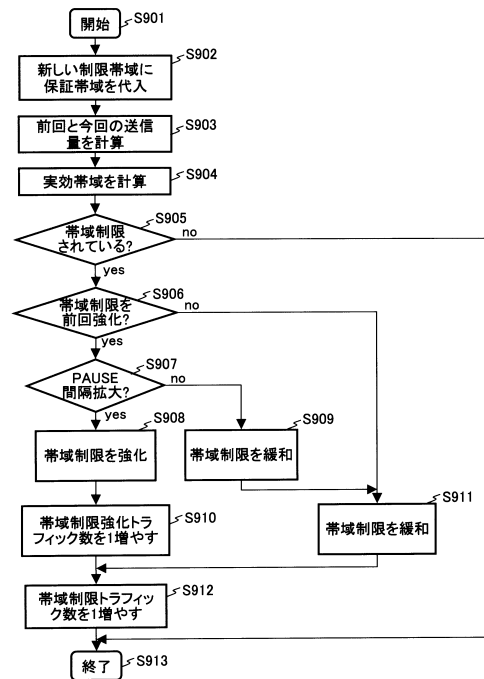
500 キュー設定表

501 トラフィックID	502 ポートID	503 キューID
1	2	3
⋮	⋮	⋮

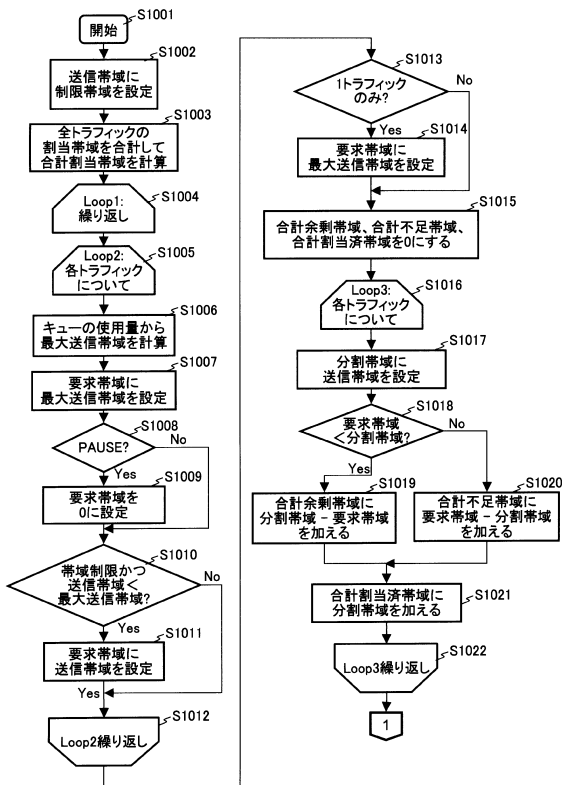
【図 8】



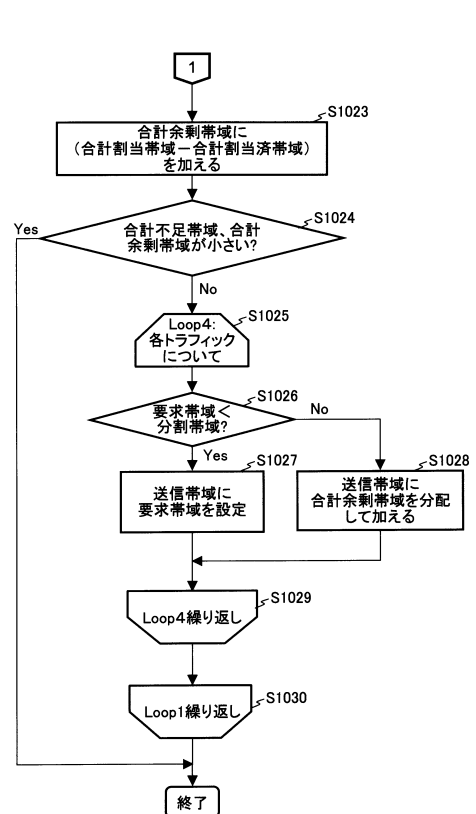
【図 9】



【図 10 A】



【図 10 B】



---

フロントページの続き

審査官 安藤 一道

(56)参考文献 特開 2 0 0 6 - 0 5 0 3 6 1 ( J P , A )  
特開 2 0 0 8 - 2 3 6 7 3 3 ( J P , A )  
特開 2 0 1 1 - 0 1 9 0 4 0 ( J P , A )  
特開 2 0 0 4 - 1 0 4 4 2 7 ( J P , A )  
特開 2 0 0 4 - 1 5 9 2 0 3 ( J P , A )  
特開 2 0 0 6 - 0 3 3 7 1 3 ( J P , A )

(58)調査した分野(Int.Cl. , D B 名)  
H 0 4 L 1 2 / 2 8  
H 0 4 L 1 2 / 8 2 5