

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3556495号
(P3556495)

(45) 発行日 平成16年8月18日(2004.8.18)

(24) 登録日 平成16年5月21日(2004.5.21)

(51) Int. Cl.⁷

H04L 12/56

F I

H04L 11/20 1 O 2 Z

H04L 11/20 1 O 2 E

請求項の数 20 (全 28 頁)

(21) 出願番号	特願平10-356012	(73) 特許権者	000003078
(22) 出願日	平成10年12月15日(1998.12.15)		株式会社東芝
(65) 公開番号	特開2000-183965(P2000-183965A)		東京都港区芝浦一丁目1番1号
(43) 公開日	平成12年6月30日(2000.6.30)	(74) 代理人	100058479
審査請求日	平成14年3月19日(2002.3.19)		弁理士 鈴江 武彦
		(74) 代理人	100084618
			弁理士 村松 貞男
		(74) 代理人	100068814
			弁理士 坪井 淳
		(74) 代理人	100092196
			弁理士 橋本 良郎
		(74) 代理人	100091351
			弁理士 河野 哲
		(74) 代理人	100088683
			弁理士 中村 誠

最終頁に続く

(54) 【発明の名称】 パケットスイッチ及びパケット交換方法

(57) 【特許請求の範囲】

【請求項1】

入側転送手段からスイッチング手段を経由して所望の出側転送手段へ、パケットを転送するパケットスイッチであって、

所定の転送先毎の輻輳状況を観測するための輻輳状況観測手段と、

前記パケットに、前記輻輳状況観測手段により観測された該パケットの転送先の輻輳状況に基づいて、優先度を付与するための優先度付与手段と、

前記スイッチング手段内部で前記パケットの衝突が発生した場合に、パケットに付与された優先度に基づいて、優先して転送すべきパケットを選択するためのパケット選択手段とを備えたことを特徴とするパケットスイッチ。

【請求項2】

前記優先度付与手段は、パケットの転送先の輻輳の度合いがより大きいほど、より低い優先度を付与することを特徴とする請求項1に記載のパケットスイッチ。

【請求項3】

前記優先度付与手段は、各入側転送手段内に設けられたものであることを特徴とする請求項1または2に記載のパケットスイッチ。

【請求項4】

前記優先度付与手段があるパケットに優先度を付与するために参照する転送先の輻輳状況が未知または無効となっている場合に、該転送先に転送する最初の1つのパケットまたは複数のパケット群に付与する優先度を一時的に高く設定することを特徴とする請求項1な

10

20

いし3のいずれか1項に記載のパケットスイッチ。

【請求項5】

前記優先度付与手段は、各入側転送手段毎に1つずつまたは複数の入側転送手段で1つ設けられた、前記輻輳状況観測手段により観測された輻輳状況に基づいて所定の転送先毎に設定される輻輳度を記憶する輻輳度テーブルを参照して、前記パケットに付与すべき優先度を設定することを特徴とする請求項1ないし4のいずれか1項に記載のパケットスイッチ。

【請求項6】

前記輻輳状況観測手段は、各出側転送手段内に設けられたものであることを特徴とする請求項1ないし5のいずれか1項に記載のパケットスイッチ。

10

【請求項7】

前記輻輳状況観測手段は、前記転送先毎の輻輳状況として、対応する出側転送手段毎、対応する出側転送手段のクラス毎、対応する出側転送手段のポート毎、対応する出側転送手段の各ポートのクラス毎、または対応する出側転送手段の各ポートの各クラスのフロー毎に、前記輻輳状況を観測する手段を含むことを特徴とする請求項6に記載のパケットスイッチ。

【請求項8】

前記輻輳状況観測手段により観測された前記輻輳状況を、前記優先度付与手段に反映させるための輻輳状況通知手段を更に備えたことを特徴とする請求項1ないし7のいずれか1項に記載のパケットスイッチ。

20

【請求項9】

前記優先度付与手段は、各入側転送手段内に設けられ、前記輻輳状況通知手段は、前記出側転送手段にパケットが到着したことを契機に、該パケットを送出した入側転送手段に対して所定の輻輳状況に関する情報を通知することを特徴とする請求項8に記載のパケットスイッチ。

【請求項10】

前記入側転送手段毎に、各入側転送手段において複数のパケットが前記スイッチング手段への転送を待っている場合に、該複数のパケットの転送順を制御するためのスケジューリング手段を更に備え、

前記スケジューリング手段は、各パケットの転送先の輻輳状況を考慮して、輻輳していない転送先行きのパケットが優先的に前記スイッチング手段に転送されるように制御することを特徴とする請求項1ないし9のいずれか1項に記載のパケットスイッチ。

30

【請求項11】

優先度の付与されたパケットを入側転送手段からスイッチング手段を経由して所望の出側転送手段へ転送し、該スイッチング手段は内部でパケットの衝突が発生した場合に該優先度を考慮して選択したパケットを出側転送手段に対して転送し、その他のパケットを該スイッチング手段内部で廃棄するパケットスイッチであって、

前記入側転送手段は、

前記衝突によるパケット廃棄が検出された場合に、該廃棄されたパケットを再送するための手段と、

40

前記再送パケットに付与すべき優先度を、もとの廃棄されたパケットに付与した優先度よりも高く設定するための手段とを含むことを特徴とするパケットスイッチ。

【請求項12】

優先度の付与されたパケットを入側転送手段からスイッチング手段を経由して所望の出側転送手段へ転送し、該スイッチング手段は内部でパケットの衝突が発生した場合に、該優先度を考慮して選択したパケットを出側転送手段に対して優先的に転送するパケットスイッチであって、

前記入側転送手段は、1つのデータが複数のパケットに分割して搭載された際の各パケットを転送する場合に、分割されたデータの後続部分に対応するパケットの優先度を、先頭部分に対応するパケットよりも高く設定するための手段を含むことを特徴とするパケット

50

スイッチ。

【請求項 13】

前記分割されたデータの先頭部分に対応するパケットの優先度を、分割されたデータを搭載したものではないパケットならば付与するであろう優先度と同等の値に設定することを特徴とする請求項 12 に記載のパケットスイッチ。

【請求項 14】

パケットを入側転送手段からスイッチング手段を経由して所望の出側転送手段へ転送するパケットスイッチであって、
前記スイッチング手段を転送されるパケットとともに該パケットを送出した入側転送手段内部の輻輳状況を出側転送手段へ転送するための手段と、前記パケットとともに通知された前記輻輳状況と、該パケットが転送された出側転送手段内部の輻輳状況とを用いて、総合的な輻輳状況を求めるための手段と、
求められた前記総合的な輻輳状況を、通過するパケットフローの量または速度を制御するためにネットワークにおいて行われる輻輳制御に用いるための手段とを備えたことを特徴とするパケットスイッチ。

10

【請求項 15】

入側転送手段から、パケット衝突の発生しないスイッチング手段を経由して所望の出側転送手段へ、パケットを転送するパケットスイッチであって、
転送先毎の輻輳状況を観測するための輻輳状況観測手段と、
各入側転送手段から送出されるパケットに前記輻輳状況観測手段により観測された該パケットの転送先の輻輳状況に基づいて優先度を付与し、各パケットをパケット衝突の発生するトポロジを有する仮想的なスイッチング網内を転送させたと仮定して、前記パケットの衝突が発生した場合にパケットに付与された優先度に基づき優先して転送すべきパケットを選択するシミュレーションを行い、このシミュレーションにより前記入側転送手段から前記出側転送手段までパケットが到達した結果と同等な結果となるように、前記パケット衝突の発生しないスイッチング手段の接続パターンを決定する接続パターン決定手段とを備えたことを特徴とするパケットスイッチ。

20

【請求項 16】

前記パケット衝突の発生しないスイッチング手段は、クロスバスイッチであり、前記仮想的なスイッチング網は、単位スイッチにより構成されたスイッチ網であることを特徴とする請求項 15 に記載のパケットスイッチ。

30

【請求項 17】

入側転送手段からスイッチング手段を経由して所望の出側転送手段へ、パケットを転送するパケットスイッチのパケット交換方法であって、
前記入側転送手段は、前記パケットにその転送先の輻輳状況に基づいた優先度を付与してこれを前記スイッチング手段へ送出し、
前記スイッチング手段は、前記入側転送手段から転送された前記パケットをその転送先に従って交換するとともに、その内部でパケットの衝突が発生した場合には各パケットに付与された優先度を考慮して選択したパケットを優先して交換し、
前記パケットの到達した前記出側転送手段は、所定の観測単位についての輻輳状況の観測結果を示す情報を、該パケットを送出した前記入側転送手段に通知することを特徴とするパケット交換方法。

40

【請求項 18】

優先度の付与されたパケットを入側転送手段からスイッチング手段を経由して所望の出側転送手段へ転送し、該スイッチング手段は内部でパケットの衝突が発生した場合に該優先度を考慮して選択したパケットを出側転送手段に対して転送し、その他のパケットを該スイッチング手段内部で廃棄するパケットスイッチのパケット交換方法であって、
前記入側転送手段は、前記衝突によるパケット廃棄が検出された場合に、該廃棄されたパケットに、もとの廃棄されたパケットに付与した優先度よりも高い優先度を付与して、これを再送することを特徴とするパケット交換方法。

50

【請求項 19】

優先度の付与されたパケットを入側転送手段からスイッチング手段を経由して所望の出側転送手段へ転送し、該スイッチング手段は内部でパケットの衝突が発生した場合に、該優先度を考慮して選択したパケットを出側転送手段に対して優先的に転送するパケットスイッチのパケット交換方法であって、

前記入側転送手段は、1つのデータが複数のパケットに分割して搭載された際の各パケットを転送する場合に、分割されたデータの先頭部分に対応するパケットに所定の優先度を付与して送出し、該パケットが転送先に到達したならば、該データの後続部分に対応するパケットに該所定の優先度より高い優先度を付与してを送出することを特徴とするパケット交換方法。

10

【請求項 20】

パケットを入側転送手段からスイッチング手段を経由して所望の出側転送手段へ転送するパケットスイッチのパケット交換方法であって、

前記入側転送手段は、その内部の輻輳状況を示す情報をパケットに付加して送出し、

前記出側転送手段は、その内部の輻輳状況と、前記パケットに付加されて通知された前記輻輳状況とを用いて、総合的な輻輳状況を求め、

求められた前記総合的な輻輳状況を、所定の輻輳制御に用いることを特徴とするパケット交換方法。

【発明の詳細な説明】**【0001】**

20

【発明の属する技術分野】

本発明は、複数のポート間でパケットを交換するパケットスイッチ及びパケット交換方法に関する。

【0002】**【従来の技術】**

近年、目覚しくデータ通信の分野が発展している。伝統的には通信網の代表といえば電話網のことであったが、企業および家庭へのパソコンの普及によりインターネットを代表とするデータ通信網の重要性が急速に高まっている。最近ではデータ通信網の低コスト性や高効率性を生かして、データ通信網上で電話サービスを実現する技術も確立されつつあり、データ通信網に電話網を吸収すべきなどという議論までなされている。

30

【0003】

さて、このデータ通信網のほとんどは、送信したいデータに宛先を示す情報を付与したものの、つまりパケット(ATMにおけるセルの意味も含む)の通信を行なうパケット交換網(ATM交換網の意味も含む)である。データ通信のトラフィック量の増大とともに、このパケットを交換するパケットスイッチの大規模化に対する要求が高まっている。

【0004】

パケットスイッチの構成方法は数多く提案されているが、大きくは2種類に分類できる。一つはやや古いタイプのものでポート数を大きくできる多段スイッチ型、もう一つは比較的新しいタイプのもので、ポート数が少ない単段スイッチ型である。

【0005】

40

後者の単段スイッチ型のパケットスイッチは、輻輳している出力ポートがあっても輻輳していない出力ポート行きのパケットの流れがそれに妨げられないようにする輻輳制御機構が備わっていることが多い。これに対して前者の多段スイッチ型のパケットスイッチは、流れ込むパケットの流量が予め決まっていることを前提にしている。したがって、重度の輻輳を想定していないものが多かった。この点は、輻輳の発生には複数の入力ポートと複数の出力ポートとが互いに関係するために、大きいポート数まで拡張することが前提の多段スイッチ型では、輻輳に対処できる機構を組み込むことが困難であった、という事情によるところも大きいと思われる。

【0006】

データ通信網の代表であるインターネットでは、各ユーザがどのような速度でパケットを

50

転送するかが予め決められているわけではないため、いつでもどこでも輻輳が発生し得る。そのため、パケットスイッチのアーキテクチャを考える際には、大規模化のための技術と同時に輻輳を扱う技術が重要といえる。輻輳を扱う技術には、パケットスイッチ内で輻輳が発生した場合に、ユーザに対してパケットの転送速度を下げてもらうように通知する技術と、輻輳が発生してもその輻輳に関係のないパケットの流れが悪くならないようにする技術とがある。

【0007】

以下、このような従来の多段スイッチ型のパケットスイッチの概略について説明する。

【0008】

多段スイッチ型パケットスイッチは、小さな単位スイッチを多段に組み合わせることで大規模なパケットスイッチを作ることが可能である (Joseph Y. Hui: "Switching and Traffic Theory for Integrated Broadband Networks", ISBN 0-7923-9061-X, Kluwer Academic Publishers, 1990)。

10

【0009】

図6に、多段スイッチ型パケットスイッチの一構成例を示す。図6は、3入力3出力の単位スイッチ145を組合せて、27入力27出力のパケットスイッチを構成した例である。

【0010】

入力ポート121から入力されたパケットは、入側転送部102からスイッチング部104を経由して所望の出側転送部106へ転送され、そこから出力ポート161へ出力される。この多段スイッチ型パケットスイッチは、単位スイッチの数と段数を増やせば、さらに大きなパケットスイッチを構成することが可能であり、また、例えば単位スイッチをより大きな8入力8出力に置き換えれば、より少ない数の単位スイッチで大規模なパケットスイッチを構築できる。

20

【0011】

多段スイッチ型パケットスイッチのよく知られている構成方法の一つが、ランダム網とルーティング網を従属接続する方法である。図6のランダム網141とルーティング網143は、互いに線対称の関係になっている結合網である。それぞれの網は、どの入力リンクからでも任意の出力リンクへ到達でき、かつ、その経路は一つしかない。ランダム網とルーティング網を従属接続することで、ある入側転送部102からある出側転送部106への経路が、ランダム網とルーティング網の境界にある単位スイッチの数(図6では9個)だけ存在することになる。

30

【0012】

入側転送部102から出力されたパケットは、ランダム網141内を他のパケットと衝突しないように、ランダムに選ばれた空いているリンクに転送され、ルーティング網143の入口へ到達する。このランダムな転送は、ルーティング網143の入口(図6では9個)に、パケットを確率的に分散させる目的で行なわれる。パケットは、ルーティング網143でその宛先に基づいて転送され、出側転送部106に到達する。

【0013】

ルーティング網143の内部においては、1本の単位スイッチの出力リンクに複数のパケットが同時に向かったためにパケットの衝突が発生する可能性がある。同じ出側転送部106行きのパケット同士であれば必ず衝突するし、異なる出側転送部106行きのパケット同士でも衝突する場合がある。各単位スイッチ145が待ち合わせバッファを持たないようなパケットスイッチでは、衝突すると唯一のパケットが転送され、残りのパケットは直ちに廃棄される。このような衝突によってパケットが廃棄された場合、該当する入側転送部102は廃棄されたパケットを再送する。

40

【0014】

このように、スイッチング部104内でパケットが廃棄されても、そのパケットが出側転送部106へ無事に転送されるまで何度も入側転送部102が再送し続ける。したがって

50

、スイッチング部 104 の内部でパケットの情報が失われることはない。パケットスイッチの外部から見ればパケットが廃棄されるのは入側転送部 102 および出側転送部 106 の内部のみである。

【0015】

このような多段スイッチ型パケットスイッチ内におけるパケット転送制御では、3種類の輻輳が発生し得る。これらの輻輳を一つずつ説明する。ここでは、説明を分かり易くするために、スイッチング部 104 として 2 入力 2 出力の単位スイッチを 3 段に組み合わせた 4 入力 4 出力のものを用いる (図 7 ~ 図 9)。この小規模のスイッチング部の場合でも、1 段目の単位スイッチ 145 から 2 段目の単位スイッチの入力まではランダム網 141、2 段目の単位スイッチ 145 の出力から 3 段目の単位スイッチ 145 はルーティング網 143、という構成は図 6 と同様であり、ここで行う 3 種類の輻輳についての議論もしくは説明が図 6 やそれ以上の大規模なスイッチング部に対しても同様に成立する (適用できる)。

10

【0016】

まず、図 7 を参照しながら、輻輳の第 1 番目の例について説明する。

【0017】

図 7 の例では、入側転送部 [A] に出側転送部 [C] 行きのパケットが存在し、入側転送部 [C] に出側転送部 [D] に行きのパケットが存在するものとする。この場合、理想的なパケットスイッチであれば、この二つのパケットは宛先が異なるため同時に転送できるはずである。ところが、実際には、ランダム網 141 が二つのパケットを例えばスイッチング部 104 の 2 段目で (偶然に) 下側の単位スイッチ 145 に転送するような状況があり得る。そのような状況が発生した場合、それら二つのパケットは同時に宛先へ到着することはできない。逆にもしそれら二つのパケットが 2 段目の上側と下側の別々の単位スイッチに分かれて転送されていたなら、この輻輳は回避されていたはずである。このような輻輳を「ルーティング網内衝突」と呼ぶことにする。

20

【0018】

このルーティング網内衝突による転送性能の劣化を解決するための方法は古くからいくつかのものが提案されていた。例えば、ルーティング網の前段にランダム網を従属接続しルーティング網内衝突が特定のパケットの集合に偏らないようにパケットを分散させる方法や、この輻輳によってパケットの再送が発生してもパケットスイッチの外部からみたスループットが低下しないようにスイッチング部のパケット転送速度を入力ポートおよび出力ポートの速度よりも高速化する方法が知られている。このようにスイッチアーキテクチャを工夫して継続的な輻輳が発生しないようにすれば、スイッチの外部からこの輻輳の影響を見えないようにすることができる。

30

【0019】

次に、図 8 を参照しながら輻輳の第 2 番目の例について説明する。

【0020】

図 8 の例では、二つの入側転送部 [A], [C] に同じ出側転送部 [C] 行きのパケットが存在する。この二つのパケットはどのような経路でランダム網 145 の内部を転送されても出側転送部 [C] に到着するまでに必ず衝突してしまう。このような輻輳を「出側転送部入口輻輳」と呼ぶことにする。

40

【0021】

この輻輳は、基本的には同じ行き先を持つパケットが継続的に集中することで発生する。例えば、図 8 において、さらに入側転送部 [B] から出側転送部 [D] 行きのパケットの流れがあったとすると、出側転送部 [C] 行きの過剰なパケット転送が継続する影響で出側転送部 [D] 行きのパケットの流れが悪くなってしまふ。

【0022】

従来の多段スイッチ型のようなパケットスイッチは出側転送部の数が多いため、この輻輳に対する適切な対処方法は未解決の課題であった。なぜなら出側転送部の数が多いと、どの出側転送部が混んで、どの出側転送部がすいているかの情報を一度に収集することが困

50

難だったからである。

【0023】

次に、図9を参照しながら輻輳の第3番目の例について説明する。

【0024】

この輻輳は、出側転送部(図9の例ではC)の出力ポートの速度が、そのポート行きのパケットの到着速度よりも遅いことが原因で発生する。この輻輳を「出側転送部内輻輳」と呼ぶことにする。図9の例においては、出側転送部[C]が輻輳している場合、本来、入側転送部[A]は、出側転送部[C]行きのパケットを必要以上に頻繁に転送する必要はなく、出側転送部[C]の出力ポートの速度程度で十分である。

【0025】

従来のパケットスイッチは過剰に多くのパケットを同一の出力ポートへ転送してしまうことにより、先に述べた第1の輻輳(「ルーティング網内衝突」)や第2の輻輳(「出側転送部入口輻輳」)を誘発する原因になっていた。

【0026】

以上、3種類の輻輳についてそれぞれ説明した。このように、従来の多段スイッチ型パケットスイッチは、単位スイッチを組み合わせることで大規模化することができるという利点を持っていたが、重度の輻輳に対するよい対処方法がなかった。つまり、スイッチング部内に大量の輻輳ポート行きパケットが転送されるおそれがあった。これらはスイッチング部内で廃棄される可能性が高い無駄な転送である。本来は出力ポートの転送速度を越える転送は無意味であり、それを越えて大量に転送されたパケットはスイッチング部内で廃棄され再送を繰り返すという事態が発生し得る。このような事態になると、輻輳していないポート行きのパケットも、輻輳しているポート行きのパケットとスイッチング部内で衝突して廃棄と再送を繰り返し流れが悪くなる。結果的にはパケットスイッチ全体のパケット転送効率が著しく低下するという欠点があった。

【0027】

【発明が解決しようとする課題】

以上述べたように、従来の多段スイッチ型のパケットスイッチでは、スイッチング部内に大量の輻輳ポート行きパケットが転送されるおそれがあった。これらはスイッチング部内で廃棄される可能性が高い無駄な転送である。本来は出力ポートの転送速度を越える転送は無意味であり、それを越えて大量に転送されたパケットはスイッチング部内で廃棄され再送を繰り返すという事態が発生し得る。このような事態になると、輻輳していないポート行きのパケットも、輻輳しているポート行きのパケットとスイッチング部内で衝突して廃棄と再送を繰り返し流れが悪くなる。結果的にはパケットスイッチ全体のパケット転送効率が著しく低下するという問題点があった。したがって、例えば、トラフィックが急増するインターネット等に適用するためには、ポート数を大きくできる多段スイッチ型のパケットスイッチにおける輻輳制御が要望される。

【0028】

本発明は、上記事情を考慮してなされたもので、輻輳による影響を回避し転送等の性能をより最大限に引き出すことのできるパケットスイッチ及びパケット交換方法を提供することを目的とする。

【0029】

本発明は、輻輳しているポート行きのパケットの影響により、輻輳していないポート行きのパケットの流れが妨げられないようにしたパケットスイッチ及びパケット交換方法を提供することを目的とする。

【0030】

また、本発明は、廃棄されたパケットを再送する場合にそれを廃棄されにくくさせることができるようにしたパケットスイッチ及びパケット交換方法を提供することを目的とする。

【0031】

また、本発明は、複数のパケットに分割した1つデータを一塊りに転送させることができ

10

20

30

40

50

るようにしたパケットスイッチ及びパケット交換方法を提供することを目的とする。

【0032】

【課題を解決するための手段】

本発明（請求項1）は、入側転送手段からスイッチング手段を経由して所望の出側転送手段へ、パケット（例えば、パケットデータ、またはパケットデータの転送に先立つリクエスト）を転送するパケットスイッチであって、所定の転送先毎の輻輳状況を観測するための輻輳状況観測手段と、前記パケットに、前記輻輳状況観測手段により観測された該パケットの転送先の輻輳状況に基づいて、優先度を付与するための優先度付与手段と、前記スイッチング手段内部で前記パケットの衝突が発生した場合に、パケットに付与された優先度に基づいて、優先して転送すべきパケットを選択するためのパケット選択手段とを備えたことを特徴とする。

10

【0033】

本発明によれば、パケットにその転送先の輻輳の程度に応じて優先度を付与することができ、例えば輻輳していないポート行きのパケットには相対的に高い優先度を付ける、といった制御が可能となる。そして、スイッチング手段内でパケットの衝突が発生しても、優先度を考慮したパケット選択を行うので、例えば輻輳していないポート行きのパケットの流れは相対的に高い優先度が付与されていれば妨げられることがなく、結果としてパケットスイッチの転送効率の低下を回避することができる。

【0034】

スイッチング手段内部の衝突で負けたパケットについては、例えば、待ち合わせバッファで待たせるようにしてもよいし、あるいは即座に廃棄し該当する入側転送手段から再送するようにしてもよい。入側転送手段がスイッチング手段内でパケットが廃棄されたことを検出する方法は、明示的に再送要求を受信したことによる方法、暗黙的に到着確認情報が期待する時刻までに戻ってこなかったことによる方法などが考えられる。

20

【0035】

なお、1つの出側転送手段には、1つの輻輳状況が存在する場合の他に、複数の輻輳状況が存在する場合がある。例えば、出側転送手段毎でもフロー毎やクラス毎に輻輳状況が異なるケース、1つの出側転送手段に複数の低速ポートがぶら下がっているケースなどが考えられる。

【0036】

好ましくは、前記優先度付与手段は、パケットの転送先の輻輳の度合いがより大きいほど、より低い優先度を付与する（パケットの転送先の輻輳の度合いがより小さいほど、より高い優先度を付与する）ようにしてもよい。

30

【0037】

好ましくは、前記優先度付与手段は、各入側転送手段内に設けられたものであるようにしてもよい。

【0038】

好ましくは、前記優先度付与手段があるパケットに優先度を付与するために参照する転送先の輻輳状況が未知または無効となっている場合に、該転送先に転送する最初の1つのパケットまたは複数のパケット群に付与する優先度を一時的に高く設定するようにしてもよい。

40

【0039】

好ましくは、前記優先度付与手段は、各入側転送手段毎に1つずつまたは複数の入側転送手段で1つ設けられた、前記輻輳状況観測手段により観測された輻輳状況に基づいて所定の転送先毎に設定される輻輳度を記憶する輻輳度テーブルを参照して、前記パケットに付与すべき優先度を設定するようにしてもよい。

【0040】

好ましくは、前記輻輳状況観測手段は、各出側転送手段内に設けられたものであるようにしてもよい。

【0041】

50

好ましくは、前記輻輳状況観測手段は、前記転送先毎の輻輳状況として、対応する出側転送手段毎、対応する出側転送手段のクラス毎、対応する出側転送手段のポート毎、対応する出側転送手段の各ポートのクラス毎、または対応する出側転送手段の各ポートの各クラスのフロー毎に、前記輻輳状況を観測する手段を含むようにしてもよい。

【0042】

好ましくは、前記輻輳状況観測手段により観測された前記輻輳状況を、前記優先度付与手段に反映させるための輻輳状況通知手段を更に備えるようにしてもよい。

【0043】

好ましくは、前記優先度付与手段は、各入側転送手段内に設けられ、前記輻輳状況通知手段は、前記出側転送手段にパケットが到着したことを契機に、該パケットを送出した入側転送手段に対して所定の輻輳状況に関する情報を通知するようにしてもよい。より具体的には、例えば、ACK信号に輻輳状況を載せる方法や、逆向きのパケットを作ってそれに輻輳状況を載せる方法が考えられる。

【0044】

なお、全出側転送手段の輻輳状況を全入側転送手段に定期的に通知するようにしてもよい。

【0045】

また、前記優先度付与手段は、前記入側転送手段内に設ける代わりにスイッチング手段内に設けるようにしてもよい。

【0046】

また、前記輻輳状況観測手段や前記輻輳状況通知手段は、前記出側転送手段内に設ける代わりにスイッチング手段内に設けるようにしてもよいし、前記出側転送手段内とスイッチング手段内の両方に設けるようにしてもよい。

【0047】

好ましくは、前記入側転送手段毎に、各入側転送手段において複数のパケットが前記スイッチング手段への転送を待っている場合に、該複数のパケットの転送順を制御するためのスケジューリング手段を更に備え、前記スケジューリング手段は、各パケットの転送先の輻輳状況を考慮して、輻輳していない転送先行きのパケットが優先的に前記スイッチング手段に転送されるように制御するようにしてもよい。

【0048】

また、例えば、パケットにリアルタイムあるいはベストエフォートといったクラスの属性がついている場合に、出力すべきクラスをまず決定し、そのクラスの中で輻輳状況を考慮してパケットを選ぶようにしてもよい。このような場合には、輻輳度の小さな順にパケットが転送されるとは限らないことになる。

【0049】

なお、このスケジューリング手段は、独立して実施することができる。すなわち、本発明は、入側転送手段からスイッチング手段を経由して所望の出側転送手段へパケットを転送するパケットスイッチであって、入側転送手段毎に、各入側転送手段において複数のパケットが前記スイッチング手段への転送を待っている場合に、該複数のパケットの転送順を制御するためのスケジューリング手段と、所定の転送先毎の輻輳状況を観測するための輻輳状況観測手段とを備え、各スケジューリング手段は、各パケットの転送先の輻輳状況を考慮して、輻輳していない転送先行きのパケットが優先的に前記スイッチング手段に転送されるように制御することを特徴とする。

【0050】

本発明（請求項11）は、優先度の付与されたパケット（例えば、パケットデータ、またはパケットデータの転送に先立つリクエスト）を入側転送手段からスイッチング手段を経由して所望の出側転送手段へ転送し、該スイッチング手段は内部でパケットの衝突が発生した場合に該優先度を考慮して選択したパケットを出側転送手段に対して転送し、その他のパケットを該スイッチング手段内部で廃棄するパケットスイッチであって、前記入側転送手段は、前記衝突によるパケット廃棄が検出された場合に、該廃棄されたパケットを再

10

20

30

40

50

送するための手段と、前記再送パケットに付与すべき優先度を、もとの廃棄されたパケットに付与した優先度よりも高く設定するための手段とを含むことを特徴とする。

【0051】

本発明によれば、廃棄されたパケットに付与する優先度をより高くすることにより、該パケットを廃棄されにくくさせることができる。

【0052】

なお、連続廃棄回数に応じて優先度を順次高くしていくようにしてもよい。

【0053】

入側転送手段がスイッチング手段内でパケットが廃棄されたことを検出する方法は、明示的に再送要求を受信したことによる方法のほか、暗黙的に到着確認情報が期待する時刻までに戻ってこなかったことによる方法などが考えられる。

10

【0054】

この発明は、これまでの各発明と組み合わせて実施することができる。

【0055】

本発明（請求項12）は、優先度の付与されたパケット（例えば、パケットデータ、またはパケットデータの転送に先立つリクエスト）を入側転送手段からスイッチング手段を経由して所望の出側転送手段へ転送し、該スイッチング手段は内部でパケットの衝突が発生した場合に、該優先度を考慮して選択したパケットを出側転送手段に対して優先的に転送するパケットスイッチであって、前記入側転送手段は、1つのデータが複数のパケットに分割して搭載された際の各パケットを転送する場合に、分割されたデータの後続部分に対応するパケットの優先度を、先頭部分に対応するパケットよりも高く設定するための手段を含むことを特徴とする。好ましくは、前記分割されたデータの先頭部分に対応するパケットの優先度を、分割されたデータを搭載したものではないパケットならば付与するであろう優先度と同等の値に設定するようにしてもよい。

20

【0056】

本発明によれば、一度データグラムの先頭に対応するパケットが宛先の出側転送部に到達すればその後の部分も高優先で続けて宛先に到達させることができ、結果として複数のパケットからなるデータグラムを一塊りに出側転送部へ転送させることができる。

【0057】

この発明は、これまでの各発明と組み合わせて実施することができる。

30

【0058】

本発明（請求項14）は、パケット（例えば、パケットデータ、またはパケットデータの転送に先立つリクエスト）を入側転送手段からスイッチング手段を経由して所望の出側転送手段へ転送するパケットスイッチであって、前記スイッチング手段を転送されるパケットとともに該パケットを送出した入側転送手段内部の輻輳状況を出側転送手段へ転送するための手段と、前記パケットとともに通知された前記輻輳状況と、該パケットが転送された出側転送手段内部の輻輳状況とを用いて、総合的な輻輳状況を求めるための手段と、求められた前記総合的な輻輳状況を、通過するパケットフローの量または速度を制御するためにネットワークにおいて行われる輻輳制御（輻輳制御とは、EFCI、ECN、source quench、BECN、REDなど）に用いるための手段とを備えたことを特徴とする。

40

【0059】

例えば、ある入側転送部のフローのパケットキューの長さ（例えば289とする）をパケットとともに出側転送部へ転送し、これに出側転送部のそのフローのパケットキューの長さ（例えば12とする）を合計することにより、パケットスイッチ内に存在するそのフローの待ちパケット量（ $= 289 + 12 = 301$ ）を把握できる。そして、例えば、これをユーザへ通知することで、パケットフローの速度を制御し、これによって輻輳制御を制御することができる。

【0060】

この発明は、これまでの各発明と組み合わせて実施することができる。

50

【0061】

本発明（請求項15）は、入側転送手段から、パケット衝突の発生しないスイッチング手段を経由して所望の出側転送手段へ、パケット（例えば、パケットデータ、またはパケットデータの転送に先立つリクエスト）を転送するパケットスイッチであって、転送先毎の輻輳状況を観測するための輻輳状況観測手段と、各入側転送手段から送出されるパケットに前記輻輳状況観測手段により観測された該パケットの転送先の輻輳状況に基づいて優先度を付与し、各パケットをパケット衝突の発生するトポロジを有する仮想的なスイッチング網内を転送させたと仮定して、前記パケットの衝突が発生した場合にパケットに付与された優先度に基づき優先して転送すべきパケットを選択するシミュレーションを行い、このシミュレーションにより前記入側転送手段から前記出側転送手段までパケットが到達した結果と同等な結果となるように、前記パケット衝突の発生しないスイッチング手段の接続パターンを決定する接続パターン決定手段とを備えたことを特徴とする。

10

【0062】

好ましくは、前記パケット衝突の発生しないスイッチング手段は、クロスバスイッチであり、前記仮想的なスイッチング網は、単位スイッチにより構成されたスイッチ網（例えば、ランダム網とルーティング網を接続したもの）であるようにしてもよい。

【0063】

本発明では、接続パターン決定手段がシミュレーションにより、例えば、各入側転送部のデータ転送のリクエストをトーナメントさせ、勝ち残ったリクエストに対応するデータを交換するようにクロスバスイッチ等の接続パターンを決定する。トーナメントの競合では、輻輳度テーブルの輻輳度を考慮して設定もしくは変更された優先度によって、勝ち残るリクエストを決定する。このようにして出側転送部へ向けて転送を許可するリクエストの集合が求まると、接続パターンが求まる。この接続パターンによって設定されたクロスバスイッチ等によってパケットが交換される。

20

【0064】

従来の接続パターン決定に要する計算量は、スイッチのポート数を N とすると N^2 のオーダーとなり、ポート数が多くなるとともに計算が困難になる問題点が知られているが、本発明によれば、スイッチのポート数を N とすると $N \cdot \log N$ のオーダーで出側転送部へ転送してよいリクエストを求めることができ、比較的高速に計算することができるとともに、出側転送部の輻輳状況を考慮しているので公平性も確保できるという利点がある。

30

【0065】

本発明（請求項17）は、入側転送手段からスイッチング手段を経由して所望の出側転送手段へ、パケットを転送するパケットスイッチのパケット交換方法であって、前記入側転送手段は、前記パケットにその転送先の輻輳状況に基づいた優先度を付与してこれを前記スイッチング手段へ送出し、前記スイッチング手段は、前記入側転送手段から転送された前記パケットをその転送先に従って交換するとともに、その内部でパケットの衝突が発生した場合には各パケットに付与された優先度を考慮して選択したパケットを優先して交換し、前記パケットの到達した前記出側転送手段は、所定の観測単位についての輻輳状況の観測結果を示す情報を、該パケットを送出した前記入側転送手段に通知することを特徴とする。

40

【0066】

本発明（請求項18）は、優先度の付与されたパケットを入側転送手段からスイッチング手段を経由して所望の出側転送手段へ転送し、該スイッチング手段は内部でパケットの衝突が発生した場合に該優先度を考慮して選択したパケットを出側転送手段に対して転送し、その他のパケットを該スイッチング手段内部で廃棄するパケットスイッチのパケット交換方法であって、前記入側転送手段は、前記衝突によるパケット廃棄が検出された場合に、該廃棄されたパケットに、もとの廃棄されたパケットに付与した優先度よりも高い優先度を付与して、これを再送することを特徴とする。

【0067】

本発明（請求項19）は、優先度の付与されたパケットを入側転送手段からスイッチング

50

手段を経由して所望の出側転送手段へ転送し、該スイッチング手段は内部でパケットの衝突が発生した場合に、該優先度を考慮して選択したパケットを出側転送手段に対して優先的に転送するパケットスイッチのパケット交換方法であって、前記入側転送手段は、1つのデータが複数のパケットに分割して搭載された際の各パケットを転送する場合に、分割されたデータの先頭部分に対応するパケットに所定の優先度を付与して送出し、該パケットが転送先に到達したならば、該データの後続部分に対応するパケットに該所定の優先度より高い優先度を付与してを送出することを特徴とする。

【0068】

本発明（請求項20）は、パケットを入側転送手段からスイッチング手段を経由して所望の出側転送手段へ転送するパケットスイッチのパケット交換方法であって、前記入側転送手段は、その内部の輻輳状況を示す情報をパケットに付加して送出し、前記出側転送手段は、その内部の輻輳状況と、前記パケットに付加されて通知された前記輻輳状況とを用いて、総合的な輻輳状況を求め、求められた前記総合的な輻輳状況を、所定の輻輳制御に用いることを特徴とする。

10

【0069】

なお、以上の本発明は、パケット網全体の輻輳の制御にも適用可能である。

【0070】

また、装置に係る本発明は方法に係る発明としても成立し、方法に係る本発明は装置に係る発明としても成立する。

【0071】

また、装置または方法に係る本発明は、コンピュータに当該発明に相当する手順を実行させるための（あるいはコンピュータを当該発明に相当する手段として機能させるための、あるいはコンピュータに当該発明に相当する機能を実現させるための）プログラムを記録したコンピュータ読取り可能な記録媒体としても成立する。

20

【0072】**【発明の実施の形態】**

以下、図面を参照しながら発明の実施の形態を説明する。

【0073】**（第1の実施形態）**

図1に、本発明の第1の実施形態に係る多段スイッチ型のパケットスイッチの構成例を示す。また、併せて、該パケットスイッチの動作概要を示す。

30

【0074】

本パケットスイッチの基本的な全体構成は、図6と同様で、入力ポートから入力されたパケットは、入側転送部2からスイッチング部4を経由して所望の出側転送部6へ転送され、そこから出力ポートへ出力される。本パケットスイッチの入力ポート数、出力ポート数は、任意である。また、1つの入側転送部や出側転送部に複数の入力ポートや出力ポートが対応しても構わない。本パケットスイッチのスイッチング部4の内部の網構成も任意であり、その網を構成する1つの単位スイッチの入力数、出力数も任意である。もちろん、本パケットスイッチが図6と同じ全体構成であってもよい。なお、図1では、4つの入側転送部2と4つの出側転送部6を示している。また、図1では、スイッチング部4の内部構成の詳細は省略している。

40

【0075】

ここで、本実施形態における「パケット」とは、本パケットスイッチで交換したいデータそのものであってもよいし、また、データの転送に先立ってスイッチング部の内部を転送される情報であってもよい（前者をデータ・パケットと呼び、後者のような情報を「リクエスト」と呼ぶものとする）。

【0076】

「パケット」という語句が「リクエスト」を指す場合には、まず、入側転送部2が、データの転送に先立って、宛先情報が書き込まれたリクエストをスイッチング部4へ送出手続き。このリクエストには優先度が書き込まれる。リクエストが無事に宛先の出側転送部6に

50

到達した場合もしくは到達したと判断された場合には、そのリクエストに対応するデータ・パケットが入側転送部 2 から出側転送部 6 へ転送される。リクエストが経路を確保し、その経路に従ってデータ・パケットを転送すれば、データ・パケットは途中で他のデータ・パケットと衝突することなく安全に出側転送部 6 へ転送される。データ・パケットを転送している時間に次のデータ・パケットに対応するリクエストを転送するようにすれば、パイプライン的に効率良くデータを入側転送部 2 から出側転送部 6 へ交換することができる。

【 0 0 7 7 】

本実施形態では、交換したいデータそのものとデータ・パケットの転送に先立つリクエストとを区別せずに「パケット」の語句を用いて説明するが、いずれの場合であっても本発明は有効に作用する。

10

【 0 0 7 8 】

また、本発明は、リクエストを用いる（もしくは経路を確保する）パケットスイッチにも、リクエストを用いない（もしくは経路を確保しない）パケットスイッチにも適用可能である。

【 0 0 7 9 】

これらの点は、後述する他の実施形態についても同様である。

【 0 0 8 0 】

さて、本実施形態のパケットスイッチでは、パケットを入側転送部 2 からスイッチング部 4 を経由して所望の出側転送部 6 へ転送する際に、各パケットに優先度を付与し、スイッチング部においてこの優先度に基づいて転送パケット間の優先制御を行う。つまり、この優先度は、パケットスイッチ内で局所的に使われるものである。

20

【 0 0 8 1 】

優先度は、パケットのパケットスイッチ内における到達すべき転送先、例えば、出側転送部、あるいは出力ポート、あるいは出力ポートをさらに細分化したクラスもしくはパケットの属するフロー（フローは、例えば、そのパケットの宛先アドレス、あるいはこれにポート番号やプロトコル番号等の情報を適宜組み合わせたものによって規定される；ATM 通信であればフローは仮想コネクションである）についての輻輳状況を少なくとも考慮して設定される。優先度を設定するにあたっては、他の情報、例えば、そのパケットの持つ所定の属性（例えば、パケットの属するクラスなど）や、そのパケットの持つ特別の条件（例えば、再送パケットであることなど）をも考慮するようにしてもよい。

30

【 0 0 8 2 】

優先度の段階については、特に制限はなく、何段階としても本発明は有効に作用する。

【 0 0 8 3 】

パケットに優先度を付与する場所については、それがどこであっても本発明は有効に作用するが、本実施形態では、各入側転送部 2 においてパケットに優先度を付与する場合を例にとって説明する。

【 0 0 8 4 】

輻輳状況を観測する場所については、例えばスイッチング部 4 内でも出側転送部 6 内でも本発明は有効に作用するが、後者がより適しているので、本実施形態では、各出側転送部 6 において輻輳状況を観測し、その情報を入側転送部 2 に通知する場合を例にとって説明する。

40

【 0 0 8 5 】

図 1 に示した輻輳状況を考慮して設定した優先度による制御を実現するためのパケットスイッチは、一構成例として、各入側転送部 2 毎に優先度付与部 2 1 と輻輳度テーブル 1 2 を持ち、出側転送部 4 毎に輻輳状況観測部 6 1 を持ち、各单位スイッチ 4 5 毎にパケット選択部 4 5 1 を持つ構成となっている。

【 0 0 8 6 】

なお、全入側転送部で 1 つの輻輳度テーブルを共用する構成としてもよい。

この場合には、出側転送部側から輻輳状況に関する情報の通知を受け、これを共用の輻輳

50

度テーブルに反映させる輻輳度テーブル管理部（図示せず）を持つ構成とすればよい。

【0087】

また、一部の単位スイッチ45にはパケット選択部451を設けないようにしてもよい。

【0088】

概略的には、入側転送部2の優先度付与部21は、出側転送部6からその輻輳状況観測部61によって観測された輻輳状況に関する情報の通知を受け、これを輻輳度テーブル12に反映させること（輻輳度テーブル管理機能）と、当該入側転送部2からパケットが送出される際に輻輳度テーブル12の内容を考慮してそのパケットに付与する優先度を決定すること（優先度付与機能）を行う。

【0089】

輻輳度は、通知された輻輳状況に関する情報そのものであってもよいし、通知された輻輳状況に関する情報を変換した値であってもよい。輻輳度の設定は、出側転送部単位に行なってもよいし、ポート単位に行なってもよいし、さらに細分化してクラス単位あるいはフロー単位に行なってもよい。

【0090】

優先度の設定については、基本的な考え方としては、輻輳度テーブル12からそのパケットに該当する輻輳度を参照し、輻輳度の高いものほど優先度を低く設定し、輻輳度の低いものほど優先度を高く設定するというものである。

【0091】

ただし、輻輳度と優先度との関係については、多種多様なバリエーションが考えられる。輻輳度と優先度とが線形の関係にあってもよいし非線形の関係にあってもよい。また、輻輳度の段階数と優先度の段階数を同じにしてもよいし、異なるものにしてもよい。また、前述したように、輻輳度と他の情報から優先度を求めるようにしてもよい。

【0092】

また、輻輳度から優先度を求める手順もしくは関数等を、輻輳度の設定と同じ単位（例えばポート単位）で定義してもよいし、輻輳度の設定より細かい単位（例えば、ポート単位に対してフロー単位）で設けてもよい。

【0093】

概略的には、単位スイッチ45のパケット選択部451は、パケットが衝突したときにそれらに付与された優先度を考慮して優先して転送すべきパケットを選択するものである。

【0094】

パケットの選択については、基本的な考え方としては、最も高い優先度を持つパケットを選択するというものであるが、優先度に加えて他の情報（例えば、パケットヘッダ内に記載されている何らかの情報）をも考慮して選択するようにしてもよい。

【0095】

なお、最も高い優先度を持つパケットが複数存在するなどにより、通常行う判断では1つのパケットを選択できない場合のために何らかの特別の選択基準を設けておいてもよいし、通常行う判断では1つのパケットを選択できない場合には最も優越する複数のパケットのうちからランダムに1つを選択するようにしてもよい。

【0096】

概略的には、出側転送部4の輻輳状況観測部61は、当該出側転送部4について輻輳状況を観測するものである。観測された輻輳状況に関する情報は、当該出側転送部4から入側転送部2に通知される。

【0097】

輻輳状況の観測は、出側転送部単位に行なってもよいし、ポート単位に行なってもよいし、さらに細分化してクラス単位あるいはフロー単位に行なってもよい。

【0098】

通知すべき輻輳状況に関する情報としては、輻輳状況に関する所定の測定値（例えば、単位時間に転送したパケットの個数もしくはバイト数、パケットキューのキュー長（パケット数もしくはバイト数）、パケットキューから廃棄されたパケット数もしくはバイト数な

10

20

30

40

50

ど)をそのまま用いてもよいし、所定の測定値を1つの基準値によって輻輳有りまたは無しに分類した結果を用いてもよいし、所定の測定値を複数の基準値によって分類したレベル値を用いてもよい。

【0099】

なお、輻輳状況の観測の単位(例えばポート単位)と輻輳度テーブルのエントリの単位(例えばポート単位)を一致させてもよいし、輻輳状況の観測の単位(例えばフロー単位)に比べて輻輳度テーブルのエントリの単位(例えばポート単位)をより荒くするようにしてもよい。

【0100】

以下では、本実施形態についてより詳しく説明する。

10

【0101】

最初に、パケットの衝突の発生した場合について説明する。

【0102】

スイッチング部4の内部でパケットの衝突(例えばルーティング網内衝突)が生じた場合に、その箇所にあたる単位スイッチ45では、パケット選択部451により、衝突した各パケットに付与されている優先度を考慮して、優先して転送すべきパケットを選択し、その選択されたパケットのみを転送する。

【0103】

衝突で負けた他のパケットについては、例えば、即座に廃棄し、該当する入側転送部2が後で再送する(この場合、入側転送部2は、パケットの再送を考慮して、送出したパケットを一定時間保持しておく必要がある)。再送の契機を与える方法としては、例えば、パケットを廃棄した単位スイッチが入側転送部2に再送メッセージを出す方法や、あるいはパケットが到達した出側転送部4が転送元の入側転送部2にパケット到達メッセージを出すものとして入側転送部2が到達先となる出側転送部4からパケット到達メッセージを一定時間経過しても受信しない場合にパケットが廃棄されたとみなす方法などがある。

20

【0104】

あるいは、全部または一部の単位スイッチを、パケットバッファを備える構成とし、上記のように衝突で負けたパケットを即座に廃棄する代わりに、該パケットを該単位スイッチ内のパケットバッファに蓄積して一時的に待ち合わせるようにしてもよい。

【0105】

続いて、輻輳状況の観測、通知について説明する。

30

【0106】

本実施形態では、その特徴の一つとして、パケットが転送される経路の輻輳状況を観測し、その輻輳状況に関する情報を、優先度をパケットに付与する部分(本実施形態では入側転送部2)へ伝えるようにしている。前述したように、輻輳状況を観測する場所は、本実施形態では、より適している出側転送部6としている。この点について説明する。

【0107】

出側転送部2にスイッチング部4からパケットが流れ込む速度よりも、出側転送部6からパケットを出力ポートへ出力する速度が低くなっていれば、スイッチング部4内の輻輳状況が出側転送部6の内部の輻輳(例えば、出力用バッファの蓄積パケット数もしくはバイト数の増大)に反映される。つまり、出側転送部入口輻輳が発生すれば必ず出側転送部の内部も輻輳になる。

40

【0108】

また、スイッチング部4に輻輳が発生しなくても、出側転送部6の内部のみに輻輳(前述の出側転送部内輻輳)が生じることもある。これは、一つの出側転送部6が複数の出力ポートを備えている構成で発生しやすい。

【0109】

以上から、出側転送部入口輻輳も出側転送部内輻輳も、出側転送部6だけを観測すれば検出できることになる。

【0110】

50

なお、輻輳状況は、出側転送部毎、出力ポート毎、クラス毎、パケットのフロー毎（例えば、仮想コネクション毎）、またはこれらの組合せ、で観測する。図1では、輻輳状況をフロー毎に観測し、輻輳テーブル12をフロー毎に設定する例を示している。

【0111】

さて、このように観測された輻輳状況は、本実施形態では入側転送部2へ伝えられるが、輻輳状況を入側転送部2へ伝える方法としては、いくつかの方法が考えられる。

【0112】

例えば、出側転送部6から入側転送部2へ転送されるACK/NACK信号とともに輻輳状況を送る方法が考えられる。従来のパケットスイッチにおいて、出側転送部から入側転送部へ転送される信号として、パケットが無事に出側転送部に到着したことを通知するACK (Acknowledgment: 受信応答) 信号、または、出側転送部にパケットが到着しても出側転送部がそのパケットを何らかの理由で拒否するNACK (Non-acknowledgment) 信号がある。一般的によく知られている方法によれば、これらの信号は、出側転送部からそのパケットを送出した入側転送部へ、パケットの転送経路を逆向きにたどって返送される。すなわち、このACK/NACK信号とともに輻輳状況に関する情報を送るとというのが、一つの方法である。

10

【0113】

また、他の方法としては、各出側転送部が、その観測単位（例えばポート）についてそれに一定個数（一定個数を1としてもよい）のパケットが到着する毎に、または一定の時間が経過する毎に、そのときにパケットを送出した入側転送部2（あるいは前回の通知からの間にパケットを送出した入側転送部2）に向かって、輻輳状況に関する情報を通知するパケットを生成して返送するようにしてもよい。

20

【0114】

以上のようにして輻輳状況を転送すれば、ポート数が大きなパケットスイッチであっても、容易に出側転送部6の輻輳状況を入側転送部2へ伝えることが可能となる。

【0115】

なお、上記のいずれの場合においても、観測単位を出側転送部毎より細分化された単位（例えばポート毎）としている場合に、通知すべき輻輳状況に関する情報を、当該到着したパケットと同じ転送先（例えばポート）についての情報とする代わりに、出側転送部を同じくする全ての転送先についての情報としてもよいし、全ての出側転送部についての全ての情報としてもよい。

30

【0116】

また、一定周期等の所定のタイミングで、全ての出側転送部についての全ての輻輳状況に関する情報を、全ての入側転送手段に通知するようにしてもよい。

【0117】

なお、上記のように出側転送手段で輻輳状況を観測し通知する代わりに、その直前または複数段前の単位スイッチで観測し通知するようにしてもよい。また、出側転送手段と単位スイッチの双方で輻輳状況を観測し通知するようにしてもよい。

【0118】

続いて、優先度によるスイッチ内パケット転送制御のいくつかの例について説明する。

40

【0119】

まず、最も代表的な制御、すなわち輻輳状態に応じて行う優先度制御について説明する。

【0120】

パケットの優先度はもともと例えばリアルタイム情報あるいはベストエフォート情報などといったそのパケットやそのパケットのフローが属するクラスに基づいて決定することが一般的であるが、本実施形態では、その優先度を輻輳状況に応じて設定する（あるいは、上記のようにしてもともと付与されていた優先度を変更する）。

【0121】

例えば、輻輳している到達先（輻輳ポートとする）行きのパケットは優先度を低くする。優先度の高いパケットはスイッチング部4の内部で衝突しても優先されるため、優先度の

50

高いパケットにとって輻輳しているポート行きの優先度の低いパケットは存在しないのと同じことになる。

【0122】

したがって、本実施形態により、輻輳していないポート行きのパケットの流れが輻輳ポート行きのパケットの流れによって乱されることが非常に少なくなるという効果がある。輻輳ポート行きのパケットにとっても、優先度が変化させられるだけで、その転送速度は抑制されていないため、優先度の高いパケットの転送のすき間を利用して出側転送部へ転送する試みを継続することができる。そのため、輻輳ポートへもスイッチング部4の転送能力の限りパケットを転送し続けることができる。

【0123】

本実施形態とは異なった方法として優先度を用いるのではなく輻輳ポート行きパケットの転送速度を輻輳度に応じて厳密に抑制するという複雑な制御が考えられるが、本実施形態によれば、そのような転送速度抑制の制御は不要であり、入側転送部の構造も簡単にできるという利点がある。

【0124】

なお、従来からパケットに優先度を付与して転送する方法は知られているが、その従来の方法に比べ、本実施形態は輻輳度に応じて動的に優先度を設定もしくは変更することが本質的に異なるものである。

【0125】

次に、再送パケットに対する優先度制御について説明する。

【0126】

スイッチング部4の内部の衝突で廃棄されたパケットを入側転送部2が再送する方式のパケットスイッチでは、同じパケットが連続して廃棄されてしまうと、そのパケットの後で転送を待っているパケット全体の流れが悪くなってしまうおそれがある。

【0127】

これを解決する方法として、スイッチング部4の内部でパケットが廃棄された場合に、その再送パケットの優先度をもとのパケットの優先度よりも高くするようにしてもよい。このような制御によって、再送パケットが再び廃棄されることが少なくなる利点がある。もし再送パケットが再び廃棄された場合には、それに対する再送パケットの優先度を更に高くするようにしてもよい。これによって、何度も再送され続けている不運のパケットを優先して転送することが可能となる。

【0128】

次に、プローブパケットに対する優先度制御について説明する。

【0129】

上記の再送パケットのように優先度を一時的に変更するとよい状況は他にも種々考えられる。

【0130】

例えば、出側転送部6にパケットが到達したことを契機として該出側転送部6から該パケットを送出した入側転送部2に輻輳状況に関する情報を通知する構成においては、入側転送部2が新しい宛先に初めて（または輻輳状況が変化している程度に久しぶりに）パケットを転送する場合、その宛先の輻輳状況はわからない（もしくはその宛先の実際の輻輳状況はわからない）。このように宛先の輻輳状況が未知のとき（もしくは未知に等しいようなとき）に送るパケットは、宛先の輻輳状況を知るための探查信号の意味合いが強い。

【0131】

この探查を行うパケット（プローブパケット）の優先度を高くすれば、スイッチング部4で廃棄されにくくなるため、早期に宛先の輻輳状況が判明し、その後で転送するパケットの優先度を素早く適切な値に設定することができるという利点がある。優先度を高く設定するパケットは最初の1パケットだけで十分であるが、実装上の都合などで複数のパケットの優先度が高くなったとしても効果がある。

【0132】

10

20

30

40

50

次に、分割されたデータグラムに対する優先度制御について説明する。

【0133】

優先度を一時的に変更するとよい他の例としては、複数に分割されたデータグラムをスイッチング部4が転送する場合に、そのデータグラムの最初の部分に対応するパケットの優先度を他の部分に対応するパケットよりも低く設定することである。他の部分に対応するパケットの優先度を最初の部分に対応するパケットよりも高くしてもよい。例えば、先頭のパケットの優先度を、本実施形態の方法で設定し、2番目以降のパケットの優先度を非常に高めの値にしてもよい。

【0134】

また、先頭のパケットの最初の送出では宛先の出側転送部6に到達しなかった場合には、優先度をさらに高く設定して繰り返し転送を試みるようにしてもよい。

10

【0135】

このようにすれば、一度データグラムの先頭に対応するパケットが宛先の出側転送部6に到達すればその後の部分も高優先で続けて宛先に到達させることができ、結果として複数のパケットからなるデータグラムを一塊りに出側転送部6へ転送することができる。出側転送部6が複数に分割されたデータグラムを再び一つに再構成して外部へ出力するような場合などでは、スイッチング部4でこのような転送を行うことにより、データグラムを再構成するために必要なバッファ量を少なくすることができる利点がある。

【0136】

以上のように、パケットの優先度は、輻輳度のみによって決定する方法だけでなく、例えば、パケットまたはフローの属するクラスによる優先度に、輻輳状況による変更と再送回数に応じた変更などを行なって、最終的に決定する（もしくは、輻輳度と、クラスおよび/またはその他のファクターとをパラメータとして求める）方法がある。

20

【0137】

また、さらにパケットの優先度を高度に制御することにより、スイッチング部4内のパケット転送速度を保証することも可能である。例えば、入側転送部2でフローの転送速度をモニタし、設定した速度を下回りそうになった場合に、優先度を一時的に高くする制御を行なうことで、転送速度の最小値を保証するサービスを提供することができる。このサービスは、スイッチがすいているときには、保証速度を上回る転送が可能である。また、重みをつけてその比で帯域を分割する、といったことも可能である。

30

【0138】

なお、パケットの優先度を設定もしくは変更する優先度付与部の実現方法として、入側転送部2の内部にプロセッサを配置しソフトウェアで優先度を決定するにすれば、パケットスイッチのポート構成や扱うクラス数などが運用中に変化しても柔軟に対応することができる。

【0139】

続いて、輻輳度テーブルと優先度付与部のバリエーションについて説明する。

【0140】

本実施形態では、その大きな特徴として、出側転送部6から通知された輻輳状況に応じて、対応する宛先へと向かうパケットの優先度を変更するようにしている。

40

【0141】

その実現例としては、パケットを転送することにより通知された輻輳状況に関する情報から得られる輻輳度を、そのパケットを転送した入側転送部2の輻輳度テーブル12で記憶し、そこから転送する同じ宛先行きのパケットの優先度を設定もしくは変更する方法があるが、この方法としては、その他にも、種々のバリエーションが考えられる。

【0142】

例えば、輻輳状況に関する情報を複数の場所、例えば複数の入側転送部2で共通に記憶してもよい。輻輳度テーブル12を共用することで、各入側転送部2は自らが得た輻輳状況に加えて他の入側転送部2が得た輻輳状況をも使用でき、より多くの輻輳状況がわかるため、輻輳に即座に対応できる確率が高くなるという利点がある。

50

【 0 1 4 3 】

また、他の方法として、図 2 に示すように、入側転送部 2 で輻輳状況を考慮した優先度をパケットに付与するのではなく、（入側転送部 2 では輻輳状況を考慮しない優先度をパケットに付与し）、スイッチング部 4 の内部に配置したパケット優先度変更部 4 2 で、パケットに付与されている優先度を、輻輳状況を考慮したものに変更するようにしてもよい。

【 0 1 4 4 】

この方法の場合、出側転送部 6 から転送された輻輳状況は、スイッチング部 4 の内部のパケット優先度変更部 4 2 へ通知される。パケット優先度変更部 4 2 は、その通知された輻輳状況をもとに輻輳度を求めて輻輳度テーブル 5 2 に記憶し、通過するパケットの宛先に応じてそのパケットの優先度を変更する。このパケット優先度変更部 5 2 は必ずしも一つでパケットスイッチの全ての出側転送部 6 を受け持つ必要はない。一つのスイッチング部 4 に複数のパケット優先度変更部 5 2 を用意し、それぞれが近くの出側転送部 6 を担当すれば処理が分散され実装が容易にできる。

10

【 0 1 4 5 】

この方法の場合、宛先の輻輳状況に応じた優先度をその宛先に近い場所でパケットに付与できるため、効果的にパケットに優先度を付与することができるという利点がある。

【 0 1 4 6 】

（第 2 の実施形態）

第 1 の実施形態では、輻輳度に応じて優先度を制御したが、第 2 の実施形態では、スケジューリング（複数のパケットがスイッチング部への転送を待っている場合に、どのパケットを転送するかを選択する処理）を輻輳度を考慮して行うようにしたものである。

20

【 0 1 4 7 】

本発明の第 2 の実施形態に係る多段スイッチ型のパケットスイッチの基本的な全体構成は、第 1 の実施形態のように、図 6 と同様で、入力ポートから入力されたパケットは、入側転送部からスイッチング部を経由して所望の出側転送部へ転送され、そこから出力ポートへ出力される。また、パケットスイッチの入力ポート数、出力ポート数が任意である点、1 つの入側転送部や出側転送部に複数の入力ポートや出力ポートが対応しても構わない点、スイッチング部の内部の網構成も任意である点、その網を構成する 1 つの単位スイッチの入力数、出力数も任意である点は、第 1 の実施形態と同様である。

【 0 1 4 8 】

また、本実施形態では、入側転送部の輻輳度テーブル、出側転送部の輻輳状況観測部については、第 1 の実施形態と同様のものを備えるものとする。また、第 1 の実施形態の優先度付与部のうち輻輳度テーブルの輻輳度を設定する部分を入側転送部に備えるものとする。

30

【 0 1 4 9 】

図 3 に、本実施形態の多段スイッチ型パケットスイッチの各入側転送部のスケジューリングに関する構成の一例を示す。

【 0 1 5 0 】

図 3 においては、スケジューリングするグループ（例えば、フロー単位）毎に設けられたパケットキュー 2 2、クラス間スケジューリング設定部 2 4、転送セル選択部 2 6、輻輳度テーブル 3 2 を示している。

40

【 0 1 5 1 】

なお、図 3 の例では、スケジューリングに輻輳度とクラスを用いるものとする。また、輻輳度テーブルに輻輳度と併せてクラスのフィールドを設けた構成例を示している。

【 0 1 5 2 】

以下、本実施形態についてより詳しく説明する。

【 0 1 5 3 】

入側転送部内で転送を待っているパケットは、同じ輻輳状況を共有するグループに分類してパケットキュー 2 2 を作ることにより、効率よく転送できる。例えば、このパケットキュー 2 2 は、パケットのフロー（または仮想コネクション）毎に作る。他の例としては、

50

パケットスイッチの各出側転送部の各出力ポートのクラス毎に作ることも有効である。さらに、より簡単には各出側転送部のクラス毎に作ってもよい。ただし、クラス概念を持たないパケットスイッチの場合には、クラス毎の分類は不要である。

【0154】

ここでは、フロー毎にパケットキュー22を作って管理するものとして説明する。

【0155】

輻輳状況が出側転送部から転送されてくると、入側転送部はそれをもとに各フローの輻輳度をテーブル32に記憶・更新する。各フローにクラスの属性がついていれば、その情報も同じテーブル32に記憶すると便利である。

【0156】

図3の転送セル選択部26は、転送すべきパケットを選択する場合に、この輻輳度とクラス情報のテーブル32、およびクラス間スケジューリング設定部24を参照する。クラス間スケジューリング設定部24は、例えば、リアルタイムクラス(図3のRT)の転送速度を、入側転送部からのパケット送出最大速度の80パーセントを限度として、ベストエフォートクラス(図3のBE)よりも優先して出力する、などといった、クラス間のパケット選択ポリシーを設定するためのものである。

【0157】

本実施形態では、転送セル選択部26は、まず、(i)このクラス間スケジューリング設定部24に記憶された設定内容に基づいて、クラスを選択し、その後、(ii)その選択されたクラスのフローの中から、なるべく輻輳していない宛先行きのフローを選択するようになっている。

【0158】

その際に、理想的には、最も輻輳していない宛先行きのフローを選択すべきであるが、実装上の都合により必ずしも厳密に輻輳度の順にフローを選ぶことが簡単にできるとは限らない。しかし、輻輳している宛先行きよりも輻輳していない宛先行きを選択することが多いようになれば、それだけでも十分な効果が得られる。なるべく簡単に実装する方法として、フローを輻輳度により例えば3段階に大きく分類し、最も輻輳していない段階のフローから選択することが考えられる。このようにすると、同じ段階の中では必ずしも輻輳度の順に厳密にフローが選択されるわけではないが、全体的にはおよそ輻輳度の低い順にフローが選択されることになる。

【0159】

フローが選択されると、(iii)そのフローのパケットキュー22の先頭からパケットを取り出してスイッチング部へ転送する。

【0160】

以降、上記の(i)~(iii)の手順を繰り返し実行する。

【0161】

なお、本実施形態を実装する場合には、輻輳度の有効期間を定めると好ましい。例えば、入側転送部において、あるフローの輻輳度が決定されてから一定時間新たな輻輳状況が通知されなかった場合には、その古い輻輳度を無効にする。

【0162】

このように本実施形態によれば、入側転送部が輻輳していない宛先へのパケットを優先して転送するので、輻輳して流れが悪くなっているフローの影響で、輻輳していない宛先行きのパケットの流れも悪くなるという問題を容易に解決することができる。

【0163】

また、本実施形態では、パケットを転送すると、そのパケットの転送によってそのフローの最新の輻輳状況が通知される。したがって、この輻輳状況によってそのフローの選択優先度が動的に最適な値に変化する利点がある。

【0164】

ところで、この第2の実施形態(スケジューリングにおけるパケット選択)は、第1の実施形態(衝突時におけるパケット選択)と組み合わせて実施することが可能である。

10

20

30

40

50

【 0 1 6 5 】

この場合には、さらに次のような効果を得ることができる。すなわち、入側転送部の内部に、輻輳している宛先行きパケットしかない場合には輻輳している宛先へ連続してパケットを出力してしまうが、これらのパケットは廃棄され再送される確率が高い。しかし、第2の実施形態に第1の実施形態を組み合わせた場合、パケットに輻輳度に応じた優先度が付与されるため、出力してもかまわない。このような構成は、輻輳している宛先へパケットを転送しないように抑制する制御よりも、簡単に実現できるという利点がある。転送を抑制する制御では、過剰な抑制のためにスイッチング部の利用効率が低下することがないように厳密な制御を行なう必要があるが、それと比較して本実施形態は簡単な制御でよいわけである。

10

【 0 1 6 6 】

なお、第2の実施形態に第1の実施形態を組み合わせる場合、輻輳度テーブル（および通知された輻輳状況から輻輳度を求める部分）は、パケットに優先度を付与するための部分と、スケジューリングのための部分とで、互いに独立した構成にしてもよいし、1つの輻輳度テーブルに共通化してもよい。

【 0 1 6 7 】

（第3の実施形態）

第1の実施形態、第2の実施形態では、輻輳していない出力ポート等行きのパケットの流れが、それとの関係を持たない他の部分の輻輳の影響を受けないようにする点を中心に説明してきた。

20

【 0 1 6 8 】

第3の実施形態では、輻輳の原因となるパケットを送信しているユーザに対してパケット送信速度を下げてもらうように通知するための機構について説明する。

【 0 1 6 9 】

第3の実施形態は、第1の実施形態、第2の実施形態、第1の実施形態と第2の実施形態を組み合わせたもののいずれにも適用可能であり、また独立して実施することも可能であるが（スイッチング部の構成がどのようなものであっても適用可能である）、第3の実施形態と第1の実施形態および/または第2の実施形態とを組み合わせて実施することで、総合的に輻輳対策に優れたパケットスイッチを構築することができる。

【 0 1 7 0 】

図4に、本実施形態に係る多段スイッチ型のパケットスイッチの構成例を示す。また、併せて、該パケットスイッチの動作概要を示す。

30

【 0 1 7 1 】

ユーザに対して輻輳を通知する機能は、従来、単段スイッチ型パケットスイッチにおいては提案または実用化されている。しかし、この機能を、ポート数の多い多段スイッチ型などのパケットスイッチに効果的に適用することは難しかった。まず、ユーザに対して輻輳を通知する方法として、従来の単段スイッチ型パケットスイッチで用いられている主な3つの方法について説明する。

【 0 1 7 2 】

輻輳を通知する第1番目の方法は、パケットの廃棄による方法である。例えば、インターネットで広く使用されているプロトコルTCPでは、パケットの廃棄を検出すると、輻輳を緩和するために転送速度（正確にはウィンドウサイズ）を小さくする制御が働く。TCPから転送されているパケットを扱うパケットスイッチの廃棄制御では、輻輳がひどくなるに従ってパケットの廃棄確率を増やすRED（Random Early Detection）と呼ばれている方法が優れていると言われている。パケットの廃棄によって輻輳を通知する場合、パケットキューの末尾のパケットを廃棄するよりも先頭のパケットを廃棄した方が、ユーザに輻輳を早く通知でき、輻輳を重くなる前に緩和できる可能性が高いことが知られている。

40

【 0 1 7 3 】

輻輳を通知する第2番目の方法は、輻輳経験通知による方法である。パケットヘッダの一

50

部に輻輳経験フラグ領域を設置し、輻輳を経験したパケットは、その領域がマーキングされる。これによりユーザに輻輳を通知する。この方法はATM通信の場合にはEFCI (Explicit Forward Congestion Indication) と呼ばれている。また、インターネットではECN (Explicit Congestion Notification) と呼ばれ、現在標準化の議論中である。輻輳経験通知の場合にもパケットを廃棄する方法と同様に、パケットキューの末尾のパケットにマーキングするよりはパケットキューの先頭のパケットにマーキングした方が、ユーザに輻輳を早く通知することができることが知られている。

【0174】

輻輳を通知する第3番目の方法は、BECN (Backward Explicit Congestion Notification) による方法である。ATM通信では、輻輳が発生するとスイッチはそれを通知するためにBECNセルと呼ばれる特殊なセル (ATM通信ではパケットのことをセルと呼ぶ) を生成し、送信側ユーザに向けて転送することが考えられている。ATMのABRサービスカテゴリの資源管理セルへ情報を載せる処理も送信側ユーザに向けて直接情報を転送する意味で同じである。インターネットにもSource Quenchと呼ばれる同様の制御メッセージが存在する。BECNの場合には、輻輳を下流側のユーザに対してではなく上流側のユーザに直接通知するため、特にパケットキューの先頭にあるパケットに対して何かを行なう必要はない。しかし、どのユーザに対してBECNを送出するかを決定する場合、輻輳している地点つまりパケットスイッチの出側で判定することが、公平な輻輳通知を実現する上では好ましい。というのは、輻輳通知は、輻輳地点毎に、最も輻輳の原因となっているユーザに対して、行なうべきだからである。

【0175】

以上、主な3つの輻輳通知方法について述べた。これらに共通することは、なるべくキューの出口で輻輳制御を行なった方がよいということである。これを、これまで説明してきたような各実施形態のパケットスイッチに当てはめて考えると、入側転送部よりも出側転送部においてこれらの輻輳制御を行なうことが望ましいと考えられる。この点について、従来の多段スイッチ型パケットスイッチでは、出側転送部が、入側転送部の内部の輻輳状況を正確に把握することはできず、パケットスイッチ全体の輻輳状況を把握することは難しかった。単段スイッチ型パケットスイッチのようなポート数の少ないパケットスイッチの場合であれば、各出側転送部に全ての入側転送部の輻輳状況をモニタするような結線を行なうことは可能であるが、多段スイッチ型パケットスイッチのようにポート数が多い場合には困難である。

【0176】

これを解決する方法として、図4に例示するように、パケットを入側転送部2からスイッチング部4を経由して所望の出側転送部6へ転送するパケットスイッチにおいて、パケットとともに、入側転送部2の内部の輻輳状況を出側転送部6へ通知し、その入側転送部2の内部の輻輳状況から出側転送部6は総合的な輻輳状況を判定して、輻輳制御を行なう方法が考えられる。

【0177】

図4の例では、入側転送部[D]の入側輻輳状況として、フロー毎のパケットキューの長さ (例えば289とする) をパケットとともに出側転送部[B]へ転送する。出側転送部[B]の内部の輻輳状況をそのフローのパケットキュー長 (例えば12とする) で観測するものとすれば、出側転送部[B]は、入側のパケットキュー長と出側のパケットキュー長とを合計することにより、パケットスイッチ内に存在するそのフローの待ちパケット量 ($= 289 + 12 = 301$) を把握できる。

【0178】

ユーザへの輻輳通知方法がREDであれば、その待ちパケット量をもとにパケットキューの先頭のパケットを廃棄する確率を決定すればよい。ECNであれば、その待ちパケット量をしきい値と比較してマーキングするかどうかを決定すればよい。Source Qu

en chであれば、その待ちパケット量をしきい値と比較してSource Queue hメッセージを転送するかどうかを決定すればいい。

【0179】

本実施形態は、輻輳状況としてフロー毎のキュー長だけでなく、クラス毎、転送部毎の輻輳状況なども考慮し、その結果得られる情報を利用して、有効に作用する。

【0180】

従来、ポート数が多いパケットスイッチでは、すべての入側転送部の輻輳状況を出側で観測することは不可能であったが、本実施形態によれば、出側転送部6において入側転送部2の輻輳状況を知ることができ、出側転送部6の内部の輻輳状況と合わせて、パケットが転送される経路の輻輳状況を総合的に把握することができる。ユーザへ輻輳を通知するための制御はパケットキューの出口にて行なうことが好ましく、本実施形態はその要求を満たしており、またポート数が大きな多段スイッチ型パケットスイッチに容易に適用できる利点がある。

10

【0181】

(第4の実施形態)

第1～第3の実施形態では、本発明を、パケットを入側転送部から出側転送部へ転送する際にスイッチング部内で衝突が発生するようなパケットスイッチに適用した場合を中心に説明してきたが(第1の実施形態では、衝突が発生した際に、その宛先の輻輳状況に応じた優先度により、優先して転送すべきパケットを決定する)、本発明は、内部で衝突が発生しないパケットスイッチ、例えば、クロスバ型のスイッチング部を持つパケットスイッチにも適用可能である。言い換えると、第1の実施形態(または第1の実施形態に第2および/または第3の実施形態を組み合わせた実施形態)に相当するような制御の仕組みを利用した、クロスバ型等のパケットスイッチを実現することができる。

20

【0182】

一般に、パケットスイッチは、高速で入側転送部と出側転送部とを接続するパターンを計算しなければならない。例えば、図6のような多段スイッチではパケットが自律分散的にルーティングされ衝突に勝ち残ったパケットの経路を得ることによって入側転送部と出側転送部とを接続するパターンを計算していることになる。一方、クロスバ型スイッチなどでは、このように自律分散的なアプローチではなく、集中して接続パターンを計算するエンジンがスイッチに1箇所存在する。一般に、このエンジン計算量は、スイッチのポート数をNとすると N^2 のオーダーとなり、ポート数が多くなるとともに計算が困難になる問題点が知られている。本発明は、この接続パターン計算エンジンのアルゴリズムに適用することが可能である。

30

【0183】

図5に、本実施形態に係るクロスバスイッチ型のパケットスイッチの構成例を示す。また、併せて、該パケットスイッチの動作概要を示す。図5では、第1の実施形態に対応するものを示してある。

【0184】

なお、本実施形態のパケットスイッチの基本的な全体構成は、第1の実施形態のように、図6と同様で、入力ポートから入力されたパケットは、入側転送部2からクロスバスイッチ3を経由して所望の出側転送部6へ転送され、そこから出力ポートへ出力される。また、パケットスイッチの入力ポート数、出力ポート数が任意である点、1つの入側転送部や出側転送部に複数の入力ポートや出力ポートが対応しても構わない点、スイッチング部の内部の網構成も任意である点、その網を構成する1つの単位スイッチの入力数、出力数も任意である点は、第1の実施形態と同様である。

40

【0185】

なお、優先度付与部、輻輳状況観測部、輻輳度テーブル、パケット選択部などの機能は第1の実施形態と同様であるが、これらは、接続パターン計算エンジン7のアルゴリズムの中に作り込むものとする。ただし、例えば輻輳状況観測部で輻輳状況の観測のためにカウンタを用いるとする場合などのように、個別のハードウェアを用いる場合には、当該部分

50

については、これまでの実施形態で示した箇所に実装されることになる。なお、スイッチング部内で行う制御・処理のように接続パターン計算エンジン7のアルゴリズムの中に作り込むことを必須とするもの以外は、これまでの実施形態で示したように実装するようにしてもよい。

【0186】

さて、入側転送部2のデータ転送のリクエストを受けて、接続パターン計算エンジン7がクロスバスイッチ3の接続パターンを計算する。また、図5に例示するように、出側転送部6から輻輳状況に関する情報が接続パターン計算エンジン7へ送られる。

【0187】

この接続パターン計算エンジン7は、その内部でソフトウェア的に各入側転送部2のデータ転送のリクエストをトーナメントさせ(すなわち、シミュレーションを行い)、勝ち残ったリクエストに対応するデータを交換するようにクロスバスイッチ3へ接続パターンを伝える。このトーナメント表は、予め定めたトポロジ(例えば図6で示されるようなトポロジ)を用いればよい。

10

【0188】

図6で示されるようなトポロジを用いるものとした場合、ランダム網のトポロジを用いてリクエストをランダムな順序に並べ替え、ルーティング網のトポロジを用いてトーナメントを行う。

【0189】

トーナメントの競合では、輻輳度テーブル92の輻輳度を考慮して設定もしくは変更された優先度によって、勝ち残るリクエストを決定する。

20

【0190】

このようにして出側転送部6へ向けて転送を許可するリクエストの集合が求まると、接続パターンが計算され、その接続パターンによって設定されたクロスバスイッチ3によってデータが交換される。

【0191】

この方式では、スイッチのポート数をNとすると $N \cdot \log N$ のオーダーで出側転送部6へ転送してよいリクエストを求めることができ、比較的高速に計算することができるとともに、出側転送部6の輻輳状況を考慮しているので公平性も確保できるという利点がある。

【0192】

なお、第1の実施形態と第2の実施形態および/または第3の実施形態を組み合わせた構成も、同様にして実現可能である。

30

【0193】

なお、以上の各実施形態では、本発明を多段スイッチ型パケットスイッチに適用する場合を中心に説明したが、本発明は単段スイッチ型パケットスイッチにも適用可能である。また、本発明は、パケット網全体の輻輳に対処する方法としても利用可能である。

【0194】

なお、各実施形態の構成のうちハードウェアを必須としない各部分は、ソフトウェアとしても実現可能である。

【0195】

また、本実施形態の制御機能は、コンピュータに所定の実行させるための(あるいはコンピュータを所定の実行手段として機能させるための、あるいはコンピュータに所定の機能を実現させるための)プログラムを記録したコンピュータ読取り可能な記録媒体としても実施することもできる。

40

【0196】

本発明は、上述した実施の形態に限定されるものではなく、その技術的範囲において種々変形して実施することができる。

【0197】

【発明の効果】

本発明によれば、転送先の輻輳状況に応じた優先度をパケットに付与し、この優先度を考

50

慮してパケット衝突時の処理を行うようにしたので、各パケットの転送先の輻輳状況に応じたパケット転送制御を行うことができる。例えば、輻輳していないポート行きのパケットに相対的に高い優先度を付与するようにすれば、輻輳しているポート行きのパケットの影響により、輻輳していないポート行きのパケットの流れが妨げられないようにすることができる。

【0198】

また、本発明によれば、廃棄されたパケットにより高い優先度を付与するようにすれば、廃棄されたパケットを再送する場合にそれを廃棄されにくくすることができる。

【0199】

また、本発明によれば、1つのデータの先頭部分に対応するパケットを後続のパケットより低い優先度で転送先に到達させることにより、複数のパケットに分割した1つデータを一塊りに転送させることができるようにすることができる。

【図面の簡単な説明】

【図1】本発明の第1の実施形態に係るパケットスイッチの構成例を示す図

【図2】同実施形態に係るパケットスイッチの他の構成例を示す図

【図3】本発明の第2の実施形態に係るパケットスイッチの要部構成の一例を示す図

【図4】本発明の第3の実施形態に係るパケットスイッチの構成例を示す図

【図5】本発明の第4の実施形態に係るパケットスイッチの構成例を示す図

【図6】多段スイッチ型パケットスイッチの構成例を示す図

【図7】ルーティング網内衝突の一例を示す図

【図8】出側転送部入口輻輳の一例を示す図

【図9】出側転送部内輻輳の一例を示す図

【符号の説明】

2 ... 入側転送部

3 ... クロスバスイッチ

4 ... スイッチング部

6 ... 出側転送部

7 ... 接続パターン計算エンジン

1 2 , 5 2 , 9 2 ... 輻輳度テーブル

2 1 ... 優先度付与部

2 2 ... パケットキュー

2 4 ... クラス間スケジューリング設定部

2 6 ... 転送セル選択部

3 2 ... 輻輳度テーブル

4 2 ... パケット優先度変更部

4 5 ... 単位スイッチ

6 1 ... 輻輳状況観測部

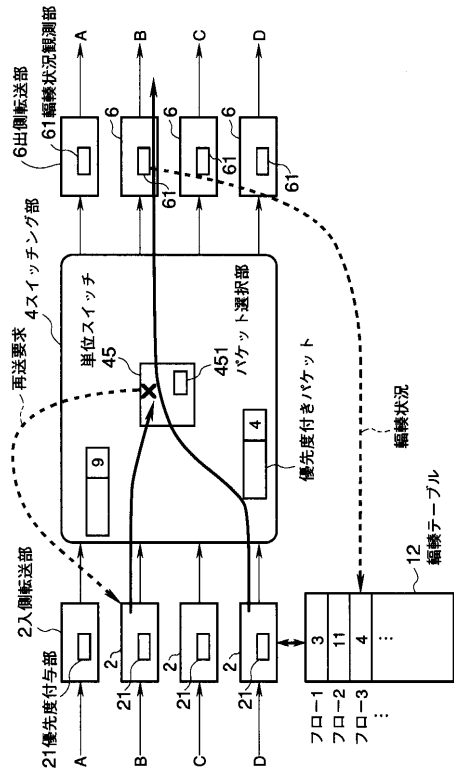
4 5 1 ... パケット選択部

10

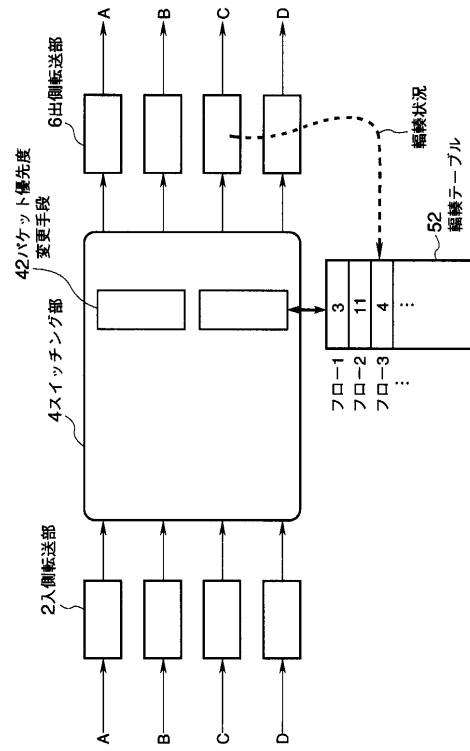
20

30

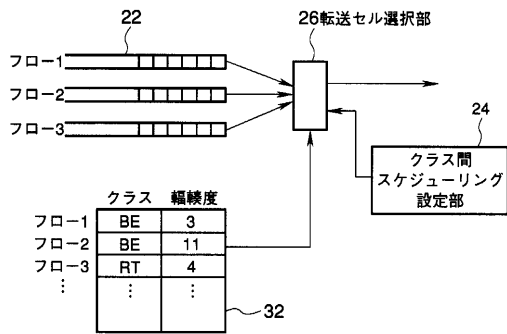
【 図 1 】



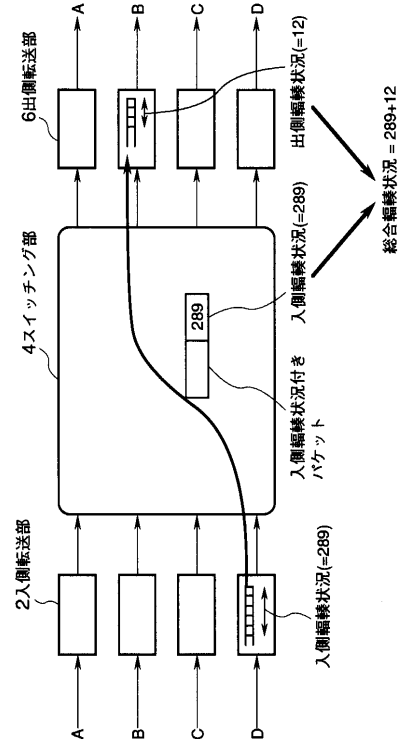
【 図 2 】



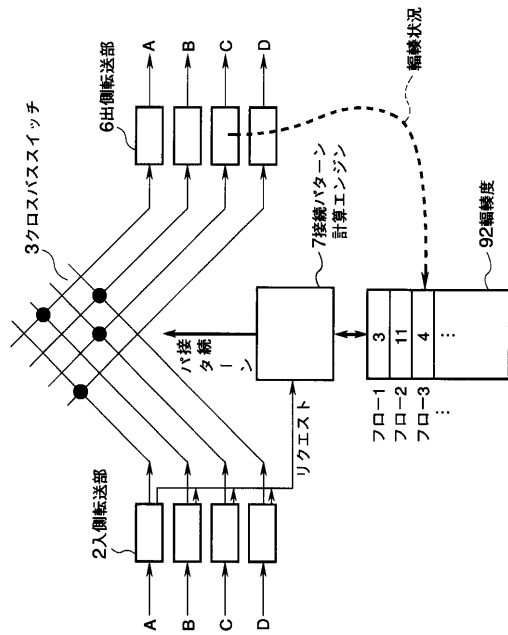
【 図 3 】



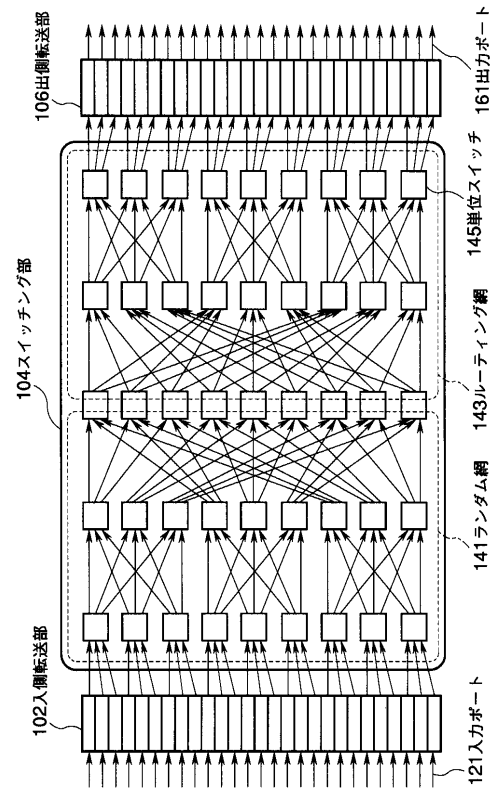
【 図 4 】



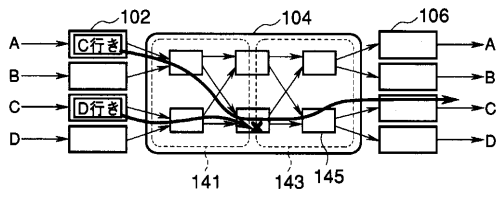
【 図 5 】



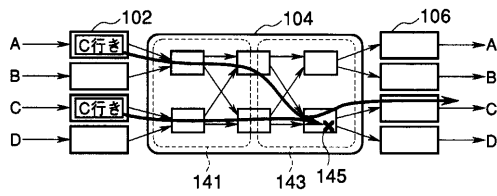
【 図 6 】



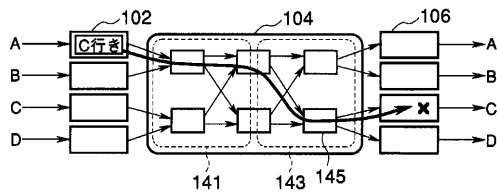
【 図 7 】



【 図 8 】



【 図 9 】



フロントページの続き

(74)代理人 100070437

弁理士 河井 将次

(72)発明者 下條 義満

神奈川県川崎市幸区小向東芝町1番地 株式会社東芝研究開発センター内

(72)発明者 中北 英明

神奈川県川崎市幸区小向東芝町1番地 株式会社東芝研究開発センター内

審査官 石井 研一

(56)参考文献 特開平09-205441(JP,A)

特開平08-288965(JP,A)

特開平07-177179(JP,A)

特開平07-107095(JP,A)

(58)調査した分野(Int.Cl.⁷,DB名)

H04L 12/56