

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3673025号

(P3673025)

(45) 発行日 平成17年7月20日(2005.7.20)

(24) 登録日 平成17年4月28日(2005.4.28)

(51) Int. Cl.⁷

F I

H04L 12/28

H04L 11/20

G

H04Q 3/00

H04Q 3/00

請求項の数 8 (全 72 頁)

(21) 出願番号	特願平8-213569	(73) 特許権者	000003078
(22) 出願日	平成8年8月13日(1996.8.13)		株式会社東芝
(65) 公開番号	特開平9-149051		東京都港区芝浦一丁目1番1号
(43) 公開日	平成9年6月6日(1997.6.6)	(74) 代理人	100058479
審査請求日	平成14年8月14日(2002.8.14)		弁理士 鈴江 武彦
(31) 優先権主張番号	特願平7-238867	(74) 代理人	100084618
(32) 優先日	平成7年9月18日(1995.9.18)		弁理士 村松 貞男
(33) 優先権主張国	日本国(JP)	(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】 パケット転送装置

(57) 【特許請求の範囲】

【請求項1】

入力されたパケットを一時的に蓄積する複数の入力バッファと、これら入力バッファを制御する制御手段と、前記各入力バッファから出力されるパケットを転送する少なくとも1つの出力ポートを具備したパケット転送装置であって、

前記各入力バッファは、入力したパケットを一時的にクラス毎に蓄積するクラス毎の蓄積手段と、

前記制御手段から指示されたクラスのパケットを前記蓄積手段から前記出力ポートへ向けて出力する出力手段とを具備し、

前記制御手段は、

前記蓄積手段の複数の入力バッファ全体における蓄積状況をクラス毎に把握し、この蓄積状況に基づいてパケットを出力すべきクラスを決定し、この決定したクラスの指定を含む指示を前記複数の入力バッファに送信することを特徴とするパケット転送装置。

【請求項2】

前記出力ポートは、前記入力バッファから出力されたパケットを一時的に蓄積する出力バッファを具備し、

前記出力手段は、前記出力バッファ内部のパケットの蓄積状況に応じてパケットを前記出力ポートへ向けて出力することを特徴とする請求項1記載のパケット転送装置。

【請求項3】

10

20

前記出力手段は、前記蓄積手段から出力されたパケットを多重化し、この多重化されたパケットを前記出力ポートに出力する多重化手段を具備し、

前記制御手段は、前記蓄積状況に加えて、前記多重化手段の複数の入力バッファ全体における多重化状況を把握し、前記蓄積状況および前記多重化状況に基づいてパケットを出力すべきクラスを決定することを特徴とする請求項1記載のパケット転送装置。

【請求項4】

前記入力バッファは、複数の入力バッファにまたがる同一クラス内で、各パケットの属する各仮想コネクション間でパケットの転送が公平になるように前記蓄積手段から出力するパケットを選択することを特徴とする請求項1記載のパケット転送装置。

【請求項5】

入力されたパケットを一時的に蓄積する複数の入力バッファと、これら入力バッファを制御する制御手段と、前記各入力バッファから出力されるパケットを転送する少なくとも1つの出力ポートを具備したパケット転送装置であって、

前記各入力バッファは、

入力したパケットを一時的に蓄積する蓄積手段と、

前記制御手段からの指示を受けて、次のフェーズで前記蓄積手段から前記出力ポートへ向けて出力すべきパケットを選択する選択手段と、

この選択手段で選択されたパケットを前記出力ポートに向けて出力する出力手段とを具備し、

前記制御手段は、

前記複数の入力バッファ全体における前記選択手段で選択されたパケットの出力状況を把握し、この出力状況に基づいて新たなパケットを選択するよう前記複数の入力バッファに指示することを特徴とするパケット転送装置。

【請求項6】

前記各入力バッファの選択手段は、それぞれ、

各パケットの属する各仮想コネクション間でパケットの転送が公平になるように前記蓄積手段から出力するパケットを選択することを特徴とする請求項5記載のパケット転送装置。

【請求項7】

入力したパケットを一時的に蓄積するバッファと、このバッファを制御する制御手段と、前記バッファから出力されたパケットを転送する少なくとも1つの出力ポートを具備したパケット転送装置であって、

前記制御手段は、

前記バッファ内に蓄積されたパケットを複数の集合に分けて管理する管理手段と、

前記バッファに入力したパケットを各パケットの属するフロー間で公平になるように前記複数の集合のいずれか1つに振り分ける振り分け手段と、

前記管理手段で管理されている前記複数の集合のうちの1つの集合に属するパケットを前記出力ポートに向けて出力するよう前記バッファに対し指示する指示手段と、

を具備したことを特徴とするパケット転送装置。

【請求項8】

前記振り分け手段は、前記バッファに入力したパケットを、そのパケットの属するフローの識別情報と、各フローに対して定められた重み及び入力したパケットの長さの少なくとも一方とに基づいて前記複数の集合のいずれか1つに振り分けることを特徴とする請求項7記載のパケット転送装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、例えば、ATM通信網におけるセル多重化装置およびセルバッファ装置に関する。

【0002】

10

20

30

40

50

また、セルスイッチはセル多重化装置が複数個重ね合わされていることから、本発明のセル多重化装置の構成をセルスイッチのひとつの出力ポートに対応する構成としてセルスイッチに適用することが可能である。

【0003】

本発明のセルバッファ装置は、セルに限らず一般的なバッファ装置にも適用可能である。

【0004】

さらに、本発明はATM通信網のみならず、パケット交換網にも適用可能である。

【0005】

【従来の技術】

現在、ATM (Asynchronous Transfer Mode) 通信方式に関する研究が、世界中の通信技術の研究者らによって精力的に行なわれている。ATM通信方式は情報をセルという固定長のパケットにより伝送交換する。ATM通信方式ではスイッチノード内のハードウェアによるセルスイッチにより、高速なセルの交換が可能で単位時間当たりの情報転送能力は既存の通信網を越えるものを実現可能である。

10

【0006】

ATM通信方式は、セルのヘッダのVPI (Virtual Path Identifier) とVCI (Virtual Channel Identifier) と呼ばれる識別情報により、ひとつの物理伝送路に論理的に複数のコネクション (Virtual Connection: VC) を設定することができる。網内の各スイッチノードにおいては、各VCに対して予めルートが定められており、スイッチノードはセルのコネクション識別子VPIとVCIからセルを出力すべき出方路を求める。VPIとVCIは、スイッチノード間の物理伝送路で一意に割り当てられるため、スイッチノードは、通過するセルのVPIとVCIの値を書き換える能力を持つ。

20

【0007】

これまで、ATM網において品質を保証されたVCは、CBR (Constant Bit Rate: 固定ビットレート) コネクションかVBR (Variable Bit Rate: 可変ビットレート) コネクションが中心であった。CBRコネクションは、セルの伝送速度 (セルレートまたは帯域ともいう。単位時間当たりの伝送セル数) が一定で予めわかっているトラヒックを伝送するVCであり、VBRコネクションはセルの伝送速度が一定ではないが、その最大値 (ピークレート) と平均値 (平均レート) などのトラヒックの性質が予めわかっているVCである。

30

【0008】

基本的には、1本の物理伝送路に複数のVCを十分な品質を保ちつつ多重化する場合には、全てのVCのピークレートの和が物理伝送路の帯域以下になっていればよい。この手法をピークレート割り当てと呼ぶ。CBRコネクションのみをピークレート割り当てした場合には、物理伝送路の十分に高い利用効率が達成可能である。VBRコネクションの場合には、ピークレート割り当てでは、物理伝送路の利用効率を高くできない。そこで予め分かっているトラヒックの性質より、統計的多重化効果を用いて品質を保ちつつ利用効率を上げる技術がさかんに検討されている。

【0009】

ところが、計算機間のATM通信を考えると、平均レートを始めとしたトラヒックの性質が予め予測できないという性質や、瞬間的に大量のセルを送信するが、送信しないときには全くセルを送信しないというバースト性と呼ばれる性質がある。そのため、CBRやVBRの様に品質を保ちつつ網の利用効率を上げることは難しい。つまり計算機間で転送されるデータは、ピークレート割り当てなどで品質を保ちようとするとも網の利用効率が著しく低下し、VBRの様に統計的多重化効果を用いるとトラヒックのバースト性のためにセルスイッチのある出力ポートに同時に大量のセルが到着し、セルスイッチのバッファ量が十分でないと、バッファ溢れによるセル廃棄が発生してしまう。また、セル廃棄が発生すると複数のセルで構成されているパケット単位で再送が発生し、これにより、実効的なスループットが低下する。

40

50

【 0 0 1 0 】

そこで近年、端末とスイッチノード間でフロー制御をかけて、セルの転送品質（特にセル廃棄に関する品質）を保証し、かつ、網の利用効率を上げる ABR（Available Bit Rate）というサービスクラスが提案され、検討が進んでいる。ABR コネクションは、経由するスイッチノードがセル輻輳に落ちりそうになった場合には、スイッチノードがセル輻輳の発生前に、送信端末にセルの送迎を迎えるように要求する。端末へのトラフィック制御情報は、主に RMセル（リソース管理セル）と呼ばれるセルを用いて行なう。ABRにおけるスイッチノードから送信端末へのトラフィック制御には無視できない遅延時間が存在する。そのため、トラフィック制御が有効に作用するまでセルを廃棄しないようにセルスイッチは大容量のバッファを実装する必要がある。

10

【 0 0 1 1 】

CBR、VBR、ABRの他のサービスクラスとして、UBR（Unspecified Bit Rate）というサービスクラスが存在する。このクラスは、端末が出力するトラフィック特性を詳細に網に申告することを必要としない。そのかわり、網はその転送品質について一切の保証をしない、いわゆるベストエフォート（Best Effort）サービスのクラスである。

【 0 0 1 2 】

前述したように、計算機間のデータはバースト性を持っているため、UBRコネクションのセル廃棄率を満足できるものとするためにはセルスイッチに大容量のセルバッファを実装する必要があると考えられている。

20

【 0 0 1 3 】

幸いなことに計算機間のトラフィックは転送の遅延時間、遅延揺らぎに関する要求がCBRやVBRと比較して厳しくはない場合が多い。容量の大きなバッファをセルスイッチに実装することによりセルの伝送遅延時間、遅延揺らぎが増大するが、それを許容できるアプリケーションは決して少なくないと考えられる。

【 0 0 1 4 】

特にABR、UBRサービスの場合、網の輻輳を回避する手段が必要となると考えられる。輻輳回避手段のひとつとして、EFCI（Explicit Forward Congestion Indication）と呼ばれる方法がある。セルのヘッダ内に、そのセルの輻輳経験を書き込むEFCIビットがあり、網内のセルスイッチは、その輻輳状況に応じてEFCIをマークする。端末はその情報を利用することにより輻輳を回避することが可能となる。

30

【 0 0 1 5 】

次に、VC間公平キューイングについて説明する。

【 0 0 1 6 】

まず、ABRサービスにおける耐故障性から要求されるVC間公平キューイングについて述べる。

【 0 0 1 7 】

ABRサービスは、網からのトラフィック制御情報に従って送信端末がセルの送迎を制御することによって成立する。もし、ある端末が故障し（もしくは故意に）、網からのトラフィック制御情報を無視すれば、網の輻輳からの脱出が困難になるおそれがある。

40

【 0 0 1 8 】

このような問題は、セル多重化装置またはセルスイッチにおいてVC間で公平なキューイングを行ない、VC間の公平なセル多重化スケジューリングとVC間の公平なバッファ割り当てを実現することにより解決できる。VC間公平キューイングを行なえばVC間の相互作用は最小限に抑えられ、もしトラフィック制御情報を無視するVCがあっても、そのVCのみが輻輳に落ちり、その他のVCにおよぼす影響を低減できるからである。

【 0 0 1 9 】

また、UBRサービスにおける公平性からもVC間公平キューイングが要求される。

【 0 0 2 0 】

50

図48に、UBRサービスにおける不公平な帯域割り当ての例を示す。この例では網に端末A、B、C、Dの4つの端末と、ひとつのファイルサーバが存在し、それぞれがセルスイッチ（またはセル多重化装置）X、Y、Zとリンク1、2、3、4、5、6、7により接続されている。

【0021】

ここで、各セルスイッチが入力リンクに設定されているVC数に関わらずセルを多重化したとすると、各セルスイッチは入力ポートを平等に扱う。例えば、セルスイッチXは、リンク2とリンク3からのセルをリンク1へ多重化する場合に、両方のリンクを平等に扱うため、リンク1の帯域を1.0とすると、リンク2とリンク3にそれぞれ0.5の帯域を与える。

10

【0022】

同様に考えていくと最終的に各端末に与えられる帯域は、端末Aが最も大きく0.5、端末C、Dは最も少なく0.125である。全ての端末が平等にファイルサーバにアクセスすることが理想だとすれば、これは理想とはほど遠い。

【0023】

またさらに、網が輻輳に陥った場合、網は端末に帯域を小さくするようにEFCI情報を送信する。このとき、網内にEFCI情報に忠実に従う端末と、無視する端末が存在する場合、EFCIに忠実な端末のみがセルの送出を抑制するため、EFCIを無視する端末は不当に多くの帯域資源を確保することになってしまうという問題点もある。

【0024】

これらの問題は、セルスイッチが入力リンクに多重されているVC数を見捨てて各入力リンクを平等に扱っていることによる。セルスイッチがVC間公平キューイングを行ない、VC数に応じた帯域の分配を行なうことにより、各端末は公平にファイルサーバにアクセスすることが可能となる。

20

【0025】

次に、クラス間の優先制御を行なうセル多重化装置について説明する。

【0026】

図49に従来のクラス間の優先制御を行なうセル多重化装置の構成を示す。

【0027】

図49は、N本の入力ポートから入力したセルを1本の出力ポートに多重化するセル多重化装置である。

30

【0028】

出力バッファは、内部のクラス毎のセル蓄積部にクラス毎にセルを蓄積し、出力ポートへはクラスの優先度に応じてセルをセル選択部が選択し出力する。

【0029】

出力バッファの各蓄積部に蓄積しているセル数としきい値とを比較してバックプレッシャ信号を生成する。入力バッファもまたクラス毎のセル蓄積部が設けられており、バックプレッシャ信号により出力が許可されているクラスの中からクラスの優先度に応じてセルを出力バッファへ転送する。

【0030】

この構成では出力バッファに複数のクラス毎セル蓄積部を設ける必要があるため、クラス数が増えると出力バッファの実現が困難になるという欠点があった。特に、出力バッファの入力速度は入力ポートの速度のN倍（Nは入力ポート数）である必要があり、クラス数、入力ポート数が大きな場合は、このような複雑な機能を実現することは困難であった。

40

【0031】

図49を、セルスイッチの各出力ポートに関する部分と考えると、入力ポートにバッファを持つセルスイッチについても同じ問題点が存在する。

【0032】

次に、セル多重化装置について説明する。

【0033】

50

図 5 0 に示すような入力バッファを持つセル多重化装置の問題点を説明する。

【 0 0 3 4 】

図 5 0 は、N 本の入力ポートから入力されたセルを 1 本の出力ポートに多重化するセル多重化装置である。

【 0 0 3 5 】

出力バッファは、入力ポートからのセルを一時的に蓄積し、出力ポートの速度に応じてセルを出力する。

【 0 0 3 6 】

出力バッファに蓄積しているセル数としきい値とを比較してバックプレッシャ信号を生成する。入力バッファは出力ポートに対応してひとつのキューを持ち、バックプレッシャ信号により出力が許可されているときのみセルを出力バッファへ転送する。

10

【 0 0 3 7 】

この様なセル多重化装置は設定されている V C 数に関わらず各入力バッファを公平に扱ってしまうため、A B R サービスにおける耐故障性や、U B R サービスにおける公平性に問題があった。

【 0 0 3 8 】

図 5 0 を、セルスイッチの各出力ポートに関する部分と考えると、入力ポートにバッファを持つセルスイッチについても同じ問題点が存在する。

【 0 0 3 9 】

次に、V C 間公平キューイングについて説明する。

20

【 0 0 4 0 】

図 5 1 に V C 間公平キューイングの実現方法である V C 毎の F I F O の構成を示す。

【 0 0 4 1 】

図 5 1 は、入力リンクより入力したセルのコネクション識別情報をバッファポインタ管理部に通知し、バッファポインタ管理部よりセルの書き込み位置を示す書き込みポインタを得てセルを一時的にセルバッファに蓄積し、バッファポインタ管理部からの、読みだしセルを示す読みだしポインタに基づいてセルをセルバッファから読みだし出力リンクへ出力するセルバッファ装置である。バッファポインタ管理部は、蓄積しているセルのセルバッファ上の位置を示すバッファポインタを管理する。

【 0 0 4 2 】

30

バッファポインタ管理部は V C テーブルによりバッファポインタを V C 毎に管理する。図 5 1 は広く知られているポインタチェーンによる管理方法を示している。

【 0 0 4 3 】

セルが入力すると空きチェーンの先頭からバッファポインタを書き込みポインタとして一つ取りだし、入力セルに対応する V C のチェーンの末尾につける。セルを出力する時は、セルを蓄積している V C を公平に選択し、その先頭のバッファポインタを読みだしポインタとして取りだす。セル出力時に V C を公平に選択する方法として、ラウンドロビンが知られている。

【 0 0 4 4 】

V C をラウンドロビンで選択するためには、V C テーブル上で、前回出力した V C の次から、バッファに蓄積しているセルがあるかどうかを順に検索しなければならない。検索は、多い場合にはほぼ全ての V C を 1 セル周期で行なう必要があるが、設定すべき V C 数は数千以上に及ぶ場合があり、実現は困難であった。

40

次に、キュー長監視について説明する。

【 0 0 4 5 】

図 5 2 に、例えば図 5 0 の様な、入力バッファを持つセル多重化装置における、入力バッファのキュー長の変化の一例を示す。

【 0 0 4 6 】

横軸は時間、縦軸は上が入力バッファのキュー長、下がバックプレッシャ信号である。バックプレッシャ信号が出力禁止を指示しているとき (T up)、入力バッファからの出力ス

50

ループットはゼロであるためキュー長はバッファの入力レート (R_i) で増加する。バックプレッシャ信号が出力可能を指示しているとき、入力バッファからの出力スループットは最大であり、通常キュー長は減少する。

【0047】

この様にキュー長はバックプレッシャ信号に同期して振動する。従来、バッファにセルが到着したとき、キュー長がしきい値を越えていれば入力したセルを受け付けず廃棄し、越えていなければ受け付けることが広く行なわれている。この方法は通常のバッファでは問題ないが、この様にバッファの出力がバックプレッシャ信号に制御されている場合に問題となる。キュー長の最大値は入力レート (R_i) とバックプレッシャ信号が出力禁止を指示している時間 (T_{up}) によって決定され、特にバックプレッシャ信号の出力禁止時間 (T_{up}) は、他の入力バッファへの入力トラヒックなどに大きく影響されるからである。

10

【0048】

バックプレッシャ信号により制御されるバッファ装置の場合は、単純にキュー長としきい値と比較してセルの廃棄を決定する方法では外部条件によりセルの廃棄が過敏に影響され安定したセルの廃棄が困難であった。キュー長から輻輳を判断する場合にも同じ問題点があった。

【0049】

次に、優先制御を行うセルバッファ装置について説明する。

【0050】

図53に、優先制御を行なうセルバッファ装置の構成を示す。

20

【0051】

遅延に関して低優先のセルを一時的に蓄積する、クラス1セル蓄積部と、高優先のセルを一時的に蓄積するクラス2セル蓄積部とあって、それぞれの出力はクラス多重化FIFOにより多重化される。クラス管理部は、クラス2セル蓄積部セル数 N_a およびクラス多重化FIFOセル数 N_b を入力し、どちらかのクラスへ転送指示を与える。転送指示を受けたクラスのセル蓄積部は1セルだけクラス多重化FIFOへ転送する。

【0052】

このセルバッファ装置の優先制御の性能を高めるためにはクラス多重化FIFOの蓄積セル数を少なくすることが必要であるが、スループットを低下させないためには、クラス多重化FIFOを努めて空にしない (= アンダフローしない) ように転送を指示する必要がある。さらに、どちらか一方のセル蓄積部が空になっている時に、空のセル蓄積部に対して転送を指示してしまった場合 (空指示と呼ぶ) にはアンダフローしてしまう可能性があるため、転送指示は蓄積セルが存在するセル蓄積部に対して行なう必要がある。

30

【0053】

以上の方針から、クラス管理部がクラス2セル蓄積部へ転送指示を行なう条件は、

$$(N_b > 1) \text{ and } (N_a \geq 1)$$

で与えられる。クラス1セル蓄積部への転送指示は、

$$(N_b > 1) \text{ and } (N_a = 0)$$

である。

【0054】

ここで、クラス毎のセル蓄積部やクラス多重化FIFOからクラス管理部へのセル数情報の伝送や、クラス管理部からクラス毎セル蓄積部への転送指示情報の伝送、さらにはクラス毎セル蓄積部からセル多重化FIFOへのセル転送に遅延時間が存在する場合を考える。このとき従来は次のような問題点が存在していた。

40

【0055】

これらの遅延時間により、クラス多重化FIFOがアンダフローするかどうかクラス管理部で正確に判定できなかった。そこでクラス多重化FIFO内セル数をこの遅延時間に応じた大きめのしきい値と比較してアンダフローしないための条件としていた。この結果、クラス多重化FIFOの平均キュー長が大きくなり優先制御の性能が低下した。

【0056】

50

また、高優先であるクラス2セル蓄積部情報の伝送が遅れることにより、クラス2セル蓄積部が空になったことを正確に判断できない。したがって空になっているのに空でないと誤り、空指示をしてしまうか、空でないのに空であると誤り、高優先であるクラス2のセルよりも低優先のクラス1のセルが優先して出力してしまうような指示を行なってしまっていた。結果としてセルバッファ装置のスループットが低下したり優先制御機能の性能が低下する問題点があった。

【0057】

以上、図53の優先制御を行なうセルバッファ装置に関して問題点を述べた。この問題はこのような優先制御を行なうセルバッファ装置のみならず、さまざまなセルバッファ装置において発生し得る。

10

【0058】

つまり、複数のバッファ間のセルの転送を管理部が指示するセルバッファ装置において、管理部へのバッファ情報の伝送、管理部からバッファへの転送指示情報の伝送、バッファ間のセルの伝送に遅延時間が存在する場合、セルバッファ装置の性能やスループットが低下してしまうという問題点があった。

【0059】

【発明が解決しようとする課題】

以上述べたように、従来のセル多重化装置には次のような第1、第2の2つの問題点があった。

【0060】

20

まず、出力バッファに複数のクラス毎セル蓄積部を設ける必要があるため、クラス数が多くなると出力バッファの実現が困難になるという欠点があった。特に、出力バッファの入力速度は入力ポートの速度のN倍(Nは入力ポート数)である必要があり、クラス数、入力ポート数が大きな場合は、このような複雑な機能を実現することは困難であった。

【0061】

次に、設定されているVC数に関わらず各入力バッファを公平に扱ってしまうため、ABRサービスにおける耐故障性や、UBRサービスにおける公平性に問題があった。

【0062】

また、従来のセルバッファ装置には次のような第3、第4の2つの問題点があった。

【0063】

30

まず、VCをラウンドロビンで選択するためには、VCテーブル上で、前回出力したVCの次から、バッファに蓄積しているセルがあるかどうかを順に検索しなければならない。検索は、多い場合にはほぼ全てのVCを1セル周期で行なう必要があるが、設定すべきVC数は数千以上に及ぶ場合があり、実現は困難であった。

【0064】

これに対し近年、"Efficient Fair Queueing Deficit Round Robin"(pp.231-242、ACM SIGCOMM '95、1995年8月)に記載されているようなアルゴリズムが提案されている。この文献で提案されているDRRというアルゴリズムはアクティブリストという概念を用いてセルを蓄積しているVCを管理するため、セル出力時のVCテーブルの検索が不要である。従って、設定すべきVC数が増えても高速にVC間公平キューイングの為の出力すべきセルの選択が実現できる。しかしながら、このDRRアルゴリズムでは、パケットを各VC毎に分類して蓄積するため各VCに設定された重みの値に応じた量のパケットが一度に出力されてしまい、その出力トラヒックのバースト性が非常に高くなってしまいうという欠点があった。

40

【0065】

次に、バックプレッシャ信号により制御されるバッファ装置の場合は、単純にキュー長としきい値と比較してセルの廃棄を決定する方法では外部条件によりセルの廃棄が過敏に影響され安定したセルの廃棄が困難であった。キュー長から輻輳を判断する場合にも同じ問題点があった。

50

【 0 0 6 6 】

さらに、複数のバッファ間のセルの転送を管理部が指示するセルバッファ装置においては、次のような第5の問題点があった。すなわち、管理部へのバッファ情報の伝送、管理部からバッファへの転送指示情報の伝送、バッファ間のセルの伝送に遅延時間が存在する場合、セルバッファ装置の性能やスループットが低下してしまうという問題点があった。

【 0 0 6 7 】

本発明は、上述の点に鑑みてなされたもので、その目的とするところは、入力ポート数が大きな場合においても実現が容易なクラス間の優先制御を行なうセル多重化装置を提供することにある。

【 0 0 6 8 】

本発明の他の目的は、A B R サービスにおける耐故障性や、U B R サービスにおける公平性を実現するV C 毎公平キューイングを行なうことができるセル多重化装置を提供することにある。

【 0 0 6 9 】

本発明の別の目的は、設定可能なV C 数の上限に依存することなく容易に実現することのできるV C 間公平キューイングを行なうセル多重化装置を提供することにある。

【 0 0 7 0 】

本発明のさらに別の目的は、キュー長の監視結果が外部条件の影響を受け難く、安定したキュー長の監視が容易であるセル多重化装置を提供することにある。

【 0 0 7 1 】

本発明のもう一つの目的は、複数のバッファ間のセルの転送を管理部が指示する際に、管理部へのバッファ情報の伝送、管理部からバッファへの転送指示情報の伝送、バッファ間のセルの伝送に遅延時間が存在する場合においても、従来のように性能やスループットが低下しないセル多重化装置を提供することにある。

【 0 0 7 2 】

【 課題を解決するための手段 】

本発明の packets 転送装置（請求項 1 ~ 4）は、入力された packets を一時的に蓄積する複数の入力バッファと、これら入力バッファを制御する制御手段と、前記各入力バッファから出力される packets を転送する少なくとも 1 つの出力ポートを具備した packets 転送装置であって、

前記各入力バッファは、

入力した packets を一時的にクラス毎に蓄積するクラス毎の蓄積手段と、

前記制御手段から指示されたクラスの packets を前記蓄積手段から前記出力ポートへ向けて出力する出力手段とを具備し、

前記制御手段は、

前記蓄積手段の複数の入力バッファ全体における蓄積状況をクラス毎に把握し、この蓄積状況に基づいて packets を出力すべきクラスを決定し、この決定したクラスの指定を含む指示を前記複数の入力バッファに送信することにより、

制御部が複数の入力バッファ全体の状況を把握して（少なくとも全体の状況を把握していればよいが、個々の入力バッファ毎に把握していても構わない）指示を出すため、各入力バッファのクラス毎の蓄積状況の違いによりクラス間の優先関係が崩れてしまうことが生じず、また、入力ポートの数すなわち入力バッファの数が増えても容易にクラス間の優先制御を実現できる。

【 0 0 7 3 】

また、本発明の packets 転送装置（請求項 5、6）は、複数の入力バッファと、これら入力バッファを制御する制御手段と、前記各入力バッファから出力される packets を転送する少なくとも 1 つの出力ポートを具備した packets 転送装置であって、

前記各入力バッファは、

入力した packets を一時的に蓄積する蓄積手段と、

前記制御手段からの指示を受けて、次のフェーズで前記蓄積手段から前記出力ポートへ向

10

20

30

40

50

けて出力すべきパケットを選択する選択手段と、
この選択手段で選択されたパケットを前記出力ポートに向けて出力する出力手段とを具備し、

前記制御手段は、

前記複数の入力バッファ全体における前記選択手段で選択されたパケットの出力状況を把握し、この出力状況に基づいて新たなパケットを選択するよう前記複数の入力バッファに指示することにより、

あるフェーズ（その時間的長さは不定）内で出力されるべきパケット集合（カゴ）という概念を導入し、このパケット集合に入れるパケットの選択と、このパケット集合に入っているパケットのみを出力する（その他のパケットは出力が許可されない）という処理を各入力バッファが分担し、パケット集合に新たなパケットを入れるタイミングを指示する処理を制御部（カゴ管理部）が分担することにより、例えば、入力バッファ方式のセルスイッチにおいて、V C間公平キューイングを容易に実現できる。

10

【0074】

また、本発明のパケット転送装置（請求項7、8）は、入力したパケットを一時的に蓄積するバッファと、このバッファを制御する制御手段と、前記バッファから出力されたパケットを転送するパケット転送装置であって、

前記制御手段は、

前記バッファ内に蓄積されたパケットを複数の集合に分けて管理する管理手段と、

前記バッファに入力したパケットを各パケットの属するフロー間で公平になるように前記複数の集合のいずれか1つに振り分ける振り分け手段と、

20

前記管理手段で管理されている前記複数の集合のうちの1つの集合に属するパケットを出力するよう前記バッファに対し指示する指示手段と、

を具備することにより、制御手段（バッファポインタ管理部）でパケットが入力されたときに複数の集合への振り分けが行われており、パケットを出力するときには1つの集合に属するポインタを出すだけで済むので、設定すべきV C数が増えても高速にV C間公平キューイングのための出力すべきパケットの選択が実現できる。また、本発明の前記指示手段により、請求項5の選択手段の動作を実現することができる。さらに、本発明は、出力バッファ方式のセルスイッチでV C間公平キューイングをする際の出力バッファとして用いることができる。

30

【0077】

【発明の実施の形態】

以下、本発明の実施形態について図面を参照して説明する。

【0078】

0. 用語の定義

以下、本発明の実施の形態を説明する上で用いられる用語の定義を行う。

【0079】

（1）パケットとは、固定長のセルを含む概念で、セルの上位概念と解することができる。また、以下の説明でパケットというときは、可変長、固定長を問わない。また、セルに限定して説明している場合でも、それは単に説明を簡単にするためのもので、特に、これに限るものではない。

40

【0080】

（2）パケット（セル）多重化装置とパケット（セル）バッファ装置の違い：

パケットバッファ装置：一つの以上の入力ポートと一つ以上の出力ポート及び一つ以上のバッファを具備し、入力ポートから入力されたパケットを必要に応じて一時的にバッファに蓄積し、出力ポートへ出力する装置。特にパケットを宛先に応じた出力ポートへ出力するパケットバッファ装置は、パケットスイッチと呼ばれる。

【0081】

パケット多重化装置：パケットバッファ装置の中でも特にパケットの出力ポートが一つのものでパケット多重化装置と呼ぶ（パケットバッファ装置を出力ポート毎に分けて考える

50

と、入力ポートが共通な複数のパケット多重化装置の重ね合わせてとしてみる事ができる)。以下の説明において、パケット多重化装置について述べていることは、パケットバッファ装置にも容易に適用できる。

【0082】

(3) パケット転送装置：パケットバッファ装置と、その一種であるパケット多重化装置の双方を含む。また、以下の説明において、多重化装置あるいはバッファ装置を例にとり説明している場合も、特にこれに限るものではなく、広くパケット転送装置として適用できる。

【0083】

1. クラス間優先制御(第1の実施形態、第2の実施形態)

10

まず、本発明に係るクラス間の優先制御を行なうセル多重化装置の実施形態について説明する。

【0084】

1.1 第1の実施形態

図1に、第1の実施形態に係るクラス間の優先制御を行なうセル多重化装置の構成例を示す。

【0085】

図1のセル多重化装置は、入力ポート#1~#Nのそれぞれに対応して設けられ、各入力ポート#1~#Nから入力したセルを一時的に蓄積する複数の入力バッファ10と、入力バッファ10の出力したセルを多重化し出力ポートへ出力する出力バッファ11と、全ての入力バッファを管理するクラス管理部12とを備える。

20

【0086】

入力バッファ10は、入力したセルをクラス(クラス1~クラス3)毎に蓄積するクラスのセル蓄積部13(クラス1セル蓄積部13a、クラス2セル蓄積部13b、クラス3セル蓄積部13c)と、セル蓄積部13から出力したセルを多重化し出力バッファ11へ出力するクラス多重化FIFO14を備えている。

【0087】

クラス管理部12は、セル蓄積部13のクラス毎セル蓄積部セル数と、クラス多重化FIFO14のクラス多重化FIFOセル数を得て、予め定められたアルゴリズムより転送クラス指示を決定して全ての入力バッファ10へ同報する。

30

【0088】

入力バッファ10に同報された転送クラス指示で指示されたクラスのセル蓄積部13は、セルをクラス多重化FIFO14へ転送し、クラス多重化FIFO14は、出力バッファ11内の輻輳状態に応じたバックプレッシャ信号の制御に従って出力バッファ11へセルを出力する。

【0089】

クラス毎セル蓄積部13は通常はFIFO(First In First Out)メモリで構成される。

【0090】

図1は3クラス(クラス1~クラス3)の場合を示したが、第1の実施形態はクラス数に関わらず有効に作用する。

40

【0091】

また、図1を、セルスイッチの、各出力ポートに関する部分と考えるとセルスイッチについても第1の実施形態を適用できる。

【0092】

図1の出力バッファ11は、単段のFIFO15から構成されているが、どんな状況においてもクラス多重化FIFO14のセルをある有限の時間内で出力ポートに出力できることを保証できれば、どんな形態でもよい。

【0093】

例えば、バッファを全く持たずにクラス多重化FIFO14のセルを調停して出力ポートに出力する構成でもよい。また例えば、後述の図3または図4に示す構成でもかまわない

50

。

【 0 0 9 4 】

図 3 に示した出力バッファ 1 6 は、F I F O が複数段（図 3 では 2 段）に構成されており、それぞれ前段の F I F O にバックプレッシャ信号が接続される。

【 0 0 9 5 】

すなわち、図 3 の出力バッファ 1 6 は、入力側（前段）の複数の F I F O 1 8 に、入力バッファ 1 0 の各出力に対応して設けられる入力リンク # 1 ~ # N が接続され、それら F I F O 1 8 の各出力が出力側（後段）の F I F O 1 9 に接続されて構成され、後段の F I F O 1 9 から、その輻輳状態に応じて出力されるバックプレッシャ信号は、前段の複数の F I F O 1 8 のそれぞれに接続され、前段の複数の F I F O 1 8 のそれぞれから、その輻輳状態に応じて出力されるバックプレッシャ信号は、複数のクラス多重化 F I F O 1 4 のそれぞれに接続される。

10

【 0 0 9 6 】

図 4 に示した出力バッファ 1 7 は、図 3 と同様に F I F O が複数段（図 4 では 2 段）に構成されている。図 3 と異なるのは、入力側の複数の F I F O 1 8 のキュー長情報（または全ての F I F O のキュー長情報でもよい）をバックプレッシャ生成部 2 0 へ転送し、バックプレッシャ生成部 2 0 はキュー長情報からバックプレッシャ信号を生成し、全ての入力リンクへ同報することである。

【 0 0 9 7 】

例えば、キュー長を全て合計してしきい値と比較し、しきい値を越えていればバックプレッシャ信号により出力バッファへの入力を禁止する。バックプレッシャ生成部 2 0 はクラス管理部 1 2 に存在しても良い。

20

【 0 0 9 8 】

第 1 の実施形態のクラス間の優先制御を行うセル多重化装置によれば、高いスループットが必要な出力バッファにてクラスを識別する必要がないため、容易に実現することが可能となる。入力バッファはクラス毎にセルを管理する必要があるが、入力バッファに必要なスループットは低いため実現が容易となる。

【 0 0 9 9 】

図 1 のセル多重化装置の構成をさらに詳細に説明する。

【 0 1 0 0 】

クラス管理部 1 2 に入力されるクラス毎セル蓄積部セル数は、クラス毎に、全ての入力バッファ 1 0 のクラス毎セル蓄積部セル数の合計でよい。また、クラス多重化 F I F O セル数も、全ての入力バッファ 1 0 の合計でよい。

30

【 0 1 0 1 】

合計の演算はクラス管理部 1 2 の内部で行なっても良いし、クラス管理部 1 2 の外部で行なっても良い。

【 0 1 0 2 】

セル数情報はクラス管理部 1 2 において比較的小さなしきい値と大小を判定する。しきい値が固定値であるか動的に変化する値であるかはクラス管理部 1 2 のアルゴリズムに依存するが、もししきい値がある固定値の場合はセル数の代わりにその比較結果をクラス管理部に入力してもよい。または、ある値以上はどのような値でもクラス管理部 1 2 の処理に影響を与えないのなら、それを利用してその値以上は一つの符号に符号化して入力しても良い。例えば 4 ビットで、セル数を (0 0 0 0) = 0、(0 0 0 1) = 1、(0 0 1 0) = 2、(0 0 1 1) = 3、...、(1 1 0 1) = 1 3、(1 1 1 0) = 1 4、(1 1 1 1) = 1 5 以上、という 1 6 段階に符号化するなどとしてもよい。これによりクラス管理部 1 2 への入力情報を圧縮することができ、実装が容易になるという利点がある。

40

【 0 1 0 3 】

セル数の合計を扱うことにより、クラス管理部 1 2 の実装を簡単化することができる。転送クラス指示も全ての入力バッファに同じ情報を同報することから、クラス管理部は入力バッファを個々に認識する必要がない。よって入力ポート数が多くなっても本質的な実装

50

の困難度は変化しないという利点がある。

【 0 1 0 4 】

クラス管理部 1 2 が転送指示を決定するアルゴリズムは例えば次のようなものである。

【 0 1 0 5 】

1 . クラス多重化 F I F O セル数より、クラス多重化 F I F O がアンダフローしないかどうかを判断し、アンダフローする可能性があるとは判断した場合のみアンダフローしないように転送指示を行なう。アンダフローとは、クラス多重化 F I F O セル数がゼロになり、本来出力できるはずのクラス毎セル蓄積部のセルが有効に出力できない状態をいう。

【 0 1 0 6 】

2 . クラス毎セル蓄積部セル数より、入力バッファ 1 0 にセルが蓄積されているクラスを転送指示の候補とする。 10

【 0 1 0 7 】

3 . 転送指示を行なうことができる場合に、転送指示の候補となっているクラスの中で最も優先度の高いクラスを求め、そのクラスの転送を指示する。

【 0 1 0 8 】

クラス管理部 1 2 は出力バッファの状態とバックプレッシャ信号を直接考慮せずに転送クラス指示を行なう。そのため、転送クラス指示によるセルの転送スループットが、入力バッファ 1 0 からのセル出力スループットを、上回ることがあり、その差を一時的に吸収するためにクラス多重化 F I F O 1 4 が用意されている。クラス管理部 1 2 はクラス多重化 F I F O セル数を知ること、転送クラス指示によるセル転送のスループットが、入力バッファ 1 0 からのセル出力スループットを、長時間連続して越えないように制御する。 20

【 0 1 0 9 】

クラス毎セル蓄積部 1 3 からクラス多重化 F I F O 1 4 へ実際に転送したセルを全ての入力バッファ 1 0 で合計したセル数（クラス毎転送セル数）をクラス管理部 1 2 が知ることにより、より詳細なクラス間の優先制御を行なうことができる。

【 0 1 1 0 】

この場合、クラス管理部 1 2 が転送指示を決定するアルゴリズムは例えば次のようなものである。

【 0 1 1 1 】

1 . クラス多重化 F I F O 内セル数より、クラス多重化 F I F O がアンダフローしないかどうかを判断し、アンダフローする可能性があるとは判断した場合のみアンダフローしないように転送指示を行なう。アンダフローとは、クラス多重化 F I F O セル数がゼロになり、本来出力できるはずのクラス毎セル蓄積部 1 3 のセルが有効に出力できない状態をいう。 30

【 0 1 1 2 】

2 . クラス毎セル蓄積部セル数より、入力バッファにセルが蓄積されているクラスを転送指示の候補とする。

【 0 1 1 3 】

3 . クラス毎転送セル数より、クラス毎に予め定められたスループット以上のスループットを得ているクラスは、その程度に応じて転送指示する優先度を下げる。または転送指示の候補としない。 40

【 0 1 1 4 】

4 . クラス毎転送セル数より、クラス毎に予め定められたスループット以下のスループットしか得ていないクラスは、その程度に応じて転送指示する優先度を上げる。

【 0 1 1 5 】

5 . 転送指示を行なうことができる場合に、転送指示の候補となっているクラスの中で最も優先度の高いクラスを求め、そのクラスの転送を指示する。

【 0 1 1 6 】

1 . 2 第 2 の実施形態

次に、第 2 の実施形態について説明する。

【 0 1 1 7 】

図 2 に、第 2 の実施形態に係るクラス間の優先制御を行なうセル多重化装置の他の構成例を示す。

【 0 1 1 8 】

図 2 のセル多重化装置は、入力ポート # 1 ~ # N のそれぞれに対応して設けられ、各入力ポート # 1 ~ # N から入力したセルを一時的に蓄積する複数の入力バッファ 1 0 0 と、入力バッファ 1 0 0 の出力したセルを多重化し出力ポートへ出力する出力バッファ 1 1 1 と、全ての入力バッファ 1 0 0 を管理するクラス管理部 1 1 2 とを備える。

【 0 1 1 9 】

入力バッファ 1 0 0 は、入力したセルをクラス (クラス 1 ~ クラス 3) 毎に蓄積するクラス毎のセル蓄積部 1 1 3 (クラス 1 セル蓄積部 1 1 3 a、クラス 2 セル蓄積部 1 1 3 b、クラス 3 セル蓄積部 1 1 3 c) と、セル蓄積部 1 1 3 から出力したセルを多重化し出力バッファ 1 1 1 へ出力するクラス多重化 F I F O 1 1 4 を備えている。

10

【 0 1 2 0 】

クラス管理部 1 1 2 は、セル蓄積部 1 1 3 のクラス毎セル蓄積部セル数と、出力バッファ内セル数を得て、予め定められたアルゴリズムより転送クラス指示を決定して全ての入力バッファ 1 0 0 へ同報する。

【 0 1 2 1 】

入力バッファ 1 0 0 に同報された転送クラス指示で指示されたクラスのセル蓄積部 1 1 3 は、セルをクラス多重化 F I F O 1 1 4 へ転送し、クラス多重化 F I F O 1 1 4 は、出力バッファ 1 1 1 内の輻輳状態に応じたバックプレッシャ信号の制御に従って出力バッファ 1 1 1 へセルを出力する。

20

【 0 1 2 2 】

クラス毎セル蓄積部 1 1 3 は通常は F I F O で構成される。

【 0 1 2 3 】

図 2 は 3 クラス (クラス 1 ~ クラス 3) の場合を示したが、第 2 の実施形態はクラス数に関わらず有効に作用する。

【 0 1 2 4 】

図 2 を、セルスイッチの、各出力ポートに関する部分と考えるとセルスイッチについても第 2 の実施形態を適用できる。

30

【 0 1 2 5 】

図 2 の出力バッファ 1 1 1 は、どんな状況においてもクラス多重化 F I F O のセルをある有限の時間内で出力ポートに出力できることを保証できれば、どんな形態でもよい。

【 0 1 2 6 】

例えば、図 4 に示す構成でもかまわない。

【 0 1 2 7 】

なお、この場合、バックプレッシャ生成部 2 0 はクラス管理部 1 1 2 に存在してもよい。出力バッファ内セル数 N_m は、バックプレッシャ生成部 2 0 にて加算し、図 2 のクラス管理部 1 1 2 へ出力する。

【 0 1 2 8 】

第 2 の実施形態のクラス間の優先制御を行うセル多重化装置によれば、高いスループットが必要な出力バッファにてクラスを識別する必要がないため、容易に実現することが可能である。入力バッファ 1 0 0 はクラス毎にセルを管理する必要があるが、入力バッファ 1 0 0 に必要なスループットは低いため実現が容易である。

40

【 0 1 2 9 】

図 2 のセル多重化装置の構成をさらに詳細に説明する。

【 0 1 3 0 】

クラス管理部 1 1 2 に入力されるクラス毎セル蓄積部セル数は、クラス毎に、全ての入力バッファ 1 0 0 のクラス毎セル蓄積セル数の合計でよい。

【 0 1 3 1 】

50

合計の演算はクラス管理部 1 1 2 の内部で行なっても良いし、クラス管理部 1 1 2 の外部で行なっても良い。

【 0 1 3 2 】

セル数情報はクラス管理部 1 1 2 において比較的小さなしきい値と大小を判定する。しきい値が固定値であるか動的に変化する値であるかはクラス管理部 1 1 2 のアルゴリズムに依存するが、もししきい値がある固定値の場合はセル数の代わりにその比較結果をクラス管理部 1 1 2 に入力してもよい。または、ある値以上はどのような値でもクラス管理部 1 1 2 の処理に影響を与えないのなら、それを利用してその値以上は一つの符号に符号化して入力しても良い。例えば 4 ビットで、セル数を (0 0 0 0) = 0、(0 0 0 1) = 1、(0 0 1 0) = 2、(0 0 1 1) = 3、...、(1 1 0 1) = 1 3、(1 1 1 0) = 1 4、(1 1 1 1) = 1 5 以上、という 1 6 段階に符号化するなどとしてもよい。これによりクラス管理部への入力情報を圧縮することができ、実装が容易になるという利点がある。

10

【 0 1 3 3 】

セル数の合計を扱うことにより、クラス管理部 1 1 2 の実装を簡単化することができる。転送クラス指示も全ての入力バッファ 1 0 0 に同じ情報を同報することから、クラス管理部 1 1 2 は入力バッファ 1 0 0 を個々に認識する必要がない。よって入力ポート数が多くなっても本質的な実装の困難度は変化しないという利点がある。

【 0 1 3 4 】

クラス管理部 1 1 2 が転送指示を決定するアルゴリズムは例えば次のようなものである。

【 0 1 3 5 】

1 . 出力バッファ 1 1 1 内セル数より、出力バッファ 1 1 1 がアンダフローしないかどうかを判断し、アンダフローする可能性がある場合のみアンダフローしないように転送指示を行なう。アンダフローとは、出力バッファセル数がゼロになり、本来出力できるはずのクラス毎セル蓄積部 1 1 3 のセルが有効に出力できない状態をいう。

20

【 0 1 3 6 】

2 . クラス毎セル蓄積部セル数より、入力バッファ 1 0 0 にセルが蓄積されているクラスを転送指示の候補とする。

【 0 1 3 7 】

3 . 転送指示を行なうことができる場合に、転送指示の候補となっているクラスの中で最も優先度の高いクラスを求め、そのクラスの転送を指示する。

30

【 0 1 3 8 】

クラス管理部 1 1 2 は出力バッファ 1 1 1 の状態を直接考慮して転送クラス指示を行なう。そのため出力バッファ 1 1 1 のバッファ量を十分用意すれば、バックプレッシャ信号が入力バッファからのセル出力を抑制することがなく、クラス多重化 F I F O 1 1 4 は基本的に不要である。クラス管理部 1 1 2 は出力バッファ内セル数を知ることにより、転送クラス指示によるセル転送のスループットが、出力バッファ 1 1 1 からのセル出力スループットを、長時間連続して越えないように制御する。

【 0 1 3 9 】

このように入力バッファ 1 0 0 へのバックプレッシャ信号およびクラス多重化 F I F O 1 1 4 は基本的に不要であるが、出力バッファ 1 1 1 のバッファ量が十分でない場合には有効に作用する。これらを備えていれば無い場合に比べ出力バッファ 1 1 1 のバッファ量を少なくできる。

40

【 0 1 4 0 】

クラス毎セル蓄積部 1 1 3 からクラス多重化 F I F O 1 1 4 へ実際に転送したセルを全ての入力バッファで合計したセル数(クラス毎転送セル数)をクラス管理部 1 1 2 が知ることにより、より詳細なクラス間の優先制御を行なうことができる。

【 0 1 4 1 】

この場合、クラス管理部 1 1 2 が転送指示を決定するアルゴリズムは例えば次のようなものである。

【 0 1 4 2 】

50

1. 出力バッファ 1 1 1 内セル数より、出力バッファ 1 1 1 がアンダフローしないかどうかを判断し、アンダフローする可能性があるとは判断した場合のみアンダフローしないように転送指示を行なう。アンダフローとは、出力バッファセル数がゼロになり、本来出力できるはずのクラス毎セル蓄積部 1 1 3 のセルが有効に出力できない状態をいう。

【 0 1 4 3 】

2. クラス毎セル蓄積部セル数より、入力バッファ 1 0 0 にセルが蓄積されているクラスを転送指示の候補とする。

【 0 1 4 4 】

3. クラス毎転送セル数より、クラス毎に予め定められたスループット以上のスループットを得ているクラスは、その程度に応じて転送指示する優先度を下げる。または転送指示の候補としない。

10

【 0 1 4 5 】

4. クラス毎転送セル数より、クラス毎に予め定められたスループット以下のスループットしか得ていないクラスは、その程度に応じて転送指示する優先度を上げる。

【 0 1 4 6 】

5. 転送指示を行なうことができる場合に、転送指示の候補となっているクラスの中で最も優先度の高いクラスを求め、そのクラスの転送を指示する。

【 0 1 4 7 】

1. 3 第 1、第 2 の実施形態に係るセル多重化装置の利点

以上、説明したように、上記第 1、第 2 の実施形態のクラス間の優先制御を行なうセル多重化装置によれば、高いスループットが必要な出力バッファにてクラスを識別する必要がないため、入力ポート数が大きな場合においても実現が容易であるという利点がある。

20

【 0 1 4 8 】

2. 公平キューイング (第 3 ~ 第 7 の実施形態)

次に、本発明に係る V C 間公平キューイングを行なうセル多重化装置の実施形態について説明する。

【 0 1 4 9 】

2. 1 第 3 の実施形態 (出力バッファ型のセル多重化装置)

図 5 に、第 3 の実施形態に係る V C 間公平キューイングを行なうセル多重化装置の構成例を示す。

30

【 0 1 5 0 】

図 5 のセル多重化装置は、出力バッファ型と呼ばれ、入力ポート # 1 ~ # N のそれぞれに対応して設けられ、入力ポート # 1 ~ # N から入力したセルを、入力ポート数 N 倍の速度に速度変換する複数の速度変換回路 2 0 0 と、速度変換回路 2 0 0 の出力したセルを一時的に蓄積し出力ポートへ出力する出力バッファ 2 0 1 とを備える。

【 0 1 5 1 】

出力バッファ 2 0 1 は、入力したセルを V C 毎に分離する V C 分離部 2 0 1 a と、各 V C 毎に設けられる複数の V C 毎 F I F O 2 0 1 b と、複数の V C 毎 F I F O 2 0 1 b のそれぞれから出力されるセルを選択して出力ポートへ出力するセル選択部 2 0 1 c を備えている。

40

【 0 1 5 2 】

V C 毎 F I F O 2 0 1 b は十分良好なセル廃棄率を得るためには大容量でなければならない。またノンブロッキング条件より、V C 毎 F I F O 2 0 1 b の入力速度は入力ポートの N 倍でなければならない。入力ポートの速度が大きい場合や、入力ポート数が多い場合は、このようなスループットを満足する高速大容量のメモリは安価ではなく、かつ、このようなスループットを満足しつつ複雑な V C 毎のキュー管理を行なうことも困難である。

【 0 1 5 3 】

2. 2 第 4 の実施形態 (「カゴ」という概念を用いたセル多重化装置)

次に、第 4 の実施形態について説明する。

【 0 1 5 4 】

50

図 6 に、第 4 の実施形態に係るカゴと呼ぶ概念を用いたセル多重化装置の構成例を示す。

【 0 1 5 5 】

図 6 に示したセル多重化装置は、入力ポート # 1 ~ # N のそれぞれに対応して設けられ、入力ポート # 1 ~ # N から入力したセルを一時的に蓄積する複数の入力バッファ 2 1 0 と、各入力バッファ 2 1 0 の出力したセルを多重化し出力ポートへ出力する出力バッファ 2 1 1 とを備え、出力バッファ 2 1 1 内の輻輳状態に応じてバックプレッシャ信号により各入力バッファ 2 1 0 のセル出力を制御する。

【 0 1 5 6 】

このセル多重化装置は、出力許可済セル集合を管理するカゴ管理部 2 1 2 を備える。

【 0 1 5 7 】

入力バッファ 2 1 0 に蓄積されているセルのうち、出力を許可されたセルの集合（出力許可済セル集合）をカゴと呼ぶ。入力バッファ 2 1 0 から出力するセルはカゴ 2 1 3 から選択する。

【 0 1 5 8 】

カゴ管理部 2 1 2 は、カゴ 2 1 3 に含まれるセル数を入力し、予め定められたアルゴリズムよりカゴ 2 1 3 へのセル転送指示を決定して全ての入力バッファ 2 1 0 へ同報する。

【 0 1 5 9 】

入力バッファ 2 1 0 は、セル転送指示により、カゴ 2 1 3 として確定されたセル以外の入力バッファ 2 1 0 に蓄積されているセルで、同一の時刻までに出力すべきセルの集合 2 1 4 をカゴ 2 1 3 に加える。

【 0 1 6 0 】

図 6 を、セルスイッチの、各出力ポートに関する部分と考えるとセルスイッチについても第 4 の実施形態を適用できる。

【 0 1 6 1 】

図 6 の出力バッファ 2 1 1 は、どんな状況においてもカゴのセルをある有限の時間内で出力ポートに出力できることを保証できれば、どんな形態でもよい。

【 0 1 6 2 】

例えば、バッファを全く持たずにカゴ 2 1 3 のセルを調停して出力ポートに出力する構成でもよい。また例えば図 3 または図 4 に示す構成でもかまわない。

【 0 1 6 3 】

なお、図 4 において、バックプレッシャ生成部 2 0 はカゴ管理部 2 1 2 に存在しても良い。

【 0 1 6 4 】

第 4 の実施形態の V C 間公平キューイングを行なうセル多重化装置によれば、高いスループットが必要な出力バッファは簡単な構成になっているため、容易に実現することが可能である。入力バッファはカゴによりセルを管理する必要があるが、入力バッファに必要なスループットは低いため、実現の容易性は高い。

【 0 1 6 5 】

図 6 のセル多重化装置の構成をさらに詳細に説明する。

【 0 1 6 6 】

カゴ管理部 2 1 2 に入力されるカゴ内セル数は、全ての入力バッファ 2 1 0 のカゴ内セル数の合計でよい。

【 0 1 6 7 】

合計の演算はカゴ管理部 2 1 2 の内部で行なっても良いし、カゴ管理部 2 1 2 の外部で行なっても良い。

【 0 1 6 8 】

セル数情報はカゴ管理部 2 1 2 において比較的小さなしきい値と大小を判定する。しきい値が固定値であるか動的に変化する値であるかはカゴ管理部 2 1 2 のアルゴリズムに依存するが、もししきい値がある固定値の場合はセル数の代わりにその比較結果をカゴ管理部 2 1 2 に入力してもよい。または、ある値以上はどのような値でもカゴ管理部 2 1 2 の処

10

20

30

40

50

理に影響を与えないのなら、それを利用してその値以上は一つの符号に符号化して入力しても良い。例えば4ビットで、セル数を(0000) = 0、(0001) = 1、(0010) = 2、(0011) = 3、...、(1101) = 13、(1110) = 14、(1111) = 15以上、という16段階に符号化するなどとしてもよい。これによりカゴ管理部212への入力情報を圧縮することができ、実装が容易になるという利点がある。

【0169】

セル数の合計を扱うことにより、カゴ管理部212の実装を簡単化することができる。カゴへのセル転送指示も全ての入力バッファ210に同じ情報を同報することから、カゴ管理部212は入力バッファ210を個々に認識する必要がない。よって入力ポート数が多くなっても本質的な実装の困難度は変化しないという利点がある。

10

【0170】

カゴ管理部212がカゴ213へのセル転送指示を決定するアルゴリズムは例えば次のようなものである。

【0171】

すなわち、カゴ内セル数 N_k より、カゴ213内のセル数がアンダフローしないかどうかを判断し、アンダフローする可能性がある場合のみアンダフローしないように転送指示を行なう。アンダフローとは、カゴ内セル数 N_k がゼロになり、本来出力できるはずの入力バッファ210のセルが有効に出力できない状態をいう。

【0172】

入力バッファ210は、カゴ213へのセル転送指示に応じて、同一の時刻までに出力すべきセルの集合214をカゴ213に加える。例えばVC毎に設定する値 N_x で重み付けされたVC間公平キューイングを行なう場合において、同一の時刻までに出力すべきセルとは、各入力バッファ210内のカゴ213に入っていないセルでVC毎に最も古い N_x 個のセルを集めたものである。

20

【0173】

転送指示は全ての入力バッファ210に同報されるため、カゴ213には複数の入力バッファ210を通してVC間で公平にセルが入る。カゴ213のセルをカゴ213の外セルよりも優先して出力バッファ211に転送することにより、複数の入力バッファ210を通してVC間で公平にセルを出力バッファ211へ転送することが可能である。

【0174】

出力バッファ211が、カゴ213から入力したセルをどんな場合でもある時間内に出力できる構成になっていれば、VC間で公平にカゴ213へ転送されたセルはその時間の遅延揺らぎを持って、出力ポートから出力される。

30

【0175】

図6は、例として、入力ポート#1、#2、#Nの入力バッファ210にそれぞれ2本、1本、3本のVCがキューイングされている状態を示している。全てのVCの重みが同じ時、VC間公平キューイングを行なえば、各入力バッファ210からの出力は2対1対3の割合になる必要がある。本実施形態のセル多重化装置は、カゴ管理部212から全ての入力バッファ210へ同報されるカゴ213へのセル転送指示によって、カゴ213内のセル数は各入力バッファで2対1対3の割合になり、結果的に出力バッファからの出力も2対1対3の割合になる。

40

【0176】

カゴ213へのセルの入力方法には2通りの方法がある。ここでは、転送を指示されたとき、各VC毎に最も古い N_x 個(VC毎に設定する値)のセルをカゴに入れる場合を例にとって説明する。

【0177】

ひとつの方法は、入力バッファ210に入力したセルがカゴ213に入れる条件を満たしていても最初はカゴ213の中へ入れない方法である。セルはカゴ管理部212からセル転送が指示された場合にのみカゴ213に転送される。これにより、カゴ内セル数はカゴ管理部212がセル転送を指示した場合以外は減少する。

50

【 0 1 7 8 】

別の方法は、入力バッファ 2 1 0 に入力したセルがカゴ 2 1 3 に入れる条件を満たしている時には最初からカゴ 2 1 3 の中へ入れてしまう方法である。つまり転送指示の時点で $N \times$ 個のセルをカゴ 2 1 3 に転送した VC は次の転送指示があるまでセルをカゴ 2 1 3 に入れることができないが、転送指示の時点でセルをカゴ 2 1 3 に全く転送しなかった VC 、または転送セル数が $N \times$ 個に満たない VC は、 $N \times$ 個になるまで入力セルを転送指示がなくてもカゴに 2 1 3 入れることができる。これにより、カゴ内セル数は転送指示と転送指示の間においても増加することがあるが、ほぼ全てのアクティブな VC が $N \times$ 個のセルを転送した段階で減少する。なお、本実施形態はどちらの方法においても有効に作用する。

【 0 1 7 9 】

第 4 の実施形態のセル多重化装置の輻輳状態の判断は、 VC 毎に、トラヒックにより変動する情報を予め定められた方法により監視すれば良い。そして、この輻輳判断を元に VC 毎に端末にトラヒック制御情報を通知すればよい。

【 0 1 8 0 】

監視すべき変動する情報は、例えば蓄積セル数、一定時間あたりの入力セル数とその目標値との関係、一定個数のセルが入力する時間とその目標値との関係である。蓄積セル数が多い VC 、一定時間あたりの入力セル数が目標値からかけ離れて多い VC 、一定個数のセルが入力する時間がその目標値からかけ離れて短い VC を輻輳している VC と判断する。

【 0 1 8 1 】

輻輳している VC のセルのヘッダにある $EFCI$ をマークすることや、通過する RM セルを書き換えることにより、端末にトラヒック制御情報を通知することができる。

【 0 1 8 2 】

2.3 第 5 の実施形態（「カゴ」という概念を用いたセル多重化装置の他の実施形態）次に、第 5 の実施形態について説明する。

【 0 1 8 3 】

図 7 に、第 5 の実施形態に係るカゴと呼ぶ概念を用いたセル多重化装置の他の構成例を示す。

【 0 1 8 4 】

図 7 に示したセル多重化装置は、入力ポート # 1 ~ # N のそれぞれに対応して設けられ、各入力ポート # 1 ~ # N から入力したセルを一時的に蓄積する複数の入力バッファ 2 2 0 と、入力バッファ 2 2 0 の出力したセルを多重化し出力ポートへ出力する出力バッファ 2 2 1 とを備え、出力バッファ 2 2 1 内の輻輳状態に応じてバックプレッシャ信号により入力バッファ 2 2 0 のセル出力を制御する。

【 0 1 8 5 】

このセル多重化装置は、出力許可済セル集合を管理するカゴ管理部 2 2 2 を備える。

【 0 1 8 6 】

入力バッファ 2 2 0 に蓄積されているセルのうち、出力を許可されたセルの集合（出力許可済セル集合）をカゴと呼ぶ。入力バッファから出力するセルはカゴ 2 2 3 から選択する。

【 0 1 8 7 】

カゴ管理部 2 2 2 は、出力バッファ内セル数を入力し、予め定められたアルゴリズムよりカゴ 2 2 3 へのセル転送指示を決定して全ての入力バッファ 2 2 0 へ同報する。

【 0 1 8 8 】

入力バッファ 2 2 0 は、セル転送指示により、カゴ 2 2 3 以外の入力バッファ 2 2 0 に蓄積されているセルで同一の時刻までに出力すべきセルの集合 2 2 4 をカゴ 2 2 3 に加える。

【 0 1 8 9 】

図 7 を、セルスイッチの、各出力ポートに関する部分と考えるとセルスイッチについても本実施形態を適用できる。

【 0 1 9 0 】

10

20

30

40

50

図7の出力バッファ221は、どんな状況においてもカゴ223のセルをある有限の時間内で出力ポートに出力できることを保証できれば、どんな形態でもよい。

【0191】

例えば、図4に示す構成でもかまわない。この場合、バックプレッシャ生成部20はカゴ管理部222に存在しても良い。出力バッファ内セル数はバックプレッシャ生成部20にて加算し、カゴ管理部222へ出力する。

【0192】

第5の実施形態に係るV C間公平キューイングを行なうセル多重化装置によれば、高いスループットが必要な出力バッファは簡単な構成になっているため、容易に実現することが可能である。入力バッファはカゴによりセルを管理する必要があるが、入力バッファに必要なスループットは低いため、実現の容易性は高い。

10

【0193】

図7のセル多重化装置の構成をさらに詳細に説明する。

【0194】

セル数情報はカゴ管理部222において比較的小さなしきい値と大小を判定する。しきい値が固定値であるか動的に変化する値であるかはカゴ管理部222のアルゴリズムに依存するが、もししきい値がある固定値の場合はセル数の代わりにその比較結果をカゴ管理部222に入力してもよい。または、ある値以上はどのような値でもカゴ管理部222の処理に影響を与えないのなら、それを利用してその値以上は一つの符号に符号化して入力してもよい。例えば4ビットで、セル数を(0000) = 0、(0001) = 1、(0010) = 2、(0011) = 3、...、(1101) = 13、(1110) = 14、(1111) = 15以上、という16段階に符号化するなどとしてもよい。これによりカゴ管理部222への入力情報を圧縮することができ、実装が容易になるという利点がある。

20

【0195】

カゴ管理部222がカゴ223へのセル転送指示を決定するアルゴリズムは例えば次のようなものである。

【0196】

すなわち、出力バッファ内セル数より、出力バッファ221内のセル数がアンダフローしないかどうかを判断し、アンダフローする可能性がある場合のみアンダフローしないように転送指示を行なう。アンダフローとは、出力バッファ内セル数がゼロになり、本来出力できるはずの入力バッファ220のセルが有効に出力できない状態をいう。

30

【0197】

入力バッファ220は、カゴ223へのセル転送指示に応じて、同一の時刻までに出力すべきセルの集合224をカゴ223に加える。例えばV C毎に設定する値N_xで重み付けされたV C間公平キューイングを行なう場合において、同一の時刻までに出力すべきセルの集合とは、各入力バッファ内のカゴ223に入っていないセルでV C毎に最も古いN_x個のセルを集めたものである。

【0198】

転送指示は全ての入力バッファに同報されるため、カゴ223には複数の入力バッファを通してV C間で公平にセルが入る。カゴ223のセルをカゴの外のセルよりも優先して出力バッファ221に転送することにより、複数の入力バッファ220を通してV C間で公平にセルを出力バッファ221へ転送することが可能である。出力バッファ221が、カゴ223から入力したセルをどんな場合でもある時間内に出力できる構成になっていれば、V C間で公平にカゴ223へ転送されたセルはその時間の遅延揺らぎを持って、出力ポートから出力される。

40

【0199】

図7は、例として、入力ポート#1、#2、#Nの入力バッファ220にそれぞれ2本、1本、3本のV Cがキューイングされている状態を示している。全てのV Cの重みが同じ時、V C間公平キューイングを行なえば、各入力バッファ220からの出力は2対1対3の割合になる必要がある。

50

【 0 2 0 0 】

図7のセル多重化装置は、カゴ管理部222から全ての入力バッファ220へ同報されるカゴ223へのセル転送指示によって、カゴ内のセル数は各入力バッファで2対1対3の割合になり、結果的に出力バッファ221からの出力も2対1対3の割合になる。

【 0 2 0 1 】

カゴ223へのセルの入力方法には2通りの方法がある。ここでは、転送を指示されたとき、各VC毎に最も古いNx個（VC毎に設定する値）のセルをカゴに入れる場合を例にとって説明する。

【 0 2 0 2 】

ひとつの方法は、入力バッファ220に入力したセルがカゴ223に入れる条件を満たしていても最初はカゴの中へ入れない方法である。セルはカゴ管理部222からセル転送が指示された場合にのみカゴ223に転送される。これにより、カゴ内セル数はカゴ管理部222がセル転送を指示した場合以外は減少する。

10

【 0 2 0 3 】

別の方法は、入力バッファ220に入力したセルがカゴ223に入れる条件を満たしている時には最初からカゴ223の中へ入れてしまう方法がある。つまり転送指示の時点でNx個のセルをカゴ223に転送したVCは次の転送指示があるまでセルをカゴ223に入れることができないが、転送指示の時点でセルをカゴ223に全く転送しなかったVC、または転送セル数がNx個に満たないVCは、Nx個になるまで入力セルを転送指示がなくてもカゴ223に入れることができる。これにより、カゴ内セル数は転送指示と転送指示の間においても増加することがあるが、ほぼ全てのアクティブなVCがNx個のセルを転送した段階で減少する。なお、本実施形態はどちらの方法においても有効に作用する。

20

【 0 2 0 4 】

第5の実施形態のセル多重化装置の輻輳状態の判断は、VC毎に、トラヒックにより変動する情報を予め定められた方法により監視すれば良い。そして、セル多重化装置は、この輻輳判断を元にVC毎に端末にトラヒック制御情報を通知すればよい。

【 0 2 0 5 】

監視すべき変動する情報は、例えば蓄積セル数、一定時間あたりの入力セル数とその目標値との関係、一定個数のセルが入力する時間とその目標値との関係である。蓄積セル数が多いVC、一定時間あたりの入力セル数が目標値からかけ離れて多いVC、一定個数のセルが入力する時間がその目標値からかけ離れて短いVCを輻輳しているVCと判断する。

30

【 0 2 0 6 】

輻輳しているVCのセルのヘッダにあるEFCIをマークすることや、通過するRMセルを書き換えることにより、端末にトラヒック制御情報を通知することができる。

【 0 2 0 7 】

2.4 第4、第5の実施形態に係るセル多重化装置の利点

以上説明したように、上記第4～第5の実施形態に係る出力許可済セル集合（カゴ）の概念を用いたセル多重化装置によれば、出力許可済セル集合に入れるセルを制御することにより各入力バッファの出力スループットを調整するため、VC間公平キューイングを行なうことができ、ABRサービスにおける耐故障性や、UBRサービスにおける公平性を実現することが可能である。

40

【 0 2 0 8 】

2.5 第6の実施形態（セルグループFIFOを用いたセルバッファ装置）次に、本発明に係るVC間公平キューイングを行なうセルバッファ装置の実施形態について説明する。

【 0 2 0 9 】

図8に、第6の実施形態に係るセルグループFIFOを用いたセルバッファ装置の構成例を示す。

【 0 2 1 0 】

図8は、入力リンクより入力したセルのコネクション識別情報をバッファポインタ管理部

50

230に通知し、バッファポインタ管理部230よりセルの書き込み位置を示す書き込みポインタを得てセルを一時的にセルバッファ231に蓄積し、バッファポインタ管理部230からの、読みだしセルを示す読みだしポインタに基づいてセルをセルバッファ231から読みだし出力リンクへ出力するセルバッファ装置である。

【0211】

バッファポインタ管理部230は、蓄積しているセルのセルバッファ231上の位置を示すバッファポインタを管理する。

【0212】

バッファポインタ管理部230は、バッファポインタの集合である複数のセルグループをFIFO管理するセルグループFIFO232aと、出力待ちセルグループFIFO232bと、セルグループ選択部233と、空きバッファポインタ管理部234からなる。

10

【0213】

バッファポインタ管理部230は、セル入力時には、空きバッファポインタ管理部234より空きのバッファポインタを得て前記書き込みポインタとするとともに、セルグループ選択部233がコネクション識別情報より書き込みポインタをセルグループFIFO232aの先頭のセルグループより順にVC毎に予め定められた重みに従って決定される数だけ入るようにセルグループ指示を行ない、セルグループFIFO232aはセルグループ指示に従って書き込みポインタを指示されたセルグループに入力する。

【0214】

セル出力時には、セルグループFIFO232aの先頭よりセルグループを出力し、さらにそのセルグループよりバッファポインタを出力して前記読みだしポインタとするとともに、その読みだしポインタを空きバッファポインタ管理部234に戻す。

20

【0215】

図8では、セルグループFIFO232aから出力され、出力を待っているセルグループが複数になる可能性がある場合の構成を示している。これらのセルグループFIFO232aから出力されたセルグループは出力待ちセルグループFIFO232bに入力される。

【0216】

出力待ちセルグループFIFO232bのセルグループは、セルバッファ231の外部の何らかの管理部により出力を許可されたセルグループである。例えば、図6のカゴ231に相当する。

30

【0217】

読みだしポインタは、出力待ちセルグループFIFO232bの先頭のセルグループから出力されたバッファポインタである。

【0218】

例えば全てのVCの重みを同じとする。セルグループFIFO232aの先頭のセルグループのセルは、各VCのキューを考えるとキューの先頭(ただし出力待ちセルグループFIFO232b内のセルは除く)のセルである。セルグループFIFO232aの先頭から2番目のセルグループのセルは、各VCのキューの2番目のセルである。3番目以降のセルグループのセルも同様である。これらのセルを出力する場合はセルグループFIFO232aの先頭のセルグループから順に出力するためVC間で公平に出力することになる。

40

【0219】

今まで全く到着していなかったVCのセルが新たにこのバッファ装置に到着した場合は、セルグループFIFO232aの先頭のセルグループに入力され、他のVCのキューの2番目以降のセルよりも優先して出力される。

【0220】

この様に、第6の実施形態に係るセルバッファ装置によれば、VC間で公平にキューイングを行なうことが可能でありながらセルの入力、出力時に図50に示した従来例のように検索動作が不要であるという利点がある。

50

【 0 2 2 1 】

第 6 の実施形態に係るセルバッファ装置のセルグループ F I F O 2 3 2 a の実現方法は、例えばポインタチェーンによる方式やリングバッファによる方式が考えられる。

【 0 2 2 2 】

2 . 6 第 7 の実施形態（セルグループ F I F O を用いたセルバッファ装置の他の実施形態）

次に、第 7 の実施形態について説明する。

【 0 2 2 3 】

図 9 に、第 7 の実施形態に係るセルグループ F I F O を用いたセルバッファ装置の他の構成例を示す。

10

【 0 2 2 4 】

図 9 は、入力リンクより入力したセルの接続識別情報をバッファポインタ管理部 2 4 0 に通知し、バッファポインタ管理部 2 4 0 よりセルの書き込み位置を示す書き込みポインタを得てセルを一時的にセルバッファ 2 4 1 に蓄積し、バッファポインタ管理部 2 4 0 からの、読みだしセルを示す読みだしポインタに基づいてセルをセルバッファ 2 4 1 から読みだし出力リンクへ出力するセルバッファ装置である。

【 0 2 2 5 】

バッファポインタ管理部 2 4 0 は、蓄積しているセルのセルバッファ 2 4 1 上の位置を示すバッファポインタを管理する。

【 0 2 2 6 】

バッファポインタ管理部 2 4 0 は、バッファポインタの集合であるセルグループを F I F O 管理するセルグループ F I F O 2 4 2 と、セルグループ選択部 2 4 3 と、空きバッファポインタ管理部 2 4 4 からなる。

20

【 0 2 2 7 】

バッファポインタ管理部 2 4 0 は、セル入力時には、空きバッファポインタ管理部 2 4 4 より空きのバッファポインタを得て前記書き込みポインタとするとともに、セルグループ選択部 2 4 3 が接続識別情報より書き込みポインタをセルグループ F I F O 2 4 2 の先頭のセルグループより順に V C 毎に予め定められた重みに従って決定される数だけ入るようにセルグループ指示を行ない、セルグループ F I F O 2 4 2 はセルグループ指示に従って書き込みポインタを指示されたセルグループに入力する。

30

【 0 2 2 8 】

セル出力時には、セルグループ F I F O 2 4 2 の先頭のセルグループよりバッファポインタを出力して前記読みだしポインタとするとともに、その読みだしポインタを空きバッファポインタ管理部 2 4 4 に戻す。

【 0 2 2 9 】

図 9 では、セルグループ F I F O 2 4 2 内で出力を待っているセルグループを出力待ちセルグループ F I F O 2 4 2 b として記してある。

【 0 2 3 0 】

出力待ちセルグループ F I F O 2 4 2 b のセルグループは、セルバッファ 2 4 1 の外部の何らかの管理部により出力を許可されたセルグループである。例えば、図 6 のカゴに相当

40

【 0 2 3 1 】

読みだしポインタは、出力待ちセルグループ F I F O 2 4 2 b の先頭のセルグループから出力されたバッファポインタである。

【 0 2 3 2 】

例えば全ての V C の重みを同じとする。セルグループ F I F O 2 4 2 の先頭のセルグループ（すなわち、出力待ちセルグループ F I F O 2 4 2 b 内のセルグループ）のセルは、各 V C のキューを考えるとキューの先頭のセルである。セルグループ F I F O 2 4 2 の先頭から 2 番目のセルグループのセルは、各 V C のキューの 2 番目のセルである。3 番目以降のセルグループのセルも同様である。これらのセルを出力する場合はセルグループ F I F

50

0242の先頭のセルグループから順に出力するためVC間で公平に出力することになる。

【0233】

今まで全く到着していなかったVCのセルが新たにこのバッファ装置に到着した場合は、セルグループFIFO0242の先頭のセルグループに入力され、他のVCのキューの2番目以降のセルよりも優先して出力される。

【0234】

この様に、第7の実施形態に係るセルバッファ装置によれば、VC間で公平にキューイングを行なうことが可能でありながらセルの入力、出力時に図51に示した従来例のように検索動作が不要であるという利点がある。

10

【0235】

第7の実施形態に係るセルバッファ装置のセルグループFIFO0242の実現方法は、例えばポインタチェーンによる方式やリングバッファによる方式が考えられる。

【0236】

2.7 第6の実施形態に係るセルバッファ装置で用いられるデータ構造

次に、第6の実施形態に係るセルバッファ装置(図8参照)で用いられるデータの構造について説明する。

【0237】

図10、図11に、図8で説明したバッファ装置をポインタチェーン方式で実現した場合のデータ構造の一例を示す。

20

【0238】

大きく分けて、図10に示すように、VCテーブル250、セルグループFIFO232a、出力待ちセルグループFIFO232b、空きバッファポインタチェーン251、また、図11に示すように、空きセルグループチェーン252で構成される。

【0239】

セルグループFIFO232aには、それを管理するセルグループFIFO管理データ253、出力待ちセルグループFIFO232bには、それを管理する出力待ちセルグループFIFO管理データ254がある。

【0240】

セルグループと空きバッファポインタチェーン251は、バッファポインタのチェーンであり、セルグループFIFO232aと空きセルグループチェーン252はセルグループ管理データ255のチェーンである。

30

【0241】

出力待ちセルグループFIFO232bは、出力待ちセルグループFIFO管理データ254にセルグループを指すポインタのリスト(Ptr1、Ptr2、Ptr3、...)がある。

【0242】

セルグループFIFO232a、出力待ちセルグループFIFO232bはリングバッファ方式でもかまわない。

【0243】

セルが入力した場合、空きバッファポインタチェーン251から取り出した書き込みポインタをどのセルグループに入力するかを決定しなければならない。そのため、まずVCテーブル250の該当する領域を読み出す。VCテーブル250のNxはそのVCの重み、Ncは作業変数、QlenはそのVCの蓄積セル数、PtrはそのVCのセルが蓄積されているセルグループの末尾へのポインタである。

40

【0244】

はじめに、Qlenをチェックして、そのVCのセルが現在バッファ装置に蓄積されているかどうかを調べる。もし、Qlenがゼロの場合は、セルグループFIFO232aの先頭のセルグループに入れる。Qlenが1以上の場合でも、Ptrが出力待ちセルグループFIFO232b内のセルグループを指している場合は、セルグループFIFO232aの先頭のセルグループに入れる。

50

【 0 2 4 5 】

その他の場合は N_c に 1 . 0 を加え N_x と比較し、 N_x の方が大きければ P_{tr} で指しているセルグループに入れ、そうでない場合は、 P_{tr} で指しているセルグループの次のセルグループに入れる。

【 0 2 4 6 】

ここで、セルが入力した時の N_c の更新手順について説明する。前述のように、その V_C に関してセルグループ $F I F O 2 3 2 a$ の先頭のセルグループに初めてバッファポインタを入力する場合は、 $N_c := 1 . 0$ とする。その他の場合は、 N_c を $N_c := N_c + 1 . 0$ と更新する。その結果、もし $N_x \geq N_c$ ならば、そのセルグループにバッファポインタを入れる。逆に、もし $N_x < N_c$ ならば、セルグループ $F I F O 2 3 2 a$ の次のセルグループにバッファポインタを入れると同時に、 $N_c := N_c - N_x$ と書き変える（つまり、 $0 < N_c \leq N_x$ である）。

10

【 0 2 4 7 】

もし、セルグループ $F I F O 2 3 2 a$ の最後のセルグループの次のセルグループにバッファポインタを入れる必要がある場合は、空きセルグループチェーン 2 5 2 の先頭よりセルグループ管理データ 2 5 5 をとりだし、セルグループ $F I F O 2 3 2 a$ の末尾に入れ、そのセルグループにバッファポインタを入れれば良い。出力待ちセルグループ $F I F O$ 管理データ 2 5 4 は、セルグループを指すポインタ P_{tr1} 、 P_{tr2} 、 P_{tr3} 、... がシフトレジスタになっている。 P_{tr1} の示すセルグループが $F I F O$ の先頭であり、このセルグループから出力したバッファポインタが読みだしポインタとなる。読みだしポインタは空きバッファポインタチェーン 2 5 1 の末尾に入力する。 P_{tr1} の示すセルグループのバッファポインタが空になったら、そのセルグループは空きセルグループチェーン 2 5 2 の末尾に入力する。そして $n \geq 2$ の全ての n についてポインタ $P_{tr}(n)$ を $P_{tr}(n - 1)$ にシフトし、新しい P_{tr1} の指すセルグループからバッファポインタを出力する。

20

【 0 2 4 8 】

前述したように、セルの入力時に V_C テーブル 2 5 0 の P_{tr} が出力待ちのセルグループを指しているかどうか判定しなければならない。したがって出力待ちセルグループ $F I F O$ 管理データ 2 5 4 は、検索が容易に可能な構成であることが必要である。

【 0 2 4 9 】

セルグループ $F I F O 2 3 2 a$ から出力待ちセルグループ $F I F O 2 3 2 b$ に新たにセルグループを転送する場合には、出力待ちセルグループ $F I F O$ 管理データ 2 5 4 のセルグループ数をインクリメントして (m とする)、 $P_{tr}(m)$ がセルグループ $F I F O 2 3 2 a$ から出力したセルグループを指すようにする。

30

【 0 2 5 0 】

なお、第 6 の実施形態に係るセルバッファ装置は、図 6 や図 7 に示すようなカゴの概念を用いた入力バッファに容易に適用することができる。その場合、出力待ちセルグループ $F I F O 2 3 2 b$ がカゴに相当する。また、図 5 の出力バッファにも適用することができる。

【 0 2 5 1 】

2 . 8 第 7 の実施形態に係るセルバッファ装置で用いれるデータ構造
次に、第 7 の実施形態に係るセルバッファ装置（図 9 参照）で用いられるデータの構造について説明する。

40

【 0 2 5 2 】

図 1 2、図 1 3 に、図 9 で説明したバッファ装置をポインタチェーン方式で実現した場合のデータ構造の一例を示す。

【 0 2 5 3 】

大きく分けて、図 1 2 に示すように、 V_C テーブル 2 6 0、セルグループ $F I F O 2 4 2$ 、空きバッファポインタチェーン 2 6 1、また、図 1 3 に示すように、空きセルグループチェーン 2 6 2 で構成される。

【 0 2 5 4 】

50

セルグループ F I F O 2 4 2 の内部には、出力待ちセルグループ F I F O 2 4 2 b が存在し、セルグループ F I F O 2 4 2 を管理するセルグループ F I F O 管理データ 2 6 3、出力待ちセルグループ F I F O 2 4 2 b を管理する出力待ちセルグループ F I F O 管理データ 2 6 4 がある。

【 0 2 5 5 】

セルグループと空きバッファポインタチェーン 2 6 1 は、バッファポインタのチェーンであり、セルグループ F I F O 2 4 2 と空きセルグループチェーン 2 6 2 はセルグループ管理データ 2 6 5 のチェーンである。

【 0 2 5 6 】

出力待ちセルグループ F I F O 2 4 2 b は、セルグループ F I F O 2 4 2 の先頭のいくつかのセルグループである。つまり図 1 2 において、セルグループ F I F O 管理データ 2 6 3 の先頭ポインタが指すセルグループから出力待ちセルグループ F I F O 管理データ 2 6 4 の末尾ポインタが指すセルグループまでが出力待ちセルグループ F I F O 2 4 2 である。

10

【 0 2 5 7 】

セルグループ F I F O 2 4 2、出力待ちセルグループ F I F O 2 4 2 b はリングバッファ方式でもかまわない。

【 0 2 5 8 】

セルが入力した場合、空きバッファポインタチェーン 2 6 1 から取り出した書き込みポインタをどのセルグループに入力するかを決定しなければならない。そのため、まず V C テーブル 2 6 0 の該当する領域を読み出す。V C テーブル 2 6 0 の N_x はその V C の重み、 N_c は作業変数、 Q_{len} はその V C の蓄積セル数、 P_{tr} はその V C のセルが蓄積されているセルグループの末尾へのポインタである。

20

【 0 2 5 9 】

はじめに、 Q_{len} をチェックして、その V C のセルが現在バッファ装置に蓄積されているかどうかを調べる。もし、 Q_{len} がゼロの場合は、セルグループ F I F O 2 4 2 の先頭のセルグループに入れる。

【 0 2 6 0 】

その他の場合は N_c に 1 . 0 を加え N_x と比較し、 N_x の方が大きければ P_{tr} で指しているセルグループに入れ、そうでない場合は、 P_{tr} で指しているセルグループの次のセルグループに入れる。

30

【 0 2 6 1 】

ここで、セルが入力した時の N_c の更新手順について説明する。前述のように、その V C に関してセルグループ F I F O 2 4 2 の先頭のセルグループに初めてバッファポインタを入力する場合は、 $N_c := 1 . 0$ とする。その他の場合は、 N_c を $N_c := N_c + 1 . 0$ と更新する。その結果、もし $N_x \geq N_c$ ならば、そのセルグループにバッファポインタを入れる。逆に、もし $N_x < N_c$ ならば、セルグループ F I F O 2 4 2 の次のセルグループにバッファポインタを入れると同時に、 $N_c := N_c - N_x$ と書き変える（つまり、 $0 < N_c \leq N_x$ である）。

【 0 2 6 2 】

セルが入力した時、セルグループがセルグループ F I F O 2 4 2 になかった場合、もし、セルグループ F I F O 2 4 2 の最後のセルグループの次のセルグループにバッファポインタを入れる必要があれば、空きセルグループチェーン 2 6 2 の先頭よりセルグループ管理データ 2 6 5 をとりだし、セルグループ F I F O 2 4 2 の末尾に入れ、そのセルグループにバッファポインタを入れれば良い。

40

【 0 2 6 3 】

出力待ちセルグループ F I F O 管理データ 2 6 4 は、出力待ちセルグループ F I F O 2 4 2 b の末尾を示すポインタを持っている。先頭はセルグループ F I F O 管理データ 2 6 3 の先頭のセルグループと同じである。

【 0 2 6 4 】

50

出力待ちセルグループ F I F O 2 4 2 b の先頭のセルグループ (セルグループ F I F O の先頭と同じセルグループ) からバッファポインタを出力し、このバッファポインタが読みだしポインタとなる。読みだしポインタは空きバッファポインタチェーン 2 6 1 の末尾に

【 0 2 6 5 】

セルグループ F I F O 2 4 2 の先頭のセルグループのバッファポインタが空になったら、セルグループ F I F O 2 4 2 からセルグループ管理データ 2 6 5 を出力し、空きセルグループチェーン 2 6 2 の末尾に

【 0 2 6 6 】

出力待ちセルグループ F I F O 2 4 2 b に新たにセルグループを転送する場合には、出力待ちセルグループ F I F O 管理データ 2 6 4 の末尾を示すポインタをセルグループ F I F O 2 4 2 上の次のセルグループを指すように変更する。 10

【 0 2 6 7 】

なお、第 7 の実施形態に係るセルバッファ装置は、図 6 や図 7 に示すようなカゴの概念を用いた入力バッファに容易に適用することができる。その場合、出力待ちセルグループ F I F O 2 4 2 b がカゴに相当する。また、図 5 の出力バッファにも適用することができる。

【 0 2 6 8 】

2 . 9 第 6 、第 7 の実施形態に係るセルバッファ装置の利点
以上説明したように、上記第 6 ~ 第 7 の実施形態に係る V C 間公平キューイングを行なうセルバッファ装置によれば、V C テーブルを検索する処理が必要なく、したがって設定可能な V C 数の上限に依存することなく容易に実現することができる。 20

【 0 2 6 9 】

2 . 1 0 フロー間重み付き公平キューイング
以上の V C 間公平キューイングを行なうセルバッファ装置は、固定長のパケットであるセルを扱う場合を例にあげて説明した。異なる長さのパケットを同時に扱うことが可能なパケットバッファ装置にも以下に説明するように入力したパケットをフロー毎に定めた重みに基づいて公平に出力することができる。

【 0 2 7 0 】

ここでのフローとは、以下の例に示すようにある基準によって識別されるパケットの集合である。 30

【 0 2 7 1 】

・ I P 網のいわゆるフロー (送信元アドレス、送信元レポート、宛先アドレス、宛先ポートの組で識別)、もしくは A T M 網の V C C または V P C。

【 0 2 7 2 】

・ I P 網のサービスクラス (g u a r a n t e e d 、 c o n t r o l l e d - l o a d 、 b e s t - e f f o r t など)、もしくは A T M 網のサービスカテゴリ (C B R 、 V B R 、 A B R など) の違い。

【 0 2 7 3 】

・ プロトコルの違い (T C P / I P 、 D E C n e t 、 S N A 、 A p p l e T a l k など) 40
・ アプリケーションの違い (f t p 、 t e l n e t など)
・ 同じパケットバッファ装置を共用する異なった組織 (企業別など)
これらを本発明のパケットバッファ装置で異なったフローとして扱うことによりトラヒックを互いに分離することができる。

【 0 2 7 4 】

はじめに、図 1 4 を用いて、一般的なフロー間重み付き公平キューイングを説明する。このパケットバッファ装置の目的は、装置に蓄積しているパケットをそのフローの重みに従って公平に出力することである。その動作は図 1 4 のようなフロー毎の F I F O を考えるとわかりやすい。図には、フロー 1、フロー 2、フロー 3 の 3 つのフローのキュー (F I F O) があり、それぞれのキューにはいろいろな長さのパケットが蓄積されている。図中 50

の packets に添えられている数字は、バイト数で表したその packets の長さである。

【0275】

この packets バッファ装置に packets が入力した場合の処理は、入力した packets が属するフローを判定しフロー毎のキューに蓄積することである。

【0276】

一方 packets の出力については、フロー間で公平に行なうことが必要であり、どのキューの packets をどの順で出力するかが重要である。ひとつの方法を図14の状況に合わせて説明する。説明を簡単にするためここでは各フローの重みを同じとしている。また公平の単位をここでは仮に500バイトと定めた。

【0277】

まずフロー1のキューの先頭から300バイトの長さの packets と200バイトの長さの packets の合わせて500バイト分の packets を出力する。次にフロー2のキューの先頭から250バイトの長さの packets を2個続けて出力する。そしてフロー3のキューの先頭から500バイトの長さの packets をひとつ出力する。このように出力すると、この時点では各フロー1バイトの長さの packets をひとつ出力する。このような出力にすると、この時点では各フローから等しく500バイトずつの packets を出力しているので公平に出力したといってよい。続いて再びフロー1から順に各フローから500バイトずつ packets を出力することを繰り返せばよい。この様にすれば packets はフロー間で公平に出力される。

【0278】

以上の方法では500バイトを単位として公平になっているが、それより小さな単位では必ずしも公平になっていない。しかし厳密な公平性の実現は packets バッファ装置に非常に大きな処理能力を要求することが知られており、現実問題として上述の程度の公平性が十分と判断される場合が多い。

【0279】

さて、上述の説明ではフローの数が3つであったが、 packets バッファ装置が莫大な数のフローを扱う場合、キューに packets を蓄積しているフローをどの様に探索し、公平に packets を出力するのが課題になる。本発明は、 packets グループという概念を用いてこの課題を解決する。

【0280】

本発明は、前述の図14のZ、A、B、C、Dの様に全てのフローのキューを一定の長さに区切る。このとき、各区切りを単位として packets を出力することで公平な packets 出力を実現できることに着目する。本発明はこれらの区切り毎に packets をグループ化して管理する。この各グループを packets グループと呼ぶ。ひとつの packets グループには複数のフローの packets が混在する。 packets の出力はまず出力側に最も近い packets グループZの packets を出力し、次に packets グループAの packets を出力する。さらにB、C、Dというように順に packets グループの packets を出力することにより、本発明の packets バッファ装置はフロー間で公平に packets を出力できる。

【0281】

ところで図14において区切りBのフロー3を見ると、400バイトの packets の次に350バイトの packets が蓄積されている。図では同じ区切りに属する packets は合計500バイトを上限としているため、 packets を packets バッファ装置内部で分割しないとすると400バイトの packets と350バイトの packets の両方を区切り、Bに入れることはできない。そうかといって、後から packets バッファ装置に到着した350バイトの packets を区切りBに入れなければ区切りBには100バイトのすき間が発生してしまう。この様な場合の扱いは後で述べる。

【0282】

以上本発明の packets グループという概念について説明した。

【0283】

続いて本発明の packets バッファ装置の一例を図15を用いてより具体的に説明する。

10

20

30

40

50

【0284】

本発明の構成要素として、パケットグループ、パケットグループFIFO、フローテーブル、出力待ちパケットグループがある。

【0285】

パケットグループは前述した様に蓄積しているパケットをグループ化したものである。図15ではA、B、C、Dがパケットグループである。各パケットグループに属するパケットには順序がついており、パケットグループはいわばFIFO構造になっている。

【0286】

本発明は図15のように、パケットのFIFOであるパケットグループをさらにFIFOで管理する。このパケットグループのFIFOをパケットグループFIFOと呼ぶ。

10

【0287】

フローテーブルはフロー毎の情報を記憶する。フローテーブルの N_x はそのフローの重みの設定値、 N_c は作業変数、蓄積量はそのフローの蓄積パケットの合計量、Ptrはそのフローの最後に到着したパケットが属するパケットグループへのポインタである。蓄積量はパケットバッファ装置がそのフローのパケットを蓄積しているかどうかを示す目的を持ち、例えばパケット入力時に入力パケット長を加え、パケット出力時に出力パケット長を引く変数とすることによりフロー毎に蓄積パケットの合計バイト数を保持する。

【0288】

出力待ちパケットグループは、パケットグループFIFOから出力されたパケットグループである(図ではパケットグループZ)。パケットバッファ装置からパケットを出力する場合、出力待ちパケットグループに属するパケットを出力する。

20

【0289】

パケットが入力した場合、そのパケットをどのパケットグループに入れるかを決定しなければならない。

【0290】

まずはじめに、そのフローのパケットが現在バッファ装置に蓄積されているかどうかを調べる。フロー識別情報よりフローテーブルのそのフローの蓄積量を参照し、それがゼロかどうかで判断できる。現在蓄積されていなければ、パケットグループFIFOの先頭のパケットグループ(図15のパケットグループA)に入力パケットを入れる。

【0291】

蓄積量がゼロより大きな場合でも、Ptrが出力待ちパケットグループ(図15のパケットグループZ)を指している場合は、やはりパケットグループFIFOの先頭のパケットグループAに入力パケットを入れる。

30

【0292】

その他の場合は、入力パケット長と、 N_x 、 N_c 、Ptrによって決定する。

【0293】

各フローの N_x は、そのフローについて、ひとつのパケットグループが管理する平均的なパケットの量を示している。 N_x はそのフローの最大パケット長よりも大きな値に設定する。

【0294】

N_c は、 $N_x - N_c$ を計算することにより、Ptrで示しているパケットグループに、そのフローが今後入れることのできる残りの量を示す。もし $N_x - N_c$ (入力パケット長)ならば、Ptrで指しているパケットグループに入力パケットを入れ、逆に $N_x - N_c <$ (入力パケット長)ならば、パケットグループFIFO上で、Ptrが目指すパケットグループの次のパケットグループに入力パケットを入れる。

40

【0295】

パケットグループが決定すると、次に N_c とPtrの値を更新する。パケットグループFIFOの先頭のパケットグループに、そのフローに関して1つめのパケットを入れた場合は、 $N_c :=$ (入力パケット長)とする。その他の場合は、まず、 $N_c := N_c +$ (入力パケット長)とし、その結果もし $N_x < N_c$ ならさらに $N_c := N_c - N_x$ と書き変える(

50

つまり、 $0 < N_c < N_x$ である)。Ptrはその入力 packets を入れた packet グループへのポインタを書き込む。

【0296】

本発明では、ひとつの packet をふたつの packet グループが管理することはない。そのため入力 packet の長さが $N_x - N_c$ を越えている場合には、入力 packet を Ptr の指す packet グループに入れず次の packet グループに入れる。しかし N_c の計算上ではひとつの packet をふたつの packet グループに分割して入れる。つまり、 $N_c := N_c + (\text{入力 packet 長})$ とし、 N_x を越えた部分があれば、その $N_c - N_x$ を次の packet グループの N_c とする。このアルゴリズムにより平均的に、ひとつの packet グループあたり N_x の packet が属することになる。

10

【0297】

図15を用いて packet 入力動作を具体的に説明する。図15はフロー1で長さ250バイトの packet が packet バッファ装置に到着した場面である。このときフロー1の Ptr は packet グループ B を指している。フローテーブルの更新は、上述の手续にしたがって N_c を $N_c := N_c + (\text{入力 packet 長}) = 200 + 250 = 450$ とする。この値は N_x 以下のなので、この入力 packet は Ptr が指している packet グループ B に入れる。

【0298】

もし、仮にフロー1の入力 packet 長が250ではなく350であった場合は $N_c := 200 + 350 = 550$ となり N_x を越える。したがってさらに、 $N_c := N_c - N_x = 550 - 500 = 50$ と更新し、この入力 packet を packet グループ B の次の packet グループ C に入れる。Ptr は packet グループ C に変更する。

20

【0299】

packet グループ FIFO は先入れ先だしのキューであり、キューイングしている packet グループ数を必要なだけ増やすことができる。例えば、図15においてフロー2で packet 長が400の packet が到着したと仮定すると、空きの packet グループ (packet グループ E) が packet グループ D の次に追加され、入力 packet はその packet グループ E に入れることになる。

【0300】

以上のように packet 入力時の処理は、フローテーブルの該当箇所を変更することと、packet グループに packet を入れることと、必要があれば packet グループ FIFO に空きの packet グループを追加することとでよい。フローテーブルを複数のフローに渡って検索することなく入力時の処理を行なうことが可能であり、フロー数が大きくなっても処理の複雑さは変化しない。

30

【0301】

一方、packet の出力時に必要な処理はさらに簡単である。つまり出力待ち packet グループから packet を出力すればよい。必要に応じて packet グループ FIFO から packet グループを出力することをその前に行なう。これらの処理もフローテーブルを複数のフローに渡って検索することなく可能であり、フロー数が大きくなっても処理の複雑さは変化しない。

【0302】

前述したように各フローの重み N_x は、そのフローに関してひとつの packet グループあたりの平均的な packet の量を示している。そして本発明の packet バッファ装置は packet グループを単位として packet を出力するため、各フローの出力スリットは N_x に比例する。

40

【0303】

この packet バッファ装置によれば、フローテーブルを検索する処理が必要なく、したがって設定可能なフロー数の上限に依存することなく容易に可変長 packet を扱い重み付き公平な packet バッファ装置を実現することができる。

【0304】

本発明の各 packet グループは packet を到着順に管理する。つまり同じ packet グループ

50

ブに属するパケットの出力順位は全てのフローを通して到着順になるということである。この性質はパケットの遅延揺らぎを低減する効果がある。

【0305】

また、本発明のパケットバッファ装置は、従来の技術で説明したDRRというアルゴリズムと比較すると、VC数が増えても高速にVC間公平キューイングの為の出力セルの選択が実現できるうえ、重みの値が大きくなっても出力トラヒックのバースト性が増加することがない。本発明でバースト性が増加しないのは、入力したパケットを複数の集合（パケットグループ）に振り分け管理する際に、ひとつの集合内ではパケットをフロー（VC）とは無関係に管理し、パケットをその集合（パケットグループ）から出力する際はフロー（VC）とは無関係に出力するからである。

10

【0306】

3. バックプレッシャを考慮したキュー監視（第8の実施形態）

次に、本発明の第8の実施形態に係るバックプレッシャを考慮したキュー監視を行うセルバッファ装置について説明する。

【0307】

図16に、例えば、前述の第1、第2、第4、第5の実施形態で説明したような入力バッファにおけるキュー長の変化の一例を示す。

【0308】

第8の実施形態に係るセルバッファ装置は、バックプレッシャ信号により出力が出力可能と出力禁止の2状態に制御されているセルバッファ装置であって、バックプレッシャ信号が出力可能から出力禁止に変化した時から次に出力可能から出力禁止に変化するまでの間で、最初のセル入力時の蓄積セル数のみを監視することを特徴とする。

20

【0309】

監視の目的は、輻輳の判定や、ひとつのバッファ領域を複数のキューが共用する場合においてバッファ領域の独占を防ぐための入力セル廃棄の判定などである。

【0310】

図17にバックプレッシャ考慮のキュー監視部を持つセルバッファ装置の構成例を示す。図17において、セルバッファ装置は、複数のキュー Q_i ($i = 1 \sim N$)、キューテーブル記憶部280、バックプレッシャ識別子(BPID)記憶部281、キュー監視部282から構成される。

30

【0311】

キューテーブル記憶部280に記憶されるキューテーブルの Q_{len_i} はそのキューの蓄積セル数、 Cnt_i は輻輳判定変数（初期値ゼロ）、 $BPID_t_i$ は最後のセル入力時のバックプレッシャ識別子(BPID)である。

【0312】

セルがキューに入力すると入力セル情報がキュー監視部282に転送される。また、BPIDはバックプレッシャ信号が出力可能から出力禁止に変化した時に図16の様にインクリメントされる。キュー監視部282はキュー長の監視を図18に示すアルゴリズムに基づいて行なう。

【0313】

図18にバックプレッシャを考慮したバッファのキュー長監視アルゴリズムの一例を示す。これはキュー長を監視して輻輳の有無を判断する例を示している。

40

【0314】

キュー Q_i にセルが入力されて、セル入力処理が開始されると、まず、 Q_{len_i} をインクリメントする（ステップS1）。

【0315】

次に、BPIDと $BPID_t_i$ を比較することにより、キュー長を監視すべきかそうでないかを判断する（ステップS2）。すなわち、バックプレッシャ信号により変動するキュー長が短くなった時に監視する。

【0316】

50

B P I D = B P I D t i の場合、キュー長を監視すべきであり、まず、 $C n t i \geq 1$ であるか否かを判断する（ステップ S 3）。このとき、 $C n t i \geq 1$ 場合は直前まで輻輳状態だと判断し、輻輳状態か否かを判断するためのしきい値 Q_{th} にヒステリシスを持たせた値（ $Q_{th} - H$ ）と $Q_{len i}$ 比較する（ステップ S 4）。

【0317】

一方、 $C n t i \geq 1$ でない場合は、 $Q_{len i} > Q_{th}$ を判定する（ステップ S 5）。ステップ S 4、ステップ S 5でキュー長 $Q_{len i}$ を $Q_{th} - H$ あるいは Q_{th} と比較した結果、これらのしきい値を越えているとき（すなわち、 $Q_{len i} > Q_{th} - H$ 、あるいは、 $Q_{len i} > Q_{th}$ であるとき）、 $C n t i := Q_{len i} - Q_{th} + H$ を行ない（ステップ S 6、ステップ S 7）、そうでなければ $C n t i := 0$ とする（ステップ S 8、ステップ S 9）。さらに、ス

10

【0318】

ステップ S 2で、 $B P I D = B P I D t i$ の場合、キュー長を監視すべきではなく、そのままステップ S 10に進む。

【0319】

ステップ S 10では、 $C n t i \geq 1$ を判定する。もし、そうであれば輻輳状態であり、ステップ S 11に進み、 $C n t i$ から「1」を減算し（ $C n t i := C n t i - 1$ ）、入力セルを廃棄するなら（ステップ S 12） $Q_{len i}$ から「1」を減算する（ $Q_{len i} := Q_{len i} - 1$ ）。 $C n t i \geq 1$ でないなら（ステップ S 10）輻輳状態ではない。

【0320】

最後に、 $B P I D t i$ に $B P I D$ の値を代入し、 $B P I D t i$ を更新する（ステップ S 13）。

20

【0321】

このように、図 18 に示したフローチャートによれば、バックプレッシャー信号が出力可能から出力禁止に変化した時から、次に、出力可能から出力禁止に変化する時までの最初のセル入力時の蓄積数のみを監視していることになる。

【0322】

なお、 H は、ヒステリシスに関するパラメタを示す。ヒステリシスはキュー長を監視して行なわれるセル廃棄や輻輳状態などの判定結果の振動を低減する作用がある。 $0 \leq H < Q_{th}$ であり、ヒステリシスが不要である場合はゼロに設定する。

30

【0323】

図 18 では、しきい値 Q_{th} とヒステリシスパラメタ H は全てのキューで同じ値を用いる場合を示しているが、キュー毎に異なる値を設定してもかまわない。

【0324】

ここで、しきい値 Q_{th} の決定方法を述べる。まず、バックプレッシャー信号によって制御されていないバッファを考える。このバッファは一定レート R_i で入力し、一定レート R_0 で出力するとする。 $R_i > R_0$ のとき、この負荷（ $R_i - R_0$ ）が加わっても t_T 時間廃棄を起こさないために必要なバッファ量 Q_{th} は、 $Q_{th} = (R_i - R_0) \times t_T$ により求められる。

【0325】

これに対し、この第 8 の実施形態では、図 16 に示すように、入力レートは R_i で一定であるが、入力バッファからの出力レートは T_{up} 時間（入力バッファからの出力禁止時間）についてはゼロ、 T_{dn} 時間（入力バッファからの出力可能時間）については R_e である。 T_{up} 、 T_{dn} の変化が少ないと仮定すると、先ほどと同じく、この負荷が加わっても t_T 時間廃棄を起こさないために必要なバッファ量 Q_{th} は、

40

【数 1】

$$Q_{th} = \frac{R_i T_{up} + (R_i - R_e) T_{dn}}{T_{up} + T_{dn}} t_T \quad \dots (1)$$

$$= (R_i - \frac{T_{dn}}{T_{up} + T_{dn}} R_e) t_T \quad \dots (2)$$

【0326】

となる。ここで、出力ポートからの平均の出力レート R_0 は、

10

【数2】

$$R'_0 = \frac{T_{dn}}{T_{up} + T_{dn}} R_e$$

【0327】

であるから、 $Q_{th} = (R_i - R_0) \times t_T$ となり、この第8の実施形態のセルバッファ装置は、バックプレッシャのない単純なバッファ装置の場合のしきい値と同じしきい値を設定することにより、ほぼ同様の効果を得られると考えられる。

【0328】

20

なお、輻輳状態を判定する代わりに、セルの廃棄判定を行なう場合でも全く同じアルゴリズムが適用できる。

【0329】

以上説明したように、上記第8の実施形態に係るバックプレッシャ信号を考慮するセルバッファ装置によれば、キュー長の監視結果が外部条件の影響を受け難く、安定したキュー長の監視が容易である。

【0330】

4. 遅延を考慮したセル数情報の修正（第9の実施形態）

次に、本発明に係るセル数情報を修正するセルバッファ装置の実施形態について説明する。

30

【0331】

4.1 第9の実施形態（遅延を考慮したセルバッファ装置）

図19に、第9の実施形態に係る遅延を考慮したセルバッファ装置の構成を示す。

【0332】

図19に示したセルバッファ装置は、大きく分けて、ひとつ以上のセルバッファB（Ba、Bb、...）を複数段（図19では2段）接続してなるセルバッファ網300と、各セルバッファBのセル数を入力してセルバッファBへ転送指示を行なう管理部301とから構成される。

【0333】

管理部301は、転送指示の履歴302を保持し、新たな転送指示を決定する際にこの転送指示履歴302を使用することを特徴とする。

40

【0334】

前記管理部301は、セルバッファBがセル数を管理部301へ送信した時刻から、管理部が決定する転送指示がそのセルバッファBに作用する時刻までに、そのセルバッファBに作用する転送指示の回数を前記転送指示履歴302より求め、その転送指示の回数とセル数から転送指示を決定する。

【0335】

図20は、本発明の遅延を考慮したセルバッファ装置におけるセル数情報の修正原理を説明する図である。

【0336】

50

入力リンク # 1、# 2 から入力したセルは、それぞれ前段のセルバッファ B a1、B a2に一時的に蓄積される。管理部 3 0 1 の転送指示に従いセルはセルバッファ B b に転送され、セルバッファ B b はセルを出力リンクへ出力する。

【 0 3 3 7 】

管理部 3 0 1 には、B a2のセル数 N_{a2} と、B b のセル数 N_b が伝送され、B a1、B a2に対して転送指示を出力する。

【 0 3 3 8 】

今、B a2からセル数情報が管理部 3 0 1 に到着し、その情報から転送指示を決定して B a2に転送を指示して、この転送指示により B a2のセル数が変化するまでの遅延時間を D セル周期とする。また、B b からセル数情報が管理部に到着し、その情報から転送指示を決定して B a1、B a2に転送を指示して、この転送指示により B a1または B a2から B b へセルが転送され B b のセル数が変化するまでの遅延時間も、説明を簡単にするため同じ D セル周期であるとする。

10

【 0 3 3 9 】

ここで、図 2 0 を、セルバッファ B a1、B a2が図 5 3 のクラス 1 セル蓄積部、クラス 2 セル蓄積部に相当し、セルバッファ B b がセル多重化 F I F O に相当する優先制御を行なうセルバッファ装置と考える。B a1、B a2、B b のセルバッファ間のセル転送の方針は図 5 3 と同じである。つまり B b の蓄積セル数を少なくすることが必要であるが、B b を努めて空にしない (= アンダフローしない) ように転送指示をする。さらに、転送指示は空指示をしないように行なう (B a2に蓄積セルが存在する場合にのみ B a2に対して転送指示を行なう)。

20

【 0 3 4 0 】

さて、時刻 t_2 で B a2から出力されたセル数情報 N_{a2} が管理部 3 0 1 へ到着し、時刻 t_3 に管理部 3 0 1 が転送指示を行なって時刻 t_4 にその転送指示により B a2からセルが出力したとする。管理部 3 0 1 が知りたい情報は、転送指示が実際に B a2に作用する時刻 t_4 における B a2のセル数である。この時点を図 2 0 に白丸で示した。

【 0 3 4 1 】

B a2が時刻 t_2 から t_4 の時間 (= D セル周期) で M_{a2} 回の転送指示を受け付けたとすると、この間に B a2は M_{a2} セル出力したことになり、時刻 t_4 における B a2のセル数は $N_{a2} - M_{a2}$ である (B a2へセルが入力しなかった場合)。

30

【 0 3 4 2 】

通常、管理部 3 0 1 から B a2までの転送指示の伝送遅延は一定だから、B a2が受信する転送指示回数 M_{a2} は、 $t_4 - t_2 = D = t_3 - t_1$ となるような t_1 と t_3 間で管理部 3 0 1 が出力した転送指示の回数と等しい。従って、転送指示履歴 3 0 2 より転送指示回数 M_{a2} を知ることにより、時刻 t_2 から t_4 における B a2の減少セル数を正確に求めることができる。

【 0 3 4 3 】

B b に関しても同様に考えることができ、時刻 t_2 から t_4 における B b のセル数の増加量 (B a2からの転送セルによる増加量) を正確に求めることができる。

【 0 3 4 4 】

時刻 t_2 から t_4 での B a2への入力セル数は、管理部 3 0 1 は知ることができない。B a2に対して空指示しないようにするためには、B a2のセル数を少なめに見積もればよい。従って、入力リンクから B a2へのセル入力、時刻 t_2 から t_4 においてなかったと仮定する。

40

【 0 3 4 5 】

また時刻 t_2 から t_4 での B b からのセル出力は、管理部 3 0 1 は知ることができない。B b がアンダフローしないようにするためには、B b のセル数を少なめに見積もる必要がある。従って、B b からのセル出力は、時刻 t_2 から t_4 において常に出力し続けたと仮定する (D セル周期であるので、D セルと仮定する)。さらに、B a1からの入力セルもゼロと仮定する。

50

【0346】

まとめると、修正された N_{a2} 、 N_b をそれぞれ N_{a2} 、 N_b とすると、

$$N_{a2} = N_{a2} + 0 - M_{a2}$$

$$N_b = N_b + M_{a2} - D$$

である。よって管理部301が B_{a2} へ転送指示を行なう条件は、

$$(N_b > Th) \text{ and } (N_{a2} > 0)$$

で与えられる。 B_{a1} への転送指示は、

$$(N_b > Th) \text{ and } (N_{a2} = 0)$$

の条件で行なう。 Th は通常は「1」である。以上のようにセル数を修正することにより、アンダフローや空指示を努めて少なくしつつ転送を指示することができる。

10

【0347】

4.2 一般的なバッファ網におけるセル数の修正方法

次に、一般的なバッファ網におけるセル数の修正方法について説明する。

【0348】

基本原理の説明には図21を用いる。図21には、すでに述べたようにセルバッファ網300と管理部301がある。セルバッファ網300は、着目しているセルバッファ B_x を含む。セルバッファ B_x は物理的には複数のバッファであっても良い。管理部はこのセルバッファ B_x からセル数情報を得て転送指示を与える。まず本発明を説明する際に使用する用語の定義を行う。

【0349】

4.2.1 用語の定義

時刻は、特に断らない限り現在時刻をゼロとする相対時刻とし、セル周期経過毎に1ずつ増加するものとする。このような時刻の定義は他の表現方法に変換可能であり一般性がある。また、ある時刻 t_1 からある時刻 t_2 までの期間を p などと表すことがあり、これを $p = [t_1, t_2)$ と表す。この期間 p は時刻 t_1 を含み、 t_2 は含まない。

20

【0350】

セルバッファ B_x の蓄積セル数を N_x で表す。 N_x は現在時刻 ($t = 0$) に管理部301に到着した蓄積セル数である。

【0351】

セルバッファ B_x に対し管理部301が時刻 t ($t \leq 0$) に出力した転送指示を $I_x(t)$ で表す。転送を指示した場合は $I_x(t) = 1$ であり、転送を指示しない場合は $I_x(t) = 0$ である。

30

【0352】

時刻 t ($t \leq 0$) に管理部301に到着したセルバッファ B_x のHOLセル数を $H_x(t)$ で表す。

【0353】

HOLセル数とは、1回の転送指示でセルバッファから一度に転送される可能性のあるセル数である。特に、現在時刻に管理部301に到着したHOLセル数情報 H_x を $H_x(0)$ と表す。

【0354】

本発明の管理部301は、時刻 t ($t \leq 0$) において得たセルバッファ B_x のHOLセル数 $H_x(t)$ および自らがセルバッファ B_x へ出力した転送指示 $I_x(t)$ を履歴として記憶する。

40

【0355】

ある期間 p に管理部がセルバッファ B_x に転送を指示した回数を $S_{mx}(p)$ と表す。 $p = [t_1, t_2)$ のとき、 $M_x[t_1, t_2)$ と表す。つまり、

【数3】

$$S m_x (p) = M_x [t_1, t_2) = \sum_{t=t_1}^{t_2-1} I_x (t) \quad \text{ただし、} p = [t_1, t_2)$$

【 0 3 5 6 】

図 2 3 に、セルバッファ B x から管理部 3 0 1 へのセル数情報の流れと、管理部 3 0 1 からセルバッファ B x への転送指示の流れを表したタイムチャートを示す。

【 0 3 5 7 】

この図の様に、時刻 t f、t t、t h を定める。時刻 t f におけるセルバッファ B x の蓄積セル数 N x および時刻 t h におけるセルバッファ B x の H O L セル数 H x は、管理部 3 0 1 が現在時刻においてセルバッファ B x に関し直接知り得る最新の情報であるとする。このとき管理部 3 0 1 は、時刻 t t におけるセルバッファ B x の蓄積セル数を予測したいものとする。

10

【 0 3 5 8 】

時刻 t f、t t、t h に対して、時刻 t 1、t 3、t 2 を定める。時刻 t 1 に管理部 3 0 1 が出力した転送指示がセルバッファ B x に届く時刻を t f とし、時刻 t 2 に管理部 3 0 1 が出力した転送指示がセルバッファ B x に届く時刻を t t とする。同様に時刻 t 3 に管理部 3 0 1 が出力した転送指示がセルバッファ B x に届く時刻を t h とする。

【 0 3 5 9 】

時刻 t 1、t 2、t 3 は現在時刻より過去の時刻であるとする。つまり、 $t 1 < t 2 \leq 0$ 、 $t 3 \leq 0$ とする。

20

【 0 3 6 0 】

4 . 2 . 2 . 基本原理

ここで本発明の基本原理を説明する。本発明の管理部 3 0 1 は転送指示の履歴から蓄積セル数の変化を予測する。前述のセルバッファ B x について、時刻 t t の蓄積セル数を求めたい場合、それより過去にの時刻 t f における蓄積セル数 N x に対して、転送指示の履歴を用いてこの期間に転送指示により転送されたセル数を予測して修正する。修正は、1 回の転送指示により転送されるセル数が 1 以上であることに基づく。特に時刻 t h 以降の最初の転送指示により転送されるセル数については H O L セル数 H x 以上であることに基づく。

30

【 0 3 6 1 】

4 . 2 . 3 . 3 つの関数

以上の基本原理より修正に使用する 3 つの関数、f ()、g ()、S c () を定める。これらの関数は、時刻 t f から時刻 t t においてセルバッファ B x が転送指示によるセルの転送でどのくらい変化するかを求めるものである。t f、t t と比較して t h がどのような値であるかにより 3 通りの場合に分類でき、適用する関数が異なる。

【 0 3 6 2 】

(1) 関数 f ()

関数 f () は、 $t h = t f < t t$ の場合に使用する。この時のタイムチャートを図 2 2 に示す。p 1 [t 1、t 2) とする。H x = 0 かつ S m x (p 1) = 0 のときは、期間 p 1 の最初の転送指示では H x セル転送され、それ以降の転送指示では 1 セルずつ転送されると考える。時刻 t f から時刻 t t において変化するセルバッファ B x のセル数を計算する関数 f () は、H x と期間 p 1 の転送指示回数 S m x (p 1) を引数として、次のように求めることができる。

40

【 0 3 6 3 】

【 数 4 】

$$f(Sm_x(p1), H_x) \triangleq \begin{cases} Sm_x(p1) + H_x - 1 \cdots \text{ただし、} Sm_x(p1) \neq 0 \\ \quad \text{かつ } H_x \neq 0 \\ Sm_x(p1) \quad \cdots \text{その他の場合} \end{cases}$$

【0364】

(2) 関数 $g(\quad)$

関数 $g(\quad)$ は、 $t_h < t_f < t_t$ の場合に使用する。この時のタイムチャートを図23に示す。 $p1 = [t1, t2)$ とし、 $p2 = [t3, t1)$ とする。

10

【0365】

この場合、時刻 t_h 以降に到着した初めての転送指示が期間 $p2$ に管理部301から出力されたか否かで場合分けが必要である。期間 $p2$ に管理部301から転送指示が出力されていないから、期間 $p1$ において H_x の値を用いることができる。時刻 t_f から時刻 t_t において変化するセルバッファ B_x のセル数を計算する関数 $g(\quad)$ は、 H_x と期間 $p1$ 、 $p2$ の転送指示回数 $Sm_x(p1)$ 、 $Sm_x(p2)$ を引数として、次のように求めることができる。

【0366】

【数5】

$$g(Sm_x(p1), Sm_x(p2), H_x) \triangleq \begin{cases} Sm_x(p1) + H_x - 1 \cdots \text{ただし、} Sm_x(p2) = 0 \\ \quad \text{かつ } Sm_x(p1) \neq 0 \\ \quad \text{かつ } H_x \neq 0 \\ Sm_x(p1) \quad \cdots \text{その他の場合} \end{cases}$$

20

【0367】

(3) 関数 $Sc(\quad)$

関数 $Sc(\quad)$ は、 $t_f < t_t = t_h$ の場合に使用する。この時のタイムチャートを図24に示す。 $p1 = [t1, t2)$ とする。管理部301が知り得る最新の蓄積セル数は時刻 t_f のときの値であるが、この場合は HOL セル数に関しては時刻 t_f から時刻 t_t までの全ての期間における情報が管理部301にすでに到着している。

30

【0368】

転送指示がセルバッファ B_x へ到着した時点の HOL セル数がわかれば、転送されるセル数もわかることから、時刻 t_f から時刻 t_t において変化するセルバッファのセル数を計算する関数 $Sc(\quad)$ は、 $I_x(t)$ 、 $H_x(t)$ を用い、期間 $p1$ を引数として、次のように求めることができる。

【0369】

【数6】

$$Sc_x(p1) \triangleq \sum_{t=t1}^{t2-1} I_x(t) H_x(t-t2) \quad \cdots \text{ただし } p1 = [t1, t2)$$

40

【0370】

以上説明してきたようなセル数の修正方法を、例えば、図1に示したようなクラス間の優先制御を行なうセル多重化装置に適用することは容易である。

【0371】

さらに、例えば、図6に示したような VC 間公平キューイングを行なうセル多重化装置に、セル数の修正方法を適用することも容易である。

50

【 0 3 7 2 】

4.3 第9の実施形態に係るセルバッファ装置の利点

以上説明したように、上記第9の実施形態によれば、セルバッファとセルバッファを管理する管理部との間に遅延のあるセルバッファ装置は、転送指示を決定する際に転送指示の履歴を用いて、セルバッファから伝送されてきたセル数を、転送指示がセルバッファに作用する時刻のセル数に修正するため、遅延時間を原因とした性能の劣化が少ない。

【 0 3 7 3 】

5. セルスイッチ (第10の実施形態)

5.1 第10の実施形態 (セルスイッチ)

次に、第10の実施形態として、ここまで説明してきた本発明に係るセル多重化装置、セルバッファ装置を組み合わせたセルスイッチについて説明する。 10

【 0 3 7 4 】

図25は、複数の入力ポート# i ($i = 1 \sim N$) から入力したセルを一時的に蓄積し、バックプレッシャに応じて出力する複数の入力バッファIBと、内部の輻輳状態に応じてバックプレッシャ信号を出力するバックプレッシャ付きセルスイッチ400と、入力バッファIBとの間で情報を転送するカゴ管理部KMとからなる。

【 0 3 7 5 】

この構成例において扱うセルは、大きくは高優先セル、低優先セルの2つの優先度を持ち、さらに低優先クラスは複数のクラスにわかれている。通常はCBR、VBRなどのリアルタイム情報を伝送するセルを高優先セルとし、ABR、UBRなどのノンリアルタイム情報を伝送するセルを低優先セルとする。 20

【 0 3 7 6 】

まず、バックプレッシャ付きセルスイッチについて説明する。

【 0 3 7 7 】

カゴスイッチの特徴は入力バッファIBとカゴ管理部KMにある。バックプレッシャ付きセルスイッチ400は入力したセルをある定められた時間内に出力できればよく、そのアーキテクチャは以下に示すものだけではない。例えば、電子情報通信学会技術報告SSE93-6“バッファ容量拡張可能なATMスイッチ：XATOM”にて示されているATMスイッチでもよい。

【 0 3 7 8 】

図25に示すバックプレッシャ付きセルスイッチ400を一例として説明する。1段目単位スイッチSE1、2段目単位スイッチSE2、出力バッファOBからなり、SE1とSE2を並列リンクにて2段デルタ網接続し、SE2の出力したセルをOBへ転送して出力ポートへ出力する。SE1からOBの間は k 倍速で動作する。 k の値はセルスイッチがノンブロッキングになるように定める。例えばSE1、SE2が8入力8出力の時、全体で16入力16出力とする場合には $k = 2$ とすればよい。 30

【 0 3 7 9 】

図25のバックプレッシャ付きセルスイッチ400は、高優先セルと低優先セルの2段階の優先度を扱う構成である。

【 0 3 8 0 】

低優先セルの交換にのみ着目すると基本構成は1995年電子情報通信学会総合大会B-589“ATM2段スイッチ網のバッファ容量拡張法”にて示されているものとほぼ等しい。低優先セルは、SE1、SE2で交換、コピーされる。 40

【 0 3 8 1 】

低優先セルがOBで輻輳した場合には低優先セル用バックプレッシャ信号を用いてIBのその出力ポート行きのセル出力を抑制する。このバックプレッシャ信号により基本的にはOBおよびSE2のバッファを溢れさせないようにすることが可能である。しかし、マルチキャストセルの影響のためこれらのバッファが溢れそうになる場合があり、それに備えて低優先セル用オーバーフローバックプレッシャ信号がOBからSE2へ出力し、低優先セル非常用バックプレッシャ信号がSE2からIBへ出力する。IBが低優先セルを蓄積 50

するのでSE2、OBでは廃棄されない。

【0382】

OBは、低優先セル用オーバーフローバックプレッシャ信号をSE2へ出力する代わりに、これを低優先セル非常用バックプレッシャ信号としてIBへ出力する構成でもよい。この場合は、SE2が低優先セル非常用バックプレッシャ信号を出力する必要はない。しかし、先に説明したように、低優先セル非常用バックプレッシャ信号をSE2から出力した方がSE2内でのバッファ共用化効果を有効に活用できるという利点がある。

【0383】

高優先セルは、低優先セルと同様、SE1、SE2で交換、コピーされる。高優先セルに関してはバックプレッシャ付きセルスイッチからIBへのバックプレッシャ信号はない。OBはその高優先セルの輻輳を反映した高優先セル用バックプレッシャ信号によりSE2のセル出力を制御する。IBは高優先セルを蓄積しない。SE2内に空きバッファが無ければ廃棄されることがある。

10

【0384】

次に、出力バッファについて説明する。

【0385】

図26に、カゴスイッチの出力バッファOBの構成例を示す。

【0386】

SE2からOBへ転送されたセルは識別部410にて優先度を識別され各優先度に対応したキューに入れられる。各キューのセルは選択部411により優先度に従って選択され出力ポートへ出力される。OBは、低優先セルのキュー長があるしきい値を越えた場合に低優先セル用バックプレッシャ信号を用いてIBの低優先セルの、そのOB行きの出力を止める。高優先セルのキュー長があるしきい値を越えた場合に高優先セル用バックプレッシャ信号を用いてSE2の高優先セルの、そのOB行きの出力を止める。低優先セルのマルチキャストセルの影響により低優先セルのキュー長が伸び続けた場合に備えて低優先セル用オーバーフローバックプレッシャ信号を用いてSE2の低優先セルの、そのOB行きの出力を止める。

20

【0387】

次に、単位スイッチについて説明する。

【0388】

図27に、カゴスイッチに用いる1段目単位スイッチSE1の構成例を示す。SE1は一般的によく知られているスイッチLSIである。入力リンクより入力したセルに対し、そのセルヘッダを識別部415で識別して所定の交換処理を行い出力リンクへ出力する。スイッチ全体がk倍速で動作していることから内部のセルバッファの蓄積セル数は多くならない。これについては1994年電子情報通信学会秋季大会B-439“並列リンクのあるスイッチ網の検討”にて知られている。

30

【0389】

図28に、カゴスイッチに用いる2段目単位スイッチSE2の構成例を示す。SE2は内部のキューが出力リンク毎に優先度毎のキューになっている。各優先度毎のキューの出力は選択部421により選択される。高優先セルが存在するときには、低優先セルの有無に関わらずその出力リンクに高優先セルを出力する。ただし、選択部421に入力している高優先セル用バックプレッシャ信号が、高優先セルの出力の禁止を指示している場合は高優先セルを出力せず、低優先セルを出力しても良い。低優先セルは低優先セル用オーバーフローバックプレッシャ信号により出力が禁止される場合がある。

40

【0390】

SE2は、複数の優先度、複数の出力リンク間で物理的なバッファ領域を共有することにより、SE2内部のバッファを効率的に使用することが可能である。バッファ状態管理部422は各キュー毎のキュー長、および優先度毎に全ての出力リンク行きのキュー長を合計した結果などをモニタする。

【0391】

50

例えば、高優先セル用キューは各キュー毎にしきい値を持ち、そのしきい値を越えた場合にそのキューの入力セルを廃棄する。低優先セル用キューは全てのキュー長の合計をモニタし、しきい値を越えた場合は低優先セル非常用バックプレッシャ信号により低優先セルの入力を禁止するようにIBに伝える。

【0392】

SE1はSE2と同じ単位スイッチを使用してもかまわない。

【0393】

図25のバックプレッシャ付きセルスイッチ400は2段の単位スイッチSE1とSE2によりセルを交換するが、内部の経路は各入力ポートと出力ポートのペアにひとつずつしか定めない。そのため各単位スイッチはセルヘッダを見てセルをコピーしつつ交換する際に複雑な経路の選択アルゴリズムは必要ない。出力ポート数が32ポート程度以下であるなら、セルのヘッダに出力ポートのビットマップ情報を付加することが容易に可能である。SE1、SE2はそのビットマップ情報に基づいてセルを交換する。

10

【0394】

次に、入力バッファIBについて説明する。

【0395】

図29に、カゴスイッチの入力バッファIBの構成例を示す。

【0396】

IBに入力した高優先セルは、ユニキャスト、マルチキャストとも優先的にSE1へ出力する。そのため高優先セルのセル遅延揺らぎは非常に小さい。またIBに高優先セル用のバッファは存在しない。

20

【0397】

本発明の特徴は低優先セルの処理方法にある。以降、低優先セルについてのみ説明する。

【0398】

図29において、入力バッファIBは、主に、入力したユニキャストセルを一時的にそのセルの出力ポート毎に蓄積する複数の出力ポート毎ユニキャストセル管理部430と、入力したマルチキャストセルを一時的に蓄積するマルチキャストセル管理部431、およびこれらの管理部のセルを選択して入力バッファから出力する出力ポート選択スケジューラ432を持つ。

【0399】

出力ポート毎ユニキャストセル管理部430は、入力したセルを一時的にそのセルのクラス毎に蓄積する複数のクラス毎ユニキャストセル管理部433と、そのクラス毎ユニキャストセル管理部433から出力したセルを多重化するクラス多重化FIFO434を持つ。

30

【0400】

マルチキャストセル管理部431は、入力したセルを一時的にそのセルのクラス毎に蓄積する複数のクラス毎マルチキャストセル管理部435と、そのクラス毎マルチキャストセル管理部435から出力したセルを多重化するクラス多重化FIFO436を持つ。

【0401】

出力ポート毎ユニキャストセル管理部430、マルチキャストセル管理部431のそれぞれのクラス多重化FIFO434、436より出力するセルは出力ポート選択スケジューラ432により選択され、入力バッファIBから出力される。

40

【0402】

出力ポート選択スケジューラ432は外部からの低優先セル用バックプレッシャ信号、低優先セル非常用バックプレッシャ信号を考慮してセルを選択する。

【0403】

クラス毎ユニキャストセル管理部433は、蓄積されているセルのうち、カゴ管理部KMから出力を許可されたセルの集合を管理するカゴと呼ばれるキュー（カゴ440）と、それ以外のセルを管理する前カゴと呼ばれるキュー（前カゴ441）を持つ。

【0404】

50

前カゴ441に入ったセルは単位カゴという単位で管理される。図29において前カゴ441、カゴ440内の縦一列が単位カゴを示す。入力したセルは単位カゴ毎にV C (C G)間で公平に管理される。前カゴ441のセルはカゴ管理部K Mの指示により、単位カゴを単位としてカゴ440に転送される。

【0405】

カゴ440のセルはカゴ管理部K Mの指示により、セルを単位として転送されクラス多重化F I F O 4 3 4で多重化される。

【0406】

クラス毎マルチキャストセル管理部435の動作は後述する。

【0407】

I Bの全ての出力ポート、全てのクラスのキューでセルバッファを共有することにより、有限のバッファ量を有効に活用することが可能である。

【0408】

I Bに到着したセルがI Bから出力されるためには、セルが到着したことをI BからK Mへ通知し、K MがI Bからの出力を指示することが必要である。従ってI BとK Mの間の情報伝送に遅延がある場合はその遅延時間だけI Bからのセルの出力が遅延する。しかし遅延時間の大きさは数セル周期程度と考えられるため、トラヒック特性に及ぼす影響は少ないと考えられる。

【0409】

ユニキャストセルを出力ポート毎に優先制御するために、I Bには出力ポート毎に、クラス毎ユニキャストセル管理部433とクラス多重化F I F O 4 3 4がある。図1のセル多重化装置で説明した優先制御と基本的に同じ方法を使用する。

【0410】

図1はセル多重化装置であるが、カゴスイッチはセルスイッチでありI Bはセルスイッチの入力バッファである。そのためクラス多重化F I F OはひとつのI Bに全ての出力ポートに対応して複数存在し、そのセルは出力ポート選択スケジューラ434により選択されS E 1へ出力される。

【0411】

図1のセル多重化装置のクラス毎セル蓄積部13は、図29のI Bのクラス毎ユニキャストセル管理部433に相当し、図1のセル多重化装置のクラス多重化F I F O 1 4は、図29のI Bのクラス多重化F I F O 4 3 4に相当する。

【0412】

図1のセル多重化装置のクラス管理部12は、図25のK Mに対応する。クラス毎ユニキャストセル管理部433からクラス多重化F I F O 4 3 4へのセルの転送をカゴ管理部K Mが指示することにより優先制御を行なうことができる。

【0413】

マルチキャストセルの優先制御は後述する。

【0414】

ユニキャストセルをクラス毎にV C間公平キューイングするために、クラス毎ユニキャストセル管理部433に前カゴ441とカゴ440がある。図6のセル多重化装置で説明したV C間公平キューイングと基本的に同じ方法を使用する。

【0415】

図6はセル多重化装置であり、かつクラスも存在しないが、I Bはセルスイッチの入力バッファでありクラスも存在する。前カゴ441、カゴ440は各I Bについて、全ての出力ポート、全てのクラスに対応して存在し、そのセルはクラス多重化F I F O 4 3 4により多重化された後、出力ポート選択スケジューラ432により選択されS E 1へ出力される。

【0416】

図6のセル多重化装置のカゴ213は、図29のI Bのカゴ440に相当し、図6のセル多重化装置の同一の時刻までに出力すべきセル集合214が、図29のI Bの単位カゴに

10

20

30

40

50

相当する。図 6 のセル多重化装置のカゴ管理部 1 2 は、図 2 5 のカゴ管理部 K M に相当する。図 2 9 の前カゴ 4 4 1 からカゴ 4 4 0 への単位カゴの転送をカゴ管理部 K M が指示することにより V C 間公平キューイングを行なうことができる。

【 0 4 1 7 】

マルチキャストセルの V C 間公平キューイングは後述する。

【 0 4 1 8 】

ユニキャストセルを I B 内で V C 間公平キューイングするために、クラス毎ユニキャストセル管理部 4 3 3 に前カゴ 4 4 1 とカゴ 4 4 0 がある。図 1 0 のセルグループ F I F O 2 3 2 a のポインタチェーンによるデータ構造例で説明した方法と基本的に同じ方法を使用する。

10

【 0 4 1 9 】

図 1 0 のセルグループ F I F O 2 3 2 a が、図 2 9 の I B の前カゴ 4 4 1 に相当し、図 1 0 の出力待ちセルグループ F I F O 2 3 2 b が、図 2 9 の I B のカゴ 4 4 0 に相当する。また、図 1 0 のセルグループは図 2 9 の I B の単位カゴに相当する。

【 0 4 2 0 】

次に、図 3 0 を参照して、入力バッファ I B のキューイング処理について説明する。

【 0 4 2 1 】

セルは、まず前カゴ 4 4 1 に蓄積され、その後単位カゴを単位としてカゴ 4 4 0 に転送される。カゴ 4 4 0 のセルはセルを単位としてカゴ 4 4 0 から出力される。それぞれのきっかけについてはカゴ管理部 K M が指示を行なう。

20

【 0 4 2 2 】

この構成例では、全ての V C は、あるコネクショングループ (C G) に所属する。コネクショングループの概念を持つことにより複数の V C にひとつのリソース (帯域やセルバッファ) を共有させることが可能となる。例えば、離れたふたりのユーザ間の通信に対して 1 本の V C 分の課金しか行なわないときに、そのふたりのユーザ間に複数の V C を設定する自由度を与えることが可能となる。もしくは、V C 数に応じて公平な帯域の割当を行なう網において、あるふたりのユーザが多数の V C を設定して不当に大きな帯域を得ようとするのを防止することが可能となる。

【 0 4 2 3 】

前カゴ 4 4 1 においてセルをどの単位カゴにいれるかを決定するために V C テーブルと C G テーブルを使用する。

30

【 0 4 2 4 】

V C テーブルを用いて入力したセルのコネクション識別子 (例えば V C 4) からコネクショングループ識別子 (例えば C G 5) を知ることができる。C G テーブルを用いてコネクショングループ識別子から C G 毎に定められた値を知ることができる。この構成例における C G 毎の値は N_x 、 N_c 、 P_{tr} 、 Q_{len} 、 C_{nt1} 、 C_{nt2} 、 C_{nt3} 、 $B P I D t$ である。

【 0 4 2 5 】

N_x 、 N_c 、 P_{tr} 、 Q_{len} は、図 1 0 のセルグループ F I F O のポインタチェーンによるデータ構造例で説明した。 N_x はその C G の重み、 N_c は作業変数、 Q_{len} はその C G の蓄積セル数、 P_{tr} はその C G のセルが蓄積されている単位カゴの末尾へのポインタである。

40

【 0 4 2 6 】

一方、 Q_{len} 、 C_{nt1} 、 C_{nt2} 、 C_{nt3} 、 $B P I D t$ は、図 1 5 および図 1 6 のバックプレッシャ考慮キュー監視部を持つセルバッファ装置で説明した。

【 0 4 2 7 】

Q_{len} はその C G の蓄積セル数、 C_{nt1} は廃棄すべきセル数、 $B P I D t$ は最後のセル入力時のバックプレッシャ識別子である。よりきめの細かいキュー監視を行なうため C_{nt1} と同じ機能を持つ変数 C_{nt2} 、 C_{nt3} が存在する。それぞれキュー長監視のしきい値が異なる。 C_{nt2} はセルヘッダ内にある C L P (C e l l L o s s P r i o r i t y : セ

50

ル廃棄優先)ビットが低優先(CLP=1)であるセル用の廃棄すべきセル数、Cnt3はこのキュー長監視において輻輳かどうかを判定するために使用する値である。輻輳かどうかの情報は、例えばセルヘッダ内にあるEFCIをマークするかどうかや、RMセルの輻輳表示ビットをマークするかどうかを決定するために用いる。

【0428】

あるクラスでVC間公平キューイングを行わない設定をすることも可能である。その場合は、そのVCが所属するCGのNxを十分大きな値に設定すればよい。

【0429】

さて、図25に示すバックプレッシャ付きセルスイッチは、IBからOBまでのセル伝送の遅延時間が大きいことなどの影響により、低優先セル用バックプレッシャ信号が出力可能と出力禁止を繰り返す周期が比較的大きい。従ってIBのCG毎のキュー長Qlenを観測すると、主に低優先セルバックプレッシャ信号に同期して図14の様に振動を繰り返す。この様に、バックプレッシャ信号により振動するキュー長の監視には、図15、図16で説明したバックプレッシャ考慮キュー長監視アルゴリズムが適している。

10

【0430】

前カゴ441のとカゴ440に蓄積されているセルのクラス全体の合計セル数についてもバックプレッシャ考慮キュー長監視により入力セルの廃棄やEFCIのマークなどを行なって、一部のクラスがIBのバッファ領域の大半を独占することを防ぐことが可能である。

【0431】

図31に、カゴスイッチに用いる入力バッファIBのカゴ管理部KMとのインタフェースの一例を示す。

20

【0432】

KMから各出力ポートのそれぞれに対し、以下に示すような各出力ポート毎の情報が入力される。

【0433】

- ABR入力セル数
- 指示クラス
- 前カゴからカゴへ単位カゴの転送指示(図31の転送指示T1)
- カゴからクラス多重化FIFOへセルの転送指示(図31の転送指示T2) ABR入力セル数はABR処理部450へ転送される。ABR処理部450は、ABR入力セル数を知ることによりそのIBとペアの出力ポートに流入するトラヒックをモニタすることができる。その情報を用いてIBを通過するRMセルのペイロードの書き換えなどを行なう。

30

【0434】

カゴ440からクラス多重化FIFO434へセルの転送をKMから指示されると、IBは指示クラスのユニキャストのカゴ440からセルをひとつ取りだしクラス多重化FIFO434に転送する。また前カゴ441からカゴ440へ単位カゴの転送指示があると、IBは指示クラスのユニキャストの前カゴ441から単位カゴをひとつ取りだしカゴ440に転送するとともに、カゴ440からセルをひとつ取りだしクラス多重化FIFO434に転送する。

40

【0435】

単位カゴの転送指示をそのクラスのセルの転送指示を行なう時のみに限定することにより、転送指示の情報量を減らし、より容易に実装することができる。

【0436】

IB内で前カゴ441、カゴ440、クラス多重化FIFO434のセル数が変化するのは、IBへセルが入力した時、KMが転送指示を行なった時、IBからセルが出力した時である。従って、それに応じてセル数情報をKMへ出力する。ひとつのIBがKMへ出力する情報は、以下のものである。

【0437】

各出力ポート毎に、入力したセルについて(ユニキャストの場合はひとつの出力ポートに

50

ついでのみ出力)、

- 入力セルの有無とクラス

各出力ポート毎に、転送の指示クラスについて、

- (転送後の)前カゴセル数(図31のCE1)

- (転送後の)カゴセル数(図31のCE2)

- クラス多重化FIFOへの転送セル数(図31のCE3)

各出力ポート毎に、ユニキャストのみの、

- クラス多重化FIFOセル数(図31のCE4)

図31を用いて説明する。入力したセルの情報はIBからKMへ入力セルの有無とクラスとして伝送される。また、ユニキャストとマルチキャストの前カゴセル数が加算され、クラス選択部L5により転送の指示クラスで選択されたものが前カゴセル数として伝送される。さらに、ユニキャストとマルチキャストのカゴセル数が加算され、クラス選択部L6において指示クラスで選択されたものがカゴセル数となり伝送される。そして転送指示によりカゴからクラス多重化FIFO434へ移動したセル数はユニキャストとマルチキャストとが加算され、クラス選択部L7において指示クラスで選択されて伝送される。クラス多重化FIFOセル数はユニキャストについてのみKMへ伝送される。マルチキャストのクラス多重化FIFOセル数は、クラス毎マルチキャストセル管理部435からのセルの出力を制御するためのバックプレッシャ信号として使用する。

10

【0438】

マルチキャストのクラス多重化FIFO436は複数の出力バッファOBからのバックプレッシャ信号により出力が制御されていることなどから、クラス多重化FIFOセル数がゼロでなくともセルスイッチの出力ポートがアンダフローする可能性がある。ユニキャストセルがその影響を受けてアンダフローしないように、この構成例ではユニキャストのクラス多重化FIFOセル数とは加算していない。

20

【0439】

ここで、マルチキャストセルの優先制御およびVC間公平キューイングについて説明する。

【0440】

マルチキャストVCをスループットの視点で考える。一般的に、ある時点のセルスイッチの状態を考えると、その複数の出力ポートには、輻輳している出力ポートもあれば、負荷が低い出力ポートもある。つまり、VCが使用できるスループットは出力ポートによって異なることになる。マルチキャストVCは、ひとつの入力ポートから複数の出力ポートへ同じスループットが同時に出力されるから、そのマルチキャストVCが得ることのできるスループットは出力先の中で最も輻輳している出力ポートのスループットになるのがよい。従って、IBからマルチキャストセルを出力する場合は、最も輻輳している出力ポートのスループットで出力することが望ましい。

30

【0441】

ある出力ポートのあるクラスに与えられるVC(CG)あたりのスループットを、IBは前カゴ441からカゴ440への単位カゴの転送指示の頻度より知ることができる。また、ある出力ポートのあるクラスに与えられるスループットを、IBはカゴ440からクラス多重化FIFO434へのセルの転送指示の頻度より知ることができる。さらにある出力ポートに与えられるスループットを、IBは低優先セル用バックプレッシャ信号の頻度より知ることができる。これらの転送指示やバックプレッシャ信号は出力ポート毎に与えられる。

40

【0442】

最も輻輳している出力ポートのスループットで出力するためには、マルチキャストセルを、そのセルの出力ポートの全てに関してこれらの信号が指示または許可した時点で転送すればよい。複数の出力ポートへの指示または許可は同時である必要はない。

【0443】

IBにおいて、低優先マルチキャストVCを管理する場合、本来そのクラス毎かつ出力ポ

50

ートパターン毎に管理することが望ましいのであるが、一般的にはその組合せの数は大きく、また通信中にその出力ポート数が増減する可能性があるため、その様な実装は現実的ではない。図25に示すようなカゴスイッチでは全てのマルチキャストコネクションをクラス毎にのみ管理し、出力ポートパターン毎での管理はしない。ある程度のHOL(Head Of Line)ブロッキングが発生することになる。

【0444】

図32は、クラス毎マルチキャストセル管理部435を説明するための図である。この管理部は、出力ポート表、VC間公平キュー、単位カゴセル数、カゴセル数、出力許可済みセル数を保持し、管理するものである。

【0445】

クラス毎ユニキャストセル管理部433が前カゴとカゴでセルを管理していたのに対して、クラス毎マルチキャストセル管理部435は前カゴに相当するVC間公平キュー460のみによりセルを管理する。ただ本発明のクラス毎マルチキャストセル管理部435は、管理の計算上では前カゴ、カゴの概念を用いてユニキャストと同様のセル管理を行なう。

【0446】

クラス毎マルチキャストセル管理部435は、出力ポート毎に次の情報を入力する。

【0447】

- ・前カゴからカゴへの単位カゴ転送指示
- ・カゴからクラス多重化FIFOへのセル転送指示

また出力ポート毎に次の情報を出力する。

【0448】

- ・前カゴセル数
- ・カゴセル数
- ・転送セル数

VC間公平キュー460は単位カゴのFIFOである。VC間公平キュー460に入力したセルは、各VC(CG)についてひとつの単位カゴに予め定められた数までしか入力できない。この構成は図30で説明したユニキャストの前カゴの構成と同様である。クラス多重化FIFO436からバックプレッシャ信号が入力しており、状況によってVC間公平キュー460からのセル出力を禁止する。

【0449】

出力ポート表は、出力ポート別に、前カゴセル数と先頭単位カゴへのポインタを記憶する。前カゴセル数は単位カゴセル数の合計、先頭単位カゴは計算上の前カゴの先頭の単位カゴへのポインタである。

【0450】

単位カゴセル数はVC間公平キュー460の単位カゴに対応して出力ポート毎のデータを記憶する。基本的には対応する単位カゴに蓄積しているマルチキャストセルの出力ポート毎のセル数を記憶する。例えば、図32のVC間公平キュー460の先頭の単位カゴには5セル蓄積されている。そのうち、ポート#1とポート#Nに出力するセルが3セル、ポート#1とポート#2に出力するセルが2セルあるので、単位カゴセル数のポート#1には5、ポート#2には2、ポート#Nには3が記憶される。

【0451】

カゴセル数は計算上のカゴのセル数である。また、出力許可済みセル数は計算上の出力を許可されたセル数である。

【0452】

ある出力ポートに対して単位カゴ転送を指示された場合、その出力ポートの先頭単位カゴの単位カゴセル数をカゴセル数に加算し、単位カゴセル数をゼロにする。さらにその出力ポートのカゴセル数から1を減算して出力許可済みセル数に1を加算する。そして出力ポート表の先頭単位カゴのポインタをVC間公平キュー460の次の単位カゴへひとつ進める。このときの転送セル数は1セルである。

【0453】

10

20

30

40

50

先頭単位カゴのポインタを進めることにより、ポインタはV C間公平キュー460の先頭の単位カゴから離れていく。この距離に上限を設け、ある距離になったら単位カゴ転送指示を無視してポインタを進めないように制御する。これは、単位カゴの転送が指示された時刻と実際にセルを出力する時刻との差を大きくしないようにする効果がある。

【0454】

ある出力ポートに対してK Mからセル転送が指示された場合、その出力ポートのカゴセル数から1を減算して出力許可済みセル数に1を加算する。このときの転送セル数は1セルである。

【0455】

クラス多重化F I F O 436からのバックプレッシャ信号が出力を許可している時、セルをV C間公平キュー460から出力することができる。その際、そのセルの全ての出力ポートについて出力許可済みセル数を1ずつ減算する。出力許可済みセル数がゼロの出力ポートが存在する場合は減算できないので、そのセルはまだ出力できない。

10

【0456】

出力許可済みセル数は単位カゴ転送指示およびセル転送指示により増加する。ある出力ポートの出力許可済みセル数がある値以上になった場合には、その出力ポートについて、クラス毎マルチキャストセル管理部435から出力する前カゴセル数とカゴセル数をゼロとし、単位カゴ転送指示とセル転送指示を無視する。これは実際のセルの転送スループット以上の転送指示を無視する効果がある。

【0457】

次に、図33～図39を参照して、クラス毎マルチキャストセル管理部の動作について説明する。

20

【0458】

図32の状態から、出力ポート#1に対して単位カゴの転送をK Mから指示された場合の動作を図33に示す。ポート#1の先頭単位カゴの単位カゴセル数をカゴセル数に加算し、単位カゴセル数ゼロにする。さらにポート#1のカゴセル数から1を減算して出力許可済みセル数に1を加算する。K Mへの転送セル数は出力ポート#1に関して1セルである。

【0459】

出力ポート表の先頭単位カゴのポインタをV C間公平キュー460の次の単位カゴへひとつ進める。

30

【0460】

さらに、出力ポート#1に対してセルの転送をK Mから指示された場合の動作を図34に示す。ポート#1のカゴセル数から1を減算して出力許可済みセル数に1を加算する。K Mへの転送セル数は出力ポート#1に関して1セルである。

【0461】

図34のV C間公平キュー460の先頭のセルの出力ポートは#1と#Nである。出力許可済みセル数は出力ポート#1に関してはゼロでない(=2)が、出力ポート#Nに関してはゼロであるのでまだセルを出力することはできない。

【0462】

セルが到着した時の動作を図35に示す。到着したセルの出力ポートは#1と#Nである。またこのセルは今までにV C間公平キュー460に蓄積していなかったV Cのセルであるとする。この場合セルはV C間公平キュー460の先頭の単位カゴに入力される。

40

【0463】

セルを入力したとき、そのセルの全ての出力ポートの単位カゴセル数を1ずつ加算する。加算する単位カゴは出力ポート毎に決定する。基本的にはセルを入力した単位カゴに1を加えるが、出力ポート表の先頭単位カゴがセルを入力した単位カゴよりも後ろにある場合は出力ポート表の先頭単位カゴが指している単位カゴセル数に1を加える。従って、図35では出力ポート#1については先頭単位カゴが指している2番目の単位カゴの単位カゴセル数に1を加算する。出力ポート#Nについてはセルを入力した単位カゴである1番目

50

の単位カゴの単位カゴセル数に 1 を加算する。

【 0 4 6 4 】

さらに、またセルが到着した時の動作を図 3 6 に示す。新たに到着したセルの出力ポートは # 1 と # 2 である。またこのセルは現在 V C 間公平キュー 4 6 0 にセルが蓄積している V C のセルであるとする。図 3 6 の場合このセルは V C 間公平キュー 4 6 0 の末尾に新たに単位カゴを設け、その単位カゴに入力した。

【 0 4 6 5 】

さらに、図 3 7 に出力ポート # N に対して単位カゴの転送を K M から指示された場合の動作を示す。ポート # N の先頭単位カゴの単位カゴセル数をカゴセル数に加算し、単位カゴセル数ゼロにする。さらにポート # N のカゴセル数から 1 を減算して出力許可済みセル数に 1 を加算する。K M への転送セル数は出力ポート # N に関して 1 である。出力ポート表の先頭単位カゴのポインタを V C 間公平キュー 4 6 0 の次の単位カゴへひとつ進める。

【 0 4 6 6 】

V C 間公平キュー 4 6 0 の先頭のセルの出力ポートが # 1 と # N であり、その # 1 と # N の出力許可済みセル数がそれぞれ 1 以上（具体的には 2 と 1 ）になったため、クラス多重化 F I F O 4 3 6 が許可すればそのセルを出力できる。

【 0 4 6 7 】

図 3 8 にセルを出力した場合の状態を示す。出力セルの出力ポートに対応する出力許可済みセル数を 1 ずつ減算する。

【 0 4 6 8 】

図 3 9 に、以上のようにして連続して単位カゴ、カゴの転送が指示され、合計 6 セル出力したときの状態を示す。V C 間公平キュー 4 6 0 の先頭の単位カゴはこの時点で空になり、その単位カゴおよび単位カゴセル数は消滅する。必然的に出力ポート表の先頭単位カゴがこの消失した単位カゴを指していることはない。次に、図 4 0 を参照して、出力ポート選択スケジューラ 4 3 2 について説明する。

【 0 4 6 9 】

出力ポート選択スケジューラ 4 3 2 に入力する、出力ポート毎のユニキャスト出力待ちセルとマルチキャスト出力待ちセルをセル選択部 4 7 0 が選択して出力ポート選択スケジューラ 4 3 2 から出力する。セル選択部 4 7 0 の選択の候補となるセルは、出力待ちセルのうち許可されたセルである。

【 0 4 7 0 】

ユニキャストセルは、ユニキャストセル出力許可信号生成部 4 7 1 が出力するユニキャストセル出力許可信号により出力が許可される。

【 0 4 7 1 】

ユニキャストセル出力許可信号生成部 4 7 1 には、バックプレッシャ信号と非常用バックプレッシャ信号が入力する。

【 0 4 7 2 】

バックプレッシャ信号は出力ポート毎である。また非常用バックプレッシャ信号は対応するバックプレッシャ信号が予め定まっている。図 2 5 の場合は、非常用バックプレッシャ信号は S E 2 から出力しており、その S E 2 の出力リンクに接続する出力バッファ O B からのバックプレッシャ信号に対応している。

【 0 4 7 3 】

バックプレッシャ信号が出力禁止を示している時、ユニキャストセル出力許可信号はその出力ポートのセルの出力を許可しない。非常用バックプレッシャ信号が出力禁止を示している場合は、その非常用バックプレッシャ信号に対応するバックプレッシャ信号が全て出力禁止を示しているものと見なして出力を許可しない。

【 0 4 7 4 】

マルチキャストセルは、マルチキャストセル出力許可信号生成部 4 7 2 が出力するマルチキャストセル出力許可信号により出力が許可される。

【 0 4 7 5 】

10

20

30

40

50

マルチキャストセル出力許可信号生成部 4 7 2 には、ユニキャストセル出力許可信号とマルチキャストセル用クラス多重化 F I F O 4 3 6 の先頭のセルの出力ポート情報が入力する。

【 0 4 7 6 】

マルチキャストセル出力許可信号生成部 4 7 2 は、マルチキャストセルの出力ポートが全て（同時にでなくてもよい）許可されたかどうかを監視する。どれかの出力ポートが許可されていない場合はマルチキャストセル出力許可信号を出力禁止とする。マルチキャストセルの全ての出力ポートが許可された時、最後に許可された出力ポートのユニキャストセル出力許可信号が出力可能を示している時間はマルチキャストセルの出力を許可する。その時間内にそのマルチキャストセルがセル選択部 4 7 0 により選択されなかった場合、およびマルチキャストセルがセル選択部 4 7 0 により選択され出力して新しいマルチキャストセルの出力ポート情報が到着した場合は、再び最初から出力ポートが全て許可されたかどうかを監視する。

10

【 0 4 7 7 】

セル選択部 4 7 0 は、こうして出力が許可された出力待ちセルを公平に選択する。

【 0 4 7 8 】

上述したように、マルチキャストセルはその出力ポートの全てのバックプレッシャ信号が同時に出力を許可していなくても出力される。すなわち、そのマルチキャストセルの出力時には、そのセルの出力ポートのバックプレッシャ信号のいくつかが出来禁止を示している可能性がある。この様な状況においても S E 2 のバッファが溢れないように低優先セル非常用バックプレッシャ信号が有効に作用する。

20

【 0 4 7 9 】

上述の構成においては、低優先セル非常用バックプレッシャ信号は出力禁止を示した場合、対応する出力ポート行きのユニキャストセルおよびマルチキャストセルの両方の出力を禁止する。他の構成としては、マルチキャストセル出力許可信号生成部 4 7 2 において、ある出力ポートが一度許可されても低優先セル非常用バックプレッシャ信号が出来禁止を指示した時に許可を取り消す処理を行うなどとしてもよい。また、簡単には、低優先非常用バックプレッシャ信号が出来禁止と指示した場合は、全ての低優先マルチキャストセルの出力を禁止する構成でもよい。

【 0 4 8 0 】

本発明の入力バッファ I B は、図 1 5 で説明したバックプレッシャ考慮キュー監視を行なう。この監視は、バックプレッシャ信号によりバックプレッシャ識別子を更新し、その値に基づいて動作する。ユニキャストセルの出力ポート毎のバックプレッシャ識別子は出力ポート選択スケジューラ 4 3 2 のユニキャストセル出力許可信号により更新し、マルチキャストセルのバックプレッシャ識別子は出力ポート選択スケジューラ 4 3 2 のマルチキャストセル出力許可信号により更新する。

30

【 0 4 8 1 】

次に、図 4 1 を参照して、図 2 5 に示したカゴスイッチに用いるカゴ管理部 K M について説明する。

【 0 4 8 2 】

図 4 1 は、カゴ管理部 K M の構成例を示したものである。

40

【 0 4 8 3 】

K M は出力ポート毎に、

- ・ A B R 入力セル数合計。

【 0 4 8 4 】

- ・ クラス毎に、前カゴセル数合計、カゴセル数合計、転送セル数。

【 0 4 8 5 】

- ・ クラス多重化 F I F O セル数合計。

【 0 4 8 6 】

といった情報を保持し、さらに、出力ポート毎に、転送指示履歴管理部 4 8 0、クラス選

50

択スケジューラ 481 を具備している。

【0487】

ABR入力セル数合計は、ABRクラスについて各出力ポート行きの全IBの合計入力セル数である。一定時間内で積算し、その周期毎にその出力ポートとペアとなるIBへ伝送する。所定の一定時間に達する前にABR入力セル数合計があるしきい値以上になった場合は割り込み的にIBへその旨を通知する構成でもかまわない。

【0488】

前カゴセル数合計は、IBへ入力したセルの有無とクラスの情報より、前カゴに入力したセルの合計を加算する。またKMが転送を指示したクラスの前カゴセル数を全てのIBについて合計したもので更新する。そのクラスはKMの転送指示履歴管理部480が記憶している。

10

【0489】

カゴセル数合計は、KMが転送を指示したクラスのカゴセル数を全てのIBについて合計したもので更新する。転送を指示したクラスはKMの転送指示履歴管理部480が記憶している。

【0490】

転送セル数は、KMが転送を指示したクラスのクラス多重化FIFOへの転送セル数を全てのIBについて合計したもので更新する。転送を指示したクラスはKMの転送指示履歴管理部480が記憶している。

【0491】

クラス多重化FIFOセル数合計は、クラス多重化FIFOセル数を全てのIBについて合計したもので更新する。

20

【0492】

クラス選択スケジューラ381は、前カゴセル数合計、カゴセル数合計、転送セル数合計、クラス多重化FIFOセル数合計などの情報よりIBに転送指示を行なう。過去の転送指示の履歴は転送指示履歴管理部480が記憶し、またクラス選択スケジューラ481はその履歴を利用して上記のセル数情報を修正する。

【0493】

図42、図43、図44に、カゴ管理部KMのクラス選択スケジューラ481で実行するスケジューリングアルゴリズムの一例を示す。このアルゴリズムは、クラスの最大帯域を制限し最小帯域を保証しつつ優先度による優先制御を行なうものである。

30

【0494】

グローバルな変数として、現在時刻を示す変数nowがある。各クラスiについて、最小帯域 $1/D_i$ 、出力義務時刻 T_{di} 、 T_{di} の上限を与えるパラメタ D_{xi} 、および最大帯域 $1/U_i$ 、出力禁止時刻 T_{ui} 、 T_{ui} の下限を与えるパラメタ U_{xi} を設定する(図42のステップS100)。

【0495】

各クラスiについて、出力すべき場合は現在時刻nowが $T_{di} \leq now$ である時であり、出力が禁止される場合は現在時刻nowが $now < T_{ui}$ である時である。また優先順位に従って出力すべき場合は現在時刻nowが $T_{ui} \leq now < T_{di}$ である時である。 T_{di} はその上限値を定められており、 $T_{di} \leq now + D_{xi}$ である。また T_{ui} はその下限値を定められており、 $now - U_{xi} \leq T_{ui}$ である。ただし $U_{xi} \geq 0$ 、 $D_{xi} \geq D_i$ である。

40

【0496】

初期設定は、nowをゼロに設定し、各クラスiに対して T_{di} を $D_i - U_i$ とし、 T_{ui} をnowつまりゼロにする(図42のステップS101)。

【0497】

このアルゴリズムはセル周期毎に次の処理を行なう。

【0498】

まず、クラスがアクティブになった時の処理を行なう。具体的にはセルが無かったクラス

50

i にセルが到着したとき、 T_{di} を now と T_{di} の大きい方に更新し、 T_{ui} を now と T_{ui} の大きい方に更新する (図 4 3 のステップ S 1 0 2)。

【0499】

次にセルを転送してもよいかどうかを判定して (図 4 3 のステップ S 1 0 3)、もしそうであれば転送クラス選択および転送クラス指示を行なう (図 4 3 のステップ S 1 0 4)。

【0500】

転送クラス選択は次のような処理である。

【0501】

転送すべきセルがある各クラス i の中で、

1. $T_{di} < now$ となっているクラスがあれば、その中で最も $now - T_{di}$ の大きいクラスを選択する (同じ値の場合は優先度で判断)。 10

【0502】

2. $T_{di} \leq now$ となっているクラスがなければ、 $T_{ui} \leq now < T_{di}$ を満足するクラスの中から優先度の最も高いクラスを選択する。

【0503】

選択されたクラスを j とすると、図 4 4 のステップ S 1 0 5 に進み、次のような転送クラス指示の処理を行なう。

【0504】

1. クラス j の転送指示を行なう。

【0505】

2. T_{uj} を $\max(T_{uj}, now - U_{xj})$ と更新。 20

【0506】

3. T_{dj} を $T_{dj} + D_j$ と更新。

【0507】

4. T_{uj} を $T_{uj} + U_j$ と更新。

【0508】

5. T_{dj} を $\min(T_{dj}, now + D_{xj})$ と更新。

【0509】

最後に、転送セル流量モニタの処理を行なった後 (図 4 4 のステップ S 1 0 6)、 now をインクリメントして (図 4 4 のステップ S 1 0 7)、セル周期毎の処理を終了する。 30

【0510】

図 4 4 のステップ S 1 0 6 の転送セル流量モニタの処理は、クラス j の転送セル数が x_j の時、クラス j に関して次の処理を行なう。

【0511】

1. T_{uj} を $\max(T_{uj}, now - U_{xj})$ と更新。

【0512】

2. T_{dj} を $T_{dj} + (x_j - 1) D_j$ と更新。

【0513】

3. T_{uj} を $T_{uj} + (x_j - 1) U_j$ と更新。

【0514】

4. T_{dj} を $\min(T_{dj}, now + D_{xj})$ と更新。 40

【0515】

図 4 3 のステップ S 1 0 3 のセルを転送指示してもよいかどうかの判定は、クラス多重化 F I F O セル数合計が、あるしきい値以下かどうかで行なう。あるしきい値とは通常は「1」である。このクラス多重化 F I F O セル数合計は転送指示履歴管理部 4 8 0 の転送指示履歴により修正された値である。

【0516】

図 4 3 のステップ S 1 0 4 における、クラスに転送すべきセルがあるかどうかの判定は、前カゴセル数合計またはカゴセル数合計のどちらか一方がゼロでないことを判定して行なう。この前カゴセル数合計とカゴセル数合計は転送指示履歴管理部 4 8 0 の転送指示履歴 50

により修正された値である。

【0517】

転送指示はカゴ440からクラス多重化FIFO434へのセルの転送指示（図31の転送指示T2）である。その際に前カゴからカゴへの単位カゴの転送指示を同時に行なうかどうかの判定は、前カゴセル数合計がゼロでなく、かつ、カゴセル数合計があるしきい値以下かどうかを判定して行なう。あるしきい値とは通常は1である。この前カゴセル数合計とカゴセル数合計は転送指示履歴管理部480の転送指示履歴により修正された値である。

【0518】

本アルゴリズムは、 D_i を大きな数に設定することにより、最小保証帯域をゼロに近付けることができる。例えば全てのクラスの D_i を大きな数に設定すれば全てのクラスを完全優先で優先制御することができる。

10

【0519】

また、クラスのひとつを、優先度は最も低いが最小保証帯域を持つIC（Idle Cell）クラスとすることにより、出力ポートの使用率を制御することができる。IBはICクラスの転送指示があった場合にアイドルセルをクラス多重化FIFOへ挿入する（このIBはひとつだけでよい）。

【0520】

アイドルセルは他のセルと同様にセルスイッチの出力ポートに向けて交換され、そこでUnassigned Cellに置き換えられる。KMは、ICクラスはいつも転送すべきセルがある状態と考えスケジューリングを行なえば良い。

20

【0521】

ところで、図42～図44に示したアルゴリズムでは説明を簡単にするために、変数now、Tdi、Tuiといった時刻を示す変数は、表現できる値が上限を持たないとしている。実装時にはこれらの変数は有限のビット長を持つレジスタなどで表現するため長時間経過すると同じ時刻が循環して現われる。アルゴリズムを正常に動作させるためには次のようなことを行なえば良い。これらの周期をMaxとする。Maxは十分大きい必要がある。

【0522】

このとき、 $x := f(x)$ などの計算は、

$$x := f(x) \bmod \text{Max}$$

と置き換える。つまりMaxでモジュロをとる。また、 $y \leq z$ の比較演算は、

$$(z \leq y - \text{Max}/2) \text{ or } ((y \leq z) \text{ and } (z \leq y + \text{Max}/2))$$

に置き換え比較する。

30

【0523】

例えば、 $\text{now} := \text{now} + 1$ という計算は、

$$\text{now} := (\text{now} + 1) \bmod \text{Max}$$

と置き換える。Tdi \leq nowという比較演算は、

$$(\text{now} \leq \text{Tdi} - \text{Max}/2) \text{ or } ((\text{Tdi} \leq \text{now}) \text{ and } (\text{now} \leq \text{Tdi} + \text{Max}/2))$$

に置き換え比較すれば良い。

40

【0524】

時刻の変数は周期的に監視し、必要に応じてnowからMax/2以上離れ過ぎないように調整を行なうことが必要であるが、Maxを大きくとることによって監視の周期を長くとることができ、実装が容易となる。

【0525】

さらに、他のクラス選択スケジューリングとしては次のようなものを有効である。

【0526】

すでに説明したように、CG毎（VC毎）に設定された重みに従ってVC間公平キューイングを行ったセルバッファ装置は、入力バッファIBのCGテーブル（図30）またはVCテーブルのNxに重みを書き込む。

【0527】

50

ここでは、クラス選択スケジューリングを併用することによって、そのクラスに設定されているVC毎の重みを係数倍できることを説明する。

【0528】

クラスA、クラスBがあり、クラスAに設定されているあるVCの重み N_{xa} とし、クラスBに設定されているあるVCの重みを N_{xb} とする。このとき双方のクラスに出力待ちセルが存在するときにクラス選択スケジューリングにより単位カゴの転送指示回数をクラスXとクラスYで R_a 対 R_b の比になるようにスケジューリングすることにより、それぞれの重みを $N_{xa} \times R_a$ 、 $N_{xb} \times R_b$ と係数倍することが可能である。3クラス以上の場合も同様である。所望の比になるように、スケジューリングを行うアルゴリズムとして例えば、SCFQ (Self-Clocked Fair Queueing) などが知られている。

10

【0529】

具体的な例で説明する。例えば、VC1からVC7の7本のVCに重みをそれぞれ1.24、3.72、2.18、934、562、1370、4360に設定したいとする。クラスAにVC1、2を設定し、 N_x をそれぞれ1.24、3.72とする。クラスBにVC3を設定し、 N_x を2.18とする。クラスCにVC4、5を設定し、 N_x をそれぞれ9.34、5.62とする。クラスDにVC6、7を設定し、 N_x をそれぞれ1.37、4.36とする。単位カゴの転送指示回数を出力待ちセルが存在するクラス間でクラスA：クラスB：クラスC：クラスD = 1：1：100：1000の比になるようにクラス選択スケジューリングを行えば、目的の重み1.24：3.72：2.18：934：562：1370：4360でセルが出力ポートに多重化される。この方法は N_x の値として大きな数値を設定するとVC間公平キューイングのセル遅延揺らぎが増大してしまうという現象を回避する効果がある。

20

【0530】

図42～図44に示したようなクラス選択スケジューリングのアルゴリズムでは、各クラスの転送セル流量をモニタしていた。これを利用して階層的なクラス選択スケジューリングを行うことが可能である。

【0531】

クラスA1、A2、A3、...からなるクラス群 G_a と、別のクラスB1、B2、B3、...からなるクラス群 G_b とがあるとき、出力待ちセルが存在するクラス群間で転送セル流量の比を R_a 対 R_b と設定してクラス選択スケジューリングを行う。クラス群が3以上の場合も同様である。各クラス群内のクラス選択スケジューリングと、クラス群間の選択スケジューリングとは独立である。

30

【0532】

これにより、あるクラス群間の剰余帯域は、そのクラス群内のクラスのみが使用することになり、別のクラス群が使用できることを防止できる。別のクラス群がこの剰余帯域を使用できるのは、あるクラス群の出力待ちセルが一つも存在しない場合だけである。このような階層的なクラス選択スケジューリングは、例えば、企業毎にクラス群を割り当てるなどの応用が考えられる。

【0533】

5.2 優先制御およびVC間公平キューイングを行うセルスイッチのセル数修正方法次に、優先制御およびVC間公平キューイングを行うセルスイッチ(カゴスイッチ)(図25、図29、図31)のカゴ管理部KMにおけるセル数の修正方法を図45を参照して説明する。

40

【0534】

修正の原理は、図19、図20などを用いて説明したカゴスイッチの場合も同じ方法でセル数を修正することにより優先制御の性能や出力ポートのスループットの向上などを実現することができる。

【0535】

本発明のセル数修正方法の一応用例として、図45に示すカゴスイッチに適用する。カゴ

50

スイッチでは、管理部 301 がクラス i の前カゴ (B_{pi}) へ転送を指示すると、そのクラスの前カゴのセルがカゴ (B_{ki}) へ転送されるとともに、その一部がクラス多重化 FIFO (B_m) へ転送される。管理部 301 がクラス i のカゴからクラス多重化 FIFO へ転送を指示すると、そのクラスのカゴのセルがクラス多重化 FIFO へ転送される。

【0536】

図45の様に各セルバッファのセル数情報および遅延時間を定める。

【0537】

H_{pi} は転送指示 I_{pi} により前カゴからカゴへ転送される HOL 数、 H_{pmi} は同じく転送指示 I_{pi} により前カゴからクラス多重化 FIFO へ転送される HOL セル数である。

【0538】

ここで一例として、各遅延時間の値が等しい場合について考える。つまり、 $D_{Op} = D_{Ik} = D_{Ok} = D_{Im} = D_{ipm} = D$ 、および、 $D_{Tpm} = D_{Tkm}$ とする。また $p = [-D, 0)$ とする。このとき、クラス i の前カゴとカゴの修正されたセル数 (N_{pi} 、 N_{ki})、およびクラス多重化 FIFO の修正されたセル数 (N_m) は次のように計算できる。

【0539】

【数7】

$$N'_{pi} = N_{pi} - f(S_{m_{pi}}(p), H_{pi}) \quad \dots (22)$$

$$\begin{aligned} N'_{ki} = & N_{ki} \\ & + f(S_{m_{pi}}(p), H_{pi}) \\ & - f(S_{m_{pi}}(p), H_{pmi}) \\ & - f(S_{m_{ki}}(p), H_{ki}) \quad \dots (23) \end{aligned}$$

$$\begin{aligned} N'_m = & N_m \\ & + \sum_i f(S_{m_{pi}}(p), H_{pmi}) \\ & + \sum_i f(S_{m_{ki}}(p), H_{ki}) \\ & - D \quad \dots (24) \end{aligned}$$

【0540】

5.3 第10の実施形態に係るセルスイッチの説明の続き

図46に、カゴスイッチの入力バッファ IB とカゴ管理部 KM との接続例を示す。

【0541】

ここでは、例えば、クラス数を 62 クラス、セルスイッチサイズを 16 入力 16 出力とする。

【0542】

クラス番号は、{ 0、1、2、...、60、61、UC、IC } と 6 ビットにコーディングする。

【0543】

クラス UC は非割り当てセル (Unassigned Cell)、IC はアイドルセル (Idle Cell) を示し、これらの使い方は以下に記述する。

【0544】

前カゴセル数、カゴセル数、クラス多重化 FIFO セル数は、{ 0、1、2、...、13、

10

20

30

40

50

14、15以上}と4ビットにコーディングする。

【0545】

入力バッファIBからカゴ管理部KMへの転送情報を以下に示す。

【0546】

(1) 入力セル 出力ポート毎に、入力したセルについて(ユニキャストの場合はひとつの出力ポートについてのみ出力する)、

・入力セルの有無とクラス... 6ビット

(入力セルがない場合はクラスUCを送信する)

(3) 転送指示クラス 出力ポート毎に転送の指示クラスについて、

・転送後の前カゴセル数... 4ビット

・転送後のカゴセル数... 4ビット

・クラス多重化FIFOへの転送セル数... 2ビット

(最大でユニキャストとマルチキャストの2セルが転送される)

(4) クラス多重化FIFO 出力ポート毎にユニキャストのみの、

・クラス多重化FIFOセル数... 4ビット

合計すると、ひとつの出力ポートに関する情報は20ビットである。

【0547】

カゴ管理部KMから全ての入力バッファIBへ同報する転送情報を以下に示す。

【0548】

(2) 出力ポート毎に、

・前カゴからカゴへ単位カゴの転送指示(下記のクラスについて)... 1ビット

・カゴからクラス多重化FIFOへセルの転送指示と、指示クラス... 6ビット

(転送指示をしない場合はUCクラスを指定。ICクラスが指定された場合はIBはアイドルセルを転送(挿入)する。アイドルセルは他のクラスのセルと同様に出力ポートへ交換転送され、セルスイッチからの出力時にUnassigned Cellに置き換えられるスループット調整用のセルである。)

合計すると、ひとつの出力ポートに関する情報は7ビットである。

【0549】

KMは、ABRトラヒック制御のため出力ポート毎に、対応するIBへ次の信号を送信する。

【0550】

(5) 全IBの入力セルのうち、各出力ポート行きのABR入力セル数... 5ビット

(16入力16出力のセルスイッチの場合)

ひとつの出力ポートに関する情報は5ビットである。

【0551】

図46は、カゴ管理部KMをふたつのLSIにて実現する構成例を示したものである。ふたつのLSIのうち一方は出力ポート#1から#8を管理し、もう一方は出力ポート#9から#16を管理する。

【0552】

IBからKMへの転送情報(前述の(1)(3)(4))は、1出力ポートあたり20ビットであるから、8ポート分で160ビットである。高速差動伝送を2ペア用いて平行伝送を行なえば1セル周期に80ビットの伝送速度でよい。

【0553】

KMからIBへの同報情報(前述の(2))は、1出力ポートあたり7ビットであるから、8ポート分で56ビットである。通常の伝送ドライバを用いて、8本の平行伝送を行なえば1セル周期に7ビットの伝送速度でよい。

【0554】

KMからIBへのABRトラヒック制御のための情報(前述の(5))は、1出力ポートあたり5ビットである。

【0555】

10

20

30

40

50

セル伝送速度が622Mbpsであるセルスイッチを考えた場合、これらの伝送速度は現在の技術において十分実現可能である。

【0556】

IBはKMへセル数情報をLSB（最下位桁）から伝送することにより、セル数の加算の回路規模を少なくすることができる。加算回路は全てのIBから伝送されてくるセル数情報を（出力ポート毎クラス毎に）下位桁から順に加算することができるためである。より具体的には、加算回路は入力した桁の数字とそれまでに計算した下位桁からの桁上がりの和を計算して、加算結果のその桁の値を得ることができる。加えて次の桁の計算のための桁上がりを持する。

【0557】

IBからKMへの信号の中にクラス多重化FIFOセル数がある。バックプレッシャ付きセルスイッチが出力ポート毎の蓄積セル数（図25の出力バッファOBのセル数）を出力できる場合には、KMは、クラス多重化FIFOセル数の代わりにこの出力ポート毎蓄積セル数を用いても良い。

【0558】

これは、図2のクラス間の優先制御を行なうセル多重化装置や、図7のVC間公平キューイングを行なうセル多重化装置と同じ効果を持つ。

【0559】

図47に、ABRサービスを実現するための情報の経路を示す。

【0560】

図47は、ポート#1からポート#Nに入力するABRクラスのセルが、バックプレッシャ付きセルスイッチ400によりポート#iに交換多重化される部分について特に抜き出して示してある。実際は全ての出力ポートに関して、この図の#iと同じ構成になっている。

【0561】

ABRサービスは、セルスイッチの出力ポートに使用可能な帯域があれば、その帯域をABRクラスのVCで公平に分配して使用するサービスクラスである。トラヒック制御には受信端末から送信端末へ伝送されるRMセルを用いる。

【0562】

セルスイッチ400は、出力ポート#iにおいてABRが使用可能な帯域を知っており、それと現在のABRの入力トラヒックとの比（または差）を知ることにより、送信端末に対してトラヒック制御を行なうことができる。現在のABRトラヒックは図47に示したポート#i行きの入力セル数の一定周期毎の合計から知ることができる。入力セル数から入力帯域を求め、これが出力ポートにおいてABRが使用可能な帯域を越えていれば過負荷であるから、ポート#iに入力するRMセルにVCの帯域を減らす指示を書き込めばよい。逆に、ポート#i行きの入力セル数から求めた入力帯域が、出力ポートにおいてABRが使用可能な帯域を下回っていれば、低負荷であるから、ポート#iに入力するRMセルにVCの帯域を増やす指示を書き込めばよい。入力帯域と出力帯域の比または差を知ることにより、過負荷、低負荷の程度を知ることができ、細やかな制御を行なうことが可能である。

【0563】

この処理により輻輳に陥る確率は少なくなるが、さらにIBにおいてキュー長を監視して輻輳を検出した場合はEFCIをマークしたり、RMセルを書き換えたりすることにより、より安全にトラヒックを制御することが可能である。特にマルチキャストの場合は、入力帯域が出力帯域より小さくとも過負荷となる可能性があり、その場合はキュー長監視によるトラヒック制御が重要になる。

【0564】

このように、第10の実施形態に係るカゴスイッチは、VC間公平キューイングを行なうため、キュー長監視によるトラヒック制御においてもVC間で公平な帯域分配を実現することが可能である。

10

20

30

40

50

【 0 5 6 5 】

【 発明の効果 】

以上説明したように、本発明のクラス間の優先制御を行なうパケット転送装置（第1、第2の実施形態（請求項1、2、3、4、））によれば、制御部（クラス管理部）が複数の入力バッファ全体の状況を把握して（少なくとも全体の状況を把握していればよいが、個々の入力バッファ毎に把握していても構わない）指示を出すため、各入力バッファのクラス毎の蓄積状況の違いによりクラス間の優先関係が崩れてしまうことが生じず、また、入力ポートの数すなわち入力バッファの数が増えても容易にクラス間の優先制御を実現できる。

【 0 5 6 6 】

また、複数の入力バッファにまたがってV C（フロー）間公平キューイングをするために、例えば、制御手段が全ての入力バッファの全てのV C（フロー）について、その入力ポート番号とV C番号、パケットの到着状況、出力状況等を把握してパケットの出力順序を制御するのでは、装置構成が複雑になり、スループットの向上が難しい。そこで、本発明のパケット転送装置（第3、第4、第5の実施形態（請求項5、6））によれば、あるフェーズ（その時間的長さは不定）内で出力されるべきパケット集合（カゴ）という概念を導入し、このパケット集合に入れるパケットの選択と、このパケット集合に入っているパケットのみを出力する（その他のパケットは出力が許可されない）という処理を各入力バッファが分担し、パケット集合に新たなパケットを入れるタイミングを指示する処理をカゴ管理部が分担することにより、例えば、入力バッファ方式のセルスイッチにおいて、V C（フロー）間公平キューイングを容易に実現できる。

【 0 5 6 7 】

また、本発明のパケット転送装置（第6、第7の実施形態（請求項7、8））によれば、制御手段（バッファポインタ管理部）がパケットが入力されたときに複数の集合への振り分けが行われており、パケットを出力するときには1つの集合に属するポインタを出すだけで済むので、例えば設定すべきV C数が増えても高速にV C（フロー）間公平キューイングのための出力すべきパケットの選択が実現できる。

【 図面の簡単な説明 】

【 図 1 】 本発明の第1の実施形態に係るクラス間の優先制御を行なうセル多重化装置の構成例を示した図。

【 図 2 】 本発明の第2の実施形態に係るクラス間の優先制御を行なうセル多重化装置の他の構成例を示した図。

【 図 3 】 出力バッファの他の構成例を示した図。

【 図 4 】 出力バッファのさらに他の構成例を示した図。

【 図 5 】 本発明の第3の実施形態に係るV C毎F I F Oを備えた出力バッファ型セル多重化装置の構成例を示した図。

【 図 6 】 本発明の第4の実施形態に係る：V C間公平キューイングを行なうセル多重化装置の構成例を示した図。

【 図 7 】 本発明の第5の実施形態に係るV C間公平キューイングを行なうセル多重化装置の構成例を示した図。

【 図 8 】 本発明の第6の実施形態に係るセルグループF I F Oを用いるセルバッファ装置の構成例を示した図。

【 図 9 】 本発明の第7の実施形態に係るセルグループF I F Oを用いるセルバッファ装置の構成例を示した図。

【 図 1 0 】 第6の実施形態に係るセルグループF I F Oのポインタチェーンによるデータ構造例を説明するための図。

【 図 1 1 】 第6の実施形態に係るセルグループF I F Oのポインタチェーンによるデータ構造例を説明するための図で、空きセルグループチェーンの構成を示している。

【 図 1 2 】 第7の実施形態に係るセルグループF I F Oのポインタチェーンによるデータ構造例を説明するための図。

10

20

30

40

50

- 【図13】第7の実施形態に係るセルグループFIFOのポインタチェーンによるデータ構造例を説明するための図で、空きセルグループチェーンの構成を示している。
- 【図14】パケットグループを用いたフロー間公平キューイングを説明するための図。
- 【図15】パケットグループを用いたフロー間公平キューイングを用いたパケットバッファ装置の動作を説明するための図。
- 【図16】本発明の第8の実施形態にかかるバックプレッシャ考慮キュー監視を説明する図である。
- 【図17】第8の実施形態に係るバックプレッシャ考慮キュー監視部を持つセルバッファ装置の構成例を示した図。
- 【図18】第8の実施形態に係るバックプレッシャ考慮キュー監視アルゴリズムを示したフローチャート。 10
- 【図19】本発明の第9の実施形態に係る遅延を考慮したセルバッファ装置の構成例を示した図。
- 【図20】遅延を考慮したセル数情報の修正方法の基本原則を説明するための図。
- 【図21】セル数修正方法の他の基本原則を説明するための図。
- 【図22】図21のセル数修正方法に用いる基本的な関数 $f()$ を説明するための図。
- 【図23】図21のセル数修正方法に用いる基本的な関数 $g()$ を説明するための図。
- 【図24】図21のセル数修正方法に用いる基本的な関数 $S_c()$ を説明するための図。
- 【図25】本発明の第10の実施形態に係るカゴスイッチの全体の構成例を示した図。
- 【図26】図25の出力バッファの構成例を示した図。 20
- 【図27】図25の1段目単位スイッチの構成例を示した図。
- 【図28】図25の2段目単位スイッチの構成例を示した図。
- 【図29】図25の入力バッファの構成例を示した図。
- 【図30】入力バッファのキューイング処理を説明するための図。
- 【図31】図25の入力バッファのカゴ管理部とのインタフェースを示した図。
- 【図32】図29のクラス毎マルチキャストセル管理部の動作を説明するための図。
- 【図33】図29のクラス毎マルチキャストセル管理部の動作を説明するための図。
- 【図34】図29のクラス毎マルチキャストセル管理部の動作を説明するための図。
- 【図35】図29のクラス毎マルチキャストセル管理部の動作を説明するための図。
- 【図36】図29のクラス毎マルチキャストセル管理部の動作を説明するための図。 30
- 【図37】図29のクラス毎マルチキャストセル管理部の動作を説明するための図。
- 【図38】図29のクラス毎マルチキャストセル管理部の動作を説明するための図。
- 【図39】図29のクラス毎マルチキャストセル管理部の動作を説明するための図。
- 【図40】図29の出力ポート選択スケジューラの構成例を示した図。
- 【図41】図25のカゴ管理部の構成例を示した図。
- 【図42】図41のカゴ管理部のクラス選択スケジューラで実行する、最大帯域を制限し最小帯域を保証するスケジューリングアルゴリズムの具体例を示したフローチャート。
- 【図43】図41のカゴ管理部のクラス選択スケジューラで実行する、最大帯域を制限し最小帯域を保証するスケジューリングアルゴリズムの具体例を示したフローチャート。
- 【図44】図41のカゴ管理部のクラス選択スケジューラで実行する、最大帯域を制限し最小帯域を保証するスケジューリングアルゴリズムの具体例を示したフローチャート。 40
- 【図45】セル数情報の修正方法の応用例であるカゴスイッチの動作を説明するための図。
- 【図46】図25の入力バッファとカゴ管理部との接続の具体例を示した図。
- 【図47】ABRサービスのための情報の経路の具体例を示した図。
- 【図48】従来のUBRサービスにおける不公平な帯域割り当てについて説明するための図。
- 【図49】従来のクラス間の優先制御を行うセル多重化装置の構成を説明するための図。
- 【図50】従来の入力バッファを持つセル多重化装置の構成を説明するための図。
- 【図51】従来のVCテーブルの検索が必要なセルバッファ装置の構成を説明するための 50

図。

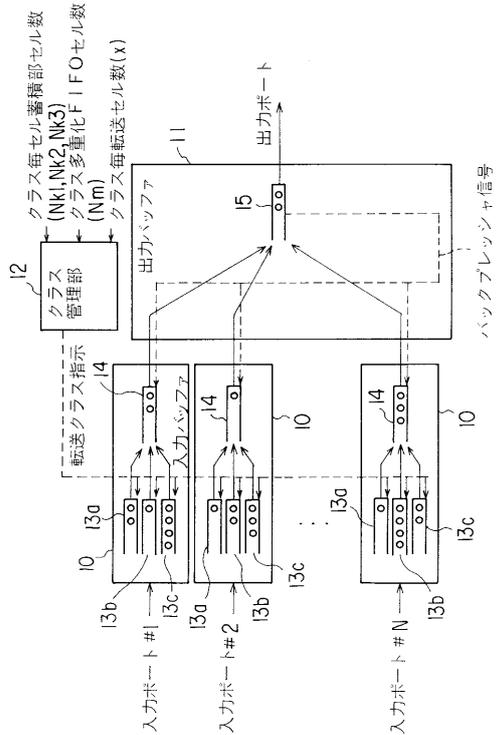
【図52】図50に示すような従来の入力バッファを持つセル多重化装置における入力バッファのキュー長の変化についての説明するための図。

【図53】従来の優先制御を行うセル多重化装置の構成を説明するための図。

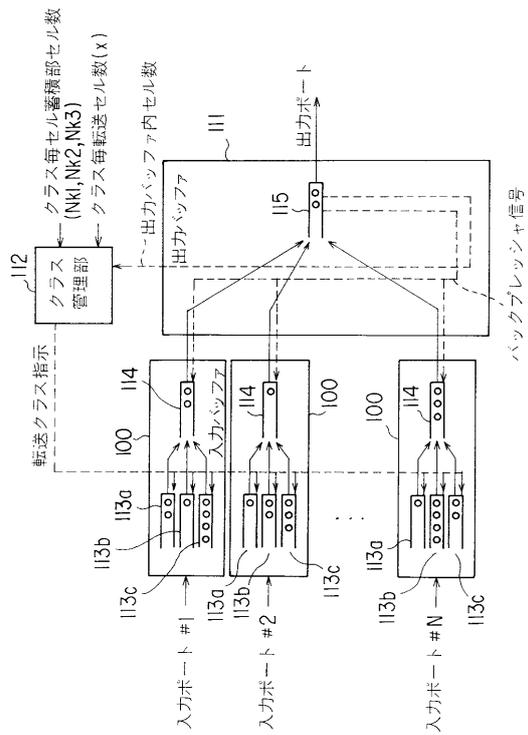
【符号の説明】

10...入力バッファ、11...出力バッファ、12...クラス管理部、212、222...カゴ管理部、230、240...バッファポインタ管理部、282...キュー監視部、302...転送指示履歴、KM...カゴ管理部、400...バックプレッシャ付きセルスイッチ。

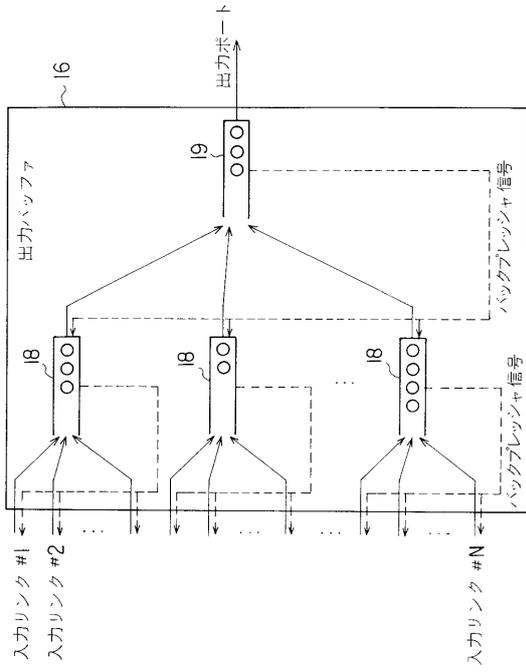
【図1】



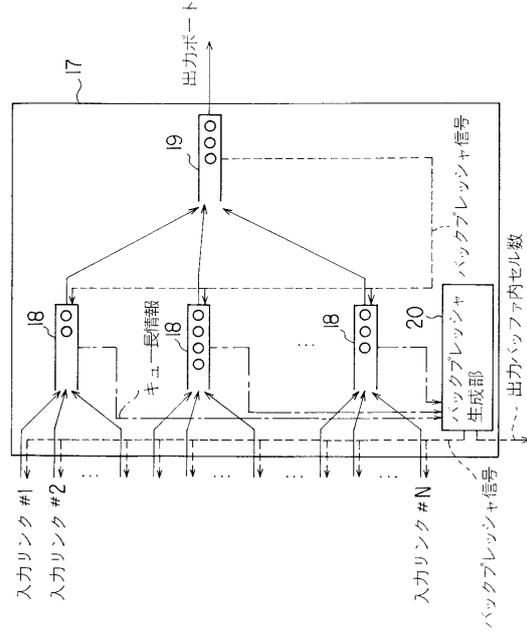
【図2】



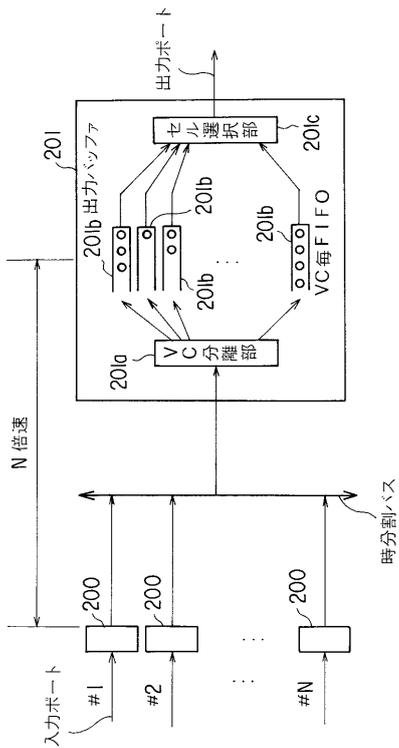
【 図 3 】



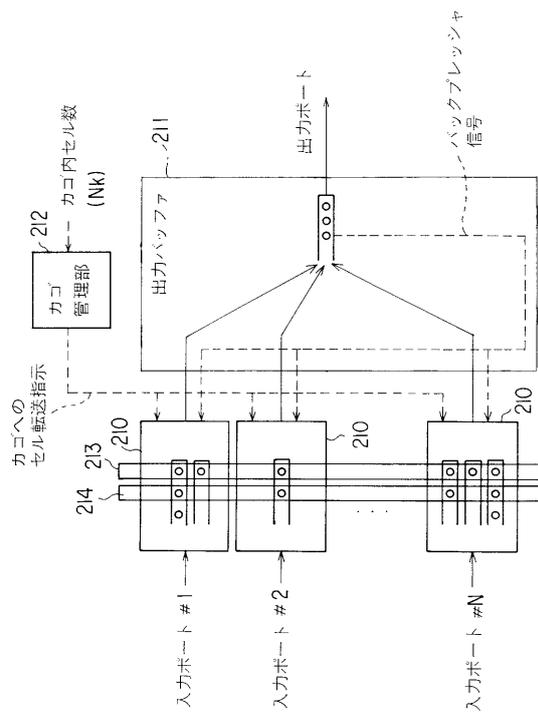
【 図 4 】



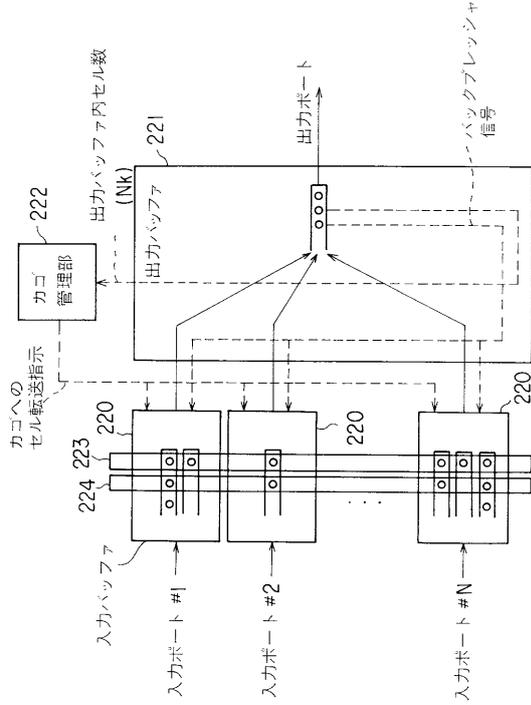
【 図 5 】



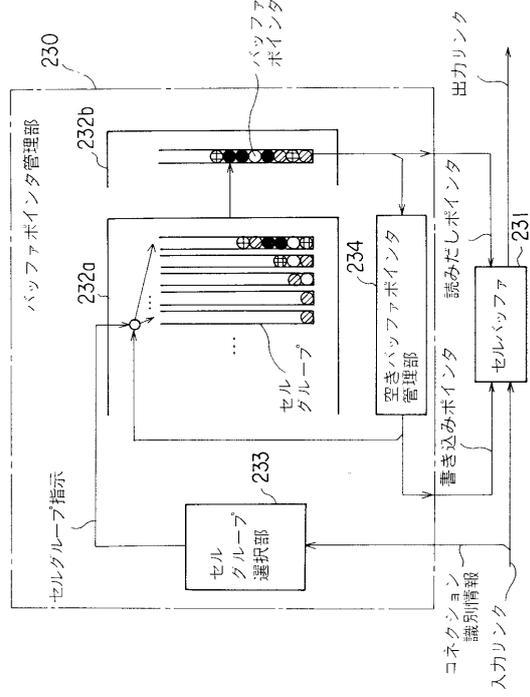
【 図 6 】



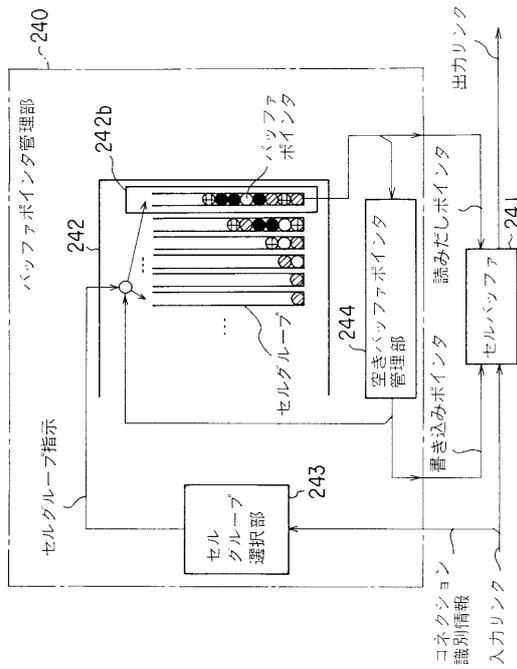
【 図 7 】



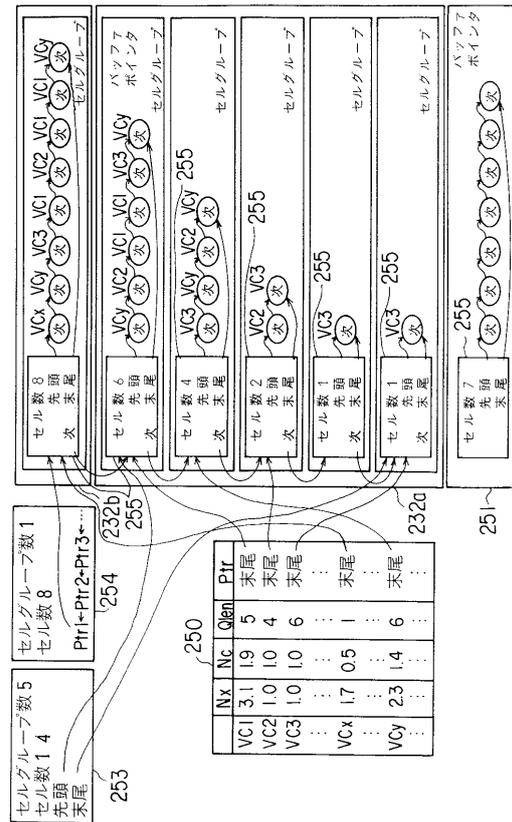
【 図 8 】



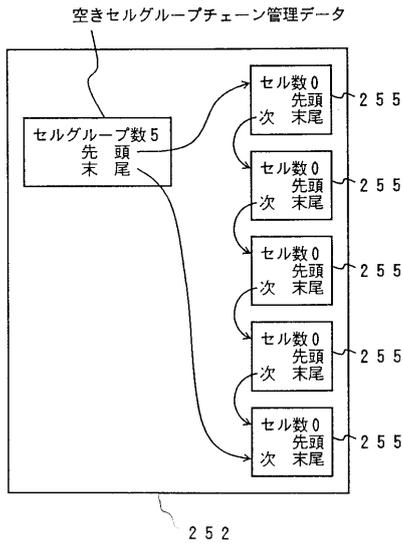
【 図 9 】



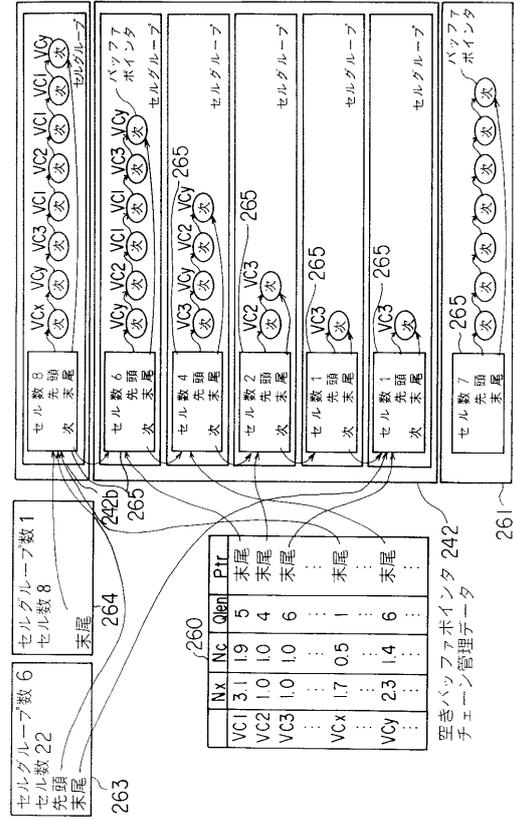
【 図 10 】



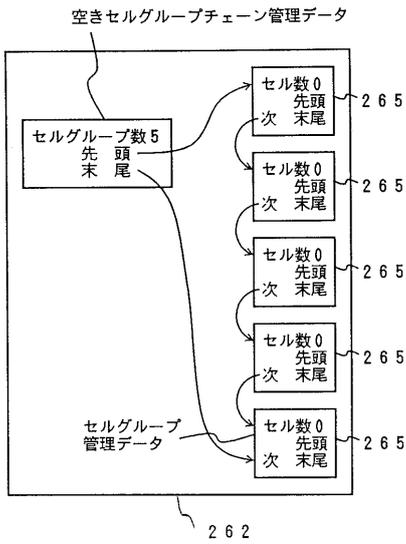
【図11】



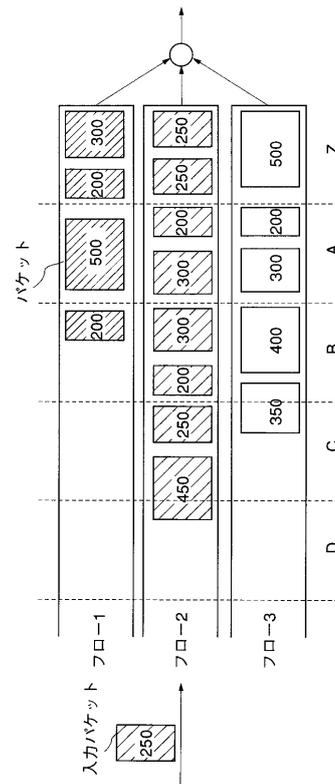
【図12】



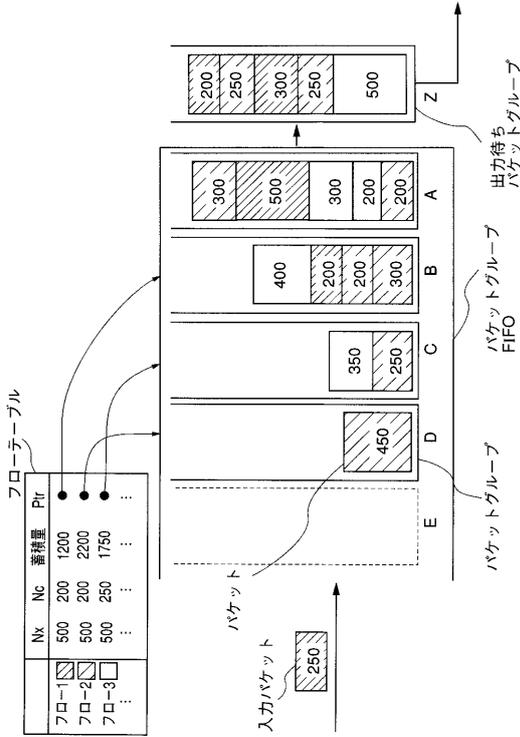
【図13】



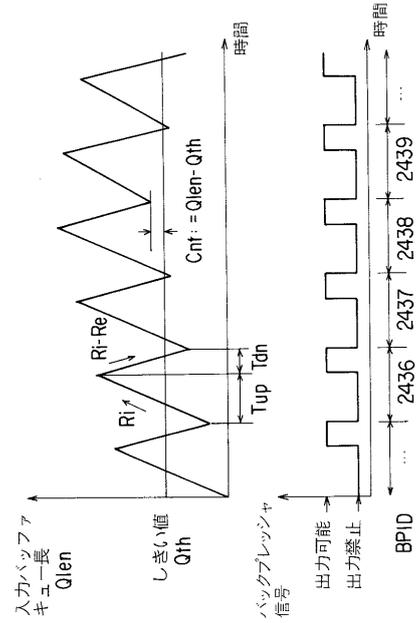
【図14】



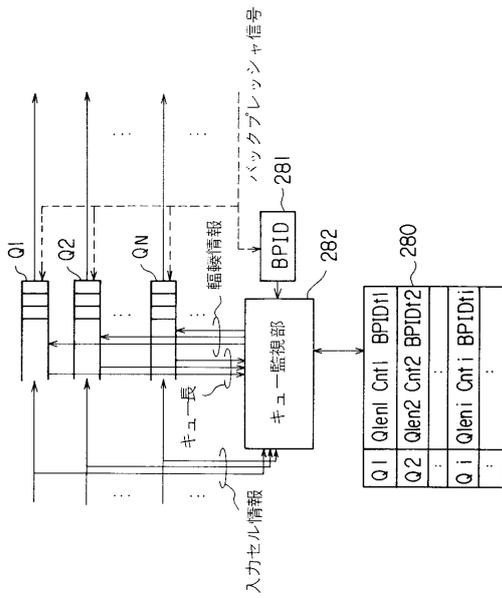
【図15】



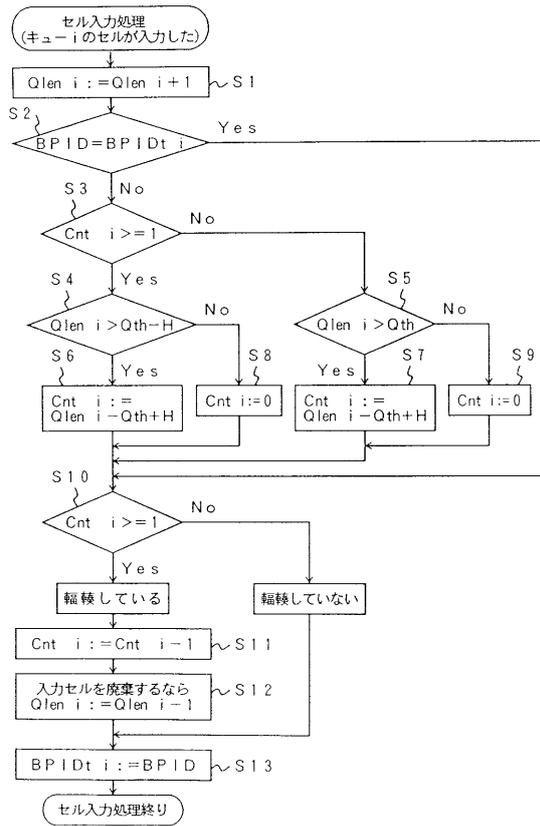
【図16】



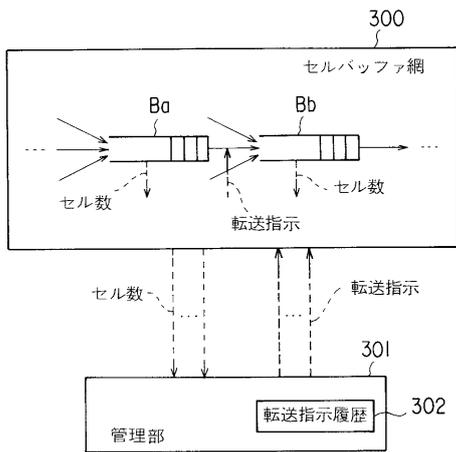
【図17】



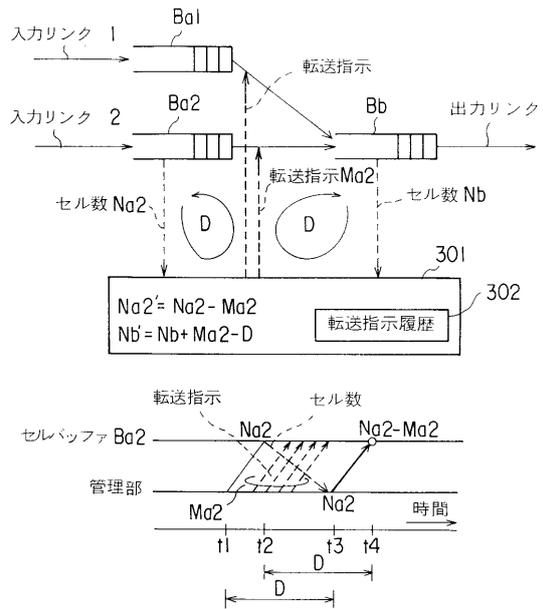
【図18】



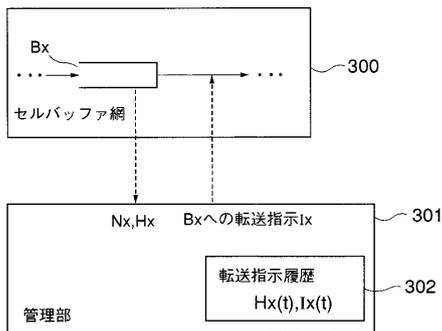
【図19】



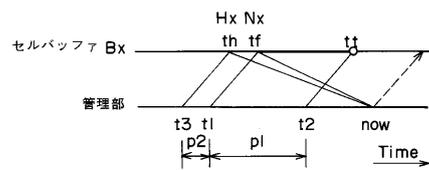
【図20】



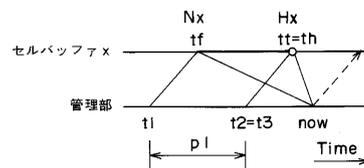
【図21】



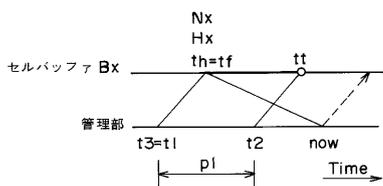
【図23】



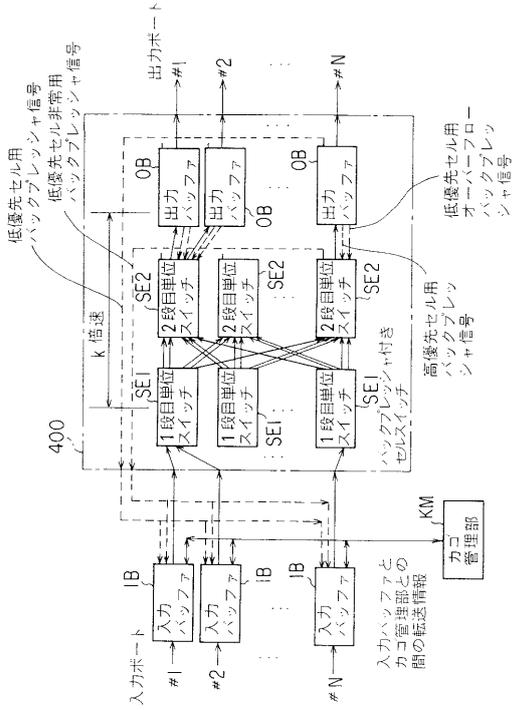
【図24】



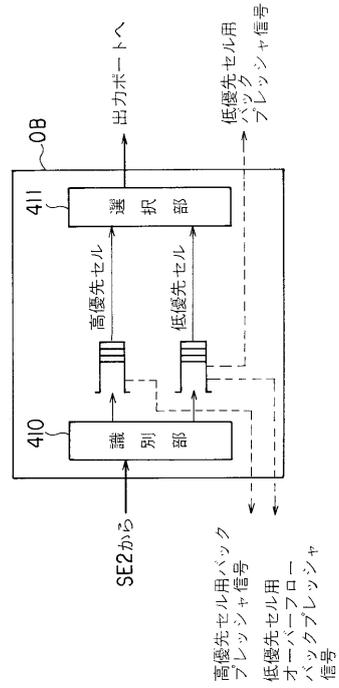
【図22】



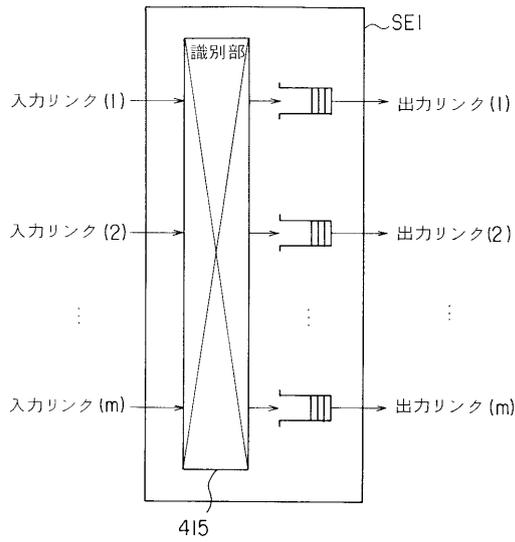
【 図 2 5 】



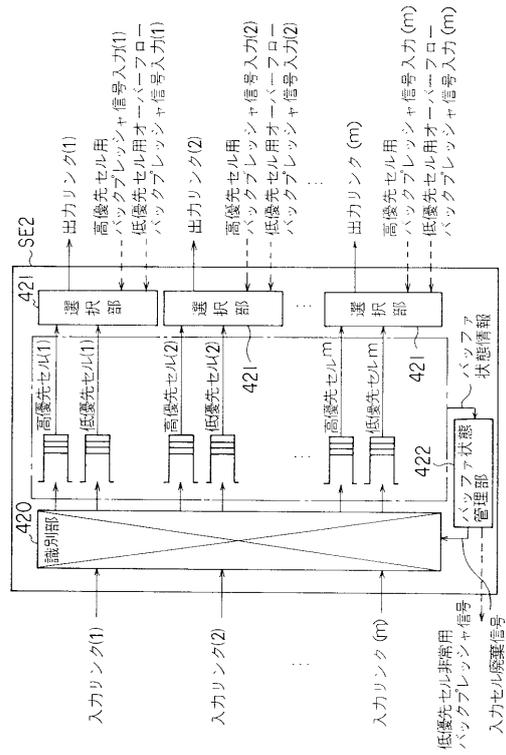
【 図 2 6 】



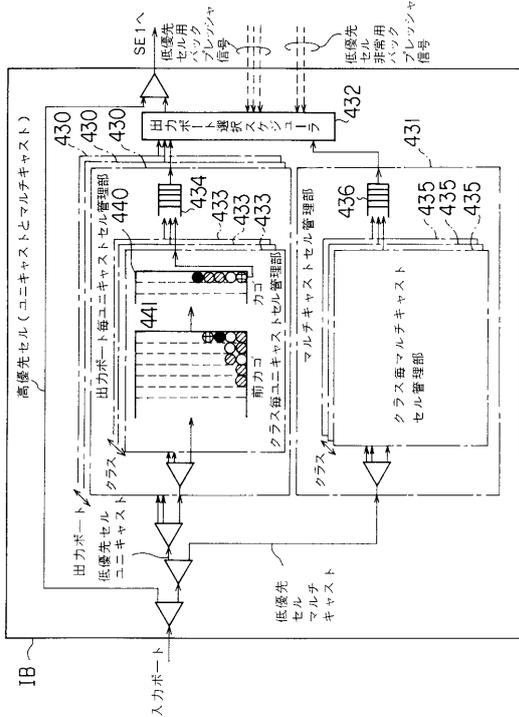
【 図 2 7 】



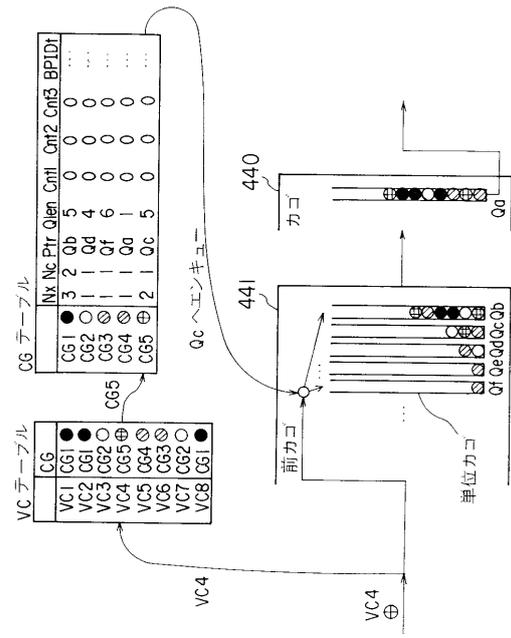
【 図 2 8 】



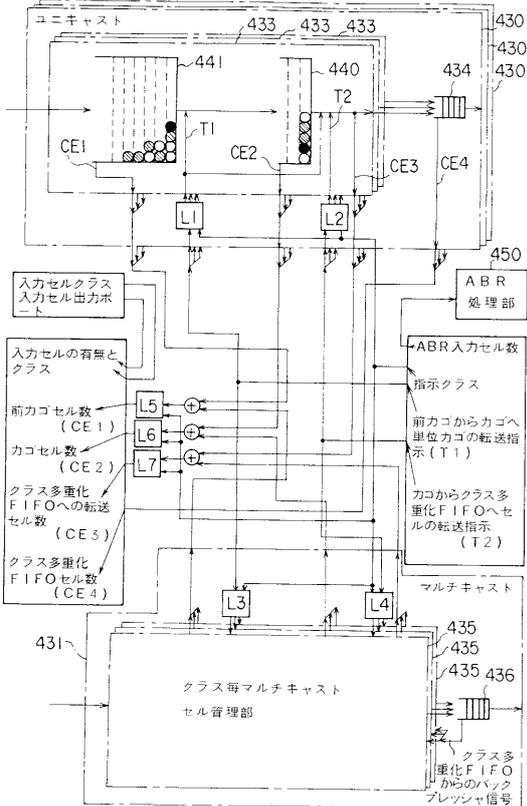
【図 29】



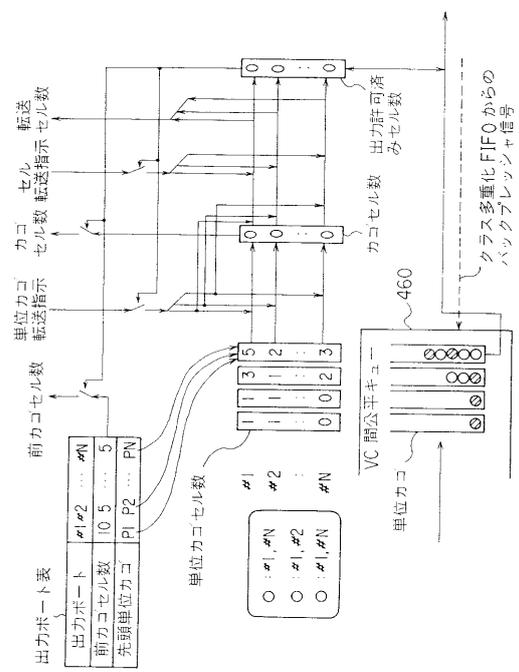
【図 30】



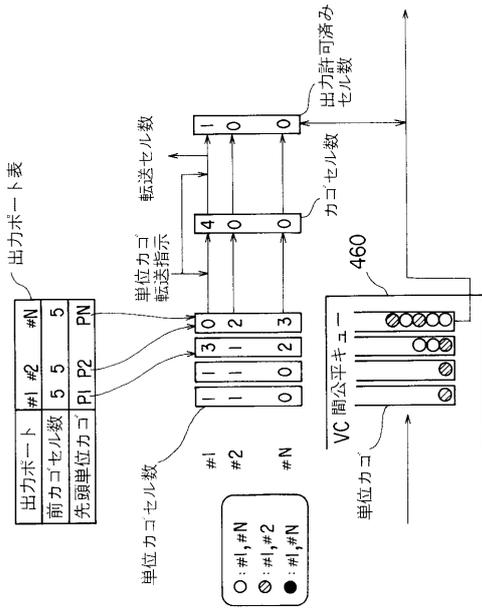
【図 31】



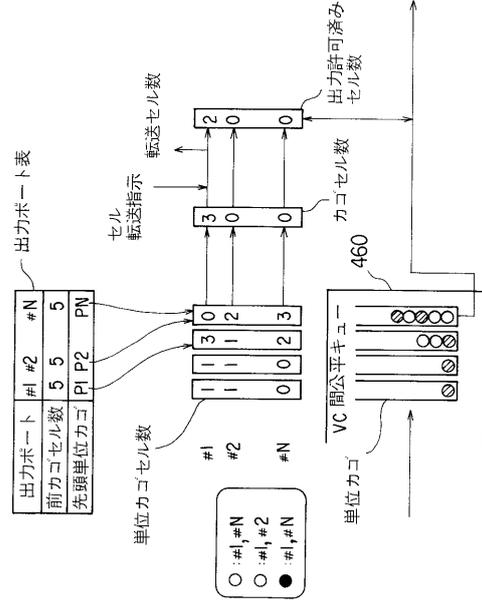
【図 32】



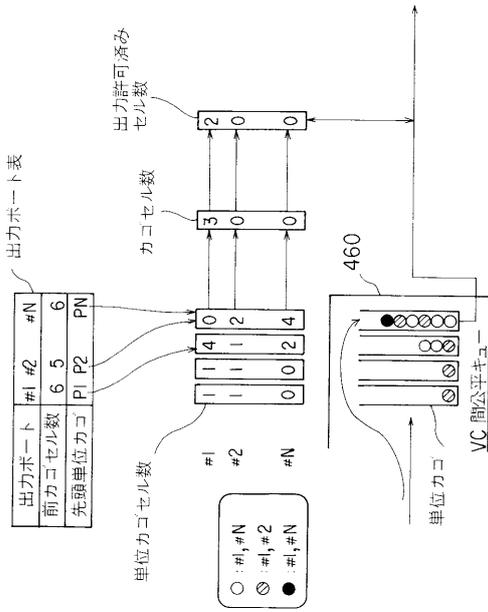
【 図 3 3 】



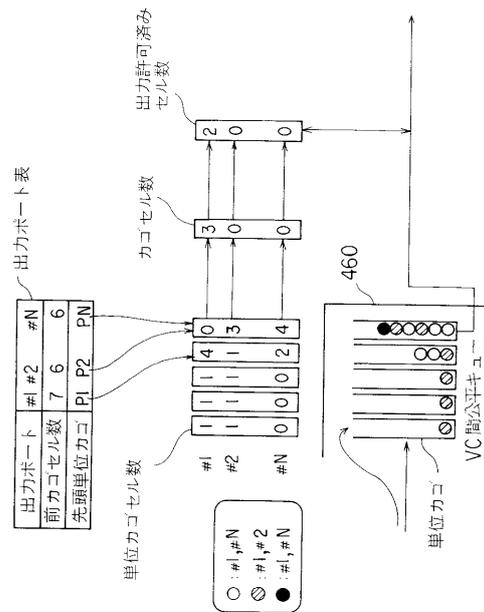
【 図 3 4 】



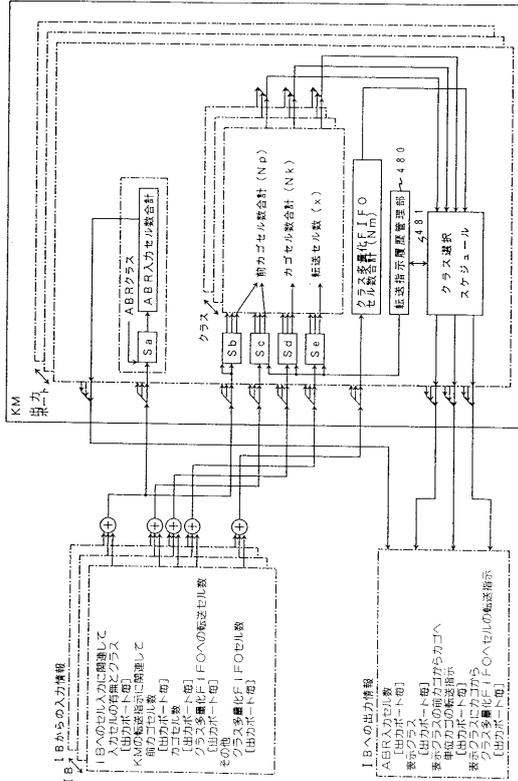
【 図 3 5 】



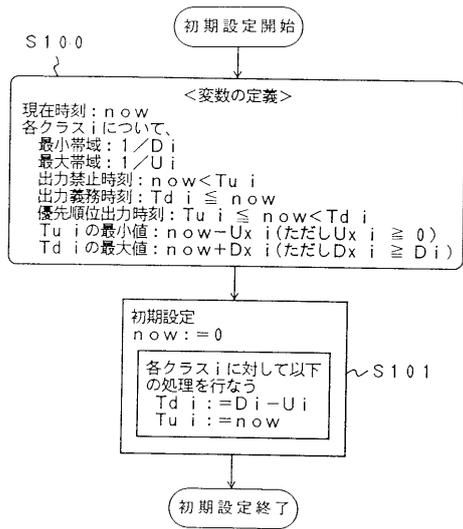
【 図 3 6 】



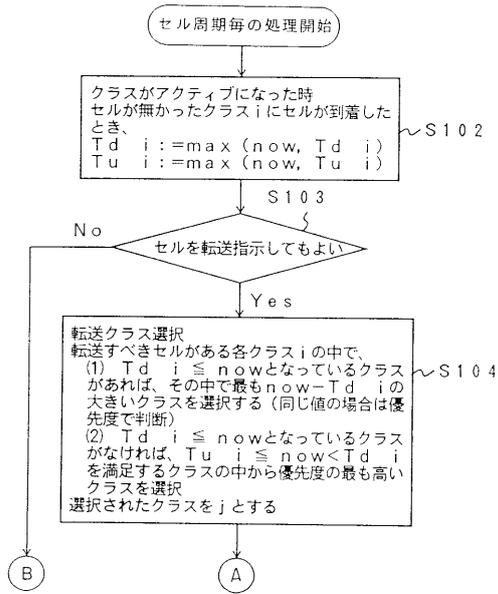
【図41】



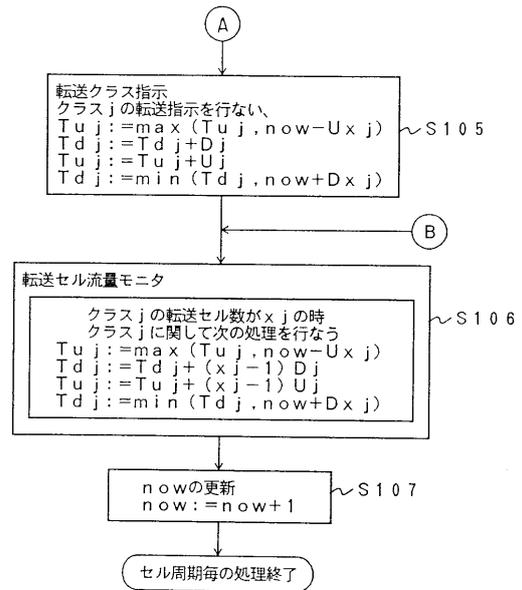
【図42】



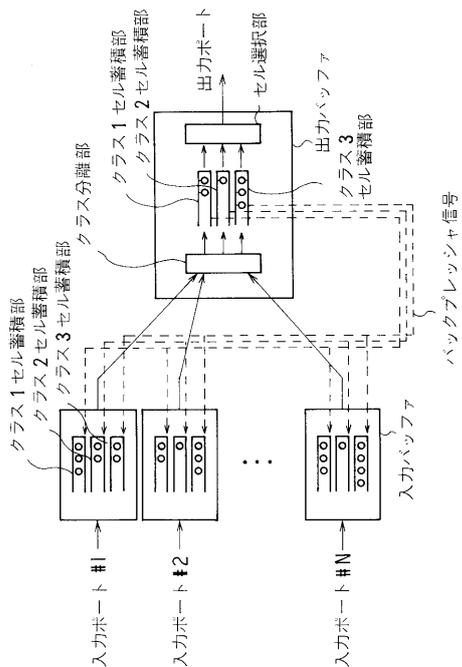
【図43】



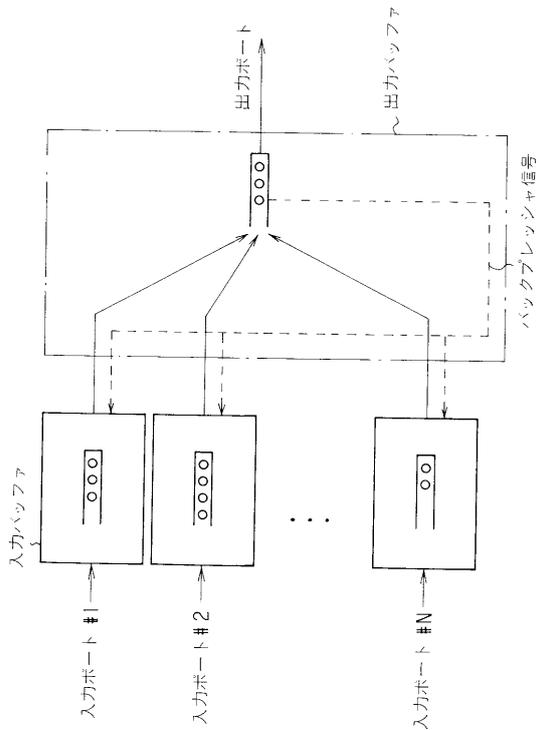
【図44】



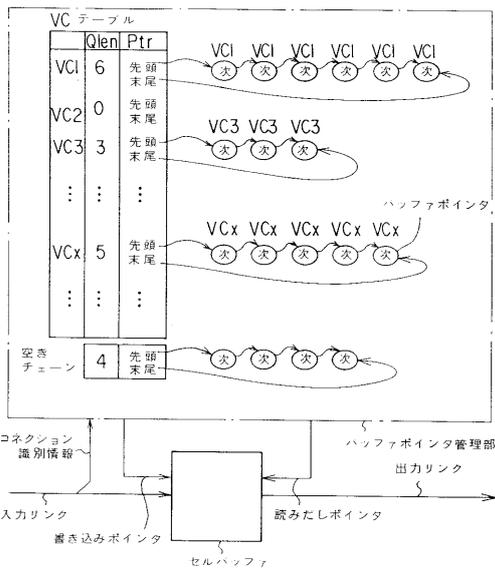
【 図 4 9 】



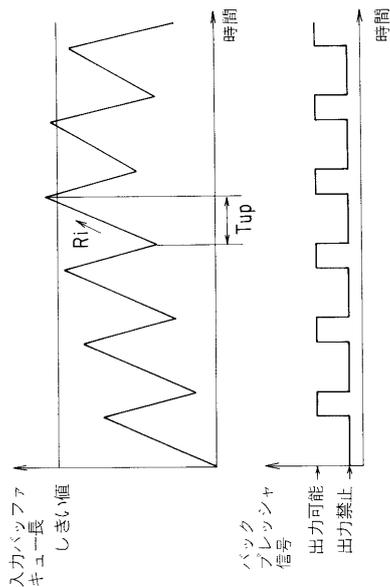
【 図 5 0 】



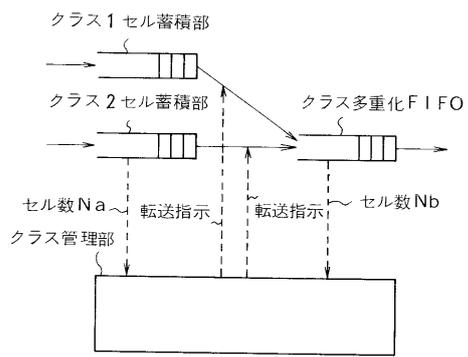
【 図 5 1 】



【 図 5 2 】



【 図 5 3 】



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

審査官 清水 稔

(56)参考文献 特開平06-224935(JP,A)

特開平07-193583(JP,A)

(58)調査した分野(Int.Cl.⁷, DB名)

H04L 12/28

H04Q 3/00