

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3953819号

(P3953819)

(45) 発行日 平成19年8月8日(2007.8.8)

(24) 登録日 平成19年5月11日(2007.5.11)

(51) Int. Cl.

H04L 12/56 (2006.01)

F I

H04L 12/56 200C

請求項の数 10 (全 20 頁)

(21) 出願番号	特願2002-601 (P2002-601)	(73) 特許権者	596092698
(22) 出願日	平成14年1月7日(2002.1.7)		ルーセント テクノロジーズ インコーポ
(65) 公開番号	特開2002-232469 (P2002-232469A)		レーテッド
(43) 公開日	平成14年8月16日(2002.8.16)		アメリカ合衆国, 07974-0636
審査請求日	平成16年11月2日(2004.11.2)		ニュージャージー, マレイ ヒル, マウン
(31) 優先権主張番号	60/260807		テン アヴェニュー 600
(32) 優先日	平成13年1月10日(2001.1.10)	(74) 代理人	100064447
(33) 優先権主張国	米国 (US)		弁理士 岡部 正夫
(31) 優先権主張番号	10/008425	(74) 代理人	100085176
(32) 優先日	平成13年11月13日(2001.11.13)		弁理士 加藤 伸晃
(33) 優先権主張国	米国 (US)	(74) 代理人	100106703
			弁理士 産形 和央
		(74) 代理人	100096943
			弁理士 臼井 伸一

最終頁に続く

(54) 【発明の名称】 スケジューリング装置およびスケジューリング方法

(57) 【特許請求の範囲】

【請求項1】

複数のデータパケットフローについてデータパケットの送信をスケジューリングするスケジューリング装置において、

前記データパケットフローには、通信リンクの送信容量 r のうちの所定シェアが割り当てられ、前記データパケットフローは複数のバンドルにまとめられ、

前記バンドルには、前記通信リンクの処理容量のうちのサービスシェアが割り当てられ、

通信リンクを通じての送信はサービスフレームに分けられ、1サービスフレームは、あらゆるバックログデータパケットフローに少なくとも1つの送信機会を提供し、バックログデータパケットフローは、複数のパケットキューのうちの対応するパケットキューに格納された少なくとも1つのデータパケットを有するデータパケットフローであり、

前記スケジューリング装置は、

サービスフレームの期間を決定する手段と、

各データパケットフローが、常に少なくとも該データパケットフローに割り当てられたサービスシェアを受けるとともに、各バンドルが、1サービスフレームの全期間中連続してバックログのあるバンドル内に少なくとも1つのデータパケットフローがある場合に、少なくとも該バンドルに割り当てられたサービスシェアを受けるとを保証する手段とを有し、

前記保証する手段は、

10

20

各バンドル I について、同じバンドル I にまとめられた前記データパケットフローのそれぞれに割り当てられる前記サービスシェアの総和に係る累積シェア Σ を管理する手段と、

各バンドル I について、該バンドル I に割り当てられるサービスシェア R_I と、該バンドルの前記累積シェア Σ との間のサービス比を計算する手段と、

前記複数のバンドルのそれぞれについて計算されたサービス比を用いて、前記複数のデータパケットフローのそれぞれに割り当てられる前記サービスシェアを修正する手段とを有することを特徴とするスケジューリング装置。

【請求項 2】

重み付けラウンドロビン (WRR) アルゴリズム、不足ラウンドロビン (DRR) アルゴリズム、および余剰ラウンドロビン (SRR) アルゴリズムを含む群から選択されるアルゴリズムが、データパケットの送信をスケジューリングするために用いられることを特徴とする請求項 1 記載のスケジューリング装置。

10

【請求項 3】

前記サービスフレームの期間を決定する手段は、

グローバルフレームカウンタ $FRMCNT$ と、

前記複数のバンドルの各バンドル I に対する開始フラグ I と、

前記複数のデータパケットフローの各データパケットフロー i に対するフレームフラグ FF_i とを有し、

全てのデータパケットについて、前記開始フラグと前記グローバルフレームカウンタとを比較することにより、データパケットが現サービスフレーム内でサービスを受ける権利がまだあるか否かが決定されることを特徴とする請求項 1 記載のスケジューリング装置。

20

【請求項 4】

前記グローバルフレームカウンタ $FRMCNT$ は、全てのサービスフレームの終わりにおいて値がトグルされるブールフラグであり、バンドル I の開始フラグ I は、該バンドル内の最初のデータパケットフローにバックログが生じたときに、グローバルフレームカウンタ $FRMCNT$ に等しくセットされることを特徴とする請求項 3 記載のスケジューリング装置。

【請求項 5】

データパケットフロー i のフレームフラグ FF_i は、該フローにバックログが生じるか、または、該フローが現サービスフレーム内で最後に処理されたフローであるときに、グローバルフレームカウンタ $FRMCNT$ とは異なる値にセットされることを特徴とする請求項 3 記載のスケジューリング装置。

30

【請求項 6】

サービスフレームの終了と次のサービスフレームの開始は、処理されるべき次のデータパケットフロー i のフレームフラグ FF_i がグローバルフレームカウンタ $FRMCNT$ とは異なる値を有するときに、同時に検出されることを特徴とする請求項 3 記載のスケジューリング装置。

【請求項 7】

バンドル I の累積シェア Σ の値は、バックログを生じているバンドル I のデータパケットフローのサービスシェアの総和に等しいことを特徴とする請求項 1 記載のスケジューリング装置。

40

【請求項 8】

バンドル I の累積シェア Σ の値は、該バンドルの最初のデータパケットフローがサービスフレーム内で最初にサービスされるときにセットされ、たとえバンドル I の 1 個または複数のデータパケットフローのバックログ状態が該サービスフレームの期間中に変化しても、同じサービスフレームの期間全体で不変のまま保持されることを特徴とする請求項 1 記載のスケジューリング装置。

【請求項 9】

現在シェア Σ は、バンドル I 内に蓄積されたデータパケットフローのサービスシェア

50

の合計を維持し、そしてバンドル内の1または複数のデータパケットフローのバックログ状態が変化するときに変化し、前記現在シェアの値は全ての新たなサービスフレームの開始時に前記累積シェア Σ の値を設定するのに使用されることを特徴とする請求項8記載のスケジューリング装置。

【請求項10】

複数のデータパケットフローについてデータパケットの送信をスケジューリングするスケジューリング方法において、

前記データパケットフローには、出通信リンクの送信容量のうちの所定シェアが割り当てられ、前記データパケットフローは複数のバンドルにまとめられ、

前記バンドルには、前記出通信リンクの送信容量 r のうちのサービスシェアが割り当てられ、

通信リンクを通じての送信はサービスフレームに分けられ、1サービスフレームは、あらゆるバックログデータパケットフローに少なくとも1つの送信機会を提供し、バックログデータパケットフローは、複数のパケットキューのうちの対応するパケットキューに格納された少なくとも1つのデータパケットを有するデータパケットフローであり、

前記方法は、

サービスフレームの期間を決定するステップと、

各データパケットフローが、常に少なくとも該データパケットフローに割り当てられたサービスシェアを受けるとともに、各バンドルが、1サービスフレームの全期間中連続してバックログのあるバンドル内に少なくとも1つのデータパケットフローがある場合に、少なくとも該バンドルに割り当てられたサービスシェアを受けることを保証するステップと、

各バンドル I について、同じバンドル I にまとめられた前記データパケットフローのそれぞれに割り当てられる前記サービスシェアの総和に係する累積シェア Σ_I を管理するステップと、

各バンドル I について、該バンドル I に割り当てられるサービスシェア R_I と、該バンドルの前記累積シェア Σ_I との間のサービス比を計算するステップと、

前記複数のバンドルのそれぞれについて計算されたサービス比を用いて、前記複数のデータパケットフローのそれぞれに割り当てられる前記サービスシェアを修正するステップとを有することを特徴とするスケジューリング方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、パケットスケジューラに関し、特に、個々のデータソースへ、および、それらのデータソースの集合(aggregate)へのデータ転送レートを保証するパケットスケジューリング装置および方法に関する。

【0002】

【従来の技術】

Integrated Service (Intserv) [後記の[付録]に掲げた文献1参照]やDifferentiated Service (Diffserv) [2]のような精緻なサービス品質(QoS)フレームワークがますます広まるとともに、フレキシブルな帯域管理の可能なパケットスケジューラが重要視されている。従来のいくつかのパケットスケジューラは、すぐれた最悪の場合の遅延性能を提供するとともに、正確な帯域保証を提供する[3, 4, 5, 6, 7, 8]が、それらのコストはかなり高い[6, 7, 8]。

【0003】

この理由のため、産業界は、重み付けラウンドロビン(WRR: Weighted Round Robin)スケジューラ[9, 10, 11]に相当の関心を示している。WRRスケジューラは、ぎりぎりまでの(タイトな)遅延限界を必ずしも達成するものではないが、最小限の複雑さで頑強(ロバスト)な帯域保証を提供する。このようなスケジューラのさまざまな実例が文献に現れている。よく知られた例には、不足ラウンドロビン(DRR: Deficit Round

10

20

30

40

50

Robin) アルゴリズム [1 0] および余剰ラウンドロビン (S R R : Surplus Round Robin) [1 1] アルゴリズムがある。

【 0 0 0 4 】

【 発明が解決しようとする課題 】

W R R スケジューラの上記の実例は、異種のデータパケットフローの帯域保証を差別化することには成功している。しかし、それらが現在規定されている形態に関しては、新しいサービス品質フレームワークのすべての帯域要求を満たすには十分でない。一般に、ネットワークノードにおけるフレキシブルな帯域管理は、帯域が、個々のフローに割り当てられることが可能であるだけでなく、それらのフローの集合にも割り当てられることが可能であるような、階層的スケジューリング構造の展開を必要とする。従来の W R R スケジューラでは、帯域分離を達成するための階層構造の重ね合わせは、基本のスケジューラの単純さを損なう。

10

【 0 0 0 5 】

データパケットネットワークのためのスケジューリングの技術分野においては、基本スケジューラの単純さを損なわずに階層的帯域分離を達成する改良されたスケジューラ装置が必要とされている。

【 0 0 0 6 】

【 課題を解決するための手段 】

本発明の目的は、階層的帯域分離のためのスケジューラ装置を規定することである。この新規な技術の目標は、完全に透過的な方法で、すなわち、追加のスケジューリング構造を使用せずに、帯域保証を、個々のデータパケットフローに提供するとともに、それらのフローの集合 (「バンドル」という) にも提供することである。

20

【 0 0 0 7 】

各バンドルについて、スケジューラは、そのバンドルに名目上割り当てられる帯域と、そのバンドルで現在積滞している (バックログのある) フローの帯域割当ての総和との間の比を決定する。スケジューラは、その比を用いて、個々のフローへの帯域の分配を規制するタイムスタンプ増分を修正する。このようにして、そのバンドルに対する比が大きいほど、そのバンドル内のそれぞれのバックログフローが受ける帯域は大きくなる。スケジューラは、常に、フローが積滞している限り、与えられたデータパケットフローの帯域要求を満たし、バンドル内に少なくとも1つのバックログフローがある限り、与えられたバンドルの帯域要求を満たす。

30

【 0 0 0 8 】

具体的には、本発明のスケジューリング装置は、複数のデータパケットフローについてデータパケットの送信をスケジューリングする。前記データパケットフローには、通信リンクの送信容量 r のうちの所定シェアが割り当てられ、前記データパケットフローは複数のバンドルにまとめられる。前記バンドルには、前記通信リンクの処理容量のうちのサービスシェアが割り当てられる。通信リンクを通じての送信はサービスフレームに分けられる。1サービスフレームは、あらゆるバックログデータパケットフローに少なくとも1つの送信機会を提供する。バックログデータパケットフローは、複数のパケットキューのうちの対応するパケットキューに格納された少なくとも1つのデータパケットを有するデータパケットフローである。本発明のスケジューリング装置は、(1) サービスフレームの期間を決定する手段と、(2) 各データパケットフローが、十分な数の連続するサービスフレームにわたり連続してそのデータパケットフローにバックログがある場合に、常に少なくともその割り当てられたサービスシェアを受けるとともに、各バンドルが、十分な数の連続するサービスフレームにわたり1サービスフレームの全期間中連続してバックログのあるバンドル内に少なくとも1つのデータパケットフローがある場合に、少なくともその割り当てられたサービスシェアを受けるとを保証する手段とを有する。前記保証する手段は、(A) 各バンドル I について、同じバンドル I にまとめられた前記データパケットフローのそれぞれに割り当てられる前記サービスシェアの総和に関係する累積シェア I を管理する手段と、(B) 各バンドル I について、該バンドル I に割り当てられるサービ

40

50

スシェア R_I と、該バンドルの前記累積シェア Σ との間のサービス比を計算する手段と、(C)前記複数のバンドルのそれぞれについて計算されたサービス比を用いて、前記複数のデータパケットフローのそれぞれに割り当てられる前記サービスシェアを修正する手段とを有する。

【0009】

本発明はまた、複数のデータパケットフローについてデータパケットの送信をスケジューリングする方法に関する。前記データパケットフローには、出通信リンクの送信容量のうちの所定シェアが割り当てられ、前記データパケットフローは複数のバンドルにまとめられる。前記バンドルには、前記出通信リンクの送信容量 r のうちのサービスシェアが割り当てられる。通信リンクを通じての送信はサービスフレームに分けられる。1サービスフレームは、あらゆるバックログデータパケットフローに少なくとも1つの送信機会を提供する。バックログデータパケットフローは、複数のパケットキューのうちの対応するパケットキューに格納された少なくとも1つのデータパケットを有するデータパケットフローである。この方法は、(1)サービスフレームの期間を決定するステップと、(2)各データパケットフローが、十分な数の連続するサービスフレームにわたり連続してそのデータパケットフローにバックログがある場合に、常に少なくともその割り当てられたサービスシェアを受けるとともに、各バンドルが、十分な数の連続するサービスフレームにわたり1サービスフレームの全期間中連続してバックログのあるバンドル内に少なくとも1つのデータパケットフローがある場合に、少なくともその割り当てられたサービスシェアを受けるとを保証するステップと、(3)各バンドル I について、同じバンドル I にまとめられた前記データパケットフローのそれぞれに割り当てられる前記サービスシェアの総和に係する累積シェア Σ を管理するステップと、(4)各バンドル I について、該バンドル I に割り当てられるサービスシェア R_I と、該バンドルの前記累積シェア Σ との間のサービス比を計算するステップと、(5)前記複数のバンドルについて計算されたサービス比を用いて、前記複数のデータパケットフローのそれぞれに割り当てられる前記サービスシェアを修正するステップとを有する。

【0010】

【発明の実施の形態】

図1に、複数の通信スイッチ $101-1 \sim 101-p$ が通信リンクにより相互に接続された例示的なパケットネットワークを示す。いくつかのデータソース $102-1 \sim 102-q$ が通信スイッチに接続される。ネットワークコネクションが、それぞれのデータソースから、対応する宛先 $103-1 \sim 103-g$ へと確立され、データパケットは、各データソースから対応する宛先へ送信される。

【0011】

図2は、パケットネットワークの通信スイッチ $101-1$ の例示的なブロック図である。図示のように、通信スイッチは、スイッチファブリック 250 と、複数の通信リンクインタフェース $200-1 \sim 200-s$ を有する。各通信リンクインタフェースは、複数の入力リンクを1つの出力リンクに接続し、データパケットを入力リンクから出力リンクへ転送する。通信スイッチ $101-1$ は、このような通信リンクインタフェース 200 を1個だけ含むことも、複数含むことも可能である。例えば、入力通信リンクインタフェース $200-1$ はスイッチファブリック 250 の前に配置され、その場合、入力通信リンクインタフェース $200-1$ の入力リンク $201-1 \sim 201-r$ は通信スイッチ $101-1$ の入力リンクであり、入力通信リンクインタフェース $200-1$ の出力リンク 203 はスイッチファブリック 250 に接続される。

【0012】

第2の例として、出力通信リンクインタフェース $200-j$ はスイッチファブリック 250 の出力に配置され、出力通信リンクインタフェース $200-j$ の入力リンクはスイッチファブリック 250 の複数の出力リンク 204 であることが可能であり、出力通信リンクインタフェース $200-j$ の出力リンクは通信スイッチ $101-1$ の出力リンク $202-j$ である。注意すべき点であるが、1つの特定のリンクを通じて、または、複数の相異なる

10

20

30

40

50

るリンクを通じて受信されるパケットは、同じ長さであることもそうでないこともある。例えば、スイッチファブリック 250 が非同期転送モード (ATM) スイッチであり、図 1 のネットワークが ATM ネットワークである場合、すべてのパケットは同じ長さを有する。本発明の以下の説明では、1 つの特定のリンクを通じて、または、複数の相異なるリンクを通じて受信されるパケットは、長さが必ずしも同じでないと仮定する。

【0013】

図 6 に関連して後の段落で説明するように、図 2 の通信リンクインタフェース 200 のそれぞれは一般に、少なくともパケット受信器、スケジューラ、およびパケット送信器を有する。前述のように、スケジューラは、重み付けラウンドロビン (WRR) スケジューラ [9, 10, 11] であることが可能であり、さらにこれは、不足ラウンドロビン (DRR : Deficit Round Robin) [10] または余剰ラウンドロビン (SRR : Surplus Round Robin) [11] のアルゴリズムに従って実装されることが可能である。

10

【0014】

不足ラウンドロビン (DRR) アルゴリズムは、その最小限の実装複雑さと、割り当てられるサービスシェアに比例してフローをサービスする際のその効率とにより、可変サイズのパケットに対する WRR スケジューラの最も普及している実例の 1 つである。WRR パラダイムに従って、DRR アルゴリズムは、サービスシェア s_i を各設定フロー i に関連づける。サービスシェアは、すべての設定フローにわたるそれらの和がサーバの容量 r を超えないときに、最小の保証サービスレートとなる。

【数 1】

20

$$\sum_{i=1}^V \rho_i \quad (1)$$

式 (1) の限界 (ただし、 V は設定フローの総数である) は、 s_i より低くない長期レートでフロー i がサービスを受けることを保証する。

【0015】

DRR アルゴリズムは、サーバのアクティビティをサービスフレームに分割する。本発明は、基準タイムスタンプ増分 T_Q を用いて仮想時間領域でフレーム期間を表現するアルゴリズムの定式化を参照する。この定式化は、最初に [10] で提示された DRR の定義と機能的に等価であるが、本発明の記述にはより適している。

30

【0016】

1 フレーム内で、各設定フロー i には、次式のような量 Q_i の情報ユニットの送信の権利が与えられる。

$$Q_i = s_i \cdot T_Q \quad (2)$$

【0017】

スケジューラは、フレームあたり 1 回だけバックログフローを処理するため、一発でそのフレームに対するそれらのフローのサービス期待値を充足する。各フロー i は、パケットのキュー (フローキュー) と、長さ l_i^k の新たなパケット p_i^k がフローキューの先頭に到着するごとに更新される次式のタイムスタンプ F_i とを管理する。

40

$$F_i^k = F_i^{k-1} + l_i^k / s_i \quad (3)$$

【0018】

スケジューラは、フローへのサービスを、そのタイムスタンプが T_Q より小さくとどまる限り続ける。タイムスタンプ F_i が基準タイムスタンプ増分 T_Q を超えると、スケジューラは、フロー i の処理の終了を宣言し、 T_Q をフロー i のタイムスタンプから減算し、サービスすべき別のバックログフローを探す。その結果、 T_Q の減算後、 F_i の値は、フロー i に対するサービスクレジット (貸し) を表現する。一般に、タイムスタンプは、バックログフローのサービスクレジットを後のフレームに持ち越すため、スケジューラは、長

50

期間では(すなわち、複数のフレームにわたって)、割り当てられたサービスシェアに比例してサービスを分配することができる。

【0019】

フロー*i*がアイドルになると、スケジューラは直ちに別のフローに移り、それへのサービスの提供を開始する。フロー*i*が短時間で再びバックログになると、フロー*i*は、サーバからの次の処理を受けるためには、次のフレームの開始を待たなければならない。フローがアイドルになると、同じフローが将来のフレームで再びバックログになった場合のサービスの損失を避けるために、そのフローのタイムスタンプはゼロにリセットされる。構成により、アイドルのフローのタイムスタンプは常に T_Q より小さいため、タイムスタンプのリセットは、他のフローに不利になるような余分のサービスクレジットを生じることはない。

10

【0020】

構成により、サービスフレームの始めには、フロー*i*のタイムスタンプ F_i の値は0と L_i / μ_i の間にある。ただし、 L_i はフロー*i*のパケットの最大サイズである。タイムスタンプの初期値のゆらぎは、フロー*i*が1フレームで送信する情報ユニットの量のゆらぎを引き起こし、これは区間 $(Q_i - L_i, Q_i + L_i)$ 内にある。このため、すべての設定フローが永続的に積滞しているときでさえ、サーバが1フレームで送信する情報ユニットの総量は一定でない。

【0021】

DRRスケジューラは、[10]では、FIFO順に巡回されるバックログフローの単一の連結リストにより実装された。次式のように、基準タイムスタンプ増分 T_Q が、最小サービスシェアのフローに対する最大サイズのパケットにより定まるタイムスタンプ増分より小さくない場合には、バックログフローを単一FIFOキューに配置すると、 $O(1)$ の実装複雑さとなる。

20

$$T_Q = L_{max} / \mu_{min} \quad (4)$$

【0022】

式(4)の条件が満たされない場合、スケジューラのアルゴリズム複雑さは、連続するパケット送信間に行われるべき基本オペレーションの最悪の場合の数とともに爆発的に増大する(基本オペレーションには、連結リストにおけるフローの取り出しおよび挿入、タイムスタンプ更新、タイムスタンプと基準タイムスタンプ増分の比較、が含まれる)。実際、スケジューラは、基準タイムスタンプ増分の反復減算によりタイムスタンプが $[0, T_Q)$ 区間内に入るまで、いくつかの連続するフレームの間、与えられたフローへのサービスを拒否しなければならない可能性がある。図3Aに、DRRにおいて、新たなパケットがフロー*i*のキューの先頭に到着したときに、そのフロー*i*を処理しそのタイムスタンプを更新するための規則を指定する例示的な擬似コードのリストを示す。

30

【0023】

余剰ラウンドロビン(SRR)の記述は[11]に与えられている。このアルゴリズムは、DRRと同じパラメータおよび変数を有するが、異なるイベントがタイムスタンプの更新を引き起こす。すなわち、フロー*i*は、パケット p_i^k の送信が完了すると、その結果のフローのバックログ状態とは独立に、新たなタイムスタンプ F_i^k を受け取る。フレームの終わりは、1パケットの送信後に常に検出され、その前には検出されない。タイムスタンプは、(DRRの場合のクレジット(貸し)の代わりに)現フレームの期間中にフローによって累積されたデビット(借り)を次フレームに持ち越す。

40

【0024】

DRRに比べてSRRの利点は、SRRは、フローのフレームの終わりを決定するためにキュー先頭パケットの長さを前もって知る必要がないことである。逆に、悪意のあるフローが競争相手から帯域を盗むことを防ぐために、このアルゴリズムは、アイドルになったフローのタイムスタンプをリセットすることができない。アイドルフローのヌルでないタイムスタンプは、そのフローがアイドルになったのと同じフレームの終了によって結局は不用(obsolete)となる。理想的には、タイムスタンプは、不用になったらすぐリセット

50

されるべきである。しかし、数十万、さらには数百万のフローを処理するスケジューラでは、同時に不用になる可能性のあるすべてのタイムスタンプを即座にリセットすることは現実には不可能である。

【0025】

本発明についてのこの記述は、アイドルフローのタイムスタンプの不用性を何らチェックしないSRRアルゴリズムの実装に焦点をさぼる。この実装では、新たに積滞したフローは常に、そのタイムスタンプの値がいかにか古くとも、そのタイムスタンプの最近の値からアクティビティを再開する。この仮定の効果は、新たに積滞したフローが、前に長時間の間に累積したデビットの結果として、サーバが最初にやってくる時にその受けるべきサービスの一部をあきらめなければならないことがあることである。図3Bに、SRRにおいて、サーバがパケット p_i^k の送信を完了したときにフロー i を処理しそのタイムスタンプを更新するための規則を指定する例示的な擬似コードのリストを示す。

10

【0026】

説明を簡単にするため、本明細書の残りの部分では、重み付けラウンドロビン(WRR)の名称は一般に、DRRまたはSRRを総称的に(それらを区別する特徴への明示的な言及がなければ)指すために用いるものとする。

【0027】

WRRスケジューラは、個々のパケットフローへの帯域の分配だけを制御することができるという点で、本質的に「1層」である。複数のスケジューリング層を重ね合わせることで、個々のフローに帯域を割り当てることができるだけでなく、フローの集合を生成しそれに従って帯域を分離することができる、階層的でフレキシブルな構造を実装すると、これらのスケジューラの単純さが損なわれる。

20

【0028】

帯域分離の参照モデルを図4に示す。割り当てられるフローのセット $(1, 1) \sim (K, V_K)$ は、バンドルと呼ばれる K 個のサブセット $401-1 \sim 401-K$ に分割される。各バンドル I は、 V_I 個のフローを集め、割り当てられたサービスレート R_I を有する。スケジューラの論理編成は、2層化された階層を反映する。すなわち、まず、バンドルの総割当てに従って帯域をバンドルに分配した後、フローの帯域割当てと、それぞれのバンドル内の他のフローのバックログ状態とに基づいてフローをサービスする。スケジューラは、少なくとも1つのフローが積滞している限り、各バンドルを、対応するフローのバックログ状態とは独立に扱う。

30

【0029】

スケジューラは、いかなる種類の遅延保証もサポートしようとせずに、フロー集合と、集合内の個々のフローの両方に対して、厳格な帯域保証を強制することを目標とする(スケジューリング階層において厳格な遅延保証を提供するためのいくつかのフレームワークがすでに利用可能である[7, 12]が、それらはすべて、スケジューラの複雑さを相当に増大させる高度なアルゴリズムによるものである)。フロー集合の帯域要求をサポートするためには、バンドルのレート割当てに対して次の条件が常に成り立たなければならない。

【数2】

$$\sum_{I=1}^K R_I \leq r \quad (5)$$

40

【0030】

同様に、各バンドル I 内で、関連するフローの帯域要求を満たすためには、次の限界が満たされなければならない。

【数3】

50

$$\sum_{i \in I} \rho_i \leq R_I \quad (6)$$

【 0 0 3 1 】

[7 , 1 2] で提示されたフレームワークによって示唆されるスケジューリング解では、フローとリンクサーバの間に、バンドルを処理する本格的な（そして高価な）スケジューリング層を導入することになる。一般に、本格的な階層スケジューラの実装コストは、バンドルの数について線形に増大する。その理由は、各バンドルがそれぞれ基本のフロー別スケジューラを必要とするからである。これに対して、本発明では、スケジューリング階層におけるバンドル要求を強制する層は純粋に仮想的であり、単一の基本スケジューラ上に重ね合わされる。したがって、個々のフローを処理する構造のコストは、設定されるバンドルの数とは独立であるため、スケジューリング階層の実装における相当の節約になる。

10

【 0 0 3 2 】

本発明によれば、改良されたWRRスケジューラは、WRRスケジューラの基本構造の実質的な変更を必要とせず、帯域を階層的に分離する。この新規な技術によれば、スケジューラは、完全に透過的な方法で、すなわち、追加のスケジューリング構造を使用せずに、帯域保証を、個々のフローに提供するとともに、フローの集合（バンドル）にも提供する。本発明は、この目標を、単に、スケジューラがタイムスタンプを操作する方法を改良することによって達成する。その結果得られる「ソフト」なスケジューリング階層は、無視できるほどの複雑さしか有しない。にもかかわらず、帯域保証を、個々のフローに提供するとともに、バンドルにも提供するのに有効である。

20

【 0 0 3 3 】

図5に、階層的帯域分離のソフト強制のための改良されたWRRスケジューラによって利用されるキュー、状態テーブル、レジスタ、およびパラメータの機能図である。図4および図5をあわせて参照すると、スケジューラは、複数のバンドル401-1~401-Kにまとめられた、複数のデータパケットフロー $i_1 \sim j_N$ (501)を処理する。フローキュー502は、それぞれのデータパケットフロー501についてデータパケットを格納する。フローキュー502は、先入れ先出し(FIFO)キューとして実装可能である。各フローキュー502は、対応するフロー別状態テーブル503を有し、これは、対応するフロー501に対するいくつかの変数を格納する。

30

【 0 0 3 4 】

例えば、フロー i_1 の場合、フロー別状態テーブル503は、タイムスタンプ F_{i_1} 、最小保証サービスレート ρ_{i_1} 、フレームフラグ FF_{i_1} 、そのフローを含むバンドルIへのポインタ、ならびに、対応するフローキュー502の先頭および末尾のポインタ（図示せず）を含む。バンドルの状態は、バンドル別状態テーブル504で管理される。例えば、バンドルIの場合、バンドル別状態テーブル504は、総帯域割当て R_I 、現在シェア ρ_I 、累積シェア $\rho_{I,acc}$ 、および開始フラグ $start_I$ を格納する。フローポインタのFIFOキュー505は、スケジューラがフロー501のパケットを送信するために巡回する順序を決定する。レジスタ506は、フローポインタのFIFOキューの先頭および末尾のポインタを格納する。スケジューラは、テーブル507でグローバル状態情報を管理する。このテーブルは、フレームカウンタFRMCNTおよび基準タイムスタンプ増分 T_Q を管理する。

40

【 0 0 3 5 】

図6は、スケジューラが利用される入力通信リンクインタフェース200の例示的なブロック図である。通信リンクインタフェース200は、データパケット受信器601、スケジューラ602、およびパケット送信器609を有する。例として、スケジューラは、コントローラ603、バンドル別状態RAM604、およびレジスタ605を、すべて同一

50

のチップ606上に有するように示されている。パケットRAM607およびフロー別状態RAM608は、別々のチップ上に位置するように示されている。明らかに、動作容量およびその他の特性に依存して、スケジューラ602は、他の構成でも実装可能である。

【0036】

コントローラ603は、本発明の方法を実装したプログラムを格納し実行する。通信リンクインタフェース200の動作を制御するプログラムの実施例を図7～図8にフローチャート形式で示す。図5および図6をあわせて参照すると、フローキュー502内のパケットはパケットRAM607に格納される。フロー別状態テーブル503はフロー別状態RAM608に格納される。フローポイントのFIFOキュー505の先頭および末尾ポイント506ならびにグローバル状態テーブル507はレジスタ605に格納される。

10

【0037】

スケジューラ602の動作の概略は以下の通りである。パケット受信器601は、入力リンク201-1～201-rから、データパケットフロー501のデータパケットを受信する。パケット受信器601は、各パケット(図示せず)のヘッダに含まれるフロー識別フィールドの内容を用いて、そのパケットのデータパケットフロー501を識別する。データパケットフロー501の識別により、関連するフローキュー502およびバンドル401が識別される。スケジューラ602は、バンドル別状態テーブル504のエントリを用いて、各バンドルIに対するサービス比を計算する。バンドルIのサービス比は、公称帯域割当て R_I と、そのバンドルの累積シェア Σ_I との間の比として定義される(累積シェア Σ_I は、バンドルIのバックログフローの帯域割当てを累積したものである)。構成により、バンドルIのサービス比は決して1より小さくならない。

20

【0038】

サービス比には、個々のフローに対するタイムスタンプ増分の計算が関係する。これは、バンドルIの各バックログフローが、同じバンドルの他のフローがフレームの最初にアイドルであると検出されたときのそのフレームの期間中に受けるサービスの量の増大分を決定する。サービス比により引き起こされるタイムスタンプ増分の修正は、スケジューラが、バンドルの総帯域割当てと、個々のフローの帯域割当てとの両方を依然として保証することができるようになされる。

【0039】

[動作の詳細]

本発明の改良されたWRRスケジューラは、従来技術のWRRスケジューラ上に仮想スケジューリング層を重ね合わせることによって、サービス帯域の階層的分離を達成する。本発明のスケジューラは、複数のバンドルをサポートし、各バンドルは、データパケットフローの集合である。基礎となるWRRスケジューラのタイムスタンプに基づく定式化において(説明のこの点では、DRRとSRRの区別はまだ重要ではない)、本発明の技術は、タイムスタンプ更新の簡単な修正により、個々のフローの帯域保証を保存すると同時に、バンドルの総要求を満たす。

30

【0040】

それぞれの設定されるバンドルIについて、スケジューラは、公称帯域割当て R_I および累積シェア Σ_I を管理する。公称帯域割当て R_I は、そのバンドル内のすべてのフローの帯域割当ての総和より小さくなることはない。累積シェア Σ_I は、バンドルIのバックログフローの帯域割当ての総和を追跡し、したがって、決して R_I より大きくならない。

40

【数4】

$$\Phi_I = \sum_{i \in B_I} \rho_i \quad (7)$$

式(7)において、 B_I は、フレームの最初に積滞している、バンドルIのフローの集合である。

50

【 0 0 4 1 】

スケジューラは、バンドル I の累積レート割当て R_I と累積シェア Φ_I の間のサービス比を用いて、そのバンドルのフローに付与される実際のサービス帯域を修正する。具体的には、サービス比は、バンドルのフローに対するタイムスタンプ増分の定義に寄与する。フロー i のパケット p_i^k に関連するタイムスタンプ更新は次式の通りである。

【数 5】

$$F_i^k = F_i^{k-1} + \frac{l_i^k}{\rho_i} \cdot \frac{\Phi_I}{R_I} \quad (8)$$

10

【 0 0 4 2 】

式 (8) のタイムスタンプ割当て規則が実際にバンドルの帯域保証を強制することを確認するには、バンドル I 内の個々のフローが 1 フレームの期間中に受けることが期待されるサービスの量を計算することが必要である。原理的には、サービス比 R_I / Φ_I がバンドル I について 1 より大きいとき、式 (8) のタイムスタンプ割当ての結果、バンドル I の個々のフローは、それらの公称帯域割当ての R_I / Φ_I 倍のサービス帯域を得る。

【 0 0 4 3 】

この計算は、次の 2 つの仮定に基づく。

(1) バンドルの累積シェア Φ_I は、バンドル内のフローのバックログダイナミクスとは無関係に、1 フレームの間中は不変である。

20

(2) そのフレームの期間中にサーバにアクセス可能なフローのセットは、そのフレームの最初に積滞していたフローのみを含む (もしバンドル内の一部のフローが、フレームの開始後に積滞した場合、それらのフローは、サーバにアクセス可能になる前に、新たなフレームの開始まで待たなければならない) 。

【 0 0 4 4 】

基準タイムスタンプ増分 T_Q は、式 (8) と組み合わせて、バンドル I のフロー i が 1 フレーム期間中に受けることが期待される基準サービス量 Q_i を設定する。

【数 6】

$$T_Q = \frac{Q_i}{\rho_i} \cdot \frac{\Phi_I}{R_I} \quad (9)$$

30

【 0 0 4 5 】

すると、バンドル I 内のすべてのフローのサービス量の総計は、バンドルのサービス量 Q_I となる。

【数 7】

$$Q_I = \sum_{i \in B_I} Q_i = \frac{\sum_{i \in B_I} \rho_i}{\Phi_I} \cdot R_I \cdot T_Q = R_I \cdot T_Q \quad (10)$$

40

【 0 0 4 6 】

式 (10) における Q_I の表式は、式 (2) におけるフロー量 Q_i の表式と同一であるため、フレームの最初にバンドル内で積滞しているフローのセットの組成とは無関係に、式 (8) のタイムスタンプ更新がバンドル I の帯域保証を保存することが証明される。

【 0 0 4 7 】

バンドル I の累積シェアがフレーム期間中に変化しないという仮定を保持すると、次式のように、式 (8) のタイムスタンプ更新規則は、フレーム期間中に決してアイドルにならないバンドル I の任意の 2 個のフロー i 、 j について、サービス割合を保存することも示

50

すことができる。

【数 8】

$$\frac{Q_i}{Q_j} = \frac{\rho_i \cdot \frac{R_I}{\Phi_I} \cdot T_Q}{\rho_j \cdot \frac{R_I}{\Phi_I} \cdot T_Q} = \frac{\rho_i}{\rho_j} \quad (11)$$

10

【0048】

帯域分離を有するWRRアルゴリズムの詳細を指定するには、式(10)および(11)の結果を生成する仮定についての議論と、それら仮定のアルゴリズム的な意味の評価とを必要とする。スケジューラが1フレーム期間中に計算するすべてのタイムスタンプ増分において一定値の累積シェア Φ_I を使用することは、バンドルIのフローにサービスを一貫して分配するための共通の基準を提供する。フレームの開始後になってはじめて積滞したフローをサービスフレームから排除することも、同一の目的を有する。

【0049】

タイムスタンプ増分は、システムがフローに対して、関連するパケットの送信について課す料金とみなすことができる。送信のコストは、それが実行されるときのパンドル内で利用可能な帯域に依存する。タイムスタンプ増分を、バンドル内の帯域資源のコストと整合させるためには、タイムスタンプ増分は、資源が使用されるときに、すなわち、対応するパケットの送信時に、計算されなければならない。もしスケジューラが、増分を前もって計算するならば、バンドルの状態は、パケットの送信が起こる前に急激な変化を受ける可能性があるため、課金メカニズムが帯域の分配と整合しなくなる。

20

【0050】

考えている2つのWRRアルゴリズムのうち、SRRは、送信とタイムスタンプ増分の間の整合性に対する要求に最もよく適合するものである。その理由は、SRRは、ちょうど送信したばかりのパケットの長さを用いてタイムスタンプを更新し、対応するフローのフレーム内状態を決定するからである。これに対して、DRRでは、スケジューラは、新たなキュー先頭パケットの長さを用いて、タイムスタンプ更新およびフレーム内状態チェックを実行するが、これは、そのパケットが実際に送信されるずっと前である可能性がある。DRRサーバが最終的にパケットを配信するとき、バンドルの累積シェア、したがってそのバンドル内の帯域のコストは、最後のタイムスタンプ更新以来相当に変化してしまっている可能性がある。

30

【0051】

SRRにおける帯域分離のためのメカニズムの導入は簡単である。最小帯域保証 R_I および累積シェア Φ_I に加えて、各バンドルIは、現在シェア Φ_I および開始フラグ Φ_I を管理する。現在シェアは、次式のように、バンドル内のバックログフローのサービスシェアの総和を瞬間的に追跡する。

40

【数 9】

$$\phi_I(t) = \sum_{i \in B_I(t)} \rho_i \quad \forall t \quad (12)$$

【0052】

現在シェアは、バンドルの1つのフローがそのバックログ状態を変えるたびに更新される。一般に、現在シェア Φ_I の更新は、累積シェア Φ_I の即時更新を引き起こさない。実際

50

、スケジューラは、開始フラグ I と、スケジューラがフレーム境界ごとにトグルする 1 ビットのグローバルなフレームカウンタ $FRMCNT$ とにおける整合しない値を検出した場合にのみ、バンドルの累積シェア I を更新する（スケジューラは、バンドル I の 1 つのフローのサービスを開始するごとに、 I と $FRMCNT$ を比較する）。

【0053】

これらの 2 つのビットが相違することは、将来のタイムスタンプ計算で用いられる累積シェアの更新を引き起こす（ I I ）とともに、バンドルの開始フラグをトグルする（ I $FRMCNT$ ）。他方、これらの 2 つのビットがすでに等しい場合、ちょうど完了したサービスは、現フレームの期間中にバンドルが受けた最初のサービスではないことが確実であり、バンドルパラメータに対して何らのアクションもとられてはならない。バンドルの最初のフローが積滞すると、開始フラグは $FRMCNT$ に等しくセットされる。

I $FRMCNT$ (13)

【0054】

フレームの終わりを識別するため、各フロー i はフレームフラグ FF_i を管理する。フロー i のフレームフラグは、そのフローがバックログフローのリスト 505 の末尾にキューイングされるときには必ず、 $FRMCNT$ の補数に等しくセットされる。スケジューラは、フレームカウンタに一致しないフレームフラグを見つけると、新たなフレームの開始を宣言し、フレームカウンタをトグルする。1 パケットの送信を完了した後に実行される一連のオペレーションは、図 9 の擬似コードにまとめられている。

【0055】

図 7 ~ 図 8 に、本発明に従ってパケット送信のスケジューリングを制御するために図 6 のスケジューリング装置を動作させる方法をフローチャート形式で示す。図 7 ~ 図 8 のフローチャートおよび図 9 の擬似コードは、 SRR が、基礎となるスケジューリングアルゴリズムであるという仮定に基づいている。機能性に関する限り、 SRR の代わりに DRR を用いても問題はない。同様に、図 5 の装置は、バックログフローの単一の $FIFO$ キューを用いてソフトスケジューリング階層を実装している。フレーム内とフレーム外のフローの明確な分離が可能な他の任意のキューイング構造も同様に使用可能である。

【0056】

以下の記述では、図 4、図 5、図 6、および図 7 ~ 図 8 を参照する。図 4（図 5、図 6）で最初に現れる要素の参照符号は 4（5、6）から始まるが、図 7 ~ 図 8 のステップは、ステップ番号の前に S を付けることによって（例えば、 $S510$ ）示される。

【0057】

図 7 において、 $S510$ で、コントローラ 603 は、新たに受信したデータパケットがあるかどうかをチェックする。 $S510$ で、新たに受信したデータパケットがなく、 $S520$ で、バックログフローがある場合、制御は $S680$ に移る。他方、 $S510$ で、新たに受信したデータパケットがなく、 $S520$ で、バックログフローがない場合、コントローラ 603 は、新たなパケットが受信されるまで、ステップ $S510$ と $S520$ の間を循環する。 $S510$ において、新たに受信したパケットの存在が受信器 601 で検出されると、 $S550$ で、コントローラ 603 は、それらのパケットのうちの 1 つを選択する。その後、 $S560$ で、コントローラ 603 は、データパケットのフローを識別し、最終的に（ $S570$ で）、パケットを適当なフローキュー 502 に格納する。 $S580$ で、識別されたフローのキューの長さがゼロでない場合、 $S970$ で、そのフローのキュー長をインクリメントし、制御は $S680$ に移る。

【0058】

他方、 $S580$ で、識別されたフローのキューの長さがゼロである場合、 $S585$ で、そのフローのフレームフラグ（フロー別状態テーブル 503 内）が $FRMCNT$ （グローバル状態テーブル 507 内）の補数に等しくセットされる。その後、コントローラは、 $S590$ で、バックログフローの総数を増大させ、 $S600$ で、識別されたフローのキュー長を増大させる。 $S610$ で、識別されたフローのバンドルが、フロー別状態テーブル 503 のバンドルポインタを用いて識別される。識別されたバンドルの現在シェアが 0 である

10

20

30

40

50

場合、S 6 3 0で、コントローラ6 0 3は、バンドルの開始フラグ（バンドル別状態テーブル5 0 4に格納される）を、グローバルフレームカウンタFRMCNTの値に等しくセットする。その後、制御はS 6 3 5に移る。他方、S 6 2 0で、バンドルの現在シェアが0でない場合、制御は直接S 6 2 0からS 6 3 5に移る。S 6 3 5で、識別されたバンドルの現在シェアが、新たに積滞したフローの帯域割当てだけ増大させられる。そして、S 6 4 0で、フローはフローポインタのFIFOキュー5 0 5の末尾に付加され、その後、制御はS 6 8 0に移る。

【0 0 5 9】

S 6 8 0で、送信器6 0 9が古いパケットの送信でビジーであり、したがって新たなパケットの送信のために利用可能でない場合、制御はS 5 1 0に戻る。そうでない場合、S 7 0 0で、サービス後処理のために待機しているサービスされたばかりのフローの利用可能性がチェックされる。S 7 0 0で、サービスされたばかりのフローが利用可能でない場合、S 7 1 0で、バックログフローがあるかどうか判定される。バックログフローが存在しない場合、制御はS 5 1 0に戻り、存在する場合、S 7 2 0で、フローポインタのFIFOキュー5 0 5の先頭5 0 6にあるフロー5 0 1がサービスのために選択され、S 7 3 0で、選択されたフローのフローキュー5 0 2内の最初のデータパケットが送信器6 0 9に送られる。S 7 4 0で、サービスのために選択されたフロー5 0 1のフレームフラグがグローバルフレームカウンタFRMCNTに等しい場合、制御はS 5 1 0に戻り、そうでない場合、S 7 5 0で、グローバルフレームカウンタFRMCNTはトグルされ、制御はS 5 1 0に戻る。

【0 0 6 0】

S 7 0 0で、サービス後処理のために待機しているサービスされたばかりのフロー5 0 1が利用可能である場合、コントローラはS 7 6 0に進み、対応するフローキュー5 0 2の長さをデクリメントする。S 7 7 0で、サービスされたばかりのフロー5 0 1のバンドル4 0 1が、フロー別状態テーブル5 0 3のバンドルポインタを用いて識別される。S 7 8 0で、グローバルフレームカウンタFRMCNTが、サービスされたばかりのフロー5 0 1のバンドル4 0 1の開始フラグに等しいかどうか判定される。開始フラグがグローバルフレームカウンタに等しくない場合、コントローラは、（S 7 9 0で）バンドルの累積シェアをバンドルの現在シェアに等しいとセットし、（S 8 0 0で）バンドルの開始フラグをグローバルフレームカウンタに等しいとセットする。その後、制御は、フロータイムスタンプの更新のためにS 8 1 0に進む。S 7 8 0で、サービスされたばかりのフローのバンドルの開始フラグがグローバルフレームカウンタに等しい場合、サービスされたばかりのフローのタイムスタンプの更新のために、制御は直接S 8 1 0に移る。

【0 0 6 1】

S 8 2 0で、コントローラ6 0 3は、サービスされたばかりのフローのキュー長がゼロに等しいかどうかを判定する。S 8 2 0で、サービスされたフロー5 0 1のフローキュー5 0 2が空であると判定された場合、S 8 8 0で、コントローラは、サービスされたばかりのフローのタイムスタンプがグローバル状態テーブル5 0 7の基準タイムスタンプ増分以上であるかどうかをチェックする。タイムスタンプが基準タイムスタンプ以上である場合、制御はS 8 9 0に移り、タイムスタンプは正当な範囲（0, T_Q）内にリセットされる。その後、制御はS 9 0 0に移る。S 8 8 0で、サービスされたばかりのフローのタイムスタンプが基準タイムスタンプ増分より小さいと判定された場合、制御は直接S 9 0 0に移る。S 9 0 0で、サービスされたばかりのフローへのポインタが、フローポインタのFIFOキュー5 0 5の先頭5 0 6から取り出される。その後、S 9 1 0で、サービスされたばかりのフローのバンドルの現在シェアが、アイドルになった、サービスされたばかりのフローの帯域割当てだけ減少させられる。S 9 2 0で、コントローラ6 0 3は、バックログフローの総数をデクリメントした後、S 7 1 0に進む。

【0 0 6 2】

S 8 2 0で、サービスされたフロー5 0 1のフローキュー5 0 2が空でない判定された場合、S 8 3 0で、コントローラは、サービスされたばかりのフローのタイムスタンプが

10

20

30

40

50

グローバル状態テーブル507の基準タイムスタンプ増分以上であるかどうかをチェックする。タイムスタンプが基準タイムスタンプより小さい場合、制御はS840に移り、サービスされたばかりのフローのフレームフラグがトグルされる。その後、S850で、サービスされたばかりのフローのタイムスタンプは正当な範囲(0, T_Q)内にリセットされる。S860で、コントローラ603は、サービスされたばかりのフローへのポインタを、フローポインタのFIFOキュー505の先頭506から取り出す。S870で、同じポインタが、FIFOキュー505の末尾506に再びキューイングされる。その後、制御はS710に移る。S830で、サービスされたばかりのフローのタイムスタンプが基準タイムスタンプ増分以上であると判定された場合、制御は直接S710に移る。

【0063】

図9のリストは、本発明に従ってバンドルIのフローiを選択しサービス後処理する方法を擬似コード形式で記述したものである。

【0064】

上記の実施例は、本発明に従って、重み付けラウンドロビン上にソフトスケジューリング階層を重ね合わせるために使用可能な諸原理の単なる例示である。当業者であれば、ここに明示的に記載していなくとも、特許請求の範囲により規定される本発明の技術思想および技術的範囲から離れずに、それらの原理を実現するさまざまな構成を考えることが可能である。

【0065】

[付録]

文献

[1] R. Braden, D. Clark, and S. Shenker, "Integrated Services in the Internet Architecture: an Overview", Request for Comments (RFC) 1633, IETF, June 1994.

[2] K. Nichols, V. Jacobson, and L. Zhang, "A Two-bit Differentiated Services Architecture for the Internet", Request for Comments (RFC) 2638, IETF, July 1999.

[3] D. Stiliadis and A. Varma, "Design and Analysis of Frame-based Fair Queueing: A New Traffic Scheduling Algorithm for Packet-Switched Networks", Proceedings of ACM SIGMETRICS '96, pp.104-115, May 1996.

[4] D. Stiliadis and A. Varma, "Efficient Fair Queueing Algorithms for Packet-Switched Networks", IEEE/ACM Transactions on Networking, Vol.6, No.2, pp.175-185, April 1998.

[5] J. C. R. Bennett and H. Zhang, "WF2Q: Worst-Case-Fair Weighted Fair Queueing", Proceedings of IEEE INFOCOM '96, pp.120-128, March 1996.

[6] D. Stiliadis and A. Varma, "A General Methodology for Designing Efficient Traffic Scheduling and Shaping Algorithms", Proceedings of IEEE INFOCOM '97, Kobe, Japan, April 1997.

[7] I. Stoica, H. Zhang, and T. S. E. Ng, "A Hierarchical Fair Service Curve Algorithm for Link-Sharing, Real-Time, and Priority Services", Proceedings of ACM SIGCOMM '97, September 1997.

[8] D. C. Stephens, J. C. R. Bennett, and H. Zhang, "Implementing Scheduling Algorithms in High-Speed Networks", IEEE Journal on Selected Areas in Communications, Vol.17, No.6, June 1999, pp.1145-1158.

[9] M. Katevenis, S. Sidiropoulos, and C. Courcoubetis, "Weighted Round Robin Cell Multiplexing in a General-Purpose ATM Switch", IEEE Journal on Selected Areas in Communications, vol.9, pp.1265-79, October 1991.

[10] M. Shreedhar and G. Varghese, "Efficient Fair Queueing Using Deficit Round Robin", IEEE/ACM Transactions on Networking, vol.4, no.3, pp.375-385, June 1996.

[11] H. Adishesu, G. Parulkar, and G. Varghese, "A Reliable and Scalable Str

10

20

30

40

50

iping Protocol", Proceedings of ACM SIGCOMM '96, August 1996.

[1 2] J. C. R. Bennett and H. Zhang, "Hierarchical Packet Fair Queueing Algorithms", Proceedings of ACM SIGCOMM '96, pp.143-156, August 1996.

【 0 0 6 6 】

【 発明の効果 】

以上述べたごとく、本発明によれば、基本スケジューラの単純さを損なわずに階層的帯域分離が達成される。

【 0 0 6 7 】

特許請求の範囲の発明の要件の後に括弧で記載した番号がある場合は、本発明の一実施例の対応関係を示すものであって、本発明の範囲を限定するものと解釈すべきではない。

10

【 図面の簡単な説明 】

【 図 1 】 データソース、通信スイッチ、およびデータ宛先を含む例示的なパケットネットワークを示す図である。

【 図 2 】 図 1 のパケットネットワークで用いられる例示的な通信スイッチを示す図である。

【 図 3 】 図 3 A は、新たなパケットがフローキューの先頭に到着したときに不足ラウンドロビン (D R R) アルゴリズムで用いられる擬似コードの例を示し、図 3 B は、サーバが 1 パケットの送信を完了したときに余剰ラウンドロビン (S R R) アルゴリズムで用いられる擬似コードの例を示す図である。

【 図 4 】 本発明による、スケジューラの 2 層論理編成を示す図である。

20

【 図 5 】 本発明のスケジューラによって利用されるキュー、状態テーブル、レジスタ、およびパラメータの機能図である。

【 図 6 】 図 5 の装置の具体的実装の例示的ブロック図である。

【 図 7 】 本発明に従ってパケットの送信をスケジューリングする方法を記述した例示的な流れ図である。

【 図 8 】 本発明に従ってパケットの送信をスケジューリングする方法を記述した例示的な流れ図である。

【 図 9 】 1 パケットの送信を完了した後にスケジューラによって使用される擬似コードの例を示す図である。

【 符号の説明 】

30

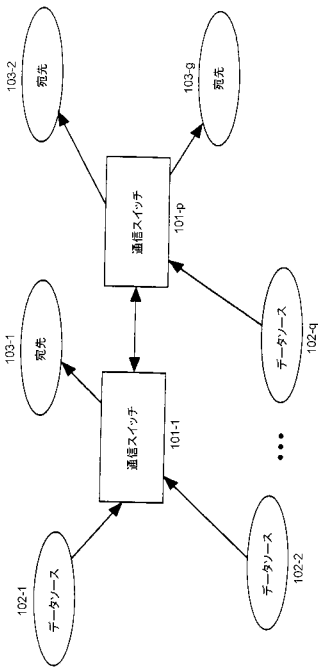
- 1 0 1 通信スイッチ
- 1 0 2 データソース
- 1 0 3 宛先
- 2 0 0 通信リンクインタフェース
- 2 0 1 入力リンク
- 2 0 2 出力リンク
- 2 5 0 スイッチファブリック
- 4 0 1 バンドル
- 5 0 1 データパケットフロー
- 5 0 2 フローキュー
- 5 0 3 フロー別状態テーブル
- 5 0 4 バンドル別状態テーブル
- 5 0 5 フローポインタの F I F O キュー
- 5 0 6 レジスタ
- 5 0 7 グローバル状態テーブル
- 6 0 1 パケット受信器
- 6 0 2 スケジューラ
- 6 0 3 コントローラ
- 6 0 4 バンドル別状態 R A M
- 6 0 5 レジスタ

40

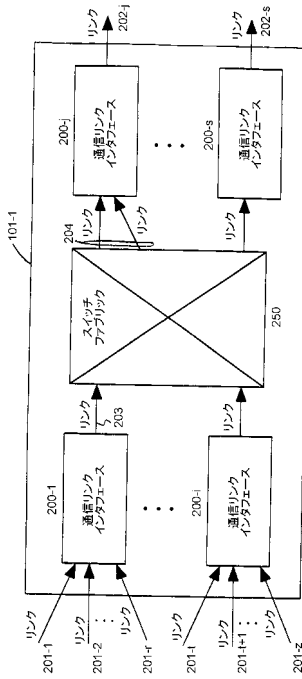
50

- 606 チップ
- 607 パケットRAM
- 608 フロー別状態RAM
- 609 パケット送信器

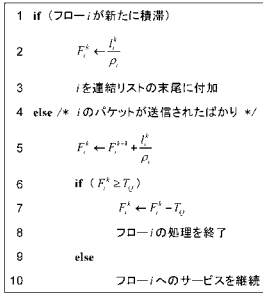
【図1】



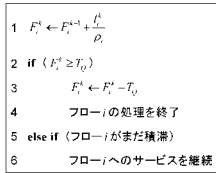
【図2】



【 図 3 】

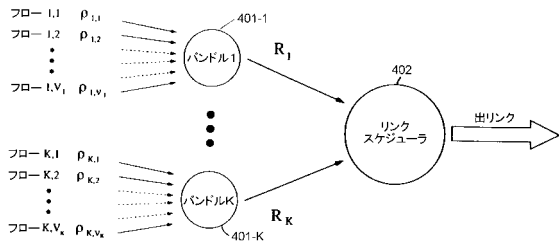


A

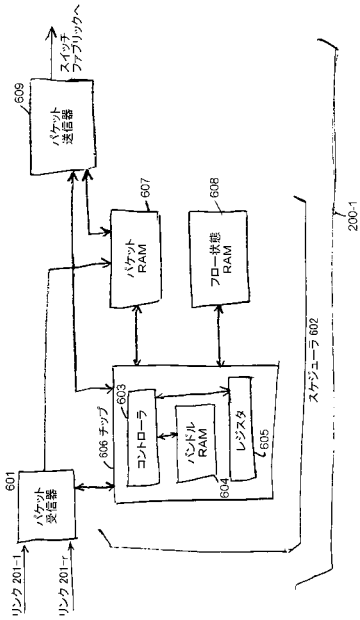


B

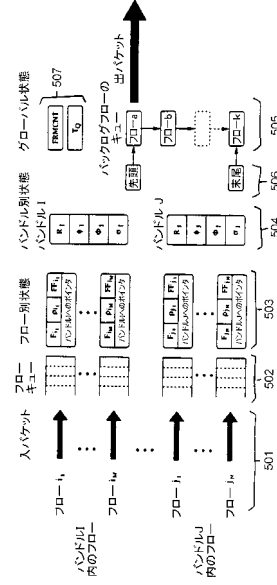
【 図 4 】



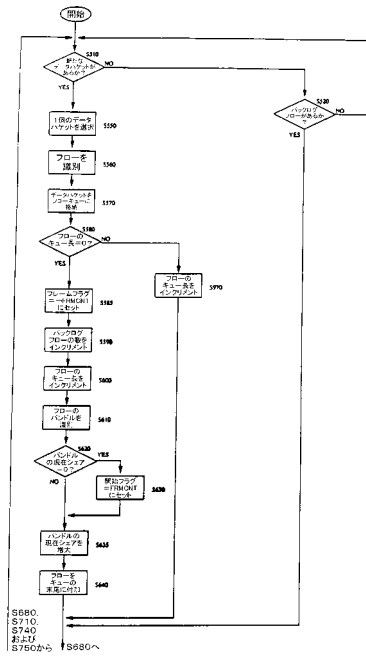
【 図 6 】



【 図 5 】



【 図 7 】



フロントページの続き

- (74)代理人 100091889
弁理士 藤野 育男
- (74)代理人 100101498
弁理士 越智 隆夫
- (74)代理人 100096688
弁理士 本宮 照久
- (74)代理人 100102808
弁理士 高梨 憲通
- (74)代理人 100104352
弁理士 朝日 伸光
- (74)代理人 100107401
弁理士 高橋 誠一郎
- (74)代理人 100106183
弁理士 吉澤 弘司
- (74)代理人 100081053
弁理士 三俣 弘文
- (72)発明者 ファビオ エム チウスイ
アメリカ合衆国、07740 ニュージャージー州、ロング ブランチ、Apt. 4D、オーシャ
ン アベニュー N 300
- (72)発明者 ロバート ティー クランシー
アメリカ合衆国、06437 コネチカット州、ギルフォード、ホーソン ロード 113
- (72)発明者 ケヴィン ディー ドラッカー
アメリカ合衆国、08801 ニュージャージー州、アナンデル、ウェストゲート ドライブ
18
- (72)発明者 アンドレーエ フランシオン
アメリカ合衆国、07701 ニュージャージー州、レッド バンク、ウェスト フロント スト
リート 537
- (72)発明者 ナサー イー イディ
アメリカ合衆国、07724 ニュージャージー州、イートンタウン、ウェッジウッド サークル
65

審査官 衣鳩 文彦

- (56)参考文献 特開2000-31997(JP,A)
F M Chiussi et al. , A distributed scheduling architecture for scalable packet switches
 , IEEE Journal on Selected Areas , 米国 , IEEE , 2000年12月 , Volume:18 Issue:12 , 第
2665 - 2683頁
Kai-Yeung Siu et al. , Virtual queueing techniques for UBR+ service in ATM with fair ac
cess and minimum bandwidth guarantee , Global Telecommunications Conference, 1997. GLOB
ECOM '97 , 米国 , IEEE , 1997年11月 , VOL: 2 , 第1081 - 1085頁

- (58)調査した分野(Int.Cl. , DB名)
H04L 12/56