

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第3715098号
(P3715098)

(45) 発行日 平成17年11月9日(2005.11.9)

(24) 登録日 平成17年9月2日(2005.9.2)

(51) Int. Cl.⁷

H04L 12/56

F I

H04L 12/56 200F

請求項の数 6 (全 9 頁)

(21) 出願番号	特願平10-27573	(73) 特許権者	596092698
(22) 出願日	平成10年2月9日(1998.2.9)		ルーセント テクノロジーズ インコーポ
(65) 公開番号	特開平10-313324		レーテッド
(43) 公開日	平成10年11月24日(1998.11.24)		アメリカ合衆国, 07974-0636
審査請求日	平成11年7月8日(1999.7.8)		ニュージャージー, マレイ ヒル, マウン
審査番号	不服2002-19979 (P2002-19979/J1)		テン アヴェニュー 600
審査請求日	平成14年10月11日(2002.10.11)	(74) 代理人	100064447
(31) 優先権主張番号	60/037844		弁理士 岡部 正夫
(32) 優先日	平成9年2月7日(1997.2.7)	(74) 代理人	100085176
(33) 優先権主張国	米国 (US)		弁理士 加藤 伸晃
(31) 優先権主張番号	08/972424	(74) 代理人	100106703
(32) 優先日	平成9年11月18日(1997.11.18)		弁理士 産形 和央
(33) 優先権主張国	米国 (US)	(74) 代理人	100081053
			弁理士 三俣 弘文

最終頁に続く

(54) 【発明の名称】 通信ネットワークにおけるパケットの配送装置とその方法

(57) 【特許請求の範囲】

【請求項1】

(A) それぞれのソースからパケットを受領し、宛先に配送する前にこの受領したパケットを一時的に記憶する複数の結線毎の待ち行列(20)と、

(B) 保証されたあらかじめ割り当てられたレートをもって、前記(A)複数の結線毎の各待ち行列(20)からのパケットをサービスする重み付き公平待ち行列スケジューリング手段(40)と、

(C) 前記待ち行列内のパケットの存否を検出する手段と、

(D) 前記パケットの不存在を検出すると、その結果発生する余剰のバンド幅を再分配する状態依存性スケジューリング手段(200)と、

からなり、

前記(D)状態依存性スケジューリング手段(200)は、前記待ち行列(20)の性能特性に応じた状態に従って、これらの待ち行列(20)にサービスし、

その結果、それぞれの待ち行列(20)のパケットをルーティングする遅延特性と分離特性が保存される

ことを特徴とする通信ネットワークにおけるパケットの配送装置。

【請求項2】

前記待ち行列(20)は、最大量のパケットを有する待ち行列を含み、前記(D)状態依存性スケジューリング手段(200)は、最長の待ち行列に最初にサービスする

ことを特徴とする請求項1の装置。

【請求項 3】

前記待ち行列(20)は、サービスされる為に最も長時間待機している最大遅延のペケットを有する待ち行列を含み、
前記(D)状態依存性スケジューリング手段(200)は、最大遅延の待ち行列に最初にサービスすることを特徴とする請求項1の装置。

【請求項 4】

前記待ち行列(20)は、そのバッファメモリがオーバーフローする可能性が最も大きい待ち行列を含み、
前記(D)状態依存性スケジューリング手段(200)は、バッファがオーバーフローする最も可能性の高い待ち行列に最初にサービスすることを特徴とする請求項1の装置。

10

【請求項 5】

前記接続に対し、最悪の場合の公平尺度が満足されていることを特徴とする請求項1の装置。

【請求項 6】

(A)それぞれのソースからペケットを受領し、宛先に配送する前に、この受領したペケットを一時的に記憶する待ち行列(20)を形成するステップと、
(B)保証されたあらかじめ割り当てられたレートをもって前記複数の結線毎の待ち行列(20)からのペケットにサービスするステップと、
(C)前記待ち行列内のペケットの存否を検出するステップと、
(D)前記ペケットの不存在を検出すると、その結果発生する余剰のバンド幅を再分配するステップと、

20

からなり、

前記(D)再分配ステップは、前記待ち行列の性能特性に応じた状態変化に従ってこれらの待ち行列にサービスし、
それぞれの待ち行列のペケットをルーティングする遅延特性と分離特性が保存されることを特徴とする通信ネットワークにおけるペケットの配送方法。

【発明の詳細な説明】**【0001】**

30

【発明の属する技術分野】

本発明は、ペケット通信システムに関し、特に、このペケット通信システムのルータと交換機内で実行される公平待ち行列システム(fair-queuing system)に関する。

【0002】**【従来の技術】**

沢山の研究が、理想的な流れ(fluid flow)モデル(即ち複数のソースから送信されるデータペケットは、無限に分解できるものとし、かつ複数のソースは、そのデータを例えば1本の物理的通信リンクに同時に送信できるようなモデル)をできるだけエミュレート(模擬)するペケット通信ネットワーク用の待ち行列システムの開発に費やされている。ペケットを無限に分割できることは、実際にはありえないことである。

40

【0003】

ペケットネットワークにおいては、通常ペケットがリンクを介して伝送されると、すべてのペケットを伝送しなければならない、即ちその間に別のペケットの伝送を阻止することはできない。ペケットネットワーク内で、サービスの質(Quality of Service(QoS))を保証するような要望があるために、データペケット交換器、あるいはルータ内でトラフィックスケジューリング方法の実現が必要とされている。

【0004】

このスケジューリング方法の機能とは、交換器の各出力用リンクにとってその出力リンクを共有する通信期間に属するペケットの中から次のサイクルに送信すべきペケットを選択することである。このような選択は、個々のトラフィック期間に対し、保証されるサービ

50

スの質 (QoS) (例えば、最大遅延の上限) を満足するように行う必要がある。この方法の実行は、ハードウェアあるいはソフトウェアで行われるが、速度のことを考慮するとスケジューリングは通常 ATM 交換器および高速ルータ内のハードウェア内で実行される。

【0005】

多くの様々な種類のスケジューリング方法が、汎用プロセッサシェアリング (Generalized Processor Sharing (GPS)) として知られる理論的スケジューリング原理を近似すべく提案されている。このスケジューリング原理は、「流体」モデルを用いて定義される原理である。このような GPS によりリンク上で通信できる各期間に割り当てられたバンド幅をきっちりと制御できる。しかし、ある期間に伝送されるパケットは、これ以上分割することができないので複数のソースからのデータは、パケットの間 (境界) にインターリーブ (挟み込む) しなければならない。このため GPS 原理は、実際のパケット交換ネットワークでは、実行することができない。

10

【0006】

単純な FIFO 技術, ラウンドロビン技術, 公平待ち行列技術により、個々の待ち行列をサービスすることは公知である。しかし、この流体システムを最もよく近似する重み付けされた公平待ち行列 (Weighted fair-queuing (WFQ)) 系が開発されている。特に、A. Demers, S. Keshav, S. Shenker 著の "Analysis and Simulation of a Fair Queuing Algorithm" Internetworking: Research and Experience, pp. 3-26, vol. 1, 1990 の文献によれば、流体フロー GPS システムをシミュレートし、この GPS システム内の出発の順番に対し、パケットのスケジューリング決定にバイアスを掛けることにより GPS をエミュレートする公平待ち行列系が記載されている。

20

【0007】

重み付け公平待ち行列においては、スケジューリング方法により制御される出力リンクを共有する各トラフィック期間 i に、その期間の予約バンド幅に対応する ϕ_i が割り当てられている。この値 ϕ_i は、リンク上の期間 i の予約バンド幅が次式で与えられるように計算される。

【数1】

$$\frac{\phi_i}{\sum_{j=1}^v \phi_j}$$

30

ここで、分母はリンクを共有する全ての 期間の ϕ_i の値の和を計算している。

【0008】

図1に従来技術に係る重み付け公平待ち行列 WFQ システム 10 を示す。この重み付け公平待ち行列 WFQ システム 10 は、複数の結線用 (結線毎に用意される) 行列 $20a, \dots, 20i$ を有し、各待ち行列はデータ端末のようなソースからの入力トラフィックである情報パケットを、共有メモリ 25 内の別々の場所にパケットを一時的に記憶する。様々な種類のトラフィック、例えばオーディオ、ビデオ、データ等を収納する様々な種類の待ち行列が存在する。

40

【0009】

さらにまた、重み付け公平待ち行列 WFQ システム 10 は、この待ち行列からパケットを割り当てられたレートと全く同一のレートでもって重み付け公平待ち行列スケジューラであるレート比例型サーバー 40 に転送するシェーパ 30a, ..., 30i を具備する。この重み付け公平待ち行列スケジューラは、重みはそれぞれ各結線用待ち行列 $20a, \dots, 20i$ に関連づけられ、その結果スケジューラによりこれらの待ち行列の1つに提供されるサービスは、待機中のパケットがある場合には、常に重みに比例することになる。

【0010】

例えば、リンクの容量 (バンド幅) を $C = 10$ パケット / 秒と仮定する。同時にスケジュー

50

ーラも3つの待ち行列にサービスするものとし、待ち行列Q1は重み $WQ1 = 20\%$ 、待ち行列Q2は重み $WQ2 = 30\%$ 、待ち行列Q3は重み $WQ3 = 50\%$ とする。すると、この3個の全ての待ち行列が待機パケットを有している場合には、Q1は2パケット/秒の保証されたバンド幅を、Q2は3パケット/秒の保証されたバンド幅を、Q3は5パケット/秒の保証されたバンド幅を受領する。

【0011】

しかし、例えばQ3が待機中のパケットを持っていない場合には、過剰(余分)のバンド幅は5パケット/秒に等しくなる。WFQシステムにおいては、この過剰バンド幅は、待機中のパケットを有する待ち行列の重みに比例して再分配される。上記の実施例では、Q3は待機パケットを有していないので過剰のバンド幅はQ1とQ2に比例して分配され、その結果Q1は4パケット/秒、Q2は6パケット/秒の瞬間的なバンド幅を受領することになる。それぞれのシェーパ-30から出ていった各パケットは、直接、レート比例型サーバー40に転送され、このレート比例型サーバー40はこのパケットをさらに出力リンク50に転送する。このレート比例型サーバー40は重み付け公平待ち行列の変形例である。

10

【0012】

このようなWFQ系においては、端末間の遅延が保証されるような利点がある。例えば、各パケットはストリーム内の各パケットフローに対しあるレートが保証され、ストリーム間の隔離が確保され、例えば誤動作をするソースは他のストリームのフローには影響を及ぼさない。さらにまた、容量の利用が十分でない場合には、例えばフローが特にバースト型であり、送信しない空き時間が存在する場合には、このWFQシステムは作業保存特性(work-conservation property)を保存するために未使用のバンド幅の再分配を容易にする。

20

【0013】

現在複数の待ち行列の間で未使用のバンド幅の容量を再分配する特性は、流体フローモデル(fluid-flow model)に特有な方法、例えば特定の待ち行列の重みに従って実行される。かくしてパケットの待ち行列がないアイドル状態の時には、過剰の(余分)のバンド幅は長期の要件に基づくその重みに比例して予備接続(backlogged connections)に再分配される。

【0014】

全ての公平待ち行列システムがGPSの近似を受け継ぐが、このGPSの欠点は、GPSが状態依存性バンド幅の共有を厳しく制限してしまうことである。GPSにおける状態の依存性のみが予備接続の数になる。従ってこれ以上の自由度は存在せず、接続の長期間のニーズに基づいて設定された保証レートによってのみ決定される。

30

【0015】

バンド幅の共有に対するこの制限は、公平待ち行列の特性を維持するため、およびリーキーパケット制御のトラフィックソースに対する最悪の遅延限界を保証するためには、必要以上に厳しいものである。したがって、GPSの余剰バンド幅共有を次善策としてエミュレートするような公平待ち行列システムに対するニーズは存在しない。

【0016】

【発明が解決しようとする課題】

したがって本発明の目的は、GPSをエミュレートする重み付け公平待ち行列システムにおいて、残りのトラフィックフローの瞬時のニーズを反映するような状態に依存する方法で未使用のバンド幅の再分配を効率よく達成する方法を提供することである。

40

【0017】

【課題を解決するための手段】

本発明は、適応型再分配系を実行する重み付け公平待ち行列への改良型アプローチである。このような再分配系においては、接続毎の各フローは、適応しながら再分配される余剰のバンド幅とリンクバンド幅の特定の共有が保証される。この再分配系により、最悪の場合の端末間の遅延限界を与えるような公平な待ち行列機能を保存でき、そしてこの再分配

50

系は、余剰のバンド幅が存在しない場合でも公平待ち行列のように働く。この余剰のバンド幅は個々の条件に従って再分配される。

【 0 0 1 8 】

状態依存性の条件の例は、1) 現在最も長い遅延を有するフローにサービスする最長遅延優先条件 (Longest delay first (LDF)) と、2) 許容された最大遅延と現在の遅延との差が最小のフローにサービスする対オーバーフロー最短時間条件 (Least time to overflow (LTO)) と、3) 最悪のケースの到着が発生したときに、バッファが最初にオーバーフローを引き起こすようなフローにサービスするリーキーバケットのオーバーフローへの最短時間条件 (Least time to overflow with leaky buckets (LTO-LB)) である。

10

【 0 0 1 9 】

この適応型バンド幅再分配系は、常に各接続に必要とされる最小限の保証を与え、最悪の場合の公平さを満足させている。

【 0 0 2 0 】

このLDF方式は、余剰のバンド幅を用いて遅延分配の変動を低減し、これにより音声ソース、ビデオソースに対しプレイアウト(使い切った)バッファサイズを縮小できるメリットがある。ビデオトレースとボイストラフィックのシミュレーションによれば、この方式は最悪の場合の保証を犠牲にすることなくGPSよりは、性能がよい。

【 0 0 2 1 】

許容最大遅延からの変動を考慮にいれていないので、遅延制限が小さいフロー(例、音声)は、遅延制限が大きいフローの存在の元では、余剰のバンド幅をほとんど得ることができない。割当重みの不正確さは、大きな遅延制限のフローよりは、これらの小さな遅延制限のフローの方がはるかに大きな損失を受けることになる。

20

【 0 0 2 2 】

LTO方式は、最も早くオーバーフローする可能性のあるフローが、最も瞬時のバンド幅を必要とするという仮定の元で、余剰バンド幅を割り当てることによりパケットの喪失を最小にしようとしている。このようにする際に、許容最大遅延からの各フローの現行の変動を考慮にいれている。CBRと音声ソースとビデオソースを異なる遅延限界で組み合わせると、この方式は全てのクラスに対しパケットの喪失を低減できかつ各クラスの遅延の変動を低減することができる。

30

【 0 0 2 3 】

【 発明の実施の形態 】

図2は適応型余剰バンド幅再分配を用いた公平待ち行列方法を表すブロック図。本発明のシステムは、図1に示した従来のシステムに加えて状態依存性スケジューラ200を有し、この状態依存性スケジューラ200へは各結線用待ち行列20a, ..., 20iからパケットが入力され、状態依存性スケジューラ200の出力はマルチプレクサ350で出力リンク50と多重化される。パケットは一時的に結線用待ち行列20a, ..., 20i内に記憶される。重みがこれらの各待ち行列に割り当てられ、この重みは、各待ち行列に到着したパケットに割り当てられなければならない出力バンド幅の一部を表す。

【 0 0 2 4 】

シェーパ30a, ..., 30iは、結線用待ち行列20a, ..., 20iからのパケットを割り当てられたレートと等しいレートでレート比例型サーバ40に転送する。このような構成により、余剰のバンド幅は、再分配用に利用できることになる。本発明によれば、他のソース即ち他の待ち行列からのトラフィックは、適応型バンド幅再分配メカニズムである状態依存性スケジューラ200にも入力される。結線毎の待ち行列、シェーパ、スケジューラ等のハードウェア構成は、当業者には公知である。

40

【 0 0 2 5 】

上述したようにシェーパ30a, ..., 30iは、パケットを割り当てられたレートに等しいレートで、パケットを状態依存性スケジューラであるレート比例型サーバ40に送る。シェーパ30a, ..., 30iよりそれぞれの待ち行列の接続「i」に時間間隔

50

の間提供されるサービスは、 $S_i(\cdot, t)$ で表す。パケットはシェーパ-30から無限の容量を有するサーバ-40に転送されるものとする。レート比例型サーバ-40により提供されるサービスは、 $R_i(\cdot, t)$ で表す。

【0026】

サービスを受ける資格のなくなったパケットは、シェーパ-内の対応する接続待ち行列内に保存され、一方全ての資格のあるパケットは、レート比例型サーバ-40内でサービスを受けるために待機する。サービスはパケットがそこにある限りレート比例型サーバ-40から与えられる。状態依存性スケジューラ200の待ち行列の全てが空の時には、状態依存性スケジューラ200は結線用待ち行列 $20a, \dots, 20i$ からパケットを選択するよう起動され、シェーパ-30の状態に影響を及ぼさない接続に提供されるサービスでもって伝送する。図2に示すように、シェーパ-30によりタイムインターバルの間待ち行列接続「 i 」へ提供されるサービスは、 $D_i(\cdot, t)$ で示し、状態依存性スケジューラ200により提供されるサービスは、 $Z_i(\cdot, t)$ で示す。

10

【0027】

シェーパ-30は当業者に公知のカレンダー（暦）待ち行列を用いて実現できる。シェーピング用のカレンダー待ち行列のメカニズムは、D. Stiliadis と A. Verma 著の A General Methodology for Designing Scheduling and Shaping Algorithms, in Proceedings of IEEE INFOCOM '97 を参照のこと。

【0028】

重み付き公平待ち行列サーバであるレート比例型サーバ-40は、公知の公平待ち行列メカニズムを用いて実現できる、これに関しては、D. Stiliadis と A. Varma の Traffic Scheduling System and Method for Packet-Switched Networks, 米国特許出願第08/634,904号（1996年4月15日提出）を参照のこと。

20

【0029】

ある時点でRPSスケジューラ200内にパケットが存在しないと決定された場合には、これはフリーの（だれでもが使用可能な）バンド幅があることを意味する。このためパケットは、状態依存の方法に基づいて状態依存性スケジューラ200によりサービスされる。この状態依存の方法により、決定はシステムに関連するある変数の現在の状態に基づいて行われる。状態依存の決定の2つの例は、最長遅延優先（Longest Delay First）とオーバフローへの最短時間（Least Time to Overflow）である。

30

【0030】

最長遅延優先方式においては、状態依存性スケジューラ200は結線用待ち行列 $20a, \dots, 20i$ の中からこの保証されたのと等しいレートでもってこれらの待ち行列がサービスされた場合には、遅延が最長となるパケットを有する待ち行列を選択する。これを実行するために、スケジューラは各待ち行列に保証レートにより除算された待ち行列のサイズに等しい遅延値を割り当てる。この値は、この待ち行列の最後のパケットが受ける遅延を表す。状態依存性スケジューラ200は、一組の数の中から最大値を選択するために何らかのメカニズムを用いて、最大値を有する待ち行列を選択する。この選択メカニズムは公知である。

【0031】

オーバフローへの最短メカニズムにおいては、状態依存性スケジューラ200は結線用待ち行列 $20a, \dots, 20i$ の中から最短時間でオーバフローする可能性のある待ち行列を選択する。このメカニズムは、最大サイズは待ち行列に関係していると見なしている。パケットが到着したり、あるいは待ち行列によりサービスされる毎にオーバフローするパケット（packet-to-overflow）の変数は、最大容量を超えることなく待ち行列に付加されるパケットの数をカウントする。

40

【0032】

第2の変数（オーバフローする時間（time-to-overflow）と称する）は、オーバフローする予測時間を表し、関連レートでオーバフローするパケット変数（packets-to-overflow）を除算することにより計算される。状態依存性スケジューラ200は、最小のオーバ

50

フローする時間変数 (time-to-overflow) を有する待ち行列を伝送用を選択する。この最小値は、公知のメカニズムにより決定できる。

【 0 0 3 3 】

上記の変数は、状態依存性の変数を計算する 2 種類の方法について説明したが、他の類似の方法も用いることができる。状態依存性の変数に基づいて、状態依存性スケジューラ 2 0 0 は、選択された待ち行列からパケットをリンク 7 5 を介して、マルチプレクサ 3 5 0 により多重化されて次の宛先に転送される。

【 0 0 3 4 】

最悪の場合の公平性の特徴は、本発明のシステムで満足されることが分かる。最悪の場合の公平性においては、異なる割当レートで異なる長さの待ち行列をサービスすることは、第 1 の待ち行列からのパケットが時間 t_1 でサービスされ、第 2 の待ち行列のパケットが時間 t_2 でサービスされるように、インタリーブされる (間に入れられる)。

【 0 0 3 5 】

本発明においては、最悪の場合の公平性は満足され、その結果、最悪の場合の時間間隔 $t_2 - t_1$ は、接続の数の関数ではない値により制限されるか、あるいはそれよりも少なく、しかし時間 t においてサービスされる待ち行列のような最長の待ち行列のパケットサイズと、最長の待ち行列の割当てレートの関数である。

【 0 0 3 6 】

最悪の場合の公平性の特徴は、本発明の重み付け公平待ち行列システム内の状態依存性スケジューラ 2 0 0 により満足される。

【 図面の簡単な説明 】

【 図 1 】 重み付け公平待ち行列とレート比例型サーバーの特徴を表すデータフロー図

【 図 2 】 本発明による重み付け公平待ち行列スケジューリングシステムを表す図

【 符号の説明 】

- 1 0 従来技術に係る重み付け公平待ち行列 W F Q システム
- 2 0 結線用待ち行列
- 2 5 共有メモリ
- 3 0 シェーパ
- 4 0 レート比例型サーバー
- 5 0 出力リンク
- 7 5 リンク
- 1 0 0 本発明に係る重み付け公平待ち行列 W F Q システム
- 2 0 0 状態依存性スケジューラ
- 3 5 0 マルチプレクサ

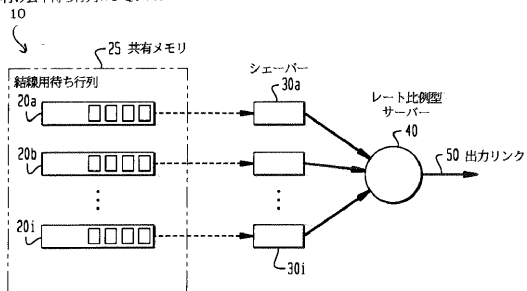
10

20

30

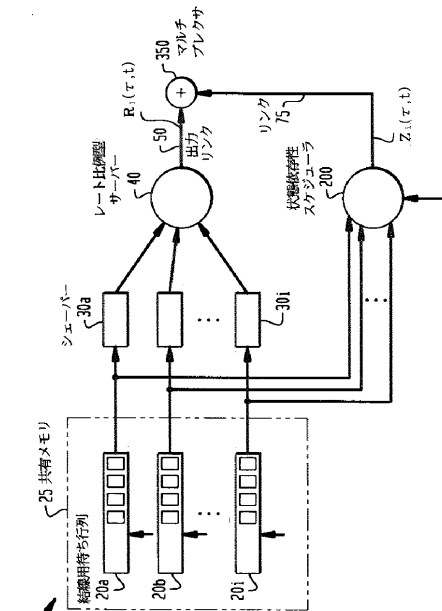
【 図 1 】

従来技術に係る
重み付け公平待ち行列WFQシステム



【 図 2 】

本発明に係る
重み付け公平待ち行列WFQシステム



フロントページの続き

- (72)発明者 ニコラス ジー . ダフォールド
アメリカ合衆国、07030 ニュージャージー、ホボーケン、#50、パークアベニュー 93
3
- (72)発明者 ティ、ブイ、ラクシマン
アメリカ合衆国、07724 ニュージャージー、イートンタウン、ヴィクトリアドライブ 18
8
- (72)発明者 ディミトリオス ステリアリス
アメリカ合衆国、07030 ニュージャージー、ホボーケン、パークアベニュー 106

合議体

審判長 山本 春樹

審判官 野元 久道

審判官 衣鳩 文彦

(56)参考文献 特開平10-51472(JP,A)

(58)調査した分野(Int.Cl.⁷, DB名)
H04L12/56