

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5521038号
(P5521038)

(45) 発行日 平成26年6月11日(2014.6.11)

(24) 登録日 平成26年4月11日(2014.4.11)

(51) Int.Cl.		F I	
HO 4 L 12/801	(2013.01)	HO 4 L 12/801	
HO 4 L 12/70	(2013.01)	HO 4 L 12/70	1 O O Z
HO 4 L 12/28	(2006.01)	HO 4 L 12/28	2 O O D
HO 4 M 3/00	(2006.01)	HO 4 M 3/00	B

請求項の数 12 (全 17 頁)

(21) 出願番号	特願2012-518078 (P2012-518078)	(73) 特許権者	391030332
(86) (22) 出願日	平成22年6月29日 (2010.6.29)		アルカテルルーセント
(65) 公表番号	特表2012-531867 (P2012-531867A)		フランス国、75007・パリ、アブニ
(43) 公表日	平成24年12月10日 (2012.12.10)		ユ・オクターブ・グレアール、3
(86) 国際出願番号	PCT/EP2010/059163	(74) 代理人	110001173
(87) 国際公開番号	W02011/000810		特許業務法人川口国際特許事務所
(87) 国際公開日	平成23年1月6日 (2011.1.6)	(72) 発明者	ラウテンシユラガー、ボルフラム
審査請求日	平成24年2月24日 (2012.2.24)		ドイツ国、74343・ザクセンハイム、
(31) 優先権主張番号	09290501.7		レッシングシユトラーセ・19
(32) 優先日	平成21年6月29日 (2009.6.29)		
(33) 優先権主張国	欧州特許庁 (EP)	審査官	衣鳩 文彦

最終頁に続く

(54) 【発明の名称】 トラフィック負荷を管理する方法

(57) 【特許請求の範囲】

【請求項1】

パケット交換網(100)において多数のアプリケーションストリームから集約されたパケットトラフィックを受信するネットワークノード(2)のトラフィック負荷を管理する方法であって、

a) ネットワークノード(2)の伝送容量Bに適合するアプリケーションストリームの最大数であるパケットトラフィックの粒度を推定するステップと、

b) 推定された粒度およびネットワークノード(2)のトラフィック負荷に基づいてドロップ確率 P_d を計算するステップと、

c) 計算されたドロップ確率 P_d を輻輳制御のために提供するステップと

を含む方法において、

ステップa)が、

伝送容量Bに対するネットワークノード(2)のトラフィック負荷の時間平均率として容量利用率 x を決定するステップであって、 $0 < x < 1$ および時間平均化は第1の時間スケールであるステップと、

容量利用率 x の時間平均値 m_1 および容量利用率 x の2乗 x^2 の時間平均値 m_2 を決定するステップであって、時間平均化は、第1の時間スケールよりも長い第2の時間スケールであるステップと、

ネットワークノード(2)の伝送容量Bに適合する、アプリケーションストリームの前記推定された最大数として、 $N = m_1 / (m_2 - (m_1)^2)$ を計算するステップと

10

20

を含むことを特徴とする、方法。

【請求項 2】

第 1 の時間スケールが、ネットワークノード (2) のバッファ保持時間に匹敵することを特徴とする、請求項 1 に記載の方法。

【請求項 3】

第 1 の時間スケールが、 $500 \mu s$ から $100 ms$ の範囲にあることを特徴とする、請求項 1 に記載の方法。

【請求項 4】

第 2 の時間スケールが、 $1 s$ よりも大きいことを特徴とする、請求項 1 に記載の方法。

10

【請求項 5】

第 2 の時間スケールが、第 1 の時間スケールよりも少なくとも 100 倍大きいことを特徴とする、請求項 1 に記載の方法。

【請求項 6】

ステップ b) が、

A アプリケーションストリームの平均を有するパケットトラフィックが k アプリケーションストリームからなる確率 $P(k)$ が、 $P(k) = A^k e^{-A} / k!$ としてポアソン分布に従うと仮定するステップと、

N がネットワークノード (2) の伝送容量 B に適合するアプリケーションストリームの最大数である場合、 $k > N$ について確率 $P(k)$ の合計として、オーバーフロー確率 P_o を計算するステップと、

20

ドロップ確率 P_d をオーバーフロー確率 P_o よりも小さいと仮定するステップとを含む

ことを特徴とする、請求項 1 に記載の方法。

【請求項 7】

計算されたドロップ確率 P_d に従って受信されているパケットトラフィックのパケットをドロップするおよび / またはこのパケットにマーキングするステップをさらに含む

ことを特徴とする、請求項 1 に記載の方法。

【請求項 8】

計算されたドロップ確率 P_d に従って輻輳通知を起動するステップをさらに含む

ことを特徴とする、請求項 1 に記載の方法。

30

【請求項 9】

輻輳通知によってトリガされ、パケットトラフィックのレートを下げるステップをさらに含む

ことを特徴とする、請求項 8 に記載の方法。

【請求項 10】

パケット交換網 (100) において多数のアプリケーションストリームから集約されたパケットトラフィックを受信するネットワークノード (2) であって、ネットワークノード (2) の伝送容量 B に適合するアプリケーションストリームの最大数であるパケットトラフィックの粒度を推定し、推定された粒度およびネットワークノード (2) のトラフィック負荷に基づいてドロップ確率 P_d を計算し、計算されたドロップ確率 P_d を輻輳制御のために提供するように構成された制御ユニット (4) を含むネットワークノード (2) において、

40

制御ユニット (4) が、パケットトラフィックの粒度の上記推定のために、さらに、

伝送容量 B に対するネットワークノード (2) のトラフィック負荷の時間平均率として容量利用率 x を決定し、 $0 < x < 1$ であって、時間平均化は第 1 の時間スケールに基づく、

容量利用率 x の時間平均値 m_1 および容量利用率 x の 2 乗 x^2 の時間平均値 m_2 を決定

50

し、時間平均化は第1の時間スケールよりも長い第2の時間スケールに基づく、

ネットワークノード(2)の伝送容量Bに適合する、アプリケーションストリームの前記推定された最大数として、 $N = m_1 / (m_2 - (m_1)^2)$ を計算する
ように構成されることを特徴とする、ネットワークノード(2)。

【請求項11】

パケット交換網(100)において多数のアプリケーションストリームから集約されたパケットトラフィックを受信するネットワークノード(2)のトラフィック負荷を管理するためのコンピュータプログラムであって、ネットワークノード(2)によって実行されるとき、

a) ネットワークノード(2)の伝送容量Bに適合するアプリケーションストリームの最大数であるパケットトラフィックの粒度を推定するステップと、

b) 推定された粒度およびネットワークノード(2)のトラフィック負荷に基づいてドロップ確率 P_d を計算するステップと、

c) 計算されたドロップ確率 P_d を輻輳制御のために提供するステップと
を実行するコンピュータプログラムにおいて、

ステップa)が、

伝送容量Bに対するネットワークノード(2)のトラフィック負荷の時間平均率として容量利用率 x を決定し、 $0 < x < 1$ および時間平均化が第1の時間スケールに基づくステップと、

容量利用率 x の時間平均値 m_1 および容量利用率 x の2乗 x^2 の時間平均値 m_2 を決定し、時間平均化が、第1の時間スケールよりも長い第2の時間スケールに基づくステップと、

ネットワークノード(2)の伝送容量Bに適合する、アプリケーションストリームの前記推定された最大数として、 $N = m_1 / (m_2 - (m_1)^2)$ を計算するステップとを含むことを特徴とする、コンピュータプログラム。

【請求項12】

第1の時間スケールが、1msから10msの範囲にあることを特徴とする、請求項1に記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、パケット交換網において多数の短期間のアプリケーションストリームから集約されたパケットトラフィックを受信するネットワークノードのトラフィック負荷を管理する方法、ならびに上記方法を実行するためのネットワークノードおよびコンピュータプログラム製品に関する。

【背景技術】

【0002】

RFC2309は、パケットトラフィックを処理する現在のルータで一般的に使用される輻輳通知アルゴリズムであるランダム初期検知(Random Early Detection = RED)アルゴリズムについて記載している。詳細には、REDアルゴリズムは、ルータまたはスイッチのようなネットワークノードにおいてトラフィック負荷の管理に使用される。

【先行技術文献】

【非特許文献】

【0003】

【非特許文献1】「Bandwidth Dimensioning in Packet-based Aggregation Networks」、Lautenschlager, W. および Frohberg, W., Alcatel-Lucent Bell Labs, Murray Hill, NJ, USA, The 13th International Telecommunications Network Stra

10

20

30

40

50

tegy and Planning Symposium 2008 (略してNetworks 2008)、Budapest 2008年9月28日 - 2008年10月2日、1 - 18頁、ISBN: 978 - 963 - 8111 - 68 - 5、<http://ieeexplore.ieee.org/>

【非特許文献2】ITU/ITC Teletraffic Engineering Handbook、ITU-D Study Group 2、Question 16/2、Geneva、2004年3月、<http://www.com.dtu.dk/teletraffic/>

【非特許文献3】Floyd、S. および Jacobson、V.、Random Early Detection Gateways for Congestion Avoidance、IEEE/ACM Transactions on Networking, V. 1 N. 4, 1993年8月、397 - 413頁

10

【発明の概要】

【発明が解決しようとする課題】

【0004】

ネットワークノードのトラフィック負荷の改善された管理を提供することが、本発明の目的である。

【課題を解決するための手段】

【0005】

本発明の第1の目的は、a) ネットワークノードの伝送容量Bに適合するアプリケーションストリームの最大数としてパケットトラフィックの粒度を推定するステップと、b) 推定された粒度およびネットワークノードのトラフィック負荷に基づいてドロップ確率 P_d を計算するステップと、c) 輻輳制御のために計算されたドロップ確率 P_d を提供するステップとを含む、パケット交換網において多数のアプリケーションストリームから集約されたパケットトラフィックを受信するネットワークノードのトラフィック負荷を管理する方法によって達成される。本発明の第2の目的は、パケット交換網において多数のアプリケーションストリームから集約されたパケットトラフィックを受信するネットワークノードによって達成され、このネットワークノードは、ネットワークノードの伝送容量Bに適合するアプリケーションストリームの最大数としてパケットトラフィックの粒度を推定し、推定された粒度およびネットワークノードのトラフィック負荷に基づいてドロップ確率 P_d を計算し、計算されたドロップ確率 P_d を輻輳制御のために提供するように構成された制御ユニットを含む。本発明の第3の目的は、パケット交換網において多数のアプリケーションストリームから集約されたパケットトラフィックを受信するネットワークノードのトラフィック負荷を管理するためのコンピュータプログラム製品によって達成され、このコンピュータプログラム製品は、ネットワークノードによって実行されるとき、ネットワークノードの伝送容量Bに適合するアプリケーションストリームの最大数としてパケットトラフィックの粒度を推定するステップと、推定された粒度およびネットワークノードのトラフィック負荷に基づいてドロップ確率 P_d を計算するステップと、計算されたドロップ確率 P_d を輻輳制御のために提供するステップとを行う。

20

30

【0006】

アプリケーションストリームは、同時発生アプリケーションストリーム、すなわちその存続期間が互いと重なり合うアプリケーションストリームである。ネットワークノードは、エンドポイント間、例えば送信元と宛先との間のパケットトラフィックの伝送における中間ノードである。ネットワークノードは、送信元と宛先との間の伝送機能を表す。ネットワークノードは、パケットを処理するために、例えばルーティングするために、パケットトラフィックを受信する。パケットトラフィックの伝送手段として、ネットワークノードは、ある一定の伝送容量Bを所有する。ネットワークノードがルータまたはスイッチである場合、ネットワークノードの伝送容量は、パケットがルータまたはスイッチを介してそれぞれルーティングまたは交換されるレートによって限定される。ネットワークノードが伝送リンク、より正確にはそのエントリポイントである場合、伝送容量は、単にリンク

40

50

の伝送速度である。パケットトラフィックの量がゼロである場合、ネットワークノードによって表されるリンクの負荷、すなわちリンク負荷もまたゼロ、すなわちその最小値である。パケットトラフィックの量が伝送容量に等しい場合、ネットワークノードのリンク負荷は1、すなわちその最大値である。パケットトラフィックの量が伝送容量より大きい場合、リンク負荷は1より大きい、すなわちネットワークノードは過負荷状態である。

【0007】

パケットトラフィックは、多数の短期間のアプリケーションストリームから多重化される。本発明は統計的方法に基づいているので、本方法は、統計的に有意な数の同時発生アプリケーションストリームにより良く機能する。考慮されるタイムスケールにわたって多数のアプリケーションストリームが含まれているが、ある一定の時点では、1つまたは少数のアプリケーションストリームのパケットのみがネットワークノードに到着している場合があることに注意されたい。「短期間」という用語は、アプリケーションストリームの継続時間が有限の長さであり、好ましくは少なくとも、輻輳制御に含まれる、例えばネットワークノードなどのシステムの典型的なサービス時間よりもかなり短いことを意味する。特定のアプリケーションストリームは、エンドユーザアプリケーションの通信イベントを表す。これは、ビットレート、サイズ、および/または継続時間によって特徴付けられる。アプリケーションストリームは、例えばパケットベースの集約ネットワークを介して上記ノードに接続された多数のエンドユーザによってランダムに、また互いとは無関係に、開始される。同時発生アプリケーションストリームの平均数および平均ビットレートは、パケットトラフィックの粒度と呼ばれ、すなわち時間の所与の瞬間において開始されたがまだ終了されていないアプリケーションストリームの数である。

【0008】

一時的に高くなるパケットの負荷を吸収することができるように、ネットワークノードは、バッファを含むこと、すなわち新しく到着した受信パケットが、そのパケットの処理される順番が来るまでバッファされることが好ましい。時間の所与の瞬間にバッファが満杯である場合、受信パケットはバッファされることが可能ではなく、ドロップされなければならない、すなわちパケットは失われる。

【0009】

イーサネット（登録商標）/IPのような無接続方式のパケット伝送網では、これらのネットワークにおける広く用いられている輻輳緩和は、弾性的トラフィックの概念である（IP=Internet Protocol、インターネットプロトコル）。トラフィック源は、現在の送信レートを下げる要求とともに輻輳状態について知らされる。この本発明の説明では、「現在の」という用語は、現在の時間、好ましくは現時点に開始される、現時点を含んだ、所定の非ゼロの時間間隔の間続く時間、または現時点に終わる、現時点を含んだ、所定の非ゼロの時間間隔の間続いた時間を意味する。理論的に言えばこれは、すべてのソースが公平に共有されて、100%のリソース利用に近い平衡、およびごくわずかな低損失をもたらす。実際に、最も一般的な弾性的トラフィックの実行は、TCPプロトコルである（TCP=Transmission Control Protocol、伝送制御プロトコル）。TCPは、失われたパケットを再送するために、接続エンドポイントへのパケットの到着の記録をとる。同時に、記録されたパケットの損失は、暗黙の輻輳通知として解釈される。適正なTCP実行は、それに応じてその送信レートを下げる。本発明の諸実施形態は、オーバフローがまれな例外となるように弾性的トラフィック（例えばTCP）を管理するために使用されることが可能である。本発明の諸実施形態は、既存のTCP/IP網に導入されることが可能であるルーティングノードまたはスイッチングノードの構成可能な機能を提供する。利点は、TCP接続の流れが非常に滑らかになり、待ち行列のジッタが大いに縮小され、バッファがより小さくなることである。本発明の諸実施形態は、例えばルータまたはスイッチなど、TCPプロトコルに従うパケット伝送機器における輻輳処理を提供する。

【0010】

中間ルータまたはスイッチにおけるパケットのドロップによる上述の暗黙的輻輳通知は

10

20

30

40

50

、特定のドロップ方式によって実行されることが可能である。直接的実行は、パケットが到着するがバッファが満杯であるときに必ずパケットのドロップが発生する単純なFIFOとしてバッファを調べることである(FIFO=First In - First Out、先入れ先出し)。どのパケットがドロップされるかによって、末尾ドロップ(Tail Drop)、先頭ドロップ(Head Drop)、またはランダムドロップ(Random Drop)の方式は区別されることが可能である。残念ながら、こうした前述の単純なオーバフロードロップ方式は、いくつかの深刻な欠点を有する：いわゆるグローバル同期の危険がある、すなわち、影響を受けるすべての接続が、同期してその送信レートを下げ、再確立し、過負荷期間と利用期間とを交互に繰り返す結果となる。第2に、不公平なリソース割り当ての恐れがある。この種の誤った振る舞いの基となる根本的原因は、一般に受け入れられているTCP理論ではランダムに分散されるパケット損失の仮定とは対照的であるバッファのオーバフローの場合のパケットのドロップをバースト様にクラスタ化することである可能性が高い。上述の問題に対して十分に確立され、実行されている緩和策が、十分に確立されたランダム初期検知アルゴリズムである。現在のバッファのオーバフローの場合の単純なランダムドロップではなく、REDは平均待ち行列のサイズに依拠する。パケットスイッチの平均待ち行列サイズは、差し迫るオーバフローの初期表示として使用される。平均待ち行列サイズが一定の閾値を超えてバッファの満杯状態に向かう場合、ランダムパケットドロップが開始されて、理想的には耐え難いバッファのオーバフローが発生する前の早い時期に、TCPエンドポイントにやがて起こるオーバフローを通知する。このプロセスは、パケット損失の危険なバースト化を回避するためである。本発明の諸実施形態は、REDの拡張を行う。

【0011】

REDアルゴリズムは、制御メトリックとして平均待ち行列サイズを計算する。最近の研究は、この測定が、意図された定常状態の待ち行列サイズではなく経時的なオーバフロー状態の割合を指し示すことを明らかにしている。バッファの充填は、一般的に、「ほぼ空」と「ほとんど満杯」の間で変動しており、「ほぼ空」に重点を有するが、低い確率で間のどこか(単に2つの端点の間の過渡事象)となる。この観点から、REDで使用される「平均待ち行列サイズ」は、経時的なオーバフロー状態の割合を示す幾分人工的な尺度である。言い換えれば、REDにおける定常状態の待ち行列サイズの仮定は有効ではない。REDは実際に機能するが、そのあらゆる好ましくない結果(クラスタ化されたパケットの損失、大規模な待ち行列のジッタなど)を伴ってバッファのオーバフローを本当には回避しない。さらにREDは、特定の転送プロセスには必要とされず、負荷測定デバイスとして誤用されるだけのバッファのディメンショニングにさらなる負担を与える。REDは基本的にバッファのオーバフローによってトリガされるが、本発明は、バッファのオーバフローに依拠しない、REDに代わるものを意味する。

【0012】

本発明の諸実施形態は、滑らかに流れるTCP接続、より優れたリソースの利用、より少ないジッタを提供する。本発明の諸実施形態は、面倒な分類/優先順位付けなしに、バッファ空間要件を縮小し、サービスの共存を可能にする。RFC2309と比べると、本発明の諸実施形態は、輻輳通知が時間とともにより良く広められる。本発明の諸実施形態は、TCPおよび上位(アプリケーション)層におけるクラスタ化された損失の重大な影響を回避する。

【0013】

REDの前述の縮小を避けるための1つの直接的考えは、輻輳通知を単に現在のトラフィック負荷にリンクさせることである。特定のネットワーク装置において平均負荷が容量の限界に近づいている場合、ランダムドロップにより、エンドポイントにおいてTCP伝送機による負荷軽減を起動することができる。この直接的手法が機能しない重大なポイントは、所与の容量の耐えられる負荷は、トラフィックの揮発性(volatility)によって決まり、これは、時間について、またあらゆる種類のネットワークについて決して均一ではないことである。先の手法に比べると、本発明は、トラフィックの揮発性に配

10

20

30

40

50

慮する。

【0014】

本発明は、面倒な分類および優先順位付けの方式なしに、多様なサービスのサービス品質に達することができる新しいパケット伝送網のパラダイムに設定される。本発明が設定される大域的考えは、極めて例外的な場合にのみ過負荷が発生するようにトラフィックを統計的に管理することである。

【0015】

提案する発明は、現在のトラフィック負荷およびトラフィック揮発性によりドロップ確率を決定する。REDとは違い、本発明の諸実施形態により決定されるドロップ確率は、バッファ負荷状態に依拠しない。本発明の諸実施形態によって決定されるドロップ確率は、バッファ空間とは無関係の大きさになる。このように、損失を滑らかに分散し(TCPフレンドリ)、待ち行列のジッタを低くして、バッファのオーバフローが十分に回避されることが可能である。副次的効果として、バッファ空間は小さく維持されることが可能である。

10

【0016】

理論的には、パケットトラフィックの粒度もまた、例えばパケットの送信元アドレスおよび宛先アドレスなど、パケットプロトコルを調べることによって決定されることが可能である。しかしながらこれは、すべての中間ネットワークノードにおいて、どれが多くのリソースをバインドし、大量のデータ比較を必要とする時間のかかる手順であるかを見付け出すステートフル接続を必要とする。

20

【0017】

従属請求項によって示される本発明の諸実施形態によって、さらなる利点が得られる。

【0018】

ステップb)は、次のステップ：平均Aアプリケーションストリームを有するパケットトラフィックが、kアプリケーションストリームからなる確率P(k)は、 $P(k) = A^k e^{-A} / k!$ としてポアソン分布に従うと仮定するステップ、そのトラフィック量がネットワークノードの容量Bよりも小さい同時発生アプリケーションストリームの推定最大数がNである場合、 $k > N$ について確率P(k)の合計としてオーバフロー確率 P_o_v を計算するステップ、ドロップ確率 P_d はオーバフロー確率 P_o_v よりも小さいと仮定するステップ、を含むことが可能である。Nは、そのトラフィック量がネットワークノードの容量Bよりも小さいアプリケーションストリームの推定最大数である。

30

【0019】

<http://ieeexplore.ieee.org/>から検索できる、文献「Bandwidth Dimensioning in Packet-based Aggregation Networks」、Lautenschlager, W.およびFrohberg, W.、Alcatel-Lucent Bell Labs、Murphy Hill, NJ, USA、The 13th International Telecommunications Network Strategy and Planning Symposium 2008 (略してNetworks 2008)、Budapest 2008年9月28日 - 2008年10月2日、1 - 18頁、ISBN: 978 - 963 - 8111 - 68 - 5によれば、トラフィックの揮発性は、基本的に、現在のトラフィック負荷を構成する同時発生アプリケーションストリーム数によって決まる。所与の負荷が多数の狭いアプリケーションストリームによって形成される場合には、その変動は低くなる。これは、(発生の可能性が低い)多数の追加ストリームを要求して、特定の伝送容量を過負荷にする。反対の場合、すなわち、少数の膨大なアプリケーションストリームがある場合、予想される変動は高くなる。ただ1つの追加アプリケーションストリームでさえリンクを過負荷にする可能性があり、これはいつでも起こる可能性がある。この影響は、次のように数学的に説明されることが可能である：同時発生アプリケーションストリームがランダムに発生するシステムでは、現在のストリーム数の確率分布は、ポアソン分布に従う：

40

50

【数 1】

$$P(k) = \frac{A^k e^{-A}}{k!} \quad \text{式 (1)}$$

k は現在のストリーム数、P (k) は、A 同時発生アプリケーションストリームの平均転送トラフィックを有するリンク上のその数 k がわかる確率とする。「現在の」という用語は、アプリケーションストリームの平均継続時間よりも実質的に短い現在の期間を指す。厳密に言えば任意の時点では、伝送リンク、例えばネットワークノードは、100%パケットで占有されるか、またはアイドル状態 / 占有されていない状態となるので、概して本発明は現在の期間を指し、現在の時点を指さない。好ましくは短い期間を考えるとのみ、パケットトラフィックの粒度が明らかになる。

10

【0020】

「リンク」という用語は、例えば、データ線またはデータリンクなど、受信パケットトラフィックを処理するネットワークノードと関連するパケット伝送機能を指す。リンクは、例えば、ビット / 秒 (= bps) で測定される、限られた伝送容量 B を有する。ネットワーク要素の特定の容量 B が最大 N の同時発生アプリケーションストリームを処理できる場合、オーバーフロー確率 P_{ov} は、 $k > N$ では式 (1) の確率の合計である：

【数 2】

$$P_{ov} = \sum_{k=N+1}^{\infty} P(k) \quad \text{式 (2)} \quad 20$$

オーバーフローの場合には、すべてのパケットが失われるとは限らず、わずかなオーバーシュート剰余があり、実際の損失確率 P_d (= ドロップ確率) は、オーバーフロー確率 P_{ov} よりもわずかに小さい。より詳細な導出は、Lautenschlager および Frohberg の上述の文献に掲載されている。 P_d の対応する式は、以下の式 (14) に示す。

【0021】

アプリケーションストリームは、パケット伝送網で宣言されず、均一サイズでもない。そこでトラフィック負荷もリンク容量も、前述の考えによれば負荷分布を拡散するための決定的パラメータである粒度がわからない。上記の式 (1) および (2) による損失の確率の予測は、推定を利用しなければならない。

30

【0022】

本発明の一実施形態によれば、ステップ a) は、次のステップ：伝送容量 B に対するネットワークノードのトラフィック負荷の時間平均率として容量利用率 x を決定するステップであって、 $0 < x < 1$ および時間平均が第 1 の時間スケールであるステップ、容量利用率 x の時間平均値 m_1 および容量利用率 x の 2 乗 x^2 の時間平均値 m_2 を決定するステップであって、時間平均が第 1 の時間スケールよりも長い第 2 の時間スケールであるステップ、およびネットワークノードの伝送容量 B に適合するアプリケーションストリームの上述の推定最大数として、 $N = m_1 / (m_2 - (m_1)^2)$ を計算するステップ、を含む。

40

【0023】

所与の集約パケットフローの粒度は、次の手段によって推定される。第 1 の平均化ステップにおいて、ネットワークノードに到着するパケットトラフィックは、ネットワークノードの、詳細には、パケットトラフィックが伝送されるリンク上のネットワークノードのデータリンクの利用可能容量 B に関連して設定される。本発明の一実施形態によれば、第 1 の平均化ステップは、関連ネットワークノードのバッファ保持時間に匹敵する第 1 の時間スケールで行われる。この第 1 の時間スケールは、およそ、パケットの伝送に要する時間、または 2 つの連続したパケット間の時間的距離とすることができる。本発明の一実施形態によれば、第 1 の時間スケールは、500 μ s から 100 ms の範囲、好ましくは 1 から 10 ms の範囲にある。結果として生じる「容量利用率」 x は、0 から 1 の値である。

50

【 0 0 2 4 】

容量利用率 x は、相対トラフィック負荷とも呼ばれる。トラフィック負荷 r は、時間単位あたりのデータユニットの率として、例えばビット/秒の単位で求められるデータレートとして測定される。相対トラフィック負荷 x は、伝送容量 B に絶対トラフィック負荷 r を関連させる 0 から 1 の範囲の無次元量である。

【 0 0 2 5 】

第 2 の平均化ステップでは、第 2 の時間スケールで時間について平均化することによって容量利用率 x から第 1 のモーメント m_1 が導出される。まず容量利用率 x を 2 乗して、次にこの 2 乗した値を第 2 の時間スケールで時間について平均化することによって、容量利用率 x から第 2 のモーメント m_2 が導出される。「モーメント」という用語は、数理統計学において明確に規定された量である。 m_1 および m_2 を作り出すために時間について平均化することは、同一パラメータを用いて行われる。この第 2 の時間平均化ステップは、極めてアプリケーションに依拠するが、上記の第 1 の時間平均化ステップと対照的に、第 2 の時間スケールは数分またはより大きい範囲にある。第 2 の時間スケールは、含まれるネットワーク数に依拠することが可能である。本発明の一実施形態によれば、第 2 の時間スケールは、1 秒より大きい範囲にある。本発明の一実施形態によれば、第 2 の時間スケールは、第 1 の時間スケールより少なくとも 100 倍大きい。

10

【 0 0 2 6 】

上記の第 1 および第 2 のモーメントの推定に、例えば、1 カ月内の毎日の負荷曲線の対応する時間間隔について平均化するなど、他の平均化方法が適用されることも可能である。

20

【 0 0 2 7 】

モーメント m_1 および m_2 から、同時発生アプリケーションストリームの推定最大数 N が計算され、すなわち、容量 B (例えばビット/秒で求められる) は、ストリームの整数に変換される：

【 数 3 】

$$N = \frac{m_1}{m_2 - m_1^2} \quad \text{式 (3)}$$

【 0 0 2 8 】

式 (3) の導出について、次に説明する。総データレート r を有する現在のトラフィックフローは、それぞれ (未知の) データレート b_r の多数のアプリケーションストリームのオーバレイであると仮定される。さらに、アプリケーションストリームが大規模 (無限大に近い) ユーザ群からランダムに、互いと無関係に到着すると仮定される。ITU / ITC Teletraffic Engineering Handbook、ITU-D Study Group 2、Question 16 / 2、Geneva、2004 年 3 月、<http://www.com.dtu.dk/teletraffic/> から、この場合、現在の同時発生ストリーム数は、ポアソン分布に従う乱数 k であることがわかる：

30

【 数 4 】

$$P(k) = \frac{\lambda^k e^{-\lambda}}{k!} \quad \text{式 (4)}$$

40

ここで、強度 λ は、同時発生ストリームの平均数に等しい (「必要トラフィック」としても知られ、「平均保持時間あたりの発呼数」としても知られる)。強度 λ は、次のように計算されることが可能である：

【 数 5 】

$$\lambda = \frac{E[r]}{b_r} = E\left[\frac{r}{b_r}\right] \quad \text{式 (5)}$$

50

r は現在のトラフィックレート、 $E[r]$ は r の期待値とし、 b_r は単一アプリケーションストリームの未知のデータレートとする。同時に、利用可能な容量 B は、データレート b_r の最大 N アプリケーションストリームで搬送することができる：

【数 6】

$$N = \frac{B}{b_r} \quad \text{式(6)}$$

【0029】

式(5)および式(6)から、以下のように導出されることが可能である：

【数 7】

$$\lambda = \frac{E[r]B}{B b_r} = E\left[\frac{r}{B}\right] \cdot N \quad \text{式(7)}$$

ここで、 r/B は、(負荷対容量の)容量利用率 x である。

【0030】

ポアソン分布の標準偏差は、以下のとおりとわかっている：

$$\sigma^2 = \lambda \quad \text{式(8)}$$

【0031】

一方、観測される同時発生アプリケーションストリームの数は、 $k^* = r/b_r$ である。したがって、観測から導出される標準偏差は以下のものである：

【数 8】

$$\sigma^2 = E\left[\left(\frac{x}{b_r} - E\left[\frac{r}{b_r}\right]\right)^2\right] = \frac{E[r^2] - E^2[r]}{b_r^2} = \frac{E[r^2] - E^2[r]}{B^2} N^2 = \left(E\left[\left(\frac{r}{B}\right)^2\right] - E^2\left[\frac{r}{B}\right]\right) \cdot N^2 \quad \text{式(9)}$$

【0032】

式(7)、(8)、および(9)から、以下のように導出されることが可能である：

【数 9】

$$N = \frac{E\left[\frac{r}{B}\right]}{E\left[\left(\frac{r}{B}\right)^2\right] - E^2\left[\frac{r}{B}\right]} \quad \text{式(10)}$$

【0033】

有限アプリケーションストリームを仮定すると、期待値は、時間についての平均値に置き換えられることが可能である：

【数 10】

$$E\left[\frac{x}{B}\right] \cong m_1, \quad \text{および} \quad E\left[\left(\frac{x}{B}\right)^2\right] \cong m_2 \quad \text{式(11), 式(12)}$$

$$N = \frac{m_1}{m_2 - m_1^2} \quad \text{式(13)}$$

【0034】

帯域幅容量 B 内の同時発生アプリケーションストリームの推定最大許容数 N はわかっているので、期待パケットドロップ確率 P_d は、次の用に計算されることが可能である：

10

20

30

40

50

【数 1 1】

$$P_d = \sum_{k=[N]}^{\infty} \frac{A^k e^{-A}}{k!} \left(1 - \frac{N}{k+1} \right) \quad \text{式(14)}$$

転送トラフィック $A = N \cdot m_1$ とする。 式 (1 5)

【 0 0 3 5 】

所与の負荷レベル m_1 におけるドロップ確率 P_d は、負荷自体だけでなく、数 N で表される、トラフィックの粒度、すなわち同時発生ストリームの最大許容数にも依拠することは、明らかである。

10

【 0 0 3 6 】

検討されるトラフィックを供給される容量 B のネットワークノードは、ほぼ推定確率 P_d でパケットをドロップすると予想されることが可能である。残念ながら、実際のドロップは、時間についてよく分散されていないが、詳細には現在の負荷 x が許容限度を超えるときに、短いバーストにクラスタ化される。実際のドロップ率をクラスタ化することは、TCPで広く利用されているにもかかわらず、輻輳通知には不適切とする。

【 0 0 3 7 】

本発明の別の実施形態によれば、実際のドロップの代わりに、推定ドロップ確率 P_d (式 (1 4) 参照) が輻輳通知に使用されることが可能である。これは、上記の第 2 の平均化演算の時間スケールである程度一定である。

20

【 0 0 3 8 】

上記解決法はさらに、TCPエンドポイントによるトラフィック適応への反応が遅れるという欠点を有する。これを克服するために、本発明の別の実施形態は、次のように発見的 (ヒューリスティック) 手法を導入する: 平均負荷 m_1 の代わりに、現在の負荷 x_c 、または、 m_1 と x_c の両方の組合せが式 (1 5) に使用され、すなわち、

【数 1 2】

$$A = N \cdot x_c \quad \text{または} \quad A = N \cdot \sqrt{m_1 x_c}, \quad \text{式(16), 式(17)}$$

m_1 の代わりに現在の負荷 x_c を使用することは、トラフィック変化の力学 (dynamic) を輻輳通知信号に再び導入するが、依然としてクラスタ化を回避し、また輻輳時のトラフィック粒度の影響力を重視する。 m_1 と x_c の組合せは、長期平均 m_1 から、現在の負荷 x_c の任意の例外的な大きい偏差の影響を飽和させるのに有益である。

30

【 0 0 3 9 】

式 (1 4) - (1 7) から導出されるヒューリスティック関数 $P_d = f (N , m_1 , x)$ は、LautenschlagerおよびFrohbergの損失確率の計算に基づく (上記参照)。この関数 $f (N , m , x)$ は、比較的平らな面を構成し、したがって、補間テーブルによって実行されることが可能である。関数 $f (N , m , x)$ は、現在の負荷 x_c だけでなく、過去の平均負荷 $m = m_1$ も考慮に入れて、例外的な不測の偏差の場合に著しくドロップすることを回避する。さらに、ヒューリスティックは、もともとのREDアルゴリズムに含まれるように、閾値メカニズムおよびスケーリング係数を含む。

40

【 0 0 4 0 】

ヒューリスティックは、推定損失がランダムドロップ確率 P_d によって予想されるという仮定に基づく。バーストにクラスタ化されるオーバーフロー損失以外は、ランダムドロップは、Floyd、S.およびJacobson、V.、Random Early Detection Gateways for Congestion Avoidance、IEEE/ACM Transactions on Networking, V. 1 N. 4, 1993年8月、397 - 413頁に説明されているように、滑らかに分散されることが可能である。したがって、これによりTCPおよびアプリケーションに適したパケット損失プロファイルになる。

50

【 0 0 4 1 】

計算されたドロップ確率 P_d は、輻輳制御のために提供される。提供されたドロップ確率 P_d によって、輻輳通知が起動されることが可能である。本発明の別の実施形態によれば、この方法はさらに、計算されたドロップ確率 P_d に従って受信されているパケットトラフィックのパケットをドロップするステップおよび/またはマーキングするステップを含む。ネットワークノードでパケットをドロップするステップは、ネットワークノードのトラフィック負荷を軽減する効果を有することができる。パケットをドロップするステップおよび/またはマーキングするステップの後に、輻輳通知が起動されることが可能である。ドロップ確率 P_d は、パケットの宛先アドレスにより、ネットワークノードによって単にルーティングされないが、ネットワークノードによって特別な方法で処理されるパケットのパーセンテージを指定する。パケットの特別な処置は、パケットがネットワークノードによってドロップされること、例えばパケットが削除されることを意味することが可能である。本発明の諸実施形態は、REDの改善として使用されることが可能であり、ドロップされるパケットは、暗黙的な輻輳通知の効果を有する。ネットワークノードによるパケットの特別な処置は、パケットが印を付けられる（例えば輻輳マーキングなど、フラグまたはビットを設定する）、カウントされる、特別な宛先アドレスにルーティングされるなどを意味することもまた可能である。パケットをドロップすることおよび/またはパケットにマーキングすることは、知られている輻輳回避処理を起動して、データパケットの伝送をいつ送信するか、または遅らせるかを決定する。アプリケーションストリームのエンドポイントは、輻輳通知によって通知されることが可能であり、これらのエンドポイントは、送信元によって送信されるパケットの量を減少させるよう要求される。

10

20

【 0 0 4 2 】

本発明の別の実施形態は、輻輳プライシングまたはアカウントリングに予想されるドロップ確率 P_d を使用する。この場合、2つの異なるネットワークドメイン間の相互接続ゲートウェイにおいて、送信網は、送信網が受信網に送り込む輻輳の程度を明らかにされる。

【 0 0 4 3 】

本発明の別の実施形態によれば、輻輳通知は、パケットトラフィックのレートを下げるためのトリガとして使用される。計算されたドロップ確率 P_d によって、トラフィックの送信元は、現在の送信レートを下げる要求とともに、輻輳状態について知らされる。理想的には、これは、すべての送信元が公平に共有されて、100%のリソース利用、および無視できる低損失に近い平衡につながる。実際には、最も一般的な弾性的トラフィックの実装は、TCPプロトコルである。TCPは、失ったパケットを再送するために、接続エンドポイントへのパケットの到着を記録する。同時に、記録されたパケットの損失は、暗黙の輻輳通知として解釈される。それに応じて正確なTCPの実行は、その送信レートを下げる。本発明の諸実施形態は、TCPの実行と協働する。

30

【 0 0 4 4 】

本発明のこれらの特徴ならびにさらなる特徴、および利点は、添付の図面と関連して次の例示的实施形態の詳細な説明を読むことによってより良く理解されるであろう。

【 図面の簡単な説明 】

40

【 0 0 4 5 】

【 図 1 a 】 本発明の基となる発見的手法を説明する図である。一般に、相対的トラフィック負荷 x の関数として、また N をパラメータとする関数 $f(N, m)$ によって導出される、ドロップ確率 P_d を示す。

【 図 1 b 】 本発明の基となる発見的手法を説明する図である。長期平均 m_1 または現在（短期）の値 x_c のどちらかによって、相対的トラフィック負荷 x が図 1 a の関数にどのように適用されるかを示す。

【 図 2 】 本発明の一実施例によるネットワークノードのブロック図である。

【 発明を実施するための形態 】

【 0 0 4 6 】

50

図1 aおよび1 bは、ドロップ確率 P_d がネットワークノードにおける相対的トラフィック負荷 x の測定からどのように導出されることが可能であるかの一例である。図1 aは、ドロップ確率 P_d を $0 < x < 1$ で容量利用率 x の関数とする曲線の概形を示す。この概形は、同時発生アプリケーションストリームの5つの異なる推定最大数 N に対して、すなわち、 $N = 1, 3, 10, 30$ および 100 に対して、5つの数値的に決定された P_d の曲線を示す。図1 aに示す関数 $P_d = f(N, m)$ は、式(14)によって選択された N の値に対して数値的に取得された。

【0047】

図1 bは、容量利用率 x を時間 t の関数とする曲線の概形を示す。第1に、概形は、現在の容量利用率 x_c の非常に変動する値を示す。現在の容量利用率 x_c は、ネットワークのノードの容量 B に対するネットワークノードのトラフィック負荷 r の時間平均率であり、時間平均化は、例えばミリ秒など、第1の時間スケールに基づく。第2に、概形は、容量利用率 x の時間平均値である一定の値 m_1 を示し、平均化は、数分の時間スケールに基づく。

10

【0048】

ドロップ確率 P_d の計算にどの容量利用率 x が使用されるかによって、ドロップ確率 P_d の著しく異なる値が得られる。また、結果として生じるドロップ確率 P_d は、同時発生アプリケーションストリームの推定最大数 N に著しく依拠する。

【0049】

$N = 10$ のドロップ確率 P_d は、例示的に図1 aにおいて、相対トラフィック負荷 x の3つの異なる値について決定される：平均値 $x = m_1 = 0.35$ は(1点鎖線矢印に従う)ドロップ確率 $P_d = 2.5 \cdot 10^{-4}$ を示し、相対トラフィック負荷の最小値 $x_{c, min} = 0.18$ は(破線矢印に従う)ドロップ確率 $P_d = 2 \cdot 10^{-6}$ を示し、相対トラフィック負荷の最大値 $x_{c, max} = 0.65$ は(点線矢印に従う)ドロップ確率 $P_d = 1.5 \cdot 10^{-2}$ を示す。図1 aおよび1 bは、計算されるドロップ確率 P_d がパケットトラフィックの推定粒度 N および相対トラフィック負荷 x にどれだけ依拠するかを示す。所与の相対負荷レベル x 、ここでは m_1 または x_c におけるドロップ確率 P_d は、負荷 x そのものにだけでなく、トラフィックの粒度 N 、すなわち同時発生ストリームの最大許容数にも依拠することが明らかになる。

20

【0050】

図2は、本発明の一実施形態によるネットワークノードを示すブロック図である。図2は、例えばインターネットなどのパケット交換網100において受信リンク6と多数の送出リンク230とを有する、例えばルータまたはスイッチなどのネットワークノード2を示す。パケット交換網は、コネクションレスパケット伝送網とも呼ばれる。受信リンク6では、多数の送信元から集約されたパケットトラフィックがルータ2に、すなわち入力インタフェース21を介してデータリンク上に届く。ルータ2は、制御ユニット4、ルーティングユニット23、およびルーティングテーブル25を備える。受信パケットは、入力インタフェース21から接続210を介してルーティングユニット23へ送信される。まず、パケットは、ルーティングユニット23のバッファ231に入れられる。パケットの変わり目である場合、またパケットがドロップされるよう選択されない場合(以下参照)、ルーティングユニット23は、受信パケットからルーティングに関する情報、例えばパケットのパケットヘッダから宛先アドレスなどを抜き出し、ルーティングテーブル25で対応するルーティングデータを調べ、多数の出力リンク230の1つまたは多数においてルーティングデータに従ってパケットを転送する。

30

40

【0051】

制御ユニット4は、第1のモジュール42と、2乗デバイス(squaring device)44と、第1の平均化デバイス46aおよび第2の平均化デバイス46bと、粒度計算機48と、確率計算機49とを含む。制御ユニット4は、1つまたは多数の相互にリンクされたコンピュータ、すなわち、ハードウェアプラットフォーム、ハードウェアプラットフォームに基づいたソフトウェアプラットフォーム、および、ソフトウェアおよ

50

びハードウェアプラットフォームによって形成されるシステムプラットフォームによって実行されるいくつかのアプリケーションプログラムからなる。制御ユニット4の機能は、こうしたアプリケーションプログラムを実行することによって提供される。アプリケーションプログラムまたはこれらのアプリケーションプログラムの選択された部分は、システムプラットフォームで実行されるとき、次に説明する確率計算サービスを提供するコンピュータソフトウェア製品を構成する。さらに、このようなコンピュータソフトウェア製品は、これらのアプリケーションプログラムまたは上記アプリケーションプログラムの選択された部分を格納する記憶媒体によって構成される。

【0052】

制御ユニット4は、ドロップ確率 P_d をルーティングユニット23に提供する。ドロップ確率 P_d に応じて、ルーティングユニット23は、対応するパーセンテージの受信パケットを、ドロップする220、またはマーキングする、または他の何らかの方法で輻輳状態に関してこれに通知する、すなわちそれらのパケットを削除する220。一例として、ドロップ確率が0.05である場合、ルーティングユニット23は、統計的手法で5パーセントの受信パケットをドロップする。ルーティングユニット23は、それ自体でどのパケットをドロップするかを決定することができる。好ましくは、ドロップされるパケットの選択は、乱数発生器によって実行される。

【0053】

ルータ2の第1の分配ノード24では、受信パケットトラフィックの信号が、2つの接続に分配される。パケットトラフィックは、ルーティングユニット23に転送される。また、パケットトラフィックは、第1のモジュール42に送信される41。第1のモジュール42への別の入力、ルータ2の現在の容量Bであり、例えばルーティングユニット23の処理容量である。第1のモジュール42は、到着パケットトラフィックの量（例えばビット/秒で測定する）を利用可能容量B（例えばビット/秒で測定する）と関連して設定し、この比を平均化する平均化デバイスである。平均化は、ルータ2のバッファ231の保持時間に匹敵する時間スケールで行われる。例えばこれは、1から100msの範囲とすることができる。到着パケットトラフィックの量および利用可能容量Bの時間平均率は、容量利用率 x と呼ばれ、これは、ミリ秒の分解能で0から1の範囲の値である。

【0054】

第2の分配ノード43では、容量利用率 x は、3つの接続に分配される。容量利用率 x は、接続43aを介して第1の平均化デバイス46aに送信される。また、容量利用率 x は、接続43bを介して2乗デバイス44に送信される。また、容量利用率 x は、接続43cを介して確率計算機49に送信される。

【0055】

第1の平均化デバイス46aは、時間について容量利用率 x を平均化する。平均化は、秒から分以上の範囲の時間スケールで行われる。例えば、この時間スケールは1秒、10秒、10分のうちの3分とすることができる。容量利用率 x の時間平均値は、 m_1 と呼ばれる。数量 m_1 は、接続47aで粒度計算機48へ、接続47bで確率計算機49へ転送される。

【0056】

2乗デバイス44は、容量利用率 x を2乗し、2乗された容量利用率 x^2 を第2の平均化デバイス46bへ転送する。

【0057】

第2の平均化デバイス46bは、2乗された容量利用率 x^2 を時間について平均化する。平均化は、第1の平均化デバイス46aの平均化と同一の時間スケールで行われる。好ましくは、第1の平均化デバイス46aおよび第2の平均化デバイス46bは、同一の時間平均化装置であり、単に異なる入力平均化される。2乗された容量利用率 x^2 の時間平均化された値は、 m_2 と呼ばれる。数量 m_2 は、接続50で粒度計算機48へ転送される。

【0058】

10

20

30

40

50

粒度計算機 4 8 は、ルータ 2 の容量 B を下回る同時発生アプリケーションストリームの推定最大数として数量 $N = m_1 / (m_2 - (m_1)^2)$ を、受信された数量 m_1 および m_2 から計算する。粒度計算機 4 8 は、計算された数量 N を接続 5 2 で確率計算機 4 9 へ転送する。

【 0 0 5 9 】

確率計算機 4 9 は、受信された数量 m_1 (単に m とも呼ばれる)、受信された数量 N、および受信された容量利用率 x から、確率 $P = f(N, m, x)$ を計算する。確率計算機 4 9 は、計算された確率 P をルーティングユニット 2 3 に送信し、ルーティングユニット 2 3 は、詳細にはどのパケット部分にマーキングするまたはこれをドロップするべきかという輻輳通知に確率 P を使用する。

10

【 0 0 6 0 】

ルーティングユニット 2 3 はパケットをドロップし、これが TCP 接続の受信者への暗黙的輻輳通知となるか、例えばパケットにマーキングすることによって、または明示的メッセージによって、通信エンドポイントの 1 つに明示的輻輳通知を送信する。このように、制御ユニット 4 からルーティングユニット 2 3 に提供される計算されたドロップ確率 P_d は、輻輳通知を制御する、すなわち、輻輳通知は、計算されたドロップ確率 P_d によって決まる。計算されたドロップ確率 P_d がゼロである場合、輻輳通知は起動されない。計算されたドロップ確率 P_d が 1 である場合、輻輳通知は確実に起動される。計算されたドロップ確率 P_d がゼロと 1 の間である場合、輻輳通知は計算されたドロップ確率 P_d の値次第で起動される。

20

【 図 1 a 】

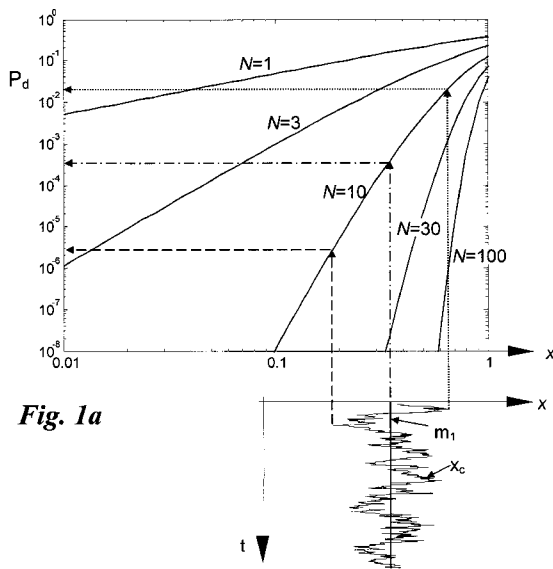


Fig. 1a

Fig. 1b

【 図 1 b 】

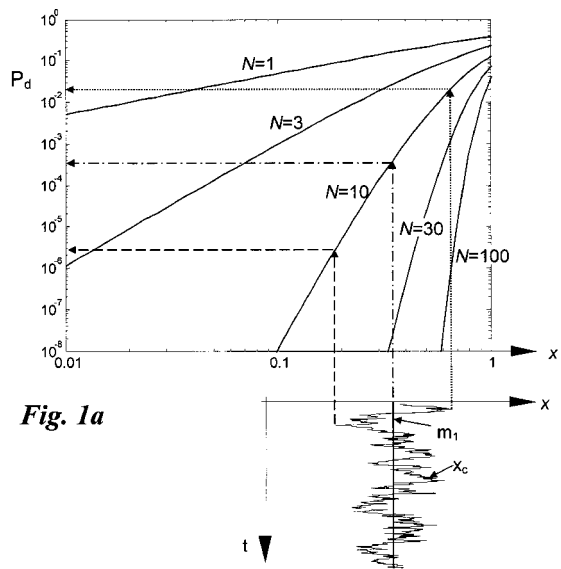


Fig. 1a

Fig. 1b

【 図 2 】

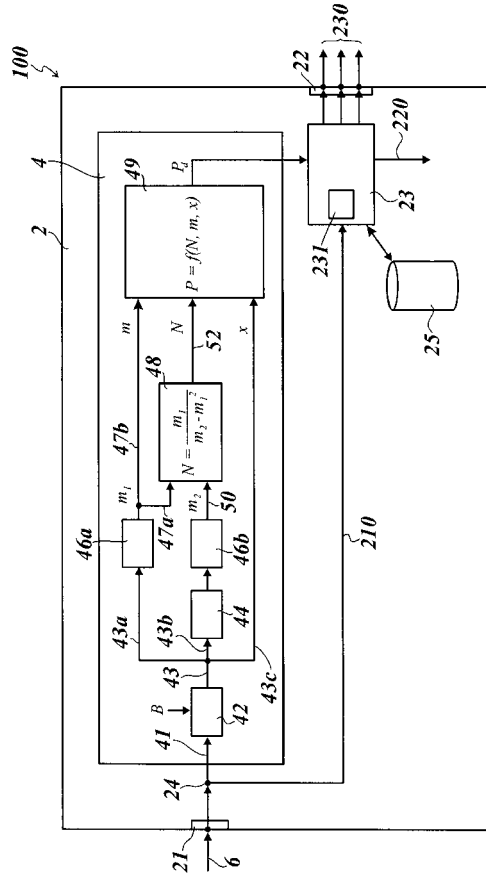


Fig. 2

フロントページの続き

(56)参考文献 石橋 圭介 他, 使用率統計情報を用いたTCPフローレベル性能劣化検出法, 電子情報通信学会技術研究報告, 2003年11月14日, 第103巻, 第443号, p.45~48

(58)調査した分野(Int.Cl., DB名)

H04L 12/00 ~ 12/955

H04M 3/00