

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4828865号
(P4828865)

(45) 発行日 平成23年11月30日(2011.11.30)

(24) 登録日 平成23年9月22日(2011.9.22)

(51) Int.Cl. F I
H O 4 L 12/56 (2006.01) H O 4 L 12/56 1 O O Z

請求項の数 8 (全 36 頁)

(21) 出願番号	特願2005-155178 (P2005-155178)	(73) 特許権者	596092698
(22) 出願日	平成17年5月27日(2005.5.27)		アルカテルルーセント ユーエスエー
(65) 公開番号	特開2005-341589 (P2005-341589A)		インコーポレーテッド
(43) 公開日	平成17年12月8日(2005.12.8)		アメリカ合衆国 07974 ニュージャ
審査請求日	平成20年5月22日(2008.5.22)		ーシー, マレイ ヒル, マウンテン アヴ
(31) 優先権主張番号	60/575350		ェニュー 600-700
(32) 優先日	平成16年5月28日(2004.5.28)	(74) 代理人	100094112
(33) 優先権主張国	米国 (US)		弁理士 岡部 譲
(31) 優先権主張番号	11/106410	(74) 代理人	100064447
(32) 優先日	平成17年4月14日(2005.4.14)		弁理士 岡部 正夫
(33) 優先権主張国	米国 (US)	(74) 代理人	100085176
			弁理士 加藤 伸晃
		(74) 代理人	100106703
			弁理士 産形 和央

最終頁に続く

(54) 【発明の名称】 トラフィック・パターン可変性と独立な効率的で堅牢なルーティング

(57) 【特許請求の範囲】

【請求項1】

リンクによって相互接続され、少なくとも1つの入口点および少なくとも1つの出口点を有するノードのネットワークを通過してデータをルーティングする方法であって、

(a) 前記入口点と前記出口点との間でのデータのルーティングに関するサービス需要を有する、パスに関する要求を受け取る工程と、

(b) 前記入口点と前記出口点との間の1つまたは複数の中間ノードの集合を選択する工程と、

(c) 前記ネットワークの帯域幅に基づいて、前記入口点から1つまたは複数の中間ノードの前記集合の各ノードに送られる前記データのめいめいの分数を決定する工程とを含み、前記入口点および前記出口点に対応するトラフィック行列が、前記トラフィック行列を制約する行合計上限および列合計上限を有し、前記めいめいの分数の前記決定が、前記トラフィック行列の実際のトラフィックのボリュームを考慮に入れずに、前記トラフィック行列を制約する前記行合計上限および前記列合計上限の少なくとも1つを使用して行われ、さらに、

(d) 前記入口点から1つまたは複数の中間ノードの前記集合の各ノードへ、前記決定されためいめいの分数で前記データをルーティングする工程と、

(e) 1つまたは複数の中間ノードの前記集合の各ノードから前記出口点へ前記データをルーティングする工程とを含む、方法。

【請求項2】

工程(c)での前記めいめいの分数の前記決定が、主解および双対解を有する線形プログラムを解くことによって行われ、前記リンクに沿ったフローが、前記主解で増やされ、前記リンクの重みが、前記双対解で乗法的な形で更新される、請求項1に記載の発明。

【請求項3】

前記主解が、

【数1】

$$\sum_{P \in P_f} x(P) = \alpha_j R_i + \alpha_i C_j \quad \forall i, j \in N, i \neq j$$

$$\sum_{P: e \in P} x(P) \leq u_e \quad \forall e \in E$$

10

のもとで

【数2】

$$\sum_{i \in N} \alpha_i$$

を最大化すること、という線形計画法定式化によって表され、前記双対解が、

【数3】

$$\sum_{i: i \neq k} R_i SP(i, k) + \sum_{j: j \neq k} C_j SP(k, j) \geq 1 \quad \forall k \in N$$

$$w(e) \geq 0 \quad \forall e \in E$$

20

のもとで

【数4】

$$\sum_{e \in E} u_e w(e)$$

を最小化すること、という線形計画法定式化によって表され、

30

Nが、ソース・ノードi、宛先ノードj、および中間ノードkを含む、前記ネットワーク内のすべてのノードの集合を表し、

Eが、前記ネットワーク内のすべてのリンクeの集合を表し、

Pが、ノードiからノードjへの所与のパスを表し、

x(P)が、パスP上のトラフィックを表し、

α_i が、ノードiに送られるトラフィックの配分比を表し、

α_j が、ノードjに送られるトラフィックの配分比を表し、

R_i が、任意の時にノードiが前記ネットワークに送るトラフィックの最大総帯域幅を表し、

C_j が、任意の時にノードjが前記ネットワークから受け取るトラフィックの最大総帯域幅を表し、

40

u_e が、リンクeの使用可能な容量を表し、

w(e)が、リンクeの重みの集合を表し、

SP(i, k)が、ノードiからノードkへの重みw(e)の下での最短パスを表し、

SP(k, j)が、ノードkからノードjへの重みw(e)の下での最短パスを表す、

請求項2に記載の発明。

【請求項4】

工程(c)での前記めいめいの分数の前記決定が、不等であるめいめいの分数をもたらす、請求項1に記載の発明。

【請求項5】

50

工程(c)での前記めいめいの分数の前記決定が、前記出口点の識別を考慮に入れずに行われる、請求項1に記載の発明。

【請求項6】

工程(e)での1つまたは複数の中間ノードの前記集合の各ノードから前記出口点への前記データのルーティングが、前記入口点および前記出口点に対応するトラフィック行列に関する他の情報の知識なしで、前記出口点の識別に基づいて実行される、請求項1に記載の発明。

【請求項7】

リンクによって相互接続され、少なくとも1つの入口点および少なくとも1つの出口点を有するノードのネットワークを通してデータをルーティングする装置であって、

(i)前記入口点と前記出口点との間でのデータのルーティングに関するサービス需要を有する、パスに関する要求と、(ii)前記要求に関連する前記データと、を受け取るように適合された入力モジュールと、

前記要求の前記パスを決定するように適合された処理モジュールとを含み、前記処理モジュールは、(a)前記入口点と前記出口点との間の1つまたは複数の中間ノードの集合を選択し、(b)前記ネットワークの帯域幅に基づいて、前記入口点から1つまたは複数の中間ノードの前記集合の各ノードに送られる前記データのめいめいの分数を決定する、ことによって前記パスを決定し、前記入口点および前記出口点に対応するトラフィック行列が、前記トラフィック行列を制約する行合計上限および列合計上限を有し、前記めいめいの分数の前記決定が、前記トラフィック行列の実際のトラフィックのボリュームを考慮に入れず、前記トラフィック行列を制約する前記行合計上限および前記列合計上限の少なくとも1つを使用して行われ、さらに、

前記要求の前記パスに従って、前記入力モジュールからのパケットをルータの出力モジュールに転送するように適合されたルータを含み、前記ルータは、(c)前記入口点から1つまたは複数の中間ノードの前記集合の各ノードへ、前記決定されためいめいの分数で前記データをルーティングし、(d)1つまたは複数の中間ノードの前記集合の各ノードから前記出口点へ前記データをルーティングするように適合される、装置。

【請求項8】

リンクによって相互接続され、少なくとも1つの入口点および少なくとも1つの出口点を有するノードのネットワークであって、

(a)前記入口点と前記出口点との間でのデータのルーティングに関するサービス需要を有する、パスに関する要求を受け取り、

(b)前記入口点と前記出口点との間の1つまたは複数の中間ノードの集合を選択し、
(c)前記ネットワークの帯域幅に基づいて、前記入口点から1つまたは複数の中間ノードの前記集合の各ノードに送られる前記データのめいめいの分数を決定し、前記入口点および前記出口点に対応するトラフィック行列が、前記トラフィック行列を制約する行合計上限および列合計上限を有し、前記めいめいの分数の前記決定が、前記トラフィック行列の実際のトラフィックのボリュームを考慮に入れず、前記トラフィック行列を制約する前記行合計上限および前記列合計上限の少なくとも1つを使用して行われ、さらに、

(d)前記入口点から1つまたは複数の中間ノードの前記集合の各ノードへ、前記決定されためいめいの分数で前記データをルーティングし、

(e)1つまたは複数の中間ノードの前記集合の各ノードから前記出口点へ前記データをルーティングする、ネットワーク。

【発明の詳細な説明】

【技術分野】

【0001】

本願は、参照によって本明細書に組み込まれる、2004年5月28日出願の同時係属の米国仮特許出願第60/575350号の優先権を主張するものである。

本発明は、遠隔通信システムにおけるルーティングに関し、具体的には、保証されたサービス・レベルを有するルーティングのためのネットワークのノードを通るパスの判定に

10

20

30

40

50

関する。

【背景技術】

【0002】

インターネットなどのパケットベース通信ネットワークでは、パケット・フローと称するデータ・パケットの各ストリームが、ネットワークを介して、ソースから宛先へのネットワーク・パス上で転送される。各ネットワーク・パスは、リンクの集合によって相互接続されたノードの集合によって定義される。ノードに、1つまたは複数のルータが含まれる場合があり、ルータは、コンピュータの間でのデータ転送を処理する、ネットワーク内のデバイスである。

【0003】

通信システムは、異なるサイズのネットワークが相互接続されるように構成することができ、その代わりにまたはそれに加えて、同等のサイズのネットワークが相互接続される1つまたは複数のピア構造を含めることができる。パケット・ネットワークは、入口点 (ingress point) および出口点 (egress point) と称するノードを介して別のパケット・ネットワークに接続することができる。用語入口点および出口点は、別のパケット・ネットワークに接続された、パケット・ネットワークのノードを指すことができ、あるいは、その代わりに、他のパケット・ネットワークの接続するノードを指すことができる。複数の他のパケット・ネットワークの間でパケットを転送する大容量のパケット・ネットワークを、一般に「バックボーン」ネットワークと称する。

【0004】

図1に、パケット・ネットワーク102から104の間の通信を可能にする、リンク101を介して相互接続されたノードn1からn9を有する従来技術のバックボーン・ネットワーク100を示す。バックボーン・ネットワーク100の入口点の1つが、ノードn1であり、このノードは、ソース(すなわちパケット・ネットワーク102)からパケットを受け取り、このバックボーン・ネットワークの出口点の1つが、ノードn4であり、このノードは、宛先(すなわちパケット・ネットワーク104)にパケットを送る。バックボーン・ネットワーク100は、内部ルーティング・プロトコルをサポートして、ネットワーク・トポロジ情報を配布し、ノードn1からn9を介するベストエフォート・ルーティング(たとえば、宛先ベースの最短パス・ルーティング)に基づいて、入口点と出口点の間でパケットをルーティングする。集中ネットワーク管理システム105を使用して、(i)バックボーン・ネットワーク100を通る仮想回線(virtual circuit)またはパケット・フローを提供し、(ii)リンク101の容量および使用率を監視し、(iii)プロビジョニングされるパスの計算およびインストールを調整することができる。転送テーブルが、各ノードによって使用されて、受け取られたパケットが、その宛先に向かって次のノードに転送される。さらに、集中ネットワーク管理システム105を使用して、ネットワーク・トポロジ情報を収集し、配布することもできる。

【0005】

内部ルーティング・プロトコルは、ソース/宛先対の間でバックボーン・ネットワークのノードを通るパスに沿ったパケットの転送を決定するのに使用される。ノードによって受け取られたパケットは、内部ルーティング・プロトコルに従って構成された転送テーブルまたは明示的経路プロビジョニングによってインストールされた経路に基づいて、他のノードに転送される。内部ルーティング・プロトコルは、ノードの間でのネットワーク・トポロジおよびリンク状態情報(「ネットワーク・トポロジ情報」)の交換を指定して、ノードが対応する転送テーブルを構成できるようにすることもできる。さらに、いくつかのルーティング・プロトコルは、リンク「コスト」をノード間の各リンクに関連付ける。このリンク・コストは、たとえば、平均リンク使用率またはリンクによって生成される収入ならびにネットワークでのリンクの重要さに関連付けることができる。リンク状態情報またはリンク帯域幅(たとえば、接続性または使用可能な帯域幅)がルータの間で交換される時に、ネットワーク内の各ノードは、そのネットワークのトポロジの完全な記述を有する。「ベストエフォート」ルーティングに広く使用されている内部ルーティング・プロ

10

20

30

40

50

トコルの例が、OSPF (Open Shortest Path First) プロトコルである。

【0006】

ルーティング・プロトコルは、接続性を提供するほかに、トラフィック管理も可能にすることができる。たとえば、MPLS (Multi-Protocol Label Switched) 標準規格は、バックボーン・ネットワークでのそのようなルーティング・プロトコルを可能にする。MPLS 標準規格は、プロビジョニングされるサービス・レベル (保証されたサービス品質 (QoS) とも称する) を有する仮想回線 (パケット・フロー) を有するネットワークに使用することができる。

【0007】

プロビジョニングされるサービス・レベルは、たとえば、バックボーン・ネットワークを通るパケット・フローのパスの保証される最低帯域幅とすることができる。入口点と出口点の間のサービスの保証されるレベルを有するこのパスを、ネットワーク・トンネル・パス (Network Tunnel Path, NTP) と称する場合がある。当業者に明白であるように、NTP の特定の実装が、異なるネットワークのタイプについて存在する。NTP の例として、仮想回線を、TCP/IP ネットワーク内のパケット・フローについて確立することができ、仮想回線を、非同期転送モード (ATM) ネットワーク内のセルについて確立することができ、LSP (Label-switched path) を、MPLS ネットワーク内のパケットについて確立することができる。NTP のルーティングが計算されたならば、RSVP (IP ネットワークおよび MPLS ネットワークの Reservation Protocol) または LDP (MPLS ネットワークの Label Distribution Protocol) などのシグナリング・プロトコルのパケットを使用して、リンク帯域幅を予約し、NTP を確立することができる。NTP は、バックボーン・ネットワークのノードの間の特定のパスに沿った明示的な経路としてプロビジョニングすることができ、すなわち、NTP がパケット・フローについてプロビジョニングされる時に、その NTP の入口点と出口点の間のすべての中間ノードを指定でき、これを、そのフローの各パケットが通過する。

【0008】

MPLS ネットワークでは、パケットが、そのパケットが入口点で受け取られる時に追加情報をパケットに付加するか、パケットから形成することによってカプセル化される。ラベルと称するこの追加情報は、バックボーン・ネットワークのルータによって、パケットを転送するのに使用される。図 2 に、ラベル 201 をパケット 202 に付加された、そのようなカプセル化されたパケット 200 を示す。ラベルは、パケット・ヘッダ内の情報を要約する。この要約は、ヘッダ・フィールドに基づくものとしてことができ、入口点のアドレスを識別する起点 (ソース) アドレス・フィールド (s) 210 と、出口点のアドレスを識別する終点 (宛先) アドレス・フィールド (t) 211 を含む。いくつかの場合に、ラベルを、単に、受け取られたパケットのヘッダ内の起点アドレス・フィールドおよび終点アドレス・フィールドを識別するか他の形でこれらに関連するポイントとすることができる。ラベルに、1 つまたは複数のサービス・レベル・フィールド (bd) 212 も含めることができる。サービス・レベル・フィールド 212 は、必要な最小帯域幅など、仮想回線について望まれるサービス・レベル (「デマンド」と称する) を識別することができる。いくつかのネットワークで、サービス・レベル・フィールドは、ラベル自体から暗示される。MPLS 標準規格バージョン、内部ルーティング・プロトコル・バージョン、最大遅延、または他のタイプのサービス・レベル・パラメータなど、他のフィールド 213 をラベル 201 に含めることができる。その代わりに、ラベル 201 を、パケット 202 のパケット・ヘッダ (PH) 214 に挿入することができ、したがって、図 2 に示されたフィールドの順番は、例示にすぎない。バックボーン・ネットワークは、ラベルを使用して、類似する LSP を有するカプセル化されたパケットをクラス (equivalence class) にグループ化することができ、equivalence class を転送する方法を使用して、LSP のルーティングの計算を単純にすることができる。

10

20

30

40

50

【 0 0 0 9 】

転送テーブルを生成するために、ネットワーク・ノードを通る好ましいパスの集合が、計算され、好ましいパスの集合を計算するのに、重みを使用することができる。各好ましいパスは、ノードの間の最小の総重み（パスの総重みは、パス内のすべてのリンクの重みの合計である）を有し、この重みが、当技術分野で最短パス・ルーティングと称する技法で使用される。好ましいパスの結果の集合は、最短路木（shortest-path tree、SPT）を用いて定義することができる。ルーティング情報（たとえば、ソース/宛先対、ソース・ポート、および宛先ポート）を有する転送テーブルが、SPTから生成される。その後、このルーティング情報が、受け取ったパケットをSPTの最短パスに沿ってその宛先に転送するのに使用される。SPTは、その教示が参照によって本明細書に組み込まれる、E. Dijkstra、「A Note: Two Problems In Connection With Graphs」、Numerical Mathematics、vol. 1、1959年、269～271頁に記載のダイクストラのアルゴリズムなどのアルゴリズムを使用して計算することができる。

10

【 0 0 1 0 】

LSPのルーティングを生成するのにルータによって使用される一般的な最短パス・ルーティング・アルゴリズムが、最小ホップ・アルゴリズムである。最小ホップ・アルゴリズムでは、各ルータが、入口点と出口点の間で、パケットのストリーム（パケット・フロー）のバックボーン・ネットワークを通るパスを計算する。各ルータは、最小の個数（「最小」）の実現可能なリンク（「ホップ」）を有する、入口点から出口点までパケット・フローをルーティングするパスを構成する（実現可能なリンクとは、そのパケット・フローをルーティングするのに十分な容量を有するリンクである）。最短パス・ルーティングなどの従来技術のルーティング方式では、宛先アドレスだけに基づいてパケットを転送し、静的でトラフィック固有の独立リンク重みを使用して、ルーティング・テーブルのパスを計算する。ある入口点/出口点对の間の最短パス上の一部のリンクが、輻輳する場合があるが、代替パス上の他のリンクは、あまり使用されない。

20

【 0 0 1 1 】

RSVPまたはLDPなどのシグナリング機構を使用して、パケット・フローに関するネットワークを介する接続の予約および確立の両方を行うことができる。シグナリング機構は、バックボーン・ネットワークをトラバースするLSPのサービス品質属性を指定することができる。複数のLSPの最短パス・ルーティングによって引き起こされるリンク輻輳は、わずかに長いだけである代替のあまり利用されていないパスにLSPに十分なサービス・レベル（サービス品質保証）が存在する可能性がある場合であっても、シグナリング機構による予約要求の拒絶を引き起こす場合がある。最短パス・ルーティングが使用される時には、使用可能なネットワーク・リソースが、必ず効率的に使用されるわけではない。

30

【 0 0 1 2 】

BGP（Border Gateway Protocol）は、自律システム間ルーティング・プロトコルである。自律システムとは、共通の監督の下で共通のルーティング・ポリシーを有するネットワークまたはネットワークのグループである。自律システム間ルーティング・プロトコルは、自律システムの間でデータをルーティングするのに使用される。BGPは、インターネットのルーティング情報を交換するのに使用され、インターネット・サービス・プロバイダ（ISP）の間で使用されるプロトコルである。大学および会社などの顧客ネットワークは、通常、ネットワーク内でルーティング情報を交換するのに、RIP（Routing Information Protocol）またはOSPF（Open Shortest Path First）などのIGP（Interior Gateway Protocol）を使用する。顧客は、ISPに接続し、ISPは、BGPを使用して、顧客経路およびISP経路を交換する。BGPは、自律システムの間で使用することができ、あるいは、サービス・プロバイダは、BGPを使用して、自律システム内で経路を交換することができる。

40

50

【特許文献1】米国仮特許出願第60/575350号

【非特許文献1】E. Dijkstra, 「A Note: Two Problems In Connection With Graphs」, Numerical Mathematics, vol. 1, 1959年, 269~271頁

【非特許文献2】N. G. Duffield, P. Goyal, A. G. Greenberg, P. P. Mishra, K. K. Ramakrishnan, J. E. van der Merwe, 「A flexible model for resource management in virtual private network」, ACM SIGCOMM 1999, 1999年8月

【非特許文献3】F. Shahrokhi and D. Matula, 「The Maximum Concurrent Flow Problem」, Journal of ACM, 37(2): 318~334, 1990年

【非特許文献4】R. K. Ahuja, T. L. Magnanti, and J. B. Orlin, 「Network Flows: Theory, Algorithms, and Applications」, Prentice Hall, 1993年2月

【発明の開示】

【発明が解決しようとする課題】

【0013】

ネットワークでの主な問題は、BGPによって導入されるトラフィック変動である。極端なネットワーク・トラフィック変動が、さまざまな理由から発生する可能性がある。たとえば、複数の他のプロバイダとトラフィックを交換する大規模なインターネット・サービス・プロバイダの場合に、キャリアの間のトラフィック交換は、通常、長い時間期間にわたる総トラフィック・ボリュームおよび多分ピーク・レート限度（通常は物理リンク容量によって決定される）によって指定される。さまざまなネットワーク出口点への入口点に入るトラフィックの実際の分布は、先験的に既知でない可能性があり、経時的に変化し得る。これは、分布が、異なる宛先プレフィックスへのトラフィックの本質的な変動などの多数の要因と、キャリアによってローカルに行われるかキャリアが制御を有しない他の自律システム内で行われる変更起因するルーティング変更によって決定される。トラフィック分布の本質的な変化は、特別な出来事に応答する急激な混雑の突然の出現など、多数の要因によって引き起こされる可能性がある。トラフィック分布に影響し得るローカル・ルーティング変化の例が、IGP重み変更と「hot-potato」ルーティングの組合せであり、これは、プレフィックスの組に宛てられたトラフィックが、そうでなければ選択するはずのネットワーク出口点を変更する可能性がある。「hot-potato」ルーティングは、ネットワークのノードが、最終的な所定の宛先に移動する前にパケットを保管するバッファを持たず、ルーティングされる各パケットが、その最終宛先に達するまでコンスタントに転送されるようになる、ルーティングの形態である。したがって、パケットは、「hot potato（熱いジャガイモ）」のようにはね回り、ネットワークを移動し続けなければならないので、時々、宛先から遠くに向かって移動する。もう1つの例が、MED（Multi-Exit Discriminator）が使用される時のBGPの変化である。MEDは、経路の「外部メトリック」とも称するが、複数の入口点を有する自律システムへの好ましいパスに関する、外部の近傍への提案である。ローカル・ルーティング変更は、キャリアの制御の下にあり、したがって計画された場合に限ってトラフィック・パターンを変更するが、他の自律システム内のルーティング変更が、下流自律システムに影響する時に、予測不能なトラフィック・シフトが発生し得る。hot-potatoルーティングの広い使用に起因して、ある自律システム内のIGP重み変化（新しいリンクの追加、保守、トラフィック・エンジニアリングなどに起因する可能性がある）が、トラフィック・パターンの大きいシフトを引き起こす可能性がある。IGPコストの変化は、プレフィックスのかなりの比率に関するBGP経路に影響し、影響されるプレフィックスは、トラフィックのかなりの比率を計上する可能性がある。したがって、トラフィックの大きいシフトが、ネットワークの他の場所での変化に起因して、あ

10

20

30

40

50

るキャリアで発生する可能性がある。

【 0 0 1 4 】

高いトラフィック変動性を考慮しなければならないもう1つの理由が、ピアリング合意に達するユーザまたはキャリアが、さまざまなサイトへのトラフィックの特性をうまく表せない場合があることである。送信または受信のいずれかの総集約帯域幅だけを推定することは、はるかに簡単である。したがって、正確なトラフィック行列を知ることには頼る必要を無くし、その代わりに、トラフィック行列の部分的な詳細だけを使用することが好ましい。また、トラフィック行列が既知の時であっても、しばしば、トラフィック分布の変化を検出することは困難である。

【 0 0 1 5 】

ネットワーク輻輳は、通常、容量の消失（ルータまたはリンクの故障時）または容量需要の増加（トラフィックの大幅な増加によって引き起こされる）のいずれかに起因して発生する。これらの制御不能な事象に応じて、キャリアは、ネットワーク輻輳を避けるようにドメイン内ルーティングを適応させるか、ルーティング変更を頼らずに、発生し得る異なるトラフィック・パターンおよび障害パターンに対処するために、先験的に十分な容量をとっておく必要があり、繰り返してそうする。動作の複雑さおよびコストに起因し、かつ、変更が正しく実施されない場合のネットワーク不安定性の危険性に起因して、頻繁なドメイン内ルーティング変更を避けることが好ましい。さらに、上で述べたように、ある自律システム内の変更が、他の自律システム内のトラフィック変化のカスケードを引き起こし、これによって多数のインターネット・パスの総合的な安定性に影響する可能性がある。大規模なルーティング変化を避けることのトレードオフは、障害パターンまたは変化するトラフィック・パターンに対処するために行わなければならないかなりの容量オーバープロビジョニングである。理想的には、プロバイダは、（ i ）構成パラメータのトラフィック依存の動的適応を必要とせず、（ i i ）障害の後の動的容量再割振りを最小にし、（ i i i ）オーバープロビジョニングの必要が最小である、ほぼ固定されたルーティング方式を使用することを好む。

【 0 0 1 6 】

トラフィック行列が先験的に未知であるもう1つの応用例が、ネットワークベースの仮想プライベート・ネットワーク（VPN）サービスの企業顧客への提供である。この場合に、各顧客とのサービス・レベル契約によって、VPNに属する各サイトによって送信または受信できるトラフィックの量が指定される。このシナリオでは、ユーザは、トラフィック行列を知らず、キャリアに、総トラフィック・ボリュームおよびピーク・レートだけを指定する。提供されるVPNトラフィックのすべてをネットワークにトランスポートし、大きすぎる遅れを導入せずにそのトラフィックを搬送することが、キャリアの仕事である。各サイトから他のサイトへの実際のトラフィック分布は、通常、未知であり、時刻によって変化し得る。キャリア・ネットワークは、トラフィック・パターン変化の時またはノード障害もしくはリンク障害の時にネットワーク輻輳を経験せずに、提供されるVPNトラフィックのすべてを搬送する仕事を割り当てられる。

【 0 0 1 7 】

グリッド・コンピューティング用のネットワークは、さらに、トラフィック変動が極端になる可能性があり、トラフィック行列が先験的に既知ではないさらなるシナリオを提供する。グリッド・コンピューティングでは、複雑な計算タスクが、地理的に分散する可能性があり、ネットワークによって接続された異なるコンピューティング・ノードの間で区分される。グリッド・コンピューティング・ノードの間の通信パターンは、非常に予測不能であり、高いバースト・レートも経験する可能性がある。トラフィック行列が先験的に既知でないので、1つのオプションは、基礎になるネットワークで容量を動的に予約することであるが、この手法は、多数のグリッド・コンピューティング応用例には低速すぎる。宛先の高い変化性およびトラフィックのバースト的性質のゆえに、ネットワークのオーバープロビジョニングが、ほとんどの時間の非常に低い容量使用量につながる。

【 0 0 1 8 】

10

20

30

40

50

トラフィック・パターンが制御不能に変化し得る時によりサービスを提供するために、キャリアは、ネットワーク輻輳を避けるためにドメイン内ルーティングを素早く繰り返して適応させるか、ルーティング変更に頼らずに発生し得る異なるトラフィック・パターンに対処するために先験的に十分な容量をとっておかなければならない。サービス・プロバイダは、(i) 動作の複雑さおよびコストと、(i i) リンク・メトリック変化が正しく実施されない場合のネットワーク不安定性の危険性に起因して、頻繁なドメイン内ルーティング変更を避けることを好む。さらに、ある自律システム内の B G P アプリケーション内の変化が、他の自律システムのトラフィック変化のカスケードを引き起こし、これによって多数のインターネット・パスの総合的安定性に影響する可能性がある。ルーティング変更を避けることのトレードオフは、ルーティングを固定したままで変化するトラフィック・パターンに対処するために行うことができる大きい容量オーバ・プロビジョニングである。理想的には、プロバイダは、構成パラメータのトラフィック依存の動的適応を必要とせず、容量の必要において儉約的な、固定ルーティング方式を使用したがる。

10

【 0 0 1 9 】

さらに、IP-over-OTN (Optical Transport Network) では、ルータが、通常はIPルータ・ポートより安価でない光クロスコネクタ (OXC) からなる再構成可能な交換光バックボーンまたはOTNを介して接続される。OXCは、波長分割多重 (WDM) リンクを使用して、メッシュ・トポロジで相互接続される。そのようなOXCからなるコア光バックボーンは、光レイヤでの交換、グルーミング、および回復の機能を引き受ける。IPTraffic・フローは、光レイヤ回路 (「ライトパス」と称する) で搬送されるので、過渡的トラフィックに関するルータ・ポートのバイパスによって、IP-over-OTNの光バックボーンを介してIPルータを相互接続することによって獲得される巨大な規模の経済の基礎が作られる。過渡的トラフィックをルータから光スイッチに移動することによって、増加するトラフィックに伴ってルータのPoP (Point-of-Presence) 構成を更新する要件が、最小になる。というのは、光スイッチが、通常はルータより増やされたポート数を有することに起因して、よりスケラブルであるからである。IP-over-OTNアーキテクチャでは、ルータ・ライン・カードが、通常は、光スイッチ・カードより高価であり、したがって、ネットワーク・コストは、通常、トラフィックのほとんどを光レイヤに保つことによって減らされる。また、光スイッチは、通常、ルータより信頼性が高いので、そのアーキテクチャは、通常はより堅牢であり、信頼性がある。ルータは、交換光バックボーンを介して相互接続されるので、ルーティング処理は、トラフィックを光レイヤに保つことと、データ・トラフィックの効率的で満足な多重化を達成するためにパケット・グルーミングに中間ルータを使用することの間の妥協である。

20

30

【 0 0 2 0 】

高速回復機能を有する帯域幅保証されたパスの動的プロビジョニングは、MPLS (Multi-Protocol Label Switched) ネットワークおよび光メッシュ・ネットワークなどの多数のネットワークの望ましいネットワーク・サービス機能である。光ネットワークでは、光トランスポート・ネットワークが、それぞれが異なる厳密な信頼性要件を有するさまざまなタイプのトラフィックを搬送するので、高速回復も望ましい。パケット化された音声、クリティカルな仮想プライベート・ネットワーク (VPN) トラフィック、または他のサービス品質 (QoS) 保証などのサービスに必要な信頼性を提供するために、類似する高速回復機能をMPLSネットワークで使用することができる。

40

【 0 0 2 1 】

ネットワーク内の接続は、パス・レベルまたはリンク・レベルで保護することができる。リンク回復 (ローカル回復または高速回復とも称する) について、接続の各リンクを、保護されるリンクを排除した1つまたは複数の事前にプロビジョニングされる迂回パスの組によって保護することができる。リンクの障害時に、故障したリンクのトラフィックを、迂回パスに切り替える。したがって、リンク回復は、リンク障害を避けてルーティング

50

するローカル機構を提供する。パス回復では、接続の主なまたは機能するパスが、ソースから宛先への「別個の」バックアップ・パスによって保護される。動作するパスのリソースのいずれかの障害時に、トラフィックが、ソース・ノードによってバックアップ・パスに切り替えられる。リンク回復は、通常は、パス回復よりはるかに高速の回復サービスである。というのは、回復が、ローカルにアクティブ化され、パス回復と異なって、障害情報を、ネットワークを介してソースまで伝搬させる必要がないからである。

【0022】

サービス回復は、光ネットワークの重要な要件である。ノード（光スイッチ）またはリンク（光ファイバ）などのネットワーク要素が故障した場合に、その障害は、1つまたは複数の特定の波長パスを故障させ、影響されるトラフィック・フローを、非常に短い間隔（たとえば50ms）内で代替パスを使用して回復しなければならない。比較的高速の回復時間を達成するために、プロビジョニングは、波長パスごとに、ネットワークを通る2つのパスすなわち、主（アクティブ）パスと副（バックアップ）パスを識別する。バックアップ・パスは、主パスに関して、リンクが互いに疎（link disjoint、アクティブ・パスとバックアップ・パスがリンクを共用しない）またはノードが互いに疎（node disjoint、アクティブ・パスとバックアップ・パスがノードもリンクも共用しない）である。バックアップ・パスのリンクの容量は、対応する主パス（たとえば波長）に排他的に割り当てることができ、あるいは、ネットワーク帯域幅使用の効率のために、容量を、望まれる回復のタイプに応じて、異なる主パスのバックアップ・パスのリンクの間で共用することができる。光ネットワーク容量設計は、通常、可能な共用を有する互いに疎な副パスをルーティングする回復の必要を計上する。

【0023】

現代の非常に動的で変化するトラフィック環境での堅牢なネットワーク・ルーティングは、（i）異なる宛先に送りたがる、トラフィックにおいて予測不能ではあるが「よい」サービスを求めるユーザに対処し、（ii）洗練されたトラフィック予測機構およびトラフィック管理機構に頼らずに「ベストエフォート・ネットワークングをよりよく」するためにネットワーク内で行われる必要がある「オーバープロビジョニング」の量を最小にし、（iii）ほとんど静的ルーティング構成を用い、動的ルーティング調整なしでネットワークを効率的に運営し、これによってネットワークの入口ルータと出口ルータの間のトラフィック・フローの劇的な変化によって引き起こされる輻輳を避けるために、インターネット・サービス・プロバイダが使用しなければならないネットワーク・ルーティング方法の知識がなければ行うことができない。これらの目標の達成は、困難であり、その代わりに、ネットワーク・ルーティングを変化したトラフィック需要に適應させるトラフィック管理方式を実施するという管理の複雑さを回避するために非常にオーバープロビジョニングされたネットワークがもたらされる。

【課題を解決するための手段】

【0024】

本発明は、キャリアのドメイン内の最短パスIGP（Interior Gateway Protocol）ルーティングを、トラフィックが、やはりキャリアのドメイン内にある1つまたは複数の事前に決定された中間ノードを通過することを保証した後にトラフィックを宛先にルーティングする変更されたルーティング方式に置換するという発想に基づく方式を提供する（中間ノードの割当は、パケット再シーケンシングの問題を避けるためにフロー・レベルで行われる）。出口ノードは、それでも、BGP（Border Gateway Protocol）によって決定される自律システム・パスおよびhot-potatoルーティングなどの補助キャリア・ルーティング・ポリシーに基づいて選択される。本発明の1実施形態と一貫する方式は、直接最短パスのIGPパス選択を、1つまたは複数の先験的に割り当てられた中間ノードを通るパスに変更する。MPLSネットワークでは、この1つまたは複数の事前に決定された中間ノードを通るルーティングを、指定された確率に従ってフローが割り当てられる1つまたは複数の中間ノードの選択された組と入口ノードとの間のMPLS LSPの事前に構成された組を使用して達成する

ことができる。純IPネットワークでは、このルーティングを、まず1つまたは複数の事前に決定された中間ノードにパケットをトンネリングすることによって達成することができる。1つまたは複数の中間ノードの事前に決定された選択を用いるこのルーティングは、エッジ-リンク容量制約の対象になる、許容可能なすべてのトラフィック・パターンを処理するのに十分であり、さらに、ルータ障害および光レイヤ・リンク障害に対する保護を提供する。さらに、ルーティング適応は、トラフィック行列が変化する時に必要ではなく、この方式は、帯域幅効率が良い。

【0025】

本発明は、さらに、IP-over-OTNまたは他の回路交換ネットワークに適用される時に、1つの中間ルータだけでのパケット・グルーミングを用いて光レイヤでパケットをルーティングでき、大きく変化するトラフィックに関するパケット交換の望みの統計的多重化特性を提供できるルーティング方式を提供する。

10

【0026】

1実施形態で、本発明は、リンクによって相互接続され、少なくとも1つの入口点および少なくとも1つの出口点を有するノードのネットワークを通過してデータをルーティングする方法を提供し、この方法には、(a)前記入口点と前記出口点との間でのデータのルーティングに関するサービス需要を有する、パスに関する要求を受け取る工程と、(b)前記入口点と前記出口点との間の1つまたは複数の中間ノードの集合を選択する工程と、(c)前記ネットワークの帯域幅に基づいて、前記入口点から1つまたは複数の中間ノードの前記集合の各ノードに送られる前記データのめいめいの分数を決定する工程と、(d)前記入口点から1つまたは複数の中間ノードの前記集合の各ノードへ、前記決定されたためいめいの分数で前記データをルーティングする工程と、(e)1つまたは複数の中間ノードの前記集合の各ノードから前記出口点へ前記データをルーティングする工程とが含まれる。

20

【0027】

もう1つの実施形態で、本発明は、リンクによって相互接続され、少なくとも1つの入口点および少なくとも1つの出口点を有するノードのネットワークを通過してデータをルーティングする装置であって、入力モジュールと、処理モジュールと、ルータとを含む装置を提供する。入力モジュールは、(i)前記入口点と前記出口点との間でのデータのルーティングに関するサービス需要を有する、パスに関する要求と、(ii)前記要求に関連する前記データとを受け取るように適合される。処理モジュールは、前記要求の前記パスを決定するように適合され、前記処理モジュールは、(a)前記入口点と前記出口点との間の1つまたは複数の中間ノードの集合を選択することと、(b)前記ネットワークの帯域幅に基づいて、前記入口点から1つまたは複数の中間ノードの前記集合の各ノードに送られる前記データのめいめいの分数を決定することとによって前記パスを決定する。ルータは、前記入力モジュールからのパケットを前記要求の前記パスに従ってルータの出力モジュールに転送するように適合され、前記ルータは、(c)前記入口点から1つまたは複数の中間ノードの前記集合の各ノードへ、前記決定されたためいめいの分数で前記データをルーティングし、(d)1つまたは複数の中間ノードの前記集合の各ノードから前記出口点へ前記データをルーティングするように適合される。

30

40

【0028】

もう1つの実施形態で、本発明は、複数の命令を保管されたコンピュータ可読媒体であって、前記複数の命令は、プロセッサによって実行される時に、前記プロセッサに、リンクによって相互接続され、少なくとも1つの入口点および1つの出口点を有するノードのネットワークを通過してデータをルーティングする方法であって、(a)前記入口点と前記出口点との間でのデータのルーティングに関するサービス需要を有する、パスに関する要求を受け取る工程と、(b)前記入口点と前記出口点との間の1つまたは複数の中間ノードの集合を選択する工程と、(c)前記ネットワークの帯域幅に基づいて、前記入口点から1つまたは複数の中間ノードの前記集合の各ノードに送られる前記データのめいめいの分数を決定する工程と、(d)前記入口点から1つまたは複数の中間ノードの前記集合の

50

各ノードへ、前記決定されたためいめいの分数で前記データをルーティングする工程と、(e) 1つまたは複数の中間ノードの前記集合の各ノードから前記出口点へ前記データをルーティングする工程とを含む方法を実行させる、コンピュータ可読媒体を提供する。

【 0 0 2 9 】

もう1つの実施形態で、本発明は、リンクによって相互接続され、少なくとも1つの入口点および1つの出口点を有するノードのネットワークを通してデータをルーティングするシステムを提供し、前記システムは、(i) 前記入口点と前記出口点との間でのデータのルーティングに関するサービス需要を有する、パスに関する要求と、(i i) 前記要求に関連する前記データとを受け取る手段を含む。前記システムは、さらに、(a) 前記入口点と前記出口点との間の1つまたは複数の中間ノードの集合を選択することと、(b) 前記ネットワークの帯域幅に基づいて、前記入口点から1つまたは複数の中間ノードの前記集合の各ノードに送られる前記データのめいめいの分数を決定することとによって前記要求の前記パスを決定する手段を含む。前記システムは、さらに、(c) 前記入口点から1つまたは複数の中間ノードの前記集合の各ノードへ、前記決定されたためいめいの分数で前記データをルーティングし、(d) 1つまたは複数の中間ノードの前記集合の各ノードから前記出口点へ前記データをルーティングすることによって前記要求の前記パスに従ってパケットを転送する手段を含む。

【 発明を実施するための最良の形態 】

【 0 0 3 0 】

図3に、本発明による、サービス・レベル保証を有するルーティングの方法の例示的実施形態を使用する、相互接続されたノード n_1 から n_{10} のネットワーク 300 を示す。このルーティング方法では、LSP (label - switched path) などのネットワーク・トンネル・パスの要求に関する、ネットワーク 300 を通るパスを決定する。ノード n_1 から n_{10} のそれぞれに、1つまたは複数のルータが含まれ、このルータが、本発明のルーティング方法に従って決定されたパスから構成された転送テーブルに基づいてパケットを転送する。この例示的ルーティング方法は、要求されたLSPのパケットを2フェーズでルーティングし、これによって、着信トラフィックが、まず、所定の比率で1つまたは複数の中間ノードの組に送られ、次に、その中間ノードから最終宛先に送られて、ネットワーク・スループットが最大になる。そのようなルーティング方法は、必ずしも、ネットワークを通る最短パスまたは最小ホップ数に沿って、要求されたLSPのパケットをルーティングしない場合がある。

【 0 0 3 1 】

本発明の例示的実施形態を、本明細書で、LSP要求などの関連するサービス・レベルを有するパス要求と共にMPLS標準規格を使用するネットワークに関して説明するが、本発明は、それに制限されない。本発明は、入口点と出口点の間の保証されたサービス・レベルを有するネットワーク・トンネル・パス (NTP) に関する要求が受け取られる場合など、他のコンテキストでも使用することができる。NTPは、たとえば、TCP/IPネットワーク内のパケット・フローの仮想回線、非同期転送モード (ATM) ネットワークのセルの接続、およびLSP (MPLSネットワーク内のパケットの) とすることができる。本発明は、OXCからなる再構成可能交換光バックボーンを介して接続され、これによってコア光バックボーンが光レイヤでの交換、グルーミング、および回復の機能を利用するルータのコンテキストでのIP-over-OTN (または他の回路交換ネットワーク) にも特に役立つ可能性がある。

【 0 0 3 2 】

ネットワーク 300 などの相互接続されたノードのネットワークは、 $G (N , E)$ と定義され、ここで、 N は、ノード n_1 から n_{10} の集合、 E は、ノードを相互接続するリンク (弧) または (有向) 辺の集合である。本明細書に記載の例示的実施形態では、サービス・レベルなどの使用可能なリソースの値が、リンクまたはパスの帯域幅容量であるが、その代わりにまたはそれに加えて、他の実施形態でのサービス・レベル値に、遅延、パケット消失の確率、収入、または他のサービス品質パラメータなどの1つまたは複数のリン

ク・パラメータを含めることができる。当技術分野で既知のように、これらのさまざまなサービス・レベル値の1つまたは複数を、有効帯域幅と称する量によって表すことができる。リンクの集合Eのリンク e_{ij} は、リンク e_{ij} によって接続されるノード n_i および n_j を表す2つの添字 i および j ($0 < i, j \leq N$)を有する。一般性を失わずに、各リンク e_{ij} は、有向である(パケット・フローが、ノード n_i からノード n_j に向かう)。

【0033】

図3のソース S_1 、 S_2 、および S_3 は、ノード n_1 、 n_2 、 n_3 、 n_5 、および n_9 内のルータにパケット・フローを集合的に供給するパケット・ネットワークとすることができ、これらのノードは、たとえば他のキャリアなどの外部ネットワークへの接続の潜在的な入口点である。同様に、宛先 D_1 、 D_2 、および D_3 は、ノード n_3 、 n_4 、 n_8 、 n_9 、および n_{10} 内のルータからパケット・フローを集合的に受け取るパケット・ネットワークとすることができ、これらのノードは、外部ネットワークへの接続の潜在的な出口点である。あるネットワーク内のどのノードでも、入口点および/または出口点になることができることを理解されたい。ソース S_1 、 S_2 、および S_3 は、は、入口点に接続され、宛先 D_1 、 D_2 、および D_3 は、出口点に接続される。ソース/宛先対は、(S_1 , D_1)、(S_1 , D_2)、(S_1 , D_3)、(S_2 , D_1)、(S_2 , D_2)、(S_2 , D_3)、(S_3 , D_1)、(S_3 , D_2)、および(S_3 , D_3)と定義され、各ノードは、1つまたは複数のソースおよび/あるいは1つまたは複数の宛先をサポートすることができる。ノード n_1 から n_{10} は、現在のネットワーク・トポロジ情報およびリンク状況情報(以下では「ネットワーク・トポロジ」と称する)を有するか、これへのアクセスを有することができる、このネットワーク・トポロジは、配布プロトコルを使用して(たとえば、OSPFプロトコルに従う制御パケットによって)ネットワークを介して供給し、配布することができる。

【0034】

ソース S_1 、 S_2 、および S_3 は、ネットワーク300内の新しいLSPまたは現在提供されているLSPに関するパケットを生成し、このパケットには、入口点/出口点对を識別するフィールド(たとえば、ソース S_1 、 S_2 、または S_3 のいずれかのアドレスと、宛先 D_1 、 D_2 、および D_3 のいずれかのアドレス)が含まれる。たとえばRsvpまたはLDPのシグナリング・パケットを使用して、帯域幅などのサービス品質(QoS)属性またはQoS保証をネットワーク要素(たとえばルータまたはノード)に通信することができるが、LSPのパケットには、QoS属性またはQoS保証に対応する1つまたは複数のサービス・レベル・パラメータの値も含めることができる。ネットワーク300を介して転送されるLSPのこれらのパケットは、MPLS標準規格に従うものとしてことができ、図2に示し、これに関して説明したものに類似するフォーマットを有することができる。

【0035】

図3に示されたネットワーク300について、9つの潜在的な入口点/出口点对(ソース/宛先対)が存在する。次の議論では、ノード n_i および n_j を相互接続する各リンク(i, j)(本明細書では e_{ij} とも称する)が、残存帯域幅(residual bandwidth)と称する関連する使用可能な容量 u_{ij} (または u_e)を有する。あるリンクの残存帯域幅 u_e は、そのリンクの総帯域幅と、現在そのリンクに割り当てられているLSPの帯域幅需要の合計との間の差である。ネットワークは、リンクの残存容量に関する情報を交換することができ(QoS PF(QoS shortest-path first)ネットワークなど)、この情報を、経路の分散計算に使用することができる。残存帯域幅は、一般に、たとえば $kbit/秒$ または $Mbit/秒$ 単位で表すことができ、あるいは、リンクの総容量の比率として表すことができる。ノード n_i および n_j を相互接続する各リンク(i, j)は、関連するリンク・コスト c_{ij} (または c_e)すなわち、特定のリンクの相対的な使用量、重要性、または他のコストに対応することができる関連するスカラー重みも有することができる。リンク・コストは、特定の入口点/出口点

10

20

30

40

50

対のスカラ重みを指す場合もある。リンク・コストを特定のリンクに割り当てて、ルーティング・アルゴリズムが、たとえば遅延、帯域幅提供のコスト、他のトラフィック・エンジニアリング考慮事項、または他の物理リンク・レイヤ考慮事項のゆえに特定のリンクを通る経路を好むか嫌うようにすることができる。

【0036】

一般に、要求は、入口点 o と出口点 t の間の、 $b d$ (「需要」 $b d$) の要求サービス・レベルを有するパスを提供し、ルーティングするために、ネットワーク 300 に達する。図 3 の例示的ネットワークでは、これを、要求された帯域幅 $b d$ $M b / 秒$ を有するソース / 宛先対の間のパス、たとえば (S_1, D_1) の提供を求める LSP 要求または他の形の NTP 要求とすることができる。LSP 要求は、将来の LSP 要求による帯域幅に関する需要の特性の先験的知識なしに、一時に 1 つが到着する可能性がある。さらに $(i) QoS$ 属性または QoS 保証の特性、 (ii) 接続到着、ホールド時間、または出発、および (iii) 他のトラフィック・エンジニアリング情報の先見的知識は、必ずしも入手可能でない。需要 $b d$ は、パケット・フローのパケットが変化する帯域幅必要に関する統計学的プロセスを表すことができるので、「等価」帯域幅値または「有効」帯域幅値とすることができる。当技術分野で既知のように、サービス・レベル (たとえば QoS) の属性または要件を、等価帯域幅値または有効帯域幅値に変換することができる。等価帯域幅値または有効帯域幅値は、たとえばピーク・パケット・レートおよび平均パケット・レート、到着時間およびホールド時間、ならびに接続持続時間などに基づく統計学的変数を近似する決定的な値である。

【0037】

本発明によるルーティング方法は、入口点 / 出口点对の間のネットワークを通る 1 つまたは複数のパスに沿った LSP を評価し、ルーティングする。集合 P は、ネットワーク $G(N, E)$ に含まれる特定の (区別可能な) ノード入口点 / 出口点对の集合であり、ネットワーク $G(N, E)$ は、潜在的なソース / 宛先対 (S_1, D_1) 、 (S_1, D_2) 、...、 (S_3, D_3) である。集合 P の要素を、 (s, d) と表し (すなわち、 $(s, d) \in P$)、 s および d は、それぞれ、ソース・ネットワークおよび宛先ネットワークに対応する。要素 (s, d) の間に、複数の LSP を提供することができる。

【0038】

ネットワーク 300 への LSP 要求は、集中ネットワーク管理システム (図 3 に図示せず) を介して、または分散プロトコルに従ってネットワーク 300 のノード n_1 から n_1 に供給される制御メッセージによってのいずれかで実施することができる。集中ネットワーク管理システムおよび / または各ネットワーク・ルータは、要求された LSP に対応するネットワークを通して提供されるパスを LSP 要求が決定する例示的なルーティング方法を実施する。集中ネットワーク管理システムおよび / または各ネットワーク・ルータによるプロビジョニングによって、RSVP 制御 (たとえば、RSVP シグナリング・プロトコルの QoS 要求) が、たとえば必要な帯域幅または他のタイプのサービス・レベルを有する 1 つまたは複数の接続 (パケット・フロー) を確立できるようになる。

【0039】

ノード - 弧接続行列 M は、 $(n \times e)$ 行列 (n は、集合 N の要素数と等しく、 e は、集合 E の要素数と等しい) と定義され、各行が、集合 N の異なるノード n に対応し、各列が、集合 E の異なるリンク e に対応する。各列は、ノード n_i と n_j の間の対応するリンク e_{ij} の 2 つの非 0 項目 (i, j) を有する。リンク e_{ij} に対応する列は、行 i の「+1」の値と、行 j の「-1」の値と、他のすべての行に対応する各位置の「0」の値を有する。

【0040】

ネットワークの入口 (または出口) ノードに入る (または出る) トラフィックの総量は、そのノードのすべての外部入口 (または出口) リンク (たとえば、顧客ネットワークまたは他のキャリアへのライン・カード) の総容量によって上限を定められる。所与のノード i について、ノード i から出るトラフィックの総量に対する上限 (たとえば、帯域幅ま

10

20

30

40

50

たは他のサービス・レベル)が、 R_i であり、ノード*i*に入るトラフィックの総量に対する上限(たとえば、帯域幅または他のサービス・レベル)が、 C_i である。これらのリンク容量限度は、ルータのシャーシ内に物理的に置かれたハードウェアの最大容量などの要因に基づいてモデル化されるが、ネットワーク内のトラフィックのポイントツーポイント行列を制約する。これらの制約だけが、そのネットワークによって搬送されるトラフィックの既知の態様である場合があり、これを知ることは、トラフィック行列の行および列の合計上限を知ることと同等である、すなわち、最大の可能な行の合計は、最大の可能な発信トラフィックを示し、最大の可能な列の合計は、最大の可能な着信トラフィックを示す。したがって、ネットワークのすべての許容可能なトラフィック行列 $T = \langle t_{ij} \rangle$ は、

10

【数1】

$$\sum_{j:j \neq i}^n t_{ij} = R_i \quad \forall i \in N \quad \text{かつ} \quad (1)$$

$$\sum_{j:j \neq i} t_{ji} = C_i \quad \forall i \in N \quad (2)$$

【0041】

20

前述の式(1)および(2)では、同等性(ではなく)を検討することで十分である。というのは、そのすべての行または列の合計が所与の上限より小さいすべての行列 $T' = T(R, C)$ を、非負(非対角)項目を有する行列 T'' の加算によって、行列 $T = T' + T'' = T(R, C)$ に変換できるからである。 $T(R, C)$ は、すべての可能なトラフィック行列の集合を表す。したがって、 T をルーティングするすべてのルーティング方式は、 T' もルーティングすることができる。

所与の R_i および C_i の値について、行および列の合計だけによって指定されるすべてのそのような行列の集合 $T(R, C)$ を、次の式(3)によって表すことができる。

【数2】

$$T(R, C) = \{ \langle t_{ij} \rangle \text{ ただし } \sum_{j:j \neq i} t_{ij} = R_i \text{ かつ } \sum_{j:j \neq i} t_{ji} = C_i \quad \forall i \} \quad (3)$$

30

トラフィック分布 T を、 $T(R, C)$ の任意の行列とすることができ、経時的に変更することに留意されたい。本発明のある実施形態と一貫するルーティング・アーキテクチャで、 T に関して行う必要がある唯一の仮定が、行および列の合計上限だけによって T が指定されることであることが望ましい。したがって、本発明の1実施形態と一貫するルーティング戦略は、望ましくは、(i) $T(R, C)$ のすべての行列のルーティングを許容しなければならず、(ii) 既存の接続の再構成を必要としてはならない、すなわち $T(R, C)$ に属する限りトラフィック行列 T の変化を忘れなければならず、(iii) 帯域幅効率が良くなければならない、すなわち、ノード*i* からノード*j* への需要の $\min(R_i, C_j)$ 量の提供の普通の戦略よりはるかに多い帯域幅を使用してはならない。

40

【0042】

VPNの帯域幅要件を指定する方法の1つの既知のモデルが、その教示が参照によって本明細書に組み込まれる、N. G. Duffield、P. Goyal、A. G. Greenberg、P. P. Mishra、K. K. Ramakrishnan、J. E. van der Merwe、「A flexible model for resource management in virtual private network」、ACM SIGCOMM 1999、1999年8月に記載のhoseモデルである。このモデルでは、トラフィック行列が、部分的にのみ指定され、VPN端点*i* ごとに、 R_i および C_i だけが指定され、 R_i は、*i* が任意の時にネットワークに送り込むトラ

50

フィックの最大総帯域幅であり、 C_i は、 i が任意の時にネットワークから受け取るトラフィックの最大総帯域幅である。 VPN 用に予約されるネットワーク容量は、 R_i および C_i の値と一貫するすべての可能なトラフィック・パターンに十分なものでなければならない。

【0043】

本発明のある実施形態と一貫するルーティング方式は、洗練されたトラフィック・エンジニアリング機構または追加のネットワーク・シグナリングを必要とせずに、ネットワークが任意の（多分素早く変化する）トラフィック需要を満たすことを可能にする。実際に、ネットワークは、トラフィック分布の変化を検出する必要さえない。必要になる可能性がある、トラフィックに関する唯一の知識は、ネットワークの端で外部インターフェースに接続されるすべてのライン・カードの総容量によって課せられる限度である。

10

【0044】

図4を参照すると、本発明の1実施形態と一貫する例示的2フェーズ・ルーティング方式が、物理ビューおよび論理ビューの両方で示されている。フェーズ1(401)では、任意のノード*i*でネットワークに入るトラフィックの所定の分数 α_k が、トラフィックの最終的な宛先と独立に、1つまたは複数の中間ノード*k*に分配される。フェーズ2(402)では、各ノード*k*が、異なる宛先宛のトラフィックを受け取り、受け取ったトラフィックをめいめいの宛先にルーティングする。このルーティング方式を実施する方法の1つが、一部がフェーズ1トラフィックを搬送し、残りがフェーズ2トラフィックを搬送する、ノードの間の固定帯域幅トンネルを形成することである。この2フェーズ・ルーティング戦略が働くのは、これらのトンネルに必要な帯域幅が、トラフィック行列の個々の項目ではなく、 R および C だけに依存するからである。フェーズ1で、 $\alpha_1, \alpha_2, \dots, \alpha_n$ が、次の式(4)を満足するものであることに留意されたい。

20

【数3】

$$\sum_{i=1}^n \alpha_i = 1 \quad (4)$$

【0045】

この2フェーズ・ルーティング方法を、これから詳細に説明する。最大発信トラフィック R_i を有する所与のノード*i*について、ノード*i*は、各 $k \in N$ について、フェーズ1中に中間ノード*k*にこのトラフィックの $\alpha_k R_i$ 量を送る。したがって、フェーズ1の結果としてのノード*i*からノード*k*への需要は、 $\alpha_k R_i$ である。フェーズ1の終わりに、ノード*k*が、各ノード*i*から $\alpha_k R_i$ を受け取っている。行上限の合計が、列上限の合計と等しくなければならないので、すべてのソース*i*からノード*k*で受け取られる総トラフィックは、

30

【数4】

$$\sum_{i=1}^n \alpha_k R_i = \sum_{j=1}^n \alpha_k C_j$$

40

である。同一宛先へのトラフィックが所定の比率で分割されると仮定すると、フェーズ1の後に、ノード*k*で受け取られたトラフィックのうちで、ノード*j*宛のトラフィックは、 $\alpha_k t_{ij}$ である。したがって、フェーズ2中にノード*k*からノード*j*にルーティングされる必要がある総トラフィックすなわち、ノード*k*からノード*j*へのトラフィック需要は、次の式(5)で示されるものである。

【数5】

$$\sum_{i \in N} \alpha_k t_{ij} = \alpha_k C_j \quad (5)$$

50

したがって、フェーズ1で、 k は本質的に j と同一であり、フェーズ2で、 k が本質的に i なので、フェーズ1および2でのルーティングの結果としてのノード i からノード j への総需要は、 $(\sum_j R_{ij} + \sum_i C_j)$ であり、これは、行列 $T = T(R, C)$ の知識なしで導出することができる。次の3つの特性が、この2フェーズ・ルーティング方式の特性を表す。

【0046】

(i) ルーティングが、トラフィック変動に気付かない。フェーズ1および2中にルーティングされる必要がある需要は、特定のトラフィック行列 $T = T(R, C)$ に依存するのではなく、 T (すなわち集合 $T(R, C)$)を制約する行および列の合計上限だけに依存する。

10

【0047】

(ii) ルーティングされる需要が、トラフィック行列に独立である。フェーズ1および2でのルーティングの結果としてのノード i と j の間の総需要は、 $t_{ij}' = \sum_j R_{ij} + \sum_i C_j$ であり、特定の行列 $T = T(R, C)$ に依存しない。

(iii) 提供される容量が、完全に使用される。行列 $T = T(R, C)$ のそれぞれについて、このルーティング方式は、フェーズ1および2で関連するポイントツーポイント需要を完全に利用する。

【0048】

特性(ii)は、この方式が、行および列の合計上限と配分比 $\alpha_1, \alpha_2, \dots, \alpha_n$ だけに依存し、特定の行列 $T = T(R, C)$ に依存しない変換された行列 $T' = \langle t_{ij}' \rangle$ を効果的にルーティングすることによってトラフィック行列 $T = T(R, C)$ の変動性を扱い、これによって、このルーティング方式が、トラフィック分布の変化に気付かないことを暗示する。

20

【0049】

トラフィック分布が行または列の合計上限に従うことの保証は、行または列の合計上限を、ノードで外部インターフェースに接続されるライン・カード容量の合計と等しくなるようにし、これによって、物理レイヤで厳しい形で制約を実施することによって達成することができる。その代わりに、*Differentiated Service* (differentiated service) タイプのポリシー設定方式(ネットワークに入るトラフィックが、ネットワークの境界で分類され、多分条件付けられ、異なる挙動の集団に割り当てられる)によって、各入口ノードでネットワークに入る総トラフィックをレート制限し、各ノードがオーバーサブスクライブ(over-subscribe)されないことを保証することができる。

30

【0050】

したがって、本発明の1実施形態と一貫するルーティング方式で、フェーズ1中の各ソース・ノードでのルーティング決定は、ネットワーク側の状態情報(たとえば、他方のピアリング・ポイントのトラフィックがどのように変化しているか)を全く必要とせず、フェーズ2中のルーティング決定は、パケット宛先だけに基づく。さらに、このネットワークは、入口点/出口点がオーバーサブスクライブされない限り、すべてのトラフィック分布を満たすことができ、他のキャリアに接続されるライン・カードのハード・レート保証によって、またはノードでネットワークに入るトラフィックのレート制限用の*differentiated service* タイプのポリシー設定方式の実施によって、輻輳を避けることができる。さらに、このルーティング方式は、トラフィック分布の変化に気付かず、これに対して堅牢であり、エンドツーエンド帯域幅保証の提供が、リアルタイムでのネットワークの再構成を必要としない。

40

【0051】

図5の流れ図に示されているように、本発明の1実施形態と一貫するルーティング・アーキテクチャは、次の例示的な方法で実施することができる。ステップ501で、この方法は、自律システム間ピアリング合意および/または他のキャリアに接続された各ノードのライン・カードのレートを使用して、行(または列)の上限 R_i (または C_i)の計算

50

を開始する。次に、ステップ502で、トラフィック配分比 $\alpha_1, \alpha_2, \dots, \alpha_n$ を計算する（下で詳細に説明する、必要なネットワーク帯域幅を最適化する例示的アルゴリズムを使用して）。次に、ステップ503で、ノード対 i, j のそれぞれについて、フェーズ1のノード i から1つまたは複数の中間ノードへの帯域幅 R_i の1組と、フェーズ2の1つまたは複数の中間ノードからノード j への帯域幅 C_j の1組の、2組の接続（たとえば、MPLS LSP、IPトンネル、または光レイヤ回路）を提供する。次に、ステップ504で、フェーズ1および2に従ってトラフィックをルーティングする（上で詳細に説明したように、ソース・ノードと中間ノードでのローカル動作だけを必要とする）。次に、ステップ505で、differentiated serviceタイプのポリシー設定機構を使用して、各ノードでネットワークに入る総トラフィックをレート制限する。次に、ステップ506で、たとえば新しいピアリング合意または既存のピアリング合意に対する変更の結果として、行（または列）の上限 R_i （または C_i ）が変化したかどうかを判定する。上限が変化していない場合には、この方法は、ステップ504に戻って、ルーティング動作を継続する。上限が変化した場合には、ステップ507で、 α_i 配分比を再最適化し、ステップ508で、フェーズ1および2中のルーティングに関するLSP（または光レイヤ回路あるいはIPトンネル）の帯域幅をそれぞれ相応に調整し、その後、ステップ504に戻る。

【0052】

前述の方法では、同一の接続内でトラフィックが分割される場合に、同一のエンドツーエンド接続に属するパケットが、出口ノードに順序はずれで到着する可能性がある。この状況は、この方式のフェーズ1でフローごとの分割を使用することによって回避することができる。それに加えておよび/またはその代わりに、下でさらに説明するように、トラフィックのソース・ノードおよび/または宛先ノードに依存するように、トラフィック分割比 α_i を一般化することができる。

【0053】

入口/出口トラフィックに関するリンク容量および制約 R_i, C_i を有するネットワークで、ネットワーク内のすべてのリンクの最大使用率を最小にするようにルーティングすることが望ましい。リンクの使用率は、リンクのトラフィックをその容量で割ったものと定義することができる。 $\rho_i = T_i / R_i$ が、その項目に α_i を乗じられた T_i のすべてのトラフィック行列の集合を表す場合に、線形プログラムを使用して、 ρ_i のすべての行列をルーティングできる最大の乗数 λ （スループット）を見つけることができる。

【0054】

等しい分割比すなわち、 $\alpha_i = 1/n$ $i = 1, \dots, N$ の場合に、ノード i と j の間の需要は、 $(R_i + C_j) / n$ であり、この問題は、その教示が参照によって本明細書に組み込まれる、F. Shahrokhi and D. Matula、「The Maximum Concurrent Flow Problem」、Journal of ACM、37(2): 318~334、1990年で説明されているように、最大同時フロー（maximum concurrent flow）問題になる。

【0055】

本発明の1実施形態での例示的なリンク・フローベースの線形計画法定式化を、これから説明するが、この定式化では、フローが、主問題の解で増やされ、重みが、対応する双対問題に対する解で乗法的な形で更新される。主問題および双対問題とそれらの解は、次のように特性を表すことができる。

【0056】

1. 主問題が、 n 個の変数および m 個のリソース制約を有する場合に、双対問題は、 m 個の変数および n 個のリソース制約を有する。したがって、双対問題の制約行列は、主問題の制約行列の転置行列である。

2. 主制約と双対変数の間に1対1対応がある、すなわち、双対問題の変数は、主問題の不等式と対になり、主変数と双対制約も同様である。

10

20

30

40

50

3. 双対問題の目的関数は、主制約の右辺によって決定され、主問題の目的関数と双対制約の右辺も同様である。

【0057】

下の例示的な線形計画法定式化では、コモディティ・インデックス (commodity index) k が与えられ、用語「コモディティ」は、ソースと宛先の間のフローを指し、コモディティ k のソース・ノードは $s(k)$ 、宛先ノードは $d(k)$ 、コモディティ k に対応するフローの量は $f(k)$ によって表される。ベクトル $x^k(e)$ は、ネットワーク内のリンク e のコモディティ k のフローの量を表し、 $\delta^-(i)$ および $\delta^+(i)$ は、それぞれノード i での着信エッジおよび発信エッジの集合を表す。式(6)および(7)ならびに不等式(8)の制約を有する例示的なリンク・フローベースの線形計画法定式化は、次のように表される。

10

【数6】

$$\sum_{e \in \delta^-(i)} x^k(e) = \sum_{e \in \delta^+(i)} x^k(e) \quad \forall i \neq s(k), d(k), \forall k \quad (6)$$

$$\sum_{e \in \delta^+(i)} x^k(e) = a_{s(k)} C_{d(k)} + a_{d(k)} R_{s(k)} \quad i = s(k), \forall k \quad (7)$$

20

$$\sum_k x^k(e) \leq u_e \quad \forall e \in E \quad (8)$$

のもとで

【数7】

$$\sum_{i \in N} \alpha_i$$

30

を最大化すること。

上の線形プログラムに、2組の決定変数すなわち、トラフィック分離比 α_i および、 $x^k(e)$ によって表される、コモディティ k のリンク e のフローが含まれる。コモディティ k に関する需要が、 $a_{s(k)} C_{d(k)} + a_{d(k)} R_{s(k)}$ によって与えられることに留意されたい。上の線形プログラムの最適解の α_i 値は、 α_i^* によって表され、最適の目的関数値は、 λ^* によって表され、

【数8】

$$\lambda^* = \sum_i \alpha_i^*$$

40

である。 $\alpha_i^* < 1$ の場合に、この問題は実行可能である、すなわち、所与の需要をネットワーク上でルーティングすることができる。 α_i^* 値を α_i^* 倍だけ減らして、実際の分割比を得ることができ、それに沿って需要がルーティングされる明示的なパスを、上の問題の解から決定することができる。値 $\alpha_i^* < 1$ の場合に、この問題は実行不能である。この場合には、出口(または入口)制約 $R_i(C_i)$ を、 $1/\alpha_i^*$ で割ることによってスケール・ダウンすることができ、その場合に、この問題は、所与のリンク容量の下でのルーティングに関して実行可能になる。その代わりに、 $1/\alpha_i^*$ を乗ずることによってリンク容量をスケール・アップして、すべての所与の需要のルーティングに対処することができる

50

【 0 0 5 8 】

本発明の1実施形態の例示的なパス・フローベースの線形計画法定式化（高速組合せFPTAS（Fully Polynomial Time Approximation Scheme）アルゴリズムを開発することができる）を、これから説明する。次の例の定式化では、 $P_{i,j}$ が、ノード*i*からノード*j*へのすべてのパスの集合を表し、 $x(P)$ が、パス*P*上のトラフィックを表す。式（9）および不等式（10）の制約を有する例示的なパス・フローベースの線形計画法定式化は、次のように表される。

【数9】

$$\sum_{P \in P_{ij}} x(P) = \alpha_j R_i + \alpha_i C_j \quad \forall i, j \in N, i \neq j \quad (9)$$

10

$$\sum_{P: e \in P} x(P) \leq u_e \quad \forall e \in E \quad (10)$$

のもとで

【数10】

$$\sum_{i \in N} \alpha_i$$

20

を最大化すること。

【 0 0 5 9 】

ネットワークは、一般に、指数関数的な個数のパス（ネットワークのサイズに関して）を有することができるので、前述の（主）線形プログラムは、多分、指数関数的な個数の変数を有する可能性があり、その双対（下で詳細に示す）は、指数関数的な個数の制約を有する可能性がある。したがって、これらのプログラムは、中程度から大きいサイズのネットワークでの実行によく適さない可能性がある。それでも、そのような主/双対定式化は、下で説明するように、この問題の高速多項式時間組合せアルゴリズムに有用である。

【 0 0 6 0 】

高速組合せ近似アルゴリズムを使用して、任意の $\epsilon > 0$ について、最適目的関数の $(1 + \epsilon)$ 倍までの分割比を計算することができる。 ϵ の値は、解の所望の度合の最適性を提供するように選択することができる。このアルゴリズムは、FPTAS方式であることが好ましく、入力サイズおよび $1/\epsilon$ の多項式である時間で動作する。このアルゴリズムでは、各ステップで主解および双対解が維持されるので、主目的関数値と双対目的関数値の比を計算することによって、最適性ギャップを推定することができる。

30

【 0 0 6 1 】

上で式（9）および不等式（10）に示された線形プログラムの双対定式化では、変数 $w(e)$ を不等式（10）で各リンク容量制約に関連付け、変数 $\alpha_{i,j}$ を式（9）の需要制約に関連付ける。 $SP(i, j)$ は、次の式（11）に示されているように、重み $w(e)$ の下での最短パス $P \in P_{i,j}$ を表す。

40

【数11】

$$SP(i, j) = \min_{P \in P_{ij}} \sum_{e \in P} w(e) \quad (11)$$

単純化および双対変数 $\alpha_{i,j}$ の消去の後に、双対線形計画法定式化を、不等式の制約（12）から（13）を用いて、次のように書くことができる。

【数 1 2】

$$\sum_{i:i \neq k} R_i SP(i, k) + \sum_{j:j \neq k} C_j SP(k, j) \geq 1 \quad \forall k \in N \quad (12)$$

$$w(e) \geq 0 \quad \forall e \in E \quad (13)$$

のもとで

【数 1 3】

$$\sum_{e \in E} u_e w(e)$$

10

を最小化すること。

【0 0 6 2】

所与のノード k について、 $V(k)$ は、不等式 (12) の制約の左辺を表す。重み $w(e)$ を与えられれば、2つの最短パス計算すなわち、ノード k をルートとし、すべての宛先に達する最短路木の計算およびすべての他のノードからノード k に達する逆の最短路木の計算によって、 $V(k)$ を多項式時間で計算できることに留意されたい。

重み $w(e)$ を与えられれば、次の不等式 (14) が満足される場合に限り、双対問題の実行可能な解が存在する。

【数 1 4】

20

$$\min_{k \in N} V(k) \geq 1. \quad (14)$$

【0 0 6 3】

このアルゴリズムは、等しい初期重み $w(e) =$ (量 は、 に依存し、後で導出される) から始まる。次に、下のステップ (1 から 5) を、双対実行可能制約が満足されるまで繰り返す。

(1) 図 6 に示されているように、その $V(k)$ が最小になるノード

【数 1 5】

 \bar{k}

30

を計算し、これによって、リンク

【数 1 6】

 \bar{k} を識別し、すべての i について、ノード i からノード

【数 1 7】

 \bar{k} へのパス P_i を識別し、すべての j についてノード

【数 1 8】

 \bar{k}

40

からノード j へのパス Q_i を識別する。

【0 0 6 4】

(2) $e \in E$ のそれぞれについて、 $N_p(e)$ を、 P_i がリンク e を含むノード i の集合と定義し、 $N_q(e)$ を、 Q_j がリンク e を含むノード j の集合と定義する。分数 を、次の式 (15) を使用して計算する。

【数 19】

$$a = \min_{e \in E} \frac{u_e}{\sum_{i \in N_P(e)} R_i + \sum_{j \in N_Q(e)} C_j} \quad (15)$$

【0065】

(3) すべての i について、量 R_i のフローをパス P_i で送り、すべての j について、量 C_j のフローをパス Q_j で送り、リンク e で送られる総フロー $w(e)$ を、すべての $e \in E$ について計算する。リンク e のフローを、 $w(e)$ だけ増分する。

(4) すべての $e \in E$ について、 $w(e) \left(1 + \frac{w(e)}{u_e} \right)$ として重み $w(e)$ を更新する。

(5) ノード

10

【数 20】

$$\alpha_k$$

に関連する分割比

【数 21】

$$\bar{k}$$

を、 \bar{k} だけ増分する。

【0066】

20

前述の手順が終了する時に、双対実行可能制約が満足される。しかし、各ステージの元の（残余ではない）リンク容量がこの手順で使用されるので、各リンクの主容量制約を侵害している場合がある。これを直すために、分割比を均一にスケールダウンし、その結果、容量制約に従わせることができる。

【0067】

上述の例示的方法を実施するのに使用できる例示的アルゴリズムの擬似コードを、下で示す。この擬似コードでは、配列 $flow(e)$ によって、リンク e のトラフィックが記憶される。変数 G は、0 に初期化され、双対制約が満足されないままである限り、1 未満になる。ループ全体が終了した後に、各リンク e の容量制約を侵害した倍率が、計算され、配列 $scale(e)$ に格納される。最後に、 α_k の値を、最大の容量侵害倍率で割り、結果の値を、最適値として出力する。

30

【0068】

この例示的アルゴリズムに関する 2 つの定理を次に示す。

定理 1 : $L = (n - 1) \left(\sum_{i \in N} R_i + \sum_{j \in N} C_j \right)$ であり、 L' が、 R_i' および C_j' の最小の非 0 の値であり、 α および β の値が、下で示すアルゴリズムの近似係数保証に関係すると考えるならば、任意の所与の $\epsilon > 0$ について、このアルゴリズムは、次の式 (16) および (17) に関して最適値の $(1 + \epsilon)$ 倍以内の目的関数値を有する解を計算する。

【数 22】

40

$$\delta = \frac{1 + \epsilon}{L' \left[(1 + \epsilon) \frac{L}{L'} \right]^{1/\epsilon}}, \quad (16)$$

$$\epsilon = 1 - \frac{1}{\sqrt{1 + \epsilon'}} \quad (17)$$

定理 2 : 定理 1 に従って所望の近似係数保証を提供するように選択された任意の所与の $\epsilon > 0$ について、このアルゴリズムは、入力サイズおよび $1/\epsilon$ の多項式であり、

【数 2 3】

$$O\left(\frac{nm}{\epsilon}(m+n \log n) \log_{1+\epsilon} \frac{L}{L'}\right)$$

である。

次の例示的擬似コードを使用して、上で示した例示的アルゴリズムを実施することができる。

```

    k  0  k  N ;
w(e)      e  E ;
flow(e)   0  e  E ;
G  0 ;
while G < 1 do
    リンク・コスト w(e)   i, j  Nの下で i から j へのコスト SP(i, j) の
最短パスを計算する ;
    V(k)      i  k  R_i SP(i, k) +   j  k  C_j SP(k, j) ;
    G  min_k  N V(k) ;
    if G  1  break ;

```

【数 2 4】

$$\bar{k}$$

が、 $g(k)$ が最小になるノードであるものとする；

P_i が、すべての i について、 i から

【数 2 5】

$$\bar{k}$$

への最短パスであるものとする；

Q_j が、すべての j について、

【数 2 6】

$$\bar{k}$$

から j への最短パスであるものとする；

$N_P(e) = \{i : e \text{ を含む } P_i\}$ for all e ;

$N_Q(e) = \{j : e \text{ を含む } Q_j\}$ for all e ;

【数 2 7】

$$\alpha \leftarrow \min_{e \in E} \frac{u_e}{\sum_{j \in N_P(e)} R_j + \sum_{j \in N_Q(e)} C_j} ;$$

すべての i についてパス P_i で R_i フローを送り、すべての j についてパス Q_j で C_j フローを送り、すべての e についてリンク e での結果の容量使用量 $f(e)$ を計算する；

$flow(e) = flow(e) + f(e)$ for all e ;

$w(e) = w(e) (1 + f(e) / u_e)$ for all e ;

10

20

30

40

【数 2 8】

$$\alpha_k \leftarrow \alpha_k + \alpha;$$

```
end while
scale(e) = flow(e) / u_e for all e ∈ E;
scale_max = max_{e ∈ E} scale(e);
α_k = α_k / scale_max for all k ∈ N;
最適トラフィック分割比として α_k を出力する;
```

【0069】

10

定理 1 および 2 の証明および基礎になる補助定理は、次の通りである。

双対重み $w(e)$ の集合に対して、 $D(w)$ が双対目的関数値を表し、 $V(k)$ が、すべてのノード $k \in N$ に関する不等式 (12) に示された双対プログラム制約の左辺の最小値を表すものとする、この双対プログラムを解くことは、 $D(w) / V(k)$ が最小になる重み $w(e)$ の集合を見つけることと同等である。 $D(w) / V(k)$ の最適目的関数値を、 $\Gamma(w)$ によって表す、すなわち、 $\Gamma(w) = \min_{k \in N} D(w) / V(k)$ である。上の while ループの反復 t の初めの重み関数を、 w_{t-1} によって表し、 f_{t-1} は、反復 $t-1$ の終りまでの f_j ($j \in N$) (主目的関数) の値である。上で定義したように、 $L = (n-1) \sum_{i \in N} R_i + \sum_{j \in N} C_j$ であり、 L' は、 R_i および C_j の最小の非 0 の値である。このアルゴリズムは、反復 N の後に終了する。

20

補助定理 1: アルゴリズムのすべての反復 $t = 1, \dots, K$ の終りに、次の不等式 (17.1) が満足される。

【数 2 9】

$$\Gamma(w_t) \leq \delta L \prod_{j=1}^t \left[1 + \frac{\varepsilon}{\theta} (f_j - f_{j-1}) \right]$$

証明: $V(k)$ が最小になるノードは、

【数 3 0】

$$k = \bar{k}$$

30

であり、反復 t 中にフローが増やされる対応するパスは、上で定義したように、 P_i, Q_j によって表される。重みが、 $w_t(e) = w_{t-1}(e) (1 + \alpha(e) / u_e)$ $e \in E$ として更新され、 $\alpha(e)$ は、反復 t 中にリンク e で送られる総フローである。これを使用すると、次の式 (17.2) に示されているように、 $D(w_t)$ を導出することができる。

【数 3 1】

$$\begin{aligned}
D(w_t) &= \sum_{e \in E} u_e w_t(e) \\
&= \sum_{e \in E} u_e w_{t-1}(e) + \varepsilon \sum_{e \in E} \Delta(e) w_{t-1}(e) \\
&= D(w_{t-1}) + \varepsilon \sum_{e \in E} w_{t-1}(e) \left[\sum_{i \in N_P(e)} \alpha R_i + \sum_{j \in N_Q(e)} \alpha C_j \right] \\
&= D(w_{t-1}) + \varepsilon \alpha \left[\sum_i R_i \sum_{e \in P_i} w_{t-1}(e) + \sum_j C_j \sum_{e \in Q_j} w_{t-1}(e) \right] \\
&= D(w_{t-1}) + \varepsilon \alpha \Gamma(w_{t-1}) \quad . \quad (17.2)
\end{aligned}$$

10

最初の反復まで戻って各反復に上で導出された式を使用すると、 $D(w_t)$ を、次の式 (17.3) のように定義することができる。

【数 3 2】

$$D(w_t) = D(w_0) + \varepsilon \sum_{j=1}^t (f_j - f_{j-1}) \Gamma(w_{j-1}) \quad . \quad (17.3)$$

20

ここで、重み関数 $w_t - w_0$ を考慮すると、 $D(w_t - w_0) = D(w_t) - D(w_0)$ であることがわかり、パス P_i 、 Q_j のどれもが、多くとも $n-1$ ホップの長さなので、 $(w_0)_i (n-1) R_i + (w_0)_j (n-1) C_j = L$ であることもわかる。したがって、 $(w_t - w_0)_i (n-1) R_i + (w_t - w_0)_j (n-1) C_j = L$ である。は、最適の双対目的関数値なので、次の不等式 (17.4) および (17.5) が成り立つ。

【数 3 3】

$$\begin{aligned}
D(w_t) &= \sum_{e \in E} u_e w_t(e) \\
&= \sum_{e \in E} u_e w_{t-1}(e) + \varepsilon \sum_{e \in E} \Delta(e) w_{t-1}(e) \\
&= D(w_{t-1}) + \varepsilon \sum_{e \in E} w_{t-1}(e) \left[\sum_{i \in N_P(e)} \alpha R_i + \sum_{j \in N_Q(e)} \alpha C_j \right] \\
&= D(w_{t-1}) + \varepsilon \alpha \left[\sum_i R_i \sum_{e \in P_i} w_{t-1}(e) + \sum_j C_j \sum_{e \in Q_j} w_{t-1}(e) \right] \\
&= D(w_{t-1}) + \varepsilon \alpha \Gamma(w_{t-1}) \quad . \quad (17.2)
\end{aligned}$$

30

不等式 (17.5) を式 (17.3) と比較することによって、次の不等式 (17.6) を導出することができる。

40

【数 3 4】

$$\Gamma(w_t) \leq \delta L + \frac{\varepsilon}{\theta} \sum_{j=1}^t (f_j - f_{j-1}) \Gamma(w_{j-1}) \quad . \quad (17.6)$$

【0070】

補助定理 1 の特性を、不等式 (17.6) および反復回数に対する数学的帰納法を使用して証明することができる。 $w_0(e) = u_e$ 、 $(w_0)_i = L$ なので、この帰納法の基礎になる事例 (反復 $t = 1$) が成り立つことに留意されたい。次に、リンク

50

容量制約が侵害されないことを保証するために、このアルゴリズムが終了した時に主解の目的関数値 f_K をスケーリングする必要がある倍率の推定を行うことができる。

補助定理 2 : このアルゴリズムが終了する時に、主実現可能性を保証するために、主解を、多くとも次の値の倍率だけスケーリングしなければならない。

【数 3 5】

$$\log_{1+\varepsilon} \frac{1+\varepsilon}{\delta L'}$$

【 0 0 7 1】

10

証明：リンク e およびそれに関連する重み $w(e)$ を考慮すると、 $w(e)$ の値は、フローがエッジ e で増やされる時に更新される。リンク e 上のフロー増大のシーケンス（反復ごとの）は、 $\Delta_1, \Delta_2, \dots, \Delta_r$ であり、 $r \leq K$ である。リンク e でルーティングされる総フローは、その容量を κ 倍だけ超える、すなわち、

【数 3 6】

$$\sum_{t=1}^r \Delta_t = \kappa u_e$$

である。このアルゴリズムは、 $w(e) \leq 1$ の時に終了し、双対重みは、各反復の後に多くとも $1 + \varepsilon$ 倍だけ更新されるので、 $w_K(e) \leq 1 + \varepsilon$ である。上で述べた増大の直前に、少なくとも L' の係数を有する重み $w(e)$ が、 $w_K(e)$ の合計成分の 1 つである。したがって、 $L' w_K(e, f) \leq 1 + \varepsilon$ であり、 $w_K(e, f)$ の値は、次式 (17.7) によって与えられる。

20

【数 3 7】

$$w_K(e, f) = \delta \prod_{t=1}^r \left(1 + \frac{\Delta_t}{u_e} \varepsilon\right). \tag{17.7}$$

$x \geq 0$ かつすべての t について $(1 + x)^t \leq (1 + x)^{\sum_{i=1}^t x_i}$ であるという事実と、設定 $x = \varepsilon / u_e$ かつ $x_i = (\Delta_t / u_e) \varepsilon$ を使用すると、次の不等式 (17.8) および (17.9) が成り立つ。

30

【数 3 8】

$$\begin{aligned} \frac{1+\varepsilon}{L'} &\geq w_K(e, f) \geq \delta \prod_{t=1}^r (1 + \varepsilon)^{\Delta_t / u_e} \\ &\geq \delta (1 + \varepsilon)^{\sum_{t=1}^r \Delta_t / u_e} \\ &\geq \delta (1 + \varepsilon)^\kappa, \end{aligned} \tag{17.8}$$

$$\kappa \leq \log_{1+\varepsilon} \frac{1+\varepsilon}{\delta L'}. \tag{17.9}$$

40

定理 1 の証明：補助定理 1 および $1 + x \leq e^x$ ($x > 0$) という事実を使用すると、次の不等式 (17.10) を導出することができる。

【数 3 9】

$$\begin{aligned}\Gamma(w_i) &\leq \delta L \prod_{j=1}^i e^{\frac{\varepsilon}{\theta}(f_j - f_{j-1})} \\ &\leq \delta L e^{\varepsilon f_i / \theta}\end{aligned}\quad (17.10)$$

前述のステップの単純化は、 j に関する合計 ($f_j - f_{j-1}$) のテレスコーピック・キャンセルーション (telescopic cancellation) を使用する。このアルゴリズムは、反復 K の後に終了するので、 $(w) = 1$ である。したがって、次の不等式 (17.11) および (17.12) が成り立つ。

【数 4 0】

$$1 \leq \Gamma(w_K) \leq \delta L e^{\varepsilon f_i / \theta} \quad (17.11)$$

$$\frac{\theta}{f_K} \leq \frac{\varepsilon}{\ln(1/\delta L)} \quad (17.12)$$

補助定理 2 から、スケーリングの後の実現可能な主解の目的関数値は、少なくとも次の値である。

【数 4 1】

$$\frac{f_K}{\log_{1+\varepsilon} \frac{1+\varepsilon}{\delta L'}}$$

【0 0 7 2】

主解の近似係数は、多くとも、主解と双対解の間の (比) ギャップである。不等式 (17.12) を使用すると、このギャップを、次の不等式 (17.13) によって与えることができる。

【数 4 2】

$$\begin{aligned}\frac{\theta}{f_K} &\leq \frac{\varepsilon \log_{1+\varepsilon} \frac{1+\varepsilon}{\delta L'}}{\ln(1/\delta L)} \\ &\leq \frac{\varepsilon \ln \frac{1+\varepsilon}{\delta L'}}{\ln(1+\varepsilon) \ln(1/\delta L)}.\end{aligned}\quad (17.13)$$

量

【数 4 3】

$$\ln \frac{1+\varepsilon}{\delta L'} / \ln(1/\delta L)$$

は、

【数 4 4】

$$\delta = \frac{1+\varepsilon}{L'} / \left[(1+\varepsilon) \frac{L}{L'} \right]^{1/\varepsilon}$$

に関して $1 / (1 - \delta)$ と等しい。この値を使用すると、近似係数は、次の不等式 (17.14) によって上限を与えられる。

【数 4 5】

$$\frac{\varepsilon}{\ln(1+\varepsilon)} \frac{1}{(1-\varepsilon)} \leq \frac{\varepsilon}{(\varepsilon - \varepsilon^2/2)(1-\varepsilon)} \leq \frac{1}{(1-\varepsilon)^2}. \quad (17.14)$$

10

$1 + \delta = 1 / (1 - \delta)^2$ とし、 δ について解くことによって、定理 1 で述べた δ の値が得られる。

定理 2 の証明：まず、このアルゴリズムの各反復の実行時間は、ノード

【数 4 6】

 k

およびそれに関連するパス P_i, Q_j がフローの増大のために選択されている間と考えられる。このノードおよびパスの選択は、参照によって本明細書に組み込まれる、R. K. Ahuja, T. L. Magnanti, and J. B. Orlin, 「Network Flows: Theory, Algorithms, and Applications」、Prentice Hall、1993 年 2 月で説明されているように、フィボナッチ・ヒープと共にダイクストラの最短パス・アルゴリズムを使用して $O(nm + n^2 \log n)$ 時間で実施できる全対最短パス計算を伴う。1 つの反復中の他のすべての演算が、この全対最短パス計算が占める時間に吸収され (一定の倍率まで)、反復ごとに合計 $O(n(m + n \log n))$ の時間につながる。

20

【0073】

次に、各反復で、フローがパス P_i, Q_j に沿って増やされ、値が、その反復中にリンク e で送られる総フロー $f(e)$ が多くとも u_e になるものであるという事実に関して、このアルゴリズムが停止する前の反復の回数を推定する。したがって、少なくとも 1 つのリンク e で、 $f(e) = u_e$ であり、 $w(e)$ が、 $1 + \delta$ 倍だけ増える。

30

ここで、固定された $e \in E$ の重み $w(e)$ を検討する。 $w_0(e) = \frac{f(e)}{L'}$ かつ $w_k(e) = (1 + \delta)^k \frac{f(e)}{L'}$ なので、この重みを反復に関連付けることができる回数の最大値は、次式 (18) によって定義することができる。

【数 4 7】

$$\log_{1+\delta} \frac{1+\varepsilon}{\delta L'} = \frac{1}{\varepsilon} (1 + \log_{1+\varepsilon} \frac{L}{L'}) = O\left(\frac{1}{\varepsilon} \log_{1+\varepsilon} \frac{L}{L'}\right). \quad (18)$$

40

合計 m 個の $w(e)$ があるので、反復の総数は、

【数 4 8】

$$O\left(\frac{m}{\varepsilon} \log_{1+\varepsilon} \frac{L}{L'}\right)$$

を上限とする。これに反復当たりの実行時間を乗じると、全体的なアルゴリズム実行時間を、

【数 4 9】

$$O\left(\frac{nm}{\epsilon}(m+n \log n) \log_{1+\epsilon} \frac{L}{L'}\right)$$

として得ることができる。log (L / L') が、log n と R_i および C_j 値を表すのに使用されるビット数の多項式であることに留意されたい。

【0074】

次に、上で示したルーティング方式の容量性能を検討すると、ルーティング方式の「ルーティング保証」を、まず定義しなければならず、次に、T(R, C) のすべての行列をルーティングできるすべての方式のクラスの最良の可能な方式のルーティング保証と比較しなければならない。トラフィック行列に対するリンク容量および上限 R_i, C_i を有するネットワークについて、最適目的関数を、* と表し、これは、本発明の1実施形態と一貫するルーティグ方式によって * - T(R, C) のすべての行列をルーティングできることの保証を提供する、上で説明した線形問題定式化の出力である。図7からわかるように、行合計

10

【数 5 0】

$$\hat{R}$$

および列合計

20

【数 5 1】

$$\hat{C}$$

のルーティング保証の概略は、すべてのルーティング方式によって許される の最高の可能な値が

【数 5 2】

$$\hat{\lambda}$$

であるとするならば

【数 5 3】

30

$$\lambda^* \leq \hat{\lambda}$$

であり、ルーティング方式の効率は、量

【数 5 4】

$$\lambda^* / \hat{\lambda} (\leq 1)$$

によって測定することができる。

値

40

【数 5 5】

$$\hat{\lambda}$$

は、計算が困難である可能性がある。それでも、単一の行列 T = T(R, C) が存在すると仮定し、最大の乗数 λ(T) を計算する、たとえば、λ(T) · T を所与のリンク容量を有するネットワーク内で実行可能にルーティングできるように最大並列フロー定式化を使用すると、

【数 5 6】

$$\hat{\lambda} \leq \lambda(T)$$

50

であり、したがって、次の不等式 (1 9) が満足される。

【数 5 7】

$$\frac{\lambda^*}{\lambda(T)} \leq \frac{\lambda^*}{\hat{\lambda}} \leq 1 \quad (19)$$

【 0 0 7 5】

したがって、任意のトラフィック行列 $T = T(R, C)$ について、量 $\lambda^*/\lambda(T)$ は、ルーティング方式の効率の下限である。ルーティング効率に関するより近い下限を得るために、 $\lambda(T)$ が最小である行列 $T = T(R, C)$ を識別しなければならないが、1つの行列がルーティングすべきかなりの帯域幅をとるので、これは計算が困難である可能性がある。次の例示的なヒューリスティック手法を使用して、ルーティングすべき最大の帯域幅をとる行列を近似することができ、ここで、 $C(T)$ は、行列 $T = T(R, C)$ をルーティングするための最小帯域幅を表す。 $C(T)$ を最小にする行列 $T = T(R, C)$ は、線形計画法定式化を使用することによって、多項式時間で計算可能である。この問題は、容量制約を有しないので、ノード i からノード j へのトラフィックが、単一の最短パスに沿って分割不能にルーティングされると仮定することができる。 d_{ij} が、ノード i からノード j への最短パスのホップ・カウントを表すならば、ルーティングに最大の帯域幅をとるトラフィック行列 $T = T(R, C)$ の判定という問題は、式 (2 0) および (2 1) と不等式 (2 2) の制約を有する、次の例示的な線形プログラムとして定式化することができる。

【数 5 8】

$$\sum_{j \in N, j \neq i} t_{ij} = R_i \quad \forall i \in N \quad (20)$$

$$\sum_{j \in N, j \neq i} t_{ji} = C_i \quad \forall i \in N \quad (21)$$

$$t_{ij} \geq 0 \quad \forall i, j \in N \quad (22)$$

のもとで、

$\sum_{i, j \in N} d_{ij} t_{ij}$ を最大化すること。
必要な帯域幅は、この線形プログラムの目的関数であり、行および列の合計上限 ($T(R, C)$ を定義する) が、制約を形成する。この定式化を使用して、行列 T および値 $\lambda(T)$ を計算して、本発明の 1 実施形態と一貫するルーティング方式の効率の下限を提供することができる。

【 0 0 7 6】

本発明の 1 実施形態と一貫する 2 フェーズ・ルーティング方式の分割比の 2 つの変形が、次のようなルーティング方式の一般化をもたらす。

(I) ソースと宛先の両方に依存する分割比：この方式では、ノード i から発し、宛先がノード j であるトラフィックの分数

【数 5 9】

$$\alpha_k^{ij}$$

が、中間ステージでノード k にルーティングされると仮定する。第 2 フェーズでノード i と j の間に必要な容量を計算する。第 1 フェーズで、ノード i と j の間に必要な容量は、次の不等式 (2 3) によって定義される。

10

20

30

40

【数 6 0】

$$\sum_k \alpha_j^{ik} t_{ik} \leq \max_k \alpha_j^{ik} \sum_k t_{ik} \leq \max_k \alpha_j^{ik} R_i . \quad (23)$$

第 2 フェーズについて、ノード i と j の間に必要な容量は、次の不等式 (2 4) によって定義される。

【数 6 1】

$$\sum_k \alpha_i^{kj} t_{kj} \leq \max_k \alpha_i^{kj} \sum_k t_{kj} \leq \max_k \alpha_i^{kj} C_j . \quad (24) \quad 10$$

したがって、この 2 つのフェーズでノード i と j の間に必要な容量は、次の不等式 (2 5) によって定義することができる。

【数 6 2】

$$C_{ij} \geq \alpha_j^{ik} R_i + \alpha_i^{mj} C_j \quad \forall k \quad \forall m . \quad (25)$$

(I I) ソースだけに依存する分割比：この方式では、

【数 6 3】

$$\alpha_k^i$$

が、ノード i からノード k に入るトラフィックの分数を表す。フェーズ 1 でのノード i からノード j へのトラフィックの量は、

【数 6 4】

$$\alpha_j^i R_i$$

によって与えられ、これは、フェーズ 1 でのノード i と j の間の必要な容量である。ノード i と j の間の必要な容量は、次の不等式 (2 6) によって定義することができる。

【数 6 5】

$$\sum_k \alpha_k^i t_{kj} \leq \max_k \alpha_k^i \sum_k t_{kj} \leq \max_k \alpha_k^i C_j . \quad (26) \quad 30$$

したがって、ノード i と j の間の必要な総容量 C_{ij} は、次の不等式 (2 7) によって与えられる。

【数 6 6】

$$C_{ij} \geq \alpha_j^i R_i + \alpha_k^i C_j \quad \forall k . \quad (27) \quad 40$$

どちらの場合でも、制約が、線形であり、トラフィック行列の個々の要素と独立であり、行合計および列合計だけに依存することに留意されたい。

【 0 0 7 7】

上で示した実施形態を使用するシミュレーションを、キャリア・バックボーン・ネットワークを表す 2 つのネットワーク・トポロジに対して、そのサイズ範囲で実行した。図 8 からわかるように、第 1 のネットワークは、ノード n 1 から n 1 5 を含み、2 8 個の両方向リンクを有する 1 5 ノード・ネットワークである。第 2 のネットワークは、3 3 個の両方向リンクを有する 2 0 ノード・ネットワークである (図示せず) 。異なる実行について、各ネットワーク・リンクの容量を、集合 { 0 C - 3 , 0 C - 1 2 , 0 C - 4 8 , 0 C - 1 9 2 } から選択した。結果について、 R_i および C_i が等しいと仮定し、1 に正規

化した、すなわち、 $R_i = C_i = 1$ i である。行列 T は、上の 2 つのトポロジすなわち、(I) ソースと宛先の両方に依存する分割比および (I I) ソースだけに依存する分割比のそれぞれについて計算した。量 (T) は、均等 ($e q u a l$ と表す) および不等 ($u n e q u a l$ と表す) の両方のトラフィック分割比の 値の上限であり、 (T) $u n e q u a l$ $e q u a l$ である。

【 0 0 7 8 】

図 9 および 1 0 に、5 つの異なる実行に関する上の 3 つの 値のプロットを示すが、値の相対的な順序付けは、期待通りである。図からわかるように、本発明の 1 実施形態と一貫する方法のルーティング効率は、プロットされた事例のすべてで、1 . 0 に非常に近い。両方のネットワークのプロットについて、ルーティング効率は、0 . 9 から 0 . 9 9 10 まで変化し、したがって、本発明の 1 実施形態と一貫する方法が、最適に近い動作が可能であることが示される。この結果から、分割比 i が不等になることを許容される時に、ネットワーク・スループットが増加することも示される。1 5 ノード・トポロジ実行について、スループットの増加比率 ($u n e q u a l - e q u a l$) / $e q u a l$ は、1 0 % から 8 5 % まで変化する。2 0 ノード・トポロジでは、増加率が、2 % から 5 3 % まで変化する。この結果から、次の 2 つの結論が引き出される：(1) 本発明の 1 実施形態と一貫するルーティング方式は、トラフィック分布から選択された単一の行列よりかなり低くはないネットワーク・スループットで、トラフィックの不確かさ (定義されたトラフィック変動モデルの下の) がある状態で効率的にルーティングすることができる場合があり、(2) トラフィック分割比を不等にすることを許容することによって、ネットワー 20 ク・スループットを、均等分割比の場合より大きく高めることができる。

【 0 0 7 9 】

したがって、本発明の 1 実施形態と一貫するルーティング戦略が、動的ルーティング変更を必要とせず、高い容量オーバープロビジョニングを必要とせずに、ネットワークの極端なトラフィック変動性の処理に関する複数の既知の問題に対処できることがわかった。ルーティング適合なしでトラフィック変動を処理する能力は、より安定した堅牢なインターネット挙動につながるができる。本発明の 1 実施形態と一貫するルーティング方式を使用することによって、サービス・プロバイダが、(i) 単一行列のルーティングのネットワーク・スループットにかなり近いネットワーク・スループットで、(i i) リアル・タイムでトラフィック変動を検出する必要もネットワークを再構成する必要もなしに、 30 すべてのトラフィック分布 (定義されたモデルの下の) をルーティングできるようになる。したがって、本発明は、最短パス・ルーティングより実施がそれほど複雑ではない単純なネットワーク・ルーティング方式を提供し、この方式は、次の追加の有利な特性を有する：(i) この方式は、入口 - 出口リンクの容量制約の中で許容可能なすべてのトラフィック・パターンを効率的に処理でき、(i i) この方式は、ルーティング・パラメータ (リンク重みまたはルーティング・ポリシーなど) の動的再構成を必要とせずに、高いトラフィック変動性の下でネットワーク輻輳を避けることができ、(i i i) この方式は、帯幅効率がよい可能性があり、この方式は入口 - 出口容量制約を受けるすべての可能なトラフィック・パターンを処理できるが、この方式の容量要件は、単一の悪いトラフィック・ 40 パターンに対処するのに必要な方式の容量要件に近いものとすることができる。

【 0 0 8 0 】

本発明の 1 実施形態と一貫するルーティングの方法は、ネットワーク・サービス・レベル容量のより効率的な使用率、ネットワーク・ノードでのルータの輻輳の低減、およびネットワークのより高いパケット・スループットという長所の 1 つまたは複数を提供することができる。この方法は、集中ネットワーク管理システムによって、ネットワークの各ノードによって、またはこの両方によって、要求された L S P について実施することができる。結果をネットワーク・ノードに分配する集中ネットワーク管理システムを使用する実施形態は、新しいパスのプロビジョニングの調整に好ましい可能性がある。ネットワークの各ノードでの分散実施形態は、集中ネットワーク管理システムが存在しない時、および / または要求された L S P が、ネットワークを介してルーティングされる制御パケットを 50

用いて実施される分散要求である場合に、好ましい可能性がある。

【0081】

本発明の1実施形態と一貫するルーティングの方法のさまざまな機能は、回路要素を用いて実施することができ、あるいは、デジタル領域で、ソフトウェア・プログラムの処理ステップとして実施することもできる。そのようなソフトウェアは、たとえば、デジタル信号プロセッサ、マイクロコントローラ、または汎用コンピュータで使用することができる。

【0082】

本明細書で使用される用語「ルータ」が、単一のハードウェア・デバイスもしくはスイッチ・ファブリックなどの複数の相互接続されたハードウェア・デバイス、ソフトウェア要素とハードウェア要素の組合せ、またはソフトウェア・プログラムを指すことができることを理解されたい。

10

【0083】

本発明は、方法およびこれらの方法を実践する装置の形で実施することができる。本発明は、フロッピ・ディスク、CD-ROM、ハード・ドライブ、または他の機械可読記憶媒体などの有形の媒体内で実施されるプログラム・コードの形で実施することもでき、この場合に、そのプログラム・コードがコンピュータなどの機械にロードされ、これによって実行される時に、その機械が、本発明を実践する装置になる。本発明は、たとえば記憶媒体に保管されるか、機械にロードされかつ/またはこれによって実行されるか、電気配線を介する、光ファイバを介する、または電磁放射を介するなどの媒体を介して送信される、プログラム・コードの形で実施することもでき、この場合に、プログラム・コードがコンピュータなどの機械にロードされ、これによって実行される時に、その機械が、本発明を実践する装置になる。汎用プロセッサで実施される時に、プログラム・コード・セグメントは、プロセッサと組み合わされて、特定の論理回路に類似して動作する独自の装置を提供する。

20

【0084】

本明細書に記載のルーティングの例示的方法の工程が、必ずしも記載の順序で実行される必要がないことを理解されたく、そのような方法の工程の順序が、単に例示的であると理解されたい。同様に、本発明のさまざまな実施形態と一貫するルーティング方法で、追加の工程をそのような方法に含めることができ、ある工程を省略しまたは組み合わせることができる。

30

【0085】

請求項で表される本発明の趣旨および範囲から逸脱せずに、当業者が、本発明の性質を説明するために説明され、図示された部分の詳細、材料、および配置のさまざまな変更を行えることを、さらに理解されたい。

【図面の簡単な説明】

【0086】

【図1】他のパケット・ネットワークの間の通信を可能にするリンクを介して相互接続されたノードを有する従来技術の例示的バックボーン・ネットワークを示す図である。

【図2】図1のバックボーン・ネットワークによって、入口点から出口点へパケットをルーティングするのに使用されるカプセル化されたパケットを示す図である。

40

【図3】本発明の1実施形態と一貫するラベル交換パスをルーティングするサービス・レベル保証を有するルーティングの方法を使用する、相互接続されたノードのネットワークを示す図である。

【図4】本発明の1実施形態と一貫する例示的2フェーズ・ルーティング方式の物理ビューおよび論理ビューを示す図である。

【図5】本発明の1実施形態と一貫するルーティング・アーキテクチャの例示的な方法を示す流れ図である。

【図6】本発明の1実施形態と一貫する例示的な主 - 双対線形プログラムの1工程を示す図である。

50

【図7】本発明の1実施形態と一貫するアルゴリズムの行合計
【数67】

\hat{R}

および列合計

【数68】

\hat{C}

のルーティング保証の概略を示す図である。

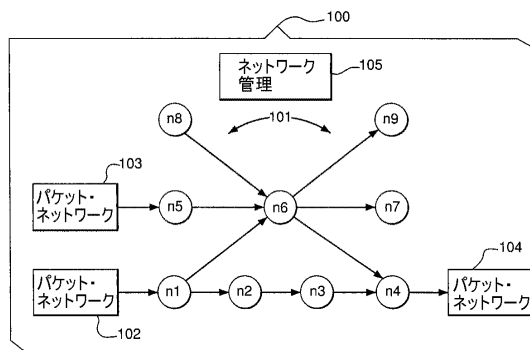
【図8】本発明の例示的实施形態のシミュレーションに使用された、キャリア・バックボーン・ネットワークを表すネットワーク・トポロジ内の28個の双方向リンクを有する例示的な15ノード・ネットワークを示す図である。

10

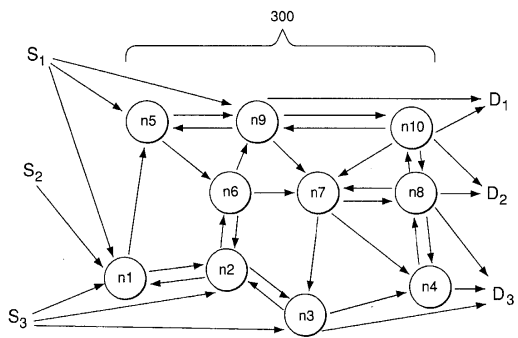
【図9】図8の例示的ネットワークのスケール係数を比較するシミュレーション結果のグラフである。

【図10】例示的な20ノード・ネットワークのスケール係数を比較するシミュレーション結果のグラフである。

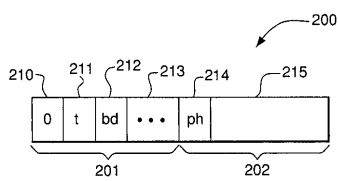
【図1】



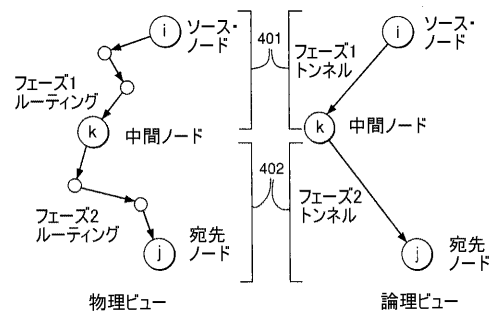
【図3】



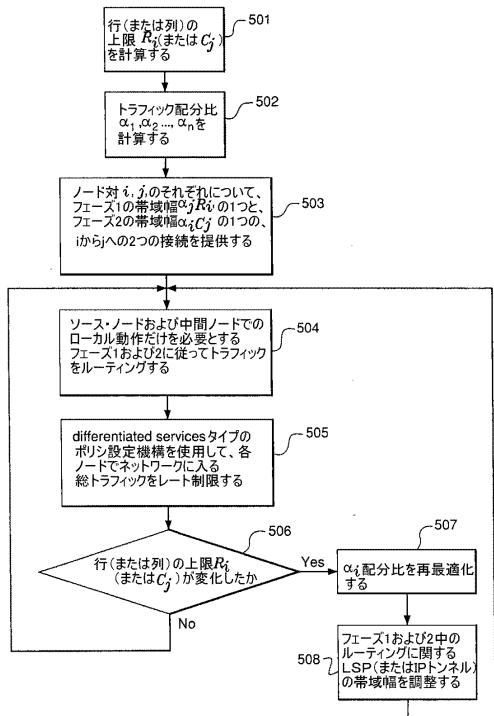
【図2】



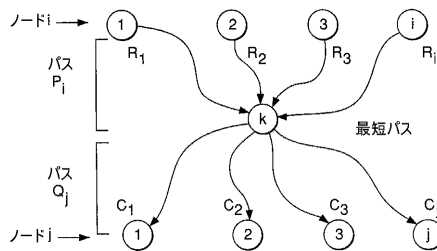
【図4】



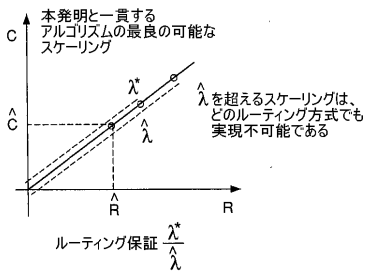
【図5】



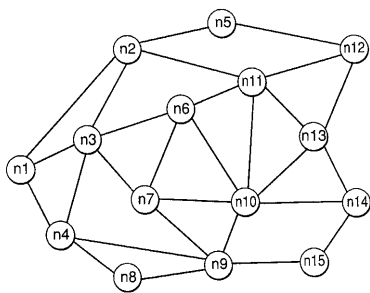
【図6】



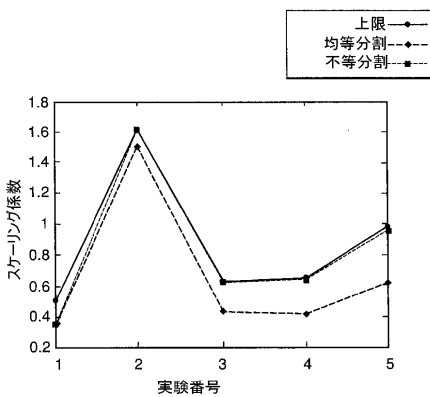
【図7】



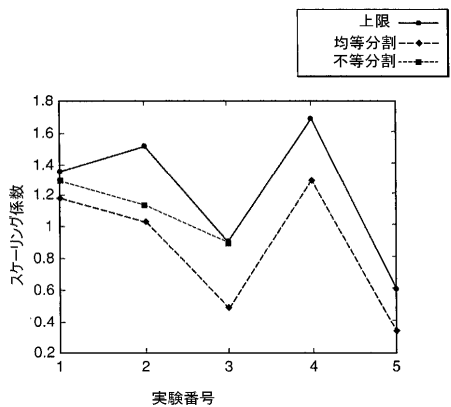
【図8】



【図10】



【図9】



フロントページの続き

- (74)代理人 100096943
弁理士 臼井 伸一
- (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 100104352
弁理士 朝日 伸光
- (74)代理人 100128657
弁理士 三山 勝巳
- (72)発明者 ムラリダラン サンパス コディアラム
アメリカ合衆国 07746 ニュージャーシィ, マールポロ, ベラ ヴィスタ コート 5
- (72)発明者 ティルネル ヴィー . ラクシュマン
アメリカ合衆国 07751 ニュージャーシィ, モーガンヴィル, レアド ドライヴ 115
- (72)発明者 スディプタ セングプタ
アメリカ合衆国 07747 ニュージャーシィ, アバディーン, ダンバートン ヒル コート
120

審査官 安藤 一道

- (56)参考文献 特開2004-048330(JP, A)
特開2004-147060(JP, A)
羽賀 太, MPLS網管理システムにおける経路設計支援機能の開発, 電子情報通信学会技術研究報告, 社団法人電子情報通信学会, 2001年11月15日, 第101巻 第442号, pp.25~28
武田 知典, コア網における公平性を考慮した動的負荷分散方式, 電子情報通信学会論文誌, 社団法人電子情報通信学会, 2003年 2月 1日, 第J86-B巻 第2号, pp.174~186

- (58)調査した分野(Int.Cl., DB名)
H04L 12/56