

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6117911号
(P6117911)

(45) 発行日 平成29年4月19日(2017.4.19)

(24) 登録日 平成29年3月31日(2017.3.31)

(51) Int.Cl. F I
HO 4 L 12/733 (2013.01) HO 4 L 12/733
HO 4 L 12/741 (2013.01) HO 4 L 12/741

請求項の数 27 (全 27 頁)

(21) 出願番号	特願2015-507625 (P2015-507625)	(73) 特許権者	598036300
(86) (22) 出願日	平成25年4月11日(2013.4.11)		テレフオンアクチーボラゲット エルエム
(65) 公表番号	特表2015-518345 (P2015-518345A)		エリクソン (パブル)
(43) 公表日	平成27年6月25日(2015.6.25)		スウェーデン国 スtockホルム エスー
(86) 国際出願番号	PCT/IB2013/052862		1 6 4 8 3
(87) 国際公開番号	W02013/160786	(74) 代理人	100095957
(87) 国際公開日	平成25年10月31日(2013.10.31)		弁理士 亀谷 美明
審査請求日	平成28年3月11日(2016.3.11)	(74) 代理人	100096389
(31) 優先権主張番号	13/458,675		弁理士 金本 哲男
(32) 優先日	平成24年4月27日(2012.4.27)	(74) 代理人	100101557
(33) 優先権主張国	米国 (US)		弁理士 萩原 康司
		(74) 代理人	100128587
			弁理士 松本 一騎

最終頁に続く

(54) 【発明の名称】 8 O 2 . 1 A Q のための3段折り返し C L O S の最適化

(57) 【特許請求の範囲】

【請求項 1】

3 段折り返し C L O S ネットワーク内の、計算上の複雑さ、ネットワーク運営、マルチキャストアドレッシング、及び当該ネットワークにおける障害のケースでの負荷再分配、についての改善された効率性を伴うイーサネットルーティングのための方法であって、

前記ネットワークは、各々がネットワークエレメントを表す複数のノードを含み、前記ノードはルートノードとエッジノードとを含み、各エッジノードは、ユーザに面する入力及び出力ポートのセットを含み、各ルートノードは、エッジノードを相互接続するためのポートのセットを含み、前記方法は、

前記ネットワーク内の各ノードにより、前記ルートノードから発するスパニングツリーについての転送ステートを計算し及びインストールするステップと、前記スパニングツリーにおけるデータ転送は、前記ネットワークに障害が無い場合にはマルチキャストパスについて任意ソースのアドレスを利用することと、

前記エッジノードのうちの所与の1つを、前記ルートノードのうちの所与の1つへの障害を起こした接続を有するものとして識別するステップと、

前記ネットワーク内の各ノード又はシステム管理機能により、前記所与のエッジノードへの接続性のためのプロトタイプとして、前記所与のエッジノードから発する最短パス優先 (S P F) ツリーを、前記所与のルートノードから発し前記障害を起こした接続を使用するスパニングツリーについての B - V I D (Backbone Virtual Local Area Network identifiers) のセットの各々について構築するステップと、前記 S P F ツリーは、前

10

20

記 B - V I D についての前記所与のエッジノードへのユニキャスト及びマルチキャストの接続性のためのプロトタイプとして供されることと、

前記 B - V I D の前記セットに関連付けられるサービスを共通して有するように識別される各ノードペアについての、前記所与のエッジノードと他のエッジノードとの間のペアごとの接続性のために、フィルタリングデータベースへデータ投入するステップと、

前記フィルタリングデータベース内にユニキャストステートをインストールするステップと、

トラフィックが前記所与のエッジノードへと向かうのか前記所与のエッジノードから来るのかに依存して選択されるマルチキャストアドレッシングのハイブリッドを用いて、前記フィルタリングデータベース内にマルチキャストステートをインストールするステップと、

10

前記フィルタリングデータベースに従って、前記ネットワークにおいてユニキャスト及びマルチキャストデータを転送するステップと、

を含む方法。

【請求項 2】

データフレームを転送する前記ステップは、

前記所与のエッジノードから前記他のエッジノードへ、同じサービス識別子を有するグループ内の前記エッジノードの全てに到達する、802.1ah マルチキャスト MAC (Medium Access Control) アドレッシングを用いて、データフレームの第 1 のセットを転送するステップと、

20

前記他のエッジノードの各々から前記所与のエッジノードへ、ユニキャスト又は 802.1aq マルチキャスト MAC アドレッシングを用いて、データフレームの第 2 のセットを転送するステップと、前記ユニキャスト及び前記 802.1aq マルチキャスト MAC アドレッシングは、前記グループ内の前記エッジノードの全てではなくサブセットへの通信を可能とすることと、

をさらに含む、請求項 1 の方法。

【請求項 3】

前記スパニングツリーの各々について、当該スパニングツリーに関連付けられるマスク値によって前記ノードのシステム ID が変換された際に、前記ネットワーク内の前記ノードの中から、最小のシステム ID を有するルートノードを選択するステップ、

30

をさらに含む、請求項 1 の方法。

【請求項 4】

各ノードにより、障害を起こしたルートノードの作業負荷が他の複数のルートノードにわたって分配されるように、スパニングツリールート選択の期間中のタイブレーキングのための複数のシステム ID を受信すること、

をさらに含む、請求項 3 の方法。

【請求項 5】

計算し及びインストールする前記ステップは、

前記ネットワーク内の所与の B - V I D についてスパニングツリーを形成するために、ダイクストラのアルゴリズムを用いて、各ルートノードについて、前記エッジノードへのそのパスを計算するステップと、

40

前記スパニングツリーから、前記エッジノードまでに 1 つよりも多くのホップを有する前記パスを除外するステップと、

前記スパニングツリー内の 3 ホップパスを、前記 3 ホップパスの障害の終点にあたるエッジノードと前記ルートノードとの間のリンク障害を有するものとして識別するステップと、

をさらに含む、請求項 1 の方法。

【請求項 6】

同じサービス識別子に関連付けられるエッジ - エッジノードペアを識別するステップと、

50

識別した前記エッジ - エッジノードペアに基づいて、前記ネットワーク内の各ノードについて前記フィルタリングデータベースを構築するステップと、

をさらに含む、請求項 1 の方法。

【請求項 7】

前記フィルタリングデータベースを構築する前記ステップは、

各エッジノードにより、I - S I D マルチキャストアドレスがスパニングツリーのルートと前記エッジノードのノードとしてのユニキャストアドレスとを指し示すように、前記フィルタリングデータベースにデータ投入するステップと、

各ルートノードにより、前記フィルタリングデータベース内の、識別した前記エッジ - エッジノードペアへのマルチキャストエントリ及び前記ルートノードのノードとしてのユニキャストアドレスへのマルチキャストエントリ、を相互接続するように前記フィルタリングデータベースにデータ投入するステップと、

をさらに含む、請求項 6 の方法。

【請求項 8】

前記スパニングツリーのサブセットを、前記ネットワーク内の他のスパニングツリーに対してプライオリティを有するものとして識別するステップと、

前記サブセットの外側での障害を前記サブセットの中へと拡散させることなく、前記サブセットの内側の障害を前記サブセットの外側へと拡散させるステップと、

をさらに含む、請求項 1 の方法。

【請求項 9】

前記ネットワーク内の前記障害は、障害を起こしたルートノード、自ノードに接続している 1 つよりも多くのリンクが障害を起こした部分的に断絶されたルートノード、単一のリンクが生き残っている事実上断絶されたルートノード、自ノードに接続している 1 つよりも多くのリンクが障害を起こした部分的に断絶されたエッジノード、又は、異なるルートノード及び異なるエッジノードへ各々接続している障害を起こした複数のリンク、を含む、請求項 1 の方法。

【請求項 10】

前記 3 段折り返し C L O S ネットワークは、データセンタ内のネットワークを表す、請求項 1 の方法。

【請求項 11】

3 段折り返し C L O S ネットワークのエッジノードとして機能するネットワークエレメントであって、計算上の複雑さ、ネットワーク運営、マルチキャストアドレッシング、及び当該ネットワークにおける障害のケースでの負荷再分配、についての改善された効率性を伴うイーサネットルーティングを使用し、前記エッジノードは、

ユーザに面する入力及び出力ポートの第 1 のセットと、

複数のルートノードへ連結される入力及び出力ポートの第 2 のセットと、

フィルタリングデータベースを記憶するメモリと、

前記メモリ、ユーザに面する入力及び出力ポートの前記第 1 のセット、並びに入力及び出力ポートの前記第 2 のセット、へ連結されるネットワークプロセッサと、

を備え、

前記ネットワークプロセッサは、

前記ルートノードから発するスパニングツリーについての転送ステートを計算し及びインストールし、

前記所与のエッジノードへの接続性のためのプロトタイプとして、前記所与のエッジノードから発する最短パス優先 (S P F) ツリーを、前記所与のルートノードから発し前記障害を起こした接続を使用するスパニングツリーについての B - V I D (Backbone Virtual Local Area Network identifiers) のセットの各々について構築し、前記 S P F ツリーは、前記 B - V I D についての前記所与のエッジノードへのユニキャスト及びマルチキャストの接続性のためのプロトタイプとして供されるものであり、

前記 B - V I D のセットに関連付けられるサービスを共通して有するように識別される

10

20

30

40

50

各ノードペアについての、前記所与のエッジノードと他のエッジノードとの間のペアごとの接続性のために、フィルタリングデータベースへデータ投入し、

前記フィルタリングデータベース内にユニキャストステートをインストールし、

トラフィックが前記所与のエッジノードへと向かうのか前記所与のエッジノードから来るのかに依存して選択されるマルチキャストアドレッシングのハイブリッドを用いて、前記フィルタリングデータベース内にマルチキャストステートをインストールし、

前記フィルタリングデータベースに従って、前記ネットワークにおいてユニキャスト及びマルチキャストデータを転送する、

ように構成される、ネットワークエレメント。

【請求項 1 2】

10

前記ネットワークエレメントは、管理システムへ連結され、前記管理システムが前記ネットワーク内の各ノードの代わりに各 B - V I D について前記 S P F を構築する、請求項 1 1 のネットワークエレメント。

【請求項 1 3】

前記ネットワークプロセッサは、

前記所与のエッジノードから前記他のエッジノードへ、同じサービス識別子を有するグループ内の前記エッジノードの全てに到達する、8 0 2 . 1 a h マルチキャスト M A C (Medium Access Control) アドレッシングを用いて、データフレームの第 1 のセットを転送し、

前記他のエッジノードの各々から前記所与のエッジノードへ、ユニキャスト又は 8 0 2 . 1 a q マルチキャスト M A C アドレッシングを用いて、データフレームの第 2 のセットを転送する、

20

ようにさらに構成され、

前記ユニキャスト及び前記 8 0 2 . 1 a q マルチキャスト M A C アドレッシングは、前記グループ内の前記エッジノードの全てではなくサブセットへの通信を可能とする、

請求項 1 1 のネットワークエレメント。

【請求項 1 4】

前記ネットワークプロセッサは、前記スパニングツリーの各々について、当該スパニングツリーに関連付けられるマスク値によって前記ルートノードのシステム I D が変換された際に、前記ルートノードの中から、最小のシステム I D を有するルートノードを選択する、

30

ようにさらに構成される、請求項 1 1 のネットワークエレメント。

【請求項 1 5】

各ノードは、障害を起こしたルートノードの作業負荷が他の複数のルートノードにわたって分配されるように、スパニングツリールート選択の期間中のタイブレーキングのための複数のシステム I D を受信する、請求項 1 4 のネットワークエレメント。

【請求項 1 6】

前記ネットワークプロセッサは

前記ネットワーク内の所与の B - V I D についてスパニングツリーを形成するために、ダイクストラのアルゴリズムを用いて、各ルートノードについて、前記エッジノードへのそのパスを計算し、

40

前記スパニングツリーから、前記エッジノードまでに 1 つよりも多くのホップを有する前記パスを除外し、

前記スパニングツリー内の 3 ホップパスを、前記ルートノードとの間の隣接関係の上でリンク障害を有するノードへ到達するものとして識別する、

ようにさらに構成される、請求項 1 1 のネットワークエレメント。

【請求項 1 7】

前記ネットワーク内の前記障害は、障害を起こしたルートノード、自ノードに接続している 1 つよりも多くのリンクが障害を起こした部分的に断絶されたルートノード、単一のリンクが生き残っている事実上断絶されたルートノード、自ノードに接続している 1 つよ

50

りも多くのリンクが障害を起こした部分的に断絶されたエッジノード、又は、異なるルートノード及び異なるエッジノードへ各々接続している障害を起こした複数のリンク、を含む、請求項 11 のネットワークエレメント。

【請求項 18】

各ルートノードは、バックボーンコアブリッジ (BCB) であり、各エッジノードは、バックボーンエッジブリッジ (EEB) である、請求項 11 のネットワークエレメント。

【請求項 19】

前記 3 段折り返し CLOS ネットワークは、データセンタ内のネットワークを表す、請求項 11 のネットワークエレメント。

【請求項 20】

3 段折り返し CLOS ネットワークのシステムであって、計算上の複雑さ、ネットワーク運営、マルチキャストアドレッシング、及び当該ネットワークにおける障害のケースでの負荷再分配、についての改善された効率性を伴うイーサネットルーティングを使用し、前記システムは、

ユーザに面する入力及び出力ポートのセットを各々が含む複数のエッジノードと、

前記複数のエッジノードを相互接続するためのポートのセットを各々が含む複数のルートノードと、

を含み、

前記ルートノード及びエッジノードの各々は、

フィルタリングデータベースを記憶するメモリと、

前記メモリへ連結されるネットワークプロセッサと、

を備え、

前記ネットワークプロセッサは、

前記ルートノードから発するスパニングツリーについての転送ステートを計算し及びインストールし、

前記所与のエッジノードへの接続性のためのプロトタイプとして、前記所与のエッジノードから発する最短パス優先 (SPF) ツリーを、前記所与のルートノードから発し前記障害を起こした接続を使用するスパニングツリーについての B - V I D (Backbone Virtual Local Area Network identifiers) のセットの各々について構築し、前記 SPF ツリーは、前記 B - V I D についての前記所与のエッジノードへのユニキャスト及びマルチキャストの接続性のためのプロトタイプとして供されるものであり、

前記 B - V I D のセットに関連付けられるサービスを共通して有するように識別される各ノードペアについての、前記所与のエッジノードと他のエッジノードとの間のペアごとの接続性のために、フィルタリングデータベースへデータ投入し、

前記フィルタリングデータベース内にユニキャストステートをインストールし、

トラフィックが前記所与のエッジノードへと向かうのか前記所与のエッジノードから来るのかに依存して選択されるマルチキャストアドレッシングのハイブリッドを用いて、前記フィルタリングデータベース内にマルチキャストステートをインストールし、

前記フィルタリングデータベースに従って、前記ネットワークにおいてユニキャスト及びマルチキャストデータを転送する、

ように構成される、システム。

【請求項 21】

前記ネットワーク内の各ノードの代わりに各 B - V I D について前記 SPF を構築する管理システム、

をさらに含む、請求項 20 のシステム。

【請求項 22】

前記ネットワークプロセッサは、

前記所与のエッジノードから前記他のエッジノードへ、同じサービス識別子を有するグループ内の前記エッジノードの全てに到達する、802.1ah マルチキャスト MAC (Medium Access Control) アドレッシングを用いて、データフレームの第 1 のセットを

10

20

30

40

50

転送し、

前記他のエッジノードの各々から前記所与のエッジノードへ、ユニキャスト又は 802.1aq マルチキャスト MAC アドレッシングを用いて、データフレームの第 2 のセットを転送する、

ようにさらに構成され、

前記ユニキャスト及び前記 802.1aq マルチキャスト MAC アドレッシングは、前記グループ内の前記エッジノードの全てではなくサブセットへの通信を可能とする、

請求項 20 のシステム。

【請求項 23】

前記ネットワークプロセッサは、前記スパニングツリーの各々について、当該スパニングツリーに関連付けられるマスク値によって前記ルートノードのシステム ID が変換された際に、前記ルートノードの中から、最小のシステム ID を有するルートノードを選択する、

ようにさらに構成される、請求項 20 のシステム。

【請求項 24】

障害を起こしたルートノードの作業負荷が他の複数のルートノードにわたって分配されるように、スパニングツリールート選択の期間中のタイブレーキングのために使用される複数のシステム ID で、各ノードを構成する管理システム、

をさらに含む、請求項 23 のシステム。

【請求項 25】

前記ネットワーク内の前記障害は、障害を起こしたルートノード、自ノードに接続している 1 つよりも多くのリンクが障害を起こした部分的に断絶されたルートノード、単一のリンクが生き残っている事実上断絶されたルートノード、自ノードに接続している 1 つよりも多くのリンクが障害を起こした部分的に断絶されたエッジノード、又は、異なるルートノード及び異なるエッジノードへ各々接続している障害を起こした複数のリンク、を含む、請求項 20 のシステム。

【請求項 26】

各ルートノードは、バックボーンコアブリッジ (BCB) であり、各エッジノードは、バックボーンエッジブリッジ (EEB) である、請求項 20 のシステム。

【請求項 27】

前記 3 段折り返し CLOS ネットワークは、データセンタ内のネットワークを表す、請求項 20 のシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明の実施形態は、コンピュータネットワーキングの分野に関し、より具体的には、3 段折り返し (three stage folded) CLOS ネットワークにおけるイーサネットルーティングの最適化に関する。

【背景技術】

【0002】

IEEE 802.1aq 標準 (これ以降、802.1aq ともいう) は、2012 年に公開され、イーサネットのためのルーティングソリューションを定義している。802.1aq は、最短パスブリッジングあるいは SPB としても知られている。802.1aq は、ネイティブなイーサネット基盤上での論理的なイーサネットネットワークの生成を可能とする。802.1aq は、トポロジー及びネットワーク内のノードの論理的なネットワークメンバシップの双方を広告するために、リンクステートプロトコルを採用する。データパケットは、802.1aq を実装するネットワークのエッジノードにおいて、mac-in-mac の 802.1ah フレーム又はタグ付き 802.1Q/p 802.1ad フレームのいずれかでカプセル化され、論理ネットワークの他のメンバにのみ移送される。ユニキャスト及びマルチキャストもまた、802.1aq によりサポートされる。全

10

20

30

40

50

てのそうしたルーティングは、対称的な最短パスを介して行われる。多くの等コストな最短パスがサポートされる。ネットワークにおける 802.1aq の実装は、プロバイダネットワーク、企業ネットワーク及びクラウドネットワークを含む多様なタイプのネットワークの生成及び構成を簡素化する。構成は比較的簡素化され、エラーの可能性、特に人為的な構成エラーの可能性が減少する。

【発明の概要】

【0003】

方法の一実施形態が、計算上の複雑さ、ネットワーク運営、マルチキャストアドレッシング、及び障害時の負荷再分配、についての改善された効率性を伴うイーサネットルーティングのために使用される 3 段折り返し CLOS ネットワークにおいて実装される。上記ネットワークは、ユーザに面する (user-facing) 入力及び出力ポートを有するエッジノードの配列に連結されるルートノードの配列を含む。上記ネットワーク内の各ノードは、上記ルートノードから発するスパニングツリーについての転送ステートを計算し及びインストールする。上記スパニングツリーにおけるデータ転送は、上記ネットワークに障害が無い場合には任意ソース (any source) のマルチキャストパスアドレスをマルチキャストパスのために利用する。これは、サービスに参加するノードの数とマルチキャストグループの数との積という組合せ積 (combinatorial product) ではなく、マルチキャストグループの数に線型的に比例してスケールするためである。所与のルートノードへの障害を起こした接続を有するものとして所与のエッジノードが識別されると、上記ネットワーク内の各ノードは、上記所与のルートノードから発し上記障害を起こした接続を使用するスパニングツリーのための各 B - V I D についての、上記所与のエッジノードから発する最短パス優先 (S P F) ツリーを計算し、上記 S P F ツリーは、上記 B - V I D についての上記所与のエッジノードへのユニキャスト及びマルチキャストの接続性のためのプロトタイプとして供される。上記所与のエッジノードと他のエッジノードとの間のペアごとの接続性のために、各ノード内のフィルタリングデータベースヘデータが投入され (populated)、それらは移動される (displaced) B - V I D に関連付けられるサービスを共通して有する。そして、ノードは、上記フィルタリングデータベース内にユニキャストステートをインストールし、トラフィックが上記所与のエッジノードへと向かうのか上記所与のエッジノードから来るのかに依存して選択されるマルチキャストアドレッシングのハイブリッドを用いて、上記フィルタリングデータベース内にマルチキャストステートをインストールする。そして、各ノードは、自身の転送データベースを用いて、上記ネットワークにおいてユニキャスト及びマルチキャストデータを転送する。

【0004】

上述した実施形態は、分散型のルーティングシステムであり、各ノードがフィルタリングデータベース内の自身の転送テーブルを計算する。代替的な実施形態において、転送テーブルを計算するシステム管理機能を、集中型のコントローラが実行する。そして、転送テーブルは、各ノードへ、当該ノードがデータ転送を実行するためにダウンロードされる。

【0005】

エッジノードとして機能するネットワークエレメントの一実施形態が、マルチキャストアドレッシング及び障害時の負荷再分配、についての改善された効率性を伴うイーサネットルーティングのために使用される 3 段折り返し CLOS ネットワークにおいて実装される。上記エッジノードは、ユーザに面する入力及び出力ポートの第 1 のセットと、複数のルートノードへ連結される入力及び出力ポートの第 2 のセットと、フィルタリングデータベースを記憶するメモリと、ネットワークプロセッサと、を備える。上記ネットワークプロセッサは、上記ルートノードから発するスパニングツリーについての転送ステートを計算し及びインストールするように構成され、上記スパニングツリーにおけるデータ転送は、上記ネットワークに障害が無い場合には任意ソースのマルチキャストパスアドレスをマルチキャストパスのために利用し、上記ネットワークプロセッサは、所与のルートノードへの障害を起こした接続を有するものとして所与のエッジノードを識別し、上記所与のエ

ッジノードから発するS P Fツリーを、上記所与のルートノードから発し上記障害を起こした接続を使用するスパニングツリーのための各B - V I Dについて構築するように構成され、上記S P Fツリーは、上記所与のエッジノードへのユニキャスト及びマルチキャストの接続性のためのプロトタイプとして供されるものであり、上記ネットワークプロセッサは、移動される上記B - V I Dに関連付けられるサービスを共通して有する上記所与のエッジノードと他のエッジノードとの間のペアごとの接続性のために、フィルタリングデータベースヘデータ投入し、上記フィルタリングデータベース内にユニキャストステートをインストールし、トラフィックが上記所与のエッジノードへと向かうのか上記所与のエッジノードから来るのかに依存して選択されるマルチキャストアドレッシングのハイブリッドを用いて、上記フィルタリングデータベース内にマルチキャストステートをインストールする、ように構成される。そして、ネットワーク内のユニキャスト及びマルチキャストデータを転送するために、転送データベースが使用される。

10

【 0 0 0 6 】

システムの一実施形態が、マルチキャストアドレッシング及び障害時の負荷再分配、についての改善された効率性を伴うイーサネットルーティングのために使用される3段階折り返しC L O Sネットワークにおいて実装される。上記システムは、ユーザに面する入力及び出力ポートのセットを有するエッジノード群へ連結されるルートノード群を含む。上記ルートノード及び上記エッジノードの各々は、フィルタリングデータベースを記憶するメモリと、ネットワークプロセッサと、を備える。上記ネットワークプロセッサは、上記ネットワークプロセッサは、上記ルートノードから発するスパニングツリーについての転送ステートを計算し及びインストールするように構成され、上記スパニングツリーにおけるデータ転送は、上記ネットワークに障害が無い場合には任意ソースのマルチキャストパスアドレスをマルチキャストパスのために利用し、上記ネットワークプロセッサは、所与のルートノードへの障害を起こした接続を有するものとして所与のエッジノードを識別し、上記所与のエッジノードから発するS P Fツリーを、上記所与のルートノードから発し上記障害を起こした接続を使用するスパニングツリーのための各B - V I Dについて構築するように構成され、上記S P Fツリーは、上記B - V I Dについての上記所与のエッジノードへのユニキャスト及びマルチキャストの接続性のためのプロトタイプとして供されるものであり、上記ネットワークプロセッサは、移動される上記B - V I Dに関連付けられるサービスを共通して有する上記所与のエッジノードと他のエッジノードとの間のペアごとの接続性のために、フィルタリングデータベースヘデータ投入し、上記フィルタリングデータベース内にユニキャストステートをインストールし、トラフィックが上記所与のエッジノードへと向かうのか上記所与のエッジノードから来るのかに依存して選択されるマルチキャストアドレッシングのハイブリッドを用いて、上記フィルタリングデータベース内にマルチキャストステートをインストールする、ように構成される。そして、ネットワーク内のユニキャスト及びマルチキャストデータを転送するために、転送データベースが使用される。

20

30

【図面の簡単な説明】

【 0 0 0 7 】

本発明は、限定ではなく例示の手法で添付図面の図において示され、図中で類似の参照符号は同様のエレメントを指し示す。なお、本開示における「一」実施形態又は「1つの」実施形態についての異なる言及は、必ずしも同じ実施形態を指すものではなく、それらの言及は、少なくとも1つの実施形態を意味することに留意すべきである。さらに、一実施形態に関連して特定の特徴、構造、又は特性が説明される場合、他の実施形態に関連してそのような特徴、構造、又は特性を作用させることは、明示されているか否かに係わらず、当業者の知識の範囲内であることを提示しておく。

40

【 0 0 0 8 】

【図1】3段階折り返しC L O Sネットワークの一例を示している。

【図2】ルートノードに割り当てられるシステムIDセットの一例を示している。

【図3】スパニングツリールート選択のために使用されるスプリットタイブレーキングの

50

仕組みの一例を示している。

【図4】3段折り返しCLOSネットワークにおいてスパニングツリーのための転送ステータスを計算し及びインストールする方法の一例を示している。

【図5】一実施形態に係るルートノードに障害が起きるシナリオを示している。

【図6】一実施形態に係るリンクに障害が起きるシナリオを示している。

【図7】一実施形態に係る2つのリンクに障害が起きる第1のシナリオを示している。

【図8】一実施形態に係る2つのリンクに障害が起きる第2のシナリオを示している。

【図9】一実施形態に係る2つのリンクに障害が起きる第3のシナリオを示している。

【図10】一実施形態に係る3つのリンクに障害が起きる第1のシナリオを示している。

【図11】一実施形態に係る3つのリンクに障害が起きる第2のシナリオを示している。

【図12】3段折り返しCLOSネットワークにおいて障害が発生した場合のデータフレーム転送のための方法の一実施形態を示している。

【図13】一実施形態に係る管理システムに連結されるネットワークエレメントを示すブロック図である。

【発明を実施するための形態】

【0009】

以下の説明において、多くの具体的な詳細が説明される。しかしながら、本発明の実施形態が、それらの具体的な詳細がなくとも実践され得ることが理解されるべきである。他の例では、本説明の理解を曖昧にすることのないように、周知の回路、構造、及び技法については詳細を示していない。しかしながら、そのような具体的な詳細がなくとも本発明が実践され得ることが、当業者によって理解されるであろう。ここに含まれる説明によって、当業者は、過度の実験をすることなく、適切な機能性を実装することができるであろう。

【0010】

IEEE 802.1aq 標準では、ネットワーク上のイーサネットフレームの転送を制御するために、リンクステートプロトコルが利用される。1つのリンクステートプロトコルである IS-IS (Intermediate System to Intermediate System) が、ネットワークのトポロジ及び論理的なネットワークメンバシップの双方を広告するために、802.1aq ネットワークにおいて使用される。

【0011】

802.1aq は、2つの動作モードを有する。VLAN (Virtual Local Area Network) ベースのネットワークのための第1のモードを、最短パスブリッジング VID (SPBV) という。MAC ベースのネットワークのための第2のモードを、最短パスブリッジング MAC (SPBM) という。SPBV ネットワーク及び SPBM ネットワークの双方は、データプレーンにおいて同時に等コスト転送ツリーのセット (ECT セット) を1つよりも多くサポートすることができる。ECT セットは、通常、複数の最短パス VLAN 識別子 (SPVID) に関連付けられ、SPBV のために SPVID セットを形成し、及び、SPBM のためにバックボーン VLAN ID (B-VID) と1対1で関連付けられる。

【0012】

802.1aq の MAC モードによれば、プロバイダネットワーク内のネットワークエレメントは、同じ宛て先アドレスを宛て先としつつも異なる B-VID にマッピングされた異なるフレームがネットワークを通じて異なるパス (“マルチパスインスタンス”という) 上で転送され得るように、B-VID によって分離されるマルチパス転送トラフィックを実行するように構成される。サービスに関連付けられる顧客のデータフレームは、802.1aq に従い、別個のサービス識別子 (I-SID) 及び B-VID を有するヘッダと共にカプセル化される。この分離は、ネットワークトポロジから独立してサービスがスケーリングされることを可能とする。よって、B-VID は、専らマルチパスインスタンスの識別子として使用されることができ、I-SID は、B-VID により識別されるマルチパスインスタンスによって提供されるべき特定のサービスを識別する。802

10

20

30

40

50

． 1 a qにおけるマルチパスインスタンスの実際のルーティングは、各ノードのシステム I Dに基づくタイブレーキングによって決定される。

【 0 0 1 3 】

W A N (Wide Area Network) 又はクラウドコンピューティングデータセンタといったデータセンタ内で、 8 0 2 . 1 a q をイーサネットルーティングのために使用することができる。データセンタ内のネットワークは、典型的には、C L O S ネットワークのトポロジーのように、非常に整然とした構造を有する。C L O S ネットワークは、スイッチングノードの配列を含む。C L O S ネットワークの一例は、入口ステージ、中間ステージ及び出口ステージを有する 3 段折り返し C L O S であり、当該ネットワークは、入口ステージが出口ステージにマージされるように、ノードの中間ステージをまたいで途中で折り返される。入口ステージ内のノードに入る各データフレームは、宛て先の出口ステージのノードへ到達するための中間ステージの利用可能なノードのいずれかを通じてルーティングされ得る。

【 0 0 1 4 】

障害の起きたノード又はリンクを有する C L O S ネットワークにおける接続性を維持するために、いくつかの技法が開発されてきた。8 0 2 . 1 a q によれば、ノード又はリンクの障害は、1 つ以上の周囲ノードにより観測され、ルーティングシステムによってネットワークにわたって広告される。ネットワーク内の各ノードは、障害により影響を受けるトラフィックについて新たなパスを再計算し、新たなパスを用いて転送が自動的に継続されるであろう。

【 0 0 1 5 】

任意のトポロジーへの、非最適な転送パスをもたらすスパニングツリーの適用とは異なり、障害の無い C L O S ネットワークでは、第 2 のティアノード群（即ち、中間ステージ内のノード群）から発する（rooted）複数のスパニングツリーにより提供される接続性が、最短パスツリーに基づくものと同一である。これは、ルーティング問題の計算上の複雑さの観点において簡易的であるだけでなく、よりスケーラブルなマルチキャストアドレッシングフォーマットを使用することをも可能とする。しかしながら、第 2 のティアノードが不完全な接続性を有するようにリンクに障害が起きた場合のトラフィックの分散は、課題（problematic）であり、最短パスツリーの動作への立ち戻りをもたらし、あまりスケーラブルでないマルチキャストアドレッシングの使用を要し得る。

【 0 0 1 6 】

さらに、8 0 2 . 1 a q によれば、障害が存在する場合、移動を強いられるトラフィックはまとめてフェイルオーバーパスへとシフトされる。フェイルオーバーパスが今やトラフィックの大幅な増加を処理することから、トラフィックのまとまった再分配がネットワークの安定性を減少させるかもしれない、それにより性能が大幅に劣化し得る。さらに、フェイルオーバーパスへのこのまとまったトラフィックのシフトは、フェイルオーバーパス内のリンク及びノードを、それらエレメントを実質的に障害に陥らせるほど圧倒し得る。基本的な仕様 8 0 2 . 1 a q の振る舞いを修正して、障害シナリオにおけるマルチパス選択との共通性のある仕組みであって、障害シナリオでのネットワークキャパシティの劣化の公平な分配のために提供される分散的なスパニングツリーのルート選択（root election）の仕組みを提供することにより、マルチキャストアドレスのスケーラビリティを最大化し、計算上の複雑さを最小化し、及びネットワーク内のマルチパス化設計を簡易にするマルチパス化技法を適用することが望ましいであろう。

【 0 0 1 7 】

ここで説明される実施形態は、3 段折り返し C L O S ネットワークといった、フラット化されたスイッチング階層を有するネットワークトポロジーを利用する。3 段折り返し C L O S ネットワークは、入口ステージが出口ステージにマージされるように、ノードの中間ステージをまたいで途中で折り返される 3 ステージ型の C L O S ネットワークである。マージされる入口ノード／出口ステージはエッジノード（第 1 ティアノードともいう）の配列を含み、中間ステージはルートノード（第 2 ティアノードともいう）の配列を含む。ルート

ノードへ接続される１つ以上のリンクに障害が起きた場合、不必要に長い転送パスを生じさせてまで全ての接続性のためにそのルート（root）を使用し続けたり、又は“全てのペア”の計算に立ち返ったりする代わりに、特別なソリューションを前提とすることが可能であり、B - V I Dについての障害の起きたリンク上で上記ルートと通常直接的な隣接関係を有する各エッジノードについて、最短パス優先（S P F）ツリーが生成され、当該S P Fツリーはエッジノードから発する。障害の起きたルートノード以外の全てのルートノードについてスパニングツリーが生成されるといった“全てのペア”の計算は存在せず、代わりに、ルートノード又はルートノードへ接続する１つ以上のリンクの障害について、B - V I Dごとに１つのS P Fツリーのみが生成される。よって、ルートノードに関わる障害のケースでの計算上の負荷が、大幅に削減される。

10

【 0 0 1 8 】

ここで開示される実施形態は、各E T Cセットについてデータフレームが横断するスパニングツリーのルートを選択するように、スプリットタイブレーキングの仕組みと共に増強される、8 0 2 . 1 a qのタイブレーキングパス選択技法の組合せのバリエーションを利用する。１つの実施形態において、スパニングツリーのためのルートは、B - V I Dに関連付けられるX O Rマスク値でのマスク後の最小のシステムI Dを伴うルートノードである。ネットワークは、C L O S内の第２のティアノードが一貫して、エッジノードに割り当てられるシステムI Dと比較した場合にタイブレーキングにおける最小のシステムI Dを有することを保証するように運営されてきた。C L O Sネットワーク内の各ノードは、タイブレーキングのために使用される複数のシステムI Dセット内の複数のシステムI Dを有する。それらシステムI Dセット及びマスク値は、スパニングツリーの複数のルートの分散的かつ独立的な選択を可能とし、データセンタ内の負荷分散への単純なアプローチを提供する。スパニングツリーのルートの次善（second best）の選択（例えば、X O Rマスク後の２番目に小さいシステムI Dを伴うルートノード）は、各システムI Dセットにおける異なる第２のティアノードとなる。よって、第２のティアノードについて障害が起きた場合、当該ノードから発するスパニングツリーのセットを、１つよりも多くの他の第２のティアノードにわたって分配することができる。

20

【 0 0 1 9 】

折り返しC L O Sネットワークについて、ルートノードから発するスパニングツリーは、エッジノードから発する等コストツリーと同じ接続性を提供する。しかしながら、等コストツリーに対して、スパニングツリーを使用するいくつかの利点が存在する。スパニングツリーの使用は、“任意ソース（any source）”というマルチキャストアドレッシングの使用を許容し、それによりネットワーク内のマルチキャストステートが大幅に削減される。対比において、“ソース固有（source specific）”マルチキャストアドレッシングを使用する等コストツリーは、所与のマルチキャストグループについてのマルチキャストソースノードの数をSとすると、オーダ（S）だけ多くのマルチキャストアドレスをスイッチフィルタリングデータベース内にもたらし。

30

【 0 0 2 0 】

さらに、ここで開示される実施形態は、マルチキャストの接続性を構築する中で、３つの形式のバックボーン宛て先M A Cアドレスを利用し、その３つは、8 0 2 . 1 a qマルチキャストM A Cアドレッシング（“ソース固有”マルチキャストアドレッシングあるいは（S , G）アドレッシングともいう）、8 0 2 . 1 a h M A Cアドレッシング（“任意ソース”マルチキャストアドレッシングあるいは（* , G）アドレッシングともいう）、及び既存のバックボーンユニキャストトンネリングの再利用、を含む。本実施形態では、ルートノードの部分的な断絶（partial severing）を扱うために、“スプリットホライズン”ツリーアプローチが採られる。障害を起こした隣接関係（failed adjacency）にアタッチしているノードは、そうではないノードとは異なる形式のアドレッシングを使用する。障害を起こした隣接関係にアタッチしているノードは、（* , G）アドレッシングを使用し得る一方、E C Tセット内の障害を起こした隣接関係に直接的にはアタッチしていないノードは、B U M（broadcast/unknown/multicast）データフレームを、（* , G

40

50

）アドレッシング及び（*S*，*G*）（あるいはユニキャスト）アドレッシングのハイブリッドを用いてバイキャスト（bi-cast）する。

【0021】

ここで説明されるハイブリッドアドレッシングのスタイルは、ノード内に維持されるマルチキャストステートを削減することにより、ネットワークのスケラビリティを改善する。802.1ahによれば、（***，*G*）マルチキャストのためのバックボーン宛て先MACアドレスは、固定的なOUI（Organizationally Unique Identifier）（*I*-*SID*）についての全てのソースのセットである $*$ を表す）と、*I*-*SID*などのサービス識別子（*G*を表す）との接続（concatenation）として符号化される。802.1aqによれば、（*S*，*G*）マルチキャストのためのバックボーン宛て先MACアドレスは、マルチキャストツリールート（*S*を表す）と*I*-*SID*などのサービス識別子（*G*を表す）との接続として符号化される。単一の受信機が存在するユニキャストについて、バックボーン宛て先MACアドレスは、例えば46ビットの値などの、固定長のビット値である。いかなる個別のソースも自身のトラフィックに行き当たらないというスプリットホライズンに依存して、（***，*G*）マルチキャストツリーは、任意のソースノードにより、*G*内の全ての受信機へ到達するために使用されることができる。（*S*，*G*）マルチキャストツリーは、ツリーが全ての“*S*”について同じようにルーティングされる（co-routed）わけではないことによるFDBコンフリクトに起因してスプリットホライズンが不可能な最短パスツリーのために必要とされる。（*S*，*G*）マルチキャストツリーは、ツリーが各“*S*”について個別化（personalize）されることから、単一のソースによって*G*内の受信機のセット又はサブセットへ到達するために使用されることができる。従って、（***，*G*）ツリーと等価なものを構築するために複数の（*S*，*G*）ツリーを要し、結果として、それら複数の（*S*，*G*）ツリーは、1つの（***，*G*）マルチキャストツリーと同じことを実質的に行うために、ルートノード内でより多くのステート（state）を要する。ユニキャストは、ノードグループ（*G*）内に単一の受信機が存在する場合に使用される。ユニキャストの転送パスはサービス固有のマルチキャストツリーからは独立して存在するため、ユニキャストは、ステートを追加しない。従って、（***，*G*）が使用されず又は使用不能である場合、（*S*，*G*）に対しユニキャストを使用することが有益である。

【0022】

ここでは特定のバージョンの標準が説明されるが、本発明の実施形態は現時点のバージョンの標準に基づく実装には限定されず、将来のバージョンの標準が開発されればそれらと共に作動するように適合され得る。同様に、本発明の実施形態は、ここで説明される特定のプロトコルの1つとの関係で動作する実装には限定されず、イーサネットマルチエリアルルーティングネットワークにおいて他のプロトコルもまた使用され得る。

【0023】

図1は、3段折り返し Clos ネットワーク 10 の一例としてのネットワークトポロジーの図である。折り返される入口／出口ステージ（図示した下のステージ）内のノードをエッジノード 12 といい、中間ステージ（図示した上のステージ）内のノードをルートノード 11 という。ルートノード 11 及びエッジノード 12 を、総称として“ノード”という。各ノード 11、12 は、自身の入力ポートの全てを自身の出力ポートへと相互接続するスイッチングエレメントである。当該ネットワークは、フレームにより横断されるノードの最大のが3つ（即ち、エッジからルート、そしてエッジへ）であることから3段（three stage）である。各ルートノード 11 は、エッジノード 12 の全てへ接続される。各ルートノード 11 は、エッジノード 12 との間の送受信のための複数の入力／出力ポートを含み、各エッジノード 12 もまた、ルートノード 11 との間の送受信のための複数の入力／出力ポートを含む。追加的に、各エッジノード 12 は、ネットワーク 10 の外部との間でトラフィックを送受信するための、複数のユーザに面する（user-facing）入力／出力ポート 13 をも含む。

【0024】

図1の例としてのネットワークは、キャパシティが同一であるノードから作られ、8つ

のエッジノード12を相互接続する4つのルートノード11を含む。別の実施形態では、異なる数のルートノード11及びエッジノード12が含まれてもよい。1つの実施形態において、ネットワーク10は、データセンタ内にある。1つの実施形態において、各ルートノード11は、バックボーンコアブリッジ（BCB）であり、各エッジノード12は、バックボーンエッジブリッジ（EEB）である。

【0025】

一般には、3段折り返しCLOSネットワークのノンブロッキング特性は、そのルートノード及びエッジノード内のポートの数により左右され得る。例えば、ノンブロッキングなCLOSネットワークにおけるルートノード11の最大数は、ノード別ポート数を2で割った商であり、エッジノード12内のユーザに面するポート13の数は、ルートノード

10

【0026】

簡明さのために、以下の説明は、基礎となるネットワークとして、破線のボックス15内に示したようなネットワーク10の一部を使用する。理解されることとして、本技術は、様々な数のルートノード及びエッジノードを伴う3段折り返しCLOSネットワークに適用可能である。図2の実施形態に示したように、ルートノード21がラベルA、B、C及びDを付与されており、エッジノード22がラベルw、x、y及びzを付与されている。1つの実施形態において、ルートノード21には、複数のシステムIDセット（例えば、セット1、セット2及びセット3）が割り当てられ、各システムIDセットはルートノード21ごとに区別される（distinct）システムIDを含む。異なるシステムIDセット

20

【0027】

複数のシステムIDセットの使用は、B-VIDについてスパニングツリーのルートノードを選択する際のスプリットタイブレーキングを可能とする。図3を参照しながら詳細に説明されるスプリットタイブレーキングの仕組みによれば、B-VID1、5、9（v1/5/9として示されている）についてのスパニングツリーは、ノードAから発する（rooted）。図2には示していないが、ルートノードBはB-VID2、6、10についての

30

【0028】

図3は、スパニングツリールート選択のために使用されるスプリットタイブレーキングの仕組みの一実施形態を示している。スプリットタイブレーキングの仕組みを説明する前に、802.1aqにおいて定義されているタイブレーキングの仕組みを説明することが有益である。802.1aqは、通常、ネットワーク内のトラフィックの各送信元（source）から発する、対称合同的（symmetrically congruent）な最短パスツリーのフルメッシュを生成する。1つのそうしたフルメッシュは、等コストツリー（ECT）セットとして知られている。ECTセットは、通常、B-VIDに関連付けられる。ECTセットの

40

50

生成の一部としてのパス計算が1つよりも多くの等コストパスからの選択の必要性に至った場合、802.1aqは、ノードIDの辞書順の並び替えを用いて各等コストパスについての一意なパスIDを構築し、パスIDのセットをソートし、最も小さい値を選択する。追加的に、802.1aqは、各ECTセットに関連付けられるセットの値のノードIDとのXOR演算を介して複数のECTセットを生成し、各パスIDにおけるノードIDの辞書順の並び替えを改訂し、パスIDのランク付けを再度行い、最も小さい値を再選択するための手段を仕様化している。

【0029】

ここで説明されるスプリットタイブレーキングの仕組みは、802.1aqに対する改善である。スプリットタイブレーキングの仕組みは、障害の無いシナリオのみならず1つ以上のルートノードに障害が起きた場合における、CLOSネットワーク内のルートノードのセットにわたるトラフィックのより均等な分配を可能とする。スプリットタイブレーキングの仕組みは、次のように、いくつかの設計エレメントを使用する。(1)インスタンス化されるスパニングツリーの数、ルートノードの数の何らかの倍数となるように選択される。その倍率は、障害の期間中にスプリットタイブレーカの性質を活用する(leverage)ことが望ましい場合、1よりも大きい。(2)ルートノードID(即ち、システムID)は、障害の無いCLOSネットワークにおいて各ルートノードから等しい数のスパニングツリーが発するように設計される。(3)エッジノードIDは、各エッジノードが決してスパニングツリーのルートにならないように設計される。これは、単純に、エッジノードIDの上位ビットに非ゼロの値を使用し、一方でルートノードは上位ビットにゼロを有し、それによりエッジノードIDのセットの最小値がルートノードIDのセットの最大値よりも大きくなるようにすることで、達成され得る。(4)(1)において倍率が1よりも大きい場合、スプリットタイブレーカ値は、障害の起きたルートノードのスパニングツリーのルート(root)が1つよりも多くの他のルートノードにわたって分配されるように設計される。

【0030】

図3の実施形態において、システムIDセット23内のシステムIDは、バイナリで示されている。加えて、ネットワーク内の各B-VIDには、例えばマスクセット33に示した例のように、マスク値及びシステムIDセット番号が割当てられている。システムIDセット1に属するB-VID1~4の各々についてルートノードを決定するために、B-VID1~4という対応するマスク値を用いて、システムIDセット1について変換が行われる。1つの実施形態において、(所与のマスクを用いた)変換後の最小のシステムIDを有するルートノードが、当該所与のマスクに関連付けられるB-VIDについてのスパニングツリーのルートである。1つの実施形態において、上記変換はXOR演算である。例えば、マスク値0000を伴うB-VID1に関し、マスク値0000及びセット1:0000、0001、0010、0011内のシステムIDについてXORが実行される。(変換されたシステムIDである)XORの結果は、0000、0001、0010、0011である。1つの実施形態において、最小のXOR値をもたらしたルートノードが、対応するB-VIDについてのスパニングツリーのルートとして選択される。よって、B-VID1について、4つの中でXORの結果として0000が最小であるため、ルートノードはノードAである。

【0031】

ノードAの障害のケースでは、ノードAを通過するトラフィックは、上で計算されたXORの結果に従って、他のルートノードへとリルートされ得る。例えば、次に小さいXOR値をもたらすルートノードを新たな通過ルートノードとして選択することができ、即ちノードBである。従って、上で計算したXORの結果がB-VID1についてのフェイルオーバー順序を左右する。

【0032】

同様に、B-VID5~8及びB-VID9~12についてのルート選択を、同じスプリットタイブレーキングの仕組みで実行することができる。B-VID1~4、B-VI

10

20

30

40

50

D 5 ~ 8 及び B - V I D 9 ~ 1 2 は異なるセット内にあるため、各セット内のマスク値及びシステム I D は、他のセットには非依存で構成され得る。1 つの実施形態において、管理システムは、各ルートノードをシステム I D 及びその対応するシステム I D セットと共に構成し得る。管理システムは、ネットワーク内で使用される B - V I D にマスク値をも割り当て得る。システムセット I D 番号（例えば、セット 1、セット 2 又はセット 3）及び各 B - V I D についてのマスク値は、拡張 I S - I S (augmented Intermediate System to Intermediate System) h e l l o 手続を介してノード間で交換され、又は管理システムにより構成データとして各ノードへロードされ得る。

【 0 0 3 3 】

図 4 は、3 段折り返し C L O S ネットワーク内でスパニングツリーのための転送ステートを計算し及びインストールするための方法 4 0 0 の一実施形態を示すフロー図である。1 つの実施形態において、拡張 I S - I S の h e l l o 手続を用いて、ノード間でタイブレーキング情報が交換される（ブロック 4 1 0）。拡張 I S - I S 手続を用いて、例えば各 B - V I D について、パス生成アルゴリズム（例えば、スパニングツリー又は等コストツリー）、ルート選択のためのシステム I D セット、及び当該 B - V I D についてのマスク値を交換することができる。例えば、ノードに関連付けられる I S - I S スピーカは、以下を広告し得る：B - V I D 1 について、スパニングツリーを使用し、ルート選択のためにシステム I D セット 1 を使用し、マスク値 = 0 0 である；B - V I D 2 について、スパニングツリーを使用し、ルート選択のためにシステム I D セット 1 を使用し、マスク値 = 0 1 である、など。代替的な実施形態において、タイブレーキング情報は、管理システムにより各ノードへ構成され得る。タイブレーキング情報を取得した後、スパニングツリーのルートノードが、当該タイブレーキング情報に基づいて各 B - V I D について選択される（ブロック 4 2 0）。図 3 において説明した通り、所与の B - V I D により識別されるスパニングツリーについてのルートとして、全てのルートノードの中で最小のシステム I D を有するルートノードが選択されるように、システム I D の変換によってルート選択が行われる。ルート選択の後、各ルートノードについて、エッジノード群までのそのパスが識別される（ブロック 4 3 0）。1 つの実施形態において、ダイクストラのアルゴリズムを用いてネットワーク内の各ノードによりパスが計算される。その代わりに、異なるアルゴリズムが使用されてもよい。

【 0 0 3 4 】

ダイクストラのアルゴリズムが使用される実施形態において、計算の結果を用いて、リンク障害の影響を迂回するように転送を修正することができる。計算はネットワーク内の各ノードにより分散的なやり方で行われ、計算の結果は、ノードにより、ネットワーク内のそれらの位置に依存することなく使用され得る。よって、この情報を使用するためにノードは障害の起きたリンクに隣接していなくてよく、なぜならノードはグローバルのトポロジーの知識を取得し、その情報のローカルな個別化を計算するためである。ダイクストラのアルゴリズムは、プロトタイプツリーを生成し、プロトタイプツリーにおいて、1 ホップの各パスはルートノードからエッジノードへのパスであり、それらはツリー内で維持され得る。1 つより多くのホップを有するパスは、ツリーから除外（pruned）され得る（ブロック 4 4 0）。例えば、ルートからエッジへ、さらにルートへ、という 2 ホップのパスが存在するかもしれない。それらパスは、プロトタイプツリーから除外され得る。リンク障害の標識として、3 ホップパスが識別される（ブロック 4 5 0）。3 ホップパスの始点はスパニングツリーのルートであり、終点は当該ルートにより直接的に到達可能でないエッジノードである。よって、3 ホップパスは、エッジから出発するトラフィックがスパニングツリーのルートを通過できず当該エッジノードへ至るために他のルート（root）を通過しなければならないことを示す。よって、それら 3 ホップパスをも検討から除外することができ、その B - V I D 及び障害の起きたリンクに連結しているエッジノードについて、新たなルートが必要である。3 ホップパスによって到達されるノードのリストは、ルートにもはや直接的にアタッチしていないエッジノードのセットとして、将来の計算のために別個に維持される。さらに、4 ホップ以上のパスは、ネットワークに病的に（pathol

10

20

30

40

50

ogically) 障害が起きていることを意味する。なお、上述したパスの除外及び障害検出は、例に過ぎず、最適化が存在してもよい。

【0035】

スパニングツリーの生成の後、B - V I Dに関連付けられる共通のサービス識別子を有する当該B - V I Dについてのエッジノードの各ペアが識別される(ブロック460)。エッジ - エッジノードペアが識別されると、そのそれぞれのフィルタリングデータベース(F D B)が、当該ノードペアの間でユニキャスト及びマルチキャストのデータフレームを転送するための転送エントリ(転送ステートともいう)を含むように構築される(ブロック470)。それらデータフレームは、ネットワークにより提供される対応するサービスを識別するために、ヘッダ内にI - S I Dを含むことになる。エッジノードについては、スパニングツリーのルートへのI - S I Dマルチキャストアドレスと共に、当該エッジノードのノードとしての(nodal)ユニキャストB - M A C (Backbone MAC)アドレスへのI - S I Dマルチキャストアドレスをポインティングすることにより、F D Bエントリを生成することができる。ルートノードについては、エッジノードのペアへのマルチキャストエントリに加えて当該ルートノードのノードとしてのユニキャストB - M A Cへのマルチキャストエントリを相互接続することにより、F D Bエントリを生成することができる。ブロック430~ブロック470の動作は、スパニングツリーの各ルートについて繰り返される。

【0036】

図5は、図2のネットワークの一実施形態を示しており、1つのルートノード(例えば、ノードA)が障害を起こす("F"により示されている)。図3において説明したスプリットタイプブレーキングの仕組みを用いることにより、ノードA上のデータ転送の作業負荷(workload)を、B - V I D 1、5及び9について2番目に小さいX O R結果(バイナリで0001)を有する他のノードへとシフトすることができる。従って、ノードB、C及びDが、B - V I D 1、5及び9についてそれぞれスパニングツリーのルートになる。作業負荷のシフトは、エッジノードにより使用されるマルチキャストアドレッシングを変化させない。よって、エッジノードw、x、y及びzは、802.1ahの"任意ソース"マルチキャストM A C (*, G)アドレッシング(簡明さのために(*, G)アドレッシングともいう)を用いて、他のエッジノードへのマルチキャストデータフレームの転送を継続する。

【0037】

図6~図11は、1つ以上のリンクに障害が起きる、一例としての3段折り返しC L O Sネットワークにおけるいくつかのシナリオを示している。図4のブロック450において説明したように、スパニングツリーのルートからエッジノードへのパスを計算するためにダイクストラのアルゴリズムが使用されるスパニングツリー生成の際に、リンク障害が検出され得る。図4のブロック450に示したように、スパニングツリー内の3ホップパスは、リンク障害及び影響を受けるノードを示している。理解されることとして、3段折り返しC L O Sネットワークにおいてノード障害又はリンク障害の帰結を検出するために、他の方法が使用されてもよい。

【0038】

図6は、ノードAとノードWとの間のリンクに障害が発生するシナリオを示している。ノードAについてスパニングツリーが計算されると、ノードWは、3ホップパスによりサービスされるものとして現れる。リンク障害を識別した後、障害の起きたリンクを使用する各B - V I Dについてノードwから発する最短パス優先(S P F)ツリーが構築される。I S - I Sなどのリンクステートルーティングプロトコルを介して、各ノードは、ネットワーク内のトポロジー情報を学習し、当該情報をS P Fツリーを計算するために使用する。1つの実施形態において、各ノードにより、ダイクストラのアルゴリズムを用いてS P Fツリーを計算することができる。これらノードは、トポロジー情報に基づいてネットワークの同じビューを構築する。

【0039】

1つの実施形態において、S P Fツリーの構築は、どのルートノードを通過すべきかを決定するための、前述したスプリットタイブレーキングの仕組みの使用を含む。なお、これらS P Fツリーのルートはネットワークのエッジノードであり、ネットワークのルートノードはS P Fツリーの通過ノード (transit nodes) になる。S P Fツリーが構築された後、エッジノードは、B - V I Dについてエッジ - エッジノードペアにより共有される、I - S I Dのインターセクション (即ち、共通のI - S I D) を識別する。識別されるエッジノードペアに基づいて、各ノードは、後続のデータ転送のために自身のF D Bにデータ投入を行う (populate)。

【0040】

図6の例において、B - V I D 1、5及び9について作られるS P Fが、それぞれv 1、v 5及びv 9でラベリングされたリンクとして示されている。ラベルv 1 / 5 / 9を伴うノードx、y及びzへ接続するリンクは、障害が無く、従って、ノードx、y及びzにより、それらの間でデータを転送するために、(*, G) アドレッシングを用いて継続的に使用されることができる。一方、ノードwと通信するためには、ノードx、y及びzは、v 1、v 5及びv 9でラベリングされたS P Fツリーのリンクを使用する必要がある。よって、1つの実施形態において、エッジノードによりハイブリッドマルチキャストアドレッシングが使用される：ノードx、y及びzは、ノードwへマルチキャストデータを転送するために、ユニキャスト又は802.1aqの“ソース固有”マルチキャストMAC (S, G) アドレッシング (簡明さのために (S, G) アドレッシングともいう) を使用し、x、y及びzの間で使用される接続性が拡張され、一方で、ノードwは、ノードx、y及びzへマルチキャストデータを転送するために、802.1ahマルチキャストMAC (*, G) アドレッシングを使用する。これは、ノードwからのマルチキャストは、全てのピアへ到達する必要がある、x、y及びzからのマルチキャストは、それらが互いへとどのようにマルチキャストを行うかとは別に、ノードwへ到達するために特有の扱いを必要とするからである。。

【0041】

同じ原理が、図7～図11に示した以下の例の各々に適用され、即ち、エッジノードが3ホップパスによってのみスパンニングツリーから到達可能である場合、そのノードからのS P Fツリーが計算され、(*, G) アドレッシングを用いてそのノードからマルチキャストデータが転送され、ユニキャスト又は (S, G) アドレッシングを用いてそのノードへマルチキャストデータが転送される。

【0042】

図7は、ノードwへ接続している2つのリンクに障害が発生する他のシナリオを示している。このシナリオにおいて、(ノードwから発する) B - V I D 1についてのS P FツリーはもはやノードA及びBを通過しないため、ノードCが通過ノードとして (例えば、タイブレーキングにより) 使用される。B - V I D 5及びB - V I D 9についてのS P Fツリーは、以前としてノードC及びDをそれぞれ通過することができる。v 1 / 5 / 9に関するアドレッシング方式は以前として (*, G) アドレッシングである。しかしながら、ノードwと通信するために、ノードx、y及びzは、v 1 / 5 (v 1及びv 5を表す) 並びにv 9でラベリングされたS P Fツリーのリンクを使用する必要がある。よって、1つの実施形態において、ノードx、y及びzは、ノードwへマルチキャストデータを転送するためにユニキャスト又は (S, G) アドレッシングを使用し、一方で、ノードwは、ノードx、y及びzへマルチキャストデータを転送するために (*, G) アドレッシングを使用する。

【0043】

図8は、ノードAへ接続している2つのリンクに障害が発生する他のシナリオを示している。このシナリオにおいて、ノードy及びzは、それらの間で依然として、v 1 / 5 / 9リンクを介し、802.1ahマルチキャストMAC (*, G) アドレッシングを用いて通信することができる。一方、ノードx及びwと通信するために、ノードy及びzは、v 1、v 5及びv 9でラベリングされたS P Fツリーのリンクを使用する必要がある。よ

って、1つの実施形態において、共通のルートへの障害の起きた隣接関係を有する複数のノードが存在することから、ノードy及びzは、ノードx及びwへマルチキャストデータフレームを転送するために(S, G)アドレッシングを使用し、一方で、ノードw及びxは、互いに並びにノードy及びzへマルチキャストを行うために802.1ahマルチキャストMAC(*, G)アドレッシングを使用する。

【0044】

図9は、ノードA及びWの間の第1のリンクに障害が起き、並びにノードB及びxの間の第2のリンクにも障害が起きる他のシナリオを示している。ノードAについてスパニングツリーが計算されると、ノードWは、3ホップパスによりサービスされるものとして現れる。B-VID1についてのSPFツリーはノードwがノードxと通信するための(ノードBとノードxとの間の)有効なパスをもはや有しないため、この障害の起きたパスをパスw-C-xによって置換えることができ、これは1つよりも多くのルートを通過するノードwからのSPFツリーの計算の結果として判定される。ラベルv1/5/9を伴うノードx、y及びzへのリンクは、障害が無く、従って、ノードx、y及びzにより、それらの間で(*, G)アドレッシングを用いてデータを転送するために継続的に使用されることができる。一方、ノードwと通信するために、ノードx、y及びzは、v1、v5及びv9でラベリングされたSPFツリーのリンクを使用する必要がある。よって、1つの実施形態において、ノードx、y及びzは、ノードwへマルチキャストデータを転送するためにユニキャスト又は(S, G)アドレッシングを使用し、一方で、ノードwは、ノードx、y及びzへマルチキャストデータを転送するために(*, G)アドレッシングを使用する。

【0045】

図10は、ノードwへ接続する4つのリンクのうちの3つに障害が起きる他のシナリオを示している。このシナリオにおいて、ラベルv1/5/9を伴うノードx、y及びzへのリンクは、障害が無く、従って、ノードx、y及びzにより、それらの間で802.1ahマルチキャストMAC(*, G)アドレッシングを用いてデータを転送するために継続的に使用されることができる。一方、ノードwとの通信は、ノードwとノードDとの間を接続する唯一の有効なリンクを介する。よって、(v1/5/9でラベリングされた)B-VID1、5及び9についてのノードwから発する3つのSPFツリーは、全てノードDを通過する。ノードwと通信するために、ノードx、y及びzは、v1/5/9でラベリングされたSPFツリーのリンクを使用する必要がある。よって、1つの実施形態において、ノードx、y及びzは、ノードwへマルチキャストデータを転送するためにユニキャスト又は(S, G)アドレッシングを使用し、一方で、ノードwは、ノードx、y及びzへマルチキャストデータを転送するために(*, G)アドレッシングを使用する。

【0046】

図11は、ノードAへ接続する4つのリンクのうちの3つに障害が起きる他のシナリオを示している。このシナリオにおいて、ノードwとの通信は、ノードwへ接続する唯一の有効なリンク(図中のv1/5/9でラベリングされた左端のリンク)を介する。A-zの間のリンクは、ノードzが他のいずれのエッジノードへ到達するためにもこのリンクを使用できないことから、使用されない。よって、ルートノードAは、有効なデータ転送のためにエッジノードのいずれにもサービスできないことから、ネットワークから“事実上”断絶されている。このシナリオは、4つのマルチキャストアドレスを使用する: ノードw、x及びyは、802.1aqマルチキャストMAC(S, G)アドレッシングを使用し、ノードzは、802.1ahマルチキャストMAC(*, G)アドレッシングを使用する。加えて、ノードAから発する1つのスパニングツリー、並びに、ノードw、x及びyからそれぞれ発する3つのスプリットホライズン計算が存在する。

【0047】

図12は、3段折り返しCLOSネットワークにおいて障害が発生した場合のデータフレーム転送のための方法1200を示すフロー図である。1つの実施形態において、方法1200は、ルートノードから発するスパニングツリーについての転送ステートを、各ノ

10

20

30

40

50

ードが計算し及びインストールすることから開始され（ブロック1210）、ここで、スパニングツリーにおけるデータ転送は、ネットワークに障害が無い場合にはマルチキャストパスについて802.1ahのマルチキャストアドレスを利用する。（例えば、図4にて上述したようにダイクストラのアルゴリズムによって）所与のルートノードへの障害を起こした接続を所与のエッジノードが有することが検出された場合（ブロック1220）、当該所与のエッジノードから発するSPFツリーが、上記所与のエッジノードから発生した障害の起きた接続を使用するB-VIDのセットの各々について構築される（ブロック1230）。上記所与のエッジノードと他のエッジノードとにより共有される共通のI-SIDが識別され、この情報を使用して、B-VIDのセットに関連付けられるサービスを共通して有するように識別されるエッジノードペアについて、上記所与のエッジノードと他のエッジノードとの間のペアごとの接続性のために、各ノード内のFDBヘータが投入される（ブロック1240）。FDB内にユニキャストステートがインストールされる（ブロック1250）。トラフィックが所与のエッジノードへと向かうのか所与のエッジノードから来るのかに依存して選択されるマルチキャストアドレッシングのハイブリッドを用いて、FDB内にマルチキャストステートもまたインストールされる（ブロック1260）。そして、ノードは、FDBに従って、ユニキャスト及びマルチキャストデータを転送する（ブロック1270）。データは、ハイブリッドマルチキャストアドレッシングに従って転送され、それと共に、上記所与のエッジノードは802.1ahマルチキャストMACアドレッシングを用いてSPFツリーを介して他のエッジノードヘータフレームを転送し、他のエッジノードはユニキャスト又は802.1aqマルチキャストMACアドレッシングを用いてSPFツリーを介して上記所与のエッジノードヘータフレームを転送する。

【0048】

図2～図12を参照しながら上で説明した方法400及び1200は、さらに最適化されてもよい。図2の例において、B-VID1～12について12個のスパニングツリーが生成される（3つのスパニングツリーのみが示されている）。実際には4つのルートのみが存在するため、4つのスパニングツリーの計算の結果が再利用されてもよい。同様に、2つのリンクに障害が起きる最初のシナリオ（図7）及び3つのリンクに障害が起きる最初のシナリオ（図10）において、ノードwからの計算が再利用されてもよい。

【0049】

より小規模なネットワークについて、ある程度の事前の計算と、別個のシステムとしての又はノードに統合されるより簡易な内部システム管理機能とのルーティングシステムの置換えと、を計画することも可能であり、例えば、8ポートのスイッチから作られるCLOSネットワークが32個のリンクを有する。全ての単一のリンク及びノードの障害シナリオの疎なテーブルを前もって構築することもでき、当該テーブルは、B-VIDに対し、どのルート及びどのアドレッシングを使用すべきかをマッピングする。例えば、テーブル内には（232のあり得るネットワークステートに對して）約40個のエントリがある。エッジノードについては、正確なルートノードを転送エントリで指し示すことによって、FDBエントリを生成することができる。ルートノードについては、各エッジペアについての対象のI-SIDのインターセクションを判定し、それに応じてルートFDBにデータ投入することにより、FDBエントリを生成することができる。

【0050】

上の説明は、移される負荷が生き残るノードにわたって可能な限り均等に共有されるように、障害時の負荷を拡散することを扱っている。しかしながら、全ての顧客がノンブロッキングなサービスからブロッキングなサービスを有するように遷移することになるため、均等な負荷の分配が常に望ましいわけではないかもしれない。1つの実施形態において、障害シナリオの下でのノンブロッキング性を保全するために、（スパニングツリーについてのそれぞれのB-VIDにより識別される）VLANのサブセットヘプライオリティが付与され得る。プライオリティの与えられたスパニングツリーのセットの外側での障害は、それらセットの内側には拡散されない。それらセットの内側の障害は、それらセット

10

20

30

40

50

の外側へと拡散される。このようにして、顧客のサブセットに、障害に直接的に影響されない場合にはいつでも、ノンブロッキングな振る舞いを保証することができる。

【 0 0 5 1 】

さらに、その時点の (current) トラフィックパターンの知識もまた活用されてよい。例えば十分に利用されているネットワークのような、データセンタ内では、トラフィックは均等に分散しているものと想定される。アルゴリズムを、利用度の低い (under-utilized) ネットワークを活用するように開発することができる。例えば、利用度の低いツリーは、障害時にシフトされるトラフィックのより大きなシェアを受け取ることができる。従って、ネットワークは、真にノンブロッキングではなく、提供されるその時点の負荷に基づいてノンブロッキングであるだけである。

10

【 0 0 5 2 】

さらに、バックアップノード/パスの構成の変更を、ヒットレスに (即ち、データのロスなく) 提供することができる。所与のシステム ID セット内の障害時の第 2 の又は第 3 の候補としてのノードの位置を、当該システム ID セット内でのそのシステム ID の変更を介して修正することができる。管理システムは、どのようにトラフィックがシフトされるかを修正するために、それらシステム ID を構成し及び調整することができる。その時点の転送パターンが最小の XOR タイプレカ値に基づくものと仮定すると、次の最良のものの序列を、転送パターンに影響を与えることなくサービス中に修正することができる。

【 0 0 5 3 】

20

図 1 3 は、本発明の一実施形態を実装するために使用され得る 3 段折り返し CLOS ネットワークエレメントの一例を示している。ネットワークエレメント 3 1 0 は、上述した 3 段折り返し CLOS ネットワーク内の任意のノード (エッジノード又はルートノード) であってよい。

【 0 0 5 4 】

図 1 3 に示したように、ネットワークエレメント 3 1 0 は、スイッチングファブリック 3 3 0、複数のデータカード 3 3 5、受信 (Rx) インタフェース 3 4 0、送信 (Tx) インタフェース 3 5 0 及び I/O ポート 3 5 5 を含むデータプレーンを含む。Rx 及び Tx インタフェース 3 4 0 及び 3 5 0 は、I/O ポート 3 5 5 を通じて、ネットワーク内のリンクとインタフェースする。ネットワークエレメントがエッジノードである場合、I/O ポート 3 5 5 は、ネットワークの外部との間の通信を提供するための、ユーザに面する複数のポートをも含む。データカード 3 3 5 は、インタフェース 3 4 0 及び 3 5 0 上で受信されるデータについての機能を実行し、スイッチングファブリック 3 3 0 はデータカード及び I/O カードの間でデータをスイッチングする。

30

【 0 0 5 5 】

ネットワークエレメント 3 1 0 は、制御プレーンをも含み、制御プレーンは、データトラフィックの経路制御、転送及び処理をハンドリングするように構成される制御ロジックを収容する 1 つ以上のネットワークプロセッサ 3 1 5 を含む。ネットワークプロセッサ 3 1 5 は、スパニングツリールート選択のためのスプリットタイプレカを実行し、スパニングツリーについての転送ステートを計算し及びインストールし、リンク障害の発生時に SPF ツリーを計算し、並びにデータ転送のために FDB 3 2 6 にデータ投入を行うようにも構成される。制御ロジック内に、他の処理もまた実装されてよい。

40

【 0 0 5 6 】

ネットワークエレメント 3 1 0 は、FDB 3 2 6 及びトポロジーデータベース 3 2 2 を記憶するメモリ 3 2 0 を含む。トポロジーデータベース 3 2 2 は、ネットワークのリンクステートを含む、ネットワークモデル又は類似のネットワークトポロジーの表現を記憶する。FDB 3 2 6 は、1 つ以上の転送テーブル内にネットワークエレメント 3 1 0 の転送ステートを記憶させ、それらはネットワークエレメント 3 1 0 へのインカミングのトラフィックをどこへ転送すべきかを示す。

【 0 0 5 7 】

50

1つの実施形態において、ネットワークエレメント310は、管理システム380へ連結され得る。1つの実施形態において、管理システム380は、メモリ370へ連結される1つ以上のプロセッサ360を含む。プロセッサ360は、システムID及びネットワークエレメント310の動作を構成するためのロジックを含み、当該動作は、システムIDの更新及びそれによるネットワーク内の作業の分散と、スパニングツリーのサブセットへの、ネットワークのノンブロッキング特性を少なくともそれらスパニングツリーについて維持するようなプライオリティの割り当てと、を含む。1つの実施形態において、管理システム380は、各ノードについて転送テーブルを計算し、そして転送テーブルをノードへダウンロードするという、システム管理機能を実行してもよい。システム管理機能は、（破線で示されているように）オプションであり、代替的な実施形態において、分散型のルーティングシステムは、各ノードが自身の転送テーブルを計算するように、上記計算を実行してもよい。

10

【0058】

ここで説明した実施形態の利点の1つは、ネットワーク内に合計でN個のノードのあるネットワークについて、スパニングツリーのルート選択に適合されたスプリットタイブレーカの使用が、ルートの障害時に計算の負荷が大幅に削減されることを意味し、それにより、計算の複雑さは、（生き残るルートの数） $\times O(N \ln N)$ となる。これは、複雑さが $O(N^2 \ln N)$ である802.1aqに対して、大幅な改善である。ここで説明した実施形態について、計算上の複雑さは、障害の無いシナリオ及び障害シナリオの双方で削減される。

20

【0059】

リンク障害のシナリオを解決するためのスプリットホライズン型のルーテッドツリーの使用もまた、計算上の複雑さを低減し、それにより、複雑さは、（ルートの数） $\times O(N \ln N)$ + （障害リンクに隣接するエッジの数） $\times O(N \ln N)$ となる。対比として、802.1aqについての名目上の複雑さは、 $O(N^2 \ln N)$ である。ここで説明した実施形態についてあらためて言うと、計算上の複雑さは、障害の無いシナリオ及び障害シナリオの双方で削減される。

【0060】

さらに、スプリットホライズン型のルーテッドツリーの使用は、障害リンクに隣接する他のエッジノード（例えば、図6～図11の例におけるノードw）への送信を行うエッジノードだけが、それらの通常のスパニングツリー（ $*$, G）と、当該他のエッジノード（例えば、ノードw）へのユニキャストパスと、ヘバイキャストする必要があることを意味し、一方で、当該他のエッジノード自体（例えば、ノードw）は、（ $*$, G）アドレッシングを使用し続けることができる。

30

【0061】

ルート選択のためのノードとしてのタイブレーカ値の使用は、ルートが他の手段で明示的に識別される必要の無いことを意味し、よって構成ミスをするのがより少ない。

【0062】

上述の機能は、コンピュータ読取可能なメモリに記憶されてネットワークエレメントに関連付けられたコンピュータプラットフォーム上の1つ以上のプロセッサ上で実行される、プログラム命令のセットとして実現されてもよい。しかしながら、ここで説明した全てのロジックを、複数の個別の構成要素、特定用途向け集積回路（ASIC）等の集積回路、フィールドプログラマブルゲートアレイ（FPGA）若しくはマイクロプロセッサ等のプログラマブルロジックデバイスと一緒に使用されるプログラマブルロジック、ステートマシン、又は、それらのいずれの組み合わせをも含む他のいずれの装置を使用しても具現化することができることは、当業者には明らかであろう。プログラマブルロジックは、読取専用メモリチップ、コンピュータメモリ、ディスク、又は他の記憶媒体等の有形の媒体に、一時的に又は永久的に固定されることができる。また、プログラマブルロジックは、搬送波で具現化されるコンピュータデータ信号内に固定されることもでき、それにより、プログラマブルロジックを、コンピュータバス又は通信ネットワーク等のインターフェ

40

50

イス上で送信することが可能になる。このようなすべての実施形態は、本発明の範囲に含まれることが意図される。

【0063】

図4及び図12のフロー図の動作を、図1、図2及び図13の例示的な実施形態を参照しながら説明した。しかしながら、図4及び図12のフロー図に示す動作を、図1、図2及び図13を参照して議論した実施形態以外の本発明の実施形態によっても実行することができること、また、図1、図2及び図13を参照して議論した実施形態が、図4及び図12の図を参照して議論した動作とは異なる動作を実行することができること、が理解されるべきである。図4及び図12の図は、本発明の特定の実施形態によって実行される動作の具体的な順序を示しているが、この順序は例示であること（例えば、代替の実施形態では、動作が異なる順序で行われ、何らかの動作が組み合わせられ、何らかの動作が重複してもよいこと、など）が理解されるべきである。

10

【0064】

本発明の多様な実施形態は、ソフトウェア、ファームウェア、及び/又はハードウェアの多様な組み合わせを使用して実装されてよい。従って、図示した技法は、1つ以上の電子デバイス（例えば、エンドステーション、ネットワークエレメント）に記憶されかつそこで実行されるコード及びデータを使用して実装されることができる。そのような電子デバイスは、非一時的なコンピュータ読取可能な記憶媒体（例えば、磁気ディスク、光ディスク、ランダムアクセスメモリ、読出し専用メモリ、フラッシュメモリデバイス、相変化メモリ）及び一時的なコンピュータ読取可能な伝送媒体（例えば、電気的、光学的、音響的、若しくは他の形式の伝搬信号 - 搬送波、赤外信号、デジタル信号等）、といったコンピュータ読取可能な媒体を使用して、コード及びデータを記憶し並びに（内部で及び/又はネットワークを通じて他の電子デバイスとの間で）通信する。加えて、そのような電子デバイスは、典型的に、1つ以上の記憶デバイス（非一時的な機械読取可能な記憶媒体）、ユーザ入出力デバイス（例えば、キーボード、タッチスクリーン、及び/又はディスプレイ）、並びにネットワーク接続等の、1つ以上の他のコンポーネントに連結された1つ以上のプロセッサのセットを含む。そのプロセッサのセットと他のコンポーネントとの連結は、典型的に、1つ以上のバス及びブリッジ（バスコントローラとも呼ばれる）を通じて行われる。よって、所与の電子デバイスの記憶デバイスは、典型的に、その電子デバイスの1つ以上のプロセッサのセット上での実行のためのコード及び/又はデータを記憶する。

20

30

【0065】

ここで使用されるところでは、ネットワークエレメント（例えば、ルータ、スイッチ、ブリッジ、コントローラ）とは、ネットワーク上の他の機器（例えば、他のネットワークエレメント、エンドステーション）を通信可能に相互に接続する、ハードウェア及びソフトウェアを含む、1つのネットワーキング機器のことである。いくつかのネットワークエレメントは、“マルチサービスのネットワークエレメント”であり、これらは、複数のネットワーキング機能（例えば、ルーティング、ブリッジング、スイッチング、レイヤ2統合、セッションボーダ制御、サービス品質（Quality of Service）、及び/又は加入者管理）についてのサポートを提供し、及び/又は、複数のアプリケーションサービス（例えば、データ、音声、及びビデオ）についてのサポートを提供する。加入者エンドステーション（例えば、サーバ、ワークステーション、ラップトップ、ネットブック、パームトップ、携帯電話、スマートフォン、マルチメディアフォン、VOIP（Voice Over Internet Protocol）フォン、ユーザ機器、端末、携帯型メディアプレーヤ、GPSユニット、ゲームシステム、セットトップボックス）は、インターネットを通じて提供されるコンテンツ/サービス、及び/又は、インターネットにオーバレイされる（例えば、インターネット経由でトンネリングされる）仮想プライベートネットワーク（VPN）上で提供されるコンテンツ/サービスにアクセスする。それらのコンテンツ及び/又はサービスは、典型的に、サービスプロバイダ又はコンテンツプロバイダに属する1つ以上のエンドステーション（例えば、サーバエンドステーション）により、又はピアツーピアサービスに

40

50

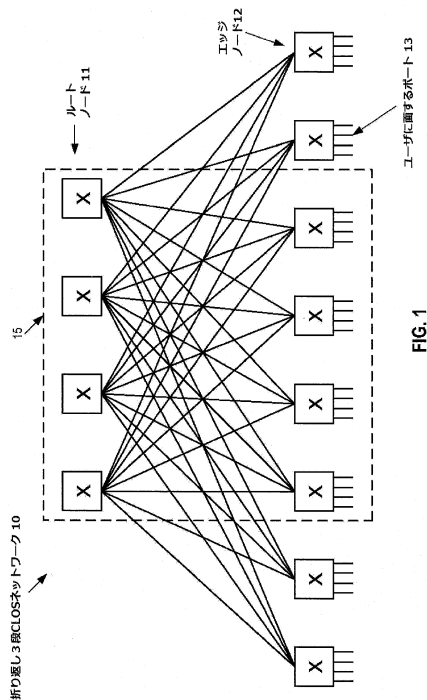
参加するエンドステーションにより提供され、例えば、パブリックウェブページ（例えば、無料コンテンツ、ストアフロント、検索サービス）、プライベートウェブページ（例えば、eメールサービスを提供する、ユーザ名/パスワードによりアクセスされるウェブページ）、及び/又は、VPN上の企業ネットワークを含み得る。典型的に、加入者エンドステーションは、（例えば、アクセスネットワークに（有線又は無線で）連結された顧客構内機器（customer premise equipment）を通じて）エッジネットワークエレメントに連結され、それらエッジネットワークエレメントは、（例えば、1つ以上のコアネットワークエレメントを通じて）他のエッジネットワークエレメントに連結され、それら他のエッジネットワークエレメントは、他のエンドステーション（例えば、サーバエンドステーション）に連結される。

10

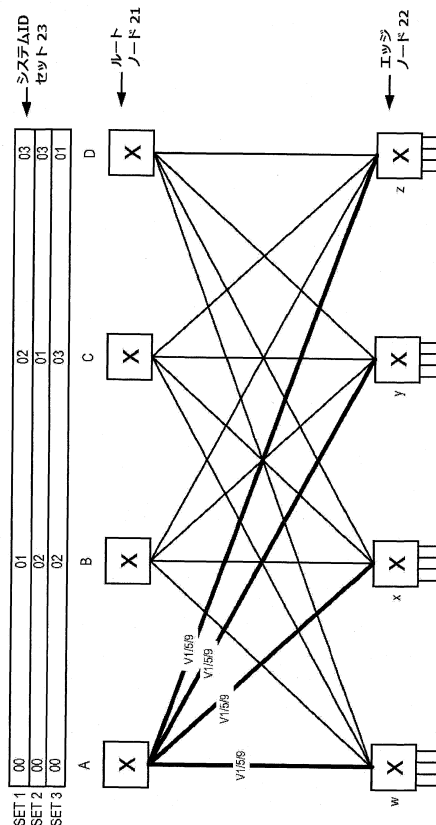
【0066】

本発明をいくつかの実施形態に関して説明してきたが、本発明が上述の実施形態に限定されるものではなく、添付の特許請求の範囲の思想および範囲内で修正及び変形と共に実践されることができ、当業者は認識するであろう。よって、本説明は、限定ではなく例示と捉えられるべきである。

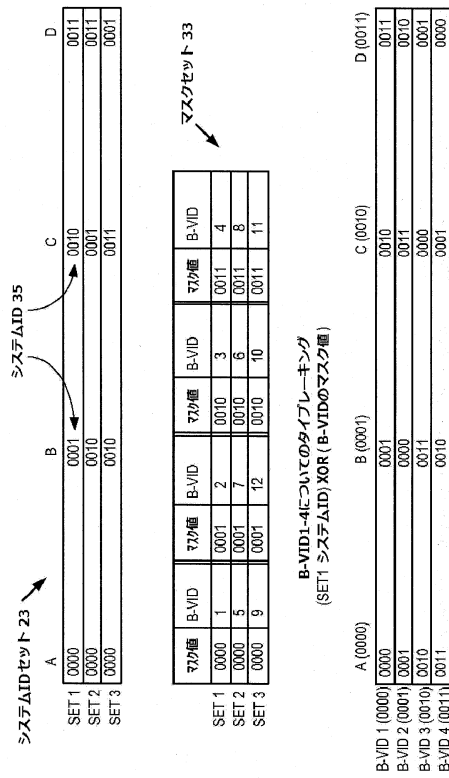
【図1】



【図2】



【図 3】



【図 4】

FIG. 3

B-VID 5-8 についてのタイプレッキング
(SET2 システムID XOR (B-VIDのマスク値))

B-VID 9-11 についてのタイプレッキング
(SET3 システムID XOR (B-VIDのマスク値))

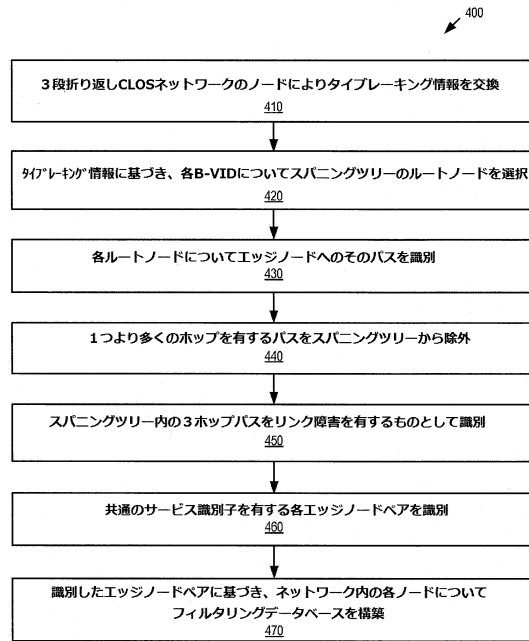


FIG. 4

【図 5】

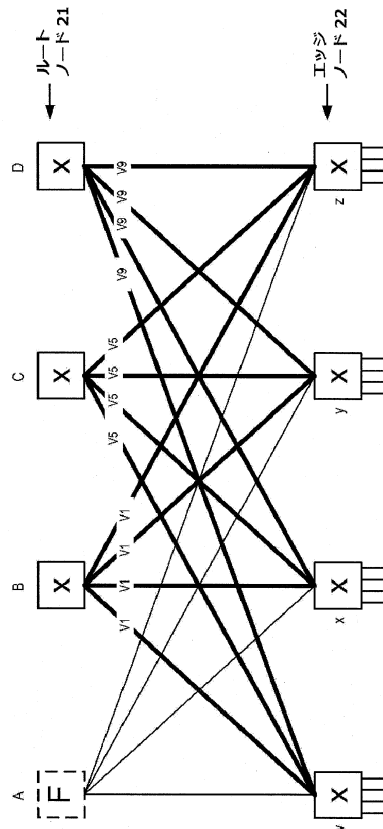


FIG. 5

【図 6】

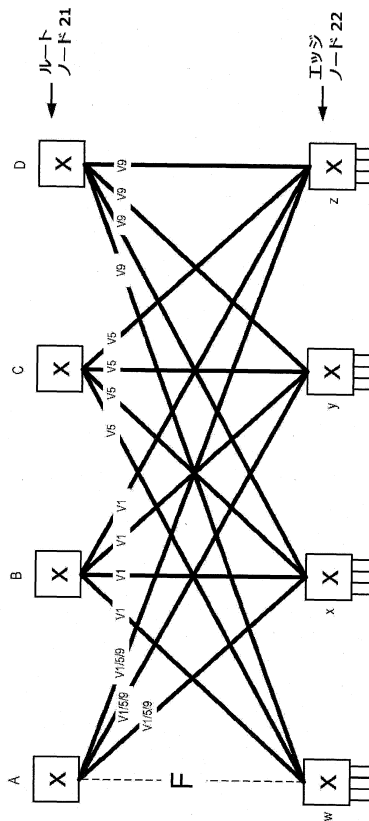


FIG. 6

【図 7】

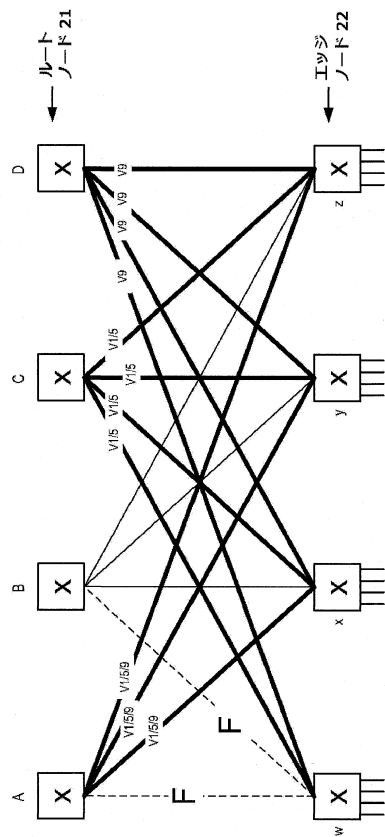


FIG. 7

【図 8】

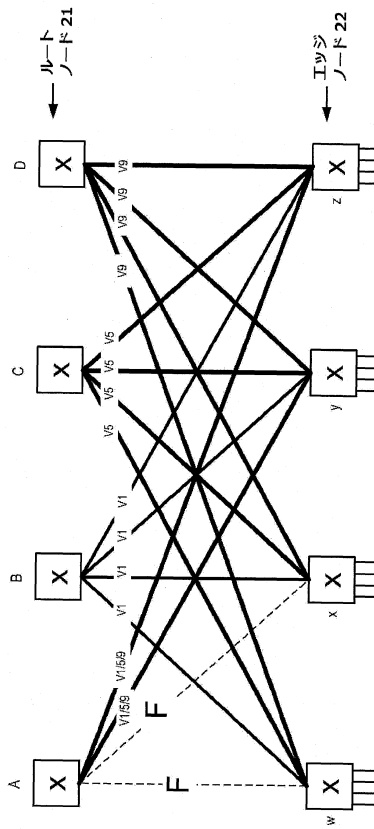


FIG. 8

【図 9】

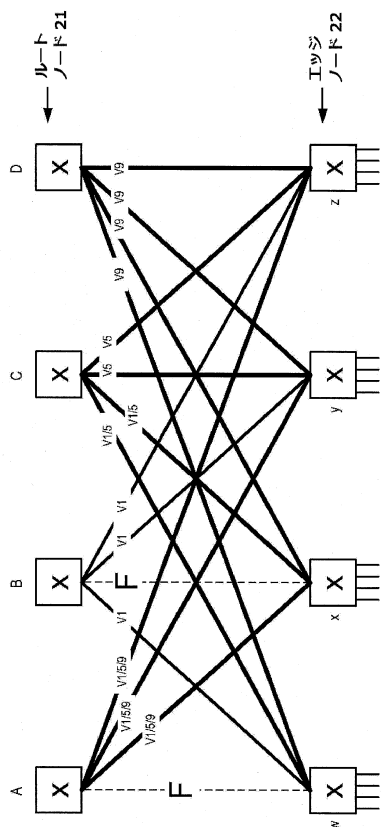


FIG. 9

【図 10】

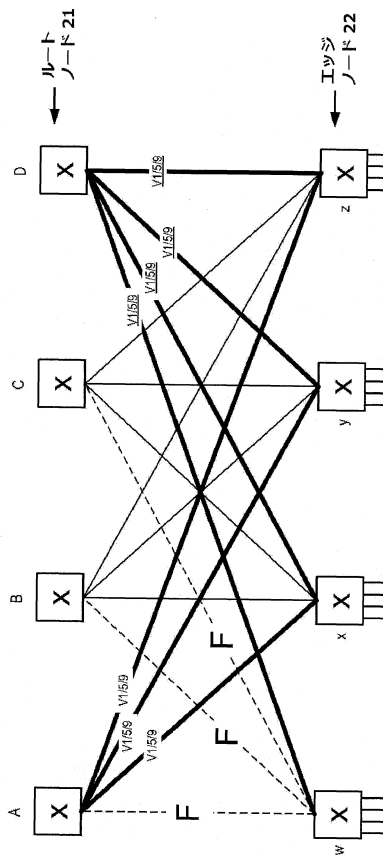


FIG. 10

【図 11】

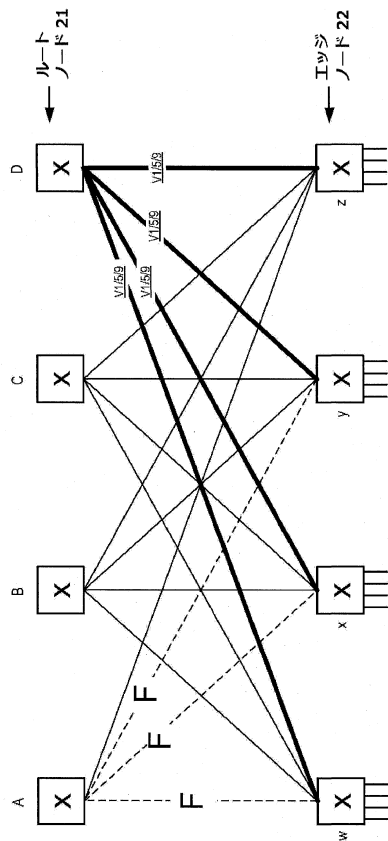


FIG. 11

【図 12】

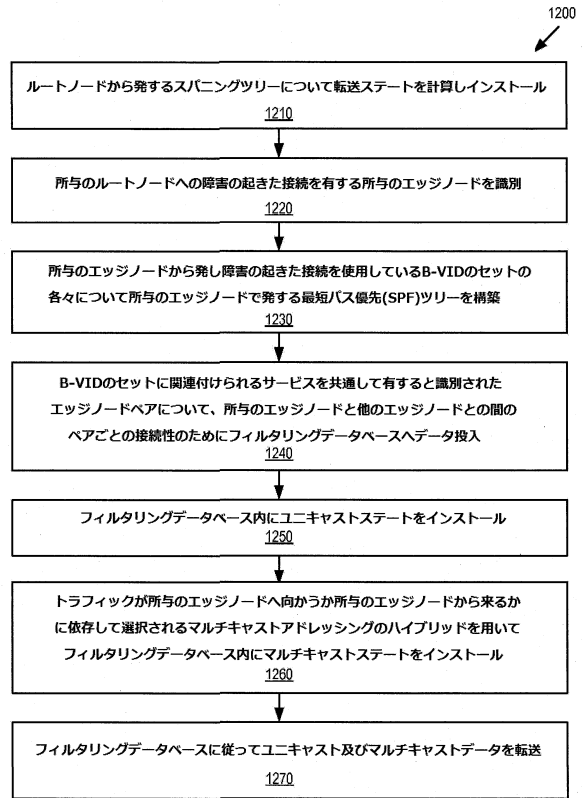


FIG. 12

【図 13】

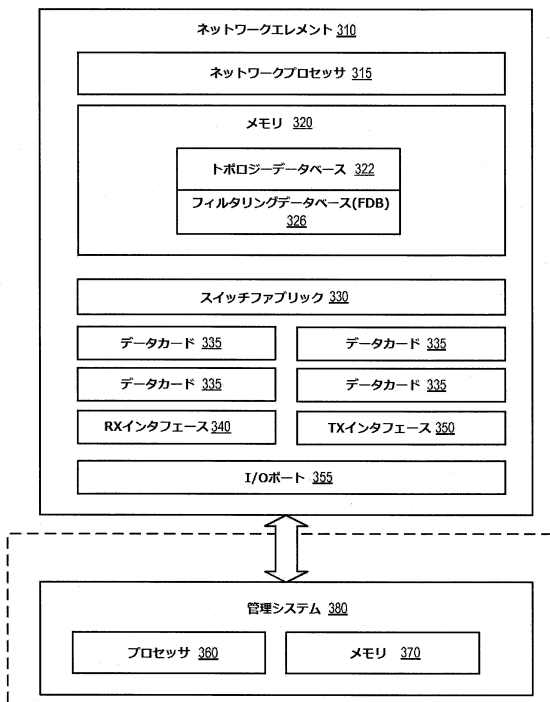


FIG. 13

フロントページの続き

(72)発明者 アラン、デイビッド イアン
アメリカ合衆国 カリフォルニア州 95112 サンノゼ 306 エス フィフティーンス
ストリート

審査官 速水 雄太

(56)参考文献 特表2008-546332(JP, A)
米国特許出願公開第2009/0213866(US, A1)
米国特許出願公開第2005/0111356(US, A1)

(58)調査した分野(Int.Cl., DB名)
H04L 12/733
H04L 12/741