

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2019-514294
(P2019-514294A)

(43) 公表日 令和1年5月30日(2019.5.30)

(51) Int.Cl.	F I	テーマコード (参考)
HO4L 12/803 (2013.01)	HO4L 12/803	5K030
HO4L 12/70 (2013.01)	HO4L 12/70	D

審査請求 未請求 予備審査請求 未請求 (全 20 頁)

(21) 出願番号 特願2018-553946 (P2018-553946)
 (86) (22) 出願日 平成29年4月12日 (2017. 4. 12)
 (85) 翻訳文提出日 平成30年12月12日 (2018. 12. 12)
 (86) 国際出願番号 PCT/US2017/027190
 (87) 国際公開番号 W02017/180731
 (87) 国際公開日 平成29年10月19日 (2017. 10. 19)
 (31) 優先権主張番号 62/321, 730
 (32) 優先日 平成28年4月12日 (2016. 4. 12)
 (33) 優先権主張国 米国 (US)
 (31) 優先権主張番号 15/485, 089
 (32) 優先日 平成29年4月11日 (2017. 4. 11)
 (33) 優先権主張国 米国 (US)

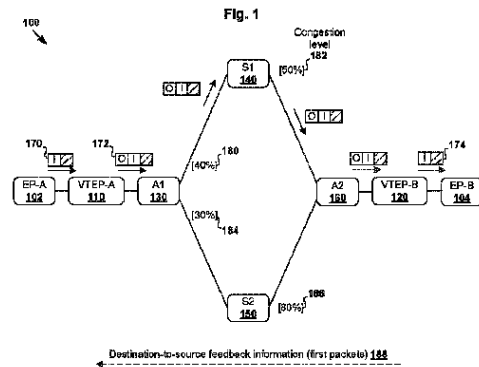
(71) 出願人 511235548
 ニシラ, インコーポレイテッド
 アメリカ合衆国 カリフォルニア州 94
 304, パロアルト, ヒルビュー ア
 ベニュー 3401
 (74) 代理人 100105957
 弁理士 恩田 誠
 (74) 代理人 100068755
 弁理士 恩田 博宣
 (74) 代理人 100142907
 弁理士 本田 淳
 (72) 発明者 ヒラ、ムケーシュ
 アメリカ合衆国 94304 カリフォル
 ニア州 パロ アルト ヒルビュー アベ
 ニュー 3401

最終頁に続く

(54) 【発明の名称】 輻輳を考慮したロードバランシングのための仮想トンネルエンドポイント

(57) 【要約】

例示的な方法は、送信元仮想トンネルエンドポイント (VTEP) のために提供されて、データセンタネットワークにおいて輻輳を考慮したロードバランシングを実行する。方法は、送信元VTEPが、送信元VTEPを送信先VTEPに接続する対応する複数の中間スイッチによって提供される複数の経路に関連付けられた輻輳状態情報を学習すること、を備える。また方法は、送信元VTEPが、送信元エンドポイントによって送信され且つ送信先VTEPに関連付けられた送信先エンドポイント宛ての複数の第2の packets を受信すること、輻輳状態情報に基づいて複数の経路から特定の経路を選択すること、を備える。さらに方法は、送信元VTEPが、特定の経路に関連付けられたタプルセットを含むヘッダ情報を有する前記複数の第2の packets の各々をカプセル化することによってカプセル化された複数の第2の packets を生成すること、カプセル化された複数の第2の packets を送信先エンドポイントに送信すること、を備える。



Legend:
 [] : Second packet with inner header information (i)
 [] : Second packet encapsulated with outer header information (o) that includes SPI

Congestion state information at VTEP-A 180

Outer source port number (source IP)	Path (path ID)	Congestion flag (congestion_flag)	Path weight (weight)
SP1	P1 (via S1)	false	w1
SP2	P2 (via S2)	true	w2

【特許請求の範囲】**【請求項 1】**

送信元仮想トンネルエンドポイント（VTEP）が、該送信元仮想トンネルエンドポイント、送信先仮想トンネルエンドポイント、送信元エンドポイント、送信先エンドポイント、複数の中間スイッチを含むデータセンタネットワークにおいて輻輳を考慮したロードバランシングを実行する方法であって、

前記送信先仮想トンネルエンドポイントからの複数の第 1 のパケットに基づいて、前記送信元仮想トンネルエンドポイントが、前記送信元仮想トンネルエンドポイントを前記送信先仮想トンネルエンドポイントに接続する対応する複数の中間スイッチによって提供される複数の経路に関連付けられた輻輳状態情報を学習すること、

前記送信元仮想トンネルエンドポイントが、前記送信元エンドポイントによって送信され且つ前記送信先仮想トンネルエンドポイントに関連付けられた送信先エンドポイント宛ての複数の第 2 のパケットを受信すること、

前記送信元仮想トンネルエンドポイントが、前記輻輳状態情報に基づいて前記複数の経路から特定の経路を選択すること、

前記送信元仮想トンネルエンドポイントが、前記特定の経路に関連付けられたタプルセットを含むヘッダ情報を有する前記複数の第 2 のパケットの各々をカプセル化することによってカプセル化された複数の第 2 のパケットを生成すること、

前記カプセル化された複数の第 2 のパケットが前記タプルセットに基づいた前記特定の経路を介して転送されるように、前記送信元仮想トンネルエンドポイントが、前記カプセル化された複数の第 2 のパケットを前記送信先エンドポイントに送信すること、を備える方法。

【請求項 2】

前記輻輳状態情報を学習することは、

前記送信先仮想トンネルエンドポイントからの特定の第 1 のパケットに基づいて前記特定の経路に関連付けられた輻輳フラグを決定することを含み、

前記特定の第 1 のパケットは、

前記特定の経路が輻輳しているかどうかを示すための前記複数の中間スイッチのうちの少なくとも 1 つからの輻輳通知を含む、請求項 1 に記載の方法。

【請求項 3】

前記輻輳状態情報を学習することは、

前記送信先仮想トンネルエンドポイントからの特定の第 1 のパケットの受信時刻、および前記特定の第 1 のパケットをトリガするための前記送信元仮想トンネルエンドポイントによって送信されたプローブパケットの送信時刻に基づいて、前記特定の経路に関連付けられたラウンドトリップタイムを決定することを含む、請求項 1 に記載の方法。

【請求項 4】

前記輻輳状態情報を学習することは、

前記送信先仮想トンネルエンドポイントからの特定の第 1 のパケットに基づいて、選択の可能性を増加または低下させるために前記特定の経路に関連付けられた重み付けを決定すること、を含む、請求項 1 に記載の方法。

【請求項 5】

前記輻輳状態情報を学習することは、

複数の外部送信元ポート番号と、対応する複数の経路との間のマッピングを学習することを含む、請求項 1 に記載の方法。

【請求項 6】

前記複数の第 2 のパケットの各々をカプセル化することは、

前記複数の外部送信元ポート番号から、前記特定の経路に関連付けられた特定の外部送信元ポート番号を決定すること、

前記特定の送信元ポート番号を含むため、前記複数の第 2 のパケットの各々において前記ヘッダ情報のタプルセットを構成すること、を含む、請求項 5 に記載の方法。

10

20

30

40

50

【請求項 7】

前記送信先エンドポイント宛ての前記現在の第 2 のパケットと以前の第 2 のパケットとの間のパケット間ギャップに基づいて、現在の第 2 のパケットである前記複数の第 2 のパケットのそれぞれをフローレットに割り当てること、

前記特定の経路に関連付けられた特定の外部送信元ポート番号に前記フローレットを関連付けるフローレット情報を記憶すること、をさらに備える請求項 6 に記載の方法。

【請求項 8】

1 組の命令を含む非一時的なコンピュータ可読記憶媒体であって、

前記 1 組の命令は、ホストのプロセッサによって実行されると、

前記プロセッサに送信元仮想トンネルエンドポイント (VTEP) を実行させて、前記送信元仮想トンネルエンドポイント、送信先仮想トンネルエンドポイント、送信元エンドポイント、送信先エンドポイント、複数の中間スイッチを含むデータセンタネットワークにおいて輻輳を考慮したロードバランシングの方法を実行し、前記ホストは、前記送信元仮想トンネルエンドポイントをサポートし、

前記方法は、

前記送信先仮想トンネルエンドポイントからの複数の第 1 のパケットに基づいて、前記送信元仮想トンネルエンドポイントが、前記送信元仮想トンネルエンドポイントを前記送信先仮想トンネルエンドポイントに接続する対応する複数の中間スイッチによって提供される複数の経路に関連付けられた輻輳状態情報を学習すること、

前記送信元仮想トンネルエンドポイントが、前記送信元エンドポイントによって送信され且つ前記送信先仮想トンネルエンドポイントに関連付けられた送信先エンドポイント宛ての複数の第 2 のパケットを受信すること、

前記送信元仮想トンネルエンドポイントが、前記輻輳状態情報に基づいて前記複数の経路から特定の経路を選択すること、

前記送信元仮想トンネルエンドポイントが、前記特定の経路に関連付けられたタプルセットを含むヘッダ情報を有する前記複数の第 2 のパケットの各々をカプセル化することによってカプセル化された複数の第 2 のパケットを生成すること、

前記カプセル化された複数の第 2 のパケットが前記タプルセットに基づいた前記特定の経路を介して転送されるように、前記送信元仮想トンネルエンドポイントが、前記カプセル化された複数の第 2 のパケットを前記送信先エンドポイントに送信すること、を備える、非一時的なコンピュータ可読記憶媒体。

【請求項 9】

前記輻輳状態情報を学習することは、

前記送信先仮想トンネルエンドポイントからの特定の第 1 のパケットに基づいて前記特定の経路に関連付けられた輻輳フラグを決定することを含み、

前記特定の第 1 のパケットは、

前記特定の経路が輻輳しているかどうかを示すための前記複数の中間スイッチのうちの少なくとも 1 つからの輻輳通知を含む、請求項 8 に記載の非一時的なコンピュータ可読記憶媒体。

【請求項 10】

前記輻輳状態情報を学習することは、

前記送信先仮想トンネルエンドポイントからの特定の第 1 のパケットの受信時刻、および前記特定の第 1 のパケットをトリガするための前記送信元仮想トンネルエンドポイントによって送信されたプローブパケットの送信時刻に基づいて、前記特定の経路に関連付けられたラウンドトリップタイムを決定することを含む、請求項 8 に記載の非一時的なコンピュータ可読記憶媒体。

【請求項 11】

前記輻輳状態情報を学習することは、

前記送信先仮想トンネルエンドポイントからの特定の第 1 のパケットに基づいて、選択の可能性を増加または低下させるために前記特定の経路に関連付けられた重み付けを決定

10

20

30

40

50

することを含み、請求項 8 に記載の非一時的なコンピュータ可読記憶媒体。

【請求項 1 2】

前記輻輳状態情報を学習することは、

複数の外部送信元ポート番号と、対応する複数の経路との間のマッピングを学習することを含む、請求項 8 に記載の非一時的なコンピュータ可読記憶媒体。

【請求項 1 3】

前記複数の第 2 のパケットの各々をカプセル化することは、

前記複数の外部送信元ポート番号から、前記特定の経路に関連付けられた特定の外部送信元ポート番号を決定すること、

前記特定の送信元ポート番号を含むため、前記複数の第 2 のパケットの各々において前記ヘッダ情報のタプルセットを構成すること、を含む、請求項 1 2 に記載の非一時的なコンピュータ可読記憶媒体。

10

【請求項 1 4】

前記送信先エンドポイント宛ての前記現在の第 2 のパケットと以前の第 2 のパケットとの間のパケット間ギャップに基づいて、現在の第 2 のパケットである前記複数の第 2 のパケットのそれぞれをフローレットに割り当てること、

前記特定の経路に関連付けられた特定の外部送信元ポート番号に前記フローレットを関連付けるフローレット情報を記憶すること、をさらに備える請求項 1 3 に記載の非一時的なコンピュータ可読記憶媒体。

【請求項 1 5】

20

送信元仮想トンネルエンドポイント (VTEP) を実装して、データセンタネットワークにおいて輻輳を考慮したロードバランシングを実行するように構成されたスイッチであって、

送信先仮想トンネルエンドポイントから複数の第 1 のパケットを受信するための 1 つまたは複数の第 1 のポートであって、前記送信元仮想トンネルエンドポイントを前記送信先仮想トンネルエンドポイントに接続する対応する複数の中間スイッチによって提供される複数の経路を介して前記送信元仮想トンネルエンドポイントに接続される前記 1 つまたは複数の第 1 のポートと、

送信元エンドポイントによって送信され且つ前記送信先仮想トンネルエンドポイントに関連付けられた送信先エンドポイント宛ての複数の第 2 のパケットを受信するための第 2 のポートと、

30

スイッチロジックであって、

前記複数の第 2 のパケットを受信することに対応して、前記複数の第 1 のパケットに基づいて前記複数の経路に関連付けられた輻輳状態情報を学習すること、

前記輻輳状態情報に基づいて前記複数の経路から特定の経路を選択すること、

前記特定の経路に関連付けられたタプルセットを含むヘッダ情報を有する前記複数の第 2 のパケットの各々をカプセル化することによってカプセル化された複数の第 2 のパケットを生成すること、

前記カプセル化された複数の第 2 のパケットが前記タプルセットに基づいた前記特定の経路を介して転送されるように、前記カプセル化された複数の第 2 のパケットを前記送信先エンドポイントに送信すること、を実行するように構成された前記スイッチロジックと、を備えるスイッチ。

40

【請求項 1 6】

前記スイッチロジックは、

前記送信先仮想トンネルエンドポイントからの特定の第 1 のパケットに基づいて前記特定の経路に関連付けられた輻輳フラグを決定することによって、前記輻輳状態情報を学習するように構成され、

前記特定の第 1 のパケットは、

前記特定の経路が輻輳しているかどうかを示すための前記複数の中間スイッチのうちの少なくとも 1 つからの輻輳通知を含む、請求項 1 1 に記載のスイッチ。

50

【請求項 17】

前記スイッチロジックは、

前記送信先仮想トンネルエンドポイントからの特定の第1のケットの受信時刻、および前記特定の第1のケットをトリガするための前記送信元仮想トンネルエンドポイントによって送信されたプローブケットの送信時刻に基づいて、前記特定の経路に関連付けられたラウンドトリップタイムを決定することによって、前記輻輳状態情報を学習するように構成される、請求項11に記載のスイッチ。

【請求項 18】

前記スイッチロジックは、

前記送信先仮想トンネルエンドポイントからの特定の第1のケットに基づいて、選択の可能性を増加または低下させるために前記特定の経路に関連付けられた重み付けを決定することによって、前記輻輳状態情報を学習するように構成される、請求項11に記載のスイッチ。

10

【請求項 19】

前記スイッチロジックは、

複数の外部送信元ポート番号と、対応する複数の経路との間のマッピングを学習することによって前記輻輳状態情報を学習するように構成される、請求項11に記載のスイッチ。

【請求項 20】

前記スイッチロジックは、

前記複数の外部送信元ポート番号から、前記特定の経路に関連付けられた特定の外部送信元ポート番号を決定すること、および

20

前記特定の送信元ポート番号を含むため、前記複数の第2のケットの各々において前記ヘッダ情報のタプルセットを構成することによって、前記輻輳状態情報を学習するように構成される、請求項19に記載のスイッチ。

【請求項 21】

前記スイッチロジックは、

前記現在の第2のケットと前記送信先エンドポイント宛ての以前の第2のケットとの間のケット間ギャップに基づいて、現在の第2のケットである前記第2のケットのそれぞれをフローレットに割り当てること、

30

前記特定の経路に関連付けられた特定の外部送信元ポート番号に前記フローレットを関連付けるフローレット情報を記憶すること、をさらに実行するように構成される、請求項20に記載のスイッチ。

【発明の詳細な説明】**【背景技術】****【0001】**

本明細書で別段の説明がない限り、このセクションに記載されたアプローチは、このセクションに含めることによって先行技術であると認められるものではない。

複数のデータセンタネットワークは、概して、多数のマルチ経路によって特徴付けられるマルチルートポロジを採用する。例えば、複数の物理サーバは、ケット転送の代替経路を提供する多数のスイッチを用いて相互に接続される。物理サーバが別の物理サーバに送信するデータを有する場合、複数の経路のうちの1つが選択されて、データが複数のケットのフローとして送信される。実際には、トラフィックは、異なる経路間で均等に分配されないことがあり、1つの経路の過剰利用と別の経路の利用不足の原因となる。ロードバランシングは、トラフィックをできるだけ均等に分散させて輻輳を減らし且つネットワークパフォーマンスを向上させるために重要である。

40

【図面の簡単な説明】**【0002】**

【図1】輻輳を考慮したロードバランシングが実行されるデータセンタネットワークの一例を示す概略図。

50

【図 2】送信元仮想トンネルエンドポイント（VTEP）がデータセンタネットワーク内で輻輳を考慮したロードバランシングを実行するための例示的なプロセスのフローチャート。

【図 3】送信元 VTEP がデータセンタネットワーク内の輻輳状態情報を学習するための第 1 の例示的なプロセスのフローチャート。

【図 4 A】送信元 VTEP がデータセンタネットワーク内の輻輳状態情報を学習するための第 2 の例示的なプロセスのフローチャート。

【図 4 B】図 4 A の例に従って学習された輻輳状態情報の一例を示す概略図。

【図 5】送信元 VTEP がデータセンタネットワークにおいてデータパケット処理を実行するための例示的なプロセスのフローチャート。

10

【発明を実施するための形態】

【0003】

以下の詳細な説明では、本明細書の一部を形成する添付の図面を参照する。図面において、類似の記号は、説明が別途指示しない限り、典型的には同様の構成要素を特定する。詳細な説明、図面、および特許請求の範囲に記載された例示的な実施形態は、限定を意味するものではない。本明細書に提示される主題の趣旨または範囲から逸脱することなく、他の実施形態が利用されてもよく、他の変更が行われてもよい。本明細書に全体的に記載され、図面に示された本開示の態様は、多種多様な異なる構成で配置され、置換され、結合され、構成され、これらのすべてが本明細書において明示的に企図されることは容易に理解されるだろう。

20

【0004】

複数のデータセンタネットワークにおけるロードバランシング（load balancing）の課題は、輻輳を考慮したロードバランシングが実行される例示的なデータセンタネットワーク 100 を示す概略図である図 1 を参照してより詳細に説明される。データセンタネットワーク 100 は、所望の実施形態に従って、示されるコンポーネントに対する追加および/または代替のコンポーネントを含むことができることを理解されたい。

【0005】

図 1 の例では、データセンタネットワーク 100 は、「VTEP - A」110 および「VTEP - B」120 等の仮想トンネルエンドポイント（virtual tunnel endpoints : VTEPs）によって提供される複数の経路を介して接続される第 1 のエンドポイント 102（「EP - A」参照）および第 2 のエンドポイント 104（「EP - B」参照）を含む。「VTEP - A」110 は、「A1」130、「S1」140、「S2」150 および「A2」160 などの複数の中間スイッチによって提供される複数の経路を介して「VTEP - B」120 に接続される。送信元「EP - A」102 から送信先「EP - B」104 に複数のデータパケットを転送する場合、複数のデータパケットは、「A1」130、「S1」140 および「A2」160 を経由する第 1 の経路、および「A1」130、「S2」150、および「A2」160 を経由する第 2 の経路のうちの一つで転送され得る。

30

【0006】

実際には、用語「仮想トンネルエンドポイント」（例えば、「VTEP - A」110 および「VTEP - B」120）は、複数のパケット転送サービス、複数のロードバランシングサービス、複数のゲートウェイサービス等を複数のエンドポイント（例えば、「EP - A」102 および「EP - B」104）に提供するように構成される任意の好適な複数のネットワーク要素に言及する。VTEP 110 / 120 は、1 つまたは複数の物理エンティティまたは仮想エンティティによって具体化されてもよい。例えば、VTEP 110 / 120 は、物理コンピューティングデバイス（例えば、エッジデバイス、物理サーバなど）によってサポートされる（supported）ハイパーバイザ（例えば、ハイパーバイザの仮想スイッチ）によって具体化されてもよい。VTEP 110 / 120 およびそれに関連するエンドポイント 102 / 104 は、同じ物理コンピューティングデバイス上に、または異なるコンピューティングデバイス上に存在してもよい。例えば、「EP - A」102 は仮想マシンであり、「VTEP - A」110 は同じ物理サーバによってサポートされる

40

50

仮想スイッチである。別の例では、「EP-A」102は、第1の物理サーバによってサポートされる仮想マシンであり、「VTEP-A」110は、第1の物理サーバに接続される第2の物理サーバまたは物理的トップ・オブ・ラック（top-of-rack：ToR）スイッチによってサポートされる。

【0007】

用語「エンドポイント」（例えば、「EP-A」102および「EP-B」104）は、一般的に、双方向性プロセス間通信フローの送信元ノード（「送信元エンドポイント」）または終端ノード（「送信先エンドポイント」）に言及する。実際には、エンドポイントは、物理コンピューティングデバイス（例えば、物理サーバ、物理ホスト）、物理コンピューティングデバイスによってサポートされる仮想化コンピューティングインスタンスなどであってもよい。仮想化されたコンピューティングインスタンスは、ワークロード、仮想マシン、アドレス指定可能なデータコンピューティングノード、隔離されたユーザ空間インスタンスなどを表してもよい。実際には、任意の適切な技術が、ハードウェア仮想化を含むが限定されない隔離された複数のユーザ空間インスタンスを提供するように用いられる。他の仮想化されたコンピューティングインスタンスは、（例えば、ハイパーバイザまたはDockerなどの別個のオペレーティングシステムを必要とせずにホストオペレーティングシステムの上で動作するか、またはオペレーティングシステムレベルの仮想化として具体される）複数のコンテナ（container）、複数の仮想プライベートサーバなどを含み得る。複数の仮想マシンは、物理コンピューティングシステムのハードウェアおよびソフトウェアコンポーネントの仮想等化物を含む完全な計算環境であってもよい。用語「ハイパーバイザ」は、一般的に、Dockerなどの複数の名前空間コンテナをサポートするシステムレベルのソフトウェアを含む複数の仮想化されたコンピューティングインスタンスの実行をサポートするソフトウェア層またはコンポーネントに言及する。

10

20

30

40

50

【0008】

用語「スイッチ」（例えば、「A1」130、「S1」140、「S2」150および「A2」160）は、概して複数のパケットを受信および転送するように構成された任意の適切なネットワーク要素に言及してもよく、レイヤ3ルータ、レイヤ2スイッチ、ゲートウェイ、ブリッジなどでもよい。ネットワークトポロジに応じて、スイッチはToRスイッチ、アグリゲイトスイッチ（aggregate switch）、スパインスイッチ（spine switch）などでもよい。簡略化のために、図1には2つの代替的な経路が示されているが、経路の数は、相互接続されたスイッチの数およびマルチルートトポロジ（例えば、リーフスパイントポロジ（leaf-spine topology）、ファットツリートポロジ（fat-tree topology）など）等のデータセンタネットワーク100に依存する。さらに、「VTEP-A」110と「VTEP-B」120を接続する追加のスイッチが図1に示すものよりも存在してもよい。

【0009】

用語「レイヤ2（layer-2）」は、一般的に、データリンクレイヤ（例えば、メディアアクセス制御（Media Access Control：MAC）またはイーサネットレイヤ）に言及し、用語「レイヤ3」は、ネットワークレイヤ（例えば、インターネットプロトコル（Internet Protocol：IP）レイヤ）に言及し、「レイヤ4」は、オープンシステム相互接続（Open System Interconnection：OSI）モデルにおけるトランスポート層（例えば、伝送制御プロトコル（Transmission Control Protocol：TCP）層）に言及するが、本明細書で説明される概念は、他のネットワークモデルにも適用可能であり得る。用語「パケット」は、一般的に、一緒に転送される1群のビットに言及し、「フレーム」、「メッセージ」、「セグメント」などの別の形態であってもよい。

【0010】

「VTEP-A」110と「VTEP-B」120との間の接続を提供するために、「トンネル」（簡略化のために図示せず）が、任意の適切なプロトコル（例えば、GENEVE（Generic Network Virtualization Encapsulation）、STT（Stateless Transport Tunneling）またはVXLAN（Virtual eXtension Local Area Network））を用いて

複数のVTEP間で確立される。用語「トンネル」は、一般的に、一对のVTEP間のエンドツーエンド且つ双方向通信の経路に言及する。この場合、「EP-A」102から複数のデータパケット(図1の170参照)を転送する前に、「VTEP-A」110はカプセル化を実行してカプセル化された複数のパケットを生成する(図1の172参照)。

【0011】

より詳細には、各データパケット170は、(図1では「I」の符号が付されている)「内部ヘッダ情報(inner header information)」およびペイロードとしてのアプリケーションデータを含む。カプセル化された後、カプセル化された各パケット172は、(図1では「O」の符号が付されている)外部ヘッダ情報(outer header information)およびペイロードとしてのデータパケット170を含む。(「外部トンネルヘッダ」としても知られている)「外部ヘッダ情報」は、外部レイヤ2ヘッダ、外部レイヤ3ヘッダ、外部レイヤ4ヘッダ等を含み得る。カプセル化が実行されて、(例えば、130~160によって形成される)ファブリック・オーバーレイ(fabric overlay)が外部トンネルヘッダに基づいて一对のVTEP間のパケット転送を実行するだけでよい。

10

【0012】

実際には、トラフィック負荷は、データセンタネットワーク100内の異なる複数の経路間で不均一に広がり、輻輳および性能低下を引き起こす可能性がある。従来、(例えば、ホップ数が等しい)等価コストで複数の経路にわたって均一にトラフィックを分散させるために、データプレーン・ロードバランシングメカニズム(data plane load balancing mechanism)として、等価コストマルチパスルーティング(equal cost multipath routing: ECMP)が一般的に用いられている。ECMPスイッチは、単純なハッシュベースのロードバランシングスキーム(hash-based load balancing scheme)を用いて、新しい各トラフィックのフローを利用可能な複数の経路のうちの一つにランダムに割り当てる。ECMPは、通常、ロードバランシングスキームを更新するために、柔軟性に欠けるカスタムシリコン(例えば、特定用途向け集積回路(ASIC))に実装される。さらに、ECMPは輻輳に依存せず、性能低下の原因となる複数の経路のオーバーサブスクリプション(oversubscription)を防止しない。

20

【0013】

例えば、図1では、「EP-A」102から「EP-B」104へ転送する複数のパケットのキュー占有レベル(queue occupancy levels)(複数の括弧内の180~186を参照)を用いて示されるように、異なる対のスイッチを接続する複数のリンクは異なる輻輳レベルを有する。「S1」140を介した第1の経路に沿ったキュー占有レベルは40%(180参照)および50%(182参照)である。「S2」150を介した第2の経路に沿ったキュー占有レベルは30%(184参照)および80%(186参照)である。ECMPは異なる輻輳レベルを考慮しないため、長時間実行されるフローは、80%(186参照)のキュー占有レベルを有する輻輳が発生する「S2」150を介した第2の経路に割り当てられる可能性がある。

30

【0014】

従来、制御プレーン・ロードバランシングメカニズムが、ECMPの欠点に対処するために用いられている。この場合、中央制御装置がデータセンタネットワーク100に配置されて、「A1」130、「S1」140、「S2」150および「A2」160から統計を収集し且つ転送ルールを「A1」130、「S1」140、「S2」150および「A2」160にプッシュして(push)、制御プレーン・ロードバランシングを実現する。しかしながら、中央制御装置が必要とされるので、制御プレーンメカニズムは、制御ループの待ち時間(latency)が長いために比較的遅く且つ高い揮発性のトラフィック(volatile traffic)を扱うことができない。

40

【0015】

従来、ホストベース(host-based)のアプローチが、ECMPの欠点に対処するために用いられている。例えば、マルチパスTCP(multipath TCP: MPTCP)と呼称される修正されたバージョンの伝送制御プロトコル(TCP)が、複数のエンドポイント間に複数

50

のサブフローを確立して、異なる複数の経路でトラフィックを分割するように用いられる。しかしながら、ホストベースのアプローチは、通常、MPTCPの場合に「EP-A」102および「EP-B」104のTCP/IPスタックを変更するなど、すべてのエンドポイントへの変更を必要とする。そのような変更は、特に、「EP-A」102および「EP-B」104が異なるオペレーティングシステムを実行している場合、または異なるエンティティによって制御されている場合には、特に困難である（場合によっては不可能である）。

【0016】

輻輳を考慮したロードバランシング (Congestion-aware load balancing)

本開示の例によれば、輻輳を考慮したロードバランシングのアプローチは、関連する「EP-A」102に気付かれない方法で「VTEP-A」110によって実施され得る。上述した従来のアプローチとは異なり、本開示の例は、MPTCPを実施するための「EP-A」102の変更、または新しいロードバランシングスキームを実施するための中間スイッチ130~160の変更を必要とせずにより具体化され得る。さらに、制御プレーンロードバランシングメカニズムとは異なり、中央制御装置を設けて輻輳監視を実行し且つ中間スイッチ130~160に転送ルールをプッシュする必要がない。

10

【0017】

より詳細には、図2は、送信元VTEP110がデータセンタネットワーク100内で輻輳を考慮したロードバランシングを実行するための例示的なプロセス200のフローチャートである。例示的なプロセス200は、205~240などの1つまたは複数のブロックによって例示される1つ以上の工程、機能、動作を含むことができる。様々なブロックは、より少数のブロックと組み合わせられ、追加のブロックに分割され、および/または所望の実施形態に応じて除去され得る。

20

【0018】

以下、「VTEP-A」110が送信元VTEPの例として用いられ、「VTEP-B」120が送信先VTEPの例として用いられ、「S1」140および「S2」150が中間スイッチの例として用いられ、「EP-A」102が「送信元エンドポイント」として用いられ、「EP-B」104が「送信先エンドポイント」として用いられる。キュー占有レベルは図1のデータセンタネットワーク100における輻輳を示す一例として用いるが、リンク利用レベル、ラウンドトリップタイム (round trip time : RTT) などの輻輳の他の適切なインジケータが用いられてもよい。

30

【0019】

図2の205において、「VTEP-A」110は、「VTEP-A」110と「VTEP-B」120とを接続する中間スイッチ130~160によって提供される複数の経路に関連する輻輳状態情報 (図1の190参照) を学習する (learn)。輻輳状態情報は、「VTEP-B」120からの送信先から送信元への (destination-to-source) フィードバック情報 (図1の188参照) を示す第1の packets に基づいて学習され得る。図2の210において、「VTEP-A」110は、「EP-A」102によって送信され且つ「EP-B」104宛ての複数の第2の (データ) packets 170を受信する。例えば、複数の第2の packets 170は、「EP-A」102上で動作するアプリケーションから「EP-B」104上で動作する別のアプリケーションまでのアプリケーションデータを含み得る。各第2の packets 170は、概して、「EP-A」102と「EP-B」104との間のプロセス間通信に関連する (図1では、「I」の符号が付された) 内部ヘッダ情報を含む。

40

【0020】

図2の220において、「VTEP-A」110は、複数の経路から特定の経路 (「選択経路」とも呼称される) を選択する。例えば、図1では、「VTEP-A」110は、複数の経路に関連する輻輳状態情報に基づいて、「A1」130、「S1」140および「A2」160を介した第1の経路を選択し得る。図2の230において、「VTEP-A」110は、220で選択された経路に関連するタプルセット (a set of tuples) を

50

含む（外部）ヘッダ情報を有する各第2の packets 170 をカプセル化することによってカプセル化された複数の第2の packets 172 を生成する。図2の240において、「VTEP-A」110は、カプセル化された第2の packets 172 を送信先「EP-B」104 に送信して、カプセル化された第2の packets 172 は、タプルセットに基づいて選択された経路を介して転送される。

【0021】

図3を用いてさらに説明されるように、「VTEP-A」110は、異なる外部送信元ポート番号（source_PN 192 を参照）を対応する経路（path_ID 194）および輻輳を示すフラグ（congestion_flag 196 を参照）に関連付ける輻輳状態情報（図1の190参照）に依存し得る。例えば、packets 転送を実行する前に、「VTEP-A」110は、source_PN 192 と path_ID 194 との間のマッピングを学習するために、経路学習（path learning）を実行し得る。「S1」140を介した第1の経路については、source_PN = SP1 および path_ID = P1 である。「S2」150を介した第2の経路については、source_PN = SP2 および path_ID = P2 である。

10

【0022】

図3を用いてさらに説明されるように、「VTEP-A」110は、source_PN 192 と congestion_flag 196 の異なる対の間のマッピングを学習してもよい。一例では、「VTEP-A」110は、データセンタネットワーク100内の輻輳状態情報のエンドツーエンド通知を可能にする明示的輻輳通知（Explicit Congestion Notification : ECN）などの既存の機能中間スイッチ（existing capabilities intermediate switches）130 ~ 160 に依存し得る。この場合、packets をドロップする代わりに、中間スイッチ130 ~ 160は、特定の経路に関連する現在のまたは目下の輻輳を「VTEP-B」120に通知するために、輻輳通知のフォームとして packets マーキングを受信する。その後、「VTEP-B」120は、「VTEP-A」110（図1の188参照）に輻輳通知を報知する。ECNの他に、他の適切なアプローチが、輻輳状態情報を学習するために用いられてもよい。例えば、図4Aおよび図4Bを用いてさらに説明されるように、「VTEP-A」110は、「VTEP-A」110と「VTEP-B」120との間のRTTを測定することができる。

20

【0023】

図1の例では、カプセル化された各第2の packets 172 は、選択された経路に関連するタプルセットを含む外部ヘッダ情報を含む。具体的には、カプセル化された各第2の packets 172 は、外部レイヤ2ヘッダ、外部レイヤ3ヘッダ、および外部レイヤ4ヘッダなどの（「0」の符号が付けされている）外部ヘッダ情報を含む。タプルセットは、送信元ポート番号、送信先ポート番号、送信元IPアドレス、送信先IPアドレス、およびプロトコルを含み得る。「S1」140を介した第1の経路（すなわち、path_ID = P1）について、外部レイヤ4ヘッダは、source_PN = SP1 の値を有する送信元ポート番号を含み得る。送信元ポート番号は、選択された経路に沿って「A1」130、「S1」140および「A2」160を介してカプセル化された複数の packets 172 が転送されるように設定される。送信先「VTEP-B」120において、カプセル化の解除が実行されて外部ヘッダ情報を除去し、データ packets 174 が「EP-B」104 に送信される。

30

40

【0024】

例示的なプロセス200を用いて、「VTEP-A」110は、異なる複数の経路に関連付けられた輻輳状態情報190を考慮して、データセンタネットワーク100内の異なる複数の経路上に仮想ネットワークラフィックを分配することができる。「VTEP-A」110は、異なる経路上にカプセル化された第2の packets 172 を送信したいときはいつでも、異なる外部送信元ポート番号を選択することができる。「VTEP-A」110と「VTEP-B」120とを接続する中間スイッチ130 ~ 160は外部ヘッダ情報に基づいてロードバランシングを実行するので、外部レイヤ4ヘッダ内の外部送信元ポ

50

ート番号はエントロピーとして用いられ、データセンタネットワーク 100 内の複数の経路（例えば、等コストパス（equal-cost paths））を利用する。

【0025】

輻輳状態情報

図3は、第1の例によるデータセンタネットワーク100内の輻輳状態情報190を学習するための送信元VTEP110の例示的な第1のプロセス300のフローチャートである。例示的なプロセス300は、310～365などの1つまたは複数のブロックによって示される1つまたは複数の工程、機能、または動作を含み得る。様々なブロックは、より少数のブロックに組み合わせられ、追加のブロックに分割され、および/または所望の実施形態に応じて除去され得る。図1の例を用いて、例示的なプロセス300は、「VTEP-A」110によって実行され得る。

10

【0026】

図3の310において、「VTEP-A」110は、source_PN192およびpath_ID194の複数の対の間のマッピングまたは関連付けを学習する。例えば、図1において、source_PN192の異なる値は、データセンタネットワーク100内の異なる重なり合わない複数のECMP経路を導く可能性がある。マッピングは、「VTEP-A」110がパケット転送のために選択することができる異なる複数の経路の先験的知識を示す。ECMPハッシング（ECMP hashing）がsource_PN192の特定の値を含むタプルセットに適用される場合、結果として関連する経路がpath_ID194によって識別される。

20

【0027】

実際には、「VTEP-A」110は、各経路で検出されたすべてのインタフェースIPに関する「トレースルート（traceroute）」スタイル情報を収集するために、（例えば、Paris tracerouteの後にモデル化される）バックグラウンドデーモンを実施して、データセンタネットワーク100内の他のすべてのVTEPに定期的なプローブパケットを送信する。「VTEP-A」110は、各プローブパケットの外部ヘッダ情報内の外部送信元ポート番号を変更して（rotate）、各ポート番号に対する経路トレースを収集することができる。次に、「VTEP-A」110は、対応するプローブ経路トレースがこれまで収集されたトレースと異なるたびに、輻輳状態情報190を更新してsource_PN192を追加または更新する。

30

【0028】

図3のブロック315～360は、path_ID194によって識別される経路と、その経路に沿った輻輳を示すcongestion_flag196とを関連付ける輻輳状態情報190を学習する「VTEP-A」110に関する。特に、図3の315および320において、「VTEP-A」110は、source_PN192の特定の値を含む外部トンネルヘッダを有する複数のパケットをカプセル化し、カプセル化された複数のパケットを送信する。

【0029】

図3の例では、「VTEP-A」110は、現在のまたは目下の輻輳などを示すためのECNマーキングなどの既存のスイッチの輻輳通知機能に依存することができる。ECNの詳細情報は、コメント番号3168のインターネットエンジニアリングタスクフォース（Internet Engineering Task Force：IETF）で見つられ、「IPへの明示的輻輳通知（ECN）の追加」と題されており、これは参照によりその全体が本明細書に組み込まれる。例としてECNが説明されているが、任意の他の適切なパケットマーキングのアプローチを使用できることを理解されたい。

40

【0030】

図3の325、330、335及び340において、「VTEP-A」110からカプセル化された複数のパケットを受信することに対応して、スイッチ（例えば、「S1」140、「S2」150）は、ECNマーキングを実行して、カプセル化された複数のパケットを転送する前にこのスイッチにおける輻輳にフラグを立てる。この場合、スイッチは

50

、輻輳通知のフォームとしてカプセル化された複数のパケットのヘッダ情報（例えば、TCPヘッダの予約されたフィールド）を変更することができるECN対応スイッチ（ECN-enabled switch）として知られている。

【0031】

例えば、図1において、「S2」150におけるキュー占有レベル（例えば、80%）が所定の閾値（例えば、 $T_Q = 60\%$ ）を超えると、「S2」150は、「VTEP-A」110から受信された複数のパケットをマークキングして、「VTEP-B」120に輻輳を通知して、第2の経路が輻輳した経路であることを示す。一方、関連するキュー占有レベル（例えば、50%）が閾値未満であるため、「S1」140が任意のパケットマークキングを実行する必要はなく、それによって第1の経路が非輻輳の経路であることを示す。

10

【0032】

図3の345および350において、送信先「VTEP-B」120は、カプセル化されたパケットを受信し、外部ヘッダ情報内の外部送信元ポート番号と送信元「VTEP-A」110（図1の188も参照）への任意の輻輳通知（例えば、ECNマークキング）との間のマッピングを通知する。図3の355および360において、「VTEP-A」110は、source_PN192と、関連するcongestion_flag196（すなわち、フラグ情報）との間のマッピングを更新する。例えば、図1において、「VTEP-A」110は、「S1」140を介した第1のパスのsource_PN=SP1についてはcongestion_flag=falseであることを判定し、「S2」150を介した第2のパスのsource_PN=SP1についてはcongestion_flag=trueであると判定する。

20

【0033】

congestion_flagは、関連する経路の選択の可能性に影響を与える重み付けアルゴリズムにおいて使用されてもよい。ロードバランシングプロセスの開始時に、「VTEP-A」110は、図3の310で経路学習機能によって発見されたすべての等コスト経路について等しい重み付けを開始することができる。続いて、図3の365で、輻輳フラグに基づいて経路の重み付け（図1の重み付け198参照）が調整される。例えば、congestion_flag=false（すなわち、クリア）であるとの判定に回答して、「S1」140を介した第1の経路に関連するweight=w1は、選択の可能性を増大するために増大される。一方、congestion_flag=true（すなわち、セット）であるとの判定に回答して、「S2」150を介した第2の経路に関連するweight=w2は、選択の可能性を低減するために低減される。

30

【0034】

輻輳を示すために使用され得る別のメトリックは、各経路について測定され且つアクティブに追跡され得るRTTである。例は、送信元VTEP110がデータセンタネットワーク100内の輻輳状態情報を学習する第2の例示的なプロセス400のフローチャートである図4Aを用いて説明される。プロセス400の例は、410~455などのより多くのブロックにより例示される1つ以上の工程、機能、または動作を含み得る。様々なブロックは、より少数のブロックに組み合わせられ、複数の追加のブロックに分割され、および/または所望の実施形態に応じて除去され得る。

40

【0035】

図4Aの410、415および420において、「VTEP-A」110は、関連するスイッチ140/150を介した各経路上に複数のプローブパケットを周期的に送信し、各プローブパケットはsource_PNおよび送信(Tx)タイムスタンプを識別する。図4Aの425、430および435において、「VTEP-B」120は、スイッチ140/150を介してsource_PNおよびTxタイムスタンプも識別する確認応答(ACK)パケットを送信することによって、各プローブパケットに回答する。図4Aの440および445で、特定の受信時刻にACKパケットを受信したことに応じて、「VTEP-A」110は、受信時刻とTxタイムスタンプとの間の差に基づいてRTT

50

を決定することができる。図4Aの450および455において、「VTEP-A」110は、RTTをsource_PNに関連付け、経路に関連する重み付けを調整する。

【0036】

図4Bは、図4Aの例を用いて学習された例示的な輻輳状態情報460を示す概略図である。例えば、「S1」140を介した第1の経路は、source_PN=SP1(462参照)、path_ID=P1(464参照)、RTT=R1(466参照)およびweight=w1(468参照)に関連付けられる。「S2」150を介した第2の経路は、source_PN=SP2、path_ID=P2、RTT=R2およびweight=w2に関連付けられる。図1の例では、第2の経路が、第1の経路よりも高い輻輳レベルを有し、R2はR1よりも大きくすべきである。この場合、w2はw1よりも小さく調整されて、第2経路の選択の可能性が低減される。実際には、タイムスタンプと肯定応答(acknowledgement)が行われることが近いほど、トランスミッタとレシーバのソフトウェアスタックによって導入されるレイテンシを含まないため、RTTは経路の実際のネットワークレイテンシをより正確に反映する。

10

【0037】

データパケット処理

本開示の例によれば、ロードバランシングは、TCPなどのトランスポート層プロトコルに関連するパケット並べ替え問題を回避または改善するために、複数のフローレットの粒度(granularity of flowlets)で実行され得る。これは、複数のパケットのフローを「フローレット(flowlet)」と呼称される複数のより小さいグループに分割することによって達成され得る。本明細書で使用される用語「フローレット」は、概して、フロー内のパケットのグループまたはバースト(burst)に言及する。

20

【0038】

図5は、仮想トンネルエンドポイント110がデータセンタネットワーク100においてデータパケット処理を実行する例示的なプロセス500のフローチャートである。例示的なプロセス500は、510~560等の1つまたは複数のブロックによって示される1つまたは複数の工程、機能、動作を含み得る。種々のブロックは、より少数のブロックに組み合わせられ、追加のブロックに分割され、および/または所望の実施形態に応じて除去され得る。

【0039】

510および515において、「EP-A」102からのデータパケット170を受信したことに応答して、「VTEP-A」110は、データパケット170が新しいフローレットまたは現在のフローレットに属するかどうかを判定する。例えば、同じフロー内の2つの連続するパケットの到着の時間間隔(すなわち、パケット間ギャップ(inter-packet))が所定の閾値(例えば、T_flowlet秒; 515参照)を超えるたびに、新しいフローレットが検出され得る。しきい値を超えないすべての後続パケットは、同じフローレットの一部であるとみなされる。2つのシナリオが以下で説明される。

30

【0040】

(a) 新しいフローレットが検出されたとき(すなわち、パケット間ギャップ>T_flowletまたはフローの第1パケットが検出されたとき)、「VTEP-A」110は、図5の520においてデータパケット170にflowlet_ID(例えば、「F1」)を割り当てる。図5の525において、その時点の輻輳状態情報190/460に基づいて、新しいフローレット(525参照)に対して、path_ID(例えば、「P1」)によって識別される経路が選択される。図5の530において、選択された経路に関連するsource_PNが決定され、選択された経路に関連するflowlet_IDとsource_PN(例えば、「SP1」)との間の関連性が記憶される。「VTEP-A」110は、フローレットの最新のデータパケット170が受信された時刻を記録するためのflowlet_timeも記憶する。

40

【0041】

(b) 既存のフローレットのデータパケット170が検出されたとき(すなわち、パケ

50

ット間ギャップ = 現在の時刻 - flowlet__time (flowlet)、 「VTEP-A」110は、図5の550および555において、現在のflowlet__ID (たとえば「F1」) および関連するsource__PN (たとえば「SP1」) を取得する。同様に、図5の560において、「VTEP-A」110は、フローレットの最新のデータパケット170が受信された時刻を記録するためのflowlet__timeも記憶する。

【0042】

(a) および (b) の両方の場合、例示的なプロセス500は、図5の535および540に続き、「VTEP-A」110は、flowlet__ID (例えば、「F1」) に関連付けられたsource__PN (例えば「SP1」) の値を有する外部送信元ポート番号を含むように構成された外部トンネルヘッダを有するデータパケット170をカプセル化し、カプセル化されたデータパケット172を送信する。外部トンネルヘッダはまた、「VTEP-A」110に関連する送信元IPアドレスと、「VTEP-B」120に関連する送信先IPアドレスを含む。送信先では、パケットが「EP-B」104に転送される前に、外部トンネルヘッダが「VTEP-B」120によって除去される。

10

【0043】

図5の515でパケット間ギャップを超過したとき、flowlet__ID = 「F2」などの別のフローレットに属する後続のデータパケット170に対して上記事項を繰り返すことができることを理解されたい。この場合、輻輳状態情報190/460が変更され、「S1」140を介した第1の経路の代わりに「S2」150を介した第2の経路 (すなわち、path__ID = P2) が選択され得る。この場合、関連するsource__PN = SP2は、カプセル化されたデータパケット172が第2の経路に沿って「A1」130、「S2」150および「A2」160によって「VTEP-B」120に転送されるような外部ヘッダ情報に含まれる。

20

【0044】

本開示の例を用いて、「VTEP-A」110は、新しいフローレットのflowlet__IDを記憶し、同じフローレットの複数のパケットのために関連するsource__PNを再利用する。実際には、閾値T_flowletは、ネットワークにおいて (例えば、RTTのオーダーにおいて) 推定RTTに基づいて設定されてもよい。複数のフローレット間のパケット間ギャップが大きいことは、(技術的には同じフローの一部である) 複数のフローレットが異なる経路を用いる場合、パケットの並べ替えを最小限に抑えることを保証する。さらに、「VTEP-A」110は、選択された経路を用いるようにするために、新しい各フローレットに対して最小の輻輳を有する外部送信元ポート番号を割り当てることができる。これにより、すべてのアクティブな経路のキューが常に低く保たれるため、ワークロードのスループットおよびレイテンシが向上する。

30

【0045】

本開示の複数の例は、(例えば、「EP-A」102および「EP-B」104として動作する) 複数のゲスト仮想マシンに気付かれない方法でネットワークフローを複数の非輻輳の経路に分割するように「VTEP-A」110を構成することによって、エンドユーザの複数のゲスト仮想マシンに気付かれないように実施される。一例では、(例えば、「VTEP-A」110として動作する) 送信元仮想スイッチと、(例えば、「VTEP-B」120として動作する) 送信先仮想スイッチの両方が用いられ得る。この場合、フローは任意に複数のフローレットに分割され、複数のフローレット間のアイドルギャップ (idle gap) に左右されない。

40

【0046】

例えば、各TCPセグメンテーションオフロード (TCP segmentation offload : TSO) セグメントは、フローレットとして扱うことができる。用いられる複数のフローが異なるために送信先仮想スイッチで複数のフローレットが順番に関係なく (out of order) 到着した場合、送信先仮想スイッチは、それらを送信先エンドポイント (例えば、ゲスト仮想マシン) に配信する前に、複数のフローレットを並べ替えることができる。これにより、

50

送信先仮想スイッチは、送信先エンドポイントの送信先プロトコルスタック（例えば、TCP/IPスタック）から順番に関係なく到達したことを隠して、送信先プロトコルスタックが、送信元エンドポイントでの送信元プロトコルスタックの遅延を防止することができる。

【0047】

本開示の例は、ハイパーバイザの仮想スイッチで具体化され、トラフィックの第1のエントリーポイントからインテリジェントな経路選択に導かれる。たとえば、仮想スイッチ以降の複数のECMP経路が存在する可能性があり、各経路は、異なる物理NICを用いる。この場合、本開示の例は、プロセッサを含むホストと、一組の命令を記憶する非一時的なコンピュータ可読記憶媒体とを含むホストを用いて具体化され得る。プロセッサによる実行に応答して、1組の命令は、図1～図5の例にしたがって、プロセッサに、データセンタネットワーク100における輻輳認識ロードバランシングを実行するために、送信元VTEP（例えば、ホストによってサポートされる仮想スイッチにおける「VTEP-A」110）を具体化させる。

10

【0048】

一例では、ホストによって実装された送信元VTEPは、送信先VTEPからの複数の第1の packets に基づいて、送信元VTEPを送信先VTEPに接続する対応する複数の中間スイッチによって提供される複数の経路に関連する輻輳状態情報を学習することができる。また、送信元エンドポイントによって送信され、送信先VTEPに関連付けられた送信先エンドポイント宛ての複数の第2の packets を受信することに応答して、送信元VTEPは、輻輳状態情報に基づいて複数の経路から特定の経路を選択し、複数の第2の packets の各々を、前記特定の経路に関連する1組のタブルを含むヘッダ情報でカプセル化することによってカプセル化された複数の第2の packets を生成し、カプセル化された複数の第2の packets が1組のタブルに基づいて特定の経路を介して転送されるように、カプセル化された複数の第2の packets を送信先エンドポイントに送信する。

20

【0049】

すべての物理NICが同じレイヤ3に接続するシナリオでは、ネクストホップ（next-hop）および経路ダイバーシティ（path diversity）が第1ホップスイッチ（例えば、TORスイッチ）を超えて開始する場合、本開示の例は、仮想スイッチソフトウェアよりも早いスピードでNICドライバ/ハードウェアまたは第1ホップスイッチにおいて具体化され得る。高度なスイッチアーキテクチャを必要とする従来のアプローチと比較して、本開示の例は、エッジハイパーバイザ（edge hypervisor）において（例えば、完全にソフトウェアで）実行され、複数の送信元と複数の送信先との間の任意の数のホップにスケールリングされ得る。この場合、図1から図5の例に従ってデータセンタネットワーク100内で輻輳認識ロードバランシングを実行する送信元VTEP（例えば、「VTEP-A」110）を具体化するために、例示的な（第1ホップ）スイッチが用いられる。スイッチは、ハードウェアロジック（例えば、ハードウェア回路）、プログラマブルロジック、またはそれらの組み合わせなど、任意の適切なスイッチロジックを含み得る。

30

【0050】

一例では、スイッチは、1つまたは複数の第1のポート、1つまたは複数の第2のポート、およびスイッチロジックを含み得る。1つまたは複数の第1のポートは、送信元VTEPを送信先VTEPに接続する対応する複数の中間スイッチによって提供される複数の経路を介して送信元VTEPに接続された送信先VTEPから複数の第1 packets を受信するように用いられてもよい。第2のポートは、送信元エンドポイントによって送信され、送信先VTEPに関連付けられた送信先エンドポイント宛である第2の packets を受信するように用いられてもよい。スイッチロジックは、複数の第1の packets に基づいて複数の経路に関連付けられた輻輳状態情報を学習するように構成され得る。複数の第2の packets を受信したことに応答して、スイッチロジックは、輻輳状態情報に基づいて複数の経路から特定の経路を選択し、特定の経路に関連付けられたタブルセットを含むヘッダ情報を有する複数の第2の packets の各々をカプセル化することによってカプセル化された

40

50

複数の第2の packets を生成し、カプセル化された複数の第2の packets がタプルセットに基づいて特定の経路を介して転送されるように、カプセル化された複数の第2の packets を送信先エンドポイントに送信するように構成され得る。

【0051】

上記で紹介した技術は、専用ハードワイヤード回路、プログラマブル回路と組み合わせたソフトウェアおよび/またはファームウェア、またはそれらの組み合わせで具体化され得る。専用ハードワイヤード回路は、例えば、1つまたは複数の特定用途向け集積回路 (ASIC)、プログラマブルロジックデバイス (PLD)、フィールドプログラマブルゲートアレイ (FPGA)、プログラマブルスイッチアーキテクチャなどの形態であってもよい。用語「プロセッサ」は、処理ユニット、ASIC、論理ユニット、またはプログラマブルゲートアレイなどを含むように広く解釈されるべきである。

10

【0052】

前述の詳細な説明は、ブロック図、フローチャート、および/または実施例の使用を介して、デバイスおよび/またはプロセスの様々な実施形態を示している。そのようなブロック図、フローチャート、および/または実施例が1つまたは複数の機能および/または動作を含む限り、そのようなブロック図、フローチャート、または実施例における各機能および/または動作が、ハードウェア、ソフトウェア、ファームウェア、またはそれらの任意の組み合わせの広い範囲によって個別におよび/または集合的に具体化可能であることは、当業者によって理解されよう。

【0053】

当業者であれば、本明細書に開示された実施形態のいくつかの態様は、全体的または部分的に、集積回路内で、1つまたは複数のコンピュータ上で実行される1つまたは複数のコンピュータプログラムとして (例えば、1つまたは複数のコンピュータシステム上で実行される1つまたは複数のプログラムとして)、1つまたは複数のプロセッサ上で実行される1つまたは複数のプログラムとして (例えば、1つまたは複数のマイクロプロセッサ上で実行される1つまたは複数のプログラムとして)、ファームウェアとして、またはそれらの仮想的な任意の組み合わせとして、透過的に具体化され得ること、および回路の構成および/またはソフトウェアおよび/またはファームウェア用のコードの記述はこの開示に照らして当業者の技術範囲内であることを当業者によって認識されよう。

20

【0054】

明細書で紹介されたソフトウェアおよび/または技術を具体化することは、非一時的なコンピュータ可読記憶媒体に記憶されてもよく、1つまたは複数の汎用または専用のプログラム可能なマイクロプロセッサによって実行されてもよい。「コンピュータ可読記憶媒体」は、本明細書で使用されるように、機械 (例えば、コンピュータ、ネットワーク装置、パーソナルデジタルアシスタント (PDA)、モバイルデバイス、製造ツール、1つまたは複数のプロセッサのセットを有する任意のデバイスなど) を含み得る。コンピュータ可読記憶媒体は、記録可能/記録不可能媒体 (例えば、読み出し専用メモリ (ROM)、ランダムアクセスメモリ (RAM)、磁気ディスクまたは光記憶媒体、フラッシュメモリデバイスなど) を含み得る。

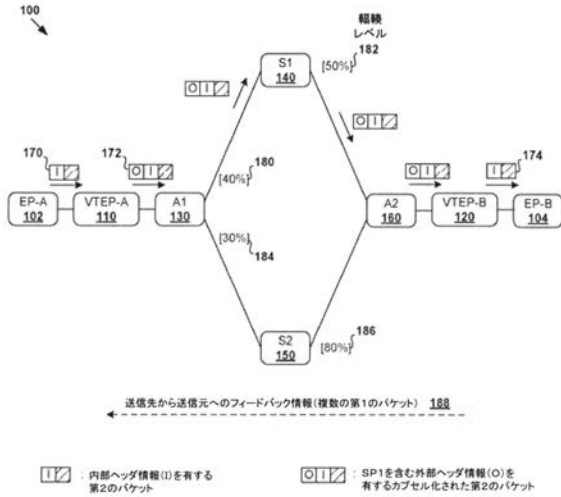
30

【0055】

図面は、例示にすぎず、図面に示される要素または手順は、本開示を実施するために必ずしも必須ではない。当業者であれば、実施例中の要素は、記載された実施例の装置に配置されるか、あるいは、実施例の装置とは異なる1つ以上の装置に代替的に配置されることを理解するであろう。記載された実施例における要素は、1つのモジュールに組み合わされるか、またはさらに複数のサブ要素に分割され得る。

40

【図1】



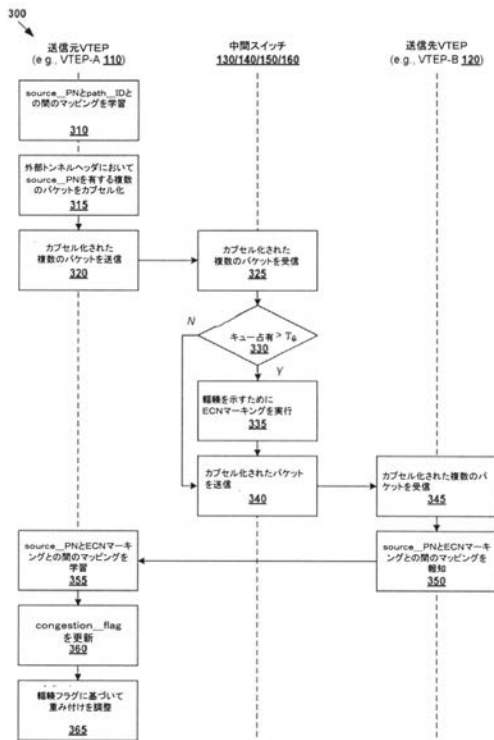
VTEP-Aにおける輻輳状態情報

外部送信元ポート番号 (source_PN)	経路 (path_ID)	輻輳フラグ (congestion_flag)	経路の重み付け (weight)
192	194	196	198
SP1	P1 (via S1)	false	w1
SP2	P2 (via S2)	true	w2

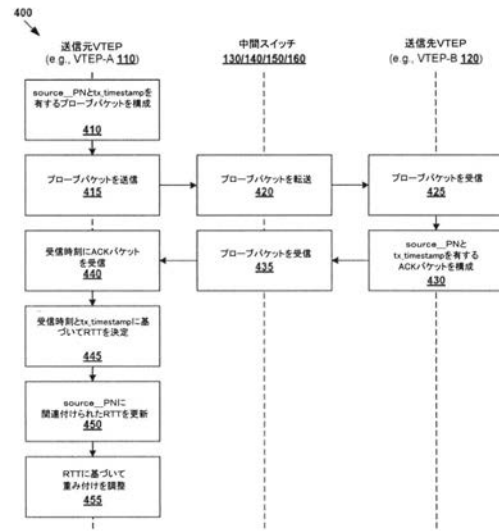
【図2】



【図3】



【図4A】

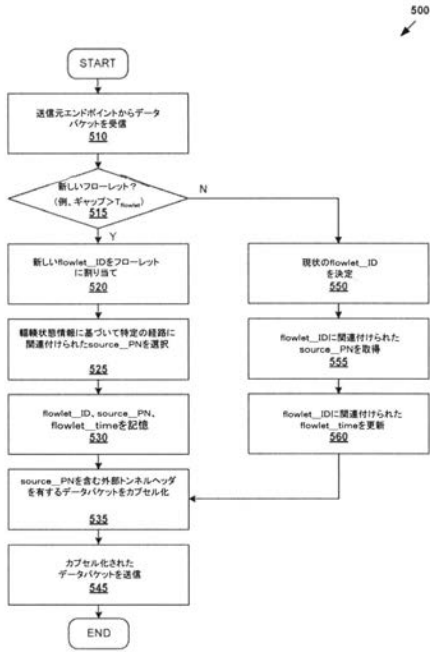


【図4B】

VTEP-Aにおける輻輳状態情報

外部送信元ポート番号 (source_PN)	経路 (path_ID)	ラウンドトリップタイム (RTT)	経路の重み付け (weight)
462	464	466	468
SP1	P1 (via S1)	R1	w1
SP2	P2 (via S2)	R2	w2

【図5】



【 国際調査報告 】

INTERNATIONAL SEARCH REPORT		International application No PCT/US2017/027190
A. CLASSIFICATION OF SUBJECT MATTER INV. H04L12/715 H04L12/729 H04L12/803 ADD.		
According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED Minimum documentation searched (classification system followed by classification symbols) H04L		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) EPO-Internal		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	MOHAMMAD ALIZADEH ET AL: "CONGA", SIGCOMM, ACM, 2 PENN PLAZA, SUITE 701 NEW YORK NY 10121-0701 USA, 17 August 2014 (2014-08-17), pages 503-514, XP058053862, DOI: 10.1145/2619239.2626316 ISBN: 978-1-4503-2836-4 the whole document -----	1-21
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/> See patent family annex.		
* Special categories of cited documents :		
"A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed		"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family
Date of the actual completion of the international search		Date of mailing of the international search report
29 May 2017		09/06/2017
Name and mailing address of the ISA/ European Patent Office, P.B. 6818 Patentlaan 2 NL - 2280 HV Rijswijk Tel (+31-70) 340-2040, Fax (+31-70) 340-3016		Authorized officer Raible, Markus

フロントページの続き

(81)指定国 AP(BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, RU, TJ, TM), EP(AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DJ, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JP, KE, KG, KH, KN, KP, KR, KW, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ

(72)発明者 カッタ、ナガ

アメリカ合衆国 9 4 1 0 9 カリフォルニア州 サンフランシスコ ポスト ストリート 7 3
7 アpartment 7 2 1

(72)発明者 ケスラシー、アイザック

アメリカ合衆国 9 4 3 0 4 カリフォルニア州 パロ アルト ヒルビュー アベニュー 3 4
0 1

(72)発明者 ガグ、アディティ

アメリカ合衆国 9 4 3 0 4 カリフォルニア州 パロ アルト ヒルビュー アベニュー 3 4
0 1

Fターム(参考) 5K030 GA13 HC13 KA07 LB06 LC11 MA07 MB02 MC07