

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第6169251号
(P6169251)

(45) 発行日 平成29年7月26日(2017.7.26)

(24) 登録日 平成29年7月7日(2017.7.7)

(51) Int. Cl.		F I	
GO6F	13/00	(2006.01)	GO6F 13/00 357Z
HO4L	12/803	(2013.01)	HO4L 12/803
HO4L	12/743	(2013.01)	HO4L 12/743
HO4L	12/951	(2013.01)	HO4L 12/951

請求項の数 15 (全 78 頁)

(21) 出願番号	特願2016-509083 (P2016-509083)	(73) 特許権者	507303550
(86) (22) 出願日	平成26年4月16日 (2014.4.16)		アマゾン・テクノロジーズ・インコーポレ ーテッド
(65) 公表番号	特表2016-520904 (P2016-520904A)		アメリカ合衆国・98108-1226・ ワシントン州・シアトル・パイオーボック ス・81226
(43) 公表日	平成28年7月14日 (2016.7.14)	(74) 代理人	100064621
(86) 国際出願番号	PCT/US2014/034426		弁理士 山川 政樹
(87) 国際公開番号	W02014/172499	(74) 代理人	100098394
(87) 国際公開日	平成26年10月23日 (2014.10.23)		弁理士 山川 茂樹
審査請求日	平成27年12月8日 (2015.12.8)	(72) 発明者	ソレンソン, ザ サード, ジェームズ・ク リストファー
(31) 優先権主張番号	13/864, 157		アメリカ合衆国・98109-5210・ ワシントン州・シアトル・テリー アヴェ ニュー ノース・410
(32) 優先日	平成25年4月16日 (2013.4.16)		最終頁に続く
(33) 優先権主張国	米国 (US)		

(54) 【発明の名称】 分散型ロードバランサにおける非対称パケットフロー

(57) 【特許請求の範囲】

【請求項1】

複数のロードバランサノードの少なくとも2つが入口サーバとして構成され、
前記複数のロードバランサノードの少なくとも2つがフロートラッカーノードとして構
成される、

前記複数のロードバランサノードと、

複数のサーバノードと、

1つまたは複数の顧客からのパケットフローを、ハッシュ化されたマルチパスルーティ
ング技術に従って、前記入口サーバへと分散させるよう構成されたルータと、

を備えた分散型ロードバランサシステムであり、

顧客のためのパケットフローにおけるパケットを前記ルータから受信し、
前記複数のサーバノードへの前記パケットフローのためのマッピングを前記入口サーバ
が有しないことを決定し、

前記パケットのソースおよび宛先アドレス情報に適用される一貫したハッシュ関数に従
って、前記パケットフローのための少なくとも1つのフロートラッカーノードを決定し、

前記パケットフローのための前記複数のサーバノードの特定の1つへの接続のマッピ
ングを、少なくとも1つのフロートラッカーノードから取得し、

前記特定のサーバノードへの前記パケットフローの1つ以上のパケットを送信する

ように各入口サーバが構成される、

分散型ロードバランサシステム。

10

20

【請求項 2】

前記パケットフローがトランスミッションコントロールプロトコル(TCP)パケットフローである、請求項 1 に記載の分散型ロードバランサシステム。

【請求項 3】

前記複数のロードバランサノードの少なくとも 2 つが、前記 1 つまたは複数の顧客への前記サーバノードからの発信パケットを送信するよう構成された出口サーバとして構成され、

前記パケットフローのための前記出口サーバを選択し、

前記パケットフローのための 1 つまたは複数の発信パケットを、前記選択した出口サーバへと送信する、

10

ように前記サーバノードが構成され、

前記顧客への前記発信パケットを送信するように、前記出口サーバが構成され、

前記パケットフローのための前記選択された出口サーバが、前記パケットフローのための前記入口サーバとは異なるロードバランサノードである、

請求項 1 に記載の分散型ロードバランサシステム。

【請求項 4】

前記サーバノードへの前記パケットの送信の前に、前記入口サーバが前記 1 つまたは複数のパケットをユーザデータグラムプロトコル(UDP)に従ってカプセル化し、前記出口サーバへの前記発信パケットの送信の前に、前記サーバノードが前記発信パケットを UDP に従ってカプセル化し、前記顧客への前記発信パケットの送信の前に、前記出口サーバが前記発信パケットから前記 UDP カプセル封じを取り外す、請求項 3 に記載の分散型ロードバランサシステム。

20

【請求項 5】

前記パケットフローのための前記出口サーバを選択し、

前記カプセル化された受信パケットを前記入口サーバから受信し、

前記パケットから前記 UDP カプセル封じを取り外し、前記パケットを前記サーバノード上のサーバへと伝達させ、

前記サーバノード上の前記サーバから前記発信パケットを取得し、

UDP に従って、前記発信パケットをカプセル化し、

前記カプセル化された発信パケットを前記出口サーバへと送信する

30

ように構成されるロードバランサモジュールを前記サーバノードが含む、

請求項 4 に記載の分散型ロードバランサシステム。

【請求項 6】

前記パケットフローのための前記複数のサーバノードの特定の 1 つへの接続のマッピングを前記少なくとも 1 つのフロートラッカーノードから取得するために、

前記パケットフローのための情報を含むメッセージを、前記入口サーバが前記パケットフローのための 1 次フロートラッカーへと送信し、

前記パケットフローのための前記情報を含むメッセージを、前記 1 次フロートラッカーが前記パケットフローのための 2 次フロートラッカーへと送信し、前記パケットフローのための前記 1 次および 2 次フロートラッカーが異なるロードバランサノードであり、

40

前記 2 次フロートラッカーが、前記パケットフローのための受信確認を、前記顧客へと送信し、

前記入口サーバが、受信確認パケットを前記顧客から受信し、前記受信確認パケットを前記 1 次フロートラッカーへと転送し、

前記 1 次フロートラッカーが、前記パケットフローを受信するための前記サーバノードとして、前記複数のサーバノードの中から前記特定のサーバノードを無作為に選択し、前記特定のサーバノードを示すメッセージを、前記 2 次フロートラッカーへと送信し、

前記 2 次フロートラッカーが、同期メッセージを生成して前記生成された同期メッセージを前記特定のサーバノードへと送信し、

前記 2 次フロートラッカーが、前記パケットフローのための接続情報を前記特定のサー

50

バノードから受信して前記 1 次フロートラッカーへの前記接続情報を含むメッセージを送信し、

前記 1 次フロートラッカーが、前記パケットフローのための前記接続情報を含むメッセージを前記入口サーバへと送信し、前記接続情報が前記パケットフローを前記特定のサーバノードへとマッピングする、

請求項 1 に記載の分散型ロードバランサシステム。

【請求項 7】

前記生成された同期メッセージを前記 2 次フロートラッカーから受信し、

前記サーバノード上のサーバが接続を受け入れることができることを決定し、

前記生成された同期メッセージに従って同期パケットを生成し、前記同期パケットを前記サーバノード上の前記サーバへと伝達し、

前記サーバノード上の前記サーバによって生成された受信確認パケットを遮断し、

前記接続情報を含むメッセージを前記 2 次フロートラッカーへと送信する、

ように構成されるロードバランサモジュールを前記サーバノードが含む、

請求項 6 に記載の分散型ロードバランサシステム。

【請求項 8】

顧客からのパケットフローにおけるパケットの受信であって、1 つまたは複数の顧客から一貫したハッシュ関数に従って前記複数のロードバランサノードへと前記パケットフローを分散させるルータからのパケットの受信、

前記パケットのソースおよび宛先アドレス情報に適用される一貫したハッシュ関数に従っての、前記パケットフローのためのフロートラッカーノードとしての役割を担うロードバランサノードの決定、

前記パケットフローのための複数のサーバノードの 1 つへの接続のマッピングの、前記パケットフローのための前記フロートラッカーノードからの取得、

前記マッピングにより示された、前記サーバノードへの前記パケットフローの 1 つまたは複数のパケットの送信、

を、複数のロードバランサノードのひとつの入口サーバによって実行すること、

を備えた方法。

【請求項 9】

前記パケットフローがトランスミッションコントロールプロトコル (TCP) パケットフローである、請求項 8 に記載の方法。

【請求項 10】

前記パケットフローの前記 1 つまたは複数のパケットの前記サーバノードへの前記送信の前に、ユーザデータグラムプロトコル (UDP) に従って前記パケットをカプセル化する前記入口サーバをさらに備えた、請求項 8 に記載の方法。

【請求項 11】

前記サーバノードによる、前記パケットフローのための出口サーバとしての前記複数のロードバランサノードの 1 つの選択であり、前記パケットフローのための前記選択した出口サーバが、前記パケットフローのための前記入口サーバとは異なるロードバランサノードである、前記複数のロードバランサノードの 1 つの選択と、

前記サーバノードによる、前記パケットフローのための 1 つまたは複数の発信パケットの、前記選択された出口サーバへの送信と、

前記出口サーバによる前記発信パケットの、前記パケットフローの前記顧客への送信と、

、

をさらに備えた、請求項 8 に記載の方法。

【請求項 12】

前記発信パケットの前記出口サーバへの前記送信の前に、ユーザデータグラムプロトコル (UDP) に従って前記発信パケットをカプセル化する前記サーバノードと、

前記発信パケットの顧客への前記送信の前に、前記発信パケットから前記 UDP カプセル封じを取り外す前記出口サーバと、

10

20

30

40

50

をさらに備えた、請求項 1 1 に記載の方法。

【請求項 1 3】

前記フロートラッカーノードが前記パケットフローのための 1 次フロートラッカーノードであり、前記一貫したハッシュ関数に従った一貫したハッシュリングにおける次のロードバランサノードが前記パケットフローのための 2 次フロートラッカーノードである、請求項 8 に記載の方法。

【請求項 1 4】

前記入口サーバによる、少なくとも 1 つのメッセージの前記パケットフローのための前記 1 次フロートラッカーノードへの送信であり、各メッセージが前記ルータから受信された前記パケットフローのパケットを含む前記送信と、

10

前記 1 次フロートラッカーノードによる、前記複数のサーバノードからの前記パケットフローのための前記サーバノードの選択と、

前記 1 次フロートラッカーノードによる、前記選択されたサーバノードを示すパケットフロー情報の前記 2 次フロートラッカーノードへの送信と、

前記 2 次フロートラッカーノードによる、前記サーバノードと前記顧客との通信による、前記パケットフローのための前記選択されたサーバノードへの前記接続の確立の円滑化と、

前記 2 次フロートラッカーノードによる、前記パケットフローのための接続情報の、前記 1 次フロートラッカーノードを通じた前記入口サーバへの送信であり、前記接続情報が前記選択されたサーバノードへの前記パケットフローのマッピングを行う前記送信と、

20

を、前記パケットフローのための複数のサーバノードの 1 つへの接続のマッピングの、前記パケットフローのための前記フロートラッカーノードからの前記取得が備える、

請求項 1 3 に記載の方法。

【請求項 1 5】

前記 2 次フロートラッカーノードによる、前記サーバノード上の前記ロードバランサモジュールへの、生成された同期メッセージの送信と、

前記サーバノード上のサーバが接続を受け入れることができることの決定、

前記生成された同期メッセージに従った同期パケットの生成、

前記同期パケットの、前記サーバノード上の前記サーバへの伝達、

前記サーバノード上の前記サーバにより生成された受信確認パケットの遮断と、

30

前記接続情報を含むメッセージの、前記 2 次フロートラッカーノードへの送信

を、前記サーバノード上の前記ロードバランサモジュールにより実行することと、

を、前記サーバノードと前記顧客との通信により、前記パケットフローのための前記選択されたサーバノードへの前記接続の確立を円滑化するロードバランサモジュールを前記サーバノードが含む、

請求項 1 4 に記載の方法。

【発明の詳細な説明】

【背景技術】

【0001】

従来のロードバランサは通常、複数のネットワークインタフェースコントローラ (NIC) を含む単一の専用装置であり、例えばその一部が顧客からのインバウンドトラフィック / 顧客へのアウトバウンドトラフィックを処理し、残りがロードバランスされたホスト装置 (例、ウェブサーバ等のサーバ) からのアウトバウンドトラフィック / そのようなホスト装置へのインバウンドトラフィックを処理するような、8 つの NIC などを含む。これら従来のロードバランサの帯域またはスループットは通常、顧客側で 40 ギガビット毎秒 (Gbps)、サーバ側で 40 Gbps の範囲である。クラウドコンピューティングサービスのようなネットワークを活用したアプリケーションやネットワークを活用したサービスの規模および範囲の拡大につれて、データセンターには、数百、または数千ものロードバランスを必要とするホスト装置 (例、ウェブサーバ) を格納する可能性が生じる。従来のロードバランサはそのような環境に対応できない。

40

50

【 0 0 0 2 】

さらに従来のロードバランサは通常、ホスト装置から収集されたデータに対して最大接続（またはmax conn）、ラウンドロビン、および/または最小接続（least conn）等の技術を適用し、どのホスト装置が接続の処理を行うかを選択する。また、従来のロードバランサは通常、ホスト装置によりフロントに配置され、それによって顧客からの接続（例、トランスミッションコントロールプロトコル（TCP）接続）を遮断し、ホスト装置およびロードバランサの間で確立されたTCP接続において顧客のトラフィックをホスト装置へ送信する、ホスト装置へのプロキシの役割を担う。したがってこれら従来のロードバランサの使用時には、ホスト装置と顧客は直接TCP接続を通じた通信を行わない。

10

【図面の簡単な説明】

【 0 0 0 3 】

【図1】少なくともいくつかの実施形態による、分散型ロードバランサシステムの実施例のブロック図である。

【図2】少なくともいくつかの実施形態による、図1の分散型ロードバランサシステムにより実装されてもよいロードバランサ方法のハイレベルフローチャートである。

【図3】少なくともいくつかの実施形態による、入口、出口およびフロートラッカーの構成要素を含むロードバランサノードの実施例を示す。

【図4】少なくともいくつかの実施形態による、分散型ロードバランサにおけるルーティングおよびパケットフローを示す。

20

【図5】少なくともいくつかの実施形態による、エッジルータへの入口ノードの提供を示す。

【図6】少なくともいくつかの実施形態による、マルチパスルーティング方法のフローチャートである。

【図7】少なくともいくつかの実施形態による、非対称パケットフローを図示する。

【図8】少なくともいくつかの実施形態による、分散型ロードバランサシステムにおけるパケットフローを示す。

【図9A】少なくともいくつかの実施形態による、分散型ロードバランサシステムにおいて接続を確立する際のパケットフローのフローチャートを提示する。

【図9B】少なくともいくつかの実施形態による、分散型ロードバランサシステムにおいて接続を確立する際のパケットフローのフローチャートを提示する。

30

【図10A】少なくともいくつかの実施形態による、分散型ロードバランサシステムにおけるパケットフローを示す。

【図10B】少なくともいくつかの実施形態による、分散型ロードバランサシステムにおけるパケットフローを示す。

【図10C】少なくともいくつかの実施形態による、分散型ロードバランサシステムにおけるパケットフローを示す。

【図10D】少なくともいくつかの実施形態による、分散型ロードバランサシステムにおけるパケットフローを示す。

【図10E】少なくともいくつかの実施形態による、分散型ロードバランサシステムにおけるパケットフローを示す。

40

【図10F】少なくともいくつかの実施形態による、分散型ロードバランサシステムにおけるパケットフローを示す。

【図10G】少なくともいくつかの実施形態による、分散型ロードバランサシステムにおけるパケットフローを示す。

【図11A】少なくともいくつかの実施形態による、ロードバランサノードの一貫したハッシュリングにおいてメンバーシップに影響を与える処理を示す。

【図11B】少なくともいくつかの実施形態による、ロードバランサノードの一貫したハッシュリングにおいてメンバーシップに影響を与える処理を示す。

【図11C】少なくともいくつかの実施形態による、ロードバランサノードの一貫したハ

50

ッシュリングにおいてメンバーシップに影響を与える処理を示す。

【図11D】少なくともいくつかの実施形態による、ロードバランサノードの一貫したハッシュリングにおいてメンバーシップに影響を与える処理を示す。

【図12】少なくともいくつかの実施形態による、ヘルスチェック間隔に従って各ロードバランサノードにより実行されてもよいヘルスチェック方法のハイレベルフローチャートである。

【図13】少なくともいくつかの実施形態による、別のロードバランサノードからのロードバランサノードのヘルスチェック方法を示す。

【図14】少なくともいくつかの実施形態による、1つまたは複数の他のロードバランサノードのヘルスチェックを行うロードバランサノードを図示する。

10

【図15】少なくともいくつかの実施形態による、サーバノードのヘルスチェックを行うロードバランサノードを示す。

【図16】少なくともいくつかの実施形態による、ロードバランサノード110により維持されてもよい別のノードのヘルス状態を図示する。

【図17】少なくともいくつかの実施形態による、各ロードバランサノードにより維持されてもよいヘルス情報を示す。

【図18A】少なくともいくつかの実施形態によるロードバランサノードの故障の処理を示す。

【図18B】少なくともいくつかの実施形態によるロードバランサノードの故障の処理を示す。

20

【図19A】少なくともいくつかの実施形態による、接続公開技術を図示する。

【図19B】少なくともいくつかの実施形態による、接続公開技術を図示する。

【図20】少なくともいくつかの実施形態による、各ロードバランサモジュールにより実行されてもよい接続公開方法のハイレベルフローチャートである。

【図21】少なくともいくつかの実施形態による、対象のロードバランサノードへの接続公開パケットにおいて受信されるアクティブな接続情報の分散方法のフローチャートである。

【図22】少なくともいくつかの実施形態による、対象のロードバランサノードへの接続公開パケットにおいて受信されるアクティブな接続情報の分散の代替方法を示す。

【図23】少なくともいくつかの実施形態による、ロードバランサノードのソフトウェアスタックアーキテクチャの実施例を示す。

30

【図24】実施形態において用いられてもよいコアパケット処理技術の態様を示す。

【図25】少なくともいくつかの実施形態による、ロードバランサノードにおけるデータフローの処理のためのマルチコアパケットプロセッサの実施例を示す。

【図26】少なくともいくつかの実施形態による、ロードバランサノードにおけるデータフローの処理のためのマルチコアパケットプロセッサの別の実施例を示す。

【図27】少なくともいくつかの実施形態による、ロードバランサノード処理による受信パケットの処理を示す。

【図28】少なくともいくつかの実施形態による、ロードバランサノード処理による発信パケットの処理を示す。

40

【図29】少なくともいくつかの実施形態による、本番環境における分散型ロードバランサを含むロードバランサシステムを示す。

【図30】少なくともいくつかの実施形態による、複数の分散型ロードバランサシステムの構成要素を、単一処理の過程でまたは単一処理として構成、実行することを可能にするメッセージバス機構を組み込む、分散型ロードバランサテストシステムを示す。

【図31】少なくともいくつかの実施形態による、メッセージバスパケットアダプタおよびパケットパイプラインを示す。

【図32】少なくともいくつかの実施形態による、メッセージバスパケットアダプタおよびパケットパイプラインを示す。

【図33A】少なくともいくつかの実施形態による、プロバイダネットワーク環境の実施

50

例を示す。

【図33B】少なくともいくつかの実施形態による、図33Aで示されるプロバイダネットワーク環境の実施例における、分散型ロードバランサの実装を示す。

【図34A】少なくともいくつかの実施形態による、分散型ロードバランサおよびサーバノードの物理ラックの実装の実施例を示す。

【図34B】少なくともいくつかの実施形態による、分散型ロードバランサおよびサーバノードの物理ラックの実装の別の実施例を示す。

【図35】少なくともいくつかの実施形態による、ネットワーク上に1つまたは2つ以上の分散型ロードバランサが実装されるネットワーク環境の実施例を示す。

【図36】いくつかの実施形態において用いられてもよいコンピュータシステムの実施例を示すブロック図である。

【0004】

本明細書内ではいくつかの実施形態および例示的な図面のために例として実施形態が記載されているが、実施形態が本明細書内で示される実施形態または図面のみに限られないことを当業者は認識すべきである。本明細書に記載の図面および詳細な説明は、本明細書で開示される特定の形式に実施形態を限定することを意図するものではなく、反対に、添付の請求項で定められる精神および範囲から逸脱することのないすべての修正、均等物および代替物を包括する意図があることを理解すべきである。本明細書で使用される見出しは構成上の目的のためのみであり、説明また請求の範囲を限定するために用いる意図はない。本願を通じて用いられる通り、「してもよい(may)」という語は、必須の意(すなわち、義務)ではなく、むしろ許容の意(すなわち、可能性を示す意味)で用いられる。同様に、「含む(include)」、「including」および「includes)」の語は含むことを意味するが、それに限定されない。

【発明を実施するための形態】

【0005】

ネットワーク環境における分散型ロードバランサ方法およびシステムの様々な実施形態が記載される。分散型ロードバランサ方法およびシステムの実施形態は、様々なネットワーク環境における分散型ロードバランサの実施形態によって実装されてもよいように記載される。例として、分散型ロードバランサの実施形態は、図33Aおよび33Bに示されるようなプロバイダネットワーク1900等のローカルネットワーク上のインターネットや宛先、通常はサーバ(例、ウェブサーバ、アプリケーションサーバ、データサーバ等)のような外部ネットワーク上の顧客間のパケットフロー、例えばトランスミッションコントロールプロトコル(TCP)技術を用いたパケットフロー等を円滑化し、維持するために用いられてもよい。本明細書に記載の実施形態は主にTCPパケットフローの処理に関するが、実施形態はTCP以外のデータ通信プロトコルや、パケットフロー処理以外の用途に適用されてもよいことに留意する。

【0006】

分散型ロードバランサは、特定の顧客と選択したサーバ(例、ウェブサーバ)との間のTCPパケットフローを円滑化し維持する役割を担ってもよい。しかし分散型ロードバランサは、従来のロードバランサにおけるやり方では、顧客からのTCPフローを遮断せず、また、サーバへのプロキシの役割も担わない。代わりに分散型ロードバランサのロードバランサノードが顧客から受信されたTCPパケットを対象のサーバヘルディングし、サーバがそのTCPスタックを用いて顧客へのTCP接続を管理する。すなわち、サーバが顧客からのTCPパケットフローを遮断する。

【0007】

また、従来のロードバランサ技術において行われるように、サーバから収集された情報に適用されるロードバランサ技術またはアルゴリズムに基づいて、どのサーバが接続要求を送信するかに関して決定を下す1つまたは複数のロードバランサノードの代わりに、ロードバランサノードが新規接続要求を受信するサーバを無作為に選択してもよく、サーバノード上にある分散型ロードバランサの構成要素がそれぞれのサーバの現状の1つまたは

10

20

30

40

50

複数の測定基準に基づき、選択したサーバが新規接続要求を受け入れるか拒否するかに関する決定をローカルで下す。したがって、どのサーバが接続要求を受け入れるかに関する決定は、1つまたは複数のロードバランサノードから、接続を処理するサーバノードへと移行する。すなわち、接続要求の送信により近い場所およびタイミングへと決定が移行される。

【0008】

顧客とサーバ間のパケットフローを円滑化し、維持するため、分散型ロードバランサの実施形態は様々な技術を用いてもよい。またそのような技術は、マルチパスルーティング技術、一貫したハッシュ技術、分散型ハッシュテーブル(DHT)技術、境界ゲートウェイプロトコル(BGP)技術、メンバーシップトラッキング、ヘルスチェック、接続公開、およびパケットのカプセル化および脱カプセル化を含むが、これらに限定されない。これらは分散型ロードバランサシステムの他の態様と同様に、図面に関連して以下に記載される。

分散型ロードバランサシステム

【0009】

図1は、少なくともいくつかの実施形態による、分散型ロードバランサシステムの実施例のブロック図である。分散型ロードバランサの実施形態は、例えば図33Aおよび33Bで示されるサービスプロバイダのプロバイダネットワーク1900等のネットワーク100において実装されてもよい。分散型ロードバランサシステムにおける顧客のパケット処理のハイレベル概観図に示される通り、ネットワーク100の1つまたは複数の顧客160は、例えばインターネット等の外部ネットワーク150を通じて、ネットワーク100の境界ルータ102に接続してもよい。境界ルータ102は、顧客160からの受信パケット(例、TCPパケット)を分散型ロードバランサの構成要素エッジルータ104へとルーティングしてもよい。エッジルータ104は受信パケットを、分散型ロードバランサシステムのロードバランサノードレイヤー上のロードバランサ(LB)ノード110へとルーティングする。少なくともいくつかの実施形態において、エッジルータ104は例えば等価コストマルチパス(ECMP)ハッシュ技術のようなフローごとにハッシュ化されたマルチパスルーティング技術に従って、ルーティングの決定を行ってもよい。次にロードバランサノード110がパケットをカプセル化し(例、ユーザデータグラムプロトコル(UDP)に従って)、ネットワーク100上のネットワークファブリック120(例、L3ネットワーク)を通じて、カプセル化されたパケットをサーバノード130上のローカルロードバランサモジュール132へとルーティングする。ファブリック120は、1つまたは複数のネットワーク装置または構成要素を含んでもよく、そのようなネットワーク装置または構成要素にはスイッチ、ルータ、およびケーブルが含まれるがこれらに限定されない。サーバノード130上では、ローカルロードバランサモジュール132がパケットを脱カプセル化し、顧客のTCPパケットをサーバ134のTCPスタックに送信する。サーバノード130上のサーバ134はその後そのTCPスタックを利用し、顧客160への接続を管理する。

【0010】

図2は少なくともいくつかの実施形態による、図1の分散型ロードバランサシステムにより実装されてもよいロードバランサ方法のハイレベルフローチャートである。分散型ロードバランサシステムの実施形態は、従来のロードバランサ技術において行われるような複数の宛先(例、ウェブサーバ)間の負荷の割り当てに関する困難問題の解決を行わなくてもよい。例えば従来のロードバランサは通常、最大接続、ラウンドロビン、および/または最小接続等の技術やアルゴリズムを用い、どのサーバが接続の処理を行うべきかを選択する。しかしこれら技術には欠点があり、特にロードバランサに関する決定を行うために用いられるデータがほぼすぐに古くなってしまいうような分散型システムにおいては良好に行うことが難しい。少なくともいくつかの分散型ロードバランサシステムの実施形態においては、従来のロードバランサにおいて行われるように接続要求を満たすために1つまたは複数のロードバランサ技術を用いてサーバノード130の選択を試みる代わりに、ロ

10

20

30

40

50

ロードバランサノードレイヤー上のロードバランサノード110が、顧客の接続のための要求を受信するサーバノード130を無作為に決定してもよい。そのサーバノード130が自身に負荷がかかり過ぎていると見なす場合は、サーバノード130はロードバランサノード110に接続要求を送信し返し、それによってサーバノード130が現在接続を処理できないことをロードバランサノード110に伝えてもよい。ロードバランサノードレイヤーはその後、接続要求を受信する別のサーバノード130を無作為に決定してもよいし、もしくはその代わりに、要求を行う顧客160に対してエラーメッセージを返信し、現在接続を確立できないことを顧客160に伝えてもよい。

【0011】

図2の10で示すように、分散型ロードバランサシステムのロードバランサノードレイヤーがソースから通信セッション（例、TCP接続）のための要求を受信する。ソースは例えば、分散型ロードバランサシステムを実装するネットワーク100への外部ネットワーク150上の顧客160であってもよい。少なくともいくつかの実施形態においては、要求はネットワーク100の境界ルータ102で顧客160から受信されてもよく、また、エッジルータ104にルーティングされてもよい。エッジルータ104は例えば顧客160からの特定の接続要求のルーティング先となるロードバランサノード110を擬似ランダムに選択するフローごとの等価コストマルチパス（ECMP）ハッシュ技術を用いて、受信パケットをロードバランサノードレイヤー上のロードバランサ（LB）ノード110にルーティングする。

【0012】

20で示すように、ロードバランサノードレイヤーが宛先ノードを無作為に選択し、選択した宛先ノードに接続要求を転送する。宛先ノードは例えば、ロードバランサによってフロントに配置された複数のサーバノード130のうちの1つであってもよい。少なくともいくつかの実施形態において、ロードバランサレイヤー上のロードバランサノード110はすべての既知のサーバノード130の中から接続要求を受信するサーバノード130を無作為に選択してもよい。しかしいくつかの実施形態においては、すべての既知のサーバノード130の中から純粋に無作為に選択する以外の方法を用いて接続要求を受信するサーバノード130を選択してもよい。例えばいくつかの実施形態においては、サーバノード130の無作為な選択の重みづけを行うために、ロードバランサノード110によってサーバノード130に関する情報が利用されてもよい。実施例として、異なるサーバノード130が異なる種類の装置であるかまたは異なるCPUによって構成されているために異なる能力や可能性を有するとロードバランサノード110が認識している場合、無作為な選択をサーバノード130の1つまたは複数の特定の種類または構成の方に（またはそれを避けるように）偏らせるために情報が用いられてもよい。

【0013】

30で示すように、宛先ノードが、通信セッションを受け入れることが可能か決定する。少なくともいくつかの実施形態において、サーバノード130上のローカルロードバランサ（LB）モジュール132がそれぞれのサーバ134の現状の1つまたは複数の測定基準に基づき、サーバノード130上のそれぞれのサーバ134が新規接続を受け入れることが可能であるかどうかを決定する。

【0014】

40において接続要求が受け入れられる場合は、その後50で示すように宛先ノードが接続を処理できることを宛先ノードからロードバランサノードレイヤーに伝える。その後60で示すようにソース（例、顧客160）と宛先ノード（例、サーバノード130上のサーバ134）との間にロードバランサノードレイヤーを通じて通信セッションが構築される。少なくともいくつかの実施形態において、サーバノード130上のサーバ134がTCPスタックを用いて顧客160への接続を管理する。

【0015】

40において接続要求が受け入れられない場合は、その後70で示すように宛先ノードがロードバランサノードレイヤーに通知し、メソッドは要素20に戻ってもよい。ロード

10

20

30

40

50

バランスノードレイヤーはその後20において別の宛先ノードを無作為に選択してもよいし、もしくはその代わりに、要求を行う顧客160に対して現在接続を確立できないことを顧客160に伝えてもよい。顧客160は、必ずそうするわけではないが、要素10においてメソッドを再び開始するために接続要求を再提出してもよいことに留意する。

【0016】

図1を再び参照する。分散型ロードバランスシステムの少なくともいくつかの実施形態は、コモディティハードウェアを用いてネットワーク100上のエッジルータ104で受信された顧客のトラフィックをネットワーク100上のサーバノード130へとルーティングしてもよい。分散型ロードバランスの少なくともいくつかの実施形態は、複数のロードバランスノード110を含むロードバランスノードレイヤーを含んでもよい。少なくともいくつかの実施形態において、各ロードバランスノード110はロードバランスノードレイヤーにおける複数の役割のうち、1つまたは複数の役割を担ってもよい。これらロードバランスノード110の役割には、入口ノード、また出口ノード、そして(所与のパケットフローのための1次フロートラッカーまたは2次フロートラッカーとしての)フロートラッカーノードの役割が含まれてもよい。少なくともいくつかの実施形態において各ロードバランスノード110は、ラック搭載型のコモディティコンピューティング装置等の個別のコンピューティング装置として、またはそういった個別のコンピューティング装置において、ロードバランスノードレイヤー上に実装されてもよい。ロードバランスノード110は一般的に、特定のパケットフローのための役割のうちただ1つ(しかし可能であれば2つまたは3つ)を担うが、少なくともいくつかの実施形態において、各ロードバランスノード110は入口ノード、また出口ノード、そして(所与のパケットフローのための1次フロートラッカーまたは2次フロートラッカーとしての)フロートラッカーノードの3つの役割それぞれを担う。しかし少なくともいくつかの実施形態において、ロードバランスノード110は特定のパケットフローのための1次フロートラッカーおよび2次フロートラッカーの両方の役割を担うことはできないことに留意する。その代わりにいくつかの実施形態においては、各ロードバランスノード110が3つの役割のうちただ1つを担ってもよい。この実施形態においては、コンピューティング装置の個別の組が、特に入口ノード、出口ノード、およびフロートラッカーノードとしてロードバランスノードレイヤー上で実装されてもよい。

【0017】

少なくともいくつかの実施形態において、パケットフローのための1次および2次フロートラッカーの決定するために、一貫したハッシュおよび一貫したハッシュリング技術が適用されてもよい。顧客からの各パケットフローは、例えば顧客のIPアドレス、顧客用ポート、サーバ(パブリック)IPアドレス、およびサーバポートから成る4タプルによって一意に識別されてもよい。この識別子は、顧客およびパブリックエンドポイントのペアを示すCPまたはCcPpとして略されてもよい。与えられた任意のTCPフロー(またはCPペア)に関連するパケットは、エッジルータ104からのハッシュ化されたマルチパス(例、ECMP)フロー分散により、入口サーバ112として動作するいずれのロードバランスノード110上にも表れることができる。パケットが入口ノードとして動作するロードバランスノード110に到達する際に、どのロードバランスノード110がパケットフローのための状態維持の役割を担うか(すなわち、1次フロートラッカーノード)を入口ノードが決定できるように、一貫したハッシュが用いられる。どのロードバランスノード110がパケットフローのための状態維持の役割を担うかを決定するために、CPペアは入口ノードにより一貫したハッシュリング内にハッシュ化されてもよい。一貫したハッシュリング内でパケットフローのためのCPペアの一貫したハッシュに従って決定されたノード110が、パケットフローのための1次フロートラッカーの役割を担うノード110である。少なくともいくつかの実施形態において、一貫したハッシュリングにおける後続ノードがパケットフローのための2次フロートラッカーの役割を担う。

【0018】

図3は少なくともいくつかの実施形態による、3つの役割すべて(入口、出口、および

10

20

30

40

50

フロートラッカー)を実装する構成要素を含む、ロードバランサ(LB)ノード110の実施例を示す。この実施例において構成要素である入口サーバ112は、1つまたは複数の顧客からインバウンドTCPパケットを受信し、TCPパケットをカプセル化されたパケットとして1つまたは複数のサーバに送信する、入口の役割を果たす。構成要素である出口サーバ114は、1つまたは複数のサーバからアウトバウンドのカプセル化されたパケットを受信し、脱カプセル化されたTCPパケットを1つまたは複数の顧客に送信する、出口の役割を果たす。構成要素であるフロートラッカー116は、1つまたは複数の顧客160と1つまたは複数のサーバ134との間に確立された1つまたは複数のパケットフローのための1次または2次フロートラッカーの役割を果たす。それぞれの顧客160から受信された接続要求に応じて顧客とサーバ134の1つとの間のTCP接続を開始するため、またはパケットフローのためにマッピング情報を取得するため、入口サーバ112はまた、ロードバランサノード110上のフロートラッカー116または別のロードバランサノード110上のフロートラッカー116と通信してもよい。

10

ロードバランサノード

【0019】

図1を再び参照する。少なくともいくつかの実施形態において、ロードバランサノードレイヤーにおけるロードバランサノード110は、ネットワーク上の1つまたは複数のルータ104から顧客のトラフィック(例えばTCPパケット等のパケット)を受信し、ファブリック120上の分散型ロードバランサシステムにより用いられるプロトコル(例、ユーザデータグラムプロトコル(UDP))に従ってパケットをカプセル化する。ロードバランサノードレイヤーはその後ファブリック120を介してカプセル化されたパケットを宛先サーバノード130に転送する。各サーバノード130はロードバランサシステムの構成要素であるローカルモジュール132を含む。モジュール132は本明細書内ではロードバランサモジュールまたは単にLBモジュールと称されてもよく、サーバノード130上のソフトウェア、ハードウェア、またはそれらの組み合わせにおいて実装されてもよい。各サーバノード130において、それぞれのロードバランサモジュール132がパケットを脱カプセル化し、通常のTCP処理のためにTCPパケットをローカルTCPスタックに送信する。少なくともいくつかの実施形態においては、ロードバランサノードレイヤーがすべての顧客サーバのTCPフローのために状態情報を維持してもよい。しかし、ロードバランサノードレイヤー上のロードバランサノード110は、TCPフローに関する一切の解釈を実行することはできない。各フローは、それぞれのサーバノード130のサーバ134と顧客160との間で管理される。分散型ロードバランサシステムはTCPパケットが正確な宛先サーバ134に確実に到着するようにする。各サーバノード130におけるロードバランサモジュール132は、ロードバランサノード110から受信した顧客の接続要求に応じてそれぞれのサーバ134が新規接続を受け入れるか拒否するかに関する決定を下す。

20

30

【0020】

少なくともいくつかの実施形態において、分散型ロードバランサシステムは一貫したハッシュ技術を、例えば、どのサーバノード130が特定のTCPパケットフローのための役割を担うかについてどの1つまたは複数のロードバランサノード110が記憶すべきかなどを決定するために用いてもよい。一貫したハッシュ技術を利用し、ロードバランサノードレイヤーにおけるロードバランサノード110は一貫したハッシュリングとして見なされてもよく、ロードバランサノード110はリングにおけるメンバーシップのトラッキングを継続し、一貫したハッシュ関数に従って、特定のパケットフローのための役割を果たすリングにおける特定のメンバーを決定してもよい。少なくともいくつかの実施形態においては、顧客160とサーバ134との間の各パケットフローのトラッキングの役割を担う2つのロードバランサノード110が存在する。これらのノード110は、1次フロートラッカー(PFT)ノードおよび2次フロートラッカー(SFT)ノードと称されてもよい。少なくともいくつかの実施形態において、1次フロートラッカーはフローのための一貫したハッシュリングにおける第1のロードバランサノード110であり、2次フロ

40

50

ートラッカーは一貫したハッシュリングにおける次のまたは後続の、1次フロートラッカーノードとは異なるロードバランサノード110である。この場合、1次フロートラッカーノードが故障した際には、その後2次フロートラッカーノードが新規1次フロートラッカーになってもよく、別のロードバランサノード110（例、一貫したハッシュリングにおける次のノード110）が2次フロートラッカーの役割を担ってもよい。少なくともいくつかの実施形態において、ロードバランサノード110は所与のケットフローのための1次フロートラッカーおよび2次フロートラッカーの両方の役割を担うことはできないことに留意する。一貫したハッシュリングにおけるこのまたは別のメンバーシップの変更については、本明細書にて後述される。少なくともいくつかの実施形態において、ロードバランサの実装のための構成情報（例、現在実装されているロードバランサノード110およびサーバノード130の1つまたは複数の信頼すべきリスト）は、例えばファブリック120を通じてロードバランサノード110に接続される1つまたは複数のサーバ装置上で実装されてもよい分散型ロードバランサシステムの構成要素である、構成サービス122によって維持されてもよい。

【0021】

少なくともいくつかの実施形態において、1次および2次フロートラッカーノードとしての役割に加えて、ロードバランサノード110は所与のフローのために、他の2つのうち1つの役割を果たしてもよい。すなわち、入口ノードの役割および出口ノードの役割である。ケットフローのための入口ノードは、エッジルータ104からそれぞれのケットフローを受信し、ケットフローを（カプセル化されたケットとして）ファブリック120を通じてサーバノード130上の選択されたサーバ134に転送する、ロードバランサノード110である。入口ノードは、実際の顧客データ（TCPデータケット）をそれぞれの宛先サーバノード130へと移動させる唯一のロードバランサノード110である。入口ノードは、宛先サーバノード130上のそれぞれのロードバランサモジュール132へのTCPフローのマッピングを維持し、顧客のトラフィックをどのロードバランサモジュール132へと転送すべきかを把握する。出口ノードは、ファブリック120を通じてサーバノード130から受信されたケットフローのための応答トラフィックを、境界ネットワークを通じてそれぞれの顧客160に転送する役割を担うロードバランサノード110である。ロードバランサモジュール132は、応答ケットをロードバランサプロトコル（例、UDP）に従ってサーバ134から得た応答ケットをカプセル化し、カプセル化された応答ケットを、ファブリック120を通じてフローのためのそれぞれの出口ノードに送信する。出口ノードはステートレスであり、単にデータケットを脱カプセル化して境界ネットワーク上の応答ケット（例、TCPケット）を境界ルータ102に送信し、外部ネットワーク150を通じてそれぞれの顧客160に送達する。

【0022】

上述のように少なくともいくつかの実施形態において、各ロードバランサノード110は異なるケットフローのための入口ノード、出口ノード、および/またはフロートラッカーノード（1次または2次いずれかのフロートラッカーとして）の役割を果たす。ロードバランサノードレイヤー上の単一のロードバランサノード110は、ノードが何のケットフローを処理しているかに応じて、役割のうちいずれか1つを担ってもよい。例えば、少なくともいくつかの実施形態において、ロードバランサノード110は、1つのケットフローのために入口ノードの役割を、また別のケットフローのために1次または2次フロートラッカーの役割を、そしてさらに別のケットフローのために出口ノードの役割を果たしてもよい。また少なくともいくつかの実施形態においてロードバランサノード110は、例えば所与のケットフローのために入口ノードおよび1次（または2次）フロートラッカーノードの役割をというように、同一のケットフローのために複数の役割を果たしてもよい。しかし少なくともいくつかの実施形態において、ロードバランサノード110は冗長化および回復を目的として、同一のケットフローのために1次および2次フロートラッカーの両方の役割を担うことはできない。

【0023】

10

20

30

40

50

上記は、各ロードバランサノード110が入口サーバ、出口サーバ、およびフロートラッカーの3つの役割のうちいずれの役割を担ってもよい実施形態である。しかしいくつかの実施形態においては、コンピューティング装置の異なるグループが、ロードバランスシステムにおける異なる役割に割り当てられてもよい。例えばいくつかの実施形態においては、個別のコンピューティング装置上でそれぞれ実装される入口ノード、出口ノードおよびフロートラッカーノードの異なる組があってもよい。別の実施例として、いくつかの実施形態においては、コンピューティング装置の別の組が出口ノードの役割のみを担う一方で、コンピューティング装置の組の1つが入口ノードおよびフロートラッカーノードの両方の役割を担ってもよい。

ロードバランサモジュール

【0024】

上述のように各サーバノード130は、ロードバランスシステムの構成要素であるローカルロードバランサモジュール132を含む。モジュール132は、サーバノード130のソフトウェア、ハードウェア、またはそれらの組み合わせにおいて実装されてもよい。少なくともいくつかの実施形態において、サーバノード130のロードバランサモジュール132は、発信パケットのカプセル化および受信パケットの脱カプセル化、ノード130上のサーバ134のためのロードバランスに関するローカルでの決定、ならびに接続公開の、3つの主要な役割を果たしてもよい。これら3つの役割は以下に簡単に記され、さらに本明細書で詳細に後述される。

【0025】

分散型ロードバランスシステムの少なくともいくつかの実施形態は、TCP接続を遮断せず、パケットのスプーフィングも行わない。ロードバランサノードレイヤーを通じて送信されるすべてのパケットのソースおよび宛先IPアドレスは、パケットフローに關与するエンドポイント（すなわち、顧客160およびサーバ134）の実際のIPアドレスである。スプーフィングの代わりに、これら実施形態は、ロードバランサノード110とサーバノード130との間で送信される、例えばUDPパケット等のすべてのパケットを、ファブリック120上でカプセル化する。フローのための入口ノードの役割を担うロードバランサノード110からサーバノード130に到着するパケットフロー内のインバウンドパケットはロードバランサノード110によってカプセル化されるため、パケットはノード130上のサーバ134のために脱カプセル化され、ローカルホストのTCPフローへと方向を変えられる必要がある。ノード130上のロードバランサモジュール132がこの脱カプセル化を行う。同様に、サーバ134からのパケットフローのための発信データパケットは、ロードバランサモジュール132によってカプセル化され、ファブリック120を通じてパケットフローのための出口ノードの役割を果たすロードバランサノード110へと送信される。

【0026】

少なくともいくつかの実施形態において、サーバノード130上のロードバランサモジュール132は、それぞれのサーバノード130上のサーバ134のためのロードバランスに関してローカルで決定を下す。特にサーバノード130上のロードバランサモジュール132は、それぞれのサーバ134が新規TCP接続要求の受信に応じて別のTCPフローを受け入れるかどうかの決定を行う。上述のように、ロードバランサノード110はロードバランサモジュール132に送信されるすべてのパケットをカプセル化するため、ロードバランサモジュール132は実際にはTCP同期(SYN)パケットを顧客160から受信することはなく、その代わりにロードバランサモジュール132はカプセル化プロトコル(例、UDP)に従ったフロートラッカー116からの接続要求メッセージを、受信する。ロードバランサモジュール132はこの接続要求メッセージを受け入れるか拒否することができる。ロードバランサモジュール132が接続要求メッセージを受け入れる場合、ロードバランサモジュール132はローカルホストに向けたSYNパケットを作成する。ローカルホストが接続を受け入れる場合、これはそれぞれの顧客接続を処理する実際のTCPスタックとなる。

10

20

30

40

50

【 0 0 2 7 】

少なくともいくつかの実施形態において、接続要求メッセージを受け入れるべきかどうかに関する決定を下すため、ロードバランサモジュール 1 3 2 はサーバノード 1 3 0 上の現在のリソース消費に関する 1 つまたは複数の測定基準を確認し、新規接続の処理のために使用可能なリソースが十分にある場合は、ロードバランサモジュール 1 3 2 が接続を受け入れる。少なくともいくつかの実施形態において、ロードバランサモジュール 1 3 2 により考慮されてもよいリソースの測定基準は、CPU 使用率、最新の帯域占有量、および確立された接続の数のうちの 1 つまたは複数を含んでもよいが、これらに限定されない。いくつかの実施形態においては、これらの測定基準の代わりに、またはこれらの測定基準に加えて他の測定基準が考慮されてもよい。例えばいくつかの実施形態においては、ロードバランサモジュールがサーバ待ち時間（すなわち、要求がサーバ接続バックログに留まる時間）を測定基準として考慮してもよく、また、サーバ待ち時間が閾値を超える場合には接続要求を拒否してもよい。これらのおよび/または他の測定基準を利用して、ロードバランサモジュール 1 3 2 はそれぞれのサーバ 1 3 4 のために、サーバ 1 3 4 が新規パケットフローを受け入れるか拒否するかに関する決定を下すことができる。少なくともいくつかの実施形態において、リソース利用率（例、N % 利用）は個別の、または組み合わせた、また、閾値（例、9 0 % 利用）と比較した測定基準から決定されてもよい。決定されたリソース利用率が閾値以上である場合、または、接続を追加することによりリソース利用率が閾値を超えることになる場合には、その後接続要求が拒否されてもよい。

10

【 0 0 2 8 】

少なくともいくつかの実施形態において、ロードバランサモジュール 1 3 2 は、接続要求メッセージを拒否すべきかどうか決定するために確率論的手法を実装してもよい。上記のようにリソース利用が閾値以上である場合にすべての接続要求を拒否する代わりに、この方法では 2 つ以上の異なる利用レベルにおいて異なる確率で接続要求を拒否してもよい。例えばリソース利用が 8 0 % の場合にロードバランサモジュール 1 3 2 が 2 0 % の確率で接続要求を拒否してもよく、リソース利用が 9 0 % の場合にロードバランサモジュール 1 3 2 が 2 5 % の確率で接続要求を拒否してもよく、リソース利用が 9 5 % の場合にロードバランサモジュール 1 3 2 が 5 0 % の確率で接続要求を拒否してもよく、そしてリソース利用が 9 8 % 以上の場合にロードバランサモジュール 1 3 2 がすべての接続要求を拒否してもよい。

20

30

【 0 0 2 9 】

少なくともいくつかの実施形態において、各接続要求メッセージには、ロードバランサモジュール 1 3 2 が接続要求メッセージを拒否した回数が含まれていてもよい。自身が閾値以上の回数に渡り拒否されたことをロードバランサモジュール 1 3 0 により受信された接続要求メッセージが示す場合、パフォーマンス測定基準がサーバノード 1 3 0 の接続要求を拒否すべきだと示していても、ロードバランサモジュール 1 3 0 は接続を受け入れてもよい。

【 0 0 3 0 】

場合によっては、接続要求メッセージの送信先であるロードバランサモジュール 1 3 2 がすべて、接続要求を拒否する可能性もある。少なくともいくつかの実施形態において、接続要求メッセージがロードバランサモジュール 1 3 2 同士の間で無期限にバウンスされるのを防ぐために、各接続要求メッセージには生存時間が与えられてもよい。この生存時間の期限が切れると、フロートラッカーノードは要求を遮断し、要求を現在伝達することができないことをそれぞれの顧客 1 6 0 に通知してもよい。

40

【 0 0 3 1 】

少なくともいくつかの実施形態において、サーバノード 1 3 0 上のロードバランサモジュール 1 3 2 はまた、ロードバランサノード 1 1 0 への接続公開を行う。少なくともいくつかの実施形態において、接続公開を行うために、各ロードバランサモジュール 1 3 2 は定期的または非定期的（例、1 秒に 1 回）サーバノード 1 3 0 上のルーティングテーブル（例、ネットスタットルーティングテーブル）を確認し、アクティブな接続（TCPF

50

ロー)のリストをロードバランサノード110に公開する。所与のパケットフローの存在について通知を受ける必要のあるロードバランサノード110は、それぞれのパケットフローのために入口ノードおよび1次および2次フロートラッカーとしての役割を担っているロードバランサノード110である。いくつかの実施形態においては、ロードバランサモジュール132は一貫したハッシュ技術を用いて、サーバノード130上のアクティブなTCPフローについて通知を受ける必要のあるロードバランサノード110のリストをフィルターにかけてもよい。例えばロードバランサモジュール132は一貫したハッシュリングに従って、どのロードバランサノード110が所与のパケットフローのために1次および2次フロートラッカーの役割を担っているかを決定してもよい。いくつかの実施形態においては、ロードバランサモジュール132は各パケットフローのためにどのロードバランサノード110が最後にデータパケットをロードバランサモジュール132に送信したかをトラッキングし、この情報を用いてどのロードバランサノード110がパケットフローのための入口ノードの役割を担っているかを決定する。これは入口ノードのみが顧客データをロードバランサモジュール132に転送するためである。いくつかの実施形態においてはロードバランサモジュール132がその後、パケットフローについて通知を受ける必要があると決定したロードバランサノード110の各々のためのメッセージを作成し、メッセージをロードバランサノード110へ送信してそれぞれのサーバノード130が1つまたは複数の接続を1つまたは複数の顧客160に対してまだ維持していることをノード110に通知する。このロードバランサモジュール132によるロードバランサノード110への接続公開はリースのロードバランサノード110への延長と見なされてもよい。ロードバランサノード110が一定の時間(例、10秒)内に特定のパケットフローを示す接続公開メッセージを受信しない場合、ロードバランサノード110はその後それぞれのパケットフローについて忘れることもできる。

ロードバランサノードへのマルチパスルーティング

【0032】

図4は少なくともいくつかの実施形態による、分散型ロードバランサにおけるルーティングおよびパケットフローの態様を示す。少なくともいくつかの実施形態において、各入口ノード(入口ノードは図4において入口サーバ112として示される)が分散型ロードバランサのために、例えば境界ゲートウェイプロトコル(BGP)を通じて1つまたは複数のパブリックエンドポイント(例、IPアドレスおよびポート)をエッジルータ104へとルーティングする能力を提供する。少なくともいくつかの実施形態においては、各入口ノードがBGPセッションを通じて自身をエッジルータ104に提供するのではなく、1つまたは複数の他の入口ノード、例えば2つの近傍ノードがエッジルータ104とのBGPセッションを確立し、図5で示されるように入口ノードを提供してもよい。

【0033】

従来のロードバランサは通常、単一のパブリックエンドポイントとしての役割しか担うことができない。それに対して、分散型ロードバランサの実施形態は複数のロードバランサノード110が単一のパブリックエンドポイントとしての役割を担うことを可能にする。これによりルータの能力に応じて、すべての入口サーバ112にルーティングされた単一のパブリックIPアドレスが1つまたは複数のエッジルータ104を通じて全帯域(例、160Gbps)の処理を行ってもよい構成が可能になる。少なくともいくつかの実施形態においてこれを達成するために、1つまたは複数のエッジルータ104がレイヤー4のフローごとにハッシュ化されたマルチパスルーティング技術、例えば等価コストマルチパス(ECMP)ルーティング技術を利用し、各々が同一のパブリックIPアドレスを提供する複数の入口サーバ112を介してトラフィックを分散してもよい。一般的に、1つまたは複数のエッジルータ104のフローハッシュの一部としてフローのためのレイヤー4のソースおよび宛先ポートを用いてすべての入口サーバ112に受信パケットを分散することにより、入口サーバ112としての役割を担う同一のロードバランサノード110へとルーティングされた各接続のためのパケットに、故障したパケットを避けさせてもよい。しかしいくつかの実施形態においては、1つまたは複数のエッジルータ104は他の

10

20

30

40

50

技術を利用して入口サーバ 1 1 2 を通じてトラフィックを分散させてもよいことに留意する。

【 0 0 3 4 】

図 4 はまた、2 つ以上の分散型ロードバランサがネットワーク 1 0 0 上に実装されてもよいことを示す。2 つ以上の分散型ロードバランサは、複数のサーバ 1 3 0 をフロントに配置し各々が異なるパブリック IP アドレスを提供するそれぞれ独立したロードバランサとして動作してもよく、またはその代わりに、図 4 に示されるように、2 つ以上の分散型ロードバランサがそれぞれ同一の IP アドレスを提供してもよく、ハッシュ技術（例、レイヤー 4 のフローごとにハッシュ化されたマルチパスルーティング技術）は 1 つまたは複数の境界ルータ 1 0 2 において、次にパケットフローをそれぞれの区分する入口サーバ 1 1 2 に分散させるエッジルータ 1 0 4 からパケットフローを隔てるために用いられてもよい。

10

【 0 0 3 5 】

図 5 は少なくともいくつかの実施形態による、エッジルータへの入口ノードの提供のための境界ゲートウェイプロトコル（BGP）の利用を示す。この実施例においては、ロードバランサの実装において入口ノード 1 1 0 A ~ 1 1 0 D の役割を担う、4 つのロードバランサノードがある。エッジルータ 1 0 4 は顧客（図示せず）からの受信パケットをロードバランサノード 1 1 0 にルーティングする。少なくともいくつかの実施形態において、エッジルータ 1 0 4 はレイヤー 4 のフローごとにハッシュ化されたマルチパスルーティング技術、例えば等価コストマルチパス（ECMP）ルーティング技術に従って、ルーティングに関する決定を下してもよい。

20

【 0 0 3 6 】

少なくともいくつかの実施形態において、エッジルータ 1 0 4 は、ロードバランサの実装において、入口ノード 1 1 0 により開始されるセッションを提供する境界ゲートウェイプロトコル（BGP）技術を通じて顧客のトラフィックを受信することが現在可能である入口ノード 1 1 0 に関して把握する。各入口ノード 1 1 0 は BGP を用いて自身をエッジルータ 1 0 4 に提供することができる。しかし、BGP は通常、収束に比較的長い時間がかかる（3 秒以上）。各入口ノード 1 1 0 が BGP を介して自身を提供するこの技術を用いる場合、入口ノード 1 1 0 の故障時には、ネットワーク時間（3 秒以上）において、エッジルータ 1 0 4 上の BGP セッションがタイムアウトするまでに相当な時間がかかる可能性があり、したがって、エッジルータ 1 0 4 が失敗による終了について把握し現在の TCP フローを入口ノード 1 1 0 に再度ルーティングするまでには相当な時間がかかる可能性がある。

30

【 0 0 3 7 】

BGP の収束問題を回避し、ノード 1 1 0 の故障時の回復を早めるため、少なくともいくつかの実施形態において、入口ノード 1 1 0 が BGP セッションを介して自身をエッジルータ 1 0 4 に提供する代わりに、ロードバランの実装において少なくとも 1 つの他の入口ノード 1 1 0 が BGP を通じてその入口ノード 1 1 0 をエッジルータ 1 0 4 に提供する役割を担う。例えば図 5 に示されるいくつかの実施形態においては、所与の入口ノード 1 1 0 の左右の近傍入口ノード 1 1 0、例えばノード 1 1 0 の番号付きリストにおける左右の近傍ノードや、例えばノード 1 1 0 により形成された一貫したハッシュリングが、所与の入口ノード 1 1 0 をエッジルータ 1 0 4 に提供してもよい。例えば図 5 では、入口ノード 1 1 0 A が入口ノード 1 1 0 B および 1 1 0 D を提供し、入口ノード 1 1 0 B が入口ノード 1 1 0 A および 1 1 0 C を提供し、入口ノード 1 1 0 C が入口ノード 1 1 0 B および 1 1 0 D を提供し、そして入口ノード 1 1 0 D が入口ノード 1 1 0 C および 1 1 0 A を提供する。入口ノード 1 1 0 は本明細書にて後述するように、互いのヘルス状態をチェックしゴシップする。記載のヘルスチェック方法を用いて、異常なノードを検知することができる。1 秒以内、例えば 1 0 0 ミリ秒（ms）以内にノード 1 1 0 間に情報を伝播することができる。入口ノード 1 1 0 が正常でないとして決定された際には、異常なノードを提供する入口ノード 1 1 0 は直ちに異常なノード 1 1 0 の提供を停止してもよい。少なくともいくつか

40

50

つかの実施形態において、入口ノード110がTCPクローズまたはBGPセッションのための類似のメッセージをエッジルータ104に送信することによりエッジルータ104とのBGPセッションを終了する。したがって、故障したノード110により確立されたBGPセッションがノード110の故障の検知についてタイムアウトするのを待つ必要なく、故障したノード110の代わりに提供を行う他の入口ノード110が、ノード110の異常の検知時にノード110を提供するエッジルータ104とのBGPセッションを遮断する時に、エッジルータ104が故障したノード110を発見することができる。ロードバランサノードの故障の処理については図18Aおよび18Bに関連して本明細書にてさらに後述される。

【0038】

10

図6は分散型ロードバランサシステムの少なくともいくつかの実施形態による、マルチパスルーティング方法のフローチャートである。900で示すように、ロードバランサの実装における入口ノード110がその近傍ノード110をエッジルータ104に提供する。少なくともいくつかの実施形態において、入口ノード110は一貫したハッシュリングのようなノード110の番号付きリストに従って、その近傍ノード110を決定してもよい。少なくともいくつかの実施形態において、入口ノード110はBGPセッションを用いてその1つまたは複数の近傍ノード110をエッジルータ104に提供する。それらBGPセッションのうち1つは、提供されるノード110の各々のために確立されたエッジルータ104へのBGPセッションである。

【0039】

20

902で示すように、エッジルータ104はフローごとにハッシュ化されたマルチパスルーティング技術、例えば等価コストマルチパス(ECMP)ルーティング技術に従って、顧客160から受信したトラフィックをアクティブな(提供された)入口ノード110に分散させる。少なくともいくつかの実施形態において、エッジルータ104はパブリックIPアドレスを顧客160に公開し、すべての入口ノード110が同一のパブリックIPアドレスをエッジルータ104に提供する。エッジルータレイヤー4のソースおよび宛先ポートをエッジルータ104のフローハッシュの一部として用い、受信パケットを入口ノード110間に分散させる。一般的に、これによって各接続のためのパケットを同一の入口ノード110に分散させる。

【0040】

30

902で示すように、入口ノードがデータフローを対象のサーバノード130に転送する。少なくともいくつかの実施形態において、入口ノード110はデータフローのための1次および2次フロートラッカーノードと対話し、データフローを対象のサーバノード130へとマッピングする。こうして各入口ノード110は、受信されたパケットを対象のサーバノード130へと適切に転送するノード110を通じ、アクティブなデータフローのマッピングを維持してもよい。

【0041】

要素906~910は入口ノード110の故障の検知およびそこからの回復に関する。906で示すように、例えば本明細書に記載のヘルスチェック技術に従って入口ノード110は入口ノード110のダウンを検知してもよい。ノード110のダウンの検知時には、近傍ノード110がそのノード110のエッジルータ104への提供を停止する。少なくともいくつかの実施形態において、これにはそれぞれのBGPセッションのためのエッジルータ104へのTCPクローズの送信が関与する。

40

【0042】

908で示すように、入口ノード110のダウンをBGPセッションの終了を通じて検知する際に、エッジルータ104はフローごとにハッシュ化されたマルチパスルーティング技術に従い、顧客160の受信トラフィックを残りの入口ノード110に再分散させる。したがって、少なくともいくつかのデータフローを入口ノード110にルーティングしてもよい。

【0043】

50

910で示すように、入口ノード110は必要に応じてマッピングを回復し、データフローを適切な対象のサーバノードに転送してもよい。入口ノード110におけるノード110の故障からの回復方法については、本明細書の別の部分でも論じる。1つの実施例として、入口ノード110は、そのための現在のマッピングがないパケットの受信の際に、一貫したハッシュリングに従ってデータフローのためのフロートラッカーノードを決定し、フロートラッカーノードからマッピングを回復するために、一貫したハッシュ関数を用いてもよい。

非対称パケットフロー

【0044】

少なくともいくつかの実施形態において、アウトバウンドトラフィックとインバウンドデータとの比が1より大きい場合に入口ノードの帯域およびCPUの使用率を効率的に利用するために、分散型ロードバランスシステムは図7で示すようにアウトバウンドパケットをサーバノード130から複数の出口ノードへと転送する。少なくともいくつかの実施形態において、各接続のために、それぞれのサーバノード130上のロードバランサモジュール132が顧客エンドポイント/パブリックエンドポイントのタプルをハッシュし、それぞれのアウトバウンドパケットフローのための出口サーバ114の役割を担うロードバランサノード110を選択するために、一貫したハッシュアルゴリズムを用いる。しかし、いくつかの実施形態においては、接続のために出口サーバ114を選択するために、他の方法および/またはデータが用いられてもよい。選択した出口サーバ114は必ずではないが、通常は接続のための入口サーバ112の役割を担うロードバランサノード110とは異なるロードバランサノード110である。少なくともいくつかの実施形態において、そのロードバランサノード110/出口サーバ114に故障がない限りは、特定の接続のためのすべてのアウトバウンドパケットは故障したパケットを回避するために、同一の出口サーバ114に転送される。

【0045】

少なくともいくつかの実施形態において、出口サーバ114の選択のためにサーバノード130によって用いられる方法およびデータは、1つまたは複数のエッジルータ104により行われる入口サーバ112の選択に用いられる方法およびデータとは異なってもよい。一般的には異なる方法およびデータの利用により、所与の接続のための出口ノードには、その接続のために入口ノードとして選択されたロードバランサノード110とは異なるロードバランサノード110が結果的に選択されてもよく、また複数のロードバランサノード110が、入口ノードの役割を担う単一のロードバランサノード110を通過する接続のための発信トラフィックを処理する出口ノードとして結果的に選択されてもよい。

【0046】

図7は少なくともいくつかの実施形態による、非対称パケットフローを図示する。外部ネットワーク150上の顧客160から入口サーバ112を通り、サーバノード130A、130B、130Cおよび130Dのそれぞれに至る接続が少なくとも1つは確立されている。少なくともいくつかの実施形態において、接続のための出口ノードを選択するため、各接続についてそれぞれのサーバノード130上のロードバランサモジュール132が、顧客エンドポイント/パブリックエンドポイントのタプルをハッシュし、それぞれのアウトバウンドパケットフローのための出口サーバ114の役割を担うロードバランサノード110を選択するために一貫したハッシュアルゴリズムを利用する。例えばサーバノード130Aが接続のために出口サーバ114Aを選択し、サーバノード130Bがある接続のために出口サーバ114Aを、また別の接続のために出口サーバ114Bを選択している。しかし、いくつかの実施形態においては他の方法および/またはデータが接続のための出口ノードの選択に用いられてもよい。

顧客接続を破棄しないロードバランサノードの故障からの回復

【0047】

どのサーバノード130が顧客のトラフィックを受信するべきが決定するためにロードバランサノード110は一貫したハッシュを用いることができるが、いくつかの接続は寿

10

20

30

40

50

命が長いために、この手法では、新規サーバノード130が一貫したハッシュメンバーシップに参加し、入口ロードバランサノード110の故障が続いて起こる場合に既存のフローを維持できない可能性がある。この場合、サーバ130のための一貫したハッシュリングが異なるメンバーシップを有することになるため、故障したノード110からのフローを引き継ぐロードバランサノード110はもともと選択したマッピングを決定することができない可能性がある。したがって少なくともいくつかの実施形態においては、接続のためのサーバノード130を選択し、選択したサーバノード130にパケットをルーティングするために、ロードバランサノード110が分散型ハッシュテーブル(DHT)技術を用いてもよい。DHTに従って、サーバノード130が特定の接続を受信するために選択された場合に、サーバノード130が正常を保ち、サーバノード130上のロードバランサモジュール132がそのアクティブな接続の状態を(例、接続公開を介して)DHTに定期的に伝達することでリースを継続的に延長すると仮定すると、DHTは接続が完了するまでマッピングを保持する。入口ノード110の故障はエッジルータ104から残りのロードバランサノード110へのパケットの分散に影響を及ぼし、その結果、ロードバランサノード110は異なる顧客接続の組からトラフィックを受信することになる。しかしDHTはすべてのアクティブな接続をトラッキングするので、アクティブなマッピングのいずれかのリースを取得するためにロードバランサノード110がDHTに問い合わせを行うことは可能である。その結果、すべてのロードバランサノード110がトラフィックを正しいサーバノード130へと渡し、それによって入口ロードバランサノード110の故障時であってもアクティブな顧客接続の故障を防ぐ。

10

20

分散型ロードバランサシステムにおけるパケットフロー

【0048】

図8は少なくともいくつかの実施形態による、分散型ロードバランサシステムにおけるパケットフローを示す。図8において実線矢印はTCPデータパケットを表し、破線矢印はUDPデータパケットを表すことに留意する。図8では、入口サーバ112が1つまたは複数の顧客160からエッジルータ104を通じてTCPデータパケットを受信する。TCPパケットの受信時に、入口サーバ112はTCPパケットフローのためのサーバノード130へのマッピングを有しているかどうかを判断する。入口サーバ112がTCPパケットフローのためのマッピングを有する場合には、サーバ112はその後TCPパケットを(例えばUDPに従って)カプセル化し、カプセル化されたパケットを対象のサーバノード130へと送信する。入口サーバ112がTCPパケットフローのためのマッピングを有していない場合には、入口サーバ112はその後、サーバノード130への接続を確立するために、さらに/またはTCPパケットフローのためのマッピングを取得するために、TCPパケットから抽出されたTCPパケットフローに関する情報を含むUDPメッセージを1次フロートラッカー116Aへと送信してもよい。図9Aおよび9Bならびに図10A~10Gは、顧客160とサーバノード130との間の接続の確立方法を示す。サーバノード130上のロードバランサモジュール132は、サーバノード130上の1つまたは複数のTCP接続のための1つまたは複数の出口サーバ114の役割を担う1つまたは複数のロードバランサノード110を無作為に選択し、1つまたは複数の出口サーバ114を通じてUDPによりカプセル化されたTCP応答データパケットを1つまたは複数の顧客160へと送信する。

30

40

【0049】

図9Aおよび9Bは少なくともいくつかの実施形態による、分散型ロードバランサシステムにおいて接続を確立する際のパケットフローのフローチャートを提示する。図9Aの200で示すように、入口サーバ112は顧客160のTCPパケットを、エッジルータ104を通じて受信する。202で入口サーバ112がTCPフローのためのサーバノード130へのマッピングを有する場合、204で示すように、入口サーバ112がその後TCPパケットをカプセル化し、それぞれのサーバノード130へと送信する。入口サーバ112は継続的に、2つ以上の顧客160から2つ以上のTCPフローのためのパケットを受信し処理してもよいということに留意する。

50

【 0 0 5 0 】

202において入口サーバ112がTCPフローのためのマッピングを有していない場合、パケットは顧客160からのTCP同期(SYN)パケットであってもよい。206で示すように、SYNパケットの受信時には、入口サーバ112がSYNパケットからデータを抽出し、例えばUDPメッセージにおいて、データを1次フロートラッカー116Aへと転送する。少なくともいくつかの実施形態において、入口サーバ112はTCPフローのための1次フロートラッカー116Aおよび/または2次フロートラッカー116Bを一貫したハッシュ関数に従って決定することができる。208では、1次フロートラッカー116Aがデータを例えばハッシュテーブル内に格納し、TCP接続のサーバノード130側のための最初のTCPシーケンス番号を生成し、データおよびTCPシーケンス番号を2次フロートラッカー116Bへと転送する。210では、2次フロートラッカー116Bもデータを格納してもよく、少なくともTCPシーケンス番号を含むSYN/ACKパケットを生成して顧客160に送信する。

10

【 0 0 5 1 】

212で示すように、入口サーバ112は顧客160からTCP受信確認(ACK)パケットを、エッジルータ104を通じて受信する。入口サーバ112はこの時点ではTCPフローのためのサーバ130ノードへのマッピングは有していないため、214で入口サーバ112がACKパケットから抽出されたデータを含むメッセージを1次フロートラッカー116Aへと送信する。216で示すように、メッセージの受信時に1次フロートラッカー116Aは格納されたデータに従ってTCPフローを確認し、承認されたACKパケットからのシーケンス番号(+1)がSYN/ACKにおいて送信された数値に一致することを確認する。1次フロートラッカー116Aはその後、TCPフローを受信するサーバノード130を選択し、選択したサーバノード130上のローカルロードバランサモジュール132のデータ、TCPシーケンス番号およびIPアドレスを含むメッセージを2次フロートラッカー116Bに送信する。218で示すように、2次フロートラッカー116BもまたデータおよびTCPシーケンス番号を確認し、SYNメッセージを生成し、生成されたSYNメッセージを選択したサーバノード130上のローカルロードバランサモジュール132へと送信する。メソッドは図9Bの要素220に続く。

20

【 0 0 5 2 】

図9Bの220で示すように、ロードバランサモジュール132は生成されたSYNメッセージに応じて、サーバノード130の1つまたは複数の測定基準を調べてサーバノード130が接続を受け入れることが可能かどうかを決定してもよい。222においてサーバノード130が接続を現在受け入れることができないとロードバランサモジュール132が決定した場合、その後224においてロードバランサモジュール132が2次フロートラッカー116Bへの伝達を行う。2次フロートラッカー116Bは、すでに格納したフローのための情報を削除してもよい。226において、2次フロートラッカー116Bは1次フロートラッカー116Aへの伝達を行う。図9Aの216で示されるように1次フロートラッカー116Aはその後、新規対象のサーバノード130を選択し、2次フロートラッカー116Bへの伝達を行ってもよい。

30

【 0 0 5 3 】

222では、サーバノード130が接続を受け入れることが可能であるとロードバランサモジュール132が決定した場合、その後図9Bの228で示すように、ローカルロードバランサモジュール132は生成されたSYNからTCP SYNパケットを構成し、TCP SYNパケットをサーバノード130上のサーバ134に送信する。TCP SYNパケットのソースIPアドレスに顧客160の実際のIPアドレスが取り込まれ、それによりサーバ134は、顧客160への直接TCP接続を受信したことを知る。ロードバランサモジュール132はTCPフローに関連する詳細を、例えばローカルハッシュテーブルに格納する。230で示すように、サーバ134はロードバランサモジュール132が遮断するSYN/ACKパケットで応答する。232で示すように、ロードバランサモジュール132はその後接続情報を含むメッセージを2次フロートラッカー116Bに

40

50

送信し、接続が受け入れられたことを示す。このメッセージの受信時に、234において2次フロートラッカー116Bがサーバ134へのマッピングを記録し、類似のメッセージを1次フロートラッカー116Aに送信し、1次フロートラッカー116Aもまたマッピング情報を記録する。236で示すように、その後1次フロートラッカー116Aがマッピングメッセージを入口サーバ112に転送する。入口サーバ112はそうして顧客160からサーバ130へのTCPフローのためのマッピングを有する。

【0054】

238において入口サーバ112はサーバノード130上のロードバランサモジュール132へのデータフローのためのいずれかのバッファされたデータパケットをカプセル化し転送する。入口サーバ112により顧客160から受信されたデータのための追加の受信データパケットはカプセル化され、ロードバランサモジュール132へと直接転送され、ロードバランサモジュール132がデータパケットを脱カプセル化し、サーバ134へと送信する。

10

【0055】

240において、ロードバランサモジュール132はデータフローのための出口サーバ114を無作為に選択する。後続のサーバ134からのアウトバウンドTCPデータパケットはロードバランサモジュール132により遮断され、UDPに従ってカプセル化され、任意に選択された出口サーバ114へと転送される。出口サーバ114は発信パケットを脱カプセル化し、TCPデータパケットを顧客160に送信する。

【0056】

20

上記の通り202において、入口サーバ112が受信されたパケットのTCPフローのためのマッピングを有していない場合、パケットは顧客160からのTCP同期(SYN)パケットであってもよい。しかし、パケットはTCP SYNパケットでなくてもよい。例えばロードバランサノード110のメンバーシップがロードバランサノード110の追加や故障により変更される場合、エッジルータ104は入口サーバ112がマッピングを有しない入口サーバ112へのTCPフローのためのデータパケットのルーティングを開始してもよい。少なくともいくつかの実施形態において、入口サーバ112がマッピングを有しないこのようなパケットの受信時には、入口サーバ112は、一貫したハッシュリングに従ってTCPフローのための1次フロートラッカー116Aおよび/または2次フロートラッカー116Bを決定するために一貫したハッシュ関数を用い、マッピングを要求するために1次フロートラッカー116Aあるいは2次フロートラッカー116Bのいずれかへの伝達を行ってもよい。TCPフローのためのフロートラッカー116へのマッピング受信時には、入口サーバ112はマッピングを格納し、TCPフローのための1つまたは複数のTCPパケットのカプセル化および正しい宛先サーバノード130への転送を開始することができる。

30

ロードバランサノードの説明

【0057】

少なくともいくつかの実施形態において、おのおののロードバランサノード110は3つの役割を有する：

* 入口 - 顧客接続における顧客160からのすべての受信パケットの受信、マッピングが把握されている場合のデータパケットのサーバノード130へのルーティング、またはマッピングが把握されていない場合のフロートラッカーへの伝達。入口ノードからの発信パケットは(例、UDPに従って)入口ノードによりカプセル化される。

40

* フロートラッキング - 接続状態(例えば、各顧客接続を伝達するためにどのサーバノード130/サーバ134が割り当てられているか)の継続的な追跡。フロートラッカーも顧客160とサーバ134との間の接続の確立に参加する。

* 出口 - サーバ134から受信されたアウトバウンドパケットの脱カプセル化および顧客160への転送。

【0058】

少なくともいくつかの実施形態において、顧客 -> サーバのマッピングが把握されてい

50

る場合に、ロードバランサノード 1 1 0 は入口の役割として、サーバ 1 3 4 へのパケットの転送をし、マッピングが把握されていない場合に、要求のフロートラッカーへの転送の役割を担う。少なくともいくつかの実施形態において、特定の顧客接続/データフローのための入口ノードの役割を担うロードバランサノード 1 1 0 はまた、顧客接続のための 1 次フロートラッカーまたは 2 次フロートラッカーのいずれかの役割も担ってもよいが、その両方の役割を果たすことはできない。

【 0 0 5 9 】

少なくともいくつかの実施形態において、ロードバランサノード 1 1 0 はフロートラッカーの役割として、すでに確立された接続の顧客 - > サーバのマッピングの維持と同様にいまだ確立されつつある接続状態の維持の役割を担う。1 次フロートラッカーおよび 2 次フロートラッカーと呼ばれる 2 つのフロートラッカーは個別の各顧客接続に参与する。少なくともいくつかの実施形態において、顧客接続に関するフロートラッカーは一貫したハッシュアルゴリズムを用いて決定されてもよい。フロートラッカーはまた、各新規顧客接続のためのサーバノード 1 3 0 の擬似ランダム選択を含むがそれに限定されないロードバランサ機能も果たす。選択されたサーバノード 1 3 0 上のローカルロードバランサモジュール 1 3 2 は、サーバ 1 3 4 が接続を処理することができないと決定された場合に接続要求を拒否してもよいことに留意する。この場合、フロートラッカーはその後別のサーバノード 1 3 0 を選択し接続要求を他のサーバノード 1 3 0 に送信してもよい。少なくともいくつかの実施形態において、所与の接続のための 1 次フロートラッカーの役割および 2 次フロートラッカーの役割は異なるロードバランサノード 1 1 0 によって果たされる。

【 0 0 6 0 】

少なくともいくつかの実施形態において、ロードバランサノード 1 1 0 は出口の役割として、ステートレスであり、サーバノード 1 3 0 から受信された受信パケットを脱カプセル化し、いくつかの検証を行い、アウトバウンド TCP データパケットをそれぞれの顧客 1 6 0 に転送する。少なくともいくつかの実施形態において、サーバノード 1 3 0 上のロードバランサモジュール 1 3 2 は所与の接続のためのロードバランサノード 1 1 0 を任意に選択してもよい。

ロードバランサノード一貫したハッシュリングトポロジ

【 0 0 6 1 】

少なくともいくつかの実施形態において、ロードバランサノード 1 1 0 は入力キー空間（顧客エンドポイント、パブリックエンドポイント）の一貫したハッシュに基づくリングトポロジを形成する。入力キー空間は利用可能なフロートラッカーノードの間で区分されてもよく、すべてのフロートラッカーノードがそのキー空間に応じた質問への回答の役割を担ってもよい。少なくともいくつかの実施形態において、一貫したハッシュリングにおける後続処理に基づいて（例、2 次フロートラッカーノードが 1 次フロートラッカーノードへの後続ノードであるか、または一貫したハッシュリングにおいて次のノードである）、データは 1 次および 2 次フロートラッカーノードに複製されてもよい。フロートラッカーノードが何らかの理由によりダウンした場合、一貫したハッシュリングにおける次のロードバランサノードが故障したノードのキー空間を獲得する。新規フロートラッカーノードが参加する際には、ノードは（例、図 1 で示すように構成サービス 1 2 2 を用いて）そのエンドポイントを登録し、それにより他のロードバランサノードはロードバランサの実装における構成変化および、その結果の一貫したハッシュリングにおける構成変化を把握してもよい。一貫したハッシュリングにおけるフロートラッカーの追加および故障の処理については、図 1 1 A ~ 1 1 D に関連して詳細に後述される。

入口ノード < - > フロートラッカーノードの通信

【 0 0 6 2 】

少なくともいくつかの実施形態において、入口ノードの役割を担うロードバランサノード 1 1 0 は、構成サービス 1 2 2 からフロートラッカーノードの役割を担うロードバランサノード 1 1 0 について把握してもよい。入口ノードは、ロードバランサの実装において変更され、その結果の一貫したハッシュリングにおいても変更されたメンバーシップのた

10

20

30

40

50

めの構成サービス 1 2 2 を監視してもよい。入口ノードがマッピングを有しない顧客 1 6 0 からのパケットを受信する際には、どのフロートラッカーノードがパケットを伝達すべきか決定するために、入口ノードは一貫したハッシュ関数を用いてもよい。少なくともいくつかの実施形態において、ハッシュ関数への入力パケットからの（顧客エンドポイント、パブリックエンドポイントの）ペアである。少なくともいくつかの実施形態において、入口ノードおよびフロートラッカーノードは、UDP メッセージを用いた通信を行う。

【 0 0 6 3 】

1 次フロートラッカーノードが入口ノードから新規パケットフローのためのメッセージを受信する際には、1 次フロートラッカーノードが無作為に TCP シーケンス番号を決定し、別のメッセージを 2 次フロートラッカーノードに転送する。2 次フロートラッカーノードは顧客のための TCP SYN / ACK メッセージを生成する。両方のフロートラッカーが顧客接続のエンドポイントのペアおよび TCP シーケンス番号を記憶し、メモリ圧迫または期限により状態が消去されるまでこの情報を保持する。

10

【 0 0 6 4 】

TCP ACK パケットを受信された入口ノードから 1 次フロートラッカーノードがメッセージを受信する際には、1 次フロートラッカーノードが、承認された TCP シーケンス番号が SYN / ACK パケットにおいて送信済みの格納された数値と一致することを確認し、要求を伝達するサーバノード 1 3 0 を選択し、メッセージを 2 次フロートラッカーノードに転送する。2 次フロートラッカーノードは、サーバノード 1 3 0 で TCP スタックを用いた実際の TCP 接続を開始するためにメッセージを選択されたサーバノード 1 3 0 上のロードバランサモジュール 1 3 2 へと送信し、その後サーバノード 1 3 0 からの確認応答を待つ。

20

【 0 0 6 5 】

2 次フロートラッカーノードがサーバノード 1 3 0 上のロードバランサモジュール 1 3 2 から接続確認を受信する際に、1 次フロートラッカーを通して入口ノードに至るリバースメッセージフローが開始され、それにより両方のフローにおいてサーバノード 1 3 0 に関連する情報が格納される。この時点より、入口ノードにおいて受信された追加の TCP データパケットがサーバノード 1 3 0 上のロードバランサモジュール 1 3 2 に直接転送される。

ロードバランサモジュール < - > ロードバランサノードの通信

30

【 0 0 6 6 】

少なくともいくつかの実施形態において、すべてのロードバランサモジュール 1 3 2 が構成サービス 1 2 2 を用いてそのエンドポイントを登録し、ロードバランサノードレイヤー内のメンバーシップ変更のために構成サービス 1 2 2 を継続的に監視する。少なくともいくつかの実施形態による、ロードバランサモジュール 1 3 2 の関数が以下に記載される：

* 接続公開 - 定期的に（例、1 秒に 1 回）または非定期的に、それぞれのサーバノード 1 3 0 上のアクティブな接続（顧客エンドポイント、パブリックエンドポイント）の組を、それら接続のために最後にパケットをロードバランサモジュール 1 3 2 へと送信した入口ノードへと同様に、それら接続のための役割を担う 1 次および 2 次フロートラッカーノードの両方にも公開する。接続公開関数はそのための役割を担うロードバランサノード 1 1 0 において、接続状態のリースを更新する。

40

* ロードバランサレイヤーにおけるメンバーシップ変更の監視。メンバーシップが変更された際に、ロードバランサモジュール 1 3 2 はこの変更情報を用いて、接続のための役割を担うロードバランサノードへとアクティブな接続を直ちに送信してもよい。

分散型ロードバランサシステムにおけるパケットフロー - 詳細

【 0 0 6 7 】

分散型ロードバランサシステムは複数のロードバランサノード 1 1 0 を含んでもよい。少なくともいくつかの実施形態において、分散型ロードバランサシステムにおける各ロードバランサノード 1 1 0 はサーバ 1 3 4 への顧客 1 6 0 の接続のために、フロートラッカ

50

ーノード、出口ノード、および入口ノードの役割を担ってもよい。分散型ロードバランスシステムはまた、各サーバノード130上のロードバランサモジュール132を含んでもよい。

【0068】

図10A~10Gは、少なくともいくつかの実施形態による、分散型ロードバランスシステムにおけるパケットフローを示す。図10A~10Gにおいて、ロードバランサノード110の間で交換されるパケットと、ロードバランサノード110とサーバノード130との間で交換されるパケットは、UDPメッセージまたはUDPによりカプセル化された顧客TCPパケットのいずれかである。少なくともいくつかの実施形態において、顧客TCPパケットはネットワーク100に脱カプセル化された形式で、ロードバランサノード110の北側で境界ルータ102との間の相互の移行(図1参照)においてのみ存在する。図10A~10Gの実線矢印はTCPパケットを表し、破線矢印はUDPパケットを表すことに留意する。

10

【0069】

少なくともいくつかの実施形態において、単一のロードバランサノード110の故障時に、分散型ロードバランスシステムは確立された接続の保持を試みてもよい。少なくともいくつかの実施形態において、1次フロートラッカーノードまたは2次フロートラッカーノードのいずれかの故障時に、接続の顧客->サーバのマッピングが残りのフロートラッカーノードにより格納できるように、1次フロートラッカーノードおよび2次フロートラッカーノードにおける接続の詳細の複製を行うことで、これは達成されてもよい。少なくともいくつかの実施形態において、ノードの故障時にはパケットの一部が喪失される可能性があるが、顧客/サーバTCPパケットの再移行により、喪失されたパケットが回復する可能性もある。

20

【0070】

顧客からの各TCP接続は、TCPフローとよばれてもよく、顧客のIPアドレス、顧客用ポート、サーバ(パブリック)IPアドレス、およびサーバポートから成る4タプルによって一意に識別される。この識別子は、顧客およびパブリックエンドポイントのペアを示すCPまたはCcpとして略されてもよい。与えられた任意のTCPフロー(またはCPペア)に関連するパケットは、上流のエッジルータ104からのハッシュ化された等価コストマルチパス(ECMP)フロー分散により、入口サーバ112として動作するどのロードバランサノード110上にも表れることができる。しかし一般的にTCPフローのためのパケットはTCPフローの方向を変えるリンクやロードバランサノード110の故障がない限りは、同一のロードバランサノード110に到着し続けてもよい。上流のルータ104からTCPフローのためのパケットを受信するロードバランサノード110はTCPフローのための入口ノードと称される。

30

【0071】

少なくともいくつかの実施形態において、パケットがTCPフローのための入口ノードの役割を担うロードバランサノード110に到着した際に、どのロードバランサノード110がTCPフローのための状態を含むか(すなわち、フロートラッカーノード)を入口ノードが決定できるように、一貫したハッシュが用いられる。TCPフローに関する状態の維持の役割をどのロードバランサノード110が担うか決定するために、CPペアは入口ノードによって一貫したハッシュリングへとハッシュされてもよい。このノードは、TCPフローのための1次フロートラッカーの役割を担う。一貫したハッシュリングにおける後続ノードは、TCPフローのための2次フロートラッカーの役割を担う。

40

【0072】

少なくともいくつかの実施形態において、すべてのロードバランサノード110は、入口ノード、1次フロートラッカーノード、および2次フロートラッカーノードとしての役割を担ってもよい。TCPフローのための一貫したハッシュの結果に応じて、TCPフローのための入口ノードの役割を担うロードバランサノード110はまた、TCPフローのための1次または2次フロートラッカーノードの役割を担ってもよい。しかし、少なくと

50

もいくつかの実施形態において、異なる物理ロードバランサノード 110 が TCP フローのための 1 次および 2 次フロートラッカーの役割を果たす。

接続の確立

【0073】

図 10A を参照すると、顧客 160 からの新規接続は顧客 TCP 同期 (SYN) パケットによって起こってもよい。ロードバランサノード 110 は実際には SYN パケットの受信時にサーバノード 130 との接続を確立せず、また、接続を受信するためのサーバノード 130 の選択も直ちには行わない。代わりに、ロードバランサノード 110 は顧客の SYN パケットからの関連データを格納し、まだ未選択のサーバノード 130 のために SYN / ACK パケットを生成する。図 10C を参照すると、顧客 160 が TCP のスリーウェイハンドシェイクにおける第 1 の ACK パケットで応答すると、ロードバランサノード 110 はサーバノード 130 を選択し、サーバノード 130 のための同等の SYN パケットを生成し、サーバノード 130 を用いて実際の TCP 接続の確立を試みる。

10

【0074】

図 10A を再度参照すると、TCP フローのための入口サーバ 112 の役割を担うロードバランサノード 110 における顧客 SYN パケットの受信時には、入口サーバ 112 が SYN パケットからデータフィールドを抽出し、TCP フローのための 1 次フロートラッカー 116A にデータを転送する。1 次フロートラッカー 116A は例えばハッシュテーブル内にデータを格納し、最初の TCP シーケンス番号 (TCP 接続のサーバ側のための) を生成し、同一のデータを次フロートラッカー 116B へ転送する。2 次フロートラッカー 116B は、サーバの TCP シーケンス番号を含む顧客 160 のための SYN / ACK パケットを生成する。

20

【0075】

図 10A では、入口サーバ 112、1 次フロートラッカー 116A、2 次フロートラッカー 116B の役割が異なるロードバランサノード 110 により果たされる。しかし場合によっては、TCP フローのための入口サーバ 112 の役割を担うロードバランサノード 110 は、TCP フローのための 1 次フロートラッカー 116A または 2 次フロートラッカー 116B いずれか (しかし両方ではない) の役割を担う同一のノード 110 であってもよい。パケットフローのための入口サーバ 112 がフローのためのフロートラッカー 116 と同一のノード 110 上であってもよい理由は、パケットフローのためのフロートラッカー 116 がパケットフローのアドレス情報に適用される一貫したハッシュ関数に従って一貫したハッシュリング上で決定される一方で、エッジルータ 104 がフローごとにハッシュ化されたマルチパスルーティング技術 (例、ECMP ルーティング技術) に従ってフローのための入口サーバ 112 を擬似ランダムに選択するからである。パケットフローのための入口サーバ 112 は、パケットフローのためのフロートラッカー 116 と同一のノード 110 上にある場合、SYN パケットからのデータは、入口サーバ 112 を実装するノード 110 から他のフロートラッカー 116 ノード 110 へのみ転送されてもよい。例えば図 10B において、1 次フロートラッカー 116A は TCP フローのための入口サーバ 112 と同一のロードバランサノード 110A 上にあるが 2 次フロートラッカー 116B は異なるロードバランサノード 110B 上にあり、したがって SYN パケットからのデータがノード 110A から (フロートラッカー 116A によって) ロードバランサノード 110B 上の 2 次フロートラッカー 116B へと転送される

30

40

【0076】

図 10C を参照すると、非 SYN パケットが入口サーバ 112 に到着した際、入口サーバ 112 はどのサーバノード 130 がパケットの転送先となるべきかを把握しているか把握していないかのいずれかである。TCP フローのための入口サーバ 112 に到着する第 1 の非 SYN パケットは、TCP 確認番号フィールドが図 10A において SYN / ACK パケットが送信されたサーバシーケンス番号 (+1) に一致する TCP スリーウェイハンドシェイクにおける第 1 の TCP 受信確認 (ACK) パケット (または後続のデータパケットである可能性もある) であるべきである。入口サーバ 112 がサーバマッピングを有

50

しない非SYNパケットを受信した際には、入口サーバ112はシーケンス番号といったACKパケットからの情報を含む、またはその代わりにACKパケット自体を含むメッセージをTCPフローのための1次フロートラッカー116Aへと転送する。少なくともいくつかの場合に、1次フロートラッカー116AはTCPフローのために格納されたデータを記憶し、承認されたシーケンス番号(+1)がSYN/ACKパケットにおいて顧客160へ送信された数値に一致することを確認する。1次フロートラッカーはその後TCPフローのためのサーバノード130を選択し、TCPフローのためにすでに格納されたデータ、サーバシーケンス番号、および選択されたサーバノード130上のロードバランサモジュール132のためのIPアドレスを含む別のメッセージを2次フロートラッカー116Bに転送する。2次フロートラッカー116Bはサーバシーケンス番号を確認し、情報を記録し、生成されたSYNメッセージを選択されたサーバノード130上のロードバランサモジュール132に送信する。TCPフローのCPエンドポイントペアはここでロードバランサモジュール132/サーバノード130へとマッピングされる。サーバノード130上のロードバランサモジュール132は、2次フロートラッカー116Bから生成されたSYNメッセージを受信する際に、サーバノード130上のサーバ134のための正しいTCP SYNパケットを作成する役割を担う。SYNパケットの増加時には、ソースIPアドレスに顧客160の実際のIPアドレスが取り込まれ、それによってサーバ134は顧客160から直接TCP接続要求を受信したことを把握する。ロードバランサモジュール132はTCPフローに関連する詳細を、例えばローカルハッシュテーブル内に格納し、TCP SYNパケットをサーバ134に送信する(例、SYNパケットをサーバ134のLinuxカーネルに注入する)。

10

20

【0077】

図10Cにおいて、入口サーバ112、1次フロートラッカー116A、および2次フロートラッカー116Bの役割はそれぞれ異なるロードバランサノード110により果たされる。しかし場合によっては、TCPフローのための入口サーバ112の役割を担うロードバランサノード110は、TCPフローのための1次フロートラッカー116Aまたは2次フロートラッカー116Bの役割を担うノードと同一のノード110である(しかし両方ではない)。例えば、図10Dにおいて、1次フロートラッカー116Aが異なるロードバランサノード110B上にある一方で、2次フロートラッカー116BはTCPフローのための入口サーバ112の役割を担う同一のロードバランサノード110A上にある。

30

【0078】

図10Eを参照すると、サーバ134(例、Linuxカーネル)はロードバランサモジュール132も遮断するSYN/ACKパケットで応答する。SYN/ACKパケットは、もともとSYN/ACKにおいて2次フロートラッカー116Bから顧客160へと伝達されたもの(図10Aを参照)とは異なるTCPシーケンス番号を含んでもよい。ロードバランサモジュール132はシーケンス番号デルタを受信および発信データパケットへと適用させる役割を担う。サーバ134からのSYN/ACKパケットはまた、ロードバランサモジュール132から2次フロートラッカー116Bへと戻るメッセージ(例、UDPメッセージ)のきっかけとなり、選択したサーバノード130/ロードバランサモジュール132/サーバ134への接続が成功したことを示す。このメッセージの受信時に2次フロートラッカー116Aは、顧客160とサーバ134との間の顧客およびパブリックエンドポイントのペア(CP)マッピングを送信されたとおりに記録し、同様にCPマッピングを記録する1次フロートラッカー116Aへと類似のメッセージを送信してもよい。1次フロートラッカー116Aはその後入口サーバ112へとCPマッピングメッセージを転送してもよく、それにより入口サーバ112は接続のためにあらゆるバッファされたデータパケットをカプセル化されたデータパケットとしてローカルサーバノード130上のロードバランサモジュール132へと転送できるようになる。

40

【0079】

図10Fを参照すると、接続のためのCPマッピングは入口サーバによって把握され、

50

したがって接続のための入口サーバ112により受信された受信TCPパケットは、(例、UDPに従って)カプセル化され、カプセル化されたデータパケットとしてサーバノード130上のローカルロードバランサモジュール132に直接転送されてもよい。ロードバランサモジュール132はデータパケットを脱カプセル化し、例えばカーネルのTCPスタック上へとTCPパケットを注入することでTCPパケットをサーバノード130上のサーバ134へと送信する。サーバ134からのアウトバウンドパケットは、サーバノード130上のロードバランサモジュール132によって遮断され、(例、UDPに従って)カプセル化され、そしてロードバランサモジュール132がこの接続のための出口サーバ114として無作為に選択する任意のロードバランサノード110へと転送される。出口サーバ114はパケットを脱カプセル化し、脱カプセル化されたデータパケットを顧客116へと送信する。選択したロードバランサノード110の出口関数はステートレスであり、そのため出口サーバの役割を担うロードバランサノード110の故障時には、異なるロードバランサノード110を接続のための出口サーバ114としてとして選択することが可能である。しかし一般的には、接続の維持のための出口サーバ114と同一のロードバランサノード110が、アウトバウンドパケットの再配置を減少させ、または取り除くために用いられる。

【0080】

図10Gを参照すると、少なくともいくつかの実施形態において、1次フロートラッカー116Aによって選択されたサーバノード130A上のロードバランサモジュール132A(図10Cを参照)は、自身に負荷がかかり過ぎていると判断した場合に、2次フロートラッカー116Bから受信された、生成されたSYNメッセージ(図10Cを参照)を拒否する選択肢を有する。少なくともいくつかの実施形態において、生成されたSYNメッセージは拒否の最大値を許す生存時間(TTL)の数値またはカウンターを含む。少なくともいくつかの実施形態において、このTTL値がゼロに達した場合、ロードバランサモジュール132Aが負荷を制限するために接続の受け入れまたは接続の破棄のいずれかを行ってもよい。ロードバランサモジュール132Aは接続の拒否を決定した場合、TTL値をディクリメントし、2次フロートラッカー116Bに拒否メッセージを送信する。2次フロートラッカー116BはCPマッピングをリセットし、リリースメッセージを同様のことを行う1次フロートラッカー116Aに送信する。1次フロートラッカー116Aは別のサーバノード130B上の新規ロードバランサモジュール132Bを選択し、2次フロートラッカー116Bに新規対象のメッセージを返信し、2次フロートラッカー116Bが新規生成されたSYNメッセージを新規に選択されたロードバランサモジュール132Bに送信する。パケットの破棄によりこれらのシーケンスが完了しない可能性があることに留意する。しかし、顧客160からの再伝達によりロードバランサモジュールの選択処理が1次フロートラッカー116Aにおいて再び行われてもよい。1次フロートラッカー116Aは必ずそうするとは限らないが、生成されたSYNパケットの前回の拒否について把握していない場合、接続のための同一のロードバランサモジュール132を選択してもよい。

【0081】

少なくともいくつかの実施形態において、TTLカウンターは継続的に接続要求をサーバノード130に送るのを阻止してもよく、これは例えばすべてのサーバノード130がビジー状態である場合に発生してもよい。少なくともいくつかの実施形態において、ロードバランサモジュール132が、それぞれのサーバノード130に代わり接続要求を拒否する時は毎回、ロードバランサモジュール132がTTLカウンターをディクリメントする。フロートラッカーノード116は、TTLカウンターがゼロではない(すなわちある特定の閾値を超える)限りTTLカウンターを監視してもよく、別のサーバノード130を選択し再度試みてよい。TTLカウンターがゼロに達する(すなわちある特定の閾値を超える)場合、接続要求は破棄され、その接続のために選択されたサーバノード130のうちの1つに接続要求を送信する試みをフロートラッカーノード116が再び行うことはない。少なくともいくつかの実施形態において、エラーメッセージがそれぞれの顧客1

10

20

30

40

50

60に送信されてもよい。

【0082】

少なくともいくつかの実施形態において、分散型ロードバランサシステムは複数のパブリックIPアドレスをサポートする。こうして、顧客160が同一の顧客用ポート番号から2つの異なるパブリックIPアドレスへの2つのTCP接続を開始することが可能になる。これらのTCP接続は顧客160の観点とは異なるが、内部では分散型ロードバランサが同一のサーバノード130への接続をマッピングしてもよく、これにより衝突が起こる。少なくともいくつかの実施形態において、可能性のある衝突を検知し処理するために、ロードバランサモジュール132は図10Cおよび10Dで示すように2次フロートラッカー116Bから生成されたSYNパケットを受信する際に、アドレス情報をアクティブな接続と比較してもよく、この接続が衝突を発生させる場合には、図10Gで示すように接続要求を拒否してもよい。

10

ロードバランサノードの故障および追加の処理

【0083】

従来のロードバランサの多くにおいては、ロードバランサの故障時には既存の接続の一部またはすべてが喪失される。少なくともいくつかの実施形態において、単一のロードバランサノード110の故障時には、接続が完全に正常に戻るまで顧客およびサーバが接続を通じてパケットの交換を継続できるように、分散型ロードバランサシステムが確立された接続のうち少なくともいくつかを維持してもよい。また分散型ロードバランサシステムは、故障時に確立の処理中であった接続の伝達を継続してもよい。

20

【0084】

分散型ロードバランサシステムの少なくともいくつかの実施形態において、単一のロードバランサノード110の故障時に既存の顧客接続を回復させることができる故障回復プロトコルが実装されてもよい。しかし複数のロードバランサノード110の故障により、顧客接続が喪失される場合がある。少なくともいくつかの実施形態において、顧客160とサーバ134との間のTCPの再伝達がロードバランサノード110の故障後の回復手段として用いられてもよい。

【0085】

可能性のあるロードバランサノード110の故障に加えて、新規ロードバランサノード110が分散型ロードバランサシステムに追加されてもよい。これら新規ノード110はロードバランサレイヤーに、またそれにより一貫したハッシュリング追加されてもよく、ロードバランサノード110の既存の顧客接続に関する役割は、必要に応じて変更に従い調整される。

30

フロートラッカーノードの故障および追加の処理

【0086】

少なくともいくつかの実施形態において、各接続が確立される(例、図10A~10Gを参照)にしたがって、接続状態の情報が1次および2次フロートラッカーと呼ばれる2つのロードバランサノード110を通じて渡される。これらは例えば(顧客IP:ポート、パブリックIP:ポート)タプルをハッシュ関数入力として用いる一貫したハッシュアルゴリズムを利用して決定されてもよい。単一のロードバランサノード110の故障時には、パケットを接続のための選択されたサーバノード130へと導くために、少なくとも1つの生存するロードバランサノード110は一貫したハッシュ関数を通じて継続的にマッピングされてもよく、また接続のために必要な状態情報を含んでもよい。また、ロードバランサノード110を一貫したハッシュリングへと追加する場合、接続のための状態情報は適切なフロートラッカーへとリフレッシュされてもよい。

40

【0087】

図11A~11Dは少なくともいくつかの実施形態による、ロードバランサノードの一貫したハッシュリングにおいてメンバーシップに影響を与えるイベントの処理を示す。これらのイベントは、新規1次フロートラッカーノードの追加、新規2次フロートラッカーノードの追加、1次フロートラッカーノードの故障、および2次フロートラッカーノード

50

の故障を含んでもよいがこれらに限定されない。

【 0 0 8 8 】

図 1 1 A は、新規 1 次フロートラッカーノードの一貫したハッシュリングへの追加処理を示す。図 1 1 A の上部列は、1 つまたは複数の顧客接続のための 1 次フロートラッカーとしてのフロートラッカー 1 1 6 A および 1 つまたは複数の同一の接続のための 2 次フロートラッカーとしてのフロートラッカーノード 1 1 6 B を示す。図 1 1 A の下部列においては、新規フロートラッカーノード 1 1 6 C が追加され、1 つまたは複数の顧客接続のための 1 次フロートラッカーとなっている。以前は 1 次フロートラッカーであったフロートラッカーノード 1 1 6 A は 2 次フロートラッカーとなり、以前は 2 次フロートラッカーであったフロートラッカーノード 1 1 6 B は、一貫したハッシュリングにおける次のフロートラッカーとなる。フロートラッカー 1 1 6 A および 1 1 6 B により維持された 1 つまたは複数の顧客接続のための状態情報は新規 1 次フロートラッカー 1 1 6 C に提供されてもよい。また、フロートラッカー 1 1 6 B は 2 次フロートラッカーの役割として以前トラッキングしていた接続を「忘れて」もよい。

10

【 0 0 8 9 】

図 1 1 B は、新規 2 次フロートラッカーノードの一貫したハッシュリングへの追加処理を示す。図 1 1 B の上部列は、1 つまたは複数の顧客接続のための 1 次フロートラッカーとしてのフロートラッカー 1 1 6 A および 1 つまたは複数の同一の接続のための 2 次フロートラッカーとしてのフロートラッカーノード 1 1 6 B を示す。図 1 1 B の下部列においては、新規フロートラッカーノード 1 1 6 C が追加され、1 つまたは複数の顧客接続のための 2 次フロートラッカーとなっている。フロートラッカーノード 1 1 6 A は 1 つまたは複数の接続のための 1 次フロートラッカーのままであり、以前は 2 次フロートラッカーであったフロートラッカーノード 1 1 6 B は一貫したハッシュリングにおける次のフロートラッカーとなる。フロートラッカー 1 1 6 A および 1 1 6 B によって維持された 1 つまたは複数の顧客接続のための状態情報は新規 2 次フロートラッカー 1 1 6 C に提供されてもよい。また、フロートラッカー 1 1 6 B は 2 次フロートラッカーの役割として以前トラッキングしていた接続を「忘れて」もよい。

20

【 0 0 9 0 】

図 1 1 C は、一貫したハッシュリングにおける 1 次フロートラッカーノードの故障の処理を示す。図 1 1 C の上部列は、1 つまたは複数の顧客接続のための 1 次フロートラッカーとしてのフロートラッカー 1 1 6 A、1 つまたは複数の同一の接続のための 2 次フロートラッカーとしてのフロートラッカーノード 1 1 6 B、および一貫したハッシュリングにおける次のフロートラッカーとしてのフロートラッカーノード 1 1 6 C を示す。図 1 1 C の下部列においては、1 次フロートラッカーノード 1 1 6 A が故障している。フロートラッカーノード 1 1 6 B は 1 つまたは複数の接続のための 1 次フロートラッカーとなり、フロートラッカーノード 1 1 6 C は 1 つまたは複数の接続のための 2 次フロートラッカーとなる。1 つまたは複数の顧客接続のための状態情報はフロートラッカー 1 1 6 B によって維持され、新規 2 次フロートラッカー 1 1 6 C へと提供されてもよい。

30

【 0 0 9 1 】

図 1 1 D は一貫したハッシュリングにおける 2 次フロートラッカーノードの故障の処理を示す。図 1 1 D の上部列は、1 つまたは複数の顧客接続のための 1 次フロートラッカーとしてのフロートラッカー 1 1 6 A、1 つまたは複数の同一の接続のための 2 次フロートラッカーとしてのフロートラッカーノード 1 1 6 B、および一貫したハッシュリングにおける次のフロートラッカーとしてのフロートラッカーノード 1 1 6 C を示す。図 1 1 D の下部列においては、2 次フロートラッカーノード 1 1 6 B が故障している。フロートラッカーノード 1 1 6 A は 1 つまたは複数の接続のための 1 次フロートラッカーのままであり、フロートラッカーノード 1 1 6 C は 1 つまたは複数の接続のための 2 次フロートラッカーとなる。1 つまたは複数の顧客接続のための状態情報はフロートラッカー 1 1 6 B により維持され、新規 2 次フロートラッカー 1 1 6 C に提供されてもよい。

40

【 0 0 9 2 】

50

少なくともいくつかの実施形態において、サーバノード130上のロードバランサモジュール132はロードバランサノード110への接続公開を行う。少なくともいくつかの実施形態において、接続公開は定期的に(例、1秒に1回)または非定期的に現在の接続状態の情報をサーバノード130から、接続のための1次および2次フロートラッカーノード両方への接続マッピングのリフレッシュまたは復元を行うフロートラッカーノードおよび入口ノードの役割を担うロードバランサノード110へとプッシュする。少なくともいくつかの実施形態において、ロードバランサモジュール132は例えば図11A~11Dで示される通り、フロートラッカーのメンバーシップ変更を検知してもよい。それに従ってロードバランサモジュール132は、メンバーシップが変更された際に接続のために変更されたかもしれない、1次および2次フロートラッカーノードにおける接続のための状態情報を追加するために、接続公開を行ってもよい。接続公開により、複数のロードバランサノードの故障時に少なくともいくつかの確立された接続が回復されてもよいことに留意する。

10

故障に関連するメッセージフロー

【0093】

少なくともいくつかの実施形態において、1次および2次フロートラッカーノードの間のプロトコルは、修正または同期機能を含んでもよい。例えば図11Aを参照すると、新規1次フロートラッカーノード116Cが一貫したハッシュリングに参加する際には、新規ノード116Cがいくつかの数(1/N)の接続のための一貫したハッシュのキー空間に要求を出し、エッジルータ104からこれらの接続に関連するトラフィックの受信を開始してもよい。しかし新規1次フロートラッカーノード116Cは接続のために格納された状態を一切有さないため、各パケットに対して、顧客160から受信された第1のパケットとして動作してもよい。1次フロートラッカーはSYNデータパケットに応じてサーバのTCPシーケンス番号の生成(例、図10Aを参照)および顧客160からの第1のACKパケットに応じてサーバノード130の選択(例、図1を参照)の役割を担い、これらの生成された数値は前の1次フロートラッカー(図11Aにおけるフロートラッカーノード116A)により選択された数値と相違してもよい。しかし少なくともいくつかの実施形態において、一貫したハッシュアルゴリズムは前の1次フロートラッカー(図11Aにおけるフロートラッカーノード116A)を2次フロートラッカーの役割に割り当て、このフロートラッカーは接続のためのすでに格納された状態をいまだに保持する。したがって少なくともいくつかの実施形態において2次フロートラッカー(図11Aにおけるフロートラッカーノード116A)は、1次フロートラッカー116Cから受信された情報における不一致を検知した時、2次フロートラッカーは接続のためのフロートラッカーとしての役割を担う2つのロードバランサノード110を同期するために更新メッセージを1次フロートラッカー116Cへと返信することができる。一貫したハッシュリングのメンバーシップにおいて他の変更がなされた後にフロートラッカーを同期するために類似の方法が用いられてもよい。

20

30

ロードバランサモジュールの説明

【0094】

少なくともいくつかの実施形態において、ロードバランサモジュール132は各サーバノード130上にある分散型ロードバランサシステムの構成要素である。ロードバランサノード132の役割は、ロードバランサノード110から受信されたパケットの脱カプセル化および脱カプセル化されたパケットのサーバノード130上のサーバ134への送信、ならびにサーバ134からの発信パケットのカプセル化およびカプセル化されたパケットのロードバランサノード110への送信を含むがそれに限定されない。

40

【0095】

少なくともいくつかの実施形態において、入口サーバ112の役割を担うロードバランサノード110からサーバノード130上のロードバランサモジュール132への受信パケットは、実際の顧客データパケットをカプセル化するステートレスプロトコル(例、UDP)パケットである。カプセル化された顧客データパケットは各々、ソースアドレスと

50

してそれぞれの顧客 160 のオリジナルの顧客 IP : ポートを、そして宛先アドレスとしてサーバ 134 パブリック IP : ポートを有する。ロードバランサモジュール 132 は顧客データパケットを脱カプセル化し、例えばパケットの方向をローカルホスト TCP フローへと変更することで、それぞれのサーバノード 130 上のサーバ 134 へと送信する。

【0096】

少なくともいくつかの実施形態において、サーバ 134 から出口サーバ 114 の役割を担うロードバランサノード 110 への発信パケットは、発信 IP パケットをカプセル化するステートレスプロトコル (例、UDP) パケットである。ロードバランサモジュール 132 は発信 IP パケットをカプセル化し、ファブリック 120 を通してカプセル化されたパケットを出口サーバ 114 へと送信する。各カプセル化された発信 IP パケットは、ソースアドレスとしてサーバ 134 パブリック IP : ポートを、そして宛先アドレスとしてそれぞれの顧客 160 の顧客 IP : ポートを有する。

ロードバランサモジュール機能

【0097】

少なくともいくつかの実施形態において、サーバノード 130 上のロードバランサモジュール 132 の機能は以下の 1 つまたは複数を含んでもよいが、それらに限定されない：

- * 1 つまたは複数のロードバランサノード 110 からの、例えば顧客 160 への接続を処理する入口サーバ 112 からの、UDP トンネルの終了。これは入口サーバ 112 から受信された受信顧客データパケットの UDP 脱カプセル化を含む。

- * 接続のための発信トラフィックを受信する出口サーバ 114 の選択。

- * それぞれのサーバ 134 上の接続における発信 IP パケットの遮断、接続のための発信 IP パケットのカプセル化、およびカプセル化されたデータパケットの出口サーバ 114 への送信。

- * フロートラッカーノード 116 が顧客 160 に SYN / ACK を送信する際にシーケンス番号がフロートラッカーノード 116 により生成されたシーケンス番号と整列させるための、受信および発信パケットにおけるシーケンス番号のマングリング。

- * 例えばそれぞれのサーバ 134 の現在の負荷を示す 1 つまたは複数の測定基準に基づく、それぞれのサーバ 134 のための接続を受け入れるか拒否するかの決定。

- * 顧客 IP : ポートアドレスの衝突を回避するためのアクティブな接続がある場合の、同一の顧客 IP : ポートアドレスからそれぞれのサーバ 134 への接続の検知および拒否。

- * 接続トラッキングおよび接続公開。

ロードバランサモジュールの構成情報

【0098】

少なくともいくつかの実施形態において、各ロードバランサモジュール 132 は構成のために、以下の情報の組の 1 つまたは複数を獲得し、ローカルに格納してもよいが、それらに限定されない：ロードバランサノード 110 エンドポイントの組、伝達する有効なパブリック IP アドレスの組、およびそれぞれのサーバ 134 が受信接続を受け入れる 1 つまたは複数のポート番号。少なくともいくつかの実施形態において、図 1 で示すようにこの情報は、分散型ロードバランサシステムの構成要素である構成サービス 122 から獲得されるか、またはそれへのアクセスが問い合わせにより更新されてもよい。いくつかの実施形態においては、他の情報獲得方法が用いられてもよい。

ロードバランサモジュールのパケット処理

【0099】

少なくともいくつかの実施形態による、インバウンドトラフィックおよびアウトバウンドトラフィックのためのロードバランサモジュール 132 の動作を以下に記載する。少なくともいくつかの実施形態において、インバウンドデータパケットがロードバランサモジュール 132 により受信される際に、データパケットが UDP パケットから脱カプセル化され、脱カプセル化された TCP パケットにおける宛先アドレスは構成された有効なパブリック IP アドレスの組に対して最初に検証される。一致がない場合、パケットは破棄さ

れるかまたは無視される。少なくともいくつかの実施形態において、シーケンス番号がSYN/ACKパケットを顧客160に送信したフロートラッカーノード116により生成され無作為に選択されたシーケンス番号に一致するように、ロードバランサモジュール132はTCPヘッダにおけるシーケンス番号を定数デルタにより調整してもよい。ロードバランサモジュール132は[顧客:パブリック]エンドポイントから[顧客/サーバ]エンドポイントへのマッピングを内部状態として記録する。

【0100】

少なくともいくつかの実施形態において、サーバ134からのアウトバウンドTCPデータパケットのために、ロードバランサモジュール132はまず内部状態を確認し、パケットが、ロードバランサモジュールが管理しているアクティブな接続のためのものであるかどうかを決定する。そうでない場合、ロードバランサモジュール132はただパケットを渡す。そうである場合、ロードバランサモジュール132は発信TCPパケットを例えばUDPに従ってカプセル化し、カプセル化されたパケットをこの接続のための出口サーバ114として選択されたロードバランサノード110へと転送する。少なくともいくつかの実施形態において、ロードバランサモジュール134は発信TCPパケットにおけるTCPシーケンス番号を定数デルタにより調整して、SYN/ACKパケットを顧客160に送信したフロートラッカーノード116により生成されたシーケンス番号を整列させてもよい。

接続のトラッキング

【0101】

少なくともいくつかの実施形態において、各サーバノード130上のロードバランサモジュール132はそれぞれのサーバ134へのすべてのアクティブな顧客接続のための接続の詳細を含むハッシュテーブルを管理する。少なくともいくつかの実施形態において、ハッシュテーブルのためのキーは(顧客IP:ポート、パブリックIP:ポート)タプルである。少なくともいくつかの実施形態において、各顧客接続のための接続状態は以下の1つまたは複数を含むが、それらに限定されない:

- * 顧客IP:ポート
- * パブリックIP:ポート
- * フロートラッカー116ノードにより提供される最初のサーバのTCPシーケンス番号。
- * サーバのTCPシーケンス番号デルタ。
- * オリジナルの1次フロートラッカーIPアドレス。
- * オリジナルの2次フロートラッカーIPアドレス。
- * 最後に検知された入口サーバ112のIPアドレス。
- * このエントリのための有効期限
- * 最長期間未使用の(LRU)/衝突指数。

【0102】

少なくともいくつかの実施形態において、各ロードバランサモジュール132はすべてのアクティブな顧客接続のための1次および2次フロートラッカーノードへの接続公開メッセージを定期的に生成する。少なくともいくつかの実施形態において、/proc/net/tcpの内容がスキャンされロードバランサモジュールのハッシュテーブルにおけるアクティブな接続と交差し、Linuxカーネルが接続のトラッキングを停止するまでフロートラッカーノードへと継続的に公開される。接続公開については本明細書にて詳細に後述される。

シーケンス番号のマングリング

【0103】

上述のように少なくともいくつかの実施形態において、ロードバランサノード110はサーバ134の代わりに顧客160SYNパケットに応じて、SYN/ACKパケットを生成する。顧客160がACKパケットを送信する後のみ(TCPスリーウェイハンドシェイク)ロードバランサモジュール110がサーバノード130上のロードバランサモ

10

20

30

40

50

ジュール132へといずれかのデータを送信する。ロードバランサモジュール132が最初に顧客接続を確立するよう指示される際は、ロードバランサモジュール132がローカルでSYNパケットを作成してサーバノード130上のサーバ134を用いてTCP接続を開始し、サーバ134に対応するSYN/ACKパケットを遮断する。通常、サーバ134(例、サーバノード130上のLinuxカーネル)がSYN/ACKパケットにおいてロードバランサノード110から受信された顧客の1つとはまったく異なるTCPシーケンス番号を選択する。こうして少なくともいくつかの実施形態において、ロードバランサモジュール132は顧客160とサーバ134との間のTCP接続のすべてのパケットにおけるシーケンス番号の補正を行ってもよい。少なくともいくつかの実施形態において、ロードバランサモジュール132はロードバランサノード110により生成されたシーケンス番号とサーバ134により作成されたシーケンス番号との間の差異を計算し、その差異をデルタ値としてTCP接続のためのハッシュテーブルエントリ内に格納する。受信データパケットが顧客160から接続に到着する際には、TCPヘッダがサーバ134により用いられるシーケンス番号と整列しない確認番号を含むため、ロードバランサモジュール132はTCPヘッダにおけるシーケンス番号の数値からデルタ値を減算する(例、2つの補数を用いて)。ロードバランサモジュールはまた、サーバ134から顧客130への接続上のアウトバウンドデータパケットにおけるシーケンス番号にデルタ値を追加する。

10

分散型ロードバランサシステムにおけるヘルスチェック

【0104】

20

分散型ロードバランサシステムの少なくともいくつかの実施形態において、各ロードバランサノード110はロードバランサの実装における正常なメンバー(すなわち、正常なロードバランサノード110およびサーバノード130)の一貫した見解を、少なくとも以下の理由により要求する:

- * ロードバランサ - ロードバランサノード110がサーバノード130の故障を検知し、顧客のトラフィックを受け入れることができる正常なサーバノード130の組において収束する必要がある。

- * 分散状態の管理 - ロードバランサは複数のロードバランサノード110で共有された/複製された状態を有する分散型システムである(例、一貫したハッシュ機構に従って)。顧客のトラフィックを正しく処理するために、各ロードバランサノード110は最終的にロードバランサの実装における正常なメンバーの一貫した見解を有する必要がある。

30

【0105】

これを達成するため、分散型ロードバランサシステムの少なくともいくつかの実施形態は、ロードバランサの実装においてノードを監視し、可能な限り迅速に異常なノードを検知するヘルスチェックプロトコルの実施形態を実装してもよい。ヘルスチェックプロトコルはロードバランサの実装においてノード間にヘルス情報を伝播してもよく、正常なノードの組においてノードの収束を可能にする方法を提供してもよい。またヘルスチェックプロトコルは、ロードバランサの実装における正常/異常なノードおよび状態の変化を報告するための機構を提供してもよい。

【0106】

40

少なくともいくつかの実施形態において、ヘルスチェックプロトコルは以下の仮定のうちの1つまたは複数に基づいてもよいが、それらに限定されない:

- * ロードバランサの実装におけるすべてのノードが把握される(すなわち、ヘルスチェックプロトコルは発見を行わなくてもよい)。

- * ノードの故障はすべてフェイルストップである。

- * ノード間のすべてのメッセージはステートレスプロトコル(例、UDP)メッセージであり、メッセージは破棄され、遅延させられ、複製され、または破損する可能性がある。メッセージの伝達の保証はない。

【0107】

少なくともいくつかの実施形態において、ロードバランサの実装におけるノード(例

50

、ロードバランサノード 1 1 0 またはサーバノード 1 3 0) は以下の条件の下で正常であると見なされてもよい：

* ノードの内部構成要素はすべてレディ状態である (顧客のトラフィックを処理する準備が完了している)。

* ノードの受信 / 発信ネットワークリンクは正常である (少なくともどの顧客のトラフィックを流すかについてのネットワークインタフェースコントローラ (NIC) に関しては)。

【 0 1 0 8 】

図 1 2 は少なくともいくつかの実施形態による、ヘルスチェック間隔に従って各ロードバランサノードにより実行されるヘルスチェック方法のハイレベルフローチャートである。1 0 0 0 で示すように、各ロードバランサ間隔において、例えば 1 0 0 ミリ秒毎に、各ロードバランサ (LB) ノード 1 1 0 は少なくとも 1 つの他の LB ノード 1 1 0 および少なくとも 1 つのサーバノード 1 3 0 のヘルスチェックを行ってもよい。1 0 0 2 で示すように、ロードバランサノード 1 1 0 はヘルスチェックに従って、そのローカルに格納されたヘルス情報を更新してもよい。1 0 0 4 で示すように、ロードバランサノード 1 1 0 はその後、少なくとも 1 つの他のロードバランサノード 1 1 0 を無作為に選択し、そのヘルス情報を選択された 1 つまたは複数のロードバランサノード 1 1 0 へと送信してもよい。少なくともいくつかの実施形態において、ノード 1 1 0 はまた、正常なロードバランサノード 1 1 0 のリストを 1 つまたは複数のサーバノード 1 3 0、例えばノード 1 1 0 によりヘルスチェックされる 1 つまたは複数の同一のサーバノード 1 3 0 へと送信してもよい。図 1 2 の要素は以下において詳細に説明される。

【 0 1 0 9 】

ヘルスチェックプロトコルの少なくともいくつかの実施形態において、ロードバランサノード 1 1 0 はそのヘルス状態を他のロードバランサノード 1 1 0 にアサートしない。その代わりに、1 つまたは複数の他のロードバランサノード 1 1 0 がそのノード 1 1 0 のヘルスチェックを行ってもよい。例えば少なくともいくつかの実施形態において、各ロードバランサノード 1 1 0 はヘルスチェックを行う 1 つまたは複数の他のノード 1 1 0 を、定期的または非定期的は無作為に選択してもよい。別の実施例として、少なくともいくつかの実施形態において、1 つまたは複数の他のロードバランサノード 1 1 0、例えば一貫したハッシュリング等のノード 1 1 0 の番号付きリスト上の所与のロードバランサノード 1 1 0 の 2 つの最近傍ノードはそれぞれ、所与のノード 1 1 0 のヘルスチェックを定期的または非定期的に行ってもよい。少なくともいくつかの実施形態において、ノード 1 1 0 のヘルスチェックは図 2 3 で示すように、ノード 1 1 0 上の NIC 1 1 1 4 へと送信されたヘルス ping の利用を含んでもよい。少なくともいくつかの実施形態において、第 2 のノード 1 1 0 が正常であると第 1 のノード 1 1 0 がヘルスチェックを通じて決定する場合、第 1 のノード 1 1 0 は、ロードバランサノード 1 1 0 のためのローカルヘルス情報に格納された、第 2 のノード 1 1 0 のためのハートビートカウンターを更新 (例、増加) してもよい。第 1 のノード 1 1 0 はそのローカルヘルス情報をロードバランサの実装における 1 つまたは複数の他のロードバランサノード 1 1 0 へと定期的または非定期的に送信してもよく、それら 1 つまたは複数の他のロードバランサノード 1 1 0 はそのローカルヘルス情報を適宜更新 (例、第 2 のノードのためのハートビートカウンターの増加により) し、その更新されたローカルヘルス情報を 1 つまたは複数の他のノード 1 1 0 へと送信してもよい。第 2 のノード 1 1 0 のためのハートビート情報はこうしてロードバランサの実装における他のノード 1 1 0 へと伝播されてもよい。第 2 のノード 1 1 0 が正常である限り、第 2 のノード 1 1 0 から到達可能な他のすべてのノード 1 1 0 はこのように、第 2 のノード 1 1 0 のハートビートカウンターが一定期間毎に、例えば、1 秒に 1 回または 1 0 秒毎に一回、増加していることを確認すべきである。第 2 のノード 1 1 0 が、そのヘルスチェックを行う 1 つまたは複数のノード 1 1 0 により、異常であると検知された場合、ノード 1 1 0 のためのハートビートはヘルスチェックを行うノード 1 1 0 により一切送信されず、ある時間閾値の経過後、ロードバランサの実装 1 1 0 における他のノード 1 1 0 が、問

10

20

30

40

50

題のノード110が異常である、またはダウンしているを見なす。

【0110】

少なくともいくつかの実施形態において、ロードバランサノード110はその内部状態の1つまたは複数の態様を確認してもよく、ノード110が何らかの理由によるその異常を検知した場合、ノード110はそのヘルスチェックを行う他のノード110からのヘルスピンギングに対して応答を停止してもよい。したがって、異常なノード110のヘルスチェックを行うノード110は、そのノード110を以上であると見なしてもよく、そのノード110の代わりにハートビートの増加を伝播しなくてもよい。

ヘルスチェックプロトコルの説明

【0111】

少なくともいくつかの実施形態において、ヘルスチェックプロトコルはハートビートカウンター技術およびゴシッププロトコル技術を活用してもよい。ヘルスチェックプロトコルは2つの主要部分 - ヘルスチェックおよびゴシップノ故障検知を有すると見なされてもよい。

【0112】

ヘルスチェック - ロードバランサの実装におけるすべてのロードバランサノード110は、実装における1つまたは複数の他のノード110のヘルスチェックを定期的または非定期的に行ってもよい。1つまたは複数の他のノードの決定方法は後述される。ヘルスチェックの中心となる概念は、ノード110が別のノード110のヘルスチェックを行い、他のノード110が正常であると決定する場合、そのチェックを行うノード110が他のノード110のハートビートカウンターを増加させまた伝播することにより、他のノード110が正常であるとアサートする。すなわち、ノード110はそのヘルス状態を他のノードにアサートせず、その代わりに、1つまたは複数の他のノード110がロードバランサの実装における各ノード110のヘルス状態をチェックしアサートする。

【0113】

ゴシップノ故障検知 - 少なくともいくつかの実施形態において、ヘルスチェックプロトコルはロードバランサの実装におけるメンバーであるロードバランサノード110の間にロードバランサノード110のヘルス情報を伝播するゴシッププロトコルを活用してもよい。ゴシッププロトコルは迅速に収束し、分散型ロードバランサシステムの目的に十分な最終的な一貫性を保証する。少なくともいくつかの実施形態において、ゴシッププロトコルの利用により各ロードバランサノード110は、ロードバランサの実装における互いのノード110のためのハートビートカウンターを、例えばハートビートリストにおいて維持する。各ロードバランサノード110は上記のように少なくとも1つの他のロードバランサノード110のヘルスチェックを定期的または非定期的に行い、ヘルスチェックを通じてチェックを行ったノード110が正常であると決定した際に、ノード110のためのハートビートカウンターを増加させる。少なくともいくつかの実施形態において、各ロードバランサノード110は定期的または非定期的に、ロードバランサの実装における少なくとも1つの他のノード110を現在のハートビートリストの送信先として無作為に選択する。別のノード110からハートビートリストを受信した際に、2つのリスト（受信されたリストおよびそのリスト）上の各ノード110のための最大のハートビートカウンターを決定し、決定された最大のハートビートカウンターをそのハートビートリストにおいて利用することで、ロードバランサノード110は受信されたリストのハートビート情報をそのハートビートリストと組み合わせる。次にこのハートビートリストは別の無作為に選択されたノード110へと送信され、選択されたノード110がそのハートビートリストの更新等を適宜行う。この技術を用い、各正常なノード110のためのハートビート情報は最終的に（例、数秒後に）すべての他のロードバランサの実装におけるロードバランサノード110へと伝播される。所与のロードバランサノード110のためにハートビートカウンターし続ける限り、それは他のノード110により正常であると見なされる。ロードバランサノード110のハートビートカウンターがヘルスチェックおよびゴシップ方法により特定の期間中に増加されない場合は、他のロードバランサノード110がその後

10

20

30

40

50

、異常であると見なされたロードバランサノード 1 1 0 上で収束してもよい。

ヘルスチェックを行うロードバランサノード

【 0 1 1 4 】

少なくともいくつかの実施形態による、別のロードバランサノード 1 1 0 により実行されてもよいロードバランサノード 1 1 0 のヘルスチェック方法を以下に記載する。図 2 3 に関連し、少なくともいくつかの実施形態において、ノード 1 1 0 のために以下の条件のうちの一つまたは複数が決定された場合、ロードバランサノード 1 1 0 は正常であると見なされてもよい：

* ノード 1 1 0 のプロセッサの閾値（例、コアパケット処理コード 1 1 0 8 の閾値）がレディ状態（内部）である。

* ノード 1 1 0 がエッジルータ 1 0 4 の IP アドレスおよび / または MAC アドレスを把握している（内部）。

* ノード 1 1 0 のすべての閾値および / またはプロトコルハンドラーがレディ状態である（内部）。

* 北側（エッジルータ 1 0 4 / 境界ネットワーク）から、また南側（サーバ 1 3 0 / 本番ネットワーク）からの受信および発信リンクがアクティブである（外部）。

* ロードバランサの実装において用いられるネットワークインタフェースコントローラ（NIC）を通じて、ノード 1 1 0 がパケットを受信およびディスパッチすることが可能である。例えば図 2 3 で示される例示的なロードバランサノード 1 1 0 の実施形態において、ノード 1 1 0 は北向きの NIC 1 1 1 4 A および南向きの NIC 1 1 1 4 B を通じてパケットの受信およびディスパッチに成功すべきである。

【 0 1 1 5 】

1 つまたは複数のこれらのヘルス条件が所与のノード 1 1 0 に当てはまらない場合、そのノード 1 1 0 は正常でないと見なされてもよい。いくつかの実施形態において、上記条件のすべてがノード 1 1 0 に当てはまる場合にのみ、ノード 1 1 0 は正常であると見なされることに留意する。

【 0 1 1 6 】

少なくともいくつかの実施形態において、上記ヘルス条件に加えて、図 2 3 において NIC 1 1 1 4 C として示され、例えば制御プレーン通信のために用いられてもよい、各ロードバランサノード 1 1 0 上の第 3 の NIC もまた、NIC へパケットを送信し、NIC からパケットを受信することでヘルスチェックを行うノード 1 1 0 によりチェックされてもよく、第 3 の NIC のチェックが失敗した場合、チェックされているノード 1 1 0 は異常であると見なされてもよい。

【 0 1 1 7 】

図 1 3 は少なくともいくつかの実施形態による、別のロードバランサノードからのロードバランサノードのヘルスチェック方法の実施例を示す。この実施例では、ロードバランサノード 1 1 0 A はロードバランサノード 1 1 0 B のヘルスチェックを行っている。ノード 1 1 0 A および 1 1 0 B はそれぞれ、北向きの NIC（図 2 3 の NIC 1 1 1 4 A）および南向きの NIC（図 2 3 の NIC 1 1 1 4 B）を有する。1 では、ノード 1 1 0 A がパケット（例、ping パケット）をその北向きの NIC からノード 1 1 0 B の北向きの NIC へとエッジルータ 1 0 4 を通じて送信する。ノード 1 1 0 B はその北向きの NIC においてパケットを受信し、上記リストにおいて与えられた条件が満たされた場合、2 においてその北向きの NIC からノード 1 1 0 A の北向きの NIC へとファブリック 1 2 0 を通じて応答を送信する。3 において、その北向きの NIC で応答を受信した際に、ノード 1 1 0 A はパケット（例、ping パケット）をその南向きの NIC からノード 1 1 0 B の南向きの NIC へとファブリック 1 2 0 を通じて送信する。ノード 1 1 0 B はその南向きの NIC においてパケットを受信し、上記リストにおいて与えられた条件が満たされた場合、4 においてその南向きの NIC からノード 1 1 0 A の南向きの NIC へとエッジルータ 1 0 4 を通じて応答を送信する。その南向きの NIC において応答を受信した際、ノード 1 1 0 A はノード 1 1 0 B を正常であると見なしてノード 1 1 0 B のローカルハー

10

20

30

40

50

トビートカウンターを増加させ、その後それは上述のようにゴシッププロトコルに従って他のノード110へと伝播されてもよい。

【0118】

上記の代わりとして、いくつかの実施形態においては、ロードバランサノード110Bはその南向きのNICを通じて、その北向きのNICで受信された、ノード110Aの南向きのNICへの第1のpingメッセージに 응답してもよく、その北向きのNICを通じて、その南向きのNICで受信されたノード110Aの北向きのNICへの第2のpingメッセージに 응답してもよい。

【0119】

また、いくつかの実施形態においては、ノード110Aはまた、ノード110Bが正常である場合に、それ自体の第3のNICからノード110Bの第3のNICへとpingを送り、ノード110Bの第3のNICからのその第3のNIC上のpingメッセージへの 응답を受信して、(図23でNIC1114Cとして示される)制御プレーン通信のために用いられるノード110Bの第3のNICのヘルスチェックを行う。pingメッセージおよび応答は1つまたは複数の制御プレーン装置170、例えばネットワークスイッチを通過してもよい。

【0120】

上記のヘルスチェック機構は、すべてのノード110BのNICと同様にすべての受信および発信リンクならびに全方向(北、南、および制御プレーンを通じて)のノード110Bのデータ経路を実行し、顧客パケットもそうするようにpingデータパケットが内部キューをトラバースしノード110Bのディスパッチを行う時に、ノード110Bの内部ヘルス状態を検証する。

ロードバランサノードへのヘルスチェックの役割の割り当て

【0121】

少なくともいくつかの実施形態において、ロードバランサの実装におけるすべてのロードバランサノード110は、例えば構成関数を通じて、および/または図1で示すように構成要素としての構成サービス122を通じて、ロードバランサの実装におけるすべての他のロードバランサノード110のリスト(例、ソートされたリスト)へのアクセスを有する。少なくともいくつかの実施形態において、各ロードバランサノード110は、各ヘルスチェック間隔でヘルスチェックを行うためにリスト上の1つまたは複数の他のノード110を無作為に選択し、正常であると決定された場合にハートビートカウンターを増加させてもよい。リストはヘルスチェック機構を通じて現在正常であると見なされているものも異常であると見なされているものも関わらずロードバランサの実装におけるすべてのロードバランサノード110を含んでよく、現在異常なノード110も正常なノード110と同様にリストから無作為に選択されてヘルスチェックを行われてもよいことに留意する。こうして、現在異常なノード110は、そのノード110のヘルスチェックを行う1つまたは複数のノード110により正常であると決定されてもよく、そのハートビートカウンターが増加させられて他のノード110へと伝播されてもよく、このようにして異常なノード110が正常な状態に戻されてもよい。

【0122】

その代わりにいくつかの実施形態においては、各ロードバランサノード110はリスト上の1つまたは複数の他のノード110のヘルスチェックの役割および、正常であると決定された場合のハートビートカウンターの増加の役割を担ってもよい。例えばいくつかの実施形態においては、各ノード110は2つの他のノード、例えばリスト上の「左」(すなわち前)そして「右」(すなわち次)の最近傍ノード110の役割を担ってもよい。リストは円環であってもよく、リストの「最後」のノード110がリストの「最初」のノード110のヘルスチェックの役割を担ってもよく、その逆もまた同様であることに留意する。いくつかの実施形態においては、2つの他のノード110が、例えばリスト上の次の2つの最近傍ノードとして他に選択されてもよい。いくつかの実施形態においては、各ノード110は、例えば3つまたは4つの他のノード110のような、リスト上の3つ以上

10

20

30

40

50

のノード110のヘルスチェックの役割を担ってもよい。少なくともいくつかの実施形態において、ノード110からチェックされている近傍ノード110が異常であると決定された場合、ノード110はその後、異常な近傍ノード110がチェックする役割を担っていたリスト上の少なくとも1つのノードのヘルスチェックの役割を担ってもよい。少なくともいくつかの実施形態において、その近傍ノード110（例、「左」および「右」の近傍ノード）のヘルスチェックに加えて、各ロードバランサノード110もまた定期的または非定期的、リング上のノード110を無作為に選択し、その無作為に選択されたノード110のヘルスチェックを行い、正常であれば無作為なノード110のハートビートを増加させ伝播させてもよい。少なくともいくつかの実施形態において、他のノード110がすでに正常であるか正常でないかのどちらに見なされているかに関わらず、番号付きリスト上のすべての他のノード110が無作為な選択およびヘルスチェックのために考慮される。

10

【0123】

少なくともいくつかの実施形態において、各ノード110は1つまたは複数の無作為に選択されたノード110の、またはその代わりにその近傍ノード110および無作為に選択されたノードのヘルスチェックを通常の間隔で行い、その間隔はヘルスチェック間隔と称されてもよい。例えば、いくつかの実施形態においては、ハートビート間隔は100ミリ秒であってもよく、それより短い、または長い間隔も利用されてもよい。また少なくともいくつかの実施形態において、各ノード110はその現在のハートビートリストを少なくとも1つの他の無作為に選択されたノード110に通常の間隔で送信または「ゴシップ」し、それはゴシップと称されてもよい。いくつかの実施形態においては、ヘルスチェック間隔およびゴシップ間隔が同じ長さであってもよいが、必ずしも同じである必要もない。

20

【0124】

図14は少なくともいくつかの実施形態による、1つまたは複数の他のロードバランサノードのヘルスチェックを行うロードバランサノードを図示する。この実施例においては、ロードバランサの実装における8つのロードバランサノード110A-110Hがある。破線の円は実装におけるすべてのノード110の番号付きリストを表す。いくつかの実施形態においては、各ノード110は各間隔でヘルスチェックを行う、リスト上の1つまたは複数の他のノード110を無作為に選択してもよい。その代わりとして、いくつかの実施形態においては、各ロードバランサノード110は番号付きリスト上の1つまたは複数の特定のノード110のチェックの役割を担ってもよい。例えばノード110Aは、図14に示された番号付きリストに従ってその2つの最近傍ノード110Bおよび110Hのヘルスチェックの役割を担ってもよい。またロードバランサノードもまた各ヘルスチェック間隔で、番号付きリストから別のノード110を無作為に選択してもよい。この実施例で示されるように、ノード110Aもヘルスチェックのためにノード110Fを無作為に選択している。ゴシップ間隔において、ノード110Aは何か別の正常なノード110、例えばノード110Dを無作為に選択し、その現在のハートビートリストを選択した他のノード110へと、例えばUDPメッセージで送信する。ノード110は、別のノード110からハートビートリストを受信した際に、それ自身のハートビートリストを適宜更新し、次のゴシップ間隔で、ハートビートリストを1つまたは複数の無作為に選択されたノード110に伝播してもよい。

30

40

サーバノードのヘルスチェック

【0125】

上述のようなロードバランサノード110のヘルスチェックに加えて、ヘルスチェックプロトコルの実施形態は、それらノード130上のロードバランサモジュール132およびサーバ134を含むサーバノード130のヘルスチェックを行ってもよい。少なくともいくつかの実施形態において、以下の条件のうちの1つまたは両方がノード130のために決定された場合、サーバノード130は正常であると見なされてもよい：

* ロードバランサモジュール132が正常である。

50

* ロス ping (例、L7ヘルス ping)への応答に成功する。

【0126】

図15は少なくともいくつかの実施形態による、サーバノードのヘルスチェックを行うロードバランサノードを示す。少なくともいくつかの実施形態において、ロードバランサの実装におけるすべてのロードバランサノード110は、ロードバランサの実装におけるすべてのサーバノード130のリストと同様に、すべての他のロードバランサの実装におけるロードバランサノード110のリストへのアクセスを有する。1つまたは複数のリストは、例えば構成関数を通じておよび/または図1で示されるように構成要素である構成サービス122を通じて取得され、更新されてもよい。少なくともいくつかの実施形態において、図15で示されるような一貫したハッシュリングを形成するために、サーバノード130は正常なロードバランサノード110に対して一貫したハッシュを行ってもよい。少なくともいくつかの実施形態において、リング内の各サーバノード130はリング内の正常なロードバランサノード110によりヘルスチェックを行う。例えば図15では、サーバノード130Aはロードバランサノード110Aおよび110Cによりヘルスチェックを行う。これら2つのノード110は、一貫したハッシュリングにおけるサーバノード130のための第1の(ノード110A)および第2の(ノード110B)ヘルスチェックノード110と称されてもよい。所与の正常なロードバランサノード110は2つ以上のサーバノード130のヘルスチェックを行ってもよいことに留意する。例えば図15において、ロードバランサノード110Aはまた、サーバノード130Bおよび130Cのヘルスチェックを行う。また所与のロードバランサノード110は、1つまたは複数の他のサーバノード130のための第1のヘルスチェックノード110および1つまたは複数のサーバノード130のための第2のヘルスチェックノード110であってもよい。例えば、図15において、ロードバランサノード110Aはサーバノード130Aおよび130Bのための第1のヘルスチェッカーノードであり、サーバノード130Cおよび130Dのための第2のヘルスチェッカーノードである。

10

20

【0127】

少なくともいくつかの実施形態において、ロードバランサノード110が故障した場合、一貫したハッシュリングにおけるメンバーシップは変更され、まだ正常でしたがって一貫したハッシュリング上にある1つまたは複数の他のロードバランサノード110が故障したノード110によってすでにヘルスチェックされたサーバノード130のヘルスチェックの役割を担ってもよい。

30

【0128】

少なくともいくつかの実施形態において、正常なノード110はそれぞれその割り当てられたサーバノード130のヘルスチェックを通常の間隔で行い、それはサーバチェック間隔と称されてもよい。少なくともいくつかの実施形態において、サーバチェック間隔は上述のゴシップ間隔以上の長さであってもよい。

【0129】

少なくともいくつかの実施形態において、サーバノード130のヘルスチェックを行うために、正常なロードバランサノード110(例、図15のノード110A)はサーバノード130(例、図15のサーバノード130A)へのヘルス ping メッセージ(例、L7 HTTPヘルス ping メッセージ)を開始する。正常である場合、サーバノード130は ping 応答をロードバランサノード110に返信する。少なくともいくつかの実施形態において、 ping メッセージはサーバノード130上のロードバランサモジュール132により受信され処理されるため、成功した場合にはヘルスチェック ping がサーバノード130上のモジュール132が正常であることを確立する。 ping への応答時に、ロードバランサノード110はサーバノード130を正常であると見なし、サーバノード130のためのハートビートカウンターを増加させる。

40

【0130】

少なくともいくつかの実施形態において、所与の正常なロードバランサノード110によりヘルスチェックを行われたすべてのサーバノード130のためのハートビートカウン

50

ターは、例えばすでにロードバランサノード 1 1 0 ハートビートカウンターのために説明された、各ノード 1 1 0 がそのハートビートリストを少なくとも 1 つの他の無作為に選択されたノード 1 1 0 へと通常の間隔（ゴシップ間隔）で送信するゴシップ技術に従って他のロードバランサノード 1 1 0 へと伝播されてもよく、また受信ノード 1 1 0 が 2 つのリスト上の最大値に従って、それ自体のハートビートリストを更新する。

故障検知およびゴシップ

【 0 1 3 1 】

少なくともいくつかの実施形態において、上記のロードバランサノード 1 1 0 のヘルスチェックおよびサーバノード 1 3 0 のヘルスチェックを通じて取得された情報は、すべてのロードバランサノード 1 1 0 がロードバランサの実装の一貫した見解を維持できるように、ロードバランサの実装におけるすべてのノード 1 1 0 へと伝播される必要があってもよい。上述の通り、少なくともいくつかの実施形態において、ロードバランサノード 1 1 0 はゴシッププロトコルに従って、このヘルス情報を交換し伝播するため、またロードバランサノード 1 1 0 およびサーバノード 1 3 0 の故障を検知するために、互いに通信してもよい。

10

【 0 1 3 2 】

少なくともいくつかの実施形態において、各ロードバランサノード 1 1 0 は通常の間隔（ゴシップ間隔と称する）で、別のロードバランサノード 1 1 0 を無作為に選択し、他のノード 1 1 0 にその正常なロードバランサノード 1 1 0 およびサーバノード 1 3 0 に関する見解をロードバランサノード 1 1 0 およびサーバノード 1 3 0 のためのハートビートカウンターとともに送信する。ロードバランサノードまたはサーバノード 1 3 0 が正常である限り、ノードはそのヘルスチェックにパスし、そのハートビートカウンターは増加し続ける。ノードのためのハートビートカウンターが特定の間隔（故障時間間隔と称されてもよい）で変化しない場合、ノードはその後ロードバランサノード 1 1 0 により故障を疑われる。ノードが故障を疑われると、ロードバランサノード 1 1 0 はノードが異常であると決定するまで特定の間隔（異常な時間間隔と称されてもよい）で待機してもよい。この異常な時間間隔により、ノードが故障したことをすべてのロードバランサノード 1 1 0 が把握するまでロードバランサノード 1 1 0 は待機することができる。

20

【 0 1 3 3 】

図 1 6 は少なくともいくつかの実施形態による、ロードバランサノード 1 1 0 により維持されてもよい別のノード（ロードバランサノード 1 1 0 またはサーバノード 1 3 0 のいずれか）のヘルス状態またはその見解を図示する。3 0 0 で示すように、ロードバランサノード 1 1 0 がまず、問題のノードが正常であるとの見解を有すると仮定する。これはノードのためのハートビートカウンターが増加しつつあることを示す。しかし、ノードのハートビートカウンターが 3 0 2 で示すように特定の間隔（故障時間間隔）で増加しない場合、ロードバランサノード 1 1 0 はその後 3 0 4 で示すように、ノードが故障したことを疑う。3 0 6 で示すように、ノードのハートビートカウンターが特定の間隔（異常な時間間隔）で増加しない場合、ロードバランサノード 1 1 0 はその後 3 0 8 で示すように、ノードを異常であると見なす。しかしノードのためのハートビートカウンターが 3 1 0 で示すように、異常な時間間隔が終了する前に増加する場合、ロードバランサノード 1 1 0 は再びノードを正常 3 0 0 であると見なす。同様に 3 1 2 で示すように、異常なノードのためのハートビートの増加を受信することでもノードは正常 3 0 0 であると見なされうる。

30

40

【 0 1 3 4 】

ノードの異常の決定には、本明細書で別に記載するように、異常なノードがロードバランサノード 1 1 0 またはサーバノード 1 3 0 のどちらであるかによって、また、ロードバランサノード 1 1 0 の異常なノードとの関係によって、1 つまたは複数のロードバランサノード 1 1 0 による異なる行動を含んでもよい。

ロードバランサノードのデータ

【 0 1 3 5 】

少なくともいくつかの実施形態において、各ロードバランサノード 1 1 0 はロードバラ

50

ンサの実装の状態に関するデータを維持してもよい。少なくともいくつかの実施形態において、このデータは各ロードバランサノード110上の、正常なロードバランサノードリスト、疑わしいロードバランサノードリスト、およびハートビートリストを含むがそれらに限定されない、1つまたは複数のデータ構成において維持されてもよい。図17は正常なロードバランサノードリスト320、疑わしいロードバランサノードリスト322、異常なロードバランサノードリスト324、およびロードバランサノードのハートビートリスト326を維持するロードバランサノード110の実施例を示す。

【0136】

少なくともいくつかの実施形態において各ロードバランサノード110は、例えばどのノード110が正常であるか、またそれによってどのノード110がゴシッププロトコルに参加するかを決定するために用いられてもよい、正常なロードバランサノード110のリストである正常なロードバランサノードリスト320を維持してもよい。リスト320上のノード110のみがゴシッププロトコルを通じたロードバランサ情報の伝播に参与し、リスト320上のノード110のみが一貫したハッシュリング内にあると見なされ、このリスト上のノード110のみがサーバノード130のヘルスチェックを行う。ノード110はこのリスト320から、そのハートビート情報の送信先となる別のノード110を無作為に選択してもよい。またハートビートカウンターは、正常なロードバランサノードリスト320上に現在あるノード110のためのみに交換される。少なくともいくつかの実施形態において、ノードNがロードバランサノード110によるヘルスチェックにパスする場合、またはロードバランサノード110がリスト320上のどれか他のロードバランサノード110からノードNに関するゴシップメッセージを受信する場合に、ロードバランサノードNは別のロードバランサノード110の正常なロードバランサノードリスト320に追加されることができる。

【0137】

少なくともいくつかの実施形態において、各ロードバランサノード110は、ハートビートカウンター（ハートビートリスト326を参照）が特定の閾値（故障時間間隔と称されてもよい）で増加しなかったロードバランサノードのリストである疑わしいロードバランサノードリスト322を維持してもよい。ロードバランサノードEがロードバランサノード110の疑わしいロードバランサノードリスト322上にある場合、ロードバランサノード110はその後ノードEについてゴシップしない。正常なリスト320上の他のどれかのロードバランサノード110が、ノード110のハートビートリスト326上でノードEのためのカウンターより高いハートビートカウンターとともに、ロードバランサノード110にノードEについてゴシップする場合、ノードEはその後疑わしいリスト322から正常なリスト320へと移行される。ノードEが特定の閾値（異常な時間間隔と称されてもよい）でロードバランサノード110の疑わしいリスト322上に留まる場合、ノードEはロードバランサノード110により異常であるが見なされ、異常なノードリスト324上に移行される。異常なノードリスト324上のノード110（この実施例のノードG）は、ノードGがノード110によるヘルスチェックにパスした際、またはノードGのための更新されたハートビートカウンターを別のノード110から受信した際に、ロードバランサノード110の正常なノードリスト320へと移行されてもよい。

【0138】

少なくともいくつかの実施形態において、各ロードバランサノード110はすべての既知のロードバランサノード110のためのハートビートリスト326を維持してもよい。各ノード110のために、このリスト326はハートビートカウンターおよび、ハートビートカウンターが最後に変更された時を示すタイムスタンプを含んでもよい。

【0139】

少なくともいくつかの実施形態において、各ロードバランサノード110はまた、図17に図示されていないすべての既知のサーバノードのためのハートビートリストを維持してもよい。このリストはロードバランサノードのハートビートリスト326に類似していてもよい。いくつかの実施形態においては、2つのリストが組み合わせられてもよい。少な

10

20

30

40

50

くともいくつかの実施形態において、サーバノード130のためのハートビート情報は、例えばゴシッププロトコルに従って、ロードバランサノード110のためのハートビート情報とともにあるいはそれに加えて、ロードバランサノード110間に伝播されてもよい。

【0140】

図17は4つの個別のリストを示すが、リストは2つ以上が単一のリストに組み合わせられてもよいことに留意する。例えば、いくつかの実施形態においては、すべてのノード110の単一のリストが各ロードバランサノード110上に維持されてもよく、ビットフラグまたは他のデータ構成は各ノードが現在正常、疑わしい、または異常のどれであることを示すために用いられてもよい。

サーバノードデータ

【0141】

少なくともいくつかの実施形態において、ノード130上のサーバノード130およびローカルロードバランサモジュール132はロードバランサノード110を含むゴシッププロトコルに参加しない。ロードバランサノード110は、ロードバランサノードヘルスチェック方法により取得された他のロードバランサノード110に関するハートビート情報およびサーバノードヘルスチェック方法により取得されたサーバノード130に関するハートビート情報を、自身らの間でのみゴシップする（特に、各ロードバランサノード110はその正常なロードバランサノードリスト320上に現在あるノードにのみゴシップする）。

【0142】

しかし各サーバノード130/ロードバランサモジュール132は、サーバノード130が、サーバノード130が顧客のトラフィックを転送する先となるロードバランサノード110（特に、出口ノード）を決定し、どのロードバランサノードが接続公開情報の送信先となるかを決定できるように、ロードバランサの実装における正常なロードバランサノード110に関する情報を必要とする場合がある。少なくともいくつかの実施形態において、この情報をサーバノード130に提供するために、ロードバランサノード110は現在正常なロードバランサノード110（例、図17の正常なロードバランサノードリスト320）を特定する情報を用いて、定期的または非定期的にはサーバノード130を更新してもよい。少なくともいくつかの実施形態において、所与のサーバノード130のヘルスチェックの役割を担うロードバランサノード110（図15を参照）は、現在正常なロードバランサノードを特定する情報をサーバ130へと提供する役割を担う。例えば図15を参照すると、ロードバランサノード110Aはその正常なロードバランサノードリスト320をサーバノード130A、130B、130Cおよび130Dへと送信してもよく、ロードバランサノード110Bはその正常なロードバランサノードリスト320をサーバノード130C、130D、および130E等へと送信してもよい。

ロードバランサノードの故障の処理

【0143】

図18Aおよび18Bは少なくともいくつかの実施形態による、ロードバランサノードの故障の処理を示す。図18Aはロードバランサの実装の実施例を示す。現在のロードバランサの実装において、4つのロードバランサノード110A~110Dがある。エッジルータ104は顧客（図示せず）からの受信パケットをロードバランサノード110へとルーティングする。少なくともいくつかの実施形態において、エッジルータ104は、レイヤー4のフローごとにハッシュ化されたマルチパスルーティング技術、例えば等価コストマルチパス（ECMP）ルーティング技術に従って、ルーティングの決定を行ってもよい。少なくともいくつかの実施形態においてエッジルータ104は、ロードバランサノード110の提供、例えばロードバランサノード110によって開始された境界ゲートウェイプロトコル（BGP）技術セッションを介した提供を通じて顧客のトラフィックを受信するために、ロードバランサの実装において現在利用可能なロードバランサノード110について把握する。しかし、少なくともいくつかの実施形態において、BGPセッション

10

20

30

40

50

を通じて自身をエッジルータ 104 に提供するロードバランサノード 110 の代わりに、ロードバランサの実装における少なくとも 1 つの他のノード 110 が、BGP を通じたノード 110 のエッジルータ 104 へ提供の役割を担う。例えば図 18A で示されるようにいくつかの実施形態においては、所与のノード 110 の左および右の近傍ノード 110 が、所与のノード 110 をエッジルータ 104 に提供する。例えばロードバランサノード 110A はノード 110B および 110D を提供し、ロードバランサノード 110B はノード 110A および 110C を提供し、ロードバランサノード 110C はノード 110B および 110D を提供する。

【0144】

図 18A の実施例において示されるように、各ロードバランサノード 110 はまた、例えば 1 つまたは複数の無作為に選択されたノード 110、ロードバランサノードの番号付きリストにより決定された 1 つまたは複数の近傍ノード 110、あるいは 1 つまたは複数の近傍ノードおよび 1 つまたは複数の無作為に選択されたノードなどの 1 つまたは複数の他のロードバランサノード 110 のヘルスチェックを定期的に行う。また各ロードバランサノード 110 は少なくとも 1 つのサーバノード 130 のヘルスチェックを定期的に行ってもよく、その正常なロードバランサノード 110 のリストをそれがヘルスチェックを行う 1 つまたは複数のサーバノードへと送信してもよい。ロードバランサノード 110 およびサーバノード 130 のためのヘルス情報は、例えばゴシッププロトコルに従ってノード 110 間に伝播されてもよい。

【0145】

図 18B は、図 18A のロードバランサの実装の実施例における、単一のロードバランサノード 110 の故障の処理を示す。この実施例において、ロードバランサノード 110B は何らかの理由で故障している。例えば、ノード 110A および 110C はノード 110B のヘルスチェックを行ってもよく、両方ともノード 110B がそのヘルスチェックに失敗していることを検知してもよい。したがって、ノード 110A および 110C はノード 110B のためのハートビートカウンターを増加させない。ノード 110A および 110B の両方からのハートビート情報は、ゴシッププロトコルに従って、他の正常なロードバランサノード 110 (この実施例においては、唯一の他のロードバランサノードはノード 110D である) へと伝播される。すべての正常なロードバランサノード 110 (この実施例においては、ノード 110A、110C、および 110D) がノード 110B の故障において収束すると直ちに、以下のイベントの 1 つまたは複数が発生してもよいがそれらに限定されない。これらのイベントは必ずしもこの順序で発生するわけではないことに留意する。

* ノード 110A および 110C はエッジルータ 104 へのノード 110B の提供を停止する。少なくともいくつかの実施形態においてこれは、ノード 110 がエッジルータ 104 を用いてノード 110B を提供するために確立した BGP セッションの終了に関連する。各ノード 110 は、提供を行う互いのノード 110 のためにエッジルータ 104 を用いて個別の BGP セッションを確立し、そのためノード 110B のための BGP セッションの終了は提供された他のノード 110 に影響を与えないことに留意する。少なくともいくつかの実施形態においてノード 110 は、BGP セッションのための TCP クローズまたは類似のメッセージをエッジルータ 104 に送信することにより、エッジルータ 104 を用いて BGP セッションを終了する。

* ノード 110B がどのノードからも提供されなくなったことを検知すると、それを受けてエッジルータ 104 が顧客データパケットのノード 110B へのルーティングを停止する。エッジルータ 104 はまた、顧客から残りの正常なロードバランサノード 110、特にノード 110 上の入口サーバ 112 へのパケットフローを再分散するために、マルチパス (例、ECMP) ハッシュを調整する。入口サーバ 112 が顧客 -> サーバマッピングを有しない、入口サーバ 112 へとルーティングされたいずれかのパケットフローのために、マッピングが顧客 -> サーバ接続のためのフロートラッカーノードから取得されてもよく、またはその代わりに新規顧客 -> サーバ接続が図 10A ~ 10G で示される技術

10

20

30

40

50

に従って確立されてもよい。

* ノード 110A および 110C はそれぞれ、互いを提供するためのエッジルータ 104 への BGP セッションを開始してもよい。ノード 110A および 110C は両方ともノード 110B と同様にロードバランサノード 110D によってエッジルータ 104 へと提供されるため、ノード 110B が故障時にノード 110A および 110C のエッジルータ 104 への提供を停止するかもしれないという事実が、エッジルータ 104 からのこれら 2 つのノード 110 へのパケットのルーティングに停止を引き起こさないことに留意する。

* 少なくともいくつかの実施形態において、ノード 110A および 110C はこの時点では近傍ノード 110 であるため、互いのヘルスチェックの役割を担ってもよい。ノード 110B は異常であると見なされてもなお、1 つまたは複数の他のノード 110 により無作為なヘルスチェックを行われてもよいことに留意する。

* 1 つまたは複数の残りの正常なロードバランサノード 110 は、以前はノード 110B によりトラッキングされていたフローのトラッキング接続の役割を担ってもよい。例えばノード 110C およびノード 110D は、図 11C および 11D で示されるように、ノード 110B が 1 次または 2 次フロートラッカーの役割を担っていた 1 つまたは複数の接続のための 1 次または 2 次フロートラッカーとしての役割を引き継いでもよい。

* 1 つまたは複数の残りの正常なロードバランサノード 110 は、ノード 110B によりすでにヘルスチェックを行われたサーバノード 130 のヘルスチェックの役割を担ってもよい。サーバノード 130 は正常なロードバランサノードリスト（この時点ではノード 110B を含まない）を用いて残りのロードバランサノード 110 により更新される。例えば図 18B において、ロードバランサノード 110A はサーバノード 130C のヘルスチェックおよび更新を開始し、ロードバランサノード 110C はサーバノード 130B のヘルスチェックおよび更新を開始する。

* エッジルータ 104 上で、故障したノード 110B からの BGP セッションが最終的にタイムアウトする。その代わりにエッジルータ 104 は、ノード 110B の故障の認識時に BGP セッションを遮断してもよい。

【0146】

2 つのロードバランサノード 110 が同時に、またはほぼ同時に故障する可能性がある。2 つの故障したロードバランサノードが互いに隣接していない場合、故障は独立となり、個別の単一のノード 110 の故障として図 18B で示す方法に従って処理されてもよい。しかし、2 つの故障したノードが互いに隣接し（例、図 18A におけるノード 110B および 110C、その後直ちにすべての正常なロードバランサノード 110（この実施例では、ノード 110A および 110D）が故障を検知し、故障上で収束する場合に、以下のイベントの 1 つまたは複数が発生してもよいがそれらに限定されない。これらのイベントは必ずしもこの順序で発生するわけではないことに留意する。

* ノード 110A はノード 110B のためのエッジルータ 104 への BGP セッションを終了する。

* ノード 110D はノード 110C のためのエッジルータ 104 への BGP セッションを終了する。

* ノード 110A および 110D はエッジルータ 104 を用いた互いを提供する BGP セッションを開始する。

* ノード 110A および 110D は互いのヘルスチェックを開始する。ノード 110A および 110D はまた、故障したノード 110 のヘルスチェックを継続してもよいことに留意する。

* 残りの正常なノード 110 は正常なロードバランサノードリストを用いてサーバノード 130 を更新する。

* トラフィックはエッジルータ 104 からノード 110B およびノード 110C へと継続して流れてもよい。これは、これら 2 つのノード 110 がエッジルータ 104 への互いの提供を継続してもよいからである。しかしこれらの BGP セッションは最終的

10

20

30

40

50

にタイムアウトし、エッジルータ104はフローを残りの提供されたノード110へと適宜、再分散させる。

* ノード110Bおよび110Cは、ノード110Bおよび110Cがいまだに正常であると判断する場合に、ノード110Aおよび110Dをそれぞれ提供する、エッジルータ104を用いたBGPセッションを終了してもよい。

接続公開

【0147】

図1を再度参照すると、少なくともいくつかの実施形態において、ロードバランサの実装におけるロードバランサノード110がサーバ130への顧客TCP接続のための状態情報を維持する。この状態情報によりロードバランサノード110が、顧客の受信トラフィックをエッジルータ104からTCP接続の役割を担うサーバノード130へとルーティングできる。サーバノード130上のロードバランサモジュール132はそれぞれのサーバ134へのアクティブなTCP接続のリストを維持する。接続公開は、サーバノード130上のロードバランサモジュール132が、それを通じてロードバランサノード110へのアクティブな顧客TCP接続のリストを公開してもよい機構である。少なくともいくつかの実施形態において、接続公開データパケットはロードモジュール132によりロードバランサノード110へと、接続公開間隔と称されてもよい通常の間隔で形成され、公開される。

10

【0148】

少なくともいくつかの実施形態において、ロードバランサノード110により維持された接続状態情報はキャッシュの形式として見なされてもよく、特定の接続のための状態情報の維持はその接続のためのロードバランサノード110上のリースの維持とみなされてもよい。キャッシュエントリが更新されない限りは、ロードバランサノード110はデータフローを処理するサーバノード130への顧客データフローをルーティングすることができない可能性がある。接続公開機構はロードバランサノード110上でサーバノード130からの現在の接続状態情報を用いて、定期的にキャッシュを、そしてその結果リースを更新し、こうしてTCPデータパケットを顧客160から適切なサーバノード130へと流し続ける。顧客160がサーバ134へのTCP接続を終了する際、その接続に関連するサーバノード130上のロードバランサモジュール132はそのアクティブな接続のリストからのリストを破棄し、したがって接続公開機構を通じたTCP接続を公開することはなくなる。こうして、その接続に関連するロードバランサノード110（特に、接続のための入口サーバ112ならびに1次および2次フロートラッカー116）上の、その接続のための接続状態情報（キャッシュエントリまたはエントリ）は更新されなくなり、接続はロードバランサノード110により破棄される。少なくともいくつかの実施形態において、接続のためのキャッシュエントリまたはエントリは、他のどれかのアクティブな接続によりメモリを要求されるまで、ロードバランサノード110上のキャッシュを保持してもよい。

20

30

【0149】

したがって接続公開機構は定期的または非定期的な、入口サーバ112ならびに1次および2次フロートラッカー116上で接続リースを延長し、顧客のトラフィックを流れさせる。また、接続公開機構は少なくともいくつかのロードバランサノード110の故障からの回復に貢献してもよい。顧客接続のための状態情報を保持する1つまたは複数のロードバランサノード110が故障した場合、接続公開によって残りのロードバランサノード110へと提供されたアクティブな接続情報は場合によっては接続の回復に用いられてもよい。

40

【0150】

接続公開機構を用いて、サーバ134と顧客160との間の接続状態のためのサーバノード130は信頼すべきソースとなる。またサーバ134への接続の終了は、サーバノード130上のロードバランサモジュール132およびロードバランサノード110によって受動的に処理される。サーバノード130とロードバランサノード110との間にハン

50

ドシェイクは要求されない。すなわちロードバランサモジュール132は、特定の接続が終了したことを積極的にノードに通知するためにメッセージをロードバランサノード110に送信する必要がない。サーバ134が接続を終了する際、サーバ134は接続のための内部状態をクリアする。ロードバランサモジュール132はサーバ134の内部状態を用いて、接続公開パケットを組み込む。接続がサーバ134の内部状態からなくなるため、接続はロードバランサノード110へと公開されない。ロードバランサノード110上の接続のためのリースはこうして終了し、ロードバランサノード110は接続に関して受動的に忘れる。接続に用いられたロードバランサノード110のキャッシュ内のメモリはその後必要に応じて、他の接続のために用いられることが可能である。

【0151】

いくつかの実施形態においては、ロードバランサノード110により維持された接続のためのリースにはキャッシュ内の接続のためのタイムスタンプのエントリが関係してもよい。接続のリースが接続公開パケットにより更新された際に、タイムスタンプは更新されてもよい。サーバノード130上のロードバランサモジュール132によって接続が公開されなくなったために接続のリースが更新されない場合には、タイムスタンプはその後更新されなくなる。少なくともいくつかの実施形態において、遅延ガベージコレクション方法が用いられ、メモリが必要になるまで接続のためのエントリがキャッシュ内に保持されてもよい。例えば少なくともいくつかの実施形態において、キャッシュエントリ内のタイムスタンプがリース更新時間閾値と比較されてもよい。キャッシュエントリのためのタイムスタンプが閾値より遅い時間である場合、エントリはその後古くなり再利用されてもよい。しかしいくつかの実施形態においては、古いエントリは積極的にガベージコレクションされてもよい。

接続公開受信者

【0152】

少なくともいくつかの実施形態において、各顧客TCP接続のために接続状態を維持する3つのロードバランサノード110がある - 入口サーバ112の役割を担うノード110、1次フロートラッカー116の役割を担うノード110、および2次フロートラッカー116の役割を担うノードである。一貫したハッシュリングにおいて1次フロートラッカー116ノードおよびその後続ノードを見つけるためのTCPフローに対する一貫したハッシュ関数を適用することにより、所与のTCPフローのために、例えばロードバランサノード110によって1次および2次フロートラッカー116を決定することができる。TCPフローのための入口サーバ112の役割を担うロードバランサノード110は、エッジルータ104の内部マルチパス(例、ECMP)ハッシュ関数に基づき、エッジルータ104からそのフローのためのトラフィックを受信するノード110である。ノード110の故障または追加がある場合には、入口サーバ112の役割を担うロードバランサノード110は多くのアクティブなTCPフローのために変更されてもよく、少なくともいくつかのアクティブなTCPフローのためのフロートラッカーの役割を担うロードバランサノード110が変更されてもよい(例、図11A~11Dを参照)。サーバノード130上のサーバ132へのすべてのTCPフローのために、そのサーバノード130上のロードバランサモジュール132は状態情報を維持し、入口サーバ112がロードバランサノード110からのトラフィックを受信することから、どのロードバランサノード110がそのTCPフローのための入口サーバ112であるかを示す。しかし少なくともいくつかの実施形態において、ロードバランサモジュール132は使用された一貫したハッシュ関数を把握していなくてもよい。どのロードバランサノード110がTCPフローのための1次および2次フロートラッカーの役割を担っているかを、ロードバランサモジュール132は把握しなくてもよく、また決定できなくてもよい。すなわち、少なくともいくつかの実施形態において、ロードバランサモジュール132は一貫したハッシュを行わない。

アクティブな接続情報の公開

【0153】

図19Aおよび19Bは少なくともいくつかの実施形態による、接続公開技術を図示する。図19Aはアクティブな接続情報をロードバランサノードに公開するロードバランサ(LB)モジュールを示す。少なくともいくつかの実施形態において、各ロードバランサモジュール132はサーバノード130上の各アクティブなTCPフローのための情報を収集し、接続公開パケットを形成する。所与のTCPフローのための情報は、フローのための入口サーバ112を担うロードバランサノード110を特定する情報を含む。接続公開パケットがレディ状態である際(例、接続公開間隔に達した際)、ロードバランサモジュール132はロードバランサノード110を、例えばすでに述べたようにサーバノード130のヘルスチェックを行うロードバランサノード110からサーバノード130へと定期的に送信される正常なロードバランサノード110のリストから、無作為に選択する。ロードバランサモジュール132はその後接続公開パケットを選択したノード110へと送信する。例えば図19Aでは、ロードバランサモジュール132Aがある接続公開パケットをロードバランサノード110Aに送信し、後ほど別の接続公開パケットをロードバランサノード110Bに送信する。

10

【0154】

図20は少なくともいくつかの実施形態による、各ロードバランサモジュール132により実行されてもよい接続公開方法のハイレベルフローチャートである。500で示すようにロードバランサ(LB)モジュール132は、それぞれのサーバノード130上のすべてのアクティブなTCPフローのための接続公開エントリを作成する。少なくともいくつかの実施形態において、ロードバランサモジュール132はサーバノード130上のサーバ134が処理するアクティブなTCP接続の組を、例えばサーバノード130上の/`proc/net/tcp`から取得する。すべてのアクティブなTCP接続のために、ロードバランサモジュール132は(例、ローカルで維持されたアクティブな接続のテーブル内で)TCPフローのための入口サーバ112の役割を担うロードバランサノード110を検索し、接続のためのTCPタプル(例、顧客のIPアドレス、顧客用ポート、サーバ(パブリック)IPアドレス、およびサーバポートから成る4タプル)および接続のための入口サーバ112を示す接続公開エントリを作成する。各ロードバランサモジュール132は、接続のために受信されたパケットを最後に送信したロードバランサノード110を示す、各アクティブなTCP接続のための情報を維持すること、また、この情報が各アクティブな接続のための入口ノード110を特定するために、ロードバランサモジュール132によって用いられてもよいことに留意する。

20

30

【0155】

502で示すように、ロードバランサモジュール132は接続公開パケット(1つまたは複数の接続公開エントリを含み、各アクティブなTCP接続毎に1つのエントリ)の送信先となるロードバランサノード110を無作為に選択する。少なくともいくつかの実施形態において、接続公開パケットの送信準備が完了しているとロードバランサモジュール132が決定した際に、ロードバランサモジュール110は無作為に選択されてもよい。少なくともいくつかの実施形態において、この決定は接続公開間隔に従って下される。非限定的な実施例として、接続公開間隔は100ミリ秒(ms)、または1秒であってもよい。少なくともいくつかの実施形態において、ロードバランサノード110の内の1つからすでに受信された正常なロードバランサノード110のリストから、ロードバランサモジュール110が選択される。504で示すように、ロードバランサモジュールはその後接続公開パケットを、選択されたロードバランサノード110に公開する。少なくともいくつかの実施形態において、接続公開パケットは例えばUDPパケットといったステータスパケットである。いくつかの実施形態において接続公開パケットは、対象のロードバランサノード110へのパケットの送信前に圧縮されてもよい。少なくともいくつかの実施形態において接続公開情報は、2つ以上のパケットにおける対象のロードバランサノード110へと送信されてもよい。

40

【0156】

要素504から要素500へと戻る矢印が示すように、ロードバランサモジュール13

50

2は継続的に接続公開データパケットを作成し、無作為なノード110を選択し、そしてパケットを選択したノードへと送信してもよい。上述のように、ロードバランサノード110がロードバランサノード110上の接続リースを維持するための現在のアクティブな接続情報を用いて比較的定期的に取りフレッシュされるように、これは接続公開間隔に従って行われてもよい。

【0157】

少なくともいくつかの実施形態において、接続公開データパケットはロードバランサモジュールによってロードバランサノード110へと無作為に分散されるため、接続公開パケットを受信するロードバランサノード110は、接続公開データパケットにおけるアクティブな接続情報の接続のための正しい入口/1次/2次ノード110への分散の役割を担う。図19Bならびに図21および22は、少なくともいくつかの実施形態において用いられてもよいアクティブな接続情報の分散方法を示す。

10

【0158】

図19Bは少なくともいくつかの実施形態による、ロードバランサノード110間のアクティブな接続情報の分散を示す。ロードバランサノード110がロードバランサモジュール132から接続公開パケットを受信する際、ロードバランサノード110は、フローのための入口ノードならびに1次および2次フロートラッカーノードを決定するために、その中に示された各TCPフローのための情報を分析してもよい。ロードバランサノード110が、フローのためのそれら役割のうち1つを担っている場合、ロードバランサノード110はフローのための情報を(例、その状態情報のキャッシュの更新により)消費する。少なくともいくつかの実施形態において、ロードバランサノード110はまたフローのための情報を、フローのための他の役割を担っている1つまたは複数の他のノード110へと送信される1つまたは複数のパケットフローに盛り込んでよい。接続公開パケットにより示される残りのフローのために、ロードバランサノード110はアクティブな接続情報を2つ以上のよりデータ量の少ないパケットに分割して各パケットを1つまたは複数の他のロードバランサノード110に送信する。例えば少なくともいくつかの実施形態において、1つまたは複数のフローのためにアクティブな接続情報を含むパケットは、1つまたは複数のフローのための入口サーバ112、1次フロートラッカー116A、および2次フロートラッカー116Bの役割を担うロードバランサノード110に送信されてもよい。

20

30

【0159】

図21は、少なくともいくつかの実施形態による、対象のロードバランサノード110への接続公開パケットにおいて受信されるアクティブな接続情報の分散方法のフローチャートである。520で示すように、ロードバランサノード110はロードバランサモジュール132から接続公開パケットを受信する。ロードバランサモジュール132は、例えば図19Aおよび20に関連して上述したようにパケットを生成し、パケットを受信するためのロードバランサノード110を選択した。接続公開パケットは、受信されたパケットの送信元であるサーバノード130を特定する情報(例、サーバノード130上のロードバランサモジュール132のIPアドレス)およびアクティブなTCP接続を特定するエントリのリスト(例、各接続のための顧客のIPアドレス、顧客用ポート、サーバ(パブリック)IPアドレス、およびサーバポートから成る4タプル)を含んでもよい。

40

【0160】

図21の要素522~530において、ロードバランサモジュール110は受信された接続公開パケットで示されるアクティブなTCP接続情報を繰り返し処理する。522で示すようにロードバランサノード110は、それぞれのTCPフローのための入口ノード110ならびに1次および2次フロートラッカーノード110を決定するために、パケット内の次のTCPフローのためのエントリを分析する。少なくともいくつかの実施形態において、ロードバランサノード110は接続公開エントリから入口ノード110を特定する。少なくともいくつかの実施形態において、TCPフローのための1次および2次フロートラッカーノード110は一貫したハッシュ関数に従って決定されてもよい。524に

50

においてロードバランサノード110が検証中のTCPフローのための役割のうちの1つを担う場合、その後526においてロードバランサノード110がフローのための情報を、例えばその状態情報のキャッシュの更新によって消費する。528で示すように、ロードバランサノード110はTCPフローのための接続公開エントリを、構成中であり別のロードバランサノード110に送信される予定のケットに追加してもよい。530において接続公開ケット内にフローのための接続公開エントリがさらにある場合は、その後メソッドは522へと戻り、次のエントリの処理を行う。そうでなければロードバランサノードは、それぞれがオリジナルの接続公開ケットからの接続公開エントリのサブセットを含む1つまたは複数の新規構成されたケットを、532で示すようにケットのための対象のロードバランサノード110へと送信する。少なくともいくつかの実施形態において、対象のロードバランサノード110に送信されるケットは、例えばUDPデータケットといったステートレスケットである。いくつかの実施形態において、ケットを対象のロードバランサノード110に送信する前に、ケットは圧縮されてもよい。

10

【0161】

このように少なくともいくつかの実施形態において、図21の要素522～528ではフロートラッカーノード110が、受信された接続公開ケット内の接続公開エントリから522で決定された情報に従って他のノード110のうち特定の1つへとそれぞれ送信されることになる1つまたは複数のケット(例、UDPケット)を構築する。少なくともいくつかの実施形態において、別のノード110に送信されるケットは、対象のノード110が入口ノード110、1次フロートラッカーノード110、または2次フロートラッカーノード110としての役割を担うTCPフローのためのエントリを含む。いくつかの実施形態においては、所与のロードバランサノード110がTCPフローのための入口および1次フロートラッカーノードの両方の役割を担ってもよく、またはTCPフローのための入口および2次フロートラッカーノードの両方の役割を担ってもよいことに留意する。

20

【0162】

図22は少なくともいくつかの実施形態による、対象のロードバランサノード110への接続公開ケットにおいて受信されるアクティブな接続情報の分散の代替方法を示す。550で示すように、ロードバランサノード110はロードバランサモジュール132から接続公開ケットを受信する。この方法では、552で示すように、ロードバランサモジュール110上の処理によってケット内の接続公開エントリの分析が行われ、受信されたケットの1つまたは複数のよりデータ量の少ないケットへの分割が適宜行われる。ロードバランサモジュール110はこの処理の間、ローカルでフロー情報を消費しない。接続公開ケットが1つまたは複数のケットに分割されると、ケットがその後554～560で示すように処理される。554においてケットのための対象のノード110がこのロードバランサノード110である場合、ロードバランサノード110はその後556で示すようにローカルでケットを消費する。そうでなければケットは対象のロードバランサノード110へと送信される。560において処理すべきケットがさらであれば、その後メソッドは554に戻る。そうでなければメソッドは完了する。

30

【0163】

このようにしてロードバランサモジュール132から接続公開ケットを受信するロードバランサノード110は、接続公開ケットを特定の他のロードバランサノード110に特有の2つ以上のよりデータ量の少ないケットに分割し、ロードバランサノード110により現在処理中であるいずれかのTCPフローのためのフロー情報を内部で消費しながらケットを適宜分散してもよい。その間、他のロードバランサノード110もまた、接続公開ケットをロードバランサモジュール132から受信し、接続公開エントリを複数のよりデータ量の少ないケットに分割し、そしてよりデータ量の少ないケット対象のノード110に送信して、ノード110間にアクティブな接続情報を分散してもよい。

40

接続公開トリガ

【0164】

50

少なくともいくつかの実施形態において、接続公開はロードバランサモジュール132上で1つまたは複数の異なるイベントによってトリガされてもよい。上述のようにいくつかの実施形態において、接続公開パケットはロードバランサノード110上のTCP接続のためのリースを更新するべく、接続公開間隔、例えば100msまたは1秒間隔に従って生成され、無作為に選択されたロードバランサノード110へと送信されてもよい。いくつかの実施形態においては、ロードバランサノード110のメンバーシップにおける変更は、即時の接続公開イベントをトリガしてもよい。少なくともいくつかの実施形態において、ロードバランサモジュール132はそれぞれのサーバノード130のヘルスチェックを行うロードバランサノード110の1つより送信された正常なロードバランサノード110のリストから変更について学習してもよい。リストに従った変更（削除または追加のいずれか）の検知時には、変更に影響されたTCP接続がロードバランサノード110によって迅速に回復できるように、ロードバランサモジュール132は接続公開パケットを生成しロードバランサノード110へと送信してもよい。

10

パケットループの阻止

【0165】

接続公開パケットの処理の間にロードバランサレイヤーのメンバーシップが変更された場合、接続公開パケットループが発生してもよい。第1のノード110はロードバランサモジュール132から接続公開パケットを受信し、よりデータ量の少ないパケットを第2のノード110へと送信してもよい。しかしメンバーシップが変更された場合には、パケットが第1のノード110へ移行するべきであると第2のノード110が決定してもよく、その結果パケットを第1のノード110へと転送してもよい。少なくともいくつかの実施形態において、このループの発生を阻止するためにロードバランサモジュール132から受信された接続公開パケットと、ロードバランサノード110から受信されたそれらとで異なるポート番号が用いられてもよく、またロードバランサノード110は他のロードバランサノード110から受信された接続公開パケットの再分散を行わない。

20

接続公開パケット分散の代替方法

【0166】

上記の接続公開方法において、ロードバランサモジュール132は接続公開パケットの送信先であるロードバランサノード110を無作為に選択する。しかしいくつかの実施形態においては、ロードバランサノード110の選択に他の方法が用いられてもよい。例えばいくつかの実施形態において、ロードバランサノード132は、1つまたは複数のアクティブなTCPフローの処理を行う特定の入口ノード110をそれぞれ対象とする1つまたは複数の接続公開データパケットを構築してもよく、また、1つまたは複数のパケットを1つまたは複数の対象の入口ノード110へと送信してもよい。1つまたは複数の入口ノード110はその後、アクティブな接続情報を接続のための1次および2次フロートラッカーへと再分散する。別の実施例として、いくつかの実施形態においては、接続公開パケットの単一の、無作為に選択されたノード110への送信の代わりに、各接続公開パケットがロードバランサモジュール132によって2つ以上の正常なノード110へと、またはすべての正常なノード110へと送信されてもよい。

30

ロードバランサノードアーキテクチャ

40

【0167】

図23は少なくともいくつかの実施形態によるロードバランサノード110のためのソフトウェアスタックアーキテクチャの実施例を示し、限定的な意図を持たない。このソフトウェアスタックアーキテクチャの実施例においては、Java Native Interface (JNI (商標)) 1104技術を用いて、ロードバランササーバネイティブコード1106およびコアパケット処理コード1108、例えばIntel (商標) Dataplane Development Kit (DPDK) 技術コードを含んでもよいネイティブコードのレイヤーを管理する、単一のJava (商標) 技術処理1102内でロードバランサノード110が動作する。ネイティブコードは2つのネットワークインタフェースコントローラ (NIC 1114Aおよび1114B) へのインターフェー

50

スとなってもよい。第1のNIC (NIC 1114A)は「北」、すなわちエッジルータ104向きに面してもよい。第2のNIC (NIC 1114B)は「南」、すなわちサーバノード130向きに面してもよい。少なくともいくつかの実施形態において、NIC 1114Aおよび1114BはTCPスタックを維持しなくてもよい。したがって少なくともいくつかの実施形態は、ロードバランサノード110が制御プレーンを通じた処理との通信を行うことができ、またその逆も可能になるようにTCP接続をサポートする第3のNIC 1114Cを含んでもよい。その代わりに、いくつかの実施形態においては、第1の、北向きのNIC 1114Aおよび第2の、南向きのNIC 1114Bのみがロードバランサノード110において実装されてもよく、また第2の、南向きのNIC 1114BがTCPスタックを実装してもよい。ロードバランサノード110はTCPスタックを通じて制御プレーンを通じた処理との通信を行ってもよい。ロードバランサノード110はまた、オペレーティングシステム(OS)技術ソフトウェア1112、例えば、Linux(商標)カーネル、およびOS技術ソフトウェア1112上のJava Virtual Machine(JVM(商標))技術ソフトウェア1110レイヤーそしてJNI 1104技術を含む。

10

【0168】

少なくともいくつかの実施形態において、分散型ロードバランサシステム内のロードバランサノード110はそれぞれ多くのデータフローを高パケットレートで同時に処理しなければならない場合がある。少なくともいくつかの実施形態において、要求レベルのスループットを達成するために、ロードバランサノード110が高性能パケット処理のためのIntel(商標)Dataplane Development Kit(DPDK)技術を活用してもよい。DPDK技術により、ユーザースペースプログラムがネットワークインタフェースコントローラ(NIC)から直接パケットを読み取ること/ネットワークインタフェースコントローラ(NIC)へと直接パケットを書き込むことが可能になり、またDPDK技術がLinuxカーネルネットワークスタックの多くのレイヤーを(Linux gbe基本NICドライバを除いて)バイパスする。パケット処理のためのDPDK手法は、ビジーリングにおいてNICハードウェアに直接ポーリングを行う専用CPUコアのための割り込みハンドラーを利用した入力を拒否する。この手法は、専用CPUコアをビジーリングにおいて継続的に動作させることで、熱出力の増加と引き換えにはるかに高いパケットレートを可能にする。DPDK技術はまた、CPUコア管理、ロックフリーのキュー、メモリプールおよび同期プリミティブを含むパケット処理ツールも提供する。図24で示されるように、DPDK技術においては、専用CPUコア600が各特定のタスクのために用いられてもよく、また非停止キュー602を用いて、あるCPUコア600Aから別のCPUコア600Bへと作業が渡される。

20

30

【0169】

DPDKキュー602は高速の2のべき乗リングバッファを用いて実装されてもよく、単一および複数の生産者/消費者バリエーションのサポートを行ってもよい。複数の生産者/消費者バリエーションは、アクセスを同期するためにコンペアアンドスワップ(CAS)ループを含むため、真の意味でロックフリーではない。バッファへのポインタのみが読み取られキュー602に書き込まれるように、すべてのパケットバッファメモリはメモリプール内に事前割り振りされてもよい。メモリプールはキューとして実装されてもよく、メモリチャネルやメモリバンクをまたいでメモリを分散させるため最適化されてもよく、非一様性メモリアクセス(NUMA)により最適化された割り当てをサポートしてもよい。少なくともいくつかの実施形態において、パケットバッファにはバッファコピーを要求せずに外部ネットワークレイヤーヘッダーを追加/取り除くことができるカプセル化/脱カプセルの動作をサポートするために、各パケットバッファのヘッドルームおよびテールルームに過剰な割り当てを行うMbufパラダイムのような方法を用いてもよい。

40

【0170】

ロードバランサノード110の少なくともいくつかの実施形態において、DPDK技術を活用するコアパケット処理アーキテクチャが実装されてもよい。各ロードバランサノ

50

ド 1 1 0 はコアパケット処理アーキテクチャに従って、少なくとも 1 つの実装されたマルチコアパケットプロセッサを含んでもよい。コアパケット処理アーキテクチャはパケットフローのための、マルチコアパケットプロセッサのキューおよびコアを通じた単一の生産者 / 単一の消費者パラダイムを用いてもよい。このパラダイムでは、各キューはただ 1 つのコアへの入力を行い、各コアはただ 1 つのコアへの出力を互に行い、パケットを与える。また、マルチコアパケットプロセッサにおいてコアにより用いられたメモリは共有されない。各コアにはそれ自体の別のメモリ領域がある。したがって、コア間にはメモリやキューの共有はなく、メモリやキューの競合もなく、メモリやキューにオーナーシップの要求 (R F O) やコンペアアンドスワップ (C A S) 等の機構の共有の必要もない。図 2 5 および 2 6 にはコアパケット処理アーキテクチャに従って実装されるマルチコアパケットプロセッサの実施例が示される。

10

【 0 1 7 1 】

図 2 5 は少なくともいくつかの実施形態による、データフローの処理のために D P D K 技術を活用するコアパケット処理アーキテクチャに従って実装されたマルチコアパケットプロセッサの実施例を示す。コアパケット処理アーキテクチャは単一の生産者 / 単一の消費者パラダイムに従って、マルチコアパケットプロセッサとして実装されてもよい。少なくともいくつかの実施形態において、図 2 3 に示されるように、ロードバランサノード 1 1 0 はそれぞれ 2 つのネットワークインタフェースコントローラ (N I C) を有する - 境界ネットワーク / エッジルータ 1 0 4 に面する北向きの N I C 1 1 1 4 A および本番ネットワーク / サーバノード 1 3 0 に面する南向きの N I C 1 1 1 4 B である。少なくともいくつかの実施形態において、N I C 1 1 1 4 は 1 0 G p b s N I C であってもよい。ロードバランサノード 1 1 0 を通過するパケットの大部分がこれら 2 つの N I C (N I C 1 1 1 4 A または 1 1 1 4 B のいずれか) で受信され、処理され (例、カプセル化され、あるいは脱カプセル化され) 、一方の N I C (N I C 1 1 1 4 B または 1 1 1 4 A のいずれか) に伝達される。

20

【 0 1 7 2 】

図 2 5 を参照すると、少なくともいくつかの実施形態においてロードバランサノード 1 1 0 は 2 つの C P U コア、受信 (R X) コア 6 1 0 および伝達 (T X) コア 6 3 0 を各 N I C 1 1 1 4 のためにスピニアップする。ロードバランサノード 1 1 0 も両方の N I C 1 1 1 4 のための両方向のパケットの処理を行ういくつかの作業員コア 6 2 0 を行う。この実施例では 4 つの作業員コア 6 2 0 A ~ 6 2 0 D が用いられる。受信コア 6 1 0 は入力キューからの受信パケットのバッチを、それらが N I C 1 1 1 4 に到着する際に読み取り、各パケットのための作業の大部分を行う作業員コア 6 2 0 へとパケットを分散させ、各受信コア 6 1 0 はパケットを各作業員コア 6 2 0 のためのそれぞれの作業員入力キュー 6 1 2 に与える。少なくともいくつかの実施形態において、受信コア 6 1 0 は各受信パケットにおいてレイヤー 4 「フローハッシュ」技術 (上述のエッジルータ 1 0 4 により用いられてもよいフローごとにハッシュ化されたマルチパスルーティング技術に類似) を実行し、いずれの特定の顧客接続 (I P アドレスおよびポートにより区別される) も同一の作業員コア 6 2 0 によって確実に処理されるようにしながらパケットを作業員コア 6 2 0 へと分散させてもよい。これは、各作業員コア 6 2 0 が常に同一のパケットのサブセットを見てもよいということの意味し、ロックを要求されないように作業員コア 6 2 0 が管理する状態データにおける競合を除去する。受信されたパケットへのポインタは作業員コア 6 2 0 が継続的に新規入力のために監視する作業員キュー 6 2 2 間に分散されてもよい。作業員コア 6 2 0 は各接続のための状態の管理の役割を担い (例えば、割り当てられたサーバノード 1 3 0) 、パケットをアウトバウンドキュー 6 3 2 の 1 つに転送する前のパケット上の U D P のカプセル化または脱カプセル化を行ってもよい。伝達コア 6 3 0 は作業員コア 6 2 0 アウトバウンドキュー 6 3 2 を通って循環し、出力パケットを、キュー 6 3 2 に現れるときに対応する N I C 1 1 1 4 に書き込む。

30

40

【 0 1 7 3 】

図 2 6 は少なくともいくつかの実施形態による、データフローの処理のために D P D K

50

技術を活用するコアパケット処理アーキテクチャに従って実装されたマルチコアパケットプロセッサの別の実施例を示す。コアパケット処理アーキテクチャは単一の生産者/単一の消費者パラダイムに従って、マルチコアパケットプロセッサとして実装されてもよい。少なくともいくつかの実施形態において、高スループット顧客TCPフローの処理に加えて、ロードバランサノード110上のDPDKコアアーキテクチャもまたARP、DHCP、およびBGP等の他のプロトコルのための、北および南向きのNIC1114上でのパケットの送受信に用いられてもよい。図26で示す実施形態においては、作業員コア620Aはこれら他のプロトコルのためのパケットの処理専用である。これらパケットの処理は一般的に顧客TCPフローよりも低速で行われるため、この作業員コア620Aは「低速の」作業員コアと称されてもよく、一方顧客TCPフローのみを処理する他の作業員コア620B~620Dは高速の作業員コアと称されてもよい。北向きのおよび南向きのNIC1114において受信パケットを処理するそれぞれ受信コア610Aおよび610Bは、低速の作業員コア620Aにより処理される予定のパケットを特定し、パケットを低速の作業員コア620Aのための入力キュー622へと導いてもよい。低速の作業員コア620Aはまた、Java/JNIにより生成されたパケットのための入力キュー622、およびJava/JNIへの出力パケットのための出力キュー634を監視してもよい。低速の作業員コア620Aはまた、低速の作業員コア620Aが高速の作業員コア620B~620Dのそれぞれへとパケット例えば接続公開データパケットを送信できるように、高速の作業員コア620B~620Dのそれぞれのための入力キュー622への出力を行ってもよい。低速の作業員コア620Aはまた伝達コア630Aおよび630Bのそれぞれに流れ込むアウトバウンドキュー632を有する。

【0174】

少なくともいくつかの実施形態において、各高速の作業員コア620B~620Dの第3の入力キュー622は低速の作業員コア620Aからの出力キューである。少なくともいくつかの実施形態において、この第3の入力キュー622は例えば、それぞれが接続状態情報を含む接続公開パケットの受信および処理のために、高速の作業員キュー620B~620Dによって用いられてもよい。これら接続公開パケットの少なくともいくつかのためには、伝達コア630への出力はなくてもよい。代わりに、データパケットにおける接続状態情報が、例えばそれぞれの高速の作業員コア620が維持する1つまたは複数のパケットフローのための格納された状態の更新により、高速の作業員コア620から消費されてもよい。こうして、高速の作業員コア620B~620Dへの入力を行う低速の作業員コア620Aからの出力キューが、入力キュー622以外の、受信コア610から直接の、高速の作業員コアの格納された状態の更新のためのパスを提供してもよい。

【0175】

少なくともいくつかの実施形態において、図25および26のマルチコアパケットプロセッサは受信パケットをフィルターにかけ、有効なパケットのみを処理し、出力してもよい。例えば、少なくともいくつかの実施形態において、受信コア610は作業員コア620のいずれにもサポートされていないプロトコルのパケットをフィルターにかけて取り除いてもよく、その結果作業員コア620へとパケットを送信しなくてもよい。少なくともいくつかの実施形態において、作業員コア620はパケットの処理時に、それぞれがまずそれぞれの作業員入力キュー622から読み取られたパケットを分析し、パケットがさらなる処理のために受け入れられるべきか、また伝達コア630への出力を行うべきかを決定してもよく、また単に受け入れられた伝達コア630へのパケットの処理および出力を完了してもよい。受け入れられなかったパケットは破棄されてもよい。例えば作業員コア620は各パケットのアドレス情報を確認し、ロードバランサされている有効なアドレスを対象とするパケットの受け入れのみを行い、他のパケットは破棄してもよい。

境界ゲートウェイプロトコル(BGP)データの処理

【0176】

少なくともいくつかの実施形態において、コアアーキテクチャの内部および外部の、BGP顧客に関連するパケットフローは以下のように処理されてもよい。NIC1114A

10

20

30

40

50

および 1114B は Linux カーネルに向かわないので、エッジルータ 104 への TCP 接続は図 26 で示されるようにコアアーキテクチャにより遮断され、出力キュー 634 を通じて BGP パケットを Java スペースへと渡す低速の作業者コア 622A によって処理される。これら TCP データパケットは BGP 顧客へと伝達される前にロードバランサノード 110 上の 1 つまたは複数のモジュールによりさらに処理される。この処理には TCP 接続を管理し、パケットを TCP ストリームへと効率的に変換するための Linux カーネルによる処理が含まれる。この設計により、標準 Java TCP ソケットライブラリを用いた BGP 顧客の書き込みが可能になる。

【0177】

図 27 は少なくともいくつかの実施形態による、ロードバランサ (LB) ノード処理 650 による受信 BGP TCP データパケットの処理を示す。エッジルータ 104 からのパケットが北向きの NIC 640 に到着し、受信コア 652 のための入力キュー 640 に入る。受信コア 652 はキュー 640 からのパケットを読み取り、パケットを BGP パケットとして特定し、パケットを低速の作業者コア 656 のための入力キュー 654 上に配置する。低速の作業者コア 656 はパケットを検証し、それを JNI 出力キュー 658 上に配置する。JNI パケット受信装置 660 はキュー 658 からのパケットを、JNI を介して読み取り、ソース / 宛先アドレスをマングルし、パケットを raw ソケット 644 に書き込む。Linux カーネル 646 は生のパケットを受信し、TCP プロトコルに従ってそれを処理し、ペイロードデータを TCP ソケット Input Stream へと追加する。パケットからのデータはその後 BGP 顧客 662 内の Java TCP ソケットに伝達される。

【0178】

図 28 は少なくともいくつかの実施形態による、ロードバランサ (LB) ノード処理 650 による、発信 BGP TCP データパケットの処理を示す。BGP 顧客 662 はデータを Linux カーネル 646 の Java TCP ソケットへと書き込む。Linux カーネル 646 は TCP プロトコルに従ってデータを処理し、データを 1 つまたは複数の TCP パケットへと変換する。少なくともいくつかの実施形態において、1 つまたは複数の TCP パケットは 127 . x . x . x IP テーブルルールに適合する。1 つまたは複数の TCP パケットは出力キュー 648、例えば Netfilter LOCAL_OUT キュー上に配置される。JNI を通じてキュー 648 を監視する JNI パケット受信装置 670 の Java スレッドは 1 つまたは複数の TCP パケットを受信し、各 NF_STOLEN に印を付けてカーネル 646 にそれらを忘れさせる。Java スレッドはソース / 宛先アドレスをマングルし、1 つまたは複数のパケットを低速の作業者コア 656 のための JNI 入力キュー 672 へと JNI を通じて追加する。低速の作業者コア 656 はその JNI 入力キュー 672 から 1 つまたは複数の TCP パケットを受信してパケットを北向きの NIC 640 伝達コア 666 のためのアウトバウンドキュー 664 上に配置する。伝達コア 666 はその入力キュー 664 から 1 つまたは複数の TCP パケットを読み取り、それらを北向きの NIC 640 に書き込む。TCP パケットは、NIC 640 によってエッジルータに送られる。

分散型ロードバランサのシミュレーションおよびテスト

【0179】

本明細書に記載のロードバランサは、多くの独立した構成要素 (例、ルータ、ロードバランサノード、ロードバランサモジュール、等) の対話を要求する分散型システムである。ノードの故障、メッセージの破棄、および遅延等のシナリオのシミュレーションと同様に、分散型構成要素、ロジック、およびプロトコルのテストを行うための、分散型ロードバランサを単一処理において動作できるようにするテストシステムの実施形態について記載する。単一処理においては、複雑なネットワークポロジ (例、本番ネットワーク) において、コードを複数のホストに展開する必要なく対話のテストを行うことができる。これを達成するための、メッセージバスと称される、複数のロードバランサ構成要素を単一処理内でまたは単一処理として構成させて実行させることができるソフトウェア機

10

20

30

40

50

構について記載する。単一処理は単一のホストシステムにおいて実行されてもよい。メッセージバス機構により、分散型ロードバランサシステムの単一処理としてのテストが、例えばロードバランサ構成要素（例、ロードバランサノードおよびロードバランサモジュール）にとっては実際の本番ネットワークで動作しているかのような単一のホストシステムにおいて可能になる。

【0180】

メッセージバスは、分散型ロードバランサが単一処理として動作することを可能にするフレームワークを提供する。処理における1つまたは複数のメッセージバスレイヤーのそれぞれが分散型ロードバランサの構成要素の間のネットワーク（例、Ethernet（登録商標））セグメントのシミュレーションを行う。分散型ロードバランサシステムのソフトウェア構成要素は、構成要素をメッセージバス環境内で動作させるために特別な方法で書かれる必要はない。代わりに、分散型ロードバランサシステムの構成要素が生成するパケットを遮断し、パケットを本物の物理ネットワークの代わりにメッセージバスレイヤーに提供された模擬ネットワーク内に導き、パケットを対象の構成要素へと伝達する構成要素（メッセージバスNICまたはパケットアダプタと称されてもよい）をメッセージバスフレームワークが提供する。メッセージバスレイヤーは、構成要素間の通信のための1つまたは複数のTCP/IPスタックを実装しない。その代わりに、メッセージバスレイヤーはホストシステムのオペレーティングシステム（OS）とのインターフェースを行い、ホストシステムのTCP/IPスタックを用いる。メッセージバスレイヤーはOSにより提供されたTCP/IPスタックを、顧客およびサーバが予期するTCPストリームの、メッセージバスが遮断し伝達する個別のパケットからの、そしてそういった個別のパケットへの変換のために活用する。

【0181】

少なくともいくつかの実施形態において、メッセージバスとのインターフェースを行うために、ロードバランサ構成要素は少なくとも1つのメッセージバスネットワークインタフェースコントローラ（NIC）を与えられていてもよい。各NICは有効なメディアアクセス制御（MAC）アドレスを有しており、それにより物理ネットワークとのやり取りの代わりに、メッセージバス模擬ネットワーク環境へとパケットを送信し、またそこからパケットを受信する。メッセージバスNICは物理ネットワークの代わりにメッセージバスに付随する仮想ネットワークインタフェースコントローラである。メッセージバスを通じて通信する必要のある各ロードバランサ構成要素は少なくとも1つのメッセージバスNICを要求する。メッセージバスNICはメッセージバスへのパイプライン出口の役割と、構成要素へのパイプライン入口の役割を担う。構成要素は複数のメッセージバスネットワークインタフェースを各メッセージバスNICへとインスタンス化することができる。

【0182】

メッセージバスネットワークインタフェースはメッセージバスNICを通じてメッセージバスに付随する構成要素のための機構である。メッセージバスネットワークインタフェースはLinux技術におけるインターフェース構成（ifconfig）インタフェースと同義であってもよく、違いとしては、メッセージバスネットワークインタフェースが物理ネットワークの代わりにメッセージバスに付随する点である。メッセージバスネットワークインタフェースはIPアドレスを有し、メッセージバスNICの最上段にある。メッセージバスネットワークインタフェースは、構成要素によってメッセージバスからのパケットの受信のために用いられることができるパケットソースインタフェースと、構成要素によってメッセージバスへのパケットの発信のために用いられることができるパケットシンクインタフェースを公開する。

【0183】

各ロードバランサノードは、パケットソースおよびパケットシンクインタフェースの実装を通じて伝達され送信される個別のネットワークパケットを処理する。メッセージバス環境における動作時には、これらのインターフェースは、レイヤー2 Ethernet

10

20

30

40

50

ヘッダを追加または削除するメッセージバスネットワークインターフェースによって実装される。(ロードバランサノードのためには、これはカーネルネットワークスタックによって実行されることになる)。図29に示す本番環境においては、パケットソースおよびパケットシンクインターフェースの実装は、実際のネットワークインターフェースにおいてパケットを受信し伝達する。図30で示されるメッセージバス環境において、パケットソースおよびパケットシンクインターフェースの実装がメッセージバスレイヤーまたはレイヤーからパケットを受信し、メッセージバスレイヤーまたはレイヤーへとパケットを伝達する。

【0184】

単純化のために、メッセージバスNICおよびメッセージバスインターフェースはメッセージバスパケットアダプタ、または単にパケットアダプタと総称されてもよい。例、図31および32を参照。

10

【0185】

図29は少なくともいくつかの実施形態による、本番環境において分散型ロードバランサ700を含むロードバランスシステムを示す。ロードバランサ700はこの説明では単純化されている。ロードバランサ700は外部ネットワーク740上の顧客742へと、ロードバランサ700を実装するデータセンター等のネットワークインストールの境界ルータ702を通じて接続してもよい。ロードバランサ700はいくつかの種類の構成要素を含む - 少なくとも1つのエッジルータ704、2つ以上のロードバランサ(LB)ノード710、それぞれが個別のサーバノード(図示せず)上で実装された2つ以上のロードバランサ(LB)モジュール732、ファブリック720を形成するルータやスイッチのような1つまたは複数のネットワーク構成要素、また、少なくともいくつかの実施形態において構成サービス722。少なくともいくつかの実施形態において、ロードバランサ700の各構成要素は、ラック搭載型のコモディティコンピューティング装置等の個別のコンピューティング装置上で実装されてもよい。

20

【0186】

図30は少なくともいくつかの実施形態による、複数の分散型ロードバランスシステムの構成要素を単一処理内または単一処理として構成させて実行させることができるメッセージバス機構を組み込む分散型ロードバランサテストシステム800を示す。図29に示されるロードバランサ700において、各ロードバランサソフトウェア構成要素は、個別のコンピューティング装置上でインストールされ、実行される(例、ロードバランサノード710上にロードバランサソフトウェア、また、サーバノード上にロードバランサモジュール732)。これらロードバランサソフトウェア構成要素を単一処理において実行させるため、各ロードバランサソフトウェア構成要素(図30においてロードバランサ(LB)ノード810およびロードバランサ(LB)モジュール832として示される)は、ロードバランサソフトウェア構成要素を出入りするパケットがまた、物理ネットワーク上で送受信される代わりに、メッセージバス機構を通じて遮断され、ルーティングされるように、構成要素のネットワーク接続性を抽出するコードを含んでもよい。

30

【0187】

少なくともいくつかの実施形態において、分散型ロードバランサテストシステム800上では、メッセージバス機構は1つまたは複数の構成要素間の通信のためのTCPスタックを実装しない。その代わりに、メッセージバス機構はホストシステムのオペレーティングシステム(OS)とのインターフェースを行い、ホストシステムのTCPスタックを用いる。少なくともいくつかの実施形態において、メッセージバス機能は、ホストシステムのOSのカーネル(例、Linuxカーネル)と、カーネルの機能であるIPテーブルを通じて、ユーザーレイヤーの下で結びついている。メッセージバス機能はカーネルレベルでIPテーブルに接続され、パケットを遮断し、ルーティングのためのメッセージバス処理へとパケットを送信する。

40

【0188】

図30において模擬エッジルータ862および模擬ファブリック864で示されるよう

50

に、物理ネットワーク構成要素（例、図29のエッジルータ704およびファブリック720）の機能は、顧客860、サーバ834、構成サービス866も可能であるように、ソフトウェアにおいてシミュレーションされてもよい。しかし、少なくともいくつかの実施形態において、模擬サーバ834ではなく実物が分散型ロードバランサテストシステム800において用いられてもよいことに留意する。図30のメッセージバスレイヤー850が物理ネットワークインフラストラクチャの代わりとなる。したがって、ロードバランサソフトウェア構成要素（ロードバランサノード810およびロードバランサモジュール832）は、図29で示すように本番ネットワーク環境で実行していないことを認識しないままロードバランサテストシステム800上で動作してもよい。

【0189】

いくつかの構成要素（例えば模擬ルータ）は、ネットワークセグメントをシミュレーションする異なるメッセージバスレイヤー850とパケットの送受信を行うために、2つ以上のメッセージバスレイヤー850に接続されてもよい。

【0190】

分散型ロードバランサテストシステム800のメッセージバスレイヤー850において実装されるメッセージバス機構は、ネットワークセグメントの「ワイヤ」をシミュレーションする。少なくともいくつかの実施形態において、メッセージバス機構は構成要素のMACアドレスに基づいて、パケットを分散型ロードバランサテストシステム800内の宛先構成要素に伝達する。こうして、各ロードバランサソフトウェア構成要素（ロードバランサノード810およびロードバランサモジュール832）は、ロードバランサソフトウェア構成要素が、分散型ロードバランサテストシステム800において他の構成要素から送信されたパケットを受信できるように、MACアドレスを1つまたは複数の接続されているメッセージバスレイヤー850に提供する。

メッセージバスパケットアダプタ

【0191】

図31および32は少なくともいくつかの実施形態による、メッセージバスパケットアダプタを示す。少なくともいくつかの実施形態において、各ロードバランサ（LB）ソフトウェア構成要素はパケットソースおよびパケットシンクインターフェースの実装を通じて伝達され送信される個々のネットワークデータパケットを処理する。図31を参照すると、これらインターフェース（パケットソースインターフェース862およびパケットシンクインターフェース864として示される）は分散型ロードバランサテストシステム800上で動作している際に、メッセージバスレイヤー850と、カーネルネットワークスタックによって実行されることになるソフトウェア構成要素870のためのレイヤー2 Ethernetヘッダを追加または削除するロードバランサソフトウェア構成要素870との間のパケットアダプタ860によって実装されてもよい。図29で示されるような本番環境では、ロードバランサソフトウェア構成要素のためのパケットソースおよびパケットシンクの実装が、構成要素が実装される物理装置の実際のネットワークインターフェース上で、パケットを受信し伝達する。

【0192】

図31を参照すると、少なくともいくつかの実施形態において、ロードバランサソフトウェア構成要素870がパケットを伝達する際に、パケットを構成要素の入力キューに追加することで最終的にパケットを宛先構成要素に伝達するために、パケットシンクインターフェース864の送信パケット方法と呼び出す実行スレッドがパケットアダプタ860内とメッセージバスレイヤー850内の関数チェーンをトラバースする。少なくともいくつかの実施形態において、ロードバランサソフトウェア構成要素870がパケットを受信する際に、ロードバランサソフトウェア構成要素870がパケットソースインターフェース862の受信パケット方法と呼び出し、その入力キューからパケットを読み取る。少なくともいくつかの実施形態において、メッセージバス機構はパケットの伝達のためにいかなる追加のスレッドも要求しない。

メッセージバスパケットパイプライン

10

20

30

40

50

【 0 1 9 3 】

図 3 2 を参照すると、少なくともいくつかの実施形態において、パケットソースインターフェース 8 6 2 およびパケットシンクインターフェース 8 6 4 のメッセージバス 8 5 0 側がパケットパイプラインの特徴を提供する。ロードバランサソフトウェア構成要素 8 7 0 がパケットシンクインターフェース 8 6 4 を通じてパケットを送信する際、パケットデータはメッセージバスレイヤー 8 5 0 に達する前に、段階のシリーズ（パケットパイプライン 8 8 0 ）をトラバースしてもよい。これらの段階はパケットを修正し、パケットを破棄し、パケットを複製し、パケットを遅延させる等してもよい。パケットがパケットパイプライン 8 8 0 をトラバースしメッセージバスレイヤー 8 5 0 が宛先構成要素 8 7 0 を選択すると、パケットが宛先構成要素 8 7 0 の入力キューに追加される前に、宛先構成要素 8 7 0 に関連するパイプライン段階の第 2 のシリーズ（パケットパイプライン 8 8 2 ）もまたトラバースされてもよい。

10

プロバイダネットワーク環境の実施例

【 0 1 9 4 】

このセクションでは、分散型ロードバランサ方法および機器の実施形態が実装されてもよいプロバイダネットワーク環境の実施例を説明する。しかし、これらプロバイダネットワーク環境の実施例は制限を意図するものではない。

【 0 1 9 5 】

図 3 3 A は少なくともいくつかの実施形態による、プロバイダネットワーク環境の実施例を示す。プロバイダネットワーク 1 9 0 0 は、顧客に仮想化リソースのインスタンス 1 9 1 2 のアクセス、購入、レンタル、またはその他の取得を可能にする、1 つまたは複数の仮想化サービス 1 9 1 0 を通じて顧客にリソース仮想化を提供してもよい。仮想化リソースのインスタンス 1 9 1 2 には、1 つまたは複数のデータセンター内のプロバイダネットワークまたはネットワーク内の装置に実装された計算およびストレージリソースが含まれるが、それらに限定されない。プライベート IP アドレス 1 9 1 6 はリソースのインスタンス 1 9 1 2 に関連してもよい。プライベート IP アドレスは、プロバイダネットワーク 1 9 0 0 上のリソースのインスタンス 1 9 1 2 の内部ネットワークアドレスである。いくつかの実施形態においては、プロバイダネットワーク 1 9 0 0 はまた、顧客がプロバイダ 1 9 0 0 から取得してもよい、パブリック IP アドレス 1 9 1 4 および / またはパブリック IP アドレスの範囲（例、インターネット Protocol バージョン 4 (IPv 4) またはインターネット Protocol バージョン 6 (IPv 6) アドレス）を提供してもよい。

20

30

【 0 1 9 6 】

従来から、仮想化サービス 1 9 1 0 を通じたプロバイダネットワーク 1 9 0 0 により、サービスプロバイダの顧客（例、顧客ネットワーク 1 9 5 0 A を運用する顧客）が顧客に割り当てられた少なくともいくつかのパブリック IP アドレス 1 9 1 4 を顧客に割り当てられた特定のリソースのインスタンス 1 9 1 2 と動的に結びつけることができる。また、プロバイダネットワーク 1 9 0 0 により顧客は、顧客に割り当て済みのある仮想化されたコンピューティングリソースのインスタンス 1 9 1 2 にすでにマッピングされていたパブリック IP アドレス 1 9 1 4 を、同様に顧客に割り当てられた別の仮想化されたコンピューティングリソースのインスタンス 1 9 1 2 へと再度マッピングすることができる。サービスプロバイダから提供された、仮想化されたコンピューティングリソースのインスタンス 1 9 1 2 およびパブリック IP アドレス 1 9 1 4 を用いて、顧客ネットワーク 1 9 5 0 A の運用者などのサービスプロバイダの顧客は例えば、インターネットのような中間ネットワーク 1 9 4 0 上で、顧客特有のアプリケーションを実装し、顧客のアプリケーションを提示してもよい。中間ネットワーク 1 9 4 0 上の他のネットワークエンティティ 1 9 2 0 はその後顧客ネットワーク 1 9 5 0 A によって公開された宛先パブリック IP アドレス 1 9 1 4 へのトラフィックを生成してもよい。トラフィックはサービスプロバイダデータセンターへとルーティングされ、データセンターで、ネットワーク基盤を通じて、現在宛先パブリック IP アドレス 1 9 1 4 にマッピングされている仮想化されたコンピューティ

40

50

ングリソースのインスタンス1912のプライベートIPアドレス1916へとルーティングされる。同様に、仮想化されたコンピューティングリソースのインスタンス1912からの応答トラフィックはネットワーク基盤を通じて中間ネットワーク1940に戻り、ソースエンティティ1920へとルーティングされてもよい。

【0197】

プライベートIPアドレスは、本明細書で用いられている通り、プロバイダネットワーク内のリソースのインスタンスの内部ネットワークアドレスを参照する。プライベートIPアドレスはプロバイダネットワーク内でのみルーティング可能である。外部のプロバイダネットワークから始まるネットワークトラフィックはプライベートIPアドレスに直接ルーティングされない。代わりにトラフィックは、リソースのインスタンスへとマッピングされたパブリックIPアドレスを用いる。プロバイダネットワークは、パブリックIPアドレスからプライベートIPアドレスへのマッピングやその逆を行うための、ネットワークアドレス変換(NAT)または類似の機能を提供するネットワーク装置または機器を含んでもよい。

10

【0198】

パブリックIPアドレスは、本明細書で用いられている通り、サービスプロバイダまたは顧客のいずれかによってリソースのインスタンスに割り当てられた、インターネット上でルーティング可能なネットワークアドレスである。パブリックIPアドレスへとルーティングされたトラフィックは例えば1:1ネットワークアドレス変換(NAT)を通じて変換され、リソースのインスタンスのそれぞれのプライベートIPアドレスへと転送される。

20

【0199】

いくつかのパブリックIPアドレスは、プロバイダネットワークインフラストラクチャから特定のリソースのインスタンスへと割り当てられてもよい。これらのパブリックIPアドレスは標準パブリックIPアドレス、または単に標準IPアドレスと称されてもよい。少なくともいくつかの実施形態において、標準IPアドレスのリソースのインスタンスのプライベートIPアドレスへのマッピングは、すべてのリソースのインスタンスタイプ向けのデフォルトの起動構成である。

【0200】

少なくともいくつかのパブリックIPアドレスは、プロバイダネットワーク1900の顧客に割り振られるか、またはそのような顧客により取得されてもよい。顧客はその後割り振られたパブリックIPアドレスを、顧客に割り振られた特定のリソースのインスタンスに割り当ててもよい。これらのパブリックIPアドレスは、顧客パブリックIPアドレス、または単に顧客のIPアドレスと称されてもよい。標準IPアドレスと同様にプロバイダネットワーク1900によりリソースのインスタンスに割り当てられる代わりに、顧客のIPアドレスは、例えばサービスプロバイダから提供されたAPIを通じて、顧客によりリソースのインスタンスに割り当てられてもよい。標準IPアドレスとは異なり、顧客のIPアドレスは顧客アカウントに割り当てられるものであり、必要に応じてまたは希望があれば、他のリソースのインスタンスへのそれぞれの顧客による再度マッピングすることが可能である。顧客のIPアドレスは顧客のアカウントに関連し、特定のリソースのインスタンスには関連せず、顧客がそのIPアドレスのリリースを選択するまでは、自身でそれを管理する。従来の固定IPアドレスとは異なり、顧客のIPアドレスにより、顧客のパブリックIPアドレスを顧客のアカウントに関連するいずれかのリソースのインスタンスへと再度マッピングすることで、顧客はリソースのインスタンスまたはアベイラビリティゾーンの問題をマスクすることができる。顧客のIPアドレスにより、例えば、顧客のリソースのインスタンスまたはソフトウェアの問題について顧客のIPアドレスを置き換えのリソースのインスタンスへと再度マッピングすることで、解決に向け動くことができる。

30

40

【0201】

図33Bは少なくともいくつかの実施形態による、図33Aに示されるような、プロバ

50

イダネットワーク環境の実施例における分散型ロードバランサの実装を示す。プロバイダネットワーク1900は、例えば仮想ストレージサービスのようなサービス1910を顧客1960に提供してもよい。顧客1960は、例えばサービス1910への1つまたは複数のAPIを通じて、サービス1910にアクセスし、プロバイダネットワーク1900の本番ネットワーク部分における複数のサーバノード1990上に実装されたリソース（例、ストレージリソースまたは計算リソース）の利用形態を得てもよい。サーバノード1990はそれぞれ、ローカルロードバランサ（LB）モジュール1992と同様に、サーバ（図示せず）、例えばウェブサーバまたはアプリケーションサーバを実装してもよい。1つまたは複数の分散型ロードバランサ1980は、境界ネットワークと本番ネットワークとの間のロードバランサレイヤーにおいて実装されてもよい。境界ルータ1970は、顧客1960からのパケットフロー上のパケット（例、TCPデータパケット）をインターネット等の中間ネットワーク1940を通じて受信し、パケットを1つまたは複数の分散型ロードバランサ1980の1つまたは複数のエッジルータへと境界ネットワークを通じて転送してもよい。データパケットは、1つまたは複数の分散型ロードバランサ1980の1つまたは複数のエッジルータによって公開されたパブリックIPアドレスを対象としてもよい。各分散型ロードバランサ1980のエッジルータは、それぞれの分散型ロードバランサ1980のロードバランサノード間にパケットフローを分散させてもよい。少なくともいくつかの実施形態において、入口ノードの役割を担う各ロードバランサノードはエッジルータに同一のパブリックIPアドレスを提供し、エッジルータはフローごとにハッシュ化されたマルチパスルーティング技術、例えば等価コストマルチパス（ECMP）ハッシュ技術に従って、顧客1960からのパケットフローを入口サーバ間に分散させる。パケットフローのための対象のサーバノード1990を決定するため、また、サーバと顧客1960との間の接続を円滑化するために、ロードバランサノードは本明細書に記載の接続プロトコルを用いてもよい。接続が確立されると、フロートラッカーノードが接続のための状態を維持する一方で、入口ノードがフローのために受信されたパケットをカプセル化して本番ネットワーク上の対象のサーバノード1990へと発信する。サーバノード1990上のロードバランサモジュール1992は、サーバノード1960上のそれぞれのサーバが接続を受け入れるかどうかについて決断を下してもよい。ロードバランサモジュールは入口ノードからのパケットを受信して脱カプセル化し、脱カプセル化されたデータパケット（例、TCPデータパケット）をサーバノード1990上のそれぞれのサーバへと送信する。ロードバランサモジュール1992はまた、パケットフローのための出口ノードとしてロードバランサノードを選択し、フローのための発信パケットをカプセル化し、選択した出口ノードに本番ネットワークを通じて送信してもよい。次に出口ノードがパケットを脱カプセル化し、脱カプセル化されたパケットを、それぞれの顧客1960への伝達のための境界ネットワークへと送信する。

【0202】

図34Aは少なくともいくつかの実施形態による、分散型ロードバランサおよびサーバノードの物理ラックの実装の実施例を示すが、制限を意図するものではない。少なくともいくつかの実施形態において、分散型ロードバランサ上の様々な構成要素が、ラック搭載型のコモディティコンピューティング装置上で、またはラック搭載型のコモディティコンピューティング装置として実装されてもよい。ラック190は、それぞれがロードバランサノードの役割を担う複数のコンピューティング装置（LBノード110A～110F）および、それぞれがサーバノードの役割を担う複数のコンピューティング装置（サーバノード130A～130L）を含んでもよい。ラック190はまた、少なくとも1つのエッジルータ104、ファブリック120を形成する1つまたは複数のラック搭載型ネットワーク装置（ルータ、スイッチ等）、および1つまたは複数の他の構成要素180（他のネットワーク装置、パッチパネル、電源、冷却システム、バス等）を含む。図33Aおよび33Bのプロバイダネットワーク1900を実装するデータセンターまたはセンター等のネットワーク100のインストールは、1つまたは複数のラック190を含んでもよい。

【0203】

10

20

30

40

50

図34Bは少なくともいくつかの実施形態による、分散型ロードバランサおよびサーバノードの物理ラックの実装の別の実施例を示すが、制限を意図するものではない。図34Bはスロット搭載型コンピューティング装置、例えばブレードサーバとして実装された、ラック190上のLBノード110およびサーバノード130を示す。

【0204】

図35は少なくともいくつかの実施形態による、ネットワーク上で1つまたは2つ以上の分散型ロードバランサが実装されてもよく、別途実装されたサーバノードを持つネットワーク環境の実施例を示す。この実施例では、2つの分散型ロードバランサ1980Aおよび1980Bが示される。分散型ロードバランサ1980はそれぞれ境界ネットワークを通じて顧客1960からのパケットフローを受信し、複数のサーバノード1990間にパケットフローを分散させるために本明細書に記載のロードバランサ方法を実行してもよい。いくつかの実装においては、各分散型ロードバランサ1980は図34Aおよび34Bで示されるラック190に類似のラック実装であってもよいが、ロードバランサラック内にインストールされたサーバノードは含まれない。サーバノード1990は、データセンター内の1つまたは複数の個別のラック上にインストールされたブレードサーバ等のラック搭載型コンピューティング装置であってもよい。いくつかの実装においては、サーバノード1990は、プロバイダネットワークにより提供され、それぞれが異なる1つまたは複数の分散型ロードバランサ1980によりフロントに配置された、2つ以上の異なるサービスを実装してもよい。

例示的システム

【0205】

少なくともいくつかの実施形態において、本明細書に記載の分散型ロードバランサ方法および機器の一部またはすべてを実装するサーバは、図36に示されるコンピュータシステム2000のような、1つまたは複数のコンピュータアクセス可能な媒体を含む、またはそのような媒体にアクセスするよう構成されている汎用コンピュータシステムを含んでもよい。示された実施形態においては、コンピュータシステム2000は、入力/出力(I/O)インターフェース2030を通じてシステムメモリ2020に接続された1つまたは複数のプロセッサ2010を含む。コンピュータシステム2000は、I/Oインターフェース2030に接続されたネットワークインターフェース2040をさらに含む。

【0206】

様々な実施形態において、コンピュータシステム2000は、1つのプロセッサ2010を含むユニプロセッサシステム、または、いくつかのプロセッサ2010(例、2つ、4つ、8つ、または別の適切な数)を含むマルチプロセッサシステムであってもよい。プロセッサ2010は、命令を実行する能力があるいずれかの適切なプロセッサであってもよい。例えば様々な実施形態において、プロセッサ2010は、x86、PowerPC、SPARC、またはMIPSISA等の各種の命令セットアーキテクチャ(ISA)、または他のいずれかの適切なISAを実装する汎用または埋め込みプロセッサであってもよい。マルチプロセッサシステムにおいては、プロセッサ2010のそれぞれが一般に同一のISAを実装してもよいが、必ずしもそうである必要はない。

【0207】

システムメモリ2020は、1つまたは複数のプロセッサ2010によりアクセス可能な命令およびデータを格納するよう構成されてもよい。様々な実施形態において、システムメモリ2020は、スタティックランダムアクセスメモリ(SRAM)、シンクロナスダイナミックRAM(SDRAM)、不揮発性/フラッシュ型メモリ、またはその他各種メモリ等のいずれかの適切なメモリ技術を用いて実装されてもよい。示された実施形態においては、ロードバランサ方法および機器のために上述された方法、技術、およびデータ等の、1つまたは複数の所望の関数を実装するプログラム命令およびデータが、システムメモリ2020内にコード2024およびデータ2026として格納されていることが示される。

【0208】

1つの実施形態において、I/Oインターフェース2030は、プロセッサ2010、システムメモリ2020、およびネットワークインターフェース2040または他の周辺インターフェースを含む装置内のいずれかの周辺装置の間のI/Oトラフィックを調整するように構成されてもよい。いくつかの実施形態においては、1つの構成要素(例、システムメモリ2020)からのデータ信号を、別の構成要素(例、プロセッサ2010)による利用に適したフォーマットへと変換するために、I/Oインターフェース2030はかかる必要なプロトコル、タイミングまたは他のデータ媒体変換を行ってもよい。いくつかの実施形態においては、I/Oインターフェース2030は例えば、周辺構成要素相互接続(PCI)バス標準またはユニバーサルシリアルバス(USB)標準のバリエーション等

10

の様々な種類の周辺バスに付随する装置のためのサポートを含んでもよい。いくつかの実施形態においては、I/Oインターフェース2030の関数は、例えばノースブリッジおよびサウスブリッジ等の、2つ以上の個別の構成要素へと分割されてもよい。また、いくつかの実施形態においては、システムメモリ2020へのインターフェース等の、I/Oインターフェース2030の一部またはすべての機能がプロセッサ2010に直接組み込まれてもよい。

【0209】

ネットワークインターフェース2040は、コンピュータシステム2000と、例えば図1~35で示すような他のコンピュータシステムまたは装置等の、1つまたは複数のネットワーク2050に付随する他の装置2060との間でのデータ交換が可能であるように構成されてもよい。様々な実施形態において、ネットワークインターフェース2040

20

は、例えばEthernetネットワークのような種類の、いずれかの適切な有線または無線の一般的なデータネットワークを通じた通信をサポートしてもよい。また、ネットワークインターフェース2040は、アナログ音声ネットワークまたはデジタルファイバ通信ネットワークのような遠隔通信/電話網を通じた通信、Fibre Channel SAN等のストレージエリアネットワークを通じた通信、またはその他いずれかの適切な種類のネットワークおよび/またはプロトコルを通じた通信をサポートしてもよい。

【0210】

いくつかの実施形態においてシステムメモリ2020は、分散型ロードバランスシステムの実施形態を実装するための図1~35について上述したように、プログラム命令およびデータを格納するよう構成されたコンピュータアクセス可能な媒体の1つの実施形態

30

であってもよい。しかし他の実施形態においては、異なる種類のコンピュータアクセス可能な媒体において、プログラム命令および/またはデータが受信され、送信され、または格納されてもよい。一般的に、コンピュータアクセス可能な媒体は、磁気または光学式媒体等の非一時的記憶媒体またはメモリ媒体を含んでもよい。例、I/Oインターフェース2030を通じてコンピュータシステム2000に接続されるディスクまたはDVD/CD。コンピュータアクセス可能な非一時的記憶媒体はまた、コンピュータシステム2000のいくつかの実施形態においてシステムメモリ2020または別の種類のメモリとして含まれてもよいRAM(例えば、SDRAM、DDRSDRAM、RDRAM、SRAM等)、ROM等のあらゆる揮発性または不揮発性媒体を含んでもよい。さらにコンピュータ

40

アクセス可能な媒体は、ネットワークインターフェース2040を通じて実装されてもよいネットワークおよび/または無線接続等の通信媒体を通じて伝達される電気、電磁、またはデジタル信号等の伝送媒体または信号を含んでもよい。

【0211】

本開示の実施形態は、以下の節を考慮して説明することができる。

1. 複数のロードバランサノードの少なくとも2つが入口サーバとして構成され、前記複数のロードバランサノードの少なくとも2つがフロートラッカーノードとして構成される、

前記複数のロードバランサノードと、
複数のサーバノードと、

1つまたは複数の顧客からのパケットフローを、ハッシュ化されたマルチパスルーティ

50

ング技術に従って、前記入口サーバへと分散させるよう構成されたルータと、
を備えた分散型ロードバランサシステムであり、
顧客のためのパケットフローにおけるパケットを前記ルータから受信し、
前記複数のサーバノードへの前記パケットフローのためのマッピングを前記入口サーバ
が有しないことを決定し、
前記パケットのソースおよび宛先アドレス情報に適用される一貫したハッシュ関数に従
って、前記パケットフローのための少なくとも1つのフロートラッカーノードを決定し、
前記パケットフローのための前記複数のサーバノードの特定の1つへの接続のマッピ
ングを、少なくとも1つのフロートラッカーノードから取得し、
前記特定のサーバノードへの前記パケットフローの1つ以上のパケットを送信する 10
ように各入口サーバが構成される、
分散型ロードバランサシステム。
2．前記パケットフローがトランスミッションコントロールプロトコル(TCP)パケ
ットフローである、第1節に記載の分散型ロードバランサシステム。
3．前記複数のロードバランサノードの少なくとも2つが、前記1つまたは複数の顧客
への前記サーバノードからの発信パケットを送信するよう構成された出口サーバとして構
成され、
前記パケットフローのための前記出口サーバを選択し、
前記パケットフローのための1つまたは複数の発信パケットを、前記選択した出口サー
バへと送信する、 20
ように前記サーバノードが構成され、
前記顧客への前記発信パケットを送信するように、前記出口サーバが構成され、
前記パケットフローのための前記選択された出口サーバが、前記パケットフローのため
の前記入口サーバとは異なるロードバランサノードである、
第1節に記載の分散型ロードバランサシステム。
4．前記サーバノードへの前記パケットの送信の前に、前記入口サーバが前記1つまた
は複数のパケットをユーザデータグラムプロトコル(UDP)に従ってカプセル化し、前
記出口サーバへの前記発信パケットの送信の前に、前記サーバノードが前記発信パケ
ットをUDPに従ってカプセル化し、前記顧客への前記発信パケットの送信の前に、前記出口
サーバが前記発信パケットから前記UDPカプセル封じを取り外す、第3節に記載の分散 30
型ロードバランサシステム。
5．前記パケットフローのための前記出口サーバを選択し、
前記カプセル化された受信パケットを前記入口サーバから受信し、
前記パケットから前記UDPカプセル封じを取り外し、前記パケットを前記サーバノ
ード上のサーバへと伝達させ、
前記サーバノード上の前記サーバから前記発信パケットを取得し、
UDPに従って、前記発信パケットをカプセル化し、
前記カプセル化された発信パケットを前記出口サーバへと送信する
ように構成されるロードバランサモジュールを前記サーバノードが含む、
第4節に記載の分散型ロードバランサシステム。 40
6．前記パケットフローのための前記複数のサーバノードの特定の1つへの接続のマッ
ピングを前記少なくとも1つのフロートラッカーノードから取得するために、
前記パケットフローのための情報を含むメッセージを、前記入口サーバが前記パケ
ットフローのための1次フロートラッカーへと送信し、
前記パケットフローのための前記情報を含むメッセージを、前記1次フロートラッカー
が前記パケットフローのための2次フロートラッカーへと送信し、前記パケットフロー
のための前記1次および2次フロートラッカーが異なるロードバランサノードであり、
前記2次フロートラッカーが、前記パケットフローのための受信確認を、前記顧客へと
送信し、
前記入口サーバが、受信確認パケットを前記顧客から受信し、前記受信確認パケットを 50

前記 1 次フロートラッカーへと転送し、

前記 1 次フロートラッカーが、前記パケットフローを受信するための前記サーバノードとして、前記複数のサーバノードの中から前記特定のサーバノードを無作為に選択し、前記特定のサーバノードを示すメッセージを、前記 2 次フロートラッカーへと送信し、

前記 2 次フロートラッカーが、同期メッセージを生成して前記生成された同期メッセージを前記特定のサーバノードへと送信し、

前記 2 次フロートラッカーが、前記パケットフローのための接続情報を前記特定のサーバノードから受信して前記 1 次フロートラッカーへの前記接続情報を含むメッセージを送信し、

前記 1 次フロートラッカーが、前記パケットフローのための前記接続情報を含むメッセージを前記入口サーバへと送信し、前記接続情報が前記パケットフローを前記特定のサーバノードへとマッピングする、

第 1 節に記載の分散型ロードバランサシステム。

7 . 前記生成された同期メッセージを前記 2 次フロートラッカーから受信し、

前記サーバノード上のサーバが接続を受け入れることができることを決定し、

前記生成された同期メッセージに従って同期パケットを生成し、前記同期パケットを前記サーバノード上の前記サーバへと伝達し、

前記サーバノード上の前記サーバによって生成された受信確認パケットを遮断し、

前記接続情報を含むメッセージを前記 2 次フロートラッカーへと送信する、

ように構成されるロードバランサモジュールを前記サーバノードが含む、

第 6 節に記載の分散型ロードバランサシステム。

8 . 顧客からのパケットフローにおけるパケットの受信であって、1 つまたは複数の顧客から一貫したハッシュ関数に従って前記複数のロードバランサノードへと前記パケットフローを分散させるルータからのパケットの受信、

前記パケットのソースおよび宛先アドレス情報に適用される一貫したハッシュ関数に従って、前記パケットフローのためのフロートラッカーノードとしての役割を担うロードバランサノードの決定、

前記パケットフローのための複数のサーバノードの 1 つへの接続のマッピングの、前記パケットフローのための前記フロートラッカーノードからの取得、

前記マッピングにより示された、前記サーバノードへの前記パケットフローの 1 つまたは複数のパケットの送信、

を、複数のロードバランサノードのひとつの入口サーバによって実行すること、

を備えた方法。

9 . 前記パケットフローがトランスミッションコントロールプロトコル (T C P) パケットフローである、第 8 節に記載の方法。

1 0 . 前記パケットフローの前記 1 つまたは複数のパケットの前記サーバノードへの前記送信の前に、ユーザデータグラムプロトコル (U D P) に従って前記パケットをカプセル化する前記入口サーバをさらに備えた、第 8 節に記載の方法。

1 1 . 前記サーバノードによる、前記パケットフローのための出口サーバとしての前記複数のロードバランサノードの 1 つの選択であり、前記パケットフローのための前記選択した出口サーバが、前記パケットフローのための前記入口サーバとは異なるロードバランサノードである、前記複数のロードバランサノードの 1 つの選択と、

前記サーバノードによる、前記パケットフローのための 1 つまたは複数の発信パケットの、前記選択された出口サーバへの送信と、

前記出口サーバによる前記発信パケットの、前記パケットフローの前記顧客への送信と、

をさらに備えた、第 8 節に記載の方法。

1 2 . 前記発信パケットの前記出口サーバへの前記送信の前に、ユーザデータグラムプロトコル (U D P) に従って前記発信パケットをカプセル化する前記サーバノードと、

前記発信パケットの顧客への前記送信の前に、前記発信パケットから前記 U D P カプセ

10

20

30

40

50

ル封じを取り外す前記出口サーバと、

をさらに備えた、第 11 節に記載の方法。

13. 前記フロートラッカーノードが前記パケットフローのための 1 次フロートラッカーノードであり、前記一貫したハッシュ関数に従った一貫したハッシュリングにおける次のロードバランサノードが前記パケットフローのための 2 次フロートラッカーノードである、第 8 節に記載の方法。

14. 前記入口サーバによる、少なくとも 1 つのメッセージの前記パケットフローのための前記 1 次フロートラッカーノードへの送信であり、各メッセージが前記ルータから受信された前記パケットフローのパケットを含む前記送信と、

前記 1 次フロートラッカーノードによる、前記複数のサーバノードからの前記パケットフローのための前記サーバノードの選択と、

前記 1 次フロートラッカーノードによる、前記選択されたサーバノードを示すパケットフロー情報の前記 2 次フロートラッカーノードへの送信と、

前記 2 次フロートラッカーノードによる、前記サーバノードと前記顧客との通信による、前記パケットフローのための前記選択されたサーバノードへの前記接続の確立の円滑化と、

前記 2 次フロートラッカーノードによる、前記パケットフローのための接続情報の、前記 1 次フロートラッカーノードを通じた前記入口サーバへの送信であり、前記接続情報が前記選択されたサーバノードへの前記パケットフローのマッピングを行う前記送信と、

を、前記パケットフローのための複数のサーバノードの 1 つへの接続のマッピングの、前記パケットフローのための前記フロートラッカーノードからの前記取得が備える、

第 13 節に記載の方法。

15. 前記 2 次フロートラッカーノードによる、前記サーバノード上の前記ロードバランサモジュールへの、生成された同期メッセージの送信と、

前記サーバノード上のサーバが接続を受け入れることができることの決定、

前記生成された同期メッセージに従った同期パケットの生成、

前記同期パケットの、前記サーバノード上の前記サーバへの伝達、

前記サーバノード上の前記サーバにより生成された受信確認パケットの遮断と、

前記接続情報を含むメッセージの、前記 2 次フロートラッカーノードへの送信

を、前記サーバノード上の前記ロードバランサモジュールにより実行すること、

を、前記サーバノードと前記顧客との通信により、前記パケットフローのための前記選択されたサーバノードへの前記接続の確立を円滑化するロードバランサモジュールを前記サーバノードが含む、

第 14 節に記載の方法。

16. 1 つまたは複数の顧客からのパケットフローを一貫したハッシュ関数に従って前記複数のロードバランサノードへと分散させるルータからパケットが受信されるように、顧客のためのパケットフロー内の前記パケットを受信し、

前記パケットのソースおよび宛先アドレス情報に適用される一貫したハッシュ関数に従って、前記パケットフローのためのフロートラッカーノードとしての役割を担う、複数のロードバランサノードの 1 つを決定し、

前記パケットフローのための複数のサーバノードの 1 つへの接続のマッピングを、前記パケットフローのための前記フロートラッカーノードから取得し、

前記パケットフローの 1 つまたは複数のパケットを、前記マッピングにより示された前記サーバノードへと送信する、

ように各入口サーバが構成された、

複数のロードバランサノードのおのおのに入口サーバおよびフロートラッカーを実装するためにコンピュータにより実行可能なプログラム命令を格納するコンピュータアクセス可能な非一時的記憶媒体。

17. 前記パケットフローのための出口サーバとして前記複数のロードバランサノードの 1 つを選択し、前記パケットフローのための選択した出口サーバが前記パケットフロー

10

20

30

40

50

のための入口サーバとは異なるロードバランサノードであり、

前記パケットフローのための1つまたは複数の発信パケットを前記選択された出口サーバへと送信する、

ロードバランサモジュールから受信されたパケットフロー上の発信パケットを、前記パケットフローの顧客へと送信するように各出口サーバが構成される、

ように各ロードバランサモジュールが構成され、

それぞれのロードバランサノード上に出口サーバ、複数の前記サーバノードのそれぞれの上にロードバランサモジュールを実装するために前記プログラム命令がさらにコンピュータにより実行可能である、第16節に記載のコンピュータアクセス可能な非一時的記憶媒体。

10

18. 各入口サーバがさらに、前記パケットフローの1つまたは複数のデータパケットの前記サーバノードへの前記送信の前に、ユーザデータグラムプロトコル(UDP)に従って前記パケットをカプセル化するように構成され、

入口サーバから受信した前記パケットから前記UDPカプセル封じを取り外し、前記パケットを前記それぞれのサーバノード上のサーバへと伝達し、

前記それぞれのサーバノード上の前記サーバからの前記発信パケットを遮断し、

出口サーバへの前記発信パケットの前記送信の前に、UDPに従って前記発信パケットをカプセル化する、

ように各ロードバランサモジュールがさらに構成され、

前記パケットフローの顧客への前記発信パケットの送信の前に、前記発信パケットから前記UDPカプセル封じを取り外すように、前記各出口サーバがさらに構成される、

20

第17節に記載のコンピュータアクセス可能な非一時的記憶媒体。

19. フロートラッカーノードが前記パケットフローのための1次フロートラッカーノードであり、一貫したハッシュ関数に従った一貫したハッシュリング内の次のロードバランサノードが、前記パケットフローのための2次フロートラッカーノードであり、前記パケットフローのため前記フロートラッカーノードから、前記パケットフローのための複数の1つへの接続のマッピングを取得するために、

前記入口サーバが少なくとも1つのメッセージを前記パケットフローのための前記1次フロートラッカーノードへと送信するように構成され、各メッセージが前記ルータから受信された前記パケットフローのパケットを含み、

30

前記パケットフローのための前記サーバノードを前記複数のサーバノードから選択し、

前記選択されたサーバノードを示すパケットフロー情報を、前記2次フロートラッカーノードへと送信する、

ように、前記1次フロートラッカーノードが構成され、

前記サーバノードと前記顧客との間の通信により、前記パケットフローのための前記選択したサーバノードへの前記接続の確立を円滑化し、

前記パケットフローのための接続情報を前記入口サーバへと前記1次フロートラッカーノードを通じて送信し、前記接続情報が前記選択されたサーバノードへの前記パケットフローのマッピングを行う、

ように、前記2次フロートラッカーノードが構成される、

40

第16節に記載のコンピュータアクセス可能な非一時的記憶媒体。

20. 前記プログラム命令がさらにコンピュータにより実行可能であり、複数のサーバノードのそれぞれの上にロードバランサモジュールを実装するため、また、前記サーバと前記顧客との間の通信により、前記パケットフローのための前記選択されたサーバノードへの前記接続の確立を円滑化するために、

前記2次フロートラッカーノードが生成された同期メッセージを前記サーバノード上の前記ロードバランサモジュールへと送信するように構成され、

前記サーバノード上の前記サーバが接続を受け入れることができることを決定し、

前記生成された同期メッセージに従って同期パケットを生成し、

同期パケットを前記サーバノード上の前記サーバへと伝達し、

50

前記サーバノード上の前記サーバにより生成された受信確認パケットを遮断し、前記接続情報を含むメッセージを前記2次フロートラッカーノードに送信する、ように、前記サーバノード上の前記ロードバランサモジュールが構成される、第19節に記載のコンピュータアクセス可能な非一時的記憶媒体。

結論

【0212】

様々な実施形態は、上記のコンピュータアクセス可能な媒体に関する説明に従って実装された命令および/またはデータの受信、送信または格納をさらに含んでもよい。一般的に、コンピュータアクセス可能な媒体は、ネットワークおよび/または無線接続により伝達される通信媒体を介した電気、電磁、またはデジタル信号等の伝送媒体または信号と同様に、磁気または光学式媒体等の記憶媒体またはメモリ媒体、例、ディスクまたはDVD/CD-ROM、RAM(例えば、SDRAM、DDR、RDRAM、SRAM等)、ROM等の揮発性または不揮発性媒体を含んでもよい。

10

【0213】

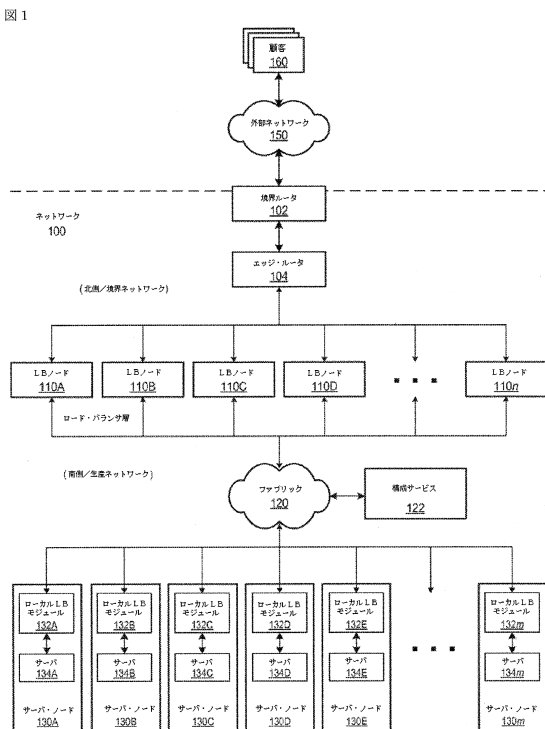
図面に示され、本明細書に記載された様々な方法は、方法の例示的な実施形態を表す。方法はソフトウェア、ハードウェア、またはそれらの組み合わせにおいて実装されてもよい。方法の順序は変更されてもよく、また、様々な要素が追加され、再整理され、組み合わせられ、省略され、修正される等してもよい。

【0214】

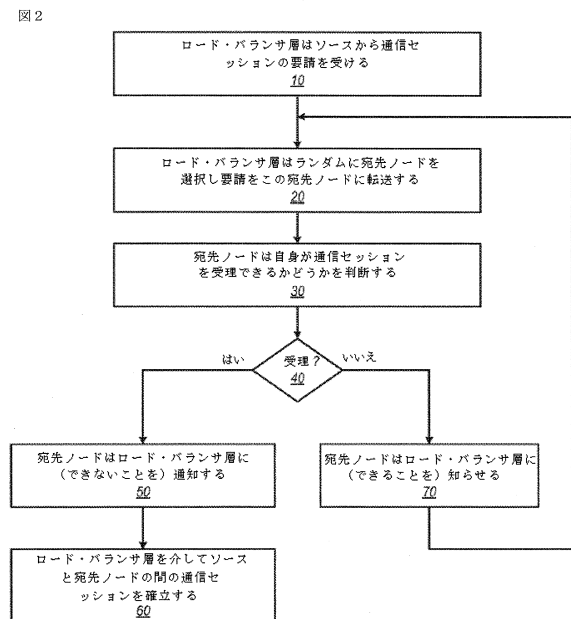
当業者には明らかであるように、本開示を利用して様々な修正や変更が加えられてもよい。本開示はそのようなすべての修正や変更を包含することを意図しており、したがって、上記記載は制限的な意味ではなく説明的な意味を持つと見なされるべきである。

20

【図1】

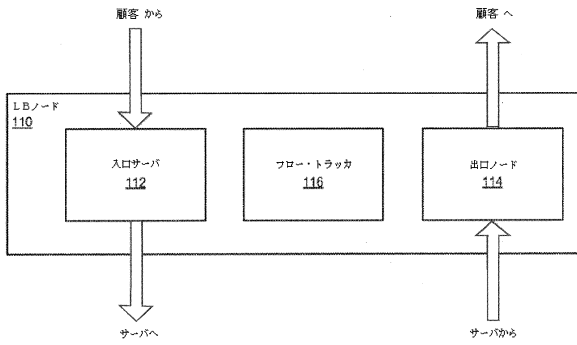


【図2】



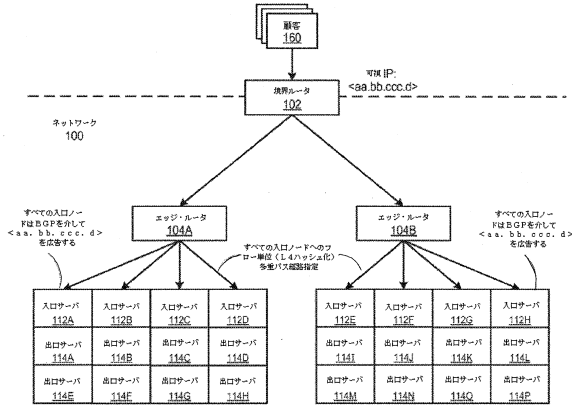
【 図 3 】

図 3



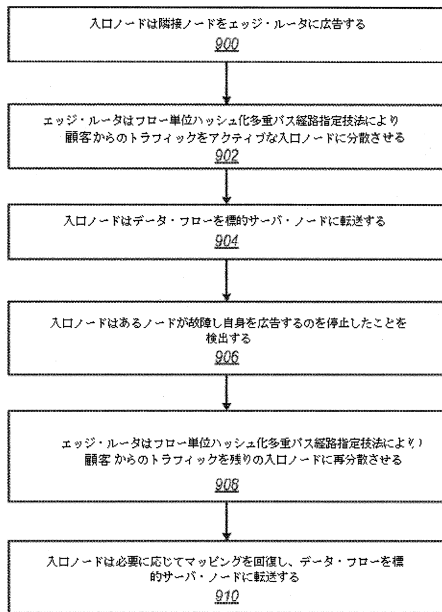
【 図 4 】

図 4



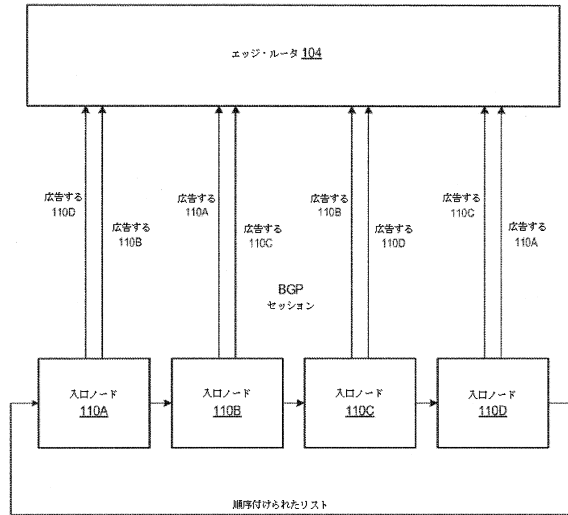
【 図 6 】

図 6



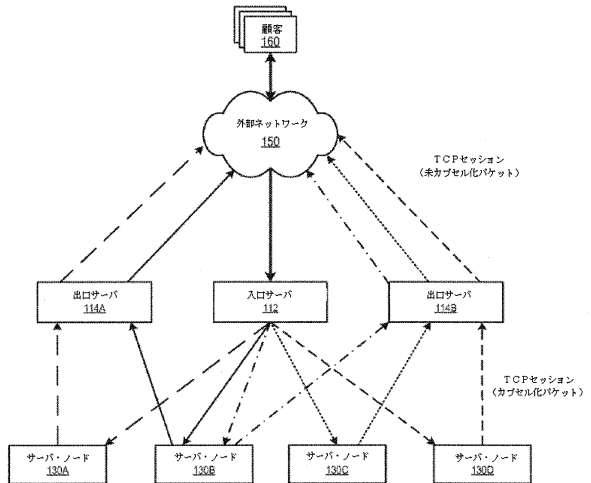
【 図 5 】

図 5

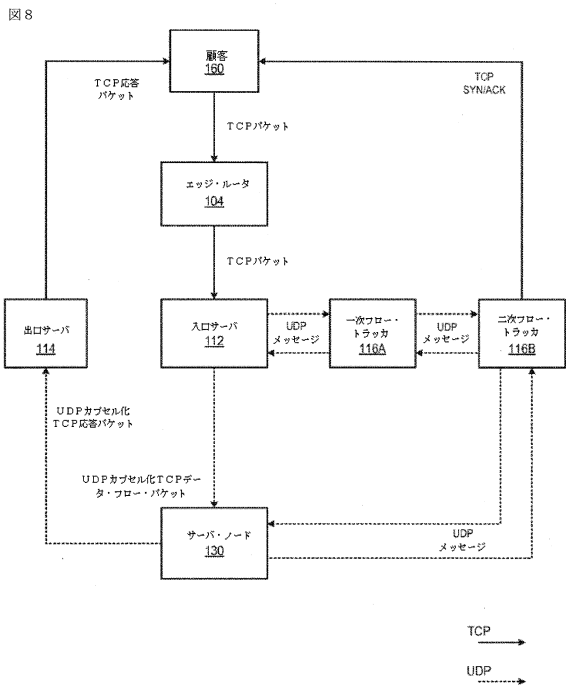


【 図 7 】

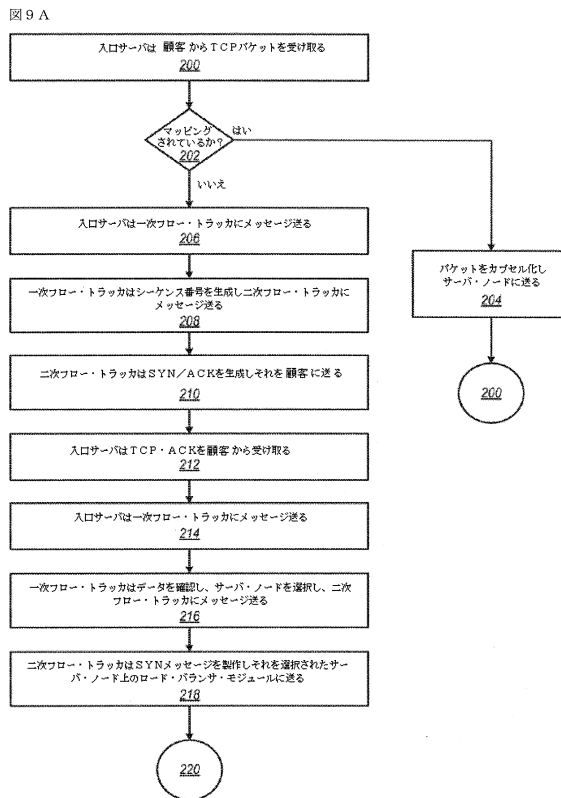
図 7



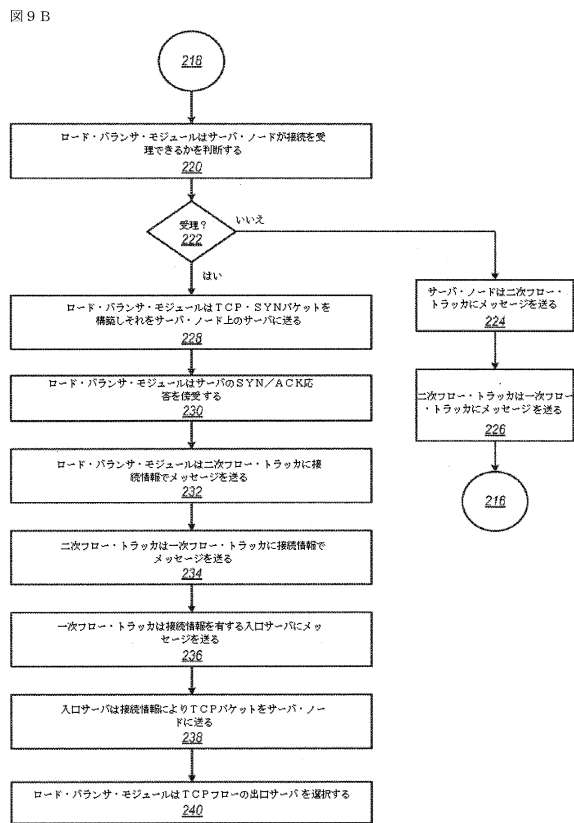
【 図 8 】



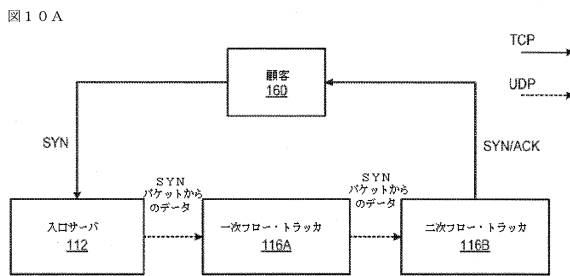
【 図 9 A 】



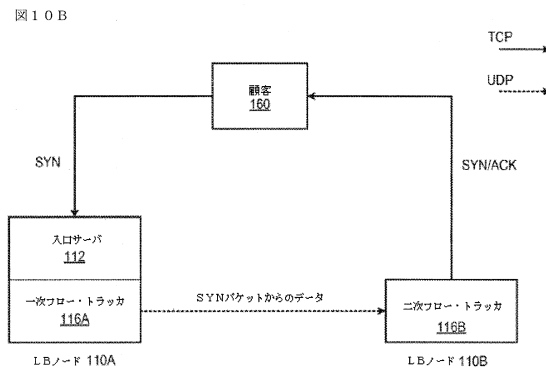
【 図 9 B 】



【 図 10 A 】

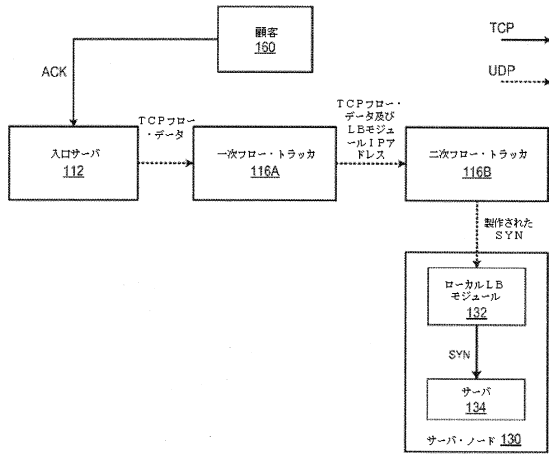


【 図 10 B 】



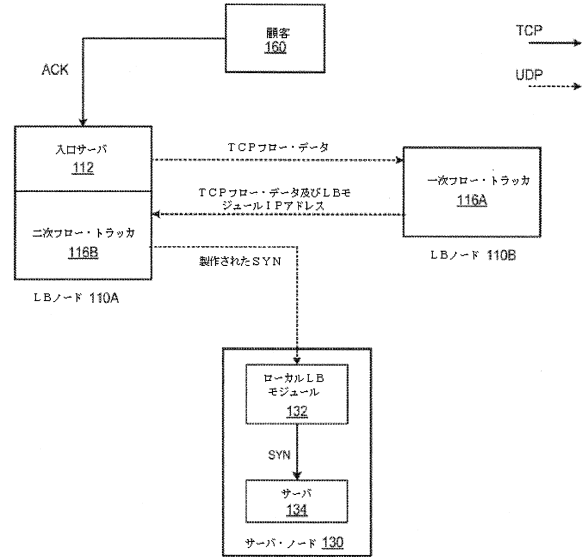
【図10C】

図10C



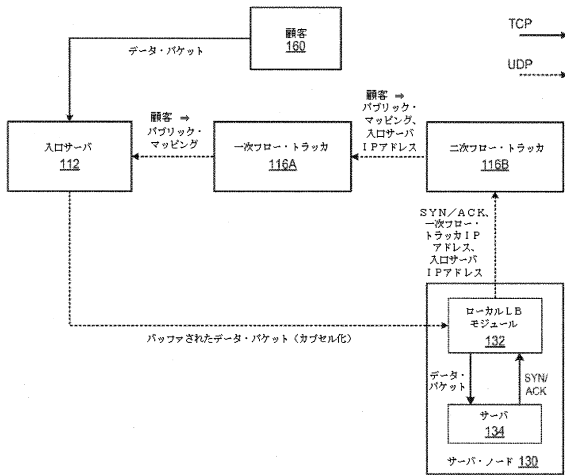
【図10D】

図10D



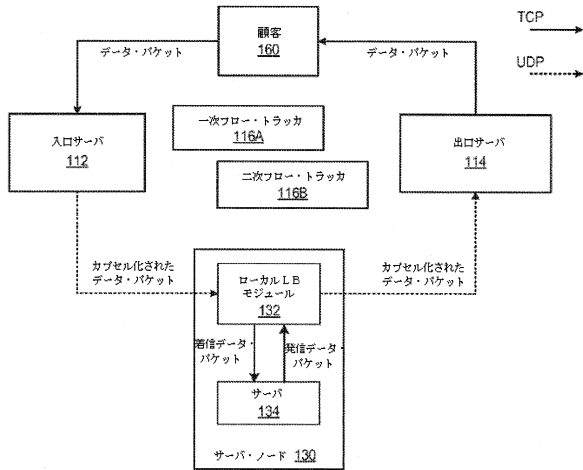
【図10E】

図10E



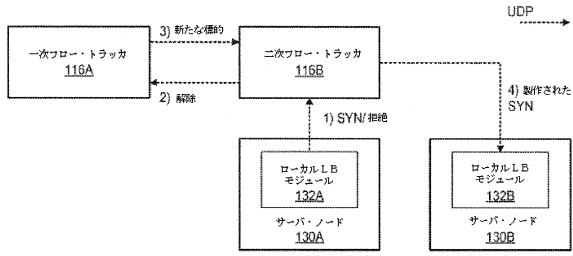
【図10F】

図10F



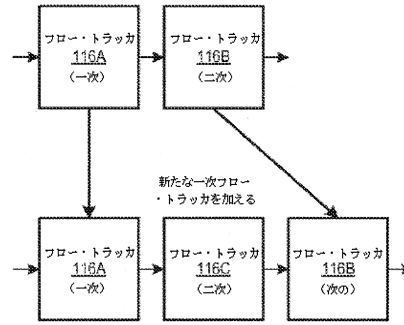
【図10G】

図10G



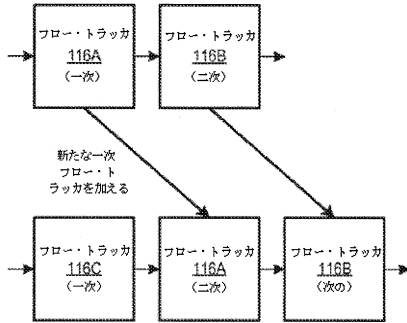
【図11B】

図11B



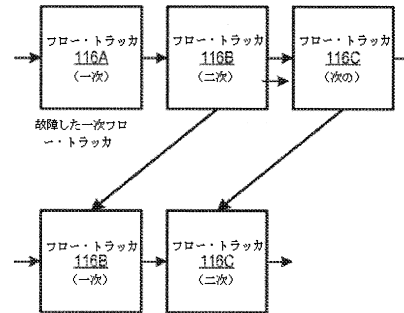
【図11A】

図11A



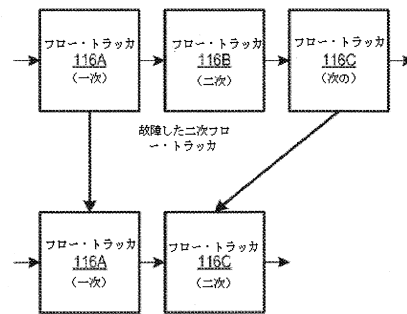
【図11C】

図11C



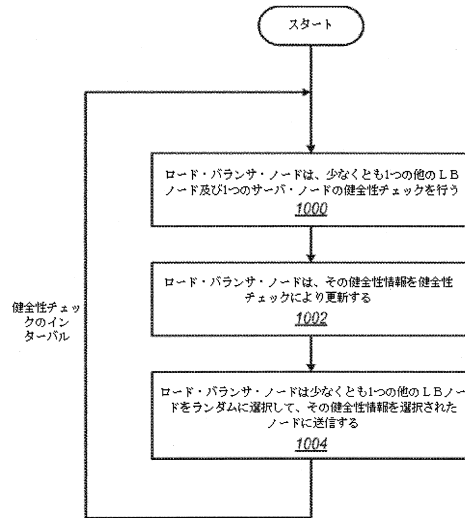
【図11D】

図11D



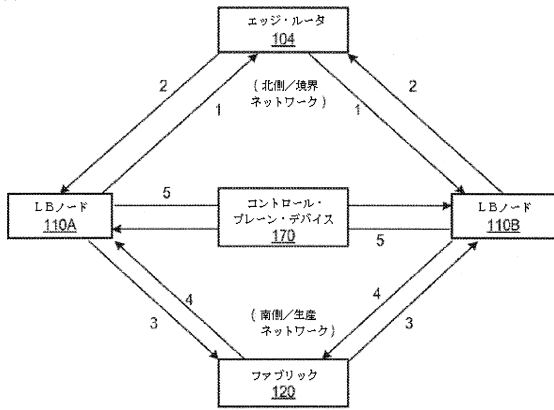
【図12】

図12



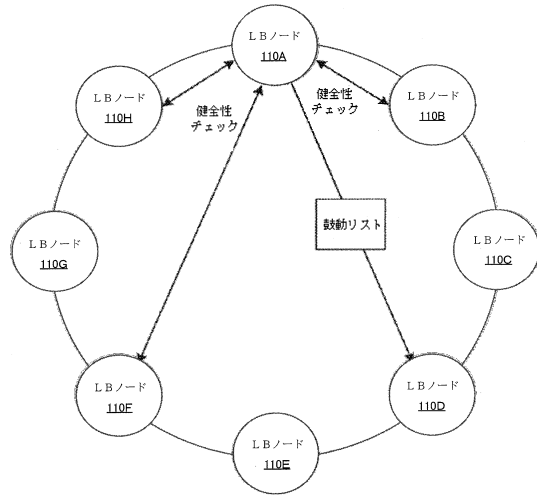
【 図 13 】

図 13



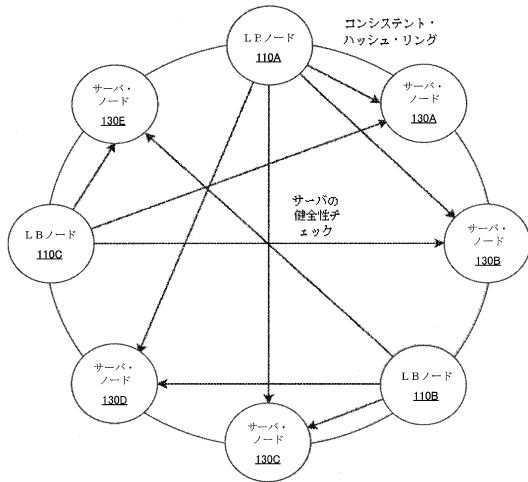
【 図 14 】

図 14



【 図 15 】

図 15



【 図 17 】

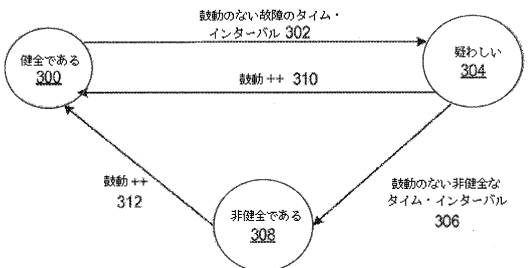
図 17

ロード・バランサ・ノード 110

健全な LB ノードのリスト 320 <ノード A> <ノード B> <ノード C> <ノード D> ...	新しい LB ノードのリスト 322 <ノード E> <ノード F> ...	非健全な LB ノードのリスト 324 <ノード G> ...	LB ノード鼓動のリスト 326 <ノード A、カウンタ、時間> <ノード B、カウンタ、時間> <ノード C、カウンタ、時間> <ノード D、カウンタ、時間> <ノード E、カウンタ、時間> <ノード F、カウンタ、時間> <ノード G、カウンタ、時間> ...
--	--	--	--

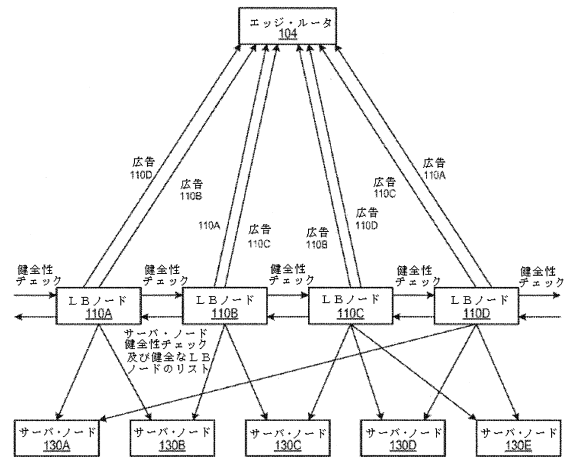
【 図 16 】

図 16



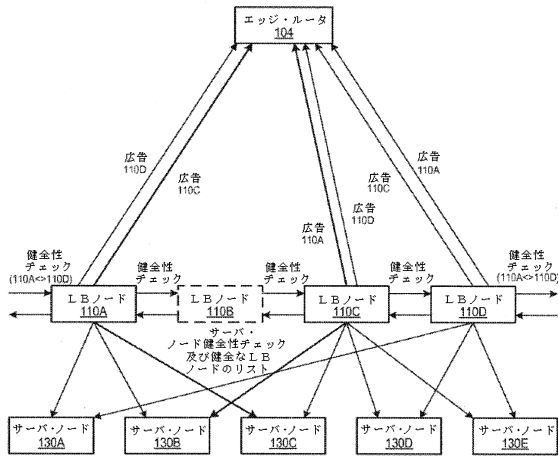
【 図 18 A 】

図 18 A



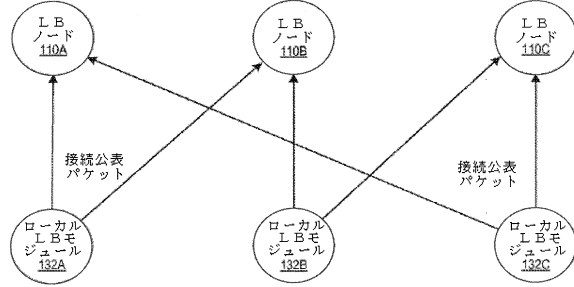
【図18B】

図18B



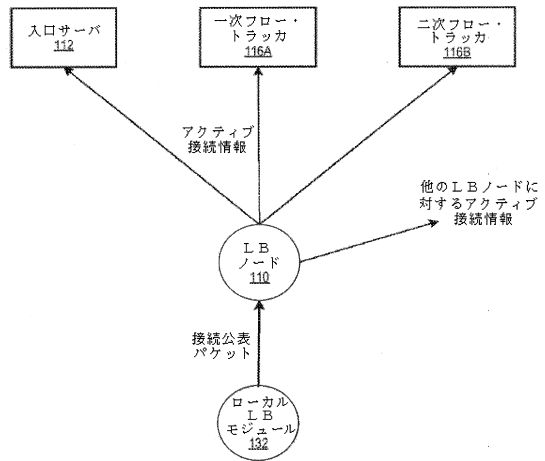
【図19A】

図19A



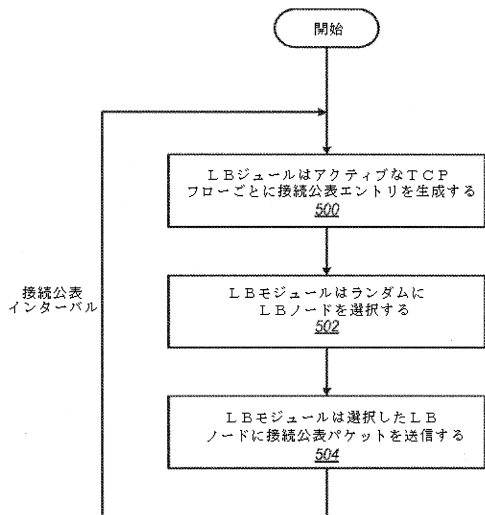
【図19B】

図19B



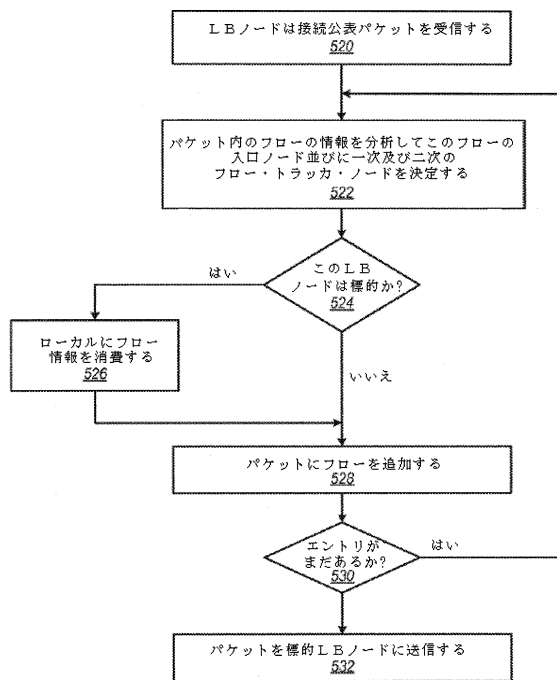
【図20】

図20

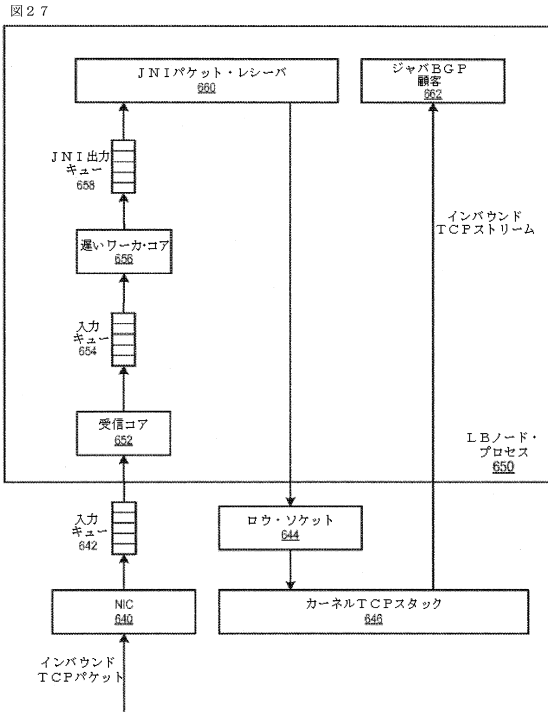


【図21】

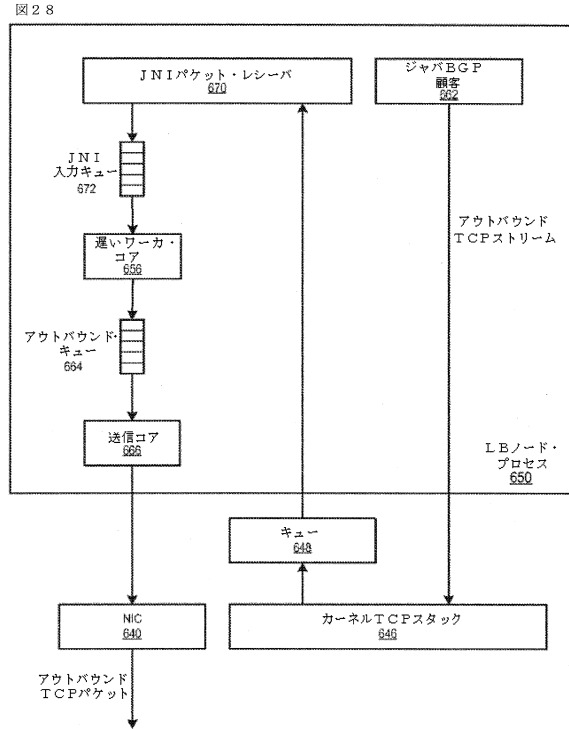
図21



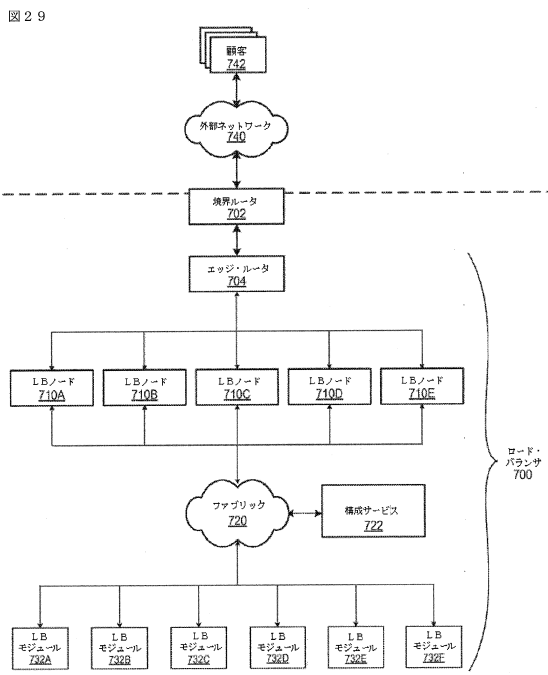
【 図 27 】



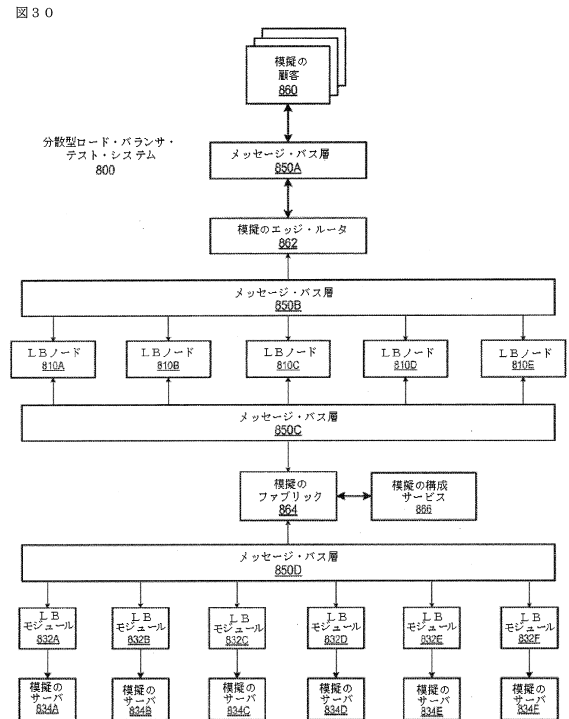
【 図 28 】



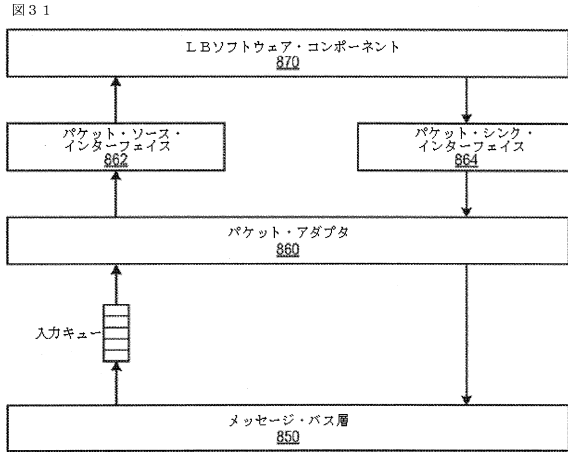
【 図 29 】



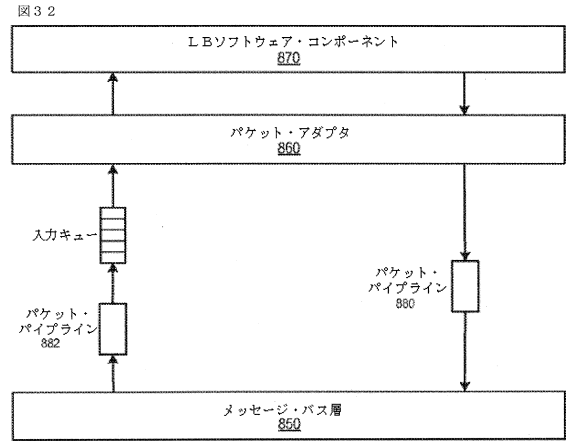
【 図 30 】



【図31】

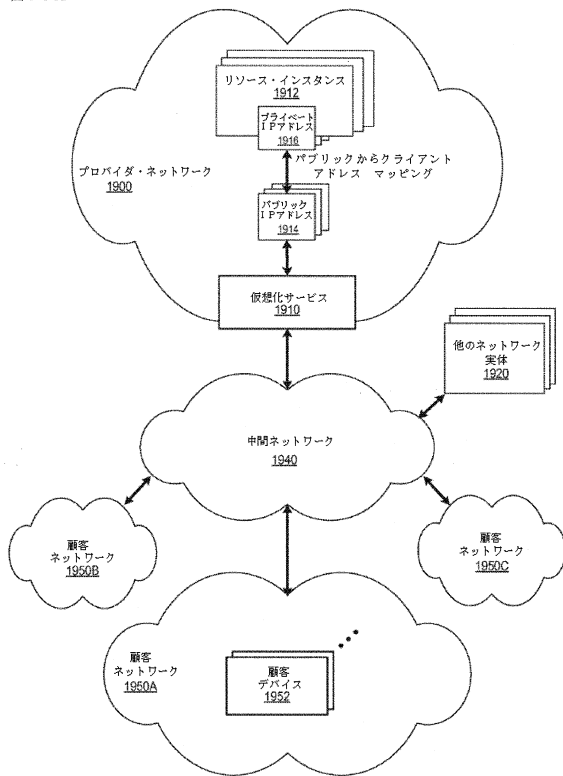


【図32】



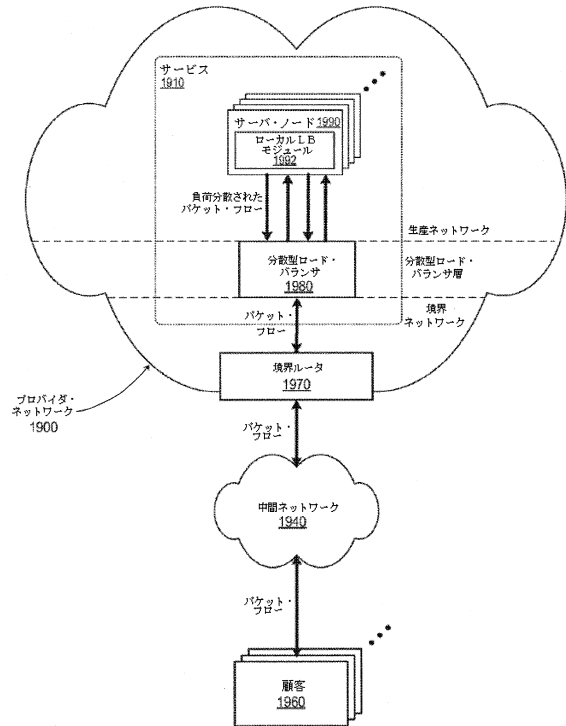
【図33A】

図33A



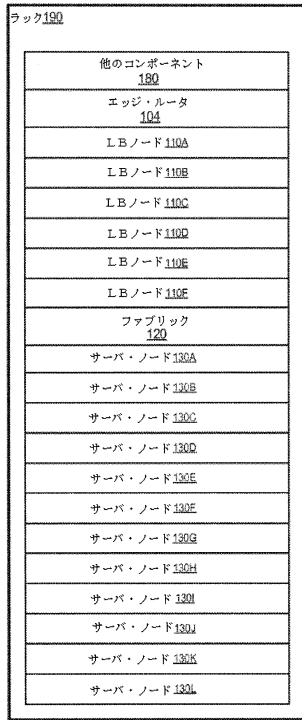
【図33B】

図33B



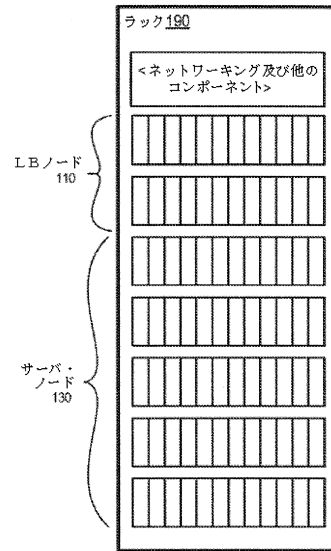
【図34A】

図34A



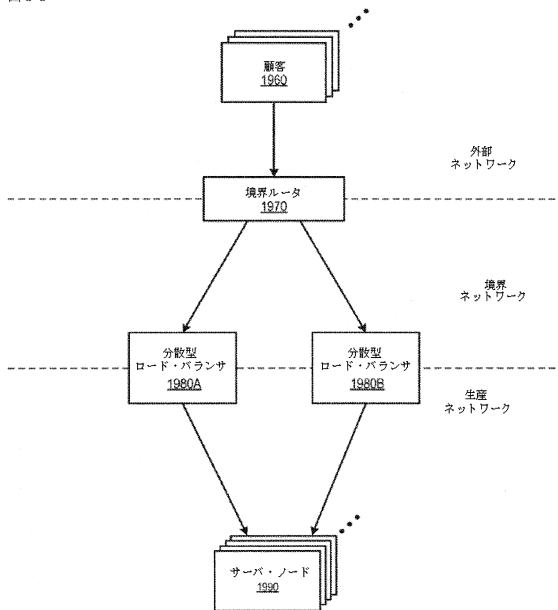
【図34B】

図34B



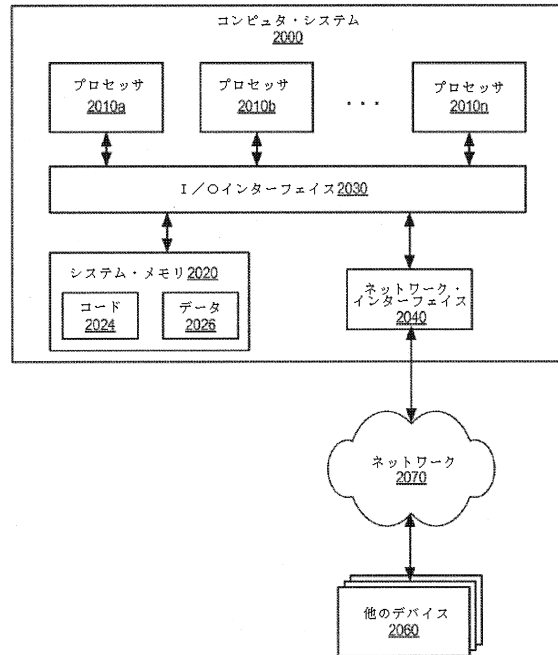
【図35】

図35



【図36】

図36



フロントページの続き

- (72)発明者 ローレンス, ダグラス・スチュワート
アメリカ合衆国・98109-5210・ワシントン州・シアトル・テリー アヴェニュー ノース
・410
- (72)発明者 スリニヴァサン, ヴェンカトラガヴァン
アメリカ合衆国・98109-5210・ワシントン州・シアトル・テリー アヴェニュー ノース
・410
- (72)発明者 ヴァイジャ, アクシャイ・スハス
アメリカ合衆国・98109-5210・ワシントン州・シアトル・テリー アヴェニュー ノース
・410
- (72)発明者 チャン, ファン
アメリカ合衆国・98109-5210・ワシントン州・シアトル・テリー アヴェニュー ノース
・410

審査官 寺谷 大亮

- (56)参考文献 特開2005-025756(JP, A)
特開2003-131961(JP, A)
特表2010-531020(JP, A)
特開2000-029831(JP, A)
特開2012-074928(JP, A)
米国特許出願公開第2010/0036903(US, A1)

(58)調査した分野(Int.Cl., DB名)

G06F 13/00
H04L 12/743
H04L 12/803
H04L 12/951
G06F 9/46