

(19) 日本国特許庁 (JP)

(12) 公表特許公報 (A)

(11) 特許出願公表番号

特表2012-528552

(P2012-528552A)

(43) 公表日 平成24年11月12日 (2012. 11. 12)

(51) Int.Cl.		F I		テーマコード (参考)
<b>H04L 12/56</b>	<b>(2006.01)</b>	H04L 12/56	H	5B089
<b>G06F 13/00</b>	<b>(2006.01)</b>	G06F 13/00	357Z	5K030

審査請求 未請求 予備審査請求 未請求 (全 31 頁)

(21) 出願番号 特願2012-513344 (P2012-513344)  
 (86) (22) 出願日 平成22年5月28日 (2010. 5. 28)  
 (85) 翻訳文提出日 平成24年1月30日 (2012. 1. 30)  
 (86) 国際出願番号 PCT/US2010/036758  
 (87) 国際公開番号 W02010/138937  
 (87) 国際公開日 平成22年12月2日 (2010. 12. 2)  
 (31) 優先権主張番号 61/182, 063  
 (32) 優先日 平成21年5月28日 (2009. 5. 28)  
 (33) 優先権主張国 米国 (US)  
 (31) 優先権主張番号 12/578, 608  
 (32) 優先日 平成21年10月14日 (2009. 10. 14)  
 (33) 優先権主張国 米国 (US)

(71) 出願人 500046438  
 マイクロソフト コーポレーション  
 アメリカ合衆国 ワシントン州 9805  
 2-6399 レッドモンド ワン マイ  
 クロソフト ウェイ  
 (74) 代理人 110001243  
 特許業務法人 谷・阿部特許事務所  
 (72) 発明者 アルバート グリーンバーグ  
 アメリカ合衆国 98052 ワシントン  
 州 レッドモンド ワン マイクロソフト  
 ウェイ マイクロソフト コーポレーシ  
 ョン エルシーエーインターナショナル  
 パテント内

最終頁に続く

(54) 【発明の名称】 アジャイルデータセンタネットワークアーキテクチャ

## (57) 【要約】

本特許出願は特にデータセンタで用いられ得るアジャイルネットワークアーキテクチャに関する。1つの実施例はレイヤ3のインフラストラクチャのマシンを接続する仮想レイヤ2のネットワークを提供する。

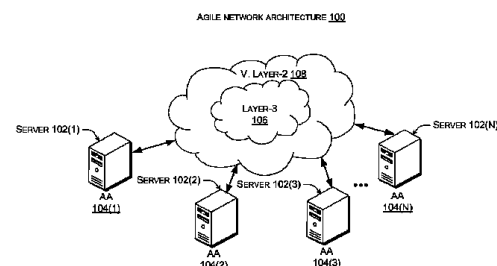


FIG. 1

**【特許請求の範囲】****【請求項 1】**

アプリケーションアドレス（１０４）を複数のマシンに割り当て、ローカルアドレス（２０６）をレイヤ３のインフラストラクチャ（１０６）のコンポーネントに割り当てることによって前記マシンを接続する仮想レイヤ２のネットワーク（１０８）を提供するステップを備えることを特徴とする方法。

**【請求項 2】**

個々のマシン間でターンアラウンドパスを早期に使用するステップを更に備えることを特徴とする請求項 1 記載の方法。

**【請求項 3】**

前記マシンは、サーバ又は仮想マシンを備えることを特徴とする請求項 1 記載の方法。

**【請求項 4】**

前記仮想レイヤ２のネットワークを提供するステップは、複数の仮想レイヤ２のネットワークを提供するステップを備えることを特徴とする請求項 1 記載の方法。

**【請求項 5】**

レイヤ３のコンポーネントのローカルアドレスを有する第１のマシンと第２のマシンの間でパケットを、前記第１のマシンと前記第２のマシンの間の前記レイヤ３のインフラストラクチャの個々のパスに従ってカプセル化するステップを更に備えることを特徴とする請求項 1 記載の方法。

**【請求項 6】**

前記マシンのうちの２つの間で前記レイヤ３のインフラストラクチャの個々のパスをランダムに選択するステップを更に備えることを特徴とする請求項 1 記載の方法。

**【請求項 7】**

前記個々のパスを選択するためにバリエーションロードバランシング（ＶＬＢ）を利用するステップを更に備えることを特徴とする請求項 6 記載の方法。

**【請求項 8】**

周期的に又はネットワークイベントに応答して前記個々のパスを再選択するステップを更に備えることを特徴とする請求項 6 記載の方法。

**【請求項 9】**

サーバ（３１６（１））であって、  
コンピュータ可読命令を実行する少なくとも１つのプロセッサと、  
前記少なくとも１つのプロセッサによって実行可能でかつ他のサーバ（３１６（Ｎ））への配信のためのパケットを受信し、中間スイッチ（３１０）を介して前記配信用の前記パケットをカプセル化するように構成されたアジャイルエージェント（３２０）と、  
を備えることを特徴とするサーバ。

**【請求項 10】**

前記アジャイルエージェントは前記中間スイッチを複数の中間スイッチからランダムに選択するように構成されていることを特徴とする請求項 9 記載のサーバ。

**【請求項 11】**

前記アジャイルエージェントは前記配信用のパスを選択し、通信障害の指示を受信すると、複数の中間スイッチから新しい中間スイッチを含む新しいパスを再選択するように構成されていることを特徴とする請求項 10 記載のサーバ。

**【請求項 12】**

前記サーバは複数の仮想マシンをサポートするように構成され、前記アジャイルエージェントは２つの仮想マシン間で前記パケットを配信するためのパスを選択するように構成されていることを特徴とする請求項 10 記載のサーバ。

**【発明の詳細な説明】****【技術分野】****【0001】**

本発明は、アジャイルデータセンタネットワークアーキテクチャに関する。

10

20

30

40

50

## 【背景技術】

## 【0002】

従来のデータセンタネットワークアーキテクチャはそれらのアジリティ（データセンタネットワークのいずれかのサーバを何らかのサービスに割り当てるためのそれらの能力）を弱体化可能なくつかの設計上の欠点に悩まされている。第一に、従来のネットワークの構成は典型的に、実質的にツリー状であり、比較的高価な機器からなる。これはネットワークの他のどこかで予備能力が利用可能であるときでさえ、計算ホットスポットの混雑及びその発展の結果である。第二に、従来のデータセンタネットワークは1つのサービスにおけるトラフィック輻輳がその周囲の他のサービスに影響を与えることを防止することにほとんど効果がない。1つのサービスがトラフィック輻輳を経験すると、同一のネットワークのサブツリーを共有する全てのサービスが巻き添え被害（ダメージ）を被ることは普通である。第三に、従来のデータセンタネットワークのルーティング設計は、通常、サーバにトポロジ的に重要なインターネットプロトコル（IP）アドレスを割り当て、サーバをVLAN（Virtual Local Area Network）のうちで分けることによりスケールを達成する。しかしながら、これはサーバがサービスの中で再割り当てされるとき大きな設定負担を作り出し、よってデータセンタのリソースを更に使ってしまう。更に、人間の関与がそれらの再設定で通常要求され、これによりこの処理の速度を制限してしまう。最後に、従来のデータセンタネットワークを設定する際の難しさ、そのようなネットワークで使用される機器のコスト等の他の事項は、これらのネットワークのアジリティを否定的に影響を与えることができる。

10

20

## 【発明の概要】

## 【0003】

本特許出願は、特にデータセンタ内で使用され得るアジャイル（agile）ネットワークアーキテクチャに関する。1つの実施例は、レイヤ3のインフラストラクチャ（インフラ）のサーバ等のマシンを接続する仮想レイヤ2のネットワークを提供する。

## 【0004】

他の実施例は複数のスイッチを介して通信可能に接続された複数の計算装置を含んでいる。個々の計算装置をアプリケーションアドレスで関連付けすることができる。1つの計算装置をソースとして動作するように設定することができ、他の計算装置を送り先として動作するように設定することができる。パケットを送り先計算装置のアプリケーションアドレスに送信するようにソース計算装置を設定することができる。この実施例は、パケットを捕まえ、送り先計算装置と関連付けられたロケーションアドレスを識別し、そのロケーションアドレスにパケットを送信するために個々のスイッチを選択するように設定されるアジャイルエージェントを含むことができる。

30

## 【0005】

上記した各実施例は紹介目的のために提供され、請求項に記載した構成の全てを包含及び/又は限定しない。

## 【図面の簡単な説明】

## 【0006】

添付図面は本願で伝えられるコンセプトの実施例を示している。添付図面と関連した次の詳細な説明を参照することにより図示の実施例の特徴を更に容易に理解することができる。各図面の同類の参照番号は実現可能であれば同類の要素を示すために用いられている。更に、各参照番号の最も左の数字は図面と、参照番号が最初に紹介されている関連記載とを表す。

40

## 【0007】

【図1】本コンセプトのいくつかの実施に応じたアジャイルネットワークアーキテクチャの例を示す図である。

【図2】本コンセプトのいくつかの実施に応じたアジャイルネットワークアーキテクチャの例を示す図である。

【図3】本コンセプトのいくつかの実施に応じたアジャイルネットワークアーキテクチャ

50

の例を示す図である。

【図 4】本コンセプトのいくつかの実施に応じたアジャイルネットワークアーキテクチャの例を示す図である。

【図 5】本コンセプトのいくつかの実施に応じたアジャイルネットワークアーキテクチャの例を示す図である。

【図 6】本コンセプトのいくつかの実施に応じたアジャイルネットワークアーキテクチャの例を示す図である。

【図 7】本コンセプトのいくつかの実施に応じたアジャイルネットワークデータセンタレイアウトの例を示す図である。

【図 8】本コンセプトのいくつかの実施に応じたアジャイルネットワークデータセンタレイアウトの例を示す図である。

【図 9】本コンセプトのいくつかの実施に応じたアジャイルネットワークデータセンタレイアウトの例を示す図である。

【図 10】本コンセプトのいくつかの実施に応じて達成可能なアジャイルネットワークのフローチャート図である。

【発明を実施するための形態】

【0008】

(概要)

本特許出願は特にデータセンタ内で使用され得るアジャイルネットワークアーキテクチャに関する。クラウドサービスは潜在的に 10 から何十万ものサーバを保有する大規模なデータセンタの構築を押し進めている。それらのデータセンタは多くの動的な数の異なるサービス (Web アプリケーション、電子メール、マップリデュースクラスタ等) を同時にサポートすることができる。クラウドサービスデータセンタの実施例はスケールアウト設計: 必要に応じてサービス間で急いで再割り当てされ得るリソース (例えば、サーバ) の大きなプールを介して達成される信頼性及び性能に依存することができる。データネットワークのいずれかサーバに何らかのサービスを割り当てるための能力にはデータセンタネットワークのアジリティを考慮することができる。大きなコストと関連付けられ得るデータセンタの利益を効果的に引き出すためには、ネットワークアジリティは高価になる可能性がある。ネットワークアジリティなしではデータセンタサーバリソースは立ち往生となり、よって無駄なお金となる。

【0009】

(第 1 のアジャイルネットワークアーキテクチャ例)

紹介目的のためにアジャイルネットワークアーキテクチャ 100 を例示する図 1 及び図 2 が参照される。アジャイルネットワークアーキテクチャ 100 はサーバ 102 (1)、102 (2)、102 (3) 及び 102 (N) 等の複数のサーバサイド計算装置を含むことができる。

【0010】

用語サーバ及びマシンはデータを送信又は受信することができるいずれかの装置に言及していることを理解されるべきである。例えば、それらの用語は物理サーバ、(例えば、仮想化技術を用いて) サーバ上で動作する仮想マシン、1つのオペレーティングシステムを動作する計算装置、1以上のオペレーティングシステムを動作する計算装置、異なるオペレーティングシステム (例えば、マイクロソフトウインドウズ (登録商標)、Linux、FreeBSD) を動作する計算装置、サーバ以外の計算装置 (例えば、ラップトップ、アドレスابلパワーサプライ (addressable power supply))、又は計算装置の一部 (例えば、ネットワーク接続ディスク、ネットワーク接続メモリ、記憶サブシステム、記憶エリアネットワーク (SAN)、グラフィックス処理ユニット、計算アクセラレータ、量子計算装置) のいずれかに言及していることを理解されるべきである。

【0011】

アジャイルネットワークアーキテクチャ 100 はサーバの数に関連してスケラビリティ

10

20

30

40

50

(拡張性)を進展させることができる。スケラビリティを達成することができる1つの方法は、アプリケーションを利用してサーバ102(1)~102(N)に対するイーサネットのようなフラットアドレッシングを作り出すことによって達成され得る。イーサネットのレイヤ2セマンティクスは、フラットアドレッシングをサポートするネットワーク状態を達成することと関連付けることができ、フラットアドレッシングはサーバがLAN(ローカルエリアネットワーク)上に存在したように、いずれかのインターネットプロトコル(IP)アドレスをいずれかのネットワークポートに接続されたいずれかのサーバに割り当てることができる。

#### 【0012】

この場合において、アプリケーションアドレス(AA)104(1), 104(2), 104(3), 104(N)を各サーバ102(1), 102(2), 102(3), 102(N)に各々割り当てることができる。サーバの観点からすれば、いずれかのサーバ関連付けされたアプリケーションアドレス104(1), 104(2), 104(3), 104(N)を介していずれかの他のサーバと通信することができる。これは何らかの様式でアプリケーションアドレスを配置することができる場合に、サーバ102(1), 102(2), 102(3), 102(N)を備えるLANに対して有効であるそれら全てを含む、レイヤ2の機能として考えることができる。しかしながら、以下に説明されるように、いくつかの実施例では、アジャイルネットワークアーキテクチャの根本のインフラは106で示されるようにレイヤ3であっても良い。よって、それらの実施例は、レイヤ3のインフラ106上に(又は利用して)仮想レイヤ2のネットワーク108を作り出すことができる。それらは同一のレイヤ3のインフラ106上に作り出した1以上の仮想レイヤ2のネットワーク108であっても良く、また各サーバはそれらの1以上の仮想レイヤ2のネットワーク108に属することができる。

#### 【0013】

図2はアジャイルネットワークアーキテクチャ100にインターネット204を介して接続されている外部クライアント202を示している。アジャイルネットワークアーキテクチャ100は、アプリケーションアドレス104(1)~104(N)の知識を有する外部クライアントなしで1以上のサーバ102(1)~102(N)に割り当てられているグローバル又はロケーションアドレス206と外部クライアントが通信することを可能にする。それらのコンセプトは図3~図5の説明において以下に詳細に説明される。

#### 【0014】

(第2のアジャイルネットワークアーキテクチャ例)

図3は上記のコンセプトを実施可能なアジャイルネットワークアーキテクチャ300を例示している。この場合において、外部クライアント302はアジャイルシステム304とインターネット306及び/又は他のネットワークを介して通信することができる。この実施例では、アジャイルシステム304は、308で概して示された複数のルータ308(1)~308(N)、310で概して示された複数の中間スイッチ310(1), 310(2)及び310(N)、312で概して示された複数の集合スイッチ312(1), 312(2)及び312(N)、314で概して示された複数のトップオブラック(TOR又はToR)スイッチ314(1), 314(2)及び314(N)、並びに316で概して示された複数のサーバ316(1), 316(2), 316(3), 316(4), 316(5)及び316(N)を含んでいる。図面のページの空間制限のため、6つのサーバ316(1)~316(N)だけがここでは示されているが、アジャイルシステム304は、何千、何万、何十万、又はそれ以上のサーバを提供することができる。なお、簡潔及び図面のページの空間制限のため、コンポーネント間の接続(すなわち、通信パス)は図3~図8に示されているとは限らない。

#### 【0015】

サーバ316(1)及び316(2)はサーバラック318(1)としてTORスイッチ314(1)と関連付けられている。同様に、サーバ316(3)及び316(4)はサーバラック318(2)としてTORスイッチ314(2)と関連付けられ、サーバ3

10

20

30

40

50

1 6 ( 5 ) 及び 3 1 6 ( N ) はサーバラック 3 1 8 ( N ) として T O R スイッチ 3 1 4 ( N ) と関連付けられている。これは上記したように図面ページの空間制限のためであり、度々、サーバラックは 1 0 以上のサーバを含む。更に、個々のサーバはアジャイルエージェントと関連付け可能である。例えば、サーバ 3 1 6 ( 1 ) はアジャイルエージェント 3 2 0 ( 1 ) と関連付けられる。同様の関係はサーバ 3 1 6 ( 2 ) ~ 3 1 6 ( N ) とアジャイルエージェント 3 2 0 ( 2 ) ~ 3 2 0 ( N ) との間に各々示される。

【 0 0 1 6 】

アジャイルエージェント 3 2 0 ( 1 ) ~ 3 2 0 ( N ) の機能は以下に詳細に説明される。短くは、アジャイルエージェントは個々のサーバ間での通信を容易にすることができる。特にこの例では、アジャイルエージェントはコンピュータ読み取り可能な命令としてサーバ上に保存された論理モジュールとして考えても良い。他の実施例は、サーバ群を提供するアジャイルエージェント 3 2 0 がスイッチ、例えば、T O R スイッチ又は中間スイッチに配置される構成を含んでも良い。スイッチに配置されたときには、アジャイルエージェントは、パケットがサーバ 3 1 6 から中間スイッチ 3 1 0 に向けてネットワークを逆送するようにパケットを処理する。そのような構成では、アジャイルエージェント 3 2 0 はパケット転送パス上のカスタムハードウェアと、その転送パス又はスイッチの制御プロセッサで実施するソフトウェア命令との組み合わせを用いて実施されても良い。

【 0 0 1 7 】

アジャイルシステム 3 0 4 は 3 つのディレクトリサービスモジュール 3 2 2 ( 1 ) ~ 3 2 2 ( N ) を更に含んでいる。ディレクトリサービスモジュールの図示の数はそのアジャイルシステムに必須ではなく、他の実施例は 2、3 以上のディレクトリサービスモジュール ( 及び / 又は他の図示のコンポーネント ) を使用することができる。ディレクトリサーバの機能は以下に詳細に説明される。短くは、ディレクトリサービスモジュールは、他の情報のうちの、アジャイルエージェント 3 2 0 ( 1 ) ~ 3 2 0 ( N ) ( 及び / 又は他のコンポーネント ) によって利用され得るアプリケーションアドレス対ロケーションアドレスマッピング ( 順方向又は逆方向のマッピングのいずれか一方又は両方 ) を含むことができ、これによりアジャイルシステム 3 0 4 内の通信を容易にする。この場合において、ディレクトリサービスモジュール 3 2 2 ( 1 ) ~ 3 2 2 ( N ) は特定のサーバ 3 1 6 ( 1 )、3 1 6 ( 3 ) 及び 3 1 6 ( 5 ) と関連付けされる。他の設定では、そのディレクトリサービスモジュールは、データセンタ制御サーバ、スイッチ等の他のコンポーネントと共に、及び / 又は専用の計算装置上で存在することができる。

【 0 0 1 8 】

アジャイルシステム 3 0 4 は 2 つの論理グループ化として考えても良い。第 1 の論理グループ化は 3 2 6 で示されたようにロケーション又はグローバルアドレスを担うリンクステートネットワークである。第 2 の論理グループ化は 3 2 8 で示されたようにアプリケーションアドレスを所有するサーバの代替可能なプールである。短くは、リンクステートネットワーク 3 2 6 のコンポーネントは、サーバ 3 2 8 のプールのどのサーバがどのアプリケーションアドレスを現在用いているかを追跡するために情報を交換することを必要としない。また、サーバの見地からは、サーバはサーバプール 3 2 8 内のいずれかの他のサーバと他のサーバのアプリケーションアドレスを介して通信することができる。この処理は、サーバに対して透過的であるような方式でアジャイルエージェント、ディレクトリサービス、及び / 又は他のコンポーネントによって容易にされる。他の方法としては、その処理は、サーバ上の他のコンポーネントがその処理に気づくかもしれないけれども、サーバ上で動作しているアプリケーションに対して透過的であっても良い。

【 0 0 1 9 】

ルータ 3 0 8、中間スイッチ 3 1 0、集合スイッチ 3 1 2、T O R スイッチ 3 1 4 及びサーバ 3 1 6 ( 1 ) ~ 3 1 6 ( N ) は、レイヤ 3 の技術を用いる等で通信可能に接続されても良い。個々のサーバの見地からは、他のサーバとの通信はレイヤ 2 の通信 ( 例えば、仮想レイヤ 2 ) として現れる。しかしながら、サーバラック 3 1 8 ( 1 ) のソースサーバ 3 1 6 ( 1 ) からサーバラック 3 1 8 ( 2 ) の送り先サーバ 3 1 6 ( 3 ) へのようなラッ

10

20

30

40

50

ク間通信はレイヤ 3 のインフラを介して実際に起きる。例えば、アジャイルエージェント 3 2 0 ( 1 ) は通信 ( 例えば、サーバ 3 1 6 ( 3 ) のアプリケーションアドレスにアドレス指定されたパケット ) を捕らえ、その伝送を容易にすることができる。

#### 【 0 0 2 0 】

アジャイルエージェント 3 2 0 ( 1 ) はディレクトリサービスモジュール 3 2 2 ( 1 ) ~ 3 2 2 ( N ) のうちの 1 以上にアクセスして、サーバ 3 1 6 ( 3 ) で関連付けられたロケーションアドレスに対するアプリケーションアドレスのマッピングを得ることができる。例えば、そのマッピングしたロケーションアドレスは T O R スイッチに対してでも良い。アジャイルエージェントはロケーションアドレスを有するパケットをカプセル化することができる。そのアジャイルエージェントはそのカプセル化パケットを送信又は中継するように個々 ( 又はセット ) の集合及び / 又は中間スイッチを選択することができる。この選択処理の機能は以下に更に詳細に説明される。T O R スイッチ 3 1 4 ( 2 ) でそのカプセル化パケットを受け取ると、その T O R スイッチはそのパケットをデカプセル化してそのパケットをサーバ 3 1 6 ( 3 ) に送信することができる。代替実施形態では、ロケーションアドレスはサーバ 3 1 6 ( 3 ) 又はサーバ 3 1 6 ( 3 ) 上で動作する仮想マシンと関連付けられても良く、パケットは送り先サーバそれ自身上でデカプセル化されても良い。それらの実施形態において、アプリケーションアドレスが他のホストがアプリケーションとの通信のために使用しているアドレスである L A N によって接続されているアプリケーションに対して示し続けるために、サーバ又は仮想マシンに割り当てられたロケーションアドレスはサーバ上で動作する他のアプリケーションから隠されても良い。

#### 【 0 0 2 1 】

代替の実施形態において、パケットはレイヤ 3 / レイヤ 2 の境界を横切る際に他のコンポーネントによってカプセル化されても良い。例えば、デカプセル化を行うことができるコンポーネントの例はハイパーバイザ及び / 又は仮想マシンモニタのルートパーティションを含んでも良い。

#### 【 0 0 2 2 】

この設定はサーバがサーバプール 3 2 8 の多数の中に加えられることを可能にし、また、サーバの見地からは、他のサーバはそれらが同一のサブネットワーク上にあるように現れることができる。代替的に又は追加的に、リンクステートネットワーク 3 2 6 のコンポーネントはサーバアプリケーションアドレスを気にする必要がない。更に、サーバが追加又は削除された時等のアドレス情報が変化したときはいつでもディレクトリサーバは多数の異なる種類のコンポーネントを更新するよりもむしろ単に更新されても良い。

#### 【 0 0 2 3 】

要約すれば、レイヤ 2 のセマンティックは、サーバが L A N 上に存在したように、いずれかのネットワークポートに接続されたいずれかのサーバにいずれかの I P アドレスを割り当てることができるフラットアドレッシングをサポートするネットワーク状態を達成することと関連付けられても良い。また、リンクステートネットワーク 3 2 6 内のコンポーネント ( すなわち、スイッチ ) はそのリンクステートネットワーク内の他のコンポーネントに気づいても良いが、サーバプール 3 2 8 のコンポーネントに気づく必要はない。更に、T O R スイッチはそれら各々のラック内のサーバについて知っても良いが、他のラックのサーバについては知る必要がない。更に、アジャイルエージェントはサーバアプリケーションアドレス ( A A ) パケットを捕らえて、A A の送り先計算装置と関連付けられたロケーションアドレス ( L A ) を識別することができる。そして、アジャイルエージェントは L A にパケットを送信するために個々のスイッチ ( 又はスイッチのセット ) を選択することができる。この場合において、その個々のスイッチは利用可能なスイッチの 1 以上を含むことができる。

#### 【 0 0 2 4 】

また、この設定はサービスに関係した他のサーバの機能を促進する。例えば、ディレクトリサービスモジュール 3 2 2 ( 1 ) ~ 3 2 2 ( N ) に含まれても良いようなデータセンタ管理ソフトウェアはいずれかのサービスにいずれかのサーバ 3 1 6 ( 1 ) ~ 3 1 6 ( N

）を割り当て、そのサービスが期待するどのようなＩＰアドレスを有するサーバでも設定することができる。各サーバのネットワーク設定はＬＡＮを介して接続された場合と同一とすることができ、リンクローカルブロードキャストのような機能をサポートすることができる。サービス間の通信分離の目的は、サービス及び通信グループを定義する簡単で一貫したアプリケーションプログラムインターフェース（ＡＰＩ）を提供することと関連付けられても良い。この点において、ディレクトリサービスはサービス（例えば、顧客）と関連付けられたサーバのグループを定義することができる。１つのグループ内のサーバ間での完全接続を許可しても良く、異なるグループのどのサーバが通信することを許可されるべきかを管理するためにアクセスコントロールリスト（ＡＣＬ）のような方針を規定しても良い。

10

#### 【００２５】

更に、上記の設定はトラフィック管理にそれ自身役立つ。説明の目的のために、第１の顧客がアジャイルシステム３０４のサーバによって実行されるべきサービスに対して比較的高いレートを支払い、それによってサービス契約の比較的高い品質を得ると仮定する。更に、第２の顧客が比較的低いレートを支払い、それによってサービス契約の比較的低い品質を受けると仮定する。そのような場合において、第１の顧客のためにトラフィックを運営するように中間スイッチ３１０（１）～３１０（Ｎ）の比較的高いパーセンテージ、又は全てを割り当てることができ、第２の顧客に対してより少ない数のスイッチを割り当てることができる。別の言い方をすれば、第１の顧客にスイッチの第１のサブセットを割り当てることができ、第２の顧客にスイッチの第２のサブセットを割り当てることができる。第１のセット及び第２のセットを互いに排他的に又は重複させても良い。例えば、いくつかの実施例では、個々のスイッチを特定顧客に対して専用に、又は複数の顧客に割り当てることができる。例えば、中間スイッチ３１０（１）を両方の顧客に割り当てることができ、中間スイッチ３１０（２）及び３１０（Ｎ）を排他的に第１の顧客に割り当てることができる。

20

#### 【００２６】

要約すると、以下に詳細に説明されるように、アジャイルネットワークアーキテクチャ３００は次の１以上の目的、すなわちサーバ間の一樣な高い能力、サービス間のパフォーマンス分離、イーサネットレイヤ２のセマンティック、及び／又はサービス間の通信分離に関連付けられても良い。サーバ間の一樣な高い能力の目的は、ネットワークを流れるトラフィックの割合が、先ず、送信するサーバ及び受信するサーバのネットワークインターフェースカードの利用可能な能力による場合を除いて、制限されないネットワーク状態を達成することと関連付けられても良い。よって、開発者の見地からは、その目的を達成することによって、ネットワークトポロジは、サーバを１つのサービスに加えるときにはもはや第一の関心事ではないかもしれない。サービス間のパフォーマンス分離の目的は、各サービスが個別の物理スイッチによって接続されたかのように、１つのサービスのトラフィックが何らかの他のサービスによって運営されたトラフィックによって影響を受けないネットワーク状態を達成することと関連付けられても良い。イーサネットレイヤ２のセマンティックの目的は、サーバがＬＡＮ上に存在したように、ほとんどのＩＰアドレスが何らかのネットワークポートに接続されたいずれかのサーバに割り当てられ得るフラットアドレスリングをサポートするネットワーク状態を達成することと関連付けられても良い。よって、データセンタ管理ソフトウェアはいずれかのサーバを何らかのサービスに割り当て、サービスが期待するどんなＩＰアドレスを有するサーバでも設定することができる。

30

40

#### 【００２７】

各サーバのネットワーク設定はＬＡＮを介して接続された場合と同一とすることができ、リンクローカルブロードキャストのような機能をサポートすることができる。サービス間の通信分離の目的は、サービス及び通信グループを定義する簡単で一貫したＡＰＩを提供することと関連付けられても良い。この点において、サーバのグループを定義するディレクトリシステム（すなわち、例えば、ディレクトリサービスモジュール３２２（１）～３２２（Ｎ）を介して）は提供されても良い。１つのグループ内のサーバ間での完全接続

50



を許可しても良く、異なるグループのどのサーバが通信することを許可されるべきかを管理するために方針を規定しても良い。

【 0 0 2 8 】

説明したアジャイルネットワークアーキテクチャを利用することによって、次のネットワークの特徴の 1 以上と関連付けられているデータセンタネットワークを備えることができる。( 1 ) サービスインスタンスがネットワークのどこにでも配置されることを可能にするフラットアドレッシング、( 2 ) ネットワークパスを越えて一様にトラフィックを分散するようにランダム化を使用するロードバランシング(例えば、バリエーションロードバランシング( V L B : V a l i a n t L o a d B a l a n c i n g ) )、( 3 ) イーサネットレイヤ 2 のセマンティックを達成し、大きなサーバプールヘスケーリングを行う新しいエンドシステムベースのアドレス解決サービス。

10

【 0 0 2 9 】

上記した目的を達成するために、少なくともいくつかの実施形態において、次のアジャイルネットワークアーキテクチャの設計原理の 1 以上を様々な実施例で用いることができる。

【 0 0 3 0 】

( 多数のパスダイバシティを有するトポロジの利用 )

「メッシィ( m e s h y )」トポロジを利用することにより、サーバの個々のセット間の多数のパスを提供することができる。例えば、サーバラック 3 1 8 ( 1 ) の複数のサーバとサーバラック 3 1 8 ( N ) の複数のサーバとの間の通信は T O R スイッチ 3 1 4 ( 1 ) から集合スイッチ 3 1 2 ( 1 ) ~ 3 1 2 ( 2 ) のいずれかを介して中間スイッチ 3 1 0 ( 1 ) ~ 3 1 0 ( N ) のいずれかに達することができる。中間スイッチから、その通信は集合スイッチ 3 1 2 ( 2 ) ~ 3 1 2 ( N ) のいずれかを介して T O R スイッチ 3 1 4 ( N ) に達することができる。

20

【 0 0 3 1 】

この設定はいくつかの利益を得る結果にすることができる。例えば、多数のパスの存在は、明白なトラフィックエンジニアリング又はパラメータのチューニングの必要性なしにネットワークの輻輳状態の減少及び / 又は解消を可能にする。更に、多数のパスは「スケールアウト( s c a l e - o u t )」ネットワーク設計を可能にさせる。言い換えると、より低コストのスイッチを加えることにより更なる能力を加えることができる。逆に、従来の階層的なネットワーク設計は階層の高いレベルで 1 又は非常に僅かのリンクのトラフィックを集中させる。その結果、従来のネットワークはトラフィックの高い密度に対処するために高価な「ビッグアイロン( b i g i r o n : 大型コンピュータ )」スイッチの購入を必要とする可能性がある。

30

【 0 0 3 2 】

更に、「メッシィ」トポロジを利用することにより、多数のパスはリンク又はスイッチが機能しなくなったとき正常な劣化を可能にさせる。例えば、所定のレイヤで「 n 」のスイッチを有する説明したアジャイルデータセンタネットワークアーキテクチャに応じて実施されたアジャイルネットワークは、スイッチが機能しなくなったとき、その能力の 5 0 % を失う可能性がある従来のネットワークと比較して、その能力の 1 / n だけ失う可能性がある。アジャイルデータネットワークアーキテクチャに応じて実施されたアジャイルネットワークは、完全な二分トポロジを潜在的に利用することができる。

40

【 0 0 3 3 】

( アドレス不安定さに対するランダム化 )

データセンタはそれらの負荷、それらのトラフィック、及びそれらの故障パターンにおいて非常に大きい不安定さを持つ可能性がある。よって、リソースの複数の大きなプールを作り出すことができる。負荷をランダムにそれらに分散することができ、最悪の場合を平均的な場合に改善するために、最良の場合のいくつかの性能を、引き換えにすることができる。少なくともいくつかの実施形態において、広範囲のパスダイバシティと関連付けられたトポロジ(例えば、図 3 に明らかにされたように)が利用されても良い。 V L B 技

50

術等のロードバランシング技術を用いてトポロジを越えてワークフローのルート決めすることができる。短くは、V L B 技術はデータ伝送を担うために用いられるパス又は複数のパスをランダムに選択することを含むことができ、ここでパスは一連のリンク及び / 又はスイッチからなる。続いて、パスを再選択し、ここで再選択は初期のパスを構成するスイッチ又はリンクの 1 以上を変えることを必要とする。その再選択は、特定のバイト / パケット数を送信 / 受信した後に、及び / 又は選択したパス、スイッチ又はリンクと関連付けられた伝送状況が生じたことに応答するように周期的に起きても良い。例えば、パケット遅延又は他の通信障害が検出されたならば、その選択処理が繰り返されても良い。この原理のアプリケーションを通して、一様な能力及びパフォーマンス分離目的を合致することができる。

10

#### 【 0 0 3 4 】

特に、データセンタトラフィックマトリックスにおけるアドレスの不安定さ及び不確かさに対して、ネットワークパス間でフローをランダムにハッシュ ( h a s h ) するためにロードバランシング技術 ( 例えば、V L B ) を利用することができる。このアプローチに対する目的は、ホーストラフィックモデルにおいてのように、ネットワークの入出制限の影響下にある任意のトラフィック変動のための帯域幅保証を提供することである。短くは、ホースモデルは、所定のパスに亘るデータ伝送レートがそのパスの最も遅い又は最も規制された部分を越えることができないことを特定する。

#### 【 0 0 3 5 】

フロー単位 ( フローの大部分のパケットがパスを再選択したときを除いてネットワークを介して同一のパスに続いていることを意味する ) で、V B L のようなロードバランシング技術を用いることは、有利である可能性がある。それは、フローのパケットが並べ直される機会、又は送り先で判断された待ち時間を急に変更する体験、及び / 又はフロー内の M T U ( M a x i m u m T r a n s m i s s i o n U n i t ) の違いのためにパス M T U 発見プロトコルの混乱動作を減少することができるからである。トラフィックのいくつかの種類 ( 例えば、パケットの並べ直しによって悪影響を与えない種類 ) 及びいくつかの環境 ( 例えば、全てのパスに亘って非常に均一な遅延を有する環境 ) は、パケット精度 ( 潜在的な異なるパスがパケットのシーケンスにおけるパケット毎に使用されることを意味する ) で V L B のようなロードバランシングを用いることを好むかもしれない。フローの共通に受け入れられた定義のいくつかを用いることができる。例えば、I P 5 タプルフロー ( 5 - t u p l e f l o w ) 、 I P 2 タプルフロー ( 2 - t u p l e f l o w ) 、又は 2 つのサブネット又はアドレスレンジ間のパケットのセット。

20

30

#### 【 0 0 3 6 】

アジャイルデータセンタネットワークを提供するコンテキストにおいて、入出制限はサーバラインカード速度に対応することができる。ハイバイセクション帯域幅トポロジ ( h i g h b i s e c t i o n b a n d w i d t h t o p o l o g y ) ( 例えば、折り畳み C l o s t ポロジ ) との組み合わせにおいて、非干渉パケットスイッチ切替ネットワーク ( 非ブロッキング回路スイッチ切替ネットワークの対応 ) を作り出し、サーバ入出ポート速度を越える負荷を維持していない) トラフィックパターンのためのホットスポットフリーパフォーマンスを提供するためにロードバランシング技術を利用することができる。この点において、いくつかの実施例では、T C P ( T r a n s m i s s i o n C o n t r o l P r o t o c o l ) のエンドツーエンド輻輳制御機構はホースモデルを実行して過動作サーバポート速度を避けるために利用されても良い。この原理は図 3 に示された論理トポロジに導入することができ、それは 3 つの異なるスイッチ層、T O R 3 1 4 、集合 3 1 2 及び中間 3 1 0 から構成することができる。1 つのサーバから他のサーバへのフローは、ランダム中間スイッチを介し、複数の T O R 及び集合スイッチを越えてランダムパスをとることができる。V L B 等のロードバランシング技術は、利用を円滑にして持続するトラフィック輻輳を解消するためにデータセンタの中間スイッチ構造のコンテキストに利用されても良い。

40

#### 【 0 0 3 7 】

50

## (ロケーションからの名称分離)

ロケーションから名称を分離することは、新しい機能を実施するために用いられ得る自由度を作り出すことができる。この原理はデータセンタネットワークにおいてアジリティを可能にし、アドレスとロケーションとの間を結びつけることが起き得るフラグメンテーションを減少することによって利用を改善するために利用することができる。この原理のアプリケーション及び以下に説明されるエンドシステムを包括する原理を通して、レイヤ2のセマンティクス目的を合致することができる。よって、ネットワークトポロジに関係なく、またそれらのアプリケーション又はネットワークスイッチを再設定することなくIPアドレスを割り当てることが開発者に可能になる。

## 【0038】

10

ネットワークアジリティ(サーバ上のサービスをサポートすること、サーバプールの動的成長及び縮小、及び負荷マイグレーション)を拡張するために、AAと称した名称及びLAと称したロケーションを分離するIPアドレッシングスキームを用いることができる。ディレクトリサービスモジュール322(1)~322(N)として示されたようなアジャイルディレクトリサービスは、拡張性及び信頼性ある方法でAAとLAとの間のマッピングを管理するために定義され得る。アジャイルディレクトリサービスは、個々のサーバ上でのネットワーキングスタックにおいて動作するシム層によって含まれても良い。図3に表された実施例では、このシム層をアジャイルエージェント320(1)~320(N)として示している。

## 【0039】

20

## (エンドシステムの包括)

データセンタサーバ上のオペレーティングシステムを含むソフトウェアは、通常、データセンタ内での利用のために広範囲に修正される。例えば、新しい又は修正したソフトウェアは仮想化のためのハイパーバイザ又は複数のサーバに亘ってデータを保存するためにプロブファイルシステムを作り出すことができる。スイッチ上のソフトウェアを変更するよりむしろ、このソフトウェアのプログラム可能性を利用することができる。更に、スイッチ又はサーバのハードウェアへの変更を避ける又は制限することができ、使われていないアプリケーションが修正されずに残すことができる。現在利用可能な低コストのスイッチASIC(application specific integrated circuit)の制限内で動作するようにサーバ上でソフトウェアを用いることによって、今日作り開発され得る設計を作り出すことができる。例えば、ARP(Address Resolution Protocol)パケットのブロードキャストにより作り出される拡張性の問題は、スイッチ上でソフトウェア又はハードウェアの変更を介してARPを制御する試みよりむしろ、サーバ上のARPLクエストを捕らえ、それらをディレクトリシステムに対してルックアップリクエストに変換することによって、減少及び/又は解消され得る。

## 【0040】

30

図4はアジャイルエージェント320(1)を詳細に例示している。この場合において、アジャイルエージェント320(1)はユーザモード406及びカーネルモード408を含むサーバマシン402(1)上で動作する。そのサーバマシンはユーザモードでユーザモードエージェント410を含む。カーネルモードはTCPコンポーネント412、IPコンポーネント414、エンカプスレータ416、NIC418及びルーティング情報キャッシュ420を含む。サーバマシンはディレクトリサービス322(1)を含む、及び/又はそれと通信することができる。ディレクトリサービスはサーバロールコンポーネント422、サーバヘルスコンポーネント424、及びネットワークヘルスコンポーネント426を含むことができる。アジャイルエージェント320(1)はユーザモードエージェント410、エンカプスレータ416、及びルーティング情報キャッシュ420を含んでも良い。エンカプスレータ416はARPを捕らえ、それをユーザモードエージェント410に送信することができる。ユーザモードエージェントはディレクトリサービス322(1)に問い合わせることができる。カーネルモードコンポーネントにユーザモ

40

50

ドエージェントを含むように、又はルーティングテーブルルックアップの期間のようなARP以外の機構を介して又はIPテーブル又はIPチェイン等の機構を介してディレクトリルックアップを含むように、それらのブロックの他の配置が可能であることが理解されるべきである。

#### 【0041】

図3のアジャイルネットワークアーキテクチャにおいて、エンドシステム制御は新しい機能を速く注入する機構を提供することができる。よって、アジャイルエージェントはロードバランシングに使用されるランダム化を制御することによってきめ細かいパス制御を提供することができる。加えて、名称及びロケーションの分離を実現するために、アジャイルエージェントはイーサネットのARP機能をアジャイルディレクトリサービスへの問い合わせ(クエリ)で置き換えることができる。アジャイルディレクトリサービス自身は、スイッチよりむしろ、サーバ上に実現され得る。このアジャイルディレクトリサービスは、サーバ到達可能性、グループ化、アクセス制御、リソース配分(例えば、中間スイッチの能力)、分離(例えば、中間スイッチの非重複化)、及び動的成長及び縮小のきめ細かい制御を可能にする。

#### 【0042】

(ネットワーク技術の利用)

ネットワークスイッチにおけるロバスト実装を有する1以上のネットワーク技術を利用することは、アジャイルネットワークの設計を簡単にし、そのようなネットワークを開発するためにオペレータ意欲を増加させることができる。例えば、少なくともいくつかの実施形態において、リンクステートルーティングプロトコルは、サーバから障害を隠すためにネットワークスイッチ上に実装されても良く、また、アジャイルディレクトリサービスでの負荷を減少させる手助けに利用され得る。それらのプロトコルはアジャイルネットワークのためにトポロジ及びルートを維持するために利用されても良く、それはアジャイルディレクトリサービスとネットワーク制御プレーンとの間の結合を減少させることができる。スイッチ上のエニキャスト(anycast)アドレスを定義するルーティング設計を通して、上記したアジャイルアーキテクチャはサーバからスイッチの障害を隠すためにECMP(Equal Cost Multi-Path)を利用することができる。更に、これはディレクトリシステム上の負荷を減少させることができる。また、多数のパスの使用をサポートする他のルーティングプロトコルは適切である。

#### 【0043】

(仮想レイヤ2のネットワーク例に関する実施詳細)

(トポロジのスケールアウト)

従来のネットワークは、典型的に、ネットワークの最高レベルでトラフィックを少しのスイッチに集中させる。これはそれらの装置の能力に二等分の帯域幅を両方制限し、それらが機能しなくなったときネットワークに影響をかなり与える可能性がある。しかしながら、それらの問題を避けるために、トラフィックの不安定さに対処するためにランダム化を用いる原理によって運営されるアジャイルネットワークトポロジを利用することができる。この点について、ネットワーク装置をスケールアウトするアプローチをとることができる。これは、図3に示されたように、早送りに専門化され得る低複雑スイッチの比較的ブロードなネットワークに、結果的になる可能性がある。これは、中間スイッチ310(1)~310(N)と集合スイッチ312(1)~312(N)との間のリンクが完全な二分グラフを形成することができる折り畳み Clos ネットワークの例である。従来のトポロジにあるように、TORは2つの集合スイッチに接続することができる。しかしながら、いずれかの2つの集合スイッチの間の多数のパスは、nの中間スイッチがあるならば、それらのいずれかの故障が1/nだけ - 帯域幅の正常な低下と称され得る望ましいプロパティ、二分帯域幅を減少させることを意味する。更に、オーバサブスクリプション(oversubscription)がないようにクロネットワーク等のネットワークを設計することができる。例えば、図3において、Dの数のインターフェースポートを有する集合スイッチ及び中間スイッチを用いることができる。それらのスイッチは、スイッチの

10

20

30

40

50

各層間の能力がリンク能力の  $D * D / 2$  倍であるように接続され得る。

【 0 0 4 4 】

ネットワークの上段又は「スパイン ( spine )」で中間スイッチを介して戻って行くことによってネットワークがサーバラインカードで入出限度の影響下にある潜在的に全ての可能なトラフィックマトリックスに対する帯域幅保証を提供することができることに於いて、ロードバランシング (例えば、V L B ) に対して C l o s ネットワークのようなネットワークは格別によく適合されても良い。ルーティングは単純で回復力 (例えば、ランダムパスはランダム中間ノード及び取り払われたランダムパスを必要としても良い) を持っていて良い。

【 0 0 4 5 】

上記したアジャイルアーキテクチャは従来のネットワークアーキテクチャで達成されるよりも大きなパス制御を提供することができる。特に、複数の中間ノードを区分することができ、異なる区分に専用化されたトラフィッククラスは高い帯域幅全体をいくつかのトラフィッククラスに割り当てる。輻輳指示は、I E E E ( I n s t i t u t e o f E l e c t r i c a l a n d E l e c t r o n i c s E n g i n e e r s ) 8 0 2 . 1 Q a u 輻輳制御におけるように、E C N ( E x p l i c i t C o n g e s t i o n N o t i f i c a t i o n ) 又は類似の機構を介して送り元に信号で戻されても良い。よって、E C N 信号を集める送り元は、(上記のパスを再選択すると呼ばれた) ネットワークを介して代替パスを選択するために用いられるソースパケット内のフィールドを変えることによって応答することができる。

【 0 0 4 6 】

(アジャイルルーティング)

ロケーションから名称を分離する原理を実施するために、アジャイルネットワークは2つのIPアドレスファミリを用いることができる。図3はそのような分離を示している。ネットワークインフラは複数のL A に関して動作することができる。スイッチ及びインターフェース ( 3 1 0 ( 1 ) ~ 3 1 0 ( N ) , 3 1 2 ( 1 ) ~ 3 1 2 ( N ) , 及び 3 1 4 ( 1 ) ~ 3 1 4 ( N ) ) は割り当てられた複数のL A であることができる。スイッチはそれらのL A を担うリンクステートIPルーティングプロトコルを動作させることができる。

【 0 0 4 7 】

サーバ 3 1 6 ( 1 ) ~ 3 1 6 ( N ) 上で動作するようなアプリケーションは複数のL A に気付かず、複数のA A に気付くことができる。この分離をいくつかの利益と関連付けることができる。第一に、パケットはA A に直接送信されるよりむしろ適当なL A に送られても良い (スイッチはそれらを送信するためにホスト毎にルーティングエントリを維持する必要がない)。これは複数のA A を複数のL A に変換するアジャイルディレクトリサービスが、どのサービスが通信のために許可されるべきかについての方針を実行することができることを意味している。第二に、低コストスイッチは度々、全てのL A ルータを保持することができる小さなルーティングテーブル (例えば、12 K エントリ) を有するが、A A の数だけ負担が掛かる。このコンセプトは特に、スイッチが保持することができるルーティングエントリの数より多いネットワークが作られることを可能にすることにおいて価値がある可能性がある。第三に、いずれかのA A がトポロジとは関係なくいずれかのサーバに割り当てられ得るので分離はアジリティを可能にする。第四に、A A から分離して複数のL A を割り当てる自由は、複数のL A がトポロジ的に重要な方法で要約され得るように割り当てられ、更に、スイッチが担う必要があるルーティング状態の量を制限することを意味している。一方、データセンタ内で動作しているサービス又はデータセンタのオペレータは、どんな方法においても、アプリケーションアドレスを割り当てる能力を妨げないことを望んでいる。

【 0 0 4 8 】

本発明の代替的な実施形態はL A 及びA A アドレスのための他の種類のデータを用いても良い。例えば、L A アドレスはIP v 4 で可能であり、A A アドレスはIP v 6 で可能であり、又は逆も同様である。又はIP v 6 はA A 及びL A の両方のアドレスのために使

10

20

30

40

50

用可能であり、又はIEEE 802.1 MACアドレスはAAアドレスとして使用可能である。一方、IPアドレス(v4又はv6)はLAアドレスのために用いられ、又は逆も同様である。また、アドレスは、例えば、IPアドレスを有するVLANタグ又はVRF識別子のように、異なる種類のアドレスと一緒に組み合わせることによって作り出されても良い。

#### 【0049】

次の記載は、下層のネットワーク構造を仮想化し、レイヤ2 LAN及びその上の何かのそれらのグループの他のサーバ316(1)~316(N)に接続され、ホストが比較的大きいデータセンタワイドレイヤ2 LANの一部であるアジャイルネットワークのサーバ316(1)~316(N)にイリュージョンを作り出すために、トポロジ、ルーティング

10

#### 【0050】

(アドレス決定及びパケット転送)

少なくともいくつかの実施例では、サーバ316(1)~316(N)が信じることを可能にするためそれらは同一のサービスで他のサーバと1つの大きなVLANを共有し、大きなイーサネットを悩ますブロードキャストARPスケーリング障害を解消し、以下に記載した解決策が提供される。なお、予め、次の解決策は下位互換があり、存在するデータセンタアプリケーションに対して透過であることが可能である。

#### 【0051】

20

(パケット転送)

AAは通常、ネットワークのルーティングプロトコルに示されても良い。よって、パケットを受信するサーバのために、パケットのソースは先ず、そのパケットをカプセル化することができ、ホストのために外側のヘッダの送り先をLAに設定する。LAアドレスを保持する装置に到達すると、パケットはデカプセル化され、送り先のサーバに配信される。1つの実施形態において、送り先のサーバのためのLAは、その送り先サーバが設置されている元のTORに割り当てられる。パケットがその送り先TORに到着すると、通常のレイヤ2の配信規則によれば、そのTORスイッチがそのパケットをデカプセル化して内側のヘッダの送り先AAに基づいてパケット配信することができる。代替的に、LAは物理的送り先サーバ又はそのサーバ上で動作する仮想マシンと関連付けられても良い。

30

#### 【0052】

(アドレス決定)

複数のAAアドレスが複数のサーバと同一のLANに存在することを信じるようにその複数のサーバを設定することができ、それでアプリケーションが初めてAAにパケットを送信するときホスト上のカーネルネットワークスタックは送り先のAAのためのブロードキャストARPLクエストを生成することができる。ソースサーバのネットワークスタックで動作するアジャイルエージェントはそのARPLクエストを捕らえてそれをアジャイルディレクトリサービスへのユニキャスト問い合わせ(クエリ)に変換することができる。アジャイルディレクトリサービスがその問い合わせに答えるとき、アジャイルディレクトリサービスはパケットが送られるべきLAを提供することができる。また、アジャイルディレクトリサービスはパケットを中継するために用いられ得る中間スイッチ又は中間スイッチのセットを提供することができる。

40

#### 【0053】

(ディレクトリサービスによるインターサービスアクセス制御)

サーバは、AAのためのパケットを送る必要がないTORのLAを得ることができないならば、パケットをAAに送信することができなくても良い。よって、アジャイルディレクトリサービス322(1)~322(N)は通信方針を実行することができる。アジャイルディレクトリサービスはルックアップリクエストを処理するとき、どのサーバがそのリクエストを作っているか、ソース及び送り先の両方が属するサービス、及びそれらのサービスの分離方針を知る。その方針が「拒否(deny)」であるならば、アジャイルディ

50

レクトリサービスはL Aを提供することを単に拒否することができる。上記したアジャイルネットワークアーキテクチャの有利な点はインターサービス通信が許可されたとき、IPゲートウェイに迂回されることなくパケットが送信サーバから受信サーバへ直接流れることができることである。これは、従来のアーキテクチャでの2つのV L A Nの接続とは異なっている。

#### 【0054】

(インターネットとのインタラクション)

度々、データセンタによって処理されるトラフィックのほぼ20%がインターネットへの又はインターネットからのトラフィックである。よって、データセンタネットワークがそれらの大きなボリュームを処理することができることは有利な点である。先ず、上記したアジャイルネットワークアーキテクチャが仮想レイヤ2ネットワークを実施するためにレイヤ3の構造を利用することは不思議に思えるかもしれないが、この1つの有利な点は、いくつかの従来提案されたネットワーク環境で要求されたようなヘッダの書き直しのためにゲートウェイサーバを経由することが強制されることなく、このアーキテクチャを有するアジャイルデータセンタネットワークを作り出すことが可能なスイッチの高速シリコンを横切って外部のトラフィックが直接流れることができることである。

#### 【0055】

インターネットから直接到達可能であることを必要とするサーバ(例えば、フロントエンドウェブサーバ)には2つのアドレス、L A及びA Aを割り当てることができる。L Aはインターネットワーク通信のために用いられても良い。A Aはバックエンドサーバを用いたイントラデータセンタ通信のために用いられても良い。L Aは、B G P (B o r d e r G a t e w a y P r o t o c o l)を介して伝えられかつ外部から到達可能であるプールから引き出されても良い。そのため、インターネットからのトラフィックはサーバに直接到達することができる。そのサーバから外部の送り先へのパケットはコアルータに向けてルート決めされ、E C M Pによって利用可能なリンク及びコアルータを介して拡がっても良い。

#### 【0056】

(ブロードキャストの処理)

上記したアジャイルネットワークアーキテクチャは下位互換性のアプリケーションに対してレイヤ2セマンティクスを提供することができる。これはブロードキャスト及びマルチメディアをサポートすることを含むことができる。アジャイルネットワークアーキテクチャのアプローチはA R P及びD H C P (D y n a m i c H o s t C o n f i g u r a t i o n P r o t o c o l)のようなブロードキャストの大部分の共通のソースを完全に除去することである。A R Pはアジリティエージェント320のA R Pパケットを捕らえて上記したようにアジャイルディレクトリサービスから情報を調べた後、応答を与えることによって処理されても良く、D H C Pパケットは従来のD H C Pリレーエージェント及びD H C Pサーバに送られたユニキャストを用いてT O Rで捕らえられても良い。他のブロードキャストパケットを処理するために、セット内の他のホストによって送信されたブロードキャストパケットを受信可能であるべきホストの各セットはIPマルチキャストアドレスを割り当てられても良い。このアドレスはディレクトリサービスによって割り当てられても良く、アジリティエージェントはそのディレクトリシステムを問い合わせることによってそのアドレスを知ることができる。

#### 【0057】

ブロードキャストアドレスに送信されたパケットは代わりにそのサービスのマルチキャストアドレスに進行するように修正されても良い。アジャイルネットワークアーキテクチャのアジャイルエージェントは混乱を防止するために限界ブロードキャストトラフィックのレートを定めても良い。アジャイルエージェントは、サーバが最新の時間間隔(例えば、過去1秒及び過去60秒)に亘って送信したブロードキャストパケットのレートの見積もりを維持し、サーバが各間隔の期間にブロードキャストパケットの設定数よりも多く送信することを防止することができる。許可されたパケット数を越えた分のパケットは次の

間隔まで引っ込められるか又は遅延させても良い。また、ネイティブIPマルチキャストがサポートされても良い。

#### 【0058】

複数のスイッチがレイヤ3のルータとして動作する実施形態の潜在的な有利な点は、マルチキャストグループに属する全てのマシン又はマシンに対してマルチキャストグループにアドレス指定したパケットの配信を実施することが特に容易であることである。PIM-BIDIR等の存在するIPマルチキャストルーティングプロトコルのいずれかは複数のスイッチに設定されても良い。これは、それらにマルチキャストグループに属する各ホスト又はマシンでエンドポイントを有するマルチキャスト分散ツリーを計算することをさせる。ホスト上のアジャイルエージェント、マシン、又はサーバは、通常、IGMPをそのデフォルトゲートウェイに送信することによって適切なマルチキャストグループの一部であるとしてホスト、マシン、又はサーバを登録する。それ故、そのマルチキャストルーティングプロトコルはホスト、マシン、又はサーバをそのマルチキャストグループのための分散ツリーに加えることを処理する。レイヤ2で動作するスイッチはマルチキャストグループ毎にVLANのような様々な機構を用いることができ、或いは、ネットワークを通して存在するパケットを、アジャイルエージェントのホスト、マシン又はサーバが受信すべきでないパケットを取り除く各ホスト、マシン、又はサーバ上のアジャイルエージェントで満たさせることができる。

#### 【0059】

(マルチパスルーティングでのランダム化)

上記したアジャイルネットワークアーキテクチャは、少なくともいくつかの実施形態では、2つの関係した機構、VLB及びECMP(Equal Cost Multipath)を用いて不安定さに対処するためにランダム化を用いる原理を活用/利用することができる。両方の目的は似ており、持続的な輻輳状態を減少又は防止するためにVLBは複数の中間ノードに亘ってランダムにトラフィックを分散し、ECMPは等コストパスを越えるトラフィックを発する。以下に詳細に説明されるように、VLB及びECMPは、各々がその他において制限を克服するために用いられる点で相補的である。両方の機構は、パケットの送り元がネットワークを越えるパスの選択に影響を与えるために使用可能な制御を提供することができる。アジャイルエージェントはそれらの制御が輻輳状態を避けるために利用されることを可能にする。

#### 【0060】

図5は、図3で紹介されたアジャイルネットワークアーキテクチャ300のサブセットを示している。図5はサーバ通信に対するサーバの更なる詳細を提供する。この例はサーバ316(5)と通信するサーバ316(1)を含む。送信サーバ316(1)及び送り先サーバ316(5)は、VLANとして機能し10.128/9のアプリケーションアドレスを有するサーバプール328内で動作することができる。中間スイッチ310(1)~310(N)はリンクステートネットワーク328内に存在する。

#### 【0061】

アジャイルネットワークアーキテクチャ300は、VLBの利益がランダムに選択された中間ノードを外れるようにパケットを強制することによって達成されることを可能にする。この場合に、送り元のアジャイルエージェント320(1)は各パケットを中間スイッチ310(1)~310(N)に対してカプセル化することによりこれを実施することができる。その中間スイッチはそのパケットを送り先のTOR(この場合、314(1))に送る。ここで、パケットは先ず、中間スイッチのうちの、310(2)のような1つに配信され、そのスイッチによってデカプセル化され、TOR314(1)のLAに配信され、再びデカプセル化され、そして、最後に送り先サーバ316(5)に送信されても良い。

#### 【0062】

アジャイルエージェント320(1)は活性状態の中間スイッチ310(1)~310(N)のアドレスを知っているならば、パケットを送信するときそれらの中からランダム



に選択することができる。しかしながら、これは、中間スイッチが機能しなくなったとき潜在的な何十万のアジャイルエージェントの更新を必要とする可能性がある。代わって、同一のL Aアドレスが多数の中間スイッチに割り当てられても良い(この場合、L Aアドレス10.0.0.5)。そのアジャイルディレクトリサービス(図3に示された)は、このエニキャストアドレスを1以上のルックアップ結果の一部としてアジャイルエージェント320(1)に戻すことができる。ECMPはそのエニキャストアドレスにカプセル化されたパケットを活性状態の中間スイッチ310(1)~310(N)のうちの1つへ配信する処理を行うことができる。もしスイッチが機能しなくなったならば、ECMPは再動作してアジャイルエージェントを通知する必要性を除去することができる。

#### 【0063】

しかしながら、ECMPはスケーリング限界を有しても良い。従来のスイッチは、今日、16wayのECMPをサポートすることができ、また、256wayのECMPスイッチが利用可能、又はまもなく利用可能である。ECMPを使用するよりも利用可能なパスが存在することが生じたならば、VLBカプセル化は補償することができる。1つの解決法はいくつかのエニキャストアドレスを定義することであり、個々のエニキャストアドレスはECMPが蓄積できるように多くの中間スイッチ310(1)~310(N)と関連付けされる。送り元は負荷を分散するためにエニキャストアドレスに亘ってハッシュを行うことができ、スイッチが機能しなくなると、個々のサーバが通知される必要性がないようにそのエニキャストアドレスがディレクトリシステムによって他のスイッチに再割り当てされても良い。説明のために、この態様はディレクトリシステムによって提供されたネットワーク制御機能として考慮されても良い。

#### 【0064】

上記したVLBベースの気がつかないルーティングが、折り畳み Clos ネットワークトポロジについて純粋なOSPF/ECMP機構を用いて実施されても良い。そのような設定は中間スイッチにおいてデカプセル化のサポートを必要としない。例えば、Nが各TOR上のアップリンクの数であるならば、集合スイッチはセットにグループ化されても良い。いくつかの場合には、それらのセットの各々はNのスイッチを正確に含むことができる。各TORは1つのセット内にNのスイッチ全てにアップリンクを有するか又は1つのセット内にスイッチのどれもアップリンクを有しないとすることができる。TORのこの書き込みで、TOR間のルート決定のためにOSPF及び/又はECMPのようなプロトコルが用いられるときでさえ、サーバの入/出制限の影響下にある任意のトラフィックのための帯域幅保証が保持し続けることを示すことができる。

#### 【0065】

TOR間のルート決定のためのOSPF又はECMPの使用は、集合スイッチの同一のセット内における2つのTOR間のパケットのようないくつかのパケットが、中間スイッチを介して進行しないパスを利用するようにさせることができる。よって、ソースと送り先との間の最も短いパスに従い、同一の集合スイッチ又は複数の集合スイッチに接続された同一のTOR下、又は複数のTOR下のサーバ間のトラフィックの早期のターンアラウンドを可能にするようにそれらのパスを「アーリターンアラウンドパス(early turnaround paths)」と称することができる。それらのトラフィックフローはコア集合/中間ネットワークのいずれにも入ることを必要としない。

#### 【0066】

アーリターンアラウンドパスを用いることの潜在的な利益は、トラフィック(例えば、外部)の他のクラスのためのコアにおける能力を制限から解くことを含むことである。自由にされた能力は、存在するアプリケーションがクロス(cross)、例えば、TORトラフィックを最小化するために書き込まれているとき、ほぼ「平均」の場合にあることができる。他の方法を見ると、これはコアが何らかの要因に基づいて設定されることを可能にし、サーバ間のトラフィックのために旨く動作することができる。また、アーリターンアラウンドパスの使用は装置の広範囲が中間スイッチとして使用されることを可能にし、それらのスイッチについて結果的に低コストにすることができる。

10

20

30

40

50

## 【 0 0 6 7 】

## ( 輻輳状態に対する対処 )

E C M P 及び V L B の両方を用いて、大きなフローが同一のリンク及び中間スイッチ各々にハッシュされる機会があっても良く、それは輻輳状態を生じさせるかもしれない。もし、これが起きるならば、送信アジャイルエージェントは、E C M P がネクストホップ、すなわちパケットが通るべき次のスイッチを選択するために使用するフィールドの値を変更することによって、そのフローがアジャイルネットワークを介して利用するパスを変更することができる。この点において、アジャイルエージェントは、周期的に大きなフローをリハッシュする、又は激しい輻輳イベント（例えば、フルウィンドウ損失）又は E C N が T C P によって検出されたとき、又は閾値のバイト / パケット数を送信 / 受信した後のような状況を簡単な機構で検出して処理することができる。

10

## 【 0 0 6 8 】

## ( ホスト情報の維持 )

上記したアジリティネットワークアーキテクチャに応じて実施されるネットワークシステムは、データセンタ負荷のために設計されたスケーラブル（拡張可能な）、信頼がある、及び / 又は高性能なストア又はディレクトリシステムを使用することができる。アジリティネットワークアーキテクチャに応じて実施されるネットワークは、一様な高能力、パフォーマンス分離、レイヤ 2 のセマンティック、及びサービス間の通信分離、それらの 4 つのプロパティの 1 以上を所有することができる。また、そのネットワークは正常な劣化を示すことができ、ここで、そのネットワークは故障後にどんな能力が残っていても使用することを継続することができる。よって、そのネットワークは故障にあっても信頼性 / 回復力がある。この点について、そのようなネットワークのディレクトリシステムは 2 つの潜在的なキー機能、( 1 ) A A 対 L A マッピングのためのルックアップ及び更新、及び ( 2 ) 例えば、ライブ仮想マシンマイグレーション等の待ち時間検出動作をサポートすることができる反応型キャッシュ更新機構を提供することができる。

20

## 【 0 0 6 9 】

## ( 要求の特徴付け )

ディレクトリシステムのためのルックアップ負荷は頻繁で集中的である。サーバは、A A 対 L A マッピングのためのルックアップを生成する各フローで短時間の期間に何千又は何万の他のサーバまでと通信することができる。更新のために、負荷を故障及びサーバスタートアップイベントによって駆動することができる。多くの故障は一般に、サイズのには小さく、大きな相関性がある故障はおそらくまれである。

30

## 【 0 0 7 0 】

## ( パフォーマンス要求 )

負荷のバースト性特質は、多くの接続を素早く確立するようにルックアップが高い処理能力及び低応答時間を要求しても良いことを暗示する。ルックアップは初めてサーバと通信するために要求された時間を増加させるので、その応答時間はできるだけ小さく保持されるべきである。例えば、1 / 1 0 0 秒が妥当な値である。しかしながら、更新のために、潜在的なキー要求は信頼性であって良く、応答時間はあまり気に掛けなくても良い。更に、更新は通常、前もって予定されるので、高い処理能力はバッチ処理の更新によって達成されても良い。

40

## 【 0 0 7 1 】

## ( 一貫性の考慮 )

従来のレイヤ 2 のネットワークにおいて、A R P は A R P 中断のために最終的一貫性を提供することができる。加えて、ホストは余計な A R P を発することによってその到着を伝えることができる。極端な例として、上記したアジリティネットワークアーキテクチャに応じて実施されたネットワークにおけるライブ仮想マシン ( V M ) マイグレーションを考慮する。V M マイグレーションはステイルマッピング ( A A 対 L A ) の高速更新を利用することができる。V M マイグレーションの潜在的な目的は、ロケーション変更を越えて進行している通信を保持することであることができる。それらの考慮は、信頼性ある更新

50

機構を提供することができる限り、A A対L Aマッピングの弱い又は最終的な一貫性が受け入れ可能であることを暗示している。

#### 【0072】

(アジャイルディレクトリシステム又はサービス設計)

パフォーマンスパラメータ及びルックアップの負荷パターンは、更新の各々から異なっているとしても良い。よって、図6に示された2層になったアジャイルディレクトリサービスアーキテクチャ600を考慮する。この場合において、アジャイルディレクトリサービスアーキテクチャ600は、アジャイルエージェント602(1)~602(N)、ディレクトリサービスモジュール604(1)~604(N)、及びRSM(Replicated State Machine:複製ステートマシン)サーバ606(1)~606(N)を含んでいる。特に、この例では、個々のディレクトリサービスモジュールは専用のコンピュータ608(1)~608(N)各々上で実施される。他の実施例では、ディレクトリサービスモジュールは、他のシステム機能を実行するコンピュータ上で明白になっても良い。この実施例では、ディレクトリサービスモジュールの数は一般に全体のシステムサイズと比較して少ない。例えば、1つの実施例は100Kサーバ(すなわち、図3のサーバ316(1)~316(N))に対してほぼ50~100のディレクトリサービスモジュールを用いることができる。この範囲は説明の目的で与えられ、重要ではない。

10

#### 【0073】

ディレクトリサービスモジュール604(1)~604(N)は、A A対L Aマッピングをキャッシュすることができる読み取り最適化された、複製ディレクトリサービスとして考慮され得る。そのディレクトリサービスモジュール604(1)~604(N)はアジャイルエージェント602(1)~602(N)、及び少量(例えば、ほぼ5~10サーバ)の読み取り最適化された、RSMサーバ606(1)~606(N)と通信することができる。RSMサーバ606(1)~606(N)はA A対L Aマッピングの強力な一貫して信頼性あるストア(保存部)を提供することができる。

20

#### 【0074】

ディレクトリサービスモジュール604(1)~604(N)は少ない待ち時間、高い処理能力、及び高ルックアップレートのための高い利用可能性を確かに行うことができる。一方、RSMサーバ606(1)~606(N)は、少なくともいくつかの実施形態において、更新の少ないレートのために、Paxos合意アルゴリズム等を用いて強力な一貫性及び耐久性を確かに行うことができる。

30

#### 【0075】

個々のディレクトリサービスモジュール604(1)~604(N)はRSMサーバ606(1)~606(N)で保存されたA A対L Aマッピングをキャッシュすることができる。キャッシュした状態を用いてアジャイルエージェント602(1)~602(N)からのルックアップに対して独立して応答することができる。強力な一貫性は必要条件でなくとも良いので、ディレクトリサービスモジュールは定期的に(例えば、30秒毎)RSMサーバとそのローカルマッピングをゆっくりと同期させることができる。高利用可能性及び少ない待ち時間を同時に達成するために、アジャイルエージェント602(1)~602(N)はランダムに選択されたディレクトリサービスモジュール604(1)~604(N)の数k(例えば、2)に対してルックアップを送信することができる。多くの応答が受信されたならば、アジャイルエージェントは単に最も速い応答を選択してそれをそのキャッシュに保存する。

40

#### 【0076】

また、ディレクトリサービスモジュール604(1)~604(N)は、ネットワーク設定(プロビジョニング)システムからの更新を処理することができる。一貫性及び耐久性のために、更新は1つのランダム選択のディレクトリサービスモジュールに送信されても良く、RSMサーバ606(1)~606(N)にライトスルーされても良い。特に、更新について、ディレクトリサービスモジュールは先ず、RSMに更新を送ることができる。

50

る。確実に、その R S M は個々の R S M サーバに対してその更新を複製し、そして、ディレクトリサービスモジュールにアクナリッジメントをもって応答することができ、それはそのアクナリッジメントを元のクライアントに順に送り戻すことができる。

#### 【0077】

一貫性を拡張するための潜在的最適化として、ディレクトリサービスモジュール 604 (1) ~ 604 (N) は、少数の他のディレクトリサービスモジュールにアクナリッジメントを行った更新を任意に広めることができる。元のクライアントは中断 (例えば 2 秒) 内にアクナリッジメントを受信しなかったならば、そのクライアントは同一の更新を他のディレクトリサービスモジュールに送信して、信頼性及び / 又は利用可能性と引き換えに応答時間を提供することができる。

10

#### 【0078】

また、ディレクトリサービスの他の実施形態は可能である。例えば、DHT (D i s t r i b u t e d H a s h T a b l e) がディレクトリサーバと、DHT にエントリーとして保存された A A / L A マッピングとを用いて構成されても良い。また、パフォーマンスが前述した実施形態と同然でないか、又は強力な一貫性がなくて良いけれども、アクティブディレクトリ (A c t i v e D i r e c t o r y) 又はライトウエイトディレクトリシステム (L i g h t w e i g h t D i r e c t o r y S y s t e m) 等の他の存在するディレクトリシステムを用いても良い。

#### 【0079】

(最終的一貫性の確保)

20

A A 対 L A マッピングはディレクトリサービスモジュール及びアジャイルエージェントのキャッシュでキャッシュされ得るので、アップデートは一貫性がないことに導く可能性がある。サーバ及びネットワークを浪費することなく一貫性がないことを解決するために、同時にスケラビリティ (拡張性) 及びパフォーマンスの両方を確実にするために反応型キャッシュ更新機構を用いることができる。キャッシュ更新プロトコルはキー観察を利用することができる。ステイルホストマッピングは、そのマッピングがトラフィックを配信するために用いられるときだけに、修正される必要がある。特に、ステイルマッピングが用いられるとき、いくつかのパケットはステイル L A (送り先サーバをもはや管理しない T O R 又はサーバ) に到着することができる。その T O R 又はサーバはそのような配信不可能なパケットをディレクトリサービスモジュールに送り、例えば、ユニキャストを介して、ソースサーバのキャッシュにおいてステイルマッピングを選択的に修正するようにディレクトリサービスモジュールを動作させることができる。更新の他の実施形態において、ディレクトリサービスは影響を受けたサーバと通信するために許可されている全てのサーバグループに対する更新をマルチキャストによって行っても良い。

30

#### 【0080】

(更なる実施例)

(ロードバランシングの最適性)

上記したように、V L B のようなロードバランシング技術は不安定さに対処するためにランダム化を用いることができ、トラフィックパターン (最良の場合及び最悪の場合の両方を含む) を平均の場合に変えることによって最良の場合のトラフィックパターンのいくつかのパフォーマンスを潜在的に犠牲にする。このパフォーマンス損失は、より最適なトラフィックエンジニアリングシステム下にあるよりも高いいくつかのリンクの利用としてそれ自身を明らかにすることができる。しかしながら、活性状態にあるデータセンタ負荷についての評価は、更に複雑なトラフィックエンジニアリングスキームと比較したとき比較的小さい能力損失と関連づけられ得る、V L B のようなロードバランシング技術の簡潔さ及び普遍性を示している。

40

#### 【0081】

(レイアウト設定)

図 7 ~ 図 9 は説明したアジャイルネットワークアーキテクチャに応じて実施されるデータセンタネットワークの 3 つの可能なレイアウト設定を示している。図 7 ~ 図 9 において

50

、図面ページ上の空間制限のためにTORは関連したサーバなしに示されている。

【0082】

図7はオープンフロアブランデータセンタレイアウト700を示している。データセンタレイアウト700はTOR702(1)~702(N)、集合スイッチ704(1)~704(N)、及び中間スイッチ706(1)~706(N)を含んでいる。図7において、TOR702(1)~702(N)は、センタ「ネットワークケージ」708を取り囲むとして示され、(例えば、銅及び/又はファイバケーブル等を用いて)接続されても良い。集合スイッチ704(1)~704(N)及び中間スイッチ706(1)~706(N)各々はネットワークケージ708内側でごく接近してレイアウトされても良く、それらの間接続のために銅ケーブルの使用を可能にする(銅ケーブルは低コストで極太であり、ファイバに対して低距離範囲用である)。ネットワークケージ708内側のケーブルの数を(例えば、約4分の1)、例えば、QSFP(Quad Small Form Pluggable)規格のような適切な規格を用いて1つのケーブルに10Gリンクの数(例えば、4)と一緒に構築することによってそれらのコスト(例えば、約2分の1)と共に減らすことができる。

【0083】

オープンフロアブランデータセンタレイアウト700において、中間スイッチ706(1)~706(N)は、ネットワークケージ708内の中央に配置され、集合スイッチ704(1)~704(N)は、中間スイッチ706(1)~706(N)とTORスイッチ702(1)~702(N)(及び関連したサーバ)との間に介在されている。

【0084】

オープンフロアブランデータセンタレイアウト700は望んだように拡張可能である。例えば、追加のサーバラックは、サーバラックを作り出すためにTOR702(1)~702(N)を有するサーバの形で計算装置を関係付けることによって加えられても良い。そして、サーバラックはネットワークケージ708の集合スイッチ704(1)~704(N)に接続されても良い。他のサーバラック及び/個々のサーバはオープンフロアブランデータセンタレイアウトによって提供されたサービスを妨げることなく削除されても良い。

【0085】

図8はモジュラー化されたコンテナベースのレイアウト800を示している。レイアウト800はTOR802(1)~802(N)、集合スイッチ804(1)~804(N)、及び中間スイッチ806(1)~806(N)を含んでいる。この場合において、中間スイッチ806(1)~806(N)はそのレイアウトのデータセンタインフラ808に含まれる。集合スイッチ及びTORスイッチはそのデータセンタインフラに接続されているプラグ着脱可能なコンテナとして関連付けられても良い。例えば、集合スイッチ804(1)及び804(2)は、データセンタインフラ808に接続され得るプラグ着脱可能なコンテナ810(1)内でTOR802(1)及び802(2)と関連付けられても良い。同様に、集合スイッチ804(3)及び804(4)は、プラグ着脱可能なコンテナ810(2)内でTOR802(3)及び802(4)と関連付けられ、集合スイッチ804(5)及び804(N)は、プラグ着脱可能なコンテナ810(N)内でTOR802(5)及び802(N)と関連付けられる。

【0086】

図7のように、図8において、サーバラックを作り出すためにTORと関連付けられているサーバは、図面ページの空間制限のために示されていない。更に、空間制限のために、2つの集合スイッチ及び2つの中間スイッチだけがプラグ着脱可能なコンテナ毎に示されている。勿論、他の実施例はそれらのコンポーネントのいずれか又は両方より多い又は少ないコンポーネントを用いることができる。また、他の実施例は本明細書に示された3より多い又は少ないプラグ着脱可能なコンテナを用いることができる。興味ある1つの特徴は、レイアウト800が1のケーブルの束を各プラグ着脱可能なコンテナ810(1)~810(N)からデータセンタスパイン(すなわちデータセンタインフラ808)にも

たらすためにそれ自身を貸すことができることである。要約すると、データセンタインフラ 808 はレイアウト 800 が個々のプラグ着脱可能なコンテナ 810 (1) ~ 810 (N) 加えるか、又は除去することによってサイズとして拡張又は縮小することを可能にする。

#### 【0087】

図 9 は「インフラ減少」及び「コンテナ化」データセンタレイアウト 900 を示している。レイアウトは多くのコンテナ 908 (1) ~ 908 (N) に配置された TOR 902 (1) ~ 902 (N)、集合スイッチ 904 (1) ~ 904 (N)、及び中間スイッチ 906 (1) ~ 906 (N) を含んでいる。例えば、TOR 902 (1) ~ 902 (2)、集合スイッチ 904 (1) ~ 904 (2)、及び中間スイッチ 906 (1) はコンテナ 908 (1) に配置されている。

10

#### 【0088】

コンテナ 908 (1) ~ 908 (N) は「インフラ減少」及び「コンテナ化」データセンタレイアウト 900 の実現を可能にする。このレイアウト 900 はコンテナ 908 (1) 及び 908 (3) の個々のペア間でケーブル束 910 (1) を繋ぐことと関連付けられても良い。他のケーブル束 910 (2) はコンテナ 908 (2) 及び 908 (N) の個々のペア間で繋ぐことができる。個々のケーブル束 910 (1), 910 (2) は、コンテナ 908 (1) 内の集合スイッチ 904 (1), 904 (2) をコンテナ 908 (3) 内の中間スイッチ 906 (3) と接続するリンク、また逆の場合も同じリンクを担うことができる。

20

#### 【0089】

要約すると、個々のコンテナ 908 (1) ~ 908 (N) は複数のスイッチを含むことができる。それらのスイッチは、相補型プラグ着脱可能なコンテナに配置されている TOR 902 (1) ~ 902 (N)、集合スイッチ 904 (1) ~ 904 (N)、及び中間スイッチ 906 (1) ~ 906 (N) を含むことができる。相補型プラグ着脱可能なコンテナのペアは、第 1 のプラグ着脱可能なコンテナの集合スイッチを第 2 のプラグ着脱可能なコンテナの中間スイッチに接続することによって、また逆の場合もケーブル束を介して同じによって結合されても良い。例えば、コンテナ 908 (1) はケーブル束 910 (1) を介してコンテナ 908 (3) に接続されても良い。特に、ケーブル束は、コンテナ 908 (1) の集合スイッチ 904 (1) 及び 904 (2) をコンテナ 908 (3) の中間スイッチ 906 (3) と接続することができる。同様に、束 910 (1) はコンテナ 908 (3) の集合スイッチ 904 (5) 及び 904 (6) をコンテナ 908 (1) の中間スイッチ 906 (1) と接続することができる。

30

#### 【0090】

少なくともいくつかの実施形態において、アジャイルネットワークアーキテクチャは次のコンポーネント (1) ~ (5) から構成することができる。(1) トポロジと一緒に接続されたスイッチのセット、(2) スwitchの 1 以上に各々接続されたサーバのセット、(3) サーバがパケットを他のサーバに送信することを望んだときリクエストが作られ、サーバ (又はサーバの代理アジャイルエージェント) がスイッチのトポロジを越えることができるように送信することを望むパケットを加えるか、又はカプセル化することにおいて使用する情報で応答するディレクトリシステム、(4) パケットがスイッチによって引っ込められてリンクに送られるほどいずれかのリンク上での利用の増大を減少又は防止するネットワークにおいて輻輳状態を制御する機構、(5) ディレクトリサービスと通信し、必要に応じてパケットをカプセル化し、アドレス指定し、又はデカプセル化し、そして、必要に応じて輻輳制御に加わるサーバ上のモジュール。

40

#### 【0091】

少なくとも 1 つの実施形態において、機能 (1) ~ (5) を提供する各サーバについてアジャイルエージェントが存在することができる。(1) パケットを送り先に送り、システム内のこのサーバに登録する等のために利用されるカプセル化情報を取り出すためにアジャイルディレクトリサービスと通信する、(2) 必要に応じて代替手段のセットの中か

50

ら（例えば、中間スイッチの中から）ランダム選択を行い、それらの選択をキャッシュする、（３）パケットをカプセル化／デカプセル化する、（４）ネットワークから輻輳表示を検出して応答する。よって、少なくともいくつかの実施形態において、それらの機能はネットワーク内の複数のサーバと複数のスイッチとの間で分散されても良い。例えば、デフォルトのルート決定はパケットを（中間スイッチ等の）スイッチのセット、及びパケットが横断する中間スイッチ上のパケット毎に実施された上記のリストの機能を管理するために用いられても良い。

#### 【００９２】

少なくともいくつかの実施形態において、本明細書で説明されたアジャイルネットワークアーキテクチャを実施することは、ネットワーク内の各スイッチがネットワーク内の他のスイッチにパケットを送信することができるように、データセンタ内でスイッチのセットの中からネットワークを作り出すことを含むことである。他のサーバと通信するためにサーバによって用いられたアドレスと同じ種類のアドレスをそれらのスイッチ又はそのネットワークがそれら自身の中からパケットを管理するために使用することが必要でない。例えば、MACアドレス、IPv4アドレス、及び／又はIPv6アドレスは全て適切であれば良い。

#### 【００９３】

アジャイルネットワークの少なくとも１つの実施形態において、データセンタ内のスイッチのセットの中における１つの考慮は、それらの各々をIPアドレス、IPv4又はIPv6のいずれかを用いて設定されること、及びそれらを１以上の標準レイヤ３ルーティングプロトコルを動作させるために一般的な例OSPF(Open-Shortest Path First)、IS-IS(Intermediate System-Intermediate System)又はBGP(Border Gateway Protocol)を用いて設定することである。そのような実施形態の利益は、ディレクトリシステムがトポロジに対して反応してサーバに大部分の変更を知らせる必要がないように、スイッチ間でパケットを送るためにネットワークの能力を維持するそのルーティングプロトコルによって作り出されるネットワークの制御プレーンを用いて、ネットワークとディレクトリシステムとの間の接続が減少されることである。

#### 【００９４】

代替的又は追加的に、ディレクトリシステムはネットワークのトポロジを監視し（例えば、スイッチ及びリンクの調子（ヘルス）を監視すること）、トポロジが変更するとサーバに提供するカプセル化情報を変更することができる。また、ディレクトリシステムは、前に応答を送信したサーバにそれらの応答がもはや有効でないことを通知しても良い。その代替に亘って第１の実施形態の潜在的な利点は、ディレクトリシステムがトポロジに対して反応してサーバに大部分の変更を知らせる必要がないように、スイッチ間でパケットを送るためにネットワークの能力を維持するそのルーティングプロトコルによって作り出されるネットワークの制御プレーンを用いて、ネットワークとディレクトリシステムとの間の接続が減少されることである。要約すれば、ネットワークパフォーマンスに関係した１以上のパラメータを監視することによってパケット配信遅延を減少又は避けることができる。そのパラメータは、特定のパスに亘る通信障害等のネットワークイベントを表しても良い。

#### 【００９５】

１つの実施形態において、ネットワークの複数のスイッチは、LAアドレスのサブネットから引き出されたIPv4のアドレスを用いて設定される。そのスイッチはOSPFルーティングプロトコルを動作させるように設定される。そのスイッチのアドレスはOSPFプロトコルによってスイッチの中で割り振られる。番号が付けられていないOSPFのインターフェース拡張はOSPFプロトコルによって分配される情報量を減らすために用いられても良い。各トップオブブラック(TOR)スイッチのサーバ向き部分はVLAN(Virtual Local Area Network)の一部であるようにそのスイッチ上で設定される。AA空間を含むサブネットはサーバ向きのVLANに割り当てられた

10

20

30

40

50

ようにスイッチに設定される。このVLANのアドレスはOSPFには分配されず、VLANは通常、トランキングされない。サーバへ行くパケットはそのサーバが接続されたTORに対してカプセル化される。このTORはそれらカプセル化パケットを受信するとそのパケットをデカプセル化し、そして、サーバの送り先アドレスに基づいてサーバ向きのVLANにそれらを送る。そのためそのサーバは通常のLANのように、パケットを受信する。

#### 【0096】

他の実施形態においては、TORスイッチのサーバ向きのVLANにAAサブネットを設定する代わりに、各TORに特有のLAサブネットはそのサーバ向きのVLANに割り当てられる。このLAサブネットはOSPFによって分配される。TORに接続されたサーバは少なくとも2つのアドレスを用いて設定される。LAサブネットから引き出されるLAアドレスはサーバ向きのVLANに割り当てられ、それはVLANの一部とAAアドレスである。サーバ行きのパケットはサーバ上に設定されているLAに対してカプセル化される。サーバ上のモジュールはカプセル化パケットを受信するとデカプセル化して、パケットに含まれたAAアドレスに基づいた送り先であるサーバ上の仮想マシン又は処理にローカル的に配信することができる。

10

#### 【0097】

他の実施形態において、TORスイッチはレイヤ2のスイッチとして動作しても良く、一方、集合レイヤスイッチはレイヤ3のスイッチとして動作しても良い。この設計は、潜在的に安価なレイヤ2スイッチがTORスイッチとして使用されることを可能にしても良く（多くのTORスイッチが存在する）、一方、レイヤ3の機能は比較的少ない数の集合レイヤスイッチで実施されても良い。この設計において、デカプセル化機能はレイヤ2のスイッチ、レイヤ3のスイッチ、送り先サーバ、又は送り先仮想マシンで実行されても良い。

20

#### 【0098】

いくつかの実施形態において、追加のアドレスはスイッチ上に設定、又はOSPFのようなルーティングプロトコルを介して分配されても良い。それらのアドレスは通常、トポロジ的に重要である（すなわち、LA）。アドレスは通常、パケットをインフラサービス - すなわちサーバ、スイッチ、又は追加のサービスとして知られていることを提供するネットワーク装置に送るために用いられる。そのようなサービスの例はロードバランサ（それらはF5からのBigIPのようなハードウェアベース又はソフトウェアベースのロードバランサであって良い）、S-NAT(Source Network Address Translators)、ディレクトリシステムの一部であるサーバ、DHCPサービスを提供するサーバ、又は（インターネット又は他のデータセンタ等の）他のネットワークへのゲートウェイを含む。

30

#### 【0099】

1つの実施形態において、各スイッチはBGPプロトコルを用いてルートリフレクタライアントとして設定されても良い。追加のアドレスはそのルートリフレクタ上でそれらを設定し、BGPがそれらをスイッチに分配することを許可にすることによってスイッチへ分配される。この実施形態は、追加のアドレスを加える又は除去することがそのスイッチのルーティングプロセッサを過負荷にするOSPF再計算をさせないという利益を有する。

40

#### 【0100】

他の実施形態において、ネットワーク内の輻輳状態を制御する機構は複数のサーバそれぞれ自身上で実施される。適切な機構はTCP(Transport Control Protocol)のようなものであり、ここで送り先にサーバによって送られるトラフィックは、ネットワークが担うことができるらしいレートにサーバによって制限される。TCPのようなプロトコルの使用の改善は次に説明される。代替的な実施形態においては、スイッチ上のQoS(Quality of Service)機構が輻輳状態制御のために用いられても良い。そのような機構の例はWFQ(Weighted Fair Qu

50



e u i n g ) 及びその派生物、R E D ( R a n d o m E a r l y D e t e c t i o n ) 、 R S V P 、 X C P ( e X p l i c i t C o n t r o l P r o t o c o l ) 、 及び R C P ( R a t e C o n t r o l P r o t o c o l ) を含む。

#### 【0101】

少なくとも1つの実施形態において、サーバ上のモジュールはアジャイルネットワークから受信されているパケットを観察し、受信したパケットから得た又は推測した情報に基づいてパケットの送信又はパケットのカプセル化を変更する。アジャイルエージェントは(1)パケットが送信されるレートを減らすためにパケットの送信を変更する、又は(2)ネットワークを通して異なるパスを利用するようにパケットのカプセル化を変更する、ことによってネットワークにおける輻輳状態を減らすことができ、それは、パケットのカプセル化及びアドレス指定を先ず選択するとき作った可能な選択肢の中からのランダム選択のいずれか又は全てを作り変えることによって達成されても良い。

10

#### 【0102】

アジャイルエージェントがすることができる観察及びその反応の例は、次の(1)~(5)を含む。(1)アジャイルエージェントはTCPパケットのフルウィンドウの損失を検出するならば、パケットがネットワークを通して利用するパスを再ランダム化する。これは、特に、パケットによって利用されたパスの変更が並べ直されたパケットを送り先によって受信されることが起きないようにフローを前に送られた全てのパケットがネットワークから抜け出たと思われると同時に異なる(望ましくは輻輳してない)パスにフローを置くようにすると都合が良い。(2)アジャイルエージェントはパケットによって利用されるパスを周期的に再ランダム化することができる。(3)アジャイルエージェントは1つのフローによって達成されている実効レートを計算し、そのレートが期待閾値より低いならば再ランダム化することができる。(4)アジャイルエージェントはECN(E x p l i c i t C o n g e s t i o n N o t i f i c a t i o n)マークの受信パケットを監視してそのレートを減らすか又はその送り先へのパケットのパスを再ランダム化することができる。(5)スイッチは輻輳状態に入った又は入ろうとしているリンクを検出するようにロジックを実行し(例えば、IEEE QCN及び802.1auにおけるように)、通知を上流スイッチ及び/又はサーバに送ることができる。それらの指示を受信するアジャイルエージェントはそれらのパケットのレートを減らすか又はパケットのパスを再ランダム化することができる。

20

30

#### 【0103】

説明した複数の実施形態の1つの有利な点は、VM(仮想マシン)が1つのサーバから他のサーバへ再配置されても良く、一方、同一のIPアドレスの使用を保持しても良いので、それらが仮想マシン(VM)のライブマイグレーションを可能にすることである。ディレクトリシステムは、移動の間に再配置されるサーバのVMのIPアドレスに向かうパケットを管理するために単に更新されても良い。ロケーションにおける物理的な変更は行われている通信を妨げる必要はない。

#### 【0104】

少なくとも1つの実施形態において、ネットワークの能力の一部は、好ましいサービスが多数又は少数のパス、又はサービスの他のセットによって使用されるパスから分離されたパスのセットに広がったそれらのパケットを有するように、分割割合の様な計算によってネットワークに亘って動作するサービスのセットに予約され又は優先的に割り当てられても良い。プリファレンス又はQoSの多数のクラスはこの同じ技術を用いて作り出されても良い。

40

#### 【0105】

図10は、本コンセプトの少なくともいくつかの実施例と一致するアジャイルネットワークング技術又は方法1000のフローチャートを示している。方法1000が説明される順番は、限定として構成されるつもりではなく、その方法又は代替方法を実施するために様々な説明のブロックを組み合わせても良い。更に、方法は、計算装置がその方法を実施することができるように、適切なハードウェア、ソフトウェア、ファームウェア、又は

50

それらの組み合わせで実施されても良い。1つの場合において、その方法は、計算装置のプロセッサによる実行が計算装置にその方法を行わせるように、命令のセットとしてコンピュータ可読記憶媒体に保存される。他の場合において、その方法はA S I Cによる実行のためにA S I Cのコンピュータ可読記憶媒体に保存される。

【0106】

1002では、その方法はパケットを送り先に送るために利用されるカプセル化情報を得る。

【0107】

1004では、その方法は複数のスイッチ等の利用可能なハードウェアを通してパスを選択する。

【0108】

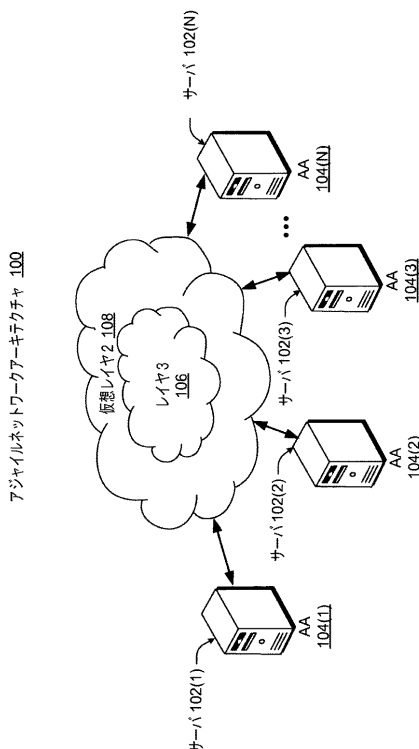
1006では、その方法はパスを介した配信のためにパケットをカプセル化する。

【0109】

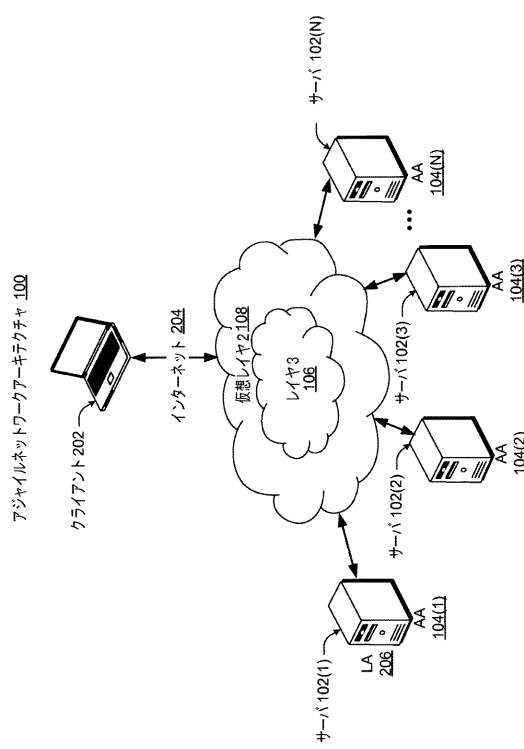
1008では、その方法は輻輳状態の表示に対して監視する。例えば、その方法はネットワークパフォーマンスに関係したパラメータを監視することができる。例えば、TCPは、パケット伝送レートに関係した更新、及び/又は輻輳状態に関係するネットワークパラメータとして動作可能なネットワークコンポーネント上の負荷を提供することができる。その方法は輻輳状態が検出されたときパスを再選択するか又は他の動作をとることができる。

10

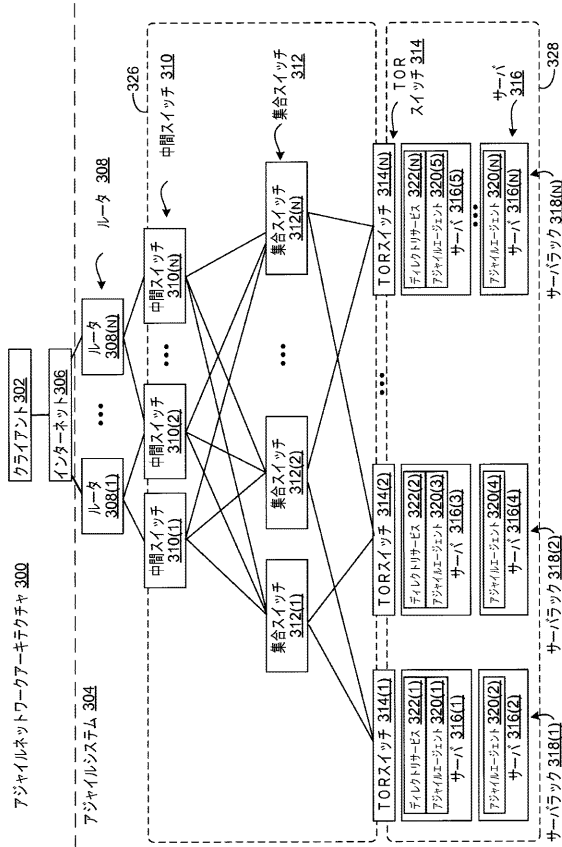
【図1】



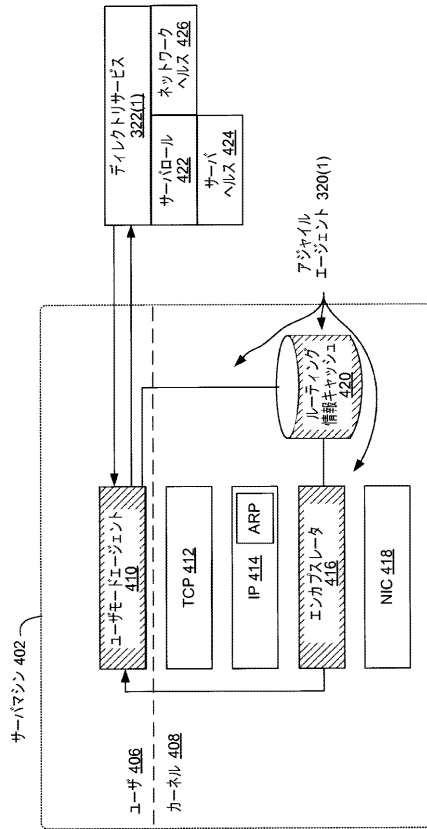
【図2】



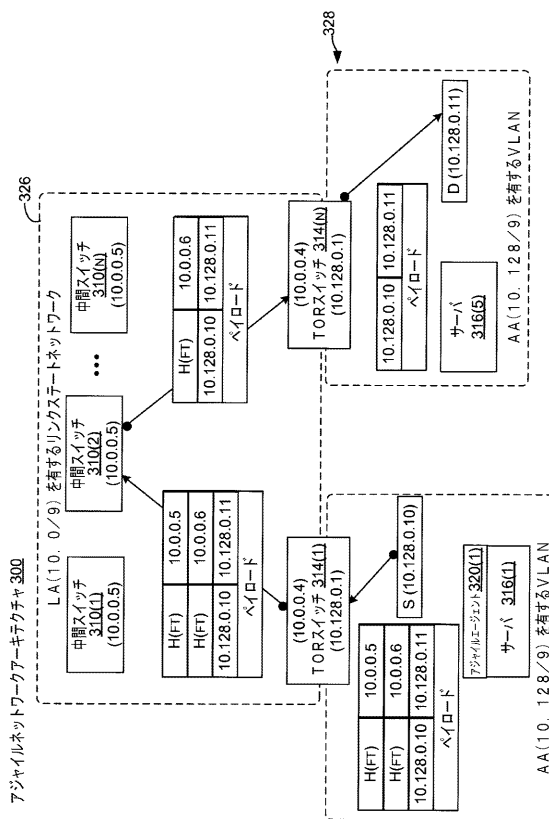
【 図 3 】



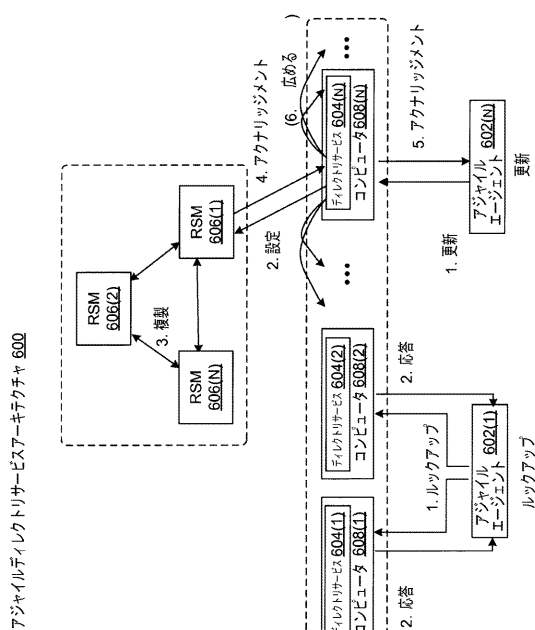
【 図 4 】



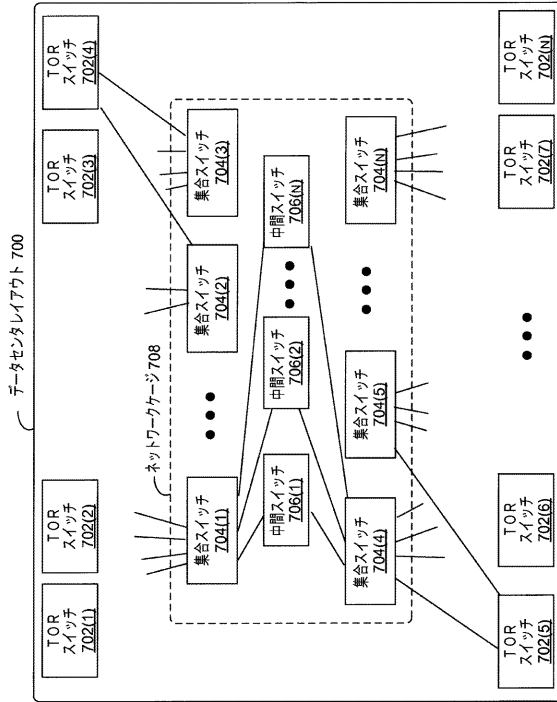
【 図 5 】



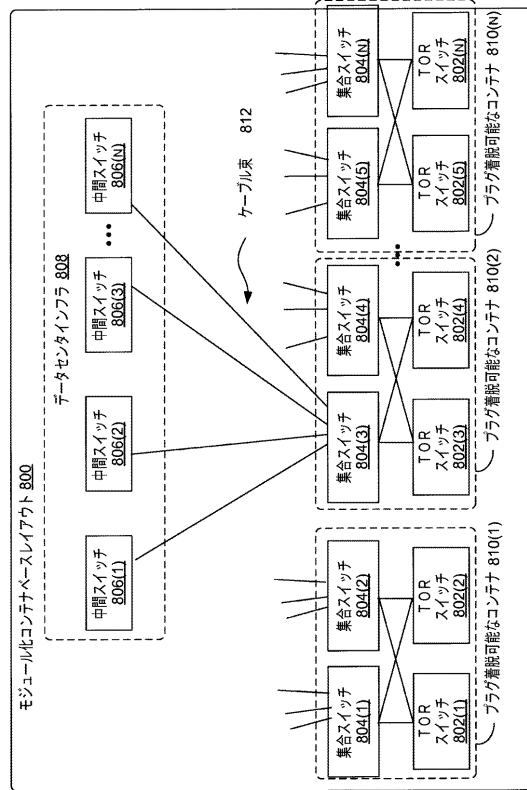
【 図 6 】



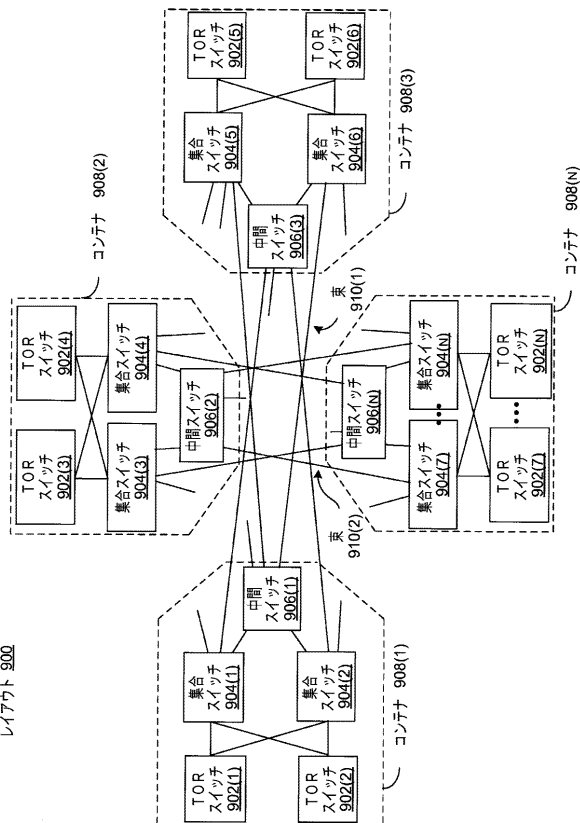
【図 7】



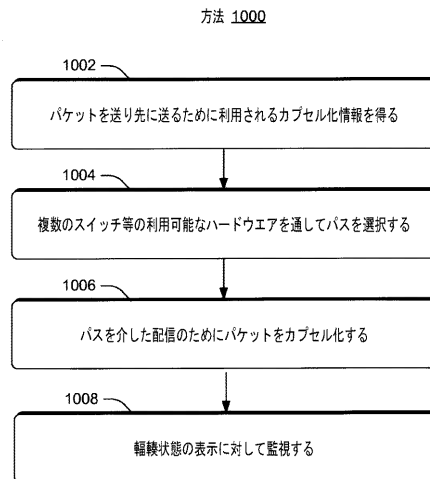
【図 8】





【図 9】



【図 10】



## 【 国際調査報告 】

<b>INTERNATIONAL SEARCH REPORT</b>		International application No. <b>PCT/US2010/036758</b>
<b>A. CLASSIFICATION OF SUBJECT MATTER</b>		
<i>H04L 29/06(2006.01)i, H04L 12/56(2006.01)i, H04L 12/28(2006.01)i</i>		
According to International Patent Classification (IPC) or to both national classification and IPC		
<b>B. FIELDS SEARCHED</b>		
Minimum documentation searched (classification system followed by classification symbols) H04L 29/06; G06F 15/16		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Korean utility models and applications for utility models Japanese utility models and applications for utility models		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) eKOMPASS(KIPO internal) & Keywords:virtual routing, encapsulate, layer-3, layer2		
<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	GIHWAN CHO, "An Efficient Location and Routing Scheme for Mobile Computing Environments", IEEE JSAC. JUNE 1995, VOL. 13. NO. 5. See page 868, left column, lines 11-20; page 871, right column, lines 1-24.	9-10
Y	JESSE M. GORDON, "Hypercube Message Routing in the Presence of Faults", ACM, 1988 See abstract; page 319, lines 20-30.	9-10
A	US 2009-0063706 A1 (GOLDMAN JOEL et al.) 05 March 2009 See whole documents.	1-8
A	BRAIN VETTER, "An Experimental Study of Insider Attacks for OSPF Routing Protocol", IEEE International Conference on Network Protocols, 1997. See whole documents.	9-12
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> See patent family annex.		
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search 29 DECEMBER 2010 (29.12.2010)		Date of mailing of the international search report <b>03 JANUARY 2011 (03.01.2011)</b>
Name and mailing address of the ISA/KR  Korean Intellectual Property Office Government Complex-Daejeon, 139 Seonsa-ro, Seo-gu, Daejeon 302-701, Republic of Korea Facsimile No. 82-42-472-7140		Authorized officer Lee Seoung Young Telephone No. 82-42-481-8591 

**INTERNATIONAL SEARCH REPORT**

Information on patent family members

International application No.

**PCT/US2010/036758**

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
US 2009-0063706 A1	05.03.2009	None	

## フロントページの続き

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

(特許庁注：以下のものは登録商標)

1. L i n u x
2. イーサネット
3. F r e e B S D

(72)発明者 パラントップ ラヒリ

アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ  
マイクロソフト コーポレーション エルシーエー - インターナショナル パテンツ内

(72)発明者 デイビッド エー・マルツ

アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ  
マイクロソフト コーポレーション エルシーエー - インターナショナル パテンツ内

(72)発明者 バルヴィーン ケー・パテル

アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ  
マイクロソフト コーポレーション エルシーエー - インターナショナル パテンツ内

(72)発明者 スディプタ セングプタ

アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ  
マイクロソフト コーポレーション エルシーエー - インターナショナル パテンツ内

(72)発明者 ナヴェンデュ ジェイン

アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ  
マイクロソフト コーポレーション エルシーエー - インターナショナル パテンツ内

(72)発明者 チャンホン キム

アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ  
マイクロソフト コーポレーション エルシーエー - インターナショナル パテンツ内

Fターム(参考) 5B089 GA11 HB19 KB03 KC23 MA03

5K030 GA19 HA08 HC13 HD03 HD06 HD09 KA01 KA07