

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号
特許第6752815号
(P6752815)

(45) 発行日 令和2年9月9日 (2020. 9. 9)

(24) 登録日 令和2年8月21日 (2020. 8. 21)

(51) Int. Cl.

F I

HO 4 L 12/715 (2013. 01)

GO 6 F 9/50 (2006. 01)

HO 4 L 12/70 (2013. 01)

HO 4 L 12/715

GO 6 F 9/50 1 5 O A

GO 6 F 9/50 1 5 O C

HO 4 L 12/70 D

請求項の数 15 (全 34 頁)

(21) 出願番号	特願2017-554456 (P2017-554456)	(73) 特許権者	502303739
(86) (22) 出願日	平成28年9月9日 (2016. 9. 9)		オラクル・インターナショナル・コーポレ イション
(65) 公表番号	特表2018-537004 (P2018-537004A)		アメリカ合衆国カリフォルニア州9406 5レッドウッド・シティ、オラクル・パ ークウェイ500
(43) 公表日	平成30年12月13日 (2018. 12. 13)		
(86) 国際出願番号	PCT/US2016/050992	(74) 代理人	110001195
(87) 国際公開番号	W02017/065905		特許業務法人深見特許事務所
(87) 国際公開日	平成29年4月20日 (2017. 4. 20)	(72) 発明者	ザヒド、フェロツ
審査請求日	令和1年6月24日 (2019. 6. 24)		ノルウェー、エヌー1325 リュサケー ル、ピィ・オウ・ボックス・134
(31) 優先権主張番号	62/240, 986	(72) 発明者	グラン、アーンスト・ガンナー
(32) 優先日	平成27年10月13日 (2015. 10. 13)		ノルウェー、エヌー1325 リュサケー ル、ピィ・オウ・ボックス・134
(33) 優先権主張国・地域又は機関	米国 (US)		
(31) 優先権主張番号	62/242, 211		
(32) 優先日	平成27年10月15日 (2015. 10. 15)		
(33) 優先権主張国・地域又は機関	米国 (US)		

最終頁に続く

(54) 【発明の名称】 マルチテナントクラスタ環境における効率的なネットワーク分離および負荷バランシングのためのシステムおよび方法

(57) 【特許請求の範囲】

【請求項 1】

マルチテナントクラスタ環境においてネットワーク分離をサポートするための方法であって、

前記マルチテナントクラスタ環境内の1つ以上のテナントをサポートするステップと、
前記1つ以上のテナントの各々を複数のパーティションのうちのあるパーティションに関連付けするステップと、

前記複数のパーティションの各々を複数のノードのうちの1つ以上のノードに関連付けるステップとを備え、前記複数のノードの各々は、複数のスイッチのうちのリーフスイッチに関連付けられ、前記複数のスイッチは、複数のリーフスイッチと、別のレベルにおける少なくとも1つのスイッチとを備え、前記方法はさらに、

前記複数のパーティションの各々を複数のポリシーパラメータのうちのあるポリシーパラメータでマーキングするステップと、

前記複数のノードの各ノードにパーティション化順序を割り当てるステップとを備え、前記パーティション化順序は、少なくとも各ノードに関連付けられた前記パーティションにマーキングされた前記ポリシーパラメータに基づき、前記方法はさらに、

少なくとも前記複数のパーティションのうちの前記パーティションの前記マーキングに基づいて、前記マルチテナントクラスタ環境で使用される1つ以上のリニアフォーディングテーブルを生成するステップを備える、方法。

【請求項 2】

前記複数のパーティションの各々をグローバルなポリシーパラメータでマーキングするステップをさらに備える、請求項 1 に記載の方法。

【請求項 3】

前記複数のポリシーパラメータは、
厳格なパラメータと、
厳格な仮想レーンパラメータと、
ベストエフォートパラメータとを備える、請求項 1 または 2 に記載の方法。

【請求項 4】

厳格なパラメータでマーキングされたパーティションに関連付けられるノードのパーティション化順序は、厳格な仮想レーンパラメータでマーキングされたパーティションに関連付けられるノードのパーティション化順序より前であり、厳格な仮想レーンパラメータでマーキングされたパーティションに関連付けられるノードのパーティション化順序は、ベストエフォートパラメータでマーキングされたパーティションに関連付けられるノードのパーティション化順序より前である、請求項 3 に記載の方法。

【請求項 5】

1 つ以上のリニアフォワーディングテーブルを生成するステップは、
前記複数のリーフスイッチの各々について、各ノードの前記パーティション化順序に従って前記複数のノードを順序付けるステップと、
前記ノードの順序で複数のエンドノードをルーティングするステップとを備え、前記ルーティングするステップは、

各ノードについて少なくとも 1 つの下りポートおよび少なくとも 1 つの上りポートを選択するステップを備える、請求項 4 に記載の方法。

【請求項 6】

前記ノードの順序で前記複数のエンドノードをルーティングするステップは、厳格なパラメータでマーキングされたパーティションに関連付けられるそれらのノードに対して厳格な分離を提供する、請求項 5 に記載の方法。

【請求項 7】

前記ノードの順序で前記複数のエンドノードをルーティングするステップは、厳格な仮想レーンパラメータでマーキングされたパーティションに関連付けられるそれらのノードに対して仮想レーン上で分離を提供する、請求項 5 または 6 に記載の方法。

【請求項 8】

マルチテナントクラスタ環境においてネットワーク分離をサポートするシステムであって、

1 つ以上のプロセッサと、

命令を格納したメモリとを備え、前記命令は、前記 1 つ以上のプロセッサによって実行されたときに前記 1 つ以上のプロセッサにステップを実行させ、前記ステップは、

前記マルチテナントクラスタ環境内の 1 つ以上のテナントをサポートするステップと

、
前記 1 つ以上のテナントの各々を複数のパーティションのうちのあるパーティションに関連付けるステップと、

前記複数のパーティションの各々を複数のノードのうちの一つ以上のノードに関連付けるステップとを備え、前記複数のノードの各々は、複数のスイッチのうちの一つ以上のスイッチに関連付けられ、前記複数のスイッチは、複数のリーフスイッチと、別のレベルにおける少なくとも一つのスイッチとを備え、前記ステップはさらに、

前記複数のパーティションの各々を複数のポリシーパラメータのうちのあるポリシーパラメータでマーキングするステップと、

前記複数のノードの各ノードにパーティション化順序を割り当てるステップとを備え、前記パーティション化順序は、少なくとも各ノードに関連付けられた前記パーティションにマーキングされた前記ポリシーパラメータに基づき、前記ステップはさらに、

少なくとも前記複数のパーティションのうちの前記パーティションの前記マーキング

10

20

30

40

50

に基づいて、前記マルチテナントクラスタ環境で使用される１つ以上のリニアフォワーディングテーブルを生成するステップを備える、システム。

【請求項 9】

前記１つ以上のプロセッサは、前記複数のパーティションの各々をグローバルなポリシーパラメータでマーキングするステップを備えるさらに他のステップを実行するように動作する、請求項 8 に記載のシステム。

【請求項 10】

前記複数のポリシーパラメータは、
厳格なパラメータと、
厳格な仮想レーンパラメータと、

10

ベストエフォートパラメータとを備える、請求項 8 または 9 に記載のシステム。

【請求項 11】

厳格なパラメータでマーキングされたパーティションに関連付けられるノードのパーティション化順序は、厳格な仮想レーンパラメータでマーキングされたパーティションに関連付けられるノードのパーティション化順序より前であり、厳格な仮想レーンパラメータでマーキングされたパーティションに関連付けられるノードのパーティション化順序は、ベストエフォートパラメータでマーキングされたパーティションに関連付けられるノードのパーティション化順序より前である、請求項 10 に記載のシステム。

【請求項 12】

１つ以上のリニアフォワーディングテーブルを生成するステップは、

20

前記複数のリーフスイッチの各々について、各ノードの前記パーティション化順序に従って前記複数のノードを順序付けるステップと、

前記ノードの順序で複数のエンドノードをルーティングするステップとを備え、前記ルーティングするステップは、

各ノードについて少なくとも１つの下りポートおよび少なくとも１つの上りポートを選択するステップを備える、請求項 11 に記載のシステム。

【請求項 13】

前記ノードの順序で前記複数のエンドノードをルーティングするステップは、厳格なパラメータでマーキングされたパーティションに関連付けられるそれらのノードに対して厳格な分離を提供する、請求項 12 に記載のシステム。

30

【請求項 14】

前記ノードの順序で前記複数のエンドノードをルーティングするステップは、厳格な仮想レーンパラメータでマーキングされたパーティションに関連付けられるそれらのノードに対して仮想レーン上で分離を提供する、請求項 12 または 13 に記載のシステム。

【請求項 15】

コンピュータプログラムであって、前記コンピュータプログラムは、マルチテナントクラスタ環境においてネットワーク分離をサポートするためのプログラム命令を備え、前記プログラム命令は、実行されたときにコンピュータに請求項 1 ～ 7 のいずれか 1 項に記載の方法を実行させる、コンピュータプログラム。

【発明の詳細な説明】

40

【技術分野】

【0001】

著作権表示：

この特許文献の開示の一部は、著作権保護の対象となる題材を含んでいる。著作権所有者は、特許商標庁の包袋または記録に掲載されるように特許文献または特許情報開示を誰でも複製できることに対して異議はないが、その他の点ではすべてのいかなる著作権をも保有する。

【0002】

発明の分野：

本発明は、一般にコンピュータシステムに関し、特にマルチテナントクラスタ環境に関

50

する。

【背景技術】

【0003】

背景：

マルチテナンシは、利用可能なシステムリソースの利用度を高くすることを約束し、サービスプロバイダのために費用効果の高い動作を維持することに役立つ。しかし、マルチテナント高性能コンピューティング (high-performance computing : H P C) インフラストラクチャは、固有の課題をもたらし、当該固有の課題は、性能分離をテナントに提供することにもネットワークファブリック全体にわたる効率的な負荷バランシングを達成することにも関連付けられる。

10

【発明の概要】

【課題を解決するための手段】

【0004】

概要：

マルチテナントクラスタ環境においてネットワーク分離をサポートするためのシステムおよび方法である。例示的な方法は、上記マルチテナントクラスタ環境内の1つ以上のテナントをサポートし得る。上記方法は、上記1つ以上のテナントの各々を複数のパーティションのうちのあるパーティションに関連付け得て、また、上記複数のパーティションの各々を複数のノードのうちのある1つ以上のノードに関連付け得て、上記複数のノードの各々は、複数のスイッチのうちのあるリーフスイッチに関連付けられ、上記複数のスイッチは、複数のリーフスイッチと、別のレベルにおける少なくとも1つのスイッチとを備える。上記方法は、上記複数のパーティションの各々を複数のポリシーパラメータのうちのあるポリシーパラメータでマーキングし得る。上記方法は、上記複数のノードの各ノードにパーティション化順序を割り当て得て、上記パーティション化順序は、少なくとも各ノードに関連付けられた上記パーティションにマーキングされたポリシーパラメータに基づく。最後に、上記方法は、少なくとも上記複数のパーティションのうちのある上記パーティションの上記マーキングに基づいて、上記マルチテナントクラスタ環境で使用される1つ以上のリニアフォワーディングテーブルを生成し得る。

20

【図面の簡単な説明】

【0005】

30

【図1】実施形態に係るマルチテナントクラスタ環境の図である。

【図2】本開示の実施形態を実施することができるネットワーク環境におけるツリートポロジの図である。

【図3】本開示の実施形態を実施することができるマルチテナントクラスタ環境におけるルーティングの図である。

【図4】本開示の実施形態を実施することができるマルチテナントクラスタ環境におけるルーティングの図である。

【図5】実施形態に係るマルチテナントクラスタ環境においてパーティション認識ルーティングをサポートすることの図である。

【図6】実施形態に係るマルチテナントクラスタ環境においてパーティション認識ルーティングをサポートすることの図である。

40

【図7】実施形態に係るマルチテナントクラスタ環境においてパーティション認識ルーティングをサポートすることの図である。

【図8】実施形態に係るマルチテナントクラスタ環境においてパーティション認識ルーティングをサポートすることの図である。

【図9】実施形態に係るマルチテナントクラスタ環境においてネットワーク分離をサポートすることの図である。

【図10】実施形態に係るマルチテナントクラスタ環境においてネットワーク分離をサポートすることの図である。

【図11】実施形態に係るマルチテナントクラスタ環境において重み付けされたパーティ

50

ション認識ルーティングをサポートすることの図である。

【図 1 2】実施形態に係るマルチテナントクラスタ環境において重み付けされたパーティション認識ルーティングをサポートすることの図である。

【図 1 3】実施形態に係るマルチテナントクラスタ環境においてネットワーク分離をサポートするための方法のフローチャートである。

【図 1 4】本発明のさまざまな実施形態に係るコンピュータシステムを示す一般化された概略図である。

【発明を実施するための形態】

【 0 0 0 6 】

詳細な説明：

以下の詳細な説明では、本発明は、限定としてではなく一例として添付の図面の図に示される。なお、本開示における「ある」または「1つの」または「いくつかの」実施形態への言及は、必ずしも同一の実施形態に言及するものではなく、このような言及は少なくとも1つを意味する。具体的な実現例を記載しているが、当該具体的な実現例は例示の目的で提供されているに過ぎないということが理解される。本発明の範囲および精神から逸脱することなく他の構成要素および構成を使用してもよいということを当業者は認識するであろう。

【 0 0 0 7 】

図面および詳細な説明全体にわたって、共通の参照番号を使用して同様の要素を示す。したがって、要素が他の場所に記載されている場合には、図で使用される参照番号は、当該図を対象とした詳細な説明の中で言及されてもされなくてもよい。

【 0 0 0 8 】

本発明の以下の説明では、高性能ネットワークの一例としてインフィニバンド（商標）（Infiniband：IB）ネットワークを使用する。他のタイプの高性能ネットワークが制限なく使用されてもよいということが当業者に明らかであろう。また、以下の説明では、ファブリックトポロジの一例としてファットツリートポロジを使用する。他のタイプのファブリックトポロジが制限なく使用されてもよいということが当業者に明らかであろう。

【 0 0 0 9 】

インフィニバンド（商標）

インフィニバンド（商標）（IB）は、インフィニバンド（商標）トレードアソシエーションによって開発された開かれた標準的な無損失ネットワーク技術である。当該技術は、特にHPCアプリケーションおよびデータセンタを対象とした、高スループットおよび低遅延通信を提供するシリアルポイントツーポイント全二重相互接続に基づく。

【 0 0 1 0 】

インフィニバンド（商標）アーキテクチャ（InfiniBand Architecture：IBA）は、二層トポロジ分割をサポートする。下位層では、IBネットワークはサブネットと称され、サブネットは、スイッチおよびポイントツーポイントリンクを使用して相互接続された一組のホストを含み得る。上位レベルでは、IBファブリックは、ルータを使用して相互接続され得る1つ以上のサブネットを構成する。

【 0 0 1 1 】

サブネット内において、ホスト同士は、スイッチおよびポイントツーポイントリンクを使用して接続される。また、1つのマスタ管理エンティティ、すなわちサブネットマネージャ（subnet manager：SM）があり、当該サブネットマネージャは、サブネット内の指定されたサブネットデバイス上に存在している。サブネットマネージャは、IBサブネットを構成、起動および維持することを担う。また、サブネットマネージャ（SM）は、IBファブリックにおいてルーティングテーブル計算を実行することを担ってもよい。ここで、たとえば、IBネットワークのルーティングは、ローカルサブネット内の全てのソース・宛先ペア間で負荷バランシングが適切になることを目指している。

【 0 0 1 2 】

サブネット管理インターフェイスを介して、サブネットマネージャは、サブネット管理

10

20

30

40

50

パケット (subnet management packet : S M P) と称される制御パケットをサブネット管理エージェント (subnet management agent : S M A) とやりとりする。サブネット管理エージェントは、どの I B サブネットデバイス上にも存在している。S M P を使用して、サブネットマネージャは、ファブリックを発見し、エンドノードおよびスイッチを構成し、S M A から通知を受け取ることができる。

【 0 0 1 3 】

一般に、マスタサブネットマネージャ以外の全ての他のサブネットマネージャは、フォールトトレランスのために待機モードで動作する。しかし、マスタサブネットマネージャが機能しなくなる状況では、待機中のサブネットマネージャによって新たなマスタサブネットマネージャが決められる。また、マスタサブネットマネージャは、サブネットの定期的なスweepを実行して、いかなるトポロジ変化も検出し、それに応じてネットワークを再構成する。

10

【 0 0 1 4 】

さらに、サブネット内のホストおよびスイッチは、ローカル識別子 (local identifier : L I D) を使用してアドレス指定可能であり、単一のサブネットは、4 9 1 5 1 個の L I D に制限される。サブネット内で有効なローカルアドレスである L I D に加えて、各 I B デバイスは、その不揮発性メモリに焼き付けられた 6 4 ビットのグローバル一意識別子 (global unique identifier : G U I D) を有し得る。G U I D を使用して、I B 層 3 (L 3) アドレスであるグローバル識別子 (global identifier : G I D) を形成することができる。G I D は、6 4 ビットのサブネット識別子 (I D) を 6 4 ビットの G U I D と連結して I P v 6 のような 1 2 8 ビットのアドレスを形成することによって作成することができる。たとえば、I B ファブリックに接続されたポートにさまざまなポート G U I D が割り当てられてもよい。

20

【 0 0 1 5 】

S M は、ネットワーク初期化時にルーティングテーブル (すなわち、サブネット内のノードの各ペア間の接続 / 経路) を計算し得る。さらに、接続性および最適性能を保証するために、トポロジが変化するたびにルーティングテーブルが更新され得る。通常動作中、S M は、ネットワークの定期的な軽いスweepを実行して、トポロジ変化をチェックし得る。軽いスweep中に変化が発見されるか、またはネットワーク変化を知らせるメッセージ (トラップ) が S M によって受信されると、S M は、発見された変化に従ってネットワークを再構成し得る。

30

【 0 0 1 6 】

たとえば、S M は、リンクがダウンしたとき、デバイスが追加されたとき、またはリンクが除去されたときなど、ネットワークトポロジが変化したときにネットワークを再構成してもよい。再構成ステップは、ネットワーク初期化中に実行されるステップを含み得る。さらに、再構成は、ネットワーク変化が生じたサブネットに限定されるローカルスコープを有し得る。また、ルータを用いた大きなファブリックのセグメント化は、再構成スコープを制限する可能性がある。

【 0 0 1 7 】

実施形態によれば、I B ネットワークは、ルータを使用して相互接続された 1 つ以上のサブネットで構成され得る。サブネット内において、ホスト同士は、スイッチおよびポイントツーポイントリンクを使用して接続される。各 I B サブネット内に、1 つのマスタ管理エンティティ、すなわちサブネットマネージャ (S M) があってもよく、当該サブネットマネージャは、任意の指定されたサブネットデバイス上に存在し、I B サブネットを構成、起動および維持する。

40

【 0 0 1 8 】

サブネット管理インターフェイスを介して、S M は、サブネット管理パケット (S M P) と呼ばれる制御パケットを、どの I B デバイス上にも存在しているサブネット管理エージェント (S M A) とやりとりする。S M P を使用して、S M は、ファブリックを発見し、エンドノードおよびスイッチを構成し、S M A から通知を受け取ることができる。また

50

、S Mは、サブネットの定期的な軽いスweepを実行して、いかなるトポロジ変化も検出し、それに応じてネットワークを再構成することができる。

【0019】

実施形態によれば、I Bネットワークにおけるサブネット内ルーティングは、スイッチに記憶されたりニアフォーワードイングテーブル (linear forwarding table: L F T) に基づき得る。L F Tは、使用時にルーティングメカニズムに従ってS Mによって計算される。サブネットにおいて、エンドノード上の全てのH C Aポートおよび全てのスイッチは、ローカル識別子 (L I D) を使用してアドレス指定される。L F Tへの各入力は、宛先L I D (D L I D) および出力ポートで構成される。テーブルへのL I D当たり1つの入力のみがサポートされる。パケットがスイッチに到達すると、当該スイッチのフォーワードイングテーブルにおいてD L I Dを検索することによってその出力ポートが決定される。パケットは所与のソース・宛先ペア (L I D ペア) 間でネットワーク内の同一のパスを通るので、ルーティングは決定論的である。

10

【0020】

実施形態によれば、パーティション化は、ネットワークファブリックを共有しているシステムの論理グループの分離を提供するためにI Bによってサポートされるセキュリティメカニズムである。ファブリック内のノード上の各H C Aポートは、1つ以上のパーティションのメンバであり得る。パーティションメンバーシップは、S Mの一部であり得る集中型パーティションマネージャによって管理される。S Mは、各ポートについてのパーティションメンバーシップ情報を16ビットのパーティションキー (P _ K e y) のテーブルとして構成し得る。また、S Mは、L I Dに関連付けられるP _ K e y 値を含むパーティション実行 (partition enforcement) テーブルを用いてスイッチおよびルータを構成し得る。

20

【0021】

実施形態によれば、ノード間での通信のために、管理キューペア (Q P 0 および Q P 1) 以外に、キューペア (Q P) およびエンドツーエンドコンテキスト (E E C) が特定のパーティションに割り当てられ得る。次いで、送信されたどのI BトランスポートパケットにもP _ K e y 情報が追加され得る。パケットがH C Aポートまたはスイッチに到達すると、そのP _ K e y 値は、S Mによって構成されたテーブルに対して有効化され得る。無効なP _ K e y 値が見つければ、パケットはすぐに廃棄される。このようにして、パーティションを共有しているポート間でのみ通信が可能になる。

30

【0022】

I Bパーティションの一例が図1に示されており、図1は、実施形態に係るマルチテナントクラスタ環境の図である。図1に示される例では、ノードA ~ E 101 ~ 105は、インフィニバンドファブリック100を使用して、それぞれのホストチャネルアダプタ111 ~ 115を介して通信する。ノードA ~ Eは、パーティション、すなわちパーティション1 110、パーティション2 120およびパーティション3 130の中に配置されている。パーティション1は、ノードA 101とノードD 104とを備える。パーティション2は、ノードA 101とノードB 102とノードC 103とを備える。パーティション3は、ノードC 103とノードE 105とを備える。このようなパーティションの配置のために、ノードD 104およびノードE 105は、パーティションを共有していないので通信することができない。一方、たとえば、ノードA 101およびノードC 103は、両方ともがパーティション2 120の一部であるので通信することができる。

40

【0023】

実施形態によれば、パーティションは、ネットワークファブリックを共有しているシステムの論理グループの分離を実施するためのセキュリティメカニズムとして提供されることができる。I Bパーティションは、同様の分離特徴をイーサネット (登録商標) 802.1Q V L A Nとして提供することができる。ファブリック内のノード上の各H C Aポートは、1つ以上のパーティションのメンバであり得る。パーティションメンバーシップ

50

は、S Mの一部であり得る集中型パーティションマネージャによって管理することができる。S Mは、各ポートについてのパーティションメンバーシップ情報を16ビットのパーティションキー(P _ K e y)のテーブルとして構成し得る。また、S Mは、P _ K e y値を含むパーティション実行テーブルを用いてスイッチおよびルータを構成し得る。

【0024】

実施形態によれば、ノード間での通信のために、管理キューペア(Q P 0およびQ P 1)以外に、キューペア(Q P)およびエンドツーエンドコンテキスト(E E C)が特定のパーティションに割り当てられ得る。次いで、送信されたどのトランスポートパケットにもP _ K e y情報が追加され得る。パケットがH C Aポートまたはスイッチに到達すると、そのP _ K e y値は、S Mによって構成されたテーブルに対して有効化され得る。無効なP _ K e y値が見つければ、パケットは廃棄される。このようにして、パーティションを共有しているポート間でのみ通信が可能になる。

10

【0025】

実施形態によれば、I Bは、V Lを使用して各物理リンクを複数の仮想チャネルに分割することができる層状アーキテクチャである。各V Lは、それ自体のバッファリング、フロー制御および輻輳管理リソースを有し得る。一組の区別されたトラフィッククラス、すなわちS L、を介してQ o Sを提供することができる。S Lは、ネットワーク内でパケットが受信することができるサービスのクラスを表わす。各S Lは、構成されたS L - V Lマッピングテーブルに基づいてリンク上のV Lにマッピングされる。I Bは、16個までのV Lをサポートする。しかし、最後のV Lは、サブネット管理トラフィックのために取っておかれて、通常はユーザアプリケーションによって使用されることはない。

20

【0026】

I Bシステムにおけるマルチテナンシ

ネットワークングの観点から、マルチテナンシは、ネットワークリソースの利用度を高くして、サービスプロバイダのために費用効果の高い動作を維持することに役立つことができる。しかし、マルチテナントインフラストラクチャは、いくつかの重要なセキュリティ問題も課し、最も困難なものの中の1つは、テナントに対して性能分離を提供することに関連する。各テナントは、システム内の他のテナントのワークロードの影響を受けない予測可能なネットワーク性能を備えているべきである。I Bシステムにおけるネットワーク分離は、パーティション化によって提供することができる。

30

【0027】

実施形態によれば、パーティションは、グループのメンバが同一グループの他のメンバとのみ通信できるようなポートの論理グループである。ホストチャネルアダプタ(H C A)およびスイッチにおいて、パーティションメンバーシップ情報を使用してパケットをフィルタリングして分離を実施することができる。無効なパーティション化情報を有するパケットは、受信ポートに到達するとすぐに削除され得る。しかし、H P Cシステムで使用するルーティングアルゴリズムは、一般に、ネットワーク内のこのようなパーティションを認識しない。したがって、異なるパーティションに属するトラフィックフローがネットワークファブリック内のリンクを共有する可能性がある。

【0028】

マルチテナントI Bシステムでは、パーティションを使用してテナントクラスタを作成することができる。パーティション実行を実施すると、ノードは、異なるテナントクラスタに属する他のノードとは通信できなくなる。このようにして、感染したまたは不正なテナントノードが存在していてもシステムのセキュリティを保証することができる。

40

【0029】

一般に、I Bルーティングは、スイッチに記憶されたリニアフォワーディングテーブル(L F T)に基づき得る。L F Tは、パーティション化情報を考慮に入れることなくサブネットマネージャ(S M)によって計算される。したがって、中間ネットワークリンクは、さまざまなパーティションに属するトラフィックを搬送し得る。このように中間リンクを共有することは、パーティション干渉につながるおそれがある。その結果、テナントは

50

予測不可能なネットワーク性能を経験することになる。さらに、パーティション化されたサブネットでは、ルーティングアルゴリズムのバランシング特徴も影響を受ける。その理由は、たとえばパーティション境界を横断するリンクがユーザトラフィックに利用されなくても、これらのリンクが他の機能的リンクと同じようにルーティングされる（したがって、バランシングの際に考慮される）からである。バランシングの質が低下することにより、有効帯域幅および準最適ネットワーク利用度が低下する可能性がある。

【0030】

I Bは、一般に、他のパーティションにおけるノードにかかわらず、あるシェアの利用可能な帯域幅を各パーティションに保証するために使用できるサービス品質（QoS）特徴を提供する。次いで、各パーティションにサービスレベル（service level：SL）と呼ばれる利用可能な区別されたトラフィッククラスを割り当てることによって帯域幅保証が提供される。次いで、各SLは、SL-VLマッピングテーブルに従ってリンク上の利用可能な15個の仮想レーン（VL）のうちの1つにマッピングされる。

10

【0031】

SLをパーティションに割り当てる際に問題が生じる可能性がある。なぜなら、当該システムは15個のVLのみを利用してネットワーク内に個別のパーティションを作成することができる一方、I Bネットワークは一般に多数のパーティションを有し得る（たとえば、各ポートは32,768個までのパーティションのメンバであり得る）からである。さらに、既存のI Bハードウェアでは、（サブネット管理のために取っておかれるものを含む）9個のVLのみをサポートすることが一般的である。さらに、SLは希少なリソースであるので、他の目的、たとえばネットワーク内でフォールトトレランスまたはサービス差別化を提供する目的のためにできるだけ多くのSLを空けておくことが望ましいであろう。

20

【0032】

ファットツリー（F T r e e）トポロジおよびルーティング

実施形態によれば、I BベースのHPCシステムのうちのいくつかは、ファットツリートポロジを利用して、有用な特性であるファットツリーの提供を活用する。これらの特性は、各ソース・宛先ペア間で複数のパスを利用できることに起因した全二分割帯域幅および固有のフォールトトレランスを含む。ファットツリーの背後にある最初の考え方は、ツリーがトポロジのルートの方に移動するにつれて、利用可能な帯域幅が増加していき、ノード間のリンクが太くなっていくことを利用するというものであった。リンクが太くなることは、上位のスイッチでの輻輳の回避に役立つことができ、二分割帯域幅が維持される。

30

【0033】

図2は、本開示の実施形態を実施することができるネットワーク環境におけるツリートポロジの図である。図2に示されるように、ネットワークファブリック200において1つ以上のエンドノード201~204を接続することができる。ネットワークファブリック200は、複数のリーフスイッチ211~214と複数のスパイン（spine）スイッチまたはルートスイッチ231~234とを含むファットツリートポロジに基づき得る。また、ネットワークファブリック200は、スイッチ221~224などの1つ以上の中間スイッチを含み得る。

40

【0034】

また、図2に示されるように、エンドノード201~204の各々は、マルチホームのノード、すなわち複数のポートを介してネットワークファブリック200の2つ以上の部分に接続される単一のノードであってもよい。たとえば、ノード201はポートH1およびH2を含み得て、ノード202はポートH3およびH4を含み得て、ノード203はポートH5およびH6を含み得て、ノード204はポートH7およびH8を含み得る。

【0035】

また、各スイッチは、複数のスイッチポートを有し得る。たとえば、ルートスイッチ231はスイッチポート1~2を有し得て、ルートスイッチ232はスイッチポート3~4

50

を有し得て、ルートスイッチ 2 3 3 はスイッチポート 5 ~ 6 を有し得て、ルートスイッチ 2 3 4 はスイッチポート 7 ~ 8 を有し得る。

【 0 0 3 6 】

実施形態によれば、ファットツリールーティングメカニズムは、I B ベースのファットツリートポロジのための最も人気のあるルーティングアルゴリズムのうちの 1 つである。ファットツリールーティングメカニズムは、O F E D (オープン・ファブリック・エンタープライズ・ディストリビューション : I B ベースのアプリケーションを構築およびデプロイするための標準的なソフトウェアスタック) サブネットマネージャ、すなわち O p e n S M、でも実現される。

【 0 0 3 7 】

ファットツリールーティングメカニズムは、ネットワークファブリック内のリンク全体にわたって最短パス経路を均一に分散させる L F T を生成することを目指している。当該メカニズムは、インデックス化順序でファブリックを横断し、エンドノード、したがって対応する経路、のターゲット L I D を各スイッチポートに割り当てる。同一のリーフスイッチに接続されたエンドノードでは、インデックス化順序は、エンドノードが接続されているスイッチポート (すなわち、ポート番号付けシーケンス) に左右され得る。各ポートについて、当該メカニズムは、ポート使用カウンタを維持することができ、このポート使用カウンタを使用して、新たな経路が追加されるたびに最も使われていないポートを選択することができる。同一の 2 個のスイッチを接続する複数のポートがある場合には、当該ポートはポートグループを形成する。その場合、新たな経路を追加するために最小負荷ポートグループの最小負荷ポートが選択される。

【 0 0 3 8 】

上記のように、パーティション化されたサブネットでは、共通のパーティションのメンバーではないノード同士は通信することができない。これは、事実上、ファットツリールーティングアルゴリズムによって割り当てられた経路のうちのいくつかがユーザトラフィックに使用されないことを意味する。ファットツリールーティングメカニズムが他の機能的パスと同じようにそれらの経路について L F T を生成する場合に問題が生じる。この挙動は、リンク上のバランシングの質の低下を生じさせる可能性がある。なぜなら、ノードがインデックス化順序でルーティングされるからである。ルーティングがパーティションに気付かずになされるので、ファットツリールーティングが行われたサブネットでは、一般にパーティション間の分離が不十分である。

【 0 0 3 9 】

図 3 は、本開示の実施形態を実施することができるマルチテナントクラスタ環境におけるルーティングの図である。より具体的に、図 3 は、負荷バランシングの質の低下および不十分な分離の問題について詳細に説明する。

【 0 0 4 0 】

図 3 は、4 個のスイッチ、すなわちルートスイッチ 3 2 5 ~ 3 2 6 およびリーフスイッチ 3 2 0 ~ 3 2 1 と、3 個の重なり合うパーティションにおける 6 個のエンドノード、すなわちノード A ~ F 3 0 1 ~ 3 0 6 とを有する 2 レベルファットツリートポロジを示す。パーティション 1 は、ノード B 3 0 2 とノード C 3 0 3 とを備える。パーティション 2 は、ノード A 3 0 1 とノード C 3 0 3 とノード D 3 0 4 とノード F 3 0 6 とを備える。最後に、パーティション 3 は、ノード D 3 0 4 とノード E 3 0 5 とを備える。

【 0 0 4 1 】

実施形態によれば、パーティション 1 および 3 は、それぞれリーフスイッチ 3 2 0 および 3 2 1 (すなわち、単一リーフスイッチパーティション) 内に完全に限定されている。このため、パーティション 1 および 3 内のノード間の通信は、トラフィックをルートスイッチ 3 2 5 または 3 2 6 に移動させることなく、それらの対応するリーフスイッチを介して行われる。このトポロジがファットツリールーティングメカニズムによってルーティングされると、リーフスイッチ 3 2 0 および 3 2 1 に接続されたノードに向かう経路がルー

トスイッチに割り当てられ、そのため、リーフスイッチ間のフローはそれらの宛先に到達することができる。負荷バランスのために、AおよびCに向かう経路がルートスイッチ325に割り当てられる（リンクpとして図示）一方、ルートスイッチ326はノードBの方にトラフィックをルーティングする（リンクqとして図示）。同様に、リーフスイッチ321では、リーフスイッチ間パーティション2内のノードDおよびFに向かうトラフィックがルートスイッチ325を介してルーティングされ（リンクrとして図示）、ノードEに向かうトラフィックがルートスイッチ326を介してルーティングされる（リンクsとして図示）。

【0042】

実施形態によれば、ルートスイッチ上でのエンドポートの選択は、ノード識別子を有する小さな円として図に示されている。（ファットツリールーティングメカニズムを使用して）パーティション化情報を考慮に入れることなくルーティングがなされるので、サブネット内のパスは適切なバランスがとられない。リンクpおよびrはオーバーサブスクライブされる（oversubscribed）一方、リーフスイッチ内のフローはリンクqまたはsを使用することはない。ノードBおよびEに向かう割り当てられた経路は、（管理トラフィックが比較的低い場合を除いて）利用されない。なぜなら、これらのノードは両方とも、パーティション化によりそれらのリーフスイッチの外部からのいかなる通信も受信できないからである。このバランス問題は、パーティションの通信がトポロジ内のレベルのうちのいくつかのレベルのみに制限される場合にはファットツリーでも生じる。

【0043】

ここで、本開示の実施形態を実施することができるマルチテナントクラウド環境におけるルーティングの図である図4を参照する。より具体的に、図4は、ファットツリー内での不十分な分離に関連する問題について詳細に説明する。

【0044】

図4は、4個のスイッチ、すなわちルートスイッチ425～426およびリーフスイッチ420～421と、8個のエンドノード、すなわちノードA～G 401～408とを有する2レベルファットツリートポロジを示す。同様に、エンドノードは2個のパーティションに分割されている。パーティション1は、ノードA 401とノードB 402とノードG 407とノードH 408とを備える。パーティション2は、ノードC 403とノードD 404とノードE 405とノードF 406とを備える。

【0045】

パーティションの各々は、2個のリーフスイッチの各々に接続された2個のノードを有する。ファットツリールーティングメカニズムは、図に示されるように、ルートスイッチ425および426に下りポートを割り当てる。ファットツリールーティングメカニズムの性質のために、各ルートスイッチは、両方のパーティションに属するノードの方にトラフィックをルーティングし、これは分離を不十分にし、パーティション化された環境では望ましくないものである。たとえば、ノードAおよびCに向かうトラフィックは、共有リンクpに沿ってルーティングされる。異なるパーティションのノード間での中間リンクの共有は、それらの間に干渉を引き起こすおそれがある。ネットワークがパーティション間の完全な分離を提供するのに十分なリソースをルートレベルで有しているにもかかわらず、ファットツリールーティングメカニズムは望ましい分離を提供しない。

【0046】

パーティション認識ファットツリー（p F T r e e）ルーティング

実施形態によれば、パーティション認識ファットツリールーティングメカニズム（本明細書では別名でp F T r e eと称される）は、I Bネットワークにおけるマルチテナンシに関連する所望の目的を達成することができる。たとえば、p F T r e eメカニズムは、ツリー内のリンク全体にわたって均一に経路を分散させることによって、バランスのとれたリアフォワーディングテーブルをファットツリートポロジに提供することができる。また、p F T r e eは、リンク上の経路をバランスのとれた状態に維持しながら、異なるパーティションに属するパス間の干渉を取り除くことができる。

【 0 0 4 7 】

実施形態によれば、p F T r e eメカニズムは、サブネットについてのパーティション化情報を使用して、パーティション内のノードが、他のパーティションでのワークロードの影響を受けない予測可能なネットワーク性能を受け取ることができることを確実にすることができる。トポロジが（負荷バランシングに関して妥協することなく）各レベルにおいてパーティション分離を提供するのに十分な利用可能なリンクを持たない状況では、p F T r e eは、V Lを割り当てて干渉の影響を軽減することができる。

【 0 0 4 8 】

実施形態によれば、p F T r e eメカニズムは、各エンドノードに関連付けられるL I Dについて全ての関連するスイッチ上でL F Tをセットアップするように再帰的に機能し得る。これは、以下の疑似コード（本明細書ではリスト1と称される）で示される。

10

【 0 0 4 9 】

【 数 1 】

```

1: for each sw ∈ leafSwitches[] do
2:   Load partitioning information
3:   Filter leaf-switch only partitions
4:   Sort end nodes in partitioning specific order
5:   for each cn ∈ endNodes[] do
6:     Get lid of cn
7:     Get partition key of the cn.hca_port
8:     Set LFT[lid] ← cn:hca port on sw
9:   ROUTEDOWNGOINGBYASCENDINGTD() on sw
10:  end for
11: end for
12: ASSIGNVIRTUALLANESTD()

```

20

【 0 0 5 0 】

実施形態によれば、ROUTEDOWNGOINGBYASCENDINGTD()は、以下の疑似コード（本明細書ではリスト2と称される）で例示される。

30

【 0 0 5 1 】

【数 2】

```

1:  Get least-loaded ports from sw.UpGroups[] as uplist[]
2:  selected port  $\leftarrow$  upList.get_port_max_guid()
3:  for each port in upList[] do
4:    r_sw  $\leftarrow$  port:get_remote_switch()
5:    if r_sw is marked with partition_key then
6:      selected_port  $\leftarrow$  port
7:      break
8:    end if
9:  end for
10: if r_sw is not marked then
11:   Mark it with partition_key in DWN direction
12: end if
13: Set LFT[lid]  $\leftarrow$  selected_port on r_sw
14: ROUTEUPGOINGBYDESCENDINGTD() on sw
15: ROUTEDOWNGOINGBYASCENDINGTD() on r_sw

```

【 0 0 5 2 】

実施形態によれば、ROUTEUPGOINGBYDESCENDINGTD()は、以下の疑似コード（本明細書ではリスト3と称される）で例示される。

【 0 0 5 3 】

【数 3】

```

1:  Get least-loaded ports from sw.DownGroups[] as dwnlist[]
2:  selected port  $\leftarrow$  dwnList.get_port_max_guid()
3:  for each port in dwnList[] do
4:    r_sw  $\leftarrow$  port:get_remote_switch()
5:    if r_sw is marked with partition_key then
6:      selected_port  $\leftarrow$  port
7:      break
8:    end if
9:  end for
10: if r_sw is not marked then
11:   Mark it with partition_key in UP direction
12: end if
13: Set LFT[lid]  $\leftarrow$  selected_port on r_sw
14: ROUTEUPGOINGBYDESCENDINGTD() on r_sw

```

【 0 0 5 4 】

実施形態によれば、ASSIGNVIRTUALLANESTD()は、以下の疑似コード（本明細書では別名でリスト4と称される）で例示される。

【 0 0 5 5 】

【 数 4 】

```

1:  vlanes_needed ← 1
2:  max_vlanes ← get_max_lanes()
3:  strict ← get_is_strict()
4:  for each partition in partition_tbl do
5:    check if any intermediate communication link in this
      partition share a switch with a partition that has not
      been assigned a virtual lane
6:    if require a separate vl then
7:      if vlanes_needed = max_vlanes and strict = false then
8:        vlanes_needed ← 1
9:      else
10:        error: routing failed
11:        return
12:      end if
13:      vlanes_needed++
14:      partition.vlane ← vlanes_needed
15:    end if
16:  end for

```

【 0 0 5 6 】

実施形態によれば、単一リーフスイッチパーティション（すなわち、もっぱら単一のリーフスイッチ内で通信することができるパーティション）をフィルタリングして除去した後、当該メカニズムは、各リーフスイッチについて、（たとえば、固有のパーティション化順序番号を有する各パーティションを介して）接続されたエンドノードをパーティション化に特有の順序でソートし得る（上記のリスト1の第4行）。この順序付けは、リーフスイッチにおける利用可能な上りポートの数を考慮に入れてノードがそれらのパーティションに従ってルーティングされることを確実なものとする助けとなるであろう。次いで、*p F T r e e* メカニズムは、*ROUTEDOWNGOINGBYASCENDINGTD*（上記のリスト1の第9行）などの関数を呼び出して、リスト2に示されるようにツリーを上って次のレベルのポートを選択して、*L I D*をルーティングし得る。

【 0 0 5 7 】

実施形態によれば、ポートの選択は、最小数の既に割り当てられている経路に基づく。これは、利用可能なパス全体にわたって負荷を分散させることを確実なものとする助けとなるであろう。しかし、いくつかのポートが同一の負荷で利用可能である場合には、当該関数は、これらの最小負荷ポートで繰返し処理を実行して、ルーティングされているノードのパーティションキーで既にマーキングされているスイッチに接続されているポートを選択し得る（リスト2の第3～9行）。マーキングされているスイッチがない場合（これは、このパーティションの第1のノードがルーティングされていることを意味し得る）、当該システムは、デフォルトで最も高いグローバル意識別子（*G U I D*）を有するポートの選択を実行し得る（リスト2の第2行）。パーティションにとって初めてスイッチが選択されると、当該スイッチはパーティションキーで下り方向にマーキングされる（リスト2の第11行）。

【 0 0 5 8 】

実施形態によれば、スイッチにおける L I D に対して下りポートが設定された後、当該メカニズムは、（リスト 3 の ROUTEUPGOINGBYDESCENDINGTD）を呼び出すツリーを下ることによって全ての接続された下りスイッチ上のそれに対して上りポートを割り当て得る。やはり、上りポートの選択は、まず負荷基準に基づき、次いで上り方向へのリモートスイッチのパーティションマーキングに基づき得る。次いで、全ての L F T が設定されるまで、次のレベルまでツリーを上ることによってプロセスが繰返され得る。なお、スイッチは複数のパーティションキーでマーキングされてもよい。p F T r e e メカニズムは、各パーティションについてルーティングされるノードの個数を記憶するテーブルを各スイッチについて維持し得る。マーキングされたパーティションを有するいくつかのスイッチがノードのルーティングに利用できる場合には、このカウンタを使用してポートの選択を決定することができる。パーティションについて最大数の既にルーティングされているノードを有するスイッチを選択することができる。

10

【 0 0 5 9 】

実施形態によれば、パーティション分離基準を保持しながらルーティングテーブルが生成されると、当該メカニズムは、リンクのうちのいくつか異なるパーティション内のノードに向かうフローに使用されているか否かを確認することに進み得る。そのような場合、当該メカニズムは、干渉しているパーティションに V L を割り当てて分離を提供することができる。V L 割り当てメカニズムの一例はリスト 4 に示されている。

20

【 0 0 6 0 】

実施形態によれば、V L 割り当てメカニズムは、パーティションで繰返し処理を実行して、パーティション内のノードによって使用されるいずれかの中間通信リンクが、別個の V L を割り当てられていない別のパーティションと中間リンクを共有しているか否かを確認し得る。このような状況に遭遇すると、新たな V L が割り当てられ得る。p F T r e e ルーティングメカニズムは、V L 選択のために 2 つのモード、すなわち厳格モードおよび通常モード、をサポートし得る。

【 0 0 6 1 】

実施形態によれば、厳格モードでは、p F T r e e ルーティングに必要とされる V L の数がシステム内の利用可能な V L を超えると、ルーティングは失敗する可能性がある（リスト 4 の第 10 行）。

30

【 0 0 6 2 】

実施形態によれば、通常モードでは、アルゴリズムはパーティションへの V L の割り当てを V L₁ から再開し得る（リスト 4 の第 8 行）。

【 0 0 6 3 】

実施形態によれば、I B ベースのファットツリーネットワークのための効率的なパーティション認識ルーティングメカニズム（別名で p F T r e e と称される）が提供される。p F T r e e メカニズムは、ファットツリートポロジのためのパーティションのネットワーク全体にわたる分離を提供することができる。また、p F T r e e は、スイッチのためにバランスのとれた L F T を生成する。ネットワークリソースが十分であると仮定すると、p F T r e e は、もっぱら物理リンクレベルでパーティションを分離することができる。たとえば、ファットツリーが 2 個の重なり合わない等しいサイズのパーティションを有している場合には、p F T r e e は、ルーティング自体に基づいて、中間ネットワークリンクを 2 個の等しいサイズの論理サブネットワークに分割することができる。さらに、ネットワークが完全なパーティション分離を提供するのに十分な利用可能なリソースを持たない場合には、p F T r e e は、物理的分離と連動して機能する相補的な V L ベースの分離スキームを利用することができる。

40

【 0 0 6 4 】

実施形態によれば、p F T r e e ルーティングメカニズムは、2 つの主な目的を達成することを目指している。第一に、当該メカニズムは、ツリー内のリンク全体にわたって均一に経路を分散させることによってファットツリートポロジのためにバランスのとれた L

50

F Tを生成することができる。第二に、p F T r e eは、リンク上の経路をバランスのとれた状態に維持しながら、異なるパーティションに属するパス間の干渉を取り除く。p F T r e eは、サブネットについてのパーティション化情報を使用することができ、パーティション内のノードが、他のパーティションでのワークロードの影響を受けない予測可能なネットワーク性能を受け取ることを確実にする。トポロジが（負荷バランシングに関して妥協することなく）各レベルにおいてパーティション分離を提供するのに十分な利用可能なリンクを持たない場合、p F T r e eは、V Lを使用して干渉の影響を軽減することができる。

【 0 0 6 5 】

実施形態によれば、p F T r e eメカニズムは、各エンドノードに関連付けられるL I Dについて全ての関連するスイッチ上でL F Tをセットアップするように再帰的に機能し得る。単一リーフスイッチパーティションをフィルタリングして除去した後、当該メカニズムは、各リーフスイッチについて、接続されたエンドノードをパーティション化に特有の順序でソートし得る。この順序付けは、リーフスイッチで利用可能な上りポートの数を考慮に入れてノードがそれらのパーティションに従ってルーティングされることを確実にする。各レベルでのポートの選択は、負荷が利用可能なパス全体にわたって分散されることを確実にするように、最小数の既に割り当てられている経路に基づき得る。しかし、いくつかのポートが同一の負荷で利用可能である場合、当該関数は、これらの最小負荷ポートで繰返し処理を実行して、ルーティングされているノードのパーティションキーで既にマーキングされているスイッチに接続されているポートを選択する。マーキングされているスイッチがない場合（すなわち、特定のパーティションの第1のノードのルーティング）、p F T r e eは、最も高いグローバル意識別子（G U I D）を有するポートのデフォルト選択の状態になり得る。パーティションにとって初めてスイッチが選択されると、当該スイッチはパーティションキーでマーキングされ得る。このようにして、当該メカニズムは、十分なパスがバランシングに利用可能であると仮定して、1つのパーティションに属するノードが同一のスイッチおよび対応するリンクを介してルーティングされることを確実なものとする助けとなるであろう。パーティション分離基準を保持しながらルーティングテーブルが生成されると、当該メカニズムは、リンクのうちのいくつかは異なるパーティション内のノードに向かうフローに使用されているか否かを確認することに進み得る。そのような場合、当該メカニズムは、干渉しているパーティションにV Lを割り当てて分離を提供することができる。

【 0 0 6 6 】

ここで、実施形態に係るマルチテナントクラスタ環境においてパーティション認識ルーティングをサポートすることの図である図5～図8を参照する。

【 0 0 6 7 】

実施形態によれば、p F T r e eルーティングにおけるポート選択メカニズムは、オーバーサブスライブされたファットツリーネットワークの単純なセクションによって図5～図8に示されている。

【 0 0 6 8 】

ここで、図5を参照する。図5は、4個のスイッチ、すなわちルートスイッチ525～526およびリーフスイッチ520～521と、8個のエンドノード、すなわちノードA～G 501～508とを有する2レベルファットツリートポロジを示す。同様に、エンドノードは2個のパーティションに分割されている。パーティション1は、ノードA 501とノードD 504とノードG 507とノードH 508とを備える。パーティション2は、ノードB 502とノードC 503とノードE 505とノードF 506とを備える。

【 0 0 6 9 】

実施形態によれば、図5に示されるように、例示的なセクションは、2個のリーフスイッチ（520および521）で構成され、各リーフスイッチは、4個のエンドノードと、リーフスイッチの次のレベルの2個のスイッチ、すなわちルートスイッチ（525および

10

20

30

40

50

526)とに接続されている。下り方向の割り当てられた経路の数を表わす変数 `down` および `max`、ならびに、各リンク上で適切なバランシングを確実にするようにルーティングされ得るノードの最大数もそれぞれ図に示されている。

【0070】

実施形態によれば、ルーティングすべき4個のエンドノードを有する各リーフスイッチに2個の上りポートがあると仮定すると、アップリンクの各々は、リンクのバランスをとることを確実にするために2個のエンドノードを下方にルーティングすべきである(すなわち、`max = 2`)。

【0071】

実施形態に従って、リーフスイッチ520についての、最初の2個のノード、すなわちノードAおよびノードB、のルーティングが図6に示されている。ルーティングメカニズムは、ルートスイッチ525を選択して、ノードAの方にトラフィックをルーティングし、図に「(パーティション1)」と示されているノードAのパーティションキーで当該スイッチをマーキングし得る。同様に、ノードBのために、ルートスイッチ526が選択され、図に「(パーティション2)」と示されているノードBのパーティションキーでマーキングされ得る。また、2個の下りリンクの各々上の単一のルーティングされるノードを数えるように変数 `down` も更新される。

【0072】

実施形態によれば、図7に示されるように、ノードCおよびDのために、対応するパーティションキーで既にマーキングされているスイッチが選択され得る。結果として生じる経路は、ルートスイッチ525を介して同一のリンクによって、第1のパーティションに属するノード、すなわちノードAおよびD、の方に流れる。同様に、第2のパーティションのノード、すなわちノードBおよびC、がルートスイッチ526を介して下り方向にルーティングされ得る。この経路の分離により、2個のパーティションのトラフィックフロー間の干渉が回避される。なお、各リンクに沿って下り方向にルーティングされるノードの数は `max` 変数を超えることはなく、これはルーティングが依然としてバランスがとれていることを意味する。

【0073】

最後に、実施形態に従って、図8は、リーフスイッチ521に接続されたエンドノードのためのルーティングを示す。ここでも、第2のレベルのスイッチが第1のリーフスイッチのルーティングからのパーティションキーで既にマーキングされているので、対応するスイッチが選択されて、ノード、すなわちノードE、F、GおよびH、の各々をルーティングすることができる。図に示されるように、最終的なルーティングは、ルーティングに基づいて中間ネットワークリンクを2個の等しいサイズの論理サブネットワークに分割することによって2個のパーティションを分離することができる。

【0074】

拡張型 p F T r e e

実施形態によれば、ネットワークがもっぱら物理リンクレベルでパーティションを分離するのに十分なリソースを持たない場合、`p F T r e e` ルーティングアルゴリズムは、`VL`を使用してパーティション間の干渉を減少させる。しかし、対応する`SLA`(サービス水準合意)または`QoS`要件によっては、パーティションが異なれば分離ニーズも異なる可能性がある。たとえば、ネットワーク内のパーティションのうちのいくつかは、重大な動作を実行しているかもしれず、いかなる場合でも完全な物理的分離を必要とするかもしれない。同様に、多くのネットワークでは、`VL`の利用可能性によっては、いくつかのパーティションは別のパーティションと`VL`を共有しなければならないかもしれない。これは通信集中的なワークロードにとっては望ましくないであろう。上記の`p F T r e e`アルゴリズムは、ルーティングにおける上記のパーティション面の要件を指定することができず、全てのパーティションは、同様の`QoS`要件を想定して等しい優先度で処理される。

【0075】

ここで、実施形態に係るマルチテナントクラスタ環境においてネットワーク分離をサバ

10

20

30

40

50

ートすることの図である図9を参照する。より具体的に、図9は、3個の異なるテナントパーティション(さまざまなシェーディングによって図示)に9個のノード(すなわち、ノードA~I 901~909)を有するファットツリーネットワークの一例を示す。ノードA 901およびE 905はパーティション1に属し、ノードB 902、F 906、G 907およびI 909はパーティション2に属し、ノードC 903、D 904およびH 908はパーティション3に属している。また、ファットツリーネットワークは、ルートスイッチ925および926と、リーフスイッチ920, 921および922を含む。

【0076】

実施形態によれば、パーティション1は非常に高いQoS要件を有し得て、このパーティションでのワークロードがいかなるパーティション間の干渉の影響も受けないことが極めて重要である。しかし、所与のファットツリーネットワークが3つの異なるテナントパーティションを有しながらルートスイッチを2つだけ有している、すなわちルートスイッチ925およびルートスイッチ926を有しているので、もっぱら物理レベルでこれらのパーティションを分離することは不可能である。上記のように、このような場合、p F T r e e ルーティングアルゴリズムは、V Lを使用してパーティションの分離を進めることができる。デフォルトのp F T r e e アルゴリズム(上記)を使用して得られるルーティングも図に示されており、ここではスイッチの直下の小さなノード円を使用して宛先ノードに向かうフローを示している。パーティション1のノードAに向かうトラフィックは、現在のところ、ルートスイッチ925とリーフスイッチ920との間のリンクを、異なるパーティション(すなわち、パーティション3)に属するノードCに向かうトラフィックと共有していることが分かる。同様に、パーティション1(高いQoSを必要とするパーティション)のEに向かうトラフィックは、パーティション2のノードFに向かうトラフィックとリンクを共有している。ルーティングメカニズムは、分離を保つために、これらのリンクの各々上で別個の仮想レーンを利用することができる。しかし、V Lを使用することにより干渉が減少するとしても、このような干渉が完全に排除されることはない。

【0077】

分離ポリシー

実施形態によれば、p F T r e e ルーティングメカニズムは、パーティションワイズの(partition-wise)分離ポリシーおよびグローバルな分離ポリシーを含むように拡張可能である。各パーティションについて、分離ポリシーは、どのようにしてパーティション内のノードが他のパーティションに属するノードとネットワークリソースを共有できるようにするかを決定し得る。グローバルなポリシーは、所与のネットワークについて全てのパーティションワイズの分離ポリシーを満たすことができない場合にルーティングが機能しなくなるかベストエフォート型分離を継続するかを判断し得る。

【0078】

実施形態によれば、拡張型p F T r e e ルーティングメカニズムのためのさまざまなポリシーパラメータを提供することができる。各パーティションは、3個のパーティションワイズのポリシーパラメータのうちの1つでマーキングされ得る。パーティションをp h y - 分離(本明細書では厳格なパラメータとも称される)でマーキングすることにより、ルーティングアルゴリズムが特に当該パーティションのためにネットワークリソースを取っておくことを保証することができ、当該パーティション内のノードは、異なるパーティション内のその他のノードといかなるリンクも共有しない。パーティションをv l a n e - 分離というパラメータ(本明細書では厳格な仮想レーンパラメータとも称される)でマーキングすることにより、マーキングされたパーティションは、別個のV Lのみを使用して他のパーティションとネットワークリソースを共有することができる。パーティションをd e f - 分離(本明細書ではベストエフォートパラメータとも称される)スキームでマーキングすることにより、マーキングされたパーティションのベストエフォート型分離が実現される。

【0079】

実施形態によれば、ポリシーパラメータは、グローバルなポリシーパラメータも含み得る。グローバルなポリシーパラメータ、すなわち厳格なパラメータおよびベストエフォートパラメータ、は、所与のサブネットにおいてパーティションワイズのポリシーパラメータを満たすことができない場合にルーティングメカニズムが機能しなくなるかベストエフォート型ルーティングに戻るかを規定し得る。たとえば、ネットワークが所望の分離を提供するための十分なリンクまたはVLを持たない場合、パーティション構成ファイルを使用してポリシーパラメータがルーティングメカニズムに提供されてもよい。

【0080】

拡張型 p F T r e e メカニズム

実施形態によれば、拡張型 p F T r e e ルーティングメカニズム（本明細書では別名で「p F T r e e - E x t」とも称される）は、ファブリックを再帰的に横断して、各エンドノードに関連付けられる L I D について全てのスイッチにおいて L F T をセットアップすることによって、（上記の）当初の p F T r e e と同様の態様で動作する。しかし、p F T r e e とは異なって、p F T r e e - E x t は、経路を割り当てる際に規定のグローバルな分離ポリシーおよびパーティションワイズの分離ポリシーも考慮に入れることができる。

【0081】

実施形態に従って、p F T r e e - E x t ルーティングメカニズムの疑似コード（本明細書では別名でリスト5とも称される）を以下に示す。

【0082】

【数5】

Ensure: The LFTs are generated for the switches conforming isolation policies

```

1: global_param ← get_global_isolation_policy()
2: partitions_info ← get_partition_information()
3: ORDERCOMPUTENODES()
4: for each sw ∈ leafSwitches[] do
5:   for each cn ∈ computeNodes[] do
6:     Get lid of cn
7:     Get partition_key of the cn.hca_port
8:     Set LFT[lid] ← cn:hca_port on sw
9:     ROUTEDOWNGOINGBYASCENDING() on sw
10:  end for
11: end for
12: ASSIGNVIRTUALLANES()
13: VALIDATEPOLICIES()
```

【0083】

実施形態によれば、当該メカニズムは決定論的であり、経路は宛先ノードから開始して後向きに計算される。当該メカニズムは、まず、コンピュータノードをパーティションに特有の順序でソートし得る（リスト5の第3行）。パーティションに特有の順序は、メカニズムのさらに高速な実行を確実にすることができる。なぜなら、ノードは、一旦順序付けられると、各々の下りおよび上りポート上に最大カウンタを維持することなく繰返しルーティングできるからである。以下のORDERCOMPUTENODESのための疑似コードに示されるように、各リーフスイッチについて、ORDERCOMPUTENODESは、まず、エンドノードをそれらのパーティションポリシー優先度の低い方から順にソートする（リスト6の第4行（以

10

20

30

40

50

下を参照))。p h y - 分離パラメータでマーキングされたパーティションに属するノードが最初に付加され得て、v l a n e - 分離によるパーティションが2番目に付加され得る。最後に、d e f - 分離のポリシーパラメータ値を有するパーティションノードがコンピュータノードのリストに付加される。次いで、当該メカニズムは、ノードのパーティション化情報を使用してルーティング順序を生成し、当該ルーティング順序では、1つのパーティションに属するノードは、反復ルーティングの際にネットワーク内の同一の上りリンクを提案するインデックスを得る傾向がある。これは、パーティションキーテーブルを使用して、利用可能な上りポートの数を、パーティションに属する第1のノードをルーティングするように選択されたインデックスに付加することによってなされる(リスト6の第14~28行)。しかし、このようなインデックスが既に取得されている場合、または当該インデックスがコンピュータレイ境界を超える場合、第1のフリーインデックスが選択されて、その後の選択のためにパーティションキーでマーキングされ得る(リスト6の第24行)。ORDERCOMPUTENODESのための疑似コード(本明細書では別名でリスト6とも称される)をここに示す。

【0084】

【数 6】

Require: List of switches and attached compute nodes

Ensure: The compute nodes are ordered for the pFTree-Ext routing algorithm

```

1: for each sw in leafswitches[] do
2:   num_up_ports ← count(sw → upPorts[])
3:   num_cns ← count(sw → computeNodes[])
4:   Sort nodes in increasing order of partition isolation policy      10
     (phy > vlane > def)
5:   if num_cns ≥ num_up_ports then
6:     return
7:   end if
8:   index_arr[] = array(num_cns)
9:   taken[] = array(num_cns)
10:  pkey_tbl[] = map()
11:  id ← 0                                                              20
12:  for each cn in sw → computeNodes[] do
13:    pkey ← cn → get_partition_key()
14:    if pkey not found in pkey_tbl then
15:      if taken[id] ≠ false then
16:        id ← get_free_id()
17:      end if
18:      index_arr[cn[i]] ← id
19:      taken[id] = true                                                30
20:      insert pkey in pkey_tbl
21:    else {pkey is already in pkey_tbl}
22:      id ← id(pkey) + num_up_ports
23:      if id ≥ num_cns or taken[id] = true then
24:        id ← get_free_id()
25:      end if
26:      index_arr[cn[i]] ← id                                          40
27:      taken[id] = true
28:      update pkey_tbl
29:    end if
30:  end for
31:  Sort sw → computeNodes[] with respect to index_arr[]
32: end for

```

【0085】

実施形態によれば、ノードが適切に順序付けられると、ROUTEDOWNGOINGBYASCENDINGのための以下の疑似コード（本明細書では別名でリスト7とも称される）に例示されるよう

に、pFTree - Extメカニズムは、ROUTEDOWNGOINGBYASCENDING (リスト5の第9行) を呼び出して、ツリーを上って次のレベルのポートを選択して、LIDを下り方向にルーティングし得る。

【0086】

【数7】

```

Require: A switch sw, an end node lid and partition_key
1: Sort sw.upPorts[] with increasing load and then GUID
2: Get least loaded ports as leastLoadedList[]
3: partition_param ← get_isolation_policy(partition_key)
4: selected_port ← null
5: for each port in leastLoadedList[] do
6:   r_sw ← port.get_remote_switch()
7:   if r_sw is marked with partition_key then
8:     selected_port ← port
9:     break
10:  end if
11: end for
12: if selected_port = null then
13:  while selected_port = null do
14:    port ← sw.upPorts[].get_next()
15:    r_sw ← port.get_remote_switch()
16:    if r_sw is marked with a partition with isolation policy >
       partition_param then
17:      continue
18:    end if
19:    selected_port ← port
20:  end while
21: end if
22: Set LFT[lid] ← selected_port on r_sw
23: if r_sw is not marked then
24:  Mark it with partition_key in DWN direction
25: end if
26: ROUTEUPGOINGBYDESCENDING() on sw
27: ROUTEDOWNGOINGBYASCENDING() on r_sw

```

【0087】

実施形態によれば、ポートの選択は、まず、ソートされた利用可能な上りポートから得られる最小負荷ポートリストに基づく(リスト7の第1~2行)。当該関数は、これらの最小負荷ポートで繰返し処理を実行して、ルーティングされているノードのパーティションキーで既にマーキングされているスイッチに接続されているポートを選択する(リスト7の第5~11行)。マーキングされているスイッチがないことが分かると、アルゴリズムは、上りポートで繰返し処理を実行して、LIDに適した経路を見つける。上りポート

リストは、ポートに対する現在の負荷の低い方から順にソートされる。同一の負荷を有するポートでは、決定論的なままであるようにするために、ソートはそれらのグローバル一意識別子 (GUID) の高い方から順になされる。さらに、当該関数は、ルーティングされるノードよりも高い分離ポリシーパラメータを有するパーティションキーで既にマーキングされているポートを選択しない (リスト 7 の第 16 ~ 17 行)。最後に、ポートが選択されると、対応するスイッチがパーティションキーで下り方向にマーキングされる (リスト 7 の第 24 行)。

【 0 0 8 8 】

実施形態によれば、スイッチにおける LID に対して下りポートが設定されると、p F T r e e - E x t メカニズムは、ROUTEUPGOINGBYDESCENDING (本明細書では別名でリスト 8 と称される) を呼び出すツリーを下ることによって全ての接続された下りスイッチ上でそれに対して上りポートを割り当てる。ROUTEUPGOINGBYDESCENDING のための疑似コードをここに示す。

【 0 0 8 9 】

【 数 8 】

```
Require: A switch sw, an end node lid and partition_key
1: Get least-loaded ports from sw.dwnPorts[] as dwnlist[]
2: selected_port ← dwnlist.get port_max_guid()
3: for each port in dwnlist[] do
4:   r_sw ← port.get_remote_switch()
5:   if r_sw is marked with partition_key then
6:     selected_port ← port
7:     break
8:   end if
9: end for
10: if r_sw is not marked then
11:   Mark it with partition_key in UP direction
12: end if
13: Set LFT[lid] ← selected_port on r_sw
14: ROUTEUPGOINGBYDESCENDING() on r_sw
```

【 0 0 9 0 】

実施形態によれば、上りポートの選択は、まず負荷基準に基づき、次いで今回は上り方向へのリモートスイッチのパーティションマーキングに基づく。次いで、全ての L F T が設定されるまで、次のレベルまでツリーを上ることによってプロセスが繰返される。なお、スイッチは複数のパーティションキーでマーキングされてもよい。p F T r e e - E x t メカニズムは、各パーティションについてルーティングされるノードの個数を記憶するテーブルを各スイッチについて維持する。マーキングされたパーティションを有するいくつかのスイッチがノードのルーティングに利用できる場合には、このカウンタを使用してポートの選択を決定する。パーティションについて最大数の既にルーティングされているノードを有するスイッチが選択される。

【 0 0 9 1 】

実施形態によれば、パーティション分離基準を保持しながらルーティングテーブルが生成されると、p F T r e e - E x t メカニズムは、リンクのうちのいくつか異なるパーティション内のノードに向かうフローに使用されているか否かを確認することに進む。そのような場合、p F T r e e - E x t メカニズムは、干渉しているパーティションに V L

を割り当てて分離を提供することができる。V L 割り当て関数、すなわちASSIGNVIRTUAL ANES、のための疑似コード（本明細書では別名でリスト9と称される）を以下に示す。

【0092】

【数9】

```
Require: The pFTree-Ext routing tables have been generated
Require: Switches have been marked with the partition keys
Require: Global policy parameter, strict or best-effort
Ensure: A partitions marked with vl-isolation has a separate VL
Ensure: No two partitions with the same SL share a link
1: vlans_needed ← 1
2: max_vlans ← get_max_lanes()
3: strict ← get_is_strict()
4: for each partition in partition_tbl do
5:   check if the isolation policy of the partition is vl-isolation
   and
       any intermediate communication link in this partition share a
       switch with a partition that has not been assigned a virtual
       lane
6:   if require a separate vl then
7:     if vlans_needed = max_vlans and global_param = strict then
8:       vlans_needed ← 1
9:     else
10:      error: routing failed
11:      return
12:    end if
13:    vlans_needed++
14:    partition.vlane ← vlans_needed
15:  end if
16: end for
```

10

20

30

【0093】

実施形態によれば、仮想レーン割り当て関数は、全てのパーティションで繰返し処理を実行して、パーティションがV L - 分離ポリシーパラメータでマーキングされているか否かおよびパーティション内のノードによって使用されるいずれかの中間通信リンクが、別個のV L を割り当てられていない別のパーティションと中間リンクを共有しているか否かを確認し得る。そうであれば、新たなV L が割り当てられる。また、V L 割り当て関数は、2つのモード、すなわち厳格モードおよびベストエフォートモード、を有するグローバルなポリシーパラメータを使用し得る。厳格モードでは、p F T r e e - E x t ルーティングメカニズムに必要とされるV L の数がシステム内の利用可能なV L の数を超えると、ルーティングは失敗する（リスト9の第10行）。ベストエフォートモードでは、仮想レーン割り当て関数はパーティションへのV L の割り当てをV L₁から再開し得る（リスト9の第8行）。

40

【0094】

実施形態によれば、p F T r e e - E x t ルーティングメカニズムは、全ての利用可能なV L ではなくV L の特定のグループを考慮に入れるように容易に変更することができる。同様に、より高い分離ポリシーを有するパーティションがV L を共有しにくくするため

50

に、全ての利用可能なVLが使用されると、(VL₁を選択する代わりに)割り当てられたパーティションの優先度を選択のために減少させることによってVLリストを順序付けることができる。VLが割り当てられた後、pFTree-Extルーティングアルゴリズムは、全てのパーティションワイズのポリシーおよびグローバルなポリシーが満たされているか否かを確認する(リスト5の第13行)。

【0095】

実施形態によれば、pFTree-Extメカニズムは、分離ポリシーをルーティングアルゴリズムに組み込むことができる。各リーフスイッチについてルーティングの前にエンドノードをパーティションに特有の順序でソートするpFTreeとは異なって、pFTree-Extルーティングメカニズムは、まず、エンドノードをそれらのパーティション優先度の順序でソートする。phy-分離でマーキングされたパーティション内のエンドノードは、最大の優先度を得る。その後、当該メカニズムは、上記のようにエンドノードをパーティションに特有の順序でソートすることに取りかかる。最も高いパーティション優先度を有するノードが最初にルーティングされることを確実にするために、付加的なソートが事前になされる。

【0096】

実施形態によれば、pFTree-Extメカニズムはさらに、新たなノードをルーティングするためのポートの選択方法を変更することによって分離ポリシーをルーティングアルゴリズムに組み込むことができる。たとえば、いくつかの候補ポートの中から下りポートを選択するために、pFTree-Extは、ポートに対する現在の負荷を確認することに加えて、対応するスイッチが現在ルーティングされているノードのパーティションよりも高い優先度を有するパーティションのキーで既にマーキングされている任意のポートグループを除去する。

【0097】

実施形態によれば、および加えて、利用可能なネットワークリソースがパーティションワイズのポリシーパラメータを満たすことを許可しない場合、pFTree-Extルーティングメカニズムは、機能しなくなるかまたはグローバルなポリシーパラメータに従って続行するかのいずれかであり得る。その場合、当初のpFTreeルーティングアルゴリズムは、利用可能なVLを考慮に入れるのみである。

【0098】

ここで、実施形態に係るマルチテナントクラスタ環境においてネットワーク分離をサポートすることの図である図10を参照する。より具体的に、図10は、3個の異なるテナントパーティション(さまざまなシェーディングによって図示)に9個のノード(すなわち、ノードA 901~ノードI 909)を有するファットツリーネットワークの一例を示す。ノードA 901およびE 905はパーティション1に属し、ノードB 902、F 906、G 907およびI 909はパーティション2に属し、ノードC 903、D 904およびH 908はパーティション3に属している。また、ファットツリーネットワークは、ルートスイッチ925および926と、リーフスイッチ920、921および922を含む。

【0099】

実施形態によれば、図10は、pFTree-Extメカニズムを使用したサブネットルーティングを示しており、当該メカニズムでは、パーティション1(すなわち、ノードA 901およびノードE 905)はphy-分離などの高い優先度でマーキングされており、ルーティングメカニズムが特に当該パーティションのためにネットワークリソースを取っておくことを保証することができ、当該パーティション内のノードは、異なるパーティション内のその他のノードといかなるリンクも共有しない。結果として生じるルーティングが図10に示されている。パーティション1がphy-分離などの高い優先度でマーキングされているので、パーティション1のノード(すなわち、ノードAおよびE)はいずれも、その他のパーティションとリンクを共有しない。しかし、このようなポリシーはパーティション2および/またはパーティション3には適用されなかったため、これ

らのパーティションはスイッチ 926 からの全ての下りリンクを共有する。

【0100】

重み付けされた p F T r e e ルーティングメカニズム

実施形態によれば、p F T r e e ルーティングメカニズムの第 2 の拡張機能は、サブネットにおけるトラフィック特性の重みを説明することができる。これは、重み付けされた p F T r e e ルーティングメカニズム (p F T r e e - W t) と称されることができる。p F T r e e - W t は、各コンピュータノードに関連付けられる重みの概念に基づく。これらの重みを使用して、経路を計算する際に既知のまたは学習したトラフィック特性を考慮に入れる。パーティション化にかかわらず、ノードの重みは、ルーティングテーブルを計算する際にノードに向かうフローが受ける優先度を反映する。たとえば、考えられる構成は、ネットワーク内でどれぐらいのトラフィックをノードが受けることが分かっているかに応じて [1 , 1 0 0] の範囲内の重みをノードに割り当てるというものであり得る。このようなスキームでは、ほとんどトラフィックを受けないノード（たとえば、主にトラフィックジェネレータ）に対しては重み = 1 を割り当て、リンク容量に近いトラフィックを受けるノードに対しては重み = 1 0 0 を割り当てることができる。その結果、中間の値、すなわち $1 < x < 100$ は、ネットワーク内でノードが受けることが見込まれるトラフィックの割合を反映することができる。

10

【0101】

実施形態によれば、コンピュータノードについての管理情報が利用できない場合には、ポートデータカウンタベースのスキームを使用して重みを計算することができる。たとえば、O F E D (オープン・ファブリック・エンタープライズ・ディストリビューション) では、データカウンタを読取るために ibdatacount と呼ばれるユーティリティが提供される。全てのノードについて初期重みが等しい状態でネットワークをセットアップした後、所定の期間が経過すると新たな重みを学習することができる。

20

【0102】

実施形態によれば、B がある期間にわたって測定される全てのノードの受信帯域幅セットを表わすとすると、各ノードの重みは、以下のように、線形変換を使用して [a , b] の範囲内で割り当てられることができる。

【0103】

【数 1 0】

$$W(x) = (x - a) \frac{b - a}{\max(B) - \min(B)} + a, \forall x \in B$$

30

【0104】

実施形態によれば、p F T r e e - W t ルーティングメカニズムを使用して、各コンピュータノードは、重みというパラメータを割り当てられることができる。ポートに対する負荷が上りおよび下り方向へのノードに向かう割り当てられた経路の数を表わす当初の p F T r e e ルーティングとは異なって、p F T r e e - W t ルーティングスキームにおけるポートに対する負荷は、当該ポートから各方向にルーティングされるコンピュータノードの累積重みである。各リーフスイッチについて、1 つのパーティション内のノードは、ルーティングの前にそれらの重みによってもソートされる。コンピュータノードをルーティングするためにスイッチにおける下りポートが選択されると、p F T r e e - W t は、対応するコンピュータノードの重みを付加することによって、選択されたポートに対する現在の負荷を更新する。同様に、上りリンクでは、各ポートに対して上り方向の負荷が維持される。ポート選択基準は p F T r e e ルーティングと同様であり、ノードのパーティションも考慮する。しかし、ポートカウンタとは異なって、p F T r e e - W t における各レベルでのポートの選択は、全ての利用可能なポートに対する最小累積重みに基づく。いくつかのポートが同一の負荷で利用可能である場合、当該メカニズムは、これらの最小負荷ポートで繰返し処理を実行して、ルーティングされているノードのパーティションキーで既にマーキングされているスイッチに接続されているポートを選択する。ルーティン

40

50

グテーブルが生成されると、p F T r e e - W t は、V L 割り当てを実行して、ネットワーク内でリンクを共有している異なるパーティションに関連付けられたノードに異なるV L が割り当てられることを確実にすることができる。

【 0 1 0 5 】

図 1 1 は、実施形態に係るマルチテナントクラスタ環境において重み付けされたパーティション認識ルーティングをサポートすることの図である。具体的に、図 1 1 は、4 個のスイッチ、すなわちルートスイッチ 1 1 2 5 ~ 1 1 2 6 およびリーフスイッチ 1 1 2 0 ~ 1 1 2 1 と、8 個のエンドノード、すなわちノード A ~ G 1 1 0 1 ~ 1 1 0 8 とを有する 2 レベルファットツリートポロジを示す。同様に、エンドノードは 2 個のパーティションに分割されている。パーティション 1 は、ノード A 1 1 0 1 とノード D 1 1 0 4 とノード G 1 1 0 7 とノード H 1 1 0 8 とを備える。パーティション 2 は、ノード B 1 1 0 2 とノード C 1 1 0 3 とノード E 1 1 0 5 とノード F 1 1 0 6 とを備える。

【 0 1 0 6 】

実施形態によれば、図 1 1 における各ノードは、重みを割り当てられている。ノード A 1 1 0 1 は 1 0 0 という重みを割り当てられている一方、残りのノードは 1 という重みを割り当てられている。p F T r e e - W t を使用したリーフスイッチ 1 1 2 0 への下りルーティングが図 1 1 に示されている。リーフスイッチ 1 1 2 0 に接続されたノードをルーティングする際、スイッチ 1 1 2 5 および 1 1 2 6 にそれぞれ接続された 2 個の上りポートが利用可能である。ノード A は、1 0 0 に等しい重みを有しているため、それらのリンクのうちの 1 つ、すなわちスイッチ 1 1 2 5 スwitch 1 1 2 0 を割り当てられる一方、他の 3 個のノードは他のリンク、すなわちスイッチ 1 1 2 6 スwitch 1 1 2 0 を共有する。その理由は、他の 3 個のノードの重みの合計が 1 0 0 よりも低い 3 に過ぎないからである。たとえ選択されたスイッチがパーティションキーでマーキングされたとしても、重み付けされたパーティション認識ルーティングにより、依然としてサブネットにおいてパーティションを分離することはできない。

【 0 1 0 7 】

図 1 2 は、実施形態に係るマルチテナントクラスタ環境において重み付けされたパーティション認識ルーティングをサポートすることの図である。具体的に、図 1 2 は、4 個のスイッチ、すなわちルートスイッチ 1 1 2 5 ~ 1 1 2 6 およびリーフスイッチ 1 1 2 0 ~ 1 1 2 1 と、8 個のエンドノード、すなわちノード A ~ G 1 1 0 1 ~ 1 1 0 8 とを有する 2 レベルファットツリートポロジを示す。同様に、エンドノードは 2 個のパーティションに分割されている。パーティション 1 は、ノード A 1 1 0 1 とノード D 1 1 0 4 とノード G 1 1 0 7 とノード H 1 1 0 8 とを備える。パーティション 2 は、ノード B 1 1 0 2 とノード C 1 1 0 3 とノード E 1 1 0 5 とノード F 1 1 0 6 とを備える。

【 0 1 0 8 】

実施形態によれば、図 1 2 における各ノードは、重みを割り当てられている。ノード A 1 1 0 1 は 1 0 0 という重みを割り当てられている一方、残りのノードは 1 という重みを割り当てられている。p F T r e e - W t を使用したリーフスイッチ 1 1 2 1 への下りルーティングが図 1 2 に示されている。リーフスイッチ 1 1 2 0 へのルーティングとは異なって、リーフスイッチ 1 1 2 1 に接続された各ノードは、同一の重み（すなわち、1）を有している。このため、パーティションはリンク上で分離されたままである可能性がある。同一のパーティションに属するノード G および H は、スイッチ 1 1 2 5 スwitch 1 1 2 1 間のリンクを介して下り方向にルーティングされることができる。等しい重みを有するパーティション 2 のノード E および F は、スイッチ 1 1 2 6 スwitch 1 1 2 1 間のリンクを介して下り方向にルーティングされることができる。

【 0 1 0 9 】

実施形態によれば、p F T r e e - W t は、パーティションが分離した状態をできるだけ保持しながらリンク上の重み付けされた負荷バランシングを満足させる。なお、最終的なルーティングは、図 1 2 に示されるように、2 個のパーティションのノードによって共有されるリンクを 1 つだけ有する。

【 0 1 1 0 】

図 1 3 は、実施形態に係るマルチテナントクラスタ環境においてネットワーク分離をサポートするための方法のフローチャートである。

【 0 1 1 1 】

ステップ 1 3 0 1 において、当該方法は、マルチテナントクラスタ環境内の 1 つ以上のテナントをサポートし得る。

【 0 1 1 2 】

ステップ 1 3 0 2 において、当該方法は、1 つ以上のテナントの各々を複数のパーティションのうちのあるパーティションに関連付け得る。

【 0 1 1 3 】

ステップ 1 3 0 3 において、当該方法は、複数のパーティションの各々を複数のノードのうちの 1 つ以上のノードに関連付け得て、複数のノードの各々は、複数のスイッチのうちのリーフスイッチに関連付けられ、複数のスイッチは、複数のリーフスイッチと、別のレベルにおける少なくとも 1 つのスイッチとを備える。

【 0 1 1 4 】

ステップ 1 3 0 4 において、当該方法は、複数のパーティションの各々を複数のポリシーパラメータのうちのあるポリシーパラメータでマーキングし得る。

【 0 1 1 5 】

ステップ 1 3 0 5 において、当該方法は、複数のノードの各ノードにパーティション化順序を割り当て得て、パーティション化順序は、少なくとも各ノードに関連付けられるパーティションにマーキングされたポリシーパラメータに基づく。

【 0 1 1 6 】

ステップ 1 3 0 6 において、当該方法は、少なくとも複数のパーティションのうちのパーティションのマーキングに基づいて、マルチテナントクラスタ環境で使用される 1 つ以上のリニアフォワーディングテーブルを生成し得る。

【 0 1 1 7 】

図 1 4 は、本発明の実施形態を実現することができる例示的なコンピュータシステム 1 4 0 0 を示すブロック図である。コンピュータシステム 1 4 0 0 は、バス 1 4 9 0 を介して電氣的に結合され得るハードウェア要素を備えるものとして示されている。ハードウェア要素は、1 つ以上の中央処理ユニット 1 4 1 0 と、1 つ以上の入力デバイス 1 4 2 0 (たとえば、マウス、キーボードなど)と、1 つ以上の出力デバイス 1 4 3 0 (たとえば、ディスプレイデバイス、プリンタなど)とを含んでいてもよい。

【 0 1 1 8 】

コンピュータシステム 1 4 0 0 は、1 つ以上の記憶装置 1 4 4 0 も含んでいてもよい。一例として、記憶装置 1 4 4 0 は、ランダムアクセスメモリ(「RAM」)および/またはリードオンリメモリ(「ROM」)などの、ディスクドライブ、光学式記憶装置、ソリッドステート記憶装置であってもよく、これらはプログラム可能、フラッシュ更新可能などであってもよい。記憶装置は、フロッピー(登録商標)ディスク、光ディスク、DVD、CD-ROM、マイクロドライブ、および光磁気ディスクを含む任意のタイプのディスク、ROM、RAM、EPROM、EEPROM、DRAM、VRAM、フラッシュメモリデバイス、磁気もしくは光カード、(分子メモリICを含む)ナノシステム、または、命令および/もしくはデータを格納するのに適した任意のタイプの媒体もしくはデバイスを含み得るが、それらに限定されるものではない。

【 0 1 1 9 】

また、コンピュータシステム 1 4 0 0 は、コンピュータ読取可能記憶媒体読取装置 1 4 5 0 と、通信システム 1 4 6 0 (たとえば、モデム、ネットワークカード(無線または有線)、赤外線通信デバイス、ブルートゥース(登録商標)デバイス、セルラー通信デバイスなど)と、上記の RAM および ROM デバイスを含み得るワーキングメモリ 1 4 8 0 とを含んでいてもよい。いくつかの実施形態では、コンピュータシステム 1 4 0 0 は、デジタル信号プロセッサ、特殊用途プロセッサなどを含み得る処理加速ユニット 1 4 7 0 も含

10

20

30

40

50

んでいてもよい。

【 0 1 2 0 】

コンピュータ読取可能記憶媒体読取装置 1 4 5 0 はさらに、コンピュータ読取可能な記憶媒体に接続可能であり、ともに（および任意に、記憶装置 1 4 4 0 と組み合わせて）コンピュータ読取可能な情報を一時的におよび／またはより永久的に収容するためのリモートの、ローカルな、固定式のおよび／またはリムーバブルな記憶装置プラス記憶媒体を包括的に表わす。本発明の特徴は、機械読取可能な媒体のいずれか 1 つに格納されて、ソフトウェアおよび／またはファームウェアに組み込まれてもよく、当該ソフトウェアおよび／またはファームウェアは、コンピュータシステムのハードウェアを制御して、本発明の結果を利用する他のメカニズムとコンピュータシステムとを対話させるためのものである。このようなソフトウェアまたはファームウェアは、アプリケーションコード、デバイスドライバ、オペレーティングシステムおよび実行環境／コンテナを含み得るが、それらに限定されるものではない。通信システム 1 4 6 0 は、上記のネットワーク、システム、コンピュータおよび／または他の構成要素とデータをやりとりすることを可能にし得る。

10

【 0 1 2 1 】

コンピュータシステム 1 4 0 0 は、ソフトウェア要素も備えていてもよく、当該ソフトウェア要素は、オペレーティングシステム 1 4 8 8 および／または他のコード 1 4 8 4 を含むワーキングメモリ 1 4 8 0 内に現在のところ置かれているものとして示されている。コンピュータシステム 1 4 0 0 の代替的な実施形態は上記のものに対する多数の変形例を有していてもよいということが理解されるべきである。たとえば、カスタマイズされたハードウェアも使用されてもよく、および／または、特定の要素がハードウェア、（アプレットなどのポータブルソフトウェアを含む）ソフトウェア、またはそれら両方で実現されてもよい。さらに、ネットワーク入出力デバイスおよびデータ取得デバイスなどの他のコンピュータシステムとの接続も行われてもよい。

20

【 0 1 2 2 】

コンピュータシステム 1 4 0 0 のソフトウェアは、本明細書に記載されているアーキテクチャのさまざまな要素のいずれかまたは全ての機能を実現するためのコード 1 4 8 4 を含んでいてもよい。たとえば、システム 1 4 0 0 などのコンピュータシステムに格納され、および／または、システム 1 4 0 0 などのコンピュータシステムによって実行されるソフトウェアが、上記のような本発明の機能および／または他の構成要素を提供してもよい。これらの構成要素のうちのいくつかのソフトウェアによって実現可能な方法については、より詳細に上記している。

30

【 0 1 2 3 】

また、本発明の特徴は、たとえば特定用途向け集積回路（application specific integrated circuit: A S I C）などのハードウェアコンポーネントを使用してハードウェアで実現されてもよい。本明細書に記載されている機能を実行するためのハードウェア状態機械の実現は、当業者に明らかであろう。

【 0 1 2 4 】

本発明のさまざまな実施形態を上記してきたが、それらは限定としてではなく一例として提示されているということが理解されるべきである。本発明の精神および範囲から逸脱することなく形態および詳細のさまざまな変更を行ってもよいことは当業者に明らかであろう。

40

【 0 1 2 5 】

特定の機能の実行およびそれらの関係を示す機能的構成要素を用いて本発明を上記してきた。これらの機能的構成要素の範囲は、説明の便宜上、本明細書では多くの場合任意に規定されている。特定の機能およびそれらの関係が適切に実行される限り、代替的な範囲が規定されてもよい。したがって、いかなるこのような代替的な範囲も本発明の範囲および精神の範囲内である。

【 0 1 2 6 】

本発明の上記の説明は、解説および説明の目的で提供されている。本発明の上記の説明

50

は、網羅的であることを意図したものではなく、本発明を開示されている厳密な形態に限定することを意図したものでもない。本発明の幅および範囲は、上記の例示的な実施形態のいずれによっても限定されるべきではない。多くの変更例および変形例が当業者に明らかであろう。当該変更例および変形例は、開示されている特徴のいずれかの関連性のある組み合わせを含む。本発明の原理およびその実際の適用を最良に説明するために実施形態が選択されて説明されており、それによって、さまざまな実施形態についておよび意図される特定の使用に適したさまざまな変更例とともに当業者は本発明を理解することができる。本発明の範囲は以下の特許請求の範囲およびそれらの等価物によって規定されることが意図される。

【図 1】

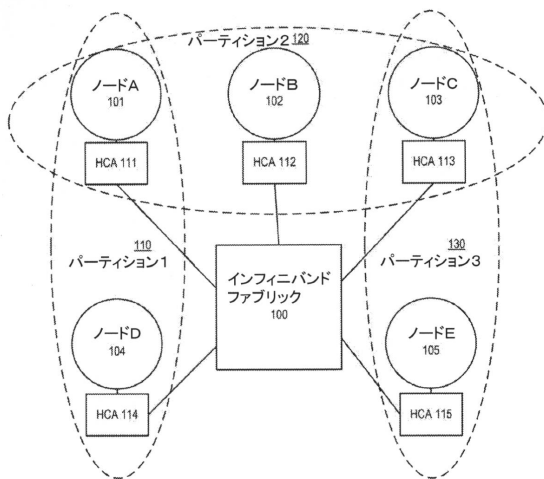


FIGURE 1

【図 2】

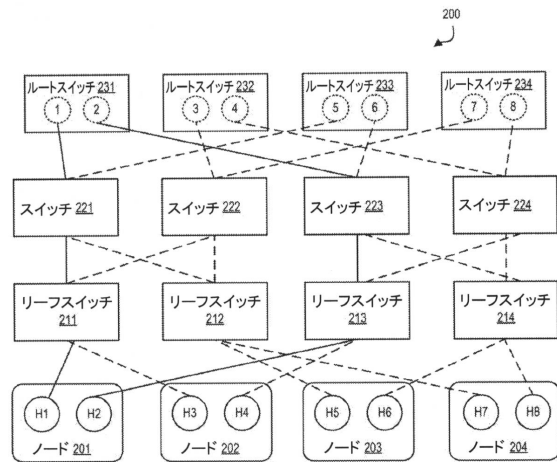


FIGURE 2

【図 3】

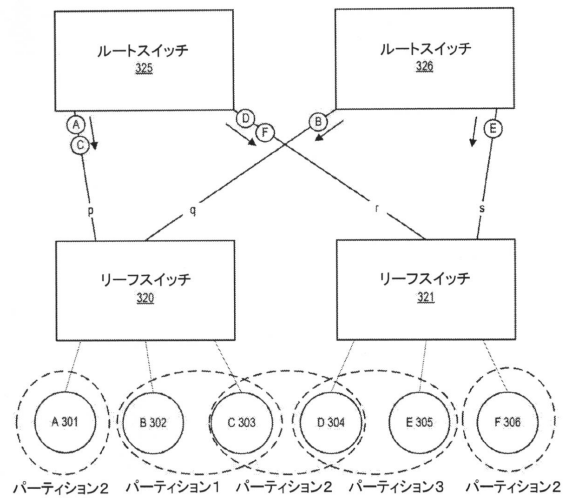


FIGURE 3

【図 4】

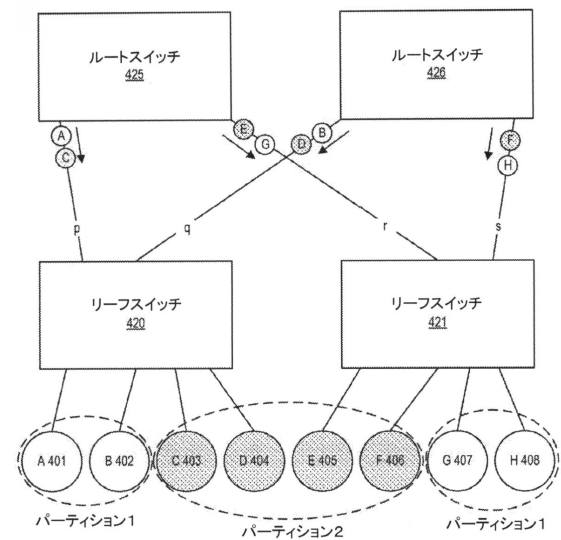


FIGURE 4

【図 5】

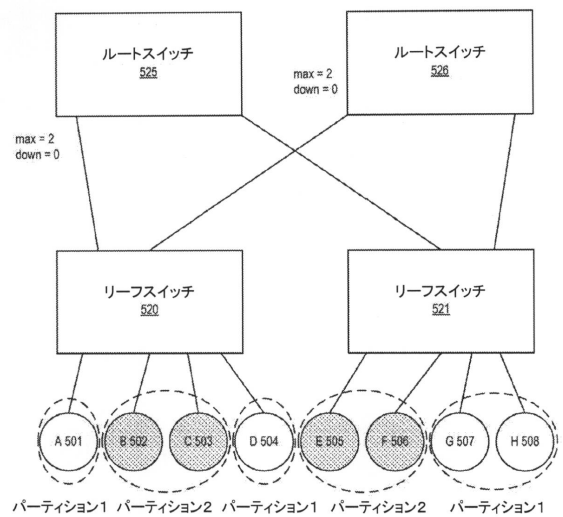


FIGURE 5

【図 6】

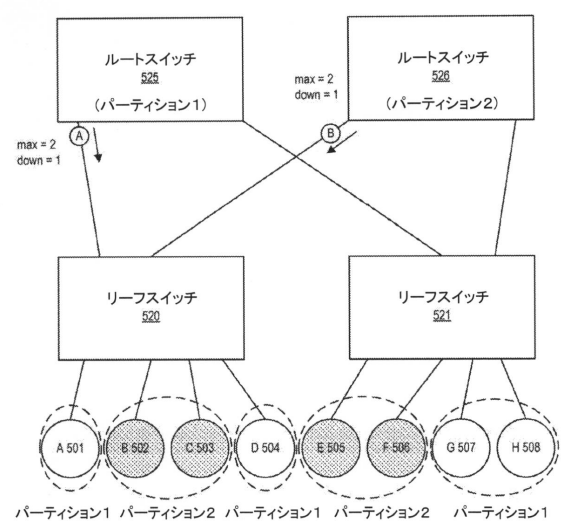


FIGURE 6

【図 7】

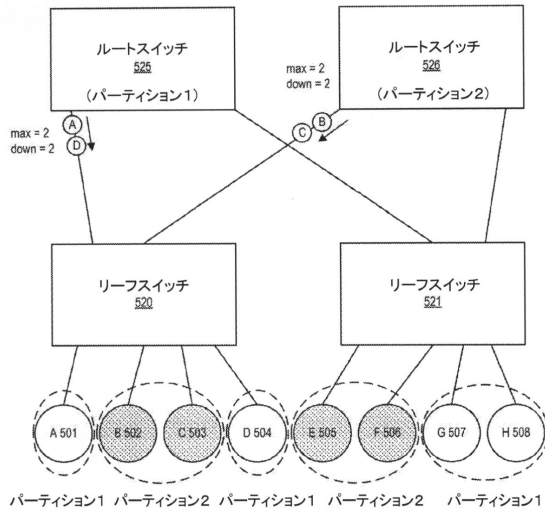


FIGURE 7

【図 8】

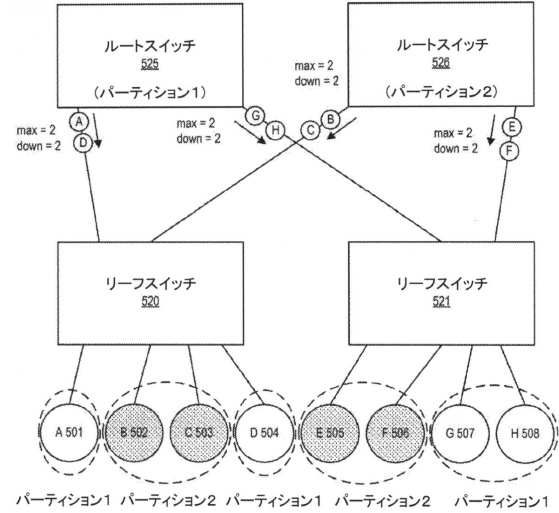


FIGURE 8

【図 9】

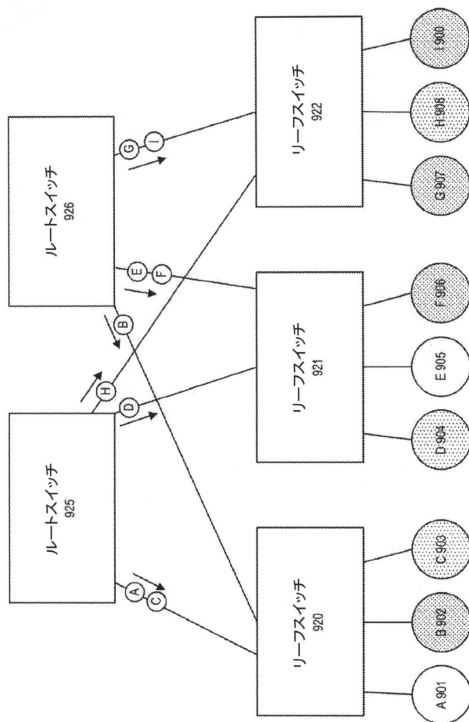


FIGURE 9

【図 10】

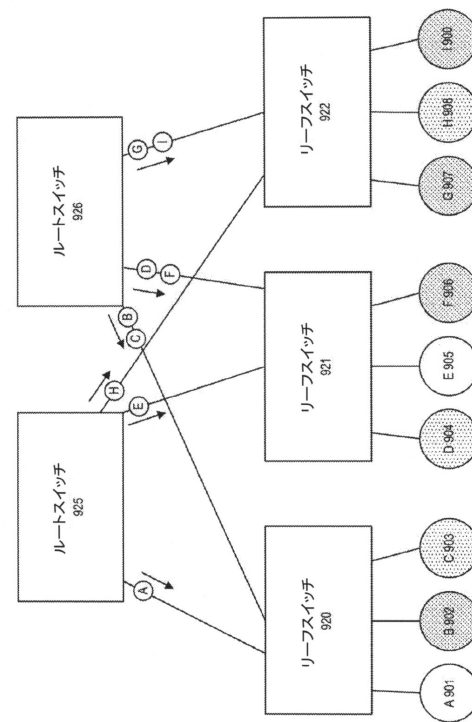


FIGURE 10

【図 1 1】

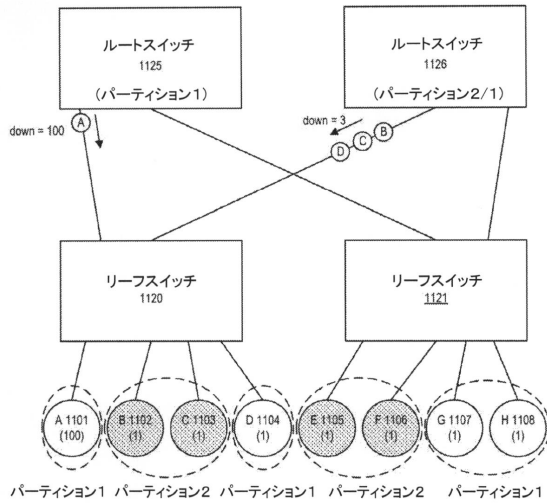


FIGURE 11

【図 1 2】

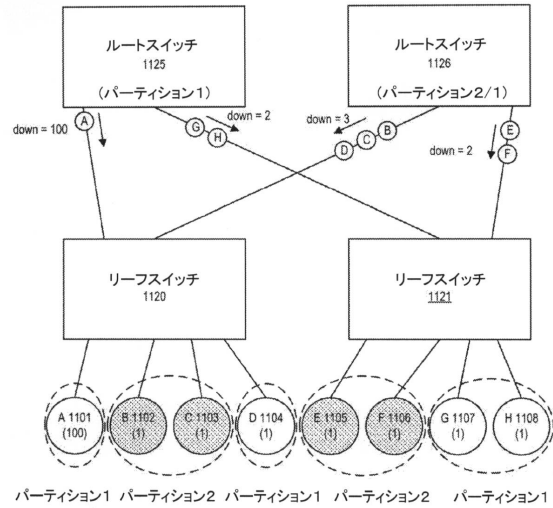


FIGURE 12

【図 1 3】

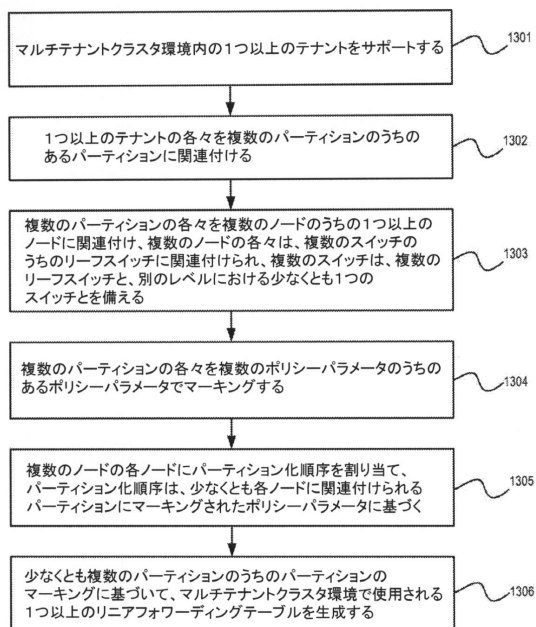


FIGURE 13

【図 1 4】

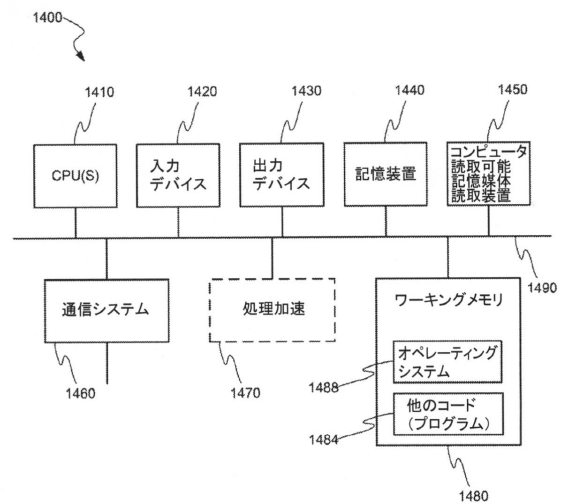


FIGURE 14

フロントページの続き

(31)優先権主張番号 15/182,383

(32)優先日 平成28年6月14日(2016.6.14)

(33)優先権主張国・地域又は機関
米国(US)

(31)優先権主張番号 15/182,397

(32)優先日 平成28年6月14日(2016.6.14)

(33)優先権主張国・地域又は機関
米国(US)

(72)発明者 ボグダンスキー、バルトシュ

ノルウェー、エヌ - 0 2 7 5 オスロ、ホフ・テラッセ、15、エイチ・0 2 0 3

(72)発明者 ヨンセン、ビョルン・ダグ

ノルウェー、エヌ - 0 6 8 7 オスロ、ビルベルクグレンダ、9

審査官 安藤 一道

(56)参考文献 国際公開第2015/020665(WO, A1)

国際公開第2013/071330(WO, A1)

国際公開第2015/130372(WO, A1)

国際公開第2015/031371(WO, A1)

米国特許出願公開第2013/0121149(US, A1)

国際公開第2016/069927(WO, A1)

(58)調査した分野(Int.Cl., DB名)

H04L 12/715

G06F 9/50

H04L 12/70