

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6843112号
(P6843112)

(45) 発行日 令和3年3月17日 (2021.3.17)

(24) 登録日 令和3年2月25日 (2021.2.25)

(51) Int. Cl.

F I

HO 4 L 12/70 (2013.01)
 GO 6 F 9/455 (2006.01)
 GO 6 F 9/50 (2006.01)
 HO 4 L 12/46 (2006.01)

HO 4 L 12/70 D
 GO 6 F 9/455 1 5 O
 GO 6 F 9/50 1 2 O Z
 HO 4 L 12/46 V

請求項の数 9 (全 37 頁)

(21) 出願番号 特願2018-501257 (P2018-501257)
 (86) (22) 出願日 平成28年11月18日 (2016.11.18)
 (65) 公表番号 特表2018-536302 (P2018-536302A)
 (43) 公表日 平成30年12月6日 (2018.12.6)
 (86) 国際出願番号 PCT/US2016/062882
 (87) 国際公開番号 WO2017/091475
 (87) 国際公開日 平成29年6月1日 (2017.6.1)
 審査請求日 令和1年7月19日 (2019.7.19)
 (31) 優先権主張番号 62/259,321
 (32) 優先日 平成27年11月24日 (2015.11.24)
 (33) 優先権主張国・地域又は機関
 米国 (US)
 (31) 優先権主張番号 62/259,831
 (32) 優先日 平成27年11月25日 (2015.11.25)
 (33) 優先権主張国・地域又は機関
 米国 (US)

(73) 特許権者 502303739
 オラクル・インターナショナル・コーポレ
 イション
 アメリカ合衆国カリフォルニア州9406
 5レッドウッド・シティー, オラクル・パ
 ークウェイ500
 (74) 代理人 110001195
 特許業務法人深見特許事務所
 (72) 発明者 タソウラス, エバンジェロス
 ノルウェー、1325 リュサケール、ピ
 イ・オウ・ボックス・134
 (72) 発明者 ザヒド, フェロツ
 ノルウェー、1325 リュサケール、ピ
 イ・オウ・ボックス・134

最終頁に続く

(54) 【発明の名称】 無損失ネットワークにおける効率的な仮想化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】

無損失相互接続ネットワークにおける効率的な仮想化をサポートするためのシステムであって、

1つ以上のマイクロプロセッサと、

1つ以上のスイッチと、

複数のホストチャネルアダプタとを含み、前記複数のホストチャネルアダプタの各々は少なくとも1つの仮想機能、少なくとも1つの仮想スイッチおよび少なくとも1つの物理機能を含み、前記複数のホストチャネルアダプタは前記1つ以上のスイッチを介して相互接続されており、前記システムはさらに、

複数のハイパーバイザを含み、前記複数のハイパーバイザの各々は、前記複数のホストチャネルアダプタのうち少なくとも1つのホストチャネルアダプタに関連付けられており、前記システムはさらに、

複数の仮想マシンを含み、前記複数の仮想マシンの各々は、少なくとも1つの仮想機能に関連付けられており、

前記複数のホストチャネルアダプタは、予めボピュレートされたローカル識別子 (L I D) アーキテクチャを備えた仮想スイッチまたは動的 L I D 割当てアーキテクチャを備えた仮想スイッチのうち1つ以上と共に配置されており、

L I D スペースは、複数の物理的 L I D (p L I D) および複数の仮想 L I D (v L I D) を含み、

10

20

前記仮想スイッチの各々には、前記複数の p L I D のうち 1 つの p L I D が割当てられており、前記割当てられた 1 つの p L I D は関連付けられた物理機能の p L I D に対応しており、

前記複数の仮想マシンの各々には前記複数の v L I D のうち 1 つの v L I D が割当てられており、

1 つ以上のリニアフォーディングテーブル (L F T) は、前記仮想スイッチの各々に割当てられた前記複数の p L I D に少なくとも基づいて計算され、前記複数の p L I D に少なくとも基づいている前記 1 つ以上の L F T の各々は、前記 1 つ以上のスイッチのうちの一のスイッチに関連付けられており、

1 つ以上の二次的 L F T は、前記複数の仮想マシンの各々に割当てられた前記複数の v L I D に少なくとも基づいて計算され、前記複数の v L I D に少なくとも基づいている前記二次的 L F T の各々は、前記 1 つ以上のスイッチのうちの一のスイッチに関連付けられている、システム。

【請求項 2】

無損失相互接続ネットワークにおける効率的な仮想化をサポートするための方法であって、

1 つ以上のマイクロプロセッサを含む 1 つ以上のコンピュータにおいて、1 つ以上のスイッチと、複数のホストチャネルアダプタと、複数のハイパーバイザと、複数の仮想マシンとを設けるステップを含み、

前記複数のホストチャネルアダプタの各々は少なくとも 1 つの仮想機能、少なくとも 1 つの仮想スイッチおよび少なくとも 1 つの物理機能を含み、前記複数のホストチャネルアダプタは前記 1 つ以上のスイッチを介して相互接続されており、

前記複数のハイパーバイザの各々は、前記複数のホストチャネルアダプタのうち少なくとも 1 つのホストチャネルアダプタに関連付けられており、

前記複数の仮想マシンの各々は、少なくとも 1 つの仮想機能に関連付けられており、前記方法はさらに、

予めポピュレートされたローカル識別子 (L I D) アーキテクチャを備えた仮想スイッチまたは動的 L I D 割当てアーキテクチャを備えた仮想スイッチのうち 1 つ以上を備えた前記複数のホストチャネルアダプタを配置するステップと、

前記仮想スイッチの各々に複数の物理的 L I D (p L I D) のうち 1 つの p L I D を割当てるステップとを含み、前記割当てられた p L I D は関連付けられた物理機能の p L I D に対応しており、前記方法はさらに、

前記複数の仮想マシンの各々に複数の仮想 L I D (v L I D) のうち 1 つの v L I D を割当てるステップを含み、

L I D スペースは前記複数の p L I D および前記複数の v L I D を含み、前記方法は、さらに、

前記仮想スイッチの各々に割当てられた前記複数の p L I D に少なくとも基づいて 1 つ以上のリニアフォーディングテーブル (L F T) を計算するステップをさらに含み、前記複数の p L I D に少なくとも基づいている前記 1 つ以上の L F T の各々は、前記 1 つ以上のスイッチのうちの一のスイッチに関連付けられており、前記方法は、さらに、

前記複数の仮想マシンの各々に割当てられた前記複数の v L I D に少なくとも基づいて、1 つ以上の二次的 L F T を計算するステップをさらに含み、前記複数の v L I D に少なくとも基づいている前記二次的 L F T の各々は、前記 1 つ以上のスイッチのうちの一のスイッチに関連付けられている、方法。

【請求項 3】

前記二次的 L F T の各々は前記複数の仮想マシンにトラフィックを転送する、請求項 2 に記載の方法。

【請求項 4】

仮想マシンがマイグレートされると、前記二次的 L F T のうち 1 つ以上が更新される、請求項 2 または 3 に記載の方法。

【請求項 5】

前記 1 つ以上のスイッチの各々および前記複数のホストチャネルアダプタの各々には p L I D が割当てられる、請求項 2 ~ 4 のいずれか 1 項に記載の方法。

【請求項 6】

サブネットマネージャは、再構成中に前記無損失相互接続ネットワーク上で経路演算を実行し、前記経路演算は当該経路演算から各々の v L I D を除外している、請求項 2 ~ 5 のいずれか 1 項に記載の方法。

【請求項 7】

前記無損失相互接続ネットワークはインフィニバンドファブリックを含み、各々の p L I D 値は、インフィニバンドパケットのローカルルートヘッダにおける標準 S L I D フィールドおよび標準 D L I D フィールドを用いて表わされ、各々の v L I D 値は、拡張を表わす追加の 2 ビット以上と組合わせて、前記標準 S L I D フィールドと前記標準 D L I D フィールドとの組合せを用いて表わされる、請求項 2 ~ 6 のいずれか 1 項に記載の方法。

10

【請求項 8】

前記 1 つ以上のスイッチはリーフスイッチを含む、請求項 2 ~ 7 のいずれか 1 項に記載の方法。

【請求項 9】

無損失相互接続ネットワークにおける効率的な仮想化をサポートするための命令が格納されているコンピュータ読取り可能プログラムであって、前記命令が 1 つ以上のコンピュータによって読出されて実行されると、前記 1 つ以上のコンピュータに請求項 2 ~ 8 のいずれか 1 項に記載の方法を実行させる、コンピュータ読取り可能プログラム。

20

【発明の詳細な説明】**【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資料を含む。この特許文献または特許開示は特許商標庁の特許ファイルまたは記録に記載されているため、著作権所有者は、何人によるその複製複製に対しても異議はないが、その他の場合には如何なるときもすべての著作権を保有する。

【0002】

30

発明の分野：

本発明は、概して、コンピュータシステムに関し、特に、S R - I O V v S w i t c h アーキテクチャを用いてコンピュータシステム仮想化およびライブマイグレーションをサポートすることに関する。

【背景技術】**【0003】**

背景：

導入されるクラウドコンピューティングアーキテクチャがより大規模になるのに応じて、従来のネットワークおよびストレージに関する性能および管理の障害が深刻な問題になってきている。クラウドコンピューティングファブリックのための基礎としてインフィニバンド（登録商標）（InfiniBand：I B）技術などの高性能な無損失相互接続を用いることへの関心がますます高まってきている。これは、本発明の実施形態が対応するように意図された一般領域である。

40

【発明の概要】**【課題を解決するための手段】****【0004】**

概要：

サブネットにおいて仮想マシンマイグレーションをサポートするためのシステムおよび方法がこの明細書中に記載される。一実施形態に従うと、方法は、1 つ以上のマイクロプロセッサを含む 1 つ以上のコンピュータにおいて、1 つ以上のスイッチを設けることがで

50

き、当該1つ以上のスイッチは少なくともリーフスイッチを含み、当該1つ以上のスイッチの各々は複数のポートを含み、当該方法はさらに、複数のホストチャネルアダプタを設けることができる。複数のホストチャネルアダプタの各々は、少なくとも1つの仮想機能、少なくとも1つの仮想スイッチおよび少なくとも1つの物理機能を含む。複数のホストチャネルアダプタは当該1つ以上のスイッチを介して相互接続されている。当該方法はさらに、複数のハイパーバイザを設けることができる。当該複数のハイパーバイザの各々は当該複数のホストチャネルアダプタのうち少なくとも1つのホストチャネルアダプタに関連付けられている。当該方法はさらに、複数の仮想マシンを設けることができる。複数の仮想マシンの各々は、少なくとも1つの仮想機能に関連付けられている。当該方法はさらに、予めボピュレートされたローカル識別子 (local identifier: L I D) アーキテクチャを備えた仮想スイッチ、または動的 L I D 割当てアーキテクチャを備えた仮想スイッチのうち1つ以上を備えた複数のホストチャネルアダプタを配置することができる。当該方法は、各々の仮想スイッチを L I D に割当てることができ、割当てられた L I D は関連付けられた物理機能の L I D に対応している。当該方法は、仮想スイッチの各々に割当てられた L I D に少なくとも基づいて、1つ以上のリニアフォワーディングテーブル (linear forwarding table: L F T) を計算することができる。1つ以上の L F T の各々は、1つ以上のスイッチのうちの一のスイッチに関連付けられている。

【0005】

一実施形態に従うと、方法は、1つ以上のマイクロプロセッサを含む1つ以上のコンピュータにおいて、1つ以上のマイクロプロセッサと、少なくともリーフスイッチを含む1つ以上のスイッチとを設けることができ、当該1つ以上のスイッチの各々は複数のポートを含み、さらに、複数のホストチャネルアダプタを設けることができ、ホストチャネルアダプタの各々は少なくとも1つの仮想機能、少なくとも1つの仮想スイッチおよび少なくとも1つの物理機能を含み、複数のホストチャネルアダプタは1つ以上のスイッチを介して相互接続されており、さらに、複数のハイパーバイザを設けることができ、複数のハイパーバイザの各々は、複数のホストチャネルアダプタのうち少なくとも1つのホストチャネルアダプタに関連付けられており、さらに、複数の仮想マシンを設けることができ、複数の仮想マシンの各々は少なくとも1つの仮想機能に関連付けられている。当該方法は、予めボピュレートされたローカル識別子 (L I D) アーキテクチャを備えた仮想スイッチまたは動的 L I D 割当てアーキテクチャを備えた仮想スイッチのうち1つ以上を備えた複数のホストチャネルアダプタを配置することができる。当該方法は、仮想スイッチの各々に複数の物理的 L I D (physical L I D: p L I D) のうち1つの p L I D を割当てることができ、割当てられた p L I D は関連付けられた物理機能の p L I D に対応している。当該方法はまた、複数の仮想マシンの各々に複数の仮想 L I D (virtual L I D: v L I D) のうち1つの v L I D を割当てることができ、L I D スペースは複数の p L I D および複数の v L I D を含んでいる。

【0006】

一実施形態に従うと、各々の p L I D 値は、インフィニバンドパケットのローカルルートヘッダにおける標準 S L I D フィールドおよび標準 D L I D フィールドを用いて表わすことができる。同様に、各々の v L I D 値は、拡張を表わす追加の2ビット以上と組合わせて、標準 S L I D フィールドと標準 D L I D フィールドとの組合せを用いて表わすことができる。

【図面の簡単な説明】

【0007】

【図1】一実施形態に従ったインフィニバンド環境の一例を示す図である。

【図2】一実施形態に従った、ネットワーク環境におけるツリートポロジの一例を示す図である。

【図3】一実施形態に従った例示的な共有ポートアーキテクチャを示す図である。

【図4】一実施形態に従った例示的な v S w i t c h アーキテクチャを示す図である。

【図5】一実施形態に従った例示的な v P o r t アーキテクチャを示す図である。

【図 6】一実施形態に従った、L I D が予めポピュレートされた例示的な v S w i t c h アーキテクチャを示す図である。

【図 7】一実施形態に従った、動的 L I D 割当てがなされた例示的な v S w i t c h アーキテクチャを示す図である。

【図 8】一実施形態に従った、動的 L I D 割当てがなされかつ L I D が予めポピュレートされている v S w i t c h を備えた例示的な v S w i t c h アーキテクチャを示す図である。

【図 9】一実施形態に従った、拡張されたローカルルートヘッダを示す図である。

【図 10】一実施形態に従った、2つの例示的なリニアフォワーディングテーブルを示す図である。

10

【図 11】一実施形態に従った、無損失相互接続ネットワークにおける効率的な仮想化のサポートの例を示す図である。

【図 12】一実施形態に従った、無損失相互接続ネットワークにおける効率的な仮想化のサポートの例を示す図である。

【図 13】一実施形態に従った、無損失相互接続ネットワークにおける効率的な仮想化のサポートの例を示す図である。

【図 14】一実施形態に従った、無損失相互接続ネットワークにおける効率的な仮想化のサポートの例を示す図である。

【図 15】一実施形態に従った潜在的な仮想マシンマイグレーションを示す図である。

【図 16】一実施形態に従ったスイッチタプルを示す図である。

20

【図 17】一実施形態に従った再構成プロセスを示す図である。

【図 18】一実施形態に従った、無損失相互接続ネットワークにおける効率的な仮想化をサポートするための方法を示すフローチャートである。

【図 19】一実施形態に従った、無損失相互接続ネットワークにおける効率的な仮想化をサポートするための方法を示すフローチャートである。

【発明を実施するための形態】

【0008】

詳細な説明：

本発明は、同様の参照番号が同様の要素を指している添付図面の図において、限定のためではなく例示のために説明されている。なお、この開示における「ある」または「1つの」または「いくつかの」実施形態への参照は必ずしも同じ実施形態に対するものではなく、そのような参照は少なくとも1つを意味する。特定の実現例が説明されるが、これらの特定の実現例が例示的な目的のためにのみ提供されることが理解される。当業者であれば、他の構成要素および構成が、この発明の範囲および精神から逸脱することなく使用され得ることを認識するであろう。

30

【0009】

図面および詳細な説明全体にわたって同様の要素を示すために、共通の参照番号が使用され得る。したがって、ある図で使用される参照番号は、要素が別のところで説明される場合、そのような図に特有の詳細な説明において参照される場合もあり、または参照されない場合もある。

40

【0010】

無損失相互接続ネットワークにおける効率的な仮想化をサポートするためのシステムおよび方法がこの明細書中に記載される。

【0011】

この発明の以下の説明は、高性能ネットワークについての一例として、インフィニバンド (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何ら限定されることなく使用され得ることが、当業者には明らかであるだろう。以下の説明ではまた、ファブリックトポロジーについての一例として、ファットツリートポロジーを使用する。他のタイプのファブリックトポロジーが何ら限定されることなく使用され得ることが当業者には明らかであるだろう。

50

【 0 0 1 2 】

現代（たとえばExascale（エクサスケール）時代）におけるクラウドの要求を満たすために、仮想マシンがリモート・ダイレクト・メモリ・アクセス（Remote Direct Memory Access：R D M A）などの低オーバーヘッドネットワーク通信パラダイムを利用できることが望ましい。R D M AはO Sスタックをバイパスし、ハードウェアと直接通信することで、シングルルートI / O仮想化（Single-Root I/O Virtualization：S R - I O V）ネットワークアダプタのようなパススルー技術が使用可能となる。一実施形態に従うと、高性能な無損失相互接続ネットワークにおける適用可能性のために、仮想スイッチ（virtual switch：v S w i t c h）S R - I O Vアーキテクチャを提供することができる。ライブマイグレーションを実際に選択できるようにするためにネットワーク再構成時間が重要となるので、ネットワークアーキテクチャに加えて、スケーラブルであるとともにトポロ

10

【 0 0 1 3 】

一実施形態に従うと、さらには、v S w i t c hを用いる仮想化された環境のためのルーティング戦略を提供することができ、ネットワークトポロジー（たとえばファットツリートポロジー）のための効率的なルーティングアルゴリズムを提供することができる。動的な再構成メカニズムは、ファットツリーにおいて課されるオーバーヘッドを最小限にするためにさらに調整することができる。

【 0 0 1 4 】

本発明の一実施形態に従うと、仮想化は、クラウドコンピューティングにおける効率的なリソース利用および融通性のあるリソース割当てに有益であり得る。ライブマイグレーションは、アプリケーションにトランスペアレントな態様で物理サーバ間で仮想マシン（virtual machine：V M）を移動させることによってリソース使用を最適化することを可能にする。このため、仮想化は、ライブマイグレーションによる統合、リソースのオン・デマンド・プロビジョニングおよび融通性を可能にし得る。

20

【 0 0 1 5 】

インフィニバンド（登録商標）

インフィニバンド（I B）は、インフィニバンド・トレード・アソシエーション（InfiniBand™ Trade Association）によって開発されたオープン標準無損失ネットワーク技術である。この技術は、特に高性能コンピューティング（high-performance computing：H P C）アプリケーションおよびデータセンタを対象とする、高スループットおよび少ない待ち時間の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（serial point-to-point full-duplex interconnect）に基づいている。

30

【 0 0 1 6 】

インフィニバンド・アーキテクチャ（InfiniBand Architecture：I B A）は、2層トポロジー分割をサポートする。低層では、I Bネットワークはサブネットと呼ばれ、1つのサブネットは、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る一組のホストを含み得る。より高いレベルでは、1つのI B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 7 】

1つのサブネット内で、ホストは、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接続され得る。加えて、サブネットにおける指定されたデバイス上に存在する、1つのマスター管理エンティティ、すなわちサブネットマネージャ（subnet manager：S M）があり得る。サブネットマネージャは、I Bサブネットを構成し、起動し、維持する役割を果たす。加えて、サブネットマネージャ（S M）は、I Bファブリックにおいてルーティングテーブル計算を行なう役割を果たし得る。ここで、たとえば、I Bネットワークのルーティングは、ローカルサブネットにおけるすべての送信元と宛先とのペア間の適正な負荷バランスングを目標とする。

40

【 0 0 1 8 】

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

50

パケット (subnet management packet : S M P) と呼ば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) と交換する。サブネット管理エージェントは、すべての I B サブネットデバイス上に存在する。S M P を使用することにより、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 から通知を受信することができる。

【 0 0 1 9 】

一実施形態によれば、I B ネットワークにおけるサブネット内のルーティングは、スイッチに格納された L F T に基づき得る。L F T は、使用中の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。サブネットでは、エンドノード上のホストチャネルアダプタ (Host Channel Adapter : H C A) ポートおよびスイッチが、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 における各エントリは、宛先 L I D (destination L I D : D L I D) と出力ポートとからなる。テーブルにおける L I D ごとに 1 つのエントリのみがサポートされる。パケットがあるスイッチに到着すると、その出力ポートは、そのスイッチのフォワーディ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判断される。所与の送信元 - 宛先ペア (L I D ペア) 間のネットワークにおいてパケットは同じ経路を通るため、ルーティングは決定論的である。

【 0 0 2 0 】

一般に、マスターサブネットマネージャを除く他のすべてのサブネットマネージャは、耐故障性のために待機モードで作動する。しかしながら、マスターサブネットマネージャが故障した状況では、待機中のサブネットマネージャによって、新しいマスターサブネットマネージャが取り決められる。マスターサブネットマネージャはまた、サブネットの周期的なスイープ (sweep) を行なってあらゆるトポロジ変化を検出し、それに応じてネットワークを再構成する。

【 0 0 2 1 】

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 (L I D) を用いてアドレス指定され得るとともに、単一のサブネットは 4 9 1 5 1 個のユニキャスト L I D に制限され得る。サブネット内で有効なローカルアドレスである L I D の他に、各 I B デバイスは、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 は、I B レイヤー 3 (L 3) アドレスであるグローバル識別子 (global identifier : G I D) を形成するために使用され得る。

【 0 0 2 2 】

S M は、ネットワーク初期化時間に、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 / ルート) を計算し得る。さらに、トポロジが変化するたびに、ルーティングテーブルは、接続性および最適性能を確実にするために更新され得る。通常動作中、S M は、トポロジ変化をチェックするためにネットワークの周期的なライトスイープ (light sweep) を実行し得る。ライトスイープ中に変化が発見された場合、または、ネットワーク変化を信号で伝えるメッセージ (トラップ) を S M が受信した場合、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。

【 0 0 2 3 】

たとえば、S M は、リンクがダウンした場合、デバイスが追加された場合、またはリンクが除去された場合など、ネットワークトポロジが変化する場合に、ネットワークを再構成し得る。再構成ステップは、ネットワーク初期化中に行なわ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制限されるローカルスコープを有し得る。また、ルータを用いる大規模ファブリックのセグメント化は、再構成スコープを制限し得る。

【 0 0 2 4 】

一実施形態によれば、I B ネットワークは、ネットワークファブリックを共有するシステムの論理グループの分離をもたらすためにセキュリティメカニズムとしてパーティショニングをサポートし得る。ファブリックにおけるノード上の各 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

10

20

30

40

50

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、各ポートに関する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(partition key: P__キー)のテーブルとして構成することができる。S Mはまた、これらのポートを介してデータトラフィックを送信または受信するエンドノードに関連付けられたP__Key情報を含むパーティション実施テーブルを用いて、スイッチポートおよびルータポートを構成することができる。加えて、一般的な場合には、スイッチポートのパーティションメンバーシップは、(リンクに向かう)出口方向に向かってポートを介してルーティングされたL I Dに間接的に関連付けられたすべてのメンバーシップの集合を表わし得る。

【0025】

一実施形態によれば、ノード間の通信のために、管理キューペア(Q P 0およびQ P 1)を除き、キューペア(Queue Pair: Q P)およびエンドツーエンドコンテキスト(End-to-End context: E E C)を特定のパーティションに割り当てることができる。次に、P__キー情報を、送信されたすべてのI Bトランスポートパケットに追加することができる。パケットがH C Aポートまたはスイッチに到着すると、そのP__キー値を、S Mによって構成されたテーブルに対して確認することができる。無効のP__キー値が見つかった場合、そのパケットは直ちに廃棄される。このように、通信は、パーティションを共有するポート間でのみ許可される。

【0026】

一実施形態に従ったインフィニバンド環境100の例を示す図1に、インフィニバンドファブリックの一例を示す。図1に示す例では、ノードA 101~E 105は、インフィニバンドファブリック120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。一実施形態に従うと、さまざまなノード(たとえばノードA 101~E 105)はさまざまな物理デバイスによって表わすことができる。一実施形態に従うと、さまざまなノード(たとえばノードA 101~E 105)は仮想マシンなどのさまざまな仮想デバイスによって表わすことができる。

【0027】

インフィニバンドにおける仮想マシン

過去10年の間に、ハードウェア仮想化サポートによってC P Uオーバーヘッドが実質的に排除され、メモリ管理ユニットを仮想化することによってメモリオーバーヘッドが著しく削減され、高速S A Nストレージまたは分散型ネットワークファイルシステムの利用によってストレージオーバーヘッドが削減され、シングルルートI / O仮想化(Single Root Input/Output Virtualization: S R - I O V)のようなデバイス・パススルー技術を使用することによってネットワークI / Oオーバーヘッドが削減されてきたことに応じて、仮想化された高性能コンピューティング(High Performance Computing: H P C)環境の将来見通しが大幅に改善されてきた。現在では、クラウドが、高性能相互接続ソリューションを用いて仮想H P C(virtual H P C: v H P C)クラスタに対応し、必要な性能を提供することができる。

【0028】

しかしながら、インフィニバンド(I B)などの無損失ネットワークと連結されたとき、仮想マシン(V M)のライブマイグレーションなどのいくつかのクラウド機能は、これらのソリューションにおいて用いられる複雑なアドレス指定およびルーティングスキームのせいで、依然として問題となる。I Bは、高帯域および低レイテンシを提供する相互接続ネットワーク技術であり、このため、H P Cおよび他の通信集約型の作業負荷に非常によく適している。

【0029】

I BデバイスをV Mに接続するための従来のアプローチは直接割り当てされたS R - I O Vを利用することによるものである。しかしながら、S R - I O Vを用いてI Bホストチャネルアダプタ(H C A)に割り当てられたV Mのライブマイグレーションを実現することは難易度の高いものであることが判明した。各々のI Bが接続されているノードは、3つの異なるアドレス(すなわちL I D、G U I DおよびG I D)を有する。ライブマイグレ

ーションが発生すると、これらのアドレスのうち1つ以上が変化する。マイグレーション中のVM (VM-in-migration) と通信する他のノードは接続性を失う可能性がある。これが発生すると、IBサブネットマネージャ (Subnet Manager : SM) にサブネット管理 (Subnet Administration : SA) 経路記録クエリを送信することによって、再接続すべき仮想マシンの新しいアドレスを突きとめることにより、失われた接続を回復させるように試みることができる。

【0030】

IBは3つの異なるタイプのアドレスを用いる。第1のタイプのアドレスは16ビットのローカル識別子 (LID) である。少なくとも1つの固有のLIDは、SMによって各々のHCAポートおよび各々のスイッチに割当てられる。LIDはサブネット内のトラフィックをルーティングするために用いられる。LIDが16ビット長であるので、65536個の固有のアドレス組合せを構成することができ、そのうち49151個 (0x0001 - 0xBFFF) だけをユニキャストアドレスとして用いることができる。結果として、入手可能なユニキャストアドレスの数は、IBサブネットの最大サイズを定義することとなる。第2のタイプのアドレスは、製造業者によって各々のデバイス (たとえば、HCAおよびスイッチ) ならびに各々のHCAポートに割当てられた64ビットのグローバル意識別子 (GUID) である。SMは、HCAポートに追加のサブネット固有GUIDを割当ててもよく、これは、SR-IOVが用いられる場合に有用となる。第3のタイプのアドレスは128ビットのグローバル識別子 (GID) である。GIDは有効なIPv6ユニキャストアドレスであり、少なくとも1つが各々のHCAポートに割当てられている。GIDは、ファブリックアドミニストレータによって割当てられたグローバルに固有の64ビットプレフィックスと各々のHCAポートのGUIDアドレスとを組み合わせることによって形成される。

【0031】

ファットツリー (Fat Tree : F T r e e) トポロジーおよびルーティング

一実施形態によれば、IBベースのHPCシステムのいくつかは、ファットツリートポロジーを採用して、ファットツリーが提供する有用な特性を利用する。これらの特性は、各送信元宛先ペア間の複数経路の利用可能性に起因する、フルバイセクション帯域幅および固有の耐故障性を含む。ファットツリーの背後にある初期の概念は、ツリーがトポロジーのルート (root) に近づくにつれて、より利用可能な帯域幅を用いて、ノード間のより太いリンクを採用することであった。より太いリンクは、上位レベルのスイッチにおける輻輳を回避するのに役立てることができ、バイセクション帯域幅が維持される。

【0032】

図2は、一実施形態に従った、ネットワーク環境におけるツリートポロジーの例を示す。図2に示すように、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が接続され得る。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と複数のスパインスイッチまたはルート (root) スwitch 231~234とを含むファットツリートポロジーに基づき得る。加えて、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。

【0033】

また、図2に示すように、エンドノード201~204の各々は、マルチホームノード、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のうち2つ以上の部分に接続される単一のノードであり得る。たとえば、ノード201はポートH1およびH2を含み、ノード202はポートH3およびH4を含み、ノード203はポートH5およびH6を含み、ノード204はポートH7およびH8を含み得る。

【0034】

加えて、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、ルートスイッチ232はスイッチポート3~4を有し、ルートスイッチ233はスイッチポート5~6を有し、ルートスイッチ234はスイッチポート7~8を有し得る。

【0035】

ー実施形態によれば、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ーに関して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また、O F E D (Open Fabric Enterprise Distribution: I B ベースのアプリケーションを構築しデプロイするための標準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において実現される。

【0036】

ファットツリールーティングメカニズムの目的は、ネットワークファブリックにおけるリンクにわたって最短経路ルートを均一に広げるL F Tを生成することである。このメカニズムは、索引付け順序でファブリックを横断し、エンドノードの目標L I D、ひいては 10
対応するルートを各スイッチポートに割当てて。同じリーフスイッチに接続されたエンドノードについては、索引付け順序は、エンドノードが接続されるスイッチポートに依存し得る(すなわち、ポートナンバリングシーケンス)。各ポートについては、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、新しいルートが追加されるたびに、ポート使用カウンタを使用して使用頻度が最小のポートを選択することができる。

【0037】

ー実施形態に従うと、パーティショニングされたサブネットでは、共通のパーティションのメンバではないノードは通信することを許可されない。実際には、これは、ファットツリールーティングアルゴリズムによって割当てられたルートのうちのいくつかがユーザ 20
トラフィックのため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、それらのルートについてのL F Tを、他の機能的経路と同じやり方で生成する場合、問題が生じる。この動作は、リンク上でバランシングを劣化させるおそれがある。なぜなら、ノードが索引付けの順序でルーティングされているからである。パーティションに気づかずにルーティングが行なわれるため、ファットツリーでルーティングされたサブネットにより、概して、パーティション間の分離が不良なものとなる。

【0038】

ー実施形態に従うと、ファットツリーは、利用可能なネットワークリソースでスケール 30
リングすることができる階層ネットワークトポロジーである。さらに、ファットツリーは、さまざまなレベルの階層に配置された商品スイッチを用いて容易に構築される。さらに、k - a r y - n - t r e e、拡張された一般化ファットツリー(Extended Generalized Fat-Tree: X G F T)、パラレルポート一般化ファットツリー(Parallel Ports Generalized Fat-Tree: P G F T)およびリアルライフファットツリー(Real Life Fat-Tree: R L F T)を含むファットツリーのさまざまな変形例が、一般に利用可能である。

【0039】

また、k - a r y - n - t r e eは、nレベルのファットツリーであって、 k^n エンドノードと、 $n \cdot k^{n-1}$ スイッチとを備え、各々が2kポートを備えている。各々のスイッチは、ツリーにおいて上下方向に同数の接続を有している。X G F Tファットツリーは、スイッチのための異なる数の上下方向の接続と、ツリーにおける各レベルでの異なる数の接続とをともに可能にすることによって、k - a r y - n - t r e eを拡張させる。P 40
G F T定義はさらに、X G F Tトポロジーを拡張して、スイッチ間の複数の接続を可能にする。多種多様なトポロジーはX G F TおよびP G F Tを用いて定義することができる。しかしながら、実用化するために、現代のH P Cクラスタにおいて一般に見出されるファットツリーを定義するために、P G F Tの制限バージョンであるR L F Tが導入されている。R L F Tは、ファットツリーにおけるすべてのレベルに同じポートカウントスイッチを用いている。

【0040】

入出力(Input/Output: I / O) 仮想化

ー実施形態に従うと、I / O仮想化(I/O Virtualization: I O V)は、基礎をなす物理リソースに仮想マシン(V M)がアクセスすることを可能にすることによって、I / O 50
を利用可能にすることができる。ストレージトラフィックとサーバ間通信とを組合せると

、シングルサーバの I / O リソースにとって抗し難い高い負荷が課され、結果として、データの待機中に、バックログが発生し、プロセッサがアイドル状態になる可能性がある。I / O 要求の数が増えるにつれて、I O V により利用可能性をもたらすことができ、最新の CPU 仮想化において見られる性能レベルに匹敵するように、(仮想化された) I / O リソースの性能、スケーラビリティおよび融通性を向上させることができる。

【0041】

一実施形態に従うと、I / O リソースの共有を可能にして、VM からリソースへのアクセスが保護されることを可能にし得るような I O V が所望される。I O V は、VM にエクスポートされる論理装置を、その物理的な実装から分離する。現在、エミュレーション、準仮想化、直接的な割当て (direct assignment : DA)、およびシングルルート I / O 仮想化 (SR - I O V) などのさまざまなタイプの I O V 技術が存在し得る。

10

【0042】

一実施形態に従うと、あるタイプの I O V 技術としてソフトウェアエミュレーションがある。ソフトウェアエミュレーションは分離されたフロントエンド / バックエンド・ソフトウェアアーキテクチャを可能にし得る。フロントエンドは VM に配置されたデバイスドライバであり得、I / O アクセスをもたらすためにハイパーバイザによって実現されるバックエンドと通信し得る。物理デバイス共有比率は高く、VM のライブマイグレーションはネットワークダウンタイムのわずか数ミリ秒で実現可能である。しかしながら、ソフトウェアエミュレーションはさらなる不所望な計算上のオーバーヘッドをもたらしてしまう。

20

【0043】

一実施形態に従うと、別のタイプの I O V 技術として直接的なデバイスの割当てがある。直接的なデバイスの割当てでは、I / O デバイスを VM に連結する必要があるが、デバイスは VM 間では共有されない。直接的な割当てまたはデバイス・パススルーは、最小限のオーバーヘッドでほぼ固有の性能を提供する。物理デバイスはハイパーバイザをバイパスし、直接、VM に取付けられている。しかしながら、このような直接的なデバイスの割当ての欠点は、仮想マシン間で共有がなされないため、1 枚の物理ネットワークカードが 1 つの VM と連結されるといったように、スケーラビリティが制限されてしまうことである。

【0044】

30

一実施形態に従うと、シングルルート I O V (Single Root IOV : SR - I O V) は、ハードウェア仮想化によって、物理装置がその同じ装置の複数の独立した軽量のインスタンスとして現われることを可能にし得る。これらのインスタンスは、パススルー装置として VM に割当てることができ、仮想機能 (Virtual Function : VF) としてアクセスすることができる。ハイパーバイザは、(1 つのデバイスごとに) 固有の、十分な機能を有する物理機能 (Physical Function : PF) によってデバイスにアクセスする。SR - I O V は、純粹に直接的に割当てする際のスケーラビリティの問題を軽減する。しかしながら、SR - I O V によって提示される問題は、それが VM マイグレーションを損なう可能性があることである。これらの I O V 技術の中でも、SR - I O V は、ほぼ固有の性能を維持しながらも、複数の VM から単一の物理デバイスに直接アクセスすることを可能にする手段を用いて PCI Express (PCIe) 規格を拡張することができる。これにより、SR - I O V は優れた性能およびスケーラビリティを提供することができる。

40

【0045】

SR - I O V は、PCIe デバイスが、各々のゲストに 1 つの仮想デバイスを割当てることによって複数のゲスト間で共有することができる複数の仮想デバイスをエクスポートすることを可能にする。各々の SR - I O V デバイスは、少なくとも 1 つの物理機能 (PF) と、1 つ以上の関連付けられた仮想機能 (VF) とを有する。PF は、仮想マシンモニタ (virtual machine monitor : VMM) またはハイパーバイザによって制御される通常の PCIe 機能であるのに対して、VF は軽量の PCIe 機能である。各々の VF はそれ自体のベースアドレス (base address : BAR) を有しており、固有のリクエスト ID

50

が割当てられている。固有のリクエストIDは、I/Oメモリ管理ユニット(I/O memory management unit: IOMMU)がさまざまなVFへのノからのトラフィックストリームを区別することを可能にする。IOMMUはまた、メモリを適用して、PFとVFとの間の変換を中断する。

【0046】

しかし、残念ながら、直接的デバイス割当て技術は、仮想マシンのトランスペアレントなライブマイグレーションがデータセンタ最適化のために所望されるような状況においては、クラウドプロバイダにとって障壁となる。ライブマイグレーションの本質は、VMのメモリ内容がリモートハイパーバイザにコピーされるという点である。さらに、VMがソースハイパーバイザにおいて中断され、VMの動作が宛先において再開される。ソフトウェアエミュレーション方法を用いる場合、ネットワークインターフェイスは、それらの内部状態がメモリに記憶され、さらにコピーされるように仮想的である。このため、ダウンタイムは数ミリ秒にまで減らされ得る。

【0047】

しかしながら、SR-IOVなどの直接的デバイス割当て技術が用いられる場合、マイグレーションはより困難になる。このような状況においては、ネットワークインターフェイスの内部状態全体は、それがハードウェアに結び付けられているのでコピーすることができない。代わりに、VMに割当てられたSR-IOV VFが分離され、ライブマイグレーションが実行されることとなり、新しいVFが宛先において付与されることとなる。インフィニバンドおよびSR-IOVの場合、このプロセスがダウンタイムを数秒のオーダーでもたらず可能性がある。さらに、SR-IOV共有型ポートモデルにおいては、VMのアドレスがマイグレーション後に変化することとなり、これにより、SMにオーバーヘッドが追加され、基礎をなすネットワークファブリックの性能に対して悪影響が及ぼされることとなる。

【0048】

インフィニバンドSR-IOVアーキテクチャ - 共有ポート

さまざまなタイプのSR-IOVモデル(たとえば共有ポートモデル、仮想スイッチモデルおよび仮想ポートモデル)があり得る。

【0049】

図3は、一実施形態に従った例示的な共有ポートアーキテクチャを示す。図に示されるように、ホスト300(たとえばホストチャネルアダプタ)はハイパーバイザ310と対話し得る。ハイパーバイザ310は、さまざまな仮想機能330、340および350をいくつかの仮想マシンに割当て得る。同様に、物理機能はハイパーバイザ310によって処理することができる。

【0050】

一実施形態に従うと、図3に示されるような共有ポートアーキテクチャを用いる場合、ホスト(たとえばHCA)は、物理機能320と仮想機能330、350、350との間において単一の共有LIDおよび共有キュー対(Queue Pair: QP)のスペースがあるネットワークにおいて単一のポートとして現われる。しかしながら、各々の機能(すなわち、物理機能および仮想機能)はそれら自体のGIDを有し得る。

【0051】

図3に示されるように、一実施形態に従うと、さまざまなGIDを仮想機能および物理機能に割当てることができ、特別のキュー対であるQP0およびQP1(すなわちインフィニバンド管理パケットのために用いられる専用のキュー対)が物理機能によって所有される。これらのQPはVFにも同様にエクスポートされるが、VFはQP0を使用することが許可されておらず(VFからQP0に向かって入来するすべてのSMPが廃棄され)、QP1は、PFが所有する実際のQP1のプロキシとして機能し得る。

【0052】

一実施形態に従うと、共有ポートアーキテクチャは、(仮想機能に割当てられることによってネットワークに付随する)VMの数によって制限されることのない高度にスケラ

10

20

30

40

50

ブルなデータセンタを可能にし得る。なぜなら、ネットワークにおける物理的なマシンおよびスイッチによってL I Dスペースが消費されるだけであるからである。

【 0 0 5 3 】

しかしながら、共有ポートアーキテクチャの欠点は、トランスペアレントなライブマイグレーションを提供することができない点であり、これにより、フレキシブルなV M配置についての可能性が妨害されてしまう。各々のL I Dが特定のハイパーバイザに関連付けられており、かつハイパーバイザ上に常駐するすべてのV M間で共有されているので、マイグレートしているV M（すなわち、宛先ハイパーバイザにマイグレートする仮想マシン）は、そのL I Dを宛先ハイパーバイザのL I Dに変更させなければならない。さらに、Q P Oアクセスが制限された結果、サブネットマネージャはV Mの内部で実行させることができなくなる。

10

【 0 0 5 4 】

インフィニバンドS R - I O Vアーキテクチャモデル - 仮想スイッチ (v S w i t c h)

図4は、一実施形態に従った例示的なv S w i t c hアーキテクチャを示す。図に示されるように、ホスト400（たとえばホストチャネルアダプタ）はハイパーバイザ410と対話することができ、当該ハイパーバイザ410は、さまざまな仮想機能430、440および450をいくつかの仮想マシンに割り当てることができる。同様に、物理機能はハイパーバイザ410によって処理することができる。仮想スイッチ415もハイパーバイザ401によって処理することができる。

20

【 0 0 5 5 】

一実施形態に従うと、v S w i t c hアーキテクチャにおいては、各々の仮想機能430、440、450は完全な仮想ホストチャネルアダプタ（virtual Host Channel Adapter: v H C A）であり、これは、ハードウェアにおいて、V Fに割り当てられたV Mに、I Bアドレス一式（たとえばG I D、G U I D、L I D）および専用のQ Pスペースが割り当てられていることを意味する。残りのネットワークおよびS Mについては、H C A 400は、仮想スイッチ415を介して追加のノードが接続されているスイッチのように見えている。ハイパーバイザ410はP F 420を用いることができ、（仮想機能に付与された）V MはV Fを用いる。

【 0 0 5 6 】

一実施形態に従うと、v S w i t c hアーキテクチャは、トランスペアレントな仮想化を提供する。しかしながら、各々の仮想機能には固有のL I Dが割り当てられているので、利用可能な数のL I Dが速やかに消費される。同様に、多くのL I Dアドレスが（すなわち、各々の物理機能および各々の仮想機能ごとに1つずつ）使用されている場合、より多くの通信経路をS Mによって演算しなければならず、それらのL F Tを更新するために、より多くのサブネット管理パケット（S M P）をスイッチに送信しなければならない。たとえば、通信経路の演算は大規模ネットワークにおいては数分かかる可能性がある。L I Dスペースが49151個のユニキャストL I Dに制限されており、（V Fを介する）各々のV Mとして、物理ノードおよびスイッチがL I Dを1つずつ占有するので、ネットワークにおける物理ノードおよびスイッチの数によってアクティブなV Mの数が制限されてしまい、逆の場合も同様に制限される。

30

40

【 0 0 5 7 】

インフィニバンドS R - I O Vアーキテクチャモデル - 仮想ポート (v P o r t)

図5は、一実施形態に従った例示的なv P o r tの概念を示す。図に示されるように、ホスト300（たとえばホストチャネルアダプタ）は、さまざまな仮想機能330、340および350をいくつかの仮想マシンに割り当てることができるハイパーバイザ410と対話することができる。同様に、物理機能はハイパーバイザ310によって処理することができる。

【 0 0 5 8 】

一実施形態に従うと、ベンダーに実装の自由を与えるためにv P o r t概念は緩やかに

50

定義されており（たとえば、当該定義では、実装がSRIOV専用とすべきであるとは規定されていない）、vPortの目的は、VMがサブネットにおいて処理される方法を標準化することである。vPort概念であれば、空間ドメインおよび性能ドメインの両方においてよりスケラブルであり得る、SR-IOV共有のポートのようなアーキテクチャおよびvSwitchのようなアーキテクチャの両方、または、これらのアーキテクチャの組合せが規定され得る。また、vPortはオプションのLIDをサポートするとともに、共有のポートとは異なり、SMは、vPortが専用のLIDを用いていなくても、サブネットにおいて利用可能なすべてのvPortを認識する。

【0059】

インフィニバンドSR-IOVアーキテクチャモデル - LIDが予めボピュレートされたvSwitch

10

一実施形態に従うと、本開示は、LIDが予めボピュレートされたvSwitchアーキテクチャを提供するためのシステムおよび方法を提供する。

【0060】

図6は、一実施形態に従った、LIDが予めボピュレートされた例示的なvSwitchアーキテクチャを示す。図に示されるように、いくつかのスイッチ501~504は、ネットワーク切替環境600（たとえばIBサブネット）内においてインフィニバンドファブリックなどのファブリックのメンバ間で通信を確立することができる。ファブリックはホストチャネルアダプタ510、520、530などのいくつかのハードウェアデバイスを含み得る。さらに、ホストチャネルアダプタ510、520および530は、それぞれ、ハイパーバイザ511、521および531と対話することができる。各々のハイパーバイザは、さらに、ホストチャネルアダプタと共に、いくつかの仮想機能514、515、516、524、525、526、534、535および536と対話し、設定し、いくつかの仮想マシンに割当てることができる。たとえば、仮想マシン1 550はハイパーバイザ511によって仮想機能1 514に割当てることができる。ハイパーバイザ511は、加えて、仮想マシン2 551を仮想機能2 515に割当て、仮想マシン3 552を仮想機能3 516に割当てることができる。ハイパーバイザ531は、さらに、仮想マシン4 553を仮想機能1 534に割当てることができる。ハイパーバイザは、ホストチャネルアダプタの各々の上で十分な機能を有する物理機能513、523および533を介してホストチャネルアダプタにアクセスすることができる。

20

30

【0061】

一実施形態に従うと、スイッチ501~504の各々はいくつかのポート（図示せず）を含み得る。いくつかのポートは、ネットワーク切替環境600内においてトラフィックを方向付けるためにリニアフォワーディングテーブルを設定するのに用いられる。

【0062】

一実施形態に従うと、仮想スイッチ512、522および532は、それぞれのハイパーバイザ511、521、531によって処理することができる。このようなvSwitchアーキテクチャにおいては、各々の仮想機能は完全な仮想ホストチャネルアダプタ（vHCA）であり、これは、ハードウェアにおいて、VFに割当てられたVMに、IBアドレス一式（たとえばGID、GUID、LID）および専用のQPスペースが割当てられていることを意味する。残りのネットワークおよびSM（図示せず）については、HCA510、520および530は、仮想スイッチを介して追加のノードが接続されているスイッチのように見えている。

40

【0063】

一実施形態に従うと、本開示は、LIDが予めボピュレートされたvSwitchアーキテクチャを提供するためのシステムおよび方法を提供する。図5を参照すると、LIDは、さまざまな物理機能513、523および533に、さらには、仮想機能514~516、524~526、534~536（その時点でアクティブな仮想マシンに関連付けられていない仮想機能であっても）にも、予めボピュレートされている。たとえば、物理機能513はLID1が予めボピュレートされており、仮想機能1 534はLID10

50

が予めポピュレートされている。ネットワークがブートされているとき、L I DはS R - I O V v S w i t c h対応のサブネットにおいて予めポピュレートされている。V FのすべてがネットワークにおけるV Mによって占有されていない場合であっても、ポピュレートされたV Fには、図5に示されるようにL I Dが割当てられている。

【0064】

一実施形態に従うと、多くの同様の物理的なホストチャネルアダプタが2つ以上のポートを有することができ（冗長性のために2つのポートが共用となっている）、仮想H C Aも2つのポートで表わされ、1つまたは2つ以上の仮想スイッチを介して外部I Bサブネットに接続され得る。

【0065】

一実施形態に従うと、L I Dが予めポピュレートされたv S w i t c hアーキテクチャにおいては、各々のハイパーバイザは、それ自体のための1つのL I DをP Fを介して消費し、各々の追加のV Fごとに1つ以上のL I Dを消費することができる。I Bサブネットにおけるすべてのハイパーバイザにおいて利用可能なすべてのV Fを合計すると、サブネットにおいて実行することが可能なV Mの最大量が得られる。たとえば、サブネット内の1ハイパーバイザごとに16個の仮想機能を備えたI Bサブネットにおいては、各々のハイパーバイザは、サブネットにおいて17個のL I D（16個の仮想機能ごとに1つのL I Dと、物理機能のために1つのL I D）を消費する。このようなI Bサブネットにおいては、単一のサブネットについて理論上のハイパーバイザ限度は利用可能なユニキャストL I Dの数によって規定されており、（49151個の利用可能なL I Dをハイパーバイザごとに17個のL I Dで割って得られる）2891であり、V Mの総数（すなわち限度）は（ハイパーバイザごとに2891個のハイパーバイザに16のV Fを掛けて得られる）46256である（実質的には、I Bサブネットにおける各々のスイッチ、ルータまたは専用のS Mノードが同様にL I Dを消費するので、これらの数は実際にはより小さくなる）。なお、v S w i t c hが、L I DをP Fと共有することができるので、付加的なL I Dを占有する必要がないことに留意されたい。

【0066】

一実施形態に従うと、L I Dが予めポピュレートされたv S w i t c hアーキテクチャにおいては、ネットワークが一旦ブートされると、すべてのL I Dについて通信経路が計算される。新しいV Mを始動させる必要がある場合、システムは、サブネットにおいて新しいL I Dを追加する必要はない。それ以外の場合、経路の再計算を含め、ネットワークを完全に再構成させ得る動作は、最も時間を消費する要素となる。代わりに、V Mのための利用可能なポートはハイパーバイザのうちの1つに位置し（すなわち利用可能な仮想機能）、仮想マシンは利用可能な仮想機能に付与されている。

【0067】

一実施形態に従うと、L I Dが予めポピュレートされたv S w i t c hアーキテクチャはまた、同じハイパーバイザによってホストされているさまざまなV Mに達するために、さまざまな経路を計算して用いる能力を可能にする。本質的には、これは、L I Dを連続的にすることを必要とするL M Cの制約によって拘束されることなく、1つの物理的なマシンに向かう代替的な経路を設けるために、このようなサブネットおよびネットワークがL I Dマスク制御ライク（L I D-Mask-Control-like：L M Cライク）な特徴を用いることを可能にする。V Mをマイグレートしてその関連するL I Dを宛先に送達する必要がある場合、不連続なL I Dを自由に使用できることは特に有用となる。

【0068】

一実施形態に従うと、L I Dが予めポピュレートされたv S w i t c hアーキテクチャについての上述の利点と共に、いくつかの検討事項を考慮に入れることができる。たとえば、ネットワークがブートされているときに、S R - I O V v S w i t c h対応のサブネットにおいてL I Dが予めポピュレートされているので、（たとえば起動時の）最初の経路演算はL I Dが予めポピュレートされていなかった場合よりも時間が長くなる可能性がある。

10

20

30

40

50

【 0 0 6 9 】

インフィニバンド S R - I O V アーキテクチャモデル - 動的 L I D 割当てがなされた v S w i t c h

一実施形態に従うと、本開示は、動的 L I D 割当てがなされた v S w i t c h アーキテクチャを提供するためのシステムおよび方法を提供する。

【 0 0 7 0 】

図 7 は、一実施形態に従った、動的 L I D 割当てがなされた例示的な v S w i t c h アーキテクチャを示す。図に示されるように、いくつかのスイッチ 5 0 1 ~ 5 0 4 は、ネットワーク切替環境 7 0 0 (たとえば I B サブネット) 内においてインフィニバンドファブリックなどのファブリックのメンバ間で通信を確立することができる。ファブリックは、
10
ホストチャネルアダプタ 5 1 0、5 2 0、5 3 0 などのいくつかのハードウェアデバイスを含み得る。ホストチャネルアダプタ 5 1 0、5 2 0 および 5 3 0 は、さらに、ハイパーバイザ 5 1 1、5 2 1 および 5 3 1 とそれぞれ対話することができる。各々のハイパーバイザは、さらに、ホストチャネルアダプタと共に、いくつかの仮想機能 5 1 4、5 1 5、5 1 6、5 2 4、5 2 5、5 2 6、5 3 4、5 3 5 および 5 3 6 と対話し、設定し、いくつかの仮想マシンに割当てることができる。たとえば、仮想マシン 1 5 5 0 はハイパーバイザ 5 1 1 によって仮想機能 1 5 1 4 に割当てることができる。ハイパーバイザ 5 1 1 は、加えて、仮想マシン 2 5 5 1 を仮想機能 2 5 1 5 に割当て、仮想マシン 3 5 5 2 を仮想機能 3 5 1 6 に割当てることができる。ハイパーバイザ 5 3 1 はさらに、仮想マシン 4 5 5 3 を仮想機能 1 5 3 4 に割当てることができる。ハイパーバイザは、
20
ホストチャネルアダプタの各々の上において十分な機能を有する物理機能 5 1 3、5 2 3 および 5 3 3 を介してホストチャネルアダプタにアクセスすることができる。

【 0 0 7 1 】

一実施形態に従うと、スイッチ 5 0 1 ~ 5 0 4 の各々はいくつかのポート (図示せず) を含み得る。いくつかのポートは、ネットワーク切替環境 7 0 0 内においてトラフィックを方向付けるためにリニアフォワーディングテーブルを設定するのに用いられる。

【 0 0 7 2 】

一実施形態に従うと、仮想スイッチ 5 1 2、5 2 2 および 5 3 2 は、それぞれのハイパーバイザ 5 1 1、5 2 1 および 5 3 1 によって処理することができる。このような v S w i t c h アーキテクチャにおいては、各々の仮想機能は完全な仮想ホストチャネルアダプタ (v H C A) であり、これは、ハードウェアにおいて、V F に割当てられた V M に、I B アドレス一式 (たとえば G I D、G U I D、L I D) および専用の Q P スペースが割当てられていることを意味する。残りのネットワークおよび S M (図示せず) については、H C A 5 1 0、5 2 0 および 5 3 0 は、仮想スイッチを介して、追加のノードが接続されているスイッチのように見えている。

【 0 0 7 3 】

一実施形態に従うと、本開示は、動的 L I D 割当てがなされた v S w i t c h アーキテクチャを提供するためのシステムおよび方法を提供する。図 7 を参照すると、L I D には、さまざまな物理機能 5 1 3、5 2 3 および 5 3 3 が動的に割当てられており、物理機能 5 1 3 が L I D 1 を受取り、物理機能 5 2 3 が L I D 2 を受取り、物理機能 5 3 3 が L I D 3 を受取る。アクティブな仮想マシンに関連付けられたそれらの仮想機能はまた、動的に割当てられた L I D を受取ることもできる。たとえば、仮想マシン 1 5 5 0 がアクティブであり、仮想機能 1 5 1 4 に関連付けられているので、仮想機能 5 1 4 には L I D 5 が割当てられ得る。同様に、仮想機能 2 5 1 5、仮想機能 3 5 1 6 および仮想機能 1 5 3 4 は、各々、アクティブな仮想機能に関連付けられている。このため、これらの仮想機能に L I D が割当てられ、L I D 7 が仮想機能 2 5 1 5 に割当てられ、L I D 1 1 が仮想機能 3 5 1 6 に割当てられ、L I D 9 が仮想機能 1 5 3 4 に割当てられている。L I D が予めポピュレートされた v S w i t c h とは異なり、アクティブな仮想マシンにその時点で関連付けられていない仮想機能は L I D の割当てを受けない。

【 0 0 7 4 】

10

20

30

40

50

一実施形態に従うと、動的 L I D 割当てがなされていれば、最初の経路演算を実質的に減らすことができる。ネットワークが初めてブートしており、V M が存在していない場合、比較的少数の L I D を最初の経路計算および L F T 分配のために用いることができる。

【 0 0 7 5 】

一実施形態に従うと、多くの同様の物理的なホストチャネルアダプタが 2 つ以上のポートを有することができ（冗長性のために 2 つのポートが共用となっている）、仮想 H C A も 2 つのポートで表わされ、1 つまたは 2 つ以上の仮想スイッチを介して外部 I B サブネットに接続され得る。

【 0 0 7 6 】

一実施形態に従うと、動的 L I D 割当てがなされた v S w i t c h を利用するシステムにおいて新しい V M が作成される場合、どのハイパーバイザ上で新しく追加された V M をブートすべきであるかを決定するために、自由な V M スロットが発見され、固有の未使用のユニキャスト L I D も同様に発見される。しかしながら、新しく追加された L I D を処理するためのスイッチの L F T およびネットワークに既知の経路が存在しない。新しく追加された V M を処理するために新しいセットの経路を演算することは、いくつかの V M が毎分ごとにブートされ得る動的な環境においては望ましくない。大規模な I B サブネットにおいては、新しい 1 セットのルートの演算には数分かかる可能性があり、この手順は、新しい V M がブートされるたびに繰返されなければならないだろう。

【 0 0 7 7 】

有利には、一実施形態に従うと、ハイパーバイザにおけるすべての V F が P F と同じアップリンクを共有しているので、新しいセットのルートを演算する必要はない。ネットワークにおけるすべての物理スイッチの L F T を繰返し、（ V M が作成されている）ハイパーバイザの P F に属する L I D エントリから新しく追加された L I D にフォーワーディングポートをコピーし、かつ、特定のスイッチの対応する L F T ブロックを更新するために単一の S M P を送信するだけでよい。これにより、当該システムおよび方法では、新しいセットのルートを演算する必要がなくなる。

【 0 0 7 8 】

一実施形態に従うと、動的 L I D 割当てアーキテクチャを備えた v S w i t c h において割当てられた L I D は連続的である必要はない。各々のハイパーバイザ上の V M 上で割当てられた L I D を L I D が予めポピュレートされた v S w i t c h と動的 L I D 割当てがなされた v S w i t c h とで比較すると、動的 L I D 割当てアーキテクチャにおいて割当てられた L I D が不連続であり、そこに予めポピュレートされた L I D が本質的に連続的であることが分かるだろう。さらに、v S w i t c h 動的 L I D 割当てアーキテクチャにおいては、新しい V M が作成されると、次に利用可能な L I D が、V M の生存期間の間中ずっと用いられる。逆に、L I D が予めポピュレートされた v S w i t c h においては、各々の V M は、対応する V F に既に割当てられている L I D を引継ぎ、ライブマイグレーションのないネットワークにおいては、所与の V F に連続的に付与された V M が同じ L I D を得る。

【 0 0 7 9 】

一実施形態に従うと、動的 L I D 割当てアーキテクチャを備えた v S w i t c h は、いくらかの追加のネットワークおよびランタイム S M オーバーヘッドを犠牲にして、予めポピュレートされた L I D アーキテクチャモデルを備えた v S w i t c h の欠点を解決することができる。V M が作成されるたびに、作成された V M に関連付けられた、新しく追加された L I D で、サブネットにおける物理スイッチの L F T が更新される。この動作のために、1 スイッチごとに 1 つのサブネット管理パケット（ S M P ）が送信される必要がある。各々の V M がそのホストハイパーバイザと同じ経路を用いているので、L M C のような機能も利用できなくなる。しかしながら、すべてのハイパーバイザに存在する V F の合計に対する制限はなく、V F の数は、ユニキャスト L I D の限度を上回る可能性もある。このような場合、当然、アクティブな V M 上で V F のすべてが必ずしも同時に付与されることが可能になるわけではなく、より多くの予備のハイパーバイザおよび V F を備えるこ

10

20

30

40

50

とにより、ユニキャスト L I D 限度付近で動作する際に、断片化されたネットワークの障害を回復および最適化させるための融通性が追加される。

【 0 0 8 0 】

インフィニバンド S R - I O V アーキテクチャモデル - 動的 L I D 割当てがなされかつ L I D が予めポピュレートされた v S w i t c h

図 8 は、一実施形態に従った、動的 L I D 割当てがなされて L I D が予めポピュレートされた v S w i t c h を備えた例示的な v S w i t c h アーキテクチャを示す。図に示されるように、いくつかのスイッチ 5 0 1 ~ 5 0 4 は、ネットワーク切替環境 8 0 0 (たとえば I B サブネット) 内においてインフィニバンドファブリックなどのファブリックのメンバ間で通信を確立することができる。ファブリックはホストチャネルアダプタ 5 1 0、5 2 0、5 3 0 などのいくつかのハードウェアデバイスを含み得る。ホストチャネルアダプタ 5 1 0、5 2 0 および 5 3 0 は、それぞれ、さらに、ハイパーバイザ 5 1 1、5 2 1 および 5 3 1 と対話することができる。各々のハイパーバイザは、さらに、ホストチャネルアダプタと共に、いくつかの仮想機能 5 1 4、5 1 5、5 1 6、5 2 4、5 2 5、5 2 6、5 3 4、5 3 5 および 5 3 6 と対話し、設定し、いくつかの仮想マシンに割当てることができる。たとえば、仮想マシン 1 5 5 0 は、ハイパーバイザ 5 1 1 によって仮想機能 1 5 1 4 に割当てることができる。ハイパーバイザ 5 1 1 は、加えて、仮想マシン 2 5 5 1 を仮想機能 2 5 1 5 に割当てることができる。ハイパーバイザ 5 2 1 は、仮想マシン 3 5 5 2 を仮想機能 3 5 2 6 に割当てることができる。ハイパーバイザ 5 3 1 は、さらに、仮想マシン 4 5 5 3 を仮想機能 2 5 3 5 に割当てることができる。ハイパーバイザは、ホストチャネルアダプタの各々の上において十分な機能を有する物理機能 5 1 3、5 2 3 および 5 3 3 を介してホストチャネルアダプタにアクセスすることができる。

【 0 0 8 1 】

一実施形態に従うと、スイッチ 5 0 1 ~ 5 0 4 の各々はいくつかのポート (図示せず) を含み得る。これらいくつかのポートは、ネットワーク切替環境 8 0 0 内においてトラフィックを方向付けるためにリニアフォワーディングテーブルを設定するのに用いられる。

【 0 0 8 2 】

一実施形態に従うと、仮想スイッチ 5 1 2、5 2 2 および 5 3 2 は、それぞれのハイパーバイザ 5 1 1、5 2 1、5 3 1 によって処理することができる。このような v S w i t c h アーキテクチャにおいては、各々の仮想機能は、完全な仮想ホストチャネルアダプタ (v H C A) であり、これは、ハードウェアにおいて、V F に割当てられた V M に、I B アドレス一式 (たとえば G I D、G U I D、L I D) および専用の Q P スペースが割当てられていることを意味する。残りのネットワークおよび S M (図示せず) については、H C A 5 1 0、5 2 0 および 5 3 0 は、仮想スイッチを介して、追加のノードが接続されているスイッチのように見えている。

【 0 0 8 3 】

一実施形態に従うと、本開示は、動的 L I D 割当てがなされ L I D が予めポピュレートされたハイブリッド v S w i t c h アーキテクチャを提供するためのシステムおよび方法を提供する。図 7 を参照すると、ハイパーバイザ 5 1 1 には、予めポピュレートされた L I D アーキテクチャを備えた v S w i t c h が配置され得るとともに、ハイパーバイザ 5 2 1 には、L I D が予めポピュレートされて動的 L I D 割当てがなされた v S w i t c h が配置され得る。ハイパーバイザ 5 3 1 には、動的 L I D 割当てがなされた v S w i t c h が配置され得る。このため、物理機能 5 1 3 および仮想機能 5 1 4 ~ 5 1 6 には、それらの L I D が予めポピュレートされている (すなわち、アクティブな仮想マシンに付与されていない仮想機能であっても L I D が割当てられている)。物理機能 5 2 3 および仮想機能 1 5 2 4 にはそれらの L I D が予めポピュレートされ得るとともに、仮想機能 2 5 2 5 および仮想機能 3 5 2 6 にはそれらの L I D が動的に割当てられている (すなわち、仮想機能 2 5 2 5 は動的 L I D 割当てのために利用可能であり、仮想機能 3 5 2 6 は、仮想マシン 3 5 5 2 が付与されているので、1 1 という L I D が動的に割当てら

10

20

30

40

50

れている)。最後に、ハイパーバイザ 3 5 3 1 に関連付けられた機能（物理機能および仮想機能）にはそれらの L I D を動的に割当てることができる。これにより、結果として、仮想機能 1 5 3 4 および仮想機能 3 5 3 6 が動的 L I D 割当てのために利用可能となるとともに、仮想機能 2 5 3 5 には、仮想マシン 4 5 5 3 が付与されているので、9 という L I D が動的に割当てられている。

【 0 0 8 4 】

L I D が予めポピュレートされた v S w i t c h および動的 L I D 割当てがなされた v S w i t c h がともに（いずれかの所与のハイパーバイザ内で独立して、または組合わされて）利用されている、図 8 に示されるような一実施形態に従うと、ホストチャネルアダプタごとの予めポピュレートされた L I D の数はファブリックアドミニストレータによって定義することができ、（ホストチャネルアダプタごとに） $0 \leq$ 予めポピュレートされた V F \leq 総 V F の範囲内になり得る。動的 L I D 割当てのために利用可能な V F は、（ホストチャネルアダプタごとに）V F の総数から予めポピュレートされた V F の数を減じることによって見出すことができる。

【 0 0 8 5 】

一実施形態に従うと、多くの同様の物理的なホストチャネルアダプタが 2 つ以上のポートを有することができ（冗長性のために 2 つのポートが共用となっている）、仮想 H C A も 2 つのポートで表わされ、1 つまたは 2 つ以上の仮想スイッチを介して外部 I B サブネットに接続され得る。

【 0 0 8 6 】

v S w i t c h スケーラビリティ

一実施形態に従うと、v S w i t c h アーキテクチャを用いる場合の問題は L I D スペースが制限されていることである。L I D スペースに関するスケーラビリティの問題を克服するために、以下の 3 つの代替例（各々を以下にさらに詳細に説明する）を独立して用いるかまたは組合わせることができる：すなわち、複数のサブネットを用いること；後方互換性のある L I D スペース拡張を導入すること；および、軽量の v S w i t c h を形成するために v P o r t アーキテクチャと v S w i t c h アーキテクチャとを組合わせること；である。

【 0 0 8 7 】

一実施形態に従うと、複数の I B サブネットを用いることができる。L I D は、層 2 アドレスであり、サブネット内において固有でなければならない。I B トポロジが複数のサブネット上にわたっている場合、L I D はそれ以上制限事項とはならないが、V M を異なるサブネットにマイグレートする必要がある場合、その L I D アドレスは変更することができる。なぜなら、そのアドレスが新しいサブネットにおいて既に使用されているかもしれないからである。複数のサブネット上にわたっていることで、単一のサブネットトポロジの L I D 制限を解決することができるが、これはまた、ルーティングプロセスに付加的なオーバーヘッドおよび待ち時間を付加するサブネット間ルーティングのために層 3 G I D アドレスを用いなければならないことを意味している。なぜなら、サブネットの端に位置するルータによって層 2 ヘッダを変更しなければならないからである。また、現在のハードウェア実装、ソフトウェア実装および緩い I B A（インフィニバンド・アーキテクチャ）規格の下では、複数のサブネット上にわたっているクラスタのために最適化されたルーティング経路を提供するために、個々のサブネットの S M はグローバルトポロジを認識することができなくなっている。

【 0 0 8 8 】

一実施形態に従うと、I B A における後方互換性のある L I D スペース拡張を導入することができる。L I D ビットの数、たとえば 2 4 ビットまたは 3 2 ビットを増やすことにより、不十分な L I D スペースを増やす場合に問題が生じる可能性がある。このような量だけ L I D スペースを増やすことにより、後方互換性に破断が生じる可能性がある。なぜなら、I B ローカルルートヘッダ（Local Route Header：L R H）がオーバーホールされなければならない、レガシーハードウェアが新しい基準では機能することができな

10

20

30

40

50

くなるからである。一実施形態に従うと、後方互換性を維持しながらも、依然として新しいハードウェアが拡張機能を利用できるように、L I Dスペースを拡張することができる。L R Hは、0として送信されて受信機には無視される予備の7ビットを有する。送信元L I D (Source LID: S L I D) についてのL R Hにおけるこれらの予備ビットと宛先L I D (Destination LID: D L I D) についての2ビットとのうち2つを利用することにより、L I Dスペースを18ビットに拡張する(L I Dスペースを4倍にすること)ことができ、物理的装置に割当てられた物理的L I D (p L I D) およびV Mに割当てられた仮想L I D (v L I D) を用いたスキームを作成することができる。

【0089】

一実施形態に従うと、追加の2ビットが0として送信されると、L I DがI B A (48 KユニキャストL I Dおよび16 KマルチキャストL I D) においてその時点で定義されるとおりに用いられ、スイッチは、パケットの転送のためにそれらの主要なL F Tを検索することができる。他の場合、L I Dはv L I Dであり、192 Kのサイズを有する二次的L F Tに基づいて転送することができる。v L I DがV Mに属しており、V Mが、p L I Dを有する物理ノードとアップリンクを共有しているので、v L I Dは、ネットワークを構成(たとえば初期構成)または(たとえばトポロジー変更後に)再構成する際に、経路演算段階から除外することができるが、スイッチにおける二次的L F Tテーブルは、上述のとおり更新することができる。S Mがネットワークをブートし発見すると、S Mはハードウェアのすべてが拡張されたL I Dスペースをサポートするかどうかを識別することができる。そうでなければ、S Mはレガシー互換モードでフォールバックすることができ、V Mはp L I DスペースからのL I Dを占有するはずである。

【0090】

図9は、一実施形態に従った、拡張されたローカルルートヘッダを示す。図に示されるように、ローカルルートヘッダ内では、仮想レーン(virtual lane: V L) 900は4ビットを含み、リンクバージョン(link version: L v e r) 901は4ビットを含み、サービスレベル(service level: S L) 902は4ビットを含み、L I D拡張フラグ(LID extension flag: L E X T F) 903は1ビットを含み、第1の予備ビット(R 1) 904は1ビットを含み、リンク次ヘッダ(link next header: L N H) 905は2ビットを含み、宛先ローカルI D (destination local ID: D L I D) 906は16ビットを含み、D L I Dプレフィックス拡張(DLID prefix extension: D P F) 907は2ビットを含み、S L I Dプレフィックス拡張(SLID prefix extension: S P F) 908は2ビットを含み、第2の予備ビット(R 2) 909は1ビットを含み、パケット長(packet length: P k t L e n) 910は11ビット含み、送信元ローカルI D (source local ID: S L I D) 911は16ビットを含む。一実施形態に従うと、両方の予備ビット904および909はゼロに設定することができる。

【0091】

一実施形態に従うと、上述のとおり、図9に示されるL R Hは、宛先ローカルI D 906および送信元ローカルI D 911についてのプレフィックス拡張として7つの(元の)予備ビットのうちの4つを利用する。これにより、利用時に、L I D拡張フラグに関連付けて、スイッチにおける二次的L F Tを介してルーティングされ得るv L I Dに関連付けてL R Hが使用されることが信号で伝えられる。代替的には、拡張907および908がゼロとして送信され(受信機によって無視され)ると、L I Dは、p L I Dに関連付けられ、I B Aにおいてその時点で定義されるとおり用いられる。

【0092】

図10は、一実施形態に従った、2つの例示的なリニアフォワーディングテーブルを示す。図10に示されるように、リニアフォワーディングテーブル916はp L I Dに関連付けられたフォワーディングテーブルである。L F Tは、エントリ912(D L I D = 0)によって索引付けされたエントリ0からエントリ913(D L I D = 48 K - 1)によって索引付けされたエントリ48 K - 1)にわたっている。この場合、L F Tにおける各エントリは、規格16ビットD L I Dによって索引付けされ、標準I Bポート番号を含んで

いる。対照的に、リニアフォーディングテーブル 917 は $v L I D$ に関連付けられた二次的フォーディングテーブルである。 $L F T$ は、エントリ 914 ($18 \text{ ビット } D P F + D L I D = 0$) によって索引付けされたエントリ 0) からエントリ 915 ($18 \text{ ビット } D P F + D L I D = 256 K - 1$) によって索引付けされたエントリ $256 K - 1$) にわたっている。この場合、各エントリは、拡張された $18 \text{ ビット } D P F + D L I D$ によって索引付けされ、標準 $I B$ ポート番号を含んでいる。

【 0093 】

一実施形態に従うと、軽量の $v S w i t c h$ アーキテクチャを形成するためにハイブリッドアーキテクチャを用いることができる。マイグレートされた $V M$ と共に $L I D$ をマイグレートすることができる $v S w i t c h$ アーキテクチャは、 $L I D$ が変化するであろう共有の $L I D$ のスキームとは対照的にマイグレーションの後にピアとの接続性を再構築するために付加的なシグナリングについての要件が存在しないので、サブネット管理に対して十分にスケールリングする。他方で、共有の $L I D$ スキームは、 $L I D$ スペースに対して十分にスケールリングする。ハイブリッド $v S w i t c h +$ 共有型 $v P o r t$ モデルは、 $S M$ がサブネットにおける利用可能な $S R - I O V$ 仮想機能を認識する場合、実現することができるが、特定の $V F$ が専用の $L I D$ を受取り得る一方で、他のものはそれらの $G I D$ に基づいて共有 $L I D$ の態様でルーティングされている。 $V M$ ノード役割についての何らかの情報があれば、(たとえば、ルートを計算し、ネットワークにおける負荷バランスングを実行している間に別々に考慮されるようにするために)、多数のピアを備えたポピュラーな $V M$ (たとえばサーバ) には専用の $L I D$ が割当てられ得る一方で、多くのピアと対話しないかまたはステートレスなサービスを実行する(マイグレートされる必要がなく、再生成され得る)他の $V M$ は $L I D$ を共有することができる。

【 0094 】

$v S w i t c h$ ベースのサブネットのためのルーティング戦略

一実施形態に従うと、より高い性能を得るために、ルーティングアルゴリズムは、ルートを計算する際に $v S w i t c h$ アーキテクチャを考慮に入れることができる。ファットツリーにおいては、 $v S w i t c h$ は、 $v S w i t c h$ が対応するリーフスイッチへの上りリンクを 1 つだけ有するという独特な特性によってトポロジー発見プロセスにおいて識別することができる。 $v S w i t c h$ が識別されると、ルーティング機能は、各 $V M$ からのトラフィックがネットワークにおける他のすべての $V M$ に向かう経路を発見することができるように、すべてのスイッチのための $L F T$ を生成することができる。各 $V M$ はそれ自体のアドレスを有しており、このため、各 $V M$ は、同じ $v S w i t c h$ に付与された他の $V M$ からは独立してルーティングすることができる。これにより、結果として、トポロジーにおける $v S w i t c h$ に向かうとともに各々が特定の $V M$ へのトラフィックを担持している独立した複数の経路を生成するルーティング機能が得られる。このアプローチの 1 つの欠点として、 $V M$ 分配が $v S w i t c h$ の間で均一でない場合、より多くの $V M$ を備えた $v S w i t c h$ には潜在的により大きなネットワークリソースが割当てられる点がある。しかしながら、 $v S w i t c h$ から対応するリーフスイッチまでの単一の上りリンクは、依然として、特定の $v S w i t c h$ に付与されたすべての $V M$ によって共有されるボトルネックリンクのままである。結果として、準最適にネットワークが利用される可能性がある。最も単純で最速のルーティング戦略は、すべての $v S w i t c h - v S w i t c h$ の対の間に経路を生成して、対応する $v S w i t c h$ に割当てられるのと同じ経路を備えた $V M$ をルーティングすることである。予めポピュレートされた $L I D$ 割当てスキームと動的 $L I D$ 割当てスキームとがあれば、各々の $v S w i t c h$ は、 $S R - I O V$ アーキテクチャにおける $P F$ によって定義された $L I D$ を有する。 $v S w i t c h$ についてのこれらの $P F - L I D$ は、ルーティングの第 1 段階で $L F T$ を生成するために用いることができ、第 2 段階では、 $V M$ の $L I D$ を生成された $L F T$ に追加することができる。予めポピュレートされた $L I D$ スキームにおいては、 $V F - L I D$ へのエントリは対応する $v S w i t c h$ の出力ポートをコピーすることによって追加することができる。同様に、新しい $V M$ がブートされた場合の動的 $L I D$ 割当ての場合、 $V M$ の $L I D$ と対応する $v S w$

i t c hによって決定された出力ポートとを備えた新しいエントリがすべてのL F Tにおいて追加される。この戦略についての問題点は、v S w i t c hを共有する別々のテナントに属するV Mが、ネットワークにおいて同じ完全な経路を共有しているせいで、それらの中で固有に干渉する可能性がある点である。高いネットワーク利用率を維持しながらもこの問題を解決するために、仮想化されたサブネットのための重み付けされたルーティングスキームを用いることができる。

【 0 0 9 5 】

一実施形態に従うと、v S w i t c hベースの仮想化サブネットのための重み付けされたルーティングスキームを利用することができる。このようなメカニズムにおいては、v S w i t c h上の各V Mには、ルートを計算する際にバランスを取るために考慮に入れることができるパラメータ重みが割当てられる。重みパラメータの値は、そのv S w i t c hにおけるV Mに割付けられたリーフスイッチリンク容量に対するv S w i t c hの割合を反映している。たとえば、単純な構成により、各V Mに、 $1 / \text{num_vms}$ に等しい重みが割当てられてもよく、この場合、 num_vms は、対応するv S w i t c hハイパーバイザ上のブートされたV Mの数である。別の可能な実現例は、最も重要なV Mに対して、これらV Mに向かって流れるトラフィックに優先順位を付けるために、より高い割合のv S w i t c h容量を割当てることであり得る。しかしながら、v S w i t c h毎のV Mの累積的な重みはすべてのv S w i t c h上で等しくなり得るので、トポロジにおけるリンクは、実際のV M分配によって影響されことなくバランスを取ることができる。同時に、スキームは、トポロジにおける中間リンクで同じv S w i t c h V M間における干渉をなくした上で、各V Mがネットワークにおいて独立してルーティングされ得る多重通路を可能にする。当該スキームは、V Mがその割当てられた容量を上回るのを確実に防止するために、V M率の上限ごとに、各v S w i t c h上での実施と組み合わせることができる。加えて、ネットワークにおいて複数のテナントグループが存在している場合、テナント認識型ルーティングのような技術は、テナント間でネットワーク全体を分離させるために、提案されたルーティングスキームと統合することができる。

【 0 0 9 6 】

一実施形態に従うと、以下に、I Bベースのファットツリートポロジについての重み付けされたルーティングを記載する。ファットツリールーティングアルゴリズムとして、v S w i t c h F a t T r e eは、サブネットにおける各V Mに関連付けられたL I DのためのすべてのスイッチにおけるL F Tを設定するために、ファットツリートポロジを再帰的に横断する。このメカニズムは決定論的であり、すべてのルートについての後方計算が宛先ノードから開始される宛先ベースのルーティングをサポートする。

【 0 0 9 7 】

仮想化されたサブネットについての重み付けされたファットツリールーティングアルゴリズム

【 0 0 9 8 】

【数 1】

```

1: procedure ROUTEVIRTUALIZEDNODES
2:   for all  $s \in \text{leafSwitches}[]$  do
3:     sort vswitches in the increasing order of connected virtual
       machines
4:     for all  $v \in \text{vSwitches}[]$  do
5:        $\text{num\_vms} \leftarrow \text{GETTOTALVMS}(v)$ 
6:        $\text{vm\_weight} \leftarrow 1/\text{num\_vms}$ 
7:       for all  $\text{vm} \in \text{vSwitches}[]$  do
8:          $\text{vm.weight} \leftarrow \text{vm\_weight}$ 
9:          $s.\text{LFT}[\text{vm.LID}] \leftarrow v.\text{port}$ 
10:         $\text{ROUTEDOWNGOINGBYGOINGUP}(s, \text{vm})$ 
11:      end for
12:    end for
13:  end for
14: end procedure
15: procedure ROUTEDOWNGOINGBYGOINGUP( $s, \text{vm}$ )
16:   $p \leftarrow \text{GETLEASTLOADEDPORT}(s.\text{UpGroups}[])$ 
17:   $r\text{Switch} \leftarrow p.\text{Switch}$ 
18:   $r\text{Switch}.\text{LFT}[\text{vm.LID}] \leftarrow p$ 
19:   $p.\text{Dwn} += \text{vm.weight}$ 
20:   $\text{ROUTEUPGOINGBYGOINGDOWN}(s, \text{vm})$ 
21:   $\text{ROUTEDOWNGOINGBYGOINGUP}(r\text{Switch}, \text{vm})$ 
22: end procedure
23: procedure ROUTEUPGOINGBYGOINGDOWN( $s, \text{vm}$ )
24:   for all  $g \in s.\text{DownGroups}[]$  do
25:     skip g if the LFT(vm.LID) is part of this group
26:      $p \leftarrow \text{GETLEASTLOADEDPORT}(g)$ 
27:      $r\text{Switch} \leftarrow p.\text{Switch}$ 
28:      $r\text{Switch}.\text{LFT}[\text{vm.LID}] \leftarrow p$ 
29:      $p.\text{Up} += \text{vm.weight}$ 
30:      $\text{ROUTEUPGOINGBYGOINGDOWN}(r\text{Switch}, \text{vm})$ 
31:   end for
32: end procedure

```

【0099】

－実施形態に従うと、 vSwitchFatTree ルーティングメカニズムは以下のよう
 に作用する。各々のVMには、比例した重みが割当てられる。この比例した重みは、
 vSwitch ノードの（たとえば、定数1として得られる）重みをその上で実行される
 VMの総数で割ることによって計算される。さまざまな重み付けスキームを実現すること
 もできる。たとえば、VMタイプに基づいて重みを割当てるための実現例を選ぶことがで
 きる。しかしながら、簡潔にするために、この説明は比例重み付けスキームに焦点を合
 わせている。各々のリーフスイッチのために、ルーティングメカニズムは、接続されたVM
 （行3）に基づいて減少する順序で、接続された vSwitch をソートする。この順序

は、より高い重みが付けられたVMが最初にルーティングされることを確実にするので、リンクに割当てられたルートのバランスを取ることができる。ルーティングメカニズムは、すべてのリーフスイッチおよびそれらの対応するvSwitchを通過し、各々のVMからツリー内を横断して、ROUTEDOWNGOINGBYGOINGUP（行10）をコールすることによって、ツリー内においてVMに向かう経路を再帰的に割当てる。各々のスイッチにおける下りポートは、利用可能な上りポート群のすべての中で最少累積の下り重み（downward weight）に基づいて選択されている（ROUTEDOWNGOINGBYGOINGUP；行16）。下りポートが選択されると、当該メカニズムは、ルーティングされているVMの重みによって、対応するポートについての下り累積重みを増やすことができる（ROUTEDOWNGOINGBYGOINGUP；行19）。下りポートが設定された後、ルーティングメカニズムは、ツリーを下降していくことによってすべての接続された下りスイッチ上において、VMに向かうルートのために上りポートを割当てることことができる（ポートについての対応する上り重み（upward weight）を更新する）（ROUTEUPGOINGBYGOINGDOWN；行20）。次いで、当該プロセスはツリーにおける次のレベルまで上っていくことによって繰返される。すべてのVMがルーティングされると、（擬似コードに図示されない）トポロジにおいてvSwitch経路とvSwitch経路との間でバランスを取るように等しい重み付けがなされているにも関わらず、アルゴリズムはまた、VMと同じ方法でvSwitchの物理的LIDをルーティングする。これは、最小限の再構成方法がライブマイグレーションの文脈において用いられる際にバランスを取るのを向上させるのに望ましい。また、vSwitchのベースとなる物理的LID上のルーティング経路は、再構成を必要とすることなく、新しいVMを迅速にデプロイするために予め定められた経路として用いることができる。しかしながら、一定の期間にわたって、全体的なルーティング性能は、元のvSwitchFatTreeルーティングの間にわずかに減少するだろう。性能の低下を制限するために、ある性能しきい値を超えたとき、vSwitchFatTreeに基づいた再構成をオフラインで実行してもよい。

【0100】

一実施形態に従うと、上述のルーティングメカニズムは、正規のノレガシーなルーティングメカニズムに勝るさまざまな改善を提供することができる。トポロジにおけるvSwitchまたはVMを考慮に入れていない当初のファットツリールーティングアルゴリズムとは異なり、vSwitchFatTreeは、vSwitchに印付けをして、vSwitchに接続された他のVMからは独立して各々のVMをルーティングする。同様に、vSwitch間で不均一なVM分配を行なうために、各々のVMには、vSwitch上で割付けられているリンクの割合に対応する重みが割当てられている。重みは、ファットツリーにおける経路配分のバランスを取るためのポートカウンタを維持するのに用いられる。スキームはまた、一般化された重み付けされたファットツリールーティングを可能にする。この場合、各々のVMには、ネットワークにおけるそのトラヒックプロファイルまたは役割の優先順位に基づいて重みを割当てることことができる。

【0101】

図11から図14は、一実施形態に従った、無損失相互接続ネットワークにおいて効率的な仮想化をサポートする例を示す。具体的には、図11は、4つのスイッチとして、ルートスイッチ925および926、リーフスイッチ920および921、さらには、4つの仮想スイッチVS1 931、VS2 941、VS3 951およびVS4 961を備えた2レベルのファットツリートポロジを示す。4つの仮想スイッチVS1 931、VS2 941、VS3 951およびVS4 961には、それぞれ、4つのホストノハイパーバイザ930、940、950、960が関連付けられており、この場合、4つの仮想スイッチは、8つの仮想マシンVM1 932、VM2 933、VM3 942、VM4 943、VM5 952、VM6 953、VM7 954、およびVM8 962のために接続性を提供する。

【0102】

vSwitchFatTreeルーティングをさらに詳しく説明するために、図11に示されるように、4つのエンドノード(vSwitch)を備えた単純な仮想化されたファットツリートポロジーについて検討する。リーフスイッチ920、VS1およびVS2に接続されたvSwitchの各々は、実行中の2つのVM(VS1についてはVM1およびVM2、ならびにVS2についてはVM3およびVM4)を有する。第2のリーフスイッチ921は、3つのVM(VM5、VM6、VM7)を備えたVS3を有し、1つのVMがホストvSwitch VS4の上で実行中である。各々のリーフスイッチは、両方のルートスイッチ925および926に接続されているため、ルートを介して各々のVMに向かうルートを設定するのに利用可能な2つの代替経路が存在している。VS1に接続されたVMのためのルーティングは、ルートスイッチからの選択された下り経路を示す円を用いて、図12に示される。VM1は925 920を用いてルーティングされ、VM2は926 920からルーティングされている。対応する下り負荷カウンタは、選択されたリンク上で更新されて、各々のVMのために0.5を追加する。同様に、図13に示されるように、VS2のためのルートを追加した後、VM3およびVM4は、リンク925 920およびリンク926 920を介してそれぞれルーティングされる。リーフスイッチ920に接続されたすべてのVMをルーティングした後、たとえVMが個々にルーティングされていたとしても、両方のリンク上の下り負荷の合計が等しくなることに留意されたい。リーフスイッチ921に接続されたvSwitch上のVM分配は異なっており、このため、1つのVMを備えたvSwitch(VS4)は最初にルーティングされることとなるだろう。ルート925 921がVM8に割付けられ、VS3に接続された3つのすべてのVMが926 921からルーティングされて、両方の下りリンク上で累積された負荷のバランスが取られた状態を維持するようにする。図14に示される最後のルーティングでは、トポロジーにVM分配がなされていると想定して、可能な限り、VMに向かう独立したルートと共に、各々のリンク上で負荷のバランスが取られている。

【0103】

仮想マシンライブマイグレーション上での最小限のオーバーヘッド再構成

一実施形態に従うと、Iterative Reconfiguration(反復再構成)と略され得る動的な再構成メカニズムは、VMがマイグレートされたときに、必要に応じて、ルートの切替えおよび更新をすべてを繰返す。しかしながら、サブネットにおける既存のLFT(すなわち、既に計算されたLFTであって、サブネット内の各スイッチに存在しているLFT)に応じて、スイッチのサブセットだけを実際に更新する必要がある。

【0104】

図15は、一実施形態に従った潜在的な仮想マシンマイグレーションを示す。より具体的には、図15は、ネットワークトポロジーにもかかわらず、対応するリーフスイッチだけがLFT更新を必要としているリーフスイッチ内のVMのマイグレーションの特別な事例を示している。

【0105】

図15に示されるように、サブネットは、いくつかのスイッチ、すなわち、スイッチ1 1301~スイッチ12 1312を含み得る。これらのスイッチのうちのいくつかは、スイッチ1 1301、スイッチ2 1302、スイッチ11 1311、スイッチ12 1312などのリーフスイッチを含み得る。サブネットは、付加的に、いくつかのホスト/ハイパーバイザ1330、1340、1350および1360、いくつかの仮想スイッチVS1 1331、VS2 1341、VS3 1351およびVS4 1361を含み得る。さまざまなホスト/ハイパーバイザは、仮想機能を介して、VM1 1332、VM2 1333、VM3 1334、VM4 1342、VM5 1343およびVM6 1352などのサブネット内の仮想マシンをホストすることができる。

【0106】

一実施形態に従うと、VM3が(太字矢印によって示されるように)付随しているハイパーバイザ1330からハイパーバイザ1340における自由な仮想機能にマイグレートする場合、リーフスイッチ1 1301におけるLFTだけが更新される必要がある。な

ぜなら、両方のハイパーバイザが同じリーフスイッチに接続されており、局所的な変更がネットワークの残りの部分に影響を及ぼさないからである。たとえば、最初のルーティングアルゴリズムは、ハイパーバイザ1360からハイパーバイザ1330に向かうトラフィックが実線（すなわち、12 9 5 3 1）によって印付けされた第1の経路を追従すると判断する。同様に、ハイパーバイザ1360からハイパーバイザ1340に向かうトラフィックは、破線（すなわち、12 10 6 4 1）によって印付けされた第2の経路を追従する。VM3がマイグレートされ、ネットワークを再構成するためにItercが用いられる場合、VM3に向かうトラフィックは、マイグレーションの前にハイパーバイザ1330に向かう第1の経路を追従し、マイグレーションの後、ハイパーバイザ1340に向かう第2の経路を追従することとなるだろう。この状況においては、ファットツリールーティングアルゴリズムが最初のルーティングのために用いられたと想定すると、Iterc法は、スイッチの総数の半分（6 / 12）を更新するだろう。しかしながら、マイグレートされたVMを接続されたままにしておくためにリーフスイッチを1つだけ更新する必要がある。

【0107】

一実施形態に従うと、VMマイグレーションの後にスイッチ更新の回数を制限することによって、ネットワークをより速く再構成することができ、従来のルーティング更新の際に必要とされる時間およびオーバーヘッドを減らすことができる。これは、トポロジーに依存しないスカイライン技術（topology-agnostic skyline technique）に基づいて、FTreeMinRCと称される、ファットツリー上でのVMマイグレーションをサポートするためのトポロジー認識型高速再構成方法によって達成することができる。

【0108】

ファットツリーにおけるサブツリーおよびスイッチタプル

一実施形態に従うと、以下の記述は、例示的なファットツリーネットワークとしてXGFTを用いて、最小限のオーバーヘッドネットワーク再構成方法であるFTreeMinRCを利用する。しかしながら、ここで提示される概念は、PGFTおよびRLFTにとっても有効である。XGFT($n; m_1, \dots, m_n; w_1, \dots, w_n$)は、 $n+1$ レベルのノードを備えたファットツリーである。レベルは0から n で表わされ、計算ノードがレベル n にあり、スイッチが他のすべてのレベルにある。子がない計算ノードを除いては、レベル i 、 $0 \leq i \leq n-1$ におけるすべてのノードは、 m_i の子ノードを有する。同様に、親がないルートスイッチを除いては、レベル i 、 $1 \leq i \leq n$ における他のすべてのノードは $w_i + 1$ の親ノードを有する。

【0109】

【数2】

XGFT($n+1; m_1, \dots, m_{n+1}; w_1, \dots, w_{n+1}$)は、新しい最上位

レベルにおける $\prod_{i=0}^{n+1} w_i$ の付加的スイッチと、XGFT($n+1; m_1, \dots, m_n;$

w_1, \dots, w_n)の m_n の別個のコピーを接続することによって再帰的に構築される。この定義を用いることにより、以下の特性が適用される： $n > 0$ の場合、 $n+1$ レベルである各々のXGFTは、 m_n サブツリーから構成される（すなわち、1レベルであるXGFTにおける n レベルの各サブツリーの場合（ $1 > n$ ）、 $n+1$ レベルである1つのイミディエイト・スーパーツリーが存在する。これは m_n の n レベルサブツリーを接続するものである）。同様に、ネットワーク接続性の観点から、XGFTにおける各々のサブツリーを別個のXGFTと見なすことができ、サブツリーにおける最上位レベルのスイッチがそのイミディエイト・スーパーツリーに向かうそのスカイラインを定義する。

【0110】

一実施形態に従うと、 $n+1$ レベルであるXGFTにおける各々のスイッチは固有の n タプル（ $1, x_1, x_2, \dots, x_n$ ）によって表わすことができる。左端のタプル値

(1) はツリーが位置するレベルを表わしており、残りの値 (x_1, x_2, \dots, x_n) は、他のスイッチに対応するツリーにおけるスイッチの位置を表わしている。特に、レベル 1 におけるスイッチ A ($1, a_1, \dots, a_1, \dots, a_n$) は、 $i = 1 + 1$ である場合を除いて、すべての値について $a_i = b_i$ であるとき、かつそのときに限り、レベル $1 + 1$, ($1 + 1, b_1, \dots, b_1, b_{1+1}, \dots, b_n$) におけるスイッチ B に接続される。

【0111】

図 16 は、一実施形態に従ったスイッチタブルを示す。より具体的には、当該図は、例示的なファットツリーである X G F T (4; 2, 2, 2, 2; 2, 2, 2, 1) のために実現された Open SM のファットツリールーティングアルゴリズムによって割付けられるようなスイッチタブルを示している。ファットツリー 1400 は、スイッチ 1401 ~ 1408、1411 ~ 1418、1421 ~ 1428 および 1431 ~ 1438 を含み得る。ファットツリーが (リーフレベルにおける列 3 まで、ルートレベルにおける列 0 として印付けされた) $n = 4$ のスイッチレベルを有しているので、ファットツリーは、各々が $n = n - 1 = 3$ スwitchレベルである $m_1 = 2$ の第 1 レベルサブツリーで構成されている。これは、図において、レベル 1 から 3 までのスイッチを囲んでいる破線によって規定される 2 つのボックスによって示されている。各々の第 1 レベルのサブツリーが 0 または 1 の識別子を受取る。第 1 レベルのサブツリーの各々は、各々がリーフスイッチを上回っている $n = n - 1 = 2$ のスイッチレベルである $m_2 = 2$ の第 2 レベルのサブツリーから構成されている。これは、図において、レベル 2 から 3 までのスイッチを囲んでいる点線によって規定される 4 つのボックスによって示されている。各々の第 2 レベルのサブツリーは 0 または 1 の識別子を受取る。同様に、リーフスイッチの各々は、図において、鎖線によって規定される 8 つのボックスによって示されるサブツリーと見なすこともできる。これらのサブツリーの各々は 0 または 1 の識別子を受取る。

【0112】

一実施形態に従うと、図に例示されているように、4 つの数字のタブルなどのタブルは、さまざまなスイッチに割当てることができ、タブルの各々の数字は、タブルにおける各々の値の位置についての特定のサブツリー対応を示している。たとえば、(スイッチ 1 __ 3 と参照され得る) スwitch 1413 は、レベル 1 におけるその位置と 0 番目の第 1 レベルのサブツリーとを表わしているタブル 1 . 0 . 1 . 1 に割当てることができる。

【0113】

ライブマイグレーションの文脈における F T r e e M i n R C を用いたファットツリー認識型の最小再構成

一実施形態に従うと、スイッチタブルは、トポロジにおけるサブツリーに対応するスイッチの位置についての情報を符号化する。F T r e e M i n R C は、ライブ V M マイグレーションの場合における迅速な再構成を可能にするためにこの情報を用いることができる。タブル情報は、V M がマイグレートされたときに S M によって再構成される必要のあるスイッチの数が最も少ないスカイランを発見するために用いることができる。特に、V M がファットツリートポロジにおける 2 つのハイパーバイザ間でマイグレートされると、更新される必要のある最小数のスイッチを表わしているスカイラインは、マイグレーションに参与しているすべてのサブツリーのうちすべての最上位レベルのスイッチによって形成されている。

【0114】

一実施形態に従うと、V M がライブマイグレートされると、スイッチ印付けメカニズムを両方のリーフスイッチから開始することができる。この場合、送信元ハイパーバイザと宛先ハイパーバイザとが接続され、スイッチのタブルを比較する。タブル同士が一致する場合、メカニズムは、V M がリーフスイッチ内でマイグレートされていると判断することができる。これにより、再構成のために対応するリーフスイッチだけに印が付けられる。しかしながら、タブルが一致していなければ、送信元リーフスイッチおよび宛先リーフスイッチの両方からの上りリンクがトレースされる。1 レベル上に位置するスイッチは、リ

ーフレベルのサブツリーが接続されているイミディエイト・スーパーツリーのうち最上位レベルのスイッチであり、ツリーを下方へと横切る際にリーフスイッチに到達する前に生じる可能性のある唯一のホップである。次いで、当該メカニズムは、送信元リーフスイッチタプルおよび宛先リーフスイッチタプルを新しくトレースされたスイッチと比較することができ、その時点のレベルを反映させるためにタプル値を調整した後、その時点のツリーのサブツリーに対応する値がワイルドカードにされる。さらに、（対応するサブツリーのための最上位レベルのスイッチである）トレースされたスイッチは更新されるべく印付けされ、送信元スイッチタプルおよび宛先スイッチタプルの両方からの比較がトレースされたすべてのスイッチのタプルと一致する場合、トレースが停止される。他の場合には、メカニズムが両端から共通の先祖スイッチを特定するまで、同じ手順が繰返される。最悪の場合、ファットツリートポロジーのルートスイッチに到達した後、メカニズムを停止することができる。すべての上り経路のトレースがリーフレベルから開始されており、かつ、連続したサブツリーのスカイラインスイッチに印付けされているので、メカニズムがマイグレーションによって影響される最上位のサブツリーに到達した場合、当該メカニズムは、その途中で、下位レベルスイッチに向かう潜在的なトラフィックゲートウェイであるすべてのスイッチや、ライブマイグレーションに関与するハイパーバイザを既に選択してしまっている。これにより、当該メカニズムは、ネットワークのうちライブマイグレーションによって影響を受けた部分のスカイラインを形成するすべてのスイッチに印を付けた。

10

【 0 1 1 5 】

20

ー実施形態に従うと、スイッチ印付けメカニズムは、物理的接続の観点から更新される必要のある、最小数のスイッチを発見する。しかしながら、これらのスイッチのすべてが再構成によって影響を受けた L I D に対するルーティングアルゴリズムによって計算されたアクティブな経路を含むとは限らない可能性もある。このため、アクティブなルートを含んでいるスイッチには更新手順において優先順位が付けられる一方で、スイッチのうち二次ルートを有する残りのスイッチは後で更新することができる。

【 0 1 1 6 】

ー実施形態に従うと、ファットツリールーティングメカニズムは、常に、同じルートスイッチを介して所与の宛先にトラフィックをルーティングする。トポロジーにおいてはルートスイッチとエンドノードとの間に単一の経路だけが存在しているので、所与のエンドノードを表わすために選択されたルートスイッチが位置特定されると、エンドノードにトラフィックをルーティングするために用いられる中間スイッチを見出すことができる。アクティブなルートを発見するために、経路は、関与するハイパーバイザの送信元 L I D から宛先 L I D にまで、またはこれの逆の態様でトレースすることができる。再構成のために既に選択されていたスイッチのサブセットであるスイッチに印を付けることができ、それらのスイッチの L F T 更新に優先順位を付けることができる。その後、すべての L F T を有効に維持するために、残りの選択されたスイッチを更新することができる。

30

【 0 1 1 7 】

図 1 7 は、一実施形態に従った再構成プロセスを示す。ファットツリー 1 4 0 0 は、スイッチ 1 4 0 1 ~ 1 4 0 8、1 4 1 1 ~ 1 4 1 8、1 4 2 1 ~ 1 4 2 8 および 1 4 3 1 ~ 1 4 3 8 を含む得る。ファットツリーが（リーフレベルにおける列 3 まで、ルートレベルにおける列 0 として印付けされた） $n = 4$ スイッチレベルを有しているので、ファットツリーは、各々が $n = n - 1 = 3$ のスイッチレベルである $m_1 = 2$ の第 1 レベルのサブツリーから構成されている。これら第 1 レベルのサブツリーの各々は、各々がリーフスイッチを上回っている $n = n - 1 = 2$ のスイッチレベルである $m_2 = 2$ の第 2 レベルのサブツリーから構成されている。同様に、リーフスイッチの各々もサブツリーと見なすことができる。

40

【 0 1 1 8 】

ー実施形態に従うと、図 1 7 は、タプル 3 . 0 . 0 . 0 および 3 . 0 . 1 . 1 を備えたリーフスイッチに接続された 2 つのハイパーバイザ間で V M がマイグレートされている状

50

況を示す。これらの2つのタプルは、選択されたリーフスイッチから上方向への経路をメカニズムがトレースする際の比較についての基準として用いられる。この例においては、共通の先祖スイッチがレベル1上で発見される。レベル0はルートレベルであり、レベル3はリーフレベルである。表示されたタプル情報を有するスイッチ間のリンクはメカニズムの実行中ずっとトレースすることができるリンクであり、それらの同じスイッチすべてに更新するための印を付けることができる。強調表示された5つのスイッチ（スイッチ1431、1421、1411、1423および1434）およびそれらの間のリンクは、アクティブなルートを表わしており、それらのLFT更新に優先順位を付けることができる。

【0119】

一実施形態に従うと、ライブマイグレーションをサポートする仮想化されたデータセンタにおける最小限のオーバーヘッドに迅速な接続性を提供するために、FTReeMinRCは、スイッチに送信される必要のあるLFT更新の回数を最小限にする。

【0120】

図18は、一実施形態に従った、無損失相互接続ネットワークにおける効率的な仮想化をサポートするための方法のフローチャートである。ステップ1810において、当該方法は、1つ以上のマイクロプロセッサを含む1つ以上のコンピュータにおいて、少なくともリーフスイッチを含む1つ以上のスイッチを設けることができ、当該1つ以上のスイッチの各々は複数のポートを含み、さらに、複数のホストチャネルアダプタを設けることができ、ホストチャネルアダプタの各々は少なくとも1つの仮想機能、少なくとも1つの仮想スイッチおよび少なくとも1つの物理機能を含み、複数のホストチャネルアダプタは1つ以上のスイッチを介して相互接続されており、さらに、複数のハイパーバイザを設けることができ、複数のハイパーバイザの各々は、複数のホストチャネルアダプタのうち少なくとも1つのホストチャネルアダプタに関連付けられており、さらに、複数の仮想マシンを設けることができ、複数の仮想マシンの各々は少なくとも1つの仮想機能に関連付けられている。

【0121】

ステップ1820において、当該方法は、予めポピュレートされたローカル識別子(local identifier: LID)アーキテクチャを備えた仮想スイッチまたは動的LID割当てアーキテクチャを備えた仮想スイッチのうち1つ以上を備えた複数のホストチャネルアダプタを配置することができる。

【0122】

ステップ1830において、当該方法は、各々の仮想スイッチにLIDを割当てることができ、割当てられたLIDは関連付けられた物理機能のLIDに対応している。

【0123】

ステップ1840において、当該方法は、仮想スイッチの各々に割当てられたLIDに少なくとも基づいて1つ以上のリニアフォワーディングテーブルを計算することができ、1つ以上のLFTの各々は、1つ以上のスイッチのうちの1つのスイッチに関連付けられている。

【0124】

図19は、一実施形態に従った、無損失相互接続ネットワークにおける効率的な仮想化をサポートするための方法のフローチャートである。ステップ1910において、当該方法は、1つ以上のマイクロプロセッサを含む1つ以上のコンピュータにおいて、1つ以上のマイクロプロセッサと、1つ以上のスイッチとを設けることができ、1つ以上のスイッチは少なくともリーフスイッチを含み、1つ以上のスイッチの各々は複数のポートを含み、さらに、複数のホストチャネルアダプタを設けることができ、ホストチャネルアダプタの各々は少なくとも1つの仮想機能、少なくとも1つの仮想スイッチおよび少なくとも1つの物理機能を含み、複数のホストチャネルアダプタは1つ以上のスイッチを介して相互接続されており、さらに、複数のハイパーバイザを設けることができ、複数のハイパーバイザの各々は、複数のホストチャネルアダプタのうち少なくとも1つのホストチャネルア

10

20

30

40

50

ダブタに関連付けられており、さらに、複数の仮想マシンを設けることができ、複数の仮想マシンの各々は、少なくとも1つの仮想機能に関連付けられている。

【0125】

ステップ1920において、当該方法は、予めポピュレートされたローカル識別子(LID)アーキテクチャを備えた仮想スイッチまたは動的LID割当てアーキテクチャを備えた仮想スイッチのうち1つ以上を備えた複数のホストチャネルアダプタを配置することができる。

【0126】

ステップ1930において、当該方法は、仮想スイッチの各々に複数のpLIDのうち1つのpLIDを割当てることができる。割当てられたpLIDは関連付けられた物理機能のpLIDに対応している。

10

【0127】

ステップ1940において、当該方法は、複数の仮想マシンの各々に複数のvLIDのうち1つのvLIDを割当てることができ、LIDスペースは、複数のpLIDおよび複数のvLIDを含む。

【0128】

本発明の多くの特徴は、ハードウェア、ソフトウェア、ファームウェアまたはそれらの組合せにおいて、それらを用いて、またはそれらの支援により、実行可能である。したがって、本発明の特徴は、(たとえば、1つ以上のプロセッサを含む)処理システムを用いて実現され得る。

20

【0129】

この発明の特徴は、ここに提示された特徴のうちのいずれかを行なうように処理システムをプログラミングするために使用可能な命令を格納した記憶媒体またはコンピュータ読取り可能媒体であるコンピュータプログラム製品において、それを使用して、またはその助けを借りて実現され得る。記憶媒体は、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、フラッシュメモリ装置、磁気カードもしくは光カード、ナノシステム(分子メモリICを含む)、または、命令および/もしくはデータを格納するのに好適な任意のタイプの媒体もしくは装置を含み得るものの、それらに限定されない。

30

【0130】

この発明の特徴は、機械読取り可能媒体のうちのいずれかに格納された状態で、処理システムのハードウェアを制御するために、および処理システムがこの発明の結果を利用する他の機構とやり取りすることを可能にするために、ソフトウェアおよび/またはファームウェアに取込まれ得る。そ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、装置ドライバ、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ものの、それらに限定されない。

【0131】

この発明の特徴はまた、たとえば、特定用途向け集積回路(application specific integrated circuit: ASIC)などのハードウェアコンポーネントを使用して、ハードウェアにおいて実現されてもよい。ここに説明された機能を行なうようにハードウェアステートマシンを実現することは、関連技術の当業者には明らかであろう。

40

【0132】

加えて、この発明は、この開示の教示に従ってプログラミングされた1つ以上のプロセッサ、メモリおよび/またはコンピュータ読取り可能記憶媒体を含む、1つ以上の従来の汎用または特殊デジタルコンピュータ、コンピューティング装置、マシン、またはマイクロプロセッサを使用して都合よく実現され得る。ソフトウェア技術の当業者には明らかであるように、この開示の教示に基づいて、適切なソフトウェアコーディングが、熟練したプログラマによって容易に準備され得る。

【0133】

50

この発明のさまざまな実施形態が上述されてきたが、それらは限定のためではなく例示のために提示されたことが理解されるべきである。この発明の精神および範囲から逸脱することなく、形状および詳細のさまざまな変更を行なうことができることは、関連技術の当業者には明らかであろう。

【 0 1 3 4 】

この発明は、特定された機能およびそれらの関係の実行を示す機能的構築ブロックの助けを借りて上述されてきた。説明の便宜上、これらの機能的構築ブロックの境界は、この明細書中ではしばしば任意に規定されてきた。特定された機能およびそれらの関係が適切に実行される限り、代替的な境界を規定することができる。このため、そのようないかなる代替的な境界も、この発明の範囲および精神に含まれる。

10

【 0 1 3 5 】

この発明の前述の説明は、例示および説明のために提供されてきた。それは、網羅的であるよう、またはこの発明を開示された形態そのものに限定するよう意図されてはいない。この発明の幅および範囲は、上述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ない。多くの変更および変形が、当業者には明らかになるだろう。これらの変更および変形は、開示された特徴の関連するあらゆる組合せを含む。実施形態は、この発明の原理およびその実用的応用を最良に説明するために選択され説明されたものであり、それにより、考えられる特定の使用に適したさまざまな実施形態についての、およびさまざまな変更例を有するこの発明を、当業者が理解できるようにする。この発明の範囲は、請求項およびそれらの同等例によって定義されるよう意図されている。

20

【 図 1 】

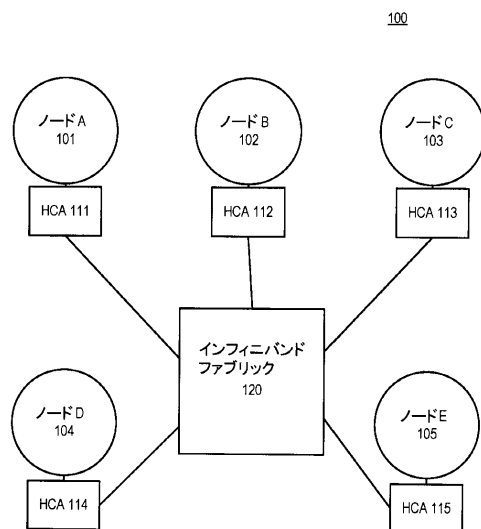


FIGURE 1

【 図 2 】

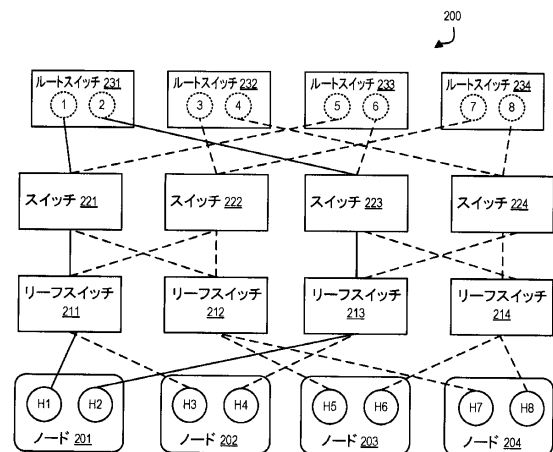


FIGURE 2

【図 3】

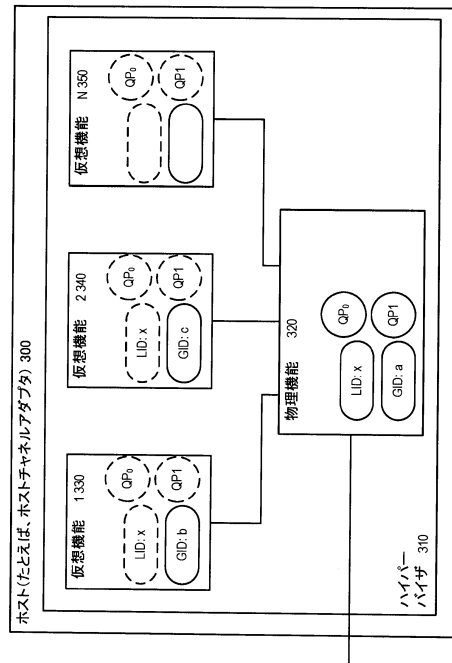


FIGURE 3

【図 4】

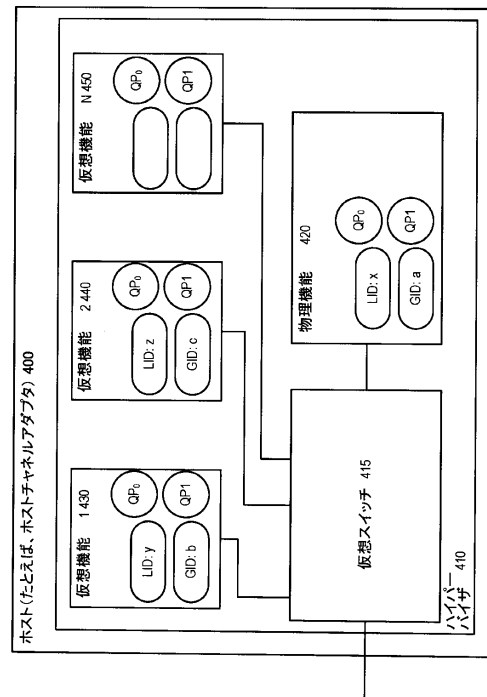


FIGURE 4

【図 5】

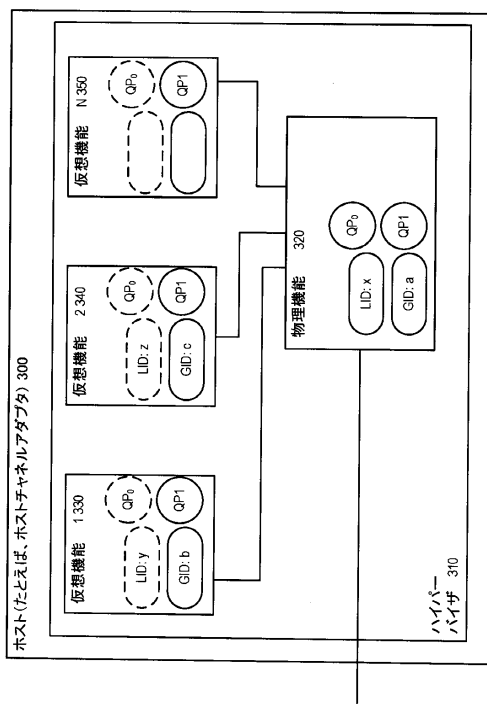


FIGURE 5

【図 6】

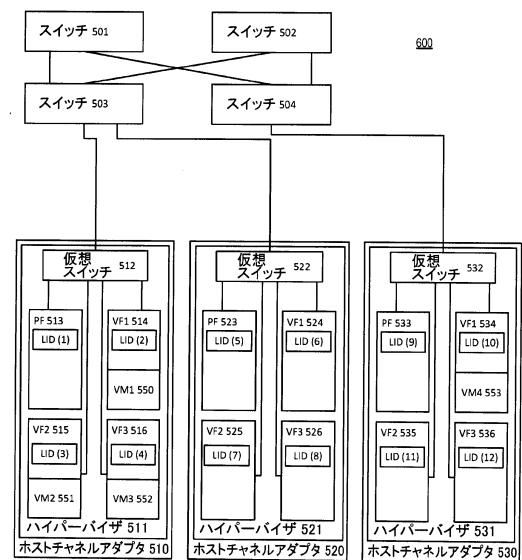


FIGURE 6

【図 7】

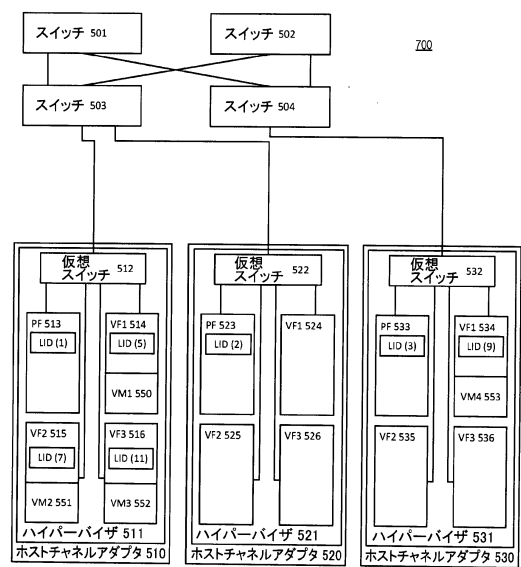


FIGURE 7

【図 8】

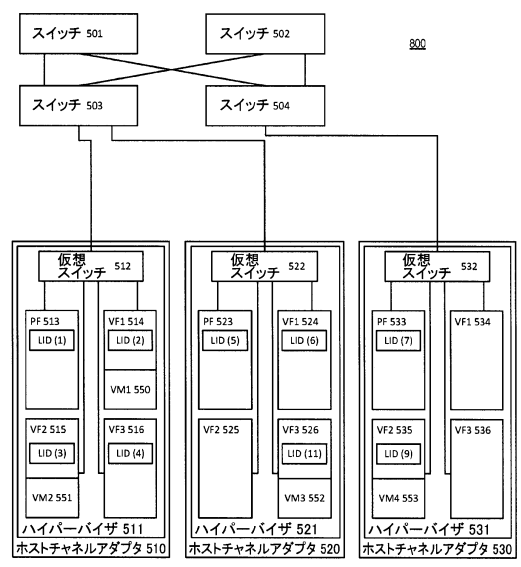


FIGURE 8

【図 9】

900	仮想レーン (VL) - 4ビット
901	リンクバージョン (LVer) - 4ビット
902	サービスレベル (SL) - 4ビット
903	LID拡張フラグ (LEXTF) - 1ビット
904	予備 (R1) - 1ビット - ゼロ
905	リンク次ヘッダ (LNH) - 2ビット
906	宛先ローカルID (DLID) - 16ビット
907	DLIDプレフィックス拡張 (DPF) - 2ビット
908	SLIDプレフィックス拡張 (SPF) - 2ビット
909	予備 (R2) - 1ビット - ゼロ
910	パケット長 (PktLen) - 11ビット
911	送信元ローカルID (SLID) - 16ビット

FIGURE 9

【図 10】

912	16ビットDLID=0によって索引付けされた エントリ0 (標準IBポート番号を含む)
...	
913	16ビットDLID=48K-1によって索引付けされた エントリ48K-1 (IBポート番号を含む)

914	18ビットDPF+DLID=64Kによって索引付け されたエントリ0 (標準IBポート番号を含む)
...	
915	18ビットDPF+DLID=256K-1によって 索引付けされたエントリ256K-1 (標準IBポート番号を含む)

FIGURE 10

【図 1 1】

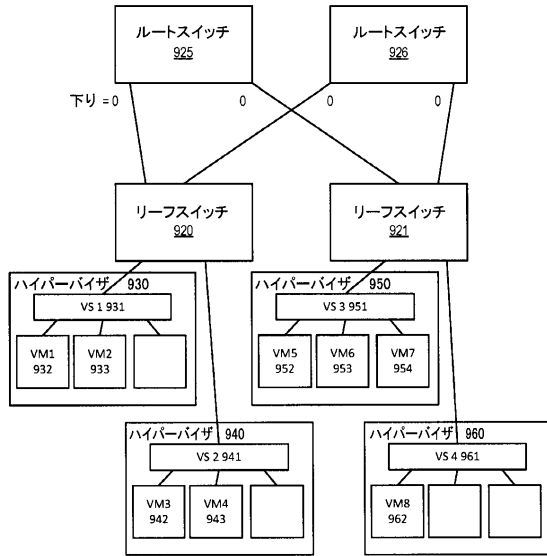


FIGURE 11

【図 1 2】

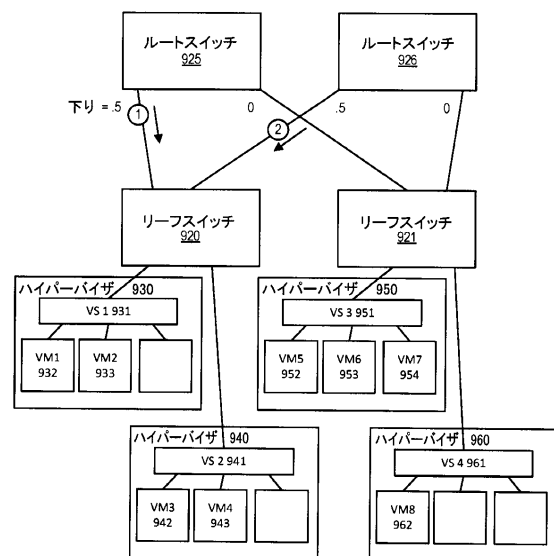


FIGURE 12

【図 1 3】

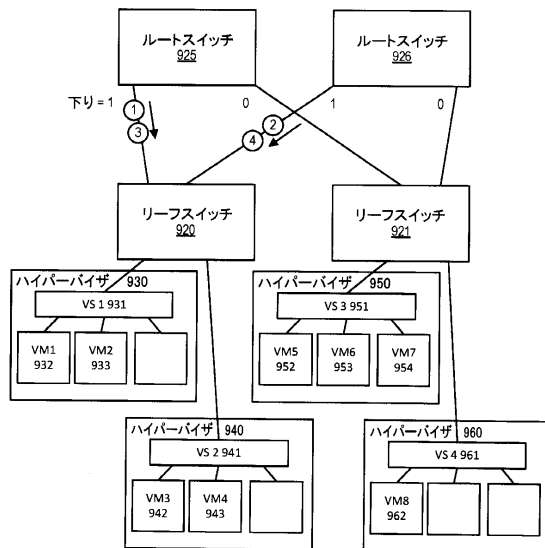


FIGURE 13

【図 1 4】

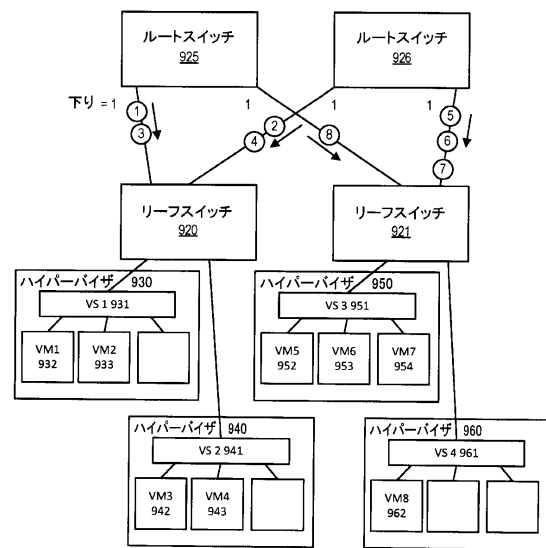


FIGURE 14

【図 15】

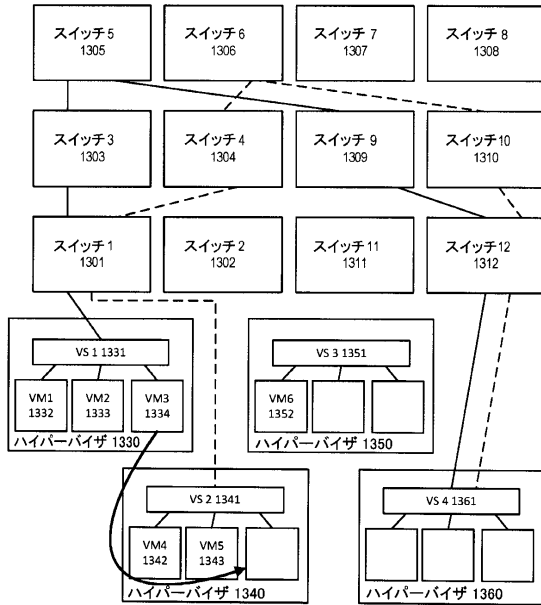


FIGURE 15

【図 16】

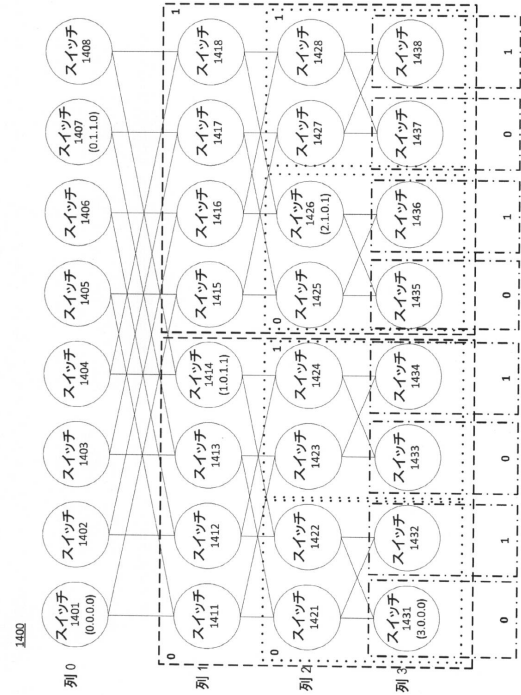


FIGURE 16

【図 17】

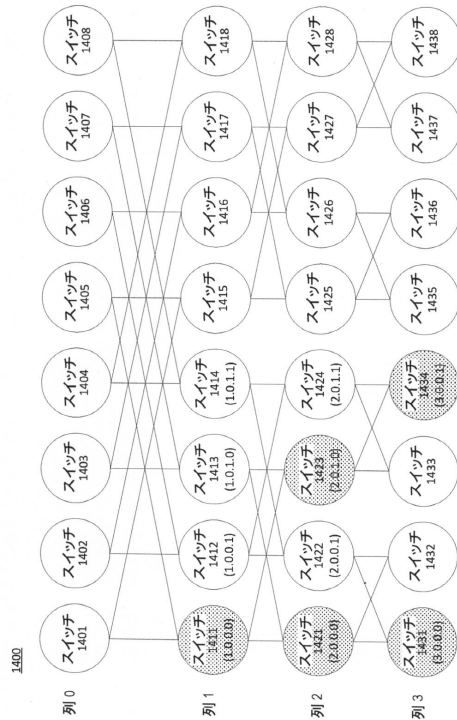


FIGURE 17

【図 18】

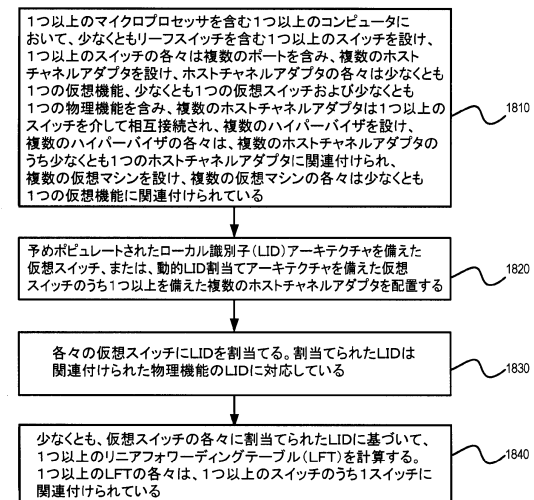


FIGURE 18

【図 19】

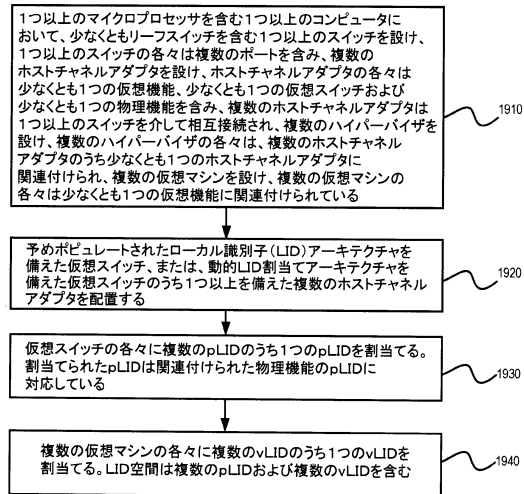


FIGURE 19

フロントページの続き

(31)優先権主張番号 62/261,103

(32)優先日 平成27年11月30日(2015.11.30)

(33)優先権主張国・地域又は機関
米国(US)

(31)優先権主張番号 15/210,595

(32)優先日 平成28年7月14日(2016.7.14)

(33)優先権主張国・地域又は機関
米国(US)

(31)優先権主張番号 15/210,599

(32)優先日 平成28年7月14日(2016.7.14)

(33)優先権主張国・地域又は機関
米国(US)

(72)発明者 ヨンセン、ビョルン・ダグ

ノルウェー、0687 オスロ、ビルベルクグレンダ、9

(72)発明者 グラン、アーヌスト・ガンナー

ノルウェー、1325 リュサケール、ピィ・オウ・ボックス・134

審査官 中川 幸洋

(56)参考文献 米国特許第06055532(US,A)

特表2015-514271(JP,A)

Evangelos Tasoulas, et al., Towards the InfiniBand SR-IOV vSwitch Architecture, 2015 IEEE International Conference on Cluster Computing, IEEE, 2015年9月8日

(58)調査した分野(Int.Cl., DB名)

H04L 12/70

G06F 9/455

G06F 9/50

H04L 12/46