

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7109527号
(P7109527)

(45)発行日 令和4年7月29日(2022.7.29)

(24)登録日 令和4年7月21日(2022.7.21)

(51)国際特許分類	F I		
G 0 6 F 15/173(2006.01)	G 0 6 F	15/173	6 8 5 M
G 0 6 F 9/455(2006.01)	G 0 6 F	9/455	1 5 0
G 0 6 F 9/50 (2006.01)	G 0 6 F	9/50	1 2 0 Z
H 0 4 L 61/00 (2022.01)	H 0 4 L	61/00	
H 0 4 L 61/10 (2022.01)	H 0 4 L	61/10	

請求項の数 8 外国語出願 (全37頁)

(21)出願番号	特願2020-212507(P2020-212507)	(73)特許権者	502303739
(22)出願日	令和2年12月22日(2020.12.22)		オラクル・インターナショナル・コーポレーション
(62)分割の表示	特願2018-501253(P2018-501253)の分割		アメリカ合衆国カリフォルニア州94065レッドウッド・シティー, オラクル・パークウェイ500
原出願日	平成28年11月18日(2016.11.18)	(74)代理人	110001195弁理士法人深見特許事務所
(65)公開番号	特開2021-73551(P2021-73551A)	(72)発明者	タソウラス, エバンジェロス
(43)公開日	令和3年5月13日(2021.5.13)		ノルウェー, 1325 リュサケール、
審査請求日	令和3年1月12日(2021.1.12)	(72)発明者	ギド, フェロツ
(31)優先権主張番号	62/259,321		ノルウェー, 1325 リュサケール、
(32)優先日	平成27年11月24日(2015.11.24)	(72)発明者	ビョルン・ダグ
(33)優先権主張国・地域又は機関	米国(US)		ノルウェー, 0687 オスロ、ビルベ
(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			

最終頁に続く

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、

複数のスイッチを含み、前記複数のスイッチは少なくとも2つのレベルに配置されており、

前記複数のスイッチの各々は、複数のリニアフォワードテーブル(linear forwarding table: LFT)のうち1つのLFTに関連付けられており、前記システムはさらに、

複数のホストチャンネルアダプタを含み、前記複数のホストチャンネルアダプタは前記複数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ンネルアダプタのうち少なくとも1つのホストチャン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、それぞれのハイパーバイザを介して前記複数のホストチャンネルアダプタのうち1つのホストチャンネルアダプタに接続しており、

前記複数の仮想マシンのうち1つの仮想マシンは、第1のホストチャンネルアダプタにおける第1のハイパーバイザから第2のホストチャンネルアダプタにおける第2のハイパーバ

イザに対するライブマイグレーションを実行し、前記ライブマイグレーション中に、前記 1 つの仮想マシンのローカル識別子 (local identifier: L I D) が更新され、

前記第 1 のホストチャンネルアダプタから前記第 2 のホストチャンネルアダプタに対する前記 1 つの仮想マシンのマイグレーションの結果、前記複数の L F T のセットが更新され、前記複数の L F T のセットは前記第 1 のホストチャンネルアダプタおよび前記第 2 のホストチャンネルアダプタの最も近い共通の先祖スイッチの判定に基づいて規定されており、

前記仮想マシンの各々には重みパラメータが割当てられており、

前記重みパラメータの各々は前記複数の L F T を計算するのに用いられ、

前記複数の L F T を計算することは、

前記仮想マシンのために、利用可能な上りポート群の間で最少累積の下り重みに基づいてスイッチにおける下りポートを選択することと、

前記仮想マシンの割当てられた前記重みパラメータによって、前記選択されたポートについて、累積された下り重みを増やすこととを含む、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、複数のスイッチと、複数のホストチャンネルアダプタと、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、

前記複数のスイッチは少なくとも 2 つのレベルに配置されており、前記複数のスイッチの各々は、複数のリニアフォワーディングテーブル (L F T) のうち 1 つの L F T に関連付けられており、

前記複数のホストチャンネルアダプタは前記複数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ンネルアダプタのうち少なくとも 1 つのホストチャン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、それぞれのハイパーバイザを介して前記複数のホストチャンネルアダプタのうち 1 つのホストチャンネルアダプタに接続しており、前記方法はさらに、

第 1 のホストチャンネルアダプタにおける第 1 のハイパーバイザから第 2 のホストチャンネルアダプタにおける第 2 のハイパーバイザに対する前記複数の仮想マシンのうち 1 つの仮想マシンのライブマイグレーションを実行するステップを含み、前記ライブマイグレーション中に、前記 1 つの仮想マシンのローカル識別子 (L I D) が更新され、前記方法はさらに、

前記第 1 のホストチャンネルアダプタから前記第 2 のホストチャンネルアダプタに対する前記 1 つの仮想マシンの前記ライブマイグレーションの結果、前記複数の L F T のセットを更新するステップを含み、前記複数の L F T のセットは前記第 1 のホストチャンネルアダプタおよび前記第 2 のホストチャンネルアダプタの最も近い共通の先祖スイッチの判定に基づいて規定されており、

前記仮想マシンの各々には重みパラメータが割当てられており、

前記重みパラメータの各々は前記複数の L F T を計算するのに用いられ、

前記複数の L F T を計算することは、

前記仮想マシンのために、利用可能な上りポート群の間で最少累積の下り重みに基づいてスイッチにおける下りポートを選択することと、

前記仮想マシンの割当てられた前記重みパラメータによって、前記選択されたポートについて、累積された下り重みを増やすこととを含む、方法。

【請求項 3】

前記第 1 のホストチャンネルアダプタおよび前記第 2 のホストチャンネルアダプタの前記最も近い共通の先祖スイッチの前記判定は、

前記第 1 のホストチャンネルアダプタから前記第 2 のホストチャンネルアダプタに向かって上方向に第 1 の経路をトレースするステップと、

10

20

30

40

50

前記第 2 のホストチャネルアダプタから前記第 1 のホストチャネルアダプタに向かって上方向に第 2 の経路を同時にトレースするステップと、

印付けメカニズムによって、前記第 1 の経路と前記第 2 の経路とが合流する地点において前記最も近い共通の先祖スイッチに印付けするステップとを含む、請求項 2 に記載の方法。

【請求項 4】

前記第 1 の経路および前記第 2 の経路に沿った各々のスイッチは前記印付けメカニズムによって印付けされる、請求項 3 に記載の方法。

【請求項 5】

前記複数の L F T のセットは、印付けされた前記最も近い共通の先祖スイッチにおける各々の L F T と、前記第 1 の経路および前記第 2 の経路に沿って印付けされた各々のスイッチの各々の L F T とを含む、請求項 4 に記載の方法。

10

【請求項 6】

前記重みパラメータの各々は、仮想スイッチとリーフスイッチとの間の各リンクのトラフィックの割合を反映している、請求項 2 ~ 5 のいずれか 1 項に記載の方法。

【請求項 7】

前記複数の L F T の各々は宛先ノードから計算が開始される、請求項 6 に記載の方法。

【請求項 8】

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 ~ 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、著作権保有者は、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。

30

【0002】

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、SR - IOV 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：IB）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。例示的な方法は、1 つ以上のマイクロプロセッサを含

50

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、複数のホストチャンネルアダプタを設けることができる。複数のホストチャン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。複数のホストチャン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ンネルアダプタのうち少なくとも1つのホストチャン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、予めポ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ン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ドイングテーブル(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。

10

【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、複数のホストチャンネルアダプタを設けることができ、ホストチャン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、複数のホストチャンネルアダプタは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、複数のホストチャンネルアダプタのうち少なくとも1つのホストチャン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つ以上を備えた複数のホストチャン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。

20

30

【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。

40

【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。

50

【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。 20

【図 16】一実施形態に従ったスイッチテーブルを示す図である。

【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 (I B) 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 50

には明らかであるだろう。

【 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：RDMA）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。RDMAはOS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：SR-IOV）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：vSwitch）SR-IOV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ジーに依存しない動的な再構成メカニズムを提供することができる。

10

【 0 0 1 3 】

一実施形態に従うと、さらには、vSwitch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：VM）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（IB）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：HPC）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：IBA）は、2層トポロジー分割をサポートする。低層では、IB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。

40

【 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：SM）があり得る。サブネットマネージャは、IB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（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行なう役割を果たし得る。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。

【 0 0 1 8 】

50

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、サブネット管理パケット (subnet management packet : SMP) 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 : SMA) と交換する。サブネット管理エ

ージェントは、すべてのIBサブネットデバイス上に存在する。SMP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、SMAから通知を受信することができる。

【0019】

一実施形態によれば、IB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LFTに基づき得る。LFTは、使用中のルーティングメカニズムに従って、SM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ンネルアダプタ (Host Channel Adapter : HCA) ポートおよびスイッチが、ローカル識別子 (LID) を使用してアドレス指定される。LFTにおける各エントリは、宛先LID (destination LID : DLID) と出力ポートとからなる。テーブルにおけるLID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DLID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 (LID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。

10

【0020】

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ー変化を検出し、それに応じてネットワークを再構成する。

20

【0021】

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 (LID) 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49151個のユニキャストLID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LIDの他に、各IBデバイスは、64ビットのグローバル一意識別子 (global unique identifier : GUID) を有し得る。GUIDは、IBレイヤー3 (L3) アドレスであるグローバル識別子 (global identifier : GID) を形成するために使用され得る。

30

【0022】

SM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、トポロジーが変化するたびに、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、SMは、トポロジー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SMが受信した場合、SM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。

40

【0023】

たとえば、SM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ーが変化する場合に、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。

【0024】

一実施形態によれば、IB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

50

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HCA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、SM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。SM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。SM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LID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。

10

【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(QP0およびQP1)を除き、キューペア(Queue Pair: QP)およびエンドツーエンドコンテキスト(End-to-End context: EEC)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IB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HCAポートまたはスイッチに到着すると、そのP__キー値を、SM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A101~E105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A101~E105)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A101~E105)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。

20

【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CPU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SAN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: SR-IOV)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: HPC)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HPC(virtual HPC: vHPC)クラスタに対応し、必要な性能

30

を提供することができる。

40

【0028】

しかしながら、インフィニバンド(IB)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(VM)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。IB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、HPC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】

IBデバイスをVM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SR-IOVを利用することによるものである。しかしながら、SR-IOVを用いてIBホストチ

50

チャンネルアダプタ (HCA) に割当てられた VM 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 IB 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 LID、GUID および GID) 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 ($0 \times 0001 - 0 \times BFFF$) 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】

10

20

30

40

50

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 2 3 1 はスイッチポート 1 ~ 2 を有し、ルートスイッチ 2 3 2 はスイッチポート 3 ~ 4 を有し、ルートスイッチ 2 3 3 はスイッチポート 5 ~ 6 を有し、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 5 】

一実施形態によれば、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一つである。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。

10

【 0 0 3 6 】

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、ひいては対応するルートを各スイッチポートに割当てる。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。

【 0 0 3 7 】

20

一実施形態に従うと、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はユーザ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。

【 0 0 3 8 】

30

一実施形態に従うと、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 、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。

【 0 0 3 9 】

また、 k - a - r - y - n - t - r - 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 $2k$ ポートを持っている。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とも可能にすることによって、 k - a - r - y - n - t - r - e を拡張させる。P 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。

40

【 0 0 4 0 】

50

入出力 (Input/Output : I / O) 仮想化

一実施形態に従うと、I / O 仮想化 (I/O Virtualization : IOV) は、基礎をなす物理リソースに仮想マシン (VM) がアクセスすることを可能にすることによって、I / O 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、IOV 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 IOV が所望される。IOV 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 - IOV) などのさまざまなタイプの IOV 技術が存在し得る。

【0042】

一実施形態に従うと、あるタイプの IOV 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な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。

【0043】

一実施形態に従うと、別のタイプの IOV 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】

一実施形態に従うと、シングルルート IOV (Single Root IOV : SR - IOV) 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 - IOV は、純粋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 - IOV 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 IOV 技術の中でも、SR - IOV 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 - IOV 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。

【0045】

SR - IOV 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 - IOV 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

10

20

30

40

50

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 P C I e 機能であるのに対して、V F は軽量の P C I e 機能である。各々の V F はそれ自体のベースアドレス (base address : B A R) を有しており、固有のリクエスト I D が割当てられている。固有のリクエスト I D は、I / O メモリ管理ユニット (I / O memory

management unit : I O M M U) がさまざまな V F への / 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。I O M M U はまた、メモリを適用して、P F と V F との間の変換を中断する。

【 0 0 4 6 】

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、V M 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、V M がソースハイパーバイザにおいて中断され、V M 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。

10

【 0 0 4 7 】

しかしながら、S R - I O V 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、V M に割当てられた S R - I O V V F 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 V F 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 S R - I O V 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、S R - I O V 共有型ポートモデルにおいては、V M 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、S M 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。

20

【 0 0 4 8 】

インフィニバンド S R - I O V アーキテクチャ - 共有ポート

30

さまざまなタイプの S R - I O V 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。

【 0 0 4 9 】

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 3 0 0 (たとえばホストチャネルアダプタ) はハイパーバイザ 3 1 0 と対話し得る。ハイパーバイザ 3 1 0 は、さまざまな仮想機能 3 3 0、3 4 0 および 3 5 0 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 3 1 0 によって処理することができる。

【 0 0 5 0 】

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 H C A) は、物理機能 3 2 0 と仮想機能 3 3 0、3 5 0、3 5 0 との間において単一の共有 L I D および共有キュー対 (Queue Pair : Q P) 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 G I D を有し得る。

40

【 0 0 5 1 】

図 3 に示されるように、一実施形態に従うと、さまざまな G I D 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 Q P 0 および Q P 1 (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 Q P は V F にも同様にエクスポーズされるが、V F は Q P 0 を使用することが許可されておらず (V F から Q P 0 に向かって入来するすべての S M P が廃棄され)

50

、 Q P 1 は、 P F が所有する実際の Q P 1 のプロキシとして機能し得る。

【 0 0 5 2 】

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) V M の数によって制限されることのない高度にスケラ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 0 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。

【 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 4 0 0 (たとえばホストチャネルアダプタ) はハイパーバイザ 4 1 0 と対話することができ、当該ハイパーバイザ 4 1 0 は、さまざまな仮想機能 4 3 0、4 4 0 および 4 5 0 をいくつかの仮想マシンに割当てることができる。同様に、物理機能はハイパーバイザ 4 1 0 によって処理することができる。仮想スイッチ 4 1 5 もハイパーバイザ 4 0 1 によって処理することができる。

【 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 4 3 0、4 4 0、4 5 0 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 については、 H C A 4 0 0 は、仮想スイッチ 4 1 5 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 4 1 0 は P F 4 2 0 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 4 9 1 5 1 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。

【 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 3 0 0 (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 3 3 0、3 4 0 および 3 5 0 をいくつかの仮想マシンに割当てることができるハイパーバイザ 4 1 0 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 3 1 0 によって処理すること

10

20

30

40

50

ができる。

【0058】

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vPort概念は緩やかに定義されており(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、SRIOV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。

10

【0059】

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

一実施形態に従うと、本開示は、LIDが予めポピュレートされた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が予めポ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1550はハイパーバイザ511によって仮想機能1514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2551を仮想機能2515に割当て、仮想マシン3552を仮想機能3516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4553を仮想機能1534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が予めポピュレートされた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514~5

50

16、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、予めポピュレートされている。たとえば、物理機能513はLID1が予めポピュレートされており、仮想機能1534はLID10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、LIDはSR-IOV vSwitch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。VFのすべてがネットワークにおけるVM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VFには、図5に示されるようにLID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HCA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IBサブネットに接続され得る。

10

【0065】

一実施形態に従うと、LIDが予めポピュレートされたvSwitch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LIDをPFを介して消費し、各々の追加のVFごとに1つ以上のLIDを消費することができる。IB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VF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VM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IB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LID（16個の仮想機能ごとに1つのLIDと、物理機能のために1つのLID）を消費する。このようなIB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LIDの数によって規定されており、（49151個の利用可能なLIDをハイパーバイザごとに17個のLIDで割って得られる）2891であり、VM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VFを掛けて得られる）46256である（実質的には、IB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SMノードが同様にLID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、vSwitchが、LIDをPFと共有することができるので、付加的なLIDを占有する必要がないことに留意されたい。

20

【0066】

一実施形態に従うと、LIDが予めポピュレートされたvSwitch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LIDについて通信経路が計算される。新しいVM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LID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、VM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。

30

【0067】

一実施形態に従うと、LIDが予めポピュレートされたvSwitch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VM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、LIDを連続的にすることを必要とするLMC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LIDマスク制御ライク（LID-Mask-Control-like：LMCライク）な特徴を用いることを可能にする。VMをマイグレートしてその関連するLIDを宛先に送達する必要がある場合、不連続なLIDを自由に使用できることは特に有用となる。

40

【0068】

一実施形態に従うと、LIDが予めポピュレートされたvSwitch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、SR-IOV vSwitch対応のサブ

50

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。

【0069】

インフィニバンドSR - IOV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。

【0070】

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501～504は、ネットワーク切替環境700（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、ホストチャンネルアダプタ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。ホストチャンネルアダプタ510、520および530は、さらに、ハイパーバイザ511、521および531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、ホストチャン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、ホストチャン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513、523および533を介してホストチャンネルアダプタにアクセスすることができる。

【0071】

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700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ーワーディングテーブルを設定するのに用いられる。

【0072】

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および531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ンネルアダプタ（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 510、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。

【0073】

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513、523および533が動的に割当てられており、物理機能513がL I D1を受取り、物理機能523がL I D2を受取り、物理機能533がL I D3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1 550がアクティブであり、仮想機能1 514に関連付けられているので、仮想機能514にはL I D5が割当てられ得る。同様に、仮想機能2 515、仮想機能3 516および仮想機能1 534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 D7が仮想機能2 515に割当てられ、L I D11が仮想機能3 516に割当てられ、L I D9が仮想機能1 534に割当てられている。L

10

20

30

40

50

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 】

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。

10

【 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。

20

【 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。

30

40

【 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

50

計に対する制限はなく、VFの数は、ユニキャストLID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VM上でVF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VFを備えることにより、ユニキャストLID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。

【0080】

インフィニバンドSR-IOVアーキテクチャモデル-動的LID割当てがなされかつLIDが予めポピュレートされたvSwitch

図8は、一実施形態に従った、動的LID割当てがなされてLIDが予めポピュレートされたvSwitchを備えた例示的な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800（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。ホストチャネルアダプタ510、520および530は、それぞれ、さらに、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1550は、ハイパーバイザ511によって仮想機能1514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2551を仮想機能2515に割当てることができる。ハイパーバイザ521は、仮想マシン3552を仮想機能3526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4553を仮想機能2535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。

【0081】

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。これらいくつかのポートは、ネットワーク切替環境800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ーワーディングテーブルを設定するのに用いられる。

【0082】

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。

【0083】

一実施形態に従うと、本開示は、動的LID割当てがなされLIDが予めポピュレートされたハイブリッド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7を参照すると、ハイパーバイザ511には、予めポピュレートされたLIDアーキテクチャを備えたvSwitchが配置され得るとともに、ハイパーバイザ521には、LIDが予めポピュレートされて動的LID割当てがなされたvSwitchが配置され得る。ハイパーバイザ531には、動的LID割当てがなされたvSwitchが配置され得る。このため、物理機能513および仮想機能514~516には、それらのLID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LIDが割当てられている）。物理機能523および仮想機能1524にはそれらのLIDが予めポピュレートされ得るとともに、仮想機能2525および仮想機能3526にはそれらのLIDが動的に割当てられている（すなわち

10

20

30

40

50

、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、または組み合わせられて）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） $0 < =$ 予めポピュレートされた $V F < =$ 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。

10

【 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。

20

【 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み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み合わせること；である。

30

【 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 トポロジー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、VM 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ーの 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ーを認識することができなくなっている。

40

【 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 ビットの数を、たとえば 2 4 ビットまたは 3 2 ビットに増やすことにより、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

50

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くなり、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 1 8 ビットに拡張する (L I D スペースを 4 倍にする) 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。

10

【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 (4 8 K ユニキャスト L I D および 1 6 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、1 9 2 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。

20

【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) 9 0 0 は 4 ビットを含み、リンクバージョン (link version : L v e r) 9 0 1 は 4 ビットを含み、サービスレベル (service level : S L) 9 0 2 は 4 ビットを含み、L I D 拡張フラグ (LID extension flag : L E X T F) 9 0 3 は 1 ビットを含み、第 1 の予備ビット (R 1) 9 0 4 は 1 ビットを含み、リンク次ヘッダ (link next header : L N H) 9 0 5 は 2 ビットを含み、宛先ローカル I D (destination local ID : D L I D) 9 0 6 は 1 6 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) 9 0 7 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) 9 0 8 は 2 ビットを含み、第 2 の予備ビット (R 2) 9 0 9 は 1 ビットを含み、パケット長 (packet length : P k t L e n) 9 1 0 は 1 1 ビットを含み、送信元ローカル I D (source local ID : S L I D) 9 1 1 は 1 6 ビットを含む。一実施形態に従うと、両方の予備ビット 9 0 4 および 9 0 9 はゼロに設定することができる。

30

【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 9 0 6 および送信元ローカル I D 9 0 8 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 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 9 0 7 および 9 0 8 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。

40

【0092】

図 1 0 は、一実施形態に従った、2 つの例示的なリニアフォーワーディングテーブルを示す。図 1 0 に示されるように、リニアフォーワーディングテーブル 9 1 6 は p L I D に関連付けられたフォーワーディングテーブルである。L F T は、エントリ 9 1 2 (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 9 1 3 (D L I D = 4 8 K - 1 によっ

50

て索引付けされたエントリ 48K - 1) にわたっている。この場合、LFTにおける各エントリは、規格 16ビットDLIDによって索引付けされ、標準IB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vLID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。LFTは、エントリ 914 (18ビットDPF + DLID = 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ビットDPF + DLID = 256K - 1) によって索引付けされたエントリ 256K - 1) にわたっている。この場合、各エントリは、拡張された 18ビットDPF + DLIDによって索引付けされ、標準IBポート番号を含んでいる。

【0093】

一実施形態に従うと、軽量のvSwitch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VMと共にLIDをマイグレートすることができるvSwitchアーキテクチャは、LIDが変化するであろう共有のLID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アップする。他方で、共有のLIDスキームは、LIDスペースに対して十分にスケールアップする。ハイブリッドvSwitch + 共有型vPortモデルは、SMがサブネットにおける利用可能なSR-IOV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VFが専用のLIDを受取り得る一方で、他のものはそれらのGIDに基づいて共有LIDの態様でルーティングされている。VM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シ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VM(たとえばサーバ)には専用のLID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VMはLIDを共有することができる。

【0094】

vSwitch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vSwitch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、vSwitchは、vSwitch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。vSwitchが識別されると、ルーティング機能は、各VM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VM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LFTを生成することができる。各VMはそれ自体のアドレスを有しており、このため、各VMは、同じvSwitchに付与された他のVM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vSwitchに向かうとともに各々が特定のVM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、VM分配がvSwitchの間で均一でない場合、より多くのVMを備えたvSwitch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、vSwitch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vSwitchに付与されたすべてのVM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vSwitch - vSwitchの対の間に経路を生成して、対応するvSwitchに割当てられるのと同じ経路を備えたVM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LID割当てスキームと動的LID割当てスキームとがあれば、各々のvSwitchは、SR-IOVアーキテクチャにおけるPFによって定義されたLIDを有する。vSwitchについてのこれらのPF LIDは、ルーティングの第1段階でLFTを生成するために用いることができ、第2段階では、VMのLIDを生成されたLFTに追加することができる。予めポピュレートされたLIDスキームにおいては、VF LIDへのエントリは対応するv

10

20

30

40

50

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、それらの間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、トポロジーにおけるリンクは、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、トポロジー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ー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ー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 】

10

20

30

40

50

```

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:          $s.LFT[\text{vm.LID}] \leftarrow v.port$ 
10:        ROUTEDOWNGOINGBYGOINGUP( $s, \text{vm}$ )
11:      end for
12:    end for
13:  end for
14: end procedure
15: procedure ROUTEDOWNGOINGBYGOINGUP( $s, \text{vm}$ )
16:   $p \leftarrow \text{GETLEASTLOADEDPORT}(s.UpGroups[])$ 
17:   $rSwitch \leftarrow p.Switch$ 
18:   $rSwitch.LFT[\text{vm.LID}] \leftarrow p$ 
19:   $p.Dwn += \text{vm.weight}$ 
20:  ROUTEUPGOINGBYGOINGDOWN( $s, \text{vm}$ )
21:  ROUTEDOWNGOINGBYGOINGUP( $rSwitch, \text{vm}$ )
22: end procedure
23: procedure ROUTEUPGOINGBYGOINGDOWN( $s, \text{vm}$ )
24:   for all  $g \in s.DownGroups[]$  do
25:     skip g if the LFT(vm.LID) is part of this group
26:      $p \leftarrow \text{GETLEASTLOADEDPORT}(g)$ 
27:      $rSwitch \leftarrow p.Switch$ 
28:      $rSwitch.LFT[\text{vm.LID}] \leftarrow p$ 
29:      $p.Up += \text{vm.weight}$ 
30:     ROUTEUPGOINGBYGOINGDOWN( $rSwitch, \text{vm}$ )
31:   end for
32: end procedure

```

【 0 0 9 9 】

一実施形態に従うと、`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

10

20

30

40

50

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 `vSwitch FatTree` 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、`vSwitch FatTree` 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、`vSwitch FatTree` 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フィ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/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】

10

20

30

40

50

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タについて検討する。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、トポロジータに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り返す。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タにもかかわらず、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ITRC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、ITRC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。

10

【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。

20

【0108】

ファットツリーにおけるサブツリーおよびスイッチダブル

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。

30

【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。

40

【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

50

タプル $(l, x_1, x_2, \dots, x_n)$ によって表わすことができる。左端のタプル値 (l) 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 l におけるスイッチ $A(l, a_1, \dots, a_l, \dots, a_n)$ は、 $i = l + 1$ 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 $l + 1$, $(l + 1, b_1, \dots, b_l, b_{l+1}, \dots, b_n)$ 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プルを示す。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 $XGFT(4; 2, 2, 2, 2; 2, 2, 2, 1)$ のために実現された $OpenSM$ 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プルを示している。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つの数字のタプルなどのタプルは、さまざまなスイッチに割当てることができ、タプルの各々の数字は、タプ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witch1413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プル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FTreeMinRC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、スイッチタプルは、トポロジー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。FTreeMinRCは、ライブVM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。タプル情報は、VMがマイグレートされたときにSM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、VMがファットツリートポロジーにおける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、VM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、スイッチのタプルを比較する。タプル同士が一致する場合、メカニズムは、VM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、タプ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

10

20

30

40

50

イ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ダブルおよび宛先リーフスイッチダブ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ダブル値を調整した後、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、送信元スイッチダブルおよび宛先スイッチダブ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ダブルと一致する場合、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。

【0115】

一実施形態に従うと、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LID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もつ残りのスイッチは後で更新することができる。

【0116】

一実施形態に従うと、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LIDから宛先LID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LFT更新に優先順位を付けることができる。その後、すべてのLFT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。

【0117】

図17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。

【0118】

一実施形態に従うと、図17は、ダブル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VMがマイグレートされている状

10

20

30

40

50

況を示す。これらの2つのタブルは、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。表示されたタブ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、FreeMinRC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一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。

【0134】

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

【0135】

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

30

40

50

【図面】

【図 1】

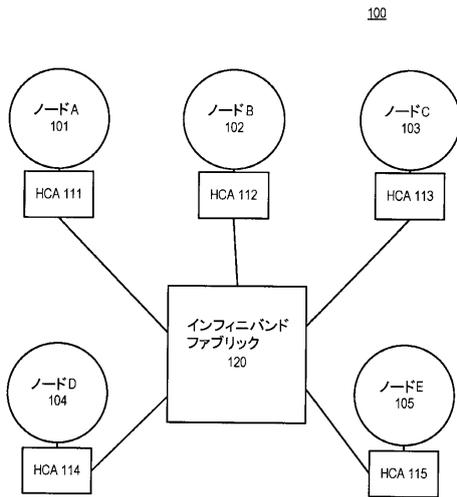


FIGURE 1

【図 2】

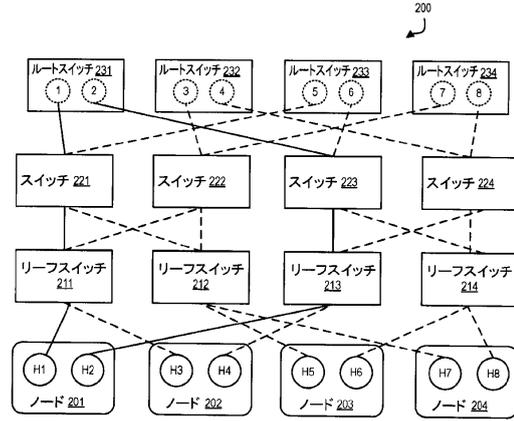


FIGURE 2

【図 3】

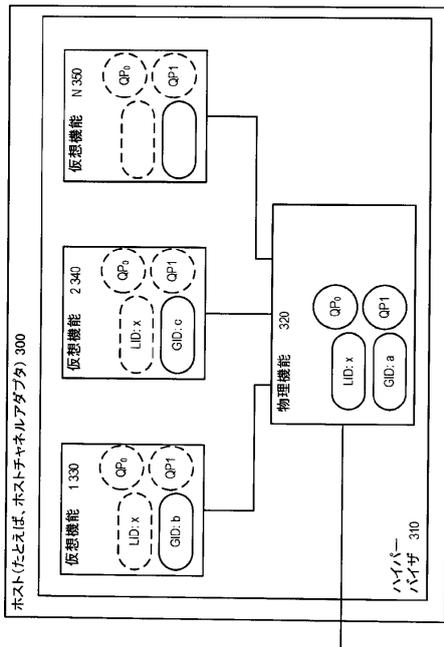


FIGURE 3

【図 4】

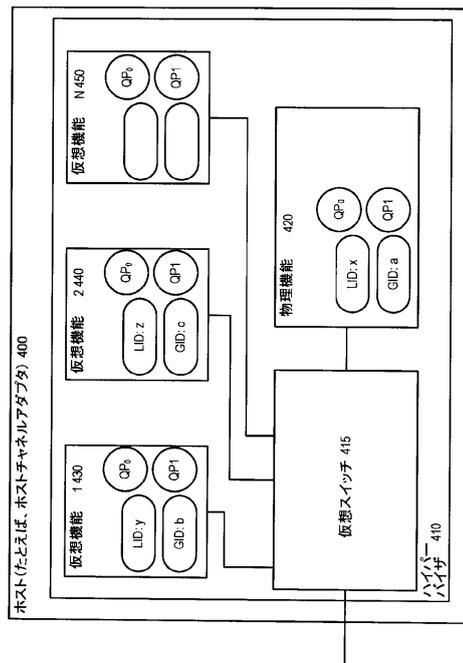


FIGURE 4

10

20

30

40

50

【図 5】

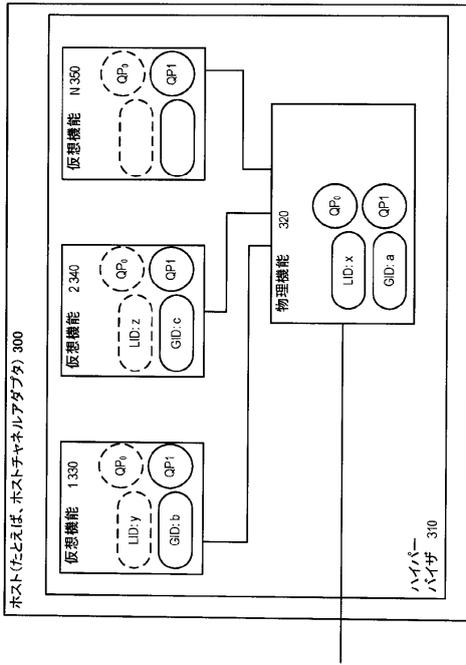


FIGURE 5

【図 6】

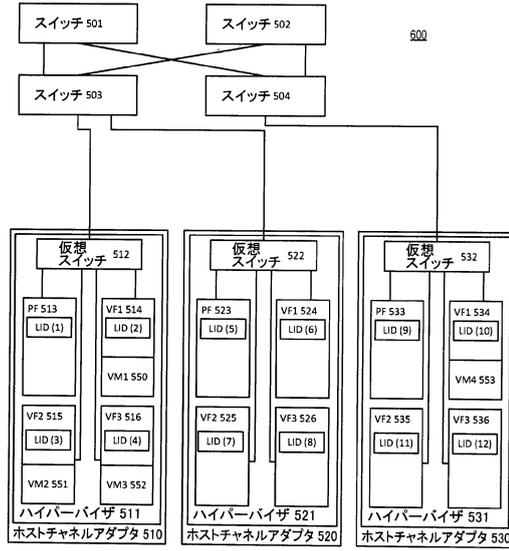


FIGURE 6

【図 7】

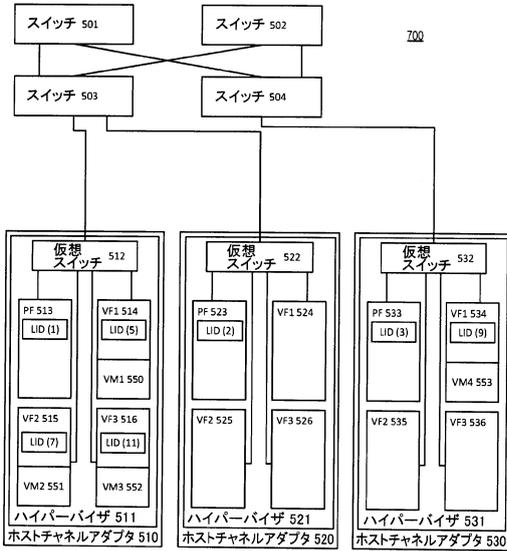


FIGURE 7

【図 8】

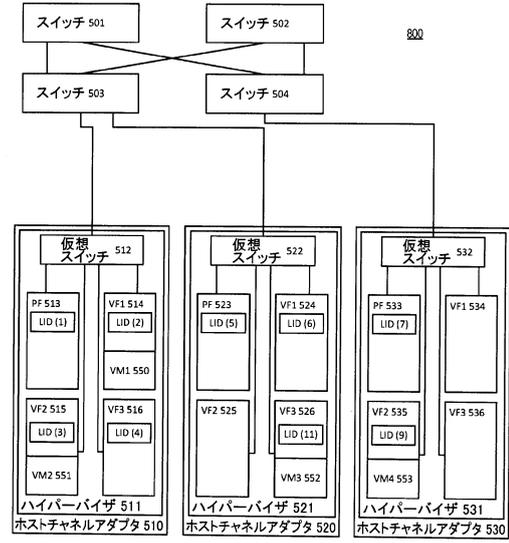


FIGURE 8

10

20

30

40

50

【 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 1 0 】

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 】

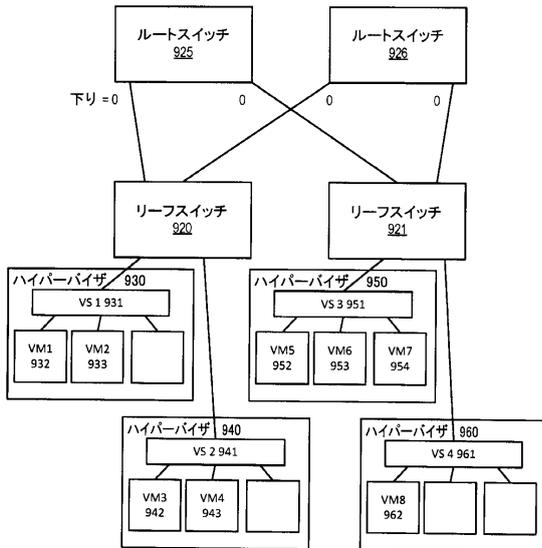


FIGURE 11

【 図 1 2 】

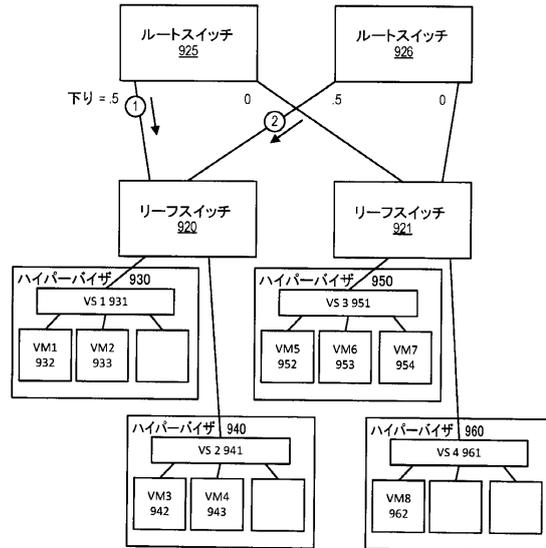


FIGURE 12

10

20

30

40

50

【 図 1 3 】

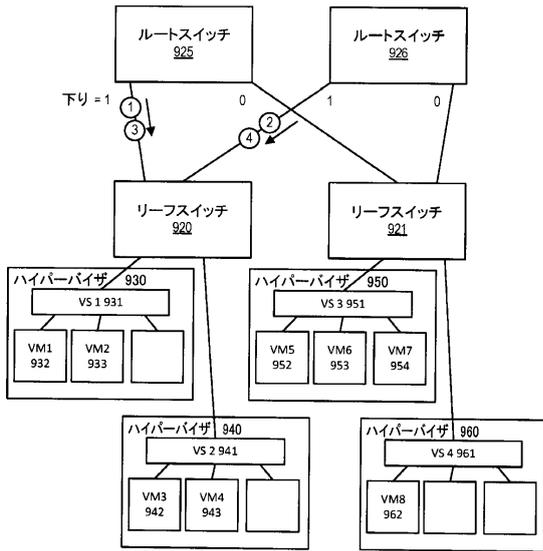


FIGURE 13

【 図 1 4 】

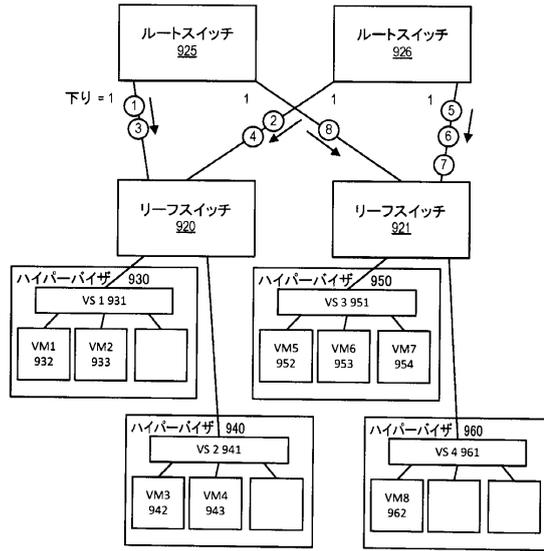


FIGURE 14

【 図 1 5 】

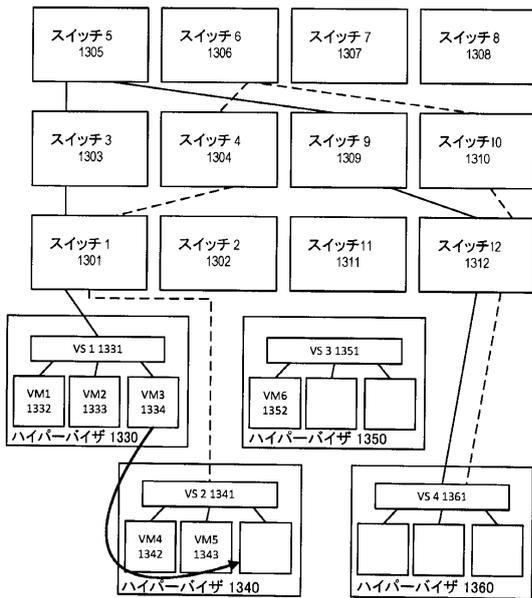


FIGURE 15

【 図 1 6 】

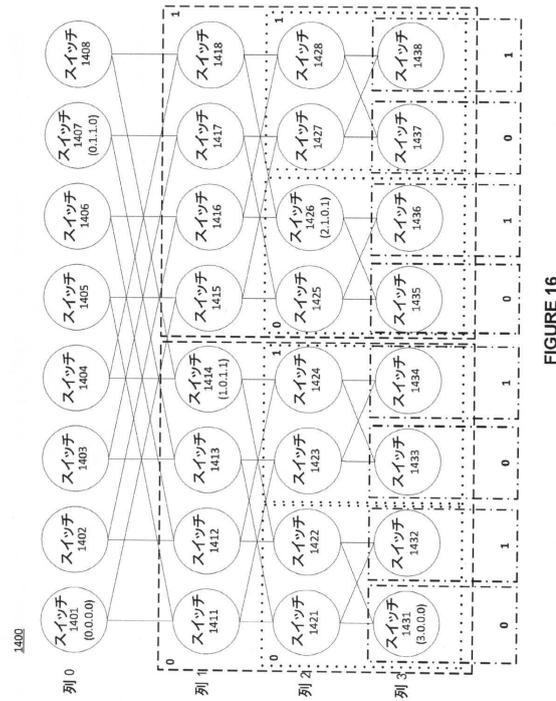


FIGURE 16

10

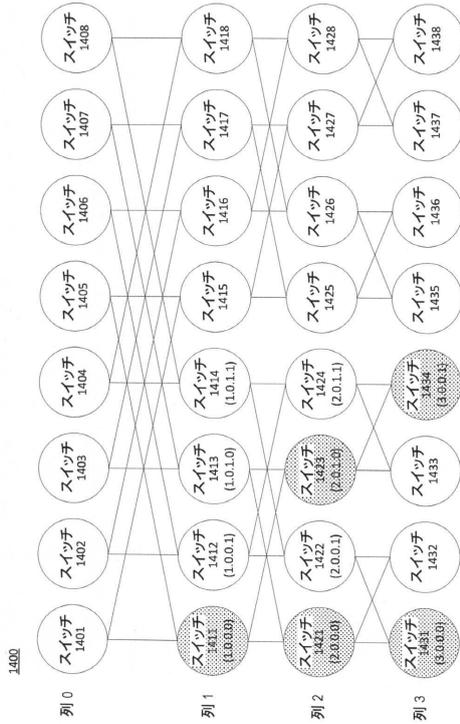
20

30

40

50

【 図 17 】



【 図 18 】

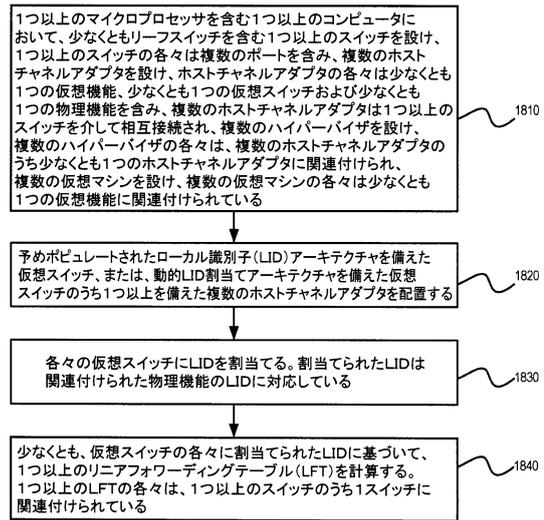


FIGURE 17

FIGURE 18

【 図 19 】

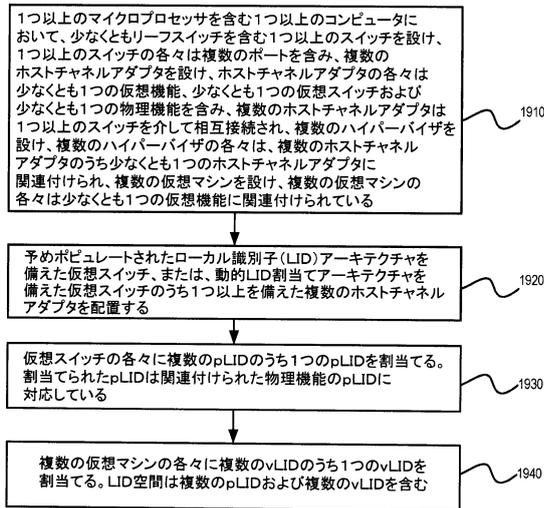


FIGURE 19

10

20

30

40

50

フロントページの続き

米国(US)

(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)

ルクグレンダ、9

(72)発明者 グラン、アーンスト・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 漆原 孝治

(56)参考文献 特開2011-028408(JP, A)

特開2013-069260(JP, A)

Evangelos Tasoulas, Towards the InfiniBand SR-IOV vSwitch Architecture, 2015 IEEE INTERNATIONAL CONFERENCE ON CLUSTER COMPUTING, 2015年09月08日, pp.371-380

(58)調査した分野 (Int.Cl., DB名)

G06F 15/173

G06F 9/455

G06F 9/50

H04L 61/00

H04L 61/10