

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6920281号  
(P6920281)

(45) 発行日 令和3年8月18日(2021.8.18)

(24) 登録日 令和3年7月28日(2021.7.28)

(51) Int.Cl. F I  
**G O 6 F 9/455 (2006.01)** G O 6 F 9/455 1 5 0  
**G O 6 F 15/177 (2006.01)** G O 6 F 15/177 A

請求項の数 7 (全 41 頁)

(21) 出願番号	特願2018-510841 (P2018-510841)	(73) 特許権者	502303739
(86) (22) 出願日	平成29年1月26日 (2017.1.26)		オラクル・インターナショナル・コーポレイション
(65) 公表番号	特表2019-509529 (P2019-509529A)		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
(43) 公表日	平成31年4月4日 (2019.4.4)	(74) 代理人	110001195
(86) 国際出願番号	PCT/US2017/015169		特許業務法人深見特許事務所
(87) 国際公開番号	W02017/132399	(72) 発明者	ヨンセン, ビョルン・ダグ
(87) 国際公開日	平成29年8月3日 (2017.8.3)		ノルウェー、0687 オスロ、ビルベルクグレンダ、9
審査請求日	令和1年12月2日 (2019.12.2)	(72) 発明者	ホウ, ハーラル
(31) 優先権主張番号	62/287, 712		ノルウェー、1326 リュサケール、ピイ・オウ・ボックス・384
(32) 優先日	平成28年1月27日 (2016.1.27)		
(33) 優先権主張国・地域又は機関	米国 (US)		
(31) 優先権主張番号	15/415, 668		
(32) 優先日	平成29年1月25日 (2017.1.25)		
(33) 優先権主張国・地域又は機関	米国 (US)		

最終頁に続く

(54) 【発明の名称】 高性能コンピューティング環境において仮想マシンの仮想マシンファブリックプロファイルを規定するためのシステムおよび方法

(57) 【特許請求の範囲】

【請求項1】

仮想マシン (VM) の仮想マシンファブリックプロファイルを作成する方法であって、  
プロセッサを含むノードが、仮想マシン識別子 (VM - id) を生成するステップと、  
前記ノードが、仮想ホストチャネルアダプタ (vHCA) を特定するvHCAインスタンスIDを生成するステップと、

前記ノードが、グローバル意識別子のプールからの仮想グローバル意識別子 (vGUID) を、前記vHCAの仮想ポートの現在のファブリックアドレスとして割り当てるステップと、

前記ノードが、管理パーティションを規定するP\_\_Keyと前記vGUIDとの間の第1の関係を作成するステップとを含み、前記P\_\_Keyと前記vGUIDとの間の前記第1の関係を、前記vGUIDが前記現在のファブリックアドレスとして割り当てられている前記仮想ポートを、前記P\_\_Keyによって規定される前記管理パーティションのメンバとして規定し、

前記ノードが、前記VM - idと前記vHCAインスタンスIDとの間の第2の関係を作成するステップを含み、前記第2の関係により、前記vHCAインスタンスIDを、前記VM - idへのアクセスを通して取出すことができ、

前記ノードが、前記VM - idと前記vGUIDとの間の第3の関係を作成するステップを含み、前記第3の関係により、前記vGUIDを、前記VM - idへのアクセスを通して取出すことができ、

前記ノードが、前記VM - id、前記vHCAインスタンスID、前記vGUID、前記第1の関係、前記第2の関係、および前記第3の関係を、前記VM - id、前記vHCA、前記vGUID、前記第1の関係、前記第2の関係、および前記第3の関係を保存するフォーマットで、永続化するステップを含む、方法。

【請求項2】

前記VM - id、前記vHCAインスタンスID、前記vGUID、前記第1の関係、前記第2の関係、前記第3の関係、および前記P\_\_Keyは、サブネットのサブネットマネージャのキャッシュに置かれ、前記P\_\_Keyによって規定される前記管理パーティションは前記サブネット内に定められる、請求項1に記載の方法。

【請求項3】

前記キャッシュに置かれた前記VM - id、前記vHCAインスタンスID、前記vGUID、前記第1の関係、前記第2の関係、前記第3の関係、および前記P\_\_Keyは、コピーである、請求項2に記載の方法。

【請求項4】

前記キャッシュは揮発性メモリである、請求項2または3に記載の方法。

【請求項5】

前記サブネットマネージャが、前記キャッシュから、前記vHCAインスタンスIDおよび前記vGUIDを取出すステップと、

前記サブネットマネージャが、前記vHCAインスタンスIDおよび前記vGUIDを物理ホストチャネルアダプタに送るステップとをさらに含む、請求項2～4のいずれか一項に記載の方法。

【請求項6】

仮想マシン（VM）の仮想マシンファブリックプロファイルを作成し前記仮想マシンファブリックプロファイルにアクセスするためのシステムであって、

プロセッサを含むノードと、

前記ノードにアクセス可能なメモリとを備え、

前記ノードは、

仮想マシン識別子（VM - id）を生成するように動作し、

仮想ホストチャネルアダプタ（vHCA）を特定するvHCAインスタンスIDを生成するように動作し、

グローバル意識別子のプールからの仮想グローバル意識別子（vGUID）を、前記vHCAの仮想ポートの現在のファブリックアドレスとして割り当てるように動作し、

管理パーティションを規定するP\_\_Keyと前記vGUIDとの間の第1の関係を作成するように動作し、前記P\_\_Keyと前記vGUIDとの間の前記第1の関係を、前記vGUIDが前記現在のファブリックアドレスとして割り当てられている前記仮想ポートを、前記P\_\_Keyによって規定される前記管理パーティションのメンバとして規定し、前記ノードはさらに、

前記VM - idと前記vHCAインスタンスIDとの間の第2の関係を作成するように動作し、前記第2の関係により、前記vHCAインスタンスIDを、前記VM - idへのアクセスを通して取出すことができ、前記ノードはさらに、

前記VM - idと前記vGUIDとの間の第3の関係を作成するように動作し、前記第3の関係により、前記vGUIDを、前記VM - idへのアクセスを通して取出すことができ、前記ノードはさらに、

前記VM - id、前記vHCAインスタンスID、前記vGUID、前記第1の関係、前記第2の関係、および前記第3の関係を、前記VM - id、前記vHCA、前記vGUID、前記第1の関係、前記第2の関係、および前記第3の関係を保存するフォーマットで、永続化するように動作する、システム。

【請求項7】

コンピュータシステムに請求項1～5のいずれか一項に記載の方法を実行させるためのコンピュータプログラム。

10

20

30

40

50

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

## 著作権に関する注意

本特許文献の開示の一部には、著作権保護の対象となるものが含まれている。著作権者は、この特許文献または特許開示の何者かによる複製が、特許商標庁の特許ファイルまたは記録にある限り、それに対して異議を唱えないが、そうでなければ、いかなる場合もすべての著作権を留保する。

## 【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、具体的には仮想マシンの仮想マシンファブリックプロファイルを規定することに関する。

## 【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<sup>TM</sup> (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。

## 【発明の概要】

## 【課題を解決するための手段】

## 【0004】

## 概要

本明細書に記載されているのは、仮想マシンの仮想マシンファブリックプロファイルを規定するためのシステムおよび方法である。ある例示的な実施形態は、仮想マシン識別子、仮想ホストチャネルアダプタインスタンスID、および仮想グローバル意識別子を提供することができる。仮想マシン識別子、仮想ホストチャネルアダプタインスタンスID、および仮想グローバル意識別子は、仮想マシン識別子にアクセスすることにより、仮想ホストチャネルアダプタインスタンスIDおよび仮想グローバル意識別子を取り出すことができるように、相互にマッピングすることができる。さらに、管理パーティションを規定するP\_Keyと仮想グローバル意識別子との間の関係を作成することができ、上記P\_Keyと仮想グローバル意識別子との間の関係は、仮想グローバル意識別子を、P\_Keyによって規定される管理パーティションのメンバとして規定する。

## 【図面の簡単な説明】

## 【0005】

【図1】一実施形態に従うインフィニバンド<sup>TM</sup>環境の一例を示す図である。

【図2】一実施形態に従う分割されたクラスタ環境の一例を示す図である。

【図3】一実施形態に従う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4】一実施形態に従う例示的な共有ポートアーキテクチャを示す図である。

【図5】一実施形態に従う例示的なvSwitchアーキテクチャを示す図である。

【図6】一実施形態に従う例示的なvPortアーキテクチャを示す図である。

【図7】一実施形態に従うLIDが予めポピュレートされている例示的なvSwitchアーキテクチャを示す図である。

【図8】一実施形態に従う動的LID割当てがなされた例示的なvSwitchアーキテクチャを示す図である。

【図9】一実施形態に従う、vSwitchに動的LID割当てがなされかつLIDが予めポピュレートされている、例示的なvSwitchアーキテクチャを示す図である。

10

20

30

40

50

【図10】一実施形態に従う例示的なマルチサブネットインフィニバンド<sup>T M</sup>ファブリックを示す図である。

【図11】一実施形態に従う、例示的な物理および論理サブネットリソースを含む例示的なインフィニバンド<sup>T M</sup>ファブリックおよびサブネットを示す図である。

【図12】一実施形態に従う、異なる管理パーティションのメンバとしての例示的なサブネットリソースを含む例示的なインフィニバンド<sup>T M</sup>ファブリックおよびサブネットを示す図である。

【図13】一実施形態に従う、管理パーティションおよびリソースドメイン双方を含む、階層管理スキームのメンバとしての例示的なサブネットリソースを含む例示的なインフィニバンド<sup>T M</sup>ファブリックおよびサブネットを示す図である。

10

【図14】一実施形態に従う、ファブリックレベルリソースドメインに対応付けられた管理パーティションのメンバおよび対応付けられたリソースに対して相互アクセス権を割当てるための方法のフローチャートを示す図である。

【図15】一実施形態に従う、VMファブリックプロファイル情報を格納するための例示的なデータベース構造を示す図である。

【図16】一実施形態に従う、VMファブリックプロファイルをサブネットリソースが利用できるようにするためのフローチャートを示す図である。

【図17】一実施形態に従う、仮想マシンの仮想マシンファブリックプロファイルを作成するためのフローチャートを示す図である。

【発明を実施するための形態】

20

【0006】

#### 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。

【0007】

30

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。

【0008】

本明細書に記載されているのは、仮想マシンの仮想マシンファブリックプロファイルを規定するためのシステムおよび方法である。

【0009】

この発明の以下の説明は、高性能ネットワークの一例として、インフィニバンド<sup>T M</sup> (IB) ネットワークを使用する。以下の説明全体にわたり、インフィニバンド<sup>T M</sup>の仕様(インフィニバンド仕様、IB仕様、またはレガシーIB仕様など、さまざまな呼ばれ方がある)を引用することがある。このような引用は、2015年3月に発表され、<http://www.inifinibandta.org>から入手可能な、本明細書にその全体を引用により援用するInfiniBand Trade Association Architecture Specification, Volume 1, Version 1.3を引用することであると理解される。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。

40

【0010】

50

### インフィニバンド™

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。

#### 【0011】

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。

10

#### 【0012】

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。

20

#### 【0013】

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。

#### 【0014】

一実施形態に従うと、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りニアフォーワーディングテーブル (linear forwarding table) (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ー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。

30

40

#### 【0015】

一般に、マスタ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、新しいマスタサブネットマネージャが取り決められる。マスタ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。

#### 【0016】

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

50

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。

#### 【0017】

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。

10

#### 【0018】

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。

20

#### 【0019】

一実施形態に従ったインフィニバンド<sup>TM</sup>環境100の一例を示す図1に、インフィニバンド<sup>TM</sup>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<sup>TM</sup>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。

#### 【0020】

##### インフィニバンド<sup>TM</sup>におけるデータパーティション

一実施形態に従うと、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を分離するためのセキュリティメカニズムとしての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HCAポートは、1つ以上のパーティションのメンバである可能性がある。一実施形態に従うと、本開示は、IBサブネット内に規定することができる2種類のパーティションとして、データパーティション(以下の段落でより詳細に説明する)と管理パーティション(本開示において後に詳細に説明する)とを提供する。

30

#### 【0021】

データ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、各ポートに関するデータ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\_Key)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情報を含むデータ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、スイッチポートのデータ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。

40

#### 【0022】

一実施形態に従うと、データパーティションはポートの論理グループであり、あるグループのメンバは同じ論理グループの他のメンバとしか通信できない。ホストチャネルアダ

50

プタ (HCA) およびスイッチにおいて、データパーティションメンバーシップ情報を用いてパケットをフィルタリングすることにより、分離を実施することができる。無効なパーティショニング情報を有するパケットは、当該パケットが入口ポートに達すると直ちにドロップすることができる。パーティショニングされたIBシステムにおいて、データパーティションを用いることにより、テナントクラスタを作成できる。データパーティションを適所で実施すると、ノードは異なるテナントクラスタに属する他のノードと通信することができない。このようにして、欠陥があるまたは悪意があるテナントノードが存在していても、システムのセキュリティを保證することができる。

#### 【0023】

一実施形態に従うと、ノード間の通信のために、マネージメントキューペア (QPO 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) を特定のデータパーティションに割り当てることができる。次に、P\_Key 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\_Key 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\_Key 値が見つかった場合、そのパケットは直ちに廃棄される。このようにして、通信は、データパーティションを共有するポート間でのみ許可される。

#### 【0024】

一実施形態に従い、データがパーティショニングされたクラスタ環境の一例を示す図2に、IBパーティションの一例が示される。図2に示す例では、ノードA101~E105は、インフィニバンド<sup>TM</sup>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。ノードA~Eは、データパーティション、すなわち、データパーティション1 130、データパーティション2 140、およびデータパーティション3 150に配置されている。データパーティション1はノードA101とノードD104とを含む。データパーティション2はノードA101とノードB102とノードC103とを含む。データパーティション3はノードC103とノードE105とを含む。データパーティションのこの配置により、ノードD104およびノードE105は、1つのデータパーティションを共有していないので、通信することができない。一方、たとえばノードA101およびノードC103は、どちらもデータパーティション2 140のメンバなので、通信することができる。

#### 【0025】

##### インフィニバンド<sup>TM</sup>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、必要な性能を提供することができる。

#### 【0026】

しかしながら、インフィニバンド<sup>TM</sup> (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。

#### 【0027】

10

20

30

40

50

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ホストチャネルアダプタ(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。

10

#### 【0028】

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のために用いられる。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。

20

#### 【0029】

##### 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

30

一実施形態に従うと、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つことができ、バイセクション帯域幅が維持される。

#### 【0030】

図3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3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)スイッチ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。

40

#### 【0031】

また、図3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。

50

## 【 0 0 3 2 】

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 3 】

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、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。

10

## 【 0 0 3 4 】

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。

20

## 【 0 0 3 5 】

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、共通のデータ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。データ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、データパーティション間の分離が不良なものとなる。

30

## 【 0 0 3 6 】

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 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 7 】

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 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

40

50

を用いている。

【 0 0 3 8 】

入出力 ( I / O ) 仮想化

－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 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み合わせると、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 C P U 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。

10

【 0 0 3 9 】

－実施形態に従うと、I / O リソースの共有を可能にして、V M 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、V M にエクスポーズされる論理装置を、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 : D A )、およびシングルルート I / O 仮想化 ( S R - 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。

【 0 0 4 0 】

－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 V M 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、V M 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

【 0 0 4 1 】

－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 V M に連結する必要があるが、デバイスは V M 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、V M 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 V M 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。

30

【 0 0 4 2 】

－実施形態に従うと、シングルルート I O V ( Single Root I O V : S R - 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 V M に割当てることができ、仮想機能 ( Virtual Function : V F ) 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 : P F ) によってデバイスにアクセスする。S R - 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、S R - I O V によって提示される問題は、それが V M 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、S R - I O V は、ほぼ固有の性能を維持しながらも、複数の V M 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 P C I E x p r e s s ( P C I e ) 規格を拡張することができる。これにより、S R - 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

【 0 0 4 3 】

S R - I O V は、P C I e 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

50

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/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。

10

## 【0044】

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。

## 【0045】

20

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<sup>T M</sup>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す可能性がある。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。

30

## 【0046】

インフィニバンド<sup>T M</sup> 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。

## 【0047】

図4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。

40

## 【0048】

一実施形態に従うと、図4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。

## 【0049】

図4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<sup>T M</sup>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

50

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。

【0050】

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LIDスペースが消費されるだけであるからである。

【0051】

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VM配置についての可能性が妨害されてしまう。各々のLID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VM間で共有されているので、マイグレートしているVM（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LIDを宛先ハイパーバイザのLIDに変更させなければならない。さらに、QP0アクセスが制限された結果、サブネットマネージャはVMの内部で実行させることができなくなる。

【0052】

インフィニバンド<sup>T M</sup>SR - IOVアーキテクチャモデル - 仮想スイッチ (vSwitch) (ch)

図5は、一実施形態に従った例示的なvSwitchアーキテクチャを示す。図に示されるように、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。

【0053】

一実施形態に従うと、vSwitch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: 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については、HCA400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PF420を用いることができ、（仮想機能に付与された）VMはVFを用いる。

【0054】

一実施形態に従うと、vSwitch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LIDが割当てられているので、利用可能な数のLIDが速やかに消費される。同様に、多くのLID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SMによって演算しなければならず、それらのLFTを更新するために、より多くのサブネット管理パケット（SMP）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。LIDスペースが49151個のユニキャストLIDに制限されており、（VFを介する）各々のVMとして、物理ノードおよびスイッチがLID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VMの数が制限されてしまい、逆の場合も同様に制限される。

【0055】

インフィニバンド<sup>T M</sup>SR - IOVアーキテクチャモデル - 仮想ポート (vPort)

図6は、一実施形態に従った例示的なvPortの概念を示す。図に示されるように、

10

20

30

40

50

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。

【0056】

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

【0057】

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。

【0058】

図 7 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ー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り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 に割り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 に割り当て、仮想マシン 3 552 を仮想機能 3 516 に割り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 に割り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

【0059】

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。

40

【0060】

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 に割り当てられた 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 SM (図示せず) については、HCA 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。

【0061】

50

一実施形態に従うと、本開示は、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は、さまざまな物理機能5 1 3、5 2 3および5 3 3に、さらには、仮想機能5 1 4 ~ 5 1 6、5 2 4 ~ 5 2 6、5 3 4 ~ 5 3 6（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5 1 3はL I D 1が予めポピュレートされており、仮想機能1 5 3 4はL I D 1 0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7に示されるようにL I Dが割当てられている。

10

## 【0062】

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。

## 【0063】

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、実際これらの数はより小さくなる）。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。

20

30

## 【0064】

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。

40

## 【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。

## 【0066】

50

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い場合よりも時間が長くなる可能性がある。

【 0 0 6 7 】

インフィニバンド<sup>T M</sup> 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 6 8 】

図8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<sup>T M</sup>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、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、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 6 9 】

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ー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 0 】

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 1 】

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8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

10

20

30

40

50

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 2 】

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 3 】

一実施形態に従うと、多くの同様の物理的なホストチャンネルアダプタが 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 4 】

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 5 】

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 6 】

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 7 】

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

10

20

30

40

50

めに、1スイッチごとに1つのサブネット管理パケット(SMP)が送信される必要がある。各々のVM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、LMC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VFの合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。

【0078】

インフィニバンド<sup>T M</sup>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9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<sup>T M</sup>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。

【0079】

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ワーディングテーブルを設定するのに用いられる。

【0080】

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。

【0081】

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9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

10

20

30

40

50

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。

10

#### 【0082】

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、ホストチャン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、（ホストチャン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、（ホストチャン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。

20

#### 【0083】

一実施形態に従うと、多くの同様の物理的なホストチャンネルアダプタが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。

#### 【0084】

##### インフィニバンド<sup>T M</sup> - サブネット間通信

一実施形態に従うと、1つのサブネット内にインフィニバンド<sup>T M</sup>ファブリックを提供することに加え、本開示の実施形態は、2つ以上のサブネットにまたがるインフィニバンド<sup>T M</sup>ファブリックを提供することもできる。

30

#### 【0085】

図10は、一実施形態に従う例示的なマルチサブネットインフィニバンド<sup>T M</sup>ファブリックを示す。この図に示されるように、サブネットA 1 0 0 0の内部の多数のスイッチ1 0 0 1 ~ 1 0 0 4は、サブネットA 1 0 0 0（たとえばI Bサブネット）内におけるインフィニバンド<sup>T M</sup>ファブリックなどのファブリックのメンバ間の通信を提供することができる。このファブリックは、たとえばチャンネルアダプタ1 0 1 0などの多数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1 0 1 0は、ハイパーバイザ1 0 1 1と対話することができる。ハイパーバイザは、対話の相手であるホストチャンネルアダプタとともに、多数の仮想機能1 0 1 4をセットアップすることができる。加えて、ハイパーバイザは、仮想マシンを仮想機能各々に割当てることができる。たとえば、仮想マシン1 1 0 1 5は仮想機能1 1 0 1 4に割当てられる。ハイパーバイザは、その対応付けられたホストチャンネルアダプタに、各ホストチャンネルアダプタ上の物理機能1 0 1 3などの十分な機能を有する物理機能を通して、アクセスすることができる。

40

#### 【0086】

さらに図10を参照して、一実施形態に従うと、多数のスイッチ1 0 2 1 ~ 1 0 2 4は、サブネットB 1 0 4 0（たとえばI Bサブネット）内におけるインフィニバンド<sup>T M</sup>ファブリックなどのファブリックのメンバ間の通信を提供することができる。このファブリックは、たとえばホストチャンネルアダプタ1 0 3 0などの多数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1 0 3 0は、ハイパーバイザ1 0 3 1と対話することができる。ハイパーバイザは、対話の相手であるホストチャンネルアダプタとともに、多

50

数の仮想機能 1 0 3 4 をセットアップすることができる。加えて、ハイパーバイザは、仮想マシンを仮想機能各々に割当てることができる。たとえば、仮想マシン 2 1 0 3 5 は仮想機能 2 1 0 3 4 に割当てられる。ハイパーバイザは、その対応付けられたホストチャンネルアダプタに、各ホストチャンネルアダプタ上の物理機能 1 0 3 3 などの十分な機能を有する物理機能を通して、アクセスすることができる。なお、各サブネット（すなわちサブネット A およびサブネット B）内に示されているホストチャンネルアダプタは 1 つだけであるが、各サブネット内に複数のホストチャンネルアダプタおよびそれらに対応するコンポーネントが含まれていてもよいことが、理解されるはずである。

#### 【 0 0 8 7 】

一実施形態に従うと、各ホストチャンネルアダプタはさらに、仮想スイッチ 1 0 1 2 および仮想スイッチ 1 0 3 2 などの仮想スイッチに対応付けられていてもよく、上記のように各 H C A は異なるアーキテクチャモデルでセットアップされてもよい。図 1 0 のサブネットはどちらも L I D が予めピュレートされている v S w i t c h のアーキテクチャモデルを使用するものとして示されているが、これは、このようなサブネット構成すべてが同様のアーキテクチャモデルに従う必要があることを示唆しようとしているのではない。

#### 【 0 0 8 8 】

一実施形態に従うと、各サブネット内の少なくとも 1 つのスイッチがルータに対応付けられていてもよい。たとえば、サブネット A 1 0 0 0 内のスイッチ 1 0 0 2 はルータ 1 0 0 5 に対応付けられ、サブネット B 1 0 4 0 内のスイッチ 1 0 2 1 はルータ 1 0 0 6 に対応付けられている。

#### 【 0 0 8 9 】

一実施形態に従うと、サブネット A 内の仮想マシン 1 などの発信ソースにおけるトラフィックを、サブネット B 内の仮想マシン 2 などの異なるサブネットを宛先としてそれに向ける場合、トラフィックは、サブネット A 内のルータ、すなわち、ルータ 1 0 0 5 に向ければよく、そうすると、ルータ 1 0 0 5 はこのトラフィックをルータ 1 0 0 6 とのリンクを介してサブネット B に送ることができる。

#### 【 0 0 9 0 】

##### ファブリックマネージャ

上述のように、インフィニバンド<sup>TM</sup> ファブリックなどのネットワークファブリックは、ファブリックの各サブネット内の相互接続されたルータを用いることで、複数のサブネットにまたがることができる。一実施形態に従うと、ファブリックマネージャ（図示せず）はホスト上で実現することができる。このホストは、ネットワークファブリックのメンバであり、ファブリック内で使用されることにより、ファブリックの一部である物理リソースおよび論理リソース双方を管理することができる。たとえば、特に、ファブリックリソースを発見すること、物理サーバ間の接続性を制御すること、リアルタイムのネットワーク統計を収集して観察すること、災害復旧（disaster recovery）、およびサービス品質（quality of service: Q o S）設定を設定することなどの管理タスクを、ユーザは、ファブリックマネージャを通して実行することができる。一実施形態に従うと、ファブリックマネージャは、ファブリック内に規定されたすべてのサブネットにまたがっていてもよい。すなわち、ファブリックマネージャは、リソースがどのサブネットのメンバであるかに関係なく、ファブリックのメンバであるかまたは少なくともファブリックに対応付けられている物理リソースおよび論理リソースを管理することができる。

#### 【 0 0 9 1 】

一実施形態に従うと、ファブリックマネージャはグラフィカルユーザインターフェイス（graphical user interface: G U I）を含み得る。ユーザは、G U I を通して管理機能を実行することができる。ファブリックマネージャ G U I は、ユーザがファブリックリソースをモニタリングし制御することを可能にする可視化ツールを取入れてもよい。たとえば、一実施形態において、ユーザは、ファブリック全体のサーバについて、サーバ接続、構成設定およびパフォーマンス統計を、ファブリックインターフェイスを通して見ることができる。ファブリックマネージャ G U I を通してモニタリングおよび/または管理する

10

20

30

40

50

ことができるファブリック機能のその他の例は、サブネット間ファブリックトポロジを  
発見すること、これらのトポロジーの視覚表現を見ること、ファブリックプロファイル（  
たとえば仮想マシンファブリックプロファイル）を作成すること、および、ファブリック  
マネージャデータベースを構築し管理することを含む。ファブリックマネージャデータベ  
ースは、ネットワークファブリックが必要としネットワークファブリックに関連する、フ  
ァブリックプロファイル、メタデータ、構成設定、およびその他のデータを格納すること  
ができる。一実施形態に従うと、ファブリックマネージャデータベースはファブリックレ  
ベルデータベースである。

【 0 0 9 2 】

加えて、ファブリックマネージャは、どのサブネットがどのルータポートを介してどの  
パーティション番号を用いて通信することを許可されるかについて、法的なサブネット間  
接続性を規定することができる。一実施形態に従うと、ファブリックマネージャは、集中  
型ファブリック管理ユーティリティである。上に挙げた例は限定を意図しているのではな  
い。

10

【 0 0 9 3 】

一実施形態に従うと、ファブリックマネージャの機能の一部はユーザによって開始され  
ることができ、その他の機能はユーザから引き離されてもよくまたは自動化されてもよい  
（たとえば、機能の一部は、起動に応じてまたはその他所定のイベント時に、ファブリッ  
クマネージャによって実行されてもよい）。

【 0 0 9 4 】

20

管理イベントの例示的な実施形態において、ユーザは、ファブリックマネージャインタ  
ーフェイスで、ネットワークファブリックデバイス向けの構成変更を開始してもよい。構  
成変更要求を受けた後に、ファブリックマネージャは、この構成変更要求が適切に実行さ  
れることを保証してもよい。たとえば、ファブリックマネージャは、要求をデバイスに伝  
え構成変更がデバイスの構成に書込まれることを保証してもよい。一実施形態において、  
物理デバイスは、構成変更が無事に終了したことをファブリックマネージャに伝える。一  
実施形態に従うと、それに応じて、ファブリックマネージャはインターフェイスを更新し  
、要求が実行されたことを視覚的に確認できるようにしてもよい。さらに、ファブリック  
マネージャは、たとえば災害復旧のためまたはそれ以外の目的のために、デバイスの構成  
をファブリックマネージャデータベースに永続化（persist）してもよい。

30

【 0 0 9 5 】

一実施形態に従うと、ファブリックマネージャは、G U I のいくつかの、すべての、ま  
たはより多くの機能を含む、コマンドラインインターフェイスなどのその他のインターフ  
ェイスを有し得る。

【 0 0 9 6 】

ファブリックレベルリソースドメイン

上述のように、ファブリックマネージャは、ユーザがファブリックマネージャのインタ  
ーフェイスを用いてネットワークファブリック全体にわたって管理タスクを実行できるよ  
うにする。一実施形態に従うと、ファブリックマネージャのその他の機能は、階層的なロ  
ール（role）ベースのアクセス制御を容易にすることである。一実施形態において、ロー  
ルベースのアクセス制御は、ファブリックレベルリソースドメインを通して実現される。

40

【 0 0 9 7 】

一実施形態に従うと、ロールベースのアクセス制御は、ファブリックユーザという概念  
に基づく。人間であるアドミニストレータからのアクセスおよび外部管理アプリケーション  
からのアクセスはいずれも、ファブリックインフラストラクチャまたはファブリックリ  
ソースのすべてまたはサブセットに対する法的処理を定める認証されたコンテキストを表  
わすことができる。たとえば、ユーザは、ファブリックにおいて、ユーザプロファイルに  
よって表わすことができる。すなわち、ファブリック内において、ユーザを、そのユーザ  
のプロファイルを作成し属性をこのプロファイルに割当てることによって定義することが  
できる。ユーザプロファイルには、ユーザ名属性およびパスワード属性を割当てることが

50

でき、ユーザ名は、ファブリック内で固有であるため、ユーザを一意に特定する。さらに、ユーザプロファイルは、ファブリック内のさまざまなリソースに特定のアクセスレベルを割当て、ファブリック内に規定された特定の役割に、対応付けられてもよい。一実施形態に従うと、ユーザプロファイルのセットアップは、ファブリックマネージャインターフェイスを通して実現することができる。ユーザプロファイルのすべてまたは一部を、ファブリックマネージャデータベースに格納することができる。加えて、一実施形態において、ファブリックマネージャは、マイクロソフト（登録商標）のアクティブディレクトリまたはLDAPディレクトリなどの周知のユーザディレクトリと一体化してもよく、または、たとえば遠隔認証のためのRADIUSネットワークングプロトコルと一体化してもよい。

10

**【0098】**

一実施形態に従うと、ファブリックマネージャは、ファブリックレベルリソースドメイン（本明細書では「リソースドメイン」または簡単に「ドメイン」と呼ぶ）を通して自身が発見したファブリックリソースを管理することができる。リソースドメインは、ファブリックレベルで規定されたファブリックリソースの論理的なグルーピングである。ファブリックリソースは、物理リソースおよび論理リソース双方を含む。リソースの例の一部は、特に、ファブリックデバイス（HCA、物理ノード、およびスイッチなど）、ファブリックプロファイル（仮想マシンファブリックプロファイル、およびユーザプロファイルなど）、仮想マシン、クラウド、および入出力モジュールを含む。

**【0099】**

一実施形態に従うと、ファブリックマネージャによって発見され管理されるすべてのファブリックリソースは、デフォルトでファブリック内に存在する（すなわちセットアップまたは構成する必要がない）デフォルトドメイン内にあり、これらのファブリックリソースには、ファブリックマネージャインターフェイスを通してアクセスすることができる。デフォルトドメインは、最高レベルのドメインである、すなわち、その他すべてのリソースドメインに対する親（parent）ドメインであり、その他すべてのリソースドメインはデフォルトドメイン内に存在する。デフォルトドメインはファブリックレベルアドミニストレータに対応付けられる。ファブリックレベルアドミニストレータもデフォルトで存在し、デフォルトでデフォルトドメイン内の管理上の権利を有するように構成される。

20

**【0100】**

一実施形態に従うと、リソースドメインは、階層形式のリソース管理を表わす。たとえば、デフォルトドメインを構成し管理するプロセスは、ファブリックレベルアドミニストレータしか利用できない。しかしながら、子（child）ドメインは、ファブリックレベルアドミニストレータが、デフォルトドメイン内に作成することができる。たとえば、ファブリックレベルアドミニストレータは、子ドメインを作成することができ、ドメインリソースを子ドメインに追加することができる。加えて、ファブリックレベルアドミニストレータは、ドメインレベルの「ドメイン管理」ユーザを作成することができ、ドメイン管理ユーザを子ドメインに追加する（すなわち対応付ける）ことができる。ドメイン管理ユーザをリソースドメインのメンバにすることにより、ドメイン管理ユーザが、子ドメインとそれに含まれるファブリックリソースのサブセットとを管理できるようにする。一実施形態に従うと、ドメイン管理ユーザは、子メインの外部にあるリソース（すなわちこのドメイン管理が対応付けられているレベルと同様のまたはそれよりも高いレベルのリソース）を管理することはできない。しかしながら、ドメイン管理は、リソースドメインの子ドメインとして作成されたリソースドメインに含まれているリソースを管理することができる。一実施形態に従うと、ファブリックマネージャは、リソースドメインの境界が厳密に守られることを保証するセキュリティを提供する役割を担う。

30

40

**【0101】**

図11は、リソースドメインの階層構造を示す。図示のように、デフォルトドメイン1102がネットワークファブリック1100内に存在する。ファブリックレベルアドミニストレータ1110は、ファブリックレベルリソース1112、1124、および113

50

4の管理のためのアクセス権を有する。ファブリックレベルアドミニストレータ1110はまた、新たなリソースドメインをデフォルトドメイン1102内に作成し管理することができる。リソースドメイン1120および1130と、対応するドメインレベルドメイン管理ユーザ1122および1132とは、ファブリックレベルアドミニストレータ1110によって作成されたものである。ドメイン管理ユーザ1122は、(ファブリックレベルアドミニストレータ1110によりリソースドメイン1120に割当てられた)ファブリックリソース1124を管理するためのアクセス権を有するが、(より高いレベルの)ファブリックリソース1112または(同様のレベルの)ドメインリソース1134を管理するためのアクセス権は有しない。同様に、ドメイン管理ユーザ1132は、(ファブリックレベルアドミニストレータ1110によりリソースドメイン1130に割当てられた)ファブリックリソース1134を管理するためのアクセス権を有するが、(より高いレベルの)ファブリックリソース1112または(同様のレベルの)ドメインリソース1124を管理するためのアクセス権は有しない。

10

#### 【0102】

##### 管理パーティション

一実施形態に従うと、リソースドメインは、アドミニストレーションまたは(本明細書における呼称である)「管理」パーティションにより、サブネットレベルで表わすことができる。管理パーティションは、サブネットレベルでアクセス権をサブネットリソースに与えるグループメンバーシップを表わす。一実施形態に従うと、管理パーティションのメンバは、管理パーティションに対応付けられているいずれのサブネットリソースへのアクセス権も有するという点において、特権があると考えられる。ファブリックマネージャレベルにおいて、管理パーティションは、リソースドメインと対応するドメイン管理ユーザとに対応付けられる。このようにして、マルチテナント環境においてユーザとロールの分離をサブネットレベルで保証することができる。さらに、特定のリソースドメインに対応付けられている管理パーティションのメンバであるリソースがリソースドメインのメンバにもなるように、リソースドメインメンバーシップを管理パーティションメンバーシップに対応付けることができる。

20

#### 【0103】

一実施形態に従うと、管理パーティションは、サブネットレベルにおいて、データパーティションの規定と同じやり方で規定することができるが、作成されたパーティションが管理パーティションであることを示す追加の属性を伴う。(先に詳述した)データパーティションと同様、管理パーティションは、一実施形態に従い、アドミニストレータにより、ファブリックマネージャインターフェイスを通して作成することができる。一実施形態において、ファブリックマネージャは、パーティションの作成中、「管理パーティション」フラグを任意のパラメータとしてサポートすることができる。作成したアドミニストレータによって選択される場合、ファブリックマネージャは、新たに作成されたパーティションが管理パーティションであることを示す追加の属性を含みかつファブリックマネージャおよびローカルマスタサブネットマネージャによって管理パーティションとして扱われるであろう。

30

#### 【0104】

一実施形態に従うと、ファブリックマネージャは、作成されたリソースドメインごとに対応する管理パーティションを自動的に作成し、自動的に作成したパーティションに対応するリソースドメインに対応付けるように構成されてもよい。このような実施形態において、ファブリックレベルリソースがリソースドメインに追加されると、ファブリックマネージャはまた、これらを、自動的に作成されリソースドメインに対応付けられた管理パーティションに対応付ける。このように、リソースドメインに追加されたリソースは、アドミニストレータ(たとえばファブリックレベルアドミニストレータまたはドメイン管理)によるさらに他の動作を伴うことなく、リソースドメインに追加されると互いにサブネットレベルアクセス権を有することになる。

40

#### 【0105】

50

加えて、一実施形態に従うと、ネットワークのサブネット全体は、トップレベルドメイン（たとえばデフォルトドメイン）を有するドメイン階層における特別なリソースドメインを表わすことができる。たとえば、ドメイン階層において、デフォルトドメインがトップレベルドメインを表わす場合、ネットワークファブリックの各サブネットは、ファブリックマネージャにより、デフォルトドメインの子ドメインとして認識されることができる。サブネット全体をトップレベルドメインの子ドメインとして認識することを、ファブリックマネージャのデフォルト挙動として構成してもよく、または、これらデフォルトドメインがアドミニストレータによって手作業で定められてもよい。ここでも、ユーザとロールを分離し、ドメイン境界を守り、サブネットレベルにおけるリソースの対応付けを実施するために、サブネットリソースドメイン全体に対応する管理パーティションを規定することができる。一実施形態に従うと、サブネット内に規定され当該サブネット内の各リソース（メンバであるかまたは管理パーティションに対応付けられている）を含む管理パーティションは、「ドメイングローバル」管理パーティションと名付けることができる。なぜなら、この構成の中で、サブネット内のすべてのリソースはその他すべてのリソースへのアクセス権を有するからである。

10

**【0106】**

一実施形態に従うと、管理パーティションは、ドメイン管理にとってトランスペアレントであってもよい。上述のように、ドメイングローバル管理パーティションは、ファブリックマネージャレベルでリソースドメインに対して自動的に作成することができ、その後、このドメインの範囲に割当てられたまたはこのドメインの範囲内に作成されたすべてのリソースに対応する管理パーティションに自動的に対応付けることができる。しかしながら、別の実施形態において、ドメイン管理は、関連するリソースドメイン内に種々の管理パーティションを明示的に作成することができ、その後、当該ドメイン内のリソースを、リソースドメインに対してデフォルトで作成された管理パーティションの代わりに明示的に作成された管理パーティションに対応付けることができる。

20

**【0107】**

一実施形態に従うと、ファブリックマネージャは、共有管理パーティションおよびプライベート管理パーティション双方の作成をサポートできる。デフォルトドメイン内のファブリックレベルアドミニストレータによって作成された管理パーティションは、個々のリソースドメインが利用できるようにすることができる共有パーティションであってもよい。ドメイン管理がメンバであるドメインによって作成された管理パーティションは、そのコンテキストに管理パーティションが作成されている特定のリソースドメインのみに対応付けられ当該特定のリソースドメインのみが利用できるプライベートパーティションであってもよい。

30

**【0108】**

一実施形態に従うと、H C Aおよびv H C Aのエンドポートは、データパーティションのメンバであってもよいのと同じく、管理パーティションのメンバであってもよい。しかしながら、一実施形態に従うと、管理パーティションはその他のサブネットリソースに対応付けることができるという点で、管理パーティションはデータパーティションと区別される。たとえば、データパーティションは管理パーティションに対応付けることができる。さらに、一実施形態に従うと、管理パーティションを子または親としての別の管理パーティションに対応付けることにより、これらの管理パーティションは階層的概念となり、対応付けられているリソースドメインの階層に対応させることができる。

40

**【0109】**

技術的な問題として、従来の専門用語では、H C A（およびv H C A）のエンドポートをパーティションの「メンバ」と呼ぶことができ、一実施形態に従うと、その他のリソースは管理パーティション「に対応付ける」ことができる。以下、これら2つの概念の技術的相違点について説明する。しかしながら、便宜上かつ読み易さのために、本文献では、管理パーティションについて、「メンバ」および「～に対応付ける」という表現を区別なく使用する場合がある。これらの表現が区別なく使用されていても、それとは関係なく、

50

管理パーティションへのエンドポート / H C A のメンバーシップと、管理パーティションに対応付けられたリソースとの間の技術的相違が読み手によって一貫して適用されることを意図していることが理解されるはずである。

**【 0 1 1 0 】**

一実施形態に従うと、管理パーティションは、データパーティションの規定と同じく、P \_ K e y によって規定される。しかしながら、エンドポートは、自身がメンバであるデータパーティションを認識しているものの、エンドポートは自身がメンバであるのがどの管理パーティションなのかを認識している必要はない。よって、一実施形態において、管理パーティションを規定する P \_ K e y は、メンバのエンドポートの P \_ K e y テーブルに入力されない。このように、管理パーティションが I B パケットトラフィック用に使用されなければ、管理パーティションの作成が、限りあるリソースである P \_ K e y テーブルエントリを無駄に使うことはない。しかしながら、別の実施形態において、管理パーティションは管理パーティションとしてもデータパーティションとしても機能する場合がある。このような実施形態において、管理パーティションのメンバであるエンドポートの P \_ K e y テーブルすべてが、それぞれの P \_ K e y テーブルの管理パーティションの P \_ K e y エントリを有することができる。一実施形態に従うと、データパーティションは、管理パーティションではない何らかのパーティションであると定義してもよい。

10

**【 0 1 1 1 】**

一実施形態に従うと、データパーティションは 1 つ以上の管理パーティションに対応付けることができる。たとえば、P \_ K e y 値によって規定されるデータパーティションは、自身の固有の P \_ K e y 値によって規定される管理パーティションに対応付けることができる。加えて、このデータパーティションは、別の固有の P \_ K e y 値によって規定される第 2 の管理パーティションに対応付けることができる。一実施形態に従うと、データパーティションを特定の管理パーティションに対応付けることにより、当該特定の管理パーティションのメンバであるエンドポートの最大メンバーシップレベルを規定することができる。

20

**【 0 1 1 2 】**

上述の通り、管理パーティションは、サブネットリソースにアクセス権を与えるグループメンバーシップを表わす。一実施形態に従うと、管理パーティションのどのエンドポートメンバも、当該管理パーティションへのエンドポートのメンバーシップのみに基づいて、同じ管理パーティションに対応付けられているいずれのサブネットリソースへのアクセス権も有する。したがって、管理パーティションのメンバであるエンドポートはいずれも、同じ管理パーティションに対応付けられているいずれのデータパーティションへのアクセス権も有する。なお、これは、メンバエンドポートが、対応付けられたデータパーティションのメンバであることを必ずしも意味しておらず、対応付けられたデータパーティションへのアクセス権を有しそれ故にこのデータパーティションのメンバであり得ることを意味している。

30

**【 0 1 1 3 】**

このような手法により、アドミニストレータが、たとえばデータパーティションへのアクセス権を、エンドポートに、データパーティションの P \_ K e y をこのエンドポートの P \_ K e y テーブルに手作業で含めることによって与える必要は、なくなる。一実施形態において、エンドポートがサブネット内で初期化されると、マスタサブネットマネージャは、管理パーティションの規定（たとえば P \_ K e y ）を保持するとともに規定された管理パーティションへのメンバーシップを規定しかつ規定された管理パーティションとの対応付けを規定する関係を保持するデータストア（たとえば以下で説明する管理パーティションレジストリ）にクエリすることによって、このエンドポートがメンバであるのはどの管理パーティションかを判断することができる。次に、サブネットマネージャはさらに、エンドポートがメンバである管理パーティションに対応付けられているデータパーティションがあるか否かを確かめるためにチェックすることができる。S M が、1 ) エンドポートが管理パーティションのメンバであること、および、2 ) 当該管理パーティションがデ

40

50

ータパーティションに対応付けられていることを見出した場合、SMは、対応付けられているデータパーティションのP\_KeyをエンドポートのP\_Keyテーブルに自動的に置くことにより、データパーティションへアクセス権をエンドポートの自動的に与える。このように、管理パーティションは、アドミニストレータによる手作業のパーティションマッピングと比較して、より簡単でよりスケーラブルな解決策を表わす。

【0114】

図12は、管理パーティションおよびデータパーティション双方を有する例示的なネットワークファブリックを示す。図12に示されるように、管理パーティション1230、1240、および1250がファブリック内において規定されている。ノードA 1201~E 1205が、それぞれのHCA1211~1215によってファブリックに物理的に接続されている。加えて、各HCAは、少なくとも1つの管理パーティションのメンバである。HCA1211およびHCA1214は、管理パーティション1230のメンバである。HCA1211はまた、HCA1212および1213とともに管理パーティション1240のメンバである。HCA1213はさらに、HCA1215とともに管理パーティション1250のメンバである。

10

【0115】

さらに図12を参照して、一実施形態に従うと、データパーティション1232、1242、および1252がファブリック内に規定されている。データパーティション1232は管理パーティション1230に対応付けられ、データパーティション1242は管理パーティション1240に対応付けられ、データパーティション1252は管理パーティション1250に対応付けられている。一実施形態に従うと、HCA1211およびHCA1214は、管理パーティション1230へのそれぞれのメンバーシップに基づく、データパーティション1232へのメンバーシップへのアクセス権を有する。同様に、HCA1211~1213は、管理パーティション1240へのそれぞれのメンバーシップに基づく、データパーティション1242へのメンバーシップへのアクセス権を有する。加えて、HCA1213および1215は、管理パーティション1250へのそれぞれのメンバーシップに基づく、データパーティション1252へのメンバーシップへのアクセス権を有する。

20

【0116】

一実施形態に従うと、管理パーティションは、物理HCAの仮想機能をvHCAに登録できるか否かを判断するために使用することもできる。一実施形態に従うと、vHCAは、特定の仮想マシン(VM)のために計画され構成されたホストチャネルアダプタを表わす。vHCAはVMとともにマイグレートするのに対し、VFは物理アダプタとともに留まっているという点において、vHCAは仮想機能(VF)と異なる。しかしながら、上述のように、物理HCAもvHCAも(また、より低いレベルでは、これらの(v)HCAのエンドポートも)管理パーティションのメンバになり得る。よって、一実施形態に従うと、管理パーティションへのメンバーシップを、SMが使用することにより、物理HCAからの要求である、この要求を出している物理HCAの仮想機能をvHCAに登録することを求める要求が許可可能か否かを、判断することができる。

30

【0117】

図13は、HCAおよびvHCAを管理パーティションのメンバとして有する例示的なネットワークファブリックを示す。図13に示されるように、サブネット1302はネットワークファブリック1300の一部である。HCA1310、1324、1332、および1344は、サブネット1302内でそれぞれのエンドポートを通してネットワークファブリック1300に物理的に接続されている物理HCAを表わす。HCA1310は、物理機能(PF)1312と仮想機能(VF)1314および1316とに対応付けられている。HCA1324は、PF1326とVF1328および1329とに対応付けられている。HCA1332は、PF1334とVF1336および1338とに対応付けられている。HCA1344は、PF1346とVF1348および1349とに対応付けられている。さらに、vHCA1320は、VF1314が登録され仮想マシン(V

40

50

M) 1318に対応付けられたものとして示されている(すなわち、VM1318は、ネットワークファブリック1300へのアクセス権を、vHCA1320を通して取得し、最終的には物理HCA1310を通して取得する)。vHCA1340は、VF1337が登録されVM1338に対応付けられている。

【0118】

引き続き図13を参照して、図示の通り、HCA1310および1324ならびにvHCA1320は、管理パーティション1350のメンバである。加えて、HCA1332および1344ならびにvHCA1340は、管理パーティション1360のメンバである。結果として、HCA1310および1324ならびにvHCA1320は各々管理パーティション1350のメンバであるため、vHCA1320に、HCA1310のVF1314もしくは1316を、またはHCA1324のVF1328もしくは1329を正式に登録することができる。同様に、HCA1332および1344ならびにvHCA1340が各々管理パーティション1360のメンバであるため、vHCA1340に、HCA1330のVF1336もしくは1338を、またはHCA1344のVF1348もしくは1349を、正式に登録することができる。

10

【0119】

上述のように、ファブリックレベルのファブリックデータベースは、ファブリックおよびファブリックリソースに関連する情報を保持し、ファブリックマネージャによって管理される。一実施形態に従うと、ファブリックデータベースは、ファブリックリソースのインベントリ内の「完全な知識」を有することができる(すなわち、ネットワークファブリックの一部であるすべてのリソースは、少なくともファブリックデータベースに保持されている記録によって表わされる)。さらに、ファブリック内の各リソースに対応付けられたアクセス権およびネームスペースは、ファブリックデータベースに格納されても、ファブリックデータベースに含まれる情報および関係から導出することも、可能である。

20

【0120】

たとえば、一実施形態に従うと、管理パーティションへのメンバーシップおよび/または管理パーティションに対するリソースの対応付けに関する情報は、ファブリックデータベースに格納できる。この情報を保持するテーブルと、これらのテーブルをまとめてリンクする関係とは、ファブリックデータベースのサブセットであってもよく、管理パーティションレジストリと呼ぶことができる。一実施形態に従うと、管理パーティションレジストリは、管理パーティショングループリソースを集めたものである。たとえば、管理パーティションレジストリ内の管理パーティショングループは、特定の管理パーティションのHCAメンバ(vHCAを含む)と、対応付けられたリソースとを集めたものであってもよく、このグループは、当該特定の管理パーティションを規定するP\_Keyによってリンクアップされる。加えて、管理パーティショングループのメンバおよび対応付けられたリソースを、メンバHCAもしくはvHCAそれぞれに対するGUIDもしくはvGUIDなどのキー、または、対応付けられたデータパーティションに対するP\_Keyを用いて、レジストリ内でリンクアップすることができる。管理パーティションのP\_Keyと、メンバまたは対応付けられたリソースの固有識別子との間の関係はそれぞれ、管理パーティションへのメンバーシップまたは対応付けを規定し、管理パーティションレジストリによって維持され、より高いレベルではファブリックデータベースによって維持される。

30

40

【0121】

一実施形態に従うと、管理パーティションレジストリのすべてまたは一部は、記録としてSMのキャッシュに保持されていてもよい。たとえば、特定のサブネットのリソースに対応する管理パーティションレジストリの記録を、当該特定のサブネットのサブネットマネージャ(たとえばマスタサブネットマネージャ)の常駐メモリのキャッシュに複製することができる。管理パーティションレジストリの記録は、(たとえばSMがブートするときは)SMによってファブリックデータベースから取出され(すなわちコピーされ)てもよく、または、ファブリックデータベースに永続化される前にキャッシュに置かれてもよい。キャッシュは、揮発性または不揮発性メモリであってもよい。いつレジストリ記録が

50

キャッシュに置かれるかには関係なく、管理パーティションレジストリのキャッシュされたコピーと、ファブリックレベルデータベースで発見された管理パーティションレジストリのコピーとの間には、同期が成立し得る。

**【 0 1 2 2 】**

クエリを受けるたびに管理パーティション情報を永続化状態から（すなわちファブリックデータベースから）取出すのではなく、SM上の高速キャッシュ内の管理パーティションレジストリのうちのすべてまたはサブネット関連部分を保持することにより、管理パーティション情報のルックアップがSMに課すオーバーヘッドは最小になる。これは特に、サブネットリソース間でアクセス権が自動的に割当てられるときのサブネット初期化中において重要である。

10

**【 0 1 2 3 】**

一実施形態に従うと、論理名または識別子は、リソースドメイン内のリソースに割当てることができる（たとえばファブリックレベルまたはドメインレベルの管理ユーザによって）。これらの論理名は、リソースドメインにとってプライベートなものであってもよい。ファブリックマネージャは、ファブリックデータベースを通して、ファブリック内で使用される固有の識別子（たとえばVGUIDおよびP\_Key）をファブリック内のリソースに与えられた論理または識別名にマッピングする関係を作成することができる。

**【 0 1 2 4 】**

たとえば、一実施形態に従うと、ファブリックデータベースは、リソースの記録、およびリソースのドメインメンバーシップおよび/または管理パーティションメンバーシップを格納することができる。論理名は、ファブリックマネージャによってリソースが発見されると、リソースに割当てることができる。これらの名称は、ファブリックデータベース内のファブリックリソースの固有識別子にリンクさせることができる。加えて、ファブリックマネージャは、ファブリックマネージャデータベース内の関係を通して、リソースドメインおよび管理パーティションへの各リソースのメンバーシップを追跡することができる。これらの記録および関係を用いて、ファブリックマネージャは、全く異なるリソースドメインおよび管理パーティションにおける同様の論理名を可能にする。一実施形態に従うと、論理ドメイン名スキームは、特定のドメインリソースがメンバである1つまたは複数のリソースドメインの階層を反映することができる。このような実施形態において、論理リソース名は、リソースがメンバである最高レベルのリソースドメインに固有であって

20

30

**【 0 1 2 5 】**

一実施形態に従うと、ファブリックにおけるリソースの識別子は、この識別子が何であるかに関わらず、管理パーティションの範囲内で固有であってもよい。そうすると、リソース名を、対応する管理パーティションでプレフィックスすることにより、グローバルな固有性（すなわちファブリックレベルで）を得ることができる。

**【 0 1 2 6 】**

図14は、一実施形態に従う、リソースドメインおよび管理ドメイン双方を有する例示的なネットワークファブリックを示す。図14に示されるように、ファブリックマネージャ1402はネットワークファブリック1400上で実行している。一実施形態に従うと、ファブリックマネージャ1402は、ネットワークファブリック1400のノード（図示せず）から実行することができる。ファブリックマネージャ1402は、ファブリックレベルアドミニストレータ1404によって管理され、ファブリックマネージャデータベース1414を含む。管理パーティションレジストリ1416は、論理名テーブル1418と同様、ファブリックマネージャデータベース1414の一部である。

40

**【 0 1 2 7 】**

引き続き図14を参照すると、サブネット1420がネットワークファブリック1400内に規定されている。サブネットマネージャ1422は、サブネット1420に対応付けられており、一実施形態に従うと、ネットワークファブリック1400における動作のためにサブネット1420が必要とするセマンティックランタイム動作を実行する。サブ

50

ネット1420が必要とするセットアップおよび管理タスクは、ファブリックレベルアドミニストレータ1404およびファブリックマネージャ1402によって実行されてもよい。

【0128】

ノード1444、1454、1474および1484は、サブネット1420の一部である。HCA1446は、ノード1444に対応付けられ、PF1448とVF1450および1452とを含む。同様に、HCA1456は、ノード1454に対応付けられ、PF1458とVF1460および1462とを含む。HCA1476は、ノード1474に対応付けられ、PF1478とVF1480および1482とを含む。さらに、HCA1486は、ノード1484に対応付けられ、PF1488とVF1490および1492とを含む。VM1440はノード1444上で実行しており、VM1470はノード1474上で実行している。vHCA1442は、VM1440のために計画および構成されており、VM1440に対応付けられ、HCA1446の仮想機能1452が登録されている。vHCA1472は、VM1470のために計画および構成されており、VM1470に対応付けられ、HCA1476の仮想機能1482が登録されている。

10

【0129】

一実施形態に従うと、HCA1446、1456、1476、および1486はドメインリソースとみなされ、各々の記録はファブリックマネージャデータベース1414に格納される。この記録は、ファブリック内のHCAリソースを特定するために使用されるGUIDなどの識別子を含み得る。さらに、vHCA1442および1472も、ドメインリソースとみなされ、各々の記録はファブリックマネージャデータベース1414に格納される。この記録は、vHCAを特定するために使用されるGUIDなどの識別子を含み得る。

20

【0130】

さらに図14を参照すると、一実施形態に従い、リソースドメイン1410およびリソースドメイン1412がファブリックマネージャ1402内に作成されている。一実施形態に従うと、ファブリックレベルアドミニストレータ1404は、リソースドメイン1410およびリソースドメイン1412を作成する役割を担う。加えて、ドメイン管理1406は、リソースドメイン1410に対応付けられたドメインレベルアドミニストレータである。同様に、ドメイン管理1408は、リソースドメイン1412に対応付けられたドメインレベルアドミニストレータである。一実施形態に従うと、ファブリックレベルアドミニストレータ1404は、リソースドメインの階層性に従う、それぞれのリソースドメインの管理として、ドメイン管理1406および1408を作成することができる。

30

【0131】

一実施形態に従うと、管理パーティション1424および管理パーティション1426は、サブネット1420内に規定されている。管理パーティション1424はリソースドメイン1410に対応付けられ、管理パーティション1426はリソースドメイン1412に対応付けられている。

【0132】

図14に示されるように、vHCA1442とHCA1446および1456とは、リソースドメイン1410のメンバである。一実施形態に従うと、管理パーティション1424はリソースドメイン1410に対応付けられているので、vHCA1442とHCA1446および1456とが、リソースドメイン1410のメンバとして追加されると、これらは管理パーティション1424のメンバにもなり、管理パーティション1424を規定するP\_KeyとHCA1446、1456およびvHCA1442の識別子との間の関係が、管理パーティションレジストリ1416内に作成される。一実施形態に従うと、この関係は、HCA1446、1456およびvHCA1442を、管理パーティション1424のメンバとして定義する。

40

【0133】

同様に、vHCA1472とHCA1476および1486とは、リソースドメイン1

50

412のメンバである。一実施形態に従うと、管理パーティション1426はリソースドメイン1410に対応付けられているので、vHCA1472とHCA1466および1486とが、リソースドメイン1412のメンバとして追加されると、これらは管理パーティション1426のメンバにもなり、管理パーティション1426を規定するP\_Keyと、HCA1476、1486およびvHCA1472の識別子との間の関係が、管理パーティションレジストリ1416内に生成される。一実施形態に従うと、この関係は、HCA1476、1486およびvHCA1472を、管理パーティション1426のメンバとして定義する。

#### 【0134】

上述のように、一実施形態に従うと、VM1440(vHCA1442を含む)、ノード1444(HCA1446を含む)およびノード1454(HCA1456を含む)は、リソースドメイン1410のメンバである。本発明の一実施形態において、ファブリックレベルアドミニストレータ1404は、ノード1444およびノード1454をリソースドメイン1410に追加する役割を担う。たとえば、ファブリックレベルアドミニストレータ1404は、ファブリックマネージャ1402のインターフェイスを通して、ノード1444および1454をリソースドメイン1410に追加することができる。ファブリックレベルアドミニストレータ1404がノード1444および1454をリソースドメイン1410に追加すると、ドメイン管理1406はノード1444および1454に対して管理タスクを実行することができる。しかしながら、リソースドメインの階層スキームに従うと、ドメイン管理1406は、ノード1444および1454に対して、これらリソースドメイン1410に追加される前に(すなわちこれらがより高いレベルのデフォルトドメイン(図示せず)のメンバである間)、管理タスクを実行することはできない。さらに、一実施形態に従うと、ドメイン管理1408は、ノード1444および1454に対して管理タスクを実行できない。なぜなら、ノード1444および1454は、ドメイン管理1408が対応付けられていないパラレルレベルのリソースドメインのメンバであるからである。

#### 【0135】

引き続き図14を参照すると、一実施形態に従い、管理パーティション1424および1426がサブネット1420内に規定されている。リソースドメインの階層スキームに従い、一実施形態では、管理パーティション1424および1426はファブリックレベルアドミニストレータ1404によって規定されたものである。別の実施形態では、ドメイン管理1406が管理パーティション1424を規定し、ドメイン管理1408が管理パーティション1426を規定している。一実施形態に従うと、管理パーティション1424はリソースドメイン1410に対応付けられ、管理パーティション1426はリソースドメイン1412に対応付けられる。上述のように、一実施形態に従うと、管理パーティション1424および1426はそれぞれ、サブネットレベルにおいてリソースドメイン1410および1412を表わす。それぞれのリソースドメインに対応付けられることに加えて、管理パーティション1424および1426は、それぞれドメイン管理1406および1408(すなわち、管理パーティション各々が対応付けられているリソースドメインの対応する管理ユーザ)に対応付けられる。上述のように、一実施形態に従うと、管理パーティションとドメインレベル管理との間のこの対応付けによって、サブネットレベルにおけるマルチテナント環境内のユーザとロールの分離を保証することができる。

#### 【0136】

一実施形態に従い、データパーティション1428および1430がサブネット1420内に規定されている。リソースドメインの階層スキームに従い、一実施形態では、データパーティション1428および1430がファブリックレベルアドミニストレータ1404によって規定されている。別の実施形態では、ドメイン管理1406がデータパーティション1428を規定し、ドメイン管理1408がデータパーティション1430を規定している。図14に示されるように、データパーティション1428は管理パーティション1424に対応付けられ、データパーティション1430は管理パーティション14

10

20

30

40

50

26に対応付けられている。加えて、先に説明し図14に示しているように、HCA1446、1456およびvHCA1442は、管理パーティション1424のメンバである。結果として、一実施形態に従うと、HCA1446、1456およびvHCA1442は、データパーティション1428が対応付けられている管理パーティション(すなわち管理パーティション1424)のメンバであるため、データパーティション1428へのアクセス許可を有する。

#### 【0137】

一実施形態に従い、データパーティション1428が管理パーティション1424に対応付けられると、データパーティション1428の識別子(たとえばデータパーティション1428のP\_Key)と、管理パーティション1424のP\_Keyとの間の関係が、管理パーティションレジストリ1416内に作成される。この関係は、データパーティション1428を、管理パーティション1424に対応付けられているものとして定義する。同様に、データパーティション1430が管理パーティション1426に対応付けられると、データパーティション1430の識別子(たとえばデータパーティション1430のP\_Key)と管理パーティション1426のP\_Keyとの間の関係が管理パーティションレジストリ1416内に作成される。この関係は、データパーティション1430を、管理パーティション1426に対応付けられているものとして定義する。

#### 【0138】

一実施形態に従うと、データパーティション1428に加わることを求める要求をHCA1446および1456またはvHCA1442から受けた場合、SM1422は、管理パーティションレジストリ1416をチェックし、HCA1446および1456ならびにvHCA1442が管理パーティション1424のメンバであることと、データパーティション1428が管理パーティション1424に対応付けられていることを見出すことができる。そうすると、SM1422は、HCA1446および1456ならびにvHCA1442が管理パーティション1424のメンバでありデータパーティション1428が管理パーティション1424に対応付けられていることに基づいて、HCA1446および1456ならびにvHCA1442が、データパーティション1428のメンバになることを許可することができる。HCA1446および1456ならびにvHCA1442をデータパーティション1428に加わらせるのに、ファブリックレベルアドミニストレータ1404またはドメインレベルアドミニストレータ1406からのマッピングは不要であろう。

#### 【0139】

加えて、vHCA1442には、HCA1446のVF1452もしくは1450を、または、HCA1456のVF1462もしくは1460を登録することができる。なぜなら、HCA1446および1456ならびにvHCA1442は、管理パーティション1424のメンバであるからである(vHCA1442はVF1452が登録されたものとして示されている)。ここでもまた、SM1422は、管理パーティションレジストリ1416をチェックしてHCA1446および1456ならびにvHCA1442が管理パーティション1424のメンバであることを見出すことができる。HCA1446および1456ならびにvHCA1442が管理パーティション1424のメンバであることを発見すると、SM1422は、いかなるファブリックユーザからも介入されることなく、vHCA1442に、VF1452、1450、1462、および1460のうちのいずれかを登録することができる。

#### 【0140】

##### 仮想マシンファブリックプロファイル

上述のように、効率的なハードウェアリソースの利用およびスケーラビリティを改善するために、IBファブリック等のファブリックにおいて仮想マシン(VM)を用いることができる。しかしながら、仮想マシン(VM)のライブマイグレーションは、これらのソリューションにおいて用いられるアドレス指定およびルーティングスキームのせいで、依然として問題である。一実施形態に従う方法およびシステムは、このようなVMのマイグ

10

20

30

40

50

レーションの問題を克服することを目的とし、アドレス指定スキームをサポートすることが可能な、予め定められ、非常に入手し易く、かつ物理的な場所に左右されない仮想マシンファブリックプロファイルを、容易にする。一実施形態に従うと、VMファブリックプロファイルにより、ファブリック接続性を用いてVMの集中型セットアップおよび構成管理を可能にし、かつ、SR-IOVに基づいて、VMに対し、最適化されたVMマイグレーションをサポートする。一実施形態に従うと、VMファブリックプロファイルは、仮想マシンに関する、詳細なファブリック構成情報の、単一の集中型リポジトリを表わす。ファブリックマネージャに対応付けられたデータベース（たとえばファブリックデータベース）は、VMファブリックプロファイルを構成する情報を永続化することができる。

#### 【0141】

一実施形態に従うと、VMファブリックプロファイルは、IBファブリック等のネットワークファブリックにおいて、仮想マシン識別子（VM-id）により、特定することができる。一実施形態において、VM-idは、ファブリック全体にわたって固有であってもよい、汎用一意識別子（universally unique identifier：UUID）である固有の128ビット数である。しかしながら、必要なVM-idの固有性は、異なる態様で管理される複数のVMマネージャドメインにわたる固有性のみである（たとえばVM-idは管理パーティション内で固有であってもよい）。したがって、他の実施形態において、VM-idは、少なくともこのようなドメインにわたって固有である何か他の適切なタイプのIDであってもよい。一実施形態に従うと、ファブリックレベルまたはサブネットレベルのすべての管理エンティティは、VMファブリックプロファイルのVM-idを参照することにより、VMファブリックプロファイルに関する情報をルックアップする。

#### 【0142】

図15は、VMファブリックプロファイル情報を格納するための例示的なデータベース構造を示す。図15は、従来のリレーショナルデータベース設計のいくつかのテーブルを示す。しかしながら、何らかの適切なデータ構造を用いて、VMファブリックプロファイルデータ（たとえばフラットファイルテーブル、または境界が定められた構造等）を格納することができる。図15において、アスタリスク（\*）はキー値を示す。図15は、VMファブリックプロファイルデータを、より大きなファブリックデータベース1500の一部として示しているが、その他の実施形態において、VMプロファイルデータは自身のデータベースに含まれていてもよい、または、ファブリックデータベース1500にアクセスできる別のデータベースであってもよい。

#### 【0143】

図15に示されるように、VMファブリックプロファイルのコンテンツは、ルックアップキーとして使用される仮想マシン識別子（VM-id）1502と、ファブリックの管理における使い易さのためおよび品質改善のための論理名1504と、たとえばVMファブリックプロファイルと、ファブリックに対して規定されたその他のプロファイルとを区別するために使用されるプロファイルタイプ1506と、ファブリックに対して規定されたすべてのプロファイルからなる組の中の固有idであるプロファイルID1508と、最高シーケンスが最も最近の更新を表わすプロファイルのシーケンス番号であってもよいコンテンツ更新エニューメレータ（CUE）1510とを含み得るが、これらに限定されない。図15に示されるように、このVMファブリックプロファイルのコンテンツは、VMプロファイルテーブル1512に格納することができ、VM-idは、各VMファブリックプロファイルを特定するための固有のキーの働きをする。

#### 【0144】

上述のように、仮想HCA（vHCA）を、物理HCAの仮想機能とともに使用することにより、VMへのネットワークアクセスを提供してもよい。一実施形態に従うと、各VMファブリックプロファイルは、少なくとも1つのvHCAにも対応付けられる。vHCAは、特定のVMに対して計画および構成することができ、この構成は、vHCAが構成されたVMのファブリックプロファイルとともに、含まれ格納されてもよい。そうすると、構成されたvHCAは、物理HCAによって規定することができ物理HCAとともに常

10

20

30

40

50

駐できる仮想機能と異なり、仮想マシンとともにマイグレートすることができる。さらに図15を参照すると、vHCAは、ファブリックマネージャデータベースにおいて、VM-id1502とvHCAインスタンスID1514との固有の組み合わせとして表わすことができる。加えて、この組み合わせは、VM-idキー1502を用いて、VMプロファイルテーブル1512に関連するvHCAテーブルに格納することができる。

**【0145】**

物理HCAと同様に、vHCAは複数の(仮想)ポートを有し得る。一実施形態に従うと、これらの仮想ポートは、物理ポートと全く同じように、ネットワーク環境におけるエンドポートとして働くことができる。一例として、vHCAポートを含むすべてのエンドポートに、GUID(たとえばIBネットワークで使用される64ビットGUID)を割当てることができる。このGUIDを用いて、SMのルーティングテーブルのLID宛先アドレスを要求することができる。一実施形態に従うと、仮想GUID(vGUID)は、各vHCAポートの現在のファブリックアドレスを表わし得る。一実施形態において、vGUIDは、上述のように、VMファブリックプロファイルに割当てられかつVMファブリックプロファイルとともに格納されるGUIDのリストからのvHCAポートに割当てることができる。vGUIDは、ファブリックマネージャGUIDポリシーに従ってファブリックマネージャにより所有され制御されるGUIDの専用プールからのファブリックプロファイルに割当てることができる。たとえば、フリーでファブリックの広さの固有vGUIDは、ファブリックプロファイルがファブリックマネージャに作成されたときにVMのために構成されたvHCAに割当てることができる。

**【0146】**

引き続き図15を参照して、一実施形態に従うと、vHCAポートは、ファブリックマネージャデータベースにおいて、vHCAインスタンスID1514とvGUID1522との固有の組み合わせとして表わすことができる。vHCA構成はまた、vHCAポート番号1520を含み得る。この構成は、vHCAインスタンスID\*1514キーによってvHCAテーブルに関連付けることができ最終的にはVM-idキー1502によってVMプロファイルテーブル1512に関連付けることができるvHCAポートテーブル1518に格納することができる。

**【0147】**

一実施形態に従うと、vHCAは、データパーティションおよび管理パーティション双方のメンバであってもよい。vHCAのパーティションメンバーシップは、VMファブリックプロファイルにおいて、vHCAがメンバであるパーティション(管理パーティションおよびデータパーティション双方)にvHCA記録をリンクするファブリックデータベース内の関係によって表わすことができる。一実施形態において、たとえばvGUIDキー1522は、データとネットワークファブリックにおいて規定された管理パーティションP\_Keyとを含むテーブル(図示せず)に関連付けることができる。一実施形態において、vGUIDキーは、関係により、(上記)管理パーティションレジストリにリンクされる。一実施形態において、これに代えてまたはこれに加えて、管理パーティションレジストリをvHCAインスタンスID1514にリンクする関係があってもよい。これらの関係により、ファブリックコンポーネントは、vHCAがメンバであるのはどのデータパーティションおよび管理パーティションであるかを識別することができる。

**【0148】**

図15は、専ら例示を目的として提供されており、当業者は、VMファブリックプロファイルを構成するデータの記憶装置を設計および監理する数多くのやり方が存在することを理解するであろう。さらに、上記VMファブリックプロファイルコンテンツのリストは、限定ではなく例示を意図したものであり、仮想マシンファブリックプロファイルのその他の実施形態は、より多くのまたはより少ないその他のコンテンツを含み得る。一実施形態に従うと、図15に示されるデータベースコンポーネントは、ファブリックについてその他の関連情報を保持するより大きなファブリックマネージャデータベースの一部であってもよく、その他のテーブルと相互に関連付けることにより、ファブリックマネージャお

10

20

30

40

50

よびその他のファブリックコンポーネントの機能を向上させることができる。

【0149】

一実施形態において、ユーザは、ファブリックマネージャのインターフェイスを通して仮想マシンファブリックプロファイルと対話する。たとえば、ユーザは、ファブリックマネージャを通してVMファブリックプロファイルを作成、編集、および削除し得る。一実施形態に従うと、ファブリックプロファイル情報の一部は、VMファブリックプロファイル（たとえばVMファブリックプロファイルの論理名）を作成または編集するユーザによって提供され、この情報のその他の部分は、ファブリックマネージャによって、または、ファブリックマネージャが作成されているサブネットのローカルSMによって提供される（たとえばVM - id、vHCAインスタンスID、またはvGUID）。

10

【0150】

一実施形態に従うと、仮想マシンファブリックプロファイルの作成は、VMファブリックプロファイルが対応付けられているファブリックリソースの管理特権を表わす管理コンテキストにおいて行なわれてもよい。たとえば、ドメイン管理は、VMファブリックプロファイルに対して構成されているvHCAと同じリソースドメイン（および同じ管理パーティション）のメンバであるHCAを有するノードが使用するために、VMファブリックプロファイルを作成する。作成されたVMファブリックプロファイルは、（論理）リソースでありかつそれが作成されたリソースドメインのメンバであるとみなされる。このように、同じ管理パーティションのメンバであるおかげで、VMファブリックプロファイルのvHCAは、リソースドメインのメンバでもあるHCAのVFのうちいずれかとともに登録される許可を得る。こうして管理オーバーヘッドは緩和される。

20

【0151】

一実施形態に従うと、ファブリックレベルまたはドメインレベルアドミニストレータユーザは、「仮想マシンマネージャ」（Virtual Machine Manager：VMM）と名付けられたファブリックマネージャのコンポーネントを用いることにより、VMファブリックプロファイルを設定アップし構成することができる。VMMは、VMファブリックプロファイルの作成においてファブリックREST APIを使用することができる。ユーザは、たとえばファブリックマネージャのGUIを通してVMMにアクセスすることができる。一実施形態に従うと、管理ユーザは、VMファブリックプロファイルに関連する特定のパラメータ（プロファイルに対応付けられるvHCAの論理名、プロファイルタイプ、および数など）を提供することができ、その他のパラメータ（VM - id、vGUID、およびプロファイルに対応付けられた各vHCAのvHCAインスタンスIDなど）は、VMMによって自動的に生成し割当てることができる。VMファブリックプロファイルの更新および削除などのその他のCRUDアクションは、ファブリックマネージャ、具体的にはVMMを通して、アドミニストレータユーザが利用することもできる。

30

【0152】

一実施形態に従うと、VMファブリックプロファイルの構築に必要なパラメータすべてがVMMを介して供給されると、ファブリックマネージャは、アドミニストレータユーザおよびVMMによって指定された属性を有するVMファブリックプロファイルオブジェクトのインスタンスを作成し、このVMファブリックプロファイルオブジェクトをファブリックレベルデータベースに永続化することができる。

40

【0153】

一実施形態に従うと、オペレーショナルネットワークファブリックにおいて、VMファブリックプロファイルは、SMのキャッシュに、1つまたは複数の記録として保持することができる。VMファブリックプロファイルは、（たとえばSMが起動したときに）SMによってファブリックデータベースから取出され（すなわちコピーされ）てもよく、または、ファブリックデータベースに永続化される前にキャッシュに置かれてもよい。このキャッシュは、揮発性メモリであっても不揮発性メモリであってもよい。VM - idを、キャッシュにクエリするためのキーとして用いることにより、特定のVMファブリックプロファイルの属性を取出してもよい。

50

## 【 0 1 5 4 】

クエリを受けるたびに、VMファブリックプロファイルを、永続化された状態から（すなわちファブリックデータベースから）取出すのではなく、SM上の高速キャッシュに保持することによって、VMファブリックプロファイル属性のルックアップがSMに対して課す可能性があるオーバーヘッドを最小にすることができる。これは、VMおよびホストの起動中に、どのHCA VMとペアにすることができどのHCA VMとペアにするかを確定するのにファブリックプロファイルデータが必要なときは、特に重要である。

## 【 0 1 5 5 】

図16は、VMファブリックプロファイルをサブネットリソースが利用できるようにするためのフローチャートである。

10

## 【 0 1 5 6 】

ステップ1610において、VMのセットアップパラメータおよび構成を含むVMファブリックプロファイルが規定される。

## 【 0 1 5 7 】

ステップ1620において、VMファブリックプロファイルはファブリックレベルデータベースに格納される。

## 【 0 1 5 8 】

ステップ1630において、VMファブリックプロファイルは、サブネットマネージャ上の高速メモリキャッシュを通して利用できるようにされる。

20

## 【 0 1 5 9 】

ステップ1640において、VMファブリックプロファイルデータが、高速メモリキャッシュに向けられたVM-idルックアップ要求に基づいてサブネットマネージャから戻される。

## 【 0 1 6 0 】

図17は、仮想マシンの仮想マシンファブリックプロファイルを作成するためのフローチャートである。

## 【 0 1 6 1 】

ステップ1710において、仮想マシン識別子（VM-id）が生成される。

ステップ1720において、仮想ホストチャンネルアダプタ（vHCA）を特定する仮想ホストチャンネルアダプタインスタンスIDが生成される。

30

## 【 0 1 6 2 】

ステップ1730において、仮想グローバル一意識別子（vGUID）が、グローバル一意識別子のプールから、vHCAの仮想ポートの現在のファブリックアドレスとして割当てられる。

## 【 0 1 6 3 】

ステップ1740において、管理パーティションを規定するP\_KeyとvGUIDとの間の第1の関係が作成される。このP\_KeyとvGUIDとの間の関係は、vGUIDが現在のファブリックアドレスとして割当てられている仮想ポートを、P\_Keyによって規定される管理パーティションのメンバとして規定する。

## 【 0 1 6 4 】

ステップ1750において、VM-idとvHCAインスタンスIDとの間の第2の関係が作成される。この第2の関係により、vHCAインスタンスIDを、VM-idへのアクセスを通して取出すことができる。

40

## 【 0 1 6 5 】

ステップ1760において、VM-idとvGUIDとの間の第3の関係が作成される。この第3の関係により、vGUIDを、VM-idへのアクセスを通して取出すことができる。

## 【 0 1 6 6 】

ステップ1770において、VM-id、仮想ホストチャンネルアダプタインスタンスID、vGUID、第1の関係、第2の関係、および第3の関係を、VM-id、仮想ホス

50

トチャンネルアダプタ、vGUI D、第1の関係、第2の関係、および第3の関係を保存するフォーマットで永続化する。

【0167】

本発明の多数の特徴は、ハードウェア、ソフトウェア、ファームウェア、またはこれらを組合わせたものにおいて、これを用いて、またはこれに支援されて、実施することができる。したがって、本発明の特徴は、処理システム（たとえば1つ以上のプロセッサを含む）を用いて実現し得る。

【0168】

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可能媒体であるコンピュータプログラムプロダクトにおいて、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。

【0169】

この発明の特徴は、機械読取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。

【0170】

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タ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。

【0171】

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。

【0172】

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。

【0173】

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。

【0174】

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

10

20

30

40

50

。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。

【図1】

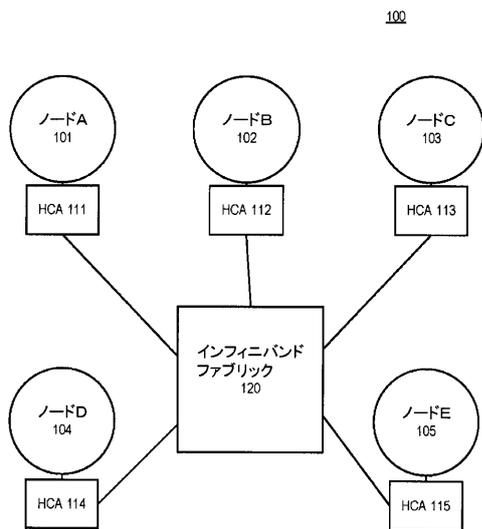


FIGURE 1

【図2】

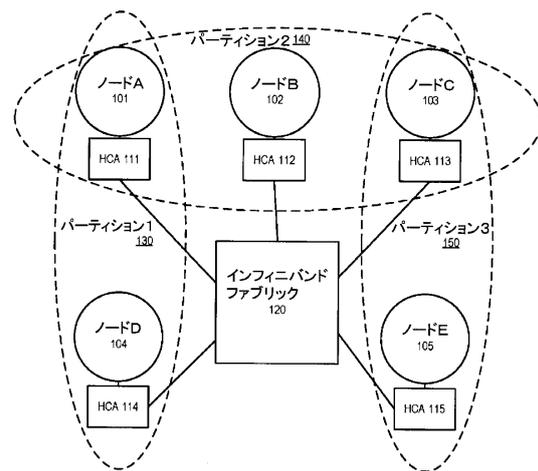


FIGURE 2

【図3】

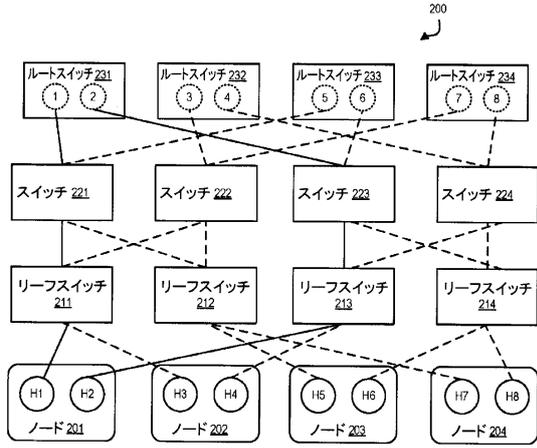


FIGURE 3

【図4】

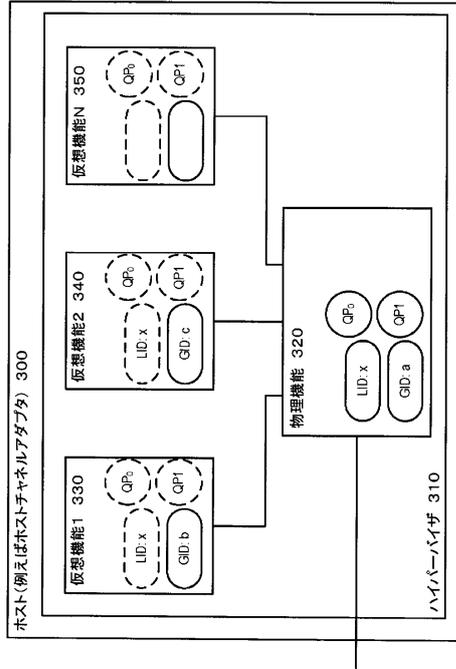


FIGURE 4

【図5】

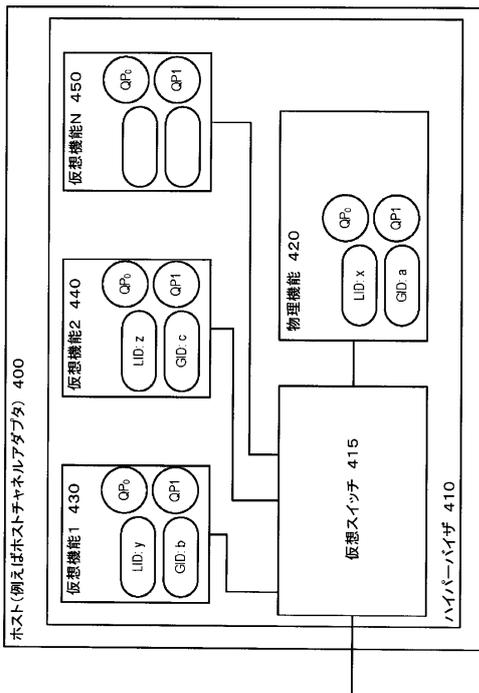


FIGURE 5

【図6】

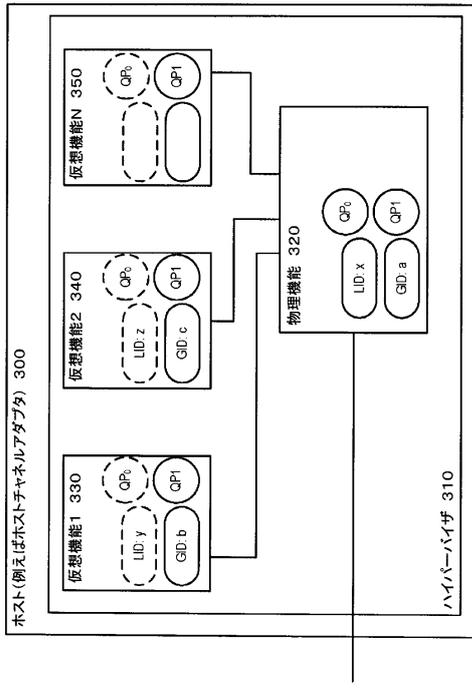


FIGURE 6

【図 7】

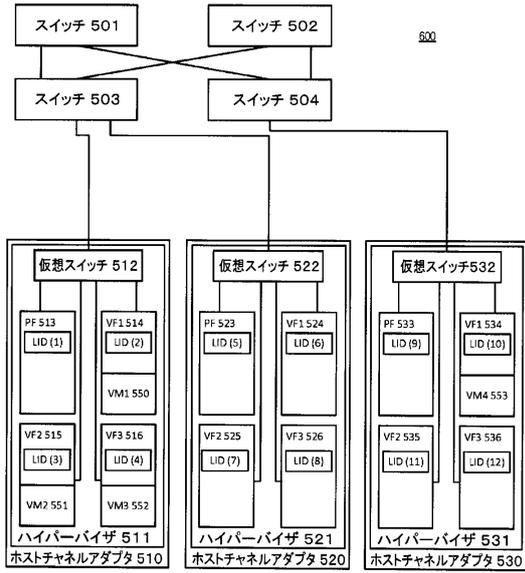


FIGURE 7

【図 8】

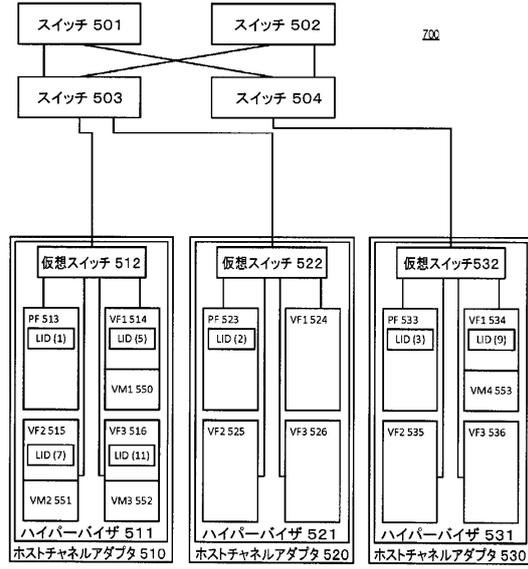


FIGURE 8

【図 9】

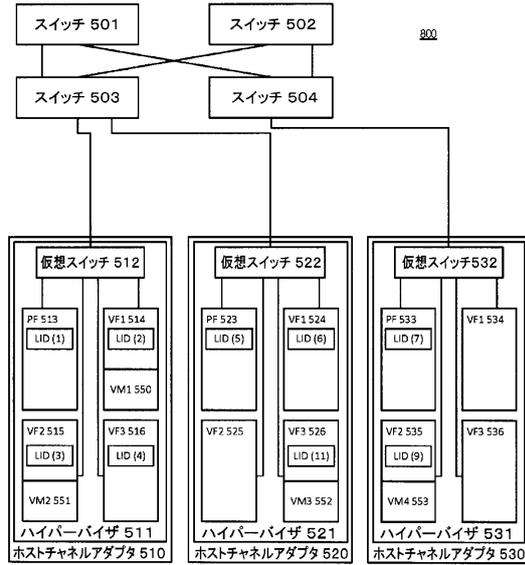


FIGURE 9

【図 10】

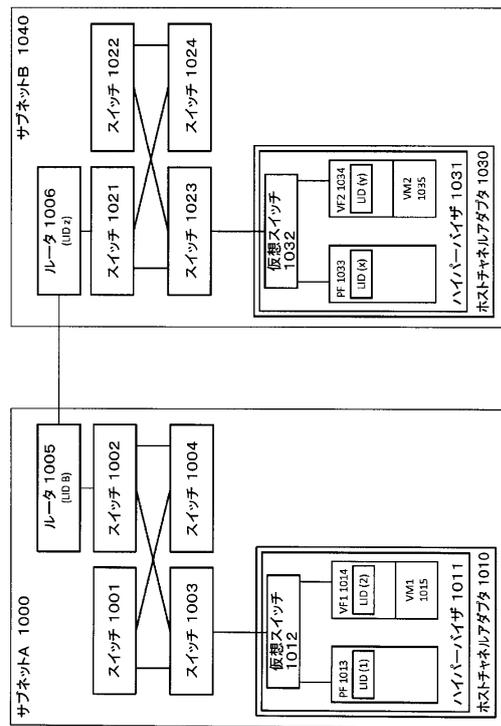


FIGURE 10

【図 1 1】

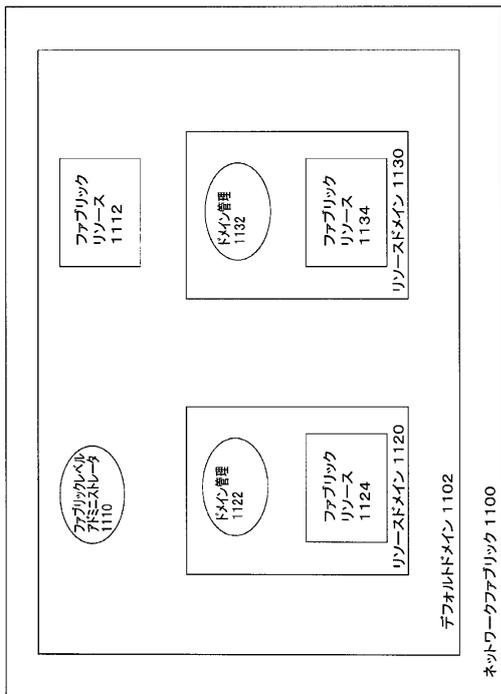


FIGURE 11

【図 1 2】

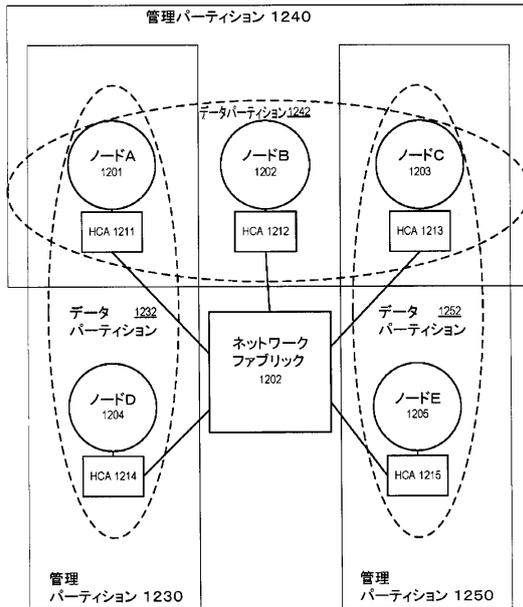


FIGURE 12

【図 1 3】

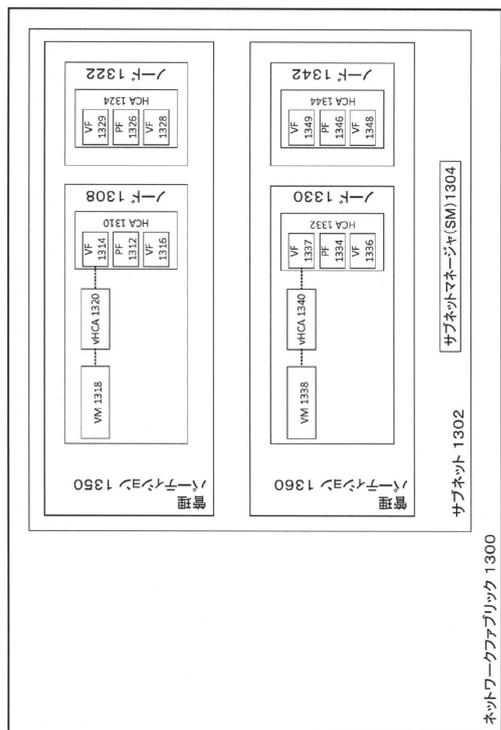


FIGURE 13

【図 1 4】

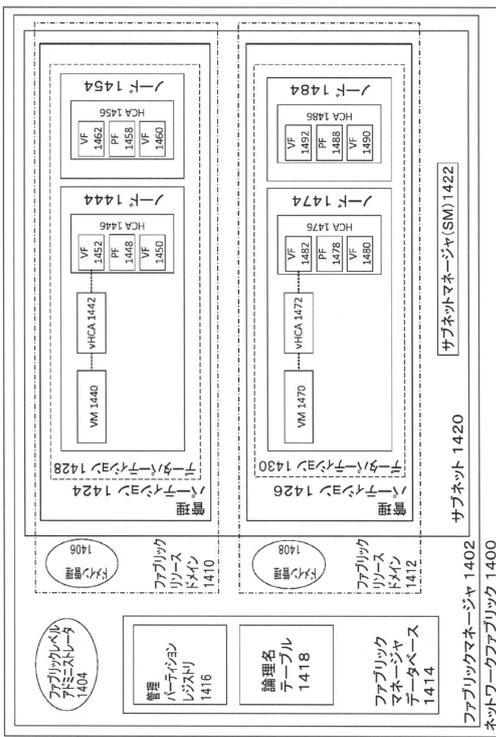


FIGURE 14

【図15】

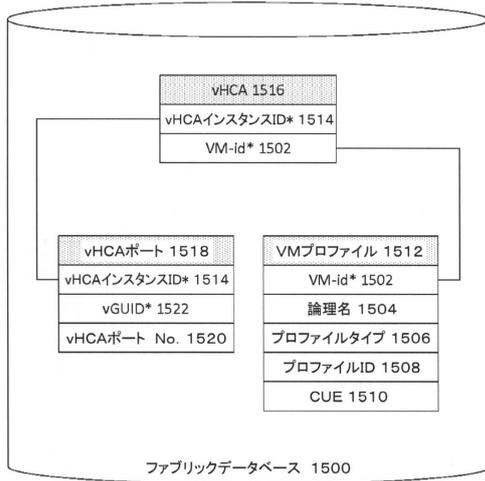


FIGURE 15

【図16】

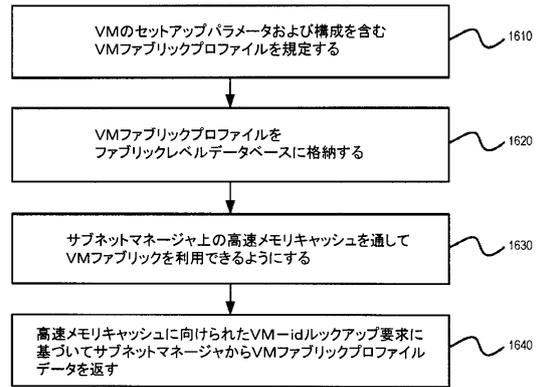


FIGURE 16

【図17】

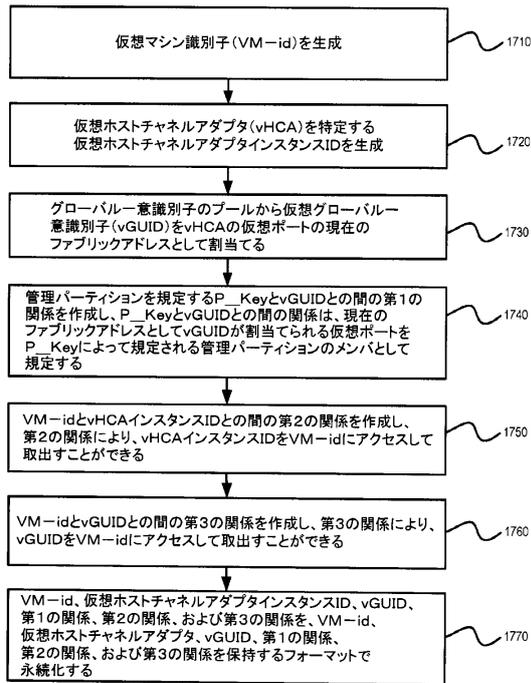


FIGURE 17

---

フロントページの続き

(72)発明者 ホレン, リネ  
ノルウェー、1900 フェツンド、ビタセン、17

審査官 井上 宏一

(56)参考文献 特表2015-518602(JP, A)  
特表2004-531175(JP, A)  
特表2015-515683(JP, A)

(58)調査した分野(Int.Cl., DB名)  
G06F 9/455 - 9/54  
G06F 15/177