

(19) 日本国特許庁 (JP)

(12) 公表特許公報 (A)

(11) 特許出願公表番号

特表2019-503598

(P2019-503598A)

(43) 公表日 平成31年2月7日 (2019. 2. 7)

(51) Int. Cl.	F I	テーマコード (参考)
H04L 12/931 (2013.01)	H04L 12/931	5K030
H04L 12/70 (2013.01)	H04L 12/70 100Z	

審査請求 未請求 予備審査請求 未請求 (全 42 頁)

(21) 出願番号 特願2018-504731 (P2018-504731) (86) (22) 出願日 平成29年1月26日 (2017. 1. 26) (85) 翻訳文提出日 平成30年1月29日 (2018. 1. 29) (86) 国際出願番号 PCT/US2017/015156 (87) 国際公開番号 W02017/132392 (87) 国際公開日 平成29年8月3日 (2017. 8. 3) (31) 優先権主張番号 62/287, 704 (32) 優先日 平成28年1月27日 (2016. 1. 27) (33) 優先権主張国 米国 (US) (31) 優先権主張番号 15/412, 972 (32) 優先日 平成29年1月23日 (2017. 1. 23) (33) 優先権主張国 米国 (US) (31) 優先権主張番号 15/415, 497 (32) 優先日 平成29年1月25日 (2017. 1. 25) (33) 優先権主張国 米国 (US)	(71) 出願人 502303739 オラクル・インターナショナル・コーポレーション 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, オラクル・パークウェイ500 (74) 代理人 110001195 特許業務法人深見特許事務所 (72) 発明者 ヨンセン, ビョルン・ダグ ノルウェー、0687 オスロ、ビルベルクグレンダ、9 (72) 発明者 スリニバサン, アルビンド アメリカ合衆国、95129 カリフォルニア州、サン・ノゼ、ハッピー・バレー・アベニュー、1075
---	--

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スケーラブルなビットマップに基づく P_Key テーブルをサポートするためのシステムおよび方法

(57) 【要約】

高性能コンピューティング環境においてスケーラブルなビットマップに基づく P_Key テーブルをサポートするためのシステムおよび方法。ある方法は、1つ以上のスイッチ、複数のホストチャネルアダプタ、および複数のエンドノードを備える少なくとも1つのサブネットを提供することができる。この方法は、複数のエンドノードを複数のパーティションの少なくとも1つに関連付けることができ、複数のパーティションの各々は P_Key 値に関連付けられる。この方法は、1つ以上のスイッチの各々を、複数のビットマップに基づく P_Key テーブルのうちのあるビットマップに基づく P_Key テーブルに関連付けることができる。この方法は、ホストチャネルアダプタの各々を、複数のビットマップに基づく P_Key テーブルのうちのあるビットマップに基づく P_Key テーブルに関連付けることができる。

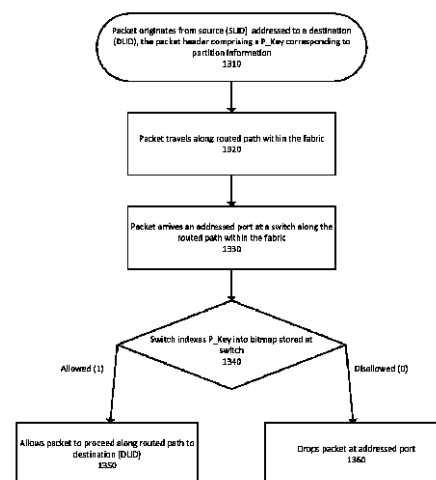


FIGURE 13

【特許請求の範囲】**【請求項 1】**

高性能コンピューティング環境においてスケーラブルなビットマップに基づく P __ K e y テーブル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、

少なくとも 1 つのサブネットとを備え、前記少なくとも 1 つのサブネットは、

1 つ以上のスイッチを含み、前記 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、前記 1 つ以上のスイッチの各々は、複数のスイッチポートを含み、前記少なくとも 1 つのサブネットはさらに、

複数のホストチャネルアダプタを含み、各ホストチャネルアダプタは、少なくとも 1 つのホストチャネルアダプタポートを含み、前記少なくとも 1 つのサブネットはさらに、

複数のエンドノードを含み、前記複数のエンドノードの各々は、前記複数のホストチャネルアダプタのうちの少なくとも 1 つのホストチャネルアダプタに関連付けられ、

前記複数のエンドノードの各々は、複数のパーティションのうちの少なくとも 1 つに関連付けられ、

前記複数のパーティションの各々は、P __ K e y 値に関連付けられ、

前記複数のスイッチポートのうちのあるスイッチポートは、複数のビットマップに基づく P __ K e y テーブルのうちのあるビットマップに基づく P __ K e y テーブルに関連付けられ、

前記複数のホストチャネルアダプタポートのうちのあるホストチャネルアダプタポートは、前記複数のビットマップに基づく P __ K e y テーブルのうちのあるビットマップに基づく P __ K e y テーブルに関連付けられる、システム。

【請求項 2】

前記 1 つ以上のスイッチの 1 つまたは前記複数のホストチャネルアダプタの 1 つで実行されるサブネットマネージャをさらに備え、

前記サブネットマネージャは、前記 1 つ以上のスイッチの各々上の複数のポートを通る許可されたトラフィックおよび許可されないトラフィックを判断する、請求項 1 に記載のシステム。

【請求項 3】

前記サブネットマネージャは、前記 1 つ以上のスイッチの各々上の複数のポートの各々を通る許可されたトラフィックおよび許可されないトラフィックの前記判断に基づいて、前記ビットマップに基づく P __ K e y テーブルの各々を構成する、請求項 2 に記載のシステム。

【請求項 4】

前記 1 つ以上のスイッチのうちのあるスイッチが、P __ K e y 値を含む少なくともヘッダを含むパケットを、アドレス指定されたポートで受信すると、前記 1 つ以上のスイッチのうちの前記あるスイッチは、許可されたまたは許可されないインジケータを受信するよう、前記パケットの前記ヘッダに含まれる前記 P __ K e y 値を、関連付けられるビットマップに基づく P __ K e y 値で索引付けする、請求項 1 に記載のシステム。

【請求項 5】

許可された値を受信すると、前記スイッチは、前記パケットが前記アドレス指定されたポートを通過することを可能にする、請求項 4 に記載のシステム。

【請求項 6】

許可されない値を受信すると、前記スイッチは、前記パケットを前記アドレス指定されたポートでドロップする、請求項 4 に記載のシステム。

【請求項 7】

前記 1 つ以上のサブネットは、2 つ以上のサブネットを含み、前記 2 つ以上のサブネットの各々は、前記 2 つ以上のサブネットの各々において少なくとも 1 つのルータによって相互接続される、先行する請求項のいずれか 1 つに記載のシステム。

【請求項 8】

高性能コンピューティング環境においてスケーラブルなビットマップに基づく P __ K e y テーブル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少なくとも 1 つのサブネットを提供することを備え、前記少なくとも 1 つのサブネットは、

1 つ以上のスイッチを含み、前記 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、前記 1 つ以上のスイッチの各々は、複数のスイッチポートを含み、前記少なくとも 1 つのサブネットはさらに、

複数のホストチャネルアダプタを含み、各ホストチャネルアダプタは、少なくとも 1 つのホストチャネルアダプタポートを含み、前記少なくとも 1 つのサブネットはさらに、

複数のエンドノードを含み、前記複数のエンドノードの各々は、前記複数のホストチャネルアダプタのうちの少なくとも 1 つのホストチャネルアダプタに関連付けられ、前記方法はさらに、

前記複数の物理ホスト仮想マシンの各々を複数のパーティションの少なくとも 1 つに関連付けることを備え、前記複数のパーティションの各々は P __ K e y 値に関連付けられ、前記方法はさらに、

前記複数のスイッチポートのうちの あるスイッチポートを、複数のビットマップに基づく P __ K e y テーブルのうちの あるビットマップに基づく P __ K e y テーブルに関連付けることと、

前記複数のホストチャネルアダプタポートのうちの あるホストチャネルアダプタポートを、前記複数のビットマップに基づく P __ K e y テーブルのうちの あるビットマップに基づく P __ K e y テーブルに関連付けることとを備える、高性能コンピューティング環境においてスケーラブルなビットマップに基づく P __ K e y テーブルをサポートするための方法。

【請求項 9】

前記 1 つ以上のマイクロプロセッサを含む前記 1 つ以上のコンピュータにおいて、

前記 1 つ以上のスイッチの 1 つまたは前記複数のホストチャネルアダプタの 1 つで実行されるサブネットマネージャを提供することと、

前記サブネットマネージャによって、前記 1 つ以上のスイッチの各々上の複数のポートを通る許可されたトラフィックおよび許可されないトラフィックを判断することとをさらに備える、請求項 8 に記載の方法。

【請求項 10】

前記サブネットマネージャによって、前記 1 つ以上のスイッチの各々上の複数のポートの各々を通る許可されたトラフィックおよび許可されないトラフィックの前記判断に基づいて、前記ビットマップに基づく P __ K e y テーブルの各々を構成することをさらに備える、請求項 9 に記載の方法。

【請求項 11】

前記 1 つ以上のスイッチのうちの あるスイッチが、P __ K e y 値を含む少なくともヘッダを含むパケットを、アドレス指定されたポートで受信すると、許可されたまたは許可されないインジケータを受信するよう、前記 1 つ以上のスイッチのうちの 前記あるスイッチによって、前記パケットの前記ヘッダに含まれる前記 P __ K e y 値を、関連付けられるビットマップに基づく P __ K e y 値で索引付けすることをさらに備える、請求項 8 ~ 10 のいずれか 1 つに記載の方法。

【請求項 12】

許可された値を受信すると、前記スイッチによって、前記パケットが前記アドレス指定されたポートを通過することを可能にすることをさらに備える、請求項 11 に記載の方法。

【請求項 13】

許可されない値を受信すると、前記スイッチによって、前記パケットを前記アドレス指定されたポートでドロップすることをさらに備える、請求項 11 に記載の方法。

【請求項 14】

前記 1 つ以上のサブネットは、2 つ以上のサブネットを含み、前記 2 つ以上のサブネットの各々は、前記 2 つ以上のサブネットの各々において少なくとも 1 つのルータによって相互接続される、請求項 8 ~ 13 のいずれか 1 つに記載の方法。

【請求項 15】

高性能コンピューティング環境においてスケーラブルなビットマップに基づく P __ K e y テーブルをサポートするための命令をそこに記憶して含む、非一時的なコンピュータ可読記憶媒体であって、前記命令は、1 つ以上のコンピュータによって読み取られ実行されると、前記 1 つ以上のコンピュータに、

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少なくとも 1 つのサブネットを提供することを含むステップを実行させ、前記少なくとも 1 つのサブネットは、

1 つ以上のスイッチを含み、前記 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、前記 1 つ以上のスイッチの各々は、複数のスイッチポートを含み、前記少なくとも 1 つのサブネットはさらに、

複数のホストチャネルアダプタを含み、各ホストチャネルアダプタは、少なくとも 1 つのホストチャネルアダプタポートを含み、前記少なくとも 1 つのサブネットはさらに、

複数のエンドノードを含み、前記複数のエンドノードの各々は、前記複数のホストチャネルアダプタのうちの少なくとも 1 つのホストチャネルアダプタに関連付けられ、前記命令は、さらに、1 つ以上のコンピュータによって読み取られ実行されると、前記 1 つ以上のコンピュータに、

前記複数のエンドノードの各々を複数のパーティションの少なくとも 1 つに関連付けることを含むステップを実行させ、前記複数のパーティションの各々は P __ K e y 値に関連付けられ、前記命令は、さらに、1 つ以上のコンピュータによって読み取られ実行されると、前記 1 つ以上のコンピュータに、

前記複数のスイッチポートのうちの あるスイッチポートを、複数のビットマップに基づく P __ K e y テーブルのうちの あるビットマップに基づく P __ K e y テーブルに関連付けることと、

前記複数のホストチャネルアダプタポートのうちの あるホストチャネルアダプタポートを、前記複数のビットマップに基づく P __ K e y テーブルのうちの あるビットマップに基づく P __ K e y テーブルに関連付けることとを含むステップを実行させる、非一時的なコンピュータ可読記憶媒体。

【請求項 16】

前記ステップは、さらに、前記 1 つ以上のマイクロプロセッサを含む前記 1 つ以上のコンピュータにおいて、

前記 1 つ以上のスイッチの 1 つまたは前記複数のホストチャネルアダプタの 1 つで実行されるサブネットマネージャを提供することと、

前記サブネットマネージャによって、前記 1 つ以上のスイッチの各々上の複数のポートを通る許可されたトラフィックおよび許可されないトラフィックを判断することを含む、請求項 15 に記載の非一時的なコンピュータ可読記憶媒体。

【請求項 17】

前記ステップは、さらに、前記サブネットマネージャによって、前記 1 つ以上のスイッチの各々上の複数のポートの各々を通る許可されたトラフィックおよび許可されないトラフィックの前記判断に基づいて、前記ビットマップに基づく P __ K e y テーブルの各々を構成することを含む、請求項 16 に記載の非一時的なコンピュータ可読記憶媒体。

【請求項 18】

前記ステップは、さらに、前記 1 つ以上のスイッチのうちの あるスイッチが、P __ K e y 値を含む少なくともヘッダを含むパケットを、アドレス指定されたポートで受信すると、許可されたまたは許可されないインジケータを受信するよう、前記 1 つ以上のスイッチのうちの 前記あるスイッチによって、前記パケットの前記ヘッダに含まれる前記 P __ K e y 値を、関連付けられるビットマップに基づく P __ K e y 値で索引付けすることを含む、請求項 15 ~ 17 のいずれか 1 つに記載の一時的なコンピュータ可読記憶媒体。

10

20

30

40

50

【請求項 19】

前記ステップは、さらに、許可された値を受信すると、前記スイッチによって、前記パケットが前記アドレス指定されたポートを通過することを可能にすることを含む、請求項 18 に記載の非一時的なコンピュータ可読記憶媒体。

【請求項 20】

前記ステップは、さらに、許可されない値を受信すると、前記スイッチによって、前記パケットを前記アドレス指定されたポートでドロップすることを含む、請求項 18 に記載の非一時的なコンピュータ可読記憶媒体。

【請求項 21】

コンピュータシステムによって実行されると、前記コンピュータシステムに請求項 8 ~ 14 のいずれか 1 つに記載の方法を実行させる、機械可読フォーマットにおけるプログラム命令を含む、コンピュータプログラム。

【請求項 22】

非一時的な機械可読データ記憶媒体に記憶される請求項 21 のコンピュータプログラムを備えるコンピュータプログラム製品。

【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。著作権所有者は、この特許文献または特許開示の、それが特許商標庁の特許ファイルまたは記録に現れているとおりの、何人による複写複製にも異議を唱えないが、それ以外の場合にはすべての著作権をどのようなものであると所有する。

【0002】

発明の分野：

本発明は、一般にコンピュータシステムに関し、特に、高性能コンピューティング環境においてスケーラブルビットマップに基づく P__K e y テーブルをサポートすることに関する。

【背景技術】

【0003】

背景：

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

【発明の概要】

【課題を解決するための手段】

【0004】

概要：

本明細書では、高性能コンピューティング環境においてスケーラブルなビットマップに基づく P__K e y テーブルをサポートするためのシステムおよび方法を説明する。ある例示的方法は、1 つ以上のスイッチを含む少なくとも 1 つのサブネットを提供し得、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、複数のスイッチポートを含み、少なくとも 1 つのサブネットはさらに、複数のホストチャネルアダプタを含み、ホストチャネルアダプタの 1 つ以上は、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の各々は、複数のホストチャネルアダプタポートを含み、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、少なくとも 1 つのサブネッ

10

20

30

40

50

トはさらに、複数の物理ホストおよびハイパーバイザを含み、複数の物理ホストおよび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の少なくとも1つのホストチャネルアダプタに関連付けられ、少なくとも1つのサブネットはさらに、複数の仮想マシンを含み、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。この方法は、複数の物理ホストおよび仮想マシンの各々を複数のパーティションの少なくとも1つに関連付けることができ、複数のパーティションの各々はP__Key値に関連付けられる。この方法は、1つ以上のスイッチポートの各々を、複数のビットマップに基づくP__Keyテーブルのうちのあるビットマップに基づくP__Keyテーブルに関連付けることができる。この方法は、ホストチャネルアダプタポートの各々を、複数のビットマップに基づくP__Keyテーブルのうちのあるビットマップに基づくP__Keyテーブルに関連付けることができる。

10

【0005】

本明細書では、高性能コンピューティング環境においてビットマップに基づくハードウェア実現例を使用してレガシーP__Keyテーブル抽象化をサポートするためのシステムおよび方法を説明する。ある例示的方法は、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少なくとも1つのサブネットを提供することができ、少なくとも1つのサブネットは、1つ以上のスイッチを含み、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、複数のスイッチポートを含み、少なくとも1つのサブネットはさらに、複数のホストチャネルアダプタを含み、各ホストチャネルアダプタは、少なくとも1つのホストチャネルアダプタポートを含み、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、少なくとも1つのサブネットはさらに、複数のエンドノードを含み、複数のエンドノードの各々は、複数のホストチャネルアダプタのうちの少なくとも1つのホストチャネルアダプタに関連付けられる。この方法は、エンドノードの各々を複数のパーティションの少なくとも1つに関連付けることができ、複数のパーティションの各々はP__Key値に関連付けられる。この方法は、1つ以上のスイッチポートの各々を、複数のビットマップに基づくP__Keyテーブルのうちのあるビットマップに基づくP__Keyテーブルに関連付けることができる。この方法は、ホストチャネルアダプタポートの各々を、複数のビットマップに基づくP__Keyテーブルのうちのあるビットマップに基づくP__Keyテーブルに関連付けることができる。この方法は、複数のビットマップに基づくP__Keyテーブルの各々を、仮想P__Key

20

30

【0006】

一実施形態によれば、複数のホストチャネルアダプタの1つ以上は、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ことができる。複数のエンドノードは、物理ホスト、仮想マシン、または物理ホストと仮想マシンとの組み合わせを含むことができ、仮想マシンは、少なくとも1つの仮想機能に関連付けられる。

【図面の簡単な説明】

【0007】

【図1】一実施形態によるインフィニバンド環境の一例を示す図である。

40

【図2】一実施形態による、パーティショニングされたクラスタ環境を示す図である。

【図3】一実施形態による、ネットワーク環境におけるツリートポロジを示す図である。

【図4】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5】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である。

【図6】一実施形態に従った例示的なvPortアーキテクチャを示す図である。

【図7】一実施形態に従った、LIDが予めボピュレートされた例示的なvSwitchアーキテクチャを示す図である。

【図8】一実施形態に従った、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】一実施形態に従った、動的LID割当てがなされかつLIDが予めボピュレート

50

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 1 0】一実施形態による例示的なマルチサブネットインフィニバンドファブリックを示す。

【図 1 1】一実施形態による、パーティション分離のための連想テーブルを有する例示的なインフィニバンドファブリックを示す。

【図 1 2】一実施形態による、パーティション分離のためのビットマップを有する例示的なインフィニバンドファブリックを示す。

【図 1 3】一実施形態による、高性能コンピューティング環境においてスケーラブルなビットマップに基づく `P _ K e y` テーブルをサポートするための方法のフローチャートである。

10

【図 1 4】一実施形態による、パーティション分離のためのビットマップと、レガシー管理エンティティのための `P _ K e y` 抽象化とを有する例示的なファブリックを示す。

【図 1 5】一実施形態による、パーティション分離のためのビットマップと、レガシー管理エンティティのための `P _ K e y` 抽象化とを有する例示的なファブリックを示す。

【図 1 6】一実施形態による、高性能コンピューティング環境においてスケーラブルなビットマップに基づく `P _ K e y` テーブルをサポートするための方法のフローチャートである。

【図 1 7】一実施形態による、高性能コンピューティング環境においてビットマップに基づくハードウェア実現例を使用してレガシー `P _ K e y` テーブル抽象化をサポートするための方法のフローチャートである。

20

【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。

【0010】

本明細書では、高性能コンピューティング環境においてビットマップに基づくハードウェア実現例を使用して構成可能なレガシー `P _ K e y` テーブル抽象化をサポートするシステムおよび方法を説明する。

40

【0011】

本発明の以下の説明は、高性能ネットワークのための例としてインフィニバンドTM (`IB`) ネットワークを使用する。以下の説明を通して、インフィニバンドTM 規格 (様々な、インフィニバンド規格、`IB` 規格、またはレガシー `IB` 規格とも呼ばれる) を参照することができる。このような参照は、引用によりその全体が本明細書に援用される、<http://www.inifinibandta.org> で入手可能なインフィニバンド (登録商標) トレード・アソシエーション・アーキテクチャ規格、第 1 巻、バージョン 1 . 3 (2 0 1 5 年 3 月リリース) を参照すると理解される。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

50

のファブリックトポロジーマトリクス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。

【0012】

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マトリクス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

【0013】

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マトリクス（たとえばファットツリートポロジーマトリクス）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。

20

【0014】

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。

【0015】

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（InfiniBandTM 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

【0016】

インフィニバンド・アーキテクチャ（InfiniBand Architecture：IBA）は、2層トポロジーマトリクス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

【0017】

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。

50

【 0 0 1 8 】

サブネット管理インターフェイスを通して、サブネットマネージャは、サブネット管理パケット (subnet management packet : S M P) と呼ばれる制御パケットを、サブネット管理エージェント (subnet management agent : S M A) と交換する。サブネット管理エージェントは、すべての I B サブネットデバイス上に存在する。S M P を使用することにより、サブネットマネージャは、ファブリックを発見し、エンドノードおよびスイッチを構成し、S M A から通知を受信することができる。

【 0 0 1 9 】

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

【 0 0 2 0 】

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

【 0 0 2 1 】

さらに、サブネット内のホストおよびスイッチは、ローカル識別子 (L I D) を用いてアドレス指定され得るとともに、単一のサブネットは 4 9 1 5 1 個のユニキャスト L I D に制限され得る。サブネット内で有効なローカルアドレスである L I D の他に、各 I B デバイスは、6 4 ビットのグローバル一意識別子 (global unique identifier : G U I D) を有し得る。G U I D は、I B レイヤー 3 (L 3) アドレスであるグローバル識別子 (global identifier : G I D) を形成するために使用され得る。

【 0 0 2 2 】

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

【 0 0 2 3 】

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

【 0 0 2 4 】

一実施形態に従ったインフィニバンド環境 1 0 0 の例を示す図 1 に、インフィニバンドファブリックの一例を示す。図 1 に示す例では、ノード A 1 0 1 ~ E 1 0 5 は、インフィ

10

20

30

40

50

ニバンドファブリック 120 を使用して、それぞれのホストチャネルアダプタ 111 ~ 115 を介して通信する。一実施形態に従うと、さまざまなノード（たとえばノード A 101 ~ E 105）はさまざまな物理デバイスによって表わすことができる。一実施形態に従うと、さまざまなノード（たとえばノード A 101 ~ E 105）は仮想マシンなどのさまざまな仮想デバイスによって表わすことができる。

【0025】

インフィニバンドにおけるパーティショニング

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。

【0026】

一実施形態によれば、パーティションは、あるグループのメンバが同じ論理グループの他のメンバとしか通信できないような、ポートの論理グループである。ホストチャネルアダプタ（HCA）およびスイッチでは、パーティションメンバーシップ情報を使用してパケットをフィルタリングして分離を実行できる。無効なパーティショニング情報を持つパケットは、そのパケットが着信ポートに到着するとすぐ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 システムでは、パーティションを使用してテナントクラスタを作成できる。パーティションの適所における実施で、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はできない。このようにして、侵害されたテナントノードまたは悪意のあるテナントノードが存在する場合でも、システムのセキュリティを保証することができる。

【0027】

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直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。

【0028】

IB パーティションのある例が、図 2 に示されており、それは、一実施形態による、パーティショニングされたクラスタ環境を示している。図 2 に示す例では、ノード A ~ E 101 ~ 105 は、インフィニバンドファブリック 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、パーティション、すなわちパーティション 1, 130、パーティション 2, 140、およびパーティション 3, 150 に配置される。パーティション 1 は、ノード A 101 およびノード D 104 を含む。パーティション 2 は、ノード A 101、ノード B 102、およびノード C 103 を含む。パーティション 3 は、ノード C 103 およびノード E 105 を含む。パーティションの配置のため、ノード D 104 とノード E 105 とは、これらのノードがパーティションを共有しないので、通信することができない。一方、例えば、ノード A 101 とノード C

10

20

30

40

50

103とは、これらのノードが両方ともパーティション2, 140のメンバであるため、通信が許可される。

【0029】

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、必要な性能を提供することができる。

10

【0030】

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。

20

【0031】

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。

30

【0032】

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

40

50

【 0 0 3 3 】

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、I B ベースの H P C 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。

【 0 0 3 4 】

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 2 0 0 において、1 つ以上のエンドノード 2 0 1 ~ 2 0 4 が接続され得る。ネットワークファブリック 2 0 0 は、複数のリーフスイッチ 2 1 1 ~ 2 1 4 と複数のスパインスイッチまたはルート (root) スwitch 2 3 1 ~ 2 3 4 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 2 0 0 は、スイッチ 2 2 1 ~ 2 2 4 などの 1 つ以上の中間スイッチを含み得る。

【 0 0 3 5 】

また、図 3 に示すように、エンドノード 2 0 1 ~ 2 0 4 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 2 0 0 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 2 0 1 はポート H 1 および H 2 を含み、ノード 2 0 2 はポート H 3 および H 4 を含み、ノード 2 0 3 はポート H 5 および H 6 を含み、ノード 2 0 4 はポート H 7 および H 8 を含み得る。

【 0 0 3 6 】

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 7 】

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 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。

【 0 0 3 8 】

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。

【 0 0 3 9 】

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ー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

10

20

30

40

50

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。

【0040】

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リングすることができる階層ネットワークトポロジである。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 -ary- n -tree、拡張された一般化ファットツリー (Extended Generalized Fat-Tree: XGFT)、パラレルポート一般化ファットツリー (Parallel Ports Generalized Fat-Tree: PGFT) およびリアルライフファットツリー (Real Life Fat-Tree: RLFT) 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。

10

【0041】

また、 k -ary- n -tree 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。XGFT 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 -ary- n -tree を拡張させる。PGFT 定義はさらに、XGFT トポロジを拡張して、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は XGFT および PGFT 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 HPC 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、PGFT の制限バージョンである RLFT が導入されている。RLFT 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。

20

【0042】

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。

30

【0043】

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。

【0044】

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。

40

【0045】

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

50

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。

【0046】

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。

10

【0047】

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。

20

30

【0048】

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。

40

【0049】

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

50

ることとなる。

【0050】

インフィニバンドSR - IOVアーキテクチャ - 共有ポート

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

【0051】

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。

10

【0052】

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。

【0053】

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管理パケットのために用いられる専用のキュー対）が物理機能によって所有される。これらのQPはVFにも同様にエクスポートされるが、VFはQP0を使用することが許可されておらず（VFからQP0に向かって入来するすべてのSMPが廃棄され）、QP1は、PFが所有する実際のQP1のプロキシとして機能し得る。

20

【0054】

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。

【0055】

30

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。

【0056】

インフィニバンドSR - IOVアーキテクチャモデル - 仮想スイッチ（vSwitch）

40

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。

【0057】

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 Adapt

50

er: 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。

【0058】

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。

【0059】

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6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ができる。

【0060】

一実施形態に従うと、ベンダーに実装の自由を与えるために v P o r t 概念は緩やかに定義されており（たとえば、当該定義では、実装が S R I O V 専用とすべきであるとは規定されていない）、v P o r t の目的は、V M 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。v P o r t 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ーラブルであり得る、S R - I O V 共有のポートのようなアーキテクチャおよび v S w i t c h 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、v P o r t はオプションの L I D をサポートするとともに、共有のポートとは異なり、S M は、v P o r t が専用の L I D 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 v P o r t を認識する。

【0061】

インフィニバンド S R - I O V アーキテクチャモデル - L I D が予めボピュレートされた v S w i t c h

一実施形態に従うと、本開示は、L I D が予めボピュレートされた 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。

【0062】

図7は、一実施形態に従った、L I D が予めボ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 6 0 0（たとえば I B サブネット）内においてインフィニバンドファブリックなどのファブリックのメンバ間で通信を確立することができる。ファブリックはホストチャネルアダプタ 5 1 0、5 2 0、5 3 0 などのいくつかのハードウェアデバイスを含み得る。さらに、ホストチャネルアダプタ 5 1 0、5 2 0 および 5 3 0 は、それぞれ、ハイパーバイザ 5 1 1、5 2 1 および 5 3 1 と対話することができる。各々のハイパーバイザは、さらに、ホストチャネルアダプタと共に、いくつかの仮想機能 5 1 4、5 1 5、5 1 6、5 2 4、5 2 5、5 2 6、5 3 4、5 3 5 および 5 3 6 と対話し、設定し、

10

20

30

40

50

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 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 に割り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 に割り当て、仮想マシン 3 5 5 2 を仮想機能 3 5 1 6 に割り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 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。

【0063】

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 6 0 0 内においてトラフィックを方向付けるためにリニアフォワーディングテーブルを設定するのに用いられる。

10

【0064】

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

20

【0065】

一実施形態に従うと、本開示は、L I D が予めボピュレートされた 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、予めボ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 が予めボピュレートされており、仮想機能 1 5 3 4 は L I D 1 0 が予めボ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ボ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、ボピュレートされた V F には、図 7 に示されるように L I D が割り当てられている。

30

【0066】

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。

【0067】

一実施形態に従うと、L I D が予めボピュレートされた 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 1 6 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 1 7 個の L I D（1 6 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、（4 9 1 5 1 個の利用可能な L I D をハイパーバイザごとに 1 7 個の L I D で割って得られる）2 8 9 1 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 2 8 9 1 個のハイパーバイザに 1 6 の V F を掛けて得られる）4 6 2 5 6 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

40

50

たは専用の S M ノードが同様に L I D を消費するので、これらの数は実際にはより小さくなる)。なお、v S w i t c h が、L I D を P F と共有することができるので、付加的な L I D を占有する必要がないことに留意されたい。

【 0 0 6 8 】

一実施形態に従うと、L I D が予めボピュレートされた 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。

10

【 0 0 6 9 】

一実施形態に従うと、L I D が予めボピュレートされた 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 - M a s k - C o n t r o l - l i k e : 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。

20

【 0 0 7 0 】

一実施形態に従うと、L I D が予めボピュレートされた 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 が予めボ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 が予めボ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。

【 0 0 7 1 】

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

30

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

【 0 0 7 2 】

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。

40

【 0 0 7 3 】

50

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 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。

【0074】

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 *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、*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。

【0075】

一実施形態に従うと、本開示は、動的 *LID* 割当てがなされた *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 8 を参照すると、*LID* 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 *LID* 1 を受取り、物理機能 523 が *LID* 2 を受取り、物理機能 533 が *LID* 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 *LID* 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 *LID* 5 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 *LID* が割当てられ、*LID* 7 が仮想機能 2 515 に割当てられ、*LID* 11 が仮想機能 3 516 に割当てられ、*LID* 9 が仮想機能 1 534 に割当てられている。*LID* が予めポピュレートされた *vSwitch* 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 *LID* の割当てを受けない。

【0076】

一実施形態に従うと、動的 *LID* 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、*VM* が存在していない場合、比較的少数の *LID* を最初の経路計算および *LFT* 分配のために用いることができる。

【0077】

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る（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。

【0078】

一実施形態に従うと、動的 *LID* 割当てがなされた *vSwitch* を利用するシステムにおいて新しい *VM* 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 *VM* をブートすべきであるかを決定するために、自由な *VM* スロットが発見され、固有の未使用のユニキャスト *LID* も同様に発見される。しかしながら、新しく追加された *LID* を処理するためのスイッチの *LFT* 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 *VM* 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 *VM* 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 *IB* 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 *VM* 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。

【0079】

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 *VF* が *PF* 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 *LFT* を繰返し、（*VM* が作成されている）ハイパーバイザの *PF* に属する *LID* エントリから新しく追加された *LID* にフォワーディング

10

20

30

40

50

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 8 0 】

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 が予めボピュレートされた 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、そこに予めボピュレートされた 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 が予めボピュレートされた 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 8 1 】

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、予めボピュレートされた L I D アーキテクチャモデルを備えた v S w i t c h の欠点を解決することができる。V M が作成されるたびに、作成された V M に関連付けられた、新しく追加された L I D で、サブネットにおける物理スイッチの L F T が更新される。この動作のために、1 スイッチごとに 1 つのサブネット管理パケット (S M P) が送信される必要がある。各々の V M がそのホストハイパーバイザと同じ経路を用いているので、L M C のような機能も利用できなくなる。しかしながら、すべてのハイパーバイザに存在する V F の合計に対する制限はなく、V F の数は、ユニキャスト L I D の限度を上回る可能性もある。このような場合、当然、アクティブな V M 上で V F のすべてが必ずしも同時に付与されることが可能になるわけではなく、より多くの予備のハイパーバイザおよび V F を備えることにより、ユニキャスト L I D 限度付近で動作する際に、断片化されたネットワークの障害を回復および最適化させるための融通性が追加される。

【 0 0 8 2 】

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

図 9 は、一実施形態に従った、動的 L I D 割当てがなされて L I D が予めボピュレートされた v S w i t c h を備えた例示的な v S w i t c h アーキテクチャを示す。図に示されるように、いくつかのスイッチ 5 0 1 ~ 5 0 4 は、ネットワーク切替環境 8 0 0 (たとえば I B サブネット) 内においてインフィニバンドファブリックなどのファブリックのメンバ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 5 1 0、5 2 0、5 3 0 などのいくつかのハードウェアデバイスを含み得る。ホストチャネルアダプタ 5 1 0、5 2 0 および 5 3 0 は、それぞれ、さらに、ハイパーバイザ 5 1 1、5 2 1 および 5 3 1 と対話することができる。各々のハイパーバイザは、さらに、ホストチャネルアダプタと共に、いくつかの仮想機能 5 1 4、5 1 5、5 1 6、5 2 4、5 2 5、5 2 6、5 3 4、5 3 5 および 5 3 6 と対話し、設定し、いくつかの仮想マシンに割当てることができる。たとえば、仮想マシン 1 5 5 0 は、ハイパーバイザ 5 1 1 によって仮想機能 1 5 1 4 に割当てることができる。ハイパーバイザ 5 1 1 は、加えて、仮想マシン 2 5 5 1 を仮想機能 2 5 1 5 に割当てることができる。ハイパーバイザ 5 2 1 は、仮想マシン 3 5 5 2 を仮想機能 3 5 2 6 に割当てることができる。ハイパーバイザ 5 3 1 は、さらに、仮想マシン 4 5 5 3 を仮想機能 2 5 3 5 に割当てることができる。ハイパーバイザは、ホストチャネルアダプタの各々の上において十分な機能を有する物理機能 5 1 3、5 2 3 および 5 3 3 を介してホストチャネルアダプタにアクセスすることができる。

【 0 0 8 3 】

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

【 0 0 8 4 】

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

【 0 0 8 5 】

一実施形態に従うと、本開示は、動的 L I D 割当てがなされ L I D が予めボピュレートされたハイブリッド v S w i t c h アーキテクチャを提供するためのシステムおよび方法を提供する。図 9 を参照すると、ハイパーバイザ 5 1 1 には、予めボピュレートされた L I D アーキテクチャを備えた v S w i t c h が配置され得るとともに、ハイパーバイザ 5 2 1 には、L I D が予めボピュレートされて動的 L I D 割当てがなされた v S w i t c h が配置され得る。ハイパーバイザ 5 3 1 には、動的 L I D 割当てがなされた v S w i t c h が配置され得る。このため、物理機能 5 1 3 および仮想機能 5 1 4 ~ 5 1 6 には、それらの L I D が予めボ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 が予めボ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。

【 0 0 8 6 】

L I D が予めボピュレートされた 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 9 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ボピュレートされた 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） $0 < =$ 予めボピュレートされた 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 の総数から予めボピュレートされた V F の数を減じることによって見出すことができる。

【 0 0 8 7 】

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

【 0 0 8 8 】

インフィニバンド - サブネット間通信（ファブリックマネージャ）

一実施形態によれば、本開示の実施形態は、単一サブネット内にインフィニバンドファブリックを提供することに加えて、2 つ以上のサブネットにまたがるインフィニバンドフ

10

20

30

40

50

アプリケーションを提供することもできる。

【0089】

図10は、一実施形態による例示的なマルチサブネットインフィニバンドファブリックを示す。図に示すように、サブネットA 1000内では、ある数のスイッチ1001～1004が、サブネットA 1000（例えば、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において通信を提供することができる。ファブリックは、例えば、チャンネルアダプタ1010などのある数のハードウェアデバイスを含むことができる。ホストチャンネルアダプタ1010は、次いで、ハイパーバイザ1011と対話することができる。ハイパーバイザは、次いで、それが対話するホストチャンネルアダプタと関連して、ある数の仮想機能1014をセットアップすることができる。ハイパーバイザは、加えて、仮想マシンを仮想機能の各々に割り当てることができ、仮想マシン1 10105が仮想機能1 1014に割り当てられるなどする。ハイパーバイザは、各ホストチャンネルアダプタ上において、物理機能1013など、十分な機能を有する物理機能を介して、それらの関連付けられるホストチャンネルアダプタにアクセスできる。サブネットB 1040内では、ある数のスイッチ1021～1024が、サブネットB 1040（例えば、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において通信を提供することができる。ファブリックは、例えば、チャンネルアダプタ1030などのある数のハードウェアデバイスを含むことができる。ホストチャンネルアダプタ1030は、次いで、ハイパーバイザ1031と対話することができる。ハイパーバイザは、次いで、それが対話するホストチャンネルアダプタと関連して、ある数の仮想機能1034をセットアップすることができる。ハイパーバイザは、加えて、仮想マシンを仮想機能の各々に割り当てることができ、仮想マシン2 1035が仮想機能2 1034に割り当てられるなどする。ハイパーバイザは、各ホストチャンネルアダプタ上において、物理機能1033など、十分な機能を有する物理機能を介して、それらの関連付けられるホストチャンネルアダプタにアクセスできる。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）内には1つのホストチャンネルアダプタしか示されていないが、複数のホストチャンネルアダプタおよびそれらの対応するコンポーネントを各サブネット内に含めることができることを理解されたい。

【0090】

一実施形態によれば、ホストチャンネルアダプタの各々は、仮想スイッチ1012および仮想スイッチ1032などの仮想スイッチにさらに関連付けることができ、各HCAは、上述したように、異なるアーキテクチャモデルでセットアップすることができる。図10内の両方のサブネットは、事前にボビュレートされたLIDアーキテクチャモデルを有するvswitchを使用しているように示されているが、これはすべてのそのようなサブネット構成が同様のアーキテクチャモデルに従わなければならないことを意味するものではない。

【0091】

一実施形態によれば、各サブネット内の少なくとも1つのスイッチはルータに関連付けられることができ、サブネットA 1000内のスイッチ1002はルータ1005に関連付けられ、サブネットB 1040内のスイッチ1021はルータ1006に関連付けられるなどする。

【0092】

一実施形態によれば、少なくとも1つのデバイス（例えば、スイッチ、ノードなど）をファブリックマネージャ（図示せず）に関連付けることができる。ファブリックマネージャは、例えば、サブネット間ファブリックトポロジを発見し、ファブリックプロファイル（例えば、仮想マシンファブリックプロファイル）を作成し、仮想マシンファブリックプロファイルを構築するための基礎を形成する仮想マシン関連データベースオブジェクトを構築するために使用することができる。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パーティション番号を使用してどのルータポートを介して通信することが許可されているかに関して、法的なサブネット間接続を定義することができる。

【 0 0 9 3 】

一実施形態によれば、サブネット A 内の仮想マシン 1 などの発信元でのトラフィックが、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の宛先にアドレス指定されている場合、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 にアドレス指定され、ルータ 1 0 0 5 は次いでそのトラフィックをルータ 1 0 0 6 とのそのリンクを介してサブネット B に渡すことができる。

【 0 0 9 4 】

スケーラブルなビットマップに基づく P __ K e y テーブル

一実施形態によれば、仮想マシンを使用を提供された拡張されたインフィニバンドネットワークのため、可能性のあるパーティションの数が劇的に増加している。しかしながら、P __ K e y 管理およびルーティングを扱う現在の方法では、トラフィックがファブリックを通過する際にかなりのオーバーヘッド時間が追加される。従来、インフィニバンド規格では、S M が索引付けされたテーブルとしてアクセスできる 1 6 ビットの P __ K e y 値の配列として、P __ K e y テーブルを定義している。この規格のハードウェア実現例は、ワイヤ速度パケットレートで I B パケットのパーティションチェックを実行するためにルックアップを実行するために連想メモリを使用することを意味する。実際には、これはハードウェアにより実現される P __ K e y テーブルの可能なサイズを、1 6 ビットの P __ K e y 値が表す 6 4 K 値空間よりも桁違いに小さく制限する。

【 0 0 9 5 】

一実施形態によれば、インフィニバンド規格は、パーティションメンバシップを 1 6 ビット P __ K e y を介して定義し、ファブリック内の各ポートで P __ K e y 値の固定テーブルを使用してパーティション分離を実施する。サブネットマネージャは、異なるテーブルエントリで異なる P __ K e y 値をプログラミングできる。そのパーティション関連付けに基づいて P __ K e y でそのヘッダにおいてマークされたパケットがポートに到着すると、下位のハードウェアは、入来するパケットを、チェックするハードウェアに関連付けられる P __ K e y テーブル内のすべての値と比較することができる（すなわち、テーブルで連想ルックアップを実行する）。しかしながら、このルックアップは多数のパーティションではうまくスケールせず、テーブル内である数のパーティション上でルックアップを実行している間に不要なオーバーヘッドが発生する可能性がある。

【 0 0 9 6 】

図 1 1 は、一実施形態による、パーティション分離のための連想テーブルを有する例示的なインフィニバンドファブリックを示す。図 1 1 に示すように、1 つ以上のエンドノード 1 1 4 1 ~ 1 1 4 4 を、ネットワークファブリック 1 1 0 0 において接続することができる。ネットワークファブリック 1 1 0 0 は、複数のリーフスイッチ 1 1 1 1 ~ 1 1 1 4 と、複数のスパインスイッチまたはルートスイッチ 1 1 3 1 ~ 1 1 3 4 とを含む、ファットツリートポロジに基づくことができる。加えて、ネットワークファブリック 1 1 0 0 は、スイッチ 1 1 2 1 ~ 1 1 2 4 のような 1 つ以上の中間スイッチを含むことができる。また、図 1 1 に示すように、エンドノード 1 1 4 1 ~ 1 1 4 4 の各々はマルチホームノードであってもよい。

【 0 0 9 7 】

一実施形態によれば、ノード 1 1 4 1 ~ 1 1 4 4 の各々は、ある数のパーティションのうちの 1 つ以上に属することができ、各パーティションは、ある P __ K e y 値に関連付けられる。上述したように、P __ K e y 値は、ファブリック内で送信される各パケットに含めることができ、サブネット内におけるパーティション実施の程度を変化させることができる。P __ K e y 値の各々は、ある数のテーブル、例えばテーブル 1 1 0 1 ~ 1 1 0 8（ここでは「ポート P __ K e y テーブル」とも呼ばれる）に格納され、各テーブルは、例えば、スイッチのポート、またはホストノードに属する H C A ポートなど、ファブリックの他のメンバのポートに関連付けられる。各テーブル 1 1 0 1 ~ 1 1 0 8 は、サブネットマネージャによって設定することができる。

【 0 0 9 8 】

一実施形態によれば、各ポート P__K e y テーブルは、各エントリが 16 ビット P__K e y 値を含むことができる N 個のエントリの配列を含むことができ、N の値は、（ポートのステータスおよび能力を発見する一環として S M によって検索される）ポートの能力である。ポートがパーティションのメンバであるはずである（つまり、フルまたは制限付き）ことを意味するポリシー情報を S M がポートが持っている場合、S M はポート P__K e y テーブルにおいて次に利用可能なエントリを見つけて、関連する 16 ビット P__K e y 値をそのエントリに格納できる。P__K e y 値が成功裡に格納される（S M A によって S M に確認される）と、S M は、ポートハードウェアがポートを通過するすべてのパケットについて新たな P__K e y テーブルエントリの一致を含むことを予期することができ、パケット P__K e y 値が更新された P__K e y テーブルエントリにおける新たな値と一致するパケットを受け入れる（すなわち、ポートを通過させる）ことができる。

10

【0099】

一実施形態によれば、サブネット内のポートがもはやパーティションのメンバではないと想定される場合、S M は関連する P__K e y テーブルエントリをクリアする（0 を書き込む）ことができる。P__K e y 値が成功裡にクリアされる（S M A によって S M に確認される）と、S M は、ポートハードウェアが、丁度ポート P__K e y テーブルからクリアされたものに対応するパケット P__K e y 値を有する、そのポートを通過するすべてのパケットについて P__K e y 一致に失敗することを予期し得、次いでパケットをドロップし得る。

20

【0100】

一実施形態によれば、S M は、既存の値を最初にクリアすることなく、有効な P__K e y テーブルエントリを新たな値と置き換えることができる。この場合の期待される効果は、エントリの内容が最初にクリアされ、次いでその後新たな値で更新された場合と同じである。

【0101】

すでに説明したように、図 11 に示す実施形態は、多数のパーティションを持たないサブネットの範囲内で十分にスケーリングすることができる。しかしながら、例えば、エンドノードおよび 1 つ以上の v S w i t c h アーキテクチャの仮想化を利用するサブネットなど、任意の所与のサブネット内に多数のパーティションが存在すると、上述の連想テーブルルックアップはオーバーヘッド時間の大幅な増加を招き得、I B 規格によって規定された関連するリンク速度要件が満たされるべきである場合、ハードウェアで実現することが可能でさえないかもしれない。

30

【0102】

一実施形態によれば、I B ファブリック内の各ポートに関連付けられる P__K e y テーブルを利用する代わりに、例示的な方法およびシステムは、ビットマップを利用して、可能なすべての P__K e y 値の表現、16 ビット（すなわち 64 k の可能な値）をハードウェアにおいて実現できる。そのような方法およびシステムでは、ハードウェア実現例内において、各可能な P__K e y 値を単一ビットとして表現することができ、単一ビットの値（例えば、1 または 0）は、P__K e y 値が許可される（すなわち、パケットは、パーティションチェックを実行するポートを通過することを許可される）か、またはポートに関連付けられない（すなわち、I B パケットは、次いで、I B 規格によって規定された関連するポートタイプについての P__K e y 実施規則に従ってそのパーティション番号に対する制限付きおよびフルメンバエントリの両方についてのビットマップエントリとパケット P__K e y を関連させることに基づいて、パーティション実施規則に従ってチェックされ得る）かを定義できる。

40

【0103】

図 12 は、一実施形態による、パーティション分離のためのビットマップを有する例示的なインフィニバンドファブリックを示す。図 12 に示すように、1 つ以上のエンドノード 1241 ~ 1244 を、ネットワークファブリック 1200 において接続することができる。ネットワークファブリック 1200 は、複数のリーフスイッチ 1211 ~ 1214

50

と、複数のスパインスイッチまたはルートスイッチ 1 2 3 1 ~ 1 2 3 4 とを含む、ファットツリートポロジに基づくことができる。加えて、ネットワークファブリック 1 2 0 0 は、スイッチ 1 2 2 1 ~ 1 2 2 4 のような 1 つ以上の中間スイッチを含むことができる。また、図 1 2 に示すように、エンドノード 1 2 4 1 ~ 1 2 4 4 の各々はマルチホームノードであってもよい。

【 0 1 0 4 】

一実施形態によれば、ノード 1 2 4 1 ~ 1 2 4 4 の各々は、ある数のパーティションのうちの 1 つ以上に属することができ、各パーティションは、ある P __ K e y 値に関連付けられる。上述したように、P __ K e y 値は、ファブリック内で送信される各パケットに含めることができ、サブネット内におけるパーティション実施の程度を変化させることができる。各 P __ K e y 値は、ビットマップ 1 2 0 1 ~ 1 2 0 8 のようなある数のビットマップで参照することができ、各ビットマップは、例えば、スイッチのポート、またはホストノードに属する（関連付けられる）H C A ポートなど、ファブリックの他のメンバのポートに関連付けられる。各ビットマップ 1 2 0 1 ~ 1 2 0 8 は、サブネットマネージャによって設定することができる。

10

【 0 1 0 5 】

一実施形態によれば、P __ K e y を連想テーブルルックアップの代わりにビットマップで表すことにより、これにより、柔軟性が増大し、ファブリックにおけるパーティション分離およびパケットルーティングのオーバーヘッド時間が短縮できる。どのノードも任意のパーティションのメンバになることができる。これにより、任意のスイッチポートが、S M が構成した P __ K e y だけをそのスイッチポートを介して許可されるようにすることができる（そうでない場合、パケットはドロップされ得る）。ビットマップに基づく P __ K e y テーブルを利用することにより、方法およびシステムはほんのわずかな数のパーティションまたは数百のものを許容することができる。そのような方法およびシステムは、スイッチを通過することになっているパーティションのみを許可することを確実にすることができる。このような方法およびシステムは、より多くの数のパーティションが利用可能であっても、より大きなパーティション分離を可能にすることができる。

20

【 0 1 0 6 】

一実施形態によれば、ビットマップスキームを使用することにより、P __ K e y 値の連想ルックアップを有する代わりに、ハードウェアはパケットヘッダから P __ K e y 値を取ることができ、連想ルックアップ機構に値を送る代わりに、P __ K e y 値を取り、それを（例えば、6 4 K ビットの）ビット配列に索引付する。次いで、I B パケットが、I B 規格によって規定された関連するポートタイプについての P __ K e y 実施規則に従ってそのパーティション番号に対する制限付きおよびフルメンバエントリの両方についてのビットマップエントリとパケット P __ K e y を相関させることに基づいて、パーティション実施規則に従ってチェックされ得る。これは、例えば、スケーリングの増加（すなわち、パフォーマンスの損失なしにより多数のパーティション P __ K e y を可能にすること）に至り得る。

30

【 0 1 0 7 】

一実施形態によれば、上述のビットマップスキームは、エンドノードでパーティション分離を実施することに加えて、エンドノードに向けて接続されるポートに対する第 1 のリーフスイッチ、第 1 のリーフスイッチとルートスイッチとの間の任意のスイッチ、および最後のリーフスイッチなど、中間スイッチでパーティション実施を追加的にサポートすることができる。スイッチからスイッチへのリンクが、ファブリックのルーティングに応じて、多くの異なるノードによって使用される。従来のシステムでは、パーティション実施は主にエンドノード上で使用され、エンドノードに直接接続されるポートに対する最後のリーフスイッチで実施されたかも知れないが、リーフスイッチよりも高いレベルでのスイッチでは実施されない（つまり、パーティション実施はスイッチ間のリンク上では使用されなかった）。しかしながら、許可された P __ K e y および許可されていない P __ K e y に関連付けられるパーティションビットマップを使用することで得られる速度および効率

40

50

により、パーティション実施がこれらの中間リンクで起こり得る。各スイッチポートは許可されているパーティションがそのポートを介して送信されることを許すにすぎないことを確実にすることによって、中間リンクで通過できないポートに遭遇するパケットは、（エンドノードまたはリーフスイッチにおける代わりに）はるかによりすぐにドロップされ得、最終的にドロップされるパケットに対するリンク帯域幅使用率に関してパフォーマンスの低下を少なくすることができる。

【0108】

一実施形態によれば、SMは、パーティションに応じて、完全なファブリックのどこでどのトラフィックがルーティングされることが許可されるかを、ビットマップで判断して設定することができる。SMは、次いで、ビットマップを使用してパーティション実施をセットアップできる。SMは、どのエンドノードがどのパーティションのメンバであるかに関するポリシー構成を受信することができる。このポリシー情報に基づいて、（ビットマップを介して）SMはエンドノード上でパーティション実施を設定できる。SMは、接続を切り替えるようにスイッチを構成することもできる。これを行った後、SMはビットマップを使用してファブリック内のトラフィックフローを制限できる。SMは、ビットマップ内のパーティション毎の値を1または0、パスまたはドロップ（すなわち、許可する、または許可しない）にセットする。これにより、許可されていないトラフィックを、宛先ノードまたは宛先ノードの前の最後のリーフスイッチにおいてではなく、トラフィックのルーティングにおいて早期にドロップできる。

【0109】

一実施形態によれば、各ビットマップは、各ポートに対するP__Key実施のための64Kビット配列を含むビット配列とすることができる。このようなビットマップは、任意の種類のIBポート：スイッチ外部ポート、スイッチ管理ポート（すなわち、スイッチポート0）、TCAポート（例えば、IBイーサネット（登録商標）ブリッジ/ゲートウェイノード実装用）、HCAポートおよびルータポートに関連付けられ得る。次いで、IBパケットを、IB規格によって規定された関連するポートタイプについてのP__Key実施規則に従ってそのパーティション番号に対する制限付きおよびフルメンバエントリの両方についてのビットマップエントリとパケットP__Keyを相関させることに基づいて、パーティション実施規則に従ってチェックし得る。

【0110】

一実施形態によれば、単一のパーティション番号を2つのビット配列索引に関連付けることができ、1つは制限付きメンバシップを表すものであり、1つはフルメンバシップを表すものである。したがって、ポートがフルメンバシップのみを表すか、フルおよび制限付きメンバシップの両方を表すと考えられるかに応じて、フルメンバ索引のみまたはフルおよび制限付きメンバ索引の両方を設定して、対応するメンバシップトラフィックを許可する必要がある。しかしながら、制限付きパーティションメンバが別の制限付きメンバと通信することが許可されていない一方で、フルメンバがフルメンバおよび制限付きメンバの両方と通信することが許可されているという事実のため、ルックアップ実装はこれを考慮し、パケットを開始しているエンドノードが、制限付きメンバのパケットを、それがパーティションの制限付きメンバであるにすぎない場合に送信することが許可されるにすぎず、しかしながら、それは、対応して、フルメンバパケットを受信することが許可されるにすぎないことになるようにしなければならない。

【0111】

図13は、一実施形態による、高性能コンピューティング環境においてスケーラブルなビットマップに基づくP__Keyテーブルをサポートするための方法のフローチャートである。

【0112】

ステップ1310において、パケットは、送信元（例えば、送信元ノードまたは送信元仮想マシン）から発信され、宛先（例えば、宛先ノードまたは宛先仮想マシン）に宛てられ得る。パケットは、他の情報の中でもとりわけ、送信元と宛先との間の共通パーティシ

ョンなどのパーティション情報に対応する P__K e y 値を含むことができるヘッダを含むことができる。

【0113】

ステップ1320において、パケットは、算出されたルートに沿ってインフィニバンドファブリックのようなファブリックを横断することを開始できる。ルートは、例えば、サブネットマネージャによって算出することができる。

【0114】

ステップ1330において、パケットは、ファブリック内のルーティングされた経路に沿ってスイッチのアドレス指定されたポートに到達することができる。

【0115】

ステップ1340で、アドレス指定されたポートが位置するスイッチは、パケット P__K e y をビットマップの内容と照合することができる（すなわち、IBパケットは、IB規格によって規定された関連するポートタイプについての P__K e y 実施規則に従ってそのパーティション番号に対する制限付きおよびフルメンバエントリの両方についてのビットマップエントリとパケット P__K e y を関連させることに基づいて、パーティション実施規則に従ってチェックされ得る）。

【0116】

ステップ1350で、チェックの結果が「許可」である場合、これは、パケットがルーティングされた経路に沿って宛先へ（DLID）進むことを許可する。

【0117】

ステップ1360で、チェックの結果が「不許可」である場合、スイッチはアドレス指定されたポートでパケットをドロップし、パケットがルーティングされた経路に沿って進まないようにすることができる。

【0118】

図13に示され、上述された実施形態は、ファブリック内の1つのスイッチのみでのパケット P__K e y のチェックを開示するが、ルーティングされた経路に沿った各スイッチおよびノードは、各アドレス指定されたポートに関連付けられるビットマップを使用して P__K e y のチェックを実行して、パケットが許可されているかいないかを判断することができることが理解されるべきである。

【0119】

ビットマップに基づくHW実現例を使用する構成可能なレガシー P__K e y テーブル抽象化

一実施形態によれば、上記のビットマップに基づく実現例は、レガシーサブネットマネージャ実現例によって直接使用することができないという点で問題がある。これは、インフィニバンド規格（インフィニバンド（登録商標）トレード・アソシエーション・アーキテクチャ規格、第1巻、バージョン1.3（2015年3月リリース））で定義されているサブネットマネージャなどの、現在定義されている管理エンティティが、ビットマップに基づく P__K e y 実現例と対話するように定義されていないためであり、それらは、代わりに、IB規格によって定義されたレガシー P__K e y テーブルに基づくスキームと対話するように定義される。

【0120】

一実施形態によれば、DRAMにおけるマッピングテーブルを、マッピングテーブルを実装するソフトウェアベースのSMAの使用によって達成することができる。このマッピングテーブルを用いて、ビットマップに基づく P__K e y テーブルのレガシー準拠のビューを提供することが可能である。そのようなレガシー準拠のビューは、仮想 P__K e y テーブル、または構成可能なレガシー P__K e y テーブル抽象化と呼ぶことができる。

【0121】

一実施形態によれば、仮想 P__K e y テーブルのサイズは、その目的のために割り当てられ得るDRAMの量によって制限される。特定の実施形態では、アクセス可能なDRAMの量は、仮想 P__K e y テーブルに関連付けられるスイッチまたは他のデバイス（エン

10

20

30

40

50

ドノードなど)内に含まれる。ある実施形態では、DRAMの外部ソースを利用して、より大きな仮想P__Keyテーブルを管理エンティティによってアクセス可能にすることができる。

【0122】

一実施形態によれば、各仮想P__Keyテーブルは、構成パラメータ(またはヘッダファイル定数)で定義されたサイズを有することができ、関連するSMA属性インスタンスに関連付けられる対応するメモリデータ構造(例えば、埋め込みプロセッササブシステムメモリデータ構造(EPS))によって実現することができる。

【0123】

一実施形態によれば、管理エンティティ(例えば、サブネットマネージャ)が現在クリアされている(値ゼロを含む)(仮想の)P__KeyテーブルエントリにP__Key値を格納する度毎に、SMA実現例はHWアクセスを呼び出して、P__Key値によって索引付けされたHWビット配列を1にセットでき、次いで、仮想P__Keyテーブルエントリ(例えば、EPSメモリデータ構造)を更新して、SMによって指定されるP__Keyテーブルエントリ番号が指定されたP__Key値を含むようにすることができる。

【0124】

一実施形態によれば、SMが仮想P__KeyテーブルにおけるエントリからP__Key値をクリアする度毎に、SMA実現例はメモリデータ構造からの古いP__Key値を使用し、次いで、仮想P__Keyテーブルエントリ(すなわち、EPSメモリデータ構造)もクリアされる前に、対応するHWビット配列索引をHWアクセスインターフェイスを介してクリアすることができる。

【0125】

一実施形態によれば、SMが、現在クリアされていない(仮想)P__Keyテーブルエントリに新たなP__Key値を書き込む度毎に、SMA実現例はまずEPSメモリデータ構造からの古いP__Key値を使用し、次いで、対応するHWビット配列索引をHWアクセスインターフェイスを介してクリアすることができる。古い値がクリアされた後、SMA実現例は、新たに追加されたP__Key値に対応するビット配列エントリを設定し、次いで、メモリデータ構造を更新して、指定されたP__Keyテーブルエントリが指定されたP__Key値を含むようにすることができる。

【0126】

図14は、一実施形態による、パーティション分離のためのビットマップと、レガシー管理エンティティのためのP__Key抽象化とを有する例示的なファブリックを示す。図に示す実施形態では、スイッチ1400は、ある数のポート1420を備えることができ、これらのポートの各々は、上述したように、P__Key実施のためのビットマップに基づくP__Keyテーブル1410に関連付けることができる。スイッチ1400は、加えて、メモリ(例えばDRAM)1430を備えることができる。

【0127】

一実施形態によれば、サブネットマネージャ1440(例えばレガシーサブネットマネージャ)は、ビットマップに基づくP__Keyテーブル1410を照会/設定することができないことがあり得る。このような状況において、スイッチのメモリ1430内にP__Keyテーブル抽象化1450(すなわち「仮想P__Keyテーブル」)を設けることができる。次いで、サブネットマネージャは、サブネットマネージャが連想メモリ/ルックアップに基づいてレガシーP__Keyテーブルと対話するであろうのと同じ態様で、仮想P__Keyテーブルを照会し、それと対話することができる。より具体的には、SM1440は、スイッチ1400に関連付けられるSMA1441とインターフェイスすることができる。このインターフェイスは、P__Key管理のために、IB規格で定義されたSMA属性および方法を利用することができるため、SMは基礎となる実装についての知識がない。

【0128】

一実施形態によれば、仮想P__Keyテーブルに含まれる各P__Keyは16ビットを

10

20

30

40

50

含むため、そして仮想 P__K e y テーブルに含まれるべき可能性のある P__K e y 値の数（すなわち、最大 6 4 K P__K e y まで）のため、仮想 P__K e y テーブルのサイズは、スイッチに含まれる、および / またはスイッチによって仮想 P__K e y テーブルに割り当てられるメモリ 1 4 3 0 の量によって制限され得る。例えば、3 6 ポートスイッチの各ポートについて 6 4 K P__K e y を表現するには、必要なメモリサイズは 4 メガバイトより大きくなり得る。

【 0 1 2 9 】

図 1 5 は、一実施形態による、パーティション分離のためのビットマップと、レガシー管理エンティティのための P__K e y 抽象化とを有する例示的なファブリックを示す。図に示す実施形態では、スイッチ 1 4 0 0 は、ある数のポート 1 4 2 0 を備えることができ、これらのポートの各々は、上述したように、P__K e y 実施のためのビットマップに基づく P__K e y テーブル 1 4 1 0 に関連付けることができる。外部メモリ 1 5 3 0 は、スイッチ 1 4 0 0 によってアクセス可能である。

10

【 0 1 3 0 】

一実施形態によれば、サブネットマネージャ 1 4 4 0（例えばレガシーサブネットマネージャ）は、ビットマップに基づく P__K e y テーブル 1 4 1 0 を照会 / 設定することができないことがあり得る。このような状況において、外部メモリ 1 5 3 0 内に P__K e y テーブル抽象化 1 4 5 0（すなわち「仮想 P__K e y テーブル」）を設けることができる。次いで、サブネットマネージャは、サブネットマネージャが連想メモリ / ルックアップに基づいてレガシー P__K e y テーブルと対話するであろうのと同じ態様で、仮想 P__K e y テーブルを照会し、それと対話することができる。より具体的には、S M 1 4 4 0 は、スイッチ 1 4 0 0 に関連付けられる S M A 1 4 4 1 とインターフェイスしてそのような照会を実行することができる。このインターフェイスは、P__K e y 管理のために、I B 規格で定義された S M A 属性および方法を利用することができるため、S M は基礎となる実装についての知識がない。

20

【 0 1 3 1 】

一実施形態によれば、仮想 P__K e y テーブルに含まれる各 P__K e y は 1 6 ビットを含むため、そして仮想 P__K e y テーブルに含まれるべき可能性のある P__K e y 値の数（すなわち、最大 6 4 K P__K e y まで）のため、仮想 P__K e y テーブルのサイズは、外部メモリ 1 5 3 0 の量によって制限され得る。上述した状況とは異なり、外部メモリは、しばしば、スイッチ自体に含まれるメモリの量よりもはるかに大きくなり得る。そのような状況では、仮想 P__K e y テーブルの潜在的なサイズは、図 1 4 に関して上述したもののよりもはるかに大きい。そのような場合、たとえば、各スイッチポートに対して 6 4 K の P__K e y を表現することが可能であろう。

30

【 0 1 3 2 】

図 1 6 は、一実施形態による、高性能コンピューティング環境においてスケーラブルなビットマップに基づく P__K e y テーブル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 1 6 1 0 において、この方法は、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少なくとも 1 つのサブネットを提供することができ、少なくとも 1 つのサブネットは 1 つ以上のスイッチを含み、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、複数のスイッチポートを含み、少なくとも 1 つのサブネットはさらに、複数のホストチャネルアダプタを含み、各ホストチャネルアダプタは、少なくとも 1 つのホストチャネルアダプタポートを含み、少なくとも 1 つのサブネットはさらに、複数のエンドノードを含み、複数のエンドノードの各々は、複数のホストチャネルアダプタのうちの少なくとも 1 つのホストチャネルアダプタに関連付けられる。

40

【 0 1 3 3 】

ステップ 1 6 2 0 において、この方法は、複数の物理ホストおよび仮想マシンの各々を複数のパーティションの少なくとも 1 つに関連付けることができ、複数のパーティションの各々は P__K e y 値に関連付けられる。

50

【 0 1 3 4 】

ステップ 1 6 3 0 において、この方法は、スイッチポートの各々を、複数のビットマップに基づく P __ K e y テーブルのうちのあるビットマップに基づく P __ K e y テーブルに関連付けることができる。

【 0 1 3 5 】

ステップ 1 6 4 0 において、この方法は、ホストチャネルアダプタポートの各々を、複数のビットマップに基づく P __ K e y テーブルのうちのあるビットマップに基づく P __ K e y テーブルに関連付けることができる。

【 0 1 3 6 】

図 1 7 は、一実施形態による、高性能コンピューティング環境においてビットマップに基づくハードウェア実現例を使用してレガシー P __ K e y テーブル抽象化をサポートするための方法のフローチャートである。この方法は、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少なくとも1つのサブネットを提供することができ、少なくとも1つのサブネットは、1つ以上のスイッチを含み、1つ以上のスイッチは少なくとも1つのリーフスイッチを含み、1つ以上のスイッチの各々は、複数のスイッチポートを含み、少なくとも1つのサブネットはさらに、複数のホストチャネルアダプタを含み、各ホストチャネルアダプタは、少なくとも1つのホストチャネルアダプタポートを含み、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、少なくとも1つのサブネットはさらに、複数のエンドノードを含み、複数のエンドノードの各々は、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、少なくとも1つのサブネットはさらに、複数の仮想マシンを含み、複数の仮想マシンの各々は少なくとも1つ仮想機能に関連付けられる。

【 0 1 3 7 】

ステップ 1 7 2 0 において、この方法は、複数のエンドノードの各々を複数のパーティションの少なくとも1つに関連付けることができ、複数のパーティションの各々は P __ K e y 値に関連付けられる。

【 0 1 3 8 】

ステップ 1 7 3 0 において、この方法は、1つ以上のスイッチのスイッチポートの各々を、複数のビットマップに基づく P __ K e y テーブルのうちのあるビットマップに基づく P __ K e y テーブルに関連付けることができる。

【 0 1 3 9 】

ステップ 1 7 4 0 において、この方法は、ホストチャネルアダプタポートの各々を、複数のビットマップに基づく P __ K e y テーブルのうちのあるビットマップに基づく P __ K e y テーブルに関連付けることができる。

【 0 1 4 0 】

ステップ 1 7 5 0 において、この方法は、複数のビットマップに基づく P __ K e y テーブルの各々を、仮想 P __ K e y テーブルに関連付けることができる。

【 0 1 4 1 】

一実施形態によれば、高性能コンピューティング環境においてビットマップに基づくハードウェア実現例を使用してレガシー P __ K e y テーブル抽象化をサポートするためのシステムは、1つ以上のマイクロプロセッサと、少なくとも1つのサブネットとを備え、少なくとも1つのサブネットは、1つ以上のスイッチを含み、1つ以上のスイッチは少なくとも1つのリーフスイッチを含み、1つ以上のスイッチの各々は、複数のスイッチポートを含み、少なくとも1つのサブネットはさらに、複数のホストチャネルアダプタを含み、各ホストチャネルアダプタは、少なくとも1つのホストチャネルアダプタポートを含み、少なくとも1つのサブネットはさらに、複数のエンドノードを含み、複数のエンドノードの各々は、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、複数のエンドノードの各々は、複数のパーティションのうち少なくとも1つに関連付けられ、複数のパーティションの各々は、P __ K e y 値に関連付けられ、スイッチポートのうち1つ以上は、複数のビットマップに基づく P __ K e y テーブルのう

ちのあるビットマップに基づく P__K e y テーブルに関連付けられ、ホストチャネルアダプタポートのうちの 1 つ以上は、複数のビットマップに基づく P__K e y テーブルのうちのあるビットマップに基づく P__K e y テーブルに関連付けられ、複数のビットマップに基づく P__K e y テーブルの各々は、仮想 P__K e y テーブルに関連付けられる。

【 0 1 4 2 】

一実施形態によれば、上記のシステムは、さらに、複数のエンドノードの 1 つを介して動作するサブネットマネージャをさらに備え、サブネットマネージャは、1 つ以上のスイッチの各々上の複数のポートを通る許可されたトラフィックおよび許可されないトラフィックを判断する。

【 0 1 4 3 】

一実施形態によれば、上記のシステムにおいて、サブネットマネージャは、関連付けられる仮想 P__K e y テーブルを介して、1 つ以上のスイッチの各々上の複数のポートの各々を通る許可されたトラフィックおよび許可されないトラフィックの判断に基づいて、ビットマップに基づく P__K e y テーブルの各々を構成する。

【 0 1 4 4 】

一実施形態によれば、上記システムにおいて、仮想 P__K e y テーブルの各々は、それぞれのスイッチおよびホストチャネルアダプタにおいてメモリ上にホストされる。

【 0 1 4 5 】

一実施形態によれば、上記システムにおいて、各仮想 P__K e y テーブルは、それぞれのスイッチおよびホストチャネルアダプタにおけるメモリ量に基づいてサイズが制限される。

【 0 1 4 6 】

一実施形態によれば、上記のシステムにおいて、各仮想 P__K e y テーブルは、それぞれのスイッチおよびホストチャネルアダプタから外部のメモリ上にホストされる。

【 0 1 4 7 】

一実施形態によれば、上記のシステムは、2 つ以上のサブネットを備え、2 つ以上のサブネットの各々は、2 つ以上のサブネットの各々における少なくとも 1 つのルータポートによって相互接続される。

【 0 1 4 8 】

一実施形態によれば、高性能コンピューティング環境においてビットマップに基づくハードウェア実装を使用してレガシー P__K e y テーブル抽象化をサポートするための方法は、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 1 つのサブネットを提供することを備え、少なくとも 1 つのサブネットは、1 つ以上のスイッチを含み、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、複数のスイッチポートを含み、少なくとも 1 つのサブネットはさらに、複数のホストチャネルアダプタを含み、各ホストチャネルアダプタは、少なくとも 1 つのホストチャネルアダプタポートを含み、少なくとも 1 つのサブネットはさらに、複数のエンドノードを含み、複数のエンドノードの各々は、複数のホストチャネルアダプタのうちの少なくとも 1 つのホストチャネルアダプタに関連付けられ、方法はさらに、複数のエンドノードの各々を複数のパーティションの少なくとも 1 つに関連付けることを備え、複数のパーティションの各々は P__K e y 値に関連付けられ、方法はさらに、1 つ以上のスイッチポートの各々を、複数のビットマップに基づく P__K e y テーブルのうちのあるビットマップに基づく P__K e y テーブルに関連付けることと、ホストチャネルアダプタポートの各々を、複数のビットマップに基づく P__K e y テーブルのうちのあるビットマップに基づく P__K e y テーブルに関連付けることと、複数のビットマップに基づく P__K e y テーブルの各々を、仮想 P__K e y テーブルに関連付けることとを備える。

【 0 1 4 9 】

一実施形態によれば、上記の方法はさらに、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、複数のエンドノードの 1 つを介して動作するサブネットマネージャを提供することと、サブネットマネージャによって、1 つ以上のスイッチの各々

10

20

30

40

50

上の複数のポートを通る許可されたトラフィックおよび許可されないトラフィックを判断することとをさらに備える。

【0150】

一実施形態によれば、上記の方法はさらに、サブネットマネージャによって、関連付けられる仮想 P__Key テーブルを介して、1つ以上のスイッチの各々上の複数のポートの各々を通る許可されたトラフィックおよび許可されないトラフィックの判断に基づいて、ビットマップに基づく P__Key テーブルの各々を構成することをさらに備える。

【0151】

一実施形態によれば、上記の方法において、仮想 P__Key テーブルの各々は、それぞれのスイッチおよびホストチャネルアダプタにおいてメモリ上にホストされる。

10

【0152】

一実施形態によれば、上記の方法において、各仮想 P__Key テーブルの各々は、それぞれのスイッチおよびホストチャネルアダプタにおけるメモリ量に基づいてサイズが制限される。

【0153】

一実施形態によれば、上記の方法において、仮想 P__Key テーブルの各々は、それぞれのスイッチおよびホストチャネルアダプタから外部のメモリ上にホストされる。

【0154】

一実施形態によれば、上記の方法において、1つ以上のサブネットは、2つ以上のサブネットを備え、2つ以上のサブネットの各々は、2つ以上のサブネットの各々における少なくとも1つのルータによって相互接続される。

20

【0155】

一実施形態によれば、高性能コンピューティング環境においてビットマップに基づくハードウェア実装を使用してレガシー P__Key テーブル抽象化をサポートするための命令をそこに記憶して含む、非一時的なコンピュータ可読記憶媒体であって、命令は、1つ以上のコンピュータによって読み取られ実行されると、1つ以上のコンピュータに、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少なくとも1つのサブネットを提供することを含むステップを実行させ、少なくとも1つのサブネットは、1つ以上のスイッチを含み、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、複数のスイッチポートを含み、少なくとも1つのサブネットはさらに、複数のホストチャネルアダプタを含み、各ホストチャネルアダプタは、少なくとも1つのホストチャネルアダプタポートを含み、少なくとも1つのサブネットはさらに、複数のエンドノードを含み、複数のエンドノードの各々は、複数のホストチャネルアダプタのうちの少なくとも1つのホストチャネルアダプタに関連付けられ、命令は、さらに、1つ以上のコンピュータによって読み取られ実行されると、1つ以上のコンピュータに、複数のエンドノードの各々を複数のパーティションの少なくとも1つに関連付けることを含むステップを実行させ、複数のパーティションの各々は P__Key 値に関連付けられ、命令は、さらに、1つ以上のコンピュータによって読み取られ実行されると、1つ以上のコンピュータに、1つ以上のスイッチポートの各々を、複数のビットマップに基づく P__Key テーブルのうちのあるビットマップに基づく P__Key テーブルに関連付けることと、ホストチャネルアダプタポートの各々を、複数のビットマップに基づく P__Key テーブルのうちのあるビットマップに基づく P__Key テーブルに関連付けることと、複数のビットマップに基づく P__Key テーブルの各々を、仮想 P__Key テーブルに関連付けることとを含むステップを実行させる。

30

40

【0156】

一実施形態によれば、上記の非一時的なコンピュータ可読記憶媒体において、上記のステップは、さらに、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サブネットマネージャを提供することを含み、サブネットマネージャは複数のエンドノードの1つを介して動作し、上記ステップはさらに、サブネットマネージャによって、1つ以上のスイッチの各々上の複数のポートを通る許可されたトラフィックおよび許可されないト

50

ラフィックを判断することを含む。

【0157】

一実施形態によれば、上記の非一時的なコンピュータ可読記憶媒体において、上記のステップは、さらに、サブネットマネージャによって、関連付けられる仮想P__Keyテーブルを介して、1つ以上のスイッチの各々上の複数のポートの各々を通る許可されたトラフィックおよび許可されないトラフィックの判断に基づいて、ビットマップに基づくP__Keyテーブルの各々を構成することをさらに含む。

【0158】

一実施形態によれば、上記の非一時的なコンピュータ可読記憶媒体において、仮想P__Keyテーブルの各々は、それぞれのスイッチおよびホストチャネルアダプタにおいてメモリ上にホストされる。

10

【0159】

一実施形態によれば、上記の非一時的なコンピュータ可読記憶媒体において、各仮想P__Keyテーブルの各々は、それぞれのスイッチおよびホストチャネルアダプタにおけるメモリ量に基づいてサイズが制限される。

【0160】

一実施形態によれば、上記の非一時的なコンピュータ可読記憶媒体において、仮想P__Keyテーブルの各々は、それぞれのスイッチおよびホストチャネルアダプタから外部のメモリ上にホストされる。

【0161】

20

一実施形態によれば、コンピュータプログラムは、コンピュータシステムによって実行されると、コンピュータシステムに上記の方法を実行させる、機械可読フォーマットにおけるプログラム命令を含む。

【0162】

一実施形態によれば、非一時的な機械可読データ記憶媒体に記憶される上記のコンピュータプログラムを備えるコンピュータプログラム製品。

【0163】

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。

30

【0164】

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。

40

【0165】

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。

【0166】

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 int

50

egrated 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。

【0167】

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。

10

【0168】

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。

【0169】

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。

20

【0170】

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。

30

【図 1】

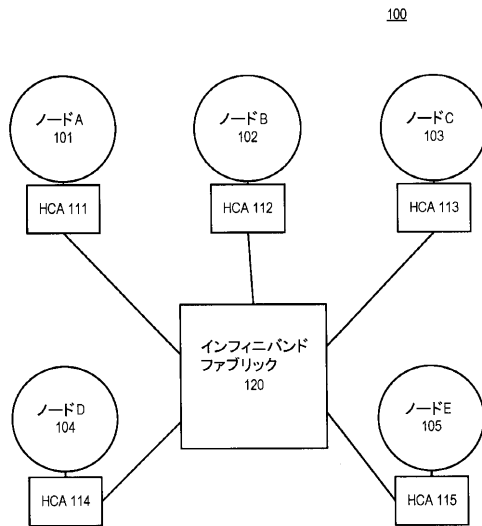


FIGURE 1

【図 2】

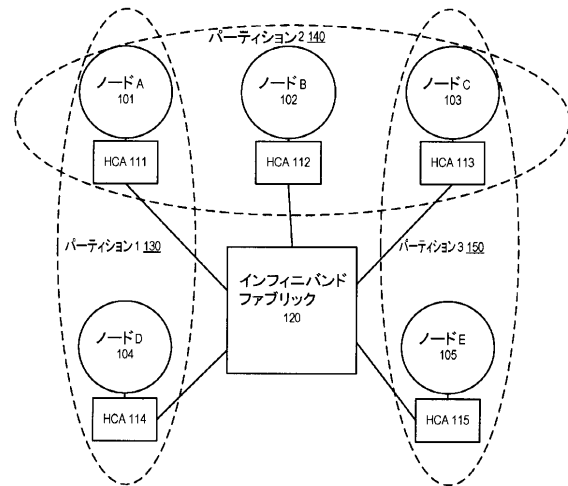


FIGURE 2

【図 3】

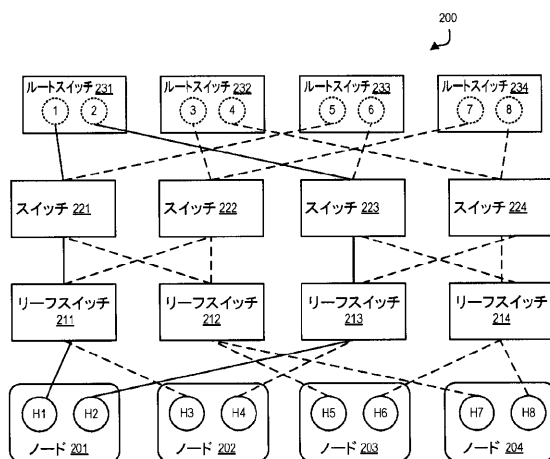


FIGURE 3

【図 4】

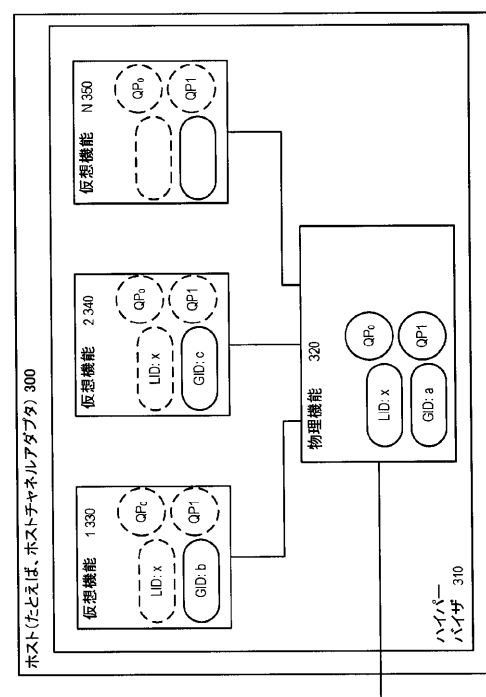


FIGURE 4

【図 5】

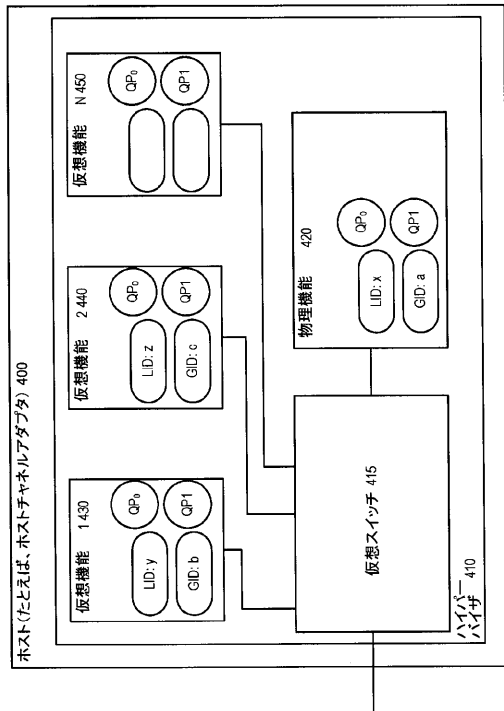


FIGURE 5

【図 6】

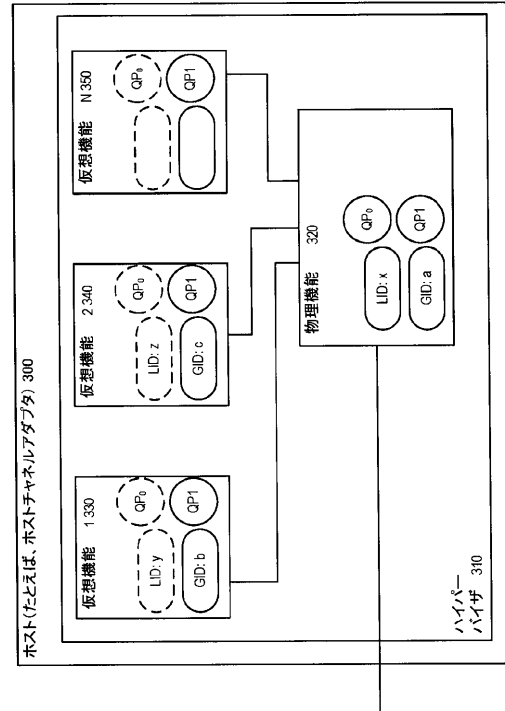


FIGURE 6

【図 7】

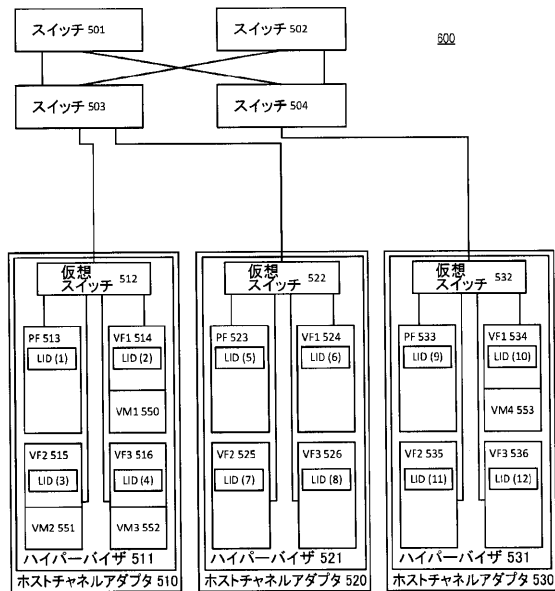


FIGURE 7

【図 8】

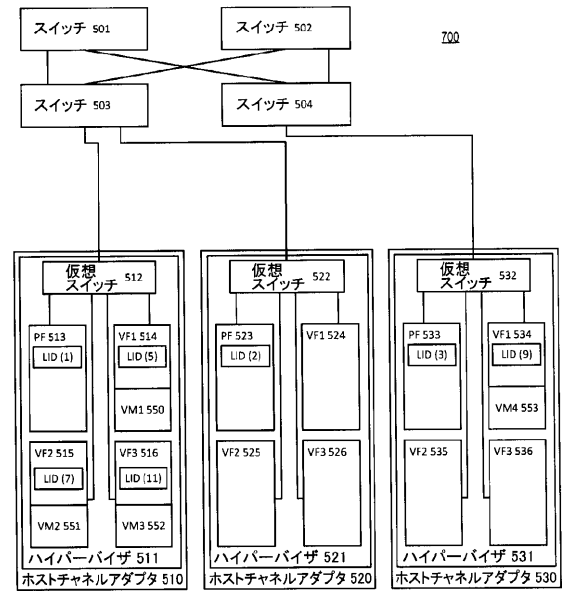


FIGURE 8

【図 9】

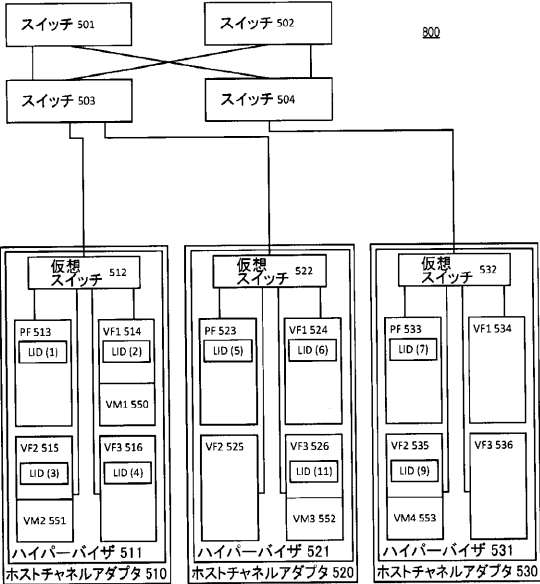


FIGURE 9

【図 10】

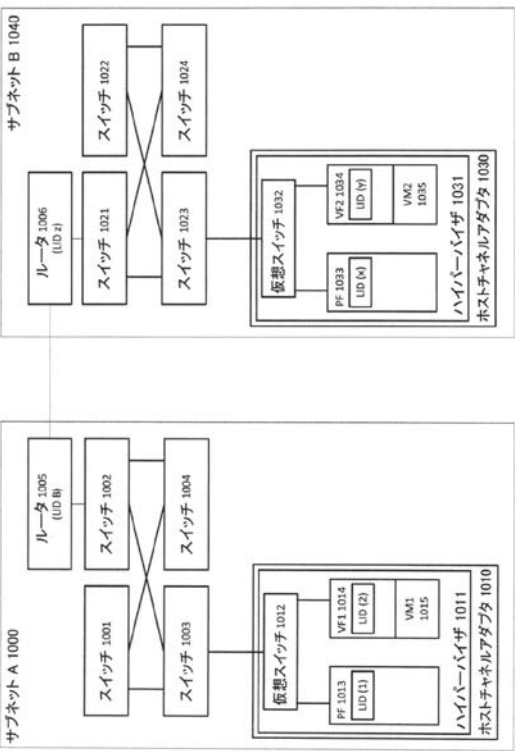


FIGURE 10

【図 11】

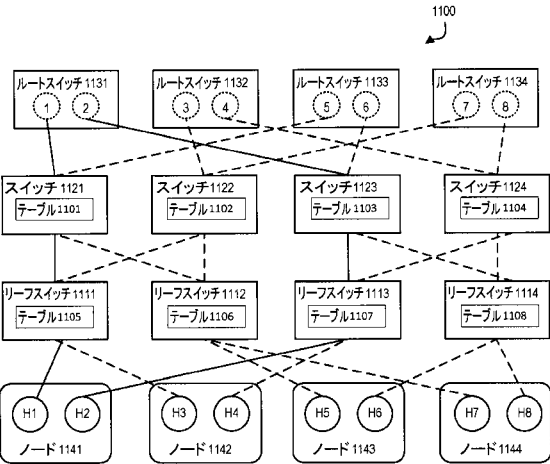


FIGURE 11

【図 12】

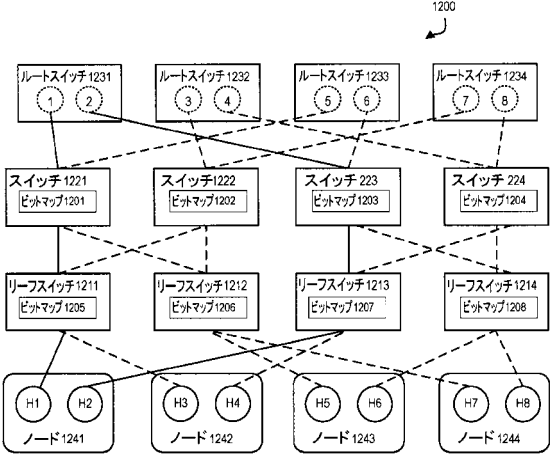
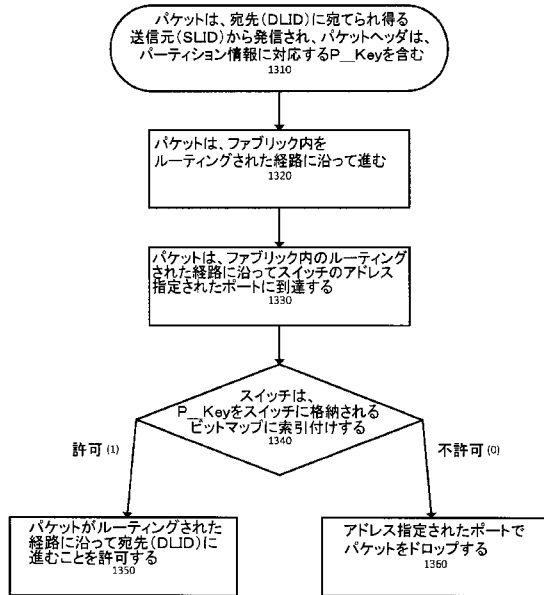


FIGURE 12

【図 13】



【図 14】

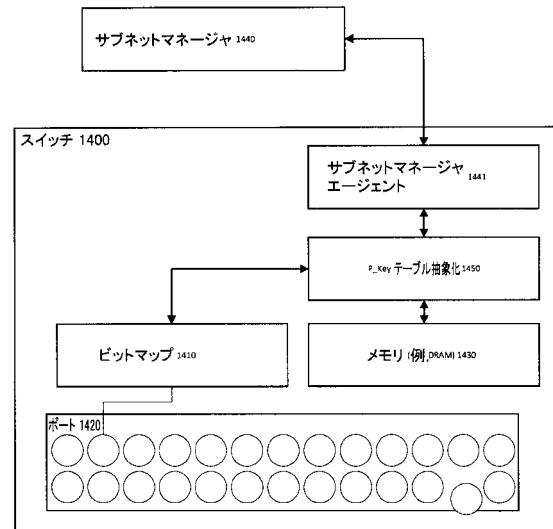
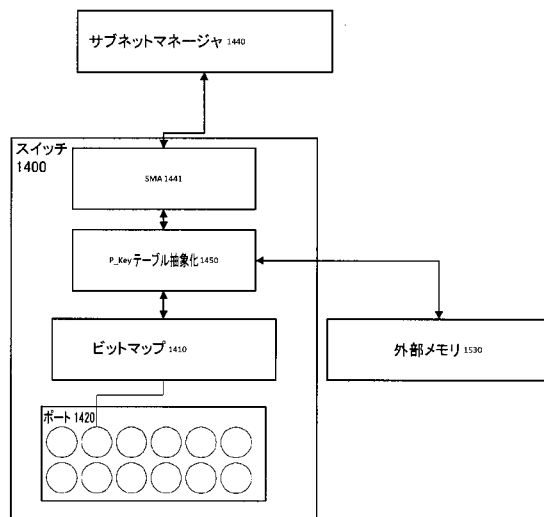
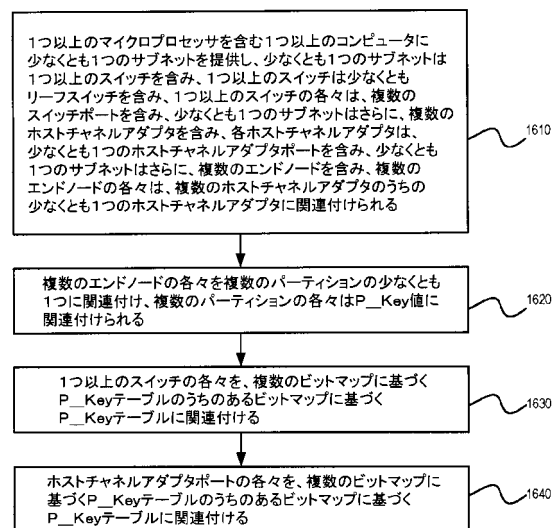


FIGURE 13

【図 15】



【図 16】



【 図 17 】

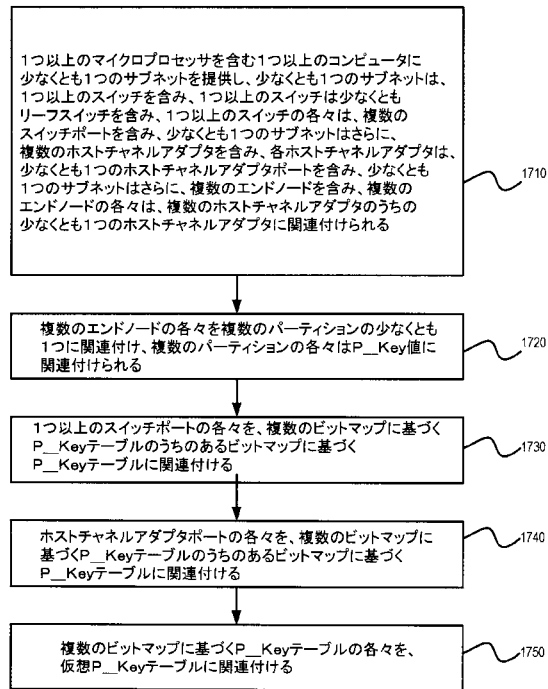


FIGURE 17

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No

PCT/US2017/015156

A. CLASSIFICATION OF SUBJECT MATTER

INV. H04L12/931 G06F9/44
ADD.

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

H04L G06F

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

EPO-Internal, WPI Data

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	<p>FEROZ ZAHID ET AL: "Partition-Aware Routing to Improve Network Isolation in Infiniband Based Multi-tenant Clusters", 2015 15TH IEEE/ACM INTERNATIONAL SYMPOSIUM ON CLUSTER, CLOUD AND GRID COMPUTING, 1 May 2015 (2015-05-01), pages 189-198, XP055322241, DOI: 10.1109/CCGrid.2015.96 ISBN: 978-1-4799-8006-2 paragraph [II.B.Infiniband.Architecture.B.Partitioning]</p> <p style="text-align: center;">----- -/--</p>	1-22

☒ Further documents are listed in the continuation of Box C.☐ See patent family annex.

* Special categories of cited documents :

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier application or patent but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art

"&" document member of the same patent family

Date of the actual completion of the international search

28 March 2017

Date of mailing of the international search report

05/04/2017

Name and mailing address of the ISA/

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040,
Fax: (+31-70) 340-3016

Authorized officer

Lefebvre, Laurent

INTERNATIONAL SEARCH REPORT

International application No

PCT/US2017/015156

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	<p>VISHNU A ET AL: "Performance Modeling of Subnet Management on Fat Tree InfiniBand Networks using OpenSM", PARALLEL AND DISTRIBUTED PROCESSING SYMPOSIUM, 2005. PROCEEDINGS. 19TH IEEE INTERNATIONAL DENVER, CO, USA 04-08 APRIL 2005, PISCATAWAY, NJ, USA, IEEE, 4 April 2005 (2005-04-04), pages 296b-296b, XP010785940, DOI: 10.1109/IPDPS.2005.339 ISBN: 978-0-7695-2312-5 paragraph [2.1.Infiniband.Subnet.Management] paragraph [3.1.Subnet.Discovery]</p> <p>-----</p>	1-22
A	<p>FRANCESCO FUSCO ET AL: "Real-time creation of bitmap indexes on streaming network data", THE VLDB JOURNAL ; THE INTERNATIONAL JOURNAL ON VERY LARGE DATA BASES, SPRINGER, BERLIN, DE, vol. 21, no. 3, 30 July 2011 (2011-07-30), pages 287-307, XP035056143, ISSN: 0949-877X, DOI: 10.1007/S00778-011-0242-X paragraph [4.2.1.Bitmap.Indexes]</p> <p>-----</p>	1,8,15, 21,22

フロントページの続き

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

(72)発明者 ナラシンハマーシー, プラブフナンダン

アメリカ合衆国、 9 4 0 6 5 カリフォルニア州、 レッドウッド・ショアーズ、 オラクル・パーク
ウェイ、 5 0 0、 エム/エス・ 5 ・オウ・ピィ・ 7

(72)発明者 ホレン, リネ

ノルウェー、 1 9 0 0 フェツンド、 ビタセン、 1 7

Fターム(参考) 5K030 GA03 HA08 HC13 HD03 KX17 MA14 MD07