

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6122483号  
(P6122483)

(45) 発行日 平成29年4月26日 (2017. 4. 26)

(24) 登録日 平成29年4月7日 (2017. 4. 7)

(51) Int. Cl.

F I

G 0 6 F 9/46 (2006. 01)

G 0 6 F 9/46 3 5 0

G 0 6 F 9/50 (2006. 01)

G 0 6 F 9/46 4 6 2 Z

G 0 6 F 13/10 (2006. 01)

G 0 6 F 13/10 3 3 0 C

請求項の数 12 (全 21 頁)

(21) 出願番号 特願2015-503440 (P2015-503440)  
 (86) (22) 出願日 平成25年3月25日 (2013. 3. 25)  
 (65) 公表番号 特表2015-518602 (P2015-518602A)  
 (43) 公表日 平成27年7月2日 (2015. 7. 2)  
 (86) 国際出願番号 PCT/US2013/033754  
 (87) 国際公開番号 W02013/148600  
 (87) 国際公開日 平成25年10月3日 (2013. 10. 3)  
 審査請求日 平成28年1月19日 (2016. 1. 19)  
 (31) 優先権主張番号 61/615, 731  
 (32) 優先日 平成24年3月26日 (2012. 3. 26)  
 (33) 優先権主張国 米国 (US)  
 (31) 優先権主張番号 61/693, 703  
 (32) 優先日 平成24年8月27日 (2012. 8. 27)  
 (33) 優先権主張国 米国 (US)

(73) 特許権者 502303739  
 オラクル・インターナショナル・コーポレ  
 イション  
 アメリカ合衆国カリフォルニア州9406  
 5レッドウッド・シティー, オラクル・パ  
 ークウェイ500  
 (74) 代理人 110001195  
 特許業務法人深見特許事務所  
 (72) 発明者 ヨンセン, ビョルン・ダグ  
 ノルウェー、エヌー0687 オスロ、ビ  
 ルベルクグレンダ、9

最終頁に続く

(54) 【発明の名称】 拡張ホストチャネルアダプタ (HCA) モデルに基づいてバーチャルマシンのライブマイグレーションをサポートするためのシステムおよび方法

(57) 【特許請求の範囲】

【請求項 1】

ネットワークにおいてバーチャルマシンライブマイグレーション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、

前記1つ以上のマイクロプロセッサに関連付けられるファブリックアダプタとを含み、前記ファブリックアダプタは複数のバーチャルホストチャネルアダプタ (HCA) に関連付けられており、前記バーチャルホストチャネルアダプタ (vHCA) の各々は、別個のキューペア (QP) スペースに関連付けられており、前記システムはさらに、

第1のホストから第2のホストにライブマイグレーションを行なうよう動作する少なくとも1つのバーチャルマシンを含み、前記少なくとも1つのバーチャルマシンには、前記キューペア (QP) スペースにおけるキューペア (QP) に関連付けられる前記バーチャルホストチャネルアダプタ (vHCA) が付与されており、

前記キューペア (QP) は、前記ライブマイグレーションについてピアQPに信号送信するよう動作し、マイグレーションの後に前記ピアQPにアドレス情報を提供する、システム。

【請求項 2】

前記少なくとも1つのバーチャルマシンを管理する少なくとも1つのバーチャルマシンモニタをさらに含み、前記バーチャルマシンの各々は、プライベート仮想アドレススペースに関連付けられる、請求項1に記載のシステム。

10

20

## 【請求項 3】

前記ファブリックアダプタは、仮想スイッチモデルおよび拡張ホストチャネルアダプタ（HCA）モデルの少なくとも1つを実現する、請求項1 または2に記載のシステム。

## 【請求項 4】

中断状態は、マイグレートされたQPとマイグレートされたQPに接続されるリモートピアQPとの両方に適用可能である、請求項1 ～3のいずれか1項に記載のシステム。

## 【請求項 5】

特別なヘッダ情報および/または非送信請求メッセージが、ピア同士の間での通信および通信状態の更新をサポートするために使用可能である、請求項1 ～4のいずれか1項に記載のシステム。

10

## 【請求項 6】

無効なvHCAを目標とする入力パケットがホストにて受け取られると、前記ホストは、ローカル非同期イベントを生成し、前記入力パケットのソースに非送信請求パケットを送信可能であり、前記入力パケットのソースに関連付けられるキューペア（QP）は、自動的に中断状態への遷移を実行可能である、請求項1 ～5のいずれか1項に記載のシステム。

## 【請求項 7】

マイグレートされたvHCAから送られるパケットは、マイグレートされた前記vHCAが新しい位置で再び作動状態になる前に更新可能であるマイグレーションカウンタフィールドを包含可能である、請求項1 ～6のいずれか1項に記載のシステム。

20

## 【請求項 8】

マイグレートされたvHCAは、非中断ペンディング状態のQPで開始可能である、請求項1 ～7のいずれか1項に記載のシステム。

## 【請求項 9】

LID値をHCAポートに割り当て可能なサブネットマネージャをさらに含み、前記LID値は、前記HCAポートに既に関連付けられる任意の他のLID値から独立して、関連するHCAポートへのパケット転送を促進するよう前記サブネットマネージャが使用することを決定する、請求項1 ～8のいずれか1項に記載のシステム。

## 【請求項 10】

前記サブネットマネージャは、アドレッシングのためのユニークLIDを必要とするVMにadmin LIDを割り当て可能である、請求項9に記載のシステム。

30

## 【請求項 11】

ネットワークにおいてバーチャルマシンライブマイグレーションをサポートするための方法であって、

ファブリックアダプタを複数のバーチャルホストチャネルアダプタ（HCA）に関連付けるステップを含み、前記バーチャルホストチャネルアダプタ（vHCA）の各々は、別個のキューペア（QP）スペースに関連付けられており、前記方法はさらに、

少なくとも1つのバーチャルマシンを介して、第1のホストから第2のホストにライブマイグレーションを行なうステップを含み、前記少なくとも1つのバーチャルマシンには、前記キューペア（QP）スペースにおけるキューペア（QP）に関連付けられる前記バーチャルホストチャネルアダプタ（vHCA）が付与されており、前記方法はさらに、

40

前記キューペア（QP）を介して、前記ライブマイグレーションについてピアQPに信号送信し、マイグレーションの後に前記ピアQPにアドレス情報を提供するステップを含む、方法。

## 【請求項 12】

コンピュータ読取可能なプログラムであって、実行されると、システムに、

ファブリックアダプタを複数のバーチャルホストチャネルアダプタ（HCA）に関連付けるステップを行わせ、前記バーチャルホストチャネルアダプタ（vHCA）の各々は、別個のキューペア（QP）スペースに関連付けられており、さらに、

少なくとも1つのバーチャルマシンを介して、第1のホストから第2のホストにライブ

50

マイグレーションを行なうステップを行わせ、前記少なくとも 1 つのバーチャルマシンには、前記キューペア ( Q P ) スペースにおけるキューペア ( Q P ) に関連付けられる前記バーチャルホストチャネルアダプタ ( v H C A ) が付与されており、さらに、

前記キューペア ( Q P ) を介して、前記ライブマイグレーションについてピア Q P に信号送信し、マイグレーションの後に前記ピア Q P にアドレス情報を提供するステップを行わせる、コンピュータ読取可能なプログラム。

【発明の詳細な説明】

【技術分野】

【 0 0 0 1 】

著作権表示

この特許文書の開示の一部は、著作権の保護下にある内容を含む。著作権所有者は、特許商標庁の特許ファイルまたはレコードに現れるので、誰でも当該特許文書または特許開示を複製することについて異議はないが、そうでなければ如何なる場合でもすべての著作権を留保する。

【背景技術】

【 0 0 0 2 】

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、特にコンピュータシステム仮想化をサポートすることに関する。

【 0 0 0 3 】

背景

より大きなクラウドコンピューティングアーキテクチャが導入されるにつれて、従来のネットワークおよびストレージに関連付けられる性能および管理上のボトルネックが重要な問題になっている。インフィニバンド ( I n f i n i B a n d ( I B ) ) 技術は、クラウドコンピューティングファブリックの基礎として、展開が増加している。これは、本発明の実施例が対応することが意図される一般的な領域である。

【発明の概要】

【課題を解決するための手段】

【 0 0 0 4 】

概要

本願明細書では、ネットワークにおいてバーチャルマシンライブマイグレーションをサポートするためのシステムおよび方法が記載される。ファブリックアダプタは、複数のバーチャルホストチャネルアダプタ ( H C A ) に関連付けられ得、上記バーチャルホストチャネルアダプタ ( v H C A ) の各々は、別個のキューペア ( Q P ) スペースに関連付けられている。少なくとも 1 つのバーチャルマシンが第 1 のホストから第 2 のホストにライブマイグレーションを行なうよう動作し、上記少なくとも 1 つのバーチャルマシンには、上記キューペア ( Q P ) スペースにおけるキューペア ( Q P ) に関連付けられる上記バーチャルホストチャネルアダプタ ( v H C A ) が付与されており、上記キューペア ( Q P ) は、上記ライブマイグレーションについてピア Q P に信号送信するよう動作し、マイグレーションの後に上記ピア Q P にアドレス情報を提供する。

【図面の簡単な説明】

【 0 0 0 5 】

【図 1】本発明の実施例に従った仮想インターフェイスアーキテクチャ ( V I A ) ネットワークデバイスの図を示す図である。

【図 2】本発明の実施例に従った、仮想インターフェイス ( V I ) に関連付けられる異なる状態の図を示す図である。

【図 3】本発明の実施例に従った、キューペア ( Q P ) に関連付けられる異なる状態の図を示す図である。

【図 4】本発明の実施例に従った、ネットワークおよびルーティングの観点からのバーチャルマシン ( V M ) のライブマイグレーションの図を示す図である。

10

20

30

40

50

【図 5】本発明の実施例に従った、バーチャルマシン（VM）ライブマイグレーションの前の仮想化環境の図を示す図である。

【図 6】本発明の実施例に従った、バーチャルマシン（VM）ライブマイグレーションの後の仮想化環境の図を示す図である。

【図 7】本発明の実施例に従った、仮想化環境におけるバーチャルマシンのライブマイグレーションをサポートするための例示的なフローチャートを示す図である。

【図 8】本発明の実施例に従った、インフィニバンドアーキテクチャ（IBA）ネットワークデバイスの図を示す図である。

【図 9】本発明の実施例に従った、マイグレートされたキューペア（QP）とリモートピア QP との間の通信をサポートする図を示す図である。

10

【図 10】本発明の実施例に従った、拡張 vHCA モデルに基づいてバーチャルマシンのライブマイグレーションをサポートするための例示的なフローチャートを示す図である。

【図 11】本発明の実施例に従った、仮想化環境においてローカル識別子（LID）の割り当てをサポートする図を示す図である。

【発明を実施するための形態】

【0006】

#### 詳細な説明

本発明は、限定目的ではなく例示目的で、添付の図面の図において示される。当該図面においては、同様の参照符号が同様の要素を示す。なお、この開示における「ある」、「1つ」または「いくつか」の実施例への参照は、必ずしも同じ実施例に対してなされるものではなく、このような参照は、少なくとも 1 つの実施例を意味する。

20

【0007】

以下の本発明の記載は、高性能ネットワークについての例として、インフィニバンド（IB）ネットワークを使用する。他のタイプの高性能ネットワークが限定なしで使用され得るということは、当業者には明らかであろう。さらに、以下の本発明の記載は、仮想化モデルについての例として、Xen 仮想化モデルを使用する。他のタイプの仮想化モデルが限定なしで使用され得るということは、当業者には明らかであろう。

【0008】

本願明細書において、ネットワークにおいてバーチャルマシン（virtual machine（VM））ライブマイグレーション（live migration）をサポートすることができるシステムおよび方法が記載される。

30

【0009】

本発明の実施例に従うと、仮想化は、クラウドコンピューティングでの効率的なリソース活用および弾力的なリソースアロケーションに有益になり得る。ライブマイグレーションは、物理サーバ同士の間でバーチャルマシン（VM）をアプリケーショントランスペアレントな態様で移動させることによって、リソースの使用を最適化することを可能にする。したがって仮想化は、コンソリデーション、リソースのオンデマンドプロビジョニング、およびライブマイグレーションによる弾性を可能にし得る。

【0010】

#### インフィニバンド（IB）アーキテクチャ

40

IB アーキテクチャはシリアル・ポイント・ツー・ポイント全二重技術である。IB ネットワークは、サブネットとも称され得る。サブネットは、スイッチとポイント・ツー・ポイントリンクとを使用して相互接続されるホストのセットからなる。IB サブネットは、サブネットにおけるすべてのスイッチ、ルータおよびホストチャネルアダプタ（HCA）のコンフィギュレーションを含み、ネットワークを初期化および起動することを担う、少なくとも 1 つのサブネットマネージャ（subnet manager（SM））を含み得る。

【0011】

IB は、リモートダイレクトメモリアクセス（remote direct memory access（RDMA））および従来の送信 / 受信セマンティックスの両方を提供するために、伝送サービスの豊富なセットをサポートする。使用される伝送サービスから独立して、IB HCA は

50

キューペア ( Q P ) を使用して通信する。 Q P は、通信セットアップの間に作り出され、 Q P ナンバー、 H C A ポート、宛先 L I D、キューサイズおよび伝送サービスといった提供される初期属性のセットを有し得る。 H C A は、多くの Q P を扱い得、各 Q P は、送信キュー ( send queue ( S Q ) ) および受信キュー ( receive queue ( R Q ) ) のようなキューのペアからなっており、このような 1 つのペアが、通信に参加する各エンドノードに存在する。送信キューは、リモートノードに転送されるワーク要求を保持する一方、受信キューは、リモートノードから受け取ったデータで何を行なうべきかについての情報を保持する。 Q P に加えて、各 H C A は、送信キューおよび受信キューのセットに関連付けられる 1 つ以上の完了キュー ( completion queue ( C Q ) ) を有する。 C Q は、送信および受信キューにポスティングしたワーク要求についての完了通知を保持する。通信の複雑性はユーザから隠されているが、 Q P 状態情報は H C A に維持される。

10

#### 【 0 0 1 2 】

各物理的な I B デバイスには、 L I D およびグローバルユニーク識別子 ( globally unique identifier ( G U I D ) ) の 2 つのアドレスが割り当てられる。 L I D は、サブネット内で I B パケットをルーティングするよう使用され、 G U I D は、物理的 I B デバイスを一意に示すハードウェアアドレスである。 6 4 ビット G U I D は、ローカルの ( 6 4 ビット ) サブネットプレフィックスと組み合わされてグローバル識別子を形成する。 1 2 8 ビットグローバル識別子は、 I B サブネット同士の間で I B パケットをルーティングするよう使用される。

#### 【 0 0 1 3 】

20

##### 入出力 ( I / O ) 仮想化

I / O 仮想化 ( I / O Virtualization ( I O V ) ) は、存在する物理リソースに 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 の可用性を提供し得る。ストレージトラフィックおよびサーバ間通信の組合せによって、単一のサーバの I / O リソースを圧倒し得るような増加した負荷が課され、これにより、データを待っている際に、バックログおよびアイドル状態になるプロセッサが引き起こされる。 I / O 要求の数の増加により、 I O V は可用性を提供し得るとともに、現代の C P U 仮想化において見られる性能のレベルに合致するよう、 ( 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フレキシビリティを向上し得る。

#### 【 0 0 1 4 】

30

エミュレーション、準仮想化、直接割当 ( direct assignment ( D A ) ) およびシングルルート I / O 仮想化 ( single root-I / O virtualization ( S R - I O V ) ) といった、異なるタイプの I O V 技術が存在し得る。これらの I O V 技術のうち、 S R - I O V は、ほぼネイティブの性能を維持しつつ、複数の V M から単一の物理デバイスへの直接アクセスを可能にする手段により、 P C I E x p r e s s ( P C I e ) 規格を拡張し得る。したがって S R - I O V は、良好な性能およびスケーラビリティを提供し得る。

#### 【 0 0 1 5 】

S R - I O V は、各ゲストに 1 つの仮想デバイスを割り当てることによって複数のゲスト同士の間で共有され得る複数の仮想デバイスを P C I e デバイスが公開することを可能にする。各 S R - I O V デバイスは、少なくとも 1 つの物理的機能 ( physical function ( P F ) ) と、 1 つ以上の関連付けられる仮想機能 ( virtual function ( V F ) ) とを有する。 P F は、バーチャルマシンモニタ ( virtual machine monitor ( V M M ) ) またはハイパーバイザによって制御される通常の P C I e 機能であり、 V F は、軽量の P C I e 機能である。各 V F は、自身のベースアドレス ( base address ( B A R ) ) を有しており、ユニークリクエスト I D が割り当てられている。ユニークリクエスト I D は、 I / O メモリ管理ユニット ( I / O memory management unit ( I O M M U ) ) が、異なる V F に対するトラフィックストリーム同士を区別することを可能にする。 I O M M U はさらに、 P F と V F との間でメモリおよび割込の変換を適用する。

40

#### 【 0 0 1 6 】

たとえば共有ポートモデルおよび仮想スイッチモデルといった、異なるタイプの S R -

50

IOVモデルが存在し得る。共有ポートモデルでは、すべてのVFは、単一のポートアドレスおよび単一のQPネームスペースを共有し得、単一のHCAポートのみがネットワークに公開される。仮想スイッチモデルでは、各VFは、ユニークポートアドレスおよびユニークQPネームスペースを含む仮想HCAであり、デバイス上のVFごとに1つのHCAがネットワークに公開される。したがって、より複雑なハードウェアであっても、仮想スイッチモデルはIOVを簡素化し得る。

【0017】

SR - IOV対応のデバイスの使用によって、ほぼネイティブの性能およびスケラビリティの向上が与えられ得る。その一方、SR - IOVは、ライブマイグレーションおよびチェックポイント/リスタートメカニズムと完全には互換性がないわけではない。

10

【0018】

ライブマイグレーションのためのハードウェアソリューション

本発明の実施例に従って、ハードウェアソリューションは、アプリケーショントランスペアレントな態様で、フレキシブルで高性能かつスケラブルな仮想化されたIOアーキテクチャに基づき、物理サーバ同士間でのバーチャルマシン（VM）のライブマイグレーションをサポートし得る。さらに、ハードウェアソリューションは、マイグレーションのトランスペアレンシーを増加し得、IBとSR - IOVとの間の相互運用性を改善し得る。さらに、当該システムは、仮想インターフェイスアーキテクチャ（virtual interface architecture（VIA））およびミリネット（Myrinet）のような他の高速ロスレスネットワークワーキング技術に適用可能であり得る。

20

【0019】

ハードウェアソリューションを使用して、IBホストチャネルアダプタ（HCA）は、QPについて自身のリソースプールを有する仮想エンドポイントとして、各物理的な機能（VF）を区別し得る。さらに、QPNのようなQP属性は、マイグレーションの後に再使用され得る。さらに、IB HCAは、通信を中断および再開するためにQP状態に基づいて、VMライブマイグレーションをサポートし得る。

【0020】

さらに、システムは、マイグレートされたQPとマイグレートされたQPのピアとの間で休止するよう利用され得るハンドシェイクメカニズムを回避し得る。したがって、当該システムは、バーチャルマシン（VM）のライブマイグレーションが、IBにおける信頼性のある接続の再試行タイムアウトに相当するレベルにまでサービスダウンタイムを低減することを可能にする。

30

【0021】

さらにシステムは、厳格なタイミング要件を有するシナリオおよびアプリケーションを含む、VMライブマイグレーションをサポートする一般的なソリューションを示し得る。たとえば高い可用性のクラスタにおいて、欠陥のある物理サーバによってホストされたVMは、限られた時間内に別のサーバにマイグレートされ得る。マイグレーションは、欠陥が発見された後に行なわれ得るが、当該欠陥のあるサーバが回復不能なエラーになる前に行なわれ得る。ソフトウェアアプローチは、オーバーヘッドに苦しんでいるので、この低いサービスダウンタイムの要求を満たすのに不十分である。さらに、ソフトウェアアプローチは開発コストの点で高価である。

40

【0022】

仮想エンドポイント

本発明の実施例に従うと、仮想インターフェイスアーキテクチャ（VIA）ハードウェアは、各VFから作り出されたVIを区別するよう、PCI仮想スイッチを統合し得る。仮想スイッチの後ろでは、複数の独立した仮想エンドポイントがネットワークマネージャによって発見可能であり得る。

【0023】

図1は、本発明の実施例に従った仮想インターフェイスアーキテクチャ（VIA）ネットワークデバイスの図を示す。図1に示されるように、仮想インターフェイスアーキテク

50

チャ（VIA）ネットワークデバイス100は、ハードウェアレイヤー110、バーチャルマシンモニタ120、およびたとえばVM101-102といった1つ以上のバーチャルマシンを含み得る。ここで、VM101は、カーネルエージェント107に基づいて、固定されたメモリ105に関連付けられるユーザプロセス103をサポートし得る。VM102は、カーネルエージェント108に基づいて、固定されたメモリ106に関連付けられるユーザプロセス104をサポートし得る。

#### 【0024】

仮想スイッチ130はハードウェアレイヤー110上に提供され得る。仮想スイッチ130は、たとえばVF111-119といった1つ以上の仮想機能（VF）に関連付けられ得る。VF111-119の各々は、1つ以上の仮想インターフェイス（VI）を作成および管理し得る。

10

#### 【0025】

本発明の実施例に従うと、仮想スイッチモデルを使用して、VF111-119の各々が、たとえばVM101またはVM102といったゲストVMに割り当てられ得る専用のPCIエンドポイントとして分離され得る。さらに、VF111-119の各々は、その専用のVIネームスペース121-129を有し得、完全なVIエンドポイントとして考えられ得る。VIネームスペースの分離は、仮想化された環境において保護およびフレキシビリティを提供し得る。各VFは、IO動作についてVMの保護および分離プロパティを拡張するパーティション内のVIリソースにアクセスするのが制限され得る。さらに、VI属性はマイグレーションの後に再使用され得る。

20

#### 【0026】

たとえば、ハードウェアレイヤー110は、IBアーキテクチャのような高速ネットワークアーキテクチャに基づき得る。VF111-119の各々は、ユニークなポートアドレスおよびユニークなQPネームスペースを含むバーチャルホストチャネルアダプタ（virtual host channel adaptor（VHCA））であり得る。さらに、IBデバイス上のVFごとに、1つのVHCAがネットワークに公開され得る。さらに、各VHCAは、自身のQPネームスペースに関連付けられ得、IBアーキテクチャにおけるQP131および132のような1つ以上のキューペア（QP）を作成および管理し得る。図1に示されるように、キューペア（QP）131-132の各々は、セNDERキューおよびレシーバークューを含み得る。

30

#### 【0027】

さらに、VIAデバイス100は、IOメモリ管理ユニット（IOMMU）の機能性を活用し得る。VI属性は、IOMMUによって変換（translate）されるIO仮想メモリを介してアクセスされるメモリ領域に格納され得る。新しいVFが、VMに再付与されると、ハイパーバイザは、IOMMUを変更して、明示的に新しいVIを再作成することなしに物理メモリの新しいセットにIO仮想メモリを再マッピングし得る。共有されたVIネームスペースモデルと比較して、仮想スイッチモデルは、リソースマイグレーションを緩和するとともに、VI属性を再マッピングするためのマッピングテーブルの使用を回避して、動作中のアプリケーションについてトランスペアレンシーを維持する。

40

#### 【0028】

##### キューペア（QP）中断状態

図2は、本発明の実施例に従った、仮想インターフェイス（VI）に関連付けられる異なる状態の図を示す。図2に示されるように、VIのライフサイクルは、アイドル状態201、ペンディング接続状態202、接続済状態203、エラー状態204および中断状態205といった複数の状態を含み得る。

#### 【0029】

中断状態205は、進行中の通信が、たとえばライブマイグレーションの間、VIを休止状態にするために一時的に停止されることを示す。接続済状態203と中断状態205との間の遷移は、ソフトウェアが開始した要求またはVIイベント要求によって行なわれ得る。たとえば、マイグレーションの間に、ソフトウェアインターフェイスは、接続済状

50

態 203 から中断状態 205 まで、マイグレートする V I を修正し得る。

【0030】

V I が中断状態 205 に入ると、V I は、如何なる以前に開始されたメッセージ送信も完了し得、それら进行处理することなく、ハードウェアセンドキュー上の如何なる残存するメッセージもフラッシュする。ハードウェアは、残存するメッセージが、フラッシュする前に、ユーザプロセスについてのセンドキューによってキャッシュされることを保証し得る。中断状態 205 においては、如何なる新しいメッセージも、センドキューおよびレシーバキューの両方においてキューに入れられ得、処理されることなくポストされたままである。

【0031】

さらに、入力メッセージが、中断された V I を目標としている場合、中断状態 205 にピア V I を移行させるよう、中断された V I によって SuspendRequest イベントが生成され得る。このような場合、ResumeRequest イベントを受け取った後、以前の開始されたメッセージが再送信され得、不完全な送信時にはタイムアウトエラーが生成され得ない。

【0032】

新しい V F がマイグレートされた V M に付与された後およびマイグレートされた V M が新しい位置で再び開始される前に、ストップ・アンド・コピーステージ (stop-and-copy stage) が、I O M M U における新しいセットの物理メモリに I O 仮想メモリを再登録し得る。その後、ソフトウェアインターフェイスは、マイグレートされた V M の V I 状態を接続済状態 203 に修正し得る。その後、システムは、ピア V I へ転送される ResumeRequest イベントをトリガし得る。その後、ピア V I は、ResumeRequest イベントを生成することなく接続済状態 203 に移行し得る。したがって、中断状態 205 は、V M のライブマイグレーションの間にスムーズな遷移を提供し得、V I は休止状態にあり得る。

【0033】

本発明の実施例に従うと、V M ライブマイグレーションの間にダウンタイムを低減するために、最小遅延のアプローチは、すべての未完了の動作が完了する直前に、V I を中断または V F を分離することである。したがって、延長によりタイムアウトが向上し得、中断状態 205 のハンドリングの部分として、信頼性のあるトランスポートプロトコルによって与えられるハンドリングが再試行され得る。延長された再試行およびタイムアウトロジックは、ピア V I またはマイグレートされる V I のいずれかについて致命的なエラーを引き起こすことなく V I を中断状態 205 に移行し得る。

【0034】

したがって、マイグレーションの後に V I の可用性に対してハードなリアルタイムの制約を課すことなく、V I はマイグレートされ得る。このアプローチはさらに、ネットワークの遅延によりイベント報告 (SuspendRequest および ResumeRequest) が V I 同士の間で生成されない場合をカバーし得る。さらに、複数の V I を備えた場合において、各 V I は上記の動作を同時に行ない得る。

【0035】

図 3 は、本発明の実施例に従った、キューペア (Q P) に関連付けられる異なる状態の図を示す。図 3 に示されるように、Q P のライフサイクル 300 は、リセット状態 301、初期化 (I N T) 状態 302、レディ・ツー・レシーブ (R T R) 状態 303、レディ・ツー・センド (R T S) 状態 304、S Q エラー (S Q E) 状態 305、エラー状態 306、S Q ドレイン (S Q D) 状態 307、および中断 (S U S) 状態 308 といった複数の状態を含む。

【0036】

中断状態である S U S 308 は、マイグレーションの間に Q P を休止するよう進行中の通信を一時的に停止させ得る。R T S 304 と S U S 308 との間の遷移は、ソフトウェアが開始した要求またはイベント要求によってトリガされ得る。マイグレーションの間に、ソフトウェアインターフェイスは、マイグレートする Q P を R T S 304 状態から S U S 308 状態まで移行し得る。Q P は、中断状態に入ると、如何なる以前に開始されたメ

10

20

30

40

50



ッセージ送信も完了し得、S Qにおける如何なる残存するメッセージを、それら进行处理することなく、フラッシュし得る。ハードウェアはさらに、残存するメッセージが、フラッシュする前にユーザプロセスバッファによってキャッシュされることを保証し得る。中断状態において、如何なる新しいメッセージも、S QおよびR Qの両方にキューに入れられ得るが、処理されることなくポストされたままである。

#### 【0037】

VMライブマイグレーションのネットワーキングおよびルーティングの観点

図4は、本発明の実施例に従った、ネットワークおよびルーティングの観点からのバーチャルマシン(VM)のライブマイグレーションの図を示す。図4に示されるように、仮想化環境400は、たとえばvHCA A-C401-403といった複数のvHCAを含み得る。vHCA A-C401-403の各々は、仮想化についてトランスピアレンシーを提供し得る完全なIBエンドポイントとして機能し得る。

10

#### 【0038】

さらに、仮想化環境400は、ネットワークマネージャ410を含み得る。ネットワークマネージャは、IBスイッチA-B431-432の後ろにvHCA A-C401-403を発見し得、異なるLIDおよびvGUIDに基づき、vHCA A-C401-403の各々を認識し得る。たとえば、vHCA A401には、LID A411およびvGUID A421が割り当てられ得、vHCA A402にはLID B412およびvGUID B422が割り当てられ得、vHCA C403には、LID C413およびvGUID C423が割り当てられ得る。

20

#### 【0039】

本発明の実施例に従うと、VMマイグレーションの間、ネットワークマネージャ410は、たとえばvHCA B402のようなIBエンドポイントがIBスイッチA431の後ろでダウン(すなわち中断)し、ネットワークの他方に(vHCA B404として)IBスイッチB432の後ろに、たとえばLID B412およびvGUID B422のような同じアドレスが新しい位置に引き継がれた状態で再び現れることを観察し得る。

#### 【0040】

マイグレーションの後、ネットワークマネージャ410は、そのシフトの後にLIDの新しい再構成を反映するようルーティングテーブル420を更新し得る。さらに、全ネットワークのリルーティングタイムを低減するために、LIDをローカルに修正する新しいルーティングアルゴリズムが実現され得る。さらに、ネットワークマネージャ410はさらに、新しいLID構成はデッドロックがないことを保証し得る。

30

#### 【0041】

さらに、高速ネットワークでは、ローカルアドレスリソースが制限されている場合がある。たとえば、16ビットLIDを使用し得るIBネットワークは、各vHCAに自身のLIDが割り当てられるとともに1000ノードクラスタにおける各ノードが複数のVMをホストする場合、すぐにアドレススペースを使い果たし得る。

#### 【0042】

本発明の実施例に従うと、システムは、仮想化環境におけるスケーラビリティの問題を扱うために、異なるアプローチを使用し得る。たとえば、システムは、32ビットまでLIDアドレッシングスペースを拡張し得るか、または、vHCA同士の間のルーティングを行なうようvGUIDをLIDに組み合わせ得る。第1のソリューションは、ネットワークマネージャにおける如何なるアーキテクチャの修正も必要としないが、より古いハードウェアとの後方互換性がない。

40

#### 【0043】

バーチャルマシンの観点

図5は、本発明の実施例に従った、バーチャルマシン(VM)ライブマイグレーションの前の仮想化環境の図を示す。図5に示されるように、IBネットワーク500は、たとえばホストA-B501-502といった複数のホストを含み得る。ホストA501は、HCA517を使用するVMM A515を含み、たとえばDomU 513といったゲ

50

ストドメインを管理する特権ドメイン（または管理ドメイン）である Dom0 511 をサポートし得る。さらに、ホスト B 502 は、HCA 518 を使用する VMM B 516 を含み、DomU 514 を管理する Dom0 512 をサポートし得る。

【0044】

さらに、ホスト A 501 上の VMa 503 には、たとえば QPa 507 といったキューペアに接続される VF 505 といった VF が付与され得る。さらに、VMa 503 は、QPb 508 に関連付けられる VMb 504 と通信し得る。

【0045】

本発明の実施例に従うと、たとえばホスト A 501 上の Dom0 511 といった特権ドメインは、たとえばマイグレーションの間に分離されるべき VF 505 のようなデバイスを必要としない。システムは単に、たとえば QPa 507 といったアクティブな QP がストップ・アンド・コピーステージで中断されることを必要とし得る。

【0046】

図5に示されるように、ホスト A 501 上の特権ドメインである Dom0 511 は、ホスト B 102 への VMa 503 のマイグレーションを（図6における VMa 604 として）開始し得る。VMa 503 のマイグレーションの前に、システムは、それを中断状態にする状態遷移動作を行ない得る。さらに、QPb 508 から入力パケットを受け取った後、中断された QPa 507 は、QPb 508 に同様に中断状態に入るように命令し得る。

【0047】

その後、たとえばリモートダイレクトメモリアクセス（RDMA）接続 520 に基づき、I/O 仮想メモリを含むダーティページ（メモリ）をソースサーバと宛先サーバとの間で同期させることによって、VMマイグレーションが行なわれ得る。

【0048】

図6は、本発明の実施例に従った、バーチャルマシン（VM）ライブマイグレーションの後の仮想化環境の図を示す。図6に示されるように、IBネットワーク 600 は、たとえばホスト A - B 601 - 602 とした複数のホストを含み得る。（さらに図5と同じ番号を有する図6における要素は同じ機能を果たす）。

【0049】

ホスト B 602 へのマイグレーションの後に、VMa 603 には、たとえば QPa' 607 のようなキューペアに接続されるたとえば VF 616 のような新しい VF が付与され得る。ここで、QPa' 607 は、図5における QPa 507 と論理上同じ QP であり、VMa 503（または VMa 603）の仮想メモリにおいて同じコンテキスト情報を有する。

【0050】

さらに、VMa 603 が停止から解放される前に、特権ドメインである Dom0 612 は、QPa' 607 を RTS 状態に戻し得る、QPa' 607 に対する状態遷移を行ない得る。さらに、QPa' 607 は、ホスト A 601 上の QPb 608 との通信を再開し得る。

【0051】

したがって、マイグレーションについてのサービスダウンタイムは、ストップ・アンド・コピーステージの期間と同等になり得る。さらに、マイグレーションを行なうよう RDMA 動作を使用することは、サービスダウンタイムおよびトータルマイグレーションタイムを低減し得る。

【0052】

図7は、本発明の実施例に従った、仮想化環境におけるバーチャルマシンのライブマイグレーションをサポートするための例示的なフローチャートを示す。図7に示されるように、ステップ 701 では、システムは、各々が別個のキューペア（QP）スペースに関連付けられる複数の仮想機能（VF）に仮想スイッチを関連付け得る。さらに、ステップ 702 では、システムは、少なくとも1つのバーチャルマシンに上記仮想機能（VF）を付

10

20

30

40

50

与し得る。上記仮想機能（VF）は、仮想インターフェイス（VI）に関連付けられる。次いで、ステップ703では、上記の少なくとも1つのバーチャルマシンは、上記仮想機能（VF）が付与された状態で第1のホストから第2のホストにライブマイグレーションを行ない得る。

【0053】

拡張HCAモデル

図8は、本発明の実施例に従った、インフィニバンドアーキテクチャ（IBA）ネットワークデバイスの図を示す。図8に示されるように、インフィニバンドアーキテクチャ（IBA）ネットワークデバイス800は、ハードウェアレイヤー810、バーチャルマシンモニタ820、およびたとえばVM801 - 802といった1つ以上のバーチャルマシンを含み得る。VM801は、カーネル807およびユーザライブラリ837に基づき、固定されたメモリ805に関連付けられるアプリケーション803をサポートし得る。VM802は、カーネルエージェント808およびユーザライブラリ838に基づき、固定されたメモリ806に関連付けられるアプリケーション804をサポートし得る。

【0054】

ハードウェアレイヤー810は、IBA HCA830のような、IBサブネット840に接続するファブリックアダプタを含み得る。さらに、IBA HCA830は、たとえばvHCA 811 - 819といった複数のvHCAインスタンスが関連付けられ得る。

【0055】

本発明の実施例に従うと、たとえば物理HCA830といったファブリックアダプタは、仮想スイッチモデル、拡張HCAモデルまたは両方のモデルの組合せを実現することができる。さらに、拡張HCAモデルを使用して、スイッチなしの物理HCA830のみが存在し得る。バーチャルHCA811 - 819は、物理HCAモデルの拡張を示し、バーチャルHCAインスタンス811 - 819の1つ以上は、物理HCA830ポートを介して、たとえばIBサブネット840といったファブリックからアクセス可能であり得る。

【0056】

本発明の実施例に従うと、ハードウェアレイヤー810は、VMライブマイグレーションを行なうための割り当てられたIBリソースを保存するよう、以下の機能をサポートし得る。

【0057】

1. QPコンテキスト（たとえばQPナンバー）は、各vHCAがプライベートQPスペースを有することを可能にすることによって保存され得る。

【0058】

2. メモリ領域ナンバーは、各vHCAがプライベートメモリ領域ナンバースペースを有することを可能にすることによって保存され得る。

【0059】

3. 登録したメモリは、各vHCAが、ローカルの物理的なホストのローカルメモリへのアクセスのためにプライベート仮想アドレススペースを実現することを可能にすることによって、保存され得る。

【0060】

プライベートQPスペースは、上記のセクションで論じられた仮想スイッチベースのモデルを用いるか、または拡張HCAモデルを用いることのいずれかによって実現され得る。HCA規格を拡張する拡張HCAモデルは、各々がプライベートQPスペースを有し得る複数のvHCAインスタンス811 - 819の定義を可能にする。図8に示されるように、vHCA811 - 819の各々は、たとえばQP821 - 829のようなキューペア（QP）をサポートし得る。当該キューペア（QP）は、パケットを送信および受信するためのたとえばアプリケーション803または804といったアプリケーションに関連付けられるセNDER/レシーババッファによって使用され得る。

【0061】

さらに、拡張HCA規格は、より複雑なファブリックトポロジに依存しないという長所

10

20

30

40

50

を有し得る。さらに、プライベートメモリ領域ナンバースペースは、プライベートQ Pスペースに関する同じ代替例に基づいて、v H C Aごとに実現され得る。

【0062】

本発明の実施例に従うと、ローカルの物理メモリについてのプライベート仮想アドレススペースは、v H C Aスキームが仮想スイッチモデルまたは拡張H C Aモデルに基づくかどうかに関わらず実現され得る。

【0063】

プライベート仮想アドレススペースを実現する1つの方法は、サーバプラットフォームを使用することによる。サーバプラットフォームでは、物理的なH C Aによって観察されたローカルメモリアドレススペースが仮想メモリを示す。この仮想メモリは、具体的にS R I O V V Fごとに、すなわち、v H C A 8 1 1 - 8 1 9ごとにホストプラットフォームによって実現され得る。

10

【0064】

プライベート仮想アドレススペースを実現する別の方法は、たとえばV M M 8 2 0といったプラットフォームハイパーバイザによって制御されるシステムM M Uのセット内の物理的なH C Aのメモリ管理ユニット(M M U)を含むことによる。したがって、如何なる仮想アドレスも変更することなく、仮想アドレスから物理アドレスへのマッピングが必要に応じて更新される状態で、V M固有のアドレススペースが異なる物理的なホスト同士の間でマイグレートされ得る。(たとえばプライベート仮想アドレススペースは、物理的なH C Aとローカルの物理メモリとの間にシステムM M Uを有することによって実現することができる)。

20

【0065】

本発明の実施例に従うと、接続状態は、マイグレートされているV Mのローカルメモリに保存され得る。さらに、接続状態は、マイグレーションの後、新しいv H C Aについて、如何なる新しいリソースも割り当てる必要なしで再確立(または再使用)され得る。

【0066】

さらに、(図3に示されるような)中断状態は、マイグレートされたQ PとマイグレートされたQ Pに接続されるリモートピアQ Pとの両方に適用され得る。したがって、システムは、V Mマイグレーションの間にタイムアウトが発生せず、システムが、マイグレートされたV Mのピアによって実現されるI B接続関連のタイムアウト値と無関係に、マイグレーション動作が任意の量の時間をかけることを可能にし得ることを保証し得る。さらに、接続に関連付けられるI Bアドレスがマイグレーションの結果変化することを可能にするために、システムは、中断状態におけるQ Pについてアドレス情報を更新し得る。

30

【0067】

リモートQ Pピアとの通信

本発明の実施例に従うと、システムは、ピア同士間の通信をサポートするための特別のヘッダ情報および/または非送信請求(non-solicited)メッセージを使用し得る。システムは、リモートピア同士間の休止動作を行なうことなくV Mライブマイグレーションが開始されることを可能にする。さらに、システムは、リモートピアがマイグレーションに関して学習すること、および非同期イベントドリブンの態様で中断状態に入ることを可能にする。

40

【0068】

たとえば、v H C Aルートヘッダ(V R H)が、サポートされ得、I Bネットワークにおけるグローバルルートヘッダ(G R H)ベースのアドレッシングの代わりに、V R Hベースのv H C Aアドレッシングが使用され得る。V R Hを含むパケットは、たとえばV R Hにおける宛先v H C Aタグ(D V T)が0でなく、かつヘッダにおけるD V Tが、V R HにおけるD V Iによって識別されたv H C AタグテーブルエントリにおけるD V Tと同じ場合、V R Hにおいて指定されたように宛先v H C Aインデックス(D V I)に送達され得る。システムは、V R Hの存在を規定するワイヤプロトコル拡張が存在しない場合、L I Dポリシーを介して、どのパケットがV R Hを含むと予想されるかを特定し得る。代替

50

的には、システムは、すべてのvHCAに共通であるQPナンバー範囲を規定し得、異なる範囲がVRHまたは他のプロトコルポリシーの使用に関連付けられ得る。さらに、VRHは、vHCAインデックスがVRHから抽出されるまでQPコンテキストが知られていないので、QPコンテキストを介して特定され得ない。

#### 【0069】

図9は、本発明の実施例に従った、マイグレートされたキューペア(QP)とリモートピアQPとの間の通信をサポートする図を示す。図9に示されるように、仮想化環境900は、各々がたとえばSW911-913のうちの1つのようなソフトウェアスタックを含み得るホストA-C901-903といった複数のホストを含み得る。(ソフトウェアスタックは、カーネルスペース、ユーザスペースまたはその両方に存在し得る)。

10

#### 【0070】

さらに、ホストA901上のQP a910は、QP a910がホストC903に(QP a930として)マイグレートする前に、ホストB902上のQP b920と通信状態にあり得る。QP a910は、その時間の間にQP b920がQP a910にパケット(UD、UCまたはRC)を送信し得るマイグレーションを可能にするよう中断状態にセットされ得る。ホストA901上の無効なvHCA931を目標とするこのパケットは、GRHまたはVRH/LIDのいずれかのマップベースのアドレッシングを使用し得る。

#### 【0071】

無効なvHCA931を目標とするパケットを受け取った後、ホストA901は、ローカルの非同期イベントを生成し得、随意にvHCA932上のQP b920に非送信請求応答パケットを送信し得る。さらにQP b920は、非送信請求応答パケットを受け取ると、当該非送信請求応答から情報を含むローカルの非同期イベントを生成し得る。その後、QP b920は、QP b920が、有効にされた信頼性のある接続された(reliable connected(RC))QPまたは信頼性のない接続(unreliable connection(UC))QPである場合、(図3に示されるように)自動的に中断状態への遷移を行ない得る。

20

#### 【0072】

本発明の実施例に従うと、IBAによって提供されるバーブズインターフェイス(verbs interface)は、キューペア(QP)を中断状態におよび中断状態から移すためにModifyQPおよびModifyAddressHandleオプションを含み得る。さらに、中断状態がRC QPについてリセットされると、再試行カウントおよびタイマーが初期値にリセットされ得る(すなわち中断状態に入る前に消費された如何なる「再試行バジェット(retry budget)」も「忘れられ(forgotten)」得る)。中断状態にある間、パス情報はRCまたはUC

30

QPについて更新され得る。QPおよびアドレスハンドルは、たとえばSW A-C911-913のようなローカルのソフトウェアが明示的に中断状態をリセットするまで、中断状態のままであり得る。ここで、中断状態をリセットする前に、ローカルソフトウェアは、宛先LIDおよびVRH情報を含み得る更新されたパス情報を得ることができる。

#### 【0073】

本発明の実施例に従うと、そのvGUIDによって識別されたvHCAポートがIBネットワークにおいて作動状態になるたびに、vHCAポートは、適切な物理的なHCAポートについて、vGUIDおよびvHCA情報を介して報告され得る。VM/vHCAマイグレーションの場合には、宛先ノードの物理的なHCAポートのvGUIDおよびvHCA情報ならびにSMA情報が、マイグレートされるvHCAが作動状態になると、更新され得る。VMのマイグレーションの後に、中断されたQPまたはアドレスハンドルを有するリモートピアは、関連するパス情報を更新した後、関連情報を観察し得、中断状態をリセットし得る。

40

#### 【0074】

さらに、マイグレートされたVMおよびvHCAが新しい位置にて再び開始される場合、通信ピアは、当該マイグレートされたvHCAを認識していない場合がある。たとえば、QP a930のようなマイグレートされたRC QPがそのピアであるQP b920に要求を送信すると、QP a930は、QP b920が同じ時間フレーム中にマイグレート

50

されなかったので、正しいパス情報を有し得る。さらに、Q P b 9 2 0 が更新されたパス情報を有していなければ、マイグレートされたQ P a 9 3 0 からの入力要求に対する承認および応答が、古い宛先へ送信され得る。その後、当該古い宛先は、Q P b 9 2 0 に非送信請求応答を送信し得、これにより、パス情報が既に更新されたことをQ P b 9 2 0 に認識させる。

【0075】

さらに、たとえばホストC 9 0 3 のような新しい位置にて開始する、v H C A 9 3 3 のようなマイグレートされたv H C A は、送出要求がタイムアウトするのを回避するために、タイムリーな態様でパス情報を更新するようピアに依存し得る。たとえば、Q P b 9 2 0 が正しいパス情報を有していない場合、マイグレートされたv H C A 9 3 3 からの送出要求は、承認および応答が戻ってこないという事実により、タイムアウトし得る。

10

【0076】

マイグレートされたQ P a 9 3 0 からのオリジナルの要求パッケージが、マイグレーションが行われたことをピアQ P b 9 2 0 に通知し得ることを保証するために、システムは、マイグレーション情報を、要求パッケージにおけるヘッダコンテンツから検知できるようにし得る。たとえば、ソースアドレス情報における任意の関連する変更を検知するのをピアに依存する代わりに、V R H は、マイグレーションカウントフィールドを含み得る。マイグレーションカウントフィールドは、マイグレートされたv H C A 9 3 0 が新しい位置であるホストC 9 0 3 で再び作動状態になる前に更新され得る。さらに、ピアQ P 状態に記録されるものに対するマイグレーションカウントにおける変化によって、ピアノードに関する条件を識別する非同期イベントが引き起こされ得る。したがって、V R H マイグレーションカウントフィールドを介して明示的にリモートピアに通知することによって、システムは、ほとんどの場合、マイグレートされたH C A 側上のQ P についてタイムアウトを回避することができる。

20

【0077】

本発明の実施例に従うと、パス情報を更新するのに完全にピアに対してリアルタイムに依存することを回避するために、マイグレートされたv H C A は、非中断ペンディング状態のQ P で開始し得る。非中断ペンディングQ P 状態は、送出要求が可能にされるという例外を有する中断状態に類似する。さらに、v H C A がマイグレーションの後で再び開始された後に生成される要求についての応答または承認がリモートピアから受け取られた後、非中断ペンディング状態はリセットされ得、Q P は正常なものとして挙動し得る。

30

【0078】

さらに、システムの実現例(H W、F WまたはS W)は、タイムアウトおよび再試行を使用してパッケージのロスに対応する一方、Q P 自体のタイムアウト挙動は中断状態におけるものと同じであり得る(すなわち、タイムアウトによりエラー状態にならない)。さらに、非中断ペンディング状態は、マイグレートされないv H C A / Q P について中断状態からの遷移を適切に行なうよう使用され得る。

【0079】

代替的には、V R H が使用されない場合、マイグレーションの信号送信は非送信請求メッセージを介して達成され得る。たとえば、マイグレーションの後のリスタートに続いて、正常なメッセージとの間でインターリーブされた非送信請求メッセージが、場合によっては送られ得る。さらに、リモートピアQ P からの要求または応答の受信により、このような非送信請求メッセージを送信する必要性が除去され得る。ここで、非送信請求メッセージの使用により、Q P / 接続に対する如何なる他の通信アクティビティとも無関係に、信号送信が行われることを可能にする。

40

【0080】

したがって、付加的なヘッダ情報および/または非送信請求メッセージを使用して、システムは、マイグレートされたQ P が、行われるマイグレーションについてそのピアQ P に自動的に信号送信することを可能にし、含まれる新しいアドレス情報をそのピアQ P に提供することを可能にする。アドレス情報の関連する更新は、サブネットマネージャの関

50

連情報の通知および報告と並行して行われ得、接続を完全な作動状態にするためにレイテンシをさらに低減し得る。

【0081】

図10は、本発明の実施例に従った、拡張vHCAモデルに基づいてバーチャルマシンのライブマイグレーションをサポートするための例示的なフローチャートを示す。図10に示されるように、ステップ1001では、システムは、複数のバーチャルホストチャンネルアダプタ(HCA)にファブリックアダプタを関連付け得、上記バーチャルホストチャンネルアダプタ(vHCA)の各々は、別個のキューペア(QP)スペースに関連付けられ得る。ここで、ファブリックアダプタは、仮想スイッチモデル、拡張ホストチャンネルアダプタ(HCA)モデル、または両方のモデルの組合せを実現することができる。さらに、ステップ1002では、少なくとも1つのバーチャルマシンが、第1のホストから第2のホストにライブマイグレーションを行なうよう動作し得、上記少なくとも1つのバーチャルマシンには、上記キューペア(QP)スペースにおけるキューペア(QP)が関連付けられる上記のバーチャルホストチャンネルアダプタ(vHCA)が付与される。その後、ステップ1003では、上記キューペア(QP)は、ライブマイグレーションについてピアQPに信号送信し得、マイグレーションの後に上記ピアQPにアドレス情報を提供し得る。

10

【0082】

スケーラブルアドレッシングスキーム

本発明の実施例に従うと、レイヤー3(L3)アドレスの一部であるVMに関連付けられる仮想GUIDが、ライブマイグレーションを通じて保存され得、仮想GUIDが、ピアに通信する観点と、サブネットマネージャの観点とから一意にVMを識別するよう使用され得る。

20

【0083】

マイグレートされているリモート接続ピアに関連付けられるアドレス情報を更新する本質的な必要性が存在するソフトウェアアプローチと異なり、拡張HCAモデルは、VMに関連付けられるアドレス情報がマイグレーションに亘って保存されることを可能にする。したがって、接続ピアについてアドレス情報を更新する本質的な必要性が存在せず、システムはIBレイヤー2(L2)またはローカルルートヘッダ(LRH)を保存し得るので、L3またはグローバルルートヘッダ(GRH)アドレスに加えてアドレス情報が保存され得る。

30

【0084】

本発明の実施例に従うと、システムは、IBサブネットにおいてスケーラビリティおよびフレキシビリティを達成することができる。たとえば、IBサブネットは、48KのLID値に限定され得る。システムは、関連するHCAポート(物理的なポートまたはvHCAポート)へのパケット転送を促進するよう利用することをIBサブネットにおけるサブネットマネージャが決定するどのようなLID値にでも基づきHCAがアドレス指定されることを可能にするignore LMC機能をサポートする。たとえば、HCAポートは、専用のLID、またはベースとなる2つの隣接するLIDの範囲(HCAポートLIDマスク、すなわちLMC値、によって規定されるような)についての権限を有することを強制され得ない。

40

【0085】

さらに、システムは、物理的なポートについて静的に割り付けられたLID/LID範囲としてサブネットマネージャが使用しないLID値を示す範囲のadmin LIDの使用をサポートし得る。さらに、サブネットマネージャは、VM作成時または動的にVMライフタイム中のいずれかに、アドレッシングのためにユニークなLIDを必要とするVMにadmin LIDを割り当て得る。

【0086】

図11は、本発明の実施例に従った、仮想化環境においてローカル識別子(LID)の割り当てをサポートする図を示す。図11に示されるように、仮想化環境1100は、物

50

理的なH C Aポート1 1 1 1 - 1 1 1 2およびバーチャルH C Aポート1 1 2 1 - 1 1 2 3のようなさまざまなネットワークデバイスについてL I Dを割り当てることを担い得るサブネットマネージャ( S M ) 1 1 0 1を含み得る。

【 0 0 8 7 】

S M 1 1 0 1は、ある範囲のa d m i n L I D 1 1 2 0が保存され得るポリシーインターフェイスをサポートし得る。各a d m i n L I Dは、もしポリシー入力を介して命令が与えられなければ、S M 1 1 0 1がいずれのポートにもこの値を割り当て得ないという点で特別である。

【 0 0 8 8 】

さらに、S M 1 1 0 1は、ポリシーインターフェイスをサポートし得る。当該ポリシーインターフェイスでは、v G U I Dが1つ以上のこのようなa d m i n L I Dに関連付けられ得るとともに、v G U I Dマイグレーションは、対応するL I Dについてのルーティングが、関連する物理的なH C Aポートに対してセットアップされることを示す。さらに、v G U I Dがa d m i n L I Dに関連付けられていれば、S M 1 1 0 1は、i g n o r e L M C機能をサポートしない物理的なH C Aポートにv G U I Dに関連付けることを許可し得ない。

【 0 0 8 9 】

本発明の実施例に従うと、a d m i n L I Dのルーティングは、ルーティングアルゴリズムと、その時のS M 1 0 1についてアクティブなポリシーに基づき得る。a d m i n L I Dに関連付けられるv G U I Dについてのパス情報は、S M 1 0 1によって規定されるポリシーに従って、割り当てられるa d m i n L I Dを反映し得る。

【 0 0 9 0 】

したがって、当該システムは、i g n o r e L M C機能およびa d m i n L I D機能を使用して、I Bサブネット内でのアドレッシングについてスケーラビリティおよびフレキシビリティを、V Mライブマイグレーション動作から独立して提供し得る。一般にこれらの機能は、I Bサブネットの内の任意の物理的またはバーチャルH C Aポートに到達するのに使用され得る1つ以上の付加的なL I Dの動的な生成を促進するよう用いられ得る。

【 0 0 9 1 】

さらに、動的なサービス品質( quality of service ( Q o S ) )、高い可用性( high availability ( H A ) )、および/または他のマルチアドレススペースのプロパティがランタイムにて動的に確立され得る。これらの機能は、サブネット再初期化を必要とすることなく、またはI Bサブネット内での進行中の通信に対して如何なる他のネガティブなインパクトも引き起こすことなく、任意の特定の物理的または仮想のH C Aポートに対して適用され得る。

【 0 0 9 2 】

本発明は、1つ以上のプロセッサ、メモリ、および/または本開示の教示に従ってプログラムされたコンピュータ読取可能な記録媒体を含む1つ以上の従来の汎用または専用デジタルコンピュータ、コンピューティングデバイス、マシン、またはマイクロプロセッサを用いて簡便に実施され得る。ソフトウェア技術の当業者には明らかであるように、適切なソフトウェアコーディングは、熟練したプログラマによって本開示の教示に基づき容易に用意され得る。

【 0 0 9 3 】

いくつかの実施例では、本発明は、本発明の処理のいずれかを実行するようコンピュータをプログラムするのに用いられ得る命令を格納した記録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を含む。当該記録媒体は、フロッピー(登録商標)ディスク、光ディスク、D V D、C D - R O M、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、R O M、R A M、E P R O M、E E P R O M、D R A M、V R A M、フラッシュメモリ素子、磁気または光学カード、ナノシステム(分子メモリI Cを含む)、または命令および/またはデータを格納するのに好適な任意

10

20

30

40

50



のタイプの媒体もしくは装置を含み得るが、これらに限定されない。

【 0 0 9 4 】

本発明の上記の記載は、例示および説明目的で与えられている。網羅的であることまたは開示されたそのものの形態に本発明を限定することを意図したものではない。当業者にとっては、多くの修正例および変形例が明確であろう。上記の実施例は、本発明の原理およびその実際の適用をもっともよく説明するために選択および記載されたものであり、これにより他の当業者が、特定の使用に好適なさまざまな修正例を考慮して、さまざまな実施例について本発明を理解するのが可能になる。本発明の範囲は、添付の特許請求の範囲およびそれらの均等物によって定義されることが意図される。

【 図 1 】

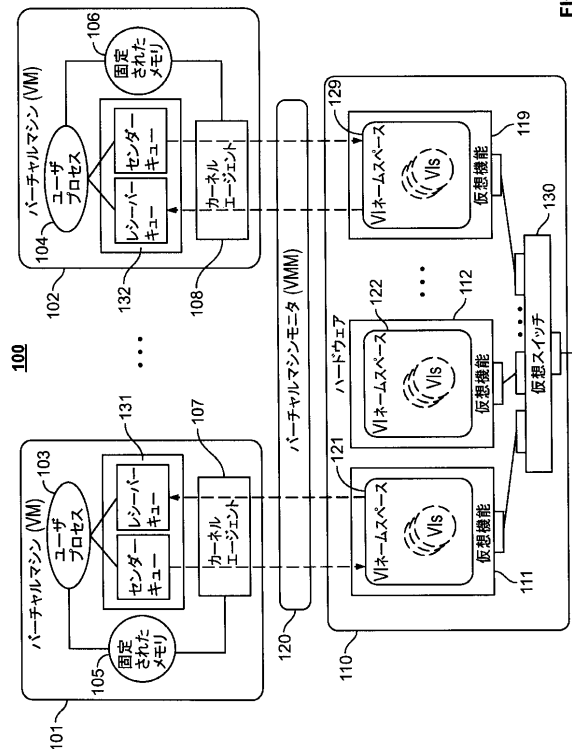


FIG. 1

【 図 2 】

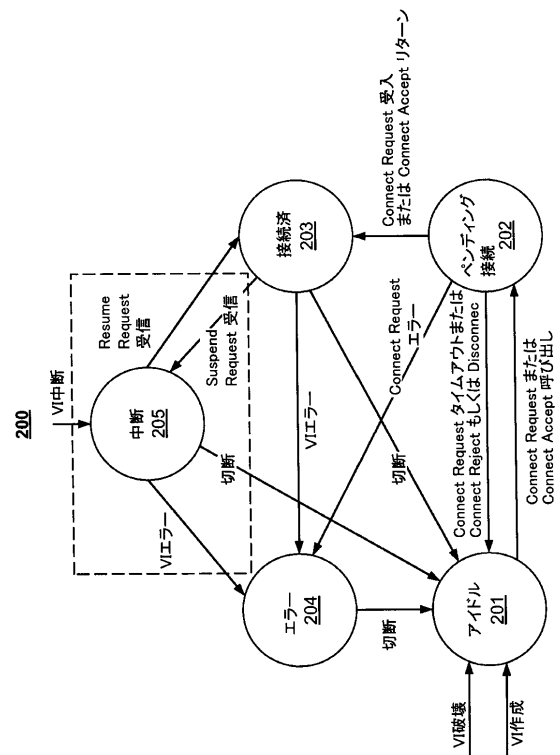


FIG. 2



【図 7】

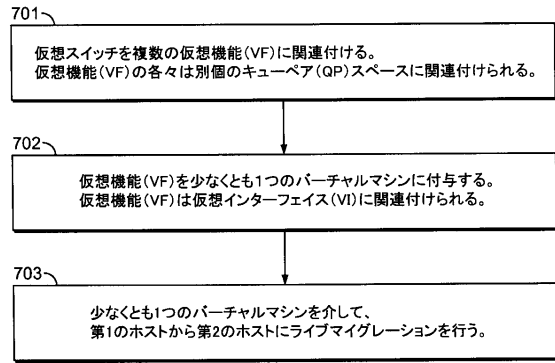


FIG. 7

【図 8】

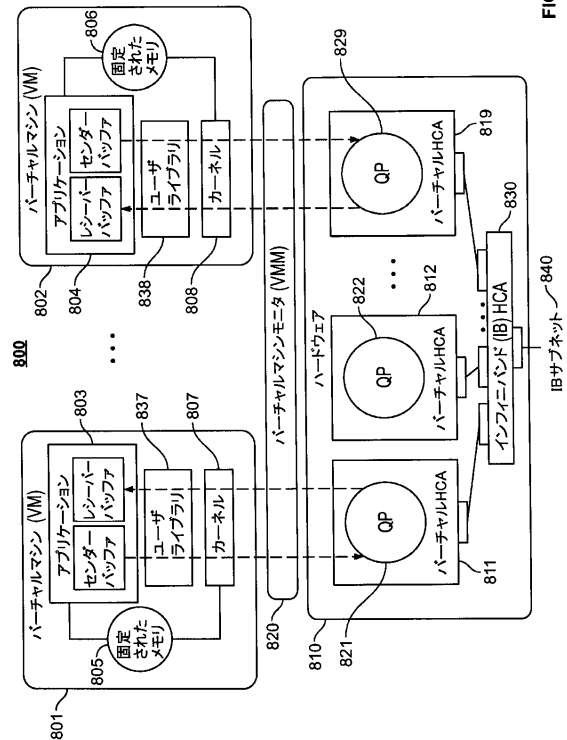


FIG. 8

【図 9】

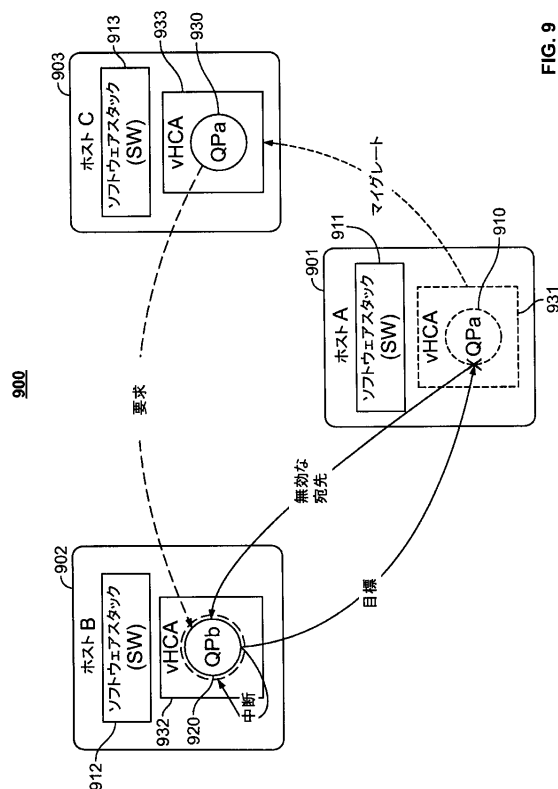


FIG. 9

【図 10】

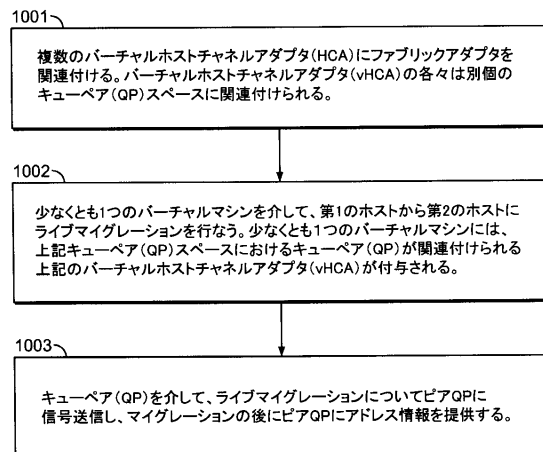


FIG. 10

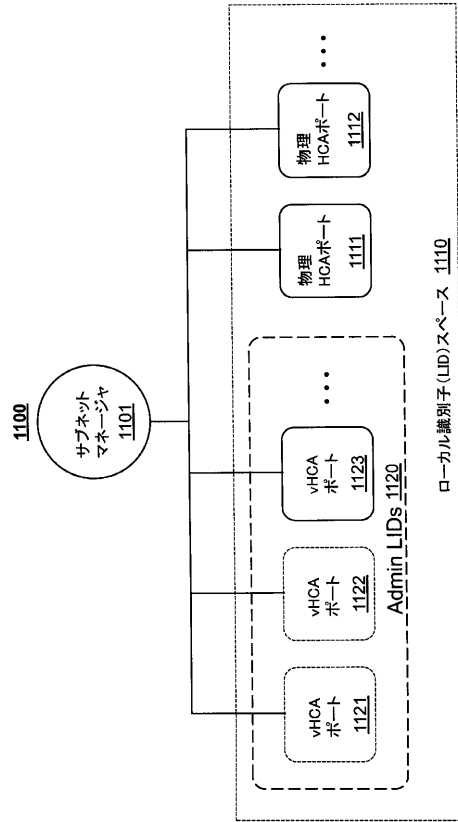


FIG. 11

## フロントページの続き

- (31)優先権主張番号 13/838,275  
(32)優先日 平成25年3月15日(2013.3.15)  
(33)優先権主張国 米国(US)  
(31)優先権主張番号 13/837,922  
(32)優先日 平成25年3月15日(2013.3.15)  
(33)優先権主張国 米国(US)  
(31)優先権主張番号 13/838,121  
(32)優先日 平成25年3月15日(2013.3.15)  
(33)優先権主張国 米国(US)  
(31)優先権主張番号 13/838,502  
(32)優先日 平成25年3月15日(2013.3.15)  
(33)優先権主張国 米国(US)

- (72)発明者 グアイ、ウェイ・リン  
マレーシア、11950 ペナン、バヤン・レパス、ティンカット・ブキット・ジャンブル、1、  
パークビュー・タワー、25-12-1

審査官 多賀 実

- (56)参考文献 特表2011-519089(JP,A)  
特開2011-186967(JP,A)  
米国特許出願公開第2008/0189432(US,A1)  
米国特許出願公開第2012/0042034(US,A1)  
池田 宗広 外4名,「Linuxカーネル Hacks」,株式会社オライリー・ジャパン,  
2011年 7月22日,第1版,pp.225-241

- (58)調査した分野(Int.Cl.,DB名)  
G06F 9/46-9/54  
G06F 13/10