

(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2019-503595
(P2019-503595A)

(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)
HO4L 12/721 (2013.01)	HO4L 12/721 Z	5K030
HO4L 12/931 (2013.01)	HO4L 12/931	5K033
HO4L 12/28 (2006.01)	HO4L 12/28 200Z	

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2018-502700 (P2018-502700)
 (86) (22) 出願日 平成29年1月25日 (2017. 1. 25)
 (85) 翻訳文提出日 平成30年1月19日 (2018. 1. 19)
 (86) 国際出願番号 PCT/US2017/014959
 (87) 国際公開番号 WO2017/132268
 (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, 720
 (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 15/413, 149
 (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)

(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 ホレン, リネ
 ノルウェー、1900 フェツンド、ピタセン、17
 (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

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ける仮想ルータポートにわたる SMP 接続性チェックのためのルータ SMA 抽象化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ける仮想ルータにわたる SMP 接続性チェックをサポートするためのシステムおよび方法を提供する。一実施形態に従うと、SMA モデル拡張は、ローカルルータポートにアドレス指定されたパケット（すなわち、SMP）を送信する可能性を考慮に入れている。パケットがアドレス指定されている SMA はパケットを受取り得るとともに、要求された情報が（たとえば、サブネットにわたる物理リンクによって接続された）リモートノード上にあることを定義する新しい属性を適用し得る。一実施形態に従うと、SMA は、（SMP を受取って別の要求を送信する）プロキシとして動作し得るか、または、SMA は、元のパケットを変更してサブネット間パケットとして送出し得る。

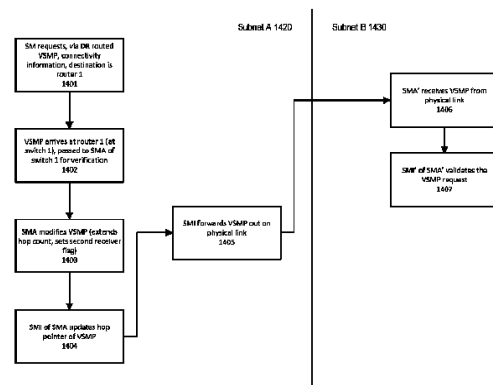


FIGURE 14

【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ける仮想ルータポートにわたる S M P ベースの接続性チェック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、

第 1 のサブネットとを含み、前記第 1 のサブネットは、

複数のスイッチを含み、前記複数のスイッチは少なくともリーフスイッチを含み、前記複数のスイッチの各々は複数のスイッチポートを含み、前記第 1 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、前記複数のホストチャネルアダプタの各々は、少なくとも 1 つのホストチャネルアダプタポートを含み、前記第 1 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少なくとも 1 つのホストチャネルアダプタに関連付けられており、前記第 1 のサブネットはさらに、

サブネットマネージャを含み、前記サブネットマネージャは、前記複数のスイッチおよび前記複数のホストチャネルアダプタのうち 1 つの上で実行しており、

前記複数のスイッチのうち 1 スイッチ上の前記複数のスイッチポートのうち 1 スイッチポートはルータポートとして構成されており、

前記ルータポートとして構成された前記 1 スイッチポートは仮想ルータに論理的に接続されており、

前記サブネットマネージャは、前記ルータポートにアドレス指定された要求パケットを、前記複数のスイッチのうち前記 1 スイッチ上にあるとともに前記ルータポートとして構成された前記複数のスイッチポートのうち前記 1 スイッチポートに送信し、前記要求パケットは、前記ルータポートを越えて接続性情報を要求し、

前記複数のスイッチのうち前記 1 スイッチ上に常駐するサブネット管理エージェントは前記要求パケットを変更する、システム。

【請求項 2】

前記仮想ルータは少なくとも 2 つの仮想ルータポートを含み、

前記ルータポートとして構成された前記スイッチポートは、前記少なくとも 2 つの仮想ルータポートのうち第 1 の仮想ルータポートに論理的に接続される、請求項 1 に記載のシステム。

【請求項 3】

前記サブネット管理エージェントは、アドレスされた前記ルータポートにおいて前記要求パケットが受取られると、前記要求パケットを確認する、請求項 1 または 2 に記載のシステム。

【請求項 4】

第 2 のサブネットをさらに含み、前記第 2 のサブネットは、

前記第 2 のサブネットの複数のスイッチを含み、前記第 2 のサブネットの前記複数のスイッチは、前記第 2 のサブネットの少なくともリーフスイッチを含み、前記第 2 のサブネットの前記複数のスイッチの各々は、前記第 2 のサブネットの複数のスイッチポートを含み、前記第 2 のサブネットはさらに、

前記第 2 のサブネットの複数のホストチャネルアダプタを含み、前記第 2 のサブネットの前記複数のホストチャネルアダプタの各々は、前記第 2 のサブネットの少なくとも 1 つのホストチャネルアダプタポートを含み、前記第 2 のサブネットはさらに、

前記第 2 のサブネットの複数のエンドノードを含み、前記第 2 のサブネットの前記複数のエンドノードの各々は、前記第 2 のサブネットの前記複数のホストチャネルアダプタのうち前記第 2 のサブネットの少なくとも 1 つのホストチャネルアダプタに関連付けられており、前記第 2 のサブネットはさらに、

前記第 2 のサブネットのサブネットマネージャを含み、前記第 2 のサブネットの前記

10

20

30

40

50

サブネットマネージャは、前記第 2 のサブネットの前記複数のスイッチおよび前記第 2 のサブネットの前記複数のホストチャネルアダプタのうち 1 つの上で実行しており、

前記第 2 のサブネットの別の複数のスイッチのうち 1 スイッチ上における前記第 2 のサブネットの前記複数のスイッチポートのうち前記第 2 のサブネットの 1 スイッチポートは、前記第 2 のサブネットのルータポートとして構成されており、

前記第 2 のサブネットの前記ルータポートとして構成された前記第 2 のサブネットの前記 1 スイッチポートは前記第 2 のサブネットの仮想ルータに論理的に接続されており、前記第 2 のサブネットの前記仮想ルータは、前記第 2 のサブネットの少なくとも 2 つの仮想ルータポートを含み、

前記第 1 のサブネットは、物理リンクを介して前記第 2 のサブネットと相互接続されている、請求項 1 から 3 のいずれか 1 項に記載のシステム。 10

【請求項 5】

前記要求パケットは、DR (ダイレクティブルーティング) ルーティングされたパケットであり、

前記サブネット管理エージェントによって実行された変更は、

1 つ以上のホップによって、前記要求に関連付けられたホップ・カウンタを拡張することと、

リモートエンド宛先を示す受取りフラグを設定することを含む、請求項 1 から 4 のいずれか 1 項に記載のシステム。 20

【請求項 6】

前記サブネット管理エージェントに関連付けられたサブネット管理インターフェイス (SMI) はさらに前記要求パケットを変更し、前記 SMI による前記変更はホップポイント属性を更新することを含み、

前記 SMI は、前記物理リンクを介して前記第 2 のサブネットの前記ルータポートに前記要求パケットを転送する、請求項 4 または 5 に記載のシステム。

【請求項 7】

前記要求パケットは LID ルーティングされたパケットであり、

前記サブネット管理エージェントによって実行された前記変更は、

1 つ以上のホップによって、前記要求に関連付けられたホップ・カウンタを拡張することと、 30

リモートエンド宛先を示す受取りフラグを設定することと、

前記 LID ルーティングされたパケットのアドレスに 1 ホップ DR (ダイレクティブルーティング) ルーティングされた経路を追加することを含む、請求項 1 から 4 のいずれか 1 項に記載のシステム。

【請求項 8】

前記サブネット管理エージェントに関連付けられたサブネット管理インターフェイス (SMI) はさらに前記要求パケットを変更し、前記 SMI による前記変更はホップポイント属性を更新することを含み、

前記 SMI は、前記物理リンクを介して前記第 2 のサブネットの前記ルータポートに前記要求パケットを転送する、請求項 4 または 5 に記載のシステム。 40

【請求項 9】

高性能コンピューティング環境における仮想ルータポートにわたる SMP ベースの接続性チェックをサポートするための方法であって、

1 つ以上のマイクロプロセッサを含む 1 つ以上のコンピュータにおいて、第 1 のサブネットを提供するステップを含み、前記第 1 のサブネットは、

複数のスイッチを含み、前記複数のスイッチは少なくともリーフスイッチを含み、前記複数のスイッチの各々は複数のスイッチポートを含み、前記第 1 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、前記複数のホストチャネルアダプタの各々は少なくとも 1 つのホストチャネルアダプタポートを含み、前記第 1 のサブネットはさらに 50

、
 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少なくとも1つのホストチャネルアダプタに関連付けられており、前記第1のサブネットはさらに、

サブネットマネージャを含み、前記サブネットマネージャは、前記複数のスイッチおよび前記複数のホストチャネルアダプタのうちの1つの上で実行しており、前記方法はさらに、

前記複数のスイッチのうち1スイッチ上において前記複数のスイッチポートのうち1スイッチポートをルータポートとして構成するステップと、

前記ルータポートとして構成された前記1スイッチポートを仮想ルータに論理的に接続するステップとを含み、前記仮想ルータは少なくとも2つの仮想ルータポートを含み、前記方法はさらに、

前記ルータポートにアドレス指定された要求パケットを、前記サブネットマネージャにより、前記複数のスイッチのうち前記1スイッチ上にあり前記ルータポートとして構成された前記複数のスイッチポートのうち前記1スイッチポートに送信するステップを含み、前記要求パケットは、前記ルータポートを越えて接続性情報を要求し、前記方法はさらに

、
 前記複数のスイッチのうち前記1スイッチ上に常駐するサブネット管理エージェントによって、前記要求パケットを変更するステップを含む、方法。

【請求項10】

前記仮想ルータは少なくとも2つの仮想ルータポートを含み、

前記仮想ルータは少なくとも2つの仮想ルータポートを含み、

前記ルータポートとして構成された前記スイッチポートは、前記少なくとも2つの仮想ルータポートのうち第1の仮想ルータポートに論理的に接続される、請求項9に記載の方法。

【請求項11】

前記サブネット管理エージェントは、アドレスされた前記ルータポートにおいて前記要求パケットが受取られると、前記要求パケットを確認する、請求項9または10に記載の方法。

【請求項12】

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、第2のサブネットを提供するステップをさらに含み、前記第2のサブネットは、

前記第2のサブネットの複数のスイッチを含み、前記第2のサブネットの前記複数のスイッチは、前記第2のサブネットの少なくともリーフスイッチを含み、前記第2のサブネットの前記複数のスイッチの各々は、前記第2のサブネットの複数のスイッチポートを含み、前記第2のサブネットはさらに、

前記第2のサブネットの複数のホストチャネルアダプタを含み、前記第2のサブネットの前記複数のホストチャネルアダプタの各々は、前記第2のサブネットの少なくとも1つのホストチャネルアダプタポートを含み、前記第2のサブネットはさらに、

前記第2のサブネットの複数のエンドノードを含み、前記第2のサブネットの前記複数のエンドノードの各々は、前記第2のサブネットの前記複数のホストチャネルアダプタのうち前記第2のサブネットの少なくとも1つのホストチャネルアダプタに関連付けられており、前記第2のサブネットはさらに、

前記第2のサブネットのサブネットマネージャを含み、前記第2のサブネットの前記サブネットマネージャは、前記第2のサブネットの前記複数のスイッチおよび前記第2のサブネットの前記複数のホストチャネルアダプタのうちの1つの上で実行しており、前記方法はさらに、

前記第2のサブネットの別の複数のスイッチのうちの1スイッチ上で、前記第2のサブネットの前記複数のスイッチポートのうち前記第2のサブネットの1スイッチポートを前記第2のサブネットのルータポートとして構成するステップを含み、

10

20

30

40

50

前記第 2 のサブネットの前記ルータポートとして構成された前記第 2 のサブネットの前記 1 スイッチポートは、前記第 2 のサブネットの仮想ルータに論理的に接続されており、前記第 2 のサブネットの前記仮想ルータは、前記第 2 のサブネットの少なくとも 2 つの仮想ルータポートを含み、

前記第 1 のサブネットは、物理リンクを介して前記第 2 のサブネットに相互接続されている、請求項 9 から 11 のいずれか 1 項に記載の方法。

【請求項 13】

前記要求パケットは、DR (ダイレクティブルーティング) ルーティングされたパケットであり、

前記サブネット管理エージェントによって実行される変更は、

1 つ以上のホップによって前記要求に関連付けられたホップ・カウンタを拡張することと、

リモートエンド宛先を示す受取りフラグを設定することを含む、請求項 9 から 12 のいずれか 1 項に記載の方法。

【請求項 14】

前記サブネット管理エージェントに関連付けられたサブネット管理インターフェイス (SMI) はさらに、前記要求パケットを変更し、前記 SMI による前記変更はホップポイント属性を更新することを含み、

前記 SMI は、前記物理リンクを介して、前記要求パケットを前記第 2 のサブネットの前記ルータポートに転送する、請求項 12 または 13 に記載の方法。

【請求項 15】

前記要求パケットは LID ルーティングされたパケットであり、

前記サブネット管理エージェントによって実行された前記変更は、

1 つ以上のホップによって前記要求に関連付けられたホップ・カウンタを拡張することと、

リモートエンド宛先を示す受取りフラグを設定することと、

前記 LID ルーティングされたパケットのアドレスに 1 ホップ DR ルーティングされた経路を追加することを含む、請求項 9 から 12 のいずれか 1 項に記載の方法。

【請求項 16】

前記サブネット管理エージェントに関連付けられたサブネット管理インターフェイス (SMI) はさらに、前記要求パケットを変更し、前記 SMI による前記変更はホップポイント属性を更新することを含み、

前記 SMI は、前記物理リンクを介して前記第 2 のサブネットの前記ルータポートに前記要求パケットを転送する、請求項 12 または 13 に記載の方法。

【請求項 17】

高性能コンピューティング環境における仮想ルータポートにわたる SMP ベースの接続性チェックをサポートするための命令が格納されている非一時的なコンピュータ読取り可能記憶媒体であって、前記命令が、1 つ以上のコンピュータによって読出されて実行されると、前記 1 つ以上のコンピュータに以下のステップを実行させ、前記以下のステップは

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、

第 1 のサブネットを提供するステップを含み、前記第 1 のサブネットは、

複数のスイッチを含み、前記複数のスイッチは少なくともリーフスイッチを含み、前記複数のスイッチの各々は複数のスイッチポートを含み、前記第 1 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、前記複数のホストチャネルアダプタの各々は、少なくとも 1 つのホストチャネルアダプタポートを含み、前記第 1 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少なくとも 1 つのホストチャネルアダプタに関連付けられており、

10

20

30

40

50

前記第 1 のサブネットはさらに、

サブネットマネージャを含み、前記サブネットマネージャは、前記複数のスイッチおよび前記複数のホストチャンネルアダプタのうちの一つの上で実行しており、前記以下のステップはさらに、

前記複数のスイッチのうち 1 スイッチ上に、前記複数のスイッチポートのうち 1 スイッチポートをルータポートとして構成するステップと、

前記ルータポートとして構成された前記 1 スイッチポートを仮想ルータに論理的に接続するステップとを含み、前記仮想ルータは少なくとも 2 つの仮想ルータポートを含み、前記以下のステップはさらに、

前記ルータポートにアドレス指定された要求パケットを、前記サブネットマネージャによって、前記複数のスイッチのうち前記 1 スイッチ上にあり前記ルータポートとして構成された前記複数のスイッチポートのうち前記 1 スイッチポートに送信するステップを含み、前記要求パケットは、前記ルータポートを越えて接続性情報を要求し、前記以下のステップはさらに、

前記複数のスイッチのうち前記 1 スイッチ上に常駐するサブネット管理エージェントによって、前記要求パケットを変更するステップを含む、非一時的なコンピュータ読取り可能記憶媒体。

【請求項 18】

前記仮想ルータは少なくとも 2 つの仮想ルータポートを含み、

前記仮想ルータは少なくとも 2 つの仮想ルータポートを含み、

前記ルータポートとして構成された前記スイッチポートは、前記少なくとも 2 つの仮想ルータポートのうち第 1 の仮想ルータポートに論理的に接続される、請求項 17 に記載の非一時的なコンピュータ読取り可能記憶媒体。

【請求項 19】

前記サブネット管理エージェントは、アドレス指定された前記ルータポートにおいて前記要求パケットが受取られると、前記要求パケットを確認する、請求項 17 または 18 に記載の非一時的なコンピュータ読取り可能記憶媒体。

【請求項 20】

前記以下のステップはさらに、

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第 2 のサブネットをさらに提供するステップを含み、前記第 2 のサブネットは、

前記第 2 のサブネットの複数のスイッチを含み、前記第 2 のサブネットの前記複数のスイッチは、前記第 2 のサブネットの少なくともリーフスイッチを含み、前記第 2 のサブネットの前記複数のスイッチの各々は、前記第 2 のサブネットの複数のスイッチポートを含み、前記第 2 のサブネットはさらに、

前記第 2 のサブネットの複数のホストチャンネルアダプタを含み、前記第 2 のサブネットの前記複数のホストチャンネルアダプタの各々は、前記第 2 のサブネットの少なくとも 1 つのホストチャンネルアダプタポートを含み、前記第 2 のサブネットはさらに、

前記第 2 のサブネットの複数のエンドノードを含み、前記第 2 のサブネットの前記複数のエンドノードの各々は、前記第 2 のサブネットの前記複数のホストチャンネルアダプタのうち前記第 2 のサブネットの少なくとも 1 つのホストチャンネルアダプタに関連付けられており、前記第 2 のサブネットはさらに、

前記第 2 のサブネットのサブネットマネージャを含み、前記第 2 のサブネットの前記サブネットマネージャは、前記第 2 のサブネットの前記複数のスイッチおよび前記第 2 のサブネットの前記複数のホストチャンネルアダプタのうちの一つの上で実行しており、前記以下のステップはさらに、

前記第 2 のサブネットの別の複数のスイッチうち 1 スイッチ上において、前記第 2 のサブネットの前記複数のスイッチポートのうち前記第 2 のサブネットの 1 スイッチポートを前記第 2 のサブネットのルータポートとして構成するステップを含み、

前記第 2 のサブネットの前記ルータポートとして構成された前記第 2 のサブネットの前

10

20

30

40

50

記 1 スイッチポートは前記第 2 のサブネットの仮想ルータに論理的に接続されており、前記第 2 のサブネットの前記仮想ルータは、前記第 2 のサブネットの少なくとも 2 つの仮想ルータポートを含み、

前記第 1 のサブネットは、物理リンクを介して前記第 2 のサブネットに相互接続されている、請求項 17 から 19 のいずれか 1 項に記載の非一時的なコンピュータ読取り可能記憶媒体。

【請求項 21】

機械読取り可能フォーマットのプログラム命令を含むコンピュータプログラムであって、前記プログラム命令が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 9 から 16 のいずれか 1 項に記載の方法を実行させる、コンピュータプログラム。

10

【請求項 22】

非一時的な機械読取り可能データ記憶媒体に格納された、請求項 21 に記載のコンピュータプログラムを含むコンピュータプログラムプロダクト。

【発明の詳細な説明】

【技術分野】

【0001】

著作権表示：

この特許文献の開示の一部は、著作権保護の対象となる資料を含む。この特許文献または特許開示は特許商標庁の特許ファイルまたは記録に記載されているため、著作権保有者は、何人によるその複製複製に対しても異議はないが、その他の場合には如何なるときもすべての著作権を保有する。

20

【0002】

発明の分野：

本発明は、概して、コンピュータシステムに関し、特に、高性能コンピューティング環境におけるデュアルポート仮想ルータポートにわたる S M P 接続性チェックのための S M A 抽象化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。

30

【発明の概要】

【課題を解決するための手段】

【0004】

概要：

高性能コンピューティング環境における仮想ルータポートにわたる S M P 接続性チェックをサポートするためのシステムおよび方法がここに記載される。例示的な方法は、1 つ以上のマイクロプロセッサを含む 1 つ以上のコンピュータにおいて、第 1 のサブネットを提供し得る。第 1 のサブネットは、複数のスイッチを含み、複数のスイッチは少なくともリーフスイッチを含む。複数のスイッチの各々は複数のスイッチポートを含む。第 1 のサブネットは、複数のホストチャネルアダプタを含み、各々の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む。第 1 のサブネットは、複数のエンドノードを含む。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ている。第 1 のサブネットはさらに、サブネットマネージャを含む。サブネットマネージャは、複数のスイッチおよび複

40

50

数のホストチャネルアダプタのうちの1つの上で実行している。当該方法は、複数のスイッチのうち1スイッチ上において複数のスイッチポートのうち1スイッチポートをルータポートとして構成し得る。当該方法は、ルータポートとして構成されたスイッチポートを仮想ルータに論理的に接続し得る。仮想ルータは少なくとも2つの仮想ルータポートを含む。当該方法は、ルータポートにアドレス指定されたリクエストパケットを、サブネットマネージャによって、ルータポートとして構成された複数のスイッチのうち1スイッチ上における複数のスイッチポートのうち1スイッチポートに送信し得る。要求パケットは、ルータポートを越えて接続性情報を要求する。最後に、当該方法は、複数のスイッチのうち1スイッチ上に常駐するサブネット管理エージェントによって、要求パケットを変更し得る。

10

【0005】

一実施形態に従うと、(第1のサブネットまたは第2のサブネットの)複数のホストチャネルアダプタのうち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得る。(第1のサブネットまたは第2のサブネットの)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合せを含み得る。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ている。

【図面の簡単な説明】

【0006】

【図1】一実施形態に従ったインフィニバンド環境の一例を示す図である。

20

【図2】一実施形態に従った、パーティショニングされたクラスタ環境の一例を示す図である。

【図3】一実施形態に従った、ネットワーク環境におけるツリートポロジの一例を示す図である。

【図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が予めポピュレートされた例示的なvSwitchアーキテクチャを示す図である。

【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。

30

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

【図10】一実施形態に従った、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である。

【図11】一実施形態に従った、高性能コンピューティング環境における2つのサブネット間の相互接続を示す図である。

【図12】一実施形態に従った、高性能コンピューティング環境におけるデュアルポート仮想ルータ構成を介する2つのサブネット間の相互接続を示す図である。

40

【図13】一実施形態に従った、高性能コンピューティング環境におけるデュアルポート仮想ルータをサポートするための方法を示すフローチャートである。

【図14】一実施形態に従った、DRルーティングされたVSMPパケットを示すフローチャートである。

【図15】一実施形態に従った、DRルーティングされたVSMPパケットに対する応答を示すフローチャートである。

【図16】一実施形態に従った、LIDルーティングされたVSMPパケットを示すフローチャートである。

【図17】一実施形態に従った、LIDルーティングされたVSMPパケットに対する応答を示すフローチャートである。

50

【発明を実施するための形態】

【0007】

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。

10

【0008】

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。

【0009】

高性能コンピューティング環境における仮想ルータポートにわたるSMP接続性チェックをサポートするためのシステムおよび方法がこの明細書中に記載される。

【0010】

この発明の以下の説明は、高性能ネットワークについての一例として、インフィニバンド（IB）ネットワークを使用する。以下の記載全体にわたり、インフィニバンド（InfiniBand™）規格（インフィニバンド規格、IB規格またはレガシーIB規格ともさまざまに称される）を参照することができる。このような参照は、全体がこの明細書中に引用により援用されている、<http://www.infinibandta.org>から入手可能な、2015年3月に公開されたインフィニバンドトレードアーキテクチャ規格；第1号；バージョン1.3を参照するものと理解される。他のタイプの高性能ネットワークが何ら限定されることなく使用され得ることが、当業者には明らかであるだろう。以下の説明ではまた、ファブリックトポロジーについての一例として、ファットツリートポロジーを使用する。他のタイプのファブリックトポロジーが何ら限定されることなく使用され得ることが当業者には明らかであるだろう。

20

【0011】

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行うことができるようにするためにネットワーク再構成時間が重要となるので、ネットワークアーキテクチャに加えて、スケーラブルであるとともにトポロジーに依存しない動的な再構成メカニズムを提供することができる。

30

40

【0012】

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

【0013】

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

50

なり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。

【 0 0 1 4 】

インフィニバンド (登録商標)

インフィニバンド (IB) は、インフィニバンド・トレード・アソシエーション (InfiniBand™ Trade Association) によって開発されたオープン標準無損失ネットワーク技術である。この技術は、特に高性能コンピューティング (high-performance computing : HPC) アプリケーションおよびデータセンタを対象とする、高スループットおよび少ない待ち時間の通信を提供するシリアルポイントツーポイント全二重相互接続 (serial point-to-point full-duplex interconnect) に基づいている。

【 0 0 1 5 】

インフィニバンド・アーキテクチャ (InfiniBand Architecture : IBA) は、2層トポロジー分割をサポートする。低層では、IBネットワークはサブネットと呼ばれ、1つのサブネットは、スイッチおよびポイントツーポイントリンクを使用して相互接続される一組のホストを含み得る。より高いレベルでは、1つのIBファブリックは、ルータを使用して相互接続され得る1つ以上のサブネットを構成する。

【 0 0 1 6 】

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

【 0 0 1 7 】

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

【 0 0 1 8 】

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

【 0 0 1 9 】

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

10

20

30

40

50

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決められる。マスターサブネットマネージャはまた、サブネットの周期的なスイープ (sweep) を行なってあらゆるトポロジー変化を検出し、それに応じてネットワークを再構成する。

【 0 0 2 0 】

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 (g l o b a l i d e n t i f i e r : G I D) を形成するために使用され得る。

10

【 0 0 2 1 】

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

20

【 0 0 2 2 】

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ーが変化する場合に、ネットワークを再構成し得る。再構成ステップは、ネットワーク初期化中に行なわれるステップを含み得る。さらに、再構成は、ネットワーク変化が生じたサブネットに制限されるローカルスコープを有し得る。また、ルータを用いる大規模ファブリックのセグメント化は、再構成スコープを制限し得る。

【 0 0 2 3 】

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、インフィニバンドファブリック 1 2 0 を使用して、それぞれのホストチャネルアダプタ 1 1 1 ~ 1 1 5 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 1 0 1 ~ E 1 0 5) 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 1 0 1 ~ E 1 0 5) 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。

30

【 0 0 2 4 】

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、I B ネットワークは、ネットワークファブリックを共有するシステムの論理グループの分離をもたらすためにセキュリティメカニズムとしてパーティショニングをサポートし得る。ファブリックにおけるノード上の各 H C A ポートは、1 つ以上のパーティションのメンバであり得る。パーティションメンバーシップは、S M の一部であり得る集中型パーティションマネージャによって管理される。S M は、各ポートに関するパーティションメンバーシップ情報を、1 6 ビットのパーティションキー (partition key : P _ K e y) のテーブルとして構成することができる。S M はまた、これらのポートを介してデータトラフィックを送信または受信するエンドノードに関連付けられた P _ K e y 情報を含むパーティション実施テーブルを用いて、スイッチポートおよびルータポートを構成することができる。加えて、一般的な場合には、スイッチポートのパーティションメンバーシップは、(リンクに向かう) 出口方向に向かってポートを介してルーティングされた L I D に間接的に関連付けられたすべてのメンバーシップの集合を表わし得る。

40

【 0 0 2 5 】

50

一実施形態に従うと、パーティションは、1グループのメンバが単に同じ論理グループの他のメンバと通信することができる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 (host channel adapter : H C A) およびスイッチにおいて、分離を実施するためにパーティションメンバーシップ情報を用いてパケットをフィルタリングすることができる。無効なパーティショニング情報を備えたパケットは、パケットが受信ポートに到達すると直ちに削除され得る。パーティショニングされた I B システムにおいては、パーティションはテナントクラスタを作成するために用いることができる。パーティションが適所で実施されていれば、ノードは、異なるテナントクラスタに属する他のノードと通信することができない。このようにして、安全性が損なわれたテナントノードまたは悪意あるテナントノードが存在する場合であっても、システムのセキュリティを保証することができる。

10

【 0 0 2 6 】

一実施形態に従うと、ノード間における通信のために、管理キュー対 ((Queue Pair : Q P) Q P 0 および Q P 1) を除いて、キュー対 (Queue Pair : Q P) およびエンドツーエンドコンテキスト (End-to-End context : E E C) を特定のパーティションに割り当てることができる。次いで、 P _ K e y 情報は、送信されたすべての I B トランスポートパケットに追加することができる。パケットが H C A ポートまたはスイッチに到達すると、その P _ K e y 値は S M によって構成されたテーブルと突き合わせて確認することができる。無効な P _ K e y 値が発見された場合、パケットは直接廃棄される。このようにして、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。

20

【 0 0 2 7 】

一実施形態に従った、パーティショニングされたクラスタ環境の一例を示す I B パーティションの例が図 2 に示される。図 2 に示される例において、ノード A 1 0 1 ~ E 1 0 5 は、インフィニバンドファブリック 1 2 0 を用いて、それぞれのホストチャネルアダプタ 1 1 1 ~ 1 1 5 を介して通信する。ノード A ~ E は、複数のパーティション、すなわちパーティション 1 1 3 0、パーティション 2 1 4 0 およびパーティション 3 1 5 0 に配置される。パーティション 1 はノード A 1 0 1 およびノード D 1 0 4 を含む。パーティション 2 はノード A 1 0 1、ノード B 1 0 2 およびノード C 1 0 3 を含む。パーティション 3 はノード C 1 0 3 およびノード E 1 0 5 を含む。パーティションの配置のせいで、ノード D 1 0 4 およびノード E 1 0 5 は、これらのノードがパーティションを共有していないために通信することが許可されない。一方で、たとえば、ノード A 1 0 1 およびノード C 1 0 3 は、これらのノードがともにパーティション 2 1 4 0 のメンバであるので、通信することが許可される。

30

【 0 0 2 8 】**インフィニバンドにおける仮想マシン**

過去 1 0 年の間に、ハードウェア仮想化サポートによって C P U オーバーヘッドが実質的に排除され、メモリ管理ユニットを仮想化することによってメモリオーバーヘッドが著しく削減され、高速 S A N ストレージまたは分散型ネットワークファイルシステムの利用によってストレージオーバーヘッドが削減され、シングルルート I / O 仮想化 (Single Root Input/Output Virtualization : S R - I O V) のようなデバイス・パススルー技術を使用することによってネットワーク I / O オーバーヘッドが削減されてきたことに応じて、仮想化された高性能コンピューティング (High Performance Computing : H P C) 環境の将来見通しが大幅に改善されてきた。現在では、クラウドが、高性能相互接続ソリューションを用いて仮想 H P C (virtual H P C : v H P C) クラスタに対応し、必要な性能を提供することができる。

40

【 0 0 2 9 】

しかしながら、インフィニバンド (I B) などの無損失ネットワークと連結されたとき、仮想マシン (V M) のライブマイグレーションなどのいくつかのクラウド機能は、これらのソリューションにおいて用いられる複雑なアドレス指定およびルーティングスキームのせいで、依然として問題となる。 I B は、高帯域および低レイテンシを提供する相互接

50

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。

【0030】

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ホストチャンネルアダプタ(HCA)に割当てられたVMのライブマイグレーションを実現することは難易度の高いものであることが判明した。各々のIBが接続されているノードは、3つの異なるアドレス(すなわちLID、GUIDおよびGID)を有する。ライブマイグレーションが発生すると、これらのアドレスのうち1つ以上が変化する。マイグレーション中のVM(VM-in-migration)と通信する他のノードは接続性を失う可能性がある。これが発生すると、IBサブネットマネージャ(Subnet Manager: SM)にサブネット管理(Subnet Administration: SA)経路記録クエリを送信することによって、再接続すべき仮想マシンの新しいアドレスを突きとめることにより、失われた接続を回復させるように試みることができる。

10

【0031】

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。

20

30

【0032】

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

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

40

【0033】

図3は、一実施形態に従った、ネットワーク環境におけるツリートポロジーの例を示す。図3に示すように、ネットワークファブリック200において、1つ以上のエンドノード201~204が接続され得る。ネットワークファブリック200は、複数のリーフスイッチ211~214と複数のスパインスイッチまたはルート(root)スイッチ231~234とを含むファットツリートポロジーに基づき得る。加えて、ネットワークファブリック200は、スイッチ221~224などの1つ以上の中間スイッチを含み得る。

【0034】

また、図3に示すように、エンドノード201~204の各々は、マルチホームノード、すなわち、複数のポートを介してネットワークファブリック200のうち2つ以上の部

50

分に接続される単一のノード、であり得る。たとえば、ノード 201 はポート H1 および H2 を含み、ノード 202 はポート H3 および H4 を含み、ノード 203 はポート H5 および H6 を含み、ノード 204 はポート H7 および H8 を含み得る。

【0035】

加えて、各スイッチは複数のスイッチポートを有し得る。たとえば、ルートスイッチ 231 はスイッチポート 1 ~ 2 を有し、ルートスイッチ 232 はスイッチポート 3 ~ 4 を有し、ルートスイッチ 233 はスイッチポート 5 ~ 6 を有し、ルートスイッチ 234 はスイッチポート 7 ~ 8 を有し得る。

【0036】

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、IB 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、OFED (Open Fabric Enterprise Distribution: IB 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 OpenSM において実現される。

10

【0037】

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 LFT 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 LID、ひいては対応するルートを各スイッチポートに割当てる。同じリーフスイッチに接続されたエンドノードについては、索引付け順序は、エンドノードが接続されるスイッチポートに依存し得る (すなわち、ポートナンバリングシーケンス)。各ポートについては、メカニズムはポート使用カウンタを維持することができ、新しいルートが追加されるたびに、ポート使用カウンタを使用して使用頻度が最小のポートを選択することができる。

20

【0038】

一実施形態に従うと、パーティショニングされたサブネットでは、共通のパーティションのメンバではないノードは通信することを許可されない。実際には、これは、ファットツリールーティングアルゴリズムによって割当てられたルートのうちのいくつかはユーザトラフィックのために使用されないことを意味する。ファットツリールーティングメカニズムが、それらのルートについての LFT を、他の機能的経路と同じやり方で生成する場合、問題が生じる。この動作は、リンク上でバランシングを劣化させるおそれがある。なぜなら、ノードが索引付けの順序でルーティングされているからである。パーティションに気づかずにルーティングが行なわれ得るため、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。

30

【0039】

一実施形態に従うと、ファットツリーは、利用可能なネットワークリソースでスケールリングすることができる階層ネットワークトポロジーである。さらに、ファットツリーは、さまざまなレベルの階層に配置された商品スイッチを用いて容易に構築される。さらに、 k -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。

40

【0040】

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 トポロジーを拡張して、スイッチ間の複数の接続を可能にする。多種多様なトポロジーは 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

50

ットツリーを定義するために、P G F Tの制限バージョンであるR L F Tが導入されている。R L F Tは、ファットツリーにおけるすべてのレベルに同じポートカウントスイッチを用いている。

【0041】

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 : I O V) は、基礎をなす物理リソースに仮想マシン (V M) がアクセスすることを可能にすることによって、I / O を利用可能にすることができる。ストレージトラフィックとサーバ間通信とを組合せると、シングルサーバのI / Oリソースにとって抗し難い高い負荷が課され、結果として、データの待機中に、バックログが発生し、プロセッサがアイドル状態になる可能性がある。I / O 要求の数が増えるにつれて、I O Vにより利用可能性をもたらすことができ、最新のC P U 仮想化において見られる性能レベルに匹敵するように、(仮想化された) I / O リソースの性能、スケーラビリティおよび融通性を向上させることができる。

10

【0042】

一実施形態に従うと、I / O リソースの共有を可能にして、V M からリソースへのアクセスが保護されることを可能にし得るようなI O V が所望される。I O V は、V M にエクスポー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 : D A) 、およびシングルルートI / O 仮想化 (S R - I O V) などのさまざまなタイプのI O V 技術が存在し得る。

20

【0043】

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

【0044】

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

30

【0045】

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

40

50

【 0 0 4 6 】

S R - I O V は、P C I e デバイスが、各々のゲストに 1 つの仮想デバイスを割り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 S R - I O V デバイスは、少なくとも 1 つの物理機能 (P F) と、1 つ以上の関連付けられた仮想機能 (V F) とを有する。P F は、仮想マシンモニタ (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 機能であるのに対して、V F は軽量の P C I e 機能である。各々の V F はそれ自体のベースアドレス (base address : B A R) を有しており、固有のリクエスト I D が割り当てられている。固有のリクエスト I D は、I / O メモリ管理ユニット (I / O memory management unit : I O M M U) がさまざまな V F への / からのトラフィックストリームを区別することを可能にする。I O M M U はまた、メモリを適用して、P F と V F との間の変換を中断する。

10

【 0 0 4 7 】

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

20

【 0 0 4 8 】

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

30

【 0 0 4 9 】

インフィニバンド S R - I O V アーキテクチャ - 共有ポート

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

【 0 0 5 0 】

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

40

【 0 0 5 1 】

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 H C A) は、物理機能 3 2 0 と仮想機能 3 3 0、3 5 0、3 5 0 との間において単一の共有 L I D および共有キュー対 (Queue Pair : Q P) のスペースがあるネットワークにおいて単一のポートとして現われる。しかしながら、各々の機能 (すなわち、物理機能および仮想機能) はそれら自体の G I D を有し得る。

【 0 0 5 2 】

図 4 に示されるように、一実施形態に従うと、さまざまな G I D を仮想機能および物理

50

機能に割当てることができ、特別のキュー対である Q P 0 および Q P 1 (すなわちインフィニバンド管理パケットのために用いられる専用のキュー対) が物理機能によって所有される。これらの Q P は V F にも同様にエクスポートされるが、V F は Q P 0 を使用することが許可されておらず (V F から Q P 0 に向かって入来するすべての S M P が廃棄され)、Q P 1 は、P F が所有する実際の Q P 1 のプロキシとして機能し得る。

【 0 0 5 3 】

一実施形態に従うと、共有ポートアーキテクチャは、(仮想機能に割当てられることによってネットワークに付随する) V M の数によって制限されることのない高度にスケラブルなデータセンタを可能にし得る。なぜなら、ネットワークにおける物理的なマシンおよびスイッチによって L I D スペースが消費されるだけであるからである。

10

【 0 0 5 4 】

しかしながら、共有ポートアーキテクチャの欠点は、トランスペアレントなライブマイグレーションを提供することができない点であり、これにより、フレキシブルな V M 配置についての可能性が妨害されてしまう。各々の L I D が特定のハイパーバイザに関連付けられており、かつハイパーバイザ上に常駐するすべての V M 間で共有されているので、マイグレートしている V M (すなわち、宛先ハイパーバイザにマイグレートする仮想マシン) は、その L I D を宛先ハイパーバイザの L I D に変更させなければならない。さらに、Q P 0 アクセスが制限された結果、サブネットマネージャは V M の内部で実行させることができなくなる。

20

【 0 0 5 5 】

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

図 5 は、一実施形態に従った例示的な v S w i t c h アーキテクチャを示す。図に示されるように、ホスト 4 0 0 (たとえばホストチャネルアダプタ) はハイパーバイザ 4 1 0 と対話することができ、当該ハイパーバイザ 4 1 0 は、さまざまな仮想機能 4 3 0、4 4 0 および 4 5 0 をいくつかの仮想マシンに割当てることができる。同様に、物理機能はハイパーバイザ 4 1 0 によって処理することができる。仮想スイッチ 4 1 5 もハイパーバイザ 4 0 1 によって処理することができる。

【 0 0 5 6 】

一実施形態に従うと、v S w i t c h アーキテクチャにおいては、各々の仮想機能 4 3 0、4 4 0、4 5 0 は完全な仮想ホストチャネルアダプタ (virtual Host Channel Adapter: v H C A) であり、これは、ハードウェアにおいて、V F に割当てられた V M に、I B アドレス一式 (たとえば G I D、G U I D、L I D) および専用の Q P スペースが割当てられていることを意味する。残りのネットワークおよび S M については、H C A 4 0 0 は、仮想スイッチ 4 1 5 を介して追加のノードが接続されているスイッチのように見えている。ハイパーバイザ 4 1 0 は P F 4 2 0 を用いることができ、(仮想機能に付与された) V M は V F を用いる。

30

【 0 0 5 7 】

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。

40

【 0 0 5 8 】

50

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

図6は、一実施形態に従った例示的なvPortの概念を示す。図に示されるように、ホスト300（たとえばホストチャンネルアダプタ）は、さまざまな仮想機能330、340および350をいくつかの仮想マシンに割り当てることができるハイパーバイザ410と対話することができる。同様に、物理機能はハイパーバイザ310によって処理することができる。

【0059】

一実施形態に従うと、ベンダーに実装の自由を与えるためにvPort概念は緩やかに定義されており（たとえば、当該定義では、実装がSRIOV専用とすべきであるとは規定されていない）、vPortの目的は、VMがサブネットにおいて処理される方法を標準化することである。vPort概念であれば、空間ドメインおよび性能ドメインの両方においてよりスケーラブルであり得る、SR - IOV共有ポートのようなアーキテクチャおよびvSwitchのようなアーキテクチャの両方、または、これらのアーキテクチャの組合せが規定され得る。また、vPortはオプションのLIDをサポートするとともに、共有ポートとは異なり、SMは、vPortが専用のLIDを用いていなくても、サブネットにおいて利用可能なすべてのvPortを認識する。

10

【0060】

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

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

20

【0061】

図7は、一実施形態に従った、LIDが予めポピュレートされた例示的なvSwitchアーキテクチャを示す。図に示されるように、いくつかのスイッチ501～504は、ネットワーク切替環境600（たとえばIBサブネット）内においてインフィニバンドファブリックなどのファブリックのメンバ間で通信を確立することができる。ファブリックはホストチャンネルアダプタ510、520、530などのいくつかのハードウェアデバイスを含み得る。さらに、ホストチャンネルアダプタ510、520および530は、それぞれ、ハイパーバイザ511、521および531と対話することができる。各々のハイパーバイザは、さらに、ホストチャンネルアダプタと共に、いくつかの仮想機能514、515、516、524、525、526、534、535および536と対話し、設定し、いくつかの仮想マシンに割り当てることができる。たとえば、仮想マシン1 550はハイパーバイザ511によって仮想機能1 514に割り当てることができる。ハイパーバイザ511は、加えて、仮想マシン2 551を仮想機能2 515に割り当て、仮想マシン3 552を仮想機能3 516に割り当てることができる。ハイパーバイザ531は、さらに、仮想マシン4 553を仮想機能1 534に割り当てることができる。ハイパーバイザは、ホストチャンネルアダプタの各々の上で十分な機能を有する物理機能513、523および533を介してホストチャンネルアダプタにアクセスすることができる。

30

【0062】

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

40

【0063】

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

50

スイッチのように見えている。

【0064】

一実施形態に従うと、本開示は、L I Dが予めポピュレートされたv S w i t c hアーキテクチャを提供するためのシステムおよび方法を提供する。図7を参照すると、L I Dは、さまざまな物理機能5 1 3、5 2 3および5 3 3に、さらには、仮想機能5 1 4 ~ 5 1 6、5 2 4 ~ 5 2 6、5 3 4 ~ 5 3 6（その時点でアクティブな仮想マシンに関連付けられていない仮想機能であっても）にも、予めポピュレートされている。たとえば、物理機能5 1 3はL I D 1が予めポピュレートされており、仮想機能1 5 3 4はL I D 1 0が予めポピュレートされている。ネットワークがブートされているとき、L I DはS R - I O V v S w i t c h対応のサブネットにおいて予めポピュレートされている。V FのすべてがネットワークにおけるV Mによって占有されていない場合であっても、ポピュレートされたV Fには、図7に示されるようにL I Dが割当てられている。

10

【0065】

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。

【0066】

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

20

30

【0067】

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

40

【0068】

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

50

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。

【 0 0 6 9 】

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

【 0 0 7 0 】

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

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

【 0 0 7 1 】

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。

【 0 0 7 2 】

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

【 0 0 7 3 】

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

【 0 0 7 4 】

一実施形態に従うと、本開示は、動的 L I D 割当てがなされた v S w i t c h アーキテクチャを提供するためのシステムおよび方法を提供する。図 8 を参照すると、L I D には、さまざまな物理機能 5 1 3、5 2 3 および 5 3 3 が動的に割当てられており、物理機能 5 1 3 が L I D 1 を受取り、物理機能 5 2 3 が L I D 2 を受取り、物理機能 5 3 3 が L I D 3 を受取る。アクティブな仮想マシンに関連付けられたそれらの仮想機能はまた、動的に割当てられた L I D を受取ることもできる。たとえば、仮想マシン 1 5 5 0 がアクテ

10

20

30

40

50

ィブであり、仮想機能 1 5 1 4 に関連付けられているので、仮想機能 5 1 4 には L I D 5 が割当てられ得る。同様に、仮想機能 2 5 1 5、仮想機能 3 5 1 6 および仮想機能 1 5 3 4 は、各々、アクティブな仮想機能に関連付けられている。このため、これらの仮想機能に L I D が割当てられ、L I D 7 が仮想機能 2 5 1 5 に割当てられ、L I D 1 1 が仮想機能 3 5 1 6 に割当てられ、L I D 9 が仮想機能 1 5 3 4 に割当てられている。L I D が予めポピュレートされた v S w i t c h とは異なり、アクティブな仮想マシンにその時点で関連付けられていない仮想機能は L I D の割当てを受けない。

【 0 0 7 5 】

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

10

【 0 0 7 6 】

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

【 0 0 7 7 】

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

20

【 0 0 7 8 】

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

30

【 0 0 7 9 】

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

40

【 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 は、いくらかの追加のネットワークおよびランタイム S M オーバーヘッドを犠牲にして、予めポピュレートされた L I D アーキテクチャモデルを備えた v S w i t c h の欠点を解決する

50

ことができる。VMが作成されるたびに、作成されたVM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VMがそのホストハイパーバイザと同じ経路を用いているので、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VM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。

10

【0081】

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

図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が予めポピュレートされたv S w i t c hを備えた例示的なv S w i t c hアーキテクチャを示す。図に示されるように、いくつかのスイッチ501~504は、ネットワーク切替環境800(たとえばIBサブネット)内においてインフィニバンドファブリックなどのファブリックのメンバ間で通信を確立することができる。ファブリックはホストチャネルアダプタ510、520、530などのいくつかのハードウェアデバイスを含み得る。ホストチャネルアダプタ510、520および530は、それぞれ、さらに、ハイパーバイザ511、521および531と対話することができる。各々のハイパーバイザは、さらに、ホストチャネルアダプタと共に、いくつかの仮想機能514、515、516、524、525、526、534、535および536と対話し、設定し、いくつかの仮想マシンに割当てることができる。たとえば、仮想マシン1 550は、ハイパーバイザ511によって仮想機能1 514に割当てることができる。ハイパーバイザ511は、加えて、仮想マシン2 551を仮想機能2 515に割当てることができる。ハイパーバイザ521は、仮想マシン3 552を仮想機能3 526に割当てることができる。ハイパーバイザ531は、さらに、仮想マシン4 553を仮想機能2 535に割当てることができる。ハイパーバイザは、ホストチャネルアダプタの各々の上において十分な機能を有する物理機能513、523および533を介してホストチャネルアダプタにアクセスすることができる。

20

30

【0082】

一実施形態に従うと、スイッチ501~504の各々はいくつかのポート(図示せず)を含み得る。これらいくつかのポートは、ネットワーク切替環境800内においてトラフィックを方向付けるためにリニアフォワーディングテーブルを設定するのに用いられる。

【0083】

一実施形態に従うと、仮想スイッチ512、522および532は、それぞれのハイパーバイザ511、521、531によって処理することができる。このような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VMに、IBアドレス一式(たとえばG I D、G U I D、L I D)および専用のQ Pスペースが割当てられていることを意味する。残りのネットワークおよびS M(図示せず)については、H C A 510、520および530は、仮想スイッチを介して、追加のノードが接続されているスイッチのように見えている。

40

【0084】

一実施形態に従うと、本開示は、動的L I D割当てがなされL I Dが予めポピュレートされたハイブリッドv S w i t c hアーキテクチャを提供するためのシステムおよび方法を提供する。図9を参照すると、ハイパーバイザ511には、予めポピュレートされた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521には、L I Dが予めポピュレートされて動的L I D割当てがなされたv S w i t c h

50

が配置され得る。ハイパーバイザ 5 3 1 には、動的 L I D 割当てがなされた v S w i t c h が配置され得る。このため、物理機能 5 1 3 および仮想機能 5 1 4 ~ 5 1 6 には、それらの L I D が予めポピュレートされている（すなわち、アクティブな仮想マシンに付与されていない仮想機能であっても L I D が割当てられている）。物理機能 5 2 3 および仮想機能 1 5 2 4 にはそれらの L I D が予めポピュレートされ得るとともに、仮想機能 2 5 2 5 および仮想機能 3 5 2 6 にはそれらの L I D が動的に割当てられている（すなわち、仮想機能 2 5 2 5 は動的 L I D 割当てのために利用可能であり、仮想機能 3 5 2 6 は、仮想マシン 3 5 5 2 が付与されているので、1 1 という L I D が動的に割当てられている）。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 5 】

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

【 0 0 8 6 】

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 7 】

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）

一実施形態に従うと、単一のサブネット内において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、本開示の実施形態はまた、2 つ以上のサブネットにわたるインフィニバンドファブリックを提供することができる。

【 0 0 8 8 】

図 1 0 は、一実施形態に従った、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。図に示されるように、サブネット A 1 0 0 0 内において、いくつかのスイッチ 1 0 0 1 ~ 1 0 0 4 は、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サブネット A 1 0 0 0（たとえば I B サブネット）内における通信を提供することができる。ファブリックは、たとえばチャネルアダプタ 1 0 1 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 1 0 1 0 は、さらに、ハイパーバイザ 1 0 1 1 と対話し得る。ハイパーバイザは、さらに、それが対話するホストチャネルアダプタと共に、いくつかの仮想機能 1 0 1 4 を設定し得る。ハイパーバイザは、加えて、仮想機能 1 1 0 1 4 に割当てられている仮想マシン 1 1 0 1 5 などの仮想マシンを仮想機能の各々に割当て得る。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、物理機能 1 0 1 3 などの十分な機能を有する物理機能によって、それらの関連付けられたホストチャネルアダプタにアクセスし得る。サブネット B 1 0 4 0 内では、いくつかのスイッチ 1 0 2 1 ~ 1 0 2 4 は、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サブネット B 1 0 4 0（たとえば I B サブネット）内において通信を提供し得る。ファブリックは、たとえばチャネルアダプタ 1 0 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 1 0 3 0 は、さらに、ハイパーバイザ 1 0 3 1 と対話し得る。ハイパーバイザは、さらに、それが対話するホストチャネルアダプタと共に

、いくつかの仮想機能 1 0 3 4 を設定し得る。ハイパーバイザは、加えて、仮想機能 2 1 0 3 4 に割当てられている仮想マシン 2 1 0 3 5 などの仮想マシンを仮想機能の各々に割当て得る。ハイパーバイザは、各々のホストチャンネルアダプタ上で、物理機能 1 0 3 3 などの十分な機能を有する物理機能によってそれらの関連付けられたホストチャンネルアダプタにアクセスし得る。なお、1つのホストチャンネルアダプタだけが各々のサブネット（すなわちサブネット A およびサブネット B）内において図示されているが、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が各々のサブネット内に含まれ得ることが理解されるはずである。

【 0 0 8 9 】

一実施形態に従うと、上述のとおり、ホストチャンネルアダプタの各々は、加えて、仮想スイッチ 1 0 1 2 および仮想スイッチ 1 0 3 2 などの仮想スイッチに関連付けられ得るとともに、各々の H C A はさまざまなアーキテクチャモデルで設定され得る。図 1 0 内における両方のサブネットは予めピュレートされた L I D アーキテクチャモデルを備えた v S w i t c h を用いるものとして示されているが、これは、このようなすべてのサブネット構成が同様のアーキテクチャモデルに従うべきであることを示唆するよう意図したものであるのではない。

10

【 0 0 9 0 】

一実施形態に従うと、サブネット A 1 0 0 0 内のスイッチ 1 0 0 2 がルータ 1 0 0 5 に関連付けられ、サブネット B 1 0 4 0 内のスイッチ 1 0 2 1 がルータ 1 0 0 6 に関連付けられているなどのように、各々サブネット内の少なくとも1つのスイッチがルータに関連付けられ得る。

20

【 0 0 9 1 】

一実施形態に従うと、少なくとも1つのデバイス（たとえばスイッチ、ノード...など）は、ファブリックマネージャ（図示せず）に関連付けられ得る。ファブリックマネージャを用いることにより、たとえば、サブネット間ファブリックトポロジーを発見し、ファブリックプロファイル（たとえば仮想マシンファブリックプロファイル）を作成し、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のデータベースオブジェクトを構築することができる。加えて、ファブリックマネージャは、どのサブネットが、どのルータポートを介して、どのパーティション番号を用いて通信することが許可されているかについて、法的なサブネット間接続性を定義し得る。

30

【 0 0 9 2 】

一実施形態に従うと、サブネット A 内における仮想マシン 1 などの送信元におけるトラフィックが、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における宛先にアドレス指定されると、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得る、次いで、ルータ 1 0 0 5 は、ルータ 1 0 0 6 とのリンクを介してサブネット B にトラフィックを渡し得る。

【 0 0 9 3 】

仮想デュアルポートルータ

一実施形態に従うと、デュアルポートルータの抽象化は、サブネットとサブネットとの間のルータ機能が、通常の L R H ベースの切替えを実行することに加えて、グローバル・ルート・ヘッダ（global route header : G R H）とローカル・ルート・ヘッダ（local route header : L R H）との変換を実行する能力を有するスイッチハードウェア実装に基づいて定義されることを可能にするための単純な方法を提供することができる。

40

【 0 0 9 4 】

一実施形態に従うと、仮想デュアルポートルータは、対応するスイッチポートの外部に論理的に接続され得る。この仮想デュアルポートルータは、サブネットマネージャなどの標準的な管理エンティティにインフィニバンド規格対応の概念を提供することができる。

【 0 0 9 5 】

一実施形態に従うと、デュアルポート・ルータモデルは、各々のサブネットが、サブネットに対する入口経路におけるアドレスマッピングおよびパケットの転送を十分に制御す

50

るような態様で、かつ、不正確に接続されたサブネットのいずれかの内部におけるルーティングおよび論理接続性に影響を及ぼすことなく、さまざまなサブネットが接続可能であることを示唆している。

【0096】

一実施形態に従うと、不正確に接続されたファブリックを伴う状況においては、仮想デュアルポートルータ抽象化を用いることにより、リモートサブネットに対して意図せず物理的に接続性している状態でサブネットマネージャおよびIB診断ソフトウェアなどの管理エンティティが正確に機能することをも可能にし得る。

【0097】

図11は、一実施形態に従った、高性能コンピューティング環境における2つのサブネット間の相互接続を示す。仮想デュアルポートルータを備えた構成に先だって、サブネットA1101におけるスイッチ1120は、スイッチ1120のスイッチポート1121を通じて、物理的接続1110を介し、サブネットB1102におけるスイッチ1130に、スイッチ1130のスイッチポート1131を介して接続され得る。このような実施形態においては、各々のスイッチポート1121および1131は、スイッチポートおよびルータポートの両方の機能を果たし得る。

【0098】

一実施形態に従うと、この構成に関する問題は、インフィニバンドサブネットにおけるサブネットマネージャなどの管理エンティティが、スイッチポートおよびルータポートの両方である物理ポートを区別することができない点である。このような状況においては、SMは、スイッチポートを、そのスイッチポートに接続されたルータポートを有するものとして取扱うことができる。しかしながら、スイッチポートが、別のサブネットマネージャとともに、たとえば物理リンクを介して別のサブネットに接続される場合、サブネットマネージャは物理リンク上で発見メッセージを送出することができる。しかしながら、このような発見メッセージは他のサブネットにおいて許可することはできない。

【0099】

図12は、一実施形態に従った、高性能コンピューティング環境におけるデュアルポート仮想ルータ構成を介する2つのサブネット間の相互接続を示す。

【0100】

一実施形態に従うと、構成後、あるサブネットマネージャが担当しているサブネットの端部を表わしている適切なエンドノードを当該サブネットマネージャが認識するように、デュアルポート仮想ルータ構成が提供され得る。

【0101】

一実施形態に従うと、サブネットA1201におけるスイッチ1220において、スイッチポートは、仮想リンク1223を介して仮想ルータ1210におけるルータポート1211に接続（すなわち、論理的に接続）され得る。仮想ルータ1210（たとえばデュアルポート仮想ルータ）は、スイッチ1220の外部にあるものとして示されており、実施形態においては、スイッチ1220内に論理的に含まれ得るものであるが、第2のルータポート（ルータポートII1212）を含み得る。一実施形態に従うと、2つの端部を有し得る物理リンク1203は、ルータポートII1212と、サブネットB1202における仮想ルータ1230に含まれるルータポートII1232とを介して、物理リンクの第1の端部を介するサブネットA1201を、物理リンクの第2の端部を介するサブネットB1202と、接続することができる。仮想ルータ1230は、加えて、ルータポート1231を含み得る。ルータポート1231は、仮想リンク1233を介して、スイッチ1240上のスイッチポート1241に接続（すなわち、論理的に接続）され得る。

【0102】

一実施形態に従うと、サブネットA上のサブネットマネージャ（図示せず）は、サブネットマネージャが制御するサブネットのエンドポイントとして、仮想ルータ1210上においてルータポート1211を検出し得る。デュアルポート仮想ルータ抽象化は、サブネットA上のサブネットマネージャが、（たとえば、インフィニバンド規格に準拠して定義

10

20

30

40

50

されるような)通常の状態ではサブネット A を処理することを可能にし得る。サブネット管理エージェントレベルにおいて、SM が通常のスイッチポートを認識するように、デュアルポート仮想ルータ抽象化が提供され得るとともに、SMA レベルにおいては、スイッチポートに接続された別のポートが存在する抽象化が行なわれ、このポートはデュアルポート仮想ルータ上のルータポートとなる。ローカル SM においては、従来のファブリックトポロジーを使用し続けることができる (SM は当該ポートをトポロジーにおける標準スイッチポートとして認識する)。このため、SM はルータポートをエンドポートとして認識する。物理的接続は、2 つの異なるサブネットにおけるルータポートとしても構築される 2 つのスイッチポート間において構成され得る。

【0103】

一実施形態に従うと、デュアルポート仮想ルータはまた、物理リンクが同じサブネットにおける他の何らかのスイッチポートに、または、別のサブネットに接続されるよう意図されていなかったスイッチポートに、間違っ て接続され得るという問題を解決することもできる。したがって、この明細書中に記載される方法およびシステムはまた、サブネットの外部に存在するものも表わしている。

【0104】

一実施形態に従うと、サブネット A などのサブネット内においては、ローカル SM がスイッチポートを決定し、次いで、そのスイッチポートに接続されるルータポート (たとえば、仮想リンク 1223 を介してスイッチポート 1221 に接続されるルータポート 1211) を決定する。SM が、当該 SM が管理しているサブネットの端部としてルータポート 1211 を認識するので、SM は、この点 (たとえばルータポート EE1212) を越えて発見メッセージおよび / または管理メッセージを送信することができない。

【0105】

一実施形態に従うと、上述のデュアルポート仮想ルータは、デュアルポート仮想ルータが属するサブネット内において、デュアルポート仮想ルータ抽象化が管理エンティティ (たとえば SM または SMA) によって完全に管理されるという利点を提供する。ローカル側でのみ管理を許可することにより、システムは、外部の独立した管理エンティティを提供する必要がなくなる。すなわち、サブネット接続に対するサブネットの両サイドが、それ自体のデュアルポート仮想ルータを構成する役割を果たし得る。

【0106】

一実施形態に従うと、離れた宛先 (すなわち、ローカルサブネットの外側) にアドレス指定されている SMP などのパケットが上述のデュアルポート仮想ルータとして構成されていないローカルターゲットポートに到達する状況においては、ローカルポートは、それがルータポートではないことを規定するメッセージを戻すことができる。

【0107】

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。

【0108】

図 13 は、一実施形態に従った、高性能コンピューティング環境におけるデュアルポート仮想ルータをサポートするための方法を示す。ステップ 1310 において、当該方法は、1 つ以上のマイクロプロセッサを含む 1 つ以上のコンピュータにおいて、第 1 のサブネットを提供し得る。第 1 のサブネットは複数のスイッチを含み、複数のスイッチは少なくともリーフスイッチを含む。複数のスイッチの各々は複数のスイッチポートを含む。第 1 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む。各々の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む。第 1 のサブネットはさらに、複数のエンドノードを含む。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ている。第 1 のサブネットはさらに、複数のスイッチおよび複数のホストチャネルアダプタのうち 1 つの上で実

10

20

30

40

50

行しているサブネットマネージャを含む。

【0109】

ステップ1320において、当該方法は、複数のスイッチのうち1スイッチ上において複数のスイッチポートのうち1スイッチポートをルータポートとして構成し得る。

【0110】

ステップ1330において、当該方法は、ルータポートとして構成されたスイッチポートを仮想ルータに論理的に接続し得る。仮想ルータは少なくとも2つの仮想ルータポートを含む。

【0111】

SMPベースの接続性チェックを可能にするためのルータSMA抽象化

一実施形態に従うと、サブネット管理パケット（Subnet Management Packet：SMP）は、パケットがルータポートを越えて送信されるであろうことを示唆するアドレス指定情報を有することが許可されていない。しかしながら、（仮想）ルータポートの遠隔側での物理的接続性の発見（すなわち、遠隔接続性のローカルな発見）を可能にするために、SMA属性の新しいセットが定義され得る。この場合、このような各々の新しい属性は、標準的なSMA属性またはベンダー特有のSMA属性とともに「遠隔情報」を表わしている。

【0112】

一実施形態に従うと、ルータSMAが「遠隔」情報/属性を表わす属性を処理すると、次いで、対応するSMP要求は、元の要求の送信側に完全にトランスペアレントな態様で外部の物理リンク上で送信することができる。

【0113】

一実施形態に従うと、ローカルSMAは、受信要求とは無関係に遠隔からの発見を実行して関連情報をローカルにキャッシュすることを選択し得るか、または、単純なプロキシのように動作して、「遠隔」情報/属性を指定する要求を受取るたびに外部リンクに対する対応する要求を生成し得る。

【0114】

一実施形態に従うと、「遠隔」属性を要求するSMPがローカルサブネット側（すなわち、ローカルスイッチポートに論理的に接続されている仮想ルータポート）から、または外部リンク（すなわち、仮想ルータの遠隔側）から、受取られたかどうかを追跡することによって、SMA実装は、ローカルサブネットにおける元の要求を表現する観点、または、ピアルータポートからのプロキシ要求を表現する観点から、遠隔要求がどの程度まで有効であるかを追跡することができる。

【0115】

一実施形態に従うと、IB規格に準拠したSMの場合、ルータポートはサブネットにおけるエンドポートである。このため、（発見に用いられてサブネットを構成する）低レベルSMPはルータポートにわたって送信することができない。しかしながら、サブネット間トラフィックのためのルートを維持するために、ローカルSMまたはファブリックマネージャは、ローカルリソースのいずれかの構成を構築する前に、物理リンクの遠隔側での物理的接続性を観察可能である必要がある。しかしながら、遠隔接続性を認識したいとの要望に関して、SMは、物理リンクの遠隔側を構成することを許可することができない（すなわち、SMの構成はそれ自体のサブネット内に含まれていなければならない）。

【0116】

一実施形態に従うと、SMAモデル拡張は、ローカルルータポートにアドレス指定されているパケット（すなわちSMP）を送信する可能性を考慮に入れている。パケットがアドレス指定されているSMAは、パケットを受取り得るとともに、次いで、要求された情報が（たとえば、サブネットにわたる物理リンクによって接続されている）リモートノード上にあることを定義している新しい属性を適用し得る。

【0117】

一実施形態に従うと、SMAは、プロキシとして動作する（SMPを受取って別の要求

10

20

30

40

50

を送信する)ことができるか、または、SMAは元の packets を変更して、これをサブネット間 packets として送信することができる。SMAは、packets におけるアドレス情報を更新することができる。これにより、SMPに対する1ホップ・ダイレクトルート経路の追加が更新される。次いで、リモートノード(ルータポート)によってSMPを受取ることができる。SMPは、リモートエンド上のノードが同じ態様で(たとえば仮想ルータとして)構成されているか、または、基本的なスイッチポート(たとえば、レガシースイッチ実装に対する物理的接続性)として構成されているかどうかとは無関係に、機能し得る。次いで、受信ノードが認識し得る要求 packets は基本的な要求となり、通常の態様で応答することとなる。実際には、サブネットを越えて発信された要求は受信ノードに対してトランスペアレント(不可視)である。

10

【0118】

一実施形態に従うと、これは、抽象化を利用することによって、ローカルサブネットによるリモートサブネットの発見を可能にする。提供されるSMA抽象化は、遠隔側で、ローカルな(すなわち、実行中のサブネット発見からは離れている)サブネットから問い合わせられていたことを認識することなく、(たとえば物理リンクにわたって)リモートサブネットから情報を検索することを可能にする。

【0119】

アドレス指定スキーム

一実施形態に従うと、IB規格と準拠させたままにするために、SMP packets は、サブネットの境界によって連結されている(すなわち、SMは、それが関連付けられているサブネットの外にある情報を「認識する」かまたは発見することが許可されていない)。しかしながら、仮想ルータポートのリモートエンドからの接続性情報などの情報(すなわち、サブネット境界を越えた1「ホップ」)を検索する必要性が依然として存在する。この情報は、ベンダー特有のSMP(vendor specific SMP: VSMP)と称され得る特別なタイプのSMPに符号化することができる。

20

【0120】

一実施形態に従うと、VSMPは、概して、一般的なSMP(サブネット向けのSMP)と同様のアドレス指定方式を利用することができる。DR(Directed Routing(ダイレクティッドルーティング)):この場合、SMPは、スイッチ間を移動しているときにどのポートから生じているかを明確に示すことができる)およびLIDの両方のルーティングを用いることができる。しかしながら、ルータポートのリモートエンドに適用され得る属性のために、属性修飾子におけるシングルビットを用いてローカルポート対リモートポートを示すことができる。属性のリモートビットが設定されている場合、付加的な処理がルータポートを表わすSMAにおいて生じ得る。

30

【0121】

一実施形態に従うと、アドレス指定方式の重要な局面は、システムの構成または配線が間違っただけでなされていたとしても、ルータポートのリモート・ピア・ポートが要求に応答することができるという点である。たとえば、仮想ルータが、リモートサブネットの汎用のスイッチポートにおいて、物理リンクを介して、たとえばリモートサブネットに、接続されている場合、リモートサブネットの汎用スイッチポートを処理するSMAは、リモート属性がサポートされていないことを示す状態値とともに要求に応答することができ、さらに、当該応答は要求しているSMに到達することができる。

40

【0122】

一実施形態に従うと、SMP(すなわち、VSMP)は、DR経路を用いることによって、またはルータポートのLIDを介して、ローカルルータポートに送信することができる。要求された情報がルータポートのリモート・ピア・ポートのためのものである場合、属性修飾子のリモートフラグ/ビットを設定することができる。SMAは、このようなVSMPを受取ると、packets を変更し、リモートエンドをアドレス指定する付加的なDRステップを追加することができる。

【0123】

50

一実施形態に従うと、パケット属性の一部（たとえば16ビット属性修飾子のビット）はVSMPがローカルであるかリモートであるかどうかを信号で伝えるために用いることができる。たとえば、この部分を1の値に設定することにより、リモート・ピア・ポートがVSMPのための最終宛先であることを示唆することができる。属性のこの部分はリモートフラグとも称され得る。

【0124】

一実施形態に従うと、リモートフラグに加えて、どの宛先インスタンスが、最終宛先への経路に沿ってVSMPを処理するべきかを示すことができる追加のフラグが存在し得る。パケット属性のうち2つの付加的部分（たとえば属性修飾子の2ビット）は、この目的のために用いることができる。第1の受取りフラグと称される第1の部分（たとえばビット20）は、（元の要求の宛先アドレスと一致するはずである）ローカルルータポートによって受取られたとき、パケットの処理が予想されることを示し得る。第1の受信側における予想される処理がすべて実行されると、パケット属性の第2の部分（たとえば属性修飾子のビット21）を設定して、パケットをリモートエンドに転送することができる。この第2の部分は第2の受取りフラグと称され得る。

10

【0125】

DRルーティングされたパケット

一実施形態に従うと、一例として、DRルーティングされたパケットは例示的なフローに追従し得る。ソースノードは、LIDAにおいて、ルータ1を宛先ノードとして指定する要求パケットを開始することができる。例示的なDRルーティングパケット構成は以下のとおりである。

20

【0126】

- ・MADHdr.Class = 0x81 (DRルーティングされたSMP)
- ・MADHdr.Method = 0x1 (取得)
- ・LRH.SLID = LRH.DLID = 0xffff (許可LID)
- ・MADHdr.DrSLID = MADHdr.DrDLID = 0xffff
- ・MADHdr.AttrID = <VSMP attrID>
- ・MADHdr.AttrMod.remote = 1
- ・MADHdr.AttrMod.first_receiver = 1
- ・MADHdr.InitPath = DR経路からRtr1 (LID B)
- ・MADHdr.HopCnt = N
- ・MADHdr.HopPtr = 0

30

要求パケットがルータ1に到達すると、当該要求パケットは、対応するSMAに渡され得るとともに、さらにこの対応するSMAが、要求が有効であることを確認し得る。要求の有効性を確認した後、SMAはパケットを変更し得る。この変更は、1つの余分なホップを用いてDR経路を拡張することと、第2の受取りフラグを設定することとを含み得る。このような例示的な構成変更を以下に示す。

【0127】

- ・MADHdr.HopCnt = N+1 (1つの余分なホップを用いてDR経路を拡張)
- ・MADHdr.InitPath[N+1] = (仮想ルータ外部ポート番号(すなわち、2))
- ・MADHdr.AttrMod.second_receiver = 1

40

一実施形態に従うと、SMAのサブネット管理インターフェイス(subnet management interface: SMI)層は、SMAが物理リンク上でVSMPを転送する前に、余分なホップを用いてDRを更新することができる。

【0128】

一実施形態に従うと、物理リンクの遠隔側でVSMPを受取るSMAのために、VSMPによって用いられるアドレス指定方式は、正常にIB規定されたパケットの方式のように見える場合もある。

【0129】

一実施形態に従うと、物理リンクの遠隔側でVSMPを受取るSMAがVSMPの最終

50

宛先である場合、V S M Pは受取られたS M Aによって確認することができる。確認はフラグ設定に対して入力ポートをチェックすることを含み得る。リモートフラグならびに第1の受取りフラグおよび第2の受取りフラグが設定される場合、パケットは物理ポート上で(すなわち、仮想ルータで構成されたポートの外部リンク側から)受取ることができる。遠隔にある第1の受取りフラグだけが設定される場合、パケットは、仮想リンク上で(すなわち、仮想ルータポートの内部スイッチ側から)到達し得る。確認が失敗すれば、状態が適切なエラーメッセージに設定され得る。

【0130】

図14は、一実施形態に従った、D RルーティングされたV S M Pパケットを示すフローチャートである。

10

【0131】

ステップ1401において、一実施形態に従うと、サブネットA内において、サブネットAのサブネットマネージャなどのエンティティは、D RルーティングされたV S M Pを介して接続性情報を要求することができる。D RルーティングされたV S M Pの宛先は、サブネットA内におけるルータ1などのルータであり得る。

【0132】

ステップ1402において、D RルーティングされたV S M Pがルータ1に到達し得る。ルータ1は、一実施形態においては、上述のとおり、スイッチ1などのスイッチに含まれ得る。D RルーティングされたV S M Pは、確認のためにスイッチ1のS M Aに渡され得る。

20

【0133】

ステップ1403において、一実施形態に従うと、S M AはD RルーティングされたV S M Pを変更し得る。この変更は、(第2のサブネットに対する物理リンクにわたるホップのための)ホップ・カウンタを拡張することと、第2の受取りフラグを設定することを含み得る。

【0134】

ステップ1404において、一実施形態に従うと、S M I (S M Iはスイッチ/ルータのS M Aに関連付けられている)は、V S M Pのホップポイントを更新し得る。

【0135】

ステップ1405において、一実施形態に従うと、S M Iは、物理リンク上で、サブネットA 1420とサブネットB 1430との間のサブネット境界にわたってD R V S M Pを転送することができる。この場合、物理リンクの第1の端部は、サブネットA内のルータに接続され得るとともに、物理リンクの第2の端部は、サブネットB内のルータに接続され得る。

30

【0136】

ステップ1406において、一実施形態に従うと、サブネットB 1430内のルータ2のS M A は物理リンクからV S M Pを受取り得る。

【0137】

ステップ1407において、一実施形態に従うと、(S M A に関連付けられた)S M I がV S M P要求を確認し得る。

40

【0138】

一実施形態に従うと、応答エンティティ(たとえば、物理リンクの遠隔側のルータ)は、S M P 応答を完了し、応答を示す方向属性を設定し得る。このような例示的な構成を以下に示す。

【0139】

・MADHdr.Method = 0x81 (GetResp)

・MADHdr.Direction = 1 (応答を示す)

一実施形態に従うと、S M I層は、ホップポイントをデクリメントし、物理リンク上で、ローカル側に応答を転送し得る。ローカルルータにおけるS M Aは、要求に対してなされた変更を元に戻し得る。ローカル側ルータにおけるS M Iは、ホップをデクリメントす

50

ることと、仮想リンク上で（すなわち、スイッチポートと仮想ルータポートとの間の内部仮想リンク上で）スイッチに応答を送出することを含む通常の処理を実行し得る。次のホップのために、S M I が再びホップをデクリメントし、要求が本来到達していた物理スイッチポート上でパケットを送出し得る。

【 0 1 4 0 】

図 1 5 は、一実施形態に従った、D R ルーティングされた V S M P パケットに対する応答を示すフローチャートである。

【 0 1 4 1 】

ステップ 1 5 0 1 において、一実施形態に従うと、ルータ 2 の S M A は、V S M P に応答情報（たとえば接続性情報）を記入し、応答を示すために方向性フラグを設定し得る。

10

【 0 1 4 2 】

ステップ 1 5 0 2 において、一実施形態に従うと、S M I は、応答のホップポイントをデクリメントし得るとともに、たとえば、2つの端部を有する物理リンク上でサブネット B 1 5 3 0 からサブネット A 1 5 2 0 に応答を転送し返し得る。物理リンクの第 1 の端部は、サブネット A 内のルータに接続され得るとともに、物理リンクの第 2 の端部は、サブネット B 内のルータに接続され得る。

【 0 1 4 3 】

ステップ 1 5 0 3 において、一実施形態に従うと、ルータ 1 における S M A は応答を受取って、V S M P 上でなされた変更を元に戻し得る。

20

【 0 1 4 4 】

ステップ 1 5 0 4 において、一実施形態に従うと、ルータ 1 における S M I は、V S M P に対する応答を処理し、物理スイッチに対してリンク（たとえば、仮想ルータ 1 と物理スイッチとの間の仮想リンク）上で応答を送信し得る。

【 0 1 4 5 】

ステップ 1 5 0 5 において、一実施形態に従うと、S M I は、ホップポイントカウンタをデクリメントし、V S M P が元々受取られていた物理スイッチポート上でパケットを送出し得る。

【 0 1 4 6 】

L I D ルーティングされたパケット

30

一実施形態に従うと、一例として、L I D ルーティングされたパケットは例示的なフローに追従し得る。ソースノードは、L I D A において、宛先ノードとしてルータ 1 を指定する要求パケットを開始させ得る。このような例示的なパケットは以下の構成を有し得る。

【 0 1 4 7 】

- ・MADHdr.Class = 0x01 (L I D ルーティングされた S M P)
- ・MADHdr.Method = 0x1 (取得)
- ・LRH.SLID = L I D A
- ・LRH.DLID = L I D B
- ・MADHdr.AttrID = <VSMP attrID>
- ・MADHdr.AttrMod.remote = 1
- ・ADHdr.AttrMod.first_receiver = 1

40

一実施形態に従うと、要求パケットがルータ 1 に到達すると、当該要求パケットは、対応する S M A に渡され、当該対応する S M A は、要求が有効であることを確認し得る。要求の有効性を確認した後、S M A はパケットを変更し得る。この変更は、端部における単一の D R 経路を追加することを含み得る。結果として、これは、アドレスの合計が L I D ルーティングされたパケットと D R ルーティングされたホップとの組合せになり得ることを意味する。例示的な構成は以下のとおりである。

【 0 1 4 8 】

- ・MADHdr.Class = 0x81 (D R ルーティングされた S M P)

50

- ・MADHdr.HopCnt = 1
- ・MADHdr.HopPtr = 0
- ・MADHdr.Direction = 0 (アウトバウンド)
- ・MADHdr.InitPath[1] = (仮想ルータ外部ポート番号(すなわち、2))
- ・MADHdr.DrSLID = LRH.SLID (すなわち、LRH.SLIDは元のリクエストからのL I D Aを含む)
- ・MADHdr.DrDLID = 0xffff
- ・MADHdr.AttrMod.second_receiver = 1

一実施形態に従うと、ルータにおけるS M I層は、アウトバウンドの packets を正常に処理し、これを物理リンク上で送出し得る。例示的な構成は以下のとおりである。

10

【0149】

- ・LRH.SLID = LRH.DLID = 0xffff
- ・MADHdr.HopPtr = 1 (1ずつインクリメント)

一実施形態に従うと、宛先ルータにおける(物理リンクの他方の端部における)S M Iは、それが要求についての宛先であると判断し、関連付けられたS M Aに対してV S M Pを渡し得る。

【0150】

図16は、一実施形態に従った、D RルーティングされたV S M P packets を示すフローチャートである。

【0151】

ステップ1601において、一実施形態に従うと、サブネットA内において、サブネットAのサブネットマネージャなどのエンティティは、L I DルーティングされたV S M Pを介して接続性情報を要求し得る。L I DルーティングされたV S M Pの宛先は、ルータ1などの、サブネットA内のルータであり得る。

20

【0152】

ステップ1602において、L I DルーティングされたV S M Pはルータ1に到達し得る。ルータ1は、一実施形態においては、上述のとおり、スイッチ1などのスイッチに含まれ得る。L I DルーティングされたV S M Pは、確認のためにスイッチ1のS M Aに渡され得る。

【0153】

ステップ1603において、一実施形態に従うと、S M AはD RルーティングされたV S M Pを変更し得る。この変更は、単一のホップD Rルーティングされた経路をアドレスの最後に追加することと、(第2のサブネットに対する物理リンクにわたるホップのために)ホップ・カウンタを拡張することと、第2の受取りフラグを設定することとを含み得る。

30

【0154】

ステップ1604において、一実施形態に従うと、S M I (S M Iはスイッチ/ルータのS M Aに関連付けられている)は、V S M Pのホップポイントを更新し得る。

【0155】

ステップ1605において、一実施形態に従うと、S M Iは、物理リンク上で、サブネットA 1620とサブネットB 1630との間におけるサブネット境界にわたって、元々L I DにルーティングされていたV S M Pを転送し得る。この場合、物理リンクの第1の端部はサブネットA内のルータに接続され得るとともに、物理リンクの第2の端部は、サブネットB内のルータに接続され得る。

40

【0156】

ステップ1606において、一実施形態に従うと、サブネットB 1630内のルータ2のS M Aは、このとき、D RルーティングされたV S M Pを物理リンクから受取り得る。

【0157】

ステップ1607において、一実施形態に従うと、(S M Aに関連付けられた)S M

50

I は V S M P 要求を確認し得る。

【 0 1 5 8 】

一実施形態に従うと、応答フローのために、ルータ 2 における S M A が V S M P を確認し得る。この確認は、フラグ設定に対して入力ポートをチェックすることを含み得る。リモートフラグならびに第 1 の受取りフラグおよび第 2 の受取りフラグが設定される場合、パケットは物理ポート上で（すなわち、仮想ルータで構成されたポートの外部リンク側から）受取ることができる。遠隔側である第 1 の受取りフラグだけが設定される場合、パケットは、仮想リンク上で（すなわち、仮想ルータポートの内部スイッチ側から）到達し得る。確認が失敗すれば、状態が適切なエラーメッセージに設定され得る。

【 0 1 5 9 】

一実施形態に従うと、ルータ 2 における S M A は、S M P 応答を完了し、応答を示す方向属性を設定し得る。このような例示的な構成を以下に示す。

【 0 1 6 0 】

- ・MADHdr.Method = 0x81 (GetResp)
- ・MADHdr.Direction = 1 (応答を示す)
- ・LRH.SLID = 0xffff
- ・LRH.DLID = ReqMAD.MADHdr.SLID = 0xffff

ルータ 2 における S M I は、ホップポイントをデクリメントし、物理リンク上で応答をルータ 1 に転送し返し得る。次いで、ルータ 1 における S M A は、応答を物理スイッチポートから送信元に転送し返す前に、元の要求上でなされた変更を元に戻し得る。このように元に戻した後の例示的な構成を以下に示す。

【 0 1 6 1 】

- ・MADHdr.Class = 0x01
- ・LRH.DLID = MADHdr.DrSLID (すなわち、元のリクエストからの L I D A を含む)
- ・LRH.SLID = local LID = B

一実施形態に従うと、加えて、S M A が V S M P の D R 特有のフィールドをクリアし得るため、応答は、当該応答が元のリクエストに到達したときに元の V S M P と完全に一致しているように見えることとなる。

【 0 1 6 2 】

一実施形態に従うと、応答が送信元に戻ってくると、送信元は、正常なルータ L I D 応答を、それがルータ 1 において完全に処理されてしまっているかのように認識するだろう。

【 0 1 6 3 】

図 1 7 は、一実施形態に従った、L I D ルーティングされた V S M P パケットに対する応答を示すフローチャートである。

【 0 1 6 4 】

ステップ 1 7 0 1 において、一実施形態に従うと、ルータ 2 の S M A は、V S M P に対する応答情報（たとえば接続性情報）を記入し得るとともに、応答を示すために方向性フラグを設定し得る。

【 0 1 6 5 】

ステップ 1 7 0 2 において、一実施形態に従うと、S M I は、応答のホップポイントをデクリメントし、たとえば、2 つの端部を有する物理リンク上で、サブネット B 1 7 3 0 からサブネット A 1 7 2 0 に応答を転送し返し得る。物理リンクの第 1 の端部は、サブネット A 内のルータに接続され得るとともに、物理リンクの第 2 の端部は、サブネット B 内のルータに接続され得る。

【 0 1 6 6 】

ステップ 1 7 0 3 において、一実施形態に従うと、ルータ 1 における S M A は、応答を受取り得るとともに、D R ルーティングされたホップを取り除くことを含む、V S M P 上でなされた変更を元に戻し得る。

10

20

30

40

50

【0167】

ステップ1704において、一実施形態に従うと、ルータ1におけるSMIは、VSM Pに対する応答を処理し得るとともに、VSM Pが元々受取られていた物理スイッチポート上でLIDルーティングされた応答を送出し得る。

【0168】

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

10

【0169】

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。

20

【0170】

この発明の特徴はまた、たとえば、特定用途向け集積回路(application specific integrated circuit: ASIC)などのハードウェアコンポーネントを使用して、ハードウェアにおいて実現されてもよい。ここに説明された機能を行なうようにハードウェアス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。

【0171】

加えて、この発明は、この開示の教示に従ってプログラミングされた1つ以上のプロセッサ、メモリおよび/またはコンピュータ読取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。

30

【0172】

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

【0173】

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

40

【0174】

この発明の前述の説明は、例示および説明のために提供されてきた。それは、網羅的であるよう、またはこの発明を開示された形態そのものに限定するよう意図されてはいない。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

50

理およびその実用的応用を最良に説明するために選択され説明されたものであり、それにより、考えられる特定の使用に適したさまざまな実施形態についての、およびさまざまな変更例を有するこの発明を、当業者が理解できるようにする。この発明の範囲は、請求項およびそれらの同等例によって定義されるよう意図されている。

【 図 1 】

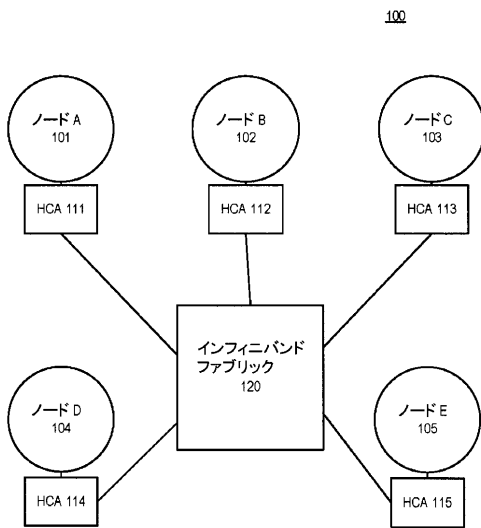


FIGURE 1

【 図 2 】

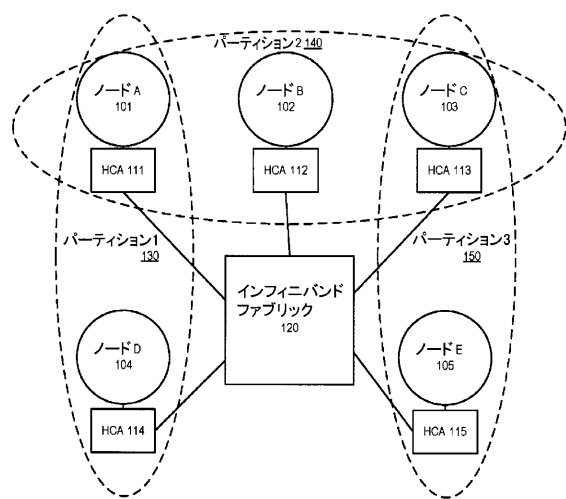


FIGURE 2

【 図 3 】

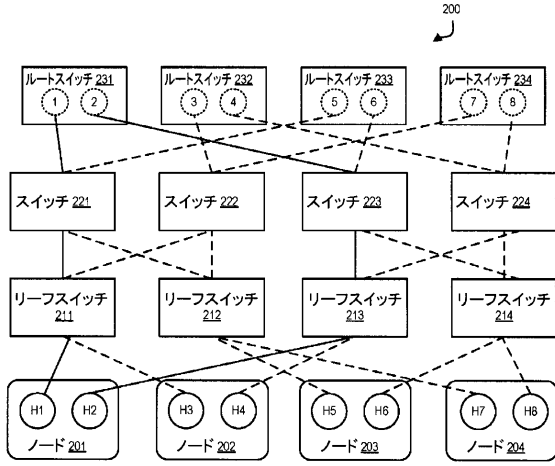


FIGURE 3

【 図 4 】

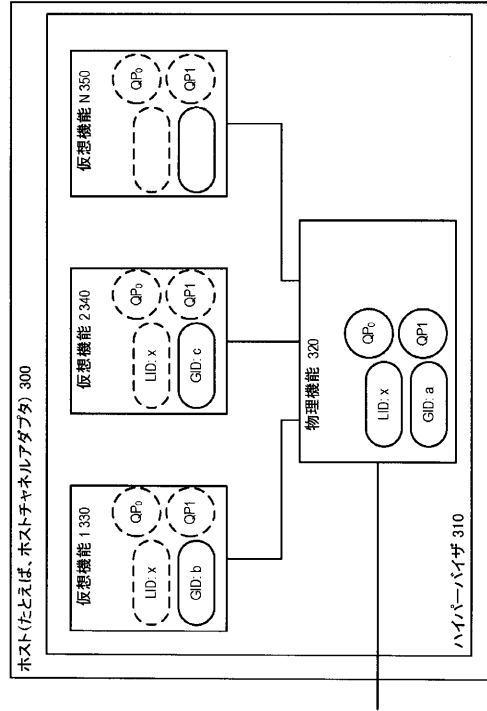


FIGURE 4

【 図 5 】

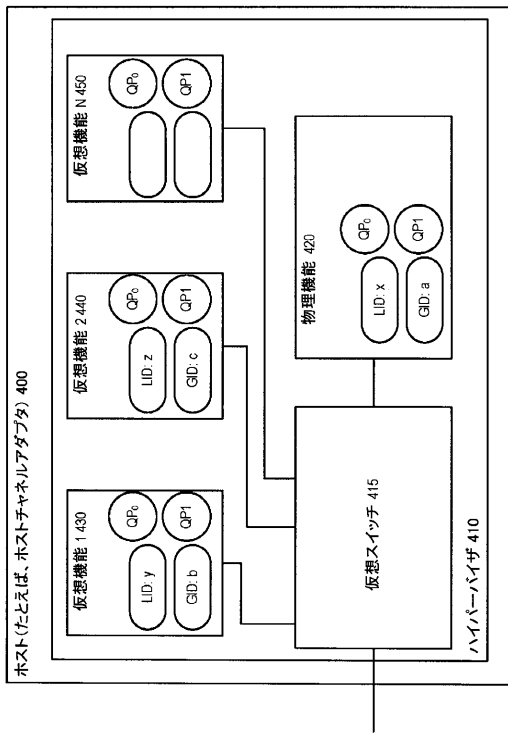


FIGURE 5

【 図 6 】

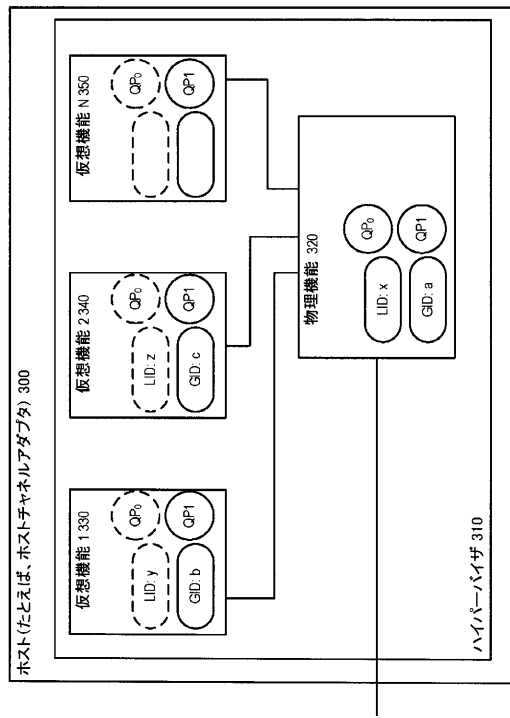


FIGURE 6

【 図 7 】

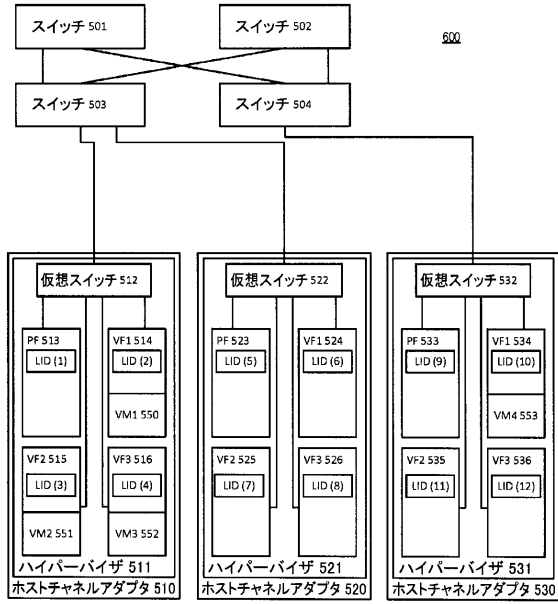


FIGURE 7

【 図 8 】

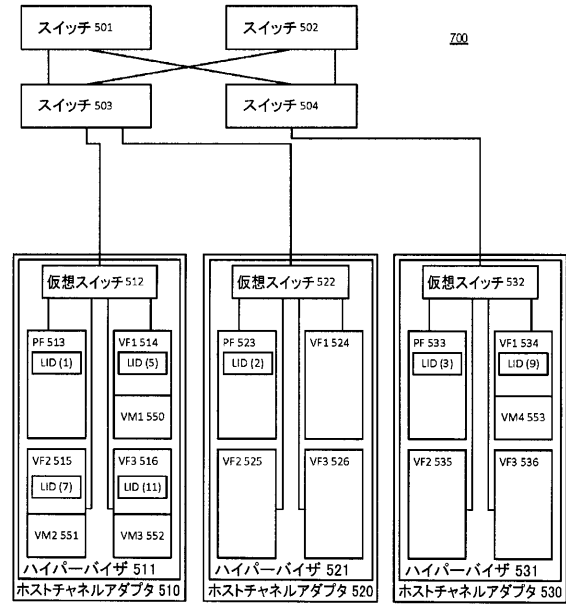


FIGURE 8

【 図 9 】

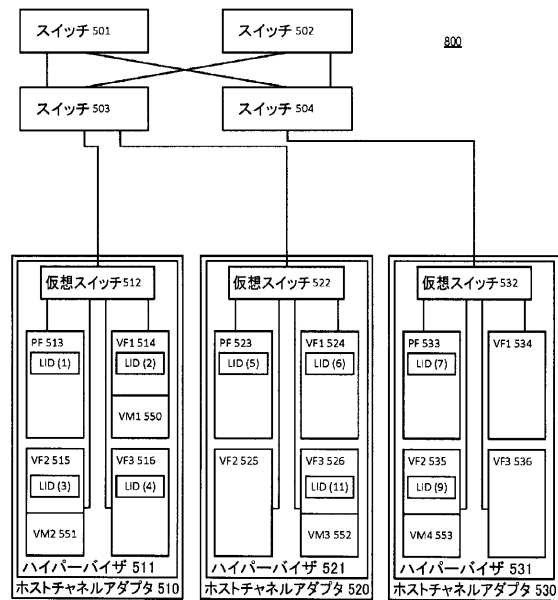


FIGURE 9

【 図 10 】

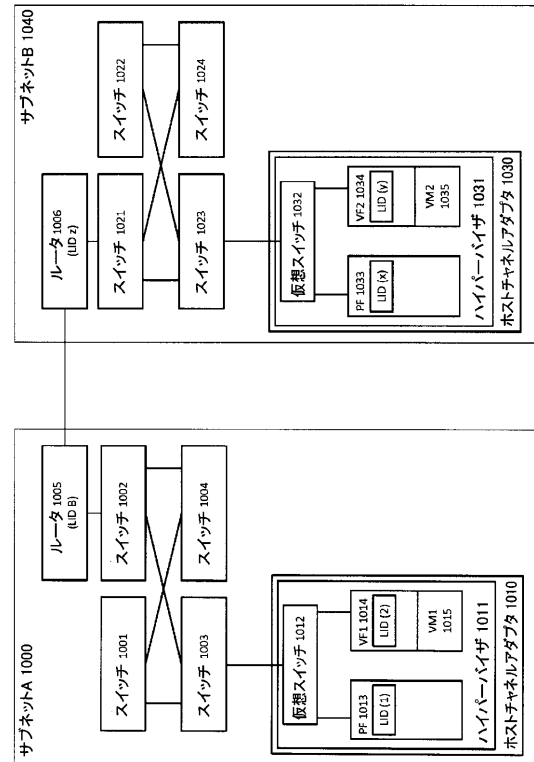


FIGURE 10

【 図 1 1 】

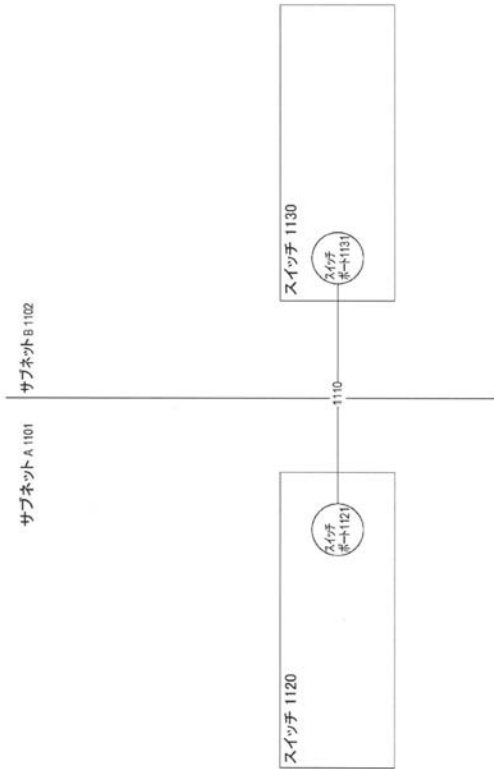


FIGURE 11

【 図 1 2 】

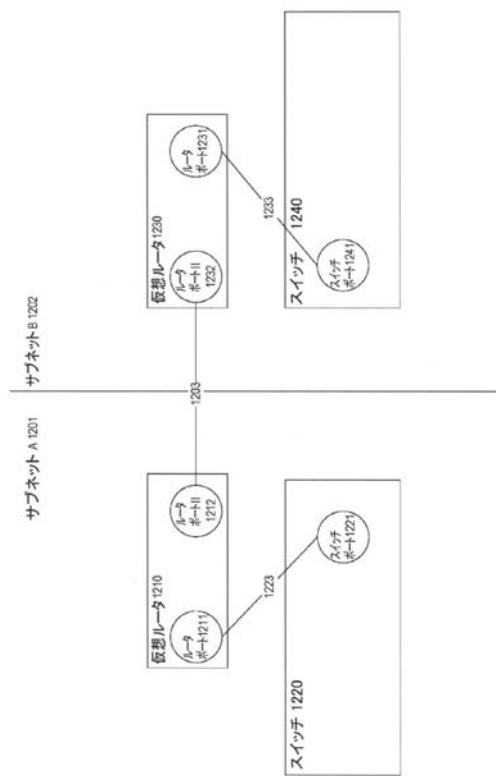


FIGURE 12

【 図 1 3 】

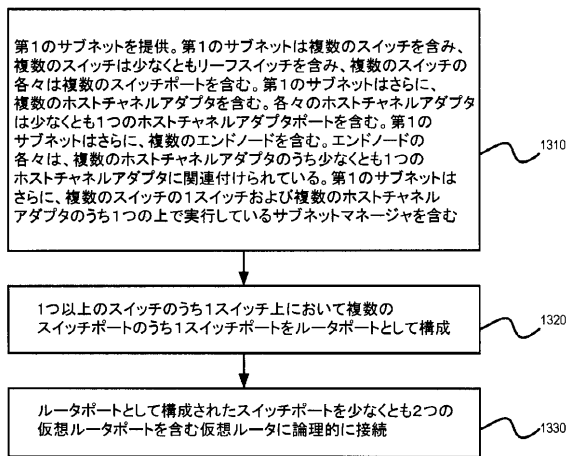


FIGURE 13

【 図 1 4 】

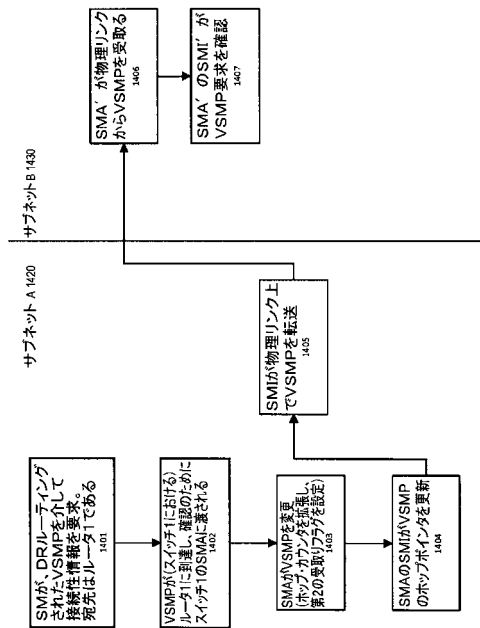


FIGURE 14

【 図 1 5 】

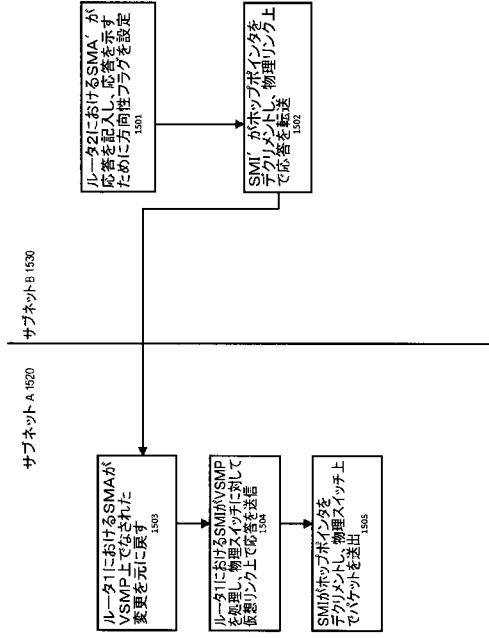


FIGURE 15

【 図 1 6 】

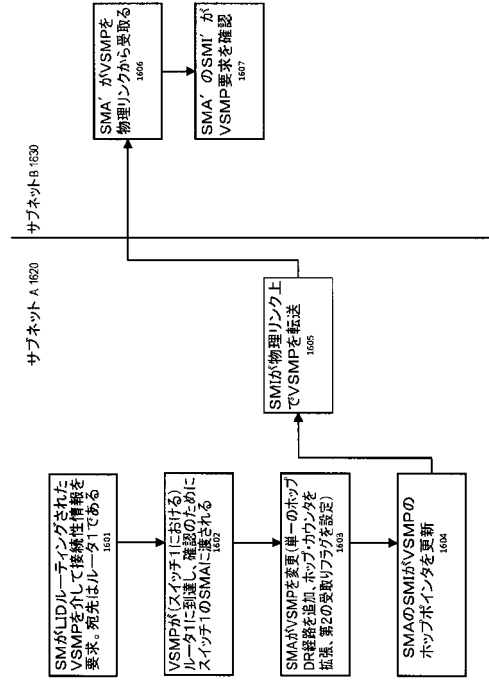


FIGURE 16

【 図 1 7 】

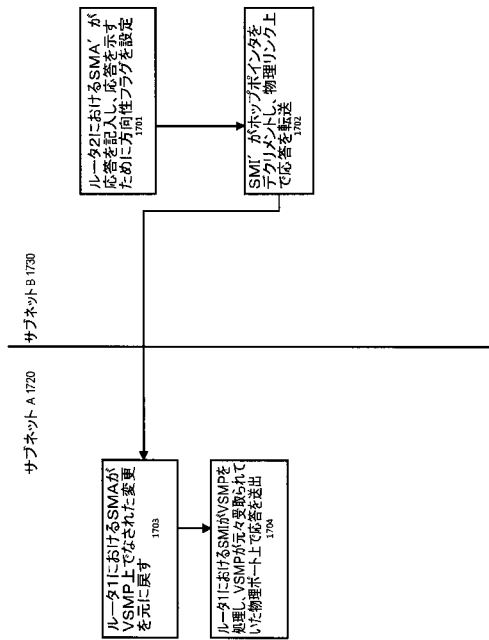


FIGURE 17

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No PCT/US2017/014959

A. CLASSIFICATION OF SUBJECT MATTER INV. H04L12/931 H04L12/733 H04L12/713 H04L12/751 H04L12/721 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 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, INSPEC, 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	US 9 219 718 B2 (JOHNSEN BJÖRN DAG [NO] ET AL) 22 December 2015 (2015-12-22) figures 1-3 column 1, line 39 - line 43 column 1, line 61 - column 2, line 4 column 2, line 23 - line 67 column 3, line 27 - line 43 column 5, line 4 - line 8 column 5, line 24 - line 48 column 6, line 23 - line 63 -----	1-22
A	US 2009/213753 A1 (BURROW STEPHEN R [US] ET AL) 27 August 2009 (2009-08-27) paragraphs [0029] - [0032], [0052] - [0057] -----	1-22
A	US 2014/269686 A1 (SRINIVASAN ARVIND [US] ET AL) 18 September 2014 (2014-09-18) the whole document -----	1-22
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> 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		Date of mailing of the international search report
25 April 2017		04/05/2017
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 Tyszka, Krzysztof

INTERNATIONAL SEARCH REPORT

Information on patent family members

International application No

PCT/US2017/014959

Patent document cited in search report	Publication date	Patent family member(s)	Publication date	
US 9219718	B2	22-12-2015	CN 103597795 A	19-02-2014
			EP 2716003 A1	09-04-2014
			JP 5965478 B2	03-08-2016
			JP 2014517406 A	17-07-2014
			US 2012307682 A1	06-12-2012
			US 2012311122 A1	06-12-2012
			US 2012311123 A1	06-12-2012
			US 2012311143 A1	06-12-2012
			US 2012311182 A1	06-12-2012
			US 2012311332 A1	06-12-2012
			US 2012311333 A1	06-12-2012
			US 2012311670 A1	06-12-2012
			US 2014241208 A1	28-08-2014
			WO 2012167268 A1	06-12-2012

US 2009213753	A1	27-08-2009	NONE	

US 2014269686	A1	18-09-2014	NONE	

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 ボグダンスキー, バルトシュ

ノルウェー、0275 オスロ、ホフ・テラッセ、15、エイチ・0203

Fターム(参考) 5K030 GA11 HB11 HD03 KX17 LB05 MA14

5K033 AA05 CB01 CB08 DB03