

(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2017-538201
(P2017-538201A)

(43) 公表日 平成29年12月21日(2017.12.21)

(51) Int. Cl.		F I		テーマコード (参考)
G06F 13/10	(2006.01)	G06F 13/10	330C	
G06F 9/50	(2006.01)	G06F 9/46	465Z	
G06F 9/46	(2006.01)	G06F 9/46	350	

審査請求 未請求 予備審査請求 未請求 (全 25 頁)

(21) 出願番号	特願2017-522887 (P2017-522887)	(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
(86) (22) 出願日	平成27年10月28日 (2015.10.28)		
(85) 翻訳文提出日	平成29年4月27日 (2017.4.27)		
(86) 国際出願番号	PCT/US2015/057860		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
(87) 国際公開番号	W02016/069773	(74) 代理人	110001195 特許業務法人深見特許事務所
(87) 国際公開日	平成28年5月6日 (2016.5.6)		
(31) 優先権主張番号	62/072,847	(72) 発明者	タソウラス, エバンジェロス ノルウェー、エヌー1325 リュサケール、 ピー・オー・ボックス・134
(32) 優先日	平成26年10月30日 (2014.10.30)		
(33) 優先権主張国	米国 (US)	(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、 ベルベルクグレンダ、9
(31) 優先権主張番号	62/075,000		
(32) 優先日	平成26年11月4日 (2014.11.4)		
(33) 優先権主張国	米国 (US)		
(31) 優先権主張番号	62/076,336		
(32) 優先日	平成26年11月6日 (2014.11.6)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 ダイナミッククラウドにサブネット管理 (SA) クエリキャッシングを提供するためのシステムおよび方法

(57) 【要約】

システムおよび方法はクラウド環境においてサブネット管理をサポートすることができる。クラウド環境における仮想マシンマイグレーション中に、サブネットマネージャは効率的なサービスを遅延させる障害ポイントになる可能性がある。システムおよび方法は、マイグレーション後に仮想マシンに複数のアドレスを確実に保持させることによって、この障害ポイントを軽減することができる。当該システムおよび方法はさらに、マイグレートされた仮想マシンとの通信を回復させるときに、仮想マシンが利用することができるローカルキャッシュに、クラウド環境内の各々のホストノードを関連付けることを可能にし得る。

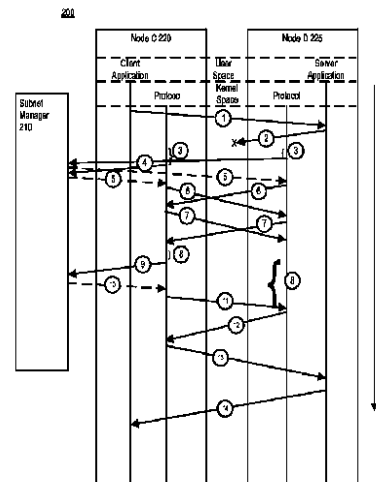


FIGURE 3

【特許請求の範囲】**【請求項 1】**

クラウド環境においてサブネット管理をサポートするための方法であって、
少なくとも第 1 のハイパーバイザおよび第 1 のホストチャネルアダプタに関連付けられた第 1 のホストノードを含む複数のホストノードを前記クラウド環境内に設けるステップと、

複数のアドレスに関連付けられた第 1 の仮想マシンを前記第 1 のホストノード上に設けるステップと、

前記第 1 の仮想マシンを、前記クラウド環境内の前記複数のホストノードのうち前記第 1 のホストノードから、設けられた第 2 のホストノードにまでマイグレートするステップとを含み、前記第 2 のホストノードは、少なくとも第 2 のハイパーバイザおよび第 2 のホストチャネルアダプタに関連付けられており、

前記複数のホストノードの各々はローカルキャッシュを含み、各々のローカルキャッシュは 1 つ以上のパス記録を含む、方法。

【請求項 2】

前記第 1 の仮想マシンを、前記クラウド環境内の前記第 1 のホストノードから前記設けられた第 2 のホストノードにまでマイグレートする前記ステップは、

前記第 1 のハイパーバイザから前記第 1 の仮想マシンを分離するステップを含み、前記第 1 のハイパーバイザから前記第 1 の仮想マシンを分離する前記ステップは、前記第 1 の仮想マシンに関連付けられた第 1 の仮想機能を前記第 1 の仮想マシンから分離するステップを含み、マイグレートする前記ステップはさらに、

前記第 1 の仮想マシンに関連付けられた前記複数のアドレスを前記第 2 のホストノードに与えるステップと、

前記第 2 のハイパーバイザに関連付けられた第 2 の仮想機能に前記複数のアドレスを割当てするステップと、

前記第 1 の仮想マシンを、前記第 1 のホストノードから前記第 2 のホストノード上の第 2 の仮想マシンにまでマイグレートするステップと、

前記第 1 の仮想マシンに関連付けられた前記複数のアドレスに前記第 2 の仮想マシンをエクスポートするステップとを含む、請求項 1 に記載の方法。

【請求項 3】

前記クラウド環境内において、前記第 1 の仮想マシンを前記第 1 のホストノードから前記設けられた第 2 のホストノードにまでマイグレートするステップの後、前記第 2 の仮想マシンと第 3 の仮想マシンとの間に通信を確立するステップをさらに含み、前記第 3 の仮想マシンは、前記複数のホストノードのうち第 3 のホストノード上に設けられ、

前記クラウド環境内において前記第 1 の仮想マシンを前記第 1 のホストノードから前記設けられた第 2 のホストノードにまでマイグレートするステップの前に、前記第 1 の仮想マシンと前記第 3 の仮想マシンとが通信していた、請求項 2 に記載の方法。

【請求項 4】

前記第 2 の仮想マシンと第 3 の仮想マシンとの間に通信を確立する前記ステップは、

前記第 3 のホストノードに関連付けられたローカルキャッシュ内に、少なくとも前記第 1 の仮想マシンに関連付けられた前記複数のアドレスを含む第 1 のパス記録を記憶するステップと、

前記第 1 の仮想マシンが前記第 1 のホストノードから前記設けられた第 2 のホストノードにまでマイグレートするとき、前記第 3 のホストノードによって、通信における途切れを検出するステップと、

前記第 3 のホストノードに関連付けられた前記ローカルキャッシュにおける前記第 1 のパス記録を前記第 3 のホストノードによって検索するステップと、

前記少なくとも第 1 のパス記録に基づいて、前記第 2 の仮想マシンと第 3 の仮想マシンとの間に通信を確立するステップとを含む、請求項 3 に記載の方法。

【請求項 5】

10

20

30

40

50

前記第 1 のパス記録は、前記第 1 の仮想マシンに関連付けられるとともに前記第 3 のホストノードによってサブネットマネージャに送信される前記複数のアドレスに関するクエリに基づいて作成され、前記サブネットマネージャは前記クラウド環境に関連付けられている、請求項 4 に記載の方法。

【請求項 6】

前記第 1 のパス記録が作成された後、前記第 1 の仮想マシンに関連付けられた前記複数のアドレスに関するクエリはいずれも、それ以上、前記第 3 のホストノードによって前記サブネットマネージャに送信されない、請求項 5 に記載の方法。

【請求項 7】

前記第 1 の仮想マシンに関連付けられるとともに前記第 3 のホストノードによって前記サブネットマネージャに送信される前記複数のアドレスに関する前記クエリに応じて、前記サブネットマネージャはマーク付けされたパス記録を戻し、前記マーク付けされたパス記録はパス・キャッシング可能マークを含み、前記パス・キャッシング可能マークは、前記第 1 のパス記録が通信における前記途切れの間も持続することを示している、請求項 5 または 6 に記載の方法。

10

【請求項 8】

前記複数のホストノードのうち別のホスト上に設けられた別の仮想マシンに関連付けられた別の複数のアドレスに関する、前記サブネットマネージャに対する別のクエリに応じて、前記サブネットマネージャはマーク付けされた別のパス記録を戻し、前記マーク付けされた別のパス記録はパス・キャッシング可能マークを含み、前記パス・キャッシング可能マークは、前記別のパス記録が通信における別の途切れの間も持続することを示している、請求項 7 に記載の方法。

20

【請求項 9】

前記第 2 の仮想マシンと前記第 3 の仮想マシンとの間の前記通信はインフィニバンド・プロトコルに基づいている、請求項 3 から 8 のいずれかに記載の方法。

【請求項 10】

前記クラウド環境内において前記第 1 の仮想マシンを前記第 1 のホストノードから前記設けられた第 2 のホストノードにまでマイグレートするステップの後、前記第 2 の仮想マシンと第 3 のエンティティとの間に通信を確立するステップをさらに含み、前記第 3 のエンティティは、物理的なホスト、記憶装置、または、マイグレートされた第 1 の仮想マシンと前記インフィニバンド・プロトコルによって以前から通信している別のエンティティのうちの 1 つであり、さらに、

30

前記第 3 のエンティティに関連付けられたローカルキャッシュ内に第 1 のパス記録を記憶するステップを含み、前記第 1 のパス記録は、少なくとも、前記第 1 の仮想マシンに関連付けられた前記複数のアドレスを含み、さらに、

前記第 1 の仮想マシンが前記第 1 のホストノードから前記設けられた第 2 のホストノードにまでマイグレートするときに、前記第 3 のエンティティによって、通信における途切れを検出するステップと、

前記第 3 のエンティティに関連付けられた前記ローカルキャッシュにおける前記第 1 のパス記録を前記第 3 のエンティティによって検索するステップと、

40

少なくとも前記第 1 のパス記録に基づいて、前記第 2 の仮想マシンと前記第 3 のエンティティとの間に通信を確立するステップとを含み、

前記クラウド環境内において前記第 1 の仮想マシンを前記第 1 のホストノードから前記設けられた第 2 のホストノードにまでマイグレートするステップの前に、前記第 1 の仮想マシンと前記第 3 のエンティティとが通信していた、請求項 2 に記載の方法。

【請求項 11】

クラウド環境においてサブネット管理をサポートするシステムであって、

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上で実行されるプロセッサとを含み、前記プロセッサは、

50

少なくとも第 1 のハイパーバイザおよび第 1 のホストチャネルアダプタに関連付けられた第 1 のホストノードを含む複数のホストノードを前記クラウド環境内に設けるステップと、

複数のアドレスに関連付けられた第 1 の仮想マシンを前記第 1 のホストノード上に設けるステップと、

前記第 1 の仮想マシンを、前記クラウド環境内の前記複数のホストノードのうち前記第 1 のホストノードから、設けられた第 2 のホストノードにまでマイグレートするステップとを含む複数のステップを実行するように動作し、前記第 2 のホストノードは、少なくとも第 2 のハイパーバイザおよび第 2 のホストチャネルアダプタに関連付けられており、

前記複数のホストノードの各々はローカルキャッシュを含み、各々のローカルキャッシュは 1 つ以上のパス記録を含む、システム。

10

【請求項 1 2】

前記プロセッサはさらに、

前記第 1 のハイパーバイザから前記第 1 の仮想マシンを分離するステップを実行するように動作し、前記第 1 のハイパーバイザから前記第 1 の仮想マシンを分離する前記ステップは、前記第 1 の仮想マシンに関連付けられた第 1 の仮想機能を前記第 1 の仮想マシンから分離するステップを含み、前記プロセッサはさらに、

前記第 1 の仮想マシンに関連付けられた前記複数のアドレスを前記第 2 のホストノードに与えるステップと、

前記複数のアドレスを、前記第 2 のハイパーバイザに関連付けられた第 2 の仮想機能に割当てるステップと、

20

前記第 1 の仮想マシンを、前記第 1 のホストノードから前記第 2 のホストノード上の第 2 の仮想マシンにまでマイグレートするステップと、

前記第 1 の仮想マシンに関連付けられた前記複数のアドレスに前記第 2 の仮想マシンをエクスポートするステップとを実行するように動作する、請求項 1 1 に記載のシステム。

【請求項 1 3】

前記プロセッサはさらに、

前記クラウド環境内において、前記第 1 の仮想マシンを前記第 1 のホストノードから前記設けられた第 2 のホストノードにまでマイグレートするステップの後、前記第 2 の仮想マシンと前記複数のホストノードのうち第 3 のホストノード上に設けられた第 3 の仮想マシンとの間に通信を確立するステップを実行するように動作し、

30

前記クラウド環境内において前記第 1 の仮想マシンを前記第 1 のホストノードから前記設けられた第 2 のホストノードにまでマイグレートするステップの前に、前記第 1 の仮想マシンと前記第 3 の仮想マシンとが通信していた、請求項 1 2 に記載のシステム。

【請求項 1 4】

前記プロセッサはさらに、

前記第 3 のホストノードに関連付けられたローカルキャッシュ内に第 1 のパス記録を記憶するステップを実行するように動作し、前記第 1 のパス記録は、少なくとも前記第 1 の仮想マシンに関連付けられた前記複数のアドレスを含み、前記プロセッサはさらに、

前記第 1 の仮想マシンが前記第 1 のホストノードから前記設けられた第 2 のホストノードにまでマイグレートするときに、前記第 3 のホストノードによって、通信における途切れを検出するステップと、

40

前記第 3 のホストノードに関連付けられた前記ローカルキャッシュにおける前記第 1 のパス記録を前記第 3 のホストノードによって検索するステップと、

少なくとも前記第 1 のパス記録に基づいて、前記第 2 の仮想マシンと第 3 の仮想マシンとの間に通信を確立するステップとを実行するように動作する、請求項 1 3 に記載のシステム。

【請求項 1 5】

前記第 1 のパス記録は、前記第 1 の仮想マシンに関連付けられるとともに前記第 3 のホストノードによってサブネットマネージャに送信される前記複数のアドレスに関するクエ

50

りに基づいて作成され、前記サブネットマネージャは前記クラウド環境に関連付けられている、請求項 14 に記載のシステム。

【請求項 16】

前記第 1 のパス記録が作成された後、前記第 1 の仮想マシンに関連付けられた前記複数のアドレスに関するクエリはいずれも、それ以上、前記第 3 のホストノードによって前記サブネットマネージャに送信されない、請求項 15 に記載のシステム。

【請求項 17】

前記第 2 の仮想マシンと前記第 3 の仮想マシンとの間の前記通信はインフィニバンド・プロトコルに基づいている、請求項 13 から 16 のいずれかに記載のシステム。

【請求項 18】

クラウド環境においてサブネット管理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機械読取可能な記憶媒体であって、前記命令は、実行されると、システムに、

少なくとも第 1 のハイパーバイザおよび第 1 のホストチャネルアダプタに関連付けられた第 1 のホストノードを含む複数のホストノードを前記クラウド環境内に設けるステップと、

複数のアドレスに関連付けられた第 1 の仮想マシンを前記第 1 のホストノード上に設けるステップと、

前記クラウド環境内において、前記第 1 の仮想マシンを、前記複数のホストノードのうちの前記第 1 のホストノードから、設けられた第 2 のホストノードにまでマイグレートするステップとを含む複数のステップを実行させ、前記第 2 のホストノードは、少なくとも

前記複数のホストノードの各々はローカルキャッシュを含み、各々のローカルキャッシュは 1 つ以上のパス記録を含む、非一時的な機械読取可能な記憶媒体。

【請求項 19】

前記複数のステップはさらに、

前記第 1 のハイパーバイザから前記第 1 の仮想マシンを分離するステップを含み、前記第 1 のハイパーバイザから前記第 1 の仮想マシンを分離する前記ステップは、前記第 1 の仮想マシンに関連付けられた第 1 の仮想機能を前記第 1 の仮想マシンから分離するステップを含み、前記複数のステップはさらに、

前記第 1 の仮想マシンに関連付けられた前記複数のアドレスを前記第 2 のホストノード

に与えるステップと、

前記第 2 のハイパーバイザに関連付けられた第 2 の仮想機能に前記複数のアドレスを割当てるステップと、

前記第 1 の仮想マシンを、前記第 1 のホストノードから前記第 2 のホストノード上の第 2 の仮想マシンにまでマイグレートするステップと、

前記第 2 の仮想マシンを、前記第 1 の仮想マシンに関連付けられた前記複数のアドレスにエクスポートするステップとを含む、請求項 18 に記載の非一時的な機械読取可能な記憶媒体。

【請求項 20】

前記複数のステップはさらに、

前記クラウド環境内において、前記第 1 の仮想マシンを、前記第 1 のホストノードから前記設けられた第 2 のホストノードにまでマイグレートするステップの後、前記第 2 の仮想マシンと第 3 の仮想マシンとの間に通信を確立するステップを含み、前記第 3 の仮想マシンは、前記複数のホストノードのうち第 3 のホストノード上に設けられ、

前記クラウド環境内において前記第 1 の仮想マシンを前記第 1 のホストノードから前記設けられた第 2 のホストノードにまでマイグレートするステップの前に、前記第 1 の仮想マシンと前記第 3 の仮想マシンとが通信していた、請求項 19 に記載の非一時的な機械読取可能な記憶媒体。

【請求項 21】

前記複数のステップはさらに、

10

20

30

40

50

前記第 3 のホストノードに関連付けられたローカルキャッシュ内に第 1 のパス記録を記憶するステップを含み、前記第 1 のパス記録は、少なくとも前記第 1 の仮想マシンに関連付けられた前記複数のアドレスを含み、前記複数のステップはさらに、

前記第 1 の仮想マシンが前記第 1 のホストノードから前記設けられた第 2 のホストノードにまでマイグレートするときに、前記第 3 のホストノードによって、通信における途切れを検出するステップと、

前記第 3 のホストノードに関連付けられた前記ローカルキャッシュにおける前記第 1 のパス記録を前記第 3 のホストノードによって検索するステップと、

少なくとも前記第 1 のパス記録に基づいて、前記第 2 の仮想マシンと第 3 の仮想マシンとの間に通信を確立するステップとを含む、請求項 20 に記載の非一時的な機械読取可能な記憶媒体。

10

【請求項 22】

前記第 1 のパス記録は、前記第 1 の仮想マシンに関連付けられるとともに前記第 3 のホストノードによってサブネットマネージャによって送信される前記複数のアドレスに関するクエリに基づいて作成され、前記サブネットマネージャは前記クラウド環境に関連付けられており、

前記第 1 のパス記録が作成された後、前記第 1 の仮想マシンに関連付けられた前記複数のアドレスに関するクエリはいずれも、それ以上、前記第 3 のホストノードによって前記サブネットマネージャに送信されない、請求項 21 に記載の非一時的な機械読取可能な記憶媒体。

20

【発明の詳細な説明】

【技術分野】

【0001】

著作権表示：

この特許文献の開示の一部は、著作権保護の対象となる題材を含んでいる。著作権の所有者は、特許商標庁の包袋または記録に掲載されるように特許文献または特許情報開示を誰でも複製できることに対して異議はないが、その他の点ではすべての如何なる著作権をも保有する。

【0002】

発明の分野：

30

本発明は概してコンピュータシステムに関し、特にクラウド環境に関する。

【背景技術】

【0003】

背景：

インフィニバンド (InfiniBand) サブネットにおいては、サブネットマネージャ (Subnet Manager : S M) が潜在的な障害となる。インフィニバンド・サブネットのサイズが大きくなると、ホスト間のパスの数が多項式的に増加し、多くの同時発生的なパス決定の要求を受取ったときに S M がタイムリーにネットワークのために機能することができなくなるかもしれない。このスケーラビリティの問題は、ダイナミック仮想化クラウド環境においてさらに増幅される。インフィニバンド相互接続された仮想マシン (Virtual Machine : V M) がライブマイグレートする (live migrate) と、 V M アドレスが変更される。これらのアドレス変更により、結果として、通信ピアが新しいパス特徴を解釈するためにサブネット管理 (Subnet Administration : S A) パス記録クエリを S M に送信したときに、 S M に負荷が追加されることとなる。

40

【発明の概要】

【課題を解決するための手段】

【0004】

概要：

システムおよび方法は、クラウド環境においてサブネット管理をサポートすることができる。クラウド環境における仮想マシンマイグレーション (migration) 中に、サブネッ

50

トマネージャは効率的なサービスを遅延させる障害ポイントになる可能性がある。システムおよび方法は、マイグレーション後に仮想マシンに確実に複数のアドレスを保持させることによって、この障害ポイントを軽減させることができる。当該システムおよび方法は、マイグレートされた仮想マシンとの通信を回復させるときに、クラウド環境内の各々のホストノードを、仮想マシンが利用できるローカルキャッシュに関連付けることをさらに可能にし得る。

【図面の簡単な説明】

【0005】

【図1】一実施形態に従った、クラウド環境におけるVMライブマイグレーションをサポートする例を示す図である。

【図2】一実施形態に従った、2つのホスト間の接続を確立するプロトコルの例を示す図である。

【図3】一実施形態に従った、接続の途絶時に2つのノード間で継続中の通信の例を示す図である。

【図4】一実施形態に従った、クラウド環境においてSAパス・キャッシングをサポートする例を示す図である。

【図5】本発明の一実施形態に従った、クラウド環境においてSAパス・キャッシングをサポートする例を示す図である。

【図6】一実施形態に従った、クラウド環境においてサブネット管理をサポートするための方法の例を示す図である。

【発明を実施するための形態】

【0006】

詳細な説明：

本発明は、同様の参照符号で同様の要素を示している添付の図面において、限定によってではなく例示によって説明されている。この開示において「ある」または「1つの」または「いくつかの」実施形態を言及している場合、このような言及は必ずしも同じ実施形態に対するものではなく、「少なくとも1つ」という意味である。

【0007】

本発明の以下の記載は、高性能ネットワークプロトコルについての一例としてインフィニバンド（IB）ネットワークプロトコルを用いている。他のタイプの高性能ネットワークプロトコルが制限なしに使用可能であることは当業者にとって明らかであるだろう。

【0008】

クラウド環境において仮想マシン（virtual machine：VM）マイグレーション・サブネット管理（subnet administration：SA）パス・キャッシングをサポートすることができるシステムおよび方法がこの明細書中に記載される。

【0009】

一実施形態に従うと、高性能コンピューティング（High Performance Computing：HPC）を提供することができるクラウドコンピューティングを提供することができる。HPC・アズ・ア・サービス（HPC-as-a-service）をコンピューティングクラウド内に提供することができる。このHPC・アズ・ア・サービスは仮想HPC（virtual HPC：vHPC）クラスタを考慮に入れたものであって、高性能相互接続ソリューションを用いてこの仮想HPCクラスタに対応し得る。

【0010】

一実施形態に従うと、各々のIBサブネットはサブネットマネージャ（SM）を利用することができる。各々のSMは、ネットワーク初期化、トポロジー発見、パス演算、ならびに、スイッチおよびホストチャネルアダプタ（Host Channel Adapter：HCA）上のIBポートの構成のために機能し得る。大規模なサブネットでは、ノード間で利用可能なパスが多項式的に大きくなる可能性があり、パス決定のための多くの同時発生的な要求を受取ったときに、SMが潜在的な障害になる可能性がある。このスケーラビリティの問題は、ダイナミックな仮想化クラウド環境においては、IB相互接続されている仮想マシンが

10

20

30

40

50

ライブマイグレートしたときにさらに増幅される。

【0011】

効率的な仮想化をサポートするために、高帯域および低レイテンシを維持しながらも、IBホストチャネルアダプタ(HCA)はシングルルートI/O仮想化(Single Root I/O Virtualization: SR-IOV)をサポートすることができる。IB接続されたノードは各々、3つのさまざまなアドレスを有する。ライブマイグレーションが発生すると、パススルーされたインターフェイス(the passed through interface)を取外すことでダウンタイムが発生するにもかかわらず、IBアドレスのうち1つ以上を変更することができる。マイグレーション中のVMと通信する他のノードは接続性を失い、サブネット管理(SA)パス記録クエリをSMに送信することによって再接続するための新しいアドレスを発見しようと試みる。基盤ネットワークにおいて結果として生じるSM宛ての通信は有効になり得る。大規模なネットワークにおいては、SM宛てのメッセージがこのようにVMのマイグレートによって溢れかえると、SMに対する負荷が増えるのに応じて、全体的なネットワークレイテンシが大きくなる可能性がある。

10

【0012】

一実施形態に従うと、さらに、VMマイグレーションによって引き起こされてSMによって受取られるSA要求の量を減らすことによって、SMに対する負荷を減らすことが望ましい。方法およびシステムは、VMがマイグレーション後にその同じアドレスを保持することができるシステムを実現することによって、これを達成することができる。加えて、2つのノード間の初期接続が確立された後に、SAパス・キャッシングメカニズムを用いてSAクエリの数を激減させることができる。

20

【0013】

一実施形態に従うと、インフィニバンドは一般に3つの異なるタイプのアドレスを用いている。まず、16ビットのローカル識別子(Local Identifier: LID)が挙げられる。少なくとも1つの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り当てられている。LIDを用いて、サブネット内のトラフィックをルーティングすることができる。LIDが16ビット長であるので、65536個の固有のアドレスの組合せを構成することができ、そのうち49151(0x0001-0xBFFF)個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、利用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。

30

【0014】

第2のタイプのアドレスは、一般に、製造業者によって各々の装置(たとえば、HCAおよびスイッチ)ならびに各々のHCAポートに割り当てられる64ビットのグローバル固有識別子(Global Unique Identifier: GUID)である。SMは、SR-IOV VFが使用可能にされたときに有用になり得る追加のサブネット固有のGUIDをHCAポートに割り当て得る。

【0015】

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(Global Identifier: GUID)である。GUIDは、一般に、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ポートおよび各々のスイッチに割り当てられている。GUIDは、ファブリックアドミニストレータによって割り当てられたグローバルに固有の64ビットプレフィックスと各々のHCAポートのGUIDアドレスとを組み合わせることによって形成される。

40

【0016】

本発明の以下の記載は、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ネットワークを用いる。他のタイプの高性能ネットワークを制限なしで使用できることは当業者にとって明らかであるだろう。また、本発明の以下の記載は、仮想化モデルについての一例としてKVM仮想化モデルを用いている。他のタイプの仮想化モデル(たとえばXen)を制限なしで使用できることは当業者にとって明らかであるだろう。

【0017】

50

本発明の以下の記載は、加えて、OpenStack、OpenSMおよびRDSELinux（登録商標）カーネルモジュールを利用する。OpenStackは、データセンターを通じて処理、記憶およびネットワークリソースのプールを制御する相互に関係付けられたプロジェクトのグループを含むクラウド・コンピューティング・ソフトウェアプラットフォームである。OpenSMは、OpenIBの上で実行させることができるインフィニバンド対応のサブネットマネージャおよびアドミニストレーションである。リライアブル・データグラム・ソケット（Reliable Datagram Socket：RDS）は、データグラムを提供するための高性能で低レイテンシかつ信頼性のある非接続プロトコルである。他の同様のプラットフォームを制限なしで利用できることは当業者にとって明らかであるだろう。

10

【0018】

本発明の一実施形態に従うと、仮想化は、クラウドコンピューティングにおける効率的なリソース利用および融通性のあるリソース割当てに有益であり得る。ライブマイグレーションは、アプリケーションにトランスペアレントな態様で物理サーバ間で仮想マシン（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。これにより、シングルルートI/O仮想化（SR-IOV）法を利用する仮想化は、ライブマイグレーションによって統合、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。

【0019】

IBアーキテクチャは、直列なポイント・ツー・ポイントの全二重技術である。IBネットワークはサブネットとも称することができ、この場合、サブネットは、スイッチおよびポイント・ツー・ポイントリンクを用いて相互接続された1セットのホストで構成されている。IBサブネットは、サブネットにすべてのスイッチ、ルータおよびホストチャネルアダプタ（host channel adaptor：HCA）の構成を含むネットワークを初期化して提供する役割を果たし得る少なくとも1つのサブネット・マネージャ（SM）を含み得る。

20

【0020】

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をとともに提供するために、1組のリッチな伝送サービスをサポートする。用いられる伝送サービスに依存せずに、IB HCAはキュー対（queue pair：QP）を用いて通信を行う。QPは、通信設定中に作成され、提供されるQP数、HCAポート、宛先LID、キューサイズおよび伝送サービスなどの1セットの初期属性を有し得る。HCAは多くのQPを処理することができ、各々のQPは、送信キュー（send queue：SQ）および受信キュー（receive queue：RQ）などの1対のキューからなる。通信に参与している各々のエンドノードにはこのような1つの対が存在している。送信キューは、リモートノードに転送されるべき作業要求を保持し、受信キューは、リモートノードから受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している。QPに加えて、各々のHCAは、1セットの送信キューおよび受信キューに関連付けられている1つ以上の完了キュー（completion queue：CQ）を有している。CQは、送信キューおよび受信キューに掲示された作業要求についての完了通知を保持する。たとえば通信の複雑さがユーザからは見えないようになっていたとしても、QP状態情報はHCAに保持されている。

30

40

【0021】**ネットワークI/O仮想化：**

一実施形態に従うと、I/O仮想化（I/O virtualization：IOV）を用いることにより、I/Oリソースを共有して、さまざまな仮想マシンからリソースへのアクセスの保護を提供することができる。IOVは、仮想マシンにエクスポーズされ得る論理装置をその物理的な実装から分離することができる。このような1タイプのIOVとして、直接的な装置の割当てがある。

【0022】

一実施形態に従うと、直接的な装置の割当ては、VM間で装置を共有せずにI/O装置

50

をVMに連結することを必要とし得る。直接的な割当て（または装置のパススルー）は最小限のオーバーヘッドで固有の性能に近いものを提供することができる。物理装置は、ハイパーバイザをバイパスしてVMに直接取付けられており、ゲストOSは未改良のドライバを用いることができる。不利な面としては、共有がなされないせいでスケーラビリティが制限されてしまう点であり、1枚の物理ネットワークカードが1つのVMと連結されている。

【0023】

一実施形態に従うと、シングルルートIOV（Single Root IOV：SR-IOV）は、ハードウェア仮想化によって、物理装置がその同じ装置の複数の独立した軽量のインスタンスとして現われることを可能にし得る。これらのインスタンスは、パススルー装置としてVMに割当てることができ、仮想機能（Virtual Function：VF）としてアクセスすることができる。SR-IOVは、純粋に直接的に割当てする際のスケーラビリティの問題を軽減する。

10

【0024】

残念ながら、SR-IOVなどの直接的な装置の割当て技術によれば、実現されたシステムがデータセンタ最適化のためにトランスペアレントなライブマイグレーション（VMマイグレーション）を使用する場合、クラウドプロバイダに対して問題が提起される可能性がある。ライブマイグレーションの本質は、仮想マシンのメモリ内容が遠隔のハイパーバイザにまでコピーされるという点にある。さらに、仮想マシンをソース・ハイパーバイザにおいて一時停止させると、その動作は、それがコピーされた宛先において再開される。基礎をなすシステムが（SR-IOVなどの）直接的な装置の割当てを利用する場合、ネットワークインターフェイスがハードウェアに接続されている場合にはそのネットワークインターフェイスの完全な内部状態をコピーすることができない。VMに割当てられたSR-IOV VFが分離されると、ライブマイグレーションが実行されることとなり、その宛先において新しいVFが取付けられることとなる。

20

【0025】

IB VFを用いるVMがライブマイグレートされている状況においては、VMの3つのアドレスをすべて変更することにより、基礎をなすネットワークファブリックおよびSMに対する明らかな影響を取り入れることができる。VMが異なるLIDを有する異なる物理ホストに移されるので、LIDが変化する。宛先において異なるVFが取付けられると、SMによってソースVFに割当てられている仮想GUID（virtual GUID：vGUID）が同様に变化し得る。次に、vGUIDを用いてGUIDが形成されるので、GUIDも変化するだろう。結果として、マイグレートされたVMが、突然、新しい1セットのアドレスに関連付けられ、マイグレートされたVMの通信ピアが、マイグレートされたVMとの失われた接続を回復しようとして、同時発生したSAパス記録クエリバーストをSMに送信し始めることができる。これらのクエリは、SMに対する余分なオーバーヘッドと、副次的作用としての追加のダウンタイムとをもたらす可能性がある。マイグレートされたノードがネットワークにおける他の多くのノードと通信する場合、SMが障害となり、全体的なネットワーク性能を阻害する可能性がある。

30

【0026】

一実施形態に従うと、この明細書中に記載される方法およびシステムは、クラウドプロバイダに提示されているSR-IOVなどの直接的な装置の割当て技術を用いて、仮想マシンのライブマイグレーションに関連付けられた問題を低減および/または排除することができる。当該方法およびシステムは、IB VFを用いるVMがライブマイグレートされている状況において提示される問題を克服することができる。

40

【0027】

仮想マシン（VM）ライブマイグレーション

図1は、一実施形態に従った、クラウド環境におけるVMのライブマイグレーションをサポートする例を示す。図1に示されるように、インフィニバンド（IB）サブネット100は、異なるハイパーバイザ111～113をサポートする複数のホストノードA10

50

1 ~ C 1 0 3 を含み得る。

【 0 0 2 8 】

加えて、各々のハイパーバイザ 1 1 1 ~ 1 1 3 は、さまざまな仮想マシン (V M) がその上で実行されることを可能にする。たとえば、ホストノード A 1 0 1 上のハイパーバイザ 1 1 1 は V M A 1 0 4 をサポートすることができ、ホストノード B 上のハイパーバイザ 1 1 2 は V M B 1 0 5 をサポートすることができる。 V M A および V M B が実行されているノード同士は通信することができる。

【 0 0 2 9 】

さらに、ホストノード A 1 0 1 ~ C 1 0 3 の各々は 1 つ以上のホストチャネルアダプタ (H C A) 1 1 7 ~ 1 1 9 に関連付けることができる。図 1 に示されるように、ホストノード A 1 0 1 上の H C A 1 1 7 は、 V M A 1 0 4 によって使用され得る Q P a 1 0 8 などのキュー対 (Q P) を利用することができ、ホストノード B 1 0 2 上の H C A 1 1 8 は、 V M B 1 0 5 によって使用され得る Q P b 1 0 7 を利用することができる。

10

【 0 0 3 0 】

本発明の一実施形態に従うと、入出力仮想化 (input/output virtualization : I O V) を用いて、 V M に I / O リソースを提供し、複数の V M から共有の I / O リソースへのアクセスの保護を提供することができる。 I O V は論理装置を分離することができ、この論理装置は、その物理的実装から V M にエクスポーズされている (exposed)。たとえば、シングルルート I / O 仮想化 (S R - I O V) は、 I B ネットワークにわたる仮想化において高性能を達成するための I / O 仮想化アプローチである。

20

【 0 0 3 1 】

また、 I B サブネット 1 0 0 は、ネットワーク初期化、スイッチおよび H C A 上の I B ポートの構成、トポロジ発見およびパス演算のために機能し得るサブネットマネージャ 1 1 0 を含み得る。

【 0 0 3 2 】

図 1 に示されるように、 V M B 1 0 5 は、 (たとえば、ハイパーバイザ 1 1 1 上の V M A 1 0 5 と通信しながら) ハイパーバイザ 1 1 2 からハイパーバイザ 1 1 3 へとマイグレートさせることができる。

【 0 0 3 3 】

マイグレーションの後、新しい V M B 1 0 6 は、宛先ホストノード C 1 0 3 における新しい 1 セットのアドレスに対して突然エクスポーズされる可能性がある。さらに、ピア V M (たとえば、 V M A 1 0 4) は、失われた接続を回復させようとしている間に、 S M 1 1 0 にサブネット管理 (S A) パス記録クエリを送信し始めることができる (V M B はまた、新しいホストノード上で実行されると、 S A パス要求を S M に送信することもできる)。これは以下の事実起因している。すなわち、一般に、 V M B がホストノード B からホストノード C にマイグレートする場合などのように V M がマイグレートすると、これに応じて、 V M のアドレス (L I D 、 G U I D 、 G I D) が変化する。というのも、これら V M のアドレス (L I D 、 G U I D 、 G I D) は、概して、 S R - I O V を用いているときにはハードウェアに接続されているからである。サブネットマネージャへのこれらの S A パスクエリは、著しいダウンタイムと、インフィニバンド S M 1 1 0 に対する余分なオーバーヘッドとをもたらす可能性がある。多くのマイグレーションが大規模なデータセンタにおいてかなり短い時間フレーム内で行なわれる場合、または、マイグレートされたノードがネットワークにおける他の多くのノードと通信している場合、 S M 1 1 0 が障害になる可能性がある。なぜなら、この S M 1 1 0 がタイムリーに応答することができないかもしれないからである。

30

40

【 0 0 3 4 】

本発明の一実施形態に従うと、当該システムは、 V M B 1 0 4 がマイグレートして I B アドレス情報が変化した場合に、関与しているホストノード A 1 0 1 ~ C 1 0 3 によって生成される S A クエリの量を減らすことができる。

【 0 0 3 5 】

50

図 1 に示されるように、当該システムは、最初に、たとえば、VM B 1 0 4 から仮想機能 (virtual function: VF) 1 1 5 を分離することによって、ハイパーバイザ 1 1 2 から VM B 1 0 4 を分離することができる。次いで、システムは、たとえば、これらのアドレスを、ホストノード C 1 0 3 上のハイパーバイザ 1 1 3 上において次に利用可能な仮想機能 (すなわち VF 1 1 6) に割当てることによって、VM B 1 0 4 に関連付けられているアドレス情報 1 2 0 を宛先ホストノード C 1 0 3 に与えることができる。最後に、VM B 1 0 4 が VM B 1 0 6 としてハイパーバイザ 1 1 3 にマイグレートされた後、当該システムは、(たとえば、QPb 1 0 9 を介する) ピア VM との通信を回復させるために、VM B 1 0 6 をアドレス情報 1 2 0 にエクスポートすることができる。

10

【0036】

これにより、宛先ホストノード C 1 0 3 へのマイグレーションの後、新しい VM B 1 0 6 を元の 1 セットのアドレスにエクスポートすることができ、ピア VM A 1 0 4 が、SAパス記録クエリを SM 1 1 0 に送信する必要がなくなる。

【0037】

一実施形態に従うと、システムは、IB SR-IOV VF が取付けられている VM の VM ライブマイグレーションをサポートすることができる。リモートダイレクトメモリアクセス (RDMA) は、VM のマイグレーション後に通信を回復させるために、リアル・データグラム・ソケット (RDS) プロトコルなどのプロトコルによって利用することができる。

20

【0038】

一実施形態に従うと、システムは、OpenStack、OpenSM および RDS Linux カーネルモジュールを利用することができる。加えて、LID tracker と称され得るプログラムは、各々の VM に関連付けられた IB アドレスを常に把握しておくために用いることができ、マイグレーションプロセスを編成することができる。

【0039】

一実施形態においては、プログラムは、OpenSM のオプション honor_guid2lid_file を使用可能にし得る。次いで、OpenSM によって生成されたファイル guid2lid は、プログラムによってパースされ、昇順などの順序で GUID によってソートされ得る。LID は、1 から開始して GUID に割当てられる。GUID に割当てられた各々の LID は、物理ホストのためのベース LID とも称することができる。

30

【0040】

一実施形態においては、ベース LID が割当てられると、IB 対応の OpenStack 演算ノードの各々は、VM を実行させるためにスキャンすることができる。実行されていることが判明した各々の VM には、49151 (最上位のユニキャスト LID) から開始して降順で LID が割当てられ得る。VM に割当てられたこれらの LID は浮動 LID とも称され得る。

【0041】

一実施形態においては、浮動 LID は、VM が実行されている OpenStack 演算ノードにおけるベース LID と置換えることができる。ハイパーバイザは LID を VM と共有する。いくつかの実施形態においては、1 つの VM をハイパーバイザごとに実行させることができ、他の VM がその時点で実行されていないハイパーバイザに対して、或る VM がマイグレートされ得る。他の実施形態においては、複数の VM をハイパーバイザ上で実行させることができ、他の VM が宛先ハイパーバイザ上でその時点で実行されているかどうかにかかわらず、或る VM を別のハイパーバイザにマイグレートすることができる。

40

【0042】

一実施形態においては、VM_x のためのマイグレーションが OpenStack API などの API から指示されると、SR-IOV VF を VM から分離することができる。装置の取外しが完了し、マイグレーションが進行中である場合、OpenStack は、プログラムに対し、Hypervisor_y などの 1 つのハイパーバイザから Hype

50

r v i s o r_zなどの宛先ハイパーバイザにまでVM_xが移動していることを通知することができる。次いで、プログラムは、Hyper v i s o r_yのL I DをそのベースL I Dに戻すことができ、Hyper v i s o r_zはVM_xに関連付けられた浮動L I Dを得ることができる。プログラムはまた、宛先ハイパーバイザであるHyper v i s o r_zにおいて、VM_xに関連付けられたv G U I Dを次に利用可能なS R - I O V V Fに割当てることができる。マイグレーション中、VMにはネットワーク接続性がない。

【0043】

一実施形態に従うと、これらの変更は再始動によって加えることができる。さらに、マイグレーションが完了すると、Open S t a c kは次に利用可能なS R - I O V V FをHyper v i s o r_z上のVM_xに追加することができ、VMはそのネットワーク接続性を復帰させることができる。VMは、マイグレーションの前に有していた同じI Bアドレス(L I D、v G U I DおよびG I D)にエクスポートすることができる。VMの視点からは、マイグレートするのに必要な時間だけI Bアダプタが分離されており、かつ、アドレスが変化しなかったため同じI Bアダプタが再び取付けられたように、見えている。

10

【0044】

サブネット管理(S A)パス・キャッシング

一実施形態に従うと、エンドノードにおけるローカルのS Aパス・キャッシングメカニズムは、初期接続が2つのノード間で確立された後にS Aクエリを減らすかまたは排除することができる。キャッシングスキームは一般的なものであり得るとともに、使用可能にされると、ライブマイグレーションの実施の有無にかかわらず、S Mに対する負荷を軽減させることができる。

20

【0045】

図2は、一実施形態に従った、2つのホスト間の接続を確立するプロトコルの例を示す。より特定的には、図2は、2つのホスト間の接続を確立するためにR D Sなどのプロトコルを用いることを例示している。

【0046】

一実施形態に従うと、接続を確立する前に、I P・オーバー・I B(I P o v e r I B: I P o I B)をすべての通信ピアにおいて設定することができる。R D Sなどのプロトコルは、特有のI BポートのI P o I Bアドレスを用いてポートのG I Dアドレスを決定することができる。G I Dアドレスが決定された後には、プロトコルは、パス記録照合を実行してI B通信を確立するのに十分な情報を有することができる。

30

【0047】

図2に示されるように、インフィニバンド・サブネット200内では、サブネットマネージャ210は、ノードC220とノードD225との間、より特定的には、ノードC上のクライアント側アプリケーションとノードD上のサーバ側アプリケーションとの間、にパス通信を提供することができる。図2においては、上位層アプリケーションのクライアント側がノードCにおいて実行され、アプリケーションのサーバ側はノードDにおいて実行される。アプリケーションのクライアント側は、R D Sソケットなどのソケットを作成して、アプリケーションのサーバ側と通信するよう試みる(ステップ1)ことができる。R D Sなどのプロトコルは、S Aパス記録要求をノードCからS Mに送信する(ステップ2)ことができる。サブネットマネージャは、プロトコルに応答を提供する(ステップ3)ことができる。この応答は、クライアント側アプリケーションのターゲットについてのアドレス情報を含み得る。サブネットマネージャから応答を受取った後、プロトコルは、接続要求を送信することによってノードDとの接続を開始するよう試みる(ステップ4)ことができる。接続が成功すれば、プロトコルは、たとえば、両側でのR D M A _ _ C M _ _ E V E N T _ _ E S T A B L I S H E Dイベントによって通信チャネルを確立する(ステップ5)ことができる。この時点において、上位層アプリケーションが通信可能となる(ステップ6)。

40

【0048】

50

初期接続において何かエラーがあれば、クライアント側（ノードC）のプロトコルは、ランダムなバックオフメカニズムとの接続を確立するために再試行しようとすることができる。サーバは、クライアントが通信しようとする意図していることにまだ気づいていない。接続が確率された後に何か不具合があれば、RDS側（アプリケーション視点からはクライアントおよびサーバ）がともに積極的にピアと再接続しようとするだろう。接続プロセスにおけるランダムなバックオフメカニズムは、両側が接続に関与している場合には、競合状態を回避するのに有用である。

【0049】

図3は、一実施形態に従った、接続が途絶えた場合における2つのノード間の継続中の通信の例を示す。

10

【0050】

図3においては、インフィニバンド・サブネット200内では、サブネットマネージャ210は、ノードC220とノードD225との間にパス通信を提供することができ、接続が途絶えた（ステップ2）場合にノードCとノードDとの間に通信を継続させている（ステップ1）。接続の途絶えは、たとえば、ノード上で実行されているアプリケーションのうち1つのアプリケーションのライブマイグレーションに付随して起こる可能性がある。両方のプロトコルエンドは、接続がダウンしていると判断することができ、再接続しようとする前にいくらかのランダム時間（すなわちバックオフ時間）だけ待機する（ステップ3）ことができる。再接続を試みる前に両側が待機する時間は、図3に示されるように、同じであってもよく、または異なってもよい。ノードは、SMにSAパス記録要求を送信する（ステップ4）ことによって再接続しようとすることができる。SAパス記録応答が受取られた（ステップ5）後、接続要求を送信する（ステップ6）ことができる。

20

【0051】

図3に示される場合においては、ステップ3において2つのノードによって選択されたバックオフ時間はほとんど同じであった。このため、ノードDがノードCよりもわずかに速くSAパス記録応答を得て、ステップ6において最初に接続を開始しようとしていたとしても、ノードCが接続要求自体を送信するより以前には接続要求はノードCに到達していなかった。この場合、両方のプロトコルエンドは未処理の接続要求を有している。次いで、ノードがそれらのピアから接続要求を受取ると、これらのノードは接続要求を拒否する（ステップ7）だろう。ステップ8において、2つのノードは、再接続しようとする前に、ランダムなバックオフ時間をもう一度選択した。このとき、ノードDによって選択されたランダムなバックオフ時間はノードCによって選択されたものより著しく長い。結果として、ノードCが優先されることとなり、このノードCが、接続確立プロセスを繰返し、SAパス記録要求を送信し（ステップ8）、サブネットマネージャから応答を受取り（ステップ10）、ノードDに接続要求を送信し（ステップ11）、そして、ノードDがノードCとの接続自体を開始しようとする前に接続要求がノードDに到達する。図3に示される状況においては、ノードDが入接続を受付ける（ステップ12）。次いで、ステップ13および14において、上位層アプリケーションのために通信を再開することができる。

30

40

【0052】

図3から推測すると、サブネットマネージャがVMマイグレーションの場合にSAパス要求を大量に受取る可能性がある（通信の遮断）ことが明らかになる。何千ものノードを有する大規模なサブネットにおいては、各々のノードから追加のSAクエリが1つだけ送信された場合でも、SMには、最終的に何千ものメッセージが殺到する可能性がある。ライブマイグレーションがダイナミックIBベースのクラウドにおいて行なわれる場合、多くの過剰なSAクエリが送信される可能性がある。SAクエリの量は、ネットワークにおけるノードの数が増加するのに応じて多項式的に増加する。これらの開示された方法およびシステムは、サブネットにおけるノードによってサブネットマネージャに送信されるSAクエリの数を減らすことができるキャッシングメカニズムを提供する。

50

【 0 0 5 3 】

図 4 は、本発明の一実施形態に従った、クラウド環境における S A パス・キャッシングをサポートする例を示す。図 4 に示されるように、インフィニバンド (I B) サブネット 4 0 0 は、サブネットマネージャ (S M) 4 1 0 および複数のホストノード A 4 0 1 および B 4 0 2 を含み得る。

【 0 0 5 4 】

送信元ホストノード A 4 0 1 (たとえば、 V M A 4 1 1) が宛先ホストノード B 4 0 2 (たとえば、 V M B 4 1 2) と通信しようとして初めて試みたとき、送信元ホストノード A 4 0 1 は S M 4 1 0 に S A パス記録要求を送信することができる。次いで、送信元ホストノードは、ローカルキャッシュ 4 2 1 を用いてパス情報 (たとえば、パス記録 4 2 2) を記憶することができる。

10

【 0 0 5 5 】

さらに、送信元ホストノード A 4 0 1 が同じ宛先ホストノード B 4 0 2 に再接続しようとして試みるときに、送信元ホストノード A 4 0 1 は、サブネットマネージャに要求を送信するのではなく、ローカルキャッシュ 4 2 1 におけるキャッシングテーブル内の宛先ホストノードのアドレスを検索することができる。

【 0 0 5 6 】

パス情報が発見された場合、送信元ホストノード A 4 0 1 は、パス記録 4 2 2 によって示されるように、パス 4 2 0 を用いて宛先ホストノード B 4 0 2 に接続することができる。この場合、 S A クエリは S M 4 1 0 に送信されない。それ以外の場合、送信元ホストノード A 4 0 1 は、必要なパス情報を取得するために S M 4 1 0 に S A パス記録要求を送信することができる。

20

【 0 0 5 7 】

図 5 は、本発明の一実施形態に従った、クラウド環境における S A パス・キャッシングをサポートする例を示す。より特定的には、図 5 は、インフィニバンド環境のサブネット内における S A パス・キャッシングをサポートする例を示す。

【 0 0 5 8 】

図 5 に示されるように、インフィニバンド (I B) サブネット 5 0 0 は、異なるハイパーバイザ 5 1 1 および 5 1 2 をサポートする複数のホストノード A 5 0 1 および B 5 0 2 を含み得る。加えて、各々のハイパーバイザ 5 1 2 および 5 1 3 は、さまざまな仮想マシン (V M) がその上で実行されることを可能にする。たとえば、ホストノード A 1 0 1 上のハイパーバイザ 5 1 1 は、 V M A 5 0 4 をサポートすることができ、ホストノード B 上のハイパーバイザ 5 1 2 は V M B 5 0 5 をサポートすることができる。

30

【 0 0 5 9 】

さらに、ホストノード A 5 0 1 および B 5 0 2 の各々は 1 つ以上のホストチャネルアダプタ (H C A) 5 1 7 および 5 1 8 に関連付けることができる。図 5 に示されるように、ホストノード A 5 0 1 上の H C A 5 1 7 は、 V M A 5 0 4 が使用可能な Q P a 5 0 8 などのキュー対 (Q P) を利用することができる、ホストノード B 5 0 2 上の H C A 5 1 8 は、 V M B 5 0 5 が使用可能な Q P b 5 0 7 を利用することができる。

40

【 0 0 6 0 】

一実施形態に従うと、各々のホストノードはまたメモリ 5 3 0 および 5 4 0 をサポートすることができ、メモリ 5 3 0 および 5 4 0 の各々は (ローカルキャッシュなどの) キャッシュ 5 3 5 および 5 4 5 を含み得るとともに、各々のキャッシュは、キャッシュテーブルに記憶することができる 1 つ以上のパス記録 5 3 7 および 5 4 7 を含み得る。

【 0 0 6 1 】

また、 I B サブネット 5 0 0 は、サブネットマネージャ 5 1 0 を含み得る。サブネットマネージャ 5 1 0 は、ネットワーク初期化、スイッチおよび H C A 上の I B ポートの構成、トポロジー発見およびパス演算のために機能し得る。

【 0 0 6 2 】

一実施形態に従うと、送信元ホストノード A 5 0 1 (たとえば、 V M A 5 0 4) が宛

50

先ホストノード B 5 0 2 (たとえば、VM B 5 0 5) と最初に通信しようとしたとき、送信元ホストノード A 5 0 1 は S M 5 1 0 に S A パス記録要求を送信することができる。次いで、送信元ホストノードは、ローカルキャッシュ 5 3 5 を用いてパス情報 (たとえば、パス記録 5 3 7) を記憶することができる。

【 0 0 6 3 】

さらに、送信元ホストノード A 5 0 1 が同じ宛先ホストノード B 5 0 2 に再接続しようとしたとき、送信元ホストノード A 5 0 1 は、サブネットマネージャに要求を送信するのではなく、キャッシュ 5 3 5 における宛先ホストノードのアドレスを検索することができる。

【 0 0 6 4 】

一実施形態に従うと、パス情報が発見されれば、送信元ホストノード A 5 0 1 は、パス記録 5 3 7 に示されたパスを用いることによって宛先ホストノード B 5 0 2 に接続することができるが、この場合、S A クエリは S M 5 1 0 に送信されない。それ以外の場合、送信元ホストノード A 5 0 1 は、S M 4 1 0 に S A パス記録要求を送信して、必要なパス情報を取得することができる。

【 0 0 6 5 】

一実施形態に従うと、ホストノード A 5 0 1 がサブネットマネージャ 5 1 0 に S A パス記録要求を送信する状況においては、受取られた応答は、キャッシングフラグを含み得る。このキャッシングフラグは、宛先ホストノード B 5 0 2 (D G I D) の所与の G I D アドレスに関連付けられたパス特徴を記憶するために (キャッシュ 5 3 5 内における) ローカルキャッシングテーブルを用いるようにホストノード A 5 0 1 に指示し得る。

【 0 0 6 6 】

図 6 は、一実施形態に従った、クラウド環境におけるサブネット管理をサポートするための方法の例を示す。例示的な方法 6 0 0 はステップ 6 0 1 から開始され得る。ステップ 6 0 1 において、少なくとも第 1 のハイパーバイザおよび第 1 のホストチャネルアダプタに関連付けられた第 1 のホストノードを含む複数のホストノードをクラウド環境内に設ける。ステップ 6 0 2 において、当該方法は、引き続き、複数のアドレスに関連付けられた第 1 の仮想マシンを第 1 のホストノード上に設け得る。ステップ 6 0 3 において、当該方法は、引き続き、クラウド環境内において、第 1 の仮想マシンを、複数のホストノードのうちの第 1 のホストノードから、設けられた第 2 のホストノードにまでマイグレートする。第 2 のホストノードは、少なくとも第 2 のハイパーバイザおよび第 2 のホストチャネルアダプタに関連付けられている。複数のホストノードの各々はローカルキャッシュを含む。各々のローカルキャッシュは 1 つ以上のパス記録を含む。

【 0 0 6 7 】

一実施形態に従うと、第 1 の仮想マシンをマイグレートするステップは、ステップ 6 0 4 において、第 1 のハイパーバイザから第 1 の仮想マシンを分離するステップを含み得る。第 1 のハイパーバイザから第 1 の仮想マシンを分離するステップは、第 1 の仮想マシンに関連付けられた第 1 の仮想機能を第 1 の仮想マシンから分離するステップを含む。ステップ 6 0 5 において、当該方法は、次に、第 1 の仮想マシンに関連付けられた複数のアドレスを第 2 のホストノードに与える。ステップ 6 0 6 において、当該方法は、第 2 のハイパーバイザに関連付けられた第 2 の仮想機能に複数のアドレスを割当てることができる。ステップ 6 0 7 において、当該方法は、第 1 の仮想マシンを第 1 のホストノードから第 2 のホストノード上の第 2 の仮想マシンにまでマイグレートすることができる。ステップ 6 0 8 において、当該方法は、第 2 の仮想マシンを、第 1 の仮想マシンに関連付けられた複数のアドレスにエクスポートするステップで終了し得る。

【 0 0 6 8 】

一実施形態に従うと、S A パス記録キャッシングメカニズムは、R D S プロトコルなどのプロトコルにおいて実現することができ、キャッシングテーブルは、各々のノードのメモリに記憶することができる。以下の擬似コードに示されるプログラムを用いることができる：

10

20

30

40

50

【 0 0 6 9 】

【表 1 - 1】

```

1: private bool SA PathCachingEnabled
2: private list SA PathRecordCacheTable
3:
4: procedure RDSMODULEINITIALIZATION
5: // The Caching table is initialized
6: SA PathRecordCacheTable = empty
7:
8: // The system does not know yet if SA Path Caching is
9: // enabled by the SM, so we assume not.
10: SA PathCachingEnabled = False
11: end procedure
12:
13: procedure (RE-)CONNECTIONESTABLISHMENT(DGID)
14: struct PathRecord DST Path = NULL
15:
16: // Use the cache only if the SA Path Caching is
17: // enabled by the SM
18: if SA PathCachingEnabled then
19: if DGID in SA PathRecordCacheTable.
20: DGIDs then
21: DST Path = Cached PathRecord
22: end if
23: end if
24:
25: // If DST Path is NULL at this point, either the
26: // cache is disabled by the SM, or the path
27: // characteristics for the host with the given DGID
28: // have never been retrieved. In any case, a
29: // send a PathRecord Query can be sent to the SM.
30: if DST Path == NULL then
31: SendAnewSA PathRecordQueryToTheSM
32: WaitForTheReply

```

【 0 0 7 0 】

【表 1 - 2】

```

33: DST Path = PathRecordResponse
34:
35: // If caching is enabled by the SM the reply will
36: // have the reserved field in the PathRecord set
37: // to 1. If not, the reserved field is 0
38: if DST Path ! Reserved Field != 0 then
39: SA PathCachingEnabled = True
40:
41: // Insert the DST Path in the caching table
42: SA PathRecordCacheTable.append(
43: DST Path)
44: end if
45: end if
46: connect to(DST Path)
47: end procedure

```

10

20

【0071】

一実施形態に従うと、送信元ホスト (source host : S H o s t) が宛先ホスト (destination host : D H o s t) と通信しようとして最初に試みるときに、S H o s t はサブネットマネージャに S A パス記録要求を送信することができる。応答によってキャッシングフラグが立てられた場合、S H o s t は、ローカルキャッシングテーブルを用いて、D H o s t (D G I D) の所与の G I D アドレスに関連付けられたパス特徴を記憶することができる。さらに、S H o s t は、キャッシングがサブネットマネージャによってサポートされていることをこのとき認識し、このため、次に S H o s t がいずれかの D H o s t に接続または再接続しようとして試みるときに、S H o s t はキャッシングテーブルを最初に検索することとなる。所与の D H o s t についてのパス情報が発見された場合、S H o s t は、サブネットマネージャに送信される S A クエリの送信が妨害される可能性があり、代わりに、S H o s t は、そのキャッシングテーブル内の情報を用いて D H o s t に接続しようとして試みることができる。

30

【0072】

遮断された接続 (ステップ 2) に関して、図 3 を再び参照すると、上述のキャッシングメカニズムが使用可能にされるシステムにおいては、サブネットマネージャに S A クエリが送信される必要はない。図 3 に記載される場合においては、ステップ 4、5、9 および 10 が省かれており、このため、接続がより高速に回復され、サブネットマネージャに対する負荷 (たとえば S A パスの要求および応答) がより低下する。

【0073】

この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、それを使用して、またはその助けを借りて実行され得る。従って、この発明の特徴は、(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。

40

【0074】

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実行するように処理システムをプログラミングするために使用可能な命令を格納した記憶媒体またはコンピュータ読取可能媒体であるコンピュータプログラム製品において、それを使用して、またはその助けを借りて実現され得る。記憶媒体は、フロッピー (登録商標) ディスク、光ディスク、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

50

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適した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。

【0075】

この発明の特徴は、機械読取可能な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メカニズムとやりと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組込まれ得る。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。

10

【0076】

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実行する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るだろう。

【0077】

加えて、この発明は、この開示の教示に従ってプログラミングされた1つ以上のプロセッサ、メモリおよび/またはコンピュータ読取可能記憶媒体を含む、1つ以上の従来の汎用または専用の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、適切なソフトウェアコーディングが、この開示の教示に基づいて、熟練したプログラマによって容易に準備され得る。

20

【0078】

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が関連技術の当業者にとって明らかになるであろう。

【0079】

この発明は、特定された機能およびそれら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しばしば、この明細書中では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内に含まれる。

30

【0080】

この発明の上述の説明は、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状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であるだろう。これらの変更および変形は、開示された特徴の関連するあらゆる組合せを含む。実施形態は、この発明の原理およびその実用的応用を最良に説明するために選択され説明されたものであり、それにより、企図される特定の用途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、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よう意図されている。

40

【 図 1 】

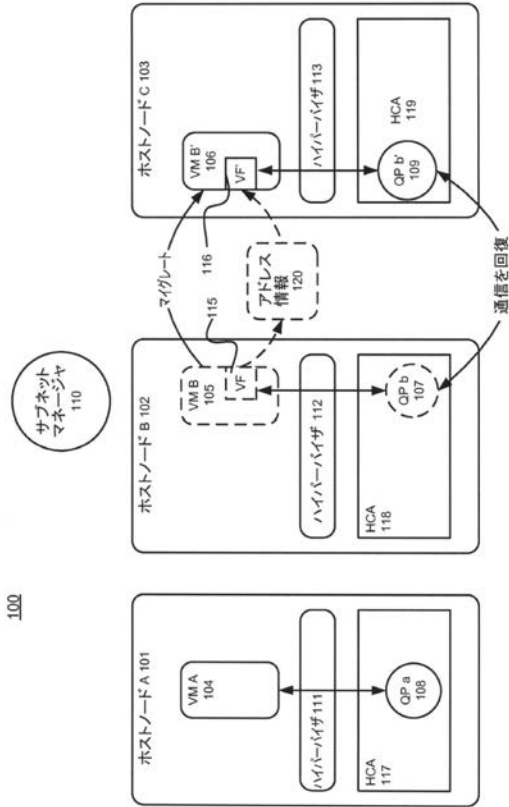


FIGURE 1

【 図 2 】

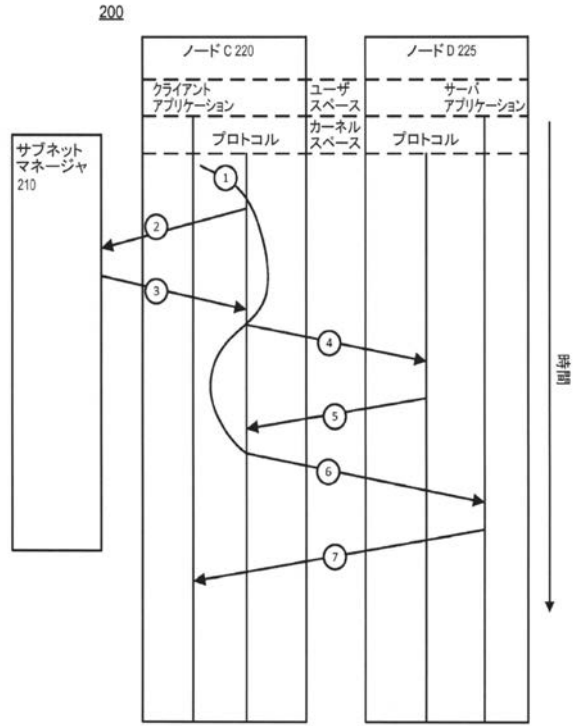


FIGURE 2

【 図 3 】

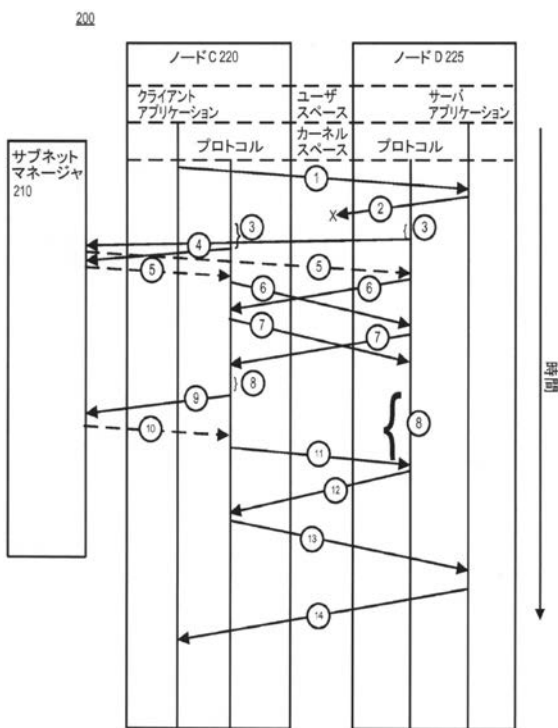


FIGURE 3

【 図 4 】

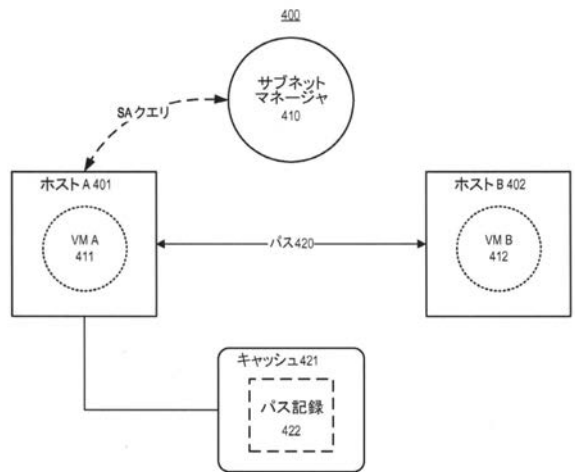


FIGURE 4

【 図 5 】

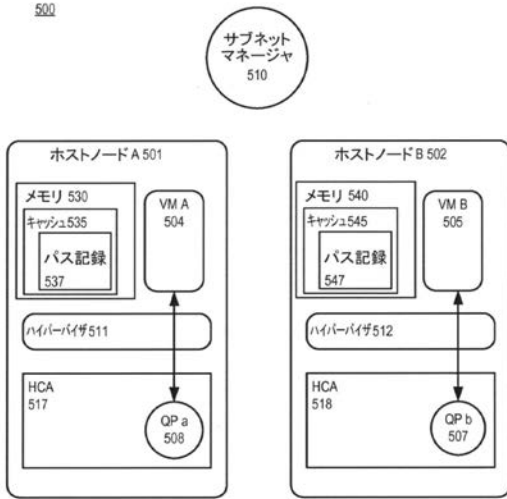


FIGURE 5

【 図 6 】



FIGURE 6

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No PCT/US2015/057860

A. CLASSIFICATION OF SUBJECT MATTER INV. G06F9/48 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) G06F		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) EPO-Internal, WPI Data		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	CLARK C ET AL: "Live migration of virtual machines", PROCEEDINGS OF THE SYMPOSIUM ON NETWORKED SYSTEMS DESIGN AND IMPLEMENTATION, USENIX ASSOCIATION, BERKELEY, CA, US, 2 May 2005 (2005-05-02), - 4 May 2005 (2005-05-04), pages 273-286, XP002443245, figure 1 page 274, left-hand column, line 1 - line 13 page 275, left-hand column, line 10 - line 17 page 276, left-hand column, line 1 - line 27 page 276, left-hand column, line 47 - right-hand column, line 33 page 282, left-hand column, line 21 - line 30 page 285, left-hand column, line 16 - line -/--	1-22
<input checked="" 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
28 January 2016		04/02/2016
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
		Wirtz, Hanno

2

INTERNATIONAL SEARCH REPORT

International application No

PCT/US2015/057860

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
	30 -----	
X	US 2008/189432 A1 (ABALI BULENT [US] ET AL) 7 August 2008 (2008-08-07)	1-3,5-9, 11-13, 15-20,22
A	claims 1,2 figure 3 paragraph [0009] - paragraph [0010] paragraph [0015] paragraph [0045] - paragraph [0055] -----	4,10,14, 21
A	US 2008/186990 A1 (ABALI BULENT [US] ET AL) 7 August 2008 (2008-08-07) claims 1-5 figure 2 paragraph [0012] - paragraph [0016] paragraph [0021] - paragraph [0022] -----	1-22
X,P	EVANGELOS TASOULAS ET AL: "A Novel Query Caching Scheme for Dynamic InfiniBand Subnets", 2015 15TH IEEE/ACM INTERNATIONAL SYMPOSIUM ON CLUSTER, CLOUD AND GRID COMPUTING, 4 May 2015 (2015-05-04), - 7 May 2015 (2015-05-07), pages 199-210, XP055244069, DOI: 10.1109/CCGrid.2015.10 ISBN: 978-1-4799-8006-2 the whole document -----	1-22

INTERNATIONAL SEARCH REPORT

Information on patent family members

International application No

PCT/US2015/057860

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
US 2008189432	A1	07-08-2008	NONE

US 2008186990	A1	07-08-2008	NONE

フロントページの続き

- (31)優先権主張番号 62/121,294
(32)優先日 平成27年2月26日(2015.2.26)
(33)優先権主張国 米国(US)
(31)優先権主張番号 62/133,179
(32)優先日 平成27年3月13日(2015.3.13)
(33)優先権主張国 米国(US)
(31)優先権主張番号 14/924,281
(32)優先日 平成27年10月27日(2015.10.27)
(33)優先権主張国 米国(US)

(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, 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, KN, KP, KR, 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, UA, UG, US

- (72)発明者 グラン, アーンスト・ガンナー
ノルウェー、エヌ - 1 3 2 5 リュサケール、ピィ・オウ・ボックス・1 3 4