

(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2016-524760
(P2016-524760A)

(43) 公表日 平成28年8月18日(2016.8.18)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 9/50 (2006.01)	G06F 9/46 462Z	
G06F 9/46 (2006.01)	G06F 9/46 350	
G06F 12/00 (2006.01)	G06F 12/00 535B	
	G06F 12/00 545A	

審査請求 有 予備審査請求 未請求 (全 34 頁)

(21) 出願番号 特願2016-518066 (P2016-518066)
 (86) (22) 出願日 平成26年6月10日 (2014.6.10)
 (85) 翻訳文提出日 平成27年12月7日 (2015.12.7)
 (86) 国際出願番号 PCT/US2014/041724
 (87) 国際公開番号 WO2014/201012
 (87) 国際公開日 平成26年12月18日 (2014.12.18)
 (31) 優先権主張番号 13/914,104
 (32) 優先日 平成25年6月10日 (2013.6.10)
 (33) 優先権主張国 米国 (US)

(71) 出願人 507303550
 アマゾン・テクノロジーズ・インコーポレ
 ーテッド
 アメリカ合衆国・89507・ネバダ州・
 レノ・ピーオーボックス 8102
 (74) 代理人 100064621
 弁理士 山川 政樹
 (74) 代理人 100098394
 弁理士 山川 茂樹
 (72) 発明者 ジェンキンス, ジョージ, オリヴァー
 アメリカ合衆国・98109-5210・
 ワシントン州・シアトル・テリー アヴェ
 ニュ ノース・410

最終頁に続く

(54) 【発明の名称】 クラウドコンピューティング環境における分散ロック管理

(57) 【要約】

分散ロックマネージャ (DLM) は、仮想化されたコンピューティングリソース及び/または仮想コンピューティングサービスをクライアントに提供する分散コンピューティングシステムにおいて実装することができる。ロックは、クライアントからの作成及び管理を実行する要求にตอบสนองして、DLMによって作成かつ管理することができる。DLMの構成要素は、クライアントアプリケーション構成要素が相互に通信するか、またはロックによって保護された共有リソースにアクセスする以外のネットワーク上で相互に通信することができる。例えば、DLM構成要素は、クラウドコンピューティング環境の制御プレーンネットワーク上で通信することができ、アプリケーション構成要素は、クラウドコンピューティング環境のデータプレーンネットワーク上で通信することができる。DLMは、クライアントにAPIを公開することができ、クライアントが、様々なロック管理操作を実施するために同じノード上のDLM構成要素へのローカル呼び出しを行うことを可能にする。ロック値の意味は、クライアントアプリケーションにおけるそれらの使

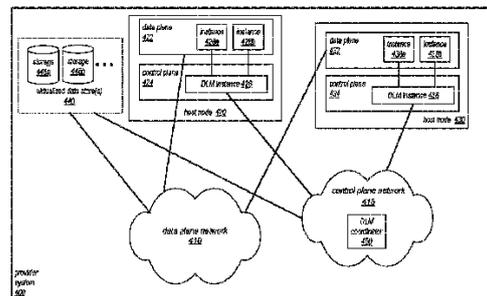


FIG. 4

【特許請求の範囲】**【請求項 1】**

ネットワーク上で相互に連結され、かつ仮想コンピューティングサービスを1つ以上のクライアントに集合的に提供する、複数のコンピューティングノードであって、前記コンピューティングノードの各々が、少なくとも1つのプロセッサとメモリとを備える、複数のコンピューティングノードと、

各々が前記コンピューティングノードのうちの2つ以上のそれぞれのコンピューティングノード上で実行している、2つ以上の仮想コンピュータインスタンスであって、前記仮想コンピュータインスタンスの各々が、クライアントのために分散アプリケーションのアプリケーション構成要素を実装するように構成され、前記仮想コンピュータインスタンスが、前記ネットワークの少なくとも一部上で相互に通信するように構成される、2つ以上の仮想コンピュータインスタンスと、

10

各々が前記2つ以上のコンピューティングノードのそれぞれのコンピューティングノード上で実行している、2つ以上の構成要素を含む分散ロックマネージャであって、前記分散ロックマネージャ構成要素が、それぞれの共有リソースの1つ以上のロックの状態を共有するように構成され、前記1つ以上のロックの状態を共有するために、前記分散ロックマネージャ構成要素が、前記仮想コンピュータインスタンスが相互に通信する前記ネットワークの前記少なくとも一部とは物理的または論理的に別個のネットワーク上で相互に通信するように構成される、分散ロックマネージャと、を備え、

前記アプリケーション構成要素のうちの1つが、前記分散アプリケーションによってアクセスされる共有リソースのロックのためのロック管理操作の実施を開始するために、前記アプリケーション構成要素を実装する前記仮想コンピュータインスタンスを実行している前記コンピューティングノード上で実行している前記分散ロックマネージャへの呼び出しを起動するように構成される、システム。

20

【請求項 2】

前記分散ロックマネージャ構成要素が相互に通信する前記ネットワークが、サービスプロバイダコンピューティング環境の制御プレーンネットワークを備え、前記仮想コンピュータインスタンスが相互に通信する前記ネットワークの前記少なくとも一部が、前記サービスプロバイダコンピューティング環境のデータプレーンネットワークを備える、請求項1に記載の前記システム。

30

【請求項 3】

前記分散ロックマネージャ構成要素への前記呼び出しに応答して、前記分散ロックマネージャ構成要素が、

前記ロック管理操作を実施し、かつ

前記分散ロックマネージャ構成要素が実行している前記コンピューティングノード以外のコンピューティングノード上で実行している少なくとも1つの分散ロックマネージャ構成要素へ前記ロックの結果として生じた状態を通信するように構成される、請求項1に記載の前記システム。

【請求項 4】

前記2つ以上のコンピューティングノードが、前記クライアントのための仮想プライベートネットワークを実装する、請求項1に記載の前記システム。

40

【請求項 5】

1つ以上のコンピュータによって、

複数のコンピューティングノードのうちの所与のコンピューティングノード上で実行している分散ロックマネージャの構成要素によって、共有リソースのロックのためのロック管理操作を実施する要求を受信することであって、前記受信することが、前記所与のコンピューティングノード上で実行しているコンピュータインスタンスから前記要求を受信することを含み、前記コンピュータインスタンスが、前記コンピュータインスタンスに割り当てられたネットワークリソース能力を使用して前記共有リソースにアクセスする、受信することを実施することと、

50

前記要求されたロック管理操作を実施することであって、前記実施することが、前記分散ロックマネージャ構成要素が、前記コンピュータインスタンスに割り当てられた前記ネットワークリソース能力以外のネットワークリソース能力を使用して、前記ロックの状態情報を共有するために、前記複数のコンピューティングノードのうちの別のコンピューティングノード上で実行している別の分散ロックマネージャ構成要素と通信することを含む、実施することと、を含む、方法。

【請求項 6】

前記コンピュータインスタンスが、分散アプリケーションのアプリケーション構成要素を実装し、

前記方法が、前記アプリケーション構成要素が、前記コンピュータインスタンスに割り当てられた前記ネットワークリソース能力を使用して、前記分散アプリケーションの少なくとも 1 つの他の構成要素と通信することをさらに含む、請求項 5 に記載の前記方法。

10

【請求項 7】

前記複数のコンピューティングノードが、1 つ以上の仮想コンピューティングサービスを実装する、請求項 5 に記載の前記方法。

【請求項 8】

前記ロック管理操作を前記実施することが、前記ロックを前記共有リソースと関連付けることを含み、前記通信することが、前記他の分散ロックマネージャ構成要素へ前記関連付けを通信することを含む、請求項 5 に記載の前記方法。

20

【請求項 9】

前記ロック管理操作を前記実施することが、前記ロックの値を変更することを含む、請求項 5 に記載の前記方法。

【請求項 10】

前記ロックの前記値が、前記共有リソースの所有者を識別する、請求項 9 に記載の前記方法。

【請求項 11】

前記分散ロックマネージャの構成要素が、前記ロックの前記変更された値を永続的データストアに書き込むことをさらに含む、請求項 9 に記載の前記方法。

【請求項 12】

前記ロック値を前記変更することが、前記ロック値をアトミックに変更することを含む、請求項 9 に記載の前記方法。

30

【請求項 13】

前記コンピュータインスタンスに割り当てられた前記ネットワークリソース能力が、ネットワーク接続リソースまたは入出力処理能力を含む、請求項 5 に記載の前記方法。

【請求項 14】

前記実施することが、前記所与のコンピューティングノードの前記ロックの状態情報をキャッシュすることをさらに含む、請求項 5 に記載の前記方法。

【請求項 15】

1 つ以上のプロセッサと、

1 つまたはメモリであって、前記 1 つ以上のプロセッサ上で実行されると、前記 1 つ以上のプロセッサに、

40

複数のコンピューティングノードのうちの所与のコンピューティング上で実行している分散ロックマネージャの構成要素によって、共有アクセスが制御されることになる実体と関連付けられたロックのためのロック管理操作を実施する要求を受信することであって、前記受信することが、前記所与のコンピューティングノード上で実行しているリソースインスタンスから、前記要求を受信することを含み、前記複数のコンピューティングノードが、分散ロックサービスを実装し、前記リソースインスタンスが、前記ロックと関連付けられた前記実体にアクセスし、前記要求が、前記分散ロックマネージャによってサポートされた 1 つ以上のロック管理操作を定義するアプリケーションプログラミングインターフェースに準拠する、受信することと、

50

前記要求されたロック管理操作を実施することであって、前記実施することが、前記分散ロックマネージャ構成要素が、前記リソースインスタンスが前記ロックと関連付けられた前記エンティティにアクセスするネットワーク接続性以外のネットワーク接続性を使用して、前記ロックの状態情報を共有するために、前記複数のコンピューティングノードのうち別のコンピューティングノード上で実行している前記分散ロックマネージャの別の構成要素と通信することを含む、実施することと、を実施させるプログラム命令を記憶している、1つ以上のメモリと、を備える、システム。

【発明の詳細な説明】

【技術背景】

【0001】

多数の企業及び他の組織は、（例えば、ローカルネットワークの一部として）共存している、またはそうではなく、（例えば、1つ以上のプライベートまたは公共中間ネットワークを介して接続された）複数の別個の地理的場所に存在しているコンピューティングシステムを用いる等、その事業活動を支援するために無数のコンピューティングシステムを相互接続するコンピュータネットワークを運用する。例えば、単一の組織によりかつそのために運用されるプライベートデータセンター、及びコンピューティングリソースを顧客またはクライアントに提供するために事業として事業体により運用される公共データセンター等、相当数の相互接続されたコンピューティングシステムを収容するデータセンターが一般的になってきている。いくつかの公共データセンター運用者は、様々なクライアントによって所有されるハードウェアのためのネットワークアクセス、電力、及び安全な設置施設を提供する一方、他の公共データセンター運用者は、そのクライアントによる使用のために利用可能になったハードウェアリソースも含む「フルサービス」の施設を提供する。しかしながら、典型的なデータセンターの規模及び範囲が増加したことに伴い、物理的コンピューティングリソースをプロビジョニング、管理、及び運営するタスクはますます複雑になっている。

【0002】

汎用ハードウェアの仮想化技術の出現は、様々なニーズを擁する多数のクライアントのための大規模コンピューティングリソースを運営する点で利点をもたらし、様々なコンピューティングリソースが、複数のクライアントによって効率的かつ安全に共有されることを可能にしている。例えば、仮想化技術は、単一の物理的コンピューティングマシンによってホストされた1つ以上の仮想マシンを各ユーザに提供することによって、単一の物理的コンピューティングマシンが、複数のユーザの間で共有されることを可能にすることができ、各々のそのような仮想マシンは、ユーザが所与のハードウェアコンピューティングリソースの唯一の運用者かつ管理者であるという錯覚をユーザに提供すると同時に、様々な仮想マシン間でアプリケーションの分離及びセキュリティも提供する、別個の論理的コンピューティングシステムとして作用するソフトウェアシミュレーションである。その上、いくつかの仮想化技術は、複数の別個の物理的コンピューティングシステムにまたがる、複数の仮想プロセッサを備える単一の仮想マシン等、2つ以上の物理的リソースにまたがる仮想リソースを提供することが可能である。別の例として、仮想化技術は、複数のデータ記憶装置にわたって分散され得る仮想化されたデータストアを各ユーザに提供することによって、データストレージハードウェアが複数のユーザの間で共有されることを可能にすることができ、各々のそのような仮想化されたデータストアは、ユーザがデータストレージリソースの唯一の運用者かつ管理者であるという錯覚をユーザに提供する別個の論理的データストアとして作用する。

【0003】

分散コンピューティングをサポートするいくつかのシステムにおいて、分散アプリケーションのアプリケーション構成要素またはプロセスは（時々）、様々なタイプの共有リソースにアクセスすることができる。いくつかのそのようなシステムにおいて、それらの共有リソースへのアクセスを制御及び/または同期するために、分散ロックマネージャが使用される。

10

20

30

40

50

【図面の簡単な説明】

【0004】

【図1】仮想コンピューティングシステムを提供する分散コンピューティングシステムにおいて分散ロックマネージャを提供するための方法の一実施形態を示すフロー図である。

【図2】少なくともいくつかの実施形態に従う、例示的プロバイダネットワーク環境を示すブロック図である。

【図3】いくつかの実施形態に従う、例示的データセンターを示すブロック図である。

【図4】いくつかの実施形態に従う、制御プレーンにおいて分散ロックマネージャを実装する例示的データセンターを示すブロック図である。

【図5】仮想化されたリソースを提供する分散コンピューティングシステムにおいて分散ロックマネージャを実装するための方法の一実施形態を示すフロー図である。

【図6】クライアントアプリケーションまたはプロセスが、クラウドコンピューティング環境において実装される分散ロックマネージャの機能にアクセスするための方法の一実施形態を示すフロー図である。

【図7】いくつかの実施形態に従う、制御プレーンにおいて分散ロックマネージャを実装し、その環境の外側で実行しているプロセスに分散ロックマネージャを公開する例示的データセンターを示すブロック図である。

【図8】クラウドコンピューティング環境の外側で実行するクライアントアプリケーションまたはプロセスが、クラウドコンピューティング環境内に実装される分散ロックマネージャの機能にアクセスするための方法の一実施形態を示すフロー図である。

【図9】異なる実施形態に従う、本明細書に記載される技法の一部または全てを実装する例示的コンピュータシステムを示すブロック図である。

【発明を実施するための形態】

【0005】

本明細書において、実施形態はいくつかの実施形態及び例示的図面によって説明されるが、当業者は、実施形態が説明される実施形態または図面に限定されないことを認識するであろう。図面及びそれらに対する詳細説明は、開示される特定の形式に実施形態を限定することを意図せず、しかし対照的に、その意図は、添付の請求項によって定義される趣旨及び範囲に含まれる全ての変形、同等物、及び代替を網羅することであることを理解されたい。本明細書に使用される見出しは、構成上の目的のみであり、記載または請求項の範囲を限定するために使用されることを意図するものではない。本明細書全体で使用されるように、「～することができる」(may)は、必須の意味(すなわち、～しなければならない(must))ではなく、許容の意味(すなわち、～する可能性を有する意味)において使用される。同様に、「含む」(include、including、includes)は、含むがそれらに限定されないことを意味する。

【0006】

本明細書において、仮想化されたコンピューティングリソースをクライアントに提供するシステムにおいて分散ロック管理を実装するためのシステム及び方法の様々な実施形態が説明される。いくつかの実施形態において、分散ロックマネージャ(DLM)の構成要素は、共有リソースにアクセスする、分散アプリケーションのアプリケーション構成要素が、相互に、及び/または共有リソースと通信する別のネットワーク(またはネットワークの一部)とは物理的または論理的に別個のネットワーク(またはネットワークの一部)上で共有リソースのロックを管理するために相互に通信することができる。例えば、いくつかの実施形態において、DLM構成要素は、クラウドコンピューティング環境の制御プレーンネットワーク上で相互に通信ことができ、一方、分散アプリケーションのアプリケーション構成要素は、クラウドコンピューティング環境のデータプレーンネットワーク上で相互に、及び/または共有リソースと通信することができる。いくつかの実施形態において、アプリケーションプログラミングインターフェース(API)は、クラウドコンピューティング環境内部で実行しているクライアントアプリケーション(またはその構成要素)に公開されることに加えて、クラウドコンピューティング環境の外側で実行

10

20

30

40

50

しているクライアントアプリケーション（またはその構成要素）に公開され得る。そのような実施形態において、クライアントアプリケーション（またはその構成要素）のうちの何れかは、（例えば、それらのクライアントアプリケーションによって共有リソースのロック上のロック管理操作を開始するために）DLMのロック機構に参与することができる。

【0007】

分散アプリケーションを実装及び/または仮想化されたコンピューティングリソースをクライアントアプリケーションに提供するようなクラスタ環境では、リソース間に発生する可能性がある様々なレベルの共有が存在し得ることに留意されたい。これらのシステムは、（例えば、コンピューティングノードの障害または他のフェイルオーバー状況に
10 応答して）移動させる必要があるリソースを管理するために、様々な機構を実装することができる。例えば、10の仮想化されたインスタンスが仮想プライベートクラウド（VPC）内部で実行しているシステムでは、（例えば、1つのインスタンスからそれを切断し、別のインスタンスに再接続するために）異なるインスタンス間でネットワークインターフェース（例えば、エラスティック（Elastic）ネットワークインターフェースまたはENI、あるいは別のタイプの仮想ネットワークインターフェース）を移動させることが
15 所望される場合が発生し得る。

【0008】

いくつかの既存のシステムは、共有リソースへのアクセスを制御及び同期するために、分散ロックマネージャを使用するが、これらのシステムは典型的に、ネットワーク接続性を必要とする。例えば、いくつかのクラスタ化技術は、クラスタが使用するロックを管理するために使用されるDLMの相互接続機構として、ネットワークを使用する。様々な
20 実施形態において、本明細書に記載されるシステムは、そのようなネットワーク接続性を必要としない分散ロックマネージャ（DLM）を実装することができる。代わりに、DLMは、クラウドコンピューティング環境の制御プレーンにおいて実装することができ、様々なAPI呼び出しを使用して、（例えば、これらの分散アプリケーションによってアクセス可能であるリソースの1つ以上のロックを管理するために）クラウドコンピューティング環境内のリソースインスタンス上で実行している分散アプリケーションによってアクセス
25 することができる。いくつかの実施形態において、例えば、顧客が、顧客のアプリケーションが他のアプリケーションから分離される（すなわち、他のアプリケーションまたは顧客には見えず、及び/または他のアプリケーションとは異なる仮想マシン上で稼動する）ことを望む場合、分散アプリケーションのアプリケーション構成要素が実行するリ
30 スインスタンスは、仮想プライベートクラウド（VPC）として構成することができる。いくつかの実施形態において、仮想プライベートクラウド内で実行している各仮想マシンには、その独自のプライベートIPアドレスを与えることができる。いくつかの実施形態において、これらのプライベートIPアドレス（「仮想プライベートIPアドレス」とも呼ぶことができる）は、これらがクライアントの独自のプライベートネットワーク内のIP
35 アドレスのいずれとも重ならないように選択することができる。いくつかの実施形態において、VPC内部のリソースインスタンスは、公共サービスAPIを使用してサービスリクエストをDLMへ送信することができるが、一方、他の実施形態において、VPC内部
40 のリソースインスタンスは、プライベートサービスAPIを使用してサービスリクエストをDLMへ送信することができる。

【0009】

クラスタ化されたシステムは一般に、少なくともいくつかの共有リソースを含むので、これらの環境において、DLMは非常に有用であり得る。上記のように、分散ロックマネージャは、様々なタイプの共有リソースへのアクセスを制御及び/または同期するために、分散アプリケーションによって使用され得る。様々な実施形態において、ロックは、任意のタイプのリソース（例えば、ディスクボリューム、ENIまたは別の種類の仮想ネットワークインターフェース、データベース行、またはファイル）、または一般に、複数の
45 プロセスによるアクセスが制御されなければならない任意の実体（例えば、アプリケーシ

10

20

30

40

50

ョン構成要素またはリソースインスタンスを含む)と関連付けられ得る。

【0010】

いくつかの実施形態において、DLMは、個々のリソース/エンティティがロックされ得るように、またはリソース/エンティティの集合体(例えば、2つ以上の「子」リソースを含む「親」リソース)がロックされ得るように、任意の数のレベルを有するロック階層を実装することができる。いくつかの実施形態において、ロックは、より下位のリソース(例えば、子リソース)上でロックを取得できる前に、より上位のリソース(例えば、親リソース)上で取得されなければならない。

【0011】

様々な実施形態において、本明細書に記載される分散ロックマネージャは、様々なロックモードをサポートすることができ、その各々は、関連付けられたリソース/エンティティ(またはそのロック)の共有プロパティを示すことができる。例えば、異なる実施形態において、以下のロックモードのうちのいずれか、または全てがサポートされ得る。

・ヌルロックモード。このモードのロックを保持するプロセスが、関連付けられた共有リソースに対する関心を表すことを可能にするが、このプロセスまたは他のプロセスによる関連付けられた共有リソースへの一切のアクセスを防止しない。

・並列読み出しロックモード。プロセスが関連付けられた共有リソースを読み出すことを可能にし(しかし変更は可能にしない)、かつ他のプロセスが関連付けられた共有リソースを読み取る、または変更することを可能にするが、関連付けられた共有リソースに対する排他的アクセスを防止する。

・並列書き込みロックモード。プロセスが関連付けられた共有リソースを読み出す、または変更することを可能にし、かつ他のプロセスが関連付けられた共有リソースを読み出す、または変更することを可能にするが、関連付けられた共有リソースに対する排他的アクセスを防止する。

・保護読み出しロックモード。プロセスが関連付けられた共有リソースを読み出すことを可能にし(しかし変更は可能にしない)、かつ他のプロセスが関連付けられた共有リソースを読み出すことを可能にする(しかし変更は可能にしない)。

・保護書き込みロックモード。プロセスが関連付けられた共有リソースを読み出す、または変更することを可能にし、かつ並列読み出しアクセスを有する他のプロセスが、関連付けられた共有リソースを読み出すことを可能にする(しかし変更することは可能にしない)。

・排他ロックモード。プロセスが関連付けられた共有リソースを読み出す、または変更することを可能にするが、他のプロセスによる関連付けられた共有リソースへの一切のアクセスを防止する。

【0012】

いくつかの実施形態において、共有リソースのロックを取得することは、(例えば、DLMの構成要素へ)ロックを取得する要求を通信すること、及び/またはロック上の待ち行列に入ること(例えば、要求するプロセスのインジケータを、ロックを取得するために待機しているプロセスの待ち行列に追加すること)を含むことができる。ロック要求は、異なる実施形態において、同期(例えば、プロセスは、ロックが与えられるまで待機することができる)または非同期(例えば、プロセスは、ロックを必要としない他の操作を継続することができるが、ロックが与えられる場合、ロックが与えられると中断され得る)であり得る。

【0013】

いくつかの実施形態において、ロックデータ構造(またはDLMによって作成及び/または管理されるロックの他の表現)は、ロック値を含むことができ、その意味は、アプリケーションにおけるその使用によって確立され得る。いくつかのそのような実施形態において、DLMは、ロック値の意味が何であるかを知らない場合がある(または知る必要がない場合がある)。いくつかの実施形態において、ロック値は、ロックと関連付けられた共有リソースを所有するプロセスまたはアプリケーションを特定することができる。いく

10

20

30

40

50

つかの実施形態において、ロック値は、ロックの、または関連付けられた共有リソースのバージョン識別子を表現する単調増分値であり得る。いくつかの実施形態において、共有リソースにアクセスする前に、プロセス（またはアプリケーション）は、関連付けられたロックの値を読み出すことができる。値が、プロセス（またはアプリケーション）がロックを最後に取得したときから変化していない場合、プロセス（またはアプリケーション）は、プロセス（またはアプリケーション）がロックを最後に取得したとき以来、共有リソースが別のプロセス（またはアプリケーション）によって変更されていないことを知ることができる。いくつかの実施形態において、本明細書に記載されるDLMは、デッドロック検出をサポートまたは提供することができる。

【0014】

上述のように、いくつかの実施形態において、DLMは、DLM構成要素が（例えば、クラウド内で実行しているクライアントアプリケーションに割り当てられた以外のネットワーク接続性または処理能力を使用する）制御プレーンネットワーク上で相互に通信するような様式において、分散コンピューティング環境（例えば、仮想コンピューティングリソース及び/または仮想コンピューティングサービスをクライアントに提供するクラウドコンピューティング環境）において実装することができる。そのような実施形態において、DLM構成要素間の通信は、クライアントアプリケーションの使用を目的とするネットワークリソースを消費しない、または、さもなければクライアントアプリケーションの構成要素間の通信に干渉しない場合がある。

【0015】

本明細書に記載されるシステム及び方法は、異なる実施形態において、ネットワーク環境内部の1つ以上のコンピューティングシステム上で、またはそれらによって実装され得る。本明細書に記載される分散ロックマネージャの実施形態が実装され得る例示的なコンピュータシステムが、図9に示される。これらの分散ロックマネージャを実装するための様々なシステム及び方法の実施形態が、概して、本明細書において、インターネット等の中間ネットワークを介して、サービスプロバイダのプロバイダネットワーク上に実装された仮想化されたリソース（例えば、仮想化されたコンピューティング及びストレージリソース）をクライアントに提供するサービスプロバイダの状況において、記載される。図2～4、7及び9（ならびにそれらの対応する説明）は、本明細書に記載されるシステム及び方法の実施形態が実装され得る例示的環境を示して記載し、限定することを目的としない。少なくともいくつかの実施形態において、プロバイダネットワークを介してサービスプロバイダのクライアントに提供されるリソースのうち少なくともいくつかは、他のクライアント（単数または複数）と共有されるマルチテナントハードウェア上、及び/または特定のクライアント専用のハードウェア上に実装された仮想化されたコンピューティングリソースであり得る。各仮想化されたコンピューティングリソースは、リソースインスタンスと呼ぶことができる。リソースインスタンスは、例えば、サービスプロバイダのクライアントに賃借またはリースされ得る。例えば、サービスプロバイダのクライアントは、リソースインスタンスを取得及び構成するため、並びに、例えば、仮想化されたプライベートネットワーク等のリソースインスタンスを含む仮想ネットワーク構成を確立及び管理するために、サービスへのAPIを介して、プロバイダネットワークの1つ以上のサービスにアクセスすることができる。

【0016】

いくつかの実施形態において、リソースインスタンスは、例えば、複数のオペレーティングシステムが、ホストコンピュータ上で同時に稼動することを可能にするハードウェア仮想化技術に従って、すなわち、ホスト上の仮想マシン（VM）として実装され得る。ホスト上のハイパーバイザ、または仮想マシンモニタ（VMM）は、ホスト上のVMに仮想プラットフォームを提示することができ、VMの実行を監視する。各VMには1つ以上のプライベートIPアドレスを提供することができ、ホスト上のVMMは、ホスト上のVMのプライベートIPアドレスを認識することができる。そのようなハードウェア仮想化技術を採用するシステムの実施例が、図3に示され、以下に詳細に説明する。

10

20

30

40

50

【 0 0 1 7 】

いくつかの実施形態において、VMMは、クライアントデータパケットをカプセル化し、プロバイダネットワーク内部の異なるホスト上のクライアントリソースインスタンス間のネットワークサブストレート上でルーティングするために、インターネットプロトコル（IP）トンネリング技術を使用することができる。プロバイダネットワークは、ルータ、スイッチ、ネットワークアドレストランスレータ（NAT）等のネットワークデバイス、並びにデバイス間の物理的接続を含む物理的ネットワークサブストレートを含むことができる。プロバイダネットワークは、カプセル化されたパケット（すなわち、オーバーレイネットワーク上でルーティングするためのオーバーレイネットワークアドレス情報を含むがこれに限定されない、オーバーレイネットワークメタデータを用いてタグ付けされているクライアントパケット）が、トンネルまたはオーバーレイネットワーク経路を介して、ネットワークサブストレートを通過させられ得るオーバーレイネットワークを提供するために、IPトンネリング技術を採用することができる。IPトンネリング技術は、ネットワークサブストレート上にオーバーレイネットワークを作成するためのマッピング及びカプセル化システムを提供することができ、オーバーレイネットワーク層（公共IPアドレス）及びネットワークサブストレート層（プライベートIPアドレス）のために別々の名前空間を提供することができる。少なくともいくつかの実施形態において、オーバーレイネットワーク層内のカプセル化されたパケットは、それらのトンネルサブストレートターゲット（プライベートIPアドレス）が何であるべきかを決定するために、マッピングディレクトリで確認することができる。IPトンネリング技術は、物理的ネットワークサブ

10

20

ストレイト上にオーバーレイされた仮想ネットワークトポロジを提供することができ、クライアントに提示されるインターフェース（例えば、サービスAPI）は、オーバーレイネットワークに付加されるので、クライアントリソースインスタンスが、パケットが送信されるIPアドレスを提供すると、IPアドレスは、IPオーバーレイアドレスがどこにあるかを決定することができるマッピングサービスと通信することによって、仮想空間内で稼動する。オーバーレイネットワーク技術の例示的使用が、図3に示され、以下に詳細に説明する。

【 0 0 1 8 】

ホスト上のクライアントリソースインスタンスは、伝送制御プロトコル（TCP）等のステートフルプロトコルに従って、及び/またはユーザデータグラムプロトコル（UDP）等のステートレスプロトコルに従って、同じホスト上または異なるホスト上の他のクライアントリソースインスタンスと通信することができる。しかしながら、クライアントパケットは、送信側VMMによって、オーバーレイネットワークプロトコルに従ってカプセル化され、受信側VMMによってカプセル化が解除される。ホスト上のVMMは、ホスト上のクライアントリソースインスタンスから、別のクライアントリソースインスタンスのIPアドレスに向けられたクライアントパケット（例えば、TCPまたはUDPパケット）を受信すると、オーバーレイネットワーク（またはIPトンネリング）プロトコルに従ってクライアントパケットをカプセル化またはタグ付けし、配信のためにカプセル化されたパケットをオーバーレイネットワーク上へ送信する。カプセル化されたパケットは次いで、IPトンネリング技術に従って、オーバーレイネットワークを介して別のVMMヘル

30

40

ピングされ得る。他のVMMは、パケットからオーバーレイネットワークカプセル化を解除し、ターゲットクライアントリソースインスタンスを実装するホスト上の適切なVMへ、クライアントパケット（例えば、TCPまたはUDPパケット）を配信する。すなわち、いくつかの実施形態において、サービスプロバイダのコンピューティング環境（例えば、サービスプロバイダのデータセンター）には単一の基底の物理的ネットワークが存在し得るが、本明細書に記載されるカプセル化は、それが、あたかも各クライアントアプリケーション（または1つ以上のクライアントアプリケーションが実行する各クライアントリソースインスタンス）が、その独自の仮想ネットワーク上で稼動しているように見えることを可能にすることができる（例えば、複数のクライアントアプリケーションのデータパケットは、同じ物理的ネットワーク上を移動することができるが、クライアントアプ

50

リケーションの各々に向かうトラフィックは、プライベートネットワーク上を移動しているかのように見えることができる)。

【0019】

いくつかの実施形態において、オーバーレイネットワークは、コネクションレス(またはステートレス)IPプロトコルに従って実装されたステートレスネットワークであり得る。いくつかのそのような実施形態において、送信側VMMは、ルーティング及び配信のために、カプセル化されたパケットをオーバーレイネットワークへ送信するが、パケットの配信に関する確認応答(ACK)または他の応答を受信しない。他の実施形態において、VMMは、カプセル化されたパケットの配信に関するACKまたは他の応答を受信することができる。

10

【0020】

分散ロックマネージャの構成要素がそれらによって管理されるロックを採用するアプリケーションと同じデータプレーン(単数または複数)内部で実行及び/または通信する既存のシステムとは異なり、本明細書に記載されるシステムのいくつかの実施形態では、分散ロックマネージャが、分散コンピューティングシステム(例えば、仮想コンピューティングリソース及び/または仮想コンピューティングサービスをクライアントに提供するシステム)のプロトコルプレーン層内に組み込むことができ、この制御プレーン層は、システム内の仮想化されたリソース上で実行しているクライアントソフトウェアとは異なるネットワーク可用性を有することができる。いくつかの実施形態において、クライアントソフトウェアは、(例えば、APIを通じて)サービスとしてクライアントに公開され得る、制御プレーン層の、またはDLMのネットワーク可用性を認識しない場合がある。

20

【0021】

上記のように、DLMを含む既存のシステムにおいて、相互に通信するためにDLMの構成要素によって使用される接続機構は、ロックと関連付けられたリソースがクライアントアプリケーションによってアクセスされるのと同じ接続機構である。しかしながら、本明細書に記載されるシステムのいくつかの実施形態においては、そうではない。代わりに、クライアントアプリケーションは、本明細書に記載されるロック及びロック機構に關与するためには、ローカルAPI呼び出しを行うことだけが必要であり得る。いくつかの実施形態において、DLMによって管理されるロックは、APIの観点から、ロックの「ハンドル」と考えられ得る、固有の識別子(または特定の名前空間を含む固有の識別子)を有することができる。上記のように、いくつかの実施形態において、クライアントアプリケーションは、相互に通信するために、DLM構成要素によって利用されるネットワーク接続について、何も知る必要がない場合がある。代わりに、各々が特定のロック管理操作に対応する、1つ以上のAPIを起動することによって、ロックを操作することができる。いくつかの実施形態において、ロックの使用法、及び/またはロックの通知法を決定するのは、クライアントアプリケーションの責任であり得る。いくつかの実施形態において、プロセスまたはアプリケーションがロックを作成すると、1つ以上の他のプロセスまたはアプリケーション(時には「関与者」と呼ばれる)は、(例えば、観察者として、または他のロックモードに従って)ロックにサブスクライブすることができる。

30

【0022】

いくつかの既存のDLMとは異なり、本明細書に記載されるDLMは、それらが実装される分散システムの制御プレーンへのフックを含むことができ、ユーザが、異なる共有レベルを含むロックを作成し、ロックを使用するアプリケーションのネットワーク接続性に依存することなく、それらのロックを管理することを可能にする。

40

【0023】

いくつかの実施形態において、本明細書に記載される分散ロックマネージャは、オンプレミスのホスト(例えば、クライアントネットワーク上で実行しているホスト)と、プロバイダネットワーク(例えば、公共またはプライベートクラウドコンピューティング環境)内で実行しているホストとの間の境界にまたがることができる。いくつかの実施形態において、本明細書に記載される技法を実装することは、DLMロックが、(例えば、顧客

50

の施設にある)クライアントネットワーク上に存在するホストによって、及び/またはクラウドコンピューティング環境内部に存在するホストによって、取得され、操作され、管理されることを可能にする、一連のソフトウェア構成要素を構築することを含むことができる。ロックそれ自体は、相対的に単純であり得て、どのリソースがロックされるべきか、いつロックされるべきかを決定することは、クライアントアプリケーションの設計者に一任され得る。様々な実施形態において、ロックと関連付けられたリソースは、ディスクボリューム、仮想ネットワークアダプタ、ファイル、データベース内部のレコードまたは他の項目、あるいはアプリケーション開発者が、複数のプロセスによる非同期のアクセスから保護することを希望する任意の他のリソースを含むことができる。上記のように、ロックは、階層にグループ化され得る。例えば、ファイルロックと関連付けられたレコードロック、または一連のロックを共有するホスト群が存在し得る。いくつかの実施形態において、ロックの状態が変化すると(例えば、排他的書き込みアクセスのために取得されると)、そのロックのサブスクリバの全ては、一貫したある時点のロックのビューを有するようになる。すなわち、ある時点で(かつプロセスが排他的ロックを保持する場合のみ)、1つのプロセスだけがロックの状態を変更することができ、任意の所与の時点で、ロックの関与者の全員が同じ値を見るとき、ロック状態の変化はアトミックであり得る。いくつかの実施形態において、ロックそれ自体は、いずれのサブスクリバが不在であっても、ロックが永続する点で耐久性があり得る。例えば、ディスクボリューム上のロックを共有している3つのホストが存在し(例えば、2つがクラウドコンピューティング環境、1つがオンプレミス)、そのうちの1つは、排他的書き込みのロックを保持している場合、全ての3つのホストに障害が発生、またはシャットダウンしてから、次いで、再起動すると、ロックの状態は、ロックマネージャによって維持することができ、排他的書き込みロックは依然として、それを取得したホストによって維持することができる。例えば、ロック状態は、ロックを使用するアプリケーション構成要素が実行しているホスト以外の、ロック機構に関与している制御プレーンのどこかで維持され得る。いくつかの実施形態において、ローカルプロセス(例えば、クラウドコンピューティング環境内及び/またはクライアントネットワーク上の様々なリソースインスタンス内で実行しているクライアントアプリケーション)は、DLMによってそれらの代わりに作成された各ロックのロック状態のキャッシュされたビューを維持することができ、クラウドコンピューティング環境内の永続的データストア内(例えば、データベース内)の各ロックのロック状態のコピーも存在し得る。

10

20

30

【0024】

いくつかの実施形態において、ロックを表示及び/または変更するためにロックグループに加わるためには、またはロックにアクセスするためには、承認を必要とするセキュリティアクセスプロパティを有することができるという意味で、ロックはセキュアでもあり得る。前述のように、ロックは、その意味がロックを利用するアプリケーションまたは仮想コンピューティングサービスによって決定される値を含むことができる。いくつかの実施形態において、それらの独自のロック規則に従うことは、クライアントアプリケーションの責任であり得ることに留意されたい。例えば、クラウドコンピューティング環境内で実行しているアプリケーションが、ロックをディスクボリュームと関連付ける場合、クラウドコンピューティング環境は、この関連付けが存在することを知らなくても、または知る必要がなくてもよい。代わりに、関連付けが何であるか、及びロックがどのようにリソースと関連付けられているかを知ることが、アプリケーションに一任され得る。いくつかの実施形態において、ロックと対応するリソースとの間の関連付けは、クラウドコンピューティング環境の他の実体(例えば、クラウドコンピューティング環境の管理構成要素)には可視ではない場合さえあり得る。前述のように、これもまた、ロックの使用法及びそれらの値の意味を決定することは、クライアントアプリケーション(またはアプリケーションの設計者)に一任され得る(例えば、ロックと、保護する、及び/またはアクセスを制御する共有リソースとの間の関係を定義することは、クライアントアプリケーションに一任され得る)。例えば、ファイルが渡されるクラスタ化されたアプリケーション

40

50

では、ファイルと関連付けられたロックは、任意の所与の時点で、どのアプリケーション構成要素またはプロセスがファイルを所有するかを示すことができる。

【0025】

いくつかの実施形態において、値は、上述のように、単調増加するバージョン番号であり得る。いくつかの実施形態において、本明細書に記載される技法は、ソフトウェアライブラリによって実装され、並びに/あるいは、オープンソースまたは専有（クローズ）オペレーティングシステム、及び/またはスマートフォンまたは他のモバイルデバイスのオペレーティングシステムを含む、様々なオペレーティングシステムのために開発されているコマンド行ツールによって起動され得る。

【0026】

仮想コンピューティングシステムを提供する分散コンピューティングシステム内で分散ロックマネージャをクライアントに提供するための方法の一実施形態が、図1のフロー図によって示される。110に示されるように、この実施例において、方法は、クライアントにAPIを公開する分散ロックマネージャ(DLM)を実装しているクライアントに仮想コンピューティングサービス(例えば、仮想コンピューティングリソースを使用して実装されたサービス)を提供する分散システムを含むことができる。方法はまた、120のように、DLMの構成要素が、リソースインスタンスが分散アプリケーションの一部(例えば、プロセスまたは他の構成要素)を実装する、同じノード上で実行しているリソースインスタンスから、共有リソースのロックのためのロック管理操作を実施する要求を受信することも含むことができる。

【0027】

この実施例に示されるように、方法は、130のように、DLM構成要素が、要求された操作を実施し、データプレーンネットワーク(アプリケーションに割り当てられている一部)とは物理的または論理的に別個の制御プレーンネットワーク上でロックの状態情報(例えば、操作を実施することから生じる変更されたロック値または他の状態情報)を共有するために別のDLM構成要素と通信することを含むことができる。方法はまた、140のように、DLM(例えば、要求を受信し、及び/または要求された操作を実施したDLM構成要素)が、状態情報をローカルに(例えば、このDLM構成要素が実行しているノード上)キャッシュすること、及び状態情報が、(例えば、要求を受信した、及び/または要求された操作を実施したDLM構成要素によって、あるいはDLMの別の構成要素によって)永続的ストレージに書き込まれることも含み得る。

【0028】

例示的プロバイダネットワーク環境

本項では、本明細書に記載される方法の実施形態が実装され得る、例示的プロバイダネットワーク環境を説明する。ただし、これらの例示的プロバイダネットワーク環境は限定することを目的としない。

【0029】

図2は、少なくともいくつかの実施形態に従う、例示的プロバイダネットワーク環境を示す。プロバイダネットワーク200は、クライアントが、1つ以上のデータセンターにおけるプロバイダネットワーク(単数または複数)内部のデバイス上に実装された、演算及びストレージリソースを含むがこれに限定されない、仮想化されたリソースのインスタンス212を購入、賃借、またはその他取得することを可能にする1つ以上の仮想化サービス210を経由して、リソース仮想化をクライアントに提供することができる。プライベートIPアドレス216は、リソースインスタンス212と関連付けることができ、プライベートIPアドレスは、プロバイダネットワーク200上のリソースインスタンス212の内部ネットワークアドレスである。いくつかの実施形態において、プロバイダネットワーク200はまた、クライアントがプロバイダ200から取得することができる、公共IPアドレス214及び/または公共IPアドレス範囲(例えば、インターネットプロトコルバージョン4(IPv4)またはインターネットプロトコルバージョン6(IPv6)のアドレス)も提供することができる。

【 0 0 3 0 】

従来、プロバイダネットワーク 200 は、仮想化サービス 210 を介して、サービスプロバイダのクライアント（例えば、クライアントネットワーク 250 A を運用するクライアント）が、クライアントに割り振られた、または割り当てられた少なくともいくつかの公共 IP アドレス 214 を、クライアントに割り振られた特定のリソースインスタンス 212 と動的に関連付けることを可能にすることができる。プロバイダネットワーク 200 はまた、クライアントが、クライアントに割り当てられた 1 つの仮想化されたコンピューティングリソースインスタンス 212 に以前にマッピングされた、公共 IP アドレス 214 を、これもクライアントに割り当てられている別の仮想化されたコンピューティングリソースインスタンス 212 へ、再マッピングすることを可能にすることもできる。サービスプロバイダによって提供された仮想化されたコンピューティングリソースインスタンス 212 及び公共 IP アドレス 214 を使用して、クライアントネットワーク 250 A の運用者等のサービスプロバイダのクライアントは、例えば、インターネット等の中間ネットワーク 240 上で、クライアント専用アプリケーションを実装し、クライアントのアプリケーションを提示することができる。中間ネットワーク 240 上の他のネットワーク実体 220 は、次いで、クライアントネットワーク 250 A によって発行された送信先公共 IP アドレス 214 へのトラフィックを生成することができ、トラフィックは、サービスプロバイダのデータセンターヘルディングされ、データセンターで、ネットワークサブストレートを介して、送信先公共 IP アドレス 214 に現在マッピングされている仮想化されたコンピューティングリソースインスタンス 212 のプライベート IP アドレス 216

10

20

【 0 0 3 1 】

本明細書に使用されるように、プライベート IP アドレスは、プロバイダネットワークにおけるリソースインスタンスの内部ネットワークアドレスを指す。プライベート IP アドレスは、プロバイダネットワーク内部のみでルーティング可能である。プロバイダネットワークの外側から発信されているネットワークトラフィックは、プライベート IP アドレスへ直接ルーティングされず、その代わりに、トラフィックは、リソースインスタンスにマッピングされている公共 IP アドレスを使用する。プロバイダネットワークは、公共 IP アドレスからプライベート IP アドレスへ、及びこの逆のマッピングを実施するために、ネットワークアドレス変換（NAT）または類似の機能性を提供するネットワークデバイスまたは機器を含むことができる。

30

【 0 0 3 2 】

本明細書に使用されるように、公共 IP アドレスは、サービスプロバイダによって、またはクライアントによってのいずれかで、リソースインスタンスに割り振られる、インターネットをルーティング可能なネットワークアドレスである。公共 IP アドレスヘルディングされたトラフィックは、例えば、1:1 ネットワークアドレス変換（NAT）を介して変換され、リソースインスタンスのそれぞれのプライベート IP アドレスへ転送される。

40

【 0 0 3 3 】

いくつかの公共 IP アドレスは、プロバイダネットワーク基礎構造によって、特定のリソースインスタンスに割り振ることができ、これらの公共 IP アドレスは、標準公共 IP アドレス、または簡単に標準 IP アドレスと呼ぶことができる。少なくともいくつかの実施形態において、リソースインスタンスの標準 IP アドレスからプライベート IP アドレスへのマッピングは、リソースインスタンスタイプ全てのデフォルトの起動構成である。

【 0 0 3 4 】

少なくともいくつかの公共 IP アドレスは、プロバイダネットワーク 200 のクライアントによって割り当てられる、または取得され得て、クライアントは次いで、それらの割

50

り当てられた公共IPアドレスを、クライアントに割り当てられた特定のリソースインスタンスに割り振ることができる。これらの公共IPアドレスは、クライアント公共IPアドレス、または簡単にクライアントIPアドレスと呼ぶことができる。標準のIPアドレスの場合のように、プロバイダネットワーク200によってリソースインスタンスに割り振られるのではなく、クライアントIPアドレスは、例えば、サービスプロバイダによって提供されるAPIを介して、クライアントによってリソースインスタンスに割り振られる場合がある。標準のIPアドレスとは異なり、クライアントIPアドレスは、クライアントアカウントに割り当てられ、必要または所望に応じて、それぞれのクライアントによって、他のリソースインスタンスに再マッピングされ得る。クライアントIPアドレスは、特定のリソースインスタンスではなく、クライアントのアカウントと関連付けられ、クライアントは、クライアントがそれを解放することを選択するまで、そのIPアドレスを制御する。従来の静的IPアドレスとは異なり、クライアントIPアドレスは、クライアントが、クライアントの公共IPアドレスを、クライアントのアカウントと関連付けられたいずれかのリソースインスタンスに再マッピングすることによって、リソースインスタンスまたは可用性ゾーンの障害をマスクすることを可能にする。クライアントIPアドレスは、例えば、クライアントが、クライアントIPアドレスを交換リソースインスタンスに再マッピングすることによって、クライアントのリソースまたはソフトウェアで問題を回避するように設計することを可能にする。

10

20

30

40

50

【0035】

図3は、少なくともいくつかの実施形態に従う、例示的データセンター（例えば、IPTunneling技術を使用してネットワークサブストレート上にオーバーレイネットワークを実装するデータセンター）を示す。この実施例に示されるように、プロバイダデータセンター300は、ルータ、スイッチ、ネットワークアドレストランスレータ（NAT）等のネットワークデバイス312を含むネットワークサブストレートを含むことができる。少なくともいくつかの実施形態は、カプセル化されたパケットがトンネルを使用してネットワークサブストレート310を通過され得るオーバーレイネットワークを提供するために、インターネットプロトコル（IP）トンネリング技術を採用することができる。IPTunneling技術は、ネットワーク（例えば、図3のデータセンター300内のローカルネットワーク）上にオーバーレイネットワークを作成するためのマッピング及びカプセル化システムを提供することができる。オーバーレイ層（公共IPアドレス）及びネットワークサブストレート310層（プライベートIPアドレス）に別々の名前空間を提供することができる。オーバーレイ層内のパケットは、それらのトンネルサブストレートターゲット（プライベートIPアドレス）が何であるべきかを決定するために、マッピングディレクトリ（例えば、マッピングサービス330によって提供される）と照合することができる。IPTunneling技術は、仮想ネットワークトポロジ（オーバーレイネットワーク）を提供し、クライアントに提示されるインターフェース（例えば、サービスAPI）は、オーバーレイネットワークに付加されるので、クライアントが、クライアントがパケットを送信したいIPアドレスを提供すると、IPアドレスは、IPオーバーレイアドレスがどこにあるか知っているマッピングサービス（例えば、マッピングサービス330）と通信することによって、仮想空間内で稼動する。

【0036】

少なくともいくつかの実施形態において、IPTunneling技術は、IPオーバーレイアドレス（公共IPアドレス）をサブストレートIPアドレス（プライベートIPアドレス）にオーバーレイし、2つの名前空間の間のトンネルでパケットをカプセル化し、パケットを、トンネルを介して正しいエンドポイントに配信ことができ、カプセル化がパケットから解除される。図3には、ホスト320A上の仮想マシン（VM）324Aから中間ネットワーク350上のデバイスへの例示的なオーバーレイネットワークトンネル334A、及びホスト320B上のVM324Bと、ホスト320C上のVM324Cとの間の例示的なオーバーレイネットワークトンネル334Bが示される。いくつかの実施形態において、パケットは、送信前にオーバーレイネットワークパケット形式にカプセル化

することができ、オーバーレイネットワークパケットは、受信後に解除することができる。他の実施形態において、パケットをオーバーレイネットワークパケットにカプセル化する代わりに、オーバーレイネットワークアドレス（公共IPアドレス）を、送信前にパケットのサブストレートアドレス（プライベートIPアドレス）に埋め込むことができ、受信後にパケットアドレスから解除することができる。実施例として、オーバーレイネットワークは、公共IPアドレスとして、32ビットIPv4（インターネットプロトコルバージョン4）アドレスを使用して実装することができ、IPv4アドレスは、プライベートIPアドレスとして、サブストレートネットワーク上で使用される128ビットIPv6（インターネットプロトコルバージョン6）アドレスの一部として埋め込むことができる。

10

【0037】

図3を参照すると、本明細書に記載される分散ロックマネージャの実施形態を実装することができる少なくともいくつかのネットワークは、複数のオペレーティングシステムが、ホストコンピュータ（例えば、図3のホスト320A及び320B）上で同時に、すなわち、ホスト320上の仮想マシン（VM）324として、稼動することを可能にするハードウェア仮想化技術を含むことができる。VM324は、例えば、ネットワークプロバイダのクライアントに賃借またはリースされ得る。ハイパーバイザ、または仮想マシンモニタ（VMM）322は、ホスト320上で、ホスト上のVM324に仮想プラットフォームを提示し、VM324の実行を監視する。各VM324には、1つ以上のプライベートIPアドレスを提供することができ、ホスト320上のVMM322は、ホスト上のVM324のプライベートIPアドレスを認識することができる。マッピングサービス330は、ローカルネットワーク上のIPアドレスに対応しているルータまたは他のデバイスの全てのネットワークIPプレフィックス及びIPアドレスを認識することができる。これは、複数のVM324に対応しているVMM322のIPアドレスを含む。マッピングサービス330は、例えば、サーバーシステム上で集中管理されてもよく、または代替として、ネットワーク上の2つ以上のサーバーシステムまたは他のデバイス間で分散されてもよい。ネットワークは、例えば、データセンター300ネットワーク内部の異なるホスト320上のVM324間でデータパケットをルーティングするために、例えば、マッピングサービス技術及びIPTunnelling技術を使用することができ、内部ゲートウェイプロトコル（IGP）は、そのようなローカルネットワーク内部のルーティング情報を交換するために使用することができることに留意されたい。

20

30

【0038】

加えて、プロバイダデータセンター300ネットワーク（時には自律システム（AS）と呼ばれる）のようなネットワークは、パケットをVM324からインターネットの送信先へ、及びインターネットの発信元からVM324へルーティングするために、マッピングサービス技術、IPTunnelling技術、及びルーティングサービス技術を使用することができる。外部ゲートウェイプロトコル（EGP）またはボーダゲートウェイプロトコル（BGP）は一般に、インターネット上の発信元と送信先との間のインターネットルーティングのために使用されることに留意されたい。図3は、少なくともいくつかの実施形態に従う、リソース仮想化技術を提供し、かつインターネットトランジットプロバイダに接続するエッジルータ314を介して完全なインターネットアクセスを提供するネットワークを実装する例示的なプロバイダデータセンター300を示す。プロバイダデータセンター300は、例えば、ハードウェア仮想化サービスを介して、仮想コンピューティングシステム（VM324）を実装する能力、及びストレージ仮想化サービスを介して、ストレージリソース318上に仮想化されたデータストア316を実装する能力をクライアントに提供することができる。

40

【0039】

いくつかの実施形態において、データセンター300は、例えば、データセンター300内のホスト320上のVM324からインターネット送信先へ、及びインターネット発信元からVM324へパケットをルーティングするために、仮想化されたリソースへ及び

50

仮想化されたリソースからのトラフィックをルーティングするために、IPトンネリング技術、マッピングサービス技術、及びルーティングサービス技術を実装することができる。インターネットの発信元及び送信先は、例えば、中間ネットワーク340に接続されたコンピューティングシステム370、及び（例えば、ネットワーク350を中間トランジットプロバイダに接続するエッジルータ（単数または複数）314を介して）中間ネットワーク340に接続するローカルネットワーク350に接続されたコンピューティングシステム352を含むことができる。プロバイダデータセンター300ネットワークはまた、例えば、データセンター300内のホスト320上のVM324から同じホスト上またはデータセンター300内の他のホスト320上の他のVM324へ、データセンター300内のリソース間でパケットをルーティングすることもできる。

10

【0040】

データセンター300を提供するサービスプロバイダはまた、データセンター300に類似するハードウェア仮想化技術を含み、中間ネットワーク340にも接続することができる、追加のデータセンター（単数または複数）360を提供することができる。パケットは、例えば、データセンター300内のホスト320上のVM324から、別の類似のデータセンター360内の別のホスト上の別のVMへ、及びこの逆等、データセンター300から他のデータセンター360へ転送することができる。

【0041】

上記に複数のオペレーティングシステムが、ホスト上の仮想マシン（VM）としてホストコンピュータ上で同時に稼動することを可能にし、VMは、ネットワークプロバイダのクライアントに賃借またはリースされ得る、ハードウェア仮想化技術を説明したが、ハードウェア仮想化技術はまた、同様な様式でネットワークプロバイダのクライアントに対する仮想化されたリソースとして、例えば、ストレージリソース318等、他のコンピューティングリソースを提供するためにも使用され得る。

20

【0042】

公共ネットワークは、広義には、複数のエンティティ間へのオープンアクセス及び相互接続を提供するネットワークとして定義され得ることに留意されたい。インターネット、またはワールドワイドウェブ（WWW）は、公共ネットワークの一例である。共有ネットワークは、広義には、アクセスが一般に限定されない公共ネットワークとは対照的に、アクセスが2つ以上のエンティティに限定されるネットワークとして定義され得る。共有ネットワークは、例えば、1つ以上のローカルエリアネットワーク（LAN）及び/またはデータセンターネットワーク、あるいは広域ネットワーク（WAN）を形成するために相互接続される2つ以上のLANまたはデータセンターネットワークを含むことができる。共有ネットワークの例として、企業ネットワーク及び他のエンタープライズネットワークを挙げることができるが、これらに限定されない。共有ネットワークは、ローカルエリアを対象とするネットワークからグローバルネットワークまでの範囲のいずれかであり得る。共有ネットワークは、少なくともいくつかのネットワーク基礎構造を公共ネットワークと共有することができ、共有ネットワークは、他のネットワーク（単数または複数）と共有ネットワークとの間のアクセスを制御して、公共ネットワークを含むことができる、1つ以上の他のネットワークに連結され得ることに留意されたい。共有ネットワークはまた、インターネットのような公共ネットワークとは対照的に、プライベートネットワークとして考えることもできる。実施形態において、共有ネットワークまたは公共ネットワークのいずれかが、プロバイダネットワークとクライアントネットワークとの間の中間ネットワークとして役目を果たすことができる。

30

40

【0043】

いくつかの実施形態において、本明細書に記載されるDLMは、図2または図3に示され、前述した例示的なプロバイダネットワーク環境のうちの1つ等、分散コンピューティング環境（例えば、仮想コンピューティングリソース及び/またはサービスをクライアントに提供するクラウドコンピューティング環境）において実装され得る。いくつかの実施形態において、そのようなシステムに実装されたDLMの構成要素は、（例えば、クラウ

50

ドコンピューティング環境で実行しているクライアントアプリケーションに割り当てられた以外の、あるいは相互に通信する及び/またはロックによって保護された共有リソースにアクセスするためにクライアントアプリケーションの構成要素によって使用される、ネットワーク接続性及び/または処理能力を使用して)制御プレーンネットワーク上で相互に通信することができる。

【0044】

いくつかの実施形態において、複数のリソースインスタンスが、クライアントの代わりに分散アプリケーションを実装するためにクラウドコンピューティング環境において実行している場合がある。前述のように、クラウドコンピューティング環境は、各アプリケーション(及び/または各仮想プライベートネットワーク)がその独自の名前空間を有することができる、マルチテナント環境であり得る。いくつかの実施形態において、各クライアントは、ネットワーク接続及び/または処理能力(帯域幅)のその独自の割当を有することができる。例えば、データプレーンネットワークのネットワーク接続性及び/または処理能力は、様々なクライアントの使用のためにプロビジョニング(例えば、専用または予約)され得る。いくつかの実施形態において、DLMの1つ以上の構成要素(またはインスタンス)はまた、リソースインスタンスのうちの一つが実行している各ノード上でも実行されている場合があり、これらの構成要素は、相互に通信するためにクライアントアプリケーションに割り当てられた以外のネットワーク接続及び/または処理能力を使用することができる。例えば、様々な実施形態において、ノードあたり1つのDLM構成要素(またはインスタンス)、または各ノード上の顧客あたり1つのDLM構成要素(またはインスタンス)が存在し得る。

10

20

【0045】

いくつかの実施形態において、クライアントアプリケーションの構成要素が、クライアントアプリケーションの構成要素がそれらの通常の作業(例えば、ロック管理以外の作業)の一部として、相互に通信するネットワーク接続とは異なる接続機構上から同じノード上のDLM構成要素/インスタンスへのローカルAPI呼び出しを行うことができる。様々なロック管理操作(例えば、ロックを作成、ロックの1つ以上のプロパティを指定、ロックのリストを取得、ロックをサブスクライブ、ロックを取得、ロックを解放、またはロックを消去する操作)を起動するためにローカルAPI呼び出しを行うことによって、クライアントアプリケーションの構成要素は、DLMによって管理されるロックに関与することができる。

30

【0046】

いくつかの実施形態において、本明細書に記載されるクライアントアプリケーション及び他のプロセスの全てを実行している物理的コンピュータが存在する一方、クライアントアプリケーションは、物理的コンピュータ上で仮想マシンとして稼動している場合がある。例えば、これらの仮想マシンの作成を管理する、これらの仮想マシンのリソースをプロビジョニングする、及び/またはクライアント及び/またはそれらのアプリケーションの代わりに他の管理タスク(例えば、リソース使用量、顧客アカウント、サービスの請求等を監視する)を実施するように構成されるクラウドコンピューティング環境の内部プロセスは、クラウドコンピューティング環境内の制御プレーン層(またはハイパーバイザ)内で実行することができる。対照的に、クライアントアプリケーション(例えば、アプリケーション構成要素を実装する各リソースインスタンス)は、クラウドコンピューティング環境のデータプレーン層内で実行することができる。これらの層の下には、いくつかの実施形態において、各ホストノードに対して(または複数のホストノードに対して)1つの物理的ネットワークカードだけが存在し得るが、各リソースインスタンスは、あたかもそれがその独自のネットワーク(例えば、仮想ネットワーク)を有するように実行することができる。いくつかの実施形態において、各リソースインスタンスは、その独自のデータプレーンネットワーク接続(単数または複数)を有することができるが、これらのデータプレーンネットワーク接続に依存する必要なく、ローカルAPI呼び出し(例えば、同じノード上のDLM構成要素への呼び出し)を行うことができる。

40

50

【 0 0 4 7 】

いくつかの実施形態において、DLMはまた、制御プレーン（ハイパーバイザ）層上で稼動しているプロセスとして実装することができる。したがって、それは、クライアントプロセスが認識せず、アクセスを有さないネットワークアクセスを有することができる。そのような実施形態において、DLMは、データプレーンのいかなるリソース（例えば、コンピュートインスタンスまたはネットワーク接続/大域幅）も消費しない場合があり、それらのリソースに対してクライアントアプリケーションと競争しない場合があり、制御プレーンのリソースだけを消費する場合がある。様々な実施形態において、DLM構成要素（またはインスタンス）は、様々な基底のネットワーク及びネットワーク機構のいずれかを使用して、作成したロックのロック状態情報を共有することができる。一実施形態において、DLM構成要素は、ロック状態情報が共有される機構として、クラスタ通信（例えば、InfiniBand（登録商標）アーキテクチャ仕様に準じる通信リンク）のために設計された高速相互接続を採用することができる。例えば、ロックの状態に変化が生じると、ロックを変更したDLM構成要素（または別のDLM構成要素）は、ロック状態が変化したことを、1つ以上の他のDLM構成要素（例えば、ロックのサブスクリバである任意のクライアントアプリケーション構成要素と同じノード上で実行しているDLM構成要素）に通知することができ、及び/または変更されたロック状態値をサブスクリバに通信する。様々な実施形態において、ロックを変更するDLM構成要素（または別のDLM構成要素）は、永続的データストア（例えば、クラウドコンピューティング環境内）内のロックのロック状態情報のコピーを更新することに責任を負う場合がある。

10

20

【 0 0 4 8 】

図4は、いくつかの実施形態に従う、制御プレーンにおいて分散ロックマネージャを実装する例示的なサービスプロバイダシステムを示すブロック図である。いくつかの実施形態において、プロバイダシステム（図4においてプロバイダシステム400として示される）は、図3に示されるプロバイダデータセンター300に類似し得る。例えば、図4の各ホストノード上の制御プレーンは、図3に示されるハイパーバイザまたは仮想マシンモニタの機能性のうちのいくつかまたは全てを実装することができる。同様に、図4の各ホストノード上のデータプレーンで実行しているインスタンスは、図3に示された仮想マシンの機能性のうちのいくつかまたは全てを実装する仮想コンピュートインスタンスであり得る。

30

【 0 0 4 9 】

より具体的には、示された実施例において、プロバイダシステム400は、複数の仮想化されたデータストア（単数または複数）440と、ホストノード420及び430（その各々がデータプレーン部分と制御プレーン部分とを含む）と、データプレーンネットワーク410と、制御プレーンネットワーク415（異なる実施形態において、データプレーンネットワーク410と異なる物理的ハードウェア上に実装される場合または実装されない場合がある）とを含むことができる。

【 0 0 5 0 】

この実施例において、インスタンス428a~428b、及び438a~438bは、システムの制御プレーンのそれぞれの部分（424及び434として示される）で実行し、1つ以上のクライアントアプリケーションまたはプロセスを実装することができ、そのうちの少なくともいくつかは、ロック（例えば、分散ロックマネージャ（DLM）によって管理されるロック）によって保護される共有リソースにアクセスするように構成される。この実施例において、アプリケーション/プロセスを実装するために、これらのインスタンスは、データプレーンネットワーク410上で相互に、及び/または他のアプリケーション構成要素（例えば、仮想化されたデータストア（単数または複数）440のストレージデバイス445）と通信するように構成することができる。

40

【 0 0 5 1 】

図4に示される実施例において、ホストノード420上で実行しているインスタンス428a~428bは、様々なロック管理操作を開始するために、DLMインスタンス42

50

6 への API 呼び出しを行うように構成することができ、ホストノード 430 上で実行しているインスタンス 438a ~ 438b は、様々なロック管理操作を開始するために、DLM インスタンス 436 への API 呼び出しを行うように構成することができる。この実施例において、DLM (DLM インスタンス 426 及び 436 を含み、かつ仮想化されたデータストア (単数または複数) 440 へのアクセスを有する) は、システムの制御プレーンに内 (例えば、424 及び 434 として制御プレーンのそれぞれの部分内) で実行し、その構成要素は、インスタンス 428a ~ 428b、及び 438a ~ 438b のために、共有リソースの 1 つ以上のロックを管理するように、制御プレーンネットワーク 415 上で相互に通信するように構成することができる。

【0052】

図 4 に示されるように、いくつかの実施形態において、サービスプロバイダシステムは、(例えば、制御プレーンネットワーク 415 上の) 制御プレーン内に DLM コーディネータ構成要素 (例えば、DLM コーディネータ 450) を含むことができる。例えば、DLM コーディネータ構成要素 (制御プレーンのクラウドマネージャ構成要素のサブ構成要素であり得る) は、分散ロックサービスを提供するために一体に機能する際、DLM インスタンスのアクティビティのうち少なくともいくつかを管理及び / または調整することができる。様々な実施形態において、DLM インスタンスは、制御プレーンネットワーク 415 上で相互に及び / または DLM コーディネータ 450 と通信することができる。例えば、DLM コーディネータ 415 は、例えば、分散ロックマネージャが (例えば、最新状態の DLM インスタンスの各々によってローカルにキャッシュされたロック状態情報の全てを保つために) DLM によって管理されたロックの状態の一貫した (または最終的に一貫した) ビューを維持することを促進するために、DLM インスタンス 426 及び 436 へ、及び / または 426 と 436 との間でメッセージを仲介することができる。他の実施形態において、DLM インスタンスは (少なくとも時々)、ロック状態情報を共有するために、及び / または DLM インスタンスの各々にローカルに記憶されたロック状態情報が最新状態に保たれていることを保証するために、制御プレーンネットワーク 415 上で直接相互に通信することができる。いくつかの実施形態において、DLM コーディネータ 450 は、ロック状態情報の永続的データストアを維持するように構成することができる (図示せず)。

【0053】

仮想化されたリソースを提供する分散コンピューティングシステム内において分散ロックマネージャを実装するための方法の一実施形態が、図 5 のフロー図によって示される。510 に示されるように、この実施例において、方法は、分散コンピューティングシステムが、クライアントから仮想コンピューティングサービスの要求を受信することを含むことができる。要求に回答して、方法は、520 のように、分散コンピューティングシステムが、分散アプリケーションを実装するために、それぞれのコンピューティングノード上でクライアントのための 2 つ以上のリソースインスタンスをプロビジョニングすることと、これらのインスタンスが、データプレーンネットワーク上で相互に通信するように、システム内でそれらを構成することとを含むことができる。例えば、システムは、様々なコンピュータインスタンスまたは実行プラットフォームインスタンスをプロビジョニングすることができ、いくつかの実施形態において、データプレーンネットワーク上のプロビジョニングされたコンピューティングリソース、プロビジョニングされた記憶能力、プロビジョニングされたネットワーク接続及び / またはプロビジョニングされた処理能力 (例えば、帯域幅) を含むことができる。データプレーンネットワークは、分散アプリケーションの構成要素がその作業 (例えば、ロック管理以外の作業) を実行するために、相互に通信するネットワークであってもよい。

【0054】

この実施例に示されるように、方法は、分散コンピューティングシステムが、分散ロックマネージャ (DLM) を実装するために、コンピューティングノードの各々の上でリソースインスタンスをプロビジョニングすることと、これらが制御プレーンネットワーク上

10

20

30

40

50

で相互に通信するなど、530のようにシステム内のこれらのDLMインスタンスを構成することを含むことができる。例えば、システムは、様々なコンピュータインスタンスまたは実行プラットフォームインスタンスをプロビジョニングすることができ、いくつかの実施形態において、制御プレーンネットワーク上のプロビジョニングされたコンピューティングリソース、プロビジョニングされた記憶能力、プロビジョニングされたネットワーク接続及び/またはプロビジョニングされた処理能力（例えば、帯域幅）を含むことができる。制御プレーンネットワークは、クライアントプロセス以外のプロセス（例えば、リソース使用量、クライアント請求、認証サービス、及び/または分散コンピューティングシステムの他の管理タスク）が実行及び/または相互に通信するネットワークであり得ることに留意されたい。

10

【0055】

540に示されるように、方法は、ロック管理操作の実施を開始するために、ローカルDLMインスタンス（例えば、呼び出しが公共またはプライベートネットワーク接続上を移動する必要がないように、アプリケーション構成要素と同じコンピューティングノード上で実行しているDLMインスタンス）へのAPI呼び出しを使用する分散アプリケーション（例えば、分散アプリケーションのアプリケーション構成要素）を含むことができる。例えば、アプリケーション構成要素は、ロックの作成、ロックの取得、ロックの解放、またはこれら及び他の操作を定義するAPIに従う別のロック管理操作を開始することができる。いくつかの実施形態において、方法は、550のように、ローカルDLMインスタンスが、要求されたロック管理操作を実施することと、他のDLMインスタンスとロック状態情報（例えば、変更されたロック値または操作を実施したことから生じる他の状態情報）を共有することを含むことができる。本明細書に記載されるように、DLMの構成要素は、いくつかの実施形態において、分散コンピューティングシステムのデータプレーンネットワーク上ではなく、システムの制御プレーンネットワーク上で相互に通信することができる。方法はまた、560のように、DLMインスタンス（単数または複数）が、ロックの状態情報に対する何らかの変化について、ロックのサブスクリバに通知することも含むことができる。例えば、各DLMインスタンスは、それらがサブスクリバされているロックのロック情報における何らかの変化について同じノード上のプロセスに通知すること、及び/またはロックの他のサブスクリバが実行しているノード上で実行している他のDLMインスタンスへのロック状態情報の変化を伝播することに責任を負う場合がある。

20

30

【0056】

本明細書に記載される分散ロックマネージャのアプリケーションプログラミングインターフェースは、以下を含むが、これらに限定されない、共有リソースのロックのための様々なロック管理操作を定義することができる。

- ・ロックを作成すること（これは、ロックとリソースの関連付けを含まない場合があり、呼び出し側の責任となる場合があることに留意されたい）。
- ・ロックの所定のプロパティの値を設定すること（例えば、ロックの共有プロパティまたはロックモードを指定する）。
- ・クライアントがサブスクリブすることができるロックのリストを取得すること。
- ・ロックの状態情報を表示すること（例えば、そのようなAPIが、ロックの複数のサブスクリバが、ロックの状態を表示することを可能にすることができ、状態がアトミックに変化するため、サブスクリバは、ロック状態の一貫してビューを得るようになる）。
- ・DLMがロックのハンドルまたはロックの状態情報を返すことに応答して、所与のロックのロックグループのメンバーになる（すなわち、所与のロックをサブスクリブする）ことを要求すること。

40

【0057】

いくつかの実施形態において、ロックは、エラスティックネットワークインターフェース（ENI）または別のタイプの仮想ネットワークインターフェースを管理するために使用することができる。例えば、いくつかの既存システムにおいて、ネットワークインター

50

フェースを移動することが可能である場合、これによって、いくつかの既存のシステムのように、システムのロックを管理するための媒体としてネットワーク自体を使用しようとするのがより困難である。すなわち、ネットワークインターフェースが、特定のアプリケーションまたはリソースインスタスの唯一のネットワークインターフェースである場合、かつそれを移動するために切断されなければならない場合、アプリケーションまたはリソースインスタスは、ロックマネージャから封鎖され得る。しかしながら、本明細書に記載されるシステムにおいて、分散ロックマネージャは、ネットワークインターフェースのロックを管理するために（例えば、ネットワークインターフェースまたはロックの所有権とのロックの関連付けを管理するために）同じ物理的または論理的ネットワークを使用しない場合がある。その代わりに、ロックマネージャは、その接続機構として、クラウドコンピューティング環境の制御プレーンネットワークを使用することができ、制御プレーンネットワークは、クライアントアプリケーションから隠すことができる。一実施例において、ENIは、インスタスの障害に回答して、（そのIPアドレス及び接続されたそのクライアントと共に）別のホストへ移動され得る。ENIを別のホストへ移動させると、それに伴いそのIPアドレスをもたすため、クライアントは、1つのIPアドレスだけを知ることが必要であり、クライアントは、ENIが接続されている場所を知る必要はない。従来のDLMソリューションでは、そのENIがそれ自体データプレーン接続機構であった場合、それを移動させることはできない（移動させると、そのデータプレーンへの接続がすべて失われるため）。

10

20

【0058】

別の例において、クラスタ化されたファイルシステムタイプのアプリケーションは、ディスクボリューム上にロックを作成することができ、予備のデータベースを起動させて稼働状態に保つことができる（例えば、一次ホスト上のデータベースをミラリングする）。この例において、一次ホストに障害が発生すると、対応するディスクボリュームを、一次ホストから切断し、予備のうちの1つ（例えば、二次ホスト）に再接続することができる。一般に、分散コンピューティングシステムに共有リソースが存在する場合は必ず、ロックの関与者に、関与者のうちの1つが特定の容量のロックを保持していることを示し、及び/またはそのロックの所有者の識別子を他の関与者に通信するためにロックを使用することが可能であるため、ロックをそれらのリソースと関連付けることが有用であり得る。

30

【0059】

クライアントアプリケーションまたはプロセスが、クラウドコンピューティング環境において実装される分散ロックマネージャの機能にアクセスするための方法の一実施形態が、図6のフロー図によって示される。610に示されるように、この実施例において、方法は、クラウドコンピューティング環境で実行しているクライアントアプリケーションまたはプロセスが、共有リソースのロックを作成するためにローカルDLMインスタス（例えば、同じコンピューティングノード上で実行しているDLMインスタス）へのAPI呼び出しを行うことと、ロックが作成されていることに回答して、新しく作成されたロックの識別子（例えば、ロックハンドル）を受信することを含むことができる。異なる実施形態において、アプリケーションは、分散アプリケーション、または別のアプリケーションまたはプロセスによってもアクセス可能であるリソースにアクセスする、単一のノード上で稼働しているアプリケーションまたはプロセスであり得る。いくつかの実施形態において、ロックを作成することは、ロックの1つ以上のプロパティ（例えば、共有プロパティ）の値を設定するために、1つ以上の追加のAPI呼び出しを行うことも含むことができる。

40

【0060】

この実施例に示されるように、方法は、620のように、クライアントアプリケーション/プロセスが、共有リソースのロックを取得するために、ローカルDLMインスタスへのAPI呼び出しを行うことを含むことができる。例えば、要求は、ロックの識別子（またはハンドル）を含むことができ、ロックを作成した同じアプリケーション/プロセスまたはロックに参加もする異なるアプリケーション/プロセスから受信され得る。クライ

50

アントアプリケーションは、任意の他のDLMインスタンスの存在または場所に対する可視性を有さない場合、または任意の他のDLMインスタンスとの通信能力を有さない場合があることに留意されたい。

【0061】

630からの肯定終了として示される、ロックが別のアプリケーションプロセスによって保持されている場合、方法は、635のように、クライアントアプリケーション/プロセスが、それが解放される（またはその他取得のために有効になる）までロック上の待ち行列に入ることまたはロックをポーリングすることを含むことができる。ロックが別のアプリケーション/プロセスによって保持されていない場合（630からの否定終了として示される）またはロックを保持する別のアプリケーション/プロセスによって解放された後）、方法は、640のように、クライアントアプリケーション/プロセスにロックが承諾されることと、共有リソースにアクセスすることを含むことができる。いくつかの実施形態において、ロックを承諾することは、永続的データストア内のロック状態のコピーをアトミックに更新すること、及び/または他のDLMインスタンス（これらの全ては制御プレーンネットワーク上で相互に通信する）のロック状態のローカルにキャッシュされたコピーを更新することを含むことができる。いくつかの実施形態において、アプリケーション/プロセスは、クラウドコンピューティング環境内のデータプレーンネットワーク上で共有リソースにアクセスする。

10

【0062】

図6に示されるように、何らかの時点で（例えば、アプリケーションまたはプロセスが、共有リソースに対するアクセスを必要としないとき）、方法は、650のように、クライアントアプリケーション/プロセスが、ロックを解放するために、ローカルDLMインスタンスへのAPI呼び出しを行うことを含むことができる。その後、方法は、660のように、別のクライアントアプリケーションまたはプロセスが、ロックの状態を表示及び/または変更するために、1つ以上のローカルDLMインスタンスへの1つ以上のAPI呼び出しを行うことを含むことができる。例えば、別のクライアントアプリケーションまたはプロセスは、ロックをクエリする（例えば、共有リソースがロックされているかを決定、及び/または現在の所有者を決定するために）、あるいはロックを取得（そしてその後解放）するために、API呼び出しを行うことができる。

20

【0063】

いくつかの実施形態において、本明細書に記載される分散ロックマネージャは、ロックが、クラウドコンピューティング環境の制御プレーンの外側で拡張されることを可能にすることができる。例えば、いくつかの実施形態において、顧客は、クラウドコンピューティング環境の外側でロックを拡張するために、（例えば、クライアントDLMエージェントをダウンロードしてインストールするために）その独自の施設上でソフトウェアプロセスをダウンロード及びインストールすることが可能であってもよい。そのような実施形態において、クライアントネットワーク上のホストコンピューティングノード上で稼働しているアプリケーションは、様々なロック管理操作を実施するために、クライアントDLMへのAPI呼び出しを行うことができる。いくつかの実施形態において、クライアントがAPI呼び出しを行う機構は、セキュアチャンネルであり、クライアントは、クラウドコンピューティング環境の外側でエージェントからのAPI呼び出しを行うことができる。

30

40

【0064】

一例において、顧客は、クラウドコンピューティング環境内部にアプリケーションを構築することができるが、アプリケーションコントローラ（ロックマネージャを含む）をオンプレミス（例えば、クライアントネットワークのローカルにあるホストノード上で実行しているマシン上）に所在させたい場合がある。より具体的な例として、顧客は、クラウドコンピューティング環境内で使用するために10のリソースインスタンスをプロビジョニングすることを要求する場合があるが、顧客は、クラウドコンピューティング環境の外側に所在するリソース（例えば、そのローカルマシン上のファイル）と関連付けられたロックを有することも希望する場合がある。本明細書に記載される分散ロックマネージャは

50

、顧客が、そのファイルのロックを作成し、ロック（またはファイル）の所定の共有プロパティを設定することを可能にすることができ、そのロックに対する全てのサブスクリバが、ロックの一貫した状態情報を表示することを可能にすることができる。

【0065】

別の例において、顧客は、ローカルマシン上で稼動しているアプリケーションを有する場合があるが、ローカルマシン上の障害の場合に使用されるように、クラウドコンピューティング環境内のいくつかのリソースインスタンスをプロビジョニングすることを所望する場合がある。この例では、アプリケーションと関連付けられたロックが存在する場合があり、障害の場合に、ロックは、クラウドコンピューティング環境内のリソースインスタンスに移動され得る。

10

【0066】

DLMがクラウドコンピューティング環境とローカルクライアントネットワークとの間に及ぶいくつかの実施形態において、クライアントアプリケーション及び/またはDLMは、VCP（仮想プライベートクラウド）内部に実装され得ることに留意されたい。

【0067】

図7は、いくつかの実施形態に従う、制御プレーンに分散ロックマネージャを実装し、その環境の外側（例えば、サービスプロバイダシステムの外側）で実行しているプロセスに対して、分散ロックマネージャを公開する例示的サービスプロバイダシステムを示すブロック図である。様々な実施形態において、プロバイダシステム700は、図3に示されるプロバイダデータセンター300及び/または図4に示されるプロバイダシステム400に類似し得る。例えば、図7の各ホストノード上の制御プレーンは、図3に示されたハイパーバイザまたは仮想マシンモニタの機能性のうちのいくつかまたは全てを実装することができる。同様に、図7の各ホストノード上のデータプレーンで実行しているインスタンスは、図7に示された仮想マシンの機能性のうちのいくつかまたは全てを実装する仮想コンピュートインスタンスであり得る。

20

【0068】

より具体的には、示された実施例において、プロバイダシステム700は、複数の仮想化されたデータストア（単数または複数）740と、ホストノード720及び730（その各々がデータプレーン部分と制御プレーン部分とを含む）と、データプレーンネットワーク710と、制御プレーンネットワーク715（異なる実施形態において、データプレーンネットワーク710と異なる物理的ハードウェア上に実装される場合または実装されない場合がある）とを含むことができる。

30

【0069】

この実施例において、インスタンス728a~728b、及び738a~738bは、システムの制御プレーンのそれぞれの部分（724及び734として示される）で実行し、1つ以上のクライアントアプリケーションまたはプロセスを実装することができ、そのうちの少なくともいくつかは、ロック（例えば、分散ロックマネージャ（DLM）によって管理されるロック）によって保護される共有リソースにアクセスするように構成される。この実施例において、アプリケーション/プロセスを実装するために、これらのインスタンスは、データプレーンネットワーク710上で相互に、及び/または他のアプリケーション構成要素（例えば、仮想化されたデータストア（単数または複数）740のストレージデバイス745）と通信するように構成することができる。

40

【0070】

図7に示される実施例において、ホストノード720上で実行しているインスタンス728a~728bは、様々なロック管理操作を開始するために、DLMインスタンス726へのAPI呼び出しを行うように構成することができ、ホストノード730上で実行しているインスタンス738a~738bは、様々なロック管理操作を開始するために、DLMインスタンス736へのAPI呼び出しを行うように構成することができる。この実施例において、DLM（DLMインスタンス726及び736を含み、かつ仮想化されたデータストア（単数または複数）740へのアクセスを有する）は、システムの制御プレ

50

ーンにおいて（例えば、724及び734として制御プレーンのそれぞれの部分において）実行し、その構成要素は、インスタンス728a～728b、及び738a～738bのために、共有リソースの1つ以上のロックを管理するように、制御プレーンネットワーク715上で相互に通信するように構成することができる。

【0071】

図7に示されるように、いくつかの実施形態において、サービスプロバイダシステムは、（例えば、制御プレーンネットワーク715上の）制御プレーンにDLMコーディネータ構成要素（例えば、DLMコーディネータ750）を含むことができる。例えば、DLMコーディネータ構成要素（制御プレーンのクラウドマネージャ構成要素のサブ構成要素であり得る）は、分散ロックサービスを提供するために一体に機能する際、DLMインスタンスのアクティビティのうち少なくともいくつかを管理及び/または調整することができる。様々な実施形態において、DLMインスタンスは、制御プレーンネットワーク715上で相互に及び/またはDLMコーディネータ750と通信することができる。例えば、DLMコーディネータ715は、例えば、分散ロックサービスが（例えば、最新状態のDLMインスタンスの各々によってローカルにキャッシュされたロック状態情報の全てを保つために）DLMによって管理されたロックの状態の一貫した（または最終的に一貫した）ビューを維持することを促進するために、DLMインスタンス726及び736へ、及び/または726と736との間でメッセージを仲介することができる。他の実施形態において、DLMインスタンスは（少なくとも時々）、ロック状態情報を共有するために、及び/またはDLMインスタンスの各々にローカルに記憶されたロック状態情報が最新状態に保たれていることを保証するために、制御プレーンネットワーク715上で直接相互に通信することができる。いくつかの実施形態において、DLMコーディネータ750は、ロック状態情報（図示せず）の永続的データストアを維持するように構成することができる。

10

20

【0072】

図4に示された実施例とは異なり、図7に示されたプロバイダシステム700のDLMに対するインターフェースは、プロバイダシステム700の外側で実行しているプロセスに公開され得る。この実施例において、クライアントネットワーク770上のホストノード775上で実行している様々なアプリケーション及び/またはプロセスは、様々なAPI760を通じて、（例えば、それらのアプリケーション/プロセス及びホストノード720またはホストノード730上で実行しているアプリケーション/プロセスによって共有リソースのロック上でロック管理操作を開始するために）DLMのロック機構に参与することが可能であり得る。本明細書に記載されるように、いくつかの実施形態において、DLMクライアントエージェントは、ホスト775上で実行しているアプリケーションまたはプロセスが、DLMクライアントエージェントへのローカルAPI呼び出しを使用して、DLMによって管理されたロックにアクセスすることを可能にするために、ホストノード775上でインスタンス化され得る。そのような実施形態において、ローカルDLMクライアントエージェントが様々なローカルAPI呼び出しに応答して、これらの構成要素との通信を処理するので、プロバイダシステムの外側で実行しているアプリケーション/プロセスは、プロバイダシステム700内部に実装されたDLM構成要素のネットワークアドレス、ネットワーク接続、及び/またはネットワークリソースについて全く知る必要なく、これらのロックの管理をサブスクライブ及び/または関与することができる。

30

40

【0073】

図7に示された実施例において、分散ロックマネージャ（または分散ロックマネージャの構成要素によって提供された分散ロックサービス）に対して、クライアントネットワーク770上のホストノード775上で実行しているアプリケーション/プロセスによって行われたAPI呼び出し（例えば、API760に準拠するAPI呼び出し）は、DLMコーディネータ750によって仲介され得て、API呼び出しがルーティングされるべき特定のホストノード（または、より具体的には、特定のホストノード上のDLMインスタンス）を決定することができる。他の実施形態において、クライアントネットワーク77

50

0上のホストノード775上で実行しているアプリケーション/プロセスによって行われたAPI呼び出しは、DLMコーディネータ750を通じてルーティングされることなく、特定のホストノード上の特定のDLMインスタンスへ向けることができる。例えば、いくつかの実施形態において、通信チャンネル（例えば、制御プレーンネットワーク715上のオーバーレイネットワークトンネル、または制御プレーンネットワーク715上の別の種類の通信チャンネル）は、（例えば、DLMコーディネータ750または別の制御プレーン構成要素によって仲介された通信を通じて）ホストノード775上で実行しているアプリケーション/プロセスと、特定のロックをサブスクライブしているアプリケーション/プロセスの結果としての特定のDLMインスタンスとの間に確立され得る。サブスク립ションが承諾された、及び/または通信チャンネルが確立された後、ホストノード775上で実行しているアプリケーション/プロセスと特定のDLMインスタンス（例えば、API760に準拠する様々なAPI呼び出し）との間のその後の通信は、DLMコーディネータ750を通じてではなく、この通信チャンネル上で（例えば、直接）発生することができる。

10

20

30

40

50

【0074】

クラウドコンピューティング環境の外側で実行するクライアントアプリケーションまたはプロセスが、クラウドコンピューティング環境において実装される分散ロックマネージャの機能にアクセスするための方法の一実施形態が、図8のフロー図によって示される。810に示されるように、この実施例において、方法は、クラウドコンピューティング環境の外側で実行しているクライアントプロセスが、クラウドコンピューティング環境内で実行している1つ以上のプロセスと共有されるリソースのロックをサブスクライブするために、クラウドコンピューティング環境内で実行する分散ロックマネージャのローカルエージェントへのAPI呼び出しを行うことを含むことができる。例えば、クライアントプロセスは、呼び出しプロセスと同じコンピューティングノード上で実行しているクライアントDLMエージェントへのAPI呼び出しを行うことができる。その呼び出しに回答して、方法は、820のように、クラウドコンピューティング環境の外側で実行しているクライアントプロセスが、ロックの識別子（例えば、ロックハンドル）またはロックの値を受信することを含むことができ、その後、共有リソースのロックを取得するために、分散ロックマネージャのローカルエージェントへのAPI呼び出しを行うことができる。いくつかの実施形態において、これによって、（例えば、ロックのロック状態情報を他のサブスクライバと共有するために）ローカルクライアントDLMエージェントと、クラウドコンピューティング環境で実行しているDLMの構成要素との間の通信を開始することができる。

【0075】

この実施例に示されるように、830からの肯定終了として示される、ロックが別のアプリケーションプロセスによって保持されている場合、方法は、835のように、クライアントプロセスが、それが解放される（またはその他取得のために有効になる）までロック上の待ち行列に入ることまたはロックをポーリングすることを含むことができる。ロックは、クラウドコンピューティング環境の外側で実行している別のプロセスによって、またはクラウドコンピューティング環境の内側で実行しているプロセスによって保持され得ることに留意されたい。ロックが別のプロセスによって保持されていない場合（830からの否定終了として示される）またはロックを保持する別のプロセスによって解放された後）、方法は、840のように、クライアントプロセスにロックを承諾することと、共有リソースにアクセスすることを含むことができる。いくつかの実施形態において、ロックを承諾することは、永続的データストア内のロック状態のコピーをアトミックに更新すること、並びに/あるいはローカルクライアントDLMエージェント及び/またクラウドコンピューティング環境内の他のDLM構成要素（これらの全ては制御プレーンネットワーク上で相互に通信する）のロック状態のローカルにキャッシュされたコピーを更新することを含むことができる。いくつかの実施形態において、ロックを承諾することは、DLMが、プロセスにロックが承諾されたことを示すロック状態値を返すことを含むことがで

きる。

【0076】

図8に示されるように、何らかの時点で（例えば、プロセスが、共有リソースに対するアクセスを必要としないとき）、方法は、850のように、クライアントプロセスが、ロックを解放するために、ローカルクライアントDLMエージェントへのAPI呼び出しを行うことを含むことができる。いくつかの実施形態において、これによって、（例えば、ロックの状態情報を他のサブクライアントと共有するために）ローカルクライアントDLMエージェントと、クラウドで実行しているDLMの1つ以上の構成要素との間の通信を開始することができる。その後、方法は、860のように、別のプロセスが、ロックの状態を表示及び/または変更するために、クラウドコンピューティング環境の外側で実行しているローカルクライアントDLMエージェントまたはクラウドコンピューティング環境の内側で実行しているローカルDLMインスタンスへの1つ以上のAPI呼び出しを行うことを含むことができる。すなわち、他のプロセスは、そのプロセスに対してローカルであるDLMエージェントまたはインスタンス（例えば、クラウドコンピューティング環境内またはクラウドコンピューティング環境外側の、そのプロセスと同じホストノード上で実行しているDLMエージェントまたはインスタンス）へのAPI呼び出しを行うことができる。例えば、別のクライアントアプリケーションまたはプロセスは、ロックをクエリする（例えば、共有リソースがロックされているかを決定、及び/または現在の所有者を決定するために）、あるいはロックを取得（そしてその後解放）するために、API呼び出しを行うことができる。

10

20

【0077】

例示的システム

少なくともいくつかの実施形態において、本明細書に記載されるような分散ロックマネージャを実装するための技法のいくつかまたは全てを実装するサーバーが、図9に示されるようなコンピュータシステム900等、非一時的コンピュータアクセス可能媒体を含む、またはこれにアクセスするように構成される汎用コンピュータシステムを含むことができる。示される実施形態において、コンピュータシステム900は、入出力（I/O）インターフェース930を介して、システムメモリ920に連結される1つ以上のプロセッサ910を含む。コンピュータシステム900は、I/Oインターフェース930に連結されたネットワークインターフェース940をさらに含む。

30

【0078】

様々な実施形態において、コンピュータシステム900は、1つのプロセッサ910を含む単一プロセッサシステム、または複数（例えば、2、4、8、または別の適切な数）のプロセッサ910を含むマルチプロセッサシステムである。プロセッサ910は、命令を実行することが可能である任意の適切なプロセッサであり得る。例えば、様々な実施形態において、プロセッサ910は、x86、PowerPC、SPARC、またはMIPS ISA、あるいは任意の他の適切なISA等、様々な命令セットアーキテクチャ（ISA）のうちのいずれかを実装する汎用または埋め込みプロセッサであり得る。マルチプロセッサシステムにおいて、プロセッサ910の各々は一般に、同じISAを実装することができるが、必ずしも実装しなくてもよい。

40

【0079】

システムメモリ920は、プロセッサ910によってアクセス可能な命令及びデータを記憶するように構成され得る。様々な実施形態において、システムメモリ920は、静的ランダムアクセスメモリ（SRAM）、同期動的RAM（SDRAM）、不揮発性/フラッシュタイプメモリ、または任意の他の種類のメモリ等、任意の適切なメモリ技術を使用して実装され得る。示された実施形態において、分散ロックマネージャを実装するための上記の方法、技法、及びデータ等、1つ以上の所望される機能を実装するプログラム命令及びデータは、コード925及びデータ926として、システムメモリ920内部に記憶されて示される。

【0080】

50

一実施形態において、I/Oインターフェース930は、ネットワークインターフェース940または他の周辺インターフェースを含めて、プロセッサ910と、システムメモリ920と、デバイス内の任意の周辺機器との間で、I/Oトラフィックを調整するように構成され得る。いくつかの実施形態において、I/Oインターフェース930は、1つの構成要素（例えば、システムメモリ920）からのデータ信号を、別の構成要素（例えば、プロセッサ910）によって使用されるために適切な形式に変換するために、任意の必要なプロトコル、タイミング、または他のデータ変換を実施することができる。いくつかの実施形態において、I/Oインターフェース930は、例えば、周辺構成要素相互接続（PCI）バス基準またはユニバーサルシリアルバス（USB）基準等、多種の周辺バスを通じて接続されたデバイスのサポートを含むことができる。いくつかの実施形態において、I/Oインターフェース930の機能は、例えば、ノースブリッジ及びサウスブリッジ等、2つ以上の別々の構成要素に分割され得る。また、いくつかの実施形態において、システムメモリ920へのインターフェース等、I/Oインターフェース930の機能のいくつかまたは全ては、プロセッサ910内に直接組み入れることができる。

10

20

30

40

50

【0081】

ネットワークインターフェース940は、データが、コンピュータシステム900と、例えば、図面に示される他のコンピュータシステムまたはデバイス等、ネットワーク（単数または複数）950に接続された他のデバイス960との間で交換されることを可能にするように構成され得る。様々な実施形態において、ネットワークインターフェース940は、例えば、イーサネット（登録商標）ワークのタイプ等、任意の適切な有線または無線汎用データネットワークを介して、通信をサポートすることができる。加えて、ネットワークインターフェース940は、ファイバチャンネルSAN等のストレージエリアネットワークを介して、あるいは任意の他の適切な種類のネットワーク及び/またはプロトコルを介して、アナログ音声ネットワークまたはデジタルファイバー通信ネットワーク等の電気通信網/テレフォニーネットワークを介して、通信をサポートすることができる。

【0082】

いくつかの実施形態において、システムメモリ920は、本明細書に記載される分散ロクマネージャの様々な実施形態を実装するために、図1～12の上記のプログラム命令及びデータを記憶するように構成されたコンピュータアクセス可能媒体の一実施形態であり得る。しかしながら、他の実施形態において、プログラム命令及び/またはデータは、受信、送信、または異なる種類のコンピュータアクセス可能媒体上に記憶され得る。一般に、コンピュータアクセス可能媒体は、磁気または光学式媒体等の非一時的記憶媒体、例えば、I/Oインターフェース930を介して、コンピュータシステム900に連結されたディスクまたはDVD/CDを含むことができる。非一時的コンピュータアクセス可能記憶媒体は、システムメモリ920または別のタイプのメモリとしてコンピュータシステム900のいくつかの実施形態に含むことができる、RAM（例えば、SDRAM、DDR SDRAM、RDRAM、SRAM等）、ROM等の任意の揮発性または不揮発性媒体も含むことができる。さらに、コンピュータアクセス可能媒体は、ネットワークインターフェース940を介して実装されてもよいように、ネットワーク及び/または無線リンク等の通信媒体を介して伝達される、電気、電磁気、またはデジタル信号等の伝送媒体または信号を含むことができる。

【0083】

前述の説明は、以下の付記の観点からさらに理解され得る。

1. ネットワーク上で相互に連結され、かつ仮想コンピューティングサービスを1つ以上のクライアントに集合的に提供する、複数のコンピューティングノードであって、コンピューティングノードの各々が、少なくとも1つのプロセッサとメモリとを備える、複数のコンピューティングノードと、

各々がコンピューティングノードのうち2つ以上のそれぞれのコンピューティングノード上で実行している、2つ以上の仮想コンピュータインスタンスであって、仮想コンピュータインスタンスの各々が、クライアントのために分散アプリケーションのアプリケー

ション構成要素を実装するように構成され、仮想コンピュータインスタンスが、ネットワークの少なくとも一部上で相互に通信するように構成される、2つ以上の仮想コンピュータインスタンスと、

各々が2つ以上のコンピューティングノードのそれぞれのコンピューティングノード上で実行している、2つ以上の構成要素を含む分散ロックマネージャであって、分散ロックマネージャ構成要素が、それぞれの共有リソースの1つ以上のロックの状態を共有するように構成され、1つ以上のロックの状態を共有するために、分散ロックマネージャ構成要素が、仮想コンピュータインスタンスが相互に通信するネットワークの少なくとも一部とは物理的または論理的に別個のネットワーク上で相互に通信するように構成される、分散ロックマネージャと、を備え、

10

アプリケーション構成要素のうちの一つが、分散アプリケーションによってアクセスされる共有リソースのロックのためのロック管理操作の実施を開始するために、アプリケーション構成要素を実装する仮想コンピュータインスタンスを実行しているコンピューティングノード上で実行している分散ロックマネージャへの呼び出しを起動するように構成される、システム。

2. 分散ロックマネージャ構成要素が相互に通信するネットワークが、サービスプロバイダコンピューティング環境の制御プレーンネットワークを備え、仮想コンピュータインスタンスが相互に通信するネットワークの少なくとも一部が、サービスプロバイダコンピューティング環境のデータプレーンネットワークを備える、付記1に記載のシステム。

3. 分散ロックマネージャ構成要素への呼び出しに 응답して、分散ロックマネージャ構成要素が、

20

ロック管理操作を実施し、かつ

分散ロックマネージャ構成要素を実行しているコンピューティングノード以外のコンピューティングノード上で実行している少なくとも一つの分散ロックマネージャ構成要素へロックの結果として生じた状態を通信するように構成される、付記1に記載のシステム。

4. 2つ以上のコンピューティングノードが、クライアントのための仮想プライベートネットワークを実装する、付記1に記載のシステム。

5. 1つ以上のコンピュータによって、

複数のコンピューティングノードのうちのととのコンピューティングノード上で実行している分散ロックマネージャの構成要素によって、共有リソースのロックのためのロック管理操作を実施する要求を受信することであって、該受信することが、所与のコンピューティングノード上で実行しているコンピュータインスタンスから要求を受信することを含み、コンピュータインスタンスが、コンピュータインスタンスに割り当てられたネットワークリソース能力を使用して共有リソースにアクセスする、受信することを実施することと、

30

要求されたロック管理操作を実施することであって、該実施することが、分散ロックマネージャ構成要素が、コンピュータインスタンスに割り当てられたネットワークリソース能力以外のネットワークリソース能力を使用して、ロックの状態情報を共有するために、複数のコンピューティングノードのうちのととのコンピューティングノード上で実行している別の分散ロックマネージャ構成要素と通信することを含む、実施することと、を含む、方法。

40

6. コンピュータインスタンスが、分散アプリケーションのアプリケーション構成要素を実装し、

方法が、アプリケーション構成要素が、コンピュータインスタンスに割り当てられたネットワークリソース能力を使用して、分散アプリケーションの少なくとも一つの他の構成要素と通信することをさらに含む、付記5に記載の方法。

7. 複数のコンピューティングノードが、分散ロックサービスを実装する、付記5に記載の方法。

8. 複数のコンピューティングノードが、1つ以上の仮想コンピューティングサービスを実装する、付記5に記載の方法。

50

9．ロック管理操作を該実施することが、ロックを共有リソースと関連付けることを含み、該通信することが、他の分散ロックマネージャ構成要素へ関連付けを通信することを含む、付記5に記載の方法。

10．ロック管理操作を該実施することが、ロックの値を変更することを含む、付記5に記載の方法。

11．ロックの値が、共有リソースの所有者を識別する、付記10に記載の方法。

12．分散ロックマネージャの構成要素が、ロックの変更された値を永続的データストアに書き込むことをさらに含む、付記10に記載の方法。

13．ロック値を該変更することが、ロック値をアトミックに変更することを含む、付記10に記載の方法。

14．コンピュータインスタンスに割り当てられたネットワークリソース能力が、ネットワーク接続リソースまたは入出力処理能力を含む、付記5に記載の方法。

15．該実施することが、所与のコンピューティングノードのロックの状態情報をキャッシュすることをさらに含む、付記5に記載の方法。

16．1つ以上のコンピュータ上で実行されると、1つ以上のコンピュータに、

複数のコンピューティングノードのうち所与のコンピューティングノード上で実行している分散ロックマネージャの構成要素によって、共有アクセスが制御されることになる実体と関連付けられたロックのためのロック管理操作を実施する要求を受信することであって、該受信することが、所与のコンピューティングノード上で実行しているリソースインスタンスから要求を受信することを含み、複数のコンピューティングノードが、分散ロックサービスを実装し、リソースインスタンスが、ロックと関連付けられた実体にアクセスし、要求が、分散ロックマネージャによってサポートされる1つ以上のロック管理操作を定義するアプリケーションプログラミングインターフェースに準拠する、受信すること、

要求されたロック管理操作を実施することであって、該実施することが、分散ロックマネージャ構成要素が、リソースインスタンスがロックと関連付けられたエンティティにアクセスするネットワーク接続以外のネットワーク接続を使用して、ロックの状態情報を共有するために、複数のコンピューティングノードのうち別のコンピューティングノード上で実行している分散ロックマネージャの別の構成要素と通信することを含む、実施すること、を実施させるプログラム命令を記憶する、非一時的コンピュータ可読記憶媒体。

17．要求が、ロックを作成する要求を含み、要求されたロック管理操作を該実施することが、ロックを作成すること、要求が受信されたリソースインスタンスへロックの識別子を返すこと、該通信することが、分散ロックマネージャの他の構成要素へロックの識別子を通信することを含む、付記16に記載の非一時的コンピュータ可読記憶媒体。

18．要求が、ロックをサブスクライブする要求、またはロックのプロパティの値を設定する要求を含む、付記16に記載の非一時的コンピュータ可読記憶媒体。

19．要求が、ロックを取得する要求、またはロックを解放する要求を含み、要求された操作を該実施することが、ロックのロック値を変更することを含み、該通信することが、変更されたロック値を分散ロックマネージャの他の構成要素に通信することを含む、付記16に記載の非一時的コンピュータ可読記憶媒体。

20．複数のコンピューティングノードが、ネットワーク上で相互に連結され、仮想コンピューティングサービスを1つ以上のクライアントに集合的に提供し、

1つ以上のコンピュータ上で実行されると、プログラム命令が、1つ以上のコンピュータに、

ロックのためのロック管理操作を実施する第2の要求を受信することをさらに実施させ、第2の要求が、ネットワーク上で相互に連結され、1つ以上のクライアントに仮想コンピューティングサービスを集合的に提供する、複数のコンピューティングノードのうちコンピューティングノード以外のコンピューティングサービスから受信され、要求が、分散ロックマネージャによってサポートされた1つ以上のロック管理操作を定義するアプリ

10

20

30

40

50

ケーションプログラミングインターフェースに準拠する、付記 16 に記載の非一時的コンピュータ可読記憶媒体。

21. ロックと関連付けられたエンティティが、仮想ネットワークインターフェースを備える、付記 16 に記載の非一時的コンピュータ可読記憶媒体。

【0084】

様々な実施形態は、コンピュータアクセス可能媒体に関する前述の説明に従って実装された命令及び/またはデータを受信すること、送信すること、または記憶することをさらに含むことができる。一般的に、コンピュータアクセス可能媒体として、磁気または光学式媒体（例えば、ディスクまたはDVD/CD-ROM）、RAM（例えば、SDRAM、DDR、RDRAM、SRAM等）、ROM等の揮発または不揮発媒体等の記憶媒体またはメモリ媒体、並びにネットワーク及び/またはワイヤレス連結等の通信媒体を介して伝達される、伝送媒体または、電気、電磁気、またはデジタル信号等の信号を挙げることができる。

10

【0085】

本明細書の図面及び説明において例示された様々な方法は、方法の例示の実施形態を表す。方法は、ソフトウェア、ハードウェア、またはこれらの組み合わせにおいて実装され得る。方法の順序は変更されてもよく、様々な要素が追加、順番を変更、組み合わせ、省略、修正等されてもよい。

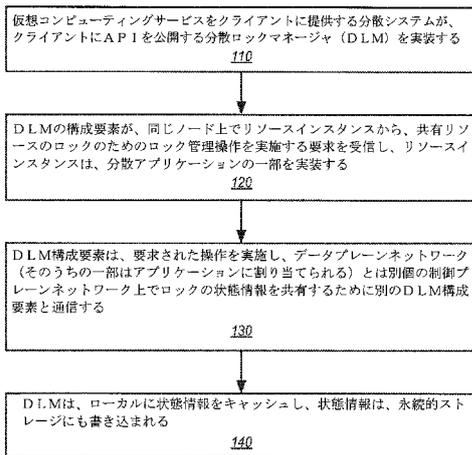
【0086】

本開示の利点を有する当業者には明らかであるように、様々な修正及び変更が行われ得る。全てのそのような修正及び変更を網羅すること、したがって、上記の説明は、制限的意味ではなく、例示として解釈されることが意図される。

20

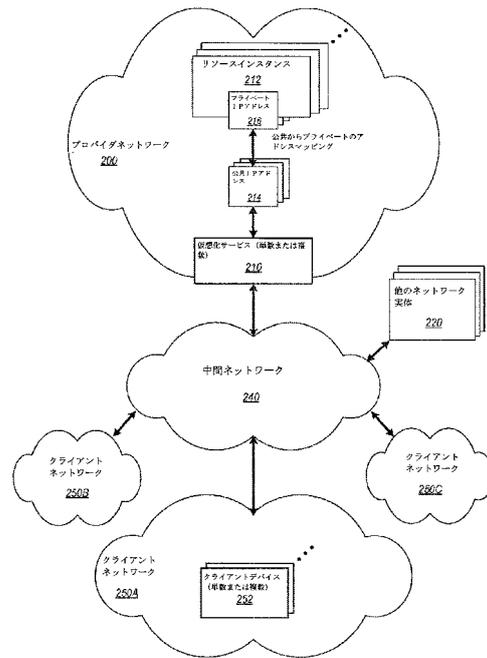
【図1】

図1



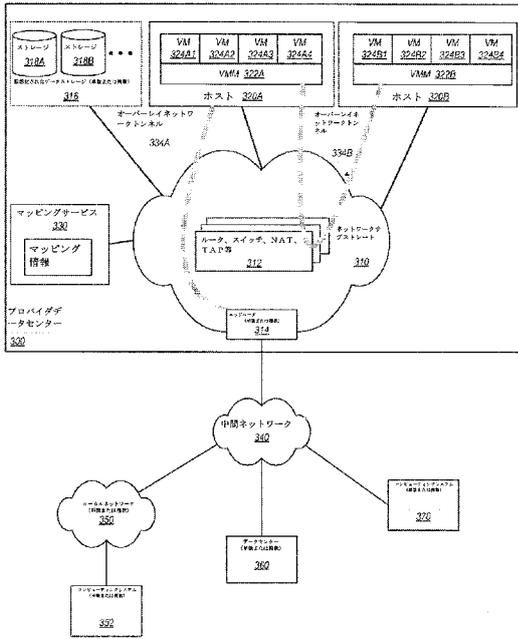
【図2】

図2



【図3】

図3



【図4】

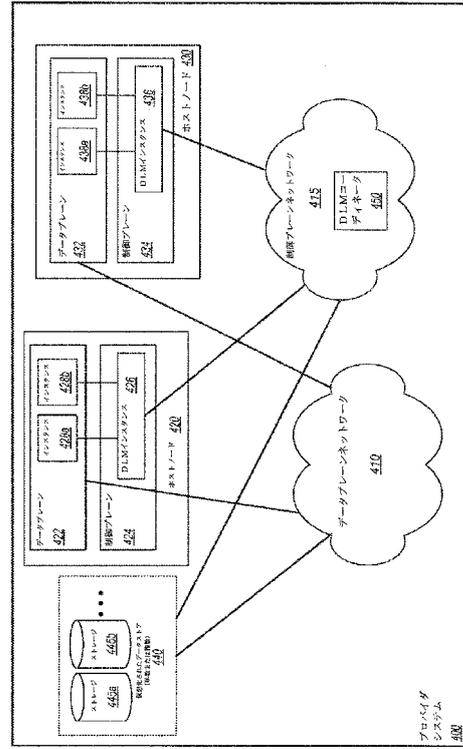
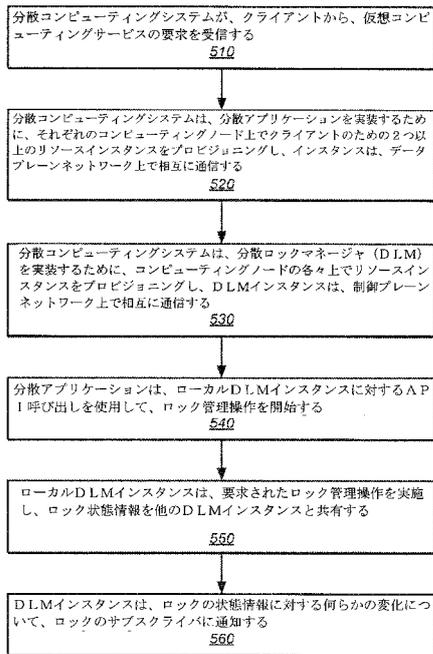


図4

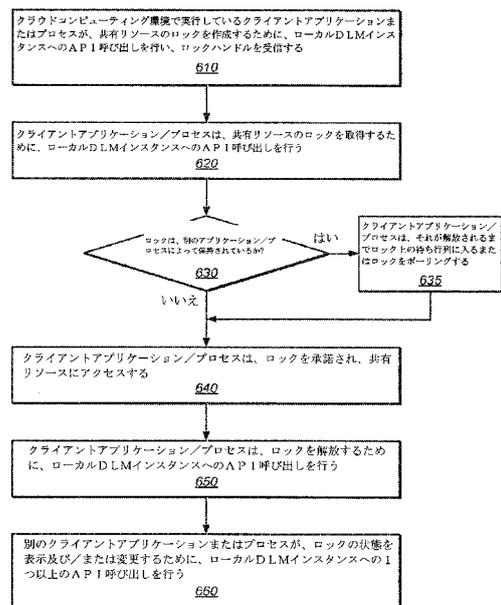
【図5】

図5



【図6】

図6



【 図 7 】

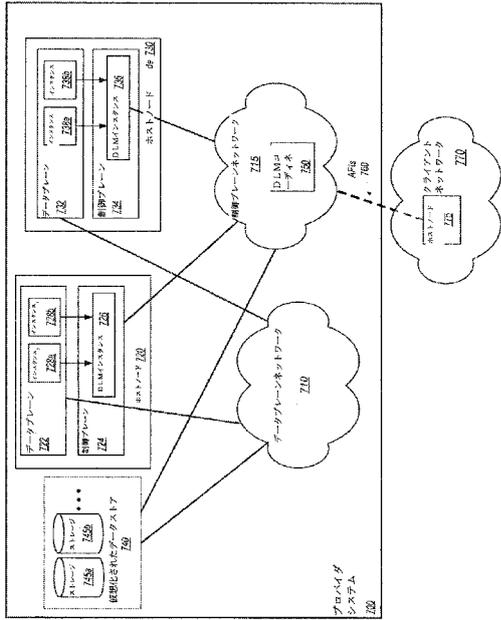
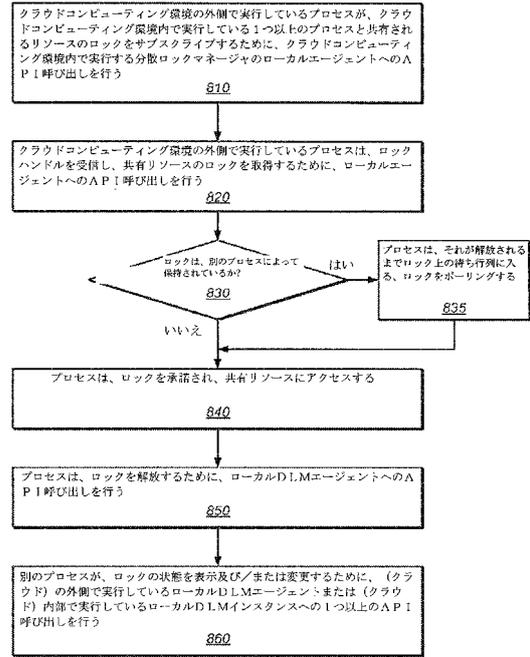


図 7

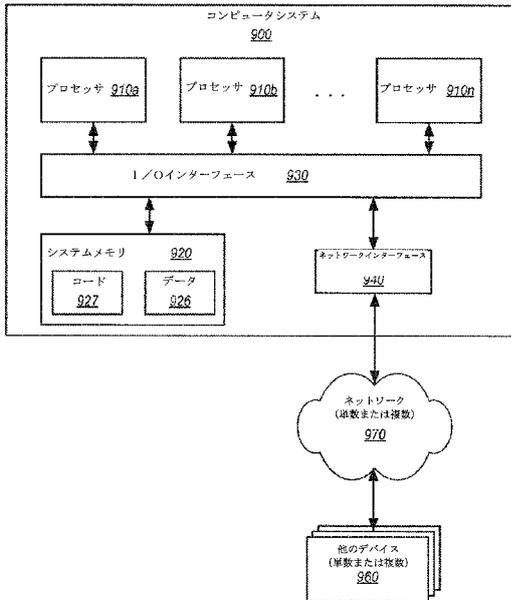
【 図 8 】

図 8



【 図 9 】

図 9



【 国際調査報告 】

INTERNATIONAL SEARCH REPORT		International application No. PCT/US 14/41724
A. CLASSIFICATION OF SUBJECT MATTER IPC(8) - G06F 7/00 (2014.01) CPC - G06F 17/30008 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) IPC(8): G06F 7/00 (2014.01) CPC: G06F 17/30008 Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched USFC: 707/704 or 707/713 or 707/887 or 707/609; IPC(8): G06F 7/00 (2014.01); CPC: G06F 17/30008 or G06F 9/526 or G06F 9/52 or G06F 9/466 or G06F 17/30067 Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) Patbase; Google (Web); distributed lock manager, lock, virtual, instance, resource, services, cloud, DLM, network, LAN, WAN		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	US 2012/0310881 A1 (Shadmon) 06 December 2012 (06.12.2012) entire document (especially para [0046]-[0049], [0053]-[0059], [0067]-[0075], [0081]-[0090])	1-15
Y	US 8,458,517 B1 (Vermeulen et al.) 04 June 2013 (04.06.2013) (col. 3, ln 50-65; col. 4, ln 32-33; col. 5, ln 52-55; col. 8, ln 59 to col. 9, ln 3; col. 14, ln 9-13; col. 12, ln 35-37)	1-15
Y	US 2008/0109807 A1 (Rosenbluth) 08 May 2008 (08.05.2008) para [0042]-[0043]	2
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/>		
* 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 "Z" document member of the same patent family		
Date of the actual completion of the international search 24 September 2014 (24.09.2014)		Date of mailing of the international search report <div style="text-align: center; font-size: 1.2em; font-weight: bold;">22 OCT 2014</div>
Name and mailing address of the ISA/US Mail Stop PCT, Attn: ISA/US, Commissioner for Patents P.O. Box 1450, Alexandria, Virginia 22313-1450 Facsimile No. 571-273-3201		Authorized officer: Lee W. Young PCT Helpdesk: 571-272-4300 PCT OSP: 571-272-7774

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, 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, LT, 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

【要約の続き】

用によって確立することができる。

【選択図】図4