

(19)日本国特許庁(JP)

(12)公表特許公報(A)

(11)公表番号

特表2025-504416
(P2025-504416A)

(43)公表日 令和7年2月12日(2025.2.12)

(51)国際特許分類		F I			
H 0 4 L	41/085 (2022.01)	H 0 4 L	41/085		
H 0 4 L	41/0895(2022.01)	H 0 4 L	41/0895		
G 0 6 F	9/50 (2006.01)	G 0 6 F	9/50	1 5 0 Z	

審査請求 未請求 予備審査請求 未請求 (全65頁)

(21)出願番号	特願2024-541878(P2024-541878)	(71)出願人	502303739 オラクル・インターナショナル・コーポ レイション
(86)(22)出願日	令和4年12月20日(2022.12.20)		アメリカ合衆国 9 4 0 6 5 カリフォル ニア州 レッドウッド ショアーズ, メー ル ストップ 5 オーピー7 オラクル パ ークウェイ 5 0 0
(85)翻訳文提出日	令和6年8月19日(2024.8.19)	(74)代理人	110001195 弁理士法人深見特許事務所
(86)国際出願番号	PCT/US2022/082073	(72)発明者	ブラール, ジャグウィンダー
(87)国際公開番号	WO2023/136964		アメリカ合衆国、9 4 0 6 5 カリフォ ルニア州、レッドウッド・シティー、オ ラクル・パークウェイ、5 0 0、オラクル ・インターナショナル・コーポレイシ ョン内
(87)国際公開日	令和5年7月20日(2023.7.20)		
(31)優先権主張番号	63/298,685		
(32)優先日	令和4年1月12日(2022.1.12)		
(33)優先権主張国・地域又は機関	米国(US)		
(31)優先権主張番号	18/050,392		
(32)優先日	令和4年10月27日(2022.10.27)		
(33)優先権主張国・地域又は機関	米国(US)		
(81)指定国・地域	AP(BW,CV,GH,GM,KE,LR,LS,MW,MZ 最終頁に続く		最終頁に続く

(54)【発明の名称】 グラフィック処理装置ワークロードの物理トポロジネットワーク局所性情報の公開

(57)【要約】

本明細書では、グラフィカル処理装置ベースのワークロードを実行するために、クラスタネットワークに含まれるホストマシンの局所性情報を利用する技術について論じられる。複数のホストマシンの各ホストマシンについて、ホストマシンの局所性情報がそこに記憶される。ホストマシンの局所性情報は、ホストマシンを含むラックを識別する。ワークロードの実行を要求する要求を受信したことに応答して、複数のホストマシンのうちの1つまたは複数のホストマシンが、ワークロードを実行するために使用可能であると識別される。1つまたは複数のホストマシンの各々について、ホストマシンの局所性情報が取得される。さらに、1つまたは複数のホストマシンのリンク情報が識別される。1つまたは複数のホストマシンの局所性情報とリンク情報は、要求に応じて提供される。

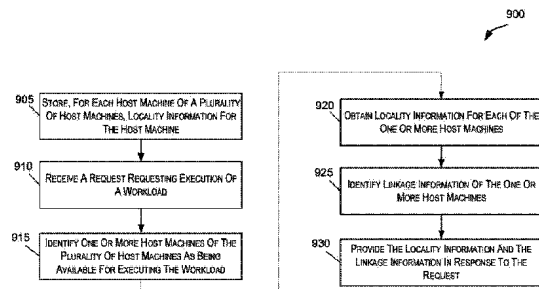


FIG. 9

【特許請求の範囲】**【請求項 1】**

複数のホストマシンの各ホストマシンについて、前記ホストマシンを含むラックを識別する前記ホストマシンの局所性情報を記憶することと、

ワークロードの実行を要求する要求を受信したことに応答して、

前記複数のホストマシンのうちの1つまたは複数のホストマシンを、前記ワークロードを実行するために使用可能であると識別することと、

前記1つまたは複数のホストマシンの各々について、前記ホストマシンの前記局所性情報を取得することと、

前記1つまたは複数のホストマシンのリンク情報を識別することと、

前記要求に応じて、前記1つまたは複数のホストマシンの前記局所性情報および前記リンク情報を提供することと、を含む、方法。

【請求項 2】

前記ホストマシンの前記局所性情報は、前記ホストマシンを含む前記ラックのラック識別子を含む情報を含む、請求項 1 に記載の方法。

【請求項 3】

前記ホストマシンの前記局所性情報は、前記ラックに関連付けられているトップオブラックスイッチの識別子を含む、請求項 1 に記載の方法。

【請求項 4】

前記1つまたは複数のホストマシンの前記リンク情報は、前記1つまたは複数のホストマシンによって形成される論理トポロジに対応する、請求項 1 に記載の方法。

【請求項 5】

前記1つまたは複数のホストマシンによって形成される前記論理トポロジはリングトポロジである、請求項 4 に記載の方法。

【請求項 6】

前記ホストマシンの前記局所性情報はインスタンスメタデータサービスから取得され、前記局所性情報は、前記ホストマシンに関連付けられるネットワーク仮想化装置を介して前記インスタンスメタデータサービスによって前記ホストマシンに記憶される、請求項 1 に記載の方法。

【請求項 7】

顧客が、前記1つまたは複数のホストマシンの前記局所性情報および前記リンク情報を受信することと、

前記顧客が、前記ワークロードを実行するための前記1つまたは複数のホストマシンのサブセットを選択することと、をさらに含み、前記選択は、前記ワークロードに関連付けられる1つまたは複数の制約に基づいて実行され、

前記1つまたは複数のホストマシンの前記サブセット上で前記ワークロードを実行することをさらに含む、請求項 1 に記載の方法。

【請求項 8】

1つまたは複数の制約は、ワークロードに関連付けられるレイテンシ閾値に対応する第1の制約と、前記ワークロードを実行する際の所望の可用性の程度に関連付けられる第2の制約とを含む、請求項 7 に記載の方法。

【請求項 9】

前記1つまたは複数のホストマシンの前記サブセットを前記選択することは、前記1つまたは複数のホストマシンの前記リンク情報にさらに基づく、請求項 8 に記載の方法。

【請求項 10】

コンピュータ実行可能命令を記憶する1つまたは複数のコンピュータ可読非一時的な媒体であって、前記コンピュータ実行可能命令は、1つまたは複数のプロセッサによって実行されると、

複数のホストマシンの各ホストマシンについて、前記ホストマシンを含むラックを識別する前記ホストマシンの局所性情報を記憶することと、

10

20

30

40

50

ワークロードの実行を要求する要求を受信したことに応答して、

前記複数のホストマシンのうちの1つまたは複数のホストマシンを、前記ワークロードを実行するために使用可能であると識別することと、

前記1つまたは複数のホストマシンの各々について、前記ホストマシンの前記局所性情報を取得することと、

前記1つまたは複数のホストマシンのリンク情報を識別することと、

前記要求に応じて、前記1つまたは複数のホストマシンの前記局所性情報および前記リンク情報を提供することと、を行わせる、1つまたは複数のコンピュータ可読非一時的な媒体。

【請求項11】

10

前記ホストマシンの前記局所性情報は、前記ホストマシンを含む前記ラックのラック識別子を含む情報を含む、請求項10に記載のコンピュータ実行可能命令を記憶する1つまたは複数のコンピュータ可読非一時的な媒体。

【請求項12】

前記ホストマシンの前記局所性情報は、前記ラックに関連付けられているトップオブラックスイッチの識別子を含む、請求項10に記載のコンピュータ実行可能命令を記憶する1つまたは複数のコンピュータ可読非一時的な媒体。

【請求項13】

前記1つまたは複数のホストマシンの前記リンク情報は、前記1つまたは複数のホストマシンによって形成される論理トポロジに対応する、請求項10に記載のコンピュータ実行可能命令を記憶する1つまたは複数のコンピュータ可読非一時的な媒体。

20

【請求項14】

前記1つまたは複数のホストマシンによって形成される前記論理トポロジは、リングトポロジである、請求項13に記載のコンピュータ実行可能命令を記憶する1つまたは複数のコンピュータ可読非一時的な媒体。

【請求項15】

前記ホストマシンの前記局所性情報はインスタンスメタデータサービスから取得され、前記局所性情報は、前記ホストマシンに関連付けられるネットワーク仮想化装置を介して前記インスタンスメタデータサービスによって前記ホストマシンに記憶される、請求項10に記載のコンピュータ実行可能命令を記憶する1つまたは複数のコンピュータ可読非一時的な媒体。

30

【請求項16】

請求項10に記載のコンピュータ実行可能命令を記憶する1つまたは複数のコンピュータ可読非一時的な媒体であって、1つまたは複数のプロセッサによって実行されると、

顧客が、前記1つまたは複数のホストマシンの前記局所性情報および前記リンク情報を受信することと、

前記ワークロードを実行するための前記1つまたは複数のホストマシンのサブセットを選択することと、を行わせる命令をさらに含み、前記選択は前記ワークロードに関連付けられる1つまたは複数の制約に基づいて実行され、

前記1つまたは複数のホストマシンの前記サブセット上で前記ワークロードを実行すること、を行わせる命令をさらに含む、1つまたは複数のコンピュータ可読非一時的な媒体

40

【請求項17】

1つまたは複数の制約は、ワークロードに関連付けられるレイテンシ閾値に対応する第1の制約と、前記ワークロードを実行する際の所望の可用性の程度に関連付けられる第2の制約とを含む、請求項16に記載のコンピュータ実行可能命令を記憶する1つまたは複数のコンピュータ可読非一時的な媒体。

【請求項18】

1つまたは複数のプロセッサと、

命令を含むメモリと、を含むコンピューティング装置であって、

50

前記命令は、前記1つまたは複数のプロセッサで実行されると、前記コンピューティング装置に少なくとも

複数のホストマシンの各ホストマシンについて、前記ホストマシンを含むラックを識別する前記ホストマシンの局所性情報を記憶することと、

ワークロードの実行を要求する要求を受信したことに応答して、

前記複数のホストマシンのうちの1つまたは複数のホストマシンを、前記ワークロードを実行するために使用可能であると識別することと、

前記1つまたは複数のホストマシンの各々について、前記ホストマシンの前記局所性情報を取得することと、

前記1つまたは複数のホストマシンのリンク情報を識別することと、

前記要求に応じて、前記1つまたは複数のホストマシンの前記局所性情報および前記リンク情報を提供することと、を行わせる、コンピューティング装置。

【請求項19】

前記ホストマシンの前記局所性情報は、前記ホストマシンを含む前記ラックのラック識別子を含む情報を含む、請求項18に記載のコンピューティング装置。

【請求項20】

前記ホストマシンの前記局所性情報は、前記ラックに関連付けられているトップオブラックスイッチの識別子を含む、請求項18に記載のコンピューティング装置。

【発明の詳細な説明】

【技術分野】

【0001】

関連出願の相互参照

本出願は、2022年1月12日に出願された米国仮出願第63/298,685号および2022年10月27日に出願された米国非仮出願第18/050,392号に対する優先権を主張する。前述の出願の全内容は、あらゆる目的のために、その全体が参照により本明細書に組み込まれる。

【0002】

分野

本開示は、クラウドインフラストラクチャのリモートダイレクトメモリアクセス(RDMA)対応クラスタネットワークに含まれるホストマシンの局所性情報を、グラフィカル処理装置ベースのワークロードの実行に利用することに関する。

【背景技術】

【0003】

背景

組織は、オンプレミスのハードウェアとソフトウェアの購入、更新、保守にかかるコストを削減するために、ビジネスアプリケーションとデータベースをクラウドに移行し続けている。高性能コンピュータアプリケーションは、特定の成果または結果を達成するために、使用可能なすべてのコンピューティング能力を一貫して消費する。このようなアプリケーションには、専用のネットワーク性能、高速記憶、高いコンピューティング能力、大量のメモリが必要であるが、これらは、今日のコモディティクラウドを構成する仮想化インフラストラクチャでは不足しているリソースである。

【0004】

クラウドインフラストラクチャサービスプロバイダは、これらのアプリケーションの要件に対応するために、より新しく高速なグラフィック処理装置(GPU)を提供している。GPUワークロードは通常、1つまたは複数のホストマシン上で実行される。通常、このようなワークロードでは、期待されるレベルのスループットを達成できない。この問題の要因の1つは、フローエントロピー(例えば、等コストマルチパス(ECMP)フローエントロピー)の欠如である。さらに、ホストマシン(つまり、ホスト)が局所ネットワーク近隣に他のホストが存在するかどうかに関係なくトラフィックを交換するという事実によって、問題はさらに悪化する。本明細書で説明する実施形態は、これらの問題および

10

20

30

40

50

その他の問題に対処する。

【発明の概要】

【0005】

概要

本開示は、一般的に、グラフィカル処理装置ベースのワークロードの実行において、クラウドインフラストラクチャのRDMA対応クラスタネットワークに含まれるホストマシンの局所性情報を利用することに関する。本明細書では、方法、システム、1つまたは複数のプロセッサによって実行可能なプログラム、コード、または命令を記憶する非一時的なコンピュータ可読記憶媒体などを含む、さまざまな実施形態が説明される。これらの例示的な実施形態は、開示を制限または定義するためではなく、理解を助けるための例を示すために言及されている。追加の実施形態は詳細な説明セクションで説明されており、ここでさらに詳しく説明されている。

10

【0006】

本開示の一実施形態は、複数のホストマシンの各ホストマシンについて、ホストマシンを含むラックを識別するホストマシンの局所性情報を記憶することと、ワークロードの実行を要求する要求を受信したことに応答して、複数のホストマシンのうちの1つまたは複数のホストマシンをワークロードの実行に使用可能であると識別することと、1つまたは複数のホストマシンの各々について、ホストマシンの局所性情報を取得することと、1つまたは複数のホストマシンのリンク情報を識別することと、要求に応じて1つまたは複数のホストマシンの局所性情報およびリンク情報を提供することと、を含む方法に関する。

20

【0007】

本開示の一態様は、1つまたは複数のデータプロセッサと、1つまたは複数のデータプロセッサ上で実行されるとコンピューティング装置に本明細書に開示される1つまたは複数の方法の一部または全部を実行させる命令を含む非一時的なコンピュータ可読記憶媒体とを含むコンピューティング装置を提供する。

【0008】

本開示の別の態様は、1つまたは複数のデータプロセッサに本明細書に開示される1つまたは複数の方法の一部または全部を実行させるように構成される命令を含む、非一時的なマシン可読記憶媒体に有形に具体化されるコンピュータプログラム製品を提供する。

【0009】

前述の内容、ならびに他の特徴および実施形態は、以下の明細書、特許請求の範囲、および添付の図面を参照することでより明らかになるであろう。

30

【0010】

本開示の特徴、実施形態、および利点は、添付の図面を参照して以下の詳細な説明を読むことによってよりよく理解される。

【図面の簡単な説明】

【0011】

【図1】特定の実施形態によるクラウドサービスプロバイダインフラストラクチャによってホストされる仮想またはオーバーレイクラウドネットワークを示す分散環境の高レベル図である。

40

【図2】特定の実施形態によるCSP I内の物理ネットワーク内の物理構成要素の簡略化されるアーキテクチャ図を示す図である。

【図3】特定の実施形態による、ホストマシンが複数のネットワーク仮想化装置(NVD)に接続されるCSP I内の配置例を示す図である。

【図4】特定の実施形態によるマルチテナントをサポートするためのI/O仮想化を提供するためのホストマシンとNVDとの間の接続を示す図である。

【図5】特定の実施形態によるCSP Iによって提供される物理ネットワークの簡略化されるブロック図を示す図である。

【図6】特定の実施形態による、CLOSネットワーク配置を組み込んだクラウドインフラストラクチャの簡略化されるブロック図を示す図である。

50

【図 7】 特定の実施形態による、クラウドインフラストラクチャに含まれるラックの例示的な構成を示す図である。

【図 8 A】 特定の実施形態による、局所性情報なしで構築される論理トポロジを示す図である。

【図 8 B】 特定の実施形態による、局所性情報を用いて構築される別の論理トポロジを示す図である。

【図 9】 特定の実施形態による、要求をプロビジョニングする際に実行されるステップを示す例示的なフローチャートを示す図である。

【図 10】 特定の実施形態による、顧客の要求を処理する際に実行されるステップを示す例示的なフローチャートを示す図である。

10

【図 11】 特定の実施形態による、顧客のワークロード要求をプロビジョニングする際に実行されるステップを示す例示的なフローチャートを示す図である。

【図 12】 少なくとも 1 つの実施形態による、クラウドインフラストラクチャをサービスシステムとして実装するための 1 つのパターンを示すブロック図である。

【図 13】 少なくとも 1 つの実施形態による、クラウドインフラストラクチャをサービスシステムとして実装するための別のパターンを示すブロック図である。

【図 14】 少なくとも 1 つの実施形態による、クラウドインフラストラクチャをサービスシステムとして実装するための別のパターンを示すブロック図である。

【図 15】 少なくとも 1 つの実施形態による、クラウドインフラストラクチャをサービスシステムとして実装するための別のパターンを示すブロック図である。

20

【図 16】 少なくとも 1 つの実施形態による例示的な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】

【0012】

詳細な説明

以下の説明では、説明の目的で、特定の実施形態を完全に理解できるように具体的な詳細が記載されている。しかし、これらの具体的な詳細がなくても、さまざまな実施形態を実践できることは明らかである。図および説明は制限的なものではない。本明細書で使用されている「例示的」という語は、「例、事例、または説明として役立つ」という意味である。本明細書で「例示的」として説明される実施形態または設計は、必ずしも他の実施形態または設計よりも好ましいまたは有利であると解釈されるものではない。

30

【0013】

クラウドインフラストラクチャのアーキテクチャ例

クラウドサービスという用語は、一般的に、クラウドサービスプロバイダ (CSP) が提供するシステムおよびインフラストラクチャ (クラウドインフラストラクチャ) を使用して、CSP によってユーザまたは顧客にオンデマンド (例えば、加入モデルを介して) で提供されるサービスを指すために使用される。通常、CSP のインフラストラクチャを構成するサーバとシステムは、顧客独自のオンプレミスサーバおよびシステムとは別である。したがって、顧客はサービスのために別途ハードウェアおよびソフトウェアリソースを購入することなく、CSP が提供するクラウドサービス自体を利用できる。クラウドサービスは、サービス提供に使用されるインフラストラクチャの調達に顧客が投資することなく、アプリケーションやコンピューティングリソースへの簡単かつスケーラブルなアクセスを加入顧客に提供するように設計されている。

40

【0014】

さまざまなタイプのクラウドサービスを提供するクラウドサービスプロバイダが複数存在する。クラウドサービスには、Software-as-a-Service (SaaS)、Platform-as-a-Service (PaaS)、Infrastructure-as-a-Service (IaaS) など、さまざまなタイプやモデルがある。

【0015】

顧客は、CSP が提供する 1 つまたは複数のクラウドサービスに加入することができる

50

。顧客は、個人、組織、企業など、あらゆるエンティティになり得る。C S P が提供するサービスに顧客が加入または登録すると、その顧客のテナントまたはアカウントが作成される。顧客は、次いでこのアカウントを介して、アカウントに関連付けられている1つまたは複数の加入済みクラウドリソースにアクセスできるようになる。

【0016】

前述のように、*infrastructure as a service (I a a S)* は、クラウドコンピューティングサービスの特定のタイプである。I a a S モデルでは、C S P は、顧客が独自のカスタマイズ可能なネットワークを構築し、顧客リソースを展開するために使用できるインフラストラクチャ（クラウドサービスプロバイダインフラストラクチャまたはC S P I と呼ばれる）を提供する。したがって、顧客のリソースとネットワークは、C S P によつて提供されるインフラストラクチャによって分散環境でホストされる。これは、顧客のリソースとネットワークが顧客が提供するインフラストラクチャによってホストされる従来のコンピューティングとは異なる。

10

【0017】

C S P I は、さまざまなホストマシン、メモリリソース、および物理ネットワークを形成するネットワークリソースを含む相互接続される高性能コンピューティングリソースを含み得、これはサブストレートネットワークまたはアンダーレイネットワークとも呼ばれる。C S P I のリソースは、1つまたは複数のデータセンターに分散し得、そのデータセンターは1つまたは複数の地理的リージョンにまたがって地理的に分散し得る。仮想化ソフトウェアは、これらの物理リソースによって実行され、仮想化される分散環境を提供し得る。仮想化により、物理ネットワーク上にオーバーレイネットワーク（ソフトウェアベースネットワーク、ソフトウェア定義ネットワーク、または仮想ネットワークとしても知られる）が作成される。C S P I 物理ネットワークは、物理ネットワーク上に1つまたは複数のオーバーレイネットワークまたは仮想ネットワークを作成するための基盤を提供する。仮想ネットワークまたはオーバーレイネットワークは、1つまたは複数の仮想クラウドネットワーク（V C N）を含むことができる。仮想ネットワークは、ソフトウェア仮想化技術（例えば、ハイパーバイザ、ネットワーク仮想化装置（N V D）によって実行される機能（例えば、スマートN I C）、トップオブブラック（T O R）スイッチ、N V Dによって実行される1つまたは複数の機能を実装するスマートT O R、およびその他のメカニズム）を使用して実装され、物理ネットワーク上で実行できるネットワーク抽象化の層が作成される。仮想ネットワークには、ピアツーピアネットワーク、I P ネットワークなど、さまざまな形式をとり得る。仮想ネットワークは通常、層3 I P ネットワークまたは層2 V L A N のいずれかである。この仮想またはオーバーレイネットワークの方法は、仮想またはオーバーレイ層3ネットワークと呼ばれることがよくある。仮想ネットワーク用に開発されるプロトコルの例には、I P - i n - I P（またはGeneric Routing Encapsulation（G R E））、Virtual Extensible LAN（V X L A N - I E T F R F C 7 3 4 8）、Virtual Private Networks（V P N）（例えば、MPLS Layer-3 Virtual Private Networks（R F C 4 3 6 4））、V M w a r e の N S X、G E N E V E（Generic Network Virtualization Encapsulation）などがある。

20

30

【0018】

I a a S の場合、C S P によって提供されるインフラストラクチャ（C S P I）は、パブリックネットワーク（例えば、インターネット）を介して仮想化されるコンピューティングリソースを提供するように構成することができる。I a a S モデルでは、クラウドコンピューティングサービスプロバイダがインフラストラクチャ構成要素（例えば、サーバ、記憶装置、ネットワークノード（例えば、ハードウェア）、展開ソフトウェア、プラットフォーム仮想化（例えば、ハイパーバイザ層）など）をホストできる。場合によっては、I a a S プロバイダが、これらのインフラストラクチャ構成要素に付随するさまざまなサービス（例えば、課金、監視、ログ記録、セキュリティ、負荷分散、クラスタリングなど）も提供することがある。したがって、これらのサービスはポリシー駆動型であり得るため、I a a S ユーザは、アプリケーションの可用性と性能を維持するために負荷分散を

40

50

駆動するポリシーを実装できる可能性がある。C S P Iは、顧客が可用性の高いホスト型分散環境で幅広いアプリケーションとサービスを構築および実行できるようにするインフラストラクチャと補完的なクラウドサービスのセットを提供する。C S P Iは、顧客のオンプレミスネットワークなどのさまざまなネットワークの場所から安全にアクセスできる柔軟な仮想ネットワークで、高性能なコンピューティングリソースと機能、および記憶容量を提供する。顧客がC S P Iが提供するIaaSサービスに加入または登録すると、その顧客用に作成されるテナントはC S P I内の安全で分離されるパーティションとなり、顧客はここで自分のクラウドリソースを作成、編成、管理できる。

【0019】

顧客は、C S P Iが提供するコンピューティング、メモリ、およびネットワークリソースを使用して独自の仮想ネットワークを構築できる。これらの仮想ネットワークには、コンピューティングインスタンスなどの1つまたは複数の顧客リソースまたはワークロードを展開できる。例えば、顧客はC S P Iが提供するリソースを使用して、仮想クラウドネットワーク（VCN）と呼ばれる1つまたは複数のカスタマイズ可能なプライベート仮想ネットワークを構築できる。顧客は、コンピューティングインスタンスなどの1つまたは複数の顧客リソースを顧客VCNに展開できる。コンピューティングインスタンスは、仮想マシン、ベアメタルインスタンスなどの形式をとることができる。したがって、C S P Iは、顧客が可用性の高い仮想ホスト環境で幅広いアプリケーションとサービスを構築および実行できるようにするインフラストラクチャと補完的なクラウドサービスのセットを提供する。顧客は、C S P Iによって提供される基礎となる物理リソースを管理または制御しないが、オペレーティングシステム、記憶、および展開されるアプリケーションを制御し、選択したネットワーク構成要素（例えば、ファイアウォール）を限定的に制御することもできる。

【0020】

C S P Iは、顧客およびネットワーク管理者がC S P Iリソースを使用してクラウドに展開されるリソースを構成、アクセス、および管理できるようにするコンソールを提供することができる。特定の実施形態では、コンソールは、C S P Iにアクセスして管理するために使用できるWebベースのユーザーインターフェースを提供する。いくつかの実装では、コンソールはC S P Iによって提供されるWebベースのアプリケーションである。

【0021】

C S P Iは、単一テナントまたはマルチテナントアーキテクチャをサポートすることができる。単一テナントアーキテクチャでは、ソフトウェア（例えば、アプリケーション、データベース）またはハードウェア構成要素（例えば、ホストマシン、サーバ）が単一の顧客またはテナントにサービスを提供する。マルチテナントアーキテクチャでは、ソフトウェアまたはハードウェア構成要素が複数の顧客またはテナントにサービスを提供する。したがって、マルチテナントアーキテクチャでは、C S P Iリソースは複数の顧客またはテナント間で共有される。マルチテナントの状況では、各テナントのデータが分離され、他のテナントから見えないようにするために、C S P I内で予防措置が講じられ、安全対策が講じられる。

【0022】

物理ネットワークにおいて、ネットワークエンドポイント（「エンドポイント」）とは、物理ネットワークに接続され、接続先のネットワークと往復して通信するコンピューティング装置またはシステムを指す。物理ネットワーク内のネットワークエンドポイントは、局所エリアネットワーク（LAN）、広域ネットワーク（WAN）、またはその他のタイプの物理ネットワークに接続され得る。物理ネットワーク内の従来のエンドポイントの例には、モデム、ハブ、ブリッジ、スイッチ、ルータ、その他のネットワーク装置、物理コンピュータ（またはホストマシン）などが含まれる。物理ネットワーク内の各物理装置は、装置との通信に使用できる固定のネットワークアドレスを有する。この固定ネットワークアドレスは、層2アドレス（例えば、MACアドレス）、固定層3アドレス（例えば、IPアドレス）などになり得る。仮想化環境または仮想ネットワークでは、エンドポイ

10

20

30

40

50

ントには、物理ネットワークの構成要素によってホストされる（例えば、物理ホストマシンによってホストされる）仮想マシンなどのさまざまな仮想エンドポイントを含めることができる。仮想ネットワーク内のこれらのエンドポイントは、オーバーレイ層2アドレス（例えば、オーバーレイMACアドレス）およびオーバーレイ層3アドレス（例えば、オーバーレイIPアドレス）などのオーバーレイアドレスによってアドレス指定される。ネットワークオーバーレイは、ネットワーク管理者がソフトウェア管理（例えば、仮想ネットワークの制御プレーンを実装するソフトウェア経路）を使用して、ネットワークエンドポイントに関連付けられているオーバーレイアドレスを移動できるようにすることで柔軟性を実現する。したがって、物理ネットワークとは異なり、仮想ネットワークでは、ネットワーク管理ソフトウェアを使用して、オーバーレイアドレス（例えば、オーバーレイIPアドレス）を1つのエンドポイントから別のエンドポイントに移動できる。仮想ネットワークは物理ネットワーク上に構築されるため、仮想ネットワーク内の構成要素間の通信には、仮想ネットワークと基盤となる物理ネットワークの両方が関与する。このような通信を容易にするために、C S P Iの構成要素は、仮想ネットワーク内のオーバーレイアドレスをサブストレートネットワーク内の実際の物理アドレスにマッピングするマッピング、およびその逆のマッピングを学習して記憶するように構成されている。これらのマッピングは、次いで、通信を容易にするために使用される。顧客トラフィックは、仮想ネットワークでのルーティングを容易にするためにカプセル化される。

10

【0023】

したがって、物理アドレス（例えば、物理IPアドレス）は物理ネットワーク内の構成要素に関連付けられ、オーバーレイアドレス（例えば、オーバーレイIPアドレス）は仮想ネットワーク内のエンティティに関連付けられる。物理IPアドレスとオーバーレイIPアドレスはどちらも実際のIPアドレスの一種である。これらは、仮想IPアドレスが複数の実際のIPアドレスにマップされる仮想IPアドレスとは異なる。仮想IPアドレスは、仮想IPアドレスと複数の実際のIPアドレス間の1対多のマッピングを提供する。

20

【0024】

クラウドインフラストラクチャまたはC S P Iは、世界中の1つまたは複数のリージョンにある1つまたは複数のデータセンターで物理的にホストされる。C S P Iには、物理ネットワークまたはサブストレートネットワーク内の構成要素と、物理ネットワーク構成要素上に構築される仮想ネットワーク内の仮想化構成要素（例えば、仮想ネットワーク、コンピューティングインスタンス、仮想マシン）が含まれ得る。特定の実施形態では、C S P Iは、レルム、リージョン、および可用性ドメインに編成され、ホストされる。リージョンは通常、1つまたは複数のデータセンターを含むローカライズされる地理的領域である。リージョンは一般に互いに独立しており、例えば、国や大陸をまたいで非常に離れた距離で隔てられている場合もある。例えば、第1のリージョンはオーストラリアにあり、別のリージョンは日本にあり、さらに別のリージョンはインドにあり得る、といった具合である。C S P Iリソースはリージョン間で分割されるため、各リージョンには独自の独立したC S P Iリソースのサブセットが存在する。各リージョンは、コンピューティングリソース（例えば、ベアメタルサーバ、仮想マシン、コンテナ、関連インフラストラクチャ）、記憶リソース（例えば、ブロックボリューム記憶、ファイル記憶、オブジェクト記憶、アーカイブ記憶）、ネットワークリソース（例えば、仮想クラウドネットワーク（VCN）、負荷分散リソース、オンプレミスネットワークへの接続）、データベースリソース、エッジネットワークリソース（例えば、DNS）、およびアクセス管理および監視リソースなど、コアインフラストラクチャサービスのセットとリソースを提供し得る。通常、各リージョンには、レルム内の他のリージョンに接続する複数のパスがある。

30

40

【0025】

一般に、アプリケーションは、最も頻繁に使用されるリージョン（すなわち、そのリージョンに関連付けられるインフラストラクチャ上に展開）に展開される。これは、近くのリソースを使用する方が遠くのリソースを使用するよりも高速であるためである。アプリ

50

ケーションは、大規模な気象システムや地震などのリージョン全体のイベントのリスクを軽減するための冗長性、法的管轄区域、課税ドメイン、その他のビジネス基準や社会基準などのさまざまな要件を満たすためなど、さまざまな理由で異なるリージョンに展開されることもある。

【0026】

リージョン内のデータセンターは、さらに編成され、可用性ドメイン（AD）に細分化されることができる。可用性ドメインは、リージョン内に配置される1つまたは複数のデータセンターに対応し得る。リージョンは1つまたは複数の可用性ドメインで構成できる。このような分散環境では、CSPリソースは、仮想クラウドネットワーク（VCN）などのリージョン固有、またはコンピューティングインスタンスなどの可用性ドメイン固有になる。

10

【0027】

リージョン内のADは互いに分離されており、耐障害性であり、同時に障害が発生する可能性が非常に低いように構成されている。これは、ADがネットワーク、物理ケーブル、ケーブルバス、ケーブルエントリポイントなどの重要なインフラストラクチャリソースを共有しないことで実現されるため、あるリージョン内の1つのADで障害が発生しても、同じリージョン内の他のADの可用性に影響が及ぶ可能性は低くなる。同じリージョン内のADは、低レイテンシ、高帯域幅のネットワークで相互に接続できるため、他のネットワーク（例えば、インターネット、顧客のオンプレミスネットワークなど）への高可用性接続を提供し、高可用性と災害復旧の両方のために複数のADに複製されるシステムを構築できる。クラウドサービスでは、高可用性を確保し、リソース障害から保護するために複数のADを使用する。IaaSプロバイダが提供するインフラストラクチャが拡大するにつれて、追加の容量とともに、より多くのリージョンとADが追加され得る。可用性ドメイン間のトラフィックは通常、暗号化される。

20

【0028】

特定の実施形態では、リージョンはレルムにグループ化される。レルムはリージョンの論理的な集合である。レルムは互いに分離されており、データを共有しない。同じレルム内のリージョンは相互に通信できるが、異なるレルム内のリージョンは通信できない。CSPで顧客のテナントまたはアカウントは単一のレルムに存在し、そのレルムに属する1つまたは複数のリージョンに分散できる。通常、顧客がIaaSサービスに加入すると、レルム内の顧客指定のリージョン（「ホーム」リージョンと呼ばれる）にその顧客のテナントまたはアカウントが作成される。顧客は、レルム内の1つまたは複数の他のリージョンにわたって顧客のテナントを拡張できる。顧客は、顧客のテナントが存在するレルム外のリージョンにアクセスすることはできない。

30

【0029】

IaaSプロバイダは、複数のレルムを提供することができ、各レルムは特定の顧客またはユーザのセットに対応する。例えば、商業顧客向けに商業レルムが提供され得る。別の例として、特定の国内の顧客に対して、その国専用のレルムが提供され得る。さらに別の例として、政府などに政府レルムが提供され得る。例えば、政府レルムは特定の政府向けに作成され、商業レルムよりもセキュリティレベルが高くなっている場合がある。例えば、Oracle Cloud Infrastructure（OCI）は現在、商用リージョン用のレルムと、政府クラウドリージョン用の2つのレルム（例えば、FedRAMP認定およびIL5認定）を提供している。

40

【0030】

特定の実施形態では、ADは1つまたは複数の障害ドメインに分割できる。障害ドメインは、反親和性を提供するためにAD内のインフラストラクチャリソースをグループ化したものである。障害ドメインを使用すると、コンピューティングインスタンスを分散して、インスタンスが単一のAD内の同じ物理ハードウェア上に存在しないようにすることができる。これは反親和性として知られている。障害ドメインとは、単一の障害点を共有するハードウェア構成要素（例えば、コンピュータ、スイッチ）のセットを指す。コンピュ

50

ーティングプールは、論理的に障害ドメインに分割される。このため、1つの障害ドメインに影響するハードウェア障害またはコンピューティングハードウェアのメンテナンスイベントは、他の障害ドメインのインスタスには影響しない。実施形態に応じて、各ADの障害ドメインの数は異なり得る。例えば、特定の実施形態では、各ADには3つの障害ドメインが含まれる。障害ドメインは、AD内の論理データセンターとして機能する。

【0031】

顧客がIaaSサービスに加入すると、CSPからのリソースが顧客向けにプロビジョニングされ、顧客のテナントに関連付けられる。顧客は、プロビジョニングされるこれらのリソースを使用してプライベートネットワークを構築し、これらのネットワーク上にリソースを展開できる。CSPによってクラウドでホストされる顧客ネットワークは、仮想クラウドネットワーク(VCN)と呼ばれる。顧客は、顧客に割り当てられるCSP Iリソースを使用して、1つまたは複数の仮想クラウドネットワーク(VCN)を設定できる。VCNは、仮想またはソフトウェア定義のプライベートネットワークである。顧客のVCNに展開される顧客リソースには、コンピューティングインスタス(例えば、仮想マシン、ベアメタルインスタス)およびその他のリソースが含まれ得る。これらのコンピューティングインスタスは、アプリケーション、負荷分散装置、データベースなどのさまざまな顧客ワークロードを表し得る。VCNに展開されるコンピューティングインスタスは、インターネットなどのパブリックネットワークを介してパブリックにアクセス可能なエンドポイント(「パブリックエンドポイント」)、同じVCNまたは他のVCN(例えば、顧客の他のVCN、顧客に属さないVCN)内の他のインスタス、顧客のオンプレミスデータセンターまたはネットワーク、サービスエンドポイント、およびその他のタイプのエンドポイントと通信できる。

【0032】

CSPは、CSP Iを使用してさまざまなサービスを提供することができる。場合によっては、CSP Iの顧客自身がサービスプロバイダのように行動し、CSP Iリソースを使用してサービスを提供できる。サービスプロバイダは、識別情報(例えば、IPアドレス、DNS名、ポート)によって特徴付けられるサービスエンドポイントを公開し得る。顧客のリソース(例えば、コンピューティングインスタス)は、特定のサービスに対してサービスによって公開されているサービスエンドポイントにアクセスすることで、特定のサービスを利用できる。これらのサービスエンドポイントは通常、インターネットなどのパブリック通信ネットワークを介してエンドポイントに関連付けられているパブリックIPアドレスを使用してユーザがパブリックにアクセスできるエンドポイントである。パブリックにアクセス可能なネットワークエンドポイントは、パブリックエンドポイントと呼ばれることもある。

【0033】

特定の実施形態では、サービスプロバイダは、サービスのエンドポイント(サービスエンドポイントと呼ばれることもある)を介してサービスを公開することができる。サービスの顧客は、次いで、このサービスエンドポイントを使用してサービスにアクセスできるようになる。特定の実装では、サービスに提供されるサービスエンドポイントに、そのサービスを利用しようとする複数の顧客がアクセスできる。他の実装では、専用のサービスエンドポイントが顧客に提供され得、その顧客のみがその専用サービスエンドポイントを使用してサービスにアクセスできるようになる。

【0034】

特定の実施形態では、VCNが作成されると、VCNに割り当てられるプライベートオーバーレイIPアドレスの範囲(例えば、10.0/16)であるプライベートオーバーレイクラスレスインタードメインルーティング(CIDR)アドレス空間に関連付けられる。VCNには、関連のサブネット、ルートテーブル、ゲートウェイが含まれる。VCNは単一のリージョン内に存在するが、リージョンの可用性ドメインの1つまたは複数またはすべてにまたがることができる。ゲートウェイは、VCN用に構成される仮想インターフェースであり、VCNとVCN外部の1つまたは複数のエンドポイントとの間のトラフ

10

20

30

40

50

ックの通信を可能にする。異なるタイプのエンドポイントとの通信を可能にするために、VCNに対して1つまたは複数の異なるタイプのゲートウェイを構成できる。

【0035】

VCNは、1つまたは複数のサブネットなどの1つまたは複数のサブネットワークに分割することができる。したがって、サブネットは、VCN内で作成できる構成単位またはサブディビジョンである。VCNには1つまたは複数のサブネットを含めることができる。VCN内の各サブネットは、そのVCN内の他のサブネットと重複せず、VCNのアドレス空間内のアドレス空間サブセットを表す、連続した範囲のオーバーレイIPアドレス（例えば、10.0.0.0/24および10.0.1.0/24）に関連付けられる。

【0036】

各コンピューティングインスタンスは、コンピューティングインスタンスがVCNのサブネットに参加できるようにする仮想ネットワークインターフェースカード（VNIC）に関連付けられている。VNICは、物理的なネットワークインターフェースカード（NIC）の論理表現である。一般に、VNICはエンティティ（例えば、コンピューティングインスタンス、サービス）と仮想ネットワーク間のインターフェースである。VNICはサブネット内に存在し、1つまたは複数の関連のIPアドレスと、関連のセキュリティルールまたはポリシーを持つ。VNICはスイッチ上の層2ポートに相当する。VNICは、コンピューティングインスタンスとVCN内のサブネットに接続される。コンピューティングインスタンスに関連付けられているVNICにより、コンピューティングインスタンスはVCNのサブネットの一部になることができ、コンピューティングインスタンスは、コンピューティングインスタンスと同じサブネット上のエンドポイント、VCN内の異なるサブネットのエンドポイント、またはVCN外のエンドポイントと通信（例えば、パケットの送受信）できるようになる。したがって、コンピューティングインスタンスに関連付けられているVNICによって、コンピューティングインスタンスがVCNの内外のエンドポイントに接続する方法が決定する。コンピューティングインスタンスが作成され、VCN内のサブネットに追加されると、コンピューティングインスタンスのVNICが作成され、そのコンピューティングインスタンスに関連付けられる。コンピューティングインスタンスのセットで構成されるサブネットの場合、サブネットにはコンピューティングインスタンスのセットに対応するVNICが含まれ、各VNICはコンピューティングインスタンスのセット内のコンピューティングインスタンスに接続される。

【0037】

各コンピューティングインスタンスには、コンピューティングインスタンスに関連付けられているVNICを介してプライベートオーバーレイIPアドレスが割り当てられる。このプライベートオーバーレイIPアドレスは、コンピューティングインスタンスの作成時にコンピューティングインスタンスに関連付けられているVNICに割り当てられ、コンピューティングインスタンスとの間のトラフィックのルーティングに使用される。特定のサブネット内のすべてのVNICは、同じルートテーブル、セキュリティリスト、およびDHCP任意選択肢を使用する。前述のように、VCN内の各サブネットは、そのVCN内の他のサブネットと重複せず、VCNのアドレス空間内のアドレス空間サブセットを表す、連続した範囲のオーバーレイIPアドレス（例えば、10.0.0.0/24および10.0.1.0/24）に関連付けられている。VCNの特定のサブネット上のVNICの場合、VNICに割り当てられるプライベートオーバーレイIPアドレスは、サブネットに割り当てられるオーバーレイIPアドレスの連続した範囲からのアドレスである。

【0038】

特定の実施形態では、コンピューティングインスタンスには、任意選択で、プライベートオーバーレイIPアドレスに加えて、例えばパブリックサブネット内にある場合は1つまたは複数のパブリックIPアドレスなどの追加のオーバーレイIPアドレスが割り当てられ得る。これらの複数のアドレスは、同じVNIC上、またはコンピューティングインスタンスに関連付けられている複数のVNIC上に割り当てられる。しかし、各インスタ

10

20

30

40

50

ンスには、インスタンスの起動時に作成され、インスタンスに割り当てられるオーバーレイプライベートIPアドレスに関連付けられているプライマリVNICがあり、このプライマリVNICは削除できない。セカンダリVNICと呼ばれる追加のVNICは、プライマリVNICと同じ可用性ドメイン内の既存のインスタンスに追加できる。すべてのVNICはインスタンスと同じ可用性ドメイン内にある。セカンダリVNICは、プライマリVNICと同じVCN内のサブネット、または同じVCN内または別のVCN内の別のサブネットにあり得る。

【0039】

コンピューティングインスタンスがパブリックサブネット内にある場合、任意選択でパブリックIPアドレスを割り当てることができる。サブネットは、作成時にパブリックサブネットまたはプライベートサブネットのいずれかとして指定できる。プライベートサブネットとは、サブネット内のリソース（例えば、コンピューティングインスタンス）と関連するVNICにパブリックオーバーレイIPアドレスを持たせることができないことを意味する。パブリックサブネットとは、サブネット内のリソースと関連するVNICにパブリックIPアドレスを持たせることができることを意味する。顧客は、サブネットを単一の可用性ドメイン内、またはリージョンもしくはレルム内の複数の可用性ドメイン間に存在するように指定できる。

10

【0040】

上述のように、VCNは1つまたは複数のサブネットに分割され得る。特定の実施形態では、VCN用に構成される仮想ルータ（VR）（VCN VRまたは単にVRと呼ばれる）によって、VCNのサブネット間の通信が可能になる。VCN内のサブネットの場合、VRはそのサブネットの論理ゲートウェイを表し、サブネット（つまり、そのサブネット上のコンピューティングインスタンス）がVCN内の他のサブネット上のエンドポイントおよびVCN外の他のエンドポイントと通信できるようにする。VCN VRは、VCN内のVNICとVCNに関連付けられる仮想ゲートウェイ（「ゲートウェイ」）間のトラフィックをルーティングするように構成される論理エンティティである。ゲートウェイについては、図1を参照して以下でさらに詳しく説明する。VCN VRは、層3/IP層の概念である。一実施形態では、VCNに対して1つのVCN VRがあり、VCN VRにはIPアドレスでアドレス指定されるポートが無制限に存在する可能性があり、VCNの各サブネットに対して1つのポートがある。このように、VCN VRは、VCN VRが接続されているVCN内のサブネットごとに異なるIPアドレスを持つ。VRは、VCN用に構成されるさまざまなゲートウェイにも接続される。特定の実施形態では、サブネットのオーバーレイIPアドレス範囲からの特定のオーバーレイIPアドレスが、そのサブネットのVCN VRのポート用に予約される。例えば、それぞれ10.0/16と10.1/16のアドレス範囲が関連付けられる2つのサブネットを持つVCNを考える。アドレス範囲が10.0/16であるVCN内の第1のサブネットの場合、この範囲のアドレスがそのサブネットのVCN VRのポート用に予約される。場合によっては、範囲の第1のIPアドレスがVCN VR用に予約され得る。例えば、オーバーレイIPアドレス範囲が10.0/16のサブネットの場合、IPアドレス10.0.0.1がそのサブネットのVCN VRのポート用に予約され得る。同じVCN内のアドレス範囲が10.1/16である第2のサブネットの場合、VCN VRにはIPアドレスが10.1.0.1であるその第2のサブネット用のポートが存在し得る。VCN VRには、VCN内のサブネットごとに異なるIPアドレスがある。

20

30

40

【0041】

いくつかの他の実施形態では、VCN内の各サブネットには、VRに関連付けられる予約済みまたはデフォルトのIPアドレスを使用してサブネットによってアドレス指定可能な独自の関連VRが存在し得る。予約済みまたはデフォルトのIPアドレスは、例えば、そのサブネットに関連付けられているIPアドレスの範囲の第1のIPアドレスであり得る。サブネット内のVNICは、このデフォルトまたは予約済みのIPアドレスを使用して、サブネットに関連付けられているVRと通信（例えば、パケットの送受信）できる。

50

このような実施形態では、VRはそのサブネットの入口/出口ポイントになる。VCN内のサブネットに関連付けられるVRは、VCN内の他のサブネットに関連付けられている他のVRと通信できる。VRは、VCNに関連付けられているゲートウェイとも通信できる。サブネットのVR機能は、サブネット内のVNICに対してVNIC機能を実行する1つまたは複数のNVD上で動作されているか、またはNVDによって実行される。

【0042】

VCNに対してルートテーブル、セキュリティルール、およびDHCP任意選択肢を構成できる。ルートテーブルはVCNの仮想ルートテーブルであり、ゲートウェイまたは特別に構成されたインスタンスを介して、VCN内のサブネットからVCN外部の宛先にトラフィックをルーティングするためのルールが含まれている。VCNのルートテーブルをカスタマイズして、VCNとの間でパケットを転送/ルーティングする方法を制御できる。DHCP任意選択肢とは、インスタンスの起動時に自動的に提供される構成情報を指す。

10

【0043】

VCNに対して構成されるセキュリティルールは、VCNのオーバーレイファイアウォールルールを表す。セキュリティルールには、入口ルールと出口ルールを含めることができ、VCN内のインスタンスに出入りできるトラフィックのタイプ（例えば、プロトコルおよびポートに基づく）を指定できる。顧客は、特定のルールがステートフルかステートレスかを選択できる。例えば、顧客は、ソースCIDR0.0.0.0/0、宛先TCPポート22でステートフル入口ルールを設定することで、どこからでもインスタンスのセットにSSHトラフィックを着信させることを許可できる。セキュリティルールは、ネットワークセキュリティグループまたはセキュリティリストを使用して実装できる。ネットワークセキュリティグループは、そのグループ内のリソースにのみ適用されるセキュリティルールのセットで構成される。一方、セキュリティリストには、セキュリティリストを使用するサブネット内のすべてのリソースに適用されるルールが含まれる。VCNには、デフォルトのセキュリティルールを含むデフォルトのセキュリティリストが提供され得る。VCNに構成されるDHCP任意選択肢は、インスタンスの起動時にVCN内のインスタンスに自動的に提供される構成情報を提供する。

20

【0044】

特定の実施形態では、VCNの構成情報は、VCN制御プレーンによって決定され、記憶される。VCNの構成情報には、例えば、VCNに関連付けられているアドレス範囲、VCN内のサブネットおよび関連情報、VCNに関連付けられている1つまたは複数のVR、VCN内のコンピューティングインスタンスおよび関連のVNIC、VCNに関連付けられるさまざまな仮想化ネットワーク機能（例えば、VNIC、VR、ゲートウェイ）を実行するNVD、VCNの状態情報、およびその他のVCN関連情報に関する情報が含まれ得る。特定の実施形態では、VCN配信サービスは、VCN制御プレーンによって記憶される構成情報またはその一部をNVDに公開する。配信される情報は、NVDによって記憶され、VCN内のコンピューティングインスタンスとの間でパケットを転送するために使用する情報（例えば、転送テーブル、ルーティングテーブル）を更新するために使用され得る。

30

40

【0045】

特定の実施形態では、VCNおよびサブネットの作成はVCN制御プレーン（CP）によって処理され、コンピューティングインスタンスの起動はコンピューティング制御プレーンによって処理される。コンピューティング制御プレーンは、コンピューティングインスタンスの物理リソースを割り当ての役割を担い、次いで、VCN制御プレーン呼び出してVNICを作成し、コンピューティングインスタンスに接続する。VCN CPは、パケット転送およびルーティング機能を実行するように構成されるVCNデータプレーンにVCNデータマッピングも送信する。特定の実施形態では、VCN CPは、VCNデータプレーンに更新を提供する役割を担う配信サービスを提供する。VCN制御プレーンの例は、図12、13、14、および15（参照1216、1316、1416、および

50

1516を参照)にも示されており、以下で説明する。

【0046】

顧客は、C S P Iによってホストされるリソースを使用して1つまたは複数のV C Nを作成することができる。顧客のV C Nに展開されるコンピューティングインスタンスは、さまざまなエンドポイントと通信し得る。これらのエンドポイントには、C S P IによってホストされるエンドポイントとC S P I外部のエンドポイントが含まれ得る。

【0047】

C S P Iを使用してクラウドベースのサービスを実装するためのさまざまな異なるアーキテクチャが図1、2、3、4、5、12、13、14、および15に示されており、以下で説明する。図1は、特定の実施形態によるC S P Iによってホストされるオーバーレイまたは顧客V C Nを示す分散環境100の高レベル図である。図1に示されている分散環境には、オーバーレイネットワーク内の複数の構成要素が含まれている。図1に示す分散環境100は単なる例であり、請求される実施形態の範囲を不当に制限することを意図するものではない。多くのパリエーション、代替、変更が可能である。例えば、いくつかの実装では、図1に示されている分散環境には、図1に示されているものよりも多いた

10

【0048】

図1に示す例に示すように、分散環境100は、顧客が加入して仮想クラウドネットワーク(V C N)を構築するために使用できるサービスおよびリソースを提供するC S P I 101を含む。特定の実施形態では、C S P I 101は加入顧客にI a a Sサービスを提供する。C S P I 101内のデータセンターは、1つまたは複数のリージョンに編成され得る。一例としてのリージョン「Region US」102を図1に示す。顧客はリージョン102に顧客V C N 104を構成した。顧客はV C N 104上にさまざまなコンピューティングインスタンスを展開できる。コンピューティングインスタンスには、仮想マシンまたはベアメタルインスタンスが含まれ得る。インスタンスの例としては、アプリケーション、データベース、負荷分散装置などが含まれる。

20

【0049】

図1に示す実施形態では、顧客V C N 104は2つのサブネット、すなわち「サブネット1」と「サブネット2」で構成され、各サブネットは独自のC I D R I Pアドレス範囲を有する。図1では、サブネット1のオーバーレイI Pアドレス範囲は10.0/16、サブネット2のアドレス範囲は10.1/16である。V C N仮想ルータ105は、V C N 104のサブネット間およびV C N外部の他のエンドポイントとの通信を可能にするV C Nの論理ゲートウェイを表す。V C N V R 105は、V C N 104内のV N I CとV C N 104に関連付けられるゲートウェイ間のトラフィックをルーティングするように構成されている。V C N V R 105は、V C N 104の各サブネットにポートを提供する。例えば、V R 105は、サブネット1にI Pアドレス10.0.0.1のポートを提供し、サブネット2にI Pアドレス10.1.0.1のポートを提供できる。

30

【0050】

各サブネットには複数のコンピューティングインスタンスが展開され得、コンピューティングインスタンスは仮想マシンインスタンスおよび/またはベアメタルインスタンスであり得る。サブネット内のコンピューティングインスタンスは、C S P I 101内の1つまたは複数のホストマシンによってホストされ得る。コンピューティングインスタンスは、コンピューティングインスタンスに関連付けられているV N I Cを介してサブネットに参加する。例えば、図1に示すように、コンピューティングインスタンスC 1は、コンピューティングインスタンスに関連付けられているV N I Cを介してサブネット1の一部になる。同様に、コンピューティングインスタンスC 2は、C 2に関連付けられているV N I Cを介してサブネット1の一部になる。同様に、仮想マシンインスタンスまたはベアメタルインスタンスである可能性のある複数のコンピューティングインスタンスがサブネット1の一部になり得る。各コンピューティングインスタンスには、関連のV N I Cを介し

40

50

て、プライベートオーバーレイIPアドレスとMACアドレスが割り当てられる。例えば、図1では、コンピューティングインスタンスC1のオーバーレイIPアドレスは10.0.0.2、MACアドレスはM1であるが、コンピューティングインスタンスC2のプライベートオーバーレイIPアドレスは10.0.0.3、MACアドレスはM2である。コンピューティングインスタンスC1およびC2を含むサブネット1の各コンピューティングインスタンスには、IPアドレス10.0.0.1を使用してVCN VR105へのデフォルトルートがある。これは、サブネット1のVCN VR105のポートのIPアドレスである。

【0051】

サブネット2には、仮想マシンインスタンスおよび/またはベアメタルインスタンスなど、複数のコンピューティングインスタンスを展開できる。例えば、図1に示すように、コンピューティングインスタンスD1とD2は、それぞれのコンピューティングインスタンスに関連付けられているVNICを介してサブネット2の一部になる。図1に示す実施形態では、コンピューティングインスタンスD1はオーバーレイIPアドレス10.1.0.2とMACアドレスMM1を持ち、コンピューティングインスタンスD2はプライベートオーバーレイIPアドレス10.1.0.3とMACアドレスMM2を持つ。サブネット2内の各コンピューティングインスタンス(コンピューティングインスタンスD1およびD2を含む)には、IPアドレス10.1.0.1を使用するVCN VR105へのデフォルトルートがある。これは、サブネット2のVCN VR105のポートのIPアドレスである。

【0052】

VCN A104には、1つまたは複数の負荷分散装置も含まれ得る。例えば、サブネットに負荷分散装置が提供され、サブネット上の複数のコンピューティングインスタンス間でトラフィックを負荷分散するように構成され得る。VCN内のサブネット間でトラフィックの負荷を分散するために、負荷分散装置が提供され得る。

【0053】

VCN 104上に展開される特定のコンピューティングインスタンスは、さまざまな異なるエンドポイントと通信できる。これらのエンドポイントには、CSP I200によってホストされるエンドポイントとCSP I200外部のエンドポイントが含まれ得る。CSP I101によってホストされるエンドポイントには、特定のコンピューティングインスタンスと同じサブネット上のエンドポイント(例えば、サブネット1内の2つのコンピューティングインスタンス間の通信)、異なるサブネット上にあるが同じVCN内のエンドポイント(例えば、サブネット1内のコンピューティングインスタンスとサブネット2内のコンピューティングインスタンス間の通信)、同じリージョン内の異なるVCN内のエンドポイント(例えば、サブネット1内のコンピューティングインスタンスと同じリージョン106または110内のVCN内のエンドポイント間の通信、サブネット1内のコンピューティングインスタンスと同じリージョン内のサービスネットワーク110内のエンドポイント間の通信)、または異なるリージョン内のVCN内のエンドポイント(例えば、サブネット1内のコンピューティングインスタンスと異なるリージョン108内のVCN内のエンドポイント間の通信)が含まれ得る。CSP I101によってホストされているサブネット内のコンピューティングインスタンスは、CSP I101によってホストされていない(つまり、CSP I101の外部にある)エンドポイントとも通信できる。これらの外部エンドポイントには、顧客のオンプレミスネットワーク116内のエンドポイント、他のリモートクラウドホストネットワーク118内のエンドポイント、インターネットなどのパブリックネットワーク経由でアクセス可能なパブリックエンドポイント114、およびその他のエンドポイントが含まれる。

【0054】

同じサブネット上のコンピューティングインスタンス間の通信は、ソースコンピューティングインスタンスと宛先コンピューティングインスタンスに関連付けられているVNICを使用して容易になる。例えば、サブネット1のコンピューティングインスタンスC1

10

20

30

40

50

は、サブネット1のコンピューティングインスタンスC2にパケットを送信したい場合がある。ソースコンピューティングインスタンスから発信され、宛先が同じサブネット内の別のコンピューティングインスタンスであるパケットの場合、パケットは最初にソースコンピューティングインスタンスに関連付けられているVNICによって処理される。ソースコンピューティングインスタンスに関連付けられているVNICによって実行される処理には、パケットヘッダーからパケットの宛先情報を決定すること、ソースコンピューティングインスタンスに関連付けられているVNICに構成されているポリシー（例えば、セキュリティリスト）を識別すること、パケットの次のホップを決定すること、必要に応じてパケットのカプセル化/カプセル化解除機能を実行すること、そして、パケットを次のホップに転送/ルーティングして、パケットを目的の宛先に容易に通信するようにすることが含まれ得る。宛先コンピューティングインスタンスがソースコンピューティングインスタンスと同じサブネットにある場合、ソースコンピューティングインスタンスに関連付けられているVNICは、宛先コンピューティングインスタンスに関連付けられているVNICを識別し、パケットをそのVNICに転送して処理するように構成される。次に、宛先コンピューティングインスタンスに関連付けられているVNICが実行され、パケットが宛先コンピューティングインスタンスに転送される。

10

【0055】

サブネット内のコンピューティングインスタンスから同じVCN内の別のサブネット内のエンドポイントにパケットを通信する場合、その通信は、ソースおよび宛先のコンピューティングインスタンスとVCN VRに関連付けられているVNICによって容易になる。例えば、図1のサブネット1内のコンピューティングインスタンスC1がサブネット2内のコンピューティングインスタンスD1にパケットを送信したい場合、パケットは、まずコンピューティングインスタンスC1に関連付けられているVNICによって処理される。コンピューティングインスタンスC1に関連付けられているVNICは、デフォルトルートまたはVCN VRのポート10.0.0.1を使用してパケットをVCN VR 105にルーティングするように構成されている。VCN VR 105は、ポート10.1.0.1を使用してパケットをサブネット2にルーティングするように構成されている。次に、パケットはD1に関連付けられているVNICによって受信および処理され、VNICはパケットをコンピューティングインスタンスD1に転送する。

20

【0056】

VCN 104内のコンピューティングインスタンスからVCN 104外のエンドポイントにパケットを通信する場合、その通信は、ソースコンピューティングインスタンスに関連付けられているVNIC、VCN VR 105、およびVCN 104に関連付けられるゲートウェイによって容易になる。VCN 104には、1つまたは複数のタイプのゲートウェイが関連付けられ得る。ゲートウェイは、VCNと別のエンドポイント間のインターフェースであり、別のエンドポイントはVCNの外部にある。ゲートウェイは層3/IP層の概念であり、VCNがVCN外部のエンドポイントと通信できるようにする。したがって、ゲートウェイは、VCNと他のVCNまたはネットワーク間のトラフィックフローを容易にする。さまざまなタイプのエンドポイントとのさまざまなタイプの通信を容易にするために、VCNにさまざまなタイプのゲートウェイを構成できる。ゲートウェイに応じて、通信はパブリックネットワーク（例えば、インターネット）経由またはプライベートネットワーク経由で行われ得る。これらの通信にはさまざまな通信プロトコルが使用され得る。

30

40

【0057】

例えば、コンピューティングインスタンスC1は、VCN 104の外部のエンドポイントと通信したい場合がある。パケットは、まずソースコンピューティングインスタンスC1に関連付けられているVNICによって処理され得る。VNIC処理により、パケットの宛先がC1のサブネット1の外部にあることが決定される。C1に関連付けられているVNICは、パケットをVCN 104のVCN VR 105に転送し得る。次に、VCN VR 105はパケットを処理し、処理の一環として、パケットの宛先に基づいて、VCN

50

104に関連付けられる特定のゲートウェイをパケットの次のホップとして決定する。その後、VCN VR105は、パケットを特定の識別されるゲートウェイに転送し得る。例えば、宛先が顧客のオンプレミスネットワーク内のエンドポイントである場合、パケットはVCN VR105によって、VCN104用に構成される動的ルーティングゲートウェイ(DRG)ゲートウェイ122に転送され得る。その後、パケットはゲートウェイから次のホップに転送され、最終的な目的の宛先へのパケットの通信が容易になり得る。

【0058】

VCNには、さまざまなタイプのゲートウェイを構成できる。VCN用に構成できるゲートウェイの例を図1に示し、以下に説明する。VCNに関連付けられるゲートウェイの例も、図12、13、14、および15に示されており(例えば、参照番号1234、1236、1238、1334、1336、1338、1434、1436、1438、1534、1536、および1538で参照されるゲートウェイ)、以下で説明する。図1に示す実施形態に示すように、動的ルーティングゲートウェイ(DRG)122は、顧客VCN104に追加されるか、または顧客VCN104に関連付けられ得、顧客VCN104と別のエンドポイントとの間のプライベートネットワークトラフィック通信のパスを提供する。ここで、別のエンドポイントは、顧客のオンプレミスネットワーク116、CSP101の別のリージョン内のVCN108、またはCSP101によってホストされていない他のリモートクラウドネットワーク118であり得る。顧客オンプレミスネットワーク116は、顧客のリソースを使用して構築される顧客ネットワークまたは顧客データセンターであり得る。顧客のオンプレミスネットワーク116へのアクセスは、通常、非常に制限されている。顧客オンプレミスネットワーク116と、CSP101によってクラウドに展開またはホストされている1つまたは複数のVCN104の両方を持つ顧客の場合、顧客はオンプレミスネットワーク116とクラウドベースのVCN104が相互に通信できるようにしたい場合がある。これにより、顧客は、CSP101によってホストされる顧客のVCN104とオンプレミスネットワーク116を含む拡張ハイブリッド環境を構築できるようになる。DRG122はこの通信を可能にする。このような通信を可能にするために、通信チャンネル124が設定され、チャンネルの1つのエンドポイントは顧客のオンプレミスネットワーク116にあり、もう1つのエンドポイントはCSP101にあり、顧客のVCN104に接続されている。通信チャンネル124は、インターネットなどのパブリック通信ネットワークまたはプライベート通信ネットワークを介することができる。インターネットなどのパブリック通信ネットワークを介したIPsecVPN技術や、パブリックネットワークの代わりにプライベートネットワークを使用するOracleのFastConnect技術など、さまざまな通信プロトコルを使用できる。通信チャンネル124の1つのエンドポイントを形成する顧客オンプレミスネットワーク116内の装置または機器は、図1に示すCPE126などの顧客構内機器(CPE)と呼ばれる。CSP101側では、エンドポイントはDRG122を実行するホストマシンであり得る。特定の実施形態では、リモートピアリング接続(RPC)をDRGに追加して、顧客が1つのVCNを別のリージョンの別のVCNとピアリングできるようにすることができる。このようなRPCを使用すると、顧客VCN104はDRG122を使用して別のリージョンのVCN108に接続できる。DRG122は、Microsoft Azureクラウド、Amazon AWSクラウドなど、CSP101でホストされていない他のリモートクラウドネットワーク118との通信にも使用できる。

【0059】

図1に示すように、顧客VCN104に対してインターネットゲートウェイ(IGW)120が構成され、VCN104上のコンピューティングインスタンスがインターネットなどのパブリックネットワークを介してアクセス可能なパブリックエンドポイント114と通信できるようになる。IGW120は、VCNをインターネットなどのパブリックネットワークに接続するゲートウェイである。IGW120は、VCN104などのVCN内のパブリックサブネット(パブリックサブネット内のリソースがパブリックオーバーレイIPアドレスを持つ)が、インターネットなどのパブリックネットワーク114上の

パブリックエンドポイント 112 に直接アクセスできるようにする。I G W 1 2 0 を使用すると、V C N 1 0 4 内のサブネットまたはインターネットから接続を開始できる。

【 0 0 6 0 】

ネットワークアドレス変換 (N A T) ゲートウェイ 1 2 8 は、顧客の V C N 1 0 4 用に構成することができ、専用のパブリックオーバーレイ I P アドレスを持たない顧客の V C N 内のクラウドリソースがインターネットにアクセスできるようにし、これらのリソースを直接の着信インターネット接続 (例えば、L 4 - L 7 接続) に公開することなく、これを実現する。これにより、V C N 1 0 4 のプライベートサブネット 1 など、V C N 内のプライベートサブネットがインターネット上のパブリックエンドポイントにプライベートアクセスできるようになる。N A T ゲートウェイでは、プライベートサブネットからパブリックインターネットへの接続のみを開始でき、インターネットからプライベートサブネットへの接続は開始できない。

10

【 0 0 6 1 】

特定の実施形態では、サービスゲートウェイ (S G W) 1 2 6 を顧客 V C N 1 0 4 用に構成することができ、V C N 1 0 4 とサービスネットワーク 1 1 0 内のサポートされているサービスエンドポイントとの間のプライベートネットワークトラフィックのパスを提供する。特定の実施形態では、サービスネットワーク 1 1 0 は C S P によって提供され、さまざまなサービスを提供し得る。このようなサービスネットワークの例としては、顧客が利用できるさまざまなサービスを提供する Oracle のサービスネットワークがある。例えば、顧客 V C N 1 0 4 のプライベートサブネット内のコンピューティングインスタンス (例えば、データベースシステム) は、パブリック I P アドレスやインターネットへのアクセスを必要とせずに、サービスエンドポイント (例えば、オブジェクト記憶装置) にデータをバックアップできる。特定の実施形態では、V C N は 1 つの S G W のみを持つことができ、接続はサービスネットワーク 1 1 0 からではなく、V C N 内のサブネットからのみ開始できる。V C N が別の V C N とピアリングされている場合、通常、他の V C N のリソースは S G W にアクセスできない。FastConnect または V P N 接続を使用して V C N に接続されているオンプレミスネットワーク内のリソースも、その V C N 用に構成されるサービスゲートウェイを使用できる。

20

【 0 0 6 2 】

特定の実装では、S G W 1 2 6 は、サービスクラスレスインタードメインルーティング (C I D R) ラベルの概念を使用する。これは、対象のサービスまたはサービスグループのすべてのリージョンのパブリック I P アドレス範囲を表す文字列である。顧客は、サービスへのトラフィックを制御するために S G W および関連ルートルールを構成するときに、サービス C I D R ラベルを使用する。顧客は、将来サービスのパブリック I P アドレスが変更される場合でもセキュリティルールを調整する必要なく、セキュリティルールを構成するときに任意選択でこれを利用できる。

30

【 0 0 6 3 】

局所ピアリングゲートウェイ (L P G) 1 3 2 は、顧客 V C N 1 0 4 に追加できるゲートウェイであり、V C N 1 0 4 が同じリージョン内の別の V C N とピアリングできるようにする。ピアリングとは、トラフィックがインターネットなどのパブリックネットワークを通過したり、顧客のオンプレミスネットワーク 1 1 6 を経由してトラフィックをルーティングしたりすることなく、V C N がプライベート I P アドレスを使用して通信することを意味する。好ましい実施形態では、V C N は、確立するピアリングごとに個別の L P G を持つ。局所ピアリングまたは V C N ピアリングは、異なるアプリケーションまたはインフラストラクチャ管理機能間のネットワーク接続を確立するために使用される一般的な方法である。

40

【 0 0 6 4 】

サービスネットワーク 1 1 0 内のサービスのプロバイダなどのサービスプロバイダは、異なるアクセスモデルを使用してサービスへのアクセスを提供することができる。パブリックアクセスモデルに従って、サービスは、インターネットなどのパブリックネットワー

50

クを介して顧客VCN内のコンピューティングインスタンスによってパブリックにアクセス可能なパブリックエンドポイントとして公開される場合もあれば、SGW126を介してプライベートにアクセス可能な場合もある。特定のプライベートアクセスモデルに従って、サービスは顧客のVCN内のプライベートサブネット内のプライベートIPエンドポイントとしてアクセス可能になる。これはプライベートエンドポイント(PE)アクセスと呼ばれ、サービスプロバイダが顧客のプライベートネットワーク内のインスタンスとしてサービスを公開できるようにする。プライベートエンドポイントリソースは、顧客のVCN内のサービスを表す。各PEは、顧客のVCN内で顧客が選択したサブネット内のVNIC(1つまたは複数のプライベートIPを持つPE-VNICと呼ばれる)として現れる。したがって、PEは、VNICを使用してプライベート顧客VCNサブネット内でサービスを提供する方法を提供する。エンドポイントはVNICとして公開されるため、ルーティングルール、セキュリティリストなど、VNICに関連付けられているすべての機能がPE-VNICで使用できるようになる。

10

【0065】

サービスプロバイダは、PEを介したアクセスを可能にするためにサービスを登録することができる。プロバイダは、サービスの可視性を顧客テナントに制限するポリシーをサービスに関連付けることができる。プロバイダは、特にマルチテナントサービスの場合、単一の仮想IPアドレス(VIP)の下に複数のサービスを登録できる。同じサービスを表すこのようなプライベートエンドポイントが(複数のVCNに)複数存在し得る。

【0066】

20

プライベートサブネット内のコンピューティングインスタンスは、次いで、PE-VNICのプライベートIPアドレスまたはサービスDNS名を使用してサービスにアクセスできる。顧客VCN内のコンピューティングインスタンスは、顧客VCN内のPEのプライベートIPアドレスにトラフィックを送信することでサービスにアクセスできる。プライベートアクセスゲートウェイ(PAGW)130は、サービスプロバイダVCN(例えば、サービスネットワーク110内のVCN)に接続できるゲートウェイリソースであり、顧客サブネットのプライベートエンドポイントとの間のすべてのトラフィックの入口/出口ポイントとして機能する。PAGW130を使用すると、プロバイダは内部IPアドレスリソースを利用せずにPE接続の数を拡張できる。プロバイダは、単一のVCNに登録されている任意の数のサービスに対して1つのPAGWのみを構成する必要がある。プロバイダは、1人または複数の顧客の複数のVCN内のプライベートエンドポイントとしてサービスを表すことができる。顧客の観点から見ると、PE-VNICは顧客のインスタンスに接続されているのではなく、顧客が対話したいサービスに接続されているように見える。プライベートエンドポイント宛てのトラフィックは、PAGW130経由でサービスにルーティングされる。これらは、顧客からサービスへのプライベート接続(C2S接続)と呼ばれる。

30

【0067】

PE概念は、トラフィックがFastConnect/IPsecリンクと顧客VCN内のプライベートエンドポイントを通過できるようにすることで、サービスのプライベートアクセスを顧客のオンプレミスネットワークとデータセンターに拡張するためにも使用できる。LPG132と顧客のVCN内のPEの間でトラフィックを流すことによって、サービスのプライベートアクセスを顧客のピアVCNに拡張することもできる。

40

【0068】

顧客はサブネットレベルでVCN内のルーティングを制御できるため、VCN104などの顧客のVCN内のどのサブネットが各ゲートウェイを使用するかを指定できる。VCNのルートテーブルは、特定のゲートウェイを介してVCNからのトラフィックを許可するかどうかを決定するために使用される。例えば、特定のインスタンスでは、顧客VCN104内のパブリックサブネットのルートテーブルは、IGW120を介して非局所トラフィックを送信し得る。同じ顧客VCN104内のプライベートサブネットのルートテーブルは、CSPサービス宛てのトラフィックをSGW126を介して送信し得る。残りの

50

トラフィックはすべてNATゲートウェイ128経由で送信され得る。ルートテーブルは、VCNから送信されるトラフィックのみを制御する。

【0069】

VCNに関連付けられるセキュリティリストは、ゲートウェイを介して受信接続でVCNに入るトラフィックを制御するために使用される。サブネット内のすべてのリソースは同じルートテーブルとセキュリティリストを使用する。セキュリティリストは、VCNのサブネット内のインスタンスに出入りできる特定のタイプのトラフィックを制御するために使用できる。セキュリティリストルールは、入口（受信）ルールと出口（送信）ルールを含み得る。例えば、入口ルールでは許可されるソースアドレス範囲を指定し、出口ルールでは許可される宛先アドレス範囲を指定することができる。セキュリティルールでは、特定のプロトコル（例えば、TCP、ICMP）、特定のポート（例えば、SSHの場合は22、Windows RDPの場合は3389）などを指定できる。特定の実装では、インスタンスのオペレーティングシステムが、セキュリティリストルールに準拠した独自のファイアウォールルールを適用し得る。ルールはステートフル（例えば、接続が追跡され、応答トラフィックに対する明示的なセキュリティリストルールがなくても応答が自動的に許可される）またはステートレスになり得る。

10

【0070】

顧客VCNからのアクセス（すなわち、VCN104上に展開されるリソースまたはコンピューティングインスタンスによるアクセス）は、パブリックアクセス、プライベートアクセス、または専用アクセスに分類できる。パブリックアクセスとは、パブリックIPアドレスまたはNATを使用してパブリックエンドポイントにアクセスするアクセスモデルを指す。プライベートアクセスにより、プライベートIPアドレス（例えば、プライベートサブネット内のリソース）を持つVCN104内の顧客ワークロードは、インターネットなどのパブリックネットワークを経由せずにサービスにアクセスできるようになる。特定の実施形態では、CSP101により、プライベートIPアドレスを持つ顧客VCNワークロードがサービスゲートウェイを使用してサービス（のパブリックサービスエンドポイント）にアクセスできるようになる。したがって、サービスゲートウェイは、顧客のVCNと顧客のプライベートネットワークの外部にあるサービスのパブリックエンドポイントとの間に仮想リンクを確立することにより、プライベートアクセスモデルを提供する。

20

30

【0071】

さらに、CSP1は、FastConnectパブリックピアリングなどの技術を使用して専用のパブリックアクセスを提供し得る。顧客のオンプレミスインスタンスは、FastConnect接続を使用して、インターネットなどのパブリックネットワークを経由せずに、顧客VCN内の1つまたは複数のサービスにアクセスできる。CSP1は、FastConnectプライベートピアリングを使用した専用プライベートアクセスも提供し得、プライベートIPアドレスを持つ顧客のオンプレミスインスタンスがFastConnect接続を使用して顧客のVCNワークロードにアクセスできるようになる。FastConnectは、パブリックインターネットを使用して顧客のオンプレミスネットワークをCSP1およびそのサービスに接続するためのネットワーク接続の代替手段である。FastConnectは、インターネットベースの接続と比較して、より高い帯域幅任意選択肢と、より信頼性が高く一貫性のあるネットワークエクスペリエンスを備えた専用のプライベート接続を作成するための、簡単、柔軟、かつ経済的な方法を提供する。

40

【0072】

図1および上記の付随する説明は、例示的な仮想ネットワーク内のさまざまな仮想化される構成要素について説明している。前述のように、仮想ネットワークは、基盤となる物理ネットワークまたはサブストレートネットワーク上に構築される。図2は、特定の実施形態による仮想ネットワークのアンダーレイを提供するCSP1200内の物理ネットワーク内の物理構成要素の簡略化されるアーキテクチャ図を示している。図に示すように、CSP1200は、クラウドサービスプロバイダ(CSP)によって提供される構成要素

50

とリソース（例えば、コンピューティング、メモリ、ネットワークリソース）で構成される分散環境を提供する。これらの構成要素とリソースは、加入している顧客、つまりCSPが提供する1つまたは複数のサービスに加入している顧客にクラウドサービス（例えば、IaaSサービス）を提供するために使用される。顧客が加入しているサービスに基づいて、CSP I 200のリソースのサブセット（例えば、コンピューティング、メモリ、ネットワークリソース）が顧客に対してプロビジョニングされる。顧客は、次いで、CSP I 200が提供する物理的なコンピューティング、メモリ、およびネットワークリソースを使用して、独自のクラウドベース（つまり、CSP I ホスト）のカスタマイズ可能なプライベート仮想ネットワークを構築できる。前述のように、これらの顧客ネットワークは仮想クラウドネットワーク（VCN）と呼ばれる。顧客は、コンピューティングインスタンスなどの1つまたは複数の顧客リソースをこれらの顧客VCNに展開できる。コンピューティングインスタンスは、仮想マシン、ベアメタルインスタンスなどの形式になり得る。CSP I 200は、顧客が可用性の高いホスト環境で幅広いアプリケーションとサービスを構築および実行できるようにするインフラストラクチャと補完的なクラウドサービスのセットを提供する。

10

20

30

40

50

【0073】

図2に示す実施例では、CSP I 200の物理構成要素には、1つまたは複数の物理ホストマシンまたは物理サーバ（例えば、202、206、208）、ネットワーク仮想化装置（NVD）（例えば、210、212）、トップオブラック（TOR）スイッチ（例えば、214、216）、物理ネットワーク（例えば、218）、および物理ネットワーク218内のスイッチが含まれる。物理ホストマシンまたはサーバは、VCNの1つまたは複数のサブネットに参加するさまざまなコンピューティングインスタンスをホストおよび実行できる。コンピューティングインスタンスには、仮想マシンインスタンスおよびベアメタルインスタンスが含まれ得る。例えば、図1に示されているさまざまなコンピューティングインスタンスは、図2に示されている物理ホストマシンによってホストされ得る。VCN内の仮想マシンコンピューティングインスタンスは、1台のホストマシンまたは複数の異なるホストマシンによって実行され得る。物理ホストマシンは、仮想ホストマシン、コンテナベースのホストまたは機能などもホストできる。図1に示されているVNICおよびVCN VRは、図2に示されているNVDによって実行され得る。図1に示すゲートウェイは、図2に示すホストマシンおよび/またはNVDによって実行され得る。

【0074】

ホストマシンまたはサーバは、ホストマシン上に仮想化環境を作成して有効にするハイパーバイザ（仮想マシンモニタまたはVMMとも呼ばれる）を実行することができる。仮想化または仮想化環境により、クラウドベースのコンピューティングが容易になる。ホストマシン上のハイパーバイザによって、ホストマシン上で1つまたは複数のコンピューティングインスタンスが作成、実行、管理され得る。ホストマシン上のハイパーバイザにより、ホストマシンの物理的なコンピューティングリソース（例えば、コンピューティング、メモリ、ネットワークリソース）を、ホストマシンによって実行されるさまざまなコンピューティングインスタンス間で共有できるようになる。

【0075】

例えば、図2に示すように、ホストマシン202および208は、それぞれハイパーバイザ260および266を実行する。これらのハイパーバイザは、ソフトウェア、ファームウェア、ハードウェア、またはそれらの組み合わせを使用して実装できる。通常、ハイパーバイザは、ホストマシンのオペレーティングシステム（OS）の上に配置され、次いで、ホストマシンのハードウェアプロセッサ上で実行されるプロセスまたはソフトウェア層である。ハイパーバイザは、ホストマシンの物理コンピューティングリソース（例えば、プロセッサ/コアなどの処理リソース、メモリリソース、ネットワークリソース）を、ホストマシンによって実行されるさまざまな仮想マシンコンピューティングインスタンス間で共有できるようにすることで、仮想化環境を提供する。例えば、図2では、ハイパーバイザ260はホストマシン202のOSの上に配置され、ホストマシン202のコンピ

ューティングリソース（例えば、処理、メモリ、およびネットワークリソース）を、ホストマシン202によって実行されるコンピューティングインスタンス（例えば、仮想マシン）間で共有できるようにする。仮想マシンには独自のオペレーティングシステム（ゲストオペレーティングシステムと呼ばれる）を持たせることができ、これはホストマシンのOSと同じでも異なってもかまわない。ホストマシンによって実行される仮想マシンのオペレーティングシステムは、同じホストマシンによって実行される別の仮想マシンのオペレーティングシステムと同じでも異なってもかまわない。したがって、ハイパーバイザを使用すると、ホストマシンの同じコンピューティングリソースを共有しながら、複数のオペレーティングシステムを相互に並行して実行できるようになる。図2に示されているホストマシンには、同じタイプのハイパーバイザが搭載されている場合も、異なるタイプのハイパーバイザが搭載されている場合もある。

10

【0076】

コンピューティングインスタンスは、仮想マシンインスタンスまたはベアメタルインスタンスであり得る。図2では、ホストマシン202上のコンピューティングインスタンス268とホストマシン208上のコンピューティングインスタンス274が仮想マシンインスタンスの例である。ホストマシン206は、顧客に提供されるベアメタルインスタンスの例である。

【0077】

場合によっては、ホストマシン全体が単一の顧客にプロビジョニングされ得、そのホストマシンによってホストされる1つまたは複数のコンピューティングインスタンス（仮想マシンまたはベアメタルインスタンス）のすべてが同じ顧客に属する。他の場合には、ホストマシンが複数の顧客（つまり、複数のテナント）間で共有され得る。このようなマルチテナントシナリオでは、ホストマシンは異なる顧客に属する仮想マシンコンピューティングインスタンスをホストし得る。これらのコンピューティングインスタンスは、異なる顧客の異なるVCNのメンバーであり得る。特定の実施形態では、ベアメタルコンピューティングインスタンスは、ハイパーバイザのないベアメタルサーバによってホストされる。ベアメタルコンピューティングインスタンスがプロビジョニングされると、単一の顧客またはテナントが、ベアメタルインスタンスをホストするホストマシンの物理CPU、メモリ、およびネットワークインターフェースの制御を維持し、ホストマシンは他の顧客またはテナントと共有されない。

20

30

【0078】

前述のように、VCNの一部である各コンピューティングインスタンスは、コンピューティングインスタンスがVCNのサブネットのメンバーになることを可能にするVNICに関連付けられている。コンピューティングインスタンスに関連付けられているVNICは、コンピューティングインスタンスとの間のパケットまたはフレームの通信を容易にする。VNICは、コンピューティングインスタンスが作成されるときにコンピューティングインスタンスに関連付けられる。特定の実施形態では、ホストマシンによって実行されるコンピューティングインスタンスの場合、そのコンピューティングインスタンスに関連付けられているVNICは、ホストマシンに接続されるNVDによって実行される。例えば、図2では、ホストマシン202は、VNIC276に関連付けられる仮想マシンコンピューティングインスタンス268を実行し、VNIC276は、ホストマシン202に接続されるNVD210によって実行される。別の例として、ホストマシン206によってホストされるベアメタルインスタンス272は、ホストマシン206に接続されるNVD212によって実行されるVNIC280に関連付けられる。さらに別の例として、VNIC284はホストマシン208によって実行されるコンピューティングインスタンス274に関連付けられ、VNIC284はホストマシン208に接続されるNVD212によって実行される。

40

【0079】

ホストマシンによってホストされるコンピューティングインスタンスの場合、そのホストマシンに接続されるNVDは、コンピューティングインスタンスがメンバーであるVC

50

Nに対応するVCN VRも実行する。例えば、図2に示す実施形態では、NVD 210は、コンピューティングインスタンス268がメンバーであるVCN VR 277を実行する。NVD 212は、ホストマシン206および208によってホストされるコンピューティングインスタンスに対応するVCN VR 283を実行することもできる。

【0080】

ホストマシンには、ホストマシンを他の装置に接続できるようにする1つまたは複数のネットワークインターフェースカード(NIC)が含まれ得る。ホストマシン上のNICは、ホストマシンを別の装置に通信可能に接続できるようにする1つまたは複数のポート(またはインターフェース)を提供し得る。例えば、ホストマシンは、ホストマシンとNVDに提供されている1つまたは複数のポート(またはインターフェース)を使用してNVDに接続できる。ホストマシンも、別のホストマシンなどの他の装置に接続され得る。

10

【0081】

例えば、図2では、ホストマシン202は、ホストマシン202のNIC 232によって提供されるポート234とNVD 210のポート236との間に延びるリンク220を使用してNVD 210に接続されている。ホストマシン206は、ホストマシン206のNIC 244によって提供されるポート246とNVD 212のポート248との間に延びるリンク224を使用してNVD 212に接続される。ホストマシン208は、ホストマシン208のNIC 250によって提供されるポート252とNVD 212のポート254との間に延びるリンク226を使用してNVD 212に接続される。

20

【0082】

NVDは、次いで、通信リンクを介してトップオブザラック(TOR)スイッチに接続され、トップオブザラックスイッチは物理ネットワーク218(スイッチファブリックとも呼ばれる)に接続される。特定の実施形態では、ホストマシンとNVD間のリンク、およびNVDとTORスイッチ間のリンクはイーサネット(登録商標)リンクである。例えば、図2では、NVD 210および212は、リンク228および230を使用して、それぞれTORスイッチ214および216に接続されている。特定の実施形態では、リンク220、224、226、228、および230はイーサネットリンクである。TORに接続されるホストマシンとNVDの集合は、ラックと呼ばれることもある。

30

【0083】

物理ネットワーク218は、TORスイッチが相互に通信できるようにする通信ファブリックを提供する。物理ネットワーク218は、多層ネットワークにすることができる。特定の実装では、物理ネットワーク218はスイッチの多層 Clos ネットワークであり、TORスイッチ214および216は多層およびマルチノードの物理スイッチングネットワーク218のリーフレベルノードを表す。2層ネットワーク、3層ネットワーク、4層ネットワーク、5層ネットワーク、および一般に「n」層ネットワークを含むがこれらに限定されない、さまざまな Clos ネットワーク構成が可能である。Clos ネットワークの例を図5に示し、以下に説明する。

【0084】

ホストマシンとNVDの間では、1対1構成、多対1構成、1対多構成など、さまざまな接続構成が可能である。1対1構成の実装では、各ホストマシンは独自の個別のNVDに接続される。例えば、図2では、ホストマシン202は、ホストマシン202のNIC 232を介してNVD 210に接続されている。多対1構成では、複数のホストマシンが1つのNVDに接続される。例えば、図2では、ホストマシン206および208は、それぞれNIC 244および250を介して同じNVD 212に接続されている。

40

【0085】

1対多構成では、1台のホストマシンが複数のNVDに接続される。図3は、ホストマシンが複数のNVDに接続されているCSP 300内の例を示している。図3に示すように、ホストマシン302は、複数のポート306および308を含むネットワークインターフェースカード(NIC) 304を備えている。ホストマシン300は、ポート30

50

6 およびリンク 3 2 0 を介して第 1 の N V D 3 1 0 に接続され、ポート 3 0 8 およびリンク 3 2 2 を介して第 2 の N V D 3 1 2 に接続される。ポート 3 0 6 および 3 0 8 はイーサネットポートであり得、ホストマシン 3 0 2 と N V D 3 1 0 および 3 1 2 間のリンク 3 2 0 および 3 2 2 はイーサネットリンクであり得る。N V D 3 1 0 は、次いで、第 1 の T O R スイッチ 3 1 4 に接続され、N V D 3 1 2 は第 2 の T O R スイッチ 3 1 6 に接続される。N V D 3 1 0 と 3 1 2 および T O R スイッチ 3 1 4 と 3 1 6 間のリンクは、イーサネットリンクであり得る。T O R スイッチ 3 1 4 および 3 1 6 は、多層物理ネットワーク 3 1 8 内の層 0 スイッチング装置を表す。

【 0 0 8 6 】

図 3 に示す配置は、物理スイッチネットワーク 3 1 8 からホストマシン 3 0 2 への 2 つの別個の物理ネットワークパスを提供する。第 1 のパスは、T O R スイッチ 3 1 4 を経由して N V D 3 1 0 を経由してホストマシン 3 0 2 へ、第 2 のパスは T O R スイッチ 3 1 6 を経由して N V D 3 1 2 を経由してホストマシン 3 0 2 へである。別々のパスにより、ホストマシン 3 0 2 の可用性が向上する（高可用性と呼ばれる）。いずれかのパスに問題がある場合（例えば、いずれかのパスのリンクがダウンしている場合）または装置に問題がある場合（例えば、特定の N V D が機能していない場合）、ホストマシン 3 0 2 との通信に他のパスが使用され得る。

【 0 0 8 7 】

図 3 に示す構成では、ホストマシンは、ホストマシンの N I C によって提供される 2 つの異なるポートを使用して 2 つの異なる N V D に接続される。他の実施形態では、ホストマシンには、ホストマシンを複数の N V D に接続できるようにする複数の N I C が含まれ得る。

【 0 0 8 8 】

図 2 に戻ると、N V D は、1 つまたは複数のネットワークおよび / または記憶仮想化機能を実行する物理装置または構成要素である。N V D は、1 つまたは複数の処理装置（例えば、C P U、ネットワーク処理装置（N P U）、F P G A、パケット処理パイプライン）、キャッシュを含むメモリ、およびポートを備えた任意の装置であり得る。さまざまな仮想化機能は、N V D の 1 つまたは複数の処理装置によって実行されるソフトウェア / ファームウェアによって実行され得る。

【 0 0 8 9 】

N V D は、さまざまな異なる形式で実装することができる。例えば、特定の実施形態では、N V D は、オンボードに組み込みプロセッサを備えたスマート N I C またはインテリジェント N I C と呼ばれるインターフェースカードとして実装される。スマート N I C は、ホストマシン上の N I C とは別の装置である。図 2 では、N V D 2 1 0 および 2 1 2 は、それぞれホストマシン 2 0 2、およびホストマシン 2 0 6 および 2 0 8 に接続されるスマート N I C として実装され得る。

【 0 0 9 0 】

しかし、スマート N I C は N V D 実装の一例にすぎない。他にもさまざまな実装が可能である。例えば、他のいくつかの実装では、N V D または N V D によって実行される 1 つまたは複数の機能は、1 つまたは複数のホストマシン、1 つまたは複数の T O R スイッチ、および C S P I 2 0 0 の他の構成要素に組み込まれるか、またはそれらによって実行され得る。例えば、N V D はホストマシンに実装され得、N V D によって実行される機能はホストマシンによって実行される。別の例として、N V D は T O R スイッチの一部であり得るか、T O R スイッチが N V D によって実行される機能を実行するように構成され、T O R スイッチがパブリッククラウドに使用されるさまざまな複雑なパケット変換を実行できるようにし得る。N V D の機能を実行する T O R は、スマート T O R と呼ばれることもある。さらに他の実装では、ベアメタル（B M）インスタンスではなく仮想マシン（V M）インスタンスが顧客に提供され、N V D によって実行される機能は、ホストマシンのハイパーバイザ内に実装され得る。いくつかの他の実装では、N V D の機能の一部が、一連のホストマシン上で実行される集中型サービスにオフロードされ得る。

10

20

30

40

50

【 0 0 9 1 】

図 2 に示すようにスマートNICとして実装される場合など、特定の実施形態では、NVDは、1つまたは複数のホストマシンおよび1つまたは複数のTORスイッチに接続できるようにする複数の物理ポートを備えることができる。NVD上のポートは、ホスト側ポート（「サウスポート」とも呼ばれる）、またはネットワーク側またはTOR側ポート（「ノースポート」とも呼ばれる）に分類できる。NVDのホスト側ポートは、NVDをホストマシンに接続するために使用されるポートである。図2のホスト側ポートの例には、NVD 210のポート236、NVD 212のポート248および254が含まれる。NVDのネットワーク側ポートは、NVDをTORスイッチに接続するために使用されるポートである。図2のネットワーク側ポートの例には、NVD 210のポート256およびNVD 212のポート258が含まれる。図2に示すように、NVD 210は、NVD 210のポート256からTORスイッチ214まで延びるリンク228を使用してTORスイッチ214に接続される。同様に、NVD 212は、NVD 212のポート258からTORスイッチ216まで延びるリンク230を使用してTORスイッチ216に接続される。

10

【 0 0 9 2 】

NVDは、ホスト側ポートを介してホストマシンからパケットおよびフレーム（例えば、ホストマシンによってホストされるコンピューティングインスタンスによって生成されるパケットおよびフレーム）を受信し、必要なパケット処理を実行した後、NVDのネットワーク側ポートを介してパケットおよびフレームをTORスイッチに転送することができる。NVDは、NVDのネットワーク側ポートを介してTORスイッチからパケットとフレームを受信し、必要なパケット処理を実行した後、NVDのホスト側ポートを介してパケットとフレームをホストマシンに転送し得る。

20

【 0 0 9 3 】

特定の実施形態では、NVDとTORスイッチの間に複数のポートと関連するリンクが存在し得る。これらのポートとリンクは集約され、複数のポートまたはリンクのリンクアグリゲータグループ（LAGと呼ばれる）を形成し得る。リンク集約により、2つのエンドポイント間（例えば、NVDとTORスイッチ間）の複数の物理リンクを1つの論理リンクとして扱うことができる。特定のLAG内のすべての物理リンクは、同じ速度で全二重モードで動作し得る。LAGは、2つのエンドポイント間の接続の帯域幅と信頼性を向上させるのに役立つ。LAG内の物理リンクの1つがダウンした場合、トラフィックはLAG内の他の物理リンクの1つに動的かつ透過的に再割り当てされる。集約される物理リンクは、個々のリンクよりも高い帯域幅を提供する。LAGに関連付けられる複数のポートは、単一の論理ポートとして扱われる。トラフィックは、LAGの複数の物理リンク間で負荷分散できる。2つのエンドポイント間に1つまたは複数のLAGを構成できる。2つのエンドポイントは、NVDとTORスイッチの間、ホストマシンとNVDの間などにあり得る。

30

【 0 0 9 4 】

NVDは、ネットワーク仮想化機能を実装または実行する。これらの機能は、NVDによって実行されるソフトウェア/ファームウェアによって実行される。ネットワーク仮想化機能の例には、パケットのカプセル化およびカプセル化解除機能、VCNネットワークを作成するための機能、VCNセキュリティリスト（ファイアウォール）機能などのネットワークポリシーを実装するための機能、VCN内のコンピューティングインスタンスとの間のパケットのルーティングおよび転送を容易にする機能などが含まれるが、これらに限定されない。特定の実施形態では、パケットを受信すると、NVDは、パケットを処理し、パケットをどのように転送またはルーティングするかを決定するためのパケット処理パイプラインを実行するように構成される。このパケット処理パイプラインの一部として、NVDは、VCN内のcisに関連付けられているVNICの実行、VCNに関連付けられる仮想ルーター（VR）の実行、仮想ネットワーク内での転送またはルーティングを容易にするためのパケットのカプセル化とカプセル化解除、特定のゲートウェイ（例えば

40

50

、局所ピアリングゲートウェイ)の実行、セキュリティリスト、ネットワークセキュリティグループの実装、ネットワークアドレス変換(NAT)機能(例えば、ホストごとにパブリックIPをプライベートIPに変換する)、スロットル機能、およびその他の機能など、オーバーレイネットワークに関連付けられる1つまたは複数の仮想機能を実行し得る。

【0095】

特定の実施形態では、NVD内のパケット処理データパスは、各々が一連のパケット変換段階から構成される複数のパケットパイプラインを含み得る。特定の実装では、パケットを受信すると、パケットは解析され、単一のパイプラインに分類される。その後、パケットはドロップされるか、NVDのインターフェースを介して送信されるまで、段階的に線形に処理される。これらの段階は、基本的な機能パケット処理構築ブロック(例えば、ヘッダーの検証、スロットルの適用、新しい層2ヘッダーの挿入、L4ファイアウォールの適用、VCNカプセル化/カプセル化解除など)を提供するため、既存の段階を組み合わせる新しいパイプラインを構築したり、新しい段階を作成して既存のパイプラインに挿入することで新しい機能を追加したりできる。

10

【0096】

NVDは、VCNの制御プレーンとデータプレーンに対応する制御プレーン機能とデータプレーン機能の両方を実行できる。VCN制御プレーンの例は、図12、13、14、および15(参照1216、1316、1416、および1516を参照)にも示されており、以下で説明する。VCNデータプレーンの例は、図12、13、14、および15(参照1218、1318、1418、および1518を参照)に示されており、以下で説明する。制御プレーン機能には、データの転送方法を制御するネットワークの構成(例えば、ルートとルートテーブルの設定、VNICの構成)に使用される機能が含まれる。特定の実施形態では、オーバーレイからサブストレートへのすべてのマッピングを一元的に計算し、それらをNVDおよびDRG、SGW、IGWなどのさまざまなゲートウェイなどの仮想ネットワークエッジ装置に公開するVCN制御プレーンが提供される。ファイアウォールルールも同じメカニズムを使用して公開できる。特定の実施形態では、NVDはそのNVDに関連するマッピングのみを取得する。データプレーン機能には、制御プレーンを使用して設定される構成に基づいてパケットを実際にルーティング/転送する機能が含まれる。VCNデータプレーンは、顧客のネットワークパケットが基盤ネットワークを通過する前にカプセル化することによって実装される。カプセル化/カプセル化解除機能はNVDに実装されている。特定の実施形態では、NVDは、ホストマシンに出入りするすべてのネットワークパケットを傍受し、ネットワーク仮想化機能を実行するように構成される。

20

30

【0097】

前述のように、NVDはVNICやVCN VRなどのさまざまな仮想化機能を実行する。NVDは、VNICに接続される1つまたは複数のホストマシンによってホストされるコンピューティングインスタンスに関連付けられているVNICを実行し得る。例えば、図2に示すように、NVD210は、NVD210に接続されるホストマシン202によってホストされるコンピューティングインスタンス268に関連付けられているVNIC276の機能を実行する。別の例として、NVD212は、ホストマシン206によってホストされるベアメタルコンピューティングインスタンス272に関連付けられているVNIC280を実行し、ホストマシン208によってホストされるコンピューティングインスタンス274に関連付けられているVNIC284を実行する。ホストマシンは、異なる顧客に属する異なるVCNに属するコンピューティングインスタンスをホストすることができ、ホストマシンに接続されるNVDは、コンピューティングインスタンスに対応するVNICを実行する(つまり、VNIC関連機能を実行する)ことができる。

40

【0098】

NVDは、コンピューティングインスタンスのVCNに対応するVCN仮想ルータも実行する。例えば、図2に示す実施形態では、NVD210は、コンピューティングインス

50

タンス 268 が属する VCN に対応する VCN VR 277 を実行する。NVD 212 は、ホストマシン 206 および 208 によってホストされるコンピューティングインスタンスが属する 1 つまたは複数の VCN に対応する 1 つまたは複数の VCN VR 283 を実行する。特定の実施形態では、その VCN に対応する VCN VR は、その VCN に属する少なくとも 1 つのコンピューティングインスタンスをホストするホストマシンに接続されるすべての NVD によって実行される。ホストマシンが異なる VCN に属するコンピューティングインスタンスをホストしている場合、そのホストマシンに接続される NVD は、それらの異なる VCN に対応する VCN VR を実行し得る。

【0099】

VNIC および VCN VR に加えて、NVD はさまざまなソフトウェア（例えば、デーモン）を実行し、NVD によって実行されるさまざまなネットワーク仮想化機能を容易にする 1 つまたは複数のハードウェア構成要素を含み得る。簡潔にするために、これらのさまざまな構成要素は、図 2 に示すように「パケット処理構成要素」としてグループ化されている。例えば、NVD 210 はパケット処理構成要素 286 を含み、NVD 212 はパケット処理構成要素 288 を含む。例えば、NVD のパケット処理構成要素には、NVD のポートおよびハードウェアインターフェースと対話して、NVD によって受信され、NVD を使用して通信されるすべてのパケットを監視し、ネットワーク情報を記憶するように構成されるパケットプロセッサが含まれ得る。ネットワーク情報には、例えば、NVD によって処理されるさまざまなネットワークフローを識別するネットワークフロー情報と、フローごとの情報（例えば、フローごとの統計）が含まれ得る。特定の実施形態では、ネットワークフロー情報は VNIC ごとに記憶され得る。パケットプロセッサは、パケットごとの操作を実行できるほか、ステートフル NAT および L4 ファイアウォール（FW）を実装することもできる。別の例として、パケット処理構成要素には、NVD によって記憶される情報を 1 つまたは複数の異なるレプリケーションターゲットストアに複製するように構成されるレプリケーションエージェントが含まれ得る。さらに別の例として、パケット処理構成要素には、NVD のログ機能を実行するように構成されるログエージェントが含まれ得る。パケット処理構成要素には、NVD の性能と健全性を監視するソフトウェアや、NVD に接続される他の構成要素の状態と健全性を監視するソフトウェアも含まれ得る。

【0100】

図 1 は、VCN、VCN 内のサブネット、サブネット上に展開されるコンピューティングインスタンス、コンピューティングインスタンスに関連付けられている VNIC、VCN の VR、および VCN 用に構成されるゲートウェイのセットを含む、仮想ネットワークまたはオーバーレイネットワークの例の構成要素を示している。図 1 に示されているオーバーレイ構成要素は、図 2 に示されている 1 つまたは複数の物理構成要素によって実行またはホストされ得る。例えば、VCN 内のコンピューティングインスタンスは、図 2 に示す 1 つまたは複数のホストマシンによって実行またはホストされ得る。ホストマシンによってホストされるコンピューティングインスタンスの場合、そのコンピューティングインスタンスに関連付けられている VNIC は通常、そのホストマシンに接続される NVD によって実行される（つまり、VNIC 機能は、そのホストマシンに接続される NVD によって提供される）。VCN の VCN VR 機能は、その VCN の一部であるコンピューティングインスタンスをホストまたは実行するホストマシンに接続されているすべての NVD によって実行される。VCN に関連付けられるゲートウェイは、1 つまたは複数の異なるタイプの NVD によって実行され得る。例えば、特定のゲートウェイはスマート NIC によって実行され得る一方、他のゲートウェイは 1 つまたは複数のホストマシンまたは NVD の他の実装によって実行され得る。

【0101】

前述のように、顧客 VCN 内のコンピューティングインスタンスは、さまざまな異なるエンドポイントと通信できる。エンドポイントは、ソースコンピューティングインスタンスと同じサブネット内にある場合もあれば、ソースコンピューティングインスタンスと同

10

20

30

40

50

じVCN内にあるが別のサブネット内にある場合もあり、あるいはソースコンピューティングインスタンスのVCN外にあるエンドポイントと通信できる。これらの通信は、コンピューティングインスタンスに関連付けられているVNIC、VCN VR、およびVCNに関連付けられるゲートウェイを使用して容易になる。

【0102】

VCN内の同じサブネット上の2つのコンピューティングインスタンス間の通信では、ソースコンピューティングインスタンスと宛先コンピューティングインスタンスに関連付けられているVNICを使用して通信が容易になる。ソースと宛先のコンピューティングインスタンスは、同じホストマシンでホストすることも、異なるホストマシンでホストすることもできる。ソースコンピューティングインスタンスから発信されるパケットは、ソースコンピューティングインスタンスをホストしているホストマシンから、そのホストマシンに接続されるNVDに転送され得る。NVDでは、パケットはパケット処理パイプラインを使用して処理される。これには、ソースコンピューティングインスタンスに関連付けられているVNICの実行が含まれ得る。パケットの宛先エンドポイントは同じサブネット内にあるため、ソースコンピューティングインスタンスに関連付けられているVNICを実行すると、パケットは宛先コンピューティングインスタンスに関連付けられているVNICを実行するNVDに転送され、NVDはパケットを処理して宛先コンピューティングインスタンスに転送する。ソースおよび宛先コンピューティングインスタンスに関連付けられているVNICは、同じNVD上で実行される場合（例えば、ソースおよび宛先コンピューティングインスタンスの両方が同じホストマシンによってホストされている場合）と、異なるNVD上で実行される場合（例えば、ソースおよび宛先コンピューティングインスタンスが異なるNVDに接続される異なるホストマシンによってホストされている場合）がある。VNICは、NVDによって記憶されるルーティング/転送テーブルを使用して、パケットの次のホップを決定し得る。

【0103】

サブネット内のコンピューティングインスタンスから同じVCN内の異なるサブネット内のエンドポイントにパケットを通信する場合、ソースコンピューティングインスタンスから発信されるパケットは、ソースコンピューティングインスタンスをホストしているホストマシンから、そのホストマシンに接続されるNVDに通信される。NVDでは、パケットはパケット処理パイプラインを使用して処理される。これには、1つまたは複数のVNICの実行と、VCNに関連付けられているVRが含まれ得る。例えば、パケット処理パイプラインの一部として、NVDはソースコンピューティングインスタンスに関連付けられているVNICに対応する機能を実行または呼び出す（VNICを実行するとも呼ばれる）。VNICによって実行される機能には、パケット上のVLANタグの確認が含まれ得る。パケットの宛先がサブネット外にあるため、次にVCN VR機能が呼び出され、NVDによって実行される。次に、VCN VRは、宛先コンピューティングインスタンスに関連付けられているVNICを実行するNVDにパケットをルーティングする。次に、宛先コンピューティングインスタンスに関連付けられているVNICがパケットを処理し、パケットを宛先コンピューティングインスタンスに転送する。ソースおよび宛先コンピューティングインスタンスに関連付けられているVNICは、同じNVD上で実行される場合（例えば、ソースおよび宛先コンピューティングインスタンスの両方が同じホストマシンによってホストされている場合）と、異なるNVD上で実行される場合（例えば、ソースおよび宛先コンピューティングインスタンスが異なるNVDに接続される異なるホストマシンによってホストされている場合）がある。

【0104】

パケットの宛先がソースコンピューティングインスタンスのVCNの外部にある場合、ソースコンピューティングインスタンスから発信されるパケットは、ソースコンピューティングインスタンスをホストしているホストマシンから、そのホストマシンに接続されるNVDに通信される。NVDは、ソースコンピューティングインスタンスに関連付けられているVNICを実行する。パケットの宛先エンドポイントはVCN外にあるため、パケ

10

20

30

40

50

ットはそのV C NのV C N V Rによって処理される。N V DはV C N V R機能を呼び出し、その結果、パケットがV C Nに関連付けられる適切なゲートウェイを実行するN V Dに転送され得る。例えば、宛先が顧客のオンプレミスネットワーク内のエンドポイントである場合、パケットはV C N V Rによって、V C N用に構成されるD R Gゲートウェイを実行するN V Dに転送され得る。V C N V Rは、ソースコンピューティングインスタンスに関連付けられているV N I Cを実行するN V Dと同じN V D上で実行されることも、別のN V Dによって実行されることもある。ゲートウェイは、スマートN I C、ホストマシン、またはその他のN V D実装であり得るN V Dによって実行され得る。その後、パケットはゲートウェイによって処理され、次のホップに転送され、パケットが目的の宛先エンドポイントに通信しやすくする。例えば、図2に示す実施形態では、コンピューティングインスタンス268から発信されるパケットは、リンク220を介して(N I C 232を使用して)ホストマシン202からN V D 210に通信され得る。N V D 210では、ソースコンピューティングインスタンス268に関連付けられているV N I Cであるため、V N I C 276が呼び出される。V N I C 276は、パケット内のカプセル化される情報を調べ、パケットを転送するための次のホップを決定して、パケットを目的の宛先エンドポイントに通信しやすくし、決定した次のホップにパケットを転送するように構成されている。

10

【0105】

V C Nに展開されるコンピューティングインスタンスは、さまざまなエンドポイントと通信できる。これらのエンドポイントには、C S P I 200によってホストされるエンドポイントとC S P I 200外部のエンドポイントが含まれ得る。C S P I 200によってホストされるエンドポイントには、同じV C Nまたは他のV C N内のインスタンスが含まれ得る。これらのV C Nは、顧客のV C Nである場合もあれば、顧客に属さないV C Nである場合もある。C S P I 200によってホストされるエンドポイント間の通信は、物理ネットワーク218を介して実行され得る。コンピューティングインスタンスは、C S P I 200によってホストされていないエンドポイント、またはC S P I 200の外部にあるエンドポイントと通信する場合もある。これらのエンドポイントの例には、顧客のオンプレミスネットワークまたはデータセンター内のエンドポイント、またはインターネットなどのパブリックネットワーク経由でアクセス可能なパブリックエンドポイントが含まれる。C S P I 200外部のエンドポイントとの通信は、さまざまな通信プロトコルを使用して、パブリックネットワーク(例えば、インターネット)(図2には示されていない)またはプライベートネットワーク(図2には示されていない)を介して実行できる。

20

30

【0106】

図2に示すC S P I 200のアーキテクチャは単なる一例であり、限定するものではない。代替実施形態では、変形、代替、および修正が可能である。例えば、いくつかの実装では、C S P I 200は、図2に示されているものよりも多いか少ないシステムまたは構成要素を備えていてもよく、2つ以上のシステムを組み合わせてもよく、またはシステムの構成または配置が異なってもよい。図2に示されているシステム、サブシステム、およびその他の構成要素は、ハードウェア、またはそれらの組み合わせを使用して、それぞれのシステムの1つまたは複数の処理装置(例えば、プロセッサ、コア)によって実行されるソフトウェア(例えば、コード、命令、プログラム)で実装され得る。ソフトウェアは、非一時的な記憶媒体(例えば、メモリ装置)に記憶され得る。

40

【0107】

図4は、特定の実施形態によるマルチテナントをサポートするためのI/O仮想化を提供するためのホストマシンとN V Dとの間の接続を示す。図4に示すように、ホストマシン402は仮想化環境を提供するハイパーバイザ404を実行する。ホストマシン402は、顧客/テナント#1に属するV M 1 406と顧客/テナント#2に属するV M 2 408の2つの仮想マシンインスタンスを実行する。ホストマシン402は、リンク414を介してN V D 412に接続される物理N I C 410を含む。各コンピューティングインスタンスは、N V D 412によって実行されるV N I Cに接続される。図4の実施形態

50

では、VM1 406はVNIC - VM1 420に接続され、VM2 408はVNIC - VM2 422に接続されている。

【0108】

図4に示すように、NIC410は、論理NIC A416と論理NIC B418の2つの論理NICを含む。各仮想マシンは、独自の論理NICに接続され、それを使用して動作するように構成される。例えば、VM1 406は論理NIC A416に接続され、VM2 408は論理NIC B418に接続される。ホストマシン402は複数のテナントによって共有される1つの物理NIC410のみを含むが、論理NICがあるため、各テナントの仮想マシンは独自のホストマシンとNICがあるものと認識する。

【0109】

特定の実施形態では、各論理NICに独自のVLAN IDが割り当てられる。したがって、テナント#1の論理NIC A416には特定のVLAN IDが割り当てられ、テナント#2の論理NIC B418には別のVLAN IDが割り当てられる。VM1 406からパケットが通信されると、テナント#1に割り当てられるタグがハイパーバイザによってパケットに接続され、その後、パケットはリンク414を介してホストマシン402からNVD412に通信される。同様に、VM2 408からパケットが通信されると、テナント#2に割り当てられるタグがハイパーバイザによってパケットに接続され、その後、パケットはリンク414を介してホストマシン402からNVD412に通信される。したがって、ホストマシン402からNVD412に通信されるパケット424には、特定のテナントおよび関連するVMを識別する関連タグ426が含まれる。NVD 20
では、ホストマシン402から受信したパケット424について、パケットに関連付けられているタグ426を使用して、パケットがVNIC - VM1 420によって処理されるか、VNIC - VM2 422によって処理されるかが決定される。その後、パケットは対応するVNICによって処理される。図4に示す構成により、各テナントのコンピューティングインスタンスは、独自のホストマシンとNICを所有していると認識できるようになる。図4に示すセットアップでは、マルチテナントをサポートするためのI/O仮想化が提供される。

【0110】

図5は、特定の実施形態による物理ネットワーク500の簡略化されるブロック図を示す。図5に示す実施形態は、Closネットワークとして構成されている。Closネットワークは、高い二分帯域幅と最大限のリソース使用率を維持しながら接続の冗長性を提供するように設計される特定のタイプのネットワークトポロジである。Closネットワークは、非ブロッキング、マルチ段階、またはマルチ層のスイッチングネットワークの一種であり、段階または層の数は2、3、4、5などになり得る。図5に示す実施形態は、層1、層2、および層3を含む3層ネットワークである。TORスイッチ504は、Closネットワーク内の層0スイッチを表す。1つまたは複数のNVDがTORスイッチに接続される。層0スイッチは、物理ネットワークのエッジ装置とも呼ばれる。層0スイッチは、リーフスイッチとも呼ばれる層1スイッチに接続される。図5に示す実施形態では、セットの「n」個の層0TORスイッチがセットの「n」個の層1スイッチに接続され、一緒にポッドを形成する。ポッド内の各層0スイッチはポッド内のすべての層1スイッチに相互接続されているが、ポッド間のスイッチの接続はない。特定の実装では、2つのポッドがブロックと呼ばれる。各ブロックは、セットの「n」個の層2スイッチ（スパインスイッチと呼ばれることもある）によってサービスが提供されるか、またはこれらのスイッチに接続される。物理ネットワークトポロジには複数のブロックが存在し得る。層2スイッチは、次に「n」個の層3スイッチ（スーパースパインスイッチと呼ばれることもある）に接続される。物理ネットワーク500を介したパケットの通信は、通常、1つまたは複数の層3通信プロトコルを使用して実行される。通常、TOR層を除く物理ネットワークのすべての層はn方向の冗長性を備えているため、高可用性が実現する。物理ネットワークのスケーリングを可能にするために、ポッドとブロックにポリシーを指定して、物理ネットワーク内のスイッチ間の可視性を制御することができる。 50

10

20

30

40

50

【0111】

Closネットワークの特徴は、1つの層0スイッチから別の層0スイッチ（または層0スイッチに接続されるNVDから層0スイッチに接続される別のNVD）に到達するまでの最大ホップ数が固定されていることである。例えば、3層のClosネットワークでは、ソースNVDとターゲットNVDがClosネットワークのリーフ層に接続されている場合、パケットが1つのNVDから別のNVDに到達するには最大7ホップが必要である。同様に、4層のClosネットワークでは、ソースNVDとターゲットNVDがClosネットワークのリーフ層に接続されている場合、パケットが1つのNVDから別のNVDに到達するには最大9ホップが必要である。したがって、Closネットワークアーキテクチャは、ネットワーク全体で一貫したレイテンシを維持する。これは、データセンター内およびデータセンター間の通信にとって重要である。Closトポロジは水平方向に拡張でき、コスト効率に優れている。ネットワークの帯域幅/スループット容量は、さまざまな層にスイッチを追加し（例えば、リーフスイッチとスパインスイッチを追加）、隣接する層のスイッチ間のリンク数を増やすことで簡単に増やすことができる。

10

【0112】

特定の実施形態では、CSP内の各リソースには、クラウド識別子（CID）と呼ばれる固有の識別子が割り当てられる。この識別子はリソースの情報の一部として含まれており、例えば、コンソールやAPIなどを通じてリソースを管理するために使用できる。CIDの構文の例は次のとおりである。

```
ocid1. RESOURCE TYPE . REALM .[REGION][.FUTURE USE]. UNIQUE ID
```

20

ocid1：CIDのバージョンを示すリテラル文字列。

resource type：リソースのタイプ（例えば、インスタンス、ボリューム、VCN、サブネット、ユーザ、グループなど）。

realm：リソースが存在するレルム。値の例としては、商用レルムの場合は「c1」、政府クラウドレルムの場合は「c2」、連邦政府クラウドレルムの場合は「c3」などがある。各レルムには独自のドメイン名があり得る。

region：リソースが存在するリージョン。リージョンがリソースに適用されない場合、この部分は空白になることがある。

future use：将来の使用のために予約されている。

30

unique ID：IDの固有の部分。形式は、リソースまたはサービスのタイプによって異なり得る。

【0113】

図6は、特定の実施形態による、CLOSネットワーク配置を組み込んだクラウドインフラストラクチャ600のブロック図を示す。クラウドインフラストラクチャ600には、複数のラック（例えば、ラック1 610、ラック2、620）が含まれる。各ラックには複数のホストマシン（以下、ホストとも呼ばれる）が含まれる。ラック1 610には、つまり、ホスト1-A 612とホスト1-B 614の2つのホストマシンが含まれていることが示されている。ラック2 620には、つまり、ホスト2-A 622とホスト2-B 624の2つのホストマシンが含まれていることが示されている。図6の説明（すなわち、各ラックに2台のホストマシンが含まれる）は、例示を目的としたものであり、限定するものではないことが理解される。例えば、クラウドインフラストラクチャには2台以上のラックが含まれ得、各ラックには2台以上のホストマシンが含まれ得る。さらに、各ラックには同じ数のホストが存在するという制約がないことに留意されたい。むしろ、ラックは、別のラックに含まれるホストマシンの数よりも多い場合も少ない場合もあり得る。

40

【0114】

各ホストマシンには、複数のグラフィック処理装置（GPU）が含まれる。例えば、ホストマシン1-A 612には、N個のGPU（例えば、GPU1、613）が含まれる。さらに、図6の説明では、各ホストマシンに同じ数のGPU、つまりN個のGPUが含ま

50

れているが、これは例示を目的としたものであり、限定するものではなく、つまり、各ホストマシンには異なる数のGPUを含めることができることが理解される。各ラックには、ラック内のホストマシンでホストされているGPUと通信可能に結合されるトップオブラック(TOR)スイッチが含まれている。例えば、ラック1 610には、ホストマシンホスト1-A、612およびホスト1-B、614に通信可能に結合されるTORスイッチ(つまり、TOR1)616が含まれるが、ラック2 620には、ホストマシンホスト2-A、622およびホスト2-K、624に通信可能に結合されるTORスイッチ(つまり、TOR2)626が含まれる。図6に示されているTORスイッチ(つまり、TOR1 616およびTORM626)には各々、ラックに含まれる各ホストマシンでホストされているN個のGPUにTORスイッチを通信可能に結合するために使用されるN個のポートが含まれていることが理解される。図6に示されているTORスイッチとGPUの結合は、例示を目的としており、限定するものではない。例えば、いくつかの実施形態では、TORスイッチは複数のポートを有し得、各ポートは各ホストマシン上のGPUに対応し、つまり、ホストマシン上のGPUは通信リンクを介してTORの固有のポートに接続され得る。

10

【0115】

各ラックのTORスイッチは、複数のスパインスイッチ、例えばスパインスイッチ1、630およびスパインスイッチP640に通信可能に結合されている。例えば、図6に示すように、TOR1 616は、2つのリンクを介してスパインスイッチ1 630に接続され、さらに別の2つのリンクを介してスパインスイッチP640にそれぞれ接続される。特定のTORスイッチからスパインスイッチに送信される情報は、本明細書ではアップリンクを介して行われる通信と呼ばれ、一方、スパインスイッチからTORスイッチに送信される情報は、本明細書ではダウンリンクを介して行われる通信と呼ばれる。いくつかの実施形態によれば、TORスイッチとスパインスイッチはCLOSネットワーク配置(例えば、多段スイッチングネットワーク)で接続され、各TORスイッチはCLOSネットワーク内の「リーフ」ノードを形成する。

20

【0116】

いくつかの実施形態によれば、ホストマシンに含まれるGPUはマシン学習に関連するタスクを実行する。このような設定では、単一のタスクが、複数のホストマシンおよび複数のラックにまたがる可能性のある多数のGPU(例えば、64個のGPU)にわたって実行/分散され得る。これらすべてのGPUは同じタスク(つまり、ワークロード)で動作しているため、すべてが時間同期される方法で相互に通信する必要がある。さらに、どの時点においても、GPUはコンピューティングモードまたは通信モードのいずれかの状態にあり、つまり、GPUはほぼ同時に相互に通信する。ワークロードの速度は、最も遅いGPUの速度によって決定する。

30

【0117】

通常、ソースGPUから宛先GPUにパケットをルーティングするために、等コストマルチパス(ECMP)ルーティングが利用される。ECMPルーティングでは、送信者から受信者へのトラフィックのルーティングに使用できる等コストパスが複数ある場合、選択技術を使用して特定のパスを選択する。したがって、トラフィックを受信するネットワーク装置(例えば、TORスイッチやスパインスイッチ)では、選択アルゴリズムを使用して、ネットワーク装置から後続の装置にトラフィックを転送するために使用される送信リンクが選択される。この送信リンクの選択は、送信者から受信者へのパス内の各ネットワーク装置で行われる。ハッシュベースの選択は、広く使用されているECMP選択技術であり、ハッシュは、例えば、パケットの4組(例えば、ソースポート、宛先ポート、ソースIP、宛先IP)に基づき得る。

40

【0118】

ECMPルーティングはフローを考慮したルーティング技術であり、各フロー(つまり、データパケットのストリーム)はフローの持続期間中、同じパスにハッシュされる。したがって、フロー内のパケットは、特定の送信ポート/リンクを使用してネットワーク装

50

置から転送される。これは通常、フロー内のパケットが順番に到着することを保証するため、つまりパケットの順序変更が不要であることを確認するために行われる。しかし、ECMPルーティングは帯域幅（またはスループット）を認識しない。つまり、TORスイッチとスパインスイッチは、並列リンク上のフローの統計的なフロー認識型（スループット非認識型）ECMP負荷分散を実行する。

【0119】

標準のECMPルーティング（つまり、フローのみを認識するルーティング）では、ネットワーク装置が2つの別々の着信リンクを介して受信したフローが同じ送信リンクにハッシュされ、それによりフロー衝突が発生し得るという問題がある。例えば、2つのフローが2つの別々の100G着信リンクを介して着信し、各フローが同じ100G送信リンクにハッシュされる状況を考慮する。このような状況では、着信帯域幅が200Gであるのに対し、発信帯域幅が100Gであるため、輻輳（つまり、フロー衝突）が発生し、パケットがドロップされる。図6に示すように、2つのフローがある。ホストマシンホスト1-A、612の第1のGPUからTORスイッチ616に向けられるフロー1641と、ホストマシン614上の別のGPUからTORスイッチ616に向けられるフロー2643である。2つのフローは別々のリンクでTORスイッチに向けられることに留意されたい。図6に示されているすべてのリンクの容量（帯域幅）は100Gであると想定されている。TORスイッチ616がECMPルーティングアルゴリズムを実行する場合、2つのフローがハッシュされて、TORの同じ送信リンク（例えば、TORスイッチ616をスパインスイッチ630に接続するリンク650）を使用し得る。この場合、2つのフロー間に衝突が発生し（「X」マークで表される）、パケットがドロップされる。

【0120】

このような衝突シナリオは、プロトコルに関係なく、すべてのタイプのトラフィックで一般的に問題になる。例えば、TCPは、パケットがドロップされ、送信者がドロップされるパケットの確認応答を受信しなかった場合に、パケットが再送信されるという点でインテリジェントである。しかし、リモートダイレクトメモリアクセス（RDMA）タイプのトラフィックの場合、状況は悪化する。RDMAネットワークでは、さまざまな理由によりTCPが使用されない（例えば、TCPの性能が高くない）。RDMAネットワークでは、インフィニバンド上のRDMAや統合イーサネット上のRDMA（RoCE）などのプロトコルが使用される。RoCEには輻輳制御アルゴリズムがあり、送信者が輻輳の発生やパケットのドロップを識別すると、送信者はパケットの送信速度を低下させる。ドロップされるパケットの場合、ドロップされるパケットだけでなく、ドロップされるパケットの周囲のいくつかのパケットも再送信されるため、使用可能な帯域幅がさらに消費され、性能が低下する。

【0121】

フロー衝突問題は、厳格な時間同期要件のため、GPUにとって重大な問題である。例えば、前述のように、GPUは、すべてのGPUが時間同期した方法で相互に通信するマシン学習タスク（つまり、ワークロード）を実行し得る。マシン学習タスクやその他のタイプのタスクの場合、GPU間の通信を可能にするために、ホストマシンの論理トポロジ（例えば、リングトポロジ、ツリートポロジなど）が構築される。GPUは、複数レベルまたは多次元であり得る論理トポロジを使用して相互に接続する。いくつかの実装では、ワークロードを実行するために、アプリケーションはGPUを相互接続するための仮想（または論理）トポロジを構築する。通常、このようなアプリケーションは、ホストマシンの基礎となる物理トポロジを認識しないため、ランダム（つまり、任意の）方法で論理トポロジを構築しようとする。このようにランダムに構築される論理トポロジにより、GPUホストマシンは、局所ネットワーク近隣に他のGPUホストが存在するかどうかに関係なく、トラフィックを交換することになる。そのため、トラフィックの輻輳が発生する可能性が高まり、GPUスループットが低下する。例えば、特定のマシン学習タスクを実行するためにホストマシンのペアが必要な場合を考慮する。この場合、ホストマシンのペアのランダム選択が実行され、ホストマシンの1つが第1の局所近隣（例えば、第1ラック

)にあり、もう1つのホストマシンが別の局所近隣(第1の局所近隣とは異なる)、例えば第2ラックにある場合、マシン学習タスクの実行によって、一定のレイテンシ(例えば、第1のホストマシンと第2のホストマシン間の通信で発生する遅延)が発生し、トラフィックの輻輳の可能性も高まる可能性がある。対照的に、マシン学習タスクを実行するために選択されるホストマシンのペアが同じ局所近隣に存在する場合(例えば、同じラックに存在する場合)、ホストマシン間の通信によるレイテンシが最小限に抑えられ、トラフィックの輻輳を回避する可能性が高くなることが理解される。

【0122】

以下に、上記の問題を克服するための技術について説明する。具体的には、本明細書で説明する技術は、論理トポロジを構築するプロセスでGPUの局所性情報を活用し、それによって unnecessary のトラフィックの輻輳を回避する。さらに、本開示の実施形態は、顧客がワークロードを近くのホストマシンに「配置」することによって、アプリケーションからサービスまでのレイテンシを削減することを可能にする。さらに、顧客は局所性情報を使用して、より高い反親和性が得られるようにワークロードを配置するため、リソースの共有結果を減らすことでより高い回復力を得ることができる。

10

【0123】

図7を参照すると、特定の実施形態によるラック700の例示的な構成が示されている。図7に示すように、ラック700には、つまり、ホスト1-A710とホスト1-B720の2つのホストマシンが含まれている。ラック700には2台のホストマシンが含まれているように示されているが、ラック700にはさらに多くのホストマシンが含まれていてもよいことは理解される。各ホストマシンは、さらに複数のGPUと複数のCPUを含み得る。

20

【0124】

ホストマシン(すなわち、ホスト1-A710およびホスト1-B720)は、トップオブラックスイッチ、すなわちTOR1 750を介してネットワークファブリックに通信可能に結合される。このネットワークファブリックは、本明細書ではラック700のフロントエンドネットワークとも呼ばれ、外部ネットワークに対応する場合もある。ホストマシン、すなわちホスト1-A710は、ネットワークインターフェースカード(NIC)730およびTOR1 750に結合されるNVD735(すなわち、ネットワーク仮想化装置)を介してフロントエンドネットワークに接続される。ホストマシン、すなわちホスト1-B720は、ネットワークインターフェースカード(NIC)740およびネットワーク仮想化装置745を介してフロントエンドネットワークに接続され、ネットワーク仮想化装置745はTOR1 750に結合される。ホストマシンホスト1-A710とホスト1-B720は、反対側のQoS対応バックエンドネットワークに接続されている。QoS対応のバックエンドネットワークは、図6に示すようにGPUクラスタネットワークに対応し得る。ホストマシン1-A710は、別のNIC765を介してTOR2スイッチ760に接続され、TOR2スイッチ760はホストマシンをバックエンドネットワークに通信可能に結合する。同様に、ホストマシン1-B720はNIC780を介してTOR2スイッチ760に接続され、TOR2スイッチ760はホストマシンをバックエンドネットワークに通信可能に結合する。

30

40

【0125】

いくつかの実施形態によれば、ホストマシンはネットワークの物理的なトポロジを認識しない、すなわち、特定のホストマシンはネットワーク内の他のホストマシンの物理的な場所/位置を認識しない。例えば、図6を参照すると、ホストマシン1-A612は、ホストマシン1-B614が実際には同じラック(つまり、ラック1 610)に含まれており、同じTORスイッチ(つまり、TOR1 616)の背後に配置されていることを認識していない。しかし、ネットワーク制御プレーンは、ホストマシンの全体的な物理トポロジを認識する。1つの実装では、ネットワーク制御プレーンは、トラフィックの局所性を実現し、unnecessary のトラフィックの輻輳を回避するために、このような局所性情報をホストマシンに公開する。そうすることで、図8Aおよび8Bを参照して以下に示すように

50

、GPUワークロードの性能に大きな影響を与える。

【0126】

いくつかの実施形態では、ネットワーク制御プレーンは、インスタンスメタデータサービス (IMDS) を利用して、メタデータ情報 (例えば、局所性情報) をホストマシンに公開 (および記憶) する。このようなメタデータ情報は、ホストマシンに関連付けられているそれぞれのスマートNICを介してホストマシンに公開され得る。例えば、図7を参照すると、局所性情報はスマートNIC730を介してホスト1-A710に (IMDS経由で) 公開され、一方、局所性情報はスマートNIC740を介してホスト1-B720に (IMDS経由で) 公開される。局所性情報には、ホストマシンのラック情報を示すメタデータが含まれ得ることが理解される。例えば、ラック情報は、特定のホストマシンが属するラックID (例えば、ラック1610) を示すだけでなく、ラックに関連付けられているTORスイッチ (例えば、TOR1スイッチ616) を示す場合もある。各ホストマシンはIMDSにクエリを実行して、ホストマシンに関連付けられているメタデータ情報を取得できることが理解される。以下に説明するように、公開される局所性情報は、より高いGPUワークロードスループットを実現するための最適な論理トポロジを構築するために使用される。

10

【0127】

図8Aを参照すると、特定の実施形態による、局所性情報なしで構築される論理トポロジが示される。図8Aに示されている論理トポロジは、図6の4つのホストマシン、すなわちホストマシン1-A612、ホストマシン1-B614、ホストマシン2-A622、およびホストマシン2-B624に対応している。図6に示すように、ホストマシン1-Aと1-Bはラック1610に属し (つまり、同じTORスイッチTOR1616の後ろに配置されている)、一方、ホストマシン2-Aと2-Bはラック2620に属し (つまり、同じTORスイッチTOR2626の後ろに配置されている) ことに留意されたい。

20

【0128】

図8Aに示す論理トポロジは、局所性情報なしで構築されるリングトポロジであり、すなわち、リングトポロジはランダム (任意) な方法で構築される。図8Aに示すように、リングは、805というラベルの付いたリンクを介して、ホスト1-Aがホスト2-Bに直接接続される (つまり、ホスト1-Aの8つのGPUが、ホスト2-Bの対応する8つのGPUに接続される) ように構築されている。さらに、ホスト2-Bは810というラベルの付いたリンクを介してホスト1-Bに接続され、ホスト1-Bは815というラベルの付いたリンクを介してホスト2-Aに接続され、ホスト2-Aは820というラベルの付いたリンクを介してホスト1-Aに接続される。

30

【0129】

図8Aに示すように構築される論理トポロジでは、ECMPトラフィック分散によりネットワークフローの衝突が発生しやすくなる。ホスト1-A612とホスト1-B614は同じTORスイッチ、つまりTOR1616の背後に配置されているが、ホスト1-Bから発信されホスト1-A宛てのトラフィックは、次のルートを通過することに留意されたい。トラフィックは、最初にホスト1-Bからホスト2-A (つまり、仮想リンク815) にルーティングされ、次にホスト2-Aからホスト1-A (つまり、仮想リンク820) にルーティングされる。したがって、発信元ホスト (ホスト1-B) と同じラック (ラック1) に配置されている宛先ホスト (ホスト1-A) 宛てのトラフィックは、TORスイッチを介してラック外のホストマシン (ホスト2-A) に不必要にルーティングされ、さらにラック内の宛先ホストマシンにリダイレクトされる。トラフィックをホストマシン2-Bからホストマシン2-Aにルーティングする場合にも、同様の状況が発生する。個々のホストマシンは自身の局所性情報を認識していないため、リングトポロジが任意に構築されることに留意されたい。図8Aの任意に構築される論理トポロジにより、フロー衝突の可能性が高まり、GPUワークロードのスループットが低下する (例えば、レイテンシの増加、ジッター損失など)。

40

50

【 0 1 3 0 】

図 8 B は、特定の実施形態による、局所性情報を使用して構築される別の論理トポロジを示す。図 8 B に示す論理トポロジは、図 6 の 4 つのホストマシン、すなわちホストマシン 1 - A 6 1 2、ホストマシン 1 - B 6 1 4、ホストマシン 2 - A 6 2 2、およびホストマシン 2 - B 6 2 4 に対応し、局所性情報に基づいて構築される。ホストマシン 1 - A と 1 - B はラック 1 6 1 0 に属し（つまり、同じ TOR スイッチ TOR 1 6 1 6 の背後に配置されている）、ホストマシン 2 - A と 2 - B はラック 2 6 2 0 に属し（つまり、同じ TOR スイッチ TOR 2 6 2 6 の背後に配置されている）ことに留意されたい。

【 0 1 3 1 】

図 8 B に示す論理トポロジは、局所性情報を用いて構築されるリングトポロジであり、すなわち、リングトポロジは、例えば IMDS から取得される局所性情報に基づいて構築される。いくつかの実施形態では、論理トポロジは、ラックに含まれるホストマシンの 1 つであり得る構成ホストマシンによって構築され得る。図 8 B に示すように、リングは、8 6 5 というラベルの付いたリンクを介して、ホスト 1 - B がホスト 1 - A に直接接続される（つまり、ホスト 1 - B の 8 つの GPU が、ホスト 1 - A の対応する 8 つの GPU に接続される）ように構築されている。さらに、ホスト 1 - A は 8 5 0 というラベルの付いたリンクを介してホスト 2 - B に接続され、ホスト 2 - B は 8 5 5 というラベルの付いたリンクを介してホスト 2 - A に接続され、ホスト 2 - A は 8 6 0 というラベルの付いたリンクを介してホスト 1 - B に接続される。

【 0 1 3 2 】

図 8 B に示す論理トポロジは、同じラック内にある（つまり、同じ TOR スイッチの背後に配置されている）2 台のホストマシン間のトラフィックがラックの外部に不必要にルーティングされることを回避する。例えば、ホスト 1 - B からホスト 1 - A にトラフィックをルーティングする場合を考慮すると、トラフィックはホストマシン 1 - B からホストマシン 1 - A に直接送信できる（仮想リンク 8 6 5 経由）。これは、図 8 A に関して実行されるルーティングとは対照的である。図 8 A では、ホスト 1 - B で発信され、ホスト 1 - A 宛てのトラフィックが最初にホスト 2 - A に向けられ（つまり、TOR スイッチ経由でラックの外部にルーティングされ）、さらに宛先ホストマシン、つまりホスト 1 - A に送信される。このように、局所性情報に基づいて論理トポロジを構築すると、フロー衝突の可能性が低減され、GPU ワークロードのスループットが向上する。

【 0 1 3 3 】

図 9 は、特定の実施形態による、要求をプロビジョニングする際に実行されるステップを示す例示的なフローチャートを示している。図 9 に示されている処理は、それぞれのシステム、ハードウェア、またはそれらの組み合わせの 1 つまたは複数の処理装置（例えば、プロセッサ、コア）によって実行されるソフトウェア（例えば、コード、命令、プログラム）で実装され得る。ソフトウェアは、非一時的な記憶媒体（例えば、メモリ装置）に記憶され得る。図 9 に示され、以下で説明される方法は、例示を目的としており、限定するものではない。図 9 は、特定のシーケンスまたは順序で発生するさまざまな処理ステップを示しているが、これは限定するものではない。特定の代替実施形態では、ステップは異なる順序で実行され得、また、一部のステップは並行して実行され得る。

【 0 1 3 4 】

プロセスはステップ 9 0 5 で開始され、複数のホストマシン（すなわち、クラスタに含まれるホストマシン）の各ホストマシンについて、そのホストマシンの局所性情報が記憶される。ホストマシンの局所性情報には、ホストマシンを含むラックを識別する情報が含まれる。いくつかの実施形態では、インスタンスメタデータサービスを利用して、ホストマシンに関連付けられるネットワーク仮想化装置（NVD）を介してホストマシンに局所性情報を記憶することができる。局所性情報は、ホストマシンを含むラックの識別子（ラック ID）、ラックに関連付けられている TOR スイッチの識別子などを示す情報に対応し得ることに留意されたい。

【 0 1 3 5 】

10

20

30

40

50

ステップ 910 では、制御プレーンは、ワークロードの実行を要求する要求（例えば、顧客からの要求）を受信する。ワークロードは、ホストマシンに関連付けられる GPU を使用して実行される 1 つまたは複数のプロセスに対応することに留意されたい。次に、プロセスはステップ 915 に移動し、複数のホストマシンから 1 つまたは複数のホストマシンがワークロードの実行に使用可能であると識別される。使用可能なホストマシンの識別は、いくつかの方法で実行できる。例えば、一実施形態では、制御プレーンは、ホストマシンによって処理される現在の負荷（つまり、ワークロード）を維持することができる。制御プレーンは、各ホストマシンの容量と、ホストマシンによって処理される現在のワークロードの量に基づいて、顧客のワークロードを実行するために使用可能な 1 つまたは複数のホストマシンを選択できる。さらに、別の実施形態では、制御プレーンは、顧客ごと

10

【0136】

次に、プロセスはステップ 920 に進み、ステップ 915 で識別される 1 つまたは複数のホストマシンの各々について、局所性情報が取得される。（各ホストマシンの）局所性情報は、例えば、対応するホストマシンのインスタンスメタデータサービスによって記憶され得ることに留意されたい。ステップ 925 では、プロセスは 1 つまたは複数のホストマシンのリンク情報を識別する。1 つまたは複数のホストマシンのリンク情報は、1 つまたは複数のホストマシンによって形成される論理トポロジ（例えば、図 8 B に示す論理トポロジ）に対応することが理解される。その後、プロセスはステップ 930 に移動し、要求に応じて 1 つまたは複数のホストマシンの局所性情報とリンク情報が提供される。いくつかの実施形態によれば、顧客は、局所性情報およびリンク情報を取得すると、ワークロードを実行するために 1 つまたは複数のホストマシンのサブセットを選択できる。ホストマシンの選択は、ワークロードに関連付けられる 1 つまたは複数の制約に基づいて実行され得ることに留意されたい。制約に関する詳細については、図 10 および 11 を参照した後で説明する。ホストマシンのサブセットを選択すると、選択したホストマシン上でワークロードが実行され得る。

20

【0137】

図 10 は、特定の実施形態による、顧客の要求を処理する際に実行されるステップを示す例示的なフローチャートを示す。図 10 に示される処理は、それぞれのシステム、ハードウェア、またはそれらの組み合わせの 1 つまたは複数の処理装置（例えば、プロセッサ、コア）によって実行されるソフトウェア（例えば、コード、命令、プログラム）で実装され得る。ソフトウェアは、非一時的な記憶媒体（例えば、メモリ装置）に記憶され得る。図 10 に示され、以下で説明される方法は、例示を目的としており、限定するものではない。図 10 は、特定のシーケンスまたは順序で発生するさまざまな処理ステップを示しているが、これは限定するものではない。特定の代替実施形態では、ステップは一部の異なる順序で実行され得、また、一部のステップは並行して実行され得る。

30

【0138】

プロセスはステップ 1005 で開始され、複数のホストマシン（つまり、クラスタに含まれるホストマシン）の各ホストマシンについて、そのホストマシンの局所性情報が記憶される。ホストマシンの局所性情報には、ホストマシンを含むラックを識別する情報が含まれる。いくつかの実施形態では、インスタンスメタデータサービスを利用して、ホストマシンに関連付けられるネットワーク仮想化装置（NVD）を介してホストマシンに局所性情報を記憶することができる。局所性情報は、ホストマシンを含むラックの識別子（ラック ID）、ラックに関連付けられている TOR スイッチの識別子などを示す情報に対応し得ることに留意されたい。

40

【0139】

ステップ 1010 では、制御プレーンは顧客からの要求を受信し、その要求は顧客に割り当てられる 1 つまたは複数のホストマシンの情報を要求する。ステップ 1015 では、制御プレーンは、顧客に割り当てられる 1 つまたは複数のホストマシンを識別する。いく

50

つかの実施形態では、制御プレーンは、異なる顧客に割り当てられるホストマシンのマッピングを維持することができる。このように、制御プレーンはマッピングを利用して、顧客に割り当てられる1つまたは複数のホストマシンを識別できる。

【0140】

次に、プロセスはステップ1020に進み、ステップ1015で識別される1つまたは複数のホストマシンの局所性情報が取得される。(各ホストマシンの)局所性情報は、例えば、対応するホストマシンのインスタンスメタデータサービスによって記憶され得ることに留意されたい。ステップ1025では、プロセスは1つまたは複数のホストマシンのリンク情報を識別する。1つまたは複数のホストマシンのリンク情報は、1つまたは複数のホストマシンによって形成される論理トポロジ(例えば、図8Bに示す論理トポロジ)に対応することが理解される。その後、プロセスはステップ1030に移動し、ここで、要求に応じて、1つまたは複数のホストマシンの局所性情報とリンク情報が(例えば、顧客に)提供される。前述のように、顧客はステップ1030で提供される情報を利用して、ワークロードに関連付けられる1つまたは複数の制約に従ってワークロードを実行するためのホストマシンを選択できる。

10

【0141】

図11は、特定の実施形態による、顧客のワークロード要求をプロビジョニングする際に行われるステップを示す例示的なフローチャートを示している。図11に示される処理は、それぞれのシステム、ハードウェア、またはそれらの組み合わせの1つまたは複数の処理装置(例えば、プロセッサ、コア)によって実行されるソフトウェア(例えば、コード、命令、プログラム)で実装され得る。ソフトウェアは、非一時的な記憶媒体(例えば、メモリ装置)に記憶され得る。図11に示され、以下で説明される方法は、例示を目的としており、限定するものではない。図11は、特定のシーケンスまたは順序で発生するさまざまな処理ステップを示しているが、これは限定するものではない。特定の代替実施形態では、ステップは異なる順序で実行され得、また、一部のステップは並行して実行され得る。

20

【0142】

プロセスはステップ1105で開始され、複数のホストマシン(つまり、クラスタに含まれるホストマシン)の各ホストマシンについて、そのホストマシンの局所性情報が記憶される。ホストマシンの局所性情報には、ホストマシンを含むラックを識別する情報が含まれる。いくつかの実施形態では、インスタンスメタデータサービスを利用して、ホストマシンに関連付けられるネットワーク仮想化装置(NVD)を介してホストマシンに局所性情報を記憶することができる。局所性情報は、ホストマシンを含むラックの識別子(ラックID)、ラックに関連付けられているTORスイッチの識別子などを示す情報に対応し得ることに留意されたい。

30

【0143】

ステップ1110では、顧客からの要求が受信される。この要求では、顧客がワークロードを実行するために所望するホストマシンの数を要求する。要求には、ワークロードに関連付けられる1つまたは複数の制約に対応するメタデータが含まれ得ることが理解される。いくつかの実施形態では、1つまたは複数の制約には、レイテンシ閾値に関連付けられる第1の制約が含まれ得、つまり、顧客は、レイテンシが所定の閾値よりも低くなるようにワークロードを実行することを所望する可能性がある。第2の制約は、反親和性制約に対応し得る。このような制約は、ホストマシンの可用性をある程度確保したいという顧客の要望に対応する。つまり、ワークロードを実行するために選択されるホストマシンの少なくとも一部は、異なるラックに配置する必要がある。このような制約は通常、ラックの障害の問題に対処するために顧客によって組み込まれる。

40

【0144】

ステップ1115では、ステップ1110で受信した要求から1つまたは複数の制約が抽出される。その後、プロセスはステップ1120に移動し、制御プレーンは、顧客のワークロードを実行するために1つまたは複数のホストマシン(複数のホストマシンから)

50

を割り当てる。ワークロードを実行するために割り当てられるホストマシンは、ステップ 1 1 1 5 で識別される 1 つまたは複数の制約に従って決定できることが理解される。いくつかの実施形態によれば、ステップ 1 1 2 0 のプロセスは、1 つまたは複数の制約およびネットワーククラスタのトポロジ（例えば、ラックに含まれるホストマシンの数、システムで使用可能なラックの数など）に従って、ワークロードを実行するためにホストマシンを割り当てる。各ラックに合計 N 台のホストマシンが含まれており、ユーザが要求するホストマシンの数が K 台であると考慮する。K が N より小さい場合（つまり、 $K < N$ ）、1 つの実装では、K 台のホストマシンすべてを 1 つのラックから割り当てることができる。つまり、K 台のホストマシンすべてを 1 つの TOR スイッチの背後に配置することができる。このような割り当てにより、レイテンシが最小限に抑えられることに留意されたい。しかし、ユーザが特定のレベルの反親和性、つまり可用性を実現することを所望する場合もあり、その場合、ユーザに割り当てられるホストマシンの一部は、ネットワーククラスタの他のラックから選択され得るため、このような割り当てはユーザの承認を待って提供され得ることが理解される。別の例としては、K が N より大きい（ $K > N$ ）場合に対応し得る。つまり、ユーザが要求するホストマシンの数が、ラックに含まれるホストマシンの数よりも多い場合である。この場合、1 つの実装では、最小数の TOR スイッチが使用されるように、すなわちラック内ルーティングが最小化されるように割り当てを実行することができる。例えば、図 6 および 8 B を参照し、制約が最小のレイテンシを達成することである場合を考慮すると、この場合、マシン学習タスクを実行するために 2 台のホストマシンが必要な場合、1 つの実装では、ホスト 1 - A とホスト 1 - B（両方とも図 6 に示すようにラック 1、6 1 0 に存在する）がマシン学習タスクを実行するために割り当てられ得る。このようなタスクの実行では、トラフィックがラック 1 内に局所に封じ込められ、それにより低レイテンシが実現されることに留意されたい。対照的に、別の例では、制約が高可用性の達成である場合、この場合は、ホスト 1 - A（ラック 1 内）とホスト 2 - B（ラック 2 内）がマシン学習タスクの実行用に割り当てられ得る。

10

20

【0 1 4 5】

次に、プロセスはステップ 1 1 2 5 に進み、ステップ 1 1 2 0 で割り当てられる 1 つまたは複数のホストマシンの各々について、局所性情報が取得される。（各ホストマシンの）局所性情報は、対応するホストマシンのインスタンスメタデータサービスによって記憶され得ることに留意されたい。このプロセスでは、さらに、1 つまたは複数のホストマシンのリンク情報を識別する。1 つまたは複数のホストマシンのリンク情報は、1 つまたは複数のホストマシンによって形成される論理トポロジ（例えば、図 8 B に示す論理トポロジ）に対応することが理解される。1 つの実装では、アプリケーションはホストマシンの局所性情報に基づいてホストマシンの論理トポロジを構築でき、各ホストマシンはインスタンスメタデータサービスから局所性情報を取得することに留意されたい。取得される局所性情報および顧客に割り当てられる 1 つまたは複数のホストマシンのリンク情報は、ステップ 1 1 2 5 で顧客に提供される。その後、プロセスはステップ 1 1 3 0 に移動し、顧客のワークロードは顧客に割り当てられる 1 つまたは複数のホストマシン上で実行される。

30

【0 1 4 6】

したがって、本開示の技術は、顧客がワークロードを例えば、ラック内の近隣のホストマシンの近くに「配置」することによって、アプリケーションからサービスへのレイテンシを削減することをプロビジョニングする。さらに、顧客は局所性情報を使用して、より高い反親和性を実現し、そのためより高い回復力を獲得するような方法でワークロードを配置できる。一実施形態によれば、GPU ワークロードを実行するためにホストマシンを割り当てるプロセスで局所性情報（例えば、ラック情報）を利用する機能により、ホストマシンをランダムな方法でユーザ要求に割り当てる、つまりホストマシンの局所性情報を考慮せずまたは顧客に提供しない単純な技術の性能と比較して、GPU ワークロードの性能（例えば、スループット）が 6 倍向上する。

40

【0 1 4 7】

50

図9～11に示されるフローチャートに対する修正は、本開示の範囲内であることが理解される。例えば、1つの実装では、ワークロードを実行するように構成される1つまたは複数のホストマシンの局所性情報が顧客に提供され得る。顧客は、1つまたは複数のホストマシンの局所性情報を利用して、それぞれのユースケースのリンクトポロジを作成できる（例えば、1つまたは複数のホストマシンに関連付けられる親和性、反親和性、またはその他の制約を取得できる）。

【0148】

クラウドインフラストラクチャの実施形態の例

前述のように、*infrastructure as a service (IaaS)* は、クラウドコンピューティングの特定のタイプである。*IaaS* は、パブリックネットワーク（例えば、インターネット）経由で仮想化されるコンピューティングリソースを提供するように構成できる。*IaaS* モデルでは、クラウドコンピューティングプロバイダがインフラストラクチャ構成要素（例えば、サーバ、記憶装置、ネットワークノード（例えば、ハードウェア）、展開ソフトウェア、プラットフォーム仮想化（例えば、ハイパーバイザ層）など）をホストできる。場合によっては、*IaaS* プロバイダが、これらのインフラストラクチャ構成要素に付随するさまざまなサービス（例えば、課金、監視、ログ記録、セキュリティ、負荷分散、クラスタリングなど）も提供することができる。したがって、これらのサービスはポリシー駆動型であり得るため、*IaaS* ユーザは、アプリケーションの可用性と性能を維持するために負荷分散を駆動するポリシーを実装できる可能性がある。

【0149】

場合によっては、*IaaS* 顧客はインターネットなどの広域ネットワーク（WAN）を介してリソースおよびサービスにアクセスし、クラウドプロバイダのサービスを使用してアプリケーションスタックの残りの要素をインストールできる。例えば、ユーザは *IaaS* プラットフォームにログインして仮想マシン（VM）を作成し、各VMにオペレーティングシステム（OS）をインストールし、データベースなどのミドルウェアを展開し、ワークロードとバックアップ用の記憶バケットを作成し、そのVMにエンタープライズソフトウェアをインストールすることもできる。顧客は、次いでプロバイダのサービスを使用して、ネットワークトラフィックの分散、アプリケーションの問題のトラブルシューティング、性能の監視、災害復旧の管理など、さまざまな機能を実行できる。

【0150】

ほとんどの場合、クラウドコンピューティングモデルではクラウドプロバイダの参加が必要になる。クラウドプロバイダは、*IaaS* の提供（例えば、提供、レンタル、販売）を専門とするサードパーティサービスである場合もあるが、そうである必要はない。エンティティはプライベートクラウドを展開し、独自のインフラストラクチャサービスのプロバイダになることを選択することもできる。

【0151】

いくつかの例では、*IaaS* 展開は、新しいアプリケーション、またはアプリケーションの新しいバージョンを準備されるアプリケーションサーバなどに配置するプロセスである。また、サーバの準備プロセス（例えば、ライブラリ、デーモンなどのインストール）も含まれ得る。これは多くの場合、ハイパーバイザ層（例えば、サーバ、記憶装置、ネットワークハードウェア、仮想化）の下でクラウドプロバイダによって管理される。したがって、顧客は、（OS）、ミドルウェア、および/またはアプリケーションの展開（例えば、オンデマンドで起動できるセルフサービス仮想マシン）の処理を担当し得る。

【0152】

いくつかの例では、*IaaS* プロビジョニングは、使用するためにコンピュータまたは仮想ホストを取得し、必要なライブラリまたはサービスをそれらにインストールすることを指し得る。ほとんどの場合、展開にはプロビジョニングは含まれず、最初にプロビジョニングを実行する必要がある。

【0153】

場合によっては、*IaaS* プロビジョニングには2つの異なる課題がある。まず、何か

10

20

30

40

50

を実行する前に、インフラストラクチャの初期セットをプロビジョニングするという最初の課題がある。第2に、すべてがプロビジョニングされた後、既存のインフラストラクチャを進化させる（例えば、新しいサービスの追加、サービスの変更、サービスの削除など）という課題がある。場合によっては、インフラストラクチャの構成を宣言的に定義できるようにすることで、これら2つの課題に対処できることがある。つまり、インフラストラクチャ（例えば、必要な構成要素やそれらの対話方法）は、1つまたは複数の構成ファイルによって定義できる。したがって、インフラストラクチャの全体的なトポロジ（例えば、どのリソースがどのリソースに依存しているか、また、各々のリソースがどのように連携するか）を宣言的に記述できる。場合によっては、トポロジが定義されると、構成ファイルに記述されているさまざまな構成要素を作成および/または管理するワークフローを生成できる。

10

【0154】

いくつかの例では、インフラストラクチャは相互接続される多くの要素を持ち得る。例えば、コアネットワークとしても知られる、1つまたは複数の仮想プライベートクラウド（VPC）（例えば、潜在的なオンデマンドで構成可能なおよび/または共有コンピューティングリソースのプール）が存在し得る。いくつかの例では、ネットワークのセキュリティの設定方法と1つまたは複数の仮想マシン（VM）を定義するためにプロビジョニングされる1つまたは複数のセキュリティグループが存在する場合もある。負荷分散装置、データベースなどの他のインフラストラクチャ要素もプロビジョニングされ得る。より多くのインフラストラクチャ要素が所望されたりおよび/または、追加されたりするにつれて、インフラストラクチャは段階的に進化し得る。

20

【0155】

場合によっては、継続的展開技術を使用して、さまざまな仮想コンピューティング環境にわたってインフラストラクチャコードの展開を可能にすることができる。さらに、説明した技術により、これらの環境内でのインフラストラクチャ管理を可能にすることができる。いくつかの例では、サービスチームは、1つまたは複数の、多くの場合は多数の異なる運用環境（例えば、さまざまな地理的な場所、場合によっては世界中）に展開する必要があるコードを作成できる。しかし、一部の例では、コードが展開されるインフラストラクチャを最初にセットアップする必要がある。場合によっては、プロビジョニングを手動で実行したり、プロビジョニングツールを使用してリソースをプロビジョニングしたり、および/またはインフラストラクチャがプロビジョニングされた展開ツールを使用してコードを展開したりすることもできる。

30

【0156】

図12は、少なくとも1つの実施形態によるIaaSアーキテクチャの例示的なパターンを示すブロック図1200である。サービスオペレータ1202は、仮想クラウドネットワーク（VCN）1206およびセキュアホストサブネット1208を含み得るセキュアホストテナント1204に通信可能に結合できる。いくつかの例では、サービスオペレータ1202は、1つまたは複数のクライアントコンピューティング装置を使用し得る。これらの装置は、ポータブルハンドヘルド装置（例えば、iPhone（登録商標）、携帯電話、iPad（登録商標）、コンピューティングタブレット、パーソナルデジタルアシスタント（PDA））またはウェアラブル装置（例えば、Google Glass（登録商標）ヘッドマウントディスプレイ）であり得、Microsoft Windows Mobile（登録商標）などのソフトウェア、および/またはiOS、Windows Phone、Android、BlackBerry 8、Palm OSなどのさまざまなモバイルオペレーティングシステムを実行し、インターネット、電子メール、ショートメッセージサービス（SMS）、BlackBerry（登録商標）、またはその他の通信プロトコルが有効になっている。あるいは、クライアントコンピューティング装置は、例えば、Microsoft Windows（登録商標）、Apple Macintosh（登録商標）、および/またはLinux（登録商標）オペレーティングシステムのさまざまなバージョンを実行するパーソナルコンピュータおよび/またはラップトップコンピュータを含む、汎用パーソナルコンピュータにすることができる。クライアントコンピューター

40

50

ィング装置は、さまざまな市販のUNIX（登録商標）またはUNIXライクなオペレーティングシステム（Google Chrome OSなどのさまざまなCNU/Linuxオペレーティングシステムを含むがこれに限定されない）のいずれかを実行するワークステーションコンピュータであり得る。代替的に、または追加的に、クライアントコンピューティング装置は、シンクライアントコンピュータ、インターネット対応ゲームシステム（例えば、Kinect（登録商標）ジェスチャ入力装置の有無にかかわらずMicrosoft Xboxゲームコンソール）、および/または個人用メッセージング装置など、VCN1206および/またはインターネットにアクセスできるネットワークを介して通信できるその他の電子装置であり得る。

【0157】

10

VCN1206には、SSH VCN1212に含まれるLPG1210を介してセキュアシェル（SSH）VCN1212に通信可能に結合できる局所ピアリングゲートウェイ（LPG）1210を含めることができる。SSH VCN1212にはSSHサブネットワーク1214が含まれ得、SSH VCN1212は、制御プレーンVCN1216に含まれるLPG1210を介して制御プレーンVCN1216に通信可能に結合できる。また、SSH VCN1212は、LPG1210を介してデータプレーンVCN1218に通信可能に結合できる。制御プレーンVCN1216とデータプレーンVCN1218は、IaaSプロバイダが所有および/または動作できるサービステナント1219に含めることができる。

【0158】

20

制御プレーンVCN1216には、境界ネットワーク（例えば、企業イントラネットと外部ネットワーク間の企業ネットワークの一部）として機能する制御プレーン非武装地帯（DMZ）層1220が含まれることができる。DMZベースのサーバは責任範囲が制限されており、セキュリティ侵害の抑制に役立ち得る。さらに、DMZ層1220には、1つまたは複数の負荷分散装置（LB）サブネットワーク1222、アプリサブネットワーク1226を含むことができる制御プレーンアプリ層1224、データベース（DB）サブネットワーク1230（例えば、フロントエンドDBサブネットワークおよび/またはバックエンドDBサブネットワーク）を含むことができる制御プレーンデータ層1228を含めることができる。制御プレーンDMZ層1220に含まれるLBサブネットワーク1222は、制御プレーンアプリ層1224に含まれるアプリサブネットワーク1226および制御プレーンVCN1216に含まれ得るインターネットゲートウェイ1234と通信可能に結合することができ、アプリサブネットワーク1226は、制御プレーンデータ層1228に含まれるDBサブネットワーク1230およびサービスゲートウェイ1236およびネットワークアドレス変換（NAT）ゲートウェイ1238と通信可能に結合することができる。制御プレーンVCN1216には、サービスゲートウェイ1236とNATゲートウェイ1238を含めることができる。

30

【0159】

制御プレーンVCN1216には、アプリサブネットワーク1226を含むことができるデータプレーンミラーアプリ層1240を含めることができる。データプレーンミラーアプリ層1240に含まれるアプリサブネットワーク1226には、コンピューティングインスタンス1244を実行できる仮想ネットワークインターフェースコントローラ（VNIC）1242を含めることができる。コンピューティングインスタンス1244は、データプレーンミラーアプリ層1240のアプリサブネットワーク1226を、データプレーンアプリ層1246に含めることができるアプリサブネットワーク1226に通信可能に結合できる。

40

【0160】

データプレーンVCN1218には、データプレーンアプリ層1246、データプレーンDMZ層1248、およびデータプレーンデータ層1250が含まれ得る。データプレーンDMZ層1248には、データプレーンアプリ層1246のアプリサブネットワーク1226およびデータプレーンVCN1218のインターネットゲートウェイ1234に通信可能に結合できるLBサブネットワーク1222を含めることができる。アプリサブネットワーク1226は、データプレーンVCN1218のサービスゲートウェイ1236およびデータプレ

50

ーンVCN1218のNATゲートウェイ1238に通信可能に結合できる。データプレーンデータ層1250には、データプレーンアプリ層1246のアプリサブネット1226と通信可能に結合できるDBサブネット1230も含まれ得る。

【0161】

制御プレーンVCN1216およびデータプレーンVCN1218のインターネットゲートウェイ1234は、パブリックインターネット1254に通信可能に結合できるメタデータ管理サービス1252に通信可能に結合できる。パブリックインターネット1254は、制御プレーンVCN1216およびデータプレーンVCN1218のNATゲートウェイ1238に通信可能に結合できる。制御プレーンVCN1216およびデータプレーンVCN1218のサービスゲートウェイ1236は、クラウドサービス1256と通信可能に結合できる。

10

【0162】

いくつかの例では、制御プレーンVCN1216またはデータプレーンVCN1218のサービスゲートウェイ1236は、パブリックインターネット1254を経由せずに、クラウドサービス1256へのアプリケーションプログラミングインターフェース(API)呼び出しを行うことができる。サービスゲートウェイ1236からクラウドサービス1256へのAPI呼び出しは一方方向にすることができる。つまり、サービスゲートウェイ1236はクラウドサービス1256へのAPI呼び出しを行うことができ、クラウドサービス1256は要求されるデータをサービスゲートウェイ1236に送信することができる。しかし、クラウドサービス1256は、サービスゲートウェイ1236へのAPI呼び出しを開始できない場合がある。

20

【0163】

いくつかの例では、セキュアホストテナント1204は、そうでなければ分離され得るサービステナント1219に直接接続することができる。セキュアホストサブネット1208は、LPG1210を介してSSHサブネット1214と通信できる。LPG1210により、分離されるシステム上で双方向通信が可能になる。セキュアホストサブネット1208をSSHサブネット1214に接続すると、セキュアホストサブネット1208はサービステナント1219内の他のエンティティにアクセスできるようになる。

【0164】

制御プレーンVCN1216は、サービステナント1219のユーザが所望するリソースを設定またはプロビジョニングすることを可能にする。制御プレーンVCN1216でプロビジョニングされる必要なリソースは、データプレーンVCN1218に展開されるか、または使用され得る。いくつかの例では、制御プレーンVCN1216はデータプレーンVCN1218から分離され得、制御プレーンVCN1216のデータプレーンミラーアプリ層1240は、データプレーンミラーアプリ層1240およびデータプレーンアプリ層1246に含まれ得るVNIC1242を介して、データプレーンVCN1218のデータプレーンアプリ層1246と通信できる。

30

【0165】

いくつかの例では、システムのユーザ、または顧客は、メタデータ管理サービス1252に要求を通信できるパブリックインターネット1254を介して、作成、読み取り、更新、または削除(CRUD)動作などの要求を行うことができる。メタデータ管理サービス1252は、インターネットゲートウェイ1234を介して制御プレーンVCN1216に要求を通信できる。要求は、制御プレーンDMZ層1220に含まれるLBサブネット1222によって受信され得る。LBサブネット1222は、要求が有効であると決定でき、この決定に応じて、制御プレーンアプリ層1224に含まれるアプリサブネット1226に要求を送信できる。要求が検証され、パブリックインターネット1254への呼び出しが必要な場合、パブリックインターネット1254への呼び出しは、パブリックインターネット1254への呼び出しを行うことができるNATゲートウェイ1238に送信され得る。要求によって記憶が望まれる可能性のあるメモリは、DBサブネット1230に記憶できる。

40

50

【 0 1 6 6 】

いくつかの例では、データプレーンミラーアプリ層 1 2 4 0 は、制御プレーン V C N 1 2 1 6 とデータプレーン V C N 1 2 1 8 との間の直接通信を容易にすることができる。例えば、データプレーン V C N 1 2 1 8 に含まれるリソースに、構成の変更、更新、またはその他の適切な修正を適用することが望まれ得る。制御プレーン V C N 1 2 1 6 は、V N I C 1 2 4 2 を介して、データプレーン V C N 1 2 1 8 に含まれるリソースと直接通信し、それによって、構成の変更、更新、またはその他の適切な修正を実行できる。

【 0 1 6 7 】

いくつかの実施形態では、制御プレーン V C N 1 2 1 6 およびデータプレーン V C N 1 2 1 8 は、サービステナント 1 2 1 9 内に含まれることができる。この場合、システムのユーザまたは顧客は、制御プレーン V C N 1 2 1 6 またはデータプレーン V C N 1 2 1 8 のいずれかを所有または動作することはできない。代わりに、I a a S プロバイダは、制御プレーン V C N 1 2 1 6 とデータプレーン V C N 1 2 1 8 を所有または動作することができ、その両方がサービステナント 1 2 1 9 に含まれ得る。この実施形態は、ユーザまたは顧客が他のユーザまたは他の顧客のリソースと対話することを妨げる可能性のあるネットワークの分離を可能にすることができる。また、この実施形態により、システムのユーザまたは顧客は、記憶のために、望ましいレベルのセキュリティを備えていない可能性のあるパブリックインターネット 1 2 5 4 に依存することなく、データベースを非公開に記憶できるようになる。

【 0 1 6 8 】

他の実施形態では、制御プレーン V C N 1 2 1 6 に含まれる L B サブネット 1 2 2 2 は、サービスゲートウェイ 1 2 3 6 からの信号を受信するように構成することができる。この実施形態では、制御プレーン V C N 1 2 1 6 およびデータプレーン V C N 1 2 1 8 は、パブリックインターネット 1 2 5 4 を呼び出すことなく、I a a S プロバイダの顧客によって呼び出されるように構成できる。I a a S プロバイダの顧客は、顧客が使用するデータベースが I a a S プロバイダによって制御され、パブリックインターネット 1 2 5 4 から分離され得るサービステナント 1 2 1 9 に記憶され得るため、この実施形態を所望し得る。

【 0 1 6 9 】

図 1 3 は、少なくとも 1 つの実施形態による、I a a S アーキテクチャの別の例のパターンを示すブロック図 1 3 0 0 である。サービスオペレータ 1 3 0 2 (例えば、図 1 2 のサービスオペレータ 1 2 0 2) は、仮想クラウドネットワーク (V C N) 1 3 0 6 (例えば、図 1 2 の V C N 1 2 0 6) およびセキュアホストサブネット 1 3 0 8 (例えば、図 1 2 のセキュアホストサブネット 1 2 0 8) を含み得るセキュアホストテナント 1 3 0 4 (例えば、図 1 2 のセキュアホストテナント 1 2 0 4) に通信可能に結合できる。V C N 1 3 0 6 には、S S H V C N 1 3 1 2 に含まれる L P G 1 3 1 0 を介してセキュアシェル (S S H) V C N 1 3 1 2 (例えば、図 1 2 の S S H V C N 1 2 1 2) に通信可能に結合できる局所ピアリングゲートウェイ (L P G) 1 3 1 0 (例えば、図 1 2 の L P G 1 2 1 0) を含めることができる。S S H V C N 1 3 1 2 には、S S H サブネット 1 3 1 4 (例えば、図 1 2 の S S H サブネット 1 2 1 4) が含まれ得、S S H V C N 1 3 1 2 は、制御プレーン V C N 1 3 1 6 に含まれる L P G 1 3 1 0 を介して制御プレーン V C N 1 3 1 6 (例えば、図 1 2 の制御プレーン V C N 1 2 1 6) に通信可能に結合されることができる。制御プレーン V C N 1 3 1 6 は、サービステナント 1 3 1 9 (例えば、図 1 2 のサービステナント 1 2 1 9) に含めることができ、データプレーン V C N 1 3 1 8 (例えば、図 1 2 のデータプレーン V C N 1 2 1 8) は、システムのユーザまたは顧客が所有または動作する可能性のある顧客テナント 1 3 2 1 に含めることができる。

【 0 1 7 0 】

制御プレーン V C N 1 3 1 6 は、L B サブネット 1 3 2 2 (例えば、図 1 2 の L B サブネット 1 2 2 2) を含むことができる制御プレーン D M Z 層 1 3 2 0 (例えば、図 1 2 の制御プレーン D M Z 層 1 2 2 0)、アプリサブネット 1 3 2 6 (例えば、図 1 2 のアプリ

10

20

30

40

50

サブネット 1 2 2 6) を含むことができる制御プレーンアプリ層 1 3 2 4 (例えば、図 1 2 の制御プレーンアプリ層 1 2 2 4)、データベース (DB) サブネット 1 3 3 0 (例えば、図 1 2 の DB サブネット 1 2 3 0 と同様) を含むことができる制御プレーンデータ層 1 3 2 8 (例えば、図 1 2 の制御プレーンデータ層 1 2 2 8) を含むことができる。制御プレーン DMZ 層 1 3 2 0 に含まれる LB サブネット 1 3 2 2 は、制御プレーンアプリ層 1 3 2 4 に含まれるアプリサブネット 1 3 2 6 および制御プレーン VCN 1 3 1 6 に含まれ得るインターネットゲートウェイ 1 3 3 4 (例えば、図 1 2 のインターネットゲートウェイ 1 2 3 4) と通信可能に結合することができ、アプリサブネット 1 3 2 6 は、制御プレーンデータ層 1 3 2 8 に含まれる DB サブネット 1 3 3 0 およびサービスゲートウェイ 1 3 3 6 (例えば、図 1 2 のサービスゲートウェイ) およびネットワークアドレス変換 (NAT) ゲートウェイ 1 3 3 8 (例えば、図 1 2 の NAT ゲートウェイ 1 2 3 8) と通信可能に結合することができる。制御プレーン VCN 1 3 1 6 には、サービスゲートウェイ 1 3 3 6 と NAT ゲートウェイ 1 3 3 8 を含めることができる。

10

【 0 1 7 1 】

制御プレーン VCN 1 3 1 6 には、アプリサブネット 1 3 2 6 を含むことができるデータプレーンミラーアプリ層 1 3 4 0 (例えば、図 1 2 のデータプレーンミラーアプリ層 1 2 4 0) を含めることができる。データプレーンミラーアプリ層 1 3 4 0 に含まれるアプリサブネット 1 3 2 6 には、コンピューティングインスタンス 1 3 4 4 (例えば、図 1 2 のコンピューティングインスタンス 1 2 4 4 と同様) を実行できる仮想ネットワークインターフェースコントローラ (VNIC) 1 3 4 2 (例えば、1 2 4 2 の VNIC) を含めることができる。コンピューティングインスタンス 1 3 4 4 は、データプレーンミラーアプリ層 1 3 4 0 に含まれる VNIC 1 3 4 2 とデータプレーンアプリ層 1 3 4 6 に含まれる VNIC 1 3 4 2 を介して、データプレーンミラーアプリ層 1 3 4 0 のアプリサブネット 1 3 2 6 と、データプレーンアプリ層 1 3 4 6 (例えば、図 1 2 のデータプレーンアプリ層 1 2 4 6) に含まれ得るアプリサブネット 1 3 2 6 との間の通信を容易にすることができる。

20

【 0 1 7 2 】

制御プレーン VCN 1 3 1 6 に含まれるインターネットゲートウェイ 1 3 3 4 は、パブリックインターネット 1 3 5 4 (例えば、図 1 2 のパブリックインターネット 1 2 5 4) に通信可能に結合できるメタデータ管理サービス 1 3 5 2 (例えば、図 1 2 のメタデータ管理サービス 1 2 5 2) に通信可能に結合できる。パブリックインターネット 1 3 5 4 は、制御プレーン VCN 1 3 1 6 に含まれる NAT ゲートウェイ 1 3 3 8 と通信可能に結合できる。制御プレーン VCN 1 3 1 6 に含まれるサービスゲートウェイ 1 3 3 6 は、クラウドサービス 1 3 5 6 (例えば、図 1 2 のクラウドサービス 1 2 5 6) と通信可能に結合できる。

30

【 0 1 7 3 】

いくつかの例では、データプレーン VCN 1 3 1 8 は顧客テナント 1 3 2 1 に含まれることができる。この場合、IaaS プロバイダは、各顧客に制御プレーン VCN 1 3 1 6 を提供することができ、また、IaaS プロバイダは、各顧客に対して、サービステナント 1 3 1 9 に含まれる固有のコンピューティングインスタンス 1 3 4 4 を設定することができる。各コンピューティングインスタンス 1 3 4 4 は、サービステナント 1 3 1 9 に含まれる制御プレーン VCN 1 3 1 6 と、顧客テナント 1 3 2 1 に含まれるデータプレーン VCN 1 3 1 8 との間の通信を可能にすることができる。コンピューティングインスタンス 1 3 4 4 は、サービステナント 1 3 1 9 に含まれる制御プレーン VCN 1 3 1 6 でプロビジョニングされるリソースを、顧客テナント 1 3 2 1 に含まれるデータプレーン VCN 1 3 1 8 に展開したり、その他の方法で使用したりすることを可能にする。

40

【 0 1 7 4 】

他の例では、IaaS プロバイダの顧客は、顧客テナント 1 3 2 1 内に存在するデータベースを有し得る。この例では、制御プレーン VCN 1 3 1 6 には、アプリサブネット 1 3 2 6 を含めることができるデータプレーンミラーアプリ層 1 3 4 0 を含めることができ

50

る。データプレーンミラーアプリ層 1340 はデータプレーン VCN 1318 内に存在できるが、データプレーンミラーアプリ層 1340 はデータプレーン VCN 1318 内に存在できない場合がある。つまり、データプレーンミラーアプリ層 1340 は顧客テナント 1321 にアクセスできる可能性があるが、データプレーンミラーアプリ層 1340 はデータプレーン VCN 1318 に存在しないか、IaaS プロバイダの顧客によって所有または動作されていない可能性がある。データプレーンミラーアプリ層 1340 は、データプレーン VCN 1318 への呼び出しを行うように構成できるが、制御プレーン VCN 1316 に含まれるエンティティへの呼び出しを行うように構成することはできない。顧客は、制御プレーン VCN 1316 にプロビジョニングされるリソースをデータプレーン VCN 1318 に展開したり、その他の方法で使用したりすることを所望する場合があります、
データプレーンミラーアプリ層 1340 は、顧客の所望する展開、またはリソースのその他の使用を容易にすることができる。

10

【0175】

いくつかの実施形態では、IaaS プロバイダの顧客は、データプレーン VCN 1318 にフィルタを適用することができる。この実施形態では、顧客はデータプレーン VCN 1318 がアクセスできるものを決定でき、顧客はデータプレーン VCN 1318 からパブリックインターネット 1354 へのアクセスを制限できる。IaaS プロバイダは、フィルタを適用したり、データプレーン VCN 1318 から外部のネットワークやデータベースへのアクセスを制御できない場合がある。顧客テナント 1321 に含まれるデータプレーン VCN 1318 に顧客によるフィルタと制御を適用すると、データプレーン VCN
1318 を他の顧客およびパブリックインターネット 1354 から分離するのに役立ち得る。

20

【0176】

いくつかの実施形態では、クラウドサービス 1356 は、サービスゲートウェイ 1336 によって呼び出され、パブリックインターネット 1354、制御プレーン VCN 1316、またはデータプレーン VCN 1318 に存在しない可能性があるサービスにアクセスすることができる。クラウドサービス 1356 と制御プレーン VCN 1316 またはデータプレーン VCN 1318 間の接続は、存在または継続的ではない場合がある。クラウドサービス 1356 は、IaaS プロバイダが所有または動作する別のネットワーク上に存在し得る。クラウドサービス 1356 は、サービスゲートウェイ 1336 からの呼び出しを受信するように構成され、パブリックインターネット 1354 からの呼び出しを受信しないように構成され得る。一部のクラウドサービス 1356 は他のクラウドサービス 1356 から分離され得、制御プレーン VCN 1316 は、制御プレーン VCN 1316 と同じリージョンにない場合があるクラウドサービス 1356 から分離され得る。例えば、制御プレーン VCN 1316 は「リージョン 1」に配置され得、クラウドサービス「展開 12」はリージョン 1 と「リージョン 2」に配置され得る。リージョン 1 にある制御プレーン VCN 1316 に含まれるサービスゲートウェイ 1336 によって展開 12 への呼び出しが行われた場合、その呼び出しはリージョン 1 の展開 12 に送信され得る。この例では、制御プレーン VCN 1316、つまりリージョン 1 の展開 12 は、リージョン 2 の展開 12 と通信可能に結合されていないか、または通信していない可能性がある。

30

40

【0177】

図 14 は、少なくとも 1 つの実施形態による、IaaS アーキテクチャの別の例のパターンを示すブロック図 1400 である。サービスオペレータ 1402 (例えば、図 12 のサービスオペレータ 1202) は、仮想クラウドネットワーク (VCN) 1406 (例えば、図 12 の VCN 1206) およびセキュアホストサブネット 1408 (例えば、図 12 のセキュアホストサブネット 1208) を含み得るセキュアホストテナント 1404 (例えば、図 12 のセキュアホストテナント 1204) に通信可能に結合できる。VCN 1406 には、SSH VCN 1412 に含まれる LPG 1410 を介して SSH VCN 1412 (例えば、図 12 の SSH VCN 1212) に通信可能に結合できる LPG 1410 (例えば、図 12 の LPG 1210) を含めることができる。SSH VCN 14

50

12には、SSHサブネット1414（例えば、図12のSSHサブネット1214）が含まれ得、SSH VCN1412は、制御プレーンVCN1416（例えば、図12の制御プレーンVCN1216）に含まれるLPG1410を介して制御プレーンVCN1416に通信可能に結合され得、データプレーンVCN1418（例えば、図12のデータプレーン1218）に含まれるLPG1410を介してデータプレーンVCN1418に通信可能に結合され得る。制御プレーンVCN1416およびデータプレーンVCN1418は、サービステナント1419（例えば、図12のサービステナント1219）に含めることができる。

【0178】

制御プレーンVCN1416には、負荷分散装置（LB）サブネット1422（例えば、図12のLBサブネット1222）を含むことができる制御プレーンDMZ層1420（例えば、図12の制御プレーンDMZ層1220）、アプリサブネット1426（図12のアプリサブネット1226と同様）を含むことができる制御プレーンアプリ層1424（例えば、図12の制御プレーンアプリ層1224）、DBサブネット1430を含むことができる制御プレーンデータ層1428（例えば、図12の制御プレーンデータ層1228）が含まれ得る。制御プレーンDMZ層1420に含まれるLBサブネット1422は、制御プレーンアプリ層1424に含まれるアプリサブネット1426および制御プレーンVCN1416に含まれ得るインターネットゲートウェイ1434（例えば、図12のインターネットゲートウェイ1234）に通信可能に結合することができ、アプリサブネット1426は、制御プレーンデータ層1428に含まれるDBサブネット1430およびサービスゲートウェイ1436（例えば、図12のサービスゲートウェイ）およびネットワークアドレス変換（NAT）ゲートウェイ1438（例えば、図12のNATゲートウェイ1238）に通信可能に結合することができる。制御プレーンVCN1416には、サービスゲートウェイ1436とNATゲートウェイ1438を含めることができる。

【0179】

データプレーンVCN1418は、データプレーンアプリ層1446（例えば、図12のデータプレーンアプリ層1246）、データプレーンDMZ層1448（例えば、図12のデータプレーンDMZ層1248）、およびデータプレーンデータ層1450（例えば、図12のデータプレーンデータ層1250）を含むことができる。データプレーンDMZ層1448には、データプレーンアプリ層1446の信頼できるアプリサブネット1460および信頼できないアプリサブネット1462、およびデータプレーンVCN1418に含まれるインターネットゲートウェイ1434に通信可能に結合できるLBサブネット1422を含めることができる。信頼できるアプリサブネット1460は、データプレーンVCN1418に含まれるサービスゲートウェイ1436、データプレーンVCN1418に含まれるNATゲートウェイ1438、およびデータプレーンデータ層1450に含まれるDBサブネット1430に通信可能に結合できる。信頼できないアプリサブネット1462は、データプレーンVCN1418に含まれるサービスゲートウェイ1436およびデータプレーンデータ層1450に含まれるDBサブネット1430に通信可能に結合できる。データプレーンデータ層1450には、データプレーンVCN1418に含まれるサービスゲートウェイ1436に通信可能に結合できるDBサブネット1430を含めることができる。

【0180】

信頼できないアプリサブネット1462には、テナント仮想マシン（VM）1466（1）-（N）に通信可能に結合できる1つまたは複数のプライマリVNIC1464（1）-（N）を含めることができる。各テナントVM1466（1）-（N）は、それぞれの顧客テナント1470（1）-（N）に含まれ得るそれぞれのコンテナ出口VCN1468（1）-（N）に含まれ得るそれぞれのアプリサブネット1467（1）-（N）に通信可能に結合することができる。それぞれのセカンダリVNIC1472（1）-（N）は、データプレーンVCN1418に含まれる信頼できないアプリサブネット14

10

20

30

40

50

62とコンテナ出口VCN1468(1)-(N)に含まれるアプリサブネット間の通信を容易にすることができる。各コンテナ出口VCN1468(1)-(N)には、パブリックインターネット1454(例えば、図12のパブリックインターネット1254)に通信可能に結合できるNATゲートウェイ1438を含めることができる。

【0181】

制御プレーンVCN1416に含まれ、データプレーンVCN1418に含まれるインターネットゲートウェイ1434は、パブリックインターネット1454に通信可能に結合できるメタデータ管理サービス1452(例えば、図12のメタデータ管理システム1252)に通信可能に結合できる。パブリックインターネット1454は、制御プレーンVCN1416に含まれるNATゲートウェイ1438およびデータプレーンVCN1418に含まれるNATゲートウェイ1438に通信可能に結合できる。制御プレーンVCN1416に含まれ、データプレーンVCN1418に含まれるサービスゲートウェイ1436は、クラウドサービス1456と通信可能に結合できる。

10

【0182】

いくつかの実施形態では、データプレーンVCN1418は顧客テナント1470と統合することができる。この統合は、コード実行時にサポートを所望する場合など、いくつかの場合で、IaaSプロバイダの顧客にとって役立つ、または望ましい場合がある。顧客は、破壊的であったり、他の顧客リソースと通信したり、その他の望ましくない影響を引き起こす可能性のあるコードを実行するように提供し得る。これに応じて、IaaSプロバイダは、顧客からIaaSプロバイダに提供されるコードを実行するかどうかを決定できる。

20

【0183】

いくつかの例では、IaaSプロバイダの顧客は、IaaSプロバイダに一時的なネットワークアクセスを許可し、データプレーン層アプリ1446に接続する機能を要求することができる。関数を実行するコードはVM1466(1)-(N)で実行され得、コードはデータプレーンVCN1418上の他の場所で実行されるように構成されない可能性がある。各VM1466(1)-(N)は1つの顧客テナント1470に接続できる。VM1466(1)-(N)に含まれるそれぞれのコンテナ1471(1)-(N)は、コードを実行するように構成されてもよい。この場合、二重の分離(例えば、コンテナ1471(1)-(N)がコードを実行し、コンテナ1471(1)-(N)が少なくとも信頼できないアプリサブネット1462に含まれるVM1466(1)-(N)に含まれ得る)が存在し得、これにより、誤ったコードまたは望ましくないコードがIaaSプロバイダのネットワークに損害を与えたり、別の顧客のネットワークに損害を与えたりするのを防ぐことに役立つことができる。コンテナ1471(1)-(N)は、顧客テナント1470に通信可能に結合され、顧客テナント1470からデータを送信または受信するように構成され得る。コンテナ1471(1)-(N)は、データプレーンVCN1418内の他のエンティティからデータを送信または受信するように構成されない場合がある。コードの実行が完了すると、IaaSプロバイダはコンテナ1471(1)-(N)を強制終了するか、その他の方法で処分することができる。

30

【0184】

いくつかの実施形態では、信頼できるアプリサブネット1460は、IaaSプロバイダが所有または動作し得るコードを実行し得る。この実施形態では、信頼できるアプリサブネット1460は、DBサブネット1430と通信可能に結合され、DBサブネット1430でCRUD動作を実行するように構成され得る。信頼できないアプリサブネット1462はDBサブネット1430と通信可能に結合され得るが、この実施形態では、信頼できないアプリサブネットはDBサブネット1430で読み取り動作を実行するように構成され得る。各顧客のVM1466(1)-(N)に含まれ、顧客からのコードを実行する可能性のあるコンテナ1471(1)-(N)は、DBサブネット1430と通信可能に結合されていない可能性がある。

40

【0185】

50

他の実施形態では、制御プレーンVCN1416とデータプレーンVCN1418は、通信可能に直接結合されない場合がある。この実施形態では、制御プレーンVCN1416とデータプレーンVCN1418の間に直接通信が行われ得ない場合がある。しかし、少なくとも1つの方法を通じて間接的に通信が行われ得る。IaaSプロバイダによって、制御プレーンVCN1416とデータプレーンVCN1418間の通信を容易にできるLPG1410が確立され得る。別の例では、制御プレーンVCN1416またはデータプレーンVCN1418は、サービスゲートウェイ1436を介してクラウドサービス1456を呼び出すことができる。例えば、制御プレーンVCN1416からクラウドサービス1456への呼び出しには、データプレーンVCN1418と通信できるサービスの要求が含まれ得る。

10

【0186】

図15は、少なくとも1つの実施形態による、IaaSアーキテクチャの別の例のパターンを示すブロック図1500である。サービスオペレータ1502（例えば、図12のサービスオペレータ1202）は、仮想クラウドネットワーク（VCN）1506（例えば、図12のVCN1206）およびセキュアホストサブネット1508（例えば、図12のセキュアホストサブネット1208）を含み得るセキュアホストテナント1504（例えば、図12のセキュアホストテナント1204）に通信可能に結合できる。VCN1506には、SSH VCN1512に含まれるLPG1510を介してSSH VCN1512（例えば、図12のSSH VCN1212）に通信可能に結合できるLPG1510（例えば、図12のLPG1210）を含めることができる。SSH VCN1512には、SSHサブネット1514（例えば、図12のSSHサブネット1214）が含まれ得、SSH VCN1512は、制御プレーンVCN1516（例えば、図12の制御プレーンVCN1216）に含まれるLPG1510を介して制御プレーンVCN1516に通信可能に結合され得、データプレーンVCN1518（例えば、図12のデータプレーン1218）に含まれるLPG1510を介してデータプレーンVCN1518に通信可能に結合され得る。制御プレーンVCN1516およびデータプレーンVCN1518は、サービステナント1519（例えば、図12のサービステナント1219）に含めることができる。

20

【0187】

制御プレーンVCN1516は、LBサブネット1522（例えば、図12のLBサブネット1222）を含むことができる制御プレーンDMZ層1520（例えば、図12の制御プレーンDMZ層1220）、アプリサブネット1526（例えば、図12のアプリサブネット1226）を含むことができる制御プレーンアプリ層1524（例えば、図12の制御プレーンアプリ層1224）、DBサブネット1530（例えば、図14のDBサブネット1430）を含むことができる制御プレーンデータ層1528（例えば、図12の制御プレーンデータ層1228）を含むことができる。制御プレーンDMZ層1520に含まれるLBサブネット1522は、制御プレーンアプリ層1524に含まれるアプリサブネット1526および制御プレーンVCN1516に含まれ得るインターネットゲートウェイ1534（例えば、図12のインターネットゲートウェイ1234）に通信可能に結合することができる。アプリサブネット1526は、制御プレーンデータ層1528に含まれるDBサブネット1530およびサービスゲートウェイ1536（例えば、図12のサービスゲートウェイ）およびネットワークアドレス変換（NAT）ゲートウェイ1538（例えば、図12のNATゲートウェイ1238）に通信可能に結合することができる。制御プレーンVCN1516には、サービスゲートウェイ1536とNATゲートウェイ1538を含めることができる。

30

40

【0188】

データプレーンVCN1518は、データプレーンアプリ層1546（例えば、図12のデータプレーンアプリ層1246）、データプレーンDMZ層1548（例えば、図12のデータプレーンDMZ層1248）、およびデータプレーンデータ層1550（例えば、図12のデータプレーンデータ層1250）を含むことができる。データプレーンD

50

MZ層1548には、データプレーンアプリ層1546の信頼できるアプリサブネット1560（例えば、図14の信頼できるアプリサブネット1460）および信頼できないアプリサブネット1562（例えば、図14の信頼できないアプリサブネット1462）、およびデータプレーンVCN1518に含まれるインターネットゲートウェイ1534に通信可能に結合できるLBサブネット1522を含めることができる。信頼できるアプリサブネット1560は、データプレーンVCN1518に含まれるサービスゲートウェイ1536、データプレーンVCN1518に含まれるNATゲートウェイ1538、およびデータプレーンデータ層1550に含まれるDBサブネット1530に通信可能に結合できる。信頼できないアプリサブネット1562は、データプレーンVCN1518に含まれるサービスゲートウェイ1536およびデータプレーンデータ層1550に含まれるDBサブネット1530に通信可能に結合できる。データプレーンデータ層1550には、データプレーンVCN1518に含まれるサービスゲートウェイ1536に通信可能に結合できるDBサブネット1530を含めることができる。

10

【0189】

信頼できないアプリサブネット1562には、信頼できないアプリサブネット1562内に存在するテナント仮想マシン（VM）1566（1）-（N）と通信可能に結合できるプライマリVNIC1564（1）-（N）を含めることができる。各テナントVM1566（1）-（N）は、それぞれのコンテナ1567（1）-（N）内でコードを実行し、コンテナ出口VCN1568内に含まれ得るデータプレーンアプリ層1546内に含まれ得るアプリサブネット1526に通信可能に結合される。それぞれのセカンダリVNIC1572（1）-（N）は、データプレーンVCN1518に含まれる信頼できないアプリサブネット1562とコンテナ出口VCN1568に含まれるアプリサブネット間の通信を容易にすることができる。コンテナ出口VCNには、パブリックインターネット1554（例えば、図12のパブリックインターネット1254）に通信可能に結合できるNATゲートウェイ1538を含めることができる。

20

【0190】

制御プレーンVCN1516に含まれ、データプレーンVCN1518に含まれるインターネットゲートウェイ1534は、パブリックインターネット1554に通信可能に結合できるメタデータ管理サービス1552（例えば、図12のメタデータ管理システム1252）に通信可能に結合できる。パブリックインターネット1554は、制御プレーンVCN1516に含まれるNATゲートウェイ1538およびデータプレーンVCN1518に含まれるNATゲートウェイ1538に通信可能に結合できる。制御プレーンVCN1516に含まれ、データプレーンVCN1518に含まれるサービスゲートウェイ1536は、クラウドサービス1556と通信可能に結合できる。

30

【0191】

いくつかの例では、図15のブロック図1500のアーキテクチャによって示されるパターンは、図14のブロック図1400のアーキテクチャによって示されるパターンの例外とみなされる場合があり、IaaSプロバイダが顧客と直接通信できない場合（例えば、切断されるリージョン）には、IaaSプロバイダの顧客にとって望ましい場合がある。各顧客のVM1566（1）-（N）に含まれるそれぞれのコンテナ1567（1）-（N）は、顧客がリアルタイムでアクセスできる。コンテナ1567（1）-（N）は、コンテナ出口VCN1568に含まれ得るデータプレーンアプリ層1546のアプリサブネット1526に含まれるそれぞれのセカンダリVNIC1572（1）-（N）を呼び出すように構成できる。セカンダリVNIC1572（1）-（N）は、NATゲートウェイ1538に呼び出しを送信することができ、NATゲートウェイ1538は、呼び出しをパブリックインターネット1554に送信することができる。この例では、顧客がリアルタイムでアクセスできるコンテナ1567（1）-（N）は、制御プレーンVCN1516から分離でき、データプレーンVCN1518に含まれる他のエンティティからも分離できる。また、コンテナ1567（1）-（N）は、他の顧客のリソースから分離され得る。

40

50

【 0 1 9 2 】

他の例では、顧客はコンテナ 1 5 6 7 (1) - (N) を使用してクラウドサービス 1 5 5 6 を呼び出すことができる。この例では、顧客はコンテナ 1 5 6 7 (1) - (N) 内でクラウドサービス 1 5 5 6 からサービスを要求するコードを実行し得る。コンテナ 1 5 6 7 (1) - (N) は、この要求をセカンダリ V N I C 1 5 7 2 (1) - (N) に送信でき、セカンダリ V N I C 1 5 7 2 (1) - (N) は、要求を N A T ゲートウェイに送信でき、N A T ゲートウェイは、要求をパブリックインターネット 1 5 5 4 に送信できる。パブリックインターネット 1 5 5 4 は、インターネットゲートウェイ 1 5 3 4 を介して、制御プレーン V C N 1 5 1 6 に含まれる L B サブネット 1 5 2 2 に要求を送信できる。要求が有効であると決定される場合、L B サブネットは、サービスゲートウェイ 1 5 3 6 を介してクラウドサービス 1 5 5 6 に要求を送信できるアプリサブネット 1 5 2 6 に要求を送信できる。

10

【 0 1 9 3 】

図に示される I a a S アーキテクチャ 1 2 0 0、1 3 0 0、1 4 0 0、1 5 0 0 は、図に示されるもの以外の構成要素を有し得ることに留意されたい。さらに、図に示される実施形態は、本開示の実施形態を組み込むことができるクラウドインフラストラクチャシステムのいくつかの例にすぎない。いくつかの他の実施形態では、I a a S システムは、図に示されているよりも多くの構成要素またはより少ない構成要素を有し、2 つ以上の構成要素を組み合わせたり、構成要素の構成または配置が異なったりし得る。

【 0 1 9 4 】

特定の実施形態では、本明細書で説明する I a a S システムには、セルフサービス、加入ベース、弾力的にスケラブル、信頼性、可用性が高く、安全な方法で顧客に提供されるアプリケーション、ミドルウェア、およびデータベースサービスオフリングのスイートが含まれ得る。このような I a a S システムの一例としては、本譲受人が提供する Oracle Cloud Infrastructure (O C I) が挙げられる。

20

【 0 1 9 5 】

図 1 6 は、さまざまな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 1 6 0 0 を示す。システム 1 6 0 0 は、上記のいずれかのコンピュータシステムを実装するために使用できる。図に示すように、コンピュータシステム 1 6 0 0 には、バスサブシステム 1 6 0 2 を介して多数の周辺サブシステムと通信する処理装置 1 6 0 4 が含まれている。これらの周辺サブシステムには、処理加速装置 1 6 0 6、I / O サブシステム 1 6 0 8、記憶サブシステム 1 6 1 8、および通信サブシステム 1 6 2 4 が含まれ得る。記憶サブシステム 1 6 1 8 には、有形のコンピュータ可読記憶媒体 1 6 2 2 とシステムメモリ 1 6 1 0 が含まれる。

30

【 0 1 9 6 】

バスサブシステム 1 6 0 2 は、コンピュータシステム 1 6 0 0 のさまざまな構成要素とサブシステムが意図したとおりに相互に通信できるようにするメカニズムを提供する。バスサブシステム 1 6 0 2 は単一のバスとして概略的に示されているが、バスサブシステムの代替の実施形態では複数のバスを利用し得る。バスサブシステム 1 6 0 2 は、メモリバスまたはメモリコントローラ、周辺バス、およびさまざまなバスアーキテクチャのいずれかを使用する局所バスを含む、いくつかのタイプのバス構造のいずれかになり得る。例えば、このようなアーキテクチャには、業界標準アーキテクチャ (I S A) バス、マイクロチャネルアーキテクチャ (M C A) バス、拡張 I S A (E I S A) バス、ビデオエレクトロニクス標準協会 (V E S A) 局所バス、および周辺機器相互接続 (P C I) バスが含まれ得、これらは I E E E P 1 3 8 6 . 1 標準に従って製造されるメザニンバスとして実装できる。

40

【 0 1 9 7 】

処理装置 1 6 0 4 は、1 つまたは複数の集積回路 (例えば、従来のマイクロプロセッサまたはマイクロコントローラ) として実装することができ、コンピュータシステム 1 6 0 0 の動作を制御する。処理装置 1 6 0 4 には 1 つまたは複数のプロセッサが含まれ得る。

50

これらのプロセッサには、単一コアプロセッサまたはマルチコアプロセッサが含まれ得る。特定の実施形態では、処理装置 1604 は、各処理装置に単一またはマルチコアプロセッサが含まれる 1 つまたは複数の独立した処理装置 1632 および / または 1634 として実装され得る。他の実施形態では、処理装置 1604 は、2 つのデュアルコアプロセッサを 1 つのチップに統合して形成されるクアッドコア処理装置として実装されることもできる。

【0198】

さまざまな実施形態において、処理装置 1604 は、プログラムコードに応じてさまざまなプログラムを実行することができ、複数の同時実行プログラムまたはプロセスを維持することができる。任意の時点で、実行されるプログラムコードの一部またはすべてが、プロセッサ 1604 および / または記憶サブシステム 1618 に存在し得る。適切なプログラミングにより、プロセッサ 1604 は上記のさまざまな機能を提供できる。コンピュータシステム 1600 には、さらに、デジタル信号プロセッサ (DSP)、専用プロセッサなどを含むことができる処理加速装置 1606 が含まれ得る。

10

【0199】

I/O サブシステム 1608 には、ユーザインターフェース入力装置とユーザインターフェース出力装置が含まれ得る。ユーザインターフェース入力装置には、キーボード、マウスやトラックボールなどのポインティング装置、ディスプレイに組み込まれたタッチパッドまたはタッチスクリーン、スクロールホイール、クリックホイール、ダイヤル、ボタン、スイッチ、キーパッド、音声コマンド認識システムを備えたオーディオ入力装置、マイク、その他のタイプの入力装置が含まれ得る。ユーザインターフェース入力装置には、例えば、Microsoft Kinect (登録商標) モーションセンサーなどのモーションセンシング装置および / またはジェスチャ認識装置が含まれ得る。Microsoft Kinect (登録商標) モーションセンサーを使用すると、ユーザはジェスチャや音声コマンドを使用した自然なユーザインターフェースを通じて、Microsoft Xbox (登録商標) 360 ゲームコントローラーなどの入力装置を制御および入力装置と対話できる。ユーザインターフェース入力装置には、ユーザの目の動き (例えば、写真を撮影しているときおよび / またはメニューを選択しているときの「まばたき」) を検出し、その目のジェスチャを入力装置 (例えば、Google Glass (登録商標)) への入力として変換する Google Glass (登録商標) まばたき検出器などの目のジェスチャ認識装置も含まれ得る。さらに、ユーザインターフェース入力装置には、ユーザが音声コマンドを通じて音声認識システム (例えば、Siri (登録商標) ナビゲータ) と対話できるようにする音声認識センシング装置が含まれ得る。

20

30

【0200】

ユーザインターフェース入力装置には、3次元 (3D) マウス、ジョイスティックまたはポインティングスティック、ゲームパッドおよびグラフィックタブレット、ならびにスピーカ、デジタルカメラ、デジタルビデオカメラ、ポータブルメディアプレーヤ、ウェブカメラ、画像スキャナ、指紋スキャナ、バーコードリーダ 3D スキャナ、3D プリンタ、レーザ距離計、視線追跡装置などのオーディオ / ビジュアル装置も含まれ得るが、これらに限定されない。さらに、ユーザインターフェース入力装置には、例えば、コンピュータ断層撮影、磁気共鳴画像化、位置放射断層撮影、医療用超音波検査装置などの医療用画像化入力装置が含まれ得る。ユーザインターフェース入力装置には、例えば、MIDI キーボード、デジタル楽器などのオーディオ入力装置も含まれ得る。

40

【0201】

ユーザインターフェース出力装置には、ディスプレイサブシステム、インジケータライト、またはオーディオ出力装置などの非視覚ディスプレイなどが含まれ得る。表示サブシステムは、ブラウン管 (CRT)、液晶ディスプレイ (LCD) やプラズマディスプレイを使用するフラットパネル装置、投影装置、タッチスクリーンなどであり得る。一般に、「出力装置」という用語の使用は、コンピュータシステム 1600 からユーザまたは他のコンピュータに情報を出力するためのあらゆるタイプの装置およびメカニズムを含むこと

50

を意図している。例えば、ユーザインターフェース出力装置には、モニタ、プリンタ、スピーカ、ヘッドフォン、自動車ナビゲーションシステム、プロッタ、音声出力装置、モデムなど、テキスト、グラフィック、オーディオ/ビデオ情報を視覚的に伝えるさまざまなディスプレイ装置が含まれ得るが、これらに限定されない。

【0202】

コンピュータシステム1600は、現在システムメモリ1610内に配置されているように示されるソフトウェア要素を含む記憶サブシステム1618を含むことができる。システムメモリ1610には、処理装置1604でロードおよび実行可能なプログラム命令、およびこれらのプログラムの実行中に生成されるデータが記憶され得る。

【0203】

コンピュータシステム1600の構成およびタイプに応じて、システムメモリ1610は揮発性（例えば、ランダムアクセスメモリ（RAM））および/または不揮発性（例えば、読み取り専用メモリ（ROM）、フラッシュメモリなど）であってもよい。RAMには通常、処理装置1604によってすぐにアクセス可能であり、または現在処理装置1604によって動作および実行されているデータおよび/またはプログラムモジュールが含まれている。いくつかの実装では、システムメモリ1610には、静的ランダムアクセスメモリ（SRAM）や動的ランダムアクセスメモリ（DRAM）などの複数の異なるタイプのメモリが含まれ得る。いくつかの実装では、起動時などコンピュータシステム1600内の要素間で情報を転送するのに役立つ基本ルーチンを含む基本入出力システム（BIOS）が、通常、ROMに記憶され得る。限定ではなく例として、システムメモリ1610には、クライアントアプリケーション、Webブラウザ、中間層アプリケーション、リレーショナルデータベース管理システム（RDBMS）などを含み得るアプリケーションプログラム1612、プログラムデータ1614、およびオペレーティングシステム1616も示されている。一例として、オペレーティングシステム1616には、Microsoft Windows（登録商標）、Apple Macintosh（登録商標）、および/またはLinuxオペレーティングシステムのさまざまなバージョン、市販のさまざまなUNIX（登録商標）またはUNIXライクなオペレーティングシステム（さまざまなGNU/Linuxオペレーティングシステム、Google Chrome（登録商標）OSなどを含むがこれらに限定されない）、および/またはiOS、Windows（登録商標）Phone、Android（登録商標）OS、BlackBerry（登録商標）OS、Palm（登録商標）OSオペレーティングシステムなどのモバイルオペレーティングシステムが含まれ得る。

【0204】

記憶サブシステム1618は、いくつかの実施形態の機能を提供する基本的なプログラミングおよびデータ構造を記憶するための有形のコンピュータ可読記憶媒体も提供することができる。プロセッサによって実行されると上記の機能を提供するソフトウェア（プログラム、コードモジュール、命令）は、記憶サブシステム1618に記憶され得る。これらのソフトウェアモジュールまたは命令は、処理装置1604によって実行され得る。記憶サブシステム1618は、本開示に従って使用されるデータを記憶するためのリポジトリも提供することができる。

【0205】

記憶サブシステム1600には、コンピュータ可読記憶媒体1622にさらに接続できるコンピュータ可読記憶媒体リーダ1620も含まれ得る。コンピュータ可読記憶媒体1622は、システムメモリ1610と一緒に、また任意選択で組み合わせて、リモート、局所、固定、および/または取り外し可能な記憶装置に加えて、コンピュータ可読情報を一時的および/またはより永続的に含有、記憶、送信、および取得するための記憶媒体を包括的に表すことができる。非一時的なコンピュータ可読記憶媒体には、揮発性メモリ記憶装置および/または不揮発性記憶装置を含む物理的に有形のメモリまたは記憶装置が含まれ得る。非一時的なコンピュータ可読記憶媒体の例には、磁気記憶媒体（例えば、ディスクやテープ）、光学記憶媒体（例えば、DVD、CD）、さまざまなタイプのRAM、ROM、またはフラッシュメモリ、ハードドライブ、フロッピー（登録商標）ドライブ、

10

20

30

40

50

取り外し可能なメモリドライブ（例えば、USBドライブ）、その他のタイプの記憶装置が含まれる。

【0206】

コードまたはコードの一部を含むコンピュータ可読記憶媒体1622には、情報の記憶および/または送信のための任意の方法または技術で実装される揮発性および不揮発性、取り外し可能および取り外し不可能な媒体など、記憶媒体および通信媒体を含むがこれらに限定されない、当該技術分野で既知または使用されている任意の適切な媒体も含まれ得る。これには、RAM、ROM、電子的に消去可能なプログラマブルROM（EEPROM）、フラッシュメモリまたはその他のメモリ技術、CD-ROM、デジタル多用途ディスク（DVD）、またはその他の光学式記憶装置、磁気カセット、磁気テープ、磁気ディスク記憶装置またはその他の磁気記憶装置、またはその他の有形のコンピュータ可読媒体などの有形のコンピュータ可読記憶媒体が含まれ得る。

10

【0207】

一例として、コンピュータ可読記憶媒体1622には、取り外し不可能な不揮発性磁気媒体から読み取りまたは書き込みを行うハードディスクドライブ、取り外し可能な不揮発性磁気ディスクから読み取りまたは書き込みを行う磁気ディスクドライブ、およびCDROM、DVD、Blu-Ray（登録商標）ディスクなどの取り外し可能な不揮発性光ディスク、またはその他の光媒体から読み取りまたは書き込みを行う光ディスクドライブが含まれ得る。コンピュータ可読記憶媒体1622には、Zip（登録商標）ドライブ、フラッシュメモリカード、ユニバーサルシリアルバス（USB）フラッシュドライブ、セキュアデジタル（SD）カード、DVDディスク、デジタルビデオテープなどが含まれ得るが、これらに限定されない。コンピュータ可読記憶媒体1622には、フラッシュメモリベースのSSD、エンタープライズフラッシュドライブ、ソリッドステートROMなどの不揮発性メモリに基づくソリッドステートドライブ（SSD）、ソリッドステートRAM、動的RAM、静的RAM、DRAMベースのSSD、磁気抵抗RAM（MRAM）SSD、およびDRAMとフラッシュメモリベースのSSDの組み合わせを使用するハイブリッドSSDなど、揮発性メモリに基づくSSDも含まれ得る。ディスクドライブおよびそれに関連するコンピュータ可読媒体は、コンピュータシステム1600のコンピュータ可読命令、データ構造、プログラムモジュール、およびその他のデータの揮発性記憶を提供し得る。

20

30

【0208】

通信サブシステム1624は、他のコンピュータシステムおよびネットワークへのインターフェースを提供する。通信サブシステム1624は、コンピュータシステム1600から他のシステムとの間でデータを送受信するためのインターフェースとして機能する。例えば、通信サブシステム1624は、コンピュータシステム1600がインターネットを介して1つまたは複数の装置に接続できるようにする。いくつかの実施形態では、通信サブシステム1624には、無線音声および/またはデータネットワーク（例えば、携帯電話技術、3G、4G、EDGE（地球規模の進化のためのデータレートの向上）などの高度なデータネットワーク技術、WiFi（IEEE802.11ファミリー標準、またはその他のモバイル通信技術、またはそれらの任意の組み合わせ）、グローバルポジショニングシステム（GPS）受信機構成要素、および/またはその他の構成要素を使用）にアクセスするための無線周波数（RF）トランシーバー構成要素が含まれ得る。いくつかの実施形態では、通信サブシステム1624は、無線インターフェースに加えて、または無線インターフェースの代わりに、有線ネットワーク接続（例えば、イーサネット）を提供できる。

40

【0209】

いくつかの実施形態では、通信サブシステム1624は、コンピュータシステム1600を使用する可能性のある1人または複数のユーザに代わって、構造化データフィールドおよび/または非構造化データフィールド1626、イベントストリーム1628、イベント更新1630などの形式で入力通信を受信することもできる。

50

【0210】

一例として、通信サブシステム1624は、Twitter（登録商標）フィード、Facebook（登録商標）アップデート、Rich Site Summary（RSS）フィードなどのWebフィード、および/または1つまたは複数のサードパーティ情報ソースからのリアルタイム更新など、ソーシャルネットワークおよび/またはその他の通信サービスのユーザからデータフィード1626をリアルタイムで受信するように構成できる。

【0211】

さらに、通信サブシステム1624は、連続データストリームの形式でデータを受信するように構成することもできる。連続データストリームには、明示的な終了がなく、本質的に連続的または無制限である可能性のあるリアルタイムイベントおよび/またはイベント更新1630のイベントストリーム1628が含まれ得る。連続データを生成するアプリケーションの例としては、例えば、センサーデータアプリケーション、金融ティッカー、ネットワーク性能測定ツール（例えば、ネットワーク監視およびトラフィック管理アプリケーション）、クリックストリーム分析ツール、自動車交通監視などが含まれ得る。

10

【0212】

通信サブシステム1624は、構造化データフィードおよび/または非構造化データフィード1626、イベントストリーム1628、イベント更新1630などを、コンピュータシステム1600に結合される1つまたは複数のストリーミングデータソースコンピュータと通信している可能性のある1つまたは複数のデータベースに出力するように構成することもできる。

20

【0213】

コンピュータシステム1600は、ハンドヘルドポータブル装置（例えば、iPhone（登録商標）携帯電話、iPad（登録商標）コンピューティングタブレット、PDA）、ウェアラブル装置（例えば、Google Glass（登録商標）ヘッドマウントディスプレイ）、PC、ワークステーション、メインフレーム、キオスク、サーバラック、またはその他のデータ処理システムを含む、さまざまなタイプのいずれかであり得る。

【0214】

コンピュータとネットワークは常に変化する性質のため、図に示されているコンピュータシステム1600の説明は、特定の例としてのみ意図されている。図に示されているシステムよりも多くの構成要素またはより少ない構成要素を持つ他の多くの構成が可能である。例えば、カスタマイズされるハードウェアも使用され得、および/または特定の要素がハードウェア、ファームウェア、ソフトウェア（タブレットを含む）、またはそれらの組み合わせで実装され得る。さらに、ネットワーク入出力装置などの他のコンピューティング装置への接続も使用できる。本明細書で提供される開示および教示に基づいて、当業者は、さまざまな実施形態を実装するための他の方法および/または手法を理解するであろう。

30

【0215】

特定の実施形態について説明したが、さまざまな修正、変更、代替構成、および同等のものも本開示の範囲内に含まれる。実施形態は、特定のデータ処理環境内での動作に限定されず、複数のデータ処理環境内で自由に動作することができる。さらに、実施形態は特定の一連のトランザクションおよびステップを使用して説明されているが、本開示の範囲は、説明されている一連のトランザクションおよびステップに限定されないことは、当業者には明らかであるはずである。上述の実施形態のさまざまな特徴および態様は、個別にまたは組み合わせで使用することができる。

40

【0216】

さらに、実施形態は特定のハードウェアとソフトウェアの組み合わせを使用して説明されているが、他のハードウェアとソフトウェアの組み合わせも本開示の範囲内にあることを認識すべきである。実施形態は、ハードウェアのみ、ソフトウェアのみ、またはそれらの組み合わせを使用して実装することができる。本明細書で説明するさまざまなプロセスは、同じプロセッサまたは異なるプロセッサの任意の組み合わせで実装できる。したがっ

50

て、構成要素またはモジュールが特定の動作を実行するように構成されていると説明されている場合、そのような構成は、例えば、動作を実行するための電子回路を設計することによって、動作を実行するためのプログラマブル電子回路（例えば、マイクロプロセッサ）をプログラミングすることによって、またはそれらの任意の組み合わせによって実現できる。プロセスは、プロセス間通信の従来技術を含むがこれに限定されないさまざまな技術を使用して通信することができ、異なるプロセスのペアが異なる技術を使用することも、同じプロセスのペアが異なる時間に異なる技術を使用することもできる。

【0217】

したがって、本明細書および図面は、限定的な意味ではなく、例示的な意味としてみなされるべきである。しかし、請求項に規定されるより広い趣旨および範囲から逸脱することなく、追加、控除、削除、およびその他の修正および変更が可能であることは明らかである。したがって、特定の開示実施形態が説明されているが、これらは限定することを意図するものではない。さまざまな変更および同等のものが、以下の請求項の範囲内にある。

10

【0218】

開示される実施形態を説明する文脈（特に以下の請求項の文脈）における用語「a」および「an」および「the」および類似の参照対象の使用は、本明細書で別途記載されていない限り、または文脈によって明らかに矛盾しない限り、単数形と複数形の両方を含むものと解釈される。「含む」、「有する」、「含む」、および「含有する」という用語は、特に明記されていない限り、オープンエンドの用語（つまり、「含むが、これらに限定されない」という意味）として解釈される。「接続」という用語は、何かが介在している場合でも、部分的または全体的に含まれている、接続されている、または結合されていると解釈される。本明細書における値の範囲の記載は、本明細書で特に明記されていない限り、その範囲内に含まれる各個別の値を個別に参照するための簡略な方法として機能することをただ目的としており、各個別の値は、本明細書で個別に記載されるかのように明細書に組み込まれる。本書で説明するすべての方法は、本書で別途記載されている場合、または文脈によって明らかに矛盾している場合を除き、任意の適切な順序で実行できる。本明細書で提供されるあらゆる例、または例示的な言語（例えば、「など」）の使用は、実施形態をよりわかりやすく説明することのみを意図しており、特に断りがない限り、開示の範囲を制限するものではない。明細書中のいかなる文言も、請求されていない要素が開示の実施に必須であることを示すものとして解釈されるべきではない。

20

30

【0219】

「X、Y、またはZのうちの少なくとも1つ」という語句などの離接語は、特に明記されていない限り、一般的に、項目、用語などがX、Y、またはZのいずれか、またはそれらの任意の組み合わせ（例えば、X、Y、および/またはZ）であり得ることを示すために使用される文脈内で理解されることが意図されている。したがって、このような離接語は、一般的に、特定の実施形態がXの少なくとも1つ、Yの少なくとも1つ、またはZの少なくとも1つを各々存在させる必要があることを暗示することを意図しておらず、また暗示すべきではない。

【0220】

本開示の好ましい実施形態は、本開示を実施するために知られている最良のモードを含めて、本明細書に記載されている。これらの好ましい実施形態の変形は、前述の説明を読めば、当業者には明らかになるであろう。当業者であれば、必要に応じてこのような変形を採用することができ、本開示は、本明細書に具体的に記載されるものとは別の方法で実施されてもよい。したがって、本開示には、適用法によって許可される限り、本明細書に添付される請求項に記載される主題のすべての修正および同等物が含まれる。さらに、本明細書に別途記載がない限り、上記要素のあらゆる可能なバリエーションの任意の組み合わせが本開示に包含される。

40

【0221】

本明細書に引用されている出版物、特許出願、特許を含むすべての参考文献は、各参考

50

文献が個別に具体的に参照により組み込まれることが示され、本明細書にその全体が記載された場合と同程度に、参照により本明細書に組み込まれる。前述の明細書では、本開示の態様は、その特定の実施形態を参照して説明されているが、当業者であれば、本開示はそれに限定されないことは理解されるであろう。上述の開示のさまざまな特徴および態様は、個別にまたは組み合わせて使用することができる。さらに、実施形態は、本明細書により広い趣旨および範囲から逸脱することなく、本明細書で説明されているもの以外の任意の数の環境および用途で利用することができる。したがって、明細書および図面は、制限的なものではなく、例示的なものとしてみなされる必要がある。

【 図面 】

【 図 1 】

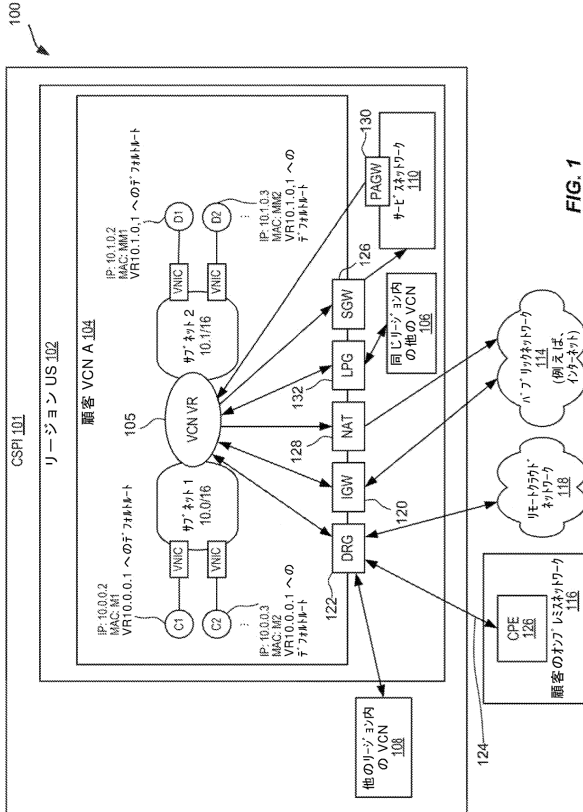


FIG. 1

【 図 2 】

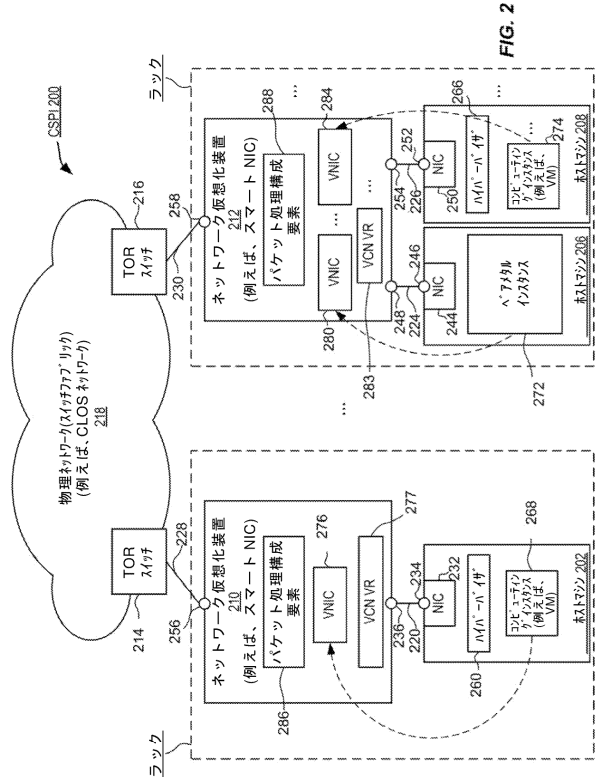


FIG. 2

10

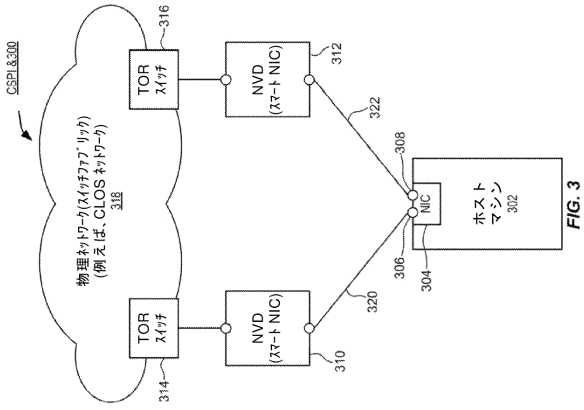
20

30

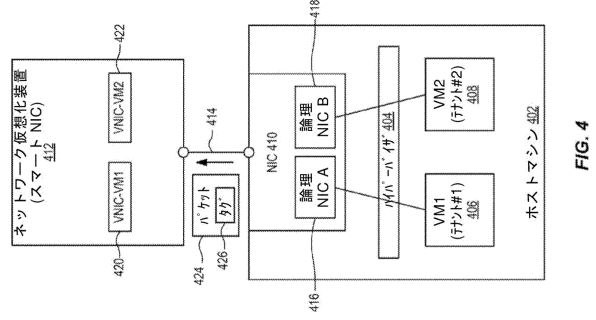
40

50

【 図 3 】

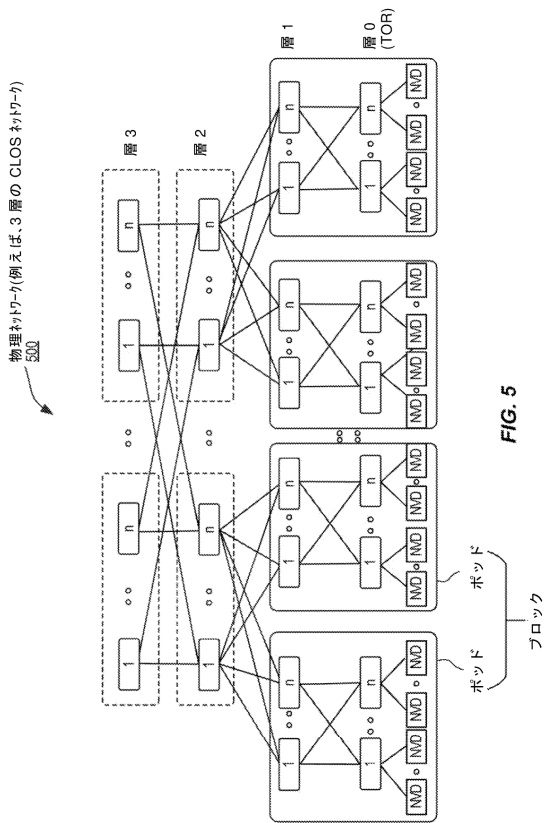


【 図 4 】

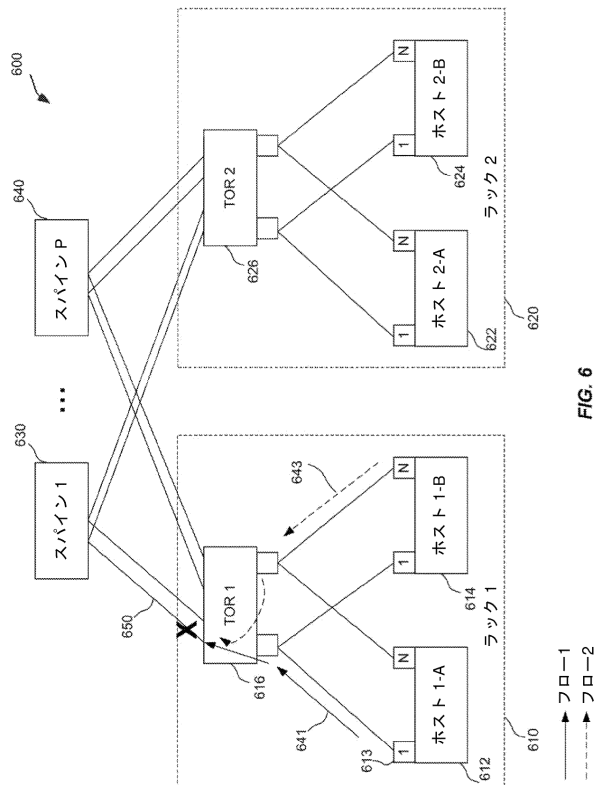


10

【 図 5 】



【 図 6 】



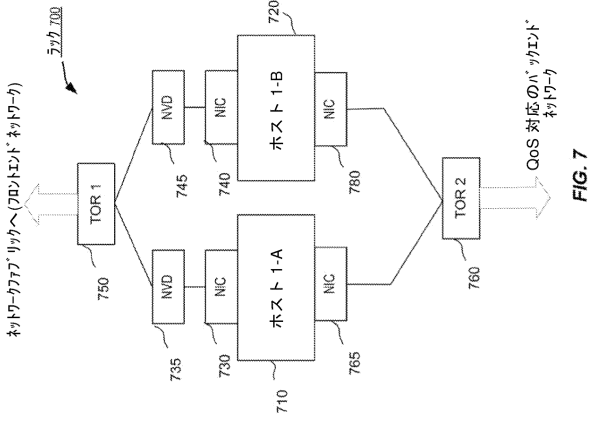
20

30

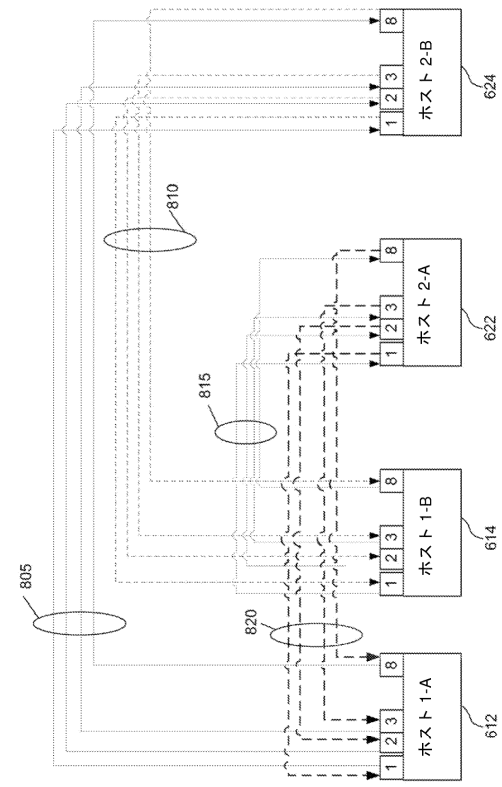
40

50

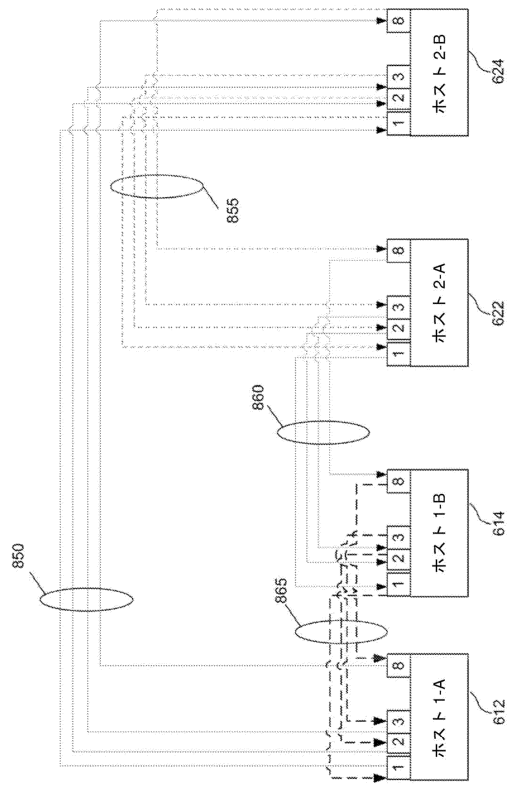
【 図 7 】



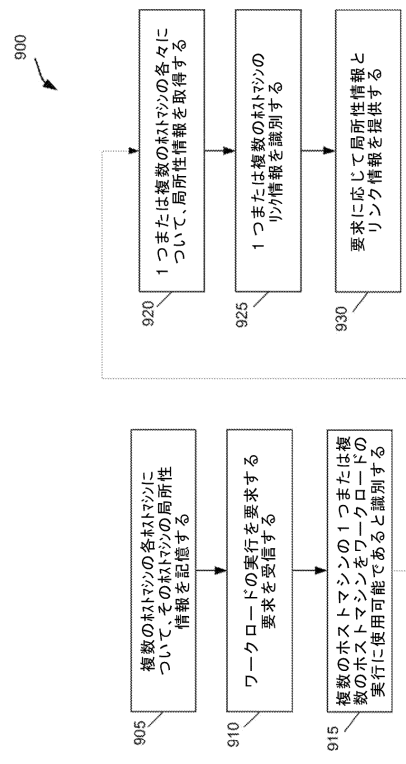
【 図 8 A 】



【 図 8 B 】



【 図 9 】



10

20

30

40

50

【 図 1 0 】

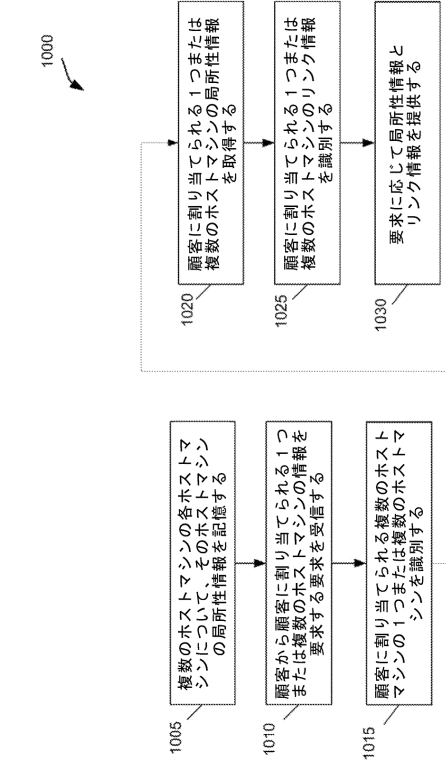


FIG. 10

【 図 1 1 】

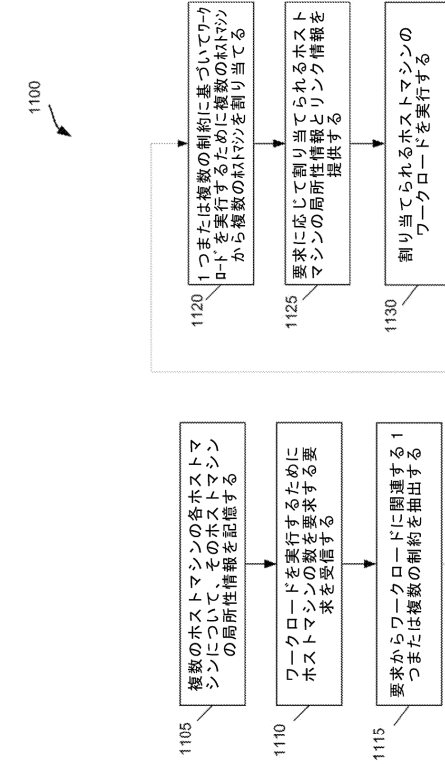


FIG. 11

【 図 1 2 】

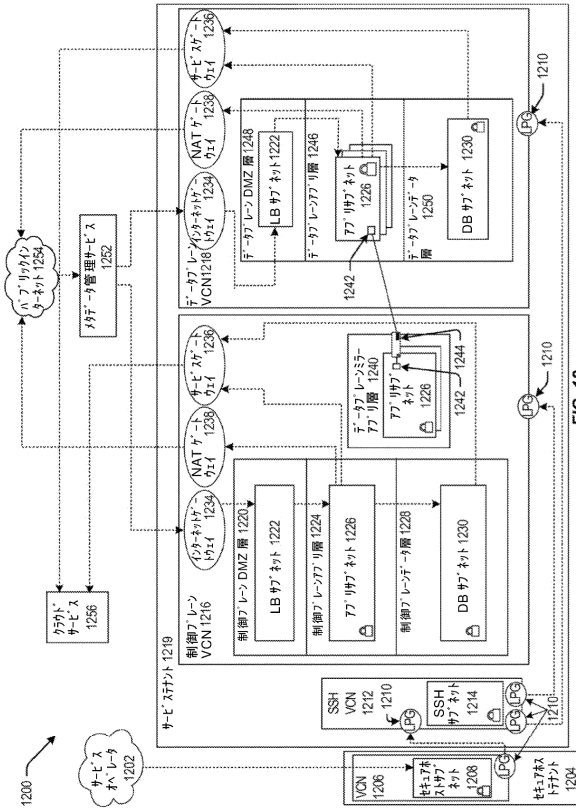


FIG. 12

【 図 1 3 】

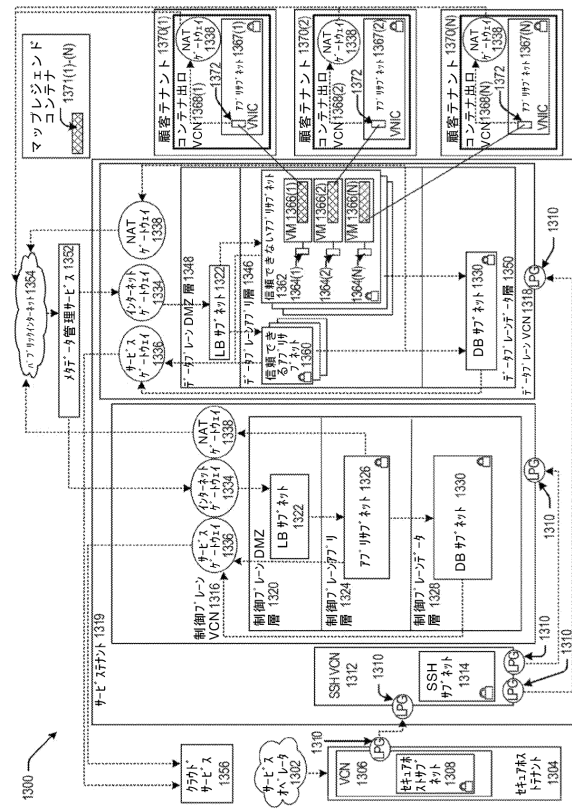


FIG. 13

【図 14】

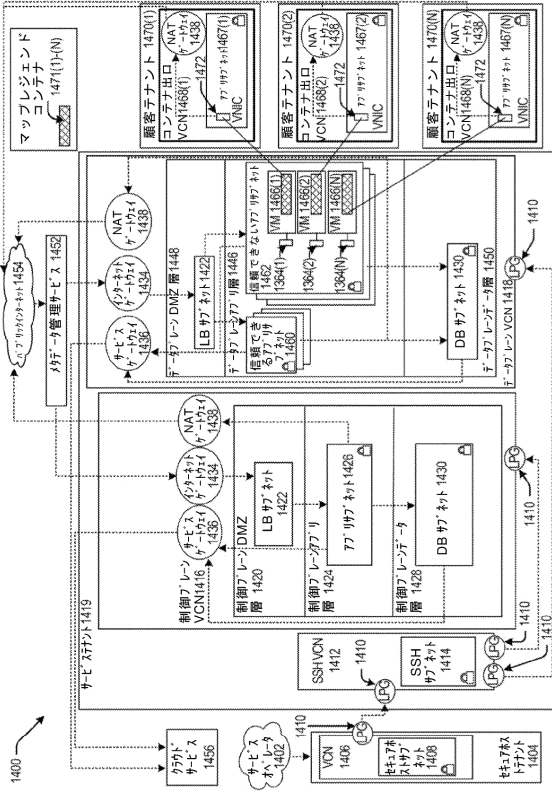


FIG. 14

【図 15】

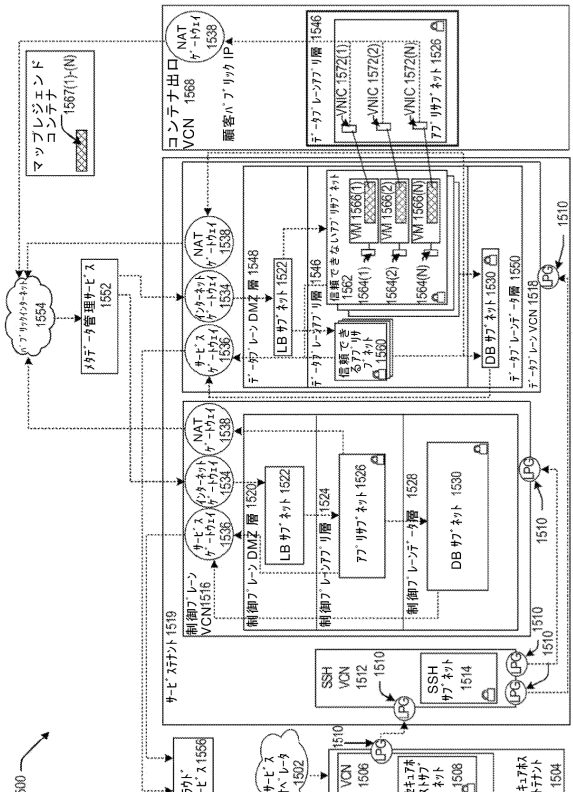


FIG. 15

10

20

【図 16】

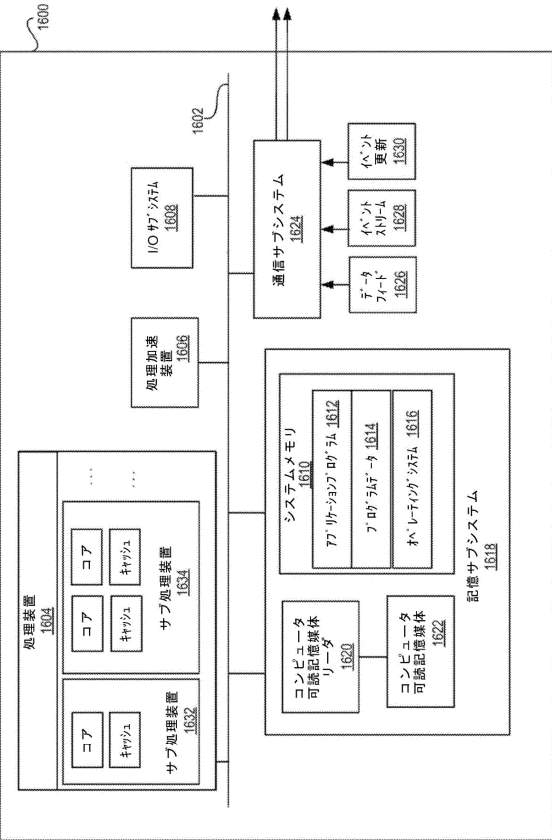


FIG. 16

30

40

50

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No
PCT/US2022/082073

A. CLASSIFICATION OF SUBJECT MATTER INV. G06F9/50 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	US 9 451 393 B1 (CULLEN PATRICK BRIGHAM [US] ET AL) 20 September 2016 (2016-09-20) column 4, line 43 - line 67 column 5, line 4 - line 7 column 9, line 22 - line 67 column 10, line 1 - line 46 column 11, line 5 - line 10 column 13, line 5 - line 12 column 6, line 1 - line 42 column 5, line 63 - line 67 -----	1-20
A	WO 2021/108358 A1 (AMAZON TECH INC [US]) 3 June 2021 (2021-06-03) paragraphs [0047] - [0049], [0098] - [0104], [0112], [0113] -----	1-20
<input type="checkbox"/> Further documents are listed in the continuation of Box C.		<input checked="" type="checkbox"/> See patent family annex.
* Special categories of cited documents :		
"A" document defining the general state of the art which is not considered to be of particular relevance	"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	
"E" earlier application or patent but published on or after the international filing date	"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	
"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)	"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	
"O" document referring to an oral disclosure, use, exhibition or other means	"&" document member of the same patent family	
"P" document published prior to the international filing date but later than the priority date claimed		
Date of the actual completion of the international search 30 March 2023	Date of mailing of the international search report 11/04/2023	
Name and mailing address of the ISA/ European Patent Office, P.B. 5618 Patentlaan 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040, Fax: (+31-70) 340-3016	Authorized officer Buzgan, C	

1

Form PCT/ISA/210 (second sheet) (April 2005)

10

20

30

40

50

INTERNATIONAL SEARCH REPORT

Information on patent family members

International application No
PCT/US2022/082073

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
US 9451393	B1	20-09-2016	NONE

WO 2021108358	A1	03-06-2021	CN 114902183 A
			12-08-2022
			EP 4049139 A1
			31-08-2022
			US 10979534 B1
			13-04-2021
			WO 2021108358 A1
			03-06-2021

10

20

30

40

50

フロントページの続き

,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,D
E,DK,EE,ES,FI,FR,GB,GR,HR,HU,IE,IS,IT,LT,LU,LV,MC,ME,MK,MT,NL,NO,PL,PT,RO,RS,SE,SI,SK,S
M,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,CV,CZ,DE,DJ,DK,DM,DO,DZ,EC,EE,EG,
ES,FI,GB,GD,GE,GH,GM,GT,HN,HR,HU,ID,IL,IN,IQ,IR,IS,IT,JM,JO,JP,KE,KG,KH,KN,KP,KR,KW,KZ,L
A,LC,LK,LR,LS,LU,LY,MA,MD,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,ST,SV,SY,TH,TJ,TM,TN,TR,TT,TZ,UA,UG,US,UZ,VC,V
N,WS,ZA,ZM,ZW

(特許庁注：以下のものは登録商標)

1 . W I N D O W S P H O N E

2 . i O S

(72)発明者 ベッカー , デイビッド

アメリカ合衆国、 9 4 0 6 5 カリフォルニア州、 レッドウッド・シティー、 オラクル・パークウ
エイ、 5 0 0、 オラクル・インターナショナル・コーポレーション内

(72)発明者 コクマード , ハロルド・ドネル

アメリカ合衆国、 9 4 0 6 5 カリフォルニア州、 レッドウッド・シティー、 オラクル・パークウ
エイ、 5 0 0、 オラクル・インターナショナル・コーポレーション内