

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号  
特許第6144346号  
(P6144346)

(45) 発行日 平成29年6月7日 (2017.6.7)

(24) 登録日 平成29年5月19日 (2017.5.19)

(51) Int. Cl.

F I

G O 6 F 9/50 (2006.01)

G O 6 F 9/46 (2006.01)

G O 6 F 9/46 4 6 2 Z

G O 6 F 9/46 3 5 0

請求項の数 13 (全 24 頁)

(21) 出願番号	特願2015-528707 (P2015-528707)	(73) 特許権者	506223509
(86) (22) 出願日	平成25年8月23日 (2013.8.23)		アマゾン・テクノロジーズ、インコーポレイテッド
(65) 公表番号	特表2015-529918 (P2015-529918A)		アメリカ合衆国、98108-1226
(43) 公表日	平成27年10月8日 (2015.10.8)		ワシントン州 シアトル ビーオー ボックス 81226
(86) 国際出願番号	PCT/US2013/056524	(74) 代理人	100108855
(87) 国際公開番号	W02014/032031		弁理士 蔵田 昌俊
(87) 国際公開日	平成26年2月27日 (2014.2.27)	(74) 代理人	100103034
審査請求日	平成27年4月23日 (2015.4.23)		弁理士 野河 信久
(31) 優先権主張番号	13/593,226	(74) 代理人	100075672
(32) 優先日	平成24年8月23日 (2012.8.23)		弁理士 峰 隆司
(33) 優先権主張国	米国 (US)	(74) 代理人	100153051
			弁理士 河野 直樹

最終頁に続く

(54) 【発明の名称】 仮想計算機インスタンスのスケーリング

(57) 【特許請求の範囲】

【請求項 1】

コンピュータ実装方法であって、  
実行可能な命令を有して構成される、1つ以上の計算システムの制御下で、  
配置サービスによって、ユーザから、仮想計算機がスケーラブルであるように要求するアプリケーションプログラミングインターフェース (API) 要求を受信することと、ここで、前記 API 要求は、要求された前記仮想計算機を設定するために使用される仮想計算機イメージを規定し、

前記配置サービスによって、要求された前記仮想計算機がスケーラブルであること、および、サービスプロバイダ環境で動作しているホスト計算装置が、規定された前記仮想計算機イメージのためのリソース容量の増加に対応することができるように、要求された前記仮想計算機に追加の計算リソースを追加する容量を含むこと、を決定することと、

前記配置サービスによって、要求された前記仮想計算機がスケーラブルである場合、規定された前記仮想計算機イメージにおける増加した作業負荷に対処するのに十分な容量を有する、または、要求された前記仮想計算機がスケーラブルではない場合、規定された前記仮想計算機イメージの過剰な容量がほとんどまたは全くない、要求された前記仮想計算機を前記ホスト計算装置上に設定することと、

前記サービスプロバイダ環境中で動作しているスケーリングサービスから、仮想計算機に割り当てられた計算リソースを調節するという第1の要求を受信することと、

前記第1の要求を受信したことに応答して、前記仮想計算機への1つ以上の計算リソー

スの割り当てを調節することと、  
を含む、コンピュータ実装方法。

【請求項 2】

前記スケーリングサービスは、

前記仮想計算機を動作させることと関連した 1 つ以上のメトリクスを監視することと

、

前記 1 つ以上のメトリクスが少なくとも 1 つの指定閾値を超えていることを検出することと、

前記仮想計算機に割り当てられた計算リソースを調節するという第 2 の要求を、前記  
ホスト計算装置に伝送することと、

10

をさらに実施する、請求項 1 のコンピュータ実装方法。

【請求項 3】

前記ホスト計算装置上で実行している監視サービスは、

前記仮想計算機と関連した 1 つ以上のメトリクスを監視する、

前記 1 つ以上のメトリクスが閾値を超えていることを検出する、および

前記 1 つ以上の計算リソースの割り当てを調節する、

ように構成される、請求項 1 のコンピュータ実装方法。

【請求項 4】

前記スケーリングサービスは、

顧客から、前記仮想計算機に割り当てられた前記計算リソースを調節するという第 2  
の要求を受信することと、

20

前記仮想計算機に割り当てられた前記計算リソースを調節するという前記第 2 の要求  
を伝送することと、

をさらに実施する、請求項 1 のコンピュータ実装方法。

【請求項 5】

前記計算リソースの前記割り当ての前記調節をすることは、

利用可能な計算リソースが所定の限界に達するまで、前記仮想計算機に前記 1 つ以上  
の計算リソースを割り当て続けることと、

1 つ以上の追加の仮想計算機に前記仮想計算機の作業負荷を分配するように、前記 1  
つ以上の追加の仮想計算機を設定することと、

30

を含む、請求項 1 のコンピュータ実装方法。

【請求項 6】

前記仮想計算機は、作業負荷を実行することと関連した 1 つ以上のメトリクスを報告す  
る、ゲストエージェントをさらに含む、請求項 1 のコンピュータ実装方法。

【請求項 7】

前記スケーリングサービスからの前記第 1 の要求に応答して、前記仮想計算機に割り当  
てられた前記ホスト計算装置の前記 1 つ以上の計算リソースに少なくとも部分的に基づき  
、前記仮想計算機のユーザに料金を請求すること

をさらに含む、請求項 1 のコンピュータ実装方法。

40

【請求項 8】

顧客が、料金に対して、前記仮想計算機に割り当てられる前記ホスト計算装置の計算リ  
ソースを取得することを可能にする、電子市場を提供することをさらに含み、前記料金は  
、取得された計算リソースの需要と供給に少なくとも部分的に基づく、  
請求項 1 のコンピュータ実装方法。

【請求項 9】

計算システムであって、

少なくとも 1 つのプロセッサと、

前記プロセッサによって実行されるとき、前記計算システムに、

配置サービスによって、ユーザから、仮想計算機がスケーラブルであるように要求す  
るアプリケーションプログラミングインターフェース (API) 要求を受信させ、ここで

50

、前記API要求は、要求された前記仮想計算機を設定するために使用される仮想計算機イメージを規定し、

前記配置サービスによって、要求された前記仮想計算機がスケーラブルであること、および、サービスプロバイダ環境で動作しているホスト計算装置が、規定された前記仮想計算機イメージのためのリソース容量の増加に対応することができるように、要求された前記仮想計算機に追加の計算リソースを追加する容量を含むこと、を決定することと、

前記配置サービスによって、サービスプロバイダ環境で動作している前記ホスト計算装置上に、要求された前記仮想計算機がスケーラブルである場合、規定された前記仮想計算機イメージにおける増加した作業負荷に対処するのに十分な容量を有する、または、要求された前記仮想計算機がスケーラブルではない場合、規定された前記仮想計算機イメージの過剰な容量がほとんどまたは全くない、要求された前記仮想計算機を設定させ、

10

動作している仮想計算機に割り当てられた計算リソースを調節することに関するウェブサービス要求を受信させ、かつ

前記ウェブサービス要求を受信したことに応答して、1つ以上の計算リソースを前記仮想計算機に割り当てさせる、  
命令を含む、メモリと、  
を備える、計算システム。

【請求項10】

前記ウェブサービス要求は、前記仮想計算機に割り当てられた前記計算リソースが調節される条件の1セットを指定する、1つ以上のメトリクスを含む、請求項9に記載の計算システム。

20

【請求項11】

前記仮想計算機に前記1つ以上の計算リソースを割り当てるための前記命令は、更に、前記計算システムに、利用可能な計算リソースが所定の限界に達するまで、1つ以上の追加の計算リソースを前記仮想計算機に割り当て続けさせ、かつ

前記メモリは、さらに、

前記プロセッサによって実行されるとき、前記計算システムに、前記仮想計算機と1つ以上の追加の仮想計算機に前記仮想計算機の作業負荷を分配するように、前記1つ以上の追加の仮想計算機を設定させる、

命令を含む、

30

請求項9に記載の計算システム。

【請求項12】

前記メモリは、

前記プロセッサによって実行されるとき、前記計算システムに、

ユーザが、前記ウェブサービス要求に応答して、前記計算リソースを割り当てることによってスケーリングされることが可能である前記仮想計算機の種類を選択したことを決定させ、かつ

前記ユーザによって選択された前記仮想計算機の前記種類に少なくとも部分的に基づき、前記ユーザに料金を請求させる、

命令をさらに含む、請求項9に記載の計算システム。

40

【請求項13】

前記メモリは、

前記プロセッサによって実行されるとき、前記計算システムに、

顧客が、前記計算システム上の前記1つ以上の計算リソースの需要と供給に少なくとも部分的に基づき、前記仮想計算機のために前記計算システムの1つ以上の追加の計算リソースを購入することを可能にする、電子市場を提供させる

命令をさらに含む、請求項9に記載の計算システム。

【発明の詳細な説明】

【背景技術】

【0001】

50

ますます多くのアプリケーションおよびサービスが、インターネットなどのネットワーク上で利用可能となっているため、ますます多くのコンテンツ、アプリケーション、および/またはサービスの提供者が、クラウドコンピューティングなどの技術に取り掛かっている。クラウドコンピューティングは概して、ウェブサービスなどのサービスを通して電子リソースへのアクセスを提供するためのアプローチであり、この場合、それらのサービスをサポートするために使用されるハードウェアおよび/またはソフトウェアは、任意の所与の時間におけるサービスの要求を満たすために、動的にスケーラブルである。ユーザまたは顧客は典型的に、クラウドを通したリソースへのアクセスをレンタルする、リースする、または別様にその代金を支払い、したがって、必要とされるハードウェアおよび/またはソフトウェアを購入および維持する必要はない。

10

#### 【0002】

この背景では、多くのクラウド計算プロバイダは、複数のユーザが基礎的ハードウェアおよび/またはソフトウェアリソースを共有することを可能にするために、仮想化を利用する。仮想化は、計算サーバ、記憶装置、または他のリソースが、特定のユーザに関連した（例えば、特定のユーザによって所有される）複数の分離したインスタンスに分割されることを可能にすることができる。クラウド計算プロバイダは通常、その顧客のそれぞれに1つ以上の仮想計算機を割り当て、仮想計算機は、それらの顧客に対するアプリケーションおよび/または他の作業負荷を実行するために使用される。しかしながら、要求の増加または他の理由により、顧客の処理負荷が仮想計算機の容量を超え始めるとき、多くの問題および不便が生じる可能性がある。

20

#### 【図面の簡単な説明】

#### 【0003】

本開示に従った種々の実施形態が、図面を参照して記載される。

【図1】種々の実施形態に従って、追加のCPUを割り当てることによって、仮想計算機インスタンスをスケールアップする一実施例を示す。

【図2】種々の実施形態に従って、サービスプロバイダによって展開される自動スケーリングサービスの一実施例を示す。

【図3A】種々の実施形態に従って、ホスト計算機上の仮想計算機インスタンスを自動的にスケーリングするための例示的な過程を示す。

【図3B】種々の実施形態に従って、ユーザからの要求を受信したことに応答して、仮想計算機インスタンスをスケーリングする例示的な過程を示す。

30

【図4】種々の実施形態に従って、仮想計算機インスタンスを自動的にスケーリングし、追加の仮想計算機インスタンスを割り当てるための例示的な過程を示す。

【図5】種々の実施形態に従って利用され得る、例示的な計算装置の一般的なコンポーネントのセットの論理配列を示す。

【図6】種々の実施形態に従った態様を実装するための環境の一実施例を示す。

#### 【発明を実施するための形態】

#### 【0004】

以下の記載において、種々の実施形態は、添付の図面の図の制限のためではなく、一例として示される。本開示における種々の実施形態の参照は、必ずしも同一の実施形態に対してではなく、本開示における種々の実施形態の参照は、少なくとも1つを意味する。特定の実装および他の詳細が考察されるが、これは例示目的でのみ行われることを理解されたい。当業者は、特許請求の範囲に記載された主題の範囲および精神から逸脱することなく、他のコンポーネントおよび構成が使用され得ることを認識するであろう。

40

#### 【0005】

本開示の種々の実施形態に従ったシステムおよび方法は、スケーリング計算リソースのための従来のアプローチにおいて経験される上記または他の欠陥のうちの1つ以上を克服し得る。具体的に、種々の実施形態は、自動的に、仮想計算機インスタンスに追加の計算リソース（例えば、プロセッサ、メモリ、ホームネットワーキングデバイスなど）を割り当てる、および/または種々のユーザ指定閾値またはユーザ要求に従って、仮想計算機イ

50

ンスタンスからの計算リソースの割り当てを解除するためのアプローチを提供する。効果的に、これは、仮想計算機インスタンスが、仮想計算機が提供するリソースに対してオンデマンドでまたは実際の要求に従って、サイズおよび容量を「拡大」または「縮小」することを可能にする。

#### 【0006】

種々の実施形態に従って、1つのアプローチは、その顧客の代わりにアプリケーションおよび仮想計算機インスタンスをホストする、共有された計算リソース環境（例えば、「クラウド」計算プロバイダ）のサービスプロバイダによって実装されてもよい。アプリケーションおよび仮想計算機インスタンスは、ビспロバイダによって所有および運営される物理的なリソース（例えば、ホストサーバおよび他のネットワークリソース）上でホストされる。一実施形態に従って、サービスプロバイダは、顧客から仮想計算機イメージを受信し、仮想計算機イメージに少なくとも部分的に基づき、顧客に対する1つ以上の仮想計算機インスタンスを設定する。これらの仮想計算機インスタンスは次いで、サービスプロバイダの物理的な計算リソースを使用して、顧客の種々のアプリケーションおよび/または他のサービスを実行することができる。

#### 【0007】

一実施形態に従って、各仮想計算機インスタンスは、ホスト計算機（例えば、計算装置）上に設定される。各ホスト計算機は、1つ以上の仮想計算機インスタンスをホストすることができる。少なくとも1つの実施形態において、ホスト計算機はさらに、ホスト計算機のハードウェアデバイスドライバへのアクセスを提供し、かつ1つ以上の仮想計算機インスタンスが直接または仮想化抽象化を通してデバイスにアクセスすることを可能にする、ハイパーバイザおよびサービスホスティング層を含む。

#### 【0008】

一実施形態に従って、いったん仮想計算機インスタンスがホスト計算機上に設定されると、サービスプロバイダは、顧客から（例えば、アプリケーションプログラミングインターフェース（API）を介して）、追加のリソースを仮想計算機インスタンスに割り当てるといふ、またはインスタンスからのリソースの割り当てを解除するといふ要求を受信してもよい。さらに、APIは、顧客が、CPU利用などの基礎的なリソースの種々の操作メトリクスに関連する、仮想計算機インスタンスに対する1つ以上の顧客によって定義された閾値を指定することを可能にし得る。加えて、顧客は、仮想計算機インスタンスをスケールアップまたはスケールダウンするための決定に関連し得る、顧客のサービスまたはアプリケーションに関連した、種々のランタイム操作メトリクスを指定することが可能になる。これらの操作メトリクスおよび閾値は、顧客が、仮想計算機インスタンスに割り当てられたリソースがスケールアップまたはダウンされるべきである条件を指示することを可能にする。

#### 【0009】

一実施形態に従って、サービスプロバイダのシステム上のサービスは、仮想計算機インスタンスの実行中に操作メトリクスを監視する。同一または代替の実施形態において、サービスは、仮想計算機インスタンス内で実行しているゲストエージェントから、操作メトリクスを受信してもよい。結果として、操作メトリクスは、サーバから、および/または仮想計算機インスタンス内から生成されてもよい。サービスが、メトリクスのうちの1つ以上が、所定の期間にわたって顧客によって定義された閾値などの閾値を超えていることを検出する場合、それは、リソース（例えば、中央処理装置（CPU）、メモリ、他のハードウェアデバイス）を追加または除去することによって、仮想計算機インスタンスのスケールアップまたはダウンを開始してもよい。例えば、サービスが、この1時間にわたって少なくとも10秒間、90%を超えるCPU容量で動作していたことを検出する場合、それは、追加のCPU容量を仮想計算機インスタンスに割り当ててもよい（例えば、追加のCPUまたはCPUコアを割り当てる、より強力なCPUに切り替えるなど）。別の例として、サービスが、仮想計算機インスタンスが特定の期間にわたって10%未満のCPU容量で動作していたことを検出する場合、それは、仮想計算機インスタンスから、ある量のCPU容量の割り当てを解除（例えば、低減）してもよい。いくつかの実施形態

10

20

30

40

50

では、仮想計算機インスタンスのスケーリングは、顧客側の任意の手動の関与を必要とせず、自動的に実施され得る。他の実施形態では、仮想計算機インスタンスのスケーリングは、ユーザ（例えば、仮想計算機の所有者）からのインスタンスをスケーリングするという要求を受信したことに応答して、実施され得る。

#### 【0010】

一実施形態に従って、仮想計算機インスタンスは、単一の仮想計算機インスタンスがもはや顧客の作業負荷を適切にサポートすることができなくなるまで、自動的にスケールアップされ得る。いったんこの限界が達せられると、サービスは、作業負荷に対応するために、追加の仮想計算機インスタンスを自動的に割り当て始めてもよい。加えて、サービスは、すでに記載された方法で、変動する要求を満たすように、追加のVMインスタンスのそれぞれを自動的にスケールアップまたはダウンし続けてもよい。いくつかの実施形態では、操作メトリクスおよびユーザによって定義された閾値は、冗長性、アベイラビリティ、耐久性などの要件を部分的に含む場合があるため、ある特定量のCPUまたはRAMなど、他のリソース要件を満たすための単一サーバ上のリソース能力が十分である場合でさえ、異なる物理的なサーバ上でホストされる複数のVMへの「スケールアウト」が生じ得る。

#### 【0011】

種々の実施形態に従って、ホスト内のVMインスタンスのスケーリングアップ（またはダウン）の管理は、ウェブベースのグラフィカルユーザインターフェース（GUI）、顧客によって定義された閾値、または閉じた測定/アクションループにおける完全に自動的な推論を使用して実施され得る。この自動スケーリングは、いくつかの請求および/または支払いモデルを可能にし得る。例えば、顧客は、RAMの1GHz時間（もしくは他の所定の期間）あたりおよび/またはGB/時間の料金を請求し、個々の装置リソース（CPU、RAM、ネットワーク）を分離し、別個に請求するなどの汎用スケラブルVMインスタンスに対する料金を請求されてもよい。

#### 【0012】

種々の実施形態では、ウェブサービスは、ユーザ（例えば、顧客）が仮想計算機インスタンスのスケーリングを要求すること、またはこれらのVMインスタンスがリソース容量を拡大もしくは収縮するときに制御する種々の閾値を指定することを可能にするために使用され得る。ウェブサービスは、クエリおよびシンプルオブジェクトアクセスプロトコル（SOAP）APIの両方を含むことができる。しかしながら、ウェブサービスがSOAPベースのAPIコールに限定されず、インターネットなどのネットワークを使用して実施される任意の遠隔の手順/機能/方法の実行を含むことができることに留意されたい。

#### 【0013】

種々の実施形態では、ウェブサービスは、サイズ変更可能な計算容量（例えば、リソースセンター中の追加のサーバインスタンス）を提供するサービスプロバイダによって展開され得る。この計算容量は、顧客のソフトウェアシステムを構築およびホストするために使用され得る。サービスプロバイダは、APIまたはウェブツールおよびユーティリティを使用して、これらのリソースへのアクセスを提供することができる。したがって、ユーザは、リソースを追加または除去し、メトリクス、冗長性、アベイラビリティに基づきスケーリングするなどのために、サービスプロバイダによって露出されるAPI機能性にアクセスすることができる。

#### 【0014】

図1は、種々の実施形態に従った、追加のCPU、仮想CPU（VCPU）、物理CPU（PCPU）、物理CPUのコア、またはその部分であって、本明細書において概して「CPU」と称される部分を割り当てることによって、仮想計算機インスタンスをスケールアップすることの一実施例100を示す。例示された実施形態では、ホスト計算装置101は、仮想計算機インスタンス103および104を管理するハイパーバイザ102を含む。ハイパーバイザ102は、1つ以上のゲストオペレーティングシステムの実行を管理し、異なるオペレーティングシステムの複数のインスタンスが基礎的ハードウェアリソ

10

20

30

40

50

ースを共有することを可能にする。従来、ハイパーバイザは、動作しているゲストオペレーティングシステムの機能を有して、サーバハードウェア上にインストールされ、この場合、ゲストオペレーティングシステム自体は、サーバとしての機能を果たす。種々の実施形態では、少なくとも2タイプのハイパーバイザ102：タイプ1（ベアメタル）ハイパーバイザおよびタイプ2（ホスト型）ハイパーバイザがあり得る。タイプ1ハイパーバイザは、ハードウェアリソース上で直接実行し、1つ以上のゲストオペレーティングシステムを管理および制御し、それはハイパーバイザの上部で実行する。タイプ2ハイパーバイザは、オペレーティングシステム内で実行され、ハードウェアリソース上の第3のレベルで、概念的に1つ以上のゲストオペレーティングをホストする。いずれのタイプのハイパーバイザも、本明細書に記載される実施形態に従って実装され得る。ハイパーバイザ102は、ホストドメイン（またはサービス層もしくは仮想化層など）および1つ以上のゲストドメインなど、多くのドメイン（例えば、仮想計算機）をホストすることができる。一実施形態では、ホストドメイン（例えば、Dom-0）は、作成される第1のドメインであり、ハードウェアデバイスおよびハイパーバイザ102上で動作している他のドメインの全てを管理するのを助ける。例えば、ホストドメインは、1つ以上のゲストドメイン（例えば、Dom-U）の作成、破壊、移動、保存、または復元を管理することができる。種々の実施形態に従って、ハイパーバイザ102は、CPU、入力/出力（I/O）メモリ、およびハイパーバイザメモリなどのハードウェアリソースへのアクセスを制御する。例示された実施形態では、ハイパーバイザ102は、仮想計算機インスタンスにリソースを割り当てるか、またはリソースの割り当てを解除することによって、仮想計算機インスタンスのスケーリングを実施する、自動スケーリングサービス114を含む。あるいは、スケーリングサービスは、Dom-0中、またはホスト計算装置に対して外部に存在することができ、ホスト計算装置は、外部スケーリングサービスから受信されるコマンドを実行するために、シンエージェントを含んでもよい。

#### 【0015】

一実施形態に従って、ホスト計算装置101のハードウェアリソースは、物理メモリ116、1つ以上の中央処理装置（CPU）（107、108、109、110）、および任意の他のハードウェアリソースまたはデバイス111を含む。物理メモリ116は、ソリッドステートドライブ（SSD）、磁気ディスク記憶装置（HDD）、ランダムアクセスメモリ（RAM）などが含まれるがこれらに限定されない、任意のデータ記憶装置を含むことができる。種々の実施形態では、他のハードウェアリソース111としては、ネットワークインターフェースコントローラ（NIC）、グラフィックスプロセッシングユニット（GPU）、周辺入力/出力（I/O）デバイスなどが含まれ得るがこれらに限定されない。

#### 【0016】

一実施形態に従って、各仮想計算機インスタンス（103、104）は、少なくとも1つのユーザ112、113（例えば、サービスプロバイダの顧客）と関連付けられ得る。各仮想計算機インスタンスは、ユーザに代わって、少なくとも1つのアプリケーション（105、106）または他のサービスを実行することができる。例示された実施形態に従って、仮想計算機インスタンス103は、CPUのうちの1つ以上のセット（例えば、107および108）が割り当てられ、仮想計算機インスタンス104は、CPUの別のセット（例えば、110）が割り当てられる。種々の実施形態では、CPUは、実際の物理CPUであってもよく、あるいは仮想計算機に割り当てられる仮想CPU容量であってもよい。

#### 【0017】

種々の実施形態では、ユーザ（112、113）は、ユーザの仮想計算機と関連した種々の操作メトリクスに対する1つ以上の閾値を指定することが可能になる。図面に示されるように、仮想計算機インスタンス103上で実行しているアプリケーション105上の処理負荷が、ユーザの仮想計算機と関連した種々の操作メトリクスに対する1つ以上の閾値を超えると、システムは、増加した要求を満たすために、追加のCPU109を仮想

10

20

30

40

50

計算機インスタンス 103 に割り当てることができる。同様に、処理負荷が減少するとき、システムは、仮想計算機インスタンス 103 に割り当てられた CPU 容量を低減してもよい。

#### 【0018】

代替の実施形態では、仮想計算機インスタンスのスケールリングは、仮想計算機インスタンスに割り当てられるリソースの量を増加または減少させるという顧客からの要求を受信すると実施され得る。例えば、顧客は、サービスプロバイダによって提供される API を起動し、サービスプロバイダに追加の CPU 容量を仮想計算機インスタンスに割り当てよう要求してもよい。要求を受信したことに応答して、スケールリングサービス 114 は、追加の CPU 容量を仮想計算機に割り当てることができる。

10

#### 【0019】

種々の実施形態に従って、システムはまた、メモリ 116 および / または他のハードウェアリソース (例えば、NIC、GPU 容量など) を割り当てるか、またはその割り当てを解除することによって、仮想計算機インスタンス (103、104) をスケールリングすることができる。例えば、仮想計算機インスタンス 103 がメモリ容量の 90% に近づいている場合、システムは、追加のメモリ (例えば、物理メモリ、仮想メモリ) を仮想計算機インスタンス 103 に割り当ててもよい。

#### 【0020】

一実施形態に従って、仮想計算機インスタンスのスケールリングは、それを 1 つの仮想計算機インスタンスタイプから別のインスタンスタイプに変更することを含むことができ、この場合、各インスタンスタイプは、所定のリソースのセットと関連付けられる。例えば、所定の閾値を超えると、サービスは、顧客に割り当てられる仮想計算機を、「小型」インスタンスタイプ (例えば、1.7 GB RAM および 160 GB の記憶装置) から「中型」インスタンスタイプ (例えば、3.75 GB RAM および 410 GB の記憶装置) に変更してもよい。代替の実施形態では、仮想計算機インスタンスのスケールリングは、例えば、仮想計算機中で実行しているユーザアプリケーションまたはサービスによって要求されるように、かつ定義されたメトリクスおよび閾値に従って、例えば、あらゆる任意の増加量で、あらゆる任意の量の CPU、メモリ、または他のリソース容量を追加することによって、円滑に連続して実施され得る。

20

#### 【0021】

図 2 は、種々の実施形態に従って、サービスプロバイダによって展開される自動スケールリングサービスの一実施例 200 を示す。例示された実施形態では、サービスプロバイダ 201 は、サービスプロバイダがその顧客にリースを提供するホストサーバ (219、220) などの計算リソースのセットを所有および運営する。少なくとも 1 つの実施形態に従って、サービスプロバイダ 201 は、各ユーザ (例えば、顧客) が 1 つ以上の仮想計算機インスタンス (209、210、211、212) と関連する、共有されたリソース実行環境を作成する。仮想計算機インスタンスは、計算リソース 214 上で動作し、ネットワーク (例えば、インターネット) 上で、種々のデバイス上のユーザによってアクセス可能である。本開示全体を通して使用されるように、ネットワークは、インターネットまたは他の広域ネットワーク (WAN)、セルラーネットワーク、ローカルエリアネットワーク (LAN)、ストレージエリアネットワーク (SAN)、イントラネット、エクストラネットなどを含まれるがこれらに限定されない、互いと通信することが可能であるデバイスの任意の有線または無線ネットワークであってもよい。サービスプロバイダのホストサーバ (219、220) などの計算リソースは、データセンター、サーバファーム、コンテンツデリバリーネットワーク (CDN)、ポイントオブプレゼンス (POP) などのリソースの任意の物理または論理グループ中に位置し得る。

30

40

#### 【0022】

一実施形態に従って、サービスプロバイダは、ユーザ (例えば、顧客) が仮想計算機インスタンス (209、210、211、212) にアクセスし、それらを管理することを可能にするための、1 つ以上のアプリケーションプログラミングインターフェース (AP

50



I)を露出する。例えば、API 208は、ユーザに対する1つ以上の仮想計算機インスタンスを設定するために使用される仮想計算機イメージを提示するために、ユーザによって採用され得る。同様に、本明細書に記載される種々の実施形態に従って、API 208は、閾値が関連する1つ以上のユーザによって定義された閾値(215、216、217、218)およびメトリクスを指定するために採用され得る。例えば、1つの閾値は、すでに記載されたように、仮想計算機インスタンスに割り当てられるCPUの動作容量と関連付けられてもよい。別の閾値は、仮想計算機に割り当てられる利用可能なメモリの量と関連付けられてもよい。別の閾値は、特定の期間にわたって仮想計算機インスタンス上で実行しているアプリケーションによって処理されている要求の平均数であってもよい。加えて、APIは、追加のリソース容量を仮想計算機インスタンス(複数を含む)に割り当てるか、または仮想計算機インスタンス(複数を含む)からリソース容量の割り当てを解除するという要求を提示するために、顧客によって使用され得る。

10

#### 【0023】

一実施形態に従って、いったんユーザが閾値(215、216、217、218)を指定すると、自動監視およびスケーリングサービス213は、メトリクスが定義された閾値を超えたときを検出するために、ランタイム実行メトリクスを監視することができる。一実施形態では、自動スケーリングサービス213は、仮想計算機インスタンス(209、210、211、212)のそれぞれからランタイム情報を収集し、各VMインスタンスからリソースを割り当てるか、または割り当てを解除するための決定を行う、集中型サービスである。代替の実施形態では、自動スケーリングサービス213は、各ホスト計算機上で動作しているサービスとして実装され得、ホスト計算機上で仮想計算機インスタンスをスケーリングすることに責任があり得る。

20

#### 【0024】

いくつかの実施形態では、ホスト計算機は、スケーリングエージェント(221、222)を含む。スケーリングエージェントは、種々のメトリクスを集中型外部スケーリングサービス213に報告してもよく、かつ中央スケーリングサービス213からのコマンドを受信し、それらを実行してもよい。一実施形態に従って、仮想計算機インスタンスのいくつかは、種々のメトリクスをスケーリングエージェントに報告する、ゲストエージェント224を含んでもよく(例えば、メトリクスは、仮想計算機ならびにユーザ指定のメトリクス内から感知されるように、メモリプレッシャー、CPUプレッシャーなどを示す)、スケーリングエージェントは次いで、メトリクスを自動スケーリングサービス213に報告する。スケーリングサービス221は次いで、リソース容量中で仮想計算機インスタンスをスケールアップまたはダウンする決定を行ってもよい。

30

#### 【0025】

一実施形態に従って、ユーザのサービスに対する作業負荷または要求が、単一仮想計算機インスタンスがもはや作業に適切に対応するには十分ではない、ある特定の限界に達する場合、自動スケーリングサービス213は、ユーザに対する新しい仮想計算機インスタンスを設定し始めることができる。加えて、自動スケーリングサービス213は、すでに記載されたように、各インスタンスから計算リソースを追加および/または除去することによって、各個々の仮想計算機インスタンスのスケーリングアップおよびダウンを管理し続けることができる。いくつかの実施形態では、例えば冗長性要件を満たすために、作業の全てが単一インスタンスによって対応され得る場合でさえ、作業負荷を支持するために複数のVMインスタンスを必要とする閾値が定義されてもよい。この場合、サービス213は、ユーザ指定のサイジングポリシーを満たすために、2つ以上のVMインスタンスに割り当てられたリソースを同時に調節してもよい。

40

#### 【0026】

一実施形態に従って、仮想計算機インスタンスの自動スケーリングは、仮想計算機の利用の料金を顧客に請求するために使用され得る、多くの異なる請求モデルを可能にすることができる。一実施形態では、顧客は、自動的にスケーラブルな仮想計算機インスタンスを利用するための割増料金が請求されてもよい。例えば、一部の顧客は、その日のある特

50

定の時間中、またはある特定の場においてのみ、容量の増加を必要とする場合がある。それらの顧客に対して、オンデマンドで必要とされる容量を自動的に追加し、要求が低下した後、容量を低減することができる、自動スケーリングサービスを利用することが、費用の観点から有利であり得る。他の顧客は、顧客のサービスに対する要求を前もって容易に知らない場合があり、自動スケーリングサービスを活用することは、顧客のアプリケーションが、要件がよく理解される前にアプリケーションに過剰なリソース容量を費やすことなく、要求を満たすことを確認する、アプローチを提供することができる。別の実施形態では、顧客は、所与の期間にわたって利用されたりリソース毎に請求されてもよい（例えば、利用されたCPU時間あたり、1時間あたりのメモリのGBあたりなど）。

#### 【0027】

一実施形態に従って、サービスプロバイダ201はさらに、ホストサーバ(219、220)上に種々の仮想計算機インスタンス(209、210、211、212)を設定することに責任がある、配置サービス223を使用することができる。配置サービスは、仮想計算機インスタンスがスケーラブル仮想計算機であるかどうかを決定することができる。配置サービス223が、仮想計算機インスタンスがスケーラブルであると決定する場合、サービスは、ランタイムまたはオンデマンドで必要とされ得るリソース容量の増加に対応することができるように、過剰な容量を有するホストサーバ上に、仮想計算機を設定することができる。例えば、顧客が割増料金で自動的にスケーラブルな仮想計算機を購入する場合、配置サービスは、VMの作業負荷の増加に対応するのに十分な容量を有するホスト計算機上に、VMを配置してもよい。仮想計算機がスケーラブルではない場合、配置サービスは、過剰なまたは確保された容量がほとんどまたは全くないホスト計算機上に仮想計算機を設定してもよい。

#### 【0028】

一実施形態に従って、サービスプロバイダ201はさらに、顧客が顧客の仮想計算機に対してホスト計算装置の追加のリソースを購入する（例えば、割り当てる）ことを可能にする、電子市場を提供することができる。追加のリソースの料金は、ホスト計算装置上の1つ以上のリソースの需要と供給に少なくとも部分的に基づき得る。例えば、ホスト計算機上に利用可能な大量のCPU容量があり、要求が低いままであると予測される場合、そのホスト計算機上の仮想計算機に追加のCPUを割り当てるための料金は、低い可能性がある。同様に、利用可能なCPU容量が少量である場合、追加のCPUに対する料金はより高くなる可能性がある。このようにして需要と供給に基づき価格変動を可能にすることによって、サービスプロバイダは、リソース利用を最適化し、そのネットワークにわたって作業負荷のより効率的な分散を提供することが可能である。

#### 【0029】

図3Aは、種々の実施形態に従って、ホスト計算機上の仮想計算機インスタンスを自動的にスケーリングするための例示的な過程300を示す。この図は、特定の配列の機能動作を示し得るが、過程は必ずしも図示される特定の順序または動作に限定されない。当業者は、この図または他の図で表現される種々の動作が種々の方法で変更、再配置、同時に実施、または応用され得ることを理解するであろう。さらに、種々の実施形態の範囲から逸脱することなく、ある特定の動作または動作配列が、過程に追加または過程から削除され得ることを理解されたい。加えて、本明細書に含有される過程の例示は、コード実行の実際の配列を指定するよりもむしろ、当業者に過程の流れのアイデアを示すよう意図され、それは、異なる流れもしくは配列として実装、性能に対する最適化、または別様に種々の方法で修正されてもよい。

#### 【0030】

動作302において、仮想計算機インスタンスが顧客に対して設定される。仮想計算機インスタンスは、顧客を代表して、共有されたリソース計算環境のサービスプロバイダによって設定され得る。一実施形態に従って、顧客に対して設定される仮想計算機インスタンスは、特定のサービスを提供するアプリケーションを実行する。例えば、顧客は、いくつかの仮想計算機を使用して、データベースサーバとしての1つの仮想計算機インスタン

10

20

30

40

50

ス、フロントエンド（例えば、プレゼンテーション論理）サーバとして機能する別個の仮想計算機インスタンス、およびミドルウェア計算サーバとして機能する第3の仮想計算機インスタンスを使用して、サービスを展開してもよい。仮想計算機インスタンスを設定するとき、ユーザは、1つ以上の仮想計算機をスケーリングするための顧客によって定義された閾値を指定してもよい。一実施形態では、顧客は、特定の操作メトリクスに対する種々の値および閾値を指定するために、サービスプロバイダによって提供されるAPIを使用することができる。例えば、ユーザは、インスタンスが1分超の間60%のCPU容量で動作している場合、仮想計算機インスタンスのサイズを増大させるべきであると指定してもよい。別の実施形態では、ユーザは、仮想計算機インスタンスとは無関係に閾値のセットを提供し、後に、それが開始されたとき、または別様にすでに動作している後の時間に、これらの閾値のセットをインスタンスと関連付けることができてもよい。

10

#### 【0031】

動作303において、自動スケーリングサービスは、作業負荷の実行中、仮想計算機インスタンスの1つ以上の操作メトリクスを監視する。例えば、ホスト計算機上に存在するエージェント過程は、CPU利用、オープン接続数、IPパケット数、要求数などの種々のランタイム情報を連続的に収集してもよい。収集された情報は、顧客によって指定された命令に従って、各仮想計算機インスタンスをスケールアップまたはダウンするための決定を行うことができる、中央サービスに報告され得る。代替の実施形態では、サービスは、ホスト計算機内でホストされ得、収集されたメトリクスは、報告される必要がない。他の実施形態では、仮想計算機インスタンスは、仮想計算機のスケーリングに関連したユーザによって指定されたメトリクスを報告する、エージェントを含んでもよい。

20

#### 【0032】

動作304において、サービスは、1つ以上のメトリクスが顧客によって定義された閾値を超えていることを検出する。例えば、サービスは、仮想計算機インスタンスのCPU使用量が、顧客によって指定された最小タイムフレームに対する使用量閾値を超えていることを検出してもよい。

#### 【0033】

動作305において、サービスは、種々のリソースの容量を増加または減少させるために、仮想計算機インスタンスをスケーリングすることができる。一実施形態では、処理負荷が増加した場合、スケーリングサービスは、追加の計算リソースを仮想計算機インスタンスに割り当てる。例えば、スケーリングサービスは、さらなるCPU（または仮想単位のCPU容量の）を仮想計算機インスタンスに追加してもよい。別の実施形態では、スケーリングサービスは、仮想計算機インスタンスからリソースの一部分の割り当てを解除してもよく、および/またはリソースのその部分を他の仮想計算機インスタンスに移動させてもよい。いくつかの実施形態では、割り当て解除のために選択された部分またはサブセットは、顧客によって定義された閾値内にメトリクスを戻すために決定される。

30

#### 【0034】

図3Bは、種々の実施形態に従って、ユーザからの要求を受信したことに応答して、仮想計算機インスタンスをスケーリングする例示的な過程310を示す。動作311において、仮想計算機インスタンスは、すでに記載されたように、ユーザに対してホスト計算機上に設定される。いったん設定されると、仮想計算機インスタンスは、ユーザの代わりに作業負荷を実行することができる。動作312において、サービスプロバイダは、仮想計算機インスタンスに割り当てられた計算リソース容量を増加または減少させるという要求を受信する。例えば、ユーザは、作業負荷の増加により、仮想計算機インスタンスがより多くのCPU容量を必要とすることを決定してもよい。ユーザは次いで、追加のCPUを仮想計算機インスタンスに割り当てるために、APIを起動してもよい。動作313において、スケーリングサービスは、要求に応答して、追加の計算リソースを仮想計算機インスタンスに割り当てるか、または仮想計算機から計算リソースの割り当てを解除する。

40

#### 【0035】

図4は、種々の実施形態に従って、仮想計算機インスタンスを自動的にスケーリングし

50

、追加の仮想計算機インスタンスを割り当てるための例示的な過程 400 を示す。

【0036】

動作 401 において、すでに記載されたように、仮想計算機インスタンスがユーザに対して設定される。仮想計算機インスタンスは次いで、1つ以上の予め指定された操作メトリクスに関して監視される。動作 402 において、サービスは、仮想計算機インスタンスに対する1つ以上の操作メトリクスが、顧客によって定義された閾値を横切っている（超えているか、または下回っている）ことを検出してもよい。動作 403 において、スケーリングサービスは、追加の計算リソースを仮想計算機インスタンスに割り当てることによって、仮想計算機インスタンスを自動的にスケーリングする。例えば、スケーリングサービスは、追加のメモリ容量またはCPU容量を仮想計算機インスタンスに追加してもよい。

10

【0037】

動作 404 において、スケーリングサービスは、仮想計算機インスタンスが、サービスの必要とされる作業負荷を適切に満たすようにスケーリングすることができないことを決定してもよい。例えば、それは、仮想計算機インスタンスがサービスプロバイダによって許可される最大サイズまで拡大していることを決定してもよい。動作 405 において、スケーリングサービスは、作業負荷に対応するために、1つ以上の追加の仮想計算機インスタンスを自動的に設定し始めてもよい。追加の仮想計算機インスタンスのそれぞれもまた、すでに記載された方法でスケーリングされ得る（動作 405）。

【0038】

20

図5は、例示的な計算装置500の一般的なコンポーネントのセットの論理配列を示す。本実施例において、装置は、メモリデバイスまたは素子504中に記憶され得る命令を実行するためのプロセッサ502を含む。当業者に明らかとなるように、装置は、プロセッサ502による実行のためのプログラム命令のための第1のデータ記憶装置、画像またはデータのための別個の記憶装置、他のデバイスと情報を共有するための着脱式メモリなど、多くの種類のメモリ、データ記憶装置、または非一時的コンピュータ可読記憶媒体を含むことができる。装置は典型的に、タッチスクリーンまたは液晶ディスプレイ（LCD）などのなんらかの種類の表示素子506を含むが、ポータブルメディアプレーヤーなどのデバイスが、音声スピーカーを通してなど、他の手段を介して情報を伝達してもよい。考察されたように、多くの実施形態における装置は、ユーザからの従来の入力を受信することが可能な少なくとも1つの入力素子508を含む。この従来の入力としては例えば、プッシュボタン、タッチパッド、タッチスクリーン、ホイール、ジョイスティック、キーボード、マウス、キーパッド、またはユーザが装置にコマンドを入力することができる任意の他のそのようなデバイスもしくは素子が含まれ得る。しかしながら、いくつかの実施形態では、従来の入力装置は、ボタンを全く含まなくてもよく、ユーザが装置と接触している必要なく、装置を制御することができるように、視覚および音声によるコマンドの組み合わせによってのみ制御されてもよい。いくつかの実施形態では、図5の計算装置500は、Wi-Fi（登録商標）、Bluetooth（登録商標）、RF、有線、または無線通信システムなど、種々のネットワーク上で通信するための、1つ以上のネットワークインターフェース素子508を含むことができる。多くの実施形態における装置は、インターネットなどのネットワークと通信することができ、他の計算装置と通信することもできる。

30

40

【0039】

本開示の実施形態が以下の節を考慮して説明され得る。

1. 仮想計算機をスケーリングするためのコンピュータ実装方法であって、

実行可能な命令を有して構成される、1つ以上のコンピュータシステムの制御下で、

少なくとも1つの顧客に対する仮想計算機インスタンスを設定することであって、仮想計算機インスタンスは、ホスト計算装置上に設定される、仮想計算機インスタンスを設定することと、

アプリケーションプログラミングインターフェース（API）を介して顧客から、仮

50

想計算機インスタンスと関連した顧客によって定義された閾値を受信することと、

仮想計算機インスタンスの実行中、仮想計算機インスタンスと関連した1つ以上のメトリクスを監視することと、

要求を受信したことに応答して、仮想計算機への1つ以上の計算リソースの割り当てを調節することと、

1つ以上のメトリクスおよび顧客によって定義された閾値に少なくとも部分的に基づき、1つ以上の計算リソースの仮想計算機インスタンスへの割り当てを調節することであって、計算リソースは、処理リソース、ネットワークリソース、またはメモリリソースのうちの少なくとも1つを含む、割り当てを調節することと、  
を含む、方法。

10

2. 追加の計算リソースが仮想計算機インスタンスに割り当てることができないことを決定することと、

追加の計算リソースが仮想計算機インスタンスに割り当てることができないことを決定したことに応答して、顧客に対する第2の仮想計算機インスタンスを設定することと、  
をさらに含む、付記1に記載の方法。

3. 顧客と関連したアカウントに、仮想計算機インスタンスを動作するための料金を請求することであって、料金は、仮想計算機インスタンスに割り当てられた調節された計算リソースに少なくとも部分的に基づく、料金を請求することと、  
をさらに含む、付記1に記載の方法。

4. コンピュータ実装方法であって、

20

実行可能な命令を有して構成される、1つ以上のコンピュータシステムの制御下で、

サービスプロバイダ環境中で動作しているホスト計算装置によって、仮想計算機をホスト計算装置上に設定させることと、

サービスプロバイダ環境中で動作しているスケーリングサービスから、仮想計算機に割り当てられたリソースを調節するという要求を受信することと、

要求を受信したことに応答して、仮想計算機への1つ以上の計算リソースの割り当てを調節することと、  
を含む、コンピュータ実装方法。

5. スケーリングサービスは、

仮想計算機を動作させることと関連した1つ以上のメトリクスを監視することと、

30

1つ以上のメトリクスが少なくとも1つの指定閾値を超えていることを検出することと、

1つ以上のメトリクスが指定閾値を超えていることを検出したことに応答して、仮想計算機に割り当てられたリソースを調節するという要求を、ホストコンピュータ装置に伝送することと、

をさらに実施する、付記4に記載のコンピュータ実装方法。

6. ホスト計算装置上で実行している監視サービスは、

仮想計算機と関連した1つ以上のメトリクスを監視する、

1つ以上のメトリクスが閾値を超えていることを検出する、および

1つ以上のメトリクスが閾値を超えていることを検出したことに応答して、1つ以上の計算リソースの割り当てを調節する、

40

ように構成される、付記4に記載のコンピュータ実装方法。

7. スケーリングサービスは、

アプリケーションプログラミングインターフェース(API)を介して顧客から、仮想計算機に割り当てられたリソースを調節するという要求を受信することと、

要求を受信したことに応答して、仮想計算機に割り当てられたリソースを調節するという要求を伝送することと、

をさらに実施する、付記4に記載のコンピュータ実装方法。

8. 所定の限界が達せられるまで、仮想計算機に1つ以上の追加のリソースを割り当て続けることと、

50

1 つ以上の追加の仮想計算機にわたって仮想計算機の作業負荷を分配するように、1 つ以上の追加の仮想計算機を設定することと、  
をさらに含む、付記 4 に記載のコンピュータ実装方法。

9 . 割り当てを調節することは、

1 つ以上の追加の仮想計算機が、仮想計算機によって提供される少なくとも 1 つのサービスに関連した冗長性、アベイラビリティ、または耐久性のうちの少なくとも 1 つを満たすために必要とされることを決定することと、

1 つ以上の追加の仮想計算機を設定することと、  
をさらに含む、付記 4 に記載のコンピュータ実装方法。

10 . 仮想計算機は、作業負荷を実行することに関連した 1 つ以上のメトリクスを報告する、ゲストエージェントをさらに含む、付記 4 に記載のコンピュータ実装方法。 10

11 . スケーリングサービスからの要求に応答して、仮想計算機に割り当てられたホスト計算装置の 1 つ以上のリソースに少なくとも部分的に基づき、ユーザに請求書を送付すること  
をさらに含む、付記 4 に記載のコンピュータ実装方法。

12 . 顧客が、料金に対して、仮想計算機に割り当てられるホスト計算装置のリソースを取得することを可能にする、電子市場を提供することをさらに含み、料金は、1 つ以上のリソースの需要と供給に少なくとも部分的に基づく、  
付記 4 に記載のコンピュータ実装方法。

13 . 配置サービスによって、仮想計算機がスケーラブルであるように要求する API 要求を受信することと、 20

配置サービスによって、ホスト計算装置が追加のリソースを仮想計算機に追加する容量を含むことを決定することと、

配置サービスによって、仮想計算機をホスト計算装置上に設定することと、  
をさらに含む、付記 4 に記載のコンピュータ実装方法。

14 . 計算システムであって、

少なくとも 1 つのプロセッサと、

プロセッサによって実行されるとき、計算システムに、

ユーザに対する仮想計算機を設定する、サービスプロバイダ環境中で動作している仮想計算機に割り当てられたリソースを調節することに関するウェブサービス要求を受信させ、かつ 30

要求を受信したことに応答して、仮想計算機をホストするサーバに、1 つ以上の計算リソースを仮想計算機に割り当てさせる、  
命令を含む、メモリと、  
を備える、計算システム。

15 . ウェブサービス要求は、仮想計算機に割り当てるための計算リソースを指定する、1 つ以上の入力パラメータを含有する、付記 14 に記載の計算システム。

16 . ウェブサービス要求は、どのサーバが仮想計算機に割り当てられた計算リソースを調節するように命令されているかに応答して、条件の 1 セットを指定する、1 つ以上の入力パラメータを含有する、付記 14 に記載の計算システム。 40

17 . メモリは、実行時に、計算システムに、

仮想計算機を動作させることと関連した 1 つ以上のメトリクスを監視させ、

1 つ以上のメトリクスが少なくとも 1 つの指定閾値を超えていることを検出させ、

1 つ以上のメトリクスが指定閾値を超えていることを検出したことに応答して、仮想計算機に割り当てられたリソースを調節するという要求を、ホストコンピュータ装置に伝送させる、

命令をさらに含む、付記 14 に記載の計算システム。

18 . スケーリングサービスは、

アプリケーションプログラミングインターフェース (API) を介して顧客から、仮想計算機をスケーリングするという要求を受信することと、 50

要求を受信したことに応答して、仮想計算機をスケーリングするという命令を送送することと、  
をさらに実施する、付記 14 に記載の計算装置。

19. メモリは、プロセッサによって実行されるとき、計算装置に、  
所定の限界が達せられるまで、1 つ以上の追加のリソースを仮想計算機に割り当て続けさせ、かつ

1 つ以上の追加の仮想計算機にわたって仮想計算機の作業負荷を分配するように、1 つ以上の追加の仮想計算機を設定させる、  
命令をさらに含む、付記 14 に記載の計算装置。

20. メモリは、プロセッサによって実行されるとき、計算装置に、  
ユーザが、ウェブサービス要求に応答して、計算リソースを割り当てることによってスケーリングされることが可能である種類の仮想計算機インスタンスを選択したことを決定させ、かつ

ユーザによって選択された仮想計算機の種類に少なくとも部分的に基づき、ユーザに請求書を送付させる、  
命令をさらに含む、付記 14 に記載の計算装置。

21. メモリは、プロセッサによって実行されるとき、計算装置に、  
顧客が、計算装置上の 1 つ以上のリソースの需要と供給に少なくとも部分的に基づき、計算装置の 1 つ以上の追加のリソースを購入することを可能にする、電子市場を提供させる  
命令をさらに含む、付記 14 に記載の計算装置。

22. ユーザに対する仮想計算機を設定することは、  
配置サービスによって、仮想計算機がスケーラブルであるように要求するアプリケーションプログラミングインターフェース (API) 要求を受信することと、  
配置サービスによって、仮想計算機をホスト計算装置上に設定することと、  
をさらに含む、付記 14 に記載の計算装置。

23. ホスト計算装置上でユーザに対する仮想計算機を設定することであって、仮想計算機は、作業負荷を実行することが可能である、仮想計算機を設定することと、

仮想計算機をスケーリングするための命令を受信することであって、命令は、スケーリングサービスからホスト計算装置に受信され、スケーリングサービスは、ホスト計算装置に対して外部に存在する、命令を受信することと、

命令を受信したことに応答して、1 つ以上の計算リソースの仮想計算機への割り当てを調節することであって、1 つ以上の計算リソースは、ホスト計算装置のハイパーバイザによって割り当て可能である、割り当てを調節することと、  
を含む、動作のセットを実施するための、1 つ以上のプロセッサによって実行可能な命令の 1 つ以上の配列を記憶する、非一時的コンピュータ可読記憶媒体。

24. 仮想計算機は、少なくとも 1 つの顧客の代わりに、共有されたリソース計算環境サービスプロバイダによって設定され、スケーリングサービスは、サービスプロバイダによって展開される、付記 23 に記載の非一時的コンピュータ可読記憶媒体。

25. スケーリングサービスは、  
仮想計算機によって実行される作業負荷と関連した 1 つ以上のメトリクスを監視することと、

1 つ以上のメトリクスが少なくとも 1 つの指定閾値を超えていることを検出することと、

1 つ以上のメトリクスが指定閾値を超えていることを検出したことに応答して、仮想計算機インスタンスをスケーリングするという命令を送送することと、  
をさらに実施する、付記 23 に記載の非一時的コンピュータ可読記憶媒体。

26. スケーリングサービスは、  
アプリケーションプログラミングインターフェース (API) を介して顧客から、仮想計算機をスケーリングするという要求を受信することと、

10

20

30

40

50

要求を受信したことに応答して、仮想計算機をスケーリングするという命令を送送することと、  
をさらに実施する、付記 2 3 に記載の非一時的コンピュータ可読記憶媒体。

【 0 0 4 0 】

考察されたように、記載された実施形態に従って、種々の環境で異なるアプローチが実装され得る。例えば、図 6 は、種々の実施形態に従った態様を実装するための環境 6 0 0 の一実施例を示す。理解されるように、ウェブベースの環境が説明目的で使用されるが、種々の実施形態を実装するために、異なる環境が必要に応じて使用されてもよい。システムは、電子クライアント装置 6 0 2 を含み、それは、適切なネットワーク 6 0 4 上で要求、メッセージ、または情報を送信および受信し、かつ装置のユーザに情報を送り返すように動作可能な任意の適切な装置を含むことができる。電子クライアント装置の例としては、パーソナルコンピュータ、携帯電話、携帯型メッセージ装置、ラップトップコンピュータ、セットトップボックス、携帯情報端末、電子ブックリーダーなどが含まれる。ネットワークとしては、イントラネット、インターネット、セルラーネットワーク、ローカルエリアネットワーク、もしくは任意のそのようなネットワーク、またはそれらの組み合わせが含まれる、任意の適切なネットワークが含まれ得る。電子クライアントシステムのために使用されるコンポーネントは、選択されたネットワークおよび/または環境の種類によって少なくとも部分的に決まってもよい。任意の適切なネットワークを介して通信するためのプロトコルおよびコンポーネントはよく知られており、本明細書では詳述されない。ネットワーク上での通信は、有線または無線接続、およびそれらの組み合わせを介して可能となり得る。本実施例では、ネットワークは、環境が要求を受信し、それに応答してコンテンツを提供するためのウェブサーバ 6 0 6 を含むようなインターネットを含むが、当業者に明らかとなるように、他のネットワークに対して、同様の目的を果たす代替の装置が使用され得る。

【 0 0 4 1 】

例示的な環境は、少なくとも 1 つのアプリケーションサーバ 6 0 8 およびデータストア 6 1 0 を含む。いくつかのアプリケーションサーバ、層、または他の素子、過程、もしくはコンポーネントがあってもよく、それらは連鎖または別様に構成されてもよく、適切なデータストアからデータを取得するなどのタスクを実施するために相互に作用することができることが理解されるべきである。本明細書で使用される場合、用語「データストア」は、データを記憶、アクセス、および検索することが可能な任意の装置または装置の組み合わせを指し、それは、任意の標準、分散、またはクラスター環境において、データサーバ、データベース、データ記憶装置、およびデータ記憶媒体の任意の組み合わせおよび数を含んでもよい。アプリケーションサーバは、クライアント装置に対する 1 つ以上のアプリケーションの態様を実行するために必要に応じてデータストアと統合し、アプリケーションに対するデータアクセスおよびビジネスロジックの大部分に対応するための、任意の適切なハードウェアおよびソフトウェアを含むことができる。アプリケーションサーバは、データストアと連携したアクセス制御サービスを提供し、ユーザに転送されるテキスト、グラフィックス、オーディオ、および/またはビデオなどのコンテンツを生成することができ、それらは本実施例において、HTML、XML、または別の適切な構造化言語の形態で、ウェブサーバによってサーバに提供されてもよい。全ての要求および応答の処理、ならびにクライアント装置 6 0 2 とアプリケーションサーバ 7 0 8 との間のコンテンツの配信は、ウェブサーバ 6 0 6 によって処理され得る。本明細書で考察される構造化コードが、本明細書の他の部分で考察される任意の適切なデバイスまたはホスト計算機上で実行され得るため、ウェブおよびアプリケーションサーバは必要とされず、ただの例示的なコンポーネントにすぎないことを理解されたい。

【 0 0 4 2 】

データストア 6 1 0 は、特定の態様に関連するデータを記憶するためのいくつかの別々のデータ表、データベース、または他のデータ記憶装置機構および媒体を含むことができる。例えば、例示されるデータストアは、生成データ 6 1 2 およびユーザ情報 6 1 6 を記

10

20

30

40

50



憶するための機構を含み、それらは、生成側に対するコンテンツを提供するために使用され得る。データストアはまた、ログまたはセッションデータ 614 を記憶するための機構を含むことが示される。ページイメージ情報およびアクセス権情報など、データストア中に記憶される必要があり得る多くの他の態様が存在し得、それらは、必要に応じて上記に列挙された機構のいずれかに、またはデータストア 610 中の追加の機構に記憶され得ることを理解されたい。データストア 610 は、それと関連したロジックを通して、アプリケーションサーバ 608 から命令を受信し、それに応答してデータを取得、更新、または別様に処理するように動作可能である。一実施例では、ユーザは、ある特定の種類の項目に対する検索要求を提示してもよい。この場合、データストアは、ユーザの身元を検証するためにユーザ情報にアクセスしてもよく、その種類の項目に関する情報を取得するために、カタログ詳細情報にアクセスすることができる。情報は次いで、ユーザがユーザデバイス 602 上のブラウザを介して見ることができるウェブページ上の結果リストなどで、ユーザに戻され得る。対象となる特定の項目に関する情報は、専用ページまたはブラウザのウィンドウで見ることができる。

#### 【0043】

各サーバは典型的に、そのサーバの一般的な管理および動作に対する実行可能なプログラム命令を提供する、オペレーティングシステムを含み、典型的に、サーバのプロセッサによって実行されるとき、サーバがその目的とする機能を実施することを可能にする命令を記憶する、コンピュータ可読媒体を含む。オペレーティングシステムおよびサーバの一般的な機能性に対する好適な実装は既知であるか、または市販されており、特に本開示を考慮すると、当業者によって容易に実装される。

#### 【0044】

一実施形態における環境は、1つ以上のコンピュータネットワークまたは直接接続を使用して、通信リンクを介して相互接続される、いくつかのコンピュータシステムおよびコンポーネントを利用する、分散計算環境である。しかしながら、コンピュータシステムは、図6に示されるよりも少ないまたは多いコンポーネントを有するシステムにおいて、同等に動作し得ることが、当業者によって理解されるであろう。したがって、図6のシステム 600 の描写は、本開示の範囲に限定するものではなく、本質的に例示的なものとして見なされるべきである。

#### 【0045】

本明細書で考察または示唆される種々の実施形態は、幅広い種類の動作環境で実装され得、それらは場合によっては、多くのアプリケーションのいずれかを動作させるために使用され得る1つ以上のユーザコンピュータ、計算装置、または処理装置を含むことができる。ユーザまたはクライアント装置は、標準的な動作システムを起動するデスクトップまたはラップトップコンピュータなどの多くの汎用パーソナルコンピュータ、ならびに携帯電話ソフトウェアを起動し、多くのネットワークおよびメッセージプロトコルをサポートすることが可能な、携帯電話、無線、および手持ち式装置のいずれかを含むことができる。コンピュータシステムはまた、開発およびデータベース管理などの目的で、種々の市販のオペレーティングシステムおよび他の既知のアプリケーションのいずれかを起動する多くのワークステーションを含むことができる。これらのデバイスはまた、ダミー端子、シンクライアント、ゲームシステム、およびネットワークを介して通信することが可能な他のデバイスなど、他の電子デバイスを含むことができる。

#### 【0046】

ほとんどの実施形態は、TCP/IP、OSI、FTP、UPnP、NFS、CIFS、およびAppleTalkなど、種々の市販のプロトコルのいずれかを使用して、通信をサポートするために、当業者によく知られている少なくとも1つのネットワークを利用する。ネットワークは、例えば、ローカルエリアネットワーク、広域ネットワーク、仮想プライベートネットワーク、インターネット、イントラネット、エクストラネット、公衆交換電話網、赤外線ネットワーク、無線ネットワーク、およびそれらの任意の組み合わせであってもよい。

## 【 0 0 4 7 】

ウェブサーバを利用する実施形態では、ウェブサーバは、HTTPサーバ、FTPサーバ、CGIサーバ、データサーバ、Java（登録商標）サーバ、およびビジネスアプリケーションサーバを含む、種々のサーバまたは中間階層アプリケーションのいずれかを起動することができる。サーバ（複数を含む）はまた、Java（登録商標）、C、C#、もしくはC++などの任意のプログラミング言語、またはPerl、Python、もしくはTCLなどの任意のスクリプト言語、ならびにそれらの組み合わせで書かれる1つ以上のスクリプトまたはプログラムとして実装されてもよい、1つ以上のウェブアプリケーションを実行することなどによって、ユーザデバイスからの要求に应答して、プログラムまたはスクリプトを実行することが可能であってもよい。サーバ（複数を含む）はまた、Oracle（登録商標）、Microsoft（登録商標）、Sybase（登録商標）、およびIBM（登録商標）から市販されているものが含まれるがこれらに限定されない、データベースサーバを含んでもよい。

10

## 【 0 0 4 8 】

環境は、上記で考察されるような種々のデータストならびに他のメモリおよび記憶装置を含むことができる。これらは、ネットワーク上のコンピュータのうちの1つ以上にローカルな（および/もしくは常駐する）、またはコンピュータのいずれかまたは全部から遠隔の記憶媒体など、種々の位置に存在することができる。実施形態の特定のセットにおいて、情報は、当業者にはよく知られているストレージエリアネットワーク（「SAN」）に存在してもよい。同様に、コンピュータ、サーバ、または他のネットワーク装置による機能を実施するための任意の必要なファイルが、必要に応じて、ローカルおよび/またはリモートで記憶されてもよい。システムがコンピュータ化された装置を含む場合、コンピュータ化された各装置は、バスを介して電氣的に連結されてもよいハードウェア要素を含むことができ、要素は、例えば、少なくとも1つの中央処理装置（CPU）、少なくとも1つの入力装置（例えば、マウス、キーボード、コントローラー、タッチスクリーン、またはキーパッド）、および少なくとも1つの出力装置（例えば、表示装置、プリンター、またはスピーカー）を含む。コンピュータ化された装置を含むシステムはまた、ディスクドライブ、光学式記憶装置、およびランダムアクセスメモリ（「RAM」）またはリードオンリーメモリ（「ROM」）などのソリッドステート記憶装置などの1つ以上の記憶装置、ならびにリムーバブルメディアデバイス、メモリカード、フラッシュカードなどを含んでもよい。

20

30

## 【 0 0 4 9 】

コンピュータ化された装置を含むシステムはまた、上記に記載されるようなコンピュータ可読記憶媒体リーダー、通信装置（例えば、モデム、ネットワークカード（無線または有線）、赤外線通信装置など）、および作業メモリを含むことができる。コンピュータ可読記憶媒体リーダーは、遠隔、ローカル、固定、および/またはリムーバブル記憶装置、ならびにコンピュータ可読情報を一時的および/またはより永久的に含有、記憶、伝送、および検索するための記憶媒体を表す、コンピュータ可読記憶媒体と接続され得るか、またはコンピュータ可読記憶媒体を受信するように構成され得る。システムおよび種々の装置はまた、典型的に、クライアントアプリケーションまたはウェブブラウザなど、オペレーティングシステムおよびアプリケーションプログラムを含む、少なくとも1つの作業メモリ装置内に位置する多くのソフトウェアアプリケーション、モジュール、サービス、または他の要素を含む。代替の実施形態が上記に記載される実施形態からの多くの変更を有してもよいことを理解されたい。例えば、カスタマイズされたハードウェアもまた使用されてもよく、および/または特定の要素がハードウェア、ソフトウェア（アプレットなどの高移植性ソフトウェアを含む）、または両方で実装されてもよい。さらに、ネットワーク入力/出力装置などの他の計算装置への接続が採用されてもよい。

40

## 【 0 0 5 0 】

コードまたはコードの部分を含むための記憶媒体およびコンピュータ可読媒体は、所望の情報を記憶するために使用され得、かつシステム装置によってアクセスされ得る、

50

RAM、ROM、EEPROM（登録商標）、フラッシュメモリもしくは他のメモリ技術、CD-ROM、デジタル多用途ディスク（DVD）、もしくは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、もしくは他の磁気記憶装置、または任意の他の媒体を含む、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、または他のデータなどの情報の記憶および/または伝送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、リムーバブルおよび非リムーバブル媒体などであるがこれらに限定されない、記憶媒体および通信媒体を含む、当該技術分野で既知であるか、または使用される任意の適切な媒体を含むことができる。本明細書に提供される本開示および教示に基づき、当業者は、種々の実施形態を実装するための他の手段および/または方法を理解するであろう。

10

#### 【0051】

結果的に、本明細書および図面は、制限的な意味よりもむしろ例示的な意味でみなされるものとする。しかしながら、特許請求の範囲に記載されるような本発明のより広い精神および範囲から逸脱することなく、種々の修正および変更がそれに行われてもよいことが明らかとなるであろう。

以下に、本願出願当初の特許請求の範囲に記載された発明を付記する。

#### [C1]

コンピュータ実装方法であって、

実行可能な命令を有して構成される、1つ以上のコンピュータシステムの制御下で、

サービスプロバイダ環境中で動作しているホスト計算装置によって、仮想計算機を前記ホスト計算装置上に設定させることと、

20

前記サービスプロバイダ環境中で動作しているスケーリングサービスから、前記仮想計算機に割り当てられたリソースを調節するという要求を受信することと、

前記要求を受信したことに応答して、前記仮想計算機への1つ以上の計算リソースの割り当てを調節することと、

を含む、コンピュータ実装方法。

#### [C2]

前記スケーリングサービスは、

前記仮想計算機を動作させることと関連した1つ以上のメトリクスを監視することと

30

、  
前記1つ以上のメトリクスが少なくとも1つの指定閾値を超えていることを検出することと、

前記1つ以上のメトリクスが前記指定閾値を超えていることを検出したことに応答して、前記仮想計算機に割り当てられたリソースを調節するという前記要求を、前記ホストコンピュータ装置に伝送することと、

をさらに実施する、[C1]のコンピュータ実装方法。

#### [C3]

前記ホスト計算装置上で実行している監視サービスは、

前記仮想計算機と関連した1つ以上のメトリクスを監視する、

前記1つ以上のメトリクスが閾値を超えていることを検出する、および

40

前記1つ以上のメトリクスが前記閾値を超えていることを検出したことに応答して、前記1つ以上の計算リソースの前記割り当てを調節する、  
ように構成される、[C1]のコンピュータ実装方法。

#### [C4]

前記スケーリングサービスは、

アプリケーションプログラミングインターフェース（API）を介して顧客から、前記仮想計算機に割り当てられたリソースを調節するという要求を受信することと、

前記要求を受信したことに応答して、前記仮想計算機に割り当てられたリソースを調節するという前記要求を伝送することと、

をさらに実施する、[C1]のコンピュータ実装方法。

50

[ C 5 ]

所定の限界が達せられるまで、前記仮想計算機に前記 1 つ以上の追加のリソースを割り当て続けることと、

1 つ以上の追加の仮想計算機にわたって前記仮想計算機の作業負荷を分配するように、1 つ以上の追加の仮想計算機を設定することと、  
をさらに含む、[ C 1 ] のコンピュータ実装方法。

[ C 6 ]

前記仮想計算機は、前記作業負荷を実行することと関連した 1 つ以上のメトリクスを報告する、ゲストエージェントをさらに含む、[ C 1 ] のコンピュータ実装方法。

[ C 7 ]

前記スケーリングサービスからの要求に応答して、前記仮想計算機に割り当てられた前記ホスト計算装置の前記 1 つ以上のリソースに少なくとも部分的に基づき、前記ユーザに請求書を送付すること

をさらに含む、[ C 1 ] のコンピュータ実装方法。

[ C 8 ]

顧客が、料金に対して、前記仮想計算機に割り当てられる前記ホスト計算装置のリソースを取得することを可能にする、電子市場を提供することをさらに含み、前記料金は、前記 1 つ以上のリソースの需要と供給に少なくとも部分的に基づく、

[ C 1 ] のコンピュータ実装方法。

[ C 9 ]

配置サービスによって、前記仮想計算機がスケーラブルであるように要求する A P I 要求を受信することと、

前記配置サービスによって、前記ホスト計算装置が追加のリソースを前記仮想計算機に追加する容量を含むことを決定することと、

前記配置サービスによって、前記仮想計算機を前記ホスト計算装置上に設定することと、  
をさらに含む、[ C 1 ] のコンピュータ実装方法。

[ C 1 0 ]

計算システムであって、

少なくとも 1 つのプロセッサと、

前記プロセッサによって実行されるとき、前記計算システムに、

ユーザに対する仮想計算機を設定する、サービスプロバイダ環境中で動作している仮想計算機に割り当てられたリソースを調節することに関するウェブサービス要求を受信させ、かつ

前記要求を受信したことに応答して、前記仮想計算機をホストするサーバに、1 つ以上の計算リソースを前記仮想計算機に割り当てさせる、

命令を含む、メモリと、

を備える、計算システム。

[ C 1 1 ]

前記ウェブサービス要求は、どの前記サーバが前記仮想計算機に割り当てられた前記計算リソースを調節するように命令されているかに応答して、条件の 1 セットを指定する、1 つ以上の入力パラメータを含有する、[ C 1 0 ] に記載の計算システム。

[ C 1 2 ]

前記メモリは、前記プロセッサによって実行されるとき、前記計算装置に、

所定の限界が達せられるまで、前記 1 つ以上の追加のリソースを前記仮想計算機に割り当て続けさせ、かつ

1 つ以上の追加の仮想計算機にわたって前記仮想計算機の作業負荷を分配するように、1 つ以上の追加の仮想計算機を設定させる、

命令をさらに含む、[ C 1 0 ] に記載の計算装置。

[ C 1 3 ]

10

20

30

40

50

前記メモリは、前記プロセッサによって実行されるとき、前記計算装置に、  
 前記ユーザが、前記ウェブサービス要求に応答して、前記計算リソースを割り当てるこ  
 とによってスケーリングされることが可能である種類の前記仮想計算機インスタンスを選  
 択したことを決定させ、かつ

前記ユーザによって選択された前記仮想計算機の前記種類に少なくとも部分的に基づき  
 、前記ユーザに請求書を送付させる、  
 命令をさらに含む、[ C 1 0 ] に記載の計算装置。

[ C 1 4 ]

前記メモリは、前記プロセッサによって実行されるとき、前記計算装置に、  
 顧客が、前記計算装置上の前記１つ以上のリソースの需要と供給に少なくとも部分的に  
 基づき、前記計算装置の１つ以上の追加のリソースを購入することを可能にする、電子市  
 場を提供させる

命令をさらに含む、[ C 1 0 ] に記載の計算装置。

[ C 1 5 ]

前記ユーザに対する前記仮想計算機を設定することは、  
 配置サービスによって、前記仮想計算機がスケーラブルであるように要求するアプリケ  
 ーションプログラミングインターフェース ( A P I ) 要求を受信することと、  
 前記配置サービスによって、前記仮想計算機を前記ホスト計算装置上に設定することと

をさらに含む、[ C 1 0 ] に記載の計算装置。

10

20

【図 1】

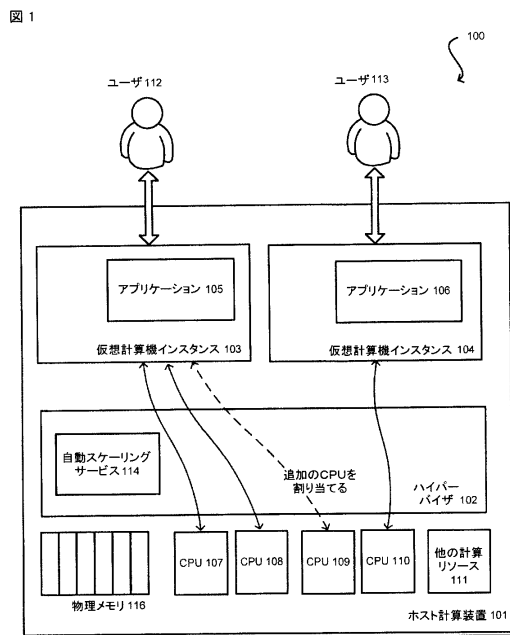


FIGURE 1

【図 2】

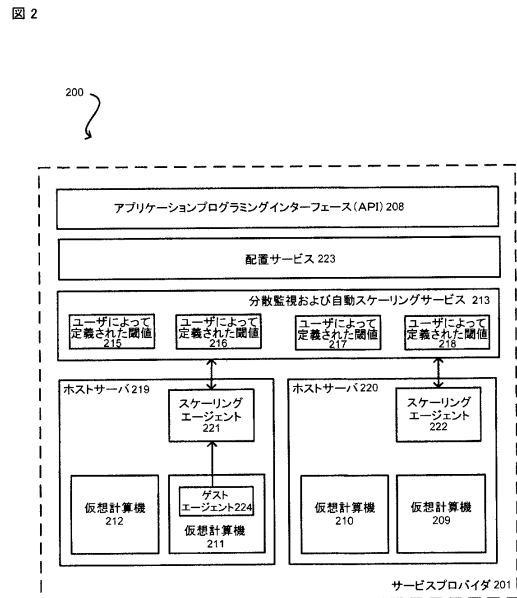


FIGURE 2

【図 3 A】

図 3A

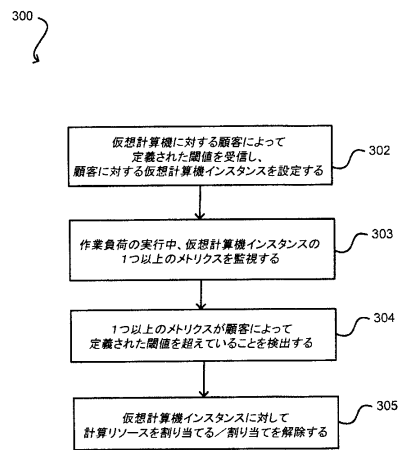


FIGURE 3A

【図 3 B】

図 3B

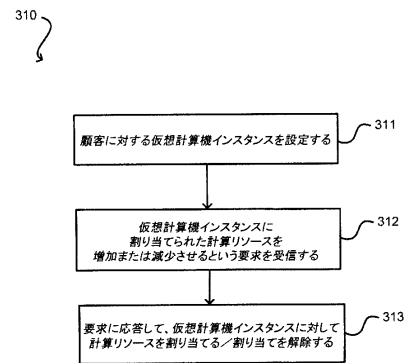


FIGURE 3B

【図 4】

図 4

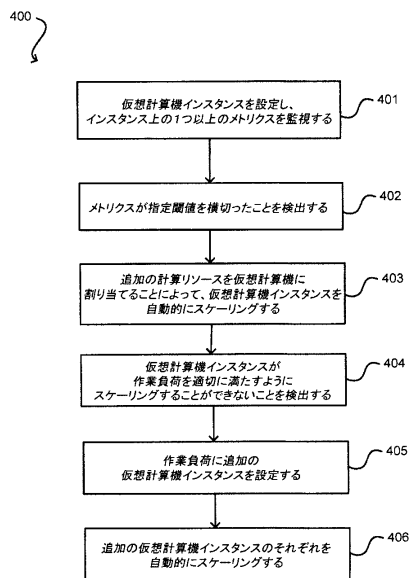


FIGURE 4

【図 5】

図 5

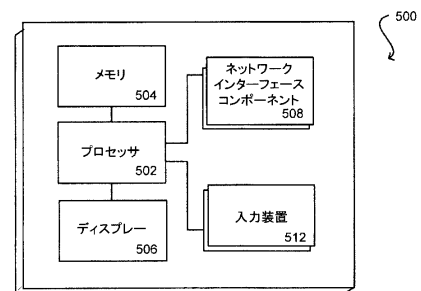


FIGURE 5

## 【図 6】

図 6

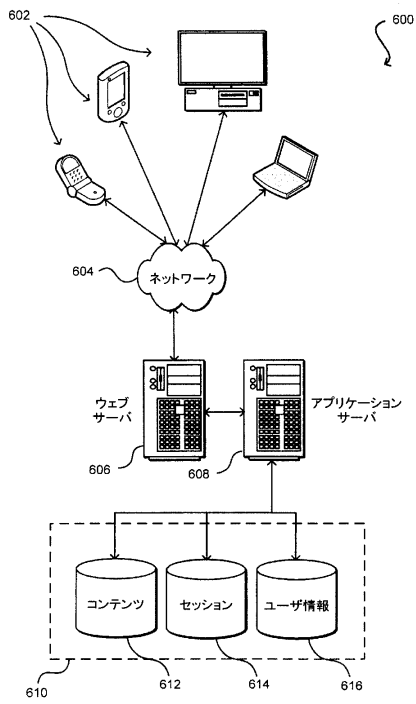


FIGURE 6

## フロントページの続き

- (74)代理人 100140176  
弁理士 砂川 克
- (74)代理人 100124394  
弁理士 佐藤 立志
- (74)代理人 100112807  
弁理士 岡田 貴志
- (74)代理人 100111073  
弁理士 堀内 美保子
- (72)発明者 マー、マイケル・デービッド  
アメリカ合衆国、ワシントン州 48109、シアトル、テリー・アベ・ノース 410
- (72)発明者 コワルスキー、マルチン・ピー.  
アメリカ合衆国、ワシントン州 48109、シアトル、テリー・アベ・ノース 410

審査官 田中 幸雄

- (56)参考文献 特開2011-203910(JP,A)  
特開2012-99062(JP,A)  
特開2011-170679(JP,A)  
特開2007-200347(JP,A)  
寺田亜紀, ネット世界の「雲」をつかめ! クラウドの成分分析と導入メリット, G-CLOUD Magazine, 日本, (株)技術評論社, 2011年 3月15日, 第2号, 128-131ページ  
荒井康宏ほか, 徹底解説! Amazon Web Services WebエンジニアのためのAWSガイド, G-CLOUD Magazine, 日本, (株)技術評論社, 2010年 9月10日, 10-19ページ  
波戸邦夫ほか, インタークラウドに向けたシステムアーキテクチャの提案, 電子情報通信学会技術研究報告, 日本, 社団法人電子情報通信学会, 2011年 2月24日, Vol.110 No.449, 151-156ページ

(58)調査した分野(Int.Cl., DB名)

G06F 9/50  
G06F 9/46