

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2015-537307

(P2015-537307A)

(43) 公表日 平成27年12月24日(2015.12.24)

| (51) Int. Cl.         | F I            | テーマコード (参考) |
|-----------------------|----------------|-------------|
| G06F 9/50 (2006.01)   | G06F 9/46 465D | 5K030       |
| H04L 12/721 (2013.01) | H04L 12/721 Z  |             |
| H04L 12/743 (2013.01) | H04L 12/743    |             |
| G06F 9/46 (2006.01)   | G06F 9/46 350  |             |
| G06F 9/54 (2006.01)   | G06F 9/46 480D |             |

審査請求 有 予備審査請求 未請求 (全 22 頁)

(21) 出願番号 特願2015-543294 (P2015-543294)  
 (86) (22) 出願日 平成26年1月23日 (2014.1.23)  
 (85) 翻訳文提出日 平成26年12月25日 (2014.12.25)  
 (86) 国際出願番号 PCT/CN2014/071233  
 (87) 国際公開番号 W02015/024368  
 (87) 国際公開日 平成27年2月26日 (2015.2.26)  
 (31) 優先権主張番号 201310367864.2  
 (32) 優先日 平成25年8月22日 (2013.8.22)  
 (33) 優先権主張国 中国 (CN)

(71) 出願人 514311221  
 インスパイア・エレクトロニック・インフォメーション・インダストリー・コーポレーション・リミテッド  
 中華人民共和国、250014 シャンドング、ジナン、ハイテク・ゾーン、シュンヤ・ロード、No. 1036  
 (74) 代理人 100095267  
 弁理士 小島 高城郎  
 (74) 代理人 100124176  
 弁理士 河合 典子  
 (74) 代理人 100146950  
 弁理士 南 俊宏

最終頁に続く

(54) 【発明の名称】 コンポーネント指向ハイブリッドクラウドオペレーティングシステムのアーキテクチャ及びその通信方法

(57) 【要約】

コンポーネント指向ハイブリッドクラウドオペレーティングシステム及びその通信方法を提供し、階層、オブジェクト及びメッセージモデルに基づきハイブリッドアーキテクチャを構築し、かつコンポーネント指向思想でコンポーネント及びそのプロセス環境を管理し、コンポーネント処理用クラスタに対して高効率ルーティング、読取り書き込み分離及び負荷バランスを行い、クラウドオペレーティングシステムに対する開放性と互換性、疎結合及び拡張性の要求を満足し、既存のクラウドオペレーティングシステムの自主管理問題、コンポーネントのスケールアウト問題及び状態コンポーネントの高可用性問題を解決し、開放性と交換性、拡張性、疎結合のコンポーネント指向クラウドオペレーティングシステムアーキテクチャを完備させ、かつ通信方法によりクラウドオペレーティングシステムのスケーラビリティ及び高可用性を保障する。

【選択図】 図 1

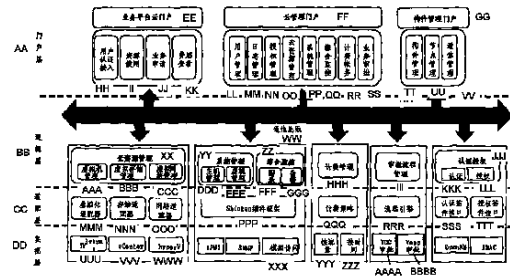


図 1 / FIG. 1

- AA PORTAL LAYER
- BB LOGICAL LAYER
- CC ADAPTATION LAYER
- DD IMPLEMENTATION LAYER
- EE SERVICE PLATFORM CLOUD PORTAL
- FF CLOUD MANAGEMENT PORTAL
- GG COMPONENT MANAGEMENT PORTAL
- HH USER AUTHENTICATION ACCESS
- II RESOURCE USE
- JJ SERVICE APPLICATION
- KK COST CHECK
- LL USER MANAGEMENT
- MM LOG MANAGEMENT
- NN AUTHORIZATION MANAGEMENT
- OO XX CLOUD RESOURCE MANAGEMENT
- PP YY SYSTEM MANAGEMENT
- QQ ZZ COMPREHENSIVE MONITORING
- RR BILLING AND ACCOUNTING
- SS SERVICE APPROVAL
- TT COMPONENT MANAGEMENT
- UU KPI MANAGEMENT
- VV COMMUNICATION MANAGEMENT
- WW COMMUNICATION BUS
- AAA VIRTUAL MACHINE MANAGEMENT
- BBB VIRTUAL STORAGE MANAGEMENT
- CCC VIRTUAL NETWORK MANAGEMENT
- DDD HOST MANAGEMENT
- EEE VIRTUAL MACHINE MANAGEMENT
- FFF CHART
- GGG ALARM
- HHH BILLING MANAGEMENT
- III APPROVAL FLOW MANAGEMENT
- JJJ AUTHENTICATION AND AUTHORIZATION
- KKK AUTHENTICATION
- LLL AUTHORIZATION
- MMM VIRTUALIZATION ADAPTER
- NNN STORAGE ADAPTER
- OOO NETWORK ADAPTER
- PPP SHIMKEN PLUG-IN FRAMEWORK
- QQQ BILLING POLICY
- RRR FLOW ENGINE
- SSS AUTHENTICATION PLUG-IN INTERFACE
- TTT AUTHORIZATION PLUG-IN INTERFACE
- UUU VIRTUAL
- VVV VCENTER
- WWW HYPERV
- XXX SIMULATED ACCESS
- YYY IN TRAFFIC
- ZZZ IN TIME
- AAAA VDC APPROVAL
- BBBB YAPP APPROVAL

## 【特許請求の範囲】

## 【請求項 1】

コンポーネント指向ハイブリッドクラウドオペレーティングシステムであって、前記システムは、階層モデル、オブジェクトモデル及びメッセージモデルに基づいてハイブリッドアーキテクチャを構築し、かつ、構成コンポーネント及びその処理環境を管理するためにコンポーネント指向思想を適用し、これをベースに、前記システムは、コンポーネント処理用クラスタに対して高効率のルーティング、読取り書込み分離及び負荷バランスとを行い、クラウドオペレーティングシステムに対する開放性と互換性、疎結合と拡張性という要求を満足し、かつ、既存のクラウドオペレーティングシステムの自主管理問題、コンポーネントのスケールアウト問題及びステートフルコンポーネントの高可用性問題を解決し；

10

階層モデルの観点から、前記システムは、上から下にポータル層、論理層、アダプテーション層及び実装層に区分され、各層は相対的に独立し、各層において各々標準インタフェースを定義することによりシステムの開放性が強化され、かつ、各層に異なる機能を適応させることにより互換性が強化され；

オブジェクトモデルの観点から、前記クラウドオペレーティングシステムは、クラウドポータル、クラウド管理ポータル、クラウドリソース管理、監視管理、計量と課金、サービス承認、及び、権限と認証の各機能モジュールにより構成され、各機能モジュールは、Restメッセージベースの呼出しにより通信を行い、互いに自由に組み合わせられ、かつ、要求により分散され配置されることができ、要求により新規機能モジュールが付加価値として開発されることができ、かつ、プラットフォームの拡張性が論理層における異なる機能モジュール間での相互操作を実現することにより強化され；

20

最小形態でインストールされた前記クラウドオペレーティングシステムは、クラウドポータル、クラウド管理ポータル及びクラウドリソース管理の各機能モジュールのみにより構成され、これをベースに、監視、課金、承認又は他の機能モジュールは、要求によりカスタム化されて拡張され、Restメッセージベースの呼出しにより通信を行い、要求により分散され配置されかつ付加価値として開発され、前記システムの拡張性が論理層における異なる機能モジュール間での相互操作を実現することにより強化され、

メッセージベースの通信方法が非同期呼出しをサポートするために導入され、メッセージ通信インターフェースJMSがRestメッセージを転送することによりシステムアーキテクチャをさらに疎結合とするために用いられ、これをベースに、コンポーネント指向設計が適用され、かつ、コンポーネント管理モジュールがコンポーネントのメタデータを管理するとともにコンポーネントの操作状態を監視することを担当し、

30

要求により拡張、分散及び配置を実現できるにも拘わらず、オブジェクトアーキテクチャは、RPC（リモート・プロセス・コール）同期通信方法に属しており、送信端は、受信端からのリターンを待った後にのみ実行を継続でき、かつ、双方のプロセスが密に結合しているため、前記システムの拡大化と複雑化に伴い、コンポーネント間の連携関係が複雑過ぎるようになり、この問題に対し、オブジェクトアーキテクチャに基づいてメッセージベースの通信方法が導入され、メッセージ通信インターフェースJMSを用いてRestメッセージを転送することにより、送信端と受信端のライフタイムを異ならせ、非同期呼出しをサポートし、かつ、前記システムアーキテクチャをさらに疎結合としており、ポータル層とクラウドリソース層の間における仮想マシンの少なくともターンオン、シャットダウン及びサスペンドの操作は、非同期方式にて実現され、ポータルが、コマンドを送信した後、応答を待つ必要がなくなりターンすることができ、これにより、ユーザの相互動作効果を向上させ；

40

上記ハイブリッドアーキテクチャに基づく前記クラウドオペレーティングシステムは、開放性、互換性及び拡張性の要求を満足させることができ、これをベースに、コンポーネント指向設計思想に基づき、前記コンポーネント管理モジュールは、コンポーネントについてのメタデータ情報を管理することを担当し、少なくともログイン、削除、修正及びクエリの操作をサポートし、その場合、

50

1つのコンポーネントは、1つの三要素セットであり、ネーム、サービスセット、アクセスアドレス、及び、記述を含み；

1つのサービスは、1つの四要素セットであり、ネーム、タイプ、メッセージプロトコル、パラメタリスト、キーネーム、機能記述、及び、非機能記述を含み；

コンポーネントの記述及び管理に加えて、前記コンポーネント管理モジュールはさらに、その処理環境を監視し、コンポーネントのスケラビリティ及び可用性を確保するための基本サービスを提供し、かつ、クラウドオペレーティングシステムの自主管理能力を完備させ、かつ、コンポーネントがログインした時、前記システムは、ユーザ名user、パスワードpsw、及び、固有のコンポーネントIDを割当て、その後のコンポーネント処理用クラスタのアクセスのプロセスは、

1) 処理用クラスタが、アドレスがurlであるシステムバスに対してアクセス要求を開始し、システムバスがアクセスノードのユーザ名、パスワード及びIDを認証し、認証が通ると、接続が確立され、そのコードは、

Connection = ConnectionFactory.createConnection(user,psw,url) であり；

2) 書込み操作トピックを設定し、コンポーネントの各処理ノードが当該トピックから書込み操作をサブスクライブし、

write\_topic = session.createTopic(id+"WRITE\_TOPIC");

write\_topic\_consumer = session.createConsumer(write\_topic);

3) コンポーネントがwriteTopicListenerのonMessage方法において書込み処理を実現し、かつシステムバスにログインし、

write\_topic\_consumer.setMessageListener(writeTopicListener);

4) コンポーネントの処理ノード数numに従って読取り操作の待ち行列セットを設定し、各処理ノードが1つの待ち行列に対応した読取り操作をサブスクライブし、

read\_queue = session.createMultiQueue(num, id+"READ\_QUEUE");

read\_queue\_consumer = session.createConsumer(read\_queue);

5) readQueueListenerのonMessage方法において特定の読取り処理機能を実現し、かつシステムバスにログインし、

read\_queue\_consumer.setMessageListener(readQueueListener);

前記書込み読取り分離待ち行列セットに基づき、コンポーネント管理モジュールにおいて各コンポーネントのサービスタイプを区別し、サービスタイプを冪等又は非冪等に設定し、冪等操作はステートレス操作に属し、同一状態下で実行した結果が毎回同一であり、非冪等操作はステートフル操作に属し、同一状態下で実行した結果が毎回異なり、ルータは、サービスタイプによってルーティングを行い、非冪等操作は固有の書込み待ち行列に送信され、そして冪等操作は負荷バランスストラテジーによって別々の読取り待ち行列に送信され、

読取り操作の負荷バランスのプロセスは、

1) ノードの処理能力を計算し、

ノードiのCPUクロック周波数、メモリ容量及びI/O帯域幅を各々 $C_i$ 、 $M_i$ 及び $B_i$ とすると、クラスタの各種リソースはノードの各種リソースの総計であり、つまり $C = C_i$ 、 $M = M_i$ 、 $B = B_i$ であり、

ノードiのCPU重み付け値は $W_i^{CPU} = C_i / C$ であり、メモリ容量重み付け値は $W_i^{RAM} = M_i / M$ であり、I/O帯域幅重み付け値は $W_i^{IO} = B_i / B$ であり、

コンポーネントサービスに必要なリソース割合を、それぞれ $p^{CPU}$ 、 $p^{RAM}$ 、 $p^{IO}$ とすると、ノードiの処理能力は、

$W_i = p^{CPU}W_i^{CPU} + p^{RAM}W_i^{RAM} + p^{IO}W_i^{IO}$  であり、

2) 書込み読取り操作の重み付け値によって各ノードの負荷を計算し、

読取り待ち行列 $L_r$ の書込み待ち行列 $L_w$ に対する書込み読取り操作のオーバーヘッド比率をaとすると、ノードiの負荷は $L_i = L_r^i + aL_w^i$  であり、各ノードの負荷状態は $S_i = L_i / W_i$  であり、

3) 負荷が最も少ないノードを選択してルーティングを行い、

10

20

30

40

50

書込み操作をパイプライン処理で行うことにより、データ入力の効率を向上させ、そのプロセスでは、先ずノード1にデータが書き込まれ、64KBのデータフラグメントが書き込まれた後、データを受信し続けると同時に、入力された64KBのデータをノード2に転送し、ノード2からノードnまでは、ノードnに最後の64KBを超えないデータフラグメントが書き込まれるまで、同じ方法でデータを受信しかつ転送し、

前記通信方法に基づいて、ノード監視モジュールはさらに、コンポーネントの処理ノードの参加、離脱、無効及び復帰というイベントを、コンポーネント管理モジュールに送信し、その場合、ノード参加イベントは、コンポーネントに処理ノードを追加することを指し、ノード離脱イベントは、コンポーネントから処理ノードを取り消すことを指し、ノード無効イベントは、コンポーネントの1つの処理ノードを使用不可とすることを指し、ノード復帰イベントは、コンポーネントの使用不可の処理ノードを復帰させて使用可とすることを指す、

10

コンポーネント指向ハイブリッドクラウドオペレーティングシステム。

#### 【請求項2】

コンポーネント指向ハイブリッドクラウドオペレーティングシステムの通信方法であって、高可用性コンポーネントクラスタ通信と、スケールアウト型コンポーネントクラスタ通信とを含み、

前記高可用性コンポーネントクラスタ通信方法は、コンポーネント管理モジュールが各コンポーネントクラスタのために読取り書込み分離待ち行列セットを構築し、前記コンポーネント管理モジュールがさらにコンポーネントのサービスタイプを冪等又は非冪等に区別し、ルータがサービスタイプによってルーティングを行い、非冪等操作は固有の書込み待ち行列に送信され、冪等操作は負荷バランスストラテジーによって別々の読取り待ち行列に送信され、これをベースに、ノード監視モジュールが処理ノードの参加、離脱、無効及び復帰というイベントを前記コンポーネント管理モジュールに送信し、前記コンポーネント管理モジュールがさらに、ノードの変化イベントによって待ち行列の構造を調整し、この方法は、読取り書込みの分離及び負荷バランスによりステータフルなコンポーネントクラスタの通信性能を向上させ、ノードの変化によって待ち行列の構造を調整することによりコンポーネントの高可用性を確保し、前記コンポーネント管理モジュールはさらに、監視モジュールにより送信されたノードの変化イベントに基づいて待ち行列の構造を調整するものであり、このプロセスは、

20

30

1) ノードが参加する時、待ち行列セットにおいて当該ノードのために1つの読取り操作の待ち行列を設定し、当該ノードはその待ち行列から読取り操作をサブスクライブし、かつ当該書込みトピックから書込みトピックをサブスクライブし、

2) ノードが離脱する時、当該ノードに対応する読取り待ち行列を削除し、当該書込みトピックに対するサブスクリプションを閉じ、

3) ノードを無効化する時、当該ノードの読取り待ち行列に読取り要求を送信することを停止し、書込みトピックに当該ノードのための書込み操作を保存し、

4) ノードが復帰する時、書込みトピックにおける書込み操作を同期し、当該ノードに対応する読取り待ち行列に読取り操作の要求を送信するために復帰し、

ステータフルなクラスタの負荷バランス方法は、各ノードに対し完全な均衡方式で操作を分散するが、それはステータフル状態において処理結果の不一致をもたらすことになり、上記方法は、読取り書込み分離待ち行列を設定しかつノードの能力に関係して負荷をバランスさせることによりこの問題を回避でき、ステータフルなコンポーネントクラスタの通信性能を向上させ、さらに、ノード変化イベントによって待ち行列の構造を調整することによりステータフルなコンポーネントクラスタの高可用性を確保し、

40

スケールアウトしたコンポーネントにより構成されたクラスタの通信方法は、前記コンポーネント管理モジュールがノード数によってコンポーネントの各サービスのためにポリゴン形待ち行列の構造を確立し、Hash値に従って二分探索によるルーティングを行い、アルゴリズムの効率は、キーワードのスケールより遙かに小さいHashパケットの数によって決定され、ルーティングの効率を向上させ、さらに、ノード数が変化する場合、

50

一部のノードデータ状態のみを調整すればよく、動的拡張の効率を向上させ、

上述した通信方法に基づき、初期化段階において、各ノードにより処理されるべき複数の H a s h インターバルをノードの処理能力によって分割し、ルーティングテーブルを形成し、比率によるデータ状態の分散を実現し、負荷バランスを達成し、

前記コンポーネント管理モジュールは、ノード数によってコンポーネントの各サービスのためにポリゴン形待ち行列セットを設定し、各待ち行列は1つの H a s h 値インターバルすなわちパケットに対応し、H a s h 値に従って二分探索によるルーティングを行い、ルーティングアルゴリズムは、

- 1) キーの H a s h 値  $h$  を計算し、
- 2) パケットの下限  $i$  を 1 に、パケットの上限  $j$  をパケット総数  $m$  に初期化し
- 3) 繰返し、
- 4) 中間パケット  $t = (i+j)/2$  を計算し、パケット  $t$  が  $h$  を含むか否かをチェックし、
- 5)  $h$  が現在のパケットの下限より小さい場合、上限  $j$  を  $j = t-1$  と更新し、
- 6)  $h$  が現在のパケットの上限より大きい場合、下限  $i$  を  $i = t+1$  と更新し、
- 7) それ以外の場合、インターバルが位置するノードに戻ってルーティングを行うものであり、

上述したアルゴリズムは、パケットをノードとして二分探索することに相当し、複雑度は  $O(\log_2 n)$  であり、このアルゴリズムの効率は、パケットの数によって決定され、パケットの数がキーワードのスケールより遙かに小さいため、ルーティング効率を向上させ、加えて、パケットストラテジーが適用されるため、ノードの数が増える場合は一部のパケットのキーワードスケール及びノードに対応するデータ状態のみを調整するだけでよく、動的拡張の効率を向上させる、

コンポーネント指向ハイブリッドクラウドオペレーティングシステムの通信方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、クラウドコンピューティング分野に関し、特に、コンポーネント指向ハイブリッドクラウドオペレーティングシステムのアーキテクチャ及びその通信方法に関する。

【背景技術】

【0002】

クラウドコンピューティングの出現に伴って、従来のデータセンタは急速にクラウドデータセンタに転換しつつある。データセンタの初期形態から発展形態への進化の4段階である、物理リソースの統合、アプリケーションと仮想化の連携、自動化管理、及び、データセンタ協調において、クラウドオペレーティングシステム (C O S : Cloud Operating System) は重要な役割を果たしており、上位に対してはアプリケーションとのインタフェース、下位に対してはハードウェアの管理という中間機能を有し、大量の異種デバイスを、クラウドアプリケーションに動的にスケジューリングされるように論理リソースプールに融合させ、端末のためのサービスを実現する。

【0003】

クラウドデータセンタ環境は、動的、不均質、大規模で単一ポイント、かつ容易に無効化されるという特徴を有する。従って、C O S は、広範かつ互換性のあるオープンアーキテクチャを適用する必要がある。サードパーティのソフトウェア及びハードウェアとの互換性を考慮するとともに二次的開発を考慮しつつ、完全な標準インタフェース A P I を提供する必要がある。クラウドコンピューティング環境の機能に対する動的変更要求に関しては、C O S は、拡張可能なコンポーネントベースの設計を適用することが必要である。その設計とは、仮想化やリソーススケジューリング等の基本的なコンポーネントに基づいて、操作とメンテナンスの管理、計量と課金、セルフサービス等のコンポーネントの付加価値開発や、要求に応じた配置を簡便に行えるものである。加えて、C O S は、クラウドコンピューティングが追求する規模拡張とサービス持続の目的を達成するために、スケーラビリティと高可用性の設計を適用する必要もある。

**【発明の概要】****【発明が解決しようとする課題】****【0004】**

クラウドデータセンタによるCOSについての疎結合、拡張可能、スケーラビリティ及び高可用性の要求に関して、従来のOSの単一モジュールアーキテクチャを適用することは、COSモジュール間の高効率の呼出しを実現するかもしれないが、結合が密であり、構造が複雑であり、かつシステムの拡張が難しい。階層化アーキテクチャを適用することは、各モジュール間の組織構造及び依存関係を明確として、COSの信頼性、移植性及びメンテナンス性を向上させるかもしれないが、ソフトウェアのスタック階層が非常に深いためにカーネルが大きすぎる上、モジュール間の結合度が依然として比較的高く、分散処理環境の構築に適さない。また、上述したアーキテクチャに基づいたオープンソースソフトウェアOpenStack及びCloudStackは、メッセージ待ち行列に基づく疎結合のクラウド管理アーキテクチャを構築するが、コンポーネント指向設計が欠落しており、コンポーネントのライフサイクルを制御できず、構成モジュールがそれ自体でスケーラビリティ及び高可用性を考慮する必要があり、このことは、それらのモジュールの開発と配置の負担及びランニングコストを増大させることになる。

10

**【0005】**

高凝集性と疎結合性という原則に従い、コンポーネントに対する管理はCOSレベルから強化すべきであり、かつ、その拡張性、スケーラビリティ及び高可用性が確保されるべきである。その場合に直面する主な問題点は、以下の通りである。

20

1. 現在のクラウドオペレーティングシステムは、自己包含性が欠落しており、構成するコンポーネントを記述しかつ管理することができず、その処理環境を動的に監視することもできない。

2. コンポーネントに対する高可用性処理用クラスタに関して、従来の通信プロトコルは、コンポーネントのステートレスな仮定に基づいて設計され、読取り書込みの分離機構及び負荷バランスストラテジーが欠落し、かつ、ステートフルなコンポーネント処理用クラスタに対して高可用性及び高性能のサポートを実現できない。

3. コンポーネントに対するスケールアウト処理用クラスタに関して、従来のツリーベースのルーティングアルゴリズムの効率は、キーワードのスケール拡大による影響を受け、一方、Hash（ハッシュ）ベースのルーティングアルゴリズムは、ノードが変化する時に大量のデータ移動をもたらすとともに、異なる機能をもつノードの負荷バランスのためのデータ分散方法が欠落している。

30

**【0006】**

従って、COSにおいて、コンポーネント及びその処理用クラスタに対する管理及び監視（スーパーバイザ）の機構をいかにして提供するか、並びに、メッセージの高効率ルーティングと負荷バランスをいかにして実現するかは、COSアーキテクチャにおいて早急に解決しなければならない課題となっている。

**【0007】**

本発明の目的は、コンポーネント指向ハイブリッドクラウドオペレーティングシステムのアーキテクチャ及びその通信方法を提供することである。

40

**【課題を解決するための手段】****【0008】**

本発明の目的は以下の方法で実現される。

コンポーネント指向ハイブリッドクラウドオペレーティングシステムは、階層モデル、オブジェクトモデル及びメッセージモデルに基づいてハイブリッドアーキテクチャを構築し、かつ、構成コンポーネント及びその処理環境を管理するためにコンポーネント指向思想を適用し、これをベースに、コンポーネント処理用クラスタに対する高効率のルーティング、読取り書込み分離及び負荷バランスを行い、クラウドオペレーティングシステムに対する開放性と互換性、疎結合と拡張性という要求を満足し、かつ、既存のクラウドオペレーティングシステムの自主管理問題、コンポーネントのスケールアウト問題及びステー

50

トフルコンポーネントの高可用性問題を解決するものである。

【0009】

階層モデルの観点から、上記システムは、上から下にポータル層、論理層、アダプテーション層及び実装層に区分され、各層は相対的に独立し、各層において標準インタフェースをそれぞれ定義することによりシステムの開放性が強化され、かつ、各層に異なる機能を適応させることにより互換性が強化される。

【0010】

オブジェクトモデルの観点から、上記クラウドオペレーティングシステムは、クラウドポータル、クラウド管理ポータル、クラウドリソース管理、監視管理、計量と課金、サービス承認、及び、権限と認証のための複数の機能モジュールにより構成され、各機能モジュールは、Restメッセージベースの呼出しにより通信を行い、互いに自由に組合わされる。そして、要求により分散され配置されることができ、要求により新規モジュールが付加価値として開発されることができ、また、プラットフォームの拡張性は、論理層における異なるモジュール間での相互操作を実現することにより強化される。

10

さらに、実施形態においては、最小形態でインストールされた上記クラウドオペレーティングシステムは、クラウドポータル、クラウド管理ポータル及びクラウドリソース管理の各機能モジュールのみにより構成される。これをベースに、監視、課金、承認又は他の機能モジュールは、要求によりカスタム化されて拡張され、Restメッセージベースの呼出しにより通信を行い、要求により分散され配置されかつ付加価値として開発される。そして上記システムの拡張性は、論理層における異なる機能モジュール間での相互操作を実現することにより強化される。

20

メッセージベースの通信方法は、非同期呼出しをサポートするために導入され、メッセージ通信インターフェースJMSがRestメッセージを転送することによりシステムアーキテクチャをさらに疎結合とするために用いられる。これをベースに、コンポーネント指向設計が適用され、かつ、コンポーネント管理ポータルがコンポーネントのメタデータを管理するとともにコンポーネントの操作状態を監視することを担当する。

【0011】

要求により拡張、分散及び配置を実現できるにも拘わらず、オブジェクトアーキテクチャは、RPC（リモート・プロセス・コール）同期通信方法に属しており、送信端は、受信端からのリターンを待った後のみ実行を継続でき、かつ、双方のプロセスが密に結合しているため、上記システムの拡大化と複雑化に伴い、コンポーネント間の連携関係が複雑過ぎるようになる。この問題に対し、オブジェクトアーキテクチャに基づいてメッセージベースの通信方法が導入される。メッセージ通信インターフェースJMSを用いてRestメッセージを転送することにより、送信端と受信端のライフタイムを異ならせ、非同期呼出しをサポートし、かつ、上記システムアーキテクチャをさらに疎結合としている。

30

そして実施形態においては、ポータル層とクラウドリソース層の間における仮想マシンのターンオン、シャットダウン、サスペンド等の操作は、非同期方式にて実現され、ポータルが、コマンドを送信した後、応答を待つ必要がなくリターンすることができる。これにより、ユーザの相互動作効果を向上させる。

【0012】

40

上述したハイブリッドアーキテクチャに基づくクラウドオペレーティングシステムは、開放性、互換性及び拡張性の要求を満足させることができる。これをベースとして、コンポーネント指向設計思想に基づき、コンポーネント管理ポータルは、コンポーネントについてのメタデータ情報を管理することを担当し、ログイン、削除、修正及びクエリ等の操作をサポートする。

【0013】

その場合、1つのコンポーネントは、1つの三要素セットであり、ネーム、サービスセット、アクセスアドレス、及び、記述を含む。

【0014】

1つのサービスは、1つの四要素セットであり、ネーム、タイプ、メッセージプロトコ

50

ル、パラメタリスト、キーネーム、機能記述、及び、非機能記述を含む。

【 0 0 1 5 】

コンポーネントの記述及び管理に加えて、コンポーネント管理モジュールはさらに、その処理環境を監視し、コンポーネントのスケラビリティ及び可用性を確保するための基本サービスを提供し、かつ、クラウドオペレーティングシステムの自主管理能力を完備させる。そして、コンポーネントがログインした時、システムは、ユーザ名user、パスワードpsw、及び、固有のコンポーネントIDidを割当て、その後のコンポーネント処理用クラスタのアクセスのプロセスは、次の通りとなる。

【 0 0 1 6 】

1) 処理用クラスタが、アドレスがurlであるシステムバスに対してアクセス要求を開始し、システムバスがアクセスノードのユーザ名、パスワード及びIDを認証し、認証が通ると、接続が確立される。そのコードは以下の通りである。

```
Connection = ConnectionFactory.createConnection(user,psw,url);
```

2) 書き込み操作トピックを設定し、コンポーネントの各処理ノードが当該トピックから書き込み操作をサブスクライブする。

```
write_topic = session.createTopic(id+"WRITE_TOPIC");
```

```
write_topic_consumer = session.createConsumer(write_topic);
```

3) コンポーネントがwriteTopicListenerのonMessage方法において書き込み処理を実現し、かつシステムバスにログインする。

```
write_topic_consumer.setMessageListener(writeTopicListener);
```

4) コンポーネントの処理ノード数numに従って読取り操作の待ち行列セットを設定し、各処理ノードが1つの待ち行列に対応した読取り操作をサブスクライブする。

```
read_queue = session.createMultiQueue(num, id+"READ_QUEUE");
```

```
read_queue_consumer = session.createConsumer(read_queue);
```

5) readQueueListenerのonMessage方法において特定の読取り処理機能を実現し、かつシステムバスにログインする。

```
read_queue_consumer.setMessageListener(readQueueListener);
```

【 0 0 1 7 】

上述した書き込み読取り分離の待ち行列セットに基づき、コンポーネント管理モジュールにおいて各コンポーネントのサービスタイプを区別し、サービスタイプを冪等又は非冪等に設定する。冪等操作はステートレス操作に属し、同一状態下で実行した結果が毎回同一である。非冪等操作はステートフル操作に属し、同一状態下で実行した結果が毎回異なる。ルータは、サービスタイプによってルーティングを行い、非冪等操作は固有の書き込み待ち行列に送信され、そして冪等操作は負荷バランスストラテジーによって別々の読取り待ち行列に送信される。

【 0 0 1 8 】

読取り操作の負荷バランスのプロセスは、以下の通りである。

1) ノードの処理能力を計算する。

ノードiのCPUクロック周波数、メモリ容量及びI/O帯域幅を各々 $C_i$ 、 $M_i$ 及び $B_i$ とすると、クラスタの各種リソースはノードの各種リソースの総計であり、つまり $C = \sum C_i$ 、 $M = \sum M_i$ 、 $B = \sum B_i$ である。

ノードiのCPU重み付け値は $W_i^{CPU} = C_i / C$ であり、メモリ容量重み付け値は $W_i^{RAM} = M_i / M$ であり、I/O帯域幅重み付け値は $W_i^{IO} = B_i / B$ である。

コンポーネントサービスに必要なリソース割合を、それぞれ $p^{CPU}$ 、 $p^{RAM}$ 、 $p^{IO}$ とすると、ノードiの処理能力は、

$$W_i = p^{CPU} W_i^{CPU} + p^{RAM} W_i^{RAM} + p^{IO} W_i^{IO}$$

2) 書き込み読取り操作の重み付け値によって各ノードの負荷を計算する。

読取り待ち行列 $L_r$ の書き込み待ち行列 $L_w$ に対する書き込み読取り操作のオーバーヘッド比率をaとすると、ノードiの負荷は $L_i = L_r^i + aL_w^i$ であり、各ノードの負荷状態は $S_i = L_i / W_i$ である。

10

20

30

40

50



3) 負荷が最も少ないノードを選択してルーティングを行う。

書込み操作をパイプライン処理で行うことにより、データ入力の効率を向上させる。そのプロセスでは、先ずノード1にデータが書き込まれ、64KBのデータフラグメントが書き込まれた後、データを受信し続けると同時に、入力された64KBのデータをノード2に転送し、ノード2からノードnまでは、ノードnに最後の64KBを超えないデータフラグメントが書き込まれるまで、同じ方法でデータを受信しかつ転送する。

【0019】

上記通信方法に基づいて、ノード監視モジュールはさらに、コンポーネントの処理ノードの参加、離脱、無効及び復帰というイベントを、コンポーネント管理モジュールに送信し、その場合、ノード参加イベントは、コンポーネントに処理ノードを追加することを指す。ノード離脱イベントは、コンポーネントから処理ノードを取り消すことを指す。ノード無効イベントは、コンポーネントの1つの処理ノードを使用不可とすることを指す。ノード復帰イベントは、コンポーネントの使用不可の処理ノードを復帰させて使用可とすることを指す。

10

【0020】

本発明の実施形態はまた、コンポーネント指向ハイブリッドクラウドオペレーティングシステムの通信方法を提供し、それは高可用性コンポーネントのクラスタ通信方法と、スケールアウトしたコンポーネントのクラスタ通信方法とを含む。

【0021】

その場合、高可用性コンポーネントのクラスタ通信方法は、コンポーネント管理モジュールが各コンポーネントクラスタのために読取り書込み分離の待ち行列セットを構築し、コンポーネント管理モジュールがさらにコンポーネントのサービスタイプを冪等又は非冪等に区別し、ルータがサービスタイプによってルーティングを行い、非冪等操作は固有の書込み待ち行列に送信され、冪等操作は負荷バランスストラテジーによって別々の読取り待ち行列に送信される。そしてこれをベースに、ノード監視モジュールがコンポーネントの処理ノードの参加、離脱、無効及び復帰というイベントをコンポーネント管理モジュールに送信し、コンポーネント管理モジュールがさらに、ノードの変化イベントによって待ち行列の構造を調整する。

20

この方法は、読取り書込みの分離及び負荷バランスによりステートフルなコンポーネントクラスタの通信性能を向上させるとともに、ノードの変化によって待ち行列の構造を調整することによりコンポーネントの高可用性を確保する。コンポーネント管理モジュールはさらに、監視モジュールにより送信されたノードの変化イベントに基づいて待ち行列の構造を調整する。このプロセスは以下のとおりである。

30

【0022】

1) ノードが参加する時、待ち行列セットにおいて当該ノードのために1つの読取り操作の待ち行列を設定し、当該ノードはその待ち行列から読取り操作をサブスクライブし、かつ当該書込みトピックから書込みトピックをサブスクライブする。

2) ノードが離脱する時、当該ノードに対応する読取り待ち行列を削除し、当該書込みトピックに対するサブスクリプションを閉じる。

3) ノードを無効化する時、当該ノードの読取り待ち行列に読取り要求を送信することを停止し、書込みトピックに当該ノードのための書込み操作を保存する。

40

4) ノードが復帰する時、書込みトピックにおける書込み操作を同期し、当該ノードに対応する読取り待ち行列に読取り操作の要求を送信するために復帰する。

【0023】

ステートレスなクラスタの負荷バランス方法は、各ノードに対し完全な均衡方式で操作を分散するが、それはステートフル状態において処理結果の不一致をもたらすことになる。上記方法は、読取り書込み分離の待ち行列を設定しかつノードの能力に関係して負荷をバランスさせることによりこの問題を回避でき、ステートフルのコンポーネントクラスタの通信性能を向上させる。さらに、ノード変化イベントによって待ち行列の構造を調整することによりステートフルのコンポーネントクラスタの高可用性を確保する。

50

## 【0024】

スケールアウトしたコンポーネントにより構成されたクラスタの通信方法は、コンポーネント管理モジュールが、ノード数によってコンポーネントの各サービスのためにポリゴン形待ち行列の構造を確立し、H a s h 値に従って二分探索によるルーティングを行う。アルゴリズムの効率は、キーワードのスケールより遙かに小さいH a s h バケットの数によって決定され、ルーティングの効率を向上させる。さらに、ノード数が増加する場合、一部のノードデータ状態のみを調整すればよく、動的拡張の効率を向上させる。

上述した通信方法に基づき、初期化段階において、各ノードにより処理されるべき複数のH a s h インターバルをノードの処理能力によって分割し、ルーティングテーブルを形成し、比率によるデータ状態の分散を実現し、負荷バランスを達成する。

10

## 【0025】

コンポーネント管理モジュールは、ノード数によってコンポーネントの各サービスのためにポリゴン形待ち行列セットを設定し、各待ち行列は1つのH a s h 値インターバルすなわちパケットに対応し、H a s h 値に従って二分探索によるルーティングを行う。ルーティングアルゴリズムは、以下の通りである。

- 1) キーのH a s h 値 $h$ を計算する。
- 2) パケットの下限 $i$ を1に、パケットの上限 $j$ をパケット総数 $m$ に初期化する。
- 3) 繰返す。
- 4) 中間パケット $t = (i+j)/2$ を計算し、パケット $t$ が $h$ を含むか否かをチェックする。
- 5)  $h$ が現在のパケットの下限より小さい場合、上限 $j$ を $j = t-1$ と更新する。
- 6)  $h$ が現在のパケットの上限より大きい場合、下限 $i$ を $i = t+1$ と更新する。
- 7) それ以外の場合、インターバルが位置するノードに戻ってルーティングを行う。

20

## 【0026】

上述したアルゴリズムは、パケットをノードとして二分探索することに相当する。複雑度は $O(\log_2 n)$ である。このアルゴリズムの効率は、パケットの数によって決定される。パケットの数がキーワードのスケールより遙かに小さいため、ルーティング効率を向上させる。加えて、パケットストラテジーが適用されるため、ノードの数が増加する場合は一部のパケットのキーワードスケール及びノードに対応するデータ状態のみを調整するだけでよく、動的拡張の効率を向上させる。

## 【発明の効果】

30

## 【0027】

本発明の有益な効果は以下の通りである。既存の技術に比べ、本発明で提供されたコンポーネント指向ハイブリッドアーキテクチャは、コンポーネント化されたクラウドオペレーティングシステムアーキテクチャを完全なものとする。そのクラウドオペレーティングシステムは、開放性と互換性、拡張性と疎結合、及び、スケーラビリティと高可用性を備えており、それらは、スケールアウトしたコンポーネントクラスタ及び高可用性クラスタの通信方法により確保される。

## 【図面の簡単な説明】

## 【0028】

【図1】図1は、コンポーネント指向ハイブリッドC O S アーキテクチャの構成図である。

40

【図1 C o n t .】図1 Cont . ( 図1の続き ) は、コンポーネント指向ハイブリッドC O S アーキテクチャの構成図である。

【図2】図2は、ステートフルコンポーネントクラスタをサポートする高可用性通信アーキテクチャの構成図である。

【図3】図3は、読取り書込み分離待ち行列セットの構成図である。

【図4】図4は、データを書込み時のシーケンスを示す図である。

【図5】図5は、ノード状態の遷移を示す図である。

【図6】図6は、待ち行列の調整プロセスを示すフロー図である。

【図7】図7は、スケールアウトしたコンポーネント処理用クラスタの通信アーキテクチャ

50

の構成図である。

【図 8】図 8 は、ポリゴン形待ち行列セットを備えたルーティングを示す概略構成図である。

【図 9】図 9 は、ポリゴン形待ち行列セットを備えたルーティングアルゴリズムのプロセスを示すフロー図である。

【発明を実施するための形態】

【0029】

図面及び実施例とともに、本発明の実施形態を以下に詳細に説明する。これにより、本発明が技術的課題を解決しかつ技術的効果を達成するためにいかにして技術手段を適用するのかについての実現プロセスを十分に理解できかつこれにより実施できる。矛盾を生じない限りにおいて、本明細書に開示に実施形態及びそれらの実施形態の多様な特徴は、いずれも開示される発明の保護範囲に含まれるものである。

10

【0030】

[1] コンポーネント指向ハイブリッドクラウドオペレーティングシステムのアーキテクチャ

階層モデルの観点から、上記システムは、上から下にポータル層、論理層、アダプテーション層及び実装層に区分され、各層は相対的に独立し、各層において標準インタフェースをそれぞれ定義することによりシステムの開放性を強化し、かつ、各層に異なる機能を適応させることにより互換性を強化する。実施形態においては、ポータル層を取り除くことにより、ユーザインタフェースUIと機能論理の分離を実現でき、例えば、論理層が外部に対して統一標準の Rest API を提供することによりポータルすなわちサードパーティの二次的開発をサポートする。また、論理的機能層を抽象化することにより、多様な具体的な機能実装における互換性を実現でき、例えば、クラウドリソース管理モジュールが、仮想化されたアダプタにより多様な仮想化されたインフラストラクチャをサポートする。監視管理モジュールは、アダプテーションフレームにより多様な監視プロトコルと互換性がある。プロセス管理モジュールは、プロセスエンジンにより異なる承認プロセスのカスタム化をサポートする。そして、課金管理及び権限と認証モジュールは、保存されたフックインタフェースによりプラグインアクセスをサポートする。

20

【0031】

オブジェクトモデルの観点から、クラウドオペレーティングシステムは、クラウドポータル、クラウド管理ポータル、クラウドリソース管理、監視管理、計量と課金、サービス承認、権限と認証、等の機能モジュールにより構成されている。各機能コンポーネントが Rest メッセージベースの呼出しにより通信を行い、自由に組合せることができ、要求により分散されかつ配置される。さらに、要求により新規モジュールを付加価値的に開発でき、論理層における異なるモジュール間で相互動作を実現することによりプラットフォームの拡張性を強化する。実施形態では、最小形態でインストールされたクラウドオペレーティングシステムは、クラウドポータル、クラウド管理ポータル及びクラウドリソース管理のモジュールのみから構成され、それに基づき、監視、課金、承認又は他のモジュールを、要求によりカスタム化し拡張してもよい。

30

【0032】

要求により拡張、分散及び配置を実現できるにも拘わらず、オブジェクトアーキテクチャは、RPC (リモート・プロセス・コール) 同期通信方法に属しており、送信端は、受信端からの応答を待った後にのみ実行を継続することができ、かつ、双方のプロセスが密に結合しているため、システムの拡大化と複雑化に伴い、コンポーネント間の連携関係が複雑過ぎるようになる。この問題に対し、オブジェクトアーキテクチャに基づいてメッセージベースの通信方法を導入し、メッセージ通信インターフェース JMS を使用して Rest メッセージを転送することにより、送信端と受信端のライフタイムが異なってもよく、非同期呼出しをサポートし、かつ、システムアーキテクチャをさらに疎結合とする。実施形態において、ポータル層とクラウドリソース層の間における仮想マシンのターンオン、ターンオフ、サスペンド、ターンオフ等の操作は非同期方式により実現され、ポ

40

50

タルは、コマンドの送信後、応答を待つ必要がなくリターンできる。このことは、ユーザの相互動作効果を向上させる。

【 0 0 3 3 】

上述したハイブリッドアーキテクチャに基づくクラウドオペレーティングシステムは、開放性、互換性及び拡張性の要求を満足させることができる。これをベースとして、コンポーネント指向設計思想に基づき、コンポーネント管理ポータルは、コンポーネントについてのメタデータ情報を管理することを担当し、ログイン、削除、修正及びクエリ等の操作をサポートする。

【 0 0 3 4 】

この場合、1つのコンポーネントは、1つの三要素セット{ネーム、サービスセット、アクセスアドレス、記述}である。 10

1つのサービスは、1つの四要素セット{ネーム、タイプ、メッセージプロトコル、パラメタリスト、キーネーム、機能記述、及び、非機能記述}である。

【 0 0 3 5 】

コンポーネントに対する記述及び管理に加えて、コンポーネント管理モジュールはさらに、その処理環境を監視し、コンポーネントのスケラビリティ及び可用性を確保するための基礎的サービスを提供し、クラウドオペレーティングシステムの自主管理能力を完備させる。

【 0 0 3 6 】

[ 2 ] 高可用性コンポーネントのクラスタ通信方法 20

本発明による高可用性コンポーネントのクラスタの実施形態は、図 2 に示す通りである。これは、主に以下の各モジュールを含む。

【 0 0 3 7 】

コンポーネント管理モジュールは、コンポーネント情報に従って、コンポーネントのメッセージ待ち行列セットを設定し、削除し、かつ調整することを担当する。

【 0 0 3 8 】

メッセージルータは、ルーティング情報に従って、待ち行列セットにメッセージを分散することを担当する。

【 0 0 3 9 】

ノード監視モジュールは、コンポーネント処理用クラスタの各ノードの参加及び離脱を検出し、処理ノードについてのリソースコンフィギュレーション情報を獲得することを担う。 30

【 0 0 4 0 】

コンポーネント処理用クラスタは、コンポーネントの具体的なサービス機能を実現することを担当し、いくつかの処理ノードから構成されている。

【 0 0 4 1 】

コンポーネントクライアントは、コンポーネントサービスに対する使用要求を開始することを担当する。

【 0 0 4 2 】

上述したアーキテクチャに基づき、コンポーネントがログオンしたとき、システムは、コンポーネントに対してユーザ名user、パスワードpsw及び固有のコンポーネントIDidを割当てて。その後、コンポーネント処理用クラスタのアクセス過程は、次の通りとなる。 40

【 0 0 4 3 】

1) 処理用クラスタは、アドレスがurlであるシステムバスに対してアクセス要求を開始し、システムバスがアクセスノードのユーザ名、パスワード及びIDを認証し、認証が通ると、接続が確立される。そのコードは以下の通りである。

```
Connection = ConnectionFactory.createConnection(user,psw,url);
```

【 0 0 4 4 】

2) 書き込み操作トピックを設定し、コンポーネントの各処理ノードが当該トピックから 50

書込み操作をサブスクライブする。

```
write_topic = session.createTopic(id+"WRITE_TOPIC");
write_topic_consumer = session.createConsumer(write_topic);
```

【 0 0 4 5 】

3) コンポーネントがwriteTopicListenerのonMessage方法において書込み処理を実現し、かつシステムバスにログインする。

```
write_topic_consumer.setMessageListener(writeTopicListener);
```

【 0 0 4 6 】

4) コンポーネントの処理ノード数numに従って読取り操作の待ち行列セットを設定し、各処理ノードが1つの待ち行列に対応した読取り操作をサブスクライブする。

```
read_queue = session.createMultiQueue(num, id+"READ_QUEUE");
read_queue_consumer = session.createConsumer(read_queue);
```

【 0 0 4 7 】

5) readQueueListenerのonMessage方法において特定の読取り処理機能を実現し、かつシステムバスにログインする。

```
read_queue_consumer.setMessageListener(readQueueListener);
```

【 0 0 4 8 】

上述した書込み読取り分離待ち行列セットに基づき、コンポーネント管理モジュールにおいて各コンポーネントのサービスタイプを区別し、サービスタイプを冪等又は非冪等に設定する。冪等操作はステートレス操作に属し、同一状態下で実行した結果が毎回同一である。非冪等操作はステートフル操作に属し、同一状態下で実行した結果が毎回異なる。図3に示すように、ルータは、サービスタイプに従ってルーティングを行い、非冪等操作は固有の書込み待ち行列に送信され、そして冪等操作は負荷バランスストラテジーによって別々の読取り待ち行列に送信される。

【 0 0 4 9 】

その場合、読取り操作負荷バランスプロセスのフローは、次の通りである。

1) ノードの処理能力を計算する。

ノード $i$ のCPUクロック周波数、メモリ容量及びI/O帯域幅を各々 $C_i$ 、 $M_i$ 及び $B_i$ とすると、クラスタの各種リソースはノードの各種リソースの総計である。つまり $C = \sum C_i$ 、 $M = \sum M_i$ 、 $B = \sum B_i$ である。

そして、ノード $i$ のCPU重み付け値は $W_i^{CPU} = C_i / C$ であり、メモリ容量重み付け値は $W_i^{RAM} = M_i / M$ であり、I/O帯域幅重み付け値は $W_i^{IO} = B_i / B$ である。

コンポーネントサービスに必要なリソース割合を、それぞれ $p^{CPU}$ 、 $p^{RAM}$ 、 $p^{IO}$ とすると、ノード $i$ の処理能力は、

$$W_i = p^{CPU} W_i^{CPU} + p^{RAM} W_i^{RAM} + p^{IO} W_i^{IO}$$

【 0 0 5 0 】

2) 書込み読取り操作の重み付け値によって各ノードの負荷を計算する。

読取り待ち行列 $L_r$ の書込み待ち行列 $L_w$ に対する書込み読取り操作のオーバーヘッド比率を $a$ とすると、ノード $i$ の負荷は $L_i = L_r + aL_w$ であり、各ノードの負荷状態は $S_i = L_i / W_i$ である。

【 0 0 5 1 】

3) 負荷が最も少ないノードを選択してルーティングを行う。

書込み操作をパイプライン処理を行うことにより、データ入力の効率を向上させる。図4に示すように、そのプロセスでは、先ずノード1にデータが書き込まれ、64KBのデータフラグメントが書き込まれた後、データを受信し続けると同時に、入力された64KBのデータをノード2に転送し、ノード2からノード $n$ までは、ノード $n$ に最後の64KBを超えないデータフラグメントが書き込まれるまで、同じ方法でデータを受信しかつ転送する。

【 0 0 5 2 】

上記通信方法に基づいて、ノード監視モジュールはさらに、コンポーネントの処理ノードの参加、離脱、無効及び復帰というイベントを、コンポーネント管理モジュールに送信

10

20

30

40

50

する。ノードの状態遷移関係は、図5に示す通りである。その場合、ノード参加イベントは、コンポーネントに処理ノードを追加することを指す。ノード離脱イベントは、コンポーネントから処理ノードを取り消すことを指す。ノード無効イベントは、コンポーネントの1つの処理ノードを使用不可とすることを指す。ノード復帰イベントは、コンポーネントの使用不可の処理ノードを復帰させて使用可とすることを指す。

#### 【0053】

コンポーネント管理モジュールはさらに、監視モジュールにより送信されたノードの変化イベントに基づいて待ち行列の構造を調整する。このフローは図6に示す通りである。

1) ノードが参加する時、待ち行列セットにおいて当該ノードのために1つの読取り操作の待ち行列を設定し、当該ノードはその待ち行列から読取り操作をサブスクライブし、かつ当該書込みトピックから書込みトピックをサブスクライブする。

2) ノードが離脱する時、当該ノードに対応する読取り待ち行列を削除し、当該書込みトピックに対するサブスクリプションを閉じる。

3) ノードを無効化する時、当該ノードの読取り待ち行列に読取り要求を送信することを停止し、書込みトピックに当該ノードのための書込み操作を保存する。

4) ノードが復帰する時、書込みトピックにおける書込み操作を同期し、当該ノードに対応する読取り待ち行列に読取り操作の要求を送信するために復帰する。

#### 【0054】

ステートレスなクラスタの負荷バランス方法は、各ノードに対し完全な均衡方式で操作を分散するが、それはステートフル状態において処理結果の不一致をもたらすことになる。上記方法は、読取り書込み分離待ち行列を設定しかつノードの能力に関係して負荷をバランスさせることによりこの問題を回避でき、ステートフルのコンポーネントクラスタの通信性能を向上させる。さらに、ノード変化イベントによって待ち行列の構造を調整することによりステートフルのコンポーネントクラスタの高可用性を確保する。

#### 【0055】

[3] スケールアウトしたコンポーネントにより構成されたクラスタの通信方法

本発明の実施形態は、図7に示す通り、スケールアウトしたコンポーネントにより構成されたクラスタの通信方法を提供する。コンポーネント管理モジュールは、ノード数に従ってコンポーネントの各サービスのためにポリゴン形待ち行列の構造を確立する。各待ち行列は、1つのHash値インターバル(パケット)に対応する。そして、Hash値に従って二分探索によるルーティングが行われる。ルーティングアルゴリズムは、図9に示す通りである。

#### 【0056】

1) キーのHash値 $h$ を計算する。

2) パケットの下限 $i$ を1に、パケットの上限 $j$ をパケット総数 $m$ に初期化する。

3) 繰返す。

4) 中間パケット $t = (i+j)/2$ を計算し、パケット $t$ が $h$ を含むか否かをチェックする。

5)  $h$ が現在のパケットの下限より小さい場合、上限 $j$ を $j = t-1$ と更新する。

6)  $h$ が現在のパケットの上限より大きい場合、下限 $i$ を $i = t+1$ と更新する。

7) それ以外の場合、インターバルが位置するノードに戻ってルーティングを行う。

#### 【0057】

上述したアルゴリズムは、パケットをノードとして二分探索することに相当する。複雑度は $O(\log_2 n)$ である。このアルゴリズムの効率は、パケットの数によって決定される。パケットの数がキーワードのスケールより遙かに小さいため、ルーティング効率を向上させる。加えて、パケットストラテジーが適用されるため、ノードの数が増える場合、一部のパケットのキーワードスケール及びノードに対応するデータ状態のみを調整するだけでよく、動的拡張の効率を向上させる。

#### 【0058】

上記通信方法に基づき、本発明の実施形態は、データ状態の分散方法を提供する。初期化段階において、各ノードにより処理される必要のあるHashインターバルが、ノード

処理能力に従って分割される。ノード処理能力の計算方法は上述した通りである。この方法は、ノード処理能力に従って各ノードの待ち行列により処理される必要のある Hash インターバルスケールを分割することができ、比率によってデータ状態を分散することができ、負荷バランスを達成できる。

【 0 0 5 9 】

本発明の他の特徴及び利点は、その後の明細書で説明されかつ一部は明細書から自明である。あるいは本発明を実行することにより理解されるであろう。本発明の目的及び他の利点は明細書、請求の範囲及び図面に示された構成により実現されかつ得ることができる。

【 図 1 】

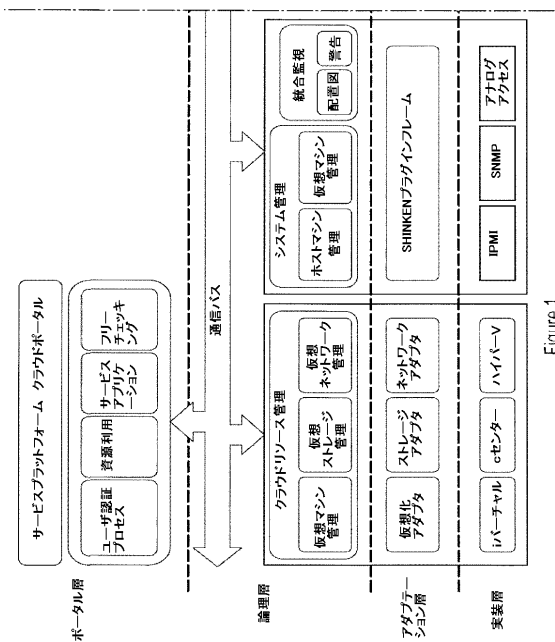


Figure 1

【 図 1 Cont . 】

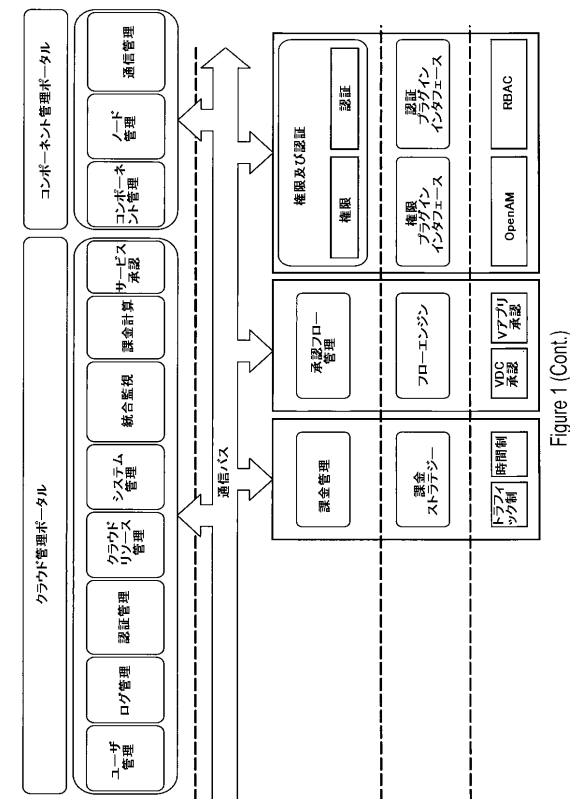


Figure 1 (Cont.)

【 図 2 】

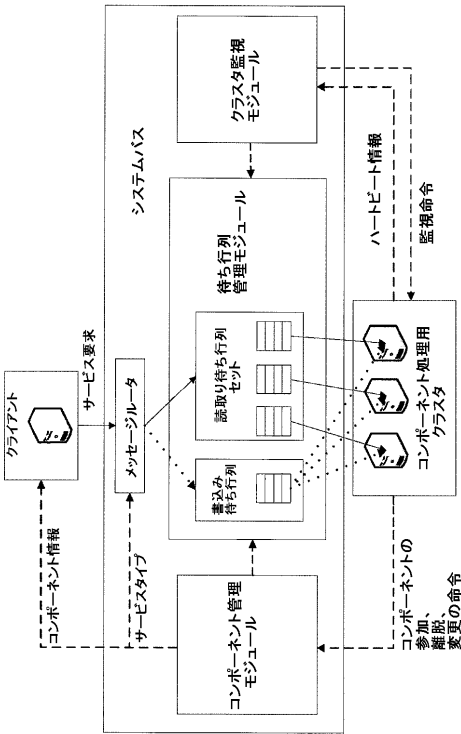


Figure 2

【 図 3 】

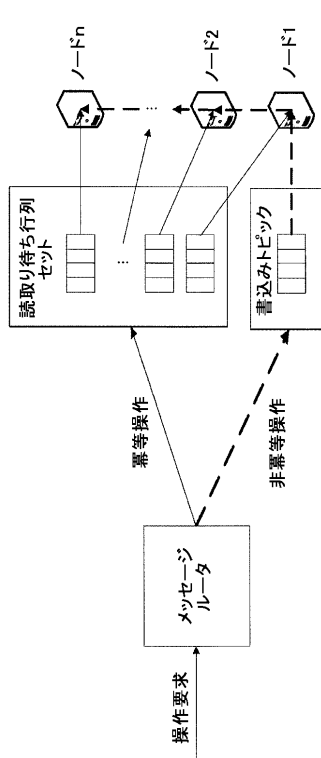


Figure 3

【 図 4 】

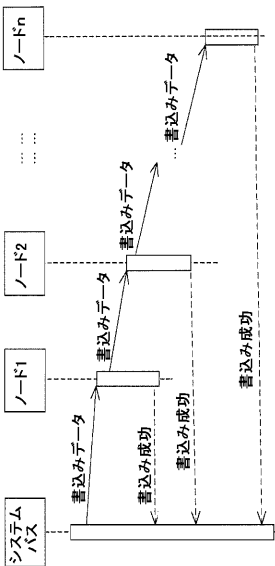


Figure 4

【 図 5 】

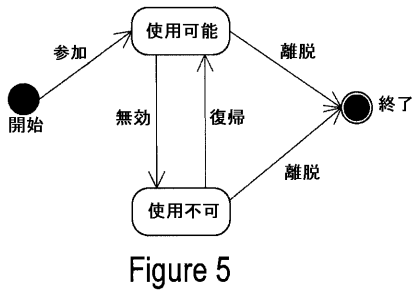


Figure 5

【 図 6 】

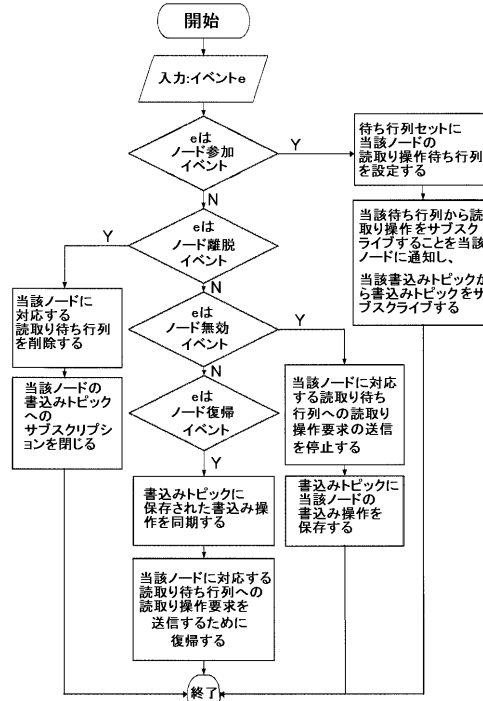


Figure 6



【 図 7 】

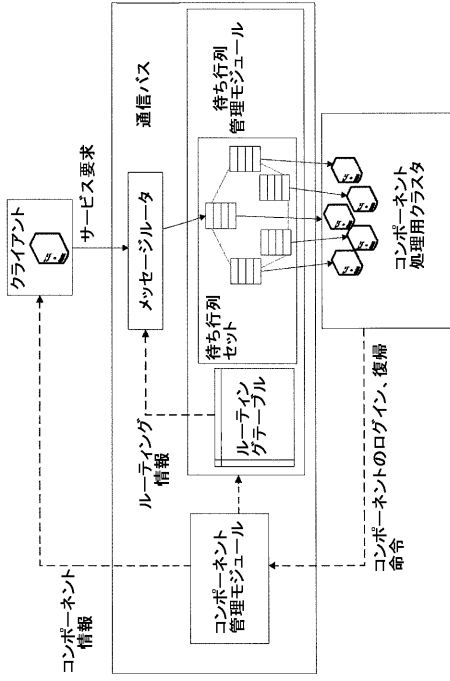
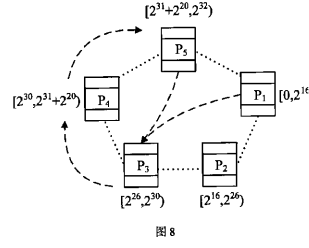


Figure 7

【 図 8 】



【 図 9 】

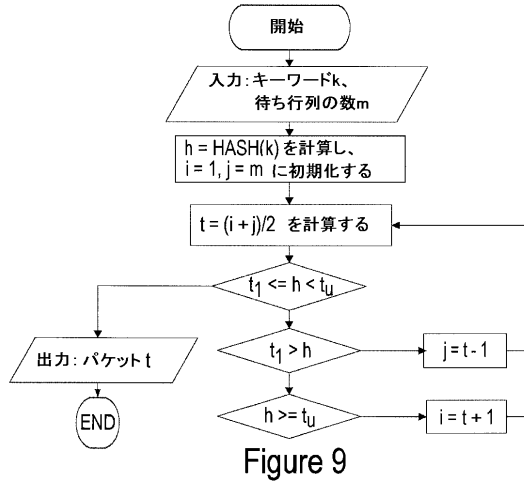


Figure 9

## 【 国际調查報告 】

| <b>INTERNATIONAL SEARCH REPORT</b>   |  | International application No.<br>PCT/CN2014/071233 |
|--|--|--|
| <b>A. CLASSIFICATION OF SUBJECT MATTER</b>   |  |  |
| H04L 29/08 (2006.01) i; G06F 9/455 (2006.01) n<br>According to International Patent Classification (IPC) or to both national classification and IPC  |  |  |
| <b>B. FIELDS SEARCHED</b>  |  |  |
| Minimum documentation searched (classification system followed by classification symbols)  |  |  |
| H04L; 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)   |  |  |
| CNABS, CNTXT, CNKI, VEN: cloud operation system, cloud operating system, cloud os, cos,<br>component+, object, module, layer   |  |  |
| <b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>  |  |  |
| Category*  | Citation of document, with indication, where appropriate, of the relevant passages   | Relevant to claim No.                              |
| PX   | CN 103442049 A (INSUR ELECTRONIC INFORMATION IND CO., LTD.), 11 December 2013 (11.12.2013), claims 1 and 2   | 1, 2   |
| A  | LI, Mingjie, INSPUR IN-CLOUD OS V2.0 PRODUCT ANALYSIS(VOLUME 1), Science & Technology Inspur, 2012 No.4, 30 April 2012 (30.04.2012), pages 30 and 31   | 1, 2   |
| A  | CN 102521022 A (MICROSOFT CORP.) 27 June 2012 (27.06.2012), the whole document   | 1, 2   |
| A  | CN 102970332 A (JIANGSU INTERNET THINGS RES DEV CENT) 13 March 2013 (13.03.2013), the whole document   | 1, 2   |
| <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:<br>"A" document defining the general state of the art which is not considered to be of particular relevance<br>"E" earlier application or patent but published on or after the international filing date<br>"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)<br>"O" document referring to an oral disclosure, use, exhibition or other means<br>"P" document published prior to the international filing date but later than the priority date claimed | "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention<br>"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<br>"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<br>"&" document member of the same patent family |  |
| Date of the actual completion of the international search<br>15 April 2014 (15.04.2014)  | Date of mailing of the international search report<br>06 May 2014 (06.05.2014)   |  |
| Name and mailing address of the ISA/CN<br>State Intellectual Property Office of the P. R. China<br>No. 6, Xitucheng Road, Jimenqiao<br>Haidian District, Beijing 100088, China<br>Facsimile No. (86-10) 62019451   | Authorized officer<br><br>NIE, Jincheng<br><br>Telephone No. (86-10) 62412155  |  |

**INTERNATIONAL SEARCH REPORT**  
**Information on patent family members**

International application No.  
PCT/CN2014/071233

| Patent Documents referred in the Report | Publication Date | Patent Family    | Publication Date  |
|---|------------------|------------------|-------------------|
| CN 103442049 A                          | 11 December 2013 | None             |                   |
| CN 102521022 A                          | 27 June 2012     | CA 2815306 A1    | 03 May 2012       |
|   |                  | US 2012110570 A1 | 03 May 2012       |
|   |                  | JP 2013541113 A  | 07 November 2013  |
|   |                  | AU 2011320899 A1 | 02 May 2013       |
|   |                  | EP 2633400 A1    | 04 September 2013 |
|   |                  | WO 2012057955 A1 | 03 May 2012       |
| CN 102970332 A                          | 13 March 2013    | None             |                   |

| 国际检索报告  |  | 国际申请号<br>PCT/CN2014/071233          |
|---|--|-------------------------------------|
| A. 主题的分类<br>H04L 29/08(2006.01)i; G06F 9/445(2006.01)n<br>按照国际专利分类(IPC)或者同时按照国家分类和IPC两种分类   |  |                                     |
| B. 检索领域<br>检索的最低限度文献(标明分类系统和分类号)<br>H04L; G06F<br>包含在检索领域中的除最低限度文献以外的检索文献<br>在国际检索时查阅的电子数据库(数据库的名称, 和使用的检索词(如使用))<br>CNABS, CNTXT, CNKI, VEN; 云操作系统, 云OS, 构件, 组件, 组成, 对象, 模型, 层, cloud operation system, cloud operating system, cloud OS, CCS, component, object, module, layer  |  |                                     |
| C. 相关文件   |  |                                     |
| 类 型*  | 引用文件, 必要时, 指明相关段落  | 相关的权利要求                             |
| PX  | CN 103442049A ((浪潮电子信息产业股份有限公司)) 2013年 12月 11日 (2013-12-11)<br>权利要求1-2               | 1-2                                 |
| A   | 李明杰. "浪潮云海OS V2.0产品解析(上)"<br>科技浪潮, 卷 2012年第4期, 2012年 4月 30日 (2012-04-30),<br>第30-31页 | 1-2                                 |
| A   | CN 102521022A ((微软公司)) 2012年 6月 27日 (2012-06-27)<br>全文                               | 1-2                                 |
| A   | CN 102970332A ((江苏物联网研究发展中心)) 2013年 3月 13日 (2013-03-13)<br>全文                        | 1-2                                 |
| <input type="checkbox"/> 其余文件在C栏的续页中列出。 <input checked="" type="checkbox"/> 见同族专利附件。  |  |                                     |
| * 引用文件的具体类型:<br>"A" 认为不特别相关的表示了现有技术一般状态的文件<br>"E" 在国际申请日的当天或之后公布的在先申请或专利<br>"L" 可能对优先权要求构成怀疑的文件, 或为确定另一篇引用文件的公布日而引用的或者其他特殊理由而引用的文件(如具体说明的)<br>"O" 涉及口头公开、使用、展览或其他方式公开的文件<br>"P" 公布日先于国际申请日但迟于所要求的优先权日的文件<br>"T" 在申请日或优先权日之后公布, 与申请不相抵触, 但为了理解发明之理论或原理的在后文件<br>"X" 特别相关的文件, 单独考虑该文件, 认定要求保护的发明不是新颖的或不具有创造性<br>"Y" 特别相关的文件, 当该文件与另一篇或者多篇该类文件结合并且这种结合对于本领域技术人员为显而易见时, 要求保护的发明不具有创造性<br>"&" 同族专利的文件 |  |                                     |
| 国际检索实际完成的日期<br>2014年 4月 15日   |  | 国际检索报告邮寄日期<br>2014年 5月 06日          |
| ISA/CN的名称和邮寄地址<br>中华人民共和国国家知识产权局(ISA/CN)<br>中国北京市海淀区蓟门桥西土城路6号<br>100088 中国<br>传真号 (86-10)62019451   |  | 受权官员<br>袁锦程<br>电话号码 (86-10)62412155 |

表 PCT/ISA/210 (第2页) (2009年7月)

国际检索报告  
关于同族专利的信息

国际申请号

PCT/CN2014/071233

| 检索报告引用的专利文件 |            | 公布日<br>(年/月/日) | 同族专利 |              | 公布日<br>(年/月/日) |
|-------------|------------|----------------|------|--------------|----------------|
| CN          | 103442049A | 2013年 12月 11日  | 无    |              | 无年             |
| CN          | 102521022A | 2012年 6月 27日   | CA   | 2815306A1    | 2012年 5月 03日   |
|             |            |                | US   | 2012110570A1 | 2012年 5月 03日   |
|             |            |                | JP   | 2013541113A  | 2013年 11月 07日  |
|             |            |                | AU   | 2011320899A1 | 2013年 5月 02日   |
|             |            |                | EP   | 2633400A1    | 2013年 9月 04日   |
|             |            |                | WO   | 2012057955A1 | 2012年 5月 03日   |
| CN          | 102970332A | 2013年 3月 13日   | 无    |              | 无年             |

表 PCT/ISA/210 (同族专利附件) (2009年7月)

## フロントページの続き

(81) 指定国 AP(BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, RU, TJ, TM), EP(AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JP, KE, KG, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US

(72) 発明者 ワング、エンドング

中華人民共和国、250014 シャンドング、ジナン、ハイ - テク・ゾーン、シュンヤ・ロード、No. 1036

(72) 発明者 チャング、ドング

中華人民共和国、250014 シャンドング、ジナン、ハイ - テク・ゾーン、シュンヤ・ロード、No. 1036

(72) 発明者 リウ、ズヘングウェイ

中華人民共和国、250014 シャンドング、ジナン、ハイ - テク・ゾーン、シュンヤ・ロード、No. 1036

(72) 発明者 クイ、カイユアン

中華人民共和国、250014 シャンドング、ジナン、ハイ - テク・ゾーン、シュンヤ・ロード、No. 1036

(72) 発明者 リウ、ジュンベング

中華人民共和国、250014 シャンドング、ジナン、ハイ - テク・ゾーン、シュンヤ・ロード、No. 1036

(72) 発明者 ガオ、フェイ

中華人民共和国、250014 シャンドング、ジナン、ハイ - テク・ゾーン、シュンヤ・ロード、No. 1036

(72) 発明者 リウ、チェングピング

中華人民共和国、250014 シャンドング、ジナン、ハイ - テク・ゾーン、シュンヤ・ロード、No. 1036

(72) 発明者 ガオ、フェイ

中華人民共和国、250014 シャンドング、ジナン、ハイ - テク・ゾーン、シュンヤ・ロード、No. 1036

(72) 発明者 ズフ、ボウ

中華人民共和国、250014 シャンドング、ジナン、ハイ - テク・ゾーン、シュンヤ・ロード、No. 1036

Fターム(参考) 5K030 GA03 HB08 LB05 LE03