

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5557840号  
(P5557840)

(45) 発行日 平成26年7月23日 (2014. 7. 23)

(24) 登録日 平成26年6月13日 (2014. 6. 13)

(51) Int. Cl.

F I

G 0 6 F 12/00 (2006. 01)

G 0 6 F 13/00 (2006. 01)

G 0 6 F 12/00 5 3 1 M

G 0 6 F 12/00 5 4 5 A

G 0 6 F 13/00 5 2 0 D

G 0 6 F 13/00 3 5 1 N

G 0 6 F 13/00 3 5 1 M

請求項の数 18 (全 31 頁)

(21) 出願番号 特願2011-529547 (P2011-529547)  
 (86) (22) 出願日 平成21年9月30日 (2009. 9. 30)  
 (65) 公表番号 特表2012-504807 (P2012-504807A)  
 (43) 公表日 平成24年2月23日 (2012. 2. 23)  
 (86) 国際出願番号 PCT/EP2009/062714  
 (87) 国際公開番号 W02010/037794  
 (87) 国際公開日 平成22年4月8日 (2010. 4. 8)  
 審査請求日 平成24年8月30日 (2012. 8. 30)  
 (31) 優先権主張番号 61/102, 408  
 (32) 優先日 平成20年10月3日 (2008. 10. 3)  
 (33) 優先権主張国 米国 (US)

(73) 特許権者 598036300  
 テレフオンアクチーボラゲット エル エ  
 ム エリクソン (パブル)  
 スウェーデン国 スtockホルム エスー  
 1 6 4 8 3  
 (74) 代理人 100076428  
 弁理士 大塚 康德  
 (74) 代理人 100112508  
 弁理士 高柳 司郎  
 (74) 代理人 100115071  
 弁理士 大塚 康弘  
 (74) 代理人 100116894  
 弁理士 木村 秀二  
 (74) 代理人 100130409  
 弁理士 下山 治

最終頁に続く

(54) 【発明の名称】 分散データベースの監視メカニズム

(57) 【特許請求の範囲】

【請求項 1】

データの少なくとも1つのパーティションのレプリカを格納するように各ノードが構成された複数のノードを有する分散データベースシステムを扱う方法であって、

格納されるデータをp個のパーティションに分割するステップ (S - 0 0 5) と、

各パーティションをr個のレプリカに複製するステップ (S - 0 1 0) と、

各パーティションについて、前記r個のレプリカを前記複数のノードの中から選択された対応するr個のノードに分配するステップ (S - 0 1 5、S - 0 2 0) と、

相互にアドレスを指定するために使用可能である他のノードの識別子のリストを各ノードに設定するステップ (S - 0 1 7、S - 0 2 0) と、

前記複数のノードの中から2つ以上のノードを起動するステップ (S - 0 3 0、S - 0 3 5、S - 0 6 0、S - 0 7 0) と、

各アクティブノードにおいて、各レプリカの最新の更新と、レプリカ状態と、各レプリカを担当するローカルリソースの状態と、各レプリカの接続状態とのうちから選択された少なくとも1つのイベントを監視するステップ (S - 0 5 0、S - 0 5 5、S - 0 9 0、S - 0 9 5) と、

前記複数のノードの中のノードの起動又は停止の際に、前記アクティブノードの中から選択された少なくとも1つのノードにおいて、前記少なくとも1つのイベントに関する情報を各アクティブノードから収集し、各レプリカについての前記収集されたイベントに依存して前記アクティブノード内の各レプリカの優先順位を決めるためのルールを適用し、

10

20

各パーティションについて最高のレプリカ優先度を有する特定のノード内のレプリカを選択するステップであって、当該選択されたレプリカが現在のマスタレプリカであり、当該特定のノードが当該パーティションについての現在のマスタノードである、ステップと、

前記分散データベースシステム内のデータの読み出し／書き込みを行うことの、ノードにおいて受信された（S - 150；S - 160）任意のリクエストについて、当該データが属するパーティション（11、12）と当該パーティションについての現在のマスタレプリカを担当する現在のマスタノード（2105）とを決定し、当該リクエストを当該現在のマスタノードヘルレーティングする（S - 151、S - 152、S - 153；S - 161）ステップと

を有することを特徴とする方法。

10

【請求項2】

各パーティションについて、前記r個のレプリカを対応するr個のノードに分配する前記ステップは、他の基準が同一のレプリカ優先度を生み出す場合に適用されるデフォルトレプリカ優先度を各レプリカに設定するステップを含むことを特徴とする請求項1に記載の方法。

【請求項3】

前記アクティブノードの中から選択された少なくとも1つのノードによって、パーティションについての所与のデフォルトレプリカ優先度が設定されていることに関する情報を少なくとも1つのアクティブノードから収集するステップと、

前記アクティブノードの中から選択された少なくとも1つのノードによって、各パーティションについて最高のデフォルトレプリカ優先度を有する特定のノード内のレプリカを選択するステップであって、当該選択されたレプリカが現在のマスタレプリカであり、当該特定のノードが当該パーティションについての現在のマスタノードである、ステップとを含むことを特徴とする請求項2に記載の方法。

20

【請求項4】

前記複数のノードの中のノードの起動又は停止の際に、前記アクティブノードが起動された順序を決定するステップをさらに有することを特徴とする請求項1に記載の方法。

【請求項5】

最初に起動されたアクティブノードは、

前記少なくとも1つのイベントに関する前記情報を各アクティブノードから収集するステップと、

30

各レプリカについての前記収集されたイベントに依存して、前記アクティブノード内の各レプリカの優先順位を決めるための前記ルールを適用するステップと、

各パーティションについて前記最高のレプリカ優先度を有する前記特定のノード内の前記レプリカを選択するステップであって、当該選択されたレプリカが現在のマスタレプリカであり、当該特定のノードが当該パーティションについての現在のマスタノードである、ステップと、

各パーティションについて選択された現在のマスタレプリカと当該マスタレプリカを保持する現在のマスタノードとに関して他のアクティブノードへ通知するステップと  
の実行を担当するシステムマスタモニタであるとみなされることを特徴とする請求項4に記載の方法。

40

【請求項6】

最初に起動されたアクティブノードは、

パーティションについての所与のデフォルトレプリカ優先度が設定されていることに関する情報を少なくとも1つのアクティブノードから収集するステップと、

各パーティションについて最高のデフォルトレプリカ優先度を有する特定のノード内のレプリカを選択するステップであって、当該選択されたレプリカが現在のマスタレプリカであり、当該特定のノードが当該パーティションについての現在のマスタノードである、ステップと、

各パーティションについて選択された現在のマスタレプリカと当該現在のマスタレプ

50

リカを保持する現在のマスタノードとに関して他のアクティブノードへ通知するステップと  
の実行を担当するシステムマスタモニタであるとみなされることを特徴とする請求項4に  
記載の方法。

【請求項7】

各パーティションについて各アクティブノードにおいて前記現在のマスタレプリカの内容を、当該パーティションについての前記現在のマスタレプリカ(112、121、132、142)を担当する前記現在のマスタノード(105)からコピーするステップをさらに有することを特徴とする請求項1に記載の方法。

【請求項8】

各アクティブノードにおいてコピーされた各レプリカについて、行われた前記最新の更新と、レプリカ状態と、前記レプリカを担当するローカルリソースの状態と、前記レプリカの接続状態とのうちの少なくとも1つを作成するステップをさらに有することを特徴とする請求項7に記載の方法。

【請求項9】

複数のノード(1、2、3、4)を有する分散データベースシステムであって、  
格納されるデータは複数のパーティションに分割され、  
各パーティションは、前記複数のノードから選択された対応する数のノードに分散され  
る複数のレプリカに複製され、

前記複数のノードのうちの少なくとも1つのノードはアクティブであり、  
前記少なくとも1つのノードは、

少なくとも1つのパーティション(112、123、143)のレプリカ(2101)  
を格納するとともに、相互にアドレスを指定するために使用可能である他のノードの識別子(152)を格納するためのデータ記憶装置(15)と、

前記分散データベースシステムの他のノード(1、3、4)と通信するとともに、前記分散データベースシステムにおける読み出し/書き込み動作を要求するクライアント(5)と通信するための入出力部(30)と、

各レプリカの最新の更新(2103)と、レプリカ状態(212、223、243)と、各レプリカを担当するローカルリソースの状態と、各レプリカ(112、123、143)の接続状態(312、323、343)とのうちから選択された少なくとも1つのイベントを監視するための監視部(60)と、

前記データ記憶装置、前記監視部及び前記入出力部と連携して、前記少なくとも1つのイベントに関する情報を前記分散データベースシステムの各アクティブノードから収集し、各レプリカについての前記収集されたイベントに依存して前記アクティブノード内の各レプリカの優先順位を決めるためのルールを適用し、各パーティションについて最高のレプリカ優先度を有する特定のノード内のレプリカを選択する動作であって、当該選択されたレプリカが現在のマスタレプリカであり、当該特定のノードが当該パーティションについての現在のマスタノードある、動作を実行し、前記分散データベースシステム内のデータの読み出し/書き込みを行うことの、受信された任意のリクエストについて、当該データが属するパーティションと当該パーティションについての現在のマスタレプリカを担当する現在のマスタノードとを決定し、当該リクエストを当該現在のマスタノードへルーティングするための制御部(20)と  
を含むことを特徴とする分散データベースシステム。

【請求項10】

各ノードの前記データ記憶装置(15)は、他の基準が同一のレプリカ優先度を生み出す場合に適用されるデフォルトレプリカ優先度を示すようにレプリカごとに設定されたインジケータ(2102)を格納するように構成されることを特徴とする請求項9に記載の分散データベースシステム。

【請求項11】

各ノードの前記処理部(20)、前記監視部(60)、前記データ記憶装置(15)及

10

20

30

40

50

び前記入力部（３０）は、

パーティションについての所与のデフォルトレプリカ優先度が設定されていることに関する情報を少なくとも１つのアクティブノードから収集し、

各パーティションについて最高のデフォルトレプリカ優先度を有する特定のノード内のレプリカを選択し、当該選択されたレプリカが現在のマスタレプリカであり、当該特定のノードが当該パーティションについての現在のマスタノードである

ようにさらに構成されることを特徴とする請求項１０に記載の分散データベースシステム。

【請求項１２】

各ノードの前記処理部（２０）、前記監視部（６０）、前記データ記憶装置（１５）及び前記入力部（３０）は、前記分散データベースシステムの各アクティブノードから、前記アクティブノードが起動された順序を決定するための情報（１１０４、２１０４、３１０４、４１０４）を収集するように構成されることを特徴とする請求項９に記載の分散データベースシステム。

【請求項１３】

最初に起動されたアクティブノードはシステムマスタモニタであるとみなされ、前記システムマスタモニタの前記処理部（２０）、前記監視部（６０）、前記データ記憶装置（１５）及び前記入力部（３０）は、

前記少なくとも１つのイベントに関する前記情報を各アクティブノードから収集し、各レプリカについての前記収集されたイベントに依存して、前記アクティブノード内の各レプリカの優先順位を決めるための前記ルールを適用し、

各パーティションについて前記最高のレプリカ優先度を有する前記特定のノード内の前記レプリカを選択し、当該選択されたレプリカが現在のマスタレプリカであり、当該特定のノードが当該パーティションについての現在のマスタノード（２１０５、４１０５）であり、

各パーティションについて選択された現在のマスタレプリカと当該マスタレプリカを保持する現在のマスタノードとに関して他のアクティブノードへ通知するように構成されることを特徴とする請求項１２に記載の分散データベースシステム。

【請求項１４】

最初に起動されたアクティブノードはシステムマスタモニタであるとみなされ、前記システムマスタモニタの前記処理部（２０）、前記監視部（６０）、前記データ記憶装置（１５）及び前記入力部（３０）は、

パーティションについての所与のデフォルトレプリカ優先度が設定されていることに関する情報を少なくとも１つのアクティブノードから収集し、

各パーティションについて最高のデフォルトレプリカ優先度を有する特定のノード内のレプリカを選択し、当該選択されたレプリカが現在のマスタレプリカであり、当該特定のノードが当該パーティションについての現在のマスタノード（２１０５、４１０５）であり、

各パーティションについて選択された現在のマスタレプリカと当該マスタレプリカを保持する現在のマスタノードとに関して他のアクティブノードへ通知するようにさらに構成されることを特徴とする請求項１２に記載の分散データベースシステム

【請求項１５】

各ノードの前記処理部（２０）、前記監視部（６０）、前記データ記憶装置（１５）及び前記入力部（３０）は、各パーティションについて各アクティブノードにおいて前記現在のマスタレプリカの内容を、当該パーティションについての前記現在のマスタレプリカ（１１２、１２１、１３２、１４２）を担当する前記現在のマスタノード（２１０５、４１０５）からコピーするようにさらに構成されることを特徴とする請求項９に記載の分散データベースシステム。

【請求項１６】

各ノードの前記処理部(20)、前記監視部(60)、前記データ記憶装置(15)及び前記入力部(30)は、各アクティブノードにおいてコピーされた各レプリカについて、行われた前記最新の更新と、レプリカ状態と、各レプリカを担当するローカルリソースの状態と、前記レプリカの接続状態とを作成するようにさらに構成されることを特徴とする請求項15に記載の分散データベースシステム。

【請求項17】

入出力部と処理部とを有するコンピュータの内部メモリにロード可能なコンピュータプログラムであって、前記コンピュータで動作する場合に請求項1乃至8の何れか1項に記載の方法を実行するように構成された実行可能なコードを含むことを特徴とするコンピュータプログラム。

10

【請求項18】

コンピュータにおいて読み出し可能であり、請求項17に記載のコンピュータプログラムを含む記録媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は分散データベースに関し、特に通信ネットワークのための共通データベースとして利用可能な地理的に分散したデータベースに関する。さらに具体的に、本発明は改良された分散データベースシステムだけでなく、このような分散データベースシステムを扱う方法を含む。

20

【背景技術】

【0002】

本発明は通信ネットワークの分野におけるアプリケーションの多様性について共通集中データベースが有する問題を解決する。通信システムの様々に異なる世代をサポートし、有線システムまたは無線システムである通信ネットワークのほとんどは、従来、加入情報および加入者データだけでなく、通信ネットワーク内または第三者のサービスネットワーク内に存在するが当該通信ネットワークの加入者がアクセス可能である様々なアプリケーションについてのサービスデータを格納するために、1つ以上の共通集中データベースを利用する。

【0003】

30

通信ネットワークが成長するにつれて、通信システムの新たな世代が現われ、既存の共通集中データベースは通信ネットワーク内のすべての通信システムのニーズに常に適合できるとは限らず、またはこのニーズに常に適切であるとは限らない。それにも係わらず、通信ネットワークは、ここで利用される任意の特定のデータベースシステムで満たされるべき極めて類似する要件を共有する。

【0004】

従来の中央管理通信データベースは一般に、少なくとも以下の特徴をサポートする必要がある。回復性および高可用性、一貫性、高性能および低レンテシ、大容量、拡張性、地理的冗長性、柔軟な配置およびデータモデル、単一アクセス(各地理的位置に1つ)、ならびに単一障害点の不存在。

40

【0005】

この点において、従来の中央管理通信データベースについての地理的冗長性は一般に、メインノードに加えて複製されたノードを有するものとして理解されてきた。ここで、メインノードは動作中であり、複製されたノードはメインノードが何らかの理由でダウンした場合に動作を再開するためにスタンバイしている。

【0006】

近年、純粋な集中データベースは、いくつかある欠点の中で特に、信号転送の観点で極めて高コストである。実際に、一部のデータが他のデータよりも頻繁にアクセスされるように様々なアプリケーションデータがするため、そのリソースの使用は十分にバランスが取られていない。集中データベースが位置する地理的位置に関して、経済的な観点だけで

50

なく負荷および障害のリスクの観点の両方において、信号転送のコストについて選択が極めて重要になりうる。明らかに、データベースクライアントとデータベース自体との間の信号経路が長くなると、ノードがその最終宛先へ信号をさらに送出することができなくなるリスクが大きくなる。同様に、この信号パスが長くなると、通信ネットワークがサポートする負荷が大きくなり、実行時間が長くなる。これとは別に、通信ネットワークは、他のネットワーク事業者に属している様々なアクセスネットワークを通じてしばしば通信される様々な遠くの領域を通じて広がる。このシナリオでは、この信号経路が長くなると、より多くの通信事業者が影響を受け、より多くのコストが引き出されうる。

【0007】

他方で、他の技術を調査すると、分散データベースは、場合によっては中央データベース管理システムの制御の下で、物理的または論理的に分散された複数のデータベースとしてみなされてもよく、必ずしもすべての記憶装置が共通の処理部に接続される必要はない。よって、分散データベースは同一の物理的位置に位置する複数のコンピュータで構築されてもよいし、相互接続されたコンピュータのネットワーク上に点在してもよい。

【0008】

一般的に言うと、データベースインスタンスの分散はデータ分散自体の結果である必要はなく、高可用性システムおよび地理的冗長性を得るためのデータレプリケーションのために有益である。特に、従来の集中データベースシステムを、中央管理または分散された共通バックエンドへアクセスする複数の通信データベースフロントエンドで置き換えるいわゆるデータレイヤアーキテクチャ技術は、通信データベースについての上記の要件を満たし、データ分散およびレプリケーションに利用可能な分散データベースの例示の適用性である。

【0009】

分散データベースの様々なデータベースインスタンスにおけるデータのレプリケーションは、様々なデータベースインスタンスにおいて現存するレプリカを最新に維持するために複雑な管理を必要とする。さらに、通信ネットワークの分野の従来の分散データベースのクライアントは、近くのデータベースインスタンスにおける任意のデータベース関連の動作が常に実行可能であるとは限らない。よって、当該データベース関連の動作に必要な信号経路を常に最短化できるとは限らない。さらに、通信ネットワークの分野における従来の分散データベースのクライアントは、任意のデータベースインスタンスが利用不可能であるデータへのアクセスで問題を経験しうる。

【発明の概要】

【発明が解決しようとする課題】

【0010】

本発明は上述の欠点を少なくとも軽減し、データの少なくとも1つのパーティションのレプリカを格納するようにそれぞれが構成された複数のノードを有する改良された分散データベースシステムと、当該分散データベースシステムを扱う方法とを提供する。

【課題を解決するための手段】

【0011】

本発明の第1側面では、データの少なくとも1つのパーティションのレプリカを格納するようにそれぞれが構成された複数のノードを有する分散データベースシステムを扱う新たな方法が提供される。

【0012】

この方法は、格納されるデータを  $p$  個のパーティションに分割するステップと、各パーティションを  $r$  個のレプリカに複製するステップと、各パーティションについて、前記  $r$  個のレプリカを前記複数のノードの中から選択された対応する  $r$  個のノードに分配するステップと、相互にアドレスを指定するために使用可能である他のノードの識別子のリストを各ノードに設定するステップと、前記複数のノードの中から2つ以上のノードを起動するステップと、各アクティブノードにおいて、各レプリカの最新の更新と、レプリカ状態と、各レプリカを担当するローカルリソースの状態と、各レプリカの接続状態とのうちか

ら選択された少なくとも1つのイベントを監視するステップと、前記複数のノードの中のノードの起動又は停止の際に、前記アクティブノードの中のどのノードが各パーティションについての現在のマスタノードであり、当該パーティションについての現在のマスタレプリカを担当するとみなされるかを決定するステップと、前記分散データベースシステム内のデータの読み出し/書き込みを行うことの、ノードにおいて受信された任意のリクエストについて、当該データが属するパーティションと当該パーティションについての現在のマスタレプリカを担当する現在のマスタノードとを決定し、当該リクエストを当該現在のマスタノードヘルディングするステップとを有する。

【0013】

前記アクティブノードの中のどのノードが各パーティションについての前記マスタノードであり、それ故当該パーティションについての現在のマスタレプリカを担当しているとみなされるかを決定する際に、前記方法は、前記アクティブノードの中の少なくとも1つのノードにおいて、各レプリカの最新の更新と、レプリカ状態と、各レプリカを担当するローカルリソースの状態と、各レプリカの接続状態とのうちから選択された前記少なくとも1つのイベントに関する情報を各アクティブノードから収集するステップと、各レプリカについての前記収集されたイベントに依存して、前記アクティブノード内の各レプリカの優先順位を決めるための事前に設定されたルールを適用するステップと、各パーティションについて最高のレプリカ優先度を有する特定のノード内のレプリカを選択するステップであって、このレプリカが前記マスタレプリカであるとみなされ、この特定のノードが前記パーティションについての前記マスタノードであるとみなされるステップとを含んでもよい。

【0014】

本発明は、前記アクティブノードの中のどのノードが各パーティションについてのマスタノードであり前記パーティションについての現在のマスタレプリカを担当するとみなされるかを決定する際に、2つの主な実施形態、すなわち動作モードを提供する。

【0015】

第1動作モードでは、処理の観点からすべてのノードは似たもの同士のノードであるため、すべてのノードが同じ情報を処理でき、各パーティションについて、現在のマスタレプリカを担当する同じマスタノードを決定することに到りうる。この実施形態のもとでは、前記少なくとも1つのアクティブノードではなく、分散データベースシステム内の各ノードが以下のステップを実行するように構成されてもよい。すなわち、前記アクティブノードの中の少なくとも1つのノードにおいて、各レプリカの最新の更新と、レプリカ状態と、各レプリカを担当するローカルリソースの状態と、各レプリカの接続状態とのうちから選択された前記少なくとも1つのイベントに関する情報を各アクティブノードから収集するステップと、各レプリカについての前記収集されたイベントに依存して、前記アクティブノード内の各レプリカの優先順位を決めるための事前に設定されたルールを適用するステップと、各パーティションについて最高のレプリカ優先度を有する特定のノード内のレプリカを選択するステップであって、このレプリカが前記マスタレプリカであるとみなされ、この特定のノードが前記パーティションについての前記マスタノードであるとみなされるステップとである。

【0016】

第2動作モードでは、前記複数のノードの中の任意のノードの起動または停止の際に、前記方法は前記アクティブノードが起動された順序を決定するステップをさらに有してもよい。この場合に、前記アクティブノードが起動された順序を決定するようにすべてのノードが構成され、その結果、最初に起動されたアクティブノードは以下のステップの実行を担当するいわゆるシステムマスタモニタノードであるとみなされる。すなわち、各レプリカの最新の更新と、レプリカ状態と、各レプリカを担当するローカルリソースの状態と、各レプリカの接続状態とのうちから選択された前記少なくとも1つのイベントに関する情報を各アクティブノードから収集するステップと、各レプリカについての前記収集されたイベントに依存して、前記アクティブノード内の各レプリカの優先順位を決めるた

めの事前に設定されたルールを適用するステップと、各パーティションについて最高のレプリカ優先度を有する特定のノード内のレプリカを選択するステップであって、このレプリカが前記マスタレプリカであるとみなされ、この特定のノードが前記パーティションについての前記マスタノードであるとみなされるステップと、各パーティションについて選択された前記マスタレプリカと当該マスタレプリカを保持する前記マスタノードとに関して他のアクティブノードへ通知するステップとである。

【0017】

前記収集されたイベントに依存して、事前に設定されたルールの適用が2つ以上のレプリカについて同一の優先度を生み出す場合に特に有用であるが、各パーティションについて、前記 $r$ 個のレプリカに対応する $r$ 個のノードに分配する方法の前記ステップは、他の基準が同一のレプリカ優先度を生み出す場合に適用されるデフォルトレプリカ優先度を各レプリカに設定するステップを含んでもよい。この場合に、前記アクティブノードの中のどのノードが各パーティションについてのマスタノードであり、当該パーティションについての現在のマスタレプリカを担当するとみなされるかを決定する前記ステップは、前記少なくとも1つのアクティブノードにおいて、各パーティションについての所与のデフォルトレプリカ優先度が設定されていることに関する情報を少なくとも1つのアクティブノードから収集するステップと、各パーティションについて最高のデフォルトレプリカ優先度を有する特定のノード内のレプリカを選択するステップであって、このレプリカが前記マスタレプリカであるとみなされ、この特定のノードが前記パーティションについての前記マスタノードであるとみなされるステップとを含んでもよい。特に、上記の第1動作モードに従って前記方法が動作する場合に、前記少なくとも1つのノードではなく、前記分散データベースシステム内の各ノードがこれらのステップを実行するように構成されてもよい。

【0018】

しかしながら、上記第2動作モードに従って前記方法が動作する場合に、最初に起動されたアクティブノードはシステムマスタモニタであるとみなされ、前記アクティブノードの中のどのノードが各パーティションについての現在のマスタノードであり、当該パーティションについての現在のマスタレプリカを担当するとみなされるかを決定するステップは、当該システムマスタモニタノードにおいて、各パーティションについての所与のデフォルトレプリカ優先度が設定されていることに関する情報を少なくとも1つのアクティブノードから収集するステップと、各パーティションについて最高のデフォルトレプリカ優先度を有する特定のノード内のレプリカを選択するステップであって、このレプリカが前記マスタレプリカであるとみなされ、この特定のノードが前記パーティションについての前記マスタノードであるとみなされるステップと、各パーティションについて選択された前記マスタレプリカと当該マスタレプリカを保持する前記マスタノードとに関して他のアクティブノードへ通知するステップとを含んでもよい。

【0019】

一般的に言うと、すべての読み出し/書き込み動作は関連するパーティションについてのマスタレプリカ上で実行される。従って、マスタレプリカの内容と、同一のパーティションについての他のレプリカの内容とは所定の時点において異なりうる。各パーティションについての様々なレプリカの間での一貫性を維持するために、前記方法は、各パーティションについて各アクティブノードにおいて前記現在のマスタレプリカの内容を、当該パーティションについての前記現在のマスタレプリカを担当する前記現在のマスタノードからコピーするステップをさらに有してもよい。コピーするステップが行われる場合に、前記方法は、各アクティブノードにおいてコピーされた各レプリカについて、行われた前記最新の更新と、レプリカ状態と、前記レプリカを担当するローカルリソースの状態と、前記レプリカの接続状態とのうちの少なくとも1つを作成するステップをさらに有してもよい。これは特に、前記マスタレプリカを担当する現在のマスタノードに障害が生じ、ダウンまたは利用不可能になり、もしくは非アクティブノードになる場合に将来の別のマスタレプリカを選択するために有利である。



## 【 0 0 2 0 】

他方で、本発明の第2側面に従うと、複数のノードを有する分散データベースシステムであって、各ノードはデータの少なくとも1つのパーティションのレプリカを格納するように構成される改良された分散データベースシステムが提供される。この分散データベースシステムにおいて、各ノードは、格納されるデータの少なくとも1つのデータパーティションのレプリカを格納するとともに、相互にアドレスを指定するために使用可能である他のノードの識別子を格納するためのデータ記憶装置と、前記分散データベースシステムの他のノードと通信するとともに、前記分散データベースシステムにおいて読み出し/書き込み動作を要求するクライアントと通信するための入出力部と、各レプリカの最新の更新と、レプリカ状態と、各レプリカを担当するローカルリソースの状態と、各レプリカの接続状態とのうちから選択された少なくとも1つのイベントを監視するための監視部と、前記データ記憶装置、前記監視部及び前記入出力部と連携して、前記分散データベースシステム内のアクティブノードの中のどのノードが各パーティションについての現在のマスタノードであり、当該パーティションについての現在のマスタレプリカを担当するとみなされるかを決定し、前記分散データベースシステム内のデータの読み出し/書き込みを行うことの、受信された任意のリクエストについて、当該データが属するパーティションと当該パーティションについての現在のマスタレプリカを担当する現在のマスタノードとを決定し、当該リクエストを当該現在のマスタノードヘルディングするための制御部とを含む。

10

## 【 0 0 2 1 】

20

上記方法と整合して、前記分散データベースシステムが第1動作モードに従って動作するか第2動作モードに従って動作するかにかかわらず、各ノードの前記処理部、前記監視部、前記データ記憶装置及び前記入出力部は、各レプリカの最新の更新と、レプリカ状態と、各レプリカを担当するローカルリソースの状態と、各レプリカの接続状態とのうちから選択された少なくとも1つのイベントに関する情報を前記分散データベースシステムの各アクティブノードから収集し、各レプリカについての前記収集されたイベントに依存して、前記アクティブノード内の各レプリカの優先順位を決めるための事前に設定されたルールを適用し、各パーティションについて最高のレプリカ優先度を有する特定のノード内のレプリカを選択し、このレプリカが前記マスタレプリカであるとみなされ、この特定のノードが前記パーティションについての前記マスタノードであるとみなされるように構成されてもよい。

30

## 【 0 0 2 2 】

特に、前記収集されたイベントに依存して、事前に設定されたルールの適用が2つ以上のレプリカについて同一の優先度を生み出す場合の上記方法に有利に整合して、各ノードの前記データ記憶装置は、他の基準が同一のレプリカ優先度を生み出す場合に適用されるデフォルトレプリカ優先度を示すようにレプリカごとに設定されたインジケータを格納するように構成されてもよい。この場合に、各ノードの前記処理部、前記監視部、前記データ記憶装置及び前記入出力部は、パーティションについての所与のデフォルトレプリカ優先度が設定されていることに関する情報を少なくとも1つのアクティブノードから収集し、各パーティションについて最高のデフォルトレプリカ優先度を有する特定のノード内のレプリカを選択し、このレプリカが前記マスタレプリカであるとみなされ、この特定のノードが前記パーティションについての前記マスタノードであるとみなされるようにさらに構成されてもよい。

40

## 【 0 0 2 3 】

また、上記方法に整合して、分散データベースシステムが第2動作モードに従って動作する場合に特に適用可能であるが、各ノードの前記処理部、前記監視部、前記データ記憶装置及び前記入出力部は、前記分散データベースシステムの各アクティブノードから、前記アクティブノードが起動された順序を決定するための情報を収集するように構成されてもよい。

## 【 0 0 2 4 】

50

この場合に、最初に起動されたアクティブノードはシステムマスタモニタであるとみなされてもよく、前記システムマスタモニタの前記処理部、前記監視部、前記データ記憶装置及び前記入力部は、各レプリカの最新の更新と、レプリカ状態と、各レプリカを担当するローカルリソースの状態と、各レプリカの接続状態とのうちから選択された少なくとも1つのイベントに関する情報を各アクティブノードから収集し、各レプリカについての前記収集されたイベントに依存して、前記アクティブノード内の各レプリカの優先順位を決めるための事前に設定されたルールを適用し、各パーティションについて最高のレプリカ優先度を有する特定のノード内のレプリカを選択し、このレプリカが前記マスタレプリカであるとみなされ、この特定のノードが前記パーティションについての前記マスタノードであるとみなされ、各パーティションについて選択されたマスタレプリカと当該マスタレプリカを保持するマスタノードとに関して他のアクティブノードへ通知するように構成されてもよい。

10

#### 【0025】

さらに、上述のデフォルト優先度が第2動作モードの下でのパーティションについてのレプリカの優先順位を決めるために関連する場合に、前記処理部、前記監視部、前記データ記憶装置及び前記入力部は、パーティションについての所与のデフォルトレプリカ優先度が設定されていることに関する情報を少なくとも1つのアクティブノードから収集し、各パーティションについて最高のデフォルトレプリカ優先度を有する特定のノード内のレプリカを選択し、このレプリカが前記マスタレプリカであるとみなされ、この特定のノードが前記パーティションについての前記マスタノードであるとみなされ、各パーティションについて選択された前記マスタレプリカと当該マスタレプリカを保持する前記マスタノードとに関して他のアクティブノードへ通知するようにさらに構成されてもよい。

20

#### 【0026】

各パーティションについての様々なレプリカの間での一貫性を維持し、それゆえ上記の方法と整合するために、各ノードの前記処理部、前記監視部、前記データ記憶装置及び前記入力部は、各パーティションについて各アクティブノードにおいて前記現在のマスタレプリカの内容を、当該パーティションについての前記現在のマスタレプリカを担当する前記現在のマスタノードからコピーするようにさらに構成される。さらに、現在のマスタレプリカが何らかの理由で非アクティブになる場合にマスタレプリカを担当する別のマスタノードをさらに選択するために、各ノードの前記処理部、前記監視部、前記データ記憶装置及び前記入力部は、各アクティブノードにおいてコピーされた各レプリカについて、行われた前記最新の更新と、レプリカ状態と、各レプリカを担当するローカルリソースの状態と、前記レプリカの接続状態とのうちの少なくとも1つを作成するようにさらに構成されてもよい。

30

#### 【0027】

他方で、本発明の第3側面に従って本発明はコンピュータプログラムで実施されてもよい。コンピュータプログラムは入出力部と処理部とを有するコンピュータの内部メモリにロード可能なコンピュータプログラムであって、上記の方法を実行するように構成された実行可能なコードを含む。特に、この実行可能なコードはコンピュータにおいて読み出し可能な記録媒体に記録されてもよい。

40

#### 【0028】

本発明の機能、目的および利点は添付の図面と併せて本明細書を読むことによって明らかになるだろう。

#### 【図面の簡単な説明】

#### 【0029】

【図1A】格納されるデータを $p$ 個のパーティションに分割し、各パーティションを $r$ 個のレプリカに複製するために実行される動作シーケンスの簡略図。

【図1B】各パーティションについて、 $r$ 個のレプリカを複数のノードの中から選択された対応する $r$ 個のノードに分配し、相互にアドレスを指定するために使用可能である他のノードの識別子のリストを各ノードに設定する動作シーケンスの簡略図。

50

【図 1 C】データの少なくとも 1 つのパーティションのレプリカを格納するように各ノードが構成された複数のノードを有する分散データベースシステムを扱う方法を実行するために、図 1 A および図 1 B に説明された動作と連動して実行される動作シーケンスの簡略図。

【図 2】本発明の一部の実施形態を説明するために有用なデータを有する分散データベース内の複数のノードの例示の構成の簡略図。

【図 3】分散データベースに含まれる複数のノードの中のノードの例示の実装を説明する図。

【図 4】本発明の一部の実施形態を説明するためにクラスタとして構成された有用なデータを有する分散データベース内の複数のノードの中のノードの例示の構成の簡略図。

10

【図 5】分散データベース内のデータの読み出し/書き込みを行うことの、ノードにおいて受信された任意のリクエストを、このようなデータが属するパーティションについての現在のマスタレプリカを担当する現在のマスタノードヘルレーティングするために実行される動作シーケンスの簡略図。

【図 6】アクティブノードの中のどのノードが、他のノードの調整と各レプリカについてのマスタノードがどれであるかの決定とを担当するコントローラシステムモニタであるとみなされるかを決定するために本発明の実施形態に従って提供される例示の状態マシンを説明する図。

【図 7】分散データベースシステム内の複数のノードの起動の際にアクティブノードの中のどのノードがコントローラシステムモニタであるとみなされるかを決定するために、図 6 に示される状態マシンのサポートのもので実行される例示の動作シーケンスを説明する図。

20

【図 8】以前にみなされていたコントローラシステムモニタの非アクティブ化の際にアクティブノードの中のどのノードがコントローラシステムモニタであるとみなされるかを決定するために、図 6 に示される状態マシンのサポートのもので実行される別の例示の動作シーケンスを説明する図。

【発明を実施するための形態】

【0030】

データの少なくとも 1 つのパーティションのレプリカをそれぞれが格納する複数のノードを有する改良された分散データベースシステムと当該分散データベースシステムを扱う方法との現時点で好適な実施形態を以下に説明する。

30

【0031】

通信データベースシステムは複数の地理的に分散されたノードを含んでもよく、各ノードは複数のデータ記憶部を含んでもよく、各ノード内の各データ記憶部はデータの部分集合の特定のレプリカ、すなわちパーティションのレプリカを割り振ってもよい。図 1 A および図 1 B に説明されるように、ノード 1 ~ 4 のデータ記憶部の間でのデータ集合 10 の例示の分配は本発明に従って提供される複数のステップに従って実行されうる。

【0032】

図 1 A に示されるように、データ集合 10 はステップ S - 005 の間に複数のパーティション 11 ~ 14 に分割され、各パーティションはデータ集合 10 の特定の部分集合を有する。次いで、各パーティションについて、ステップ S - 010 の間に複数のレプリカが生成される。各パーティションについてのレプリカの個数はすべてのパーティションについて同じである必要はない。よって、図 1 A に示される例のように、4 つのレプリカ 111 ~ 114 がパーティション 11 について生成され、3 つのレプリカ 121 ~ 123 とレプリカ 141 ~ 143 とがパーティション 12 と 14 とについてそれぞれ生成され、ただ 2 つのレプリカ 131、132 がパーティション 13 について生成される。

40

【0033】

図 1 B に示されるように、これらのレプリカは、必要となる地理的な分散を決定する予備的なステップ S - 015 の間にパーティションごとにグループ分けされてもよい。これとは別に、データベースシステムを構成するノードの地理的な分配を決定する際に、各ノ

50

ードはステップS - 017の間にアドレス目的に利用可能な識別子が割り当てられる。図1Bの例示の説明は別個の識別子N - 1ID、N - 2ID、N - 3IDおよびN - 4IDを有する4つのノード1~4で構成される。

#### 【0034】

次いで、ステップS - 020の間に、各パーティションについて生成されたレプリカは、データベースシステムを構成するノードに分配されてもよい。図1Bに説明される例では、第1パーティションの第1レプリカ111はノード1に格納され、第1パーティションの第2レプリカ112はノード2に格納され、第1パーティションの第3レプリカ113はノード3に格納され、第1パーティションの第4レプリカ114はノード4に格納され、第2パーティションの第1レプリカ121はノード3に格納され、第2パーティションの第2レプリカ122はノード1に格納され、第2パーティションの第3レプリカ123はノード2に格納され、第3パーティションの第1レプリカ131はノード4に格納され、第3パーティションの第2レプリカ132はノード1に格納され、第4パーティションの第1レプリカ141はノード3に格納され、第4パーティションの第2レプリカ142はノード4に格納され、第4パーティションの第3レプリカ143はノード2に格納される。よって、各パーティションについてのレプリカを格納するためにすべてのノードが必要であるとは限らず、本発明の側面に従って提供される分散データベースシステムの各ノードにすべてのパーティションがレプリカを有さなければならないというわけではない。

#### 【0035】

これとは別に、各ノードはまたこのステップの間に、他のノードの識別子を設定されてもよい。よって、ノード1はノード2、3、4を識別する識別子151を格納する。ノード2はノード1、3、4を識別する識別子152を格納する。ノード3はノード1、2、4を識別する識別子153を格納する。ノード4はノード1、2、3を識別する識別子154を格納する。

#### 【0036】

動作中に、以下に説明される所定のイベントに依存して、所定のノード内の1つの特定のレプリカが最高の優先度を取得してもよく、それ故パーティションについてのマスタレプリカであると決定される一方で、所定のノードは当該パーティションについてのマスタノードとみなされる。しかしながら、イベントが相異なるノード内のパーティションの相異なるレプリカについて同一の優先度を生み出し、その結果、マスタレプリカが決定されない状況が存在しうる。2つ以上のレプリカに同一の優先度が与えられる曖昧さを省くために、本発明は他の基準、すなわち上記のイベントの処理結果は同一のレプリカ優先度を生み出す場合に、各レプリカを適用されるデフォルトのレプリカ優先度を設定することを提供する。図2に示されるように、各ノード1~4はそれぞれ、複数のパーティションについてレプリカ1101、2101、3101、4101と、レプリカごとのデフォルト優先度1102、2102、3102、4102を含む。

#### 【0037】

分散データベースシステムの各ノードが上述のように構成されると、分散データベースシステムはオペレータが望むようにノードごとにまたは一斉に動作に入ることの準備が整う。

#### 【0038】

図1Cは本発明の別の側面に従って上記の分散データベースシステムを扱う方法を実行するための後続の動作シーケンスを説明する。すべてのノードが同様に振舞うが、本発明の実施形態に従って以下に説明されるように、各ノードにおける完全な動作シーケンスは相異なるノードが起動される順序に依存してもよい。よって、図1Cの例はノード2がステップS - 030の間に最初に起動されるものであり、その後ステップS - 035の間にノード3が起動され、次いでステップS - 060の間にノード4であり、最後にステップS - 070の間にノード1であるシナリオが説明される。

#### 【0039】

各ノード2、3、4、1の起動の後に、各アクティブノードにおいて当該ノードが起動された開始時刻を決定するそれぞれのステップS - 040、S - 045、S - 080、S - 075が続いてもよい。このオプションのステップは、以下に説明される他のノードとの連携と各パーティションについてのマスタレプリカがどれであるかの決定とを担当するシステムマスタモニタとして動作するノードを有するデータベースシステムが動作している場合に、アクティブノードが起動された順序をさらに決定するために有用である。この点において、図2に説明されるように、各ノード1～4は、開始時刻からノードが稼働している動作時間を示すそれぞれの表示1104、2104、3104、4104を含む。

【0040】

各ノードについて開始時刻を決定するステップが実行されるかどうかに関わらず、総括的に、各ノード2、3、4、1の起動の後に、各レプリカの最新の更新と、レプリカ状態と、各レプリカを担当するローカルリソースの状態と、各レプリカの接続状態とから選択される少なくとも1つのイベントを各アクティブノード2、3、4、1において監視するステップS - 050、S - 055、S - 095、S - 090のそれぞれが続く。

【0041】

この目的のために、図2に説明される例のように、各ノード1～4はそれぞれ、複数のパーティションについて、レプリカ1101、2101、3101、4101とともに、レプリカごとに最新の更新のインジケータ1103、2103、3103、4103を含み、且つ上述のようにレプリカごとのデフォルトの優先度1102、2102、3102、4102を含む。図2に示されるデータの例とは別に、図4に示されるように各ノード1～4においてレプリカごとに他のデータが格納されてもよい。よって、図4に説明される例のように、ノード2はそれぞれ複数のパーティションについてレプリカ112、123、143とともに、レプリカごとの接続状態のインジケータ312、323、343を含み、且つレプリカ212、223、243ごとに、レプリカがマスタレプリカであるか、起動中であるがマスタではないアクティブレプリカであるか、起動していないとみなされうるように設定が行われているスタンバイレプリカであるかを示すレプリカ状態を含む。特に、各ノードはさらに、図示されてはいないが、各レプリカについてのローカルリソースを監督するプロセスがこのようなイベントを監視するステップから分離されうるように、レプリカごとにローカルリソースの状態のインジケータを含めるための、レプリカごとの記憶フィールドをさらに有してもよい。これに替えて、各アクティブノードにおいて選択の少なくとも1つのイベントを監視する上述のステップは、各レプリカについてローカルリソースを監督するステップと各ローカルリソースの状態を決定するステップとを含んでもよい。

【0042】

なおも図4を参照して、この分散データベースシステムの各ノード1～4は複数のクラスタで構成されてもよく、各クラスタは、パーティションのレプリカと、レプリカについての接続状態のインジケータと、レプリカについてのレプリカ状態と、レプリカごとの最新の更新のインジケータと、レプリカごとのデフォルト優先度とを含む。よって、図4に説明された例として、ノード2はクラスタ160、170、180を備え、各クラスタはそれぞれレプリカについてのデータおよびインジケータとともに、パーティションのレプリカ112、123、143を含む。

【0043】

図1Cに説明される動作シーケンスに戻り、イベントを監視するステップが各ノードにおいて実行されると、各ノードはステップS - 065、S - 085、S - 100、S - 110のそれぞれの間に、各ノードに知られているすべての他のノードへのいわゆる「アラライブ(Alive)」メッセージの送信を開始する。

【0044】

当業者は理解するかもしれないが、ノード4がノード1よりも前に開始されたにもかかわらず、ノード1における一部の動作はノード4における対応する動作の前に生じる。これは一般的に起こりうる。なぜなら、1つのノードにおけるプロセッサ負荷は別のノード

10

20

30

40

50

における負荷よりも高いかもしれず、結果として前者における性能が低下するからであり、またこれは様々なネットワーク経路を通じる様々な信号遅延に起因しうる。

【0045】

図1Cに特に例示として説明されるように、ノード2は、ノード1が起動される前に、ステップS-065でいわゆる「アライブ」メッセージを送信した。ノード2はステップS-100の間にノード1からの「アライブ」メッセージを受信するので、ノード2はノード1が当初の「アライブ」メッセージを受信しておらず、ノード1がすべてのノードからの完全な情報を有していないことを認識してもよい。これらの状況では、現在の例示の状況におけるノード2のようなノードはイベントを再び監視するステップS-105と、ノード2に知られているすべてのノードへ「アライブ」メッセージを再び送出するステップS-115とを実行する。さらに、これも図1Cに示されるように、ノード2においてイベントを監視するステップS-105と「アライブ」メッセージを送出するステップS-115との中間のステップS-110の間にノード2においてノード3からの「アライブ」メッセージが受信され、その結果としてノード2から送出された最新の「アライブ」メッセージは、各ノードがすべての他のアクティブノードを認識している最新であるとみなされうる。

10

【0046】

これらの「アライブ」メッセージは分散データベースシステムのノード間で交換され、分散データベースシステムのノードの起動または停止の際に、以下に説明されるように、アクティブノードのうちのどのノードが各パーティションについての現在のマスタノードであり、当該パーティションについての現在のマスタレプリカを担当するとみなされるかを決定するために有用である。

20

【0047】

上述のように、動作中に決定された所定のイベントが、パーティションについてのマスタレプリカと当該マスタレプリカを担当するマスタノードとを決定する際に考慮に入れられる。各パーティションについてのマスタレプリカがどれであるかを決定するために、以下の情報が考慮に入れられてもよい。すなわち、各パーティションについてどのレプリカが完全な情報、つまり更新レベルを有する直近に更新された内容を有するかと、明らかに稼働中のレプリカだけが適格であるので各パーティションについての各レプリカのレプリカ状態と、各パーティションについての各レプリカの接続状態と、パーティションの各レプリカについて設定されたデフォルト優先度である。さらに、デフォルト優先度は以前の基準の結果を上書きするように構成されてもよく、以前の基準がパーティションのレプリカの2つ以上について同じ結果を生み出す場合にのみ適用可能となるように構成されてもよい。

30

【0048】

さらに、動作中に決定されたこれらのイベントに依存して、アクティブノード内の各レプリカの優先順位を決めるために事前に設定されたルールが適用されてもよい。例えば、事前に設定されたルールは、接続状態が更新レベルに優先度を引き継ぐようなものであってもよいし、レプリカ状態のすぐ後でデフォルト優先度が考慮に入れられるものであってもよいし、レプリカの優先順位を決めるためのイベントに関する他の如何なる基準であつてもよい。

40

【0049】

パーティションについての各レプリカの内容だけでなく、このような内容の更新に関して、所定の時刻に各レプリカがマスタレプリカから更新するような従来のルーチンが提供されてもよい。よって、パーティションについてのすべてではないレプリカが同時に内容を更新し、且つ当該パーティションについてのすべてではないレプリカが同じ速度で更新を進行する。更新中に、分散データベースシステム内の各ノードは、レプリカの内容が完全な情報であるとみなされうるかどうかと、交換が実行された時点とを、他のノードと交換される情報がこの点において考慮するように、どれくらい更新が進行しているかを監視する必要がある。

50

## 【 0 0 5 0 】

パーティションについてのマスタレプリカと当該マスタレプリカを担当するマスタノードとの決定に関して上述のイベントを考慮に入れるために、各ノードにおいて監視されたイベントは分散データベースシステム内の他のノードへ通信される。この目的のために、各ノード 1 ~ 4 は他のノードの仮想 IP アドレスを設定されてもよく、または互いに識別してアドレスを指定するために別個のノード識別子 1 5 1 ~ 1 5 4 を利用してもよい。各ノード 1 ~ 4 は定期的に、例えば連続した遅延時間が満了した後に、他のノードへ「アライブ」メッセージを送信してもよい。有利には本発明の実施形態において、「アライブ」メッセージは、TCP で簡単に検出される一方向リンク障害の可能性を回避するために、ある種のハートビートについてUDPを利用する代わりに既知のTCPプロトコルで送信されてもよい。

10

## 【 0 0 5 1 】

さらに正確に、各「アライブ」メッセージは受信側ノードに対して送信側ノードを識別するノード識別子 N - 1 ID、N - 2 ID、N - 3 ID、N - 4 ID を含み、パーティションの各レプリカについて、レプリカが属するパーティションの識別子、レプリカ状態、更新レベル、更新時刻、接続状態およびデフォルト優先度のうちの少なくとも1つを含んでもよい。

## 【 0 0 5 2 】

これに加えて、アクティブノードが起動された順序をさらに決定するために特に有用であるように、任意のノード 1 ~ 4 から分散データベースシステム内の他のノードへ送信された各「アライブ」メッセージはそれぞれ、送信側ノードが自身の開始時刻からアクティブである動作時間の表示 1 1 0 4、2 1 0 4、3 1 0 4、4 1 0 1 を含んでもよい。

20

## 【 0 0 5 3 】

本発明は2つの主な実施形態、すなわちアクティブノードのうちのどのノードが各パーティションについてのマスタノードであり、当該パーティションについての現在のマスタレプリカを担当するとみなされるかの決定における動作モードを提供する。

## 【 0 0 5 4 】

第1動作モードでは、処理の観点からすべてのノードは似たもの同士のノードであるため、各ノードはすべての他のノードから、そこで決定された監視情報を有する「アライブ」メッセージを受信し、その結果、すべてのノードが同じ情報を処理でき、各パーティションについて、現在のマスタレプリカを担当する同じマスタノードを決定することに到りうる。この実施形態のもとでは、分散データベースシステム内の各ノードは以下のステップを実行するように構成されてもよい。すなわち、各レプリカの最新の更新、レプリカ状態、各レプリカを担当するローカルリソースの状態、および各レプリカの接続状態から選択された少なくとも1つのイベントに関する情報を各アクティブノードから収集するステップと、各レプリカについて収集されたイベントに依存して、アクティブノード内の各レプリカの優先順位を決めるための事前に設定されたルールを適用するステップと、各パーティションについて最高のレプリカ優先度を有する特定のノード内のレプリカを選択するステップであって、このレプリカがマスタレプリカであるとみなされ、この特定のノードが当該パーティションについてのマスタノードであるとみなされるステップとである。

30

40

## 【 0 0 5 5 】

特に有利にはこの第1動作モードの下で如何なる特定のレプリカについてこれらのイベントが最高のレプリカ優先度を結果として生じない場合に、分散データベースシステム内の各ノードはさらに以下のステップを実行するように構成されてもよい。すなわち、パーティションについての所与のデフォルトレプリカ優先度を設定されている情報を少なくとも1つのアクティブノードから収集するステップと、各パーティションについて最高のデフォルトレプリカ優先度を有する特定のノード内のレプリカを選択するステップであって、このレプリカがマスタレプリカであるとみなされ、この特定のノードが当該パーティションについてのマスタノードであるとみなされるステップとである。

## 【 0 0 5 6 】

50

第2動作モードでは、アクティブノードが起動された順番を決定するためにすべてのノードが「アライブ」メッセージを処理し、その結果、最初に起動されたアクティブノードは以下のステップの実行を担当するいわゆるシステムマスタモニタノードであるとみなされる。すなわち、各レプリカの最新の更新、レプリカ状態、各レプリカを担当するローカルリソースの状態、および各レプリカの接続状態から選択された少なくとも1つのイベントに関する情報を各アクティブノードから収集するステップと、各レプリカについて収集されたイベントに依存して、アクティブノード内の各レプリカの優先順位を決めるための事前に設定されたルールを適用するステップと、各パーティションについて最高のレプリカ優先度を有する特定のノード内のレプリカを選択するステップであって、このレプリカがマスタレプリカであるとみなされ、この特定のノードが当該パーティションについてのマスタノードであるとみなされるステップと、各パーティションについて選択されたマスタレプリカと当該マスタレプリカを保持するマスタノードとに関して他のアクティブノードへ通知するステップとである。

#### 【0057】

特に有利にはこの第2動作モードの下で如何なる特定のレプリカについてこれらのイベントが最高のレプリカ優先度を結果として生じない場合に、分散データベースシステム内のシステムマスタモニタノードはさらに以下のステップを実行するように構成されてもよい。すなわち、パーティションについての所与のデフォルトレプリカ優先度を設定されている情報を少なくとも1つのアクティブノードから収集するステップと、各パーティションについて最高のデフォルトレプリカ優先度を有する特定のノード内のレプリカを選択するステップであって、このレプリカがマスタレプリカであるとみなされ、この特定のノードが当該パーティションについてのマスタノードであるとみなされるステップと、各パーティションについて選択されたマスタレプリカと当該マスタレプリカを保持するマスタノードとに関して他のアクティブノードへ通知するステップとである。

#### 【0058】

上述の方法を実行するために、本発明は総括的にノード1~4のような複数のノードを有する分散データベースシステムを提供する。各ノードはデータの少なくとも1つのパーティションのレプリカを格納するように構成される。ノード2について図3に説明される例のように、各ノードは格納されるデータの少なくとも1つのデータパーティション112、123、143のレプリカ2101を格納するとともに、相互にアドレスを指定するために用いられる他のノードの識別子152を格納するためのデータ記憶装置15と、分散データベースシステムの他のノード1、3、4と通信するための入出力部30と、各レプリカの最新の更新2103、レプリカ状態212、223、243、各レプリカを担当するローカルリソースの状態、各レプリカの接続状態312、323、343から選択された少なくとも1つのイベントを監視するための監視部60と、データ記憶装置、監視部および入出力部と連動して、分散データベースシステムのアクティブノードのうちのどのノードが各パーティションについての現在のマスタノード2105であり、当該パーティションについての現在のマスタレプリカを担当するとみなされるかを決定するための処理部20とを備える。

#### 【0059】

さらに、この分散データベースシステムでは、各ノードの処理部20、監視部60、データ記憶装置15および入出力部30は、各レプリカの最新の更新1103、2103、3103、4103、レプリカ状態212、223、243、各レプリカを担当するローカルリソースの状態、および各レプリカの接続状態312、323、343から選択された少なくとも1つのイベントに関する情報を各アクティブノードから収集し、各レプリカについて収集されたイベントに依存して、アクティブノード内の各レプリカの優先順位を決めるための上述の事前に設定されたルールを適用し、各パーティションについて最高のレプリカ優先度を有する特定のノード内のレプリカを選択するステップし、このレプリカがマスタレプリカであるとみなされ、この特定のノードが当該パーティションについてのマスタノードであるとみなされるように構成される。



## 【 0 0 6 0 】

さらに、特に有利には分散データベースシステムにおいて如何なるレプリカについて以前のイベントが最高のレプリカ優先度を結果として生じない場合に、各ノードの処理部 20、監視部 60、データ記憶装置 15 および入出力部 30 は、パーティションについての所与のデフォルトレプリカ優先度を設定されている情報を少なくとも 1 つのアクティブノードから収集し、各パーティションについて最高のデフォルトレプリカ優先度を有する特定のノード内のレプリカを選択し、このレプリカがマスタレプリカであるとみなされ、この特定のノードが当該パーティションについてのマスタノードであるとみなされるように構成される。

## 【 0 0 6 1 】

第 2 動作モードに従って分散データベースシステムを動作するために、各ノードの処理部 20、監視部 60、データ記憶装置 15 および入出力部 30 は、アクティブノードが起動された順序を決定するための情報 1104、2104、3104、4104 を分散データベースシステムの各アクティブノードから収集するように構成され、その結果、最初に起動されたアクティブノードがシステムマスタモニタであるとみなされる。

## 【 0 0 6 2 】

特に、ノードが分散データベースシステムのシステムマスタモニタノードであるとみなされる場合に、当該システムマスタモニタの処理部 20、監視部 60、データ記憶装置 15 および入出力部 30 はさらに、各パーティションについて選択されたマスタレプリカと当該マスタレプリカを保持するマスタノードとに関して他のアクティブノードへ通知するように構成されてもよい。

## 【 0 0 6 3 】

分散データベースシステムを扱う上記の方法は、分散データベースシステムのクライアントがどのように当該システムに格納された情報にアクセスしうるか、特に読み出し動作または書き込み動作のためにアクセスしうるかに関する特定の議論が必要となりうる。

## 【 0 0 6 4 】

原則として、通信ネットワークのホームロケーションレジスタ、認証センタ、またはホーム加入者サーバのような分散データベースシステムのクライアントは分散データベースシステムの任意のノードから任意のデータへアクセスできる。しかしながら、様々なレプリカ内でデータの一貫性を維持するために、レプリカのうちの 1 つだけが読み出しおよび/または書き込みのリクエストを受信し、これはマスタレプリカであるだろう。上述のように、マスタレプリカ内のデータは他のレプリカへ時々刻々と更新される。

## 【 0 0 6 5 】

よって、分散データベースシステムの全内容は、分散データベースシステムを構成する任意のノードからアクセス可能である。この目的ために、各ノードは 1 つ以上のアクセスゲートウェイ（以下、AG）を含んでもよく、AG はマスタレプリカが位置するノードへデータの読み出し/書き込みを行うリクエストを転送することを担当するエンティティである。データベースプロトコルはアクセスプロトコルと異なりうるため、AG はデータの読み出し/書き込みを行うリクエストをクライアントから受信し、分散データベースシステム内の他のノードへアクセスすることを担当する。性能を最大化するために各ノード内に 2 つ以上の AG が提供されうるため、2 つ以上の AG の間でトラフィックを分配するためにロードバランサ（以下、LB）が提供されてもよい。しかしながら、このような LB はただ 1 つの AG を有するノード構成では必要ではない。

## 【 0 0 6 6 】

図 5 が説明するように、ノード 2 は LB 109a および 3 つの AG 191a ~ 193a を備えてもよく、第 1 パーティション 11 についてのレプリカ 112、第 2 パーティション 12 についてのレプリカ 123、および第 4 パーティション 14 についてのレプリカ 143 を担当してもよく、レプリカ 112 は第 1 パーティションについてのマスタレプリカであり、ノード 2 は第 1 パーティションについての当該マスタレプリカを担当するマスタノードである。一方で、ノード 3 は LB 190b および 3 つの AG 191b ~ 193b を

10

20

30

40

50

備えてもよく、第1パーティション11についてのレプリカ113、第2パーティション12についてのレプリカ121、および第4パーティション14についてのレプリカ141を担当してもよく、レプリカ121は第2パーティションについてのマスタレプリカであり、ノード3は第2パーティションについての当該マスタレプリカを担当するマスタノードである。

【0067】

よって、分散データベースシステムを扱うこの方法は、分散データベースシステム内のデータの読み出し/書き込みを行うためにノードにおいて受信された任意のリクエストについて、当該データが属するパーティションと当該パーティションについての現在のマスタレプリカを担当する現在のマスタノードとを決定するステップと、当該リクエストを当該現在のマスタノードヘルレーティングするステップとを含む。

10

【0068】

図5に説明される例のように、分散データベースシステム内のデータの読み出し/書き込みを行うためのクライアント5からのリクエストはノード2のような任意のノードで受信されうる。図5の例はステップS-150の間にLB190aにおいて受信されたデータの読み出し/書き込みを行うリクエストを説明する。このリクエストはステップS-151の間にAG193aに割り当てられ、AG193aは読み出し/書き込みを行うデータが第2パーティション12に属することを判定し、このAGはまた、当該パーティションを担当する現在のマスタノードがノード3であることを判定する。次いで、AG193aはステップS-152の間にリクエストをノード3ヘルレーティングする。このリクエストは当該ノード3に2つ以上のAGが存在するならばノード3のLB190bにおいて受信されてもよいし、ノード3が1つのAGのみを含むならば唯一のAGにおいて受信されてもよいし、この図5に示される例のように、設定手段によってノード3の特定のAG191bがAG193aに知られているならばこのようなAGにおいて受信されてもよい。リクエストを受信するノード3のAG191bはステップS-152の間にリクエストに従ってデータの読み出し/書き込みを行うために第2パーティション12についてのマスタレプリカ121にアクセスする。

20

【0069】

他方で、図5はまた、データが属するパーティションについてのマスタレプリカを保持するマスタノードにおいてリクエストが受信される場合を説明する。この例の場合では、データの読み出し/書き出しを行うリクエストはステップS-160の間にノード2のLB190aにおいて受信される。

30

【0070】

このリクエストはステップS-161の間にAG191aに割り当てられてもよく、AG191aは読み出し/書き込みを行うデータが第1パーティション11に属することを判定し、このAGはまた、当該パーティションを担当する現在のマスタノードがこのノード2であることを判定する。次いで、AG191aはステップS-162の間に、リクエストに従ってデータの読み出し/書き込みを行うために、第1パーティション11についてのマスタレプリカ112へアクセスするようにリクエストを内部的にルーティングする。

40

【0071】

ノード2のような、分散データベースシステムの任意のノードのLBは、図3に説明されるようなクライアント5との通信専用の入出力部50と、ノードの負荷および性能のバランスをとるために適切なAGを選択するように構成された処理部20のリソースとで構築されてもよい。この入出力部50は、分散データベースシステムの各ノードが備える入出力部30の一体部分であってもよいし、当該入出力部30に含まれる別個のモジュールであってもよい。

【0072】

ノード2のような、分散データベースシステムの任意のノードのAGは図3に説明されるような他のノード1、3、4との通信専用の入出力部40と、読み出し/書き込みを行

50

うデータが属するパーティションを決定し、当該パーティションについてのマスタレプリカを担当するマスタノードが現在のノードであるかまたは分散データベースシステムの別のノードであるかを判定し、当該パーティションについてのマスタレプリカを担当するマスタノードが現在のノードである場合に当該データへアクセスするためにデータ記憶装置 15 へアクセスし、それ以外の場合に現在のマスタノードであると判定された別のノードへリクエストをルーティングするように構成された処理部 20 のリソースとで構築されてもよい。

【0073】

本発明の実施形態に従うと、監視部 60 は上述のイベントを監視し蓄積する一意のユニットである、本明細書におけるいわゆるローカルシステムモニタ（以下、LSM）を含んでもよく、またはアクティブ LSM とスタンバイ LSM とを含み、前者に障害が生じた場合に後者が動作を引き継げるようにしてもよい。以下において、監視部または LSM への如何なる言及も、言及されているノード内のアクティブ LSM を意味すると解釈される。

【0074】

特に、本発明が上述の第 2 動作モードに従って動作する場合に、いわゆるシステムマスタモニタノードの監視部 60 がコントローラシステムモニタ（以下、CSM）であるとみなされるが、分散データベースシステム内の他のノードの各監視部はなおも LSM として言及される。よって、CSM は分散データベースシステム内の各ノードの LSM から受信されたイベント情報を考慮に入れ、且つアクティブノード内の各レプリカの優先順位を決めるための事前に設定されたルールを適用することによって、マスタレプリカでどれであるかを決定してもよい。CSM は各パーティションについてのマスタノードが何であるかを他のノード内の各 LSM と通信する。

【0075】

第 2 動作モードにおいてアクティブノードのうちのどのノードがシステムマスタモニタとみなされるかを決定する際に、図 6 に説明される例示の状態マシンが適用される。状態間の遷移は他のノードからの「アライブ」メッセージの受信か、タイマの満了かのいずれかに起因する。簡単のために、分散データベースシステム内の他のノードに対してシステムマスタモニタノードを言及するのではなく、以下の議論ではこれらの別個の監視部 60、すなわち当該他のノードの各 LSM に対するシステムマスタモニタノードの CSM を言及する。

【0076】

図 6 に説明される本発明の実施形態では、取りうる状態および遷移は以下でありうる。

非アクティブ。これは各 LSM が開始した際の状態である。LSM がこの状態にあるならば、ノードはリクエストにตอบสนองせず、マスタレプリカのホストになりえない。この状態へ遷移すると、LSM は残りのノードへのアライブメッセージの送信を開始する。

アクティブ。アクティブ状態では、各 LSM は他のノードに関する情報を有する「アライブ」メッセージをリッスンし、特に存在するならばマスタレプリカを有する CSM からの「アライブ」メッセージをリッスンする。各 LSM はまた、自身のノードに関する情報を送信し、自身のノード内で任意のローカル AG へマスタレプリカに関する情報を内部的に分配し、場合によっては「アライブ」メッセージの情報を残りのノードへ転送する。

ポテンシャル CSM。所定の設定可能な期間、すなわち本明細書におけるいわゆる DELAY\_\_TIME の間、LSM がこの状態に留まるならば、このような LSM は CSM になる。これはまた、DELAY\_\_TIME が満了する前にシステム内の残りのすべてのノードから「アライブ」メッセージを受信する場合にも生じる。

CSM。この状態に到達するノードは分散データベースシステム内のすべてのパーティションについてどのレプリカがマスタレプリカであるかを決定する。

【0077】

これにもかかわらず、状態、遷移、またはその両方の他の実施形態が同様に予見可能である。

【0078】

10

20

30

40

50

簡単のために任意の状態から非アクティブ状態への遷移が示されていない図6に示される状態間の遷移に関して、任意のLSMは非アクティブ状態から開始し、LSMを有するノードが分散データベースシステム内のすべての他のノードへ「アライブ」メッセージを送信するとすぐに、このLSMはアクティブ状態へ移行する。すなわち、非アクティブ状態からアクティブ状態への遷移ST-1は、非アクティブ状態のLSMを有するノードから分散データベースシステム内のすべての他のノードへ「アライブ」メッセージを送信することである。

#### 【0079】

2つ以上のサブネットワークがアクティブになることを回避することによってスプリットブレイン状況における一貫性を解決するために提供される本発明のオプションの実施形態に従うと、非アクティブ状態への遷移は、1つ以上のノードに障害が発生し、CSMを含む $(n+1)/2$ 個未満の稼働中のノードが存在するか、CSMを含まない $n/2+1$ 個未満が存在する場合に発生してもよい。ここで「n」は分散データベースシステム内の全ノード数である。例えば、3つのノードで構成される分散データベースシステムがあり、CSMのホストであるノードが隔離されるならば、そのサブネットワークは1つだけ、すなわち $(3+1)/2=2$ 個未満のノードを有する。よって、CSMのホストであるこのノードは非アクティブ状態へ移行する。他方のサブネットワークは2つのノードを有し、これは $2/2+1=2$ 個以上のノードを有することを意味し、よってこれらはアクティブに留まり、新たなCSMが選択される。このオプションの構成では、各LSMは、設定可能な期間である、本明細書におけるいわゆるINACTIVE\_\_TIMEの後に各LSMが1つ以上のノードから「アライブ」メッセージを受信しないならば、当該1つ以上のノードがダウンしていることを検出する。非アクティブへの遷移により、任意のLSMは自身の実行時間をリセットし、すなわち再び「若く」になる。これは、隔離された以前のCSMが再びCSMになり、使用されるべきでない情報を送信することを防ぐ。

#### 【0080】

このオプションの実施形態は分散データベースシステムを適切に構成することによって実行されうる。このような実施形態が望まれないならば、適切な設定パラメータがリセットされ、その結果、システムから分離された複数のノードのxxxが、ノードをアクティブ状態へ移行させ、ノード間で分離されている場合であっても2つ以上のサブネットワークで動作させることに無関係になる。

#### 【0081】

図6に示される状態間の遷移に戻り、上述のオプションの実施形態が動作するように構成されていると仮定すると、アクティブ状態にあるLSMを有する任意のノードは、他のノードから十分な「アライブ」メッセージが得られないならば、すなわちCSMを含む稼働中の $(n+1)/2$ 個未満またはCSMを含まない $n/2+1$ 個未満であるならば、非アクティブ状態へ移行できる。

#### 【0082】

そうではなく、他のノードから十分な「アライブ」メッセージが受信され、送信側ノードが後から起動したことを受信された「アライブ」メッセージ内の情報が示し、且つメッセージを送信するノードが残っていない場合に、受信側ノード内のLSMはCSMになる。すなわち、アクティブ状態からCSM状態への遷移ST-2.1は、他のノードからの十分な「アライブ」メッセージの受信であり、分散データベース内にメッセージを送信するノードが残っていないことである。特に2ノードシステムにおいて、メッセージが受信されることなくDELAY\_\_TIMEが経過するかも知れず、同様に当該ノードはCSMになる。しかしながら、他のノードから、当該ノードが後から起動したことを示す「アライブ」メッセージを受信したが、「アライブ」メッセージを受信できると予想されるノードがさらに存在する場合に、遷移ST-2.2が行われ、受信側ノード内のLSMはポテンシャルCSMになる。

#### 【0083】

以前の遷移について説明したのと同じ理由のために、ポテンシャルCSM状態から、ノー

10

20

30

40

50

ドは非アクティブ状態へ戻りうる。そうでなければ、自身のLSMが確定したCSMであることを示す、さらに予想された「アライブ」メッセージがノードから受信された場合、または他のLSMが先に起動したことを示す、さらに予想された「アライブ」メッセージがノードから受信された場合に、図6に説明される遷移ST-3.1が行われ、ポテンシャルCSM状態のLSMは、これらの予想される「アライブ」メッセージのいずれかを受信するノードにおいてアクティブ状態へ移行する。他方で、古いノードを示すさらに予想された「アライブ」メッセージをさらに受信することなくいわゆるDELAY\_\_TIMEが経過した場合、または先に起動したノードが存在しないことを示す「アライブ」メッセージが残りのノードから受信された場合に、図6に説明される遷移ST-3.2が行われ、ポテンシャルCSM状態にあるLSMは、DELAY\_\_TIMEが経過したか、これらの予想された「アライブ」メッセージのいずれかが受信されたノードにおいてCSM状態へ移行する。

10

**【0084】**

上述のように、2種類のタイマが存在する。1つは本明細書におけるいわゆるDELAY\_\_TIMEタイマであり、自身がCSMであると宣言する前に古いLSMに関して通知する、残りのノードの「アライブ」メッセージをLSMが待つ時間である。もう1つは本明細書におけるいわゆるINACTIVE\_\_TIMEタイマであり、ノードがダウンしており利用不可能であることを結論付ける前にこのようなノードからの「アライブ」メッセージをLSMが待つ時間である。

**【0085】**

20

上記の処理が終了すると、CSMとして確定された監視部60を有するノードはマスターレプリカのリストを含む「アライブ」メッセージの送信を開始する。

**【0086】**

図6に説明される状態マシンの処理に加えて、図7は、3つの例示のノードの中のどのノードがCSMとして確定された監視部60を有するノードであるかを決定するための分散データベースシステムのノード間での例示の動作シーケンスを示す。上述のように、CSMは最初に起動されたLSMである。これを決定するために、各LSMは起動時に図2に示された例のように動作時間2104、4104を有する「アライブ」メッセージを残りのノードへ送信する。システムの一部であるノードは設定によって知られている。「アライブ」メッセージを送信した後に、LSMは所定の期間であるDELAY\_\_TIMEの間、他のノードから「アライブ」メッセージが受信されるのを待つ。CSMを確立するフェーズは、この時間が経過した場合、または他のノードから十分な「アライブ」メッセージが受信されたならばその前に終了する。このフェーズの間に、この時間までに受信された情報に従ってポテンシャルCSMが割り当てられるかもしれないが、このフェーズが終了するまでそれは確定されないだろう。

30

**【0087】**

よって、図7に説明されるように、ステップS-200の間に起動される最初のノードはノード1である。LSM、すなわちノード1の監視部60が起動するとすぐに、ステップS-205の間に分散データベースシステムの他の2つのノード1、2へ「アライブ」メッセージを送信し、DELAY\_\_TIME秒のタイマを開始する。他のノードはまだ稼働していないので、対応する各LSMはこのような「アライブ」メッセージを受信しない。

40

**【0088】**

次いで、ノード2のLSMはステップS-210の間に稼働を開始し、ステップS-220の間に「アライブ」メッセージをノード1、3へ送信し、DELAY\_\_TIMER秒の自身のタイマを開始する。ノード2のLSMが起動した後ではあるが、ノード2のLSMが「アライブ」メッセージを送信する前に、ノード3のLSMはステップS-215の間に起動し、DELAY\_\_TIMER秒の自身のタイマを開始する。

**【0089】**

ノード1がノード2から「アライブ」メッセージを受信する場合に、ノード1はステッ

50

ブ S - 2 2 5 の間に自身をポテンシャル C S M として認定する。なぜなら、ノード 1 が有する情報を用いて、ノード 1 は、ノード 3 からの情報を待つものの、早くに起動した L S M だからである。この段階で、ノード 3 はステップ S - 2 3 0 の間にノード 1、2 へ「アライブ」メッセージを送信できる。

【 0 0 9 0 】

ノード 2 がノード 3 から「アライブ」メッセージを受信する場合に、ノード 2 はステップ S - 2 4 0 の間に自身をポテンシャル C S M として認定する。なぜなら、ノード 2 が有する情報を用いて、ノード 2 の L S M はノード 3 の L S M よりも早く起動したからである。同様に、ノード 3 がノード 2 から「アライブ」メッセージを受信する場合に、ノード 3 はステップ S - 2 4 5 の間にノード 2 をポテンシャル C S M として認定する。なぜなら、ノード 3 が有する情報を用いて、ノード 1、2 はノード 1 からの情報を待っており、ノード 3 にはまだ知られていないが、ノード 2 の L S M は早くに起動したからである。

【 0 0 9 1 】

ノード 1 がノード 3 から「アライブ」メッセージを受信する場合に、ノード 1 はステップ S - 2 3 5 の間に自身を確定した C S M として認定する。なぜなら、ノード 1 はシステム内のすべてのノードから「アライブ」メッセージを受信しており、この情報を用いて、ノード 1 の L S M は早くに起動したものであるからである。よって、ノード 1 はステップ S - 2 5 0 の間に、ノード 1 が確定した C S M であることを通知する「アライブ」メッセージを他のノードへ送信する。ノード 1 から最終的に「アライブ」メッセージを受信するノード 2、3 はステップ S - 2 5 5 の間にノード 1 の監視部が C S M であり、現在の状況において他のノードがこの役割を負わないことを認識し、ステップ S - 2 6 0 の間に記録する。

【 0 0 9 2 】

障害の場合に C S M の役割を再割当てできるように、「アライブ」メッセージは各ノードから残りのノードへ定期的送信される。このように、ノードの追加または削除のような、分散データベースシステムの構成の任意の変化が即座にすべてのノードにおいて知られうる。

【 0 0 9 3 】

よって、本発明の両方の実施形態、すなわち両方の動作モードに従うと、各ノードは、他のノードから受信された情報に基づいて、どのノードが C S M であるとみなされるか、すなわちアクティブノードの中で、より長い動作時間を有するものがどれかを決定してもよい。ノードが C S M になるために、このようなノードは、任意の他のポテンシャル C S M から遠隔の「アライブ」メッセージを受信するための何らかの時間が存在することを保証するために、少なくとも D E L A Y \_ T I M E を待つことができるだろう。C S M 状態に到達したノードは、曖昧さが存在しないことを保証するために残りのノードへ自身の決定を通信してもよい。しかしながら、すべてのノードが第 1 動作モードに従って振舞う場合にこの通信は必要ない。この通信は信頼性を保証するために T C P ベースであってもよい。再設定時間を低減するために、この振る舞いは、C S M が他のノードへレプリカ構成を通信する際に間接的に実現されてもよい。この点において、レプリカ状態は「アライブ」メッセージに含まれるので、ひとたび選択されると、C S M はどのレプリカ状態の設定が動作に適するかがわかる。実際に、この情報はすべてのノードに知られるかもしれないが、これらは第 2 動作モードのもとで C S M が確認するのを待ってもよい。

【 0 0 9 4 】

各ノードは、当該ノードが「アライブ」メッセージを受信した、現在のノードを有するリストを、当該メッセージ内で受信されたレプリカ状態とともに管理できるだろう。「アライブ」メッセージが受信されるごとに、送信側ノードがアクティブとして設定されてもよい。パーティションについてのマスタノードはアクティブノードの中だけから選択されてもよい。ノードは、それからメッセージを受信することなく、いわゆる I N A C T I V E \_ P E R I O D という期間が経過している場合に利用可能でないとして設定されてもよい。この時間はメッセージ受信時間の平均時間の 2 倍または 3 倍でありえ、これは最初に

前述の D E L A Y \_ T I M E に設定されうる。ノードはまた、それに送信された「アライブ」メッセージが届かない場合に利用不可能であるとして設定されてもよい。このようにノードの利用可能性は非常に高速に検出される。

【0095】

一般的に言うと、C S Mであるとみなされる監視部60を有するノードは、分散データベースシステム内の他のノードへ、自身のノード識別子、動作時間、各レプリカについてのレプリカ状態、各レプリカについてのマスタノード、状態マシンからのノード状態、(オフラインでさえも)「アライブ」メッセージが受信されているアクティブノードのリスト、現在のマスタレプリカ情報(以下、M R I)、およびM R Iを設定するC S Mについての(実行時間を含む)更新時間およびノードIDを送出してもよい。

10

【0096】

特に、いわゆるM R Iメッセージはホストノードを有するレプリカのリストを含みうる。本発明の実施形態では、M R Iメッセージはまた、M R Iについての確認を受信するために、オフラインのノードへも送信されてもよい。この観点で、オフラインのノードは実C S MノードといわゆるサブC S Mノードとの間のリンクであってもよい。後者は、マスタであると信じているがそうではないノード内で稼動しているL S Mである。いずれにせよ、上述のように、オフラインのノードのレプリカはマスタとして設定されえない。従って、この問題を回避するために、C S M処理はM R Iを送信する前にD E L A Y \_ T I M Eを待ってもよい。

【0097】

20

一般的に言うと、C S M選択メカニズムは、すべての監視プロセスが同期されるようなものである。上述のように、最も古いプロセスがC S Mになることがアイデアである。最も古いL S Mがどれかを決定する際に、複数の実施形態が予見可能である。第1実施形態では、本発明は動作時間を考慮に入れる。すなわち、各プロセスは、自身のローカル時間を用いて、起動から何秒動作しているかを決定し、この情報を「アライブ」メッセージ内で送信する。この場合に、恐らくはレイテンシが重要な役割を果たすため、以下の欠点が存在するかもしれない。受信側ノードは送信された作業時間を見るだろうが、レイテンシを見ないので、受信側ノードが送信側ノードよりも若いか古いかを正確に判断することは難しいだろう。レイテンシ時間はピングメッセージを介して測定されえ、平均が確立されうるだろう。

30

【0098】

第2実施形態では、本発明は起動時刻を考慮に入れる。この実施形態の下で、すべてのプロセスは自身の起動時刻を「アライブ」メッセージ内で送信する。この時刻は各ノードにおけるローカルマシン時刻に関連するだろう。従って、システム内のすべてのマシンを同期する必要があるかもしれない。リナックスシステムおよび一般のオペレーティングシステムはN T Pを用いて何年も前にこの問題を解決している。

【0099】

これらの第1実施形態および第2実施形態の下で、C S M選択処理の間に、図6に説明される例示の状態マシンについての非アクティブ状態にL S Mが到達するといつても、個別の動作時間または起動時刻は、劣化したC S Mが回復した際に再びC S M状態に到達することを回避するためにリセットされる必要があるかもしれない。

40

【0100】

上述のように、分散データベースシステムの任意のノードの起動または停止の際に、アクティブノードの中で、どのノードが各パーティションについてのマスタレプリカを担当する現在のマスタノードであるかを判定するステップが存在する。このマスタノードがどれであるかを判定する際に、本発明は2つの実施形態を提供する。第1実施形態では、どれがマスタノードであるかを判定するためにすべてのノードは独立して動作する。第2の実施形態では、各パーティションについてどれがマスタノードであるかを決定するためにC S Mが決定される。

【0101】

50

図 8 は、例ではノード 3 である現在みなされている C S M がダウンして、他のノード、すなわちノード 1、2 に対して利用不可能になる例示の状況を説明する。

【 0 1 0 2 】

図 8 に示されるように、C S M として動作するノード 3 はステップ S - 3 0 0 の間に他のノード 1、2 へ「アライブ」メッセージの最新の集合を送信した。このようなメッセージがそれぞれ受信されると、ノード 1、2 の両方は、新たな「アライブ」メッセージがノード 3 から受信されるまでリセットされないいわゆる非活動期間を開始する。ノード 3 からさらに「アライブ」メッセージを受信することなくノード 1、2 の両方において非活動期間が満了するため、ノード 1、2 の両方はそれぞれステップ S - 3 2 0、S - 3 3 0 の間にノード 3 を利用不可能として記録する。この例示の場合では、ノード 2 はノード 1 よりも古く、この情報はノード 2 に知られている。なぜなら、ノード 2 はステップ S - 3 1 0 の間にノード 1 から最新の定期的な「アライブ」メッセージを受信しているからである。従って、ノード 2 はステップ S - 3 4 0 の間に、図 6 に説明される状態マシンに従って、自身をポテンシャル C S M 状態へ移行する。次いで、ノード 2 はステップ S - 3 5 0 の間に、自身の定期的な「アライブ」メッセージをノード 1 とノード 3 とへ送信するが、後者はこのようなメッセージを受信できない。ノード 2 から「アライブ」メッセージを受信したノード 1 は、ノード 2 が古いことを認識し、ステップ S - 3 6 0 の間にノード 2 が現在の C S M であるという結論にいたる。ノード 2 においていわゆる D E L A Y \_ T I M E が満了した後に、ノード 2 の監視部 6 0 はステップ S - 3 7 0 の間に C S M として確定され、ノード 2 は現在のシステムマスタモニタである。

【 0 1 0 3 】

さらに、図には示されていないが、ノード 3 が再び復旧した場合に、ノード 3 は自身の定期的な「アライブ」メッセージをノード 1、2 へ送信する。ノード 1、2 においてこのようなメッセージを受信すると、これらはともに、ノード 3 を再び利用可能として記録する。上述のように、ノード 3 はオプションとして自身の動作時間または起動時刻をリセットし、その結果、状況は変化せず、D E L A Y \_ T I M E の期間にポテンシャル C S M 状態にあった後に、現在の C S M はなおもノード 2 の監視部であり、これは図 6 に示される状態マシンに従って C S M 状態へ到達する。

【 0 1 0 4 】

以前に分散データベースシステムに含まれていなかった新たなノードの追加に関して、C S M 選択の観点から、状況は、ノードがダウンし再び復旧した状況と非常に類似している。唯一の違いは、障害が発生したノードは設定テーブルに提示されているため C S M に知られているが、全く新たなノードは既存のノードに知られていない点である。この点において、C S M だけでなく他の監視プロセスも、設定されていないノードから受信された「アライブ」メッセージに注意を払わない。従って、新たなノードを導入する重要な側面は、新たなノードの識別子を図 1 B に説明されるような既存のノードの対応する設定テーブル 1 5 1 ~ 1 5 4 に含めることと、このような新たなノード内にレプリカの構成を含めることとである。これが完了すると、新たなノードの監視部は起動しうる。このように、新たなノードは他の既存のアクティブノードへ「アライブ」メッセージを送信し、他の既存のアクティブノードは新たなノードを認識するだろう。C S M は新たなノード内の新たなレプリカを認識することになるので、これに依じて C S M は必要ならばシステムを再構成する。上記の第 1 動作モードに従って、すなわち C S M を有せずにシステムが動作する観点で、すべてのノードはそれぞれ、「アライブ」メッセージ内で受信された情報を処理し、各パーティションについてのマスタレプリカを担当するマスタノードがどれであるかに関して同じ結論に到達する。原則として、「アライブ」メッセージはいわゆる M R I メッセージ内に情報を有する。それにもかかわらず、M R I メッセージは「アライブ」メッセージとは別に送信されてもよい。

【 0 1 0 5 】

他方で、稼働中の分散データベースシステムへの新たなパーティションの追加は上述の振る舞いを用いて簡単な作業でありうる。1 つのノードにただ 1 つのレプリカを有する場

10

20

30

40

50



合よりも高可用性を得るために、それぞれ2つの既存のノードへ追加される少なくとも2つのレプリカ上にパーティションを複製することを考えてもよい。上述のように、任意のノードの構成について、レプリカとは別に、レプリカごとの上述のすべての対応するインジケータが同様に構成されなければならない。よって、第2動作モードでのCSMを有する選択されたシステムマスタモニタノード、または第1動作モードでのすべてのノードは、遅かれ早かれすべてのノードから「アライブ」メッセージを受信する。ここで、新たなレプリカについての情報は当該新たなレプリカのホストとなるノードから受信される。

【0106】

高可用性を有する分散データベースシステムを効率的に提供するために、監視プロセスも高可用性を有することが期待される。従って、上述のように、2つのLSMであるアクティブLSMとスタンバイLSMとが各ノードで動作してもよい。動作中に、アクティブLSMはすべての受信メッセージ、特に「アライブ」メッセージと、任意のMRIメッセージが受信されるならこのようなMRIメッセージとをスタンバイLSMへ転送する。このようにして、アクティブLSMがクラッシュした場合に、スタンバイLSMが即座に引継ぎ、その結果、単に監視プロセスに障害が発生したからといってノードはダウンすることがない。

【0107】

本発明の目的のために、監視プロセスは場合によっては処理部20と連動する監視部60によって実行されてもよい。

【0108】

上記とは別に他の個別の実施形態が予見可能である。例えば、設定データとともに格納される代わりに、すべての動的データが動的に管理され、時々刻々と更新されてもよい。この動的データはMRI情報だけでなく、レプリカおよび通常は「アライブ」メッセージで送信されるノードに関する情報を含む。これは、リクエストに回答するプロセスで用いられるポートとは異なるポートにおけるマルチキャスト技術を用いて容易に実現されうる。

【0109】

この点において、他のノードから受信されたすべての「アライブ」メッセージだけでなく、設定されたノードへ送信された「アライブ」メッセージは、マルチキャストにより送信されてもよく、特別な処理は必要ない。同様に、ノードの任意のLSMにより受信された、またはシステムマスタモニタのCSMにより送信されたMRIメッセージはマルチキャストにより受信または送信されてもよい。

【0110】

MRIは2相コミットであるため、以下の特別な検討が必要になる。

LSMプロセスであり、MRI\_\_ACKが送信されていないならば、CSMが接続を再開しうるため問題はない(いずれにせよ、ちょうど中間に発生する可能性は非常に低い)。マルチキャストによって何も送信されない。

LSMプロセスであり、MRI\_\_ACKが送信されているならば、「確認待機中」を示すフラグを有するMRIをマルチキャストする(MRI\_\_NACKについて、確認されないだろうために何も行われてはならない)。このように、スタンバイプロセスは確認を理解することができてよい。

LSMプロセスであり、MRI\_\_CONFIRMが受信されているならば、「確認済」としてのフラグを有するMRIをマルチキャストする。スタンバイプロセスは直近のものとしてMRIを解釈してもよい。

CSMプロセスであり、MRIがどのノードへも送信されていないならば、何もマルチキャストされてはならない(アクティブになるので、スタンバイプロセスは自身の現在のMRIが精密なものに一致しないことを検出してもよく、プロセスを開始してもよい)。

CSMプロセスであり、MRIが何れかのノードへ送信されているならば、「確認のペンダント」としてそれをマルチキャストする。スタンバイが制御を得るならば、すべて

10

20

30

40

50

の L M との接続を再開してもよく（すでに開いている接続はタイムアウトしてもよく、リモート S M プロセスは最初の M R I 表示を無視してもよい）、M R I を再送する。

C S M プロセスであり、M R I \_ C O N F I R M が何れかのノードへ送信されているならば、「確認中」というフラグを有する M R I をマルチキャストする。スタンバイが制御を得るならば、再び確認するためにすべての L S M プロセスとの接続を再開してもよい（すでに確認されたものはこのメッセージを無視してもよい）。

C S M プロセスであり、M R I \_ C O N F I R M がすべてのノードへ送信されているならば、「確認済」として M R I をマルチキャストし、その結果、スタンバイはこれを新たなものとして解釈してもよいが、確認手順を開始しなくてもよい。

【 0 1 1 1 】

10

上述のように、M R I メッセージは、構文解析のための共通コードを用いるために、末端において「確認」状態についての追加のフラグをマルチキャストするために必要であってもよい。いずれにせよ、すべてのマルチキャストメッセージは、アクティブとスタンバイとの間の競合状態を解決するために、プロセスの重みを表す 1 バイトを最後の部分に含んでもよいだろう。

【 0 1 1 2 】

このセクションは、アクティブプロセスとスタンバイプロセスとがどのようにそれらの状態を検出しえるか、そしてこれらが何を行うことが想定されるかを説明する。プロセスが起動するときには常に、分散データベースシステム内の任意のマルチメディアメッセージを受信するためにリッスンを開始してもよい。何らかのマルチキャストが受信されると、スタンバイとして設定されてもよいし、受信した状態に従って何れかの内部データを更新してもよい（自身の選択された状態がアライブメッセージから受信されるので、状態マシンを処理しなくてもよい）。プロセスが起動後または以前のマルチキャスト受信パケットの後、D E L A Y \_ T I M E 期間内に何もマルチキャストメッセージを受信しないならば、「アライブ」メッセージまたは M R I メッセージについてリッスン（ポートのオープン）を開始してもよく、処理重み（これはノード内の任意のプロセスについて異なるだろう）とともに（最初の始動時に空であるだろう）自身の現在の M R I を送信してもよい。アクティブプロセスが自身のものよりも大きな重みを有するマルチキャストメッセージと空でない M R I とを受信したならば、以前のポートでのリッスンを停止してもよく、スタンバイへ低下してもよい（これは最初の競合状態を解消してもよく、これは通常動作の間に発生しないと想定される）。空の M R I を比較することの目的は、非常に速く再開される場合に制御を得るために、より大きな重みを有するプロセスを回避することである（いずれにせよ、再起動の後に D E L A Y \_ T I M E を待つことができ、それによってこれは、現在のアクティブプロセスが自身の「アライブ」メッセージを送信するのと同様、且つ隣接ノードも自身のマルチキャストされる「アライブ」メッセージを送信するのと同時に、それが起動しない限りこれは発生し得ないだろう）。

20

30

【 0 1 1 3 】

以下の表は本発明の実施形態に従う上記の分散データベースシステムのノード間で用いられるプロトコルの取りうる実装を説明する。

【 0 1 1 4 】

40

位置(バイト)	長さ(バイト)	名称	説明
0	1	Msgタイプ	0x01はアライブメッセージである。
1	1	バージョン	
2	2	ノードId	ノードId
4	4	実行時間	現在のノードの実行時間
8	2	m	ローカルで監視されるレプリカ数。利用可能なレプリカだけが送信される。
10	2	レプリカマスタ状態	送信されるレプリカがノード内のマスタとしてみなされるか否かを通知する2バイト論理和。最初のノードが最初のビット(2<SUP>0</SUP>)でありバイトはビッグインディアンであるとみなされる。 例0x5(0.0101)は送信される1番目のノードと3番目のノードとがマスタであることを意味する。
12	1	ノード状態	ノード状態イベントマシンに従って非アクティブ、アクティブ、ポテンシャルOSM、OSM(1,2,3,4)。
P+13	2	レプリカid	送信される最初のレプリカのid。
P+1	1	メモリ使用率&D状態	レプリカに使用されるメモリの割合
P+3	2	レプリカidグループ	第2レプリカのid
P+5	1	メモリ使用率	第2レプリカのメモリ使用率
...			残りの利用可能なレプリカ
P+3*(m-1)	2	レプリカid	最後のレプリカのid
P+3*(m-1)+2	1	メモリ&St	最後のDについての割合
P+3*m	2	ノードn	利用可能なノード数
Q = P + 3*m + 2	2	ノードid	最初の利用可能なノードのid
Q+2	2	ノードid	2番目の利用可能なノードのid
...			
Q + 2*(n-1)	2	ノードid	最後の利用可能なノードのid
Q + 2*n	8+5*D	MRI	マスタレプリカ情報。 Dはこのフィールド内のレプリカ数を表す。

10

20

30

40

## 【 0 1 1 5 】

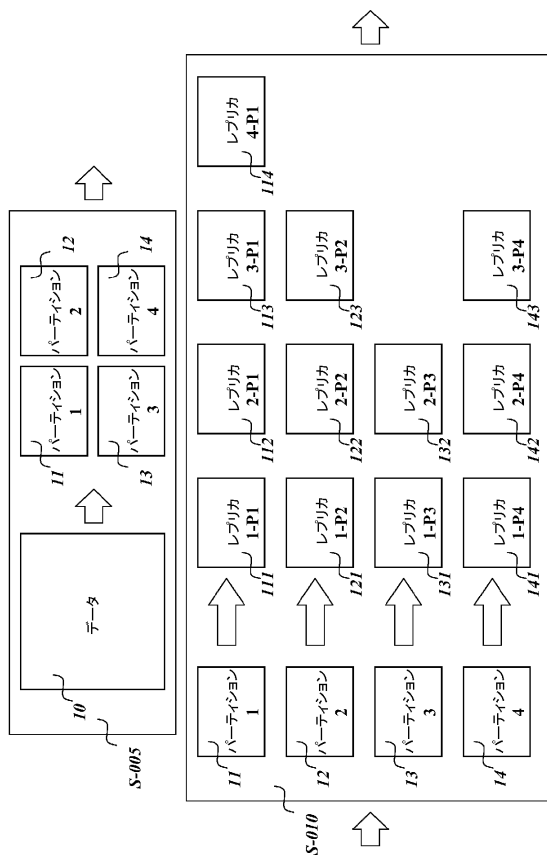
本発明は、入出力部だけでなく処理部を有するコンピュータの内部メモリにロード可能なコンピュータプログラムで実施されてもよい。このコンピュータプログラムはこの目的のためにコンピュータで実行される場合に上述の方法のステップを実行するように適応された実行可能なコードを有する。特に、実行可能なノードはコンピュータにおいて読み取り可能な媒体手段に記録されてもよい。

## 【 0 1 1 6 】

50

本発明は例示であり非限定的であることが意図される様々な実施形態に関連して上述された。当業者がこれらの実施形態を変形してもよいことが期待される。本発明の範囲は明細書と図面に連動して特許請求の範囲により規定され、特許請求の範囲に含まれるすべての変形が本明細書に含まれることが意図される。

【図 1 A】



【図 1 B】

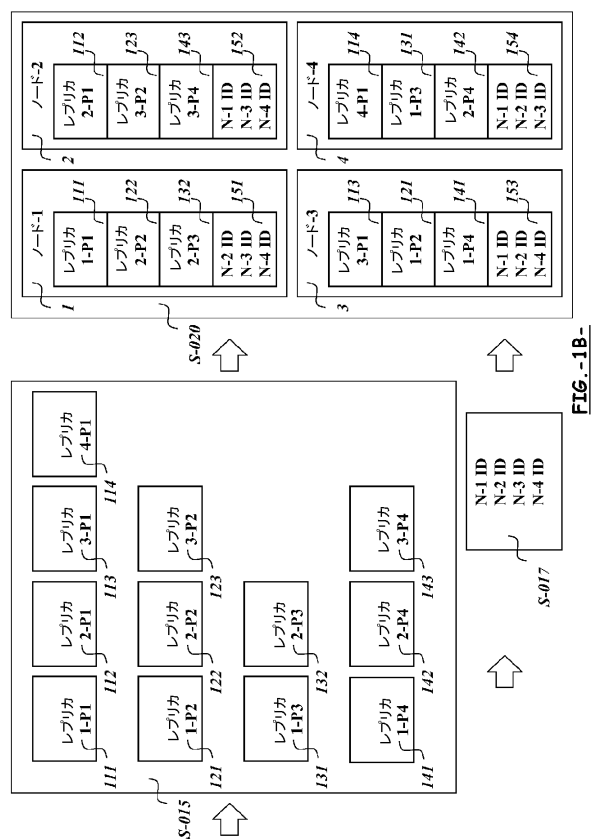


FIG. 1A-

FIG. 1B-

【図 1 C】

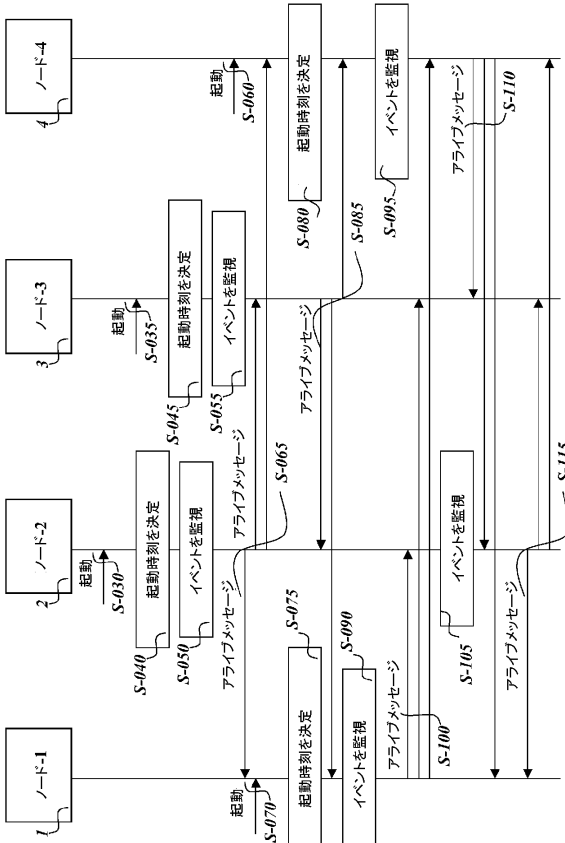


FIG. -1C-

【図 2】

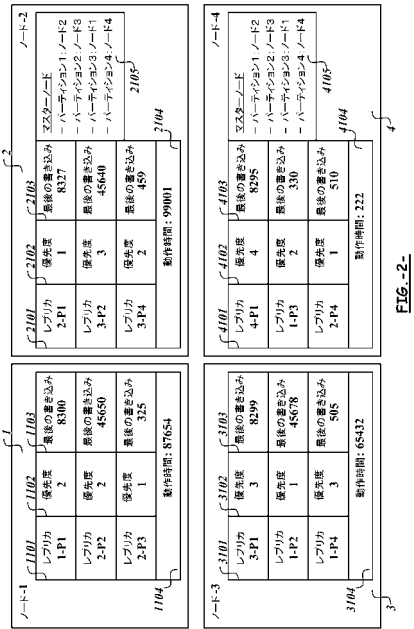


FIG. -2-

【図 3】

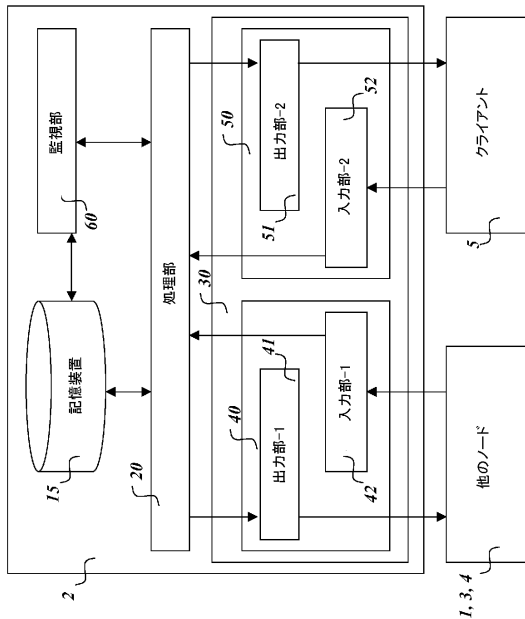


FIG. -3-

【図 4】

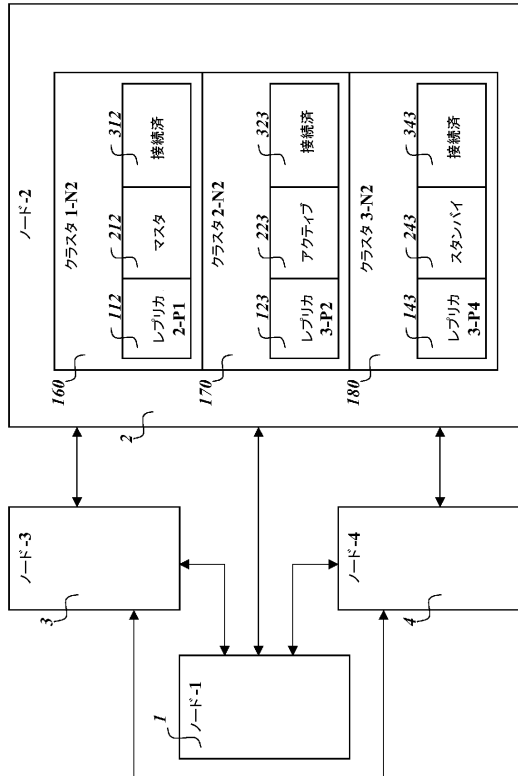
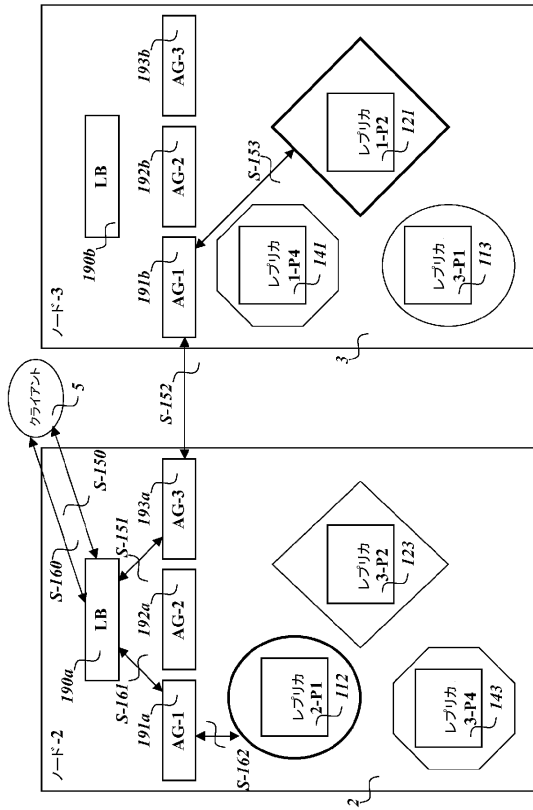
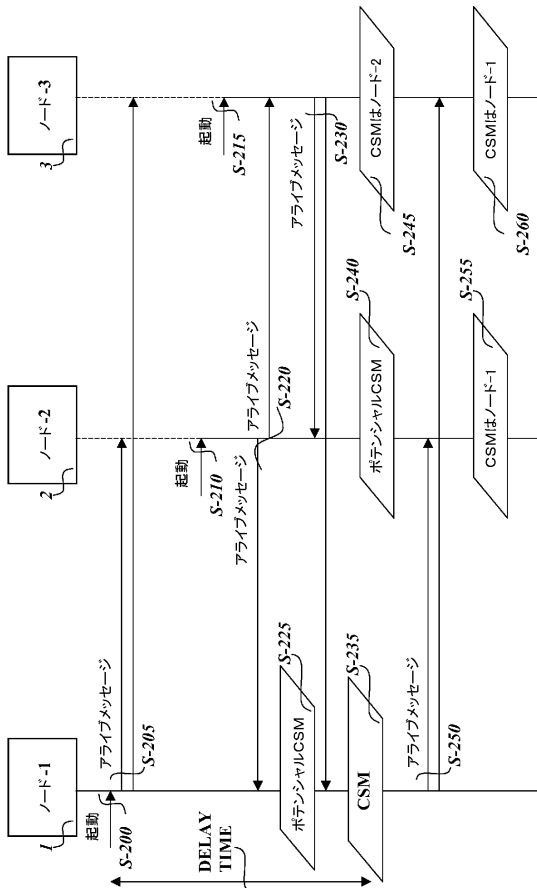


FIG. -4-

【 図 5 】



【圖 7】



【 図 6 】

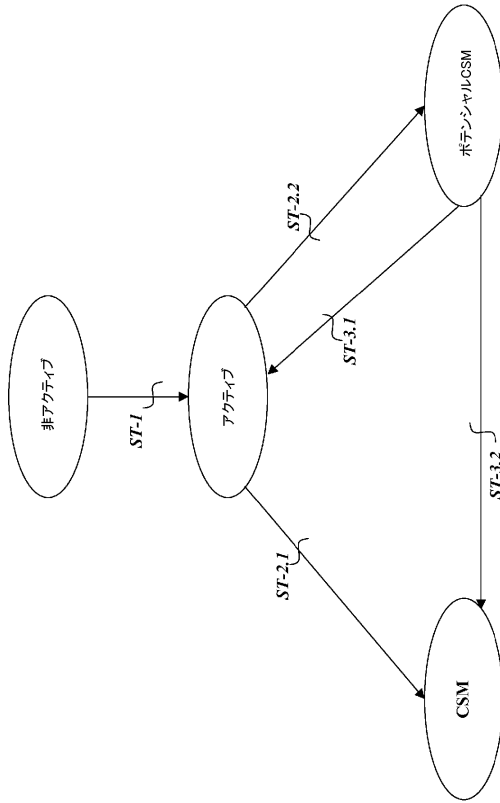


FIG.-6-

【 図 8 】

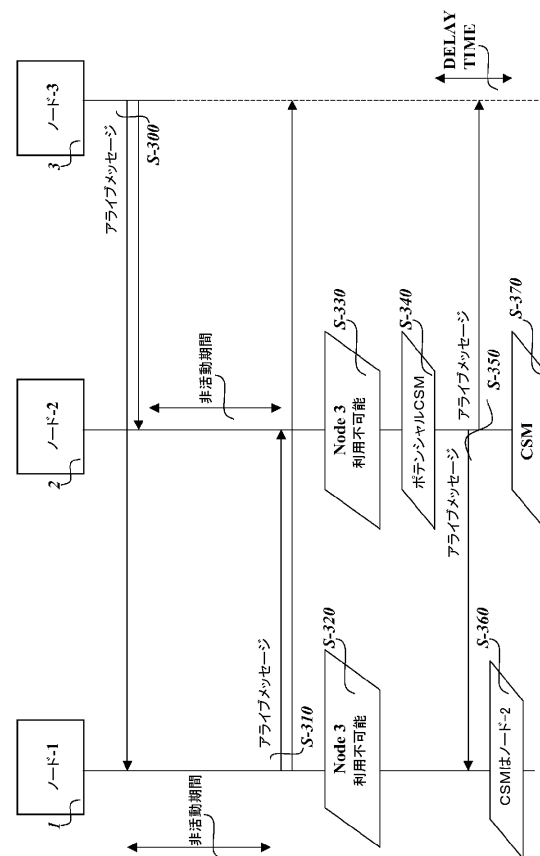


FIG.-8-

---

フロントページの続き

(74)代理人 100161399

弁理士 大戸 隆広

(72)発明者 ヘンリックセン, デニス

デンマーク国 コペンハーゲン ディーケー - 2 2 0 0 , プロフスゲーテ 2 0 , 2 . ティー  
エイチ

(72)発明者 ネヴァド ヒメネス, ホルヘ

スペイン国 マドリード イー - 2 8 4 3 6 , アルベデレテ, プラド デ ラ イグレシア  
2 , ビージェイ 4

(72)発明者 サン マルティン アリバス, マルタ

スペイン国 マドリード イー - 2 8 0 0 7 , カヴァニリエス 6 0

審査官 加内 慎也

(56)参考文献 米国特許第 0 6 5 3 9 3 8 1 ( U S , B 1 )

米国特許第 0 7 7 8 8 2 3 3 ( U S , B 1 )

Ruben de Juan-Marin ET AL , Revisiting Hot Passive Replication , Second International Co  
nference on Availability, Reliability and Security (ARES'07) , 2 0 0 7 年 4 月 1 日

(58)調査した分野(Int.Cl. , D B 名)

G 0 6 F 1 2 / 0 0

G 0 6 F 1 3 / 0 0