

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4900982号
(P4900982)

(45) 発行日 平成24年3月21日(2012.3.21)

(24) 登録日 平成24年1月13日(2012.1.13)

(51) Int.Cl. F I
G06F 11/20 (2006.01) G O 6 F 11/20 3 1 O C
G06F 9/50 (2006.01) G O 6 F 9/46 4 6 5 Z

請求項の数 15 (全 14 頁)

(21) 出願番号	特願2010-528436 (P2010-528436)	(73) 特許権者	390009531
(86) (22) 出願日	平成20年10月15日(2008.10.15)		インターナショナル・ビジネス・マシーンズ・コーポレーション
(65) 公表番号	特表2011-501254 (P2011-501254A)		INTERNATIONAL BUSINESS MACHINES CORPORATION
(43) 公表日	平成23年1月6日(2011.1.6)		アメリカ合衆国10504 ニューヨーク州 アーモンク ニュー オーチャードロード
(86) 国際出願番号	PCT/EP2008/063850	(74) 代理人	100108501
(87) 国際公開番号	W02009/050187		弁理士 上野 剛史
(87) 国際公開日	平成21年4月23日(2009.4.23)	(74) 代理人	100112690
審査請求日	平成23年8月4日(2011.8.4)		弁理士 太佐 種一
(31) 優先権主張番号	11/872, 235	(74) 代理人	100091568
(32) 優先日	平成19年10月15日(2007.10.15)		弁理士 市位 嘉宏
(33) 優先権主張国	米国 (US)		
早期審査対象出願			

最終頁に続く

(54) 【発明の名称】 サーバ・クラスタにおいてフェイルオーバを管理するための方法、フェイルオーバ・サーバ及びコンピュータ・プログラム

(57) 【特許請求の範囲】

【請求項1】

サーバ・クラスタにおいてフェイルオーバを管理するための方法であって、
前記サーバ・クラスタのうちの一つのサーバが分散型ネットワーク内の前記サーバ・クラスタにおける障害のあるサーバ(以下、障害サーバという)を検出したことに応答して、前記障害サーバを検出した前記サーバ(以下、フェイルオーバ・サーバという)が、当該フェイルオーバ・サーバのサブスクリプション・キューについてのサブスクリプション・メッセージ処理を停止するステップと、
前記フェイルオーバ・サーバが、前記障害サーバのサブスクリプション・キューを開くステップと、
前記フェイルオーバ・サーバが、前記障害サーバの前記サブスクリプション・キューを現在管理しているフェイルオーバ・サーバの識別子を含むマーカ・メッセージを、特定のメッセージング・トピックをサブスクライブしている全てのサブスクライバにパブリッシュするステップであって、前記特定のメッセージング・トピックは、サブスクライバに対して関心のある主題を示す、前記パブリッシュするステップと、
前記フェイルオーバ・サーバが、前記障害サーバのサブスクリプション・キュー内のメッセージを処理するステップと、
前記フェイルオーバ・サーバが、前記障害サーバの前記サブスクリプション・キューにおけるメッセージが前記マーカ・メッセージであるかどうかを判断するステップと、
前記フェイルオーバ・サーバが、前記障害サーバの前記サブスクリプション・キューに

おけるメッセージが前記マーカ・メッセージであることに応答して、前記障害サーバの前記サブスクリプション・キューを閉じるステップと、

前記フェイルオーバー・サーバが、当該フェイルオーバー・サーバの前記サブスクリプション・キューについての前記サブスクリプション・メッセージ処理を再開するステップとを含む、前記方法。

【請求項 2】

前記フェイルオーバー・サーバが、前記障害サーバの前記サブスクリプション・キューにおけるメッセージが前記マーカ・メッセージではないことに応答して、前記メッセージを読み取り済みとして記録するステップと、

前記フェイルオーバー・サーバが、前記障害サーバに対する関連のサブスクリプション・セッションにおける処理を遂行するステップとを更に含む、請求項 1 に記載の方法。

10

【請求項 3】

前記フェイルオーバー・サーバが、当該フェイルオーバー・サーバの前記サブスクリプション・キューにおけるメッセージが前記マーカ・メッセージであるかどうかを判断するステップと、

前記フェイルオーバー・サーバが、当該フェイルオーバー・サーバの前記サブスクリプション・キューにおけるメッセージが前記マーカ・メッセージであることに応答して、正規のオペレーションを再開するステップと

を更に含む、請求項 1 又は 2 に記載の方法。

20

【請求項 4】

前記フェイルオーバー・サーバが、当該フェイルオーバー・サーバの前記サブスクリプション・キューにおけるメッセージがマーカ・メッセージではなく且つ前記メッセージが読み取られてないことに応答して、前記フェイルオーバー・サーバに対する関連のサブスクリプション・セッションにおける処理を遂行するステップを更に含む、請求項 1 ~ 3 のいずれか一項に記載の方法。

【請求項 5】

前記フェイルオーバー・サーバが、サーバの起動に応答して、前記起動したサーバに対する一意的なサブスクリプション識別子を生成するステップと、

前記フェイルオーバー・サーバが、新しいセッションを作成したことに応答して、前記一意的なサブスクリプション識別子を前記新しいセッションに格納するステップとを更に含む、

30

前記一意的なサブスクリプション識別子は、前記新しいセッションのライフタイムの間前記新しいセッションにおいて継続する、請求項 1 ~ 4 のいずれか一項に記載の方法。

【請求項 6】

前記一意的なサブスクリプション識別子は、メッセージング・トピックに接続するために使用される、請求項 5 に記載の方法。

【請求項 7】

前記分散型ネットワークが、セッション・アフィニティを利用する分散型パブリッシュ・サブスクリプ・ネットワークである、請求項 1 ~ 6 のいずれか一項に記載の方法。

40

【請求項 8】

前記マーカ・メッセージが、前記障害サーバの前記サブスクリプション・キュー及び前記フェイルオーバー・サーバの前記サブスクリプション・キューにおいて同時に現れる、請求項 1 ~ 7 のいずれか一項に記載の方法。

【請求項 9】

前記マーカ・メッセージが、前記フェイルオーバー・サーバが前記正規のオペレーションを再開する前に、前記障害サーバの前記サブスクリプション・キュー及び前記フェイルオーバー・サーバの前記サブスクリプション・キューを同期させるために使用される、請求項 3 に記載の方法。

【請求項 10】

50

前記マーカ・メッセージが、前記障害サーバの前記サブスクリプション・キュー及び前記フェイルオーバ・サーバの前記サブスクリプション・キューを同期させるために使用される、請求項 1 ~ 8 のいずれか一項に記載の方法。

【請求項 1 1】

前記マーカ・メッセージが、フェイルオーバに参加していない他のすべてのサーバによって無視される、請求項 1 ~ 1 0 のいずれか一項に記載の方法。

【請求項 1 2】

サーバ・クラスタにおいてフェイルオーバを管理するためのフェイルオーバ・サーバであって、

バス・システムと、

前記バス・システムに接続され、命令のセットを含む記憶装置と、

前記バス・システムに接続された処理ユニットと

を含み、

前記処理ユニットに、請求項 1 ~ 1 1 のいずれか一項に記載の方法の各ステップを実行させる、前記フェイルオーバ・サーバ。

【請求項 1 3】

サーバ・クラスタにおいてフェイルオーバを管理するためのサーバであって、

バス・システムと、

前記バス・システムに接続され、命令のセットを含む記憶装置と、

前記バス・システムに接続された処理ユニットと

を含み、当該サーバが、分散型ネットワーク内の前記サーバ・クラスタにおける障害のあるサーバ（以下、障害サーバという）を検出したことに応答して、

前記処理ユニットに、請求項 1 ~ 1 1 のいずれか一項に記載の方法の各ステップを実行させる、前記サーバ。

【請求項 1 4】

コンピュータ・プログラムであって、フェイルオーバ・サーバに、請求項 1 ~ 1 1 のいずれか一項に記載の方法の各ステップを実行させるためのコンピュータ・プログラム。

【請求項 1 5】

コンピュータ・プログラムであって、分散型ネットワーク内の前記サーバ・クラスタにおける障害のあるサーバ（以下、障害サーバという）を検出したことに応答して、前記障害サーバを検出したサーバ（以下、フェイルオーバ・サーバという）に、請求項 1 ~ 1 1 のいずれか一項に記載の方法の各ステップを実行させるためのコンピュータ・プログラム

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データ処理システムの改良に関する。具体的にいえば、本発明は、セッション・アフィニティを利用する分散型ネットワーク環境において、サーバ・フェイルオーバを処理するためにコンピュータを利用する方法、システム、及びコンピュータ・プログラムに関する。

【背景技術】

【0002】

今日、ほとんどのコンピュータが何らかのタイプのネットワークに接続される。ネットワークは、コンピュータが他のコンピュータ・システムと情報を共有することを可能にする。インターネットは、コンピュータ・ネットワークの一例である。インターネットは、コンピュータのグローバル・ネットワークであり、データ転送を処理し且つ、送信ネットワークのプロトコルから受信ネットワークにより使用されるプロトコルへのメッセージの変換を処理する、ゲートウェイによって結合されたネットワークである。インターネットでは、任意のコンピュータが、プロトコルとも呼ばれる種々の言語を通して情報がインターネット上を移動することによって、他のコンピュータと通信を行なうことを可能にする

10

20

30

40

50

。一般に、インターネットは、伝送制御プロトコル/インターネット・プロトコル(TCP/IP)と呼ばれる一連のプロトコルを使用する。

【0003】

多くの最新のインターネット・アプリケーションは、異なる組織上の境界、異機種のプラットフォーム、並びにパブリッシャ(即ち、発行元)及びサブスクライバ(即ち、引用先)の大きな動的集団の全体にわたって情報を配布することを必要とする。パブリッシュ・サブスクライブ(パブ・サブ)・ネットワーク・サービスは、潜在的に無制限の数のパブリッシャおよびサブスクライバの全体にわたる情報の配布を可能にする通信機能である。パブ・サブ・システムは、多くの場合、ピア・ツー・ピア・オーバレイ・ネットワークよりも上位の通信を行なう空間的にまったく異なるノードの集合体として実現されている。

10

【0004】

そのような環境では、パブリッシャはイベントの形で情報をパブリッシュし、サブスクライバは、パブ・サブ・ネットワークにサブスクリプション・フィルタを送ることにより或るイベント又は或るパターンのイベントに対する関心事を表す能力を有する。パブ・サブ・ネットワークは、すべての活動中のサブスクリプションに各アプリケーションを動的に対抗させるためにコンテンツ・ベースの経路指定方法を使用し、イベントがそれらの登録済みの関心事と一致する場合、しかもその場合にのみ、そのイベントをサブスクライバに通知する。

【0005】

20

集中サービス(converged service)は、高いレベルの機能を提供するために多数のネットワーク・プロトコル及びプロトコル・セッション上の通信にまたがるアプリケーションである。ハイパーテキスト転送プロトコル(HTTP)及びセッション開始プロトコル(SIP)の場合、集中サービスは、HTTP及びSIPの両方からのセッション情報を結合し、1つのプロトコルによる対話がそのプロトコルの制約次第で他のプロトコルによる通信に影響を及ぼすことがある。集中サービスは、これらのプロトコルの各々にまたがって多数のプロトコル・セッションに及ぶこともある。

【0006】

コードの構造化および高い有効性サービスを簡単にするために、セッション・アフィニティと呼ばれるメカニズムが集中サービスと関連して使用される。セッション・アフィニティは、そのセッション内の要求をサーバ・クラスタ内の特定のサーバに関連付けるためのクラスタ構成された環境における機構である。この関連付けは、セッションを管理サーバにマップする経路指定メカニズムを介して達成される。集中サービスと共にセッション・アフィニティを使用するとき、集中セッション・データは、単一のアプリケーション・サーバの場合、セッションのライフタイム内では生きており、集中セッションに関する要求を処理するときアプリケーション・コードがクラスタ間通信を行なう必要のないようにし得る。

30

【0007】

しかし、多くの集中アプリケーションは、多数の集中セッション全体にわたって共通の資源又はデータ構造にアクセスすること及びそれを操作することを必要とする。セッション・アフィニティによっても、集中セッションは、クラスタ内の種々のサーバに割り当てられる。その結果、クラスタにおける関連するセッションの場所に関係なくそれらのセッションに関連する共通の情報の集中セッションを通知するための方法が必要である。例えば、サーバA、B、及びCを含む3つのサーバがクラスタ構成された環境を考察する。サーバA及びサーバCにおけるアプリケーション資源に関する通知のためにサブスクリプションが設定される。パブリッシュ要求が入ってきてサーバBに与えられる。サーバBは、クラスタにおけるどのサーバがその関心のあるサブスクリプション・セッションを含むかを知らない。サーバBはサブスクリプション・データを高い信頼性でブロードキャストできなければならない。更に、サーバBが障害を生じた場合、サーバBによって管理されたセッションがそのクラスタ内のどこで再活性化されるかはわからない。

40

50

【発明の概要】**【発明が解決しようとする課題】****【0008】**

従って、本発明の目的は、セッション・アフィニティを利用するパブ・サブ分散型ネットワーク環境においてサーバ・フェイルオーバを管理するためにコンピュータを利用する改良された方法、システム、及びコンピュータ・プログラムを提供することにある。

【課題を解決するための手段】**【0009】**

本実施例は、サーバ・クラスタにおいてフェイルオーバを管理するためにコンピュータを利用する方法、システム、及びコンピュータ・プログラムを提供する。分散型ネットワーク内のサーバ・クラスタにおける障害のあるサーバ（以下、障害サーバという）を検出したことに応答して、フェイルオーバ・サーバのサブスクリプション・メッセージ処理が停止される。障害サーバのサブスクリプション・キューが開かれる。特定のメッセージング・トピックのすべてのサブスクライバに対して、マーカ・メッセージがパブリッシュされる。マーカ・メッセージは、障害サーバのサブスクリプション・キューを現在管理しているフェイルオーバ・サーバの識別子を含む。障害サーバのサブスクリプション・キュー内のメッセージが処理される。障害サーバのサブスクリプション・キューにおけるメッセージがマーカ・メッセージであるかどうか判断される。障害サーバのサブスクリプション・キューにおけるメッセージがマーカ・メッセージであるという判断に応答して、障害サーバのサブスクリプション・キューが閉じられる。しかる後、フェイルオーバ・サーバが、当該フェイルオーバ・サーバのオリジナルのサブスクリプション・キューの処理を再開する。

【0010】

オリジナルのサブスクリプション・キューからの未検出のメッセージを処理している間、マーカ・メッセージが探索されることが望ましい。一旦そのマーカ・メッセージがオリジナルのサブスクリプション・キューにおいて見つかり、正規のオペレーションが再開されることが望ましい。

【図面の簡単な説明】**【0011】**

【図1】本発明の実施例を具体化し得るデータ処理システムのネットワークの概略図である。

【図2】本発明の実施例を具体化し得るデータ処理システムのブロック図である。

【図3】本発明の実施例に従って通常のサーバ・オペレーションのための例示的プロセスを示すフローチャートである。

【図4】本発明の実施例に従ってフェイルオーバ・サーバ・オペレーションのための例示的プロセスを示すフローチャートである。

【発明を実施するための形態】**【0012】**

次に、図面、特に図1及び図2を参照すると、本実施例を具体化し得るデータ処理環境の例示的概略図が示される。図1及び図2は単に例示的なものであって、種々の実施例を具体化し得る環境に関して如何なる制限も主張及び暗示することを意図するものではないということは明らかであろう。図示の環境に対して多くの修正を施すことが可能である。

【0013】

図1は、本実施例を具体化し得るデータ処理システムのネットワークの概略図を示す。ネットワーク・データ処理システム100は、本実施例を具体化し得るコンピュータのネットワークである。ネットワーク・データ処理システム100はネットワーク102を含み、そのネットワークは、ネットワーク・データ処理システム100内で接続される種々の装置及びコンピュータの間の通信リンクを提供するために使用される媒体である。ネットワーク102は、有線通信リンク、無線通信リンク、又は光ファイバ・ケーブルのような接続体を含み得る。

【 0 0 1 4 】

図示の例では、サーバ104及びサーバ106が記憶ユニット108と共にネットワーク102に接続している。更に、クライアント110、112、及び114もネットワーク102に接続している。しかし、ネットワーク・データ処理システム100が、図示されていない更なるサーバ、クライアント、及び他の装置を含み得るということを留意されたい。クライアント110、112、及び114は、サーバ104及び/又はサーバ106に対するクライアントである。更に、クライアント110、112、及び114は、例えば、パーソナル・コンピュータ又はネットワーク・コンピュータであってもよい。

【 0 0 1 5 】

図示の例では、サーバ104及びサーバ106はクラスタ構成されたサーバである。更に、サーバ104及び106は、サブスクライバであるクライアント110、112、及び114にパブ・サブ・ネットワーク・サービスを提供する。パブ・サブ・ネットワークは、特定のメッセージング・トピックに対するメッセージのプブリッシュをサポートする。トピックは、複数のサブスクライバに対して関心のある主題を示す。一般に、メッセージは、プブリッシュ・プロセス中にトピックに割り当てられ、しかる後、その特定のメッセージング・トピックをサブスクライブしているすべての消費者によって受け取られる。ゼロ又はそれ以上のサブスクライバ・クライアントが、特定のメッセージ・トピックに関するメッセージを受け取ることに関心があることを登録し得る。

【 0 0 1 6 】

サブスクリプションは、それが永続的及び耐久的なものとなるように構成される。サブスクリプションは、或る消費者が或るクラスのイベントを受け取ることに関心があるということを表す。例えば、Javaメッセージング・サービス(JMS)のようなメッセージング・サービス環境内のサブスクリプションは、トピックに対するイベントを、そのイベントがプブリッシュされた順序で受け取るための「仮想キュー」として作用する。永続的であるということは、クライアントがサブスクリプションからメッセージを読むことを停止するとき、そのクライアントが取りやめたサブスクリプション・キュー内に未読のメッセージが残っているということの意味する。

【 0 0 1 7 】

更に、ネットワーク・データ処理システム100は、セッション・アフィニティを利用する分散型ネットワーク環境である。セッション・アフィニティは、同じアプリケーション・サーバ・インスタンスに同じ集中セッションの一部である要求を経路指定するために負荷平衡化要素を利用する。集中セッションに対処するために、本実施例はセッション情報を伴う要求を符号化する。なお、セッション情報は、要求を経路指定するために使用することが可能である。このセッション情報の符号化は2つの方法で行うことが可能である。1つの方法は、クライアント・アプリケーションが、例えば、クッキーのようなセッション基準を集中アプリケーションから得るものである。この場合、クライアント・アプリケーションは、その要求におけるかかるクッキーを再生する。もしくは、クライアント・アプリケーションが、要求が適切なサーバに送られるように、その要求のユニフォーム・リソース・アイデンティファイア(URI)を符号化する。

【 0 0 1 8 】

図示の例では、ネットワーク・データ処理システム100は、相互に通信を行なうためにTCP/IPスイート・プロトコルを使用するネットワーク及びゲートウェイの世界的な集合体を表すネットワーク102を備えたインターネットである。そのインターネットの中心には、数千個の商業的コンピュータ・システム、政府機関のコンピュータ・システム、教育機関のコンピュータ・システム、及び、データ及びメッセージを経路指定する他のコンピュータ・システムから成る、主要ノード又はホスト・コンピュータ間的高速データ通信回線のバックボーンがある。勿論、ネットワーク・データ処理システム100は、例えば、イントラネット、ローカル・エリア・ネットワーク(LAN)、広域ネットワーク(WAN)のような数多くの様々なタイプのネットワークとして具現化することも可能である。図1は種々の実施例の一例として意図され、それらの実施例に対するアーキテク

10

20

30

40

50

チャ上の限定として意図されるものではない。

【0019】

次に、図2を参照すると、本実施例を具現化し得るデータ処理システムのブロック図が示される。データ処理システム200は、図1におけるサーバ104又はクライアント110のようなコンピュータの一例であり、そこには、本実施例のためのプロセスを具現化するコンピュータ使用可能プログラム・コード又は命令が設けられる。

【0020】

図示の例では、データ処理システム200は、インターフェース及びメモリ・コントローラ・ハブ(インターフェース/MCH)202並びにインターフェース及び入出力(I/O)コントローラ・ハブ(インターフェース/ICH)204を含むハブ・アーキテクチャを使用する。処理ユニット206、メイン・メモリ208、及びグラフィックス・プロセッサ210はインターフェース/MCH202に接続される。処理ユニット206は1つ又は複数のプロセッサを含み得るし、1つ又は複数の異種のプロセッサ・システムを使用して具現化されてもよい。グラフィックス・プロセッサ210は、例えば、アクセラレイテッド・グラフィック・ポート(AGP)を介してインターフェース/MCH202に接続されてもよい。

10

【0021】

図示の例では、ローカル・エリア・ネットワーク(LAN)アダプタ212がインターフェース/ICH204に接続され、オーディオ・アダプタ216、キーボード及びマウス・アダプタ220、モデム222、リード・オンリ・メモリ(ROM)224、ユニバーサル・シリアル・バス(USB)及び他のポート232、並びにPCI/PCIe装置234がバス238を介してインターフェース/ICH204に接続され、ハードディスクドライブ(HDD)226及びCD-ROM230がバス240を介してインターフェース/ICH204に接続される。PCI/PCIe装置は、例えば、イーサネット・アダプタ、アドイン・カード、及びノート型コンピュータ用PCカードを含み得る。PCIはカード・バス・コントローラを使用し、一方、PCIeはそれを使用しない。ROM224は、例えば、フラッシュ・バイナリ入出力システム(BIOS)であってもよく、HDD226及びCD-ROM230は、例えば、統合ドライブ・エレクトロニクス(IDE)又はシリアル・アドバンスト・テクノロジー・アタッチメント(SATA)インターフェースを使用し得る。スーパーI/O(SIO)装置236がインターフェース/ICH204に接続されてもよい。

20

30

【0022】

オペレーティング・システムが処理ユニット206において作動し、図2におけるデータ処理システム200内の種々のコンポーネントの制御を調整及び提供する。オペレーティング・システムは、Microsoft Windows Vistaのような市販のオペレーティングであってもよい。なお、Microsoft及びWindows Vistaは、米国におけるマイクロソフト社の登録商標である。Javaプログラミング・システムのようなオブジェクト指向プログラミング・システムが、そのオペレーティング・システムに関連して作動し得るし、データ処理システム200において実行されるJavaプログラム又はアプリケーションからそのオペレーティング・システムにコールを行ない得る。Java及びすべてのJavaベースの商標は、米国におけるサン・マイクロシステムズ社の登録商標である。

40

【0023】

オペレーティング・システム、オブジェクト指向プログラミング・システム、及びアプリケーション又はプログラムのための命令は、HDD226のような記憶装置に置かれ、処理ユニット206による実行のためにメイン・メモリ208にロードされる。本実施例のプロセスは、例えば、メイン・メモリ208、ROM224のようなメモリ或いは1つ又は複数の周辺装置に置かれるコンピュータ実装の命令を使用して処理ユニット206により遂行される。

【0024】

図1及び図2におけるハードウェアは具現化態様に従って変化し得る。フラッシュ・メ

50

メモリ、同等の不揮発性メモリ、又は光ディスク・ドライブ等のような他の内部ハードウェア又は周辺装置が、図1及び図2に示されたハードウェアに加えて、又はそのハードウェアの代わりに使用されてもよい。更に、本実施例のプロセスは、マルチプロセッサ・データ処理システムに適用されてもよい。

【0025】

或る実施例では、データ処理システム200は、オペレーティング・システム・ファイル及び/又はユーザ生成のデータを格納するための不揮発性メモリを提供するために一般にはフラッシュ・メモリと共に構成される携帯情報端末(PDA)であってもよい。バス・システムは、システム・バス、I/Oバス、及びPCIバスのような1つ又は複数のバスで構成されてもよい。勿論、任意のタイプの通信ファブリック又はアーキテクチャを使用して、その通信ファブリック又はアーキテクチャに接続された種々のコンポーネント又は装置の間でデータの転送を行なうバス・システムが実装されてもよい。通信ユニットは、モデム又はネットワーク・アダプタのような、データを送受信するために使用される1つ又は複数の装置を含み得る。メモリは、例えば、メイン・メモリ208又は、インターフェイス/MCH202において見られるようなキャッシュであってもよい。処理ユニットは、1つ又は複数のプロセッサ又はCPUを含み得る。図1及び図2における図示の例並びに上記の例は、アーキテクチャ上の限定を暗示することを意味しない。例えば、データ処理システム200は、PDA形式をとることの他に、タブレット・コンピュータ、ラップトップ・コンピュータ、又は電話装置であってもよい。

【0026】

本実施例は、分散型ネットワーク内のサーバ・クラスタにおいてフェイルオーバを管理するためにコンピュータを利用する方法、システム、及びコンピュータ使用可能なプログラム・コードを提供する。サーバ・クラスタにおける障害サーバを検出したことに応答して、フェイルオーバ・サーバが、自己のサブスクリプション・メッセージング・キューのサブスクリプション・メッセージ処理を停止する。しかる後、フェイルオーバ・サーバは、障害サーバのサブスクリプション・キューを開き、特定のメッセージング・トピックの全てのサブスクライバにマーカ・メッセージをパブリッシュする。マーカ・メッセージは、その障害サーバのサブスクリプション・キューをその時点で管理しているフェイルオーバ・サーバの識別子を含む。更に、マーカ・メッセージは、フェイルオーバに参加していない他のすべてのサーバによって無視される。

【0027】

更に、フェイルオーバ・サーバは、障害サーバのサブスクリプション・キュー内のメッセージを処理する。障害サーバのサブスクリプション・キューにおけるメッセージを処理している間、フェイルオーバ・サーバは、サブスクリプション・キューにおけるメッセージがマーカ・メッセージであるかどうかを判断する。障害サーバのサブスクリプション・キューにおけるマーカ・メッセージが見つかったことに応答して、フェイルオーバ・サーバは、障害サーバのサブスクリプション・キューを閉じ、当該フェイルオーバ・サーバのサブスクリプション・メッセージ処理を再開して自己のサブスクリプション・キューを処理する。

【0028】

本実施例は、分散型パブ・サブ・ネットワーク環境において具現化することが可能である。開始時における各アプリケーション・サーバは、セッションのライフタイムにおいて存続する一意的なサブスクリプション識別子を発生する。各サブスクリプションは耐久性及び永続性がある。セッションがアプリケーション・サーバにおいて作成されるときには常に、アプリケーション・サーバは、セッション属性がセッション状態複製の一部として複製されるように、一意的なサブスクリプション識別子をセッション属性に格納する。フェイルオーバが生じるとき、障害サーバのセッションが別のサーバアプリケーション上で活性化される。

【0029】

一般に、例えば、J2EEアプリケーション・サーバのようなアプリケーション・サー

10

20

30

40

50

バ環境では、フェイルオーバ・サーバ・アプリケーション・コードは、ライフサイクル・リスナを介してこの活動を信号で知らされる。しかし、ライフサイクル・リスナを介するほかに、プラットフォームによる多くの様々な方法でこの活動をアプリケーション・コードに知らせることが可能であるということに留意されたい。活動時、フェイルオーバ・サーバ・アプリケーションは、フェイルオーバしたサブスクリプションをセッションから見つける。しかる後、フェイルオーバ・サーバは、先入れ/先出し(FIFO)順に配布され且つ障害サーバのサブスクリプション・キューにおいても現れるマーカ・メッセージを直ちにすべてのサブスクリプションにパブリッシュする。

【0030】

しかる後、フェイルオーバ・サーバは、自己の一次サブスクリプションに関する処理を停止し、フェイルオーバしたサブスクリプションから回復するように進み、障害サーバのセッションに関連するメッセージだけを処理する。回復している間、フェイルオーバ・サーバは、それが処理したメッセージのマップを作成する。一旦マーカ・メッセージがその障害サーバのサブスクリプション・キューにおいてヒットされると、回復の第1段階が終了する。

10

【0031】

その後、フェイルオーバ・サーバ・アプリケーションは、それがマーカ・メッセージを見つけるまで一次キューの処理を再開する。マーカ・メッセージは、キューが現在同期していることを示す。一次キューを探索している間、フェイルオーバ・サーバはマップにあるメッセージをスキップするので、それは、回復されたセッションのための如何なるメッセージも処理しない。一旦フェイルオーバ・サーバがマーカ・メッセージを見つけると、フェイルオーバ・サーバはフェイルオーバしたサブスクリプションを終了し、正規の処理を再開する。しかし、フェイルオーバ・サーバが回復しているとき、障害サーバが動き出す場合、障害サーバが新しい一意的なサブスクリプション識別子を生成し、それによって、フェイルオーバ・サーバとの不一致を回避するので、フェイルオーバ・サーバは影響を受けない、ということに留意されたい。

20

【0032】

次に図3を参照すると、本実施例に従って、通常のサーバ・オペレーションに関する典型的なプロセスを表すフローチャートが示される。図3に示されるプロセスは、例えば、図1におけるサーバ104のようなサーバにおいて具現化することが可能である。

30

【0033】

プロセスは、サーバが起動するときに始まる(ステップ302)。起動後、サーバは、そのセッションのライフタイムの間生きている一意的なサブスクリプション識別子を生成する(ステップ304)。サーバは、その一意的なサブスクリプション識別子を使って、サブスクリプション・データに対する変更の通知をサブスクライバ・クライアントにパブリッシュするために使用されるメッセージング・トピックに接続する。

【0034】

ステップ304において一意的なサブスクリプション識別子を生成することに続いて、サーバは、その一意的なサブスクリプション識別子を使ってメッセージング・トピックに接続する(ステップ306)。ステップ306においてメッセージング・トピックに接続した後、サーバは、これが新しいセッションであるかどうかに関して決定を行なう(ステップ308)。これが新しいセッションである場合(ステップ308の「イエス」出力)、サーバは、一意的なサブスクリプション識別子を属性としてセッション・オブジェクトに格納する(ステップ310)。しかる後、プロセスはステップ312に進む。

40

【0035】

これが新しいセッションではない場合(ステップ308の「ノー」出力)、サーバは要求をサービスする(ステップ312)。サーバは要求処理を行なう(ステップ314)。更に、サーバは、特定のメッセージング・トピックに変更をパブリッシュする(ステップ316)。更に、サーバは、セッション状態複製のためにセッション・オブジェクトを複製する(ステップ318)。しかる後、プロセスはステップ312に戻り、サーバは要求

50

のサービスを継続する。

【0036】

次に、図4を参照すると、本実施例に従って、フェイルオーバ・サーバ・オペレーションのための例示的なプロセスを表すフローチャートが示される。図4に示されるプロセスは、例えば、図1におけるサーバ106のようなサーバにおいて具現化することが可能である。

【0037】

プロセスは、サーバが、例えば図1におけるサーバ104のようなサーバ・クラスタ内の別のサーバの障害を検出するときに始まる(ステップ402)。ステップ402において他のサーバの障害を検出した後、別のサーバの障害を検出したサーバ(フェイルオーバ・サーバ)はそれ自身のサブスクリプション・メッセージの処理を停止する(ステップ404)。しかる後、フェイルオーバ・サーバは、障害サーバのサブスクリプション・キューを開く(ステップ406)。更に、フェイルオーバ・サーバは、現在その障害サーバのセッションにサービスを提供しているサーバに関する一意的なサブスクリプション識別子を含むマーカ・メッセージを、すべてのサブスクリバにパブリッシュする(ステップ408)。更に、マーカ・メッセージは、フェイルオーバが生じたとき、フェイルオーバ・サーバのサブスクリプション・キュー及び障害サーバのサブスクリプション・キューの両方においてFIFO順に現れる。フェイルオーバ・サーバは、このマーカ・メッセージを使用して、後述するように両方のサーバのサブスクリプション・キューを同期させる。

【0038】

しかる後、フェイルオーバ・サーバは障害サーバのサブスクリプション・キューからメッセージを取得する(ステップ410)。しかる後、フェイルオーバ・サーバは、そのメッセージがマーカ・メッセージであるかどうかに関して判断を行なう(ステップ412)。メッセージがマーカ・メッセージではない場合(ステップ412からの「ノー」出力)、フェイルオーバ・サーバはそのメッセージを「読み取り済み」として記録する(ステップ414)。ステップ414で見られるようにメッセージを記録した後、フェイルオーバ・サーバは関連のセッションに関する処理を行なう(ステップ416)。しかる後、プロセスはステップ410に戻り、フェイルオーバ・サーバは、障害サーバのサブスクリプション・キューから別のメッセージを得る。

【0039】

ステップ412に戻ると、メッセージがマーカ・メッセージである場合(ステップ412の「イエス」出力)、フェイルオーバ・サーバは、障害サーバのサブスクリプション・キューを閉じる(ステップ418)。しかる後、フェイルオーバ・サーバは、それ自身のサブスクリプション・キューの処理を再開する(ステップ420)。そのサブスクリプション・キューの処理を再開した後、フェイルオーバ・サーバは、そのサブスクリプション・キューからメッセージを得る(ステップ422)。しかる後、フェイルオーバ・サーバは、そのメッセージがマーカ・メッセージであるかどうかに関して判断を行なう(ステップ424)。メッセージがマーカ・メッセージである場合(ステップ424の「イエス」出力)、フェイルオーバ・サーバは正規の処理を再開する(ステップ426)。しかる後、プロセスは終了する。

【0040】

メッセージがマーカ・メッセージではない場合(ステップ424の「ノー」出力)、フェイルオーバ・サーバは、メッセージが以前に読み取られたかどうかに関して判断を行なう(ステップ428)。メッセージが読み取られた場合(ステップ428の「イエス」出力)、プロセスはステップ422に戻り、フェイルオーバ・サーバはそれ自身のサブスクリプション・キューから別のメッセージを得る。メッセージが読み取られなかった場合(ステップ428の「ノー」出力)、フェイルオーバ・サーバは、関連のセッションに関する処理を行なう(ステップ430)。しかる後、そのプロセスは、再びステップ422に戻る。

【0041】

従って、本実施例は、セッション・アフィニティを利用するパブリッシュ・サブスクリプション分散型ネットワーク環境においてサーバ・フェイルオーバを処理するためにコンピュータを利用する方法、システム、及びコンピュータ使用可能プログラム・コードを提供する。本発明は、全体的にハードウェアの実施例、全体的にソフトウェアの実施例、又は、ハードウェア要素及びソフトウェア要素の両方を含む実施例の形式を取り得る。好ましい実施例では、本発明は、ファームウェア、駐在ソフトウェア、マイクロコード等を含むがそれに限定されないソフトウェアにおいて具現化される。

【0042】

更に、本発明は、コンピュータ又は任意の命令実行システムによる使用、或いは、それに関連した使用のためのプログラム・コードを提供するコンピュータ使用可能媒体又はコンピュータ可読媒体からアクセスし得るコンピュータ・プログラムの形式を取り得る。説明の便宜上、コンピュータ使用可能媒体又はコンピュータ可読媒体は、命令実行システム、装置、又はデバイスによる使用、或いは、それに関連した使用のためのプログラムを含み、格納し、通信し、又は搬送し得る任意の有形の装置であってもよい。

10

【0043】

その媒体は、電子的、磁氣的、光学的、電磁氣的、赤外線、又は半導体システム（又は装置又はデバイス）であってもよい。コンピュータ可読媒体の例は、半導体又はソリッド・ステート・メモリ、磁気テープ、取り外し可能なコンピュータ・ディスク、ランダム・アクセス・メモリ（RAM）、リード・オンリ・メモリ（ROM）、固定磁気ディスク、及び光ディスクを含む。光ディスクの現在の例は、コンパクト・ディスク・リード・オンリ・メモリ（CD-ROM）、コンパクト・ディスク・リード/ライト（CD-R/W）、又はDVDを含む。

20

【0044】

更に、コンピュータ記憶媒体はコンピュータ可読プログラム・コードを含み又はそれを格納し得るので、そのコンピュータ可読プログラム・コードがコンピュータ上で実行されるとき、このコンピュータ可読プログラム・コードの実行は、そのコンピュータに、通信リンクを介して他のコンピュータ可読プログラム・コードを伝送させ得る。この通信リンクは、例えば、制限なく、物理的な媒体又は無線媒体を使用し得る。

【0045】

プログラム・コードを格納及び/又は実行するに適したデータ処理システムは、システム・バスを介してメモリ素子に直接的又は間接的に接続された少なくとも1つのプロセッサを含むであろう。メモリ素子は、プログラム・コードの実際の実行中に使用されるローカル・メモリ、大容量記憶装置、及び、プログラム・コードの実行中にそれらのコードが大容量記憶装置から検索されるべき回数を減らすために少なくとも幾つかのプログラム・コードの一時的記憶装置を提供するキャッシュ・メモリを含み得る。

30

【0046】

入出力装置、即ちI/O装置（キーボード、ディスプレイ、ポインティング・デバイスを含むがそれらに限定されない）は、システムに、直接に又は介在するI/Oコントローラを介して接続される。

【0047】

データ処理システムが、介在する専用ネットワーク又は公衆ネットワークを介して他のデータ処理システム或いは遠隔のプリンタ又は記憶装置に接続されることを可能にするために、ネットワーク・アダプタがシステムに接続されてもよい。モデム、ケーブル・モデム、及びイーサネット・カードは数少ない現在利用可能なタイプのネットワーク・アダプタである。

40

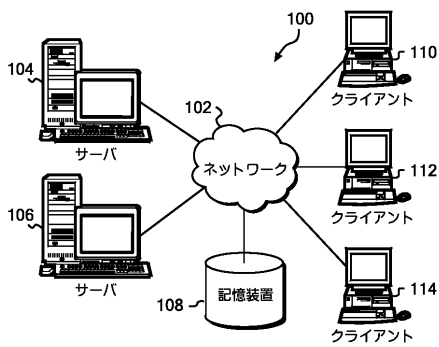
【0048】

本発明に関する記述は説明を目的として示され、網羅的であること又は開示された形式の発明に限定されることを意図するものではない。当業者には多くの修正および変更が明らかであろう。実施例は、本発明の原理及び実用的応用例を最もよく説明するために、また、当業者が、意図した特定の用途に適するように種々の修正を備えた種々の実施例のた

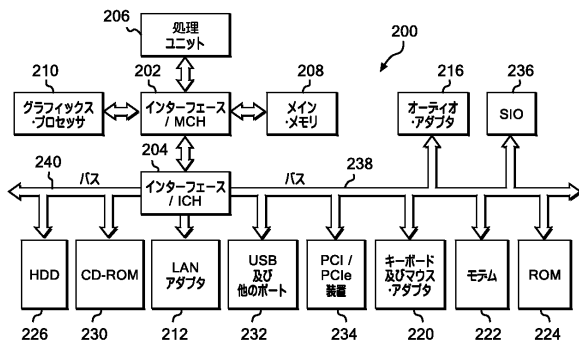
50

めに本発明を理解することを可能にするために、選択及び説明されたものである。

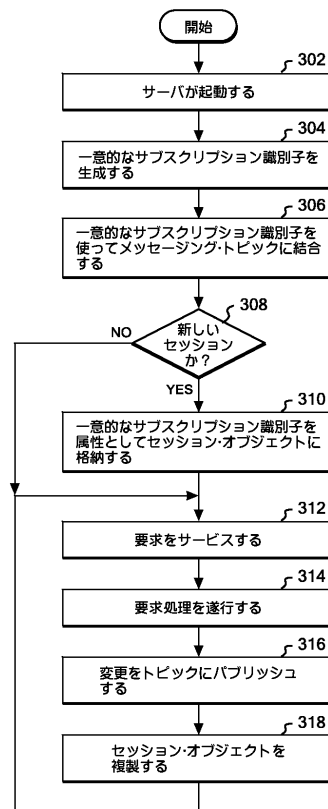
【図1】



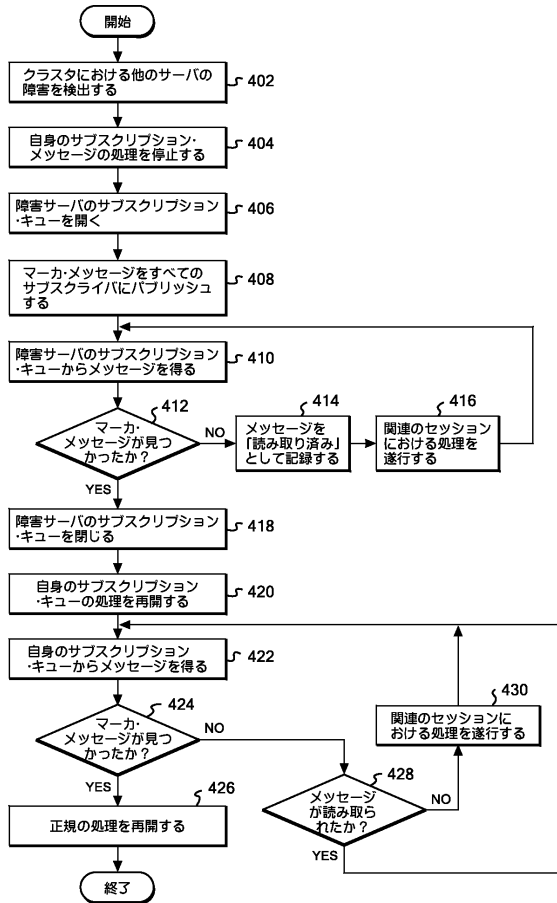
【図2】



【図3】



【図4】



フロントページの続き

- (72)発明者 ギルフィックス、マイケル
アメリカ合衆国78731、テキサス州、オースチン、キャニオンサイド・トレール 4301
- (72)発明者 ムーア、ビクター
アメリカ合衆国32024、フロリダ州、レイク・シティ、サウスウエスト・ダイアル・アベニュー 776
- (72)発明者 ロウベル ジュニア、アンソニー、ウィリアム
アメリカ合衆国27615、ノース・カロライナ州、ローリー、レスリー・ドライブ 10612
- (72)発明者 チェン、ベンソン、クアン、イー
アメリカ合衆国27713、ノース・カロライナ州、ダーハム、チャンセラーズ・リッジ・ドライブ 1014
- (72)発明者 ギルモア、マーク、デイビッド
アメリカ合衆国27713、ノース・カロライナ州、ダーハム、シルバーウッド・コート 1
- (72)発明者 タルアビブ、オフィラ
イスラエル国60946、ピザロン、モシャブ、ピー・オー・ボックス 109

審査官 高橋正徳

- (56)参考文献 特開2007-179310(JP,A)
特開2005-301442(JP,A)
国際公開第2005/091134(WO,A1)
特開2005-209191(JP,A)
特表2005-505833(JP,A)

- (58)調査した分野(Int.Cl., DB名)
G06F 11/16-11/20,
G06F 9/50