

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5551749号
(P5551749)

(45) 発行日 平成26年7月16日(2014.7.16)

(24) 登録日 平成26年5月30日(2014.5.30)

(51) Int.Cl. F 1
G06F 11/20 (2006.01) G06F 11/20 310C

請求項の数 14 (全 32 頁)

<p>(21) 出願番号 特願2012-219878 (P2012-219878) (22) 出願日 平成24年10月1日 (2012.10.1) (62) 分割の表示 特願2009-522059 (P2009-522059) の分割 原出願日 平成19年2月19日 (2007.2.19) (65) 公開番号 特開2013-33493 (P2013-33493A) (43) 公開日 平成25年2月14日 (2013.2.14) 審査請求日 平成24年10月31日 (2012.10.31) (31) 優先権主張番号 11/498,802 (32) 優先日 平成18年8月4日 (2006.8.4) (33) 優先権主張国 米国 (US)</p>	<p>(73) 特許権者 509034030 ティーエスエックス インコーポレイテッド TSX INC. カナダ, オンタリオ州 エム5エックス 1ジェイ2, トロント, ザ エクスチェン ジ タワー, キング ストリート ウェス ト 130 130 King Street Wes t, The Exchange Towe r, Toronto, Ontario M 5X 1J2, Canada (74) 代理人 100116872 弁理士 藤田 和子</p> <p style="text-align: right;">最終頁に続く</p>
--	---

(54) 【発明の名称】 フェイルオーバーシステムおよび方法

(57) 【特許請求の範囲】

【請求項 1】

一次サーバ(62-1)および少なくとも1つのバックアップサーバ(62-2)を備えており、

上記一次サーバ(62-1)が、複数の処理対象となる入力を確定するように、かつ、上記複数の処理対象となる入力を処理する上記一次サーバ(62-1)の処理に先立って、上記バックアップサーバ(62-2)に上記複数の処理対象となる入力を送信するように構成され、

上記複数の処理対象となる入力は、上記サーバ(62-1、62-2)の両方による上記入力の処理の決定性を保証する入力を含む、フェイルオーバーのためのシステム(50)。

【請求項 2】

上記複数の処理対象となる入力は、外部リソースに対する呼び出し結果を含む、請求項1に記載のシステム(50)。

【請求項 3】

上記複数の処理対象となる入力は、タイムスタンプの取得を要求して得た結果であるタイムスタンプ値を含む、請求項1または請求項2に記載のシステム(50)。

【請求項 4】

上記複数の処理対象となる入力は、共有リソース(122)をさらに含む、請求項1から3の何れか1項に記載のシステム(50)。

【請求項 5】

上記共有リソース（122）が、個々のサーバ（62）におけるランダムアクセスメモリ内に維持されている、請求項4に記載のシステム（50）。

【請求項 6】

上記複数の処理対象となる入力、上記一次サーバ（62-1）により生成されたシーケンスナンバをさらに含む、請求項1から5の何れか1項に記載のシステム（50）。

【請求項 7】

上記バックアップサーバ（62-2）は、複製エージェント（126-2）を含み、上記複製エージェント（126-2）は、上記一次サーバ（62-1）から上記バックアップサーバ（62-2）への情報のミラーリングを助けるために、上記一次サーバ（62-1）のライブラリ（102-1）と通信するように構成されている、請求項1から6の何れか1項に記載のシステム（50）。

10

【請求項 8】

上記バックアップサーバ（62-2）は、クライアント（54）からの複製メッセージを認識するように構成されている、請求項1から7の何れか1項に記載のシステム（50）。

【請求項 9】

上記バックアップサーバ（62-2）は、上記クライアント（54）からの上記複製メッセージが受信されたとき、リカバリプロトコルを実行するように構成されている、請求項8に記載のシステム（50）。

20

【請求項 10】

上記リカバリプロトコルは、複製された情報に基づいて、上記バックアップサーバ（62-2）からの返答の送信を行うことを含む、請求項9に記載のシステム（50）。

【請求項 11】

上記リカバリプロトコルは、ギャブリカバリであって、当該ギャブリカバリは、バックアップサーバ（62-2）を用いて一次単独状態にフェイルオーバーする間、上記複数の処理対象となる入力をギャブリカバリする、請求項9または10に記載のシステム（50）。

【請求項 12】

上記一次サーバ（62-1）および上記バックアップサーバ（62-2）は、同時に動作する、請求項1から11の何れか1項に記載のシステム（50）。

30

【請求項 13】

システム（50）におけるフェイルオーバーのための方法であって、
 一次サーバ（62-1）において、複数の処理対象となる入力を確定するステップと、
 上記一次サーバ（62-1）からバックアップサーバ（62-2）に上記複数の処理対象となる入力を送信するステップと、
 上記一次サーバ（62-1）において、上記複数の処理対象となる入力を送信した後に、上記複数の処理対象となる入力を処理するステップと、を含み、
 上記複数の処理対象となる入力は、上記サーバ（62-1、62-2）の両方による上記入力の処理の決定性を保証する入力を含む、フェイルオーバーのための方法。

40

【請求項 14】

ネットワーク（58）を介して相互接続された少なくとも2つのサーバ（62）のうちのいずれか1つのサーバ上で実行可能な、一組のプログラミング命令を保存するコンピュータ読み取り可能な記録媒体であって、

上記一組のプログラミング命令は、

一次サーバ（62-1）において、複数の処理対象となる入力を確定するための命令と、

上記一次サーバ（62-1）からバックアップサーバ（62-2）に上記複数の処理対象となる入力を送信するための命令と、

上記一次サーバ（62-1）において、上記複数の処理対象となる入力を送信した後に

50

、上記複数の処理対象となる入力を処理するための命令と、を含み、

上記複数の処理対象となる入力は、上記サーバ(62-1、62-2)の両方による上記入力の処理の決定性を保証する入力を含む、コンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、コンピュータおよびネットワークアーキテクチャに関する。より具体的には、フェイルオーバーシステムおよび方法に関する。

【背景技術】

【0002】

相互交流やビジネスを行う上で、社会はコンピュータおよびネットワークにますます依存している。重要なシステムにおいて求められる高レベルの可用性を確保するために、ソフトウェアおよびハードウェアの欠陥によって起こる不測のダウンタイムを最小限に抑える必要がある。

【0003】

金融サービス産業は、高可用性を有するシステムを必要とする産業の一例である。実際、今日の金融産業における多くのデータ処理活動がコンピュータシステムに支えられている。特に興味深いのは、いわゆるリアルタイムおよびニアリアルタイムオンライントランザクション処理(OLTP)のアプリケーションであり、それらは通常、多数のビジネス取引を長期間にわたって高速かつ低遅延で行う。これらのアプリケーションは、概して下記の特徴を示す：(1)高度かつ高速のデータ処理(2)信頼性の高い不揮発性データ記憶(3)高レベルの可用性、すなわち、実質的に継続可能な基本原理に基づいてサービスを支える能力。しかし、実行する場合、下記に概要を詳述するように、システム動作に与える相反する影響のため、3つの特徴全てを同時に完全に満たすことができる設計は存在しない。

【0004】

第一に、高度なデータ処理とは、タイムリーな方法で、多数のコンピュータ処理、データベース検索/アップデートなどを実行する能力のことを指す。並列処理を介して上記処理が実施され、同様の物理的機械または分散型ネットワーク上で同時に作業の統括ユニット(multiple unit of work)が実行される。あるシステムにおいて、各トランザクションの結果は既に完了したトランザクションの結果によって決まる。そのようなシステム並行の形態は、本質的に非決定性を有する。すなわち、乱調状態、オペレーティングシステムのスケジュールタスク、または可変ネットワーク遅延(variable network delay)が原因で、メッセージのシーケンスおよびスレッド実行(thread execution)が予測できず、単に複製システムに入力メッセージのコピーを送信するだけでは、それらを並列処理することができない。非決定性のシステムは非同一性のアウトプットを有するので、障害が起きた場合に代替物と代えようとしても、2つの異なるコンピュータ機器上で並列にそれらを動作させることができない。

【0005】

第二に、信頼性の高い不揮発性データ記憶とは、多数のシステムのソフトウェアまたはハードウェア機器が不測の障害に晒された場合であっても、処理済みのデータを持続して保存できる能力のことを指す。共有データにアクセスするか、または変更するとき、アトミック性(Atomic)、一貫性(Consistent)、独立性(Isolated)、および永続性(Durable)("ACID")トランザクションを使用することによって、通常、上記処理を実施することができる。作業単位が完了すると、ACIDトランザクションがデータ保全とデータ持続性とを直ちに確保することができる。コミットメントされた全ACIDを非揮発性コンピュータメモリ(ハードディスク)に書き込み、それによってデータ持続性が確保される。しかし、パフォーマンスの点と一般に全システムを減速させる点とを鑑みると、上記は非常にコストがかかる。

【0006】

10

20

30

40

50

第三に、高可用性のシステムとは、既定のコンピュータシステムにおける可用性率を時間あたり100%に可能な限り近づけることを確実にすることを目的とする。冗長ソフトウェアおよび/またはハードウェアを介してそのような可用性が実施可能であり、それらのソフトウェアまたはハードウェアが、機器障害が検出された場合に機能性を引き継ぐ。引き継ぐために、ファイルオーバがデータだけでなく処理状態をも複製する。当業者であればわかるように、非決定性システム（すなわち、同じ一連のイベントのコンピュータ処理が、それらのイベントの処理順序に従って1つ以上の結果を有し得るシステム）において、状態の複製は特に困難な問題をはらんでいる。

【0007】

高可用性のソフトウェアアプリケーションは冗長環境上に配備されることが多く、基盤となるハードウェアと共通する障害の一点を低減、および/または、除去する。2つの一般的なアプローチとしては、ホットフェイルオーバおよびウォームフェイルオーバが知られている。ホットフェイルオーバとは、複合システムにおいて同じ入力を同時に処理し、基本的に、これらシステムのうちの一つにおける障害のイベントにおいて完全な冗長化を提供することを指す。ウォームフェイルオーバとは、バックアップシステムにおけるアプリケーション（すなわち、データ）の状態を複製することを指す。このとき、バックアップシステムにおいて当該データの複製処理を行うことはないが、一次システムにおける障害のイベントにおいて、読み込み済みでスタンバイ状態のデータを処理することができるアプリケーションを有する。コールドフェイルオーバとは、単にバックアップシステムを立ち上げ、一次システムからの処理負担を想定してバックアップシステムを準備すること

10

20

【0008】

ホットフェイルオーバ設定におけるアプリケーションの二つの例としては、同時に、2種類のハードウェア機器において実行するものと、同じ入力のコピーを処理するものがある。そのうちの一つが重大な障害を起こした場合、補助同期システムは、もう一方が継続して作業負荷をサポートするように確保することができる。ウォームフェイルオーバ設定において、システムのうちの一つ、すなわち、一次に指定されたシステムが、アプリケーションを起動させる。障害の場合は、スタンバイ状態で待機している二次システム、すなわちバックアップに指定されたシステムが起動し、引き継いで上記機能性を再開する。

【0009】

ホットフェイルオーバアプローチに関する従来技術は、少なくとも2つの欠点を有する。第一に、2つのシステムの同期を維持するために、補助ソフトウェアを起動する必要がある。非決定性システムの場合では、この同期化のための注力によって、イベントの発生順序の同一性を保証するというパフォーマンスおよび高度性において、許容範囲外の（または、不要な）低下を招き得る。また、そのようなアプリケーションにおいて使用される従来技術に関する並行システムによって、通常は、多重スレッドが同時に実行できるようになるので、当該スレッドは本来的に非決定性を有する。また、非決定性は、サーバおよび地理的に離れたクライアントを含むシステムであり、当該システムにおいて、可変ネットワーク遅延によって、不測のシーケンスでサーバにメッセージが分配される。

30

【0010】

ホットフェイルオーバに関する問題を克服するのにウォームフェイルオーバを使用することができる。ウォームフェイルオーバは、冗長性のバックアップシステムにシステムデータを複製し、二次システムにアプリケーション機能性を保存することによって、非決定性システムのフェイルオーバを実行する他の方法であり得る。このアプローチは、まず安定状態にデータを修復し、それからアプリケーションを機能状態にし、最終的に中断した処理地点にアプリケーションを戻すのに要する時間という点で欠点を有する。この処理には通常数時間かかり、手動介入が必要であり、一般的にはインフライトトランザクション（in-flight transaction）を復旧させることができない。

40

【0011】

多数の特許が、上記問題のうちの少なくとも一部に対処しようとしている。US特許第5

50

、305,200では、要約すると、買い手/売り手と販売人(値付け業者)との間の協議による貿易シナリオにおけるコミュニケーションのための冗長性機構について提案している。障害発生イベントにおいて冗長性機構の作業を確保するために、冗長性が提供される。上記特許は、非決定性環境におけるオンライン・トランザクション・アプリケーションのフェイルオーバーに対処していない。簡単に言うと、US特許第5,305,200は、ネットワーク障害が起こった後で「命令が送信されたか、されていないか」という問いに対する明確な答えを提供することを目的としている。

【0012】

US特許第5,381,545では、データのアップデートをしている間に、(データベース内に)保存されたデータをバックアップするための技術を提案している。US特許第5,987,432は、地域的分布に関する世界規模の金融市場データを収集するための耐障害性市場データ相場表示装置システム(fault-tolerant market data ticker plant system)に対処している。これは決定性環境であって、解決策は、データを消費者に送信する連続した一方向のフローを提供することに焦点を当てている。US特許第6,154,847は、従来の不揮発性記憶装置上のトランザクションログと揮発性記憶装置内のトランザクションリストとを組み合わせることによって、トランザクションのロールバックを行うことについての改善方法を提供する。US特許第6,199,055は、無認証(unsecured)のコミュニケーションリンクを介して、システムとポータブルプロセッサとの間のトランザクション分配を扱うための方法を提案している。US特許第6,199,055は、遠隔装置との完全なトランザクションを確保する認証について扱っており、障害のイベントにおける遠隔装置のリセットについて扱っている。概して、上記は非決定性環境におけるオンライン・トランザクション・アプリケーションのフェイルオーバーに対処していない。

【0013】

US特許第6,202,149は、タスクを自動的に再分配してコンピュータ停止の影響を低減させるための方法および装置を提案している。装置は、1つ以上のコンピュータシステムからなる少なくとも1つの冗長性グループを含んでおり、当該コンピュータシステムはコンピュータパーティションからなる。パーティションは、各コンピュータシステムパーティションにおいて複製されたデータベーススキーマのコピーを含む。冗長性グループはコンピュータシステムおよびコンピュータシステムパーティションの状態をモニタし、モニタされたコンピュータシステムの状態に基づいて当該コンピュータシステムにタスクを割り当てる。US特許第6,202,149の問題は、バックアップシステムが処理トランザクションの負荷を想定する場合のワークフロー復旧方法について教示しておらず、代わりに、非効率のおよび/または低速になり得る、全データベースの複製を目的としている点である。さらに、そのような複製によって、インフライトにおける重要なトランザクション情報が失われる可能性がある。特に、一次システムにおいて障害が発生するか、または、一次とバックアップシステムとを相互接続させるネットワークにおいて障害が発生することによって、一次とバックアップとの間で非一貫性(inconsistent)状態が生じ得る。概して、US特許第6,202,149は、オンライントランザクションなどの処理に必要な特徴、特に、非決定性システムのフェイルオーバーに必要とされる特徴を欠いている。

【0014】

US特許第6,308,287は、機器トランザクションの障害を検知し、障害を除去し(back out)、システム障害後に回復可能になるように障害表示器を確実に保存し、この障害表示器がさらなるトランザクションに対して可用性を有するようにするための方法を提案している。当該特許は非決定性環境におけるトランザクションアプリケーションのフェイルオーバーに対処していない。US特許第6,574,750は、分配および複製されたオブジェクトに関するシステムを提供しており、この場合、オブジェクトは非決定性を有する。当該特許は、複製されたオブジェクトにおいて障害が発生した場合に、一貫性を保証するとともにロールバックを制限するための方法を提案している。この方法には、オブ

10

20

30

40

50

ジェクトが着信中のクライアント要求 (incoming client request) を受信すること、および、要求IDと、オブジェクトのレプリカによって前処理された全要求のログとを比較することについて記載されている。適合がわかると、関連のレスポンスがクライアントに返ってくる。しかし、この方法だけでは、従来技術における様々な問題を解決するには不十分である。

【0015】

他の問題として、US特許第6,574,750は同期起動チェーン (synchronous invocation chain) を想定しており、当該チェーンは高性能オンラインランザクション処理 (OLTP) アプリケーションには不適であることが挙げられる。同期起動の場合、クライアントは応答またはタイムアウトのいずれかを待ってから続行する。続いて、起動したオブジェクトが他のオブジェクトのクライアントとなって同期コールチェーンを伝搬させてもよい。結果は、広域同期オペレーションとなり得るものであり、このとき、発信元クライアントにおいて構成されたロングタイムアウトを処理および必要とするクライアントをブロックする。

10

【発明の概要】

【0016】

本発明の形態は、ネットワークを介して相互接続された少なくとも2つのサーバのうちのいずれか1つを選択可能であって、選択した方と接続する少なくとも1つのクライアントを備える、フェイルオーバーのためのシステムを提供する。通常状態において、サーバのうちの1つがクライアントに接続される場合に当該サーバが一次サーバに指定され、もう一方のサーバがクライアントに接続されない場合に当該もう一方のサーバがバックアップサーバに指定される。少なくとも1つのクライアントが一次サーバにメッセージを送信するように構成されている。サーバは、各サーバにおいて同一の少なくとも1つのサービスを用いてメッセージ処理を行うように構成されている。サービスは、サービスに関するサーバが一次サーバとして動作しているのか、バックアップサーバとして動作しているのか認識しない。サーバは、ライブラリまたは他の固有の有効コードを維持するよう構成されており、当該ライブラリまたは有効コードは、サーバが一次サーバであるかバックアップサーバであるかの表示を含む様々なタスクを実行するように構成されている。各サーバ内部にあるサービスは、その個々のライブラリに対して外部呼び出しを行うものである。一次サーバにおけるライブラリは、外部呼び出しを完了し、一次サーバにおけるサービスに外部呼び出しの結果を戻し、バックアップサーバにおけるサービスに外部呼び出しの結果を転送するように構成されている。二次サーバにおけるライブラリは、二次サーバにおけるサービスに要求されたとき、外部呼び出しを行わず、単に、一次サーバから受信したとおり外部呼び出しの結果を二次サーバにおけるサービスに転送する。

20

30

【0017】

ライブラリは、1つ以上の固有の有効コードとして実施され得る。

【0018】

サーバは、サービスがメッセージ処理の結果を保存することが可能な共有リソースを維持するようにそれぞれ構成されている。性能上の理由から、個々のサーバにおけるランダムアクセスメモリ内に共有リソースを維持することが好ましい。しかし、必ずしもランダムアクセスメモリ内に共有リソースを維持する必要はない。

40

【0019】

外部呼び出しは、(例示の非制限リストとして)タイムスタンプを要求するか、または、同サーバ上に配置された他のサービスを呼び出すか、または、物理的に離れた機械上に配置された別のサービスを呼び出すかであり得る。

【0020】

システムは電子取引システムの一部であり得るので、メッセージはセキュリティの売買注文であり得る。この場合、外部呼び出しは、セキュリティの値についての市場供給相場を要求するものであり得る。システムが電子取引システムである場合、少なくとも1つのサービスが、注文受付サービス；注文取消サービス；注文変更サービス；注文適合サービ

50

ス；予め実行された取引を行うためのサービス；またはクロス取引を行うためのサービスのうちのいずれか1つを含み得る。

【0021】

一次サーバにおけるサービスは、外部呼び出しが完全にバックアップサーバに転送されたことをバックアップサーバが認証した場合に限って、メッセージの処理が行われたことをクライアントが認証するように構成され得る。

【0022】

一次サーバにおけるサービスは、外部呼び出しがバックアップサーバに完全に転送されたという結果をバックアップサーバが認証するかどうかに関わらず、メッセージの処理が行われたことをクライアントが認証するように構成され得る。バックアップサーバが、所定の期限内に外部呼び出しが当該バックアップサーバに完全に転送されたという結果を認証しない場合、一次サーバは、バックアップサーバにおいて障害が発生したとみなし得る。

【0023】

本発明の他の形態は、ネットワークを介して相互接続された少なくとも2つのサーバのうちのいずれか1つを選択可能であって、選択した方と接続する少なくとも1つのクライアント；クライアントに接続される場合に一次サーバに指定されるサーバのうちの1つ、および、クライアントに接続されない場合にバックアップサーバに指定されるもう一方のサーバ；および、一次サーバにメッセージを送信するように構成された少なくとも1つのクライアントを有する、システムにおけるフェイルオーバーのための方法である。

【0024】

当該方法は、サーバの各々と同一で、かつ、サービスに対応するサーバが、一次サーバかバックアップサーバのいずれで動作しているかを認識しない、少なくとも1つのサービスを用いてメッセージを処理するようにサーバを構成する工程、サーバが一次サーバであるか、バックアップサーバであるかを表示するライブラリを維持するようにサーバを構成する工程、その各ライブラリに対して外部呼び出しを行うようにサービスを構成する工程、および、外部呼び出しを完了させて一次サーバにおけるサービスに外部呼び出しの結果を戻し、バックアップサーバにおけるサービスに外部呼び出しの結果を転送するように、一次サーバにおけるライブラリを構成する工程を含む。

【0025】

本発明の他の形態は、選択された少なくとも1つのクライアントと接続可能なネットワークを介して相互接続された少なくとも2つのサーバのうちのいずれか1つのサーバ上で実行可能な一連のプログラミングインストラクションを保存するコンピュータ読み取り可能な記録媒体であって、当該相互接続されたサーバが、少なくとも1つのクライアントと選択的に接続可能である記録媒体を提供する。サーバのうちの1つがクライアントと接続されたときに、当該サーバが一次サーバに指定され得、このとき、クライアントに接続されていないもう一方のサーバが、バックアップサーバに指定され得る。少なくとも1つのクライアントが一次サーバにメッセージを送信するように構成されている。プログラミングインストラクションは、

サーバの各々と同一で、かつ、サービスに対応するサーバが一次サーバまたはバックアップサーバとしてのいずれで動作するかを認識しない、少なくとも1つのサービスを用いてメッセージを処理するようにサーバを構成するためのインストラクション、

サーバが一次サーバであるか、バックアップサーバであるかを表示するライブラリを維持するようにサーバを構成するためのインストラクション、

その各ライブラリに対して外部呼び出しを行うようにサービスを構成するためのインストラクション、および、

外部呼び出しを完了させて外部呼び出しの結果を一次サーバにおけるサービスに戻し、バックアップサーバにおけるサービスに外部呼び出しの結果を転送するようにライブラリを構成するためのインストラクション、を含む。

【図面の簡単な説明】

【0026】

本発明は、添付の図面を参照し、例示することによってのみ説明される。

【図1】本発明の実施形態に係るフェイルオーバのためのシステムの概略図である。

【図2】通常状態で動作する場合の図1のシステムの概略図であって、システムにおけるサービスを実行する様々なソフトウェア部材の例示の詳細を含む。

【図3】本発明の他の実施形態に係る通常状態でのフェイルオーバのためのシステムの起動方法を示すフローチャートである。

【図4】図3の方法に関するパフォーマンス中の図2のシステムを示す。

【図5】図3の方法に関するパフォーマンス中の図2のシステムを示す。

【図6】図3の方法に関するパフォーマンス中の図2のシステムを示す。

【図7】図3の方法に関するパフォーマンス中の図2のシステムを示す。

【図8】図3の方法に関するパフォーマンス中の図2のシステムを示す。

【図9】図3の方法に関するパフォーマンス中の図2のシステムを示す。

【図10】図3の方法に関するパフォーマンス中の図2のシステムを示す。

【図11】図3の方法に関するパフォーマンス中の図2のシステムを示す。

【図12】図3の方法に関するパフォーマンス中の図2のシステムを示す。

【図13】図3の方法に関するパフォーマンス中の図2のシステムを示す。

【図14】図3の方法に関するパフォーマンス中の図2のシステムを示す。

【図15】本発明の他の実施形態に係るフェイルオーバのための方法を示すフローチャートである。

【図16】本発明の他の実施形態に係る一次単独(primary-only)状態で動作するサーバのうちの一つについての図2のシステムを示す。

【図17】本発明の他の実施形態に係る一次単独(primary-only)状態で動作するサーバのうちのもう一方についての図16のシステムを示す。

【図18】本発明の他の実施形態に係る一次単独状態でサーバのうちの一つを動作させるための方法を示すフローチャートである。

【図19】通常状態から、本発明の他の実施形態に係る一次単独状態において動作するバックアップサーバまでフェイルオーバするための方法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】

【0027】

ここで図1を参照すると、フェイルオーバのためのシステムは主に符号50で示されている。システム50は複数の遠隔クライアント54-1および54-2を備える(ここで言う「クライアント54」とは、総じて集合的な「クライアント54」を指す。図面における他の部材を指すときも同様に用語を定義する)。クライアント54をネットワーク58に接続する。ネットワーク58は、インターネットやローカルエリアネットワーク、広域ネットワークまたはその組合せなど、あらゆるタイプのコンピュータネットワークであり得る。続いて、第一サーバ62-1および第二サーバ62-2にネットワーク58を接続する。その結果、下記に詳述するように、クライアント54がネットワーク58を介してサーバ62-1および62-2にそれぞれ交信することができる。

【0028】

クライアント54は、個々のクライアント54を用いてサーバ62-2にリクエストを送信する個人および/または団体にそれぞれ属する。便宜上、そのような個人または団体をここでは取引者Tとし、クライアント54-1を使用するものを取引者T-1、クライアント54-2を使用するものを取引者T-2とする。各クライアント54は、典型的には、キーボードおよびマウス(または他の入力装置)と、モニタ(または他の出力装置)ならびに上記キーボード、マウス、およびモニタに接続されたデスクトップモジュールとを有するパーソナルコンピュータなどのコンピュータデバイスであって、1つ以上の中央演算処理装置と、揮発性メモリ(すなわち、ランダムアクセスメモリ)と、非揮発性メモリ(すなわち、ハードディスクデバイス)と、ネットワーク58を介してクライアント54

10

20

30

40

50

の通信を可能にするネットワークインターフェイスとを内蔵するパーソナルコンピュータなどのコンピュータデバイスである。しかし、クライアント54は、電子手帳、携帯電話、ノート型パソコン、e-mailページング装置などのあらゆるタイプのコンピュータデバイスであり得ることを理解されたい。

【0029】

サーバ62は、UNIX(登録商標)オペレーティングシステムを実行する、カリフォルニア州パロアルト(Palo Alto Calif)に住所を有するサンマイクロシステムズ株式会社(Sun Microsystems Inc.)製のSun Fire V480など、クライアント54からメッセージを受信および処理するように動作可能なあらゆるタイプのコンピュータデバイスであり得る。さらに、サーバ62は、約900メガヘルツで各々動作する4つの中央演算処理装置、および、4ギガバイトのランダムアクセスメモリおよびハードディスクドライブなどの非揮発性記憶装置を有する。サーバ62に適する他のタイプのコンピュータデバイスとしては、アメリカ合衆国コロラド州80537, ラブランド, サウス タフト 800(800 South Taft, Loveland, CO 80537)に住所を有するヒューレットパッカード株式会社(Hewlett-Packard Company)製のHP ProLiant BL25pサーバが挙げられる。しかし、これらの具体的なサーバは単なる例示であり、サーバ62-1および62-2について他の様々なタイプのコンピュータ環境が本発明の範囲内であることを強調したい。サーバ62-1によって受信および処理されるメッセージのタイプは特に限定されていないが、本実施形態において、サーバ62-1はオンライン取引システムを実行するので、オンラインで取引可能なセキュリティを購入する、売る、取り消すなどのリクエストを含むメッセージを処理することができるものである。より具体的には、サーバ62-1は中央適合エンジン(不図示)を維持するように動作可能であり、そこで互いに対してリクエストを実行し、さらに、注文を行う中央保存装置に対してリクエストを実行し、セキュリティ取引を処理する。

10

20

【0030】

サーバ62-2は、通常、サーバ62-1と同一の(または、少なくとも実質的に同一の)コンピュータ環境を有する。下記にさらに説明するように、ハードウェア、オペレーティングシステム、アプリケーションなどを含むコンピュータ環境が、サーバ62-2として選択される。このとき、当該サーバ62-2は、サーバ62-1の障害発生時にサーバ62-1の機能性を代替するように動作可能である。

30

【0031】

システム50はまた、サーバ62-1とサーバ62-2とを相互接続する複製リンク78を備える。本実施形態において、複製リンク78それ自体がメインリンク82とフェイルセーフリンク86とを備えており、サーバ62-1とサーバ62-2との間の通信において、より高いロバスト性を有する。

【0032】

一次サーバ62-1、バックアップサーバ62-2、および、複製リンク78の機能についてのさらなる詳細、ひいては、サーバ62-1および62-2を実行するのに使用可能な多種のハードウェアについて、下記の説明で明らかとなるだろう。

【0033】

図2において、サーバ62-1およびサーバ62-2をより詳細に示す。また、図1においてシステム50における様々な物理的な接続を実線で表すのに対し、図2においてはシステム50における様々なネットワーク上の接続を点線で示す。したがって、図2に示されたような接続は、サーバ62-1が一次サーバに指定されるとともに、サーバ62-2がバックアップサーバに指定されている通常状態で動作するシステム50を表すことを意図している。当該通常状態において、一次サーバ62-1がクライアント54からのリクエストに応答する。通常状態、および、システム50が動作可能な他の状態についての更なる詳細を以下に記載する。

40

【0034】

さらに図2を参照すると、サーバ62-1およびサーバ62-2がそれぞれ複数のソフ

50

トウェア構成要素を含んでおり、当該ソフトウェア構成要素は個々のハードウェア環境で実行され、クライアントからのリクエストに応答し、フェイルオーバ機能性を有する。

【 0 0 3 5 】

サーバ 6 2 - 1 およびサーバ 6 2 - 2 は、フェイルオーバエージェント 9 0 - 1 および 9 0 - 2 をそれぞれ備える。フェイルオーバエージェント 9 0 は相互交信し、リンク 7 8 および相互のインテグリティを定期的にテストするように動作可能である。本実施形態において、通常状態では、フェイルオーバエージェント 9 0 - 1 が、フェイルオーバエージェント 9 0 - 2 にキープアライブ信号（例えば、「アライブ状態になっているか？」）を定期的に送り、フェイルオーバエージェント 9 0 - 2 が定期的に応答する（例えば、「はい、アライブ状態です」）とされる。フェイルオーバエージェント 9 0 - 2 がそのようなリクエストに応答するとし、さらに、一次サーバ 6 2 - 1 が通常通り動作し続けるとすると、システム 5 0 は図 2 に示した通常状態を維持する。したがって、フェイルオーバエージェント 9 0 - 1 はまた、サーバ 6 2 - 1 における他のソフトウェア構成部材と交信し、通常状態が有効であることを示すように動作可能である。

10

【 0 0 3 6 】

フェイルオーバ 9 0 が、適切に、または、所望のように、リンク 7 8 を備えるメインリンク 8 2 とフェイルセーフリンク 8 6 の両方を利用するよう動作可能であることは明白である。この場合、メインリンク 8 2 とフェイルセーフリンク 8 6 のうちの少なくとも 1 つが動作している限り、システム 5 0 は通常状態を維持する。

【 0 0 3 7 】

サーバ 6 2 は、1 つ以上のクライアント 5 4 からの多種の要求を受信および処理することができる 1 つ以上のサービスをそれぞれ含む。サービスのタイプは特に限定されておらず、フェイルオーバ保護が必要とされるあらゆるタイプのサービス、アプリケーション、または処理などを含み得る。システム 5 0 がオンライン取引システムであるという、この単なる例示の実施形態において、サーバ 6 2 は注文受付サービス 9 4 および注文取消サービス 9 8 をそれぞれ含む。注文受付サービス 9 4 とは、その名前からもわかるように、特定のセキュリティの売り注文または買い注文を受け付けるためにクライアント 5 4 からの要求を受信するように構成されている。注文取消サービス 9 8 とは、その名前からもわかるように、特定のセキュリティの売り注文または買い注文を取り消すためにクライアント 5 4 からの要求を受信するように構成されている。このとき、当該セキュリティは、サービス 9 4 によってすでに受け付けられているが、当該特定の注文が実際に実行される前のものである。電子取引の分野の当業者であれば想定し得る実行可能な他のタイプのサービスは、これらに制限されるものではないが、注文適合、注文変更、取引開始、または、クロス開始を含む。本実施形態において、要件ではないものの、サービス 9 4 およびサービス 9 8 はマルチスレッドである。（ここで使用するとき、マルチスレッドとは限定的な意味で使用するものではなく、複合メッセージが同時に処理されるという並行処理の様々な形態を指す。これは、さらに、システムの非決定性特性の一因となる。例えば、複合処理、または、単一処理を用いた実行のマルチスレッドを用いて、マルチスレッドが実行され得る。）

20

30

サーバ 6 2 はまた、各々に付属する対応のサービス 9 4 およびサービス 9 8 にアクセス可能なライブラリ 1 0 2 をそれぞれ有する。各ライブラリ 1 0 2 はシーケンサ 1 0 6 およびキャッシュメモリ 1 1 0 を有する。下記に詳述するように、シーケンサ 1 0 6 が、シーケンスナンバを発生させ、ライブラリ 1 0 2 に関してサービス 9 4 または 9 8 からの要求に応答する。シーケンサ 1 0 6 - 2 は通常状態で非アクティブであり、そのような非アクティブ性は、オーバル型シーケンサ 1 0 6 - 2 を介してハッシュすることによって図 2 に示されている。（既定の特定状態において構成要素がアクティブであるか、非アクティブであるかを表示する他の構成要素において、ハッシュが使用される。）キャッシュメモリ 1 1 0 は、ライブラリ 1 0 2 によってなされる外部機能呼び出し結果の保存領域である。

40

【 0 0 3 8 】

各ライブラリ 1 0 2 はまた、システム 5 0 が並列処理を行う状態を維持する状態記録装

50

置 (state register) 114 を有する。当該状態記録装置 114 は、システム 50 が並行処理している状態を適合するために、その各フェイルオーバーエージェント 90 と絶えず交信する。図 2 において、システム 50 は通常状態で作動する。それに従って、サーバ 62 - 1 が現時点で一次サーバに指定されていることを状態記録装置 114 - 1 が表示し、一方で、サーバ 62 - 2 が現時点でバックアップサーバに指定されていることを情報記録装置 114 - 2 が表示する。しかし、下記に詳述するように、システム 50 の状態は、システム 50 における様々な構成材の動作状態によって変化し得る。

【0039】

各サーバ 62 はまた、サービス 94 および 98 に代わって外部リソースに外部呼び出しを行う機能を果たす外部リソースエージェント 118 を備えているが、当該呼び出しはライブラリ 102 を経由して行われる。外部リソースは、サービス 94 および 98 に外付けのリソースであって、各サーバ 62 上に設置されているリソースを備える。当該外部リソースは、各サーバに外付けの、オペレーティングシステムクロック (不図示) からのタイムスタンプ、および / または、リソースなどであり、電子取引システムの場合では、マーケットフィード (不図示) などである。当該マーケットフィードは、注文受付サービス 94 を経由して受け付けられた売り注文または買い注文のサブジェクトであり得る様々なセキュリティの市場価格についての最新情報を維持するものである。ここで、当業者であれば、サービス 94 および 98 が行うそのような外部リソースに対する呼び出しは、システム 50 の有する非決定性特性の一因となることがわかるだろう。通常状態において、外部リソースエージェント 118 - 1 のみがアクティブであり、外部リソースエージェント 118 - 2 は非アクティブである。オーバル型の外部リソースエージェント 118 - 2 を介してハッシュすることによって、外部リソースエージェント 118 - 2 の非アクティブ性が図 2 に示される。

【0040】

各サーバ 62 はまた、共有リソース 122 を維持しており、当該共有リソース 122 は、サービス 94 および 98 によって実行された処理ステップの結果、および / または、サービス 94 および 98 がアクセス可能であることを要求するデータを維持する。例えば、本発明に係る電子取引システムにおいて、共有リソース 122 が注文票を維持する。当該注文票は、サービス 94 によって受け付けられた注文の記録に過ぎない。それゆえ、注文サービス 94 は、例えば、共有リソース 122 において買い注文についての記録を作成してもよい。後に、注文を取り消してそれを共有リソース 122 に表示するために、注文取消サービス 98 が上記買い注文にアクセスする必要性が生じる可能性もある。同様に、サーバ 62 上で実行される適合サービス (不図示) が上記買い注文にアクセスする必要性が生じる可能性もある。当該適合サービスは、上記買い注文と、適する対応の売り注文とを市場法則に従って適合させ、上記買い注文および売り注文をアップデートし、適合が有効で取引が完了可能であることを示す。

【0041】

各サーバ 62 はまた、複製エージェント 126 を維持する。通常状態において、複製エージェント 126 - 2 のみがアクティブであって、複製エージェント 126 - 1 は非アクティブである。オーバル型エージェント 126 - 1 を介してハッシュすることによって、複製エージェント 126 - 1 の非アクティブ性が図 2 に示されている。下記に詳述するように、アクティブ複製エージェント 126 が、カウンターパートサーバ 62 におけるライブラリ 102 と交信し、一次サーバからバックアップサーバへの情報のミラーリングを助ける。

【0042】

ここで図 3 を参照すると、本発明の他の形態に係る、通常状態の間に要求を処理するための方法は、概して符号 300 に示されている。上記方法の説明を容易にするために、図 2 に示された、通常状態でのシステム 50 を用いて方法 300 を実行すると想定する。さらに、方法 300 についての詳細な説明によって、システム 50 およびその様々な構成材をよりよく理解できるだろう。しかし、便宜上の理由のみであるが、方法 300 の様々な

10

20

30

40

50

処理ステップが、システム50の特定の構成部材内で起こるように図3に示されている。そのような表示は、限定的な意味で解釈すべきものではない。しかし、当然のことながら、システム50および/または方法300は異なったものである可能性があり、本発明において互いに組み合わせて説明されたように機能する必要はなく、さらに、方法300におけるステップは、図示されたような注文において実施される必要はない。そのような変形は本発明の範囲内である。そのような変形は、ここに説明された他の方法およびシステム図にも適用される。

【0043】

まず、クライアントからメッセージを受信するステップ310から始める。メッセージタイプは特に限定されておらず、概して、サーバ上で実行するサービスのうちの1つについてのインプットの想定タイプを補完するものである。システム50上で実行されるとき、メッセージは、注文受付サービス94のためのインプットを意図する売り注文または買い注文か、または、注文取消サービス98のためのインプットを意図する取消注文であり得る。例えば、取引者T-1による買い注文が、クライアント54-1からのメッセージの内に受信され、当該メッセージが、ネットワーク58を介して注文サービス94-1に送信されることを想定する。この場合、ステップ310によると、注文受付サービス94-1が当該メッセージを受信する。図4に、ステップ310の例示の性能を示す。メッセージM(0₁)はクライアント54-1から発信され、注文受付サービス94-1においてサーバ62-1内に受信されるように示されている。表1は、注文受付メッセージM(0₁)の例示のフォーマットを示す。

【0044】

【表1】

メッセージM(0₁)

フィールド番号	フィールド名	例示の内容
1	取引者	取引者 T-1
2	セキュリティ名	ABC Co.
3	トランザクションタイプ	買い
4	量	1,000 ユニット

【0045】

より具体的には、取引者と名付けられた表1のフィールド1は、メッセージM(0₁)の発信元となる取引者が取引者T-1であると認識する。セキュリティネームと名付けられた表2のフィールド2は、取引のサブジェクトである特定のセキュリティの名前を、この例では、ABC Co.と認識する。トランザクションタイプと名付けられた表1のフィールド3は、フィールド2で認識されたセキュリティが買い注文なのか、売り注文なのか、などを認識する。この例では、トランザクションタイプは買い、であり、これは、買い注文であることを示している。量と名付けられた表1のフィールド4は、必要なセキュリティ量を認識する。この例では、量は1,000ユニットであり、ABC Co.の1,000ユニット買いを意図する。ここで、当業者であれば、注文の価格が、フィールド2のセキュリティに関する現行市場価格を基礎としていることから、表1の注文が市場注文であることがわかるだろう。

【0046】

ステップ310でメッセージを受信した後、方法300はステップ315に進み、この時点で、関連するサービスが、メッセージのさらなる処理に使用される外部データに対して、いずれかの呼び出しを行う。引き続き上記例でみていくと、ステップ315において、注文受付サービス94-1が、そのような外部呼び出しを一次ライブラリ102-1に対して行う。この例では、そのような呼び出しは下記i)、ii)を目的とすると想定さ

れている。

i) 注文が受信された時間を認識するメッセージM(0₁)において注文に割り当てられるタイムスタンプ、

ii) メッセージM(0₁)における注文で認識されるセキュリティに対する現行市場価格。

【0047】

図5に、ステップ315のパフォーマンスを点線で示す。点線130は、注文受付サービス94-1から一次ライブラリ102-1への呼び出しを示す。

【0048】

続いて、ステップ320において、一次ライブラリ102-1が呼び出しを行う。一次ライブラリ102-1はフェイルオーバーエージェント114-1から情報を得て(consult)、サーバ62-1が一次サーバに指定されており、かつ、システム50が通常状態であることを認証する。上記認証後、下記i)、ii)によって、サービス94-1が行う呼び出しに対して、一次ライブラリ102-1が応答する。

i) タイムスタンプを得るために、外部リソースエージェント118-1に対して外部呼び出しを行う。

ii) 現行市場価格を得るために外部リソースエージェント118-1にさらなる外部呼び出しを行う。

【0049】

その結果、ステップ325において、外部リソースエージェント118-1がタイムスタンプを得るためにシステムクロック(不図示)に外部呼び出しを、現行市場価格を得るために市場供給(不図示)に外部呼び出しを、それぞれ行う。

【0050】

図6に、ステップ320およびステップ325のパフォーマンスを点線で示す。点線は、外部リソースエージェント118-1を介するタイムスタンプのための呼び出しを132に、外部リソースエージェント118-1を介する市場価格のための呼び出しを134に、それぞれ示す。

【0051】

ここで、当業者であれば、外部呼び出し132および134によって、特にシステム50に非決定性特性を有するようになることがわかり、それゆえ、フェイルオーバーの発生において、復旧中に、当該復旧が取引者Tに自明であるように、システムにおける非決定性特性に対処しているフェイルオーバーシステムを提供するという点で、本発明独自の取り組みが理解できるだろう。(さらなる説明によって、両サーバ62が各メッセージに対する呼び出しを行うようにシステム50が変更されたと想定される。しかし、既定のメッセージMのいずれに関しても、市場の公正さを確保するためには、タイムスタンプに対する呼び出しが行われるその瞬間が非常に重要であって、両サーバ62が、同時に同じメッセージのためにタイムスタンプに対して呼び出しを行う確率は非常に低いと考えられる。それゆえ、各サーバ62が同じメッセージMについて異なる時間優先を割り当てることが可能であり、その結果、同じ機械処理の結果が異なることになる。フェイルオーバーの間、各サーバ62が一貫したビジネスデータを有することはなく、フェイルオーバーは無意味であり得る。)さらに読み進めると、当業者であれば、そのようなチャレンジに対する対処法および本発明の他の形態が理解できるだろう。

【0052】

ステップ330において、外部呼び出し132および外部呼び出し134の結果が一次ライブラリ102-1に戻ってくる。ステップ335において、キャッシュメモリ110-1内に呼び出し132および134の全ての結果が保存され、サーバ94-1に戻される。

【0053】

引き続き上記例でみていくと、呼び出し132の結果はタイムスタンプ2000年1月5日午後12時であると考えられる。さらに、呼び出し134の結果は\$2.00の市場

10

20

30

40

50

価格であると考えられる。表 2 および図 7 に、キャッシュメモリ 110-1 におけるこれらの結果の保存を示す。

【 0 0 5 4 】

【表 2】

ステップ335の後の、キャッシュメモリ110-1における例示の内容

記録番号	フィールド番号	フィールド名	例示の内容
1	1	メッセージ	M(O ₁)
1	2	タイムスタンプ	2000年1月5日午後12時
1	3	市場価格	\$2.00

10

【 0 0 5 5 】

ステップ 3 4 0 で、サービスが呼び出し結果を受信する。上記例で続けると、図 7 にも示されているように、表 2 で保存された呼び出し結果はサービス 9 4 - 1 に戻される。

【 0 0 5 6 】

続いて、ステップ 3 4 5 において、サービスが共有リソースのための要求を行う。上記例において、サービス 9 4 - 1 がライブラリ 1 0 2 - 1 に対して要求を行う。次に、ステップ 3 5 0 において、ライブラリ 1 0 2 - 1 が共有リソース 1 2 2 - 1 に対して、当該リソース 1 2 2 - 1 をロックするように指示を発行し、他のあらゆるサービス（例えば、サービス 9 8 - 1 またはサービス 9 4 - 1 内の他のスレッド）が共有リソース 1 2 2 - 1 にアクセスしないようにする。（下記に詳述するように、共有リソース 1 2 2 - 1 がすでにロックされている場合、共有リソース 1 2 2 - 1 が解除されるまで、方法 3 0 0 はステップ 3 4 5 で停止する。）ステップ 3 4 5 およびステップ 3 5 0 のパフォーマンスを図 8 に点線で示す。このとき、共有リソース要求を符号 1 4 0 で示す。共有リソース 1 2 2 - 1 のロックをパッドロック 1 3 8 で表す。

20

【 0 0 5 7 】

続いて、ステップ 3 5 5 において、共有リソースシーケンスナンバーが戻される。シーケンス 1 0 6 - 1 を利用するライブラリ 1 0 2 - 1 によって、このステップが実行される。当該シーケンス 1 0 6 - 1 は、メッセージ M (O₁) に関するシーケンスナンバーを生成する。上記例で続けると、シーケンスナンバー 1 が生成されることが考えられる。表 3 および図 8 に、キャッシュメモリ 110-1 内の結果の保存を示す。表 3 は表 2 をアップデートしたものであることに留意されたい。

30

【 0 0 5 8 】

【表 3】

ステップ355の後の、キャッシュメモリ110-1における例示の内容

記録番号	フィールド番号	フィールド名	例示の内容
1	1	メッセージ	M(O ₁)
1	2	タイムスタンプ	2000年1月5日午後12時
1	3	市場価格	\$2.00
1	4	シーケンスナンバ	1

40

【 0 0 5 9 】

続いて、ステップ 3 6 0 において、複製が要求される。サービス 9 4 - 1 が上記例におけるステップ 3 6 0 を実行し、当該サービス 9 4 - 1 がライブラリ 1 0 2 - 1 に複製実行の指示を送信する。ステップ 3 6 5 において、メッセージ、呼び出し結果およびシーケ

50

スナンバーの複製が開始される。上記例において、ライブラリ 102 - 1 が表 3 の内容を複製する。ステップ 365 については後に説明する。

【0060】

ステップ 370 において、呼び出し結果およびロックされた共有リソースを用いてメッセージを処理する。上記例において、サービス 94 - 1 がステップ 370 を実行する。サービス 94 - 1 は、表 3 の結果を生み出すために、表 3 の内容を用いて、サービス 94 - 1 に関する処理ステップを実行する。サービス 94 - 1 は注文受付サービスであり、メッセージ M(0₁) は買い注文を示すので、ステップ 370 において、サービス 94 - 1 は、共有リソース 122 - 1 に記録される買い注文を作り出し、例えば取引者 T - 2 からの買い注文に対して買い注文を続けて適合させるか、または、サービス 98 - 1 を用いた注文の取消など他の取引処理を行う。

10

【0061】

上記例の目的に関して、メッセージ M(0₁) が適合し得る共有リソース 122 - 1 に注文がない場合、ステップ 370 の結果は、単に、メッセージ M(0₁) に関する買い注文の詳細についての完全な記録を作り出すことである。表 4 に、ステップ 370 のパフォーマンスの例示結果を示す。

【0062】

【表 4】

ステップ370におけるパフォーマンスの例示の結果

記録番号	フィールド番号	フィールド名	例示の内容
1	1	タイムスタンプ	2000年1月5日午後12時
1	2	市場価格	\$2.00
1	3	シーケンスナンバ	1
1	4	取引者	取引者T-1
1	5	セキュリティ名	ABC Co.
1	6	トランザクションタイプ	買い

20

30

【0063】

続いて、ステップ 375 において、ステップ 370 のパフォーマンス結果が共有リソースに書き込まれ、その後、共有リソースが解除される。ステップ 370 におけるサービス 94 - 1 による表 4 の生成と、ステップ 375 における共有リソース 122 - 1 内の結果の保存を図 9 に示す。

【0064】

続いて、ステップ 380 では、ステップ 375 において結果が書き込まれたこと、および、ステップ 400 において複製が実行されたことについての認証を、サービスが認証する。上記例では、ステップ 380 において、サービス 94 - 1 が、表 4 が共有リソース 122 - 1 に書き込まれたという共有リソース 122 - 1 からの認証を待つ。同様に、ステップ 380 では、サービス 94 - 1 は、ステップ 365 によって開始された複製が完了したというステップ 400 からの認証を待つ。ステップ 365 および 400 については下記に詳述する。

40

【0065】

(代替的实施例において、ステップ 380 は、ステップ 390 に進む前に、ステップ 400 からの認証を実際に待つ必要はない。しかし、ステップ 380 は、ステップ 400 から最終的にはそのような認証を受信することを想定しており、そのような認証が得られない場合には、ステップ 380 はサーバ 62 - 2 が障害を起こしたとみなす。このとき、以下に説明するように、イベントサーバ 62 - 1 が方法 600 の実行を開始する。ここで、

50

当業者であれば、これが非同期モードの動作であり、サーバ62-2の状態を認証する速度が早いことが好まれる特定の状況において好適である可能性があることがわかるだろう。

続いて、ステップ390において、クライアントに認証が戻される。上記例において、ステップ390で、サービス94-1は、取引者T-1が要求した通りにメッセージM(0₁)が処理されたという認証メッセージをクライアント54-1に送信する。

【0066】

繰り返し述べるが、方法300におけるステップ390(すなわち、通常状態中の動作)は、ステップ380で初めて完了する。続いて、ステップ380は、ステップ365において開始される複製が完了して初めて完了する。ここで、ステップ365に戻ると、メッセージ、呼び出し結果、および、共有リソースシーケンスナンバーが複製される。本例では、ステップ360においてサービス94-1からの要求に応答するライブラリ102-1が、ステップ365を実行する。それゆえ、ライブラリ102-1が表3の内容を一括し、それを複製エージェント126-2に伝達する。

【0067】

図10に、ステップ365、370、375、395、400および390のパフォーマンスを示す(図10は、図9に示されたステップ370および375のパフォーマンスに関する表示に基づく)。ステップ365、すなわち、ライブラリ102-1のキャッシュメモリ110-1から複製エージェント126-2への表3における伝達のステップを、線142で表す。図9に関連して上述したように、図10にステップ370および375を示す。ステップ395、すなわち、メッセージ、呼び出し結果、および、共有リソースシーケンスナンバーの待機するステップを、複製エージェント126-2内部に表された表3のように楕円形で示す。ステップ400、すなわち、複製エージェント126-2から(ライブラリ102-1を介して運ばれる)サービス94-1への複製の認証に戻ることを144の線で示す。ステップ390、すなわち、サービス94-1からクライアント54-1までの認証に戻るステップを、点線146で示す。

【0068】

上記記述で、通常状態での動作中に一次サーバが行う1つのメッセージ処理の説明は実質的に終わりである。ここで、ステップ400を介するステップ310の上述に従って、一次サーバ62-1が直列、および/または略並行かのいずれかで、複合メッセージを処理し得ることを理解すべきである。例えば、サービス94-1が、あるメッセージMを扱っているとき、同様に、サービス98-1もまた、実質的に上記のように他のメッセージMを処理することができる。このとき、ライブラリ102-1が両サービス94-1、98-1と通信する。加えて、サービス94-1におけるあるスレッドが、あるメッセージMを扱うとき、サービス94-1における他のスレッドがまた、実質的に上記のように他のメッセージMを処理することができる。このとき、ライブラリ102-1がサービスにおける両スレッドと通信する。ステップ350は、サービス94-1と98-1(または、そのスレッド)との間のコンテンションを避けるように共有リソース122-1を確実にロックし、一度に共有リソース122-1と通信し得るサービスが、確実にそれらサービスのうちの一つだけであるようにする。(通信とは、リーディング、ライティング、削除に留まらず、それらを含むあらゆるタイプの機能を含み得ることに留意されたい。)回避すべきコンテンションの例として、共有リソース122-1が、与えられた注文を取り消すためにロックされているときに、注文取消サービス98-1が共有リソース122-1からの読み取り、および、リソース122-1への書き込みを行い、それによって、取り消された注文と適合サービス(不図示)とが適合できないようになることが挙げられる。

【0069】

同じ特徴から、ステップ335がシーケンサ106-1を利用し、サービス94-1または98-1(またはそのスレッド)のいずれかがメッセージMを扱っているのに関わらず、各メッセージに固有のシーケンスナンバーを生成する。したがって、共有リソース1

10

20

30

40

50

22 - 1がロックされているときに、特定のサービス94 - 1または98 - 1（またはそのスレッド）がステップ345において共有リソース122 - 1への要求を行うときがあり得る。それゆえ、ステップ345より後に進む前に、共有リソース122 - 1が解除されるまでは特定のサービス（またはそのスレッド）は停止する。

【0070】

通常状態で動作中である、プライマリサーバ62 - 1によるメッセージ処理について記述したが、方法300に関して、ステップ405およびその先のステップにおけるパフォーマンス、ならびに二次サーバ62 - 2によるメッセージ処理を説明する。

【0071】

再び図3を参照すると、ステップ405において、共有リソースシーケンスナンバーに従って、メッセージ、呼び出し結果およびシーケンスナンバーが発信（dispatch）する。上記例で続けると、この時点で、メッセージM(0₁)（すなわち、表3からの記録1におけるフィールド1の内容）をサービス94 - 2に発信し、一方で、呼び出し結果（すなわち、表3からの記録1におけるフィールド2および3の内容）およびシーケンスナンバー（すなわち、表3からの記録1におけるフィールド4の内容）を二次ライブラリ102 - 2に発信する。

【0072】

したがって、ステップ310Sでは、ステップ310でサービス94 - 1がクライアント54 - 1からメッセージM(0₁)を受信したのとほぼ同じ方法で、サービス94 - 2が複製エージェント126 - 2からメッセージM(0₁)を受信する。サービス94 - 2において、クライアントからメッセージM(0₁)を受信していた。この段階で、サービス94 - 2は、全ての面においてサービス94 - 1と実質的に同一であることが明白になるだろう（同様に、サービス98 - 2はサービス98 - 1と実質的に同一である）。サーバ62 - 2におけるサービス94 - 2は、サーバ62 - 1におけるサービス94 - 1の動作とほぼ同じ方法で動作する。換言すれば、サーバ62 - 1におけるサービス94 - 1がステップ310、315、340、345、360、370、380、および390を実行するのと同じ方法で、サービス94 - 2がステップ310S、315S、340S、345S、360S、370S、380S、および390Sを実行する。サービス94 - 1およびサービス94 - 2の両方もが、それらサービスが内部で動作させている特定のサーバが、一次サーバかバックアップサーバかのいずれに指定されているかを認識しない。これは、本発明の数ある利点のうちの一つを提示する。すなわち、サービスは、一括で2つ（またはそれ以上）のサーバのためのサービスをを進めることができ、一次サーバに指定されたサーバのための一連のサービス、および、バックアップサーバに指定されたサーバのための一連のサービスをを進める必要がないのである。

【0073】

しかし、各ライブラリ102は、個々のフェイルオーバーエージェント90および状態記録装置114と通信する中で、その個々のサーバ62が一次サーバかバックアップサーバのいずれに指定されているかを認識している。したがって、サービス94 - 2がステップ315Sを実行して呼び出しを行うとき、ライブラリ102 - 2は外部リソースエージェント118 - 2を利用せず、ステップ415において、ステップ410でライブラリ102 - 2が受信した呼び出し結果（すなわち、表3からの記録1におけるフィールド2および3の内容）を単に戻すだけである。

【0074】

図11に、ステップ405、310S、410のパフォーマンスを示す。図12に、ステップ315S、415、および340Sを示す。

【0075】

同じ特徴から、サービス94 - 2がステップ345Sを実行して共有リソースに要求するとき、共有リソース122 - 2をロックすること、および、ステップ425で共有リソースシーケンスナンバー（すなわち、表3からの記録1におけるフィールド4の内容）を戻すことによって、ライブラリ102 - 2が応答する。当該シーケンスナンバーはステッ

10

20

30

40

50

プ 4 1 0 においてライブラリ 1 0 2 - 2 が受信するものであって、シーケンサ 1 0 6 - 2 を利用するものではない。

【 0 0 7 6 】

図 1 3 に、ステップ 3 4 5 S、4 2 0、4 2 5 のパフォーマンスを示す。

【 0 0 7 7 】

同じ特徴から、サービス 9 4 - 2 がステップ 3 6 0 S を実行して複製を要求するとき、実際の複製実行によってではなく、基本的にステップ 4 0 0 を模倣して、ステップ 3 8 0 S においてサービス 9 4 - 2 に対する複製認証を戻すことによって、ライブラリ 1 0 2 - 2 がステップ 4 3 0 において応答する。ステップ 3 7 0 および 3 7 5 と実質的に、サービス 9 4 - 2 が独立して表 4 の内容を生成し、共有リソース 1 2 2 - 2 の内部に保存するようにステップ 3 7 0 S および 4 3 5 を実行する。

10

【 0 0 7 8 】

図 1 4 に、ステップ 3 7 0 S および 4 3 5 のパフォーマンスを示す。

【 0 0 7 9 】

同様に、ステップ 3 8 0 およびステップ 3 9 0 と同じ方法で、ステップ 3 8 0 S および 3 9 0 S を実行する。但し、上記ステップは、ステップ 3 9 0 S において戻された認証は、クライアント 5 4 - 1 ではなく、複製エージェント 1 2 6 - 2 に戻される。

【 0 0 8 0 】

当該段階、方法 3 0 0 における当該パフォーマンスの結論として、表 4 のように、処理メッセージ M (0₁) の結果が共有リソース 1 2 2 - 1 と共有リソース 1 2 2 - 2 との両方に保存されることがわかるだろう。ステップ 3 1 0 S、3 1 5 S、3 4 0 S、3 4 5 S、3 6 0 S、3 7 0 S、3 8 0 S、3 9 0 S のパフォーマンスとステップ 3 1 0、3 1 5、3 4 0、3 4 5、3 6 0、3 7 0、3 8 0、3 9 0 のパフォーマンスとの間の実際の待機時間は実際ごく短いこともわかるだろう。そのような待機時間は、ステップ 3 6 5 におけるネットワークの待機時間と、ステップ 3 9 5 および 4 0 5 の処理とによって決定可能であって、それらは非常に高速であり得る。あらゆるイベントにおいて、システム 5 0 は、待機時間がハードディスクにバックアップ情報を書き込むよりも数段速くなるように構成されており、これは本発明の他の利点である。

20

【 0 0 8 1 】

したがって、サービス 9 4 - 1 (および後続のサービス 9 4 - 2) を用いてセキュリティの売買注文を受け付けるためのメッセージを処理するのに、方法 3 0 0 を使用することができる。同様に、サービス 9 8 - 1 (および後続のサービス 9 8 - 2) を用いてこれらの注文を取り消すのに、方法 3 0 0 を使用することができる。サーバ 6 2 - 1 内に付加的なサービスを生成および包含することができ、当該サービスのためのロバストフェイルオーバーを備えるサーバ 6 2 - 2 上に当該サービスを容易に配置することができる。しかし、当該サービスは、サーバ 6 2 - 1 上にあるサービス用の 1 組のコードを必要とせず、サーバ 6 2 - 2 上にあるサービスのための他の組のコードを必要とする。このとき、特定のサービスのための 1 組のコードとは、両サーバにとって必要なものすべてである。おそらく、ある観点からより重要なことは、システム 5 0 が、通常はハードディスクへの書き込みに付随して起こる減速を起こすことなく、フェイルオーバー発生時の結果を実質的に保証することができる。

30

40

【 0 0 8 2 】

通常状態において、サーバ 6 2 - 2 がサーバ 6 2 - 1 において実行される処理の最新ミラーを維持するので、サーバ 6 2 - 1 がレフトオフ (left-off) であるサーバ 6 2 - 1 の処理タスクをサーバ 6 2 - 2 が想定することによって、サーバ 6 2 - 1 の障害が速やかに復旧し得る。図 1 5 は、1 対のサーバを運用するための方法 5 0 0 を表すフローチャートである。このとき、該サーバのうちの一つが一次サーバに指定されており、もう一つがバックアップサーバに指定されている。システム 5 0 を用いて実行するとき、ステップ 5 0 5 において、両サーバが有効かどうか決定される。フェイルオーバーエージェント 9 0 および状態記録装置 1 1 4 を用いてステップ 5 0 5 を実行する。ハイの場合、ステップ 5 0

50

5 はステップ 5 1 0 に進み、そこで、システム 5 0 が方法 3 0 0 に関して上述のような通常状態で動作する。両サーバが無効であると判断されない限りは、ステップ 5 0 5 およびステップ 5 1 0 が循環し続ける。両サーバが無効であると判断された場合、方法はステップ 5 2 0 に進む。ステップ 5 2 0 において、第一サーバのみが有効かどうかを決定する。たとえば、何らかの理由で、フェイルオーバーエージェント 9 0 - 1 がフェイルオーバーエージェント 9 0 - 2 との接続を確立できない場合、ステップ 5 2 0 において第一サーバのみが有効であることを決定し、方法 5 0 0 はステップ 5 3 0 に進む。当該ステップ 5 3 0 において、システム 5 0 が一次単独状態 (primary-only state) で動作する。フェイルオーバーエージェント 9 0 - 1 がフェイルオーバーエージェント 9 0 - 2 との接続が確立できない理由として考えられるのは、これらに留まるわけではないが、サーバ 6 2 - 2 が致命的な損傷を被ったか、または、リンク 7 8 が切断されたかである。

10

【 0 0 8 3 】

第一サーバが無効である場合、方法 5 0 0 はステップ 5 2 0 からステップ 5 4 0 に進む。ステップ 5 4 0 では、第二サーバのみが有効であるかどうか決定される。異なる場合、方法 5 0 0 は例外として終了する。しかし、第二サーバが有効であると判断された場合、方法 5 0 0 はステップ 5 4 0 からステップ 5 5 0 に進む。ステップ 5 5 0 において、第二サーバがさらなる処理を実行できるように、システム 5 0 がフェイルオーバーする。続いて、ステップ 5 6 0 において、第二単独状態においてさらなる処理を行うようにオペレーションが続行される。両サーバが再び有効になるまで、方法 5 0 0 はステップ 5 6 0 と 5 7 0 との間を循環する。両サーバが再び有効になった時点で、方法 5 0 0 はステップ 5 1 0 に進み、システム 5 0 は通常状態に戻る。

20

【 0 0 8 4 】

図 1 6 は、一次単独状態でのシステム 5 0 の例を示すものであり、ここではサーバ 6 2 - 1 が一次サーバに指定されているが、サーバ 6 2 - 2 はオフライン (または、リンク 7 8 の障害によって無効) である。図 1 6 において、サーバ 6 2 - 1 が一次単独状態で動作するので、状態記録装置 1 1 4 - 1 は、サーバ 6 2 - 1 が現在一次サーバに指定されており、かつ、一次単独状態で動作中であることを表示する。

【 0 0 8 5 】

図 1 7 は、二次単独状態におけるシステム 5 0 の例を示すものであり、ここではサーバ 6 2 - 2 が一次サーバに指定されているが、サーバ 6 2 - 1 はオフラインである。図 1 7 において、サーバ 6 2 - 2 が一次単独状態で動作しているので、状態記録装置 1 1 4 - 2 は、サーバ 6 2 - 2 が現在一次サーバに指定されており、かつ、一次単独状態で動作中であることを表示する。

30

【 0 0 8 6 】

図示されていないが、システム 5 0 はまた、サーバ 6 2 - 2 が一次サーバに指定されており、サーバ 6 2 - 1 がバックアップサーバに指定されている通常状態で構成され得ることに留意されたい。

【 0 0 8 7 】

図 1 8 は、サーバ 6 2 のうちのひとつのみが有効なときにメッセージを処理するための方法 6 0 0 を表すフローチャートである。方法 5 0 0 におけるステップ 5 3 0 にあるサーバ 6 2 - 1 が方法 6 0 0 を実行するか、または、方法 5 0 0 におけるステップ 5 6 0 にあるサーバ 6 2 - 2 が方法 6 0 0 を実行する。ここで、当業者であれば、方法 6 0 0 が方法 3 0 0 における一次サーバの動作を実質的に反映していることがわかるだろう。より具体的には、方法 3 0 0 におけるステップ 3 1 0 ~ 3 6 0 およびステップ 3 7 0 ~ 3 9 0 が、方法 6 0 0 におけるカウンターパート (同じ番号を有し、後ろに F が付記する) に対応する。しかし、方法 6 0 0 におけるステップ 3 6 5 F は、方法 3 0 0 におけるステップ 3 6 5 とは異なる。ステップ 3 6 5 F は、方法 3 0 0 におけるステップ 4 3 0 に対応している。ステップ 3 6 5 F では、ライブラリ 1 0 2 は複製が成されたという認証を単に複製することによって、サービス 9 4 (または 9 8) からの複製要求に回答するので、サービス 9 4 (または 9 8) がステップ 3 8 0 F においてそのような認証を受信することができ、方法

40

50

600がステップ390Fに進むことができる。

【0088】

図19は、一次サーバから、方法500におけるステップ550を実行するのに使用可能であるバックアップサーバにフェイルオーバーするための方法700を表すフローチャートである。フェイルオーバーエージェント114-2がサーバ62-1の障害(故障、または、何らかの理由でもはや有効ではないこと)を発見した場合、例えば、サーバ62-2が方法700を実行する。サーバ62-2は、サーバ62-1が障害を起こしており、当該サーバが一次サーバであるとみなすのだが、クライアント54は既にサーバ62-1と交信しているので、当該クライアント54はそのままサーバ62-1と交信し続ける。この場合、方法700はステップ710から始まる。当該ステップにおいて、複製エージェントキューが解消される。上記例において、サーバ62-1が障害を起こす前に当該サーバ62-1において処理された全メッセージ(および関連の外部呼び出し)処理を削除および複製するために、ステップ405(および後続のステップ310S、315S、340S、345S、360S、370S、380S、390S、410、415、420、425、430、および435)に従って、サーバ62-2は複製エージェント126-2に保存された全データの処理を続ける。サーバ62-1がステップ370において障害を起こした場合は、サーバ62-2はクライアントから複製メッセージを受信してもよい。該クライアントは、ギャップリカバリや他の例、例えば、出願人の同時係属出願であるUS公開公報US20050138461において説明されているタイプのリカバリなどのリカバリ
10
20
プロトコルを実行する。クライアントが、メッセージを処理するサーバ62-1からの認
証を受信することはない。この場合、サーバ62-2は、複製メッセージを認識し、複製
メッセージの再処理を試みることなく、単に同じ返答を返すように構成されている。

【0089】

続いて、ステップ720において、複製エージェントが停止する。上記例において、複製エージェント126-2は、サーバ62-1から受信されたデータのキューをもはや維持しないように停止するか、または、サービス94-2および98-2にメッセージを送信するように構成されている。ステップ730では、外部リソースエージェントおよびシーケンサがアクティブになる。上記例において、外部リソースエージェント118-2がアクティブになるので、該エージェントは、方法600におけるステップ325Fおよび
30
ステップ330Fで示された外部機能呼び出しを行うように構成されている。シーケンサ
106-2も同様にアクティブになるので、該シーケンサ106-2は、方法600の
ステップ335において示されたシーケンスナンバーを割り当てるように構成され得る。同様に、シーケンサ106-2は、方法600におけるステップ355Fにおいて示された
外部機能呼び出しを行うように構成され得る。続いて、ステップ740において、フェ
イルオーバーエージェントが一次単独状態を表示するように設定されている。上記例において、フェイルオーバーエージェント114-2が一次単独状態を表示するよう設定されている
40
ので、方法600におけるステップ320F、335F、350F、355Fおよび36
5Fに従ってライブラリ102-2が動作を認識する。続いて、ステップ720では、サ
ーバの存在(presence)がクライアントに伝えられる。上記例において、サーバ62-2
は、サーバ62-2がクライアント54からのメッセージを受容および処理する準備が
40
できたことを、ネットワーク58を介してクライアント54に伝える。これを行う方法は特
に限定されていないが、方法300の開始に先立ってサーバ62-1がクライアント54
に伝えた方法と実質的に同じ方法である。セッションプロトコルがギャップリカバリを実
行することができるので、各サイドは、カウンターパーティが受信しない可能性のある交
信を再送信し得る。この時点で、システム50は図17に示された状態である。該状態
では、サーバ62-2が一次サーバに指定されており、システム50は、一次サーバに指定
されたサーバ62-2を用いて一次単独状態で動作する準備ができています。この時点で、
方法は方法500におけるステップ560に戻り、そこで、クライアントからのメッセ
ージを受信し、方法600に従って処理が行われる。

【0090】

10

20

30

40

50

本発明の様々な特徴および構成要素の特定の組合せのみをここで述べたが、当業者であれば、開示された特徴、および構成要素、および/または、これら特徴および構成要素の代替的組合せについての所望の部分を、その希望に応じて利用し得ることが明らかであろう。例えば、システム50は2つのサーバ62-1および62-2を含んでいるが、あらゆる数のサーバが利用可能であると考えられる。ここに記載された方法の応用形を用いて、あるサーバを一次サーバに指定し、あらゆる数の付加的サーバをバックアップサーバに指定することができ、かつ、当該付加的サーバを直列または並列に接続することができる。そのような付加的サーバは、ここに開示されたサーバ62と実質的に同じコンピュータ環境および構造を有する。該付加的サーバは、あらゆる場合において、ライブラリおよび他のソフトウェア構成要素と交信する同一のサービスを有しており、これらサービスに代

10

【0091】

また、方法300が変更可能であることを理解されたい。例えば、方法300は完全に同時に動作するように構成可能であって、その場合、一次および二次共有リソースの両方が、特定のサービスによる処理の結果とともに書き込みされたことを鑑みると、一次サーバは、メッセージが処理されたことをクライアントに対して認証するのみである。ステップ380Sが実行されると直ちに、ステップ400が実行されるように方法300を変更することによって、上記を実行することができる。

【0092】

20

また、フェイルオーバーシステムを以下のように構成してもよい。

【0093】

< 1 >

ネットワークを介して相互接続された少なくとも2つのサーバのうちのいずれか1つを選択可能であって、選択した方と接続する少なくとも1つのクライアントを備えており、

上記サーバのうちの1つは、上記クライアントに接続されるときに一次サーバに指定されており、かつ、上記サーバのうちの残りは、上記クライアントに接続されないときにバックアップサーバに指定されており、

上記少なくとも1つのクライアントが、上記一次サーバにメッセージを送信するように構成されており、

30

上記サーバの各々が、異なるタイプの上記メッセージを処理する複数のサービスを介して上記メッセージの全てを処理するように構成されており、

上記サービスの各々が、上記メッセージの処理の結果に基づいて上記サーバの各々によって維持された共有リソースに対してアクセスおよびアップデートの少なくともいずれかを行なうように構成されており、

上記サーバの各々が、上記サービスと連結されたライブラリを維持しており、

上記ライブラリは、上記一次サーバによって維持されるとき、

i) 外部リソースに対して少なくとも1つの外部呼び出しを実行すること、

ii) 上記個々のメッセージに関するサービスからの要求に基づいて、各メッセージにシーケンスを付けること、

40

iii) 上記個々のメッセージに関する上記サービスに対する外部呼び出し結果および上記シーケンスの結果を戻すこと、並びに、

iv) 上記サービスによって共有されるキャッシュメモリに、上記外部呼び出し結果および上記シーケンスの結果を保存すること、

によって、上記個々のメッセージに関する上記サービスからの要求に応答するように構成されており、かつ、

上記一次サーバおよび上記バックアップサーバは、

i) 上記一次サーバにおいて受信される上記少なくとも1つのクライアントからのメッセージ、および、

ii) 上記保存された上記外部呼び出し結果、

50

が上記バックアップサーバに複製されるように、互いに接続されており、

上記ライブラリは、上記バックアップサーバによって維持されるとき、

i) 上記一次サーバが複製した上記保存された外部呼び出し結果の内容を利用して上記外部呼び出し結果をライブラリに戻すこと、
 によって、上記個々のメッセージに関するサービスからの要求に应答するように構成されており、

上記バックアップサーバは、上記一次サーバと同じシーケンスにおける上記メッセージの処理に上記キャッシュメモリを利用し、上記一次サーバおよび上記バックアップサーバにおいて上記共有リソースは実質的に同一であり、上記一次サーバが障害を起こした場合、上記バックアップサーバが上記一次サーバに指定され、実質的にトランスペアレントな方法において、上記クライアントに代わって他のさらなるメッセージを処理し続けるフェイルオーバのためのシステム。

【0094】

< 2 >

上記キャッシュメモリが揮発性メモリ内で維持される、< 1 >に記載のシステム。

【0095】

< 3 >

上記システムが電子取引システムの一部であって、上記サービスは取引エンジンに含まれており、

上記システムが、2つの他のさらなるメッセージを上記一次サーバにそれぞれ送信する2つのさらなるクライアントを備えており、

上記他のさらなるメッセージの各々が上記メッセージと実質的に同じ方法で処理される、< 1 >に記載のシステム。

【0096】

< 4 >

上記メッセージのうちの1つが買い注文を示しており、上記受信された2つの他のさらなるメッセージが、実質的に同じ回数 (different but nearly identical times) 受信されてあり、

上記2つの他のさらなるメッセージの両方が、上記買い注文に適応する売り注文を示しており、上記一次サーバが障害を起こした場合に、上記バックアップサーバが上記メッセージの処理を続けて、上記売り注文のうち注文のタイミングが早い方が上記買い注文と適合されるようになっている、< 3 >に記載のシステム。

【0097】

< 5 >

上記外部リソースが、オペレーティングシステムのタイムスタンプおよびマーケットフィールドである、< 4 >に記載のシステム。

【0098】

< 6 >

上記外部リソースが、オペレーティングシステムのタイムスタンプである、< 1 >に記載のシステム。

【0099】

< 7 >

ネットワークを介して相互接続された少なくとも2つのサーバのうちのいずれか1つを選択可能であって、選択した方と接続する少なくとも1つのクライアントを備えてあり、

上記サーバのうちの1つは、上記クライアントに接続されるときに一次サーバに指定され、かつ、上記サーバのうちの残りは、上記クライアントに接続されないときにバックアップサーバに指定されており、

上記少なくとも1つのクライアントが、上記一次サーバにメッセージを送信するように構成されており、

上記サーバは、上記クライアントに代わって複合スレッドを用いて上記メッセージを処

10

20

30

40

50

理するように構成されており、さらに、上記メッセージの処理に関して使用される上記スレッドの各々にアクセス可能な共有リソースを維持するように構成されており、

上記一次サーバは、上記1つ以上の個々のメッセージに関する外部リソースに対して少なくとも1回の外部機能呼び出しを行うことによって上記メッセージを処理するように構成されており、

上記一次サーバは、

i) 上記メッセージ、

ii) 上記メッセージに関する上記外部機能呼び出し結果、および、

iii) 上記メッセージを処理するためのシーケンス、

を、上記バックアップサーバに複製するように構成されており、

上記バックアップサーバは、上記シーケンスに従って上記一次サーバから受信される上記複製された外部機能呼び出し結果を用いて、上記メッセージを処理するように構成されており、上記サーバの各々による上記メッセージ処理の間、上記共有リソースが、上記一次サーバおよび上記バックアップサーバの両方において実質的に同一であるように構成されている、フェイルオーバのためのシステム。

【0100】

< 8 >

ネットワークを介して相互接続された少なくとも2つのサーバのうちのいずれか1つを選択可能であって、選択した方と接続する少なくとも1つのクライアントを備えており、

上記サーバのうちの1つは、上記クライアントに接続されるときに一次サーバに指定され、かつ、上記サーバのうちの残りは、上記クライアントに接続されないときにバックアップサーバに指定されており、

上記少なくとも1つのクライアントが、上記一次サーバにメッセージを送信するように構成されており、

上記サーバは、上記サーバの各々において同一である少なくとも1つのサービスであって、上記個々のサービスに関するサーバが上記一次サーバか上記バックアップサーバのいずれで動作しているかを認識しない少なくとも1つのサービスを用いて、上記メッセージを処理するように構成されており、

上記サーバは、上記サーバが上記一次サーバであるか、または、上記サーバが上記バックアップサーバであるかを示すライブラリを維持するようにさらに構成されており、

各サーバ内部にて提供される上記サービスは、その個々の上記ライブラリを介して少なくとも1回の外部呼び出しを行うように構成されており、

上記一次サーバにおける上記ライブラリは、上記外部呼び出しを完了させて当該外部呼び出しの結果を上記一次サーバにおける上記ライブラリに戻すように構成されており、上記バックアップサーバにおける上記ライブラリは、上記外部呼び出しの結果を上記バックアップサーバにおける上記サービスに戻すように構成されており、

上記一次サーバにおける上記サービスと、上記バックアップサーバにおける上記サービスとが、上記外部呼び出しの結果を用いて、上記各メッセージを処理するようにさらに構成されている、フェイルオーバのためのシステム。

【0101】

< 9 >

上記ライブラリが、1組または異なる複数の組の固有の有効コードとして実施されている、< 8 >に記載のシステム。

【0102】

< 10 >

上記サーバ各々が、上記サービスが上記メッセージの処理結果を保存することができる共有リソースを維持するように構成されている、< 8 >に記載のシステム。

【0103】

< 11 >

上記共有リソースが、個々のサーバにおけるランダムアクセスメモリ内に維持されてい

10

20

30

40

50

る、 < 10 > に記載のシステム。

【 0 1 0 4 】

< 1 2 >

上記外部呼び出しがタイムスタンプ要求である、 < 8 > に記載のシステム。

【 0 1 0 5 】

< 1 3 >

上記システムが電子取引システムの一部であって、メッセージは、セキュリティが買い注文または売り注文であるかのメッセージであって、上記外部呼び出しが、上記セキュリティの価値に関するマーケットフィールド相場要求である、 < 8 > に記載のシステム。

【 0 1 0 6 】

< 1 4 >

上記少なくとも1つのサービスが、注文受付サービス、注文取消サービス、注文変更サービス、注文適合サービス、予め実行された取引を行うためのサービス、またはクロス取引を行うためのサービス、を含んでいる、 < 8 > に記載のシステム。

【 0 1 0 7 】

< 1 5 >

上記一次サーバにおける上記サービスは、上記外部呼び出しの結果が正しく上記バックアップサーバに送信されたことを上記バックアップサーバが確認した場合、上記メッセージが処理されたことを上記クライアントが認証できるように構成されている、 < 10 > に記載のシステム。

【 0 1 0 8 】

< 1 6 >

上記一次サーバにおける上記サービスは、上記外部呼び出しの結果が正しく上記バックアップサーバに送信されたことを上記バックアップサーバが確認するか否かに関わらず、上記メッセージが処理されたことを上記クライアントが認証できるように構成されている、 < 10 > に記載のシステム。

【 0 1 0 9 】

< 1 7 >

上記一次サーバは、所定の期限内に、上記外部呼び出しの結果が正しく上記バックアップサーバに送信されたことが、上記バックアップサーバが確認しない場合、当該バックアップサーバが障害を起こしたとみなす、 < 16 > に記載のシステム。

【 0 1 1 0 】

< 1 8 >

システムにおけるフェイルオーバーのための方法であって、

ネットワークを介して相互接続された少なくとも2つのサーバのうちのいずれか1つを選択可能であって、選択した方と接続する少なくとも1つのクライアントを備えており、

上記サーバのうちの1つは、上記クライアントに接続されるときに一次サーバに指定され、かつ、上記サーバのうちの残りは、上記クライアントに接続されないときにバックアップサーバに指定されており、

上記少なくとも1つのクライアントが、上記一次サーバにメッセージを送信するように構成されており、

上記方法は、

上記サーバが上記一次サーバであるか、または、上記サーバが上記バックアップサーバであるかを示すライブラリを維持するように上記サーバを構成するステップと、

その個々の上記ライブラリを介して外部呼び出しを行うように上記サービスを構成するステップと、

上記外部呼び出しを完了させて当該外部呼び出しの結果を上記一次サーバにおけるサービスに戻すように、かつ、上記バックアップサーバにおける上記ライブラリに対して上記外部呼び出しの結果を送信するように上記一次サーバにおける上記ライブラリを構成するステップと、

10

20

30

40

50

上記バックアップサーバにおけるサービスに、上記外部呼び出しの結果を送信するように上記バックアップサーバにおける上記ライブラリを構成するステップと、

上記一次サーバにおける上記ライブラリが備える上記外部呼び出しの結果を用いて、上記メッセージを処理するように上記一次サーバにおける上記サービスを構成するステップと、

上記バックアップサーバにおける上記ライブラリが備える上記外部呼び出しの結果を用いて、上記メッセージの処理をするように上記バックアップサーバにおける上記サービスを構成するステップと、を備えており、

上記サービスの各々が、上記サーバ各々において実質的に同一であって、上記サービスの各々が、上記個々のサービスに関する上記サーバが上記一次サーバであるか上記バックアップサーバかのいずれで動作しているか認識しない、フェイルオーバのための方法。

10

【0111】

<19>

選択された少なくとも1つのクライアントと接続可能なネットワークを介して相互接続された少なくとも2つのサーバのうちのいずれか1つのサーバ上で実行可能な、一組のプログラミング命令を保存するコンピュータ読み取り可能な記録媒体であって、

上記サーバのうちの1つが、上記クライアントに接続されるときに一次サーバに指定されており、かつ、上記サーバのうちの残りが、上記クライアントに接続されないときにバックアップサーバに指定されており、上記少なくとも1つのクライアントが、上記一次サーバにメッセージを送信するように構成されており、

20

上記プログラミング命令は、

上記サーバが上記一次サーバであるか、上記サーバが上記バックアップサーバであるかを示すライブラリを維持するように、上記サーバを構成するための命令と、

上記個々のライブラリを介して外部呼び出しを行うように、上記サービスを構成するための命令と、

上記外部呼び出しを完了させて上記一次サーバにおけるサービスに対して上記外部呼び出しの結果を戻すように、かつ、上記バックアップサーバにおける上記ライブラリに対して上記外部呼び出し結果を送信するように、上記一次サーバにおける上記ライブラリを構成するための命令と、

上記バックアップサーバにおける上記サービスに対して、上記外部呼び出しの結果を送信するように上記バックアップサーバにおける上記ライブラリを構成するための命令と、

30

上記バックアップサーバにおける上記ライブラリが備える上記外部呼び出しの結果を用いて、上記メッセージを処理するように上記バックアップサーバにおける上記サービスを構成するための命令と、

上記サーバの各々において実質的に同一であるように上記サービスの各々を構成するための命令であって、上記個々のサービスに関する上記サーバが上記一次サーバか上記バックアップサーバのいずれで動作しているか、上記サービスの各々が認識しないように構成するための命令と、を含むコンピュータ読み取り可能な記録媒体。

【0112】

<20>

一次サーバおよび少なくとも1つのバックアップサーバを備えており、

上記一次サーバが、複数の処理対象 (processing inputs) となる入力を確定するように、かつ、上記入力を処理する上記一次サーバの処理に先立って、上記バックアップサーバに上記処理対象となる入力を送信するように構成されている、フェイルオーバのためのシステム。

40

【0113】

<21>

上記処理対象となる入力は、上記サーバの両方による上記入力の処理の決定性を保証する入力を含む、<20>に記載のシステム。

【0114】

50

< 2 2 >
 上記処理対象となる入力は、外部リソースに対する呼び出し結果を含む、< 2 0 >に記載のシステム。

【 0 1 1 5 】

< 2 3 >

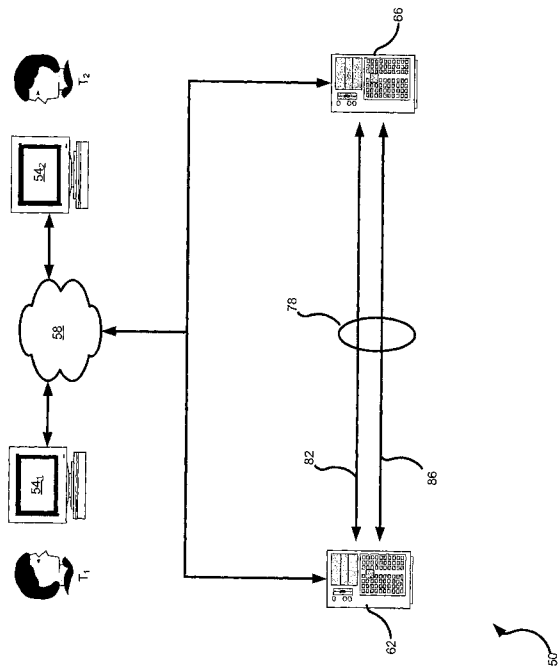
上記処理対象となる入力は、共有リソースをさらに含む、< 2 0 >に記載のシステム。

【 0 1 1 6 】

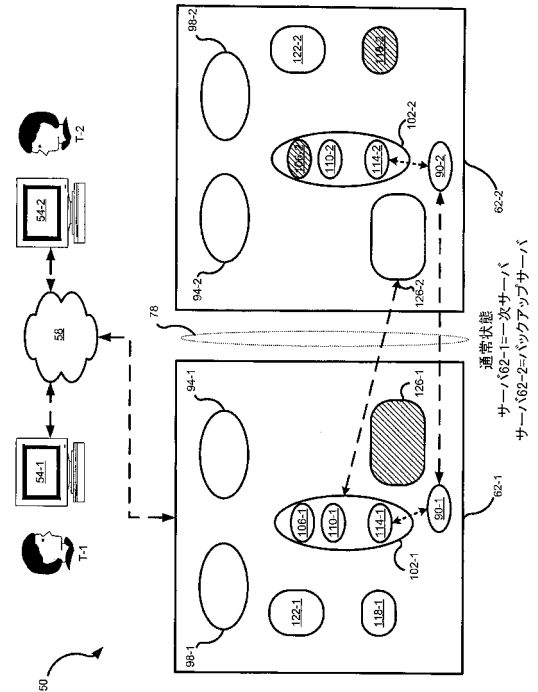
< 2 4 >

上記処理対象となる入力は、上記一次サーバにより生成されたシーケンスナンバをさらに含む、< 2 0 >に記載のシステム。

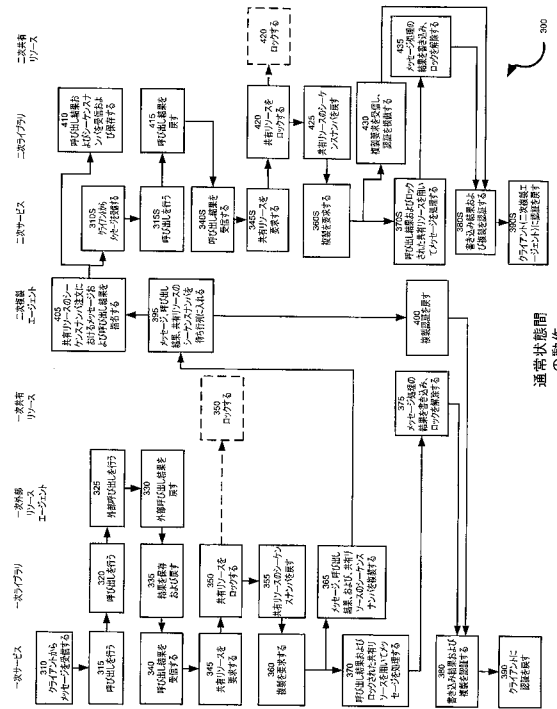
【 図 1 】



【 図 2 】

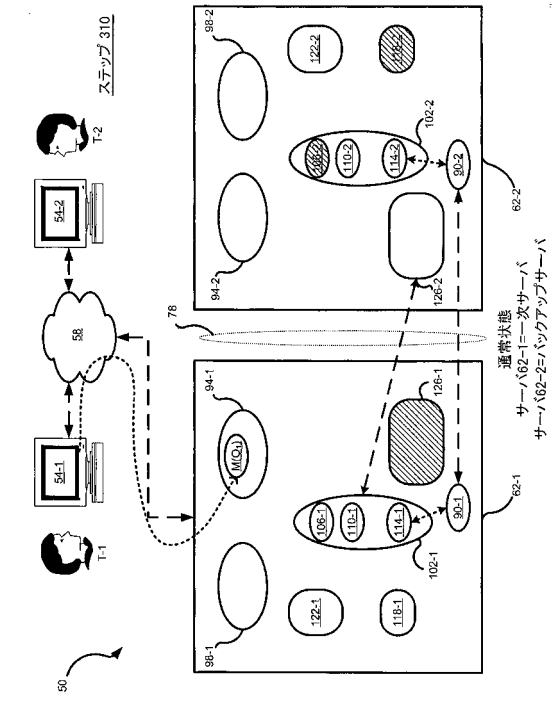


【図3】

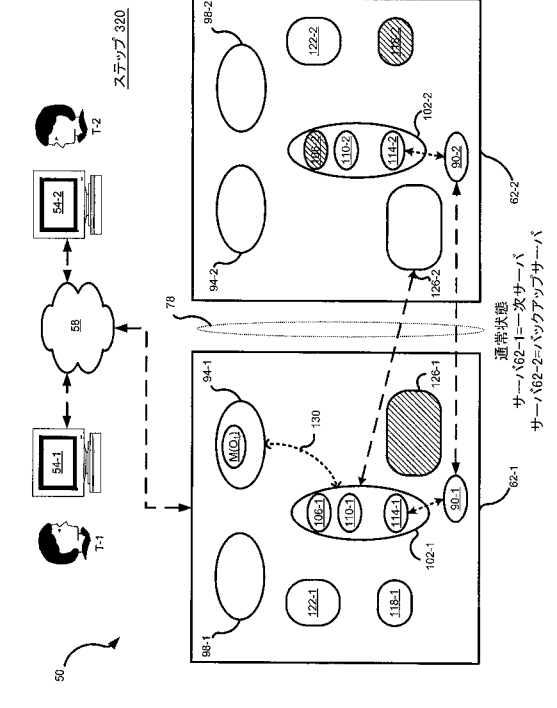


通常の動作

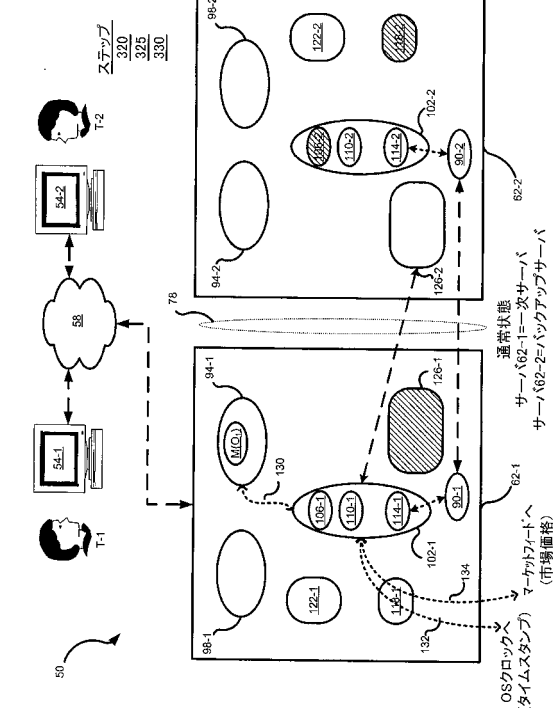
【図4】



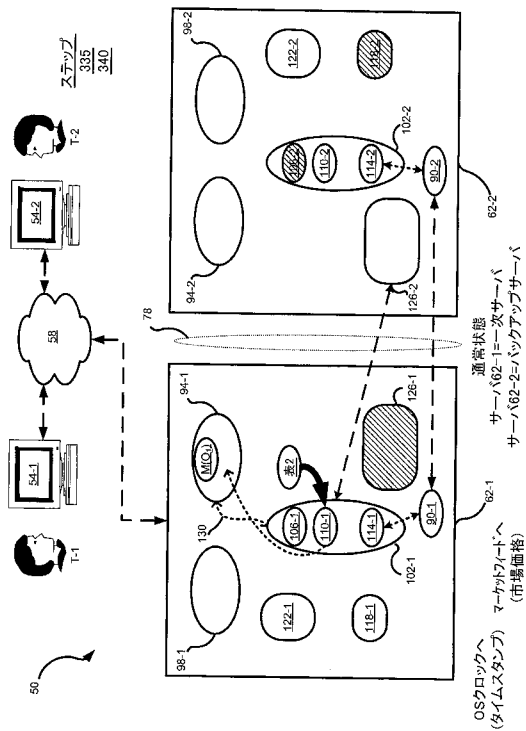
【図5】



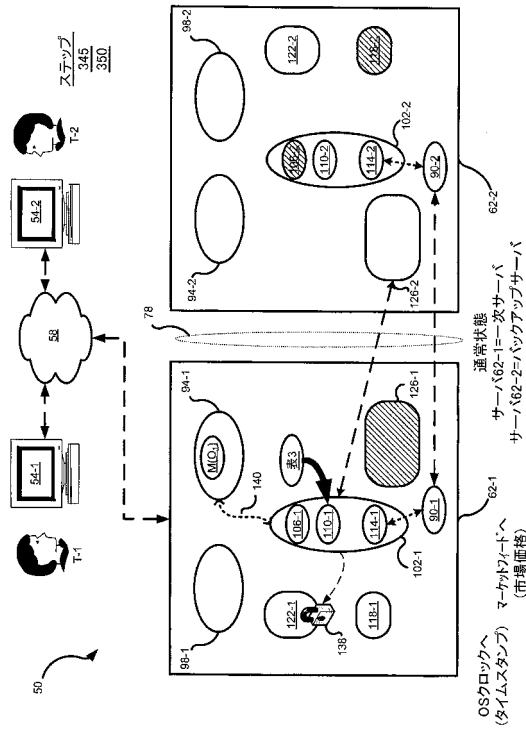
【図6】



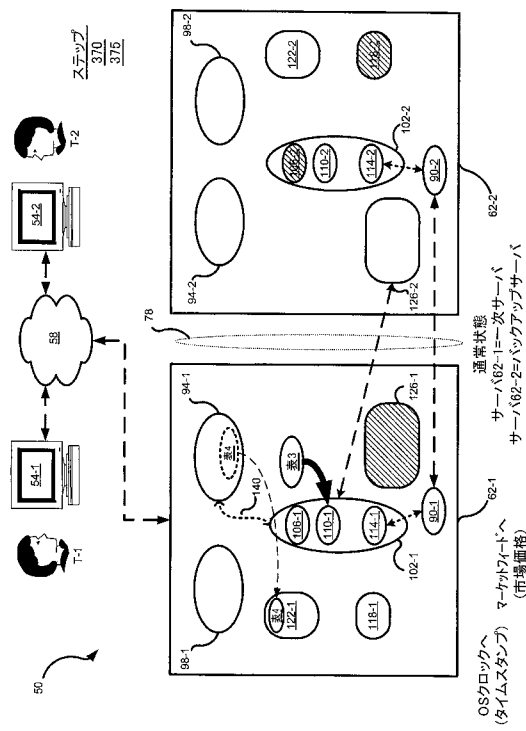
【図7】



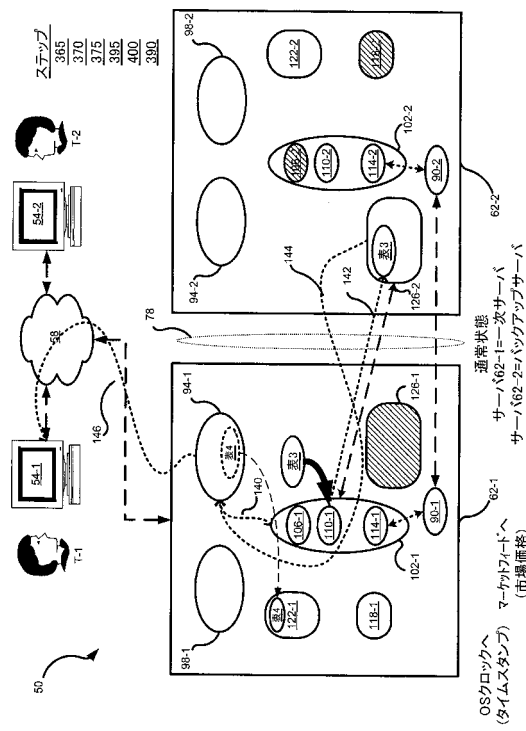
【図8】



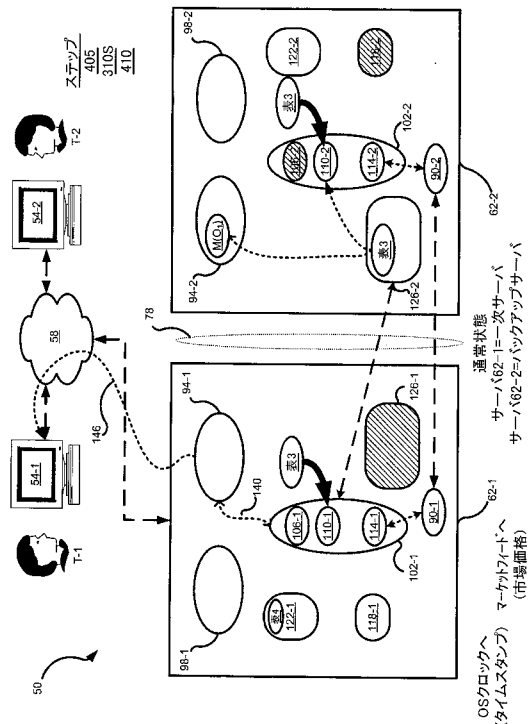
【図9】



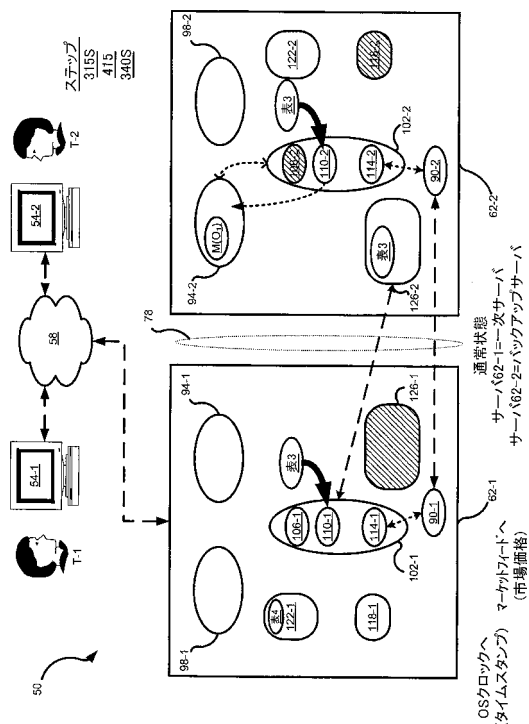
【図10】



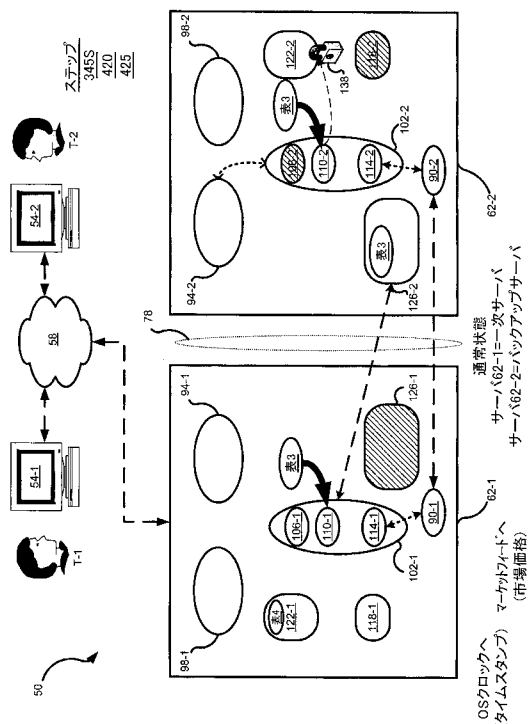
【図 1 1】



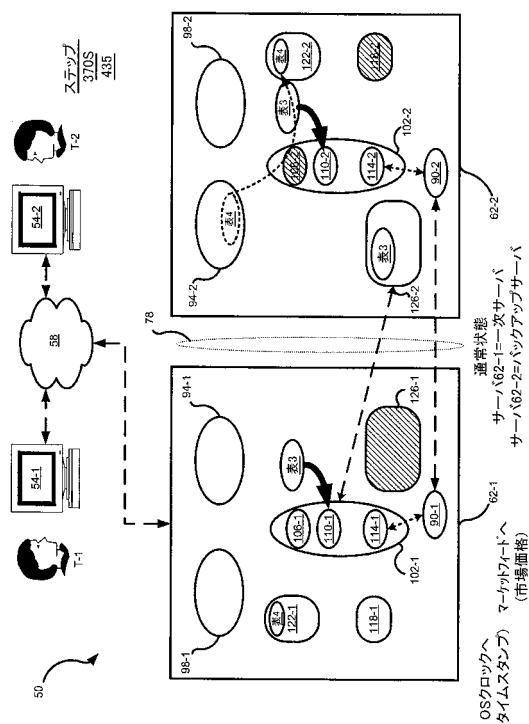
【図 1 2】



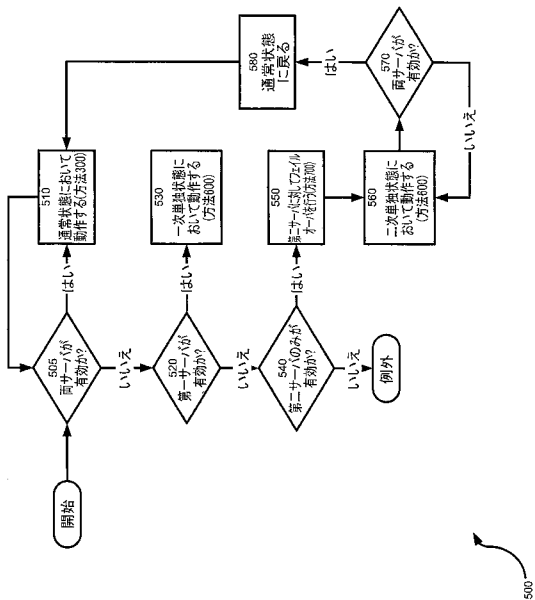
【図 1 3】



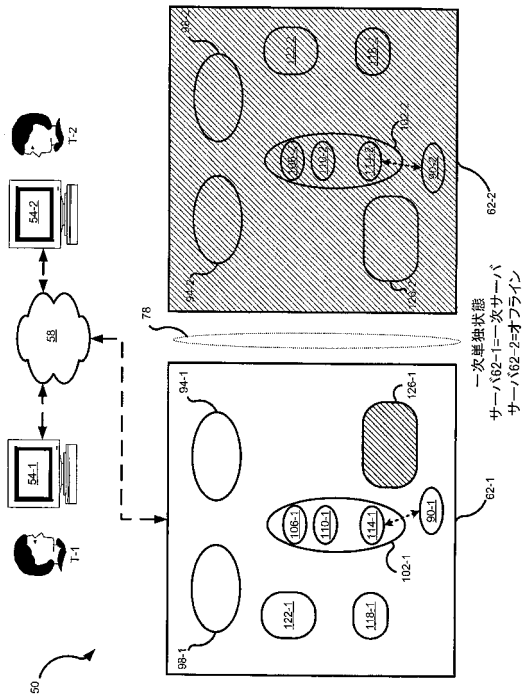
【図 1 4】



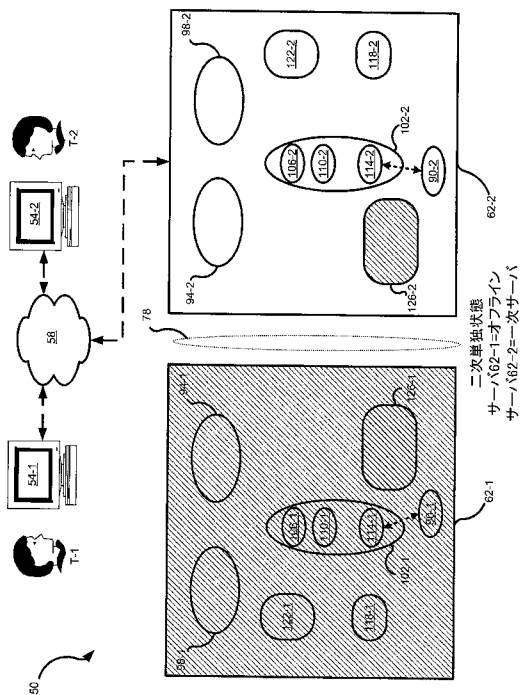
【図15】



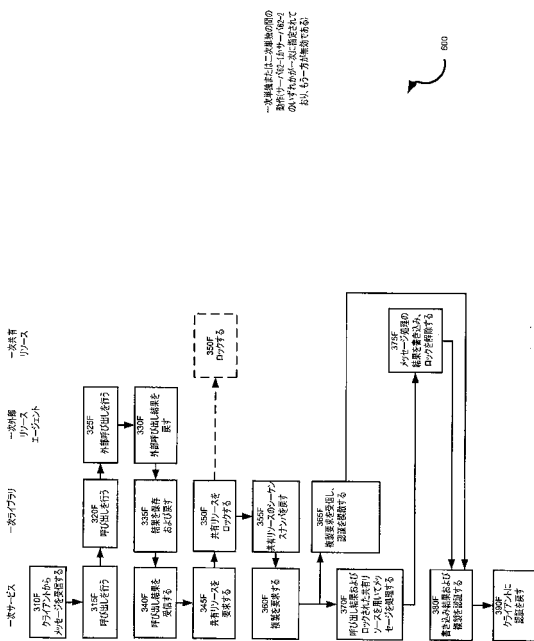
【図16】



【図17】



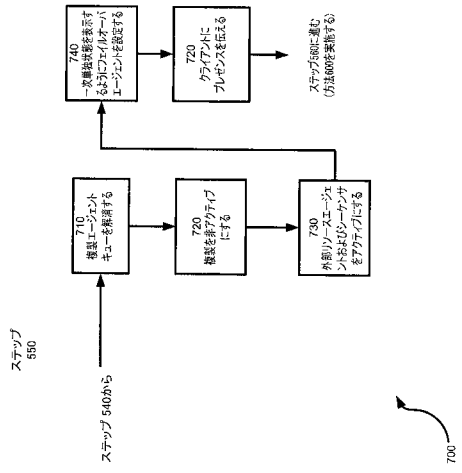
【図18】



一次単独状態は二次単独状態の
状態からサーバ62-1にサーバ62-2
のCPUリソースを一次単独状態
から二次単独状態に移す

600

【図19】



フロントページの続き

- (74)代理人 100107560
弁理士 佐野 惣一郎
- (72)発明者 モロサン, トウドル
カナダ, オンタリオ州 エム2ケイ 1シー1, トロント, シェパード アベニュー イースト
644 #1527
- (72)発明者 アレン, グレゴリー, エイ.
カナダ, オンタリオ州 エル6ジェイ 2ケイ5, オークヴィル, アンソニー ドライブ 464
- (72)発明者 パブレンコ, ヴィクター
カナダ, オンタリオ州 エル7エル 7イー8, バーリントン, オーチャード ロード 2321
- (72)発明者 ラム, ベンソン, ゼ-キット
カナダ, オンタリオ州 エル5エム 7エイ3, ミシソーガ, デル ダン ドライブ 5864

審査官 多賀 実

- (56)参考文献 特開平08-272753(JP, A)
特表平11-502659(JP, A)
特表2002-522845(JP, A)
特開2002-287999(JP, A)
特表平06-504389(JP, A)

- (58)調査した分野(Int.Cl., DB名)
G06F11/16 - 11/20