

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号
特許第4496093号
(P4496093)

(45) 発行日 平成22年7月7日(2010.7.7)

(24) 登録日 平成22年4月16日(2010.4.16)

(51) Int.Cl.	F I
GO6F 13/00 (2006.01)	GO6F 13/00 351N
GO6F 11/16 (2006.01)	GO6F 11/16 310C
GO6F 11/30 (2006.01)	GO6F 11/30 D
	GO6F 11/30 E
	GO6F 11/30 K

請求項の数 19 (全 25 頁)

(21) 出願番号	特願2005-6858 (P2005-6858)	(73) 特許権者	390009531
(22) 出願日	平成17年1月13日 (2005.1.13)		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) 公開番号	特開2005-209191 (P2005-209191A)		INTERNATIONAL BUSINESS MACHINES CORPORATION
(43) 公開日	平成17年8月4日 (2005.8.4)		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
審査請求日	平成17年1月13日 (2005.1.13)	(74) 代理人	100108501
審判番号	不服2008-7416 (P2008-7416/J1)		弁理士 上野 剛史
審判請求日	平成20年3月26日 (2008.3.26)	(74) 代理人	100112690
(31) 優先権主張番号	10/761164		弁理士 太佐 種一
(32) 優先日	平成16年1月20日 (2004.1.20)	(74) 代理人	100091568
(33) 優先権主張国	米国 (US)		弁理士 市位 嘉宏
早期審査対象出願			最終頁に続く

(54) 【発明の名称】 高可用性システムの遠隔エンタープライズ管理

(57) 【特許請求の範囲】

【請求項 1】

高可用性システムの遠隔エンタープライズ管理を可能にするためのシステムであって、ネットワークを介して遠隔エンタープライズ・サーバに通信接続された複数の高可用性システムのうちの特定の高可用性システムを有し、

前記特定の高可用性システムが、

ウェブ・アプリケーションをサポートするためのミドルウェア・スタックを走らせる一次ノードであって、第1のIPアドレスと、該第1のIPアドレスとは異なり、要求を向けるための仮想IPアドレスを割り当てられる前記1次ノードと、

前記1次ノードの前記ミドルウェア・スタックの複数の層をミラーリングするための冗長ミドルウェア・スタックを走らせる2次ノードであって、前記第1のIPアドレス及び前記仮想IPアドレスとは異なる第2のIPアドレスを割り当てられ、前記冗長ミドルウェア・スタックの前記複数の層のうちの第1の選択されたものがアクティブであって、前記冗長ミドルウェア・スタックの前記複数の層のうちの第2の選択されたものが待機状態である、前記2次ノードと、

前記アクティブなミドルウェア・スタック、すなわち前記複数の層の前記第1の選択されたものに対してアクセス可能なデータをもち、前記1次ノードと2次ノードの間で共有されるデータ複製区画と、

前記1次ノードのみが前記データ複製区画にアクセスできるように前記データ複製区画をマウントし、前記高可用性システムの前記1次ノードの状況を監視し、該状況がエラー

を示すことに応答して、前記仮想IPアドレスを前記第1のノードから前記第2のノードに転送し、前記1次ノードに対して電源を切り、前記2次ノードのみがアクセスできるように前記データ複製区画を再マウントし、前記データ複製区画内の前記データに対するアクセスを要求する前記冗長ミドルウェア・スタックの前記複数の層のうちの第2の選択されたものを活動化させる、クラスタ管理コントローラと、

前記クラスタ管理コントローラが、前記1次ノードにおいて前記エラーに反応した時を検出し、前記エラーの時点で、前記高可用性システムの複数のコンポーネントの状態を検出するための監視コントローラであって、前記エラーと、前記複数のコンポーネントの状態を、前記複数の高可用性システムの各々から受け取った別のレポートに基づき前記複数の高可用性システムを管理するようにイネーブルされた前記遠隔エンタープライズ・サーバに報告する、前記監視コントローラとを有し、

10

前記仮想IPアドレスは、前記1次ノードと前記2次ノードのうち、一度に1つのノードのみが利用可能である、

システム。

【請求項2】

高可用性システムの遠隔エンタープライズ管理を可能にするための請求項1に記載のシステムであって、前記特定の高可用性システムが、各々が前記クラスタ管理コントローラによって監視されるJ2EE準拠ミドルウェア・スタックを実装する複数のサーバをもつ前記1次ノード及び前記2次ノードをさらに備えるシステム。

【請求項3】

20

高可用性システムの遠隔エンタープライズ管理を可能にするための請求項1に記載のシステムであって、前記クラスタ管理コントローラが、前記高可用性システムの1次ノードのステータスを検出するためのハートビート・モニタをさらに備えるシステム。

【請求項4】

高可用性システムの遠隔エンタープライズ管理を可能にするための請求項1に記載のシステムであって、前記クラスタ管理コントローラが、前記高可用性システムの前記1次ノードのミドルウェア層によって提供されるサービスのステータスを検出するためのサービス監視デーモンをさらに備えるシステム。

【請求項5】

高可用性システムの遠隔エンタープライズ管理を可能にするための請求項1に記載のシステムであって、前記監視コントローラが、構成要求を前記遠隔エンタープライズ・サーバから受信し、更なるエラーに応答して前記高可用性システムを調整するように前記クラスタ管理コントローラが対応する手法の構成を調整するシステム。

30

【請求項6】

高可用性システムの遠隔エンタープライズ管理を可能にするための請求項1に記載のシステムであって、前記監視コントローラが、構成要求を前記遠隔エンタープライズ・サーバから受信し、前記要求に従って前記高可用性システムのハードウェア構成を調整するシステム。

【請求項7】

高可用性システムの遠隔エンタープライズ管理を可能にするための方法であって、ネットワークを介して遠隔エンタープライズ・サーバに通信接続された、複数の高可用性システムうちの特定の高可用性システムの少なくとも1つのコンポーネントのステータスを監視する段階と、

40

前記特定の高可用性システム内の1次ノードで、ウェブ・アプリケーションをサポートするためのミドルウェア・スタックを走らせる段階であって、前記1次ノードが第1のIPアドレスを割り当てられ、前記1次ノードはさらに、該第1のIPアドレスとは異なり、要求を向けるための仮想IPアドレスを割り当てられる、段階と、

前記特定の高可用性システム内の2次ノードで、前記1次ノードの前記ミドルウェア・スタックの複数の層をミラーリングするための冗長ミドルウェア・スタックを走らせる段階であって、前記2次ノードが前記第1のIPアドレスとは異なる第1のIPアドレスを

50

割り当てられ、前記冗長ミドルウェア・スタックの前記複数の層のうちの第1の選択されたものがアクティブであって、前記冗長ミドルウェア・スタックの前記複数の層のうちの第2の選択されたものが待機状態である、段階と、

前記1次ノードと2次ノードの間で、前記アクティブなミドルウェア・スタックである、前記複数の層の第1の選択されたものに対してアクセス可能なデータをもつデータ複製区画を共有する段階と、

前記1次ノードのみが前記データ複製区画にアクセスできるように前記データ複製区画をマウントし、エラーを示す状況にตอบสนองして、前記1次ノードから前記2次ノードに前記仮想IPアドレスを転送することによって前記特定の高可用性システムを調整するように対応し、前記1次ノードに対して電源を切り、前記2次ノードのみアクセスできるように前記データ複製区画を再マウントし、前記データ複製区画内のデータに対するアクセスを要求する前記冗長ミドルウェア・スタックの前記複数の層のうちの第2の選択されたものを活動化させる段階と、

前記エラーと、前記複数のコンポーネントの状態を、前記複数の高可用性システムの各々から受け取った別のレポートに基づき前記複数の高可用性システムを管理するようにイネーブルされた前記遠隔エンタープライズ・サーバに報告する段階を有し、

前記仮想IPアドレスは、前記1次ノードと前記2次ノードのうち、一度に1つのノードのみが利用可能である、

方法。

【請求項8】

高可用性システムの遠隔エンタープライズ管理を可能にするための請求項7に記載の方法であって、各々が、J2EE準拠ミドルウェア・スタックを実装する複数のサーバを有する、前記1次ノード及び前記2次ノードを監視することによって、前記コンポーネントの状況を監視する段階をさらに含む方法。

【請求項9】

高可用性システムの遠隔エンタープライズ管理を可能にするための請求項7に記載の方法であって、前記高可用性システムの1次ノードのステータスをハートビート・モニタによって監視する段階をさらに含む方法。

【請求項10】

高可用性システムの遠隔エンタープライズ管理を可能にするための請求項7に記載の方法であって、前記高可用性システムの前記1次ノードの前記ミドルウェア層によって提供されるサービスのステータスを、サービス監視デーモンによって検出する段階をさらに含む方法。

【請求項11】

高可用性システムの遠隔エンタープライズ管理を可能にするための請求項7に記載の方法であって、構成要求を前記遠隔エンタープライズ・サーバから受信し、

将来のエラーにตอบสนองして前記高可用性システムを調整するために前記クラスタ管理コントローラが対応する手法の構成を調整する、ことをさらに含む方法。

【請求項12】

高可用性システムの遠隔エンタープライズ管理を可能にするための請求項7に記載の方法であって、

構成要求を前記遠隔エンタープライズ・サーバから受信し、

前記要求に従って前記高可用性システムのハードウェア構成を調整する、

ことをさらに含む方法。

【請求項13】

コンピュータ読み取り可能媒体に常駐し、高可用性システムの遠隔エンタープライズ管理を可能にするためのコンピュータ・プログラムであって、

ネットワークを介して遠隔エンタープライズ・サーバに通信接続された、複数の高可用性システムうちの特定の高可用性システムの少なくとも1つのコンポーネントのステータスを監視する手段と、

10

20

30

40

50

前記特定の高可用性システム内の１次ノードで、ウェブ・アプリケーションをサポートするためのミドルウェア・スタックを走らせる手段であって、前記１次ノードが第１のＩＰアドレスを割り当てられ、前記１次ノードはさらに、該第１のＩＰアドレスとは異なり、要求を向けるための仮想ＩＰアドレスを割り当てられる、手段と、

前記特定の高可用性システム内の２次ノードで、前記１次ノードの前記ミドルウェア・スタックの複数の層をミラーリングするための冗長ミドルウェア・スタックを走らせる手段であって、前記２次ノードが前記第１のＩＰアドレス及び前記仮想ＩＰアドレスとは異なる第２のＩＰアドレスを割り当てられ、前記冗長ミドルウェア・スタックの前記複数の層のうちの第１の選択されたものがアクティブであって、前記冗長ミドルウェア・スタックの前記複数の層のうちの第２の選択されたものが待機状態である、手段と、

前記１次ノードと２次ノードの間で、前記アクティブなミドルウェア・スタックである、前記複数の層の第１の選択されたものに対してアクセス可能なデータをもつデータ複製区画を共有する手段と、

前記１次ノードのみが前記データ複製区画にアクセスできるように前記データ複製区画をマウントし、エラーを示す状況にตอบสนองして、前記１次ノードから前記２次ノードに前記仮想ＩＰアドレスを転送することによって前記特定の高可用性システムを調整するように対応し、前記１次ノードに対して電源を切り、前記２次ノードによるアクセスのため前記データ複製区画を再マウントし、前記データ複製区画内のデータに対するアクセスを要求する前記冗長ミドルウェア・スタックの前記複数の層のうちの第２の選択されたものを活動化させる手段と、

前記エラーの時点で前記高可用性システムの複数のコンポーネントの状態を検出する手段と、

、前記エラーと、前記複数のコンポーネントの状態を、前記複数の高可用性システムの各々から受け取った別のレポートに基づき前記複数の高可用性システムを管理するようにイネーブルされた前記遠隔エンタープライズ・サーバに報告する手段を有し、

前記仮想ＩＰアドレスは、前記１次ノードと前記２次ノードのうち、一度に１つのノードのみが利用可能である、

コンピュータ・プログラム。

【請求項１４】

高可用性システムの遠隔エンタープライズ管理を可能にするための請求項１３に記載のコンピュータ・プログラムであって、各々がＪ２ＥＥ準拠ミドルウェア・スタックを実装する複数のサーバをもつ１次ノード及び２次ノードを監視することによって、前記コンポーネントのステータスを監視するための手段、
をさらに有するコンピュータ・プログラム。

【請求項１５】

高可用性システムの遠隔エンタープライズ管理を可能にするための請求項１３に記載のコンピュータ・プログラムであって、
前記高可用性システムの１次ノードのステータスをハートビート・モニタによって監視するための手段、
をさらに有するコンピュータ・プログラム。

【請求項１６】

高可用性システムの遠隔エンタープライズ管理を可能にするための請求項１３に記載のコンピュータ・プログラムであって、
前記特定の高可用性システムの前記１次ノードの前記ミドルウェア層によって提供されるサービスのステータスを、サービス監視デーモンによって検出するための手段、
をさらに有するコンピュータ・プログラム。

【請求項１７】

高可用性システムの遠隔エンタープライズ管理を可能にするための請求項１３に記載のコンピュータ・プログラムであって、

構成要求を前記遠隔エンタープライズ・サーバから受信するための手段と、

10

20

30

40

50

前記高可用性システムを調整するために前記クラスタ管理コントローラが対応する手法の構成を調整するための手段と、
をさらに有するコンピュータ・プログラム。

【請求項 18】

高可用性システムの遠隔エンタープライズ管理を可能にするための請求項 13 に記載のコンピュータ・プログラムであって、

構成要求を前記遠隔エンタープライズ・サーバから受信するための手段と、

前記要求に従って前記特定の高可用性システムのハードウェア構成を調整するための手段と、

をさらに有するコンピュータ・プログラム。

10

【請求項 19】

複数の高可用性システムを遠隔的に構成するためのシステムであって、

ネットワークを介して遠隔エンタープライズ・サーバに通信接続された複数の高可用性システムのうちの特定の高可用性システムを有し、

前記高可用性システムの各々が、更に、

ウェブ・アプリケーションをサポートするためのミドルウェア・スタックを走らせる 1 次ノードであって、該ミドルウェア・スタックの複数の層が活動的であり、前記 1 次ノードが第 1 の IP アドレスを割り当てられ、さらに、該第 1 の IP アドレスとは異なり、要求を向けるための仮想 IP アドレスを割り当てられる前記 1 次ノードと、

前記 1 次ノードの前記ミドルウェア・スタックの複数の層をミラーリングするための冗長ミドルウェア・スタックを走らせる 2 次ノードであって、前記冗長ミドルウェア・スタックの前記複数の層のうちの第 1 の選択されたものがアクティブであって、前記第 1 の IP アドレス及び前記仮想 IP アドレスとは異なる第 2 の IP アドレスを割り当てられ、前記冗長ミドルウェア・スタックの前記複数の層のうちの第 2 の選択されたものが待機状態である、前記 2 次ノードと、

20

前記 1 次ノードと、前記アクティブ・ミドルウェア・スタックである、前記複数の層の第 1 の選択されたものに対してアクセス可能なデータをもつ 2 次ノードの間で、共有されるデータ複製区画と、

前記 1 次ノードのみが前記データ複製区画にアクセスできるように前記データ複製区画をマウントし、前記高可用性システムの前記 1 次ノードの状況を監視し、該状況がエラーを示すことに応答して、前記仮想 IP アドレスを前記第 1 のノードから前記第 2 のノードに転送し、前記 1 次ノードに対して電源を切り、前記 2 次ノードのみがアクセスできるように前記データ複製区画を再マウントし、前記データ複製区画内の前記データに対するアクセスを要求する前記冗長ミドルウェア・スタックの前記複数の層のうちの第 2 の選択されたものを活動化させる、クラスタ管理コントローラと、

30

前記ステータスがエラーを示すのと同時に、前記複数の高可用性システムの各々の個別のものの複数のコンポーネントについての監視された情報を検出するための監視コントローラと、

前記ネットワークに通信接続された、遠隔エンタープライズ・サーバであって、前記個別の監視コントローラから前記複数の高可用性システムの各々についての前記監視された情報を受け取り、前記監視された情報を解析し、再構成によって調整可能なエラーを示す監視された情報を提出する、前記高可用性システムの各々に対して再構成の要求を送る前記遠隔エンタープライズ・サーバとを有し、

40

前記仮想 IP アドレスは、前記 1 次ノードと前記 2 次ノードのうち、一度に 1 つのノードのみが利用可能である、

システム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、一般に、改善された高可用性 (high availability) クラ

50

スタ管理に関し、特に、高可用性システムの遠隔クラスタ管理に関する。より具体的には、本発明は、エンタープライズ・ネットワークにおける多数の高可用性システムの改善された遠隔監視及び管理に関する。

【背景技術】

【0002】

本発明は、以下の同時継続出願に関連する。

(1) 2004年1月20日出願の米国特許出願第10/761163号

負荷及び需要が常に変動し、各々の顧客の要求を処理することが最重要事項である小売業、銀行、及び他のオンライン・サービスのために、ミッションクリティカルな作動を取扱う高可用性(HA)システムが開発されてきた。一般に、HAシステムは、ネットワーク・システムのコンポーネントの計画停止又は計画外停止のいずれかを原因とするサービスの損失をなくすか又は最小限にするように設計されたシステムである。HAシステムを提供する主要な方法は、サーバのクラスタにグループ化された冗長化ハードウェア及びソフトウェア・コンポーネントによるものである。

10

【0003】

HAシステムにおいては、クラスタの1つのノードに故障が発生したとき、システムが1つのノードによって行われる処理を別のノードに移すようになっているため、冗長性が重要である。例えば、2ノードのHAクラスタにおいては、一方のノードは、典型的には1次ノードとして指定され、他方のノードは、典型的にはバックアップ・ノードとして指定される。一般に、1次ノードは、クラスタが起動されたときに最初にアプリケーションを稼働させる。さらに、一般に、バックアップ・ノードが指定され、1次ノードが故障した場合にそのアプリケーションを稼働させる。HAクラスタ・システムは、典型的には、1次ノードのポーリング(すなわち、ハートビートの検査)を定期的に行って該ノードが引き続きアクティブ状態にあるかどうかを判断するクラスタ管理処理を実施するものである。「ハートビート」が検出されない場合には、クラスタ・マネージャが、ソフトウェア処理をクラスタ内の別のサーバに移動させる。

20

【0004】

HAシステムの重要な特徴は、回復時間にある。一般に、HAシステムの回復時間は、バックアップ・ノードが故障した1次ノードからアプリケーションを引き継ぐのにかかる時間である。顧客が迅速に取引を完了できない場合には、小売業者は重要なビジネスを失うことがあるため、回復時間は、HAシステムを基盤とする販売では特に重要である。回復時間の30秒の遅れでさえ、小売業者の商取引を減少させるものとなる。

30

【0005】

HAシステムの別の重要な特徴は、フェイルオーバーの際にデータをほとんど又はまったく喪失しないことにある。特に、受託データをほとんど又はまったく喪失しないことが重要である。例えば、フェイルオーバーの際に、顧客の注文又は顧客の情報に関する貴重な情報を喪失することは不都合なことである。

【0006】

回復時間を短くし、フェイルオーバーの際にデータをほとんど又はまったく喪失しないようにするためには、まず、HAシステムの構築といった方法でハードウェア及びソフトウェアを組み合わせたことが重要である。しかしながら、HAシステムを起動した後は、該HAシステムの構成を監視して調整し、フェイルオーバー効率及び他のエラーの修正効率の改善を試みるのが重要である。

40

【0007】

ハードウェア及びソフトウェアをHAシステムとして構成するときは、多くの開発者は、新たなハードウェアを必要とすることが多い顧客環境でアプリケーションを制御するために、カスタマイズされたHAソフトウェア・サービスを開発してきた。これらの解決策は、高価なものとなることが多く、多数のプラットフォーム間でのアプリケーションの移植を可能にするオープン・ソース技術の利点を生かすものではない。さらに、サーバ・システムで利用可能な能力がフェイルオーバー効率を自動的に向上させることを期待して、

50

多くの場合、高価なサーバ・システムが選択される。

【 0 0 0 8 】

別の方法として、オープン・ソース開発者は、H Aシステムを実装する際に構成できる機能を用いて、オープン・ソース技術を拡張し続ける。例えば、L i n u xは、プラットフォームに依存しない安価なオペレーティング・システムを提供する。L i n u xの開発者は、他の開発者がオープン・ソース方式で実装することができるオペレーティング・システムに、機能を追加し続ける。「ハートビート」及びディストリビューテッド・リプリケイテッド・ブロック・デバイス (d r b d) といったこれらの機能の幾つかは、L i n u xオペレーティング・システムと共に実装されて、H Aシステムの構成を支援する。

【 発明の開示 】

10

【 発明が解決しようとする課題 】

【 0 0 0 9 】

L i n u xツールは、故障を監視し、H Aシステムで用いられるハードウェアを構成するためのフレームワークを提供するが、付加的な監視及び構成機能についての必要性がある。具体的には、H Aシステムのハードウェア及びソフトウェアの両方における故障、エラー、及び他の非理想的な状態を監視し、オープン・ソースのH Aツールが故障及びエラーを検出したことを監視する方法についての必要性がある。さらに、監視されたシステム・ステータスを遠隔的に収集し、H Aシステムの遠隔的な再構成を容易にする必要性がある。

【 0 0 1 0 】

20

さらに、典型的には、多数のH Aシステムをネットワーク内で組み合わせて、エンタープライズ・システムを形成する。各々のH Aシステムは、例えば、エンタープライズ・システム内部の異なる記憶装置についてのトランザクション要求を処理することができる。遠隔的に、エンタープライズ・システム内部の多数のH Aシステムの監視されたシステム・ステータスを収集し、該システム・ステータスを性能要件と比較し、該エンタープライズ・システム内部の各H Aシステムのハードウェア及びソフトウェアのニーズを突き止めるための方法、システム、及びプログラムについての必要性がある。

【 0 0 1 1 】

さらに、オープン・ソースのオペレーティング・システム・フレームワークを用いてH Aシステムを実装するときは、オープン・ソース対応のミドルウェア層を実装してトランザクション要求を処理することが有利になる。具体的には、(1) 遠隔エンタープライズ・コンソールとインターフェース接続するオープン・ソースに基づくクラスタ管理によって制御され、(2) エンタープライズ・ネットワーク内の多数のH Aシステムを監視し、構成することができる、J a v a (商標) 2 プラットフォーム、エンタープライズ・エディション (J 2 E E) 準拠ミドルウェア・スタックを実装することが有利になる。

30

【 課題を解決するための手段 】

【 0 0 1 2 】

本発明は、改善された高可用性クラスタ管理を提供するものであり、具体的には、オープン・ソース・フレームワークに従って実装される高可用性システムの遠隔クラスタ管理を提供するものである。さらにより具体的には、本発明は、エンタープライズ・ネットワークにおける多数の高可用性システムの改善された遠隔監視及び管理に関するものである。

40

【 0 0 1 3 】

本発明の一態様によると、多数の高可用性システムがエンタープライズ内でネットワーク化され、遠隔エンタープライズ・サーバによって全体的に管理される。各々の高可用性システム内部では、クラスタ管理コントローラが、高可用性システムの特定のコンポーネントのステータスを監視し、そのステータスがエラーを示したときには該高可用性システムを調整するように対応する。さらに、各々の高可用性システムでは、監視コントローラが、クラスタ管理コントローラが特定のコンポーネントのステータスに対応した時を検出し、該高可用性システムの多数のコンポーネントの状態を検出する。次いで、監視コント

50

ローラは、エラー及びコンポーネントの状態を遠隔エンタープライズ・サーバに報告する。遠隔エンタープライズ・サーバは、その報告に基づいて、高可用性システムを管理することができる。

【0014】

具体的には、高可用性サーバは、ハートビート・モニタ及びサービス監視デーモンなどのオープン・ソース機能によって監視されるJ2EE準拠ミドルウェア・スタックを実装する。ハートビート・モニタは、具体的には、ミドルウェア・スタックが常駐する特定のサーバのステータスを検出する。サービス監視デーモンは、具体的には、ミドルウェア・スタックによって提供されるサービスの特定のインスタンスのステータスを検出する。

【0015】

遠隔エンタープライズ・サーバは、構成変更が行われるべきであることを報告から判断し、構成要求を高可用性システムに送信する。次に、監視コントローラは、ハートビート・モニタ又はサービス監視デーモンがエラーを検出し、エラーに対して対応する方法を調整するために、高可用性システムの構成を調整する。さらに、高可用性システム内部の他のハードウェア及びソフトウェア・コンポーネントは、監視コントローラによって再構成することができる。

【0016】

遠隔エンタープライズ・サーバは、各々の高可用性システムに関する監視情報をデータベースに格納することが好ましい。さらに、エンタープライズ・サーバは、監視情報を分析し、どの高可用性システムが性能要件を満足していないかを判断することが好ましい。エンタープライズ・サーバは、ハードウェア及びソフトウェアの変更並びに構成の変更を推奨することができる。さらに、エンタープライズ・サーバは、比較性能を表示し、高可用性システムの実時間表示と、エラーが各々のシステムで検出された時とを提供する。

【0017】

発明の特性と考えられる新規な特徴が、添付の特許請求の範囲に記載される。しかしながら、発明自体並びに好ましい使用モード、そのさらなる目的及び利点は、以下の例示的な実施形態の詳細な説明を添付の図面と併せて読んだときに、最もよく理解されることになるであろう。

【発明を実施するための最良の形態】

【0018】

ここで図面、具体的には図1を参照すると、本方法、システム、及びプログラムを実装することができるシステムの一実施形態が示されている。本発明は、各種のコンピュータ・システム、サーバ・システム、及びエンタープライズ・システムを含む様々なシステムで実施することができる。

【0019】

コンピュータ・システム100は、該コンピュータ・システム100内部の情報を伝達するためのバス122又は他の通信装置と、情報を処理するために該バス122に結合された多数のプロセッサ112a~112nを含む。バス122は、ブリッジ及びアダプタによって接続され、多数のバス・コントローラによってコンピュータ・システム100内部で制御される低待ち時間のバス及びより高待ち時間のバスを含むことが好ましい。

【0020】

プロセッサ112a~112nは、通常作動の際に、ランダム・アクセス・メモリ(RAM)114などの動的記憶装置、及びリード・オンリー・メモリ(ROM)116などの静的記憶装置からアクセス可能なオペレーティング・システム及びアプリケーション・ソフトウェアの制御下でデータを処理する、IBMのPowerPC(商標)プロセッサなどの汎用プロセッサとすることができる。好ましい実施形態においては、多数のソフトウェア層は、プロセッサ112a~112nで実行されるときに、本明細書で説明される図7、図8、図9、図11、図12、図13などのフローチャートに示される動作を行う機械実行可能命令を含む。代替的には、本発明のステップは、該ステップを実施するための配線論理回路を含む特定のハードウェア・コンポーネントか、又は、プログラムされた

10

20

30

40

50

コンピュータ・コンポーネントとカスタム・ハードウェア・コンポーネントとの任意の組み合わせによって実施することができる。

【 0 0 2 1 】

本発明は、本発明に係る処理を行うようにコンピュータ・システム 1 0 0 をプログラミングするのに用いられる機械実行可能命令を格納した機械読み取り可能媒体に含ませて、コンピュータ・プログラムとして提供することができる。ここで用いられる「機械読み取り可能媒体」という用語は、実行命令をプロセッサ 1 1 2 a ~ 1 1 2 n 又はコンピュータ・システム 1 0 0 の他のコンポーネントに与えることに関係する何らかの媒体を含む。こうした媒体は、不揮発性媒体、揮発性媒体、及び伝送媒体を含む多くの形態をとることができるが、これらに限定されるものではない。不揮発性媒体の一般的な形態として、例えば、フロッピー（登録商標）・ディスク、フレキシブル・ディスク、ハード・ディスク、磁気テープ又は他の何らかの磁気媒体、コンパクト・ディスク ROM (C D - R O M) 又は他の何らかの光媒体、パンチ・カード又は孔のパターンを備えた他の何らかの物理媒体、プログラマブル ROM (P R O M)、消去可能 P R O M (E P R O M)、電氣的 E P R O M (E E P R O M)、フラッシュメモリ、他の何らかのメモリ・チップ又はカートリッジ、又は、コンピュータ・システム 1 0 0 が読み取り可能で、命令を格納するのに適した他の何らかの媒体が挙げられる。本発明の実施形態においては、不揮発性媒体の例は、図示されるようなコンピュータ・システム 1 0 0 の内部コンポーネントである大容量記憶装置 1 1 8 であるが、外部装置として構成できることも理解されるであろう。揮発性媒体として、R A M 1 1 4 などの動的メモリが挙げられる。伝送媒体として、バス 1 2 2 を構成するワイヤを含む、同軸ケーブル、銅線、又は光ファイバが挙げられる。伝送媒体は、無線データ通信又は赤外線データ通信の際に生成される波などの音波又は光波の形態をとることもできる。

【 0 0 2 2 】

さらに、本発明は、コンピュータ・プログラムとしてダウンロードすることも可能であり、この場合、プログラム命令は、搬送波又は他の伝搬媒体に統合されたデータ信号として、サーバ 1 4 0 などの遠隔コンピュータから、バス 1 2 2 に結合された通信インターフェース 1 3 2 へのネットワーク・リンク 1 3 4 a ~ 1 3 4 n の 1 つを経由して、必要なコンピュータ・システム 1 0 0 に転送することができる。通信インターフェース 1 3 2 は、例えばローカル・エリア・ネットワーク (L A N)、広域ネットワーク (W A N) に接続することができる多数のネットワーク・リンク 1 3 4 a ~ 1 3 4 n につながる双方向データ通信を提供する。サーバ・システムとして実装されるときは、コンピュータ・システム 1 0 0 は、通常、入力 / 出力コントローラに接続された多数のペリフェラル・コンポーネント・インターコネクト (P C I) バス・ブリッジを経由してアクセス可能な多数の通信インターフェースを含む。このようにして、コンピュータ・システム 1 0 0 は、多数のネットワーク・コンピュータへの接続が可能になる。

【 0 0 2 3 】

ネットワーク環境においては、コンピュータ・システム 1 0 0 は、ネットワーク 1 0 2 を通して他のシステムと通信する。ネットワーク 1 0 2 は、伝送制御プロトコル (T C P) 及びインターネット・プロトコル (I T) などの特定のプロトコルを用いて相互に通信する世界中のネットワーク及びゲートウェイの集まりを参照することができる。ネットワーク 1 0 2 は、デジタル・データ・ストリームを搬送する電気信号、電磁信号、又は光信号を用いる。様々なネットワークを通り、ネットワーク・リンク 1 3 4 a ~ 1 3 4 n と通信インターフェース 1 3 2 とを通過して、デジタル・データをコンピュータ・システム 1 0 0 に又はそこから搬送する信号は、情報を運ぶ搬送波の例示的な形態である。図示されていないが、コンピュータ・システム 1 0 0 には、通信を容易にする多数の周辺コンポーネントを含むこともできる。

【 0 0 2 4 】

コンピュータ・システム 1 0 0 がサーバ・システムとして H A クラスタに実装されるときは、他のサーバ・システムとの局所的な接続を支援するために、付加的なネットワーク

・アダプタを含むことができる。さらに、サーバ・システムとしてH A クラスタに実装されるときは、コンピュータ・システム 1 0 0 は、I B M 社の x S e r i e s (商 標) サーバなどの商用ハードウェア・サーバとして設計することができる。

【 0 0 2 5 】

当業者であれば、図 1 に示されるハードウェアは変更可能であることを認識するであろう。さらに、当業者であれば、図示された例は、本発明に関してアーキテクチャの限定を意味することを意図するものではないことを認識するであろう。

【 0 0 2 6 】

ここで図 2 を参照すると、フェイルオーバーの際にミドルウェアを効率的に移行させるための高可用性クラスタのハードウェア構成のブロック図が示されている。図示されるように、クライアント・システム 2 0 2 及び 2 0 4 が、サービス要求を転送するためにネットワークに接続される。この実施形態においては、クライアント・システム 2 0 2 及び 2 0 4 は、フェイルオーバーの際の回復時間を短くして受託データの損失を最小限にするように構成された高可用性 (H A) システム 2 0 8 から、サービスを要求する。

【 0 0 2 7 】

図示されるように、H A システム 2 0 8 は、1 次ノード 2 1 0 及び 2 次ノード 2 2 0 を含む。以下に説明するように、1 次ノード 2 1 0 及び 2 次ノード 2 2 0 は、実行時には高可用性システムとなる冗長化ハードウェア及びソフトウェアを実装することが好ましい。具体的には、以下に説明するように、1 次ノード 2 1 0 及び 2 次ノード 2 2 0 は、好ましい実施形態においては J 2 E E アプリケーションに対応する冗長化ミドルウェアを実装する。ミドルウェアとは、ウェブ・アプリケーション及びシステムを開発し、統合し、管理するソフトウェアである。以下に説明するように、ミドルウェアは、通信、処理、データの統合と、トランザクション能力管理及びシステム管理の自動化とを可能にする。

【 0 0 2 8 】

具体的には、J a v a (商 標) 2 プラットフォーム、エンタープライズ・エディション (J 2 E E) は、ウェブ・アプリケーションの作成に用いるための再利用可能なコンポーネント・モデルを提供する。J 2 E E は、標準アプリケーション・モデルと、アプリケーションのホストとして動作するための標準プラットフォームと、互換性必要条件と、J 2 E E プラットフォームのオペレーション定義とを定める。このオープン・ソース・モデルの利点は、多数の開発者が付加的なコンポーネント及び構成と共に J 2 E E モデルを実装することが可能で、さらに、すべての J 2 E E アプリケーションが J 2 E E ベースのシステム上で稼動することである。

【 0 0 2 9 】

インターナショナル・ビジネス・マシーニズ・コーポレーション (I B M (商 標)) の開発者は、J 2 E E モデルを実装するソフトウェアを開発した。このソフトウェアは、多くの場合、J 2 E E フレームワークでは規定されない隙間の部分を埋めるものである。例えば、具体的には、I B M (商 標) は、サーバのクラスタに実装されると J 2 E E アプリケーションに対応する J 2 E E 準拠ソフトウェアであるミドルウェア・スタックを開発した。一般に、ミドルウェア・スタックとして、ウェブ・サーバ、データベース・サーバ、及びユニバーサル・インターネット・アプリケーション・サーバが挙げられる。具体的には、このスタックとして、I B M D B 2 (商 標) U D B E n t e r p r i s e E d i t i o n 、 I B M H T T P S e r v e r 、 及び、I B M W e b S p h e r e (商 標) A p p l i c a t i o n S e r v e r などの製品を挙げることができる。

【 0 0 3 0 】

さらに、1 次ノード 2 1 0 及び 2 次ノード 2 2 0 は、H A クラスタの J 2 E E 準拠ミドルウェア・スタック及びハードウェアの故障及びエラーを監視する監視・構成コントローラを実装する。監視・構成コントローラの例として、J 2 E E フレームワークで稼動するソフトウェアを監視するための隙間の部分を埋め、J 2 E E フレームワークが稼動するシステムの構成を容易にする、T i v o l i (商 標) M o n i t o r i n g コントローラを実装することができる。

【0031】

1次ノード210及び2次ノード220は、各々のノードが他のノードのハートビートを迅速に検査することを可能にする信頼性の高い簡単な方法で、接続される。実施形態においては、この接続は、各々のノードにおけるネットワーク・アダプタ間に接続されるクロスケーブル218によって可能になる。具体的には、クロスケーブル218は、ハートビート・データを転送するためのイーサネット（登録商標）接続を可能にすることが好ましい。代替的には、ハートビート・データは、クロスケーブル218が故障した場合に、ネットワーク102を経由して公衆IP接続間で転送することもできる。ハートビートの通信チャネルを1次ノード210と2次ノード220との間に提供するために他のハードウェアを実装しても良いこと、及び、ネットワークを基盤とする接続に加えてシリアル接続を実装しても良いことが理解されるであろう。

10

【0032】

具体的には、ハートビート信号が1次ノード210と2次ノード220との間でクロスケーブル218を通して送られるとき、該ハートビートが機能しなくなった場合には、2次ノード220が、故障の前に1次ノード210によって提供されたサービスを引き継ぐことになる。しかしながら、以下に説明するように、本発明の利点によるとミドルウェア・コンポーネントは、ハートビート故障をさらに分析し、2次ノード220が1次ノード210によって提供されるサービスを引き継ぐ前に、該故障に関する付加的な情報を提供することができる。そのうえ、以下に説明するように、Linuxによるハートビート及びLinuxによらないハートビートは共に、クロスケーブル218を介して監視することができる。

20

【0033】

1次ノード210及び2次ノード220は、データ記憶システム214及び224にアクセスする。有利なことに、ここではd r b dパーティション230として示されるデータ・リプリケータが、1次ノード210と2次ノード220との間で実際に物理的に共有される記憶装置を必要とせずに1次ノード及び2次ノードによってアクセス可能なデータを複製するために、データ記憶装置214及び224の各々のパーティションを含む。本発明の利点によると、d r b dは、フェイルオーバーの際の1次ノード210から2次ノード220へのデータ転送を容易にするために、パーティション上で稼働するように構成される。本発明は、d r b dスクリプトによって管理されるd r b dパーティションに関して説明されるが、他の分散型データ複製システムを実装しても良いことが理解されるであろう。

30

【0034】

無停電電源装置（UPS）212及びUPS222の各々は、それぞれ1次ノード210及び2次ノード220に独立電源を供給する。好ましくは、UPS212と2次ノードとの間、及び、UPS222と1次ノード210との間にも、接続が確立される。一実施形態においては、シリアルケーブル216が1次ノード210とUPS222との間に設けられ、シリアルケーブル226が2次ノード220とUPS212との間に設けられる。しかしながら、他の形式の接続ハードウェアを実装しても良いことが理解されるであろう。

40

【0035】

本発明の利点によると、1次ノード210に故障が検出されたとき、2次ノード220は、事前に1次ノード210に向けられた要求をフェイルオーバー後に受信し始める。1次ノード210上で稼働しているハードウェア、ソフトウェア、又はネットワークの一部のみが故障することがあるため、該1次ノード210がフェイルオーバー後にデータを更新しないようにすることを保証する唯一の方法は、UPS212の電源を切ることである。有利なことに、以下に説明するように、待機ノード220へのフェイルオーバーが検出されたときは、電源を切る命令を該待機ノード220からUPS212に送るために、本明細書でさらに詳細に説明されるSTONITHがクラスタ・マネージャによって実行される。

50

【0036】

ここで図3を参照すると、本発明の方法、システム、及びプログラムに係るクラスタ・マネージャのブロック図が示されている。図示されるように、クラスタ・マネージャ322には、効率的なフェイルオーバーを実行するのに利用される、ハートビート・ツール402と、drbdスクリプト404と、mon406と、STONITH機能408とを含む多数のコンポーネントが組み込まれる。クラスタの他の状況を管理するために、他のコンポーネントをクラスタ・マネージャに組み込んでも良いことが理解されるであろう。さらに、フェイルオーバーを管理するために、付加的なコンポーネントをクラスタ・マネージャ322に組み込んでも良いことが理解されるであろう。

【0037】

ハートビート・ツール402は、J2EE準拠ハードウェア・スタックを用いるHAクラスタ内部のフェイルオーバーを管理するために構成される、Linux用のハートビート・パッケージを含むことが好ましい。具体的には、ハートビート・ツール402は、一般に、クラスタの2つのノード間で「ハートビート」要求を送信することによって作動する。図2で説明されたように、ハートビート要求は、各々のノードにおけるネットワーク・アダプタ間でクロスケーブルを通して送信することができる。サーバ・システムのクラスタ上で稼働するJ2EE準拠ミドルウェア・スタックに使用されるときは、ハートビート・ツール402によって送信されるハートビート要求は、スタックの異なる層全体に配信される。

【0038】

ハートビート要求が戻されない場合には、2次ノードは、1次ノードが故障したものと想定し、該1次ノード上で稼働していたIP、データ、及びサービスを引き継ぐことができる。2次ノードが、1次ノード上で稼働していたIP、データ、及びサービスを引き継ぐときは、ハートビート・ツール402は、待機モードで待機している該2次ノードのコンポーネントを起動し、IPアドレスを該2次ノードのコンポーネントに割り当て、他のフェイルオーバー・タスクを実行する。

【0039】

drbd404は、フェイルオーバーの際のデータ切り替えを改善するための、HAクラスタのデータを管理する関連スクリプトをもつカーネル・モジュールである。この切り替えは、drbd404によって管理されるブロック・デバイスをミラーリングすることにより行われる。drbdは、drbdモジュールを読み込み、HAクラスタの関連システム及び共有記憶装置のIPアドレスを用いて構成するスクリプトである。

【0040】

J2EE準拠ミドルウェア・スタックに使用されるときは、drbd管理ブロック・デバイスは、該ミドルウェア・スタックが稼働することができる記憶装置を提供する。当初は、1次ノードのみがdrbd管理ブロック・デバイスから読み込むか又は書き込むことができるように、クラスタが構成され、drbdパーティションがマウントされる。フェイルオーバーが発生したときは、2次ノードのみがdrbd管理ブロック・デバイスから読み込み/書き込みができるようにdrbdパーティションをマウントするために、drbd404のデータディスク・スクリプトがハートビート・ツール402によって実行される。

【0041】

mon406は、J2EE準拠ミドルウェア・スタック内部の重要なシステム・サービスを監視する監視スクリプトを定期的に行うサービス監視デーモンである。サービスが機能しなくなったか又は異常終了したことが発見された場合には、mon406は、該サービスを再開させて、ミドルウェア・スタックのすべてのコンポーネントが1次サービスの範囲内で引き続き稼働し続けることを保証する。異常終了は、例えば、プログラミング・エラー、又は、RAMを原因とする一時的なクリティカル・リソース制約などの突発的なオペレーティング・システム事象によって発生することがある。具体的には、monがサービスを再開させるときは、該monは、休止サービスとはプロセス識別子(PID

10

20

30

40

50

）は異なるが、同一の仮想IPアドレスをもつサービスの新たなインスタンスを再開させる。

【0042】

STONITH408は、フェイルオーバーの際のデータ保全を保証するために、ハートビート・ツール402によって呼び出される機能である。具体的には、STONITH408は、図2に示されるように、UPS212及び222へのシリアルケーブル構成を含む。ハートビート・ツール402がSTONITH408を呼び出すとき、その呼び出しは、シャットダウンさせるノードを指定する。STONITHは、信号を送信して、要求されたUPSの電源を切る。

【0043】

監視・構成コントローラ410は、HAクラスタ内部のハードウェア及びソフトウェアのステータスを監視するように指定された多数の監視コントローラを含む。本発明の利点によると、HAクラスタの多数のハードウェア及びソフトウェア・コンポーネントに関するステータス情報は、遠隔集中エンタープライズ・コンソールに送られる。好ましくは、監視・構成コントローラ410は、Java(商標)Management Extensions(JMX)を補完して、HAクラスタのハードウェア及びソフトウェア・コンポーネントを監視し、障害及び潜在的な問題を検出し、該クラスタを自動的に危機的な状態から回復させる。一実施形態においては、監視コントローラは、監視情報をTivoli(商標)Enterprise Console(TEC)に送るTivoli(商標)Monitoringによって使用可能にされる。

【0044】

具体的には、ハートビート・ツール402及びmon406が、ノード内部の特定のコンポーネント及びサービスの特定のインスタンスを監視する一方で、監視・構成コントローラ410は、これらのツールによって監視された状態を検出し、フェイルオーバーを開始するようにハートビート・ツール402が起動されるか又はサーバを再始動させるようにmon406が起動されるときには、システムの全ステータスを検出する。このように、監視・構成コントローラ410は、故障、エラー、及び非理想的な状態が発生したときに、ノードの多数のコンポーネントのステータスを収集することによって、オープン・ソース・ツールを補完するものである。

【0045】

本発明の1つの利点として、遠隔集中監視コンソールは、集められた情報を用いて構成変更を判断することができる。具体的には、本発明の利点によると、監視・構成コントローラ410の監視コントローラの各々は、HAクラスタの各々のハードウェア・コンポーネントと、J2EE準拠ミドルウェア・スタックの層の各々とを監視するように構成される。このように、コンソールは、ハードウェア及びミドルウェア層に関する監視情報に基づいて、どのミドルウェア層が、要求をキャッシュするためにより多くのメモリを必要とするか、要求を処理するためにより多くのスレッドを必要とするか、又は他の幾つかの方法で再構成されることを必要とするかを判断することができる。コンソールは、構成変更を監視・構成コントローラ410の構成コントローラに送信することができ、次いで、この構成コントローラは、HAクラスタの構成を調整する。一実施形態においては、構成コントローラは、HAクラスタの構成特性を管理するTivoli(商標)Configuration Managerである。

【0046】

本発明の別の利点として、エンタープライズ・システムにおいては、コンソールは、集められた情報を用いて、どのHAクラスタがハードウェア及びソフトウェアの更新を必要とするかを判断する。例えば、監視情報について、コンソールは、どの記憶装置が故障している様子の交換が必要なハードウェアを有するか、どの記憶装置が最大容量に達して更新が必要なハードウェアを有するか、及び、どの記憶装置が機能していないか又は確実に稼働していないソフトウェアを有するかを判断することができる。

【0047】

本発明のさらに別の利点として、監視・構成コントローラ 410 は、クラスタ・マネージャ 322 内部の他の監視コンポーネントと情報交換を行って、コンソールに送られるステータス情報を収集する。例えば、mon 406 が、監視状態にあるサービスのいずれかの故障を検出したときは、監視・構成コントローラ 410 は、システムにおける故障の全体像を収集することができるように、通知を遠隔集中監視コンソールに送信する。さらに、ハートビート・ツール 402 がシステムの 1 つのノードから別のノードへのフェイルオーバーを開始したときは、監視・構成コントローラ 410 は、ノード故障の統計情報を収集することができるように、通知を遠隔集中監視コンソールに送信する。

【0048】

ここで図 4 を参照すると、本発明の方法、システム、及びプログラムに係る、フェイルオーバー前の HA クラスタにおけるソフトウェア構成の一実施形態のブロック図が示されている。図示されるように、1 次ノード 210 及び 2 次ノード 220 はサーバ・システムのクラスタを表し、各々のクラスタは IP アドレスに割り当てられる。

【0049】

本発明の利点として、クラスタ・マネージャ 322 は、1 次ノード 210 及び 2 次ノード 220 上で稼働して、故障を監視し、サービスを再開し、故障が検出されたときにはフェイルオーバーを制御する。図示されるように、クラスタ・マネージャ 322 は、1 次ノード 210 と 2 次ノード 220 との間で共有される記憶装置上に配置される drbd パーティション 230 を構成する。

【0050】

1 次ノード 210 は、ミドルウェア・スタックのすべてのアクティブなコンポーネント、すなわち、ロード・バランサ 312、HTTP サーバ 314、ウェブ・アプリケーション・サーバ (WAS) 316、メッセージング・コントローラ 318、及びデータベース・サーバ 320 を含む。2 次ノード 220 は、アクティブな HTTP サーバ 334 及び WAS 336 を含むが、ロード・バランサ 332、メッセージング・コントローラ 338、及びデータベース 340 は、待機モードである。

【0051】

ロード・バランサ 312 及び 332 は、クラスタ化することもできる HTTP サーバと WAS サーバとの間で要求の負荷を平準化することが好ましい。好ましくは、ロード・バランサ 312 及び 332 は、サーバの可用性、容量、作業負荷、及び他の基準を用いて、インテリジェント負荷平準化を行う。一実施形態によると、ロード・バランサ 312 及び 332 は、IBM WebSphere (商標) Edge Server を介して実装することができる。

【0052】

図示されるように、ロード・バランサ 312 及び 332 は、Linux ベースのハートビートとは無関係のハートビートを実装することができる。代替的には、Linux ベースのハートビート・モニタ 332 及び 342 が、ロード・バランサ 312 及び 332 のステータスを監視するようにしてもよい。

【0053】

HTTP サーバ 314 及び 334 は、HTTP 要求を受信し、HTTP 要求をそれぞれ WAS 316 及び 336 の間で分散させるように設計されたサーバのクラスタを含むことができる。さらに、HTTP サーバ 314 及び 334 は、サブレット及び EJB についての要求などの他の要求を受信したときに、サブレット・コンテナ及び Enterprise Java (商標) Bean (EJB) コンテナなどのイネーブラを呼び出すことができるようにされる。一実施形態においては、HTTP サーバ 314 及び 334 は、IBM の WebSphere (商標)、特に WebSphere (商標) v. 5.0 に組み込まれた HTTP サーバを介して実装することができる。WebSphere (商標) コンポーネントの多数の複製を 1 か所から制御することができるため、WebSphere (商標) v. 5.0 は有利である。このように、構成変更は、多数のサーバ・システムに配置されたソフトウェア・コンポーネントの多数のインスタンスを生じさせる 1 つの場所

10

20

30

40

50

で行うことができる。

【0054】

本発明の利点として、HTTPサーバ314及び334は、1次ノードが稼働状態になった後にクラスタ・マネージャ322のハートビート・ツールがHTTPサーバをアクティブにする、アクティブ/アクティブ構成で稼働させられる。HTTPサーバ314及び334をアクティブ/アクティブ構成で稼働させることによって、要求負荷を2つの（又は、それ以上の）サーバにわたって分割して、クライアント要求を処理する速度を速めることができる。さらに、HTTPサーバ314及び334をアクティブ/アクティブ構成で稼働させることによって、フェイルオーバー時の起動時間が短くなる。

【0055】

WAS316及び336は、ミッションクリティカルなサービスを顧客に提供するウェブ・アプリケーションに対応可能なサーバのクラスタを含み、具体的には、これらのサーバは、J2EEアプリケーションに対応可能である。一実施形態によると、WAS316及び336は、J2EEアプリケーション及びサービスに対応するために必要なサブレット、EJB、及び他のJ2EEコンポーネントのホストとなる、IBMのWebSphere（商標）5.0によってサポートされるWebSphere（商標）Application Serverである。

【0056】

WAS316は、メッセージング・コントローラ318及びデータベース・サーバ320と情報交換を行って、メッセージング・コントローラ及びデータベースと一体化したアプリケーション・サーバ機能を提供する。本発明の利点として、WAS316及びWAS336は、アクティブ/アクティブ構成で稼働させられる。具体的には、システムを初期化するときは、メッセージング・コントローラ318及びデータベース・サーバ320が利用可能になると、クラスタ・マネージャ322のハートビート・ツールは、WAS336を立ち上げて、アクティブ/アクティブ構成を作り出す。アクティブ/アクティブ構成で稼働させることによって、要求負荷をシステムの多数のクラスタにわたって分割して、クライアント要求を処理する速度を速めることができる。さらに、アクティブ/アクティブ構成で稼働させることによって、フェイルオーバー時の起動時間が短くなる。

【0057】

メッセージング・コントローラ318及び338は、非同期要求を聴取し、これらの要求をローカル・キューに格納して、J2EEベースのシステムと通信するキューを提供するためのコントローラを含む。メッセージング・コントローラ318及び338は、IBM MQSeries（商標）、IBM WebSphere（商標）MQ、又はJava（商標）Messaging Service（JMS）を補完する他のメッセージ・コントローラを実装することができる。

【0058】

本発明の利点として、メッセージング・コントローラ318及び338は、クラスタ・マネージャ322のdrbdが永続的リソースをdrbdパーティション230のメッセージング・キューで管理し、クラスタ・マネージャ322のハートビート・ツールがフェイルオーバー時にメッセージング・コントローラ338の起動を管理する、アクティブ/待機構成で稼働する。

【0059】

データベース・サーバ320及び340は、永続的記憶装置を制御する。データベース・サーバ320及び340は、IBM DB2 UDB Enterprise Edition又は他のリレーショナル・データベース管理システムなどのデータベース制御システムを介して実装することができる。

【0060】

本発明の利点として、データベース・サーバ320及び340は、クラスタ・マネージャ322のdrbdが永続的リソースをdrbdパーティション230のデータベースで管理し、クラスタ・マネージャ322のハートビート・ツールがフェイルオーバー時にデ

10

20

30

40

50

ータベース・サーバ 340 の起動を管理する、アクティブ / 待機構成で稼働する。

【0061】

メッセージング・コントローラ 318 及び 338 並びにデータベース・サーバ 320 及び 340 が、アクティブ / 待機構成で稼働し、最小限のデータ損失で迅速にフェイルオーバーを行うために、メッセージング・コントローラ 318 及びデータベース・サーバ 320 は、キュー及びデータベースを格納するためのルートとして `drbd` パーティション 230 がマウントされる位置を指定するように構成される。さらに、クラスタ・マネージャ 322 は、メッセージング・コントローラ 318 及びデータベース・サーバ 320 の仮想 IP アドレスを用いて、`drbd` 及びハートビート・ツールを構成する。

【0062】

さらに、本発明の利点として、クラスタ・マネージャ 322 の `mon` 機能は、メッセージング・コントローラ 318 及びデータベース・サーバ 320 によって提供されるサービスなどの重要なシステム・サービスを監視する監視スクリプトを定期的に行う。サービスが機能しなくなったか又は異常終了したことが発見された場合には、`mon` は、該サービスを再開させて、ミドルウェア・スタックのすべてのコンポーネントが 1 次サービスの範囲内で引き続き稼働し続けることを保証する。

【0063】

効率的なフェイルオーバーを達成するようにミドルウェアの各々の階層を構成し、クラスタ・マネージャ 322 を通じてミドルウェアの各々の階層を制御する方法は、他のタイプのミドルウェアに応用できることに注目することが重要である。このように、J2EE 互換ミドルウェア・ソフトウェア・スタックから利用可能な機能は拡張し続けるので、各々のミドルウェア・コンポーネントを、アクティブ / アクティブ構成又はアクティブ / 待機構成のいずれかで構成し、クラスタ・マネージャ 322 によって監視し、フェイルオーバーの際に制御することができる。

【0064】

ここで図 5 を参照すると、本発明の方法、システム、及びプログラムに係る、フェイルオーバー後の HA クラスタにおけるソフトウェア構成の一実施形態のブロック図が示されている。図示されるように、フェイルオーバー後は、1 次ノード 210 は、故障ノードとして表示される。2 次ノード 220 は、すべてをアクティブなノードとして引き継ぐことになる。

【0065】

故障が検出され、2 次ノード 220 が 1 次ノード 210 を「休止」として指定するときには、ハードウェア及びソフトウェアの問題が存在する。具体的には、1 次ノード 210 が、必要な時間内にハートビート要求に応答せず、その後間もなく作動可能になる場合がある。先に説明されたように、1 次ノード 210 及び 2 次ノード 220 が共に作動可能になる状況を回避するために、クラスタ・マネージャ 322 のハートビート・ツールは、`STONITH` を呼び出して、1 次ノード 210 への UPS 電源を切ることになる。`STONITH` によって制御できる安価な UPS を実装することにより、データ保全を達成することが可能であり、1 次ノードが実際には休止していないときに生じる可能性のある HA の「スプリット・ブレイン」問題が回避される。

【0066】

次に、フェイルオーバーの際に、ロード・バランサのハートビートは、ロード・バランサ 332 の起動を管理する。起動されると、クラスタ・マネージャ 322 のハートビート・ツールは、1 次ノード 210 の仮想 IP 1 アドレスをロード・バランサ 332 に割り当てる。その結果、仮想 IP アドレスに対する要求は、負荷平準化クラスタの IP アドレスの変更が生じないようにロード・バランサ 332 に転送される。

【0067】

フェイルオーバーの際には HTTP サーバ 334 及び WAS 336 はすでにアクティブなので、クラスタ・マネージャ 322 のハートビート・ツールは、これらのコンポーネントを起動させる必要はない。しかしながら、メッセージング・コントローラ 338 及びデ

10

20

30

40

50

ータベース・サーバ340は待機状態にあるので、クラスタ・マネージャ322のハートビート・ツールは、これらの層のフェイルオーバーを管理する必要がある。まず、ハートビート・ツールは、仮想IP2アドレスを引き継ぐ。次に、ハートビート・ツールは、drbdのデータディスク・サービスを開始して、drbdミラー・パーティションを構成し、マウントする。最後に、ハートビート・ツールは、仮想IP2アドレスが設定され、ミラーdrbdパーティション230上で立ち上がるメッセージ・キュー及びデータベース・インスタンスをもつ、メッセージ・コントローラ338及びデータベース・サーバ340を起動することになる。代替的に、図示されてはいないが、仮想IP2アドレスは1度に1つのノードのみが利用可能であるため、データベース・サーバ340は、待機状態ではなく、アクティブ・モードにしておいてもよい。データベース・サーバ340は、フェイルオーバー時には、要求が着信するまでdrbdパーティション230上のデータにアクセスしようとしないので、データベース・サーバ340には仮想IP2アドレスが設定され、ミラーdrbdパーティション230は要求が着信する前にアクセス可能である。逆に、メッセージング・コントローラ338などの幾つかの層は、起動時に直接データを読み込むものであり、フェイルオーバー前には2次ノード220はdrbdパーティション230上のデータを利用できないため、フェイルオーバー前に2次ノード220上で起動された場合には破壊されることになる。

10

【0068】

ここで図6を参照すると、HAシステムにおけるJ2EE準拠ミドルウェア内部の独立系ソフトウェア・ベンダー(ISV)アプリケーションの一実施例のブロック図が示されている。図示されるように、アクティブなWAS602と、アクティブなIBM MQSeries(商標)サーバ610と、アクティブなIBM DB2サーバ614とが、drbd630とインターフェース接続しているJ2EE準拠ミドルウェア・スタックの1次ノードの一部を表している。参照数字620で示されるように、明細登録又はトランザクション完了が、アクティブなWebSphere(商標)Application Server602で受信される。ISVは、特定のタイプの着信要求を処理するように、サーブレット又はEJBをプログラムすることができる。例えば、参照数字620で示されるように、ルックアップ・サーブレット604は、項目をキャッシュ・レジスタで走査するときに、その価格を調べる価格ルックアップ(PLU)を処理するISVウェブ・アプリケーションである。ルックアップ・サーブレット604は、次に、トランザクション・サーブレット608又は別のサーブレット若しくはEJBといった別のコンポーネントによって非同期的に完了させられる保持トランザクションについての要求を送信する。しかしながら、まず、参照数字622で示されるように、その情報をMQリスナ612に送り、MQキュー632上に置いて、次の着信要求を受信するようにルックアップ・サーブレット604を解放し、トランザクションがMQキュー632を介して順序良く厳密に1度だけ記録されることを保証する。次に、参照数字624で示されるように、MDB606を呼び出して、トランザクションをMQキュー632から取り出し、参照数字626で示されるように該トランザクションをトランザクション・サーブレット608に送信する。最終的に、トランザクション・サーブレット608は、PLUを処理し、参照数字628で示されるように、DB2 634に格納するために結果をIBM DB2コントローラ616に出力する。

20

30

40

【0069】

具体的には、要求がすでにスタックの層間で移行し始めた後でフェイルオーバーが生じたとしても、該スタックは各々のトランザクションが厳密に一度だけ記録されることを保証するため、図6は、フェイルオーバーの際のHAシステムにおけるJ2EE準拠ミドルウェア・スタックの利点を示す。さらに、アクティブ層のMQSeries(商標)サーバ610及びDB2サーバ614は、1次ノードのみがアクセス可能であるが、フェイルオーバーの際には2次ノードによるアクセスのために迅速に再マウントされるdrbdパーティション630とインターフェース接続するため、図6は、フェイルオーバーの際のHAシステムにおけるJ2EE準拠ミドルウェア・スタックの利点を示す。

50

【 0 0 7 0 】

ここで図 7 を参照すると、H A クラスタの J 2 E E 準拠ミドルウェア・スタックに d r b d パーティションを構成するための処理及びプログラムの概略的な論理フローチャートが示されている。図示されるように、処理は、ブロック 7 0 0 で開始し、その後、ブロック 7 0 2 に進む。ブロック 7 0 2 は、d r b d パーティションを構成し、マウントするステップを示す。次のブロック 7 0 4 は、d r b d パーティション上のメッセージ・キュー及びデータベースをアクティブにするステップを示す。その後のブロック 7 0 6 は、フェイルオーバーの際に d r b d パーティションへのアクセスを効率的に移すために、d r b d パーティションにアクセスするメッセージング・サーバ及びデータベース・サーバの仮想 I P アドレスを記録するステップを示し、処理は終了する。

10

【 0 0 7 1 】

ここで図 8 を参照すると、H A クラスタの J 2 E E 準拠ミドルウェア・スタックの構成及びフェイルオーバーを、ハートビート・コントローラを通じて制御するための処理及びプログラムの概略的な論理フローチャートが示されている。示されるように、処理は、ブロック 8 0 0 で開始し、その後、ブロック 8 0 2 に進む。ブロック 8 0 2 は、1 次ノードのミドルウェア層をアクティブにするステップを示す。その後のブロック 8 0 4 は、2 次ノードの H T T P サーバ及び W A S ミドルウェア層をアクティブにするステップを示す。さらに、アクティブ / アクティブ構成で稼動するように指定される他のミドルウェア層がアクティブにされる。その後のブロック 8 0 6 は、2 次ノードから 1 次ノードへのハートビート要求を定期的に起動するステップを示す。ブロック 8 0 8 は、ハートビートの戻りが 2 次ノードによって検出されたかどうかの判断を示す。ハートビートの戻りが検出された場合には、処理は 8 0 6 に戻る。ハートビートの戻りが検出されない場合には、処理はブロック 8 1 0 に移行する。

20

【 0 0 7 2 】

ブロック 8 1 0 は、S T O N I T H を呼び出して、1 次ノードの電源を切るステップを示す。次のブロック 8 1 2 は、仮想 I P アドレスを 1 次ノードから引き継いで、2 次ノードの冗長化コンポーネントに割り当てるステップを示す。その後のブロック 8 1 4 は、データディスク・スクリプトを呼び出して、2 次ノードによるアクセスのために d r b d パーティションを再マウントするステップを示す。次いで、ブロック 8 1 6 は、2 次ノード上の待機ミドルウェア層と、d r b d パーティション上の起動データとをアクティブにするステップを示す。フェイルオーバーの際に、ハートビート・ツール及び他のクラスタ管理サービスによって付加的なステップを実施しても良いことが理解されるであろう。

30

【 0 0 7 3 】

ここで図 9 を参照すると、J 2 E E 準拠ミドルウェア・スタックによって提供されるサービスを監視するために m o n 機能を制御する処理及びプログラムの概略的な論理フローチャートが示されている。図示されるように、処理は、ブロック 9 0 0 で開始し、その後、ブロック 9 0 2 に進む。ブロック 9 0 2 は、ミドルウェアによって提供されるサービスを監視するためのスケジュールを設定するステップを示す。次のブロック 9 0 4 は、予定された時間監視を起動させるかどうかの判断を示す。予定された時間監視が起動されない場合には、処理はブロック 9 0 4 で繰り返される。予定された時間監視が起動された場合には、処理はブロック 9 0 6 に移行する。ブロック 9 0 6 は、予定されたサービスのステータスを監視するステップを示す。その後のブロック 9 0 8 は、サービスが何らかの方法で休止又は故障として検出されたかどうかの判断を示す。サービスが休止として検出されない場合には、処理は終了する。サービスが休止として検出された場合には、処理はブロック 9 1 0 に移行する。ブロック 9 1 0 は、新たな P I D を用いて同一のサービスを再開するステップを示し、処理は終了する。

40

【 0 0 7 4 】

ここで図 1 0 を参照すると、本発明の方法、システム、及びプログラムに従って、J 2 E E ミドルウェア・スタックが稼動し、遠隔エンタープライズ・コンソールによって管理される多数の H A システムを含むエンタープライズ・ネットワークのブロック図が示され

50

ている。図示されるように、H Aシステム1202及びH Aシステム1204が、ネットワーク102を介してH Aシステム1202及び1204を監視し、遠隔的に制御する遠隔エンタープライズ・コンソール1210に通信接続される。多数のH Aシステムを、単一又は多数の遠隔集中コンソールによって監視し、制御できることが理解されるであろう。

【0075】

本発明の利点として、H Aシステム1202及び1204の各々は、小売トランザクション及び他のミッションクリティカルな作業を処理することができる。一実施形態によると、H Aシステム1202及び1204の各々には、図4及び図5に示されたミドルウェア・スタックなどの、J2EEアプリケーションを使用可能にする冗長化J2EE準拠ミドルウェア・スタックによって、高可用性を与えることができる。具体的には、H Aシステム1202及び1204の各々は、図3に示されるように、監視・構成コントローラ410を稼働させるクラスタ・マネージャを含む。

10

【0076】

有利なことに、エラー、故障、又は非理想的な状態がH Aシステム1202及び1204のいずれかで発生したときは、監視・構成コントローラ410は、そのエラー、故障、又は非理想的な状態の時点におけるシステムの状態を検出し、次に、情報を収集して、遠隔エンタープライズ・コンソール1210に対する報告を作成する。本発明の利点として、ハートビート・モニタ又はmon機能が故障又はエラーを検出した場合には、監視・構成コントローラ410が起動して、その故障又はエラーを検出し、故障又はエラーの時点のシステム状態を判断する。

20

【0077】

遠隔エンタープライズ・コンソール1210は、監視情報をデータベースに格納することが好ましい。次に、遠隔エンタープライズ・コンソール1210は、H Aシステム1202及び1204から受信したエラー又は故障の情報を分析し、場合によっては構成変更をH Aシステムに戻して、フェイルオーバーの回避、及びフェイルオーバー効率の改善を試みる第1のコントローラを含むことが好ましい。さらに、遠隔エンタープライズ・コンソール1210は、多数のH Aシステムから受信した故障情報、エラー情報、又は他の情報を比較して、どのシステムが修復及び更新を必要とし、どのシステムが性能要件を満たしていないかを判断する第2のコントローラを含むことができる。遠隔エンタープライズ・コンソール1210は、H Aシステム1202及び1204についての性能統計を収集して、その表示を制御することができる。

30

【0078】

ここで図11を参照すると、本発明の方法、システム、及びプログラムに係るH Aクラスタ・マネージャ内部の監視コントローラを制御するための処理及びプログラムの概略的な論理フローチャートが示されている。図示されるように、処理は、ブロック1000で開始し、その後、ブロック1002に進む。ブロック1002は、故障又はエラーが、H Aシステムのミドルウェア・スタックを監視するハートビート・モニタ、mon、又は他の監視コントローラから検出されたかどうかの判断を示す。故障又はエラーが検出されない場合には、ブロック1002において処理が繰り返される。故障又はエラーが検出された場合には、処理は1004に移行する。ブロック1004は、故障又はエラーの時点で利用可能なシステム情報を収集し、分析するステップを示す。次のブロック1006は、故障又はエラーの情報、及び利用可能なシステム情報を、H Aシステムを監視する遠隔集中コンソールに送信するステップを示し、処理は終了する。

40

【0079】

ここで図12を参照すると、H Aシステムのクラスタ・マネージャを遠隔的に制御して該H Aシステムを再構成するための処理及びプログラムの概略的な論理フローチャートが示されている。図示されるように、処理は、ブロック1100で開始し、その後、ブロック1102に進む。ブロック1102は、ミドルウェア・スタックが稼働しているH Aシステムを再構成するための構成要求を、遠隔エンタープライズ・コンソールから受信した

50

かどうかの判断を示す。要求を受信しない場合には、ブロック 1 1 0 2 において処理が繰り返される。要求を受信した場合、処理はブロック 1 1 0 4 に移行する。ブロック 1 1 0 4 は、ハートビート・モニタを呼び出して、H A システムのフェイルオーバー設定を再構成するステップを示し、処理は終了する。さらに、H A システムのクラスタ・マネージャ内部の他のコントローラを呼び出して、H A システムの他のソフトウェア及びハードウェア構成を調整することができる。

【 0 0 8 0 】

ここで図 1 3 を参照すると、クラスタにおける多数の H A システムを管理するために遠隔エンタープライズ・コンソールを制御する処理及びプログラムの概略的な論理フローチャートが示されている。図示されるように、処理は、ブロック 1 3 0 0 で開始し、その後、ブロック 1 3 0 2 に進む。ブロック 1 3 0 2 は、監視情報が H A システムから受信されたかどうかの判断を示す。監視情報が受信されない場合には、ブロック 1 3 0 2 において処理が繰り返される。監視情報が受信された場合には、処理はブロック 1 3 0 4 に移行する。具体的には、遠隔エンタープライズ・コンソールは、監視情報についての要求を H A システムの各々に定期的送信することができ、各々の H A システムもまた、監視情報を自動的に送信することができる。

【 0 0 8 1 】

ブロック 1 3 0 4 は、多数の H A システムからの監視情報を格納するエンタープライズ・データベースに監視情報を追加するステップを示す。次のブロック 1 3 0 6 は、監視情報が再構成のきっかけとなる場合に、H A システムの再構成を要求するステップを示す。具体的には、遠隔エンタープライズ・コンソールは、特定の形式のエラーが監視情報中に検出されたときに要求されることになる所定の構成を含むことができる。代替的には、システム管理者が、特定の形式のエラーに対する構成のタイプを推奨するようにしても良い。その後のブロック 1 3 0 8 は、監視情報に基づいて、H A システムについての性能統計を再計算するステップを示す。具体的には、性能統計の計算は、ある種の監視エラー又は変動についてのみ起動されることになる。次のブロック 1 3 1 2 は、この H A システムの性能を、エンタープライズ・ネットワークの他の H A システムの性能、及び、該エンタープライズ・ネットワークについて設定された性能要件と比較するステップを示す。次いで、ブロック 1 3 1 4 は、性能の比較結果をチャート及びグラフで表示するステップを示す。例えば、チャートは、H A システムの位置の図形表現を描き、どのシステムが故障したかについての図形標識を提供し、他の H A システムに対する各々の H A システムの性能を示す図形標識を提供することができる。さらに、各々のシステムの実時間性能と、報告されたあらゆるエラーとを表示することができる。次のブロック 1 3 1 6 は、H A システムの脆弱性に対する修正措置を推奨するステップを示し、処理は終了する。例えば、この推奨は、どの H A システムが交換を必要とするか、どの H A システムが更新を必要とするか、及び、どの H A システムがソフトウェア更新又は微調整を必要とするかを示すものとして示すことができる。図 1 3 に示される処理は、多数の高可用性サーバから受信される監視情報に関して実施することができる処理の形式の例であり、本発明の範囲から逸脱することなく、他の同様の分析又は出力を実施できることが理解されるであろう。

【 0 0 8 2 】

本発明は、好ましい実施形態に関して具体的に示され、説明されてきたが、当業者であれば、本発明の精神及び範囲から逸脱することなく、形態及び詳細の様々な変更を行うことができることを理解するであろう。

【図面の簡単な説明】

【 0 0 8 3 】

【図 1】本発明の方法、システム、及びプログラムを実装することができるサーバ・システムを示すブロック図である。

【図 2】フェイルオーバーの際にミドルウェアを効率的に移行させるための高可用性クラスタのハードウェア構成を示すブロック図である。

【図 3】本発明の方法、システム、及びプログラムに係るクラスタ・マネージャを示すブ

10

20

30

40

50

ロック図である。

【図４】本発明の方法、システム、及びプログラムに係るフェイルオーバー前のＨＡクラスタのソフトウェア構成の一実施形態を示すブロック図である。

【図５】本発明の方法、システム、及びプログラムに係るフェイルオーバー後のＨＡクラスタのソフトウェア構成の一実施形態を示すブロック図である。

【図６】ＨＡシステムにおけるＪ２ＥＥ準拠ミドルウェア内部の独立系ソフトウェア・ベンダ・アプリケーションの実装の一実施形態を示すブロック図である。

【図７】ｄｒｂｄパーティションをＨＡクラスタのＪ２ＥＥ準拠ミドルウェア・スタックに構成するための処理及びプログラムを示す概略的なフローチャートである。

【図８】ＨＡクラスタのＪ２ＥＥ準拠ミドルウェア・スタックの構成及びフェイルオーバーを、ハートビート・コントローラを介して制御するための処理及びプログラムを示す概略的なフローチャートである。

10

【図９】Ｊ２ＥＥ準拠ミドルウェア・スタックによって提供されるサービスを監視するｍｏｎ機能を制御するための処理及びプログラムを示す概略的なフローチャートである。

【図１０】本発明の方法、システム、及びプログラムに従って、Ｊ２ＥＥミドルウェア・スタックを実行する多数のＨＡシステムを含むエンタープライズ・ネットワークを示すブロック図である。

【図１１】本発明の方法、システム、及びプログラムに従ってＨＡクラスタ内部の監視コントローラを監視するための処理及びプログラムを示す概略的なフローチャートである。

【図１２】ＨＡシステムのクラスタ・マネージャを遠隔的に制御して、該ＨＡシステムを再構成するための処理及びプログラムを示す概略的なフローチャートである。

20

【図１３】多数のＨＡシステムをクラスタ状態で管理する遠隔エンタープライズ・コンソールを制御するための処理及びプログラムを示す概略的なフローチャートである。

【符号の説明】

【００８４】

１０２：ネットワーク

２０２、２０４：クライアント・システム

２０８：高可用性システム

２１２、２２２：ＵＰＳ

２１４、２２４：データ記憶システム

30

２１６、２２６：シリアルケーブル

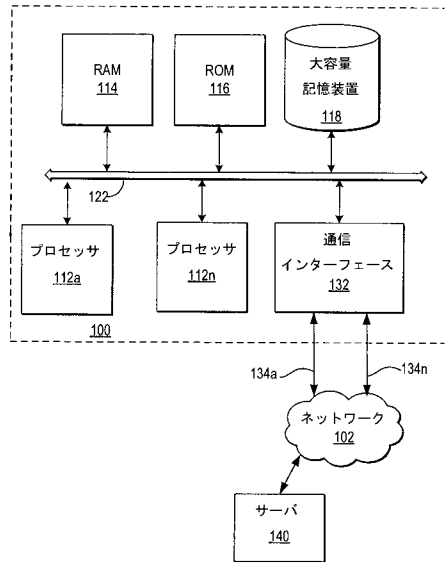
２１８：クロスケーブル

２１０：１次ノード

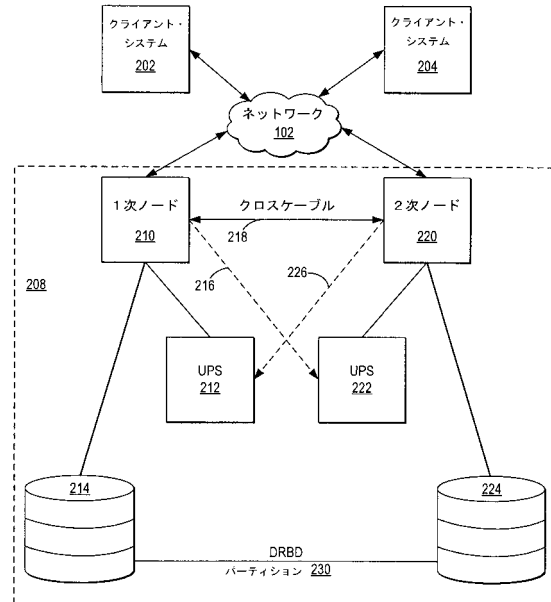
２２０：２次ノード

２３０：ＤＲＢＤパーティション

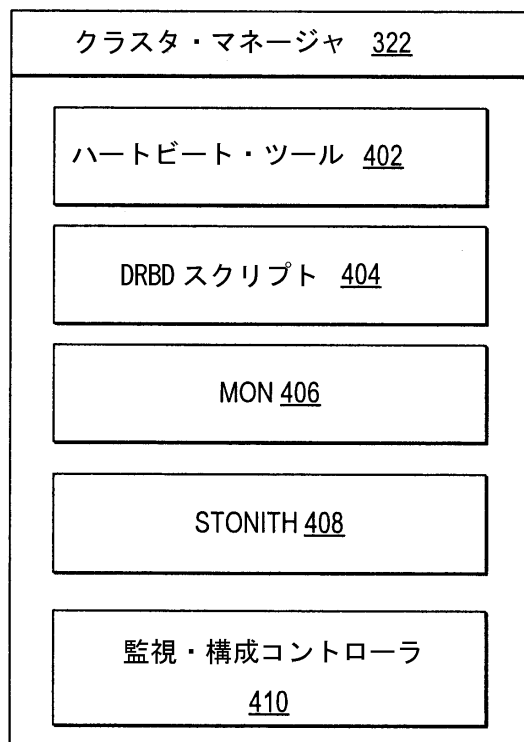
【図 1】



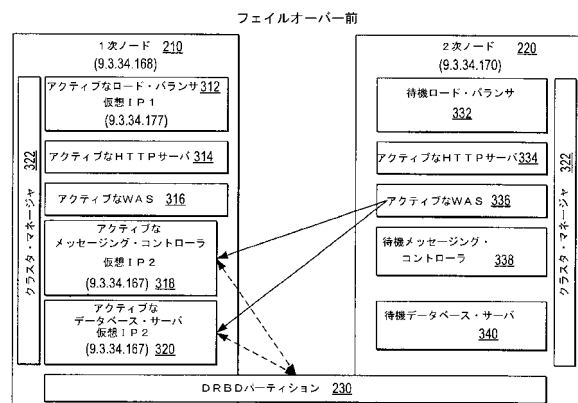
【図 2】



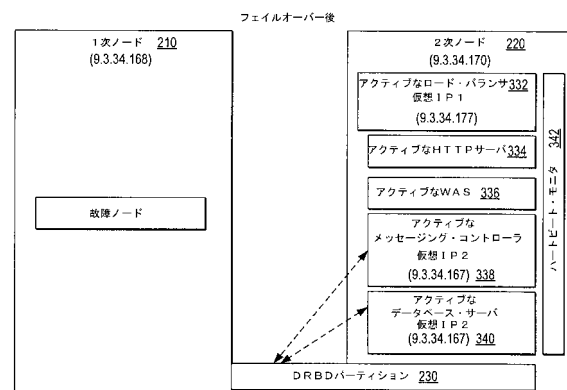
【図 3】



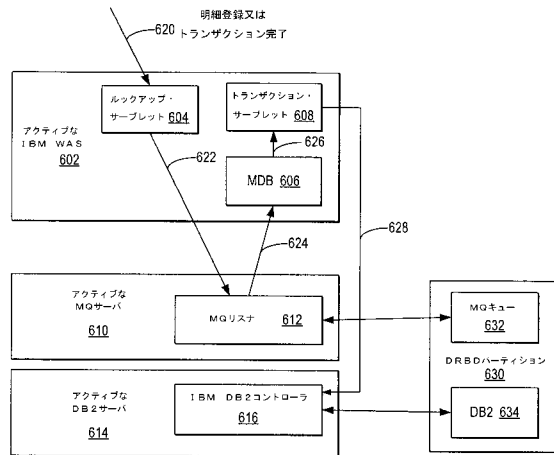
【図 4】



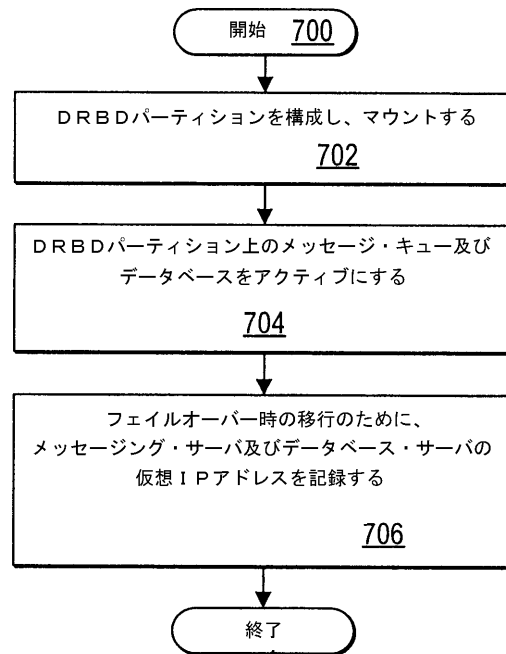
【図 5】



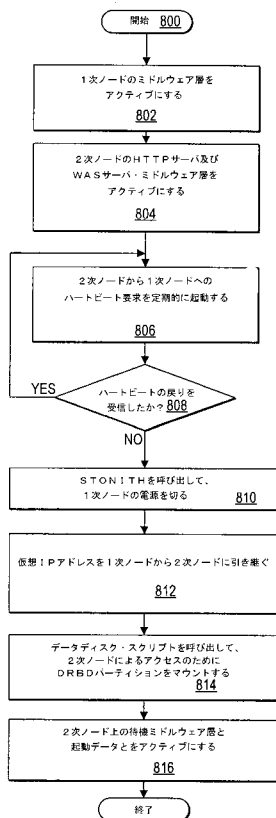
【図 6】



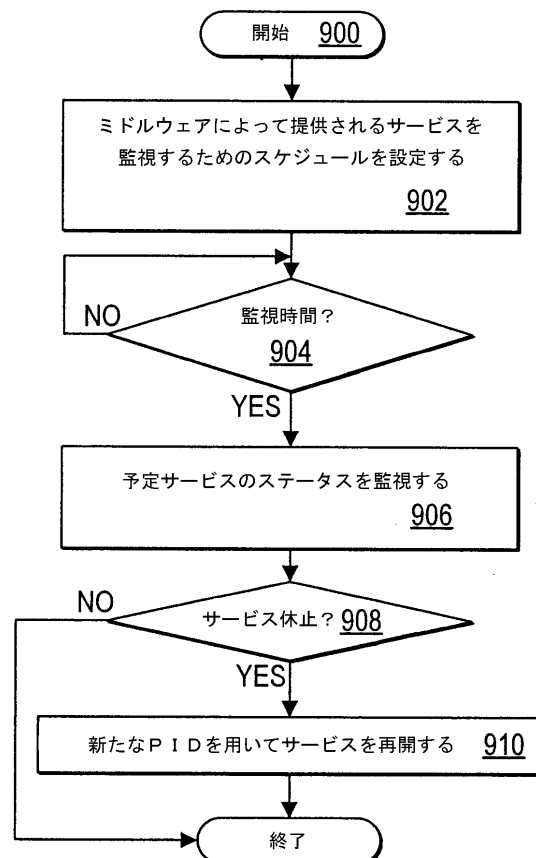
【図 7】



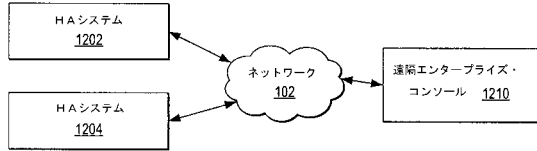
【図 8】



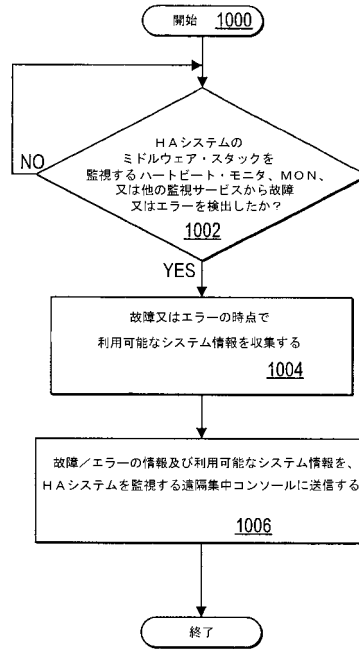
【図 9】



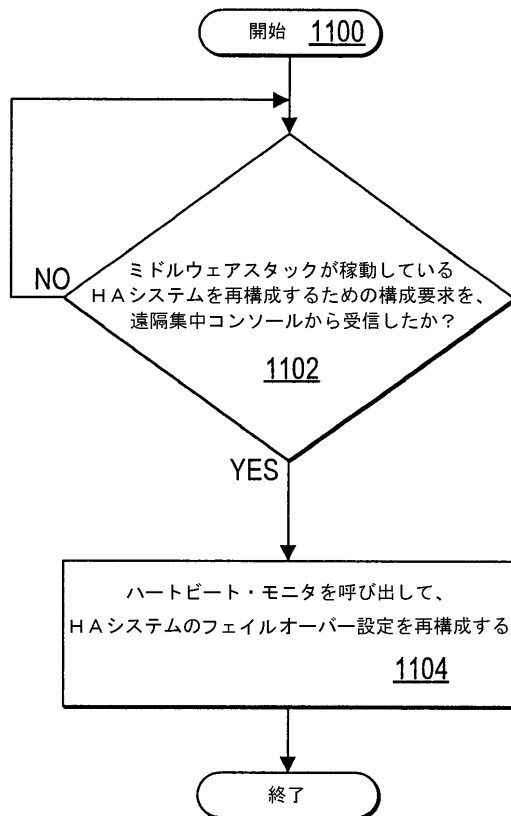
【図 10】



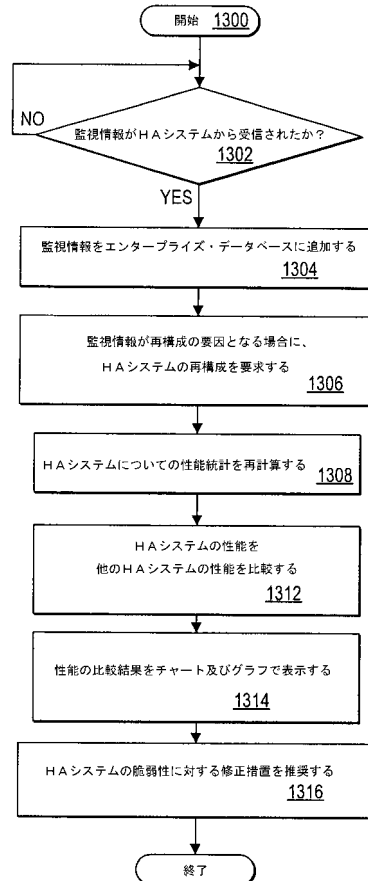
【図 11】



【図 12】



【図 13】



フロントページの続き

(74)代理人 100086243

弁理士 坂口 博

(72)発明者 フランシスコ・デ・ラ・クルーズ

アメリカ合衆国 78717 テキサス州 オースティン モースベリー・ドライブ 14513

(72)発明者 マイケル・エイ・パオリニ

アメリカ合衆国 78750 テキサス州 オースティン ウォレス・コーブ 6407

(72)発明者 ダグラス・スコット・ロザート

アメリカ合衆国 78758 テキサス州 オースティン ホビー・ホース・コート 11901
ナンバー1815

(72)発明者 ラダクリシュナン・セスラマン

アメリカ合衆国 78758 テキサス州 オースティン ストーンホロー・ドライブ 1191
5 ナンバー1418

合議体

審判長 吉岡 浩

審判官 宮司 卓佳

審判官 富吉 伸弥

(56)参考文献 特開平10-254844(JP,A)

特開2003-296141(JP,A)

特開平9-91233(JP,A)

特開平9-311842(JP,A)

特開平6-175868(JP,A)

特開2000-20336(JP,A)

特開2002-312189(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F11/16, G06F13/00