

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3930743号
(P3930743)

(45) 発行日 平成19年6月13日(2007.6.13)

(24) 登録日 平成19年3月16日(2007.3.16)

(51) Int. Cl.	F I
H O 4 L 12/56 (2006.01)	H O 4 L 12/56 G
H O 4 L 12/28 (2006.01)	H O 4 L 12/28 2 O O Z

請求項の数 16 (全 17 頁)

(21) 出願番号	特願2002-20221 (P2002-20221)	(73) 特許権者	398038580
(22) 出願日	平成14年1月29日(2002.1.29)		ヒューレット・パカード・カンパニー
(65) 公開番号	特開2002-319963 (P2002-319963A)		HEWLETT-PACKARD COMPANY
(43) 公開日	平成14年10月31日(2002.10.31)		アメリカ合衆国カリフォルニア州パロアルト
審査請求日	平成16年11月15日(2004.11.15)		ハノーバー・ストリート 3000
(31) 優先権主張番号	09/777,609	(74) 代理人	100081721
(32) 優先日	平成13年2月6日(2001.2.6)		弁理士 岡田 次生
(33) 優先権主張国	米国(US)	(74) 代理人	100105393
			弁理士 伏見 直哉
		(74) 代理人	100111969
			弁理士 平野 ゆかり

最終頁に続く

(54) 【発明の名称】 耐故障性プラットフォームにおいてネットワーク接続を提供する方法

(57) 【特許請求の範囲】

【請求項1】

アクティブ状態のプロセス、スタンバイ状態のプロセス、スタンバイ状態のプロセスをアクティブ状態に移行させるための切り換え能力を有する耐故障性プラットフォームにおいてネットワーク接続を提供する方法であって、

該アクティブ状態のプロセスは、ネットワークアドレスを有し、該方法は、

該スタンバイ状態のプロセスをアクティブ状態へ移行させ、該スタンバイ状態のプロセスの移行中に該ネットワークアドレスを該スタンバイ状態のプロセスに転送し、

該スタンバイ状態のプロセスのアクティブ状態への移行に先立って、

該スタンバイ状態のプロセスに、該アクティブ状態のプロセスのネットワーク接続の状態データを複製するステップと、

該スタンバイ状態のプロセスについて、複製されたデータで更新された対応するスタンバイネットワーク接続を維持するステップと、

アクティブ状態への該スタンバイ状態のプロセスの移行中に、該アクティブ状態のプロセス内の該ネットワーク接続を、該ネットワーク上の該接続を閉じることなく非活動化し、該スタンバイ状態のプロセスにネットワークアドレスが転送されたら、該ネットワークアドレスで該対応するスタンバイネットワーク接続を活動化するステップと、を含み、

該移行したスタンバイ状態のプロセスが該ネットワークにおける接続を再開する必要がないようにする、耐故障性プラットフォームにおいてネットワーク接続を提供する方法。

【請求項2】

前記ネットワーク接続が伝送制御プロトコル（ＴＣＰ）接続、トランザクション処理（ＯＳＩ ＴＰ）接続またはストリーム制御伝送プロトコル（ＳＣＴＰ）接続のいずれかである請求項 1 に記載の方法。

【請求項 3】

前記アクティブプロセスの前記接続を監視するステップと、該接続のアイドル状態が識別されたことに応答して前記状態データの複製を起動するステップと、を含む、請求項 1 または 2 に記載の方法。

【請求項 4】

接続の状態データを複製する前記ステップは、前記アクティブプロセスが前記接続を用いることができる間に実行される、請求項 1 から 3 のいずれか 1 項に記載の方法。

10

【請求項 5】

それぞれ異なるネットワークアドレスを有するアクティブプロセスおよびスタンバイプロセスの複数の対を維持するステップを含む、請求項 1 から 4 のいずれかに記載の方法。

【請求項 6】

前記ネットワークアドレスはＩＰアドレスである、請求項 1 から 5 のいずれかに記載の方法。

【請求項 7】

前記ネットワーク接続がＴＣＰ接続である場合に、状態データを複製するステップは、ＴＣＰ状態情報を複製することを含む請求項 2 から 6 のいずれかに記載の方法。

【請求項 8】

20

アクティブ状態のプロセス、スタンバイ状態のプロセス、および該スタンバイ状態のプロセスをアクティブ状態に移行させるための切り換え手段を有する耐故障性プラットフォームであって、

該アクティブ状態のプロセスは、ネットワークアドレスを有し、該切り換え手段は、スタンバイ状態のプロセスのアクティブ状態のプロセスへの移行の一部として、該ネットワークアドレスをスタンバイ状態のプロセスへ転送するように構成され、該プラットフォームは、

該アクティブ状態のプロセスに関連付けられた接続から状態データを抽出する第 1 の接続マネージャと、

該アクティブ状態のプロセスに関連付けられたネットワーク接続の状態データを該スタンバイ状態のプロセスに複製する複製マネージャと、

30

該スタンバイ状態のプロセスについて、複製されたデータで更新された対応するスタンバイネットワーク接続を維持する第 2 の接続マネージャとを備え、

該切り換え手段は、該スタンバイ状態のプロセスのアクティブ状態への移行に先立って、該アクティブ状態のプロセス内の該ネットワーク接続を、該ネットワーク上の該接続を閉じることなく非活動化するように構成され、ネットワークアドレスが該スタンバイ状態のプロセスに転送されたら、該ネットワークアドレスで該対応するスタンバイネットワーク接続を活動化するように構成されており、

該移行したスタンバイ状態のプロセスが該ネットワーク上で接続を再開する必要がないようにするプラットフォーム。

40

【請求項 9】

前記ネットワーク接続が伝送制御プロトコル（ＴＣＰ）接続、トランザクション処理（ＯＳＩ ＴＰ）接続またはストリーム制御伝送プロトコル（ＳＣＴＰ）接続のいずれかである請求項 8 に記載のプラットフォーム。

【請求項 10】

前記複製マネージャは、個別のソフトウェアモジュールの形をとる、請求項 8 または 9 に記載のプラットフォーム。

【請求項 11】

前記第 1 および／または第 2 の接続マネージャは、個別のソフトウェアモジュールの形をとる、請求項 8 から 10 のいずれかに記載のプラットフォーム。

50

【請求項 1 2】

各プロセスは、命令セットを形成するアプリケーションソフトウェア層と、オペレーティングシステムソフトウェア層と、伝送制御プロトコルを適用することができるソフトウェア層とを提供し、

前記システムは、前記アクティブ接続の前記状態データを複製することができるプロセスを提供し、該プロセスは、少なくとも部分的に、前記アクティブプロセスの前記アプリケーション層によって実行される、請求項 8 から 1 1 のいずれかに記載のプラットフォーム。

【請求項 1 3】

それぞれ異なるネットワークアドレスを有する、アクティブプロセスおよびスタンバイプロセスの複数の対を維持することを含む、請求項 8 から 1 2 のいずれかに記載のプラットフォーム。

【請求項 1 4】

前記ネットワークアドレスは IP アドレスである、請求項 8 から 1 3 のいずれかに記載のプラットフォーム。

r

【請求項 1 5】

前記ネットワーク接続が TCP 接続である場合に、前記複製マネージャは、TCP 状態情報を含む状態データを複製するように構成されている請求項 8 から 1 4 のいずれかに記載のプラットフォーム。

【請求項 1 6】

請求項 8 から 1 5 のいずれかに記載の耐故障性プラットフォームであって、

ネットワーク接続は、関連付けられたプロセスが前記ネットワーク接続を閉じるか、または該関連付けられたプロセスが終わる度に、該ネットワーク接続が前記ネットワーク上で閉じられる第 1 の状態と、関連付けられたプロセスが前記ネットワーク接続を閉じるか、または該関連付けられたプロセスが終わるときに、該ネットワーク接続が該ネットワーク上で閉じられない第 2 の状態とをとることができ、

前記プラットフォームは、プログラム制御下で前記接続を前記第 1 の状態と前記第 2 の状態との間で切り換え、接続状態情報を抽出し、該接続状態情報を設定することを可能にするアプリケーションプログラミングインターフェースを含み、

前記接続状態情報は、ネットワークアドレスを前記スタンバイ接続に転送することにより、前記ネットワーク上で接続が再開されることを必要とすることなく、アクティブ状態のプロセスから複製された前記状態情報で更新され維持されたスタンバイネットワーク接続が、移行したスタンバイ状態のプロセスによって用いられることを可能にするような情報である耐故障性プラットフォーム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は耐故障性コンピューティングに関し、より具体的には、たとえばインターネットプロトコル (IP) ネットワークにおいて、耐故障性プラットフォームによって確立されるネットワーク接続を維持する耐故障性コンピューティングに関する。

【0002】

【従来の技術】

インターネットがさらに普及し、一般的になってくると、それに応じて、毎日、終日稼動し続けなければならないミッション・クリティカルなアプリケーションに対応するために高い可用性を提供するインターネット装置の必要性も増す。ネットワーク内の 1 つの故障箇所、すなわちゲートウェイおよびファイアウォールのような構成要素は、高い信頼性をもって構成される必要があり、冗長な装置および種々のタイプのクラスタ化システムで強化されている。しかしながら、耐故障性の TCP/IP ベースのシステムでは、効率的なフェイルオーバー処理が実施されないと、メッセージトラフィックの損失、およびネットワ

10

20

30

40

50

ークセッションの再初期化といった問題が生じる。

【 0 0 0 3 】

たとえば、H . 3 2 3 ゲートキーパーのハードウェア故障に対処するために、同じソフトウェアを動作させる2つの冗長ゲートキーパーを採用することができる。それらのゲートキーパーはいずれも、全く同じように、IPネットワークに接続され、同じデータを受信し、同じデータを生成する。

【 0 0 0 4 】

ゲートキーパーの一方に故障が生じた場合、故障の直前にはいずれのゲートキーパーにおいてもソフトウェア状態が同じであるので、何ら切り換えの問題を生じさせることなく、他方のゲートキーパーが単独で作業を継続する。しかしながら、このような装置はコスト10
がかかり、いずれのシステムにおいてもソフトウェアが全く同じように動作するため、その解決策は完全に満足できるものではない。すなわち、ゲートキーパーの一方に現れたソフトウェア障害は、他方のゲートキーパーにも現れる。

【 0 0 0 5 】

米国特許第6 , 0 7 8 , 9 5 7号は、それぞれが自身の特定の接続を持つ1組のクラスタメンバを備えたクラスタアセンブリを提案する。各クラスタメンバは、別のクラスタメンバが動作不能になっていることを認識する手段と、動作不能となったクラスタメンバによって実施されていたタスクのうちのいくつかを再びバランスするための手段とを備える。動作不能になったクラスタメンバがクラスタマスターである場合には、他のクラスタメンバは直ちに、クラスタマスターのタスクを、別のクラスタメンバに再度割当てて、20

【 0 0 0 6 】

【 発明が解決しようとする課題 】

あるクラスタメンバが故障した場合にその接続を移すために、各クラスタメンバは、他の各クラスタメンバに、そのクラスタメンバが責任を担っていた接続に関する、記憶された不可欠な状態情報を転送する。しかしながら、このシステムは、IPネットワークと直接通信するアプリケーションを提供する耐故障性プロセスの問題には対処しない。したがって、上記特許は、クラスタ要素から他のクラスタ要素に接続が移るとき、接続される端末から見て、そのようなアプリケーションが依然として同じ段階のままであることを保証するための手段を提供しない。

【 0 0 0 7 】

本発明は主に、IPエンドポイントとの接続を有するシステム、およびそのような接続を処理するための方法を提供することを目的とする。それにもかかわらず、ここに記載される技術は、たとえばTCP（伝送制御プロトコル）接続、OSI TP（トランザクション処理）接続、またはSCTP（ストリーム制御伝送プロトコル）を用いるシステムにおける、任意のタイプの通信ネットワークまたは通信プロトコルに適用されることができることは理解されよう。30

【 0 0 0 8 】

【 課題を解決するための手段 】

簡単に言うと、本発明は、アクティブ状態のプロセス、スタンバイ状態のプロセス、およびスタンバイ状態のプロセスをアクティブ状態に移行させるための切り換え能力とを有する耐故障性プラットフォームにおいて、ネットワーク接続を提供する方法を提供する。この方法は、アクティブプロセスからスタンバイプロセスに、アクティブプロセスのネットワーク接続の状態データを複製するステップと、スタンバイプロセスについて、該複製されたデータで更新された対応するスタンバイネットワーク接続を維持するステップと、スタンバイプロセスからアクティブ状態への移行中、ネットワーク上の接続を閉じる（クローズする）ことなく、アクティブシステムのネットワーク接続を非活動化するステップと、ネットワークアドレスをスタンバイプロセスに転送するステップと、該ネットワークアドレスで対応するスタンバイ接続をアクティブ状態にするステップとを含み、それにより、移行したスタンバイプロセスが、ネットワーク上の接続を再開する必要性をなくす。40

【 0 0 0 9 】

このようにしてスタンバイ接続を維持することにより、ネットワーク接続のリモートエンドに対して高いトランスペアレンシー（透過性）を保持しつつ、フェイルオーバーの達成を可能にする。

【 0 0 1 0 】

スタンバイプロセスのネットワークアドレスを活動化するステップは、スタンバイ接続が活動化される前に実行される。アクティブプロセスの接続を監視し、状態データの複製を起動するステップは、その接続のアイドル状態中に実行される。

【 0 0 1 1 】

接続の状態データを複製するステップは、アクティブプロセスがその接続を使用することができる間に実行されることができるようにするのが有利である。

10

【 0 0 1 2 】

第2の態様によれば、本発明は、アクティブ状態のプロセス、スタンバイ状態のプロセス、およびスタンバイ状態のプロセスをアクティブ状態に移行させるための切り換え手段とを有する耐故障性プラットフォームを提供する。このプラットフォームは、アクティブ状態に関連付けられた接続から状態データを抽出する第1の接続マネージャと、アクティブプロセスに関連付けられたネットワーク接続の状態データをスタンバイプロセスに複製する複製マネージャと、スタンバイプロセスについて、該複製されたデータで更新された対応するスタンバイネットワーク接続を維持する第2の接続マネージャとを備える。最後に、切り換え手段が、スタンバイプロセスのアクティブ状態への移行の一部として、ネットワーク上の接続を閉じることなくアクティブシステムのネットワーク接続を非活動化し、ネットワークアドレスをスタンバイプロセスに転送し、該ネットワークアドレスで対応するスタンバイ接続を活動化するように構成され、それにより、移行したスタンバイプロセスが、ネットワーク上で接続を再開する必要性をなくす。

20

【 0 0 1 3 】

好ましい実施形態では、複製マネージャ、および/または第1および/または第2の接続マネージャは、個別のソフトウェアモジュールの形態をとる。

【 0 0 1 4 】

本発明の第3の態様は、上記の汎用タイプの耐故障性プラットフォームを提供し、そのプラットフォームでは、ネットワーク接続は、関連するプロセスがその接続を閉じるか、またはそのプロセスが終わる度に、ネットワーク上でそれらの接続が閉じられる第1の状態と、関連するプロセスがその接続を閉じるか、またはそのプロセスが終了するときに、ネットワーク上でそれらの接続が閉じられない第2の状態とを有する。このプラットフォームは、プログラムの制御下で、接続が第1の状態と第2の状態との間で切り換えられるようにし、さらに、接続状態情報を抽出して、該接続状態情報の設定を可能にするアプリケーションプログラミングインターフェースを有する。接続状態情報は、ネットワークアドレスをそのスタンバイ接続に転送することにより、ネットワーク上で接続を再開することを必要とすることなく、アクティブプロセスから複製された状態情報で更新された、維持されるスタンバイネットワーク接続が、移行したスタンバイプロセスによって用いられることを可能にするような情報である。

30

【 0 0 1 5 】

【 発明の実施の形態 】

本発明の他の特徴、目的および利点は、本発明の好ましい実施形態の以下に記載される説明を通して、および図面を通して、当業者には明らかになるであろう。

40

【 0 0 1 6 】

図1を参照すると、接続110および210を介していずれもIPネットワーク300に接続されることができ、2つのハードウェア装置100および200を含む耐故障性コンピュータシステムが示される。耐故障性システムはたとえば、米国特許第5,978,933号に記載されるタイプの耐故障性プラットフォーム、またはヒューレット・パッカード社によって市販されるOpenCallINプラットフォームに基づくことができる。そのようなシステムでは、故障検出のためのアプリケーション監視および故障に起因する

50

アクションは、構成要素、詳細には、当業者によく知られている高可用性（HA）コントローラ 101、201 によって管理される。図 2 は、HA コントローラ 101、201 によって用いられる状態機械と、HP OpenCall IN プラットフォームによって用いられる HA プロセスとを示す。

【0017】

図 2 では、図示される状態のうちのいくつか、すなわちブート状態 300 と、同期状態 310 と、活動化状態 320 と、停止状態 330 は、過渡（非常駐）状態である。そのプロセスは、最終的な安定状態に達する前の中間ステップとしてのみ、それらの状態を通過する。安定状態は、アクティブ状態 340、ホットスタンバイ状態 350、およびコールドスタンバイ状態 360 である。そのプロセスは、任意の状態からダウンすることがあるが、簡略化するために、ダウンへの状態遷移は図 2 には示されない。1つの装置上の OpenCall 耐故障性コントローラは、ピア装置 200 上の対応するプロセスの状態を考慮に入れて、装置 100 上で実行される HA プロセスの状態遷移を調整する。

10

【0018】

図 1 の装置 100 および 200 は、ローカルエリアネットワーク（LAN）400 によって互いにリンクされ、それぞれ、その接続に最も近い層からその接続から最も離れた層への少なく 4 つのソフトウェア層、120～150 および 220～250、すなわちインターネットプロトコル（IP）層 120、220、伝送制御プロトコル（TCP）層 130、230、オペレーティングシステム層 140、240、およびアプリケーション層 150、250 を含む。

20

【0019】

アプリケーション層 150、250 は、関連する特定用途のアプリケーションにしたがって種々のアプリケーションレベルのサービスを実行する。たとえば、1つのアプリケーションレベルの機能は、プリペイド通信に対するアプリケーションにおいて、通信時間のプリペイド度数を更新することができる。各装置 100 および 200 において、プロセスはアクティブ状態またはスタンバイ状態のいずれかであることができる。アクティブプロセスは、任意の特定の時間においてアプリケーションサービスを提供するプロセスである。スタンバイプロセスは、アクティブプロセスに障害が生じた場合に引き継ぎの役割を果たす。アクティブ装置 100 のアプリケーション 150 はデータを処理し、一方、スタンバイ装置 200 のアプリケーション 250 の状態は、それ自体がよく知られているやり方で、ローカルネットワーク 400 を通して耐故障性コントローラによって更新される。

30

【0020】

本システムでは、システムのうちの 1 つのシステムの TCP 接続のみ、たとえば、装置 100 の TCP 接続 110 のみが任意のある時点でアクティブ状態にある。他のシステム上の TCP 接続、たとえば 210 の TCP 接続は、その接続がデータを全く受信も送信もできないように構成されるという意味で、スタンバイ接続である。したがって、スタンバイ装置 200 は、物理的にはネットワークには接続されるが、インターネットプロトコルネットワーク 300 からは見えない。

【0021】

ここで、2つの装置間でのアプリケーション切り換え中に、確立された TCP 接続を保持するための技術について説明する。用いられる一般的なアプローチは、その接続を同期した状態にしておくために、アクティブ接続とスタンバイ接続との間でデータおよび状態情報を伝達するための複製機構をアクティブアプリケーション 150 に設けることである。言い換えると、アプリケーション 150 は、TCP コンテキストを、アクティブ装置 100 からスタンバイ装置 200 に移動させる。

40

【0022】

以下の説明では、リモート側、すなわちエンドポイントが、装置 100 と装置 200 との間の切り換え後にその接続を再開する必要があるなければ、TCP 接続は保持されるものと見なされる。ここで用いられる用語「アクティブ TCP 接続」、ここでは接続 110 は、パケットを送受信することのできる接続のことを指す。用語「スタンバイ TCP 接続」、こ

50

ここでは接続 210 は、活動化されるまでパケットを全く送受信することができない接続のことを指す。詳細には、これは、その接続が開始または終了されるときに、この接続を提供している TCP スタック 130、230 によって、アプリケーションのためにパケットが送信されることはなく、その TCP 接続のリモートエンドに、生存しているパケットが全く送信されないことを意味する。

【0023】

図3は、複製マネージャモジュール160、260および接続マネージャモジュール170、270を用いて、接続管理および複製機能を、アプリケーションコア150、250から分離する好ましい構造を示す。当然、他の実施形態では、接続マネージャ機能および複製マネージャ機能がアプリケーション自体に、またはオペレーティングシステムに組み込まれることができることは理解されよう。

10

【0024】

接続マネージャ170、270は、接続110、210を開始し、構成し、その状態を抽出し、その状態を更新し、かつその接続を終了するために、アプリケーション150、250によって用いられる単一インターフェースである。その際、接続マネージャ170、270は、アプリケーション150、250のコアから接続保持の実施の細部を隠すことができ、ここに記載される接続保持機能を使用することを望む他のアプリケーションによって再使用されることができる。

【0025】

また接続マネージャ170、270は、1つのアプリケーションのすべての保持される接続を管理するのに好ましい場所でもある。例として、接続マネージャは、すべてのアプリケーションコアモジュール150、250の内部の詳細に影響を及ぼすことなく、かつそれを知る必要なく、ある時点におけるすべての接続の終了、非活動化、または活動化を容易にする。

20

【0026】

複製マネージャ160、260は、スタンバイにどのようにデータが送信されるかの詳細と、そのようなデータの肯定応答受信の複雑な手順とを隠すようにして、アプリケーションのための抽象複製サービスを提供することができる。

【0027】

図4は、1つの接続が有することができる種々の状態と、アクティブ側100およびスタンバイ側200の両方において取り得る状態遷移とを示す状態図である。図4に示される状態および遷移は、アプリケーション150の視点からのTCP接続を表す。

30

【0028】

その接続の各状態は、その状態に付随した以下の特性を有することができる。

【0029】

P：ある状態において“P”があるとき、これは、その接続が保持された接続（TCP接続保持拡張のうちのいずれにも対応しない標準的なTCP接続とは対照的に）であることを示す。図4の3つのすべての状態は保持された状態であり、他のタイプの接続は示されていない。

【0030】

A：ある状態において“A”があるとき、これは、その状態がアクティブな（活動中の）接続であることを示す。アクティブ状態はデータ転送を実行することができ、F/NFフラグによって示されるように終了動作中に特定の挙動を持つ。

40

【0031】

S：ある状態において“S”があるとき、これは、その状態がスタンバイ接続であることを示す。接続が閉じられるときであっても、そのような接続上ではパケットは送信されない。

【0032】

F：ある状態において“F”があるとき、これは、TCP接続がローカル的に閉じられるときには必ず、該TCP接続はそのネットワーク上においてTCPピアで終了されること

50

を示す。状態 370 の場合のように、A インジケータと組み合わせられるときには、その接続は標準的な TCP 接続のように動作する。

【0033】

NF：ある状態において“NF”があるとき、これは、アプリケーションによって明示的に要求される際、またはプロセス終了時に、TCP 接続がローカル的に閉じられるときには、該 TCP 接続がピア TCP で終了されないことを示す。そのローカルソケットは、暗黙のうちにパーズされる。このオプションは、リモート TCP ピアが接続の終了を起動する場合には機能しない。その場合には、現在のオプション値が何であろうと、その接続は有効に終了される。状態 390 の場合のように、A インジケータと組み合わせられるときには、その接続は、同期した対応するスタンバイ接続を有するアクティブ接続である。

10

【0034】

アクティブアプリケーション 150 が開始し、スタンバイアプリケーション 250 が開始され同期する前に、そのアプリケーションが開いた（開始した）すべての保持された接続がアクティブ状態（状態 370）になる。アクティブアプリケーション 150 が終わった場合には、TCP は、ネットワーク上のリモート側への接続を閉じる。この挙動によって、TCP スタックは、リモート側への接続を閉じることができる。なぜなら、接続処理を引き継ぐためのスタンバイアプリケーションが存在しないためである。

【0035】

以下により詳細に記載されることになる同期段階中に、スタンバイアプリケーションは、保持されたスタンバイ接続を形成し、アクティブ側から接続状態情報を複製する。同期段階が正常に完了すると、アクティブアプリケーションは、そのすべての保持された接続を状態 390 に移し、保持された TCP 接続が処理終了時に閉じられないようにする。これにより、切り換え後に、スタンバイアプリケーションがその接続処理を引き継ぐことができるようになる。アプリケーションが接続を閉じることを欲した場合には、そのアプリケーションはその接続を状態 370 に戻さなければならない。

20

【0036】

スタンバイ接続は状態 380 にある。そのような接続が活動化されると、その接続は状態 370 に移される。新しいスタンバイシステムが再開され同期すると、そのような接続は状態 390 に移されることことができる。

【0037】

アクティブアプリケーションでは、socket()、connect()、bind()、listen()、accept() コールを含む標準的なソケットコールを用いて、ソケットが形成される。これが図 5(a) に示される。アクティブアプリケーション 150 はアクティブ接続 110 を管理する。2つの装置間で複製される必要がある情報の大部分は、両方の装置 100 および 200 上に設けられる TCP スタックモジュール 130、230 によって保持される。

30

【0038】

その後、アクティブ側の接続に生じている変化は、その接続の持続時間にわたって、スタンバイ側に繰返し複製される。複製は、パケットが受信または送信される度に行われる必要はない。むしろ、その状態が安定している接続持続時間中の選択された時点においてのみ、この情報を複製することがより効率的である。その接続によってトラフィックが処理されていない（すなわち、その接続がアイドル状態である）とき、接続状態は安定している。待ち状態のアウトバウンドデータが存在せず、アプリケーションによって読み取られるのを待っている受信データが存在しない場合には、TCP 接続 110 はアイドル状態であると見なされる。そのようなアイドル接続は、ここに記載される技術を用いて保持されることができる。

40

【0039】

アプリケーション 150 が、その接続状態が安定している（すなわち、その接続がアイドル状態である）と判断したとき、該アプリケーションは、TCP スタック 130 から各接続の状態を抽出し、それをスタンバイアプリケーション 250 に送信することにより、対

50

応するスタンバイ接続に接続 1 1 0 を複製することを決定することができる。スタンバイアプリケーション 2 5 0 はソケットを形成し、それをスタンバイソケットに構成し、そのソケットを、アクティブアプリケーションから受信された接続状態情報で更新する。最後に、スタンバイアプリケーションは、アクティブアプリケーションに肯定応答 A C K を送信する。このプロセスが図 5 (b) に示される。接続状態データは、アクティブ装置 1 0 0 のアプリケーション層 1 5 0 によって、O S 層 1 4 0 を介して T C P スタック 1 3 0 から取得され、その後、L A N 4 0 0 を介して、スタンバイ装置 2 0 0 のアプリケーション層 2 5 0 に送信される。

【 0 0 4 0 】

この期間中、アクティブアプリケーションは、スタンバイシステムから A C K を受信するのを待つことはなく、接続は単に状態 3 7 0 においてアクティブ状態のままであることが好ましいことに留意されたい。A C K が受信されると、アクティブアプリケーションはソケットを構成し、それを図 4 の状態 3 9 0 に設定することにより、決してネットワーク上の接続を閉じないようにする。スタンバイ 2 0 0 において接続の複製が失敗した場合に、これがアクティブ接続 1 1 0 に影響を及ぼさないことが許容可能とされる。その場合、アプリケーションは、スタンバイ肯定応答を待つて妨害されることがないので、複製が非常に速く行われる。他の状況では、確立された接続 1 1 0 の保持を促進し、その接続がスタンバイシステム 2 0 0 上で保持されることができた場合でもその接続が該アプリケーションによってのみ用いられることを確実にするために、A C K を待つことが好ましい場合がある。

【 0 0 4 1 】

アクティブアプリケーション 1 5 0 が接続 1 1 0 を閉じることを決めたとき、最初に、その接続を状態 3 7 0 に設定することにより、ネットワーク上においてその接続をピアで終了するようにその接続を構成し、その後に該接続を閉じなければならない。その後、この終了動作は、そのスタンバイ接続 2 1 0 を閉じるスタンバイアプリケーション 2 5 0 に複製される。このプロセスが図 5 (c) に示される。

【 0 0 4 2 】

その接続がリモート側によって閉じられる場合には、T C P スタック 1 3 0 は、標準的な接続の場合と同様に動作する。アクティブアプリケーション 1 5 0 は終了動作をスタンバイアプリケーション 2 5 0 に複製し、その後、その接続を閉じる。これが図 5 (d) に示される。

【 0 0 4 3 】

アクティブアプリケーション 1 5 0 が終わる場合には、スタンバイアプリケーション 2 5 0 がアクティブ状態になる前に、該アクティブアプリケーションに割り当てられた I P アドレスがスタンバイホストに移される。装置 1 1 0 は、1 つの I P アドレスのみを有することが可能である。しかしながら、各アクティブアプリケーション 1 5 0 は、そのアプリケーションのプロセスによって用いられる、自身に専用の I P アドレスを有することが好ましい。そのアプリケーションの I P アドレスは、アクティブアプリケーション 1 5 0 を提供するアクティブ装置 1 0 0 上でのみ有効である。切り換え中に、I P アドレスは、古いアクティブ装置 1 0 0 から新しいアクティブ装置 2 0 0 に移される（すなわち、その I P アドレスは古いアクティブ装置上で非活動化され、新しい装置上で活動化される）。その I P アドレスは、任意の時点で 1 つの装置においてのみアクティブ状態である。

【 0 0 4 4 】

1 つの M A C アドレスから別の M A C アドレスに I P アドレスを転送するための技術は当業者にはよく知られており、H P O p e n c a l l 製品は、これを行うための I P アドレスマネージャ構成要素を含む。そのような技術の記載は、たとえば、米国特許第 6 , 0 4 9 , 8 2 5 号に見いだすことができる。

【 0 0 4 5 】

H A プロセス状態機械は、切り換え中に、I P アドレスが装置上でアクティブ状態になった後にのみ、スタンバイプロセスがアクティブ状態になり、I P アドレスが装置上で非活

10

20

30

40

50

動化された後にのみ、アクティブプロセスがスタンバイ状態になるように、IPアドレスの移動と同期する。これは、IPアドレスマネージャとアプリケーションとの間で通信を行うことによって、または任意のスタンバイ接続を活動化する前にアプリケーション200にIPアドレス状態を検査させることによって、確実に実行されることができる。標準的なAPIコールを、IPアドレス状態を検査するために利用することができることは理解されよう。

【0046】

プロセス終了時に、カーネルは、すべてのそのファイル記述子を閉じる。すべての複製されたアクティブ接続が、ピアで接続を終了しないように設定されたので(状態390)、TCPスタック130は、その接続が閉じられていることを該接続のリモート側に示さない。

10

【0047】

同時に、高可用性コントローラによって障害を通知されているスタンバイシステムは、その装置上でIPアドレスが活動化されるまで待機する。その後、その装置は、そのスタンバイ接続を活動化し、そのリスン(listen)接続を開始する。活動化の後、新しいスタンバイシステムが再開され同期するまで、その接続は状態370に留まらなければならない場合があることに留意されたい。これらのプロセスが図6(a)に示されており、好ましい実施形態では、それらのプロセスは、図2に示される活動化状態320中に行われる。

【0048】

スタンバイソケットがアクティブ状態になると、TCP接続は、そのTCPピアと依然として同期しているか否かを判断する。最後の同期点以来アクティブTCP接続がデータを受信していた場合には、スタンバイTCPは同期から外れるであろう。その場合、該TCP接続は閉じられ、スタンバイアプリケーションによって再形成される。

20

【0049】

手動切り換え(すなわち、プロセス終了によって引き起こされない切り換え)の場合、たとえば予防保守の場合、または、ソフトウェアあるいはハードウェアアップグレードを可能にするために、アクティブアプリケーション100はそのすべての接続を非活動化し、スタンバイ状態になる。一方、スタンバイアプリケーション200はすべての接続を活動化し、アクティブ状態になる。この手動切り換えは図6(b)に示されており、好ましい実施形態では、そのプロセスは、図2に示される「停止」状態330において行われる。

30

【0050】

リスン(listen)するために用いられるソケットは、それらが接続の一部を形成しないので保持される必要がないことに留意されたい。アクティブプロセスによって確立されるリスンソケットを処理するために、スタンバイプロセスは2つの方策のうちの1つの採用することができる。第1に、スタンバイプロセスは、アクティブプロセス状態に切り替わるときに、リスンソケットを再形成することができる。これは、ソケット機能シーケンス、socket()、connect()、bind()、listen()、accept()を実行することにより行うことができる。一般に、このアプローチは、最も安全で、かつ最も簡単であると考えられる。しかしながら、代替的に、スタンバイプロセスがリスンソケットを形成し、そのリスンソケットを任意のIPアドレス(IPアドレス=INADDR_ANY)および指定されたポート番号にバインドすることができる。これは、リスンソケットを再形成するステップを節約することになるが、スタンバイプロセスは、依然としてスタンバイモードにあるとき(その装置上でアクティブ状態にあるIPアドレスについて)、接続要求の可能性を制御しなければならないであろう。

40

【0051】

TCP接続110が、切り換えが生じた時点でアイドル状態でない場合には、その接続110は終了され、スタンバイアプリケーションによって再形成される必要がある。

【0052】

以下により詳細に記載されるように、好ましい実施形態では、アプリケーションが接続特性を制御し、アクティブ側の状態を抽出し、スタンバイ側の状態を更新できるようにする

50

ために、拡張されたソケットアプリケーション・プログラミング・インターフェース（API）が設けられる。このAPIはたとえば、アプリケーションとは個別の、アプリケーションによってコールされる接続マネージャ170、270および複製マネージャ160、260のようなモジュールの形で、都合のいいように実装されることができる。図1に示されるIP層120、220またはTCP層130、230を基点としたソフトウェア階層を参照すると、このAPIは、アプリケーション150、250とオペレーティングシステム140、240との間に配置されることができる。各接続110の状態は、具体的にはTCPデータとOSデータとを含む。

【0053】

好ましい実施形態は、標準のHP-UXソケットおよび関連するコールへの追加からなる。具体的には、`getsockopt()`コールは、同様の接続/確立されたソケットを構築するのに必要とされるTCP状態情報を戻すために拡張される。これは、アクティブまたはスタンバイソケット上で実行されることができ、そのソケットにもその接続にも影響を及ぼさない読出し専用動作である。状態370から390への、または状態390から370へのソケットを設定するのに同じコールが用いられる。

【0054】

状態情報の特質は、プラットフォームに応じて変化し、ネットワーク上で接続を再開することを必要とせずに、活動化されるスタンバイプロセスによって用いられるべき情報で、スタンバイソケットが更新されることを可能にするように選択される。

【0055】

`setsockopt()`コールは、アクティブソケットから取得されるTCP状態情報でスタンバイソケットを再同期できるようにするために拡張される。これは、スタンバイソケットにおいてのみ行われることができ、すべての必要とされる層に影響を及ぼす。`setsockopt()`コールは、`socket()`後にスタンバイソケット接続を形成するために実行される。これは、アクティブ接続をスタンバイ状態にする（非活動化する）ために、アクティブ接続上でコールされる。これは、スタンバイ接続をアクティブ状態にするために、スタンバイ接続上でコールされる。

【0056】

図3に戻ると、好ましい実施形態では、接続マネージャ170は、1つの場所で、すべての開始している接続110と、その状態とを追跡する。接続マネージャ170は、同時に複数の接続110上で動作を実行することができ、たとえば、すべての接続を閉じ、すべての接続を複製し、すべての接続を活動化/非活動化することができる。接続マネージャ170は、それが複製されているか否かを示すために、接続110において特定のフラグを設定することができ、その際、接続テーブル上をただ単にループすることにより、未だ複製されていない接続を複製する。接続マネージャ170は、アプリケーション150によって用いられている接続の複製、活動化、抽出、または状態取得の間の長い期間の間に、アプリケーション150を妨害するのを避けるために、マルチスレッドであるのが好ましい。

【0057】

アプリケーション150はTCP接続状態を抽出し、それを、スタンバイアプリケーション250に送信する。その後、スタンバイアプリケーション250は、状態データを、スタンバイ装置200のTCP層230に送信する。

【0058】

ここに記載されるソフトウェア障害フェイルオーバー技術は、個別のスタンバイ装置を用いて実装することに制限されない。むしろ、記載される方法は、同じ装置上のスタンバイアプリケーションに対するアプリケーションソフトウェアのフェイルオーバーにも同様に適用することができる。アプリケーションモニタが、そのアプリケーション自体に障害があるものと判断することができる場合には、ローカルにフェイルオーバーを行うことを選択することができる。

【0059】

10

20

30

40

50

2つのプロセスを含み、各プロセスが同様のサービスアプリケーション命令セットを提供し、各プロセスがIPアドレスを有し、各プロセスがインターネットプロトコルエンドポイントを有する一連のアクティブ接続を提供することができるシステムが記載される。このシステムは、他のプロセスにおいてアクティブ接続を提供しない間に、1つのアクティブプロセスの一連のアクティブ接続を提供するための手段を備える。このシステムは、そのアクティブプロセスから他のアクティブプロセスに、アクティブ接続の状態データを含むデータと、アクティブプロセスのサービスアプリケーションの命令セットの状態データを含むデータとを複製するための手段を備える。アクティブ接続の状態データを含むデータは、更新されたスタンバイ接続を該他のプロセスが維持できるようにするためのものであり、アクティブプロセスのサービスアプリケーションの命令セットの状態データを含むデータは、該他のプロセスの命令セットが更新された状態に維持されるようにするためのものである。このシステムは、該アクティブプロセスが利用できなくなったとき、アクティブプロセスのIPアドレスを非活動化し、他のプロセスのIPアドレスを活動化するための手段を備える。

10

【0060】

本発明の特定の実施形態が記載されてきたが、本発明は、そのような記載される特定の構成に限定されるべきではない。本発明は、特許請求の範囲によってのみ制限される。特許請求の範囲自体は、請求される本発明の周辺を指示することを意図しており、本明細書によって開示される例示的な実施形態のみを請求するものと解釈されるのではなく、言語そのものが許容するのと同じ広さで解釈されることを意図している。

20

【0061】

本発明は、以下の実施態様を含む。

【0062】

(1) アクティブ状態のプロセス、スタンバイ状態のプロセス、および前記スタンバイ状態のプロセスをアクティブ状態に移行させる切り換え能力を有する耐故障性プラットフォームにおいてネットワーク接続を提供する方法であって、アクティブプロセスからスタンバイプロセスに、該アクティブプロセスのネットワーク接続の状態データを複製するステップと、前記スタンバイプロセスについて、前記複製されたデータで更新された対応するスタンバイネットワーク接続を維持するステップと、前記アクティブ状態への前記スタンバイプロセスの移行中に、前記アクティブシステム内の前記ネットワーク接続を、該ネットワーク上の該接続を閉じることなく非活動化し、前記スタンバイプロセスにネットワークアドレスを転送し、該ネットワークアドレスで前記対応するスタンバイ接続を活動化するステップと、を含み、前記移行したスタンバイプロセスが前記ネットワーク接続を再開する必要があるようにする、耐故障性プラットフォームにおいてネットワーク接続を提供する方法。

30

【0063】

(2) 前記スタンバイプロセスの前記ネットワークアドレスを活動化するステップは、前記スタンバイ接続が活動化される前に実行される、上記(1)に記載の方法。

(3) 前記アクティブプロセスの前記接続を監視するステップと、該接続のアイドル状態が識別されたことに応答して前記状態データの複製を起動するステップと、を含む、上記(1)または(2)に記載の方法。

40

(4) 接続の状態データを複製する前記ステップは、前記アクティブプロセスが前記接続を用いることができる間に実行される、上記(3)に記載の方法。

(5) それぞれ異なるネットワークアドレスを有するアクティブプロセスおよびスタンバイプロセスの複数の対を維持するステップを含む、上記(1)から(4)のいずれかに記載の方法。

(6) 前記ネットワークアドレスはIPアドレスである、上記(1)から(5)のいずれかに記載の方法。

【0064】

50

(7) アクティブ状態のプロセス、スタンバイ状態のプロセス、および該スタンバイ状態のプロセスをアクティブ状態に移行させるための切り換え手段を有する耐故障性プラットフォームであって、

アクティブプロセスに関連付けられた接続から状態データを抽出する第 1 の接続マネージャと、

前記アクティブプロセスに関連付けられたネットワーク接続の状態データをスタンバイプロセスに複製する複製マネージャと、

前記スタンバイプロセスについて、前記複製されたデータで更新された対応するスタンバイネットワーク接続を維持する第 2 の接続マネージャとを備え、

前記切り換え手段は、前記スタンバイプロセスのアクティブ状態への移行の一部として、前記アクティブシステム内の前記ネットワーク接続を、該ネットワーク上の該接続を閉じることなく非活動化し、ネットワークアドレスを前記スタンバイプロセスに転送し、該ネットワークアドレスで前記対応するスタンバイ接続を活動化するように構成されており、前記移行したスタンバイプロセスが前記ネットワーク上で前記接続を再開する必要がないようにするプラットフォーム。

【 0 0 6 5 】

(8) 前記複製マネージャは、個別のソフトウェアモジュールの形をとる、上記 (7) に記載のプラットフォーム。

(9) 前記第 1 および / または第 2 の接続マネージャは、個別のソフトウェアモジュールの形をとる、上記 (7) または (8) に記載のプラットフォーム。

(1 0) 各プロセスは、命令セットを形成するアプリケーションソフトウェア層と、オペレーティングシステムソフトウェア層と、伝送制御プロトコルを適用することができるソフトウェア層とを提供し、

前記システムは、前記アクティブ接続の前記状態データを複製することができるプロセスを提供し、該プロセスは、少なくとも部分的に、前記アクティブプロセスの前記アプリケーション層によって実行される、上記 (7) から (9) のいずれかに記載のプラットフォーム。

(1 1) それぞれ異なるネットワークアドレスを有する、アクティブプロセスおよびスタンバイプロセスの複数の対を維持することを含む、上記 (7) から (1 0) のいずれかに記載のプラットフォーム。

(1 2) 前記ネットワークアドレスは IP アドレスである、上記 (7) から (1 1) のいずれかに記載のプラットフォーム。

【 0 0 6 6 】

(1 3) アクティブ状態のプロセスをサポートする手段、スタンバイ状態のプロセスをサポートする手段、および該スタンバイ状態のプロセスをアクティブ状態に移行させるための切り換え能力を有する耐故障性プラットフォームであって、ネットワーク接続は、関連付けられたプロセスが前記ネットワーク接続を閉じるか、または該関連付けられたプロセスが終わる度に、該ネットワーク接続が前記ネットワーク上で閉じられる第 1 の状態と、関連付けられたプロセスが前記ネットワーク接続を閉じるか、または該関連付けられたプロセスが終わるときに、該ネットワーク接続が該ネットワーク上で閉じられない第 2 の状態とをとることができ、

前記プラットフォームは、プログラム制御下で前記接続を前記第 1 の状態と前記第 2 の状態との間で切り換え、接続状態情報を抽出し、該接続状態情報を設定することを可能にするアプリケーションプログラミングインターフェースを含み、前記接続状態情報は、ネットワークアドレスを前記スタンバイ接続に転送することにより、前記ネットワーク上で接続が再開されることを必要とすることなく、アクティブプロセスから複製された前記状態情報で更新され維持されたスタンバイネットワーク接続が、移行したスタンバイプロセスによって用いられることを可能にするような情報である耐故障性プラットフォーム。

【 0 0 6 7 】

(1 4) アクティブ状態のプロセス、スタンバイ状態のプロセス、および該スタンバイ状

10

20

30

40

50

態のプロセスをアクティブ状態に移行させるための切り換え能力とを有する耐故障性プラットフォームにおいてネットワーク接続を提供するための方法であって、
 アクティブプロセスからスタンバイプロセスに、該アクティブプロセスのネットワーク接続の状態データを複製するステップと、
 前記スタンバイプロセスについて、前記複製されたデータで更新された対応するスタンバイネットワーク接続を維持するステップと、
 前記スタンバイ状態のプロセスをアクティブ状態に移行させるステップとを含み、
 前記スタンバイ状態のアクティブ状態への移行中に、前記アクティブシステム内の前記ネットワーク接続を、該ネットワーク上の該接続を閉じることなく非活動化し、ネットワークアドレスを前記スタンバイプロセスに転送し、該ネットワークアドレスで前記対応するスタンバイネットワーク接続を活動化し、該移行したスタンバイプロセスが前記ネットワーク上で前記接続を再開する必要があるようにする、方法。

10

【 0 0 6 8 】

(1 5) アクティブ状態のプロセス、スタンバイ状態のプロセス、および該スタンバイ状態のプロセスをアクティブ状態に移行させるための切り換え能力とを有する耐故障性プラットフォームにおいてインターネットプロトコル (I P) ネットワーク接続を提供する方法であって、
 アクティブプロセスからスタンバイプロセスに、該アクティブプロセスのネットワーク接続の状態データを複製するステップと、
 前記スタンバイプロセスについて、前記複製されたデータで更新された対応するスタンバイネットワーク接続を維持するステップと、
 前記スタンバイ状態のプロセスをアクティブ状態に移行させるステップとを含み、
 前記スタンバイ状態のアクティブ状態への移行中に、前記アクティブシステム内の前記ネットワーク接続を、該ネットワーク上の該接続を閉じることなく非活動化し、 I P ネットワークアドレスを前記スタンバイプロセスに転送し、該 I P ネットワークアドレスで前記対応するスタンバイネットワーク接続を活動化し、前記移行したスタンバイプロセスが前記ネットワーク上で前記接続を再開する必要があるようにし、
 前記スタンバイプロセスの前記ネットワークアドレスを活動化するステップは、前記スタンバイ接続が活動化される前に実行される、方法。

20

【 0 0 6 9 】

30

【 発明の効果 】

本発明によれば、ネットワーク接続のリモートエンドに対して高いトランスペアレンシー (透過性) が保持され、フェイルオーバが達成される。

【 図面の簡単な説明 】

【 図 1 】 本発明の一実施形態に従う、システムのアーキテクチャの概略図。

【 図 2 】 本発明の一実施形態に従う、システムにおける 1 つの装置の種々の状態を示す状態図。

【 図 3 】 本発明の一実施形態に従う、複製マネージャモジュールおよび接続マネージャモジュールを含むシステムの概略図。

【 図 4 】 本発明の一実施形態に従う、システム内の接続の種々の状態を示す状態図。

40

【 図 5 】 本発明の一実施形態に従う、アクティブ接続およびスタンバイ接続を設定するために用いられるプロセス、および閉じるために用いられるプロセスを示す図。

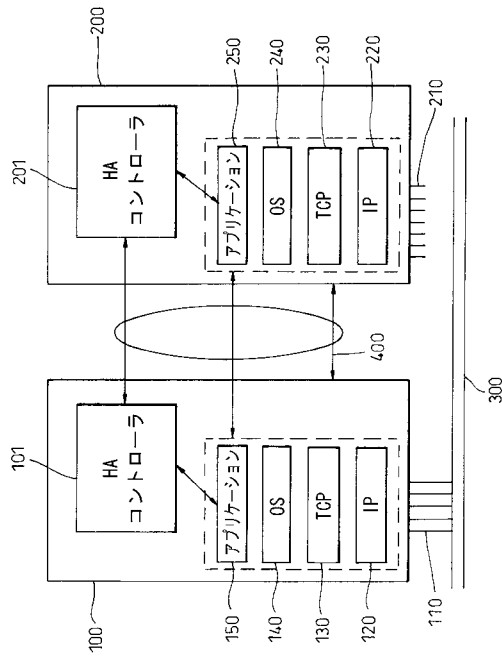
【 図 6 】 本発明の一実施形態に従う、プロセス終了時、またはアクティブシステムからスタンバイシステムへの接続の手動による移動時のプロセスを示す図。

【 符号の説明 】

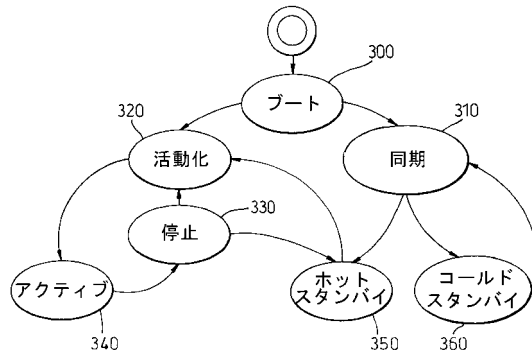
- 1 0 0、2 0 0 ハードウェア装置
- 1 1 0、2 1 0 接続
- 3 0 0 ネットワーク
- 1 6 0、2 6 0 複製マネージャ
- 1 7 0、2 7 0 接続マネージャ

50

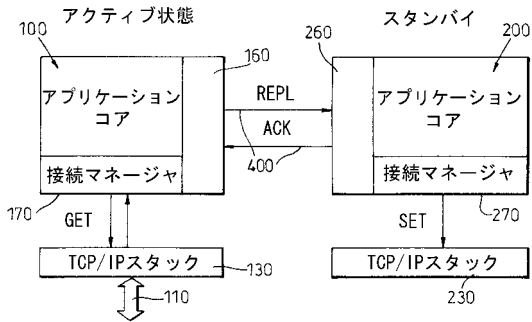
【図 1】



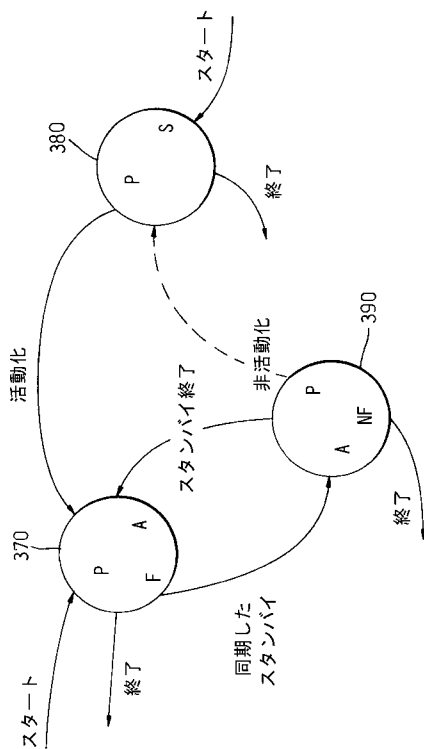
【図 2】



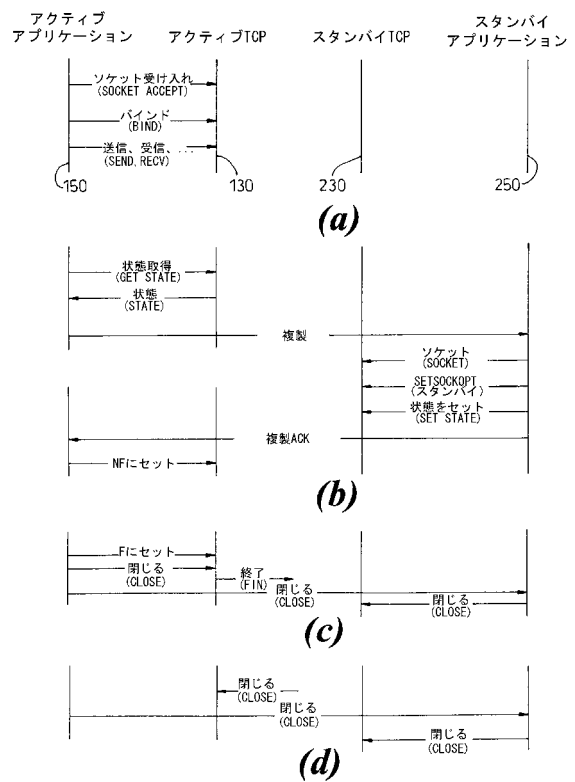
【図 3】



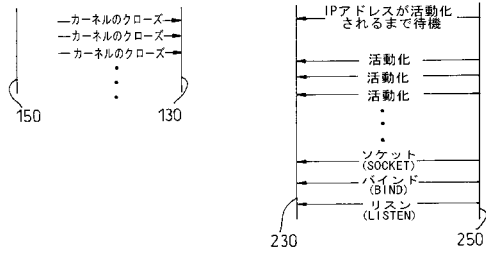
【図 4】



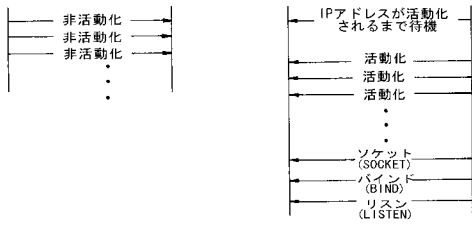
【図 5】



【図 6】



(a)



(b)

フロントページの続き

- (72)発明者 ジャック・ボウド
フランス国3 8 6 6 0 サン・ヴァンサン・デュ・マルキュゼ、シャマン・ジュ・トーティエ
- (72)発明者 ステラ・クウォン
アメリカ合衆国9 4 0 2 2 カリフォルニア州ロス・アルトス、マニエル・ロード 1 4 4 7 0
- (72)発明者 アイザック・ウォン
アメリカ合衆国9 5 1 5 9 カリフォルニア州サン・ノゼ、ピー・オー・ボックス 2 6 4 5 4
- (72)発明者 デニス・ロジャー
フランス国3 8 7 6 0 ヴァルス、リュト・ジュ・マルティネ・ダン・バ、レ・マージョエラ、ロト
2 0
- (72)発明者 アブデサター・サッシ
フランス国エフ - 3 8 1 0 0 グルノーブル、シャマン・デュ・レグリス 7

審査官 石井 研一

- (56)参考文献 国際公開第0 0 / 0 6 0 4 7 2 (WO, A 1)
特開平0 8 - 0 1 6 4 2 2 (JP, A)
国際公開第0 0 / 0 3 1 9 4 2 (WO, A 1)

- (58)調査した分野(Int.Cl., DB名)
H04L 12/56
H04L 12/28