

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3759410号

(P3759410)

(45) 発行日 平成18年3月22日(2006.3.22)

(24) 登録日 平成18年1月13日(2006.1.13)

(51) Int. Cl.

G06F 9/54 (2006.01)

F I

G06F 9/46 480Z

請求項の数 11 (全 27 頁)

(21) 出願番号	特願2000-562827 (P2000-562827)	(73) 特許権者	591067923
(86) (22) 出願日	平成11年7月21日 (1999.7.21)		ユニシス コーポレーション
(65) 公表番号	特表2002-521765 (P2002-521765A)		UNISYS CORPORATION
(43) 公表日	平成14年7月16日 (2002.7.16)		アメリカ合衆国 ペンシルバニア州 19
(86) 国際出願番号	PCT/US1999/016540		424 ブルーベル, ピー. オー. ボ
(87) 国際公開番号	W02000/007101		ックス 500 タウンシップ ライン
(87) 国際公開日	平成12年2月10日 (2000.2.10)		アンド ユニオン ミーティング ローズ
審査請求日	平成13年1月25日 (2001.1.25)		(番地なし)
(31) 優先権主張番号	09/127,167	(74) 代理人	100064746
(32) 優先日	平成10年7月29日 (1998.7.29)		弁理士 深見 久郎
(33) 優先権主張国	米国 (US)	(74) 代理人	100085132
			弁理士 森田 俊雄
		(74) 代理人	100083703
			弁理士 仲村 義平

最終頁に続く

(54) 【発明の名称】 クラスタ化された計算処理環境において実行する分散ネットワークアプリケーションの管理用リクエストを処理するための方法および装置

(57) 【特許請求の範囲】

【請求項1】

定義されたマスタノードを介した集中管理を必要とする分散ネットワークアプリケーションが集中管理をサポートしないクラスタ化された計算処理環境のノード(10、12)上で実行することを可能にするための方法であって、前記方法は、

クラスタ化された計算処理環境からの管理用リクエストをその親ノードで受取るステップと、

親ノードが分散ネットワークアプリケーションのための指定されたマスタノードであるかどうかを決定するステップと、

もし親ノードが指定されたマスタノードでなければ親ノードから指定されたマスタノードへ管理用リクエストを経路制御するステップとを含む、方法。

【請求項2】

前記経路制御するステップは、

親ノードとマスタノードとの間に接続を与える名前付きパイプのインスタンスを生成するステップと、

名前付きパイプを介して親ノードからマスタノードへ管理用リクエストを送るステップとを含む、請求項1に記載の方法。

【請求項3】

名前付きパイプを介してマスタノードから親ノードへ管理用リクエストに対する応答を送り返すステップをさらに含む、請求項2に記載の方法。

10

20

【請求項 4】

名前付きパイプを介してマスタノードで管理用リクエストを受取るステップと、分散ネットワークアプリケーションの管理アプリケーションプログラミングインターフェイス (API) を呼出して、指定されたマスタノードによるリクエストの処理を起動するステップとをさらに含む、請求項 1 に記載の方法。

【請求項 5】

コンピュータ読出可能媒体上で実現されるサーバプログラム (42) が記録されたコンピュータ読出可能媒体であって、サーバプログラムは、定義されたマスタノードを介した集中管理を必要とする分散ネットワークアプリケーションが集中管理をサポートしないクラスタ化された計算処理環境のノード (10、12) 上で実行することを可能にするプログラムコードを含み、前記サーバプログラムのプログラムコードは、クラスタ化された計算処理環境の各ノード上で実行されると、前記ノードのいずれかに、以下のステップ：

クラスタ化された計算処理環境からの管理用リクエストをそのノードで受取るステップを含み、そのノードは親ノードを定義し、

親ノードが分散ネットワークアプリケーションのための指定されたマスタノードであるかどうかを決定するステップと、

もし親ノードが指定されたマスタノードでなければ親ノードから指定されたマスタノードへ管理用リクエストを経路制御するステップとを含むステップを実行させる、サーバプログラムが記録されたコンピュータ読出可能媒体。

【請求項 6】

プログラムコードは、親ノードとマスタノードとの間に接続を与える名前付きパイプのインスタンスを生成し、名前付きパイプを介して親ノードからマスタノードへ管理用リクエストを送ることによって、ノードに前記経路制御するステップを実行させる、請求項 5 に記載のサーバプログラムが記録されたコンピュータ読出可能媒体。

【請求項 7】

プログラムコードはさらに、マスタノードに、名前付きパイプを介してマスタノードから親ノードへ管理用リクエストに対する応答を送り返させる、請求項 6 に記載のサーバプログラムが記録されたコンピュータ読出可能媒体。

【請求項 8】

定義されたマスタノードを介した集中管理を必要とする分散ネットワークアプリケーションが集中管理をサポートしないクラスタ化された計算処理環境のノード (10、12) 上で実行することを可能にするサーバプログラムが記録されたコンピュータ読出可能媒体であって、サーバプログラムは、各ノード上でクラスタ化された計算処理環境からの管理用リクエストをそのノードで傍受し、もし親ノードが指定されたマスタノードでなければ親ノードから指定されたマスタノードへ管理用リクエストを経路制御する、サーバプログラムが記録されたコンピュータ読出可能媒体。

【請求項 9】

サーバプログラムは、前記管理用リクエストを経路制御する前に親ノードが分散ネットワークアプリケーションのための指定されたマスタノードであるかどうかを決定する、請求項 8 に記載のサーバプログラムが記録されたコンピュータ読出可能媒体。

【請求項 10】

サーバプログラムは、親ノードとマスタノードとの間に接続を与える名前付きパイプのインスタンスを生成し、次に名前付きパイプを介して親ノードからマスタノードへ管理用リクエストを送ることによって管理用リクエストを指定されたマスタに経路制御する、請求項 9 に記載のサーバプログラムが記録されたコンピュータ読出可能媒体。

【請求項 11】

サーバプログラムはさらに、名前付きパイプを介してマスタノードから親ノードへ管理用リクエストに対する応答を送り返す、請求項 8 に記載のサーバプログラムが記録されたコンピュータ読出可能媒体。

【発明の詳細な説明】

10

20

30

40

50

【 0 0 0 1 】

【 著作権表示 】

この特許書類の開示の一部は、著作権保護の対象となる事項を含む。著作権者は、米国特許商標庁の特許書類または記録にある特許書類または特許開示に関して何人による複写にも異議を唱えないが、他の態様では何であれすべての著作権を留保する。

【 0 0 0 2 】

【 背景 】

発明分野

本発明は、クラスタ化された計算処理環境における分散ネットワークアプリケーションに向けられ、より特定のには、そのような環境における分散アプリケーションの管理用リクエストを処理するための方法および装置に向けられる。

10

【 0 0 0 3 】

先行技術の説明

分散ネットワークアプリケーションは、しばしば、分散アプリケーション全体の管理を行なうマスタマシンの概念を有する。この場合には、管理は、構成要素をオンラインにしたり、構成要素をオフラインにしたり、または個々の構成要素のステータスをチェックしたりすることを含む。複数のコンピュータがクラスタ内であわせて動作して全体の可用性を向上させることを可能にするソフトウェアアーキテクチャには、これらの分散ネットワークアプリケーションによって必要とされる集中管理の種類をサポートしないものがある。そのようなクラスタ化された計算処理環境の一例は、Ghormley, D.P.らの「G L U N I X : ワークステーションのネットワークのためのグローバル層 Unix (" A global layer Unix for a Network of Workstations ") 」Software Practice & Experience, GB, John Wiley & Sons Ltd. Chichester, Vol.28 no. 9, 929-961頁, XP000765510 ISSN:0038-0644に記載されている。しかしながら、Microsoft (登録商標) Cluster Server (M S C S) のアーキテクチャなどの、他のクラスタ化された計算処理環境は、そのような分散ネットワークアプリケーションの集中管理モデルに適さない。たとえば、M S C S では、クラスタの各ノード上のカスタムリソースDLLが、構成要素をオンラインにしたり、構成要素をオフラインにしたり、または構成要素のステータスをチェックしたりするなどの、そのノードのための管理操作を実現しなければならない。したがって、管理操作 (たとえば、オフライン、オンラインなど) が所与のノード上で行なわれるとき、M S C S は、そのノード上のリソースDLLを呼出して操作を行なう。これは、そのような操作のすべてがクラスタ内の単一のマスタノードによって起動および制御される必要がある、すなわち、たとえば、ノードがオンラインまたはオフラインにされるべきならば、その操作は分散アプリケーションのためのマスタノードによって起動および実行されるであろう多くの分散ネットワークの集中管理モデルに反するものである。より具体的な例として、たとえば、BEA System, Inc.のTuxedo (登録商標) トランザクションマネージャおよびメッセージングミドルウェアのために書かれた分散アプリケーションが、クラスタ化された計算処理環境においてMicrosoftのCluster Server (M S C S) ソフトウェアの制御下で実行される

20

30

【 0 0 0 4 】

A . Tuxedo

Tuxedo環境では、1つ以上の論理的マシンをあわせてグループ化してタキシードドメインを規定することができる。1個の論理的マシンが、タキシードサーバマシンを表わし、しばしば物理的コンピュータシステム (すなわち、ネットワークの1つのノード) に対応するが、いつもそうとは限らない。いくつかの場合、2個以上の論理的マシンが単一ノード上に定義可能である。

40

【 0 0 0 5 】

各論理的マシンは、B B Lプロセス、ブリッジプロセス、およびおそらくD B B Lプロセスを含む、タキシードアドミンサーバを含む。各論理的マシンはまた、タキシードTリスンプロセスを有する。B B Lプロセスは、ユーザ書込のサーバを監視し制御する。ある

50

論理的マシン上のブリッジプロセスは、Tリッスンプロセスを別の論理的マシンに接続する。ドメイン内の論理的マシンの1つがマスタとして指定される；マスタにはバックアップが構成されてもよい。DBBLは、マスタとして動作する論理的マシン上で実行される。DBBLは、さまざまなBBLを調整し制御する、すなわち、それは、構成要素をオンラインにしたり、構成要素をオフラインにしたり、または個々の構成要素のステータスをチェックしたりすることを含む、ドメインのための管理を行なう。

【0006】

論理的マシン上のユーザ書込のサーバは、1つ以上のタキシードサーバグループにグループ化される。タキシードサーバグループがトランザクションをサポートする場合、トランザクションの状態は、TLOGにログされる。各論理的マシンは、それ自身の独自のTLOGを有する。他の状態および構成情報は、バイナリTUXCONF I Gファイルに記憶される。各論理的マシンは、それ自身のTUXCONF I Gのコピーを有するが、DBBLはさまざまなTUXCONF I Gファイルの同期をとる。

10

【0007】

B. Microsoft Cluster Server (MSCS)

Microsoft (登録商標) Windows NT (登録商標), Enterprise Editionの現在のリリースでは、クラスタは、可用性の高いシステムとして動作する2つの接続されたシステム(ふつうノードと呼ばれる)のことをいう。これらの2つのシステムは、共有SCSIバス上に1つ以上の物理的ディスクドライブを有する；共有SCSIバス上の各ディスクはいずれのシステムによってもアクセス可能であるが、同時には1つのシステムによってのみ可能である。クラスタ内のシステムの1つに障害が発生すると、他のシステムが共有ディスクの制御を獲得し、障害が発生したシステムが止めたところをひろってサービスを与え続ける。Windows NTの将来のリリースは、3つ以上のノードでクラスタをサポートするものと見込まれる。

20

【0008】

MSCSは、クラスタの各ノード上で実行するクラスタサービスによって制御される。クラスタサービスは、1つ以上のリソースモニタを生じる。リソースモニタは、1つ以上のリソースを監視する。リソースとは、MSCSによって監視および制御可能であるエンティティである。リソースモニタは、リソースDLL内のエントリポイント(MSCSリソースAPI内に定義される)を呼出して、特定のリソースを監視および制御する。特に、リソースモニタは、適切なときにリソースDLL内のエントリポイントを呼出して、リソースの状態をチェックしたり、リソースをオンラインにしたり、またはリソースをオフラインにする。したがって、リソースDLLは、リソースをオンラインにしたりリソースをオフラインにしたりするのに必要とされる動作 - Tuxedo環境においてはタキシードマスタによって行なわれなければならない管理操作を実現することに注目されたい。あるノードから別のノードへリソースを移すことは、第1のノード上でリソースをオフラインにし、第2のノード上でリソースをオンラインにすることによっておこる。

30

【0009】

MSCSは、いくつかのリソースタイプを定義するリソースDLL (clusres.dll) を含む。以下は、clusres.dllによってサポートされる関連のリソースタイプのリストである

40

：
物理的ディスク - 共有SCSIバス上に位置する物理的ディスクドライブを制御する。

【0010】

IPアドレス - あるノードまたは他のノードに動的に割当可能であるIPアドレスを定義する。

【0011】

ネットワーク名 - IPアドレスリソースのための記号名を定義する。

汎用サービス - いかなるWindows NTサービスをも制御する。

【0012】

汎用アプリケーション - 正しく振舞うWindows NTコンソールアプリケーションプログラム

50

を制御する。

【 0 0 1 3 】

加えて、第三者の開発者が、リソース A P I に従うリソース D L L を開発および登録することによってカスタムリソースタイプを生成することができる。

【 0 0 1 4 】

M S C S リソースは、あわせてリソースグループにグループ化可能である。リソースグループは、フェールオーバの基本ユニットである；すなわち、所与の時点では特定のグループ内のあらゆるリソースが同じノード上で実行する。リソースグループ内のリソースは、オンライン操作およびオフライン操作の相対順序を制御する依存性を有し得る。

【 0 0 1 5 】

典型的には、リソースグループは、1つ以上の物理的ディスクリソースと、I P アドレスリソースと、ネットワーク名リソースと、汎用サービスリソース、汎用アプリケーションリソース、および/またはカスタムリソースタイプなどの、サーバアプリケーションを表わす1つ以上の付加的なリソースとを含む。それ自身のI P アドレスリソースおよびネットワーク名リソースを有するリソースグループは、仮想サーバとして知られている。

【 0 0 1 6 】

外部のクライアントには、1個の仮想サーバが独自のサーバコンピュータとしてT C P / I P クライアント/サーバタイプのアプリケーションを実行しているように見える。実際には、M S C S クラスタの単一のノード上で実行するいくつかの仮想サーバがあり、各々が異なったI P アドレスを有するだろう。さらに、仮想サーバは、M S C S クラスタのあるノードから別のノードへ移ることができ、これはクライアントには見えない(ただし、一時的割込またはサービスの減速がある)。

【 0 0 1 7 】

図1は、各々がそれぞれの仮想サーバ14、16を実行する2つのノード10、12を含む例示的なクラスタを示す。各ノード上で実行するクラスタサービス18は、たとえば、各仮想サーバ14、16のリソースモニタ(図示せず)への管理用リクエストをイニシエートすることを含め、クラスタを制御する。各ノード10、12は、たとえば、S C S I バスなどの共有バス24を介して複数個のディスク20に接続される。各ノード10、12はまた、それぞれのネットワークインターフェイスを介してローカルエリアネットワーク(L A N) 22に接続される。

【 0 0 1 8 】

ここで図2を参照して、クラスタの第2のノード12に障害が発生した場合、M S C S は、障害が発生した仮想サーバ16をクラスタの第1のノード10上で起動する。これがフェールオーバである。この時点で、仮想サーバ14、16はどちらも、第1のノード上で実行している。クライアントプログラムは引き続き正常に実行する(ただし、性能が落ちる可能性がある)。

【 0 0 1 9 】

クラスタの第2のノード12が正常動作を再開すると、M S C S は、フェールオーバした仮想サーバ16を第1のノード10上でオフラインにする。次に、M S C S は、この仮想サーバ16を第2のノード12上で再びオンラインにする。これがフェールバックである。この時点で、図1に示される構成が再開し、すなわち各ノードは1個の仮想サーバ上で実行している。

【 0 0 2 0 】

C . クラスタ化された環境における分散ネットワークアプリケーションの配備

Tuxedo環境において実行するよう設計されたアプリケーションなどの、分散ネットワークアプリケーションは、M S C S などのクラスタ化された計算処理環境内で実行するよう設定可能である。たとえば、図1を参照して、2つの論理的マシンを備えるタキシードドメインは、2つのM S C S 仮想サーバ14、16を備えるよう構成可能である。正常動作中は、1個の仮想サーバ(およびしたがって1個の論理的マシン)がクラスタのノード10、12の各々上で実行する。しかしながら、前に例示したように、管理用リクエストがTu

10

20

30

40

50

xedo環境とM S C S環境とは異なって処理されるという点で問題が生じる。Tuxedoアーキテクチャでは、構成要素をオンラインにしたりオフラインにしたりすることを含む、管理用リクエストは、指定されたマスタノードによって行なわれる。反対に、M S C S環境では、管理用リクエストは、指定されたマスタノードによってではなく、操作が向けられるノード上のリソースD L Lによって行なわれる。本発明は、この課題に対する解決を提供する。

【 0 0 2 1 】

【 発明の概要 】

本発明は、アプリケーションがクラスタ化により達成される増大した可用性を利用できるように、マスタノードを介した集中管理を必要とする分散ネットワークアプリケーション（たとえば、Tuxedoベースのアプリケーション）がそのような集中管理をサポートしないクラスタ化された環境（たとえば、M S C S）内に配備されることを可能にするための方法および装置に向けられる。より特定的には、本発明の方法および装置は、クラスタ化された計算処理環境内で実行する分散ネットワークアプリケーションの管理用リクエスト（たとえばオンライン、オフラインなど）を、リクエストが発するノードから、そのクラスタ内の分散ネットワークアプリケーションのためのマスタとして動作しているノードへ経路制御する。たとえば、ある具体的な実施例では、この発明は、M S C SリソースD L LからM S C Sで制御されたクラスタ内で実行するTuxedoアプリケーションのマスタマシンへTuxedoベースの管理用リクエストを経路制御する。

【 0 0 2 2 】

好ましい実施例では、管理用リクエストの経路制御は、クラスタの各ノード上にインストールされる、パイプサーバと以下呼ばれる複数個のサーバプログラムによって達成される。所与のノードのパイプサーバは、名前付きパイプのインスタンスを生成し、これを介して、他のノード上のパイプサーバを含む、クライアントプログラムは、管理用リクエストを含むメッセージをそのノードに送ることができる。パイプサーバは、受取ったメッセージを、分散ネットワークアプリケーションの管理A P Iに転送する。

【 0 0 2 3 】

本発明のさらなる特徴および利点は、以下に明らかとなるであろう。

前述の概要、および好ましい実施例の以下の詳細な説明は、添付の図面と関連付けて読まれるとよりよく理解される。この発明を例示する目的のために、現在のところ好ましい実施例が図面に示されるが、この発明は開示される特定の方法および手段に限られるものではないことが理解される。

【 0 0 2 4 】

【 好ましい実施例の詳細な説明 】

I . 概観

本発明は、アプリケーションがクラスタ化により達成される増大した可用性を利用できるように、マスタノードを介した集中管理を必要とする分散ネットワークアプリケーションがそのような集中管理をサポートしないクラスタ化された環境内に配備されることを可能にするための方法および装置に向けられる。より特定的には、本発明の方法および装置は、クラスタ化された計算処理環境内で実行する分散ネットワークアプリケーションの管理用リクエスト（たとえば、オンライン、オフラインなど）を、リクエストが発するノードから、そのクラスタ内の分散ネットワークアプリケーションのためのマスタとして動作するノードへ経路制御する。本発明は、フロッピィディスク、C D - R O M、ハードドライブ、またはいかなる他の機械で読出可能な記憶媒体などの有形の媒体に実現されるプログラムコード（すなわち、命令）の形をとってもよく、プログラムコードがコンピュータなどのマシンによってロードされ実行されると、マシンがこの発明を実施するための装置になる。本発明はまた、電気的配線もしくはケーブルによって、光ファイバを介して、または伝送のいかなる他の形によるなどの、何らかの伝送媒体によって伝送されるプログラムコードの形で実施可能であり、プログラムコードが、コンピュータなどのマシンによって受取られロードされ実行されると、マシンがこの発明を実施するための装置になる。汎用

プロセッサ上で実現されると、プログラムコードは、プロセッサと組合さって、特定の論理回路と同じように動作する独特の装置を提供する。

【0025】

本発明に従うと、アプリケーションがマスタノードを介した集中管理を必要とするような、複数個のノードを含むクラスタ化された計算処理環境内で実行する分散ネットワークアプリケーションの管理用リクエストは、リクエストが発するノードから、分散ネットワークアプリケーションのためのマスタとして動作するノードへ経路制御される。次に、リクエストの成功または失敗を示す応答が、マスタノードから、リクエストを発したノードまで送り返される。管理用リクエストの経路制御は、クラスタの各ノード上にインストールされる、パイプサーバと以下に呼ばれる複数個のサーバプログラムによって達成される。10
所与のノードのパイプサーバは、名前付きパイプのインスタンスを生成し、これを介して、他のノード上のパイプサーバを含む、クライアントプログラムは、管理用リクエストを含むメッセージをそのノードに送ることができる。パイプサーバは、受取ったメッセージを、分散ネットワークアプリケーションの管理APIに転送する。

【0026】

図3Aおよび図3Bは、本発明の方法の好ましい実施例を例示する流れ図を含む。図示のとおり、ステップ60で、クラスタ化された環境は、管理用リクエスト（たとえば、分散ネットワークアプリケーションの構成要素をオンラインにするか、またはこれをオフラインにするためのリクエスト）を親ノードで発行する。次に、ステップ64で、分散ネットワークアプリケーションの構成要素は、親ノード上で実行するパイプサーバに名前付きパイプを開き、ステップ66で、名前付きパイプを介してパイプサーバにリクエストメッセージを送る。20

【0027】

ステップ68で、親ノード上のパイプサーバは、親ノードが分散ネットワークアプリケーションのマスタノードとして指定されるかどうかを決定する。もしそうならば、ステップ70で、リクエストはローカルに処理される。

【0028】

ステップ68で、親ノードが指定されたマスタノードでないことが決定された場合、ステップ72で、パイプサーバは、そのマスタノードがどのマスタノードであるかを決定する。次に、ステップ74で、親ノード上のパイプサーバは、指定されたマスタノード上のパイプサーバに名前付きパイプを開く。ステップ76で、リクエストメッセージは、名前付きパイプを介してマスタノード上のパイプサーバに転送される。30

【0029】

ステップ78で、マスタノード上のパイプサーバは、これへの名前付きパイプ接続を介してメッセージを受取る。ステップ80で、マスタノード上のパイプサーバは、分散ネットワークアプリケーションの管理APIを呼出して、メッセージ内に含まれるリクエストの処理を開始する。ステップ82（図3B）で、マスタノード上のパイプサーバは、分散ネットワークアプリケーションからリクエスト処理の結果（たとえば、成功または失敗）を獲得し、ステップ84で、同じ名前付きパイプを介して親ノードのパイプサーバに応答を送り返す。ステップ88で、親ノード上のパイプサーバは、親ノード上のクラスタ環境に40
応答を転送し、処理を完了する。

【0030】

II. 例示的アプリケーション

図4は、本発明の方法を実施するための装置の好ましい実施例と、（集中管理をサポートしない）MSSCによって制御されかつ複数個のノードを含むクラスタ化された計算処理環境内に配備される（マスタノードを介した集中管理を必要とする）Tuxedo環境のために書かれた分散ネットワークアプリケーションを含む環境における本発明の例示的アプリケーションとの両方を例示するブロック図である。以下により詳細に記載されるように、タキシードアプリケーションの管理用リクエストは、リクエストが発するノードのリソースDLLからタキシードアプリケーションのためのマスタとして動作するノードへ経路制御50

される。しかしながら、本発明は、Tuxedoベースのアプリケーションおよび/またはM S C S制御のコンピュータクラスタでの使用に限られるものでなく、分散ネットワークアプリケーションがマスタノードを介した集中管理を必要とするが、アプリケーションが配備されるクラスタ化された環境がそのような集中管理をサポートしないようないかなるシステムにおいても実施可能である。

【0031】

以下に記載される実施例に従うと、Tuxedoベースの分散ネットワークアプリケーションの各論理的マシンはそれぞれのM S C S仮想サーバに対応付けられている。この例では、各論理的マシンは、Oracle(登録商標)データベースを制御する。さらなる背景として、Oracle Corporationから入手可能なデータベースソフトウェアに含まれる無料のアドオンプロダクト - Oracle(登録商標) FailSafe - が、オラクルデータベースと呼ばれるカスタムリソースタイプを定義するリソースDLLを提供する。Oracle Corporationは、Oracle FailSafe Managerと呼ばれる管理用ツールも提供して、M S C S環境のためにOracleリソースを生成および管理する。

10

【0032】

図4は、例示の構成の各仮想サーバ/論理的マシンの機能構成要素を例示する。図示のとおり、各論理的マシンごとに、以下のリソースを含むM S C Sリソースグループが生成される：

物理的ディスク30 - 共有SCSIバス(図示せず)上の1つ以上の物理的ディスク30 - これらのディスクは、タキシード論理的マシンによって使用されるTLOG、TUXCONFIG、およびオラクルデータベースを含む；

20

オラクルリソース32 - タキシード論理的マシンによって使用されるオラクルデータベースのOracle FailSafe Managerによって生成されるリソース；

タキシードTリッスン34；

タキシードアドミンサーバ36；および

タキシードサーバグループ38 - タキシード構成のGROUPSセクション内にこの論理的マシンのために構成される各サーバグループごとに1つのタキシードサーバグループリソースインスタンスがある。

【0033】

各論理的マシンについてのリソースグループは、ローカルエリアネットワーク(図示せず)上のIPアドレスと、タキシード構成のMACHINESセクション内で用いられるマシン名と一致するネットワーク名と、NETWORKセクションにおいて指定されるNLSADDR値およびNADDR値とを有する。本発明に従うと、各リソースグループはさらに、リソースDLL40とパイプサーバ42とを含み、これらはあわせて本発明の方法を実施するための装置を提供する。以下に記載されるように、リソースDLL40は、クラスタ環境とそのノード上のタキシードアプリケーションの他の構成要素との間にインターフェイスを与えるタキシード分散ネットワークアプリケーションの構成要素を表わす。

30

【0034】

この構成では、たとえば2個の論理的マシンを備えるタキシードドメインは、2個のM S C S仮想サーバを備えるよう構成される。正常動作中は、1個の仮想サーバ(およびしたがって1個の論理的マシン)がクラスタのノードの各々上で実行する。クラスタの第1のノードに障害が発生した場合、M S C Sは、障害が発生した仮想サーバ(および対応する論理的マシン)をクラスタの第2のノード上で起動する - すなわちフェールオーバー状況である。両方の仮想サーバ(および両方の論理的マシン)が次に、第2のノード上で実行する。タキシードクライアントプログラムは引き続き正常に実行する(ただし、性能が落ちる可能性がある)。クラスタの第1のノードが正常動作を再開すると、M S C Sは、第2のノード上で、フェールオーバーした仮想サーバ(および対応する論理的マシン)をオフラインにする。次にM S C Sは、第1のノード上でこの仮想サーバ(および対応する論理的マシン)をオンラインにする。これがフェールバック状況である。フェールバックの後、各ノードは再び、1個の仮想サーバ(およびしたがって1個の論理的マシン)を実行する。

40

50

以下に記載するとおり、本発明の方法は、これおよび他のシナリオに関連する管理用リクエストが、親ノードからタキシードドメインのためのマスタとして動作するノードへ経路制御されることを可能にする。

【0035】

A. リソースDLL40

リソースDLL40は、以下のリソースタイプ：タキシードサーバグループ、タキシードアドミンサーバ、パイプサーバ、およびタキシードTリッスンをサポートする。各リソースタイプは、リソースDLL40内のエントリポイント関数44、46、48および50のそれぞれの組によってサポートされる。タキシードサーバグループ・リソースタイプ44は、タキシードサーバグループ38の操作を監視しかつ制御するために用いられる。タキシードアドミンサーバ・リソースタイプは、タキシードアドミンサーバ36の操作を監視しかつ制御するために用いられる。パイプサーバ・リソースタイプは、パイプサーバプロセス42の操作を監視しかつ制御するために用いられる。タキシードTリッスン・リソースタイプは、タキシードTリッスンプロセス34の操作を監視しかつ制御するために用いられる。

10

【0036】

リソースDLLは、リソースDLLによってサポートされる各リソースタイプ（たとえば、リソースタイプ44、46、48および50）について実現されなければならないエントリポイント関数を特定するマイクロソフトリソースAPIの仕様に従う。リソースモニタ26は、これらのエントリポイント関数を呼出して、特定のリソース上の所与の関数を実行する。リソースAPIはまた、リソースモニタ26によって実現されリソースDLLから呼出されるコールバック関数を特定する。リソースDLLによって実現されなければならないエントリポイント関数のうち、2つが、タキシードアプリケーションなどの分散ネットワークアプリケーションがマスタノードによって集中的に実行されることを必要とし得る管理用リクエストのタイプの例である、すなわち、（1）オンライン、および（2）オフラインである。各リソースタイプごとに1個のオンライン関数と1個のオフライン関数とがある。オンラインの目的は、リソースの特定のインスタンスをオンラインにすることであり、オフラインの目的は、リソースの特定のインスタンスをオフラインにすることである。この発明に特に関連するのは、タキシードアドミンサーバ・リソースタイプ46のオンライン関数とオフライン関数とである。これらの関数は、タキシードアドミンAPI内に対応する関数を有するが、これはタキシードアプリケーションの指定されたマスタによって実行されなければならない。したがって、MCS環境のリソースAPIは、リソースDLLがエントリポイント関数を与えてリクエストが発するノードでオンラインリクエストおよびオフラインリクエストを処理することを必要とするが、タキシードアドミンAPIは、そのようなリクエストが指定されたマスタノードによって集中的に実行されることを必要とする。これは、本発明が解決を提供する課題の特定の例の1つである。

20

30

【0037】

この発明に従うと、リソースDLL40は、タキシードリソースを制御するためにタキシードアドミンAPIと直接インターフェイスしない。代わりに、リソースDLLは、パイプサーバAPI（以下に記載される）を用いて、リクエストをパイプサーバ42に送る。パイプサーバ42は次に、タキシードアドミンAPIを用いてタキシードリソースを監視および制御する。したがって、複数個のノードを含むが集中管理をサポートしないクラスタ化された計算処理環境において、パイプサーバ42は、分散ネットワークアプリケーションの管理用リクエストを、リクエストが発するノードからそのアプリケーションのためのマスタとして指定されるノードへ経路制御するための手段を提供する。

40

【0038】

リソースDLL40が同じノード上のタキシードアドミンAPIと直接インターフェイスせず、代わりにそのノード上のパイプサーバ42とインターフェイスする理由は、リソースDLL40はネイティブクライアントとして動作不可能であるがパイプサーバ42は可

50

能であるからである。タキシードアドミンAPIコールを実行するクライアントプログラムは、ネイティブタキシードクライアントでなければならない。ネイティブクライアントとは、タキシード構成のMACHINESセクション内に特定されるマシンの1つの上で実行するクライアントである。リソースDLL40は、リソースモニタ（たとえば、リソースモニタ26）のプロセス下で実行し、（WIN32API関数GetComputerNameによって戻される）そのコンピュータ名は、それが実行しているクラスタノードの名前と一致する。しかしながら、タキシード構成は、対応する仮想サーバのネットワーク名をマシン名として用いる。したがって、リソースDLLは、ネイティブクライアントとして動作できない。反対に、パイプサーバ42は、GetComputerNameが呼出されたとき（クラスタノード名の代わりに）ネットワーク名が戻されるようにするMCS特徴で実行する。したがって、パイプサーバ42は、ネイティブクライアントとして動作することができる。

10

【0039】

B. パイプサーバ42

本発明の好ましい実施例に従うと、名前付きパイプ機能を用いて、クライアントプロセス（たとえば、リソースDLL40）とパイプサーバ42との間、および/または2つの異なるノード上のパイプサーバ42間でメッセージをやり取りする。したがって、名前付きパイプを用いて、同じコンピュータ上のプロセス間、またはネットワーク上の異なるコンピュータ上のプロセス間で双方向の通信が可能になる。この実施例では、Microsoft Windows NTオペレーティングシステムの名前付きパイプ機能が採用される。

【0040】

20

パイプサーバ42は、管理用リクエストを含むメッセージを名前付きパイプを介してプロセス間でやり取りするためのアプリケーション・プログラム・インターフェイス（API）を与える。パイプサーバAPIは、交換されるべきメッセージのフォーマットおよび、交換を行なうための方法を特定する。

【0041】

1. パイプサーバAPI - メッセージタイプ

メッセージは、32ビットのリクエストコードを含み、特定のリクエストメッセージの内容が後に続く。リクエストコードは、リクエストのタイプを特定する。オンラインリクエストおよびオフラインリクエストを処理するために、パイプサーバAPIは、以下のリクエストコードを定義する：

30

【0042】

【数1】

```
#define PIPE_REQ_MIB          1
#define PIPE_REQ_ADM_ONLINE  2
#define PIPE_REQ_ADM_OFFLINE 3
```

【0043】

PIPE__REQ__MIBは、マネジメント・インフォメーション・ベース（MIB）にアクセスするタキシードアドミンリクエストの一般的な形を表わすメッセージタイプである。MIBは、タキシードのプログラムされた管理を行なうためにアクセスされる。リクエストコードがPIPE__REQ__MIBであれば、リクエストコードワードのすぐ後にタキシードFML32バッファが続く。FML32は、フィールド化されたバッファのタイプである。フィールド化されたバッファは、フィールドと呼ばれる属性値対を含む。属性はフィールドの識別子であり、関連付けられた値は、フィールドのデータ内容を表わす。フィールド化されたバッファと関連付けられるデータのすべては独立している（すなわち、フィールド化されたバッファの外のデータへのポインタがない）。これは、フィールド化されたバッファが、独自のプロセス間で、異なるコンピュータ上で実行するプロセス間でさえ、伝送可能であることを意味する。タキシードは、FML32バッファにデー

40

50

タを割当て、組立て、加算し、かつそこからデータを抽出するのに必要とされる関数ライブラリを提供する。タキシードアドミンAPIは、MIBアクセスリクエストを実行するために必要とされるFML32リクエストバッファの特定のフォーマットを定義する。この実施例では、PIPE_REQ_MIBリクエストタイプは、“GET”または“SET”のTA_OPERATIONで、いかなる有効なMIBアクセスリクエストをもサポートする。“GETNEXT”操作は、この実施例ではサポートされないが、他の実施例で可能である。

【0044】

PIPE_REQ_ADM_ONLINEおよびPIPE_REQ_ADM_OFFLINEなどの、PIPE_REQ_MIB以外のすべてのリクエストコードについて、リクエストメッセージは、リクエストコードおよび、リクエストコードによって特定されるリクエストの特定のタイプに関連するいかなる他の情報をも含むデータ構造を含む。たとえば、以下のC言語の型定義が採用可能である：

【0045】

【数2】

```
typedef struct _pipe_request_t {
    DWORD    requestCode;
    DWORD    packetRevision;
    DWORD    cmdInputSize;
    BOOL     returnCmdOutput;
    char     cmd[MAX_TUX_CMD_LEN+1];
    char     opts[MAX_TUX_OPTS_LEN+1];
    char     groupName[MAX_GROUPNAME_LEN+1];
} PIPEREQ;
typedef PIPEREQ      *LPPIPEREQ;
typedef const PIPEREQ *LPCPIPEREQ;
```

【0046】

requestCodeフィールドは、送られるリクエストメッセージのタイプを特定する特定のリクエストコードを含む。この実施例では、packetRevisionフィールドは、将来の使用のために確保され、単に0にセットされる。残りのフィールドは、この発明の一部ではないパイプサーバ42の他の機能に関連する。

【0047】

PIPE_REQ_ADM_ONLINEのリクエストコードは、MSCSのために、(パイプサーバが実行している仮想サーバに対応する)タキシードアドミンサーバをオンラインにするリクエストを表わす。この場合には、リクエスト構造内のすべての他のフィールドは、0にセットされる。PIPE_REQ_ADM_OFFLINEのリクエストコードは、MSCSのために、(パイプサーバが実行している仮想サーバに対応する)タキシードアドミンサーバをオフラインにするリクエストを表わす。この場合もやはり、リクエスト構造内のすべての他のフィールドは、0にセットされる。

【0048】

パイプサーバAPIはまた、特定のリクエストの結果を表示するのに用いられる応答メッセージのフォーマットを特定する。リクエストコードがPIPE_REQ_MIBであった場合、応答メッセージはタキシードFML32バッファである。タキシードアドミンAPIは、MIBアクセスリクエストを実行したとき戻されるFML32応答バッファの特定のフォーマットを定義する。PIPE_REQ_MIB以外のすべてのリクエストコードについて、応答メッセージフォーマットは、以下のCの型定義によって定義可能である：

【0049】

【数3】

```

typedef struct _pipe_reply_t {
    DWORD    packetRevision;
    DWORD    cmdOutputSize;
    int      status;
#define PIPE_REP_GOOD_STATUS          0    // See statuses below:
                                         // Success
#define PIPE_REP_ERROR_STATUS        -1    // General error
#define PIPE_REP_MASTER_UNAVAILABLE -2    // Tuxedo Master was
                                         // unavailable, so
                                         // unable to accurately
                                         // determine
                                         // looks alive status.
#define PIPE_REP_NO_RESPONSE        -3    // No reply was received
                                         // a broken pipe,
                                         // timeout, etc.

} PIPEREP;
typedef PIPEREP    *LPPIPEREP;
typedef const PIPEREP    *LPCPIPEREP;

```

10

【 0 0 5 0 】

packetRevisionフィールドは、起こり得る将来の使用のために確保され、この実施例では常に0にセットされる。cmdOutputSizeフィールドは、この発明の一部ではないパイプサーバ42の特徴に関する。statusフィールドは、リクエストの結果を示す値を保持する。statusフィールドは、パイプ切断エラーまたはタイムアウトのためにパイプサーバから応答が戻されない場合、PIPE_REP_NO_RESPONSEにセットされる。PIPE_REQ_ADM_ONLINEリクエストに対して、statusフィールドは、タキードアドミンサーバがうまくオンラインにされた場合にはPIPE_GOOD_STATUSにセットされ、または失敗の場合にはPIPE_ERROR_STATUSにセットされる。PIPE_REQ_ADM_OFFLINEリクエストに対して、statusフィールドは、タキードアドミンサーバがうまくオフラインにされた場合にはPIPE_GOOD_STATUSにセットされ、または失敗の場合にはPIPE_ERROR_STATUSにセットされる。PIPE_MASTER_UNAVAILABLEは、この発明の一部ではない特徴に関連して用いられるステータスコードである。

20

【 0 0 5 1 】

2. パイプサーバAPI - アクセス関数

上述したように、この実施例では、Windows NT名前付きパイプは、クライアントプロセス（たとえば、リソースDLL40）とパイプサーバ42との間、または異なったノード上の2つのパイプサーバ42間でメッセージをやり取りするために用いられる。所与のノード上のパイプサーバ42からサービスを要求するために、リソースDLL40などのクライアントプロセス、または別のノード上のパイプサーバ42は、以下のリスティング1に示されるステップを実行する：

30

リスティング1：

- (a) WIN32 CreateFile関数を呼出して、読出および書込の両方のために名前付きパイプを開き、
- (b) WIN32 SetNamedPipeHandleState関数を呼出して、（バイトストリームモードとは反対の）メッセージモードに名前付きパイプをセットし、
- (c) WIN32 WriteFile関数を呼出して、リクエストバッファをサーバに送り、
- (d) WIN32 ReadFile関数を呼出して、サーバから応答バッファを読出し、
- (e) WIN32 CloseHandle関数を呼出して、名前付きパイプを閉じる。

40

【 0 0 5 2 】

割込を見込んでWindows NT Overlapped I/Oがリスティング1のステップ（c）および（d）に用いられてもよい。

【 0 0 5 3 】

所与のノード上のクライアントによって、その同じノード上または異なったノード上のパ

50

パイプサーバ42にメッセージを送るために用いられる名前付きパイプの名前（クライアント側の名前）は、以下の形をとる：

```
\\serverComputerName\pipe\tuxpipeserver.serverComputerName
```

ただし、serverComputerNameは、（WIN32 GetComputerName関数によって戻される）パイプサーバが実行しているコンピュータの名前である。M S C S環境では、これは、パイプサーバ42が実行しているM S C S仮想サーバのネットワーク名値である。

【0054】

名前付きパイプのインスタンスを生成するとき、パイプサーバ42によって用いられる名前付きパイプの名前（サーバ側の名前）は、以下の形をとる：

```
\\.pipe\tuxpipeserver.serverComputerName
```

なお、サーバが名前付きパイプを生成するとき、名前の最初にコンピュータ名の代わりにピリオド（.）の使用が必要とされる。したがって、名前付きパイプの名前のserverComputerName部分は、2つ以上の仮想サーバが同じ物理的コンピュータ上で実行している場合を処理するよう名前付きパイプの名前の最後の部分に特定される。

【0055】

パイプサーバ42は、デフォルトセキュリティを用いて名前付きパイプを生成する。パイプサーバ42は、M S C Sサービスを起動するために用いられるWindows NTドメインアカウント下で実行する。これは、この同じアカウント（またはサーバシステムの管理者特権を備えるアカウント）を用いるクライアントのみが名前付きパイプを開くために必要とされるアクセス権を有することを意味する。すべての他のクライアントには、CreateFile関数でパイプのクライアント側を開くと、ERROR_ACCESS_DENIEDエラーコード（5）が戻される。

【0056】

クライアント側の処理を簡素化するために、上のステップ（a）から（e）に記載される機能をカプセル化するアクセス関数が与えられる。関数の各々は、IntrPtrパラメータを有する。パイプサーバ42にリクエストを送り応答を待つことは、時間のかかる動作であり得るので、IntrPtrパラメータを用いて、操作の割込を制御する。IntrPtrがNULLであれば、要求された操作は割込可能でない。IntrPtrがNULLでなければ、これは、以下に定義される構造を指す：

【0057】

【数4】

```
// The structure type INTR is used to control interrupting long
// running
// operations.
typedef struct _interrupt_t {
    LPVOID workerPtr;
    // workerPtr is the pointer for this worker thread, which is
    // required
    // by the ClusWorkerCheckTerminate function, or NULL if
    // ClusWorkerCheckTerminate should not be called.
    time_t stopTime;
    // stopTime is the time (as returned by the "time" function)
    // when
    // this operation should be interrupted, or 0 if there is no
    // time
    // limit.
} INTR;
typedef INTR *LPINTR;
```

【0058】

workerPtrフィールドがNULLでなければ、これは、クラスタユーティリティ関数ClusWorkerCreateでワークスレッドを生成するとき用いられるワーカ構造へのポインタである。クラスタユーティリティ関数は、M S C Sの一部として与えられる。ワークスレッドの使用は、以下により詳細に記載される。workerPtrがNULLでなければ、アクセス関数

は周期的に、クラスタユーティリティ関数ClusWorkerCheckTerminateを呼出して、ワークスレッドが終了するよう要求されたかどうかを決定する。ClusWorkerCheckTerminateがTRUEを戻した場合、アクセス関数は、進行中の操作を打ち切り、エラーを戻す。

【 0 0 5 9 】

stopTimeフィールドが0でなければ、それは、Cのランタイム時間関数によって戻されるのと同じ形(すなわち、1970からの秒数)で、この操作のデッドライン時間を表わす。stopTimeが0でないとき、アクセス関数は周期的に、デッドライン時間が来たかどうかを見るようチェックする。もしそうであれば、アクセス関数は操作を打ち切り、エラーを戻す。

【 0 0 6 0 】

関数の各々は、numRetriesパラメータを有する。0が送られた場合、要求された操作は自動的に再試行されない。numRetriesが0でなく、応答がパイプサーバから受取られないようなエラーが生じた場合(たとえば、パイプサーバが実行していないためにパイプサーバを開くことができない)、アクセス関数は1秒スリープし、自動的に、numRetriesによって示される回数まで同じリクエストを再試行する。numRetries再試行を実行した後リクエストがなおもフェールするのであれば、アクセス関数はエラーを戻す。

【 0 0 6 1 】

これらの関数の各々はまた、serverComputerNameパラメータを有する。このパラメータは、(WIN32 GetComputerName関数によって戻される)パイプサーバが実行しているコンピュータの名前である。M S C S環境では、これは、パイプサーバが実行しているM S C S仮想サーバのネットワーク名値である。他の指定がない限り、このパラメータは指定されなければならない、NULLであってはならない。

【 0 0 6 2 】

以下のアクセス関数は、前述の、PIPE__REQ__MIB、PIPE__REQ__ADM__ONLINEおよびPIPE__REQ__ADM__OFFLINEメッセージタイプに対して用いられる。他の指定がなければ、関数の各々は、成功に対しては整数()を戻し、またはエラーが生じたならば1を戻す。

【 0 0 6 3 】

getFromMib()

目的 :

タキシードMIBから所与のオブジェクトクラスの1つ以上の属性値を取出すことである。

【 0 0 6 4 】

Function Call :

【 0 0 6 5 】

【 数 5 】

```
int getFromMib(
    LPINTR IntrPtr,
    DWORD numRetries,
    LPCSTR serverComputerName,
    LPCSTR className,
    FLDID32 keyAttr,
    LPCVOID keyValue,
    ...
);
```

【 0 0 6 6 】

Parameters :

serverComputerName - NULLであってもよい。その場合、リクエストはパイプサーバに提出されず、代わりにリクエストは現在のプロセス内で実行される。

10

20

30

40

50

【 0 0 6 7 】

className - M I B クラス名を特定するために用いられる文字列である。

keyAttr - “ GET key ” として使用すべき属性を表わす F L D I D 3 2 値である。(クラスが T _ D O M A I N である場合のように) “ GET key ” がなければ、これを B A D F L D I D にセットする。

【 0 0 6 8 】

keyValue - “ GET key ” 値へのポイドポインタである ; この値は、keyAttr が B A D F L D I D であれば無視される。

【 0 0 6 9 】

getAttr1 - 取出すべき属性を表わす F L D I D 3 2 値である。 10

getValue1 - 取出された属性値を保持するのに十分大きいバッファへのポイドポインタである。

【 0 0 7 0 】

可変数のさらなる getAttrN パラメータおよび getValueN パラメータがある。最後の getAttrN は、 B A D F L D I D の値によって表示される。

【 0 0 7 1 】

setInMib ()

目的 :

タキシード M I B 内に所与のオブジェクトクラスのための 1 つ以上の属性値をセットする。 20

【 0 0 7 2 】

Function Call :

【 0 0 7 3 】

【 数 6 】

```
int setInMib(
    LPINTR IntrPtr,
    DWORD numRetries,
    LPCSTR serverComputerName,
    long flags,
    LPCSTR className,
    FLDID32 keyAttr,
    LPCVOID keyValue,
    ...
);
```

30

【 0 0 7 4 】

Parameter :

serverComputerName - N U L L であってもよい。その場合、リクエストはパイプサーバに提出されず、代わりにリクエストは現在のプロセス内で実行される。

【 0 0 7 5 】

className - M I B クラス名を特定するために用いられる文字列である。 40

flags - 使用すべき D A _ F L A G S 値である。

【 0 0 7 6 】

keyAttr - “ SET key ” として使用すべき属性を表わす F L D I D 3 2 値である。(クラスが T _ D O M A I N である場合のように) “ SET key ” がなければ、これを B A D F L D I D にセットする。

【 0 0 7 7 】

keyValue - “ SET key ” 値へのポイドポインタである ; この値は、keyAttr が B A D F L D I D であれば無視される。

【 0 0 7 8 】

setAttr1 = > セットすべき属性を表わす F L D I D 3 2 値である。 50

setValue1 = > セットすべき属性値へのポインタである。

【 0 0 7 9 】

可変数のさらなる setAttrN パラメータおよび setValueN パラメータがある。最後の setAttrN は、BADFLDID の値によって表示される。

【 0 0 8 0 】

processSimpleRequestViaPipe ()

目的 :

PIPEREQ 構造体を組立て、これを、処理のために特定されたサーバコンピュータに送る。この関数は、この実施例では、PIPE_REQ_ADM_ONLINE リクエストおよび PIPE_REQ_ADM_OFFLINE リクエストに対して用いられる。

10

【 0 0 8 1 】

Function Call :

【 0 0 8 2 】

【 数 7 】

```
int processSimpleRequestViaPipe(
    LPINTR IntrPtr,
    DWORD numRetries,
    LPCSTR serverComputerName,
    DWORD requestCode,
    LPCSTR groupName,
    BOOL logPipeNotFound,
    BOOL logBrokenPipe
);
```

20

【 0 0 8 3 】

Parameters :

requestCode - PIPEREQ 構造体に対するリクエストコードである。

【 0 0 8 4 】

groupName - この発明の一部でないリクエストタイプに対する groupName 値である。このフィールドは、PIPE_REQ_ADM_ONLINE リクエストタイプおよび PIPE_REQ_ADM_OFFLINE リクエストタイプに対しては NULL である。

30

【 0 0 8 5 】

logPipeNotFound - パイプ・ノット・ファウンド・エラー (ERROR_FILE_NOT_FOUND) がログされるべきであれば、TRUE であるフラグである。

【 0 0 8 6 】

logBrokenPipe - パイプ切断エラー (ERROR_BROKEN_PIPE) がログされるべきであれば、TRUE であるフラグである。

【 0 0 8 7 】

出力 :

この関数は、以下の値の 1 つを戻す :

40

PIPE_REP_GOOD_STATUS - 成功。

【 0 0 8 8 】

PIPE_REP_ERROR_STATUS - 一般的エラー。

【 0 0 8 9 】

PIPE_REP_MASTER_UNAVAILABLE - タキシードマスタが利用不可能であり、ステータスを正確に判断することができなかった。このステータスコードは、この発明の一部でない特徴で使用される。

【 0 0 9 0 】

PIPE_REP_NO_RESPONSE - 切断パイプ、タイムアウトなどのために応答が受取られなかった。

50

【 0 0 9 1 】

3 . パイプサーバアーキテクチャ

パイプサーバは、WIN32 コンソールアプリケーションとして実現される。スタートアップ時に、以下のものがコマンドラインパラメータとして送られる：

TUXCONFIG Path - サービスされるべきタキシード論理的マシンに対応する T U X C O N F I G ファイルのパスネーム全体である。このパラメータは必須である。

【 0 0 9 2 】

Trace Level - サーバのために活性化されるべきトレースのレベルである（オフに対しては 0、またはトレース詳細のレベルを増大させるには 1 から 4）。このパラメータはオプションである；省略されれば、トレーシングはオフにされる。

10

【 0 0 9 3 】

Number of Pipe Instances - 生成されるパイプインスタンスの数である。このパラメータはオプションである；省略されれば、デフォルト値は 1 0 である。このパラメータの最大値は 6 3 である。名前付きパイプの生成は、以下により詳細に説明される。

【 0 0 9 4 】

Worker Thread Inactivity Timeout - 特定のワークスレッドが終了する前にそのワークスレッドへのいかなるサービスリクエストもなしに経過しなければならない時間の量（ミリ秒）である。このパラメータはオプションである；省略されれば、デフォルト値は 6 0 0 0 0（すなわち、1分）である。ワークスレッドの使用は、以下により詳細に記載される。

20

【 0 0 9 5 】

パイプサーバ 4 2 は、名前付きパイプの複数のインスタンスを生成可能である。インスタンスは、同じ名前を備えた名前付きパイプの独自のコピーである。各インスタンスを用いて、異なったクライアントとの接続を同時に維持することができる。

【 0 0 9 6 】

パイプサーバ 4 2 のメインスレッドは、Number of Pipe Instancesパラメータによって特定されるとおり固定数のパイプインスタンスを生成する。メインスレッドは、Overlapped I/O操作を用いてパイプインスタンス上の接続操作を多重化する。Overlapped I/Oでは、接続操作は、WIN32 ConnectNamedPipe関数へのコールで起動される。もし操作が即座に完了しなければ、ConnectNamedPipe関数は、E R R O R _ I O _ P E N D I N G のステータスを戻す；接続が完了すると、オーバラップした構造体内のイベントに信号が与えられる。

30

【 0 0 9 7 】

各パイプインスタンスは、以下の C 言語構造によって表わされる：

【 0 0 9 8 】

【 数 8 】

```
typedef struct {
    OVERLAPPED          connectOverlap;
    HANDLE              hPipeInst;
    CLUS_WORKER        worker;
    HANDLE              workerActivateEvent;
    BOOL volatile      workerActive;
    CRITICAL_SECTION   workerCriticalSection;
    DWORD               dwState;
#define CONNECTION_PENDING_STATE    0
#define PROCESSING_STATE            1
#define RECONNECTING_STATE          2
} PIPEINST, *LPPIPEINST;
```

40

【 0 0 9 9 】

O V E R L A P P E D、H A N D L E、B O O L、C R I T I C A L _ S E C T I O N お

50

よび D W O R D は、Win32によって定義されるデータ型である。C L U S _ W O R K E R は、M S C S によって定義されるデータ型である。P I P E I N S T 構造体のフィールドは、以下のように定義される：

connectOverlap - オーバラップした接続を行うために必要とされるオーバラップした構造体である。この構造体は、オーバラップした接続が完了したとき、またはワークスレッドがサービスリクエストを完了したときを表示するよう用いられるイベントハンドルを含む。

【 0 1 0 0 】

hPipeInst - このパイプインスタンスのためのハンドルである。

worker - M S C S クラスユーティリティ関数の組によってワークスレッドを制御するために用いられるワーカ構造体である。 10

【 0 1 0 1 】

workerActivateEvent - 処理すべきサービスリクエストがあるときワークスレッドに通知するために用いられるイベントである。

【 0 1 0 2 】

workerActive - ワークスレッドが現在実行しているかどうかを表示するフラグである。

【 0 1 0 3 】

workerCriticalSection - ワークスレッドの起動および終了を調整するために用いられるクリティカルセクションである。クリティカルセクションは、Microsoft Windows NTによって定義されるオブジェクトである。 20

【 0 1 0 4 】

dwState - パイプインスタンスの現在の状態である。

C O N N E C T I O N _ P E N D I N G _ S T A T E は、ConnectNamedPipeが進行中であることを表わす。P R O C E S S I N G _ S T A T E は、クライアントと接続が確立され、パイプサーバ42がワークスレッド内のリクエストを処理する準備ができていないことを表わす。R E C O N N E C T I N G _ S T A T E は、リクエストがワークスレッドによって処理され、パイプサーバ42が、クライアントから切断され、新しいクライアントとの新しい接続を開始する準備ができていないことを意味する。

【 0 1 0 5 】

パイプサーバ42は、マルチスレッディングを用いて、それが生成する名前付きパイプを介して受取られた別個のリクエストの各々を処理する。そのようなスレッドは各々、ここではワークスレッドと呼ばれる。より詳細には、クライアントがメインスレッド内の名前付きパイプインスタンスに接続した後、リクエストは別個のワークスレッド内で処理される。これは、メインスレッドが即座に他のパイプインスタンスのサービスを再開することを可能にする。サービスリクエストは、完了するのに数秒かかることがある。したがって、メインスレッド内でサービスリクエストを実行することは、実際的でない。 30

【 0 1 0 6 】

各パイプインスタンスは、対応するワークスレッドを有する。しかしながら、システムリソースを最適化するために、ワークスレッドは、リクエストが初めて特定のパイプインスタンスについて処理されるまで、起動されない。これは、10個のパイプインスタンスがあって、クライアントからのサービスリクエストの同時レベルが6個を決して超えない場合、最後の4個のパイプインスタンスについてのワークスレッドは決して起動されないということを意味する。 40

【 0 1 0 7 】

ワークスレッドが一旦起動されると、それはさらなるサービスリクエストを待つループする。さらなるサービスリクエストが、パラメータWorker Thread Inactivity Timeoutによって定義されるタイムアウト時間期間内に特定のパイプインスタンスについて生じなければ、ワークスレッドは終了する。ワークスレッドは、必要なときに自動的に再起動される。

【 0 1 0 8 】

4. メインスレッドおよびワークスレッド処理

パイプサーバ42のメインスレッドは、リスティング2に示す処理を実行する：

リスティング2：

1) serverActiveフラグをTRUEに初期化する。(パイプサーバ42を停止するリクエストがなされたとき、このフラグはFalseにセットされる)

2) パイプサーバ42コマンドラインパラメータを処理する。

【0109】

3) パイプインスタンスを表わすPIPEINST構造体のアレイ(パイプと呼ばれる)を割当てる。

【0110】

4) ハンドル(hEventsと呼ばれる)の配列を割当てる。配列内のエントリ数は、パイプインスタンスの数よりも1つ大きい。

【0111】

5) パイプサーバ42の終了を制御するために用いられるイベントを生成する。このイベントハンドルは、hEvents配列の最終部分に記憶される。このイベントには、パイプサーバ42を停止するリクエストがなされたとき信号が与えられる。

【0112】

6) 各パイプインスタンスについて：

a) connectOverlap構造のためにイベントを生成し(かつhEventsにハンドルのコピーを記憶する)。

【0113】

b) ワーク活性化のためにイベントを生成する。

c) ワーク活性化のために使用されるクリティカルセクションを初期化する。

【0114】

d) パイプインスタンスを生成する。

e) WIN32関数ConnectNamedPipeを呼出してクライアントからの接続に対してリスニングを開始する。

【0115】

f) If ConnectNamedPipeがERROR__IO__PENDINGを戻した場合：

i) 状態をCONNECTION__PENDING__STATEに変える。

【0116】

g) Else if ConnectNamedPipeがERROR__PIPE__CONNECTEDを戻した場合：

i) 状態をPROCESSING__STATEに変える。

【0117】

7) serverActiveフラグがTRUEである間ループする：

a) hEvents配列内のハンドル上で待つWIN32関数WaitForMultipleObjectsを呼出す。これは、すべき何らかのワークを待つ。

【0118】

b) 信号を与えられたイベントがhEvents内の最終イベントであれば、パイプサーバ42は終了する。さもなければ、信号を与えられた特定のイベントは、特定のパイプインスタンスに対応する。

【0119】

c) If CONNECTION__PENDING__STATE:

i) WIN32関数GetOverlappedResultを呼出して、完了された、オーバーラップしたConnectNamedPipe操作の結果を得る。

【0120】

ii) PROCESSING__STATEに変える。

d) If PROCESSING__STATE:

i) connectOverlap構造内のイベントのためにWIN32関数ResetEventを呼出す。ワークスレッドは、それがサービスリクエストの処理を完了したときこのイベントに信号を与え

10

20

30

40

50

る。

【 0 1 2 1 】

ii) このパイプインスタンスのためにActivateWorkerThread(下記)を呼出す。

【 0 1 2 2 】

iii) RECONNECTING__STATEに変える。

e) Else if RECONNECTING__STATE:

i) WIN32関数DisconnectNamedPipeを呼出す。

【 0 1 2 3 】

ii) WIN32関数ConnectNamedPipeを呼出して、クライアントからの接続に対してリスニングを開始する。

【 0 1 2 4 】

iii) If ConnectNamedPipeがERROR__IO__PENDINGを戻した場合:

(1) 状態をCONNECTION__PENDING__STATEに変える。

【 0 1 2 5 】

iv) Else if ConnectNamedPipeがERROR__PIPE__CONNECTEDを戻した場合:

(1) 状態をPROCESSING__STATEに変える。

【 0 1 2 6 】

ActivateWorkerThread関数は、それがまだ実行していなければワークスレッドを起動し、または既存のワークスレッドを活性化する。入力パラメータは、PIPEINST構造へのポインタである。関数は、リスティング3に示すステップを実行する。

【 0 1 2 7 】

リスティング3:

1) workerCriticalSectionのためにWIN32関数EnterCriticalSectionを呼出す。

【 0 1 2 8 】

2) If workerActiveフラグがTRUE:

a) workerActivateEventに信号を与える。

【 0 1 2 9 】

b) workerCriticalSectionのためにWIN32関数LeaveCriticalSectionを呼出す。

【 0 1 3 0 】

3) Else:

a) workerCriticalSectionのためにLeaveCriticalSectionを呼出す。

【 0 1 3 1 】

b) クラスターティリティ関数ClusWorkerTerminateを呼出して、いかなる先行のワークスレッドもが完全に終了されたことを確実にする。

【 0 1 3 2 】

c) クラスターティリティ関数ClusWorkerCreateを呼出して、ワークスレッド実行関数WorkerThreadFunction(下記)を生成する。

【 0 1 3 3 】

d) workerActiveフラグをTRUEにセットする。

e) workerActivateEventに信号を与える。

【 0 1 3 4 】

TerminateWorkerThread関数は、サーバのシャットダウン中に呼出され、ワークスレッドの終了を強制する。入力パラメータは、PIPEINST構造体へのポインタである。serverActiveフラグが既にFALSEにセットされていると仮定して、関数はリスティング4に示すステップを実行する。

【 0 1 3 5 】

リスティング4:

1) workerActivateEventに信号を与えて、ワークスレッドが待っていないということを確認する。

【 0 1 3 6 】

10

20

30

40

50

2) クラスターユーティリティ関数ClusWorkerTerminateを呼出す。

WorkerThreadFunctionは、クラスターユーティリティ関数ClusWorkerCreateが呼出されたとき、新しいワークスレッド内で呼出される。WorkerThreadFunction関数は、リスティング5に示すステップを実行する。

【0137】

リスティング5：

1) workerActivateEvent上で待つWIN32関数WaitForSingleObjectを呼出す。タイムアウトの長さは、サーバパラメータWorker Thread Inactivity Timeoutによって与えられる値である。

【0138】

2) If waitForSingleObjectがタイムアウトした場合：

a) workerCriticalSectionのためにEnterCriticalSectionを呼出す。

【0139】

b) WorkerActivateEvent上で待つWaitForSingleObjectを再び呼出す。タイムアウトの長さは0であり、即時にリターンする。

【0140】

c) If waitForSingleObjectが再びタイムアウトした場合：

i) workerActiveフラグをFALSEにセットする。

【0141】

ii) workerCriticalSectionのためにLeaveCriticalSectionを呼出す。

iii) リクエストおよび応答のために割当てられるF M L 3 2バッファを解放する。

【0142】

iv) WorkerThreadFunctionからリターンし、これはワークスレッドを終わらせる。

【0143】

d) workerCriticalSectionのためにLeaveCriticalSectionを呼出す。

3) If serverActive flagがFALSE:

a) WorkerThreadFunctionからリターンし、これがワークスレッドを終了させる。

【0144】

4) 名前付きパイプインスタンスから一時バッファにリクエストを讀出す。

5) If リクエストコードがPIPE__PEQ__MIB:

a) まだ割当てられていなければ、リクエストおよび応答のためにF M L 3 2バッファを割当てる。

【0145】

b) (第2のワードで開始する)一時バッファをリクエストF M L 3 2バッファにコピーする。

【0146】

c) M I Bリクエストを処理する(以下により詳細に記載する)。

d) 応答F M L 3 2バッファを名前付きパイプインスタンスに書込む。

【0147】

e) F M L 3 2リクエストおよび応答バッファを再使用のためにクリアする。

【0148】

6) Else:

a) 一時バッファをP I P E R E Q構造体にキャストする。

【0149】

b) requestCodeフィールドに従ってリクエストを処理する(以下により詳細に記載する)。

【0150】

c) P I P E R E Q構造体を名前付きパイプインスタンスに書込む。

7) クライアントが処理を完了しパイプのクライアント側を閉じるのを待つことだけを目的に名前付きパイプインスタンスから讀出す。クライアントがパイプを閉じると、この

10

20

30

40

50

読出は、ERROR__BROKEN__PIPEエラーを受取る。

【0151】

8) connectOverlap構造体内のイベントに信号を与えて、処理が完了したことを表示する。

【0152】

5. エラーロギング

すべてのエラーログは、タキシードULOGに書込まれる。以下の条件がログされる：

a) WIN32API関数から戻されるすべての予期せぬエラー；(b) クラスタAPI関数から戻されるすべての予期せぬエラー；および(c) タキシードアドミンAPI関数から戻されるすべての予期せぬエラー。

【0153】

6. 詳細な操作

図5A～図5Bは、図4に例示される例示のアプリケーションにおけるこの発明の操作のより詳細な例示をする流れ図を含む。特に、図5A～図5Bは、ONLINEリクエストの処理を例示する。

【0154】

ステップ100で、リソースモニタ26は、リソースDLL40内のタキシードアドミンサーバ・リソースタイプのONLINEエントリポイント関数を呼出す。ONLINEエントリポイント関数は、新しいスレッドを起動して、リクエストを処理し、ERROR__IO__PENDINGステータスをリソースモニタ26に戻す。

【0155】

次に、ステップ101で、リソースDLL40は、processSimpleRequestViaPipe()関数を呼出し、PIPE__REQ__ADM__ONLINEリクエストコードを渡す。ProcessSimpleRequestViaPipe()関数は、適切なリクエストメッセージを構成し(すなわち、PIPE__REQ__ADM__ONLINEのリクエストコードでPIPEREQ構造体を組立てる)、このノード(すなわち、親ノード)上のパイプサーバ42のために名前付きパイプの正しい名前を構成する。ステップ102で、リソースDLL40は、上のListing1に示す初期メッセージ交換ステップを実行する。

【0156】

ステップ104で、パイプサーバ42のメインスレッドは、リソースDLL40から接続を受取り、ワークスレッドを活性化してリクエストを処理する。ワークスレッドは、PIPE__REQ__ADM__ONLINEリクエストコードを認識し、リクエストの処理を開始する。ステップ106で、ワークスレッドは、現在のノード(すなわち親ノード)がタキシードアプリケーションのための指定されたマスタかどうかを判断する。具体的には、ワークスレッドは、タキシードアプリケーションのT__DOMAINオブジェクトからTA__MASTER属性を獲得するgetMasterIds()関数(図示せず)を呼出す。TA__MASTER属性は、タキシードアプリケーションの現在のところ指定されているマスタの論理的マシンID(LMID)を特定する。マスタのLMIDを使用して、マスタの物理的マシンID(PMID)は、タキシードT__MACHINEオブジェクトのTA__PMID属性から獲得可能である。次に、ワークスレッドは、Win32関数GetComputerNameを呼出すgetLmid()関数(図示せず)を呼出して、親ノードの物理的マシンID(PMID)を決定する。それは次に、PMIDに対応するタキシードT__MACHINEオブジェクトからTA__LMID属性を獲得する。リクエストはこの例ではONLINEリクエストであるので、ワークスレッドは、現在の状態がINACTIVEであるということをもT__MACHINEオブジェクトのTA__STATE属性からベリファイする。

【0157】

ステップ106で、親ノードがマスタノードであるということが、LMIDおよびPMIDから決定されると、PIPE__REQ__ADM__ONLINEリクエストは、ステップ108でローカルに処理される。そのような場合には、ワークスレッドはまず、バックアップマスタノードが指定されているかどうかを決定する。もしそうなら、パイプサーバは

10

20

30

40

50

、バックアップマスタ上のパイプサーバにリクエストを送り、バックアップマスタにマスタとして自身を再指定するよう要請する。もしこれが成功であれば、ワークスレッドは次に、名前付きパイプ接続を介して新しいマスタ上のパイプサーバにオンラインリクエストを転送して、それが、親ノード（これはたとえば、失敗後に再起動しようとしているかもしれない）をオンラインにし直すためにリクエストを処理できるようにする。バックアップマスタがなければ、親ノード上のパイプサーバは、タキシードアドミンAPIを直接呼出して、そのノード上でタキシードドメインをオンラインにする。

【0158】

ステップ106で、親ノードがマスタノードでないことが決定されれば、制御はステップ110に移る。ステップ110で、ワークスレッドは、パイプサーバアクセス関数setInMIB()を呼出して、FML32バッファが親ノードのLMIDと関連付けられたT_MACHINEオブジェクト内のTA_STATE属性をACTIVEにセットするリクエストを含むPIPE_REQ_MIBメッセージを生成する（これは、ONLINEリクエストがタキシードアドミンAPIを介してなされる態様である）。setInMIB()関数は次に、マスタノード上のパイプサーバのために名前付きパイプの正しい名前を構成し、初期メッセージ交換ステップ（上記Listing1参照）を実行して、メッセージをマスタノード上のパイプサーバに転送する。

【0159】

ステップ112で、マスタノード上のパイプサーバ42のメインスレッドは、親ノード上のパイプサーバ42から接続を受取り、自分自身のワークスレッドを活性化してPIPE_REQ_MIBリクエストを処理する。ステップ114で、ワークスレッドは、PIPE_REQ_MIBリクエストを認識し、タキシードアドミンAPIに適切なコールを行なってリクエストを処理させ応答を受取る。具体的には、この例では、ワークスレッドは、タキシードアドミン関数tpinit()を呼出してタキシードアプリケーションに接続し、次にタキシードアドミンAPI関数tpacall()とそれに続くtpgetreply()を呼出してそれぞれ、リクエストを提出して処理させ、応答を受取る。

【0160】

次にステップ116で、（成功または失敗を表示する）応答FML32バッファは、同じ名前付きパイプ接続を介して親ノードに送り返される。ステップ118で、親ノード上のsetInMIB()関数は、名前付きパイプから応答を読み出し、パイプを閉じる。ステップ120で、親ノード上のワークスレッドは次に、名前付きパイプ接続を介して親ノード上のリソースDLL40に応答PIPEREPメッセージを書込む。最後に、ステップ122で、リソースDLL40のONLINEエントリポイント関数は、リソースモニタ26の適切なコールバック関数を呼出して、ONLINEリクエストのステータスを表示する。OFFLINEリクエストなどの、タキシードアプリケーションのための指定されたマスタノードによって実行されなければならない他の管理用リクエストの処理は、上記と同様である。

【0161】

「Microsoft」、「Windows」および「Windows NT」は、Microsoft Corporationの登録商標である。「BEA」および「TUXEDO」は、BEA Systems, Inc.の登録商標である。「Oracle」は、Oracle Corporationの登録商標である。前に例示したように、この発明は、アプリケーションがクラスタ化により達成される増大した可用性を利用できるように、マスタノードを介した集中管理を必要とする分散ネットワークアプリケーションがそのような集中管理をサポートしないクラスタ化された環境内に配備されることを可能にするための方法に向けられる。広い発明の概念から逸脱することなしに上述した実施例に変更がされ得ることが理解される。たとえば、パイプサーバ間の接続の好ましい形は、Windows NTの名前付きパイプであるが、プロセス間通信の他の形が採用可能である。加えて、前述したように、この発明を採用してTuxedoベースのアプリケーションがMSS制御のコンピュータクラスタ内で実行することを可能にするような特定の実施例が記載されたが、この発明は、Tuxedoベースのアプリケーションおよび/またはMSS制御のコンピュータクラスタ

10

20

30

40

50

での使用に限られるものでなく、分散ネットワークアプリケーションがマスタノードを介した集中管理を必要とするが、アプリケーションが配備されているクラスタ化された環境がそのような集中管理をサポートしないようないかなるシステムにも採用可能である。したがって、この発明は、開示される特定の実施例に限られるものでなく、前掲の特許請求の範囲によって定義されるようなこの発明の精神および範囲内にあるすべての変形を含むものと意図される。

【図面の簡単な説明】

【図 1】 2つのノードを含む先行技術のコンピュータクラスタの操作を例示するブロック図である。

【図 2】 図 1 のクラスタの操作のフェールオーバーモードを例示するブロック図である。

10

【図 3 A】 この発明の方法の好ましい実施例を例示する流れ図である。

【図 3 B】 この発明の方法の好ましい実施例を例示する流れ図である。

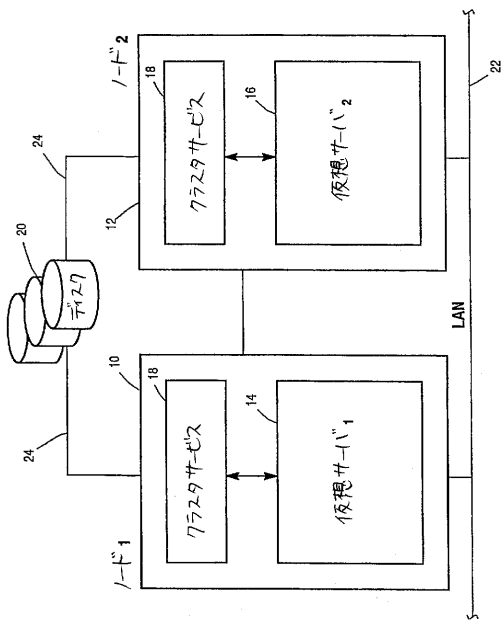
【図 4】 この発明の方法を実施するための装置の好ましい実施例、およびタキシード環境のために書かれかつ M S C S によって制御される計算クラスタ化された計算処理環境内に配備される分散ネットワークアプリケーションを含む環境におけるこの発明の例示的アプリケーションの両方を例示するブロック図である。

【図 5 A】 図 4 の例示的環境におけるこの発明の操作のさらなる詳細を例示する流れ図である。

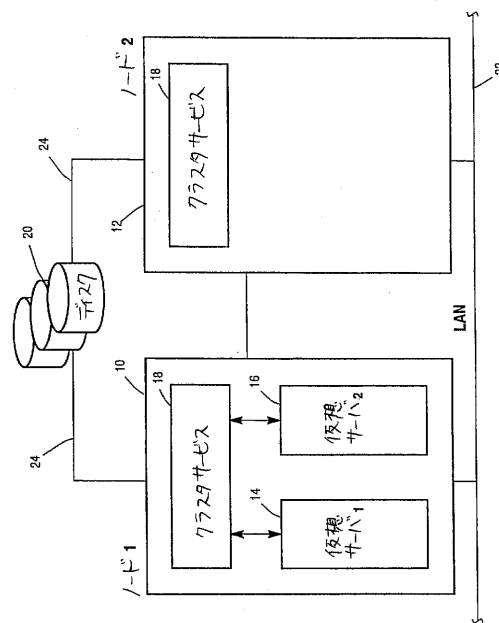
【図 5 B】 図 4 の例示的環境におけるこの発明の操作のさらなる詳細を例示する流れ図である。

20

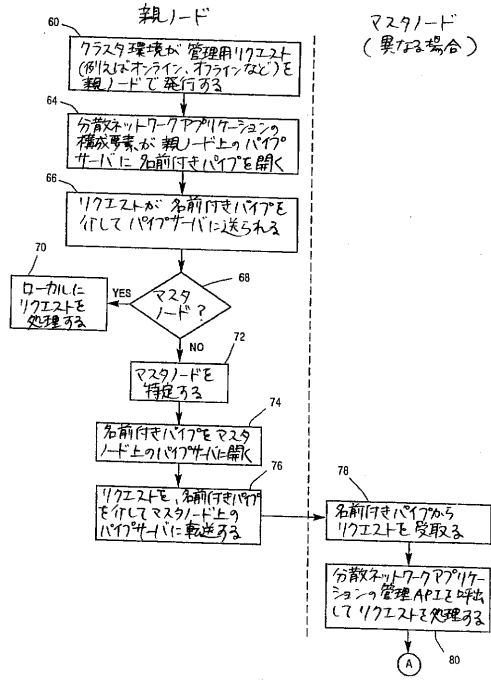
【 図 1 】



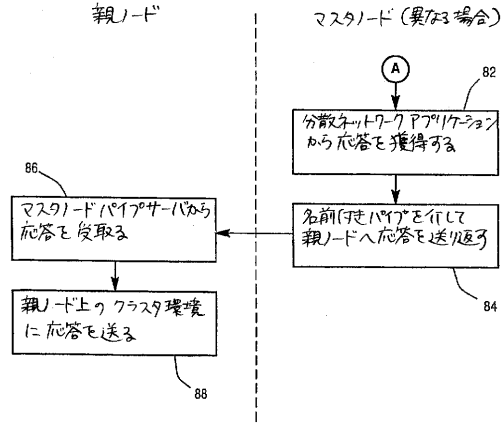
【 図 2 】



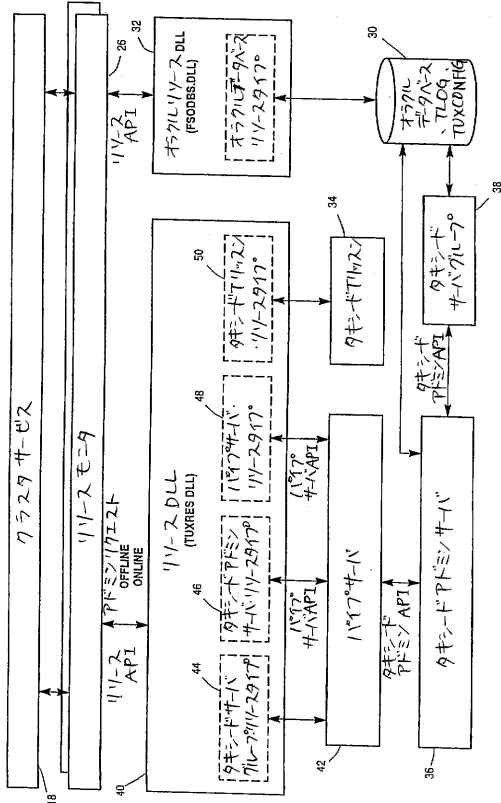
【図3A】



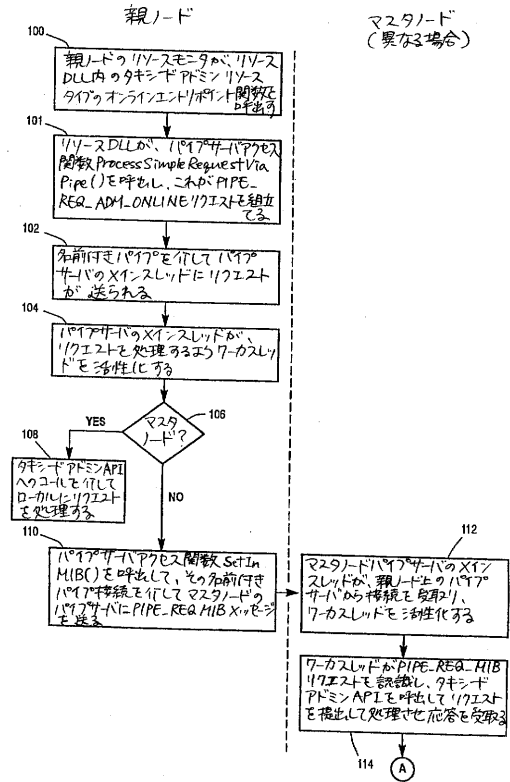
【図3B】



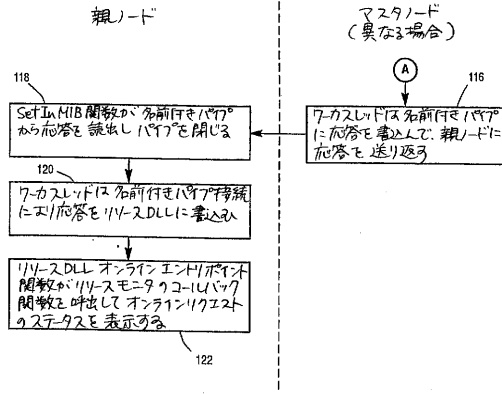
【図4】



【図5A】



【図 5 B】



フロントページの続き

(74)代理人 100091409

弁理士 伊藤 英彦

(74)代理人 100096781

弁理士 堀井 豊

(74)代理人 100096792

弁理士 森下 八郎

(72)発明者 バーデウ, スティーブン・エイ

アメリカ合衆国、19640 ペンシルバニア州、フェニックスビル、デイビッツ・ラン、1008

審査官 殿川 雅也

(56)参考文献 特開昭64-008742(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 9/46 - 9/54