

【公報種別】特許法第 17 条の 2 の規定による補正の掲載

【部門区分】第 6 部門第 3 区分

【発行日】平成 19 年 3 月 8 日 (2007.3.8)

【公表番号】特表 2002-526821 (P2002-526821A)

【公表日】平成 14 年 8 月 20 日 (2002.8.20)

【出願番号】特願 2000-566746 (P2000-566746)

【国際特許分類】

**G 0 6 F 11/20 (2006.01)**

**G 0 6 F 12/00 (2006.01)**

**G 0 6 F 12/16 (2006.01)**

【F I】

G 0 6 F 11/20 3 1 0 C

G 0 6 F 12/00 5 3 1 Z

G 0 6 F 12/00 5 4 5 B

G 0 6 F 12/16 3 1 0 C

G 0 6 F 12/16 3 1 0 J

【手続補正書】

【提出日】平成 18 年 10 月 30 日 (2006.10.30)

【手続補正 1】

【補正対象書類名】明細書

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【書類名】明細書

【発明の名称】複数のファイルサーバ間における永続的状態情報の調整

【特許請求の範囲】

【請求項 1】 ファイルサーバを動作させる方法であって、  
 共用される複数の記憶装置にてなる集合のうちの部分集合を制御するステップと、  
 上記共用される複数の記憶装置を含む通信経路を用いて、第 2 のファイルサーバとの間でメッセージを送受信するステップと、  
 上記通信経路と上記共用される記憶装置とをモニタリングするステップと、  
 上記ファイルサーバについての状態情報を永続的メモリに記憶するステップと、  
 上記モニタリングするステップと、上記永続的メモリの状態とに応答して、上記共用される複数の記憶装置の引き継ぎ動作を実行するステップとを含む、ファイルサーバを動作させる方法。

【請求項 2】 第 1 のサーバにおいて、共用されたリソースの少なくとも一部を管理するステップと、

上記共用されたリソースを用いて、上記第 1 のサーバと第 2 のサーバとの間で複数のメッセージにてなるシーケンスを送受信するステップと、

上記複数のメッセージにてなるシーケンスに応答して、上記共用されたリソースの少なくとも一部に係る上記第 1 のサーバにおける引き継ぎ動作を実行するステップとを含む方法であって、

これにより、上記送受信するステップは、上記第 1 のサーバ及び上記第 2 のサーバの両方が同時に上記引き継ぎ動作を実行することを防止することを特徴とする方法。

【請求項 3】 上記第 1 のサーバにおいて、通信経路に応答して、当該第 1 のサーバ自体の状態及び上記第 2 のサーバの状態を決定するステップと、

上記第 2 のサーバにおいて、上記通信経路に応答して、当該第 2 のサーバ自体の状態及び上記第 1 のサーバの状態を決定するステップとを含んでいて、

これにより、上記第 1 のサーバと上記第 2 のサーバとは、上記第 1 のサーバ及び上記第 2 のサーバのいずれか一方により他方が動作不可能状態にあるとみなされることがないように、それぞれ同時に互いの状態を決定することを特徴とする請求項 2 記載の方法。

【請求項 4】 上記第 1 のサーバについての状態情報を永続的メモリに記憶するステップを含んでいて、

上記第 1 のサーバは、上記永続的メモリの状態に応答してそれ自体の状態を決定する請求項 2 記載の方法。

【請求項 5】 上記第 1 のサーバから、サービス中断からの回復における上記第 1 のサーバの状態に関連する回復情報を送信するステップと、

上記回復情報に応答して上記共用されたリソースの少なくとも一部の返却動作を実行するステップとを含むことを特徴とする請求項 2 記載の方法。

【請求項 6】 上記第 1 のサーバから、サービス中断からの回復における上記サーバの状態に関連する回復情報を送信するステップを含んでいて、

上記引き継ぎ動作を実行するステップは上記回復情報に応答することを特徴とする請求項 2 記載の方法。

【請求項 7】 上記共用されたリソースは複数の記憶装置を含み、

通信経路は上記記憶装置の少なくとも一部を含んでいて、

これにより、上記記憶装置の一部へのアクセスが失われることは、上記通信経路を中断させることを特徴とする請求項 2 記載の方法。

【請求項 8】 上記第 1 のサーバが上記引き継ぎを試みていることを示す少なくとも 1 つのメッセージを、上記第 1 のサーバから上記第 2 のサーバに送信するステップと、

上記メッセージに応答して、上記第 2 のサーバの状態を変更させるステップと、

上記変更された状態においては、上記共用されたリソースへの書き込みを抑制するステップとを含むことを特徴とする請求項 2 記載の方法。

【請求項 9】 通信経路は、上記第 1 のサーバ及び上記第 2 のサーバにてなる対の間に複数の独立した通信経路を含んでいて、

上記方法は、

上記複数のメッセージにてなるシーケンスに対して番号を付与するステップを含み、上記番号を付与するステップは、(a) 上記メッセージの送信者に係るサービス中断及び永続的メモリに응答して生成番号を決定することと、(b) 上記シーケンス中の実質的にそれぞれのメッセージにおいて上記生成番号を提供することとを含み、

各受信者において、上記複数の独立した通信経路のうちの異なるものを用いて配信される複数のメッセージのための統一された順序を決定するステップと、

上記第 1 のサーバにおいて、上記共用されたリソースの状態と上記第 1 のサーバにおける永続的メモリの状態とに응答して、当該第 1 のサーバ自体の状態及び上記第 2 サーバの状態を決定するステップと、

上記シーケンス中の実質的にそれぞれのメッセージを、上記複数の独立した通信経路のうちの少なくとも 2 つにおいて送信するステップとを含み、これにより、上記対の間の通信において単一の障害地点は存在しないことを特徴とする請求項 2 記載の方法。

【請求項 10】 ファイルサーバを制御する方法であって、

複数の大容量記憶装置にてなる集合にアクセスするステップと、

上記複数の大容量記憶装置へのアクセスを有する少なくとも第 2 のファイルサーバとの間で複数のメッセージを伝送するステップと、

上記サーバ及び上記第 2 のファイルサーバについての状態情報を処理するステップとを含み、

上記複数のメッセージは、上記第 2 のファイルサーバとの間で上記状態情報を伝送するために使用され、上記複数のメッセージは、上記複数の大容量記憶装置のうちの少なくとも一部とネットワークとを含む複数の異なる通信経路を介して送信され、上記サーバのうちの一方は、上記複数の大容量記憶装置に係る制御を他方のサーバから引き継ぐことが可能であることを特徴とする、ファイルサーバを制御する方法。

【請求項 1 1】 上記通信経路に含まれた上記複数の大容量記憶装置の一部は、上記複数の大容量記憶装置に格納された 1 つ又は複数のメールボックスをさらに備えることを特徴とする請求項 1 0 記載の方法。

【請求項 1 2】 上記引き継ぎは、上記他方のサーバからのメッセージがタイムアウトしたとき、又は上記他方のサーバからのメッセージが、上記他方のサーバが状態を変化させたことを示すときに実行され、

複数の異なる通信経路に対して、異なるタイムアウトが使用され、

上記複数のメッセージは、上記複数の大容量記憶装置に係る制御の引き継ぎを両方のサーバが同時に実行しようとすることを防止するために使用されることを特徴とする請求項 1 0 記載の方法。

【請求項 1 3】 請求項 1 乃至 1 2 のいずれか 1 つに記載の方法を実行するシステム。

【請求項 1 4】 請求項 1 乃至 1 2 のいずれか 1 つに記載の方法をコンピュータに実行させるプログラムであって、記憶媒体上に記録されたプログラムを有する、コンピュータが読み取り可能な記憶媒体。

#### 【発明の詳細な説明】

##### 【0001】

#### 1. 発明の分野

この発明はコンピュータシステムに関する。

##### 【0002】

#### 2. 従来技術

コンピュータストレージシステムは、データの記録および復元のために使用される。ストレージシステムにより供給されたサービスおよびデータに対し、可能な限り最大限にサービスを利用できることが望ましい。従って、いくつかのコンピュータストレージシステムは複数のファイルサーバを備え、これらのファイルサーバでは、第 1 のファイルサーバに障害が発生したとき、非障害発生時であれば第 1 のファイルサーバにより与えられるサービスおよびデータを提供するために、第 2 のファイルサーバが利用できるようになっている。その第 2 のファイルサーバは、第 1 のファイルサーバにより管理されていたリソースを引き継ぐことにより、これらのサービスおよびデータを提供する。

##### 【0003】

この分野で知られている問題の一つは、2 つのファイルサーバがそれぞれ他方のためにバックアップを提供していたとき、2 つのファイルサーバの各々が他方における障害の発生を確実に検出でき、いかなる引き継ぎ動作が要求されてもそれを円滑に処理できることが重要となる。2 つのファイルサーバのいずれもが正常動作の他のサーバに干渉することなく、上記のことを実行できることが便利である。一つまたは双方のファイルサーバがサービス中断から復旧するとき、この問題はシステムにおいて特に深刻な問題となる。

##### 【0004】

従って、複数の独立したファイルサーバの間での比較的迅速で信頼できる引き継ぎをもたらすためのストレージシステム、およびストレージシステムを動作させるための方法を提供することが都合よい。この利点は、本発明の実施例で達成され、

本発明では、各ファイルサーバは、

(a) 他のファイルサーバへの冗長な複数の通信経路を維持し、

(b) 当該ファイルサーバ自体の状態を永続的メモリに保持し、ここで、各永続的メモリのうちの少なくとも一部は他のファイルサーバからアクセス可能であり、そして、

(c) 他のファイルサーバの状態を定期的に確認する。

##### 【0005】

#### (発明の概要)

本発明は、複数の独立したファイルサーバの間における比較的迅速かつ信頼できる引き継ぎをもたらすためのストレージシステム、およびストレージシステムを動作させるための方法を提供する。各ファイルサーバは、他のファイルサーバへの、信頼できる(例えば

冗長性を有する)通信経路を維持し、複数のファイルサーバ間においていかなる単一の障害地点が生じることとも防止する。各ファイルサーバは、その状態を信頼できる(例えば永続的な)メモリに保持し、それ自体の状態情報を確実にアップデートするとともに、状態情報が正しくないときそれを再構築する方法を提供する。ここで、各永続的メモリのうちの少なくとも一部は他のファイルサーバからアクセス可能である。各ファイルサーバは定期的に、他のファイルサーバの状態を確認し、そして、他方のファイルサーバがその担当分のサービスを明らかに提供できないときのみ、引き継ぎ動作を試みる。

【0006】

好ましい実施例では、各ファイルサーバは、他のファイルサーバが冗長な複数の通信経路を順序付けられた単一のメッセージストリームに結合できるように、冗長な複数の通信経路上でメッセージをシーケンス化する。各ファイルサーバは、他のファイルサーバが当該他のファイルサーバ自体に係る既知の最新の状態から更新されたか否かを決定するために、それ自体の状態をその永続的メモリに保持し、そしてその状態を、順序付けられたメッセージストリームと比較する。各ファイルサーバは他のファイルサーバに障害が発生したと確信したときにリソースの引き継ぎを相互に試行することを防ぐために、各ファイルサーバは、それ自体が冗長な通信経路の一部となっている共用のリソース(例えば磁気ディスク)を使用する。

【0007】

好ましい実施例では、各ファイルサーバは、複数のファイルサーバが障害発生と他のファイルサーバのリソース捕捉の試行とを繰り返す可能性を防止するために、エラーから回復するときに状態レポートを他のファイルサーバに提供する。

【0008】

【好ましい実施例の詳細な説明】

以下の説明では、この発明の好ましい実施例が好ましい処理ステップおよびデータ構造に関して述べられている。しかしながら、当業者であれば、この明細書を読めば、所定の制御プログラム下で動作する一つまたは複数の汎用プロセッサ(または特定のプロセスステップおよびデータ構造に適合した特定目的のプロセッサ)を用いて本発明の実施形態を実現可能であり、そのような装置を用いてここに説明した好ましいプロセスステップおよびデータ構造について実施することは、他の実験や別の発明を必要とはしないということを理解するであろう。

【0009】

好ましい実施例では、ファイルサーバシステムおよびその中の各ファイルサーバは、次の特許出願で述べられた発明を用いて動作する。

【0010】

発明者 Steven Kleimanによる1998年3月10日出願の出願番号第09/037,652の「Highly Available File Servers」代理人の事件番号NAP-012

【0011】

この出願の内容は、本明細書では「クラスタリングの開示内容」として参照する。

【0012】

好ましい実施例では、次の特許出願で述べられた発明を用い、RAIDストレージシステムのような冗長なアレーを形成するために、ファイルサーバシステム内の各ファイルサーバは、それに関係した複数の大容量記憶装置を制御する。

【0013】

発明者 David Hitz その他による1995年6月5日出願の出願番号第08/471,218の「A Method for Providing Parity in a Raid Sub-System Using Non-Volatile Memory」代理人の事件番号NET-004

【0014】

発明者 David Hitz その他による1995年3月31日出願の出願番号第08/454,921の「Write Anywhere File-System Layout」代理人の事件番号NET-005

【0015】

発明者 David Hitz その他による 1995 年 3 月 31 日出願の出願番号第 08/464,591 の「Method for Allocating Files in a File System Integrated with a Raid Disk Sub-System」代理人の事件番号 NET-006

【0016】

これらの出願の内容は、ともに、本明細書では「W A F L の開示内容」として参照する。

【0017】

(システム要素)

図 1 は調整された永続的状態情報を有する複数ファイルサーバシステムのブロック図を示す。システム 100 は複数のファイルサーバ 110、複数の大容量記憶装置 120、S A N (system area network) 130 および P N (public network) 140 を含む。

【0018】

好ましい実施例では、正確に 2 個のファイルサーバ 110 が存在する。各ファイルサーバ 110 は、複数の大容量記憶装置 120 に関して独立して動作することが可能である。各ファイルサーバ 110 は、クライアント装置 (不図示) からのファイルサーバ要求を受信し、それに応答して大容量記憶装置 120 に対する処理動作を実行し、また、ファイルサーバ要求への応答をクライアント装置に送信するように配置される。

【0019】

例えば、好ましい実施例では、各ファイルサーバ 110 は、上述の「クラスタリングの開示内容」に記載されたファイルサーバに似る。

【0020】

各ファイルサーバ 110 は、プロセッサ 111 と、プログラム及びデータメモリ 112 と、発生しうるサービス中断の時間中にわたって状態情報を保持する永続的メモリ 113 とを含む。好ましい実施の形態において、永続的メモリ 113 は、不揮発性の R A M を含む。

【0021】

大容量記憶装置 120 は、好ましくは、複数の磁気ディスクや磁気光ディスクや光ディスクを含む。好ましい実施の形態において、大容量記憶装置 120 は、レイド (R A I D) 構成で配置されるか、或いは、発生しうるサービス中断の時間中にわたって不変な情報を保持する他のシステム構成で配置される。

【0022】

各大容量記憶装置 120 は、大容量記憶装置用バス 121 を用いてファイルサーバのそれぞれに接続される。好ましい実施の形態において、各ファイルサーバ 110 は、独自の大容量記憶装置用バス 121 をそれぞれ有している。第 1 のファイルサーバ 110 は、複数の大容量記憶装置 120 のうちの第 1 サブセットを制御するプライマリコントローラとなり、また大容量記憶装置 120 のうちの第 2 サブセットを制御するセカンダリコントローラとなるように、複数の大容量記憶装置 120 に接続される。第 2 のファイルサーバ 110 は、複数の大容量記憶装置 120 のうちの第 2 サブセットを制御するプライマリコントローラとなり、また大容量記憶装置 120 のうちの第 1 サブセットを制御するセカンダリコントローラとなるように、複数の大容量記憶装置 120 に接続される。

【0023】

各ファイルサーバ 110 に関連付けられた大容量記憶装置用バス 121 は、ファイルサーバ 110 が大容量記憶装置 120 を制御できるように、ファイルサーバ 110 のプロセッサ 111 に接続される。他のもう 1 つの実施の形態において、ファイルサーバ 110 は、ファイバチャネルスイッチ (fiber channel switches) やスイッチドファブリック (switched fabrics) のような他の技術を用いて大容量記憶装置 120 に接続されてもよい。

【0024】

大容量記憶装置 120 は、少なくとも 1 つの指定された領域 123 をそれぞれ有する複数のメールボックスディスク 122 を含むように配置され、あるファイルサーバ 110 は

、メッセージ 1 2 4 を、他のファイルサーバ 1 1 0 によって読み出されるために上記指定された領域 1 2 3 に書き込むことができる。好ましい実施の形態においては、各メールボックスディスク 1 2 2 には、各ファイルサーバ 1 1 0 による読み出し用の少なくとも 1 つの指定された領域 1 2 3 と、各ファイルサーバ 1 1 0 による書き込み用の少なくとも 1 つの指定された領域 1 2 3 が存在する。

【 0 0 2 5 】

S A N 1 3 0 は、各ファイルサーバ 1 1 0 におけるプロセッサ 1 1 1 と永続的メモリ 1 1 3 とに接続される。S A N 1 3 0 は、第 1 のファイルサーバ 1 1 0 に設けられたプロセッサ 1 1 1 から第 2 のファイルサーバ 1 1 0 に設けられた永続的メモリ 1 1 3へメッセージ 1 2 4 を伝送するために配置される。同様に、S A N 1 3 0 は、第 2 のファイルサーバ 1 1 0 に設けられたプロセッサ 1 1 1 から第 1 のファイルサーバ 1 1 0 に設けられた永続的メモリ 1 1 3へメッセージ 1 2 4 を伝送するためにも配置される。

【 0 0 2 6 】

好ましい実施の形態において、S A N 1 3 0 は、2 つのファイルサーバ 1 1 0 間を接続するサーバネット接続経路 ( S e r v e r N e t C o n n e c t i o n ) から構成される。他のもう 1 つの実施の形態においては、永続的メモリ 1 1 3 は、ファイルサーバ 1 1 0 からは遠隔し、S A N 1 3 0 を用いてアクセス可能であるように、論理的に配置されてもよい。

【 0 0 2 7 】

P N 1 4 0 は、各ファイルサーバ 1 1 0 のプロセッサ 1 1 1 に接続される。P N 1 4 0 は、各ファイルサーバ 1 1 0 から他のファイルサーバ 1 1 0 へメッセージ 1 2 4 を伝送するために配置される。

【 0 0 2 8 】

好ましい実施の形態において、P N 1 4 0 は、直接通信チャネルや、L A N ( l o c a l a r e a n e t w o r k ) や、W A N ( w i d e a r e a n e t w o r k ) 、或いはそれらの組み合わせから構成される。

【 0 0 2 9 】

大容量記憶装置 1 2 0、S A N 1 3 0、及び P N 1 4 0 は、それぞれメッセージ 1 2 4 を伝送するために配置されるが、たとえこれらのメッセージ 1 2 4 に対するペイロードが同一のときであっても、ファイルサーバ 1 1 0 間においてこれらの伝送経路を用いて伝送されるメッセージ 1 2 4 は、実質上異なる形式を有してもよい。

【 0 0 3 0 】

( 動作方法 )

図 2 は、調整された永続的状态情報を有する複数ファイルサーバシステムの動作方法を示す図である。

【 0 0 3 1 】

状態遷移図 2 0 0 は、複数の状態と、それらの間における複数の遷移とを含んでいる。各遷移は、第 1 の状態から第 2 の状態への遷移であり、ある選択されたイベントが検出されたときに実行される。

【 0 0 3 2 】

状態遷移図 2 0 0 は、各ファイルサーバ 1 1 0 により独立して迎られる。このようにして、“この”ファイルサーバ 1 1 0 に対する 1 つの状態と、“他の”ファイルサーバ 1 1 0 に対する他のもう 1 つの状態が存在する。各ファイルサーバ 1 1 0 は、各状態から次へ進むべきそれ自体の状態が何であるかを独立して決定する。図 2 において、状態遷移図 2 0 0 は、“この”ファイルサーバ 1 1 0 に関して描かれている。

【 0 0 3 3 】

通常 ( N O R M A L ) 状態 2 1 0 において、このファイルサーバ 1 1 0 は、それに割り当てられた固有の大容量記憶装置 1 2 0 の制御権を有する。

【 0 0 3 4 】

引き継ぎ ( T A K E O V E R ) 状態 2 2 0 において、このファイルサーバ 1 1 0 は、通

常は他のファイルサーバ 1 1 0 に割り当てられている大容量記憶装置 1 2 0 の制御権を引き継ぐ。

【 0 0 3 5 】

停止 ( S T O P P E D ) 状態 2 3 0 において、このファイルサーバ 1 1 0 は、どの大容量記憶装置 1 2 0 の制御権も持たず、非動作状態にある。

【 0 0 3 6 】

再起動 ( R E B O O T I N G ) 状態 2 4 0 において、このファイルサーバ 1 1 0 は、どの大容量記憶装置 1 2 0 の制御権も持たず、サービス中断から回復中の状態にある。

【 0 0 3 7 】

( 通常状態 )

通常状態 2 1 0 において、ファイルサーバ 1 1 0 の両者は、正常に動作しており、それぞれの大容量記憶装置 1 2 0 のセットを制御する。

【 0 0 3 8 】

この状態において、各ファイルサーバ 1 1 0 は、2 つのファイルサーバ 1 1 0 間の冗長な通信経路を用いて、状態情報をメッセージ 1 2 4 において周期的に送信する。このようにして、各ファイルサーバ 1 1 0 は、次に挙げる技術により、状態情報を有するメッセージ 1 2 4 を周期的に伝送する。

【 0 0 3 9 】

・各ファイルサーバ 1 1 0 は、各ファイルサーバ 1 1 0 ごとに割り当てられた大容量記憶装置 1 2 0 にあるメールボックスディスクにメッセージをコピーすることにより、メッセージ 1 2 4 を伝送する。

好ましい実施の形態において、メッセージ 1 2 4 は、メッセージ 1 2 4 を第 1 のメールボックスディスクに書き込み、次いで第 2 のメールボックスディスクに書き込むというように、メールボックスディスクを用いて伝送される。

・各ファイルサーバ 1 1 0 は、S A N 1 3 0 を用いてメッセージ 1 2 4 を当該ファイルサーバ 1 1 0 の永続的メモリ 1 1 3 ( 或いは、当該ファイルサーバ 1 1 0 自体の永続的メモリ 1 1 3 と、他のファイルサーバ 1 1 0 の永続的メモリ 1 1 3 との両方 ) にコピーすることにより、メッセージ 1 2 4 を伝送する。

好ましい実施の形態において、メッセージ 1 2 4 は、N U M A 技術を利用した S A N 1 3 0 を用いて伝送される。

・そして、各ファイルサーバ 1 1 0 は、P N 1 4 0 を用いてメッセージ 1 2 4 を他のファイルサーバ 1 1 0 に伝送することにより、メッセージ 1 2 4 を伝送する。

好ましい実施の形態において、メッセージ 1 2 4 は、U D P や I P のような、両方のファイルサーバ 1 1 0 により利用可能な通信プロトコルへのカプセル化を行う P N 1 4 0 を用いてファイルサーバ 1 1 0 の両者へ伝送される。

【 0 0 4 0 】

各メッセージ 1 2 4 のそれぞれは、“この”ファイルサーバ 1 1 0 ( メッセージ 1 2 4 を送信しているファイルサーバ 1 1 0 である ) に対する次に挙げるような情報を含む。

【 0 0 4 1 】

・このファイルサーバ 1 1 0 のシステム I D 。

・このファイルサーバ 1 1 0 の状態表示子。

・このファイルサーバ 1 1 0 に係る現在のインスタンスと同一視される単調増加する番号からなる生成番号 ( 世代番号 ) G i 。

・このファイルサーバ 1 1 0 により伝送された現在のメッセージ 1 2 4 と同一視される単調増加する番号からなるシーケンス番号 S i 。

【 0 0 4 2 】

好ましい実施の形態において、状態表示子は、次に挙げるもののなかの 1 つであってよい。

【 0 0 4 3 】

N O R M A L ( 通常 ) : 通常の動作状態。

T A K E O V E R ( 引 き 継 ぎ ) : このファイルサーバ 1 1 0 が大容量記憶装置 1 2 0 の制御権を引き継ぐ状態。

N O - T A K E O V E R ( 引 き 継 ぎ な し ) : このファイルサーバ 1 1 0 が、当該ファイルサーバ 1 1 0 に割り当てられた大容量記憶装置 1 2 0 の制御権が、所定の引き受け側のファイルサーバ 1 1 0 によって引き継がれることを所望しない状態。

D I S A B L E ( 動 作 不 可 ) : 両方のファイルサーバ 1 1 0 に関して引き継ぎが実行不可能にされている状態。

【 0 0 4 4 】

好ましい実施の形態において、ブートアップ ( b o o t - u p ) 時にこのファイルサーバ 1 1 0 が動作開始されるとき、このファイルサーバ 1 1 0 のインスタンスが増分される。ファイルサーバ 1 1 0 が、再初期化を必要とするサービス中断の状態になると、生成番号 G i は増分され、メッセージ 1 2 4 は、それがサービス中断の前に送信された任意のメッセージ 1 2 4 に後続するものであることを示す。

【 0 0 4 5 】

同様にして、各メッセージ 1 2 4 のそれぞれは、“他の”ファイルサーバ 1 1 0 ( すなわち、メッセージ 1 2 4 を受信しているファイルサーバ 1 1 0 ) に対する次に挙げるような情報を含む。

【 0 0 4 6 】

・他のファイルサーバ 1 1 0 に係る現在のインスタンスと同一視される単調増加する番号からなる生成番号 G i 。

・他のファイルサーバ 1 1 0 から受信された最新のメッセージ 1 2 4 と同一視される単調増加する番号からなるシーケンス番号 S i 。

【 0 0 4 7 】

また、各メッセージ 1 2 4 は、当該メッセージ 1 2 4 が伝送される際に使用される状態プロトコルのバージョン番号を含む。

【 0 0 4 8 】

ファイルサーバ 1 1 0 は、複数の経路を用いて複数のメッセージ 1 2 4 を受信するので、それは、各メッセージ 1 2 4 に対して、そのメッセージ 1 2 4 が「新しい」か ( ファイルサーバ 1 1 0 がそれを以前に見なかったか ) 、又は「古い」か ( ファイルサーバ 1 1 0 がそれを以前に見たか ) を決定する。ファイルサーバ 1 1 0 は、最も最近の新しいメッセージの生成番号 G i とシーケンス番号 S i との記録を保持する。ファイルサーバ 1 1 0 は、特定のメッセージ 1 2 4 に対して、

その生成番号 G i が最も最近の新しいメッセージ 1 2 4 よりも大きいとき、又は、

その生成番号 G i が最も最近の新しいメッセージ 1 2 4 に等しく、そのシーケンス番号 S i が最も最近の新しいメッセージ 1 2 4 よりも大きいとき、  
かつ、そのときに限って、上記特定のメッセージ 1 2 4 が新しいと決定する。

【 0 0 4 9 】

複数のファイルサーバ 1 1 0 のいずれかが、メッセージ 1 2 4 が新しくないと決定するとき、そのファイルサーバ 1 1 0 はそのメッセージ 1 2 4 を無視することができる。

【 0 0 5 0 】

この状態において、各ファイルサーバ 1 1 0 は、複数のメッセージ 1 2 4 を用いて、それ自体の状態情報を周期的に保存する。従って、各ファイルサーバ 1 1 0 は、その状態情報を、それ自体の複数のメールボックスディスクと、それ自体の永続的 ( persistent ) メモリ 1 1 3 との両方に記録する。

【 0 0 5 1 】

この状態において、各ファイルサーバ 1 1 0 は、他のファイルサーバ 1 1 0 における状態変化を周期的に観測する。第 1 のファイルサーバ 1 1 0 は、少なくとも以下の 2 つの方法の中の 1 つで、第 2 のファイルサーバ 1 1 0 における状態変化を検出する。

【 0 0 5 2 】

・第 1 のファイルサーバ 1 1 0 は、第 2 のファイルサーバ 1 1 0 が タイムアウト期間 の



間に(メッセージ124を用いて)その状態情報を更新しなかったことを記録(ノート)する。

好ましい実施形態において、このタイムアウト期間は、複数のメールボックスディスクを用いる通信に対しては1秒(two-half seconds)であり、SAN130を用いる通信に対しては0.5秒(one-half second)である。しかしながら、これらのタイムアウト値を用いることに対する特定の要求は存在せず、それに代わる実施形態においては、異なるタイムアウト値、又はタイムアウト期間とは異なる技術を用いてもよい。

・第1のファイルサーバ110は、第2のファイルサーバ110がその状態を変化させたことを示すために、第2のファイルサーバ110が(1つ又はそれより多くのメッセージ124を用いて)その状態情報を更新したことを記録(ノート)する。

好ましい実施形態において、第2のファイルサーバ110は、それが、いつ、各メッセージ124に関して記述された状態の中の1つにあるかを示す。

#### 【0053】

第1のファイルサーバ110が、第2のファイルサーバ110もまた通常状態にあると決定するとき、状態210にとどまるために通常動作遷移(トランジション)211が実行される。

#### 【0054】

第1のファイルサーバ110は、それが第2のファイルサーバ110から受信する複数のメッセージ124に応答して、その決定を実行する。前述のタイムアウト期間に応答した(例えばタイムアウト期間の2乃至5倍のような)時間期間の間に、そのようなメッセージ124が存在しないならば、第1のファイルサーバ110は、第2のファイルサーバがサービス中断の状態になったと決定する。

#### 【0055】

第1のファイルサーバ110が、第2のファイルサーバ110がサービス中断の状態になったと決定するとき(すなわち、第2のファイルサーバ110が停止状態230にあるとき)、引き継ぎ状態220に入るために引き継ぎ動作遷移212が実行される。

#### 【0056】

引き継ぎ動作遷移212は、DISABLE又はNO-TAKEOVERのようなメッセージ124の状態表示子によって実行不可能にすることができる。

#### 【0057】

好ましい実施形態において、いずれのファイルサーバ110も、(a)オペレータのコマンド、(b)複数の永続的メモリ113の間の同期エラー、又は(c)複数のファイルサーバ110の間の互換性の不整合に応答して、引き継ぎ動作遷移212を実行不可能にすることができる。

#### 【0058】

引き継ぎ動作遷移212を実行するために、このファイルサーバ110はステップ213において以下の動作を実行する。

#### 【0059】

・このファイルサーバ110は、他のファイルサーバ110に、(複数のメールボックスディスク122、SAN130、及びPN140を含む)信頼性のある通信経路を含むことを用いて、メッセージ124の状態表示子TAKEOVERを送る。

・このファイルサーバ110は、他のファイルサーバ110が引き継ぎ動作遷移212を受信してそれに基づいて動作する(すなわち、複数の大容量記憶装置120へのそれ自体のアクセスをサスペンドする)機会を持つまで待機する。

・このファイルサーバ110は、通常は他のファイルサーバ110に割り当てられている複数の大容量記憶装置120に、複数のディスク予約コマンドを発行する。

・このファイルサーバ110は、他のファイルサーバ110が受動的であることを確かめるために、他の任意の適当な動作を実行する。

#### 【0060】

引き継ぎ動作が成功したときは、引き継ぎ動作遷移212は完了し、このファイルサーバ

バは引き継ぎ状態 2 2 0に入る。そうでないとき（引き継ぎが実行不可能にされているようなとき）は、このファイルサーバ 1 1 0 は通常状態 2 1 0に戻る。

【 0 0 6 1 】

（引き継ぎ状態）

引き継ぎ状態 2 2 0において、このファイルサーバ 1 1 0 は適正に動作しているが、他のファイルサーバ 1 1 0 はそうではない。このファイルサーバ 1 1 0 は、その大容量記憶装置 1 2 0 と、他のファイルサーバの大容量記憶装置 1 2 0 との双方の制御を引き継ぐ。

【 0 0 6 2 】

この状態において、このファイルサーバ 1 1 0 は、サービス中断のイベントにおいてそれ自体の状態を保存するように、永続的メモリ 1 1 3 と複数のメールボックスディスク 1 2 2 に複数のメッセージ 1 2 4 を書き込み続ける。

【 0 0 6 3 】

この状態において、このファイルサーバ 1 1 0 は、全ての大容量記憶装置 1 2 0、すなわちそれ自体のものと、通常は他のファイルサーバ 1 1 0 に割り当てられるものに対して、このファイルサーバ 1 1 0 がいくつかの大容量記憶装置 1 2 0 の制御を返却するべきであると決定するまで、制御を続ける。

【 0 0 6 4 】

好ましい実施形態において、第 1 のファイルサーバ 1 1 0 は、オペレータの制御に応答してその決定を実行する。このファイルサーバ 1 1 0 に対するオペレータは、他のファイルサーバ 1 1 0 がそのサービス中断から回復したと判断する。通常状態に入るために、返却動作遷移 2 2 1 が実行される。

【 0 0 6 5 】

それに代わる実施形態において、第 1 のファイルサーバ 1 1 0 は、それが第 2 のファイルサーバ 1 1 0 から受信する複数のメッセージ 1 2 4 に応答してその決定を実行してもよい。第 2 のファイルサーバが、それがサービス中断から回復した（すなわち、それが R E B O O T I N G 状態 2 4 0 にある）ことを示す複数のメッセージ 1 2 4 を送るとき、第 1 のファイルサーバ 1 1 0 は、返却動作遷移 2 2 1 を開始してもよい。

【 0 0 6 6 】

返却動作遷移 2 2 1 を実行するために、このファイルサーバ 1 1 0 はステップ 2 2 2 において以下の動作を実行する。

【 0 0 6 7 】

・このファイルサーバ 1 1 0 は、通常は他のファイルサーバ 1 1 0 に割り当てられている大容量記憶装置への、その複数のディスク予約コマンドを開放する。

・このファイルサーバ 1 1 0 は、複数のメールボックスディスク 1 2 2、S A N 1 3 0 及び P N 1 4 0 を用いることを含んで、メッセージ 1 2 4 の状態表示子 N O R M A L を他のファイルサーバ 1 1 0 に送る。

・このファイルサーバ 1 1 0 は、他のファイルサーバ 1 1 0 による引き継ぎ動作遷移 2 1 2 を、他のファイルサーバ 1 1 0 が通常状態 2 1 0 に入るまで実行不可能にする。このファイルサーバ 1 1 0 は、他のファイルサーバ 1 1 0 が通常状態 2 1 0 に入るまで、ステップ 2 2 2 にとどまる。

【 0 0 6 8 】

返却動作が成功したとき、返却動作遷移 2 2 1 は完了し、このファイルサーバは通常状態 2 1 0 に入る。

【 0 0 6 9 】

（停止状態）

停止状態 2 3 0 において、このファイルサーバ 1 1 0 は、どの大容量記憶装置 1 2 0 の制御権も持たず、非動作状態にある。

【 0 0 7 0 】

この状態において、このファイルサーバ 1 1 0 は、再起動の実行を当該ファイルサーバ

1 1 0 自体が決定するまで、何の動作も実行しない。

【0071】

好ましい実施の形態において、第1のファイルサーバ110は、オペレータの制御に  
答してその決定を行う。このファイルサーバ110のオペレータは、ファイルサーバ110  
がそのサービス中断から回復したと決定する。再起動状態240に入るために、再起動  
動作遷移231が実行される。

【0072】

代替の実施の形態において、第1のファイルサーバ110は、タイマー、または、他  
の再起動しようとする自動的な試みに応答してその決定を行ってもよい。このファイルサ  
ーバ110は、それ自体がそのサービス中断から回復したと決定するとき再起動を試み、  
再起動状態240に入るために、再起動動作遷移231が実行される。

【0073】

(再起動状態)

再起動状態240において、このファイルサーバ110は、どの大容量記憶装置120  
も制御せず、サービス中断からの回復中の状態にある。

【0074】

この状態において、ファイルサーバ110は、サービス中断から回復を試みる。

【0075】

もし、このファイルサーバ110がサービス中断から回復できないなら、再起動失敗遷  
移241が行われ、このファイルサーバ110は再起動状態240にとどまる。

【0076】

もし、このファイルサーバ110がサービス中断から回復でき、他のファイサーバ110  
が引き継ぎ状態220にあるなら、再起動失敗遷移241が行われ、このファイルサ  
ーバ110は再起動状態240にとどまる。この場合に、他のファイルサーバ110は、通  
常はこのファイルサーバ110に割り当てられている大容量記憶装置120を制御し、こ  
のファイルサーバ110は、サービス中断からの回復を再び試みる前に、返却動作遷移2  
21を待つ。

【0077】

もし、このファイルサーバ110がサービス中断から回復でき、それ自体が通常状態2  
10 (以下で説明される)に入るべきであると決定するなら、再起動 - 通常 (REBOOT - NORMA  
L) 遷移242が行われ、このファイルサーバ110は通常状態210に入る。

【0078】

もし、このファイルサーバ110がサービス中断から回復でき、引き継ぎ状態210 (以  
下で説明される)に入るべきであると決定するなら、再起動 - 引き継ぎ (REBOOT - T  
AKEOVER) 遷移243が行われ、このファイルサーバ110は引き継ぎ状態210に入る。

【0079】

好ましい実施の形態において、このファイルサーバ110は、以下のステップでサービ  
ス中断からの回復の試みを実行する。

【0080】

ステップ251において、このファイルサーバ110は、回復動作を始める。

【0081】

ステップ252において、このファイルサーバ110は、大容量記憶装置120のいず  
れかに書き込める (すなわち、他のファイルサーバ110が引き継ぎ状態220にある)  
かどうかを決定する。もし書き込めるなら、このファイルサーバ110は、オペレータへ  
のプロンプトを表示し、オペレータに、返却動作遷移221を他のファイルサーバ110  
に実行させるコマンドを表示して要求する。

【0082】

このファイルサーバ110は、オペレータが他のファイルサーバ110に返却動作を実

行させるよう命令し、返却動作遷移 2 2 1 が完了するまで待ち、次のステップに移る。

【 0 0 8 3 】

ステップ 2 5 3 において、このファイルサーバ 1 1 0 は、他のファイルサーバ 1 1 0 の状態を決定する。このファイルサーバ 1 1 0 は、それ自体の永続的メモリ 1 1 3 とメールボックスディスク 1 2 2 に応答してこの決定を行う。このファイルサーバ 1 1 0 は、他のファイルサーバ 1 1 0 の、再起動状態 2 2 0 に入る以前における状態（すなわち、通常状態 2 1 0 であったか、それとも引き継ぎ状態 2 2 0 であったか）を記録する。

【 0 0 8 4 】

もし、このファイルサーバ 1 1 0 が、他のファイルサーバ 1 1 0 が 通常状態 2 1 0 であると決定するなら、ステップ 2 5 4 に移る。もし、このファイルサーバ 1 1 0 が、前もって全ての大容量記憶装置 1 2 0 から引き継いだ（すなわち、他のファイルサーバ 1 1 0 は 停止状態 2 3 0、または、再起動状態 2 4 0 にある）と決定するなら、それはステップ 2 5 5 に移る。

【 0 0 8 5 】

ステップ 2 5 4 において、このファイルサーバ 1 1 0 は、それ自体の大容量記憶装置 1 2 0 を捕捉しようとするが、通常は他のファイルサーバ 1 1 0 に割り当てられている大容量記憶装置を捕捉しようとはしない。このファイルサーバ 1 1 0 は、ステップ 2 5 6 に移る。

【 0 0 8 6 】

ステップ 2 5 5 において、このファイルサーバ 1 1 0 は、それ自体の大容量記憶装置 1 2 0 と、通常は他のファイルサーバ 1 1 0 に割り当てられている大容量記憶装置 1 2 0 との両方を捕捉しようとする。このファイルサーバ 1 1 0 は、ステップ 2 5 6 に移る。

【 0 0 8 7 】

ステップ 2 5 6 において、このファイルサーバ 1 1 0 は、その永続的メモリ 1 1 3 が、実行中のファイルサーバ動作に関して最新であるかどうかを決定する。もし 最新の状態を反映していないならば、このファイルサーバ 1 1 0 は、実行中のファイルサーバ動作に関してその永続的メモリ 1 1 3 をすっかり消去する。

【 0 0 8 8 】

ステップ 2 5 7 において、このファイルサーバ 1 1 0 は、他のファイルサーバと通信できるかどうか、および、（例えばオペレータのコマンドのような）引き継ぎ動作を妨げる何かがあるかどうかを決定する。このファイルサーバ 1 1 0 は、永続的メモリ 1 1 3 とメールボックスディスク 1 2 2 に応答してその決定を行う。

【 0 0 8 9 】

ステップ 2 5 8 において、もし、このファイルサーバ 1 1 0 が、再起動状態 2 4 0 に入る前に通常状態 2 1 0 にあった（すなわち、このファイルサーバ 1 1 0 が、ステップ 2 5 4 を実行し、それ自体の大容量記憶装置 1 2 0 のみを捕捉した）なら、それは通常状態 2 1 0 に入る。

【 0 0 9 0 】

ステップ 2 5 8 において、もし、このファイルサーバ 1 1 0 が、再起動状態 2 4 0 に入る前に引き継ぎ状態 2 2 0 にあった（すなわち、このファイルサーバ 1 1 0 がステップ 2 5 5 を実行し、全ての大容量記憶装置 1 2 0 を 捕捉した）なら、それは 引き継ぎ状態 2 2 0 に入る。

【 0 0 9 1 】

（代わりの実施の形態）

好ましい実施の形態がここに開示されるが、本発明の概念、範囲、および、精神を越えない多くの変形が可能である。これらの変形は、本明細書を読んだ後で当業者に明らかとなる。

【図面の簡単な説明】

【図 1】 調整された永続的状態情報を有する複数ファイルサーバのブロック図である。

【図 2】 調整された永続的状態情報を有する複数ファイルサーバに関する動作方法を示した状態遷移図である。

【手続補正 2】

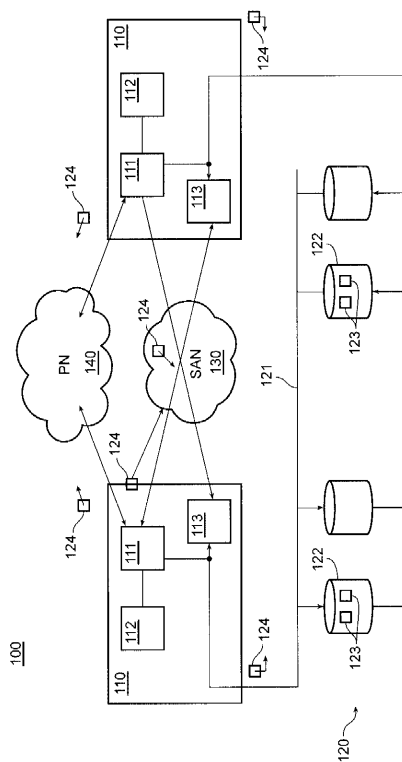
【補正対象書類名】図面

【補正対象項目名】全図

【補正方法】変更

【補正の内容】

【図 1】



【図 2】

