

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4620457号  
(P4620457)

(45) 発行日 平成23年1月26日(2011.1.26)

(24) 登録日 平成22年11月5日(2010.11.5)

(51) Int.Cl.

F I

G 0 6 F 12/00 (2006.01)

G 0 6 F 12/00 5 2 0 J

G 0 6 F 12/00 5 3 1 M

請求項の数 11 (全 15 頁)

(21) 出願番号 特願2004-512028 (P2004-512028)  
 (86) (22) 出願日 平成15年6月3日(2003.6.3)  
 (65) 公表番号 特表2005-529410 (P2005-529410A)  
 (43) 公表日 平成17年9月29日(2005.9.29)  
 (86) 国際出願番号 PCT/US2003/017497  
 (87) 国際公開番号 W02003/105026  
 (87) 国際公開日 平成15年12月18日(2003.12.18)  
 審査請求日 平成17年12月15日(2005.12.15)  
 (31) 優先権主張番号 10/165,188  
 (32) 優先日 平成14年6月7日(2002.6.7)  
 (33) 優先権主張国 米国 (US)

前置審査

(73) 特許権者 303039534  
 ネットアップ、インコーポレイテッド  
 アメリカ合衆国 カリフォルニア 940  
 89, サニーヴェール, イースト ジ  
 ャバ ドライブ 495  
 (74) 代理人 100101454  
 弁理士 山田 卓二  
 (74) 代理人 100081422  
 弁理士 田中 光雄  
 (72) 発明者 デイビッド・ヒッツ  
 アメリカ合衆国 94022 カリフォルニア  
 州ロス・アルトス、シルビアン・ウェイ 1  
 37番

最終頁に続く

(54) 【発明の名称】 複数の同時にアクティブなファイルシステム

(57) 【特許請求の範囲】

【請求項 1】

複数のアクティブな、すなわち書き込み可能なファイルシステムを含むデータ記憶を行う方法であって、

第1のアクティブなファイルシステムに基づいて第2のアクティブなファイルシステムを作成し、当初は上記第1のアクティブなファイルシステムおよび上記第2のアクティブなファイルシステムがデータを共有し、

上記第1のアクティブなファイルシステムに変更が加えられたとき、変更されたデータが上記第1のアクティブなファイルシステム内の上記第2のアクティブなファイルシステムとは共有されない場所に記録され、

上記第2のアクティブなファイルシステムに変更が加えられたとき、変更されたデータが上記第2のアクティブなファイルシステム内の上記第1のアクティブなファイルシステムとは共有されない場所に記録され、

上記アクティブなファイルシステムの各々に加えられた変更が、上記変更されたアクティブなファイルシステムと上記データを共有するアクティブなファイルシステムに反映されず、

更に各々のアクティブなファイルシステムが第1のポインタと第2のポインタを含む組織データを有し、変更されたアクティブなファイルシステムにおける変更されていないデータへの第1のポインタは複数のアクティブなファイルシステムを記憶する1以上の記憶装置内で変更されず、既存のアクティブなファイルシステムにおける新しいアクティブな

ファイルシステムへの第 2 のポインタが変更される方法。

【請求項 2】

複数のスナップショットが上記複数のアクティブなファイルシステムの 1 つから作成され、上記スナップショットの各々が、過去の 1 整合点でのそのアクティブなファイルシステムのイメージを形成する請求項 1 に記載の方法。

【請求項 3】

各スナップショットが、上記複数のアクティブなファイルシステムのためのアクティブなファイルシステムデータから分離した、ファイルシステムデータのための完全な階層を含む請求項 2 に記載の方法。

10

【請求項 4】

上記スナップショットの少なくとも 1 つが、新しいアクティブなファイルシステムに変換される請求項 2 に記載の方法。

【請求項 5】

上記スナップショットの 1 つが、上記スナップショットの 1 つを書き込み可能にすることにより変換される請求項 4 に記載の方法。

【請求項 6】

いずれかの上記アクティブなファイルシステムから上記新しいアクティブなファイルシステムへのスナップショットポインタが切断される請求項 5 に記載の方法。

20

【請求項 7】

複数のアクティブなファイルシステムを作成する方法であって、

第 1 のアクティブなファイルシステムのスナップショットを作成し、各々のアクティブなファイルシステムがスナップショットポインタと他のポインタを含む組織データを有し、上記スナップショットが当初は上記第 1 のアクティブなファイルシステムとデータを共有するステップ；および

上記スナップショットを書き込み可能にすることにより、上記スナップショットを第 2 のアクティブなファイルシステムに変換し、第 1 のアクティブなファイルシステムから第 2 のアクティブなファイルシステムへの第 1 のスナップショットポインタを切断するステップを含み、

30

上記第 1 のアクティブなファイルシステムに変更が加えられたとき、変更されたデータが、上記第 1 のアクティブなファイルシステム内の上記第 2 のアクティブなファイルシステムとは共有されない場所に記録され、

上記第 2 のアクティブなファイルシステムに変更が加えられたとき、変更されたデータが、上記第 2 のアクティブなファイルシステム内の上記第 1 のアクティブなファイルシステムとは共有されない場所に記録され、

ここに、上記第 1 のアクティブなファイルシステムに加えられた変更が上記第 2 のアクティブなファイルシステムに反映されず、上記第 2 のアクティブなファイルシステムに加えられた変更が上記第 1 のアクティブなファイルシステムに反映されず、変更された各々のアクティブなファイルシステムにおける変更されていないデータへの上記他のポインタが複数のアクティブなファイルシステムを記憶する 1 以上の記憶装置内で変更されない方法。

40

【請求項 8】

請求項 7 に記載の方法であって、さらに

上記第 1 のアクティブなファイルシステムの、当初は上記第 1 のアクティブなファイルシステムとデータを共有する新しいスナップショットを作成するステップ；および

上記新しいスナップショットを書き込み可能にすることにより、上記新しいスナップショットを第 3 のアクティブなファイルシステムに変換し、第 1 のアクティブなファイルシステムから第 3 のアクティブなファイルシステムへの第 2 のスナップショットポインタを切断するステップを含み、

ここに、上記第 1 のアクティブなファイルシステムまたは上記第 2 のアクティブなファ

50

イルシステムに加えられた変更が上記第3のアクティブなファイルシステムに反映されない

方法。

【請求項9】

上記第1のアクティブなファイルシステムまたは上記第2のアクティブなファイルシステムに変更が加えられたとき、変更されたデータが、上記第3のアクティブなファイルシステムとは共有されない場所に記録される請求項8に記載の方法。

【請求項10】

コンピュータシステム上で実行された時、上記コンピュータシステムに請求項1から請求項9のいずれか1つのステップをもたらすプログラムコード手段を含むコンピュータプログラム。

10

【請求項11】

少なくとも1つの記憶装置；

情報を送受信する少なくとも1つのコンピュータ・デバイスまたはネットワークへのインタフェース；および

上記記憶装置内の上記情報の記憶および取り出しを制御し、プログラム制御下で動作して請求項1から請求項9のいずれか1つのステップを実行するコントローラ

を含む記憶システム。

【発明の詳細な説明】

【技術分野】

20

【0001】

本発明は複数の同時に書き込み可能なファイルシステムに関する。

【背景技術】

【0002】

ファイルシステムは情報を記憶する構造を提供する。この情報は、例えば、ディスクドライブ、CD-ROMドライブなどの記憶装置上のアプリケーションプログラム、ファイルシステム情報、他のデータなど（今後、まとめて単にデータという）である。多くのファイルシステムについての1つの問題は、もし上記ファイルシステムがどういうわけか損傷すると、多量のデータが失われうることである。

【0003】

30

そのようなデータの損失を防ぐため、ファイルシステムのバックアップがしばしば作成される。ファイルシステムのバックアップを生成する1つのとても効率的な方法は、上記ファイルシステムのスナップショットを作成することである。スナップショットは整合点での上記ファイルシステムのイメージであって、整合点とは上記ファイルシステムが、自己整合的である点である。もしファイルシステムの中に記憶されたデータが正当なファイルシステムイメージを構成するならば、そのファイルシステムは自己整合的である。

【0004】

あるファイルシステム、例えばW A F L ( W r i t e A n y w h e r e F i l e s y s t e m L a y o u t ) ファイルシステムでは、ファイルシステムのスナップショットは、上記ファイルシステム内のデータの組織に関する情報をコピーすることにより作製可能である。次に、上記データ自体が上記記憶装置に保存されている限り、上記データは上記スナップショット経由でアクセス可能である。1つの仕組みが、例えばブロッックマップ経由で、このデータを保存するため、これらのファイルシステム内で提供される。

40

【発明の開示】

【発明が解決しようとする課題】

【0005】

従来は、スナップショットは読み取り専用である。読み取り専用スナップショットは、以前のバージョンのデータを再度読み出して、ファイルシステムへの損傷を修復する。これらの機能は非常に有用になりうる。しかし、これらの種類のスナップショットは、有用でありうる他の機能を提供しない。

50

## 【課題を解決するための手段】

## 【0006】

もし、スナップショットを変更したい人がスナップショットを変更するように、スナップショットに書き込めたら、好都合である。これはいくつかの効果があるだろう。

## 【0007】

- ・スナップショットに記憶されていた誤ったエントリの修正が可能になる。

## 【0008】

- ・ファイルシステムから消去されることを所望されたデータが消去可能になる。

## 【0009】

・ファイルシステム（またはファイルシステムにより維持されるデータについて）の「試験的な」バージョンに対する変更が可能になる。ファイルシステムの「試験的な」バージョンとは、上記ファイルシステムの壊滅的な誤りがその「実際の」アクティブなバージョンにおいてデータ損失を生じない上記ファイルシステムのバージョンである。

10

## 【0010】

- ・上記ファイルシステムの動作、または、上記ファイルシステムの保護の下で動作する、プログラムまたはデータベース動作に対して誤った更新を無効にすることが可能になる。

## 【0011】

書き込み可能なスナップショットは実はもう1つのアクティブなファイルシステムである。このアクティブなファイルシステムがもう1つのアクティブなファイルシステムからのデータに基づいているので、元のアクティブなファイルシステムに損害を与える危険を冒さずに、上記アクティブなファイルシステムに関する変更および修正を上記書き込み可能なスナップショットへ加えることができる。加えて、スナップショットは、単に組織情報をコピーして既存のデータを保存することにより作成できるので、書き込み可能なスナップショット（すなわち、新しいアクティブなファイルシステム）は、容易にかつわずかなシステム資源の利用で作成できる。

20

## 【0012】

これらおよび他の効果は、ここに記載した、複数の上記アクティブなファイルシステムが維持される本発明の実施の形態において提供される。複数の上記アクティブなファイルシステムの各々が、その中のもう1つの上記アクティブなファイルシステムと共有されるデータに最初にアクセスし、また、上記アクティブなファイルシステムの各々に加えられる変更は、他のアクティブなファイルシステムに反映されない。

30

## 【0013】

上記の好ましい実施の形態では、第2のアクティブなファイルシステムが第1のアクティブなファイルシステムに基づいて作成されたとき、上記第1のアクティブなファイルシステムおよび上記第2のアクティブなファイルシステムは、最初はデータを共有する。上記第1のアクティブなファイルシステムに変更が加えられたとき、変更されたデータが、上記第1のアクティブなファイルシステム内の上記第2のアクティブなファイルシステムとは共有されない場所に記録される。上記第2のアクティブなファイルシステムに変更が加えられたとき、変更されたデータが、上記第2のアクティブなファイルシステム内の上記第1のアクティブなファイルシステムとは共有されない場所に記録される。

40

## 【0014】

好ましくは、さらに別のスナップショットが複数の上記アクティブなファイルシステムのスナップショットから作られ、各スナップショットは、過去の整合点におけるそのアクティブなファイルシステムのイメージを形成する。各スナップショットは、複数の上記アクティブなファイルシステムのアクティブなファイルシステムデータから分離し、ファイルシステムデータの完全な階層を含む。これらのスナップショットの1つは、次に、上記スナップショットを書き込み可能にすることにより、および、いずれかの上記アクティブなファイルシステムから新しいアクティブなファイルシステムへのスナップショットポイントを切断することにより、上記新しいアクティブなファイルシステムに変換可能である

50

。

【 0 0 1 5 】

本発明は、また、これらの動作を実行する命令を含むメモリおよび上記動作を実行する記憶システムを含む。

【発明を実施するための最良の形態】

【 0 0 1 6 】

以下、添付の図を参照して発明の実施の形態を説明する。

【 0 0 1 7 】

以下の語句は、以下に説明するように本発明の複数の面に関係する。これらの語句の一般的意味の説明は、限定することを意図せず、説明に役立つことのみを意図している。

10

【 0 0 1 8 】

・データとは、一般的にどれかの情報。記憶装置またはファイルシステムに関連して、アプリケーションプログラムまたはデータ、マルチメディアデータ、上記記憶装置またはファイルシステムの組織データなどを含むがそれらに限定されない上記記憶装置またはファイルシステム内に記憶されたどれかの情報。

【 0 0 1 9 】

・組織データとは、一般的に、ファイルシステム内の他のデータのレイアウトを指定するデータ。W A F L ( W r i t e A n y w h e r e F i l e S y s t e m ) 設計では、上記組織データは直接にまたは間接に（すなわち、他の i ノードを経由して）上記ファイルシステム内のすべてのファイルのためのデータのブロックを指すルート i ノードを含む。W A F L 設計では、上記組織データを含むすべてのデータ（それ故、ルート i ノードおよび他の i ノード）は複数のブロックに記憶される。

20

【 0 0 2 0 】

・ i ノードとは、一般的に、情報ノード。W A F L 設計では、上記ファイルシステム内の他のブロックについてのデータを含む情報ノード。

【 0 0 2 1 】

・自己整合的である（ファイルシステム内のコンテキスト ( c o n t e x t ) の中で）とは、一般的に、ファイルシステムの組織に関するデータを含み、上記ファイルシステム内に記憶された上記データが有効なファイルシステムイメージを構成するとき、上記ファイルシステムは自己整合的であるという。

30

【 0 0 2 2 】

整合点とは、一般的に、整合点は、（ a ）ファイルシステムが自己整合的である時間、または、（ b ）整合点の時間におけるファイルシステム内の 1 組のデータをいう。

【 0 0 2 3 】

スナップショットとは、一般的に、整合点の時間において上記ファイルシステムにより維持されるデータの書かれた記録である。好ましい実施の形態では、各スナップショットは（ a ）アクティブなファイルシステムと同じフォーマットに維持され、かつ、（ b ）ファイルシステム名前空間を使用して参照可能であるけれども、これらの条件の一方を必要とすることに対して限定はない。

【 0 0 2 4 】

アクティブなファイルシステムとは、一般的に、アクティブなファイルシステムはアクセス可能で変更可能な 1 組のデータである。

40

【 0 0 2 5 】

ファイルシステム階層とは、一般的に、ファイルシステム階層とは、（ a ）名前空間へのデータの組織、または（ b ）記憶装置に維持されている、データまたはメタデータの情報を記録し、その情報にアクセスするために使用される 1 組のデータブロックおよびそれらの相互結合をいう。

【 0 0 2 6 】

スナップショットおよびアクティブなファイルシステム

【 0 0 2 7 】

50

図 1 は、本発明によりアクティブなファイルシステムに変換可能なスナップショットの作成を示す。

【 0 0 2 8 】

図 1 に示すファイルシステム 1 0 0 は 1 以上の記憶装置、例えばハードディスクドライブ、C D - R O M、または他の装置に存在する。好ましい実施の形態では、ファイルシステム 1 0 0 は W A F L システムであるが、これは必須ではない。

【 0 0 2 9 】

ファイルシステム 1 0 0 はルート i ノード 1 1 0 およびデータ 1 2 0 を他のデータとともに含む。ファイルシステム 1 0 0 内のすべての i ノードおよびデータは好ましくはブロック内に記憶されるが、これは必須ではない。

10

【 0 0 3 0 】

ルート i ノード 1 1 0 は、ファイルシステム 1 0 0 のための組織データの複数部分を記憶する。特に、ルート i ノード 1 1 0 は、データを指し、他の i ノードおよびデータを指し、このデータは、次に、ファイルシステム 1 0 0 内に記憶されたすべての情報に関するデータを指す。それ故、ファイルシステム 1 0 0 内に記憶されたすべてのデータは、ルート i ノード 1 1 0 で開始することにより到達可能である。

【 0 0 3 1 】

スナップショット 1 3 0 はファイルシステム 1 0 0 から形成された。図 1 では、スナップショット 1 3 0 の要素は、ファイルシステム 1 0 0 からそれらの要素を区別するのを助けるために破線を使用して示される。本発明の好ましい実施の形態によれば、上記スナップショットは、ファイルシステム 1 0 0 に関する整合点において、単にルート i ノード 1 1 0 をスナップショットのルート i ノード 1 4 0 にコピーすることにより、形成可能である。ある実施の形態では、追加の組織データがコピーされねばならないかもしれない。次に、ルート i ノード 1 1 0 により指されるすべてのデータおよび i ノード（およびすべての他のコピーされた組織データ）が保存される限りは、スナップショットのルート i ノード 1 4 0 はファイルシステム 1 0 0 の有効なコピーを指す。

20

【 0 0 3 2 】

スナップショットのルート i ノード 1 4 0 が作成された後で、スナップショット 1 3 0 およびファイルシステム 1 0 0 は 1 以上の記憶装置でのデータを実際に共有する。それ故、好ましくは、スナップショット 1 3 0 は、図 1 のデータ 1 2 0 のまわりの実線と破線の重なった境界により示されるように、ファイルシステム 1 0 0 のような 1 以上の記憶装置での同じ物理データ 1 2 0 を含む。すなわち、上記スナップショットと上記ファイルシステムは重なる。これは、記憶容量および他のシステム資源の効率的な使用により、スナップショット 1 3 0 の短時間での作成を可能にする。

30

【 0 0 3 3 】

ファイルシステム 1 0 0 は、好ましくは、ファイルシステム 1 0 0 のスナップショットを指すスナップショットデータ 1 5 0 を含む。特に上記スナップショットデータ内のポインタ 1 6 0 は、好ましくは、それらのスナップショットのルート i ノードを指す。

【 0 0 3 4 】

スナップショット 1 3 0 は、好ましくは、他のスナップショットを指すスナップショットデータ 1 7 0 も含む。しかし、スナップショット 1 3 0 のスナップショットデータ 1 7 0 は、スナップショット 1 3 0 が好ましくはそれ自体を指さないのので、ファイルシステム 1 0 0 のスナップショットデータ 1 5 0 とは異なる。この違いは、図 1 で、ファイルシステム 1 0 0 のスナップショットデータ 1 5 0 のまわりのスナップショット 1 3 0 の切欠により示される。

40

【 0 0 3 5 】

好ましくは、本発明によるファイルシステムのスナップショットは、アクティブなファイルシステムのアクティブなファイルシステムデータから分離したファイルシステムデータについての完全な階層を含む。この階層は、上記スナップショットのルート i ノード、および、可能ならば上記スナップショットに関してコピーされた他のノードおよびデータ

50

(図示されていない)に含まれる。

【0036】

関連するアクティブなファイルシステムに元は使用された名前空間を複製するというスナップショットに関する上記ファイルシステム階層について、要求はない。1つの好ましい実施の形態では、スナップショットのルートiノード内のファイル名(および他の組織データ)は、各スナップショットに関して記憶されねばならない組織データを最小にするため、ハッシュコードまたは他の技術を使用して圧縮可能である。しかし、代替の実施の形態では、好ましいかも知れない状況で、人間のユーザにとって比較的読みやすい形で、各スナップショットに関する元の名前空間および他の組織データを維持することは、よりすぐれている。これは、そのようなスナップショットに基づくバックアップおよびリストア作業をする人間のユーザを助ける有益な効果をもつ。

10

【0037】

図2は、アクティブなファイルシステムのスナップショットからのアクティブなファイルシステムの分岐を示す。

【0038】

ファイルシステム100はアクティブなので、上記ファイルシステム内のデータを変更するために、1つの仕組みが提供されねばならない。しかし、スナップショット130の整合性を維持するために、スナップショットのルートiノード140により指されるデータは保存されねばならない。それ故、例えば、データ120がファイルシステム100内で変更されたとき、変更されたデータ120'は1以上の記憶装置に記憶される。ファイルシステム100のルートiノード110およびすべての介在するiノードならびに組織データは、変更されたデータ120'を指すように更新される。加えて、変更されなかったデータ120は1以上の記憶装置上に保存される。スナップショットのルートiノード140はこの変更されなかったデータを指し続けるので、スナップショット130の整合性を維持する。

20

【0039】

同様に、アクティブなファイルシステム100からデータが削除されたとき、そのデータへのポインタは上記ファイルシステムから削除される。しかし、そのデータ自体は、もしスナップショット130の中に含まれれば、保存される(実際は、上記スナップショット自体が削除されたとき、このデータは削除できる。 )。

30

【0040】

実際には、ルートiノード110、他のiノード、およびファイルシステム100に対する多くの変更に関するデータへの変更は、1以上の記憶装置に書き込まれる前に、累積される。そのような変更が書き込まれた後では、ファイルシステム100は自己整合的である(すなわち整合点で)。好ましくは、スナップショットはそのような整合点でのみ作成される。

【0041】

本発明によれば、スナップショット130は、そのスナップショットを書き込み可能にすることにより、新しいアクティブなファイルシステムに変換可能である。書き込み可能なスナップショット内のデータを変更するために、変更されたデータが1以上の記憶装置に書き込まれる。上記変更されたデータを指すルートiノード140および介在するiノードならびに組織データは更新される。さらに、上記データの変更されなかったコピーは、もしファイルシステム100内にまだ含まれれば、保存される。この処理は、ファイルシステム100に変更が加えられたときに発生する処理と実質的に同一であり、保存された上記未変更データのみが、ルートiノード110が指すデータである。

40

【0042】

すなわち、第1のアクティブなファイルシステム(すなわち、ファイルシステム100)に変更が加えられたとき、変更されたデータが、第1のアクティブなファイルの、第2のアクティブなファイル(例えば、書き込み可能なスナップショット130)と共有されない場所に記録される。同様に、第2のアクティブなファイルシステムに変更が加えられ

50

たとき、変更されたデータが、第2のアクティブなファイルの、第1のアクティブなファイルと共有されない場所に記録される。結果として、第1のアクティブなファイルに加えられた変更は第2のアクティブなファイルに反映されず、第2のアクティブなファイルに加えられた変更は第1のアクティブなファイルに反映されない。

【0043】

スナップショット130は、作成されるとき、実質的にファイルシステム100と重なる。もし上記スナップショットがその作成のすぐ後に書き込み可能にされたら、上記書き込み可能なスナップショットより形成された新しいアクティブなファイルは、当初は元のアクティブなファイルシステムとそのデータのほとんどすべてを共有する。結果として、上記発明は、処理時間および記憶容量などの資源の効率的利用を用いて、新しいアクティブなファイルシステムの作成を可能にする。

10

【0044】

変更されたデータを記憶する処理および変更されなかったデータを保存する処理は、ファイルシステム100およびスナップショット130（読み取り専用か書き込み可能かに関わらず）を互いに分岐させる。この分岐は、ファイルシステム100およびスナップショット130の間の重なりを減少により、図2に示されている。

【0045】

図3は、図2のアクティブなファイルシステムおよびスナップショットの間の関係を示す。この種の図は、複数のファイルシステムおよびそれらのスナップショットの間の簡略化された図を提供する。図3では、ファイルシステム100はスナップショット130を指す。加えてファイルシステム100およびスナップショット130の両方は他のスナップショット（図示されていない）を指す。

20

【0046】

図4は、本発明によりアクティブなファイルシステムに変換可能なスナップショットの連鎖を示す。この図では、第2のスナップショットがファイルシステム100から作成される。上記第2のスナップショット作成時にスナップショット100はまだスナップショット130を指しているため、スナップショット180は、スナップショット130を指すスナップショットデータ190を含む。

【0047】

スナップショット130およびスナップショット180のいずれかまたは両方は、それらのスナップショットを書き込み可能にすることにより、アクティブなファイルシステムになることができる。1つのデータが上記アクティブなファイルシステム（すなわち、ファイルシステム100、書き込み可能なスナップショット130または書き込み可能なスナップショット180）のいずれかに書き込まれると、上記ファイルシステムは互いに分岐する。

30

【0048】

図5は、図4の上記アクティブなファイルシステムと上記スナップショットとの間の関係を示す。図5では、ファイルシステム100はスナップショット130およびスナップショット180を指す。同様に、スナップショット180はスナップショット130を指し、スナップショット130は次に他の1以上のスナップショットを指すことができる。

40

【0049】

図6は、本発明によりアクティブなファイルシステムに変換されたスナップショットを示す。この図では、スナップショット180は、書き込み可能にされることにより、アクティブなファイルシステム180'になった。この新しいアクティブなファイルシステムは変更可能であるので、もうファイルシステム100の真のスナップショットを表さない。結果として、ファイルシステム100のスナップショットデータ150の中のスナップショット180を指すスナップショットポインタは、例えば削除することにより、切断されている。

【0050】

図7は、図6のアクティブなファイルシステム、新しいアクティブなファイルシステム

50



およびスナップショットの間の関係を示す。この図では、アクティブなファイルシステム 100 はスナップショット 130 を指す。同様に、アクティブなファイルシステム 180 ' はスナップショット 130 を指す。上に説明したように、ファイルシステム 100 は、好ましくは、もはや、スナップショット 180 へのスナップショットポイントを含まない。しかし、ファイルシステム 100 は、例えば 1 つのファイルシステムから他のファイルシステムへ移れるように、ファイルシステム 180 ' へのポイントをまだ含むことができる。このファイルシステム間ポイントは、図 7 の破線で示されて、スナップショットポイントから区別される。

#### 【0051】

本発明によりアクティブなファイルシステムに変換可能なスナップショットのより複雑な連鎖を示す。図 8 では、ファイルシステム 800 はアクティブなファイルシステムである。4 つのスナップショットがこのファイルシステムから作成されている。スナップショット 810 が最も古く、スナップショット 820 は次に古く、スナップショット 830 はスナップショット 820 の次に古く、スナップショット 840 は最も新しい。スナップショット 810 より古いスナップショットが削除されていて、これにより、いずれかの他のスナップショットまたはアクティブなファイルシステムと重ならないデータにより占有される記憶容量を開放する。スナップショット 810 から 840 の各々は、書き込み可能にされることにより、アクティブなファイルシステムになることができる。

#### 【0052】

図 9 は、図 8 に示されている連鎖と、本発明によりアクティブなファイルシステムに変換されたスナップショットの 1 つとを示す。

#### 【0053】

図 9 では、スナップショット 830 は、データが変更可能、追加可能、および削除可能なアクティブなファイルシステム 830 ' に変換された。結果として、ファイルシステム 800 は、好ましくは、もう、スナップショットとしてスナップショット 830 を指さない。アクティブなファイルシステム 830 ' はスナップショット 810 および 820 を指し続けていられる。

#### 【0054】

図 10 は、複数のアクティブなファイルシステムとそれらに関連するスナップショットとの間の、本発明による他の可能な関係を示す。

#### 【0055】

追加のスナップショットが図 9 のアクティブなファイルシステムから作成されたことを除いて、図 10 の上部は図 9 に対応する。それ故、スナップショット 1000 はファイルシステム 800 から作成されていて、スナップショット 1010 はファイルシステム 830 ' から作成されている。さらに、スナップショット 810 は、1 以上の記憶装置上の容量を開放するために削除されている。

#### 【0056】

アクティブなファイルシステム 800 および 830 ' の両方とも、共通のスナップショット 820 にさかのぼれる。しかし、そのスナップショットが削除されるとき、上記アクティブなファイルシステムは、もはや、共通のスナップショットを共有しない。この状況はファイルシステム 1020 およびスナップショット 1030 から 1050 に関して発生する。この配置は、すべて 1 つの記憶装置または 1 組の記憶装置で、複数のアクティブなファイルシステムとそれらに関連するスナップショットとの間のリンクにより形成される「森」（すなわち、接続されていない木の集まり）をもつことが可能であることを示す。上記ファイルシステムおよびそれらのスナップショットが、もはや、共通のスナップショットを指していないという事実にもかかわらず、これらのスナップショットおよび上記アクティブなファイルシステムでさえ、まだ、データを共有している（すなわち、重なっている）ので、これにより本発明の効率性を維持している。

#### 【0057】

以下の説明では、新しいアクティブなファイルシステムはスナップショットから作成さ

10

20

30

40

50

れる。しかし、本発明は新しいアクティブなファイルシステムを作成するために、スナップショットの実際の作成を要求していない。むしろ、要求していることは、スナップショットの中で見出される構造のラインに沿った構造、すなわち、スナップショットのルート i ノードの中で見出される構造のラインにしたがった組織データの作成と、上記組織データにより指されたデータの保存である。

【 0 0 5 8 】

さらに、本発明は上に説明した特定の配置に限定されない。むしろ、それらの配置は、アクティブなファイルシステム、スナップショットおよび新しいアクティブなファイルシステムの間の関係の複数の可能な型を示す。他の配置は可能であり本発明の範囲に含まれる。

【 0 0 5 9 】

システム要素

【 0 0 6 0 】

図 1 1 は、本発明による複数のアクティブなファイルシステムを含む記憶システムのブロック図を示す。

【 0 0 6 1 】

システム 1 1 0 0 は、少なくとも 1 つのファイルシステムプロセッサ 1 1 1 0 (すなわち、コントローラ) およびハードディスクドライブまたは C D - R O M ドライブなどの少なくとも 1 つの記憶装置 1 1 2 0 を含む。このシステムは、また、好ましくは、情報を受信するための、少なくとも 1 つのコンピュータ・デバイスまたは、ネットワークへのインタフェース 1 1 3 0 を含む。別の実施の形態では、プロセッサ 1 1 1 0 は、インタフェース 1 1 3 0 経由で記憶システムに接続されたコンピュータ・デバイスのためのプロセッサである。

【 0 0 6 2 】

プロセッサ 1 1 1 0 は、ここで説明するように、プログラムおよびデータメモリに制御されて、上記ファイルシステムに関連するタスクを実行する。上記プログラムおよびデータメモリは、記憶装置 1 1 2 0 で演算を実行する制御プロセッサ 1 1 1 0 のための(さらに可能ならば、プロセッサ 1 1 1 0 と協調して記憶装置 1 1 2 0 を制御するための)適切なソフトウェアを含む。

【 0 0 6 3 】

好ましい実施の形態では、少なくとも 1 つのそのような記憶装置 1 1 2 0 が 1 以上のブートレコード 1 1 4 0 を含む。各ブートレコード 1 1 4 0 は、アクティブファイルシステムに関するファイルシステム階層におけるルートデータブロック(すなわち、i ノード)を指定する 2 以上(好ましくは 2)のエントリを含む。1 つのアクティブなファイルがある場合は、好ましくは 1 つのそのようなブートレコードがあり、1 より多いそのようなアクティブなファイルがある場合は、好ましくは 1 より多いそのようなブートレコードがある。

【 0 0 6 4 】

上で述べたように、1 以上のアクティブなファイルシステムが記憶装置 1 1 2 0 に存在しうる。そのような場合、上記ファイルシステム維持装置(すなわち、プログラムの制御で動作するプロセッサ 1 1 1 0)は、好ましくは、そのようなアクティブなファイルシステムごとに、1 より多いブートレコードを指定し、整然と維持する。

【 0 0 6 5 】

読み取り専用スナップショットも記憶装置 1 1 2 0 の中に存在できる。この場合、アクティブなファイルシステムからスナップショットへのポインタおよびスナップショットから他のスナップショットへのポインタは、上に説明したように、上記記憶装置の中に記憶される。

【 0 0 6 6 】

高い可用性

【 0 0 6 7 】

本発明により、複数の同時動作のファイルサーバにより使用されている複数の同時にアクティブなファイルシステムを含むファイルシステムクラスタのブロック図を示す。

【0068】

1つのファイルシステムクラスタは複数のファイルシステムプロセッサ1200および1以上のファイルシステムディスク1210を含む。好ましい実施の形態では、そのような各プロセッサ1200は、ファイルサーバとして動作するために配置されていて、例えば既知のファイルサーバプロトコルを使用して、ファイルサーバ要求を受け付け、ファイルサーバ応答を返すことができる。好ましい実施の形態では、1以上のファイルシステムディスク1210はそのような複数のディスクを含むので、全体の高い可用性クラスタに関し、どの個々のディスク1210も単一位置の故障を発生させない。好ましくは本発明と共に使用されるW A F L ( W r i t e   A n y w h e r e   F i l e   S y s t e m ) はそのような配置を組み込む。

10

【0069】

上に説明したように、複数のプロセッサ1200は、複数の並列に書き込み可能なアクティブなファイルシステムを、それらの並列に書き込み可能なアクティブなファイルシステムに関連するスナップショットとともに維持できる。上記アクティブなファイルシステムおよびスナップショットは、同じ1組のディスク1220に維持可能である。こうして、1組のプロセッサ1200および1組のディスク1220は、資源の実質的な重複の必要なしに高い可用性クラスタを提供可能である。

20

【0070】

別の実施の形態

【0071】

本発明は、複数のアクティブなファイルシステムを作成および維持するための方法だけでなく、本発明の方法を実行するソフトウェアおよび/または1以上の記憶装置などのハードウェア、そして種々の他の実施の形態において、具体化できる。

【0072】

以上の説明では、本発明の好ましい実施の形態は、好ましい処理ステップおよびデータ構造に関して説明されている。しかし、当業者は、本出願の熟読後、本発明の実施の形態が、プログラムの制御で動作する特定の処理ステップおよびデータ構造に適応可能な1以上の汎用プロセッサまたは専用プロセッサを使用して実行可能であること、そのような処理ステップおよびデータ構造はメモリ（例えば、D R A M、S R A M、ハードディスク、キャッシュなどの固定メモリおよびフロッピーディスク、C D - R O M、データテープなどの着脱式メモリ）に記憶されるか入出力され、そのようなプロセッサで実行可能な命令（例えば、直接実行可能なオブジェクトコード、コンパイル後に実行可能なソースコード、インタプリタを使用して実行可能なコードなど）を含む情報により実施可能であること、また、そのような装置を使用してここで説明されている好ましい上記処理ステップおよびデータ構造の実施は必要以上の実験またはさらなる発明を必要としないことを理解するだろう。

30

【図面の簡単な説明】

【0073】

40

【図1】本発明によりアクティブなファイルシステムに変換可能なスナップショットの作成を示す図

【図2】アクティブなファイルシステムのスナップショットからのアクティブなファイルシステムの分岐を示す図

【図3】図2におけるアクティブなファイルシステムおよびスナップショットの間の関係を示す図

【図4】本発明によりアクティブなファイルシステムに変換可能なスナップショットの連鎖を示す図

【図5】図4におけるアクティブなファイルシステムおよびスナップショットの間の関係を示す図

50

【図 6】本発明によりアクティブなファイルシステムに変換可能なスナップショットを示す図

【図 7】図 6 におけるアクティブなファイルシステム、新しいアクティブなファイルシステム、およびスナップショットの間の関係を示す図

【図 8】本発明によりアクティブなファイルシステムに変換可能なスナップショットのより複雑な連鎖を示す図

【図 9】図 8 に示されている連鎖と本発明によりアクティブなファイルシステムに変換されたスナップショットの 1 つとを示す図

【図 10】本発明による複数のアクティブなファイルシステムおよびそれらに関するスナップショットの間のより多くの別の可能な関係を示す図

10

【図 11】本発明による複数のアクティブなファイルシステムを含む記憶装システムのブロック図

【図 12】本発明により、複数の同時作動のファイルサーバにより使用される複数の同時にアクティブとなるファイルシステムを含むファイルシステムクラスタのブロック図

【符号の説明】

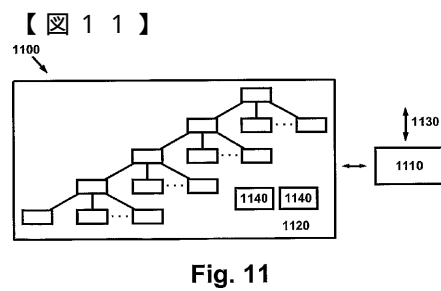
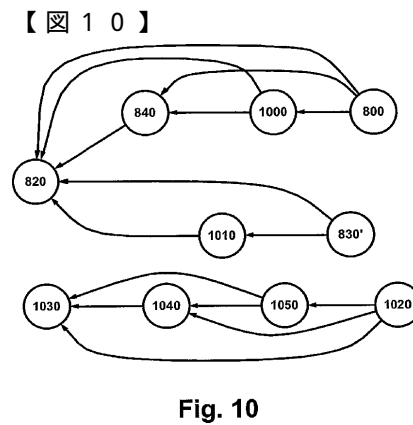
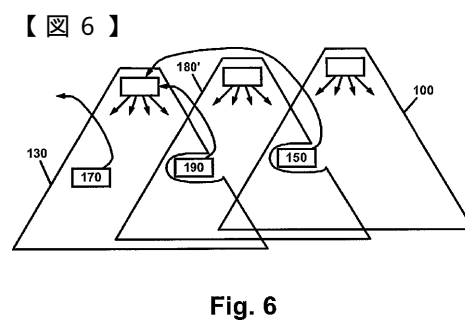
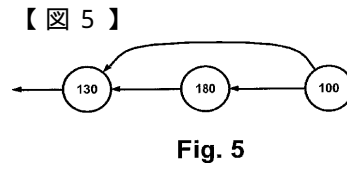
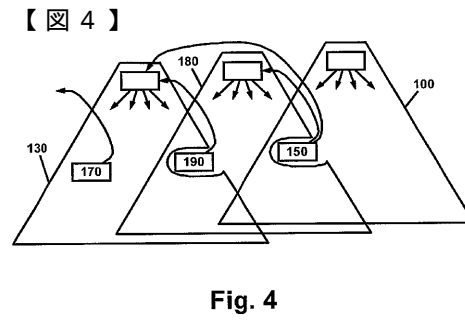
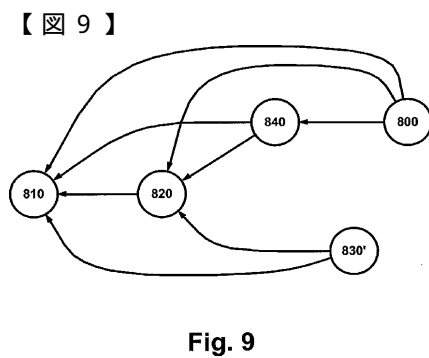
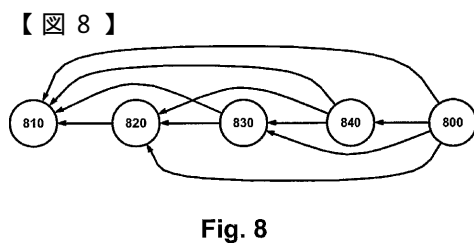
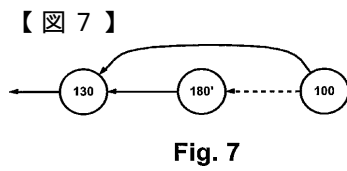
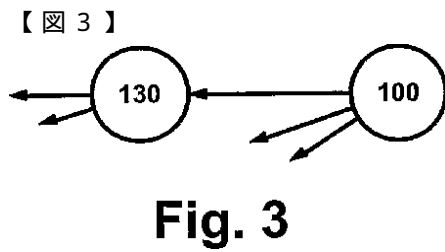
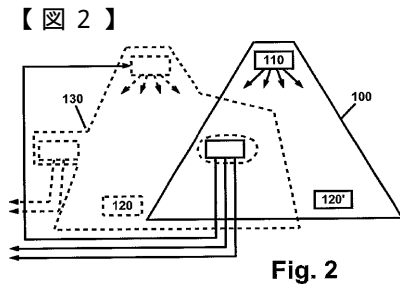
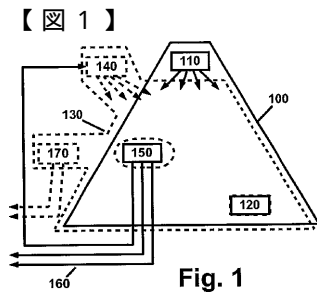
【0074】

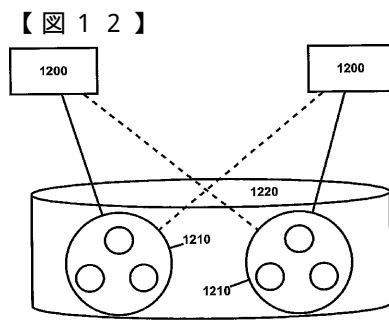
100 ファイルシステム  
 110 ルート i ノード  
 120 (変更されなかった) データ  
 120' 変更されたデータ  
 130 スナップショット  
 140 スナップショットのルート i ノード  
 150 ファイルシステム (100) のスナップショットデータ  
 160 スナップショットデータ内のポインタ  
 170 スナップショット 130 の他のスナップショットを指すスナップショットデータ  
 180 スナップショット  
 180' アクティブなファイルシステム  
 190 スナップショット (130) を指すスナップショットデータ  
 800 ファイルシステム  
 810 スナップショット  
 820 スナップショット  
 830 スナップショット  
 830' アクティブなファイルシステム  
 840 スナップショット  
 1000 スナップショット  
 1010 スナップショット  
 1020 ファイルシステム  
 1030 スナップショット  
 1040 スナップショット  
 1050 スナップショット  
 1100 システム  
 1110 ファイルシステムプロセッサ  
 1120 記憶装置  
 1130 インタフェース  
 1140 ブートレコード  
 1200 ファイルシステムプロセッサ  
 1210 ファイルシステムディスク  
 1220 ディスク

20

30

40



**Fig. 12**

---

フロントページの続き

(72)発明者 ジョン・エドワーズ

アメリカ合衆国 9 4 0 8 7 - 2 0 7 6 カリフォルニア州サニーベイル、克蘭ダノ・コート 1 1 7  
3 番

(72)発明者 ブレイク・ルイス

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・ヒルズ、ピア・フェリス 2 7 8 7 8  
番

審査官 田川 泰宏

(56)参考文献 実良 栗富, ITプロのための解説, 日経コンピュータ, 日本, 日経BP社, 2002年 5月  
20日, no. 548, p.158-164

(58)調査した分野(Int.Cl., DB名)

G06F 12/00