

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号
特許第4965210号
(P4965210)

(45) 発行日 平成24年7月4日(2012.7.4)

(24) 登録日 平成24年4月6日(2012.4.6)

(51) Int.Cl.
G06F 12/00 (2006.01)

F I
G06F 12/00 531J

請求項の数 1 (全 16 頁)

(21) 出願番号	特願2006-261837 (P2006-261837)	(73) 特許権者	390009531
(22) 出願日	平成18年9月27日 (2006. 9. 27)		インターナショナル・ビジネス・マシーンズ・コーポレーション
(65) 公開番号	特開2007-95064 (P2007-95064A)		INTERNATIONAL BUSINESS MACHINES CORPORATION
(43) 公開日	平成19年4月12日 (2007. 4. 12)		アメリカ合衆国10504 ニューヨーク州 アーモンク ニュー オーチャードロード
審査請求日	平成21年5月20日 (2009. 5. 20)		
(31) 優先権主張番号	11/236450	(74) 代理人	100108501
(32) 優先日	平成17年9月27日 (2005. 9. 27)		弁理士 上野 剛史
(33) 優先権主張国	米国 (US)	(74) 代理人	100112690
			弁理士 太佐 種一
		(74) 代理人	100091568
			弁理士 市位 嘉宏
		最終頁に続く	

(54) 【発明の名称】 ファイル・システムの稠密診断データを取得および送信するためにコンピュータに実施させる方法

(57) 【特許請求の範囲】

【請求項 1】

ファイル・システムの稠密診断データを取得および送信するためにコンピュータに実施させる方法であって、

(a) 診断ファイル・システムを生成するステップと、

(b) 前記診断ファイル・システムをマウントするステップと、

(c) 前記ファイル・システムのファイル・システム・メタデータを前記診断ファイル・システムに書き込むステップとを有し、

前記ステップ (c) が、

(c 1) 前記ファイル・システムにおいて複数のデータ単位を識別するステップと、

(c 2) 前記複数のデータ単位のうち一のデータ単位をトラバースして、当該一のデータ単位が前記ファイル・システム・メタデータを含むかどうかを判定するステップと、

(c 3) 前記一のデータ単位が前記ファイル・システム・メタデータを含むとの判定に
応答して、前記ファイル・システム・メタデータを前記一のデータ単位から前記診断ファイル・システムに書き込むステップと、

(c 4) 前記一のデータ単位がユーザ・データを含むかどうかを判定するステップと、

(c 5) 前記一のデータ単位がユーザ・データを含むとの判定に
応答して、スパース・ファイルを前記診断ファイル・システムに書き込むステップと、

(c 6) 前記複数のデータ単位のうち他のデータ単位がトラバースされていないとの判定
に
応答して、当該他のデータ単位の各々ごとに、前記ステップ (c 2) ないしステップ (

10

20

c 5) を繰り返すステップとを有し、
(d) 前記ステップ (c) が完了した段階で、前記診断ファイル・システムを圧縮して稠密診断データを形成するステップと、
(e) 前記稠密診断データをオフサイトのシステムに送信するステップとをさらに有する方法。

【発明の詳細な説明】

【技術分野】

【 0 0 0 1 】

本発明は一般に、ファイル・システムの整合性を維持および改善することに関する。より詳細には、本発明は、ファイル・システム・メタデータを稠密フォーマットで取得して、そのようなデータが専門家の検査および修復のためにオフサイトに送られ得るようにするための方法に関する。

10

【背景技術】

【 0 0 0 2 】

コンピュータ分野では、ファイル・システムは、ファイルを見つけアクセスし易くするように、コンピュータ・ファイルおよびそれらが含むデータを保存し組織化するための構造である。ファイル・システムは、ハード・ディスクまたはCD-ROMなどのストレージ装置を使用することができ、ファイルの物理的ロケーションを維持することを必要とする。代替として、ファイル・システムは、仮想ファイル・システムとすることができる。仮想ファイル・システムは、ネットワークを介した仮想データまたはデータのためのアクセス方式としてのみ存在する。

20

【 0 0 0 3 】

ファイル・システムは、2つのタイプのデータから作成される。通常大部分を占める第1のタイプのデータは、ユーザ・データである。ユーザ・データの内容タイプは、例えば、テキスト、グラフィックス、音楽、およびコンピュータ命令とすることができる。第2のタイプのデータは、ファイル・システム・メタデータである。ファイル・システム・メタデータは、ユーザ・データ以外のすべてのデータである。メタデータは、ユーザ・データを含むファイルの統計および構造を処理システムに通知する。

【 0 0 0 4 】

ファイル・システム・メタデータは、スーパーブロックを含み、スーパーブロックは、ファイル・システム概要を提供し、その他の情報へのポインタを含む。i ノードは、各ファイルに関連付けられたファイル・システム・メタデータである。i ノードは、バイト数によるファイル長、関連する装置識別子、ユーザ識別子、グループ識別子、一意とし得るi ノード番号、ファイル・モード、タイムスタンプ、および参照カウントを示す。

30

【 0 0 0 5 】

ブロックは、ファイルに割り当てられ得る最小単位のディスク・ストレージである。例えば、プロセッサは、特定のファイル・システムにおいて、1ブロックを1024バイトに設定することができる。これは、ファイルがほぼ常に1つ以上のブロックを満たし、最終ブロックはデータによって部分的にのみ占有されることを意味する。

【 0 0 0 6 】

40

ファイルの部分は、複数のブロックに存在することができ、時にはディスク・ドライブ内に散在させられる。i ノードは、i ノード内で複数のブロックをリストすることができる。より大きなファイルの場合、i ノードは、間接ブロックを含むことができ、間接ブロックは、さらなるブロックのリストをポイントすることができる。しばしば、これは、後続のより深いブロックの層をポイントする多数のレベルの間接ブロックからなるツリー状構造を生じさせる。

【 0 0 0 7 】

ファイル・システムは、非常に大きなデータ構造になる傾向にある。プロセッサがファイル・システムに変更を加える場合、プロセッサはしばしば、多くの別個の書き込み操作を必要とする。ときには、エラーまたはその他の障害、例えば、電源障害の発生などは、

50

一連の書き込みを途中で妨げることがある。

【0008】

こうした状況でプロセッサがエラーに遭遇すると、競合条件が発生することがある。競合条件は、1つの電子装置における2つのイベントが、どちらが装置の状態または出力に影響を与えるかを識別するために互いに本質的に競合するときに発生するものであり、その場合、最初に到着したイベントまたは信号が、装置の状態を制御する。ジャーナル・ファイル・システムなどのファイル・システムの文脈では、ファイル・システムを更新する場合、1) ファイルのディレクトリ・エントリを削除し、2) そのファイルのiノードをフリー・スペース・マップ内でフリー・スペースとしてマークするという2つのステップが発生する。

10

【0009】

電源が障害を起こし、ステップ1がクラッシュの直前に発生した場合、孤立したiノードが存在するようになり、実際に割り当てられているよりも多くのブロックが、ストレージに割り当てられているように見える。ステップ2がクラッシュの直前に発生した場合、まだ削除されていないiノードがフリーとしてマークされ、おそらく何か他のものによって上書きされることになる。

【0010】

特定のタイプのファイル・システムであるジャーナル・ファイル・システムは、さらなる障害モードを有する。前述の2つのステップに加えて、ジャーナル・ファイル・システムは、トランザクションのために行われた変更をコミットする第3のステップを有する。ものが正常に機能した場合、プロセッサは、トランザクションのすべてをジャーナル・ログにコミットするか、さもなければ、トランザクションの何もジャーナル・ログにコミットしない。ジャーナル・ファイル・システムでは、プロセッサは、メタデータを整合性のある状態に確立するため、ジャーナル・ログを再生することができる。しかし、ジャーナル・ファイル・システムは、ある書き込みが失敗したのに、プロセッサがトランザクションのその他の部分をジャーナル・ログに書き込んだときに発生するようなI/Oエラー処理の失敗が起こると、不整合になることがある。

20

【0011】

ファイル・システムの保守および回復を遠隔地にアウトソーシングすることに伴う危険は、送られるファイル・システムの詳細が傍受され得ることである。これは、データがインターネットを介して送信される場合に特に当てはまる。暗号化データであっても、傍受されれば、危険にさらされ得る。

30

【発明の開示】

【発明が解決しようとする課題】

【0012】

したがって、可能であればインターネットを介した機密情報の送信を回避することが長年にわたって望まれている。

【課題を解決するための手段】

【0013】

本発明の側面は、ファイル・システムの稠密診断データを取得および送信するためにコンピュータに実施させる方法を提供する。プロセッサは、ファイル・システム内の各データ単位を識別する。プロセッサは、データ単位がファイル・システム・メタデータを含むかどうかを判定する。データ単位がファイル・システム・メタデータを含むと判定された場合、プロセッサは、メタデータのデータ単位を診断ファイルに書き込む。次にプロセッサは、データ単位がユーザ・データを含むかどうかを判定する。データ単位がユーザ・データを含む場合、プロセッサは、スパース・オブジェクト(sparse object)を診断ファイルに書き込む。

40

【0014】

本発明の特色と考えられる新規な特徴が、添付の特許請求の範囲において説明される。しかし、本発明自体は、本発明の好ましい使用方法、さらなる目的、および利点と同様、

50

本発明の例示的な実施形態についての以下の詳細な説明を添付の図面と併せて読みながら参照することによって最も良く理解されよう。

【発明を実施するための最良の形態】

【0015】

ここで図面を、特に図1を参照すると、本発明の実施形態を実施し得るデータ処理システムのブロック図が示されている。示された例では、データ処理システム100は、ノース・ブリッジおよびメモリ・コントローラ・ハブ(MCH)108と、サウス・ブリッジおよび入力/出力(I/O)コントローラ・ハブ(ICH)110とを含むハブ構成を利用する。処理ユニット102、メイン・メモリ104、およびグラフィックス・プロセッサ118は、ノース・ブリッジおよびメモリ・コントローラ・ハブ108に接続する。グラフィックス・プロセッサ118は、アクセラレーテッド・グラフィックス・ポート(AGP)またはグラフィックス・プロセッサ118を介して、ノース・ブリッジおよびメモリ・コントローラ・ハブ108に接続してもよい。

10

【0016】

示された例では、ローカル・エリア・ネットワーク(LAN)アダプタ112、オーディオ・アダプタ116、キーボードおよびマウス・アダプタ120、モデム122、読み取り専用メモリ(ROM)124、ユニバーサル・シリアル・バス(USB)ポートおよびその他の通信ポート132、ならびにPCI/PCIeデバイス134が、バス138を介して、サウス・ブリッジおよびI/Oコントローラ・ハブ110に接続する。PCI/PCIeデバイスは、例えば、イーサネット(登録商標)・アダプタ、アドイン・カード、およびノートブック・コンピュータ用のPCカードを含むことができる。PCIはカード・バス・コントローラを使用するが、PCIeは使用しない。ROM124は、例えば、フラッシュ基本入力/出力システム(BIOS: basic input/output system)とすることができる。

20

【0017】

ハード・ディスク・ドライブ126、およびCD-ROMドライブ130は、バス140を介して、サウス・ブリッジおよびI/Oコントローラ・ハブ110に接続する。ハード・ディスク・ドライブ126、およびCD-ROMドライブ130は、例えば、統合ドライブ電子(IDE)またはシリアル・アドバンスド・テクノロジー・アタッチメント(SATA)インタフェースを使用することができる。スーパーI/O(SIO)デバイス136は、サウス・ブリッジおよびI/Oコントローラ・ハブ110に接続されることができる。

30

【0018】

オペレーティング・システムは、処理ユニット102上で動作し、図1のデータ処理システム100内の様々なコンポーネントを調整および制御する。オペレーティング・システムは、IBM(IBMはIBM Corporationの商標)コーポレーションから入手可能なアドバンスド・インタラクティブ・エグゼクティブ(AIX(R))などの市販のオペレーティング・システムとすることができる。AIXは、IBM Corporationの登録商標である。

【0019】

本発明の様々な実施形態は、重荷となるユーザ・データを送信する面倒を伴わずに、ファイル・システム不整合の詳細だけを送信することを可能にする。ユーザ・データは、ファイル・システム・データの大部分をなす傾向にあり、したがって、場所確保形式(place-keeping form)で以外は送信されない。ファイル・システム整合性チェック・プログラムは、ファイル・システム内の複数のデータ単位を識別する。識別するための1つの方法は、ファイル・システム整合性チェック・プログラムが、1つのデータ単位から別のデータ単位に移動しながら、ファイル・システムのツリー状構造を識別することである。いくつかの実施形態は、ブロックと同じ大きさのデータ単位を使用することができる。別の実施形態は、例えば、セクタなど、より大きなデータ単位を使用することができる。データ単位がファイル・システム・メタデータまたはメタデータを含む場合、ソフトウェアは、

40

50

メタデータの単位を、診断ファイル・システムの一部とすることができる診断ファイルに書き込みまたはコピーする。データ単位がユーザ・データを含む場合、ソフトウェアは、スパース・オブジェクトを診断ファイルに書き込みまたはコピーする。その後、診断ファイルをバックアップまたはその他の方法で統合するステップは、スパース・オブジェクトを、元のユーザ・ファイルのコンパクトではあるが完全に空で何も書き込まれていないバージョンとして扱う。スパース・オブジェクトは、ファイルのメタデータ内のヌル・ポインタとすることができる。このポインタは指示対象をポイントしないので、保存された基となるデータがなく、単に空またはすべてが「0」のデータ単位を表すにすぎない。したがって、そのようなスパース・オブジェクトは、ある意味で、ファイル・システムを構成するディスクまたはその他の媒体上のデータ単位内の「0」領域の代わりに用いられる場所確保データとして動作する。したがって、処理システムは、ソース・ファイル・システムと比べて非常に縮小された形式で各スパース・オブジェクトを送信することができる。

10

【0020】

図2は、本発明の一実施形態に従って、顧客システム201がどのように専門家システム205と対話することができるかを示している。顧客システム201は、例えば、図1の処理システム100を使用して動作することができる。管理者は、ファイル・システムが異常な挙動をしていると判定することができる。それに応じて、管理者は、稠密診断データを含む稠密診断ファイルを生成するために、顧客システム201を操作することができる。

【0021】

20

専門家システム205は、顧客システム201から診断ファイル・システムを受信することができる。専門家システム205は、診断ファイル・システムにサービスを提供して、サービス済ファイル・システムを形成する。専門家システム205は、例えば、図1の処理システム100を使用して動作することができる。サービス提供は、例えば、診断ファイル・システム上で保守を実行して、サービス済ファイル・システムを形成することを含むことができる。加えて、サービス提供は、診断ファイル・システム上で回復動作を実行することを含むことができる。請求システム207は、診断ファイル・システムにサービス提供した際に提供されたサービスに対する請求を顧客に行うことができるように、専門家システム205上の動作を監視することができる。専門家システム205は、顧客システム201にサービス済ファイル・システムを返送する。

30

【0022】

稠密ファイル・フォーマットは、スパース・ファイルが圧縮されたときに生じる結果である。管理者は、とりわけ、以下で説明されるようにユーザ・データを縮小するファイル整合性チェックを実行するため、コマンド・ライン・エディタでコマンドを入力することができる。プロセッサは、ユーザ・データに格納されている内容のタイプさえ隠蔽する程度までユーザ・データを縮小することができるので、結果的な1つ以上のファイルは稠密であり、ファイル・システム・エラーのソースにより直接的に関係するメタデータを有する。したがって、稠密診断データという用語は、本発明の様々な例示的な実施形態から生じる1つ以上の圧縮されたファイルを示すのに利用される。言い換えると、稠密診断データは、ユーザ・データを除去するかまたはユーザ・データの代わりに場所確保データを使用し、次にその結果の1つ以上のファイルを圧縮した後に生成される。

40

【0023】

顧客システム201は、例えば、インターネット203などのネットワークを介して、稠密診断データを送信する。稠密診断データは、専門家システム205に到着する。専門家システム205は、データ回復専門家の制御下にあることができる。データ回復専門家は、管理者との確立された信用関係を有していないこともある。加えて、インターネット203は、安全でないデータ送信手段として知られている。

【0024】

図3は、本発明の一実施形態による、ソース・ファイル・システムから診断ファイルへの変換を仲介するオペレーティング・システムの動作を示している。オペレーティング・

50

システム 301 は、例えば、図 1 の処理ユニット 102 などのプロセッサ上で動作することができる。ソース・ファイル・システム 303 は、例えば、図 1 のハード・ディスク・ドライブ 126 などに保存されることができる。管理者は、コマンド・ライン・エディタ 311 を使用して、コマンドを個々にオペレーティング・システム 301 に入力することができる。1 つ以上のコマンドが、診断ファイル 305 を生成することができる。加えて、ユーティリティ・プログラムが、コマンド・ライン・エディタ 311 の代わりに使用されることができ、その場合、ユーティリティ・プログラムは、オペレーティング・システムにコマンドを書き込む。

【0025】

図 4 は、従来のファイル・フォーマットの一例を示している。シーク・コマンドを使用するとき、ファイル・ポインタが、例えば、「0」ブロック 413 などの物理ディスク・アドレス上を通過するとしても、ファイル 411 に割り当てられた各ビットは、「0」に設定される。このフォーマットでは、プロセッサは、ファイルのすべての割り当てビットおよびブロックを物理的に書き込み、物理レベルにおいて圧縮は利用されない。

【0026】

図 5 は、診断ファイルに関連し得る、スパース・オブジェクトと呼ばれることもある、スパース・ファイル・フォーマットの一例を示している。このフォーマットは、各ビットの論理書き込みを実行するプロセッサを伴うが、実際には、プロセッサは、ファイル・ポインタが指示するデータだけを物理的に保存する。ファイル・ポインタは、次の物理書き込みが生じることになるロケーションである。プロセッサがファイル書き込みの間にギャップが存在することを許可する場合、例えば、ブロック・ポインタ・テーブル 414 などのメタデータ内に注記が作成される。プロセッサは、「0」からなるブロックが論理的に記録されるヌル・ポインタ 415、416、417、および 418 を設定することによって、ブロック・ポインタ・テーブル 414 に注記を作成する。言い換えると、各ヌル・ポインタは、さもなければハード・ドライブに物理的に書き込まれていた「0」からなるブロックのための一種の場所確保データとして動作する。「0」からなる大きな領域は後の読み取りにおいて回復され得るので、このフォーマットはスパース (sparse) である。さらに、ファイル・メタデータは、「0」からなる領域をユーザ・データに保存するのではなく、「0」ブロックへの参照が保持される場所である。「0」からなる領域または「0」ブロックは、スパース・オブジェクトによって表され、いくつかのファイル・システムは、例えば、プログラムがオペレーティング・システムに読み取りを要求したときなど、オペレーティング・システムの通常の動作中に、「0」からなる大きな領域をスパース・オブジェクトから読み取ることができる。

【0027】

稠密ファイル・フォーマットは、スパース・ファイルが圧縮されたときに生じる結果である。プロセッサは、各「0」ブロックを著しく圧縮することができる。いくぶん無秩序な「1」と「0」のシーケンスから構成されるデータも同様に圧縮されるが、おそらくそれほどは圧縮されない。したがって、「0」ブロックは、ファイルのすべてが圧縮形式で物理的に書き込まれるように、いくぶん無秩序なデータと一緒に圧縮される。診断データとして始まった 1 つ以上のファイルは、圧縮が完了すると、稠密診断データになることができる。

【0028】

図 6 は、本発明の一実施形態による、ファイル・システムに関する診断データを取得する最初のステップを示している。管理者は、コマンド・ラインを入力して、各ステップを実行するようオペレーティング・システムに命令することによって、これらのステップを実行するようプロセッサに命令することができる。自動的に次々と実行されるコマンドのスクリプトにすべてのステップを統合することも、等しく適している。オペレーティング・システムは、例えば、図 3 のオペレーティング・システム 301 とすることができる。プロセッサは、診断ファイル・システムを生成するためのコマンドを受け取る (ステップ 501)。診断ファイル・システムは、図 3 の診断ファイル 305 に関連することができ

る。UNIX (The Open Groupの商標) ライクなシステムでは、新規に生成されたファイル・システムをオペレーティング・システムでアクセス可能にするには、コマンドが必要なことがある。この例示的な実施形態では、そのコマンドは「mount」である。プロセッサは、診断ファイル・システムをマウントするためのmountコマンドを受け取る(ステップ503)。これで、診断ファイル・システムは、データを収集するための準備が整う。プロセッサは、診断ファイルにメタデータを抽出するためのコマンドを受け取る(ステップ505)。そのコマンドは、例えば、ファイル・システム整合性チェックfsckとすることができる。診断ファイルは、スパース・オブジェクトを含むことができる。診断ファイルは、診断ファイル・システムの構成要素であることによって、診断ファイル・システムに関連することができる。

10

【0029】

fsckは、多くのデータを生成し、それらは、プロセッサが診断ファイル・システムを生成したときに割り当てられたストレージを埋め尽くす可能性がある。したがって、プロセッサは、抽出されたメタデータが診断ファイル・システムに収まるかどうかを調べるためにテストを行う(ステップ507)。収まらないという判定を下した場合、プロセッサは、より大きな診断ファイル・システムを生成するためのコマンドを受け取ることができる(ステップ509)。プロセッサは、ステップ509の後、診断ファイル・システムのマウントおよびその他のステップを実行することができる。

【0030】

プロセッサが抽出されたメタデータは診断ファイル・システムに収まるという判定を下した場合、処理は結合子「A」から図7へと続いて行く。

20

【0031】

ファイル・システムを生成し、ファイル・システムをマウントする代替方法が存在する。例えば、診断データを記録する目的で、専用ファイル・システムが存在することができ、稠密診断データを専用ファイル・システムに書き込む動作の前に削除されることができる。

【0032】

データ処理システムが回復を行う1つの方法は、プロセッサが次にファイル・システムをマウントするとき、ファイル・システムのデータ構造全体にわたる完全なウォーク(walk)またはトラバース(traverse)を実行することである。トラバースは、その他のステップと同様、どのような不整合でも検出し、訂正することができる。この修復を実行するのに使用される1つのツールが、ファイル整合性チェックまたはfsckコマンドである。その他のファイル整合性チェック・プログラムは、「チェック・ディスク」プログラムとしても知られるマイクロソフト社の「chkdsk」を含む。コンピュータ管理者は、ファイル・システムを有するオペレーティング・システムを利用するコンピュータでfsckを使用する。そのようなオペレーティング・システムは、例えば、操作がUNIXに類似したAIXおよびLinux(LinuxはLinus Torvaldsの商標)を含む。Linuxは、リーナス・トーバルズの登録商標である。UNIXは、オープン・グループの登録商標である。その他のオペレーティング・システムのクラスは、Microsoft(MicrosoftはMicrosoft Corporationの商標)コーポレーションによる様々な世代のWindows(WindowsはMicrosoft Corporationの商標)オペレーティング・システム、およびApple(TM)コーポレーションによるMac OSを含む。ファイル・システムの例は、IBMの拡張ジャーナル・ファイル・システムまたはJFS2などのジャーナル・ファイル・システムを含む。

30

40

【0033】

Unixライクのオペレーティング・システムのファイル・システムは、いくつかの点で冗長データ用にストレージが割り当てられているかのように振舞う、冗長データを保存するための便利な方法を提供する。その方法は、サイズが少なくとも1ブロックのスパース・オブジェクトを生成することを含む。処理システムが最初にファイルを開くとき、プ

50

ロセッサは、ファイルの先頭をポイントするファイル・ポインタを生成することができる。ファイル・ポインタは、プロセッサが次にデータを書き込むべき場所を示す場所確保データとして動作し、その場合、プロセッサは、ファイルを線形であるとして扱う。言い換えると、プロセッサは、次に書き込みコマンドを実行するとき、ファイル・ポインタがある場所₁₀に書き込みを行う。書き込みコマンドを伴わずにプロセッサがファイル・ポインタを進める場合、ファイル・システムによってファイル用に予約された領域ができるが、ファイルへのブロックの対応する割り当ては行われない。代わりに、ファイル・システム・メタデータは、拡張空ヘッダの存在を数バイトでファイルに保存するが、プロセッサは、どのデータ単位も使用されているとしてマークしない。データ単位は、例えば、ブロックとすることができる。多くのブロックがこの方法で使用されているとして注記されれば、数バイトが多くのブロックを表し得るので、実際に使用されるディスク・スペースが大幅に節約される。ファイル・ポインタを前方および反対方向に動かす一般的な方法は、Unixライクなシステムで普通に利用可能なlseekシステム・コールを使用する。その他のシステムでは、その他のファイル・ポインタ移動コマンドが存在する。

【0034】

ファイル整合性チェックが探すものの1つは、データ構造における不整合である。ファイル整合性チェック・ソフトウェアは、ファイル・システム用のデータ構造を探索することによって不整合を探す。ファイル・システム用のデータ構造は、ルートとルートから延びる1つ以上のブランチとを有するツリーから構成される。例えば、ディレクトリ構造は、ルートから始まり、ルートの下に1つ以上のディレクトリを有することができる。同様に、各ファイルは、ブロックに保存されたその構成部分を有し、ブロックも、多数のブランチおよび層を有することができ、その場合、ブロックは、各分岐点におけるノードである。さらなるブロックへのポインタまたは参照を含むブロックは、そのようなポインタをメタデータとして保存する。₂₀

【0035】

ファイル整合性チェック・ソフトウェアは、データ単位からデータ単位へトラバースすることによってツリーを探索し、各データ単位は、ブロックとすることができる。トラバースとは、ファイル整合性チェック・ソフトウェアが、ブロックをその下のさらなるブランチの存在について検査し、いくつか存在する場合、ソフトウェアが、今度は各ブランチを検査し、すべてのブランチが検査され尽くすまで、それを続けることを意味する。検査とは、多くのことを意味する。一般に、検査とは、プロセッサがさらなるブランチに沿ってさらなる参照またはポインタを探すことを意味する。₃₀

【0036】

整合性のあるファイル・システムでは、各データ単位は、そのデータ単位への参照を正確に1つだけ有し、これは、プロセッサがデータ単位を1回だけトラバースすることを意味する。やはり、データ単位は、例えば、ブロックとすることができる。しかし、不整合なファイル・システムでは、ブロックが複数回参照されることがある。ソフトウェアは、トラバースすることによってファイルの単位を識別することができる。例えば、ファイル整合性チェック・ソフトウェアは、第1のファイルをトラバースし、あるブロックが第1のファイルに割り当てられていることを注記することができる。ソフトウェアは、第2のファイル₄₀をトラバースし、その同じブロックが第2のファイルにも割り当てられていることを注記することができる。これは本当は、ファイル・システムが、事故的に同じブロックを2回割り当てたもので、第1のファイルのデータは第2のファイルのデータによって上書きされる可能性がある。すべてのファイルをトラバースすることによって、ファイル整合性チェックは、多くの問題を発見することができる。

【0037】

ファイル整合性チェックは、質問/回答形式によるユーザとの対話に応じて、問題を修復することができる。しかし、この機能は、非常に複雑かつ専門的なので、ファイル・システムの管理者は、できるだけ多くのデータを回復するため、専門家からの外部支援を求めることがある。多くの専門家は高い適格性を備えているが、残念ながら、ファイル・シ₅₀

システムの管理者と専門家の間に確立されたレベルの信用がないこともある。言い換えると、ファイル・システム上のユーザ・データが危険にさらされることから保護する機構が存在することが必要である。

【 0 0 3 8 】

ソフトウェア開発者は、fsckおよびその他のファイル整合性チェッカにコードを追加することができる。追加のコードは、各ファイルに割り当てられたブロックをトラバースするときに、追加のステップを実行する。不整合を認識するのに加えて、fsckは、各ファイルのある様相のコピーを書き込むことができる。トラバースおよび書き込みの処理は、図 6のステップ 5 0 5 で実行される。コピーは、プロセッサが診断ファイル・システムを生成したときに割り当てられたストレージを埋め尽くす可能性がある。

10

【 0 0 3 9 】

図 7 は、本発明の一実施形態による、1 つ以上のファイル・システムに関する診断データをパッケージ化して送信するステップを示している。図 3 のオペレーティング・システム 3 0 1 などのオペレーティング・システムは、図 7 のステップを実行することができる。代替として、図 7 のステップは、ユーティリティ・プログラムによって実施されることができ、その場合、ユーティリティ・プログラムは、コマンドをオペレーティング・システムに書き込む。一般に、図 7 のステップは、より小さく、ネットワークまたはその他の手段を介し転送するのがより容易なデータ構造を生成する汎用的機能を達成する。最初、オペレーティング・システムは、診断ファイル・システムを単一のファイルにバックアップまたはその他の方法で統合する（ステップ 5 1 1）。その後、オペレーティング・システムは、ファイルを圧縮することができる（ステップ 5 1 3）。オペレーティング・システムが圧縮を行った場合、それは、プロセッサが、繰り返しおよび冗長性を利用してデータを要約し、短縮された要約を保存することを意味する。例えば、4 0 9 6 個の「0」からなるブロックを記述するための、4 0 9 6 個の「0」をすべて示すのとは別の方法は、単に電子的に 4 0 9 6 × 0 を保存することであり、そのほうがはるかに短い。加えて、ヌル・ポインタがすでにメタデータ内に存在していることがあるので、単にヌル・ポインタを保存することが、圧縮として動作する。ヌル・ポインタの大きなシーケンスさえも、圧縮されることができる。例えば、1 0 0 0 個のヌル・ポインタからなるシーケンスを記述するための、物理的に各ヌル・ポインタを保存するのとは別の方法は、電子的に 1 0 0 0 × n u l l を保存することである。その後、オペレーティング・システムは、圧縮されたファイルまたは稠密診断データを送信する（ステップ 5 1 5）。したがって、オペレーティング・システムは、ユーザ・データを稠密診断データに変換する。推定上、多くの小さなファイルを有する 1 0 2 4 ギガバイトの拡張ジャーナル・ファイル・システム（J F S 2）ファイル・システムは、約 1 6 0 メガバイトの稠密診断データに変換されることができ、これは 6 0 0 0 分の 1 の圧縮である。

20

30

【 0 0 4 0 】

図 8 は、拡張ファイル整合性チェック・ソフトウェアを使用して診断データを取得する詳細なステップを示している。拡張ファイル整合性チェックは、fsckの変形とすることができる。fsckは、ファイル・システム内の各データ単位をトラバースする（ステップ 6 2 1）。fsckを実行するプロセッサは、データ単位がファイル・システム・メタデータを含むかどうかを判定することができる（ステップ 6 2 3）。「含む」とは、データ単位が、ファイル・システムを構成するすべてのメタデータの少なくとも一部を含むことを意味する。fsckは、データ単位がファイル・システム・メタデータであるとの判定に回答して、メタデータの単位を書き込む（ステップ 6 2 5）。さもなければ、ステップ 6 2 3 および 6 2 5 は、ステップ 6 2 7 に続く。fsckは、データ単位がユーザ・データを含むかどうかを判定する（ステップ 6 2 7）。fsckは、データ単位がユーザ・データであるとの判定に回答して、スパス・オブジェクトを診断ファイルに書き込む（ステップ 6 2 9）。スパス・オブジェクトの書き込みは、プロセッサがメタデータ内のヌル・ポインタを書き込むことを含むことができる。ヌル・ポインタは、例えば、図 5 のブロック・ポインタ・テーブル 4 1 4 のヌル 4 1 5 とすることができる。fsckは、トラバースすべきさらなるデー

40

50

タ単位が残っているかどうかを判定する（ステップ631）。ステップ625および629の各書き込みは、図3の診断ファイル305へのものとすることができ、その場合、診断ファイルは、診断ファイル・システムに関連することができる。スパーズ・オブジェクトが非スパーズ・オブジェクトと異なるのは、ブロックが空であるかまたは有効でないデータで満たされていることを注記することを除き、ディスクへの書き込みがデータを変化させない点である。ブロック全体を占有する代わりに、fsckは、メタデータに関連するすべてのスパーズ・オブジェクトのリストを単に作成するにすぎない。代替として、fsckは、lseekシステム・コールを単に使用して、書き込みのためにオープンされたファイルのファイル位置を進める。lseekシステム・コールは、Unixライクのオペレーティング・システムで利用可能なファイル処理システム・コールである。基本的に、lseekコマンドは、現在のファイル末尾を越えてファイル位置を動かすように、オペレーティング・システムに命令する。lseekシステム・コールがブロック全体にわたってファイル位置を進めた場合、オペレーティング・システムは、ディスク・ドライブ上のストレージの物理ブロックを割り当てることなく、ブロックがスパーズ・ブロックまたはスパーズ・オブジェクトであることを記録する。スパーズ・オブジェクトは、ストレージの仮想ブロックであるが、物理的ストレージの完全なブロックのスペースを占有しない。むしろ、メタデータは、とりわけ、ファイルに関連するそのようなブロックの番号を識別する。それにも関わらず、プロセッサは、ファイルによって占有される全スペースがディレクトリ・リスティング中に決定されるとき、そのようなスパーズ・オブジェクトを計算に追加する。

【0041】

図9は、本発明の一実施形態による、AIXオペレーティング・システムのコマンド・ラインで入力されるコマンドの一例を示している。山括弧<>で括られたテキストは、ユーザ指定の文字列を示す。例えば、<新規ファイル・システム・サイズ> 711は、新規ファイル・システムに割り当てられる512バイトのブロックの個数を指定する数とすることができる。

【0042】

第1のステートメントは、crfsコマンド701を使用して、ファイル・システムを生成する。フラグは以下の通りである。

【0043】

-v jfs2 703は、仮想ファイル・システムのタイプを指定する。

【0044】

-g rootvg 705は、ファイル・システムを作成すべき既存のボリューム・グループを指定する。ボリューム・グループは、1つ以上の物理的ボリュームの集まりである。

【0045】

-m /newfs 707は、ファイル・システムが利用可能にされるディレクトリであるマウント・ポイントを指定する。

【0046】

-a size=<新規ファイル・システム・サイズ> 709は、512バイトのブロック数でJFS2のサイズを指定し、その際、指定されたサイズが物理的パーティション・サイズで均等に割り切れない場合、プロセッサは、均等に割り切れる最も近い数に切り上げる。

【0047】

診断ファイル・システムは、図3の診断ファイル305に関連することができる。

【0048】

mount /newfs 721は、ファイル・システムをマウントするようプロセッサに命令する。

【0049】

次がファイル整合性チェックであり、以下のフラグを含むfsck 722コマンドを使用する。

【0050】

-n 723は、非対話モードで、すなわち、修復についてユーザに問い合せず、指定さ

10

20

30

40

50

れたファイル・システムを変更せずにファイル整合性チェックを実行するようプロセッサに命令する。

【 0 0 5 1 】

-o metacapture=/newfs/out1 7 2 5 は、メタデータおよびスパス・オブジェクトのコピーをスパス・オブジェクトを有するファイルに書き込むようプロセッサに命令するパラメータを識別し、その場合、ファイルは、/newfs/out1またはユーザ指定のディレクトリに置かれる。

【 0 0 5 2 】

<ファイル・システム・マウント・ポイント> 7 3 5 は、プロセッサが診断を実行するファイル・システムの識別情報である。

10

【 0 0 5 3 】

プロセッサは、残りの2つのコマンド・ラインに基づいて統合を行う。backup 7 3 9 は、fsck 7 2 2 で生成されたすべてのファイルを単一のファイルに配置するようプロセッサに命令する。すなわち、診断ファイル・システムを単一のファイルに配置する。backup 7 3 9 は、以下のフラグを含む。

【 0 0 5 4 】

-f <目的ファイル名> 7 4 1 は、出力を保存する1つ以上の装置を識別する。

【 0 0 5 5 】

-0 /newfs 7 4 3 は、バックアップのためのソース・ファイル・システムを指定する。

【 0 0 5 6 】

20

compress 7 5 1 は、スパス・オブジェクトを圧縮トークンで置き換えるステップである。compress 7 5 1 がとる唯一のパラメータは、稠密診断データが保存されるファイルである<目的ファイル名>である。

【 0 0 5 7 】

統合は、ファイル整合性チェックによって提供された個々のファイルを使い、各ファイルを圧縮することによってもなされ得ることが理解される。その場合、各ファイルが稠密診断データになり、ファイルの集まりに追加される各後続ファイルが、稠密診断データに加わる。

【 0 0 5 8 】

したがって、本発明の側面は、ファイル・システム不整合を修復するために任命された専門家からさえもユーザ・データが保護された方法で、ファイル・システム不整合に関する詳細を取得するためのコンピュータ実施方法を提供する。加えて、本発明の側面は、ファイル・システム開発者が、ファイル・システム不整合に寄与するかもしれない開発中のソフトウェアにおいて問題を診断することを可能にする。さらに、スパス・オブジェクトは、スパス・オブジェクトを論理ボリュームに抽出できる、またはスパス・オブジェクトを直接見てファイル・システム問題をデバッグできるサービス・チームに送られることができる。

30

【 0 0 5 9 】

本発明の実施形態は、完全なハードウェア実施形態、完全なソフトウェア実施形態、またはハードウェア要素とソフトウェア要素の両方を含む実施形態の形を取ることができる。好ましい実施形態では、本発明は、ファームウェア、常駐ソフトウェア、マイクロコードなどの、ただしこれらに限定されないソフトウェアで実施される。

40

【 0 0 6 0 】

さらに、本発明は、コンピュータもしくは任意の命令実行システムによってまたはそれらに関連して使用するためのプログラム・コードを提供するコンピュータ使用可能またはコンピュータ可読媒体等から取得可能なコンピュータ・プログラムの形を取ることができる。この説明では、コンピュータ使用可能またはコンピュータ可読媒体等は、命令実行システム、機器、または装置によってまたはそれらに関連して使用するためのプログラムを格納、保存、伝達、伝播、または転送することができる任意の機器とすることができる。

【 0 0 6 1 】

50

媒体は、電子、磁気、光学、電磁気、赤外線、または半導体システム（または機器もしくは装置）、あるいは伝播媒体とすることができる。コンピュータ可読媒体の例は、半導体またはソリッド・ステート・メモリ、磁気テープ、着脱可能コンピュータ・ディスクレット、ランダム・アクセス・メモリ（RAM）、読み取り専用メモリ（ROM）、固定磁気ディスク、および光ディスクを含む。光ディスクの現行例は、コンパクト・ディスク・読み取り専用メモリ（CD-ROM）、コンパクト・ディスク・リード/ライト（CD-R/W）、およびDVDを含む。

【0062】

プログラム・コードを保存または実行し、あるいはその両方を行うのに適したデータ処理システムは、システム・バスを介してメモリ要素に直接的または間接的に結合される少なくとも1つのプロセッサを含む。メモリ要素は、プログラム・コードを実際に行うときに利用されるローカル・メモリ、大容量ストレージ、実行中に大容量ストレージからコードを取り出さなければならない回数を減らすためにプログラム・コードの少なくとも一部を一時的に記憶するキャッシュ・メモリを含むことができる。

10

【0063】

（キーボード、ディスプレイ、ポインティング・デバイスなどを含むが、これらに限定されない）入力/出力またはI/O装置は、直接的にまたは仲介I/Oコントローラを介してシステムに結合されることができる。

【0064】

データ処理システムが、その他のデータ処理システムまたは遠隔プリンタもしくはストレージに仲介私設または公衆ネットワークを介して結合されることができるよう、ネットワーク・アダプタもシステムに結合されることができる。モデム、ケーブル・モデム、およびイーサネット（登録商標）・カードは、現在利用可能なタイプのネットワーク・アダプタのうちのほんの2、3の実例である。

20

【0065】

本発明の説明は、例示および説明の目的で提示されたもので、網羅的であること、または本発明を開示された形態に限定することは意図されていない。当業者には、多くの変更および変形が明らかであろう。本発明の例示的な実施形態は、本発明の原理および実際の応用を最も良く説明し、企図される特定の使用に適した様々な変更を有する様々な実施形態について当業者が本発明を理解できるように選択され、説明された。

30

【図面の簡単な説明】

【0066】

【図1】本発明の実施形態が実施され得るデータ処理システムのブロック図である。

【図2】本発明の一実施形態に従って、顧客システムがどのように専門家システムと対話することができるかを示した図である。

【図3】本発明の一実施形態による、ソース・ファイル・システムから診断ファイルへの変換を仲介するオペレーティング・システムの動作を示した図である。

【図4】本発明の一実施形態による、従来のファイル・フォーマットの一例を示した図である。

【図5】本発明の一実施形態による、スパース・ファイル・フォーマットの一例を示した図である。

40

【図6】本発明の一実施形態による、ファイル・システムに関する診断データを取得する最初のステップを示した図である。

【図7】本発明の一実施形態による、ファイル・システムに関する診断データをパッケージ化して送信するステップを示した図である。

【図8】本発明の一実施形態による、拡張ファイル整合性チェック・ソフトウェアを使用して診断データを取得する詳細なステップを示した図である。

【図9】本発明の一実施形態による、AIXオペレーティング・システムのコマンド・ラインで入力されるコマンドの一例を示した図である。

【符号の説明】

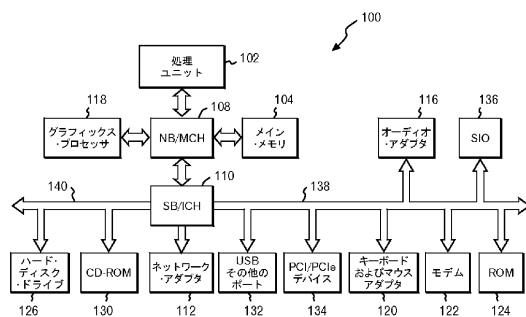
50

【 0 0 6 7 】

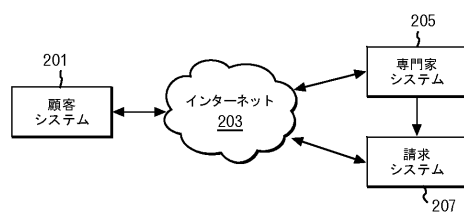
- 2 0 1 顧客システム
- 2 0 3 インターネット
- 2 0 5 専門家システム
- 2 0 7 請求システム
- 3 0 1 オペレーティング・システム
- 3 0 3 ソース・ファイル・システム
- 3 0 5 診断ファイル
- 3 1 1 コマンド・ライン・エディタ
- 4 1 1 ファイル
- 4 1 3 「 0 」ブロック
- 4 1 4 ブロック・ポインタ・テーブル
- 4 1 5 ヌル・ポインタ
- 4 1 6 ヌル・ポインタ
- 4 1 7 ヌル・ポインタ
- 4 1 8 ヌル・ポインタ

10

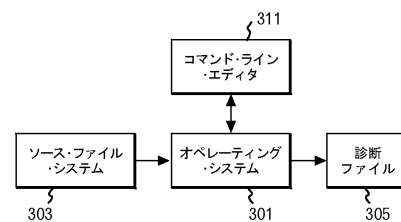
【 図 1 】



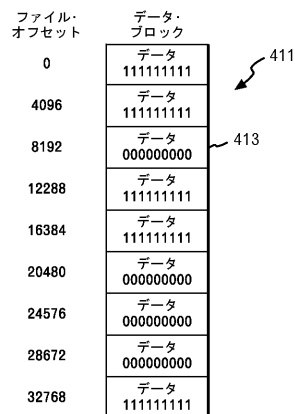
【 図 2 】



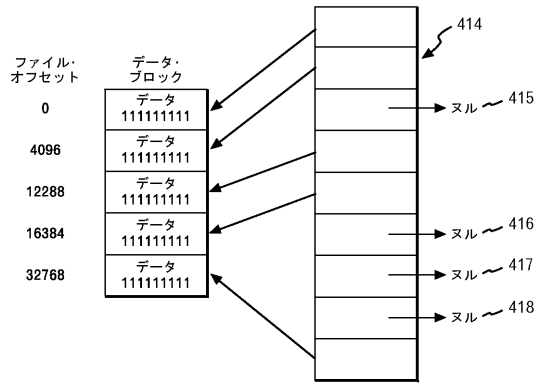
【 図 3 】



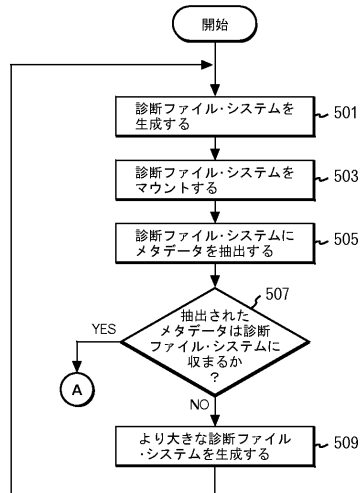
【 図 4 】



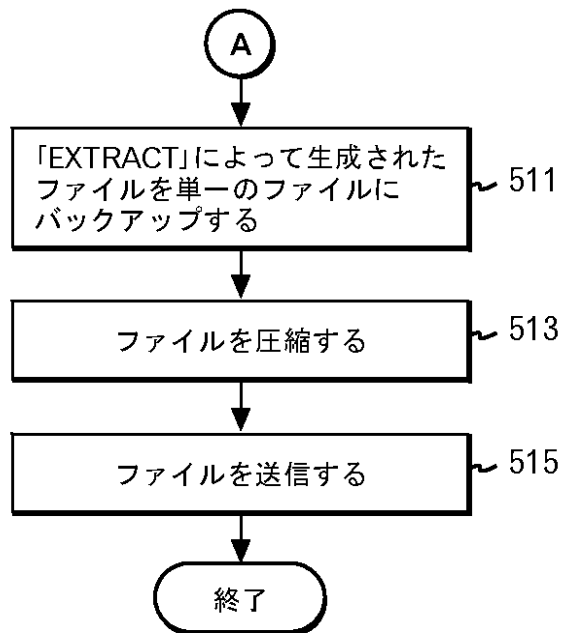
【図 5】



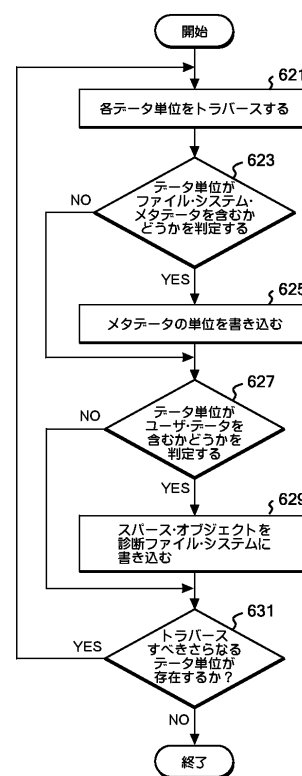
【図 6】



【図 7】



【図 8】



【図 9】

701
\$>crfs -v jfs2 -g rootvg -m /newfs -a size=<新規ファイル・システム・サイズ>
703 705 707 709
711

721
\$>mount /newfs

722
\$>fsck -n -o metacapture=/newfs/out1 <ファイル・システム・マウント・ポイント>
723 725 735

739
\$>backup -f <目的ファイル名> -o /newfs
741 743

751
\$>compress <目的ファイル名>

フロントページの続き

(74)代理人 100086243

弁理士 坂口 博

(72)発明者 ジャネット・エリザベス・アドキンス

アメリカ合衆国 7 8 7 4 6 テキサス州オースチン セッジフィールド・ドライブ 5 8 0 2

(72)発明者 マーク・アレン・グラブズ

アメリカ合衆国 7 8 6 8 1 テキサス州ラウンド・ロック ボブキャット・ドライブ 8 6 1 0

審査官 桜井 茂行

(56)参考文献 特開 2 0 0 4 - 0 3 8 9 2 9 (J P , A)

米国特許出願公開第 2 0 0 3 / 0 1 8 2 3 0 1 (U S , A 1)

(58)調査した分野(Int.Cl. , D B 名)

G 0 6 F 1 2 / 0 0

J S T P l u s (J D r e a m I I)