

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号  
特許第4988214号  
(P4988214)

(45) 発行日 平成24年8月1日 (2012. 8. 1)

(24) 登録日 平成24年5月11日 (2012. 5. 11)

(51) Int. Cl.

F I

GO 6 F 3/06 (2006. 01)

GO 6 F 12/08 (2006. 01)

GO 6 F 3/06 3 O 4 U

GO 6 F 12/08 5 5 7

GO 6 F 12/08 5 4 3 Z

GO 6 F 12/08 5 5 1 J

請求項の数 2 (全 14 頁)

(21) 出願番号	特願2006-17056 (P2006-17056)	(73) 特許権者	390009531
(22) 出願日	平成18年1月26日 (2006. 1. 26)		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2006-221623 (P2006-221623A)		I N T E R N A T I O N A L B U S I N E S S M A S C H I N E S C O R P O R A T I O N
(43) 公開日	平成18年8月24日 (2006. 8. 24)		アメリカ合衆国 1 0 5 0 4 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20年10月24日 (2008. 10. 24)		
(31) 優先権主張番号	11/053226	(74) 代理人	100108501
(32) 優先日	平成17年2月7日 (2005. 2. 7)		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
前置審査		(74) 代理人	100091568
			弁理士 市位 嘉宏
			最終頁に続く

(54) 【発明の名称】 記憶装置における喪失した書き込みの検出および回復

(57) 【特許請求の範囲】

【請求項 1】

記憶媒体にデータを書き込むための要求を受け取るステップと、  
前記記憶媒体に書き込むように要求された前記データを、キャッシュに格納するステップと、  
前記記憶媒体への前記データの書き込みを開始するステップと、  
前記キャッシュに格納された前記データが前記記憶媒体に書き込まれた前記データと同じであるかどうかを定期的に判定するステップと、  
前記キャッシュ内のエントリ数がしきい値を超えるかどうかを判定するステップを有し、  
前記エントリはキャッシュに入れられた書き込みに対応しており、  
合体された非重複書き込みアドレスのリストを生成するために、前記エントリ内の書き込みアドレスを合体するステップと、  
前記合体された非重複書き込みアドレスのリストを順序付けするステップをさらに有し、  
低順位の合体された非重複書き込みアドレスより前に、高順位の合体された非重複書き込みアドレスが、前記記憶媒体に対する書き込みエラーについて検証される、方法。

【請求項 2】

前記しきい値を超えることは、前記キャッシュが満杯の 1 0 パーセントを上回ることを示す、請求項 1 に記載の方法。

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

本発明は、記憶装置における喪失した書き込みの検出および回復のための方法に係る。

## 【背景技術】

## 【0002】

ディスク・ドライブにデータを書き込む書き込み動作は、断続的または永続的に障害が発生する可能性がある。こうした障害を検出するために、ある種のドライブは、書き込み電流が一定のしきい値よりも低い場合に検出する、ヘッド読み取りおよび書き込み前置増幅 (preamplification) 回路を使用する。この検出回路に加えられた複雑性が、信頼性の問題を提示する可能性がある。さらに、検出しきい値の設定が、すべての書き込みエラーの検出を保証するものでない場合がある。

10

## 【0003】

ある種の実施では、アクチュエータを定期的にディスクの予約領域へ移動させ、その後書き込んで、各ヘッドによって書き込まれた内容を検証することによる、ディスク・ドライブの書き込み喪失の問題がある。このすべてのヘッドに関する書き込み能力の検証は、永続問題用自己テスト (persistent problem self test: P P S T) と呼ばれることがある。このメカニズムは、書き込み問題が永続的である場合、すなわち、書き込みを喪失しているドライブが不良ヘッドに関する後続のすべての書き込みの喪失を継続する場合にのみ、喪失した書き込みを検出する。P P S T 検証メカニズムは、断続的に喪失した書き込みの検出には十分に有効ではない。加えて P P S T 検証は、連続する P P S T 書き込み能力検証の間に喪失した書き込みが原因で書き込まれなかったデータを回復させることはできない。さらに、データ破損の量を最小限にするために P P S T 検証の頻度が上がった場合、入力/出力 (I/O) 性能は許容不可能なレベルまで低下する可能性がある。

20

## 【0004】

ある種の実施では、連続する P P S T 検証間にキャッシュ内のすべての書き込みを維持することによって、P P S T 検証を増補することができる。誤った書き込みが検出された場合、データはキャッシュから直接回復される。ヘッドの定期的チェックでは断続的に喪失する書き込みが検出されない場合があるため、このメカニズムでは断続的に喪失する書き込みは保護しない。さらに、P P S T 検証のオーバーヘッドによる性能の低下を緩和するために、かなり大型で潜在的に高価な専用キャッシュが必要となる可能性がある。すべての書き込みを取り込むために必要なキャッシュは、P P S T 検証間の時間間隔の増加と共に増大する。

30

## 【0005】

断続的および永続的の両方の喪失した書き込みを保護するある種の実施では、それぞれの書き込み動作に対して書き込み検証を実行することが可能であり、書き込みが実行されるごとにディスク・ドライブが回転を完了し、書き込まれたばかりのデータを読み取って、その書き込まれたばかりのデータと書き込みバッファ内のデータとを比較する。これによってデータが損失しないことは保証されるが、ドライブの待ち時間が追加され、結果として生じる I/O 性能が受け入れ不可能なものとなる可能性がある。

40

## 【発明の開示】

## 【発明が解決しようとする課題】

## 【0006】

本発明の目的は、記憶媒体における喪失した書き込みの検出および回復のための方法を提供することである。

## 【課題を解決するための手段】

## 【0007】

記憶媒体に書き込むように要求されたデータがキャッシュに格納される。記憶媒体へのデータの書き込みが開始される。キャッシュに格納されたデータと記憶媒体に書き込まれたデータとが同じであるかどうか、定期的に判別される。

50

## 【 0 0 0 8 】

追加の実施形態では、キャッシュに格納されたデータと記憶媒体に書き込まれたデータとが同じであるという判定に応答して、格納されたデータがキャッシュから除去される。

## 【 0 0 0 9 】

さらに追加の実施形態では、キャッシュに格納されたデータと記憶媒体に書き込まれたデータとが同じでないという判定に応答して、記憶媒体は書き込み保護される。記憶媒体に対して喪失された少なくとも1つの書き込みを示すエラーが生成される。

## 【 0 0 1 0 】

他の実施形態では、記憶媒体はディスクであり、キャッシュはディスクよりも記憶容量が小さく、ディスクに比べてキャッシュにデータを書き込む方が信頼性が高く、ホスト・アプリケーションは記憶媒体からではなくキャッシュから未検証のデータを読み取る。

10

## 【 0 0 1 1 】

さらに他の実施形態では、キャッシュ内のエントリ数がしきい値を超えたかどうか判别され、エントリはキャッシュに入れられた書き込みに対応する。エントリ内の書き込みアドレスが合体され、合体済み書き込みアドレスのリストを生成する。合体済み書き込みアドレスのリストが順序付けされ、低順位の合体済み書き込みアドレスより前に、記憶媒体への書き込みエラーについて高順位の合体済み書き込みアドレスが検証される。ある種の実施形態では、しきい値を超えることはキャッシュが満杯の10パーセントを上回ることを示す。

## 【 0 0 1 2 】

20

追加の実施形態では、キャッシュ内のエントリ数がしきい値を超えたかどうか判别され、エントリは書き込みアドレスに対応する。選択されたエントリに対応するデータが、エラーなしで記憶媒体に書き込まれたかどうか判别される。ある種の実施形態では、しきい値を超えることはキャッシュが満杯の90パーセントを上回ることを示す。他の実施形態では、検証は、総合検証用自己テストが実行可能であるという判定に응答するものであり、総合検証用自己テストが実行可能でない場合、永続問題用自己テストが実行される。

## 【 0 0 1 3 】

ある種の追加の実施形態では、記憶媒体への書き込みはある1つの順序で実行され、記憶媒体への書き込みの検証は別の順序で実行される。

30

## 【 0 0 1 4 】

さらに追加の実施形態では、キャッシュに格納されたデータが記憶媒体に書き込まれたデータと同じでない場合、断続的書き込みエラーが検証される。

## 【 0 0 1 5 】

次に図面を参照するが、全体を通じて同じ参照番号は対応する部分を表す。

## 【 発明を実施するための最良の形態 】

## 【 0 0 1 6 】

以下の説明では、本明細書の一部を形成し、いくつかの実施形態を示す、添付の図面を参照する。他の実施形態が使用可能であり、構造および動作の変更が可能であることを理解されよう。

40

## 【 0 0 1 7 】

ある種の実施形態では、ディスク・ドライブが全ブロック書き込みを断続的または永続的に喪失した場合、および書き込みエラーの報告に障害が生じた場合に、ディスク・ドライブからデータを回復することができる。

## 【 0 0 1 8 】

図1は、ある種の実施形態に従ったコンピューティング環境100を示すブロック図である。ディスク・ドライブ102などの記憶装置が、ホスト・コンピュータ・デバイスなどのコンピュータ・デバイス104に結合される。ディスク・ドライブ102は、直接、またはストレージ・エリア・ネットワーク(SAN)、ローカル・エリア・ネットワーク(LAN)、イントラネット、インターネットなどの、ネットワークを介して、コンピュ

50

ータ・デバイスに結合することができる。

【0019】

コンピュータ・デバイス104は、クライアント、ストレージ・サーバ、サーバ、パーソナル・コンピュータ、ワークステーション、メインフレーム、ミッドレンジ・コンピュータ、ネットワーク・アプライアンス、パーム・トップ・コンピュータ、テレフォニ・デバイス、ブレード・コンピュータ、ハンド・ヘルド・コンピュータなどの、当分野で現在知られているものを含む、任意の好適なコンピュータ・デバイスとすることができる。図1は、ディスク・ドライブ102を記憶装置として示しているが、代替の実施形態では、記憶装置は当分野で知られた任意の他の好適なデバイスとすることができる。例示的なディスク・ドライブ102は、ある種の実施形態では、新磁気ディスク制御機構(RAID: Redundant Array of Independent Disk)アレイに含めるか、または単純ディスク束(JBOD: Just a Bunch of Disks)に含めることができる。

10

【0020】

ディスク・ドライブ102は、未検証の書き込みデータ106aを格納するキャッシュ106と、キャッシュ106に格納されたエントリに対応する非重複論理ブロック・アドレスおよび関連するカウントの順序付けリスト110を維持するリスト維持管理(maintainer)アプリケーション108と、所定の条件114のセットのうちの1つまたは複数が満たされた場合に実行される総合検証用自己テスト(CVST)アプリケーション112と、PPSTアプリケーション116とを含むことができる。ディスク・ドライブ102は、データを書き込むことが可能な1つまたは複数のディスク118を備えることができる。

20

【0021】

キャッシュ106は、任意の好適な不揮発性メモリを含むことができる。ある種の実施形態では、キャッシュへの書き込みはディスク118への書き込みよりもかなり高い信頼度で実行される。キャッシュ106に格納された未検証書き込みデータ106aは、たとえばディスク・ドライブ102がリセットされた場合であっても保持される。

【0022】

リスト維持管理アプリケーション108、CVSTアプリケーション116、およびPPSTアプリケーション116は、ソフトウェア、ファームウェア、ハードウェア、またはそれらの任意の組み合わせで実施可能である。順序付けリスト110および所定の条件114は、任意の好適なデータ構造に格納することができる。

30

【0023】

ある種の実施形態では、ディスク・ドライブ102でコンピュータ・デバイス104から書き込み要求が受け取られた場合、書き込み要求に対応するデータはキャッシュ106に格納される。ディスク118へのデータの書き込み時に何らかのエラーが生じた場合、CVSTアプリケーション112あるいはPPSTアプリケーション116またはその両方は、ディスク118に書き込まれたデータとキャッシュ106に格納されたデータとを比較することによって、エラーを判別することができる。

【0024】

図2は、ある種の実施形態に従った、コンピューティング環境100に含まれるコンポーネントおよびデータ構造を示すブロック図である。

40

【0025】

リスト維持管理アプリケーション108は順序付けリスト110を維持し、順序付けリスト110は非重複論理ブロック・アドレス(LBA)および関連するカウントを備えたエントリを含む。たとえば、順序付けリスト110の例示的エントリは、カウント=4で論理ブロック・アドレス=3への書き込みに対応することができる。これは、データが4つの論理ブロック・アドレス3、4、5、および6に書き込まれることを暗黙に示す。論理ブロック・アドレス3、4、5、および6はすでに例示的エントリに含まれているため、順序付けリスト110には論理ブロック・アドレス3、4、5、または6を含む他のエントリはない。リスト維持管理アプリケーション108は、ディスク・ドライブ102か

50

らの読み取り性能を向上させるために、リスト 1 1 0 も順序付けする。

【 0 0 2 6 】

ある種の所定の条件が満たされた場合、C V S T アプリケーション 1 1 2 は順序付けリスト 1 1 0 のサブセット 2 0 0 を選択し、対応する読み取りコマンド 2 0 2 を発行して、ディスク 1 1 8 からデータを読み取ることができる。ある種の実施形態では、選択されたサブセット 2 0 0 は、順序付けリスト 1 1 0 内のエントリの上位 5 % を含むことができる。たとえば、サブセット 2 0 0 内の選択されたエントリが論理ブロック・アドレス = 3 およびカウント = 4 に対応する場合、C V S T アプリケーション 1 1 2 は、ディスク 1 1 8 から論理ブロック・アドレス 3、4、5、6 を読み取るために読み取りコマンド 2 0 2 を発行することができる。

10

【 0 0 2 7 】

ディスク 1 1 8 内に格納されたデータ 2 0 4 とキャッシュ 1 0 6 に格納されたデータとが比較され（参照数 2 0 6）、データ 2 0 4 は実行された読み取りコマンド 2 0 2 の結果である。比較の結果を示すステータス 2 0 8 を、リスト維持管理アプリケーション 1 0 8 に送信するために生成することができる。たとえばステータス 2 0 8 が、ディスク 1 1 8 内のデータ 2 0 4 がキャッシュ 1 0 8 内のデータと一致することを示す場合、リスト維持管理アプリケーション 1 0 8 は順序付けリスト 1 1 0 内の対応するエントリを削除することができる。というのは、この一致は、削除されたエントリに対応する未検証書き込みデータ 1 0 6 a がディスク 1 1 8 に正しく書き込まれることが検証されたことを示すからである。

20

【 0 0 2 8 】

したがって図 2 は、リスト維持管理アプリケーション 1 0 8 および C V S T アプリケーション 1 1 2 が、ディスク 1 1 8 に書き込まれたデータとキャッシュ 1 0 6 に書き込まれたデータとを比較することによってディスク 1 1 8 に書き込まれたデータを検証する、ある種の実施形態を示す。

【 0 0 2 9 】

図 3 は、ある種の実施形態に従った、ディスク・ドライブ 1 0 2 で実施される喪失した書き込みから回復するための動作を示す図である。

【 0 0 3 0 】

制御はブロック 3 0 0 で開始され、ディスク・ドライブ 1 0 2 のディスク 1 1 8 にデータを書き込むための要求が受け取られる。ディスク・ドライブ 1 0 2 は（ブロック 3 0 2 で）書き込み要求に対応するデータをキャッシュ 1 0 6 に格納し、データはディスク 1 1 8 に書き込まれる。制御は、ブロック 3 0 2 からブロック 3 0 4 および 3 0 8 へ並行して進む。

30

【 0 0 3 1 】

ディスク・ドライブ 1 0 2 は（ブロック 3 0 4 で）、キャッシュ 1 0 6 がかなりの数のエントリを有するかどうかを判別する。ある種の実施形態では、キャッシュ 1 0 6 は、キャッシュが満杯の 1 0 パーセントを上回る場合に、かなりの数のエントリを有するものとみなすことができる。これを上回る場合ディスク・ドライブ 1 0 2 は（ブロック 3 0 6 で）、リスト維持管理アプリケーション 1 0 8 を実行してキャッシュ・エントリを合体する。リスト維持管理アプリケーション 1 0 8 は、非重複論理ブロック・アドレスおよび関連するカウントの順序付けリスト 1 1 0 を生成することができる。その後、制御はブロック 3 0 0 に戻る。キャッシュ 1 0 6 がかなりの数のエントリを有していない場合、ディスク・ドライブ 1 0 2 は制御をブロック 3 0 0 に戻し、ここで書き込み要求が受け取られる。

40

【 0 0 3 2 】

ディスク・ドライブ 1 0 2 は（ブロック 3 0 8 で）、キャッシュ 1 0 6 がほぼ満杯であるかどうかを判別する。たとえばある種の実施形態では、キャッシュが満杯の 9 0 % を超えている場合、キャッシュ 1 0 6 はほぼ満杯であるものとみなされる。ディスク・ドライブ 1 0 2 は（ブロック 3 1 0 で）、C V S T 処理が実行可能であるかどうかを判別する。実行可能である場合、ディスク・ドライブ 1 0 2 は（ブロック 3 1 2 で）、順序付けリス

50

ト 1 1 0 内の選択されたエントリ上で C V S T アプリケーション 1 1 2 を実行する。ある種の実施形態では、選択されたエントリは、順序付けリスト 1 1 0 内で最高順位を有するエントリ、たとえば順序付けリスト 1 1 0 内で上位 5 % のエントリとすることができる。新しい書き込み要求が受け取られると、制御はブロック 3 0 0 に戻る。ディスク・ドライブ 1 0 2 が ( ブロック 3 0 8 で )、キャッシュ 1 0 6 がほぼ満杯ではないと判別した場合も、制御はブロック 3 0 0 に戻る。

#### 【 0 0 3 3 】

ディスク・ドライブ 1 0 2 が ( ブロック 3 1 0 で )、C V S T 処理が実行不能であると判別した場合、ディスク・ドライブ 1 0 2 は ( ブロック 3 1 4 で ) 即時に P P S T アプリケーション 1 1 6 を実行し、制御はブロック 3 0 0 に戻る。P P S T アプリケーション 1 1 6 は、永続的に喪失した書き込み問題を検出することができる。ある種の実施形態では、書き込み問題が永続的である場合にのみ、喪失した書き込みの検出がある種の環境で適切であることから、C V S T 処理は実行不能である。C V S T 処理が実行不能なこうした実施形態では、C V S T 処理のある種のオーバーヘッドを招く可能性がある。

#### 【 0 0 3 4 】

ある種の実施形態では、キャッシュがほぼ満杯であるかどうかの判別 ( ブロック 3 0 8 ) あるいはキャッシュがかなりの数のエントリを有するかどうかの判別 ( ブロック 3 0 4 )、またはその両方は、ある種の実施形態でキャッシュの満杯の程度またはキャッシュ内のエントリ数を示すことが可能な、所定のしきい値との比較に基づいて実施することができる。

#### 【 0 0 3 5 】

したがって図 3 は、キャッシュがほぼ満杯の場合に C V S T アプリケーション 1 1 2 が実行される、ある種の実施形態を示す。C V S T アプリケーション 1 1 2 は、ディスク 1 1 8 への書き込みの検証のために順序付けリスト 1 1 0 からエントリを選択し、順序付けリスト 1 1 0 は、キャッシュ 1 0 6 内のエントリ数が事前に定義されたしきい値を超えた場合に、リスト維持管理アプリケーション 1 0 8 によってアクティブに管理可能である。

#### 【 0 0 3 6 】

図 4 は、ある種の実施形態に従った、リスト維持管理アプリケーション 1 0 8 で実施されるある種の動作を示す図である。

#### 【 0 0 3 7 】

制御はブロック 4 0 0 で開始され、リスト維持管理アプリケーション 1 0 8 の実行が開始される。制御は、ブロック 4 0 0 からブロック 4 0 2 および 4 0 6 へ並行して進む。

#### 【 0 0 3 8 】

リスト維持管理アプリケーション 1 0 8 は ( ブロック 4 0 2 で )、連続するかまたは重複するキャッシュに入れられた書き込みの論理ブロック・アドレスおよびカウントを含む、書き込みアドレスを合体する。その結果、リスト維持管理アプリケーション 1 0 8 は、非重複論理ブロック・アドレスおよび関連するカウントのリスト 1 1 0 を作成する。たとえば、カウント = 4 で L B A = 3 への書き込み、およびその後のカウント = 8 で L B A = 5 への書き込みは、データ・ブロックのうちの 2 つが重複することから、カウント = 1 0 で L B A = 3 の単一のアドレスに合体することができる。これにより、論理ブロック・アドレスおよび関連するカウントを含む、連続するかまたは重複する 2 つまたはそれ以上の書き込みの検証が、単一の読み取りに効果的に削減される。代替の実施形態では、2 つまたはそれ以上のほぼ連続する書き込みも合体することができる。ブロック 4 0 2 の結果、重複なしの論理ブロック・アドレスおよびカウントの削減されたリストとなる。

#### 【 0 0 3 9 】

リスト維持管理アプリケーション 1 0 8 は ( ブロック 4 0 4 で )、その後のディスク・ドライブ 1 0 2 からの読み取り性能を向上させるために順序付けリスト 1 1 0 を再度順序付けする。たとえばある種の実施形態では、順序付けリスト 1 1 0 は、同じランクにある関連するカウント・サイズおよび論理ブロック・アドレス数に基づいて論理ブロック・アドレスをランク付けすることができる。

## 【 0 0 4 0 】

ブロック 4 0 6 でリスト維持管理アプリケーション 1 0 8 は、C V S T アプリケーションからステータス 2 0 8 の通知を受け取る。受け取ったステータス 2 0 8 の通知に基づき、リスト維持管理アプリケーション 1 0 8 は（ブロック 4 0 8 で）、新しいホスト書き込みの論理ブロック・アドレスおよびカウントを順序付けリスト 1 1 0 に追加することができる。リスト維持管理アプリケーション 1 0 8 は、キャッシュ内に格納された論理ブロック・アドレスに対応するデータが、ディスク 1 1 8 に書き込まれたデータと首尾よく一致した、C V S T アプリケーション 1 1 2 によって実行された論理ブロック・アドレスのリストを削除することもできる。

## 【 0 0 4 1 】

10

したがって図 4 は、ディスク・ドライブ 1 0 2 からのその後の読み取り性能を向上させるために、リスト維持管理アプリケーション 1 0 8 が非重複論理ブロック・アドレスおよび関連するカウントの順序付けリストを維持する、ある種の実施形態を示す。

## 【 0 0 4 2 】

図 5 は、ある種の実施形態に従った、C V S T アプリケーション 1 1 2 で実施される動作を示す図である。

## 【 0 0 4 3 】

制御はブロック 5 0 0 で開始され、C V S T アプリケーション 1 1 2 が開始される。C V S T アプリケーション 1 1 2 は（ブロック 5 0 2 で）、所定の条件 1 1 4 が満たされたかどうかを判別する。所定の条件は、C V S T アプリケーション 1 1 2 によって順序付け

20

リスト 1 1 0 から選択されるエントリのサイズを決定することができる。

## 【 0 0 4 4 】

C V S T アプリケーション 1 1 2 が、所定の条件 1 1 4 が満たされていると判定した場合、C V S T アプリケーション 1 1 2 は（ブロック 5 0 4 で）、リスト維持管理アプリケーション 1 0 8 によって提供された順序付けリスト 1 1 0 のサブセット 2 0 0 を読み取りコマンド 2 0 2 に変換する。そうでない場合、C V S T アプリケーション 1 1 2 は（ブロック 5 0 2 で）、所定の条件が満たされたかどうかの判別を続ける。

## 【 0 0 4 5 】

C V S T アプリケーション 1 1 2 は（ブロック 5 0 6 で）、実行されたディスク 1 1 8 からの読み取りと、キャッシュ 1 0 6 内の未検証書き込みデータ 1 0 6 a を示す対応する

30

キャッシュ・エントリとを比較する。

## 【 0 0 4 6 】

C V S T アプリケーション 1 1 2 はブロック 5 0 8 で、キャッシュ 1 0 8 内のデータがディスク 1 1 8 に書き込まれたデータと一致するかどうかを判別する。一致する場合、C V S T アプリケーション 1 1 2 は（ブロック 5 1 0 で）、キャッシュ 1 0 6 からの対応するデータを削除し、（ブロック 5 1 2 で）ステータス 2 0 8 の通知をリスト維持管理アプリケーション 1 0 8 に送信し、リスト維持管理アプリケーション 1 0 8 は順序付けリスト 1 1 0 の修正および再順序付けを管理する。

## 【 0 0 4 7 】

C V S T アプリケーション 1 1 2 がブロック 5 0 8 で、実行された読み取りについてキャッシュ 1 0 8 内のデータがディスク 1 1 8 に書き込まれたデータと一致しないと判定した場合、C V S T アプリケーション 1 1 2 は（ブロック 5 1 4 で）、ディスク・ドライブ 1 0 2 を書き込み保護モードに設定し、エラーを報告する。ディスク・ドライブ 1 0 2 が書き込み保護モードの場合、ディスク 1 1 8 にデータを書き込むことはできない。ディスク 1 1 8 への書き込みエラーは、ディスク・ドライブ 1 0 2 が書き込み保護モードにされた後に判定することができる。正しいデータは、キャッシュ 1 0 6 に格納された未検証書き込みデータ 1 0 6 a から回復することができる。

40

## 【 0 0 4 8 】

ブロック 5 0 4 で C V S T アプリケーション 1 1 2 によって実行される読み取りコマンドの数は、複数の所定の条件 1 1 4 によって決定することができる。たとえば第 1 の例示

50

的条件において、ホストのアイドル時間が2秒の場合、C V S Tアプリケーション112は順序付けリスト110の20%を実行することができる。第2の例示的条件では、順序付けリスト110から選択されるエントリの数を、ホストの作業負荷に基づくものとしてすることができる。第3の例示的条件では、順序付けリスト110から選択されるエントリの数を、使用されるキャッシュ106の量に基づくものとしてすることができる。たとえば、キャッシュが90%満杯である場合、C V S Tアプリケーション112は順序付けリスト110の5%を選択することができる。このパーセントの数字を調整して、ディスク118の読み取りおよび書き込み性能を最適化することができる。

【0049】

したがって図5は、C V S Tアプリケーション112が所定の条件114を使用して、ディスク118への書き込み中に断続的書き込みエラーが生じたかどうかを検証する、ある種の実施形態を示す。

10

【0050】

図6は、ある種の実施形態に従った、ディスク・ドライブ102で実施される読み取り要求を処理するための動作を示す図である。

【0051】

制御はブロック600で開始され、ディスク・ドライブ102がコンピュータ・デバイス104から読み取り要求を受け取る。ディスク・ドライブ102は(ブロック602で)、読み取り要求に対応するデータがキャッシュ106内に存在する未検証書き込みデータ106aであるかどうかを判別する。存在する未検証書き込みデータ106aである場合、(ブロック604で)キャッシュ106からのデータを戻すことによって、読み取り要求が満たされる。存在する未検証書き込みデータ106aでない場合、(ブロック606で)ディスク118からのデータを戻すことによって、読み取り要求が満たされる。

20

【0052】

したがって図6は、読み取り要求に対応するデータが未検証書き込みデータ106aである場合、外部ホストからの読み取り要求がキャッシュ106から満たされる、ある種の実施形態を示す。結果として、読み取り要求に応答して戻されるデータはエラーではない。

【0053】

ある種の実施形態は、リスト維持管理アプリケーション108およびC V S Tアプリケーション112を提供し、C V S Tアプリケーション112は、リスト維持管理アプリケーション108によってC V S Tアプリケーション112に提供された順序付けリスト110を使用することによって、キャッシュ106内の書き込みデータの一部とディスク118から読み取られたデータとを比較する。ある種の実施形態では、キャッシュ106に格納される書き込み数を削減することによって性能を向上させる。さらにある種の実施形態は、書き込みキャッシュ106がほぼ満杯になった際にP P S Tを実施するP P S TアプリケーションがC V S Tアプリケーションに置き換えられた場合、断続的および永続的の両方の喪失した書き込みに対する保護も提供する。

30

【0054】

ある種の実施形態は、検出段階間の書き込みキャッシュ106のコンテンツを削減するために、リスト維持管理アプリケーション108およびC V S Tアプリケーション112を使用しながら、喪失した書き込みのチェックとしてP P S Tアプリケーション116を使用することにより、永続的に喪失した書き込みに対する完全な保護を提供することができる。このケースでは、書き込みキャッシュ106がほぼ満杯になり、エラーが検出されなければ書き込みキャッシュに入れられたエントリがフラッシュされる場合、P P S Tアプリケーション116を呼び出すことができる。

40

【0055】

ある種の実施形態では、データの回復は背景モードまたはリアルタイムで実行することができる。ある種の実施形態は、インターフェースから独立した任意のディスク・ドライブに関する任意の種類<sub>の</sub>喪失した書き込み問題に対する保護を提供する。喪失した書き込

50



みは断続的または永続的とすることが可能であり、ホスト・アプリケーションへの中断（disruption）の量を削減することが可能である。ディスク・ドライブでの喪失した書き込みを保護するある種の実施形態は、入力／出力（I/O）性能に大幅な影響を与えないことが可能である。

#### 【0056】

リスト維持管理アプリケーション108を含めることにより、キャッシュ106がほぼ満杯の状態に達するまでの期間を引き延ばし、P P S TまたはC V S T検証のオーバーヘッドによって生じる可能性のある性能の影響を遅らせる。ある種の実施形態では、C V S TおよびP P S Tは長いアイドル・サイクル中にキャッシュ106全体をフラッシュすることができる。ホストの作業負荷が少ない場合、削減の速度がホストからの書き込み追加速度よりも速い場合があり、その結果としてキャッシュ全体のフラッシュも可能である。したがってこうした状況では、キャッシュ106は決して満杯に近づくことはなく、後続のP P S TまたはC V S Tを実行する必要のない可能性がある。

#### 【0057】

ある種の実施形態では、論理ブロック・アドレスおよび関連するカウンタを単純化すること、順序付けリスト110内のエントリを再順序付けすること、およびディスク118に読み取りのキューを送信することにより、ディスク118へのあらゆる書き込みが検証される場合の検証のプロセスがより効率良くなる。ある種の実施形態では、重複する書き込みを単一の読み取りに組み合わせ、再順序付けされた読み取りをキューとして送信し、所与の時間についての読み取り数を最大限にするために異なるシリンダへのシークを最適化することにより、ディスク・ドライブの読み取り性能をさらに向上させることが可能である。

#### 【0058】

##### 追加の詳細な実施形態

説明される技法は、ソフトウェア、ファームウェア、マイクロコード、ハードウェア、あるいはそれらの任意の組み合わせ、またはそれらすべてを含む、方法、装置、または装置として実施することができる。本明細書で使用される「装置」という用語は、回路内（たとえば集積回路チップ、プログラマブル・ゲート・アレイ（PGA）、ASICなど）、あるいは、コンピュータ読み取り可能メディア（たとえばハード・ディスク・ドライブ、フロッピー・ディスク、テープなどの磁気記憶媒体）、光ストレージ（たとえばCD-ROM、DVD-ROM、光ディスクなど）、揮発性および不揮発性メモリ・デバイス（たとえば電氣的消去可能プログラマブル読み取り専用メモリ（EEPROM）、読み取り専用メモリ（ROM）、プログラマブル読み取り専用メモリ（PROM）、ランダム・アクセス・メモリ（RAM）、動的ランダム・アクセス・メモリ（DRAM）、静的ランダム・アクセス・メモリ（SRAM）、フラッシュ、ファームウェア、プログラマブル論理など）内、またはそれらすべてで実施される、プログラム命令、コード、あるいは論理またはそれらすべてを言い表す。コンピュータ読み取り可能メディア内のコードには、プロセッサなどのマシンによってアクセスし、実行することができる。ある種の実施形態では、実施形態が作成されたコードに、伝送メディアを介して、またはネットワークを介してファイル・サーバからもアクセスすることができる。こうしたケースでは、コードが実施される装置は、ネットワーク伝送回線、無線伝送メディア、空気、電波、赤外線信号を介して伝搬される信号などの、伝送メディアを含むことができる。もちろん当業者であれば、諸実施形態の範囲を逸脱することなく多くの修正が実行可能であること、および、装置が当分野で知られた任意の情報搬送メディアを含むことが可能であることを理解されよう。たとえば装置は、マシンによって実行されると結果として実行される動作を生じる命令を格納した、記憶媒体を含む。

#### 【0059】

図7は、ある種の実施形態が実施可能なシステム700を示すブロック図である。ある種の実施形態では、コンピュータ・デバイス104および記憶装置102をシステム700に従って実施することができる。システム700は、ある種の実施形態ではプロセッサ

704を含むことができる回路702を含むことができる。システム700は、メモリ706（たとえば揮発性メモリ・デバイス）およびストレージ708を含むこともできる。システム700のある種の要素は、コンピュータ・デバイス104および記憶装置102のうちの一部またはすべてで見られるか、または見られない場合がある。ストレージ708は、不揮発性メモリ・デバイス（たとえばEEPROM、ROM、PROM、RAM、DRAM、SRAM、フラッシュ、ファームウェア、プログラマブル論理など）、磁気ディスク・ドライブ、光ディスク・ドライブ、テープ・ドライブなどを含むことができる。ストレージ708は、内部記憶装置、取り付け型記憶装置、あるいはネットワーク・アクセス可能記憶装置、またはそれらすべてを含むことができる。システム700は、メモリ706にロード可能であり、プロセッサ704または回路702によって実行可能である、コード712を含むプログラム論理710を含むことができる。ある種の実施形態では、コード712を含むプログラム論理710をストレージ708に格納することができる。ある種の他の実施形態では、プログラム論理710を回路702内で実施することができる。したがって、図7ではプログラム論理710を他の要素とは別に示しているが、プログラム論理710はメモリ706あるいは回路702またはその両方で実施することができる。

10

#### 【0060】

ある種の実施形態は、人間またはコンピュータ読み取り可能コードをコンピュータ・システムに統合する自動化された処理によって、コンピュータ命令を展開するための方法を対象とすることが可能であり、コードとコンピュータ・システムと組み合わせることによって、説明された実施形態の動作を実行することができる。

20

#### 【0061】

図3～6に示された動作の少なくとも一部は、並行して、ならびに逐次、実行することができる。代替の実施形態では、動作の一部は異なる順序で実行する、修正する、または削除することができる。

#### 【0062】

さらに、ソフトウェアおよびハードウェア・コンポーネントの多くは、例示の目的で別々のモジュールにあるものとして説明してきた。こうしたコンポーネントを、より少ない数のコンポーネントに統合するか、またはより多くの数のコンポーネントに分割することが可能である。さらに、特定のコンポーネントによって実行されるものとして説明したある種の動作を、他のコンポーネントによって実行することも可能である。

30

#### 【0063】

図1～7で示されたかまたは参照されたデータ構造およびコンポーネントは、特定のタイプの情報を有するものとして説明される。代替の実施形態では、データ構造およびコンポーネントは異なる構造とすること、ならびにこれらの図面で示されたかまたは参照されたものよりも少ない、多い、または異なるフィールド、あるいは異なる機能を有することが可能である。したがって、諸実施形態の前述の説明は例示および説明の目的で提示してきた。諸実施形態を網羅するか、または開示された精密な形に限定することを意図するものではない。上記の教示に鑑みて、多くの修正および変形が可能である。

#### 【図面の簡単な説明】

40

#### 【0064】

【図1】ある種の実施形態に従ったコンピューティング環境を示すブロック図である。

【図2】ある種の実施形態に従った、コンピューティング環境に含まれるコンポーネントおよびデータ構造を示すブロック図である。

【図3】ある種の実施形態に従った、喪失した書き込みから回復するための動作を示す図である。

【図4】ある種の実施形態に従った、リスト維持管理アプリケーションで実施される動作を示す図である。

【図5】ある種の実施形態に従った、総合検証用自己テスト・アプリケーションで実施される動作を示す図である。

50

【図 6】ある種の実施形態に従った、読み取り要求を処理するための動作を示す図である。

【図 7】ある種の実施形態が実施されるシステムを示す図である。

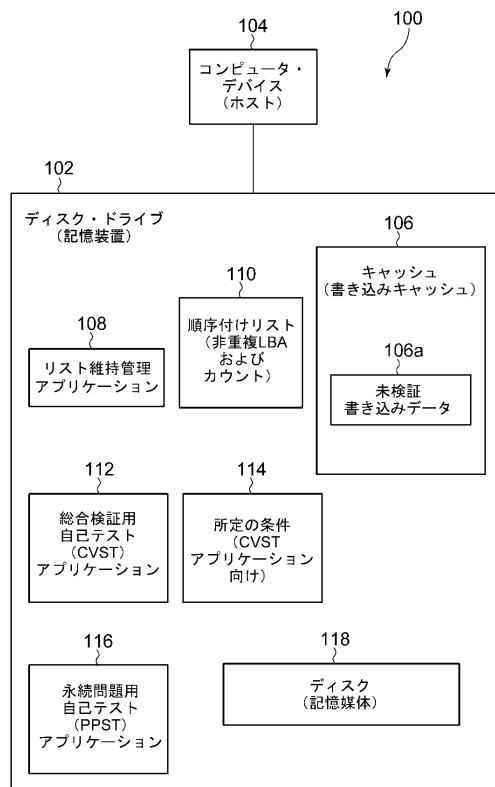
【符号の説明】

【 0 0 6 5 】

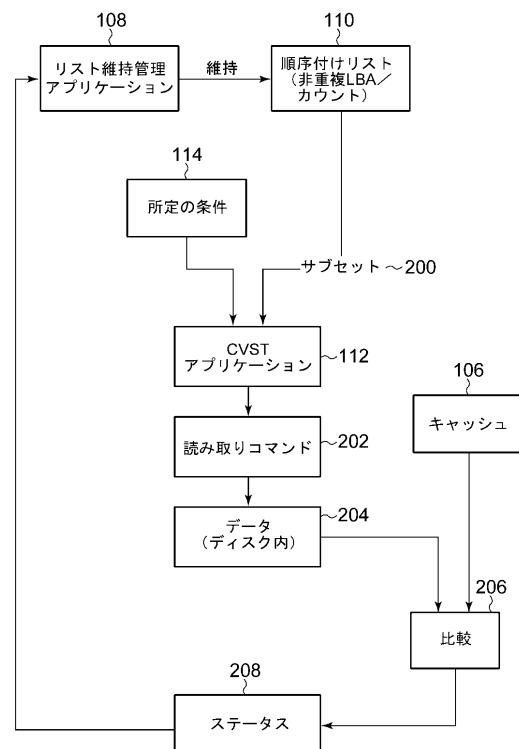
- 1 0 0      コンピューティング環境
- 1 0 2      ディスク・ドライブ（記憶装置）
- 1 0 4      コンピュータ・デバイス（ホスト）
- 1 0 6      キャッシュ（書き込みキャッシュ）
- 1 0 6 a     未検証書き込みデータ
- 1 0 8      リスト維持管理アプリケーション
- 1 1 0      順序付けリスト（非重複 L B A およびカウント）
- 1 1 2      総合検証用自己テスト（C V S T）アプリケーション
- 1 1 4      所定の条件（C V S T アプリケーション向け）
- 1 1 6      永続問題用自己テスト（P P S T）アプリケーション
- 1 1 8      ディスク（記憶媒体）

10

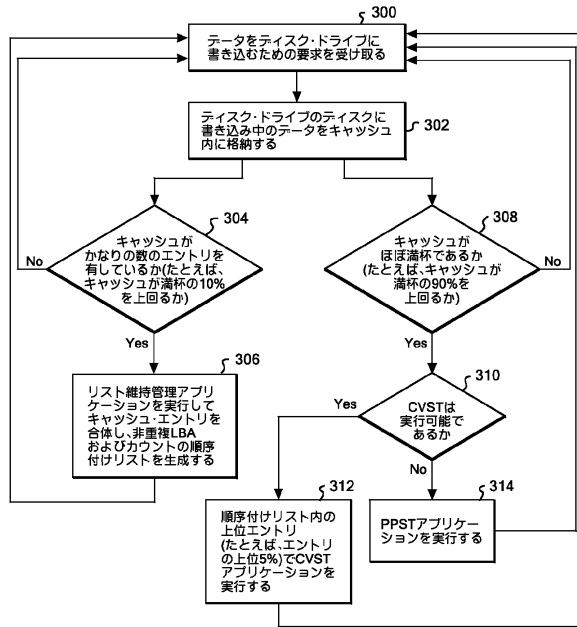
【図 1】



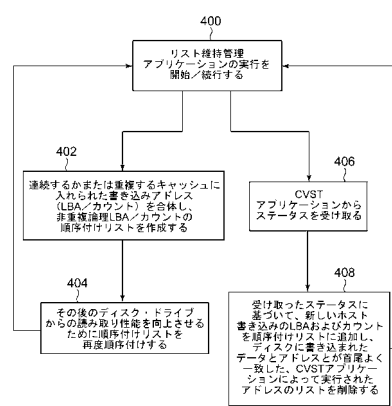
【図 2】



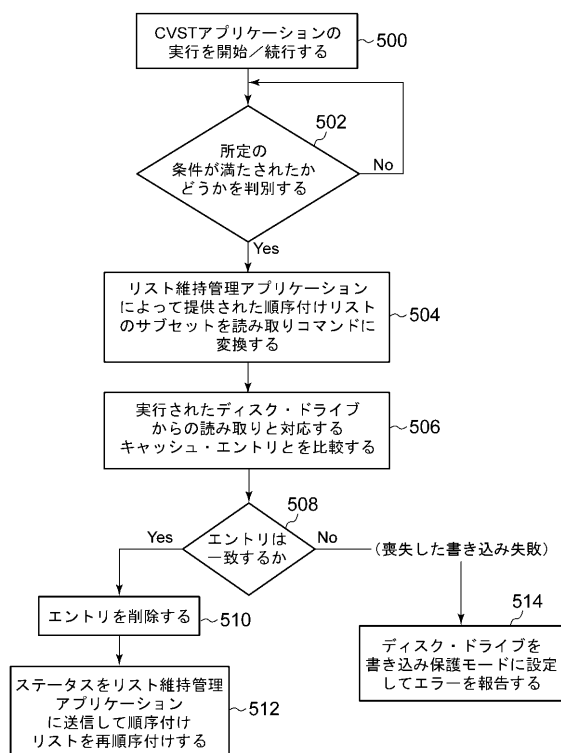
【図 3】



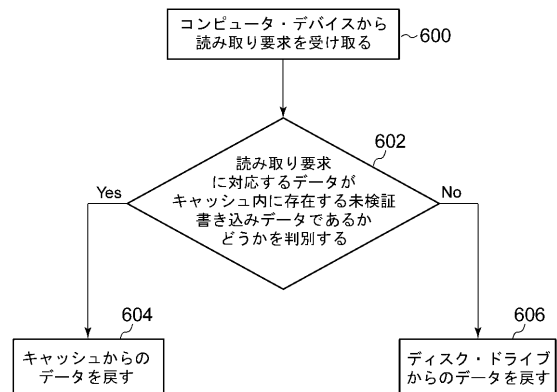
【図 4】



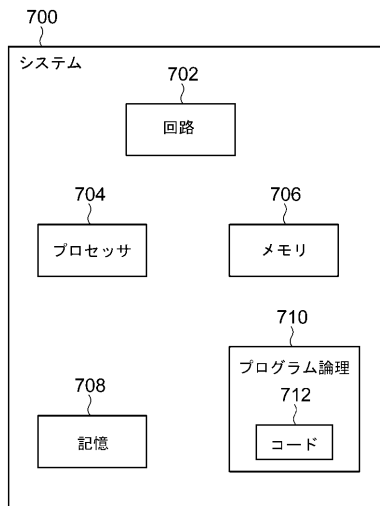
【図 5】



【図 6】



【図 7】



---

フロントページの続き

- (72)発明者 ウィリアム・ジョン・デュリカ  
アメリカ合衆国 9 5 0 3 7 カリフォルニア州モーガン・ヒル ライヴ・オーク・レーン 3 8 5 5
- (72)発明者 エム・アミン・ハッジ  
アメリカ合衆国 9 5 1 2 0 カリフォルニア州サンノゼ クエイル・クリーク・クラーク 1 2 1 1
- (72)発明者 ジョゼフ・スミス・ハイドⅡ世  
アメリカ合衆国 8 5 7 1 2 アリゾナ州トゥーソン イースト・プラシータ・デル・メスキート 5  
3 7 0
- (72)発明者 ロナルド・ジェイ・ヴェンチャーリ  
イギリス国 エス・オー 2 2 6 キュー・エル ウィンチェスター州リトルトン ザ・ホール・ウェ  
イ 2 3 ニューアーク

審査官 菅原 浩二

- (56)参考文献 特開平 0 4 - 0 1 7 0 2 0 ( J P , A )  
特開 2 0 0 1 - 1 4 2 6 5 0 ( J P , A )  
特開平 0 7 - 2 1 0 3 2 7 ( J P , A )  
特開 2 0 0 3 - 2 6 3 7 0 3 ( J P , A )

- (58)調査した分野(Int.Cl. , D B 名)  
G 0 6 F 3 / 0 6  
G 0 6 F 1 2 / 0 8