

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4163298号
(P4163298)

(45) 発行日 平成20年10月8日 (2008. 10. 8)

(24) 登録日 平成20年8月1日 (2008. 8. 1)

(51) Int. Cl.

F I

G 0 6 F 3 / 0 6 (2006. 01)

G 0 6 F 3 / 0 6 3 0 5 C

G 0 6 F 3 / 0 6 3 0 4 P

G 0 6 F 3 / 0 6 5 4 0

請求項の数 12 (全 18 頁)

(21) 出願番号	特願平10-235862	(73) 特許権者	398038580
(22) 出願日	平成10年8月21日 (1998. 8. 21)		ヒューレット・パカード・カンパニー
(65) 公開番号	特開平11-119919		HEWLETT-PACKARD COM
(43) 公開日	平成11年4月30日 (1999. 4. 30)		PANY
審査請求日	平成17年8月16日 (2005. 8. 16)		アメリカ合衆国カリフォルニア州パロアル
(31) 優先権主張番号	08/920-120		ト ハノーバー・ストリート 3000
(32) 優先日	平成9年8月26日 (1997. 8. 26)	(74) 代理人	100075513
(33) 優先権主張国	米国 (US)		弁理士 後藤 政喜
		(74) 代理人	100084537
			弁理士 松田 嘉夫
		(72) 発明者	ダグラス・エル・ヴォイグト
			アメリカ合衆国 アイダホ、ボイセ、エヌ
			・24ス 3030

最終頁に続く

(54) 【発明の名称】 記憶システムへのデータ書き込み方法

(57) 【特許請求の範囲】

【請求項 1】

コンピュータによって行われる、複数の記憶媒体を有する記憶システムへの書き込み方法であり、メモリに格納されているトランザクションログを、前記複数の記憶媒体に設けられたディスクステージングログ、および、前記ディスクステージングログとは異なるディスクログの2つの領域のうちのいずれか一方の領域に書き込む方法であって、

前記トランザクションログを格納しているメモリがページフルの状態になると、前記ディスクログに前記トランザクションログを書き込むことと、

前記メモリに格納されているトランザクションログの強制引出し要求に応じて、前記ディスクステージングログを構成する記憶媒体のうち、アクセス頻度が最低の記憶媒体（以下、最低使用頻度の記憶媒体）に前記トランザクションログを書き込むことと、

を含むことを特徴とする記憶システムへのデータ書き込み方法。

【請求項 2】

前記複数の記憶媒体は、ランダムアクセス記憶媒体であることを特徴とする請求項 1 に記載の記憶システムへのデータ書き込み方法。

【請求項 3】

前記最低使用頻度の記憶媒体を選択することは、入出力動作に基づいて当該の選択が行われることを特徴とする請求項 1 または 2 に記載の記憶システムへのデータ書き込み方法。

【請求項 4】

前記強制引出し的に書き込む要求は、前記トランザクションログを格納しているメモリがページフルの状態になる前に発生することを特徴とする請求項 1、2 または 3 に記載の記憶システムへのデータ書き込み方法。

【請求項 5】

前記トランザクションログは、前記選択された最低使用頻度の記憶媒体に非冗長に書き込まれることを特徴とする請求項 1 に記載の記憶システムへのデータ書き込み方法。

【請求項 6】

前記トランザクションログは、前記書き込まれるトランザクションログの順序を示す標識を含むことを特徴とする請求項 1 に記載の記憶システムへのデータ書き込み方法。

【請求項 7】

前記記憶媒体はそれぞれ、前記最低使用頻度の記憶媒体が選択された際の書き込みのためにのみ使用されるように確保された領域を含むことを特徴とする請求項 1 に記載の記憶システムへのデータ書き込み方法。

【請求項 8】

前記確保された領域は、少なくとも 2 つのサブ領域を含み、

最低使用頻度の記憶媒体が選択される最初の事象の発生時に前記サブ領域の 1 つに対して書き込みが行なわれ、

最低使用頻度記憶媒体が選択される次の事象の発生時に他方のサブ領域への書き込みが行なわれ、

それによって、同じ最低使用頻度記憶媒体が二度連続して選択された場合であっても、次の連続する書き込みにおいて直前に書き込まれたサブ領域に対して重ね書きが生じないことを特徴とする請求項 7 に記載の記憶システムへのデータ書き込み方法。

【請求項 9】

記憶システムであって、

(a) データ記録を保持する第 1 のメモリと、

(b) 前記第 1 のメモリに接続され、ディスクステージングログ、および、前記ディスクステージングログとは異なるディスクログが設けられる複数の記憶媒体と、

(c) 前記第 1 のメモリの状態を検出する手段と、

(d) 前記第 1 のメモリがページフルの状態であることが検出されると、前記ディスクログに前記データ記録を書き込む手段と、

(e) 前記第 1 のメモリに保持されている前記データ記録の強制引出し要求に応じて、前記ディスクステージングログを構成する記憶媒体のうち、アクセス頻度が最低の記憶媒体（以下、最低使用頻度の記憶媒体）に前記データ記録を書き込む手段と、を含むことを特徴とする記憶システム。

【請求項 10】

前記第 1 のメモリがページフルの状態であることが検出された時に、前記ディスクログに前記データ記録を書き込む手段は、前記複数の記憶媒体上の前記データ記録の冗長性を維持することを特徴とする請求項 9 に記載の記憶システム。

【請求項 11】

前記第 1 のメモリに保持されている前記データ記録の強制引出し要求に応じて、前記最低使用頻度の記憶媒体に前記データ記録を書き込む手段は、前記複数の記憶媒体上において前記データ記録の冗長性を維持しないことを特徴とする請求項 9 に記載の記憶システム。

【請求項 12】

前記最低使用頻度の記憶媒体は、前記データ記録の書き込みのために確保された少なくとも 2 つのサブ領域を含むことを特徴とする請求項 9 に記載の記憶システム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は一般的にはデータ記憶システムに関し、特にディスクアレイ記憶システムのためのトランザクションログ（transaction log）管理に関する。

10

20

30

40

50

【 0 0 0 2 】

【従来の技術】

コンピュータシステムの速度、信頼性および処理能力は絶えず進歩し続けている。その結果、コンピュータはより複雑で高度なアプリケーションを処理することができる。コンピュータの改良にともなって、大量記憶および入出力（I/O）装置の性能に対する要求も高くなる。したがって、進歩し続けるコンピュータシステムに性能上つりあう大量記憶システムを設計することが常に必要とされている。

【 0 0 0 3 】

本発明は特にディスクアレイ型の大量記憶装置に関する。ディスクアレイデータ記憶システムは単一の大量記憶システムを形成するように構成され、統合された多数の記憶ディスクドライブ装置を有する。大量記憶システムには、コスト、性能および利用可能性という3つの主要な設計評価基準がある。メガバイトあたりのコストが低く、入出力性能が高く、データの利用可能性が高いメモリ装置を製作することが最も望ましい。“利用可能性”とは、記憶システムに記憶されたデータにアクセスする能力のことであり、またなんらかの故障があった場合に連続動作を保証する能力のことをいう。通常、データの利用可能性は冗長性を用いて提供され、この場合データあるいはデータ間の関係が複数の場所に記憶される。冗長データの記憶には、“ミラー”法および“パリティ”法の2つの一般的な方法がある。

【 0 0 0 4 】

【発明が解決しようとする課題】

ディスクアレイデータ記憶システムの設計にあたって発生する問題の1つは、システムの誤りあるいは故障の場合における記憶されたデータの正確なマッピング情報の保持の問題に係るものである。これは前記の冗長データ記憶法のいずれかあるいはその両方を用いるシステムについていえることである。したがって、ディスクアレイマッピング情報の管理において、誤りから回復する目的のためには、最近変更されたマッピング情報がディスク上に確実に記憶されるようにする必要がある場合が多い。このディスク書き込みの必要性は、(i) 時間に基づいた頻度状態の更新、(ii) ログページ・フルが状態、あるいは(iii) 特定のホストによる要求等のいくつかの理由で発生する。一般的には、最近のマッピング情報の変更は、ディスクアレイ機能の性能対して最適化されたデータ構造内におけるランダムな場所に蓄積され、さらに他のデータ構造より高速にディスクに書き込む（ポストする（post））ことができるログに順次蓄積される。この技術はトランザクション処理技術においては周知である。しかし、ポストの必要性が進行中の他のディスク読み出しあるいは書き込み動作と同時に発生してシステム内に入出力の競合が生じるという問題が発生する場合がある。かかる入出力競合は、特にポストが頻繁に発生する場合にシステム上の重要な性能を阻害することが多い。これはディスクにログを1回ポストするには複数の入出力事象が必然的に発生するためである。たとえば、通常、ログページはまず無効と表示される（すなわち、更新が必要である）。次に、ログページはディスクにコピーされ、その後有効と表示される。最後に、冗長システムでは、冗長ログページがディスクにコピーされる。

【 0 0 0 5 】

以上のことから、また増大し続ける計算速度および管理対象となる膨大な情報量から、ディスクアレイシステム等の性能の改善が常に必要とされている。

【 0 0 0 6 】

したがって、本発明はディスクアレイマッピング情報の管理システムの性能を向上させることを目的とする。また、複数の利用可能な記憶ディスクから選択された任意の最も使用頻度の低いディスクへのログ書き込みを管理および分散してログ入出力と進行中の他の入出力との間におけるディスクアクセスの競合を低減することによって、ディスクログ書き込みのシステム性能を改善し、さらにシステムの誤り又は故障時に記録を確実に回復可能にすることを目的とする。

【 0 0 0 7 】

【課題を解決するための手段】

本発明による記憶システムへの書き込み方法は、コンピュータによって行われる、複数の記憶媒体（１２）を有する記憶システム（１０）への書き込み方法であって、

（a）第１のトランザクションログを記憶媒体（１２）に書き込む第１の要求を示す第１の基準を検出することと、

（b）第１の要求に応じて、第１のトランザクションログを、最低使用頻度の記憶媒体であるかないかにかかわらず、記憶媒体（１２）に書き込むことと、

（c）第２のトランザクションログを記憶媒体（１２）に強制引出し的に書き込む要求を含む第２の要求であって、第２のトランザクションログを記憶媒体（１２）に書き込む第２の要求を示す第２の基準を検出することと、

（d）複数の記憶媒体（１２）から最低使用頻度の記憶媒体を選択することと、

（e）第２の要求に応じて、第２のトランザクションログを選択された最低使用頻度の記憶媒体へ書き込むこととを含むことを特徴とする。

また、本発明による記憶システムは、

（a）データ記録（１１０）を保持する第１のメモリ（５５）と、

（b）第１のメモリに接続された複数の記憶媒体（１２）と、

（c）第１のメモリ（５５）の状態を検出する手段（１６）と、

（d）第１のメモリにおける第１の検出された状態に応答して第１の記憶管理基準にしたがって記憶媒体（１２）にデータ記録を書き込む手段（１６）と、

（e）第１のメモリにおける第２の検出された状態に応答して第２の記憶管理基準にしたがって記憶媒体（１２）にデータ記録を書き込む手段（１６）とを含む、

第１の記憶管理基準は、最低使用頻度の記憶媒体であるかないかにかかわらず、記憶媒体（１２）への書き込みを示す基準を含み、第２の記憶管理基準は、複数の記憶媒体における最低使用頻度の記憶媒体への書き込みを示す基準を含み、第１のメモリ（５５）における第１の検出された状態は、第１のメモリの一部がデータ記録で満たされていることを示す状態を含み、第１のメモリ（５５）における第２の検出された状態は、第１のメモリにおいて所定の部分がデータ記録で満たされる前に、第１のメモリ（５５）内のデータ記録を複数の記憶媒体（１２）に強制引出し的に書き込む要求を示す状態を含むことを特徴とする。

【０００８】

一実施形態における本発明の原理によれば、ディスクドライブ等の複数の記憶媒体を有する記憶システムにおいて、第１メモリに記憶されたトランザクションログが、２つの異なるログ領域に選択的にポストされる。詳細には、第１メモリのトランザクションログのページ・フル状態が検出されると、ポストは"ディスクログ"領域に対して行なわれる。ポスト要求がトランザクションログのページ・フル状態が検出される前に発生すると、ポストはただちに"ステージングログ"領域の最も使用頻度の低いディスクに対して行なわれる。

他の原理によれば、"ディスクログ"へのポストは通常の記憶システム管理技術およびデータ冗長性技術を用いて行なわれる。一方、"ステージングログ"領域へのポストは、ステージングログのデータが記憶媒体上で冗長な状態で保持されることのないように通常の記憶システムデータ管理および冗長性技術を用いずに行なわれる。冗長性はトランザクションログがステージングログ領域にコピーされることに加えて第１メモリに残ることによって維持される。

【０００９】

ステージングログ領域は複数の記憶媒体のそれぞれに確保されたスペースを含み、かかる確保されたスペースは各記憶媒体上で論理的に分離された部分に分割される。この構成によって、ステージングログへのポストを確保された部分の間で"トグルする（toggle）"ことができる。したがって、２つの連続するステージングログポストによって、どのディスクが最も使用頻度が低いかにかかわらず確保された領域の同じ部分が重ね書きされ

10

20

30

40

50

ることではない。

【 0 0 1 0 】

他の原理によれば、記憶媒体上のログ領域にポストされるデータ記録にシーケンス番号およびディスク群番号が割り当てられる。ディスクログ領域およびステージングログ領域からのデータの回復においては、シーケンス番号およびディスク群番号が参照されて完全なトランザクションログが適正に再構築される。

【 0 0 1 1 】

本発明の他の目的、利点および機能は以下の説明から明らかになるであろう。

【 0 0 1 2 】

【発明の実施の形態】

図 1 には、本発明の分散書き込みディスクログ法を用いたデータ記憶システム 10 のブロック図を示す。図示する例では、データ記憶システム 10 は階層的ディスクアレイ 11 を含むディスクアレイデータ記憶システムである。本発明は非階層的アレイ（図示せず）にも適用可能である。ディスクアレイ 11 は R A I D (Redundant Array of Independent Disks) 記憶システムを実施するための複数の記憶ディスク 12 を含む。データ記憶システム 10 はディスクアレイ 11 に結合され記憶ディスク 12 との間のデータ転送を調整するディスクアレイコントローラ 14 を含み、さらに R A I D 管理システム 16 を含む。R A I D 管理システム 16 は本発明の分散書き込みディスクログ法を実行する手段を含む。

【 0 0 1 3 】

本明細書においては、“ディスク”とは自己の記憶故障を検出することのできる任意の不揮発性のランダムアクセス可能・書き換え可能な大量記憶装置である。ディスクには、回転磁気ディスクおよび光ディスクとソリッドステートディスクの両方あるいは（PROM、EPROM および EEPROM 等の）不揮発性電子記憶素子を含む。“ディスクアレイ”という用語は、ディスクと、ディスクを 1 つあるいはそれ以上のホストコンピュータに接続するのに要するハードウェアと、物理的ディスクの動作を制御しそれらをホスト動作環境に 1 つあるいはそれ以上の仮想ディスクとして提示するのに必要な管理ソフトウェアの集合である。“仮想ディスク”は管理ソフトウェアによってディスクアレイ中に実現される抽象的存在である。

【 0 0 1 4 】

“R A I D”という用語はその物理的記憶容量の一部が記憶容量の残りの部分に記憶されたユーザデータに関する冗長な情報の記憶に用いられるディスクアレイを意味する。この冗長情報によって、このアレイを構成するディスクの 1 つあるいはこのアレイへのアクセス経路が故障した場合にユーザデータを再生することができる。R A I D システムについては、ミネソタ州 Lino Lakes の REID Advisory Board から 1993 年 6 月 9 日に刊行された“The RAID Book: A Source Book for RAID Technology”に詳細に説明されている。R A I D システムを本発明との関係において例示するが、本発明は非 R A I D システムにも適用可能であることはいうまでもない。

【 0 0 1 5 】

ディスクアレイコントローラ 14 は、small computer system interface (S C S I) 等の 1 つあるいはそれ以上のインターフェースバス 13 を介してディスクアレイ 11 に結合されている。R A I D 管理システム 16 はインターフェースプロトコル 15 によってディスクアレイコントローラ 14 に操作的に結合されている。R A I D 管理システム 16 は図示するように別個の要素として（すなわちソフトウェアあるいはファームウェアとして）実施することもでき、あるいはディスクアレイコントローラ 14 内あるいはホストコンピュータ内に構成して、ディスクの記憶および信頼性レベルの制御、さまざまな信頼性の記憶装置レベル間でのデータ転送、および本発明の分散書き込みディスクログの実施を行なうデータ管理手段を提供することもできる。また、データ記憶システム 10 は入出力インターフェースバス 17 を介してホストコンピュータ（図示せず）に結合されている。

【 0 0 1 6 】

図示するシステムでは、ディスクアレイコントローラ 14 はディスクアレイコントローラ

10

20

30

40

50

“ A ” 1 4 A およびディスクアレイコントローラ “ B ” 1 4 B からなるデュアルコントローラとして実施される。デュアルコントローラ 1 4 A および 1 4 B は一方のコントローラが動作不能となったとき連続的なバックアップと冗長性を供給することによって信頼性を向上させる。しかし、本発明の方法は単一のコントローラあるいは他のアーキテクチャで実施することができる。実際に、本発明は完全で正確なディスクログの維持がデュアルコントローラ環境におけるよりも重要である単一コントローラアーキテクチャにおいて特に有益である。

【 0 0 1 7 】

階層的ディスクアレイ 1 1 は物理的記憶スペースと 1 つあるいはそれ以上の仮想記憶スペースを含む異なる記憶スペースとして特徴付けることができる。たとえば、ディスクアレイ 1 1 内の記憶ディスク 1 2 は、複数のディスク 2 0 のミラーグループ 1 8 および複数のディスク 2 4 のパリティグループ 2 2 に構成されるものとして概念化することができる。記憶装置のかかる諸相はマッピング技術を用いて関係付けられる。たとえば、ディスクアレイの物理的記憶スペースは記憶領域をさまざまなデータ信頼性レベルに応じて区分した仮想記憶スペースにマップすることができる。仮想記憶スペース内の領域の一部をミラーすなわち R A I D レベル 1 の第 1 の信頼性の記憶レベルに割り当て、他の領域をパリティすなわち R A I D レベル 5 の第 2 の信頼性の記憶レベルに割り当てることができる。かかる領域は同じディスクあるいは別々のディスク上に構成することができ、また任意の組み合わせのディスク上に構成することもできる。

【 0 0 1 8 】

データ記憶システム 1 0 はディスクアレイ 1 1 のマッピングに用いる仮想マッピング情報の永続的な記憶を可能とするメモリマップ記憶域 2 1 を含む。このメモリマップ記憶域はディスクアレイの外部にあり、好適にはディスクアレイコントローラ 1 4 に常駐する。メモリマッピング情報は異なるビュー (view) の間でさまざまなマッピング構成が変化するにつれてディスクアレイコントローラ 1 4 あるいは R A I D 管理システム 1 6 によって連続的あるいは定期的に更新することができる。

【 0 0 1 9 】

好適には、メモリマップ記憶域 2 1 はそれぞれディスクアレイコントローラ “ A ” 1 4 A およびディスクアレイコントローラ “ B ” 1 4 B に設けられた 2 つの不揮発性 R A M (Non-Volatile RAM) 2 1 A および 2 1 B として実施される。これら 2 つの N V R A M 2 1 A および 2 1 B はメモリマッピング情報の冗長記憶を可能とする。仮想マッピング情報はミラー冗長性技術によって N V R A M 2 1 A および N V R A M 2 1 B の両方に複製され記憶される。これによって、 N V R A M 2 1 A をオリジナルのマッピング情報の記憶にのみ使い、 N V R A M 2 1 B を冗長マッピング情報の記憶にのみ用いることができる。

【 0 0 2 0 】

図示するように、ディスクアレイ 1 1 は複数の記憶ディスク 1 2 を有する。記憶ディスク 1 2 上の冗長性の管理は R A I D 管理システム 1 6 によって統御される。ユーザすなわちホストアプリケーションプログラムから見た際、アプリケーションレベルの仮想ビューによって記憶ディスク 1 2 上の利用可能な記憶スペースを示す 1 つの大きな記憶容量を表わすことができる。 R A I D 管理システム 1 6 はこの物理的記憶スペース上での R A I D 領域の構成を動的に変更することができる。その結果、 R A I D レベル仮想ビュー内の R A I D 領域のディスクへのマッピングおよびフロントエンド仮想ビューの R A I D ビューへのマッピングは一般的にはある変化の状態ということになる。 N V R A M 2 1 A および N V R A M 2 1 B 内のメモリマップ記憶域は、 R A I D 管理システム 1 6 による R A I D 領域のディスクへのマッピングに用いられる現在のマッピング情報および 2 つの仮想ビューの間でのマッピングに用いられる情報を保持する。 R A I D 管理システムは R A I D レベルのマッピングを動的に変更するとき、メモリマップ記憶域のマッピング情報にかかる変更を反映するように更新する。

【 0 0 2 1 】

しかし、ディスクアレイに用いられる R A I D 機構すなわちデータ記憶機構にかかわりな

10

20

30

40

50

く、メモリマップ記憶域 21 は一般的にはシステムの使用時全体を通じて常に変化する状態にあることは明らかである。したがって、メモリマップログ記録がメモリに保持され、RAID 管理システム 16 によってメモリからディスクに絶えずポストされ、NVRAM 21 の損失時にかかる記録が確実に回復されるようにする。よって、本発明は複数の利用可能な記憶ディスク 12 から選択された任意の最も使用頻度の低いディスクへのログ書き込みを管理および分散してログ入出力と進行中の他の入出力との間におけるディスクアクセスの競合を低減することによってディスクログ書き込みのシステム性能を改善するものである。一般的には、これはログの新しい部分を保持するために各記憶ディスク 12 上に“ステージングログ”領域を確保することによって行なわれる。そして、トランザクションログメモリのページがいっぱいになる前にポスト要求が発生した場合、ポストは最も使用頻度の低いディスクに確保された“ステージングログ”領域にほとんど即時に実行される。続いて、ログの回復が必要である場合、すべての記憶ディスク 12 からの断片がまとめられて単一の完全なイメージが形成される。

【0022】

図 2 は、本発明の分散ログ書き込みディスクログ法を示すブロック図である。NVRAM マップ 45 はデータ記憶システム 10 に用いられるデータが記憶されるディスクアレイコントローラ 14A / 14B (図 1) 上における不揮発性のメモリマップ記憶域 21 の部分集合を表わす。ディスクマップ 50 はディスクアレイ 11 に (冗長的に) 属する NVRAM マップ 45 の従来のディスクマップイメージである。ディスクマップ 50 への NVRAM マップ 45 の定期的記憶 (ポスト) によって誤り訂正を行なうために、NVRAM マップ 45 の内容の冗長コピーをディスクに維持する手段が提供される。一般的には、ディスクマップ 50 への NVRAM マップ 45 のポストは、通常のシステム処理および入出力競合状態下で可能である際に (RAID 管理システム 16 による制御のもとに) バックグラウンド処理として実行される。よって、ディスクマップ 50 への NVRAM マップ 45 データのポストは、入出力およびディスクスペースに関する通常のシステム競合の影響を受け、したがってポストが実際にいつ発生するかについては不確定要素がある。

【0023】

RAM ログイメージ (RLI) 55 もまた不揮発性メモリ 21 の部分集合である。あるいはこれは別個の (好適には不揮発性の) メモリとすることもできる。RLI 55 は NVRAM マップ 45 内で発生するインクリメンタルな変化を迅速に記憶 / 記録するのに用いられる。一実施形態では、RLI 55 は 16 (図では N) のアドレス指定可能な 64 K バイトページを含むが、他の構成も可能である。その後 RAID 管理システム 16 からの要求があった際、RLI 55 に記憶されたインクリメンタルな変化はディスクログ 60 あるいはディスクステージングログ 65 にポストされる。

【0024】

いくつかの要因によって、RAID 管理システム 16 に RLI 55 からディスクログ 60 あるいはディスクステージングログ 65 へのデータのポスト要求を行なう。たとえば、一実施形態では、RLI が“ページ・フル”状態になった際、ディスクログ 60 への“排出 (flush)”ポスト要求が発生する。“排出”ポスト要求が発生すると、RLI 55 からのトランザクションログデータのフル・ページがディスクログ 60 に書き込まれる。一方、(i) 時間に基づく頻度の要求が発生するか、(ii) ある特定のホスト要求を受けた場合に、ディスクステージングログ 65 への“強制引出し”ポスト要求が発生する。“強制引出し”ポスト要求が発生すると、RLI 55 から 1 つあるいはそれ以上のトランザクションログデータのフル・ブロックがディスクステージングログ 65 に書き込まれる。この 1 つあるいはそれ以上のブロックには、前に完全に書き込まれておらず、また未書き込みの 1 つあるいはそれ以上のトランザクションログ記録を含む (現在のページ内の) ブロックが含まれる。ディスクログ 60 に“排出”ポストで書き込みされるページおよびディスクステージングログ 65 に“強制引出し”ポストで書き込みされるブロックをここでは RLI 55 の“未書き込み”データと称する (ただし、“排出”ポストで書き込みされたページはそれ以前にディスクステージングログ 65 に“強制引出し”ポストで書き込みされた

記録を含む場合がある)。いずれの場合にも、(ディスクログ60あるいはディスクステージングログ65への)かかるポストによって、(RLI55で捕捉された)NVRAMマップ45の変化が、ディスクマップ50が更新されていない際にNVRAM21に損失が生じた場合においての、誤り回復目的のためにディスクアレイ11に確実に記憶されることを保証する。

【0025】

ディスクログ60はディスクアレイ11(図1)上に常駐するRLI55の従来と同様のディスクイメージである。好適には、ディスクログ60はRLI55と同様に多数のデータページを記憶することができる。図示するように、ディスクログ60にはデータ記憶用のNのページが示され、また従来と同様に、連続的あるいは円形にリンクすることができる。ディスクログ60は(図1の)通常データ冗長性機構を用いてディスクアレイ11上に記憶および管理される。したがって、ディスクログ60へのRLI55における“未書き込み”内容の“排出”ポストは通常の入出力状態で発生し、ディスクアクセスおよびスペースのためのシステム入出力競合の影響を受ける。ディスクログ60は一般的にはディスクマップ50より頻繁に更新されるが、ディスクログ60は(RLI55中で捕捉された)NVRAMマップ45へのインクリメンタルな変化のみを保持するのに対して、ディスクマップ50は(最終更新時の)NVRAMマップ45の完成イメージを保持する。

【0026】

ディスクステージングログ65はディスクアレイ11(図1)のディスクの各部分を表わす確保されたステージング領域70、75、80、85、90、95、100および105を含む。上述したように、一実施形態では、ディスクステージングログ65は指定された事象の発生時あるいは“ページ・フル”状態以外の時にRLI55の内容の記憶に用いられる。しかし、このポスト基準には当業者には明らかなようにシステム設計の変更および/またはユーザによる優先的な指示に合わせて自由度を持たせることができる。いずれの場合にも、(“ページ・フル”状態以外の)RAID管理システム16によって要求された所定の事象が発生した場合、RLI55はその“未書き込みの”内容をディスクステージングログ65のディスク1~Mのうち最も使用頻度の低いディスクに“強制引出し”する。最低使用頻度のディスクはディスクアレイ11のディスク1~Mの入出力動作をモニターすることによって検出される。

【0027】

基本的には、最低使用頻度のディスクへRLI55を“強制引出し”することによって、時間の経過とともにディスクアレイ全体にトランザクションログの分散書き込みが実行される。これは、所定の単一のディスクログ60へのRLI55のページ・フル“排出”とは対照的である。パリティ冗長機構を用いる場合、ディスクログ60は実際には複数のディスクに分散させることができるが基本的には“単一の”すなわち“非分散の”ディスクログである。これは(冗長コピーを考えなければ)1つのディスクドライブ上で1つの基底アドレスのみを用いてログ全体をアドレス指定/アクセスすることができるためである。

【0028】

最低使用頻度のディスクが選択されるため、“強制引出し”ポストが(進行中の他のシステム呼び出し/書き込み入出力動作との)入出力競合が低減された状態で発生するという利点がある。したがって、ディスクマップ50あるいはディスクログ60のポストとは異なり、ディスクステージングログ65へのこの分散書き込みは一般的にはただちに(あるいは、少なくともより迅速に)完了する。さらに、“強制引出し”ポストは一般的には転送される未書き込みブロックが最小限であるため“排出”ポストより高速である。

【0029】

ディスクログ60とは異なり、ディスクステージングログ65はRLI55のインクリメンタルな変化をディスクアレイ11全体にわたって分散した非冗長的な態様で保持する。これは、ディスクステージングログ65に発生する書き込みがRAID管理システム16の通常データ冗長性機構から除外されることから非冗長的である。したがって、ディスクステージングログ65へのポストにおいて発生する入出力ステップはディスクログ60の場合

10

20

30

40

50

に比べて少なくとも１ステップ少ない。トランザクションログはステージングログ領域にコピーされることに加えて第１メモリ（ＲＬＩ５５）にも残るため、“強制引出し”ポスト後にも冗長性が維持される。

【００３０】

一実施形態では、ディスクアレイ１１のディスク１～Ｍはそれぞれ分散されたログの記憶のために確保された専用のスペース量を有する。たとえば、図示するそれぞれのディスク上には２つの６４Ｋバイトページ７０／７５、８０／８５、９０／９５および１００／１０５が確保される。ディスクステージングログポスト処理中の故障の場合に発生する可能性のある有効データの重ね書き（および損失）を防止するために、それぞれのディスク上に少なくとも２つのページが確保される。すなわち、ＲＬＩ５５はディスクステージングログ６５への書き込み（ポスト／強制引出し）を奇数および偶数ページに交互に実行する（スワッピングすなわちトグル）。たとえば、最初の書き込みでは、ＲＬＩ５５は最低使用頻度のディスクに確保された偶数番号のページ７０、８０、９０あるいは１００にポストを行なう。次に発生する書き込みにおいては、ＲＬＩ５５は最低使用頻度のディスクの奇数番号ページ７５、８５、９５あるいは１０５にポストする。これによって、システムはデータの完全性の別のレベルを保証され、続いて発生するポストにおいて最後にポストされたデータの重ね書きの可能性（すなわち、同じ最低使用頻度ディスクが連続して選択される場合の重ね書き）が防止される。

【００３１】

図３から図６は、ある時間における本発明のディスクステージングログの状態を示すブロック図である。本発明の分散書き込みディスクステージング動作の例をさらに詳細に説明するために、ＲＬＩ５５のログイメージページ５７の一部、およびディスクステージングログ６５の各ページ７０～１０５の一部を示す。すなわち、図３から図６にはそれぞれＲＬＩ５５からの異なるポストにตอบสนองしてディスクステージングログ６５の状態時間における異なるスナップショットを示す。ＲＬＩ５５のログイメージページ５７およびディスクステージングログ６５の各ページ７０～１０５は破線によって３つの５１２バイトブロック（すなわち部分）Ｂ１、Ｂ２およびＢ３に（論理的に）分割されているものとして示されている。説明を簡略化するために、それぞれの６４Ｋバイトページ中の全ブロックではなく３つのブロックのみを示す。（ＲＬＩ５５の）ログイメージページ５７を本説明および図では“ＬＩ”で示す。さらに、ディスクステージングログ６５中の各ディスクをそれぞれ“Ｄ１～ＤＭ”で示し、各ディスク中に確保される２つのページをそれぞれ“Ｐ１”あるいは“Ｐ２”で示す。

【００３２】

図３において、論理標識である事象／時刻Ｔ１は（図１のＲＡＩＤ管理システム１６の要求によって）ある特定の事象が発生してＲＬＩ５５におけるログイメージページ５７の未書き込みデータのディスクステージングログ６５への“強制引出し”ポストが開始されるある所定の時点を反映している。事象／時刻Ｔ１はさらにその所定の時点においてＲＬＩのログイメージページ５７にログデータがどれだけ“フル”かを示す位置を同定する。ポストが要求されると、論理標識である事象／時刻Ｔ１によって同定される（ＲＬＩ５５におけるログイメージページ５７の）未書き込みデータの全ブロックがポストされる。５１２バイトデータブロックは（この例では任意のシステム設計条件に対する）最小ポストサイズであるため全てのブロックがポストされる。

【００３３】

したがって、たとえば、事象／時刻Ｔ１の発生時にＲＬＩ５５はログイメージページ５７の（事象／時刻Ｔ１で示される位置）“未書き込みの”内容をディスクステージングログ６５の最低使用頻度のディスク１～Ｍにおけるページ７０～１０５の１つに（後に詳述する“トグル”状態で）ポスト（post）する。すなわち、ログイメージページ５７のブロック“１”（ＬＩＢ１）は“未書き込み”でしかも完全にフルであるためその全体がポストされ、またログイメージのブロック“２”（ＬＩＢ２）も“未書き込み”であるためこれもその全体がポストされる（ただし、ログデータは事象／時刻Ｔ１まではＬＩＢ２の一部し

が満たしていない)。ディスク 2 が最低使用頻度のディスクとして検出され、ポストがディスクステージングログ 65 の偶数ページ番号から開始されるものと(便宜上)仮定すると、ログイメージページ 57 はそのブロック内容 L I B 1 および L I B 2 をディスクステージングログ 65 におけるディスク D 2 のページ P 2 に対応するブロック B 1 および B 2 (すなわち、D 2 P 2 B 1 および D 2 P 2 B 2) にポストする。したがって、ブロック D 2 P 2 B 1 は(反転ビデオ水平線の形態で示す)すべての有効データを含み、ブロック D 2 P 2 B 2 は時刻 T 1 が示す位置までの部分的有効データを含み、ブロック D 2 P 2 B 2 の残りの部分には(クロスハッチで示す)無効データすなわち“ドント・ケア (don't care)”データが含まれる。

【 0 0 3 4 】

図 4 に示すように、第 2 の事象 / 時刻 T 2 は R A I D 管理システム 16 が R L I 5 5 によるそのデータのポストを再度要求する時点を同定する。この例では、時刻 T 1 および時刻 T 2 の間に記憶されたログイメージページ 57 のデータ(すなわち、“未書き込み”データを)ディスクステージングログ 65 にポストしなければならない(これは、まだページ・フル状態になっていないためである)。(しかし、事象 / 時刻 T 2 以前にログイメージページ 57 がすべてトランザクションデータで満たされている場合、R L I 5 5 はログイメージページ 57 の一部をディスクステージングログ 65 にポストするよりむしろその全体をディスクログ 60 にポストする。)ディスク 1 が最低使用頻度のディスクであり、書き込み入出力がフルブロックサイズでのみ発生すると仮定すると、L I B 2 がすべて D 1 P 1 B 2 に書き込まれる。このとき、奇数ページ P 1 (7 5) が書き込まれ、前述したページ“トグル”(スワッピング)データ保護技術が実行される。無効データ(すなわち、そのブロックサイズ内にあって指定された時刻 T 2 より後のデータ)をここでもクロスハッチで示す。

【 0 0 3 5 】

図 5 には R L I 5 5 に対してページ・フル状態になる前にデータのポストが要求される第 3 の事象 / 時刻 T 3 を示す。この例では、時刻 T 2 と時刻 T 3 との間に記憶されたログイメージページ 57 のデータ(“未書き込み”データを)をポストしなければならない。したがって、この例ではディスク 3 (D 3) が最低使用頻度のディスクであると仮定すると、L I B 2 のすべてが D 2 P 2 B 2 にポストされ、L I B 3 のすべてが D 3 P 2 B 3 にポストされる。この場合も、ページスワッピングを行なうために、“偶数の”ページ P 2 (9 0) がこのとき書き込まれる。

【 0 0 3 6 】

図 6 には R L I 5 5 に対してページ・フル状態になる前にデータのポストが再度要求される第 4 の事象 / 時刻 T 4 を示す。この例では、時刻 T 3 と時刻 T 4 との間に記憶されたログイメージページ 57 の“未書き込み”データをポストしなければならない。したがって、この例ではディスク 1 (D 1) が最低使用頻度のディスクであると仮定すると、L I B 3 のすべてが D 1 P 1 B 3 に強制引出しされる。

【 0 0 3 7 】

図 3 ~ 図 6 からわかるように、最低使用頻度のディスクに対して書き込みが行なわれるだけでなくディスクに対する冗長書き込みが発生しないことからシステムの入出力性能に対する全体的影響が低減される。冗長性はログデータがディスク(ディスクステージングログ 65) 上に書き込まれ、しかも R L I 5 5 にも残ることによって維持される。さらに、一実施形態において、ディスクステージングログ 65 への“強制引出し”は R L I 5 5 がページ・フル状態になる前の事象について発生し、R L I 5 5 のページ・フル状態が検出された場合、R L I 5 5 からディスクログ 60 (図 2) への“排出”が発生することに注意しなければならない。

【 0 0 3 8 】

図 7 には、分散書き込み用のデータ記録 110 の構成を示すブロック図である。R L I 5 5 (およびディスクログ 60 およびディスクステージングログ 65) 内のデータの各 512 バイトブロック(セクタ)はそれぞれ 1 つあるいはそれ以上のデータ記録 110 を有し

10

20

30

40

50

、データ記録 1 1 0 はブロックの境界にまたがっている場合もある。説明の目的上、データ記録 1 1 0 は簡略化した形態で示す。すなわち、記録に用いられる可能性のあるすべてのフィールドをここでは図示しない。しかし、データ記録 1 1 0 は少なくとも記録の長さを同定する長さ標識 1 1 5、ディスクステージングログ 6 5 からデータを回復するための記録の順序付けを同定するためのシーケンス番号 1 2 0、トランザクションログに關係付けられたディスク群を同定するディスク群識別子 1 2 5、記憶される実際のログデータを保持する本体 1 3 0、およびデータ確認用のチェックサム 1 3 5を含む。

【 0 0 3 9 】

シーケンス番号 1 2 0 はトランザクションログに新たな記録が追加されるたびに順次インクリメントされる生成番号である。チェックサム 1 3 5 はデータ記録 1 1 0 全体のチェックサムであり、トランザクションログの回復中に記録の状態を確認するために用いられる。ディスク群識別子 1 2 5 は R L I 5 5 に關係付けられたディスク群の現在のインスタンス (instance) の任意の識別子であり、トランザクションログの回復中に “ 陳腐化した ” (すなわち無効な) ディスクステージングログ 6 5 データが使用されていないことを保証するために用いられる。すなわち、回復中に、記録はそのディスク群識別子 1 2 5 がディスクの現在のインスタンスに一致する場合に有効と認識され、ディスク群識別子 1 2 5 がディスク群の現在のインスタンスに一致しない場合無効と認識される。たとえば、他のディスク群からディスクドライブがスワップされた場合に陳腐化した記録あるいは有効な記録が発生する。この場合、各記録に關係付けられたディスク群識別子によって、トランザクションログ回復処理はその新たなディスクに關係付けられたあらゆる陳腐化データを認識し、その使用を避けることができる。つまり、記録のディスク群識別子は現在のディスク群インスタンスに一致しなければならない。

【 0 0 4 0 】

図 8 は分散書き込みディスクログの処理フローを示すブロック図であり、本発明におけるログトランザクションの管理のために (図 1 の) R A I D 管理システム 1 6 内で実施される処理の相互關係を示す。これらの処理は、好適にはファームウェアで実行される。アプリケーション 1 5 0 (たとえば図 1 の R A I D 管理システム 1 6) が N V R A M マップ 4 5 (図 2) を操作する際、この操作動作を同定するデータ記録 1 1 0 がログ管理 1 5 5 の制御に追加され、R A M ログイメージ 5 5 (図 2) に記録される。記録はいくつかの鍵となる事象の 1 つが発生するまで絶えず追加される。R L I 5 5 中の現在のページがフルになると、ログ管理 1 5 5 がデータ管理 1 6 0 に制御を渡し、データ管理 1 6 0 がディスクドライバ 1 6 5 とインターフェースで連結してこのフルページの内容をディスクアレイ 1 1 のディスクログ 6 0 に冗長ポストすることによってそのフルになったページが “ 排出される ”。R L I 5 5 中の現在のページがフルではないがログ管理 1 5 5 がその事象をディスクアレイ 1 1 のディスクステージング 6 5 への “ 強制引出し ” ポストの要求として検出すると、ログ管理 1 5 5 はデータ管理プログラムを迂回してディスクドライバ 1 6 5 と直接インターフェースで連結してそのデータをディスクアレイに “ 強制引出し ” する。 “ 強制引出し ” ポストの発生後に冗長コピーは書き込まれない。

【 0 0 4 1 】

ここで、すべての図面を参照して、なんらかのメモリ故障あるいはシステム故障のためにログの回復が必要である場合、ディスクログ 6 0 およびディスクステージングログ 6 5 に記憶されたインクリメンタルなログデータの回復のためにいくつかのステップが発生する。まず、ディスクログ 6 0 のすべてのフルページが R L I 5 5 にコピーされてログデータができるだけ多く再構築される。しかし、ページ 5 9 のような非フルページの場合ディスクステージングログ 6 5 にデータが残っている可能性があり、これも R L I 5 5 にコピーしなければならない。したがって、ディスク 1 ~ M 上のディスクステージングログ 6 5 の全ページ 7 0 ~ 1 0 5 からのログデータの断片をまとめて 1 つの完全なイメージにして R L I 5 5 にコピーしなければならない。

【 0 0 4 2 】

ディスクステージングログ 6 5 からのこの回復に備えて、R L I 5 5 を走査してディスク

ログ 60 に書き込まれた最後の（最も新しい）記録を示すシーケンス番号 120 を有する記録を発見する。この走査にあたってはログの循環性とシーケンス番号のラッピングの両方を考慮しなければならない。ディスクステージングログ 65 内に最後に書き込まれたものに続く次のシーケンス番号（記録）がある場合これを発見しなければならない。したがって、ディスクステージングログ 65 を次に走査して（RLI55 に次に回復すべき記録を示す）次のシーケンス番号を有する記録を発見する。ディスクステージングログ 65 内で次の記録（シーケンス番号）が発見されると、そのディスク群識別子 125 をチェックして、その記録がディスク群の現在のインスタンスに属するものであることを確認する。さらに、その記録のチェックサム 135 を評価してその記録の完全性を判定する。完全である場合、その記録が RLI55 にコピーされ、トランザクションログ回復処理が続行される。図 6 に示す例では、ブロック D2P2B1 の最初の記録はこの第 1 ステップの回復基準を満足する。

10

【0043】

続いて、ディスクステージングログ 65 を再度走査して、次の連続シーケンス番号を有する次の記録を探す。前に発見された記録の長さ標識 115 から、次の記録が前に発見された記録の長さ標識 115 によって記述されるオフセットで始まることがわかる。したがって、そのオフセットで始まり、適当な連続シーケンス番号 120 を有し、適当なディスク群識別子 125 を有し、チェックサム 135 による有効性分析を満足する次の記録が発見されるまですべてのディスクステージングログ 65 が探索される。図 6 の例では、ブロック D2P2B1 内の第 2 の記録（視覚的には識別不能）がこの回復基準を満足する。かかる記録が発見されると、その記録が RLI55 にコピーされる。

20

【0044】

一般的には、(i)（次の連続シーケンス番号によって同定される）次の連続する記録を発見し、(ii) そのディスク群識別子を検証し、(iii) そのチェックサムを検証するこの処理全体が、かかる回復基準のそれぞれを満足するすべての記録が発見されるまでディスクステージングログ 65 全体に対して継続的に反復される。さらに例を挙げれば、たとえば図 6 において、D2P2B1 および D2P2B2（反転ビデオ水平線で示される）中に同定される有効な記録のそれぞれがまずその順序で RLI55 に回復される。すると、D1P1B2 あるいは D3P2B2 内で発見されるいかなる有効な記録も次の回復ステップを満足する。たとえば、次の有効な記録がまず D3P2B2 の前に D1P1B2 内で発見された場合、D1P1B2 内のすべての有効な記録が RLI55 に（順次一度に 1 つずつ）コピーされ、D1P1B2 内で発見されなかった残りの有効記録は続いて D3P2B2 内で発見され、D3P2B3 の有効記録に続く。一方、次の有効な記録が（D1P1B2 ではなく）D3P2B2 内で発見された場合、これらの有効記録はすべて RLI55 にコピーされ、次に D3P2B3 の有効記録が RLI55 にコピーされる。その後のみ D1P1B3 内で発見された最終有効記録が処理され、RLI55 にコピーされる。この最終記録はディスクステージングログ 65 内の他のいかなる記録にも次の連続シーケンス番号が発見されない際に、ディスクステージングログ 65 から回復されることは明らかである。

30

【0045】

図示する例では、ディスクログ 60 およびディスクステージングログ 65 から RLI55 へのログ回復がこれで完了する。したがって、図 8 に示すように、ログ管理 155 が（現在は RLI55 にある）回復された記録をアプリケーション 150（図 1 の RAID 管理システム 16）の制御に返して、ログ回復が完了し、アプリケーション 150 が NVRA M マップ 45 をトランザクションログ回復処理を開始させたシステム誤り / 故障の前の状態に戻すために RLI55 に示されるログ変更の実行に着手可能であることを示す。

40

【0046】

ディスクアレイ内の複数のディスクに対する分散書き込み動作を用いてディスクログ書き込み性能を向上させる方法および装置の実施形態を上記に説明した。当業者には当該技術分野のさまざまなソフトウェア、ファームウェアおよび / またはハードウェアのうち任意の

50

ものを用いて容易に実施されることは明らかであろう。さらに、本発明をその具体的実施形態の参照により説明したが、本発明の精神と範囲から逸脱することなく他の代替実施形態および実施方法あるいは変更形態の使用が可能であることは明らかであろう。

【0047】

以下に本発明の実施の形態を要約する。

【0048】

1. 複数の記憶媒体(12)を有する記憶システム(10)への書き込み方法であって、

- (a) 前記記憶システム(10)への書き込み要求を示す基準の検出と、
 - (b) 前記複数の記憶媒体(12)から最低使用頻度の記憶媒体の選択と、
 - (c) 前記選択された最低使用頻度の媒体へのデータの書き込みと
- を含むことを特徴とする記憶システムへのデータ書き込み方法。

10

【0049】

2. 前記複数の記憶媒体(12)はランダムアクセス記憶媒体である上記1に記載の記憶システムへのデータ書き込み方法。

【0050】

3. 前記最低使用頻度の記憶媒体は入出力動作に基づいて選択される上記1または2に記載の記憶システムへのデータ書き込み方法。

【0051】

4. 前記検出される基準は前記記憶媒体(12)へのトランザクションログ(55)の強制引出しポスト要求を含む上記1、2または3に記載の記憶システムへのデータ書き込み方法。

20

【0052】

5. 前記強制引出しポスト要求は前記トランザクションログ(55)の指定された部分(57)がデータで満たされていることが検出される前に発生する上記4に記載の記憶システムへのデータ書き込み方法。

【0053】

6. 前記データは前記選択された最低使用頻度の記憶媒体に非冗長に書き込まれる請求項1に記載の記憶システムへのデータ書き込み方法。

【0054】

7. 前記データは前記書き込まれるデータの順序を示す標識(120)を含む上記1に記載の記憶システムへのデータ書き込み方法。

30

【0055】

8. 前記記憶媒体(12)はそれぞれ前記最低使用頻度の記憶媒体が選択された際の書き込みのためにのみ使用されるよう確保された領域(65)を含む上記1に記載の記憶システムへのデータ書き込み方法。

【0056】

9. 前記確保された領域は少なくとも2つのサブ領域(70、75、80、85、95、100、105)を含み、最低使用頻度の記憶媒体を選択する第1の事象の発生時に前記サブ領域の1つに対して書き込みが行なわれ、最低使用頻度記憶媒体を選択する次の事象の発生時に他方のサブ領域への書き込みが行なわれ、それによって同じ最低使用頻度記憶媒体が二度連続して選択される場合にも次の連続する書き込みにおいて直前に書き込まれたサブ領域に対して重ね書きが生じないことを特徴とする上記8に記載の記憶システムへのデータ書き込み方法。

40

【0057】

10. 記憶システム(10)であって、
(a) データ記録(110)を保持する第1メモリ(55)と、
(b) 前記第1メモリに接続された複数の記憶媒体(12)と、
(c) 前記第1メモリ(55)の状態を検出する手段(16)と、
(d) 前記第1メモリにおける第1の検出された状態に応答して第1の記憶管理基準にし

50

たがって前記記憶システム（１０）に前記データ記録を書き込む手段（１６）と、
（ｅ）前記第１メモリにおける第２の検出された状態に応答して第２の記憶管理基準にしたがって前記記憶システム（１０）に前記データ記録を書き込む手段（１６）を含み、
前記第２の記憶管理基準は前記複数の記憶媒体における最低使用頻度の記憶媒体への書き込みを含む記憶システム。

【００５８】

１１． 前記第１メモリ（５５）における前記第１の検出された状態は、前記第１メモリの一部が前記データ記録で満たされていることを示す状態を含む上記１０記載の記憶システム。

【００５９】

１２． 前記第１の記憶管理基準は前記複数の記憶媒体（１２）上の前記データ記録の冗長性の維持を含む上記１０記載の記憶システム。

【００６０】

１３． 前記第１メモリにおける前記第２の検出された状態は、前記第１メモリにおいて所定の部分が前記データ記録で満たされる前に前記第１メモリ（５５）内の前記データ記録を前記複数の記憶媒体（１２）に強制引出し的に書き込む要求を示す状態を含む上記１０記載の記憶システム。

【００６１】

１４． 前記第２の記憶管理基準は前記複数の記憶媒体（１２）上において前記データ記録の冗長性を維持しないことを含む上記１０記載の記憶システム。

【００６２】

１５． 前記最低使用頻度の記憶媒体は前記データ記録の書き込みのために確保された少なくとも２つのサブ領域（７０、７５、８０、８５、９０、９５、１００、１０５）を含む上記１０記載の記憶システム。

【００６３】

【発明の効果】

本発明によれば、複数の利用可能な記憶ディスクから選択された任意の最も使用頻度の低いディスクへのログ書き込みを管理および分散してログ入出力と進行中の他の入出力との間におけるディスクアクセスの競合を低減することによって、ディスクログ書き込みのシステム性能を改善でき、且つ、システムの誤り又は故障時に記録を確実に回復することができる。

【００６４】

【図面の簡単な説明】

【図１】本発明の分散書き込みディスクログ法を用いたデータ記憶システムのブロック図である。

【図２】本発明の分散ログ書き込みディスクログ法を示すブロック図である。

【図３】ある時間における本発明のディスクステージングログの状態を示すブロック図である。

【図４】図３の状態以後の、他の時間における本発明のディスクステージングログの状態を示すブロック図である。

【図５】図４の状態以後の、他の時間における本発明のディスクステージングログの状態を示すブロック図である。

【図６】図５の状態以後の、他の時間における本発明のディスクステージングログの状態を示すブロック図である。

【図７】分散書き込み用のデータ記録の構成を示すブロック図である。

【図８】分散書き込みディスクログの処理フローを示すブロック図である。

【符号の説明】

- １，２，３，Ｍ ディスク
- １０ データ記憶システム
- １１ ディスクアレイ

10

20

30

40

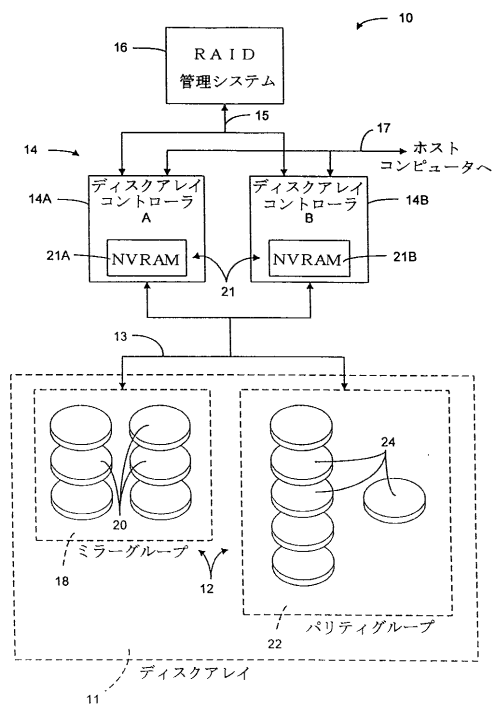
50

1 2	記憶ディスク
1 4	ディスクアレイコントローラ
1 4 A	ディスクアレイコントローラ A
1 4 B	ディスクアレイコントローラ B
1 6	R A I D 管理システム
1 8	ミラーグループ
2 1	メモリマップ記憶域
2 1 A , 2 1 B	N V R A M
2 2	パリティグループ
4 5	N V R A M マップ
5 0	ディスクマップ
5 5	R A M ログ イメージ (R L I)
6 0	ディスクログ
6 5	ディスクステージングログ
1 1 0	データ記録
1 1 5	長さ標識
1 2 0	シーケンス番号
1 2 5	ディスク群識別子
1 3 0	本体
1 3 5	チェックサム
1 5 0	アプリケーション
1 5 5	ログ管理
1 6 0	データ管理
1 6 5	ディスクドライバ

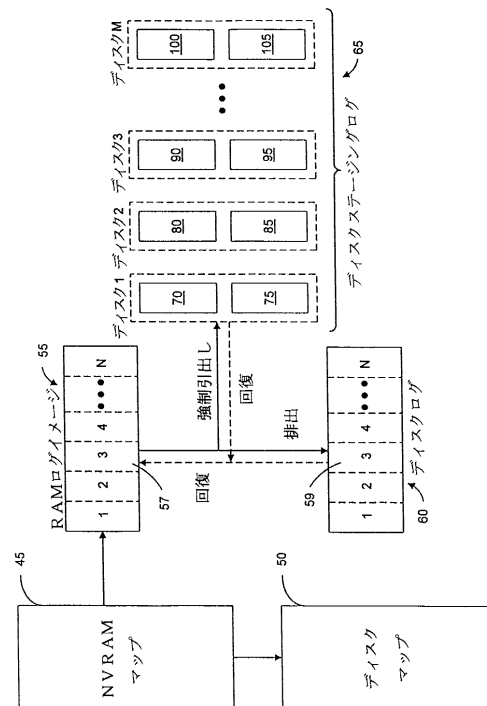
10

20

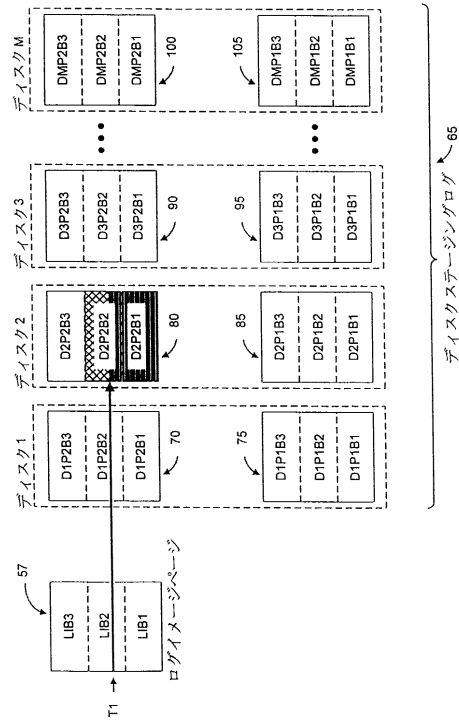
【図 1】



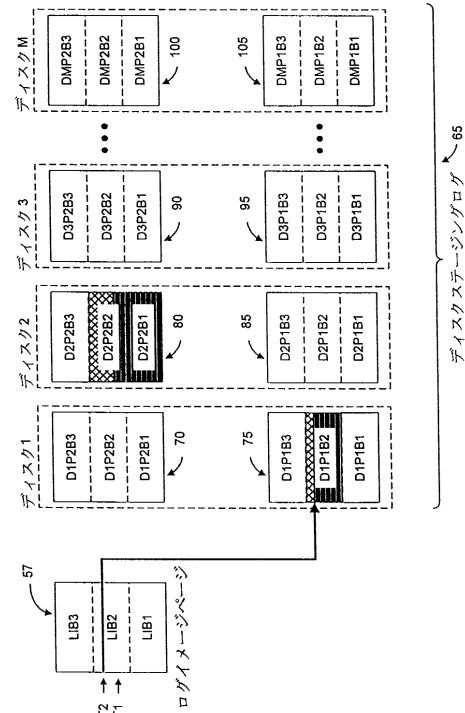
【図 2】



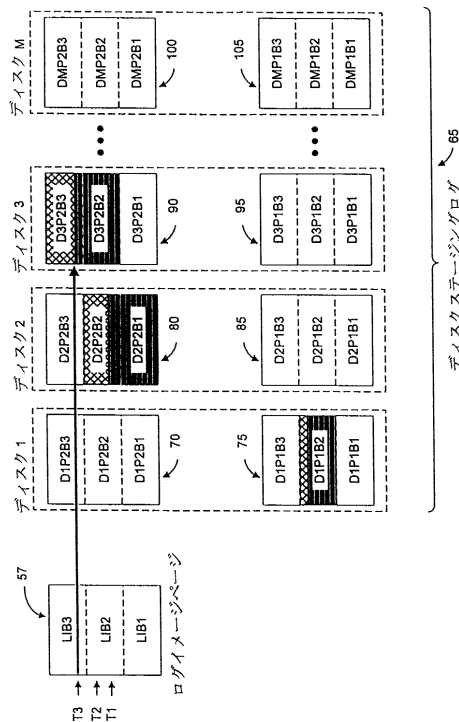
【図 3】



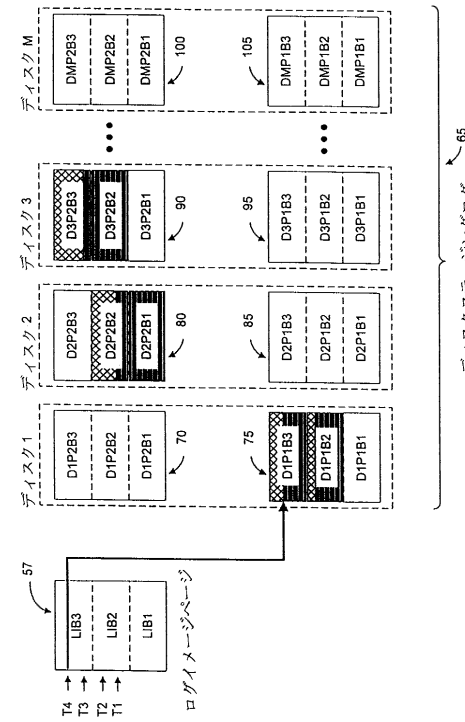
【図 4】



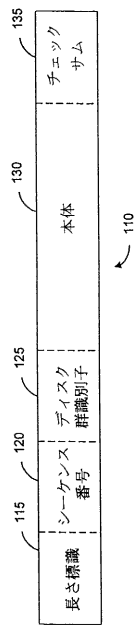
【図 5】



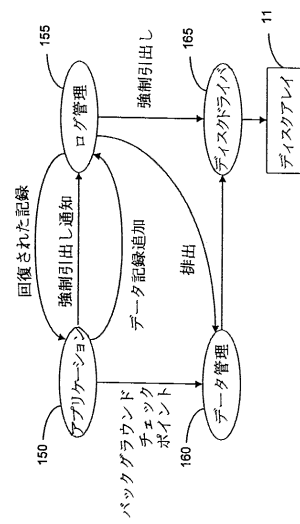
【図 6】



【図 7】



【図 8】



フロントページの続き

(72)発明者 ドン・エル・ブルケス

アメリカ合衆国 アイダホ，メリディアン，サン・ラモ・ディーアール 3100

(72)発明者 キルク・エー・ハンソン

アメリカ合衆国 アイダホ，イーグル，ウエスト・ニューフィールド・ディーアール 1129

審査官 吉 田 美彦

(56)参考文献 実開平04 - 073251 (JP, U)

特開平09 - 054658 (JP, A)

特開平05 - 081099 (JP, A)

特開平06 - 324817 (JP, A)

特開平06 - 231012 (JP, A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06