

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6206161号  
(P6206161)

(45) 発行日 平成29年10月4日 (2017. 10. 4)

(24) 登録日 平成29年9月15日 (2017. 9. 15)

(51) Int. Cl.

F I

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

G 0 6 F 3 / 0 8 (2006. 01)

G 0 6 F 1 2 / 1 6 (2006. 01)

G 0 6 F 3 / 0 6 3 0 1 A

G 0 6 F 3 / 0 8 H

G 0 6 F 3 / 0 6 3 0 4 Z

G 0 6 F 3 / 0 6 3 0 1 N

G 0 6 F 1 2 / 1 6 3 1 0 A

請求項の数 9 (全 38 頁)

(21) 出願番号 特願2013-261843 (P2013-261843)  
 (22) 出願日 平成25年12月18日 (2013. 12. 18)  
 (65) 公開番号 特開2015-118572 (P2015-118572A)  
 (43) 公開日 平成27年6月25日 (2015. 6. 25)  
 審査請求日 平成28年9月5日 (2016. 9. 5)

(73) 特許権者 000005223  
 富士通株式会社  
 神奈川県川崎市中原区上小田中4丁目1番  
 1号  
 (74) 代理人 100104190  
 弁理士 酒井 昭徳  
 (72) 発明者 斎藤 博紀  
 神奈川県川崎市中原区上小田中4丁目1番  
 1号 富士通株式会社内  
 審査官 田中 啓介

最終頁に続く

(54) 【発明の名称】 ストレージ制御装置、制御方法、および制御プログラム

(57) 【特許請求の範囲】

【請求項 1】

ストレージ装置の記憶領域に対して所定のデータサイズのブロック単位でデータを冗長化して記憶させる制御を行うストレージ制御装置であって、

データの書込要求から特定される書込先の論理アドレスを、前記ブロック単位に対応付けてグループ分けしたグループ情報を記憶する記憶部と、

受け付けたデータの読出要求に応じて、前記グループ情報に基づき前記読出要求から特定される読出先の論理アドレスを含むグループの読出回数を計数し、計数した前記グループの読出回数に基づき前記読出先の論理アドレスを含む読出要求、または、前記読出先の論理アドレスのデータに対応する冗長データの記憶先の論理アドレスを含む読出要求を発行する制御部と、

を有し、

前記記憶領域は、規定回数の読み出しがブロックに対して行われると該ブロックのデータを他のブロックへコピーするコピー制御が行われる半導体メモリの記憶領域であり、

前記制御部は、

計数した前記グループの読出回数が前記規定回数より小さい第1の回数となったことに  
 応じて、前記グループに含まれる論理アドレスのデータを対象とした前記コピー制御を起  
 動させるように、前記グループに含まれる論理アドレスを含む読出要求を発行する処理を  
 行い、

前記処理を行っているときに、データの読出要求を受け付けたことに応じて、当該読出

先の論理アドレスのデータに対応する冗長データの記憶先の論理アドレスを含む読出要求を発行する、

ことを特徴とするストレージ制御装置。

【請求項 2】

前記ストレージ制御装置は、前記半導体メモリを含む前記ストレージ装置と、当該ストレージ装置が記憶するデータの冗長データを記憶する前記ストレージ装置とは異なる他のストレージ装置の制御を行うものであり、

前記制御部は、

計数した前記グループの読出回数が前記第 1 の回数を超えた後のデータの読出要求に応じて、当該読出要求から特定される読出先の論理アドレスが前記グループに含まれる場合に、当該読出先の論理アドレスのデータに対応する冗長データの記憶先の論理アドレスを含む読出要求を、前記他のストレージ装置に発行する、

ことを特徴とする請求項 1 に記載のストレージ制御装置。

【請求項 3】

前記制御部は、

計数した前記グループの読出回数が前記第 1 の回数となったことに応じて、前記グループに含まれる論理アドレスのデータを対象とした前記コピー制御を起動させるように、前記グループに含まれる論理アドレスを含む読出要求を、前記規定回数から前記第 1 の回数を減じた第 2 の回数、または、当該読出要求に対する応答時間が所定時間を超えるまで、前記ストレージ装置に発行し、前記グループの読出回数を初期化することを特徴とする請求項 1 または 2 に記載のストレージ制御装置。

【請求項 4】

前記制御部は、

前記グループに含まれる論理アドレスを含む読出要求を、前記第 2 の回数、または、当該読出要求に対する応答時間が所定時間を超えるまで発行したことに応じて、前記グループに含まれる論理アドレスを含む読出要求の発行を停止するとともに、前記グループの読出回数を初期化することを特徴とする請求項 3 に記載のストレージ制御装置。

【請求項 5】

前記制御部は、

前記グループに含まれる論理アドレスを含む書込要求を前記ストレージ装置に発行したことに応じて、前記グループの読出回数を初期化することを特徴とする請求項 1 ～ 4 のいずれか一つに記載のストレージ制御装置。

【請求項 6】

前記ストレージ制御装置は、前記半導体メモリを含む前記ストレージ装置と、当該ストレージ装置が記憶するデータの冗長データを記憶する前記ストレージ装置とは異なる他のストレージ装置の制御を行うものであり、

前記記憶部は、

書込要求から特定される前記他のストレージ装置の書込先の論理アドレスを、前記ブロック単位に対応付けてグループ分けした前記他のストレージ装置のグループ情報を記憶しており、

前記制御部は、

受け付けたデータの読出要求に応じて、前記他のストレージ装置のグループ情報に基づき当該読出要求から特定される前記他のストレージ装置の読出先の論理アドレスを含むグループの読出回数を計数し、計数した当該グループの読出回数が前記規定回数より小さく前記第 1 の回数とは異なる第 3 の回数となったことに応じて、当該グループに含まれる論理アドレスを含む読出要求を、前記規定回数から前記第 3 の回数を減じた第 4 の回数、または、当該読出要求に対する応答時間が所定時間を超えるまで前記他のストレージ装置に発行し、当該グループの読出回数を初期化することを特徴とする請求項 3 に記載のストレージ制御装置。

【請求項 7】

前記ストレージ制御装置は、さらに、半導体メモリの記憶領域を所定のデータサイズで分割したブロック単位で書き込みを行う前記ストレージ装置とは異なる移行先のストレージ装置の制御を行うものであり、

前記制御部は、

前記ストレージ装置によりガベージコレクションが実行された場合に、受け付けた書込要求から特定される前記ストレージ装置の書込先の論理アドレスをグループ分けした複数のグループに含まれる論理アドレスを含むデータ移行要求を、前記ストレージ装置から前記移行先のストレージ装置への移行対象データの移行間隔が所定間隔を超えないように前記ストレージ装置に発行し、

発行した前記データ移行要求に含まれる論理アドレスを、前記ブロック単位に対応付けてグループ分けした前記移行先のストレージ装置のグループ情報を生成し、

前記データ移行要求に含まれる論理アドレスを含む読出要求の発行先を前記移行先のストレージ装置へ切り替える、

ことを特徴とする請求項 1 ～ 6 のいずれか一つに記載のストレージ制御装置。

【請求項 8】

ストレージ装置の記憶領域に対して所定のデータサイズのブロック単位でデータを冗長化して記憶させる制御を行うストレージ制御装置の制御方法であって、

前記記憶領域は、規定回数の読み出しがブロックに対して行われると該ブロックのデータを他のブロックへコピーするコピー制御が行われる半導体メモリの記憶領域であり、

前記ストレージ制御装置が、

受け付けたデータの読出要求に応じて、データの書込要求から特定される書込先の論理アドレスを、前記ブロック単位に対応付けてグループ分けしたグループ情報に基づき前記読出要求から特定される読出先の論理アドレスを含むグループの読出回数を計数し、計数した前記グループの読出回数に基づき前記読出先の論理アドレスを含む読出要求、または、前記読出先の論理アドレスのデータに対応する冗長データの記憶先の論理アドレスを含む読出要求を発行し、

計数した前記グループの読出回数が前記規定回数より小さい第 1 の回数となったことに応じて、前記グループに含まれる論理アドレスのデータを対象とした前記コピー制御を起動させるように、前記グループに含まれる論理アドレスを含む読出要求を発行する処理を行い、

前記処理を行っているときに、データの読出要求を受け付けたことに応じて、当該読出先の論理アドレスのデータに対応する冗長データの記憶先の論理アドレスを含む読出要求を発行する、

処理を実行することを特徴とする制御方法。

【請求項 9】

ストレージ装置の記憶領域に対して所定のデータサイズのブロック単位でデータを冗長化して記憶させる制御を行うストレージ制御装置の制御プログラムであって、

前記記憶領域は、規定回数の読み出しがブロックに対して行われると該ブロックのデータを他のブロックへコピーするコピー制御が行われる半導体メモリの記憶領域であり、

前記ストレージ制御装置に、

受け付けたデータの読出要求に応じて、データの書込要求から特定される書込先の論理アドレスを、前記ブロック単位に対応付けてグループ分けしたグループ情報に基づき前記読出要求から特定される読出先の論理アドレスを含むグループの読出回数を計数し、計数した前記グループの読出回数に基づき前記読出先の論理アドレスを含む読出要求、または、前記読出先の論理アドレスのデータに対応する冗長データの記憶先の論理アドレスを含む読出要求を発行し、

計数した前記グループの読出回数が前記規定回数より小さい第 1 の回数となったことに応じて、前記グループに含まれる論理アドレスのデータを対象とした前記コピー制御を起動させるように、前記グループに含まれる論理アドレスを含む読出要求を発行する処理を行い、

10

20

30

40

50

前記処理を行っているときに、データの読出要求を受け付けたことに応じて、当該読出先の論理アドレスのデータに対応する冗長データの記憶先の論理アドレスを含む読出要求を発行する、

処理を実行させることを特徴とする制御プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ストレージ制御装置、制御方法、および制御プログラムに関する。

【背景技術】

【0002】

従来、不揮発性のメモリとして、フラッシュメモリを採用したストレージ装置がある。フラッシュメモリでは、読み出しを行う際に浮遊ゲートに電圧をかけるため浮遊ゲートの電子量が増加する。このため、あるブロックに対して読み出しが繰り返されることにより、電子の増加量が大きくなってビットエラーが発生する、いわゆるリードディスタ urb が起きる。リードディスタ urb を防止するため、フラッシュメモリを含むストレージ装置内部では、読出回数を管理して、一定回数の読み出しが行われるとブロックのデータを別ブロックへコピーする処理を行う。

【0003】

関連する先行技術として、たとえば、フラッシュメモリの読出頻度に対応してリードをフラッシュメモリから行うか RAM ( Random Access Memory ) から行うかを制御し、所定のエリアのデータのリード頻度が一定の値を超えたとき、フラッシュメモリのデータを RAM に転送するものがある。また、ホスト PC ( Personal Computer ) からのアクセスコマンドに基づいて、リードディスタ urb によりデータが読み出される回数に制限があるフラッシュメモリからのデータ読み出し、および、フラッシュメモリをリフレッシュするための読み出し回数の更新を制御する技術がある。さらに、データ復旧処理において、フラッシュメモリから所定の時間内にデータ読出処理が完了しない場合に、データ読出処理を中止して付加データを設定し、パリティデータを用いて誤りデータを訂正した退避対象データをキャッシュメモリに書き込む技術がある。(たとえば、下記特許文献 1 ~ 3 を参照。)

【先行技術文献】

【特許文献】

【0004】

【特許文献 1】特開 2001 - 290791 号公報

【特許文献 2】特開 2008 - 181380 号公報

【特許文献 3】国際公開第 2009 / 107213 号

【発明の概要】

【発明が解決しようとする課題】

【0005】

しかしながら、従来技術によれば、フラッシュメモリを含むストレージ装置内部で行われるリードディスタ urb を防止する処理により、ストレージ装置への読出要求に対する応答性能が低下することがある。

【0006】

1 つの側面では、本発明は、ストレージ装置への読出要求に対する応答性能の低下を抑制するストレージ制御装置、制御方法、および制御プログラムを提供することを目的とする。

【課題を解決するための手段】

【0007】

本発明の一側面によれば、データが繰り返し読み出されることによりビットエラーが発生する特性を有する半導体メモリの記憶領域を所定のデータサイズで分割したブロック単位で書き込みを行うストレージ装置と接続されるストレージ制御装置であって、ストレ

10

20

30

40

50

ジ装置への書込要求に応じて、書込要求に含まれる書込先の論理アドレスを、同一グループに含まれる論理アドレスに対応するデータサイズが所定のデータサイズを超えないようにグループ分けし、ストレージ装置への読出要求に応じて、読出要求に含まれる読出先の論理アドレスを含むグループの読出回数を計数し、計数したグループの読出回数に基づきグループに含まれる論理アドレスの読み出しを制御するストレージ制御装置、制御方法、および制御プログラムが提案される。

【発明の効果】

【0008】

本発明の一態様によれば、ストレージ装置への読出要求に対する応答性能の低下を抑制することができるという効果を奏する。

10

【図面の簡単な説明】

【0009】

【図1】図1は、本実施の形態にかかるストレージ制御装置の動作例を示す説明図である。

【図2】図2は、ストレージシステムの接続例を示す説明図である。

【図3】図3は、ストレージ制御装置のハードウェア構成例を示すブロック図である。

【図4】図4は、ストレージ制御装置の機能構成例を示すブロック図である。

【図5】図5は、ブロック予測テーブルの記憶内容の一例を示す説明図である。

【図6】図6は、書込時における動作例を示す説明図である。

【図7】図7は、読出時における動作例を示す説明図である。

20

【図8】図8は、上書き書込時における動作例を示す説明図である。

【図9】図9は、RD防止コピー閾値の設定例を示す説明図（その1）である。

【図10】図10は、RD防止コピー閾値の設定例を示す説明図（その2）である。

【図11】図11は、ガベージコレクションが行われた際のブロックとグループとの乖離の一例を示す説明図である。

【図12】図12は、リフレッシュ処理の実行前後におけるブロックとグループとの関係の一例を示す説明図である。

【図13】図13は、ストレージ装置制御処理手順の一例を示すフローチャートである。

【図14】図14は、書込時処理手順の一例を示すフローチャート（その1）である。

【図15】図15は、書込時処理手順の一例を示すフローチャート（その2）である。

30

【図16】図16は、RD防止コピー引起し処理手順の一例を示すフローチャートである。

【図17】図17は、リフレッシュ処理手順の一例を示すフローチャートである。

【発明を実施するための形態】

【0010】

以下に図面を参照して、開示のストレージ制御装置、制御方法、および制御プログラムの実施の形態を詳細に説明する。

【0011】

図1は、本実施の形態にかかるストレージ制御装置の動作例を示す説明図である。図1に示すストレージシステム100は、ストレージシステム100の利用者に対して、ストレージシステム100が有する記憶領域を提供するシステムである。たとえば、ストレージシステム100は、ストレージシステム100の利用者が作成した文書ファイルやデータベースを記憶する。

40

【0012】

ストレージシステム100は、ストレージシステム100の全体を制御するストレージ制御装置101と、ストレージシステム100が提供する記憶領域を有するストレージ装置102-1、2、...、nとを有する。nは1以上の整数である。

【0013】

また、ストレージシステム100が提供する記憶領域への読出要求に対する応答時間を短縮するために、ストレージ装置102は、SSD (Solid State Driv

50

e) が適用される。以下、単に、ストレージ装置 102 と記載した場合、SSD が適用されたストレージ装置であるとする。また、本実施の形態では、ストレージ装置 102 - 1 ~ n のうち、少なくとも一つは SSD が適用されたストレージ装置である。図 1 の例では、ストレージ装置 102 - 1 が、SSD が適用されたストレージ装置であるとする。また、ストレージシステム 100 は、RAID (Redundant Arrays of Inexpensive Disks) 技術を用いて形成された仮想的なボリュームを利用者に提供する。ストレージ制御装置 101 は、RAID 技術を用いて、ストレージ装置 102 の記憶領域に対して所定のデータサイズのブロック単位でデータを冗長化して記憶させる制御を行う。RAID については、図 2 で後述する。

#### 【0014】

SSD は、記憶媒体としてフラッシュメモリを用いるドライブ装置である。具体的に、SSD は、フラッシュメモリと、受け付けた書込データやフラッシュメモリの記憶内容を一時的に記憶するキャッシュメモリと、SSD コントローラと、を有する。

#### 【0015】

フラッシュメモリとは、記憶領域を所定のデータサイズで分割したブロック単位で消去、書き込みを行うことができ、電源を切ってもデータを失わない性質を有する半導体メモリである。所定の管理データサイズ（以下、「ブロックデータサイズ」と呼称する）は、フラッシュメモリの製造時に、フラッシュメモリの仕様としてフラッシュメモリの設計者により設定される値である。

#### 【0016】

具体的に、フラッシュメモリは、1 ビットの情報を蓄積するメモリセルを複数有する。メモリセルは、シリコン基板上の P 型半導体層を挟み込むようにソースとドレインとなる 2 つの N 型半導体部分を有し、P 型半導体層の上に絶縁膜を介して浮遊ゲートを配置し、浮遊ゲート上に制御ゲートを配置した構造である。そして、メモリセルは、浮遊ゲートに蓄えられた電子量を用いて情報を記憶する。具体的に、メモリセルは、浮遊ゲートの電子量に応じてソース - ドレイン間に電流が流れるか否かが変化することを用いて、情報を記憶する。電子量に応じて電流が流れるか否かが変化する理由として、浮遊ゲートに電子があればソース - ドレイン間の抵抗が高くなりゲート電圧を高くしなければ電流が流れない。一方、浮遊ゲートに電子がなければソース - ドレイン間の抵抗が低くなり、ゲート電圧が低くても電流が流れる。ソース - ドレイン間の電流を流すために要求されるゲートへの印加電圧は、閾値電圧と呼称される。また、フラッシュメモリには、NOR 型のフラッシュメモリと NAND 型のフラッシュメモリとがある。

#### 【0017】

SSD コントローラは、フラッシュメモリおよびキャッシュメモリを制御する。たとえば、SSD コントローラは、受け付けた論理アドレスに対する書込データを一旦キャッシュメモリに格納する。そして、SSD コントローラは、キャッシュメモリに格納した書込データをブロックデータサイズで分割して、分割したデータをブロックに順次格納するとともに、格納したブロックと論理アドレスとの対応関係を示す情報を、フラッシュメモリに格納する。書込データと、書込先のブロックとの関係については、SSD の仕様によって様々であり、たとえば、SSD への書込データの発行間隔が、書込データと書込先のブロックとの関係に影響する場合もある。ここで、論理アドレスを指定する方法として、SSD においては、LBA (Logical Block Addressing) が採用される。以下、論理アドレスを、単に、LBA と称する。SSD の内部動作については、図 6 ~ 図 8 に後述する。

#### 【0018】

ここで、ストレージ制御装置 101 がストレージ装置 102 に発行する書込要求の LBA は、利用者が操作する装置による仮想ボリュームの書込要求の LBA から特定されるものである。同様に、ストレージ制御装置 101 がストレージ装置 102 に発行する読出要求の LBA は、利用者が操作する装置による仮想ボリュームの読出要求の LBA から特定されるものである。たとえば、ストレージ制御装置 101 が、ストレージ装置 102 - 1

10

20

30

40

50

とストレージ装置 102 - 2 とにより、RAID0 の仮想ボリュームを形成するとする。このとき、ストレージ制御装置 101 は、ストレージ装置 102 - 1 の LBA0、ストレージ装置 102 - 2 の LBA0、ストレージ装置 102 - 1 の LBA1、ストレージ装置 102 - 2 の LBA1、...、という順に仮想ボリュームの LBA0 ~ 3、...を形成する。そして、ストレージ制御装置 101 は、利用者が操作する装置による仮想ボリュームの書込要求を受け付けた場合、書込要求の LBA を 2 で除した際の商と余りとを用いて、書込要求の LBA と発行先のストレージ装置 102 とを特定する。

#### 【0019】

読出要求からデータ受領完了までにかかる時間は、HDD (Hard Disk Drive) が 10 ミリ秒程度であるのに対して、SSD が数ミリ秒程度である。したがって、SSD が適用されたストレージ装置は、HDD が適用されたストレージ装置より読出要求に対応する応答時間を短縮することができる。

10

#### 【0020】

しかしながら、SSD に読出要求を続けると、読出要求に対する応答が遅延する場合がある。読出要求に対する応答が遅延する場合として、ガベージコレクションの動作のタイミングやリードディスタurb (Read Disturb) によるビットエラーを防止させるため実施される動作のタイミングと、読出要求のタイミングとが重なった場合が挙げられる。

#### 【0021】

ここで、ガベージコレクションとは、小ブロックを纏まったページに集めてデータ記録域を確保する動作である。また、リードディスタurb は、データが繰り返し読み出されることによりビットエラーが発生させる。リードディスタurb は、電子を蓄える半導体メモリが有する特性である。リードディスタurb は、NAND 型のフラッシュメモリで発生しやすい。リードディスタurb によるビットエラーは、書き込んだ値が正しく読み出せなかったことによるエラーである。

20

#### 【0022】

リードディスタurb の現象をより詳細に説明する。SSD 内部の記録媒体に用いられるフラッシュメモリの書込単位となるブロックに書き込みを行った時点では、浮遊ゲートに電子がない場合、閾値電圧は低く、読み出し時に印加する電圧とのマージンが十分にある状態である。しかしながら、ブロックのいずれかのセルに読出要求を繰り返し行くと、同一ブロック内の他のセルにも電圧が印加されて浮遊ゲートに保持した電子の量が少しずつ変化して、閾値電圧と読み出し時に印加する電圧とのマージンが減っていき、ビットエラーとなる。したがって、SSD は、各ブロックの読出回数を管理しておき、規定回数の読み出しが行われるとブロックのデータを別ブロックへコピーするように制御を行う。以降、上述の制御を「RD 防止コピー」と呼称する。

30

#### 【0023】

ここで、RD 防止コピーは、SSD 内部の動作であるから、SSD にアクセスする装置は、RD 防止コピーがいつ行われるかという情報を得ることが難しい。このため、SSD が RD 防止コピーの実行中に、SSD にアクセスする装置が RD 防止コピーの対象であるブロックに読出要求を行ってしまう可能性がある。この場合、SSD は、RD 防止コピーが完了した後にデータ読み出しと完了応答を行うため、SSD にアクセスする装置は RD 防止コピーが完了するまで待たされることになる。特に、RD 防止コピーとガベージコレクションとが重なる場合には、SSD にアクセスする装置は大きく待たされることになり、たとえば、1 秒以上完了応答が得られない場合がある。

40

#### 【0024】

そこで、ストレージ制御装置 101 は、ストレージ装置 102 に発行する書込要求の LBA をブロック単位でグループ分けしたグループの読出回数を計数する。そして、ストレージ制御装置 101 は、計数したグループの読出回数に応じて、受け付けた読出要求の LBA から特定される LBA を含む読出要求、または特定される LBA のデータに対応する冗長データの記憶先の LBA を含む読出要求を発行する。これにより、ストレージ制御装

50

置 1 0 1 は、R D 防止コピーが発生するであろう L B A を予測して、応答遅延が発生するであろう L B A に対する読出要求を避けることができ、読出要求に対する応答性能の低下を抑制することができる。

【 0 0 2 5 】

以下に、書込要求の L B A をブロック単位でグループ分けしたグループの読出回数を計数し、計数したグループの読出回数に応じた読出要求の制御についての具体的な説明を行う。

【 0 0 2 6 】

上述したように、S S D は、一連の書込データをブロックデータサイズで分割して、分割したデータを順次ブロックに格納する。したがって、ストレージ制御装置 1 0 1 は、ストレージ制御装置 1 0 1 が受け付けた書込要求から特定される書込先の L B A を、同一グループに含まれる L B A に対応するデータサイズがブロックデータサイズを超えないように（ブロックデータサイズ以下で）グループ分けしたグループ情報を記憶する。ストレージ制御装置 1 0 1 は、1 つのストレージ装置 1 0 2 に対するグループ情報を、ブロック予測テーブル 1 1 1 に登録する。ストレージ制御装置 1 0 1 は、ストレージ装置 1 0 2 に書込要求を発行する前にグループ分けしてグループ情報を生成してもよいし、ストレージ装置 1 0 2 に書込要求を発行した後にグループ分けしてグループ情報を生成してもよい。また、ストレージ制御装置 1 0 1 は、ストレージ装置 1 0 2 に対して既に発行された複数の書込要求に対して、グループ情報を生成してもよい。

【 0 0 2 7 】

ストレージ制御装置 1 0 1 は、各ストレージ装置 1 0 2 に対応するブロック予測テーブル 1 1 1 を管理する。図 1 の例では、ストレージ制御装置 1 0 1 は、ストレージ装置 1 0 2 - 1 に対応するブロック予測テーブル 1 1 1 - 1 を管理する。

【 0 0 2 8 】

図 1 の例では、L B A の 1 0 0 0 個分のデータサイズが、ブロックデータサイズと一致するものとし、ストレージ制御装置 1 0 1 が、L B A 0 - 1 0 0 7 への書込要求をストレージ装置 1 0 2 - 1 に発行するとする。ストレージ制御装置 1 0 1 は、L B A 0 - 9 9 9 をグループ G 0 にグループ分けし、L B A 1 0 0 0 - 1 0 0 7 をグループ G 1 にグループ分けしたグループ情報を生成する。また、ストレージ装置 1 0 2 - 1 は、ストレージ制御装置 1 0 1 からの L B A 0 - 1 0 0 7 への書込要求を受け付けて、L B A 0 - 9 9 9 をブロック B 0 に割り当て、L B A 1 0 0 0 - 1 0 0 7 をブロック B 1 に割り当てたとする。以下の説明では、説明の簡略化のため、グループ G x を、単に「G x」と記載することがある。同様に、ブロック B x を、単に「B x」と記載することがある。x は 0 以上の整数である。

【 0 0 2 9 】

次に、ストレージ制御装置 1 0 1 は、ストレージ制御装置 1 0 1 が受け付けた読出要求 1 1 2 に応じて、ブロック予測テーブル 1 1 1 - 1 に基づき読出要求 1 1 2 から特定される読出先の L B A を含むグループの読出回数を計数する。計数する契機として、ストレージ制御装置 1 0 1 は、読出要求 1 1 2 から読出先の L B A を特定した際に計数してもよいし、読出先の L B A を含む読出要求をストレージ装置 1 0 2 - 1 に発行する前、または発行した後に計数してもよい。

【 0 0 3 0 】

図 1 は、ストレージ装置 1 0 2 - 1 の L B A 0 - 9 9 9 を読出先とする読出要求 1 1 2 が 7 9 0 0 回、ストレージ装置 1 0 2 - 1 の L B A 1 0 0 0 - 1 0 0 7 を読出先とする読出要求が 1 0 2 回の例を示す。このとき、ストレージ制御装置 1 0 1 は、G 0 の読出回数を 7 9 0 0 [ 回 ] と計数し、G 1 の読出回数を 1 0 2 [ 回 ] と計数する。また、ストレージ装置 1 0 2 - 1 は、ストレージ制御装置 1 0 1 から、読出先の L B A 0 - 9 9 9 を含む読出要求を 7 9 0 0 回受け付け、読出先の L B A 1 0 0 0 - 1 0 0 7 を含む読出要求を 1 0 2 回受け付ける。そして、ストレージ装置 1 0 2 - 1 は、B 0 の読出回数を 7 9 0 0 [ 回 ] とし、B 1 の読出回数を 1 0 2 [ 回 ] とする。



## 【 0 0 3 1 】

続けて、ストレージ制御装置 1 0 1 は、計数したグループの読出回数に基づき、読出要求 1 1 2 から特定される読出先の L B A を含む読出要求 1 1 3、または、読出先の L B A のデータに対応する冗長データの記憶先の L B A を含む読出要求 1 1 4 を発行する。図 1 中、点線で括られた読出要求 1 1 3 は、実際には発行していないことを示す。図 1 の例では、ストレージ制御装置 1 0 1 は、読出要求 1 1 3 を発行せず、読出要求 1 1 4 を発行する。

## 【 0 0 3 2 】

ここで、読出要求 1 1 3、1 1 4 について説明する。ここで、説明を容易にするため、読出要求 1 1 3 に含まれる読出先の L B A のデータを、「元データ」と呼称する。

10

## 【 0 0 3 3 】

冗長データは、たとえば、読出要求 1 1 3 に含まれる読出先の L B A が示す記憶領域に記憶されたデータと同一内容のデータである。また、冗長データは、冗長データを加工することにより、元データと同一内容のデータが得られるデータでもよい。

## 【 0 0 3 4 】

また、読出要求 1 1 4 の発行先について、ストレージ制御装置 1 0 1 は、R A I D により、読出要求 1 1 2 に含まれる L B A、または読出要求 1 1 3 に含まれる L B A から、読出要求 1 1 4 の発行先が特定できるように管理してある。図 1 では、R A I D 1 の例を示す。図 1 に示す d 1 - s r c、d 2、d 1 - r d n は、図 2 における R A I D の説明時に併せて説明を行う。

20

## 【 0 0 3 5 】

次に、計数したグループの読出回数に基づき、読出要求 1 1 3、または、読出要求 1 1 4 を発行する例について説明する。たとえば、ストレージ制御装置 1 0 1 は、ストレージ装置 1 0 2 - 1 の仕様に応じて設定された規定回数より小さい第 1 の回数を記憶しておく。第 1 の回数は、ストレージシステム 1 0 0 の管理者によって設定される。そして、ストレージ制御装置 1 0 1 は、計数したグループの読出回数が第 1 の回数になるまでは、読出要求 1 1 3 を発行する。また、ストレージ制御装置 1 0 1 は、計数したグループの読出回数が第 1 の回数になった場合、読出要求 1 1 4 を発行する。以下、グループの読出回数と比較する値を、「R D 防止コピー閾値」と呼称する。R D 防止コピー閾値は、各ストレージ装置 1 0 2 に対してそれぞれで異なる値を設定してもよい。

30

## 【 0 0 3 6 】

読出要求 1 1 3、または、読出要求 1 1 4 を発行した後、ストレージ制御装置 1 0 1 は、発行先からの応答に含まれるデータに基づき読出データを取得し、利用者に読出データを通知する。たとえば、読出要求 1 1 3 を発行した場合、ストレージ制御装置 1 0 1 は、発行先からの応答に含まれるデータと同一内容のデータを利用者が操作する装置に通知する。また、読出要求 1 1 4 を発行した場合、ストレージ制御装置 1 0 1 は、R A I D 技術を用いて、発行先からの応答に含まれる冗長データを用いて元データの内容を復元して、復元したデータを利用者が操作する装置に通知する。図 1 の例では、ストレージ制御装置 1 0 1 は、データ d 1 - r d n と同一内容のデータを、利用者が操作する装置に通知する。R A I D 技術による読出データの復元については、図 9 および図 1 0 で後述する。次に、図 2 を用いて、ストレージシステム 1 0 0 の接続例を示す。

40

## 【 0 0 3 7 】

図 2 は、ストレージシステムの接続例を示す説明図である。ストレージシステム 1 0 0 は、ホストサーバ 2 0 1 と、ストレージ制御装置 1 0 1 と、ストレージ装置 1 0 2 - 1、2、...、n と、ホットスペア 2 0 2 とを有する。ホストサーバ 2 0 1 は、ストレージ制御装置 1 0 1 に接続する。ストレージ装置 1 0 2 - 1、2、...、n と、ホットスペア 2 0 2 とは、ストレージ制御装置 1 0 1 に接続する。

## 【 0 0 3 8 】

ホストサーバ 2 0 1 は、ストレージシステム 1 0 0 の利用者が操作する装置から、ストレージ装置 1 0 2 - 1、2、... n へのアクセス要求を受け付ける装置である。ホットス

50

ア 2 0 2 は、ストレージ装置 1 0 2 - 1、2、...、n に対して予備となる装置である。

【 0 0 3 9 】

ここで、R A I D 技術について説明する。ストレージシステム 1 0 0 は、R A I D 技術により、1 つの仮想的なボリュームを形成する。ここで、仮想的なボリュームを形成するストレージ装置群を、R A I D グループに含まれるストレージ装置であると呼称する。

【 0 0 4 0 】

R A I D には、仮想的なボリュームの形成の仕方を表す R A I D レベルが存在する。R A I D レベルは、主に R A I D 0 ~ R A I D 6 までの R A I D レベルが存在する。また、R A I D には、R A I D 0 + 1 というように、R A I D レベルを組み合わせたレベルも存在する。たとえば、R A I D 0 は、冗長性を持たず、複数のストレージ装置にデータを分散する R A I D レベルである。また、R A I D 1 は、2 台のストレージ装置に同一のデータを記憶する R A I D レベルである。

10

【 0 0 4 1 】

ここで、R A I D 1 について、図 1 を用いて説明する。図 1 は、R A I D 1 の例を示す。ストレージ装置 1 0 2 - 1 が元データと冗長データとを記憶するならば、ストレージ制御装置 1 0 1 は、ストレージ装置 1 0 2 - 1 に読出要求 1 1 4 を発行する。具体的には、ストレージ装置 1 0 2 - 1 の L B A 0 - 9 9 9 は、データ d 1 - s r c を記憶し、L B A 1 0 0 0 - 1 0 0 7 は、データ d 2 を記憶しているとする。また、ストレージ装置 1 0 2 - 2 の L B A 0 - 9 9 9 は、データ d 1 と同一内容のデータ d 1 - r d n を記憶しているとする。この場合、ストレージ制御装置 1 0 1 は、ストレージ装置 1 0 2 - 2 に、読出先の L B A が L B A 0 - 9 9 9 となる読出要求 1 1 4 を発行することになる。

20

【 0 0 4 2 】

また、ストレージ装置 1 0 2 - 1 ~ 3 により R A I D 5 を形成することもできる。同様に、読出要求 1 に含まれる読出先の L B A がストレージ装置 1 0 2 - 1 であれば、ストレージ制御装置 1 0 1 は、冗長先のストレージ装置 1 0 2 - 2、3 に読出要求 1 1 4 を発行することになる。

【 0 0 4 3 】

R A I D を形成することにより、ストレージシステム 1 0 0 は、R A I D レベルが R A I D 0 以外であれば、利用者にデータの冗長性による信頼性の高い記憶領域を提供できる。また、R A I D を形成することにより、アクセスが複数のストレージ装置 1 0 2 に分散することから、ストレージシステム 1 0 0 は、利用者にアクセス応答が高速な記憶領域を提供できる。

30

【 0 0 4 4 】

また、複数台のストレージ装置 1 0 2 にデータを分散して I / O を高速化させる技術はストライピングと呼ばれ、データが X b y t e ずつ順番に R A I D グループを形成する仮想的なボリュームに配置される。X b y t e の各データを、「ストリップ ( s t r i p ) 」と呼称する。また、R A I D グループを形成する各ストレージ装置 1 0 2 に配置されたストリップを横つなかりに組み合わせた 1 つのデータセットを、「ストライプ ( s t r i p e ) 」と呼称する。

40

【 0 0 4 5 】

図 2 の例では、ストレージ制御装置 1 0 1 は、仮想的なボリュームを形成して、ホストサーバ 2 0 1 に提供する R A I D コントローラの機能を有する。仮想的なボリュームの提供を受けたホストサーバ 2 0 1 は、仮想的なボリュームの論理アドレスを用いてストレージ制御装置 1 0 1 にアクセスを行う。ストレージ制御装置 1 0 1 は、仮想的なボリュームの論理アドレスをストレージ装置 1 0 2 の L B A に変換して、変換した L B A を用いてストレージ装置 1 0 2 にアクセスする。以下の記載では、説明の簡略化のため、ホストサーバ 2 0 1 からのアクセス先の論理アドレスは、ストレージ装置 1 0 2 の L B A に変換済みであることを前提とする。

【 0 0 4 6 】

( ストレージ制御装置 1 0 1 のハードウェア )

50

図3は、ストレージ制御装置のハードウェア構成例を示すブロック図である。図3において、ストレージ制御装置101は、CPU(Central Processing Unit)301と、ROM(Read Only Memory)302と、RAM303と、を含む。また、ストレージ制御装置101は、ディスクドライブ304およびディスク305と、チャンネルアダプタ306と、I/Oコントローラ307と、SAS(Serial Attached SCSI)エキスパンダ308と、を含む。また、CPU301~ディスクドライブ304、チャンネルアダプタ306、I/Oコントローラ307はバス309によってそれぞれ接続される。

【0047】

CPU301は、ストレージ制御装置101の全体の制御を司る演算処理装置である。また、CPU301は、複数のプロセッサコアを有するマルチコアであってもよい。ROM302は、ブートプログラムなどのプログラムを記憶する不揮発性メモリである。RAM303は、CPU301のワークエリアとして使用される揮発性メモリである。

【0048】

ディスクドライブ304は、CPU301の制御に従ってディスク305に対するデータのリードおよびライトを制御する制御装置である。ディスクドライブ304には、たとえば、磁気ディスクドライブ、SSDなどを採用することができる。ディスク305は、ディスクドライブ304の制御で書き込まれたデータを記憶する不揮発性メモリである。たとえばディスクドライブ304が磁気ディスクドライブである場合、ディスク305には、磁気ディスクを採用することができる。また、ディスクドライブ304がSSDである場合、ディスク305には、半導体メモリを採用することができる。

【0049】

チャンネルアダプタ306は、ホストサーバ201に接続するアダプタである。I/Oコントローラ307は、ストレージ装置102-1、2、...、nと、ホットスペア202とのデータの入出力を制御する制御装置である。SASエキスパンダ308は、複数のSASデバイスを接続可能にする装置である。ストレージ装置102-1、2、...、nと、ホットスペア202は、SAS対応のデバイスである。

【0050】

また、ストレージ制御装置101は、ストレージシステム100の管理者等から直接操作される際を想定して、キーボード、マウスを有してもよい。

【0051】

(ストレージ制御装置101の機能)

次に、ストレージ制御装置101の機能について説明する。図4は、ストレージ制御装置の機能構成例を示すブロック図である。ストレージ制御装置101は、制御部401と、記憶部402とを有する。制御部401は、グループ分け部411と、計数部412と、アクセス制御部413とを含む。制御部401は、記憶装置に記憶された、本実施の形態にかかるストレージ装置の制御プログラムをCPU301が実行することにより、制御部401の機能を実現する。記憶装置とは、具体的には、たとえば、図3に示したROM302、RAM303、ディスク305などである。グループ分け部411~アクセス制御部413の処理結果は、RAM303、ディスク305などの記憶装置に格納される。

【0052】

また、ストレージ制御装置101は、ブロック予測テーブル111を記憶する記憶部402にアクセス可能である。記憶部402は、RAM303、ディスク305といった記憶装置に格納される。ブロック予測テーブル111の記憶内容の一例については、図5で後述する。

【0053】

グループ分け部411は、ストレージ制御装置101が受け付けた書込要求に応じて、書込要求から特定される書込先のLBAを、同一グループに含まれるLBAに対応するデータサイズがブロックデータサイズを超えないようにグループ分けする。また、グループ分け部411は、さらに、同一グループに含まれるLBAを含む書込要求のストレージ装

置 1 0 2 への発行間隔が所定間隔  $d_i$  を超えないようにグループ分けしてもよい。ここで、所定間隔  $d_i$  とは、SSD の仕様に関わる値であり、ある書込要求を受け付けた際に、フラッシュメモリに書き込む前に次の書込要求を待つ時間間隔である。SSD は、ある書込要求を受け付けてから所定間隔  $d_i$  を超えると次の書込要求を待たずにフラッシュメモリに書き込む。所定間隔  $d_i$  は、ストレージシステム 1 0 0 の管理者により設定される値である。ストレージシステム 1 0 0 の管理者は、ストレージ装置 1 0 2 の仕様を確認して、所定間隔  $d_i$  を設定する。具体的なグループ分けの一例と所定間隔  $d_i$  については、図 6 で後述する。

#### 【 0 0 5 4 】

計数部 4 1 2 は、ストレージ制御装置 1 0 1 が受け付けた読出要求に応じて、読出要求から特定される読出先の L B A を含むグループの読出回数を計数する。具体的なグループの読出回数の計数の例については、図 7 に後述する。

10

#### 【 0 0 5 5 】

また、計数部 4 1 2 は、グループの読出回数を初期化してもよい。初期化する契機について、アクセス制御部 4 1 3 の処理結果に基づくものがあるため、アクセス制御部 4 1 3 の説明の後に説明する。

#### 【 0 0 5 6 】

アクセス制御部 4 1 3 は、計数したグループの読出回数に基づき読出先の L B A を含む読出要求、または読出先の L B A のデータに対応する冗長データの記憶先の L B A を含む読出要求を発行する。また、アクセス制御部 4 1 3 は、計数したグループの読出回数が第 1 の回数となったことに応じて、グループに含まれる L B A を含む読出要求を、第 2 の回数、または読出要求に対する応答時間が所定時間を超えるまでストレージ装置 1 0 2 に発行してもよい。たとえば、アクセス制御部 4 1 3 は、計数したグループの読出回数が第 1 の回数となった後、グループに含まれる L B A を含む読出要求をストレージ装置 1 0 2 に発行する。また、アクセス制御部 4 1 3 は、計数したグループの読出回数が第 1 の回数となり、読出要求の応答をストレージ装置 1 0 2 から受け付けた後、グループに含まれる L B A を含む読出要求をストレージ装置 1 0 2 に発行してもよい。

20

#### 【 0 0 5 7 】

第 2 の回数は、ストレージ装置 1 0 2 の仕様に応じて設定された規定回数から第 1 の回数を減じた値である。所定時間は、ストレージ制御装置 1 0 1 とストレージ装置 1 0 2 の仕様によって設定される値であり、たとえば、読出要求に対する通常の応答時間の平均値に所定値を加えた値である。たとえば、ストレージ制御装置 1 0 1 は、サービス提供前に、ストレージ装置 1 0 2 に対し読出要求を何回か発行し、発行してから応答があるまでの時間の平均値を通常の応答時間として、通常の応答時間に所定値を加えた値を、所定時間として設定しておく。所定値は、ストレージシステム 1 0 0 の管理者が設定する値である。

30

#### 【 0 0 5 8 】

ここで、アクセス制御部 4 1 3 は、グループに含まれる L B A を含む読出要求を、第 2 の回数、または読出要求に対する応答時間が所定時間を超えるまでストレージ装置 1 0 2 に発行した場合、ストレージ装置 1 0 2 において R D 防止コピーが行われたと判断する。読出要求を発行することにより、ストレージ装置 1 0 2 に R D 防止コピーを行わせることができるため、以下、グループの読出回数が第 1 の回数となった後に発行する読出要求を、「R D 防止コピーを引き起こすコマンド」と呼称する。R D 防止コピーを引き起こすコマンドは、読出要求である R e a d でもよいし、記憶先が正確であるか否かを確認する V e r i f y でもよい。

40

#### 【 0 0 5 9 】

また、R D 防止コピーを引き起こすコマンドを第 2 の回数回発行することにより、R D 防止コピーが行われたと判断する方法を、以下、「第 1 の R D 防止コピー実行判断方法」と呼称する。また、応答時間が所定時間を超えるまでストレージ装置 1 0 2 に R D 防止コピーを引き起こすコマンドを発行することにより、R D 防止コピーが行われたと判断する

50

方法を、以下、「第2のRD防止コピー実行判断方法」と呼称する。

【0060】

第1のRD防止コピー実行判断方法の一例として、たとえば、RD防止コピー閾値が7700[回]であり、RD防止コピーが発生する規定回数が8000[回]であるとする。このとき、第2の回数は300[回]となる。このとき、計数したグループの読出回数が7700回となったことに応じて、ストレージ制御装置101は、RD防止コピーを引き起こすコマンドを300[回]発行した際に、RD防止コピーが行われたと判断する。第1のRD防止コピー実行判断方法を採用した場合、ストレージ制御装置101は、第2のRD防止コピー実行判断方法と比較して、コマンドの発行回数を比較すればよいため、第2のRD防止コピー実行判断方法より判断にかかる負荷を減らすことができる。

10

【0061】

第2のRD防止コピー実行判断方法を採用した場合、ストレージ制御装置101は、ストレージ装置102のRD防止コピーが発生する回数を記憶していなくてよい。したがって、第2のRD防止コピー実行判断方法を採用した場合、ストレージ制御装置101は、ストレージ装置102がRD防止コピーを実行する規定回数という仕様が不明であっても実施することができる。

【0062】

また、アクセス制御部413は、RD防止コピーを引き起こすコマンドを発行する間に、グループに含まれるLBAを含む書込要求をストレージ装置102に発行したことに応じて、グループに含まれる論理アドレスを含む読出要求の発行を停止してもよい。

20

【0063】

また、あるストレージ装置102において、計数したグループの読出回数が第1の回数となった後に、ストレージ制御装置101が読出要求を受け付けたとする。ここで、アクセス制御部413は、受け付けた読出要求に応じて、読出要求から特定されるLBAが上述のグループに含まれる場合に、読出先のLBAのデータに対応する冗長データの記憶先のLBAを含む読出要求を、他のストレージ装置102に発行してもよい。

【0064】

ここで、他のストレージ装置102は、あるストレージ装置102が記憶するデータの冗長データを記憶する装置である。たとえば、ストレージ装置102-1~3によりRAID5を形成するとする。そして、ストレージ装置102-1がデータ1を記憶しており、ストレージ装置102-2がデータ2を記憶しており、ストレージ装置102-3がデータ1およびデータ2から生成されるパリティデータを記憶するとする。このとき、ストレージ装置102が記憶するデータ1の冗長データを記憶する他のストレージ装置102は、パリティデータを記憶するストレージ装置102-3と、データ2を記憶するストレージ装置102-2とである。

30

【0065】

また、計数部412は、アクセス制御部413により、RD防止コピーが行われたと判断した際に、グループの読出回数を初期化してもよい。また、計数部412は、グループに含まれるLBAを含む書込要求をストレージ装置102に発行したことに応じて、グループの読出回数を初期化してもよい。また、計数部412は、アクセス制御部413がRD防止コピーを引き起こすコマンドを発行する間に、グループに含まれるLBAを含む書込要求をストレージ装置102に発行したことに応じて、グループの読出回数を初期化してもよい。

40

【0066】

次に、あるストレージ装置102と他のストレージ装置102とのブロックデータサイズが同一の値であり、かつ、あるストレージ装置102と他のストレージ装置102との規定回数が同一の値であるときに用いる機能について説明する。

【0067】

このとき、グループ分け部411は、他のストレージ装置102への書込要求に応じて、書込要求に含まれる書込先のLBAを、同一グループに含まれるLBAに対応するデー

50

タサイズがブロックデータサイズを超えないようにグループ分けする。また、計数部 4 1 2 は、他のストレージ装置 1 0 2 への読出要求に応じて、読出要求に含まれる読出先の L B A を含むグループの読出回数を計数する。

【 0 0 6 8 】

また、アクセス制御部 4 1 3 は、計数したグループの読出回数が、第 3 の回数となったことに応じて、グループに含まれる L B A を含む読出要求を、第 4 の回数、または、読出要求に対する応答時間が所定時間を超えるまで他のストレージ装置 1 0 2 に発行する。ここで、第 3 の回数は、ストレージ装置 1 0 2 の仕様に応じて設定された規定回数より小さく前記第 1 の回数とは異なる値である。また、第 4 の回数は、規定回数から第 3 の回数を減じた値である。ブロックデータサイズと規定回数とがそれぞれ同一の値となる例について、図 9 および図 1 0 に後述する。

10

【 0 0 6 9 】

次に、ストレージ装置 1 0 2 のデータを、S S D を適用したストレージ装置 1 0 2 に移行するときに用いる機能について説明する。データを移行する理由としては、図 1 1、図 1 2 で後述する。移行先のストレージ装置 1 0 2 は、たとえば、未使用であるホットスペア 2 0 2 である。以下、移行先のストレージ装置 1 0 2 が、ホットスペア 2 0 2 であるとして説明する。

【 0 0 7 0 】

このとき、アクセス制御部 4 1 3 は、受け付けた書込要求から特定されるストレージ装置 1 0 2 の書込先の L B A をグループ分けした複数のグループに含まれる L B A を含むデータ移行要求を、ストレージ装置 1 0 2 に発行する。このとき、アクセス制御部 4 1 3 は、ホットスペア 2 0 2 へのデータの移行間隔が所定間隔  $d_i$  を超えないように発行する。また、グループ分け部 4 1 1 は、発行したデータ移行要求に含まれる L B A を、同一グループに含まれる L B A に対応するデータサイズがホットスペア 2 0 2 のブロックデータサイズを超えないようにグループ分けする。

20

【 0 0 7 1 】

また、グループ分け部 4 1 1 は、受け付けた書込要求から特定されるホットスペア 2 0 2 の書込先の L B A を、同一グループに含まれる L B A に対応するデータサイズが当該ブロックデータサイズを超えないようにグループ分けする。また、計数部 4 1 2 は、受け付けた読出要求に応じて、受け付けた読出要求から特定されるホットスペア 2 0 2 の読出先の L B A を含むグループの読出回数を計数する。また、アクセス制御部 4 1 3 は、計数したグループの読出回数に基づき読出先の論理アドレスを含む読出要求、または、読出先の論理アドレスのデータに対応する冗長データの記憶先の論理アドレスを含む読出要求を発行する。

30

【 0 0 7 2 】

図 5 は、ブロック予測テーブルの記憶内容の一例を示す説明図である。ブロック予測テーブル 1 1 1 は、ストレージ装置 1 0 2 のうち、S S D が適用されたストレージ装置 1 0 2 ごとに有するテーブルである。たとえば、図 5 に示すブロック予測テーブル 1 1 1 - 1 は、ストレージ装置 1 0 2 - 1 のブロックと L B A 割り当ての関係とを登録するテーブルである。図 5 に示すブロック予測テーブル 1 1 1 - 1 は、レコード 5 0 1 - 1 ~ 3 を有する。

40

【 0 0 7 3 】

ブロック予測テーブル 1 1 1 は、グループ番号と、L B A と、読出回数という 3 つのフィールドを有する。グループ番号フィールドは、書込要求に含まれる書込先の論理アドレスを、同一グループに含まれる論理アドレスに対応するデータサイズがブロックデータサイズを超えないようにグループ分けした際の、グループの識別番号を記憶するフィールドである。L B A フィールドは、該当のグループにグループ分けされた L B A を記憶するフィールドである。読出回数フィールドは、該当のグループの読出回数を記憶するフィールドである。

【 0 0 7 4 】

50

たとえば、レコード501-1は、LBA0-999が一つのブロックに割り当てられたと予測したグループG0にグループ分けされており、グループG0の読出回数が7700[回]であることを示す。

【0075】

(SSDへの書込時と読出時との動作例)

次に、図6～図8を用いて、ストレージ装置102-1への書込時と読出時との動作例について説明する。ここで、本実施の形態において、SSDが適用されたストレージ装置102における1ブロックのブロックデータサイズは、説明の簡略化のため、全て512[Kバイト]であるとする。そして、1ブロックには、1000個のLBAが割当可能であるとする。したがって、1つのLBAは、512[バイト]の記憶領域を有する。

10

【0076】

図6は、書込時における動作例を示す説明図である。図6では、SSDが適用されたストレージ装置102-1への書込時において、ブロック予測テーブル111の更新例について説明する。以下の説明において、書込先のLBAが割り当て済みでない書込要求を「新規書込要求」と称する。また、書込先のLBAが割り当て済みの書込要求を「上書き書込要求」と称する。

【0077】

ここで、一般的なSSD内部の動作について説明する。SSDは、ブロック(512MB等)に複数のLBAを割り当てる。割り当て動作は、通常固定ではなく、データが書き込まれた際に順次関連付けが行われる。SSD内部にあるSSDコントローラは、SSDに転送されたデータをSSD内部にあるキャッシュメモリに一時保存して、纏まった単位でフラッシュメモリへ書き込む。SSDコントローラは、このデータ纏め作業を優先で実施するが、転送データの発行間隔が所定間隔di、たとえば10秒を超えた場合、電源遮断によるデータ損失を防止するため、以降のデータを待つことなくフラッシュメモリへ書き込む。そして、SSDコントローラは、新規書込要求に関して、一旦書き込んだブロックに書き込まない。転送されたデータが新たに発生した場合、SSDコントローラは、次のブロックに書き込む。

20

【0078】

上述したSSDの動作から、ストレージ制御装置101でも同様の管理を行うことにより、ストレージ制御装置101は、ブロックとLBA割り当てとの対応関係について予測することができる。ストレージ制御装置101は、ストレージ装置102が同一のブロックに割り当てたであろうLBA群を、同一のグループに設定する。

30

【0079】

具体的な予測例として、ホストサーバ201から、以下に示す第1の書込要求～第3の書込要求を受け付けたとして、ブロックとLBA割り当てとの関係を説明する。第1の書込要求～第3の書込要求は、いずれも新規書込要求であるとする。ストレージ制御装置101は、書込要求を受け付けた場合、ブロック予測テーブル111を参照して、受け付けた書込要求が新規書込要求か上書き書込要求かを判断するが、図6では説明を省略し、図8で説明を行う。また、まだLBAが一つも割り当てられていないストレージ装置102に対して、ストレージ制御装置101が書込要求に含まれる書込先のLBAを含める割当中のグループが、G0であるとする。

40

【0080】

第1の書込要求は、時刻t0から時刻t1までにかけて発生しており、書込先のLBAがLBA0-1007となる書込要求である。LBA0-1007のデータサイズは、516,096[バイト]に相当する。第2の書込要求は、時刻t1の2秒後である時刻t2から時刻t3までにおけるLBA6328-6345に対する書込要求である。LBA6328-6345のデータサイズは、9,216[バイト]に相当する。第3の書込要求は、時刻t3の13秒後である時刻t4から時刻t5にかけて、LBA4230-4529に対する書込要求である。LBA4230-4529のデータサイズは、153,600[バイト]に相当する。第1の書込要求～第3の書込要求を受け付けた後、ストレ

50

ジ制御装置 101 は、それぞれをストレージ装置 102 に発行する。

【0081】

時刻  $t_0$  において、ストレージ制御装置 101 は、第 1 の書込要求に含まれる書込先の LBA のうちの LBA0 - 999 を、割当中のグループである G0 に含めるようにグループ分けする。そして、G0 に LBA の空きがなくなったため、ストレージ制御装置 101 は、割当中のグループを、現在の割当中のグループの次のグループである G1 に設定する。次に、ストレージ制御装置 101 は、第 1 の書込要求に含まれる書込先の LBA のうちの LBA1000 - 1007 を、割当中のグループである G1 に含めるようにグループ分けする。

【0082】

次に、時刻  $t_2$  において第 2 の書込要求を受け付けると、ストレージ制御装置 101 は、第 1 の書込要求の完了から第 2 の書込要求の発生までの間隔が所定間隔  $d_i$  を超えないため割当中のグループを G1 のままとする。そして、ストレージ制御装置 101 は、第 2 の書込要求に含まれる書込先の LBA6328 - 6345 を、割当中のグループである G1 に含めるようにグループ分けする。ここで、10 秒経過したか否かの判断方法について、ストレージ制御装置 101 は、時刻  $t_1$  時点でタイマをスタートしておき、時刻  $t_2$  時点でのタイマが示す時間が 10 秒を超えたか否かで判断することができる。ストレージ制御装置 101 は、タイマを、ストレージ装置 102 ごとに用意しておく。

【0083】

続けて、時刻  $t_4$  において、ストレージ制御装置 101 は、第 3 の書込要求を受け付けると、第 2 の書込要求の完了から第 3 の書込要求の発生までの間隔が所定間隔  $d_i$  を超えるため、割当中のグループを次のグループ G2 に設定する。そして、ストレージ制御装置 101 は、第 3 の書込要求に含まれる書込先の LBA4230 - 4529 を、割当中のグループである G2 に含めるようにグループ分けする。

【0084】

図 7 は、読出時における動作例を示す説明図である。図 7 では、ストレージ装置 102 への読出時において、ブロック予測テーブル 111 の更新例について説明する。

【0085】

ここで、RD 防止コピーの動作は、各グループに対し、何回の読出要求を受けたかで予測することができる。したがって、ストレージ制御装置 101 は、データの読出要求を行う際に、ブロック予測テーブル 111 の LBA を参照し、該当するグループの読出回数を確認する。これにより、ストレージ制御装置 101 は、読出要求に含まれる読出先の LBA が含まれるであろうブロックの読出回数を予測することができる。ここで、RD 防止コピーが行われる読出回数は、SSD の仕様によって様々である。たとえば、SSD コントローラは、あるブロックに対して 8000 回読み出しが行われた場合に、RD 防止コピーを行う。そこで、ストレージ制御装置 101 は、SSD コントローラが RD 防止コピーを行う読出回数より小さい値を RD 防止コピー閾値として、読出回数が RD 防止コピー閾値に到達する場合に、RD 防止コピーが行われる可能性があるとして判断する。

【0086】

具体的な予測例として、図 7 では、図 6 に示した第 3 の書込要求が発生した後、ホストサーバ 201 から、以下に示す第 1 の読出要求～第 3 の読出要求があったとして、RD 防止コピーが行われることを予測する動作を説明する。また、RD 防止コピー閾値を 7900 [回] とする。

【0087】

第 1 の読出要求は、時刻  $t_6$  に発生した、読出先の LBA が LBA0 - 1007 となる読出要求である。第 2 の読出要求は、時刻  $t_7$  に発生した、読出先の LBA が LBA6328 - 6345 となる読出要求である。第 3 の読出要求は、時刻  $t_8$  に発生した、読出先の LBA が LBA6328 - 6329 となる読出要求である。第 1 の読出要求～第 3 の読出要求を受け付けた後、ストレージ制御装置 101 は、それぞれをストレージ装置 102 に発行する。また、第 2 の読出要求と第 3 の読出要求との間に、LBA1000 - 100

10

20

30

40

50



7、L B A 6 3 2 8 - 6 3 4 5 のいずれかに対する読出要求が 7 8 9 7 回あったものとする。

【 0 0 8 8 】

時刻 t 6 において、第 1 の読出要求に含まれる読出先の L B A のうちの L B A 0 - 9 9 9 は G 0 に含まれるため、ストレージ制御装置 1 0 1 は、G 0 の読出回数をインクリメントする。また、第 1 の読出要求の読出先の L B A のうちの L B A 1 0 0 0 - 1 0 0 7 は G 1 に含まれるため、ストレージ制御装置 1 0 1 は、G 1 の読出回数をインクリメントする。

【 0 0 8 9 】

時刻 t 7 において、第 2 の読出要求に含まれる読出先の L B A である L B A 6 3 2 8 - 6 3 4 5 は G 1 に含まれるため、ストレージ制御装置 1 0 1 は、G 1 の読出回数をインクリメントする。

【 0 0 9 0 】

時刻 t 8 となる前段階で G 1 に対する読出回数は、7 8 9 9 [ 回 ] になったものとする。そして、時刻 t 8 において、第 3 の読出要求に含まれる読出先の L B A である L B A 6 3 2 8 - 6 3 4 5 は G 1 に含まれるため、ストレージ制御装置 1 0 1 は、G 1 の読出回数をインクリメントする。インクリメントした結果、G 1 の読出回数が 7 9 0 0 [ 回 ] となり、R D 防止コピー閾値に到達したため、ストレージ制御装置 1 0 1 は、R D 防止コピーが行われる可能性があると判断する。

【 0 0 9 1 】

R D 防止コピーが行われる可能性があると判断した後、ストレージ制御装置 1 0 1 は、R D 防止コピーを引き起こすコマンドをストレージ装置 1 0 2 - 1 に発行する。図 7 の例では、ストレージ制御装置 1 0 1 は、時刻 t 8 から時刻 t 9 にかけて、R D 防止コピーを引き起こすコマンドをストレージ装置 1 0 2 - 1 に発行する。そして、ストレージ制御装置 1 0 1 は、R D 防止コピーが行われたと判断した場合に、G 1 の読出回数を 0 [ 回 ] に初期化する。ここで、R D 防止コピーが行われたと判断する方法としては、たとえば、上述した第 1 の R D 防止コピー実行判断方法と、第 2 の R D 防止コピー実行判断方法とがある。

【 0 0 9 2 】

図 8 は、上書き書込時における動作例を示す説明図である。図 8 では、ストレージ装置 1 0 2 への書込時において、既に書込要求を行った L B A に対して、さらに書込要求があった場合と、R D 防止コピーを引き起こすコマンドの発行中に書込要求があった場合のブロック予測テーブル 1 1 1 の更新例について説明する。

【 0 0 9 3 】

ここで、一般的な S S D 内部における上書き書込時の動作について説明する。S S D は、上書き書込を行う場合、ブロック単位でリードモディファイライト動作を行う。具体的に、S S D は、書込先の L B A を含むブロックの記憶領域を読み出して、読み出したデータに上書き書込のデータを反映して、反映したデータを書込先の L B A を含むブロックに書き込む。ブロックに書き込む際に、S S D は、対象のブロックの記憶内容を消去した後、書き込みを行う。ブロックの記憶内容が消去されて書き込みが行われると、浮遊ゲートの電子量に応じた閾値電圧と、読み出し時に印加する電圧との間に十分なマージンが発生することになる。このように、上書き書込を行うと、十分なマージンが発生するので、S S D は、書き込みを行ったブロックの読出回数を初期化する。

【 0 0 9 4 】

上述した S S D の動作から、ストレージ制御装置 1 0 1 でも同様の予測を行うことにより、ブロックの予測精度を向上することができる。具体的な更新例として、図 8 では、図 7 に示した第 3 の読出要求が発生した後、ホストサーバ 2 0 1 から、第 4 の書込要求と、第 4 の読出要求と、第 5 の書込要求とを受け付けたとして、ブロック予測テーブル 1 1 1 の更新例について説明する。

【 0 0 9 5 】

第4の書込要求は、時刻 $t_{10}$ から時刻 $t_{11}$ までにかけて発生しており、書込先のLBAがLBA0-899となる書込要求である。LBA0-899のデータサイズは、460,800[バイト]に相当する。第4の読出要求は、時刻 $t_{12}$ に発生した、読出先のLBAがLBA6238-6345となる読出要求である。第5の書込要求は、時刻 $t_{13}$ から時刻 $t_{14}$ までにかけて発生しており、書込先のLBAがLBA6238-6345となる書込要求である。また、第4の書込要求と第4の読出要求との間に、LBA1000-1007、LBA6328-6345のいずれかに対する読出要求が7899回あったものとする。

【0096】

ストレージ制御装置101は、第4の書込要求を受け付けた場合、第4の書込要求が新規書込要求か上書書込要求かを判断するため、ブロック予測テーブル111に、第4の書込要求に含まれる書込先のLBAがあるか否かを判断する。図8の例では、G0にLBA0-899が含まれることから、ストレージ制御装置101は、第4の書込要求に含まれる書込先のLBAがあり、第4の書込要求が上書書込要求であると判断する。ここで、図8では、書込要求に含まれる書込先のLBA全てがブロック予測テーブル111にあった例を示したが、書込要求に含まれる書込先のLBAの一部がブロック予測テーブル111にある場合も起こり得る。この場合、ストレージ制御装置101は、書込要求のうち、ブロック予測テーブル111にあるLBAを上書書込要求として扱い、書込要求のうち、ブロック予測テーブル111にないLBAを新規書込要求として扱う。

【0097】

時刻 $t_{10}$ において、第4の書込要求が上書書込要求であるため、ストレージ制御装置101は、第4の書込要求に含まれる書込先のLBAを含むG0の読出回数を0[回]に初期化する。また、第4の書込要求を受け付けた後、ストレージ制御装置101は、受け付けた第4の読出要求をストレージ装置102に発行する。

【0098】

時刻 $t_{12}$ において、時刻 $t_{12}$ となる前段階でG1に対する読出回数は、7899[回]になったものとする。そして、時刻 $t_{12}$ において、第4の読出要求に含まれる読出先のLBAであるLBA6328-6345はG1に含まれるため、ストレージ制御装置101は、G1の読出回数をインクリメントする。インクリメントした結果、G1の読出回数がRD防止コピー閾値に到達したため、ストレージ制御装置101は、RD防止コピーが行われる可能性があるとして判断し、RD防止コピーを引き起こすコマンドをストレージ装置102-1に発行する。

【0099】

時刻 $t_{13}$ において、RD防止コピーを引き起こすコマンドを発行中に、第5の書込要求を受け付けると、ストレージ制御装置101は、第5の書込要求に含まれる書込先のLBAが、RD防止コピーの対象のグループに含まれるかを判断する。図8の例では、第5の書込要求に含まれる書込先のLBAが、RD防止コピーの対象であるG1に含まれるため、ストレージ制御装置101は、RD防止コピーを引き起こすコマンドの発行を停止する。同時に、ストレージ制御装置101は、G1の読出回数を0[回]に初期化する。

【0100】

(RD防止コピー閾値の設定例)

次に、図9と図10とを用いて、RD防止コピー閾値の設定例を示す。図9と図10とで用いる例の前提として、ストレージ装置102-1~3が、ブロックデータサイズが同一のサイズであり、RD防止コピーが起こる回数が同一となるという、同一の特性を有するものとする。この場合、ストレージ制御装置101は、ストレージ装置102-1~3のRD防止コピー閾値を、それぞれ、7700[回]、7800[回]、7900[回]に設定しておく。また、SSDが適用されたストレージ装置102-1~3により、RAIDレベルがRAID5であるRAIDグループを形成するものとする。そして、ストレージ装置102-1のLBA1000-1007と、ストレージ装置102-2のLBA1000-1007と、ストレージ装置102-3のLBA1000-1007とにより

10

20

30

40

50

、1つのストライプが形成されたとする。

【0101】

ストレージ装置102-1のLBA1000-1007はブロックB<sub>x</sub>に対応付けられており、ストレージ制御装置101は、ストレージ装置102-1のLBA1000-1007をグループG<sub>x</sub>に対応付けたものとする。同様に、ストレージ装置102-2のLBA1000-1007はブロックB<sub>y</sub>に対応付けられており、ストレージ制御装置101は、ストレージ装置102-2のLBA1000-1007をグループG<sub>y</sub>に対応付けたものとする。また、ストレージ装置102-3のLBA1000-1007はブロックB<sub>z</sub>に対応付けられており、ストレージ制御装置101は、ストレージ装置102-3のLBA1000-1007をグループG<sub>z</sub>に対応付けたものとする。

10

【0102】

さらに、該当のストライプにおいて、ストレージ装置102-1、2には実データd<sub>1</sub>、d<sub>2</sub>が格納されており、ストレージ装置102-3には実データd<sub>1</sub>、d<sub>2</sub>によるパリティデータp(d<sub>1</sub>, d<sub>2</sub>)が格納されたとする。以下、説明の簡略化のため、実データd<sub>1</sub>、d<sub>2</sub>、パリティデータp(d<sub>1</sub>, d<sub>2</sub>)を、単に、d<sub>1</sub>、d<sub>2</sub>、p(d<sub>1</sub>, d<sub>2</sub>)と記載する。

【0103】

図9は、RD防止コピー閾値の設定例を示す説明図(その1)である。図9に示すストレージシステム100として、図9に示すブロック予測テーブル111-1は、ストレージ装置102-1のG<sub>x</sub>の読出回数が7700[回]であることを記憶する。また、図9に示すブロック予測テーブル111-2は、ストレージ装置102-2のG<sub>y</sub>の読出回数が7700[回]であることを記憶する。さらに、図9に示すブロック予測テーブル111-3は、ストレージ装置102-3のG<sub>z</sub>の読出回数が0[回]であることを記憶する。

20

【0104】

また、ストレージ制御装置101は、ストレージ装置102-1~3のRD防止処理フラグを「処理していない」に設定したとする。ここで、RD防止処理フラグは、該当のストレージ装置102が、RD防止コピーを行っているであろうか否かを判断するフラグである。RD防止処理フラグは、ストレージ装置102ごとにあるデータである。ストレージ制御装置101は、RD防止処理フラグを、たとえば、ブロック予測テーブル111に

30

【0105】

ここで、RD防止コピー閾値が各ストレージ装置102間で同一の場合に起こりうる現象について説明する。同一のRAIDグループに含まれるストレージ装置102同士は、ほぼ均等に読出が行われる可能性が高く、読出回数の値が近い値になり易い。読出回数の値が近い値になった結果、あるタイミングにおいてRD防止コピーを引き起こすコマンドを発行中の各ストレージ装置102が同時に複数台存在する状態になる可能性がある。上述の状態となり、RD防止コピーを引き起こすコマンドを発行中でないストレージ装置102の数が、データが復元可能な台数を下回ると、ストレージ制御装置101は、RD防止コピーが行われたと判断するまで読出要求に対するデータを取得できなくなる。

40

【0106】

そこで、ストレージ制御装置101は、RD防止コピー閾値をストレージ装置102間でずらして設定しておく。図9の例では、ストレージ装置102-1のG<sub>x</sub>の読出回数の値が、ストレージ装置102-1のRD防止コピー閾値である7700[回]に到達しており、ストレージ制御装置101は、ストレージ装置102-1のG<sub>x</sub>への読出を抑止す

50

る。また、読出の抑止を行った後、ストレージ制御装置 101 は、ストレージ装置 102 - 1 の R D 防止処理フラグを「処理中」に設定するとともに、ストレージ装置 102 - 1 に R D 防止コピーを引き起こすコマンドを発行する。コマンドを発行して R D 防止コピーが行われたと判断した場合、ストレージ制御装置 101 は、ストレージ装置 102 - 1 の G x の読出回数を 0 [ 回 ] に初期化するとともに、ストレージ装置 102 - 1 の R D 防止処理フラグを「処理していない」に設定する。

#### 【 0 1 0 7 】

また、読出の抑止中において、図 9 に示すように、ホストサーバ 201 からストレージ装置 102 - 1 の L B A 1000 - 1007 の読出要求があったとする。このとき、ストレージ制御装置 101 は、ストレージ装置 102 - 2 から d 2 を読み出すとともに、ストレージ装置 102 - 3 から p ( d 1 , d 2 ) を読み出す。読み出したことにより、ストレージ制御装置 101 は、ストレージ装置 102 - 2 の G y の読出回数と、ストレージ装置 102 - 3 の G z の読出回数とをインクリメントする。続けて、ストレージ制御装置 101 は、d 2 と p ( d 1 , d 2 ) とから d 1 を復元して、ホストサーバ 201 に d 1 を返却する。

#### 【 0 1 0 8 】

図 10 は、R D 防止コピー閾値の設定例を示す説明図 ( その 2 ) である。図 10 に示すストレージシステム 100 の状態は、図 9 で示したストレージシステム 100 の状態から、下記に示す第 1 の動作と第 2 の動作とが実行されたことにより、ストレージシステム 100 の状態が変化したものである。

#### 【 0 1 0 9 】

第 1 の動作は、ストレージ装置 102 - 1 に対して R D 防止コピーを引き起こすコマンドを複数発行中にストレージ装置 102 - 1 の L B A 1000 - 1007 の読出要求が 49 回発行されたという動作である。第 1 の動作により、ストレージ制御装置 101 は、グループ G y の読出回数とグループ G z の読出回数とをそれぞれ 49 回インクリメントする。

#### 【 0 1 1 0 】

第 2 の動作は、ストレージ装置 102 - 1 に対して R D 防止コピーを引き起こすコマンドを複数発行中と発行終了後との合計で、ストレージ装置 102 - 2 の L B A 1000 - 1007 の読出要求が 50 回発行されたという動作である。第 2 の動作により、ストレージ制御装置 101 は、グループ G y の読出回数を 50 回インクリメントする。

#### 【 0 1 1 1 】

第 1 の動作と第 2 の動作とにより、図 10 に示すブロック予測テーブル 111 には、次に示す値が登録される。次に示す値として、具体的に、図 10 に示すブロック予測テーブル 111 - 1 は、ストレージ装置 102 - 1 の G x の読出回数が 0 [ 回 ] であることを記憶する。また、図 10 に示すブロック予測テーブル 111 - 2 には、ストレージ装置 102 - 2 の G y の読出回数が 7800 [ 回 ] であることを登録される。さらに、図 10 に示すブロック予測テーブル 111 - 3 には、ストレージ装置 102 - 3 の G z の読出回数が 50 [ 回 ] であることが登録される。

#### 【 0 1 1 2 】

図 10 の例では、ストレージ装置 102 - 2 の G y の読出回数の値が R D 防止コピー閾値である 7800 [ 回 ] に到達したため、ストレージ制御装置 101 は、ストレージ装置 102 - 2 の G y への読出を抑止する。また、読出の抑止を行った後、ストレージ制御装置 101 は、ストレージ装置 102 - 2 の R D 防止処理フラグを「処理中」に設定するとともに、ストレージ装置 102 - 2 に R D 防止コピーを引き起こすコマンドを複数発行する。コマンドを複数発行して R D 防止コピーが行われたと判断した場合、ストレージ制御装置 101 は、ストレージ装置 102 - 2 の G y の読出回数を 0 [ 回 ] に初期化するとともに、ストレージ装置 102 - 2 の R D 防止処理フラグを「処理していない」に設定する。

#### 【 0 1 1 3 】

また、読出の抑止中において、図 10 に示すように、ホストサーバ 201 からストレージ装置 102 - 2 の LBA 1000 - 1007 の読出要求があったとする。このとき、ストレージ制御装置 101 は、ストレージ装置 102 - 1 から d1 を読み出すとともに、ストレージ装置 102 - 3 から p(d1, d2) を読み出す。読み出したことにより、ストレージ制御装置 101 は、ストレージ装置 102 - 1 の Gx の読出回数と、ストレージ装置 102 - 3 の Gz の読出回数とをインクリメントする。続けて、ストレージ制御装置 101 は、d1 と p(d1, d2) とから d2 を復元して、ホストサーバ 201 に d2 を返却する。

#### 【0114】

このように、RD 防止コピー閾値をストレージ装置 102 間でずらして設定することにより、RD 防止コピーが行われる期間をストレージ装置 102 間でずらすことができる。RD 防止コピーが行われる期間をずらすことにより、ストレージシステム 100 は、ミラーリングからのデータ取得ができない状態、またはパリティを用いたデータ復元ができない状態になることを避けることができる。

#### 【0115】

(ガベージコレクションとリフレッシュ処理との説明)

次に、図 11 を用いて、ガベージコレクションによりストレージ装置 102 が管理するブロックと、ストレージ制御装置 101 が予測するグループとが乖離した様子について説明する。また、図 12 において、リフレッシュ処理により乖離が解消される様子について説明する。ここで、図 11 と図 12 とで共通して用いる前提として、ストレージ装置 102 - 1 が、ガベージコレクションを行うものとする。

#### 【0116】

図 11 は、ガベージコレクションが行われた際のブロックとグループとの乖離の一例を示す説明図である。表 1101 は、ガベージコレクションの実行前の状態における、ストレージ装置 102 が管理するブロックと LBA 割り当てとの関係を示すテーブルである。また、表 1102 は、ガベージコレクションの実行後の状態における、ストレージ装置 102 が管理するブロックと LBA 割り当てとの関係を示すテーブルである。

#### 【0117】

ガベージコレクションの実行前の状態として、表 1101 とブロック予測テーブル 111 - 1 より、ストレージ装置 102 - 1 の LBA 0 - 999 は、ブロック B0 に対応付けられており、ストレージ制御装置 101 は、LBA 0 - 999 を G0 に含めたものとする。B0 の読出回数と G0 の読出回数とは、ともに、20 [回] であるとする。

#### 【0118】

また、ストレージ装置 102 - 1 の LBA 1000 - 1007、6328 - 6345 は、ブロック B1 に対応付けられており、ストレージ制御装置 101 は、LBA 1000 - 1007、6328 - 6345 をグループ G1 に含めたものとする。B1 の読出回数と G1 の読出回数とは、ともに、28 [回] であるとする。

#### 【0119】

また、ストレージ装置 102 - 1 の LBA 1008 - 1099、4230 - 4529 は、ブロック B2 に対応付けられており、ストレージ制御装置 101 は、LBA 1008 - 1099、4230 - 4529 をグループ G2 に含めたものとする。B2 の読出回数と G2 の読出回数とは、ともに、2367 [回] であるとする。

#### 【0120】

また、ストレージ装置 102 - 1 の LBA 1100 - 2099 は、ブロック B1000 に対応付けられており、ストレージ制御装置 101 は、LBA 1100 - 2099 をグループ G1000 に含めたものとする。B1000 の読出回数と G1000 の読出回数とは、ともに、692 [回] であるとする。

#### 【0121】

この状態で、ストレージ装置 102 - 1 が、ガベージコレクションを実行したとする。具体的には、ストレージ装置 102 - 1 の SSD コントローラは、B1 の LBA 数が 26

10

20

30

40

50

【個】であり、B2のLBA数が392【個】であるから、B1のLBAとB2のLBAとを纏めても、一つのブロックに収まるLBA数1000を超えないと判断する。そして、ストレージ装置102-1のSSDコントローラは、B2が示す記憶領域に記憶されたLBA1008-1099、4230-4529のデータを、B1の記憶領域に移行する。

#### 【0122】

ガベージコレクションの実行後において、ストレージ装置102-1のLBA1008-1099、4230-4529は、LBA1000-1007、6328-6345と併せてB1に対応付けられる。これに対し、ストレージ制御装置101は、LBA1000-1007、6328-6345をG1に、LBA1008-1099、4230-4529をG2にグループ分けしており、LBAの割当内容が乖離したことになる。

10

#### 【0123】

図12は、リフレッシュ処理の実行前後におけるブロックとグループとの関係の一例を示す説明図である。リフレッシュ処理は、LBAの割当内容が乖離したストレージ装置を移行元のストレージ装置102とし、移行元のストレージ装置102が記憶する移行対象データを移行先のストレージ装置102に移行させて、グループとLBAとの関係を再定義する処理である。移行先のストレージ装置102は、LBA割り当てがない、未使用のストレージ装置102であることが好ましい。LBA割り当てがある、使用中のストレージ装置102に移行すると、既にLBAの割当内容が乖離している可能性があるためである。図12の例では、移行先のストレージ装置がホットスペア202であるとする。また、ホットスペア202は、SSDを適用したストレージ装置102であるとする。

20

#### 【0124】

ここで、データの移行方法は、たとえば、下記に示す第1の移行方法と、第2の移行方法とがある。第1の移行方法は、ストレージ制御装置101が移行元となるストレージ装置102にデータ移行要求を発行し、移行元のストレージ装置102が、移行先のストレージ装置102に直接移行対象データを送信する方法である。また、第2の移行方法は、データ移行要求を受け付けた移行元のストレージ装置102が移行対象データをストレージ制御装置101に送信し、ストレージ制御装置101は、移行対象データを移行先のストレージ装置102に送信する方法である。本実施の形態では、第1の移行方法を採用したとする。

30

#### 【0125】

どちらの場合であっても、ストレージ制御装置101は、移行先のストレージ装置102への移行対象データのうちの一つのグループに含めるようにグループ分けしたデータ内の移行間隔が所定間隔diを超えないようにする。たとえば、一つのグループに含めるようにグループ分けしたLBA群のデータ512【Kバイト】が、256【Kバイト】、256【Kバイト】というように分割して移行先のストレージ装置102に送信されたとする。さらに、1回目の256【Kバイト】のデータが移行してから、2回目の256【Kバイト】のデータが移行するまでの時間が、所定間隔diを超えたとする。このとき、ストレージ制御装置101では一つのグループに含めるようにグループ分けされたLBA群が、移行先のストレージ装置102の2つのブロックに割り当てられてしまい、LBA割り当ての内容に乖離が発生してしまう。

40

#### 【0126】

移行間隔が所定間隔diを超えないようにするために、第1の移行方法の場合には、ストレージシステム100を構築する際に、ストレージシステム100の管理者が、ストレージ装置102のデータ転送の間隔が所定間隔di以内であることを確認しておけばよい。また、第2の移行方法の場合には、ストレージ制御装置101が、たとえば、移行元のストレージ装置102から受け付けた移行対象データをバッファリングして、移行先のストレージ装置102にブロックデータサイズごとに転送すればよい。

#### 【0127】

データ移行要求は、移行対象データが記憶された記憶領域を示すLBAと、移行対象デ

50

ータの移行先のストレージ装置 102 を特定する情報と、を含む。

【0128】

次に、図 12 の例を用いて、リフレッシュ処理の動作例を示す。ストレージ制御装置 101 は、LBA0 - 2099、4230 - 4529、6328 - 6345 をホットスペア 202 に移行するデータ移行要求を、ストレージ装置 102 - 1 に発行する。

【0129】

そして、ストレージ制御装置 101 は、まだグループ分けしていない LBA のうちの先頭の LBA を選択する。図 12 の例では、ストレージ制御装置 101 は、LBA0 を選択する。次に、ストレージ制御装置 101 は、ホットスペア 202 のブロックデータサイズである、512 [K バイト] を超えず、選択した LBA を先頭とした、LBA0 - 999 を、G0 にグループ分けする。

10

【0130】

次に、ストレージ制御装置 101 は、まだグループ分けしていない LBA のうちの先頭となる LBA1000 を選択する。そして、ストレージ制御装置 101 は、512 [K バイト] を超えず、選択した LBA を先頭とした、LBA1000 - 1999 を、G1 にグループ分けする。

【0131】

続けて、ストレージ制御装置 101 は、まだグループ分けしていない LBA のうちの先頭となる LBA2000 を選択する。そして、ストレージ制御装置 101 は、512 [K バイト] を超えず、選択した LBA を先頭とした、LBA2000 - 2099、4230 - 4529、6328 - 6345 を G2 にグループ分けする。全ての LBA をグループ分けしたため、ストレージ制御装置 101 は、リフレッシュ処理を終了する。

20

【0132】

リフレッシュ処理後、ストレージ制御装置 101 は、ホットスペア 202 を、ストレージ装置 102 - 1 が属する RAID グループに追加するとともに、ストレージ装置 102 - 1 を RAID グループから外して新たなホットスペアに設定する。また、ストレージ制御装置 101 は、移行後のホットスペア 202 の記憶内容を、ストレージ装置 102 - 1 に戻し、RAID グループを変更せずにそのままにしてもよい。

【0133】

これにより、ストレージ制御装置 101 が予測するグループと LBA 割り当ての関係と、ホットスペア 202 が管理するブロックと LBA 割り当てとの内容が一致して、RD 防止コピーが発生する LBA の予測精度が向上する。図 12 の例では、ホットスペア 202 のグループと LBA 割り当ての予測を示すブロック予測テーブル 111 - h s の内容と、ホットスペア 202 が管理するブロックと LBA 割り当てとの内容を示す表 1201 の内容が一致する。

30

【0134】

リフレッシュ処理を実行する契機としては、ストレージ制御装置 101 は、たとえば、定期的にリフレッシュ処理を実行する。また、ストレージ制御装置 101 は、LBA の割当内容が乖離したことを判断した場合にリフレッシュ処理を実行してもよい。LBA の割当内容が乖離したことを判断する方法として、ストレージ制御装置 101 は、たとえば、上述した第 1 の RD 防止コピー実行判断方法および第 2 の RD 防止コピー実行判断方法を組み合わせることにより、LBA の割当内容が乖離したことを判断する。

40

【0135】

具体的には、第 1 の RD 防止コピー実行判断方法によりストレージ装置 102 のあるグループの読出回数が RD 防止コピー閾値に到達し、ストレージ制御装置 101 は、RD 防止コピーを引き起こすコマンドを第 2 の回数分、ストレージ装置 102 に発行したとする。コマンド発行後、第 2 の RD 防止コピー実行判断方法により、ストレージ制御装置 101 は、ストレージ装置 102 からの応答が遅延しなかった場合に、ストレージ装置 102 において RD 防止コピーが行われておらず、LBA の割当内容が乖離したと判断する。LBA の割当内容が乖離したと判断した場合、ストレージ制御装置 101 は、リフレッシュ

50

処理を実行する。また、ストレージ制御装置 101 は、LBA の割当内容が乖離したと判断した回数が所定値を超えた場合に、リフレッシュ処理を実行してもよい。

#### 【0136】

具体的に、図 12 のリフレッシュ処理を行う前の状態を用いて、LBA の割当内容が乖離したことを判断する例について説明する。ここで、RD 防止コピー閾値を 7900 [回] とし、ストレージ装置 102 が RD 防止コピーを行う読出回数を 8000 [回] とする。そして、図 12 に示す表 1102、ブロック予測テーブル 111-1 の状態となった以降、リフレッシュ処理を行わず、かつ、ホストサーバ 201 から、ストレージ装置 102-1 の LBA 1008-1099 の読出要求が、5533 回あったとする。

#### 【0137】

上述した前提となる場合に、ストレージ制御装置 101 は、ブロック予測テーブル 111-1 の LBA 1008-1099 が含まれる G2 の読出回数を 5533 回インクリメントする。インクリメントした結果、G2 の読出回数は、 $2367 + 5533 = 7900$  となり、RD 防止コピー閾値に到達したため、ストレージ制御装置 101 は、ストレージ装置 102-1 に対して、RD 防止コピーを引き起こすコマンドを 100 [回] 発行する。

#### 【0138】

しかしながら、ストレージ装置 102-1 は、LBA 1008-1099 の読出要求が 5533 回あった時点で、B1 の読出回数が 5533 [回] であると管理している。したがって、ストレージ装置 102 は、RD 防止コピーを引き起こすコマンドを 100 [回] 受け付けても、B1 の読出回数が 8000 [回] に到達しないため、B1 に対する RD 防止コピーを行わない。このため、ストレージ装置 102-1 は、ストレージ制御装置 101 に対して、RD 防止コピーを引き起こすコマンドに対する応答を遅延せずに行うことになる。このように、ストレージ制御装置 101 は、RD 防止コピーを引き起こすコマンドの発行回数と、RD 防止コピーを引き起こすコマンドの発行に対する応答時間を確認することにより、LBA の割当内容が乖離したことを判断することができる。

#### 【0139】

LBA の割当内容が乖離したと判断したときにリフレッシュ処理を実行することにより、ストレージ制御装置 101 は、定期的にリフレッシュ処理を実行する場合に比べて効率的にリフレッシュ処理を実行することができる。具体的に、ストレージ制御装置 101 は、応答遅延が発生すればリフレッシュ処理を行い応答遅延を抑制することができ、応答遅延が発生しなければ移行先のストレージ装置 102 の書込回数を減らすことができる。ここで、SSD 内の浮遊ゲートに書き込める回数には上限があり、ある一定回数以上書き込みを行うと、浮遊ゲート内の絶縁膜が劣化して、メモリセルは情報を記憶することができなくなる。したがって、書込回数を減らすことにより、ストレージ制御装置 101 は、移行先のストレージ装置 102 が故障するまでの期間を延ばすことができる。

#### 【0140】

次に、図 13 ~ 図 17 を用いて、ストレージ制御装置 101 が実行するフローチャートについて説明する。

#### 【0141】

図 13 は、ストレージ装置制御処理手順の一例を示すフローチャートである。ストレージ装置制御処理は、ホストサーバ 201 のアクセスコマンドに応じてストレージ装置 102 を制御する処理である。

#### 【0142】

ストレージ制御装置 101 は、ホストサーバ 201 から、アクセスコマンドを受け付ける (ステップ S1301)。次に、ストレージ制御装置 101 は、アクセスコマンドの種別が次に示す要求のいずれに一致するかを判断する (ステップ S1302)。次に示す要求は、書込要求と、読出要求と、である。アクセスコマンドの種別が書込要求である場合 (ステップ S1302: 書込要求)、ストレージ制御装置 101 は、書込時処理を実行する (ステップ S1303)。書込時処理の詳細については、図 14 および図 15 で後述する。

10

20

30

40

50



## 【 0 1 4 3 】

アクセスコマンドの種別が読出要求である場合（ステップ S 1 3 0 2：読出要求）、ストレージ制御装置 1 0 1 は、続けて、R D 防止処理フラグが次に示す識別子のいずれに一致するかを判断する（ステップ S 1 3 0 4）。次に示す識別子は、「処理していない」と、「処理中」と、である。R D 防止処理フラグが「処理していない」である場合（ステップ S 1 3 0 4：“処理していない”）、ストレージ制御装置 1 0 1 は、ブロック予測テーブル 1 1 1 の読出先の L B A が含まれるグループの読出回数が R D 防止コピー閾値以上か否かを判断する（ステップ S 1 3 0 5）。

## 【 0 1 4 4 】

ブロック予測テーブル 1 1 1 の読出先の L B A が含まれるグループの読出回数が R D 防止コピー閾値未満である場合（ステップ S 1 3 0 5：N o）、ストレージ制御装置 1 0 1 は、読出先のストレージ装置 1 0 2 に読出要求を発行する（ステップ S 1 3 0 6）。発行した結果、ストレージ制御装置 1 0 1 は、読出要求に対するデータを得る。次に、ストレージ制御装置 1 0 1 は、ブロック予測テーブルのアクセス対象 L B A を含むグループの読出回数をインクリメントする（ステップ S 1 3 0 7）。

## 【 0 1 4 5 】

ブロック予測テーブルの読出先の L B A が含まれるグループの読出回数が R D 防止コピー閾値以上である場合（ステップ S 1 3 0 5：Y e s）、ストレージ制御装置 1 0 1 は、読出先のストレージ装置への読出要求の発行を抑止する（ステップ S 1 3 0 8）。次に、ストレージ制御装置 1 0 1 は、R D 防止コピー引起し処理を実行する（ステップ S 1 3 0 9）。R D 防止コピー引起し処理の詳細については、図 1 6 で後述する。

## 【 0 1 4 6 】

ここで、ステップ S 1 3 0 9 と、後続するステップ S 1 3 1 0 との実行の関係において、ストレージ制御装置 1 0 1 は、ステップ S 1 3 0 9 の処理開始後に、ステップ S 1 3 1 0 の処理を実行する。たとえば、ストレージ制御装置 1 0 1 は、R D 防止コピー引起し処理を実行するスレッドを起動した後、ストレージ装置制御処理を実行するスレッドと、R D 防止コピー引起し処理を実行するスレッドとをマルチスレッドにより並列に処理してもよい。また、C P U 3 0 1 がマルチコアである場合、マルチコアのうちのあるコアが、ストレージ装置制御処理を実行するスレッドを実行し、他のコアが R D 防止コピー引起し処理を実行するスレッドを実行してもよい。

## 【 0 1 4 7 】

ステップ S 1 3 0 9 の処理開始後、ストレージ制御装置 1 0 1 は、R A I D グループ内の読出先のストレージ装置 1 0 2 以外のストレージ装置 1 0 2 に、読出要求に応じた読出データを復元可能なデータを記憶する記憶領域を示す L B A の読出要求を発行する（ステップ S 1 3 1 0）。ステップ S 1 3 1 0 における読出要求は、ホストサーバ 2 0 1 からの読出要求である。また、R D 防止処理フラグが「処理中」である場合（ステップ S 1 3 0 4：“処理中”）、ストレージ制御装置 1 0 1 は、ステップ S 1 3 1 0 の処理を実行する。

## 【 0 1 4 8 】

ステップ S 1 3 1 0 の処理について、たとえば、読出先のストレージ装置 1 0 2 が R A I D 1 の R A I D グループに属する場合、ストレージ制御装置 1 0 1 は、R A I D グループのミラーリング先のストレージ装置 1 0 2 に読出要求を発行する。そして、ストレージ制御装置 1 0 1 は、ミラーリング先のストレージ装置 1 0 2 から、ホストサーバ 2 0 1 が発行した読出要求に応じたデータと同一内容の読出データを取得する。

## 【 0 1 4 9 】

また、読出先のストレージ装置 1 0 2 が R A I D 5 の R A I D グループに属する場合、ストレージ制御装置 1 0 1 は、R A I D 5 の R A I D グループの読出先のストレージ装置 1 0 2 以外の全てのストレージ装置 1 0 2 に読出要求を発行する。そして、ストレージ制御装置 1 0 1 は、パリティデータと、パリティデータを生成する際に用いたデータ群のうち、読出先のストレージ装置 1 0 2 に格納した読出データ以外のデータを取得する。次に

10

20

30

40

50

、ストレージ制御装置 101 は、取得したパリティデータと、読出先のストレージ装置 102 に格納した読出データ以外のデータとから、読出先のストレージ装置 102 に格納した読出データを復元する。

【0150】

ステップ S1307、または、ステップ S1310 の処理終了後、ストレージ制御装置 101 は、ホストサーバ 201 へ、読出要求に応じた読出データを通知する（ステップ S1311）。ステップ S1303、または、ステップ S1311 の処理終了後、ストレージ制御装置 101 は、ステップ S1301 の処理に移行する。ストレージ装置制御処理を実行することにより、ストレージ制御装置 101 は、ホストサーバ 201 のアクセスコマンドに応じてストレージ装置 102 を制御することができる。

10

【0151】

図 14 は、書込時処理手順の一例を示すフローチャート（その 1）である。また、図 15 は、書込時処理手順の一例を示すフローチャート（その 2）である。書込時処理は、ホストサーバ 201 からのアクセスコマンドの種別が書込要求である場合に実行する処理である。

【0152】

ストレージ制御装置 101 は、タイマが示す時間が所定間隔  $d_i$  を超えたか否かを判断する（ステップ S1401）。このとき、書込先のストレージ装置 102 が、LBA 割り当てがない、まだ何も書き込まれていないストレージ装置 102 である場合、タイマがスタートしていない。この場合、ストレージ制御装置 101 は、ステップ S1401：No

20

【0153】

タイマが示す時間が所定間隔  $d_i$  を超えていない場合（ステップ S1401：No）、ストレージ制御装置 101 は、書込先のストレージ装置 102 のブロック予測テーブル 111 に、書込先の LBA があるか否かを判断する（ステップ S1402）。書込先の LBA がない場合（ステップ S1402：No）、ホストサーバ 201 からの書込要求は新規書込要求であるとして、ストレージ制御装置 101 は、割当中のグループを、書込先のグループに設定する（ステップ S1403）。次に、ストレージ制御装置 101 は、書込先のグループが書込先の LBA を含むようにグループ分けする（ステップ S1404）。ステップ S1404 の処理終了後、ストレージ制御装置 101 は、グループ分けして生成したグループ情報を、ブロック予測テーブル 111 に格納する。一方、書込先の LBA がある場合（ステップ S1402：Yes）、ホストサーバ 201 からの書込要求は上書き書込要求であるとして、ストレージ制御装置 101 は、書込先の LBA が含まれるグループを、書込先のグループに設定する（ステップ S1405）。

30

【0154】

タイマが示す時間が所定間隔  $d_i$  を超えた場合（ステップ S1401：Yes）、ストレージ制御装置 101 は、書込先のストレージ装置 102 のブロック予測テーブル 111 に、書込先の LBA があるか否かを判断する（ステップ S1406）。書込先の LBA がない場合（ステップ S1406：No）、ホストサーバ 201 からの書込要求は新規書込要求であるとして、ストレージ制御装置 101 は、割当中のグループを、現在割当中のグループの次のグループに設定する（ステップ S1407）。次に、ストレージ制御装置 101 は、割当中のグループを、書込先のグループに設定する（ステップ S1408）。続けて、ストレージ制御装置 101 は、書込先のグループが書込先の LBA を含むようにグループ分けする（ステップ S1409）。ステップ S1409 の処理終了後、ストレージ制御装置 101 は、グループ分けして生成したグループ情報を、ブロック予測テーブル 111 に格納する。一方、書込先の LBA がある場合（ステップ S1406：Yes）、ストレージ制御装置 101 は、ホストサーバ 201 からの書込要求は上書き書込要求であるとして、書込先の LBA が含まれるグループを、書込先のグループに設定する（ステップ S1410）。

40

【0155】

50

ステップS 1 4 0 5、ステップS 1 4 1 0の処理において、ストレージ制御装置1 0 1は、書込先のL B Aを、同一グループに含まれるL B Aに対応するデータサイズがブロックデータサイズを超えないようにグループ分けする。書込先のL B Aがブロックデータサイズに収まらない場合、ストレージ制御装置1 0 1は、書込先のL B Aを、L B A順に複数のグループにグループ分けする。

【0 1 5 6】

ステップS 1 4 0 4、ステップS 1 4 0 5、ステップS 1 4 0 9、ステップS 1 4 1 0のうちのいずれかの処理終了後、ストレージ制御装置1 0 1は、書込先のグループの読出回数を0に初期化する(ステップS 1 4 1 1)。ホストサーバ2 0 1からの書込要求が上書書込要求である場合、ステップS 1 4 1 1を実行することにより、ストレージ制御装置1 0 1は、書込要求に含まれるL B Aが含まれるブロックの読出回数を初期化することになる。

10

【0 1 5 7】

ステップS 1 4 1 1の処理終了後、ストレージ制御装置1 0 1は、R D防止処理フラグが次に示す識別子のいずれに一致するかを判断する(ステップS 1 5 0 1)。次に示す識別子は、「処理していない」と、「処理中」と、である。R D防止処理フラグが「処理していない」である場合(ステップS 1 5 0 1: “処理していない”)、ストレージ制御装置1 0 1は、書込先のストレージ装置1 0 2へ書込データを含む書込要求を発行する(ステップS 1 5 0 2)。

【0 1 5 8】

20

一方、R D防止処理フラグが「処理中」である場合(ステップS 1 5 0 1: “処理中”)、ストレージ制御装置1 0 1は、続けて、書込先のL B AがR D防止コピーの対象のグループに含まれるか否かを判断する(ステップS 1 5 0 3)。

【0 1 5 9】

書込先のL B AがR D防止コピーの対象のグループに含まれない場合(ステップS 1 5 0 3: N o)、ストレージ制御装置1 0 1は、R D防止コピーを引き起こすコマンドの発行を一時停止する(ステップS 1 5 0 4)。次に、ストレージ制御装置1 0 1は、書込先のストレージ装置1 0 2へ書込データを含む書込要求を発行する(ステップS 1 5 0 5)。続けて、ストレージ制御装置1 0 1は、R D防止コピーを引き起こすコマンドの発行を再開する(ステップS 1 5 0 6)。

30

【0 1 6 0】

一方、書込先のL B AがR D防止コピーの対象のグループに含まれる場合(ステップS 1 5 0 3: Y e s)、ストレージ制御装置1 0 1は、R D防止コピーを引き起こすコマンドの発行を停止させる(ステップS 1 5 0 7)。次に、ストレージ制御装置1 0 1は、書込先のストレージ装置へ書込データを含む書込要求を発行する(ステップS 1 5 0 8)。続けて、ストレージ制御装置1 0 1は、R D防止処理フラグを、「処理していない」に変更する(ステップS 1 5 0 9)。

【0 1 6 1】

ステップS 1 5 0 2、ステップS 1 5 0 6、ステップS 1 5 0 9のうちのいずれかの処理終了後、ストレージ制御装置1 0 1は、タイマをスタートする(ステップS 1 5 1 0)。ステップS 1 5 1 0の処理終了後、ストレージ制御装置1 0 1は、書込時処理を終了する。書込時処理を実行することにより、ストレージ制御装置1 0 1は、ホストサーバ2 0 1からの書込要求に対応してストレージ装置1 0 2を制御することができる。

40

【0 1 6 2】

図1 6は、R D防止コピー引起し処理手順の一例を示すフローチャートである。R D防止コピー引起し処理は、ストレージ装置1 0 2にR D防止コピーを起こさせる処理である。図1 6に示す処理について、ブロック予測テーブル1 1 1とR D防止処理フラグとは、図1 3における読出先のストレージ装置1 0 2におけるブロック予測テーブル1 1 1とR D防止処理フラグとを指す。説明の簡略化のため、図1 6において、ブロック予測テーブル1 1 1とR D防止処理フラグには、「読出先のストレージ装置1 0 2における」を省略

50

して記載する。

【0163】

ストレージ制御装置101は、RD防止処理フラグを、「処理中」に変更する(ステップS1601)。次に、ストレージ制御装置101は、読出先のストレージ装置へ、RD防止コピーを引き起こすコマンドを発行する(ステップS1602)。続けて、ストレージ制御装置101は、RD防止コピーが行われたか否かを判断する(ステップS1603)。RD防止コピーが行われたか否かを判断する方法は、上述した第1のRD防止コピー実行判断方法と、第2のRD防止コピー実行判断方法とがある。

【0164】

RD防止コピーが行われていない場合(ステップS1603:No)、ストレージ制御装置101は、ステップS1602の処理に移行する。RD防止コピーが行われた場合(ステップS1603:Yes)、ストレージ制御装置101は、ブロック予測テーブルの読出先のLBAが含まれるグループの読出回数の値を0に初期化する(ステップS1604)。次に、続けて、ストレージ制御装置101は、RD防止処理フラグを、「処理していない」に変更する(ステップS1605)。ステップS1605の処理終了後、ストレージ制御装置101は、RD防止コピー引起し処理を終了する。RD防止コピー引起し処理を実行することにより、ストレージ制御装置101は、ストレージ装置102にRD防止コピーを起こさせることができる。

【0165】

図17は、リフレッシュ処理手順の一例を示すフローチャートである。リフレッシュ処理は、移行元のストレージ装置102が記憶する移行対象データを移行先のストレージ装置に移行させて、グループとLBAとの関係を再定義する処理である。

【0166】

ストレージ制御装置101は、移行元のストレージ装置102の全てのLBAを、移行先のストレージ装置102に移行するデータ移行要求を、移行元のストレージ装置に発行する(ステップS1701)。次に、ストレージ制御装置101は、移行元のストレージ装置102のブロック予測テーブル111の先頭LBAを選択する(ステップS1702)。続けて、ストレージ制御装置101は、移行先のストレージ装置102の割当中のグループを、先頭のグループに設定する(ステップS1703)。

【0167】

次に、ストレージ制御装置101は、移行先のストレージ装置102のブロック予測テーブル111の割当中のグループに、選択したLBAを先頭として1ブロック分のLBAを含むようにグループ分けする(ステップS1704)。続けて、ストレージ制御装置101は、移行元のストレージ装置102の全てのLBAをグループ分けしたか否かを判断する(ステップS1705)。

【0168】

移行元のストレージ装置102のLBAのうち、グループ分けしていないLBAがある場合(ステップS1705:No)、ストレージ制御装置101は、移行元のストレージ装置102のLBAのうち、まだグループ分けしていない先頭のLBAを選択する(ステップS1706)。続けて、ストレージ制御装置101は、割当中のグループを、現在割当中のグループの次のグループに設定する(ステップS1707)。ステップS1707の処理終了後、ストレージ制御装置101は、ステップS1704の処理に移行する。

【0169】

一方、移行元のストレージ装置102の全てのLBAをグループ分けした場合(ステップS1705:Yes)、ストレージ制御装置101は、リフレッシュ処理を終了する。リフレッシュ処理を実行することにより、ストレージ制御装置101は、ストレージ装置102が管理するグループとLBA割り当ての関係と、ストレージ制御装置101が予測するグループとLBA割り当ての関係との乖離を解消することができる。

【0170】

以上説明したように、ストレージ制御装置101によれば、書込要求のLBAをブロッ

10

20

30

40

50

ク単位でグループ分けしたグループの読出回数を計数し、計数したグループの読出回数に基づいて元データか冗長データかを読み出す読出要求をストレージ装置 102 に発行する。これにより、ストレージ制御装置 101 は、読み出しを行うと R D 防止コピーが発生するであろう L B A を予測することができる。そして、ストレージ制御装置 101 は、ホストサーバ 201 からの読出要求に対する応答性能の低下を抑制することができる。また、ストレージ制御装置 101 は、ストレージ装置 102 の内部処理とストレージ制御装置 101 からの読出要求とのタイミングが重なることを防ぎ、安定した応答時間を得ることができる。

【0171】

また、ストレージ制御装置 101 によれば、ストレージ制御装置 101 が受け付けた書込要求に応じて、書込要求の L B A をブロック単位でグループ分けしてもよい。これにより、ストレージ制御装置 101 は、ストレージ装置 102 への書込要求に追従してグループ分けが行えるため、R D 防止コピーが発生するであろう L B A の予測精度の低下を抑制することができる。

【0172】

また、ストレージ制御装置 101 によれば、ストレージ装置 102 への読出要求に応じて、読出要求に含まれる読出先の L B A を含むグループに含まれる L B A の書込要求を発行したことに応じて、グループの読出回数を初期化してもよい。S S D コントローラは、上書き書込が行われると該当するブロックの読出回数を初期化するため、ストレージ制御装置 101 は、予測精度の低下を抑制することができる。

【0173】

また、ストレージ制御装置 101 によれば、グループの読出回数が R D 防止コピー閾値以上であることに応じて、R D 防止コピーを引き起こすコマンドをストレージ装置 102 に発行して、グループの読出回数を初期化してもよい。ここで、ストレージ制御装置 101 は、ストレージ装置 102 に R D 防止コピーを引き起こすコマンドを、ストレージ装置 102 の仕様に応じて設定された規定回数から R D 防止コピー閾値を減じた回数分発行する。または、ストレージ制御装置 101 は、ストレージ装置 102 が R D 防止コピーを引き起こすコマンドを、コマンドに対する応答時間が所定時間を超えるまで発行する。これにより、ストレージ制御装置 101 は、R D 防止コピーが発生するであろう L B A に対し R D 防止コピーを事前に起こしておくことになる。したがって、ホストサーバ 201 からの読出要求を受け付ける前に R D 防止コピーが完了していれば、ホストサーバ 201 からの読出要求に対する応答遅延が発生しなくなる。これにより、ストレージ制御装置 101 は、ホストサーバ 201 からの読出要求に対する応答性能の低下を抑制することができる。

【0174】

また、ストレージ制御装置 101 によれば、R D 防止コピーを引き起こすコマンドを発行する間に、コマンドの L B A を含むグループに含まれる L B A の書込要求を発行したことに応じて、コマンドの発行を停止するとともにグループの読出回数を初期化してもよい。S S D コントローラは、上書き書込が行われると該当するブロックの読出回数を初期化するため、ストレージ制御装置 101 は、読出回数の予測精度を向上することができる。さらに、上書き書込が行われれば R D 防止コピーを行わなくてよいので、コマンドの発行を停止することにより、ストレージ制御装置 101 は、無駄なコマンドの発行を行わなくて済む。

【0175】

また、ストレージ制御装置 101 によれば、同一グループに含まれる L B A を含む書込要求のストレージ装置 102 への発行間隔が所定間隔  $d_i$  を超えないようにグループ分けしてもよい。これにより、ストレージ制御装置 101 は、ストレージ装置 102 が行うブロックと L B A 割り当てとの関係の予測精度を向上することができる。

【0176】

また、本実施の形態にかかるストレージ装置 102 の半導体メモリは、N A N D 型のフ

10

20

30

40

50

ラッシュメモリであってもよい。NAND型のフラッシュメモリは、NOR型のフラッシュメモリと比較してリードディスタurbが発生し易いため、NAND型のフラッシュメモリを制御するSSDコントローラは、RD防止コピーを行う回数が増えることになる。したがって、NAND型のフラッシュメモリを有するストレージ装置102に対して本実施の形態におけるストレージ装置の制御方法を実行することにより、多く発生する応答性能が遅延する状態を起こさないようにして、応答性能の低下を抑制することができる。

【0177】

また、ストレージ制御装置101によれば、グループの読出回数がRD防止コピー閾値以上であれば、グループに含まれるLBAの読出要求に応じて、読出先のデータを復元可能なデータ先を示すLBAの読出要求を、他のストレージ装置102に発行してもよい。これにより、ストレージ制御装置101は、RD防止コピーが発生するであろうストレージ装置102を避けて、ホストサーバ201の読出要求に応答することができ、読出要求に対する応答性能の低下を抑制することができる。

10

【0178】

また、ストレージ制御装置101によれば、同一の特性を有するストレージ装置102群によりRAIDグループを形成する場合、ストレージ装置102群の各々のRD防止コピー閾値をずらして設定してもよい。これにより、ストレージ制御装置101は、RD防止コピーが行われる期間をストレージ装置102間でずらすことができる。RD防止コピーが行われる期間をずらすことにより、ストレージ制御装置101は、ミラーリングからのデータ取得ができない状態、またはパリティを用いたデータ復元ができない状態になることを避け、読出要求に対する応答性能の低下を抑制することができる。

20

【0179】

また、ストレージ制御装置101によれば、移行元のストレージ装置102が記憶する移行対象データを移行先のストレージ装置102に移行させて、グループとLBAとの関係を再定義してもよい。これにより、ストレージ制御装置101は、ストレージ制御装置101が予測するグループとLBA割り当ての関係と移行先のストレージ装置102が管理するブロックとLBA割り当てとの内容が一致して、RD防止コピーが発生するLBAの予測精度が向上する。

【0180】

また、ストレージ制御装置101によれば、移行先のストレージ装置102に対してもグループ分けや、グループの読出回数の計数を行うことにより、移行先のストレージ装置102に対しても、RD防止コピーが発生するLBAの予測を行うことができる。

30

【0181】

なお、本実施の形態で説明したストレージ装置の制御方法は、予め用意されたプログラムをパーソナル・コンピュータやワークステーション等のコンピュータで実行することにより実現することができる。本ストレージ装置の制御プログラムは、ハードディスク、フレキシブルディスク、光ディスク等のコンピュータで読み取り可能な記録媒体に記録され、コンピュータによって記録媒体から読み出されることによって実行される。また本ストレージ装置の制御プログラムは、インターネット等のネットワークを介して配布してもよい。

40

【0182】

上述した実施の形態に関し、さらに以下の付記を開示する。

【0183】

(付記1) ストレージ装置の記憶領域に対して所定のデータサイズのブロック単位でデータを冗長化して記憶させる制御を行うストレージ制御装置であって、

データの書込要求から特定される書込先の論理アドレスを、前記ブロック単位に対応付けてグループ分けしたグループ情報を記憶する記憶部と、

データの読出要求に応じて、前記グループ情報に基づき前記読出要求から特定される読出先の論理アドレスを含むグループの読出回数を計数し、計数した前記グループの読出回数に基づき前記読出先の論理アドレスを含む読出要求、または、前記読出先の論理アドレ

50

スのデータに対応する冗長データの記憶先の論理アドレスを含む読出要求を発行する制御部と、

を有することを特徴とするストレージ制御装置。

【0184】

(付記2) 前記記憶領域は、規定回数の読み出しが行われるとブロックのデータを他のブロックへコピーするコピー制御が行われる半導体メモリの記憶領域であることを特徴とする付記1に記載のストレージ制御装置。

【0185】

(付記3) 前記ストレージ制御装置は、前記半導体メモリを含む前記ストレージ装置と、当該ストレージ装置が記憶するデータの冗長データを記憶する前記ストレージ装置とは異なる他のストレージ装置の制御を行うものであり、

10

前記制御部は、

計数した前記グループの読出回数が前記規定回数より小さい第1の回数となった後に、データの読出要求に応じて、当該読出要求から特定される読出先の論理アドレスが前記グループに含まれる場合に、当該読出先の論理アドレスのデータに対応する冗長データの記憶先の論理アドレスを含む読出要求を、前記他のストレージ装置に発行する、

ことを特徴とする付記2に記載のストレージ制御装置。

【0186】

(付記4) 前記制御部は、

計数した前記グループの読出回数が前記規定回数より小さい第1の回数となったことに  
応じて、前記グループに含まれる論理アドレスを含む読出要求を、前記規定回数から前記  
第1の回数を減じた第2の回数、または、当該読出要求に対する応答時間が所定時間を超  
えるまで前記ストレージ装置に発行し、前記グループの読出回数を初期化することを特徴  
とする付記2または3に記載のストレージ制御装置。

20

【0187】

(付記5) 前記制御部は、

前記グループに含まれる論理アドレスを含む読出要求を、前記第2の回数、または、当  
該読出要求に対する応答時間が所定時間を超えるまで前記ストレージ装置に発行する間に  
、前記グループに含まれる論理アドレスを含む書込要求を前記ストレージ装置に発行した  
ことに応じて、前記グループに含まれる論理アドレスを含む読出要求の発行を停止すると  
ともに、前記グループの読出回数を初期化することを特徴とする付記4に記載のストレ  
ージ制御装置。

30

【0188】

(付記6) 前記制御部は、

前記グループに含まれる論理アドレスを含む書込要求を前記ストレージ装置に発行した  
ことに応じて、前記グループの読出回数を初期化することを特徴とする付記1または2に  
記載のストレージ制御装置。

【0189】

(付記7) 前記記憶部は、

書込要求から特定される前記他のストレージ装置の書込先の論理アドレスを、前記ブロ  
ック単位に対応付けてグループ分けした前記他のストレージ装置のグループ情報を記憶し  
ており、

40

前記制御部は、

受け付けた読出要求に応じて、前記他のストレージ装置のグループ情報に基づき当該読  
出要求から特定される前記他のストレージ装置の読出先の論理アドレスを含むグループの  
読出回数を計数し、計数した当該グループの読出回数が前記規定回数より小さく前記第1  
の回数とは異なる第3の回数となったことに応じて、当該グループに含まれる論理アドレ  
スを含む読出要求を、前記規定回数から前記第3の回数を減じた第4の回数、または、当  
該読出要求に対する応答時間が所定時間を超えるまで前記他のストレージ装置に発行し、  
当該グループの読出回数を初期化することを特徴とする付記3に記載のストレージ制御装

50

置。

【0190】

(付記8) 前記ストレージ制御装置は、さらに、半導体メモリの記憶領域を所定のデータサイズで分割したブロック単位で書き込みを行う移行先のストレージ装置の制御を行うものであり、

前記制御部は、

受け付けた書込要求から特定される前記ストレージ装置の書込先の論理アドレスをグループ分けした複数のグループに含まれる論理アドレスを含むデータ移行要求を、前記移行先のストレージ装置へのデータの移行間隔が所定間隔を超えないように前記ストレージ装置に発行し、

10

発行した前記データ移行要求に含まれる論理アドレスを、前記ブロック単位に対応付けてグループ分けした前記移行先のストレージ装置のグループ情報を生成することを特徴とする付記1～7のいずれか一つに記載のストレージ制御装置。

【0191】

(付記9) 前記制御部は、

受け付けた書込要求に応じて、当該書込要求から特定される前記移行先のストレージ装置の書込先の論理アドレスを、同一グループに含まれる論理アドレスに対応するデータサイズが当該所定のデータサイズを超えないようにグループ分けし、

受け付けた読出要求に応じて、当該読出要求から特定される前記移行先のストレージ装置の読出先の論理アドレスを含むグループの読出回数を計数し、

20

計数した当該グループの読出回数に基づき当該読出先の論理アドレスを含む読出要求、または、当該読出先の論理アドレスのデータに対応する冗長データの記憶先の論理アドレスを含む読出要求を発行する付記8に記載のストレージ制御装置。

【0192】

(付記10) ストレージ装置の記憶領域に対して所定のデータサイズのブロック単位でデータを冗長化して記憶させる制御を行うストレージ制御装置の制御方法であって、

前記ストレージ制御装置が、

データの読出要求に応じて、データの書込要求から特定される書込先の論理アドレスを、前記ブロック単位に対応付けてグループ分けしたグループ情報に基づき前記読出要求から特定される読出先の論理アドレスを含むグループの読出回数を計数し、計数した前記グループの読出回数に基づき前記読出先の論理アドレスを含む読出要求、または、前記読出先の論理アドレスのデータに対応する冗長データの記憶先の論理アドレスを含む読出要求を発行する、

30

処理を実行することを特徴とする制御方法。

【0193】

(付記11) ストレージ装置の記憶領域に対して所定のデータサイズのブロック単位でデータを冗長化して記憶させる制御を行うストレージ制御装置の制御プログラムであって、

前記ストレージ制御装置に、

データの読出要求に応じて、データの書込要求から特定される書込先の論理アドレスを、前記ブロック単位に対応付けてグループ分けしたグループ情報に基づき前記読出要求から特定される読出先の論理アドレスを含むグループの読出回数を計数し、計数した前記グループの読出回数に基づき前記読出先の論理アドレスを含む読出要求、または、前記読出先の論理アドレスのデータに対応する冗長データの記憶先の論理アドレスを含む読出要求を発行する、

40

処理を実行させることを特徴とする制御プログラム。

【符号の説明】

【0194】

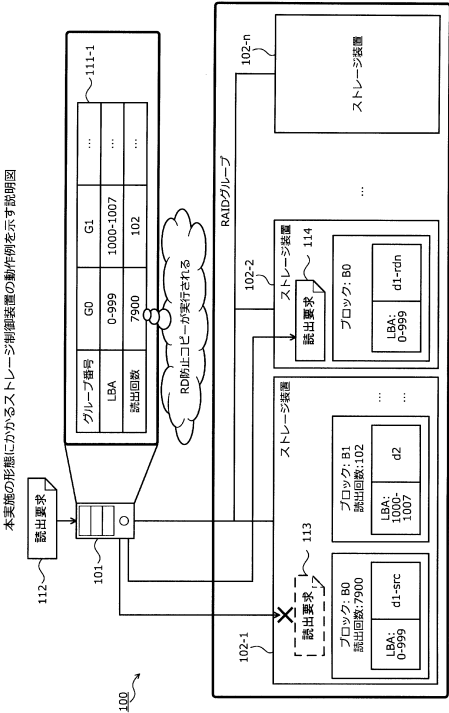
- 100 ストレージシステム
- 101 ストレージ制御装置
- 102 ストレージ装置

50

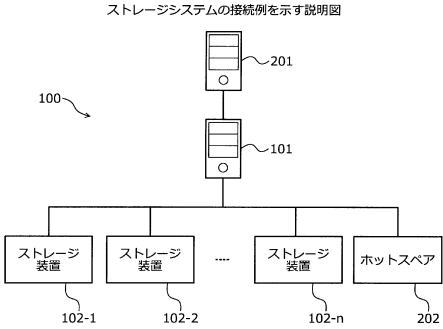


- 1 1 1    ブロック予測テーブル
- 4 0 1    制御部
- 4 0 2    記憶部
- 4 1 1    グループ分け部
- 4 1 2    計数部
- 4 1 3    アクセス制御部

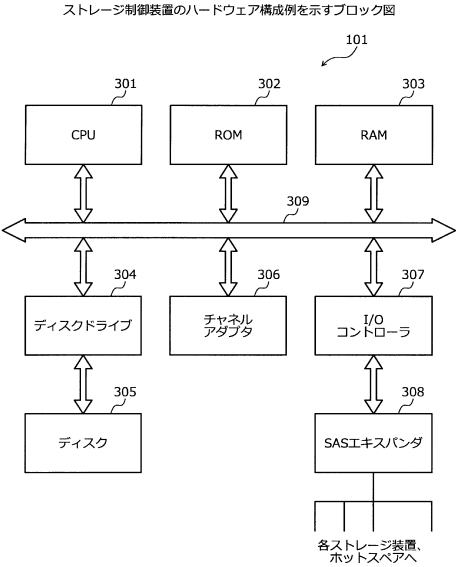
【 図 1 】



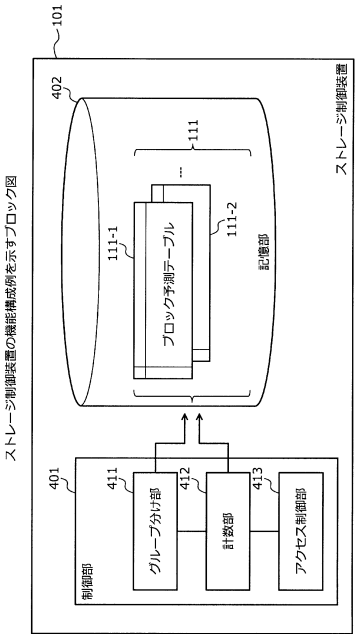
【 図 2 】



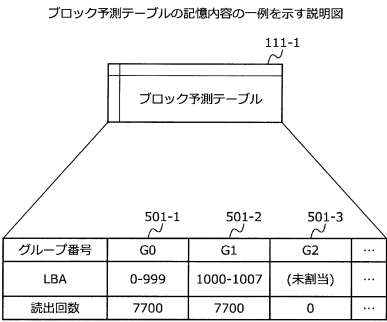
【図 3】



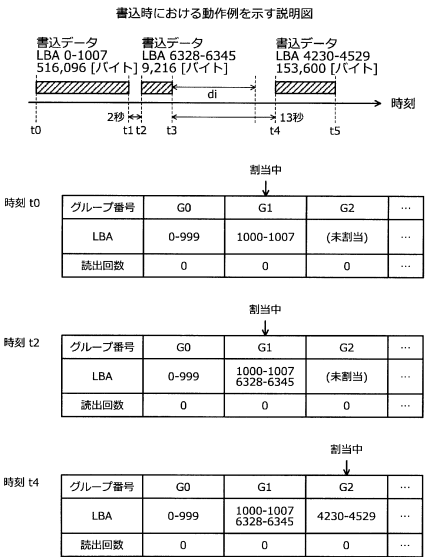
【図 4】



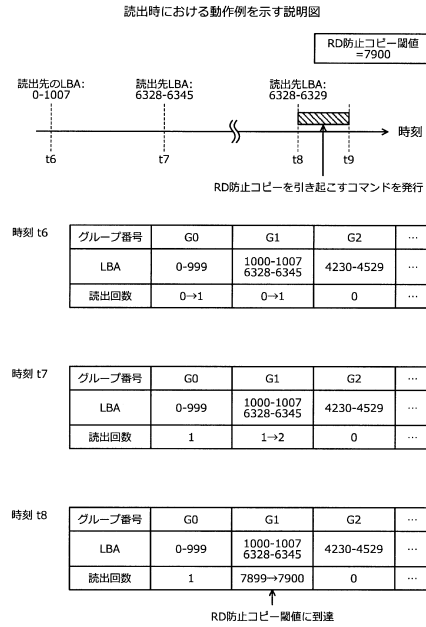
【図 5】



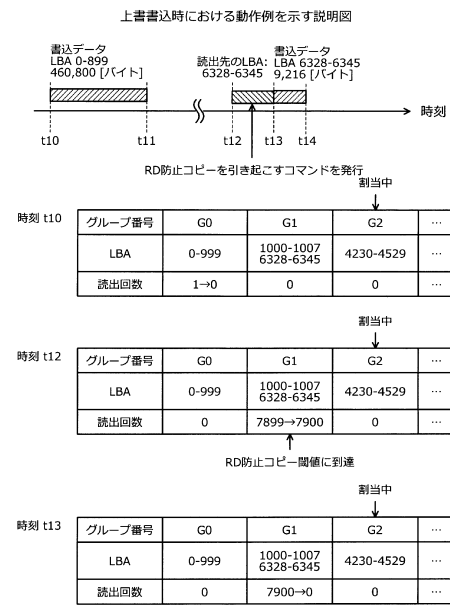
【図 6】



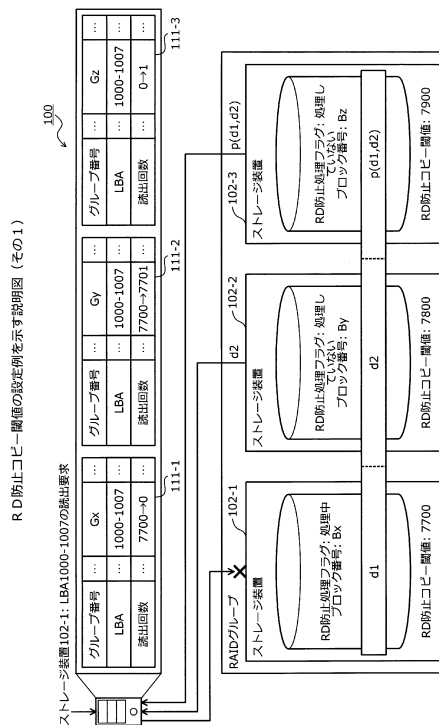
【 図 7 】



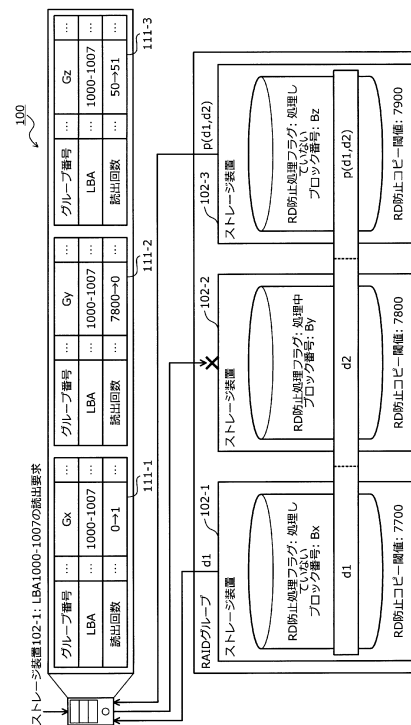
【 図 8 】



【 図 9 】

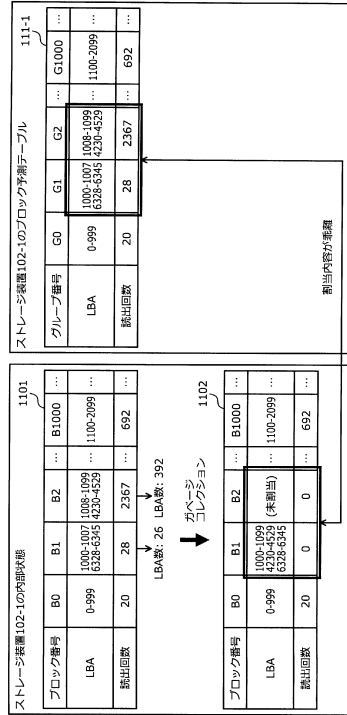


【 図 1 0 】



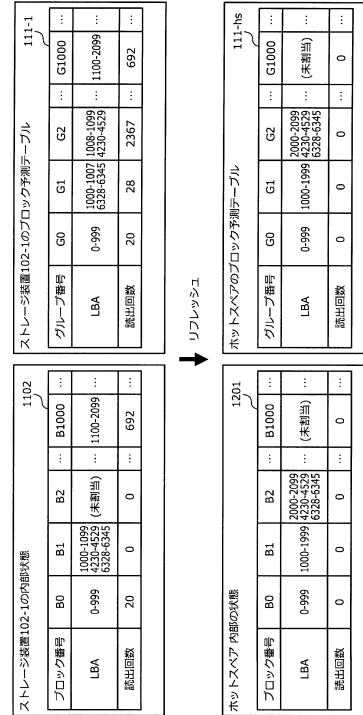
【図 1 1】

ガベージコレクションが行われた際のブロックとグループとの関係の一例を示す説明図



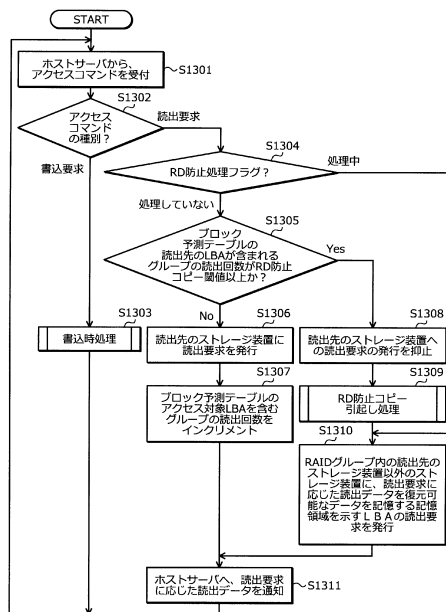
【図 1 2】

リフレッシュ処理の実行前後におけるブロックとグループとの関係の一例を示す説明図



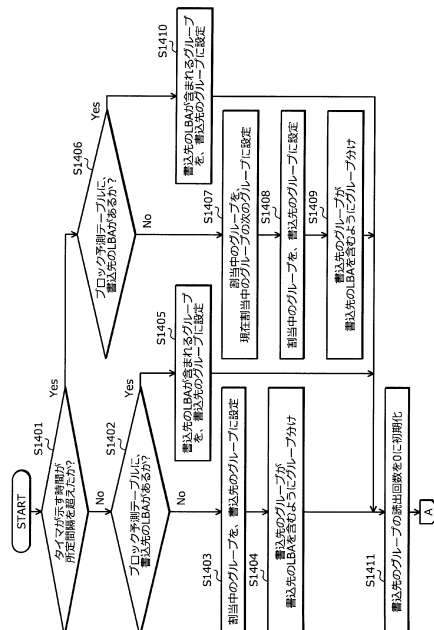
【図 1 3】

ストレージ装置制御処理手順の一例を示すフローチャート



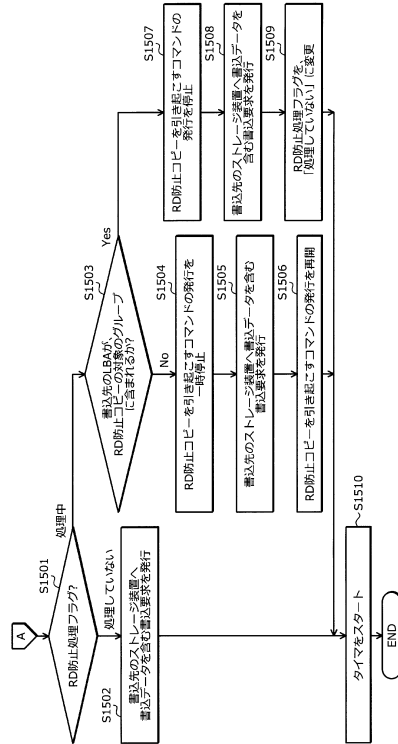
【図 1 4】

書込時処理手順の一例を示すフローチャート (その1)



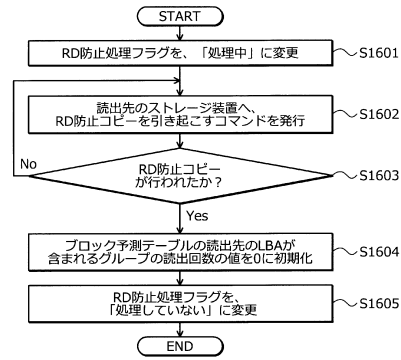
【図 15】

書込時処理手順の一例を示すフローチャート（その2）



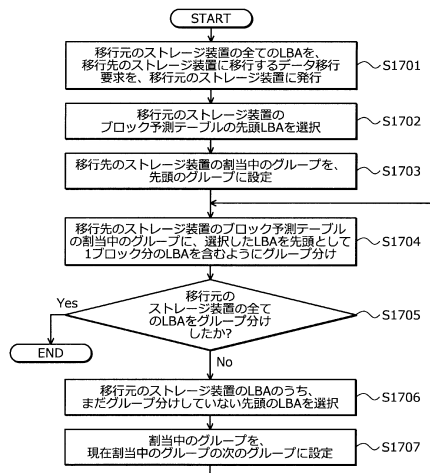
【図 16】

RD防止コピー引き起こし処理手順の一例を示すフローチャート



【図 17】

リフレッシュ処理手順の一例を示すフローチャート



---

フロントページの続き

(56)参考文献 特開 2012 - 238259 (JP, A)  
特開 2012 - 203642 (JP, A)  
米国特許出願公開第 2009 / 0193174 (US, A1)  
米国特許出願公開第 2013 / 0262750 (US, A1)  
米国特許出願公開第 2011 / 0038203 (US, A1)  
特開 2009 - 037317 (JP, A)  
特開 2008 - 287404 (JP, A)

(58)調査した分野(Int.Cl., DB名)  
G06F3/06 - 3/08  
G06F12/00 - 12/06、12/16  
G06F13/16 - 13/18