

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2010-15195
(P2010-15195A)

(43) 公開日 平成22年1月21日(2010.1.21)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 12/16 (2006.01)	G06F 12/16 320L	5B018
G06F 3/06 (2006.01)	G06F 12/16 320M	5B065
	G06F 3/06 540	
	G06F 3/06 305F	
	G06F 3/06 305C	

審査請求 未請求 請求項の数 7 O L (全 17 頁)

(21) 出願番号 特願2008-171800 (P2008-171800)
(22) 出願日 平成20年6月30日 (2008.6.30)

(71) 出願人 000003078
株式会社東芝
東京都港区芝浦一丁目1番1号
(74) 代理人 100089118
弁理士 酒井 宏明
(72) 発明者 福富 和弘
東京都港区芝浦一丁目1番1号 株式会社東芝内
(72) 発明者 佐藤 英昭
東京都港区芝浦一丁目1番1号 株式会社東芝内
(72) 発明者 菅野 伸一
東京都港区芝浦一丁目1番1号 株式会社東芝内

最終頁に続く

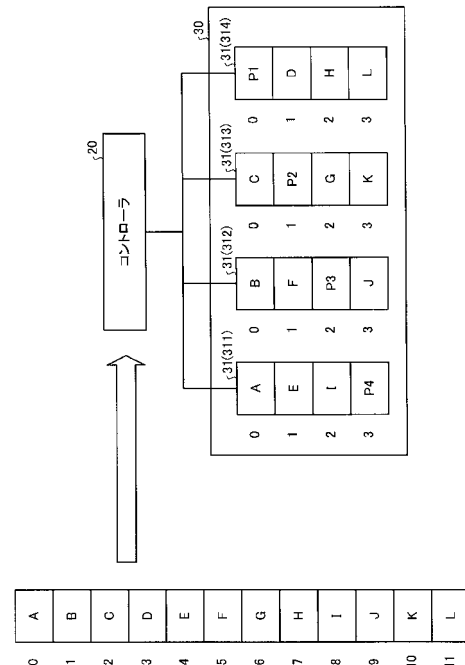
(54) 【発明の名称】 記憶制御装置及び記憶制御方法

(57) 【要約】

【課題】 R A I D構成とされた複数の不揮発性半導体記憶装置を有効に活用することが可能な記憶制御装置及び記憶制御方法を提供する。

【解決手段】 複数の不揮発性半導体記憶装置を用いて、当該不揮発性半導体記憶装置に記憶されるデータを復元可能な R A I Dを構成し、外部から入力されるデータの読み込み要求に応じ、前記 R A I Dを構成する不揮発性半導体記憶装置からデータを読み込み、この読み込み時にエラーが発生した場合には、当該読み込みエラーが発生したデータを復元し、前記読み込みエラーが発生した前記不揮発性半導体記憶装置上の領域に書き戻す。

【選択図】 図 2



【特許請求の範囲】

【請求項 1】

複数の不揮発性半導体記憶装置を接続可能なインタフェースと、
 前記複数の不揮発性半導体記憶装置を用いて、記憶対象となるデータを、当該データを復元可能な復元情報とともに記憶するための R A I D を構成する構成手段と、
 外部から入力されるデータの読み込み要求に応じ、前記 R A I D を構成する不揮発性半導体記憶装置からデータを読み込む読込手段と、
 前記読込手段による読み込み時にエラーが発生したデータを、前記復元情報に基づいて復元する復元手段と、
 前記復元手段により復元されたデータを一時的に保存する保存手段と、
 前記保存手段に保存されたデータを、当該データの読み込みエラーが発生した前記不揮発性半導体記憶装置上の領域に書き込む書込手段と、
 を備えたことを特徴とする記憶制御装置。

10

【請求項 2】

外部から入力されるデータの書き込み要求に応じ、前記 R A I D を構成する不揮発性半導体記憶装置にデータを書き込む書込手段を更に備え、
 前記構成手段は、前記書込手段による書き込み時にエラーが発生した不揮発性半導体記憶装置を、前記 R A I D の構成から除外することを特徴とする請求項 1 に記載の記憶制御装置。

20

【請求項 3】

前記構成手段は、前記書き込みエラーが発生した不揮発性半導体記憶装置以外の他の不揮発性半導体記憶装置により、前記 R A I D の構成を維持することが可能な場合に、前記書き込みエラーが発生した不揮発性半導体記憶装置を、前記 R A I D の構成から除外することを特徴とする請求項 2 に記載の記憶制御装置。

【請求項 4】

前記不揮発性半導体記憶装置は、自己の記録媒体に記憶されたデータのリフレッシュを行うことを特徴とする請求項 1 ~ 3 の何れか一項に記載の記憶制御装置。

【請求項 5】

前記不揮発性半導体記憶装置は、N A N D 型フラッシュメモリを記録媒体として有することを特徴とする請求項 1 ~ 4 の何れか一項に記載の記憶制御装置。

30

【請求項 6】

前記構成手段は、前記複数の不揮発性半導体記憶装置を、R A I D 1、5、6 又はこれらの組み合わせとした構成とすることを特徴とする請求項 1 ~ 3 の何れか一項に記載の記憶制御装置。

【請求項 7】

複数の不揮発性半導体記憶装置を用いて、記憶対象となるデータを、当該データを復元可能な復元情報とともに記憶するための R A I D を構成する記憶制御装置の記憶制御方法であって、

読込手段が、外部から入力されるデータの読み込み要求に応じ、前記 R A I D を構成する不揮発性半導体記憶装置からデータを読み込む読込工程と、

40

復元手段が、前記読込工程で読み込みエラーが発生したデータを、前記復元情報に基づいて復元する復元工程と、

保存手段が、前記復元工程で復元されたデータを一時的に保存する保存工程と、

書込手段が、前記保存工程に保存されたデータを、当該データの読み込みエラーが発生した前記不揮発性半導体記憶装置上の領域に書き込む書込工程と、

を含むことを特徴とする記憶制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、R A I D 構成とされた複数の不揮発性半導体記憶装置を制御する記憶制御装

50

置及び記憶制御方法に関する。

【背景技術】

【0002】

従来、サーバ環境等で使用されるストレージシステムでは、耐障害性・冗長性を向上させるため、複数の磁気ディスク装置を用いることで、RAID (Redundant Array of Independent/Inexpensive Disks) を構成することが行われている (例えば、非特許文献1参照)。例えば、RAID 5の構成では、3つ以上の磁気ディスク装置を使用し、データの復旧のためのパリティを当該データとともに各磁気ディスク装置に分散して記憶することで、データが破損した場合であっても、このパリティを用いることでデータを復旧することが可能である。このようなRAID構成を実現するコントローラでは、ある磁気ディスク装置からデータの読み込みが不可能になると、当該磁気ディスク装置が故障したと判断し、RAID構成から除外することが一般に行われており、除外された磁気ディスク装置を新たな磁気ディスク装置に交換することで、従前のRAID構成を回復することが可能となっている。

10

【0003】

一方、不揮発性の半導体記憶素子を記録媒体として用いたSSD (Solid State Drive) 等の不揮発性半導体記憶装置が存在しており、磁気ディスク装置と同様に補助記憶装置 (二次記憶装置) として使用することも行われている。この、不揮発性半導体記憶装置では、磁気ディスク装置のようにディスクを持たないため、データの読み書きが磁気ディスク装置に比べて高速であり、消費電力を抑えることができるため、サーバ環境での利用も期待されている。

20

【0004】

記録媒体としてNAND型フラッシュメモリが採用された不揮発性半導体記憶装置では、データの読み込み時に生じる電荷変位や自然放電等により、記憶したデータが劣化し、正常に読み込みすることができなくなる可能性がある。そのため、このような不揮発性半導体記憶装置では、記憶されているデータを読み出して誤り訂正を行った後に再びNAND型フラッシュメモリに書き戻すリフレッシュ処理を所定時間毎に行うことで、データの損失を防ぐ機構が備えられている。なお、この場合、記憶素子 (メモリセル) 自体は故障しておらず、データが劣化しているだけであるため、正常なデータを再度書き込むことにより、当該記憶素子を再び使用することが可能となる。

30

【0005】

【非特許文献1】D. Patterson, G. Gibson, and R. Katz. "A Case for Redundant Array of Inexpensive Disks (RAID)", Proceedings of the 1988 ACM SIGMOD, pp.109-116, June 1988.

【発明の開示】

【発明が解決しようとする課題】

【0006】

ところで、劣化の程度によっては上記リフレッシュ処理により全てのデータを復元できるとは限らず、この場合、復元に失敗したデータの読み込みについては、読み込みエラーが発生することになる。そのため、不揮発性半導体記憶装置を用いて上記RAIDを構成すると、従来技術と同様、読み込みエラーが発生した際にはその不揮発性半導体記憶装置が故障と判断されることになる。なお、上述したように読み込みエラーの発生したデータの格納領域に正常なデータを再度書き込むことで、故障と判断された不揮発性半導体記憶装置を復旧できる可能性があるが、上記従来技術は、磁気ディスク装置での使用が前提となるため、当該不揮発性半導体記憶装置を復旧することはできない。また、復旧可能な不揮発性半導体記憶装置であっても、RAID構成から除外され、交換の対象とされてしまう可能性があるため、従来技術では、不揮発性半導体記憶装置を有効に活用できないという問題がある。

40

【0007】

本発明は上記に鑑みてなされたものであって、RAID構成とされた複数の不揮発性半

50

導体記憶装置を有効に活用することが可能な記憶制御装置及び記憶制御方法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】

【0008】

上述した課題を解決し、目的を達成するために、本発明は、複数の不揮発性半導体記憶装置を接続可能なインタフェースと、前記複数の不揮発性半導体記憶装置を用いて、記憶対象となるデータを、当該データを復元可能な復元情報とともに記憶するためのRAIDを構成する構成手段と、外部から入力されるデータの読み込み要求に応じ、前記RAIDを構成する不揮発性半導体記憶装置からデータを読み込む読込手段と、前記読込手段による読み込み時にエラーが発生したデータを、前記復元情報に基づいて復元する復元手段と、前記復元手段により復元されたデータを一時的に保存する保存手段と、前記保存手段に保存されたデータを、当該データの読み込みエラーが発生した前記不揮発性半導体記憶装置上の領域に書き込む書込手段と、を備える。

10

【0009】

また、本発明は、複数の不揮発性半導体記憶装置を用いて、記憶対象となるデータを、当該データを復元可能な復元情報とともに記憶するためのRAIDを構成する記憶制御装置の記憶制御方法であって、読込手段が、外部から入力されるデータの読み込み要求に応じ、前記RAIDを構成する不揮発性半導体記憶装置からデータを読み込む読込工程と、復元手段が、前記読込工程で読み込みエラーが発生したデータを、前記復元情報に基づいて復元する復元工程と、保存手段が、前記復元工程で復元されたデータを一時的に保存する保存工程と、書込手段が、前記保存工程に保存されたデータを、当該データの読み込みエラーが発生した前記不揮発性半導体記憶装置上の領域に書き込む書込工程と、を含む。

20

【発明の効果】

【0010】

本発明によれば、RAID構成とされた不揮発性半導体記憶装置に読み込みエラーが発生した場合であっても、この読み込みエラーが発生したデータを復元し、当該読み込みエラーが発生した不揮発性半導体記憶装置上の領域に書き戻すことで、当該不揮発性半導体記憶装置を正常な状態に復旧させることができるため、RAID構成とされた複数の不揮発性半導体記憶装置を有効に活用することができる。

30

【発明を実施するための最良の形態】

【0011】

以下、添付図面を参照して、本発明の最良な実施形態を詳細に説明する。なお、本発明は以下の記述に限定されるものではなく、本発明の要旨を逸脱しない範囲において適宜変更可能である。

【0012】

[第1の実施形態]

図1は、第1の実施形態にかかるストレージシステム100の概略構成を示したブロック図である。同図に示したように、ストレージシステム100は、ホスト装置10と、コントローラ20と、複数の不揮発性半導体記憶装置31から構成されるストレージ装置30とを有している。

40

【0013】

ホスト装置10は、PC(Personal Computer)等であって、コントローラ20に対してデータの書き込みや読み込みを要求する指示情報を出力する。以下、データの書き込みを要求する指示情報を「書き込み要求」と呼び、データの読み込みを要求する指示情報を「読み込み要求」と呼ぶ。なお、ホスト装置10からコントローラ20に出力される書き込み要求には、少なくとも書き込み対象のデータが含まれているものとし、また、読み込み要求には、読み込み先となるストレージ装置30のアドレス情報(例えば、LBA: Logical Block Addressing)が含まれているものとする。

【0014】

コントローラ20は、ストレージ装置30を構成する複数の不揮発性半導体記憶装置3

50

1を、RAID技術を用いて管理し、これら複数の不揮発性半導体記憶装置31により論理的に構成される記憶領域に対し、データの書き込みや読み出しをホスト装置10からの要求に応じて実行する。

【0015】

具体的に、コントローラ20は、複数の不揮発性半導体記憶装置31をRAID1、5、6の何れか又はこれらの組み合わせとした構成とすることで、ストレージ装置30の耐障害性・冗長性を実現している。以下、本実施形態ではストレージ装置30をRAID5の構成とした態様について説明する。

【0016】

RAID5は、パリティと呼ばれる誤り訂正符号の記憶用に割り当てる記憶装置と、データの記憶用に割り当てる記憶装置とを、ストライプ毎に順次変更するものである。RAID5を実装するディスクアレイ装置では、耐障害性の向上、大容量化、リード処理の高速化が実現できる。

10

【0017】

図2は、RAID5で構成されたストレージ装置30の記憶領域を模式的に示した図である。なお、同図では4台の不揮発性半導体記憶装置31（不揮発性半導体記憶装置311～314）によりストレージ装置30を構成した例を示しており、当該ストレージ装置30の記憶領域には12個のデータA～Lが記憶されている。

【0018】

RAID5を構成する不揮発性半導体記憶装置31の記憶領域は、コントローラ20により、データの書き込みまたは読み込みの単位となる複数の論理ブロックに分割される。図2に示した例では、データA～Lの夫々やパリティP1～P4の夫々が格納された領域が、1つの論理ブロックを示している。

20

【0019】

ここで、パリティP1～P4は、夫々同一のストライプグループ（0～3）に記憶された複数のデータから算出される復元情報であって、この復元情報に基づいて当該復元情報の生成元となったデータを復元することが可能となっている。例えば、パリティP1はストライプグループ0に記憶されたデータA、B、Cから生成されており、データA、B、Cのうち何れか一のデータにエラーが発生した場合であっても、残りのデータとパリティP1とからエラーの発生したデータを復元することが可能である。なお、データが格納される論理ブロック及びパリティが格納される論理ブロック（以下、パリティ領域という）は、所定のルールに基づき定められているものとするが、その配置位置は図2の例に限定されないものとする。

30

【0020】

ストレージ装置30は、NAND型フラッシュメモリ等の不揮発性半導体素子を記録媒体とする複数の不揮発性半導体記憶装置31を有し、コントローラ20によるRAID管理の下、データを記憶するストレージとして機能する。なお、ストレージ装置30を構成する不揮発性半導体記憶装置31の個数は、コントローラ20が使用するRAIDの規約に応じた個数（例えば、RAID1であれば2個以上、RAID5ならば3個以上）であれば特に問わないものとする。

40

【0021】

<コントローラ20の構成>

次に、図3を参照して、コントローラ20の構成について詳細に説明する。図3は、コントローラ20の詳細構成を示したブロック図である。同図に示したように、コントローラ20は、ホスト側I/F部21と、コマンド処理部22と、ストレージ側I/F部23とを備えている。

【0022】

ホスト側I/F部21は、ホスト装置10と接続するためのインタフェース装置であって、ホスト装置10とコントローラ20（コマンド処理部22）との間で行われるデータの授受を制御する。

50

【 0 0 2 3 】

コマンド処理部 2 2 は、復元情報生成部 2 2 1、復元処理部 2 2 2、キャッシュ管理部 2 2 3 を有し、ホスト側 I / F 部 2 1 を介して入力されるホスト装置 1 0 からの要求に応じて、ストレージ装置 3 0 に対しデータの書き込みや読み込みをストレージ側 I / F 部 2 3 を介して行う。

【 0 0 2 4 】

なお、コマンド処理部 2 2 は、A S I C や C P U 等の処理装置、コントローラ 2 0 の動作を制御する所定のプログラムが格納された R O M や当該処理装置のワーク領域となる R A M 等の記憶装置を備え（何れも図示せず）、これら処理装置と記憶装置に格納されたプログラムとの協働により、復元情報生成部 2 2 1、復元処理部 2 2 2 及びキャッシュ管理部 2 2 3 の各機能部を実現する。

10

【 0 0 2 5 】

ここで、復元情報生成部 2 2 1 は、書き込み対象となるデータのパリティを生成する機能部である。なお、復元情報生成部 2 2 1 は、書き込み先となる領域に既存のデータが存在する場合、当該既存のデータと、このデータに係るパリティと、書き込み対象のデータとから、新たなパリティを生成する。

【 0 0 2 6 】

復元処理部 2 2 2 は、読み込みエラーが発生したデータについて、当該データと同一のストライプに記憶された他のデータ及びパリティを用いて、読み込みエラーが発生したデータの復元を行う機能部である。

20

【 0 0 2 7 】

また、キャッシュ管理部 2 2 3 は、不揮発性半導体記憶装置 3 1 へ書き込むデータ及び不揮発性半導体記憶装置 3 1 から読み込むデータを一時的に保存して管理するとともに、読み込みエラー発生時に復元処理部 2 2 2 により復元されたデータを一時的に保存する機能部である。

【 0 0 2 8 】

コマンド処理部 2 2 は、上記した各機能部（復元情報生成部 2 2 1、復元処理部 2 2 2、キャッシュ管理部 2 2 3）との協働により、ストレージ装置 3 0 に対するデータの書き込み又は読み込みを制御する。

【 0 0 2 9 】

具体的に、コマンド処理部 2 2 は、ホスト装置 1 0 からデータの書き込み要求を受け付けると、このデータの書き込み先が、どの不揮発性半導体記憶装置 3 1 のどの領域（論理ブロック）に対応するのかを特定する。また、復元情報生成部 2 2 1 は、書き込み対象のデータに基づいてパリティを生成する。次いで、コマンド処理部 2 2 は、特定した領域へのデータの書き込みと、当該領域に応じたパリティ用の領域へのパリティの書き込みとを、書き込み先となる不揮発性半導体記憶装置 3 1 に要求することで、書き込み対象のデータと当該データについてのパリティとをストレージ装置 3 0 に書き込む。

30

【 0 0 3 0 】

また、書き込み先となる領域に既存のデータが存在する場合には、既存のデータを新たなデータに更新することになる。この場合、コマンド処理部 2 2 は、書き込み先として特定した領域に格納されている既存のデータと、当該データに係るパリティとの読み込みを、該当する不揮発性半導体記憶装置 3 1 に要求することで、既存のデータ及びパリティを読み込む。このとき、復元情報生成部 2 2 1 は、読み込まれた既存のデータ及びパリティと、書き込み対象のデータとから新たなパリティを生成し、この生成された新たなパリティと書き込み対象のデータとを、キャッシュ管理部に保存し、不揮発性半導体記憶装置 3 1 の該当する領域に書き込むことで、データの更新を行う。

40

【 0 0 3 1 】

また、コマンド処理部 2 2 は、ホスト装置 1 0 からデータの読み込み要求を受け付けると、この読み込み先が、どの不揮発性半導体記憶装置 3 1 のどの領域（論理ブロック）に対応するのかを特定する。そして、コマンド処理部 2 2 は、特定した領域からのデータの

50

読み込みを、読み込み先となる不揮発性半導体記憶装置 3 1 に要求することで、読み込み対象のデータをストレージ装置 3 0 から読み込み、ホスト装置 1 0 に出力する。

【 0 0 3 2 】

なお、不揮発性半導体記憶装置 3 1 の記録媒体が N A N D 型フラッシュメモリの場合、データの読み込みの際に発生する電荷変位や自然放電等の理由により、記憶セル上のデータが破損する可能性がある。一般に、N A N D 型フラッシュメモリを記録媒体とする記憶装置には、破損したデータの誤りを検出し訂正する機構が設けられているが、必ずしも全ての誤りを訂正できるとは限らず、この場合、データを読み込み時にエラーが発生することになる。そのため、本実施形態では、このようなデータの破損について、復元処理部 2 2 2 が読み込みエラーの発生したデータを当該データのパリティに基づいて復元し、キャッシュ管理部 2 2 3 へ一時的に保存した後、該当部分への書き込み要求があった場合に、キャッシュ上のデータを元にして書き込むデータを生成し該当する領域に書き込むことで不揮発性半導体記憶装置 3 1 の復旧を行う。これにより、読み込みエラーの発生した不揮発性半導体記憶装置 3 1 を、正常な状態に復旧させることができる。

10

【 0 0 3 3 】

ストレージ側 I / F 部 2 3 は、不揮発性半導体記憶装置 3 1 と接続するためのインタフェース装置であって、コントローラ 2 0 (コマンド処理部 2 2) と不揮発性半導体記憶装置 3 1 の間で行われるデータの授受を制御する。なお、ストレージ側 I / F 部 2 3 は、不揮発性半導体記憶装置 3 1 毎に設けられているものとするが、これに限らず、一のストレージ側 I / F 部 2 3 と複数の不揮発性半導体記憶装置 3 1 とが接続される態様としてもよい。

20

【 0 0 3 4 】

< コントローラ 2 0 の動作 >

次に、コントローラ 2 0 の動作について説明する。まず、図 4 を参照して、ストレージ装置 3 0 にデータを書き込む際の動作について説明する。図 4 は、コントローラ 2 0 により実行される書き込み処理の手順を示したフローチャートである。なお、本処理の前提として、ストレージ装置 3 0 は R A I D 5 で構成されているものとし、データの書き込み及び読み込みはストレージ装置 3 0 のストライプ単位で行われるものとする。

【 0 0 3 5 】

まず、コマンド処理部 2 2 は、ホスト側 I / F 部 2 1 を介しホスト装置 1 0 からデータの書き込み要求を受け付けると (ステップ S 1 1)、書き込み先となる領域がどの不揮発性半導体記憶装置 3 1 のどの領域に対応するのかを特定する (ステップ S 1 2)。なお、書き込み先となる領域は一であってもよいし、複数であってもよい。

30

【 0 0 3 6 】

続いて、復元情報生成部 2 2 1 は、ステップ S 1 2 で特定した領域に既存のデータが記憶されているか否かを判定する (ステップ S 1 3)。ここで、既存のデータが存在すると判定した場合 (ステップ S 1 3 ; Y e s)、復元情報生成部 2 2 1 は、ステップ S 1 2 で特定された領域について、データとパリティとの読み込みをストレージ側 I / F 部 2 3 を介してストレージ装置 3 0 に要求することで、既存のデータと当該データに係るパリティとをストレージ装置 3 0 から読み込む (ステップ S 1 4)。

40

【 0 0 3 7 】

次いで、復元情報生成部 2 2 1 は、ステップ S 1 4 で読み込んだ既存のデータ及びパリティと、書き込み対象のデータとから新たなパリティを生成し (ステップ S 1 5)、ステップ S 1 6 の処理に移行する。

【 0 0 3 8 】

なお、ステップ S 1 4 の読み込みの際、読み込みエラーが発生した場合には、読み込みエラーの発生した既存のデータと同一のストライプに記憶された他のデータをストレージ装置 3 0 から読み込み、当該他のデータと書き込み対象のデータとから、新たなパリティの生成を行うものとする。

【 0 0 3 9 】

50

一方、ステップ S 1 3 において、既存のデータが存在しないと判定した場合（ステップ S 1 3 ; N o ）、復元情報生成部 2 2 1 は、書き込み対象のデータからパリティを生成し（ステップ S 1 5 ）、ステップ S 1 6 の処理に移行する。

【 0 0 4 0 】

続くステップ S 1 6 では、コマンド処理部 2 2 が、書き込み対象のデータをステップ S 1 2 で特定した領域に書き込むとともに、当該データの書き込み領域に応じたパリティ領域にステップ S 1 5 で生成したパリティを書き込む（ステップ S 1 6 ）。ここで、コマンド処理部 2 2 は、データ及び / 又はパリティの書き込み時に、書き込みエラーが発生したか否かを判定し、正常に書き込みが行われたと判定した場合には（ステップ S 1 7 ; N o ）、ステップ S 2 2 の処理に直ちに移行する。

10

【 0 0 4 1 】

また、ステップ S 1 7 において、書き込みエラーを検出した場合（ステップ S 1 7 ; Y e s ）、コマンド処理部 2 2 は、書き込み先となる不揮発性半導体記憶装置 3 1 に障害が発生したと判断し、当該不揮発性半導体記憶装置 3 1 を除いた残りの不揮発性半導体記憶装置 3 1 で R A I D 5 の構成を維持する縮退動作が可能か否かを判定する（ステップ S 1 8 ）。ここで、縮退動作が不可能と判定した場合（ステップ S 1 8 ; N o ）、コマンド処理部 2 2 は、書き込みが失敗したことを示す応答をホスト装置 1 0 へ出力し（ステップ S 1 9 ）、本処理を終了する。

【 0 0 4 2 】

また、ステップ S 1 8 において、縮退動作が可能と判定した場合、コマンド処理部 2 2 は、障害が発生した不揮発性半導体記憶装置 3 1 を R A I D 5 の構成から除外し（ステップ S 2 0 ）、縮退動作としたストレージ装置 3 0 に書き込み対象のデータと、ステップ S 1 5 で生成したパリティとを書き込むと（ステップ S 2 1 ）、ステップ S 2 2 の処理に移行する。

20

【 0 0 4 3 】

続くステップ S 2 2 において、コマンド処理部 2 2 は、ステップ S 1 2 で特定した全ての領域にデータを書き込んだか否かを判定し、未処理の領域が存在すると判定した場合には（ステップ S 2 2 ; N o ）、ステップ S 1 3 の処理に再び戻り、他のストライプに含まれた領域を処理対象とする。また、ステップ S 2 2 で特定した全ての領域にデータを書き込んだと判定した場合（ステップ S 2 2 ; Y e s ）、コマンド処理部 2 2 は、書き込みが終了したことを示す応答をホスト装置 1 0 へ出力し（ステップ S 2 3 ）、本処理を終了する。

30

【 0 0 4 4 】

次に、図 5 を参照して、ストレージ装置 3 0 からデータを読み込む際の動作について説明する。図 5 は、コントローラ 2 0 により実行される読み込み処理の手順を示したフローチャートである。なお、本処理の前提として、ストレージ装置 3 0 は R A I D 5 で構成されているものとし、データの書き込み及び読み込みはストレージ装置 3 0 のストライプ単位で行われるものとする。

【 0 0 4 5 】

まず、コマンド処理部 2 2 は、ホスト側 I / F 部 2 1 を介しホスト装置 1 0 からデータの読み込み要求を受け付けると（ステップ S 3 1 ）、読み込み先となる領域がどの不揮発性半導体記憶装置 3 1 のどの領域に対応するのかを特定する（ステップ S 3 2 ）。なお、読み込み先となる領域は一であってもよいし、複数であってもよい。

40

【 0 0 4 6 】

続いて、コマンド処理部 2 2 は、ステップ S 1 2 で特定した領域に対応する各不揮発性半導体記憶装置 3 1 に対し、当該領域の読み込みを要求することで、ストレージ装置 3 0 から読み込み対象のデータを読み込む（ステップ S 3 3 ）。このとき、コマンド処理部 2 2 は、ステップ S 3 3 の読み込みの際に読み込みエラーが発生したか否かを判定する（ステップ S 3 4 ）。ここで、コマンド処理部 2 2 が、正常に読み込みできたと判定した場合（ステップ S 3 4 ; N o ）、キャッシュ管理部 2 2 3 へ読み込んだデータを保存し（ステ

50

ップ S 3 8)、ステップ S 3 9 の処理に移行する。

【 0 0 4 7 】

一方、ステップ S 3 4 において、コマンド処理部 2 2 が、読み込みエラーが発生したと判定した場合 (ステップ S 3 4 ; Y e s)、復元処理部 2 2 2 は、読み込みエラーが発生したデータと同一のストライプグループに記憶された未読み込みのデータと、パリティとをストレージ装置 3 0 から読み込み、当該データ及びパリティと、ステップ S 3 3 で読み込んだデータとから読み込みエラーが発生したデータを復元する (ステップ S 3 5)。

【 0 0 4 8 】

続いて、復元処理部 2 2 2 は、ステップ S 3 5 の処理でデータを復元できたか否かを判定する。ここで、パリティが読み込めない等の理由によりデータを復元することができな
10
いと判定した場合 (ステップ S 3 6 ; N o)、復元処理部 2 2 2 は、読み込みが失敗したことを示す応答をホスト装置 1 0 に出力し (ステップ S 3 7)、本処理を終了する。また、ステップ S 3 6 において、復元処理部 2 2 2 が、データを復元できたと判定した場合 (ステップ S 3 6 ; Y e s)、この復元したデータをキャッシュ管理部 2 2 3 に保存し (ステップ S 3 8)、ステップ S 3 9 の処理に移行する。

【 0 0 4 9 】

続くステップ S 3 9 において、コマンド処理部 2 2 は、ステップ S 3 2 で特定した全ての領域からデータを読み込んだか否かを判定する (ステップ S 3 9)。ここで、未処理の領域が存在すると判定した場合 (ステップ S 3 9 ; N o)、ステップ S 3 3 の処理に再び
20
戻り、他の領域を処理対象とする。また、ステップ S 3 2 で特定した全ての領域からデータを読み込んだと判定した場合 (ステップ S 3 9 ; Y e s)、コマンド処理部 2 2 は、キャッシュ管理部 2 2 3 を参照し各領域から読み込んだデータをホスト装置 1 0 に出力し (ステップ S 4 0)、読み込み終了を示す応答をホスト装置 1 0 に出力し (ステップ S 4 1)、さらに以下に記載するようにデータの書き込み時にキャッシュ管理部 2 2 3 のデータを参照して不揮発性半導体記憶装置 3 1 へ書き込み (ステップ S 4 2)、本処理を終了する。

【 0 0 5 0 】

上記ステップ S 3 4 で読み込みエラーが発生した場合、復元したデータはキャッシュ管理部 2 2 3 にのみ保存されることになる。このため、不揮発性半導体記憶装置 3 1 とキャ
30
ッシュ管理部 2 2 3 とのデータの一貫性が保たれていない状態となるが、例えば、当該データの書き込み時にキャッシュ管理部 2 2 3 のデータを参照して不揮発性半導体記憶装置 3 1 へ書き込むことで一貫性が保持できる。これにより、不揮発性半導体記憶装置 3 1 に書き込む場合には書き込み時間がかかり、読み込み処理が完了せず、データの要求元 (ホスト装置) への出力が遅くなることを防止することができる。

【 0 0 5 1 】

以上のように、第 1 の実施形態によれば、読み込み時にエラーが発生した場合、この読み込みエラーが発生したデータを復元してキャッシュ管理部 2 2 3 へ保存し、以降、書き
40
込み要求があったときに当該読み込みエラーが発生した不揮発性半導体記憶装置 3 1 上の領域に書き込むことで、当該不揮発性半導体記憶装置 3 1 を正常な状態に復旧させることができるため、R A I D 構成とされた複数の不揮発性半導体記憶装置 3 1 を有効に活用することができる。

【 0 0 5 2 】

[第 2 の実施形態]

第 1 の実施形態では、ストレージ装置 3 0 を R A I D 5 の構成としたが、データの復元
を可能とするストレージシステムであれば特に問わず、上述したように R A I D 1 や R A I D 6 の構成を採用してもよい。以下、第 2 の実施形態として、ストレージ装置 3 0 を R A I D 1 の構成とした場合について説明する。なお、第 1 の実施形態と同様の構成要素については、同一の符号を付与し説明を省略する。

【 0 0 5 3 】

< コントローラ 4 0 の構成 >

10

20

30

40

50

まず、第2の実施形態に係るコントローラ40について説明する。コントローラ40は、ストレージ装置30を構成する二つの不揮発性半導体記憶装置31を、RAID1の技術を用いて管理し、これら複数の不揮発性半導体記憶装置31により論理的に構成された記憶領域に対し、データの書き込みや読み出しをホスト装置10からのアクセス要求に応じて行う。

【0054】

図6は、RAID1で構成されたストレージ装置30の記憶領域を模式的に示した図である。RAID1は、ミラーリングとも呼ばれ、少なくとも2台以上の不揮発性半導体記憶装置31に同一のデータを、同一のストライプ領域に同時に書き込みすることで、耐障害性・冗長性を確保している。なお、同図では2台の不揮発性半導体記憶装置31（不揮発性半導体記憶装置311、312）によりストレージ装置30を構成した例を示しており、当該ストレージ装置30の記憶領域に12個のデータA～Lが記憶されている。

10

【0055】

このRAID1の構成の場合、一方の不揮発性半導体記憶装置31についてデータの読み込みエラーが発生した場合、他方の不揮発性半導体記憶装置31から同一のデータを読み込むことで、システム自体は問題無く稼動し続けることができる。

【0056】

図7は、第2の実施形態に係るコントローラ40の詳細構成を示したブロック図である。同図に示したように、コントローラ40は、ホスト側I/F部21と、コマンド処理部41と、ストレージ側I/F部23とを備えている。

20

【0057】

ここで、コマンド処理部41は、復元処理部411、キャッシュ管理部412を有し、ホスト側I/F部21を介して入力されるホスト装置10のアクセス要求に応じて、ストレージ側I/F部23を介して接続されるストレージ装置30に対し、データの書き込みや読み込みを行う。

【0058】

なお、コマンド処理部41は、ASICやCPU等の処理装置、コントローラ20の動作を制御する所定のプログラムが格納されたROMや当該処理装置のワーク領域となるRAM等の記憶装置を備え（何れも図示せず）、これら処理装置と記憶装置に格納されたプログラムとの協働により、復元処理部411及びキャッシュ管理部412の各機能部を実現する。

30

【0059】

具体的に、コマンド処理部41は、ホスト装置10からデータの書き込み要求を受け付けると、このデータの書き込み先となる領域（論理ブロック）を各不揮発性半導体記憶装置31から夫々特定する。そして、コマンド処理部22は、書き込み対象のデータを不揮発性半導体記憶装置31の個数に応じた数だけ複製すると、各不揮発性半導体記憶装置31の特定した領域に夫々書き込む。

【0060】

また、コマンド処理部22は、ホスト装置10からデータの読み込み要求を受け付けると、アクセス対象となる一の不揮発性半導体記憶装置31から、この読み込み先に対応する領域（論理ブロック）を特定する。そして、コマンド処理部22は、アクセス対象となる不揮発性半導体記憶装置31の特定した領域からデータを読み込み、ホスト装置10に出力する。ここで、アクセス対象となる不揮発性半導体記憶装置31は、予め定められているものとしてもよいし、負荷分散などの理由により他の不揮発性半導体記憶装置31と動的に切り替わるものとしてもよい。以下、アクセス対象の不揮発性半導体記憶装置31を「主記憶装置」と呼び、他の不揮発性半導体記憶装置31を「待機記憶装置」と呼ぶ。

40

【0061】

なお、不揮発性半導体記憶装置31の記録媒体がNAND型フラッシュメモリの場合、上記したように読み込みの際に発生する電荷変位や自然放電等の理由により、記憶セル上のデータが破損する可能性がある。そのため、コマンド処理部41では、主記憶装置に発

50

生したデータの破損について、復元処理部 4 1 1 が当該データと同一のデータを待機記憶装置から読み込み、キャッシュ管理部 4 1 2 へ一時的に保存する。そして、コマンド処理部 4 1 は、読み込みエラーの発生した部分への書き込み要求があった場合に、キャッシュ上のデータを元にして書き込むデータを生成し、該当する領域に書き込むことでデータの復元を行う。つまり、待機記憶装置に記憶されるデータは、主記憶装置に記憶されたデータを復元するための復元情報として機能する。これにより、読み込みエラーの発生した主記憶装置を、正常な状態に復旧することができる。

【 0 0 6 2 】

< コントローラ 4 0 の動作 >

次に、コントローラ 4 0 の動作について説明する。まず、図 8 を参照して、ストレージ装置 3 0 にデータを書き込む際の動作について説明する。

10

【 0 0 6 3 】

図 8 は、コントローラ 4 0 により実行される書き込み処理の手順を示したフローチャートである。なお、本処理の前提として、ストレージ装置 3 0 は R A I D 1 で構成されているものとし、データの書き込みは論理ブロック単位で行われるものとする。

【 0 0 6 4 】

まず、コマンド処理部 4 1 は、ホスト側 I / F 部 2 1 を介しホスト装置 1 0 からデータの書き込み要求を受け付けると (ステップ S 5 1)、ストレージ装置 3 0 を構成する各不揮発性半導体記憶装置 3 1 について書き込み先となる領域を特定する (ステップ S 5 2)

20

【 0 0 6 5 】

続いて、コマンド処理部 4 1 は、書き込み対象のデータを不揮発性半導体記憶装置 3 1 の個数分複製すると、ステップ S 5 2 で特定した各不揮発性半導体記憶装置 3 1 の領域に、書き込み対象のデータを書き込む (ステップ S 5 3)。ここで、コマンド処理部 2 2 は、データの書き込みの際に書き込みエラーが発生したか否かを判定し、正常に書き込みが行われたと判定した場合には (ステップ S 5 4 ; N o)、ステップ S 5 8 の処理に直ちに移行する。

【 0 0 6 6 】

また、ステップ S 5 4 において、書き込みエラーを検出した場合 (ステップ S 5 4 ; Y e s)、コマンド処理部 4 1 は、書き込み先となった不揮発性半導体記憶装置 3 1 に障害が発生したと判定し、当該不揮発性半導体記憶装置 3 1 を除いた残りの不揮発性半導体記憶装置 3 1 でシステムを維持する縮退動作が可能か否かを判定する (ステップ S 5 5)。ここで、縮退動作が不可能と判定した場合 (ステップ S 5 5 ; N o)、コマンド処理部 2 2 は、書き込みが失敗したことを示す応答をホスト装置 1 0 に出力し (ステップ S 5 6)、本処理を終了する。

30

【 0 0 6 7 】

また、ステップ S 5 5 において、縮退動作が可能と判定した場合、コマンド処理部 4 1 は、障害が発生した不揮発性半導体記憶装置 3 1 を R A I D 1 の構成から除外し (ステップ S 5 7)、ステップ S 5 8 の処理に移行する。

【 0 0 6 8 】

続くステップ S 5 8 において、コマンド処理部 4 1 は、ステップ S 5 2 で特定した全ての領域にデータを書き込んだか否かを判定する (ステップ S 5 8)。ここで、未処理の領域が存在すると判定した場合には (ステップ S 5 8 ; N o)、ステップ S 5 3 の処理に再び戻り、他の領域を処理対象とする。また、ステップ S 5 2 で特定した全ての領域にデータを書き込んだと判定した場合 (ステップ S 5 8 ; Y e s)、コマンド処理部 4 1 は、書き込みが終了したことを示す応答をホスト装置 1 0 に出力し (ステップ S 5 9)、本処理を終了する。

40

【 0 0 6 9 】

次に、図 9 を参照して、ストレージ装置 3 0 からデータを読み込む際の動作について説明する。図 9 は、コントローラ 4 0 により実行される読み込み処理の手順を示したフロー

50

チャートである。なお、本処理の前提として、ストレージ装置 30 は R A I D 1 で構成されているものとし、データの書き込みは論理ブロック単位で行われるものとする。

【0070】

まず、コマンド処理部 41 は、ホスト側 I / F 部 21 を介しホスト装置 10 からデータの読み込み要求を受け付けると (ステップ S 61)、読み込み先となる領域が主記憶装置のどの領域に対応するのかを特定する (ステップ S 62)。

【0071】

続いて、コマンド処理部 41 は、主記憶装置に対し、ステップ S 62 で特定した領域のデータの読み込みを要求することで、読み込み対象となったデータをストレージ装置 30 から読み込む (ステップ S 63)。このとき、コマンド処理部 41 は、ステップ S 63 の読み込みの際に読み込みエラーが発生したか否かを判定する (ステップ S 64)。ここで、コマンド処理部 41 が、正常に読み込みできたと判定した場合 (ステップ S 64 ; N o)、キャッシュ管理部 412 へ読み込んだデータを保存し (ステップ S 68)、ステップ S 69 の処理に移行する。

10

【0072】

一方、ステップ S 64 において、コマンド処理部 41 が、読み込みエラーが発生したと判定した場合 (ステップ S 64 ; Y e s)、復元処理部 411 は、読み込みエラーの発生したデータと同一のデータを待機記憶装置から読み込む (ステップ S 65)。続いて、復元処理部 411 は、待機記憶装置からデータを読み込めたか否かを判定する (ステップ S 66)。ここで、何れの待機記憶装置からもデータを読み込めないと判定した場合 (ステップ S 66 ; N o)、復元処理部 411 は、読み込みが失敗したことを示す応答をホスト装置 10 へ出力し (ステップ S 67)、本処理を終了する。

20

【0073】

また、ステップ S 66 において、復元処理部 411 は、待機記憶装置からデータを読み込めたと判定した場合 (ステップ S 66 ; Y e s)、ステップ S 65 で読み込まれたデータをキャッシュ管理部 412 へ保存し (ステップ S 68)、ステップ S 69 の処理に移行する。

【0074】

続くステップ S 69 において、コマンド処理部 41 は、ステップ S 62 で特定した全ての領域からデータを読み込んだか否かを判定する (ステップ S 69)。ここで、未処理の領域が存在すると判定した場合には (ステップ S 69 ; N o)、ステップ S 63 の処理に再び戻り、次の領域を処理対象とする。また、ステップ S 62 で特定した全ての領域からデータを読み込んだと判定した場合 (ステップ S 69 ; Y e s)、コマンド処理部 22 は、各領域から読み込んだデータをホスト装置 10 へ出力すると (ステップ S 70)、読み込み終了を示す応答をホスト装置 10 へ出力し (ステップ S 71)、さらに以下に記載するようにデータの書き込み時にキャッシュ管理部 223 のデータを参照して不揮発性半導体記憶装置 31 へ書き込み (ステップ S 72)、本処理を終了する。

30

【0075】

なお、上記ステップ S 64 で読み込みエラーが発生した場合、復元したデータはキャッシュ管理部 412 にのみ保存されることになる。このため、不揮発性半導体記憶装置 31 とキャッシュ管理部 412 とのデータの一貫性が保たれていない状態となるが、例えば、当該データの書き込み時にキャッシュ管理部 412 のデータを参照して不揮発性半導体記憶装置 31 へ書き込むことで一貫性が保持できる。これにより、不揮発性半導体記憶装置 31 に書き込む場合には書き込み時間がかかり、読み込み処理が完了せず、データの要求元 (ホスト装置) への出力が遅くなることを防止することができる。

40

【0076】

以上のように、第 2 の実施形態によれば、読み込み時にエラーが発生した場合、この読み込みエラーが発生したデータを復元してキャッシュ管理部 412 へ保存し、以降、書き込み要求があったときに当該読み込みエラーが発生した不揮発性半導体記憶装置 31 上の領域に書き込むことで、当該不揮発性半導体記憶装置 31 を正常な状態に復旧させること

50

ができるため、R A I D構成とされた複数の不揮発性半導体記憶装置 3 1 を有効に活用することができる。

【 0 0 7 7 】

以上、発明の実施の形態について説明したが、本発明はこれに限定されるものではなく、本発明の主旨を逸脱しない範囲での種々の変更、置換、追加などが可能である。

【 図面の簡単な説明 】

【 0 0 7 8 】

【 図 1 】 ストレージシステムの構成を示した図である。

【 図 2 】 ストレージ装置を R A I D 5 とした場合の記憶領域を模式的に示した図である。

【 図 3 】 第 1 の実施形態に係るコントローラの構成を示した図である。

10

【 図 4 】 第 1 の実施形態に係る書き込み処理の手順を示したフローチャートである。

【 図 5 】 第 1 の実施形態にかかる読み込み処理の手順を示したフローチャートである。

【 図 6 】 第 2 の実施形態に係るコントローラの構成を示した図である。

【 図 7 】 ストレージ装置を R A I D 1 とした場合の記憶領域を模式的に示した図である。

【 図 8 】 第 2 の実施形態に係る書き込み処理の手順を示したフローチャートである。

【 図 9 】 第 2 の実施形態にかかる読み込み処理の手順を示したフローチャートである。

【 符号の説明 】

【 0 0 7 9 】

1 0 0 ストレージシステム

1 0 ホスト装置

20

2 0 コントローラ

2 1 ホスト側 I / F 部

2 2 コマンド処理部

2 2 1 復元情報生成部

2 2 2 復元処理部

2 2 3 キャッシュ管理部

2 3 ストレージ側 I / F 部

3 0 ストレージ装置

3 1 不揮発性半導体記憶装置

4 0 コントローラ

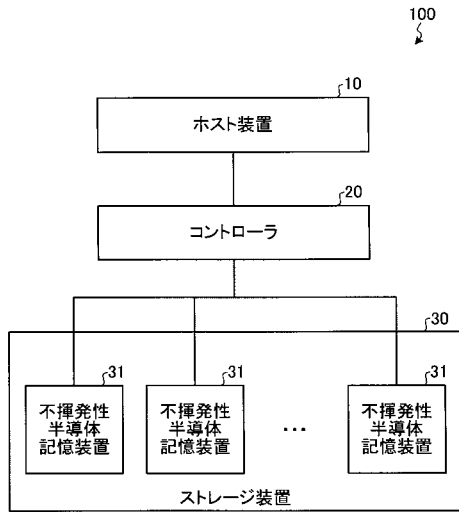
30

4 1 コマンド処理部

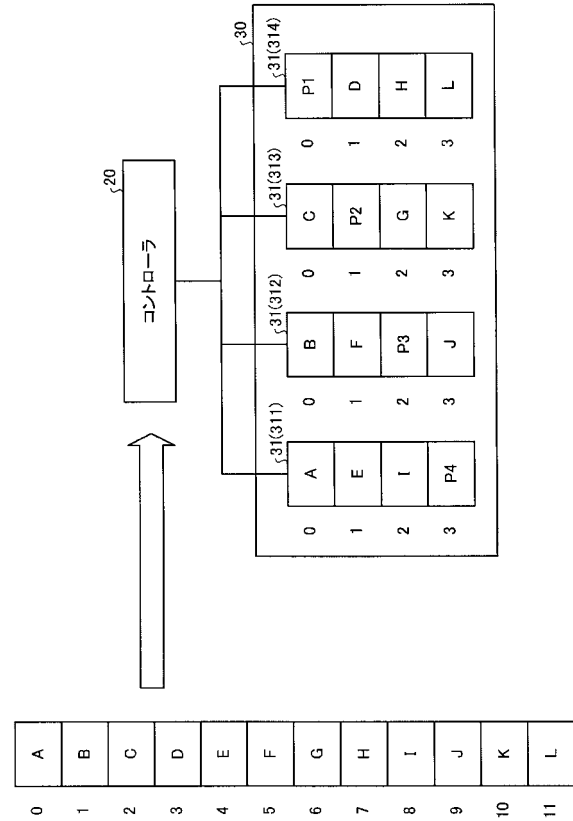
4 1 1 復元処理部

4 1 2 キャッシュ管理部

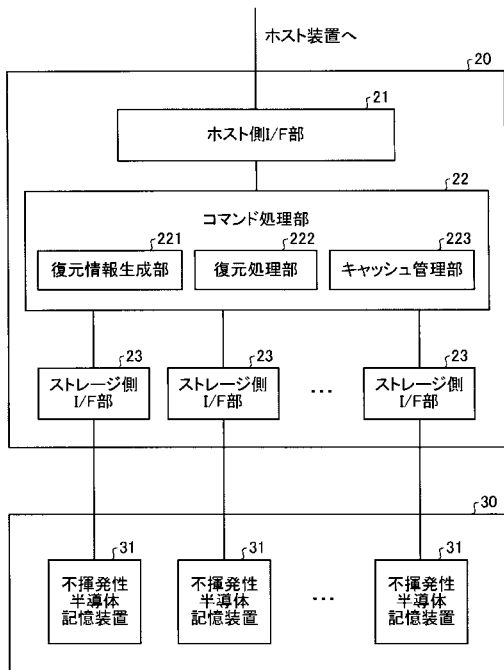
【図1】



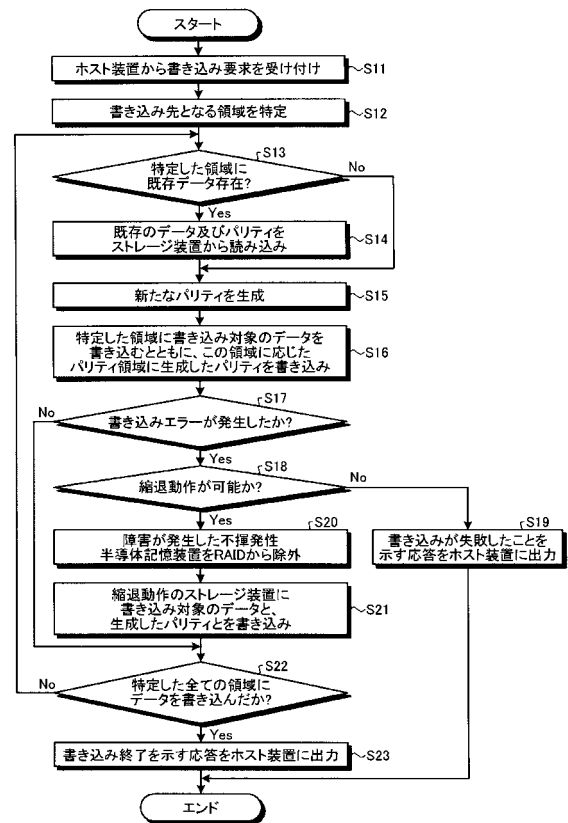
【図2】



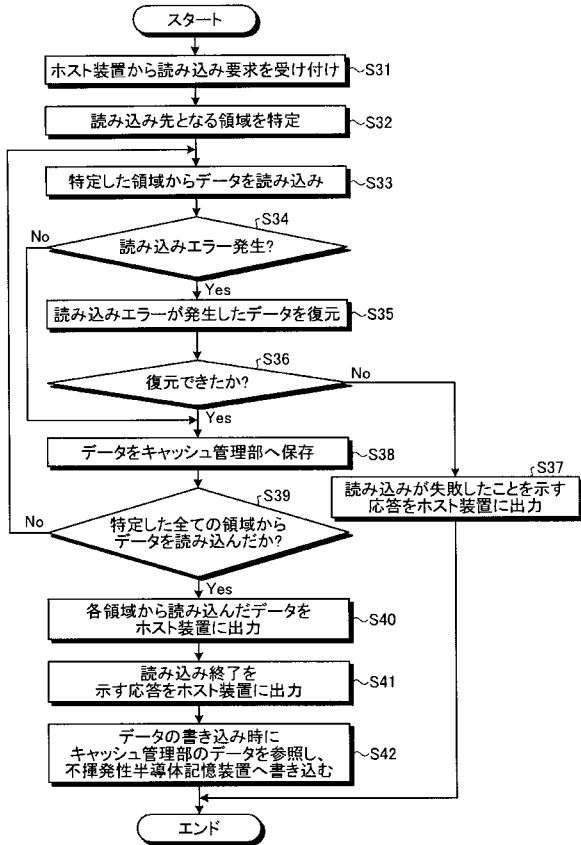
【図3】



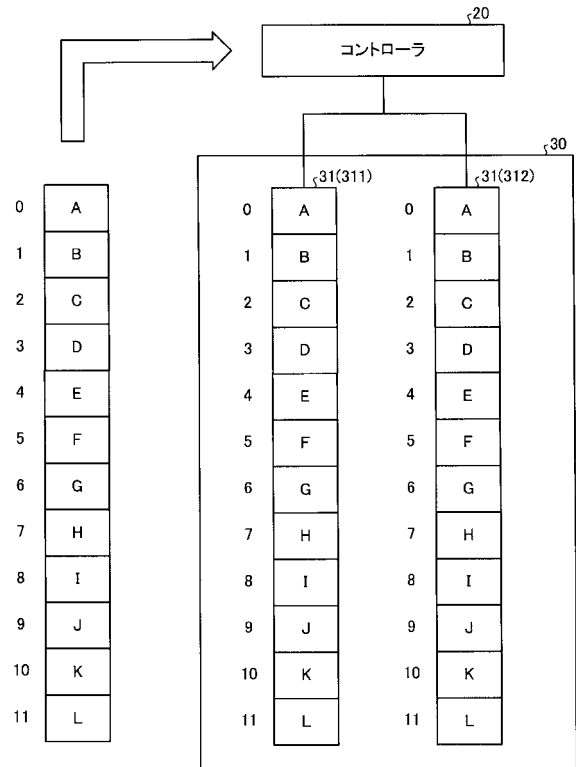
【図4】



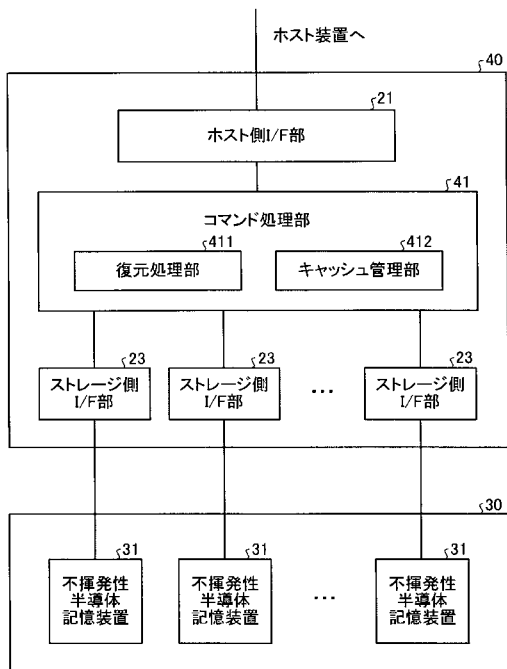
【 図 5 】



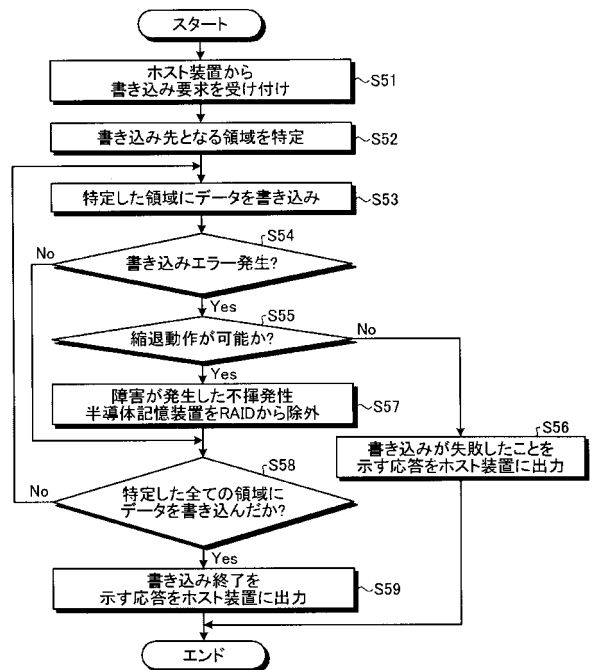
【 図 6 】



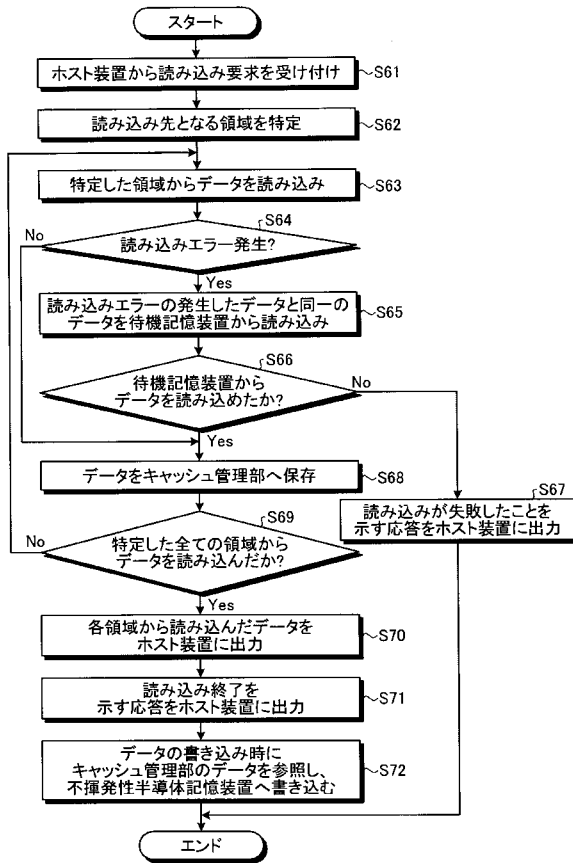
【 図 7 】



【 図 8 】



【 図 9 】



フロントページの続き

(72)発明者 浅野 滋博

東京都港区芝浦一丁目1番1号 株式会社東芝内

Fターム(参考) 5B018 GA04 GA06 HA35 KA22 MA22 NA06 QA03

5B065 BA05 CA30 EA03 EA24