

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4863749号
(P4863749)

(45) 発行日 平成24年1月25日(2012.1.25)

(24) 登録日 平成23年11月18日(2011.11.18)

(51) Int. Cl.	F I		
G06F 12/00	(2006.01)	G06F 12/00	542K
G06F 12/16	(2006.01)	G06F 12/00	514E
G06F 3/08	(2006.01)	G06F 12/16	310A
G06F 3/06	(2006.01)	G06F 3/08	H
		G06F 3/06	301J
請求項の数 23 (全 32 頁) 最終頁に続く			

(21) 出願番号 特願2006-92217(P2006-92217)
 (22) 出願日 平成18年3月29日(2006.3.29)
 (65) 公開番号 特開2007-265265(P2007-265265A)
 (43) 公開日 平成19年10月11日(2007.10.11)
 審査請求日 平成21年1月8日(2009.1.8)

(73) 特許権者 000005108
 株式会社日立製作所
 東京都千代田区丸の内一丁目6番6号
 (74) 代理人 100064414
 弁理士 磯野 道造
 (72) 発明者 田中 勝也
 神奈川県川崎市麻生区王禅寺1099番地
 株式会社日立製作所
 システム開発研究所内
 (72) 発明者 島田 健太郎
 神奈川県川崎市麻生区王禅寺1099番地
 株式会社日立製作所
 システム開発研究所内
 審査官 桜井 茂行
 最終頁に続く

(54) 【発明の名称】フラッシュメモリを用いた記憶装置、その消去回数平準化方法、及び消去回数平準化プログラム

(57) 【特許請求の範囲】

【請求項1】

記憶装置コントローラと、記憶媒体として前記記憶装置コントローラに接続される複数のフラッシュメモリ・モジュールとを有する記憶装置であって、

前記フラッシュメモリ・モジュールは、1個以上のフラッシュメモリ・チップと、前記フラッシュメモリ・チップに属するブロックの消去回数を平準化するメモリコントローラとを有し、

前記記憶装置コントローラは、

前記複数のフラッシュメモリ・モジュールを含むウエアレベリング・グループを構成し、

前記ウエアレベリング・グループに属する前記複数のフラッシュメモリ・モジュール毎に、書き込みデータ容量に基づいてフラッシュメモリ・モジュール内の平均消去回数を算出し、

前記フラッシュメモリ・モジュール毎の平均消去回数が所定の条件を満たした場合に、前記ウエアレベリング・グループに属するフラッシュメモリ・モジュール間でデータを移動する

ことを特徴とするフラッシュメモリを用いた記憶装置。

【請求項2】

前記記憶装置コントローラは、

前記フラッシュメモリ・モジュール毎の前記平均消去回数の最大値と最小値の差が所定

値以上の場合に、前記平均消去回数が最大のフラッシュメモリ・モジュールに格納されたデータと前記平均消去回数が最小のフラッシュメモリ・モジュールに格納されたデータを入換える

ことを特徴とする請求項 1 に記載のフラッシュメモリを用いた記憶装置。

【請求項 3】

前記記憶装置コントローラは、

前記ウエアレベリング・グループに属するフラッシュメモリ・モジュールにアクセスするための論理ページアドレスを前記記憶装置コントローラの内部で扱うための仮想ページアドレスに変換し、

前記データを移動した後に、前記仮想ページアドレスと前記論理ページアドレスの対応関係を変更する

10

ことを特徴とする請求項 1 に記載のフラッシュメモリを用いた記憶装置。

【請求項 4】

前記記憶装置は、複数のウエアレベリング・グループを備え、

前記記憶装置コントローラは、前記複数のウエアレベリング・グループを含む R A I D グループを構成し、

前記 R A I D グループは、フラッシュメモリ・モジュール故障時に記録データを修復するための冗長情報を含む

ことを特徴とする請求項 1 に記載のフラッシュメモリを用いた記憶装置。

【請求項 5】

20

前記 R A I D グループは、R A I D レベル 0、R A I D レベル 1、R A I D レベル 1 + 0、R A I D レベル 3、R A I D レベル 5、及び R A I D レベル 6 のうちのいずれかひとつの論理グループであり、前記 R A I D グループを構成する全てのウエアレベリング・グループの容量を等しく設定した

ことを特徴とする請求項 1 に記載のフラッシュメモリを用いた記憶装置。

【請求項 6】

前記フラッシュメモリ・モジュールのシステム運用時の実効的書込み速度と前記記憶装置の製品寿命との積をフラッシュメモリの書換え保証回数で割った商を第 1 の値とし、

前記フラッシュメモリ・モジュールの連続書込み速度と前記記憶装置の製品寿命との積をフラッシュメモリの書換え保証回数で割った商を第 2 の値とし、

30

前記ウエアレベリング・グループの容量を、第 1 の値以上、かつ、第 2 の値以下に設定した

ことを特徴とする請求項 1 又は請求項 2 に記載のフラッシュメモリを用いた記憶装置。

【請求項 7】

前記 R A I D グループは、R A I D レベル 2 又は R A I D レベル 4 の論理グループであり、R A I D グループを構成するウエアレベリング・グループの内、冗長情報格納用のウエアレベリング・グループ数を m 、データ格納用の第 1 の論理グループ数を n とするとき、冗長情報格納用のウエアレベリング・グループ容量を、データ格納用のウエアレベリング・グループ容量の 1 倍以上、かつ、 n / m 倍以下に設定した

ことを特徴とする請求項 4 に記載のフラッシュメモリを用いた記憶装置。

40

【請求項 8】

前記論理ページアドレスでアクセスするメモリ領域は、前記仮想ページアドレスでアクセスするメモリ領域より大きく設定した

ことを特徴とする請求項 3 に記載のフラッシュメモリを用いた記憶装置。

【請求項 9】

前記記憶装置コントローラは、前記ウエアレベリング・グループに属するフラッシュメモリ・モジュールにアクセスするための論理ページアドレスと前記記憶装置コントローラの内部で扱うための仮想ページアドレスの対応情報及び前記フラッシュメモリ・モジュール内のブロックの平均消去回数を記憶する

ことを特徴とする請求項 1 又は請求項 4 に記載のフラッシュメモリを用いた記憶装置。

50

【請求項 1 0】

前記記憶装置コントローラは、複数の前記 R A I D グループを構成することを特徴とする請求項 4 に記載のフラッシュメモリを用いた記憶装置。

【請求項 1 1】

前記記憶装置コントローラは、R A I D レベルの異なる複数の前記 R A I D グループを構成する

ことを特徴とする請求項 4 に記載のフラッシュメモリを用いた記憶装置。

【請求項 1 2】

前記記憶装置コントローラは、前記記憶装置が起動された時あるいは前記記憶装置に記録媒体が接続された時に、前記記録媒体がフラッシュメモリであるか否かを判定する

ことを特徴とする請求項 1 又は請求項 4 に記載のフラッシュメモリを用いた記憶装置。

【請求項 1 3】

1 個以上のフラッシュメモリ・チップと、前記フラッシュメモリ・チップに属するブロックの消去回数を平準化するためのメモリコントローラとを有するフラッシュメモリ・モジュールと、

複数の前記フラッシュメモリ・モジュールを組み合わせるウエアレベリング・グループを構成し、前記 ウエアレベリング・グループ に属するフラッシュメモリ・モジュールをアクセスするための 論理ページアドレス を前記記憶装置コントローラの内部で扱うための 仮想ページアドレス に変換し、前記 ウエアレベリング・グループ を複数組み合わせる R A I D グループを構成する記憶装置コントローラと

を備えたフラッシュメモリを用いた記憶装置の消去回数平準化方法であって、

前記記憶装置コントローラは、

前記フラッシュメモリ・モジュール内の所定メモリ空間毎の書込みデータ容量を計数管理するステップと、

前記フラッシュメモリ・モジュール毎に延べ書込みデータ容量を前記フラッシュメモリ・モジュール容量で割った平均消去回数を算出するステップと、

前記平均消去回数の最大値と最小値の消去回数差が所定値以上であるか否かを判定する第 1 の判定ステップと、

前記第 1 の判定ステップで、消去回数差が所定値以上の場合に、消去回数差が最大の前記フラッシュメモリ・モジュール間で、書込みデータ容量最大メモリ空間と書込みデータ容量最小メモリ空間との間でデータの入換え及び前記 論理ページアドレス と前記 仮想ページアドレス の対応情報を変更するステップとを含んで実行する

ことを特徴とするフラッシュメモリを用いた記憶装置の消去回数平準化方法。

【請求項 1 4】

前記記憶装置コントローラは、前記変更するステップを実行後、データの入換えをしていないフラッシュメモリ・モジュールが複数あるか否かを判定する第 2 の判定ステップと、

第 2 の判定ステップで、データの入換えをしていないフラッシュメモリ・モジュールが複数ある場合、第 1 の判定ステップに戻るステップとを含んで実行する

ことを特徴とする請求項 1 3 に記載のフラッシュメモリを用いた記憶装置の消去回数平準化方法。

【請求項 1 5】

前記記憶装置コントローラは、前記 ウエアレベリング・グループ 当たりの延べ書込みデータ容量が、所定の値に達した場合に、第 1 の判定ステップを実行する

ことを特徴とする請求項 1 3 又は請求項 1 4 に記載のフラッシュメモリを用いた記憶装置の消去回数平準化方法。

【請求項 1 6】

前記記憶装置コントローラは、前記 論理ページアドレス のメモリ空間が、前記 仮想ページアドレス のメモリ空間より大きく設定されており、

前記変更するステップが、論理ページアドレス のメモリ空間内でのデータの上書き操作

により、データの入換えを行うステップを含んで実行する

ことを特徴とする請求項 1 3 から請求項 1 5 のうちのいずれか 1 項に記載のフラッシュメモリを用いた記憶装置の消去回数平準化方法。

【請求項 1 7】

書込みデータ容量を計数管理した前記フラッシュメモリ・モジュール内の所定メモリ空間毎に、前記メモリコントローラの 1 回当たりのデータ転送容量に同じ大きさの空き領域が設けられている

ことを特徴とする請求項 1 3 に記載のフラッシュメモリを用いた記憶装置の消去回数平準化方法。

【請求項 1 8】

書込みデータ容量を計数管理した前記フラッシュメモリ・モジュール内の所定メモリ空間と同じ大きさの空き領域が前記フラッシュメモリ・モジュール内に設けられている

ことを特徴とする請求項 1 3 に記載のフラッシュメモリを用いた記憶装置の消去回数平準化方法。

【請求項 1 9】

前記記憶装置コントローラは、前記フラッシュメモリ・モジュールを新規のフラッシュメモリ・モジュールへ交換した場合、交換前のフラッシュメモリ・モジュールにおいて計数管理した所定メモリ空間毎の書込みデータ容量を、交換後のフラッシュメモリ・モジュールの所定メモリ空間毎の書込みデータ容量とするステップを含んで実行する

ことを特徴とする請求項 1 3 に記載のフラッシュメモリを用いた記憶装置の消去回数平準化方法。

【請求項 2 0】

1 個以上のフラッシュメモリ・チップと、前記フラッシュメモリ・チップに属するブロックの消去回数を平準化するためのメモリコントローラとを有するフラッシュメモリ・モジュールと、

複数の前記フラッシュメモリ・モジュールを組み合わせるウエアレベリング・グループを構成し、前記ウエアレベリング・グループに属するフラッシュメモリ・モジュールをアクセスするための論理ページアドレスを前記記憶装置コントローラの内部で扱うための仮想ページアドレスに変換し、前記ウエアレベリング・グループを複数組み合わせる R A I D グループを構成する記憶装置コントローラと

を備えたフラッシュメモリを用いた記憶装置の消去回数平準化方法であって、

前記記憶装置コントローラは、

前記フラッシュメモリ・モジュール内の所定メモリ空間毎の書込みデータ容量を計数管理するステップと、

前記フラッシュメモリ・モジュール毎に、所定時の前回の平均消去回数である第 1 の平均消去回数と該所定時からの延べ書込みデータ容量を前記フラッシュメモリ・モジュール容量で割った第 2 の平均消去回数とを加算して平均消去回数を算出するステップと、

前記平均消去回数の最大値と最小値の消去回数差が所定値以上であるか否かを判定する第 1 の判定ステップと、

前記第 1 の判定ステップで、消去回数差が所定値以上の場合に、消去回数差が最大の前記フラッシュメモリ・モジュール間で、書込みデータ容量最大メモリ空間と書込みデータ容量最小メモリ空間との間でデータの入換え及び前記論理ページアドレスと前記仮想ページアドレスの対応情報を変更するステップとを含んで実行する

ことを特徴とするフラッシュメモリを用いた記憶装置の消去回数平準化方法。

【請求項 2 1】

前記記憶装置コントローラは、前記第 1 の平均消去回数を、第 2 の平均消去回数により置換するステップを含んで実行する

ことを特徴とする請求項 2 0 に記載のフラッシュメモリを用いた記憶装置の消去回数平準化方法。

【請求項 2 2】

10

20

30

40

50

1個以上のフラッシュメモリ・チップと、前記フラッシュメモリ・チップに属するブロックの消去回数を平準化するためのメモリコントローラとを有するフラッシュメモリ・モジュールと、

複数の前記フラッシュメモリ・モジュールを組み合わせるウエアレベリンググループであるウエアレベリング・グループを構成し、前記ウエアレベリング・グループに属するフラッシュメモリ・モジュールをアクセスするための論理ページアドレスを前記記憶装置コントローラの内部で扱うための仮想ページアドレスに変換し、前記ウエアレベリング・グループを複数組み合わせるR A I Dグループを構成する記憶装置コントローラと

を備えたフラッシュメモリを用いた記憶装置の消去回数平準化プログラムであって、コンピュータに、

10

前記フラッシュメモリ・モジュール内の所定メモリ空間毎の書込みデータ容量を計数管理する処理と、

前記フラッシュメモリ・モジュール毎に延べ書込みデータ容量を前記フラッシュメモリ・モジュール容量で割った平均消去回数を算出する処理と、

前記平均消去回数の最大値と最小値の消去回数差が所定値以上であるか否かを判定する第1の判定処理と、

前記第1の判定処理で、消去回数差が所定値以上の場合に、消去回数差が最大の前記フラッシュメモリ・モジュール間で、書込みデータ容量最大メモリ空間と書込みデータ容量最小メモリ空間との間でデータの入換え及び前記論理ページアドレスと前記仮想ページアドレスの対応情報を変更する処理とを

20

実行させるためのフラッシュメモリを用いた記憶装置の消去回数平準化プログラム。

【請求項23】

前記コンピュータに、

前記変更する処理を実行後、データの入換えをしていないフラッシュメモリ・モジュールが複数あるか否かを判定する第2の判定処理と、

前記第2の判定処理で、データの入換えをしていないフラッシュメモリ・モジュールが複数ある場合、第1の判定する処理に戻る処理とを

実行させるための請求項22に記載のフラッシュメモリを用いた記憶装置の消去回数平準化プログラム。

【発明の詳細な説明】

30

【技術分野】

【0001】

本発明は、複数のフラッシュメモリ・モジュールに渡った消去回数の平準化が可能なフラッシュメモリを用いた記憶装置、その消去回数平準化方法、及び消去回数平準化プログラムに関する。

【背景技術】

【0002】

記憶装置（以下、ストレージ装置という。）は、一般的に、ランダムアクセス可能な不揮発性記憶媒体を備える。ランダムアクセス可能な不揮発性記憶媒体は、例えば、磁気ディスク、光ディスク等である。また、現在、主流ストレージ装置は、小型ディスクドライブを多数備える。

40

また、半導体技術の進歩に伴って、一括消去可能な不揮発性半導体メモリが開発されている。一括消去可能な不揮発性半導体メモリは、例えば、フラッシュメモリである。フラッシュメモリは、リードオンリメモリ（ROM）のように不揮発性でありながら、リードだけでなく、ランダムアクセスメモリ（RAM）のようにライトも可能な半導体メモリである。フラッシュメモリを記憶媒体とするストレージ装置は、小型ディスクドライブを多数備えるストレージ装置に比べ、寿命、省電力、アクセス時間等に優れている。

【0003】

ここで、フラッシュメモリについて説明する。フラッシュメモリは、特性上、データを直接書換えることができない。つまり、フラッシュメモリは、記憶しているデータを書換

50

える場合、記憶している有効なデータを退避させる。次に、記憶しているデータをブロック単位で消去する。そして、データを消去したブロックにデータを書き込む。なお、ブロックは、データを一括消去する単位の記憶領域である。

【0004】

フラッシュメモリは、例えば、データを消去した場合には、「1」に、すべて設定される。このため、データの書換え時に、ビット変換して「1」を「0」に書換えることはできる。しかし、データの消去をしないで、直接的に「0」を「1」に書換えることができない。そこで、フラッシュメモリは、データの書換え時に、一旦、ブロックの全体を消去する。このように、フラッシュメモリにおけるデータの書換えは、ブロックの消去が伴う。

10

【0005】

フラッシュメモリは、ブロック消去の回数には制限がある。例えば、ブロック当たり10万回までの消去回数が保証されている。データの書換えが集中して消去回数が増大したブロックは、データ消去ができなくなり使用不能となる問題点がある。そのため、フラッシュメモリを記憶媒体として使用するストレージ装置では、特定ブロックに消去が集中しないような消去回数平準化処理が必須となる。

【0006】

特許文献1には、消去回数平準化処理方法についての記載がある。ストレージ装置は、ホストコンピュータとフラッシュメモリ間のブロック対応関係に柔軟性を与え、ホストコンピュータがアクセスする論理ブロックによってフラッシュメモリの物理ブロックが一方的に決定されることのないようなマッピング制御方式を導入した。上記従来のストレージ装置は、ホストコンピュータがアクセスする論理ブロック毎の書込み回数と、ストレージ装置が消去する物理ブロック毎の消去回数を計数管理する。そして、書込み回数の多い論理ブロックで、かつ、消去回数の多い物理ブロックと、書込み回数が少ない論理ブロックで、かつ、消去回数の少ない物理ブロックが生じた場合に、書込み回数が多い論理ブロックには消去回数の少ない物理ブロックを対応させ、書込み回数の少ない論理ブロックには消去回数の多い物理ブロックが対応するように、マッピングを変更する。

20

【特許文献1】特開平8-16482号公報(段落0038-0046、図10)

【発明の開示】

【発明が解決しようとする課題】

30

【0007】

一般的に、フラッシュメモリ・モジュール(以下、PDEVという。)は、メモリコントローラと複数のフラッシュメモリ・チップで構成され、メモリコントローラがフラッシュメモリ・チップに対して上記従来技術と同様な消去回数平準化処理を行う。大規模ストレージ装置では、記憶媒体として多数のフラッシュメモリ・モジュールを接続して大容量化することが考えられる。このとき、各フラッシュメモリ・モジュール内ではメモリコントローラにより消去回数平準化が行われる。しかし、特定のフラッシュメモリ・モジュールのみにデータ書換えが集中した場合に、書換えが集中したフラッシュメモリ・モジュールは消去回数が増大し、書換え寿命に早く達する。特定モジュールの消去回数増大を防止するためには、複数のモジュールに渡る消去回数の平準化が必要となる。

40

【0008】

また、上記従来技術の消去回数平準化方法を、多数のフラッシュメモリ・モジュールを接続したストレージ装置へ適用しようとした場合、フラッシュメモリ・モジュール内のメモリコントローラが、フラッシュメモリ・チップの物理ブロックを隠蔽しているため、ストレージ装置内の上位の記憶制御部(以下、記憶装置コントローラという。)が、物理ブロック毎の消去回数を計数管理することができない問題がある。

【0009】

さらにまた、フラッシュメモリ・モジュール内のメモリコントローラを使用しない、つまりモジュール毎に消去回数を平準化せずに、ストレージ装置全体で従来技術の消去回数平準化方法を適用する場合は、ストレージコントローラが膨大な数の物理ブロックの消去

50

回数を一元管理しなければならないため管理負荷が大きく、ストレージ装置の性能低下の要因となるおそれがある。

【 0 0 1 0 】

本発明は、前記の問題点に鑑みてなされたものであり、フラッシュメモリ物理ブロック関連情報を使用せずに、複数のフラッシュメモリ・モジュールに渡った消去回数の平準化が可能なフラッシュメモリを用いた記憶装置、その消去回数平準化方法、及び消去回数平準化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】

【 0 0 1 1 】

本発明によるフラッシュメモリを用いた記憶装置は、記憶装置コントローラと、記憶媒体として記憶装置コントローラに接続される複数のフラッシュメモリ・モジュールとを有する記憶装置であって、フラッシュメモリ・モジュールは、1個以上のフラッシュメモリ・チップと、フラッシュメモリ・チップに属するブロックの消去回数を平準化するメモリコントローラとを有し、記憶装置コントローラは、複数のフラッシュメモリ・モジュールを含むウエアレベリング・グループを構成し、ウエアレベリング・グループに属する複数のフラッシュメモリ・モジュール毎に、書き込みデータ容量に基づいてフラッシュメモリ・モジュール内の平均消去回数を算出し、フラッシュメモリ・モジュール毎の平均消去回数が所定の条件を満たした場合に、ウエアレベリング・グループに属するフラッシュメモリ・モジュール間でデータを移動することを特徴とする。

【 0 0 1 2 】

本発明によるフラッシュメモリを用いた記憶装置の消去回数平準化方法は、1個以上のフラッシュメモリ・チップと、フラッシュメモリ・チップに属するブロックの消去回数を平準化するためのメモリコントローラとを有するフラッシュメモリ・モジュールと、複数のフラッシュメモリ・モジュールを組み合わせる第1の論理グループを構成し、第1の論理グループに属するフラッシュメモリ・モジュールをアクセスするための第1のアドレスを記憶装置コントローラの内部で扱うための第2のアドレスに変換し、第1の論理グループを複数組み合わせる第2の論理グループを構成する記憶装置コントローラとを備えたフラッシュメモリを用いた記憶装置の消去回数平準化方法であって、記憶装置コントローラは、フラッシュメモリ・モジュール内の所定メモリ空間毎の書き込みデータ容量を計数管理するステップと、フラッシュメモリ・モジュール毎に延べ書き込みデータ容量をフラッシュメモリ・モジュール容量で割った平均消去回数を算出するステップと、平均消去回数の最大値と最小値の消去回数差が所定値以上であるか否かを判定する第1の判定ステップと、第1の判定ステップで、消去回数差が所定値以上の場合に、消去回数差が最大のフラッシュメモリ・モジュール間で、書き込みデータ容量最大メモリ空間と書き込みデータ容量最小メモリ空間との間でデータの入換え及び第1のアドレスと第2のアドレスの対応情報を変更するステップとを含むことを特徴とする。

【発明の効果】

【 0 0 1 3 】

本発明によれば、複数のフラッシュメモリ・モジュールに渡って消去回数を平準化できるので、ストレージ装置の記憶媒体を長寿命化できる。

【発明を実施するための最良の形態】

【 0 0 1 4 】

以下、本発明の実施の形態について図面を参照して説明する。

《概要》

本発明によるフラッシュメモリを用いた記憶装置は、記憶装置コントローラと、記憶媒体として複数のフラッシュメモリ・モジュールとを有するフラッシュメモリを用いた記憶装置であって、フラッシュメモリ・モジュール（例えば、フラッシュメモリ・モジュール P 0 0）は、1個以上のフラッシュメモリ・チップ（例えば、フラッシュメモリ・チップ 4 0 5）と、フラッシュメモリ・チップに属するブロック（例えば、ブロック 4 0 6）の消去回数を平準化するためのメモリコントローラ（例えば、メモリコントローラ M C）と

10

20

30

40

50

を有し、記憶装置コントローラ（例えば、ストレージコントローラ S C）は、複数のフラッシュメモリ・モジュールを組み合わせて第 1 の論理グループ（例えば、ウエアレベリング・グループ W 0 0）を構成し、第 1 の論理グループに属するフラッシュメモリ・モジュールにアクセスするための第 1 のアドレス（例えば、論理ページアドレス 6 0 0）を記憶装置コントローラの内部で扱うための第 2 のアドレス（例えば、仮想ページアドレス 6 0 4）に変換し、第 1 の論理グループを複数組み合わせることで第 2 の論理グループ（例えば、R A I D グループ）を構成する。

【 0 0 1 5 】

図 1 は、本発明によるストレージ装置の実施の形態の構成を示すブロック図である。ストレージ装置 1 0 0 は、ストレージコントローラ S C 及びフラッシュメモリ・モジュール P 0 0 ~ P 3 5 を備える。

10

【 0 0 1 6 】

ストレージコントローラ S C は、チャンネルアダプタ C A 0、チャンネルアダプタ C A 1、キャッシュメモリ C M 0、キャッシュメモリ C M 1、ストレージアダプタ S A 0、ストレージアダプタ S A 1、相互接続網 N W 0、及び相互接続網 N W 1 を備える。なお、チャンネルアダプタ C A 0、チャンネルアダプタ C A 1、キャッシュメモリ C M 0、キャッシュメモリ C M 1、ストレージアダプタ S A 0、及びストレージアダプタ S A 1 は、二つずつを図示しているが、いくつ備えられていてもよい。

【 0 0 1 7 】

相互接続網 N W 0 及び相互接続網 N W 1 は、例えば、スイッチ等であり、ストレージコントローラ S C を構成する装置を相互に接続する。具体的には、相互接続網 N W 0 及び相互接続網 N W 1 は、チャンネルアダプタ C A 0、キャッシュメモリ C M 0、及びストレージアダプタ S A 0 を相互に接続する。同様に、相互接続網 N W 0 及び相互接続網 N W 1 は、チャンネルアダプタ C A 1、キャッシュメモリ C M 1、及びストレージアダプタ S A 1 を相互に接続する。

20

【 0 0 1 8 】

チャンネルアダプタ C A 0 は、図 2 で後記するが、チャンネル C 0 0、チャンネル C 0 1、チャンネル C 0 2、チャンネル C 0 3 を介して、外部の上位装置（図示省略）に接続されている。同様に、チャンネルアダプタ C A 1 は、チャンネル C 1 0、チャンネル C 1 1、チャンネル C 1 2、チャンネル C 1 3 を介して、外部の上位装置（図示省略）に接続されている。なお、上位装置は、本実施の形態のストレージ装置 1 0 0 にデータを読み書きする計算機である。ストレージ装置 1 0 0 は、上位装置や他のストレージ装置等と接続する場合、ファイバチャンネルスイッチ、F C - A L (F i b r e C h a n n e l - A r b i t r a t e d L o o p)、S A S (S e r i a l A t t a c h e d S C S I) E x p a n d e r 等を介して接続する。

30

【 0 0 1 9 】

キャッシュメモリ C M 0 は、チャンネルアダプタ C A 0 及びストレージアダプタ S A 0 から受信したデータを一時的に記憶する。同様に、キャッシュメモリ C M 1 は、チャンネルアダプタ C A 1 及びストレージアダプタ S A 1 から受信したデータを一時的に記憶する。

【 0 0 2 0 】

ストレージアダプタ S A 0 は、図 3 で後記するが、フラッシュメモリ・モジュール P 0 0 等に接続されている。具体的には、ストレージアダプタ S A 0 は、チャンネル D 0 0 を介して、フラッシュメモリ・モジュール P 0 0 ~ P 0 5 に接続されている。また、ストレージアダプタ S A 0 は、チャンネル D 0 1 を介して、フラッシュメモリ・モジュール P 1 0 ~ P 1 5 に接続されている。また、ストレージアダプタ S A 0 は、チャンネル D 0 2 を介して、フラッシュメモリ・モジュール P 2 0 ~ P 2 5 に接続されている。また、ストレージアダプタ S A 0 は、チャンネル D 0 3 を介して、フラッシュメモリ・モジュール P 3 0 ~ P 3 5 に接続されている。

40

【 0 0 2 1 】

同様に、ストレージアダプタ S A 1 は、フラッシュメモリ・モジュール P 0 0 等に接続

50

されている。具体的には、ストレージアダプタ S A 1 は、チャンネル D 1 0 を介して、フラッシュメモリ・モジュール P 0 0 ~ P 0 5 に接続されている。また、ストレージアダプタ S A 1 は、チャンネル D 1 1 を介して、フラッシュメモリ・モジュール P 1 0 ~ P 1 5 に接続されている。また、ストレージアダプタ S A 1 は、チャンネル D 1 2 を介して、フラッシュメモリ・モジュール P 2 0 ~ P 2 5 に接続されている。また、ストレージアダプタ S A 1 は、チャンネル D 1 3 を介して、フラッシュメモリ・モジュール P 3 0 ~ P 3 5 に接続されている。なお、具体的には、ストレージアダプタとフラッシュメモリ・モジュールとは、ファイバチャンネルスイッチ、FC - AL、SAS Expander等を介して接続されている。

【 0 0 2 2 】

チャンネルアダプタ C A 0 及びチャンネルアダプタ C A 1、並びに、ストレージアダプタ S A 0 及びストレージアダプタ S A 1 は、保守端末 S V P に接続されている。保守端末 S V P は、ストレージ装置 1 0 0 の管理者から入力された設定情報を、チャンネルアダプタ C A 0、チャンネルアダプタ C A 1 及び/又はストレージアダプタ S A 0、ストレージアダプタ S A 1 に送信する。なお、ストレージ装置 1 0 0 は、ストレージアダプタ S A 0 及びチャンネルアダプタ C A 0 に代わって、一つのアダプタを備えていてもよい。この場合、当該アダプタが、ストレージアダプタ S A 0 及びチャンネルアダプタ C A 0 の処理を行う。

【 0 0 2 3 】

図 2 は、チャンネルアダプタの構成を示すブロック図である。チャンネルアダプタ C A 0 は、ホストチャンネル・インターフェース 2 1、キャッシュメモリ・インターフェース 2 2、ネットワーク・インターフェース 2 3、プロセッサ 2 4、ローカルメモリ 2 5、及びプロセッサ周辺制御部 2 6 を備える。

【 0 0 2 4 】

ホストチャンネル・インターフェース 2 1 は、チャンネル C 0 0、チャンネル C 0 1、チャンネル C 0 2、及びチャンネル C 0 3 を介して、外部の上位装置（図示省略）と接続するインターフェースである。また、ホストチャンネル・インターフェース 2 1 は、チャンネル C 0 0、チャンネル C 0 1、チャンネル C 0 2、及びチャンネル C 0 3 上のデータ転送プロトコルと、ストレージコントローラ S C の内部のデータ転送プロトコルとを相互に変換する。

【 0 0 2 5 】

キャッシュメモリ・インターフェース 2 2 は、相互接続網 N W 0 及び相互接続網 N W 1 と接続するインターフェースである。ネットワーク・インターフェース 2 3 は、保守端末 S V P と接続するインターフェースである。なお、ホストチャンネル・インターフェース 2 1 とキャッシュメモリ・インターフェース 2 2 とは、信号線 2 7 によって接続されている。

【 0 0 2 6 】

プロセッサ 2 4 は、ローカルメモリ 2 5 に記憶されているプログラムを実行することによって、各種処理を行う。具体的には、プロセッサ 2 4 は、上位装置と、相互接続網 N W 0、及び相互接続網 N W 1 との間のデータ転送を制御する。

【 0 0 2 7 】

ローカルメモリ 2 5 は、プロセッサ 2 4 によって実行されるプログラムを記憶する。また、ローカルメモリ 2 5 は、プロセッサ 2 4 によって参照されるテーブルを記憶する。なお、当該テーブルは、管理者によって設定又は変更される。

【 0 0 2 8 】

この場合、管理者は、テーブルの設定又はテーブルの変更に関する情報を、保守端末 S V P に入力する。保守端末 S V P は、入力された情報をネットワーク・インターフェース 2 3 を介して、プロセッサ 2 4 に送信する。プロセッサ 2 4 は、受信した情報に基づいて、テーブルを作成又は変更する。そして、プロセッサ 2 4 は、当該テーブルを、ローカルメモリ 2 5 に格納する。

【 0 0 2 9 】

プロセッサ周辺制御部 2 6 は、ホストチャンネル・インターフェース 2 1、キャッシュメ

10

20

30

40

50

メモリ・インターフェース 22、ネットワーク・インターフェース 23、プロセッサ 24、及びローカルメモリ 25間のデータ転送を制御する。プロセッサ周辺制御部 26は、例えば、チップセット等である。なお、チャンネルアダプタ CA1も、チャンネルアダプタ CA0と同一の構成である。よって、説明を省略する。

【0030】

図3は、ストレージアダプタの構成を示すブロック図である。ストレージアダプタ SA0は、キャッシュメモリ・インターフェース 31、ストレージチャンネル・インターフェース 32、ネットワーク・インターフェース 33、プロセッサ 34、ローカルメモリ 35、及びプロセッサ周辺制御部 36を備える。

【0031】

キャッシュメモリ・インターフェース 31は、相互接続網 NW0及び相互接続網 NW1と接続するインターフェースである。ストレージチャンネル・インターフェース 32は、チャンネル D00、チャンネル D01、チャンネル D02、及びチャンネル D03と接続するインターフェースである。また、ストレージチャンネル・インターフェース 32は、チャンネル D00、チャンネル D01、チャンネル D02、及びチャンネル D03上のデータ転送プロトコルと、ストレージコントローラ SCの内部のデータ転送プロトコルとを相互に変換する。なお、キャッシュメモリ・インターフェース 31とストレージチャンネル・インターフェース 32とは、信号線 37によって接続されている。ネットワーク・インターフェース 33は、保守端末 SVPと接続するインターフェースである。

【0032】

プロセッサ 34は、ローカルメモリ 35に記憶されているプログラムを実行することによって、各種処理を行う。

【0033】

ローカルメモリ 35は、プロセッサ 34によって実行されるプログラムを記憶する。また、ローカルメモリ 35は、プロセッサ 34によって参照されるテーブルを記憶する。なお、当該テーブルは、管理者によって設定又は変更される。

【0034】

この場合、管理者は、テーブルの設定又はテーブルの変更に関する情報を、保守端末 SVPに入力する。保守端末 SVPは、入力された情報をネットワーク・インターフェース 33を介して、プロセッサ 34に送信する。プロセッサ 34は、受信した情報に基づいて、テーブルを作成又は変更する。そして、プロセッサ 34は、当該テーブルを、ローカルメモリ 35に格納する。

【0035】

プロセッサ周辺制御部 36は、キャッシュメモリ・インターフェース 31、ストレージチャンネル・インターフェース 32、ネットワーク・インターフェース 33、プロセッサ 34、及びローカルメモリ 35間のデータ転送を制御する。プロセッサ周辺制御部 36は、例えば、チップセット等である。なお、ストレージアダプタ SA1も、ストレージアダプタ SA0と同一の構成である。よって、説明を省略する。

【0036】

図4は、フラッシュメモリ・モジュールの構成を示すブロック図である。フラッシュメモリ・モジュール P00は、メモリコントローラ MC及びフラッシュメモリ MEMを備える。フラッシュメモリ MEMはデータを記憶する。メモリコントローラ MCは、フラッシュメモリ MEMに対してデータを読み書きあるいは消去する。

【0037】

メモリコントローラ MCは、プロセッサ (μP) 401、インターフェース部 (I/F) 402、データ転送部 (HUB) 403、メモリ (RAM) 404、及びメモリ (ROM) 407を備える。

【0038】

フラッシュメモリ MEMは、複数のフラッシュメモリ・チップ 405を備える。フラッシュメモリ・チップ 405は、複数のブロック 406を含み、データを記憶する。ブロッ

10

20

30

40

50

ク406は、図5で後記するが、メモリコントローラMCがデータを消去する単位である。

【0039】

ブロック406は、複数のページを含む。ページは、図5で後記するが、メモリコントローラMCがデータを読み書きする単位である。なお、ページは、有効ページ、無効ページ、未使用ページ、又は不良ページのいずれかに分類される。有効ページは、有効なデータを記憶しているページである。無効ページは、無効なデータを記憶しているページである。未使用ページは、データを記憶していないページである。不良ページは、当該ページの記憶素子が壊れている等の理由によって、物理的に使用できないページである。

【0040】

インターフェース部402は、チャンネルD00を介して、ストレージコントローラSC内のストレージアダプタSA0に接続されている。また、インターフェース部402は、チャンネルD10を介して、ストレージコントローラSC内のストレージアダプタSA1に接続されている。

【0041】

インターフェース部402は、ストレージアダプタSA0及びストレージアダプタSA1からの命令を受信する。ストレージアダプタSA0及びストレージアダプタSA1からの命令は、例えば、SCSIコマンドである。

【0042】

具体的には、インターフェース部402は、ストレージアダプタSA0及びストレージアダプタSA1からデータを受信する。そして、インターフェース部402は、受信したデータをメモリ404に格納する。また、インターフェース部402は、メモリ404に格納されているデータを、ストレージアダプタSA0及びストレージアダプタSA1へ送信する。

【0043】

メモリ404は、例えば、ダイナミック型ランダムアクセスメモリであり、高速に読み書きできる。メモリ404は、インターフェース部402が送受信するデータを一時的に記憶する。また、メモリ407は不揮発性メモリであり、プロセッサ401によって実行されるプログラムを記憶する。当該プログラムは、プロセッサ401が実行可能となるように、ストレージ装置起動時にメモリ407からメモリ404へコピーされる。また、メモリ404は、プロセッサ401によって参照されるテーブルを記憶する。当該テーブルは、例えば、フラッシュメモリMEMの論理ページアドレスと物理ページアドレスとの変換テーブルである。論理ページアドレスは、フラッシュメモリ・モジュール外から（例えばストレージアダプタSA0から）、フラッシュメモリの読み書きする単位であるページをアクセスするためのアドレスである。物理ページアドレスは、メモリコントローラMCが、フラッシュメモリの読み書きする単位であるページをアクセスするためのアドレスである。

【0044】

データ転送部403は、例えばスイッチであり、プロセッサ401、インターフェース部402、メモリ404、メモリ407及びフラッシュメモリMEMを相互に接続し、それらの間のデータ転送を制御する。

【0045】

プロセッサ401は、メモリ404に記憶されているプログラムを実行することによって、各種処理を行う。例えば、プロセッサ401は、メモリ404に記憶されているフラッシュメモリの論理ページアドレスとフラッシュメモリの物理ページアドレスとの変換テーブルを参照して、フラッシュメモリMEMにデータを読み書きする。また、プロセッサ401は、フラッシュメモリ・モジュール内のブロック406に対して、リクラメーション処理（ブロック再生処理）及びウエアレベリング処理（消去回数平準化処理）を行う。

【0046】

リクラメーション処理（ブロック再生処理）は、未使用ページが少なくなったブロック

10

20

30

40

50

を使用できるように、ブロック406内の無効ページを未使用ページに再生する処理である。すなわち、リクラメーション処理の対象となるブロック(対象ブロック)406内には、有効ページ、無効ページ、及び未使用ページが含まれ、多くの無効ページが存在しているとする。この場合に、未使用ページを増加させるには、無効ページを消去する必要がある。しかしながら、消去は、ページ単位ではなく、ブロック単位でしか消去できない。このため、有効ページを空きのあるブロックに複写し、その後、対象ブロックを消去して、ブロックを再生する処理が必要となる。具体的には、プロセッサ401は、リクラメーション処理の対象となるブロック(対象ブロック)406内の有効ページに記憶されているデータを、未使用ブロックへ複写する。そして、プロセッサ401は、データを複写した未使用ブロックの論理ブロック番号を、対象ブロックの論理ブロック番号に変更する。そして、対象ブロックのデータをすべて消去し、リクラメーション処理を完了する。

10

【0047】

例えば、プロセッサ401がブロック406にデータを書き込むと、ブロック406内の未使用ページが少なくなる。そして、ブロック406内の未使用ページが足りなくなると、プロセッサ401は、当該ブロック406へデータを書き込めなくなる。そこで、プロセッサ401は、当該ブロック406をリクラメーションすることによって、無効ページを未使用ページに再生する。

【0048】

また、ウェアレベリング処理(消去回数平準化処理)は、それぞれのブロック406のデータ消去回数を平準化する処理である。これによって、フラッシュメモリMEMの寿命を長くできる。なぜなら、フラッシュメモリMEMは、データの消去回数が増えると、寿命になるからである。フラッシュメモリMEMは、一般的に、約一万回から十万回のデータの消去が保証されている。なお、他のフラッシュメモリ・モジュールP01~P35も、フラッシュメモリ・モジュールP00と同一の構成である。よって、説明を省略する。

20

【0049】

図5は、フラッシュメモリ・モジュールのブロックの構成を示す説明図である。フラッシュメモリ・モジュールP00のブロック406は、複数のページ501を含む。ブロック406は、一般的に、数十程度のページ501(例えば、32ページ、64ページ等)を含む。

30

【0050】

ページ501は、メモリコントローラMC等がデータを読み出しあるいは書き込む単位である。例えば、NAND型フラッシュメモリでは、メモリコントローラMC等は、20~30 μ s弱/ページでデータを読み出し、0.2~0.3ms/ページでデータを書き込む。また、メモリコントローラMC等は、2~4ms/ブロックの速度でデータを消去する。

【0051】

ページ501は、データ部502及び冗長部503を含む。ページ501は、例えば、512バイトのデータ部502及び16バイトの冗長部503を含む。データ部502は、通常のデータを記憶する。

40

【0052】

冗長部503は、当該ページ501の管理情報及びエラー訂正情報を記憶する。管理情報は、オフセットアドレス及びページステータスを含む。なお、オフセットアドレスは、当該ページ501が属するブロック406内における相対的なアドレスである。また、ページステータスは、当該ページ501が有効ページ、無効ページ、未使用ページ、又は処理中のページのいずれであることを示す。エラー訂正情報は、当該ページ501のエラーを検出及び訂正するための情報であり、例えば、ハミングコードである。

【0053】

図6は、本発明によるストレージ装置の実施の形態における論理グループの構成とアドレス変換の階層を示す説明図である。なお、図6のストレージ装置のハードウェア的な構

50

成は、図1のストレージ装置と同様であるが、簡単のため、フラッシュメモリ・モジュールP00～P35と接続するストレージコントローラSCのチャンネルはD00、チャンネルD01、チャンネルD02、及びチャンネルD03のみを示している（チャンネルD10、チャンネルD11、チャンネルD12、チャンネルD13は省略）。

【0054】

本発明のストレージ装置100では、同一チャンネル上に接続したフラッシュメモリ・モジュールを組み合わせてウエアレベリング・グループ(WDEV)を構成する。例えば、チャンネルD00上のフラッシュメモリ・モジュールP00、フラッシュメモリ・モジュールP01、フラッシュメモリ・モジュールP02、及びフラッシュメモリ・モジュールP03により、ウエアレベリング・グループW00を構成する。同様に、チャンネルD01上のフラッシュメモリ・モジュールP10、フラッシュメモリ・モジュールP11、フラッシュメモリ・モジュールP12、及びフラッシュメモリ・モジュールP13により、ウエアレベリング・グループW10を構成する。チャンネルD02上のフラッシュメモリ・モジュールP20、フラッシュメモリ・モジュールP21、フラッシュメモリ・モジュールP22、及びフラッシュメモリ・モジュールP23により、ウエアレベリング・グループW20を構成する。チャンネルD03上のフラッシュメモリ・モジュールP30、フラッシュメモリ・モジュールP31、フラッシュメモリ・モジュールP32、及びフラッシュメモリ・モジュールP33により、ウエアレベリング・グループW30を構成する。

【0055】

各フラッシュメモリ・モジュールは、ストレージコントローラSCからモジュールごとに論理ページアドレスでアクセスすることができる。例えば、チャンネルD00上のフラッシュメモリ・モジュールP00、フラッシュメモリ・モジュールP01、フラッシュメモリ・モジュールP02、及びフラッシュメモリ・モジュールP03は、各モジュールの論理ページアドレス600でアクセスする。同様に、チャンネルD01上のフラッシュメモリ・モジュールP10、フラッシュメモリ・モジュールP11、フラッシュメモリ・モジュールP12、及びフラッシュメモリ・モジュールP13は各モジュールの論理ページアドレス601でアクセスし、チャンネルD02上のフラッシュメモリ・モジュールP20、フラッシュメモリ・モジュールP21、フラッシュメモリ・モジュールP22、及びフラッシュメモリ・モジュールP23は各モジュールの論理ページアドレス602でアクセスし、チャンネルD03上のフラッシュメモリ・モジュールP30、フラッシュメモリ・モジュールP31、フラッシュメモリ・モジュールP32、及びフラッシュメモリ・モジュールP33は、各モジュールの論理ページアドレス603でアクセスする。

【0056】

ストレージコントローラSCは、同じウエアレベリング・グループに属する複数のフラッシュメモリ・モジュールの論理ページアドレスを纏めて、仮想ページアドレスに変換する。例えば、ウエアレベリング・グループW00に属するフラッシュメモリ・モジュールP00～P03の論理ページアドレス600を纏めて、仮想ページアドレス604に変換する。同様に、ウエアレベリング・グループW10に属するフラッシュメモリ・モジュールP10～P13の論理ページアドレス601を纏めて、仮想ページアドレス605に変換し、ウエアレベリング・グループW20に属するフラッシュメモリ・モジュールP20～P23の論理ページアドレス602を纏めて、仮想ページアドレス606に変換し、ウエアレベリング・グループW30に属するフラッシュメモリ・モジュールP30～P33の論理ページアドレス603を纏めて、仮想ページアドレス607に変換する。

【0057】

このように、ストレージコントローラSCは、論理ページアドレスを仮想ページアドレスに変換している。このため、消去回数平準化のためフラッシュメモリ・モジュール間でデータ移動を実施されると、論理ページアドレスが変更されるが、上位系のストレージコントローラSCは、論理ページアドレスと仮想ページアドレスとのマッピングを変更できるので、矛盾無くデータにアクセスできるようになる。

【0058】

10

20

30

40

50

本発明の実施の形態のストレージ装置100では、複数のウエアレベリング・グループを組み合わせてRAIDグループ(VDEV)を構成する。図6では、4個のウエアレベリング・グループW00、ウエアレベリング・グループW10、ウエアレベリング・グループW20、及びウエアレベリング・グループW30を組み合わせて、RAIDグループV00を構成している。一つのRAIDグループを構成する各ウエアレベリング・グループの仮想ページアドレス空間における容量は、同じである。一つあるいは複数のRAIDグループ内の領域を組み合わせて、論理ボリューム608を構成する。論理ボリューム608は、ストレージコントローラSCが上位装置に対して見せる記憶空間である。

【0059】

また、チャンネルD00上のフラッシュメモリ・モジュールP04及びフラッシュメモリ・モジュールP05は、予備グループ(YDEV)Y00を構成する。同様に、チャンネルD01上のフラッシュメモリ・モジュールP14及びフラッシュメモリ・モジュールP15は、予備グループY10を、チャンネルD02上のフラッシュメモリ・モジュールP24及びフラッシュメモリ・モジュールP25は、予備グループY20を、チャンネルD03上のフラッシュメモリ・モジュールP34及びフラッシュメモリ・モジュールP35は、予備グループY30を構成する。モジュール交換の詳細は後記する。

【0060】

図7は、本発明によるストレージ装置の実施の形態におけるRAIDグループの構成を示す説明図である。RAIDグループ720は、ウエアレベリング・グループ700、ウエアレベリング・グループ701、ウエアレベリング・グループ702、及びウエアレベリング・グループ703からなる、RAIDレベル5のRAIDグループである。例えば、ウエアレベリング・グループ700は、フラッシュメモリ・モジュール730とフラッシュメモリ・モジュール731からなる。なお、RAIDは、機能によってRAIDレベル0、RAIDレベル1のようにレベル分けされている。

【0061】

RAIDグループ721は、ウエアレベリング・グループ704及びウエアレベリング・グループ705からなる、RAIDレベル1のRAIDグループである。同様にRAIDグループ722は、ウエアレベリング・グループ706及びウエアレベリング・グループ707からなる、RAIDレベル1のRAIDグループである。

【0062】

本発明のストレージ装置100では、RAIDレベル0、RAIDレベル1、RAIDレベル3、RAIDレベル5、RAIDレベル6、あるいはRAIDレベル1+0の場合に、RAIDグループを構成する各ウエアレベリング・グループの論理ページアドレス空間における容量を等しく設定する。その容量は、上限を(1)式で、下限を(2)式で与えられる計算値で設定する。すなわち、フラッシュメモリ・モジュールの連続書込み速度とシステム製品寿命との積をフラッシュメモリの書換え寿命(書換え保証回数)で割った商を第2の値(上限)とし、フラッシュメモリ・モジュールのシステム運用時の実効的書込み速度とシステム製品寿命との積をフラッシュメモリの書換え寿命(書換え保証回数)で割った商を第1の値(下限)とし、各ウエアレベリング・グループの論理ページアドレス空間における容量(第1の論理グループの容量)を、第1の値以上、かつ、第2の値以下に設定する。例えば、システム製品寿命は5年から10年程度であり、書換え寿命は1万回から10万回程度である。(2)式の実効的書込み速度とは、上位機装置からストレージ装置100への書込みアクセス割合を考慮した場合の、実効的な書込み速度である。

【0063】

10

20

30

40

【数 1】

$$\text{ウェアレベリング・グループ容量値(上限)} = \frac{\text{モジュール当たり連続書込み速度} \times \text{システム製品寿命}}{\text{書換え寿命}} \dots (1)$$

10

$$\text{ウェアレベリング・グループ容量値(下限)} = \frac{\text{モジュール当たり実効的書込み速度} \times \text{システム製品寿命}}{\text{書換え寿命}} \dots (2)$$

【0064】

上記(1)式及び(2)式で与えられる範囲内にウェアレベリング・グループの容量を設定する。ウェアレベリング・グループ内フラッシュメモリ・モジュール間で消去回数平準化を行うことにより、ストレージ装置100のシステム製品寿命期間内でフラッシュメモリ・モジュールの書換え寿命が保証される。

20

【0065】

RAIDグループ723は、ウェアレベリング・グループ708、ウェアレベリング・グループ709、ウェアレベリング・グループ710、及びウェアレベリング・グループ711からなる、RAIDレベル4のRAIDグループである。そしてウェアレベリング・グループ708、ウェアレベリング・グループ709、及びウェアレベリング・グループ710は、データ格納用のウェアレベリング・グループであり、ウェアレベリング・グループ711は、パリティ(冗長情報)格納用のウェアレベリング・グループである。パリティ格納用のウェアレベリング・グループは、他のデータ格納用のウェアレベリング・グループに対してデータの更新回数が多い。従ってRAIDレベル4では、RAIDグループ内での消去回数平準化のため、パリティ格納用のウェアレベリング・グループの論理ページアドレス空間における容量値を、データ格納用ウェアレベリング・グループ容量値より大きく設定する。例えば、RAIDグループを構成するウェアレベリング・グループの数がn個の場合、パリティ格納用のウェアレベリング・グループは、データ格納用のウェアレベリング・グループに対して、1倍以上、かつ、(n-1)倍以下の論理ページアドレス空間を持つように設定する。

30

【0066】

また、図示していないが、RAIDレベル2の場合も冗長情報格納用のウェアレベリング・グループは、データ格納用のウェアレベリング・グループより更新回数が多い。RAIDレベル2で、例えば10D4Pならば、冗長情報格納用のウェアレベリング・グループの論理ページアドレス空間における容量値をデータ格納用ウェアレベリング・グループ容量値の1倍以上、かつ、 $10/4 = 2.5$ 倍以下に設定する。また、25D5Pならば、冗長情報格納用のウェアレベリング・グループの論理ページアドレス空間における容量値を、データ格納用ウェアレベリング・グループ容量値の1倍以上、かつ、 $25/5 = 5$ 倍以下に設定する。

40

【0067】

言い換えると、RAIDレベル2又はRAIDレベル4の場合、データ格納用のウェアレベリング・グループ数をn、冗長情報格納用のウェアレベリング・グループ数をmとすると、冗長情報格納用ウェアレベリング・グループは、データ格納用ウェアレベリング・グループに対して、1倍以上、かつ、 n/m 倍以下の論理ページアドレス空間容量を持つ

50

ように設定する。

【 0 0 6 8 】

このように、ストレージコントローラ S C の R A I D グループは、ウエアレベリング・グループを組み合わせる R A I D グループを構成している。すなわち、ストレージコントローラ S C の管理単位は、ウエアレベリング・グループを考慮して R A I D グループで管理しているため、各ウエアレベリング・グループ内の論理ページアドレスと仮想ページアドレスの対応に係わらず、各ウエアレベリング・グループの仮想ページアドレスは、独立することになる。これにより、ストレージコントローラ S C は、複数、かつ、レベルの異なる R A I D グループを接続することができる。

【 0 0 6 9 】

図 8 は、フラッシュメモリ・モジュールとハードディスクドライブをストレージコントローラ S C に接続した例を示す構成図である。フラッシュメモリ・モジュール 8 1 0、フラッシュメモリ・モジュール 8 1 1、及びフラッシュメモリ・モジュール 8 1 2 でウエアレベリング・グループ 8 3 0 を構成する。フラッシュメモリ・モジュール 8 1 3、フラッシュメモリ・モジュール 8 1 4、フラッシュメモリ・モジュール 8 1 5 でウエアレベリング・グループ 8 3 1 を構成する。ウエアレベリング・グループ 8 3 0 とウエアレベリング・グループ 8 3 1 を組み合わせる、R A I D グループ 8 4 0 を構成する。

【 0 0 7 0 】

図 6 の場合と同様に、ストレージコントローラ S C は、フラッシュメモリ・モジュール 8 1 0、フラッシュメモリ・モジュール 8 1 1、及びフラッシュメモリ・モジュール 8 1 2 をアクセスするため、論理ページアドレス 8 0 0 を仮想ページアドレス 8 0 2 へ変換する。また、ストレージコントローラ S C は、フラッシュメモリ・モジュール 8 1 3、フラッシュメモリ・モジュール 8 1 4、及びフラッシュメモリ・モジュール 8 1 5 をアクセスするため、論理ページアドレス 8 0 1 を仮想ページアドレス 8 0 3 へ変換する。

【 0 0 7 1 】

ハードディスクドライブ 8 2 0 及びハードディスクドライブ 8 2 3 を組み合わせる R A I D グループ 8 4 1 を構成する。同様に、ハードディスクドライブ 8 2 1 及びハードディスクドライブ 8 2 4 を組み合わせる R A I D グループ 8 4 2 を、ハードディスクドライブ 8 2 2 及びハードディスクドライブ 8 2 5 を組み合わせる R A I D グループ 8 4 3 を、それぞれ構成する。ストレージコントローラ S C は、論理ブロックアドレス 8 0 4 あるいは論理ブロックアドレス 8 0 5 で各ハードディスクドライブをアクセスする。ハードディスクドライブからなる R A I D グループにおいては、消去回数平準化が不要なため、ウエアレベリング・グループを定義しない。フラッシュメモリ・モジュールからなる R A I D グループに対してのみ、ウエアレベリング・グループを定義し、かつ、論理ページアドレスと仮想ページアドレスの変換を行う。

【 0 0 7 2 】

ストレージコントローラ S C は、装置起動時あるいは記憶媒体を接続した際に、記憶媒体がフラッシュメモリ・モジュールかハードディスクドライブかによって、アドレス変換の要否を判断及び R A I D グループの構成方法の決定等、制御方法を変更する。

【 0 0 7 3 】

ストレージコントローラ S C は、フラッシュメモリ・モジュールからなる R A I D グループ 8 4 0 とハードディスクドライブからなる R A I D グループ 8 4 1 ~ 8 4 3 との一方あるいは両方の領域を組み合わせる、論理ボリューム 8 0 8 を構成する。フラッシュメモリ・モジュールからなる記憶領域とハードディスクドライブから記憶領域の使い分け例としては、読み出しアクセスが多く、あまり更新されないデータを主にフラッシュメモリ・モジュールに記憶し、更新する頻度が高いデータをハードディスクドライブに記憶することが考えられる。フラッシュメモリ・モジュールはハードディスクドライブに対してアクセス速度が高速なため、このようなデータのアクセス特性によって記憶領域を使い分けることにより、ストレージ装置が高性能化する。

【 0 0 7 4 】

10

20

30

40

50

次に、動作について図面を参照しながら詳細に説明する。

図9から図14において、本発明の実施の形態のストレージ装置における消去回数平準化方法を説明する。本方法は、複数のフラッシュメモリ・モジュールに渡って消去回数を平準化する点が特徴である。

【0075】

図9は、複数のフラッシュメモリ・モジュール間の消去回数平準化方法を示すフローチャートである。簡単のため、平準化の対象となるウェアレベリング・グループW00は、2個のフラッシュメモリ・モジュールP00とフラッシュメモリ・モジュールP04からなる場合を考える。

【0076】

図10は、消去回数平準化に伴うデータ入換え前の、仮想ページアドレスと論理ページアドレスの変換テーブルを示す説明図である。

【0077】

図11は、消去回数平準化に伴うデータ入換え後の、仮想ページアドレスと論理ページアドレスの変換テーブルを示す説明図である。

【0078】

図10と図11において、仮想ページアドレスと論理ページアドレスの対応関係と共に、オフセット値が示されている。本発明の実施の形態のストレージ装置では、論理ページアドレス空間の大きさ(データ長)を、対応する仮想ページアドレス空間の大きさ(データ長)より大きく設定する。そして、論理ページアドレス空間において、有効データが先頭アドレス側に書き込まれていて、末尾アドレス側に空き領域が存在する場合をオフセット値0、有効データが末尾アドレス側に書き込まれていて、先頭アドレス側に空き領域が存在する場合をオフセット値1と表示する。空き領域の大きさはフラッシュメモリのページのデータ部の整数倍(1倍以上)であり、フラッシュメモリ・モジュール内のメモリコントローラがフラッシュメモリに対して一度に書き込み可能なデータの大きさとする。

【0079】

図12は、ストレージコントローラにおいて管理するフラッシュメモリ・モジュール毎の消去回数管理テーブルを示す説明図である。ストレージコントローラSCは、フラッシュメモリ・モジュール内のデータ入換の単位となる領域ごとに、延べ書き込み容量を記録する。(3)式のように、前回の平均消去回数とモジュール内論理ページアドレス空間に対する所定期間の延べ書き込み容量をモジュールの論理ページアドレス空間容量で割った平均処理回数を加算することにより、そのモジュール内フラッシュメモリの平均消去回数を算出することができる。

【数2】

$$\text{平均消去回数} = \text{前回の値} + \frac{\sum \text{論理ページアドレス空間毎の所定期間の延べ書き込み容量値}}{\text{モジュール容量値}} \dots (3)$$

【0080】

図12の管理テーブルでは、2個の平均消去回数値を記録している。一つは前回の平準化処理実行時の平均消去回数(f00、f04)であり、もう一つは現在までの平均消去回数(e00、e04)である。論理ページアドレス空間の所定領域毎に計数管理する延べ書き込み容量値は、前回平準化処理実行時から現在までの容量値を記憶する。現在の平均消去回数値は、(3)式により自動的に計算される。このように期間を分けて延べ書き込み容量を計数管理することにより、論理ページアドレス空間領域における最近のアクセス頻度を知ることができる。また、管理テーブルには、移動フラグが設けられており、データ入換え前の場合、移動フラグには0が設定されており、データ入換え後、移動フラグには1が設定される。なお、(3)式は、期間を分けて延べ書き込み容量を計数管理した場合であるが、期間をわけずに通期での延べ書き込み容量を計数管理した場合は、(4)式のように表現できる。(3)式と(4)式とは、平均消去回数の計算結果は同一となる。

10

20

30

40

50

【数3】

$$\text{平均消去回数} = \frac{\sum \text{論理ページアドレス空間毎の
通期の延べ書込み容量値}}{\text{モジュール容量値}} \dots (4)$$

【0081】

図10又は図11のアドレス変換テーブル及び図12の平均消去回数管理テーブルは、電源断等の障害時やシステム起動時間外でも保持する必要がある。従って、ストレージコントローラSCは、各フラッシュメモリ・モジュールの所定領域に、各フラッシュメモリ・モジュールに関連するアドレス変換テーブル及び各モジュールの平均消去回数管理テーブル関連データを保存する。

10

【0082】

図9において、ストレージコントローラSCは、イベントの発生時、例えば、ウェアレベリング・グループ(WDEV)当たりの延べ書込み容量が所定値に達するか、あるいは所定期間毎に、平均消去回数平準化処理の実行を開始する。また、ウェアレベリング・グループ内のフラッシュメモリ・モジュールの移動フラグを0に設定する(ステップS901)。

【0083】

次に、ストレージコントローラSCは、図12の平均消去回数管理テーブルから、移動フラグが0で、かつ、平均消去回数の最大値と最小値をチェックする(ステップS902)。

20

【0084】

ストレージコントローラSCは、平均消去回数の最大値と最小値の消去回数差が所定値以上か否かを判定する(ステップS903)。消去回数差が所定値以上であればステップS904へ進む。消去回数差が所定値以上でなければ、終了する。

【0085】

次にストレージコントローラSCは、図12から、平均消去回数が最大であるフラッシュメモリ・モジュール(PDEV)における延べ書込み容量が最大の論理ページアドレス空間と、平均消去回数が最小であるフラッシュメモリ・モジュール(PDEV)における延べ書込み容量が最小の論理ページアドレス空間とを選択する(ステップS904)。

30

【0086】

ストレージコントローラSCは、仮想ページ 論理ページアドレス変換テーブルの状態欄を入換え中表示にする。具体的には、ストレージコントローラSCは、選択した2つの論理ページアドレス空間の間で、データ入換えと、仮想ページアドレスとのマッピング変更するため、図10の変換テーブルの状態欄に、入換え操作中を示す値を入力する(ステップS905)。ストレージコントローラSCは、入換え操作中表示がある記憶領域に対してはアクセスを一時的に保留し、データ入換え及びマッピング変更完了を待って、再アクセスする。その際、上位装置からの書込みデータはストレージコントローラSC内のキャッシュメモリに保持する。

【0087】

40

次に、ストレージコントローラSCは、上記2つの論理ページアドレス空間でデータ入換えを行う(ステップS906)。入換え操作の詳細は後記する。

【0088】

ストレージコントローラSCは、入換え操作完了後、図12の消去回数管理テーブルにおいて、入換え対象の領域が属するフラッシュメモリ・モジュール(PDEV)の、平均消去回数(前回)欄に現在の平均消去回数値を置換し、かつ、延べ書込み容量値をクリアする(ステップS907)。これにより入換え直後の平均消去回数値は、前回値と同じになる。

【0089】

ストレージコントローラSCは、図11のように、仮想ページアドレスと論理ページア

50

ドレスのマッピングとオフセット値を変更し、状態欄をクリアし、移動フラグを1に設定する(ステップS908)。図12の場合は、平準化の対象となるウエアレベリング・グループW00は、2個のフラッシュメモリ・モジュールP00とフラッシュメモリ・モジュールP04からなっているので、データ入換えがあると、フラッシュメモリ・モジュールP00とフラッシュメモリ・モジュールP04の移動フラグは、共に1に設定される。

【0090】

ストレージコントローラSCは、移動フラグ0の複数のフラッシュメモリ・モジュール(PDEV)があるか否かを判定する(ステップS909)。移動フラグ0の複数のフラッシュメモリ・モジュールがない場合、平準化処理を終了し、移動フラグ0の複数のフラッシュメモリ・モジュールがある場合は、ステップS902に戻る。図12の場合は、2個のフラッシュメモリ・モジュールP00とフラッシュメモリ・モジュールP04の移動フラグは1に設定されている場合、平準化処理を終了する。例えば、ウエアレベリング・グループ内のフラッシュメモリ・モジュールが4個以上から構成されている場合は、さらに残りの2個のフラッシュメモリ・モジュールのデータ入換えができるか否かを、ステップS902でチェックする。

10

【0091】

なお、消去回数平準化処理を実行する度に仮想ページアドレスと論理ページアドレスのマッピングが変更され、しかも、平均消去回数値が更新されるので、フラッシュメモリ・モジュールの所定領域に格納したアドレス管理テーブル(図10又は図11)と、平均消去回数管理テーブルの図12を、消去回数平準化処理を実行する度に更新する必要がある。

20

【0092】

次にステップS906におけるデータ入換え操作の詳細を説明する。

図13は、消去回数平準化に伴うデータ入換え前の仮想ページアドレスと論理ページアドレスのマッピングを示す説明図である。例として、図13の仮想ページアドレス空間のデータ領域1301と仮想ページアドレス空間のデータ領域1302のデータ入換えとマッピング変更する場合を説明する。仮想ページアドレス空間のデータ領域1301は、論理ページアドレス空間のデータ領域1303に対応している。また、仮想ページアドレス空間のデータ領域1302は、論理ページアドレス空間のデータ領域1304に対応している。論理ページアドレス空間には、データ領域間に空き領域があり、データ領域の後に空き領域がある場合がオフセット値0であり、データ領域の前に空き領域がある場合がオフセット値1と設定される。例えば、図13の1303は、データ領域の前に空き領域があるので、オフセット値1であり、1304は、データ領域の後に空き領域があるので、オフセット値0である。よって、論理ページアドレス空間の全体のメモリ容量は、仮想ページアドレス空間の全体のメモリ容量より大きく設定されている。

30

【0093】

図14は、消去回数平準化に伴うデータ入換え後の仮想ページアドレスと論理ページアドレスのマッピングを示す説明図である。仮想ページアドレス空間のデータ領域1401は、論理ページアドレス空間のデータ領域の1404に対応している。また、仮想ページアドレス空間のデータ領域1402は、論理ページアドレス空間のデータ領域の1403に対応している。1403はオフセット値0であり、1404はオフセット値0である。

40

【0094】

図15から図24は、データ入換えの例として、オフセット値0のデータ領域とオフセット値1の領域でのデータ入換えを段階的に示した図である。図15から図24において、左側がオフセット値0の論理ページアドレス空間である。ここでは、論理ページアドレス空間がE、F、G、H、-と示しているように5分割されており、EからHが有効データの書き込まれている領域を示し、-が空き領域を示す。

【0095】

図15から図24において、右側がオフセット値1の論理ページアドレス空間である。ここでは、論理ページアドレス空間が-、A、B、C、Dと示しているように5分割され

50

ており、AからDが有効データの書き込まれている領域を、-が空き領域を示す。

【0096】

図15は、データ入換え操作の入換え前の初期状態を示す説明図である。左側がオフセット値0、右側がオフセット値1の状態である。

【0097】

図16は、データ入換え操作の入換え途中の状態を示す説明図である。左側(オフセット値0)の論理ページアドレス空間Eのデータを、右側(オフセット値1)の空き領域へ上書きする。

【0098】

図17は、データ入換え操作の入換え途中の状態を示す説明図である。右側の論理ページアドレス空間Aのデータを、左側の元の論理ページアドレス空間Eへ上書きする。

10

【0099】

図18は、データ入換え操作の入換え途中の状態を示す説明図である。左側の論理ページアドレス空間Fのデータを、右側の元の論理ページアドレス空間Aへ上書きする。

【0100】

図19は、データ入換え操作の入換え途中の状態を示す説明図である。右側の論理ページアドレス空間Bのデータを、左側の元の論理ページアドレス空間Fへ上書きする。

【0101】

図20は、データ入換え操作の入換え途中の状態を示す説明図である。左側の論理ページアドレス空間Gのデータを、右側の元の論理ページアドレス空間Bへ上書きする。

20

【0102】

図21は、データ入換え操作の入換え途中の状態を示す説明図である。右側の論理ページアドレス空間Cのデータを、左側の元の論理ページアドレス空間Gへ上書きする。

【0103】

図22は、データ入換え操作の入換え途中の状態を示す説明図である。左側の論理ページアドレス空間Hのデータを、右側の元の論理ページアドレス空間Cへ上書きする。

【0104】

図23は、データ入換え操作の入換え途中の状態を示す説明図である。右側の論理ページアドレス空間Dのデータを、左側の元の論理ページアドレス空間Hへ上書きする。

【0105】

30

図24は、データ入換え操作の入換え後の最終状態を示す説明図である。左側がオフセット値0、右側がオフセット値0の状態である。

【0106】

基本的にフラッシュメモリは、物理アドレス空間において上書き操作ができない半導体デバイスである。すなわち、データを書換えたい場合、物理アドレス空間では、未使用ページにデータを書込み、元のページを無効ページに設定されるのであって、元ページのデータを上書きする操作が実施されていない。

【0107】

本発明の実施の形態によれば、論理ページアドレス空間での操作であるため、論理ページ上での上書き操作が可能となっている。このように上書き操作によるデータ入換えに基づいて消去回数平準化を行うことができる。

40

【0108】

図25は、データ入換え操作のデータ入換え前後でのオフセット値の遷移を示す説明図である。オフセット値0同士の論理ページアドレス空間でデータ入換えを行うと、データ入換え後のオフセット値は0と1になる。オフセット値0の論理ページアドレス空間とオフセット値1の論理ページアドレス空間との間でデータ入換えを行うと、データ入換え後のオフセット値は0と0になる。オフセット値1の論理ページアドレス空間同士でのデータ入換えを行うと、データ入換え後のオフセット値は1と0になる。

【0109】

次にデータ入換え操作の詳細をフローチャートを用いて説明する。

50

図26は、図15から図24で説明したオフセット値0の論理ページアドレス空間とオフセット値1の論理ページアドレス空間との間のデータ入換え手順を示すフローチャートである。ストレージコントローラSCは、オフセット値0の論理ページアドレス空間とオフセット値1の論理ページアドレス空間をデータ入換え対象とする(ステップS2601)。

【0110】

ストレージコントローラSCは、対象論理ページアドレス空間を n 分割し、 $i = 1$ に設定する(ステップS2602)。図15の例では $n = 5$ であり、分割した($n - 1$)部分に有効データが書き込まれており、残りの1部分が空き領域である。ストレージコントローラSCは、オフセット値0の論理ページアドレス空間からオフセット値1の論理ページアドレス空間へのデータ移動(ステップS2603)と、オフセット値1の論理ページアドレス空間からオフセット値0の論理ページアドレス空間へのデータ移動(ステップS2604)をし、 i に1を加算する(ステップS2605)。そして、ストレージコントローラSCは、 $i = n$ であるか否かを判定する(ステップS2606)。 $i = n$ と判定されない場合は、ステップS2603に戻る。 $i = n$ の場合に判定された場合は、データ入換えを終了する。

10

【0111】

図27は、オフセット値0の論理ページアドレス空間とオフセット値0の論理ページアドレス空間との間のデータ入換え手順を示すフローチャートである。オフセット値0の論理ページアドレス空間とオフセット値0の論理ページアドレス空間をデータ入換え対象とする。(ステップS2701)。

20

【0112】

ストレージコントローラSCは、対象論理ページアドレス空間を n 分割し、 $i = n$ に設定する(ステップS2702)。分割した($n - 1$)部分に有効データが書き込まれており、残りの1部分が空き領域である。 $i = 1$ になるまで、分割した単位ごとのデータ入換えを繰り返す(ステップS2703～ステップS2706)。そして、ステップS2706において、 $i = 1$ と判定された場合、データの入換え操作が完了する。

【0113】

図28は、オフセット値1の論理ページアドレス空間とオフセット値1の論理ページアドレス空間との間のデータ入換え手順を示すフローチャートである。オフセット値1の論理ページアドレス空間とオフセット値1の論理ページアドレス空間をデータ入換え対象とする(ステップS2801)。

30

【0114】

対象データ領域を n 分割し、 $i = 2$ に設定する(ステップS2802)。分割した($n - 1$)部分に有効データが書き込まれており、残りの1部分が空き領域である。 $i > n$ になるまで、分割した単位ごとのデータ入換えを繰り返す(ステップS2803～ステップS2806)。そして、ステップS2806において、 $i > n$ に判定された場合、データの入換え操作が完了する。

【0115】

図29から図34は、本発明の実施の形態における他の消去回数平準化方法を説明する図である。前述の方法(図13から図28)ではデータ入換え用の空き領域をフラッシュメモリ・モジュール内に分散配置したが、ここではモジュール毎に纏めて配置した方法を述べる。

40

【0116】

図29は、データ入換え前の仮想ページアドレスと論理ページアドレスのマッピングを示す説明図である。図30は、データ入換え後の仮想ページアドレスと論理ページアドレスのマッピングを示す説明図である。簡単のため、平準化対象のウェアレベリング・グループW00は、2個のフラッシュメモリ・モジュールP00とフラッシュメモリ・モジュールP04からなる場合を考える。データ入換え前の状態(図29)において、フラッシュメモリ・モジュールP00及びフラッシュメモリ・モジュールP04の論理ページアド

50

レス空間は、アドレス A C 0 以上アドレス A C 4 未満の空間がデータ領域であり、アドレス A C 4 以上の領域 (2 9 0 3、2 9 0 4) がデータ入換え用の空き領域となっている。この空き領域の大きさ (データ長) は、消去回数平準化のためにデータ入換えを行うデータ領域と同じ大きさである。

【 0 1 1 7 】

図 3 1 は、データ入換え前の仮想ページアドレスと論理ページアドレスの変換テーブルを示す説明図である。図 3 2 は、データ入換え後の仮想ページアドレスと論理ページアドレスの変換テーブルを示す説明図である。データ入換え用の空き領域を纏めて配置したため、図 1 0、図 1 1 の変換テーブルで必要であったオフセット値管理は不要となる。その代わりに、空き領域の位置を管理する必要がある。

10

【 0 1 1 8 】

図 3 3 は、データ入換え前の空き領域管理テーブルを示す説明図である。図 3 4 は、データ入換え後の空き領域管理テーブルを示す説明図である。空き領域管理テーブルは、フラッシュメモリ・モジュール毎に、空き領域の先頭論理ページアドレスと大きさ (データ長) を管理する。

【 0 1 1 9 】

図 2 9 において、仮想ページアドレス空間におけるデータ領域 2 9 0 1 とデータ領域 2 9 0 2 のデータを入換え、かつ、仮想ページアドレスと論理ページアドレスのマッピングを変更する手順を説明する。図 3 1 の論理ページアドレスと仮想ページアドレス変換テーブルから、仮想ページアドレス空間におけるデータ領域 2 9 0 1 に対応する論理ページアドレス空間は、データ領域 2 9 0 5 であり、仮想ページアドレス空間におけるデータ領域 2 9 0 2 に対応する論理ページアドレス空間は、データ領域 2 9 0 6 であることが分かる。また、図 3 3 の空き領域管理テーブルから、フラッシュメモリ・モジュール P 0 0 におけるデータ入換え用の空き領域は 2 9 0 3、フラッシュメモリ・モジュール P 0 4 におけるデータ入換え用の空き領域は 2 9 0 4 であることが分かる。

20

【 0 1 2 0 】

次に、データ領域 2 9 0 5 のデータを空き領域 2 9 0 4 に書込み、データ領域 2 9 0 6 のデータを空き領域 2 9 0 3 へ書込む。そして、図 3 0 に示すように、仮想ページアドレス空間のデータ領域 3 0 0 1 を論理ページアドレス空間のデータ領域 3 0 0 4 へ対応させ、仮想ページアドレス空間におけるデータ領域 3 0 0 2 を論理ページアドレス空間のデータ領域 3 0 0 3 へ対応させる。上記データの入換え完了後、図 3 2 に示すように、仮想ページアドレスと論理ページアドレスの変換テーブルを変更する。また、図 3 4 の空き領域管理テーブルから、フラッシュメモリ・モジュール P 0 0 におけるデータ入換え用の空き領域は 3 0 0 5、フラッシュメモリ・モジュール P 0 4 におけるデータ入換え用の空き領域は 3 0 0 6 であることが分かる。

30

【 0 1 2 1 】

本方法によれば、データ入換え用の空き領域をフラッシュメモリ・モジュール内に分散配置せずに、モジュール毎に纏めて配置して空き領域を設けているため、オフセット値の管理が不要のため、データ入換えの制御が簡単になる。

【 0 1 2 2 】

次に、フラッシュメモリ・モジュール (P D E V) に障害が発生した場合の対処方法について説明する。図 3 5 から図 3 9 により、本発明の実施の形態において、フラッシュメモリ・モジュールに障害が発生した場合のモジュール交換方法を示す。

40

【 0 1 2 3 】

図 3 5 は、モジュール交換の手順を示すフローチャートである。図 3 6 から図 3 9 は、図 3 5 のフローチャートの各ステップの構成を示す説明図である。

【 0 1 2 4 】

図 3 6 は、フラッシュメモリ・モジュールに障害発生した場合を示す説明図である。図 3 6 には、RAID グループ (V D E V) V 0 0 と、RAID グループ V 0 0 を構成するウェアレベリング・グループ (W D E V) W 0 0、ウェアレベリング・グループ W 1 0、

50

ウエアレベリング・グループW20、及びウエアレベリング・グループW30が示されている。また、ウエアレベリング・グループW00が接続されているチャンネルと同じD01上に、予備グループ(YDEV)Y00が接続されている。ウエアレベリング・グループW00内のフラッシュメモリ・モジュール(PDEV)P01に障害発生(ステップS3501)した場合を考える。

【0125】

次に、ウエアレベリング・グループ(WDEV)W00が使用できる予備グループ(YDEV)を選定する。ウエアレベリング・グループW00と同じチャンネルD01上に接続されている予備グループY00を選定する(ステップS3502)。そして予備グループY00に属するフラッシュメモリ・モジュールから、フラッシュメモリ・モジュールP01と交代させるモジュールP04を選定する(ステップS3503)。

10

【0126】

図37は、フラッシュメモリ・モジュールを入換え後の状態を示す説明図である。図37に示すように、ウエアレベリング・グループW00と予備グループY00の間で、フラッシュメモリ・モジュールP01とフラッシュメモリ・モジュールP04を入換える。障害のモジュールP01は、交換待ちの状態となる。

【0127】

次に図38は、フラッシュメモリ・モジュールを入換え後のデータ再生を示す説明図である。図38に示すように、新しくウエアレベリング・グループW00に組み入れたフラッシュメモリ・モジュールP04に対して、フラッシュメモリ・モジュールP01に書き込まれていたデータを再生し、書き込む(ステップS3504)。このとき、データ再生に使用するデータは、消去回数平準化のため、他のウエアレベリング・グループ内のフラッシュメモリ・モジュール間で分散して格納されていることに注意する。つまり、データ再生は、ウエアレベリング・グループ内で同じ仮想ページアドレスに格納されたデータについて行う。

20

【0128】

図39は、予備グループのフラッシュメモリ・モジュールを新しく交換した場合を示す説明図である。図39に示すように、交換待ちフラッシュメモリ・モジュールP01を新しいフラッシュメモリ・モジュールP06と交換し、そのP06を予備グループY00に組み入れる(ステップS3505)。以上で、フラッシュメモリ・モジュールの交換処理が終了する。

30

【0129】

なお、モジュール交換直後における図9の消去回数平準化処理の実行可否は、交換前のモジュールにおける延べ書込み容量値に基づき判定する。新規モジュールでは、その論理ページアドレス空間全領域の延べ書込み容量値がデータ再生に伴う書込み以外はゼロであり、論理ページアドレス空間内所定領域ごとの書込み頻度を知ることができない。従って交換前のモジュールの書込み容量値を用いることにより、論理ページアドレス空間における書込み頻度が分かり、消去回数平準化処理を実行できる。

【産業上の利用可能性】

【0130】

本発明は、複数のフラッシュメモリ・モジュールに渡って消去回数を平準化し、フラッシュメモリ・モジュールの長寿命化の用途に適用でき、例えば、複数のフラッシュメモリ・モジュールを使用した大容量のフラッシュメモリを用いた記憶装置、その消去回数平準化方法、及び消去回数平準化プログラムの用途に適用できる。

40

【図面の簡単な説明】

【0131】

【図1】本発明によるストレージ装置の実施の形態の構成を示すブロック図である。

【図2】チャンネルアダプタの構成を示すブロック図である。

【図3】ストレージアダプタの構成を示すブロック図である。

【図4】フラッシュメモリ・モジュールの構成を示すブロック図である。

50

- 【図5】フラッシュメモリ・モジュールのブロックの構成を示す説明図である。
- 【図6】本発明によるストレージ装置の実施の形態における論理グループの構成とアドレス変換の階層を示す説明図である。
- 【図7】本発明によるストレージ装置の実施の形態におけるRAIDグループの構成を示す説明図である。
- 【図8】フラッシュメモリ・モジュールとハードディスクドライブをストレージコントローラに接続した例を示す構成図である。
- 【図9】複数のフラッシュメモリ・モジュール間の消去回数平準化方法を示すフローチャートである。
- 【図10】消去回数平準化に伴うデータ入換え前の仮想ページアドレスと論理ページアドレスの変換テーブルを示す説明図である。 10
- 【図11】消去回数平準化に伴うデータ入換え後の仮想ページアドレスと論理ページアドレスの変換テーブルを示す説明図である。
- 【図12】ストレージコントローラにおいて管理するフラッシュメモリ・モジュール毎の消去回数管理テーブルを示す説明図である。
- 【図13】消去回数平準化に伴うデータ入換え前の仮想ページアドレスと論理ページアドレスのマッピングを示す説明図である。
- 【図14】消去回数平準化に伴うデータ入換え後の仮想ページアドレスと論理ページアドレスのマッピングを示す説明図である。
- 【図15】データ入換え操作の入換え前の初期状態を示す説明図である。 20
- 【図16】データ入換え操作の入換え途中の状態を示す説明図である。
- 【図17】データ入換え操作の入換え途中の状態を示す説明図である。
- 【図18】データ入換え操作の入換え途中の状態を示す説明図である。
- 【図19】データ入換え操作の入換え途中の状態を示す説明図である。
- 【図20】データ入換え操作の入換え途中の状態を示す説明図である。
- 【図21】データ入換え操作の入換え途中の状態を示す説明図である。
- 【図22】データ入換え操作の入換え途中の状態を示す説明図である。
- 【図23】データ入換え操作の入換え途中の状態を示す説明図である。
- 【図24】データ入換え操作の入換え後の最終状態を示す説明図である。
- 【図25】データ入換え操作のデータ入換え前後でのオフセット値の遷移を示す説明図である。 30
- 【図26】図15から図24で説明したオフセット値0の論理ページアドレス空間とオフセット値1の論理ページアドレス空間との間のデータ入換え手順を示すフローチャートである。
- 【図27】オフセット値0の論理ページアドレス空間とオフセット値0の論理ページアドレス空間との間のデータ入換え手順を示すフローチャートである。
- 【図28】オフセット値1の論理ページアドレス空間とオフセット値1の論理ページアドレス空間との間のデータ入換え手順を示すフローチャートである。
- 【図29】データ入換え前の仮想ページアドレスと論理ページアドレスのマッピングを示す説明図である。 40
- 【図30】データ入換え後の仮想ページアドレスと論理ページアドレスのマッピングを示す説明図である。
- 【図31】データ入換え前の仮想ページアドレスと論理ページアドレスの変換テーブルを示す説明図である。
- 【図32】データ入換え後の仮想ページアドレスと論理ページアドレスの変換テーブルを示す説明図である。
- 【図33】データ入換え前の空き領域管理テーブルを示す説明図である。
- 【図34】データ入換え後の空き領域管理テーブルを示す説明図である。
- 【図35】モジュール交換の手順を示すフローチャートである。
- 【図36】フラッシュメモリ・モジュールに障害発生した場合を示す説明図である。 50

【図37】フラッシュメモリ・モジュールを入換え後の状態を示す説明図である。

【図38】フラッシュメモリ・モジュールを入換え後のデータ再生を示す説明図である。

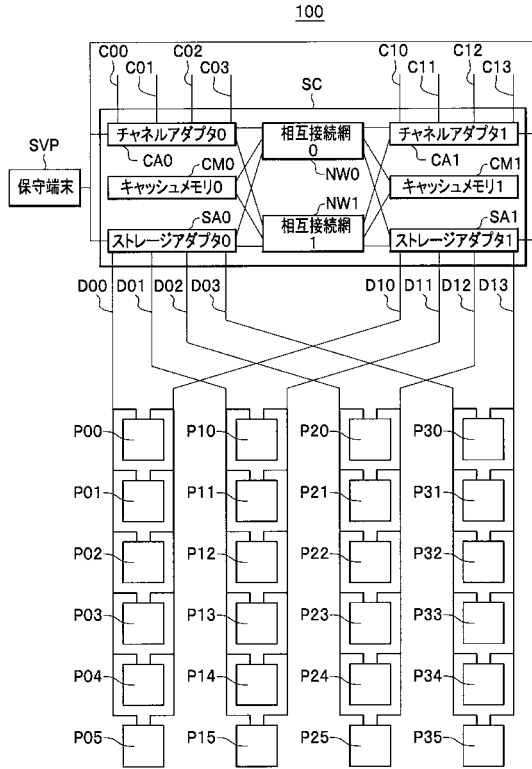
【図39】予備グループのフラッシュメモリ・モジュールを新しく交換した場合を示す説明図である。

【符号の説明】

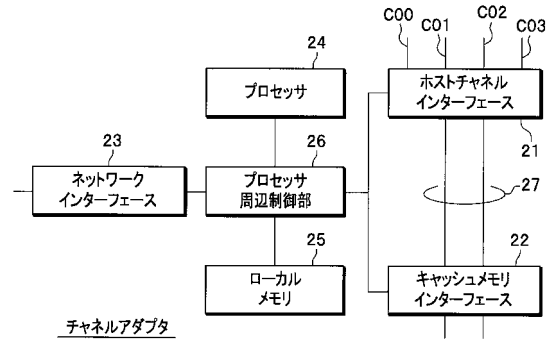
【0132】

CA0, CA1	チャンネルアダプタ	
CM0, CM1	キャッシュメモリ	
C00, C01, C02, C03, C10, C11,		
C12, C13, D00, D01, D02, D03,		10
D10, D11, D12, D13	チャンネル	
MC	メモリコントローラ	
MEM	フラッシュメモリ	
NW0, NW1	相互接続網	
P00, P01, P02, P03, P04, P05,		
P10, P11, P12, P13, P14, P15,		
P20, P21, P22, P23, P24, P25,		
P30, P31, P32, P33, P34,		
P35	フラッシュメモリ・モジュール	20
SA0, SA1	ストレージアダプタ	20
SC	ストレージコントローラ	
SV P	保守端末	
21	ホストチャンネル・インターフェース	
22	キャッシュメモリ・インターフェース	
23	ネットワーク・インターフェース	
24	プロセッサ	
25	ローカルメモリ	
26	プロセッサ周辺制御部	
27	信号線	
31	キャッシュメモリ・インターフェース	30
32	ストレージチャンネル・インターフェース	
33	ネットワーク・インターフェース	
34	プロセッサ	
35	ローカルメモリ	
36	プロセッサ周辺制御部	
37	信号線	
100	ストレージ装置	
401	プロセッサ	
402	インターフェース部	
403	データ転送部	40
404, 407	メモリ	
405	フラッシュメモリ・チップ	
406	ブロック	

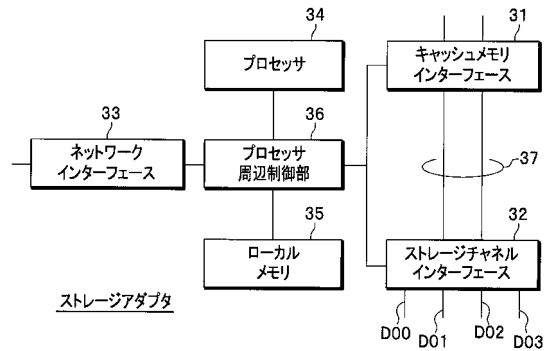
【図1】



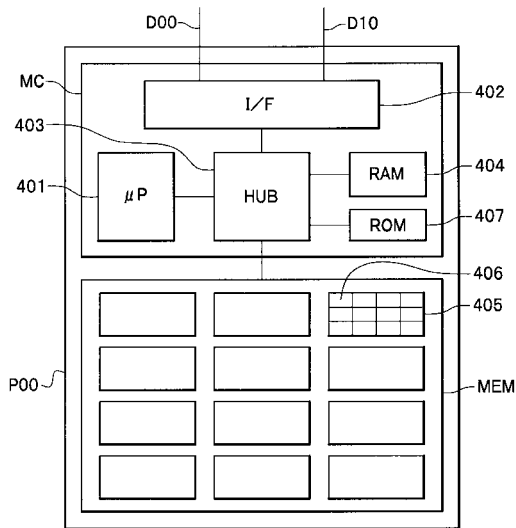
【図2】



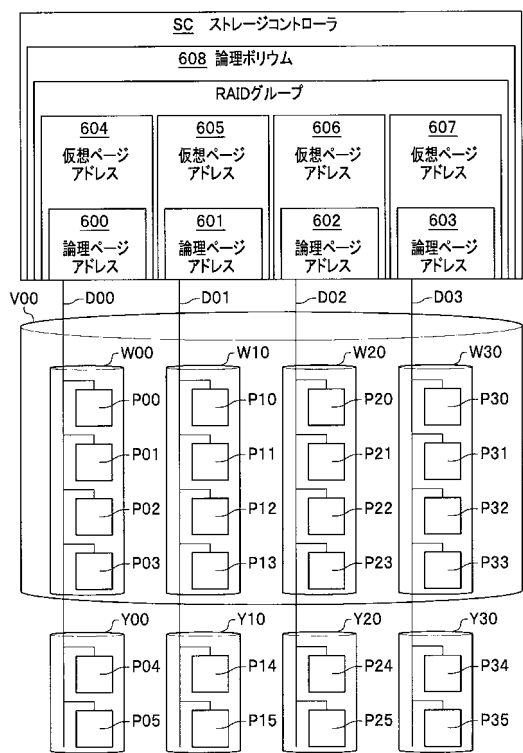
【図3】



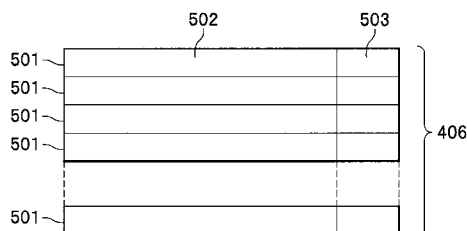
【図4】



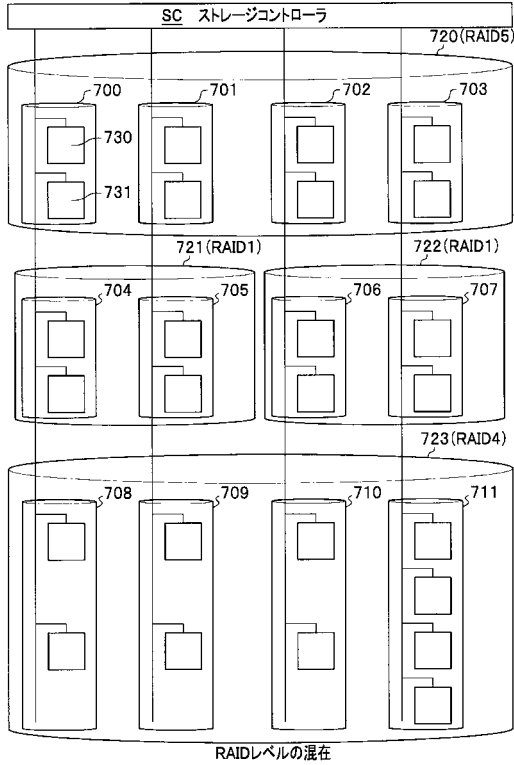
【図6】



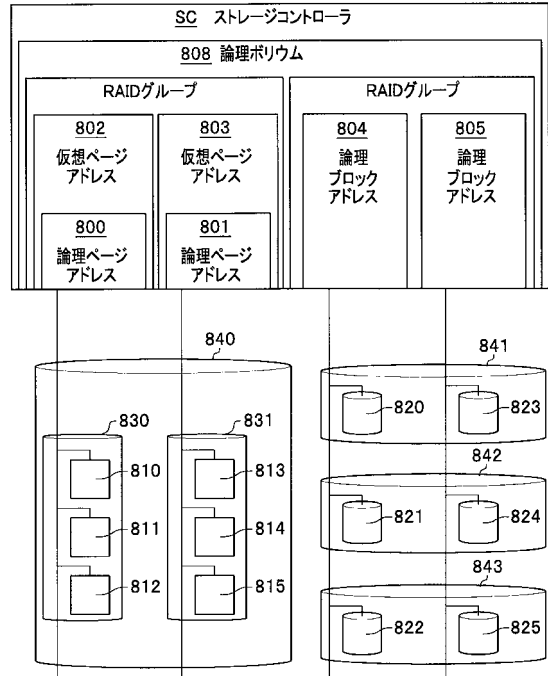
【図5】



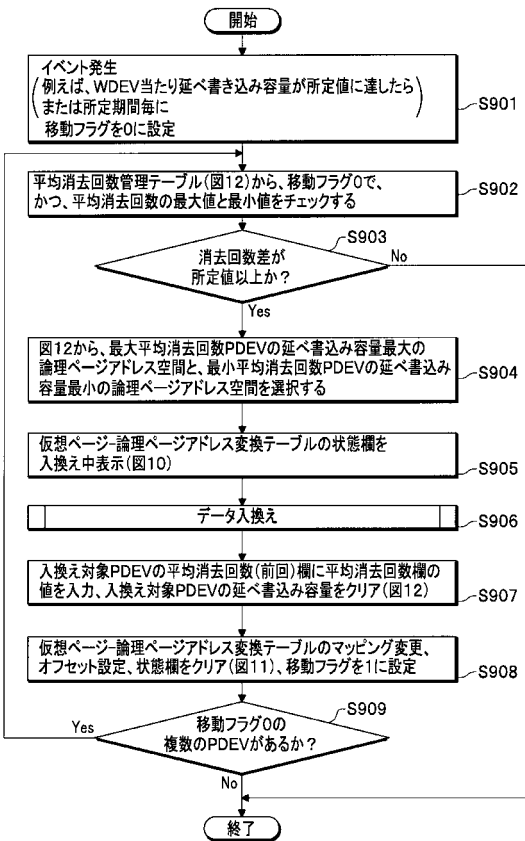
【図7】



【図8】



【図9】



【図10】

WDEV	仮想ページアドレス	データ長	PDEV	論理ページアドレス	オフセット	状態
W00	AA0		P00	AB0	0	
	AA1		P00	AB1	1	1
	AA2		P00	AB2	0	
	AA3		P00	AB3	0	
	AA4		P04	AB0	1	
	AA5		P04	AB1	1	
	AA6		P04	AB2	0	1
	AA7		P04	AB3	1	

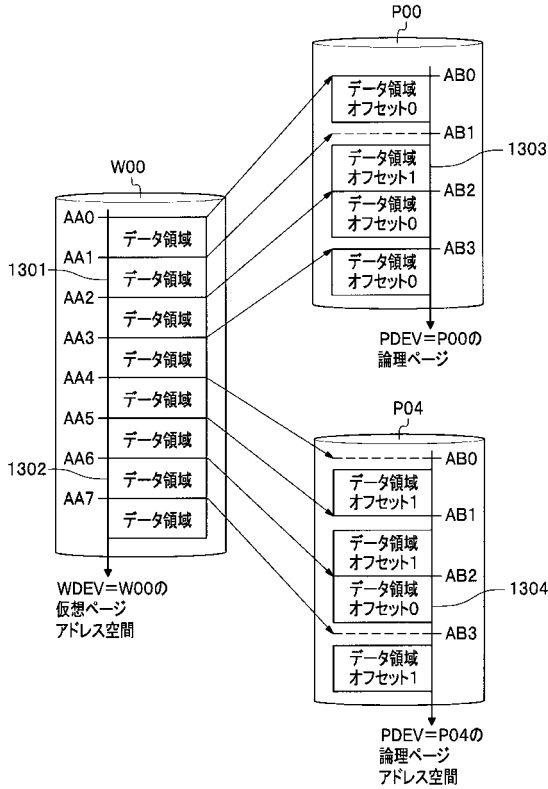
【図11】

WDEV	仮想ページアドレス	データ長	PDEV	論理ページアドレス	オフセット	状態
W00	AA0		P00	AB0	0	
	AA1		P04	AB2	0	
	AA2		P00	AB2	0	
	AA3		P00	AB3	0	
	AA4		P04	AB0	1	
	AA5		P04	AB1	1	
	AA6		P00	AB1	0	
	AA7		P04	AB3	1	

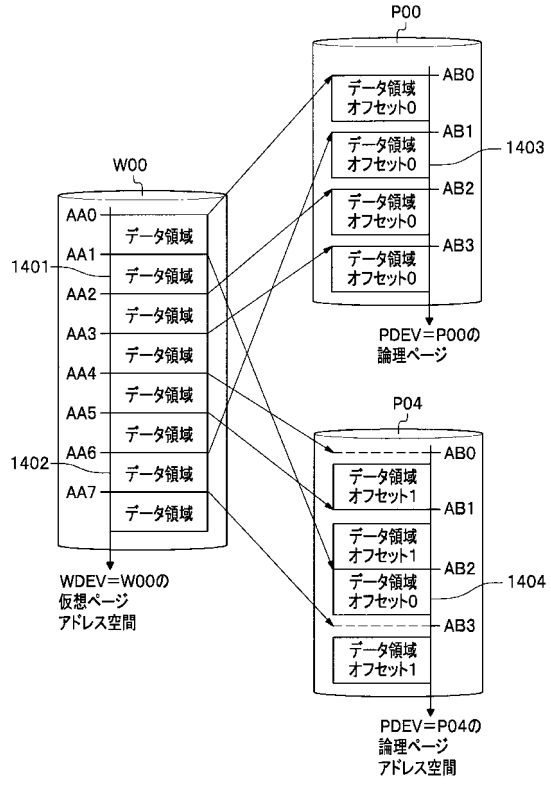
【図12】

PDEV	PDEV容量	論理ページアドレス	延べ書込み容量	平均消去回数	平均消去回数(前回)	移動フラグ
P00	s00	AB0	a000	e00=f00+(a000+a001+a002+a003)/s00	f00	0
		AB1	a001			
		AB2	a002			
P04	s04	AB0	a040	e04=f04+(a040+a041+a042+a043)/s04	f04	0
		AB1	a041			
		AB2	a042			
		AB3	a043			

【図13】



【図14】



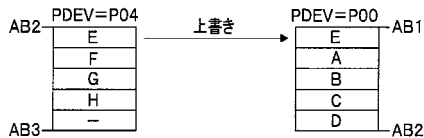
【図15】



【図20】



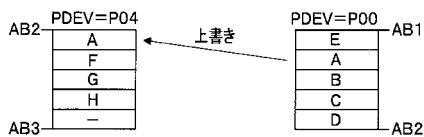
【図16】



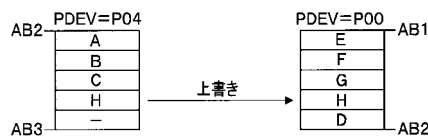
【図21】



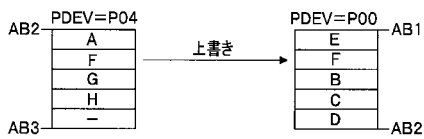
【図17】



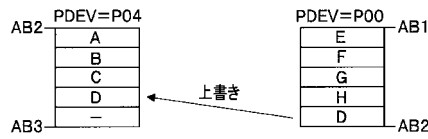
【図22】



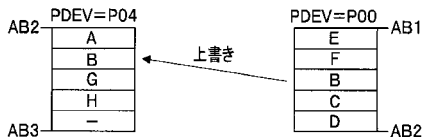
【図18】



【図23】



【図19】



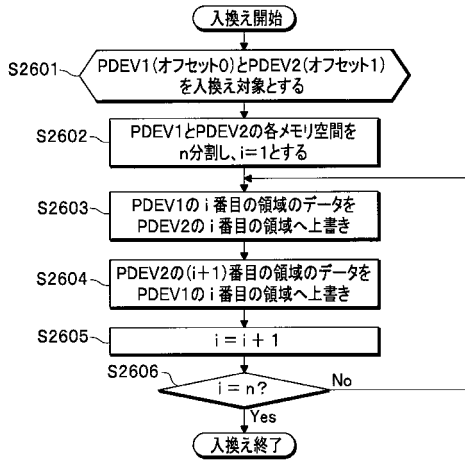
【図24】



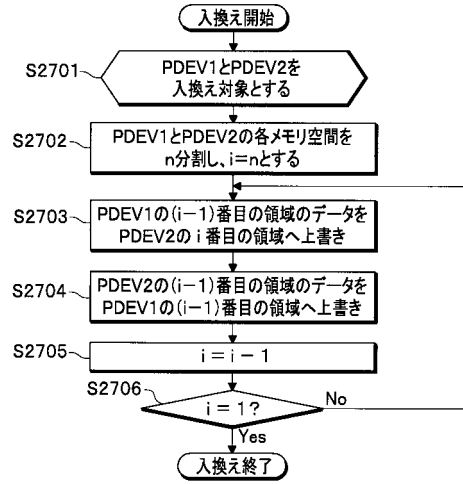
【図 25】

データ入換え前の オフセット値の組合せ		データ入換え後の オフセット値の組合せ	
0	0	0	1
0	1	0	0
1	1	1	0

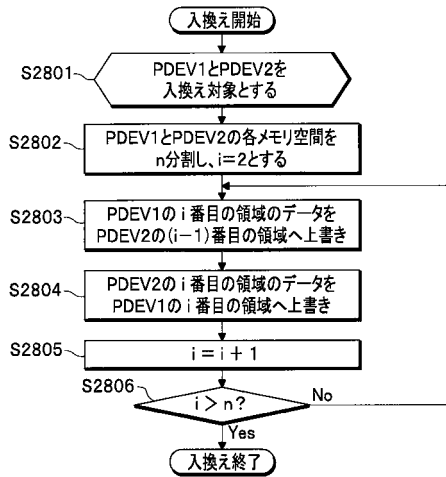
【図 26】



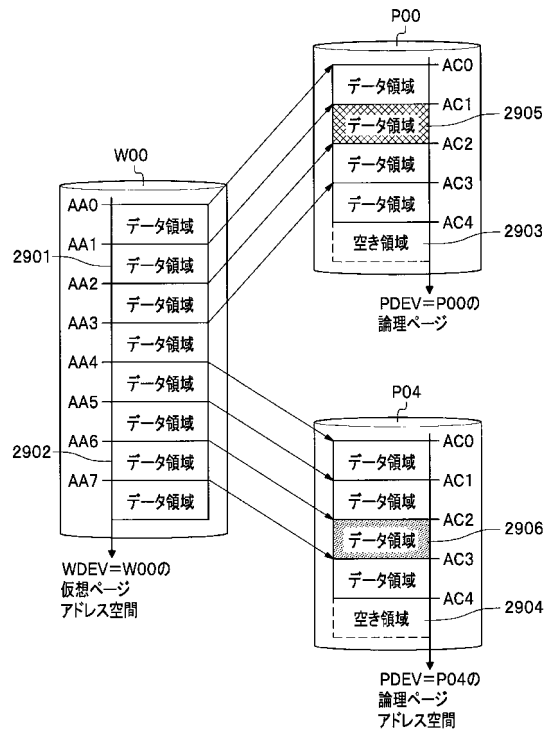
【図 27】



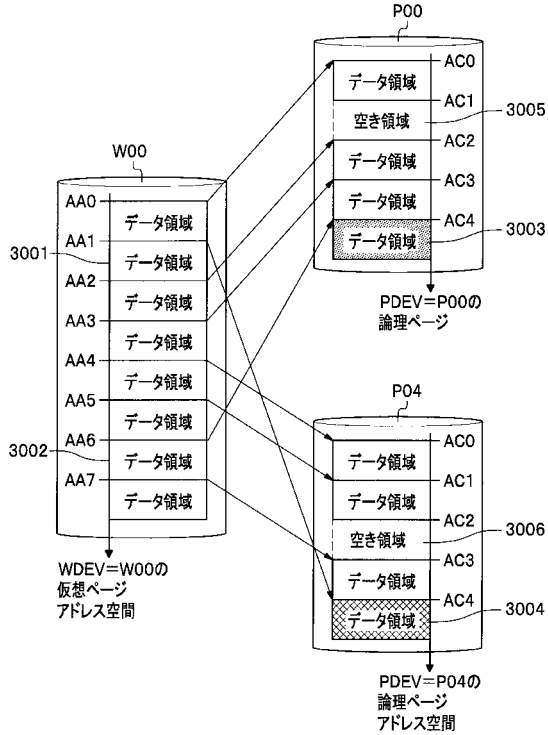
【図 28】



【図 29】



【図30】



【図31】

WDEV	仮想ページアドレス	データ長	PDEV	論理ページアドレス	状態
W00	AA0		P00	AC0	
	AA1		P00	AC1	
	AA2		P00	AC2	
	AA3		P00	AC3	
	AA4		P04	AC0	
	AA5		P04	AC1	
	AA6		P04	AC2	
AA7		P04	AC3		

【図32】

WDEV	仮想ページアドレス	データ長	PDEV	論理ページアドレス	状態
W00	AA0		P00	AC0	
	AA1		P04	AC4	
	AA2		P00	AC2	
	AA3		P00	AC3	
	AA4		P04	AC0	
	AA5		P04	AC1	
	AA6		P00	AC4	
AA7		P04	AC3		

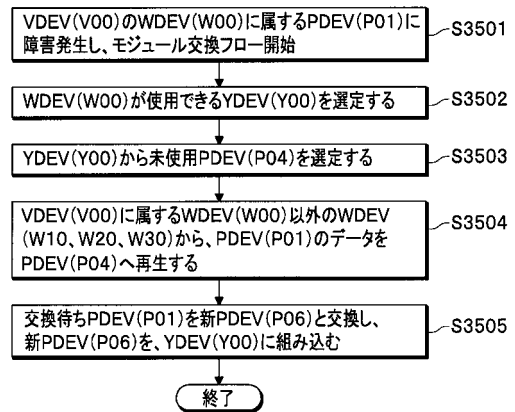
【図33】

PDEV	論理ページアドレス	データ長
P00	AC4	
P04	AC4	

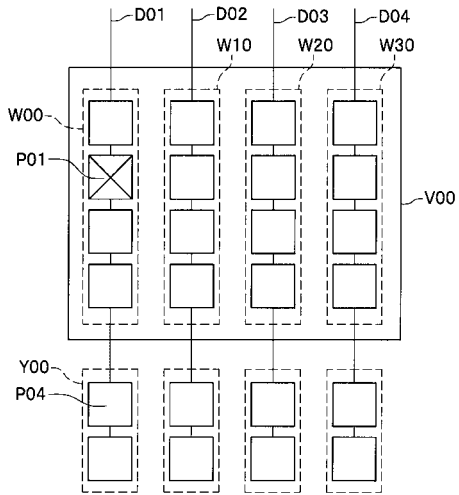
【図34】

PDEV	論理ページアドレス	データ長
P00	AC1	
P04	AC2	

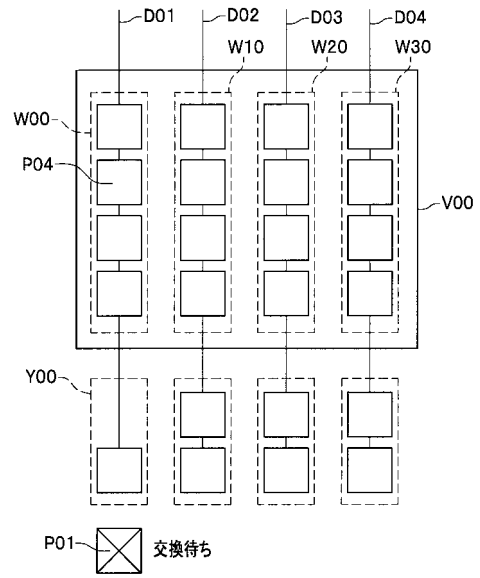
【図35】



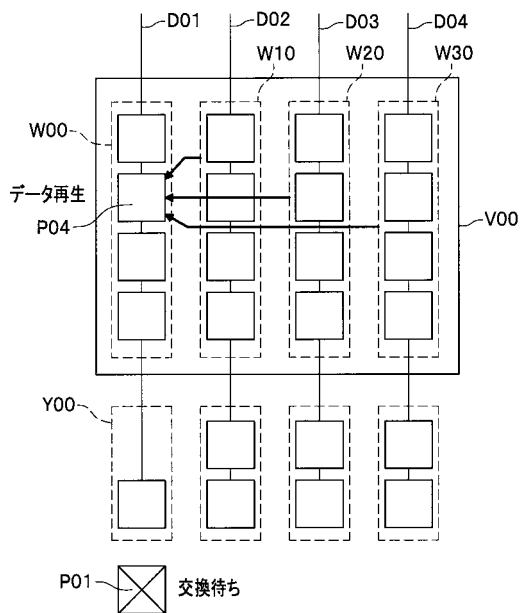
【図36】



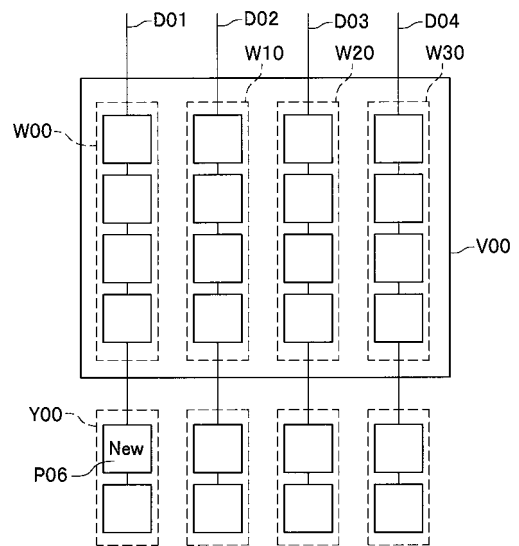
【図37】



【図38】



【図39】



フロントページの続き

(51)Int.Cl. F I
G 0 6 F 3/06 3 0 5 C

(56)参考文献 特開2000-207137(JP,A)
特開平09-218754(JP,A)
特開2004-021811(JP,A)
特開平8-16482(JP,A)

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