

(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2008-226149

(P2008-226149A)

(43) 公開日 平成20年9月25日(2008.9.25)

(51) Int. Cl.	F I	テーマコード (参考)
G06F 3/06 (2006.01)	G06F 3/06 302A	5B005
G06F 12/16 (2006.01)	G06F 12/16 310J	5B018
G06F 12/00 (2006.01)	G06F 12/16 310Q	5B065
G06F 12/08 (2006.01)	G06F 12/16 310R	5B082
	G06F 12/00 542L	

審査請求 未請求 請求項の数 17 O L (全 30 頁) 最終頁に続く

(21) 出願番号 特願2007-67142 (P2007-67142)
 (22) 出願日 平成19年3月15日 (2007. 3. 15)

(71) 出願人 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 100079108
 弁理士 稲葉 良幸
 (74) 代理人 100093861
 弁理士 大賀 真司
 (72) 発明者 水島 永雅
 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 中村 崇仁
 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

最終頁に続く

(54) 【発明の名称】 ストレージシステム及びストレージシステムのライト性能低下防止方法

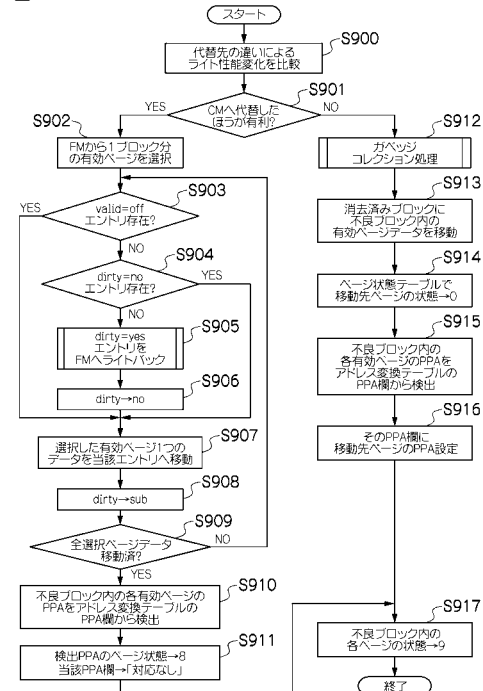
(57) 【要約】

【解決課題】 ストレージシステムのライト性能の低下を抑制する。

【解決手段】 フラッシュメモリと、キャッシュメモリと、フラッシュメモリのデータの読み出し、書き込み及び消去と、キャッシュメモリのデータの読み出し及び書き込みを制御し、フラッシュメモリ内に不良なブロックが発生したことを検出するコントローラと、データのライト処理を要求するコマンドを発行する宿主計算機とを含むストレージシステムにおいて、コントローラは、フラッシュメモリ内に不良ブロックが発生したことを検出すると、フラッシュメモリに格納された所定のデータをキャッシュメモリに移動し、その移動したデータを更新するためのコマンドを宿主計算機から受信しても、そのコマンドに基づくデータをフラッシュメモリへ書き込むことを禁止する。

【選択図】 図 1 5

図15



【特許請求の範囲】**【請求項 1】**

ページ単位でデータを書き込み、複数の前記ページから構成されるブロックを単位として前記データを消去するとともに複数の前記ブロックを有し、前記データの更新に前記ページを含む前記ブロックの消去を必要とするフラッシュメモリと、

前記フラッシュメモリに書き込むべきデータを前記フラッシュメモリよりも高速に書き込むとともに一時的に記憶するキャッシュメモリと、

前記フラッシュメモリのデータの読み出し、書き込み及び消去と、前記キャッシュメモリのデータの読み出し及び書き込みを制御し、前記フラッシュメモリ内に不良な前記ブロックが発生したことを検出するコントローラと、

前記データのライト処理を要求するコマンドを発行するホスト計算機とを備え、前記コントローラは、

前記フラッシュメモリ内に前記不良ブロックが発生したことを検出したときに、前記フラッシュメモリに格納された所定のデータを前記キャッシュメモリに移動し、その移動したデータを更新するためのコマンドを前記ホスト計算機から受信しても、そのコマンドに基づくデータを前記フラッシュメモリへ書き込むことを禁止する、

ことを特徴とするストレージシステム。

【請求項 2】

請求項 1 記載のストレージシステムであって、

前記移動したデータを更新するためのコマンドに基づいて、前記キャッシュメモリ内に移動したデータを更新する、

ことを特徴とするストレージシステム。

【請求項 3】

請求項 1 に記載のストレージシステムであって、

前記キャッシュメモリに記憶する前記データを管理するキャッシュ管理テーブルを備え、

前記キャッシュ管理テーブルは、前記データを前記フラッシュメモリへ書き込むことを禁止する禁止情報を保持する、

ことを特徴とするストレージシステム。

【請求項 4】

請求項 1 に記載のストレージシステムであって、

前記フラッシュメモリに記憶するデータの論理ページアドレスと物理ページアドレスとの対応関係を管理するアドレス変換テーブルを備え、

前記アドレス変換テーブルは、前記データの論理ページアドレスに対応する物理ページアドレスが存在しないことを表すアドレス不存在情報を保持する、

ことを特徴とする前記ストレージシステム。

【請求項 5】

請求項 1 に記載のストレージシステムであって、

前記コントローラは、

前記ホスト計算機のコマンドに基づいてデータを前記キャッシュメモリに書き込むときに前記データの論理ページアドレスと同じ論理ページアドレスの前記データが前記キャッシュメモリに存在する確率を示すヒット率を計算し、

前記不良ブロックが発生したことを検出したときに、その検出時まで計算した前記ヒット率に基づいて、前記不良ブロックの代替先として前記フラッシュメモリ内のその他のブロックを適用したときの前記ストレージシステムの第 1 のライト性能と、前記フラッシュメモリに格納された前記データの一部を前記キャッシュメモリに移動したときの前記ストレージシステムの第 2 のライト性能とを計算し、その計算結果に基づいて、前記第 2 のライト性能が前記第 1 のライト性能に勝ると判定したならば、前記フラッシュメモリに格納された前記データを前記キャッシュメモリに移動する、

ことを特徴とするストレージシステム。

10

20

30

40

50

- 【請求項 6】
請求項 1 に記載のストレージシステムであって、
前記キャッシュメモリは不揮発性のランダムアクセスメモリである、
ことを特徴とするストレージシステム。
- 【請求項 7】
請求項 6 に記載のストレージシステムであって、
前記不揮発性のランダムアクセスメモリは、相変化 R A M である、
ことを特徴とするストレージシステム。
- 【請求項 8】
請求項 1 に記載のストレージシステムであって、
前記キャッシュメモリに移動した前記データを、前記キャッシュメモリが冗長的に保持
する、
ことを特徴とするストレージシステム。 10
- 【請求項 9】
請求項 8 に記載のストレージシステムであって、
前記キャッシュメモリとは、異なるキャッシュメモリを備え、
前記冗長的な保持は、前記キャッシュメモリへ移動した前記データを、前記異なるキャ
ッシュメモリへコピーして行なう、
ことを特徴とするストレージシステム。
- 【請求項 10】
請求項 8 に記載のストレージシステムであって、
前記キャッシュメモリとは、異なる複数のキャッシュメモリを備え、
前記冗長的な保持は、前記キャッシュメモリへ移動した前記データを、前記異なる複数
のキャッシュメモリと R A I D グループを組んで行なう、
ことを特徴とするストレージシステム。 20
- 【請求項 11】
ページ単位でデータを書き込み、複数の前記ページから構成されるブロックを単位とし
て前記データを消去するとともに複数の前記ブロックを有し、前記データの更新に前記ペ
ージを含む前記ブロックの消去を必要とするフラッシュメモリと、
前記フラッシュメモリに書き込むべきデータを前記フラッシュメモリよりも高速に書き
込むとともに一時的に記憶するキャッシュメモリと、
前記フラッシュメモリのデータの読み出し、書き込み及び消去と、前記キャッシュメモ
リのデータの読み出し及び書き込みを制御し、前記フラッシュメモリ内に不良な前記ブロ
ックが発生したことを検出するコントローラと、
前記データのライト処理を要求するコマンドを発行するホスト計算機とを含むストレ
ージシステムのライト性能低下防止方法において、
前記コントローラは、
前記フラッシュメモリ内に前記不良ブロックが発生したことを検出し、
前記フラッシュメモリに格納された所定のデータを前記キャッシュメモリに移動し、
その移動したデータを更新するためのコマンドを前記ホスト計算機から受信しても、そ
の コマンドに基づくデータを前記フラッシュメモリへ書き込むことを禁止する、
ことを特徴とするストレージシステムのライト性能低下防止方法。 30 40
- 【請求項 12】
請求項 11 に記載のストレージシステムのライト性能低下防止方法であって、
前記キャッシュメモリに記憶する前記データを管理するキャッシュ管理テーブルにより
保持される禁止情報に基づいて前記データを前記フラッシュメモリへ書き込むことを禁止
する、
ことを特徴とする方法。
- 【請求項 13】
請求項 11 に記載のストレージシステムのライト性能低下防止方法であって、 50

前記コントローラは、

前記ホスト計算機のコマンドに基づいてデータを前記キャッシュメモリに書き込むときに前記データの論理ページアドレスと同じ論理ページアドレスの前記データが前記キャッシュメモリに存在する確率を示すヒット率を計算し、

前記不良ブロックが発生したことを検出したときに、その検出時まで計算した前記ヒット率に基づいて、前記不良ブロックの代替先として前記フラッシュメモリ内のその他のブロックを適用したときの前記ストレージシステムの第1のライト性能と、前記フラッシュメモリに格納された前記データの一部を前記キャッシュメモリに移動したときの前記ストレージシステムの第2のライト性能とを計算し、その計算結果に基づいて、前記第2のライト性能が前記第1のライト性能に勝ると判定したならば、前記フラッシュメモリに格納された前記データを前記キャッシュメモリに移動する、
ことを特徴とする方法。

10

【請求項14】

請求項11記載のストレージシステムのライト性能低下防止方法であって、前記キャッシュメモリは不揮発性のランダムアクセスメモリである、ことを特徴とする方法。

【請求項15】

請求項11記載のストレージシステムのライト性能低下防止方法であって、前記キャッシュメモリに移動した前記データを、前記キャッシュメモリが冗長的に保持する、
ことを特徴とする方法。

20

【請求項16】

請求項15に記載のストレージシステムのライト性能低下防止方法であって、前記冗長的な保持は、前記キャッシュメモリへ移動した前記データを、前記キャッシュメモリとは異なるキャッシュメモリへコピーして行なう、
ことを特徴とする方法。

【請求項17】

請求項15に記載のストレージシステムのライト性能低下防止方法であって、異なる複数のキャッシュメモリを備え、前記冗長的な保持は、前記キャッシュメモリへ移動した前記データを、前記キャッシュメモリとは異なる複数のキャッシュメモリとRAIDグループを組んで行なう、
ことを特徴とする方法。

30

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、電氣的に書き換え可能な不揮発性メモリを使用したストレージシステム及びストレージシステムのライト性能低下防止方法に関し、特に不揮発性メモリとしてフラッシュメモリを用い、フラッシュメモリよりも高速なランダムアクセスメモリをキャッシュメモリとして搭載したストレージシステムにおいて、不良化したフラッシュメモリブロックの代替にキャッシュメモリを使用するものに適用しても好適なものである。

40

【背景技術】

【0002】

フラッシュメモリは直接的に上書きのできない不揮発性メモリであり、書き換えを実施するためには、複数個の書き込み単位（ページ）から構成される消去単位（ブロック）を消去し、それらのページを未書き込み状態にする必要がある。

【0003】

したがって、フラッシュメモリを主たる記憶媒体とする一般的なストレージシステムでは、ホスト装置が格納データを更新したい場合、フラッシュメモリの一部にあらかじめ確保した予備領域に新データを書いて有効データとする一方で旧データを無効データとし、後に旧データを含むブロックを消去して新たな予備領域とする、といった制御を行う。な

50

お、消去するブロック内にその他の有効データが残っている場合は、消去の前にそれらを別のブロックの未書き込みのページに退避しなければならない。

【0004】

一般に、フラッシュメモリ上の予備領域は他のブロックが不良化したときの代替領域としても利用される。フラッシュメモリのブロックは10万回程度の書き換え寿命しか保証されていないため、ホスト装置がストレージシステムの格納データの更新を繰り返していると日常的に不良ブロックが発生し、その数は徐々に増加する。そして、これが予備領域を満たすほどの数になると、上記の格納データ更新制御が困難になり、そのストレージシステムは書き換え不能となる。

【0005】

特許文献1には、フラッシュメモリを記憶媒体とするデータファイルストレージシステムにおいて、キャッシュメモリを使用した以下の手続きによって、ストレージシステムの書き換え寿命を延長する方法が記載されている。(1)ホスト装置からのフラッシュメモリ向けのデータファイルを、フラッシュメモリよりかなり多くのライト・消去のサイクルに耐えうるキャッシュメモリに一時保管する。(2)ホスト装置からのライト要求に応じて、新データファイルをフラッシュメモリの代わりにキャッシュメモリに書く。(3)データファイルの識別子および各データファイルが最後にキャッシュメモリに書かれてからの時間を、タグメモリに格納する。(4)キャッシュメモリに新データファイル用の追加空間が必要になったら、タグメモリの参照によって最後の書込から最長の時間がたったデータファイルを優先してキャッシュメモリからフラッシュメモリに移動する。上記(1)~(4)の処理を行うことで、フラッシュメモリへの実際のライト回数や関連ストレスを大きく減らしている。

【特許文献1】米国特許第5936971号公報

【発明の開示】

【発明が解決しようとする課題】

【0006】

フラッシュメモリを記憶媒体とするストレージシステムにおいて、フラッシュメモリの不良ブロックの数の増加に伴って、上記の格納データ更新制御に使用する予備領域のサイズは徐々に減少する。正味の格納データ量に対して更新制御用の予備領域サイズの割合が小さいほど、ストレージシステムの格納データ更新時の作業効率(ホスト装置の書き換え単位あたりのフラッシュメモリのライト回数)は低下する。これは、上記の格納データ更新制御において、消去ブロック内に残された他の有効データを退避するときの平均作業量が多くなることが原因である。

【0007】

その結果、不良ブロックの数が増加するにつれて、ストレージシステムのライト性能は徐々に低下するという問題があった。

【0008】

本発明は、以上の点を考慮してなれたもので、ストレージシステムのライト性能の低下を抑制することができるストレージシステム及びストレージシステムのライト性能低下防止方法を提案しようとするものである。

【課題を解決するための手段】

【0009】

本発明は、ページ単位でデータを書き込み、複数の前記ページから構成されるブロックを単位として前記データを消去するとともに複数の前記ブロックを有し、前記データの更新に前記ページを含む前記ブロックの消去を必要とするフラッシュメモリと、前記フラッシュメモリに書き込むべきデータを前記フラッシュメモリよりも高速に書き込むとともに一時的に記憶するキャッシュメモリと、前記フラッシュメモリのデータの読み出し、書き込み及び消去と、前記キャッシュメモリのデータの読み出し及び書き込みを制御し、前記フラッシュメモリ内に不良な前記ブロックが発生したことを検出するコントローラと、前記データのライト処理を要求するコマンドを発行するホスト計算機とを備え、前記コント

10

20

30

40

50

ローラは、前記フラッシュメモリ内に前記不良ブロックが発生したことを検出したときに、前記フラッシュメモリに格納された所定のデータを前記キャッシュメモリに移動し、その移動したデータを更新するためのコマンドを前記ホスト計算機から受信しても、そのコマンドに基づくデータを前記フラッシュメモリへ書き込むことを禁止するストレージシステムである。

【0010】

すなわち、ページ単位でデータを書き込み、複数のページから構成されるブロックを単位としてデータを消去するとともに複数のブロックを有し、データの更新にページを含むブロックの消去を必要とするフラッシュメモリと、フラッシュメモリに書き込むべきデータをフラッシュメモリよりも高速に書き込むとともに一時的に記憶するキャッシュメモリと、フラッシュメモリのデータの読み出し、書き込み及び消去と、キャッシュメモリのデータの読み出し及び書き込みを制御し、フラッシュメモリ内に不良なブロックが発生したことを検出するコントローラと、データのライト処理を要求するコマンドを発行するホスト計算機とを含むストレージシステムのライト性能低下防止方法において、コントローラは、フラッシュメモリ内に不良ブロックが発生したことを検出し、前記キャッシュメモリに格納された所定のデータをキャッシュメモリに移動し、その移動したデータを更新するためのコマンドをホスト計算機から受信しても、そのコマンドに基づくデータをフラッシュメモリへ書き込むことを禁止することにより、フラッシュメモリの不良ブロック数の増加に伴って格納データ更新時の作業効率が低下するのを抑制することができるので、キャッシュメモリのヒット率が低い状況下において、従来よりもストレージシステムのライト性能の低下を抑制することができる。

10

20

【0011】

また、上記ストレージシステムにおけるキャッシュメモリを不揮発性のランダムアクセスメモリ、例えば、相変化RAM(Random Access Memory)とすることにより、フラッシュメモリからキャッシュメモリに代替されたデータを補助電源なしに保持することができるため、ストレージシステムの電源消費電力を削減し、データを突然の電源遮断などの障害による消失から保護するという効果を奏する。

【発明の効果】

【0012】

本発明によれば、ストレージシステムのライト性能の低下を抑制するという効果を奏する。

30

【発明を実施するための最良の形態】

【0013】

以下、本発明の実施形態について図面を参照しながら説明する。

【0014】

(1) ストレージシステムの構成

【0015】

図1は、本発明を適用したストレージシステム10のハードウェア構成を簡単に示す図である。

【0016】

ストレージシステム10は、ストレージコントローラ120及びフラッシュメモリ・モジュール(FMM)151~154、161~164、171~174、181~184を備える。ストレージコントローラ120は、チャンネルアダプタ121、122、キャッシュメモリ123、124、ストレージアダプタ125、126、共有メモリ129及び相互接続網127、128を備える。

40

【0017】

なお、ストレージコントローラ120は、チャンネルアダプタ121、122、キャッシュメモリ123、124及びストレージアダプタ125、126、共有メモリ129を有する場合を例示しているが、それらの個数は限定されるものではない

【0018】

50

相互接続網 1 2 7 及び 1 2 8 は、例えば、スイッチ等であり、ストレージコントローラ 1 2 0 を構成する各装置を相互に接続する。詳細には、相互接続網 1 2 7 及び 1 2 8 は、チャンネルアダプタ 1 2 1、キャッシュメモリ 1 2 3 及びストレージアダプタ 1 2 5、共有メモリ 1 2 9 を相互に接続する。同様に、相互接続網 1 2 7、1 2 8 は、チャンネルアダプタ 1 2 2、キャッシュメモリ 1 2 4 及びストレージアダプタ 1 2 6、共有メモリ 1 2 9 を相互に接続する。

【 0 0 1 9 】

チャンネルアダプタ 1 2 1 は、チャンネル 1 1 0、1 1 1、1 1 2、1 1 3 を介してホスト計算機 1 0 0 に接続されている。同様に、チャンネルアダプタ 1 2 2 は、チャンネル 1 1 4、1 1 5、1 1 6、1 1 7 を介してホスト計算機 1 0 0 に接続されている。ホスト計算機 1 0 0 は、パーソナルコンピュータ、ワークステーション、メインフレームコンピュータ等の計算機であり、ストレージシステム 1 0 において、ストレージコントローラ 1 2 0 にデータの読み書きを要求する。ストレージコントローラ 1 2 0 は、チャンネルアダプタ 1 2 1、1 2 2 を用いてそれらの要求を解釈し、それらの要求を満たすためにストレージアダプタ 1 2 5、1 2 6 を用いてフラッシュメモリ・モジュール 1 5 1 ~ 1 5 4、1 6 1 ~ 1 6 4、1 7 1 ~ 1 7 4、1 8 1 ~ 1 8 4 のデータを読み書きする。

10

【 0 0 2 0 】

その際、キャッシュメモリ 1 2 3、1 2 4 は、チャンネルアダプタ 1 2 1、1 2 2 又はストレージアダプタ 1 2 5、1 2 6 から受信したデータを一時的に記憶したり、必要に応じて特定の受信データを永続的に記憶したりする。キャッシュメモリ 1 2 3、1 2 4 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。共有メモリ 1 2 9 は、キャッシュメモリ 1 2 3、1 2 4 の記憶データを管理するためのテーブルを格納し、チャンネルアダプタ 1 2 1、1 2 2 又はストレージアダプタ 1 2 5、1 2 6 がそれを参照・設定することができる。共有メモリ 1 2 9 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。

20

【 0 0 2 1 】

ストレージアダプタ 1 2 5 は、フラッシュメモリ・モジュール 1 5 1 ~ 1 5 4、1 6 1 ~ 1 6 4、1 7 1 ~ 1 7 4、1 8 1 ~ 1 8 4 に接続されている。詳細には、ストレージアダプタ 1 2 5 は、チャンネル 1 4 0 を介してフラッシュメモリ・モジュール 1 5 1 ~ 1 5 4 に接続されている。また、ストレージアダプタ 1 2 5 は、チャンネル 1 4 1 を介してフラッシュメモリ・モジュール 1 6 1 ~ 1 6 4 に接続されている。また、ストレージアダプタ 1 2 5 は、チャンネル 1 4 2 を介してフラッシュメモリ・モジュール 1 7 1 ~ 1 7 4 に接続されている。また、ストレージアダプタ 1 2 5 は、チャンネル 1 4 3 を介してフラッシュメモリ・モジュール 1 8 1 ~ 1 8 4 に接続されている。

30

【 0 0 2 2 】

同様に、ストレージアダプタ 1 2 6 は、フラッシュメモリ・モジュール 1 5 1 ~ 1 5 4、1 6 1 ~ 1 6 4、1 7 1 ~ 1 7 4、1 8 1 ~ 1 8 4 に接続されている。詳細には、ストレージアダプタ 1 2 6 は、チャンネル 1 4 4 を介してフラッシュメモリ・モジュール 1 5 1 ~ 1 5 4 に接続されている。また、ストレージアダプタ 1 2 6 は、チャンネル 1 4 5 を介してフラッシュメモリ・モジュール 1 6 1 ~ 1 6 4 に接続されている。また、ストレージアダプタ 1 2 6 は、チャンネル 1 4 6 を介してフラッシュメモリ・モジュール 1 7 1 ~ 1 7 4 に接続されている。また、ストレージアダプタ 1 2 6 は、チャンネル 1 4 7 を介してフラッシュメモリ・モジュール 1 8 1 ~ 1 8 4 に接続されている。

40

【 0 0 2 3 】

チャンネルアダプタ 1 2 1、1 2 2 及びストレージアダプタ 1 2 5、1 2 6 は、保守端末 1 3 0 に接続されている。保守端末 1 3 0 は、ストレージシステム 1 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 1 2 1、1 2 2 及び / 又はストレージアダプタ 1 2 5、1 2 6 に送信する。

【 0 0 2 4 】

なお、ストレージシステム 1 0 は、ストレージアダプタ 1 2 5 及びチャンネルアダプタ 1

50

21に替わって一つのアダプタを備えていても良い。この場合、その一つのアダプタが、ストレージアダプタ125及びチャンネルアダプタ121の処理を行う。

【0025】

190～193は、RAID (Redundant Arrays of Inexpensive Disks) グループである。例えば、RAIDグループ190は、フラッシュメモリ・モジュール151、161、171、181から成る。RAIDグループ190に属するフラッシュメモリ・モジュールの一つ、例えば、フラッシュメモリ・モジュール151でエラーが発生してデータを読み出せなくなると、RAIDグループ190に属する他のフラッシュメモリ・モジュール161、171、181からデータを再生できる。

【0026】

図2はチャンネルアダプタ121のハードウェア構成を示す図である。チャンネルアダプタ121はホストチャンネル・インタフェース214、キャッシュメモリ・インタフェース215、ネットワーク・インタフェース211、プロセッサ210、ローカルメモリ213、及びプロセッサ周辺制御部212を備える。

【0027】

ホストチャンネル・インタフェース214は、チャンネルアダプタ121を、チャンネル110、111、112、113を介してホスト計算機100に接続するためのインタフェースである。ホストチャンネル・インタフェース214は、チャンネル110、111、112、113上のデータ転送プロトコルと、ストレージコントローラ120内部のデータ転送プロトコルとを相互に変換する。

【0028】

キャッシュメモリ・インタフェース215は、チャンネルアダプタ121を、相互結合網127、128に接続するためのインタフェースである。

【0029】

ネットワーク・インタフェース211は、チャンネルアダプタ121を、保守端末130に接続するためのインタフェースである。

【0030】

なお、ホストチャンネル・インタフェース214とキャッシュメモリ・インタフェース215とは、信号線216によって接続されている。

【0031】

プロセッサ210は、ローカルメモリ213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。より詳細には、プロセッサ210は、ホスト計算機100と相互結合網127、128との間のデータ転送を制御する。

【0032】

ローカルメモリ213は、プロセッサ210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213は、プロセッサ210によって参照されるテーブルを記憶する。プロセッサ210によって参照されるテーブルは、チャンネルアダプタ121の動作を制御するための設定情報を含み、管理者によって設定又は変更される。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130に入力する。保守端末130は、入力された情報を、ネットワーク・インタフェース211を介してプロセッサ210に送信する。プロセッサ210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210は、テーブルをローカルメモリ213に格納する。

【0033】

プロセッサ周辺制御部212は、ホストチャンネル・インタフェース214、キャッシュメモリ・インタフェース215、ネットワーク・インタフェース211、プロセッサ210、及びローカルメモリ213の間におけるデータ転送を制御する。プロセッサ周辺制御部212は、例えば、チップセット等である。

【0034】

なお、チャンネルアダプタ122のハードウェア構成は、チャンネルアダプタ121のハー

10

20

30

40

50

ドウェア構成と同一であるので、チャンネルアダプタ 1 2 2 のハードウェア構成についての説明を省略する。

【 0 0 3 5 】

図 3 はストレージアダプタ 1 2 5 のハードウェア構成を示す図である。ストレージアダプタ 1 2 5 は、キャッシュメモリ・インタフェース 2 2 4、ストレージチャンネルインタフェース 2 2 5、ネットワーク・インタフェース 2 2 1、プロセッサ 2 2 0、ローカルメモリ 2 2 3、及びプロセッサ周辺制御部 2 2 2 を備える。

【 0 0 3 6 】

キャッシュメモリ・インタフェース 2 2 4 は、ストレージアダプタ 1 2 5 を、相互結合網 1 2 7、1 2 8 に接続するためのインタフェースである。

【 0 0 3 7 】

ストレージチャンネルインタフェース 2 2 5 は、ストレージアダプタ 1 2 5 を、チャンネル 1 4 0、1 4 1、1 4 2、1 4 3 に接続するためのインタフェースである。ストレージチャンネルインタフェース 2 2 5 は、チャンネル 1 4 0、1 4 1、1 4 2、1 4 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 1 2 0 内部のデータ転送プロトコルとを相互に変換する。

【 0 0 3 8 】

なお、キャッシュメモリ・インタフェース 2 2 4 とストレージチャンネルインタフェース 2 2 5 とは、信号線 2 2 6 によって接続されている。

【 0 0 3 9 】

ネットワーク・インタフェース 2 2 1 は、ストレージアダプタ 1 2 5 を、保守端末 1 3 0 に接続するためのインタフェースである。

【 0 0 4 0 】

プロセッサ 2 2 0 は、ローカルメモリ 2 2 3 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。

【 0 0 4 1 】

ローカルメモリ 2 2 3 は、プロセッサ 2 2 0 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 2 3 は、プロセッサ 2 2 0 によって参照されるテーブルを記憶する。プロセッサ 2 2 0 によって参照されるテーブルは、ストレージアダプタ 1 2 5 の動作を制御するための設定情報を含み、管理者によって設定又は変更される。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 1 3 0 に入力する。保守端末 1 3 0 は、入力された情報を、ネットワーク・インタフェース 2 2 1 を介してプロセッサ 2 2 0 に送信する。プロセッサ 2 2 0 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 2 0 はテーブルをローカルメモリ 2 2 3 に格納する。

【 0 0 4 2 】

プロセッサ周辺制御部 2 2 2 は、キャッシュメモリ・インタフェース 2 2 4、ストレージチャンネルインタフェース 2 2 5、ネットワーク・インタフェース 2 2 1、プロセッサ 2 2 0 及びローカルメモリ 2 2 3 の間のデータ転送を制御する。プロセッサ周辺制御部 2 2 2 は、例えば、チップセット等である。

【 0 0 4 3 】

なお、ストレージアダプタ 1 2 6 のハードウェア構成は、ストレージアダプタ 1 2 5 のハードウェア構成と同一であるので、ストレージアダプタ 1 2 6 のハードウェア構成についての説明を省略する。

【 0 0 4 4 】

図 4 はフラッシュメモリ・モジュール 1 5 1 のハードウェア構成を示す図である。フラッシュメモリ・モジュール 1 5 1 は、メモリコントローラ 3 1 0 及びフラッシュメモリ 3 2 0 を備える。フラッシュメモリ 3 2 0 は、データを記憶する。メモリコントローラ 3 1 0 は、フラッシュメモリ 3 2 0 のデータの「読み出し」、「書き込み」、及び「消去」を制御する。

10

20

30

40

50

【 0 0 4 5 】

メモリコントローラ 3 1 0 は、プロセッサ 3 1 2、インタフェース 3 1 1、データ転送部 3 1 5、RAM 3 1 3、及び ROM 3 1 4 を備える。フラッシュメモリ 3 2 0 は、複数のフラッシュメモリ・チップ 3 2 1 を備える。

【 0 0 4 6 】

図 5 はフラッシュメモリ・チップ 3 2 1 の内部構成を示す図である。フラッシュメモリ・チップ 3 2 1 は、複数のブロック 3 3 0 を含み、それぞれのブロック 3 3 0 にデータを記憶する。ブロック 3 3 0 は、メモリコントローラ 3 1 0 がデータを消去する単位である。ブロック 3 3 0 は、複数のページ 3 4 0 を含む。ページ 3 4 0 は、メモリコントローラ 3 1 0 がデータを読み書きする単位である。フラッシュメモリ 3 2 0 では、1 ページあたり 20 μ s 程度の時間でデータを読み出し、1 ページあたり 0.2 ms 程度の時間でデータを書き込む。また、1 ブロックあたり 1.5 ms 程度の時間でデータを消去する。フラッシュメモリ 3 2 0 へのページ書き込みにかかる時間は、キャッシュメモリ 1 2 3、1 2 4 への同サイズのデータの書き込みにかかる時間よりも長い。なお、書き込みや消去はメモリセルを徐々に劣化させ、多数回（例えば、数 10 万回）の書き換えを行うとエラーが発生することがある。

10

【 0 0 4 7 】

ページ 3 4 0 は、メモリコントローラ 3 1 0 により有効ページ、無効ページ、未書込ページ、又は不良ページの何れかに分類される。有効ページは、ストレージシステム 1 0 として格納しておく必要がある有効なデータを記憶しているページ 3 4 0 である。無効ページは、ストレージシステム 1 0 として格納する必要がなくなった無効なデータ（ガベッジ）を記憶しているページ 3 4 0 である。未書込ページは、所属するブロック 3 3 0 が消去されて以来データを記憶していないページ 3 4 0 である。不良ページは、ページ 3 4 0 内の記憶素子が壊れている等の理由によって、物理的に書き換えできないページ 3 4 0 である。ページ 3 4 0 が不良ページとなる要因には次の 3 つがある。

20

【 0 0 4 8 】

1 つ目は、チップ製造段階の検査で不合格になること。2 つ目は、ページ 3 4 0 の書込においてエラーが発生すること。なお、当該不良ページは、それ以降は読み出しのみ可能となる。このページを 1 つでも含むブロック 3 3 0 は不良ブロックと呼ばれ、ブロック消去やページ書込が禁止される。3 つ目は、ブロック 3 3 0 の消去においてエラーが発生すること。なお、当該ブロック内の全ページは不良ページとなる。当該ブロックは不良ブロックと呼ばれ、ブロック消去やページ書込が禁止される。

30

【 0 0 4 9 】

インタフェース 3 1 1 は、チャンネル 1 4 0 を介してストレージコントローラ 1 2 0 内のストレージアダプタ 1 2 5 に接続されている。また、インタフェース 3 1 1 は、チャンネル 1 4 4 を介してストレージコントローラ 1 2 0 内のストレージアダプタ 1 2 6 に接続されている。インタフェース 3 1 1 は、ストレージアダプタ 1 2 5 及びストレージアダプタ 1 2 6 からの命令を受信する。ストレージアダプタ 1 2 5 及びストレージアダプタ 1 2 6 からの命令は、例えば、SCSI コマンドである。

【 0 0 5 0 】

詳細には、インタフェース 3 1 1 は、ストレージアダプタ 1 2 5 及びストレージアダプタ 1 2 6 からデータを受信する。そして、インタフェース 3 1 1 は、受信したデータを RAM 3 1 3 にバッファする。また、インタフェース 3 1 1 は、RAM 3 1 3 にバッファされているデータを、ストレージアダプタ 1 2 5 及びストレージアダプタ 1 2 6 へ送信する。

40

【 0 0 5 1 】

また、インタフェース 3 1 1 は、ハードディスクドライブとの互換性を有するインタフェース機能を有している。そのため、ストレージアダプタ 1 2 5、1 2 6 は、フラッシュメモリ・モジュール 1 5 1 ~ 1 8 4 を、ハードディスクドライブとして認識する。ストレージシステム 1 0 は、データを格納するための記録媒体として、フラッシュメモリ・モジ

50

ユーザとハードディスクドライブとを混載してもよい。

【0052】

RAM 313 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。RAM 313 は、インタフェース 311 が送受信するデータを一時的に記憶する。一方、ROM 314 は、不揮発性メモリであり、プロセッサ 312 によって実行されるプログラムを記憶する。プロセッサ 312 によって実行されるプログラムは、プロセッサ 312 が実行可能となるように、ストレージシステム 10 の起動時に ROM 314 から RAM 313 へロードされる。また、RAM 313 は、プロセッサ 312 によって参照される管理情報を記憶する。

【0053】

プロセッサ 312 によって参照される管理情報は、フラッシュメモリ 320 の論理ページアドレスと物理ページアドレスとを変換するためのアドレス変換テーブルを含む。論理ページアドレスは、フラッシュメモリ・モジュール 151 の外部から（例えば、ストレージアダプタ 125 から）、フラッシュメモリ 320 に読み書きする単位であるページを論理的に指示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 310 が、フラッシュメモリ 320 を読み書きする単位であるページを物理的にアクセスするためのアドレスである。プロセッサ 312 は、ページ対応関係の変化に応じてアドレス変換テーブルの内容を書き換える。なお、アドレス変換テーブルの具体的な例については口授する。

【0054】

さらに、この管理情報は、フラッシュメモリ 320 の物理的なページ 340 の状態を管理するためのページ状態テーブルを含む。ページ状態テーブルは、予め定義されたページ状態を符号化したものを格納する。ページ状態は 16 進数表記で以下の 4 通りを定義する。

- ・状態 = 0 ... 有効ページ
- ・状態 = 8 ... 無効ページ
- ・状態 = 9 ... 不良ページ
- ・状態 = F ... 未書込ページ

【0055】

ページ状態テーブルはブロック単位で状態を保持する。例えば、あるブロックのページ状態が “880F” のとき、その第 1・第 2 ページは無効データを含み、第 3 ページは有効データを含み、第 4 ページは未書込であることを意味する。プロセッサ 312 は、ページ状態の変化に応じてページ状態テーブルの内容を書き換える。なお、ページ情報テーブルの具体的な例については後述する。

【0056】

データ転送部 315 は、例えばスイッチであり、プロセッサ 312、インタフェース 311、RAM 313、ROM 314、及びフラッシュメモリ 320 を相互に接続し、それらの間のデータ転送を制御する。

【0057】

プロセッサ 312 は、RAM 313 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 312 は、RAM 313 に記憶されているアドレス変換テーブルを参照して、フラッシュメモリ 320 の論理ページアドレスとフラッシュメモリ 320 の物理ページアドレスとを変換し、フラッシュメモリ 320 にデータを読み書きする。また、プロセッサ 312 は、フラッシュメモリ・モジュール 151 内のブロック 330 について、ガベッジコレクション処理（ブロック再生処理）を行う。

【0058】

ガベッジコレクション処理（ブロック再生処理）は、新しいデータを書くための未書込ページが少なくなったときに、未書込ページ数を増やすために、あるブロック 330 内の無効ページを未書込ページに再生する処理である。ガベッジコレクション処理の対象となるブロック（対象ブロック）330 としては、最も多くの無効ページが存在しているブ

10

20

30

40

50

ックを選択する。未書込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が、消去はブロック単位でしかできないため、プロセッサ 3 1 2 は、有効ページのデータを未書込ページ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。このように、フラッシュメモリ・モジュール 1 5 1 に書き込まれたデータは、ストレージコントローラ 1 2 0 からの指示とは独立に、フラッシュメモリ・モジュール 1 5 1 内で移動することがある。メモリコントローラ 3 1 0 は、このデータ移動の結果を、アドレス変換テーブルやページ状態テーブルに正しく反映する。それにより、ストレージコントローラ 1 2 0 は、正しいデータをアクセスすることができる。

【 0 0 5 9 】

プロセッサ 3 1 2 は、フラッシュメモリ 3 2 0 におけるページ書き込みやガベッジコレクション処理を通じて変化する論理ページアドレスと物理ページアドレスの対応関係、ページ状態を、それぞれアドレス変換テーブル、ページ状態テーブルを用いて管理する。

10

【 0 0 6 0 】

なお、フラッシュメモリ・モジュール 1 5 1 のハードウェア構成について詳述したが、他のフラッシュメモリ・モジュール 1 5 2 ~ 1 8 4 についても、同様のハードウェア構成を有している。このため、他のフラッシュメモリモジュール 1 5 2 ~ 1 8 4 については、図示及び説明を省略する。

【 0 0 6 1 】

図 5 に示すように、各ブロックを構成する複数のページ 3 4 0 は、いずれもデータ部 3 5 0 及び冗長部 3 5 1 を含む。例えば、1 ページにつき、ページ 3 4 0 は 2 1 1 2 バイト、データ部 3 5 0 は 2 0 4 8 バイト、冗長部 3 5 1 は 6 4 バイトである。なお、本発明はそれらのページ・サイズを特に限定するものではない。

20

【 0 0 6 2 】

データ部 3 5 0 は、ユーザデータを記憶する。冗長部 3 5 1 は、ページ 3 4 0 自身に対応する論理ページアドレス、書き込み時刻、及びエラー訂正コードを記憶する。論理ページアドレスは、ストレージシステム 1 0 の起動時に R A M 3 1 3 上にアドレス変換テーブルを作成するときや、ガベッジコレクション処理のときに参照する。書き込み時刻は、ストレージシステム 1 0 の起動時に R A M 3 1 3 上にページ状態テーブルを作成するとき、ページ 3 4 0 が有効ページか無効ページであるかを知るために参照する。同じ論理ページアドレスが記録された複数のページが存在する場合、この時刻が最も遅いものが有効ページであり、それ以外は無効ページである。エラー訂正コードは、ページ 3 4 0 のエラーを検出及び訂正するための情報であり、例えば、B C H (Bose-Chaudhuri-Hocquenghem) 符号である。冗長部 3 5 1 は、通常、メモリコントローラ 3 1 0 のみがアクセス可能であり、ストレージアダプタ 1 2 5、1 2 6 からは、データ部 3 5 0 の内容のみがアクセス可能である。

30

【 0 0 6 3 】

(2) 不良ブロックの代替方法の選択とその影響

【 0 0 6 4 】

図 6 ~ 9 は、フラッシュメモリ・チップ 3 2 1 を構成する複数のブロックの中に不良ブロックが発生したときの 2 種類のブロック代替方法と、それぞれがフラッシュメモリ・モジュール 1 5 1 ~ 1 5 4、1 6 1 ~ 1 6 4、1 7 1 ~ 1 7 4 及び 1 8 1 ~ 1 8 4 へのライト性能およびキャッシュメモリ 1 2 3 のヒット確率に及ぼす影響を説明するための図である。

40

【 0 0 6 5 】

説明を簡単化するため、一般的なフラッシュメモリ・チップよりも少ないブロック数、ブロック内ページ数とする。すなわち、ユーザデータの読み書きに使用するブロック数を 7 個とし、各ブロックは 4 ページで構成する。また、ユーザデータを格納する論理ページは A ~ L の 1 2 ページ (3 ブロック分) とする。つまり、初期状態では、格納データ更新制御に使用される予備領域は 4 ブロックである。なお、図 6 ~ 8 で斜線の入ったページは無効ページ、空白のページは未書込ページである。

50

【 0 0 6 6 】

図 6 は、フラッシュメモリ・チップ 3 2 1 とキャッシュメモリ 1 2 3 を示す T 1 を示している。図 6 に示すように、T 1 は、フラッシュメモリ・チップ 3 2 1 内の 1 つのブロックが不良化し、2 つ目の不良ブロックが発生する前の状態を示している。不良化したブロックは、4 ページ全てが「b a d」と示されているブロック 3 3 0 A である。この 1 つ目の不良ブロック 3 3 0 A はフラッシュメモリ 3 2 1 内で代替される。このため、予備領域は 4 ブロックから 3 ブロックになる。また、未書込ページの残数はブロック 3 3 0 B の 4 ページになっており、ガベッジコレクションが必要な状態である。

【 0 0 6 7 】

この時点で、利用可能な 2 4 個の物理ページ (6 ブロック分) 内に 1 2 個の論理ページが配置されている。したがって、フラッシュメモリ 3 2 0 の論理ページの冗長度は $2 4 / 1 2 = 2 0 0$ パーセント (%) である。また、4 個の未書込ページを除く 2 0 個の物理ページ中に 1 2 個の有効ページが含まれるため、平均的な無効ページ含有率は $8 / 2 0 = 4 0$ % (1 ブロック辺り 1 . 6 ページ) である。したがって、ガベッジコレクション時に他ブロックへ退避すべき有効ページは平均 6 0 % (1 ブロック辺り 2 . 4 ページ) である。

10

【 0 0 6 8 】

図 7 は、フラッシュメモリ・チップ 3 2 1 とキャッシュメモリ 1 2 3 を示す T 2 を示している。図 7 に示すように、T 2 は、フラッシュメモリ・チップ 3 2 1 内の 2 つのブロック 3 3 0 A、3 3 0 C が不良化した状態を示している。2 つ目の不良化したブロックは、4 ページ全てが「b a d」と示されているブロック 3 3 0 C である。1 つ目の不良ブロック 3 3 0 A と同様に、2 つ目の不良ブロック 3 3 0 C もフラッシュメモリ 3 2 0 内のブロック 3 3 0 B で代替される。このため、予備領域は 4 ブロックから 2 ブロックになる。また、未書込ページの残数は 4 ページになっており、ガベッジコレクションが必要な状態である。

20

【 0 0 6 9 】

この時点で、利用可能な 2 0 個の物理ページ (5 ブロック分) 内に 1 2 個の論理ページが配置されている。したがって、フラッシュメモリ 3 2 0 の論理ページの冗長度は $2 0 / 1 2 = 1 6 7$ % である。また、4 個の未書込ページを除く 1 6 個の物理ページ中に 1 2 個の有効ページが含まれるため、平均的な無効ページ含有率は $4 / 1 6 = 2 5$ % (1 ブロック辺り 1 ページ) である。したがって、ガベッジコレクション時に他ブロックへ退避すべき有効ページは平均 7 5 % (1 ブロック辺り 3 ページ) である。このため、図 6 の場合に比べてガベッジコレクション時の有効ページ退避量は多くなり、このフラッシュメモリ・モジュール 1 5 1 のライト性能は低下する。なお、有効利用可能なキャッシュメモリ 1 2 3 容量は図 6 の場合と変わらないので、ストレージシステム 1 0 のキャッシュヒット確率は不変である。

30

【 0 0 7 0 】

図 8 は、フラッシュメモリ・チップ 3 2 1 とキャッシュメモリ 1 2 3 を示す T 3 を示している。図 8 に示すように、T 3 は、フラッシュメモリ・チップ 3 2 1 内の 2 つのブロック 3 3 0 A、3 3 0 C が不良化した状態を示す。1 つ目の不良ブロック 3 3 0 A はフラッシュメモリ 3 2 0 内で代替されているが、2 つ目の不良ブロック 3 3 0 C はキャッシュメモリ 1 2 3 内の一部の領域 4 0 0 で代替されている。キャッシュメモリ 1 2 3 内の一部の領域 4 0 0 内で代替される E ~ H の 4 個の論理ページデータはキャッシュメモリ 1 2 3 上にもみ存在し、フラッシュメモリ 3 2 0 から消失している。その結果、ユーザデータを格納する論理ページは A ~ D、I ~ L の 8 ページ (2 ブロック分) となる。このため、予備領域は 4 ブロックから 3 ブロックになる。また、未書込ページの残数は 4 ページになっており、ガベッジコレクションが必要な状態である。

40

【 0 0 7 1 】

この時点で、利用可能な 2 0 個の物理ページ (5 ブロック分) 内に 8 個の論理ページが配置されている。したがって、フラッシュメモリ 3 2 0 の論理ページの冗長度は $2 0 / 8 = 2 5 0$ % である。また、4 個の未書込ページを除く 1 6 個の物理ページ中に 8 個の有効

50

ページが含まれるため、平均的な無効ページ含有率は $8 / 16 = 50\%$ (1ブロック辺り2ページ)である。したがって、ガベッジコレクション時に他ブロックへ退避すべき有効ページは平均 50% (1ブロック辺り2ページ)である。このため、図6の場合に比べてガベッジコレクション時の有効ページ退避量は少なくなり、このフラッシュメモリ・モジュール151のライト性能は向上する。なお、有効利用可能なキャッシュメモリ123の容量は、4個の論理ページデータ削減されるので、図6のT1の場合に比べてストレージシステム10のキャッシュヒット確率は低下する。

【0072】

図9は、以上の説明をまとめたテーブルT4である。図9に示すように、テーブルT4は、項目名T41、T1~T3における結果T41~T43が対応付けられている。項目名T41は、フラッシュメモリの論理ページの冗長度、ガベッジコレクション時の無効ページ含有率(平均)、ガベッジコレクション時の有効ページ移動量(平均)、フラッシュメモリのライト性能及びキャッシュメモリのヒット確立という項目名が配置されている。フラッシュメモリの論理ページの冗長度は、T41(200%)、T42(167%)、T43(250%)となっている。ガベッジコレクション時の無効ページ含有率(平均)は、T41(40%)、T42(25%)、T43(50%)となっている。ガベッジコレクション時の有効ページ移動量は、T41(60%)、T42(75%)、T43(50%)となっている。また、フラッシュメモリ320のライト性能は、T42の場合は低下し、T43の場合は向上する。キャッシュメモリ123のヒット確立は、T42の場合は不変であるが、T43の場合は低下する。

10

20

【0073】

図9に示すように、不良ブロック330Cの代替先としてフラッシュメモリ320を選択すると、ストレージシステム10のライト性能は低下するという欠点がある。一方、不良ブロック330Cの代替先としてキャッシュメモリ123を選択すると、ストレージシステム10のライト性能は向上するという利点があるが、キャッシュヒット確率は低下するという欠点がある。キャッシュヒット確率が低下するということは、ミスヒットによるキャッシュメモリ123からフラッシュメモリ320へのライトバック頻度が多くなることである。上記のように、フラッシュメモリ320へのページ書き込みにかかる時間は、キャッシュメモリ123への同サイズのデータの書き込みにかかる時間よりも長い。したがって、キャッシュヒット確率が低下するということは、結果的にストレージシステム10のライト性能を低下させる。

30

【0074】

不良ブロック330Cの代替先としてキャッシュメモリ123を選択することの利点・欠点はトレードオフの関係にあり、ホスト計算機100のフラッシュメモリ320へのアクセスパターンにより、ストレージシステム10のライト性能が向上したり低下したりする。例えば、ランダムライトでは、キャッシュメモリ123の有用性が低いので、フラッシュメモリ・モジュール151のライト性能を向上させたほうが良く、不良ブロック(例えば、330C)の代替先としてキャッシュメモリ123を選択したほうが有利である。また、例えば、部分集中ライトでは、キャッシュメモリ123の有用性が高いので、有効利用可能なキャッシュメモリ123の容量を減らさずにヒット確率を維持したほうがよく、不良ブロック(例えば、330C)の代替先としてフラッシュメモリ320を選択したほうが有利である。

40

【0075】

本発明を適用したストレージシステム10は、不良ブロックの代替先としてキャッシュメモリ123を選択したときのストレージシステム10のライト性能をホスト計算機100のアクセスパターンに基づいて推定し、ストレージシステム10のライト性能が向上すると判断される場合には不良ブロックの代替先としてキャッシュメモリ123を選択する。逆に、ストレージシステム10のライト性能が低下すると判断される場合には不良ブロックの代替先として従来技術のようにフラッシュメモリ320を選択する。この処理については、後述する。

50

【 0 0 7 6 】

(3) キャッシュメモリおよびフラッシュメモリの管理手段

【 0 0 7 7 】

図 1 0、図 1 1 を参照しながら、ストレージシステム 1 0 におけるキャッシュメモリ 1 2 3 およびフラッシュメモリ・チップ 3 2 1 の管理手段を説明する。図 1 0 及び図 1 1 は、上記管理手段を説明するために、キャッシュメモリ 1 2 3、共有メモリ 1 2 9、RAM 3 1 3 及びフラッシュメモリ・チップ 3 2 1 に記憶される内容を説明するための図である。なお、キャッシュメモリ 1 2 4 や、その他のフラッシュメモリ・チップについても同様であるため、キャッシュメモリ 1 2 4 や、その他のフラッシュメモリ・チップについては説明を省略する。

10

【 0 0 7 8 】

説明を簡単化するため、図 6 ~ 図 9 の場合と同様にユーザデータの読み書きに使用するブロック数を 7 個とし、各ブロックは 4 ページで構成する。また、ユーザデータを格納する論理ページは $Ax \sim Lx$ の 1 2 ページ (3 ブロック分) とする。ここで、 x は 0 以上の整数であり、各論理ページの更新回数を示す。例えば、E 2 とは 2 回更新された論理ページ E のデータを示す。図 6 ~ 図 9 の場合と同様に、斜線の入ったページは無効ページ、空白のページは未書込ページである。また、キャッシュメモリ 1 2 3 の管理方式はフラッシュメモリ 3 2 0 のページ・サイズを単位とする 2 way - セットアソシアティブ方式とする。この方式は、1 つの way が 4 つの $index 0 \sim 3$ を持ち、論理ページ A、E、I はインデックス ($index$) = 0 の 2 エントリを利用し、論理ページ B、F、J はインデックス = 1 の 2 エントリを利用し、論理ページ C、G、K はインデックス = 2 の 2 エントリを利用し、論理ページ D、H、L はインデックス = 3 の 2 エントリを利用するものである。なお、本発明はキャッシュメモリ 1 2 3 の管理方式を特に限定するものではない。

20

【 0 0 7 9 】

図 1 0 及び図 1 1 に示すように、共有メモリ 1 2 9 には、キャッシュメモリ 1 2 3 を管理するためのキャッシュ管理テーブル 5 0 0 を含む。フラッシュメモリ・モジュール 1 5 1 内の RAM 3 1 3 には、フラッシュメモリ・チップ 3 2 1 を管理するためのアドレス変換テーブル 5 1 0、ページ状態テーブル 5 2 0 を含む。なお、図 1 0 は、第 5 ブロックが不良化し、そのブロックがフラッシュメモリ・チップ 3 2 1 内で代替されている状態を示している。また、図 1 1 は、図 1 0 の状態から第 4 ブロックの第 2 ページにデータを書き込む際にライトエラーが発生し、第 4 ブロックが 2 つ目の不良ブロックになり、キャッシュメモリ 1 2 3 の一部をその代替先としたときの状態を示している。

30

【 0 0 8 0 】

図 1 0 及び図 1 1 に示すように、キャッシュ管理テーブル 5 0 0、は、4 つのインデックス 5 0 0 1 とそれぞれ 2 つのウェイ (way) 5 0 0 2、合計 8 エントリの使用状態を管理する。各エントリについてバリッド (valid) フラグ 5 0 0 3、キー (key) レジスタ 5 0 0 4、ダーティ (dirty) レジスタ 5 0 0 5 を持つ。バリッドフラグ 5 0 0 3 はそのエントリが使用中か否かを記録し、

- ・ valid = on : 使用中
- ・ valid = off : 空き

40

と定義される。キーレジスタ 5 0 0 4 はそのエントリに格納しているキャッシュデータの論理ページを記録する。ダーティ (dirty) レジスタ 5 0 0 5 はそのエントリのキャッシュデータがフラッシュメモリ 3 2 0 上のデータよりも新しいかどうか、あるいはそのエントリがフラッシュメモリ・チップ 3 2 1 の不良ブロック代替用に使用されているかを記録し、

- ・ dirty = yes : フラッシュメモリ 3 2 0 のものより新しい更新データ保持
- ・ dirty = no : フラッシュメモリ 3 2 0 のものと同じオリジナルデータ保持
- ・ dirty = sub : 代替先として使用状態

と定義される。dirty = yes のエントリは、別の論理ページのために使用する前にフラッシュメモリ 3 2 0 へライトバックして同期をとる必要がある。dirty = sub

50

のエントリは、フラッシュメモリ 320 上にライトバック先がないので別の論理ページのために使用しないように管理する必要がある。

【0081】

図 10 及び図 11 に示すアドレス変換テーブル 510 は、上記のように論理ページアドレスと物理ページアドレスとを変換するためにそれらの対応状態を管理する。アドレス変換テーブル 510 は、LPA5101 と PPA5102 が対応して構成される。このアドレス変換テーブル 510 において、LPA は論理ページアドレス、PPA は物理ページアドレスを意味する。なお、ブロック番号 X 内のページ番号 Y の物理ページアドレスは XY で表している。

【0082】

図 10 及び図 11 に示すページ状態テーブル 520 は、前述のようにブロック単位で各ページ状態を管理する。図 10 及び図 11 において、ブロック (block) 5201 はブロック番号、ステータス (status) 5202 はそのブロックのページ状態を表す。

【0083】

以下、図 10 の状態から図 11 の状態に遷移するときのキャッシュ管理テーブル 500 アドレス変換テーブル 510、ページ状態テーブル 520 の変化について説明する。

【0084】

キャッシュメモリ 123 でのみ管理する対象論理ページとして、例えば、4 個の論理ページ E ~ H を選択する。

【0085】

キャッシュ管理テーブル 500 を参照し、論理ページ E の最新データ E2 はすでにキャッシュメモリ 123 上にあるので、データ転送せずにダートイレジスタを sub (代替状態) に設定する。

【0086】

キャッシュ管理テーブル 500 を参照し、論理ページ F の最新データ F0 はキャッシュメモリ 123 上にないので、インデックス = 1 の空きエントリにデータ転送して、バリッドフラグを「on」(使用中)、ダートイレジスタを「sub」(代替状態) に設定する。

【0087】

キャッシュ管理テーブル 500 を参照し、論理ページ G の最新データ G0 はキャッシュメモリ 123 上にないので、K0 をキャッシュしているインデックス = 2 のダートイ = 「no」(オリジナルデータ保持) のエントリにデータ転送して、バリッドフラグを「on」(使用中)、ダートイレジスタを「sub」(代替状態) に設定する。

【0088】

キャッシュ管理テーブル 500 を参照し、論理ページ H の最新データ H0 はキャッシュメモリ 123 上にないので、インデックス = 3 の空きエントリにデータ転送して、バリッドフラグを「on」(使用中)、ダートイレジスタを「sub」(代替状態) に設定する。

【0089】

図 11 に示すキャッシュ管理テーブル 600 は、以上の設定結果を示している。

【0090】

次に、アドレス変換テーブル 510 において、論理ページ E ~ H に対応する物理ページアドレス 5102 を調べ、ページ状態テーブル 520 において、その物理ページのページ状態を「8」(無効) に設定する。すなわち、ページ状態テーブル 520 の第 2 ブロックの欄を「8888」に設定する。そして、アドレス変換テーブル 510 において、論理ページ E ~ H に対応する物理ページアドレスをいずれもクリアして「対応なし」の状態に設定する。なお、「対応なし」は、アドレス変換テーブル 510 の実装上は、「対応なし」を意味する特別な値 (例えば、「FFFFFFFF (16 進数)») を定義し、その値を RAM 313 に書き込むことで「対応なし」状態を表現する。この RAM 313 に書き込

10

20

30

40

50

む値は実際に使用する論理ページアドレス範囲外の値であれば、どのような値でもよい。

【0091】

次に、不良が発生した第4ブロックに含まれる有効ページのデータを退避する。つまり、第0ページのB1は未書込ページ（例えば、第6ブロックの第0ページ）にコピーする。当該コピー先ページが有効ページであることを示すため、ページ状態テーブル520の第6ブロックの欄を「0FFF」に設定する。また、アドレス変換テーブル510のLPA5101「B」に対応するPPA5102を「60」に設定する。

【0092】

第1ページの「E2」はキャッシュメモリ123でのみ管理するので、すでに退避が完了しており、特に何もしない。

【0093】

最後に、第4ブロックが不良ブロックとなったことを示すため、ページ状態テーブル520の第4ブロックのステータス5202の欄を「9999」に設定する。

【0094】

図11に示すアドレス変換テーブル510、ページ状態テーブル520は、以上の設定結果を示している。

【0095】

(4) ストレージコントローラ120およびメモリコントローラ310の処理手順

【0096】

図12～図16を参照しながら、ストレージシステム10におけるストレージコントローラ120およびメモリコントローラ310の詳細な処理手順を説明する。以下、上記図10、図11に示したキャッシュメモリ123およびフラッシュメモリ320の管理手段に基づいて説明する。

【0097】

図12は、ホスト計算機100からのデータライト要求について、ストレージコントローラ120とメモリコントローラ310が行う処理を示すフローチャートである。以下、その手順を説明する。

【0098】

まず、ストレージコントローラ120はライト要求としてライト対象論理ページアドレスと新しいデータを受信し(S701)、当該論理ページアドレスのデータを含むエントリがキャッシュメモリ123内に存在するかを、キャッシュ管理テーブル500で調べる(S702)。

【0099】

その結果が真(存在する)ならば(S702: YES)、ストレージコントローラ120はキャッシュメモリ123内の当該データを、受信した新しいデータで更新する(S703)。そして、そのエントリのダーティレジスタ5005が「sub(代替状態)」であるかを調べる(S704)。それが真ならば(S704: YES)、そのまま処理を終える。それが偽ならば(S704: NO)、当該ダーティレジスタ5005を「yes」(更新データ保持)に設定し(S705)、処理を終える。

【0100】

一方、ステップS702の結果が偽(存在しない)ならば(S702: NO)、ストレージコントローラ120はキャッシュ管理テーブル500でバリッドフラグ5003が「off」のエントリ(空きエントリ)が存在するか調べる(S706)。空きエントリが存在すれば(S706: YES)、受信した新しいデータを当該エントリに書き込む(S710)。このとき、ストレージコントローラ120はそのエントリのバリッドフラグ5003を「on:(使用中)」に設定する。

【0101】

一方、空きエントリが存在しない場合は(S706: NO)、ストレージコントローラ120はダーティレジスタ5005が「no」のエントリ(オリジナルデータ保持エントリ)が存在するか調べる(S707)。オリジナルデータ保持エントリが存在すれば(S

10

20

30

40

50

707: YES)、そのエントリは上書きしても問題ないので、空きエントリ存在の場合と同様にステップS710へ遷移する。

【0102】

オリジナルデータ保持エントリが存在しなければ(S707:NO)、ストレージコントローラ120はダーティレジスタ5005が「yes」のエントリ(更新データ保持エントリ)の格納データをフラッシュメモリ320へライトバックし(S708、図14参照)、そのダーティレジスタ5005を「no」に設定する(S709)。その結果、そのエントリは上書きしても問題ないので、空きエントリ存在の場合と同様にステップS710へ遷移する。なお、ライトバック処理S708の対象となるエントリの選択には、例えば、LRU(Least Recently Used)アルゴリズムを適用する。

10

【0103】

ステップS710の後、ストレージコントローラ120はそのダーティレジスタ5005を「yes」(更新データ保持)に設定して(S705)、処理を終える。

【0104】

図13は、ホスト計算機100からのデータリード要求について、ストレージコントローラ120とメモリコントローラ310が行う処理を示すフローチャートである。以下、その手順を説明する。

【0105】

まず、ストレージコントローラ120はリード要求としてリード対象論理アドレスを受信し(S711)、当該論理ページアドレスのデータを含むエントリがキャッシュメモリ123内に存在するかを、キャッシュ管理テーブル500で調べる(S712)。

20

【0106】

その結果が偽(存在しない)ならば(S712:NO)、ストレージコントローラ120はキャッシュ管理テーブル500でバリッドフラグ5003が「off」のエントリ(空きエントリ)が存在するか調べる(S715)。空きエントリが存在すれば(S715:YES)、受信した論理アドレスからフラッシュメモリ・モジュール151~154、161~164、171~174及び181~184のいずれ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特定し、メモリコントローラ310がアドレス変換テーブル510によりそれに対応する物理ページアドレス5102を調べ、そのページの格納データを読み出す(S719)。そして、メモリコントローラ310は、そのデータをキャッシュメモリ123の当該エントリに転送する(S720)。このとき、ストレージコントローラ120はそのエントリのバリッドフラグ5003を「on」(使用中)、ダーティレジスタ5005を「no」(オリジナルデータ保持)に設定する。その後、ステップS712に戻る。

30

【0107】

一方、空きエントリが存在しない場合は(S715:NO)、ストレージコントローラ120はダーティレジスタ5005が「no」のエントリ(オリジナルデータ保持エントリ)が存在するか調べる(S716)。オリジナルデータ保持エントリが存在すれば(S716:YES)、そのエントリは上書きしても問題ないので、空きエントリ存在の場合と同様にステップS719へ遷移する。オリジナルデータ保持エントリが存在しなければ(S716:NO)、ストレージコントローラ120はダーティレジスタ5005が「yes」のエントリ(更新データ保持エントリ)の格納データをフラッシュメモリ320へライトバックし(S717、図14参照)、そのダーティレジスタ5005を「no」に設定する(S718)。その結果、そのエントリは上書きしても問題ないので、空きエントリ存在の場合と同様にステップS719へ遷移する。なお、ライトバック処理S717の対象となるエントリの選択には、例えば、LRUアルゴリズムを適用する。

40

【0108】

一方、ステップS712の結果が真(存在する)ならば(S712:YES)、ストレージコントローラ120はキャッシュメモリ123内の当該データを読み出し(S713)、ホスト計算機100へそれを送信して(S714)、処理を終える。

50

【 0 1 0 9 】

図 1 4 は、図 1 2 におけるメモリコントローラ 3 1 0 によるキャッシュデータのライトバック処理 S 7 0 8、図 1 3 におけるメモリコントローラ 3 1 0 によるキャッシュデータのライトバック処理 S 7 1 7、後述する図 1 5 におけるメモリコントローラ 3 1 0 によるキャッシュデータのライトバック処理 S 9 0 5 の詳細な処理手順を示すフローチャートである。以下、その手順を説明する。

【 0 1 1 0 】

まず、メモリコントローラ 3 1 0 はキャッシュメモリ 1 2 3 からライトバック対象の論理ページアドレスとそのデータを取得する (S 8 0 1)。そして、ページ状態テーブル 5 2 0 を用いて状態 = 「 F 」 (未書込) のページを選択し、そのページに取得したデータを書き込む (S 8 0 2)。そして、メモリコントローラ 3 1 0 は、書き込みにエラーが発生したかを判定する (S 8 0 3)。書き込みがエラーならば (S 8 0 3 : Y E S)、メモリコントローラ 3 1 0 は、不良ブロック代替処理 (S 8 0 4) を実施して、ステップ S 8 0 2 に戻る。不良ブロック代替処理の詳細は図 1 5 を参照して後述する。

10

【 0 1 1 1 】

一方、書き込みが成功ならば (S 8 0 3 : N O)、メモリコントローラ 3 1 0 は、アドレス変換テーブル 5 1 0 を用いてライトバック対象の論理ページアドレス 5 1 0 1 に対応する物理ページアドレス (旧アドレス) 5 1 0 2 を調べ、ページ状態テーブル 5 2 0 においてその旧アドレスが示すページの状態を 「 8 」 (無効) に設定する (S 8 0 5)。

【 0 1 1 2 】

そして、アドレス変換テーブル 5 1 0 において、メモリコントローラ 3 1 0 は、ライトバック対象の論理ページアドレス 5 0 1 2 に対応する物理ページアドレス 5 1 0 2 に、ステップ S 8 0 2 で書き込んだ物理ページのアドレス (新アドレス) を設定する (S 8 0 6)。また、メモリコントローラ 3 1 0 は、ページ状態テーブル 5 2 0 においてその新アドレスが示すページの状態を 「 0 」 (有効) に設定する (S 8 0 7)。最後に、メモリコントローラ 3 1 0 は、ガベッジコレクション処理 (S 8 0 8) を実施し、次のライトバックのために十分な未書込ページを確保して、ライトバック処理を終える。

20

【 0 1 1 3 】

図 1 5 は、図 1 4 におけるストレージコントローラ 1 2 0 とメモリコントローラ 3 1 0 による不良ブロックの代替処理 S 8 0 4 の詳細な処理手順を示すフローチャートである。以下、その手順を説明する。

30

【 0 1 1 4 】

まず、メモリコントローラ 3 1 0 は不良ブロックの代替先としてキャッシュメモリ 1 2 3 を選択した場合と、フラッシュメモリ 3 2 0 を選択した場合とで、ストレージシステム 1 0 のライト性能の変化をそれぞれ推定し、比較する (S 9 0 0)。そのライト性能の具体的方法の一例は後述する。

【 0 1 1 5 】

メモリコントローラ 3 1 0 は、比較の結果、キャッシュメモリ 1 2 3 へ代替したほうが有利か否かを判定し (S 9 0 1)、キャッシュメモリ 1 2 3 へ代替したほうが有利 (ライト性能が高い) であれば (S 9 0 1 : Y E S)、キャッシュメモリ 1 2 3 への代替処理 (S 9 0 2 ~ S 9 1 1 及び S 9 1 7) を実施し、さもなくば (S 9 0 1 : N O)、フラッシュメモリ 3 2 0 への代替処理 (S 9 1 2 ~ S 9 1 7) を実施する。

40

【 0 1 1 6 】

キャッシュメモリ 1 2 3 への代替処理としては、まず、メモリコントローラ 3 1 0 はページ状態テーブル 5 2 0 を用いて不良ブロックの発生したフラッシュメモリ・チップ 3 2 1 から 1 ブロック分の有効ページ (図 6 ~ 図 1 1 の例では 4 ページ) を選択する (S 9 0 2)。ストレージコントローラ 1 2 0 はキャッシュ管理テーブル 5 0 0 でバリッドフラグ 5 0 0 3 が 「 o f f 」 のエントリ (空きエントリ) が存在するか調べる (S 9 0 3)。

【 0 1 1 7 】

空きエントリが存在すれば (S 9 0 3 : Y E S)、メモリコントローラ 3 1 0 はステッ

50

ブ S 9 0 2 で選択した有効ページの 1 つの格納データを当該エントリに移動する (S 9 0 7)。このとき、ストレージコントローラ 1 2 0 はそのエントリのバリッドフラグ 5 0 0 3 を「 on 」 (使用中) に設定する。そして、ストレージコントローラ 1 2 0 はそのダーティレジスタ 5 0 0 5 を「 sub 」 (代替状態) に設定する (S 9 0 8)。

【 0 1 1 8 】

一方、空きエントリが存在しない場合は (S 9 0 3 : NO)、ストレージコントローラ 1 2 0 はダーティレジスタ 5 0 0 5 が「 no 」のエントリ (オリジナルデータ保持エントリ) が存在するか調べる (S 9 0 4)。オリジナルデータ保持エントリが存在すれば (S 9 0 4 : YES)、そのエントリは上書きしても問題ないので、空きエントリ存在の場合と同様にステップ S 9 0 7 へ遷移する。

10

【 0 1 1 9 】

オリジナルデータ保持エントリが存在しなければ (S 9 0 4 : NO)、ストレージコントローラ 1 2 0 はダーティレジスタ 5 0 0 5 が「 yes 」のエントリ (更新データ保持エントリ) の格納データをフラッシュメモリ 3 2 0 へライトバックし (S 9 0 5、図 1 4 参照)、そのダーティレジスタ 5 0 0 5 を「 no 」に設定する (S 9 0 6)。その結果、そのエントリは上書きしても問題ないので、空きエントリ存在の場合と同様にステップ S 9 0 7 へ遷移する。なお、ライトバック処理 S 9 0 5 の対象となるエントリの選択には、例えば、LRU アルゴリズムを適用する。ストレージコントローラ 1 2 0 は、ステップ S 9 0 2 で選択した全ての有効ページの格納データの移動が終わるまでステップ S 9 0 3 ~ S 9 0 8 を繰り返す (S 9 0 9)。

20

【 0 1 2 0 】

全ての移動が完了すれば (S 9 0 9 : YES)、メモリコントローラ 3 1 0 は不良ブロック内の各有効ページの物理ページアドレスをアドレス変換テーブル 5 1 0 の P P A 5 1 0 2 から検出する (S 9 1 0)。そして、メモリコントローラ 3 1 0 は、検出した物理ページアドレスのページ状態をページ状態テーブル 5 2 0 で「 8 」 (無効) に設定し、アドレス変換テーブル 5 1 0 の当該 P P A 5 1 0 2 をクリア (「 対応なし 」 を設定) する (S 9 1 1)。最後に、メモリコントローラ 3 1 0 は、不良ブロックを構成する各ページのページ状態を「 9 」 (不良) に設定し (S 9 1 7)、不良ブロック代替処理を終える。

【 0 1 2 1 】

一方、フラッシュメモリへの代替処理としては、まず、メモリコントローラ 3 1 0 がガベッジコレクション処理を実施し (S 9 1 2)、全てのページが未書込状態の消去済みブロックを確保する。

30

【 0 1 2 2 】

そして、不良ブロック内の有効ページをページ状態テーブル 5 2 0 から検出し、その格納データを当該消去済みブロックに移動する (S 9 1 3)。次に、メモリコントローラ 3 1 0 は、ページ状態テーブル 5 2 0 でそれぞれの移動先ページの状態を「 0 」 (有効) に設定する (S 9 1 4)。また、メモリコントローラ 3 1 0 は、不良ブロック内の各有効ページの物理ページアドレスをアドレス変換テーブル 5 1 0 の P P A 5 1 0 2 から検出する (S 9 1 5)。そして、その P P A 5 1 0 2 の欄にステップ S 9 1 3 でのそれぞれの移動先ページの物理ページアドレスを設定する (S 9 1 6)。最後に、メモリコントローラ 3 1 0 は、不良ブロックを構成する各ページのページ状態を「 9 」 (不良) に設定し (S 9 1 7)、不良ブロック代替処理を終える。

40

【 0 1 2 3 】

なお、ステップ S 9 0 2 で選択する有効ページとしては、書き換え頻度が高い論理ページのデータを格納しているものを選択することが好ましい。なぜなら、書き換えがほとんどない論理ページのデータをキャッシュメモリ 1 2 3 上で永続的に保持することは、性能面で非効率であるからである。

【 0 1 2 4 】

図 1 6 は、図 1 4 におけるメモリコントローラ 3 1 0 によるガベッジコレクション処理 S 8 0 8、図 1 5 におけるメモリコントローラ 3 1 0 によるガベッジコレクション処理 S

50

9 1 2 の詳細な処理手順を示すフローチャートである。以下、その手順を説明する。

【0 1 2 5】

まず、メモリコントローラ 3 1 0 はページ状態テーブル 5 2 0 を見て、ページ状態値 = 「F」(未書込)のページの残数が所定の値以下になったかを調べる(S 1 0 0 1)。なお、ステップ S 8 0 8 の時の所定数は「1ブロック分の総ページ数(図 6 ~ 図 1 1 の例では 4 ページ)」、ステップ S 9 1 2 の時の所定数は「2ブロック分の総ページ数(図 6 ~ 図 1 1 の例では 8 ページ)」である。ステップ S 1 0 0 1 の結果が偽(所定値より大きい)ならば(S 1 0 0 1 : NO)、メモリコントローラ 3 1 0 は、何もせずにそのまま処理を終える。

【0 1 2 6】

一方、その結果が真(所定値以下)ならば(S 1 0 0 1 : YES)、ページ状態値 = 「8」(無効)であるページを最も多く含むブロックを 1 つ選択し、その中にあるページ状態値 = 「0」(有効)である有効ページを全て検出する(S 1 0 0 2)。なお、検出した有効ページ数を N とする。

【0 1 2 7】

そしてメモリコントローラ 3 1 0 は、選択したブロック以外のブロック上にあるページ状態値 = F の未書込ページを選択し、ステップ S 1 0 0 2 で検出した有効ページの 1 つに格納されたデータを、それぞれ未書込ページの 1 つにコピー(退避)する(S 1 0 0 3)。次に、メモリコントローラ 3 1 0 は、ページ状態テーブル 5 2 0 において、ステップ S 1 0 0 3 でのそれぞれのコピー先ページのページ状態値を「0」に設定する(S 1 0 0 4)。そして、メモリコントローラ 3 1 0 は、ステップ S 1 0 0 2 で検出した有効ページのページ状態値を「8」(無効)に設定する(S 1 0 0 5)。そして、アドレス変換テーブル 5 1 0 上に設定された全ての物理ページアドレス(PPA)の中からステップ S 1 0 0 2 で検出した有効ページのアドレスを検索することにより、そのアドレスが設定された PPA 5 1 0 2 を検出し、そこにステップ S 1 0 0 3 でのコピー先ページの物理ページアドレスをコピーする(S 1 0 0 6)。

【0 1 2 8】

そして、メモリコントローラ 3 1 0 は、検出した全有効ページがコピー済みかどうかを判定する(S 1 0 0 7)。そして、検出した全ての有効ページがコピー済みでなければ(S 1 0 0 7 : NO)、ステップ S 1 0 0 3 に戻る。すなわち、検出した N 個の有効ページの全てについてステップ S 1 0 0 3 ~ S 1 0 0 6 を繰り返す。これにより、ステップ S 1 0 0 2 で選択したブロック内の全ページが無効化され、そこにあった保存すべきページデータの退避も完了する。

【0 1 2 9】

一方、検出した全ての有効ページがコピー済みであれば(S 1 0 0 7 : YES)、メモリコントローラ 3 1 0 は、ステップ S 1 0 0 2 で選択したブロックを消去する(S 1 0 0 8)。そして、メモリコントローラ 3 1 0 は、消去時にエラーが発生したか否かを判定する(S 1 0 0 9)。この消去時にエラーが発生したならば(S 1 0 0 9 : YES)、ページ状態テーブル 5 2 0 において、選択したブロックの各ページの状態を「9」(不良)に設定して(S 1 0 1 0)、ステップ S 1 0 0 2 に戻る。消去時にエラーが発生せず、消去が成功したならば(S 1 0 0 9 : NO)、ページ状態テーブル 5 2 0 において、そのブロック内の全ページのページ状態値を「F」(未書込)に設定する(S 1 0 1 1)。以上をステップ S 1 0 0 1 の結果が偽となるまで繰り返し、処理を終える。

【0 1 3 0】

(5) 不良ブロック代替先によるライト性能の変化の推定・比較方法

【0 1 3 1】

図 1 5 のステップ S 9 0 0 におけるライト性能の変化の推定・比較方法の一例を以下で説明する。

【0 1 3 2】

不良ブロックが発生した時点で、そのフラッシュメモリ・チップ 3 2 1 内で管理してい

10

20

30

40

50

る論理ブロック数をM、フラッシュメモリ・チップ321で利用可能な(つまり、不良でない)物理ブロック数をNと定義する。このとき、フラッシュメモリ・チップ321の論理ページ冗長度は N/M となる。

【0133】

また、キャッシュメモリ123の総容量をC、フラッシュメモリ・チップ321のブロックサイズをB、キャッシュメモリ123で代替されている不良ブロック数をSと定義する。このとき、キャッシュとして有効利用可能なキャッシュメモリ123の容量は「 $C - B * S$ 」(*は乗算を示す)となる。

【0134】

さらに、図12においてP730で囲まれた処理(S703~S705)に必要な時間をキャッシュメモリライト処理時間 T_c 、P740で囲まれた処理(S706~S710及びS705)に必要な時間をフラッシュメモリライト処理時間 T_f と定義する。

【0135】

キャッシュメモリライト処理時間 T_c は一定である。しかし、フラッシュメモリライト処理時間 T_f は、ライトバック処理S708に要する時間がフラッシュメモリ・チップ321の状態に依存して変化するため一定ではない。図14に示したように、ライトバック処理S708はガベッジコレクション処理S808を含む。図16に示したように、ガベッジコレクション処理S808では2ブロック間で有効ページの置換を行う。図6~図8を用いて説明したように、このときの置換ページ数の平均は論理ページ冗長度 N/M が小さいほど多くなる。ゆえに、論理ページ冗長度 N/M が小さいほどライトバック処理S708に要する時間は長くなる。したがって、フラッシュメモリライト処理時間 T_f は論理ページ冗長度 N/M の関数 $T_f(N/M)$ となる。

【0136】

次に、過去一定時間 t 内に、ホスト計算機100がライトアクセスしたサイズをA、その間のキャッシュメモリ123のヒット率をRとする。このときヒット率Rは、

- ・ Aが $C - B * S$ 以上のとき、 $R = (C - B * S) / A$
- ・ Aが $C - B * S$ 未満のとき、 $R = 1$

と表される。そして、ホスト計算機100のライトアクセス処理時間の期待値 T_w は、

- ・ Aが $C - B * S$ 以上のとき、 $T_w = R * T_c + (1 - R) T_f$
- ・ Aが $C - B * S$ 未満のとき、 $T_w = T_c$

と表される。以上より、ライト処理時間期待値 T_w はS、M、Nの関数 $T_w(S, M, N)$ となる。メモリコントローラ310は、ストレージコントローラ120の管理する「R」、「S」、「A」といった情報に基づいて、この関数 T_w が今回発生した不良ブロックの代替先によってどのように変化するかを評価する。

【0137】

キャッシュメモリ123で代替した場合、「S」は1増加、「M」は1減少、「N」は1減少するため、ライト処理時間期待値は $T_w(S+1, M-1, N-1)$ に変化すると推定する。フラッシュメモリ320で代替した場合、「S」は不変、「M」は不変、「N」は1減少するため、ライト処理時間期待値は $T_w(S, M, N-1)$ に変化すると推定する。

【0138】

以上より、 $T_w(S+1, M-1, N-1)$ が $T_w(S, M, N-1)$ 未満ならば、フラッシュメモリ320を代替先とした方が有利である。また、 $T_w(S+1, M-1, N-1)$ が $T_w(S, M, N-1)$ 以上ならば、キャッシュメモリ123を代替先とした方が有利である。このような比較により、不良ブロックの代替先を選択する。

【0139】

例えば、ランダムライトのようなキャッシュメモリ123のヒット率Rが非常に小さい状況下では、フラッシュメモリライト処理時間 T_f を短くすることによりライトアクセス処理時間の期待値 T_w を短くするのが有効であるため、キャッシュメモリ123が代替先として選択されやすくなる。また、例えば、部分集中ライトのようなキャッシュメモリ1

10

20

30

40

50

23のヒット率Rが非常に大きい状況下では、キャッシュとして有効利用可能なキャッシュメモリ123の容量を維持することによりヒット率Rを維持してライトアクセス処理時間の期待値 T_w の増加を抑えるのが有効であるため、フラッシュメモリ320が代替先として選択されやすくなる。

【0140】

なお、上記評価において、キャッシュメモリ123の総容量Cに対するキャッシュとして有効利用可能なキャッシュメモリ123の容量の割合に下限を与えたり、キャッシュメモリ123で代替される不良ブロック数Sに上限を与えたりすることにより、キャッシュメモリ123が代替先として選択されるのを制限し、キャッシュとして有効利用可能なキャッシュメモリ123の容量が少なくなり過ぎてストレージシステム10の性能が不安定になることを防止してもよい。ここでの不安定とは、ホスト計算機100からのアクセスパターンの変動に対して性能が大きく変動し、保証性能が確保できない状況になるという意味である。

10

【0141】

なお、上記評価において、メモリコントローラ310がストレージコントローラ120の管理する「R」、「S」、「A」といった情報を取得するため、ストレージコントローラ120はそれらを送信するコマンドを発行し、フラッシュメモリ・モジュール151~154、161~164、171~174及び181~184はそのコマンドを解釈するように構成されている。

20

【0142】

(6)代替データの移設によるライト性能最適化

【0143】

上記の実施形態では、不良ブロック発生時にその代替先の最適な場所を評価したが、通常の動作状況でもそれまでの不良ブロック代替先の最適な配分を評価して、その配分を調整してもよい。

【0144】

メモリコントローラ310はストレージコントローラ120の管理する「R」、「S」、「A」といった情報を常時監視しながら、関数 $T_w(S, M, N)$ の最適化を実施する。

30

【0145】

すなわち、例えばキャッシュメモリ123上の代替データを1ブロック分フラッシュメモリ320へ移設した場合、「S」は1減少、「M」は1増加、「N」は不変のため、ライト処理時間期待値は $T_w(S-1, M+1, N)$ に変化すると推定する。逆に、フラッシュメモリ320上の代替データを1ブロック分キャッシュメモリ123へ移設した場合、「S」は1増加、「M」は1減少、「N」は不変のため、ライト処理時間期待値は $T_w(S+1, M-1, N)$ に変化すると推定する。

【0146】

以上より、ライト処理時間期待値 $T_w(S-1, M+1, N)$ がライト処理時間期待値 $T_w(S, M, N)$ 未満ならば、キャッシュメモリ123上の代替データを1ブロック分フラッシュメモリ320へ移設する。また、ライト処理時間期待値 $T_w(S+1, M-1, N)$ がライト処理時間期待値 $T_w(S, M, N)$ 未満ならば、フラッシュメモリ320上の代替データを1ブロック分キャッシュメモリ123へ移設する。

40

【0147】

以上のように、フラッシュメモリ320を主たる記憶媒体とし、キャッシュメモリ123を搭載し、本発明を適用したストレージシステム10は、フラッシュメモリ320の不良ブロック数の増加に伴って格納データ更新時の作業効率が低下するのを抑制することができるため、フラッシュメモリおよびキャッシュメモリを搭載した従来技術のようなストレージシステムよりもライト性能の低下を抑制するという効果を奏する。

【0148】

以上(4)~(6)においては、キャッシュメモリ123で不良ブロックを代替する場

50

合を示したが、キャッシュメモリ 1 2 4 で代替するように構成しても良い。また、(5) や (6) の評価では、キャッシュメモリ 1 2 3 と 1 2 4 を合わせた総容量 C や両者の平均ヒット率 R 等を計算に用いても良い。

【 0 1 4 9 】

(7) キャッシュメモリデータの信頼性向上

【 0 1 5 0 】

上記の実施形態でのキャッシュメモリ 1 2 3、1 2 4 を不揮発性メモリによって構成してもよい。そうすれば、フラッシュメモリ・チップ 3 2 1 からキャッシュメモリ 1 2 3 に代替されたデータを無電源で永続的に保持することができる。

【 0 1 5 1 】

不揮発性メモリの一例は、相変化 R A M である。相変化 R A M は、ダイナミック型ランダムアクセスメモリのキャパシタ部分を、光ディスクなどに使用されている G S T (Ga-Sb-Te) と呼ばれる相変化材料に置き換えた構造のものを利用することが好ましい。この相変化 R A M はダイナミック型ランダムアクセスメモリとほぼ同程度のライト性能を持つため、このような実施形態も、上で述べたような本発明の効果を享受する。

【 0 1 5 2 】

不揮発性メモリのもう一例は、フラッシュメモリである。このフラッシュメモリはフラッシュメモリ・チップ 3 2 1 と同じ種類のものでよい。ただし、このフラッシュメモリはデータ更新作業効率を高めるために多くの予備領域を含み、フラッシュメモリ・チップ 3 2 1 よりも高速にデータを書き換えられる性能を持つ。したがって、このような実施形態も、上で述べたような本発明の効果を享受する。なお、このときの共有メモリ 1 2 9 には、R A M 3 1 3 にあるようなアドレス変換テーブル 5 1 0 やページ状態テーブル 5 2 0 をさらに作成し、キャッシュメモリ 1 2 3、1 2 4 内のデータ格納位置を管理する。

【 0 1 5 3 】

不揮発性のキャッシュメモリを搭載したストレージシステム 1 0 は、フラッシュメモリからキャッシュメモリに代替されたデータを補助電源なしに保持することができるため、ストレージシステムの電源消費電力を削減することができ、そのデータを突然の電源遮断などの障害による揮発的消失から保護するという効果を奏する。

【 0 1 5 4 】

また、上記の実施形態でのキャッシュメモリ 1 2 3 に含まれるフラッシュメモリ代替データ (ダーティ (d i r t y) = サブ (s u b) のエントリの格納データ) を、もう一つのキャッシュメモリ 1 2 4 に複写し、2 重に保持・管理してもよい。そうすれば、キャッシュメモリ 1 2 3 が故障しても、フラッシュメモリ代替データが消失しないので、ストレージシステム 1 0 の信頼性が保全される。

【 0 1 5 5 】

なお、フラッシュメモリ代替データの保持方式では上の 2 重化 (ミラーリング) 以外の方式を適用してもよい。例えば、ストレージシステム 1 0 にキャッシュメモリをさらに追加し、複数のキャッシュメモリで R A I D グループを組み、R A I D 5 などの冗長化方式を適用してフラッシュメモリ代替データを保持してもよい。

【 0 1 5 6 】

(8) 他の実施形態

【 0 1 5 7 】

上述の実施形態では本発明を、ページ単位でデータを書き込み、複数のページから構成されるブロックを単位としてデータを消去するとともに複数のブロックを有し、データの更新にページを含むブロックの消去を必要とするフラッシュメモリ 3 2 0 と、フラッシュメモリ 3 2 0 に書き込むべきデータをフラッシュメモリ 3 2 0 よりも高速に書き込むとともに一時的に記憶するキャッシュメモリ 1 2 3 と、フラッシュメモリ 3 2 0 のデータの読み出し、書き込み及び消去と、キャッシュメモリ 1 2 3 のデータの読み出し及び書き込みを制御し、フラッシュメモリ 3 2 0 内に不良なブロックが発生したことを検出するコントローラ 1 2 0 と、データのライト処理を要求するコマンドを発行するホスト計算機 1 0 0

10

20

30

40

50

とを含むストレージシステム 10 において、チャンネルアダプタ 121, 122 及びストレージアダプタ 125, 126 を含むストレージコントローラ 120 は、フラッシュメモリ 320 内に不良ブロックが発生したことを検出したときに、フラッシュメモリ 320 に格納された「E」、「F」、「G」及び「H」というデータをキャッシュメモリ 123 に移動し、その移動したデータを更新するためのコマンドをホスト計算機 100 から受信しても、そのコマンドに基づくデータをフラッシュメモリ 320 へ書き込むことを禁止する場合について適用した場合について述べたが、本発明はこれに限られず、この他種々の構成のストレージシステムに広く適用することができる。

【0158】

また、ストレージシステム 10 は、キャッシュメモリ 123 に記憶するデータを管理するキャッシュ管理テーブル 500 を共有メモリ 129 内に備え、キャッシュ管理テーブル 500 は、データをフラッシュメモリ 320 へ書き込むことを禁止する禁止情報として「sub」（ダーティレジスタ 5005 の項目）を保持する場合について説明したが、禁止情報の保持形式は、これに限られるものではない。

10

【0159】

さらに、ストレージシステム 10 は、フラッシュメモリ 320 に記憶するデータの LPA 5101 と PPA 5102 との対応関係を管理するアドレス変換テーブル 510 を RAM 313 内に備え、アドレス変換テーブル 510 は、データの LPA 5101 の「E」、「F」、「G」及び「H」に対応する PPA 5102 が存在しないことを表すアドレス不在情報として「 」を保持する場合について説明したが、アドレス不在情報の保持形式は、これに限られるものではない。

20

【産業上の利用可能性】

【0160】

本発明は、種々のストレージシステムに広く適用することができる。

【図面の簡単な説明】

【0161】

【図1】本発明に係わるストレージシステムの構成を示す図である。

【図2】本発明に係わるストレージシステムを構成するチャンネルアダプタの内部構成を示す図である。

【図3】本発明に係わるストレージシステムを構成するストレージアダプタの内部構成を示す図である。

30

【図4】本発明に係わるストレージシステムを構成するフラッシュメモリ・モジュールの内部構成を示す図である。

【図5】本発明に係わるストレージシステムを構成するフラッシュメモリ・モジュールに搭載されたフラッシュメモリ・チップの構造を示す図である。

【図6】本発明に係わるフラッシュメモリで発生した不良ブロックの代替方法の違いがもたらす影響を説明するための図である。

【図7】本発明に係わるフラッシュメモリで発生した不良ブロックの代替方法の違いがもたらす影響を説明するための図である。

【図8】本発明に係わるフラッシュメモリで発生した不良ブロックの代替方法の違いがもたらす影響を説明するための図である。

40

【図9】本発明に係わるフラッシュメモリで発生した不良ブロックの代替方法の違いがもたらす影響を説明するための図である。

【図10】本発明に係わるキャッシュメモリおよびフラッシュメモリ・チップの管理に関する説明をするための図である。

【図11】本発明に係わるキャッシュメモリおよびフラッシュメモリ・チップの管理に関する説明をするための図である。

【図12】本発明に係わるホスト計算機からのデータライト要求について、ストレージコントローラとメモリコントローラが行う処理を示すフローチャートである。

【図13】本発明に係わるデータリード要求について、ストレージコントローラとメモリ

50

コントローラが行う処理を示すフローチャートである。

【図14】本発明に係わるメモリコントローラによるキャッシュデータのライトバック処理を示すフローチャートである。

【図15】本発明に係わるストレージコントローラとメモリコントローラによる不良ブロックの代替処理を示すフローチャートである。

【図16】本発明に係わるメモリコントローラによるガベッジコレクション処理を示すフローチャートである。

【符号の説明】

【0162】

10 ... ストレージシステム	10
100 ... ホスト計算機	
120 ... ストレージコントローラ	
121, 122 ... チャンネルアダプタ	
123, 124 ... キャッシュメモリ	
125, 126 ... ストレージアダプタ	
129 ... 共有メモリ	
151 ~ 184 ... フラッシュメモリ・モジュール	
310 ... メモリコントローラ	
320 ... フラッシュメモリ	
321 ... フラッシュメモリ・チップ	20
500 ... キャッシュ管理テーブル	
510 ... アドレス管理テーブル	
520 ... ページ情報テーブル	
5001 ... インデックス	
5002 ... ウェイ	
5003 ... バリッドフラグ	
5004 ... キーレジスタ	
5005 ... データレジスタ	
5101 ... LPA (論理ページアドレス)	
5102 ... PPA (物理ページアドレス)	30
5201 ... ブロック	
5202 ... ステータス	

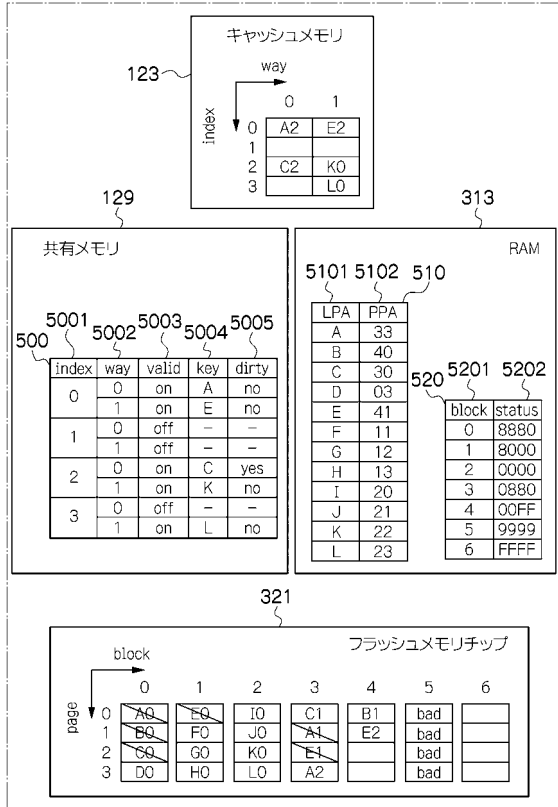
【図9】

図9

	T1	T2	T3
項目名	T1	T2	T3
フラッシュメモリの論理ページの冗長度	24/12 =200%	20/12 =167%	20/8 =250%
ガベージコレクション時の無効ページ含有率(平均)	8/20 =40%	4/16 =25%	8/16 =50%
ガベージコレクション時の有効ページ移動量(平均)	60%	75%	50%
フラッシュメモリのライト性能		低下	向上
フラッシュメモリのビット確率		不変	低下

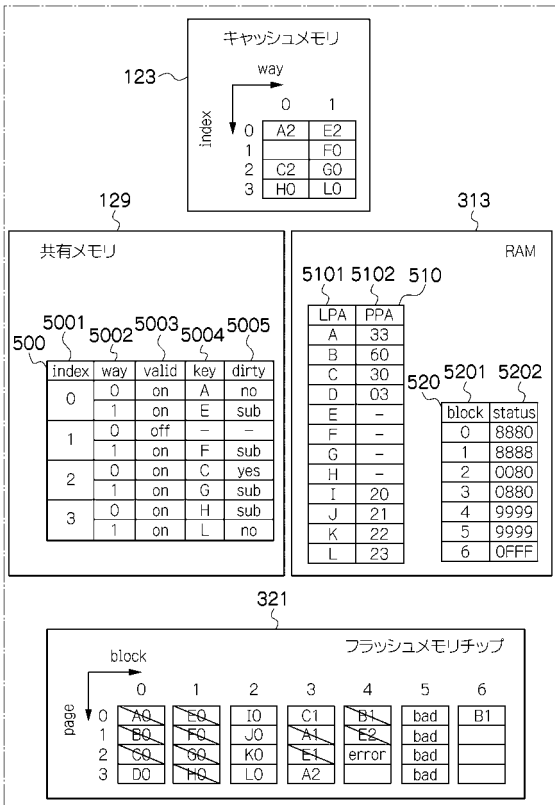
【図10】

図10



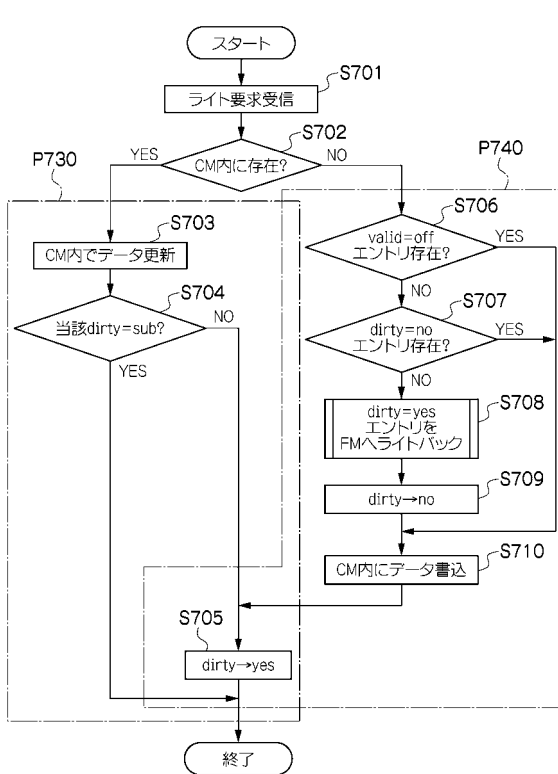
【図11】

図11



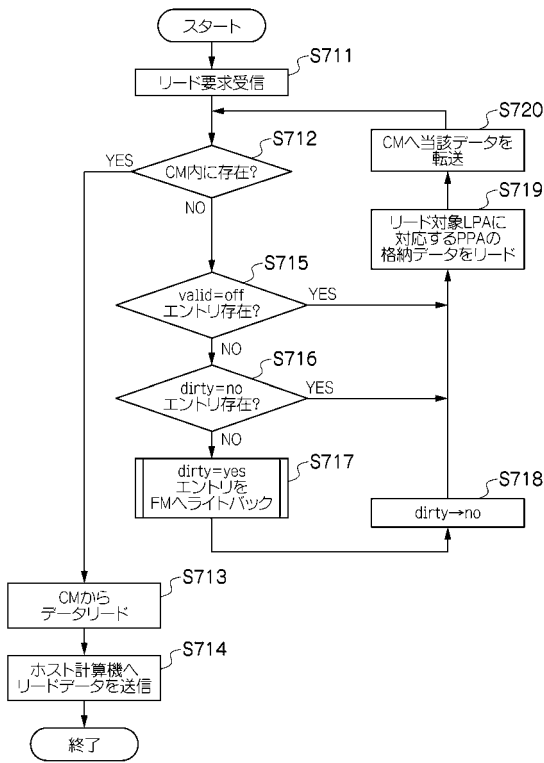
【図12】

図12



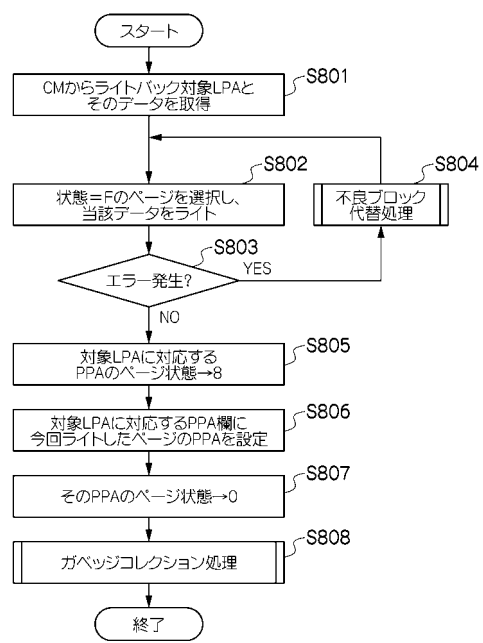
【図13】

図13



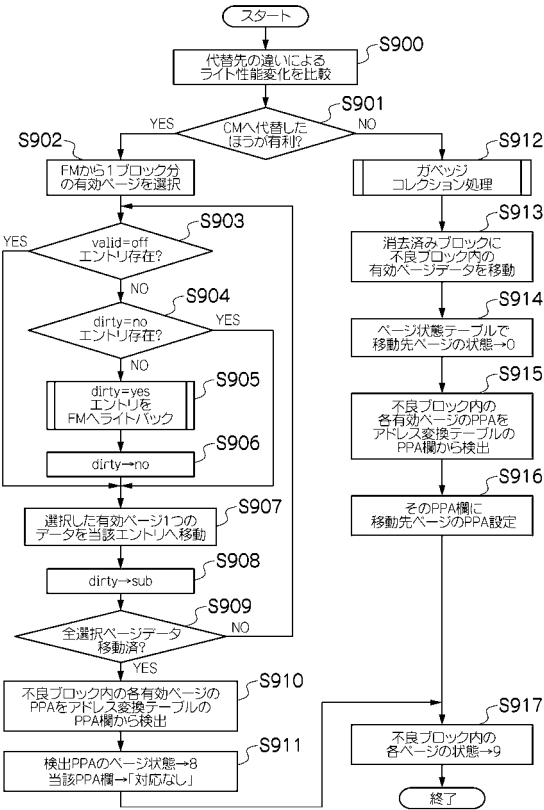
【図14】

図14



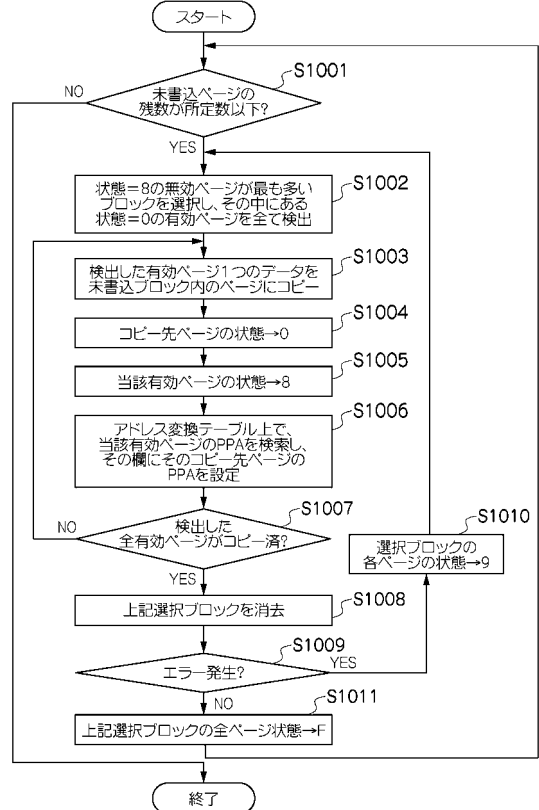
【図15】

図15



【図16】

図16



フロントページの続き

(51)Int.Cl.	F I	テーマコード(参考)
	G 0 6 F 12/16	3 2 0 L
	G 0 6 F 12/08	5 5 7
	G 0 6 F 12/08	5 4 1 B
	G 0 6 F 12/08	5 4 1 Z

Fターム(参考) 5B005 JJ11 MM11 VV13 WW14
5B018 GA04 GA06 GA10 HA04 HA35 KA14 KA18 MA22 NA01 NA06
QA16 RA11
5B065 BA01 BA09 CH02 CS01
5B082 JA07