

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6886486号  
(P6886486)

(45) 発行日 令和3年6月16日(2021.6.16)

(24) 登録日 令和3年5月18日(2021.5.18)

(51) Int.Cl. F I  
**G 0 6 F 3 / 0 6 (2006.01)**  
 G 0 6 F 3 / 0 6 3 0 6 Z  
 G 0 6 F 3 / 0 6 3 0 1 Z

請求項の数 6 (全 41 頁)

|            |                                     |           |   |
|------------|-------------------------------------|-----------|---|
| (21) 出願番号  | 特願2019-72298 (P2019-72298)          | (73) 特許権者 | 000005108                                 |
| (22) 出願日   | 平成31年4月4日(2019.4.4)                 |           | 株式会社日立製作所                                 |
| (62) 分割の表示 | 特願2017-511430 (P2017-511430)<br>の分割 |           | 東京都千代田区丸の内一丁目6番6号                         |
| 原出願日       | 平成27年4月9日(2015.4.9)                 | (74) 代理人  | 110002365<br>特許業務法人サンネクスト国際特許事務所          |
| (65) 公開番号  | 特開2019-106224 (P2019-106224A)       | (72) 発明者  | 山本 貴大<br>東京都千代田区丸の内一丁目6番6号 株<br>式会社日立製作所内 |
| (43) 公開日   | 令和1年6月27日(2019.6.27)                | (72) 発明者  | 藤本 和久<br>東京都千代田区丸の内一丁目6番6号 株<br>式会社日立製作所内 |
| 審査請求日      | 令和1年5月7日(2019.5.7)                  | (72) 発明者  | 坏 弘明<br>東京都千代田区丸の内一丁目6番6号 株<br>式会社日立製作所内  |

最終頁に続く

(54) 【発明の名称】 ストレージシステム及びデータ制御方法

(57) 【特許請求の範囲】

【請求項1】

データを記憶する複数の半導体記憶装置と、  
 仮想ボリュームに対して前記複数の半導体記憶装置に基づくプールから記憶領域を割り  
 当てるプロセッサとを備え、  
 前記複数の半導体記憶装置には、書き込み上限回数とセル当たりの記憶容量が異なる複  
 数の種類の半導体記憶装置が含まれ、セル当たりの記憶容量が大きい半導体記憶装置は、  
 書き込み上限回数が少なく、記憶容量当たりのコストが小さく、  
 前記プロセッサは、前記仮想ボリュームの記憶領域ごとに書き込み頻度を記録し、  
 前記プロセッサは、前記仮想ボリュームの記憶領域の書き込み頻度に基づいて、前記記  
 憶領域に割り当てられる半導体記憶装置の種類を変更し、  
 前記プロセッサは、前記半導体記憶装置の種類の変更では、書き込み頻度の高い記憶領  
 域を書き込み上限回数が多い種類の半導体記憶装置に割り当て、書き込み頻度の低い記憶  
 領域を書き込み上限回数が少ない種類の半導体記憶装置に割り当て、  
 前記プロセッサは、前記プール又は前記仮想ボリュームの容量と、前記プール又は前記  
 仮想ボリュームへ想定される書き込み頻度と、前記半導体記憶装置の各種類について単位  
 容量当たりの許容される書き込み頻度とに基づいて、前記記憶領域に割り当てるプールの  
 基となる前記半導体記憶装置にかかる推奨される種類ごとの記憶容量にかかる情報を算出  
 する、  
 ことを特徴とするストレージシステム。

10

20

## 【請求項 2】

請求項 1 において、

前記記憶領域に割り当てるプールの基となる前記半導体記憶装置にかかる推奨される種類ごとの記憶容量にかかる情報は、前記プール又は前記仮想ボリュームの容量のうち、当該種類の半導体記憶装置に基づく記憶領域の容量の比率であり当該種類の半導体記憶装置の使用する想定期間としての寿命が保証される推奨の比率である

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 において、

前記記憶領域に割り当てるプールの基となる前記半導体記憶装置にかかる推奨される種類ごとの記憶容量は、前記半導体記憶装置への書き込み回数とその書き込み上限回数とに基づいて半導体記憶装置を交換する場合に、ストレージシステムの運用管理コストを小さくする前記半導体記憶装置の種類ごとの容量である

ことを特徴とするストレージシステム。

## 【請求項 4】

請求項 3 において、

前記半導体記憶装置には、使用する想定期間が設定されており、

前記想定期間を使用するための半導体記憶装置が不足する場合に、前記半導体記憶装置の不足量を報知する

ことを特徴とするストレージシステム。

## 【請求項 5】

請求項 1 において、

所定の閾値よりも前記書き込み頻度の高い記憶領域を、書き込み上限回数が少ない種類の半導体記憶装置から書き込み上限回数が多い種類の半導体記憶装置に割り当てを移動させ、

前記所定の閾値は、書き込み上限回数が少ない種類の半導体記憶装置への書き込み頻度に基づいて変更される

ことを特徴とするストレージシステム。

## 【請求項 6】

データを記憶する複数の半導体記憶装置と、

仮想ボリュームに対して前記複数の半導体記憶装置に基づくプールから記憶領域を割り当てるプロセッサと、

を備えたストレージシステムにおける制御方法において、

前記複数の半導体記憶装置には、書き込み上限回数とセル当たりの記憶容量が異なる複数の種類の半導体記憶装置が含まれ、セル当たりの記憶容量が大きい半導体記憶装置は、書き込み上限回数が少なく、記憶容量当たりのコストが小さく、

前記プロセッサは、前記仮想ボリュームの記憶領域ごとに書き込み頻度を記録し、

前記プロセッサは、前記仮想ボリュームの記憶領域の書き込み頻度に基づいて、前記記憶領域に割り当てられる半導体記憶装置の種類を変更し、

前記プロセッサは、前記半導体記憶装置の種類の変更では、書き込み頻度の高い記憶領域を書き込み上限回数が多い種類の半導体記憶装置に割り当て、書き込み頻度の低い記憶領域を書き込み上限回数が多い種類の半導体記憶装置に割り当て、

前記プロセッサは、前記プール又は前記仮想ボリュームの容量と、前記プール又は前記仮想ボリュームへ想定される書き込み頻度と、前記半導体記憶装置の各種類について単位容量当たりの許容される書き込み頻度とに基づいて、前記記憶領域に割り当てるプールの基となる前記半導体記憶装置にかかる推奨される種類ごとの記憶容量にかかる情報を算出する、

ことを特徴とするストレージシステム制御方法。

## 【発明の詳細な説明】

## 【技術分野】

10

20

30

40

50

## 【0001】

本発明は、ストレージシステム及びデータ制御方法に関し、特に同一種別であって特性の異なる複数の記憶装置を利用して、記憶装置間でデータを自動的に再配置するストレージシステム及びデータ制御方法に適用して好適なものである。

## 【背景技術】

## 【0002】

近年、HDD (Hard Disk Drive) やSSD (Solid State Drive) の普及により、ストレージシステムに搭載される記憶装置の種類が多様化している。特にSSDについては、データ格納方式としてSLC (Single Level Cell)、2bit MLC (Multi Level Cell) 又はTLC (Triple Level Cell) 等があり、同一種別の記憶装置であつても、寿命 (書き込み上限回数) 及びビットコストが異なるという特性がある。

10

## 【0003】

寿命は、一般にSLCが最も長く、MLC、TLCの順に短くなる。またビットコストは、TLCが最も低く、MLC、SLCの順に高くなる。よってSSDを搭載するストレージシステムにおいては、その特性を考慮してデータを適切に配置することが重要となる。

## 【0004】

特許文献1及び2には、複数の記憶装置が搭載されたストレージシステムを管理又は制御する技術が開示されている。具体的に特許文献1には、複数の記憶装置をプールとして統合的に管理し、仮想的なボリュームをホスト装置に提供する技術が開示されている。

20

## 【0005】

また特許文献2には、HDDとSSDとを異なる階層の記憶装置として管理し、ホスト装置からのアクセス頻度に応じてデータを自動的に再配置し、高階層でのデータヒット率を向上させることにより、ホスト装置に対する応答性能を向上させる技術が開示されている。

## 【先行技術文献】

## 【特許文献】

## 【0006】

【特許文献1】米国特許第7447832号明細書

【特許文献2】米国特許第8041682号明細書

30

## 【発明の概要】

## 【発明が解決しようとする課題】

## 【0007】

しかし特許文献1及び2に記載の技術では、同一種別であって特性 (寿命及びビットコスト) の異なる複数の記憶装置 (例えばSSD) が混在したストレージシステムに対して書き込みを行った場合、特性が考慮されることなく、各記憶装置にデータが書き込まれる。

## 【0008】

SSDの特性を考慮せずに行き込みを行うと、書き込み上限回数の少ないSSDの寿命が書き込み上限回数の多いSSDよりも先に尽きることになり、寿命の尽きたSSDに対する交換回数が増え、ストレージシステムの運用管理コストが高くなるという課題がある。

40

## 【0009】

本発明は以上の点を考慮してなされたもので、同一種別であって特性の異なる記憶装置の寿命を延ばし、運用管理コストを削減し得るストレージシステム及びデータ制御方法を提案する。

## 【課題を解決するための手段】

## 【0010】

かかる課題を解決するために、本発明においては、第1の記憶装置と、第1の記憶装置よりも書き込み上限回数が少なく、かつ、単位面積当たりの記憶容量が多い第2の記憶装

50

置と、ホストに提供する仮想ボリュームに対して第1の記憶装置及び第2の記憶装置から記憶領域を割り当てるプロセッサとを備え、プロセッサは、仮想ボリュームに対して第2の記憶装置から割り当てられている記憶領域のうち、ホストからのライト頻度が予め定められたライト閾値よりも多い記憶領域に格納されているデータを第1の記憶装置の記憶領域に再配置する。

【0011】

またかかる課題を解決するために、本発明においては、第1の記憶装置と、第1の記憶装置よりも書き込み上限回数が少なく、かつ、単位面積当たりの記憶容量が多い第2の記憶装置と、ホストに提供する仮想ボリュームに対して第1の記憶装置及び第2の記憶装置から記憶領域を割り当てるプロセッサとを備えたストレージシステムのデータ制御方法であって、プロセッサが、仮想ボリュームに対して第2の記憶装置から割り当てられている記憶領域のうち、ホストからのライト頻度が予め定められたライト閾値よりも多い記憶領域に格納されているデータを第1の記憶装置の記憶領域に再配置する第1のステップと、仮想ボリュームに対して第1の記憶装置から割り当てられている記憶領域のうち、ホストからのライト頻度がライト閾値以下である記憶領域に格納されているデータを第2の記憶装置の記憶領域に再配置する第2のステップとを備える。

【発明の効果】

【0012】

本発明によれば、同一種別であって特性の異なる記憶装置の寿命を延ばし、運用管理コストを削減することができる。

【図面の簡単な説明】

【0013】

【図1】第1の実施の形態における計算機システムの全体構成図である。

【図2】ストレージシステムの論理構成図である。

【図3】TLC-MLC間で行われるページ再配置処理の概念構成図である。

【図4】TLC-TLC間で行われるページ再配置処理の概念構成図である。

【図5】共有メモリの内部構成図である。

【図6】ローカルメモリの内部構成図である。

【図7】ページ毎モニタテーブルの論理構成図である。

【図8】パリティグループ毎モニタテーブルの論理構成図である。

【図9】パリティグループ毎再配置管理テーブルの論理構成図である。

【図10】プール毎再配置管理テーブルの論理構成図である。

【図11】ホストI/O処理のフローチャートである。

【図12】デステージ処理のフローチャートである。

【図13】寿命情報採取処理のフローチャートである。

【図14】閾値決定処理のフローチャートである。

【図15】再配置処理のフローチャートである。

【図16】TLC-MLC間再配置処理のフローチャートである。

【図17】ライトリバランス処理のフローチャートである。

【図18】性能リバランス処理のフローチャートである。

【図19】新規割り当て決定処理のフローチャートである。

【図20】推奨容量算出方法の概念図である。

【図21】パラメータ設定画面の画面構成の一例である。

【図22】警告画面の画面構成の一例である。

【図23】第2の実施の形態における計算機システムの全体構成図である。

【図24】ストレージシステムの論理構成図である。

【図25】ライトデモーション処理によるページ再配置処理の概念構成図である。

【図26】パリティグループ毎再配置管理テーブルの論理構成図である。

【図27】プール毎再配置管理テーブルの論理構成図である。

【図28】再配置処理のフローチャートである。

10

20

30

40

50

- 【図29】ティア間再配置処理のフローチャートである。  
【図30】ライトデモーション処理のフローチャートである。  
【図31】新規割り当て決定処理のフローチャートである。  
【図32】ライト追加可能量及びライト削減要求量の算出方法の概念図である。  
【発明を実施するための形態】

#### 【0014】

以下、図面を参照しながら本発明の一実施の形態を説明する。なお以下に説明する実施の形態は、本発明を実現するための一例であって、本発明の技術的範囲を限定するものではない。また各図において共通の構成については、同一の参照番号を付してその説明を省略する。また各図において複数の同一部材を図示する場合は、「201A」、「201B」のようにアルファベットを付して区別する一方で、総称して呼ぶ場合には「201」のようにアルファベットを省略する。さらに本発明の実施の形態は、汎用コンピュータ上で稼動するソフトウェアに実装してもよいし、専用ハードウェアに実装してもよい。またソフトウェアとハードウェアとの組み合わせに実装してもよい。以下の説明では、管理用の情報をテーブル形式で説明するが、管理用の情報は必ずしもテーブルによるデータ構造で表現されなくてもよく、「リスト」、「DB」、「キュー」などについて単に「情報」と呼ぶことがある。また「プログラム」を主語（動作主体）として本発明の実施の形態における各処理について説明する場合がある。プログラムはプロセッサによって実行されることで定められた処理をメモリ及び通信ポート（通信制御装置）を用いながら行うため、プロセッサを主語とした説明としてもよい。プログラムの一部又は全部は専用ハードウェアで実現してもよく、またモジュール化されていてもよい。各種プログラムはプログラム配布サーバや記憶メディアによって各計算機にインストールされてもよい。

#### 【0015】

##### (1) 第1の実施の形態

第1の実施の形態は、記憶装置として半導体メモリ（SSD）のみを搭載するストレージシステムにおいて、特性の異なる半導体メモリ間（例えばTLC-MLC間）でデータを再配置することにより、半導体メモリの長寿命化を実現するとともに、ビットコストを低減しようとするものである。さらには同一種別の半導体メモリ間（例えばTLC-TLC間）でデータを再配置することにより、負荷分散を実現して長寿命化を実現しようとするものである。

#### 【0016】

##### (1-1) 計算機システムの全体構成

図1は、第1の実施の形態における計算機システム1の全体構成を示す。計算機システム1は、ホスト101、管理サーバ102、ネットワーク103及びストレージシステム104から構成される。ホスト101は、例えば一般的なサーバであり、ネットワーク103を介してストレージシステム104のポート105に接続される。ホスト101は、ストレージシステム104に対してリード要求又はライト要求を発行して、データの読み書きを実行する。

#### 【0017】

ネットワーク103は、例えばSAN（Storage Area Network）やイーサネット（登録商標）等の通信回線である。管理サーバ102は、ネットワーク103を介して、ストレージシステム104の保守I/F106又はポート105に接続される。ストレージ管理者は、管理サーバ102を操作して、ストレージシステム104を運用する上で必要な各種設定や管理を行う。

#### 【0018】

次にストレージシステム104の内部構成について説明する。ストレージシステム104は、ポート105、保守I/F106、プロセッサパッケージ107、キャッシュメモリ108、共有メモリ109、ドライブ111及びドライブ112等を備える。これらは内部ネットワーク110を介して互いに通信可能に接続される。

#### 【0019】

10

20

30

40

50

キャッシュメモリ108は、ストレージシステム104のI/O処理のスループットやレスポンスを向上させるために、データを一時的なキャッシュとして格納するための高速アクセス可能なメモリである。プロセッサパッケージ107は、ローカルメモリ113及びプロセッサ114を備えて構成される。なおここではプロセッサパッケージ107は、107A及び107Bの2つを図示しているが、数はこれに限定されない。

【0020】

プロセッサ114は、ホスト101からのリード要求及びライト要求を処理するために、ドライブ111及び112と、キャッシュメモリ108との間のデータの転送処理を実行する。共有メモリ109は、プロセッサ114がリード要求又はライト要求を処理し、またストレージシステム104の機能(ボリュームのコピー機能など)を実行する上で必要な制御用の情報を格納するメモリである。共有メモリ109は、複数のプロセッサパッケージ107(ここでは107A及び107B)のプロセッサ114により共有される情報を格納する。

10

【0021】

ローカルメモリ113は、プロセッサ114がリード要求又はライト要求を処理し、またストレージシステム104の機能を実行する上で必要な制御用の情報を格納するメモリである。ローカルメモリ113は、各プロセッサ114により占有される情報を格納する。ローカルメモリ113には、例えばプロセッサ114により実行するプログラムが格納される。

【0022】

ドライブ111及び112は、複数の半導体メモリ(SSD: Solid State Drive)から構成される記憶装置である。ここでのドライブ111は、1つのセルに2ビットの情報を格納するMLC SSD (Multi Level Cell SSD)であり、ドライブ112は、1つのセルに3ビットの情報を格納するTLC (Triple Level Cell SSD)である。

20

【0023】

MLC SSDと、TLC SSDとの違いは、上記の通りセル当たりの情報量が異なることに加えて、寿命(書き込み上限回数)及びビットコストが異なる。寿命は、一般にセル当たりの情報量が多くなるほど短くなる。よってMLC SSDの方がTLC SSDよりも長い。またビットコストは、単位面積当たりの情報量が多いほど低く抑えられるため、MLC SSDの方がTLC SSDよりも高い。

30

【0024】

すなわちMLC SSDは、寿命が長くてビットコストが高く、TLC SSDは、寿命が短くてビットコストが低いと言える。なおここではMLC SSD及びTLC SSDを利用してドライブ111及び112を構成するとしているが、これに限らず、SLC SSD (Single Level Cell SSD)を利用するとしてもよい。

【0025】

ドライブ111及び112について、以降は同一種別であって特性の異なる記憶装置であることを強調する場合にはドライブ111をMLC SSDと呼び、ドライブ112をTLC SSDと呼ぶ場合がある。また複数のMLC SSDから構成されるパリティグループをMLC PGと呼び、複数のTLC SSDから構成されるパリティグループをTLC PGと呼ぶ場合がある。

40

【0026】

(1-2) ストレージシステムの論理構成

図2は、ストレージシステム104の論理構成を示す。プロセッサ114がホスト101に提供する仮想ボリューム202A及び202Bは、ホスト101から認識される論理的な記憶領域であり、ホスト101からのリード要求又はライト要求の発行対象となる記憶領域である。

【0027】

プール204は、1個以上のプールボリューム206A~206Eにより構成される。プールボリューム206A~206Eは、それぞれMLC SSD又はTLC SSDの

50

何れかの記憶領域から構成される。一又は複数のプールボリューム 206A ~ 206E により、パリティグループ (PG: Parity Group) 205A 又は 205B が形成される。

【0028】

なおここでのパリティグループ 205A は、MLC SSD のみで構成され、パリティグループ 205B は TLC SSD のみで構成されているが、必ずしもこれに限らず、MLC SSD 及び TLC SSD が混在して 1 つのパリティグループを構成するとしてもよい。プールボリューム 206 は、パリティグループ 205 の一部又は全部の領域が切り出されて使用される。

【0029】

プロセッサ 114 は、例えばホスト 101 が仮想ボリューム 202A に対してライト要求を発行した場合、この仮想ボリューム 202A においてライト要求の対象となった仮想的な記憶領域に対し、未使用の実記憶領域をプール 204 から所定単位 (ページ) で割り当てる。

10

【0030】

なおページとは、データを書き込む際の最小単位の記憶領域である。ここでは仮想ボリューム 202A 又は 202B に対して割り当てられた仮想ページを 201A ~ 201E として図示しており、これらのページ 201A ~ 201E に対して割り当てるプールボリューム 206A ~ 206E における実ページを 207A ~ 207F として図示している。

【0031】

次回ホスト 101 から同じページ 201A に対してリード要求又はライト要求が発行された場合、プロセッサ 114 は、既に割り当てられているプールボリューム 206A の記憶領域 207A に対して I/O 処理を実行することにより、あたかもホスト 101 が仮想ボリューム 202A に対して I/O 処理を実行しているように処理することができる。

20

【0032】

すなわち仮想ボリューム 202A 又は 202B を用いて、使用する部分のみプールボリューム 206A ~ 206E の記憶領域 (ページ) 207A ~ 207F を割り当てることにより、限られた記憶容量を効率的に使用することができる。

【0033】

ここで、各仮想ボリューム 202A 又は 202B を構成するページ 201 毎に、ホスト 101 からの単位時間当たりの書き込み回数 (これをライト頻度と呼ぶ) は異なる。よって例えばライト頻度が高いページ 201A は、特性として書き込み上限回数の多い MLC SSD で構成されたパリティグループ 205A に配置することで、TLC SSD のような書き込み上限回数の少ない記憶装置の寿命を延ばすことができる。

30

【0034】

上記の再配置は、具体的にはページ 207C に格納されたデータを未使用のページ 207B にコピーし、仮想ボリューム 202A のページ 201C とプールボリューム 206B のページ 207C との対応付けを仮想ボリューム 202A のページ 201C とプールボリューム 206A のページ 207B との対応付けに変更することにより実行する。

【0035】

(1-3) ページ配置処理の概念構成

40

図 3 は、TLC - MLC 間で行われるページ再配置処理の概念構成を示す。ここでは MLC SSD で構成されたパリティグループ (以下、MLC PG と呼ぶ) 205A と、TLC SSD で構成されたパリティグループ (以下、TLC PG と呼ぶ) 205B との間で、ページ 301 を再配置する。ページ 301 の再配置先は、MLC PG 205A 及び TLC 205B 上のページ 301 毎に採取しているモニタ情報に基づいて決定される。

【0036】

具体的には、ページ 301 毎にライト回数を一定期間採取し、周期の満了後、モニタ情報に基づいて算出された MLC PG 205A 及び TLC PG 205B 毎のライト頻度と、TLC - MLC 間ライト閾値 302 とに基づいてページ 301 の再配置先を決定する

50

。

## 【0037】

図3では、TLC - MLC間ライト閾値302が2000以上のライト負荷のページ301をMLC PG205Aに再配置し(矢印303A)、2000未満のライト負荷のページ301をTLC PGに再配置する(矢印303B)。

## 【0038】

このように寿命の異なるSSDで構成されたMLC PG205AとTLC PG205Bとの間で、ライト頻度の高いページを書き込み上限回数の多いSSDで構成されたMLC PG205Aに再配置し、ライト頻度の低いページを書き込み上限回数の少ないSSDで構成されたTLC PG205Bに再配置する。

10

## 【0039】

この結果、書き込み上限回数の少ないSSDで構成されたTLC PG205Bにおけるライト回数を削減することができ、書き込み上限回数の少ないSSD(ここではTLC SSD)の寿命を延ばすことができる。

## 【0040】

図4は、同一特性のドライブ間(TLC - TLC間又はMLC - MLC間)で行われるページ再配置処理の概念構成を示す。同一特性のドライブ間で行われるページ再配置処理は、書き込み上限回数が同程度であるSSDで構成されたパリティグループ間でのライト頻度を調整することを目的として行われる。

## 【0041】

20

実際にはパリティグループ間のライト頻度を調整するライトリバランスと、ライトリバランスにより崩れたパリティグループ間のI/O頻度を調整する性能リバランスとの2種類のリバランス処理を行う。

## 【0042】

ここではTLC PG間で行うライトリバランス及び性能リバランスについて説明する。まずライトリバランスでは、寿命を保証するために削減が必要なライト頻度(これをライト削減要求量と呼ぶ)又は寿命を保証することができる範囲内で追加可能なライト頻度(これをライト追加可能量と呼ぶ)をパリティグループ毎に算出する。

## 【0043】

このときライト削減要求量が正の値のパリティグループをライト削減PG205Bとし、ライト追加可能量が正の値のパリティグループをライト追加PG205Cとする。そしてライト削減PG205Bのライト高負荷なページ401Aとライト追加PG205Cのライト低負荷なページ401Bとを互いに移動して(矢印403)、パリティグループ間のライト頻度を調整する。

30

## 【0044】

ライトリバランスの際の移動対象となるページは、モニタ情報を元にした閾値に基づいて決定される。閾値には、ライト高負荷なページを決定する閾値(これをライト削減閾値と呼ぶ)405と、ライト低負荷なページを決定する閾値(これをライト追加閾値と呼ぶ)406とがある。ライトリバランスの際、ライト削減PG205Bでは、ライト削減閾値405以上のライト頻度のページ401Aを移動対象とし、ライト追加PG205Cでは、ライト追加閾値406以下のライト頻度のページ401Bを移動対象とする。そしてそれぞれのページを移動することでライト頻度を調整する。

40

## 【0045】

移動するページ数は、1回のページ移動で移動するライト頻度を算出し、目標とするライト削減要求量を満たすために必要な移動ページ数(これを移動計画ページ数と呼ぶ)を予め算出することにより決定される。そしてこの算出された移動計画ページ数分だけページの移動(再配置)が行われる。

## 【0046】

次に性能リバランスでは、ライトリバランスにより移動したライト頻度の分だけリード頻度をライトリバランスとは逆の方向に移動することで、各パリティグループのリード頻

50

度及びライト頻度を合算したI/O頻度をパリティグループ間で調整する。具体的には、ライト削減PG205Bのリード低負荷なページ402Bとライト追加PG205Cのリード高負荷なページ402Aとを移動して(矢印404)、パリティグループ間のリード頻度を調整する。

【0047】

性能リバランスの際の移動対象となるページは、モニタ情報を元にした閾値に基づいて決定される。閾値には、リード高負荷なページを決定する閾値(これをリード削減閾値と呼ぶ)407とリード低負荷なページを決定する閾値(これをリード追加閾値と呼ぶ)408がある。

【0048】

性能リバランスの際、ライト削減PG205Bでは、リード追加閾値408以下のリード頻度のページ402Bを移動対象とし、ライト追加PG205Cでは、リード削減閾値407以上のリード頻度のページ402Aを移動対象とする。そしてそれぞれのページを移動することでリード頻度を調整する。移動するページ数は、ライトリバランスと同様、移動計画ページ数を予め算出することにより決定される。そしてこの移動計画ページ数だけページの移動(再配置)が行われる。

【0049】

(1-4)メモリの内部構成

図5は、共有メモリ109の内部構成を示す。共有メモリ109には、ページ毎モニタテーブル501、パリティグループ毎モニタテーブル502、パリティグループ毎再配置管理テーブル503、プール毎再配置管理テーブル504、キャッシュ管理テーブル505、ダイナミックマッピングテーブル506及び論理物理アドレス変換テーブル507が格納される。

【0050】

ページ毎モニタテーブル501は、I/O回数を含む各ページのモニタ情報を管理するテーブルであり、パリティグループ毎モニタテーブル502は、I/O回数を含む各パリティグループのモニタ情報を管理するテーブルである。またパリティグループ毎再配置管理テーブル503は、パリティグループ毎のページ再配置に関する制御情報を管理するテーブルである。

【0051】

プール毎再配置管理テーブル504は、プール毎のページ再配置に関する制御情報を管理するテーブルであり、キャッシュ管理テーブル505は、キャッシュメモリ108にデータを格納する際にキャッシュメモリ108にあるデータのダーティ/クリーン状態管理を保持するテーブルである。

【0052】

またプール毎再配置管理テーブル504は、仮想ボリューム202に対して書き込むデータをキャッシュする場合には、キャッシュメモリ108のアドレスと、対応する仮想ボリューム202のページ201を特定するアドレスとを対応づけて管理するテーブルである。

【0053】

ダイナミックマッピングテーブル506は、仮想ボリューム202の各ページ201と、各ページ201に割り当てられているプールボリューム206のページ207と、各ページ201のモニタ情報との対応関係を管理するテーブルである。

【0054】

論理物理アドレス変換テーブル507は、パリティグループと、プールボリュームと、プールボリュームのデータを格納する物理ディスクに対応するパリティグループのアドレスとの対応関係を管理するテーブルである。

【0055】

図6は、ローカルメモリ113の内部構成を示す。ローカルメモリ113には、ホストI/O処理プログラム601、デステージ処理プログラム602、寿命情報採取処理プロ

10

20

30

40

50

グラム603、閾値決定処理プログラム604、再配置処理プログラム604A及び新規割り当て決定処理プログラム605が格納される。これらの各種プログラムは、プロセッサ114により実行される。

【0056】

ホストI/O処理プログラム601は、ホスト101からのI/O要求を受領した場合に仮想ボリューム202に対するリード要求又はライト要求を処理するプログラムである。デステージ処理プログラム602は、キャッシュメモリ108上の物理ディスクに未反映のデータを物理ディスクに格納するプログラムである。この処理は、ホスト101からのI/O要求に対する処理とは非同期に実行される。

【0057】

寿命情報採取処理プログラム603は、ドライブ111及び112に対して所定周期でコマンドを発行して寿命情報を採取し、採取した情報を共有メモリ109に反映するプログラムである。閾値決定処理プログラム604は、所定周期で採取したモニタ情報とドライブ111及び112の寿命情報とに基づいて、ページ再配置のための閾値を算出するプログラムである。

【0058】

再配置処理プログラム604Aは、閾値決定処理プログラム604により呼び出されるプログラムであり、閾値決定処理プログラム604により決定された各種閾値に基づいて、ページを再配置するプログラムである。新規割り当て決定処理プログラム605は、ホストI/O処理プログラム601に同期して実行され、仮想ボリューム202における新規の仮想ページに対して実ページの割り当て先であるパリティグループ205を閾値に基づいて決定するプログラムである。

【0059】

(1-5) テーブル構成

図7は、ページ毎モニタテーブル501の論理構成を示す。ページ毎モニタテーブル501は、ページ番号欄701、ライトI/Oカウンタ欄702、リードI/Oカウンタ欄703、合計I/Oカウンタ欄704及び新規ライトフラグ欄705から構成される。

【0060】

ページ番号欄701には、ページ201を特定するページ番号が格納され、ライトI/Oカウンタ欄702には、一定周期のライト回数が格納される。またリードI/Oカウンタ欄703には、一定周期のリード回数が格納され、合計I/Oカウンタ欄704には、一定周期のリード回数及びライト回数の合計I/O回数が格納される。

【0061】

周期は、上記した閾値決定処理プログラム604がモニタ情報を採取する周期と同じであり、閾値決定処理プログラム604は、この一定期間のモニタ情報を処理対象とする。また新規ライトフラグ欄705には、ページが新規割り当てページか否かを示すフラグが格納される。

【0062】

図8は、パリティグループ毎モニタテーブル502の論理構成を示す。パリティグループ毎モニタテーブル502は、パリティグループ番号欄801、最大ライト頻度欄802、最小ライト頻度欄803、最大リード頻度欄804、最小リード頻度欄805、リード/ライト比率欄806、ライト追加可能量欄807、ライト削減要求量欄808、新規ライトI/Oカウンタ欄809、新規ライト比率欄810、平均I/O頻度欄811及び割り当てページ数欄812から構成される。

【0063】

パリティグループ番号欄801には、パリティグループ205を特定するパリティグループ番号が格納され、最大ライト頻度欄802には、パリティグループ内のページの最大ライト頻度が格納される。また最小ライト頻度欄803には、パリティグループ内のページの最小ライト頻度が格納される。

【0064】

10

20

30

40

50

最大リード頻度欄 804 には、パリティグループ内のページの最大リード頻度が格納され、最小リード頻度欄 805 には、パリティグループ内のページの最小リード頻度が格納される。またリード/ライト比率欄 806 には、パリティグループに対するリード回数とライト回数との比率が格納される。

【0065】

ライト追加可能量欄 807 には、寿命を保証することができる範囲内でパリティグループに追加可能なライト頻度が格納され、ライト削減要求量欄 808 には、寿命を保証するためにパリティグループから削減が必要なライト頻度が格納される。ライト追加可能量欄 807 及びライト削減要求量欄 808 には、パリティグループのライト頻度及び SSD の寿命情報に基づいて算出される値が格納され、何れかが正の値となる。

10

【0066】

ライト追加可能量欄 807 及びライト削減要求量欄 808 に格納される値は、SSD から採取できる寿命情報がライト追加可能率及びライト削減要求率である場合、下記式 1 及び 2 を計算して算出することができる。

【0067】

なおライト追加可能率とは、現状のライト頻度を 100% として、追加可能なライト頻度の割合であり、ライト削減要求率とは、現状のライト頻度を 100% として寿命を維持するために削減すべきライト頻度の割合である。

【0068】

【数 1】

$$\begin{aligned} & \text{ライト追加可能量 [IOPH (Input Output Per Hour)]} \\ & = \text{Min (パリティグループを構成する全 SSD のライト追加可能率 [\%])} \\ & \quad \times \text{パリティグループのライト頻度 [IOPH]} \\ & \dots\dots\dots (1) \end{aligned}$$

20

【0069】

【数 2】

$$\begin{aligned} & \text{ライト削減要求量 [IOPH]} \\ & = \text{Max (パリティグループを構成する全 SSD のライト削減要求率 [\%])} \\ & \quad \times \text{パリティグループのライト頻度 [IOPH]} \\ & \dots\dots\dots (2) \end{aligned}$$

30

【0070】

新規ライト I/O カウンタ欄 809 には、パリティグループに対する新規ライト回数が格納され、新規ライト比率欄 810 には、パリティグループにおける書き込み処理のうち、更新ライトと新規ライトとの比率が格納される。また平均 I/O 頻度欄 811 には、パリティグループにおける各ページの平均 I/O 頻度が格納され、割り当てページ数欄 812 には、パリティグループに割り当てられているページ数が格納される。

【0071】

図 9 は、パリティグループ毎再配置管理テーブル 503 の論理構成を示す。パリティグループ毎再配置管理テーブル 503 は、パリティグループ番号欄 901、メディアタイプ欄 902、移動元 PG 種別欄 903、移動先 PG 欄 904、ライト削減閾値欄 905、ライト追加閾値欄 906、リード削減閾値欄 907、リード追加閾値欄 908、移動計画ページ数欄 909A 及び 909B、移動実績ページ数欄 910A 及び 910B 並びに新規ライト可能量欄 911 から構成される。

40

【0072】

パリティグループ番号欄 901 には、パリティグループを特定するパリティグループ番号が格納され、メディアタイプ欄 902 には、パリティグループを構成する SSD の特性の情報が格納される。また移動元 PG 種別欄 903 には、ライト削減 PG 又はライト追加 PG の何れかを示す情報が格納される。移動元 PG 種別は、パリティグループ毎モニタテ

50

ーブル502のライト追加可能量欄807又はライト削減要求量欄808に格納される情報に基づいて決定される。

【0073】

具体的には、ライト追加可能量欄807に正の値が格納されている場合、このパリティグループの移動元PG種別欄903には、ライト追加PGを示す情報が格納される。またライト削減要求量欄808に正の値が格納されている場合、このパリティグループの移動元PG種別903には、ライト削減PGを示す情報が格納される。

【0074】

移動先PG欄904には、ページ再配置実行時の移動先のパリティグループ番号が格納される。移動先のパリティグループ番号は、再配置進捗度が最低のパリティグループを移動先として決定してもよい。再配置進捗度は、移動計画ページ数欄909A及び909Bと、移動実績ページ数欄910A及び910Bとを用いて、下記式3を計算して算出することができる。

10

【0075】

【数3】

$$\text{再配置進捗度} = \text{移動実績ページ数} / \text{移動計画ページ数} \dots \dots \dots (3)$$

【0076】

ライト削減閾値欄905には、ライト削減PGでのライトリバランス時の移動対象のページを決定するための閾値が格納され、ライト追加閾値欄906には、ライト追加PGでのライトリバランス時の移動対象のページを決定するための閾値が格納される。

20

【0077】

リード削減閾値欄907には、ライト追加PGでの性能リバランス時の移動対象ページを決定するための閾値が格納され、リード追加閾値欄908には、ライト削減PGでの性能リバランス時の移動対象ページを決定するための閾値が格納される。ライト削減閾値、ライト追加閾値、リード削減閾値及びリード追加閾値は、下記式4を計算して算出することができる。

【0078】

【数4】

$$\begin{aligned}
&\text{ライト削減閾値 [IOPH]} = \text{最大ライト頻度 [IOPH]} \times 0.7 \\
&\text{ライト追加閾値 [IOPH]} = \text{最小ライト頻度 [IOPH]} \times 1.3 \\
&\text{リード削減閾値 [IOPH]} = \text{最大リード頻度 [IOPH]} \times 0.7 \\
&\text{リード追加閾値 [IOPH]} = \text{最小リード頻度 [IOPH]} \times 1.3 \\
&\dots \dots \dots (4)
\end{aligned}$$

30

【0079】

移動計画ページ数(ライトリバランス)欄909Aには、同一特性のSSDにより構成されるパリティグループ間でライト頻度を調整するために必要なページの移動数が格納され、移動計画ページ数(性能リバランス)欄909Bには、同一特性のSSDにより構成されるパリティグループ間でリード頻度を含むI/O頻度を調整するために必要なページの移動数が格納される。ライトリバランス用の移動計画ページ数は、下記式5を計算して算出することができる。

40

【0080】

【数5】

ライトリバランス用の移動計画ページ数 [ページ]  
 = 目標移動量 [IOPH] ÷ 1回のページ移動で変化するI/O頻度 [IOPH/ページ]  
 目標移動量 [IOPH]  
 = Min (ライト追加可能量 [IOPH]、ライト削減要求量 [IOPH])  
 1回のページ移動で変化するI/O頻度 [IOPH/ページ]  
 = (ライト追加PGの移動対象ページの平均I/O頻度 [IOPH]  
 - ライト削減PGの移動対象ページの平均I/O頻度 [IOPH]) ÷ 2  
 . . . . . (5)

10

【0081】

また性能リバランス用の移動計画ページ数は、下記式6を計算して算出することができる。

【0082】

【数6】

性能リバランス用の移動計画ページ数 [ページ]  
 = 目標移動量 [IOPH] ÷ 1回のページ移動で変化するI/O頻度 [IOPH/ページ]  
 目標移動量 [IOPH]  
 = ( (ライトリバランスで移動したライト頻度 [IOPH] × 2) +  
 (ライト追加PGの平均I/O頻度 [IOPH]  
 - ライト削減PGの平均I/O頻度 [IOPH]) ) ÷ 2  
 1回のページ移動で変化するI/O頻度 [IOPH]  
 = (ライト追加PGの移動対象ページの平均I/O頻度 [IOPH]  
 - ライト削減PGの移動対象ページの平均I/O頻度 [IOPH]) ÷ 2  
 . . . . . (6)

20

【0083】

移動実績ページ数(ライトリバランス)欄910Aには、ライトリバランスのために他のパリティグループに移動したページ数が格納され、移動実績ページ数(性能リバランス)欄910Bには、性能リバランスのために他のパリティグループに移動したページ数が格納される。新規ライト可能量欄912には、パリティグループで処理することのできる新規ライト回数が格納される。新規ライト可能量は、下記式7を計算して算出することができる。

30

【0084】

【数7】

新規ライト可能量 [I/O]  
 = パリティグループの移動後のライト頻度 [IOPH]  
 × 新規割り当て比率 [%] × 再配置周期 [H]  
 . . . . . (7)

【0085】

図10は、プール毎再配置管理テーブル504の論理構成を示す。プール毎再配置管理テーブル504は、プール番号欄1001、寿命制御再配置周期欄1002、TLC-MLC間ライト閾値欄1003、新規ライト閾値欄1004、ワークロードタイプ欄1005、同種ドライブ間新規割り当てポリシ欄1006、新規割り当てバッファ欄1007及び再配置バッファ欄1008から構成される。

40

【0086】

プール番号欄1001には、プールを特定するプール番号が格納され、寿命制御再配置周期欄1002には、プールにおいて寿命制御によるページの再配置を実行する周期が格納される。またTLC-MLC間ライト閾値欄1003には、TLC PGとMLC PGとの間でページを再配置する際にどちらの種類のパリティグループにページを配置するのかを決定するための閾値が格納される。TLC-MLC間ライト閾値は、下記式8を計

50

算して算出することができる。

【0087】

【数8】

TLC-MLC間ライト閾値 [IOPH]  
 = Ave (全TLC PGの限界ライト閾値 [IOPH])  
 TLC PGの限界ライト閾値 [IOPH]  
 = IF (当該パリティグループの移動元PG種別がライト追加PG?)  
 当該パリティグループのライト頻度 [IOPH] + ライト追加可能量 [IOPH]  
 ELSE IF (当該パリティグループの移動元PG種別がライト削減PG?)  
 当該PGのライト頻度 [IOPH] - ライト削減要求量 [IOPH]  
 . . . . . (8)

10

【0088】

新規ライト閾値欄1004には、プールにおいて新規ページに対する割り当て処理を行う際、TLC PGとMLC PGのどちらから新規ページに対する実ページの割り当てを行うかを決定するための閾値が格納される。新規ライト閾値は、ページ移動時のペナルティ (= 1ページあたりのデータサイズ) としてストレージシステム104が指定してもよい。

【0089】

ワークロードタイプ欄1005には、プールに関連付けられている仮想ボリュームに対してホスト101から発行されるI/O要求の特性の情報が格納される。ワークロードタイプには、例えばWrite intensive、Read intensive、Unknown等がある。

20

【0090】

Write intensiveは、ホスト101からのI/O要求のうち、ライト比率が高いことを意味し、Read intensiveは、ホスト101からのI/O要求のうち、リード比率が高いことを意味し、Unknownは、ホスト101からのI/O要求のうち、リード/ライト比率が不明であることを意味する。ワークロードタイプは、ユーザが指定してもよいし、ストレージシステム104が自動で決定してもよい。

【0091】

同種ドライブ間新規割り当てポリシー欄1006には、新規割り当て処理でTLC PG内又はMLC PG内でどのパリティグループから新規ページに対する実ページの割り当てを行うかを決定するための情報が格納される。新規割り当てポリシーには、例えばパリティグループ間で順番に割り当てるラウンドロビン、容量の多いパリティグループから割り当てる容量優先又はライト追加可能量の多いパリティグループから割り当てる寿命優先がある。新規割り当てポリシーは、ユーザが指定してもよいし、ストレージシステム104が自動で決定してもよい。

30

【0092】

新規割り当てバッファ欄1007には、プールで新規割り当て処理のために使用するバッファのプール容量に対する比率が格納される。再配置バッファ欄1008には、プールでページの再配置処理のために使用するバッファのプール容量に対する比率が格納される。なおこのプール毎再配置管理テーブル504内には、再配置の制御情報のデフォルト値を格納するエントリが1つ存在し、ユーザからの指定がない制御情報にはデフォルト値が格納される。

40

【0093】

(1-6) フローチャート

図11は、ホストI/O処理のフローチャートを示す。このホストI/O処理は、ストレージシステム104がホスト101からのI/O要求を受領したことを契機として、プロセッサ114とホストI/O処理プログラム601との協働により実行される。説明の便宜上、処理主体をホストI/O処理プログラム601として説明する。

50

## 【 0 0 9 4 】

まずホスト I / O 処理プログラム 6 0 1 は、ホスト 1 0 1 からの I / O 処理要求を受領すると、受領した I / O 処理要求が仮想ボリューム 2 0 2 に対してデータを書き込むライト要求であるか否かを判断する ( S 1 1 0 1 ) 。

## 【 0 0 9 5 】

ライト要求である場合 ( S 1 1 0 1 : Y )、ホスト I / O 処理プログラム 6 0 1 は、ダイナミックマッピングテーブル 5 0 6 を参照して、ライト対象の仮想ページに対して実ページが割り当て済みであるか否かを判断する。そして未割り当ての場合には、未使用の実ページを割り当てる新規割り当て決定処理を実行する ( S 1 1 0 9 )。新規割り当て決定処理の詳細については後述する ( 図 1 9 ) 。

10

## 【 0 0 9 6 】

新規割り当て決定処理を実行した後、次いでホスト I / O 処理プログラム 6 0 1 は、仮想ボリューム 2 0 2 上のライト対象のアドレスに対応した領域がキャッシュメモリ 1 0 8 上に確保されているか否かを確認し、確保されていない場合にはキャッシュメモリ 1 0 8 上の領域を確保する ( S 1 1 1 0 ) 。

## 【 0 0 9 7 】

次いでホスト I / O 処理プログラム 6 0 1 は、ホスト 1 0 1 に対してライト用のデータを送信するように通知する。ホスト I / O 処理プログラム 6 0 1 は、ホスト 1 0 1 からデータが送信されると、このデータをキャッシュメモリ 1 0 8 上の確保した領域に書き込む ( S 1 1 1 1 )。そしてホスト I / O 処理プログラム 6 0 1 は、まだドライブ 1 1 1 又は 1 1 2 に書き込みが完了していない領域であることを示すダーティフラグをキャッシュ管理テーブル 5 0 5 において ON に設定する。

20

## 【 0 0 9 8 】

ここで、ダーティフラグは、キャッシュメモリ 1 0 8 にのみデータが格納されており、ドライブ 1 1 1 又は 1 1 2 にはまだ格納されていない状態を示す情報である。ダーティフラグは、キャッシュメモリ 1 0 8 の領域を管理するキャッシュ管理テーブル 5 0 5 において ON 又は OF に設定される。キャッシュ管理テーブル 5 0 5 は、後述するデステージ処理 ( 図 1 2 ) において参照される。

## 【 0 0 9 9 】

そしてダーティフラグが ON である場合にはキャッシュメモリ 1 0 8 上の領域に格納されているデータは、ドライブ 1 1 1 又は 1 1 2 に書き込まれる。ドライブ 1 1 1 又は 1 1 2 に書き込まれた後は、ダーティフラグは OFF に設定され、リード処理に対応してドライブ 1 1 1 又は 1 1 2 から読み込んだデータをキャッシュメモリ 1 0 8 に格納した場合を含めて、キャッシュ管理テーブル 5 0 5 にはクリーンフラグが ON に設定される。

30

## 【 0 1 0 0 】

このようにキャッシュ管理テーブル 5 0 5 は、キャッシュメモリ 1 0 8 上のアドレスと対応する仮想ボリューム 2 0 2 のアドレスと、キャッシュメモリ 1 0 8 上のデータの状態とを少なくとも管理する。なおキャッシュメモリ 1 0 8 のアドレスに対応する仮想ボリューム 2 0 2 のアドレスは、仮想ボリューム 2 0 2 のデータを置くためにキャッシュメモリ 1 0 8 上の領域を確保した場合にのみ、有効な値である仮想ボリューム 2 0 2 のアドレスが格納される。

40

## 【 0 1 0 1 】

以上の処理を行った後、ホスト I / O 処理プログラム 6 0 1 は、ホスト 1 0 1 に I / O 処理 ( ライト処理 ) が完了したことを通知して ( S 1 1 1 2 )、本処理を終了する。

## 【 0 1 0 2 】

ステップ S 1 1 0 1 に戻り、ホスト I / O 処理プログラム 6 0 1 は、ホスト 1 0 1 から受領した I / O 要求が仮想ボリューム 2 0 2 からデータを読み込むリード要求である場合 ( S 1 1 0 1 : N )、キャッシュ管理テーブル 5 0 5 を参照して、リード要求に対応する仮想ボリューム 2 0 2 上のアドレスに対応したデータがキャッシュメモリ 1 0 8 上に格納されているか否かを判断する ( S 1 1 0 2 ) 。

50

## 【 0 1 0 3 】

リード要求に対応する仮想ボリューム 2 0 2 上のアドレスに対応したデータがキャッシュメモリ 1 0 8 上に格納されている場合をキャッシュヒットと呼ぶ。キャッシュヒットである場合 ( S 1 1 0 2 : Y )、ホスト I / O 処理プログラム 6 0 1 は、キャッシュメモリ 1 0 8 上のデータをホスト 1 0 1 に転送するとともに ( S 1 1 0 8 )、I / O 処理 ( リード処理 ) が完了したことをホスト 1 0 1 に通知して、本処理を終了する。

## 【 0 1 0 4 】

これに対し、キャッシュヒットしなかった場合 ( S 1 1 0 2 : N )、ホスト I / O 処理プログラム 6 0 1 は、リード対象の仮想ボリューム 2 0 2 のアドレスに対応したデータを格納するための領域をキャッシュメモリ 1 0 8 上に確保する ( S 1 1 0 3 )。次いでホスト I / O 処理プログラム 6 0 1 は、ダイナミックマッピングテーブル 5 0 6 を参照して、リード対象の仮想ボリューム 2 0 2 のアドレスにプール 2 0 4 から実ページが割り当てられているか否かを確認する。

10

## 【 0 1 0 5 】

仮想ボリューム 2 0 2 に実ページが割り当てられていない場合、ホスト I / O 処理プログラム 6 0 1 は、ダイナミックマッピングテーブル 5 0 6 を参照して、デフォルト値を格納するページを用いてデフォルト値の格納ページのドライブ 1 1 1 又は 1 1 2 のアドレスを算出し、ドライブ 1 1 1 又は 1 1 2 からデフォルト値をキャッシュメモリ 1 0 8 の領域に転送する ( S 1 1 0 5 )。

## 【 0 1 0 6 】

ここで、デフォルト値の場合は、ダイナミックマッピングテーブル 5 0 6 の仮想ボリューム 2 0 2 及び論理アドレスに対応する、プールボリューム番号及び論理アドレスには、デフォルト値格納ページのあるプールボリューム番号及び論理アドレスが設定されている。デフォルト値格納ページは、プール 2 0 4 に 1 つ以上あればよい。容量効率を考えればプール 2 0 4 にデフォルト値ページは 1 又は 2 つである。

20

## 【 0 1 0 7 】

デフォルト値格納ページのアドレスと対応付けられている仮想ボリューム 2 0 2 の論理アドレスは、ホスト 1 0 1 から新規にデータの書き込みがあった際に、ホスト 1 0 1 のデータ書き込み用のページで未だどの仮想ボリューム 2 0 2 のアドレスにも対応付けられていない未使用のページと対応付け直される。

30

## 【 0 1 0 8 】

これに対し、仮想ボリューム 2 0 2 に実ページが割り当てられている場合、ホスト I / O 処理プログラム 6 0 1 は、ダイナミックマッピングテーブル 5 0 6 を参照して、プールボリューム番号及び論理アドレスを取得し、更に論理物理アドレス変換テーブル 5 0 7 を参照して、物理ドライブ番号及び物理開始アドレスを算出することで、リード対象の仮想ボリューム 2 0 2 のアドレスに対応したデータが格納されているドライブ 1 1 1 又は 1 1 2 のアドレスを算出する ( S 1 1 0 4 )。

## 【 0 1 0 9 】

次にホスト I / O 処理プログラム 6 0 1 は、算出したアドレスからデータをキャッシュメモリ 1 0 8 上の領域に転送する ( S 1 1 0 5 )。そしてリード時にページ毎モニターテーブル 5 0 1 を参照して、リード I / O カウンタの数値をカウントアップする ( S 1 1 0 6 )。

40

## 【 0 1 1 0 】

そしてホスト I / O 処理プログラム 6 0 1 は、ドライブ 1 1 1 又は 1 1 2 からキャッシュメモリ 1 0 8 上に格納したデータをホスト 1 0 1 に転送するとともに ( S 1 1 0 7 )、I / O 処理 ( リード処理 ) が完了したことをホスト 1 0 1 に通知して、本処理を終了する。

## 【 0 1 1 1 】

以上の処理により、ホスト I / O 処理を行うとともに、必要なモニタ情報を採取することができる。

50

## 【 0 1 1 2 】

図 1 2 は、デステージ処理のフローチャートを示す。このデステージ処理は、ホスト I / O 処理 ( 図 1 1 ) とは非同期で、プロセッサ 1 1 4 とデステージ処理プログラム 6 0 2 との協働により適宜実行される。説明の便宜上、処理主体をデステージ処理プログラム 6 0 2 として説明する。

## 【 0 1 1 3 】

まずデステージ処理プログラム 6 0 2 は、キャッシュ管理テーブル 5 0 5 を参照して、ダーティフラグの ON 又は OFF を確認し、ドライブ 1 1 1 又は 1 1 2 に書き込みがなされていない未反映データがキャッシュメモリ 1 0 8 上にあるか否かを判断する ( S 1 2 0 1 ) 。

10

## 【 0 1 1 4 】

デステージ処理プログラム 6 0 2 は、未反映データがキャッシュメモリ 1 0 8 上にある場合、キャッシュ管理テーブル 5 0 5 から仮想ボリューム番号及び論理アドレスを取得する。そしてこの仮想ボリューム番号及び論理アドレスを元にダイナミックマッピングテーブル 5 0 6 を参照して、プールボリューム番号及び論理アドレスを取得する。

## 【 0 1 1 5 】

このときプールボリューム番号及び論理アドレスがデフォルト値格納ページのアドレスであった場合、デステージ処理プログラム 6 0 2 は、新規データを書き込むためにダイナミックマッピングテーブル 5 0 6 から新規空きページを割り当てる。そしてこの割り当てページのプールボリューム番号及び論理アドレスをダイナミックマッピングテーブル 5 0 6 の対応する仮想ボリューム番号及び論理アドレスに対応付けて格納する。

20

## 【 0 1 1 6 】

既にページが割り当てられている場合は、デフォルト値のプールボリューム番号及び論理アドレスとは異なるプールボリューム番号及び論理アドレスの値が仮想ボリュームの論理アドレスに対応付けて格納されている。デステージ処理プログラム 6 0 2 は、プールボリューム番号及び論理アドレスを取得した後、論理物理アドレス変換テーブル 5 0 7 を参照して、ドライブ 1 1 1 又は 1 1 2 のアドレスを算出する ( S 1 2 0 2 ) 。

## 【 0 1 1 7 】

次いでデステージ処理プログラム 6 0 2 は、算出したドライブ 1 1 1 又は 1 1 2 のアドレスに対して、キャッシュメモリ 1 0 8 上の未反映データを書き込む ( S 1 2 0 3 ) 。そしてダイナミックマッピングテーブル 5 0 6 のページ番号に対応するページ毎モニタテーブル 5 0 1 を参照して、ライト I / O カウンタの数値をカウントアップする ( S 1 2 0 4 ) 。

30

## 【 0 1 1 8 】

次いでデステージ処理プログラム 6 0 2 は、ページ毎モニタテーブル 5 0 1 の新規ライトフラグ欄 7 0 5 を参照して、デステージ処理対象のページが新規割り当てページであるか否かを判断する ( S 1 2 0 5 ) 。新規割り当てページである場合 ( S 1 2 0 5 : Y ) 、デステージ処理プログラム 6 0 2 は、このページの格納先のパリティグループに対応するパリティグループ毎モニタテーブル 5 0 2 を参照して、新規ライト I / O カウンタの数値をカウントアップする ( S 1 2 0 6 ) 。

40

## 【 0 1 1 9 】

これに対し、デステージ処理プログラム 6 0 2 は、ページ処理対象のページが新規割り当てページでない場合 ( S 1 2 0 5 : N ) 、ステップ S 1 2 0 1 に移行する。ステップ S 1 2 0 1 においてデステージ処理プログラム 6 0 2 は、キャッシュメモリ 1 0 8 上に未反映データがさらにあるか否かを判断する。そしてさらなる未反映データがない場合には ( S 1 2 0 1 : N ) 、本処理を終了する。

## 【 0 1 2 0 】

以上の処理により、キャッシュメモリ 1 0 8 上のデータを非同期にドライブ 1 1 1 又は 1 1 2 に格納するとともに、必要なモニタ情報を採取することができる。

## 【 0 1 2 1 】

50

図13は、寿命情報採取処理のフローチャートである。この寿命情報採取処理は、プロセッサ114と寿命情報採取処理プログラム603との協働により一定周期で実行される。説明の便宜上、処理主体を寿命情報採取処理プログラム603として説明する。

【0122】

まず寿命情報採取処理プログラム603は、ドライブ111又は112に対して寿命採取のコマンドを発行する(S1301)。次いで寿命情報採取処理プログラム603は、寿命情報としてライト追加可能率又はライト削減要求率を受信する(S1302)。そして受信した寿命情報をローカルメモリ113に格納して(S1303)、本処理を終了する。

【0123】

以上の処理により、ドライブ111又は112から寿命情報を採取することができる。

【0124】

図14は、閾値決定処理のフローチャートである。この閾値決定処理は、プロセッサ114と閾値決定処理プログラム604との協働により一定周期で実行される。一定周期の情報は、プール毎再配置管理テーブル504の寿命制御再配置周期欄1002に格納される。説明の便宜上、処理主体を閾値決定処理プログラム604として説明する。

【0125】

まず閾値決定処理プログラム604は、全てのパリティグループについて、ページ毎のモニタ情報の集計が完了したか否かを判断する(S1401)。すなわちパリティグループ毎モニタテーブル502の各欄に情報が格納されているか否かを判断する。

【0126】

集計が完了していない場合(S1401:N)、閾値決定処理プログラム604は、寿命情報採取処理プログラム603を呼び出してドライブ111又は112から寿命情報を採取し(S1402)、ページ毎のモニタ情報をパリティグループ毎に集計する(S1403)。集計が完了している場合(S1401:Y)、閾値決定処理プログラム604は、ページ再配置のための各種閾値を算出する(S1404)。

【0127】

ここで算出される各種閾値は、パリティグループ毎モニタテーブル502、パリティグループ毎再配置管理テーブル503及びプール毎再配置管理テーブル504の各欄に格納される。そして閾値決定処理プログラム604は、再配置処理プログラム604Aを呼び出して再配置処理を実行した後(S1405)、本処理を終了する。

【0128】

以上の処理により、パリティグループ毎にモニタ情報を集計し、集計したモニタ情報に基づいて、各種閾値を算出することができる。そして算出した閾値を用いて再配置処理を実行することができる。

【0129】

図15は、再配置処理のフローチャートである。この再配置処理は、プロセッサ114と閾値決定処理プログラム604により呼び出される再配置処理プログラム604Aとの協働により実行される。説明の便宜上、処理主体を再配置処理プログラム604Aとして説明する。

【0130】

まず再配置処理プログラム604Aは、閾値決定処理プログラム604により算出されたTLC-MLC間ライト閾値に基づいて、TLC PGとMLC PGとの間でページを再配置するTLC-MLC間再配置処理を実行する(S1501)。

【0131】

このTLC-MLC間再配置処理により、TLC SSD(ドライブ112)の寿命を延ばすことができる。またビットコストを削減することができる。

【0132】

次に再配置処理プログラム604Aは、パリティグループ毎再配置管理テーブル503を参照し、ライトリバランス用の移動計画ページ数に基づいて、同一特性のSSDで構成

10

20

30

40

50

されたパリティグループ間（TLC PG間又はMLC PG間）でページを再配置するライトリバランス処理を実行する（S1502）。

【0133】

このライトリバランス処理により、同一特性のSSDで構成されたパリティグループ間でライト負荷を分散して、寿命を平準化することができる。

【0134】

次に再配置処理プログラム604Aは、パリティグループ毎再配置管理テーブル503を参照して、ライト削減PGのライトリバランス用の移動計画ページ数が0以下であるか否かを判断する（S1503）。ライト削減PGのライトリバランス用の移動計画ページ数が0以下である場合（S1503：Y）、再配置処理プログラム604Aは、ライト頻度の調整が完了していると判断して、次に性能リバランス処理を実行する（S1504）。

10

【0135】

これに対し、ライト削減PGのライトリバランス用の移動計画ページ数が0よりも大きい場合（S1503：N）、再配置処理プログラム604Aは、現在のパリティグループの構成では、ライト頻度を調整できず、寿命を維持することができないため、警告画面を表示して、寿命を保証するために追加すべきTLC SSD又はMLC SSDの容量をユーザに通知する（S1505）。

【0136】

以上の処理により、パリティグループ間でのライト頻度及びリード頻度を調整することができる。またライト頻度及びリード頻度を調整することができない場合には寿命を維持するために必要なTLC SSD又はMLC SSDの容量をユーザに通知することができる。

20

【0137】

図16は、TLC-MLC間再配置処理のフローチャートである。このTLC-MLC間再配置処理は、プロセッサ114と再配置処理プログラム604Aとの協働により実行される。説明の便宜上、処理主体を再配置処理プログラム604Aとして説明する。

【0138】

まず再配置処理プログラム604Aは、閾値決定処理プログラム604により算出されたTLC-MLC間ライト閾値に基づいて、全てのパリティグループ内の各ページをTLC PG又はMLC PGに再配置する。このとき再配置について未判定のページがあるか否かを判断する（S1601）。

30

【0139】

未判定のページがある場合（S1601：Y）、再配置処理プログラム604Aは、パリティグループ毎再配置管理テーブル503を参照して、この未判定の対象ページが所属するパリティグループのメディアタイプがTLCであるか否かを判断する（S1602）。メディアタイプがTLCである場合（S1602：Y）、再配置処理プログラム604Aは、対象ページのライト頻度がTLC-MLC間ライト閾値以上であるか否かを判断する（S1603）。

【0140】

対象ページのライト頻度がTLC-MLC間ライト閾値以上である場合（S1603：Y）、再配置処理プログラム604Aは、このページをTLC PGからMLC PGに移動する（S1604）。これに対し、対象ページのライト頻度がTLC-MLC間ライト閾値未満である場合（S1603：N）、再配置処理プログラム604Aは、何もせずにステップS1601に移行する。

40

【0141】

ステップS1602に戻り、メディアタイプがMLCである場合（S1602：N）、対象ページのライト頻度がTLC-MLC間ライト閾値未満であるか否かを判断する（S1605）。対象ページのライト頻度がTLC-MLC間ライト閾値未満である場合（S1605：Y）、再配置処理プログラム604Aは、このページをMLC PGからTL

50

C P Gに移動する(S 1 6 0 6)。

【0 1 4 2】

これに対し、対象ページのライト頻度がT L C - M L C間ライト閾値以上である場合(S 1 6 0 5 : N)、再配置処理プログラム6 0 4 Aは、何もせずにステップS 1 6 0 1に移行する。再配置処理プログラム6 0 4 Aは、全てのパリティグループ内の各ページについて判定を終えると、本処理を終了する。以上の処理により、T L C - M L C間でページを再配置して、T L C S S D (ドライブ1 1 2)の寿命を延ばすとともに、ビットコストを削減することができる。

【0 1 4 3】

図17は、ライトリバランス処理のフローチャートである。このライトリバランス処理は、プロセッサ1 1 4と再配置処理プログラム6 0 4 Aとの協働により実行される。説明の便宜上、処理主体を再配置処理プログラム6 0 4 Aとして説明する。

【0 1 4 4】

まず再配置処理プログラム6 0 4 Aは、閾値決定処理プログラム6 0 4により算出されたライトリバランス用移動計画ページ数、ライト削減閾値及びライト追加閾値に基づいて、同一特性のS S Dで構成されたパリティグループ間(T L C P G間又はM L C P G間)でページを再配置する。このとき再配置について未判定のページがあるか否かを判断する(S 1 7 0 1)。

【0 1 4 5】

未判定のページがない場合には(S 1 7 0 1 : N)、再配置処理プログラム6 0 4 Aは、本処理を終了する。これに対し、未判定のページがある場合(S 1 7 0 1 : Y)、再配置処理プログラム6 0 4 Aは、パリティグループ毎再配置管理テーブル5 0 3を参照して、この未判定の対象ページが所属するパリティグループの移動元P G種別がライト削減P Gであるか否かを判断する(S 1 7 0 2)。

【0 1 4 6】

対象ページが所属するパリティグループの移動元P G種別がライト削減P Gである場合(S 1 7 0 2 : Y)、再配置処理プログラム6 0 4 Aは、ページ毎モニタテーブル5 0 1のライトI/Oカウンタを参照して、この対象ページのライト頻度を取得する。そして取得したライト頻度がライト削減閾値以上であるか否かを判断する(S 1 7 0 3)。

【0 1 4 7】

対象ページのライト頻度がライト削減閾値未満である場合(S 1 7 0 3 : N)、再配置処理プログラム6 0 4 Aは、ステップS 1 7 0 1に移行する。これに対し、対象ページのライト頻度がライト削減閾値以上である場合(S 1 7 0 3 : Y)、再配置処理プログラム6 0 4 Aは、この対象ページの移動先のパリティグループを決定する。

【0 1 4 8】

移動先のパリティグループを決定する際、再配置処理プログラム6 0 4 Aはパリティグループ毎再配置管理テーブル5 0 3を参照して、ライト追加P Gのうち、ライトリバランス用の移動実績ページ数がライトリバランス用の移動計画ページ数未満であるパリティグループが存在するか否かを判断する(S 1 7 0 4)。

【0 1 4 9】

ライト追加P Gのうち、ライトリバランス用の移動実績ページ数がライトリバランス用の移動計画ページ数未満であるパリティグループが存在する場合(S 1 7 0 4 : Y)、再配置処理プログラム6 0 4 Aは、このパリティグループにライト負荷の高い対象ページを移動しても寿命を維持できると判断して、このパリティグループを移動先P Gとしてパリティグループ毎再配置管理テーブル5 0 3に登録し、対象ページをこのパリティグループに移動する(S 1 7 0 5)。

【0 1 5 0】

これに対し、ライト追加P Gのうち、ライトリバランス用の移動実績ページ数がライトリバランス用の移動計画ページ数未満であるパリティグループが存在しない場合(S 1 7 0 4 : N)、再配置処理プログラム6 0 4 Aは、この対象ページについては判定を終え、

10

20

30

40

50

ステップ S 1 7 0 1 に移行する。

【 0 1 5 1 】

ステップ S 1 7 0 2 に戻り、対象ページが所属するパリティグループの移動元 P G 種別がライト削減 P G でない場合 ( S 1 7 0 2 : N )、すなわち対象ページが所属するパリティグループの移動元 P G 種別がライト追加 P G である場合、再配置処理プログラム 6 0 4 A は、ページ毎モニタテーブル 5 0 1 のライト I / O カウンタ欄 7 0 2 を参照して、この対象ページのライト頻度を取得する。そして取得したライト頻度がライト追加閾値未満であるか否かを判断する ( S 1 7 0 6 )。

【 0 1 5 2 】

対象ページのライト頻度がライト追加閾値以上である場合には ( S 1 7 0 6 : N )、再配置処理プログラム 6 0 4 A は、ステップ S 1 7 0 1 に移行する。これに対し、対象ページのライト頻度がライト追加閾値未満である場合 ( S 1 7 0 6 : Y )、再配置処理プログラム 6 0 4 A は、この対象ページの移動先のパリティグループを決定する。

10

【 0 1 5 3 】

移動先のパリティグループを決定する際の処理は、上記のステップ S 1 7 0 4 及び S 1 7 0 5 と同様であるため説明を省略する。以上の処理により、同一特性の S S D で構成された異なるパリティグループ間でページを再配置してライト負荷を分散し、寿命を平準化することができる。

【 0 1 5 4 】

図 1 8 は、性能リバランス処理のフローチャートである。この性能リバランス処理は、プロセッサ 1 1 4 と再配置処理プログラム 6 0 4 A との協働により実行される。説明の便宜上、処理主体を再配置処理プログラム 6 0 4 A として説明する。

20

【 0 1 5 5 】

なお性能リバランス処理は、ライトリバランス処理 ( 図 1 7 ) で調整したライト頻度と同程度のリード頻度のページをライトリバランス処理で移動させた方向とは逆方向に移動する点で、ライトリバランス処理と異なり、他の基本的な処理内容は同様である。

【 0 1 5 6 】

まず再配置処理プログラム 6 0 4 A は、閾値決定処理プログラム 6 0 4 により算出された性能リバランス用の移動計画ページ数、ライト削減閾値及びライト追加閾値に基づいて、同一特性の S S D で構成されたパリティグループ間 ( T L C P G 間又は M L C P G 間 ) でページを再配置する。このとき再配置について未判定のページがあるか否かを判断する ( S 1 8 0 1 )。

30

【 0 1 5 7 】

未判定のページがない場合には ( S 1 8 0 1 : N )、再配置処理プログラム 6 0 4 A は、本処理を終了する。これに対し、未判定のページがある場合 ( S 1 8 0 1 : Y )、再配置処理プログラム 6 0 4 A は、パリティグループ毎再配置管理テーブル 5 0 3 を参照して、この未判定の対象ページが所属するパリティグループの移動元 P G 種別がライト削減 P G であるか否かを判断する ( S 1 8 0 2 )。

【 0 1 5 8 】

対象ページが所属するパリティグループの移動元 P G 種別がライト削減 P G である場合 ( S 1 8 0 2 : Y )、再配置処理プログラム 6 0 4 A は、ページ毎モニタテーブル 5 0 1 のリード I / O カウンタを参照して、この対象ページのリード頻度を取得する。そして取得したリード頻度がリード追加閾値未満であるか否かを判断する ( S 1 8 0 3 )。

40

【 0 1 5 9 】

対象ページのリード頻度がリード追加閾値以上である場合には ( S 1 8 0 3 : N )、再配置処理プログラム 6 0 4 A は、ステップ S 1 8 0 1 に移行する。これに対し、対象ページのリード頻度がリード追加閾値未満である場合 ( S 1 8 0 3 : Y )、再配置処理プログラム 6 0 4 A は、この対象ページの移動先のパリティグループを決定する。

【 0 1 6 0 】

移動先のパリティグループを決定する際、再配置処理プログラム 6 0 4 A はパリティグ

50

ループ毎再配置管理テーブル503を参照して、ライト追加PGのうち、性能リバランス用の移動実績ページ数が性能リバランス用の移動計画ページ数未満であるパリティグループが存在するか否かを判断する(S1804)。

【0161】

ライト追加PGのうち、性能リバランス用の移動実績ページ数が性能リバランス用の移動計画ページ数未満であるパリティグループが存在する場合(S1804:Y)、再配置処理プログラム604Aは、このパリティグループにリード負荷の低い対象ページを移動させても高負荷にならないと判断して、このパリティグループを移動先PGとしてパリティグループ毎再配置管理テーブル503に登録し、対象ページをこのパリティグループに移動する(S1805)。

10

【0162】

これに対し、ライト追加PGのうち、性能リバランス用の移動実績ページ数が性能リバランス用の移動計画ページ数未満であるパリティグループが存在しない場合(S1804:N)、再配置処理プログラム604Aは、この対象ページについては判定を終え、ステップS1801に移行する。

【0163】

ステップS1802に戻り、対象ページが所属するパリティグループの移動元PG種別がライト削減PGでない場合(S1802:N)、すなわち対象ページが所属するパリティグループの移動元PG種別がライト追加PGである場合、再配置処理プログラム604Aは、ページ毎モニタテーブル501のリードI/Oカウンタ欄703を参照して、この対象ページのリード頻度を取得する。そして取得したリード頻度がリード削減閾値以上であるか否かを判断する(S1806)。

20

【0164】

対象ページのリード頻度がリード削減閾値未満である場合には(S1806:N)、再配置処理プログラム604Aは、ステップS1801に移行する。これに対し、対象ページのリード頻度がリード削減閾値以上である場合(S1806:Y)、再配置処理プログラム604Aは、この対象ページの移動先のパリティグループを決定する。

【0165】

移動先のパリティグループを決定する際の処理は、上記のステップS1804及びS1805と同様であるため説明を省略する。以上の処理により、同一特性のSSDで構成された異なるパリティグループ間でページを再配置してリード負荷を分散し、I/O頻度を平準化することができる。

30

【0166】

図19は、新規割り当て決定処理のフローチャートを示す。この新規割り当て決定処理は、ホスト101から新規の仮想ページに対するライト要求を受領したことを契機として、プロセッサ114と新規割り当て決定処理プログラム605との協働により実行される。説明の便宜上、処理主体を新規割り当て決定処理プログラム605として説明する。

【0167】

まず新規割り当て決定処理プログラム605は、プール毎再配置管理テーブル504を参照して、新規ページの割り当て対象の仮想ボリューム202に記憶領域を提供するプール204のワークロードタイプが「Unknown」であるか否かを判断する(S1901)。

40

【0168】

プール204のワークロードタイプが「Unknown」ではない場合(S1901:N)、新規割り当て決定処理プログラム605は、更に「Write intensive」であるか否かを判断する(S1906)。ワークロードタイプが「Write intensive」である場合(S1906:Y)、新規割り当て決定処理プログラム605は、新規ページに対するライト頻度は多いと予測して、書き込み上限回数が比較的多いMLC PGを新規ページに対する実ページの割り当て先のパリティグループに設定する(S1907)。

50

## 【0169】

これに対し、ワークロードタイプが「Read intensive」である場合（S1906：N）、新規割り当て決定処理プログラム605は、新規ページに対するライト頻度は少ないものと予測して、書き込み上限回数が少ないTLC PGを新規ページに対する実ページの割り当て先のパリティグループに設定する（S1904）。

## 【0170】

ステップS1901に戻り、ワークロードタイプが「Unknown」である場合（S1901：Y）、新規割り当て決定処理プログラム605は、ホスト101からのI/O特性が分からないため、新規ページに対する将来のライト頻度を予測して、割り当て先を決定する。まずは新規ページの予測ライト頻度を算出する（S1902）。

10

## 【0171】

例えば新規割り当て決定処理プログラム605は、ホスト101からの1ページ当たりの平均ライト頻度をモニタ情報として採取して見積もることにより、予測ライト頻度を算出する。

## 【0172】

次いで新規割り当て決定処理プログラム605は、プール毎再配置管理テーブル504を参照して、予測ライト頻度が新規ライト閾値未満であるか否かを判断する（S1903）。予測ライト頻度が新規ライト閾値未満の場合（S1903：Y）、新規割り当て決定処理プログラム605は、書き込み上限回数の少ないTLC PGを割り当て先のパリティグループに設定する（S1904）。

20

## 【0173】

これに対し、予測ライト頻度が新規ライト閾値以上の場合（S1903：N）、新規割り当て決定処理プログラム605は、書き込み上限回数の多いMLC PGを割り当て先のパリティグループに設定する（S1907）。

## 【0174】

次いで新規割り当て決定処理プログラム605は、割り当て先に設定した特性のパリティグループについて、新規ページの割り当てが可能か否かを判断するため、パリティグループ毎モニタテーブル502及びパリティグループ毎再配置管理テーブル503を参照して、新規ライト可能量が新規ライトI/Oカウンタよりも大きいかなんかを割り当て先に設定した特性のパリティグループごとに判断する（S1905）。

30

## 【0175】

新規ライト可能量が新規ライトI/Oカウンタよりも大きいパリティグループが存在する場合（S1905：Y）、新規割り当て決定処理プログラム605は、このパリティグループから新規ページを割り当てるために、ページ毎モニタテーブル501を参照して、このパリティグループにおける何れかのページについて新規ライトフラグを設定し（S1909）、本処理を終了する。

## 【0176】

これに対し、新規ライト可能量が新規ライトI/Oカウンタよりも大きいパリティグループが存在しない場合（S1905：N）、新規割り当て決定処理プログラム605は、寿命を維持するにあたり新規に許容可能なライト頻度を超過しているため、ユーザに推奨容量を通知する（S1908）。

40

## 【0177】

そして新規割り当て決定処理プログラム605は、判断対象のパリティグループから新規ページを割り当てるために、何れかのパリティグループにおける何れかのページについて新規ライトフラグを設定し（S1909）、本処理を終了する。

## 【0178】

以上の処理により、ホスト101から新規ページにデータを書き込むライト要求を受領した場合、新規ページに対するライト頻度を考慮して、寿命特性に応じたSSDにより構成されたパリティグループから新規ページに対して実ページを割り当てることができる。

## 【0179】

50

( 1 - 7 ) 推奨容量算出方法の概念構成

図 20 は、ホスト 101 からのライト頻度に対して寿命を保証するにあたり推奨される書き込み上限回数異なる記憶装置の容量比率の算出方法の考え方を示す。図 20 は、プール 204 内の各ページ 207 ( 又は仮想ボリューム 202 内の各ページ 201 ) のライト頻度の分布を表す。グラフ 2007 は、左からライト頻度が多い順番に全ページ 207 を並べたときの各ページ 207 のライト頻度を示す。縦軸はライト頻度であり、横軸はページ数である。

【 0180 】

TLC - MLC 間ライト閾値 2005 は、TLC PG と MLC PG のどちらにページを配置するのかが決める閾値であり、TLC - MLC 間ライト閾値 2005 とグラフ 2007 との交点が TLC と MLC の推奨される容量比率 ( 2006 ) となる。TLC - MLC 間ライト閾値 2005 は、ユーザが指定してもよいし、ストレージシステム 104 が算出してもよい。

10

【 0181 】

ストレージシステム 104 が算出する場合、推奨容量比率は、顧客要件のライト頻度 ( = Whost ) 、プール 204 ( 又は仮想ボリューム 202 ) 容量 ( = C ) 及び各 SSD の寿命を保証するにあたり許容できるライト頻度の相関グラフ 2001、2002 を利用し、下記式 9 を計算して算出することができる。

【 0182 】

ここでグラフ 2001 は、寿命を保証するにあたり単位容量あたりに許容できるライト頻度 ( = Wtlc ) を傾きとした TLC SSD の容量と許容できるライト頻度との相関を表す。またグラフ 2002 は、寿命を保証するにあたり単位容量あたりに許容できるライト頻度 ( = Wmlc ) を傾きとした MLC SSD の容量と許容できるライト頻度との相関を表す。

20

【 0183 】

【 数 9 】

$$\begin{aligned}
 & \text{TLC SSD 推奨容量} \\
 & = ( ( ( Wmlc \times C ) - Whost ) \div ( Wmlc - Wtlc ) ) \div C \\
 & \text{MLC SSD 推奨容量} \\
 & = ( ( C - \text{TLC 推奨容量} ) \div C ) \\
 & \dots\dots\dots (9)
 \end{aligned}$$

30

【 0184 】

以上により、顧客要件ライト頻度とプール容量とを満たす TLC SSD と MLC SSD との容量比率を算出することができる。

【 0185 】

( 1 - 8 ) 画面構成

図 21 は、プール毎にパラメータを設定する際の画面構成の一例を示す。プール単位の GUI 画面 2101 は、設定対象のプール 204 を特定できるプール番号を表示する領域 2102 と、寿命制御再配置の ON / OFF を設定する領域 2103 と、寿命制御再配置を ON にした場合の詳細設定の ON / OFF を設定する領域 2104 と、詳細設定の内容を設定する領域 2105 とから構成される。本画面で設定した情報は、プール毎再配置管理テーブル 504 の各欄に格納される。

40

【 0186 】

寿命制御再配置領域 2103 の設定が OFF の場合、閾値決定処理プログラム 604 が閾値決定処理を実行することはなく、よって再配置処理プログラム 604A がページの再配置を実行することはないが、寿命制御の精度の低下を防ぐため、モニタ情報は寿命制御再配置の ON / OFF にかかわらず採取される。

【 0187 】

寿命制御再配置領域 2103 の設定が ON の場合、上記説明してきた通りページの再配

50

置が行われる。この場合、詳細設定領域 2 1 0 4 の項目設定領域 2 1 0 5 が入力可能に表示される。詳細設定領域 2 1 0 4 の設定が OFF の場合、項目設定領域 2 1 0 5 の各種パラメータには、デフォルト値又はストレージシステム 1 0 4 内で自動的に算出された値が設定される。

**【 0 1 8 8 】**

再配置周期領域 2 1 0 6 には、寿命制御のための再配置を実行する周期が設定される。この周期はユーザが指定することができる。例えばユーザが「7 days」と指定すると、7日周期でページの再配置が実行される。

**【 0 1 8 9 】**

ワークロードタイプ領域 2 1 0 7 には、ホスト 1 0 1 からの I / O 特性が設定される。この I / O 特性はユーザが指定することができる。ホスト 1 0 1 からの I / O 特性を予め把握している場合は、ユーザが I / O 特性を指定することにより、新規割り当て先の SSD の特性を明示的に指定することができる。

10

**【 0 1 9 0 】**

具体的には「Write intensive」が指定された場合、ホスト 1 0 1 の I / O 特性は、ライト高負荷であるため、書き込み上限回数の比較的多い MLC PG から新規ページに対して実ページの割り当てが行われる。

**【 0 1 9 1 】**

また「Read intensive」が指定された場合、ホスト 1 0 1 の I / O 特性は、ライト低負荷であるため、書き込み上限回数の少ない TLC PG から新規ページに対して実ページの割り当てが行われる。

20

**【 0 1 9 2 】**

またホスト 1 0 1 の I / O 特性をユーザが把握していない場合、ユーザは「Unknown」を指定する。この場合、ストレージシステム 1 0 4 が新規ページに対して実ページの割り当て先のパリティグループを自動的に決定することになる。

**【 0 1 9 3 】**

同種ドライブ間新規割り当てポリシー領域 2 1 0 8 には、新規ページに対する実ページの割り当て時に TLC PG 又は MLC PG のうち、何れの特性のパリティグループからページを割り当てるかを決定するポリシーが設定される。

**【 0 1 9 4 】**

例えばユーザにより「ラウンドロビン」が指定された場合、各パリティグループから均等にページが割り当てられる。また「容量優先」が指定された場合、容量が少なパリティグループから優先してページが割り当てられる。また「寿命優先」が指定された場合、寿命が長いパリティグループから優先してページが割り当てられる。

30

**【 0 1 9 5 】**

バッファサイズ領域 2 1 0 9 及び 2 1 1 0 には、特性が MLC PG であるパリティグループの容量に対するバッファの割合が設定される。新規割り当て用領域 2 1 0 9 には、特性が MLC PG であるパリティグループから新規ページを割り当てる際に使用するバッファが設定される。新規割り当て時に MLC PG から割り当てられるはずのライト高負荷なページが MLC PG の残容量が足りないことによって TLC PG から割り当てられることを防ぐ効果がある。

40

**【 0 1 9 6 】**

新規割り当てバッファは、再配置周期毎に確保しなおされる。このため、新規割り当てバッファは、周期内に予想されるホスト 1 0 1 から新規ページに対して書き込まれるデータ量に基づいて最適なサイズを見積もることができる。

**【 0 1 9 7 】**

再配置バッファ用領域 2 1 1 0 には、ページの再配置時に使用されるバッファが設定される。再配置バッファにより、再配置時に単位時間当たりに移動できるデータサイズを調整する。このため、再配置バッファを多くとることにより、再配置時のスループットを増やす効果がある。

50

## 【 0 1 9 8 】

図 2 2 は、ストレージシステム 1 0 4 がユーザに通知する警告画面の画面構成の一例を示す。警告画面は、ホスト 1 0 1 からのライト頻度が多く、現在の構成では目標期間よりも S S D の寿命が短くなること又は寿命を保証するにあたり過剰に S S D を搭載していることをユーザに通知することができる。

## 【 0 1 9 9 】

画面を表示する契機は、周期的に実行されるページ再配置処理の完了後にストレージシステム 1 0 4 が自動で表示するとしてもよいし、ユーザがストレージシステム 1 0 4 に対する操作により、任意のタイミングで表示させてもよい。後者の場合、任意の操作が行われたタイミングでのモニタ情報に基づいて推奨容量が算出される。

10

## 【 0 2 0 0 】

プール単位の G U I 画面 2 2 0 1 は、設定対象のプール 2 0 4 を特定できるプール番号を表示する領域 2 2 0 2 と、警告の内容を通知する領域 2 2 0 3 から構成される。警告の内容を通知する領域 2 2 0 3 は、現在のドライブ構成から追加が必要又は削減が可能な T L C 又は M L C の容量を通知する領域 2 2 0 4 及び現在のホスト 1 0 1 からの I / O 要求の情報から推奨される M L C と T L C の容量を通知する領域 2 2 0 5 から構成される。なお通知する容量の情報は、T L C と M L C の容量の比率で表してもよい。

## 【 0 2 0 1 】

( 1 - 9 ) 第 1 の実施の形態による効果

以上のように第 1 の実施の形態におけるストレージシステム 1 0 4 によれば、T L C - M L C 間でページを再配置することにより、寿命の異なる S S D から構成されたパリティグループ間でライト頻度を調整することができる。またライトリバランス処理を行うことにより、同一特性の S S D で構成されたパリティグループ間でライト頻度を調整することができる。よって寿命劣化の激しい S S D に対するライト頻度のページを寿命劣化の緩やかな S S D に移動して、S S D の保守交換回数を削減することができる。またストレージシステム 1 0 4 のコストを削減することができる。

20

## 【 0 2 0 2 】

( 2 ) 第 2 の実施の形態

第 2 の実施の形態は、記憶装置として半導体メモリ ( S S D ) だけでなく、ハードディスクドライブ ( H D D ) を搭載するストレージシステムにおいて、記憶装置 ( S S D 及び H D D ) を性能に応じた階層に分類し、ホストからのアクセス頻度に応じて適切な階層の記憶装置にデータを配置する階層制御を行う点で、第 1 の実施の形態と異なる。

30

## 【 0 2 0 3 】

( 2 - 1 ) 計算機システムの全体構成

図 2 3 は、第 2 の実施の形態における計算機システム 1 A の全体構成を示す。計算機システム 1 A は、S A S ( Serial Attached SCSI ) 規格の H D D ( ドライブ 2 3 0 1 ) を搭載している点及び S S D ( ドライブ 1 1 1 及び 1 1 2 ) をティア 1 に設定し、S A S 規格の H D D ( ドライブ 2 3 0 1 ) をティア 2 に設定するとともに、各ページの I / O 頻度に応じてティア間でデータを再配置する点で、第 1 の実施の形態と異なる。

## 【 0 2 0 4 】

例えばティア 1 からは I / O 頻度が 1 0 0 [ I O P S ] のページを割り当てるように設定されており、ティア 2 からは I / O 頻度が 1 0 [ I O P S ] のページを割り当てるように設定されているとする。一方で I / O 頻度が 5 0 [ I O P S ] のページがティア 2 から割り当てられており、I / O 頻度が 2 0 [ I O P S ] のページがティア 1 から割り当てられているとする。

40

## 【 0 2 0 5 】

この場合、ストレージシステム 1 0 4 A は全体として 1 0 (ティア 2 の上限 I O P H ) + 2 0 = 3 0 [ I O P S ] の性能しか発揮することができない。そこで I / O 頻度が 5 0 [ I O P H ] のページをティア 2 からティア 1 に移動 (再配置) すると、ストレージシステム 1 0 4 A は全体として 5 0 + 2 0 = 7 0 [ I O P S ] の性能を発揮することができる

50

ようになる。

【0206】

(2-2) ストレージシステムの論理構成

図24は、ストレージシステム104Aの論理構成を示す。ストレージシステム104Aは、プールボリューム206A~206CがMLC SSD又はTLC SSDから構成され、プールボリューム206D及び206EがSAS HDDから構成されており、各プールボリューム206A~206Eがティア1又は2に分類されている点で、第1の実施の形態と異なる。

【0207】

そして第2の実施の形態のストレージシステム104Aは、第1の実施の形態において説明してきたページの再配置に加えて、SSDから構成されるティア1においてホスト101からのライト負荷が高くなり、SSDの寿命を保証することができなくなった場合、ライト負荷の高い例えばページ207Dをティア2に移動するライトデモーション処理を実行する。これによりSSDの寿命劣化を防ぐことができる。

10

【0208】

(2-3) ページ配置処理の概念構成

図25は、ライトデモーション処理によるページ再配置処理の概念構成を示す。ここではMLC PG205A又はTLC PG205Bから構成されるティア1と、SAS HDDで構成されたパリティグループ(以下、SAS PGと呼ぶ)205Cから構成されるティア2との間で、ページ301を再配置する。ライトデモーション処理は、閾値決定処理時にティア1内でのページ再配置だけではSSDの寿命を維持することができないと判断された場合に実行される。

20

【0209】

具体的には、閾値決定処理においてライトリバランスによる移動計画ページ数を算出した後、ライト削減PGのライト削減要求量を満たすだけのページを移動することができなかった場合、移動することができなかった分のライト頻度に基づいて、ライトデモーション処理による移動計画ページ数を算出する。そしてライトリバランス実行時に、ライトデモーション処理による移動計画ページ数分だけMLC PG205A又はTLC PG205BからSAS PG205Cにライトデモーション閾値2501以上のライト頻度のページを移動する。

30

【0210】

(2-4) テーブル構成

図26は、第2の実施の形態におけるパリティグループ毎再配置管理テーブル503Aの論理構成を示す。パリティグループ毎再配置管理テーブル503Aは、パリティグループ毎再配置管理テーブル503(図9)の各欄に加えて、ライトデモーション閾値欄2601、ライトデモーション用移動計画ページ数欄2602、ライトデモーション用移動実績ページ数欄2603及びティアレベル欄2604から構成される。

【0211】

ライトデモーション閾値欄2601には、ライトデモーション処理対象のページを決定するための閾値が格納され、ライトデモーション用移動計画ページ数欄2602には、ライトデモーション処理により移動するページ数が格納される。またライトデモーション用移動実績ページ数欄2603には、ライトデモーション処理により移動したページ数が格納される。ティアレベル欄2604には、パリティグループが所属する階層順序(例えばティア1、ティア2、ティア3)が格納される。なおここでは階層順序の値が小さいティアほど高性能のドライブで構成されている。

40

【0212】

図27は、第2の実施の形態におけるプール毎再配置管理テーブル504Aの論理構成を示す。プール毎再配置管理テーブル504Aは、プール毎再配置管理テーブル504(図10)の各欄に加えて、ティア間I/O閾値欄2701から構成される。ティア間I/O閾値欄2701には、ページをどの階層に配置するのかを決定するための閾値が格納さ

50

れる。

【0213】

(2-5) フローチャート

図28は、第2の実施の形態における再配置処理のフローチャートを示す。この再配置処理は、プロセッサ114と閾値決定処理プログラム604により呼び出される再配置処理プログラム604Aとの協働により実行される。説明の便宜上、処理主体を再配置処理プログラム604Aとして説明する。

【0214】

なお前提として閾値決定処理プログラム604は、閾値決定処理(図14)においてティア間I/O閾値及びライトデモーション用移動計画ページ数を算出し、各種閾値をパリティグループ毎再配置管理テーブル503A及びプール毎再配置管理テーブル504Aにそれぞれ格納しているものとする。

10

【0215】

まず再配置処理プログラム604Aは、ティア間再配置処理を実行する(S2801)。ここでは再配置処理プログラム604Aは、プール毎再配置管理テーブル504Aのティア間I/O閾値欄2701を参照して、対象ページを配置するティアを決定する。

【0216】

その後再配置処理プログラム604Aは、第1の実施の形態における再配置処理(図15)と同様にTLC-MLC間再配置処理を実行し(S1501)、次いでライトリバランス処理を実行する(S1502)。ライトリバランス処理を実行した後、再配置処理プログラム604Aは、パリティグループ毎再配置管理テーブル503Aを参照して、ライトデモーション用移動計画ページ数が0であるか否かを判断する(S2802)。

20

【0217】

ライトデモーション用移動計画ページ数が0でない場合(S2802:N)、再配置処理プログラム604Aは、ライトリバランス処理を実行しただけではライト頻度を調整しきれなかったため、ライトデモーション処理を実行する(S2803)。その後再配置処理プログラム604Aは、性能リバランスを実行して(S1504)、本処理を終了する。

【0218】

これに対し、ライトデモーション用移動計画ページ数が0である場合(S2802:Y)、再配置処理プログラム604Aは、ライトデモーション処理を実行する必要はないため、性能リバランスを実行して(S1504)、本処理を終了する。以上の処理により、ホスト101からのライト頻度が過多であるためティア1を構成するSSDの寿命を維持することができない場合、ティア2にライト負荷が高いページを移動して、ティア1内のSSDの寿命を維持することができる。

30

【0219】

図29は、ティア間再配置処理のフローチャートを示す。このティア間再配置処理は、プロセッサ114と再配置処理プログラム604Aとの協働により実行される。説明の便宜上、処理主体を再配置処理プログラム604Aとして説明する。

【0220】

まず再配置処理プログラム604Aは、閾値決定処理プログラム604により算出されたティア間I/O閾値に基づいて、全てのパリティグループ内の各ページをティア1又は2に再配置する。このとき再配置について未判定のページがあるか否かを判断する(S2901)。

40

【0221】

未判定のページがある場合(S2901:Y)、再配置処理プログラム604Aは、パリティグループ毎再配置管理テーブル503Aを参照して、この未判定の対象ページが所属するティアがティア2であるか否かを判断する(S2902)。

【0222】

対象ページが所属するティアがティア2である場合(S2902:Y)、再配置処理プ

50

プログラム 604A は、プール毎再配置管理テーブル 504A を参照して、対象ページの I/O 頻度がティア間 I/O 閾値以上であるか否かを判断する (S2903)。

【0223】

対象ページの I/O 頻度がティア間 I/O 閾値以上である場合 (S2903:Y)、再配置処理プログラム 604A は、対象ページをティア 2 からティア 1 に移動する (S2904)。これに対し、対象ページの I/O 頻度がティア間 I/O 頻度未満である場合 (S2903:N)、再配置処理プログラム 604A は、何もせずにステップ S2901 に移行する。

【0224】

ステップ S2902 に戻り、対象ページが所属するティアがティア 1 の場合 (S2902:N)、再配置処理プログラム 604A は、プール毎再配置管理テーブル 504A を参照して、対象ページの I/O 頻度がティア間 I/O 閾値未満であるか否かを判断する (S2905)。

【0225】

対象ページの I/O 頻度がティア間 I/O 閾値未満である場合 (S2905:Y)、再配置処理プログラム 604A は、対象ページをティア 1 からティア 2 に移動する (S2906)。これに対し、対象ページの I/O 頻度がティア間 I/O 頻度以上である場合 (S2905:N)、再配置処理プログラム 604A は、何もせずにステップ S2901 に移行する。

【0226】

再配置処理プログラム 604A は、全てのパーティグループ内の各ページについて判定を終えると、本処理を終了する。以上の処理により、ホスト 101 からの I/O 頻度及び各ティアの性能に応じて、各ティアにページを再配置することができる。

【0227】

図 30 は、ライトデモーション処理のフローチャートを示す。このライトデモーション処理は、プロセッサ 114 と再配置処理プログラム 604A との協働により実行される。説明の便宜上、処理主体を再配置処理プログラム 604A として説明する。

【0228】

まず再配置処理プログラム 604A は、閾値決定処理プログラム 604 により算出されたライトデモーション用移動計画ページ数、ライト削減閾値及びライトデモーション閾値に基づいて、異なるティア間でページを再配置する。このとき再配置について未判定のページがあるか否かを判断する (S3001)。

【0229】

未判定のページがない場合 (S3001:N)、再配置処理プログラム 604A は本処理を終了する。これに対し未判定のページがある場合 (S3001:Y)、再配置処理プログラム 604A は、パーティグループ毎再配置管理テーブル 503A を参照して、この未判定の対象ページが所属するパーティグループの移動元 PG 種別がライト削減 PG であるか否かを判断する (S3002)。

【0230】

対象ページが所属するパーティグループの移動元 PG 種別がライト削減 PG でない場合 (S3002:N)、すなわちライト追加 PG である場合、再配置処理プログラム 604A は対象ページに対してライトデモーション処理を実行する必要はないと判断して、ステップ S3001 に移行する。

【0231】

これに対し、対象ページが所属するパーティグループの移動元 PG 種別がライト削減 PG である場合 (S3002:Y)、再配置処理プログラム 604A は、ページ毎モニターテーブル 501 のライト I/O カウンタ欄 702 を参照して、この対象ページのライト頻度を取得する。

【0232】

そして再配置処理プログラム 604A は、パーティグループ毎再配置管理テーブル 50

10

20

30

40

50

3 Aのライト削減閾値欄905及びライトデモーション閾値欄2601を参照して、ライト削減閾値及びライトデモーション閾値を取得する。

【0233】

そして取得した対象ページのライト頻度がライトデモーション閾値以上であり、かつ、ライト削減閾値未満であるか否かを判断する(S3003)。対象ページのライト頻度がライトデモーション閾値以上であり、かつ、ライト削減閾値未満でない場合(S3003:N)、再配置処理プログラム604AはステップS3001に移行する。

【0234】

これに対し、対象ページのライト頻度がライトデモーション閾値以上であり、かつ、ライト削減閾値未満である場合(S3003:Y)、再配置処理プログラム604Aは、対象ページが所属するパリティグループのライトデモーション用移動実績ページ数がライトデモーション用移動計画ページ数未満であるか否かを判断する(S3004)。

10

【0235】

対象ページが所属するパリティグループのライトデモーション用移動実績ページ数がライトデモーション用移動計画ページ数未満でない場合(S3004:N)、再配置処理プログラム604Aは、ステップS3001に移行する。

【0236】

これに対し、対象ページが所属するパリティグループのライトデモーション用移動実績ページ数がライトデモーション用移動計画ページ数未満である場合(S3004:Y)、再配置処理プログラム604Aは、対象ページをティア1からティア2に移動する(S3005)。再配置処理プログラム604Aは、未判定の全てのページについて判定を終えると、本処理を終了する。

20

【0237】

図31は、新規割り当て決定処理のフローチャートを示す。この新規割り当て決定処理は、ホスト101から新規ページに対するライト要求を受領したことを契機として、プロセッサ114と新規割り当て決定処理プログラム605との協働により実行される。説明の便宜上、処理主体を新規割り当て決定処理プログラム605として説明する。

【0238】

ステップS1901～S1905までは、第1の実施の形態における新規割り当て決定処理(図19)と同様であるため、ここでの説明は省略する。ステップS1905において、新規割り当て決定処理プログラム605は、割り当て先に設定した特性のティア1のパリティグループについて、新規ライト可能量が新規ライトI/Oカウンタよりも大きい  
か否かを判断する(S1905)。

30

【0239】

新規ライト可能量が新規ライトI/Oカウンタよりも大きいパリティグループが存在しない場合(S1905:N)、再配置処理プログラム604Aは、ティア1からこれ以上ページを割り当てると、SSDの寿命を維持することができなくなると判断して、新規ページに対する実ページの割り当て先をティア2に設定する(S3101)。

【0240】

ステップS1909は、第1の実施の形態における新規割り当て決定処理(図19)と同様であるため説明は省略する。以上の処理により、寿命維持のために必要な分のライト頻度をティア1から削減することができる。なお新規割り当て時にホスト101からのライト要求がシーケンシャルライトである場合、新規割り当て決定処理プログラムを実行せずに、ティア2からページを割り当てるようにしてもよい。

40

【0241】

(2-6)第2の実施の形態による効果

以上のように第2の実施の形態におけるストレージシステム104Aによれば、ホスト101からティア1に対するライト頻度をティア2に移動することにより、一定期間当たり許容できるライト頻度を越えたSSDの寿命劣化を緩やかにすることができる。よってSSDの保守交換回数を削減でき、ストレージシステム104Aのコストを削減するこ

50

とができる。

【0242】

(3) 第3の実施の形態

第3の実施の形態では、第1の実施の形態において説明したライト追加可能量及びライト削減要求量(図8)を算出する手法について説明する。第1の実施の形態においては、SSDからライト追加可能率及びライト削減要求率を採取し、これらを用いて上記式1及び2を計算することによりライト追加可能量及びライト削減要求量を算出するとしたが、ライト追加可能率及びライト削減要求率をSSDから直接採取することができない場合、以下の手法を用いてライト追加可能量及びライト削減要求量を算出することができる。

【0243】

図32は、SSDの寿命情報として磨耗指標(Wear out Indicator)を取得し、磨耗指標からライト追加可能量及びライト削減要求量を算出するための考え方を示す。グラフの縦軸は、SSDの磨耗指標(寿命率)3201であり、横軸は、SSDを使用しはじめからの経過年数(経過時間率)3202である。

【0244】

寿命率3201は、SSDの消去回数から算出される磨耗指標を意味しており、一般的にS.M.A.R.T.情報として取得できることが知られている。磨耗指標は、値が100(=L<sub>0</sub>)に達したとき、SSDの寿命を意味し、保守交換が必要となる。経過時間率は、目標寿命期間(例えば、3年や5年)を100%(=T<sub>0</sub>)とした経過時間の割合を意味する。

【0245】

直線3215は、区間 t1(=T<sub>1</sub>-T<sub>0</sub>)の寿命率の変化 L1(=L<sub>1</sub>-L<sub>0</sub>)を示しており、直線の傾きは、区間 t1のライト頻度(=W<sub>0</sub>)を表している。直線3215のライト頻度が継続すると、目標寿命期間を達成する前にSSDの寿命が尽きる。

【0246】

このため、直線3216で示すライト頻度(=W<sub>1</sub>)までライト頻度を削減し、区間 t2(=T<sub>0</sub>-T<sub>1</sub>)の間の寿命率の増加量が L(=L<sub>0</sub>-L<sub>1</sub>)となるように寿命変化の速度を調整する。また一般的にSSDに対するライト頻度は、WA(Write Amplification)により、ホスト101がSSDに対して発行したライト要求の回数よりもSSD内部のフラッシュチップに対するライト回数の方が大きくなることが知られている。SSDのWA値は、(フラッシュチップに対するライト回数÷ホスト101からのライト要求の回数)により算出される。

【0247】

以上の数値を用いてストレージシステム104は、寿命を維持するために必要なライト削減要求量を算出する。ライト削減要求量は、下記式10を計算して算出することができる。

【0248】

【数10】

$$\begin{aligned} & \text{ライト削減要求量 [IOPH]} \\ & = W_0 \times WA \times (1 - ((\Delta t_1 \div \Delta L_1) \times (\Delta L_2 \div \Delta t_2))) \\ & \dots \dots \dots (10) \end{aligned}$$

【0249】

またライト頻度が目標寿命期間を達成するにあたり余裕がある場合、ストレージシステム104は、ライト追加可能量を算出する。ライト追加可能量は、下記式11を計算して算出することができる。

【0250】

10

20

30

40

【数 1 1】

$$\begin{aligned} & \text{ライト追加可能量 [IOPH]} \\ & = W_0 \times WA \times \left( (\Delta t 1 \div \Delta L 1) \times (\Delta L 2 \div \Delta t 2) \right) - 1 \\ & \dots\dots\dots (11) \end{aligned}$$

【0 2 5 1】

以上により、一般的に取得可能なSSDの寿命情報からSSD PGに対するライト追加可能量及びライト削減要求量(図8)を算出することができる。そして算出したライト追加可能量及びライト削減要求量に基づいてページを再配置することにより、ストレージシステム104のコストを削減することができる。

10

【0 2 5 2】

なお上記説明してきた実施の形態においては、「書き込み」上限回数の少ないSSDのライト回数を削減するようにページを再配置することで、書き込み上限回数の少ないSSDの寿命を延ばす構成について説明してきたが、「書き換え」上限回数の少ないSSDについても上記構成を採用することにより、同様に寿命を延ばすことができる。SSDにおける書き換えとは、複数ページから構成されるブロックを一旦消去し、その後ブロック内の各ページにデータを書き込む一連の処理をいう。よって書き換えを1回行うと、消去と書き込みとを両方行うことになる。すなわち書き換え上限回数の少ないSSDと、書き込み上限回数の少ないSSDとを同様に扱うことで、書き換え上限回数の少ないSSDの寿命を延ばすことができる。

20

【0 2 5 3】

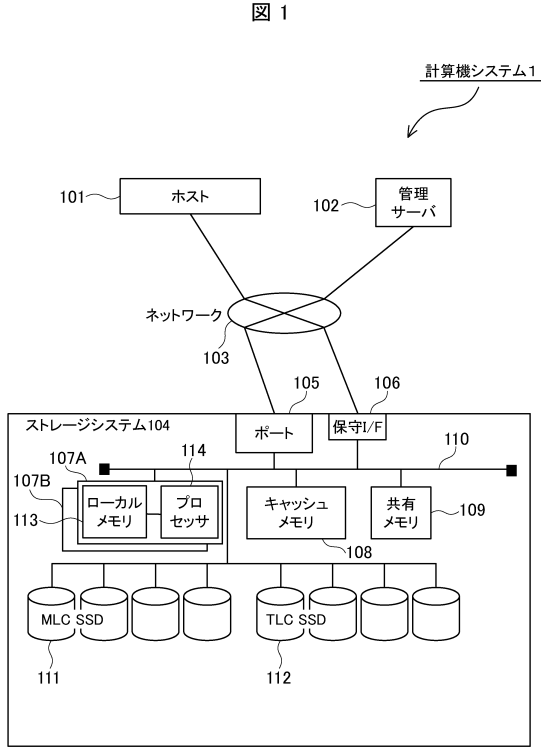
また上記説明してきた実施の形態においては、書き込み上限回数の少ないSSDの「ライト頻度」(ライト回数)を削減するようにページを再配置することで、書き込み上限回数の少ないSSDの寿命を延ばす構成について説明してきたが、ライトする「データ量」を削減するようにページを再配置することでも、同様に書き込み上限回数の少ないSSDの寿命を延ばすことができる。ライトするデータ量が大きい場合には複数ページにデータをライトする必要があり、よってライト回数も増加するためである。

【符号の説明】

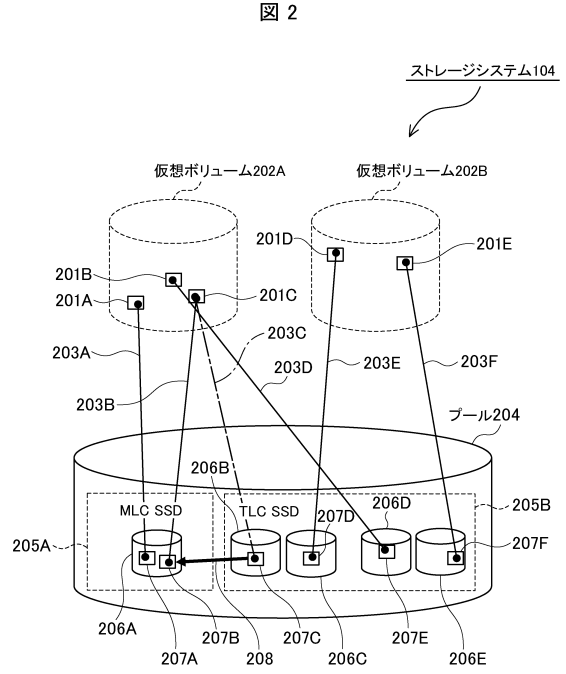
【0 2 5 4】

|                 |                 |    |
|-----------------|-----------------|----|
| 1、1A            | 計算機システム         | 30 |
| 104、104A        | ストレージシステム       |    |
| 114             | プロセッサ           |    |
| 111             | MLC SSD         |    |
| 112             | TLC SSD         |    |
| 2301            | SAS HDD         |    |
| 202             | 仮想ボリューム         |    |
| 201、301、401、207 | ページ             |    |
| 206             | プールボリューム        |    |
| 205             | パリティグループ        |    |
| 302             | TLC - MLC間ライト閾値 | 40 |

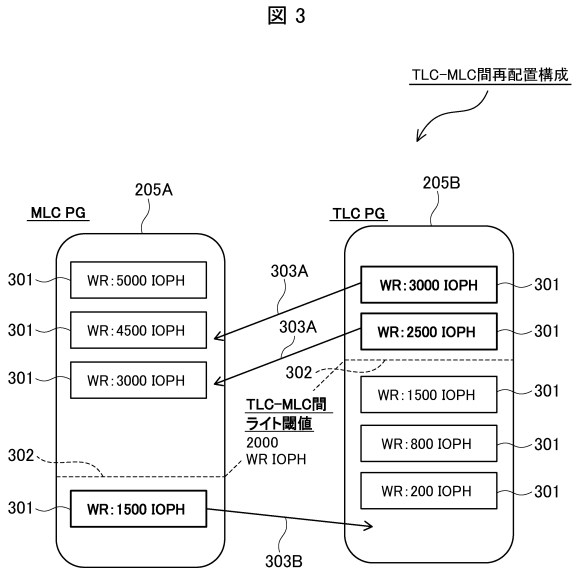
【図1】



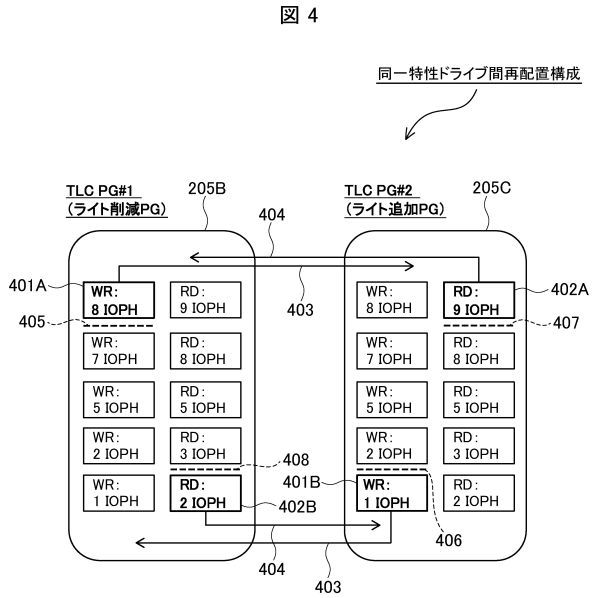
【図2】



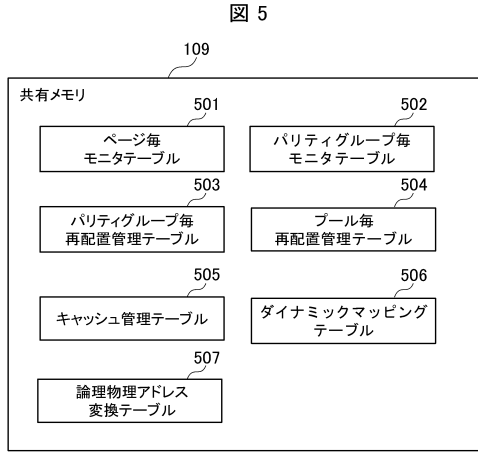
【図3】



【図4】



【図5】



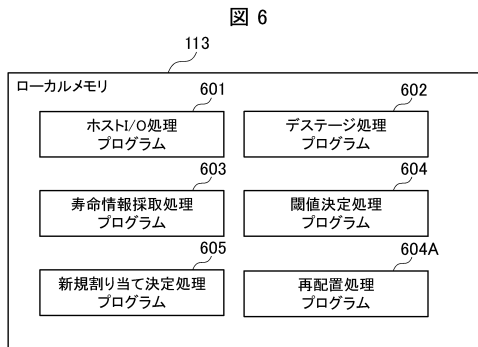
【図7】

図7

ページ毎モニターテーブル501

|     |            |      |       |      |     |
|-----|------------|------|-------|------|-----|
| 701 | ページ番号      | 0    | 1     | 2    | 3   |
| 702 | ライトI/Oカウンタ | 2000 | 12100 | 156  | 56  |
| 703 | リードI/Oカウンタ | 5010 | 6520  | 4600 | 300 |
| 704 | 合計I/Oカウンタ  | 7010 | 18620 | 4756 | 356 |
| 705 | 新規ライトフラグ   | 0    | 1     | 1    | 0   |

【図6】



【図8】

図8

パリティグループ毎モニターテーブル502

|     |              |       |       |       |
|-----|--------------|-------|-------|-------|
| 801 | パリティグループ番号   | 0     | 1     | 2     |
| 802 | 最大ライト頻度      | 22187 | 1962  | 4703  |
| 803 | 最小ライト頻度      | 5120  | 453   | 1085  |
| 804 | 最大リード頻度      | 4118  | 8018  | 17862 |
| 805 | 最小リード頻度      | 950   | 1850  | 4122  |
| 806 | リード/ライト比率    | 16/84 | 80/20 | 79/21 |
| 807 | ライト追加可能量     | 0     | 1790  | 2170  |
| 808 | ライト削減要求量     | 2447  | 0     | 0     |
| 809 | 新規ライトI/Oカウンタ | 896   | 202   | 69    |
| 810 | 新規ライト比率      | 16%   | 40%   | 6%    |
| 811 | 平均I/O頻度      | 8800  | 12000 | 9000  |
| 812 | 割り当てページ数     | 458   | 128   | 386   |

【図9】

図9

パリティグループ毎再配置管理テーブル503

|      |                     |         |         |         |
|------|---------------------|---------|---------|---------|
| 901  | パリティグループ番号          | 0       | 1       | 2       |
| 902  | メディアタイプ             | TLC     | TLC     | MLC     |
| 903  | 移動元PG種別             | ライト削減PG | ライト追加PG | ライト追加PG |
| 904  | 移動先PG               | 1       | 0       | 3       |
| 905  | ライト削減閾値             | 15531   | Invalid | Invalid |
| 906  | ライト追加閾値             | Invalid | 589     | 1411    |
| 907  | リード削減閾値             | Invalid | 5613    | 12503   |
| 908  | リード追加閾値             | 1236    | Invalid | Invalid |
| 909A | 移動計画ページ数 (ライトリバランス) | 156     | 70      | 1486    |
| 910A | 移動実績ページ数 (ライトリバランス) | 154     | 56      | 562     |
| 909B | 移動計画ページ数 (性能リバランス)  | 84      | 134     | 896     |
| 910B | 移動実績ページ数 (性能リバランス)  | 0       | 0       | 232     |
| 911  | 新規ライト可能量            | 846     | 1026    | 2394    |

【図10】

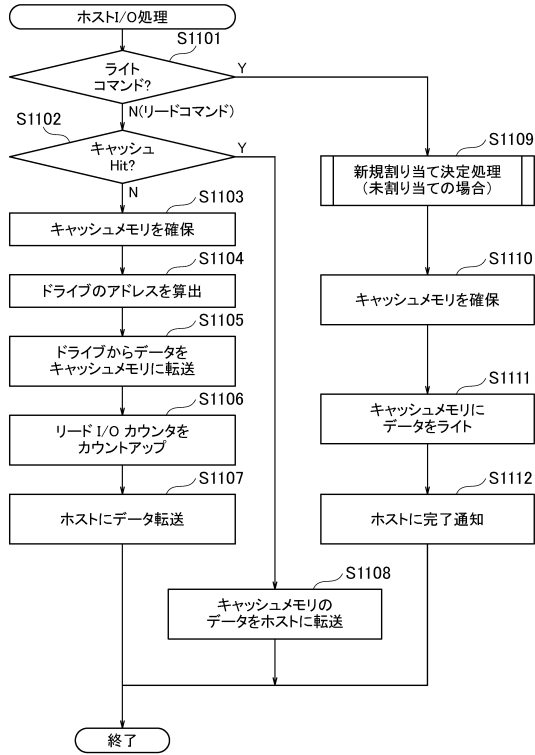
図10

プール毎再配置管理テーブル504

|      |                   |                 |         |                |
|------|-------------------|-----------------|---------|----------------|
| 1001 | プール番号             | 0               | 1       | 2              |
| 1002 | 寿命制御再配置周期 [hour]  | 168             | 24      | 72             |
| 1003 | TLC-MLC間ライト閾値     | 7040            | 9541    | 4530           |
| 1004 | 新規ライト閾値           | 7550            | 8205    | 4260           |
| 1005 | ワークロードタイプ         | Write intensive | Unknown | Read intensive |
| 1006 | 同種ドライブ間新規割り当てポリシー | 寿命優先            | 容量優先    | ラウンドロビン        |
| 1007 | 新規割り当てバッファ        | 3               | 3       | 5              |
| 1008 | 再配置バッファ           | 2               | 2       | 3              |

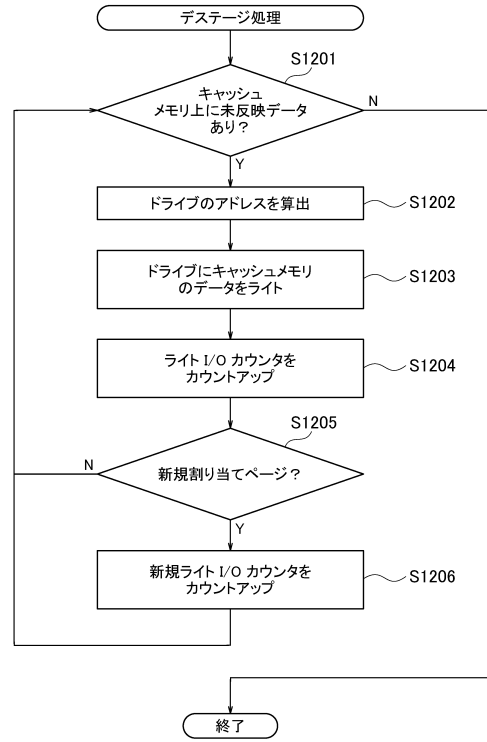
【図 1 1】

図 11



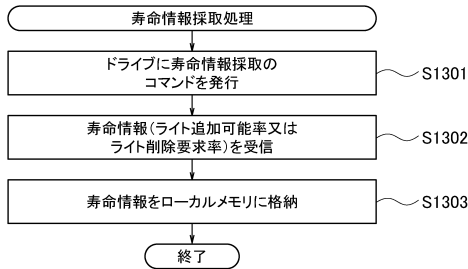
【図 1 2】

図 12



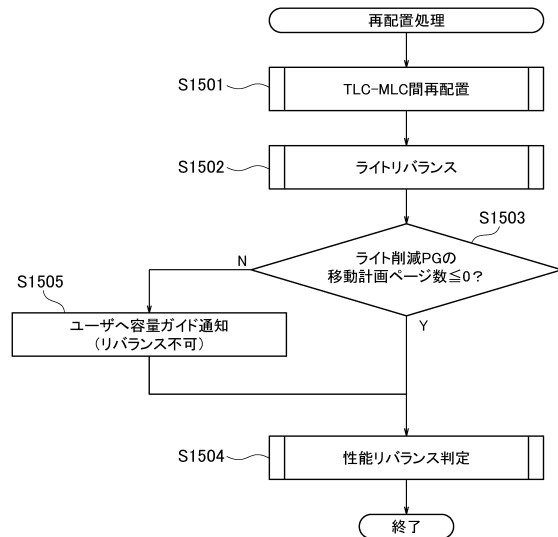
【図 1 3】

図 13



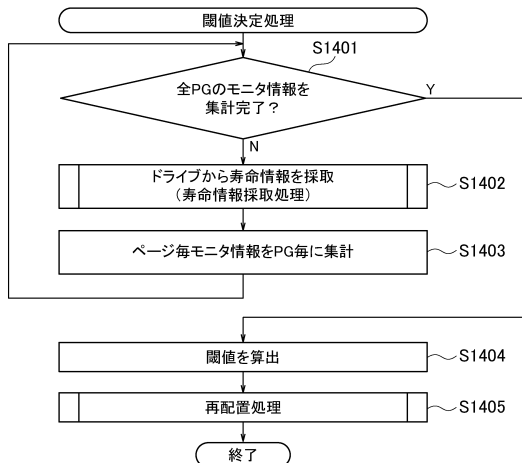
【図 1 5】

図 15



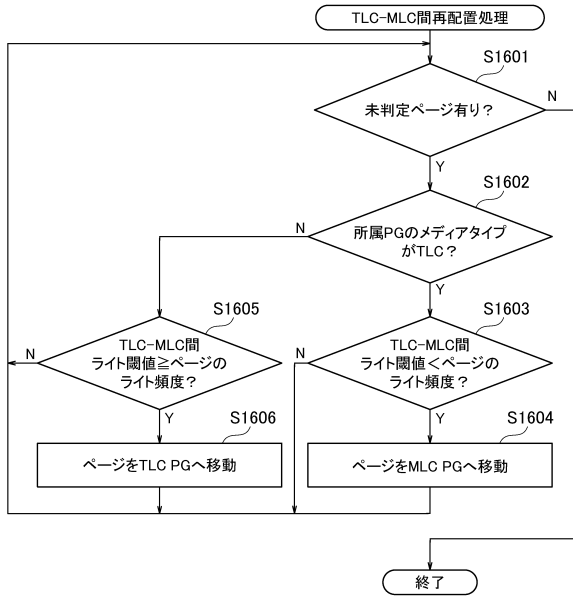
【図 1 4】

図 14



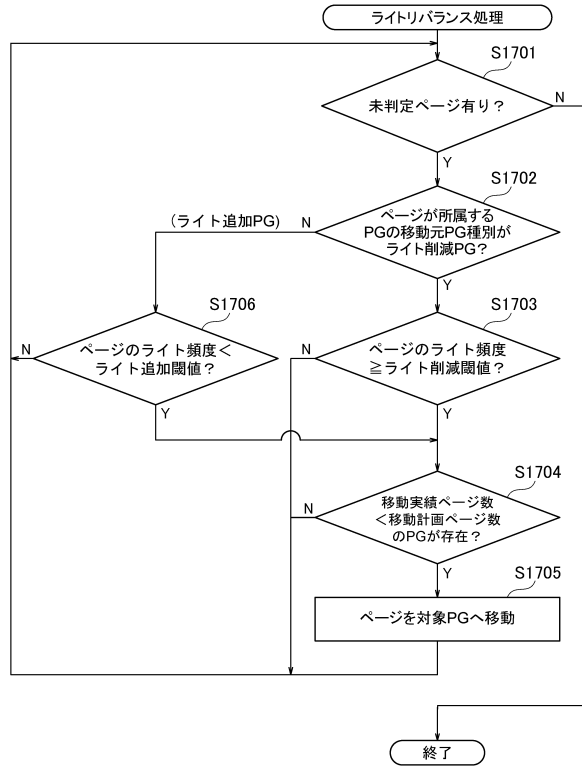
【図16】

図16



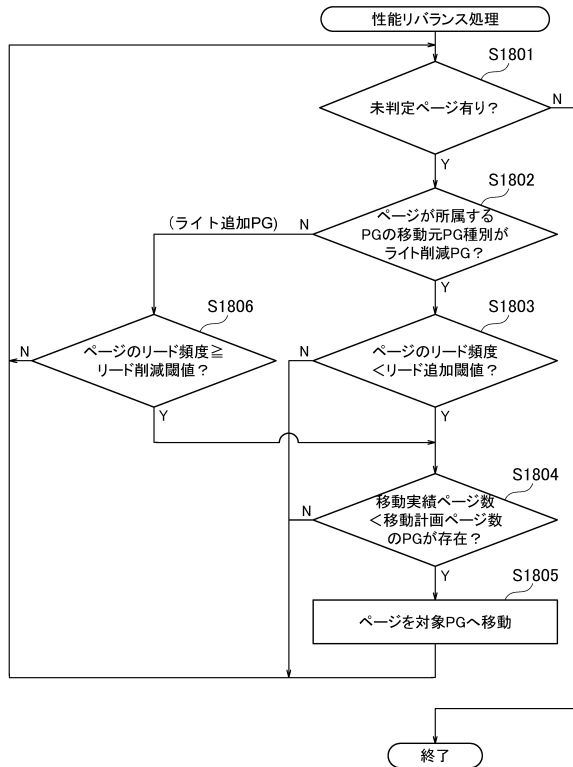
【図17】

図17



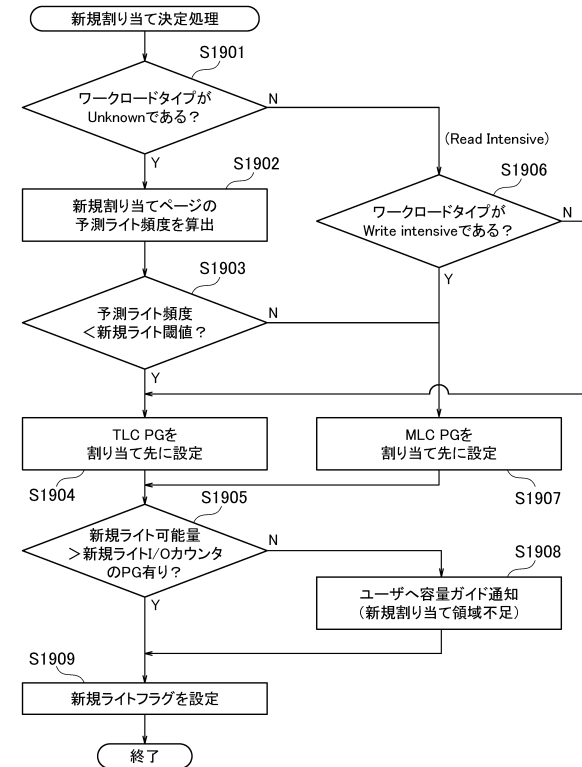
【図18】

図18



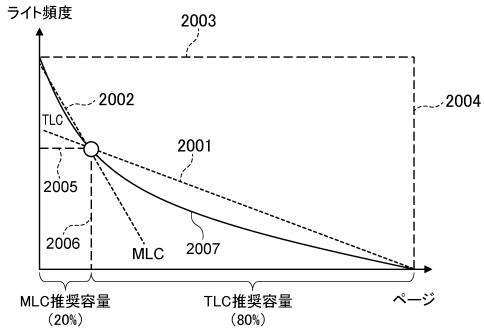
【図19】

図19



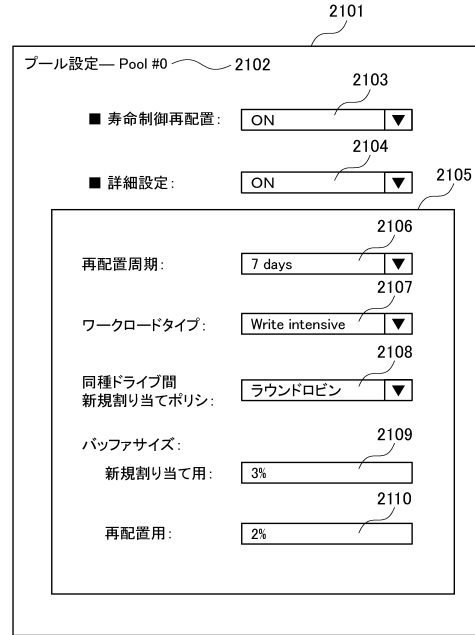
【図20】

図20



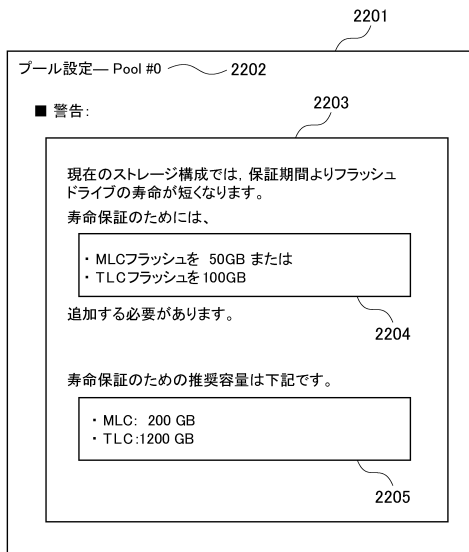
【図21】

図21



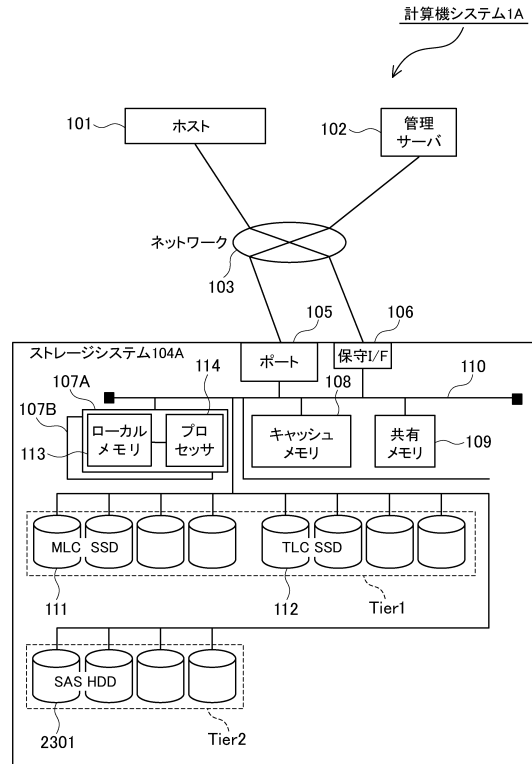
【図22】

図22

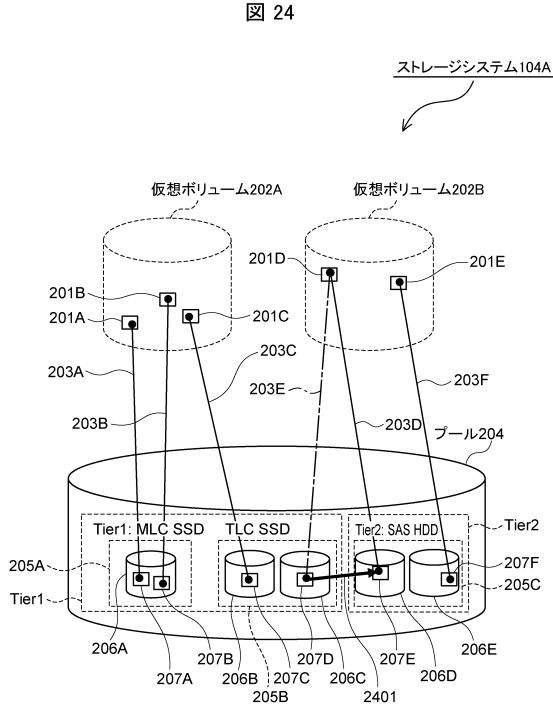


【図23】

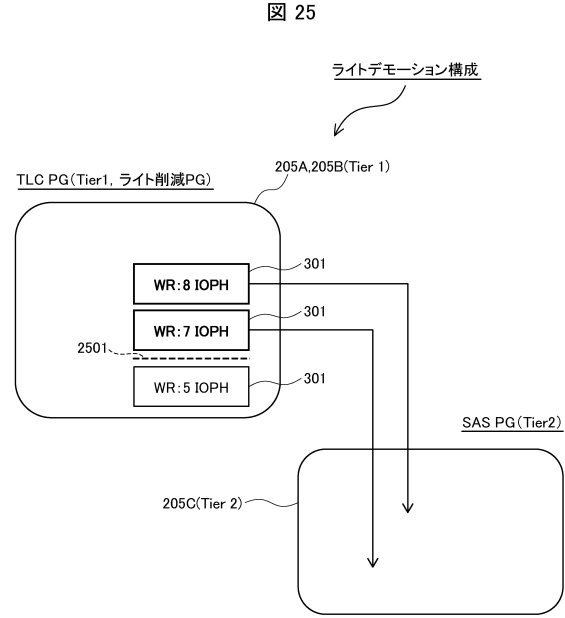
図23



【図24】



【図25】



【図26】

図 26

パリティグループ毎再配置管理テーブル503A

|      |                      |         |         |         |
|------|----------------------|---------|---------|---------|
| 901  | パリティグループ番号           | 0       | 1       | 2       |
| 902  | メディアタイプ              | TLC     | SAS     | MLC     |
| 903  | PGタイプ                | ライト削減PG | ライト追加PG | ライト追加PG |
| 904  | 移動先PG                | 1       | -       | 3       |
| 905  | ライト削減閾値              | 15531   | -       | -       |
| 906  | ライト追加閾値              | -       | -       | -       |
| 907  | リード削減閾値              | -       | -       | -       |
| 908  | リード追加閾値              | 1236    | -       | -       |
| 909  | 移動計画ページ数 (ライトリバランス)  |         |         |         |
| 910  | 移動実績ページ数 (ライトリバランス)  |         |         |         |
| 911  | 移動計画ページ数 (性能リバランス)   |         |         |         |
| 912  | 移動実績ページ数 (性能リバランス)   |         |         |         |
| 913  | 新規ライト可能量             |         |         |         |
| 2601 | ライトデモーション閾値          |         |         |         |
| 2602 | 移動計画ページ数 (ライトデモーション) |         |         |         |
| 2603 | 移動実績ページ数 (ライトデモーション) |         |         |         |
| 2604 | Tier Level           | 1       | 2       | 1       |

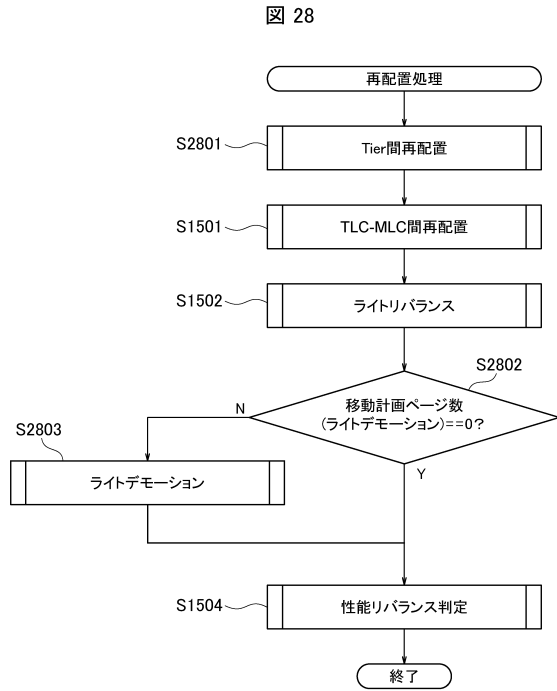
【図27】

図 27

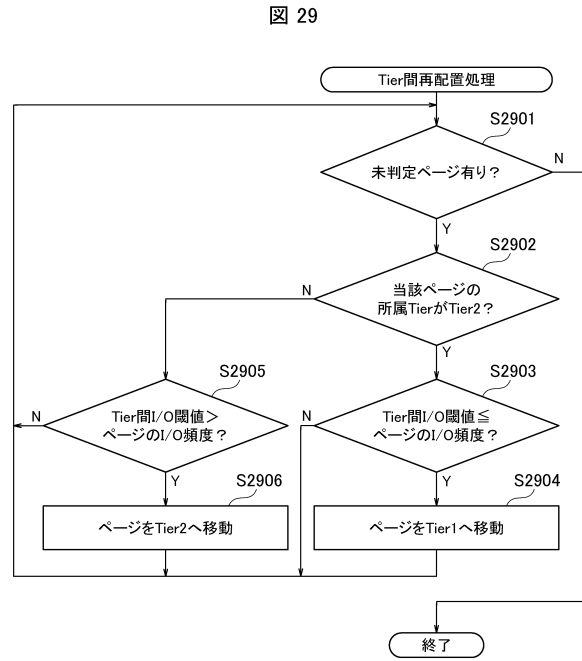
プール毎再配置管理テーブル504A

|      |                   |                 |         |                |
|------|-------------------|-----------------|---------|----------------|
| 1001 | プール番号             | 0               | 1       | 2              |
| 1002 | 寿命制御再配置周期 [hour]  | 168             | 24      | 72             |
| 1003 | TLC-MLC間ライト閾値     | 7040            | 9541    | 4530           |
| 1004 | 新規ライト閾値           | 7550            | 8205    | 4260           |
| 1005 | ワークロードタイプ         | Write intensive | Unknown | Read intensive |
| 1006 | 同種ドライブ間新規割り当てポリシー | 寿命優先            | 容量優先    | ラウンドロビン        |
| 1007 | 新規割り当てバッファ        | 3               | 3       | 5              |
| 1008 | 再配置バッファ           | 2               | 2       | 3              |
| 2701 | Tier間I/O閾値        | 5680            | 7520    | 2602           |

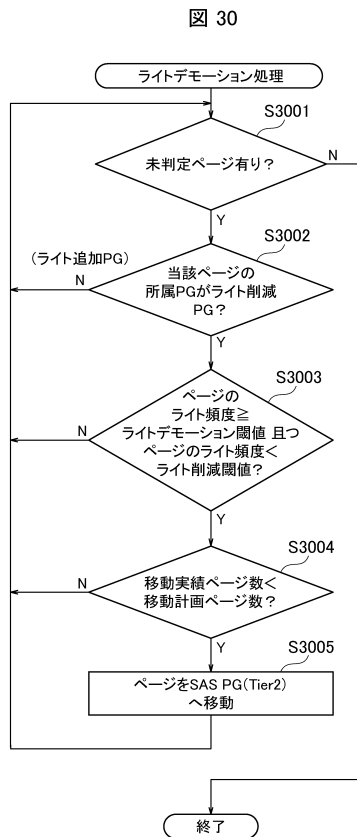
【図 28】



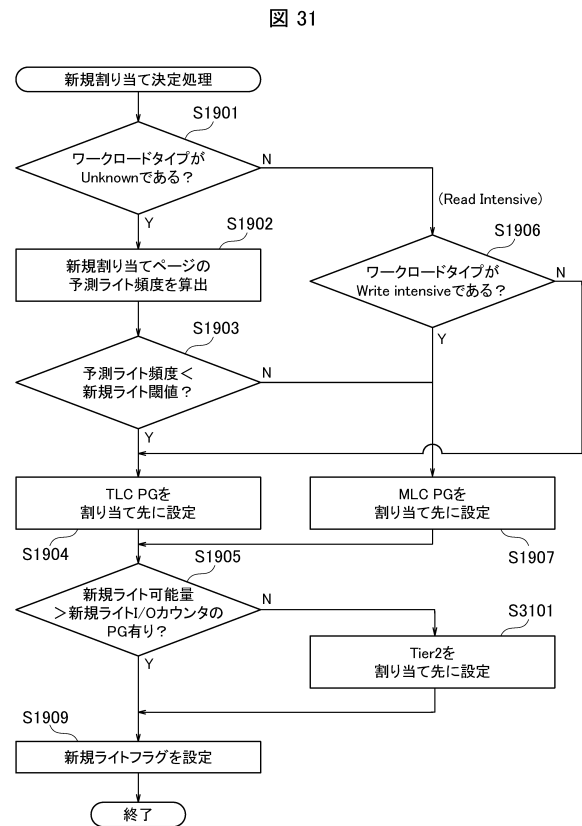
【図 29】



【図 30】



【図 31】





---

フロントページの続き

審査官 田名網 忠雄

(56)参考文献 特表2015-505078(JP,A)  
国際公開第2012/164714(WO,A1)  
特開2010-108246(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06 - 3/08  
G06F 12/00 - 12/06  
G06F 16/00 - 16/958