

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第6678230号  
(P6678230)

(45) 発行日 令和2年4月8日(2020.4.8)

(24) 登録日 令和2年3月18日(2020.3.18)

(51) Int. Cl. F I  
**G06F 3/06 (2006.01)** G O 6 F 3/06 3 O 1 Z  
**G06F 16/185 (2019.01)** G O 6 F 3/06 3 O 1 W  
 G O 6 F 16/185

請求項の数 12 (全 31 頁)

(21) 出願番号	特願2018-502866 (P2018-502866)	(73) 特許権者	000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号
(86) (22) 出願日	平成28年2月29日(2016.2.29)	(74) 代理人	110000062 特許業務法人第一国際特許事務所
(86) 国際出願番号	PCT/JP2016/056018	(72) 発明者	高岡 伸光 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
(87) 国際公開番号	W02017/149592	(72) 発明者	山本 彰 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
(87) 国際公開日	平成29年9月8日(2017.9.8)	(72) 発明者	川口 智大 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
審査請求日	平成30年6月4日(2018.6.4)		

最終頁に続く

(54) 【発明の名称】 ストレージ装置

(57) 【特許請求の範囲】

【請求項1】

ホスト計算機から書き込み要求のあったライトデータを格納するための1以上の記憶デバイスと、前記ホスト計算機に1以上のボリュームを提供するストレージコントローラと、を有し、

前記ストレージコントローラは、前記ホスト計算機から前記ボリューム内の区画に対するライト要求とライト対象データを受領すると、前記区画に、前記区画と同サイズの前記記憶デバイスの第1記憶領域を割り当て、前記割り当てられた前記第1記憶領域に前記ライト対象データを格納し、

前記ストレージコントローラは、前記ボリューム内の区画ごとに、前記ホスト計算機から最後にライト要求を受け付けた時刻である最終ライト時刻を保持しており、

前記ストレージコントローラは、前記最終ライト時刻から所定の期間以上ライト要求を受け付けていない前記区画について、前記ホスト計算機から前記区画に書き込まれたデータと同一のデータが前記記憶デバイスに格納済みの場合、前記区画に書き込まれたデータを前記記憶デバイスに格納しないようにする重複排除処理を実施し、

前記ストレージコントローラは、前記重複排除処理において、前記区画を複数の重複排除ブロックに区分し、前記重複排除ブロックごとに、前記重複排除ブロックと同サイズで前記第1記憶領域と異なる記憶領域であるデータブロックを割り当て、前記重複排除ブロックと前記データブロックとのマッピングを管理するよう構成されており、

前記ストレージコントローラは、前記重複排除ブロックのうち第2の重複排除ブロック

10

20

に前記データブロックを割り当てる時、前記第2の重複排除ブロックに書き込まれているデータが、第1の重複排除ブロックに割り当てられている第1データブロックに格納されているデータと同一の場合、前記第2の重複排除ブロックに前記第1データブロックを割り当てる、  
ことを特徴とする、ストレージ装置。

【請求項2】

前記ストレージコントローラは、前記区画ごとに状態を管理しており、

前記区画の状態として、前記重複排除処理が実施されていない状態である第1状態と、前記重複排除処理が実施された状態である第2状態と、前記第2状態から前記重複排除処理が行われる前の状態に戻された状態である第3状態とがあり、

前記ストレージコントローラは、前記第1状態の前記区画に対して前記重複排除処理を行った結果、前記区画の重複排除率が所定の閾値未満だった場合、前記区画の状態を前記第3状態に変更する

ことを特徴とする、請求項1に記載のストレージ装置。

【請求項3】

前記第3状態の前記区画は、前記ホスト計算機からのライト要求が発行されるまでは前記状態が変更されない、

ことを特徴とする、請求項2に記載のストレージ装置。

【請求項4】

前記ストレージコントローラは、前記ホスト計算機から前記第2状態または前記第3状態の前記区画に対してライト要求が発行されると、前記区画の状態を前記第1状態に変更する、

ことを特徴とする、請求項3に記載のストレージ装置。

【請求項5】

前記ストレージコントローラは、前記重複排除ブロックに書き込まれたデータの特徴量を算出し、検索テーブルに、前記算出された特徴量と前記重複排除ブロックに割り当てられた前記データブロックとのマッピングを記録し、

前記ストレージコントローラは、前記第2の重複排除ブロックの特徴量を算出すると、前記第2の重複排除ブロックの特徴量と同一の値が前記検索テーブルに格納されていない場合、前記第2の重複排除ブロックに第2データブロックを割り当て、前記第2データブロックに前記第2の重複排除ブロックに書き込まれたデータを格納する、

ことを特徴とする、請求項4に記載のストレージ装置。

【請求項6】

ホスト計算機から書き込み要求のあったライトデータを格納するための1以上の記憶デバイスと、前記ホスト計算機に1以上のボリュームを提供するストレージコントローラと、を有するストレージ装置において、

前記ストレージコントローラが、前記ホスト計算機から前記ボリューム内の区画に対するライト要求とライト対象データを受領すると、前記区画に、前記区画と同サイズの前記記憶デバイスの第1記憶領域を割り当て、前記割り当てられた前記第1記憶領域に前記ライト対象データを格納するステップと、

前記ストレージコントローラが、前記ボリューム内の区画ごとに、前記ホスト計算機から最後にライト要求を受け付けた時刻である最終ライト時刻を記録するステップと、

前記最終ライト時刻から所定の期間以上ライト要求を受け付けていない前記区画を検出するステップと、

前記検出された区画について、前記区画に書き込まれたデータのうち、前記記憶デバイスに格納済みのデータと異なるデータのみを前記記憶デバイスに格納する重複排除処理を実施するステップと、を実行し、

前記ストレージコントローラは、前記重複排除処理において、前記区画を複数の重複排除ブロックに区分し、前記重複排除ブロックごとに、前記重複排除ブロックと同サイズで前記第1記憶領域と異なる記憶領域であるデータブロックを割り当て、前記重複排除プロ

10

20

30

40

50

ックと前記データブロックとのマッピングを管理するよう構成されており、

前記ストレージコントローラが、前記重複排除ブロックのうち第2の重複排除ブロックに前記データブロックを割り当てる時、前記第2の重複排除ブロックに書き込まれているデータが、第1の重複排除ブロックに割り当てられている第1データブロックに格納されているデータと同一の場合、前記第2の重複排除ブロックに前記第1データブロックを割り当てるステップを実行する、

ことを特徴とする、ストレージ装置の制御方法。

【請求項7】

前記ストレージコントローラは、前記区画ごとに状態を管理しており、

前記区画の状態として、前記重複排除処理が実施されていない状態である第1状態と、  
前記重複排除処理が実施された状態である第2状態と、前記第2状態から前記重複排除処理が行われる前の状態に戻された状態である第3状態とがあり、

前記重複排除処理を実施するステップでは、前記ストレージコントローラが、前記第1状態の前記区画に対して前記重複排除処理を行った結果、前記区画の重複排除率が所定の閾値未満だった場合、前記区画の状態を前記第3状態に変更する、  
ことを特徴とする、請求項6に記載のストレージ装置の制御方法。

【請求項8】

前記第3状態の前記区画は、前記ホスト計算機からのライト要求が発行されるまでは前記状態が変更されない、

ことを特徴とする、請求項7に記載のストレージ装置の制御方法。

【請求項9】

前記ストレージコントローラは、前記ホスト計算機から前記第2状態または前記第3状態の前記区画に対してライト要求が発行されると、前記区画の状態を前記第1状態に変更するステップを実行する、

ことを特徴とする、請求項8に記載のストレージ装置の制御方法。

【請求項10】

ライトデータを格納するための1以上の記憶デバイスを有する計算機のプロセッサで実行されるプログラムを記録した記憶媒体であって、前記プログラムは前記プロセッサに、ボリュームに対するライト要求とライト対象データを受け付けるステップと、

前記プロセッサは、前記計算機から前記ボリューム内の区画に対するライト要求とライト対象データを受領すると、前記区画に、前記区画と同サイズの前記記憶デバイスの第1記憶領域を割り当て、前記割り当てられた前記第1記憶領域に前記ライト対象データを格納するステップと、

前記ボリューム内の区画ごとに、最後にライト要求を受け付けた時刻である最終ライト時刻を保持するステップと、

前記最終ライト時刻から所定の期間以上ライト要求を受け付けていない前記区画を検出するステップと、

前記検出された区画について、前記区画に書き込まれたデータのうち、前記記憶デバイスに格納済みのデータと異なるデータのみを前記記憶デバイスに格納する重複排除処理を実施するステップと、を  
実行させ、

前記プロセッサは、前記重複排除処理において、前記区画を複数の重複排除ブロックに区分し、前記重複排除ブロックごとに、前記重複排除ブロックと同サイズで前記第1記憶領域と異なる記憶領域であるデータブロックを割り当て、前記重複排除ブロックと前記データブロックとのマッピングを管理するよう構成されており、

前記プロセッサが、前記重複排除ブロックのうち第2の重複排除ブロックに前記データブロックを割り当てる時、前記第2の重複排除ブロックに書き込まれているデータが、第1の重複排除ブロックに割り当てられている第1データブロックに格納されているデータと同一の場合、前記第2の重複排除ブロックに前記第1データブロックを割り当てるステップを実行させる、

ことを特徴とする、プログラムを記録した記憶媒体。

10

20

30

40

50

## 【請求項 1 1】

前記プロセッサは、前記区画ごとに状態を管理しており、

前記区画の状態として、前記重複排除処理が実施されていない状態である第 1 状態と、前記重複排除処理が実施された状態である第 2 状態と、前記第 2 状態から前記重複排除処理が行われる前の状態に戻された状態である第 3 状態とがあり、

前記重複排除処理を実施するステップでは、前記プロセッサが前記第 1 状態の前記区画に対して前記重複排除処理を行った結果、前記区画の重複排除率が所定の閾値未満だった場合、前記区画の状態を前記第 3 状態に変更させる、  
ことを特徴とする、請求項 1 0 に記載のプログラムを記録した記憶媒体。

## 【請求項 1 2】

前記プロセッサに、前記第 2 状態または前記第 3 状態の前記区画に対するライト要求を受領すると、前記区画の状態を前記第 1 状態に変更するステップを実行させる、  
ことを特徴とする、請求項 1 1 に記載のプログラムを記録した記憶媒体。

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

本発明は、ストレージ装置に関する。

## 【背景技術】

## 【0002】

ストレージ装置は、データを格納する複数の記憶デバイスと、記憶デバイスを制御するストレージコントローラとを有しており、ホスト計算機に大容量のデータ格納空間を提供することを目的としている。

## 【0003】

ストレージ装置には、低コストで大量のデータを保存することが求められる。こうした要求を満たすために、ホストから受領したライトデータのサイズを縮小して記憶デバイスに記録する技術が知られている。ライトデータのサイズを縮小してから記憶デバイスに記録すると、データの保持コスト（記憶媒体のビットコスト、ストレージ装置の消費電力コスト等）を削減できる。データサイズを縮小するために、可逆圧縮アルゴリズムを用いて、データの意味を保ったままデータサイズを縮小する技術がある。この処理は「可逆圧縮」または「圧縮」と呼ばれる。

## 【0004】

データを圧縮して記憶デバイスに記録する場合、ライト時には圧縮処理、リード時には圧縮データの伸長処理のオーバーヘッドが発生するため、アクセス性能が低下することがある。つまり格納データ量の削減とアクセス性能はトレードオフの関係にある。これを避けるために、選択的にデータの圧縮を行う技術が存在する。たとえば特許文献 1 には、複数の記憶階層（tier）を管理するストレージシステムにおいて、例えば下位 tier に移動されるデータを圧縮して格納することで、アクセス性能の低下を抑止しつつ、格納データ量を削減する方法が開示されている。

## 【0005】

また格納すべきデータ量を削減するもう一つの技術として、重複排除技術がある。たとえばストレージ装置に同内容のデータが複数個存在していることをストレージ装置が検出した時、そのうちの 1 つだけをストレージ装置内の記憶デバイスに残し、残りのデータは記憶デバイスに格納しないようにする技術である。重複排除技術と可逆圧縮技術の何れが用いられても、ホストから受領したライトデータの量より記憶デバイスに格納されるライトデータの量が小さくなる。そのため重複排除技術も広義には圧縮技術の 1 つといえる。

## 【先行技術文献】

## 【特許文献】

## 【0006】

【特許文献 1】米国特許 8 3 5 9 4 4 4 号明細書

## 【発明の概要】

10

20

30

40

50

## 【発明が解決しようとする課題】

## 【0007】

データを圧縮する場合、データ内容に依存して圧縮率（あるいはデータの縮小量）が異なり得る。そのため、データに圧縮または重複排除処理が施されたが、記憶デバイスに格納されるデータ量が殆ど削減されないという事も起こり得る。その場合、アクセス性能も低下し、かつデータの保持コストも低下しないことになる。低ビットコストで高性能なストレージ装置を提供するためには、このような事態が発生することを防ぐ必要がある。

## 【課題を解決するための手段】

## 【0008】

本発明の一観点に係るストレージ装置は、ホスト計算機から書き込み要求のあったライトデータを格納するための1以上の記憶デバイスと、ホスト計算機に1以上のボリュームを提供するストレージコントローラとを有し、ボリューム内の区画ごとに、ホスト計算機から最後にライト要求を受け付けた時刻である最終ライト時刻を保持する。そしてストレージコントローラは、最終ライト時刻から所定の期間以上ライト要求を受け付けていない区画を検出すると重複排除処理を実施する。またストレージコントローラは、区画の重複排除処理の結果、重複排除率が低い区画については、区画を重複排除処理の実施されていない状態へと戻す処理を実施する。

10

## 【発明の効果】

## 【0009】

本発明によれば、低ビットコストで高性能なストレージ装置を提供することが可能になる。

20

## 【図面の簡単な説明】

## 【0010】

【図1】実施例に係るストレージ装置を含む計算機システムの論理構成図である。

【図2】論理ページの状態遷移を表した図である。

【図3】ストレージ装置の構成図である。

【図4】論理ページと物理ページのマッピング関係を説明する図である。

【図5】論理ページ管理テーブルの構成例である。

【図6】マッピングテーブルの構成例である。

【図7】プール管理テーブルの構成例である。

30

【図8】検索テーブルの構成例である。

【図9】逆参照テーブルの構成例である。

【図10】追記ポインタの構成例である。

【図11】重複排除処理部のフローチャートである。

【図12】重複排除処理の流れを表した図である。

【図13】重複排除解除部のフローチャートである。

【図14】移行処理のフローチャートである。

【図15】ライト処理のフローチャートである。

【図16】リード処理のフローチャートである。

【図17】物理ページ解放処理のフローチャートである。

40

【図18】共有データ判定・複製処理のフローチャートである。

【図19】マッピング切替処理のフローチャートである。

## 【発明を実施するための形態】

## 【0011】

以下、幾つかの実施例について、図面を用いて説明する。

## 【0012】

なお、以下の実施例において、ストレージ装置内で実行される処理について、「プログラム」を主語として説明を行う場合がある。実際には、ストレージ装置が有するプロセッサ（CPU）がプログラムを実行することによって、プログラムに記述された処理が行われるため、処理の主体はプロセッサ（CPU）であるが、説明が冗長になることを防ぐた

50

め、プログラムを主語にして処理の内容を説明することがある。また、プログラムの一部または全ては専用ハードウェアによって実現されてもよい。また、以下で説明される各種プログラムは、プログラム配布サーバや計算機が読み取り可能な記憶メディアによって提供され、プログラムを実行する各装置にインストールされてもよい。計算機が読み取り可能な記憶メディアとは、非一時的なコンピュータ可読媒体で、例えばICカード、SDカード、DVD等の不揮発性記憶媒体である。

【0013】

実施例の説明に入る前に、実施例で用いられる各種用語について説明する。

【0014】

「ボリューム」とは、ストレージ装置や記憶デバイス等のターゲットデバイスが、ホスト計算機等のイニシエータデバイスに提供する記憶空間のことを意味する。イニシエータデバイスが記憶空間上の領域に対するデータ書き込み要求を発行すると、その領域に対応付けられているターゲットデバイス上の領域にデータが格納される。本実施例に係るストレージ装置はボリュームとして、いわゆるThin Provisioning技術により形成される仮想ボリュームをホストに提供する。仮想ボリュームは、その初期状態（仮想ボリュームが定義された直後）では、記憶空間上の領域に記憶デバイスが対応付けられていない。イニシエータデバイス（ホスト）が記憶空間上の領域にデータ書き込み要求を発行した時点で、ストレージ装置はその領域に対応付けられる記憶デバイスを動的に決定する。

10

【0015】

「重複排除処理」とは、ストレージ装置内に同内容のデータが複数存在する場合、1つだけをストレージ装置に残し、それ以外のデータをストレージ装置から削除する処理である。ストレージ装置内に同内容のデータが存在するか判定する処理のことを、「重複判定」処理と呼ぶ。なお、特に断りのない限り、重複排除処理は重複判定処理を含む処理である。

20

【0016】

以下で説明する実施例に係るストレージ装置では、重複排除ブロックと呼ばれる所定サイズのデータ毎に重複判定を行う。以下の実施例では、重複排除ブロックのサイズが8KBの例について説明されるが、重複排除ブロックのサイズは8KB以外のサイズであってもよい。同内容のデータのことを「重複データ」と呼ぶ。

30

【0017】

重複判定の際、2つのデータをビット単位あるいはバイト単位で比較すると、判定処理に長時間を要することになる。そのため一般的には重複判定を行う装置は、比較対象のデータに所定の演算（たとえばハッシュ関数を用いた演算等）を行うことで、小サイズ（たとえば8バイト程度）の特徴量を生成し、それを用いて重複判定を行う。以下の実施例では、データから生成される特徴量のことを、「フィンガープリント」と呼ぶ。フィンガープリントは、FPと略記されることもある。

【0018】

以下で説明する実施例では、データAから算出されたFPの値がHであった場合、値HはデータAのFPと呼ばれる。逆にデータAのことを、「FP Hを持つデータ」と呼ぶことがある。また、データAの書き込まれる領域（重複排除ブロック）のことも、「FP Hを持つ領域（重複排除ブロック）」と呼ぶことがある。

40

【0019】

本実施例において「衝突」とは、複数の異なるデータそれぞれに対して所定の演算を施してFPを生成した時、生成されたそれぞれのFPが同一になることを意味する。ハッシュ関数などを用いて、小サイズの特徴量を算出する場合、衝突は発生し得る。

【0020】

「重複排除率」とは、重複排除処理による記憶領域消費量の削減効率を表す指標値である。たとえば重複排除率は、ストレージ装置のボリュームに書き込まれたデータ量と、ストレージ装置がデータ格納のために使用（消費）した記憶領域の量の比で表される値であ

50

る。重複排除処理では、同内容のデータが多数格納されると、1つだけがストレージ装置の記憶領域に書き込まれ（記憶領域が消費され）、それ以外のデータは記憶領域に書き込まれない（記憶領域が消費されない）ので、記憶領域消費量の削減効率が高くなる。

【実施例】

【0021】

(1) 発明の概要

まず図1、図2、そして図4を用いて、本発明の実施例に係るストレージ装置が実施する、重複排除方法の概要を説明する。図1は、本発明の実施例に係るストレージ装置が計算機に提供する仮想ボリュームの構成を表した図である。

【0022】

ストレージ装置1は、複数の記憶デバイス（図1では非図示）を有し、ホスト計算機5（以下では「計算機5」と略記する）からのライトデータを記憶デバイスに格納する。記憶デバイスは、所定サイズの記憶空間をストレージ装置1に提供するが、ストレージ装置1は、記憶デバイスの提供する記憶空間を直接計算機5には提供しない。計算機5には、記憶デバイスの有する記憶空間とは異なる、1以上の仮想的な記憶空間を提供する。この仮想的な記憶空間を「仮想ボリューム」と呼ぶ。図1では、2つの仮想ボリューム（仮想ボリューム16、仮想ボリューム20）が計算機5に提供される例が示されている。

【0023】

ストレージ装置1は、仮想ボリュームの記憶空間を、複数の所定サイズ（一例として42MB）の区画に分割して管理している。本実施例では、この区画のことを「論理ページ」と呼ぶ。各論理ページには仮想ボリューム内で一意な識別子が付されており、この識別子を論理ページ識別子（または論理ページ番号）と呼ぶ。

【0024】

仮想ボリュームは、公知のThin Provisioning技術などを用いて形成されるボリュームであり、ストレージ装置1は仮想ボリュームの論理ページに対するアクセス要求を受け付けた時点で、記憶デバイスの記憶領域を動的に論理ページに割り当てる（マップする）。言い換えると、各論理ページには、計算機5からアクセス要求を受け付けるまでは、その論理ページに記憶領域が割り当てられていない。

【0025】

ストレージ装置1はまた、論理ページに対して割り当てる記憶領域（これは記憶デバイスが提供する記憶領域である）を管理するための管理概念を有しており、これを「ストレージプール」または「プール」と呼ぶ。図1では、円筒状のオブジェクト17としてストレージプールが表現されている。ストレージ装置1は、プールの記憶領域を、論理ページと同サイズの領域（あるいは論理ページより大きいサイズの領域でも良い）に区分し、区分された領域ごとに識別子を付して管理する。この区分された領域は「物理ページ」と呼ばれ、物理ページに付される識別子は「物理ページ識別子」または「物理ページ番号」と呼ばれる。

【0026】

ストレージ装置1が計算機5から、仮想ボリュームに対するライト要求を受信すると、ライト要求に含まれているライト対象領域のアドレスを論理ページ番号に変換し、ライト対象領域を含む論理ページを特定する。特定された論理ページに物理ページが割り当てられていない場合、ストレージ装置1は、プール17内の未使用の物理ページ（まだ論理ページに割り当てられていない物理ページ）を選択し、アクセス対象の論理ページに、選択された物理ページを割り当てる（マップする）。計算機5からのライトデータは、このアクセス対象論理ページにマップされた物理ページに格納される。

【0027】

またストレージ装置1は、論理ページと、論理ページに割り当てられた物理ページとの対応関係（マッピング）を論理ページ管理テーブル126に記憶している。論理ページに対するリード要求を受け付けた時には、ストレージ装置1は論理ページ管理テーブル126を参照することで、論理ページに割り当てられた記憶領域を特定し、特定された記憶領

10

20

30

40

50

域からデータを読み出す。

【0028】

本実施例に係るストレージ装置1は、計算機5からライトデータの書き込まれた論理ページのうち、所定の条件に該当する論理ページを重複排除処理対象として扱うことと決定し、その論理ページに対して重複排除処理を行う。所定の条件とはたとえばライト頻度の低い論理ページで、具体的には、所定時間以上ライトデータの書き込みがない論理ページである。重複排除処理の行われた論理ページのことを、「重複排除論理ページ」と呼ぶ。逆に、重複排除論理ページでない論理ページのことを、「通常論理ページ」と呼ばれることがある。

【0029】

図4は、通常論理ページに物理ページが割り当てられた状態と、重複排除論理ページに物理ページが割り当てられた状態の概念図である。図4(a)は、通常論理ページに物理ページが割り当てられた状態を表している。1つの通常論理ページ161bには1つの物理ページ171cが割り当てられる。たとえば計算機5が通常論理ページ161bの先頭からkバイト目の領域361にデータをライトする要求をストレージ装置1に発行すると、ストレージ装置1は、通常論理ページ161bに割り当てられた物理ページ171cの先頭からkバイト目の領域371に、ライト要求で指定されたデータを格納する。この関係が維持されている為、ストレージ装置1は、通常論理ページと物理ページのマッピングを管理する際には、論理ページごとに、論理ページに割り当てられている物理ページの物理ページ番号だけを管理すればよい。

【0030】

図4(b)は、重複排除論理ページに物理ページが割り当てられた状態を表している。ストレージ装置1は、論理ページ内の領域を、所定サイズ(たとえば8KB)の部分領域に区分し、この部分領域ごとに物理ページの領域を割り当てる。本実施例ではこの部分領域ごとに重複判定が行われるので、これを「重複排除ブロック」と呼ぶ。一方、重複排除ブロックに割り当てられる物理ページ上の領域は、「データブロック」と呼ばれる。

【0031】

なお、通常論理ページが重複排除論理ページとして扱われるようになると、これまで通常論理ページに割り当てられていた物理ページは削除される(割り当てられていない状態にされる)。そしてこれまで物理ページに格納されていたデータは、重複排除ブロックに割り当てられる物理ページ上の領域(データブロック)に移動される。データブロックは、プール17に属する物理ページ上の領域である。重複排除ブロックとデータブロックのマッピング情報は、後述するマッピングテーブル127に格納される。

【0032】

ストレージ装置1が行う重複排除処理では、重複排除ブロックごとにデータの比較を行う。図4(b)を用いて重複排除処理について概説する。図4(b)では、161aは重複排除論理ページであり、重複排除ブロック18bには物理ページ171d内のデータブロック19cが割り当てられている。ここでストレージ装置1が論理ページ201xに対して重複排除処理を行う場合を想定する。ストレージ装置1は、論理ページ201xの重複排除ブロック21cに対して書き込まれたデータと、物理ページ171dの各データブロックに格納されているデータの内容を比較する。比較の結果、重複排除ブロック21cに対して書き込まれたデータがデータブロック19cに格納されているデータと同一だった場合、ストレージ装置1は重複排除ブロック21cにデータブロック19cを割り当てる。そしてストレージ装置1は、重複排除ブロック21cに対して書き込まれたデータを、新たに物理ページ171dに書き込むことはしない。この結果、データブロック19cは、2つの重複排除ブロック(18bと21c)に割り当てられることになる。本実施例では、データブロックが複数の重複排除ブロックに割り当てられた状態のことを、「データブロックが複数の重複排除ブロックに共有されている」または「複数の重複排除ブロックが1つのデータブロックを共有している」と表現する。

【0033】

10

20

30

40

50



これにより、ストレージ装置 1 は計算機 5 から重複排除ブロック 18 b と重複排除ブロック 21 c のいずれに対するリード要求を受け付けた場合でも、データブロック 19 c からデータ A を読み出して、計算機 5 に返送する。また、重複排除ブロック 18 b と重複排除ブロック 21 c のそれぞれに対して異なるデータブロックを割り当てる必要がないため、実質的にストレージプール 17 への格納データ量（言い換えると、記憶デバイスの記憶領域消費量）が 1 / 2 に削減される効果がある。もし n 個の重複排除ブロックが 1 つのデータブロックを共有している場合、記憶領域消費量は 1 / n に削減されることになる。

【0034】

このように、本実施例に係るストレージ装置 1 では、論理ページが、重複排除論理ページとして扱われる場合と通常論理ページとして扱われる場合がある。ストレージ装置 1 では、各論理ページがいずれの状態にあるかを管理するため、各論理ページの状態についての情報を保持する。図 2 を用いて、論理ページの各状態の説明を行う。

10

【0035】

図 2 は、ある論理ページの状態遷移を表した図である。論理ページは、P、X、Q のいずれかの状態を有する。初期状態（仮想ボリュームが定義された直後で、論理ページに物理ページが割り当てられていない）では、論理ページの状態は P である。論理ページに計算機 5 からデータのライトが行われると、ストレージ装置 1 はその論理ページの状態を X に変更する。

【0036】

状態が X の論理ページに対して、一定以上の間計算機 5 からのライト要求が到来しなかった場合、ストレージ装置 1 はその論理ページの状態を Q に変更する。状態が Q の論理ページは、重複排除論理ページであり、重複排除処理が行われる。状態が Q の論理ページ（重複排除論理ページ）に対して、計算機 5 からのライト要求が到来すると、ストレージ装置 1 は再びその論理ページの状態を X に変更する。重複排除論理ページに対してデータを格納する場合、重複判定等の処理のオーバーヘッドが増加し、アクセス性能の低下を招くためである。

20

【0037】

計算機 5 からのライト要求が到来しない場合には、重複排除論理ページは原則としてその状態（状態 Q）が維持される。ただし、その重複排除論理ページの記憶領域消費量の削減効果が薄い場合、たとえばその重複排除論理ページの各重複排除ブロックに書き込まれたデータの大半が、他のデータと異なっている場合は、ストレージ装置 1 のビットコスト削減に寄与せず、かつアクセス性能の低下を招くこともある。そのためストレージ装置 1 は、記憶領域消費量の削減効果が薄い重複排除論理ページがあった場合、その重複排除論理ページの状態を Q から P へと遷移させる。

30

【0038】

ストレージ装置 1 は、状態を Q から P に変更させるべき重複排除論理ページを決定する際、論理ページの重複排除率を用いる。なお、本実施例における“論理ページの重複排除率”とは、以下の計算式により定義されるものである。

【0039】

論理ページ内の重複排除ブロック数を P、論理ページ（重複排除論理ページ）の重複排除ブロックのうち、既に別の重複排除ブロックへ割り当てられているデータブロックが割り当てられた重複排除ブロックの数を W とする（別の表現をすると、重複排除論理ページに割り当てられたデータブロックのうち、複数の重複排除ブロックに共有されているデータブロックの数が W である）。なお論理ページ内の重複排除ブロック数 P は、論理ページサイズ（たとえば 42 MB）÷ 重複排除ブロックのサイズ（たとえば 8 KB）で求まる固定値である。このとき論理ページの重複排除率 D は、以下の計算式で表現される。

40

$$D = W \div P$$

【0040】

P は固定値のため、W が大きいと重複排除率 D も大きくなる。そのため、D が大きい（1 に近い）場合、記憶デバイスの記憶領域の消費量が少ない（重複排除による記憶領域の

50

削減効果が高い)ことを意味し、逆にDが小さい(0に近い場合)は、重複排除による記憶領域の削減効果が小さいことを意味する。本実施例ではストレージ装置1は、重複排除率が所定の閾値未満の重複排除論理ページがあった場合、その重複排除論理ページの状態をQからPへと遷移させる。

#### 【0041】

論理ページの状態がQからPに変更される場合、ストレージ装置はその論理ページのデータを、一旦ストレージ装置1内のキャッシュメモリ等の一時的な記憶領域上に読み出し、論理ページに新たな物理ページを割り当てて、割り当てられた物理ページにデータを書き戻すことで、重複排除処理が行われていない状態にする。本実施例に係るストレージ装置は、論理ページのアクセス頻度や記憶領域消費量の削減効率を考慮して、重複排除処理の対象にする論理ページを制限することで、高アクセス性能と低ビットコストの両立を図っている。

10

#### 【0042】

なお、状態Pの論理ページと状態Xの論理ページは、以下の点で異なる。状態Xの論理ページは、状態Qに変更されることがある。具体的には、状態Xのある論理ページに、計算機5からのライト要求が一定の時間以上到来しなかった場合に、その論理ページの状態はQに変更される。一方状態Pの論理ページは、状態がQに変更されることはない。状態Pの論理ページは、記憶領域消費量の削減効果が薄いと判定されたために、状態がQからPに変更された論理ページだからである(あるいは、まだ計算機5から一度もライトデータの書き込みが行われておらず、物理ページが割り当てられていない論理ページである)。そのような論理ページの状態をQに変更しても、記憶領域消費量の削減効果が小さい(またはない)ことは明らかであるため、ストレージ装置1は状態Pの論理ページの状態をQに変更することは行わない。これにより、頻りに論理ページの状態の変更が発生しないようにしている(処理オーバーヘッドの増加を抑制している)。

20

#### 【0043】

##### (2) システム構成

図3は、本実施例に係るストレージ装置1を含む計算機システムのハードウェア構成例を示している。ストレージ装置1は、ストレージコントローラ10と、ストレージコントローラ10に接続された複数の記憶デバイス15を有する。

#### 【0044】

記憶デバイス15は、ストレージ装置1が計算機5などの上位装置からのライトデータを記憶するために用いられる。記憶デバイスとしては、たとえば磁気ディスクを記憶媒体として用いるHDD(Hard Disk Drive)や、フラッシュメモリ等の不揮発性半導体メモリを記憶媒体として採用したSSD(Solid State Drive)が用いられる。本実施例では、記憶デバイス15は「ドライブ15」と表記されることもある。記憶デバイス15は一例として、SAS(Serial Attached SCSI)規格に従う伝送線(SASリンク)や、PCI(Peripheral Component Interconnect)規格に従う伝送線(PCIリンク)などによって、ストレージコントローラ10と接続される。

30

#### 【0045】

ストレージコントローラ10は少なくとも、プロセッサ(CPUとも呼ばれる)11、システムメモリ12、キャッシュメモリ14、そしてSAN(Storage Area Network)6に接続するためのインタフェース(非図示)を有する。SAN6は、一例としてファイバチャネルを用いて形成されるネットワークである。

40

#### 【0046】

プロセッサ11は、ストレージ装置1の各種制御を行う。システムメモリ12は、CPU11が実行するプログラムや、CPU11がプログラム実行の際に使用する管理情報を格納するためのものである。一方、キャッシュメモリ14は、計算機5から受領したライトデータや、記憶デバイス15から読み出されたリードデータを一時的に記憶するために用いられる。

50

## 【 0 0 4 7 】

システムメモリ 1 2 やキャッシュメモリ 1 4 には、D R A M、S R A M 等の揮発性記憶媒体が用いられるが、別の実施形態として、不揮発性メモリを用いてキャッシュメモリ 1 4 を構成してもよい。また、キャッシュメモリ 1 4 に揮発性記憶媒体が用いられる場合、ストレージ装置 1 にバッテリー等の補助電源を搭載し、停電時にキャッシュメモリ 1 4 の記憶内容を維持できるように構成されていてもよい。

## 【 0 0 4 8 】

また、別の実施形態として、ストレージ装置 1 は 2 種類のメモリ（システムメモリ 1 2 やキャッシュメモリ 1 4 ）を有さない構成でも良い。つまりストレージ装置 1 が 1 種類のメモリのみを有する構成でも良い。その場合、プログラム、管理情報、ライトデータなどは、同じメモリに格納される。

10

## 【 0 0 4 9 】

計算機 5 は、ストレージ装置 1 へのアクセス要求発行元となる装置である。計算機 5 は、P C（パーソナルコンピュータ）等の汎用のコンピュータであり、少なくともプロセッサとメモリ（非図示）を有する。プロセッサでは、ストレージ装置 1 が提供する仮想ボリュームに I / O 要求を発行するアプリケーションプログラムなどが実行される。

## 【 0 0 5 0 】

## （ 3 ）管理情報

続いてストレージ装置 1 の有する管理情報、プログラムの内容の説明を行う。ストレージ装置 1 のシステムメモリ 1 2 には少なくとも、論理ページ管理テーブル 1 2 6、マッピングテーブル 1 2 7、プール管理テーブル 1 2 8、検索テーブル 1 2 9、逆参照テーブル 1 3 0、追記ポインタ 1 3 1 の、6 種類の管理情報が格納されている。以下では、これらの各種管理情報の内容について説明する。

20

## 【 0 0 5 1 】

図 5 に、論理ページ管理テーブル 1 2 6 の構成を示す。論理ページ管理テーブル 1 2 6 は、各論理ページの状態を管理するためのテーブルで、各行（レコード）に、それぞれの論理ページの状態や属性情報が格納される。以下、論理ページ管理テーブル 1 2 6 の各カラムに格納される情報について説明する。

## 【 0 0 5 2 】

仮想ボリューム 1 2 6 1、論理ページ 1 2 6 2 にはそれぞれ、論理ページの属する仮想ボリュームの識別子、論理ページの識別子が格納される。

30

## 【 0 0 5 3 】

重複排除 1 2 6 3 には、“有効”または“無効”が格納される。あるレコードの重複排除 1 2 6 3 に、“有効”が格納されている場合、そのレコードで管理される論理ページは重複排除論理ページであることを意味し、“無効”が格納されている場合、そのレコードで管理される論理ページは通常論理ページであることを意味する。

## 【 0 0 5 4 】

物理ページ 1 2 6 4 には、論理ページに割り当てられている物理ページの識別子が格納される。論理ページに物理ページが割り当てられていない場合には、N U L L が格納される。

40

## 【 0 0 5 5 】

なお、一般に仮想ボリュームや論理ページの識別子には非負の整数値が使われるが、図 5 では説明の都合上、仮想ボリューム 1 2 6 1、論理ページ 1 2 6 2、物理ページ 1 2 6 4 の各カラムには、図 1 等に記載の仮想ボリュームに付された参照番号、論理ページに付された参照番号、物理ページに付された参照番号が格納されている。また、これ以降で説明される管理情報についても同様に、仮想ボリューム、論理ページ、物理ページ等の識別子を格納するためのカラムには、図 1 等に記載の仮想ボリューム、論理ページ、物理ページ等に付された参照番号が格納されている。

## 【 0 0 5 6 】

状態 1 2 6 5 には、論理ページの状態が格納される。論理ページの状態には、先に述べ

50

たとおり、P、Q、Xの3つがある。最終ライト時刻1266には、論理ページに対して計算機5からライト要求を受領した最新の時刻が格納される。

【0057】

排除ブロック数1267は、論理ページが重複排除論理ページである場合に有効な情報である。排除ブロック数1267には、論理ページ（重複排除論理ページ）内の重複排除ブロックのうち、重複排除処理によって、既に別の重複排除ブロックへ割り当てられているデータブロックがさらに割り当てられた重複排除ブロックの数が記録される。排除ブロック数1267が大きいほど、その論理ページの記憶領域消費量の削減効果が大きいといえる。以降、このように、重複排除処理によって既存のデータブロックが割り当てられることで、記憶領域を消費せずにデータを保持している重複排除ブロックを、削減済み重複排除ブロックと呼ぶことがある。

10

【0058】

続いてマッピングテーブル127について、図6を参照しながら説明する。マッピングテーブル127は、重複排除ブロックとデータブロックのマッピング状態を管理するためのテーブルで、各行（レコード）には、重複排除ブロックの識別子や、重複排除ブロックに割り当てられているデータブロックの識別子等の情報が格納される。以下、マッピングテーブル127の各カラムに格納される情報について説明する。

【0059】

仮想ボリューム1271、論理ページ1272にはそれぞれ、管理対象の重複排除ブロックの属する仮想ボリュームの識別子、論理ページの識別子が格納される。そして重複排除ブロック1273には、管理対象の重複排除ブロックの識別子が格納される。フィンガープリント1274には、重複排除処理において算出した、重複排除ブロックごとのフィンガープリントが格納される。

20

【0060】

物理ページ1275とデータブロック1276にはそれぞれ、管理対象の重複排除ブロックに割り当てられているデータブロックの属する物理ページの識別子、データブロックの識別子が格納される。

【0061】

なお、重複排除ブロックまたはデータブロックの識別子には、重複排除ブロックまたはデータブロックを一意に識別可能な情報であれば、任意の情報が用いられて良い。本実施例では、重複排除ブロックの識別子には、重複排除ブロックの属する仮想ボリューム内のアドレスが用いられる。一方、データブロックの識別子には、データブロックの属する物理ページ内の相対アドレス（物理ページの先頭のデータブロックの識別子を0とするアドレス）が用いられる。ただし図6では、図4に記載の重複排除ブロックとデータブロックのマッピング関係を説明するために、重複排除ブロックの識別子として、図4に記載の重複排除論理ブロックに付されたアルファベットを用い、またデータブロックの識別子としては、データブロックに付されたアルファベットを用いている。

30

【0062】

削減フラグ1277は、管理対象の重複排除ブロックが削減済み重複排除ブロックか否かを示す情報である。管理対象の重複排除ブロックが削減済み重複排除ブロックである場合には、削減フラグ1277にはTRUEが格納され、そうでなければFALSEが格納される。

40

【0063】

物理ページについての情報は、プール管理テーブル128に格納されて管理される。図7を参照しながらプール管理テーブル128の内容を説明する。

【0064】

プール管理テーブル128の各レコードには、物理ページの状態等の情報が格納される。物理ページ1281には、管理対象の物理ページの識別子が格納され、論理ページ1283には、物理ページが割り当てられている論理ページの識別子が格納される。物理ページが論理ページ（通常論理ページ）に割り当てられていない場合、あるいは物理ページ内

50

の領域が重複排除ブロックに割り当てられている場合には、論理ページ 1 2 8 3 には N U L L が格納される。

【 0 0 6 5 】

使用状況 1 2 8 2 には、物理ページの使用状態が格納される。物理ページの状態には、通常論理ページに割り当てられている状態、重複排除ブロックに割り当てられるために使用されている状態、そして未使用状態がありえる。物理ページが通常論理ページに割り当てられている場合、使用状況 1 2 8 2 には“論理ページ”が格納される。物理ページ内の領域が重複排除ブロックに割り当てられている場合には、使用状況 1 2 8 2 には“データブロック”が格納される。物理ページが論理ページにも割り当てられておらず、かつ物理ページ内の領域が重複排除ブロックに割り当てられてもいない場合、その物理ページの状態は“未使用状態”と呼ばれ、その場合には使用状況 1 2 8 2 には“未使用”が格納される。

10

【 0 0 6 6 】

ストレージ装置 1 がたとえば、ある論理ページに物理ページを割り当てる際には、プール管理テーブル 1 2 8 のレコードのうち、使用状況 1 2 8 2 が“未使用”のレコードを 1 つ特定し、特定されたレコードで管理される物理ページを論理ページに割り当てると決定する。またストレージ装置 1 は論理ページ管理テーブル 1 2 6 に、割り当ててことを決定した物理ページの識別子（物理ページ 1 2 8 1）を格納する。

【 0 0 6 7 】

なお、物理ページは、実際には 1 または複数の記憶デバイス 1 5 上の領域である。そのためストレージ装置 1 は、プール管理テーブル 1 2 8 に加えて、物理ページと、物理ページが存在する記憶デバイス 1 5 及び記憶デバイス上のアドレスとのマッピングを管理するための情報も保持する。ストレージ装置 1 は、たとえば計算機 5 からのアクセス要求を受け付けると、アクセス要求で指定されている領域に対応する論理ページを特定し、その後論理ページに割り当てられている物理ページを特定する。さらにその後ストレージ装置 1 は、このマッピング情報を参照することで、アクセス対象の記憶デバイス 1 5 上アドレスを特定する。ただし、このマッピング情報及びマッピング情報を用いたアクセス先記憶デバイス 1 5 のアドレス特定方法は、Thin Provisioning 技術を用いるストレージ装置において公知のものであるため、本実施例では詳細説明を略す。

20

【 0 0 6 8 】

続いて検索テーブル 1 2 9 について説明する。検索テーブル 1 2 9 は重複排除処理で用いられる。重複排除処理では、重複判定対象のデータと同じデータが既にストレージプール 1 7 にあるか判定する処理が行われるが、この判定処理の高速化のために、ストレージ装置 1 は重複排除ブロックごとにフィンガープリントを算出し、算出したフィンガープリントを検索テーブル 1 2 9 に格納する。

30

【 0 0 6 9 】

図 8 に検索テーブル 1 2 9 の例を示す。検索テーブル 1 2 9 のカラムのうち、フィンガープリント 1 2 9 1 にはフィンガープリントが格納される。そして重複排除ブロック 1 2 9 2 には、フィンガープリント 1 2 9 1 に格納されたフィンガープリントを持つ重複排除ブロックの位置情報が格納される。位置情報としては、重複排除ブロックの属する仮想ボリュームの識別子と、重複排除ブロックの識別子（仮想ボリューム内アドレス）の組からなる値が格納される。検索テーブル 1 2 9 の各レコードは、フィンガープリント 1 2 9 1 の値の小さい順にソートされて格納される。

40

【 0 0 7 0 】

あるレコードの重複排除ブロック 1 2 9 2 には、複数の位置情報が格納されることもある（重複データが存在するケース）。また、フィンガープリント 1 2 9 1 の値が同じレコードが複数存在することもありえる。内容の異なる複数のデータについて F P を算出した場合でも、それぞれの F P が同じになることがあり得るからである。

【 0 0 7 1 】

ストレージ装置 1 があるデータについて重複判定を行う場合、データの F P を算出し、

50

算出されたFPと同じものが、検索テーブル129のカラム“フィンガープリント1291”に格納されているレコードがあるか判定する。そのようなレコードがある場合、ストレージ装置1はさらに、そのレコードのカラム“重複排除ブロック1292”を参照して重複排除ブロック(の位置情報)を特定する。さらにストレージ装置1はマッピングテーブル127を参照して、この重複排除ブロックに割り当てられているデータブロックを特定し、当該データブロックからデータを読み出して、読み出されたデータと判定対象のデータとをバイト単位で比較することで、両者が一致しているかを判定する。

#### 【0072】

図9に、逆参照テーブル130の例を示す。逆参照テーブル130は、データブロックと重複排除ブロックのマッピング情報を管理するテーブルである。マッピングテーブル127とは異なり、逆参照テーブル130は、管理対象のデータブロックのアドレスから、そのデータブロックが割り当てられている重複排除ブロックのアドレスを特定するために用いられる。

10

#### 【0073】

逆参照テーブル130のカラムのうち、物理ページ1301とデータブロック1302にはそれぞれ、管理対象のデータブロックが属する物理ページの識別子とデータブロックの識別子が格納される。データブロックの識別子には先に述べたとおり、データブロックの属する物理ページ内の相対アドレスが用いられる。重複排除ブロック1303には、管理対象のデータブロックが割り当てられている重複排除ブロックの位置情報が格納される。重複排除ブロックの位置情報は、その重複排除ブロックの属する仮想ボリュームの識別子と重複排除ブロックの識別子の組からなる値である。

20

#### 【0074】

図10に、追記ポインタ131の例を示す。追記ポインタ131は、重複排除ブロックへのデータブロックの割り当てが必要な場合に、割り当てべきデータブロックのアドレスを管理するための情報である。追記ポインタ131は、物理ページの識別子と、データブロックの識別子の組からなる。ストレージ装置1は追記ポインタ131で指し示されているデータブロックを重複排除ブロックに割り当て、割り当てが終わった後には、追記ポインタ131の内容を、現在追記ポインタ131に格納されているアドレスの次のアドレスに更新する。物理ページの終端のデータブロックが重複排除ブロックに割り当てられた場合には、ストレージ装置1は未使用状態の物理ページを新たにデータブロック格納用に確保し、その物理ページの識別子と、その物理ページの先頭のデータブロックの識別子から成る情報で、追記ポインタ131を更新する。

30

#### 【0075】

##### (4) 処理の流れ

続いて、ストレージ装置1で行われる各種処理の流れを説明する。ストレージ装置1のシステムメモリ12には少なくとも、I/O処理部121、重複排除処理部123、重複排除解除部124、論理ページ変更部125、物理ページ解放処理部132の、5種類のプログラムが格納されている。なお、システムメモリ12に格納されているこれらのプログラムのことを、「ストレージ制御プログラム」と呼ぶこともある。以下ではこれらのプログラムによって行われる処理の説明を行う。なお、本実施例に係る図面において、参照番号の前に付されている文字列“SP”は「ステップ」を意味する。

40

#### 【0076】

まずI/O処理部121により行われる主な処理の説明を行う。I/O処理部121は、計算機5等のイニシエータデバイスに仮想ボリュームを提供し、イニシエータデバイスから受け付けたI/O要求(リード要求やライト要求)の処理を行うプログラムである。

#### 【0077】

図16を参照しながらリード処理の流れを説明する。ストレージ装置1が計算機5から仮想ボリュームに対するリード要求を受領すると、ストレージ装置1はリード要求で指定されたデータを計算機5に返送する。この処理をリード処理と呼ぶ。

#### 【0078】

50

ステップ702：ストレージ装置1が計算機5からリード要求を受領すると、I/O処理部121はリード要求に含まれるリード先アドレス(LBA)から論理ページ番号を算出することで、リード先アドレスに含まれる論理ページを特定する。また同時にI/O処理部121は、リード先アドレスに対応する論理ページ内アドレスも算出する。以下では、ここで特定された論理ページを「アクセス対象論理ページ」と呼ぶ。

【0079】

ステップ703：続いてI/O処理部121は論理ページ管理テーブル126を参照することで、アクセス対象論理ページが重複排除論理ページか判定する。この判定ではI/O処理部121は、論理ページ管理テーブル12のレコードのうち、アクセス対象論理ページに対応するレコードの状態1265が“Q”か否かを判定する。状態1265が“Q”の場合(SP703：Y)、次にステップ704が行われる。それ以外の場合には、I/O処理部121は次にステップ705を行う。

10

【0080】

ステップ704：I/O処理部121はマッピングテーブル127を参照することで、リード先アドレスで指定されている領域に割り当てられているデータブロックの位置情報(物理ページ1275とデータブロック1276の組)を特定する。そしてこの情報を用いて、リード対象データの格納されている記憶デバイス15のアドレスを求める。

【0081】

ステップ705：I/O処理部121は論理ページ管理テーブル126を参照することで、アクセス対象論理ページに割り当てられている物理ページ(の物理ページ番号)を特定する。さらに特定された物理ページ番号及びステップ702で算出した論理ページ内アドレスをもとに、リード対象データの格納されている記憶デバイス15のアドレスを求める。

20

【0082】

ステップ706：I/O処理部121は、ステップ704またはステップ705で求められた記憶デバイス15のアドレスからデータを読み出し、読み出したデータをキャッシュメモリ14に格納する。またI/O処理部121は、キャッシュメモリ14に格納されたデータを計算機5に返送し、処理を終了する。

【0083】

続いてストレージ装置1が、計算機5から仮想ボリュームに対するライト要求及びライトデータを受領した時の処理(ライト処理)の流れを、図15を用いて説明する。この処理もI/O処理部121が実行する。なお、計算機5が発行するライト要求には、ライトデータの書き込み先位置の情報(LBA及びデータ長)が含まれている。

30

【0084】

ステップ602：I/O処理部121はライト要求に含まれる書き込み先位置の情報から、ライトデータの書き込み先となる論理ページの識別子(論理ページ番号)、論理ページ内アドレスを算出する。またI/O処理部121は論理ページ管理テーブル126を参照し、ライトデータの書き込み先となる論理ページの状態1265がPでかつ物理ページ1264がNULLの場合には、この論理ページに物理ページを割り当てる処理を行う。具体的にはI/O処理部121はプール管理テーブル128を参照することで、使用状況1282が“未使用”のレコードを選択し、このレコードの論理ページ1283に、ライトデータの書き込み先となる論理ページの識別子を格納し、また使用状況1282を“論理ページ”に変更する。さらにここで選択されたレコードの物理ページ1281の値を、論理ページ管理テーブル126のレコードのうち、ライトデータの書き込み先の論理ページのレコードの物理ページ1264に格納する。

40

【0085】

ステップ603：I/O処理部121はライトデータをキャッシュメモリ14に格納する。この時、I/O処理部121はライトデータに、書き込み先論理ページの識別子と論理ページ内アドレスを付加してキャッシュメモリ14に格納する。

【0086】

50

ステップ604：I/O処理部121は、論理ページ管理テーブル126のレコードのうち、書き込み先論理ページに対応するレコードの最終ライト時刻1266に、現在時刻（ステップ604実行時点の時刻）を格納する。

【0087】

ステップ605：I/O処理部121は、論理ページ管理テーブル126のレコードの状態1265を参照し、書き込み先論理ページが状態“P”であるか否かを判定する。書き込み先論理ページが状態“P”であった場合（SP605：Y）、I/O処理部121は書き込み先論理ページの状態1265を“X”へ変更し（ステップ607）、次にステップ609を行う。

【0088】

ステップ606：I/O処理部121は、論理ページ管理テーブル126のレコードの状態1265を参照し、書き込み先論理ページが重複排除論理ページかを判定する。状態1265が“Q”だった場合、書き込み先論理ページが重複排除論理ページであったことを意味する。状態1265が“Q”だった場合（SP606：Y）、I/O処理部121は次にステップ608を行う。状態1265が“Q”でない場合（SP606：N）、ステップ608はスキップされる。

【0089】

ステップ608：I/O処理部121は重複排除解除部124を呼び出すことで、書き込み先論理ページの重複排除解除処理を行う。重複排除解除処理は重複排除解除部124が実行する。重複排除解除部124の処理は後述する。重複排除解除処理が行われることで、書き込み先論理ページは通常論理ページになる（状態1265がXになる）。

【0090】

ステップ609：I/O処理部121は、ステップ603でキャッシュメモリ14に格納したデータを記憶デバイス15に格納する。キャッシュメモリ14に一時的に格納されたデータを記憶デバイス15に格納する処理のことを「デステージ処理」と呼ぶ。I/O処理部121は、図16のステップ705と同様に、論理ページ管理テーブル126を参照することで、書き込み先論理ページに割り当てられている物理ページを特定し、さらに特定された物理ページの存在する記憶デバイス15のアドレスを特定する。そしてI/O処理部121は、この特定された記憶デバイス15のアドレスに対して、キャッシュメモリ14に格納されていたデータを書き込む。デステージ処理の後、I/O処理部121は計算機5にライト処理が完了した旨を通知し、ライト処理を終了する。

【0091】

なお、上で説明した例では、ストレージ装置1はライトスルー処理を行っている。つまりデステージ処理（ステップ609）の後で、計算機5にライト処理が完了した旨が通知される。ただし別の実施形態として、キャッシュメモリ14がライトバックキャッシュとして用いられてもよい。その場合、I/O処理部121は、キャッシュメモリ14に計算機5から受領したデータを格納した時点（ステップ603）で、計算機5にライト処理が完了した旨を通知してもよい。またこの場合、ステップ605以降の処理は、かならずしもキャッシュメモリ14にデータが格納された直後に行われなくてもよく、任意のタイミングで行われるようにしてもよい。

【0092】

続いて、論理ページ変更部125が行う処理の流れを、図14を参照しながら説明する。本実施例ではこの処理を「移行処理」と呼ぶ。移行処理は仮想ボリュームの各論理ページのうち、所定の条件を満足する通常論理ページの状態を変更する処理である。移行処理は各仮想ボリュームについて定期的に行われる。以下ではある特定の仮想ボリュームについて、ストレージ装置1が移行処理を行う時の処理の流れを説明する。

【0093】

ステップ402：論理ページ変更部125は、仮想ボリュームの論理ページのうち、まだステップ403以降の処理が行われていない論理ページを1つ選択する。ここでの論理ページの選択方法は任意であるが、たとえば論理ページ変更部125は、仮想ボリューム

10

20

30

40

50



の先頭の論理ページから順に、ステップ403以降の処理を行うとよい。

【0094】

ステップ403：論理ページ変更部125は論理ページ管理テーブル126を参照することで、選択された論理ページの状態1265が“X”か否かを判定する。状態1265が“X”の場合（ステップ403：Y）、次にステップ404が行われる。状態1265が“X”でない場合（ステップ403：N）、ステップ404～ステップ407はスキップされる。

【0095】

ステップ404：論理ページ変更部125は選択された論理ページの最終ライト時刻1266を参照し、現在時刻と最終ライト時刻1266との差を算出し、この差が所定値以上であるか判定する。この差が所定値以上の場合（ステップ404：Y）、選択された論理ページには、一定期間以上の間、計算機5からの書き込みが発生していないことを意味する。この場合には、論理ページ変更部125は次にステップ405を行う。現在時刻と最終ライト時刻1266との差が所定値未満の場合、ステップ405～ステップ407はスキップされる。

10

【0096】

ステップ405：論理ページ変更部125は重複排除処理部123を呼び出して、選択された論理ページに対して重複排除処理を行う。この処理は、通常論理ページの状態を、状態Qに変更する処理（つまり重複排除論理ページに変更する処理）である。詳細は後述する。

20

【0097】

ステップ406：論理ページ変更部125は、選択された論理ページを重複排除論理ページとして維持しておくべきか否かを判定する。具体的には論理ページ変更部125はまず、論理ページ管理テーブル126内の各レコードのうち、重複排除1263が“有効”となっているレコードの数を計数することで、ストレージ装置1内の重複排除論理ページの数特定し、この数が閾値以上であるか否かを判定する。この数が閾値未満の場合には、論理ページ変更部125は選択された論理ページは重複排除論理ページとして維持しておくべきと判断する。

【0098】

ストレージ装置1内の重複排除論理ページの数が少ない（閾値未満の）場合、重複判定対象となる重複排除ブロック数が少ない。重複判定対象となる重複排除ブロック数が少ない場合、重複排除処理による記憶領域消費量の削減効果が出ないため、ストレージ装置1は、閾値以上の数の重複排除論理ページが存在するようになるまでは、選択された論理ページを重複排除論理ページとして維持する。

30

【0099】

なお重複排除論理ページの数特定する方法は、上で挙げた方法に限定されない。別の実施形態として、ストレージ装置1は重複排除論理ページの数特定するために、システムメモリ12に重複排除論理ページの数記録する領域を予め設けておき、通常論理ページを重複排除論理ページへ、またはその逆へ変更する処理を行う度に、システムメモリ12に記録された重複排除論理ページの数増減させるようにしてもよい。

40

【0100】

重複排除論理ページの数閾値以上である場合、論理ページ変更部125はさらに、選択された論理ページの排除ブロック数1267を参照することで重複排除率を算出し、この値が所定の閾値以上か判定する。重複排除率は、“排除ブロック数1267÷論理ページ内重複排除ブロック数”で求められる。

【0101】

重複排除率が所定の閾値以上であれば、選択された論理ページは重複排除処理によって記憶領域消費量が少なくなっていることを意味するので、論理ページ変更部125は、選択された論理ページを重複排除論理ページとして維持しておくべきと判断する。

【0102】

50

一方重複排除率が閾値未満の場合には、選択された論理ページの記憶領域消費量の削減効果はあまり大きくないため、重複排除論理ページとして維持しておく必要性は低い。そのため論理ページ変更部125は、選択された論理ページを重複排除論理ページとして維持しないと判断する。選択された論理ページが、重複排除論理ページとして維持しておく必要がないと判断された場合（ステップ406：N）、次にステップ407が行われる。重複排除論理ページとして維持しておく必要があると判断された場合（ステップ406：Y）、ステップ407はスキップされる。

【0103】

ステップ407：論理ページ変更部125は重複排除解除部124を呼び出して、選択された論理ページに対して重複排除解除処理を行う。この処理は、重複排除論理ページを、状態Pの論理ページに変更する処理である。この処理の詳細は後述する。

10

【0104】

ステップ408：論理ページ変更部125は、仮想ボリュームの論理ページのうち、まだステップ403～ステップ407の処理が行われていない論理ページがあるか判定する。まだステップ403～ステップ407の処理が行われていない論理ページがある場合（ステップ408：Y）、論理ページ変更部125は再びステップ402を実行する。仮想ボリューム内の全論理ページに対してステップ403～ステップ407の処理が行われた場合には（ステップ408：N）、論理ページ変更部125は処理を終了する。

【0105】

なお、上で説明した例では、ステップ406で論理ページ変更部125は重複排除率を算出し、重複排除率に基づいて選択された論理ページを重複排除論理ページとして維持しておくか否かを判定している。ただし論理ページ変更部125は、重複排除率を算出する代わりに、排除ブロック数1267が所定の閾値以上か否かを判定することで、選択された論理ページを重複排除論理ページとして維持しておくか否か決定してもよい。排除ブロック数1267が大きければ、実質的に重複排除率も大きく、また排除ブロック数1267が小さければ、実質的に重複排除率も小さいことは明らかだからである。また、重複排除率や排除ブロック数1267以外にも、論理ページの記憶領域消費量の削減効果を推し量ることができるその他の指標値が、ステップ406の判定に用いられてもよい。

20

【0106】

続いて、ステップ405で行われる重複排除処理の流れを、図11、図12を参照しながら説明する。図11、図12は重複排除処理部123が実行する処理のフローチャートである。重複排除処理部123は、論理ページ変更部125から呼び出されることにより（図14 ステップ405）、処理を開始する。この時論理ページ変更部125は重複排除処理部123に、処理対象の論理ページの情報（具体的には論理ページの属する仮想ボリュームの識別子、論理ページ識別子）を通知する。重複排除処理部123は、この通知された論理ページの状態を変更する処理を行う。以下では論理ページ変更部125から通知された論理ページのことを、「指定された論理ページ」と呼ぶ。

30

【0107】

ステップ102：重複排除処理部123は論理ページ管理テーブル126を参照することで、指定された論理ページに割り当てられている物理ページを特定し、この物理ページに格納されているデータを読み出してキャッシュメモリ14に格納する。

40

【0108】

ステップ103：重複排除処理部123は重複排除処理を行う。ここでは、指定された論理ページ内の重複排除ブロックごとに重複排除処理を行う。以下、図12を参照しながら説明する。

【0109】

まずステップ201で重複排除処理部123は、変数kを用意し、kが指定された論理ページの先頭の重複排除ブロックを指し示すようにする（具体的にはkに、指定された論理ページの先頭の重複排除ブロックの識別子を代入する）。以下では、変数kによって指し示される重複排除ブロックのことを「選択された重複排除ブロック」と呼ぶ。

50

## 【 0 1 1 0 】

続いて重複排除処理部 1 2 3 は、ステップ 1 0 2 においてキャッシュメモリ 1 4 に格納されたデータのうち、選択された重複排除ブロックのデータを選択（特定）する（ステップ 2 0 2）。続いて重複排除処理部 1 2 3 は、特定されたデータを読み出して、このデータのフィンガープリントを算出する（ステップ 2 0 3）。

## 【 0 1 1 1 】

次にステップ 2 0 4 で、重複排除処理部 1 2 3 は検索テーブル 1 2 9 を参照し、ステップ 2 0 3 で算出されたフィンガープリントと同じ値が、フィンガープリント 1 2 9 1 に格納されているレコードを検索する。レコードがあった場合（ステップ 2 0 5 : Y）、ステップ 2 0 7 が実行される。レコードがなかった場合（ステップ 2 0 5 : N）、ステップ 2 0 6 が実行される。

10

## 【 0 1 1 2 】

ステップ 2 0 6 で重複排除処理部 1 2 3 は、ステップ 2 0 3 で算出されたフィンガープリントと、選択された重複排除ブロックの位置情報とで構成されるレコードを作成して、これを検索テーブル 1 2 9 に追加する。次いでステップ 2 0 8 が実行される。

## 【 0 1 1 3 】

ステップ 2 0 7 で、重複排除処理部 1 2 3 は、ステップ 2 0 4 の検索の結果得られたレコードのカラム“重複排除ブロック 1 2 9 2”で特定される重複排除ブロックからデータを読み出す。より詳細には、重複排除処理部 1 2 3 は、特定された重複排除ブロックに割り当てられているデータブロックを、マッピングテーブル 1 2 7 の物理ページ 1 2 7 5 とデータブロック 1 2 7 6 を参照して特定し、特定されたデータブロックからデータを読み出す。そしてこの読み出されたデータと、ステップ 2 0 2 で選択されたデータとを、バイト単位で比較し、両者が一致するか判定する。両者が一致する場合（ステップ 2 0 7 : Y）、次にステップ 2 1 0 が行われる。

20

## 【 0 1 1 4 】

もし、ステップ 2 0 4 の検索により得られた検索テーブル 1 2 9 のレコードの重複排除ブロック 1 2 9 2 に、複数の重複排除ブロックの識別子が格納されている場合は、重複排除処理部 1 2 3 は、それらの重複排除ブロックについても同様のデータの比較を行う。また、前記検索の結果、検索テーブル 1 2 9 のレコードが複数得られた場合は、重複排除処理部 1 2 3 は、各レコードにつき同様にデータの比較を行う。比較の結果、いずれの重複排除ブロックも、ステップ 2 0 2 で選択されたデータと一致しない場合には（ステップ 2 0 7 : N）、次にステップ 2 0 6 が行われる。

30

## 【 0 1 1 5 】

ステップ 2 0 8 では、重複排除処理部 1 2 3 はステップ 2 0 2 で選択されたデータを書き込むデータブロックを確保する。具体的には、重複排除処理部 1 2 3 は、追記ポインタ 1 3 1 に記録されているデータブロックのアドレスをデータ格納先として選択し、追記ポインタ 1 3 1 にデータブロックのサイズを加算する。もし、選択した追記ポインタ 1 3 1 が物理ページの最後のアドレスであった場合には、重複排除処理部 1 2 3 は、プール管理テーブル 1 2 8 を参照して未使用状態の物理ページを選択し、この物理ページをデータブロック格納用の物理ページへ変更し、追記ポインタ 1 3 1 に選択した物理ページの先頭アドレスを格納する。

40

## 【 0 1 1 6 】

次いでステップ 2 0 9 で、重複排除処理部 1 2 3 は確保したデータブロックを有する物理ページ（記憶デバイス 1 5）に、ステップ 2 0 2 で選択されたデータを書き込む。さらに、重複排除処理部 1 2 3 は、逆参照テーブル 1 3 0 に、ステップ 2 0 8 で確保したデータブロックが割り当てられる重複排除ブロックの情報を記録する。具体的には、データブロックに対応する逆参照テーブル 1 3 0 のレコードの重複排除ブロック 1 3 0 3 に、処理対象の重複排除ブロックの識別子を格納する。

## 【 0 1 1 7 】

ステップ 2 1 0 は、ステップ 2 0 2 で選択されたデータの重複データが、すでに他の重

50

重複排除ブロックに割り当てられているデータブロックに存在する場合に実行される処理である。ステップ210では重複排除処理部123は、指定された論理ページの排除ブロック数1267に1を加算する。なお、この場合、ステップ208～ステップ209と異なり、ステップ202で選択されたデータは記憶デバイス15に書き込まれない。

【0118】

ステップ211で重複排除処理部123は、ステップ207の結果一致するデータが見つかった検索テーブル129のレコードのカラム“重複排除ブロック1292”に、選択された重複排除ブロックの識別子を追加する。

【0119】

ステップ212で重複排除処理部123は、マッピングテーブル127のレコードのうち、選択された重複排除ブロックについてのレコードの更新を行う。もしステップ207の判定が肯定的な場合、つまりステップ202で選択されたデータの重複データが格納されたデータブロックがあった場合、そのデータブロックの物理ページ識別子およびデータブロックアドレスをマッピングテーブル127のレコードに登録する。この時重複排除処理部123は、このレコードのカラム“削減フラグ1277”には、重複排除ブロックが削減済み重複排除ブロックであることを示す値である“TRUE”を格納する。

【0120】

一方、ステップ202で選択されたデータの重複データが格納されたデータブロックがなかった場合（つまりステップ205またはステップ207の判定が否定的な場合）、ステップ206でデータを書き込んだ先のデータブロックの物理ページ識別子およびデータブロックの識別子をマッピングテーブル127のレコードに登録する。この登録されるレコードのカラム“削減フラグ1277”には、“FALSE”が格納される。

【0121】

また、ステップ207の判定が肯定的な場合も否定的な場合も、レコードのカラム“フィンガープリント1274”には、ステップ203で算出したフィンガープリント値が格納される。

【0122】

変数kが論理ページ内の終端アドレスに等しい場合、つまりすべての重複排除ブロックに対して、ステップ212までの処理が完了した場合（ステップ213：Y）、この処理（ステップ103）は終了する。まだステップ212までの処理が完了していない重複排除ブロックが残っている場合（ステップ213：N）、重複排除処理部123は変数kが次の重複排除ブロックを指し示すように、変数kの更新を行い（変数kに重複排除ブロックのサイズを加算する）、再びステップ202から処理を行う。以上が、ステップ103で行われる処理の内容である。

【0123】

図11の説明に戻る。ステップ104、ステップ105で、重複排除処理部123は指定された論理ページに割り当てられている物理ページの割り当てを解除する。

【0124】

ステップ104：重複排除処理部123はプール管理テーブル128のレコードのうち、割り当てられている物理ページのレコードの使用状況1282を“未使用”に変更し、論理ページ1283をNULLに変更する。

【0125】

ステップ105：重複排除処理部123は論理ページ管理テーブル126のレコードのうち、指定された論理ページのレコードの内容を変更する。具体的にはこのレコードの物理ページ1264と最終ライト時刻1266がNULLに変更され、重複排除1263には“無効”が、状態1265には“Q”が格納される。ステップ105までの処理が完了すると、重複排除処理部123は処理を終了する。

【0126】

続いて、ステップ407またはステップ608で行われる重複排除解除処理の流れを、図13を参照しながら説明する。図13は重複排除解除部124が実行する処理のフロー

10

20

30

40

50

チャートである。重複排除解除部 1 2 4 は、論理ページ変更部 1 2 5 または I / O 処理部 1 2 1 から呼び出されることにより、処理を開始する。

【 0 1 2 7 】

この時論理ページ変更部 1 2 5 または I / O 処理部 1 2 1 は、重複排除解除部 1 2 4 に処理対象の論理ページの情報（具体的には論理ページの属する仮想ボリュームの識別子、論理ページ識別子）を通知する。重複排除解除部 1 2 4 は、この通知された論理ページの状態を変更する処理を行う。以下では、重複排除解除部 1 2 4 に通知された情報により特定される論理ページのことを、「指定された論理ページ」と呼ぶ。

【 0 1 2 8 】

ステップ 3 0 2：重複排除解除部 1 2 4 はマッピングテーブル 1 2 7 を参照することで、指定された論理ページの各重複排除ブロックに割り当てられているデータブロックを特定し、特定された全てのデータブロックに格納されているデータを読み出してキャッシュメモリ 1 4 に格納する。

10

【 0 1 2 9 】

ステップ 3 0 3：重複排除解除部 1 2 4 は指定された論理ページに、未使用の物理ページを割り当てる。物理ページを割り当てる処理は、ステップ 6 0 2 で説明したものと同様である。

【 0 1 3 0 】

ステップ 3 0 4：重複排除解除部 1 2 4 は指定された論理ページに対応する、論理ページ管理テーブル 1 2 6 のレコードの内容を更新する。具体的には以下の内容更新が行われる。まずこのレコードの重複排除 1 2 6 3 は“無効”に変更され、排除ブロック数 1 2 6 7 は 0 に更新される。またステップ 3 0 3 が実行された時点で、物理ページ 1 2 6 3 には割り当てられた物理ページの識別子が格納されているので、ここでは物理ページ 1 2 6 3 の更新は行われぬ。

20

【 0 1 3 1 】

ステップ 3 0 5：重複排除解除部 1 2 4 は、ステップ 3 0 2 でキャッシュメモリ 1 4 に格納したデータを、ステップ 3 0 3 で割り当てられた物理ページにデステージする。

【 0 1 3 2 】

ステップ 3 0 6：重複排除解除部 1 2 4 はマッピングテーブル 1 2 7 のレコードのうち、指定された論理ページに含まれる各重複排除ブロックについての情報を管理しているレコードの内容を更新する。ここでは、各レコードの物理ページ 1 2 7 5 およびデータブロック 1 2 7 6 は、NULL に変更される。なお、この時点ではフィンガープリント 1 2 7 4 および削減フラグ 1 2 7 7 は更新されず、元の値を保持しつづける。

30

【 0 1 3 3 】

ステップ 3 0 7：重複排除解除部 1 2 4 は、指定された論理ページの各重複排除ブロックに割り当てられていたデータブロックを有する物理ページを特定する。具体的には、重複排除解除部 1 2 4 は、マッピングテーブル 1 2 7 のレコードのうち、指定された論理ページの識別子と“論理ページ 1 2 7 2”のカラムの情報とが一致し、かつ削減フラグ 1 2 7 7 が FALSE であるレコードを検索し、検索の結果得られたレコードから物理ページ 1 2 7 5 の情報を重複なく取り出す。ステップ 3 0 7 においては、1 個以上の物理ページが特定される場合がある。またここでは、削減フラグ 1 2 7 7 が FALSE であるレコードだけが検索されるので、削減済み重複排除ブロックは検索結果から除外される。

40

【 0 1 3 4 】

ステップ 3 0 8：重複排除解除部 1 2 4 は、ステップ 3 0 7 で特定した物理ページを未割当状態に変更する。本実施例では、データブロック格納用に使われている物理ページを未割当状態に変更する処理を、「物理ページ解放処理」と呼ぶ。また、物理ページを未割当状態にすることを、「物理ページを解放する」と表現する場合がある。

【 0 1 3 5 】

ステップ 3 0 8 の実行前時点において、解放する対象の物理ページには、依然として複数の重複排除ブロックに共有されているデータブロックがある場合がある。このようなデ

50

ータブロックを、以降では共有データブロックと呼ぶ。また、共有データブロックのデータを共有データと呼ぶ。

【0136】

物理ページ解放処理は、物理ページを未割当状態へ変更する前に、その物理ページの共有データを、別の物理ページのデータブロックへ移動（複製）する。そして、その共有データブロックが割り当てられている重複排除ブロックを、複製先のデータブロックが割り当てられた状態へ変更する。その後、対象の物理ページを未割当状態に変更する。

【0137】

物理ページ解放処理は、物理ページ解放処理部132により実行される。ステップ308において、重複排除処理部124は、ステップ307で特定した物理ページの識別子のうち、追記ポインタ131に記録されている物理ページを除いたものを、物理ページ解放処理部132に通知する。物理ページ解放処理部132は、通知された物理ページを対象に、物理ページ解放処理を実行する。

10

【0138】

物理ページ解放処理部132における処理の流れを、図17を参照しながら説明する。図17は、物理ページ解放処理部132が実行する処理のフローチャートである。

【0139】

ステップ802：物理ページ解放処理部132は、通知された物理ページの識別子の内のひとつを選択する。この選択した物理ページの識別子により特定される物理ページを、処理対象の物理ページと呼ぶ。

20

【0140】

ステップ803：物理ページ解放処理部132は、処理対象の物理ページに含まれるデータブロックを1つ選択する。以下では、ステップ803で選択したデータブロックを処理対象のデータブロックと呼ぶ。

【0141】

ステップ804：物理ページ解放処理部132は、処理対象のデータブロックに対応している逆参照テーブル130のレコードを参照し、重複排除ブロック1303に格納された重複排除ブロックの識別子より重複排除ブロックを特定する。ここで特定された重複排除ブロックは、ステップ209（重複排除処理）により処理対象のデータブロックが割り当てられた重複排除ブロックである。

30

【0142】

ステップ806：物理ページ解放処理部132は、ステップ804で特定された重複排除ブロックに対応しているマッピングテーブル127のレコードを参照し、フィンガークリント1274の値を得る。

【0143】

ステップ807：物理ページ解放処理部132は、処理対象のデータブロックが共有データブロックである場合には、そのデータを別の物理ページのデータブロックへ移動（複製）し、処理対象のデータブロックを共有している重複排除ブロックが複製先のデータブロックを共有するように変更する。ステップ807の詳細については後述する。

【0144】

ステップ808：物理ページ解放処理部132は、処理対象の物理ページに含まれる全てのデータブロックについて、ステップ804～807を実行したどうかを判定する。ステップ804～807を実行していないデータブロックがある場合（ステップ808：Y）、物理ページ解放処理部132は、処理対象のデータブロックの次のデータブロックを選択し（ステップ803）、ステップ804～807を実行する。処理対象の物理ページに含まれる全てのデータブロックについて、ステップ804～807を実行済みであれば（ステップ808：N）、ステップ809が実行される。

40

【0145】

ステップ809：物理ページ解放処理部132は、処理対象の物理ページに対応するプール管理テーブル128のレコードの使用状況1282を、“未使用”に変更する。

50

## 【 0 1 4 6 】

ステップ 8 1 0 : 物理ページ解放処理部 1 3 2 は、処理対象の物理ページに含まれる全てのデータブロックについて、それぞれ逆参照テーブル 1 3 0 の対応するレコードの重複排除ブロック 1 3 0 3 に N U L L を格納する。

## 【 0 1 4 7 】

ステップ 8 1 1 : 物理ページ解放処理部 1 3 2 は、重複排除解除部 1 2 4 から通知された全ての物理ページについて、ステップ 8 0 3 ~ 8 1 1 を実行したかどうかを判定する。ステップ 8 0 3 ~ 8 1 1 を実行していない物理ページがある場合 (ステップ 8 1 1 : Y)、物理ページ解放処理部 1 3 2 はそのような物理ページのの一つを選択し (ステップ 8 0 2)、ステップ 8 0 3 からのステップを実行する。通知された物理ページの全てについて、  
10  
ステップ 8 0 3 ~ 8 1 1 を実行済みであれば (ステップ 8 1 1 : N)、物理ページ解放処理を終了する。

## 【 0 1 4 8 】

図 1 8 を参照しながら、ステップ 8 0 7 の共有データ判定・複製処理の流れを説明する。図 1 8 は、共有データ判定・複製処理のフローチャートである。

## 【 0 1 4 9 】

ステップ 9 0 2 : 物理ページ解放処理部 1 3 2 は、内部変数「移動済み」を用意し、初期値として F A L S E を設定する。

## 【 0 1 5 0 】

ステップ 9 0 3 : 物理ページ解放処理部 1 3 2 は、検索テーブル 1 2 9 を検索し、ステップ 8 0 6 で得たフィンガープリントを“フィンガープリント 1 2 9 1”のカラムに含み、かつステップ 8 0 4 で得た重複排除ブロックの識別子を“重複排除ブロック 1 2 9 2”のカラムに含むレコードを特定する。  
20

## 【 0 1 5 1 】

ステップ 9 0 4 : 物理ページ解放処理部 1 3 2 は、特定した検索テーブル 1 2 9 のレコードの重複排除ブロック 1 2 9 2 に格納されている重複排除ブロックの識別子を選択する。  
。

## 【 0 1 5 2 】

ステップ 9 0 5 : 物理ページ解放処理部 1 3 2 は、ステップ 9 0 4 で選択した重複排除ブロックの識別子に対応するマッピングテーブル 1 2 7 のレコードを参照し、レコードに格納された各種情報を得る。  
30

## 【 0 1 5 3 】

ステップ 9 0 6 : 物理ページ解放処理部 1 3 2 は、処理対象のデータブロックが、ステップ 9 0 4 で選択した重複排除ブロックに割り当てられているか否かを判定する。割り当てられている場合には、処理対象のデータブロックは共有データブロックであることになる。判定の具体的な方法は以下のとおりである。すなわち、ステップ 9 0 5 で得られたレコードの物理ページ 1 2 7 5 およびデータブロック 1 2 7 6 で特定されるデータブロックが、処理対象のデータブロックと一致するならば、割り当てられていると判定される。一致する場合 (ステップ 9 0 6 : Y)、そのデータを別の物理ページのデータブロックへ移動 (複製) し、重複排除ブロックに複製先のデータブロックを割り当てる処理を行う。この処理は、ステップ 9 0 7 のマッピング切替処理で実行される。マッピング切替処理の流れについては後述する。一致しない場合 (ステップ 9 0 6 : N)、次にステップ 9 0 8 が実行される。  
40

## 【 0 1 5 4 】

なお、ステップ 9 0 4 で選択した重複排除ブロックが、重複排除解除処理を適用した論理ページに含まれる場合、ステップ 9 0 6 の判定は「一致しない」となる。なぜなら、重複排除解除処理のステップ 3 0 6 で、重複排除ブロックに対応するマッピングテーブル 1 2 7 のレコードの物理ページ 1 2 7 5 およびデータブロック 1 2 7 6 には、N U L L が格納されるからである。

## 【 0 1 5 5 】

ステップ908：物理ページ解放処理部132は、ステップ903で特定した検索テーブル129のレコードの重複排除ブロック1292から、ステップ904で選択した重複排除ブロックの識別子を削除する。この結果、重複排除ブロック1292に重複排除ブロックの識別子が登録されていない状態、つまり重複排除ブロック1292が空の状態となった場合は、当該レコードを検索テーブル129から削除する。

【0156】

ステップ909：物理ページ解放処理部132は、ステップ905で参照したマッピングテーブル127のレコードのフィンガープリント1274にNULLを格納する。また、同じレコードの削減フラグ1277にFALSEを格納する。

【0157】

ステップ910：物理ページ解放処理部132は、ステップ903で特定した検索テーブル129のレコードの重複排除ブロック1292に格納されている重複排除ブロックの識別子のうち、ステップ905～909を実行していないものがあるかどうかを判定する。ステップ905～909をまだ実行していない重複排除ブロックの識別子がある場合（ステップ910：Y）、物理ページ解放処理部132はそのうちのひとつを選択し（ステップ904）、ステップ905からの処理を実行する。そうでない場合（ステップ910：N）、処理を終了する。

【0158】

図19を参照しながら、ステップ907のマッピング切替処理の流れを説明する。図19は、マッピング切替処理のフローチャートである。

【0159】

マッピング切替処理は、ステップ803で選択されたデータブロックのデータを、別の物理ページのデータブロックに移動（複製）する処理である。

【0160】

ステップ1002：物理ページ解放処理部132は、ステップ902で用意した内部変数「移動済みフラグ」の値がFALSE場合（ステップ1002：N）、ステップ1003を実行する。移動済みフラグがTRUEである場合（ステップ1002：Y）、ステップ1003～ステップ1006はスキップされる。

【0161】

ステップ1003：物理ページ解放処理部132は、移動済みフラグにTRUEを設定する。

【0162】

ステップ1004：物理ページ解放処理部132は、ステップ803で選択したデータブロックのデータの複製先となるデータブロックを確保する。データブロックの確保の処理内容は、ステップ208（重複排除処理）と同様である。

【0163】

ステップ1005：物理ページ解放処理部132は、ステップ803で選択したデータブロックのデータを読み出し、ステップ1003で確保したデータブロックに書き込む。この処理は、ステップ209（重複排除処理）と同様である。

【0164】

ステップ1006：物理ページ解放処理部132は、ステップ904で選択した重複排除ブロックを、削減済み重複排除ブロックから、通常のリバース重複排除ブロックへ変更する。このため、ステップ1006では、物理ページ解放処理部132は、当該重複排除ブロックを含む論理ページの排除ブロック数（論理ページ管理テーブル126の排除ブロック数1267）を、1減じる。

【0165】

ステップ1007：物理ページ解放処理部132は、ステップ904で選択した重複排除ブロックに対応するマッピングテーブル127のレコードを更新する。変更内容は以下のとおりである。まず、物理ページ1275とデータブロック1276には、ステップ1004で確保したデータブロックの物理ページおよびデータブロックの識別子がそれぞれ

10

20

30

40

50



格納される。次に、ステップ1002において移動済みフラグがFALSEであった場合、すなわち、削減済み重複排除ブロックを通常のリダイレクトブロックへ変更した場合、削減フラグ1277の内容がFALSEに変更される。

【0166】

以上、本発明の実施例を説明したが、これらは、本発明の説明のための例示であって、本発明の範囲をこれらの実施例にのみ限定する趣旨ではない。すなわち、本発明は、他の種々の形態でも実施する事が可能である。

【0167】

上では、ストレージ装置が使用するいくつかの情報が、論理ページ管理テーブル126等のように、テーブル構造で管理される例を説明したが、テーブル構造で情報を管理する 10  
態様に限定されるわけではない。ストレージ装置はテーブル以外のデータ構造、例えばリスト構造などを用いて、情報を管理してもよい。

【0168】

また、ストレージ装置は重複排除処理に加えて圧縮処理を行ってもよい。たとえばストレージ装置に搭載される記憶デバイスとして、データを圧縮して格納する機能（圧縮機能）を有する記憶デバイスを用いることで、記憶デバイスに格納されるデータを圧縮するようにしてもよい。その際、通常論理ページに書き込まれたデータに対するアクセス性能を維持するために、重複排除論理ページに書き込まれたデータだけが、圧縮機能を有する記憶デバイスに格納されるようにするとよい。さらに別の実施形態として、ストレージ装置は、ストレージコントローラのCPUが、記憶デバイスに格納されるデータの圧縮を行う 20  
ように構成されていてもよい。

【0169】

また、上では重複排除論理ページへデータが書き込まれたとき、その重複排除論理ページを通常のリダイレクトページへ変更した上で、データをドライブへ格納する例を説明したが、重複排除論理ページを通常のリダイレクトページへ変更せずに、書き込まれたデータをドライブへ格納するようにしてもよい。例えば、重複排除論理ページへデータが書き込まれたとき、物理ページの領域を確保してデータを格納し、重複排除論理ページのアドレスと物理ページのアドレスのマッピング関係を更新するようにしてもよい。この場合、重複排除論理ページへライトされたデータの量、ライトの回数、それらの一定期間あたりの量などの指標から、重複排除論理ページを通常リダイレクトページへ戻すか否かを判断するようにしてもよい。 30

【0170】

また、上では重複排除論理ページを通常のリダイレクトページへ変更したときに、その論理ページのデータが変更前に格納されていた物理ページを未使用の状態に変更する例を説明したが、重複排除論理ページを通常のリダイレクトページへ変更した後で、任意の契機で物理ページを未使用の状態に変更するようにしてもよい。例えば、ストレージ装置は、未使用の状態の物理ページの数、論理ページへ割り当てられている物理ページの数、重複排除論理ページのデータを格納するために割り当てられた物理ページの数、重複排除論理ページの数、ストレージ装置を構成するハードウェアの処理負荷などの指標から、物理ページを未使用の状態に変更する契機を決定するようにしてもよい。また、ストレージ装置は、物理ページごとに、格納したデータの量や、重複排除ブロックへの割り当てを解除した量等の情報を 40  
記録するようにしてもよい。そして、これらの物理ページごとの情報に基づいて、未使用の状態に変更する物理ページを選択するようにしてもよい。

【0171】

また、ストレージ装置のハードウェア構成は、上で説明した構成に限定されない。たとえば、上の実施例で述べた計算機に1以上の記憶デバイスを搭載した装置を、ストレージ装置として用いてもよい。そして、上の実施例で説明した各プログラムを計算機のプロセッサに実行させれば、上の実施例で説明したストレージ装置と同じことが実現できる。

【0172】

また、上で説明した実施例では、アクセス要求元の計算機5とストレージ装置1とが異なるハードウェアである例が説明されたが、別の実施形態として、計算機5とストレージ 50

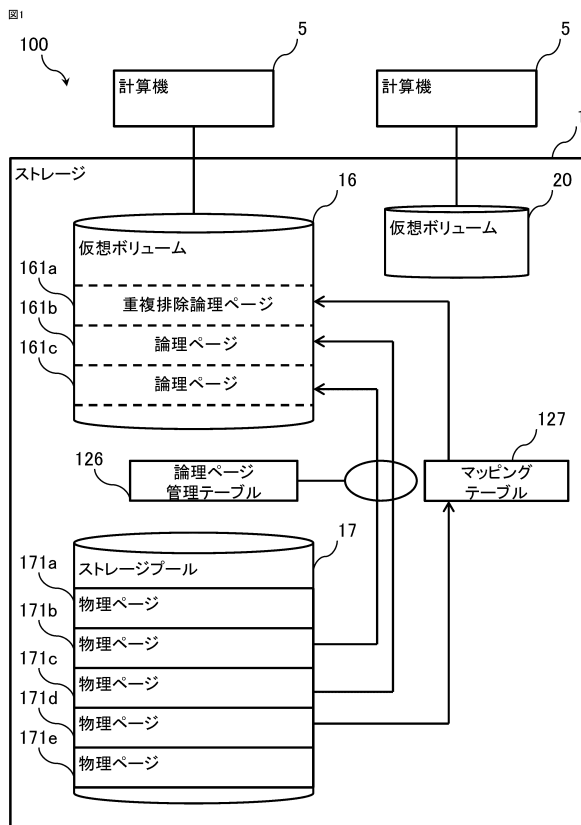
装置 1 とが単一のハードウェアとして実装されてもよい。つまり、上の実施例で説明した、ストレージ装置 1 で実行されるストレージ制御プログラムと、計算機 5 で実行されるアプリケーションプログラムとが、同一の計算機上で実行されてもよい。この場合、アプリケーションプログラムはアクセス要求元として、ストレージ制御プログラムに対して I/O 要求を発行し、ストレージ制御プログラムは I/O 要求に対する応答（I/O 要求がリード要求の場合にはリードデータ）をアクセス要求元であるアプリケーションプログラムに返却するように構成されていても良い。好ましくは、計算機上で、仮想計算機を形成するためのプログラム（ハイパーバイザ等）を実行させることで、アプリケーションプログラムを実行する仮想計算機と、I/O 処理部 121 や重複排除処理部 123 等のプログラムを実行する仮想計算機とを形成するとよい。

【符号の説明】

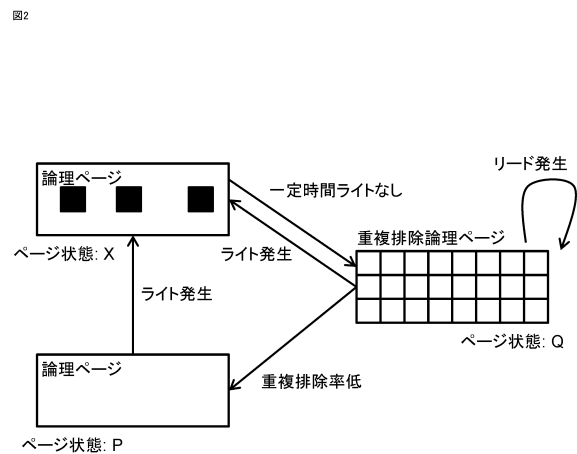
【0173】

1：ストレージ装置， 5：計算機， 6：SAN， 10：ストレージコントローラ，  
 11：CPU， 12：システムメモリ， 14：キャッシュメモリ， 15：記憶デバイス，  
 121：I/O処理部， 123：重複排除処理部， 124：重複排除解除部，  
 125：論理ページ変更部， 126：論理ページ管理テーブル， 127：マッピングテーブル，  
 128：プール管理テーブル， 129：検索テーブル， 130：逆参照テーブル，  
 131：追記ポインタ， 132：物理ページ解放処理部

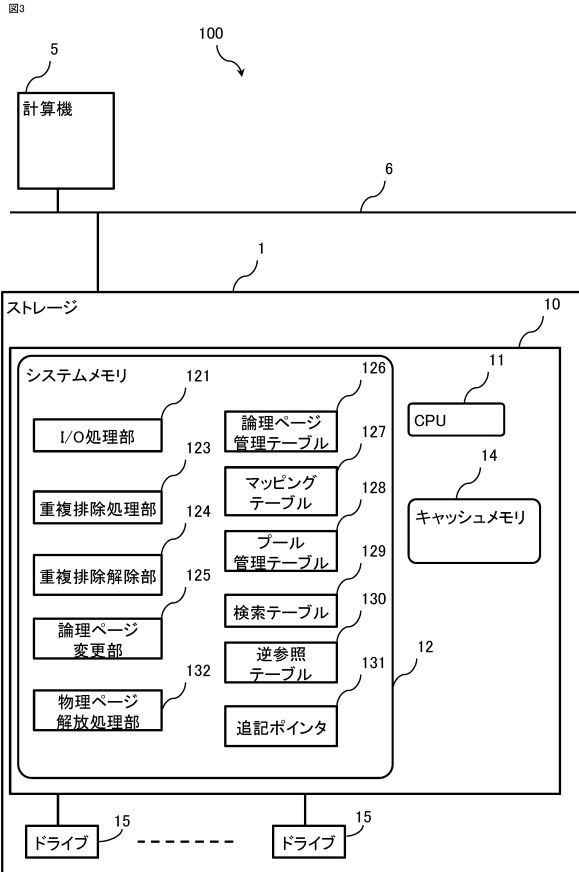
【図 1】



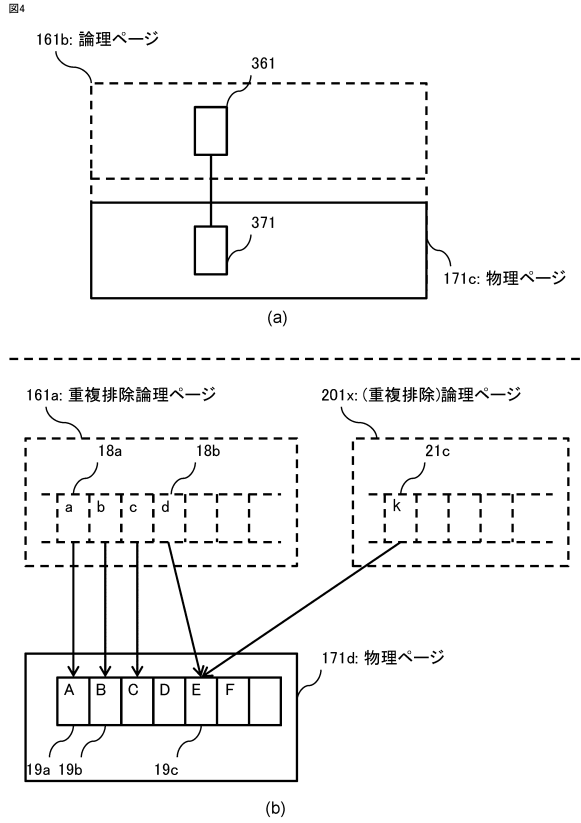
【図 2】



【図3】



【図4】



【図5】

図5  
126: 論理ページ管理テーブル

仮想ボリューム	論理ページ	重複排除	物理ページ	状態	最終ライト時刻	排除ブロック数
16	161a	有効	NULL	Q	NULL	0
16	161b	無効	171b	X	12/10 14:30	0
16	161c	無効	171c	P	12/15 20:00	0
20	201a	有効	NULL	Q	NULL	100

1261

1262

1263

1264

1265

1266

1267

【図7】

図7  
128: プール管理テーブル

物理ページ	使用状況	論理ページ
171a	未使用	NULL
171b	論理ページ	16b
171c	論理ページ	16c
171d	データブロック	NULL
171e	未使用	NULL

1281

1282

1283

【図6】

図6  
127: マッピングテーブル

仮想ボリューム	論理ページ	重複排除ブロック	フィンガープリント	物理ページ	データブロック	削減フラグ
...	...	...	...	...	...	...
16	161a	a	11111111	171d	A	FALSE
16	161a	b	22222222	171d	B	FALSE
16	161a	c	33333333	171d	C	FALSE
16	161a	d	44444444	171d	E	FALSE
...	...	...	...	...	...	...
20	201x	k	44444444	171d	E	TRUE
...	...	...	...	...	...	...

1271

1272

1273

1274

1275

1276

1277

【図8】

図8  
129: 検索テーブル

フィンガープリント	重複排除ブロック
...	...
11111111	16:a
...	...
22222222	16:b
...	...
33333333	16:c
...	...
44444444	16:d, 20:k
...	...

1291

1292

【図9】

図9

130: 逆参照テーブル

物理ページ	データブロック	重複排除ブロック
...	...	...
171d	A	16:a
171d	B	16:b
171d	C	16:c
171d	E	16:d
...	...	...

1301

1302

1303

【図10】

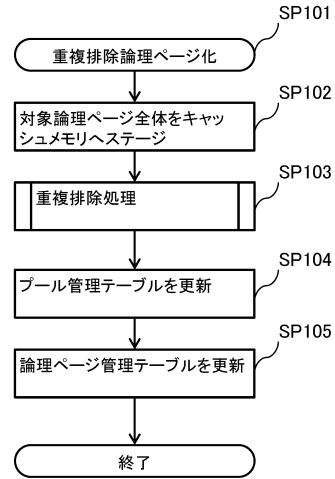
図10

131: 追記ポインタ

171d:F

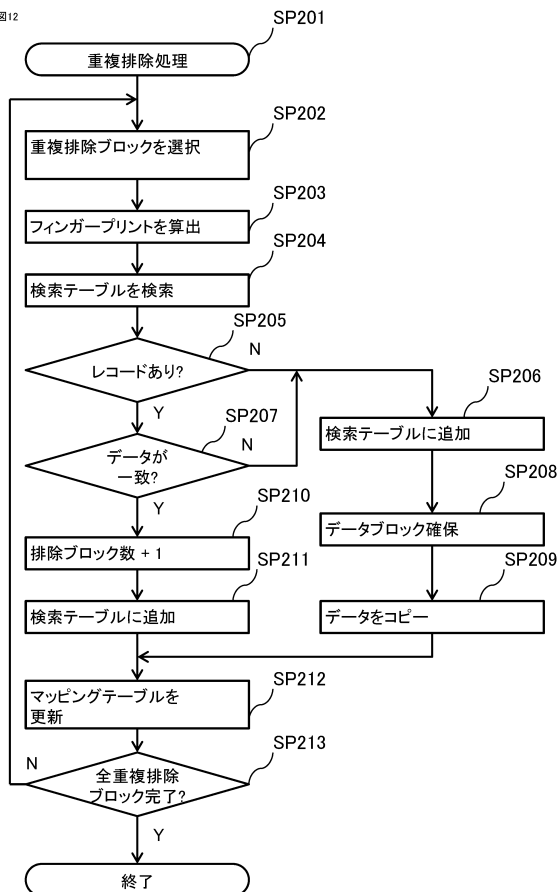
【図11】

図11



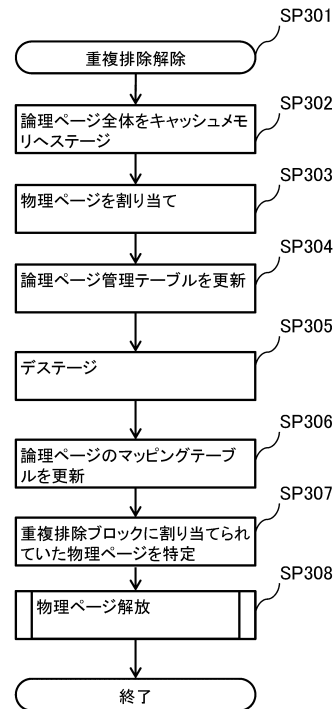
【図12】

図12



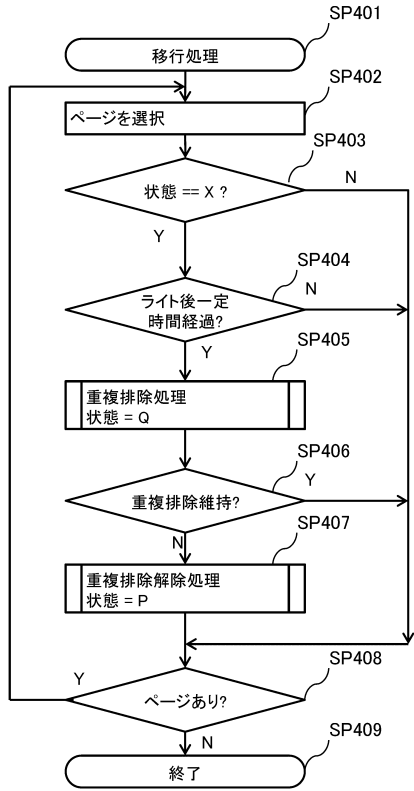
【図13】

図13



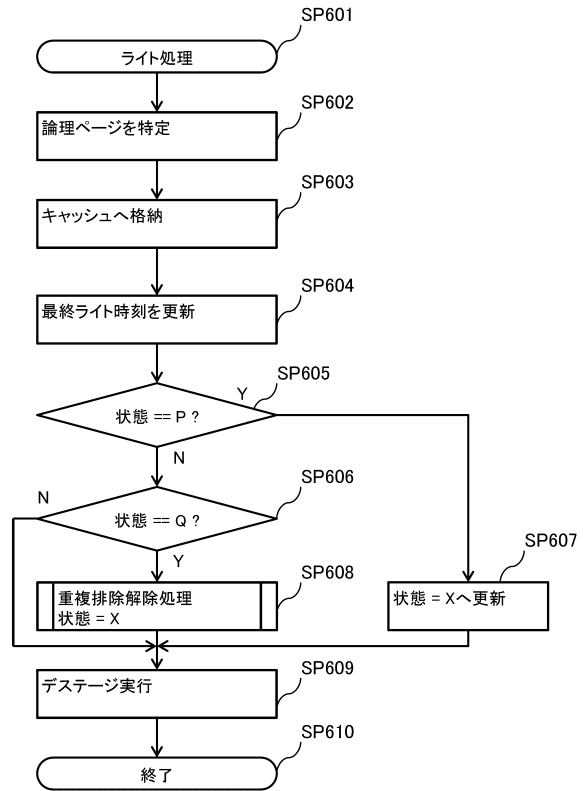
【図14】

図14



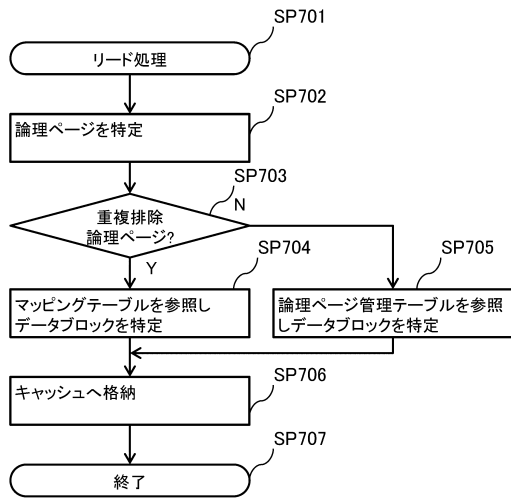
【図15】

図15



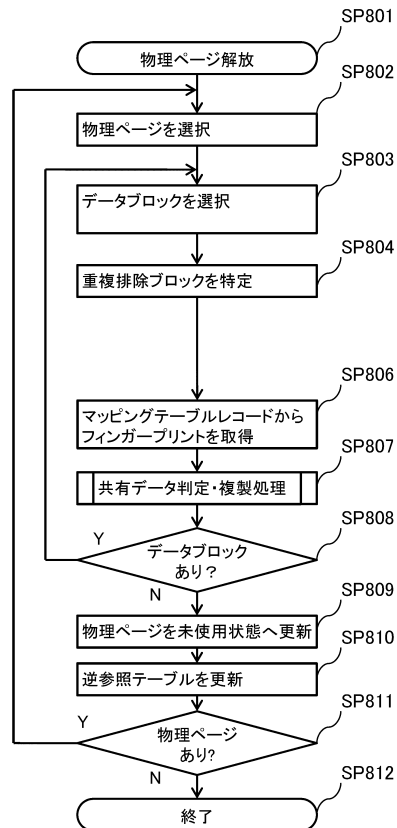
【図16】

図16



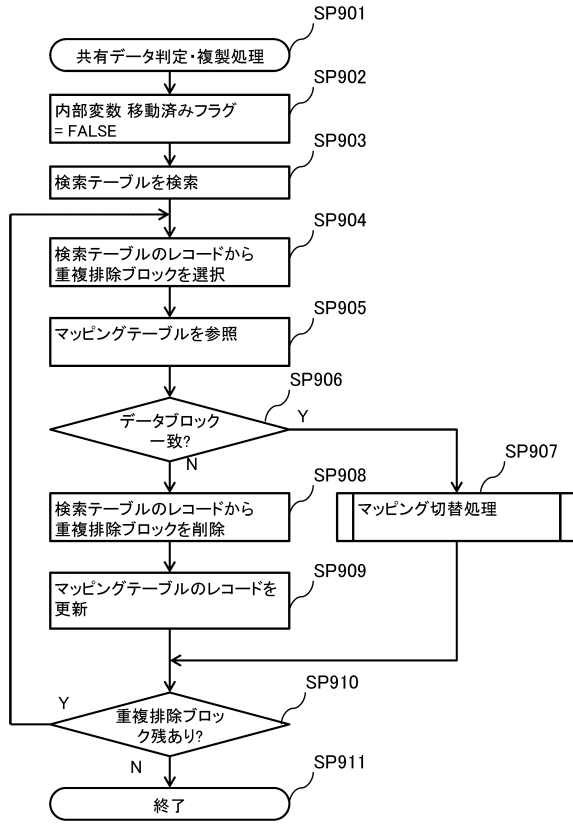
【図17】

図17



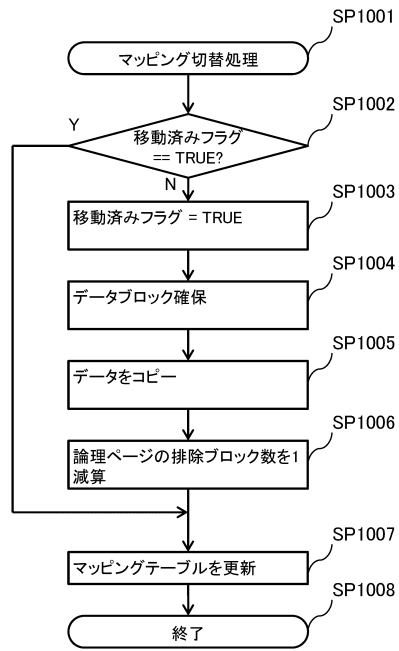
【図18】

図18



【図19】

図19



---

フロントページの続き

- (72)発明者 渡邊 恭男  
東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
- (72)発明者 吉井 義裕  
東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
- (72)発明者 松上 一樹  
東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内

審査官 田中 啓介

- (56)参考文献 特開2011-065314(JP,A)  
国際公開第2014/157243(WO,A1)  
特表2012-523023(JP,A)  
特表2015-528928(JP,A)  
特開2009-093571(JP,A)  
特開2008-234158(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F3/06-3/08  
G06F12/00-12/06  
G06F13/10-13/18  
G06F16/00-16/958