

(12) 特許協力条約に基づいて公開された国際出願

(19) 世界知的所有権機関  
国際事務局

(43) 国際公開日  
2016年2月4日(04.02.2016)



(10) 国際公開番号  
WO 2016/017002 A1

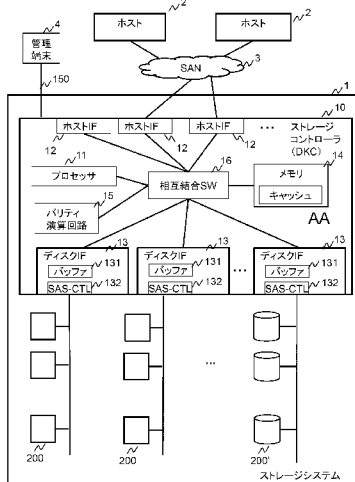
- (51) 国際特許分類:  
G06F 3/06 (2006.01)
- (21) 国際出願番号: PCT/JP2014/070224
- (22) 国際出願日: 2014年7月31日(31.07.2014)
- (25) 国際出願の言語: 日本語
- (26) 国際公開の言語: 日本語
- (71) 出願人: 株式会社日立製作所 (HITACHI, LTD.)  
[JP/JP]; 〒1008280 東京都千代田区丸の内一丁目  
6番6号 Tokyo (JP).
- (72) 発明者: 松井 章(MATSUI, Akira); 〒1008280 東京  
都千代田区丸の内一丁目6番6号 株式会社日  
立製作所内 Tokyo (JP).
- (74) 代理人: 特許業務法人第一国際特許事務所(PAT-  
ENT CORPORATE BODY DAI-ICHI KOKUSAI  
TOKKYO JIMUSHO); 〒1080014 東京都港区芝4  
丁目10番5号 Tokyo (JP).
- (81) 指定国 (表示のない限り、全ての種類の国内保  
護が可能): AE, AG, AL, AM, AO, AT, AU, AZ, BA,  
BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN,  
CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES,  
FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN,  
IR, IS, JP, KE, KG, KN, KP, KR, KZ, LA, LC, LK, LR,  
LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX,  
MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH,  
PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK,  
SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA,  
UG, US, UZ, VC, VN, ZA, ZM, ZW.
- (84) 指定国 (表示のない限り、全ての種類の広域保  
護が可能): ARIPO (BW, GH, GM, KE, LR, LS, MW,  
MZ, NA, RW, SD, SL, SZ, TZ, UG, ZM, ZW), ユーラシ  
ア (AM, AZ, BY, KG, KZ, RU, TJ, TM), ヨーロッパ  
(AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR,  
GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT,  
NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OAPI  
(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML,  
MR, NE, SN, TD, TG).

[続葉有]

(54) Title: STORAGE SYSTEM

(54) 発明の名称: ストレージシステム

81



(57) Abstract: A storage system according to the present invention comprises a controller and a plurality of storage devices, wherein (n + m) storage devices constitute a RAID group, write data from a host computer is stored in the n storage devices, and redundant data generated from n data is stored in the m storage devices. When any fault has occurred in at least one of the storage devices, the controller reads compressed data and redundant data from each of the storage devices in which no fault has occurred among the storage devices constituting the RAID group, and transmits the read compressed data to the storage device for which data restoration is to be performed.

(57) 要約: 本発明のストレージシステムは、コントローラと、複数の記憶デバイスとを備え、(n + m) 台の記憶デバイスにより RAID グループを構成し、ホスト計算機からのライトデータを n 台の記憶デバイスに、n 個のデータから生成された冗長データを m 個の記憶デバイスに格納している。少なくとも 1 台の記憶デバイスに障害が発生した場合、コントローラは RAID グループを構成する記憶デバイスの中で障害の発生していない記憶デバイスの各々から、圧縮された状態のデータと冗長データを読み出して、読み出された圧縮状態のデータを、データ復旧先の記憶デバイスに送信する。

- 1 Storage system
- 2 Host
- 4 Management terminal
- 10 Storage controller (DKC)
- 11 Processor
- 12 Host IF
- 13 Disk IF
- 14 Memory
- 15 Parity calculation circuit
- 16 Mutual connection SW
- 131 Buffer
- AA Cache

WO 2016/017002 A1

添付公開書類:

— 国際調査報告 (条約第 21 条(3))

## 明 細 書

発明の名称：ストレージシステム

### 技術分野

[0001] 本発明は、ストレージシステムにおけるデータ復旧技術に関する。

### 背景技術

[0002] ストレージシステムの中には、いわゆるRAID (Redundant Arrays of Inexpensive/Independent Disks) 技術を用いて、システムを高可用化しているものが多い。RAID技術とは、ホスト計算機などの上位装置から受け付けたライトデータからパリティなどの冗長データを算出し、ライトデータとパリティとをそれぞれ異なる記憶デバイスに分散格納する技術である。RAID技術を採用することで、一部の記憶デバイスに障害が発生し、その記憶デバイスからデータを読み出せなくなった場合でも、その他の記憶デバイスに格納されている情報を用いて、データを再生成することができる。

[0003] RAID技術ではパリティ算出のために、ストレージシステムに搭載されるコントローラ (ストレージコントローラ) の処理負荷、あるいはストレージシステム内の構成要素間 (たとえばストレージコントローラと記憶デバイス間) のデータ転送量が増加する。処理負荷やデータ転送量の増加を抑止するために、従来から様々な技術が考えられてきた。たとえば特許文献1には、ストレージコントローラと記憶デバイスの間で発生するデータ転送量を抑制するために、記憶デバイス側にパリティ生成等の機能を有するストレージシステムが開示されている。

### 先行技術文献

#### 特許文献

[0004] 特許文献1：米国特許出願公開第2013/0290629号明細書

### 発明の概要

#### 発明が解決しようとする課題

[0005] RAID技術の特徴として、先に述べたとおり、障害が発生した記憶デバイスに格納されていたデータを復旧（再生）できる点が挙げられる。ただしデータの復旧のためには、障害が発生していない記憶デバイスに格納されているデータをすべて読み出し、読み出されたデータに対して所定の演算を施すことでデータを再生成し、再生成されたデータを、新たな記憶デバイス（スペアドライブまたはスペアデバイス）に書き込む、という処理を行う必要がある、これらの処理の過程では、ストレージコントローラと記憶デバイス間で大量のデータ転送が発生する。

[0006] 特に近年、記憶デバイスの記憶容量が増大しており、記憶デバイス内のデータをすべて読み出すだけでも、長時間を要する。そのため、RAID技術を用いたデータ復旧処理も長時間化する傾向がある。データ復旧処理中は冗長性がない状態であるため、復旧処理中に別の記憶デバイスに障害が発生した場合、データ復旧が不可能になる。特許文献1には、記憶デバイス側に設けられたパリティ生成機能を用いてデータ復旧を行うことが開示されているが、記憶デバイス内のデータをすべて読み出す必要があることには変わりはなく、データ復旧処理の時間を大幅に削減することは困難である。

### 課題を解決するための手段

[0007] 本発明の一態様に係るストレージシステムは、コントローラと、複数の記憶デバイスとを備え、 $(n+m)$  台の記憶デバイスによりRAIDグループを構成し、ホスト計算機からのライトデータを  $n$  台の記憶デバイスに、 $n$  個のデータから生成された冗長データを  $m$  個の記憶デバイスに格納している。少なくとも1台の記憶デバイスに障害が発生した場合、コントローラはRAIDグループを構成する記憶デバイスの中で障害の発生していない記憶デバイスの各々から、圧縮された状態のデータと冗長データを読み出して、読み出された圧縮状態のデータを、データ復旧先の記憶デバイスに送信する。

### 発明の効果

[0008] 本発明により、データ復旧時の転送データ量を削減することができ、データ復旧時間を短縮することができる。

## 図面の簡単な説明

- [0009] [図1]本発明の実施例に係るストレージシステムの構成図である。
- [図2]FMPKの構成図である。
- [図3]RAIDグループ内の記憶領域の内容の説明図である。
- [図4]RAIDグループの記憶空間と論理ユニット(LU)の関係を説明する図である。
- [図5]デバイス管理テーブルの構成を示す図である。
- [図6]RG管理テーブルの構成を示す図である。
- [図7]LU管理テーブルの構成を示す図である。
- [図8]FMPK内の記憶領域の管理方法の説明図である。
- [図9]マッピングテーブルの構成を示す図である。
- [図10]状態管理テーブルの構成を示す図である。
- [図11]ホストからのライトデータにDIFが付加される過程を表した概念図である。
- [図12]圧縮Readコマンドのフォーマットを表した図である。
- [図13]圧縮Readコマンド3000に対する応答情報のフォーマットを表した図である。
- [図14]圧縮コピーWriteコマンドのフォーマットを表した図である。
- [図15]圧縮パリティ演算Writeコマンドのフォーマットを表した図である。
- [図16]圧縮コピーWriteコマンドに対する応答情報のフォーマットを表した図である。
- [図17]データ復旧処理の全体の流れを表したフローチャートである。
- [図18]コピー管理テーブルの構成を示す図である。
- [図19]コピー復旧処理の詳細フローチャートである。
- [図20]コレクションコピーの詳細フローチャートである。
- [図21]圧縮Readコマンドを受信した時に行われる処理の流れ図である。
- [図22]圧縮コピーWriteコマンドを受信した時に行われる処理の流れ図

である。

[図23]圧縮パリティ演算Writeコマンドを受信した時に行われる処理の流れ図である。

[図24]ステージング情報管理テーブルの構成を示す図である。

[図25]圧縮中間パリティ演算コマンドのフォーマットを表した図である。

[図26]パリティコミットコマンドのフォーマットを表した図である。

[図27]圧縮中間パリティ演算コマンドを受信した時に行われる処理の流れ図である。

[図28]パリティコミットコマンドを受信した時に行われる処理の流れ図である。

[図29]実施例2に係るストレージシステムで行われるコレクションコピーの詳細フローチャートである。

[図30]実施例3に係るストレージシステムでサポートされるストライプラインの例を表した図である。

[図31]実施例3に係るストレージシステムで行われるコレクションコピーの詳細フローチャートである。

[図32]実施例3に係るストレージシステムでサポートされる圧縮パリティ演算Writeコマンドのフォーマットを表した図である。

### 発明を実施するための形態

[0010] 以下、図面を参照して、本発明の一実施形態に係るストレージシステムについて説明する。なお、本発明は、以下に説明する実施形態に限定されるものではない。

#### 実施例 1

[0011] 図1は、実施例1に係るストレージシステム1の構成を示す。ストレージシステム1は、ストレージコントローラ(DKC)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(200、200')を有する。

[0012] 記憶デバイス200、200'は、ホスト2などの上位装置からのライト

データを格納するための最終記憶媒体である。本実施例のストレージシステム1は、最終記憶装置として、磁気ディスクを記録媒体とするHDD (Hard Disk Drive) 200' の他、フラッシュメモリなどの不揮発性半導体メモリを記憶媒体として用いるFMPK (Flash Memory Package) 200を用いることができる。ただしその他の記憶デバイスを用いることも可能である。また、記憶デバイス200、200' は一例として、SAS (Serial Attached SCSI) 規格に従って、ストレージコントローラ10 (以下、「DKC10」と略記する) との通信を行う。

- [0013] DKC10は複数のFMPKを、1又は複数のRAID (Redundant Arrays of Inexpensive/Independent Disks) グループ145として管理する。
- [0014] DKC10には、1以上のホスト2と、管理端末4が接続される。DKC100とホスト2とは、一例としてファイバチャネルを用いて形成されるSAN (Storage Area Network) 1を介して接続される。DKC10と管理端末4とは、一例としてイーサネットを用いて形成されるネットワーク150を介して接続される。
- [0015] DKC10は少なくとも、プロセッサ11、ホストインタフェース (図中では「ホストIF」と表記) 12、ディスクインタフェース (図中では「ディスクIF」と表記) 13、メモリ14、パリティ演算回路15を有する。そしてプロセッサ11、ホストIF12、ディスクIF13、メモリ14、パリティ演算回路15は、相互結合スイッチ (相互結合SW) 16を介して相互接続されている。これらの構成要素は、高性能化及び高可用性の確保のため、DKC10内に複数搭載されている。ただしこれらの構成要素が1つだけDKC10内に設けられている構成でもよい。
- [0016] ディスクIF13は少なくとも、バッファ131、記憶デバイス200、200' と通信するためのインタフェースコントローラ132 (図中では「SAS-CTL」と表記されている)、及び転送回路 (非図示) を有する。

インタフェースコントローラ132は、記憶デバイス200、200'の用いているプロトコル（一例ではSAS）を、DKC10内部で用いられている通信プロトコル（一例としてPCI-Express）に変換するためのものである。本実施例では、記憶デバイス200、200'がSAS規格に従った通信を行うため、インタフェースコントローラ132にはSASコントローラ（以下、「SAS-CTL」と略記する）が用いられる。図1では、1つのディスク1F13にSAS-CTL132が1つのみ記載されているが、1つのディスク1F13に複数のSAS-CTL132が存在する構成を採用してもよい。

- [0017] ディスク1F13は、たとえばFMPK200からデータを読み出し、読み出したデータをバッファ131に一時的に格納する。バッファ143に格納されたデータは、転送回路によって、メモリ14あるいは他のディスク1F13のバッファへと送信される。バッファ131はたとえば揮発性の半導体メモリから構成されるが、不揮発メモリを用いて構成されていても良い。
- [0018] ホスト1F12は、ディスク1F13と同様に、インタフェースコントローラとバッファ、転送回路（非図示）を少なくとも有する。インタフェースコントローラは、ホスト2とDKC10間で用いられている通信プロトコル（たとえばファイバチャネル）と、DKC10内部で用いられている通信プロトコルを変換するためのものである。バッファはホスト2からのデータを一時的に格納するためのものである（逆にメモリ14からホスト2に転送すべきデータを一時的に格納するために用いられることもある）。
- [0019] パリティ演算回路15は、RAID技術で必要とされる冗長データの生成を行うハードウェアである。パリティ演算回路15により生成される冗長データの例としては、排他的論理和（XOR）、リードソロモン符号等がある。
- [0020] プロセッサ11は、ホスト1F12から到来するI/O要求の処理を行う。メモリ14は、プロセッサ11が実行するプログラムや、プロセッサが使用するストレージシステム1の各種管理情報を記憶するために用いられる。

またメモリ14は、記憶デバイス200、200'に対するI/O対象データを一時的に記憶するためにも用いられる。以下、記憶デバイス200、200'に対するI/O対象データを一時的に記憶するために用いられる、メモリ14中の記憶領域を、「キャッシュ」と呼ぶ。メモリ14はDRAM、SRAM等の揮発性記憶媒体で構成されるが、別の実施形態として、不揮発性メモリを用いてメモリ14を構成してもよい。

[0021] ストレージシステム1は、先にも述べたとおり、FMPK200、HDD200'等の、複数種類の記憶デバイスを搭載できる。ただし以下では特に断りのない限り、ストレージシステム1にFMPK200のみが搭載されている構成を前提として説明する。

[0022] 図2を用いて、FMPK200の構成について説明する。FMPK200は、デバイスコントローラ（FMコントローラ）201と複数のFMチップ210から構成される。FMコントローラ201は、メモリ202、プロセッサ203、データの圧縮伸長を行うための圧縮伸長回路204、パリティを計算するためのパリティ演算回路205、SAS-CTL206、FM-IF207を備える。メモリ202、プロセッサ203、圧縮伸長回路204、パリティ演算回路205、SAS-CTL206、FM-IF207は、内部接続スイッチ（内部接続SW）208を介して相互接続されている。

[0023] SAS-CTL206は、FMPK200とDKC10間の通信を行うためのインタフェースコントローラである。SAS-CTL206は、伝送線（SASリンク）を介して、DKC10のSAS-CTL132に接続される。またFM-IF207は、FMコントローラ201とFMチップ210間の通信を行うためのインタフェースコントローラである。

[0024] プロセッサ203は、DKC10から到来する各種コマンドに係る処理を行う。メモリ202には、プロセッサ203が実行するプログラムや、各種管理情報が記憶される。メモリ202には、DRAM等の揮発性メモリが用いられる。ただしメモリ202に不揮発性メモリを使用しても良い。

[0025] 圧縮伸長回路204は、データの圧縮、または圧縮されたデータの伸長を

行う機能を備えたハードウェアである。またパリティ演算回路205は、DKC10が備えるパリティ演算回路15と同様の機能を持つハードウェア、つまりRAID技術で必要とされる冗長データの生成機能を有する。

[0026] FMチップ210は、不揮発性半導体メモリチップで、一例としてNAND型フラッシュメモリである。フラッシュメモリは周知のとおり、ページ単位でデータの読み出し・書き込みが行われ、またデータ消去は、複数ページの集合であるブロック単位で行われる。そして一度書き込みが行われたページは上書きが出来ず、一度書き込みが行われたページに対して再度書き込みを行うためには、当該ページを含むブロック全体を消去する必要がある。そのため、FMPK200は、FMPK200が接続されるDKC10に対しては、FMチップ210の有する記憶領域をそのまま提供することはせず、論理的な記憶空間を提供する。

[0027] なお、FMPK200は、圧縮伸長回路204により、DKC10からのライトデータを圧縮してFMチップ210に格納することができる。ただし、DKC10に対しては原則として、透過的にデータ圧縮を行う。FMPK200はDKC10に対し、所定のサイズの記憶空間（論理アドレス空間）を提供する。DKC10がFMPK200にデータを書き込む際には、この論理アドレス空間上のアドレスとライト対象データのサイズを指定したライトコマンドを発行する。一例として、DKC10がFMPK200に対し、論理アドレス空間の先頭（アドレス0）に48KBのデータを書き込むライトコマンド（及び48KBのライトデータ）を送信したとする。またFMPK200がこの48KBのデータを圧縮した結果、8KBの圧縮データになり、この8KBの圧縮データがFMチップ210に格納されたとする。この状態において、DKC10がこのデータを読み出す場合には、論理アドレス空間の先頭（アドレス0）とリードデータサイズ（たとえば48KB）を指定したリードコマンドを発行することで、先ほど格納した48KBのデータを読み出すことができる。読み出しの過程で、FMPK200は圧縮伸長回路204により、8KBの圧縮データを48KBに伸長（復元）し、伸長さ

れたデータをDKC10に返送するよう動作するからである。そのためDKC10は、（実際にはデータが圧縮されて格納されている場合でも）論理アドレス空間上にデータが非圧縮状態で格納されているかのように認識する。

[0028] 上では、データの圧縮、伸長が、圧縮伸長回路204というハードウェアで行われる例について説明したが、必ずしもデータの圧縮、伸長を、ハードウェアを用いて行わなければならないわけではない。プロセッサ203が圧縮、伸長を行うプログラムを実行することによって、データの圧縮、伸長を行うようにしてもよい。またパリティ演算回路205についても同様で、パリティ演算回路205というハードウェアでパリティの演算を行わなければならないわけではない。プロセッサ203がパリティ演算を行うプログラムを実行することによって、パリティ演算を行うようにしてもよい。

[0029] さらに、上ではFMPK200が圧縮／伸長機能、パリティ演算機能を備えている例を説明したが、圧縮／伸長機能、またはパリティ演算機能を有していないFMPK200がストレージシステム1内に存在してもよい。FMPK200が圧縮／伸長機能を有していない場合には、データは圧縮して格納されない。またFMPK200がパリティ演算機能を有していない場合には、DKC10がパリティの生成を行う。

[0030] 続いて、ストレージシステム1で用いられる記憶領域の概念について説明する。ストレージシステム1は、複数のFMPK200を1つのRAID (Redundant Arrays of Inexpensive / Independent Disks) グループとして管理する。そしてRAIDグループ内で1つ（あるいは2つ）のFMPK200に障害が発生してデータアクセスできなくなった場合に、残りのFMPK200内のデータを用いて、障害が発生したFMPK200に格納されていたデータを復旧できるようにしている。また、RAIDグループ内の一部の記憶領域（あるいは全記憶領域）を、ホスト2などの上位装置に提供する。

[0031] RAIDグループ内の記憶領域について、図3を用いて説明する。図3において、FMPK#x（xは0～3の数値）は、FMPK200がDKC1

0に提供している記憶空間を表している。DKC10は、複数（図3の例では4つ）のFMPK200から1つのRAIDグループ20を構成し、RAIDグループ20に所属する各FMPK（FMPK#0（200-0）～FMPK#3（200-3））上の記憶空間を、ストライプブロックと呼ぶ複数の固定サイズの記憶領域に分割して管理している。

[0032] また図3では、RAIDグループ20のRAIDレベル（RAID技術におけるデータ冗長化方式を表すもので、一般的にはRAID1～RAID6の種類がある）がRAID5である場合の例を表している。DKC10は、図3においてRAIDグループ20内の、「0」、「1」、「P」などのボックスがストライプブロックを表しており、1つのRAIDグループ20内の各ストライプブロックのサイズ（以下では「ストライプサイズ」と呼ばれる）は同じである。ストライプサイズとしては、たとえば64KB、256KB、512KBなどが用いられる。また、各ストライプブロックに付されている、「1」等の番号のことを、「ストライプブロック番号」と呼ぶ。

[0033] 図3で、ストライプブロックのうち、「P」と記載されているストライプブロックは、冗長データの格納されるストライプブロックであり、これを「パリティストライプ」と呼ぶ。一方、数字（0、1等）が記載されているストライプブロックは、ホスト2等の上位装置から書き込まれるデータ（冗長データではないデータ）が格納されるストライプブロックである。このストライプブロックのことは、「データストライプ」と呼ばれる。

[0034] 図3に示されたRAIDグループ20では、たとえばFMPK#3（200-3）の先頭に位置するストライプブロックはパリティストライプ301-3である。そしてDKC10がこのパリティストライプ301-3に格納される冗長データを作成する際、各FMPK200（FMPK#0（200-0）～FMPK#2（200-2））の先頭に位置するデータストライプ（ストライプブロック301-0、301-1、301-2）に格納されるデータに対して所定の演算（たとえば排他的論理和（XOR）等）を施すことによって、冗長データを生成する。

- [0035] 以下、パリティストライプと、当該パリティストライプに格納される冗長データを生成するために用いられるデータストライプのセット（たとえば図3中の要素300）のことを、「ストライプライン」と呼ぶ。実施例1に係るストレージシステム1の場合、図3に示されているストライプライン300のように、1つのストライプラインに属する各ストライプブロックは、FMPK200-0~200-3の記憶空間の同じ位置（アドレス）に存在するという規則に基づいて、ストライプラインが構成される。
- [0036] なお、先に説明したストライプブロック番号は、データストライプに付される番号で、RAIDグループ内で一意な番号である。図3に示されているように、DKC10は、RAIDグループを構成する各FMPK200の先頭に位置するデータストライプに0、1、2の番号を付し、以下順に、3、4、5...の連続番号を付して、データストライプを管理している。以下、ストライプブロック番号がn番（nは0以上の整数値）であるデータストライプのことを、「データストライプn」と表記する。
- [0037] またストレージシステム1では、RAIDグループに属する各記憶デバイス200（200'）に対し、RAIDグループ内で一意な番号を付して管理する。この一意な番号は「RAIDグループ内位置番号」または「位置番号」と呼ばれる。具体的には、RAIDグループ内の先頭ストライプラインに、データストライプ0~データストライプk（ $k > 0$ ）が含まれている場合、データストライプm（ $0 \leq m \leq k$ ）が格納される記憶デバイス200（200'）の位置番号が「m」と定められる。
- [0038] そしてRAID5のように1ストライプライン内に1つのパリティストライプが存在するRAID構成の場合、パリティストライプが格納される記憶デバイス200（200'）の位置番号が「k+1」と定められる。さらにRAID6のように、1ストライプライン内に2つのパリティストライプが存在するRAID構成の場合、パリティストライプが格納される2つの記憶デバイス200（200'）の位置番号はそれぞれ、「k+1」及び「k+2」と定められる。

- [0039] 図3に示されているRAIDグループの場合、先頭ストライプライン300内のデータストライプ0～2が格納される3つの記憶デバイス200（200'）の位置番号は0、1、2と定められる。そしてデータストライプ0～2と同一ストライプライン内の冗長データ（パリティ）の格納される記憶デバイス200（200'）の位置番号は3と定められる。
- [0040] また、DKC10は、ホスト2等の上位装置に対し、論理ユニット（LU）と呼ばれる記憶空間を、1以上提供する。RAIDグループにより形成される記憶空間（以下、「RAIDグループの記憶空間」と呼ぶ）と、論理ユニットの関係について、図4を用いて説明する。RAIDグループの記憶空間とは、図3に示されているRAIDグループ内の領域のうち、データストライプ0から順に、データストライプのみを並べることによって形成される記憶空間である。DKC10は、1つのRAIDグループの記憶空間上の連続領域を論理ユニットと定義することができる。1つのRAIDグループの記憶空間上に複数の論理ユニットを定義してもよいし、1つのRAIDグループの記憶空間全体を1つの論理ユニットとして定義することもできる。図4では、RAIDグループの記憶空間上に2つの論理ユニット（LU#0、LU#1）が定義されている例を表している。
- [0041] 一例として、ストレージシステム1が、ホスト2からLU#0の先頭から3ストライプブロック分の領域に対するデータ書き込み要求（ライトコマンド）を受け付けた場合に行われる処理の概要を説明する。なお、LU#0は図3に示されているRAIDグループに定義されているものとする。この場合、DKC10はライトコマンドとともに受信した3ストライプブロック分のデータを、1ストライプブロックごとに分割する。以下、分割されて生成された3つのデータをそれぞれ、「データ0」、「データ1」、「データ2」と呼ぶ。DKC10はデータ0～データ2を用いて冗長データ（パリティ）を生成し、データ0～データ2及び生成されたパリティをそれぞれ、異なる記憶デバイス200に格納する。LU#が図3に示されているRAIDグループ内に定義されている場合、データ0～データ2及びパリティはそれぞれ

れ、FMPK（200-0）～FMPK（200-3）に格納される。

[0042] LU#0が定義されているRAIDグループのRAIDレベルが、図3のようにRAID5である場合、パリティはデータ0～データ2の排他的論理和（XOR）を計算することで生成できる。パリティの計算は、DKC10のパリティ演算回路15を用いて行われる。ただしRAIDグループを構成する記憶デバイス200が図2のように、パリティ演算回路205等のパリティ生成機能を備えている場合には、記憶デバイス200の有するパリティ生成機能を用いてパリティを生成することもできる。

[0043] RAIDグループ、論理ユニット（LU）を管理するために、DKC10はメモリ14に、デバイス管理テーブル、RAIDグループ管理テーブル（RG管理テーブル）、LU管理テーブルという3種類の管理情報を有する。図5は、デバイス管理テーブルT1000の例を示している。

[0044] デバイス管理テーブルT1000は、ストレージシステム1に搭載された各記憶デバイス200（または200'）についての情報を管理するためのテーブルである。デバイス管理テーブルT1000内の各行（以下では、テーブル内の行のことを「レコード」と呼ぶ）に、ストレージシステム1に搭載された各記憶デバイス200（200'）の情報が格納される。デバイス管理テーブルT1000のレコードは、デバイス#（T1001）、デバイス種別（T1002）、所属RG#（T1003）、デバイスステータス（T1004）、圧縮機能サポート（T1005）、パリティ演算機能サポート（T1006）、サイズ（T1007）の項目を有する。

[0045] DKC10は、ストレージシステム1に搭載された各記憶デバイス200（または200'）に、一意な識別番号を付して管理しており、この識別番号は「デバイス番号」（または「デバイス#」）と呼ばれる。デバイス#（T1001）には、記憶デバイス200（200'）のデバイス#が格納される。デバイス種別（T1002）は、記憶デバイス200（200'）の種類についての情報を格納するための項目である。本実施例の場合、デバイス種別（T1002）には、「FMPK」または「HDD」のいずれかの情

報が格納される。あるレコードのデバイス種別（T1002）に「FMPK」が格納されている場合、当該レコードで管理される記憶デバイスがFMPK200であることを表しており、また「HDD」が格納されている場合、当該レコードで管理される記憶デバイスがHDD200'であることを表している。所属RG#（T1003）については後述する。

[0046] デバイスステータス（T1004）には、記憶デバイスの状態が格納される。デバイスステータス（T1004）に「正常」が格納されている場合、当該レコードで管理される記憶デバイスは正常に稼働していることを表している。デバイスステータス（T1004）に「閉塞」が格納されている場合、当該レコードで管理される記憶デバイスは、障害が発生した等の理由により、稼働していない（閉塞状態にある）ことを表している。

[0047] デバイスステータス（T1004）に「障害復旧中（復旧元）」または「障害復旧中（復旧先）」が格納されている場合、当該レコードで管理される記憶デバイスの所属するRAIDグループについてデータ復旧処理が行われていることを表している。詳細は後述するが、たとえば1台の記憶デバイスに障害が発生した場合、データ復旧処理では障害が発生した記憶デバイスの代わりとなるデバイス（以下、これを「スペアデバイス」と呼ぶ）を1台用意する。そしてDKC10は、障害が発生した記憶デバイスのデバイスステータス（T1004）には「障害復旧中（復旧元）」を格納する。またDKC10は、スペアデバイスとされる記憶デバイスのデバイスステータス（T1004）に「障害復旧中（復旧先）」を格納する。そして障害が発生した記憶デバイスに格納されていたデータを復旧し、スペアデバイスに書き込むことで、データ復旧を行う。

[0048] 圧縮機能サポート（T1005）、パリティ演算機能サポート（T1006）にはそれぞれ、当該レコードで管理される記憶デバイスの、圧縮機能のサポート有無、パリティ演算機能のサポート有無についての情報が格納される。圧縮機能サポート（T1005）に「サポート」が格納されている場合には、当該レコードで管理される記憶デバイスが圧縮機能を有していること

を意味し、「未サポート」が格納されている場合には、当該レコードで管理される記憶デバイスが圧縮機能を有していないことを意味する。パリティ演算機能サポート（T1006）にも同様に、「サポート」または「未サポート」のいずれかの情報が格納され、「サポート」が格納されている場合には、当該レコードで管理される記憶デバイスがパリティ演算機能を有していることを意味する。

[0049] 圧縮機能サポート（T1005）、パリティ演算機能サポート（T1006）に格納される情報は、ストレージシステム1の管理者が、管理端末を用いて設定してもよい。あるいは別の実施形態として、DKC10が各記憶デバイスに対して、各記憶デバイスが有する機能を問い合わせるコマンドを発行することによって、各記憶デバイスが圧縮機能及び／またはパリティ演算機能を有しているか問い合わせ、DKC10がこの問い合わせ結果を圧縮機能サポート（T1005）、パリティ演算機能サポート（T1006）に反映するようにしてもよい。

[0050] サイズ（T1007）には、記憶デバイスの容量、具体的には記憶デバイス200（200'）がDKC10に対して提供している記憶空間のサイズが格納される。この記憶空間のサイズについての情報は、DKC10が記憶デバイス200（200'）に対して、サイズの問い合わせを行うコマンドを発行することで、記憶デバイス200（200'）から取得可能な情報である。なお、記憶デバイス200がFMPK200のように圧縮機能を有する場合、記憶空間のサイズは、記憶デバイス200自身が有する記憶媒体（たとえばFMチップ210）の合計サイズよりも大きいこともある。

[0051] 図6は、RG管理テーブルT1100の例を示している。デバイス管理テーブルT1000と同様、RG管理テーブルT1100内の複数のレコードの各々に、RAIDグループについての情報が格納されている。

[0052] DKC10は、ストレージシステム1内に定義された各RAIDグループに一意的な識別番号を付して管理しており、この識別番号は「RAIDグループ番号」または「RG#」と呼ばれる。RG#（T1101）には、RAI

Dグループ番号 (RG#) が格納される。

所属デバイス# (T1102) は、RAIDグループに含まれている記憶デバイス200 (200') のデバイス#が記憶される。たとえばDKC10が、デバイス#が0、1、2、3、4番の記憶デバイスを用いてRAIDグループを形成した時 (仮にそのRG#が0番であったとする)、DKC10は、RG# (T1101) が0のレコードの、所属デバイス# (T1102) の欄に、「0、1、2、3、4」を格納する。また、図5のデバイス管理テーブルT1000内の各レコードの所属RG# (T1003) には、各記憶デバイスの所属しているRAIDグループ番号が格納される。

[0053] 先にも述べたが、ストレージシステム1は、RAIDグループに属する各記憶デバイス200 (200') に、位置番号を対応付けて管理する。そのため、所属デバイス# (T1102) に登録された各FMPK200には、位置番号 (T1102') が対応付けられる。図6において、レコードT1100-1で管理されているRAIDグループ (RG# (T1101) が1) にはデバイス#が8~15番のFMPK200が所属しているが、各FMPK200には、位置番号0~7の番号が対応付けられている。

[0054] RAID構成 (T1103) は、RAIDグループの構成についての情報が格納される項目であり、少なくとも、RAID技術によるデータ保護方式を表すRAIDレベル、冗長データを生成する時に用いられるデータストライプの数、生成されるパリティストライプの数の情報が格納される。図6の例では、レコードT1100-1のRAID構成 (T1103) は「RAID5 (3D+1P)」だが、これはRAIDレベルが5であること、そして3つのデータストライプから1つのパリティを生成するRAID構成であることを表している。

[0055] RG容量 (T1104)、ストライプサイズ (T1105) にはそれぞれ、RAIDグループに格納できるデータの量 (サイズ)、ストライプサイズが格納される。本実施例では、RG容量 (T1104) に格納される容量は、RAIDグループ内全データストライプの合計サイズであり、パリティの

サイズは含まれない。ただし別の実施形態として、パリティのサイズを含んだ容量が格納されるようにしてもよい。

[0056] RGステータス (T1106) には、RAIDグループの状態 (「正常」、「障害復旧中」、「障害復旧失敗」のいずれかの状態) が格納される。RAIDグループの状態の意味は、デバイスステータス (T1004) と同様であり、RGステータス (T1106) に「正常」が格納されている場合には、RAIDグループが正常に稼働していることを表している。また「障害復旧中」が格納されている場合には、RAIDグループの復旧処理が行われていることを表しており、「障害復旧失敗」が格納されている場合にはRAIDグループが閉塞状態にあることを表している。

[0057] 圧縮 (T1107)、パリティ演算 (T1108) にはそれぞれ、「実施」または「未実施」の情報が格納される。圧縮 (T1107) に「実施」が格納されている場合、当該レコードで管理されるRAIDグループでは、記憶デバイス (FMPK200) が備える圧縮機能を用いたデータ圧縮が行われ、圧縮データがFMPK200に格納されていることを表している (逆に「未実施」が格納されている場合には、FMPK200がデータ圧縮を行っていないことを意味する)。パリティ演算 (T1108) に「実施」が格納されている場合、当該レコードで管理されるRAIDグループに格納されるパリティは、FMPK200が備えるパリティ演算機能を用いて算出されることを表している。

[0058] DKC10は各RAIDグループについて、RAIDグループに所属する全記憶デバイスが圧縮機能をサポートしている場合 (所属デバイス# (T1102) で特定される全記憶デバイスの、圧縮機能サポート (T1005) が「サポート」である場合)、圧縮 (T1107) に「実施」を設定する。同様に、RAIDグループに所属する全記憶デバイスがパリティ演算機能をサポートしている場合 (所属デバイス# (T1102) で特定される全記憶デバイスの、パリティ演算機能サポート (T1006) が「サポート」である場合)、パリティ演算 (T1108) に「実施」を設定する。また別の実

施形態として、ストレージシステム1の管理者が、管理端末を用いて圧縮（T1107）、パリティ演算（T1108）に「実施」または「未実施」を設定するようにしてもよい。

[0059] 図7は、LU管理テーブルT1200の例を示している。デバイス管理テーブルT1000と同様、LU管理テーブルT1200内の各レコードそれぞれに、1つのLUについての情報が格納されている。各LUは論理ユニット番号（LUN）という一意な識別番号を有しており、LU#（T1201）にはLUの論理ユニット番号が格納される。

[0060] 先に述べたとおり、DKC10は、RAIDグループ内の連続領域をLUとして定義する。LU管理テーブルT1200には、LUが定義されているRAIDグループのRG#（T1202）、LUが定義されているRAIDグループ内のオフセットアドレス（T1203）、LUのサイズ（T1204）が格納される。

[0061] 続いて、FMPK200内の記憶領域の管理方法について、図8を用いて説明する。FMPK200は、DKC10に提供する記憶空間（論理アドレス空間）を、論理ページと呼ばれる所定サイズの領域単位で管理する。各論理ページを特定するために、各ページには一意な番号が付されている。この番号のことを「論理ページ番号」と呼ぶ。論理アドレス空間上の先頭に位置する論理ページの論理ページ番号が0番で、それ以降に続く論理ページには、連続した番号が付されている。論理ページのサイズは一例として、16セクタ（8KB）である。

[0062] また、本実施例では、FMPK200内のFMチップ210に存在する記憶領域は「物理ページ」と呼ばれる。物理ページは、フラッシュメモリにおける、アクセス（リード、ライト）の最小単位である。そのためFMコントローラ201がFMチップ210に対してデータの読み書きを行う場合には、物理ページ単位での読み書きを行う。FMPK200には複数のFMチップ210が搭載され、各FMチップ210に複数の物理ページが存在しているが、FMPK200では全FMチップ210内の各物理ページに対して、

ユニークな番号を付して管理している。この番号のことを「物理ページ番号」と呼ぶ。アクセス対象のデータが格納されている物理ページの物理ページ番号が特定されれば、当該物理ページが存在するFMチップ210、及びFMチップ210内位置が一意に特定できる。

[0063] 物理ページのサイズと論理ページのサイズは、同一でも良いし異なってもよい。本実施例では、物理ページのサイズは、528×16バイト(=8KB+256バイト)とする。論理ページのサイズ(8KB)よりも256バイト大きい理由は、各データには後述するDIF、及びECCが付加されるためである。

[0064] 先にも述べたが、FMPK200は圧縮機能を有する。FMPK200がデータを圧縮して格納する場合、FMPK200は論理ページ単位で圧縮を行う。以下、1つの論理ページ上のデータを圧縮することで生成されるデータのことを、「圧縮ページ」と呼ぶ。圧縮ページのサイズは520バイトの倍数であり、最小で520バイト、最大で(520×16)バイトである。

[0065] 圧縮により、圧縮ページのサイズは物理ページ以下になる。そのため、1つの物理ページに複数の圧縮ページを格納することが可能である。また、1つの圧縮ページが複数の物理ページに跨って格納されることもある。なお、後述するように、圧縮ページがFMチップ210に格納される際、実際には8バイトのECCが1または複数付与された状態で格納される。上で説明している圧縮ページのサイズは、ECCが付与されていない状態でのサイズである。

[0066] 図9に、FMPK200が管理するマッピングテーブルT2100の一例を示す。各行(レコード)には、各論理ページについての情報が格納される。レコードには、論理ページ番号(T2101)、物理ページ番号(T2102)、サイズ(T2103)、オフセット(T2104)の項目が含まれる。マッピングテーブルT2100は、メモリ202上に格納されている。

[0067] 論理ページ番号(T2101)には、当該レコードで管理される論理ページの論理ページ番号が格納される。物理ページ番号(T2102)には、当

該レコードで管理される論理ページがマッピングされた物理ページの物理ページ番号が格納される。

[0068] 物理ページには、圧縮されたデータ（圧縮ページ）が格納されるため、マッピングテーブルT 2 1 0 0では、圧縮ページの格納されている、物理ページ内の領域を特定する情報も管理される。その情報が、サイズ（T 2 1 0 3）とオフセット（T 2 1 0 4）である。オフセット（T 2 1 0 4）は物理ページの先頭アドレスを0としたときの、相対アドレスが格納される。そしてオフセット（T 2 1 0 4）とサイズ（T 2 1 0 3）で特定される領域に圧縮ページが格納されることを表している。

[0069] たとえば図9において、レコードT 2 1 0 0 - 2（論理ページ番号T 2 1 0 1が1 0 0 0のレコード）は、論理ページ番号が1 0 0 0の論理ページは、物理ページ番号（T 2 1 0 2）が1 0番の物理ページにマッピングされ、かつこの論理ページのデータ（圧縮ページ）は2 K Bに圧縮されていること、そして物理ページ番号（T 2 1 0 2）が1 0番の物理ページの先頭から4 K Bの位置から始まる2 K Bの大きさの領域に格納されていることを表している。

[0070] また、1つの圧縮ページが複数の物理ページに跨って格納されることもある。図9において、レコードT 2 1 0 0 - 1（論理ページ番号T 2 1 0 1が1のレコード）は、論理ページが、2つの物理ページにマッピングされている例を表している。レコードT 2 1 0 0 - 1で特定される論理ページ（論理ページ番号T 2 1 0 1が1番）は、物理ページ番号（T 2 1 0 2）が5 0 0番の物理ページと4 2番の物理ページにマッピングされている。且つ論理ページのデータは、2 K B（= 1 K B + 1 K B）に圧縮されており、物理ページ番号（T 2 1 0 2）が5 0 0番の物理ページの先頭から7 K Bの位置から始まる1 K Bの大きさの領域、及び物理ページ番号（T 2 1 0 2）が4 2番の物理ページの先頭から0 K Bの位置から始まる1 K Bの大きさの領域に跨って格納されていることを表している。

[0071] また、論理ページに対するアクセス（リード、ライト）がなかった場合に

は、物理ページへのマッピングは行われぬ。図9において、レコードT2100-0は、論理ページ（T2101が2番）に物理ページが割り当てられていない（物理ページ番号T2102がNULL）場合の、レコードの例を表している。FMPK200は、DKC10から論理ページに対するアクセスを受け付けた時点で、初めて論理ページに対して物理ページをマッピングする。

[0072] 続いて、マッピングテーブルT2100以外にFMPK200が管理している情報について説明する。図10は状態管理テーブルT2000の内容を示している。状態管理テーブルT2000は、物理容量T2001、データ圧縮T2002、論理容量T2003、接続DKC種別T2004、所属RAIDグループ構成T2005、RAIDグループ内位置T2006、所属RAIDグループ番号（所属RG#）T2007の項目を有している。

[0073] 物理容量T2001は、FMPK200の有するFMチップ210の合計記憶容量である。データ圧縮T2002には「有」または「無」のいずれかの情報が格納され、「有」が格納されている場合には、FMPK200はDKC10からのライトデータを圧縮してFMチップ210に格納する。データ圧縮T2002の設定は、DKC10がFMPK200を用いてRAIDグループを定義する時等に、DKC10（または管理者）が「有」または「無」を設定する。

[0074] 論理容量T2003は、FMPK200がDKC10に提供しているアドレス空間の容量である。FMPK200にデータが圧縮されずに格納される場合には、原則として物理容量T2001の値と論理容量T2003の値は等しい。FMPK200にデータが圧縮されて格納される場合には、論理容量T2003の値は物理容量T2001より大きくなる。DKC10（または管理者）により、データ圧縮T2002には「有」が設定されると、FMPK200は論理容量T2003に仮の値（たとえば物理容量T2001の8倍の値等）を格納し、DKC10に対して、仮の値と等しいサイズの記憶空間を提供する。またデータがFMチップ210に格納されていくにしたが

って、論理容量T2003に等しい量のデータを格納できないとFMPK200が判断した場合には、論理容量T2003のサイズを小さくする等を行ってよい。逆に、圧縮によりデータサイズが予想以上に小さくなったために、論理容量T2003よりも多い量のデータを格納できるとFMPK200が判断した場合、論理容量T2003のサイズを大きくする等の処理を行ってよい。

[0075] 接続DKC種別T2004は、FMPK200が接続されているストレージシステム1の種別（機種名等）が格納される。FMPK200がストレージシステム1に接続されたことを契機に、DKC10からストレージシステム1の種別についての情報がFMPK200に渡される。FMPK200は渡された情報を接続DKC種別T2004に格納する。

[0076] 所属RAIDグループ構成T2005、RAIDグループ内位置T2006、所属RG#（T2007）は、FMPKが所属するRAIDグループについての情報で、RG管理情報T1100に格納されているRAID構成T1003、所属デバイス#（T1102）、RG#（T1101）の情報と同様の情報が格納される。これらの情報は、DKC10がFMPK200を用いてRAIDグループを定義する際に、DKC10からFMPK200に通知される。

[0077] 続いて、保証コードについて説明する。ストレージコントローラ（DKC）10は、ホスト2から受信したライトデータをFMPK200に格納する過程で、エラー検出用の情報である検証用情報を付加して、データとこの検証用情報とをドライブ121に格納する。なお、この検証用情報は、ホスト2が論理ユニットにアクセスする際の最小アクセス単位である、1ディスクブロック（1セクタともいう。また1ディスクブロック（セクタ）のサイズは、一般的には512バイトであり、本実施例のストレージシステム1においても、1ディスクブロック（セクタ）のサイズは512バイトとする）ごとに付加される。以下ではこの検証用情報のことをDIFと呼ぶ。

[0078] さらにFMPK200内でも、データに更なる検証用情報を付加する処理

が行われる。以下では、DKC10がデータに付加する検証用情報の事を「DKC-DIF」と呼び、FMPK200内でFMPK200がデータに付加する検証用情報のことを「PK-DIF」と呼ぶ。また両者を区別せずに呼ぶ場合には、「DIF」と表記される。また、FMPK200がFMチップ210にデータを格納する際にも、更なる検証用情報を付加するが、この検証用情報のことは「ECC」と呼ばれる。

[0079] 図11を用いて、DKC-DIF、PK-DIFについて説明する。図11は、ホスト2からのライトデータにDIFが付加される過程を表した概念図である。ホスト2からライトデータ500がストレージシステム1に到来すると、DKC（ストレージコントローラ）10はFMPK200に当該ライトデータを送信する前に、DKC-DIF511を付加する。ライトデータ501は、FMPK200に送信される直前の、ライトデータの状態（DKC-DIF511が付加されている状態）を表している。DKC10は1セクタ（512バイト）毎にDKC-DIF511を付加し、FMPK200に対して、DKC-DIF511の付加されたライトデータ501を送信する。

[0080] 続いてFMPK200内でのデータの流れについて説明する。以下では特に、FMPK200でデータが圧縮される場合について説明する。この場合、SAS-CTL206に到着したデータ（ライトデータ501）は、圧縮伸長回路204へと渡される。ライトデータ502は、SAS-CTL206から圧縮伸長回路204に渡されるデータの形式を表している。SAS-CTL206が圧縮伸長回路204にデータを渡す際、1セクタ分のライトデータ毎に、PK-DIF521を付加する。

[0081] ここで、DKC-DIF及びPK-DIFに含まれる情報について説明する。1セクタ（512バイト）のデータに対して付されるDKC-DIF511のサイズは、8バイトである。DKC-DIF511には、CRC（Cyclic Redundancy Check）、RAIDグループ番号、シーケンス番号、アドレス情報が含まれる。

[0082] CRCは、データ510に所定の演算を施して生成される情報である。SAS-CTL206がDKC10から、DKC-DIF511の付加されたライトデータ501を受信すると、データ510に所定の演算を施してCRCを算出する。そして算出されたCRCと、DKC-DIF511内のCRCとが一致するか判定する（以下、この判定のことを「CRCチェック」と呼ぶ）。両者が一致しない場合、DKC10からSAS-CTL206にデータが転送されてくる過程で、障害等の要因でデータ内容が変更されたことを意味する。そのため両者が一致しなかった場合には、DKC10にエラーを返却し、ライトデータ501の書き込み処理を中断する。

[0083] アドレス情報は、データ510の書き込まれるFMPK200上の論理記憶空間上のアドレスである（またはアドレスの一部がアドレス情報に含まれる。たとえばアドレスが4バイトを超えるものである場合、アドレスの下位4バイトのみがアドレス情報として用いられる）。SAS-CTL206が、DKC10からDKC-DIF511の付加されたライトデータ501を受信する場合、それとともにライトデータ501をFMPK200に格納することを指示する命令（いわゆるWRITEコマンド）も受信する。WRITEコマンドにも、データ510の書き込み先である、FMPK200上の論理記憶空間上のアドレス情報が含まれているので、SAS-CTL206はDIFに含まれる内のアドレス情報と、WRITEコマンドに含まれているアドレス情報とが一致するか判定する。両者が一致しない場合には、FMPK200はDKC10にエラーを返却し、ライトデータ501の書き込み処理を中断する。

[0084] RAIDグループ番号は、データ510の書き込まれるFMPK200が所属するRAIDグループの番号（RG#）である。FMPK200はDKC10からあらかじめ、自身が所属するRAIDグループの番号の情報を受信している。そのため、SAS-CTL206が、DKC10からDKC-DIF511の付加されたライトデータ501を受信すると、DKC-DIF511に含まれているRAIDグループ番号と、あらかじめ受信している

RAIDグループ番号とを比較することができる。両者が一致しない場合、FMPK200はDKC10にエラーを返却し、ライトデータ501の書き込み処理を中断する。

[0085] シーケンス番号は一種の連続番号である。DKC10がFMPK200に複数セクタ分のデータを書き込む際、隣り合うデータ510に付されるDKC-DIF511には、連続したシーケンス番号が格納されている。たとえばDKC10が10セクタ分のデータを書き込む場合、先頭のデータ510に付与されるDKC-DIF511には、シーケンス番号0が格納され、次のデータ510に付与されるDKC-DIF511には、シーケンス番号1が格納されている。そのため、連続した複数セクタのデータがライト（またはリード）される時、SAS-CTL206は隣り合ったセクタのシーケンス番号が連続番号であるか判定する。連続番号が付されていない場合には、FMPK200はDKC10にエラーを返却し、ライトデータ501の書き込み処理を中断する。

[0086] PK-DIF521にも同様に、データから算出されたCRCが含まれる。PK-DIF521に含まれるCRCは、データ510とDKC-DIF511で構成される520バイトのデータから算出されるCRCである。

[0087] 図11の説明に戻る。圧縮伸長回路204に対して、DKC-DIF511、PK-DIF521の付加されたデータが渡されると、圧縮伸長回路204はデータ圧縮を行う。なお、先に述べたとおり、圧縮は1論理ページ分のデータ毎に行われる。また圧縮の際、データ510に付加されているDKC-DIF511とPK-DIF521も併せて圧縮される。つまり圧縮伸長回路204は、「データ510とDKC-DIF511とPK-DIF521」のセットを16個まとめて圧縮する。なお圧縮伸長回路204は、サイズが520バイトの倍数になるように圧縮データを生成する。

[0088] 圧縮データ530-0が生成されると、圧縮伸長回路204は圧縮データ530-0にPK-DIF531を付加する。PK-DIF531は、520バイトのデータ（圧縮されたデータ）ごとに付される。PK-DIF53

1はPK-DIF521と同様に、データ（圧縮データ530-0）から算出されたCRCを含んでいる。また圧縮伸長回路204は、圧縮を行う前に、データ510とDKC-DIF511からCRCを算出する。そして算出されたCRCとPK-DIF521に含まれるCRCとが一致するか判定する。両者が一致しない場合にはDKC10にエラーを返却し、ライトデータ501の書き込み処理を中断する。

[0089] 圧縮伸長回路204により生成された、圧縮データ530-0及びそのPK-DIF531は、FM-IF207を経由してFMチップ210に書き込まれる。FM-IF207に圧縮データ530-0及びそのPK-DIF531が到来すると、FM-IF207はPK-DIF531に含まれているCRCのチェックを行う。チェック方法は圧縮伸長回路204で行われているものと同様で、圧縮データ530-0からCRCを算出し、算出されたCRCとPK-DIF531に含まれているCRCが一致するか判定する。両者が一致しない場合、DKC10にエラーを返却し、ライトデータ501の書き込み処理を中断する。

[0090] CRCのチェックでエラーが発生しなかった場合には、FM-IF207は圧縮データ530-0に付与されているPK-DIF531を除去する。そして圧縮データ530-0から別のエラーチェックコードを生成する。このエラーチェックコードは「ECC」と呼ばれる。ECC541はPK-DIF531と同様、520バイトの圧縮データ530-0毎に付与される。そしてFM-IF207は、ECC541の付与された圧縮データ530-0をFMチップ210に書き込む。

[0091] データがFMチップ210から読み出される場合には、上で説明した処理と逆の処理が行われる。FM-IF207がFMチップ210から、ECC541の付加された圧縮データ530-0を読み出し、ECC541のチェックを行う（圧縮データ530-0から算出されたECCと、ECC541の比較を行う）。その後ECC541を圧縮データ530-0から除去して、PK-DIF531を付与し、圧縮伸長回路204へPK-DIF531

の付与された圧縮データ530-0が渡される。圧縮伸長回路204では、PK-DIF531に含まれるCRCのチェックを行い、その後圧縮データ530-0を伸長し、「データ510とDKC-DIF511とPK-DIF521」のセットを（1または複数）生成する。

[0092] 生成された「データ510とDKC-DIF511とPK-DIF521」のセットが、SAS-CTL206を経由してDKC10に転送される際、SAS-CTL206は、PK-DIF521に含まれるCRCのチェックを行い、その後「データ510とDKC-DIF511とPK-DIF521」のセットからPK-DIF521を除去し、データ510とDKC-DIF511をDKC10へと転送する。

[0093] 上で説明したデータの流は、データが圧縮伸長回路204によって圧縮される場合の例である。ただしFMPK200は、データを圧縮せずにFMチップ210に格納することもできる。その場合には、SAS-CTL206でPK-DIF521の付加されたデータ510は圧縮伸長回路204を経由せずにFM-IF207に渡される。FM-IF207では、データ510及びそのDKC-DIF511とPK-DIF521が到来すると、PK-DIF521に含まれているCRCのチェックを行う。チェック方法は上で説明したものと同様である。

[0094] その後FM-IF207は、PK-DIF521の付加されたデータ510からPK-DIF521を除去し、ECCを生成して付加する。ここでのECCは、データ510とDKC-DIF511で構成される520バイトのデータから生成されるものである。そしてECCの付加されたデータ510及びDKC-DIF511をFMチップ210に格納する。

[0095] なお、上で説明したDKC-DIF511、PK-DIF521、PK-DIF531に含まれる情報は一例であって、上で説明した以外の検証用情報が含まれるようにしてもよい。またDKC-DIF511は、FMPK200が接続されるDKC10によって付加される情報であるので、DKC10の種類（機種）によって、DKC-DIF511のフォーマットが異なる

場合がある。たとえばCRCやアドレス情報の長さがDKC10の種類によって異なっている場合があり得る。またDKC-DIF511内における、アドレス情報、シーケンス番号、CRCの並び順がDKC10の種類によって異なっていることもあり得る。本実施例に係るFMPK200は、各機種（DKC）のDKC-DIFのフォーマットについての情報（CRCやアドレス情報の格納されている位置など）を把握している。また接続されるDKCから、DKCの種類（機種）情報を受信することで、CRC、アドレス情報、シーケンス番号、RAIDグループ番号の格納位置を認識可能に構成されている。

[0096] 続いて、FMPK200がサポートするコマンドの種類と、コマンドのフォーマットについて説明する。FMPK200は、DKC10等の上位装置からコマンドを受信し、受信したコマンドに含まれている指示情報（パラメータ）の内容に従った処理（データのリード、ライト等）を行う。FMPK200は、SSDやHDD等の周知の記憶デバイスと同様、データの読み出しを指示するリードコマンド、データの書き込みを指示するライトコマンドもサポートしているが、それ以外のコマンドもサポートしている。ここではFMPK200でサポートされているコマンドのうち、本実施例で行われるデータ復旧処理で用いられるコマンドの内容について説明する。以下、FMPK200にコマンドを発行する発行元の装置は、DKC10であるという前提で説明する。

[0097] [圧縮Readコマンド]

このコマンドは、リードデータを圧縮された状態で、DKC10等のコマンド発行元に対して返却することを指示するためのコマンドである。圧縮Readコマンドに含まれているパラメータについて図12を用いて説明する。

[0098] 圧縮Readコマンド3000には、オペコード（Opcode）3001、Read開始オフセット3002、Readサイズ3003、バッファアドレス3004、転送サイズ3005、のパラメータが含まれる。オペコ

ード (Op code) 3001は、FMPK200がサポートする全コマンドに共通に含まれている情報で (ただしオペコード3001に含まれている情報の内容は、コマンドによって異なる)、FMPK200は受信したコマンドのオペコード3001を参照することで、受信したコマンドの種類を識別する。当然ながら、圧縮Readコマンド3000に含まれているオペコード3001には、圧縮Readコマンドであることが識別できる情報が格納されている。

[0099] Read開始オフセット3002、Readサイズ3003は、DKC10がリードしたいデータ (リード対象データ) が格納されている、FMPK200上の論理アドレス空間上の領域を特定するための情報である。Read開始オフセット3002には、リード対象データの格納されている、FMPK200上の論理アドレス空間上領域の先頭アドレスが、Readサイズ3003にはリード対象データのサイズが指定される。本実施例では、アドレスを特定する情報として、論理ブロックアドレス (LBA) が用いられる。ただし別の実施形態として、アドレスを特定する情報として、論理ページ番号を用いてもよい。またReadサイズ3003には、セクタ数が指定される。ただし別の実施形態として、他の単位 (たとえば論理ページ数、バイト数等) が指定されるようにしてもよい。またReadサイズ3003には、FMPK200内でデータが圧縮されて格納されている、いないにかかわらず、非圧縮時のデータサイズが指定される。

[0100] バッファアドレス3004、転送サイズ3005は、リード対象データの転送先の領域 (領域の先頭アドレスと領域のサイズ) を特定する情報であり、DKC10がFMPK200に圧縮Readコマンド3000を発行する場合には、バッファアドレス3004にはバッファ131のアドレスを指定する。

[0101] FMPK200は、DKC10から圧縮Readコマンド3000を受信すると、圧縮状態のリードデータをDKC10に転送 (DKC10の、バッファアドレス3004、転送サイズ3005で特定されるバッファ131上

領域に転送)するとともに、応答情報をDKC10に返却する。応答情報には、受信したコマンドに係る処理が正常に行われたか否かを表す情報や、サイズについての情報が含まれる。圧縮Readコマンド3000に対する応答情報のフォーマットについて、図13を用いて説明する。

[0102] 圧縮Readコマンド3000に対する応答情報には、転送結果3011、Readサイズ3012、バッファ使用サイズ3013が含まれる。転送結果3011には、「成功」または「エラー」の情報が含まれており、転送結果3011が「成功」の場合には、圧縮Readコマンド3000に係る処理が正常に行われたことを意味する。

[0103] Readサイズ3012には、リード対象データのサイズ（非圧縮時）が格納される。原則として、圧縮Readコマンド3000のReadサイズ3003と同じ値が格納されている。またバッファ使用サイズ3013には、DKC10に転送された、圧縮状態のリードデータのサイズが格納される。

[0104] [圧縮コピーWriteコマンド]

このコマンドは、DKC10が先に説明した圧縮Readコマンド3000を用いて、FMPK200から読み出したデータ（圧縮状態）を、FMPK200に格納する際に用いられる。圧縮コピーWriteコマンドに含まれているパラメータについて、図14を用いて説明する。

[0105] 圧縮コピーWriteコマンド3100には、オペコード（Opcode）3101、Write開始オフセット3102、Writeサイズ3103、転送元アドレス3104、転送サイズ3105、のパラメータが含まれる。オペコード（Opcode）3101は、圧縮Readコマンドの説明の際に述べたとおり、受信したコマンドの種類をFMPK200が識別するための情報が含まれている。

[0106] Write開始オフセット3102、Writeサイズ3103は、ライト対象データの書き込み先の領域を特定するための情報で、Write開始オフセット3102には、ライト対象データの書き込み先となる、（FMP

K200が提供する)論理アドレス空間の先頭アドレスが、Writeサイズ3103にはライト対象データのサイズが指定される。なお、圧縮コピーWriteコマンド3100が発行される時、DKC10から圧縮されたデータがFMPK200に送信されるが、ここで指定されるWrite開始オフセット3102、Writeサイズ3103は、非圧縮時のライトデータの格納される領域(論理アドレス空間上の領域)が指定される。

[0107] 転送元アドレス3104、転送サイズ3105は、FMPK200に転送すべき、圧縮状態のライト対象データが格納されている領域を特定するための情報である。通常、DKC10が圧縮コピーWriteコマンド3100をFMPK200に発行する際、ライト対象データはバッファ131に格納されている。そのため転送元アドレス3104、転送サイズ3105にはそれぞれ、バッファ131上の圧縮状態のライト対象データが格納されている領域の先頭アドレス、圧縮状態のライト対象データのサイズが指定される。

[0108] 圧縮コピーWriteコマンド3100に対する応答情報は、図16に示されているように、転送結果3011のみが含まれている。転送結果3011は、圧縮Readコマンド3000に対する応答情報に含まれている転送結果3011と同じである。つまり「成功」または「エラー」の情報が含まれている。転送結果3011が「成功」の場合には、圧縮コピーWriteコマンド3100に係る処理が正常に行われたことを意味する。

[0109] [圧縮パリティ演算Writeコマンド]

このコマンドは、DKC10が先に説明した圧縮Readコマンド3000を用いて、FMPK200から読み出したデータ(圧縮状態)を、FMPK200に送信するとともに、FMPK200に対し、送信したデータを用いて冗長データ(パリティ)を算出することを指示するためのコマンドである。なお、以下では、圧縮パリティ演算WriteコマンドとともにFMPK200に送信されるデータ(FMPK200から読み出した圧縮状態のデータ)のことを「ライト対象データ」と呼ぶ。圧縮パリティ演算Writeコマンドに含まれているパラメータについて、図15を用いて説明する。

- [0110] 圧縮パリティ演算Writeコマンド3200には、オペコード (Opcode) 3201、Write開始オフセット3202、Writeサイズ3203、バッファアドレス3204、転送サイズ3205、RAIDグループ内位置3206、のパラメータが含まれる。
- [0111] オペコード (Opcode) 3201は、受信したコマンドの種類をFMPK200が識別するための情報が含まれている。
- [0112] Write開始オフセット3202、Writeサイズ3203は、パリティ演算を行って生成されたデータ（以下、これをパリティ演算結果と呼ぶ）の格納先を特定するための情報で、Write開始オフセット3202には、パリティ演算結果の書き込み先となる、（FMPK200が提供する）論理アドレス空間の先頭アドレスが、Writeサイズ3203にはパリティ演算結果のサイズが指定される。なお、圧縮コピーWriteコマンド3100と同様、ここで指定されるWrite開始オフセット3202、Writeサイズ3203には、非圧縮時のパリティ演算結果の格納される領域（論理アドレス空間上の領域）が指定される。
- [0113] バッファアドレス3204、転送サイズ3205は、圧縮コピーWriteコマンドの転送元アドレス3104、転送サイズ3105と同様である。つまりライト対象データが格納されている領域を特定するための情報である。通常、DKC10が圧縮パリティ演算Writeコマンド3200をFMPK200に発行する際、ライト対象データはバッファ131に格納されている。そのためバッファアドレス3204、転送サイズ3205にはそれぞれ、ライト対象データが格納されているバッファ131上領域の先頭アドレス、ライト対象データデータのサイズ（圧縮状態のサイズ）が指定される。
- [0114] RAIDグループ内位置3206には、（圧縮状態の）ライト対象データが、本来格納されていたFMPK200の位置番号が格納される。たとえば図3に示されているRAIDグループ20を例にとって説明する。FMPK#0（位置番号0番）から読み出した圧縮状態のデータを、圧縮パリティ演算Writeコマンドを発行することによってFMPK#3に送信する場合

、DKC10はRAIDグループ内位置3206に、圧縮状態のデータが格納されていたFMPK#0の位置番号（つまり0番）が格納された、圧縮パリティ演算Writeコマンドを作成して、FMPK#3に作成したコマンドを発行する。

[0115] 圧縮パリティ演算Writeコマンド3200に対する応答情報は、圧縮コピーWriteコマンド3100に対する応答情報と同じである。つまり図16に示されているように、転送結果3011のみが含まれている。

[0116] なお、FMPK200は、上で説明した圧縮Readコマンド、圧縮コピーWriteコマンド、圧縮パリティ演算Writeの他に、状態管理テーブルT2000への情報設定を行うためのコマンド（以下では、「情報設定コマンド」と呼ぶ）、FMPK200の故障部位診断コマンドも用意されている。状態管理テーブルT2000への情報設定を行うためのコマンドは、所属RAIDグループ構成（T2005）等の、状態管理テーブルT2000に設定すべき情報を送信するだけであるため、詳細は省略する。また、FMPK200の故障部位の診断用コマンドは、DKC10が、障害の発生したFMPK200に対して発行するコマンドである。このコマンドを受信したFMPK200は、FMPK200がDKC10に対して提供している記憶アドレス空間の中で、障害が発生しているためDKC10がアクセス不可能（リード、ライト不可能）な論理ページの一覧を、コマンド発行元（DKC10）に返却する。

[0117] 続いて、FMPK200に障害が発生した時に、ストレージシステム1で行われるデータ復旧処理について説明する。以下では一例として、RAIDグループ番号が0のRAIDグループに属する、デバイス#が1番のFMPK200に障害が発生した場合を例にとって説明する。

[0118] 図17は、DKC10が実行するデータ復旧処理の全体の流れを表したフローチャートである。DKC10は、FMPK200にアクセス（リードまたはライト）した際に、アクセスが失敗した（エラーが返却された）場合に、図17のデータ復旧処理を開始する。ストレージシステム1のメモリ14

には、データ復旧処理を実行するためのプログラム（データ復旧プログラム）が格納されており、プロセッサ11がこのデータ復旧プログラムを実行することにより、データ復旧処理が実行される。

[0119] データ復旧処理では最初にスペアデバイスの選択が行われる。プロセッサ11は、デバイス管理テーブルT1000を参照し、所属RG#（T1003）が「未割当（スペア）」であるFMPK200を1つ選択する（S20）。以下この選択されたFMPK200を、「復旧先デバイス」または「スペアデバイス」と呼ぶ。なお、復旧先デバイスの選択の際、障害が発生したFMPK200（以下、このFMPK200を「復旧元デバイス」と呼ぶ）と同等のFMPK200を選択する。具体的には、圧縮機能サポートT1005、パリティ演算機能サポートT1006、サイズT1007の内容が同じFMPK200を選択する。

[0120] ただし、復旧元デバイスで、圧縮機能サポートT1005またはパリティ演算機能サポートT1006が「サポート」であったが、圧縮機能サポートT1005またはパリティ演算機能サポートT1006が「サポート」されているスペアデバイスが存在しない場合には、圧縮機能サポートT1005またはパリティ演算機能サポートT1006が「未サポート」であるFMPK200をスペアデバイスとして選択する。ただしこの場合、FMPK200が有する圧縮機能、パリティ演算機能を用いたデータ復旧処理を行うことができず、周知のストレージ装置で行われている、ストレージコントローラによるデータ復旧処理を行う。

[0121] また、デバイス管理テーブルT1000内の、復旧先デバイスのレコードについて、所属RG#（T1003）に、障害の発生したFMPK200の属するRAIDグループ番号と同じ番号を格納し、デバイスステータスT1004に「障害復旧中（復旧先）」を格納する。さらに障害の発生したFMPK200についてのレコードのデバイスステータスT1004に「障害復旧中（復旧元）」を格納する。復旧先デバイスとして、デバイス#が4番のFMPK200（以下、デバイス#がx（xは整数値）のFMPK200の

ことを、「FMPK#x」と表記する)が選択された場合、デバイス管理テーブルT1000は、図5に示された状態になる。

[0122] 本実施例におけるデータ復旧処理では、復旧元デバイスのデータを復旧するために、大きくは以下の2通りの方法が用いられる。1つ目の方法は、復旧元デバイスからデータを読み出し、それを復旧先デバイスに書き込む(コピーする)方法である。復旧元デバイスに障害が発生した場合、全記憶領域がアクセス不可能な状態になることは少ない。そのため、復旧元デバイスの記憶空間の中でも、DKC10からアクセス可能(リード可能)な領域が存在することがある。その場合、復旧元デバイスの記憶空間のうち、DKC10からリード可能な領域については、このリード可能な領域からデータを読み出して復旧先デバイスにコピーすることによってデータ復旧を行う。この方法のことを、以下では「コピー復旧」と呼ぶ。ただしこの方法は、DKC10からアクセス不可能(リード不可能)な領域については利用できない。

[0123] 2つ目の方法は、復旧元デバイスが所属しているRAIDグループ内の各デバイスからデータを読み出し、読み出された各データを用いて所定の演算を行うことで、復旧元デバイスに格納されていたデータを再生成する方法である。この方法のことは、以下では「コレクション(correction)」あるいは「コレクションコピー」と呼ばれる。データの再生成方法は、たとえば特許文献1に記載の演算を行えばよい。コレクションを行う場合、本実施例に係るストレージシステム1では3つの方式を採り得るが、この3つの方式については後述する。

[0124] S30では、プロセッサ11は復旧先デバイスに対して、情報設定コマンドを送信することにより、復旧先デバイスの状態管理テーブルT2000に、所属RAIDグループ構成(T2005)、RAIDグループ内位置(T2006)、データ圧縮(T2002)の設定を行う。所属RAIDグループ構成(T2005)には、復旧先デバイスが所属するRAIDグループのRAID構成と同じもの(RG管理テーブルT1100内のT1103に格納されている情報)が設定される。RAIDグループ内位置(T2006)

には、障害が発生したFMPK 200のRAIDグループ内位置（RG管理テーブルT 1100内のT 1102に格納されている情報）が設定される。またデータ圧縮（T 2002）についても、障害が発生したFMPK 200と同じ情報が設定される。つまり障害が発生したFMPK 200（の所属しているRAIDグループ）がデータ圧縮を行っていた場合には、データ圧縮（T 2002）に「有」が設定される。障害が発生したFMPK 200でデータ圧縮が行われていなかった場合には、データ圧縮（T 2002）に「無」が設定される。

[0125] S 40では、プロセッサ11は復旧元デバイスに対して、故障部位診断コマンドを発行する。このコマンドを受信した復旧元デバイスは、プロセッサ11に対して診断結果を返送する。診断結果には先に述べたとおり、アクセス不可能な論理ページの一覧が含まれている。

[0126] S 50では、プロセッサ11はコピー管理テーブルを作成する。コピー管理テーブルは、プロセッサ11がデータ復旧処理を行う際に把握しておくべき情報を集約しているテーブルである。図18を用いてコピー管理テーブルT 1500で管理される情報の内容を説明する。コピー管理テーブルT 1500には、障害RG#（T 1501）、復旧元デバイス（T 1502）、復旧先デバイス（T 1503）、コレクション方式（T 1504）、復旧デバイス容量（T 1505）、コピー方法ビットマップ（T 1506）、復旧完了済みオフセット（T 1507）、の項目が含まれる。

[0127] 障害RG#（T 1501）には、データ復旧処理による復旧対象であるRAIDグループのRG#が格納される。復旧元デバイス（T 1502）、復旧先デバイス（T 1503）にはそれぞれ、復旧元デバイスのデバイス#、復旧先デバイスのデバイス#が格納される。コレクション方式（T 1504）には、データ復旧処理で実施されるコレクション方式についての情報が格納される。本実施例に係るストレージシステム1では、以下で説明する3通りのコレクション方式が選択可能である。

[0128] 1つ目の方式（以下では「方式1」と呼ぶ）は、周知のストレージ装置で

実施されるコレクション方法と同じものである。具体的にはDKC10が、RAIDグループ内の、復旧元FMPK200以外の正常なFMPK200からデータを読み出し、読み出されたデータからDKC10内のパリティ演算回路15を用いて、復旧元デバイスに格納されていたデータを再生成する。そしてDKC10は再生成されたデータを復旧先に書き出す。

[0129] 2つ目の方式（以下では「方式2」と呼ぶ）は、復旧先FMPK200がパリティ演算機能を有している場合に適用可能な方法で、特許文献1にも開示されている方法である。具体的にはDKC10が、RAIDグループ内の復旧元FMPK200以外の正常なFMPK200からデータを読み出し、読み出されたデータを復旧先FMPK200に送信する。復旧先FMPK200では、自身の有するパリティ演算機能（パリティ演算回路205）を用いて、DKC200から送信されたデータからパリティを算出することでデータを再生成する。なお、方式1、方式2は公知の方法であるので、本実施例では説明を省略する。

[0130] 3つ目の方式（以下では「方式3」と呼ぶ）は、復旧先FMPK200がパリティ演算機能、及び圧縮機能を有している場合に適用可能な方法である。以下では、方式3によるデータ復旧が行われるという前提で説明を行う。そのため方式3の具体的な内容についても、図17以降の処理の流れを説明する過程で説明する。

[0131] コピー管理テーブルT1500の説明に戻る。コレクション方式（T1504）には、上で説明した「方式1」、「方式2」、「方式3」のいずれかが格納される。「方式1」、「方式2」、「方式3」のいずれが格納されるかは、データ復旧対象のRAIDグループに所属するFMPK200が、パリティ演算機能、データ圧縮機能を備えているか否かで決定される。FMPK200が、パリティ演算機能、データ圧縮機能の両方を備えている場合には、プロセッサ11は、コレクション方式（T1504）に「方式3」を設定する。FMPK200が、パリティ演算機能を有し、データ圧縮機能を有していない場合には、プロセッサ11はコレクション方式（T1504）に

「方式2」を設定する。FMPK200が、パリティ演算機能、データ圧縮機能のいずれも有していない場合には、プロセッサ11はコレクション方式(T1504)に「方式1」を設定する。

[0132] コピー方法ビットマップ(T1506)には、データ復旧処理を実施するデータについての情報が格納される。先にも述べたが、FMPK200に障害が発生した時、FMPK200内の全論理ページがアクセス不可能(具体的にはリード不可能)になっていないこともある。そのため、本実施例に係るストレージシステム1では、S40における診断の結果に基づいて、復旧元FMPK200の中で、DKC10がアクセス不可能な論理ページについてのみ、コレクションを行い、アクセス可能な論理ページについては、コピー復旧によるデータ復旧を行う。

[0133] 復旧元FMPK200の論理アドレス空間の大きさが、 $n$ 論理ページ分である場合、DKC10はコピー方法ビットマップ(T1506)として、 $n$ ビットのサイズのビットマップを用意する。そしてコピー方法ビットマップ(T1506)の $k$ ビット目( $1 \leq k \leq n$ )のビットは、復旧元FMPK200の論理アドレス空間の、 $k$ 番目の論理ページについて、コレクションを行うか否かを表している。S40における診断の結果、復旧元FMPK200の論理アドレス空間の、 $k$ 番目の論理ページがアクセス不可能である場合、コピー方法ビットマップ(T1506)の $k$ ビット目のビットに1が格納される(つまり $k$ 番目の論理ページについては、コレクションによるデータ復旧が行われる)。 $k$ 番目の論理ページがアクセス可能である場合は、コピー方法ビットマップ(T1506)の $k$ ビット目のビットに0が格納される(つまり $k$ 番目の論理ページについては、コピー復旧が行われる)。プロセッサ11は、S70において、このビットマップの内容に基づいて、次に行うべき処理を決定する。

[0134] 復旧完了済みオフセット(T1507)には、データ復旧が完了した、論理アドレス空間のアドレスが格納される。本実施例では復旧完了済みオフセット(T1507)に格納するアドレスとして、論理ページ番号を用いる。

ただしその他のアドレス情報（たとえばLBA等）を用いてもよい。

- [0135] また本実施例では、データ復旧は、復旧元デバイスの論理空間の先頭アドレス（論理ページ番号が0の論理ページ）から順に、データ復旧を行う。そのためS50では、プロセッサ11は復旧完了済みオフセット（T1507）に、初期値として0を格納する。そして1論理ページ分のデータ復旧が完了したら、プロセッサ11は復旧完了済みオフセット（T1507）に、データ復旧の完了した論理ページのページ数（1）を加算する。
- [0136] S50の説明に戻る。S50でプロセッサ11は、コピー管理テーブルT1500の障害RG#（T1501）、復旧元デバイス（T1502）、復旧先デバイス（T1503）、コレクション方式（T1504）、復旧デバイス容量（T1505）、コピー方法ビットマップ（T1506）、復旧完了済みオフセット（T1507）に、情報を格納する。RAIDグループ番号が0のRAIDグループに属する、デバイス#が1番のFMPK200に障害が発生した場合、プロセッサ11は障害RG#（T1501）に0を、復旧元デバイス（T1502）に1（FMPK#1）を格納する。また復旧先デバイスとしてFMPK#4が選択された場合には、復旧先デバイス（T1503）には4（FMPK#4）が格納される。
- [0137] 障害が発生したRAIDグループ（RG#0）の状態が図6に示されたものであった場合、そして復旧先デバイスであるFMPK#4の属性（特に圧縮機能サポートT1005、パリティ演算機能サポートT1006）が図5に示されたものであった場合を想定する。この場合、RG#0は圧縮を実施しており、またFMPK200の有するパリティ演算機能を用いたパリティ生成を行っている。またFMPK#4は、圧縮機能、パリティ演算機能のいずれもサポートしている。そのため、プロセッサ11はコレクション方式T1504として「方式3」を格納する。また、復旧デバイス容量（T1505）には、復旧先デバイスのサイズが格納される。
- [0138] コピー方法ビットマップ（T1506）には、上で説明したとおり、S40における診断の結果に基づいて設定が行われる。復旧元FMPK200の

論理アドレス空間の、k番目の論理ページがアクセス不可能である場合、コピー方法ビットマップ(T1506)のkビット目のビットに1が格納される。k番目の論理ページがアクセス可能である場合は、コピー方法ビットマップ(T1506)のkビット目のビットに0が格納される。

[0139] またS50では、プロセッサ11は復旧完了済みオフセット(T1507)を初期化(0を格納)する。

[0140] S60ではプロセッサ11は、コピー方法ビットマップ(T1506)の、(復旧完了済みオフセット(T1507)+1)番目のビットを選択し、S70ではデータ復旧を実施する論理ページの復旧方法を決定する。選択されたビットが0であれば(S70:コピー復旧)、コピー復旧を行う(S71)。選択されたビットが1であれば(S70:コレクション)、コレクションコピーを行う(S72)。S71、S72の処理内容については後述する。

[0141] S71またはS72の処理が完了すると、プロセッサ11はS80の処理を行う。S80ではプロセッサ11は、復旧完了済みオフセット(T1507)に、復旧されたデータのサイズを加算する。論理ページ単位での復旧の場合、1が加算される。S90でプロセッサは、復旧元デバイスの全領域について、データ復旧が完了したか判断する。データ復旧完了の判断は、復旧デバイス容量(T1505)から復旧元デバイスの終端論理ページ番号を求め(復旧デバイス容量(T1505)を論理ページサイズで除算する)、復旧完了済みオフセット(T1507)で示される論理ページ番号が、復旧元デバイスの終端論理ページ番号に達したか判断すればよい。

[0142] データ復旧が完了していない場合(S90:NO)、プロセッサ11は再びS60以降の処理を実施する。データ復旧が完了した場合(S90:YES)、プロセッサ11は管理情報の更新を行う。具体的には、デバイス管理テーブルT1000内の復旧元デバイスのレコードについて、デバイスステータス(T1004)を「閉塞」に変更する。またデバイス管理テーブルT1000内の、復旧先デバイスのレコードについて、デバイスステータス(

T1004)を「正常」に変更する。またプロセッサ11は、RG管理テーブルT1100内の、所属デバイス(T1102)の情報を変更する。具体的には所属デバイス(T1102)に登録されている情報の中から、復旧元デバイスのデバイス#を削除し、復旧先デバイスのデバイス#を追加する(S100)。

[0143] 以上がデータ復旧処理の全体の流れである。なお、上では原則として、論理ページごとにデータ復旧を行う例を説明したが、データ復旧の単位は論理ページに限られない。論理ページよりも大きい単位(たとえば論理ページの整数倍単位)でデータ復旧を行うようにしてもよいし、論理ページより小さい単位でデータ復旧を行っても、データ復旧を行うことができる。

[0144] 続いて、S71の処理(つまり復旧元デバイスから復旧先デバイスにデータをコピーすることによる復旧処理)の詳細を図19を用いて説明する。以下では、論理ページ単位でデータをコピーする場合について説明する。最初にプロセッサ11は、バッファ131にリードデータを格納するための領域として、1論理ページ分の領域を確保する。そして、復旧元デバイスに対して圧縮Readコマンドを発行する(S210)。圧縮Readコマンドのパラメータとして、復旧完了済みオフセットT1507(論理ページ番号)をLBAに変換した値をRead開始オフセット3002に設定する。またReadサイズ3003には1論理ページ分のサイズを用いる。そしてバッファアドレス3004、転送サイズ3005には、先に確保したバッファ131上の領域の情報を用いる。

[0145] 圧縮Readコマンドを発行した後、プロセッサ11は復旧元デバイスから圧縮Readコマンドに対する応答情報を受信する(S220)。応答情報に含まれる転送結果3011が「エラー」であった場合(S230:NO)、プロセッサ11はコレクションコピーを実施して(S280)、処理を終了する。コレクションコピーで行われる処理の詳細は後述する。

[0146] 応答情報に含まれる転送結果3011が「成功」であった場合(S230:YES)、圧縮コピーWriteコマンドを復旧先デバイスに発行するこ

とにより、S210、S220で読み出したデータの復旧先デバイスへの書き込みを指示する(S250)。ここで発行される圧縮コピーWriteコマンドのパラメータである、Write開始オフセット3102、Writeサイズ3103、転送元アドレス3104にはそれぞれ、S210で発行された圧縮Readコマンドのパラメータ、つまり、Read開始オフセット3002、Readサイズ3003、バッファアドレス3004と同じ値が指定される。また圧縮コピーWriteコマンドの転送サイズ3105には、S220で受信した応答情報に含まれているバッファ使用サイズ3013の値が指定される。

[0147] S260でプロセッサ11は、復旧先デバイスから圧縮コピーWriteコマンドに対する応答情報を受信する。応答情報に含まれる転送結果3011が「成功」であった場合(S270: YES)には、処理を終了する。応答情報に含まれる転送結果3011が「エラー」であった場合(S270: NO)、プロセッサ11はコレクションコピーを実施して(S280)、処理を終了する。

[0148] 続いて、S72(またはS280)の処理(コレクションコピー)の詳細を図20を用いて説明する。なお、先にも述べたが、本実施例に係るストレージシステム1では、コレクションとして3つの方式が選択可能であるが、以下では方式3におけるコレクション方式のみ説明する。つまり復旧対象RAIDグループに属するFMPK200が、パリティ演算機能、データ圧縮機能を備えていることが前提である。また、特に断りのない限り、RAID構成がRAID4やRAID5のように、1ストライプライン内のパリティストライプ数が1であるRAIDグループに係るコレクション方法を説明する。また、図19の処理と同様、以下では、論理ページ単位でコレクションコピーを行う場合について説明する。

[0149] まずプロセッサ11は、復旧対象RAIDグループの中で、まだS410以降の処理を行っていないFMPK200のうち、正常なFMPK200を1つ選択する(S400)。続いて変数r、wを用意し、両方の変数に0を

代入することで初期化を行う（S410）。なお変数 *r* は、FMPK200からのデータ読み出しに失敗した場合に、再試行を行った回数を記録するために用いられる。また変数 *w* は復旧先デバイスへのデータ書き込みに失敗した場合に、再試行を行った回数を記録するために用いられる。

[0150] S420でプロセッサ11は、リードデータを格納するための領域として、1論理ページ分の領域をバッファ131に確保し、さらにS400で選択したFMPK200に対して圧縮Readコマンドを発行する。圧縮Readコマンドのパラメータとして指定する内容は、S210で指定されたものと同様である。

[0151] S430では、プロセッサ11は圧縮Readコマンドを発行したFMPK200から応答情報を受信する。応答情報に含まれる転送結果3011が「エラー」であった場合（S440:NO）、プロセッサ11はS450の処理を行う。応答情報に含まれる転送結果3011が「成功」であった場合（S440:YES）、プロセッサ11はS480の処理を実行する。

[0152] 転送結果3011が「エラー」であった場合（S440:NO）、プロセッサ11は変数 *r* が一定値以上であるか判定し（S450）、変数 *r* がまだ一定値以上でない場合（S450:NO）、*r* に1を加算し（S460）、再びS420の処理を実行する。変数 *r* が一定値以上である場合（S450:YES）、プロセッサ11はS540の処理を行う。S540では、RAIDグループの状態（RGステータスT1106）を「障害復旧失敗」に変更し、データ復旧処理を中断する。また管理端末に、データ復旧が失敗した旨を表示する。あるいはホスト2に、データ復旧が失敗した旨を通知するようによい。

[0153] S440の判定で、転送結果3011が「成功」の場合（S440:YES）、プロセッサ11は圧縮パリティ演算Writeコマンドを復旧先デバイスに発行することにより、S420、S430で読み出したデータの復旧先デバイスへの書き込みを指示する（S480）。ここで発行される圧縮パリティ演算Writeコマンドのパラメータとして、Write開始オフセ

ット3202、Writeサイズ3203、バッファアドレス3204にはそれぞれ、S420で発行された圧縮Readコマンドのパラメータ、つまり、Read開始オフセット3002、Readサイズ3003、バッファアドレス3004と同じ値が指定される。また転送サイズ3205には、S430で受信した応答情報に含まれているバッファ使用サイズ3013の値が指定される。そしてRAIDグループ内位置3206には、S400で選択したFMPK200（つまりS420で圧縮Readコマンドを発行したFMPK200）の位置番号（RG管理テーブルT1100の位置番号（T1102'）を参照することで特定可能である）が指定される。

[0154] S490でプロセッサ11は、復旧先デバイスから圧縮パリティ演算Writeコマンドに対する応答情報を受信する。応答情報に含まれる転送結果3011が「成功」であった場合（S500：YES）には、RAIDグループを構成する全ての正常なFMPK200についてS410～S500の処理を行ったか判定し（S550）、全ての正常なFMPK200に対して処理が完了している場合（S550：YES）には、処理を終了する。RAIDグループを構成する全ての正常なFMPK200のうち、まだS410～S500の処理が行われていないFMPK200が存在する場合には（S550：NO）、プロセッサ11は再びS400以降の処理を実施する。

[0155] S500の判定で、転送結果3011が「エラー」の場合（S500：NO）、プロセッサ11は変数wが一定値以上であるか判定し（S510）、変数wがまだ一定値以上でない場合（S510：NO）、wに1を加算し（S520）、再びS420の処理を実行する。変数wが一定値以上である場合（S510：YES）、プロセッサ11はS540の処理を行う。S540では、RAIDグループの状態（RGステータスT1106）を「障害復旧失敗」に変更し、データ復旧処理を中断する。

[0156] なお、上で説明した処理は、1ストライプライン内のパリティストライプ数が1であるRAIDグループについてのデータ復旧処理であるため、S550ではRAIDグループを構成する全ての正常なFMPK200について

S 4 1 0 ~ S 5 0 0 の処理を行ったかの判定が行われている。一方、1ストライプライン内にn個のデータストライプと複数（たとえばRAID6の場合には2個）のパリティストライプが存在するRAIDグループ（一例としてRAID6）についてのデータ復旧処理が行われる場合、S 5 5 0 では、n個のFMPK 2 0 0 についてS 4 1 0 ~ S 5 0 0 の処理を行ったか判定されればよい。

[0157] 以上が、データ復旧処理において、DKC 1 0 で実施される処理の流れである。次に、DKC 1 0 がFMPK 2 0 0 に対して圧縮Readコマンド等のコマンドを発行した時に、FMPK 2 0 0 が実行する処理の流れについて説明する。なお、FMPK 2 0 0 がDKC 1 0 等の上位装置からコマンドを受信した時、FMPK 2 0 2 のメモリ202に格納されているコマンド処理用プログラムをプロセッサ203が実行することによって、そのコマンドに係る処理が行われる。

[0158] まずFMPK 2 0 0 がDKC 1 0 から、圧縮Readコマンドを受信した時に行われる処理の流れを、図21を用いて説明する。FMPK 2 0 0 がDKC 1 0 から圧縮Readコマンドを受信すると、プロセッサ203は変数u、cを用意し、これらの変数を初期化（0を代入）する（S 1 0 2 0）。なお、変数uは主として、FMPK 2 0 0 がDKC 1 0 に返却する応答情報に含まれる、Readサイズ3012の算出のために用いられ、変数cはバッファ使用サイズ3013の算出のために用いられる。

[0159] 続いてS 1 0 3 0 でプロセッサ203は、圧縮Readコマンドで指定されたリード対象データの格納されている、FMチップ210上アドレス（正確には、物理ページの物理ページ番号及び物理ページ内オフセット）の算出を行う。具体的には圧縮Readコマンドのパラメータに含まれている、Read開始オフセット3002に変数uの値を加算して算出されるアドレスから、論理ページ番号を算出する。そしてマッピングテーブルT 2 1 0 0 を参照することによって、算出された論理ページ番号がマッピングされている物理ページの物理ページ番号（T 2 1 0 2）、物理ページ内オフセット（T

2104)、そしてサイズ(T2103)を求める。1論理ページが複数の物理ページに跨って格納されている場合には、物理ページ番号(T2102)、物理ページ内オフセット(T2104)、サイズ(T2103)のセットが複数求められる。

[0160] 続いてプロセッサ203は、S1030で求められた物理ページ番号(T2102)、物理ページ内オフセット(T2104)、サイズ(T2103)で特定される領域(FMチップ210上の領域)からデータを読み出す(S1040)。先にも述べたが、FMチップ210に格納されているデータには、520バイトのデータ毎にECCが付されている。データ読み出しの過程において、FM-IF207は、ECCを用いてデータのチェックを行う。データのチェックにおいて、エラーが発生しなかった場合(データに付されているECCと、データから算出されるECCとが一致している場合)には、FM-IF207は、520バイトのデータ毎に付されているECCを除去する。代わりにFM-IF207が520バイトのデータ毎にPK-DIFを作成してデータに付加し、PK-DIFの付加されたデータを、メモリ202上に格納する。その後プロセッサ203に読み込みが成功した旨を通知する。

[0161] また、FMチップ210の最小アクセス(リード、ライト)単位は、物理ページであるので、S1040では物理ページ単位でのデータ読み出しが行われる。そのため、1論理ページが複数物理ページに跨って格納されている場合、論理ページが圧縮された状態で物理ページに格納され、かつその物理ページに他の論理ページのデータが格納されている場合などには、アクセス対象の論理ページ(あるいは論理ページのデータが圧縮された、圧縮ページ)以外の情報も一緒に読み出される。その場合には、物理ページのデータを読み出してメモリ202に格納したあと、アクセス対象の論理ページ以外の情報をメモリ202から削除する。

[0162] 一方データのチェックにおいて、エラーが発生した場合(データに付されているECCと、データから算出されるECCとが一致しない場合)には、

FM-IF207はプロセッサ203に、読み込みが失敗した旨を通知する。

[0163] プロセッサ203は、FM-IF207から読み込みが失敗した旨を受信した場合（S1050：NO）、DKC10に返却する応答情報の作成を行う（S1160）。ここで作成される応答情報は、転送結果3011に「エラー」が格納された応答情報である。そしてプロセッサ203は、作成された応答情報をDKC10に返却し（S1150）、処理を終了する。

[0164] 一方FM-IF207から読み込みが成功した旨をプロセッサ203が受信した場合（S1050：YES）、S1060以降の処理が行われる。S1060でプロセッサ203は、リードデータが圧縮されているか否かであるか判定する。リードデータが圧縮されているか否かの判定は、S1030で求められたサイズ（T2103）が論理ページのサイズと同じか否かを参照することで判定できる。リードデータが圧縮されていない場合（サイズ（T2103）が論理ページのサイズと同じ場合）、プロセッサ203は圧縮伸長回路204を用いて、リードデータを圧縮する。圧縮伸長回路204ではデータの圧縮後、圧縮されたデータにPK-DIFを付与し、メモリ202に格納する。

[0165] リードデータが圧縮されていた場合（S1060：NO）には、S1070は実行せず、S1080以降の処理が行われる。なお、以下では、S1070で圧縮されたリードデータのサイズ（S1070が実行されなかった場合には、S1040で読み出されたリードデータのサイズ）を $c'$ とする。

[0166] S1080でプロセッサ203は、 $c+c'$ が、圧縮Readコマンドのパラメータで指定されている転送サイズ3005以下であるかを判定する。 $c+c'$ が転送サイズ3005を超過している場合（S1080：NO）、プロセッサ203は、転送結果3011に「成功」が格納された応答情報を作成し（S1140）、作成された応答情報をDKC10に返却し（S1150）、処理を終了する。 $c+c'$ が転送サイズ3005以下である場合（S1080：YES）、S1090以降の処理が行われる。

[0167] 変数  $c$  には、これまでに何回か S1030～S1070 が実行されている場合、DKC10 のバッファ131 に転送済みのデータの合計量が格納されている。一方  $c'$  は、これから DKC10 のバッファ131 に転送すべきデータのサイズである。そのため  $c + c'$  が転送サイズ3005 を超過している場合に、S1040 で読み出したデータ（または S1070 で圧縮したデータ）を転送すると、DKC10 で確保したバッファ131 上領域のサイズを超えたデータが転送されることになる。そのため FMPK200 では、S1080 の判定を行うことにより、DKC10 に返送されるデータの量が、圧縮 Read コマンドのパラメータで指定された転送サイズ3005 を超過しないようにしている。また原則として、DKC10 が圧縮 Read コマンドを発行する時、バッファ131 上に大き目のサイズの領域（伸長時のリードデータサイズと同量の領域等）を確保する。

[0168] S1090 では、プロセッサ203 は、S1040 で読み出されたデータ（S1070 が実行されている場合には、S1070 で圧縮されたデータ）を DKC10 のバッファ131 に転送する。なお、図示を省略しているが、SASCTL206 は転送の過程で、データに付された PK-DIF のチェックを行う。チェック結果が正常な場合には、SASCTL206 はデータに付された PK-DIF を除去して、PK-DIF の除去されたデータを DKC10 に転送する。チェック結果が正常でない場合には、転送結果3011 に「エラー」が含まれた応答情報を DKC10 に返却して、処理を終了する。

[0169] S1090 が終了した後、プロセッサ203 は変数  $c$  に  $c'$  を加算する（S1100）。また S1110 ではプロセッサ203 は、変数  $u$  に対し、S1040 で読み出したデータのサイズ（非圧縮時のサイズ）を加算する。なお、ここで説明している処理では、S1030、S1040 の処理において、1 論理ページ分のデータを読み出しているので、S1110 では1 論理ページサイズが  $u$  に加算される。ただし別の実施形態として、S1030、S1040 の処理で、複数論理ページ分のデータを読み出す、あるいは論理ペ

ージとは関係のない単位でデータを読み出すようにしてもよい。その場合、S 1 1 1 0では読み出したデータのサイズ（ただし非圧縮時サイズ）を変数uに加算する。

[0170] S 1 1 2 0ではプロセッサ203は、変数uがReadサイズ3003未満であるか判定する。uがReadサイズ3003未満の場合（S 1 1 2 0 : Y E S）、プロセッサ203は再びS 1 0 3 0からの処理を行う。uがReadサイズ3003未満でない場合（S 1 1 2 0 : N O）、プロセッサ203は応答情報を作成する（S 1 1 4 0）。ここで作成する応答情報は、転送結果3011に「成功」が、Readサイズ3012に変数uの値が、バッファ使用サイズ3013には変数cの値が格納されたものである。その後プロセッサ203は応答情報をDKC 1 0に返却して、処理を終了する。

[0171] 次に、FMPK 2 0 0がDKC 1 0から、圧縮コピーWriteコマンドを受信した時に行われる処理の流れを、図22を用いて説明する。FMPK 2 0 0がDKC 1 0から圧縮コピーWriteコマンドを受信すると、プロセッサ203はライトデータを受信する（圧縮コピーWriteコマンドのパラメータである、転送元アドレス3104及び転送サイズ3105で特定される、バッファ131上領域からライトデータを取得する）。そして受信したライトデータを、メモリ202に格納する。さらにプロセッサ203は圧縮伸長回路204を用いて、ライトデータを伸長し、伸長されたデータもメモリ202に格納する（S 1 5 2 0）。先に述べたが、圧縮コピーWriteコマンドは、DKC 1 0が圧縮Readコマンド3000を用いて、FMPK 2 0 0から読み出したデータ（圧縮状態）をFMPK 2 0 0に格納する際に用いられる。そのため、圧縮コピーWriteコマンドによってライト対象となるデータは、圧縮された状態でFMPK 2 0 0に到来する。プロセッサ203は圧縮された状態のデータをFMチップ210に格納する前に、D I F（DKC-D I F）を用いたデータの検証（以下、D I Fを用いた検証のことを「D I Fのチェック」と呼ぶ。また、DKC-D I Fを用いたデータの検証は「DKC-D I Fのチェック」と呼び、PK-D I Fを用い

たデータの検証は「PK-DIFのチェック」と呼ぶ)を行うため、一旦データの伸長を行う。

[0172] 続いてプロセッサ203は、伸長データに付加されているPK-DIFとDKC-DIFを用いたデータの検証(チェック)を行う(S1540)。PK-DIFのチェックでは、伸長データから生成したCRCと、PK-DIFに含まれているCRCとを比較し、両者が一致しているか否かを判定する処理が行われる。

[0173] DKC-DIFのチェックでは、主に以下のチェックを行う。

a) 伸長されたデータから生成されるCRCと、伸長されたデータに付加されているDKC-DIF内のCRCが一致するかチェックする。

b) 圧縮コピーWriteコマンド(または圧縮パリティ演算Writeコマンド)のパラメータに含まれているWrite開始オフセット3102と、DKC-DIFに含まれているアドレス情報が一致するかチェックする。

c) FMPK200の状態管理テーブルT2000に格納されている所属RG#(T2007)と、DKC-DIFに含まれているRAIDグループ番号が同一かチェックする。

d) 伸長されたデータに複数の512バイトデータが含まれている場合、各512バイトデータに付加されているDKC-DIF内のシーケンス番号が、連続した番号であるかチェックする。

[0174] なお、FMPK200が接続されるストレージシステム1の種類によって、DKC-DIFのフォーマットが異なる場合がある。つまりストレージシステム1の種類によって、CRCやアドレス情報の格納されているDKC-DIF内の位置が異なる場合がある。そのためプロセッサ203は、FMPK200の状態管理テーブルT2000に格納されている、接続DKC種別T2004の内容に基づいて、DKC-DIF内のどの位置にCRC、アドレス情報、シーケンス番号が格納されているかを特定する。

[0175] S1540の処理の結果、DIFのチェック結果が正常でなかった場合(

S1550:NO)、プロセッサ203は転送結果3011に「エラー」が含まれた応答情報を作成し(S1590)、作成した応答情報をDKC10に返却して(S1580)、処理を終了する。DIFのチェック結果が正常であった場合(S1550:YES)、プロセッサ203は状態管理テーブルT2000のデータ圧縮(T2002)を参照することで、FMPK2000でデータ圧縮を行うことになっているか否か判断する(S1560)。

[0176] データ圧縮(T2002)に「有」が格納されている場合(S1560:YES)、FMチップ210には、メモリ202に格納されたライトデータ(圧縮データ)と伸長データのうち、圧縮データを書き込む処理を行う(S1561)。S1561では、プロセッサ203は圧縮データに付加されているPK-DIFのチェックを行い、チェックが終了したらPK-DIFを削除する。ただしPK-DIFのチェック結果が正常でなかった場合には、S1550、S1590と同様に、DKC10に対して転送結果3011が「エラー」である応答情報を返却して処理を終了する。その後FMチップ210に圧縮データを書き込む。なお、FMチップ210にデータ(圧縮データ)を書き込む際、プロセッサ203は未使用の物理ページ(マッピングテーブルT2100において、どの論理ページT2101にもマッピングされていない物理ページ)を選択して、この選択された物理ページにデータを書き込む。この処理は周知のフラッシュメモリで行われている処理と同じである。

[0177] FMチップ210に圧縮データを書き込む過程で、FM-IF207は、圧縮データからECCを生成し、圧縮データにECCを付加し、ECCの付加された圧縮データをFMチップ210に書き込む。なおECCの生成と付与は、先に述べたとおり、520バイトのデータ毎に行われる。FMチップ210へのデータ書き込みが終了した時点で、マッピングテーブルT2100の内容を更新する。

[0178] マッピングテーブルT2100の更新の概略は以下の通りである。Write開始オフセット3102から論理ページ番号を算出する。マッピングテ

ーブルT2100中のレコードのうち、算出された論理ページ番号と論理ページ番号(T2101)が等しいレコードが、更新されるべきレコードである。このレコードの物理ページ番号(T2102)、オフセット(T2104)、サイズ(T2103)に対して、圧縮データの書き込まれた物理ページ番号、物理ページ内のオフセット、及び圧縮データのサイズを書き込むことで、マッピングテーブルT2100の更新が行われる。

[0179] その後、プロセッサ203は転送結果3011に「成功」が含まれた応答情報を作成し(S1570)、作成した応答情報をDKC10に返却して(S1580)、処理を終了する。また、伸長されたデータ及び圧縮状態のデータの両方が、メモリ202に格納されているが、処理終了時に両方のデータを削除する。

[0180] 一方データ圧縮(T2002)に「無」が格納されている場合(S1560:NO)、FMPK200でデータ圧縮を行わないことを意味する。そのため、伸長データがFMチップ210に書き込まれる(S1562)。S1562では、プロセッサ203は伸長データに付加されているPK-DIFを削除し、FMチップ210に伸長データを書き込む。FMチップ210に伸長データを書き込む過程では、S1561と同様、FM-IF207は、伸長データからECCを生成し、伸長データにECCを付加し、ECCの付加された伸長データをFMチップ210に書き込む。またマッピングテーブルT2100の内容の更新を行う点は、S1561と同様である。

[0181] その後、プロセッサ203は転送結果3011に「成功」が含まれた応答情報を作成し(S1570)、作成した応答情報をDKC10に返却して(S1580)、処理を終了する。また、メモリ202に格納されている、伸長されたデータ及び圧縮状態のデータの両方を、処理終了時に削除する。

[0182] なお、上の処理でFMチップ210に圧縮データを書き込む処理を行う場合(S1561)、圧縮データ(圧縮ページ)は概して物理ページのサイズより小さい。そのため1物理ページに複数の圧縮データを書き込む方が、記憶領域を効率的に使用できる。しかしFMチップ210の最小書き込み単位

は物理ページであるため、圧縮コピーWriteコマンドを受信する毎に、圧縮データを物理ページに書き込む処理を行うと、その物理ページに未使用の領域が残っていても、書き込みができない。

[0183] そのため、別の実施形態として、圧縮コピーWriteコマンドを受信するたびに、S1561においてFMチップ（物理ページ）に圧縮データを書き込むのではなく、メモリ202上で圧縮データにECCを付与した時点で、DKC10に応答情報を返却して（S1570, S1580）、処理を終了してもよい。そして圧縮コピーWriteコマンドを複数回受信した結果、物理ページサイズに等しい（あるいはそれ以上の）量の圧縮データが蓄積された時点で、圧縮データを物理ページに格納する。このようにすると、FMチップ210の記憶領域を効率的に使用できる。この場合、停電などによりFMPK200への電源供給が途絶すると、メモリ202に蓄積されたデータが消失する可能性がある。そのため、FMPK200にバッテリーを備える等して、電源供給の途絶時にもメモリ202の内容が揮発しないように構成することが望ましい。また、S1562の処理を行う場合（伸長データをFMチップ210に書き込む場合）でも、複数のデータを（たとえば1物理ブロック分まとめて）FMチップ210に書き込むようにしてもよい。

[0184] さらに別の実施形態として、DKC10が圧縮コピーWriteコマンドをFMPK200に発行する際に、できるだけ物理ページサイズと等しい量のライトデータを送信するようにしてもよい。DKC10は、各圧縮データのサイズを認識しているので（圧縮Readコマンドの終了時に、応答情報としてFMPK200から圧縮データのサイズを受領するため）、S250において圧縮コピーWriteコマンドを発行する際に、できるだけ物理ページサイズと等しい量の複数の圧縮データをまとめてFMPK200に送信するとよい。このようにすると、FMPK200では圧縮コピーWriteコマンドを受信するたびに、S1561においてFMチップ（物理ページ）に圧縮データを書き込んだとしても、物理ページに効率的に圧縮データを格納できる。

- [0185] 次に、FMPK200がDKC10から、圧縮パリティ演算Writeコマンドを受信した時に行われる処理の流れを、図23を用いて説明する。以下では、圧縮パリティ演算WriteコマンドとともにFMPK200に送信されるデータ（正常なFMPK200から読み出した圧縮状態のデータ）のことを「ライト対象データ」と呼ぶ。また、圧縮パリティ演算Writeコマンドのパラメータである、「Write開始オフセット3202及びWriteサイズ3203」で特定される領域のことを「ライト対象領域」と呼ぶ。
- [0186] S2020～S2050は、図22のS1520～S1550と同じである。またS2150は、S1590と同じである。ただし図22のS1520の実行後には、伸長データと圧縮状態のデータの両方がメモリ202に格納されていたが、S2020では、伸長データのみをメモリ202に格納し、DKC10から受信した圧縮状態のデータは、メモリ202から削除してもよい。
- [0187] S2050の判定の後、DIFのチェック結果が正常であった場合（S2050：YES）、S2060以降の処理が行われる。S2060ではプロセッサ203は、今回のライト対象領域に対して、以前にも圧縮パリティ演算Writeコマンドを受信しているか判定する。この判定方法は後述する。今回のライト対象領域に対する圧縮パリティ演算Writeコマンドを初めて受信した場合（S2060：NO）、ライト対象データを格納する領域をメモリ202上に確保し、S2020で伸長されたデータを、確保されたメモリ202上領域に格納する（S2070）。その後プロセッサ203はS2090以降の処理を行う。
- [0188] プロセッサ203は、ライト対象領域と、ライト対象データ（またはライト対象データから生成されたパリティ）が格納されているメモリ202上の領域との関係を、図24に示されているようなステージング情報管理テーブルT2500で管理する。ステージング情報管理テーブルT2500は、論理ページ番号T2501で特定される論理ページに対するライト対象データ

(あるいはライト対象データを用いて生成されるパリティ)が、アドレスT2502で特定されるメモリ202上領域に格納された状態にあることを管理するテーブルである。また、回数T2503には、論理ページ番号T2501で特定される論理ページをライト対象領域として指定した圧縮パリティ演算Writeコマンドを受信した回数が記録される。

[0189] 初期状態では、アドレスT2502には無効値(NULL)が、回数T2503には0が格納されている。そしてライト対象データを格納する領域がメモリ202上に確保されると、プロセッサ203は確保されたメモリ202上のアドレスを、アドレスT2502に格納する。そして回数T2503に1を加算する。

[0190] そのため、過去に論理ページ番号T2501で特定される論理ページをライト対象領域とする圧縮パリティ演算Writeコマンドを受信していない場合、その論理ページに対応するアドレスT2502には無効値(NULL)が、回数T2503には0が格納された状態にある。S2060では、プロセッサ203は回数T2503の値を参照することにより(アドレスT2502の値を参照してもよい)、今回のライト対象領域に対して、過去に圧縮パリティ演算Writeコマンドを受信しているか判定する。回数T2503が0以外の値であれば、過去に圧縮パリティ演算Writeコマンドを受信していると判定する。

[0191] 過去に圧縮パリティ演算Writeコマンドを受信している場合(S2060: YES)、メモリ202に格納されているデータと、S2020で伸長されたデータとから、パリティが計算される(S2080)。S2080ではプロセッサ203は、メモリ202上のアドレスT2502で特定される領域に格納されているデータと、S2020で生成された伸長データとを、パリティ演算回路205によって演算する。そしてこの演算結果(以下ではこれを「中間パリティ」と呼ぶ)を、アドレスT2502で特定されるメモリ202上領域に格納し、回数T2503に1を加算する(S2080)。なお、上の説明で分かる通り、パリティ生成は伸長されたデータに対して

行われる。圧縮されたデータからパリティを計算しても、データの再生成はできないからである。

[0192] S 2 0 8 0 で、パリティ演算回路 2 0 5 によって行われる演算は、たとえば F M P K 2 0 0 が所属する R A I D グループの R A I D レベルが R A I D 5 の場合、排他的論理和 (X O R) である。一方 F M P K 2 0 0 が所属する R A I D グループの R A I D レベルが R A I D 6 の場合、排他的論理和が計算されることもあれば、リードソロモン符号 (ガロア体の多項式演算) の計算を行うこともある。いずれの計算を行うべきかは、R A I D グループ内のどの F M P K 2 0 0 のデータを再生成する必要があるか、に依存する。

[0193] またリードソロモン符号の計算を行う場合、たとえば特許文献 1 に記載されているように、データに所定の係数を乗じる積演算を行う必要がある。そしてこの乗じるべき所定の係数の値も、R A I D グループ内の位置に依存する。そのため、S 2 0 8 0 ではプロセッサ 2 0 3 は、状態管理テーブル T 2 0 0 0 に格納されている、所属 R A I D グループ構成 T 2 0 0 5、R A I D グループ内位置 T 2 0 0 6、及び圧縮パリティ演算 W r i t e コマンドのパラメータに含まれている R A I D グループ内位置 3 2 0 6 をもとに、データ再生成方法 (排他的論理和、またはリードソロモン符号)、及び積演算で用いられる係数を決定し、それを用いたパリティ生成を行う。

[0194] S 2 0 9 0 で、プロセッサ 2 0 3 はデータの復元が完了したか判定する。たとえば圧縮パリティ演算 W r i t e コマンドを受領した F M P K 2 0 0 の所属する R A I D グループの構成が、R A I D 5 (3 D + 1 P) だった場合 (ストライプラインが 4 ストライプブロックで構成される場合)、3 つの F M P K 2 0 0 から読み出されたデータの X O R を計算すると、データを再生成できる。つまり圧縮パリティ演算 W r i t e コマンドを 3 回受信し、X O R 演算を 3 回行った場合には、データの復元が完了している (アドレス T 2 5 0 2 で特定されるメモリ 2 0 2 上領域には、復元されたデータが格納されている)、と判定可能である。

[0195] そのため本実施例の F M P K 2 0 0 では、プロセッサ 2 0 3 は S 2 0 9 0

において、回数T2503に記録されている、ライト対象アドレスに対する圧縮パリティ演算Writeコマンドを受信している回数を参照することによって、データの復元が完了したか否かを判定する。たとえば圧縮パリティ演算Writeコマンドを受領したFMPK200の所属するRAIDグループの構成がRAID5 (nD+1P) の場合 (nは1以上の整数値) には、n回圧縮パリティ演算Writeコマンドを受信したか判定する。またFMPK200の所属するRAIDグループが、RAID6 (nD+2P) のように、複数のパリティストライプを格納するRAID構成であった場合も、n回圧縮パリティ演算Writeコマンドを受信したか否かを判定することで、データの復元が完了しているか否かを判定可能である。

[0196] データの復元が完了していない場合 (S2090:NO)、つまりまだ圧縮パリティ演算Writeコマンドを所定回数分受信していない場合には、プロセッサ203は、転送結果3011に「成功」が格納された応答情報を作成し (S2170)、作成された応答情報をDKC10に返却し (S2130)、処理を終了する。なおこの場合、S2070またはS2080でメモリ202に格納したデータが、停電などの要因で消失することを防ぐために、FMチップ210に格納するようにしても良い。

[0197] データの復元が完了した場合 (S2090:YES)、つまり圧縮パリティ演算Writeコマンドを所定回数分受信した場合には、プロセッサ203は、メモリ202に格納されているデータ (復元データと呼ぶ) に、DKC-DIFを付与し、さらにPK-DIFの付与を行う (S2100)。

[0198] S2110では、プロセッサ203は状態管理テーブルT2000のデータ圧縮 (T2002) を参照することで、FMPK200でデータ圧縮を行うことになっているか否かを判断する。これはS1560と同様の処理である。

[0199] データ圧縮 (T2002) に「有」が格納されている場合 (S2110:YES)、FMPK200でデータ圧縮を行うことになっているので、S2110の処理によってDIFが付与されたデータを圧縮し、その後FMチッ

プ210に圧縮されたデータを書き込む(S2111)。なお、FMチップ210へのデータ書き込み後、マッピングテーブルT2100の更新も行う。

[0200] その後プロセッサ203は、転送結果3011に「成功」が格納された応答情報を作成し(S2120)、作成された応答情報をDKC10に返却し(S2130)、処理を終了する。

[0201] データ圧縮(T2002)に「無」が格納されている場合(S2110:NO)、FMPK200でデータ圧縮を行わない。そのため、S2110においてDIFの付与されたデータをそのままFMチップ210に書き込む(S2112)。なおFMチップ210へのデータ書き込み後、マッピングテーブルT2100の更新も行う点は、S2111の処理と同様である。その後プロセッサ203は、転送結果3011に「成功」が格納された応答情報を作成し(S2120)、作成された応答情報をDKC10に返却し(S2130)、処理を終了する。

[0202] 実施例1に係るストレージシステム1では、データ復旧時に、復旧元となる記憶デバイス(障害の発生した記憶デバイス、あるいは障害が発生した記憶デバイスと同一RAIDグループに属する記憶デバイス)から圧縮状態でデータを読み出して、復旧先の記憶デバイス(スペアデバイス)にデータを送信するので、復旧用のデータの送信時間を短縮することができる。また、復旧先の記憶デバイスでコレクションを行うため、ストレージコントローラ側でデータコレクションを行う場合に比べて、ストレージコントローラの負荷を低減することができる。さらに、復旧先の記憶デバイスでデータを伸長してDIFのチェックを行うため、データ転送に係るエラーの検出も可能である。

## 実施例 2

[0203] 続いて実施例2に係るストレージシステムについて説明する。実施例2に係るストレージシステム1の構成は、実施例1に係るストレージシステム1と同じである。

[0204] 実施例1と実施例2に係るストレージシステムの違いは、データ復旧処理時、特にコレクション時にFMPK200に対して発行されるコマンドに違いがある。実施例1に係るストレージシステム1では、コレクション時に復旧先デバイスに対して圧縮パリティ演算Writeコマンドという、1種類のコマンドのみを発行していた。一方実施例2に係るストレージシステム1では、コレクション時に復旧先デバイスに対して、圧縮中間パリティ演算コマンド、そしてパリティコミットコマンドという、2種類のコマンドを発行する。以下、実施例2に係るストレージシステム1で行われる、コレクション時の処理の流れについて説明していく。

[0205] まず、圧縮中間パリティ演算コマンド、そしてパリティコミットコマンドという、2種類のコマンドについて説明する。圧縮中間パリティ演算コマンドは、DKC10が実施例1で説明した圧縮Readコマンド3000を用いて、FMPK200から読み出したデータ（圧縮状態）を、FMPK200に送信するとともに、FMPK200に対し、送信したデータとFMPK200に格納されているデータとからパリティを算出することを指示するためのコマンドであるという点で、圧縮パリティ演算Writeコマンドと類似している。

[0206] ただし圧縮パリティ演算Writeコマンドを受信したFMPK200は、最終的には生成されたデータをFMチップ210に格納する（かつマッピングテーブルT2100の更新を行う）が、圧縮中間パリティ演算コマンドがFMPK200に発行された場合、FMPK200は生成されたデータをFMチップ210に格納することは行わない。FMPK200はパリティコミットコマンドを受信した時にはじめて、生成されたデータをFMチップ210に格納する。

[0207] 圧縮中間パリティ演算コマンド、そしてパリティコミットコマンドのコマンドフォーマットを、図25、図26を用いて説明する。圧縮中間パリティ演算コマンドに含まれるパラメータは、圧縮パリティ演算Writeコマンドに含まれるパラメータと同じで、オペコード（Opcode）3201'

、Write開始オフセット3202'、Writeサイズ3203'、バッファアドレス3204'、転送サイズ3205'、RAIDグループ内位置3206'、のパラメータが含まれる。オペコード3201'の値を除き（圧縮中間パリティ演算コマンドのオペコード3201'の値は当然、圧縮パリティ演算Writeコマンド等の他のコマンドのオペコードとは異なる値である）、各パラメータの内容は圧縮パリティ演算Writeコマンドに含まれるパラメータと同じであるため、説明は省略する。

[0208] 図26はパリティコミットコマンドのコマンドフォーマットを示している。パリティコミットコマンドにはオペコード3301、Write開始オフセット3302、Writeサイズ3303が含まれる。パリティコミットコマンドを受信したFMPK200で行われる処理の流れは後述する。また圧縮中間パリティ演算コマンドとパリティコミットコマンドに対する応答情報のフォーマットは、圧縮パリティ演算Writeコマンド等の応答情報と同じものであるため、ここでの説明は省略する。

[0209] 続いて、FMPK200がDKC10から、圧縮中間パリティ演算コマンドを受信した時に行われる処理の流れを、図27を用いて説明する。S2020～S2080までの処理は、図23のS2020～S2080と同じである。つまり中間パリティを生成して、メモリ202に格納する処理が行われる。

[0210] S2080またはS2070の処理が終了した後、プロセッサ203は、転送結果3011に「成功」が格納された応答情報を作成し（S2170）、作成された応答情報をDKC10に返却し（S2130）、処理を終了する。つまり、圧縮パリティ演算Writeコマンドの処理のうち、図23のS2090、S2100、S2110、S2111、S2112、S2120が行われない点が、圧縮中間パリティ演算コマンドと圧縮パリティ演算Writeコマンドの違いである。図23のS2090、S2100、S2110、S2111、S2112、S2120に相当する処理は、FMPK200がパリティコミットコマンドを受信した時に行われる。

- [0211] 続いて、FMPK200がDKC10から、パリティコミットコマンドを受信した時に行われる処理の流れを、図28を用いて説明する。まずプロセッサ203は、コマンドのパラメータに含まれる、Write開始オフセット3302、Writeサイズ3303から、受信したパリティコミットによって処理対象となる領域（以下、この領域をコミット対象領域と呼ぶ）の論理ページ番号を算出する（S2010）。
- [0212] 続いてプロセッサ203はステージング情報管理テーブルT2500を参照し、コミット対象領域に対応するデータ（中間パリティ）が、メモリ202上に格納されているか判定する（S2060'）。これはステージング情報管理テーブルT2500内の論理ページ番号（T2501）が、S2010で特定された論理ページ番号と等しいレコードについて、アドレス（T2502）に有効値（NULL以外の値）が格納されているかを判定すればよい。
- [0213] アドレス（T2502）にNULLが格納されている場合（S2060'：NO）、これまでに、コミット対象領域に対する圧縮中間パリティ演算コマンドを受信していない（中間パリティの生成が行われていない）ことを意味する。そのためプロセッサ203は、転送結果3011に「エラー」が含まれた応答情報を作成し（S2150）、作成した応答情報をDKC10に返却して（S2130）、処理を終了する。
- [0214] アドレス（T2502）に有効値が格納されている場合（S2060'：YES）、メモリ202に格納されているデータ（コミット対象領域に対応するデータ（中間パリティ））にPK-DIFとDKC-DIFを付与する（S2100'）。これは図23のS2100と同じ処理である。以下、S2110'、S2111'、S2112'、S2120'、S2130'は、図23のS2110～S2130と同じ処理である。つまりPK-DIFとDKC-DIFの付与されたデータを、（必要があれば圧縮した後）FMチップ210に書き込み、DKC10に処理が成功した旨の応答情報を返却する。

[0215] 次に、実施例2に係るストレージシステム1で行われるコレクションコピーの処理の流れを、図29を用いて説明する。この処理は、実施例1におけるコレクションコピー（図20）とほとんど同じであるので、相違点のみを説明する。

まず実施例1におけるコレクションコピー（図20）では、DKC10のプロセッサ11は、圧縮パリティ演算Writeコマンドを復旧先デバイスに発行していた（S480）。実施例2におけるコレクションコピー処理では、圧縮中間パリティ演算コマンドを復旧先デバイスに発行する（S480'）。

[0216] また実施例2におけるコレクションコピー処理では、RAIDグループを構成する全ての正常なFMPK200についてS410～S500の処理を行った場合に、復旧先デバイスにパリティコミットコマンドを発行し（S600）、その後処理を終了する。それ以外の点は、実施例1におけるコレクションコピー処理と同じである。また、実施例1で説明した、図17、図19、図21、図22の処理は、実施例2に係るストレージシステム1でも同じである。

[0217] 実施例1におけるコレクションコピー処理では、復旧先デバイスは圧縮パリティ演算Writeコマンドを受信した回数を、領域（論理ページ番号等）毎に記憶しておき、圧縮パリティ演算Writeコマンドを受信した回数が所定回数（RAIDグループを構成する正常なFMPK200の数）に達したことを、復旧先デバイス自身が判断して、FMPK200に復元データを格納していた。実施例2におけるコレクションコピー処理では、DKC10が復旧先デバイスに、FMPK200に復元データを格納する契機を通知するため、FMPK200では圧縮パリティ演算Writeコマンド（圧縮中間パリティ演算コマンド）を受信した回数を記憶しておく必要がない。そのため、実施例2に係るFMPK200では、ステージング情報管理テーブルT2500で、必ずしも回数（T2503）を管理しておく必要はない。

### 実施例 3

- [0218] 続いて実施例3に係るストレージシステムについて説明する。実施例3に係るストレージシステム1の構成は、実施例1に係るストレージシステム1と同じである。
- [0219] 実施例1に係るストレージシステム1では、同一ストライプラインに属するストライプブロックはすべて、各記憶デバイス200(200')上の同じ位置(アドレス)に格納されることが前提であった。実施例3に係るストレージシステム1では、同一ストライプラインに属するストライプブロックがそれぞれ、各記憶デバイス200(200')上の異なる位置(アドレス)に格納される構成を許すものである。
- [0220] 実施例3に係るストレージシステム1でサポートされる、ストライプラインの構成例について、図30を用いて説明する。同一ストライプラインに属するストライプブロックが、すべて記憶デバイス200(200')上の同じ位置(アドレス)に格納されるという構成は、記憶デバイス200(200')の障害時におけるデータ復旧を可能にするための必須条件ではない。データ復旧の観点からすれば、同一ストライプラインに属する各ストライプブロックがすべて、異なる記憶デバイス200(200')に格納されるように配置されていればよい。実施例3に係るストレージシステム1では、この規則が守られるように、各ストライプライン内の各ストライプブロックを記憶デバイス200(200')に配置する。それ以外の制約はない。
- [0221] 図30において、ストライプライン300-1は、FMPK200-1, 200-2, 200-3, 200-4に跨って定義されている。そしてストライプライン300-1内の各ストライプブロックのFMPK200上の位置は、いずれも異なっている。またストライプライン300-2は、FMPK200-3, 200-4, 200-5, 200-6, 200-7に跨って定義されている。そしてストライプライン300-2内の各ストライプブロックのFMPK200上の位置は、同じものもあれば(たとえばFMPK200-3上のストライプブロックとFMPK200-7上のストライプブロック)、異なっているものもある。

- [0222] 実施例3に係るストレージシステム1では、同一ストライプラインに属する各ストライプブロックがすべて、異なる記憶デバイス200(200')に格納されるように配置されるという規則さえ守られれば、各ストライプラインは任意の記憶デバイス200(200')に存在してもよい。
- [0223] この場合、1つのFMPK200、たとえばFMPK200-3に障害が発生してアクセスできなくなったとしても(その他のFMPKは正常に動作しているとする)、FMPK200-1、200-2、200-4がアクセス可能であるため、ストライプライン300-1については、FMPK200-1、200-2、200-4に存在するストライプブロックからコレクションが可能である。同様に、ストライプライン300-2についても、FMPK200-4、200-5、200-6、200-7に存在するストライプブロックからコレクションが可能である。
- [0224] 以下で、実施例3に係るDKC10が実行するデータ復旧処理について説明する。この処理は、実施例1(または実施例2)に係るストレージシステム1で行われるものと多くの部分で共通するので、図17~24を用いて処理の流れを説明する。また、実施例1におけるデータ復旧処理と同じ処理が行われる箇所については説明を省略し、実施例1と相違する点について中心に説明する。なお、実施例1、2では、論理ページ単位にコレクションまたはコピー復旧を行っていたが、実施例3では、ストライプブロック単位にコレクションまたはコピー復旧を行う例について説明する。ただし実施例3に係るストレージシステム1でも、論理ページ単位にコレクションまたはコピー復旧を行うことは可能である。
- [0225] 先にも述べたが、実施例1または2に係るストレージシステム1では、同一ストライプラインに属するストライプブロックが、すべて記憶デバイス200(200')上の同じ位置(アドレス)に格納される。そのためデータ復旧処理もその前提で行われていた。実施例3に係るストレージシステムでは、同一ストライプラインに属するストライプブロックが、記憶デバイス200(200')上の異なる位置に存在し得るため、復旧元デバイスのデー

タをコレクションにより再生成する際、コレクションに必要なデータが、どの記憶デバイスのどのアドレスに格納されているか、特定する必要がある。

[0226] 実施例1または2に係るストレージシステム1では、複数の記憶デバイス200(200')を1つのRAIDグループという概念でまとめて管理している。そして、異なるRAIDグループ内に存在する記憶デバイスに跨ってストライプラインが定義されることはなかった。実施例3に係るストレージシステム1では、ストライプラインに属するストライプブロックは、上で述べた規則に従っている限り、任意の記憶デバイスに存在して良いため、RAIDグループという概念はない。つまり、実施例1または2における、RG管理テーブルT1100は、実施例3に係るストレージシステム1には存在しない。

[0227] 代わりにDKC10は、ストライプラインごとに、ストライプラインに属するストライプブロックの存在する記憶デバイス200(200')及び記憶デバイス内アドレス、及びストライプラインのRAID構成についての情報を管理するテーブルを有する(以下、これをストライプライン管理テーブルと呼ぶ。またストライプライン管理テーブルの内容はRG管理テーブルT1100と類似するため、詳細説明は省略する)。また実施例1に係るストレージシステム1では、デバイス管理テーブルT1000に、記憶デバイスの所属するRAIDグループの情報(所属RG#(T1003))が格納されていたが、実施例3に係るストレージシステム1では所属RG#(T1003)に代えて、デバイス管理テーブルT1000に、記憶デバイスに格納されているストライプブロックが属するストライプラインについての情報(ストライプラインの識別番号等)のリストを格納して管理する。

[0228] これにより、ストライプラインを構成する1(または2)の記憶デバイス200(200')にアクセスできなくなった場合、アクセスできなくなった記憶デバイス200(200')に属するストライプライン(1または複数)を特定し、特定されたストライプラインを構成するストライプブロックが存在する、複数の記憶デバイス200(200')及び記憶デバイス20

0 (200') 内位置を特定する。実施例3に係るデータ復旧処理では、この処理が行われる以外は、実施例1に係るデータ復旧処理と大きな違いはない。

[0229] 実施例1で説明した図17の処理（データ復旧処理の全体の流れ）は、実施例3においても同様の処理が行われるが、上でも述べたとおり、実施例1に係るストレージシステム1と実施例3に係るストレージシステム1では、管理している管理テーブルに違いがあるため、管理情報の更新（S100）の処理が異なる。また実施例3では、ストライプブロック単位でコレクションまたはコピー復旧を行うため、コピー管理テーブルT1500の、コピー方法ビットマップT1506の1ビットは、1ストライプブロックに相当する点、そしてS80の処理で、復旧完了済オフセット（T1507）に復旧されたデータサイズを加算する際、（ストライプサイズ÷論理ページサイズ）を加算する点が、実施例1で説明したものとは異なる。それ以外の点は、実施例1で説明したものと同様の処理が行われる。

[0230] また、S72で行われる処理（コレクションコピー）も、実施例1で説明したものとは若干の違いがある。以下では、実施例3に係るストレージシステム1で行われるコレクションコピーの処理の流れについて、図31を用いて説明する。

[0231] まずプロセッサ11は、復旧元デバイス内の、復旧完了済みオフセットT1507（論理ページ番号）で特定される領域が属するストライプライン（以下、復旧対象ストライプラインと呼ぶ）を特定し、及び復旧対象ストライプラインに属する全ストライプブロックが存在する記憶デバイス（FMPK200）のうち、正常な記憶デバイスのデバイス#、及び当該正常な記憶デバイス内の、復旧対象ストライプラインに属するストライプブロックが格納されているアドレスをすべて特定する（S400'）。これはストライプライン管理テーブル及びデバイス管理テーブルT1000を参照することで特定可能である。

[0232] 続いてプロセッサ11は、S400'で特定された「デバイス#、アドレ

ス」のセットのうち1つを選択する（S405'）。その後S410（変数r、wの初期化）を実行する。

[0233] 続いてプロセッサ11はS420'で、バッファ131にリードデータを格納するための領域として、1ストライプブロック分の領域を確保するとともに、S405'で選択されたデバイス#のFMPK200に対し、同じくS405'で選択されたアドレスをパラメータ（Read開始オフセット3002）として指定した圧縮Readコマンドを発行する。またこの時、圧縮ReadコマンドのReadサイズ3003には、1ストライプブロックのサイズが指定される。続いて行われるS430～S460の処理は、実施例1で説明したものと同一である。

[0234] S480'でプロセッサは、圧縮パリティ演算Writeコマンドを復旧先デバイスに発行する。実施例3に係るストレージシステム1でサポートされる圧縮パリティ演算Writeコマンドに含まれるパラメータについて、図32を用いて説明する。パラメータのうち、オペコード3201～RAIDグループ内位置3206までのパラメータは、実施例1で説明したものと同一である。

[0235] S480'で発行される、圧縮パリティ演算Writeコマンドで指定されるパラメータであるが、Write開始オフセット3202には、復旧完了済みオフセットT1507（論理ページ番号）をLBAに変換した値が設定される。またWriteサイズ3203には1ストライプブロック分のサイズが指定される。そしてバッファアドレス3204には、S420'で確保されたバッファ131上の領域の情報が指定される。また、転送サイズ3205には、S420'、S430で読み出されたデータのサイズ（これはS430で受信した、圧縮Readコマンドの応答情報（バッファ使用サイズ3013）に含まれている）が指定される。

[0236] 図32に示されているとおり、実施例3に係るストレージシステム1でサポートされる圧縮パリティ演算Writeコマンドには、リード元データアドレス3207が追加されている。実施例1で説明したとおり、圧縮パリティ

ィ演算Writeコマンドは、DKC10が圧縮Readコマンド3000を用いて読み出したデータ（圧縮状態）を、復旧先FMPK200に送信するとともに、パリティ算出を復旧先FMPK200に指示するためのコマンドである。リード元データアドレス3207には、この読み出したデータ（圧縮状態）の格納されていたアドレス（正常なストライプブロックの格納されていたアドレス）が指定される。

[0237] S480'の後に行われる、S490～S520、S540の処理は、実施例1で説明したものと同じであるので、ここでの説明は省略する。S500の判定の後、プロセッサ11は、S400'で特定されたすべてのストライプブロックに対して、S405'～S500の処理を実施したか判定する（S550'）。まだS405'～S500の処理が実施されていないストライプブロックが残っている場合（S550'：NO）、プロセッサ11は再びS405'からの処理を繰り返す。すべてのストライプブロックに対して、S405'～S500の処理を実施している場合（S550'：YES）は、処理を終了する、

[0238] 次に、実施例3に係るFMPK200が、DKC10から圧縮パリティ演算Writeコマンドを受信した時に行われる処理の流れを説明する。この処理は、実施例1において説明したものとほとんど同じであるため、図23を用いて相違点についてのみ説明する。

[0239] 実施例3に係るFMPK200で行われる、圧縮パリティ演算Writeコマンドに係る処理は、DKC-DIFのチェック（S2040）に関する処理のみが異なり、後は同じである。

[0240] 実施例3に係るFMPK200でのDKC-DIFのチェックは、主に以下のチェックを行う。

a) 伸長されたデータから生成されるCRCと、伸長されたデータに付加されているDKC-DIF内のCRCが一致するかチェックする。

b') 圧縮パリティ演算Writeコマンドのパラメータに含まれている、リード元データアドレス3207と、DKC-DIFに含まれているア

ドレス情報が一致するかチェックする。

d) 伸長されたデータに複数の512バイトデータが含まれている場合、各512バイトデータに付加されているDKC-DIF内のシーケンス番号が、連続しているかチェックする。

[0241] 上のa)、d)は、実施例1において行われるDKC-DIFのチェックと同じである。一方b')に関して、実施例1に係るFMPK200では、圧縮コピーWriteコマンド（または圧縮パリティ演算Writeコマンド）のパラメータに含まれているWrite開始オフセット3102と、DKC-DIFに含まれているアドレス情報が一致するかチェックしていた。ただし実施例3に係るストレージシステム1では、同一ストライプラインに属するストライプブロックが、記憶デバイス200（200'）上の異なる位置に存在し得るため、Write開始オフセット3102と、DKC-DIFに含まれているアドレス情報を比較することに意味がない。そのため、実施例3に係るストレージシステム1で用いられる圧縮パリティ演算Writeコマンドのパラメータには、リード元データアドレス3207が含まれている。そしてS2040で行われるDKC-DIFのチェックの際、DKC-DIFに含まれているアドレス情報とリード元データアドレス3207を比較するようにしている。その他の点は、実施例1で説明した処理と同じである。

[0242] 以上が実施例3に係るストレージシステム1で行われるデータ復旧処理の説明である。実施例3に係るストレージシステム1では、同一ストライプラインに属するストライプブロックがそれぞれ、記憶デバイス200（200'）上の異なる位置（アドレス）に格納される構成を許しているため、データ配置の自由度が高まる。

[0243] 特に記憶デバイスが圧縮機能を備える場合、格納されるデータの内容によって、記憶デバイスに格納可能なデータ量が異なる。ストライプブロックの格納される記憶デバイス及び記憶デバイス内の位置（アドレス）が固定されていると、圧縮によって記憶デバイスに格納可能な容量が増加したとしても

、増加した記憶領域を有効に活用できないことがある。

[0244] 実施例3に係るストレージシステムの場合、同一ストライプラインに属する各ストライプブロックがすべて、異なる記憶デバイス200(200')に格納されるという規則が守られる限り、任意の記憶デバイスにストライプブロックを格納することができるため、圧縮によって特定の記憶デバイスの格納可能容量が増加した場合、その記憶デバイスに多くのストライプブロックを格納する等の工夫により、圧縮によって増加した記憶領域を有効に活用しやすくなる。

[0245] 以上、本発明の実施例を説明してきたが、これは本発明の説明のための例示であって、本発明を上で説明した実施例に限定する趣旨ではない。本発明は、他の種々の形態でも実施可能である。たとえば実施例に記載のストレージシステム1では、ホスト計算機(ホスト2)からのライトデータを格納する最終記憶媒体が、フラッシュメモリを用いた記憶装置であるFMPKである構成について説明したが、最終記憶媒体はフラッシュメモリを用いた記憶装置に限定されるものではない。たとえばPhase Change RAMやResistance RAM等の不揮発メモリを採用した記憶装置であってもよい。

[0246] また、上の説明において、実施例1または2に係るストレージシステム1では、1つのRAIDグループの記憶領域が、1または複数の論理ユニットに対応付けられている構成であるという前提で説明したが、本発明はこの構成に限定されるものではない。たとえば1つの論理ボリュームが、複数のRAIDグループに対応付けられる構成を採用することもできる。

[0247] また、実施例の説明において、ホスト計算機に提供する論理ユニットの記憶領域とRAIDグループの記憶領域とは、静的にマッピングされている(定義時点で、論理ユニットの各記憶領域がマッピングされるRAIDグループ上の記憶領域が一意に決定される)ものとして説明したが、論理ユニットの記憶領域とRAIDグループ(またはストライプライン)の記憶領域の関係が固定的である構成に限定されるものではない。たとえば周知の技術であ

るThin-Provisioning技術を用いて論理ユニットを定義し、ホスト計算機から論理ユニット上の記憶領域に対してライト要求を受け付けた時点ではじめて、当該記憶領域に対してRAIDグループ（またはストライプライン）の記憶領域を割り当てる等の構成を採用することもできる。

[0248] また実施例においてプログラムとして記載されている構成物は、ハードワイヤードロジックなどを用いたハードウェアによって実現してもよい。また実施例中の各プログラムを、CD-ROM、DVD等の記憶媒体に格納して提供する形態をとることも可能である。

### 符号の説明

- [0249] 1： ストレージシステム  
2： ホスト  
3： SAN  
10： ストレージコントローラ（DKC）  
11： プロセッサ  
12： ホストIF  
13： ディスクIF  
14： メモリ  
15： パリティ演算回路  
16： 相互結合スイッチ  
20： RAIDグループ  
200： 記憶デバイス（FMPK）  
200'： 記憶デバイス（HDD）  
201： FMコントローラ  
202： メモリ  
203： プロセッサ  
204： 圧縮伸長回路  
205： パリティ演算回路  
206： SAS-CTL

207: FM-IF207

208: 内部接続スイッチ

210: FMチップ

## 請求の範囲

### [請求項1]

ホスト計算機と接続されるストレージコントローラと、前記ストレージコントローラに接続される複数の記憶デバイスとを有するストレージシステムにおいて、

前記ストレージシステムは、前記複数の記憶デバイスの中の（ $n + m$ ）台の記憶デバイスからRAIDグループを構成しており、

前記ホスト計算機からのライト要求とともに受け取った $n$ 個のデータから、該 $n$ 個のデータを復元するための $m$ 個の冗長データを生成し、前記 $n$ 個のデータと前記 $m$ 個の冗長データを、前記RAIDグループを構成する（ $n + m$ ）台の前記記憶デバイスに格納するよう構成されており、

前記RAIDグループを構成する記憶デバイスの1つに障害が発生した時、前記ストレージコントローラは、

前記複数の記憶デバイスの中から1台の復旧先記憶デバイスを選択し、

前記RAIDグループの中で障害の発生していない記憶デバイスのそれぞれから、前記データ及び冗長データを圧縮状態で読み出し、

前記読み出された圧縮状態の前記データ及び冗長データを前記復旧先記憶デバイスに送信する、

ことを特徴とする、ストレージシステム。

### [請求項2]

前記記憶デバイスは、記憶媒体とデバイスコントローラを有し、

前記RAIDグループを構成する記憶デバイスの1つに障害が発生した時、前記ストレージコントローラは前記RAIDグループの中で障害の発生していない記憶デバイスのそれぞれに圧縮Readコマンドを発行し、

前記デバイスコントローラは、前記ストレージコントローラから前記圧縮Readコマンドを受信すると、前記記憶デバイスに格納されているデータを圧縮状態で読み出して、前記ストレージコントローラ

に前記圧縮状態のデータを転送する、  
ことを特徴とする、請求項1に記載のストレージシステム。

[請求項3] 前記記憶媒体に前記データが非圧縮状態で格納されている場合、  
前記デバイスコントローラは、前記ストレージコントローラから前記圧縮Readコマンドを受信すると、前記記憶媒体に格納されているデータを圧縮した後、前記ストレージコントローラに前記圧縮状態のデータを転送する、  
ことを特徴とする、請求項2に記載のストレージシステム。

[請求項4] 前記ストレージシステムは、前記n個のデータと前記m個の冗長データを圧縮状態で、前記(n+m)台の記憶デバイスに格納するよう構成されており、  
前記デバイスコントローラは、前記ストレージコントローラから前記圧縮Readコマンドを受信すると、前記記憶媒体に格納されている圧縮状態のデータを前記ストレージコントローラに転送する、  
ことを特徴とする、請求項2に記載のストレージシステム。

[請求項5] 前記ストレージコントローラは、前記圧縮状態のデータまたは冗長データを前記復旧先記憶デバイスに送信する際、圧縮パリティ演算Writeコマンドを前記復旧先記憶デバイスに発行し、  
前記圧縮パリティ演算Writeコマンドを受信した前記復旧先記憶デバイスは、前記圧縮状態のデータまたは冗長データを伸長し、前記伸長されたデータまたは冗長データと過去に前記圧縮パリティ演算Writeコマンドを受信した際に生成された中間パリティとから、新たなパリティを算出し、前記新たなパリティを中間パリティとして前記記憶デバイスに格納する、  
ことを特徴とする、請求項2に記載のストレージシステム。

[請求項6] 前記圧縮パリティ演算Writeコマンドを受信した前記復旧先記憶デバイスは、過去に前記圧縮パリティ演算Writeコマンドを受信していなかった場合、前記圧縮パリティ演算Writeコマンドと

ともに受信した前記圧縮状態のデータ又は冗長データを伸長し、前記伸長されたデータを前記中間パリティとして前記記憶デバイスに格納する、

ことを特徴とする、請求項5に記載のストレージシステム。

[請求項7] 前記圧縮パリティ演算Writeコマンドを受信した前記復旧先記憶デバイスは、前記圧縮パリティ演算Writeコマンドを受信した回数を記憶しており、前記圧縮パリティ演算Writeコマンドを受信した回数が所定回数に達した時点で、前記新たなパリティを復元されたデータとして、前記記憶媒体に格納する、

ことを特徴とする、請求項5に記載のストレージシステム。

[請求項8] 前記復旧先記憶デバイスは、コミットコマンドを受信すると、前記記憶デバイスに格納された中間パリティを復元されたデータとして、前記記憶媒体に格納する、

ことを特徴とする、請求項5に記載のストレージシステム。

[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のライト要求とともに受け取ったデータに検証用情報を付加して、前記記憶デバイスに格納するよう構成されており、

前記圧縮パリティ演算Writeコマンドを受信した前記復旧先記憶デバイスは、前記圧縮状態のデータまたは冗長データを伸長し、前記伸長されたデータに付加されている前記検証用情報を用いてデータの検証を行い、

前記検証結果が正常でなかった場合、前記ストレージコントローラにエラーを返却する、

ことを特徴とする、請求項5に記載のストレージシステム。

[請求項10] 前記ストレージコントローラは、前記ホスト計算機からのライト要求とともに受け取ったデータに前記検証用情報を付加する際、前記検証用情報に前記記憶デバイス上のアドレス情報を格納するよう構成されており、

前記圧縮パリティ演算Writeコマンドには、前記圧縮状態のデータまたは冗長データの格納されていた、前記記憶デバイス上のアドレス情報が含まれており、

前記圧縮パリティ演算Writeコマンドを受信した前記復旧先記憶デバイスは、前記圧縮状態のデータまたは冗長データを伸長し、前記伸長されたデータに付加されている前記検証用情報に含まれている前記アドレス情報と、前記圧縮パリティ演算Writeコマンドに含まれている前記アドレス情報とが一致していない場合、前記ストレージコントローラにエラーを返却する、

ことを特徴とする、請求項9に記載のストレージシステム。

[請求項11]

前記ストレージコントローラは、前記RAIDグループを構成する記憶デバイスの1つに障害が発生した時、

前記障害が発生した記憶デバイスの記憶領域のうち、アクセス可能な記憶領域が存在する場合、前記RAIDグループの中で障害の発生していない記憶デバイスのそれぞれから前記データ及び冗長データを読み出すことに代えて、

前記アクセス可能な記憶領域からデータを圧縮状態で読み出し、

前記読み出された圧縮状態の前記データを復旧先記憶デバイスに格納する、

ことを特徴とする、請求項1に記載のストレージシステム。

[請求項12]

記憶媒体とデバイスコントローラを有し、ストレージコントローラからのデータアクセス要求を受け付ける記憶デバイスであって、

前記記憶デバイスは、前記ストレージコントローラから圧縮パリティ演算Writeコマンドと圧縮状態のデータを受信すると、

前記圧縮状態のデータを伸長し、前記伸長されたデータと、過去に前記圧縮パリティ演算Writeコマンドを受信した際に生成された中間パリティとから、新たなパリティを算出し、前記新たなパリティを中間パリティとして前記記憶デバイスに格納する、

ことを特徴とする、記憶デバイス。

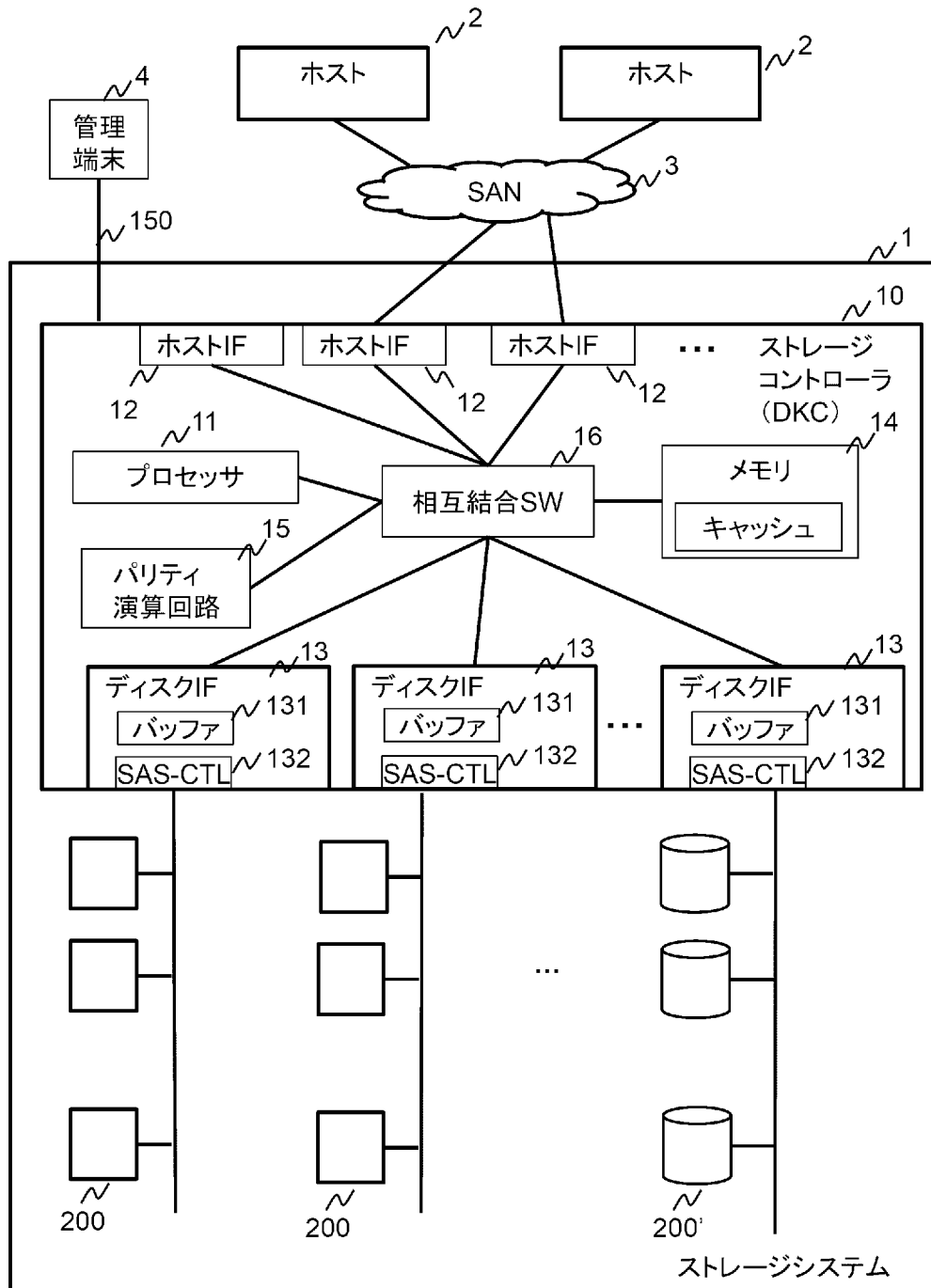
- [請求項13] 前記デバイスコントローラは、前記ストレージコントローラから前記圧縮Readコマンドを受信すると、前記記憶デバイスに格納されているデータを圧縮状態で読み出して、前記ストレージコントローラに前記圧縮状態のデータを転送する、  
ことを特徴とする、請求項12に記載の記憶デバイス。

- [請求項14] 前記圧縮パリティ演算Writeコマンドを受信した前記記憶デバイスは、前記圧縮パリティ演算Writeコマンドを受信した回数を記憶しており、前記圧縮パリティ演算Writeコマンドを受信した回数が所定回数に達した時点で、前記新たなパリティを復元されたデータとして、前記記憶媒体に格納する、  
ことを特徴とする、請求項12に記載の記憶デバイス。

- [請求項15] 前記復旧先記憶デバイスは、コミットコマンドを受信すると、前記記憶デバイスに格納された中間パリティを復元されたデータとして、前記記憶媒体に格納する、  
ことを特徴とする、請求項12に記載の記憶デバイス。

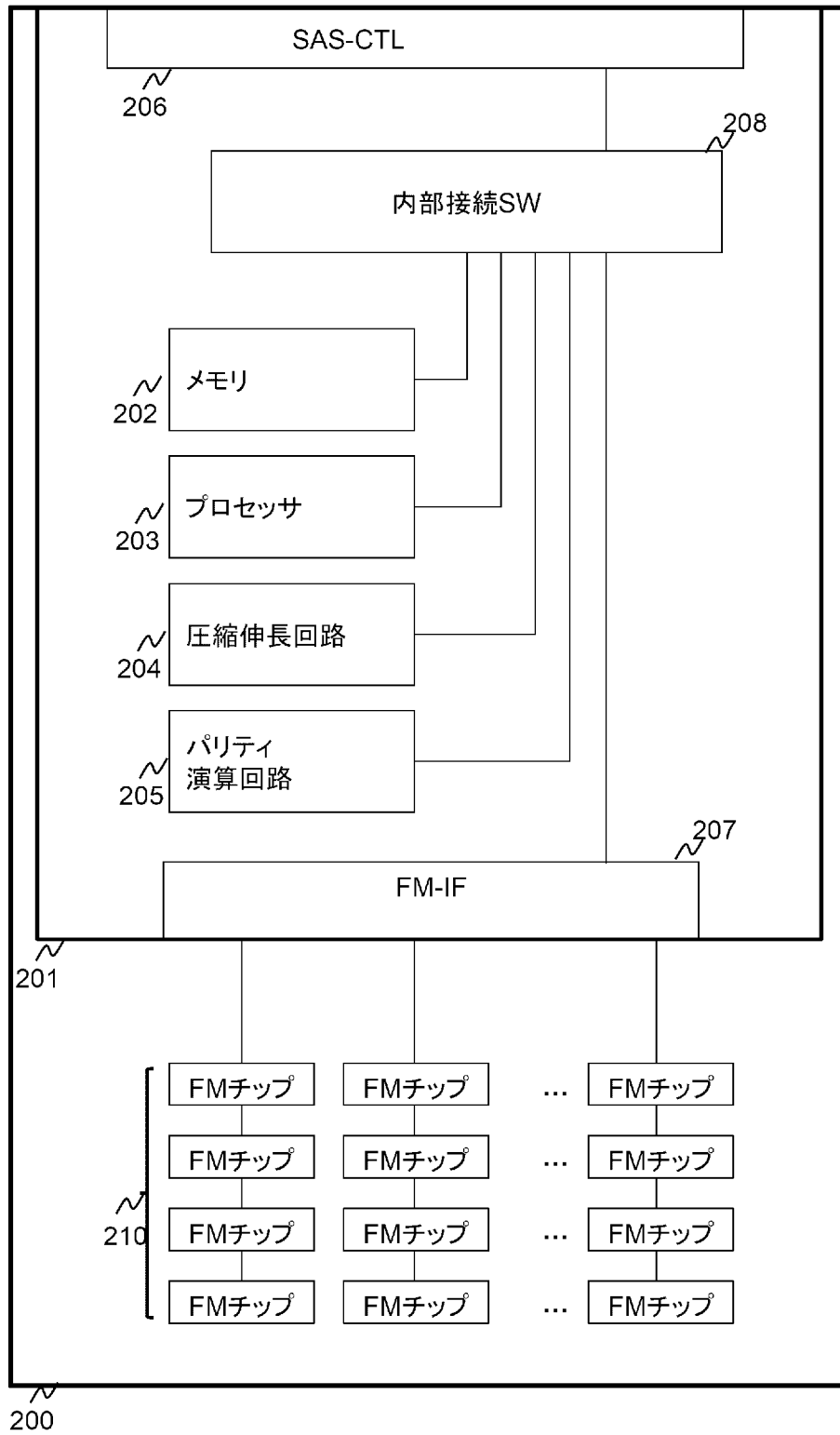
[図1]

図1



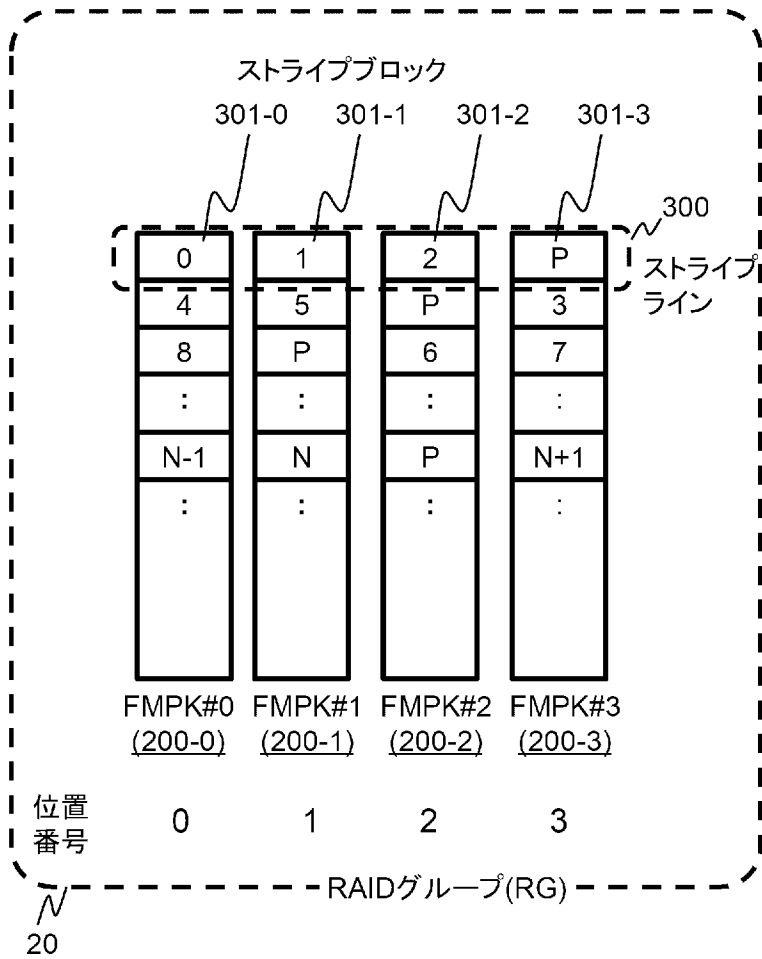
[図2]

図2



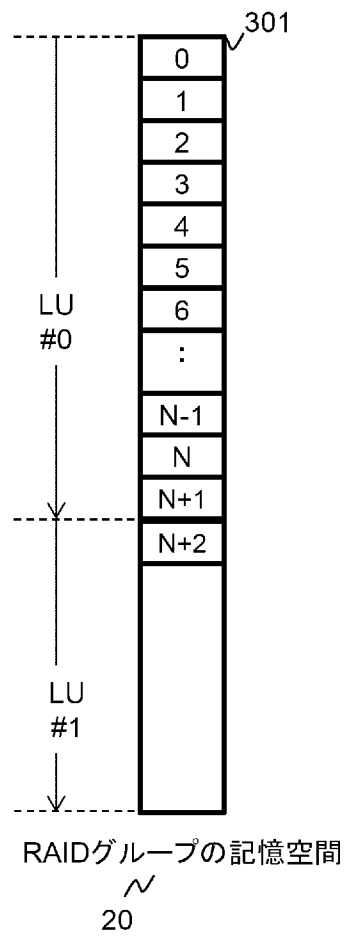
[図3]

図3



[図4]

図4



## [図5]

図5

デバイス #	デバイス 種別	所属RG#	デバイス ステータス	圧縮機能 サポート	パリティ演算 機能サポート	サイズ
0	FMPK	0	正常	サポート	サポート	10240GB
1	FMPK	0	障害復旧中 (復旧元)	サポート	サポート	10240GB
2	FMPK	0	正常	サポート	サポート	10240GB
3	FMPK	0	正常	サポート	サポート	10240GB
4	FMPK	0	障害復旧中 (復旧先)	サポート	サポート	10240GB
5	FMPK	未割当 (スペア)	正常	サポート	サポート	10240GB
6	FMPK	未割当	正常	サポート	サポート	5120GB
7	FMPK	未割当	閉塞	サポート	サポート	10240GB
8	FMPK	1	正常	サポート	サポート	:
:	:	:	:	:	:	:
16	HDD	2	正常	未サポー ト	未サポート	:
:	:	:	:	:	:	:



[図7]

図7

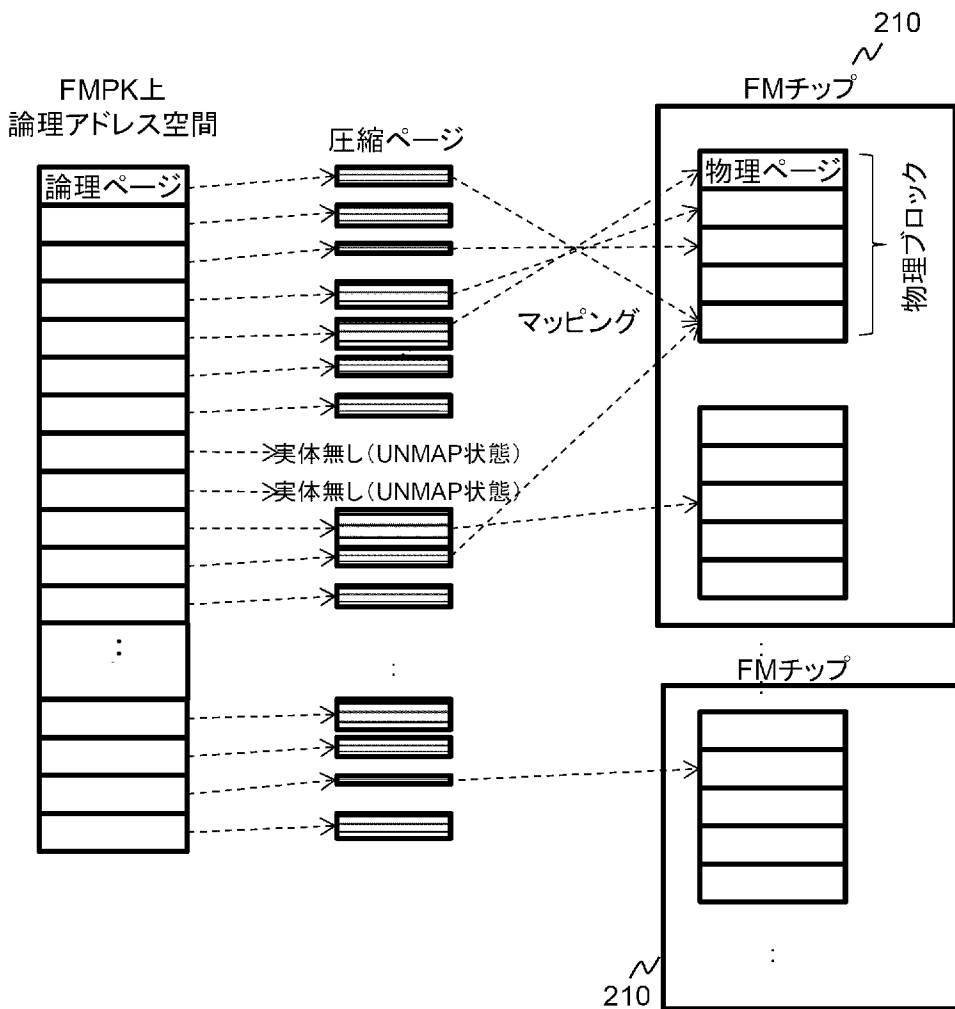
T1200  
~

LU#	RG#	RG内LU開始オフセット	LUサイズ
0	0	0x0	100GB
1	0	0xC800000	500GB
2	0	0x4B000000	50GB
:	:	:	:
5	2	0x20000	10GB
:	:	:	:

T1201
T1202
T1203
T1204

[図8]

図8



[図9]

図9

~ T2100

論理ページ 番号	物理ページ番号	サイズ (圧縮サイズ)	オフセット
0	10	4KB	0KB
1	500	1KB	7KB
	42	1KB	0KB
2	NULL	-	-
:	:	:	:
1000	10	2KB	4KB
:	:	:	:

T2100-1  
T2100-0  
T2100-2

[図10]

図10

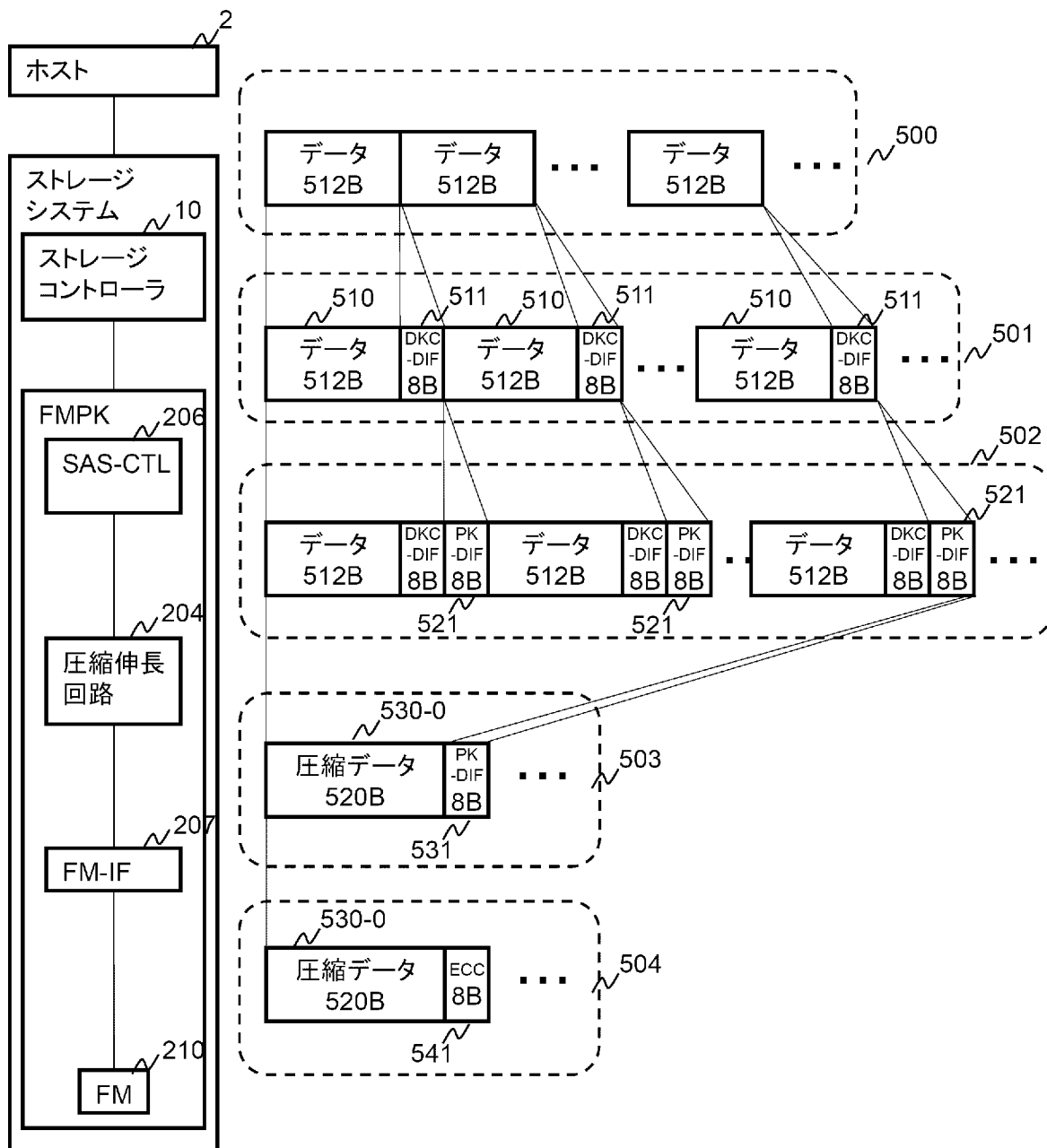
~ T2000

項目	内容
物理容量(非圧縮容量)	1634GB
データ圧縮(実施有無)	有
論理容量	10240GB
接続DKC種別	Hitachi VSP
所属RAIDグループ構成	RAID5(3D+1P)
RAIDグループ内位置	2
所属RG#	0
:	:

T2001  
T2002  
T2003  
T2004  
T2005  
T2006  
T2007

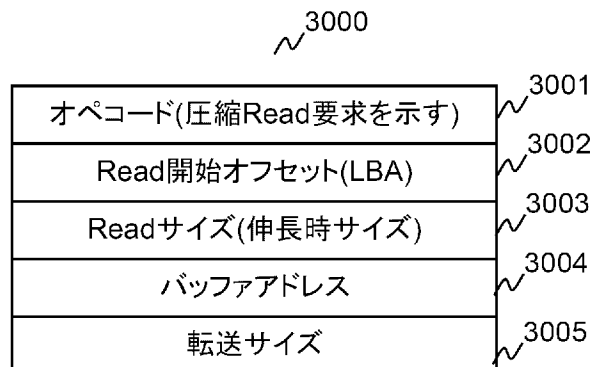
[図11]

図11



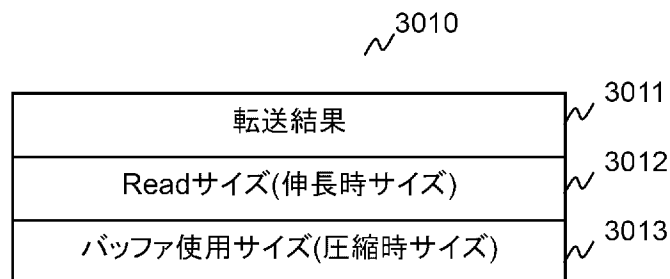
## [図12]

図12



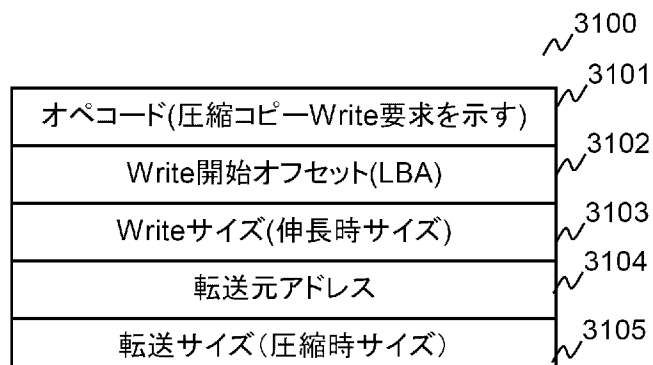
## [図13]

図13



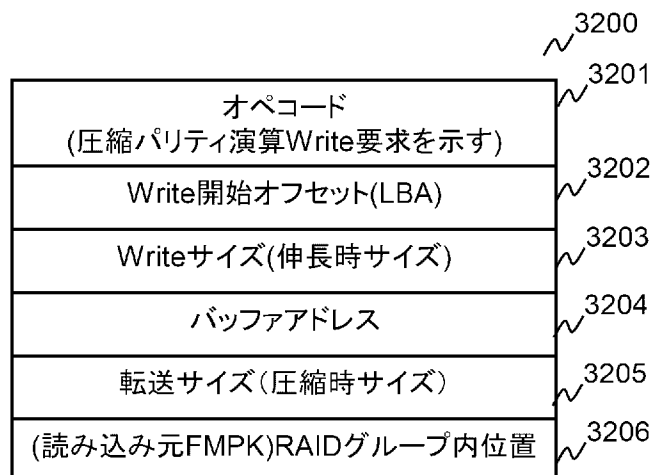
## [図14]

図14



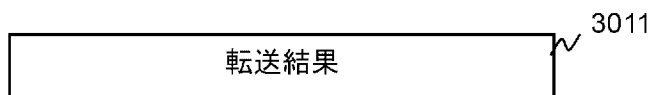
## [図15]

図15



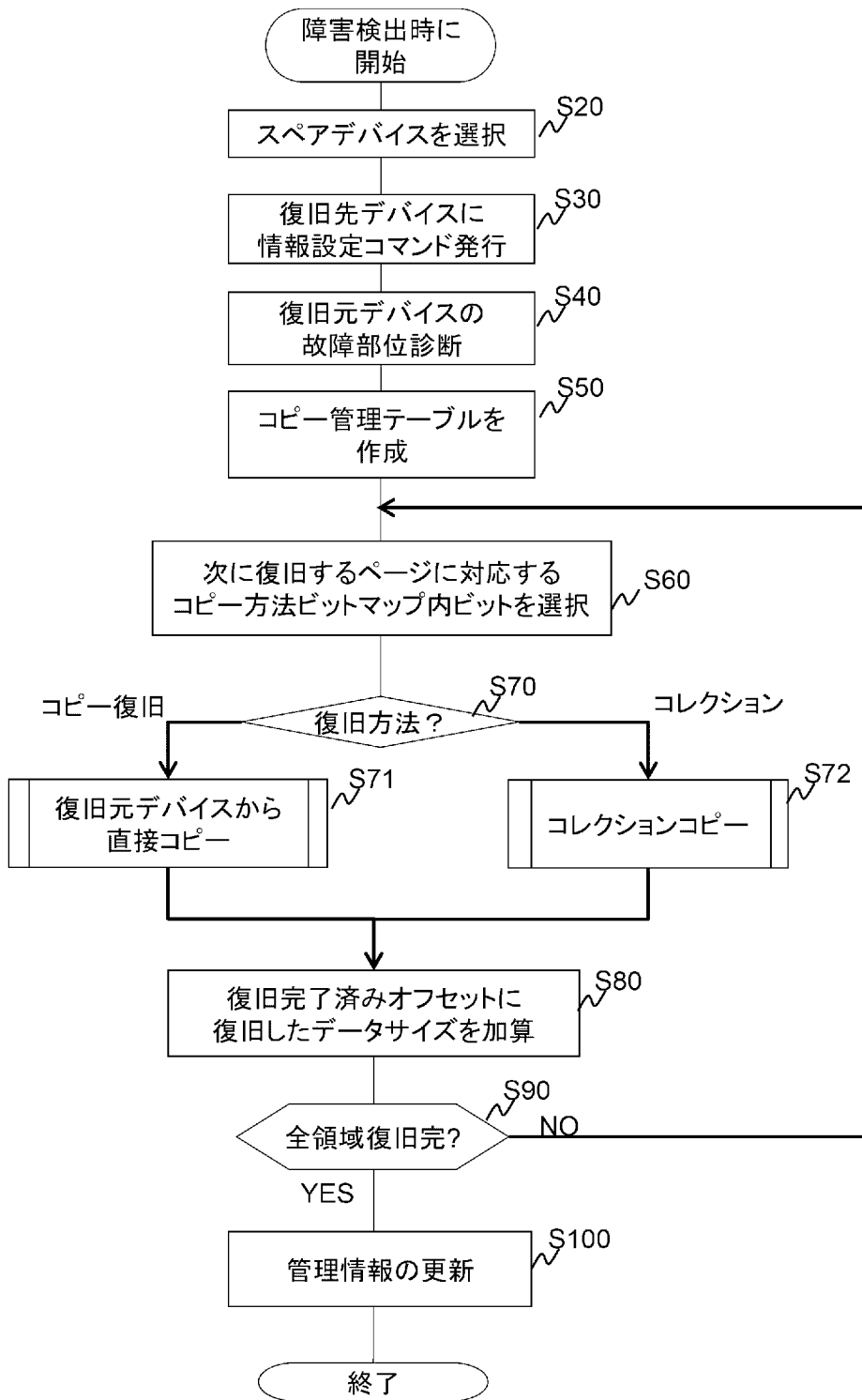
[図16]

図16



[図17]

図17



[図18]

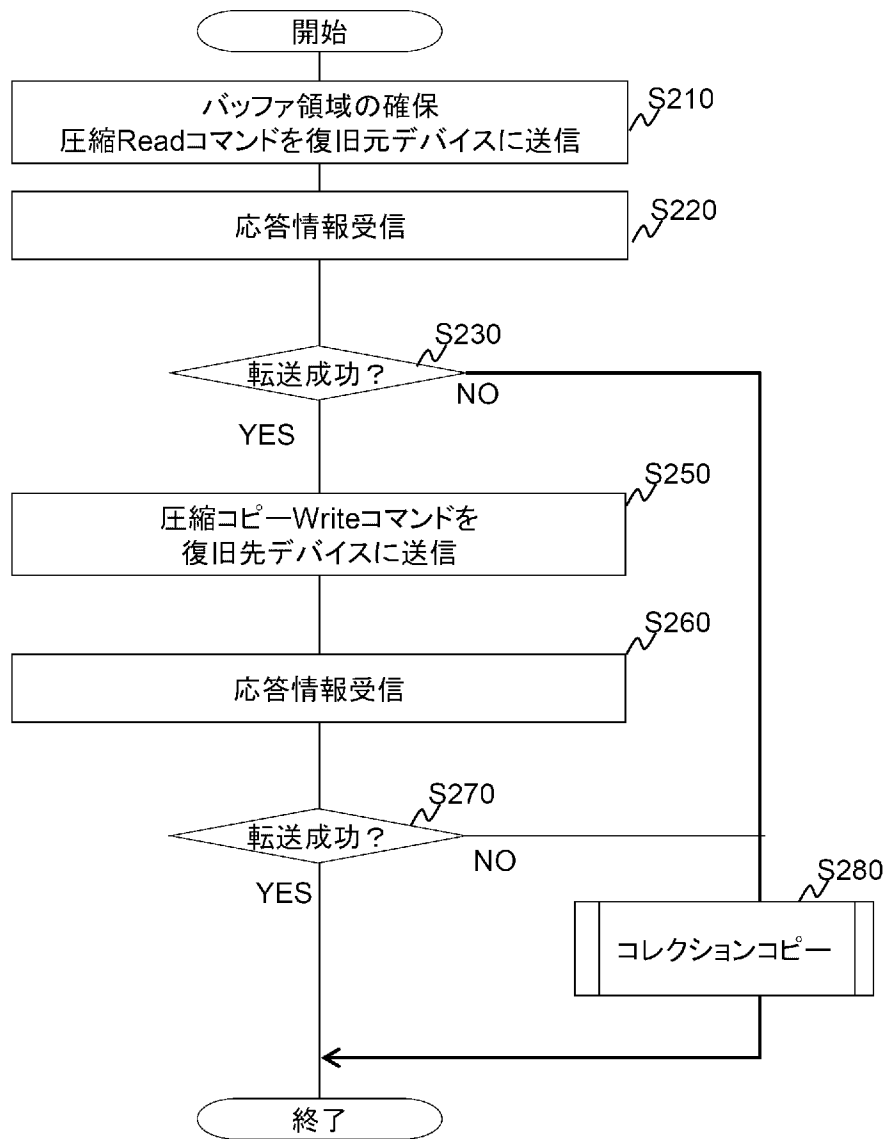
図18

~T1500

項目	内容	
障害RG#	0	~T1501
復旧元デバイス	FMPK#1	~T1502
復旧先デバイス	FMPK#4	~T1503
コレクション方式	方式3 (FMPK連携圧縮コピー)	~T1504
復旧デバイス容量	10240GB	~T1505
コピー方法ビットマップ 0:復旧元から直接コピー 1:コレクションコピー	000000100011111000...	~T1506
復旧完了済みオフセット	0x10500000	~T1507

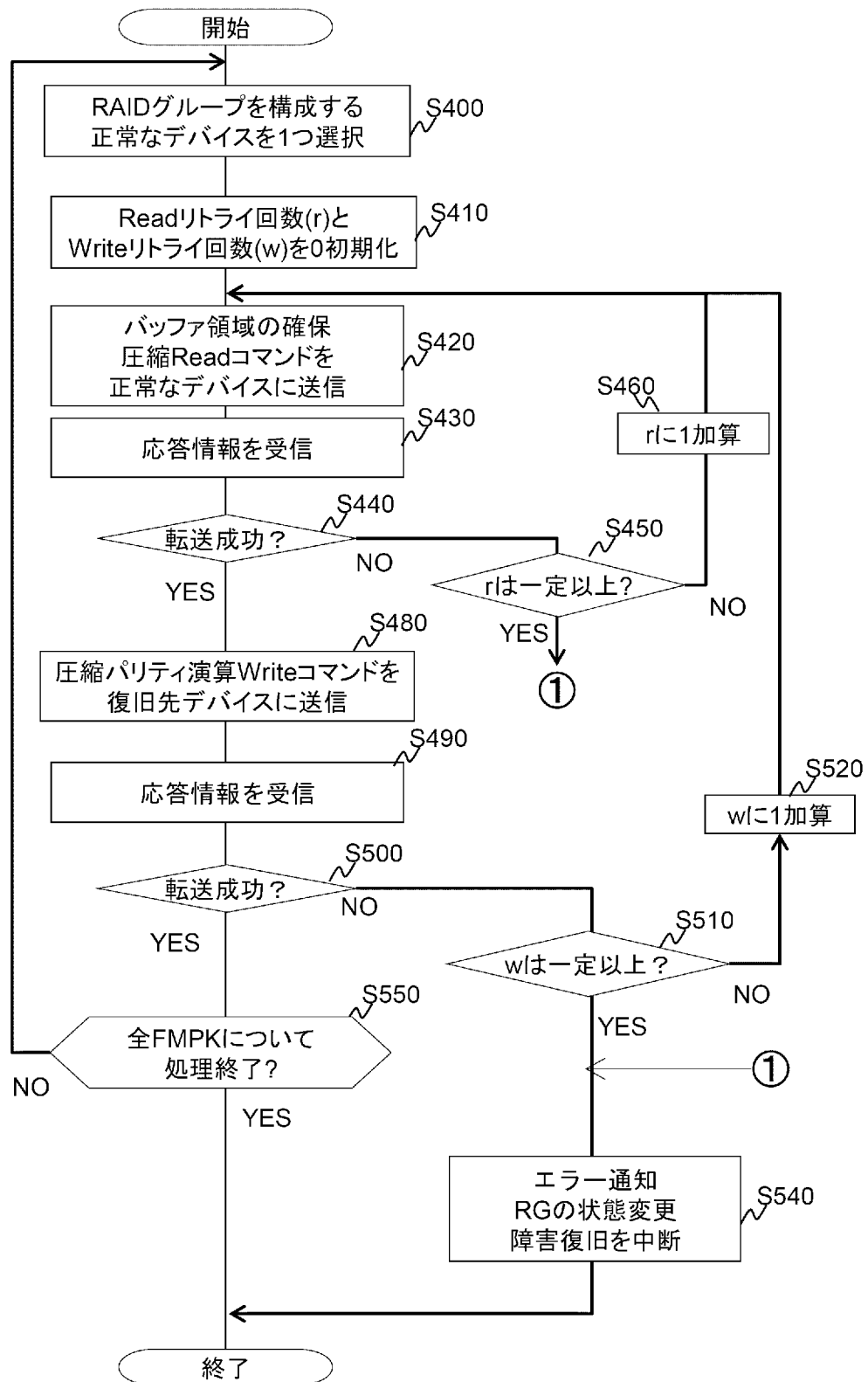
[図19]

図19



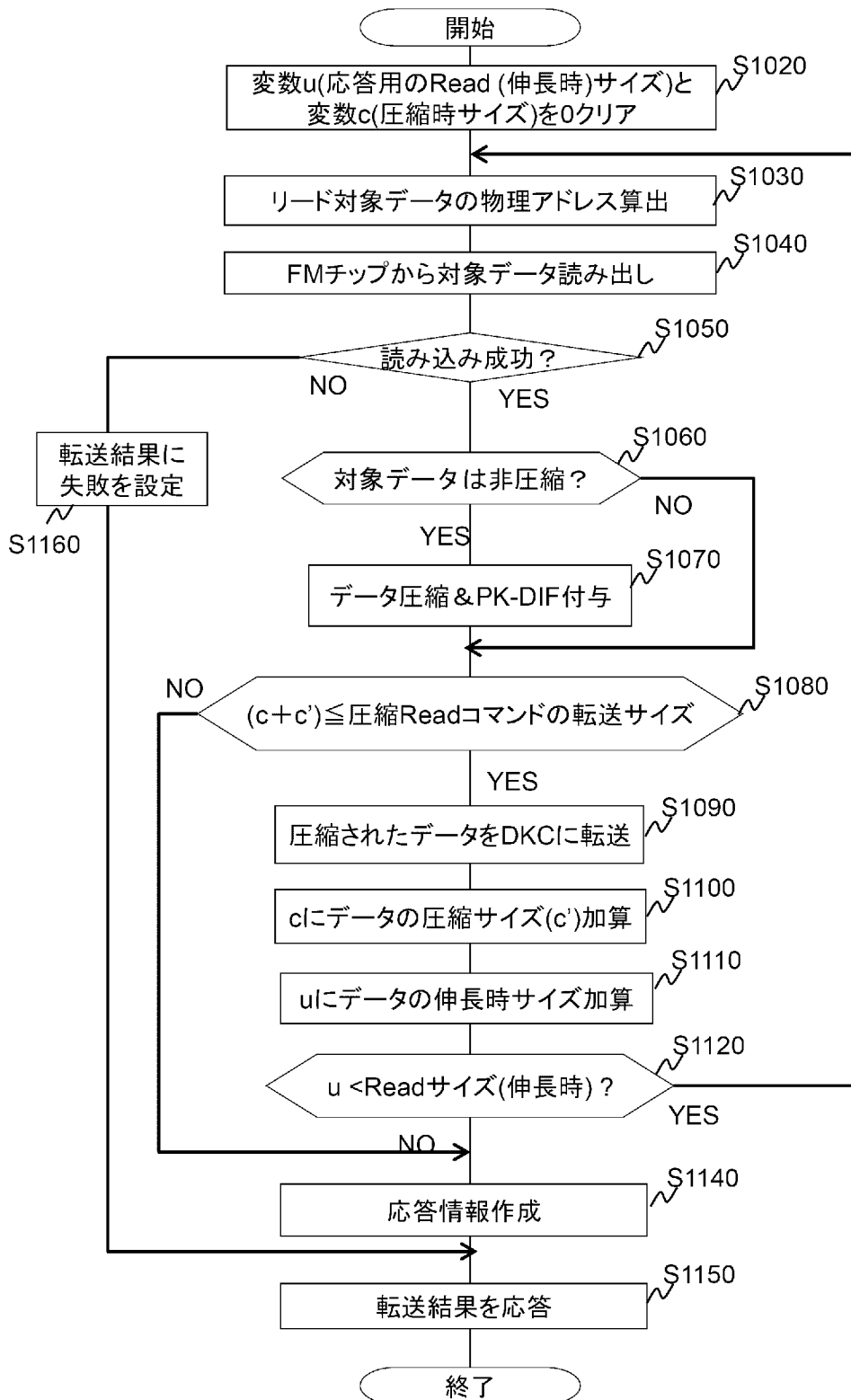
[図20]

図20



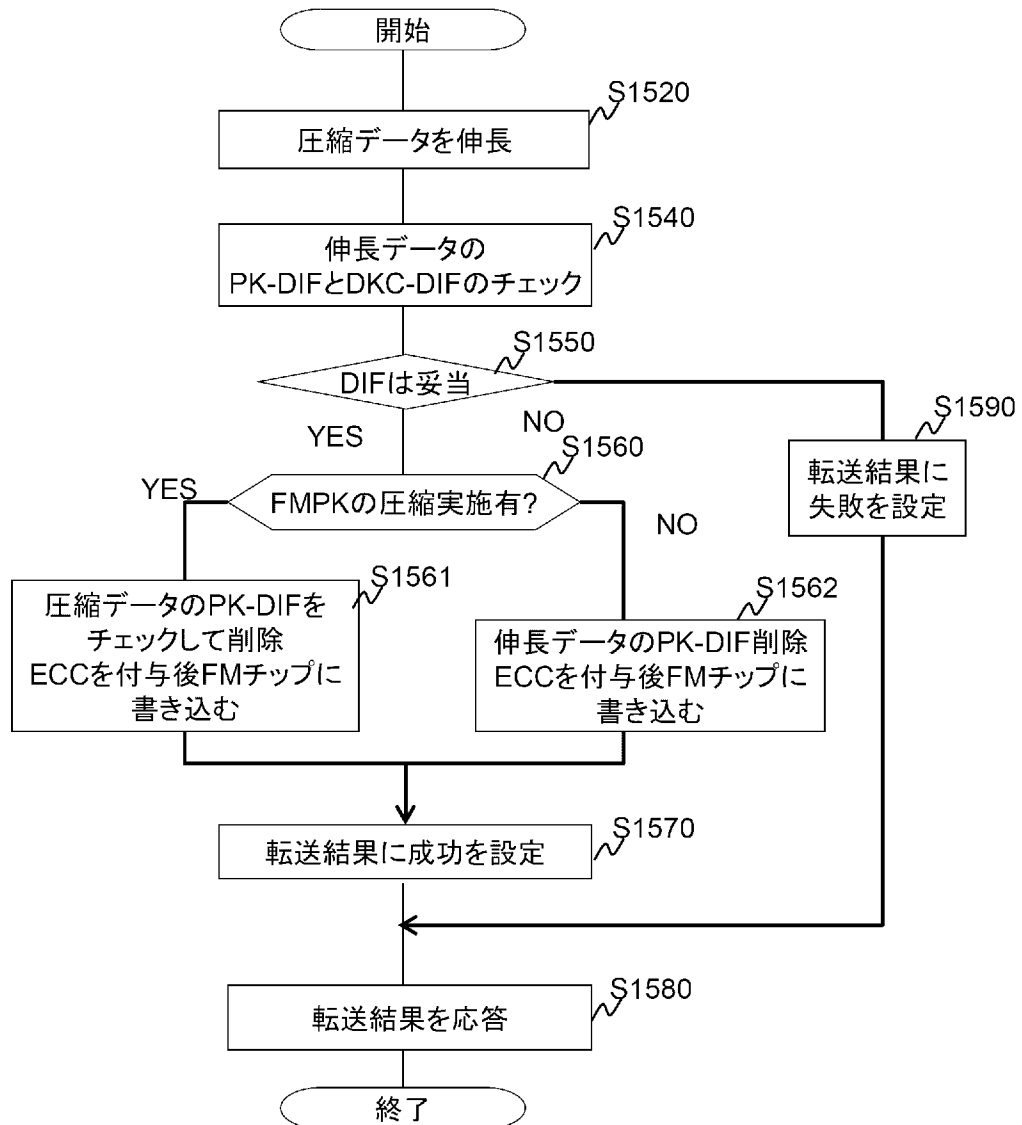
[図21]

図21



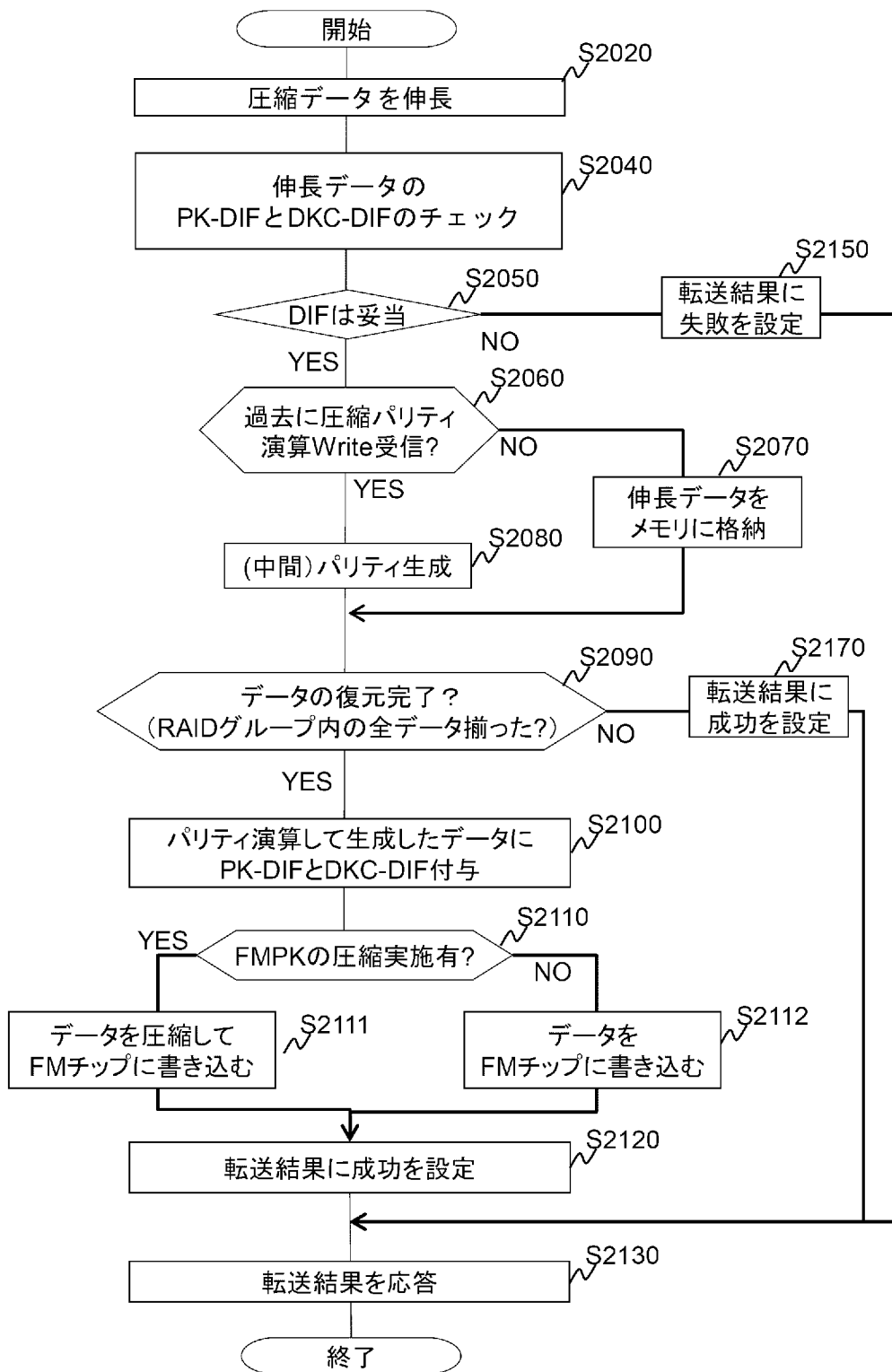
[図22]

図22



[図23]

図23



## [図24]

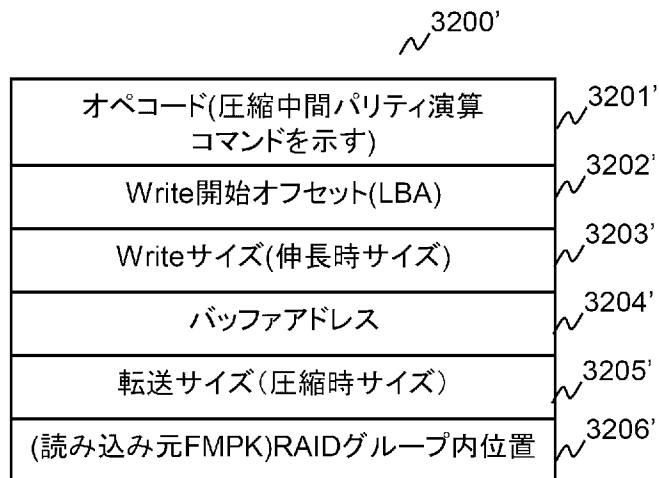
図24

~ T2500

~ T2501	~ T2502	~ T2503
論理ページ番号	アドレス	回数
0	0x80000000	1
1	0x80001000	2
2	NULL	0
:	:	

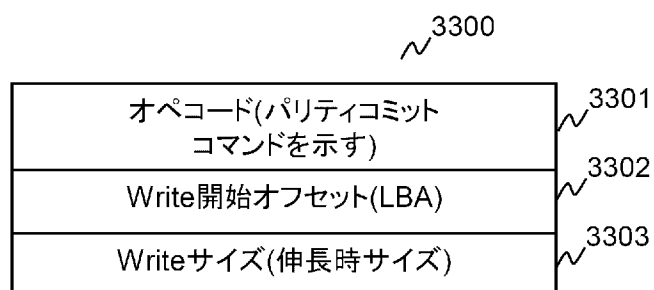
## [図25]

図25



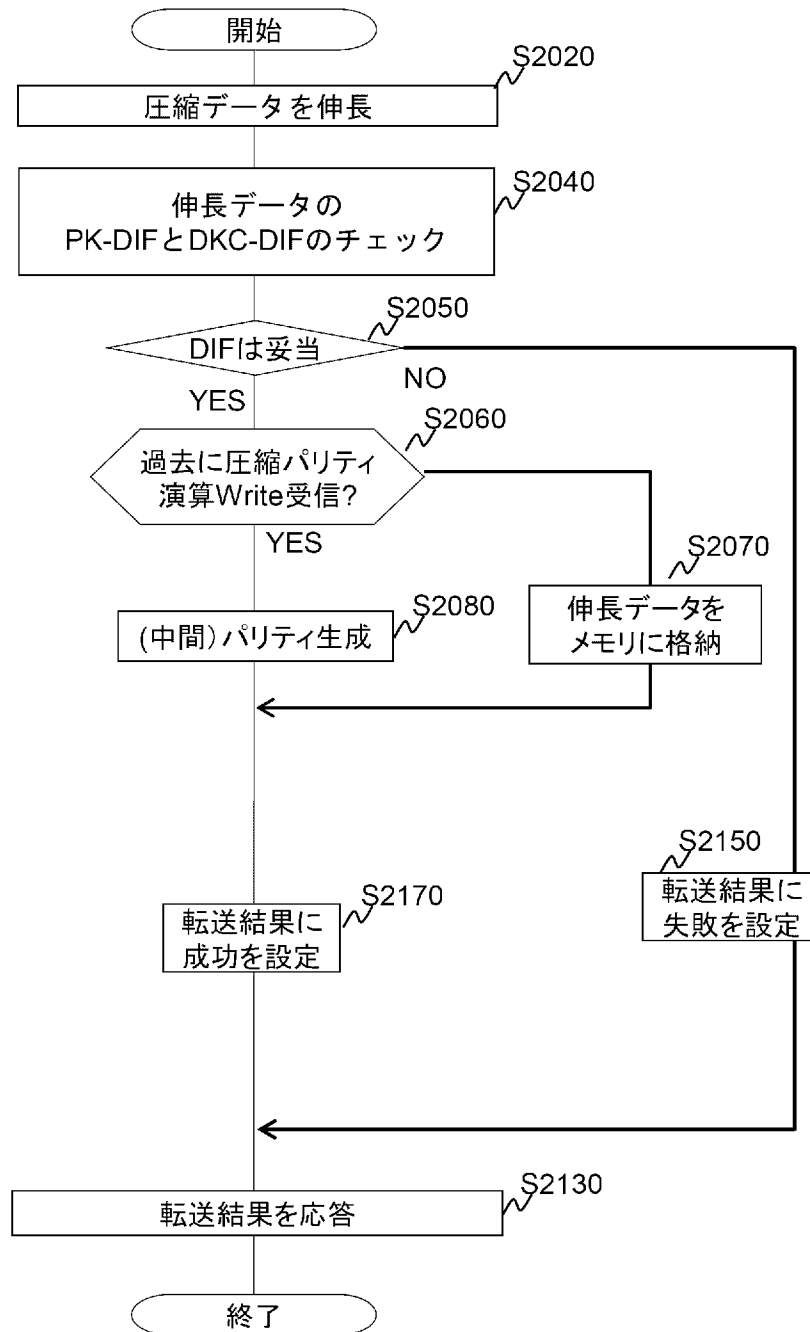
## [図26]

図26



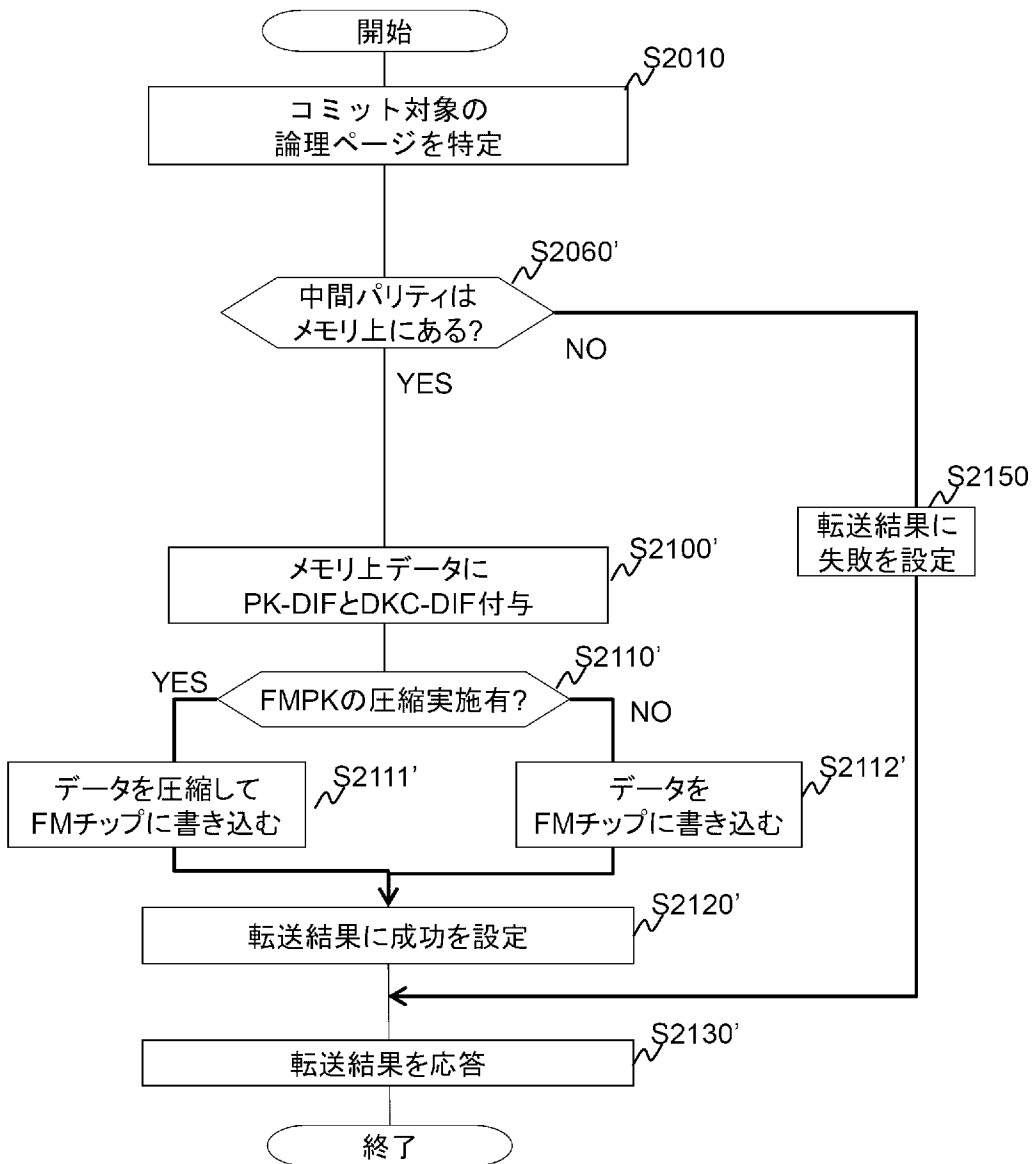
[図27]

図27



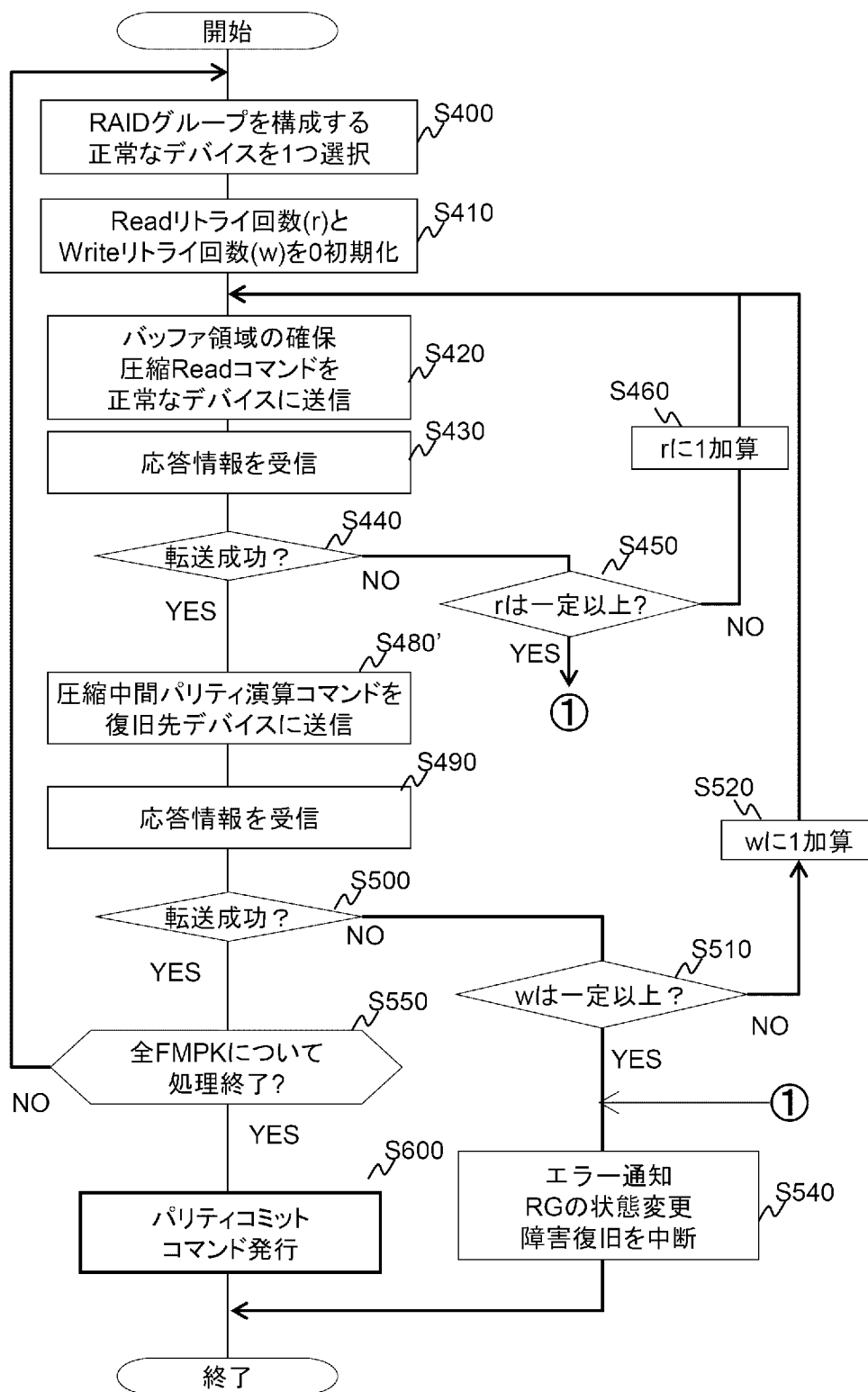
[図28]

図28



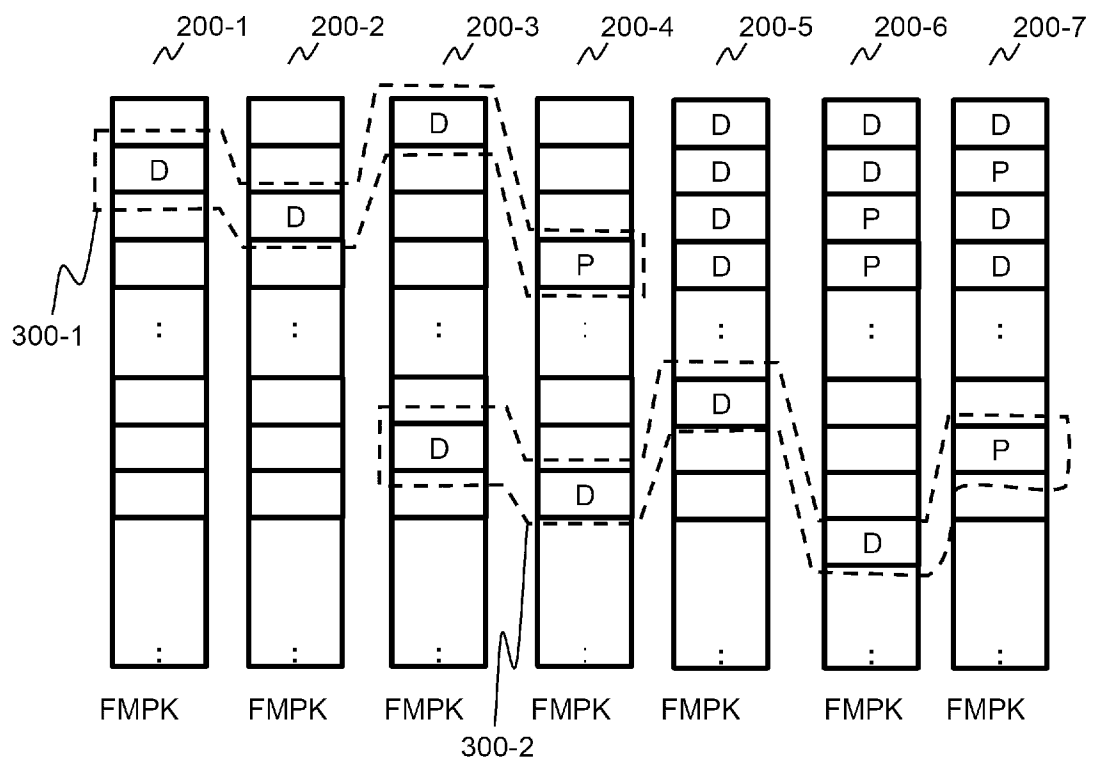
[図29]

図29



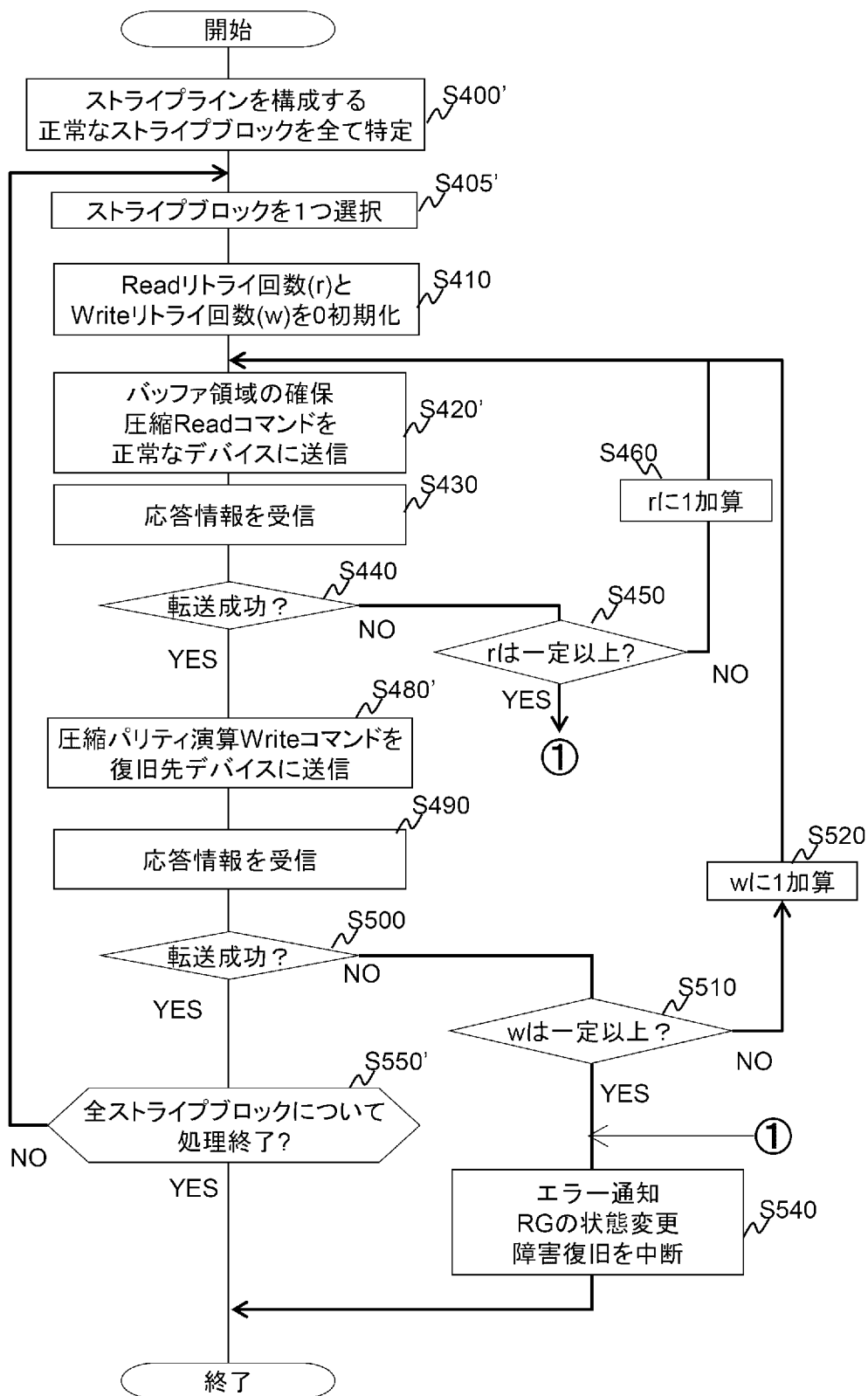
[図30]

図30



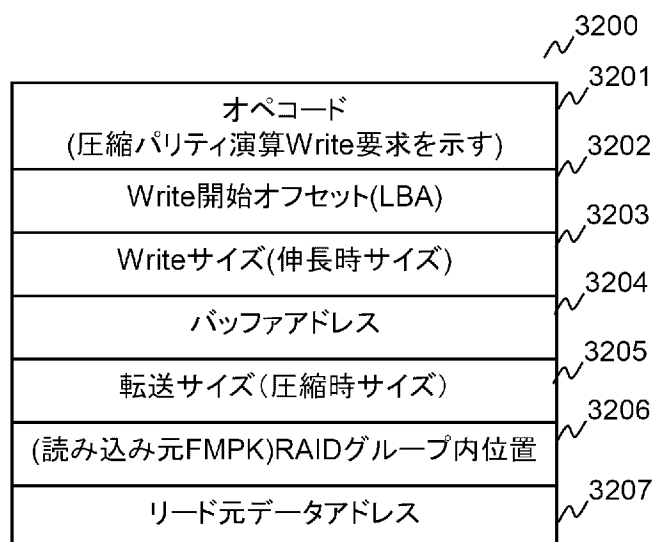
[図31]

図31



[図32]

図32



## INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2014/070224

A. CLASSIFICATION OF SUBJECT MATTER G06F3/06(2006.01) i		
According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED		
Minimum documentation searched (classification system followed by classification symbols) G06F3/06		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Jitsuyo Shinan Koho 1922-1996 Jitsuyo Shinan Toroku Koho 1996-2014 Kokai Jitsuyo Shinan Koho 1971-2014 Toroku Jitsuyo Shinan Koho 1994-2014		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y A	JP 2012-519319 A (Hitachi, Ltd.), 23 August 2012 (23.08.2012), paragraphs [0001] to [0081]; fig. 1 to 15 & US 2011/0238885 A1 & US 2012/0297244 A1 & WO 2010/137178 A1	1, 2, 4, 11 3, 5-10, 12-15
Y A	JP 5-011934 A (NEC Corp.), 22 January 1993 (22.01.1993), paragraphs [0007] to [0022]; fig. 1 to 4 (Family: none)	1, 2, 4, 11 3, 5-10, 12-15
Y	JP 7-311661 A (Fujitsu Ltd.), 28 November 1995 (28.11.1995), paragraphs [0004] to [0101]; fig. 1 to 56 & US 5859960 A1 & DE 19515661 A & KR 10-0226211 B	4
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/> See patent family annex.		
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search 03 October, 2014 (03.10.14)		Date of mailing of the international search report 14 October, 2014 (14.10.14)
Name and mailing address of the ISA/ Japanese Patent Office		Authorized officer
Facsimile No.		Telephone No.

**INTERNATIONAL SEARCH REPORT**

International application No.

PCT/JP2014/070224

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	JP 10-254634 A (Hitachi, Ltd.), 25 September 1998 (25.09.1998), paragraphs [0039] to [0041]; fig. 9 (Family: none)	11

A. 発明の属する分野の分類 (国際特許分類 (IPC)) Int.Cl. G06F3/06(2006.01)i		
B. 調査を行った分野 調査を行った最小限資料 (国際特許分類 (IPC)) Int.Cl. G06F3/06		
最小限資料以外の資料で調査を行った分野に含まれるもの 日本国実用新案公報 1922-1996年 日本国公開実用新案公報 1971-2014年 日本国実用新案登録公報 1996-2014年 日本国登録実用新案公報 1994-2014年		
国際調査で使用した電子データベース (データベースの名称、調査に使用した用語)		
C. 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
Y A	J P 2 0 1 2 - 5 1 9 3 1 9 A (株式会社日立製作所) 2 0 1 2 . 0 8 . 2 3 , 段落【0001】 - 【0081】, 図1-15 & US 2 0 1 1 / 0 2 3 8 8 8 5 A 1 & US 2 0 1 2 / 0 2 9 7 2 4 4 A 1 & WO 2 0 1 0 / 1 3 7 1 7 8 A 1	1, 2, 4, 11 3, 5-10, 12-15
Y A	J P 5 - 0 1 1 9 3 4 A (日本電気株式会社) 1 9 9 3 . 0 1 . 2 2 , 段落【0007】 - 【0022】 , 図1-4 (ファミリーなし)	1, 2, 4, 11 3, 5-10, 12-15
<input checked="" type="checkbox"/> C欄の続きにも文献が列挙されている。 <input type="checkbox"/> パテントファミリーに関する別紙を参照。		
* 引用文献のカテゴリー 「A」 特に関連のある文献ではなく、一般的技術水準を示すもの 「E」 国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの 「L」 優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献 (理由を付す) 「O」 口頭による開示、使用、展示等に言及する文献 「P」 国際出願日前で、かつ優先権の主張の基礎となる出願日の後に公表された文献 「T」 国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの 「X」 特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの 「Y」 特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの 「&」 同一パテントファミリー文献		
国際調査を完了した日 03.10.2014	国際調査報告の発送日 14.10.2014	
国際調査機関の名称及びあて先 日本国特許庁 (ISA/J P) 郵便番号100-8915 東京都千代田区霞が関三丁目4番3号	特許庁審査官 (権限のある職員) 古河 雅輝 電話番号 03-3581-1101 内線 3568	5 T 3 2 4 2

C (続き) . 関連すると認められる文献		
引用文献の カテゴリ*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
Y	JP 7-311661 A (富士通株式会社) 1995. 11. 28, 段落【0004】-【0101】, 図1-56 & US 5859960 A1 & DE 19515661 A & KR 10-0226211 B	4
Y	JP 10-254634 A (株式会社日立製作所) 1998. 09. 25, 段落【0039】-【0041】、図9 (ファミリーなし)	11