

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6233086号  
(P6233086)

(45) 発行日 平成29年11月22日 (2017.11.22)

(24) 登録日 平成29年11月2日 (2017.11.2)

(51) Int.Cl.	F I
<b>G 0 6 F 3 / 0 6 (2006.01)</b>	G O 6 F 3 / 0 6 3 O 6 Z
	G O 6 F 3 / 0 6 5 4 O
	G O 6 F 3 / 0 6 3 O 5 C
	G O 6 F 3 / 0 6 3 O 4 F

請求項の数 8 (全 31 頁)

(21) 出願番号	特願2014-30390 (P2014-30390)	(73) 特許権者	000005223
(22) 出願日	平成26年2月20日 (2014.2.20)		富士通株式会社
(65) 公開番号	特開2015-156081 (P2015-156081A)		神奈川県川崎市中原区上小田中4丁目1番1号
(43) 公開日	平成27年8月27日 (2015.8.27)	(74) 代理人	100092978
審査請求日	平成28年11月2日 (2016.11.2)		弁理士 真田 有
		(74) 代理人	100112678
			弁理士 山本 雅久
		(72) 発明者	小嵐 弘
			神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		審査官	桜井 茂行

最終頁に続く

(54) 【発明の名称】 ストレージ制御装置、ストレージシステム及び制御プログラム

(57) 【特許請求の範囲】

【請求項 1】

冗長構成がなされた複数の記憶装置及び複数の予備記憶装置と通信路を介して通信可能に接続されるストレージ制御装置であって、

前記複数の記憶装置のうち復元対象記憶装置以外の冗長用記憶装置から読み出した冗長データを用いて、前記復元対象記憶装置のデータを、前記複数の予備記憶装置のうちの第1の予備記憶装置に再構成する再構成処理部と、

前記再構成処理部による再構成を行なう際に、前記再構成処理部が前記冗長用記憶装置から読み出したデータを利用して、前記複数の予備記憶装置における、当該データを読み出した前記冗長用記憶装置に対応する第2の予備記憶装置に格納することで、前記冗長用記憶装置の複製を行なう複製処理部と  
を備えることを特徴とする、ストレージ制御装置。

【請求項 2】

リード要求受信時には、前記冗長用記憶装置と、当該冗長用記憶装置のデータを格納する前記第2の予備記憶装置とを併用することを特徴とする、請求項1記載のストレージ制御装置。

【請求項 3】

ライト要求受信時には、前記冗長用記憶装置及び当該冗長用記憶装置のデータを格納する前記予備記憶装置の双方に書き込みを行なうことを特徴とする、請求項1又は2記載のストレージ制御装置。

## 【請求項 4】

前記記憶装置に対して前記予備記憶装置を割り当てる割当処理部を備え、

前記複製処理部が、前記冗長用記憶装置から読み出したデータを、前記割当処理部が割り当てた前記予備記憶装置に格納することを特徴とする、請求項 1～3 のいずれか 1 項に記載のストレージ制御装置。

## 【請求項 5】

前記記憶装置に割り当て可能な前記予備記憶装置の数が前記記憶装置の数よりも少ない場合に、

前記割当処理部が、

前記予備記憶装置を、安定度の低い前記記憶装置から優先して割り当てることを特徴とする、請求項 4 記載のストレージ制御装置。

10

## 【請求項 6】

前記再構成処理部による再構成の完了後に、安定度の低い前記冗長用記憶装置に代えて、当該安定度の低い前記冗長用記憶装置の複製がされた前記予備記憶装置を用いて、前記冗長構成を変更する冗長構成変更部を備えることを特徴とする、請求項 1～請求項 5 のいずれか 1 項に記載のストレージ制御装置。

## 【請求項 7】

冗長構成がなされた複数の記憶装置と、

複数の予備記憶装置と、

前記複数の記憶装置のうち復元対象記憶装置以外の冗長用記憶装置から読み出した冗長データを用いて、前記復元対象記憶装置のデータを、前記複数の予備記憶装置のうちの第 1 の予備記憶装置に再構成する再構成処理部と、

20

前記再構成処理部による再構成を行なう際に、前記再構成処理部が前記冗長用記憶装置から読み出したデータを利用して、前記複数の予備記憶装置における、当該データを読み出した前記冗長用記憶装置に対応する第 2 の予備記憶装置に格納することで、前記冗長用記憶装置の複製を行なう複製処理部と  
を備えることを特徴とする、ストレージシステム。

## 【請求項 8】

冗長構成がなされた複数の記憶装置及び複数の予備記憶装置と通信路を介して通信可能に接続されるコンピュータに、

30

前記複数の記憶装置のうち復元対象記憶装置以外の冗長用記憶装置から読み出した冗長データを用いて、前記復元対象記憶装置のデータを、前記複数の予備記憶装置のうちの第 1 の予備記憶装置に再構成し、

前記再構成を行なう際に前記冗長用記憶装置から読み出したデータを利用して、前記複数の予備記憶装置における、当該データを読み出した前記冗長用記憶装置に対応する第 2 の予備記憶装置に格納することで、前記冗長用記憶装置の複製を行なう  
処理を実行させることを特徴とする、制御プログラム。

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

40

本発明は、ストレージ制御装置、ストレージシステム及び制御プログラムに関する。

## 【背景技術】

## 【0002】

情報通信技術 (Information and Communication Technology ; I C T ) システムの普及に伴い、近年、Hard Disk Drive ( H D D ) に代表される記憶装置 ( 以下、「ディスク」と総称する ) を複数使用するディスクアレイ装置が広く用いられるようになってきている。このようなディスクアレイ装置では、一般に、Redundant Arrays of Inexpensive Disks ( R A I D ) 技術を用いて、データが 2 台以上のディスクに冗長化されて記録されることにより、データの安全性が担保されている。

## 【0003】

50

ここで、RAIDとは、複数のディスクを組合せて、仮想的な1台のディスク(RAIDグループ)として管理する技術を指す。RAIDには、各ディスクへのデータ配置及び冗長性に依りて、RAID0~RAID6のレベルが存在する。

RAID装置では、RAIDを構成する複数ディスクにまたがるストライプ内にパリティデータを入れてRAIDを構成するディスクの故障に対してデータを保護する。そして、このRAID構成上にLUN(Logical Unit Number)を割り当ててサーバにディスク領域として見せて使用させている。

#### 【0004】

データが冗長化されたディスクアレイ装置において、ディスクが故障すると、故障したディスクに記憶されていたデータが再構築されて、予備ディスクなどの代替ディスクに格納される。このような処理は、一般にリビルド処理と呼ばれる。リビルド処理が実行されることで、データの冗長性が回復する。

#### 【先行技術文献】

#### 【特許文献】

#### 【0005】

【特許文献1】特開平10-293658号公報

【特許文献2】特開2005-78430号公報

#### 【発明の概要】

#### 【発明が解決しようとする課題】

#### 【0006】

しかしながら、このような従来のディスクアレイ装置において、例えば、RAID上に複数のLUNがある場合に、リビルド処理により一部のLUNの復元が完了していても、RAID内で更にもう1本ディスクが故障した場合には、RAID内の全データがロストし、復元されていたLUNのデータも失われてしまう。

また、RAID装置上に複数のLUNがある場合に、復元が完了しているLUNであっても、リビルドの全処理が終わるまでは、このリビルド処理による性能劣化の影響を受ける。

#### 【0007】

1つの側面では、本発明は、記憶装置の故障時における信頼性を向上させることを目的とする。

なお、前記目的に限らず、後述する発明を実施するための形態に示す各構成により導かれる作用効果であって、従来の技術によっては得られない作用効果を奏することも本発明の他の目的の1つとして位置付けることができる。

#### 【課題を解決するための手段】

#### 【0008】

このため、このストレージ制御装置は、冗長構成がなされた複数の記憶装置及び複数の予備記憶装置と通信路を介して通信可能に接続されるストレージ制御装置であって、前記複数の記憶装置のうち復元対象記憶装置以外の冗長用記憶装置から読み出した冗長データを用いて、前記復元対象記憶装置のデータを、前記複数の予備記憶装置のうちの第1の予備記憶装置に再構成する再構成処理部と、前記再構成処理部による再構成を行なう際に、前記再構成処理部が前記冗長用記憶装置から読み出したデータを利用して、前記複数の予備記憶装置における、当該データを読み出した前記冗長用記憶装置に対応する第2の予備記憶装置に格納することで、前記冗長用記憶装置の複製を行なう複製処理部とを備える。

#### 【発明の効果】

#### 【0009】

一実施形態によれば、記憶装置の故障時における信頼性を向上させることができる。

#### 【図面の簡単な説明】

#### 【0010】

【図1】実施形態の一例としてのストレージ装置を備えるストレージシステムのハードウェア構成を模式的に示す図である。

10

20

30

40

50

【図 2】実施形態の一例としてのストレージ装置の機能構成を示す図である。

【図 3】実施形態の一例としてのストレージ装置における L U N 管理テーブルの構成を例示する図である。

【図 4】実施形態の一例としてのストレージ装置におけるリビルド処理を説明する図である。

【図 5】実施形態の一例としてのストレージ装置におけるリビルド処理を説明する図である。

【図 6】実施形態の一例としてのストレージ装置におけるリビルド処理を説明する図である。

【図 7】実施形態の一例としてのストレージ装置における R A I D 構成変更部による R A I D 構成変更方法を示す図である。

10

【図 8】実施形態の一例としてのストレージ装置におけるリビルド処理の概要を説明するフローチャートである。

【図 9】実施形態の一例としてのストレージ装置におけるストレージ構成の決定方法をフローチャートである。

【図 10】実施形態の一例としてのストレージ装置におけるリビルド処理を説明するフローチャートである。

【図 11】実施形態の一例としてのストレージ装置における予備ディスクの解放要求の有無の確認処理の詳細を説明するフローチャートである。

【図 12】実施形態の一例としてのストレージ装置におけるリビルド処理後の処理を説明するフローチャートである。

20

【図 13】実施形態の一例としてのストレージ装置におけるリード受信時の処理を説明するフローチャートである。

【図 14】実施形態の一例としてのストレージ装置におけるライト受信時の処理を説明する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、図面を参照して本ストレージ制御装置、ストレージシステム及び制御プログラムに係る実施の形態を説明する。ただし、以下に示す実施形態はあくまでも例示に過ぎず、実施形態で明示しない種々の変形例や技術の適用を排除する意図はない。すなわち、本実施形態を、その趣旨を逸脱しない範囲で種々変形して実施することができる。又、各図は、図中に示す構成要素のみを備えるという趣旨ではなく、他の機能等を含むことができる。

30

【0012】

図 1 は実施形態の一例としてのストレージ装置 1 を備えるストレージシステム 4 のハードウェア構成を模式的に示す図である。

ストレージシステム 4 においては、ストレージ装置 1 と 1 つ以上（図 1 に示す例では 2 つ）のホスト装置 2 a , 2 b とが冗長化された複数のパスを介して接続されている。

ストレージ装置 1 は、ドライブエンクロージャ（D E : Drive Enclosure）3 0 に格納された記憶装置 3 1 を仮想化して、仮想ストレージ環境を形成する。そしてストレージ装置 1 は、仮想ボリュームを上位装置であるホスト装置 2 a , 2 b に提供する。

40

【0013】

ホスト装置 2 a , 2 b は、例えば、サーバ機能をそなえた情報処理装置であり、本ストレージ装置 1 との間において、N A S（Network Attached Storage）や S A N（Storage Area Network）のコマンドを送受信する。これらのホスト装置 2 a , 2 b は、同様の構成を有している。

以下、ホスト装置を示す符号としては、複数のホスト装置のうち 1 つを特定する必要があるときには符号 2 a , 2 b を用いるが、任意のホスト装置を指すときには符号 2 を用いる。

【0014】

50

ホスト装置 2 は、図示しない C P U (Central Processing Unit) やメモリを備え、C P U がメモリ等に格納された O S (Operating System) やプログラムを実行することで、種々の機能を実行する。

ホスト装置 2 は、例えば、ストレージ装置 1 に対して N A S におけるリード/ライト等のディスクアクセスコマンドを送信することにより、ストレージ装置 1 が提供するボリュームにデータの書き込みや読み出しを行なう。

【 0 0 1 5 】

そして、本ストレージ装置 1 は、ホスト装置 2 からボリュームに対して行なわれる入出力要求 (例えば、リードコマンドやライトコマンド) に応じて、このボリュームに対応する実ストレージに対して、データの読み出しや書き込み等の処理を行なう。なお、ホスト装置 2 からの入出力要求のことを I O コマンドという場合がある。

スイッチ 3 a , 3 b は、ホスト装置 2 a , 2 b とストレージ装置 1 のコントローラ 1 0 0 との通信を中継する中継装置である。各スイッチ 3 a , 3 b は、それぞれホスト装置 2 a , 2 b に接続されるとともに、コントローラ 1 0 0 に接続されている。

【 0 0 1 6 】

図 1 に示す例においては、コントローラ 1 0 0 に 2 つのポート 1 0 1 a , 1 0 1 b が備えられている。そして、ポート 1 0 1 a にスイッチ 3 a が、又、ポート 1 0 1 b にスイッチ 3 b が、それぞれ接続され、更に、各スイッチ 3 a , 3 b にそれぞれホスト装置 2 a , 2 b が接続されている。

ストレージ装置 1 は、図 1 に示すように、1 つ以上 (本実施形態では 1 つ) のコントローラ 1 0 0 及び 1 つ以上 (図 1 に示す例では 1 つ) のドライブエンクロージャ (Drive Enclosure : D E ) 3 0 をそなえる。

【 0 0 1 7 】

ドライブエンクロージャ 3 0 には、1 以上 (図 1 に示す例では 8 つ) の記憶装置 (物理ディスク) 3 1 a - 1 ~ 3 1 a - 4 , 3 1 b - 1 ~ 3 1 b - 4 が搭載され、これらの記憶装置 3 1 a - 1 ~ 3 1 a - 4 , 3 1 b - 1 ~ 3 1 b - 4 の記憶領域 (実ボリューム, 実ストレージ) を、本ストレージ装置 1 に対して提供する。

以下、記憶装置を示す符号としては、複数の記憶装置のうち 1 つを特定する必要があるときには符号 3 1 a - 1 ~ 3 1 a - 4 , 3 1 b - 1 ~ 3 1 b - 4 を用いるが、任意の記憶装置を指すときには符号 3 1 を用いる。又、記憶装置 3 1 をディスク 3 1 という場合がある。

【 0 0 1 8 】

また、以下、記憶装置 3 1 a - 1 , 3 1 a - 2 , 3 1 a - 3 , 3 1 a - 4 を、disk1 , disk2 , disk3 , disk4 とそれぞれ表す場合がある。更に、以下、記憶装置 3 1 b - 1 , 3 1 b - 2 , 3 1 b - 3 , 3 1 b - 4 を、disk1 ' , disk2 ' , disk3 ' , disk4 ' とそれぞれ表す場合がある。

記憶装置 3 1 は、ハードディスクドライブ (Hard disk drive : H D D )、S S D (Solid State Drive) 等の記憶装置であって、種々のデータを格納するものである。

【 0 0 1 9 】

ドライブエンクロージャ 3 0 は、例えば複数段のスロット (図示省略) を備えて構成され、これらのスロットに記憶装置 3 1 を挿入することにより、実ボリューム容量を随時変更することができる。

また、ドライブエンクロージャ 3 0 に備えられた複数の記憶装置 3 1 を用いて R A I D (Redundant Arrays of Inexpensive Disks) が構成される。図 1 に示す例においては、記憶装置 3 1 a - 1 ~ 3 1 a - 4 を用いて R A I D が構成されており、これらの記憶装置 3 1 a - 1 ~ 3 1 a - 4 が R A I D グループ 3 0 a を構成している。

【 0 0 2 0 】

この R A I D グループ 3 0 a を構成する記憶装置 3 1 a - 1 ~ 3 1 a - 4 を R A I D 構成ディスク 3 1 a という場合がある。

記憶装置 3 1 b - 1 ~ 3 1 b - 4 は、R A I D ディスクグループ内のディスク故障にそ

10

20

30

40

50

なえて予備的に設けられた予備ディスクであり、ホットスペア (Hot spare ; H S ) として用いられる。これらの記憶装置 3 1 b - 1 ~ 3 1 b - 4 が予備ディスクグループ 3 0 b を構成している。以下、記憶装置 3 1 b を予備ディスク 3 1 b という場合がある。

#### 【 0 0 2 1 】

ドライブエンクロージャ 3 0 は、コントローラ 1 0 0 のデバイスアダプタ ( Device Adapter : D A ) 1 0 3 , 1 0 3 とそれぞれ接続されている。

コントローラエンクロージャ 4 0 は、1 以上 ( 図 1 に示す例では 1 つ ) のコントローラ 1 0 0 を備える。

コントローラ 1 0 0 は、ストレージ装置 1 内の動作を制御するストレージ制御装置であり、ホスト装置 2 から送信される I O コマンドに従って、ドライブエンクロージャ 3 0 の記憶装置 3 1 へのアクセス制御等、各種制御を行なう。

10

#### 【 0 0 2 2 】

なお、図 1 に示す例においては、ストレージ装置 1 に 1 つのコントローラ 1 0 0 が備えられているが、これに限定されるものではなく、2 つ以上のコントローラ 1 0 0 を備えてもよい。すなわち、複数のコントローラ 1 0 0 により冗長化を行ない、通常は、いずれかのコントローラ 1 0 0 がプライマリとして各種制御を行ない、このプライマリコントローラ 1 0 0 の故障時には、セカンダリのコントローラ 1 0 0 がプライマリとしての動作を引き継いでもよい。

#### 【 0 0 2 3 】

コントローラ 1 0 0 はポート 1 0 1 a , 1 0 1 b を介してホスト装置 2 に接続される。そして、コントローラ 1 0 0 は、ホスト装置 2 から送信されるリード/ライト等の I O コマンドを受信し、D A 1 0 3 等を介して記憶装置 3 1 の制御を行なう。

20

コントローラ 1 0 0 は、図 1 に示すように、ポート 1 0 1 a , 1 0 1 b と複数 ( 図 1 に示す例では 2 つ ) の D A 1 0 3 , 1 0 3 とをそなえるとともに、C P U 1 1 0 , メモリ 1 0 6 , S S D 1 0 7 及び I O C ( Input Output Controller ) 1 0 8 をそなえる。

#### 【 0 0 2 4 】

以下、ポートを示す符号としては、複数のポートのうち 1 つを特定する必要があるときには符号 1 0 1 a , 1 0 1 b を用いるが、任意のポートを指すときには符号 1 0 1 を用いる。

ポート 1 0 1 は、ホスト装置 2 等から送信されたデータを受信したり、コントローラ 1 0 0 から出力するデータをホスト装置 2 等に送信する。すなわち、ポート 1 0 1 は、ホスト装置等の外部装置との間でのデータの入出力 ( I / O ) を制御する。

30

#### 【 0 0 2 5 】

ポート 1 0 1 a は、S A N を介してホスト装置 2 と通信可能に接続され、例えば、F C ( Fibre Channel ) インタフェース等のネットワークアダプタに備えられる F C ポートである。

ポート 1 0 1 b は、N A S を介してホスト装置 2 と通信可能に接続され、例えば、L A N ( Local Area Network ) インタフェースの N I C ポートである。コントローラ 1 0 0 は、これらのポート 1 0 1 により通信回線を介してホスト装置 2 等と接続され、I / O コマンドの受信やデータの送受信等を行なう。

40

#### 【 0 0 2 6 】

そして、ポート 1 0 1 a にスイッチ 3 a が、又、ポート 1 0 1 b にスイッチ 3 b が、それぞれ接続され、更に、各スイッチ 3 a , 3 b にそれぞれホスト装置 2 a , 2 b が接続されている。

すなわち、ホスト装置 2 a , 2 b は、それぞれスイッチ 3 a を介してポート 1 0 1 a に接続されるとともに、スイッチ 3 b を介してポート 1 0 1 b に接続されている。

#### 【 0 0 2 7 】

なお、図 1 に示す例においてはコントローラ 1 0 0 に 2 つのポート 1 0 1 a , 1 0 1 b が備えられているが、これに限定されるものではなく、1 つもしくは 3 つ以上のポートを備えてもよい。

50

ＤＡ１０３は、ドライブエンクロージャ３０や記憶装置３１等と通信可能に接続するためのインタフェースである。ＤＡ１０３にはドライブエンクロージャ３０の記憶装置３１が接続され、コントローラ１００は、ホスト装置２から受信したＩＯコマンドに基づき、これらの記憶装置３１に対するアクセス制御を行なう。

【００２８】

コントローラ１００は、これらのＤＡ１０３を介して、記憶装置３１に対するデータの書き込みや読み出しを行なう。又、図１に示す例においては、コントローラ１００に２つのＤＡ１０３がそなえられている。そして、コントローラ１００において、各ＤＡ１０３にドライブエンクロージャ３０が接続されている。これにより、ドライブエンクロージャ３０の記憶装置３１には、コントローラ１００からデータの書き込みや読み出しを行なうことができる。

10

【００２９】

ＳＳＤ１０７は、ＣＰＵ１１０が実行するプログラムや種々のデータ等を格納する記憶装置である。

メモリ１０６は、種々のデータやプログラムを一時的に格納する記憶装置であり、図示しないメモリ領域とキャッシュ領域とをそなえる。キャッシュ領域は、ホスト装置２から受信したデータや、ホスト装置２に対して送信するデータを一時的に格納する。メモリ領域には、ＣＰＵ１１０がプログラムを実行する際に、データやプログラムを一時的に格納・展開して用いる。

【００３０】

20

このメモリ１０６には、後述するＲＡＩＤ制御部１２によるＲＡＩＤ制御に用いられる、仮／実ボリューム変換テーブル６２やディスク構成情報６３，ＲＡＩＤ構成テーブル６４等が格納される。仮／実ボリューム変換テーブル６２は、ホスト装置２に提供する仮想ボリュームのアドレスを、記憶装置３１の物理アドレス（実アドレス）にマッピングしているテーブルである。

【００３１】

ディスク構成情報６３は、本ストレージ装置１に備えられる記憶装置３１を管理する情報である。ディスク構成情報６３において、例えば、各記憶装置３１のディスク種別や、各記憶装置３１がドライブエンクロージャ３０のどのスロットに取り付けられているか、又、どの記憶装置３１が予備ディスク３１であるか等の情報が管理される。

30

ＲＡＩＤ構成テーブル６４は、後述するＲＡＩＤ制御部１２がＲＡＩＤの管理を行なうために用いる情報であり、例えば、ＲＡＩＤ種別や、ＲＡＩＤグループ３０ａを構成する記憶装置３１ａを特定する情報等が格納される。

【００３２】

なお、これらの仮／実ボリューム変換テーブル６２やディスク構成情報６３，ＲＡＩＤ構成テーブル６４は既知であるので、その詳細な説明は省略する。

また、メモリ１０６には、後述するリビルド処理部１３がリビルド処理を行なう際に、各リビルド元ディスク３１ａから読み出されたデータが一時的に格納される。

さらに、メモリ１０６には、後述する、ＬＵＮ管理テーブル６１が格納される。ＬＵＮ管理テーブル６１の詳細は後述する。

40

【００３３】

ＩＯＣ１０８は、コントローラ１００内におけるデータ転送を制御する制御装置であり、例えば、メモリ１０６に格納されたデータをＣＰＵ１１０を介することなく転送させるＤＭＡ転送（Direct Memory Access）を実現する。

ＣＰＵ１１０は、種々の制御や演算を行なう処理装置であり、例えばマルチコアプロセッサ（マルチＣＰＵ）である。ＣＰＵ１１０は、ＳＳＤ１０７等に格納されたＯＳやプログラムを実行することにより、種々の機能を実現する。

【００３４】

図２は実施形態の一例としてのストレージ装置１の機能構成を示す図である。

この図２に示すように、ストレージ装置１は、ＩＯ制御部１１，ＲＡＩＤ制御部１２，

50

リビルド処理部 13, 割当処理部 14, ミラーリング処理部 15 及び R A I D 構成変更部 16 として機能する。

そして、コントローラ 100 の C P U 110 が、制御プログラムを実行することにより、これらの I O 制御部 11, R A I D 制御部 12, リビルド処理部 13, 割当処理部 14, ミラーリング処理部 15 及び R A I D 構成変更部 16 として機能する。

【0035】

なお、これらの I O 制御部 11, R A I D 制御部 12, リビルド処理部 13, 割当処理部 14, ミラーリング処理部 15 及び R A I D 構成変更部 16 としての機能を実現するためのプログラム(制御プログラム)は、例えばフレキシブルディスク, C D ( C D - R O M, C D - R, C D - R W 等), D V D ( D V D - R O M, D V D - R A M, D V D - R, D V D + R, D V D - R W, D V D + R W, H D D V D 等), ブルーレイディスク, 磁気ディスク, 光ディスク, 光磁気ディスク等の、コンピュータ読取可能な記録媒体に記録された形態で提供される。そして、コンピュータはその記録媒体からプログラムを読み取って内部記憶装置または外部記憶装置に転送し格納して用いる。又、そのプログラムを、例えば磁気ディスク, 光ディスク, 光磁気ディスク等の記憶装置(記録媒体)に記録しておき、その記憶装置から通信経路を介してコンピュータに提供するようにしてもよい。

【0036】

I O 制御部 11, R A I D 制御部 12, リビルド処理部 13, 割当処理部 14, ミラーリング処理部 15 及び R A I D 構成変更部 16 としての機能を実現する際には、内部記憶装置(本実施形態では S S D 107 やメモリ 106)に格納されたプログラムがコンピュータのマイクロプロセッサ(本実施形態では C P U 110)によって実行される。このとき、記録媒体に記録されたプログラムをコンピュータが読み取って実行するようにしてもよい。

【0037】

R A I D 制御部 12 は、記憶装置 31 a を用いて R A I D を実現し、R A I D を構成する記憶装置 31 の制御を行なう。すなわち、R A I D 制御部 12 は、複数の記憶装置 31 a を用いた冗長構成を設定する。

この R A I D 制御部 12 は、上述した R A I D 構成テーブル 64 を作成・管理して、R A I D 構成ディスク 31 a を用いて R A I D グループ 30 a を設定し、各種 R A I D 制御を行なう。なお、R A I D 制御部 12 による R A I D の管理は既知の手法で実現することができ、その説明は省略する。

【0038】

また、R A I D 制御部 12 は、記憶装置 31 を用いて L U N の設定・管理を行ない、ホスト装置 2 はこの設定された L U N に対してデータアクセスを行なう。R A I D 制御部 12 は、L U N 管理テーブル 61 を用いて、L U N の管理を行なう。

図 3 は実施形態の一例としてのストレージ装置 1 における L U N 管理テーブル 61 の構成を例示する図である。

【0039】

L U N 管理テーブル 61 は L U N 毎に備えられる管理情報であり、この図 3 に例示する L U N 管理テーブル 61 は L U N 1 についての情報を示す。

L U N 管理テーブル 61 は、項目と内容と備考とを関連付けて構成されている。図 3 に示す例においては、L U N 管理テーブル 61 に 12 個の項目が登録されており、これらの項目を特定するための項番 1 ~ 12 が設定されている。

【0040】

具体的には、項目として、L U N 名, 構成ディスク名リスト, 各ディスク上の位置、サイズ, 故障ディスク名, L U N の状態, 利用予備ディスク名リスト, ディスクの安定度リスト, リビルド済サイズ, 元ディスクのリード/ライト中カウンタ, 予備ディスクのリード/ライト中カウンタ, 元ディスクの I O 禁止フラグ及び予備ディスクの I O 禁止フラグが登録されており、これらに対して 1 ~ 12 の項(項番)が設定されている。

【0041】

項 1 の L U N 名には、L U N を特定する識別情報が登録され、図 3 に示す例においては“LUN1”が登録されている。項 2 の構成ディスク名リストは、その L U N を構成する記憶装置 3 1 を示す情報が登録される。すなわち、構成ディスク名リストには、R A I D 構成ディスク 3 1 a を示す情報が格納される。図 3 に示す例においては、構成ディスク名リストに disk1 , disk2 , disk3 , disk4 が登録されている。

【 0 0 4 2 】

項 3 の各ディスク上の位置、サイズは、その L U N を構成する R A I D 構成ディスク 3 1 a 上における、当該 L U N のデータの格納位置とサイズを示す情報が登録される。なお、位置としては、例えばオフセット (Offset) 値が登録される。

項 4 の故障ディスク名は、その L U N を構成する R A I D 構成ディスク 3 1 a において、故障の発生が検知された記憶装置 3 1 a を示す情報が登録され、図 3 に示す例においては“disk2”が登録されている。

【 0 0 4 3 】

項 5 の L U N の状態は、その L U N の状態を示す情報が格納され、図 3 に示す例においては、正常な状態であることを示す“正常”、リビルド中であることを示す“rebuild中”及び、リビルド済みであることを示す“rebuild済”のいずれかが格納される。

項 6 の利用予備ディスク名リストは、各 R A I D 構成ディスク 3 1 a に対応付けられている予備ディスク 3 1 b を示す情報が格納される。なお、この R A I D 構成ディスク 3 1 a に対応づけられる予備ディスク 3 1 b は、後述する割当処理部 1 4 によって設定される。

【 0 0 4 4 】

図 3 に示す例においては、項 2 の構成ディスク名リストの disk1 , disk2 , disk3 及び disk4 に対して、disk1' , disk2' , disk3' 及び disk4' が対応付けられている。すなわち、disk1 と disk1' とが対をなし、同様に、disk2 , disk3 及び disk4 が disk2' , disk3' 及び disk4' とそれぞれ対をなす。

また、R A I D 構成ディスク 3 1 a の数に対して割り当て可能な予備ディスク 3 1 b の数 ( 本数 ) が少ない場合には、後述の如く、一部の R A I D 構成ディスク 3 1 a に対してのみ予備ディスク 3 1 b が割り当てられる。図 3 においては、このように予備ディスク 3 1 b の数が少ない場合の例も表示されており、割り当てられた予備ディスク 3 1 b のみが示され、割り当てられる予備ディスク 3 1 b がいない部分には、横線 ( - ) を記している。

【 0 0 4 5 】

項 7 のディスクの安定度リストには、その L U N を構成する R A I D 構成ディスク 3 1 a を示す情報が、その安定度が高い順に登録される。図 3 に示す例においては、R A I D 構成ディスク 3 1 a が、disk1 , disk4 , disk3 の順に登録されている。なお、disk2 は故障中であるので、このディスクの安定度リストには登録されない。

項 8 のリビルド済サイズには、項 2 の故障ディスクに関するリビルド処理が実行された場合に、その復元されたデータの合計値 ( データサイズ ) が進捗状況として登録される。このリビルド済サイズの初期値は 0 であり、リビルド処理の進行に伴い、その値が大きくなる。

【 0 0 4 6 】

項 9 の元ディスクのリード / ライト中カウンタには、その L U N を構成する各 R A I D 構成ディスク ( 元ディスク ) 3 1 a のそれぞれについて、実行中のリードアクセス及びライトアクセスの数が格納される。すなわち、各 R A I D 構成ディスク 3 1 a への I O アクセス状態をリアルタイム表す。

図 3 に示す例においては、各 R A I D 構成ディスク 3 1 a のそれぞれに対して、リード中カウンタの値 n とライト中カウンタの値 n とを、( n , n ) の形式で示している。

【 0 0 4 7 】

ホスト装置 2 から L U N に対するリード要求もしくはライト要求の I O コマンドを受信すると、例えば、後述する I O 制御部 1 1 が、そのアクセス先の R A I D 構成ディスク 3 1 a についてのリード中カウンタもしくはライト中カウンタの値をインクリメントする。

又、ＩＯ制御部１１は、リードアクセスもしくはライトアクセスが完了すると、対応するカウンタの値をデクリメントする。

【００４８】

この元ディスクのリード／ライト中カウンタの値を参照することで、各ＲＡＩＤ構成ディスク３１ａがＩＯ処理を実行中であるか否かを判断することができる。

項１０の予備ディスクのリード／ライト中カウンタには、各予備ディスク３１ｂのそれぞれについて、実行中のリードアクセス及びライトアクセスの数が格納される。すなわち、各予備ディスク３１ｂへのＩＯアクセス状態を示す。

【００４９】

図３に示す例においては、各予備ディスク３１ｂのそれぞれに対して、リード中カウンタの値 $n$ とライト中カウンタの値 $n$ とを、 $(n, n)$ の形式で示している。

例えば、後述するミラーリング処理部１５が、その予備ディスク３１ｂに対してデータのライトを行なう際にライト中カウンタの値をインクリメントする。又、ミラーリング処理部１５は、ライト処理が完了すると、対応するカウンタの値をデクリメントする。

【００５０】

この予備ディスクのリード／ライト中カウンタの値を参照することで、各予備ディスク３１ｂがＩＯ処理を実行中であるか否かを判断することができる。

項１１の元ディスクのＩＯ禁止フラグは、そのＬＵＮを構成する各ＲＡＩＤ構成ディスク３１ａのそれぞれについて、ＩＯ処理が禁止されているか否かを示す情報が格納される。図３に示す例においては、“０”もしくは“１”が格納され、“１”がフラグとして設定されている場合には、そのＲＡＩＤ構成ディスク３１ａへのＩＯが禁止されていることを示す。

【００５１】

項１２の予備ディスクのＩＯ禁止フラグは、予備ディスク３１ｂのそれぞれについて、ＩＯ処理が禁止されているか否かを示す情報が格納される。図３に示す例においては、“０”もしくは“１”が格納され、“１”がフラグとして設定されている場合には、その予備ディスク３１ｂへのＩＯが禁止されていることを示す。

リビルド処理部１３は、例えば、いずれかのディスク３１ａの故障を検出した場合に、リビルド処理を実行制御する。以下、故障が検出されたディスク３１ａを故障ディスク３１ａという場合がある。この故障ディスク３１ａが復元対象記憶装置に相当する。なお、ディスク３１ａの故障は、例えば媒体エラー等の所定のエラーが閾値として設定された頻度以上で生じる場合に発生したと判断される。

【００５２】

リビルド処理は、ＲＡＩＤグループの冗長性を自動回復する処理であり、ＲＡＩＤグループに属する記憶装置３１ａが故障した際に、代理として用いられる予備ディスク（代理ディスク、第１の予備記憶装置）３１ｂに対し、故障ディスク３１ａのデータを、同一ＲＡＩＤグループにおける故障ディスク３１ａ以外の記憶装置３１ａのデータを用いて再構築する処理である。以下、同一ＲＡＩＤグループにおける故障ディスク３１ａ以外の記憶装置（冗長用記憶装置）３１ａをリビルド元ディスク３１ａという場合がある。又、このリビルド元ディスクを単に元ディスクという場合もある。

【００５３】

つまり、リビルド処理部１３は、故障ディスク３１ａの発生を検出すると、故障ディスク３１ａ以外のリビルド元ディスク３１ａのデータを用いて、故障ディスク３１ａに代わる代理ディスク（リビルド先ディスク）３１ｂに故障ディスク３１ａのデータを再構築する。

このように、リビルド処理部１３は、複数のＲＡＩＤ構成ディスク３１ａのうち故障ディスク（復元対象記憶装置）３１ａ以外のリビルド元ディスク（冗長用記憶装置）３１ａから読み出した冗長データを用いて、故障ディスク３１ａのデータを、予備ディスク（第１の予備記憶装置）に再構成する再構成処理部として機能する。

【００５４】

リビルド処理部（データ復元部）１３による故障ディスク３１ａのデータの復元は、既知の手法で実現することができる。

同一のＲＡＩＤグループ３０ａを構成する複数の記憶装置３１ａは、各記憶装置３１ａのデータが他の複数の記憶装置（冗長用記憶装置）３１ａに分散してコピーされることにより、冗長化されている。

【００５５】

リビルド処理部１３は、故障ディスク３１ａと同じＲＡＩＤグループを構成する複数の記憶装置（冗長用記憶装置）３１ａにそれぞれ格納されたデータを読み出して、代替ディスク３１ｂに格納することで、故障ディスク３１ａのデータを代替ディスク３１ｂに復元（データディスク構築）する。

10

例えば、メモリ１０６の所定の領域に、格納した故障ディスク３１ａ以外の各リビルド元ディスク３１ａから読み出したデータを格納し、このデータに対してパリティを用いたＸＯＲ（排他的論理和）演算を行なうことで、故障ディスク３１ａのデータの復元を行なう。

【００５６】

図４は実施形態の一例としてのストレージ装置１におけるリビルド処理を説明する図である。この図４に示す例においては、ＲＡＩＤグループ３０ａは、４本（３本＋１本）のＲＡＩＤ構成ディスク３１ａ－１～３１ａ－４によりＲＡＩＤ５を実現している。又、これらのＲＡＩＤ構成ディスク３１ａ－１～３１ａ－４のうち、ディスク３１ａ－２の故障が検知された場合を示す。

20

【００５７】

また、この図４に示す例においては、ＲＡＩＤグループ３０ａのディスク３１ａ－１～３１ａ－４を用いて３つのＬＵＮ１～３が形成されており、これらのうちＬＵＮ１のデータ（Data1-2）を復元する例を示している。

リビルド処理部１３は、ディスク３１ａ－２の故障を検出すると、同一ＲＡＩＤグループ３０ａ内のリビルド元ディスク３１ａ－１、３１ａ－３、３１ａ－４のデータ（Data1-1, 1-3, 1-4）を用いて、故障ディスク３１ａ－２のデータ（復元Data1-2）を作成する（図４中の破線参照）。

【００５８】

そして、この作成した故障ディスク３１ａ－２のデータ（復元Data1-2）を、故障ディスク３１ａ－２に代わる代替ディスク（リビルド先ディスク）３１ｂ－２に格納することで、故障ディスク３１ａ－２を再構築する。

30

ＲＡＩＤグループ３０ａにおいて、複数のＬＵＮが形成されている場合には、リビルド処理部１３はＬＵＮ毎にリビルド処理を行なう。又、リビルド処理部１３は、複数のＬＵＮについてリビルドを行なう場合には、これらの複数のＬＵＮの一覧（図示省略）と、当該一覧においてリビルド処理の進捗を示すポイントとを用いて、リビルド処理の進捗状況を管理する。

【００５９】

ミラーリング処理部（複製処理部）１５は、リビルド処理部１３によるリビルド処理中に、ＲＡＩＤグループ３０ａ内のリビルド元ディスク（冗長用記憶装置）３１ａからそれぞれ読み出したデータを、各リビルド元ディスク３１ａに割り当てられた予備ディスク（第２の予備記憶装置）３１ｂに格納する。これにより、ミラーリング処理部１５は、各リビルド元ディスク３１ａを予備ディスク３１ｂに複製する。

40

【００６０】

図４に示す例においては、リビルド元ディスク３１ａ－１に対して予備ディスク３１ｂ－１が対応付けられており、同様に、リビルド元ディスク３１ａ－３に対して予備ディスク３１ｂ－３が、リビルド元ディスク３１ａ－４に対して予備ディスク３１ｂ－４が、それぞれ対応付けられている。なお、このようなリビルド元ディスク３１ａへの予備ディスク３１ｂの対応付けは、後述する割当処理部１４が行なう。

【００６１】

50

ミラーリング処理部 15 は、リビルド処理部 13 によるリビルド処理中に、RAID グループ 30 a 内のリビルド元ディスク (冗長用記憶装置) 31 a から読み出されメモリ 106 に格納されたデータを、対応する予備ディスク 31 b に格納 (デッドコピー) することで、予備ディスク 31 b へのリビルド元ディスク 31 a の複製を行なう。

すなわち、図 4 に示す例においては、リビルド元ディスク 31 a - 1, 31 a - 3, 31 a - 4 の LUN 1 の各データ (Data1-1, 1-3, 1-4) が、予備ディスク 31 b - 1, 31 b - 3, 31 b - 4 にそれぞれコピーされる (矢印 P1 ~ P3 参照)。

#### 【0062】

割当処理部 14 は、リビルド元ディスク 31 a に対して予備ディスク 31 b を対応付ける。以下、リビルド元ディスク 31 a に対して予備ディスク 31 b を対応付けることを、  
「割り当てる」と表現する場合がある。

また、予備ディスク 31 b のうち、RAID 構成ディスク 31 a に対応付けられていないディスク 31 b を特に未割当予備ディスク 31 b という。

#### 【0063】

図 4 に示した例においては、リビルド元ディスク 31 a のそれぞれに予備ディスク 31 b が対応付けられている。

しかしながら、使用可能な予備ディスク 31 b の数がリビルド元ディスク 31 a の数に足りず、予備ディスク 31 b を全てのリビルド元ディスク 31 a に対応付けることができない場合もある。このような場合に、割当処理部 14 は、一部のリビルド元ディスク 31 a に対してだけ予備ディスク 31 b を割り当てる。

#### 【0064】

割当処理部 14 は、RAID サイズ (P) とスペア d i s k 数 (m) と制限値 (L) とを参照することで、リビルド時に使用する予備ディスク 31 b の数を決定する。

ここで、RAID サイズ (P) は、RAID グループ 30 a において実現される RAID を構成する記憶装置 31 a の数 (本数) であり、RAID 種類により決定される。例えば、RAID 5 の場合は 4 本 (3 本 + 1 本) である (P=4)。この RAID サイズは、RAID 構成ディスク 31 a の数でもある。

#### 【0065】

スペア d i s k 数 (m) は、使用可能な予備ディスク 31 b の数であり、例えば、ディスク構成情報 63 を参照することで確認することができる。制限値 (L) は、例えば予備ディスク 31 b の数 (m) が RAID サイズ (P) 未満の場合に、リビルド時に使用する最低限の予備ディスク 31 b の数である。例えば、制限値 L = 2 の場合には、リビルド時に 2 本の予備ディスク 31 b が用いられることを示す。制限値 (L) は、管理者等によって予め設定される。

#### 【0066】

また、割当処理部 14 は、予備ディスク数 (m) が RAID サイズ (P) に満たない場合に、リビルド元ディスク 31 a のうち、安定度が低い、すなわち、不安定なリビルド元ディスク 31 a に対して予備ディスク 31 b を優先して割り当てる。割当処理部 14 は、例えば、LUN 管理テーブル 61 の項 7 のディスクの安定度リストを参照することで、安定度の低い、すなわち、最も不安定なリビルド元ディスク 31 a を知ることができる。

#### 【0067】

図 5 は実施形態の一例としてのストレージ装置 1 におけるリビルド処理を説明する図であり、予備ディスク 31 b が一部のリビルド元ディスク 31 a にだけ対応付けられた例を示す。

この図 5 に示す例においては、RAID 構成ディスク 31 a - 2 が故障し、この RAID 構成ディスク 31 a - 2 のデータがリビルド処理により代理ディスク 31 b - 2 に復元されている (図 5 中の破線参照)。

#### 【0068】

そして、リビルド元ディスク 31 a - 1, 31 a - 3, 31 a - 4 のうち、リビルド元ディスク 31 a - 3 に対してだけ予備ディスク 31 b - 3 が対応付けられており、他のリ

10

20

30

40

50

ビルド元ディスク 3 1 a - 1 , 3 1 a - 4 に対しては予備ディスク 3 1 b が対応付けられていない。

全てのリビルド元ディスク 3 1 a に対して予備ディスク 3 1 b を割り当てることができない場合には、割当処理部 1 4 は安定度の低い一部のリビルド元ディスク 3 1 a に対してだけ予備ディスク 3 1 b が割り当てて。

【 0 0 6 9 】

このように、安定度の低い一部のリビルド元ディスク 3 1 a に対してだけ予備ディスク 3 1 b が割り当てられた場合には、ミラーリング処理部 1 5 は、この予備ディスク 3 1 b が割り当てられているリビルド元ディスク 3 1 a のデータを、当該リビルド元ディスク 3 1 a に対して割り当てられた予備ディスク 3 1 b にコピーする。

10

図 5 に示す例においては、リビルド元ディスク 3 1 a - 3 の L U N 1 のデータ (Data1-3) が、予備ディスク 3 1 b - 3 にコピーされる (矢印 P 4 参照)。

【 0 0 7 0 】

すなわち、ミラーリング処理部 1 5 は、リビルド処理部 1 3 によるリビルド処理中に、リビルド元ディスク 3 1 a - 3 から読み出されメモリ 1 0 6 に格納されたデータを、対応する予備ディスク 3 1 b - 3 に格納 (デッドコピー) する。これにより、予備ディスク 3 1 b - 3 にリビルド元ディスク 3 1 のデータの複製が行なわれる。

余分に割り当てた予備ディスク 3 1 b は、システム内で別の R A I D に故障が発生して予備ディスク 3 1 b が不足となった場合に、リビルド処理が完了していなくても、対をなす R A I D 構成ディスク 3 1 a が安定している予備ディスク 3 1 b から順番に、割り当てを解除することで未割当予備ディスク 3 1 b として利用する。

20

【 0 0 7 1 】

I O 制御部 1 1 は L U N に対する I O 制御を行なう。この I O 制御部 1 1 は、例えば、L U N に対して F C P (Fibre Channel Protocol) もしくは N A S の I O 制御を行なう。

I O 制御部 1 1 は、上述した仮 / 実ボリューム変換テーブル 6 2 を用いて、ホスト装置 2 からの I O 要求に応じて実ボリュームである記憶装置 3 1 に対する I O 処理を行なう。I O 制御部 1 1 は、ホスト装置 2 から送信された I O 要求に応じて、L U N を構成する記憶装置 3 1 に対して、データのリードやライトを行なう。

【 0 0 7 2 】

図 6 は実施形態の一例としてのストレージ装置 1 におけるリビルド処理を説明する図である。この図 6 に示す例においては、図 4 と同様に、R A I D グループ 3 0 a は 4 本の R A I D 構成ディスク 3 1 a - 1 ~ 3 1 a - 4 により R A I D 5 を実現しており、又、代理ディスク 3 1 b - 2 に故障ディスク 3 1 a - 2 のデータを再構築している。特にこの図 6 においては、図 4 に示した状態の後、L U N 1 のリビルド処理が完了し、L U N 2 のリビルド処理中の状態を示す。

30

【 0 0 7 3 】

リビルド処理部 1 3 は、リビルド元ディスク 3 1 a - 1 , 3 1 a - 3 , 3 1 a - 4 の L U N 2 のデータ (Data2-1 , 2-3 , 2-4) を用いて、故障ディスク 3 1 a - 2 の L U N 2 のデータ (復元Data2-2) を作成する。そして、この作成した故障ディスク 3 1 a - 2 の L U N 2 のデータ (復元Data2-2) を、故障ディスク 3 1 a - 2 に代わる代替ディスク 3 1 b - 2 に格納している (図 6 中の破線参照)。

40

【 0 0 7 4 】

また、これに伴い、ミラーリング処理部 1 5 が、リビルド元ディスク (冗長用記憶装置) 3 1 a - 1 , 3 1 a - 3 , 3 1 a - 4 の L U N 2 の各データ (Data2-1 , 2-3 , 2-4) を、予備ディスク 3 1 b - 1 , 3 1 b - 3 , 3 1 b - 4 にそれぞれコピーしている。

この図 6 に示す状態において、リビルド済みの L U N 1 に対してホスト装置 2 から I O 要求が行なわれると、I O 制御部 1 1 は、R A I D グループ 3 0 a の R A I D 構成ディスク 3 1 a と、予備ディスクグループ 3 0 b の予備ディスク 3 1 b との両方のディスクデータ域を用いて I O 処理を行なう。

【 0 0 7 5 】

50

具体的には、LUN 1 のデータに対するリードアクセスが行なわれると、RAID グループ 30 a の RAID 構成ディスク 31 a と、予備ディスクグループ 30 b の予備ディスク 31 b とを併用してリード処理を行なう。例えば、IO 制御部 11 は、RAID グループ 30 a の RAID 構成ディスク 31 a と、予備ディスクグループ 30 b の予備ディスク 31 b とをラウンドロビンで交互に選択してリード処理を行なう。

【0076】

リビルド済みの LUN 1 のデータは RAID グループ 30 a と予備ディスクグループ 30 b との両方に存在し、二重化されている。

リード処理については、これらのデータは RAID グループ 30 a の RAID 構成ディスク 31 a と、予備ディスクグループ 30 b の予備ディスク 31 b とのいずれか一方から行なえばよい。本ストレージ装置 1 においては、ホスト装置 2 からのリード要求に対して、これらの RAID グループ 30 a と予備ディスクグループ 30 b との両方を用いてデータリードを行なう。すなわち、複数のデータリードを並列して実行できるようにすることでデータリードのパフォーマンスを向上させるとともに、ディスクアクセス負荷を分散させることで、記憶装置 31 の寿命を延ばすことができる。

【0077】

具体的には、RAID グループ 30 a の RAID 構成ディスク 31 a と、予備ディスクグループ 30 b の予備ディスク 31 b とをラウンドロビンで交互に選択してリード処理を行なうことで、リビルド中であっても高いストレージアクセス性能を実現することができる。

一方、LUN 1 のデータに対するライトアクセスが行なわれると、RAID グループ 30 a の RAID 構成ディスク 31 a と、予備ディスクグループ 30 b の予備ディスク 31 b との両方に対してライト処理を行なう。

【0078】

このように、ライト処理については、RAID グループ 30 a の RAID 構成ディスク 31 a と、予備ディスクグループ 30 b の予備ディスク 31 b とで重複して行なう必要がある。しかしながら、一般に、ストレージ装置においては、ライト処理に比べてリード処理の割合が多く、リード：ライト = 2：1 程度であると言われているので、全体的にはライト処理を重複して行なうことによるストレージアクセス性能への影響は少ないと考えられる。

【0079】

なお、図 6 に示す状態において、リビルドが完了していない LUN 2, 3 に対してホスト装置 2 から IO コマンドが発行されると、IO 制御部 11 は、RAID グループ 30 a の RAID 構成ディスク 31 a のディスクデータ域を用いて IO 処理を行なう。

例えば、LUN 2 のデータに対するリードアクセスが行なわれると、IO 制御部 11 は、RAID グループ 30 a の RAID 構成ディスク 31 a からリード処理を行なう。又、故障ディスク 31 a - 2 に対するリード処理は、この故障ディスク 31 a - 2 以外の各リビルド元ディスク（冗長用記憶装置）31 a - 1, 31 a - 3, 31 a - 4 のデータに対して、パリティを用いた XOR 演算を行なって故障ディスク 31 a - 2 のデータの復元を行なって対応する。

【0080】

また、LUN 2 のデータに対するライトアクセスが行なわれると、IO 制御部 11 は、RAID グループ 30 a の RAID 構成ディスク 31 a に対してライト処理を行なう。又、故障ディスク 31 a - 2 に対するライト処理は、この故障ディスク 31 a - 2 以外の各リビルド元ディスク（冗長用記憶装置）31 a - 1, 31 a - 3, 31 a - 4 のデータに対して、パリティを用いた XOR 演算を行なって故障ディスク 31 a - 2 のデータの復元を行ない、この復元したデータを用いて、予備ディスクグループ 30 b の代理ディスク 31 b - 2 に対してライト処理を行なう。

【0081】

RAID 構成変更部 16 は、リビルド処理部 13 によるリビルド処理の完了後に、RA

10

20

30

40

50

ＩＤグループ３０ａにおいて不安定なリビルド元ディスク３１ａがある場合には、その不安定なリビルド元ディスク３１ａに代えて予備ディスク３１ｂを用いてＲＡＩＤ構成を組み直す。

すなわち、ＲＡＩＤ構成変更部１６は、リビルド処理部１３によるリビルド処理の完了後に、安定度の低いリビルド元ディスク３１ａに代えて、当該リビルド元ディスク３１ａの複製がされた予備ディスク３１ｂを用いて、ＲＡＩＤ構成を変更する冗長構成変更部として機能する。

【００８２】

リビルド元ディスク３１ａが不安定であるか否かの判断は、例えば、ストレージ統計情報（ログ情報）に基づいて行なう。ストレージ統計情報において所定のエラーが閾値以上検知された場合に、不安定であると判断することができる。

10

ストレージ統計情報としては、例えば、各ＲＡＩＤ構成ディスク３１ａにおける媒体エラーの発生数やシークエラーの発生数を用いることができる。このようなストレージ統計情報は、例えば、各ディスク３１ａのファーム等を参照することにより取得することができる。

【００８３】

ＲＡＩＤ構成変更部１６は、リビルド処理部１３によりすべてのＬＵＮのリビルド処理が完了したら、障害が検知された障害ＲＡＩＤを構成するＲＡＩＤ構成ディスク３１ａと、それに対をなす予備ディスク３１ｂの信頼性を確認して、より安定なＲＡＩＤ構成となるように構成を組みなおす。

20

なお、元のＲＡＩＤを構成するディスク群の各搭載位置を維持する必要がある場合には、元のＲＡＩＤのＲＡＩＤ構成ディスク３１ａを問題がなければ使用し続けてもよい。

【００８４】

図７は実施形態の一例としてのストレージ装置１におけるＲＡＩＤ構成変更部１６によるＲＡＩＤ構成変更方法を示す図である。

この図７に示す例においても、図４と同様に、ＲＡＩＤグループ３０ａは４本のＲＡＩＤ構成ディスク３１ａ－１～３１ａ－４によりＲＡＩＤ５を実現しており、又、故障ディスク３１ａ－２が代理ディスク３１ｂ－２に再構築されている。特にこの図７においては、図６に示した状態の後、ＬＵＮ２，３のリビルド処理が完了した状態を示す。

【００８５】

30

すなわち、予備ディスク３１ｂ－１，３１ｂ－３，３１ｂ－４には、ミラーリング処理部１５によりＲＡＩＤ構成ディスク３１ａ－１，３１ａ－３，３１ａ－４のデータがそれぞれ複製されている。

そして、この図７に示す例においては、ＲＡＩＤ構成ディスク３１ａ－３が不安定であると判断された状態を示す。

【００８６】

リビルド処理の完了後に、ＲＡＩＤ構成ディスク３１ａ－３が不安定であると判断された場合に、ＲＡＩＤ構成変更部１６は、このＲＡＩＤ構成ディスク３１ａ－３に代えて、ＲＡＩＤ構成ディスク３１ａ－３のデータが複製された予備ディスク３１ｂ－３を用いてＲＡＩＤを構成し直す。

40

なお、ＲＡＩＤ構成ディスク３１ａ－３から予備ディスク３１ｂ－３への切り替えは既知の種々の手法を用いて実現することができる。例えば、ＲＡＩＤ構成変更部１６は、不安定であると判断されたＲＡＩＤ構成ディスク３１ａ－３をfail状態に設定するコマンドを発行することで、ＲＡＩＤ構成ディスク３１ａ－３を予備ディスク３１ｂ－３に切り替える。

【００８７】

また、不安定であると判断されたＲＡＩＤ構成ディスク３１ａが複数ある場合には、これらの複数のＲＡＩＤ構成ディスク３１ａについても同様に予備ディスク３１ｂに切り替える。

また、ＲＡＩＤ構成変更部１６は、リビルド処理の完了後、不安定であると判断されて

50

いないRAID構成ディスク31aに対応する予備ディスク31bを未割当予備ディスク31bに戻す。これにより、ドライブエンクロージャ30のスロットに予備ディスク31bとして取り付けられた記憶装置31を、積極的に予備ディスク31bとして使用することができ、管理が容易になる。図7に示す例においては、予備ディスク31b-1, 31b-4がフリーの未割当予備ディスク31bに戻され、RAID構成ディスク31aとの対応付けが解除される。

#### 【0088】

上述の如く構成された実施形態の一例としてのストレージ装置1におけるリビルド処理の概要を、図8に示すフローチャート(ステップA1~A3)に従って説明する。

ステップA1において、割当処理部14は使用可能な予備ディスク31bの本数を確認して、ストレージ構成を決定する。なお、このストレージ構成の決定手法の詳細は、図9を用いて後述する。

#### 【0089】

ステップA2において、リビルド処理部13が故障ディスク31のリビルド処理を行なう。まず、ステップA21において、RAIDグループ30aに構成されている全てのLUNについてのリビルドが完了したかを確認する。全LUNのリビルドが完了していない場合には(ステップA21のNルート参照)、ステップA22において、リビルド処理部13が、故障ディスク31のデータを復元する。又、このリビルド処理中に、ミラーリング処理部15が、リビルド元ディスク31aから読み出されたデータを対応する予備ディスク31bにデッドコピーする。これにより、予備ディスク31bへのリビルド元ディスク31aの複製が行なわれる。

#### 【0090】

ステップA23において、処理中のLUNのリビルドが完了すると、対になるRAID構成ディスク31aと予備ディスク31bとを二重化して利用可能にする。例えば、ホスト装置2からリード要求を受信した場合には、対を形成するRAID構成ディスク31aと予備ディスク31bとからラウンドロビンで交互にリードされるようにする。又、ホスト装置2からライト処理を受信した場合には、対を形成するRAID構成ディスク31aと予備ディスク31bとの両方にライト処理が行なわれるようにする。その後、ステップA21に戻る。

#### 【0091】

なお、このリビルド処理の詳細は、図10及び図11を用いて後述する。

ステップA21において全てのLUNについての処理が完了すると(ステップA21のYルート参照)、リビルド処理後の処理、すなわち後処理(ステップA3)に移行する。

この後処理においては、ステップA31において、RAID構成変更部16が、故障が検知されたRAID構成ディスク31aに代えて予備ディスク31bを用いてRAID構成を組み直す。

#### 【0092】

また、RAIDグループ30aにおいて不安定なRAID構成ディスク31aがある場合には、RAID構成変更部16は、その不安定なRAID構成ディスク31aに代えて、当該RAID構成ディスク31aに対応する予備ディスク31bを用いてRAID構成を組み直す。

その後、ステップA32において、ステップA31で使用されなかった予備ディスク31bを解放して、未割当予備ディスク31bに戻し、処理を終了する。

#### 【0093】

なお、故障ディスク31aや、ステップA32において不安定であると判断されたRAID構成ディスク31aは、保守作業により当該ストレージ装置1から取り外され、新しい記憶装置31と交換される。

なお、この後処理の詳細は、図12を用いて後述する。

次に、実施形態の一例としてのストレージ装置1におけるストレージ構成の決定方法を

10

20

30

40

50

、図 9 に示すフローチャート（ステップ B 1 ～ B 8 ）に従って説明する。

【 0 0 9 4 】

ステップ B 1 において、割当処理部 1 4 が、R A I D サイズ（ P ）, 予備ディスク数（ m ）及び制限値（ L ）を確認する。

ステップ B 2 において、割当処理部 1 4 は、“ R A I D サイズ P = 2 ”, “ 予備ディスク数 m = 1 ” 及び “ 制限値 L = 1 ” の 3 つの条件のうち、少なくともいずれか 1 つが満たされるか否かを確認する。

【 0 0 9 5 】

確認の結果、これらの 3 つの条件のうち、少なくともいずれか 1 つが満たされる場合には（ステップ B 2 の Y E S ルート参照）、リビルド処理部 1 3 は、従来手法によるリビルド処理を行なう。すなわち、故障ディスク 3 1 a のデータを、同一 R A I D グループにおける故障ディスク 3 1 a 以外のリビルド元ディスク 3 1 a のデータを用いて再構築する。又、この際、割当処理部 1 4 によるリビルド元ディスク 3 1 a に対する予備ディスク 3 1 b を割り当てや、ミラーリング処理部 1 5 によるリビルド元ディスク 3 1 a のデータの対応する代理ディスク 3 1 b への複製は行なわれない。

【 0 0 9 6 】

ここで R A I D サイズ P は、R A I D により必要とされる記憶装置 3 1 の数を示す。従って、“ R A I D サイズ P = 2 ” は、R A I D グループ 3 0 a の R A I D 種類が R A I D 1 の二重化（ミラーリング）であり、リビルド処理が不要であることを意味する。又、“ 予備ディスク数 m = 1 ” 及び “ 制限値 L = 1 ” は、いずれも予備ディスク 3 1 b が 1 つもしくは 0 であり、割当処理部 1 4 により R A I D 構成ディスク 3 1 a に割り当て可能な予備ディスク 3 1 b がないことを意味する。

【 0 0 9 7 】

一方、確認の結果、上述した 3 つの条件のいずれも満たされない場合には（ステップ B 2 の N O ルート参照）、ステップ B 3 において、“ R A I D サイズ P > 制限値 L ” であるか否かが確認される。

R A I D サイズ P > 制限値 L である場合（ステップ B 3 の Y E S ルート参照）、ステップ B 4 において、L（例えば、図 5 に示したように L = 2）本の予備ディスク 3 1 b がリビルド処理及びミラーリング処理に使用される。

【 0 0 9 8 】

一方、R A I D サイズ P > 制限値 L でない場合（ステップ B 3 の N O ルート参照）、ステップ B 5 において、P（例えば、図 4 等 に示したように L = 4）本の予備ディスク 3 1 b がリビルド処理及びミラーリング処理に使用される。

その後、ステップ B 6 において、R A I D 制御部 1 2 は、記憶装置 3 1 a に故障が生じた R A I D グループ 3 0 a（故障 R A I D）において生き残っている R A I D 構成ディスク 3 1 a のストレージ統計情報を確認する。そして R A I D 制御部 1 2 は、L U N 管理テーブル 6 1 の項 7 のディスクの安定度リストとして格納する情報を作成し、L U N 管理テーブル 6 1 に登録する。ストレージ統計情報として取得した各 R A I D 構成ディスク 3 1 a のエラーの発生数に基づき、エラーの発生数に従って R A I D 構成ディスク 3 1 a をソートすることで、ディスクの安定度リストを作成する。

【 0 0 9 9 】

また、ステップ B 7 において、割当処理部 1 4 は、故障が検知された R A I D 構成ディスク 3 1 a を含む不安定な R A I D 構成ディスク 3 1 a を優先して、R A I D 構成ディスク 3 1 a に予備ディスク 3 1 b を割り当てる。R A I D 制御部 1 2 は、L U N 管理テーブル 6 1 の項 6 の利用予備ディスク名リストとして格納する情報を作成し、L U N 管理テーブル 6 1 に登録する。

【 0 1 0 0 】

ステップ B 8 において、R A I D 制御部 1 2 は、L U N 管理テーブル 6 1 の項 4 ～ 1 2 の各内容を初期化し、又、L U N 管理テーブル 6 1 の項 5 の L U N の状態に “ リビルド中 ” を登録して、処理を終了する。

次に、実施形態の一例としてのストレージ装置 1 におけるリビルド処理を、図 10 に示すフローチャート（ステップ C 1 ~ C 10）に従って説明する。

【0101】

ステップ C 1 において、リビルド処理部 13 は、LUN 管理テーブル 61 を参照することで LUN 数を取得する。LUN 管理テーブル 61 は LUN 毎に作成されるので、例えば LUN 管理テーブル 61 の数を参照することで、RAID グループ 30a に形成されている LUN の数を把握することができる。又、リビルド処理部 13 は、処理中の LUN を示す図示しない処理中ポインタに、最初の LUN（LUN 域）を示す情報を記憶する。この処理中ポインタを参照することで、リビルド処理の進捗状況を把握することができる。

【0102】

ステップ C 2 において、リビルド処理部 13 は、処理中ポインタを参照して、未処理の LUN 数を確認する。

未処理の LUN 数が 1 つ以上ある場合には（ステップ C 2 の NO ルート参照）、ステップ C 3 において、処理中 LUN の情報に基づき、生き残っている各リビルド元ディスク 31a を調べ、この生き残っているリビルド元ディスク 31a から、RAID ストライプ分のデータを読み出し、メモリ 106 の所定の領域に格納する。

【0103】

ステップ C 4 において、リビルド処理部 13 は、各リビルド元ディスク 31a から読み出したデータを用いて故障ディスク 31a の復元を行なう（復元データ作成）。

ステップ C 5 において、リビルド処理部 13 は、故障ディスク 31 に対してアサインされた予備ディスク 31b の所定位置（元と同じ位置）へ復元データを書き出す。又、同時に、ミラーリング処理部 15 は、生き残っているリビルド元ディスク 31 から読み出したデータを、各リビルド元ディスク 31 に対応付けた（アサインした）予備ディスク 31b の同じ位置へヘッドコピーする。

【0104】

ステップ C 6 において、リビルド処理部 13 は、LUN 管理テーブル 61 の項 8 のリビルド済サイズに、リビルドが完了したデータサイズ（リビルド済サイズ）を加算（up）する。

ステップ C 7 において、リビルド処理部 13 は、他の RAID の故障による予備ディスクの解放要求があるか否かを確認する。

【0105】

ここで、このステップ C 7 にかかる、予備ディスクの解放要求の有無の確認処理の詳細を、図 11 に示すフローチャート（ステップ C 71 ~ C 77）に従って説明する。

ステップ C 71 において、他 RAID の故障による予備ディスク 31 の解放要求があるか否かを確認する。他の RAID から予備ディスク 31b の解放要求がない場合には（ステップ C 71 の NO ルート参照）、処理を終了し、図 10 のステップ C 8 に移行する。

【0106】

他の RAID から予備ディスク 31b の解放要求がある場合には（ステップ C 71 の YES ルート参照）、ステップ C 72 において、リビルド処理部 13 は、LUN 管理テーブル 61 の項 7 のディスクの安定度リストを参照して、解放できる予備ディスク 31b があるか否かを確認する。対応する RAID 構成ディスク 31a が安定していれば予備ディスク 31b が代理ディスク 31b として使用される可能性は低い。従って、LUN 管理テーブル 61 の項 7 のディスクの安定度リストにおいて、安定度が高い RAID 構成ディスク 31a に対応付けられている予備ディスク 31b は、解放して他の記憶装置 31a の予備ディスク 31b として使用しても問題ないと考えられる。

【0107】

そこで、例えば、LUN 管理テーブル 61 の項 7 のディスクの安定度リストに、一つでも RAID 構成ディスク 31a が登録されていれば、解放できる予備ディスク 31b があると判断することができる。

確認の結果（ステップ C 73）、解放できる予備ディスク 31b がある場合には（ステ

10

20

30

40

50

ップC73のYESルート参照)、ステップC74において、リビルド処理部13は、LUN管理テーブル61の項12の予備ディスクのIO禁止フラグに、当該解放対象の予備ディスク31bのIOを禁止にするフラグを設定する。このIO禁止フラグが設定される予備ディスク31bは、LUN管理テーブル61の項7のディスクの安定度リストにおいて先頭に登録されているRAID構成ディスク31aに対応する予備ディスク31bであることが望ましい。

【0108】

これにより、当該予備ディスク31bに対するIO処理が禁止され、最終的にIOアクセスが無くなるので、当該予備ディスク31bを使用することができるようになる。

その後、ステップC75において、リビルド処理部13は、LUN管理テーブル61の項10の予備ディスクのリード/ライト中カウンタにおいて、解放対象の予備ディスク31bのカウント値を確認する。

【0109】

確認の結果(ステップC76)、解放予定の予備ディスク31bがリード中でなく、且つ、ライト中(RW中)でもない場合、すなわち、解放対象の予備ディスク31bに関してリード中カウンタもしくはライト中カウンタのいずれにおいても1以上の値が格納されていない場合には(ステップC76のNOルート参照)、ステップC77に移行する。

このステップC77において、リビルド処理部13は、LUN管理テーブル61の項6の利用予備ディスク名リストから、当該予備ディスク31bを削除してLUN管理テーブル61を更新する。これにより、リビルド処理部13は、その予備ディスク31bを割り当て可能な予備ディスク31bとしてシステムに返却して、処理を終了し、図10のステップC8に移行する。

【0110】

また、ステップC73における確認の結果、解放できる予備ディスク31bがない場合(ステップC73のNOルート参照)や、ステップC76における確認の結果、解放予定の予備ディスク31bがリード中もしくはライト中(RW中)である場合にも(ステップC76のYESルート参照)、処理を終了し、図10のステップC8に移行する。

ステップC8において、リビルド処理部13は、処理中のLUNのリビルド処理が完了したか否かを確認する。この確認の結果(ステップC9)、リビルド処理が完了していない場合には(ステップC9のNOルート参照)、ステップC3に戻る。

【0111】

一方、リビルド処理が完了した場合には(ステップC9のYESルート参照)、ステップC10に移行する。

ステップC10において、リビルド処理部13は、LUN管理テーブル61の項5のLUNの状態に“リビルド済”を設定する。又、リビルド処理部13は、未処理のLUN数をダウンして、次のLUN管理テーブル100を処理中ポインタに記憶する。その後、ステップC2に戻る。

【0112】

ステップC2において、未処理のLUN数が0である場合には(ステップC2のYESルート参照)、処理を終了する。

次に、実施形態の一例としてのストレージ装置1におけるリビルド処理後の処理(後処理)を、図12に示すフローチャート(ステップD1~D9)に従って説明する。

ステップD1において、RAID構成変更部16は、上述したリビルド処理を行なった故障ディスク31a以外のRAID構成ディスク31aについて、その統計情報を確認することで、これらのRAID構成ディスク31aに異常がないかを確認する。

【0113】

ステップD2において、RAID構成変更部16は、RAID構成ディスク31aに問題の有無を確認する。この問題の有無の確認は、例えば、媒体エラーやシークエラー等の所定のエラーの発生数と閾値とを比較することで行なう。所定のエラーの発生数が閾値を越えている場合に問題があると判断することができる。

10

20

30

40

50

確認の結果、問題があると判断された場合には（ステップD2のNOルート参照）、ステップD3において、その問題があると判断されたRAID構成ディスク31aを、failさせるリビルド元ディスク（元ディスク）としてメモリ106等に記憶して、ステップD4に移行する。

【0114】

また、ステップD2における確認の結果、問題がないと判断された場合には（ステップD2のYESルート参照）、ステップD3をスキップしてステップD4に移行する。

ステップD4において、RAID構成変更部16は、failさせるディスク31aと解放する予備ディスク31bとにIO禁止フラグを設定する。具体的には、LUN管理テーブル61の項11において、該当する記憶装置31aにIO禁止フラグを設定する。又、RAID構成変更部16は、LUN管理テーブル61の項12において、該当する予備ディスク31bにIO禁止フラグを設定する。

10

【0115】

ステップD5において、RAID構成変更部16は、LUN管理テーブル61の項9及び項10を参照して、failさせる元ディスク及び解放する予備ディスクの各リード/ライト中カウンタの値を確認する。すなわち、これらのfailさせる元ディスク及び解放する予備ディスクが使用されておらずIOアクセスがないことを確認する。

確認の結果（ステップD6）、failさせる元ディスク及び解放する予備ディスクの各リード/ライト中カウンタの値が0でない場合には（ステップD6のNOルート参照）、ステップD7において所定の時間（例えば1秒）待った後、ステップD5に戻る。

20

【0116】

ステップD4において、failさせる元ディスク及び解放する予備ディスクに対してIO禁止フラグを設定することにより、これらのディスク31への新規のディスクアクセスがなくなり、最終的には各リード/ライト中カウンタの値が0となる。

failさせる元ディスク及び解放する予備ディスクの各リード/ライト中カウンタの値が0である場合には（ステップD6のYESルート参照）、ステップD8に移行する。

【0117】

ステップD8において、RAID構成変更部16は、RAID構成テーブル64及びLUN管理テーブル61を変更して、RAID構成を変更する。

例えば、RAID構成変更部16は、RAID構成テーブル64及びLUN管理テーブル61において、failさせる元ディスク31aに代えて、ミラーリング処理部15により、当該failさせる元ディスクのデータを複製した予備ディスク31bを登録する。これにより、RAID構成変更部16は、RAIDグループ30aにおいて不安定なりビルド元ディスク31aに代えて予備ディスク31bを用いてRAID構成を組み直す。

30

【0118】

また、RAID構成変更部16は、LUN管理テーブル61の項6の利用予備ディスク名リストから、解放する予備ディスク31bを削除する。

ステップD9において、RAID構成変更部16は、failさせる元ディスク31aをfail状態に設定するコマンドを発行することでfailさせる。これにより、不要になった予備ディスク31bがシステムに返却（解放）されることになる。

40

【0119】

なお、上述したステップD2において、RAID構成変更部16は、failさせる元ディスクがないと判断された場合には、ステップD4以降の処理において、failさせる元ディスクに対して行なわれる処理が省略される。

次に、実施形態の一例としてのストレージ装置1におけるリード受信時の処理を図13に示すフローチャート（ステップE1～E13）に従って説明する。

【0120】

ステップE1において、IO制御部11は、受信したリード要求がリビルド処理中のRAIDへのリードであるか否かを確認する。確認の結果（ステップE2）、リードの要求先のRAIDグループ30aがリビルド中である場合には（ステップE2のYESルート

50

参照)、ステップE3において、I/O制御部11は、受信したリード要求がリビルド処理が完了したLUN領域へのリードであるか否かを確認する。確認の結果(ステップE4)、リードの要求先が未リビルドのLUN領域である場合や(ステップE4のNORルート参照)、リードの要求先のRAIDグループ30aがリビルド中でない場合には(ステップE2のNORルート参照)、従来手法によるリード処理を行なう。

【0121】

すなわち、I/O制御部11は、リード要求先の記憶装置31aにアクセスして、要求されたデータを読み出し、ホスト装置2へ送信する。又、リード要求先が故障ディスク31aである場合には、この故障ディスク31a以外の各リビルド元ディスク31aのデータに対して、パリティを用いたXOR演算を行なって故障ディスク31aのデータの復元を行なってホスト装置2へ送信する。

10

【0122】

一方、リードの要求先がリビルド済のLUN領域である場合には(ステップE4のYESルート参照)、ステップE5において、I/O制御部11は、アクセス対象のRAID構成ディスク31aが二重化されているかを確認する。すなわち、アクセス対象のRAID構成ディスク31aに対して予備ディスク31bが割り当てられ、この予備ディスク31bにRAID構成ディスク31aのデータの複製が格納されているかを確認する。

【0123】

確認の結果(ステップE6)、二重化されている場合には(ステップE6のYESルート参照)、ステップE7において、LUN管理テーブル61の項11及び項12のI/O禁止フラグを確認する。すなわち、リード対象のRAID構成ディスク31a及びその対応する予備ディスク31bの両方のI/O禁止フラグが“0(off)”であるかを確認する。

20

【0124】

確認の結果(ステップE8)、リード対象のRAID構成ディスク31a及びその対応する予備ディスク31bの各I/O禁止フラグが共にoffの場合には(ステップE8のYESルート参照)、ステップE9において、I/O制御部11は、RAID構成ディスク31aと対応する予備ディスク31bとのうちの一方をラウンドロビンにより交互に選択して、リード対象ディスク31を選択する。このように、RAID構成ディスク31aのデータが予備ディスク31bにより二重化されている場合に、これらのRAID構成ディスク31aと予備ディスク31bとを交互に選択してリードを行なう。これにより、各ディスク31に対するアクセスを分散し、負荷を軽減することで寿命を延ばし、信頼性を向上させることができる。

30

【0125】

一方、ステップE7における確認の結果、禁止フラグが共にoffではない場合、すなわち、RAID構成ディスク31a及び対応する予備ディスク31bの一方に、禁止フラグにonが設定されている場合には(ステップE8のNORルート参照)、ステップE10に移行する。

ステップE10において、I/O制御部11は、RAID構成ディスク31a及び対応する予備ディスク31bのうち、I/O禁止フラグにoffが設定されている方をリード対象ディスク31として選択する。

40

【0126】

その後、ステップE11において、I/O制御部11は、選択されたリード対象ディスク31について、LUN管理テーブル61の項9もしくは項10のリード/ライト中カウンタをカウントアップする。これにより、選択されたリード対象ディスク31がfailされたりフリーにされることを阻止することができ、信頼性が向上する。

また、ステップE5における確認の結果、二重化がされていない場合にも(ステップE6のNORルート参照)、このステップE11に移行する。

【0127】

ステップE12において、I/O制御部11は、リード対象ディスク31へリードを実行

50

し、このリードが完了すると、ステップ E 1 3 において、当該リード対象ディスク 3 1 について、LUN 管理テーブル 6 1 の項 9 もしくは項 1 0 のリード/ライト中カウンタをカウントダウンして、処理を終了する。

次に、実施形態の実施形態の一例としてのストレージ装置 1 におけるライト受信時の処理を図 1 4 に示すフローチャート (ステップ F 1 ~ F 1 4 ) に従って説明する。

【 0 1 2 8 】

ステップ F 1 において、I/O 制御部 1 1 は、受信したライト要求がリビルド処理中の RAID へのライトであるか否かを確認する。確認の結果 (ステップ F 2 )、ライトの要求先の RAID グループ 3 0 a がリビルド中である場合には (ステップ F 2 の YES ルート参照)、ステップ F 3 において、I/O 制御部 1 1 は、受信したライト要求がリビルド処理が完了した LUN 領域へのライトであるか否かを確認する。確認の結果 (ステップ F 4 )、ライトの要求先が未リビルドの LUN 領域である場合や (ステップ F 4 の NO ルート参照)、ライトの要求先の RAID グループ 3 0 a がリビルド中でない場合には (ステップ F 2 の NO ルート参照)、従来手法によるライト処理を行なう。

10

【 0 1 2 9 】

すなわち、I/O 制御部 1 1 は、ライト要求先の記憶装置 3 1 a にアクセスして、要求されたデータを書き込む。又、ライト要求先が故障ディスク 3 1 a である場合には、この故障ディスク 3 1 a 以外の各リビルド元ディスク 3 1 a のデータに対して、パリティを用いた XOR 演算を行なって故障ディスク 3 1 a のデータの復元を行ない、この復元したデータとつぎ合わせをしながらデータのライトを行なう。

20

【 0 1 3 0 】

一方、ライトの要求先がリビルド済の LUN 領域である場合には (ステップ F 4 の YES ルート参照)、ステップ F 5 において、I/O 制御部 1 1 は、ライトデータと RAID メンバディスクのデータからパリティデータを生成する。

その後、ステップ F 6 において、I/O 制御部 1 1 は、書き出し対象の RAID 構成ディスク 3 1 a が二重化されているかを確認する。すなわち、アクセス対象の RAID 構成ディスク 3 1 a に対して予備ディスク 3 1 b が割り当てられ、この予備ディスク 3 1 b に RAID 構成ディスク 3 1 a のデータの複製が格納されているかを確認する。

【 0 1 3 1 】

確認の結果 (ステップ F 7 )、二重化されている場合には (ステップ F 7 の YES ルート参照)、ステップ F 8 において、LUN 管理テーブル 6 1 の項 1 1 及び項 1 2 の I/O 禁止フラグを確認する。すなわち、ライト対象の RAID 構成ディスク 3 1 a 及びその対応する予備ディスク 3 1 b の両方の I/O 禁止フラグが “ 0 ( o f f ) ” であるかを確認する。

30

【 0 1 3 2 】

確認の結果 (ステップ F 9 )、ライト対象の RAID 構成ディスク 3 1 a 及びその対応する予備ディスク 3 1 b の各 I/O 禁止フラグが共に o f f の場合には (ステップ F 9 の YES ルート参照)、ステップ F 1 0 において、I/O 制御部 1 1 は、RAID 構成ディスク 3 1 a と対応する予備ディスク 3 1 b とを二重書きの対象ディスク 3 1 として選択する。

一方、ステップ F 8 における確認の結果、禁止フラグが共に o f f ではない場合、すなわち、RAID 構成ディスク 3 1 a 及び対応する予備ディスク 3 1 b の一方に、禁止フラグに o n が設定されている場合には (ステップ F 9 の NO ルート参照)、ステップ F 1 1 に移行する。

40

【 0 1 3 3 】

ステップ F 1 1 において、I/O 制御部 1 1 は、RAID 構成ディスク 3 1 a 及び対応する予備ディスク 3 1 b のうち、I/O 禁止フラグに o f f が設定されている方をライト対象ディスク 3 1 として選択する。

その後、ステップ F 1 2 において、I/O 制御部 1 1 は、選択されたライト対象ディスク 3 1 について、LUN 管理テーブル 6 1 の項 9 もしくは項 1 0 のリード/ライト中カウンタをカウントアップする。これにより、選択されたライト対象ディスク 3 1 が fail された

50

リフリーにされることを阻止することができ、信頼性が向上する。

【0134】

また、ステップF6における確認の結果、二重化がされていない場合にも（ステップF7のN0ルート参照）、このステップF12に移行する。

ステップF13において、IO制御部11は、ライト対象ディスク31へライトを実行し、このライトが完了すると、ステップF14において、当該ライト対象ディスク31について、LUN管理テーブル61の項9もしくは項10のリード/ライト中カウンタをカウントダウンして、処理を終了する。

【0135】

このように、実施形態の一例としてのストレージ装置1によれば、ミラーリング処理部15がRAID構成ディスク31aの複製を、当該RAID構成ディスク31aに対応する予備ディスク31bに作成することで、各RAID構成ディスク31aのデータを冗長化することができる。これにより、複数のRAID構成ディスク31aが故障した場合でも、予備ディスク31bを用いてRAIDを構成しなおすことにより、RAID内の全データのロストの発生を阻止することができる。これにより信頼性を向上させることができる。

10

【0136】

また、リビルド処理が完了したLUNは、リビルド処理による性能劣化の影響が緩和され、データ保護状態に速やかに戻ることができる。

リード時において、冗長化されたRAID構成ディスク31aのデータと予備ディスク31bのデータとをラウンドロビンで交互に選択してリード処理を行なうことで、リビルド中であっても高いストレージアクセス性能を実現することができる。更に、リード時に、冗長化されたRAID構成ディスク31aと予備ディスク31bとを均等に用いることで負荷を分散させることができる。

20

【0137】

RAID構成変更部16が、リビルド処理の完了時点で、RAIDグループ30a内に不安定なRAID構成ディスク31aを検知した場合に、この不安定なRAID構成ディスク31aに代えて対応する予備ディスク31bを用いてRAIDを構成し直す。これにより、RAID構成ディスク31aの故障を未然に阻止することができ、信頼性を向上させることができる。

30

【0138】

ミラーリング処理部15が、リビルド時にリビルド元ディスク31aからメモリ106に読み出されたデータを、当該リビルド元ディスク31aに対応する予備ディスク31bにデッドコピーすることで、リビルド元ディスク31aの複製を予備ディスク31bに容易に作成することができる。この際、コントローラ100のCPU110等による特別な制御が不要であり、コントローラ100の負荷の増大や処理速度の低下が発生することがない。

【0139】

そして、開示の技術は上述した実施形態に限定されるものではなく、本実施形態の趣旨を逸脱しない範囲で種々変形して実施することができる。本実施形態の各構成及び各処理は、必要に応じて取捨選択することができ、あるいは適宜組み合わせてもよい。

40

例えば、上述した実施形態においては、リビルド処理部13がいずれかの記憶装置31aの故障を検知した場合にリビルド処理を行なっている例を示しているが、これに限定されるものではない。例えば、いずれかの記憶装置31a故障の発生が予測される場合に、この故障の発生が予測される記憶装置31aを故障ディスク31aとしてリビルド処理を行なってもよく、又、予防保守の観点等から、異常のない記憶装置31aを故障ディスク31aとしてリビルド処理を行なってもよい。

【0140】

また、上述した実施形態においては、RAIDグループ30aが4本（3本+1本）のRAID構成ディスク31aによりRAID5を実現している例を示しているが、これに

50

限定されるものではない。例えば、RAID 2 ~ 4 や RAID 5 0 ( 5 + 0 )、RAID 6 , RAID 1 0 等、種々変形して実施することができる。

上述した実施形態においては、IO制御部 1 1 が、RAIDグループ 3 0 a の RAID 構成ディスク 3 1 a と、予備ディスクグループ 3 0 b の予備ディスク 3 1 b とをラウンドロビンで交互に選択してリード処理を行なっているが、これに限定されるものではない。すなわち、RAID 構成ディスク 3 1 a と、この RAID 構成ディスク 3 1 a の複製が格納された予備ディスク 3 1 b とから必ずしも交互にデータリードを行なう必要はなく、結果として均等にデータのリードを行なうことで負荷が分散できればよい。

【 0 1 4 1 】

また、上述した開示により本実施形態を当業者によって実施・製造することが可能である。

以上の実施形態に関し、更に以下の付記を開示する。

( 付記 1 )

冗長構成がなされた複数の記憶装置及び複数の予備記憶装置と通信路を介して通信可能に接続されるストレージ制御装置であって、

前記複数の記憶装置のうち復元対象記憶装置以外の冗長用記憶装置から読み出した冗長データを用いて、前記復元対象記憶装置のデータを、前記複数の予備記憶装置のうちの第 1 の予備記憶装置に再構成する再構成処理部と、

前記再構成処理部による再構成を行なう際に前記冗長用記憶装置から読み出したデータを、前記複数の予備記憶装置のうちの前記冗長用記憶装置に対応する第 2 の予備記憶装置に格納することで、前記冗長用記憶装置の複製を行なう複製処理部とを備えることを特徴とする、ストレージ制御装置。

【 0 1 4 2 】

( 付記 2 )

リード要求受信時には、前記冗長用記憶装置と、当該冗長用記憶装置のデータを格納する前記第 2 の予備記憶装置とを併用することを特徴とする、付記 1 記載のストレージ制御装置。

( 付記 3 )

ライト要求受信時には、前記冗長用記憶装置及び当該冗長用記憶装置のデータを格納する前記予備記憶装置の双方に書き込みを行なうことを特徴とする、付記 1 又は 2 記載のストレージ制御装置。

【 0 1 4 3 】

( 付記 4 )

前記記憶装置に対して前記予備記憶装置を割り当てる割当処理部を備え、

前記複製処理部が、前記冗長用記憶装置から読み出したデータを、前記割当処理部が割り当てた前記予備記憶装置に格納することを特徴とする、付記 1 ~ 3 のいずれか 1 項に記載のストレージ制御装置。

【 0 1 4 4 】

( 付記 5 )

前記記憶装置に割り当て可能な前記予備記憶装置の数が前記記憶装置の数よりも少ない場合に、

前記割当処理部が、

前記予備記憶装置を、安定度の低い前記記憶装置から優先して割り当てることを特徴とする、付記 4 記載のストレージ制御装置。

【 0 1 4 5 】

( 付記 6 )

前記再構成処理部による再構成の完了後に、安定度の低い前記冗長用記憶装置に代えて、当該安定度の低い前記冗長用記憶装置の複製がされた前記予備記憶装置を用いて、前記冗長構成を変更する冗長構成変更部を備えることを特徴とする、付記 1 ~ 付記 5 のいずれか 1 項に記載のストレージ制御装置。

## 【 0 1 4 6 】

( 付 記 7 )

冗長構成がなされた複数の記憶装置と、  
複数の予備記憶装置と、

前記複数の記憶装置のうち復元対象記憶装置以外の冗長用記憶装置から読み出した冗長データを用いて、前記復元対象記憶装置のデータを、前記複数の予備記憶装置のうちの第 1 の予備記憶装置に再構成する再構成処理部と、

前記再構成処理部による再構成を行なう際に前記記憶装置から読み出したデータを、前記複数の予備記憶装置のうちの前記冗長用記憶装置に対応する第 2 の予備記憶装置に格納することで、前記冗長用記憶装置の複製を行なう複製処理部と  
を備えることを特徴とする、ストレージシステム。

10

## 【 0 1 4 7 】

( 付 記 8 )

リード要求受信時には、前記冗長用記憶装置と、当該冗長用記憶装置のデータを格納する前記第 2 の予備記憶装置とを併用することを特徴とする、付記 7 記載のストレージシステム。

( 付 記 9 )

ライト要求受信時には、前記冗長用記憶装置及び当該冗長用記憶装置のデータを格納する前記予備記憶装置の双方に書き込みを行なうことを特徴とする、付記 7 又は 8 記載のストレージシステム。

20

## 【 0 1 4 8 】

( 付 記 1 0 )

前記記憶装置に対して前記予備記憶装置を割り当てる割当処理部を備え、

前記複製処理部が、前記冗長用記憶装置から読み出したデータを、前記割当処理部が割り当てた前記予備記憶装置に格納することを特徴とする、付記 7 ～ 9 のいずれか 1 項に記載のストレージシステム。

## 【 0 1 4 9 】

( 付 記 1 1 )

前記記憶装置に割り当て可能な前記予備記憶装置の数が前記記憶装置の数よりも少ない場合に、

30

前記割当処理部が、

前記予備記憶装置を、安定度の低い前記記憶装置から優先して割り当てることを特徴とする、付記 1 0 記載のストレージシステム。

## 【 0 1 5 0 】

( 付 記 1 2 )

前記再構成処理部による再構成の完了後に、安定度の低い前記冗長用記憶装置に代えて、当該安定度の低い前記冗長用記憶装置の複製がされた前記予備記憶装置を用いて、前記冗長構成を変更する冗長構成変更部を備えることを特徴とする、付記 7 ～ 付記 1 1 のいずれか 1 項に記載のストレージシステム。

## 【 0 1 5 1 】

40

( 付 記 1 3 )

冗長構成がなされた複数の記憶装置及び複数の予備記憶装置と通信路を介して通信可能に接続されるコンピュータに、

前記複数の記憶装置のうち復元対象記憶装置以外の冗長用記憶装置から読み出した冗長データを用いて、前記復元対象記憶装置のデータを、前記複数の予備記憶装置のうちの第 1 の予備記憶装置に再構成し、

前記再構成を行なう際に前記冗長用記憶装置から読み出したデータを、前記複数の予備記憶装置のうちの前記冗長用記憶装置に対応する第 2 の予備記憶装置に格納することで、前記冗長用記憶装置の複製を行なう  
処理を実行させることを特徴とする、制御プログラム。

50

## 【 0 1 5 2 】

( 付 記 1 4 )

リード要求受信時には、前記冗長用記憶装置と、当該冗長用記憶装置のデータを格納する前記第 2 の予備記憶装置とを併用する  
処理を前記コンピュータに実行させることを特徴とする、付記 1 3 記載の制御プログラム。

## 【 0 1 5 3 】

( 付 記 1 5 )

ライト要求受信時には、前記冗長用記憶装置及び当該冗長用記憶装置のデータを格納する前記予備記憶装置の双方に書き込みを行なう  
処理を前記コンピュータに実行させることを特徴とする、付記 1 3 又は 1 4 記載の制御プログラム。

10

## 【 0 1 5 4 】

( 付 記 1 6 )

前記記憶装置に対して前記予備記憶装置を割り当て、  
前記冗長用記憶装置から読み出したデータを、前記割当処理部が割り当てた前記予備記憶装置に格納する  
処理を前記コンピュータに実行させることを特徴とする、付記 1 3 ~ 1 5 のいずれか 1 項に記載の制御プログラム。

## 【 0 1 5 5 】

( 付 記 1 7 )

前記記憶装置に割り当て可能な前記予備記憶装置の数が前記記憶装置の数よりも少ない場合に、  
前記予備記憶装置を、安定度の低い前記記憶装置から優先して割り当てる  
処理を前記コンピュータに実行させることを特徴とする、付記 1 6 記載の制御プログラム。

20

## 【 0 1 5 6 】

( 付 記 1 8 )

前記再構成の完了後に、安定度の低い前記冗長用記憶装置に代えて、当該安定度の低い前記冗長用記憶装置の複製がされた前記予備記憶装置を用いて、前記冗長構成を変更する  
処理を前記コンピュータに実行させることを特徴とする、付記 1 3 ~ 付記 1 7 のいずれか 1 項に記載の制御プログラム。

30

## 【 符号の説明 】

## 【 0 1 5 7 】

1      ストレージ装置  
2      ホスト装置  
3 a , 3 b      スイッチ  
4      ストレージシステム  
1 1      I O 制御部  
1 2      R A I D 制御部  
1 3      リビルド処理部  
1 4      割当処理部  
1 5      ミラーリング処理部  
1 6      R A I D 構成変更部  
3 0      ドライブエンクロージャ  
3 0 a      R A I D グループ  
3 0 b      予備ディスクグループ  
3 1      記憶装置  
3 1 a , 3 1 a - 1 ~ 3 1 a - 4      記憶装置 ( R A I D 構成ディスク , リビルド元ディスク )

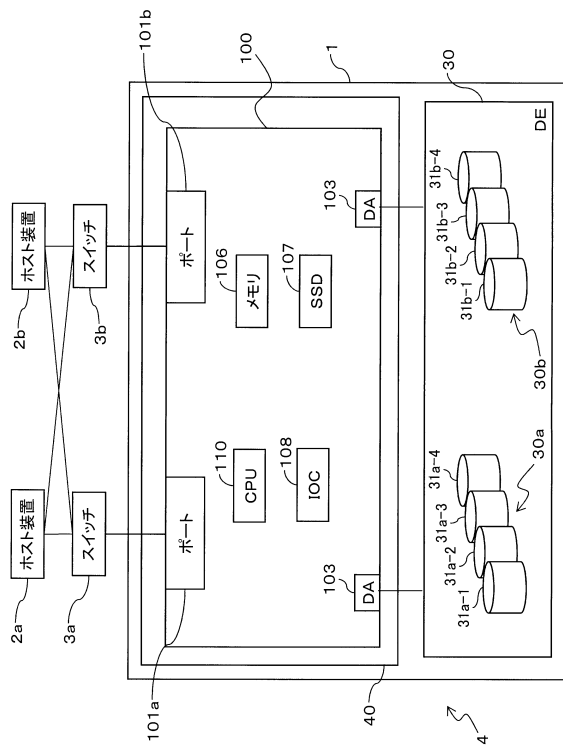
40

50

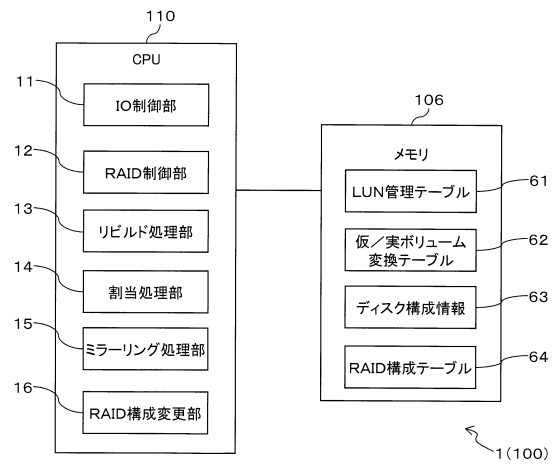
31b, 31b-1 ~ 31b-4 予備ディスク  
 40 コントローラエンクロージャ  
 61 LUN管理テーブル  
 62 仮ノ実ボリューム変換テーブル  
 63 ディスク構成情報  
 64 RAID構成テーブル  
 100 コントローラ(ストレージ制御装置)  
 101a, 101b ポート  
 103 DA  
 106 メモリ  
 107 SSD  
 108 IOC  
 110 CPU

10

【図1】



【図2】

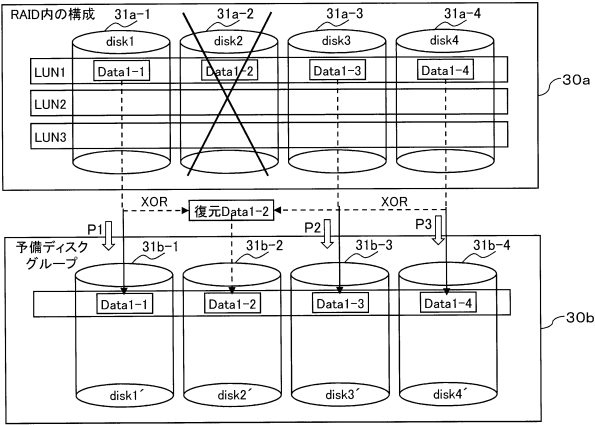


【図 3】

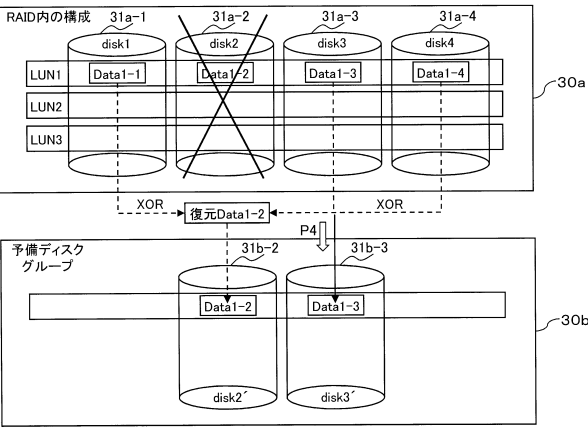
項	項目	内容	備考
1	LUN名	LUN1	
2	構成ディスク名リスト	(disk1, disk2, disk3, disk4)	
3	各ディスク上の位置・サイズ	Offset値・サイズ	
4	故障ディスク名	disk2	
5	LUNの状態	正常 or rebuild中 or rebuild済	
6	利用予備ディスク名リスト	(disk1', disk2', disk3', disk4') or (-, disk2', disk3', -)	予備ディスクが少ない時 高い順
7	ディスクの安定度リスト	disk1, disk4, disk3	初期値=0
8	リビルド済サイズ	サイズ	(U-D=n, ライト=0)
9	元ディスクのリード/ライト中カウンタ	((0, 0), (0, 0), (0, 0), (0, 0))	RAID構成の変更同期用
10	予備ディスクのリード/ライト中カウンタ	((0, 0), (0, 0), (0, 0), (0, 0))	ディスク毎のフラグ(1:禁止、 0:0, 0, 0)
11	元ディスクのIO禁止フラグ	(0, 0, 0, 0)	RAID構成の変更同期用
12	予備ディスクのIO禁止フラグ	(0, 0, 0, 0)	

61

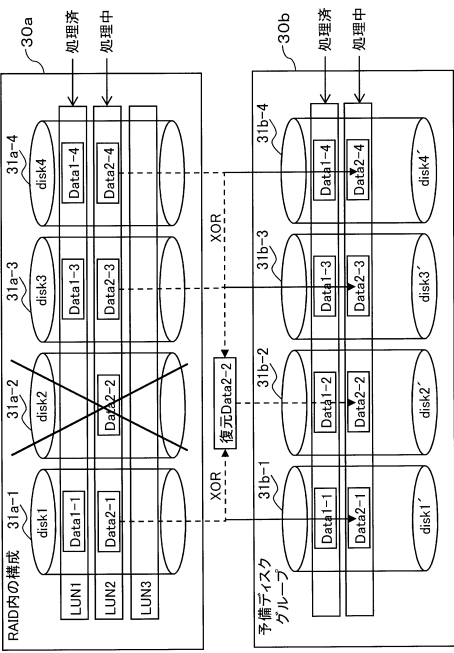
【図 4】



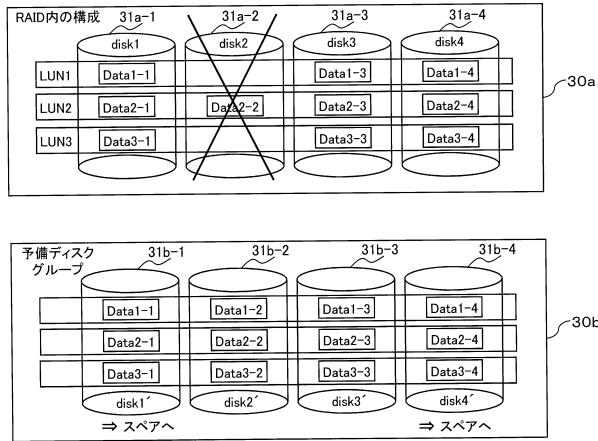
【図 5】



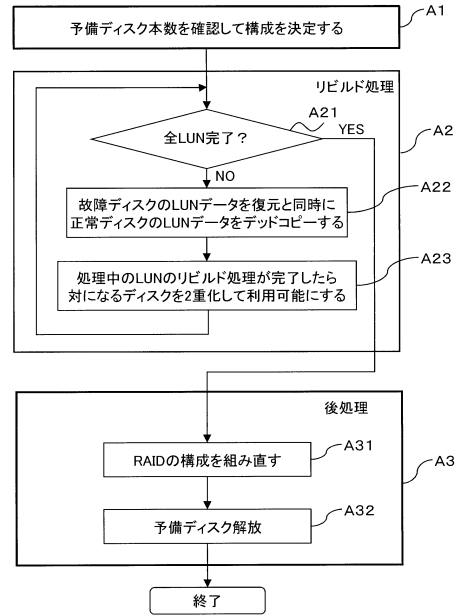
【図 6】



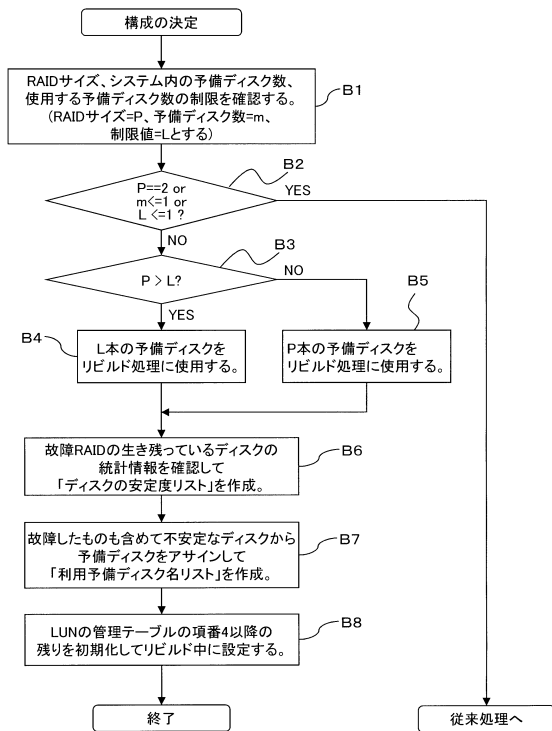
【図 7】



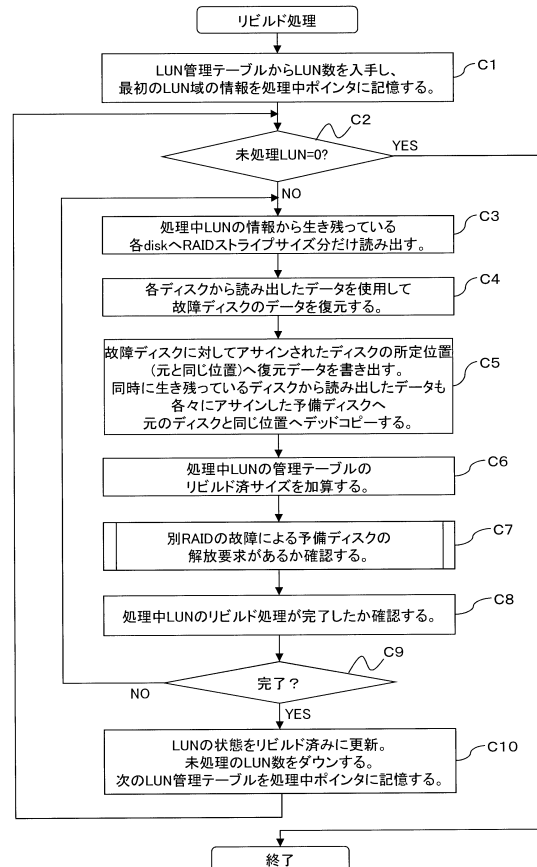
【図 8】



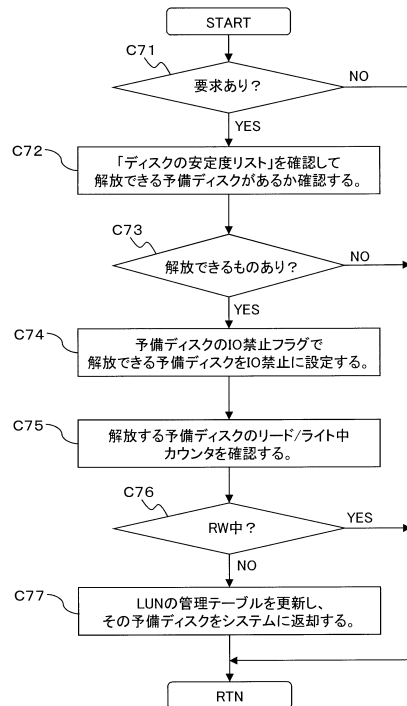
【図 9】



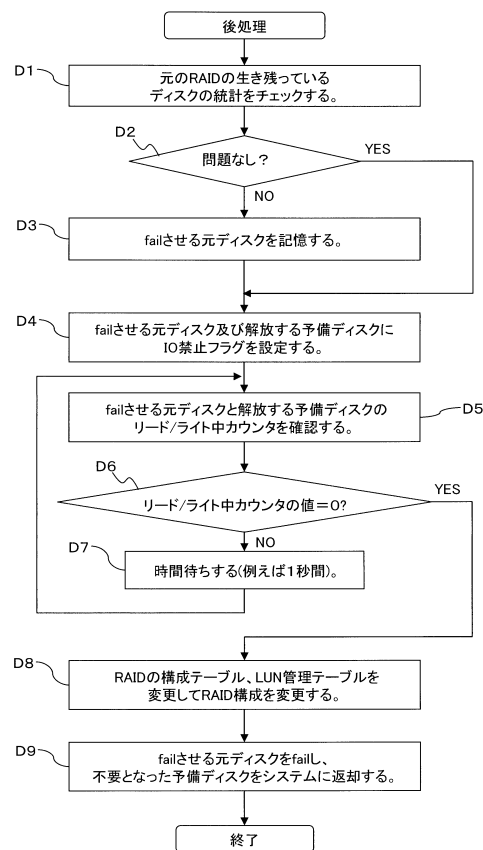
【図 10】



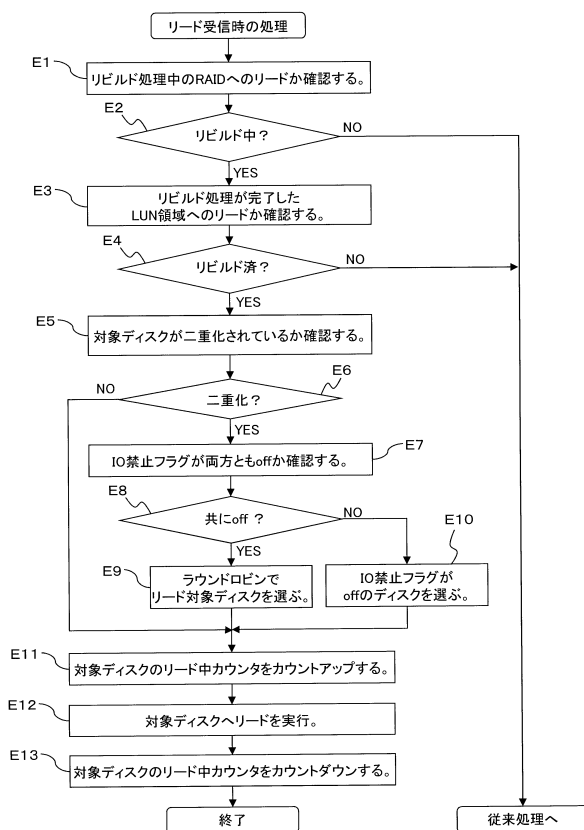
【図 1 1】



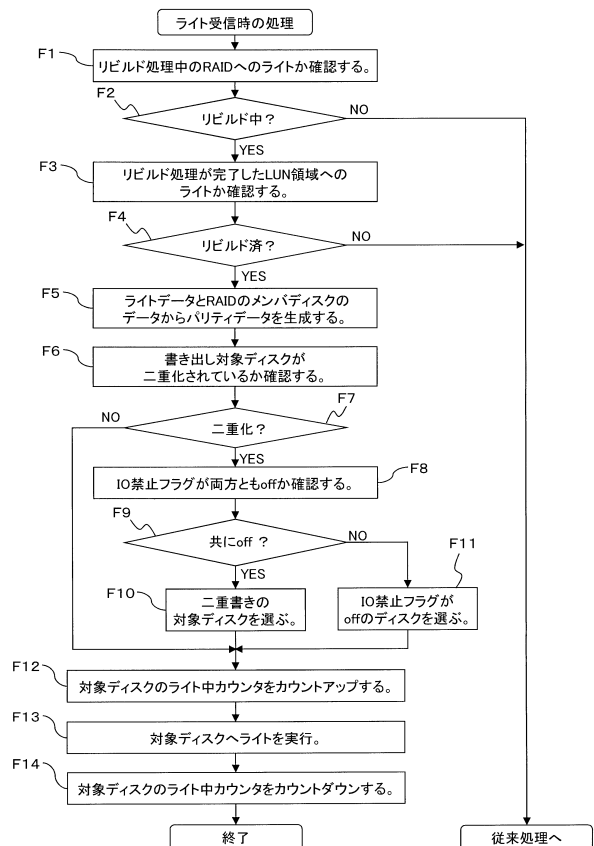
【図 1 2】



【図 1 3】



【図 1 4】



---

フロントページの続き

(56)参考文献 特開2007-233903(JP,A)  
米国特許出願公開第2007/0220313(US,A1)  
特開平11-085410(JP,A)  
特開平09-305324(JP,A)  
特開2006-164304(JP,A)  
米国特許出願公開第2008/0126839(US,A1)  
米国特許第7685463(US,B1)  
米国特許出願公開第2007/0294567(US,A1)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06 - 3/08  
G06F 13/10 - 13/14  
G06F 12/16