

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6212934号
(P6212934)

(45) 発行日 平成29年10月18日(2017.10.18)

(24) 登録日 平成29年9月29日(2017.9.29)

(51) Int.Cl.	F I
G 0 6 F 3/06 (2006.01)	G 0 6 F 3/06 3 0 1 Z
	G 0 6 F 3/06 5 4 0
	G 0 6 F 3/06 3 0 4 Z

請求項の数 8 (全 26 頁)

(21) 出願番号	特願2013-97648 (P2013-97648)	(73) 特許権者	000005223
(22) 出願日	平成25年5月7日(2013.5.7)		富士通株式会社
(65) 公開番号	特開2014-219787 (P2014-219787A)		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26年11月20日(2014.11.20)	(74) 代理人	100092152
審査請求日	平成28年2月26日(2016.2.26)		弁理士 服部 毅巖
		(72) 発明者	荻原 一隆
			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	圓道 浩史

最終頁に続く

(54) 【発明の名称】 ストレージシステム、情報処理装置の制御プログラム、およびストレージシステムの制御方法

(57) 【特許請求の範囲】

【請求項 1】

複数のストレージデバイスを有する複数のストレージユニットと、前記ストレージデバイスを管理する管理装置と、前記管理装置から割当を受けて前記ストレージデバイスと接続可能な情報処理装置と、を備えるストレージシステムであって、

前記情報処理装置は、

それぞれ異なる前記ストレージユニットに属する前記ストレージデバイスの割当を受けて第1のグループを構成する第1の構成部と、

前記第1のグループを構成するストレージデバイスの障害を検出する検出部と、

障害を検出したストレージデバイスを代替するストレージデバイスを、前記第1のグループを構成するその余のストレージデバイスが属するストレージユニットから割当を受け、前記その余のストレージデバイスが属するストレージユニットのうちの第1のストレージユニットから第1のストレージデバイスの割当を受けて前記第1のグループの構成を前記その余のストレージデバイスと前記第1のストレージデバイスとが属する第2のグループに構成する第2の構成部と、

前記その余のストレージデバイスが属するストレージユニットのうちの第2のストレージユニットから第2のストレージデバイスの割当を受けて、前記第2のストレージデバイスに前記第1のストレージデバイスを複製する複製部と、

を備えることを特徴とするストレージシステム。

【請求項 2】

10

20

前記情報処理装置は、前記第2のグループを構成するストレージデバイスが属さないストレージユニットから第3のストレージデバイスの割当を受けて、前記第1のストレージデバイスまたは前記第2のストレージデバイスを前記第3のストレージデバイスに複製して前記第2のグループの構成を前記その余のストレージデバイスと前記第3のストレージデバイスとが属する第3のグループに構成する第3の構成部を備えることを特徴とする請求項1記載のストレージシステム。

【請求項3】

前記情報処理装置は、前記第3のグループに構成した後、前記第1のストレージデバイスおよび前記第2のストレージデバイスの割当を解放することを特徴とする請求項2記載のストレージシステム。

10

【請求項4】

前記情報処理装置は、前記第2のグループを構成するストレージデバイスを特定可能な管理情報を記憶する記憶部を有し、

前記管理情報は、前記ストレージデバイスが複製を有するか否かを判別可能な複製判別情報を含む、

ことを特徴とする請求項1記載のストレージシステム。

【請求項5】

前記情報処理装置と前記ストレージユニットの接続と接続解除を切替可能なスイッチを備え、

前記管理装置は、前記スイッチの接続と接続解除の切替を制御して、前記情報処理装置に前記ストレージデバイスを割り当てる、

ことを特徴とする請求項1記載のストレージシステム。

20

【請求項6】

前記管理装置は、前記ストレージユニットの障害を検出し、検出した前記ストレージユニットの障害を前記情報処理装置に通知することを特徴とする請求項1記載のストレージシステム。

【請求項7】

複数のストレージデバイスを有する複数のストレージユニットを管理する管理装置と接続される情報処理装置の制御プログラムにおいて、

前記情報処理装置に、

それぞれ異なる前記ストレージユニットに属する前記ストレージデバイスの割当を受けて第1のグループを構成させ、

前記第1のグループを構成するストレージデバイスの障害を検出させ、

障害を検出したストレージデバイスを代替するストレージデバイスを、前記第1のグループを構成するその余のストレージデバイスが属するストレージユニットから割当を受ける場合に、前記その余のストレージデバイスが属するストレージユニットのうちの第1のストレージユニットから第1のストレージデバイスの割当を受けて前記第1のグループの構成を前記その余のストレージデバイスと前記第1のストレージデバイスとが属する第2のグループに構成させ、

30

前記その余のストレージデバイスが属するストレージユニットのうちの第2のストレージユニットから第2のストレージデバイスの割当を受けて、前記第2のストレージデバイスに前記第1のストレージデバイスを複製させる、

ことを特徴とする情報処理装置の制御プログラム。

40

【請求項8】

複数のストレージデバイスを有する複数のストレージユニットと、前記ストレージデバイスを管理する管理装置と、前記管理装置から割当を受けて前記ストレージデバイスと接続可能な情報処理装置と、を備えるストレージシステムの制御方法において、

前記情報処理装置が、

それぞれ異なる前記ストレージユニットに属する前記ストレージデバイスの割当を受けて第1のグループを構成し、

50

前記第 1 のグループを構成するストレージデバイスの障害を検出し、

障害を検出したストレージデバイスを代替するストレージデバイスを、前記第 1 のグループを構成するその余のストレージデバイスが属するストレージユニットから割当を受ける場合に、前記その余のストレージデバイスが属するストレージユニットのうちの第 1 のストレージユニットから第 1 のストレージデバイスの割当を受けて前記第 1 のグループの構成を前記その余のストレージデバイスと前記第 1 のストレージデバイスとが属する第 2 のグループに構成し、

前記その余のストレージデバイスが属するストレージユニットのうちの第 2 のストレージユニットから第 2 のストレージデバイスの割当を受けて、前記第 2 のストレージデバイスに前記第 1 のストレージデバイスを複製する、

ことを特徴とするストレージシステムの制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ストレージシステム、情報処理装置の制御プログラム、およびストレージシステムの制御方法に関する。

【背景技術】

【0002】

複数のディスクで R A I D (Redundant Arrays of Independent Disks) を構成してデータの可用性および信頼性を確保するストレージシステムが知られている。ストレージシステムでは、R A I D を構成するディスクの障害に備えて、複数の R A I D グループに共通のスペアディスクを用意したり、特定の R A I D グループに専用のスペアディスクを用意したりしてディスクの故障に備える。ストレージシステムは、ディスクの障害時にスペアディスクを交えたりビルド処理により、障害からの復旧をおこなう。

【先行技術文献】

【特許文献】

【0003】

【特許文献 1】特開 2 0 0 9 - 1 8 7 4 0 6 号公報

【特許文献 2】特開 2 0 0 5 - 1 0 0 2 5 9 号公報

【発明の概要】

【発明が解決しようとする課題】

【0004】

しかしながら、ストレージシステムは、所定数のディスクを収容可能な、ディスクボックスやディスクエンクロージャと呼ばれるストレージユニットを複数備えることで、多数のディスクをディスクプールとして管理する場合がある。

【0005】

このようなストレージシステムは、信頼性を考慮してそれぞれ異なるストレージユニットのディスクによって R A I D を構成するが、ディスクの障害によりスペアディスクを交えたりビルド処理をおこなうと、R A I D を構成する複数のディスクが同一のストレージユニットに属することがある。このとき、ストレージユニット単位の障害が発生すると、ストレージシステムは、複数のディスクで同時に障害が発生することとなり、データロストするおそれがある。

【0006】

1 つの側面では、本発明は、ストレージユニット単位の障害によるデータロストを防止できるストレージシステム、情報処理装置の制御プログラム、およびストレージシステムの制御方法を提供することを目的とする。

【課題を解決するための手段】

【0007】

上記目的を達成するために、以下に示すような、ストレージシステムが提供される。ストレージシステムは、複数のストレージデバイスを有する複数のストレージユニットと、

10

20

30

40

50

ストレージデバイスを管理する管理装置と、管理装置から割当を受けてストレージデバイスと接続可能な情報処理装置と、を備える。情報処理装置は、構成部と、検出部と、再構成部と、複製部と、を備える。構成部は、それぞれ異なるストレージユニットに属するストレージデバイスの割当を受けてグループを構成する。検出部は、グループを構成するストレージデバイスの障害を検出する。再構成部は、障害を検出したストレージデバイスを代替するストレージデバイスを、グループを構成するその余のストレージデバイスが属するストレージユニットから割当を受ける場合に、その余のストレージデバイスが属するストレージユニットのうちの第1のストレージユニットから第1のストレージデバイスの割当を受けてグループの再構成をおこなう。複製部は、その余のストレージデバイスが属するストレージユニットのうちの第2のストレージユニットから第2のストレージデバイスの割当を受けて、第2のストレージデバイスに第1のストレージデバイスを複製する。

10

【発明の効果】

【0008】

1 態様によれば、ストレージシステム、情報処理装置の制御プログラム、およびストレージシステムの制御方法において、ストレージユニット単位の障害によるデータロストを防止できる。

【図面の簡単な説明】

【0009】

【図1】第1の実施形態のストレージシステムの構成の一例を示す図である。

【図2】第2の実施形態のストレージシステムの構成の一例を示す図である。

20

【図3】第2の実施形態のディスクエンクロージャの構成の一例を示す図である。

【図4】第2の実施形態のサーバのハードウェア構成の一例を示す図である。

【図5】第2の実施形態のサーバが構成するRAIDグループの一例を示す図である。

【図6】第2の実施形態のRAID構成情報の一例を示す図である。

【図7】第2の実施形態の故障処理のフローチャートを示す図である。

【図8】第2の実施形態のディスク故障再構成処理のフローチャートを示す図である。

【図9】第2の実施形態のディスクエンクロージャ故障再構成処理のフローチャートを示す図である。

【図10】第2の実施形態のサーバが構成するRAIDグループの一例を示す図である。

【図11】第2の実施形態のサーバが構成するRAIDグループの一例を示す図である。

30

【図12】第2の実施形態のRAID構成情報の一例を示す図である。

【図13】第2の実施形態の復旧処理のフローチャートを示す図である。

【図14】第2の実施形態のサーバが構成するRAIDグループの一例を示す図である。

【図15】第3の実施形態のディスクエンクロージャ故障再構成処理のフローチャートを示す図である。

【図16】第3の実施形態のディスクエンクロージャ故障再構成処理のフローチャートを示す図である。

【図17】第3の実施形態のサーバが構成するRAIDグループの一例を示す図である。

【図18】第3の実施形態のRAID構成情報の一例を示す図である。

【図19】第3の実施形態のサーバが構成するRAIDグループの一例を示す図である。

40

【発明を実施するための形態】

【0010】

以下、実施の形態について、図面を参照しながら詳細に説明する。

〔第1の実施形態〕

まず、第1の実施形態のストレージシステムについて図1を用いて説明する。図1は、第1の実施形態のストレージシステムの構成の一例を示す図である。

【0011】

ストレージシステム1は、管理装置2と、情報処理装置3と、複数のストレージユニット4を備える。図1に示すストレージシステム1は、4つのストレージユニット4(4a, 4b, 4c, 4d)を備えるが、3または5以上のストレージユニット4を備えるもの

50

であってもよい。

【 0 0 1 2 】

ストレージユニット 4 は、複数のストレージデバイス 5 を有する。ストレージユニット 4 は、たとえば、ディスクボックスやディスクエンクロージャなどであり、複数のストレージデバイス 5 を収容する。ストレージユニット 4 は、収容する複数のストレージデバイス 5 の電源供給や冷却、所要のコントロールなどを担う。ストレージデバイス 5 は、データを格納可能なデバイスであり、たとえば、HDD (Hard Disk Drive) や SSD (Solid State Drive:フラッシュメモリドライブ) などである。

【 0 0 1 3 】

管理装置 2 は、ストレージシステム 1 におけるストレージ資源を管理し、情報処理装置 3 へのストレージデバイス 5 の割当を担う。管理装置 2 は、通信バス 7 を介してストレージユニット 4 と接続し、ストレージユニット 4、およびストレージユニット 4 が収容するストレージデバイス 5 を管理する。また、管理装置 2 は、通信バス 6 を介して情報処理装置 3 と接続し、情報処理装置 3 とストレージデバイス 5 との接続を管理する。

【 0 0 1 4 】

情報処理装置 3 は、管理装置 2 からストレージ資源の割当を受けて、割当を受けたストレージ資源にデータを格納する。情報処理装置 3 は、通信バス 8 を介してストレージユニット 4 と接続し、割当を受けたストレージ資源、すなわち割当を受けたストレージデバイス 5 と接続する。

【 0 0 1 5 】

情報処理装置 3 は、構成部 3 a と、検出部 3 b と、再構成部 3 c と、複製部 3 d を備える。構成部 3 a は、それぞれ異なるストレージユニット 4 に属するストレージデバイス 5 の割当を受けてグループを構成する。グループは、たとえば RAID グループであり、RAID の 1 つに RAID 5 などがある。図 1 に示す RAID グループの構成例では、ストレージデバイス 5 a , 5 b , 5 c , 5 d は、それぞれデータ「A」, 「B」, 「C」, 「D」を格納して RAID グループを構成する。ストレージデバイス 5 a , 5 b , 5 c , 5 d は、それぞれ、異なるストレージユニット 4 a , 4 b , 4 c , 4 d に属する。

【 0 0 1 6 】

このように、グループを構成するストレージデバイス 5 をそれぞれ異なるストレージユニット 4 から割り当ててすることで、ストレージシステム 1 は、グループを構成する複数のストレージデバイス 5 の同時故障によるデータロストの危険を低減する。

【 0 0 1 7 】

検出部 3 b は、グループを構成するストレージデバイス 5 の障害を検出する。検出部 3 b は、通信バス 8 を介して接続するストレージデバイス 5 の障害を検出することができる。また、検出部 3 b は、通信バス 6 を介して接続する管理装置 2 からの通知によりストレージデバイス 5 の障害を検出することができる。また、検出部 3 b は、管理装置 2 からの通知によりストレージデバイス 5 が属するストレージユニット 4 の障害を検出することができる。

【 0 0 1 8 】

再構成部 3 c は、検出部 3 b によるグループを構成するストレージデバイス 5 の障害検出により、障害を検出したストレージデバイス 5 をグループの構成から外す。再構成部 3 c は、あらたにストレージデバイス 5 の割当を受けてグループの再構成をおこなう。このとき、再構成部 3 c は、構成部 3 a がグループを構成したように、障害を検出していないストレージデバイス 5 とは異なるストレージユニット 4 からストレージデバイス 5 の割当を受けることが望ましい。しかしながら、再構成部 3 c は、ストレージ資源の状況によっては、障害を検出していないストレージデバイス 5 が属するストレージユニット 4 から割当を受けざるを得ない場合がある。このような場合に、再構成部 3 c は、障害を検出していないストレージデバイス 5 が属するストレージユニット 4 のうちの 1 つからストレージデバイス 5 の割当を受けてグループの再構成をおこなう。

【 0 0 1 9 】

図 1 に示すグループの構成例において、再構成部 3 c は、ストレージデバイス 5 a の障害検出を受けて、ストレージデバイス 5 b の属するストレージユニット 4 b からストレージデバイス 5 e を割り当てる。したがって、再構成部 3 c は、ストレージデバイス 5 a の障害検出後、ストレージデバイス 5 e , 5 b , 5 c , 5 d に、それぞれデータ「A 1 (A)」, 「B」, 「C」, 「D」を格納してグループを再構成する。ストレージデバイス 5 e , 5 b は、ともにストレージユニット 4 b に属し、ストレージデバイス 5 c , 5 d は、それぞれ、ストレージユニット 4 c , 4 d に属する。このとき、ストレージユニット 4 b に障害が発生すると、ストレージデバイス 5 e , 5 b にも障害が発生してデータをロストするおそれがある。

【 0 0 2 0 】

10

複製部 3 d は、障害を検出していないストレージデバイス 5 が属するストレージユニット 4 のうち再構成部 3 c が割当を受けたストレージデバイス 5 が属するストレージユニット 4 と異なるストレージユニット 4 からストレージデバイス 5 の割当を受ける。複製部 3 d は、複製部 3 d が割当を受けたストレージデバイス 5 に、再構成部 3 c が割当を受けたストレージデバイス 5 を複製する。

【 0 0 2 1 】

図 1 に示すグループの構成例において、複製部 3 d は、ストレージデバイス 5 e , 5 b がともにストレージユニット 4 b に属することから、ストレージユニット 4 b 以外のストレージユニット 4 からストレージデバイス 5 を割り当てる。この場合、複製部 3 d は、ストレージデバイス 5 c の属するストレージユニット 4 c からストレージデバイス 5 f を割り当てる。複製部 3 d は、ストレージデバイス 5 e をストレージデバイス 5 f に複製する。したがって、複製部 3 d は、ストレージデバイス 5 e , 5 f は、それぞれデータ「A 1 (A)」, 「A 2 (A)」を格納して R A I D 1 を構成する。

20

【 0 0 2 2 】

これにより、ストレージシステム 1 は、ストレージユニット 4 b , 4 c , 4 d のいずれに障害が発生しても、データをロストすることがない。また、ストレージシステム 1 は、R A I D 1 による 2 重化機会を限定するため、グループ構成時に使用するストレージデバイス 5 の数を抑制できる。したがって、ストレージシステム 1 は、低コストにして耐障害性に優れたシステムを構築可能である。

【 0 0 2 3 】

30

[第 2 の実施形態]

次に、第 2 の実施形態のストレージシステムの構成について図 2 を用いて説明する。図 2 は、第 2 の実施形態のストレージシステムの構成の一例を示す図である。

【 0 0 2 4 】

ストレージシステム 1 0 は、リソースマネージャ 1 1 と、サーバ 1 2 と、複数のディスクエンクロージャ 1 3 を備える。ディスクエンクロージャ 1 3 は、ストレージ資源として複数のディスク 1 5 を収容する。なお、図 2 に示すディスクエンクロージャ 1 3 は、6 つのディスク 1 5 を備えるが、2 以上を備えるもの（たとえば、2 4）であればいくつであってもよい。ディスク 1 5 は、データを格納可能なストレージデバイスであり、たとえば、H D D や S S D などである。

40

【 0 0 2 5 】

ディスクエンクロージャ 1 3 は、スイッチ 1 4 を備え、スイッチ 1 4 により外部機器（サーバ 1 2）とディスク 1 5 との接続および切り離しをおこなう。スイッチ 1 4 は、リソースマネージャ 1 1 の制御対象であり、通信バス 1 7 によりリソースマネージャ 1 1 と接続する。

【 0 0 2 6 】

ストレージシステム 1 0 は、ディスクエンクロージャ 1 3 を一単位にしてストレージ資源の交換あるいは増減をおこなうことができる。なお、図 2 に示すストレージシステム 1 0 は、4 つのディスクエンクロージャ 1 3（1 3 a , 1 3 b , 1 3 c , 1 3 d）を備えるが、3 または 5 以上のディスクエンクロージャ 1 3 を備えるものであってもよい。

50

【 0 0 2 7 】

リソースマネージャ 11 は、ストレージシステム 10 におけるストレージ資源を管理する管理装置であり、サーバ 12 へのディスク 15 の割当を担う。リソースマネージャ 11 は、通信バス 17 を介してディスクエンクロージャ 13 と接続し、ディスクエンクロージャ 13、およびディスクエンクロージャ 13 が収容するディスク 15 を管理する。また、リソースマネージャ 11 は、通信バス 16 を介してサーバ 12 と接続し、サーバ 12 とディスク 15 との接続を管理する。

【 0 0 2 8 】

リソースマネージャ 11 は、スイッチ 14 を制御し、サーバ 12 とディスク 15 との接続によりサーバ 12 へのディスク 15 の割当をおこなう。また、リソースマネージャ 11 は、スイッチ 14 を制御し、サーバ 12 とディスク 15 との接続解除（切り離し）によりサーバ 12 へのディスク 15 の割当解消をおこなう。なお、図 2 に示すストレージシステム 10 は、1 つのリソースマネージャ 11 を備えるが、2 以上のリソースマネージャ 11 を備えて冗長性確保あるいは負荷分散を図るものであってもよい。

【 0 0 2 9 】

サーバ 12 は、ストレージシステム 10 においてリソースマネージャ 11 からストレージ資源の割当を受ける情報処理装置である。サーバ 12 は、各ディスクエンクロージャ 13 が備えるスイッチ 14 と通信バス 18 を介して接続する。サーバ 12 は、スイッチ 14 を介して接続するディスク 15 の接続と接続解除を、ホットプラグ機能により認識できる。なお、図 2 に示すストレージシステム 10 は、3 つのサーバ 12（12a, 12b, 12c）を備えるが、任意の数のサーバ 12 を備えることができる。

【 0 0 3 0 】

サーバ 12 は、複数のディスクエンクロージャ 13 からそれぞれディスク 15 の割当を受けて、RAID（たとえば、RAID5）グループを構成する。サーバ 12 は、複数の RAID グループを構成可能であり、それぞれの RAID グループは識別情報によって区別される。このように、RAID グループを構成するディスク 15 をそれぞれ異なるディスクエンクロージャ 13 から割り当てることで、ストレージシステム 10 は、RAID グループを構成する複数のディスク 15 の同時故障によるデータロストの危険を低減する。

【 0 0 3 1 】

なお、複数のサーバ 12 と、複数のディスクエンクロージャ 13 を、通信バス 18 を介してそれぞれ接続するようにしたが、別途スイッチを設けて接続するようにしてもよい。

次に、第 2 の実施形態のディスクエンクロージャの構成について図 3 を用いて説明する。図 3 は、第 2 の実施形態のディスクエンクロージャの構成の一例を示す図である。

【 0 0 3 2 】

ディスクエンクロージャ 13 は、スイッチ 14 と、複数のディスク 15 と、コントローラ 25 と、電源部 26 と、冷却部 27 を備える。スイッチ 14 は、各ディスク 15 について外部機器との接続および切り離しをおこなう。冷却部 27 は、ディスク 15、電源部 26 を含めてディスクエンクロージャ 13 の筐体内を冷却する。電源部 26 は、コントローラ 25、冷却部 27、スイッチ 14、およびディスク 15 を含むディスクエンクロージャ 13 内の各機器に電力を供給する。

【 0 0 3 3 】

コントローラ 25 は、ディスクエンクロージャ 13 内の各機器を制御する。また、コントローラ 25 は、ディスクエンクロージャ 13 内の各機器の状態監視をおこない、ディスクエンクロージャ 13 内の各機器の故障、またはディスクエンクロージャ 13 全体としての故障を検出する。コントローラ 25 は、検出した故障をリソースマネージャ 11 に通知する。

【 0 0 3 4 】

次に、第 2 の実施形態のサーバのハードウェア構成について図 4 を用いて説明する。図 4 は、第 2 の実施形態のサーバのハードウェア構成の一例を示す図である。

サーバ 12 は、プロセッサ 101 によって装置全体が制御されている。プロセッサ 10

10

20

30

40

50

1 には、バス 1 0 6 を介して R A M (Random Access Memory) 1 0 2 と複数の周辺機器が接続されている。プロセッサ 1 0 1 は、マルチプロセッサであってもよい。プロセッサ 1 0 1 は、たとえば C P U (Central Processing Unit)、M P U (Micro Processing Unit)、D S P (Digital Signal Processor)、A S I C (Application Specific Integrated Circuit)、または P L D (Programmable Logic Device) である。またプロセッサ 1 0 1 は、C P U、M P U、D S P、A S I C、P L D のうちの 2 以上の要素の組み合わせであってもよい。

【 0 0 3 5 】

R A M 1 0 2 は、サーバ 1 2 の主記憶装置として使用される。R A M 1 0 2 には、プロセッサ 1 0 1 に実行させるオペレーティングシステム (Operating System) のプログラムやファームウェア、アプリケーションプログラムの少なくとも一部が一時的に格納される。また、R A M 1 0 2 には、プロセッサ 1 0 1 による処理に必要な各種データ (たとえば、システム制御の情報管理) が格納される。また、R A M 1 0 2 は、各種データの格納に用いるメモリと別体にキャッシュメモリを含むものであってもよい。

【 0 0 3 6 】

バス 1 0 6 に接続されている周辺機器としては、不揮発性メモリ 1 0 3、入出力インタフェース 1 0 4、および通信インタフェース 1 0 5 がある。

不揮発性メモリ 1 0 3 は、サーバ 1 2 の電源遮断時においても記憶内容を保持する。不揮発性メモリ 1 0 3 は、たとえば、E E P R O M (Electrically Erasable Programmable Read-Only Memory) やフラッシュメモリなどの半導体記憶装置や、H D D などである。また、不揮発性メモリ 1 0 3 は、サーバ 1 2 の補助記憶装置として使用される。不揮発性メモリ 1 0 3 には、オペレーティングシステムのプログラムやファームウェア、アプリケーションプログラム、および各種データが格納される。

【 0 0 3 7 】

入出力インタフェース 1 0 4 は、図示しない入出力装置と接続して入出力をおこなう。

通信インタフェース 1 0 5 は、通信バス 1 6 , 1 8 を形成するネットワークと接続することで、通信バス 1 6 , 1 8 を介して、リソースマネージャ 1 1 やディスクエンクロージャ 1 3 との間でデータの送受信をおこなう。

【 0 0 3 8 】

以上のようなハードウェア構成によって、第 2 の実施形態のサーバ 1 2 の処理機能を実現することができる。なお、サーバ 1 2 の他、リソースマネージャ 1 1、コントローラ 2 5、第 1 の実施形態に示した管理装置 2、情報処理装置 3、ストレージユニット 4 も、図 4 に示したサーバ 1 2 と同様のハードウェアにより実現することができる。

【 0 0 3 9 】

サーバ 1 2 は、たとえばコンピュータ読み取り可能な記録媒体に記録されたプログラムを実行することにより、第 2 の実施形態の処理機能を実現する。サーバ 1 2 に実行させる処理内容を記述したプログラムは、様々な記録媒体に記録しておくことができる。たとえば、サーバ 1 2 に実行させるプログラムを不揮発性メモリ 1 0 3 に格納しておくことができる。プロセッサ 1 0 1 は、不揮発性メモリ 1 0 3 内のプログラムの少なくとも一部を R A M 1 0 2 にロードし、プログラムを実行する。またサーバ 1 2 に実行させるプログラムを、図示しない光ディスク、メモリ装置、メモリカードなどの可搬型記録媒体に記録しておくこともできる。光ディスクには、D V D (Digital Versatile Disc)、D V D - R A M、C D - R O M (Compact Disc Read Only Memory)、C D - R (Recordable) / R W (ReWritable) などがある。メモリ装置は、入出力インタフェース 1 0 4 あるいは図示しない機器接続インタフェースとの通信機能を搭載した記録媒体である。たとえば、メモリ装置は、メモリリーダライタによりメモリカードへのデータの書き込み、またはメモリカードからのデータの読み出しをおこなうことができる。メモリカードは、カード型の記録媒体である。

【 0 0 4 0 】

可搬型記録媒体に格納されたプログラムは、たとえばプロセッサ 1 0 1 からの制御によ

10

20

30

40

50

り、不揮発性メモリ 103 にインストールされた後、実行可能となる。またプロセッサ 101 が、可搬型記録媒体から直接プログラムを読み出して実行することもできる。

【0041】

次に、第2の実施形態のサーバが構成するRAIDグループについて図5を用いて説明する。図5は、第2の実施形態のサーバが構成するRAIDグループの一例を示す図である。

【0042】

ディスクエンクロージャ13aは、複数のディスク15(「#A1」,「#A2」,「#A3」,・・・)を備える。ディスクエンクロージャ13bは、複数のディスク15(「#B1」,「#B2」,「#B3」,・・・)を備える。ディスクエンクロージャ13cは、複数のディスク15(「#C1」,「#C2」,「#C3」,・・・)を備える。ディスクエンクロージャ13dは、複数のディスク15(「#D1」,「#D2」,「#D3」,・・・)を備える。

10

【0043】

サーバ12(たとえば、サーバ12a)は、ディスクエンクロージャ13a,13b,13c,13dからそれぞれディスク15(「#A1」,「#B1」,「#C1」,「#D1」)の割当を受けてRAIDグループ30(たとえば、RAID5)を構成する。

【0044】

RAIDグループ30は、ディスク15(「#A1」,「#B1」,「#C1」,「#D1」)がそれぞれ異なるディスクエンクロージャ13に属する。そのため、サーバ12は、4つのディスクエンクロージャ13のうちの1つが故障しても、故障するディスク15が1つに限られる。したがって、サーバ12は、RAIDグループ30のデータへのアクセスを継続可能であり、またRAIDグループ30を再構成可能である。

20

【0045】

次に、第2の実施形態のサーバが管理するRAID構成情報について図6を用いて説明する。図6は、第2の実施形態のRAID構成情報の一例を示す図である。

RAID構成情報50は、サーバ12(たとえば、サーバ12a)が管理するRAIDグループの構成を示す情報である。サーバ12は、たとえば、不揮発性メモリ103にRAID構成情報50を保持する。RAID構成情報50は、RAIDグループID(Identification)、ブロックNo.、ステータス、ディスクエンクロージャID、ディスクIDを含む。

30

【0046】

RAIDグループIDは、サーバ12が管理するRAIDグループ30を識別するための情報である。RAID構成情報50に示すRAIDグループIDは、いずれも「#0001」であり、RAID構成情報50に示す情報は、同一のRAIDグループに属する情報である。

【0047】

ブロックNo.は、RAIDグループ30を構成するディスク15に付したシリアル番号である。ステータスは、RAIDグループ30を構成するディスク15の状態を示す。通常時のステータスは、「物理」である。ディスクエンクロージャIDは、ストレージシステム10内でディスクエンクロージャ13を一意に識別可能な識別情報である。ディスクIDは、各ディスクエンクロージャ13内でディスク15を一意に識別可能な識別情報である。したがって、サーバ12は、ディスクエンクロージャIDとディスクIDとから、ストレージシステム10内でディスク15を一意に識別できる。

40

【0048】

したがって、RAID構成情報50は、ブロックNo.「1」,「2」,「3」,「4」の4つのディスク15からRAIDグループID「#0001」のRAIDグループ30が構成されていることを示す。また、RAID構成情報50は、ステータス「物理」より、通常のRAIDグループが構成されていることを示す。また、RAID構成情報50は、ブロックNo.「1」のディスク15がディスクエンクロージャID「#A」、ディ

50

スクID「#1」であることを示す。同様に、RAID構成情報50は、ブロックNo.「2」のディスク15がディスクエンクロージャID「#B」、ディスクID「#1」であることを示す。同様に、RAID構成情報50は、ブロックNo.「3」のディスク15がディスクエンクロージャID「#C」、ディスクID「#1」であることを示す。同様に、RAID構成情報50は、ブロックNo.「4」のディスク15がディスクエンクロージャID「#D」、ディスクID「#1」であることを示す。

【0049】

次に、第2の実施形態の故障処理について図7を用いて説明する。図7は、第2の実施形態の故障処理のフローチャートを示す図である。

故障処理は、ストレージ資源の故障を検出してRAIDの再構成をおこなう処理である。故障処理は、サーバ12が定期的に行う処理である。

【0050】

【ステップS11】サーバ12のプロセッサ101（制御部）は、ディスク故障情報を取得する。ディスク故障情報は、サーバ12に割り当てのあるディスク15の故障に関する情報である。プロセッサ101は、通信バス18を介して定期または不定期にディスク15の稼働状態を監視することによりディスク15の故障を検出してディスク故障情報を生成する。プロセッサ101は、ディスク15へのポーリングまたはディスク15からの通知により、ディスク15の稼働状態を監視することができる。

【0051】

【ステップS12】制御部は、ディスクエンクロージャ故障情報を取得する。ディスクエンクロージャ故障情報は、リソースマネージャ11が管理するディスクエンクロージャ13の故障に関する情報である。リソースマネージャ11は、通信バス17を介して定期または不定期にディスクエンクロージャ13の稼働状態を監視することによりディスクエンクロージャ13の故障を検出してディスクエンクロージャ故障情報を生成する。リソースマネージャ11は、ディスクエンクロージャ13へのポーリングまたはディスクエンクロージャ13からの通知により、ディスクエンクロージャ13の稼働状態を監視することができる。なお、リソースマネージャ11は、ディスクエンクロージャ13を介して検出したディスク故障にもとづいてディスク故障情報を生成してサーバ12に通知するようにしてもよい。

【0052】

【ステップS13】制御部は、ディスク故障情報およびディスクエンクロージャ故障情報からディスク15の故障の有無を判定する。制御部は、ディスク15の故障ありと判定した場合にステップS14にすすみ、ディスク15の故障なしと判定した場合に故障処理を終了する。

【0053】

【ステップS14】制御部は、RAID構成情報を参照して故障したディスク15に係るRAIDグループの有無を判定する。制御部は、故障したディスク15に係るRAIDグループがある場合にステップS15にすすみ、故障したディスク15に係るRAIDグループがない場合に故障処理を終了する。

【0054】

【ステップS15】制御部は、故障したディスク15を構成要素とするRAIDグループの1つを特定する。

【ステップS16】制御部は、故障個所がディスク15かディスクエンクロージャ13かを判定する。制御部は、故障個所がディスク15の場合にステップS17にすすみ、故障個所がディスクエンクロージャ13の場合にステップS18にすすむ。

【0055】

【ステップS17】制御部は、ディスク故障再構成処理を実行する。ディスク故障再構成処理は、ディスク15が故障した場合にRAIDの再構成をおこなう処理である。詳細は、図8を用いて後で説明する。

【0056】

〔ステップS 1 8〕制御部は、ディスクエンクロージャ故障再構成処理を実行する。ディスクエンクロージャ故障再構成処理は、ディスクエンクロージャ1 3が故障した場合にR A I Dの再構成をおこなう処理である。詳細は、図9を用いて後で説明する。

【0 0 5 7】

〔ステップS 1 9〕制御部は、ステップS 1 5で特定したR A I Dグループの他に、故障したディスク1 5を構成要素とするR A I Dグループ、すなわち再構成対象のR A I Dグループがまだあるか否かを判定する。制御部は、再構成対象のR A I Dグループがあると判定した場合にステップS 1 5にすすみ、再構成対象のR A I Dグループがないと判定した場合に故障処理を終了する。

【0 0 5 8】

次に、第2の実施形態のディスク故障再構成処理について図8を用いて説明する。図8は、第2の実施形態のディスク故障再構成処理のフローチャートを示す図である。ディスク故障再構成処理は、故障処理のステップS 1 7でサーバ1 2が実行する処理である。

【0 0 5 9】

〔ステップS 2 1〕サーバ1 2のプロセッサ1 0 1（制御部）は、故障したディスク1 5が属するディスクエンクロージャ1 3（D E：Disk Enclosure）、すなわち故障ディスクのD Eに空きディスクがあるか否かを判定する。制御部は、リソースマネージャ1 1に照会することにより、故障ディスクのD Eに空きディスクがあるか否かを判定することができる。制御部は、故障ディスクのD Eに空きディスクがある場合にステップS 2 2にすすみ、故障ディスクのD Eに空きディスクがない場合にステップS 2 3にすすむ。

【0 0 6 0】

〔ステップS 2 2〕制御部は、故障ディスクが属するディスクエンクロージャ1 3からディスク1 5の割当を受けて、代替ディスクを獲得する。

〔ステップS 2 3〕制御部は、R A I Dグループを構成するディスクが属さないディスクエンクロージャ1 3（D E）、すなわちR A I D構成外D Eに空きディスクがあるか否かを判定する。制御部は、リソースマネージャ1 1に照会することにより、R A I D構成外D Eに空きディスクがあるか否かを判定することができる。制御部は、空きディスクがある場合にステップS 2 4にすすみ、空きディスクがない場合にステップS 2 5にすすむ。

【0 0 6 1】

〔ステップS 2 4〕制御部は、R A I D構成外D Eから代替ディスクを獲得する。

〔ステップS 2 5〕制御部は、R A I Dグループを構成するディスクが属する2以上のディスクエンクロージャ1 3（D E）、すなわちR A I D構成D Eに空きディスクがあるか否かを判定する。制御部は、リソースマネージャ1 1に照会することにより、2以上のR A I D構成D Eに空きディスクがあるか否かを判定することができる。制御部は、2以上のR A I D構成D Eに空きディスクがない場合にステップS 2 6にすすみ、2以上のR A I D構成D Eに空きディスクがある場合にステップS 2 7にすすむ。

【0 0 6 2】

〔ステップS 2 6〕制御部は、1つのR A I D構成D Eからディスク1 5の割当を受けて、代替ディスクを獲得する。

〔ステップS 2 7〕制御部は、2以上のR A I D構成D Eのうちの2つのR A I D構成D Eから1つずつディスク1 5の割当を受けて、代替ディスクを獲得する。

【0 0 6 3】

〔ステップS 2 8〕制御部は、獲得した代替ディスクを含めてR A I D再構成（第1のリビルド処理）をおこなう。このとき、ステップS 2 2，S 2 4において獲得した代替ディスクによりR A I D再構成したR A I Dグループは、R A I Dグループを構成するディスクがそれぞれ異なるディスクエンクロージャ1 3に属する。一方、ステップS 2 6において獲得した代替ディスクによりR A I D再構成したR A I Dグループは、R A I Dグループを構成するディスクが一部のディスクエンクロージャ1 3に重複して属する。そのため、このR A I Dグループは、一部のディスクエンクロージャ1 3が故障した場合に、デ

10

20

30

40

50

ータロストのおそれがある。

【 0 0 6 4 】

〔ステップS 2 9〕制御部は、獲得した代替ディスクを含めてRAID 1 併用 RAID 再構成（第2のリビルド処理）をおこなう。ステップS 2 7において獲得した代替ディスクによりRAID再構成したRAIDグループは、RAIDグループを構成するディスクが一部のディスクエンクロージャ1 3に重複して属する。制御部は、2つの代替ディスクがそれぞれの複製となるようにして、RAIDグループを再構成する。RAIDグループの再構成については、後で図1 0から図1 2を用いてRAIDグループの再構成例を挙げて説明する。

【 0 0 6 5 】

〔ステップS 3 0〕制御部は、RAID構成情報を更新してディスク故障再構成処理を終了する。

次に、第2の実施形態のディスクエンクロージャ故障再構成処理について図9を用いて説明する。図9は、第2の実施形態のディスクエンクロージャ故障再構成処理のフローチャートを示す図である。ディスクエンクロージャ故障再構成処理は、故障処理のステップS 1 8でサーバ1 2が実行する処理である。

【 0 0 6 6 】

〔ステップS 3 1〕サーバ1 2のプロセッサ1 0 1（制御部）は、RAIDグループを構成するディスクが属さないディスクエンクロージャ1 3（DE）、すなわちRAID構成外DEに空きディスクがあるか否かを判定する。制御部は、リソースマネージャ1 1に照会することにより、RAID構成外DEに空きディスクがあるか否かを判定することができる。制御部は、空きディスクがある場合にステップS 3 2にすすみ、空きディスクがない場合にステップS 3 3にすすむ。

【 0 0 6 7 】

〔ステップS 3 2〕制御部は、RAID構成外DEから代替ディスクを獲得する。

〔ステップS 3 3〕制御部は、RAIDグループを構成するディスクが属する2以上のディスクエンクロージャ1 3（DE）、すなわちRAID構成DEに空きディスクがあるか否かを判定する。制御部は、リソースマネージャ1 1に照会することにより、2以上のRAID構成DEに空きディスクがあるか否かを判定することができる。制御部は、2以上のRAID構成DEに空きディスクがない場合にステップS 3 4にすすみ、2以上のRAID構成DEに空きディスクがある場合にステップS 3 5にすすむ。

【 0 0 6 8 】

〔ステップS 3 4〕制御部は、1つのRAID構成DEからディスク1 5の割当を受けて、代替ディスクを獲得する。

〔ステップS 3 5〕制御部は、2以上のRAID構成DEのうちの2つのRAID構成DEから1つずつディスク1 5の割当を受けて、代替ディスクを獲得する。

【 0 0 6 9 】

〔ステップS 3 6〕制御部は、獲得した代替ディスクを含めてRAID再構成（第1のリビルド処理）をおこなう。このとき、ステップS 3 2において獲得した代替ディスクによりRAID再構成したRAIDグループは、RAIDグループを構成するディスクがそれぞれ異なるディスクエンクロージャ1 3に属する。一方、ステップS 3 4において獲得した代替ディスクによりRAID再構成したRAIDグループは、RAIDグループを構成するディスクが一部のディスクエンクロージャ1 3に重複して属する。そのため、このRAIDグループは、一部のディスクエンクロージャ1 3が故障した場合に、データロストのおそれがある。

【 0 0 7 0 】

〔ステップS 3 7〕制御部は、獲得した代替ディスクを含めてRAID 1 併用 RAID 再構成（第2のリビルド処理）をおこなう。ステップS 3 5において獲得した代替ディスクによりRAID再構成したRAIDグループは、RAIDグループを構成するディスクが一部のディスクエンクロージャ1 3に重複して属する。制御部は、2つの代替ディスク

10

20

30

40

50

がそれぞれの複製となるようにして、RAIDグループを再構成する。RAIDグループの再構成については、後で図10から図12を用いてRAIDグループの再構成例を挙げて説明する。

【0071】

[ステップS38] 制御部は、RAID構成情報を更新してディスクエンクロージャ故障再構成処理を終了する。

次に、第2の実施形態のRAID構成DEから代替ディスクを獲得する場合のRAID再構成について図10から図12を用いて説明する。まず、RAID構成DEの1つが故障して、代替ディスクを他のRAID構成DEから獲得しなければならない場合について図10を用いて説明する。図10は、第2の実施形態のサーバが構成するRAIDグループの一例を示す図である。

10

【0072】

ディスクエンクロージャ13aは、故障により、複数のディスク15(「#A1」,「#A2」,「#A3」,・・・)が故障した状態である。制御部は、ディスクエンクロージャ13aの故障検出により、他のディスクエンクロージャ13からディスク15「#A1」を代替するディスク15の割当を受ける。たとえば、ディスク故障再構成処理のステップS26、またはディスクエンクロージャ故障再構成処理のステップS34の場合、制御部は、ディスクエンクロージャ13bからディスク15「#B2」の割当を受ける。また、ディスク故障再構成処理のステップS27、またはディスクエンクロージャ故障再構成処理のステップS35の場合、制御部は、ディスク15「#B2」に加えて、ディスクエンクロージャ13cからディスク15「#C2」の割当を受ける。

20

【0073】

制御部は、ディスクエンクロージャ13aのディスク15「#A1」に代えてディスクエンクロージャ13bのディスク15「#B2」を加えて、RAIDグループ30をRAIDグループ30aとして再構成する。

【0074】

すなわち、制御部は、ディスクエンクロージャ13b,13c,13dからそれぞれディスク15(「#B2」,「#B1」,「#C1」,「#D1」)の割当を受けてRAIDグループ30をRAIDグループ30aとして再構成する。

【0075】

RAIDグループ30aは、ディスク15(「#B2」,「#B1」)が同一のディスクエンクロージャ13bに属する。そのため、サーバ12は、ディスクエンクロージャ13bが故障するとデータロストするおそれのある状態である。したがって、ディスク故障再構成処理のステップS26、またはディスクエンクロージャ故障再構成処理のステップS34を経てRAID再構成をおこなった場合、サーバ12は、ディスクエンクロージャ13bが故障した場合にデータロストするおそれがある。

30

【0076】

一方、サーバ12は、2つのRAID構成DEからそれぞれ代替ディスクを獲得できた場合は、データロストのおそれのないRAID再構成をおこなうことができる。2つのRAID構成DEからそれぞれ代替ディスクを獲得できた場合について図11を用いて説明する。図11は、第2の実施形態のサーバが構成するRAIDグループの一例を示す図である。

40

【0077】

制御部は、ディスク15「#B2」に加えて、ディスクエンクロージャ13cからディスク15「#C2」の割当を受けた場合、ディスク15「#B2」をディスク15「#C2」に複製する。すなわち、制御部は、ディスク15「#B2」とディスク15「#C2」とでRAID1を構成する。言い換えれば、制御部は、ディスクエンクロージャ13aのディスク15「#A1」を、ディスク15「#B2」とディスク15「#C2」とに置き換える。

【0078】

50

これにより、制御部は、ディスクエンクロージャ 13b, 13c, 13d からそれぞれディスク 15 (「#B2」, 「#C2」, 「#B1」, 「#C1」, 「#D1」) の割当を受けて、RAID1 を併用して RAID グループ 30 を RAID グループ 30b として再構成する。

【0079】

RAID グループ 30b は、ディスク 15 (「#B2」, 「#B1」) が同一のディスクエンクロージャ 13b に属し、ディスク 15 (「#C2」, 「#C1」) が同一のディスクエンクロージャ 13b に属する。しかしながら、ディスクエンクロージャ 13b, 13c のいずれか一方が故障しても、サーバ 12 は、RAID グループ 30 に対してアクセス可能である。

10

【0080】

したがって、ストレージシステム 10 は、ディスクエンクロージャ 13 (ストレージユニット単位) の障害によるデータロストを防止できる。また、ストレージシステム 10 は、通常時において、各ディスク 15 を 2 重化することを要しないから信頼性の向上とストレージ資源の効率的な利用とを両立することができる。

【0081】

RAID1 併用なしの RAID 再構成をおこなった場合、ディスク 15 が 2 台存在しているディスクエンクロージャ 13 の故障でデータロストのおそれがある。このとき、ディスク故障率を f_{hdd} 、ディスクエンクロージャ故障率を f_{de} とすると、故障率は、 $f_{hdd} \times (f_{hdd} + f_{de})$ となる。

20

【0082】

一方、RAID1 併用 RAID 再構成をおこなった場合、ディスク 15 が 2 台存在しているディスクエンクロージャ 13 の故障があってもデータロストのおそれがない。このとき、故障率は、 $f_{hdd} \times f_{hdd}$ となり、RAID1 併用なしの RAID 再構成をおこなった場合と比較して故障率を低減できる。

【0083】

次に、第 2 の実施形態の RAID1 併用 RAID 再構成後の RAID 構成情報について図 12 を用いて説明する。図 12 は、第 2 の実施形態の RAID 構成情報の一例を示す図である。

【0084】

RAID 構成情報 51 は、RAID 構成情報 50 を RAID1 併用 RAID 再構成後に更新した情報である。

30

RAID 構成情報 51 は、ブロック No. 「1」が 2 つと、ブロック No. 「2」, 「3」, 「4」が 1 つずつの合計 5 つのディスク 15 から RAID グループ ID 「#0001」の RAID グループ 30 が構成されていることを示す。また、RAID 構成情報 51 は、ステータス「RAID1」より、ブロック No. 「1」の 2 つのディスク 15 が RAID1 を構成していることを示す。ステータス「RAID1」は、ディスク 15 が複製を有することを示す。すなわち、ステータス「RAID1」は、ディスク 15 が複製を有するか否かを判別可能な複製判別情報に相当する。

【0085】

また、RAID 構成情報 51 は、ブロック No. 「1」の 1 つのディスク 15 がディスクエンクロージャ ID 「#B」、ディスク ID 「#2」であることを示す。また、RAID 構成情報 51 は、ブロック No. 「1」のもう 1 つのディスク 15 がディスクエンクロージャ ID 「#C」、ディスク ID 「#2」であることを示す。同様に、RAID 構成情報 51 は、ブロック No. 「2」のディスク 15 がディスクエンクロージャ ID 「#B」、ディスク ID 「#1」であることを示す。同様に、RAID 構成情報 51 は、ブロック No. 「3」のディスク 15 がディスクエンクロージャ ID 「#C」、ディスク ID 「#1」であることを示す。同様に、RAID 構成情報 51 は、ブロック No. 「4」のディスク 15 がディスクエンクロージャ ID 「#D」、ディスク ID 「#1」であることを示す。

40

50

【 0 0 8 6 】

次に、第2の実施形態の復旧処理について図13を用いて説明する。図13は、第2の実施形態の復旧処理のフローチャートを示す図である。復旧処理は、故障したディスク15や、故障したディスクエンクロージャ13の交換処理など、ストレージ資源のメンテナンスの終了を契機にしてサーバ12が実行する処理である。たとえば、復旧処理は、管理者による指示にもとづいて実行されるが、定期的に行われるものであってもよい。

【 0 0 8 7 】

〔ステップS41〕サーバ12のプロセッサ101（制御部）は、復旧情報を取得する。復旧情報は、サーバ12に割当のあるディスク15、またはサーバ12に割当のあるディスク15が属するディスクエンクロージャ13の故障に関する情報である。プロセッサ101は、通信バス18を介して定期または不定期にディスク15の稼働状態を監視することによりディスク15の復旧を検出してディスク15に関する復旧情報を生成する。リソースマネージャ11は、通信バス17を介して定期または不定期にディスクエンクロージャ13の稼働状態を監視することによりディスクエンクロージャ13の復旧を検出してディスクエンクロージャに関する復旧情報を生成する。なお、リソースマネージャ11は、ディスクエンクロージャ13を介して検出したディスク15の復旧にもとづいてディスクに関する復旧情報を生成してサーバ12に通知するようにしてもよい。なお、ここでいう復旧は、故障したディスク15、または故障したディスクエンクロージャ13の交換であるが、ディスク15の割当可能な状態への移行という観点から、ディスク15またはディスクエンクロージャ13の追加を含むものであってもよい。

【 0 0 8 8 】

〔ステップS42〕制御部は、復旧したディスク15を割当可能なRAIDグループの有無を判定する。制御部は、復旧したディスク15を割当可能なRAIDグループがある場合にステップS43にすすみ、復旧したディスク15を割当可能なRAIDグループがない場合に復旧処理を終了する。

【 0 0 8 9 】

〔ステップS43〕制御部は、復旧したディスク15を割当可能なRAIDグループの1つを特定する。

〔ステップS44〕制御部は、1つのディスクエンクロージャ13（同一DE）に属するRAIDグループを構成するディスク（RAID構成ディスク）の有無を判定する。制御部は、同一DEに属するRAID構成ディスクがある場合にステップS45にすすみ、同一DEに属するRAID構成ディスクがない場合にステップS49にすすむ。

【 0 0 9 0 】

〔ステップS45〕制御部は、復旧情報とRAID構成情報とにもとづいて、特定したRAIDグループに属するディスクエンクロージャ13（RAID構成DE）から代替ディスクを獲得可能か否かを判定する。制御部は、RAID構成DEから代替ディスクを獲得できる場合にステップS46にすすみ、RAID構成DEから代替ディスクを獲得できない場合にステップS49にすすむ。

【 0 0 9 1 】

〔ステップS46〕制御部は、代替ディスクを獲得する。

〔ステップS47〕制御部は、同一DEに属するRAID構成ディスクについて代替ディスクとの間でRAID1化をおこなう。

【 0 0 9 2 】

〔ステップS48〕制御部は、RAID構成情報を更新してステップS49にすすむ。

〔ステップS49〕制御部は、復旧情報とRAID構成情報とにもとづいて、特定したRAIDグループに属さないディスクエンクロージャ13（RAID構成外DE）から代替ディスクを獲得可能か否かを判定する。制御部は、RAID構成外DEから代替ディスクを獲得できる場合にステップS50にすすみ、RAID構成外DEから代替ディスクを獲得できない場合にステップS42にすすむ。

【 0 0 9 3 】

【ステップS50】制御部は、RAID構成情報のステータスを参照して、特定したRAIDグループにRAID1があるか否かを判定する。制御部は、特定したRAIDグループにRAID1がある場合にステップS51にすすみ、特定したRAIDグループにRAID1がない場合にステップS42にすすむ。

【0094】

【ステップS51】制御部は、RAID構成外DEから代替ディスクを獲得する。

【ステップS52】制御部は、RAID1を構成していたディスク15から代替ディスクにコピーバック処理をおこなう。

【0095】

【ステップS53】制御部は、RAID1を構成していたディスク15を解放する。

10

【ステップS54】制御部は、RAID構成情報を更新してステップS42にすすむ。

ここで、獲得ディスクへのコピーバックと、RAID1を構成していたディスク15の解放について図14を用いて説明する。図14は、第2の実施形態のサーバが構成するRAIDグループの一例を示す図である。

【0096】

図11に示したRAIDグループ30bは、ディスクエンクロージャ13aが故障し、RAID1を併用してRAIDグループを構成している状態である。ここで、図14に示すようにディスクエンクロージャ13aが復旧した場合、制御部は、ディスクエンクロージャ13aのディスク15（たとえば、ディスク15「#A1」）を獲得する。制御部は、RAID1を構成していたディスク15「#B2」からディスク15「#A1」にコピーバック処理をおこなう。これにより、サーバ12は、RAIDグループ30cを再構成することができる。また、制御部は、RAID1を構成していたディスク15「#B2」,「#C2」を解放対象31として、リソースマネージャ11に解放依頼をおこない、サーバ12への割当から解放する。

20

【0097】

したがって、ストレージシステム10は、ディスクエンクロージャ13（ストレージユニット単位）の復旧により、RAID1を併用していたRAIDグループを復旧することができる。

【0098】

【第3の実施形態】

30

次に、第3の実施形態のディスクエンクロージャ故障再構成処理について図15、図16を用いて説明する。図15および図16は、第3の実施形態のディスクエンクロージャ故障再構成処理のフローチャートを示す図である。第3の実施形態のディスクエンクロージャ故障再構成処理は、RAID1を併用するRAIDグループを構成するディスクエンクロージャ13の故障に対応する。

【0099】

【ステップS61】サーバ12のプロセッサ101（制御部）は、故障ディスクのステータスを確認する。制御部は、故障ディスクのステータスが「物理」である場合にステップS66にすすみ、故障ディスクのステータスが「RAID1」である場合にステップS62にすすむ。

40

【0100】

【ステップS62】制御部は、RAIDグループを構成するディスクが属するディスクエンクロージャ13（RAID構成DE）のうち、RAID1を復元可能なディスクエンクロージャ13（DE）に空きディスクがあるか否かを判定する。RAID1を復元可能なディスクエンクロージャ13は、ステータスが「RAID1」のディスク15が属するディスクエンクロージャ13と異なるディスクエンクロージャ13である。制御部は、RAID1を復元可能なディスクエンクロージャ13（DE）に空きディスクがある場合にステップS63にすすみ、空きディスクがない場合にステップS74にすすむ。

【0101】

【ステップS63】制御部は、RAID1を復元可能なディスクエンクロージャ13（

50

DE) から代替ディスクを獲得する。

[ステップS64] 制御部は、故障ディスクに代えて代替ディスクでRAID1を復元する。

【0102】

[ステップS65] 制御部は、RAID構成情報を更新してステップS74にすすむ。

[ステップS66] 制御部は、RAIDグループを構成するディスクが属さないディスクエンクロージャ13(RAID構成外DE)に空きディスクがあるか否かを判定する。制御部は、RAID構成外DEに空きディスクがある場合にステップS67にすすみ、RAID構成外DEに空きディスクがない場合にステップS68にすすむ。

【0103】

[ステップS67] 制御部は、RAID構成外DEから代替ディスクを獲得する。

[ステップS68] 制御部は、2以上のRAID構成DEに空きディスクがあるか否かを判定する。制御部は、2以上のRAID構成DEに空きディスクがない場合にステップS69にすすみ、2以上のRAID構成DEに空きディスクがある場合にステップS70にすすむ。

【0104】

[ステップS69] 制御部は、1つのRAID構成DEからディスク15の割当を受けて、代替ディスクを獲得する。

[ステップS70] 制御部は、2以上のRAID構成DEのうちの2つのRAID構成DEから1つずつディスク15の割当を受けて、代替ディスクを獲得する。

【0105】

[ステップS71] 制御部は、獲得した代替ディスクを含めてRAID再構成(第1のリビルド処理)をおこなう。このとき、ステップS67において獲得した代替ディスクによりRAID再構成したRAIDグループは、RAIDグループを構成するディスクがそれぞれ異なるディスクエンクロージャ13に属する。一方、ステップS69において獲得した代替ディスクによりRAID再構成したRAIDグループは、RAIDグループを構成するディスクが一部のディスクエンクロージャ13に重複して属する。そのため、このRAIDグループは、一部のディスクエンクロージャ13が故障した場合に、データロスのおそれがある。

【0106】

[ステップS72] 制御部は、獲得した代替ディスクを含めてRAID1併用RAID再構成(第2のリビルド処理)をおこなう。ステップS70において獲得した代替ディスクによりRAID再構成したRAIDグループは、RAIDグループを構成するディスクが一部のディスクエンクロージャ13に重複して属する。制御部は、2つの代替ディスクがそれぞれの複製となるようにして、RAIDグループを再構成する。

【0107】

[ステップS73] 制御部は、RAID構成情報を更新する。

[ステップS74] 制御部は、すべての故障ディスクについて代替ディスクを獲得したか否かを判定する。制御部は、すべての故障ディスクについて代替ディスクを獲得していない場合にステップS61にすすみ、すべての故障ディスクについて代替ディスクを獲得した場合にディスクエンクロージャ故障再構成処理を終了する。

【0108】

これにより、制御部は、故障ディスクのステータスが「RAID1」であっても、さらにRAID1を併用したRAIDグループを再構成することができる。したがって、ストレージシステム10は、ディスクエンクロージャ13(ストレージユニット単位)の繰り返しの障害があってもデータロスを防止できる。また、ストレージシステム10は、通常時において、各ディスク15を2重化することを要しないから信頼性の向上とストレージ資源の効率的な利用とを両立することができる。

【0109】

次に、第3の実施形態のRAID構成DEから代替ディスクを獲得する場合のRAID

10

20

30

40

50

再構成について図 17 および図 18 を用いて説明する。まず、RAID 構成 DE の 1 つが故障して、代替ディスクを他の RAID 構成 DE から獲得した図 11 に示す RAID グループ 30 b の状態から、さらに RAID 構成 DE の 1 つが故障した場合について図 17 を用いて説明する。図 17 は、第 3 の実施形態のサーバが構成する RAID グループの一例を示す図である。

【0110】

ディスクエンクロージャ 13 b は、故障により、複数のディスク 15 (「# B 1」, 「# B 2」, 「# B 3」, ...) が故障した状態である。制御部は、ディスクエンクロージャ 13 b の故障検出により、他のディスクエンクロージャ 13 からディスク 15 (「# B 1」, 「# B 2」) を代替するディスク 15 の割当を受ける。たとえば、ディスクエンクロージャ故障再構成処理のステップ S 67 の場合、制御部は、ディスク 15 「# B 1」を代替するため、ディスクエンクロージャ 13 c, 13 d からそれぞれディスク 15 「# C 3」, 「# D 2」の割当を受ける。また、ディスクエンクロージャ故障再構成処理のステップ S 70 の場合、制御部は、ディスク 15 「# B 2」を代替するため、ディスクエンクロージャ 13 d からディスク 15 「# D 3」の割当を受ける。

10

【0111】

制御部は、ディスクエンクロージャ 13 b のディスク 15 「# B 1」に代えてディスクエンクロージャ 13 c のディスク 15 「# C 3」を加えて、RAID グループ 30 を再構成する。

【0112】

20

すなわち、制御部は、ディスクエンクロージャ 13 c, 13 d からそれぞれディスク 15 (「# C 2」, 「# C 3」, 「# C 1」, 「# D 1」) の割当を受けて RAID グループ 30 を再構成する。また、制御部は、ディスク 15 「# C 2」をディスク 15 「# D 3」に複製して RAID 1 を構成する。また、制御部は、ディスク 15 「# C 3」をディスク 15 「# D 2」に複製して RAID 1 を構成する。

【0113】

これにより、制御部は、2 組の RAID 1 を併用して RAID グループ 30 d として RAID グループ 30 を再構成する。

RAID グループ 30 d は、ディスクエンクロージャ 13 c, 13 d のいずれか一方が故障しても、サーバ 12 が RAID グループ 30 に対してアクセス可能である。

30

【0114】

したがって、ストレージシステム 10 は、ディスクエンクロージャ 13 (ストレージユニット単位) の障害によるデータロストを防止できる。また、ストレージシステム 10 は、通常時において、各ディスク 15 を 2 重化することを要しないから信頼性の向上とストレージ資源の効率的な利用とを両立することができる。

【0115】

また、故障ディスクのステータスが「RAID 1」の場合に、制御部は、故障ディスクとペアになるステータスが「RAID 1」のディスク 15 からデータを取得し、代替ディスクへの書き込みをおこなうことができる。また、故障ディスクとペアになるステータスが「RAID 1」のディスク 15 を RAID グループの再構成に利用することで、RAID グループの再構成のためのデータリードタイムを短縮することができる。

40

【0116】

次に、第 3 の実施形態の RAID 1 併用 RAID 再構成後の RAID 構成情報について図 18 を用いて説明する。図 18 は、第 3 の実施形態の RAID 構成情報の一例を示す図である。

【0117】

RAID 構成情報 52 は、2 組の RAID 1 併用による RAID 再構成後に RAID 構成情報 51 を更新した情報である。

RAID 構成情報 52 は、ブロック No. 「1」, 「2」が 2 つと、ブロック No. 「3」, 「4」が 1 つずつの合計 6 つのディスク 15 から RAID グループ ID 「# 000

50

1」のRAIDグループ30が構成されていることを示す。また、RAID構成情報52は、ステータス「RAID1」より、ブロックNo.「1」の2つのディスク15がRAID1を構成し、ブロックNo.「2」の2つのディスク15がもう1つのRAID1を構成していることを示す。また、RAID構成情報52は、ブロックNo.「1」の1つのディスク15がディスクエンクロージャID「#C」、ディスクID「#2」であることを示す。また、RAID構成情報52は、ブロックNo.「1」のもう1つのディスク15がディスクエンクロージャID「#D」、ディスクID「#3」であることを示す。また、RAID構成情報52は、ブロックNo.「2」の1つのディスク15がディスクエンクロージャID「#C」、ディスクID「#3」であることを示す。また、RAID構成情報52は、ブロックNo.「2」のもう1つのディスク15がディスクエンクロージャID「#D」、ディスクID「#2」であることを示す。同様に、RAID構成情報52は、ブロックNo.「3」のディスク15がディスクエンクロージャID「#C」、ディスクID「#1」であることを示す。同様に、RAID構成情報52は、ブロックNo.「4」のディスク15がディスクエンクロージャID「#D」、ディスクID「#1」であることを示す。

10

【0118】

ここで、獲得ディスクへのコピーバックと、RAID1を構成していたディスク15の解放について図19を用いて説明する。図19は、第3の実施形態のサーバが構成するRAIDグループの一例を示す図である。

【0119】

20

図17に示したRAIDグループ30dは、ディスクエンクロージャ13a, 13bが故障し、RAID1を併用してRAIDグループを構成している状態である。ここで、図19に示すようにディスクエンクロージャ13a, 13bが復旧した場合、制御部は、ディスクエンクロージャ13a, 13bからそれぞれディスク15（たとえば、ディスク15「#A1」, 「#B1」）を獲得する。制御部は、RAID1を構成していたディスク15「#C2」からディスク15「#A1」にコピーバック処理をおこなう。また、制御部は、もう1つのRAID1を構成していたディスク15「#C3」からディスク15「#B1」にコピーバック処理をおこなう。これにより、サーバ12は、RAIDグループ30eを再構成することができる。また、制御部は、RAID1を構成していたディスク15「#C2」, 「#C3」, 「#D2」, 「#D3」を解放対象32として、リソースマネージャ11に解放依頼をおこない、サーバ12への割当から解放する。

30

【0120】

したがって、ストレージシステム10は、ディスクエンクロージャ13（ストレージユニット単位）の復旧により、RAID1を併用していたRAIDグループを復旧することができる。

【0121】

なお、上記の処理機能は、コンピュータによって実現することができる。その場合、管理装置2、情報処理装置3、リソースマネージャ11、サーバ12が有すべき機能の処理内容を記述したプログラムが提供される。そのプログラムをコンピュータで実行することにより、上記処理機能がコンピュータ上で実現される。処理内容を記述したプログラムは、コンピュータで読み取り可能な記録媒体に記録しておくことができる。コンピュータで読み取り可能な記録媒体としては、磁気記憶装置、光ディスク、光磁気記録媒体、半導体メモリなどがある。磁気記憶装置には、ハードディスク装置（HDD）、フレキシブルディスク（FD）、磁気テープなどがある。光ディスクには、DVD、DVD-RAM、CD-ROM/RWなどがある。光磁気記録媒体には、MO（Magneto-Optical disk）などがある。

40

【0122】

プログラムを流通させる場合には、たとえば、そのプログラムが記録されたDVD、CD-ROMなどの可搬型記録媒体が販売される。また、プログラムをサーバコンピュータの記憶装置に格納しておき、ネットワークを介して、サーバコンピュータから他のコンピ

50

ュータにそのプログラムを転送することもできる。

【 0 1 2 3 】

プログラムを実行するコンピュータは、たとえば、可搬型記録媒体に記録されたプログラムもしくはサーバコンピュータから転送されたプログラムを、自己の記憶装置に格納する。そして、コンピュータは、自己の記憶装置からプログラムを読み取り、プログラムに従った処理を実行する。なお、コンピュータは、可搬型記録媒体から直接プログラムを読み取り、そのプログラムに従った処理を実行することもできる。また、コンピュータは、ネットワークを介して接続されたサーバコンピュータからプログラムが転送されるごとに、逐次、受け取ったプログラムに従った処理を実行することもできる。

【 0 1 2 4 】

また、上記の処理機能の少なくとも一部を、DSP、ASIC、PLDなどの電子回路で実現することもできる。

【符号の説明】

【 0 1 2 5 】

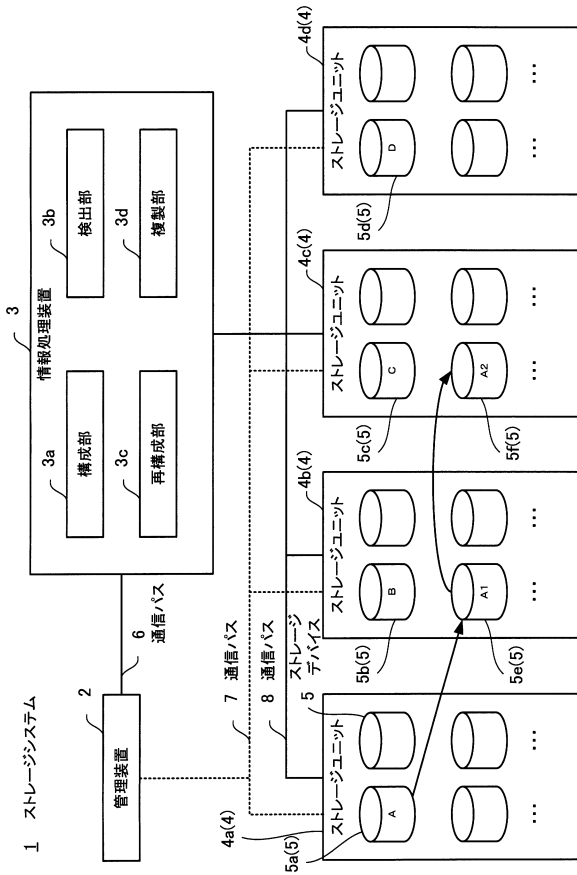
- 1、10 ストレージシステム
- 2 管理装置
- 3 情報処理装置
- 3 a 構成部
- 3 b 検出部
- 3 c 再構成部
- 3 d 複製部
- 4, 4 a, 4 b, 4 c, 4 d ストレージユニット
- 5, 5 a, 5 b, 5 c, 5 d, 5 e, 5 f ストレージデバイス
- 6, 7, 8, 16, 17, 18 通信バス
- 11 リソースマネージャ
- 12, 12 a, 12 b, 12 c サーバ
- 13, 13 a, 13 b, 13 c, 13 d ディスクエンクロージャ
- 14 スイッチ
- 15 ディスク
- 25 コントローラ
- 26 電源部
- 27 冷却部
- 101 プロセッサ
- 102 RAM
- 103 不揮発性メモリ
- 104 入出力インタフェース
- 105 通信インタフェース
- 106 バス

10

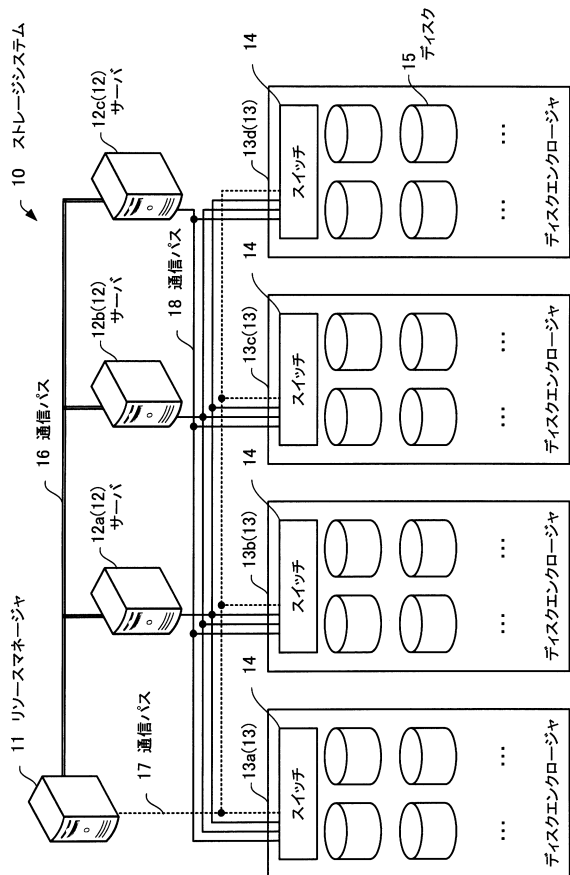
20

30

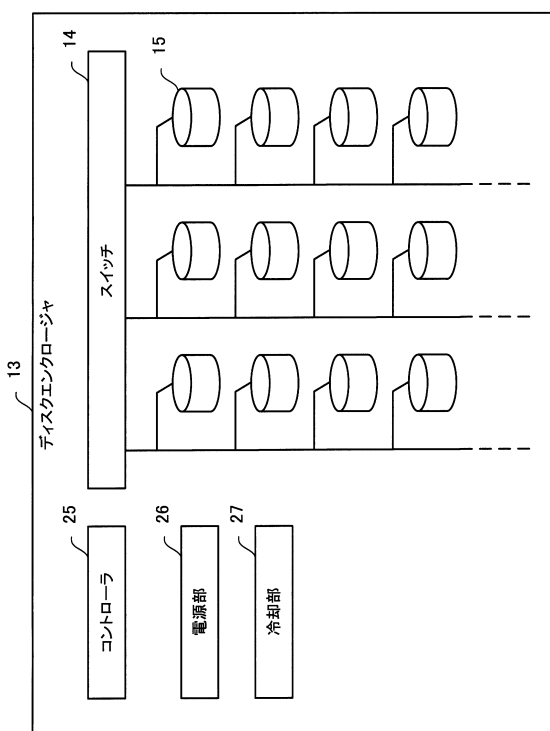
【 図 1 】



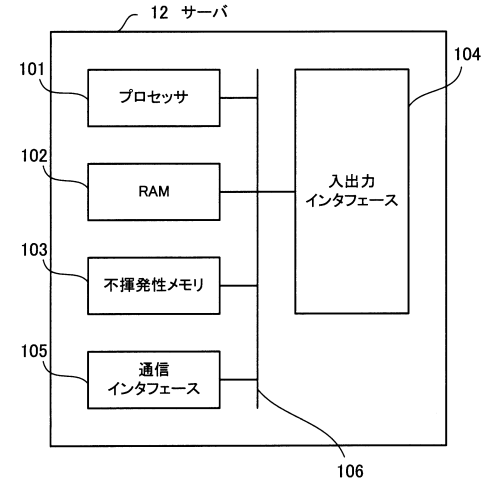
【 図 2 】



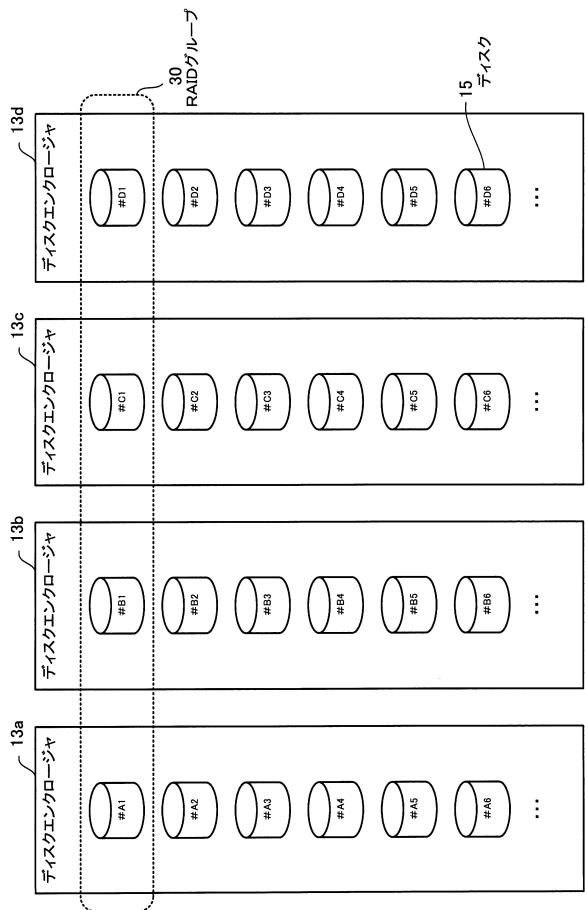
【 図 3 】



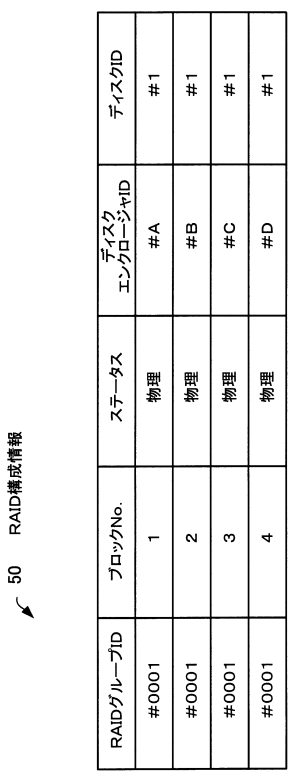
【 図 4 】



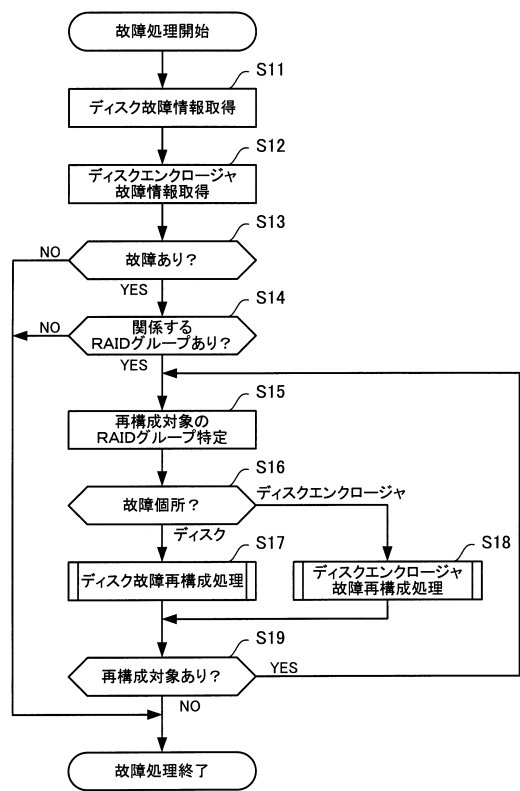
【図 5】



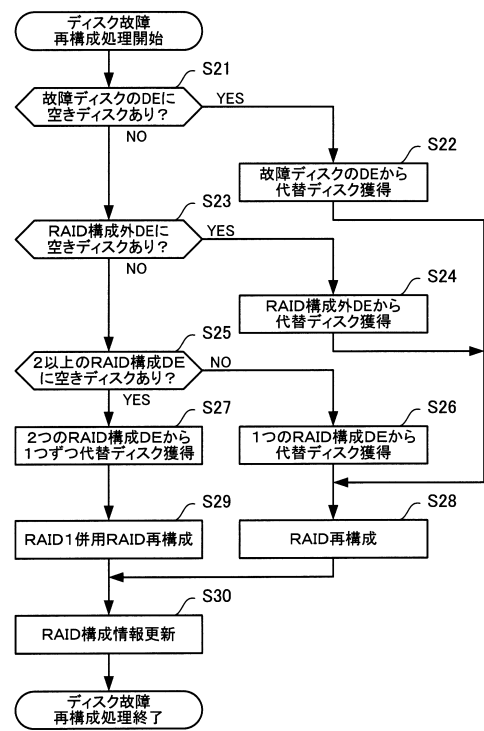
【図 6】



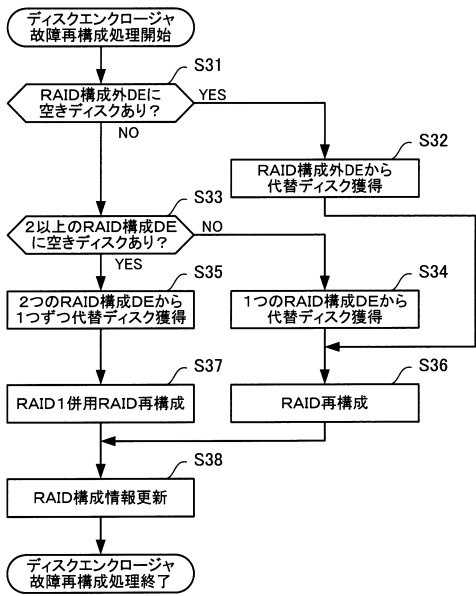
【図 7】



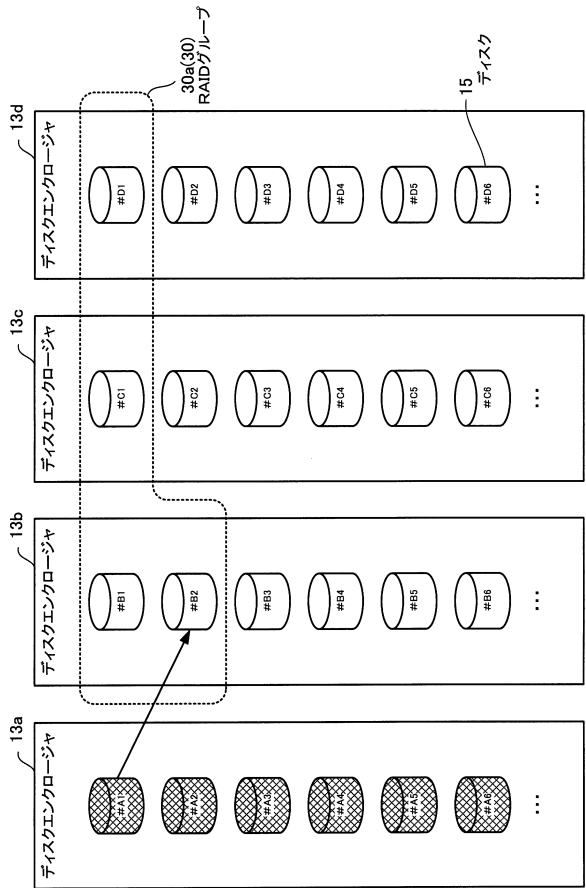
【図 8】



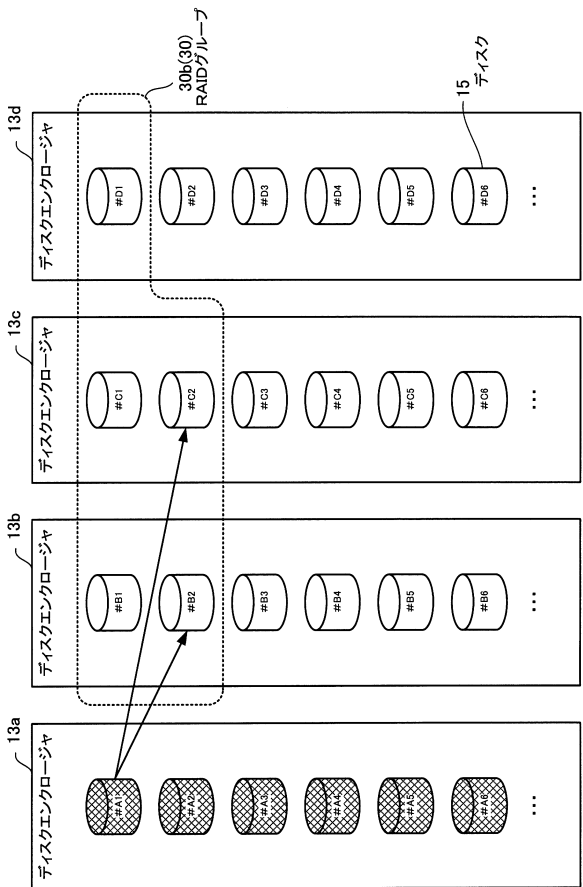
【図 9】



【図 10】



【図 11】

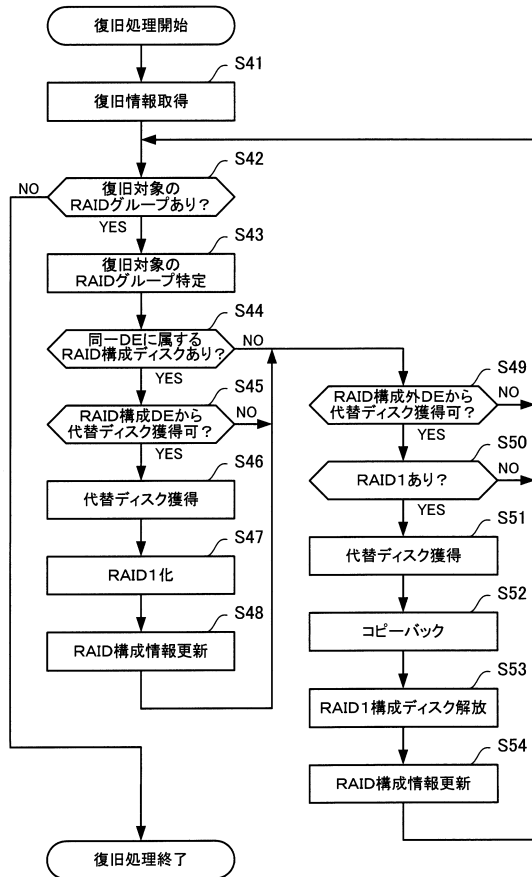


【図 12】

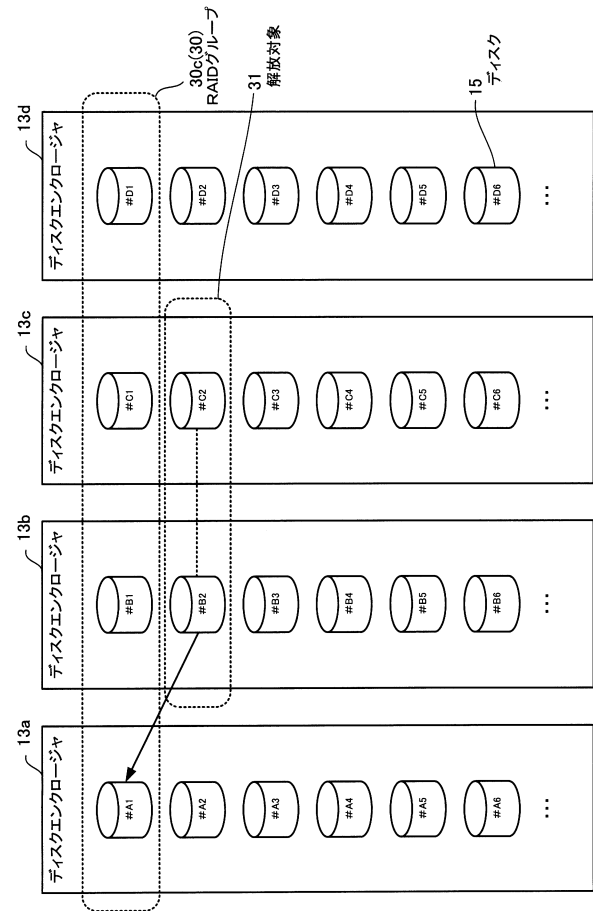
51 RAID構成情報

RAIDグループID	ブロックNo.	ステータス	ディスク エンクロージャID	ディスクID
#0001	1	RAID1	#B	#2
#0001	1	RAID1	#C	#2
#0001	2	物理	#B	#1
#0001	3	物理	#C	#1
#0001	4	物理	#D	#1

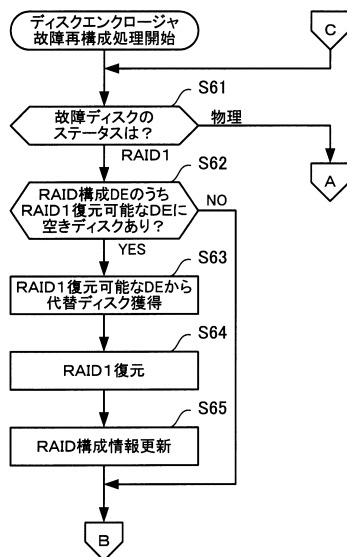
【図 13】



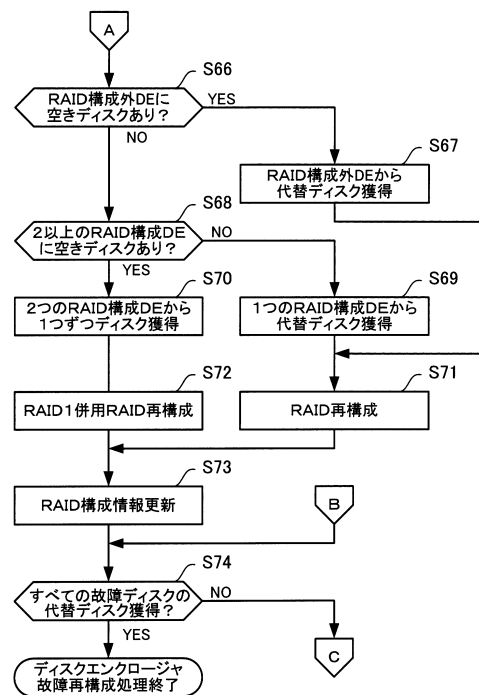
【図 14】



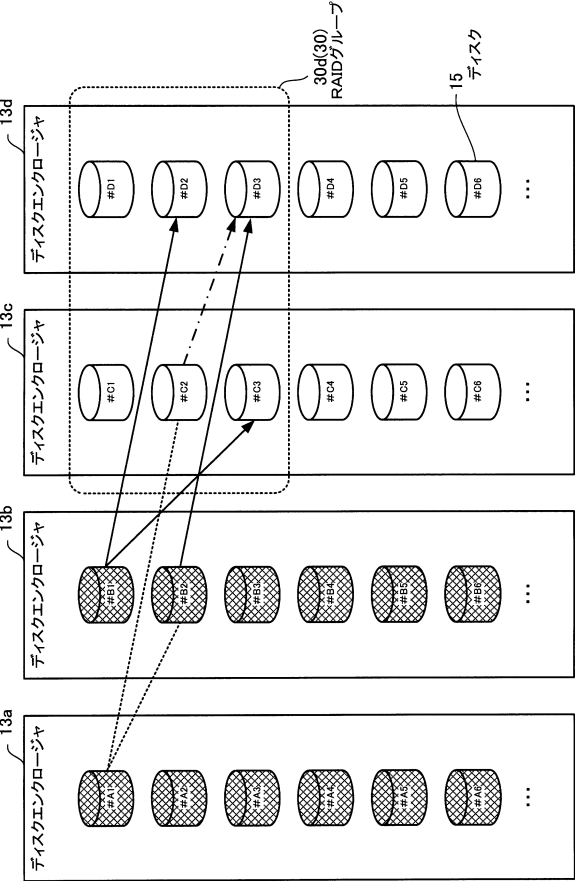
【図 15】



【図 16】



【図 17】

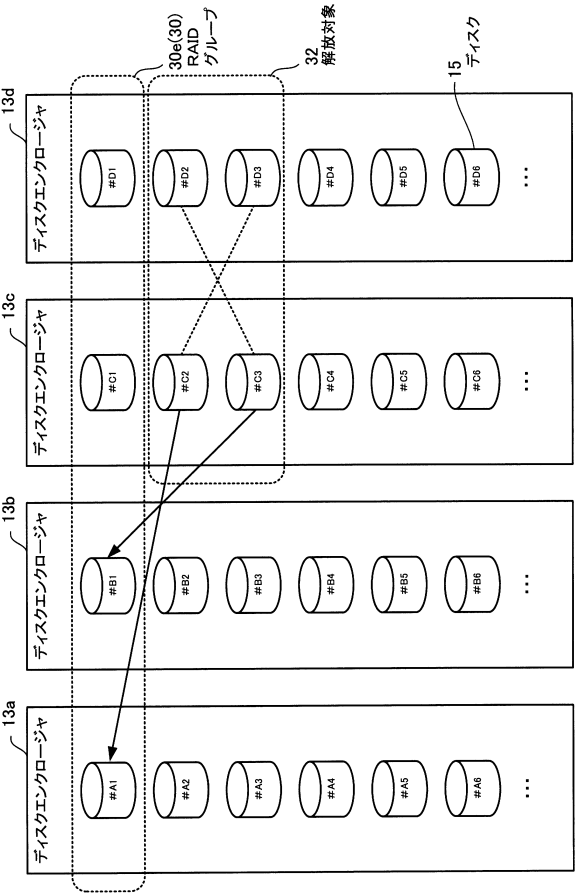


【図 18】

52 RAID構成情報

RAIDグループID	ブロックNo.	ステータス	ディスクエンクローザID	ディスクID
#0001	1	RAID1	#C	#2
#0001	1	RAID1	#D	#3
#0001	2	RAID1	#C	#3
#0001	2	RAID1	#D	#2
#0001	3	物理	#C	#1
#0001	4	物理	#D	#1

【図 19】



フロントページの続き

(56)参考文献 特開平06-230903(JP,A)
特開2005-293547(JP,A)
特開2009-187406(JP,A)
特開2009-252114(JP,A)
米国特許出願公開第2009/0024792(US,A1)
特開2010-256975(JP,A)
特開2014-056445(JP,A)

(58)調査した分野(Int.Cl., DB名)
G06F 3/06 - 3/08
11/16 - 11/20