

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2008-9767

(P2008-9767A)

(43) 公開日 平成20年1月17日(2008.1.17)

(51) Int. Cl.

G06F 3/06 (2006.01)

F I

G06F 3/06 304F

テーマコード(参考)

5B065

審査請求 未請求 請求項の数 18 O L (全 26 頁)

(21) 出願番号 特願2006-180256 (P2006-180256)
 (22) 出願日 平成18年6月29日(2006.6.29)

(71) 出願人 000005108
 株式会社日立製作所
 東京都千代田区丸の内一丁目6番6号
 (74) 代理人 100079108
 弁理士 稲葉 良幸
 (74) 代理人 100093861
 弁理士 大賀 真司
 (72) 発明者 新井 政弘
 神奈川県川崎市麻生区王禅寺1099番地
 株式会社日立製作所システム開発研究所
 内
 Fターム(参考) 5B065 BA01 CA30 CC03 EA12 EA23
 EA33

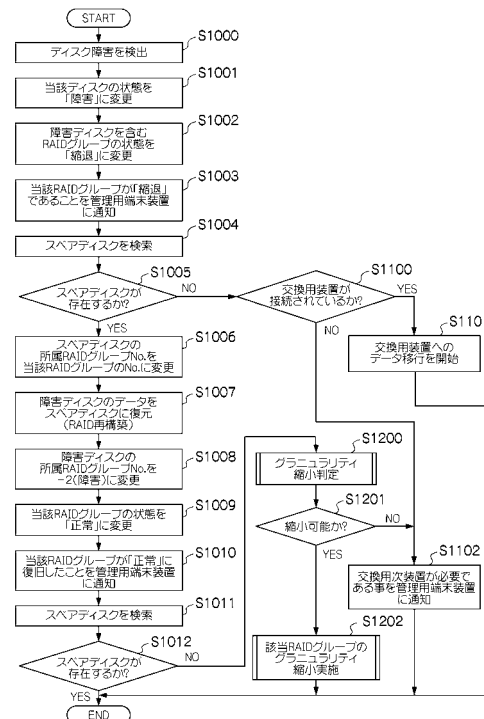
(54) 【発明の名称】 データ処理システム及びその方法並びにストレージ装置

(57) 【要約】

【課題】本発明は、多数のスペアディスクを搭載せずとも、ディスク非交換で耐障害性を維持し、小型で長期利用可能なデータ処理システム及びその方法並びにストレージ装置を提供することを目的とする。

【解決手段】上位装置と、前記上位装置に対してデータを読み書きするための第1の記憶領域を提供するストレージ装置とを有するデータ処理システムにおいて、前記ストレージ装置は、複数設けられた前記第1の記憶領域の一部領域に障害が発生した場合に、当該一部領域に格納された前記データを移行するためのスペアの第2の記憶領域が存在しなかったときに、前記複数の第1の記憶領域の他の一部領域を前記第2の記憶領域として動的に確保する確保部とを備えることとした。

【選択図】 図1



【特許請求の範囲】

【請求項 1】

上位装置と、前記上位装置に対してデータを読み書きするための第 1 の記憶領域を提供するストレージ装置とを有するデータ処理システムにおいて、

前記ストレージ装置は、

複数設けられた前記第 1 の記憶領域の一部領域に障害が発生した場合に、当該一部領域に格納された前記データを移行するためのスベアの第 2 の記憶領域が存在しなかったときに、前記複数の第 1 の記憶領域の他の一部領域を前記第 2 の記憶領域として動的に確保する確保部

を備えることを特徴とするデータ処理システム。

10

【請求項 2】

前記ストレージ装置は、

前記複数の第 1 の記憶領域の一部領域に格納された前記データを移行するための第 2 の記憶領域を検索する検索部

を備えることを特徴とする請求項 1 に記載のデータ処理システム。

【請求項 3】

前記複数の第 1 の記憶領域は、論理ボリュームを構成し、

前記確保部は、前記複数の第 1 の記憶領域のうち一の記憶領域に予め格納された前記データを、前記複数の第 1 の記憶領域のうち他の記憶領域に移行して前記論理ボリューム内に前記第 2 の記憶領域を確保する

ことを特徴とする請求項 1 に記載のデータ処理システム。

20

【請求項 4】

前記複数の第 1 の記憶領域は、パリティデータを格納する論理ボリュームを構成し、

前記確保部は、前記パリティデータを削除して前記論理ボリューム内に前記第 2 の記憶領域を確保する

ことを特徴とする請求項 1 に記載のデータ処理システム。

【請求項 5】

前記複数の第 1 の記憶領域は、複数の論理ボリュームを構成し、

前記確保部は、前記複数の論理ボリュームから前記データを格納していない未使用の論理ボリュームを前記第 2 の記憶領域として確保する

ことを特徴とする請求項 1 に記載のデータ処理システム。

30

【請求項 6】

前記複数の第 1 の記憶領域は、動的に拡張する仮想ボリュームを構成し、

前記上位装置からアクセスされた前記複数の第 1 の記憶領域のうち一の記憶領域は、動的に拡張するプールボリュームを構成し、

前記確保部は、前記プールボリュームを縮小して前記第 2 の記憶領域を確保する

ことを特徴とする請求項 1 に記載のデータ処理システム。

【請求項 7】

上位装置と、前記上位装置に対してデータを読み書きするための第 1 の記憶領域を提供するストレージ装置とを有するデータ処理システムのデータ処理方法において、

40

前記ストレージ装置には、複数設けられた前記第 1 の記憶領域の一部領域に障害が発生した場合に、当該一部領域に格納された前記データを移行するためのスベアの第 2 の記憶領域が存在しなかったときに、前記複数の第 1 の記憶領域の他の一部領域を前記第 2 の記憶領域として動的に確保する確保ステップ

を備えることを特徴とするデータ処理方法。

【請求項 8】

前記ストレージ装置には、前記複数の第 1 の記憶領域の一部領域に格納された前記データを移行するための第 2 の記憶領域を検索する検索ステップ

を備えることを特徴とする請求項 7 に記載のデータ処理方法。

【請求項 9】

50

前記複数の第 1 の記憶領域は、論理ボリュームを構成し、
前記確保ステップは、

前記複数の第 1 の記憶領域のうち一の記憶領域に予め格納された前記データを、前記複数の第 1 の記憶領域のうち他の記憶領域に移行して前記論理ボリューム内に前記第 2 の記憶領域を確保する、

ことを特徴とする請求項 7 に記載のデータ処理方法。

【請求項 10】

前記複数の第 1 の記憶領域は、パリティデータを格納する論理ボリュームを構成し、
前記確保部ステップは、

前記パリティデータを削除して前記論理ボリューム内に前記第 2 の記憶領域を確保する
ことを特徴とする請求項 7 に記載のデータ処理方法。 10

【請求項 11】

前記複数の第 1 の記憶領域は、複数の論理ボリュームを構成し、
前記確保ステップは、

前記複数の論理ボリュームから前記データを格納していない未使用の論理ボリュームを
前記第 2 の記憶領域として確保する

ことを特徴とする請求項 7 に記載のデータ処理方法。

【請求項 12】

前記複数の第 1 の記憶領域は、動的に拡張する仮想ボリュームを構成し、

前記上位装置からアクセスされた前記複数の記憶領域のうち一の記憶領域から動的に拡張する
プールボリュームを構成し、 20

前記確保ステップは、

前記プールボリュームを縮小して前記第 2 の記憶領域を確保する

ことを特徴とする請求項 7 に記載のデータ処理方法。

【請求項 13】

上位装置と、前記上位装置に対してデータを読み書きするための第 1 の記憶領域を提供する
ストレージ装置において、

複数設けられた前記第 1 の記憶領域の一部領域に障害が発生した場合に、当該一部領域に格納された前記データを移行するためのスペアの第 2 の記憶領域が存在しなかったときに、
前記複数の第 1 の記憶領域の他の一部領域を前記第 2 の記憶領域として動的に確保する
確保部 30

を備えることを特徴とするストレージ装置。

【請求項 14】

前記複数の第 1 の記憶領域の一部領域に格納された前記データを移行するための第 2 の記憶領域を検索する
検索部

を備えることを特徴とする請求項 13 に記載のストレージ装置。

【請求項 15】

前記複数の第 1 の記憶領域は、論理ボリュームを構成し、

前記確保部は、前記複数の第 1 の記憶領域のうち一の記憶領域に予め格納された前記データを、
前記複数の第 1 の記憶領域のうち他の記憶領域に移行して前記論理ボリューム内に前記第 2 の記憶領域を確保する 40

ことを特徴とする請求項 13 に記載のストレージ装置。

【請求項 16】

前記複数の第 1 の記憶領域は、パリティデータを格納する論理ボリュームを構成し、

前記確保部は、前記パリティデータを削除して前記論理ボリューム内に前記第 2 の記憶領域を確保する

ことを特徴とする請求項 13 に記載のストレージ装置。

【請求項 17】

前記複数の第 1 の記憶領域は、複数の論理ボリュームを構成し、

前記確保部は、前記複数の論理ボリュームから前記データを格納していない未使用の論 50

理ボリュームを前記第2の記憶領域として確保することを特徴とする請求項13に記載のストレージ装置。

【請求項18】

前記複数の第1の記憶領域は、動的に拡張する仮想ボリュームを構成し、前記上位装置からアクセスされた前記複数の第1の記憶領域のうち一の記憶領域は、動的に拡張するプールボリュームを構成し、前記確保部は、前記プールボリュームを縮小して前記第2の記憶領域を確保することを特徴とする請求項13に記載のストレージ装置。

【発明の詳細な説明】

【技術分野】

10

【0001】

本発明は、RAID (Redundant Array of Inexpensive Disks) を構成するデータ処理システム及びその方法並びにストレージ装置に関し、特にディスクドライブ (以下、単にディスクと呼ぶ) の一部、または、すべてを非交換に搭載するデータ処理システム及びその方法並びにストレージ装置に関する。

【背景技術】

【0002】

従来、ストレージ装置として、ディスクの障害に起因するデータ損失を防止する方策として RAID 構成を採用したものがあ (例えば、非特許文献1参照)。また、更なるディスク障害に備えるため、RAID 回復用にスペアのディスク (以下、スペアディスク) を用意し、迅速な復旧を実現している (例えば、特許文献1参照)。

20

【非特許文献1】David A.Patterson, Garth Gibson, and Randy H.Katz, "A Case for Redundant Arrays of Inexpensive Disks (RAID)", University of California Berke.

【特許文献1】特開平7-146760号公報

【発明の開示】

【発明が解決しようとする課題】

【0003】

前述のようなストレージ装置では、次のディスクの障害に備えるために、障害が発生したディスクは交換され、新たなディスクがスペアディスクとして補充されることが前提となっている。従い、ディスクを交換できないようなストレージ装置 (以下、ディスク非交換型ストレージ装置) では、新たなスペアディスクの補充ができず、耐障害性が低下してしまうという問題が生じる。

30

【0004】

一方、非交換を前提に多数のスペアディスクをあらかじめ搭載する方法も考えられるが、この方法では装置コストが高くなる上、装置自体も大きくなってしまいう問題が生じる。

【0005】

本発明は、多数のスペアディスクを搭載せずとも、ディスク非交換で耐障害性を維持し、小型で長期利用可能なデータ処理システム及びその方法並びにストレージ装置を提供することを目的とする。

40

【課題を解決するための手段】

【0006】

かかる課題を解決するため本発明は、上位装置と、前記上位装置に対してデータを読み書きするための第1の記憶領域を提供するストレージ装置とを有するデータ処理システムにおいて、前記ストレージ装置は、複数設けられた前記第1の記憶領域の一部領域に障害が発生した場合に、当該一部領域に格納された前記データを移行するためのスペアの第2の記憶領域が存在しなかったときに、前記複数の第1の記憶領域の他の一部領域を前記第2の記憶領域として動的に確保する確保部とを備えることを特徴とする。

【0007】

50

本発明のデータ処理システムによれば、既存のディスクの記憶領域からスペアの記憶領域を動的に確保することができるため、小型で長期利用可能なデータ処理システムの耐障害性を維持することができる。

【0008】

また、本発明は、上位装置と、前記上位装置に対してデータを読み書きするための第1の記憶領域を提供するストレージ装置とを有するデータ処理システムのデータ処理方法において、前記ストレージ装置には、複数設けられた前記第1の記憶領域の一部領域に障害が発生した場合に、当該一部領域に格納された前記データを移行するためのスペアの第2の記憶領域が存在しなかったときに、前記複数の第1の記憶領域の他の一部領域を前記第2の記憶領域として動的に確保する確保ステップとを備えることを特徴とする。

10

【0009】

本発明のデータ処理方法によれば、既存のディスクの記憶領域からスペアの記憶領域を動的に確保することができるため、小型で長期利用可能なデータ処理システムの耐障害性を維持することができる。

【0010】

さらに、本発明は、上位装置と、前記上位装置に対してデータを読み書きするための第1の記憶領域を提供するストレージ装置において、複数設けられた前記第1の記憶領域の一部領域に障害が発生した場合に、当該一部領域に格納された前記データを移行するためのスペアの第2の記憶領域が存在しなかったときに、前記複数の第1の記憶領域の他の一部領域を前記第2の記憶領域として動的に確保する確保部とを備えることを特徴とする。

20

【0011】

本発明のストレージ装置によれば、既存のディスクの記憶領域からスペアの記憶領域を動的に確保することができるため、小型で長期利用可能なストレージ装置の耐障害性を維持することができる。

【発明の効果】

【0012】

本発明によれば、多数のディスクを搭載せずとも、ディスクを交換することなく、スペアディスク又はスペア領域を確保することができ、小型で長期利用可能なデータ処理システム又はストレージ装置の耐障害性を維持することができる。

【0013】

また、耐障害性を維持できなくなる可能性が生じた場合（スペアディスク又はスペア領域を確保できなくなった場合）、通知により交換用ストレージ装置の準備の必要性を知ることができる。

30

【発明を実施するための最良の形態】

【0014】

以下、本発明の第1から第4の実施の形態を、図面を参照して説明する。

【0015】

1. 第1の実施の形態

図1は本発明の第1から第4の実施の形態におけるストレージ装置の構成を表した図である。ストレージ装置1は、管理用LAN5を介して管理用端末装置2と接続されている。また、ストレージ装置1は、ストレージネットワーク4を介して複数のホスト計算機3と接続されている。さらに、ストレージ装置1は交換用の他ストレージ装置との接続インターフェース7を備えている。

40

【0016】

ストレージ装置1には、ストレージコントローラ171、172が備えられており、ストレージネットワークとの間で同時に多くのデータを入出力できるようにしている。

【0017】

なお、実施の形態によっては、ストレージコントローラは1つでもよいし、3つ以上あってもよい。

【0018】

50

また、ストレージ装置 1 には、複数のディスク D 0、D 1、D 2・・・D N が備えられている。ストレージコントローラ 1 7 1、1 7 2 は接続インターフェース 6 を介して、ディスク D 0 ~ D N と接続されている。これにより、ストレージコントローラ 1 7 1、1 7 2 とディスク D 0 ~ D N は互いにデータの入出力を行うことができる。ストレージコントローラ 1 7 1、1 7 2 とディスク D 0 ~ D N との接続には、S A T A (Serial ATA) や S A S (Serial Attached SCSI)、ファイバーチャネル (Fibre Channel) 等、データ転送に適切な通信路が用いられる。

【0019】

ストレージコントローラ 1 7 1、1 7 2 は、制御プログラムが動作しており、ディスク D 0 ~ D N に対するデータの入出力を制御する。また、ストレージコントローラ 1 7 1、1 7 2 は、ディスク D 0 ~ D N によって構成される R A I D の構成を管理する。また、ストレージコントローラ 1 7 1、1 7 2 は、管理用端末装置 2 と通信を行い、種々のデータをやりとりする。

10

【0020】

ディスク D 0 ~ D N は、S A T A、S A S またはファイバーチャネル等で接続可能なディスクである。

【0021】

管理用端末装置 2 は、C P U、メモリ、記憶装置、インターフェース、入力装置および表示装置が備わるコンピュータ装置である。管理用端末装置 2 では管理プログラムが動作しており、当該管理プログラムによってストレージ装置 1 の動作状態を把握し、ストレージ装置 1 の動作を制御する。なお、管理用端末装置 2 では W e b ブラウザ等のクライアントプログラムが動作しており、ストレージ装置 1 から供給される管理プログラム (Common Gateway Interface や Java (登録商標) 等) によってストレージ装置 1 の動作を制御するようにしてもよい。

20

【0022】

表示画面 2 1 は管理用端末装置 2 に備わる表示装置に表示される画面である。

【0023】

ホスト計算機 3 は、C P U、メモリ、記憶装置およびインターフェースが備わるコンピュータ装置 (上位装置) であり、ストレージ装置 1 から供給されるデータを利用して、データベースサービスやウェブサービス等を利用可能にする。

30

【0024】

ストレージネットワーク 4 は、例えば S A S プロトコル、ファイバーチャネルプロトコルのような、データの転送に適するプロトコルで通信可能なネットワークである。

【0025】

管理用 L A N 5 は、例えば、T C P / I P (Transmission Control Protocol / Internet Protocol) によって、コンピュータ間でデータや制御情報を通信可能であり、例えばイーサネット (登録商標) が用いられる。

【0026】

図 2 は、本発明の第 1 の実施の形態のストレージ装置 1 の構成を示す概観図である。ストレージコントローラ 1 7 1、1 7 2、ディスク D 0 ~ D N、ファンおよび電源 8 は、ストレージ装置筐体 1 に内蔵されており、天板 1 1 で覆われている (図では、内部の構成を表現するために、天板 1 1 の一部を切り取って描画している)。ディスク D 0 ~ D N は、非交換に設けられている。

40

【0027】

また、前面は複数の吸気口を設けたベゼル 1 2 で覆われている。ベゼル 1 2 には、装置の起動・停止に必要な電源ボタン 1 3、1 4、および、障害時に鳴るブザーを止めるブザー停止ボタン 1 5 が備わっている。

【0028】

なお、図には示さないが、各ディスクのアクセス状況を示す、いくつかのアクセス L E D が前面に備わっていてもよい。

50

【0029】

図3は、本発明の第1の実施の形態におけるストレージコントローラ171の内部構成を示すブロック図である。なお、ストレージコントローラ172も同様である。

【0030】

ストレージコントローラ171には、CPU1901、メモリ1902、パリティジェネレータおよびデータ転送コントローラ1905、データバッファ1908、フロントエンド接続インターフェースコントローラ1906、バックエンド接続インターフェースコントローラ1907、LANインターフェースコントローラ1909が備えられている。これらは実装に適切なデータ転送路を介して、互いに接続されている。

10

【0031】

CPU1901は、メモリ1902に格納される制御プログラム1903、管理テーブル1904を読み込んでストレージ装置1のさまざまな制御を行う。

【0032】

パリティジェネレータおよびデータ転送コントローラ1905は、CPU1901の指示に基づいて、メモリ1902、フロントエンド接続インターフェースコントローラ1906、バックエンド接続インターフェースコントローラ1907、データバッファ1908との間で互いのデータの転送を行う。また、データバッファ1908を用いて、所定のデータに対するパリティ計算を実行する。

【0033】

フロントエンド接続インターフェース1906は、図1にて示した、ホスト計算機3との間でデータの入出力の制御を行い、必要に応じて、データバッファ1908にデータを格納したり、取り出したりする。

20

【0034】

バックエンド接続インターフェース1907は、ディスクD0等との間でデータの入出力を行い、必要に応じてデータバッファ1908にデータを格納したり、取り出したりする。

【0035】

データバッファ1908は、例えばDIMM(Dual Inline Memory Module)等で構成されるメモリであり、バッテリー等によって不揮発にされている。

30

【0036】

LANインターフェースコントローラ1909は、管理用端末装置2とデータや制御情報の入出力を行うインターフェースコントローラである。

【0037】

なお、第1の実施の形態におけるストレージコントローラ172は、上述で説明をしたストレージコントローラ171と同様の構成のため、説明を省略する。また、図3に示すストレージコントローラのブロック構成は一例であり、同様の機能を有せば、実装にハードウェア、ソフトウェアを問わない。また、ブロックの構成が異なってもよい。

【0038】

図4は、メモリ1902に格納される制御プログラム1903及び管理テーブル1904を示す説明図である。

40

【0039】

制御プログラム1903は、RAID制御プログラム1911、RAIDグループ設定プログラム1912、LU設定プログラム1913および管理・通知プログラム1914から構成される。

【0040】

管理テーブル1904は、ディスク情報管理テーブルT1、RAIDグループ情報管理テーブルT2、LU情報管理テーブルT3から構成される。

【0041】

そうしてCPU1901は、制御プログラム1903及び管理テーブル1904に基づ

50

いて、ホスト計算機 3 と入出力されたデータを処理し、ディスク D 0 等と入出力を行う。また、CPU 1901 は、同制御にかかわる障害復旧処理を行う。

【0042】

RAIDグループ設定プログラム 1912 は、管理用端末装置 2 からの指示や CPU 1901 の指示に基づき、RAIDグループ情報管理テーブル T 2 を用いて、RAIDグループの作成、変更または削除を行う。

【0043】

LUグループ設定プログラム 1913 は、管理用端末装置 2 からの指示や CPU 1901 の指示に基づき、LU情報管理テーブル T 3 を用いて、LUの作成、変更または削除を行う。

10

【0044】

管理・通知プログラム 1914 は、管理用端末装置 2 とデータや制御情報の授受を行う。

【0045】

ディスク情報管理テーブル T 1 は、ディスク D 0 等の各種情報を記録するテーブルである。RAIDグループ情報管理テーブル T 2 は、ディスク D 0 ~ D N から構成される RAIDグループの各種情報を記録するテーブルである。LU情報管理テーブル T 3 は、RAID上に構成される論理ディスクである LU (Logical Unit) の各種情報を記録するテーブルである。

【0046】

図 5 は以降の説明での理解を助けるために、ディスク、RAIDグループ、LUの関係を模式的に表した図である。RAIDグループは複数の物理的なディスク (図 1 の D 0 ~ D N) を用いて任意のレベルの RAID (例えば RAID - 5) を構成した際に、その RAID構成を一単位とするボリュームである。図 5 の例では、図 1 のディスク D 0 ~ D 4 をディスク 0、1、2、3、4 として用いて、RAIDが構成されており、そのボリューム (RAIDグループ) に RAIDグループ 0 のラベルが付されている。

20

【0047】

LUは、RAIDグループを複数に論理分割したボリュームである。各LUは、ホスト計算機からそれぞれディスクとして見えるので、論理ディスクとも呼ばれる。

【0048】

図 5 の例では、RAIDグループ 0 の一部を分割して LU 0 を構成している。

30

【0049】

図 6 には、データがどのように物理ディスクに格納されているかを示した説明図である。図 6 は、RAID - 5 で構成された RAIDグループ上に LU 0 のデータが格納されている様子を描いている。データはストライプブロックと呼ばれる 1 以上のセクタから成る構成単位ごとに格納される。d 1、d 2、・・・はストライプブロック単位に格納されたデータを示しており、p 1、p 2、・・・はストライプブロック横一列を単位とするストライプ列ごとに d 1 等から計算生成されるパリティを示している。

【0050】

RAID - 5 は、パリティブロックの位置は各ディスクに分散配置される。そのため、一列目ではディスク 4 に、二列目ではディスク 3 に、というふうにずらしながら格納される。

40

【0051】

図 7 は、本発明の第 1 の実施の形態におけるディスク情報管理テーブル T 1 の説明図である。ディスク情報管理テーブル T 1 には、ディスク D 0 ~ D N の各種情報が記録される。具体的には、ディスクの番号、ディスクの容量、ディスクの状態 (「正常」か「障害」か「未装着」か)、所属する RAIDグループの番号が記録される。例えば、ディスク番号 1 のディスクは、容量が 500GB であり、ディスクの状態は「正常」であり、所属する RAIDグループは 0 である。

【0052】

50

なお、所属するRAIDグループの番号は、0か正の整数ならば、当該RAIDグループの番号を意味し、-1の場合にはスペアディスクを意味し、-2の場合には障害中を意味し、-3の場合には未使用のディスクであることを意味する。

【0053】

ディスク情報管理テーブルT1の内容は、メモリ1902上に置かれると共に、管理端末装置2からも確認することができる。

【0054】

図8は、管理用端末装置2の表示画面21において、ディスク情報管理テーブルT1で管理されるディスク情報を表示した際のディスク管理画面V1の例である。ユーザーはディスク管理画面V1を通じて、ディスクの使用状況を随時知ることができる。

10

【0055】

図9は、本発明の第1の実施の形態におけるRAIDグループ情報管理テーブルT2の説明図である。RAIDグループ情報管理テーブルT2には、ディスクD0~DNを用いて構成されるRAIDグループの情報が記録される。具体的には、RAIDグループの番号、RAIDレベル、RAIDに含まれるデータディスク数(パリティ分を除いたディスク数)、RAIDの有効容量(実際にデータを書き込める容量)、当該RAIDグループの状態(「正常」か「縮退」か「破損」か)、LUへの割当済容量、未割当容量である。なお、ここでいう縮退とは、RAID回復が可能な障害をいい、具体的にはパリティディスク分台以下の障害をいう。また、破損とは、RAID回復が不可能な障害をいう。例えば、RAIDグループ0は、RAIDレベル5(RAID-5)で構築されており、データディスク数は4であり、有効容量は2000GB、状態は「正常」であり、LUへの割当済容量は900GB、未割当容量は1100GBである。

20

【0056】

RAIDグループ情報管理テーブルT1の内容は、メモリ1902上に置かれると共に、管理端末装置2からも確認することができる。

【0057】

図10は、本発明の第1の実施の形態におけるLU情報管理テーブルT3の説明図である。LU管理テーブルT3は各RAIDグループ上に構成されるLUの情報が記録される。具体的にはLUの番号、LUの割当容量、割当元のRAIDグループ、当該LUの割当開始Logical Block Address(以下、LBAという)、終了LBAである。例えば、LU番号0のLUは、容量が100GB、RAIDグループ0上に設けられており、開始LBAは00000000hであり、終了LBAは0000F000hである。LBAは、先に述べたストライプブロックに一定間隔で割り振られる一意のアドレスである。

30

【0058】

LU情報管理テーブルT3は、ディスク情報管理テーブルT1と同様に、メモリ1902上に置かれると共に、管理端末装置2からも確認することができる。

【0059】

図11は、ディスク障害が発生した際に、CPU1901が制御プログラム1903および管理テーブル1904に基づいて、障害復旧処理および次スペアディスク確保処理を示すフローチャートである。

40

【0060】

すなわち、CPU1901は、例えばデータの読み書き不可能な場合、電源が切れてしまった場合、データの読み出せない回数がある閾値を越えてしまった場合等のディスク障害を検出すると(S1000)、ディスク情報管理テーブルT1における当該ディスクの状態を「障害」に変更する(S1001)。また、CPU1901は、RAIDグループ情報管理テーブルT2において、当該ディスクの所属するRAIDグループの状態を「縮退」に変更する(S1002)。そして、CPU1901は、管理・通知プログラム1914を呼び出し、管理用端末装置2に、当該ディスクの障害および当該RAIDグループが「縮退」であることを通知する(S1003)。

【0061】

50

次に、CPU 1901は、ディスク管理情報テーブルT1を用いて、装置内のスペアディスクを検索する(S1004)。スペアディスクが存在した場合(S1005: YES)、CPU 1901は、ディスク情報管理テーブルT1における当該スペアディスクの所属RAIDグループの番号-1(スペアディスク)表示を、障害が発生しているRAIDグループの番号に変更し(S1006)、当該スペアディスク上に障害が発生したディスクの内容を復元する。すなわち、CPU 1901は、RAIDの再構築を行う(S1007)。

【0062】

再構築が完了すると、CPU 1901は、ディスク情報管理テーブルT1を用いて、障害ディスクの所属するRAIDグループの番号を-2(障害ディスク)に変更する(S1008)。また、CPU 1901は、RAIDグループ情報管理テーブルT2を用いて、当該RAIDグループの状態を「縮退」から「正常」に変更する(S1009)。

10

【0063】

次に、CPU 1901は、管理・通知プログラム1914を呼び出し、管理用端末装置2に、当該RAIDグループが正常に復旧したことを通知する(S1010)。

【0064】

続いて、CPU 1901は、既に上述で使用した1台のスペアディスクの他に、残りのスペアディスクを検索するために、2回目のスペアディスク検索を行う(S1011)。スペアディスクが存在した場合は(S1012: YES)、障害復旧処理を終了する。

【0065】

一方、ステップS1005の1回目のスペアディスク検索において、スペアディスクが存在しない場合には(S1005: NO)、CPU 1901は、外部に交換用の他ストレージ装置が接続されているかを確認する。交換用の他ストレージ装置が接続されている場合には(S1100: YES)、CPU 1901は、交換用の他ストレージ装置へのデータ移行を開始する。接続されていない場合には(S1100: NO)、CPU 1901は、管理・通知プログラム1914を呼び出し、管理用端末装置2に、交換用次装置が必要であることを通知し(S1102)、障害復旧処理を終了する。

20

【0066】

一方、ステップS1012の2回目のスペアディスク検索において、スペアディスクが存在しない場合には(S1012: NO)、CPU 1901は、ディスク情報管理テーブルT1およびRAIDグループ情報管理テーブルT2を用いてグラニュラリティ縮小可能なRAIDグループが存在するかを判定する(S1200)。

30

【0067】

グラニュラリティ縮小可能なRAIDグループが存在する場合(S1201: YES)、CPU 1901は、最初に見つかったRAIDグループもしくは、管理端末装置2によって優先度が高く設定されたRAIDグループのグラニュラリティ縮小を実施し(S1202)、障害復旧処理を終了する。

【0068】

グラニュラリティ縮小可能なRAIDグループが存在しない場合(S1201: NO)、CPU 1901は、管理・通知プログラム1914を呼び出し、管理端末装置2に、交換用次装置が必要であることを通知し(S1102)、障害復旧処理を終了する。

40

【0069】

図12は、本発明の第1の実施の形態において、CPU 1901が制御プログラム1903および管理テーブル1904に基づき、グラニュラリティ縮小判定を示すフローチャートである。グラニュラリティ縮小とは、RAID構成におけるデータディスクの台数を減少させることを目的に、ストライプブロックおよびパリティブロック内に格納されたデータの移動を行う処理である。第1の実施の形態においては、例えば、RAID-5にて4D+1Pの構成を3D+1Pに変更するようなケースを指す。以下、フローチャートを用いて説明していく。

【0070】

50

すなわち、CPU 1901は、当該RAIDグループに存在するストライプ列数 a を $a = (1 \text{ 台のディスク容量}) \div (\text{ストライプブロックサイズ})$ にて求める (S 2 0 0 0)。ここで、1台のディスク容量とは、当該RAIDグループを構成するディスク1台の物理容量を指しており、また、ストライプブロックサイズとは、CPU 1901があらかじめ保持する定数である。CPU 1901は、ストライプ列数 a は小数を含む可能性がある値であり、 a の小数点以下を切り上げ、整数とした値を変数 X に格納する (S 2 0 0 1)。

【0071】

次に、CPU 1901は、グラニュラリティ縮小後のストライプ列数 b を $b = (\text{現在の列数}) \times (\text{増加割合})$ で求める。すなわち、ストライプ列数 b は、 $b = (LU \text{ 割当済容量} / \text{ストライプブロックサイズ} / \text{データディスク数}) \times \{ (\text{データディスク数} + \text{パリティディスク数}) / (\text{データディスク数} + \text{パリティディスク数} - 1) \}$ で求める (S 2 0 0 2)。ここで、パリティディスク数とは、当該RAIDグループのRAIDレベルにより自動的に定まる値であり、例えばRAID-5ならば1である。ストライプ列数 b は小数を含む可能性がある値であり、 b の小数点以下を切り捨て、整数とした値を Y に格納する (S 2 0 0 3)。

【0072】

次にCPU 1901は、 X と Y を比較し、 $Y \geq X$ であった場合には (S 2 0 0 4 : YES)、グラニュラリティ縮小可能と判定する (S 2 0 0 5)。一方、 $Y < X$ であった場合には (S 2 0 0 4 : NO)、CPU 1901は、グラニュラリティ縮小不可能と判定する (S 2 0 0 7)。

【0073】

次に、CPU 1901は、全RAIDグループを調査したかを判定し (S 2 0 0 6)、全て調査を終えている場合は (S 2 0 0 6 : YES)、終了する。終わっていない場合は (S 2 0 0 6 : NO)、CPU 1901は、別のRAIDグループに対して上記判定を繰り返す。

【0074】

図13は、本発明の第1の実施の形態において、CPU 1901が制御プログラム1903および管理テーブル1904に基づいてグラニュラリティ縮小を行う際の手順を示すフローチャートである。

【0075】

CPU 1901は、当該RAIDグループを構成するディスクよりストライプ列1列目のデータブロックをデータバッファ1908にリードする (S 3 0 0 0)。例えば、図17に示すように、CPU 1901はデータブロック $d_1 \sim d_4$ までをデータバッファ1908にリードする。続いて、CPU 1901は、次のストライプ列のデータブロックをデータバッファ1908にリードする (S 3 0 0 1)。例えば、図18に示すように、CPU 1901はデータブロック $d_5 \sim d_8$ までをデータバッファ1908にリードする。そして、CPU 1901は、データバッファ1908に蓄えられたデータのうち、(データディスク数 - 1) このデータブロック群を用いて、新パリティ $p'x$ を生成する (S 3 0 0 2)。例えば、図19に示すように、CPU 1901はデータディスク数4 - 1のデータブロック群 $d_1 \sim d_3$ のうち新パリティ $p'1$ を生成する。次に、(データディスク数 - 1) 個のデータブロック群と新パリティ $p'x$ で構成される新ストライプ列をディスクにライトする (S 3 0 0 3)。

【0076】

現在リードしてあるストライプ列がLU割当済領域の最終ストライプであった場合には (S 3 0 0 4 : YES)、CPU 1901は、データバッファ1908にある残りのデータブロックから新パリティ $p'x$ を生成し (S 3 0 0 5)、ディスクにライトする (S 3 0 0 6)。

【0077】

次に、CPU 1901は、RAIDグループ情報管理テーブルT2を用いて、当該RAIDグループのデータディスク数を1つ減少させる (S 3 0 0 7)。また、CPU 190

10

20

30

40

50

1 は、ディスク情報管理テーブル T 1 を用いて、上記ストライプ列の書き換えで未使用となったディスクの所属 R A I D グループ番号を - 1 (スペアディスク) に変更し (S 3 0 0 8)、処理を終える。

【 0 0 7 8 】

一方、ステップ S 3 0 0 4 において、最終ストライプ列でなかった場合には (S 3 0 0 4 : N O)、C P U 1 9 0 1 は、ステップ S 3 0 0 1 に戻り、処理を繰り返す。

【 0 0 7 9 】

図 1 4 ~ 図 1 6 は、前記図 1 3 で説明したグラニュラリティ縮小処理時の、ディスク上のデータブロック d およびパリティブロック p の配置状態を示した図である。

【 0 0 8 0 】

図 1 4 は、グラニュラリティ縮小処理前のデータ配置を示している。データブロック d およびパリティブロック p は、ディスク 0 ~ ディスク 4 上に配置されているのがわかる。

【 0 0 8 1 】

図 1 5 は、C P U 1 9 0 1 がステップ S 3 0 0 0 ~ S 3 0 0 3 を実行することにより、1 列目が新ストライプ列に書き換えられた状態を示している。1 列目においてデータブロック d 1、d 2、d 3 はディスク 0、1、2 に配置され、パリティブロック p ' 1 はディスク 3 に配置されており、そして、ディスク 4 のブロックは「空き」になっている。

【 0 0 8 2 】

図 1 6 は、グラニュラリティ縮小処理が終了した際のデータブロック d およびパリティブロック p の配置状態を示した図である。データブロック d 1、d 2、d 3 およびパリティブロック p ' 1 は、ディスク 0 ~ ディスク 3 上にのみ配置され、ディスク 4 は「空き」、すなわち未使用となっていることがわかる。そして、2 列目においても、データブロック d 3、d 4、d 5 およびパリティブロック p ' 2 は、ディスク 0 ~ ディスク 3 上にのみ配置され、ディスク 4 は「空き」となっていることがわかる。3 列目以降もディスク 4 は「空き」が形成されることがわかる。これにより C P U 1 9 0 1 が、ステップ S 3 0 0 7 および S 3 0 0 8 の処理を実行することで、ディスク 4 を新たなスペアディスクにすることができる。

【 0 0 8 3 】

図 1 7 ~ 図 1 9 は、ステップ S 3 0 0 0 ~ S 3 0 0 2 におけるデータバッファ 1 9 0 8 を表した図である。

【 0 0 8 4 】

図 1 7 は、ステップ S 3 0 0 0 にて、ストライプ列一列目のデータブロック d 1、d 2、d 3、d 4 が読み込まれた状態を表している。

【 0 0 8 5 】

図 1 8 は、ステップ S 3 0 0 1 にて、次のストライプ列のデータブロック d 5、d 6、d 7、d 8 が読み込まれた状態を表している。

【 0 0 8 6 】

図 1 9 は、(データディスク数 - 1) 個のデータブロック d 1、d 2、d 3 から新パリティ p ' 1 が生成された状態を表している。

【 0 0 8 7 】

以上のように本実施の形態によるデータ処理システムでは、スペアディスクが存在しない場合に、グラニュラリティ縮小を行うことで既存のディスクからスペアディスクを確保することができる。これにより、ストレージ装置の耐障害性を長期的に維持することができる。また、既存のディスクからスペアディスクを確保するので、ストレージ装置の小型化が実現できる。

【 0 0 8 8 】

一方、耐障害性を維持できなくなる可能性が生じた場合 (スペアディスクを確保できなくなった場合) には、通知により交換用ストレージ装置の準備の必要性を知ることができる。

【 0 0 8 9 】

10

20

30

40

50

2. 第2の実施の形態

第2の実施の形態として、RAIDレベルを変更することにより、新たなスペアディスクを確保する方法がある。例えば、RAID-6 (nD+2P)をRAID-5 (nD+1P)に変更し、この差で生じる1ディスクをスペアディスクとする方法である。本実施の形態について図20～図23を用いて説明する。

【0090】

図20は、本実施の形態において、ディスク障害が発生した際に、CPU1901が制御プログラム1903および管理テーブル1904に基づいて、障害復旧処理および次スペアディスク確保処理を示すフローチャートである。

【0091】

ステップS1000～S1012、およびステップS1100～S1102については、先に示した第1の実施の形態と同じなので説明を割愛する。

【0092】

そして、ステップS1012の2回目のスペアディスク検索においてスペアディスクが存在しない場合 (S1012:NO)、CPU1901は、RAIDレベルを変更可能なRAIDグループ存在するか、RAID情報管理テーブルT2を用いて検索し、存在した場合には (S4100:YES)、該当RAIDグループのRAIDレベル変更を行う (S4101)。なお、該当RAIDグループが複数存在する場合には、最初に見つかったRAIDグループか、管理用端末装置2にて、対象優先度が高く設定されているRAIDグループに対して処理を行う。

【0093】

図21は、本実施の形態におけるCPU1901が制御プログラム1903および管理テーブル1904に基づいて行うRAIDレベル変更の手順を示すフローチャートである。

【0094】

CPU1901は、ストライプ列1列目のデータブロックdおよびパリティブロックp、qをデータバッファ1908にリードすると (S5000)、パリティブロックqをデータバッファ1908より破棄する。例えば、図22に示すように、CPU1901は、ディスク6～ディスク10に格納されたストライプ列1列目のデータブロックd1、d2、d3、パリティブロックp1、q1から、図23に示すように、パリティブロックq1を破棄する。次に、CPU1901は、データブロックdとパリティブロックpから構成される新ストライプ列をディスクにライトする (S5002)。現在データバッファ1908にリードされたストライプ列がLU割当済領域の最終ストライプ列であった場合には (S5003:YES)、CPU1901は、残りのデータブロックdとパリティブロックpをディスクにライトする (S5004)。

【0095】

後に、CPU1901は、RAIDグループ情報管理テーブルT2を用いて、当該RAIDグループのRAIDレベルの記載を変更する (S5005)。次に、CPU1901は、ディスク情報管理テーブルT1を用いて、未使用となったディスクの所属RAIDグループ番号を-1 (スペアディスク)に変更し (S5006)、処理を終える。

【0096】

一方、CPU1901は、ステップS5003において、現在データバッファ1908にリードされたストライプ列がLU割当済領域の最終ストライプ列でなかった場合には (S5003:NO)、次のストライプ列のデータブロックdおよびパリティブロックp、qをデータバッファ1908にリードし (S5007)、S5001に戻って処理を繰り返す。

【0097】

図22、図23は、本実施の形態による、ディスク上のデータブロックdおよびパリティブロックp、qの変化を示す図である。

【0098】

10

20

30

40

50

図 2 2 は R A I D レベル変更を行う前の R A I D グループの状態である。R A I D グループは R A I D - 6 のデータ構成をとっており、データブロック d と 2 種のパリティブロック p、q がディスク 6 ~ 1 0 上に分散配置されている。

【 0 0 9 9 】

図 2 3 は、図 1 9、図 2 1 に示した R A I D レベル変更の処理により、R A I D - 5 のデータ構成に変更されたあとのデータ配置を示している。R A I D - 6 (2 パリティ) から R A I D - 5 に変更されたため、パリティブロック q が破棄され、データブロック d とパリティブロック p のみで構成されている。また、データブロック d およびパリティブロック p はディスク 6 ~ ディスク 9 上に配置されており、ディスク 1 0 が「空き」すなわち未使用となっている。従い、C P U 1 9 0 1 が前述のステップ S 5 0 0 5、S 5 0 0 6 を実行することにより、ディスク 1 0 がスペアディスクに設定される。

【 0 1 0 0 】

以上のように本実施の形態によるデータ処理システムでは、スペアディスクが存在しない場合に、R A I D のレベル変更を行うことでより安定した性能が実現できるとともに、既存のディスクからスペアディスクを確保することができる。これにより、ストレージ装置の耐障害性を長期的に維持することができる。また、既存のディスクからスペアディスクを確保するので、ストレージ装置の小型化が実現できる。

【 0 1 0 1 】

一方、耐障害性を維持できなくなる可能性が生じた場合 (スペアディスクを確保できなくなった場合) には、通知により交換用ストレージ装置の準備の必要性を知ることができる。

【 0 1 0 2 】

3 . 第 3 の実施の形態

第 3 の実施の形態として、スペアの領域をディスクとして確保せず L U として確保する方法がある。第 1 及び第 2 の実施の形態では物理ディスク上に障害ディスクの内容を復元できるようにしたが、本実施の形態では、L U、すなわち論理ディスク上に障害ディスクの内容を復元できるようにする。L U として配されるスペア領域 (以下、スペア L U) を確保する方法は、第 1 及び第 2 の実施の形態の様に、既存のデータブロックを移動させずとも容易にスペア領域を確保できる点で優れている。一方、障害ディスク上にあった元の L U のアドレス (L B A) をスペア L U 上のアドレスに変換してアクセスするため、性能は、第 1 及び第 2 の実施の形態に比べ、低下する恐れがある。以下、図 2 4 ~ 図 2 7 を用いて説明する。

【 0 1 0 3 】

図 2 4 は、第 3 の実施の形態におけるメモリ 1 9 0 2 に格納される制御プログラム 1 9 0 3 および管理テーブル 1 9 0 4 を示す説明図である。

【 0 1 0 4 】

第 3 の実施の形態における制御プログラム 1 9 0 3 は、R A I D 制御プログラム 1 9 1 1、R A I D グループ設定プログラム 1 9 1 2、L U 設定プログラム 1 9 1 3、管理・通知プログラム 1 9 1 4 から構成される。

【 0 1 0 5 】

また管理テーブル 1 9 0 4 は、ディスク情報管理テーブル T 1、R A I D グループ情報管理テーブル T 2、L U 情報管理テーブル T 3、スペア L U 管理テーブル T 4、L B A - L B A 変換テーブル T 5 から構成される。

【 0 1 0 6 】

図 2 5 はスペア L U 管理テーブル T 4 の説明図である。スペア L U 管理テーブル T 4 は、L U として配されるスペア領域の情報を記録するテーブルである。具体的には、スペア L U 管理テーブル T 4 は、スペア L U の番号、容量、割当元 R A I D グループの番号、開始 L B A、終了 L B A、スペアの使用状態 (「使用済」か「未使用」か) を記録する。例えば、スペア L U 番号 0 のスペア L U は容量が 5 0 0 G B であり、割当元 R A I D グループ番号は 0 であり、開始 L B A は 0 0 0 0 0 0 0 0 h であり、終了 L B A は 0 0 0 4 F 0

10

20

30

40

50

00hであり、使用状態は「未使用」である。

【0107】

なお、スベアLU管理テーブルT4は、LU情報管理テーブルT3と共にメモリ1902におかれる。また、管理用端末装置2を通じて、その情報を確認することができる。

【0108】

図26は、LBA LBA変換テーブルT5の説明図である。LBA LBA変換テーブルT5は、元のLUのアドレス(LBA)とスベアLU上のアドレス(LBA)との対応を管理するテーブルである。具体的には、LBA-LBA変換テーブルT5は、アクセス元LUの番号、アクセス元のLBA、スベアLUの番号及びアクセス先のLBAから構成され、スベアLUのLBAを管理する。例えば、アクセス元LU1のLBA00356780hのデータに対するアクセスは、スベアLU0のLBA10004780hに対するアクセスに変換される。

10

【0109】

図27は、本実施の形態において、ディスク障害が発生した際に、CPU1901が制御プログラム1903および管理テーブル1904に基づいて行う、障害復旧処理およびスベアLUの確保処理を示すフローチャートである。

【0110】

ステップS1000~S1012、およびステップS1100~S1102については、先に示した実施の形態と同じなので説明を割愛する。

【0111】

ステップS1005の1回目のスベアディスク検索においてスベアディスクが存在しない場合(S1005:NO)、CPU1901はスベアLU管理テーブルT4を用いて、未使用のスベアLUを検索する(S6100)。未使用スベアLUが存在した場合(S6101:YES)、CPU1901は、障害ディスクのデータをスベアLUに復元し、RAID回復を行う(S6102)。

20

【0112】

次にCPU1901は、当該スベアLUのLBA-LBA変換テーブルT5を更新生成する(S6103)。そしてCPU1901は、スベアLU管理テーブルT4を用いて、当該スベアLUの使用状態を「使用済」に変更し(S6104)、ステップS1008以降の処理を行う。

30

【0113】

一方、ステップS1012の2回目のスベアディスク検索において、スベアディスクが存在しない場合(S1012:NO)、CPU1901は、ディスク情報管理テーブルT1およびRAID情報管理テーブルT2を用いて、ディスク1台分の容量以上の未使用容量を持つ、RAIDグループを検索する(S6200)。このとき、ディスク障害が発生したディスクが所属するRAIDグループに未使用容量のLUが存在していても、基本的には当該未使用容量のLUをスベアディスクとして使用せず、他のRAIDグループのLUを検索して使用する。但し、当該未使用容量のLUをスベアディスクとして使用する必要がある場合には、ディスク1台分のLUに加え、ディスク障害が発生したディスクの所属するRAIDグループを縮退してもなおも使用できる未使用容量が存在する場合に限り、未使用容量のLUをスベアディスクとして割り当てる。

40

【0114】

該当RAIDグループが存在した場合(S6201:YES)、CPU1901は、RAID情報管理テーブルT2を用いて、該当RAIDグループのLU割当済容量と未使用容量を変更し、スベアLUの容量確保を行う。具体的には、CPU1901は、LU割当済容量をディスク1台分の容量だけ増加させ、未使用容量をディスク1台分の容量だけ減少させる(S6202)。例えば、CPU1901は、スベアLUの容量として100GBの確保が必要な場合には、RAID情報管理テーブルT2のLU割当済容量を100GBだけ増加させ、未割当容量から100GBだけ減少させる。

【0115】

50

次に、CPU 1901は、スペアLU管理テーブルT4にステップS6202で確保したLUを追加登録し(S6203)、処理を終える。

【0116】

以上のように本実施の形態によるデータ処理システムでは、スペアディスクが存在しない場合に、ディスク1台分以上の未使用容量を有するRAIDグループを検索することで既存のディスクからスペアLU(スペア領域)を確保することができる。これにより、ストレージ装置の耐障害性を長期的に維持することができる。また、既存のディスクからスペア領域を確保するので、ストレージ装置の小型化が実現できる。

【0117】

さらに、耐障害性を維持できなくなる可能性が生じた場合(スペアディスクを確保できなくなった場合)、通知により交換用ストレージ装置の準備の必要性を知ることができる。

【0118】

4. 第4の実施の形態

第4の実施の形態として、容量を動的に割り当てるストレージ装置において適用する形態を説明する。

【0119】

容量を動的に割り当てるストレージ装置とは、複数のRAIDグループの容量をひとつの容量プールとして管理し、実際に書き込まれる部分にだけ動的にエクステントと呼ばれる、小容量同サイズのLUを割り当てる制御を行うストレージ装置である。ストレージ装置において、スペア領域は未使用エクステントの集合として形成される。第3の実施の形態において、スペアLUを確保するためには、ディスク1台の容量以上の未使用容量が単一のRAIDグループに必要であった。本実施の形態は、スペア領域をエクステントプール(エクステントの集合体)から確保するため、単一のRAIDグループである必要がなく、また不連続な領域からも確保できる。以下、図28~図35を用いて、本実施の形態について説明する。

【0120】

図28は、第4の実施の形態におけるメモリ1902に格納される制御プログラム1903および管理テーブル1904を示す説明図である。

【0121】

第4の実施の形態における制御プログラム1903は、RAID制御プログラム1911、RAIDグループ設定プログラム1912、LU設定プログラム1913および管理・通知プログラム1914から構成される。

【0122】

また管理テーブル1904は、ディスク情報管理テーブルT1、RAIDグループ情報管理テーブルT2、LU情報管理テーブルT3、エクステント管理テーブルT6、エクステント数管理テーブルT7、仮想LU用アクセス先エクステント対応テーブルT8およびスペアLU用エクステント対応テーブルT9から構成される。

【0123】

図29は、本実施の形態の概念を模式的に示した図である。本実施の形態のストレージ装置1は、LUの代わりに仮想的な容量を持った仮想LUをホスト計算機3に対して提供する。ホスト計算機3が仮想LUに書き込みを行うと、CPU1901は、容量プールから書き込まれたアドレスに対し、エクステントのディスパッチ(割当て)を行う。エクステントとは、RAIDグループを分割してできる小容量同サイズのLUである。エクステントの集合体は、RAIDグループの隔たりなく、1つのエクステントプールとして管理され、必要な分だけディスパッチされる。

【0124】

図30は、エクステント管理テーブルT6の説明図である。エクステント管理テーブルT6は、エクステントプールの管理実態であり、エクステントの各種情報が記録される。具体的には、エクステント管理テーブルT6は、エクステントの番号、割当元のRAID

10

20

30

40

50

グループの番号、開始 L B A、割当状態（「割当済」か「未割当」か）である。例えば、エクステント番号「0」は R A I Dグループ番号「0」上にあり、開始 L B Aは 0 0 0 0 0 0 0 0 hであり、割当状態は「割当済」である。なお、エクステントのサイズは予め決まっているので終了 L B Aの指定アドレス表示は必要ない。

【0125】

なお、エクステント管理テーブル T 6 は R A I D 情報管理テーブル T 2 と共に、メモリ 1 9 0 2 内に配される。

【0126】

図 3 1 は、エクステント数管理テーブル T 7 の説明図である。エクステント数管理テーブル T 7 は、エクステントの仮想 L U 領域への割当数、スペア領域への割当数、未使用数、総数を記録している。例えば、図 3 0 の場合、仮想 L U 領域への割当数は、200,000,000,000個であり、スペア領域への割当数は0個であり、未使用数は、700,000,000,000個であり、総数は、900,000,000,000である。

10

【0127】

なお、エクステント数管理テーブル T 7 は、エクステント管理テーブル T 6 と同様にメモリ 1 9 0 2 内に配される。

【0128】

図 3 2 は、仮想 L U に対する仮想 L U 用アクセス先エクステント対応テーブル T 8 である。仮想 L U 用アクセス先エクステント対応テーブル T 8 は、仮想 L U 上の L B A と ディスパッチされたエクステントとの対応関係を記録するテーブルである。具体的には、仮想 L U 用アクセス先エクステント対応テーブル T 8 は、仮想 L U の番号と仮想 L U の開始 L B A、エクステントの番号で構成される。例えば、仮想 L U 0 の開始 L B A 0 0 3 5 6 7 8 0 h に ディスパッチされたエクステント番号は 0 である。

20

【0129】

仮想 L U 用アクセス先エクステント対応テーブル T 8 は、エクステント管理テーブル T 6 と同様に、メモリ 1 9 0 2 内に配される。

【0130】

図 3 3 は、スペア L U に対するスペア L U 用アクセス先エクステント対応テーブル T 9 である。スペア L U 用アクセス先エクステント対応テーブル T 9 は、スペア L U の番号、スペア L U の開始 L B A、エクステントの番号で構成される。例えば、スペア L U 0 の開始 L B A 0 0 3 5 6 7 8 0 h に ディスパッチされたエクステント番号は 1 8 0 0 である。

30

【0131】

図 3 4 は、エクステントプールを容量表示した容量プールの管理画面である。容量プール管理画面 V 2 は、管理用端末装置 2 からアクセスすることができる。

【0132】

容量プール管理画面 V 2 は、仮想 L U に割り当てたエクステントと未使用のエクステントを容量表示しており、仮想 L U ごとの使用率を容量で示している。例えば、仮想 L U 0 は 5 0 0 G B 消費している。また未使用エクステントの総容量は 1 2 0 0 G B である。容量プール画面 V 2 はまた、スペア領域を確保により、未使用エクステント（未使用容量）が減少した際には、その旨を通知するメッセージウィンドウ 1 6 を備えている。

40

【0133】

図 3 5 は、C P U 1 9 0 1 が制御プログラム 1 9 0 3 および管理テーブル 1 9 0 4 に基づいて、エクステントを用いてスペア L U 用の領域を確保する手順を示すフローチャートである。なお、ステップ S 1 0 0 0 ~ S 1 0 1 2、ステップ S 6 1 0 0 ~ S 6 1 0 4 およびステップ S 1 1 0 0 ~ S 1 1 0 2 の処理は、上述と同様なので説明は割愛する。

【0134】

C P U 1 9 0 1 は、1 台分のディスク容量と（未割当数 * エクステントサイズ）とを比較し、未割当エクステントの集合体である未割当エクステント群からスペア L U を確保できるか否かを判別する（S 7 0 0 0）。ここでエクステントサイズは、C P U 1 9 0 1 においてあらかじめ定められた定数である。

50

【0135】

1台分のディスク容量が(未割当数*エクステントサイズ)と同じか、それより小さい場合、すなわち、1台分のディスク容量(未割当数*エクステントサイズ)が成り立つ場合には(S7000:YES)、CPU1901は、エクステント数管理テーブルT7を用いて、未割当数を(1台分のディスク容量÷エクステントサイズ)分だけ減少させ、スペア領域割当数を(1台分のディスク容量÷エクステントサイズ)分だけ増加させ、エクステントを割当済に変更し(S7001)、エクステントカウンタを変更する(S7002)。すなわち、CPU1901は、エクステント管理テーブルT6の割当状態を割当済に変更し、エクステント数管理テーブルT7の未使用数を減らし、スペア領域割当数を増やす。

10

【0136】

次に、CPU1901は、スペアLU用アクセス先エクステント対応テーブルT9を追加更新する(S7003)。CPU1901は、管理・通知プログラム1914を呼び出して、管理用端末装置2に、容量プールが減少した旨を通知し(S7004)、容量プール管理画面V2にこの旨が表示され、この処理が終了する。

【0137】

一方、ステップS7000において、1台分のディスク容量に見合う未使用領域が存在しない場合、すなわち、1台分のディスク容量>(未割当数*エクステントサイズ)となる場合(S7000:NO)、CPU1901は、管理・通知プログラム1914を呼び出して、管理用端末装置2に、交換用次装置が必要である旨を通知し(S7005)、この処理が終了する。

20

【0138】

以上のように本実施の形態によるデータ処理システムでは、スペアディスクが存在しない場合に、エクステントプールから動的にスペア領域を確保することができる。これにより、ストレージ装置の耐障害性を長期的に維持することができる。また、既存のディスクからスペア領域を確保するので、ストレージ装置の小型化が実現できる。

【0139】

さらに、耐障害性を維持できなくなる可能性が生じた場合(スペアディスクを確保できなくなった場合)、通知により交換用ストレージ装置の準備の必要性を知ることができる。

30

【0140】

5. 他の実施の形態

上述の第1から第4の実施の形態においては、ステップS1000~S1012までを検索部とし、複数設けられた第1の記憶領域をデータブロック及びパリティブロックとして、複数設けられたデータブロック及びパリティブロックの一部領域に障害が発生した場合に当該一部領域に格納されたデータを、第2の記憶領域であるスペアディスク又はスペア領域に移行して動的に記憶領域を確保する確保部を備えた場合について述べてきたが、本発明はこれに限らず、当該一部領域に格納されるべきデータをスペアディスク又はスペア領域に格納するようにしても良い。

【産業上の利用可能性】

40

【0141】

本発明は、1又は複数のストレージ装置や、1又は複数のストレージ装置を有するデータ処理システムに広く適用することができる。

【図面の簡単な説明】

【0142】

【図1】本発明におけるストレージ装置の構成を表した図である。

【図2】本発明の第1の実施の形態のストレージ装置の構成を示す概観図である。

【図3】本発明の第1の実施の形態におけるストレージコントローラの内部構成を示すブロック図である。

【図4】本発明の第1の実施の形態におけるメモリに格納される制御プログラムおよび管

50

理テーブルを示す説明図である。

【図 5】本発明の第 1 の実施の形態におけるディスク、RAID グループ、LU の関係を模式的に表した図である。

【図 6】本発明の第 1 の実施の形態における物理ディスクの格納状態を示す説明図である。

【図 7】本発明の第 1 の実施の形態におけるディスク情報管理テーブルの説明図である。

【図 8】本発明の第 1 の実施の形態における管理用端末装置の表示画面の画面例である。

【図 9】本発明の第 1 の実施の形態における RAID グループ情報管理テーブルの説明図である。

【図 10】本発明の第 1 の実施の形態における LU 情報管理テーブルの説明図である。

10

【図 11】本発明の第 1 の実施の形態におけるディスク障害が発生した際の、障害復旧処理および次スペアディスク確保処理を示すフローチャートである。

【図 12】本発明の第 1 の実施の形態におけるグラニュラリティ縮小判定を示すフローチャートである。

【図 13】本発明の第 1 の実施の形態におけるグラニュラリティ縮小処理を行う際の手順を示すフローチャートである。

【図 14】本発明の第 1 の実施の形態におけるグラニュラリティ縮小処理前のデータ配置を示した図である。

【図 15】本発明の第 1 の実施の形態におけるグラニュラリティ縮小処理により、1 列目が新ストライプ列に書き換えられた状態を示した図である。

20

【図 16】本発明の第 1 の実施の形態におけるグラニュラリティ縮小処理が終了した際の、データブロックおよびパリティブロックの配置状態を示した図である。

【図 17】本発明の第 1 の実施の形態におけるグラニュラリティ縮小処理にて、ストライプ列一列目のデータブロックがデータバッファに読み込まれた状態を表した図である。

【図 18】本発明の第 1 の実施の形態におけるグラニュラリティ縮小処理にて、次のストライプ列のデータブロックがデータバッファに読み込まれた状態を表した図である。

【図 19】本発明の第 1 の実施の形態におけるデータバッファ上で新パリティが生成された状態図である。

【図 20】第 2 の実施の形態において、ディスク障害が発生した際の、障害復旧処理および次スペアディスク確保処理を示すフローチャートである。

30

【図 21】第 2 の実施の形態における RAID レベル変更処理の手順を示すフローチャートである。

【図 22】第 2 の実施の形態において、RAID レベル変更処理を行う前の RAID グループの状態である。

【図 23】第 2 の実施の形態において、RAID レベル変更処理により、RAID - 5 のデータ構成に変更されたあとのデータ配置を示している。

【図 24】第 3 の実施の形態におけるメモリに格納される制御プログラムおよび管理テーブルを示す説明図である。

【図 25】第 3 の実施の形態における、スペア LU 管理テーブルの説明図である。

【図 26】第 3 の実施の形態における、LBA LBA 変換テーブルの説明図である。

40

【図 27】第 3 の実施の形態において、ディスク障害が発生した際の、障害復旧処理およびスペア LU の確保処理を示すフローチャートである。

【図 28】第 3 の実施の形態におけるメモリに格納される制御プログラムおよび管理テーブルを示す説明図である。

【図 29】第 4 の実施の形態の概念を模式的に示した図である。

【図 30】第 4 の実施の形態における、エクステント管理テーブルの説明図である。

【図 31】第 4 の実施の形態における、エクステント数管理テーブルの説明図である。

【図 32】第 4 の実施の形態における、仮想 LU 用アクセス先エクステント対応テーブルの説明図である。

【図 33】第 4 の実施の形態における、スペア LU 用アクセス先エクステント対応テーブ

50

ルの説明図である。

【図34】第4の実施の形態における、エクステントプールを容量表示した容量プール管理画面である。

【図35】第4の実施の形態における、ディスク障害が発生した際の、障害復旧処理およびスペアLUの確保処理を示すフローチャートである。

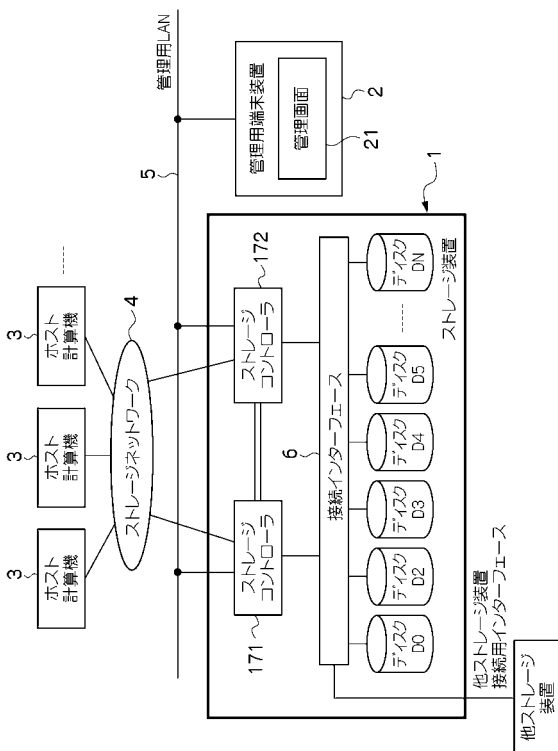
【符号の説明】

【0143】

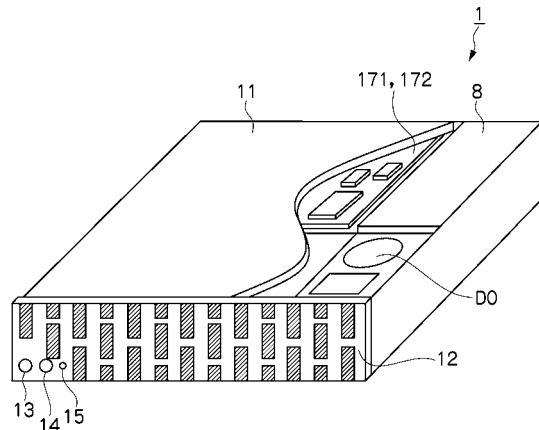
1 ストレージ装置、2 管理用端末装置、3 ホスト計算機、4 ストレージネットワーク、5 管理用LAN、6 接続インターフェース、7 交換用ストレージの接続インターフェース、8 ファンおよび電源、11 天板、12 ベゼル、13, 14 電源ボタン、15 ブザー停止ボタン、171, 172 ストレージコントローラ、1901 CPU、1902 メモリ、1905 パリティジェネレータ及びデータ転送コントローラ、1906 フロントエンド接続インターフェースコントローラ、1907 バックエンド接続インターフェースコントローラ、1908 データバッファ、1909 LANインターフェースコントローラ、V1 ディスク管理画面、V2 容量プール管理画面。

10

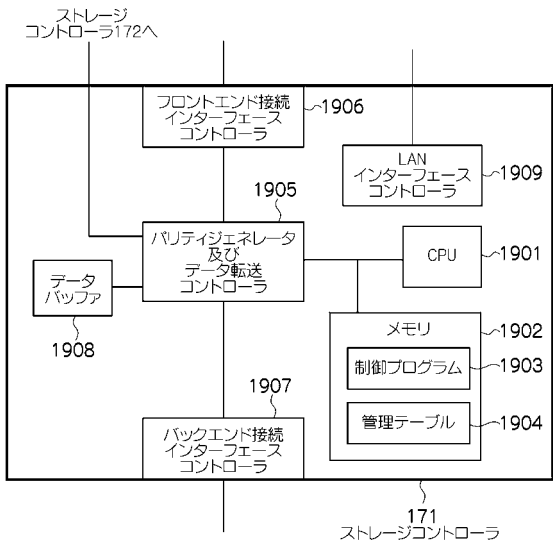
【図1】



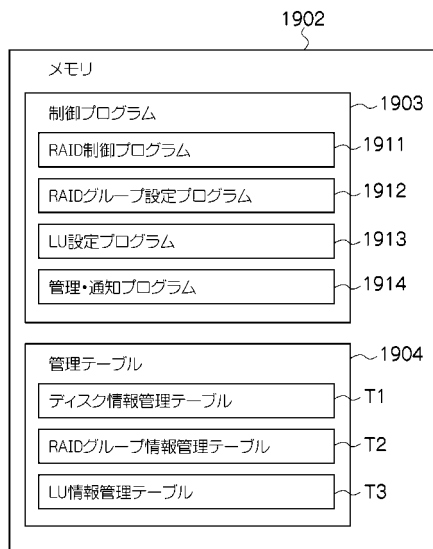
【図2】



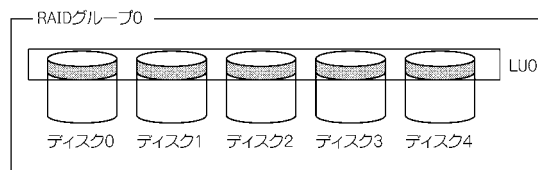
【 図 3 】



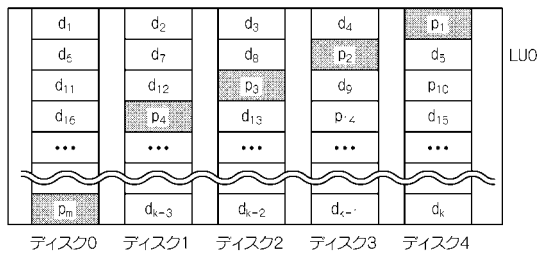
【 図 4 】



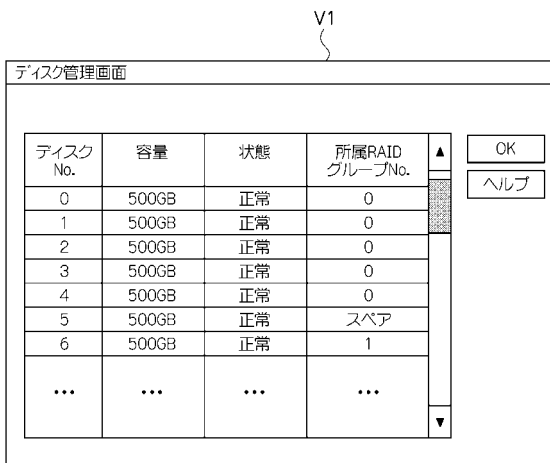
【 図 5 】



【 図 6 】



【 図 8 】



【 図 7 】

T1

ディスクNo.	容量	状態	所属RAID グループNo.
0	500GB	正常	0
1	500GB	正常	0
2	500GB	正常	0
3	500GB	正常	0
4	500GB	正常	0
5	500GB	正常	-1(スベア)
6	500GB	正常	1
...
N	-	未装着	-

【 図 9 】

T2

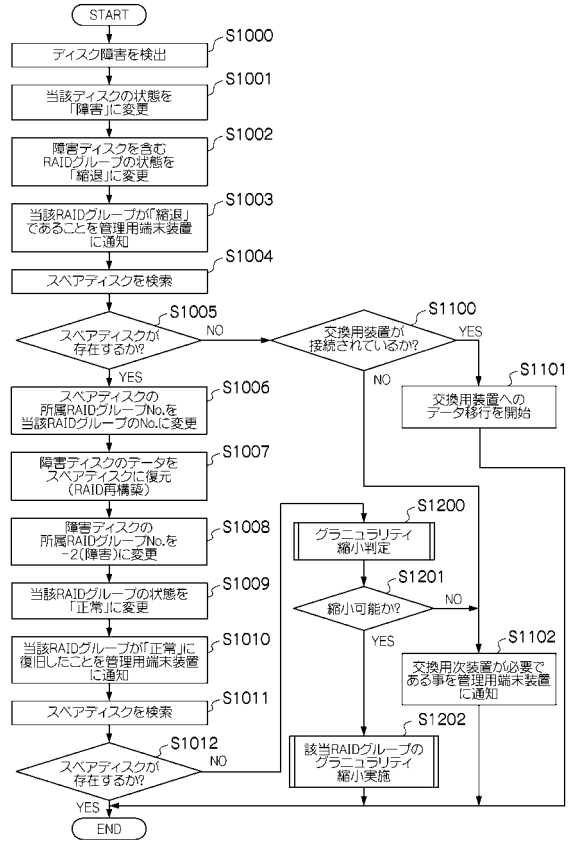
RAID グループ No.	RAID レベル	データ ディスク 数	有効 容量	状態	LU割当済 容量	未割当 容量
0	RAID-5	4	2000MB	正常	900GB	1100GB
1	RAID-6	3	1500MB	正常	1000GB	500GB
...

【 図 1 0 】

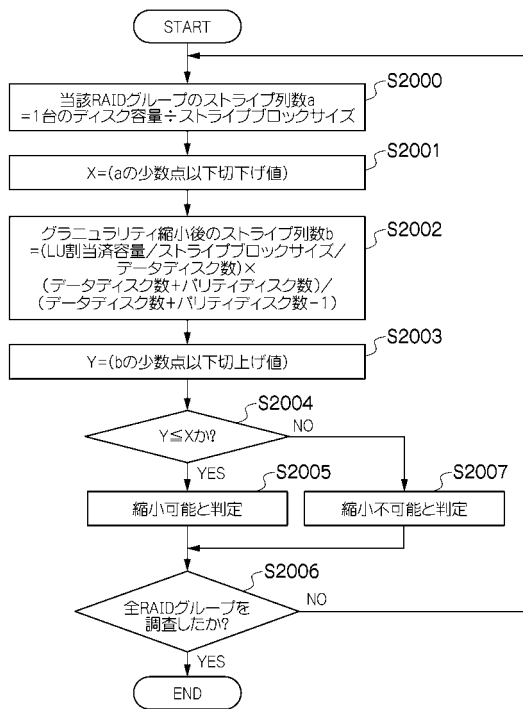
LU No.	容量	RAID グループ No.	開始LBA	終了LBA
0	100GB	0	00000000h	0000F000h
1	300GB	0	00010000h	0003F000h
2	500GB	0	00040000h	0008F000h
...

T3

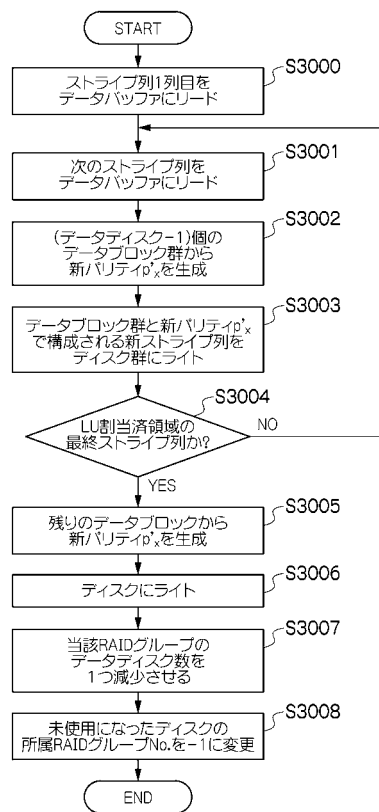
【 図 1 1 】



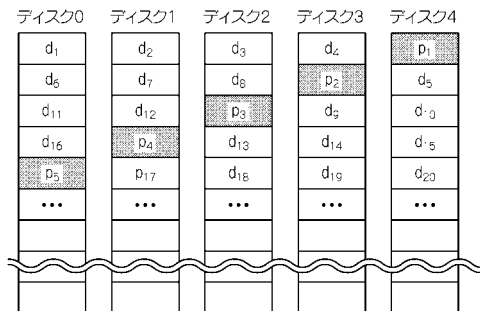
【 図 1 2 】



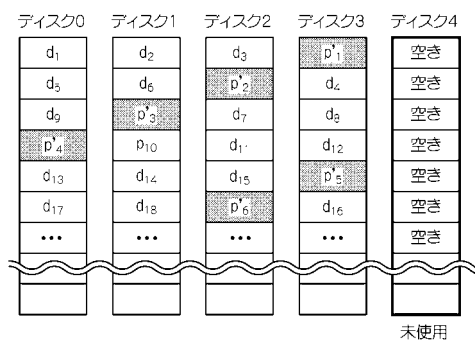
【 図 1 3 】



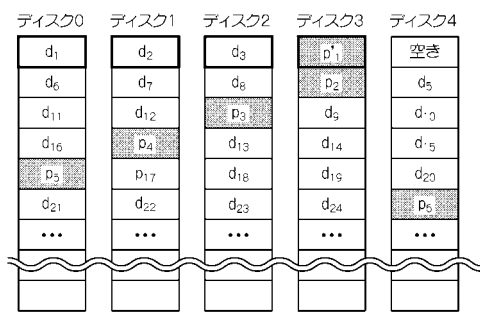
【図14】



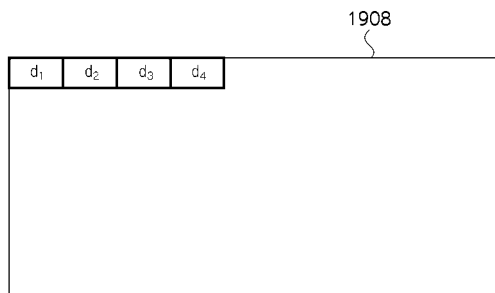
【図16】



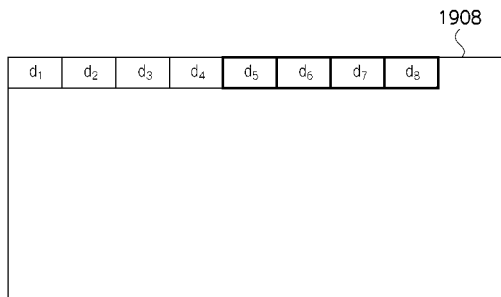
【図15】



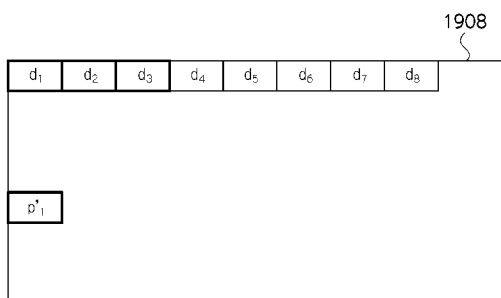
【図17】



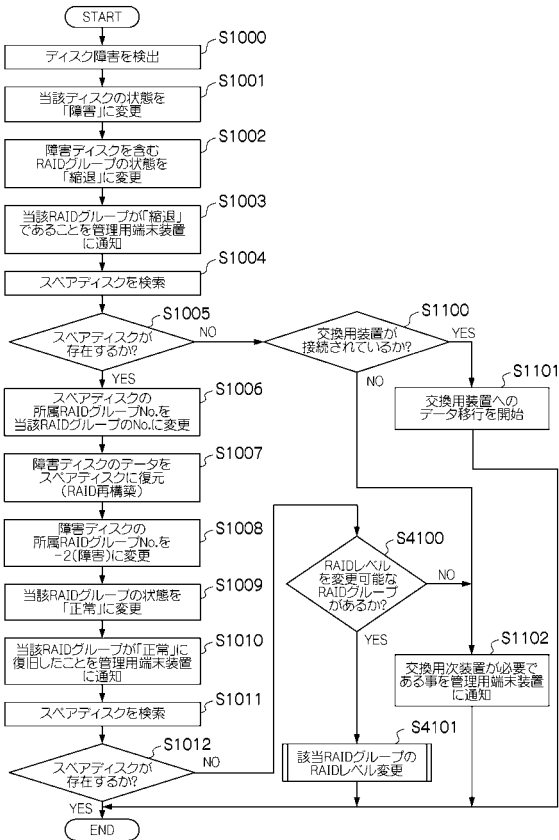
【図18】



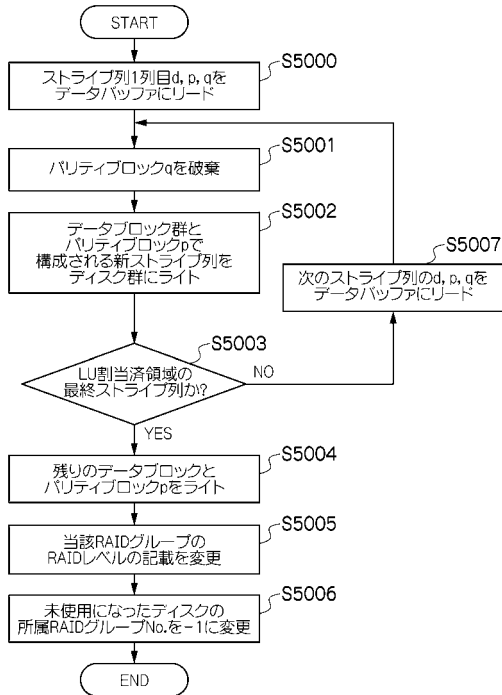
【図19】



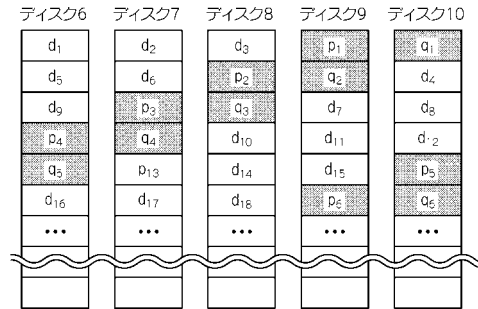
【図20】



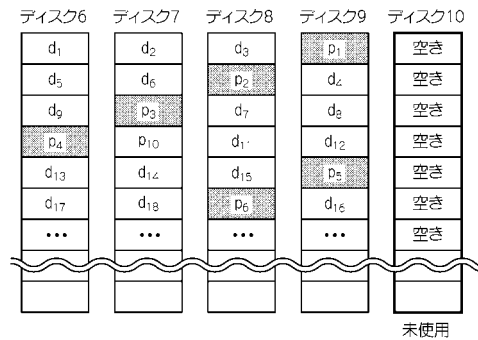
【図 2 1】



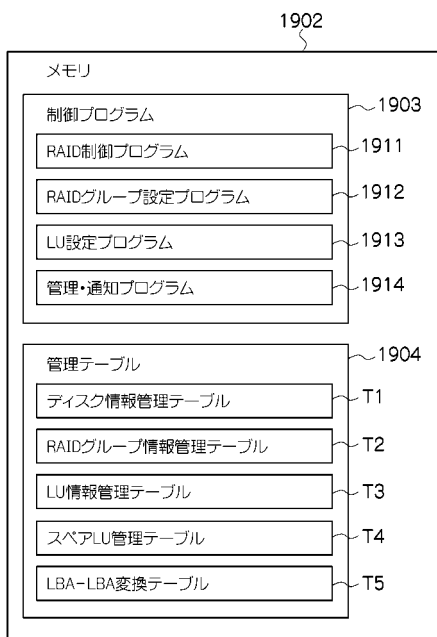
【図 2 2】



【図 2 3】



【図 2 4】



【図 2 5】

T4

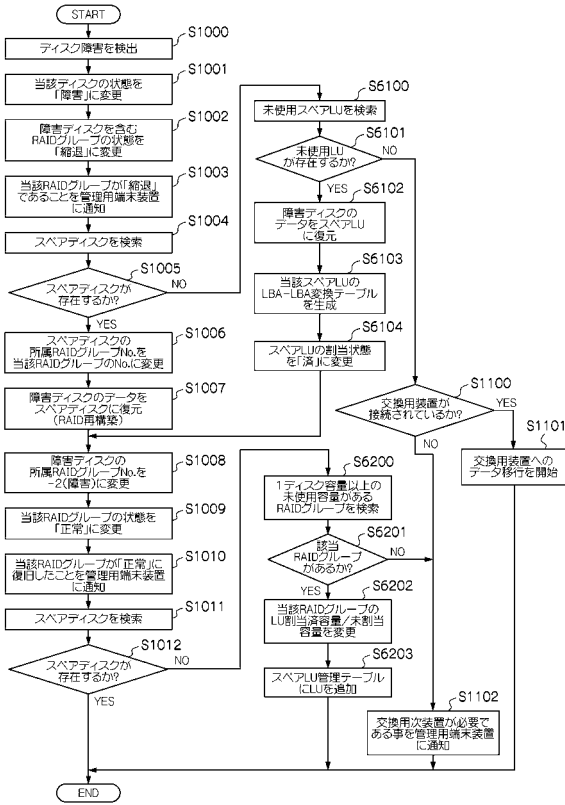
スベア LU No.	容量	RAID グループ No.	開始LBA	終了LBA	使用状態
0	500GB	0	00000000h	0004F000h	未使用
...

【図 2 6】

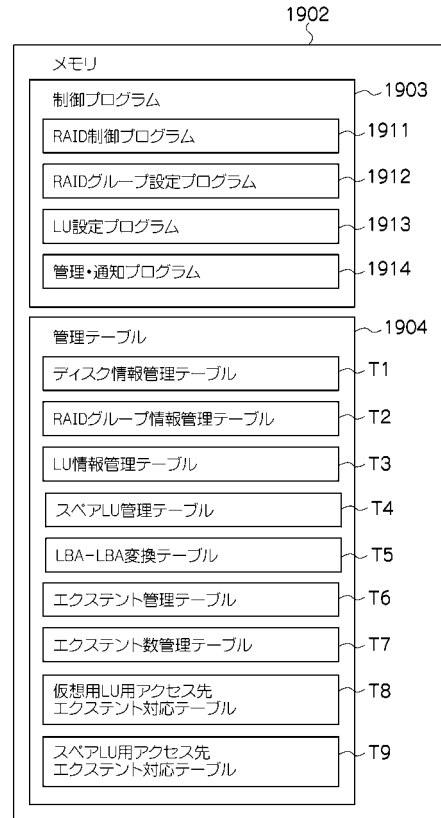
T5

アクセス元 LU No.	アクセス元LBA	スベア LU No.	アクセス先LBA
...
1	00356780h	1	10004780h
1	00356790h	1	10004790h
...

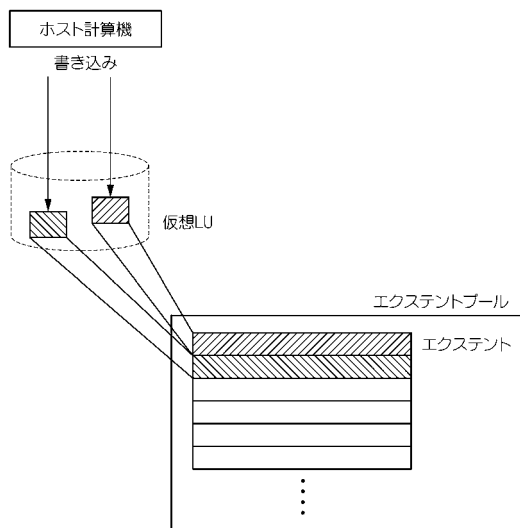
【図 27】



【図 28】



【図 29】



【図 31】

T7

ユーザー領域割当数	200,000,000,000
スペア領域割当数	0
未使用数	700,000,000,000
総数	900,000,000,000

【図 32】

T8

仮想LU No.	LBA	エクステント No.
...
0	00356780h	0
0	02466790h	1
...

【図 30】

T6

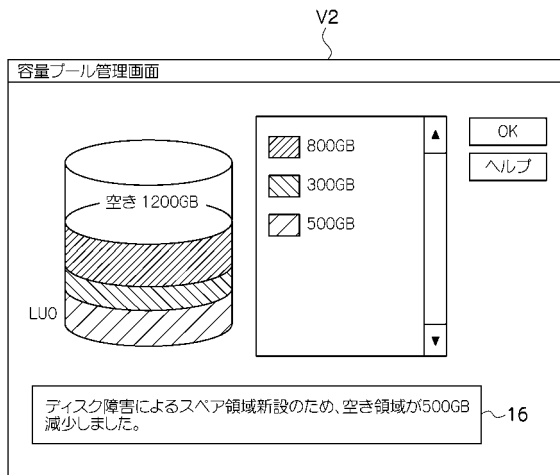
エクステント No.	RAID グループ No.	開始LBA	割当状態
0	0	00000000h	割当済
1	0	00000400h	割当済
2	0	00000800h	未割当
...

【図 33】

T9

スペアLU No.	開始LBA	エクステント No.
...
0	00356780h	1800
0	00356790h	8310
...

【図34】



【図35】

