

(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2016-161970

(P2016-161970A)

(43) 公開日 平成28年9月5日(2016.9.5)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 3/06 (2006.01)	G06F 3/06 306Z	
	G06F 3/06 540	
	G06F 3/06 305C	
	G06F 3/06 302Z	
	G06F 3/06 301J	

審査請求 未請求 請求項の数 8 O L (全 30 頁)

(21) 出願番号 特願2015-37096 (P2015-37096)
 (22) 出願日 平成27年2月26日 (2015.2.26)

(71) 出願人 000005223
 富士通株式会社
 神奈川県川崎市中原区上小田中4丁目1番1号
 (74) 代理人 100092978
 弁理士 真田 有
 (74) 代理人 100112678
 弁理士 山本 雅久
 (72) 発明者 塩沢 賢輔
 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

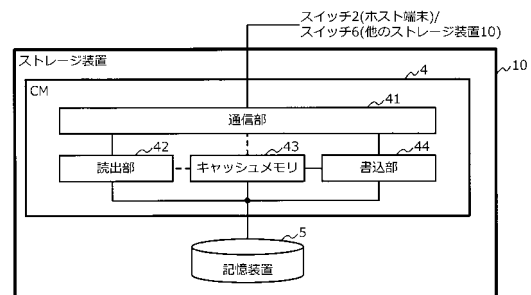
(54) 【発明の名称】 ストレージ装置、ストレージシステム、リカバリプログラム、及びリカバリ方法

(57) 【要約】

【課題】複数の記憶装置をそなえるストレージシステムにおいて、故障した記憶装置のデータの復旧処理における他の記憶装置からの情報の読み出し性能を向上させる。

【解決手段】1以上の第1記憶装置5と、前記1以上の第1記憶装置5が記憶する複数の第1ブロック情報であって、故障した第2記憶装置5が記憶する複数の第2ブロック情報の復元に用いられる前記複数の第1ブロック情報を、前記1以上の第1記憶装置5から記憶領域のアドレス順に読み出す読出部42と、前記読出部42により前記1以上の第1記憶装置5から読み出し済の第1ブロック情報を、前記複数の第2ブロック情報を段階的に復元するために、前記複数の第2ブロック情報の復元先へ出力する出力部41と、をそなえる。

【選択図】図4



【特許請求の範囲】**【請求項 1】**

1 以上の第 1 記憶装置と、

前記 1 以上の第 1 記憶装置が記憶する複数の第 1 ブロック情報であって、故障した第 2 記憶装置が記憶する複数の第 2 ブロック情報の復元に用いられる前記複数の第 1 ブロック情報を、前記 1 以上の第 1 記憶装置から記憶領域のアドレス順に読み出す読出部と、

前記読出部により前記 1 以上の第 1 記憶装置から読み出し済の第 1 ブロック情報を、前記複数の第 2 ブロック情報を段階的に復元するために、前記複数の第 2 ブロック情報の復元先へ出力する出力部と、をそなえることを特徴とする、ストレージ装置。

【請求項 2】

前記ストレージ装置とは異なる第 1 ストレージ装置、又は、前記出力部、から入力された第 1 ブロック情報に基づき、前記第 2 ブロック情報を段階的に復元する復元部をさらにそなえることを特徴とする、請求項 1 記載のストレージ装置。

【請求項 3】

前記第 1 ストレージ装置又は前記出力部から入力される第 1 ブロック情報を保持する保持部をさらにそなえ、

前記復元部は、所定のタイミングで、前記保持部が保持する 1 以上の第 1 ブロック情報に基づき、前記第 2 ブロック情報を段階的に復元することを特徴とする、請求項 2 記載のストレージ装置。

【請求項 4】

前記第 2 ブロック情報の復元先は、前記ストレージ装置とは異なる第 2 ストレージ装置にそなえられた第 1 記憶装置における空き記憶領域であり、

前記出力部は、前記読出部により前記第 1 記憶装置から読み出し済の第 1 ブロック情報を、前記第 2 ストレージ装置へ送信することを特徴とする、請求項 1 ~ 3 のいずれか 1 項記載のストレージ装置。

【請求項 5】

前記第 2 ブロック情報の復元先は、前記故障した第 2 記憶装置の代替となる第 3 記憶装置における空き記憶領域であり、

前記出力部は、前記読出部により前記第 1 記憶装置から読み出し済の第 1 ブロック情報を、前記第 3 記憶装置へ送信することを特徴とする、請求項 1 ~ 3 のいずれか 1 項記載のストレージ装置。

【請求項 6】

複数の記憶装置と、

前記複数の記憶装置のうちの複数の第 1 記憶装置の各々が記憶する複数の第 1 ブロック情報に基づき、前記複数の記憶装置のうちの故障した第 2 記憶装置が記憶する複数の第 2 ブロック情報を復元する 1 以上のストレージ装置と、をそなえ、

前記 1 以上のストレージ装置は、

複数の第 1 記憶装置の各々から、前記複数の第 1 ブロック情報を記憶領域のアドレス順に読み出し、

前記読み出しの処理において前記複数の第 1 記憶装置の各々から読み出し済の第 1 ブロック情報に基づき、前記複数の第 2 ブロック情報を段階的に復元する、ことを特徴とする、ストレージシステム。

【請求項 7】

コンピュータに、

1 以上の第 1 記憶装置が記憶する複数の第 1 ブロック情報であって、故障した第 2 記憶装置が記憶する複数の第 2 ブロック情報の復元に用いられる前記複数の第 1 ブロック情報を、前記 1 以上の第 1 記憶装置から記憶領域のアドレス順に読み出し、

前記読み出しの処理において前記 1 以上の第 1 記憶装置から読み出し済の第 1 ブロック情報を、前記複数の第 2 ブロック情報を段階的に復元するために、前記複数の第 2 ブロック情報の復元先へ出力する、

10

20

30

40

50

処理を実行させることを特徴とする、リカバリプログラム。

【請求項 8】

複数の記憶装置と、前記複数の記憶装置に対する制御を行なう 1 以上のストレージ装置とをそなえるストレージシステムにおけるリカバリ方法であって、

前記 1 以上のストレージ装置は、

前記複数の記憶装置のうちの複数の第 1 記憶装置の各々が記憶する複数の第 1 ブロック情報に基づき、前記複数の記憶装置のうちの故障した第 2 記憶装置が記憶する複数の第 2 ブロック情報を復元し、

前記復元の処理において、

前記 1 以上のストレージ装置により、

複数の第 1 記憶装置の各々から、前記複数の第 1 ブロック情報を記憶領域のアドレス順に読み出し、

前記読み出しの処理において前記複数の第 1 記憶装置の各々から読み出し済の第 1 ブロック情報に基づき、前記複数の第 2 ブロック情報を段階的に復元する、

ことを特徴とする、リカバリ方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ストレージ装置、ストレージシステム、リカバリプログラム、及びリカバリ方法に関する。

【背景技術】

【0002】

複数の HDD (Hard Disk Drive) や SSD (Solid State Drive) 等の記憶装置をそなえるストレージ装置では、種々の消失訂正技術により、記憶装置が故障した場合でもデータのリカバリを可能としている。

【0003】

図 2 1 は、ストレージ装置をそなえるストレージシステム 1 0 0 におけるディスク故障からの復旧の手法の一例を示す図である。図 2 1 に示すストレージシステム 1 0 0 は、6 つの HDD - 1 ~ HDD - 6 に 3 つのストライプ - 1 ~ ストライプ - 3 が設定されたディスクグループをそなえ、このディスクグループでは消失訂正符号 (Erasure Code) により符号化された情報が各ストライプに格納される。

【0004】

例えば図 2 1 の上段に示すように、ストライプ - 2 には、データが 4 つの HDD - 2 ~ HDD - 5 に分散して格納され (2 D 1 ~ 2 D 4 参照)、ストライプ - 2 のパリティが HDD - 6 に格納される (2 P 参照)。ここで、ストライプ - 2 における HDD - 1 は空きブロックであり、HDD - 2 ~ HDD - 6 のいずれかの故障時に代替ブロックとして用いられる。また、2 D 1 とはストライプ - “ 2 ” のデータ “ D ” のうちの “ 1 ” 番目のブロックを意味し、2 P とはストライプ - “ 2 ” のパリティ “ P ” のブロックを意味する。なお、以下の説明において、代替ブロックの “ 空き ” の表記を省略する場合がある。

【0005】

このようなストレージシステム 1 0 0 において HDD - 5 が故障した場合、図 2 1 の下段に例示するように、ストレージ装置のコントローラモジュール (CM ; Controller Module, 図示省略) 等の制御装置 (以下、CM と表記する) は、HDD - 5 に格納されていたデータを復旧するために再構築 (リビルド) を行なう。例えば CM は、ストライプ - 1 について、HDD - 1 ~ HDD - 4 から 1 D 1 ~ 1 D 4 のデータを取得し、1 D 1 ~ 1 D 4 からパリティ計算を行なって 1 P を生成し、生成した 1 P を HDD - 6 の代替ブロックに書き込む。他のストライプについても同様に、CM は、ストライプ - 2 について 2 D 1 ~ 2 D 3 及び 2 P のデータから 2 D 4 を生成し、生成した 2 D 4 を HDD - 1 の代替ブロックに書き込む。また、CM は、ストライプ - 3 について 3 D 1、3 D 2、3 D 4、及び 3 P のデータから 3 D 3 を生成し、生成した 3 D 3 を HDD - 2 の代替ブロックに書き込

10

20

30

40

50

む。

【 0 0 0 6 】

関連する技術として、データの各ブロックを消失符号化して、グループ化されたストレージノードに分配し、ストレージノードに障害が発生した場合、未使用のノードに他のノードのデータを使用して故障ノードのデータの再構築を行なう技術が知られている（例えば、特許文献 1 参照）。

【 先行技術文献 】

【 特許文献 】

【 0 0 0 7 】

【 特許文献 1 】 特開 2 0 1 0 - 7 9 8 8 6 号 公 報

10

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 0 8 】

ストレージシステム 1 0 0 において、各 HDD 上のブロック配置は CM 等による当該 HDD の空き領域の管理に応じて変化する。

【 0 0 0 9 】

従って、図 2 2 の下段に例示するように、ストライプ単位のデータ解放や再割り当てが繰り返されると、各 HDD 内のブロックには、データ又はパリティの情報がストライプ順ではなくランダムな順序で格納されることになる。例えば HDD - 3 には、HDD の記憶領域の先頭から順に、2 D 2 (ストライプ - 2)、1 D 3 (ストライプ - 1)、3 D 1 (ストライプ - 3) のブロックが格納される。

20

【 0 0 1 0 】

このような図 2 2 の下段に示す状態において HDD - 5 が故障した場合を考える。この場合、各 HDD では、ストライプ - 1 ~ ストライプ - 3 の構成ブロックの配置がランダムとなっているため、図 2 3 の下段に例示するように、CM によるリビルドの際、HDD 上のアドレス順とは異なる順序による HDD へのアクセスが発生する。

【 0 0 1 1 】

例えば HDD - 3 では、CM により、最初にストライプ - 1 のリカバリのために記憶領域の中央付近のアドレスから 1 D 3 が読み出され、次いでストライプ - 2 のリカバリのために記憶領域の先頭付近のアドレスから 2 D 2 が読み出される。そして、最後にストライプ - 3 のリカバリのために記憶領域の末尾 (最終アドレス) 付近のアドレスから 3 D 1 が読み出される。

30

【 0 0 1 2 】

また、例えば代替ブロックを持つ HDD - 1 では、CM により、最初にストライプ - 1 のリカバリのために記憶領域の末尾付近のアドレスから 1 D 1 が読み出され、次いでストライプ - 2 のリカバリのために他の HDD の情報に基づき生成された 2 D 4 が記憶領域の中央付近のアドレスに書き込まれる。そして、最後にストライプ - 3 のリカバリのために記憶領域の先頭付近のアドレスから 3 P が読み出される。

【 0 0 1 3 】

このように、ストレージシステム 1 0 0 においてストライプ単位のデータ解放や再割り当てが繰り返されると、各 HDD では、リビルドの際にリカバリ対象のストライプの順にアクセスが行なわれるため、このアクセスはランダムアクセスとなる。

40

【 0 0 1 4 】

図 2 3 の例では、ストレージシステム 1 0 0 が 3 つのストライプを管理し、各 HDD に最大 3 つのブロックが格納されるものとして説明したが、実際にはさらに多くのストライプが管理され、1 つの HDD に格納されるブロックも非常に多くなる。例えば HDD が 1 0 0 M B / s 程度の読み出し性能を持つ SAS (Serial Attached SCSI (Small Computer System Interface)) 規格に対応した記憶装置であり、ブロックが 4 K B 程度の小さいサイズである場合、図 2 3 に示す例では各 HDD のアクセスが 1 0 M B / s 程度のランダムアクセスとなってしまう。

50

【 0 0 1 5 】

以上のように、ストレージ装置を1以上そなえるストレージシステムでは、HDDに格納されたブロックがストライプ順ではない場合、リビルドの際にストレージ装置においてHDDの読み出しがランダム化され、性能劣化が生じてしまう。

【 0 0 1 6 】

なお、上述した課題は、上述の如く、消失訂正符号により符号化された情報を格納するストレージシステムにおいて生じ得るものである。このようなストレージシステムとしては、RAID (Redundant Arrays of Inexpensive Disks) 5 と、RAID 5 に対してストライプを追加のディスクにまで拡張するワイドストライプとを組み合わせた構成が挙げられる。また、RAID 5 に代えて他のRAID技術(例えばRAID 6) 或いは複数のRAID技術の組み合わせが採用された構成や、他の消失訂正符号を用いたストレージシステムにおいても、上記課題は同様に生じ得る。

10

【 0 0 1 7 】

さらに、代替ブロックを持たない、例えば通常のRAID 5 やRAID 6 等を採用したストレージシステムにおいても、HDDに格納されたブロックがストライプ順にならない場合、上記課題は同様に生じ得る。

【 0 0 1 8 】

1つの側面では、本発明は、複数の記憶装置をそなえるストレージシステムにおいて、故障した記憶装置のデータの復旧処理における他の記憶装置からの情報の読み出し性能を向上させることを目的とする。

20

【課題を解決するための手段】

【 0 0 1 9 】

1つの態様では、本件のストレージ装置は、1以上の第1記憶装置と、前記1以上の第1記憶装置が記憶する複数の第1ブロック情報を、前記1以上の第1記憶装置から記憶領域のアドレス順に読み出す読出部と、をそなえる。前記第1ブロック情報は、故障した第2記憶装置が記憶する複数の第2ブロック情報の復元に用いられる情報である。また、前記ストレージ装置は、前記読出部により前記1以上の第1記憶装置から読み出し済の第1ブロック情報を、前記複数の第2ブロック情報を段階的に復元するために、前記複数の第2ブロック情報の復元先へ出力する出力部をさらにそなえる。

30

【発明の効果】

【 0 0 2 0 】

1つの側面では、複数の記憶装置をそなえるストレージシステムにおいて、故障した記憶装置のデータの復旧処理における他の記憶装置からの情報の読み出し性能を向上させることができる。

【図面の簡単な説明】

【 0 0 2 1 】

【図1】第1実施形態の一例としてのストレージシステムの構成例を示す図である。

【図2】図1に示すストレージシステムにおけるリビルド処理の一例を説明する図である。

。

【図3】図1に示すストレージシステムにおけるリビルド処理の一例を説明する図である。

40

。

【図4】図1に示すストレージ装置の機能構成例を示す図である。

【図5】図1に示すストレージシステムにおけるリビルド処理の動作例を説明する図である。

【図6】図1に示すストレージシステムにおけるリビルド処理の動作例を説明する図である。

【図7】図1に示すストレージシステムにおけるリビルド処理の動作例を説明する図である。

【図8】図1に示すストレージシステムにおけるリビルド処理の動作例を説明する図である。

50

【図 9】図 1 に示すストレージシステムにおけるリビルド処理の動作例を説明する図である。

【図 10】図 1 に示すストレージシステムにおけるリビルド処理の動作例を説明する図である。

【図 11】第 1 実施形態に係るストレージシステムにおける全体の処理の一例を説明するフローチャートである。

【図 12】第 1 実施形態に係るストレージシステムにおけるリビルド処理の一例を説明するフローチャートである。

【図 13】第 1 実施形態に係るストレージシステムにおけるリビルド処理の一例を説明するフローチャートである。

【図 14】第 2 実施形態に係るストレージシステムにおけるリビルド処理の動作例を説明する図である。

【図 15】第 2 実施形態に係るストレージシステムの他の適用例におけるリビルド処理の動作を説明する図である。

【図 16】第 2 実施形態に係るストレージシステムにおける運用ストレージ装置のリビルド処理の一例を説明するフローチャートである。

【図 17】第 2 実施形態に係るストレージシステムにおける待機ストレージ装置のリビルド処理の一例を説明するフローチャートである。

【図 18】第 2 実施形態に係るストレージシステムにおける待機ストレージ装置のリビルド処理の他の例を説明するフローチャートである。

【図 19】第 1 及び第 2 実施形態に係るストレージ装置のハードウェア構成例を示す図である。

【図 20】第 1 及び第 2 実施形態に係るストレージシステムの他の構成例を示す図である。

【図 21】ストレージシステムにおけるディスク故障からの復旧の手法の一例を示す図である。

【図 22】ストレージシステムにおいてストライプ単位のリカバリ処理の一例を示す図である。

【図 23】図 22 の下段に示す状態においてディスクが故障した場合のリカバリ処理におけるディスクアクセスの様子を示す図である。

【発明を実施するための形態】

【0022】

以下、図面を参照して本発明の実施の形態を説明する。ただし、以下に説明する実施形態は、あくまでも例示であり、以下に明示しない種々の変形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。なお、以下の実施形態で用いる図面において、同一符号を付した部分は、特に断らない限り、同一若しくは同様の部分を表す。

【0023】

〔1〕第 1 実施形態

〔1-1〕ストレージシステムの構成例

図 1 は第 1 実施形態の一例としてのストレージシステム 1 の構成例を示す図である。図 1 に示すように、ストレージシステム 1 は例示的にスイッチ 2 及び 6、並びに 1 以上（図 1 では複数）のストレージ装置 10 - 1 ~ 10 - m（m は自然数）をそなえることができる。なお、以下の説明においてストレージ装置 10 - 1 ~ 10 - m を区別しない場合には単にストレージ装置 10 と表記する。

【0024】

ストレージシステム 1 は、ユーザに対してストレージ装置 10 の記憶領域を提供するものであり、例えばネットワーク 3 を介してユーザの使用ユーザ端末（ホスト装置）からストレージ装置 10 へのアクセスが可能となっている。ストレージシステム 1 としては、例えば複数のストレージ装置 10（筐体）をそなえるクラスタ構成のストレージシステ

10

20

30

40

50

ムであってもよいし、単一のストレージ装置 10 をそなえた構成であってもよい。図 1 の例では、ストレージシステム 1 は、m 個のストレージ装置 10 をそなえ、分散アルゴリズムによりストレージ装置 10 が相互に通信可能な分散ストレージシステムである。

【0025】

スイッチ 2 は、ストレージ装置 10 及びネットワーク 3 と接続され、ストレージシステム 1 を使用するユーザのユーザ端末とストレージ装置 10 との間の通信（クライアント通信）の切り替え制御等を行なうものである。スイッチ 6 は、ストレージ装置 10 と接続され、ストレージ装置 10 間の通信（クラスタ内部通信）の切り替え制御等を行なうものである。

【0026】

なお、ネットワーク 3 は、インターネットであってもよいし、LAN（Local Area Network）又は SAN（Storage Area Network）等のイントラネットを形成するネットワークであってもよい。また、図 1 の例ではストレージシステム 1 にスイッチ 2 及び 6 がそれぞれ 1 つずつそなえられるものとしたが、これに限定されるものではない。例えば複数のストレージ装置 10 が互いに異なる拠点にそなえられる場合には、各拠点にスイッチ 2 及びスイッチ 6 を設け、スイッチ 2 間及びスイッチ 6 間を、それぞれインターネットや LAN 又は SAN 等のイントラネットを形成するネットワーク等を介して相互に通信可能に接続してもよい。さらに、ストレージシステム 1 は、スイッチ 2 及びスイッチ 6 をクライアント通信及びクラスタ内部通信で共用のスイッチとして、いずれか一方のみそなえてもよい。

10

20

【0027】

ストレージ装置 10 は、それぞれ 1 以上（図 1 では 1 つ）の CM 4 をそなえる。また、ストレージ装置 10 - 1 は記憶装置 5 - 1 をそなえ、ストレージ装置 10 - m は記憶装置 5 - n をそなえる。なお、以下の説明において記憶装置 5 - 1 ~ 5 - n を区別しない場合には単に記憶装置 5 と表記する。ストレージ装置 10 は、図 1 の例ではそれぞれ 1 つの記憶装置 5 をそなえるものとしたが、複数の記憶装置 5 をそなえてもよい。

【0028】

CM 4 は、スイッチ 2 を介したユーザ端末からの要求、並びにスイッチ 6 を介した他のストレージ装置 10（CM 4）からの要求に応じて、記憶装置 5 の記憶領域に対する種々のアクセス制御を行なうコンピュータ（情報処理装置）の一例である。このアクセス制御には、記憶装置 5 の故障に伴うリビルド処理が含まれる。例えば CM 4 は、リビルド処理において、他のストレージ装置 10 の CM 4 とともに、分散アルゴリズムに基づく協調動作を行なうことができる。

30

【0029】

CM 4 は、例えば CPU（Central Processing Unit）4 a、メモリ 4 b、及び IF（Interface）4 c ~ 4 e をそなえる。

【0030】

CPU 4 a は、種々の制御や演算を行なう演算処理装置（プロセッサ）の一例である。CPU 4 a は、メモリ 4 b、及び IF 4 c ~ 4 e とバスで相互に通信可能に接続され、メモリ 4 b 又は図示しない ROM（Read Only Memory）等に格納されたプログラムを実行することにより、CM 4 における種々の機能を実現することができる。

40

【0031】

メモリ 4 b は、種々のデータやプログラムを格納する記憶装置である。第 1 実施形態に係るメモリ 4 b はさらに、後述するリビルド処理において記憶装置 5 から読み出した情報及び記憶装置 5 へ格納する（書き込む）情報を一時的に記憶するキャッシュメモリとして用いられる。なお、メモリ 4 b としては、例えば RAM（Random Access Memory）等の揮発性メモリが挙げられる。

【0032】

IF 4 c ~ 4 e は、それぞれスイッチ 2、記憶装置 5、スイッチ 6 との間の接続及び通信の制御等を行なう通信インタフェースである。例えば IF 4 c 及び 4 e はホストアダプ

50

タであり、I F 4 d はデバイス（ディスク）アダプタである。これらの通信インタフェース（アダプタ）としては、L A N、S A N、F C（Fibre Channel）、インフィニバンド（InfiniBand）等に準拠したアダプタが挙げられる。

【 0 0 3 3 】

記憶装置 5 は、種々のデータやプログラム等を格納するハードウェアである。記憶装置 5 としては、例えば H D D 等の磁気ディスク装置や、S S D 等の半導体ドライブ装置等の各種記憶装置が挙げられる。

【 0 0 3 4 】

〔 1 - 2 〕ストレージシステムにおけるリビルド処理について

次に、第 1 実施形態に係るストレージシステム 1 におけるリビルド処理について、図 2 及び図 3 を参照して簡単に説明する。以下、前提として、記憶装置 5 が H D D であるものとする。また、ストレージシステム 1 が 6 つの H D D - 1 ~ H D D - 6 に 3 つのストライプ - 1 ~ ストライプ - 3 が設定されたディスクグループをそなえ、このディスクグループでは消失訂正符号により符号化された情報が各ストライプに格納されるものとする（図 2 3 の上段参照）。

10

【 0 0 3 5 】

例えば図 2 3 の上段に示す状態から H D D - 5 が故障した場合に、ストライプ - 2 の消失したブロック（消失ブロック）である 2 D 4 を H D D - 1 に復旧する場合を想定する。ストレージシステム 1（ストレージ装置 1 0 の C M 4）は、図 2 の上段に示すように、ストライプ順に関係なく、各 H D D 上で物理アドレス順にシーケンシャルにブロックを読み出す。例えば H D D - 3 では、2 D 2、1 D 3、3 D 1 の順にブロックが読み出される。なお、H D D - 2 では H D D - 5 の消失ブロック（3 D 3）の復旧及び格納前であるため、1 D 2、2 D 1 の順にブロックが読み出される。

20

【 0 0 3 6 】

そして、ストレージシステム 1（C M 4）は、図 2 の中段に示すように、読み出したブロック（データ（D）ブロック及びパリティ（P）ブロック）順に、インクリメンタルに各ストライプを復旧する。以下、データブロック及びパリティブロックを区別しない場合には、これらを情報ブロックと表記する。

【 0 0 3 7 】

例えば H D D - 3 の記憶領域の先頭付近のアドレスに格納された 2 D 2 が最初に読み出されると、この 2 D 2 が H D D - 1 の新 2 D 4 のブロックに反映される（矢印（1）参照）。次いで、H D D - 6 の記憶領域の中央付近のアドレスに格納された 2 P が読み出されると、この 2 P が H D D - 1 の新 2 D 4 のブロックに反映される（矢印（2）参照）。次に、H D D - 2 の記憶領域の中央付近のアドレスに格納された 2 D 1 が読み出されると、この 2 D 1 が H D D - 1 の新 2 D 4 のブロックに反映される（矢印（3）参照）。そして、H D D - 4 の記憶領域の中央付近のアドレスに格納された 2 D 3 が読み出されると、この 2 D 3 が H D D - 1 の新 2 D 4 のブロックに反映される（矢印（4）参照）。

30

【 0 0 3 8 】

なお、図 2 にはストライプ - 2 の消失ブロックに着目して新 2 D 4 の復旧について説明したが、ストライプ - 1 及びストライプ - 3 の消失ブロックについて新 1 P 及び新 3 D 3 を復旧する処理も、各 H D D からのシーケンシャルな情報ブロックの読み出しの過程で順次実施される。

40

【 0 0 3 9 】

以上のように、第 1 実施形態に係るストレージシステム 1 は、読み出したブロックを用いて段階的に消失ブロックを復旧することができる。これにより、ストレージシステム 1 では、記憶装置の故障によるリビルド処理において、C M 4 は各記憶装置からシーケンシャルに情報ブロックを読み出すことができるため、正常な記憶装置からの情報の読み出し性能を向上することができる。従って、図 2 3 に示す各記憶装置からの完全なランダムアクセスにより情報を読み出す手法と比較して、最大で 1 0 倍以上のスループットとすることができ、リカバリ性能を大幅に向上させることができる。

50

【 0 0 4 0 】

ここで、ストレージシステム 1 が段階的にストライプを復旧することのできる理由を説明する。

【 0 0 4 1 】

ストレージシステム 1 では、各ストライプの情報ブロックが消失訂正符号により符号化されている。消失訂正符号により符号化された情報ブロックでは、消失ブロックの情報を、或るストライプにおける消失ブロック以外の正常な情報ブロックから算出（復旧）することができる。この算出手法としては、パリティが 1 つの場合、例えば正常な情報ブロックの排他的論理和（XOR）を算出するといった手法が挙げられる。

【 0 0 4 2 】

XOR 演算は、可換（Commutative）演算である。このため、ストレージシステム 1 は、各 HDD からシーケンシャルに読み出した情報ブロックを用いて段階的に（インクリメンタルに）各ストライプを復旧することができるのである。以下の説明では、XOR 演算を示す演算子として “ + ” を用い、図面では演算子として “ + ” を丸で囲んだ記号を用いる。

【 0 0 4 3 】

例えば図 2 の下段に示すように、対比例としての図 2 3 の演算では、消失ブロックを復旧するための全ての情報ブロックが揃ってから、新 $2D4 = 2D1 + 2D2 + 2D3 + 2P$ の演算が順序通りに行なわれる。一方、第 1 実施形態に係るストレージシステム 1 は、図 2 の矢印（1）及び（2）で読み出した $2D2$ 及び $2P$ の演算を行ない、その演算結果と図 2 の矢印（3）で読み出した $2D1$ との演算を行なう。そして、ストレージシステム 1 は、最後にその演算結果と図 2 の矢印（4）で読み出した $2D3$ との演算を行なう。

【 0 0 4 4 】

なお、消失ブロックの情報を算出（復旧）する算出手法としては、上述した XOR 演算に限定されるものではない。例えばストライプごとにパリティが 2 つ以上含まれる場合には、XOR 演算を用いたリードソロモン（RS ; Reed-Solomon）符号ベースの演算が行なわれてもよい。ストレージシステム 1 では、このような RS 符号を消失訂正符号として用いて符号化された情報が各ストライプに格納されている場合でも、段階的に消失ブロック（ストライプ）を復旧することができる。消失ブロックの情報を算出（復旧）する算出手法としては、消失訂正符号がガロア体（有限体）を用いる可換演算の可能な符号であれば、上述したものの以外の種々の手法が用いられてよい。

【 0 0 4 5 】

説明の簡略化のため、以下の説明では、消失ブロックの情報の算出手法として XOR 演算が用いられるものとする。

【 0 0 4 6 】

ところで、CM 4 は、各記憶装置 5 から情報ブロックを読み出す都度、当該情報ブロックのストライプにおける復旧先の代替ブロック（リカバリ対象ブロック）へ XOR 演算等による反映を行なってよい。しかし、当該代替ブロックを有する記憶装置 5 においてもシーケンシャルリードが行なわれているため、当該代替ブロックを有する記憶装置 5 では、シーケンシャルリードの最中に代替ブロックへの書き込みアクセスが頻発して読み出し性能が劣化し、リビルド処理の性能低下が生じることがある。

【 0 0 4 7 】

そこで、CM 4（又は記憶装置 5）は、図 3 に示すように、或るストライプの消失ブロックを復旧する場合、当該消失ブロックの復旧に用いる情報ブロックをキャッシュメモリ等に一時的に保持しておくことができる。

【 0 0 4 8 】

例えば CM 4 は、図 3 の矢印（1）に示すように、HDD - 3 から $2D2$ を読み出し、HDD - 1 の新 $2D4$ のための書込キャッシュに反映する。このときストライプ - 2 の $2D4$ 復旧の進捗は、 $2D4$ として $2D2$ が設定された状態である。

【 0 0 4 9 】

10

20

30

40

50

また、CM4は、図3の矢印(2)に示すように、HDD-6から2Pを読み出し、書込キャッシュに反映する。このときストライプ-2の2D4復旧の進捗は、2D4として、2D4に設定された2D2と2PとのXOR演算結果が設定された状態である。

【0050】

なお、このタイミングで書込キャッシュの容量が逼迫した場合、CM4は、書込キャッシュ内の情報ブロック(2D2+2P)をHDD-1の代替ブロックに書き込む(フラッシュする)。

【0051】

さらに、CM4は、図3の矢印(3)に示すように、HDD-2から2D1を読み出し、HDD-1の代替ブロック(新2D4)の情報を書込キャッシュにリロードする。そして、CM4は、2D1を書込キャッシュに反映する。このときストライプ-2の2D4復旧の進捗は、2D4として、書込キャッシュにリロードされた2D2+2Pと2D1とのXOR演算結果が設定された状態である。

10

【0052】

また、CM4は、図3の矢印(4)に示すように、HDD-4から2D3を読み出し、書込キャッシュに反映する。このときストライプ-2の2D4復旧の進捗は、2D4として、2D4に設定された(2D2+2P)+2D1と2D3とのXOR演算結果が設定された状態である。

【0053】

なお、このタイミングで書込キャッシュの容量が逼迫した場合、CM4は、書込キャッシュ内の情報ブロック((2D2+2P)+2D1)+2D3をHDD-1の代替ブロックに書き込む(フラッシュする)。

20

【0054】

以上により書込キャッシュを用いた段階的なリビルド処理が行なわれる。なお、ストライプ-1及びストライプ-3の消失ブロックについて新1P及び新3D3を復旧する処理も、HDD-2及びHDD-6に対応する書込キャッシュにおいて、各HDDからのシーケンシャルな情報ブロックの読み出しの過程で順次実施される。

【0055】

各記憶装置5についてシーケンシャルに情報ブロックを読み込む場合、できるだけ多くの情報ブロックを保持できるように、大容量のキャッシュメモリが用いられることが好ましい。しかし、CM4(又は記憶装置5)の各々に大容量のキャッシュメモリを搭載することは、コスト増加の観点から難しい場合がある。

30

【0056】

これに対し、上述したストレージシステム1によれば、図3に示すように、書込キャッシュの容量が逼迫した等の場合に、書込キャッシュに保持された情報ブロックをまとめてXOR演算して代替ブロックに書き込むことができる。そして、CM4は、容量の空いた書込キャッシュに情報ブロックを蓄積していき、再び書込キャッシュの容量が逼迫した等の場合に、代替ブロックから情報を再読み出しし、書込キャッシュに格納された情報ブロックと再読み出した情報とをXOR演算して代替ブロックに書き込むのである。

【0057】

このように、ストレージシステム1では、図2に示すように各ストライプをインクリメンタルに復旧することができるため、消失ブロックの復旧データである全ての情報ブロックが揃うまで書込キャッシュに情報ブロックを溜め込まなくてよい。これにより、ストレージシステム1は、大容量のキャッシュメモリをそなえなくてもよく、コスト増加を抑制することができる。

40

【0058】

〔1-3〕ストレージ装置の構成例

次に、図4を参照してストレージ装置10の構成例について説明する。ストレージ装置10(CM4)は、複数のストレージ装置10のCM4と協働して、複数の記憶装置5に対する各種制御を行なうことができる。この制御には、ユーザ端末からの書込要求に応じ

50

て書込データのパリティ演算を行ない、複数の情報ブロックを生成してストライプとして各記憶装置 5 に分散させる制御が含まれる。また、この制御には、ユーザ端末からの読出要求に応じてストライプから情報ブロックを取得して読出データを構築し、ユーザ端末へ出力する制御も含まれる。

【 0 0 5 9 】

また、例えば C M 4 は、ユーザ端末からの要求に応じてアクセス対象のストライプに対応する記憶装置 5 の情報ブロックへのアクセスを行ったり、記憶装置 5 の故障を検出した場合に他のストレージ装置 1 0 の C M 4 へ通知を行なってリビルド処理を実行することができる。これらの制御は、既知の種々の手法により行なうことが可能であり、その詳細な説明は省略する。

【 0 0 6 0 】

さらに、第 1 実施形態に係るストレージ装置 1 0 (C M 4) は、リビルド処理において図 2 及び図 3 に示すような動作を実現するため、図 4 に示すように、例示的に通信部 4 1、読出部 4 2、キャッシュメモリ 4 3、及び書込部 4 4 をそなえることができる。

【 0 0 6 1 】

通信部 4 1 は、他の C M 4 との間で通信を行なうものであり、例えばリビルド処理に関する制御情報や情報ブロック等の種々の情報を他の C M 4 との間で送受信する。例えば通信部 4 1 は、図 1 に示す C P U 4 a、I F 4 c、及び I F 4 e の少なくとも一部の機能により実現することができる。

【 0 0 6 2 】

例えば記憶装置 5 (第 2 記憶装置) の故障を検出した C M 4 は、故障した記憶装置 5 に格納されていた消失ブロック (第 2 ブロック情報) に関する情報を通信部 4 1 により他の C M 4 に通知する。この通知を受信した C M 4 は、自装置 1 0 がそなえる記憶装置 5 について、消失ブロックを復旧させる代替ブロックの有無や、消失ブロックと同じストライプの情報ブロックの有無等を判断して、消失ブロックのストライプごとに復旧先の記憶装置 5 を決定する。例えば記憶装置 5 に代替ブロックが有り、或る消失ブロックと同じストライプの情報ブロックが無い場合、C M 4 は、自装置 1 0 の記憶装置 5 が当該消失ブロックのストライプの復旧先であることを通信部 4 1 により他の C M 4 に通知する。

【 0 0 6 3 】

読出部 4 2 は、リビルド処理において、自装置 1 0 がそなえる記憶装置 5 (第 1 記憶装置) に格納された情報ブロック (第 1 ブロック情報) を記憶装置 5 の物理アドレスの先頭からシーケンシャルに読み出し、読み出した情報ブロックを順次通信部 4 1 に渡す。なお、自装置 1 0 が或るストライプの復旧先 (復元先) の記憶装置 5 をそなえる場合、読出部 4 2 は、シーケンシャルに読み出す過程で当該情報ブロックの読み出しをスキップしてよい。

【 0 0 6 4 】

このように、読出部 4 2 は、1 以上の第 1 記憶装置 5 が記憶する複数の第 1 ブロック情報であって、故障した第 2 記憶装置 5 が記憶する複数の第 2 ブロック情報の復元に用いられる複数の第 1 ブロック情報を、1 以上の第 1 記憶装置 5 から記憶領域のアドレス順に読み出すものであるといえる。

【 0 0 6 5 】

なお、通信部 4 1 は、読出部 4 2 が読み出した情報ブロックを、当該情報ブロックのストライプの復旧先であるストレージ装置 1 0 (C M 4) へ送信 (転送) し、復旧先の C M 4 は、他の C M 4 から受信した情報ブロックを書込部 4 4 に出力する。このとき復旧先の C M 4 は、書込部 4 4 によりキャッシュメモリ 4 3 の使用量を監視し、容量が逼迫した場合、通信部 4 1 により他の C M 4 に対して容量が逼迫したことを示す通知 (或いは情報ブロックの送信を抑止させる通知) を行なうことができる。この通知を受信した C M 4 (通信部 4 1) は、読出部 4 2 に対して復旧先の C M 4 に対応する情報ブロックの読み出しを中止させてもよいし、読出部 4 2 に読み出された当該情報ブロックをキャッシュメモリ 4 3 に一時的に格納 (退避) してもよい (図 4 の破線参照) 。

10

20

30

40

50

【 0 0 6 6 】

また、復旧先の C M 4 は、キャッシュメモリ 4 3 が使用可能になった場合、通信部 4 1 により他の C M 4 に対して容量が確保できたことを示す通知（或いは情報ブロックの送信を再開させる通知）を行なうことができる。この通知を受信した C M 4（通信部 4 1）は、読出部 4 2 に対して復旧先の C M 4 に対応する情報ブロックを読み出させてもよいし、キャッシュメモリ 4 3 に格納した当該情報ブロックを復旧先の C M 4 へ送信してもよい（図 4 の破線参照）。

【 0 0 6 7 】

このように、通信部 4 1 は、読出部 4 2 により 1 以上の第 1 記憶装置 5 から読み出し済の第 1 ブロック情報を、複数の第 2 ブロック情報を段階的に復元するために、複数の第 2 ブロック情報の復元先へ出力する出力部の一例であるといえる。

10

【 0 0 6 8 】

キャッシュメモリ 4 3 は、図 3 に示す書込キャッシュの一例であり、例えば図 1 に示すメモリ 4 b の少なくとも一部の記憶領域を用いることにより実現することができる。キャッシュメモリ 4 3 は、自装置 1 0 が復旧先である場合に、復旧に用いる情報ブロックが格納される記憶領域である。また、上述のように、読出部 4 2 が読み出した情報ブロックを復旧先の C M 4 へ送信できない場合、当該情報ブロックの退避用の記憶領域として用いられてもよい。

【 0 0 6 9 】

このように、キャッシュメモリ 4 3 は、他の第 1 ストレージ装置 1 0 から入力される第 1 ブロック情報を保持する保持部の一例であるといえる。

20

【 0 0 7 0 】

書込部 4 4 は、通信部 4 1 から入力された情報ブロックをキャッシュメモリ 4 3 へ書き込む。また、書込部 4 4 は、例えば定期的に、キャッシュメモリ 4 3 の使用量を監視し、使用量が閾値を超えた（容量が逼迫した）場合、その旨を通信部 4 1 へ通知する。このとき書込部 4 4 は、記憶装置 5 の代替ブロックに情報ブロックが格納されているか否かを判断する。なお、閾値は、情報ブロックのサイズやキャッシュメモリ 4 3 の容量に応じて予め設定されるものであり、閾値として例えばキャッシュメモリ 4 3 の記憶領域のサイズの 8 0 % ~ 9 0 % 等の値を設定することができる。

【 0 0 7 1 】

代替ブロックに情報ブロックが格納されていない場合、書込部 4 4 は、キャッシュメモリ 4 3 に格納された複数の情報ブロックについて X O R 演算を行ない、演算結果を記憶装置 5 の代替ブロック（空き記憶領域）に書き込み、キャッシュメモリ 4 3 をクリアする。一方、代替ブロックに情報ブロックが格納されている場合、書込部 4 4 は、代替ブロックから情報ブロックを読み出し、読み出した情報ブロックと、キャッシュメモリ 4 3 に格納された複数の情報ブロックと、について X O R 演算を行ない、演算結果を記憶装置 5 の代替ブロックに書き込み、キャッシュメモリ 4 3 をクリアする。なお、書込部 4 4 は、キャッシュメモリ 4 3 の使用量が閾値以下となった（使用可能になった）場合、その旨を通信部 4 1 へ通知する。

30

【 0 0 7 2 】

なお、書込部 4 4 によるキャッシュメモリ 4 3 の容量の監視において、使用量が閾値を超えたか否かを判断する代わりに、キャッシュメモリ 4 3 の残容量を監視し、残容量が閾値以下となったか否かを判断してもよい。或いは、情報ブロックが一定（ブロック単位の）サイズであるため、書込部 4 4 は、キャッシュメモリ 4 3 に格納した情報ブロックの数をカウントし、情報ブロックの数が所定数以上となった場合にキャッシュメモリ 4 3 の容量が逼迫したと判断してもよい。

40

【 0 0 7 3 】

また、書込部 4 4 は、例えば受信した（又はキャッシュメモリ 4 3 に格納した）情報ブロックの数をカウントし、カウント値が閾値に達した場合に、当該ストライプの復旧（リカバリ）処理が完了したと判断することができる。なお、カウントする数としては、これ

50

に限定されるものではなく、代替ブロックに反映した情報ブロックの数であってもよいし、XOR演算を行なった回数であってもよい。また、情報ブロックの数をカウントする場合、閾値を「ストライプに含まれる情報ブロックの数」 - 「当該ストライプに含まれる消失ブロックの数」とすることができる。或いは、XOR演算を行なった回数をカウントする場合、閾値を「ストライプに含まれる情報ブロックの数」 - 「当該ストライプに含まれる消失ブロックの数」 - 1としてもよい。

【0074】

このように、書込部44は、他の第1ストレージ装置10から入力された第1ブロック情報に基づき、第2ブロック情報を段階的に復元する復元部の一例であるといえる。この復元部の一例としての書込部44は、入力された第1ブロック情報及び第2ブロック情報の復元先が記憶する情報を用いて、消失訂正符号に基づく演算を行ない、演算結果を第2ブロック情報の復元先へ書き込むのである。また、この復元部の一例としての書込部44は、所定のタイミングで、キャッシュメモリ43が保持する1以上の第1ブロック情報に基づき、第2ブロック情報を段階的に復元するのである。

10

【0075】

〔1-4〕ストレージシステムにおけるリビルド処理の動作説明

次に、図5～図10を参照して、ストレージシステム1におけるリビルド処理の動作をストレージ装置10間の通信に着目して説明する。以下、ストレージシステム1が6台のストレージ装置10をそなえるものとし、便宜上、これらのストレージ装置10をノード-1～ノード-6と表記する。また、ノード-5が故障し、ストレージシステム1がノード5のHDD-5に格納された2D4、1P、3D3の3ブロックのリビルド処理を行なうものとする。

20

【0076】

図5の上段に示すように、各ノードのCM4は、HDDの先頭ブロックを読出部42により読み出し、読み出した先頭ブロック（情報ブロック）を通信部41により復旧先の代替ブロック（リカバリ対象ブロック）を有するノードへ転送する。情報ブロックを受信したノードは、キャッシュメモリ43に格納する。

【0077】

図5の例では、ノード-1が2D4（ストライプ-2）のリカバリ対象ブロックを有し、ノード-2が3D3（ストライプ-3）のリカバリ対象ブロックを有し、ノード-6が1P（ストライプ-1）のリカバリ対象ブロックを有している。この場合、ノード-2及びノード-4は、それぞれHDD-2の先頭ブロックの1D2及びHDD-4の先頭ブロックの1D4をノード-6に転送し、ノード-3は、HDD-3の先頭ブロックの2D2をノード-1に転送する。また、ノード-1及びノード-6は、それぞれHDD-1の先頭ブロックの3P及びHDD-6の先頭ブロックの3D4をノード-2に転送する。

30

【0078】

次いで、図5の下段に示すように、ノード-2及びノード-6のCM4は、書込部44により、キャッシュメモリ43の容量が枯渇したと判断して、ノード-2及びノード-6のそれぞれのリカバリ対象ブロックについて部分的なリカバリを実施する。このときノード-2及びノード-6のリカバリ対象ブロックに書き込まれる情報である3D3'及び1P'は、それぞれ $3D3' = 3P + 3D4$ 、 $1P' = 1D2 + 1D4$ となる。一例として、図6に示すように、ストライプ番号とインデックス番号（パリティを除く）をそれぞれバイナリ化して、 $3P = 000100$ 、 $3D4 = 011100$ 、 $1D2 = 001010$ 、 $1D4 = 001100$ とした場合、3D3'及び1P'はそれぞれ以下の値となる。

40

【0079】

$$3D3' = 3P + 3D4 = 000100 + 011100 = 011000$$

$$1P' = 1D2 + 1D4 = 001010 + 001100 = 000110$$

【0080】

次に、図7の上段に示すように、各ノードのCM4は、HDDの2ブロック目（2ブロック目がリカバリ対象ブロックであれば3ブロック目）を読出部42により読み出し、読

50

み出した情報ブロックを通信部 4 1 によりリカバリ対象ブロックを有するノードへ転送する。情報ブロックを受信したノードは、キャッシュメモリ 4 3 に格納する。

【 0 0 8 1 】

図 7 の例では、ノード - 1 及びノード - 3 は、それぞれ HDD - 1 の 3 ブロック目の 1 D 1 及び HDD - 3 の 2 ブロック目の 1 D 3 をノード - 6 に転送し、ノード - 2 は、HDD - 2 の 3 ブロック目の 2 D 1 をノード - 1 に転送する。また、ノード - 4 は、HDD - 4 の 2 ブロック目の 3 D 2 をノード - 2 に転送する。なお、ノード - 6 は、HDD - 6 の 2 ブロック目の 2 P を転送する前にノード - 1 のキャッシュメモリ 4 3 の容量が逼迫したため、ノード - 1 からの通知により転送を抑止（保留）している。

【 0 0 8 2 】

次いで、図 7 の下段に示すように、ノード - 1 及びノード - 6 の CM 4 は、書込部 4 4 により、キャッシュメモリ 4 3 の容量が枯渇したと判断して、ノード - 1 及びノード - 6 のそれぞれのリカバリ対象ブロックについて部分的なりカバリを実施する。このときノード - 1 及びノード - 6 のリカバリ対象ブロックに書き込まれる情報である 2 D 4 ' 及び 1 P " は、それぞれ $2 D 4 ' = 2 D 2 + 2 D 1$ 、 $1 P " = 1 P ' + 1 D 1 + 1 D 3$ となる。なお、ノード - 6 の書込部 4 4 は、リカバリ対象ブロックから 1 P ' を読み出してから、読み出した 1 P ' とキャッシュメモリ 4 3 内の 1 D 1 及び 1 D 3 との XOR 演算を行なう。

【 0 0 8 3 】

一例として、図 8 に示すように、 $2 D 2 = 0 1 0 0 1 0$ 、 $2 D 1 = 0 1 0 0 0 1$ 、 $1 D 1 = 0 0 1 0 0 1$ 、 $1 D 3 = 0 0 1 0 1 1$ とした場合、 $2 D 4 '$ 及び $1 P "$ はそれぞれ以下の値となる。

【 0 0 8 4 】

$$\begin{aligned} 2 D 4 ' &= 2 D 2 + 2 D 1 &&= 0 1 0 0 1 0 + 0 1 0 0 0 1 \\ &= 0 0 0 0 1 1 \\ 1 P " &= 1 P ' + 1 D 1 + 1 D 3 = 0 0 0 1 1 0 + 0 0 1 0 0 1 + 0 0 1 0 1 1 \\ &= 0 0 0 1 0 0 \end{aligned}$$

【 0 0 8 5 】

ここで、 $1 P "$ については、ストライプ - 1 における消失ブロック（消失パリティ）以外の情報ブロックの XOR 演算が完了しているため、ストライプ - 1 のリカバリが完了する。また、図 8 の下段に示すように、 $1 P "$ の値（ $0 0 0 1 0 0$ ）がノード - 5 の HDD - 5 における $1 P$ （ $0 0 0 1 0 0$ ）と一致していることがわかる。

【 0 0 8 6 】

そして、図 9 の上段に示すように、各ノードの CM 4 は、HDD の未読み出しの情報ブロックを読出部 4 2 により読み出し、読み出した情報ブロックを通信部 4 1 によりリカバリ対象ブロックを有するノードへ転送する。情報ブロックを受信したノードは、キャッシュメモリ 4 3 に格納する。

【 0 0 8 7 】

図 9 の例では、ノード - 4 及びノード - 6 は、それぞれ HDD - 4 の 3 ブロック目の 2 D 3 及び HDD - 6 の 2 ブロック目の 2 P をノード - 1 に転送し、ノード - 3 は、HDD - 3 の 3 ブロック目の 3 D 1 をノード - 2 に転送する。

【 0 0 8 8 】

次いで、図 9 の下段に示すように、ノード - 1 及びノード - 2 の CM 4 は、書込部 4 4 により、キャッシュメモリ 4 3 の容量が枯渇したと判断して、ノード - 1 及びノード - 2 のそれぞれのリカバリ対象ブロックについて部分的なりカバリを実施する。このときノード - 1 及びノード - 2 のリカバリ対象ブロックに書き込まれる情報である $2 D 4 "$ 及び $3 D 3 "$ は、それぞれ $2 D 4 " = 2 D 4 ' + 2 P + 2 D 3$ 、 $3 D 3 " = 3 D 3 ' + 3 D 2 + 3 D 1$ となる。なお、ノード - 1 の書込部 4 4 は、リカバリ対象ブロックから $2 D 4 '$ を読み出してから、読み出した $2 D 4 '$ とキャッシュメモリ 4 3 内の $2 P$ 及び $2 D 3$ との XOR 演算を行なう。また、ノード - 2 の書込部 4 4 は、リカバリ対象ブロックから $3 D 3$

10

20

30

40

50

'を読み出してから、読み出した3D3'とキャッシュメモリ43内の3D2及び3D1とのXOR演算を行なう。

【0089】

一例として、図10に示すように、 $2P = 000100$ 、 $2D3 = 010011$ 、 $3D2 = 011010$ 、 $3D1 = 011001$ とした場合、 $2D4''$ 及び $3D3''$ はそれぞれ以下の値となる。

【0090】

$$\begin{aligned} 2D4'' &= 2D4' + 2P + 2D3 = 000011 + 000100 + 010011 \\ &= 010100 \end{aligned}$$

$$\begin{aligned} 3D3'' &= 3D3' + 3D2 + 3D1 = 011000 + 011010 + 011001 \\ &= 011011 \end{aligned}$$

10

【0091】

ここで、 $2D4''$ 及び $3D3''$ のいずれについても、それぞれストライプ-2及びストライプ-3における消失ブロック(消失パリティ)以外の情報ブロックのXOR演算が完了しているため、ストライプ-2及びストライプ-3のリカバリが完了する。また、図10の下段に示すように、 $2D4''$ の値(010100)及び $3D3''$ の値(011011)が、それぞれノード-5のHDD-5における $2D4$ (010100)及び $3D3$ (011011)と一致していることがわかる。

【0092】

このように、第1実施形態においては、故障した記憶装置5の消失ブロック(第2ブロック情報)の復元先は、自装置10の1以上の第1記憶装置5又は自装置10とは異なる第2ストレージ装置10にそなえられた第1記憶装置5における空き記憶領域となる。この場合、通信部41は、読出部42により第1記憶装置5から読み出し済の第1ブロック情報を、復元先である第1記憶装置5へ送信するのである。

20

【0093】

〔1-5〕ストレージシステムの動作例

次に、上述の如く構成されたストレージシステム1の動作例を、図11～図13を参照して説明する。

【0094】

〔1-5-1〕全体処理の説明

30

はじめに、図11を参照して、ストレージシステム1における全体の処理について説明する。

【0095】

図11に示すように、ストレージシステム1が正常に運用されている状態において(ステップS1、ステップS2、及びステップS2のNoルート)、いずれかのストレージ装置10が記憶装置5の故障を検出すると(ステップS2のYesルート)、処理がステップS3に移行する。ステップS3では、各ストレージ装置10によりリビルドによる障害の復旧が可能か否かが判断される。

【0096】

リビルドによる障害の復旧が可能である場合(ステップS3のYesルート)、各ストレージ装置10は、リビルド処理を実施し(ステップS4)、処理がステップS1に移行する。一方、リビルドによる障害の復旧が不可能である場合(ステップS3のNoルート)、ストレージシステム1における少なくとも1つのストレージ装置10がシステムの管理者等へエラー出力を行ない(ステップS5)、処理が終了する。

40

【0097】

なお、リビルドによる障害の復旧が不可能である場合としては、少なくとも1つのストライプにおいて消失訂正符号の訂正能力を超えた情報ブロックの消失が生じた場合が挙げられる。また、エラー出力の手法としては、システムの管理者が使用する管理者端末へエラーの発生及びエラーの内容を含むメールを送信したり、管理者端末のモニタへエラー出力を行なう等、既知の種々の手法により行なうことが可能である。なお、管理者端末は、

50

例えばスイッチ 2 及びネットワーク 3、又はスイッチ 6 を介してストレージ装置 10 と相互に通信可能に接続されている。

【0098】

〔1-5-2〕リビルド処理の説明

次に、図 12 及び図 13 を参照して、ストレージシステム 1 におけるリビルド処理（図 11 のステップ S4 参照）について説明する。

【0099】

図 12 に示すように、はじめに、ストレージシステム 1 のストレージ装置 10（CM4）は、記憶装置 5 の故障が発生した CM4 の通信部 41 からの消失ブロックに関する情報の通知に基づき、CM4 間でストライプごとのリカバリ対象ブロックを決定する（ステップ S11）。例えば図 5 の上段に示す構成の場合、各ノードの CM4 は、ノード - 5 の CM4 からの通知に基づき、消失ブロックが存在するストライプ - 1、ストライプ - 2、ストライプ - 3 について、リカバリ対象ブロックに決定する。この場合、リカバリ対象ブロックは、ノード - 1（ストライプ - 2）、ノード - 2（ストライプ - 3）、ノード - 6（ストライプ - 1）の各記憶装置 5 の代替ブロックとなる。

10

【0100】

次いで、各 CM4 は、書込部 44 によりキャッシュメモリ 43 の初期化を行なう（ステップ S12）。なお、ステップ S12 の処理は、少なくともリカバリ対象ブロックを持つストレージ装置 10 の CM4 が実施すればよい。

【0101】

また、CM4 は、読出部 42 により記憶装置 5 の未読出の情報ブロックを物理アドレスの昇順で 1 つ読み出す（ステップ S13）。なお、自装置 10 の記憶装置 5 にリカバリ対象ブロックが存在する場合、読出部 42 により読み出す情報ブロックには、当該リカバリ対象ブロックは含まれない。読出部 42 は、読み出す情報ブロックが当該リカバリ対象ブロックであれば、このブロックをスキップして次のアドレスの情報ブロックを読み出す。

20

【0102】

そして、各 CM4 は、読み出した情報ブロックが当該情報ブロックに対応するストライプの復旧先（転送先）の CM4 でキャッシュメモリ 43 に格納可能か否かを判断する（ステップ S14）。なお、この判断は、復旧先の CM4 から、キャッシュメモリ 43 の容量が逼迫したことを示す通知等を受信しているか否かの判断により行なうことができる。

30

【0103】

このような通知を受信しておらず、復旧先の CM4 で情報ブロックを格納可能である場合（ステップ S14 の Yes ルート）、処理がステップ S16 に移行する。一方、このような通知を受信しており、復旧先の CM4 で情報ブロックを格納不可能である場合（ステップ S14 の No ルート）、CM4 は、復旧先が情報ブロックを格納可能になるまで待機し（ステップ S15）、格納可能になった旨の通知を受けると、処理がステップ S16 に移行する。なお、CM4 は、この待機において、読み出した情報ブロックをキャッシュメモリ 43 に退避しておいてもよい。

【0104】

ステップ S16 では、CM4 が読み出した情報ブロックを通信部 41 により復旧先の CM4 へ転送する。

40

【0105】

また、CM4 が通信部 41 により他の CM4 から情報ブロックを受信した場合（ステップ S17 及びステップ S17 の Yes ルート）、処理が図 13 のステップ S22 に移行する。なお、情報ブロックを受信したということは、自装置 10 の記憶装置 5 にリカバリ対象ブロックが存在することを意味する。

【0106】

一方、CM4 が他の CM4 から情報ブロックを受信していない場合（ステップ S17 の No ルート）、CM4 は読出部 42 により記憶装置 5 の最終ブロックまで読み出しが完了したか否かを判断する（ステップ S18）。最終ブロックまで読み出しが完了していない

50

場合（ステップ S 1 8 の N o ルート）、処理がステップ S 1 3 に移行し、最後に読み出した情報ブロックの次のアドレスの情報ブロックを読み出す。

【 0 1 0 7 】

また、ステップ S 1 8 において、最終ブロックまで読み出しが完了した場合（ステップ S 1 8 の Y e s ルート）、自装置 1 0 の記憶装置 5 が復旧先でなければ（ステップ S 1 9 及びステップ S 1 9 の N o ルート）、リビルド処理が終了する。なお、自装置 1 0 の記憶装置 5 が復旧先でないとは、自装置 1 0 の記憶装置がリカバリ対象ブロックを持っていない場合である。

【 0 1 0 8 】

一方、自装置 1 0 の記憶装置 5 が復旧先であれば（ステップ S 1 9 の Y e s ルート）、C M 4 は、自装置 1 0 の記憶装置 5 に対するリカバリ対象ブロックのリカバリが完了しているか否かを判断する（ステップ S 2 0 ）。リカバリが完了している場合（ステップ S 2 0 の Y e s ルート）、リカバリ処理が終了する。また、リカバリが完了していない場合（ステップ S 2 0 の N o ルート）、C M 4 は、通信部 4 1 により他の C M 4 から情報ブロックを受信するまで待機し（ステップ S 2 1 ）、情報ブロックを受信すると、処理が図 1 3 のステップ S 2 2 に移行する。

10

【 0 1 0 9 】

図 1 3 に示すように、ステップ S 2 2 では、C M 4 の書込部 4 4 が、通信部 4 1 により受信した情報ブロックをキャッシュメモリ 4 3 へ格納する。そして、書込部 4 4 は、キャッシュメモリ 4 3 の容量を監視するための情報として、例えばキャッシュメモリ 4 3 の使用量を示す容量情報を更新する。また、リカバリ対象ブロックのリカバリが完了したか否かを判断するための情報として、例えば受信した（又はキャッシュメモリ 4 3 に格納した）情報ブロック数のカウント値を更新する（ステップ S 2 3 ）。なお、既述のように、キャッシュメモリの使用量についても情報ブロック数のカウント値が用いられてもよく、この場合、書込部 4 4 は、情報ブロック数のカウント値のみを更新すればよい。

20

【 0 1 1 0 】

次に、書込部 4 4 は、容量情報が閾値を超えたか否かを判断し（ステップ S 2 4 ）、超えていない場合（ステップ S 2 4 の N o ルート）、処理が図 1 2 のステップ S 1 8 に移行する。

【 0 1 1 1 】

一方、容量情報が閾値を超えた場合（ステップ S 2 4 の Y e s ルート）、書込部 4 4 は、容量情報が閾値を超えたこと（キャッシュ不可の旨）を通信部 4 1 を介して他の C M 4 に通知する（ステップ S 2 5 ）。そして、書込部 4 4 は、リカバリ対象ブロックに情報が格納済みであるか否かを判断する（ステップ S 2 6 ）。

30

【 0 1 1 2 】

リカバリ対象ブロックに情報が格納済みである場合（ステップ S 2 6 の Y e s ルート）、書込部 4 4 は、リカバリ対象ブロックの情報を読み出す。そして、書込部 4 4 は、読み出した情報とキャッシュメモリ 4 3 内の情報ブロックとの X O R 演算を実行し（ステップ S 2 7 ）、処理がステップ S 2 9 に移行する。一方、リカバリ対象ブロックに情報が格納されていない場合（ステップ S 2 6 の N o ルート）、書込部 4 4 は、キャッシュメモリ 4 3 内の情報ブロックの X O R 演算を実行し（ステップ S 2 8 ）、処理がステップ S 2 9 に移行する。

40

【 0 1 1 3 】

ステップ S 2 9 では、書込部 4 4 は、ステップ S 2 7 又はステップ S 2 8 における X O R 演算結果をリカバリ対象ブロックへ書き込む。なお、キャッシュメモリ 4 3 の内容はリカバリ対象ブロックへ反映されたため、書込部 4 4 はキャッシュメモリ 4 3 をクリアする。また、書込部 4 4 は、キャッシュ可能の旨を通信部 4 1 を介して他の C M 4 に通知する（ステップ S 3 0 ）。

【 0 1 1 4 】

そして、書込部 4 4 は、カウント値が閾値に達したか否かを判断し（ステップ S 3 1 ）。

50

、達していない場合（ステップ S 3 1 の N o ルート）、つまりリカバリが完了していない場合、処理が図 1 2 のステップ S 1 8 に移行する。

【 0 1 1 5 】

一方、カウント値が閾値に達した場合（ステップ S 3 1 の Y e s ルート）、書込部 4 4 は、リカバリ対象ブロックのリカバリが完了したと判断し（ステップ S 3 2 ）、処理が図 1 2 のステップ S 1 8 に移行する。

【 0 1 1 6 】

このように、リビルド処理では、リカバリ対象ブロックを持たない C M 4 では読出処理（ステップ S 1 3 ~ S 1 8 ）が主に実行され、リカバリ対象ブロックを持つ C M 4 では読出処理と書込処理（ステップ S 1 7、ステップ S 2 0 ~ S 3 2 ）とが実施される。

10

【 0 1 1 7 】

読出処理は、シーケンシャルリードによりスループットを向上させているため、読出処理の最中に書込処理が頻発すると、記憶装置 5 のアクセス先が変化してシーケンシャルなアクセスが阻害されてしまう。そこで、リカバリ対象ブロックを持つ C M 4 は、受信した情報ブロックをキャッシュメモリ 4 3 に蓄積し、一括して書込処理を行なうことで、書込処理の実行頻度を低下させることができるのである。

【 0 1 1 8 】

なお、図 1 2 のステップ S 1 7 における情報ブロックの受信確認の処理は、ステップ S 1 3 の前（ステップ S 1 8 の N o ルートからの合流後）等、任意のタイミングで行なわれてよい。また、図 1 3 のステップ S 2 3 における情報ブロック数のカウント値の更新処理は、ステップ S 2 9 の後等に行なわれてもよい。

20

【 0 1 1 9 】

〔 2 〕 第 2 実施形態

上述した第 1 実施形態に係るストレージシステム 1 では、リカバリ対象ブロックに運用記憶装置 5 の代替ブロック（空き領域）が用いられるものとして説明したが、これに限定されるものではない。

【 0 1 2 0 】

例えばストレージシステム 1 は、図 1 4 の上段に示すように、運用中のストレージ装置 1 0 （ノード - 1 ~ ノード - 6 ）の他に、待機用のストレージ装置 1 0 （ノード - 7 ）をそなえてもよい。このとき、ストレージシステム 1 は、待機用のストレージ装置 1 0 のホットスワップディスクである H D D - 7 上に、故障した記憶装置 5 （ H D D - 5 ）のデータを復元することができる。

30

【 0 1 2 1 】

以下、第 2 実施形態に係るストレージシステム 1 について説明する。なお、ストレージシステム 1 及びストレージ装置 1 0 の構成及び機能は、特に言及しない限り、第 1 実施形態と基本的に同様とすることができる。

【 0 1 2 2 】

第 2 実施形態に係るストレージシステム 1 では、記憶装置 5 が故障した場合の代替となる記憶装置 5 （ H D D - 7 ）を予め用意し、ホットスワップとしてシステムに組み込まれている。例えばノード - 5 の H D D - 5 が故障した場合、図 1 4 の上段に示すように、運用中のノード - 1 ~ ノード - 4 及びノード - 6 は、それぞれ情報ブロックを待機用のノード - 7 へ送信する。情報ブロックは複数のノードから順次送られてくるため、情報ブロックを受信したノード - 7 は、各ストライプについてキャッシュメモリ 4 3 を用いて順次リカバリ対象ブロックへの書込処理（ X O R 演算）を行なう。

40

【 0 1 2 3 】

このように、運用中のストレージ装置 1 0 では、消失ブロックのデータを復旧するための読出処理（シーケンシャルリード）を行ない、読み出した情報ブロックを待機用のストレージ装置 1 0 へ送信するだけでよく、書込処理は発生しない。一方、待機用のストレージ装置 1 0 では、情報ブロックを受信し、受信した情報ブロックの記憶装置 5 への書込処理（ランダムライト）を行なうだけでよい。

50

【 0 1 2 4 】

以上のように、第2実施形態に係るストレージシステム1では、復旧用の情報ブロックの読出要求と、復旧先での情報ブロックの書込要求とが別々の記憶装置5に発行されるため、シーケンシャルリードが中断されず、スループットをさらに向上させることができる。

【 0 1 2 5 】

このため、第2実施形態に係るストレージシステム1では、運用中のストレージ装置10(CM4)は、例えば図4に示す機能のうち、少なくとも通信部41及び読出部42の機能をそなえていけばよい。また、待機用のストレージ装置10(CM4)は、例えば図4に示す機能のうち、少なくとも通信部41、キャッシュメモリ43、及び書込部44の機能をそなえていけばよい。

10

【 0 1 2 6 】

なお、図14に示す例では、ホットスワップディスクであるHDD-7が待機用のストレージ装置10にそなえられるものとして示したが、これに限定されるものではない。例えばHDD-7は、運用中のストレージ装置10(ノード-1~ノード-4及びノード-6のいずれか)に追加してそなえられてもよい。この場合、運用中のストレージ装置10は、図4に示す通信部41、読出部42、キャッシュメモリ43、及び書込部44の全ての機能をそなえればよい。そして、当該ストレージ装置10は、例えばHDD-7以外のHDDから読出部42により情報ブロックを読み出してキャッシュメモリ43に格納し、キャッシュメモリ43内の情報ブロックを書込部44によりHDD-7のリカバリ対象ブロックへ反映すればよい。

20

【 0 1 2 7 】

なお、図14の例において、各ストライプの消失ブロックの復旧の進捗は、各記憶装置5での読出順序に関連して、図14の下段に示すように、先頭の情報ブロック、中央付近の情報ブロック、末尾の情報ブロックの順に、段階的に(この場合3段階で)行なわれる。

【 0 1 2 8 】

また、第2実施形態に係るストレージシステム1では、上述のように運用中のストレージ装置10における記憶装置5の代替ブロックは使用されない。換言すれば、第2実施形態に係るストレージシステム1としては、代替ブロックを持たない、例えば通常のRAID5やRAID6等を採用したストレージシステムを用いることもできる。

30

【 0 1 2 9 】

例えば図15の上段に示すように、消失訂正符号としてRAID5を採用したストレージシステム1は、図14と比較して、HDD-1が2D5を格納し、HDD-2が3D5を格納し、HDD-6が1D5を格納している。このような構成であっても、復旧用の情報ブロックはノード-7へ送信されるため、各ストレージ装置10は、図14を参照した説明と同様の処理を行なうことができる。

【 0 1 3 0 】

ところで、第2実施形態においては、消失ブロックに対応する全てのストライプについて、復旧用の情報ブロックが待機用の記憶装置5(ストレージ装置10)に送信されることになる。このため、待機用のストレージ装置10(CM4)では、大量の書込処理が発生することになる。

40

【 0 1 3 1 】

そこで、待機用のストレージ装置10は、第1実施形態と比較して大容量のキャッシュメモリ43をそなえることが好ましい。これにより、1つのリカバリ対象ブロックについて或る程度の数の情報ブロックのXOR演算結果をまとめて反映できるため、1ブロック当たりの書込処理の発生頻度を低減させ、記憶装置5の処理負荷を低減させることができる。また、キャッシュメモリ43の容量の逼迫により運用中のストレージ装置10で送信待ちが発生する頻度も低減させることができ、リカバリ性能を向上させることができる。

【 0 1 3 2 】

50

又は、待機用のストレージ装置 10 は、キャッシュメモリ 43 の容量が逼迫するよりも早いタイミングで、キャッシュメモリ 43 内の情報ブロックをリカバリ対象ブロックに反映してもよい。このタイミングとしては、例えば所定期間ごとが挙げられる。或いはキャッシュメモリ 43 の使用量の閾値を第 1 実施形態よりも小さい値とすることで、キャッシュメモリ 43 内の情報ブロックをパージする頻度を上げてよい。

【0133】

これにより、待機用の記憶装置 5 における 1 ブロック当たりの書込処理の発生頻度は増加するものの、キャッシュメモリ 43 の容量が逼迫する前にキャッシュメモリ 43 内の情報ブロックがリカバリ対象ブロックに反映される。従って、運用中のストレージ装置 10 で送信待ちが発生する頻度を低減させる、或いは無くすることができ、リカバリ性能を向上させることができる。

10

【0134】

このように、第 2 実施形態においては、故障した記憶装置 5 の消失ブロック（第 2 ブロック情報）の復元先は、故障した第 2 記憶装置 5 の代替となる第 3 記憶装置 5 における空き記憶領域となる。この場合、通信部 41 は、読出部 42 により第 1 記憶装置 5 から読み出し済の第 1 ブロック情報を、第 3 記憶装置 5 へ送信するのである。

【0135】

次に、図 16 ~ 図 18 を参照して、第 2 実施形態に係るストレージシステム 1 の動作例を説明する。なお、図 16 ~ 図 18 において、第 1 実施形態に係る図 12 及び図 13 に示す符号と同一の符号を付した処理は、図 12 及び図 13 に示す処理と基本的に同様の処理であるため、重複した説明を省略する。

20

【0136】

第 2 実施形態に係るストレージシステム 1 では、運用ストレージ装置 10 と待機ストレージ装置 10 とで処理が異なる。

【0137】

例えば図 16 に示すように、運用ストレージ装置 10 の処理は、図 12 に示すストレージ装置 10 の処理から、ステップ S12、ステップ S17、及びステップ S19 ~ S21 の書込処理に係る処理を省略したものとなる。

【0138】

また、図 17 に示すように、待機ストレージ装置 10 の処理は、図 13 に示すストレージ装置 10 の処理に対して、以下の処理を追加又は変更したものとなる。

30

【0139】

例えば図 17 に示すように、待機ストレージ装置 10 の CM4 は、リビルド処理において、キャッシュメモリ 43 を初期化し（ステップ S41）、運用ストレージ装置 10 から情報ブロックを受信するまで待機する（ステップ S42）。また、ステップ S24 において容量情報が閾値を超えていない場合（ステップ S24 の No ルート）、及びステップ S31 においてカウント値が閾値に達していない場合（ステップ S31 の No ルート）、処理がステップ S42 に移行する。その他の点については、待機ストレージ装置 10 の処理は図 13 に示すストレージ装置 10 の書込処理の動作と基本的に同様である。

【0140】

なお、待機ストレージ装置 10 は、上述のようにキャッシュメモリ 43 の容量が逼迫するよりも早いタイミングでキャッシュメモリ 43 内の情報ブロックをリカバリ対象ブロックに反映してもよい。この場合、待機ストレージ装置 10 の処理は、図 18 に示すように、図 17 からステップ S24、ステップ S25、及びステップ S30 の処理を省略し、ステップ S23 の処理を情報ブロック数のカウント値の更新のみを行なうステップ S51 に置き換えたものとすることができる。

40

【0141】

以上のように、第 2 実施形態に係るストレージシステム 1 によっても、第 1 実施形態と同様の効果を奏することができるほか、CM4 が記憶装置 5 から情報ブロックをよりシームレスに読み出すことができ、スループットをさらに向上させることができる。

50

【 0 1 4 2 】

〔 3 〕 ハードウェア構成例

図 19 に例示するように、上述した第 1 及び第 2 実施形態に係るストレージ装置 10 の CM4 は、図 1 に示す CPU4a 及びメモリ 4b に加えて、記憶部 4f、インタフェース部 4g、入出力部 4h、及び読取部 4i をそなえることができる。

【 0 1 4 3 】

記憶部 4f は、種々のデータやプログラム等を格納するハードウェアである。記憶部 4f としては、例えば HDD 等の磁気ディスク装置、SSD 等の半導体ドライブ装置、フラッシュメモリや ROM 等の不揮発性メモリ等の各種装置が挙げられる。

【 0 1 4 4 】

例えば記憶部 4f は、ストレージ装置 10 (CM4) の各種機能の全部もしくは一部を実現するリカバリプログラム 40 を格納することができる。CPU4a は、例えば記憶部 4f に格納されたリカバリプログラム 40 をメモリ 4b 等の記憶装置に展開して実行することにより、故障した記憶装置 5 の消失ブロックをリカバリ (復旧) することでリビルドを行なう上述したストレージ装置 10 (CM4) の機能を実現することができる。

【 0 1 4 5 】

インタフェース部 4g は、有線又は無線による、スイッチ 2 及び 6、並びに記憶装置 5 等との間の接続及び通信の制御等を行なう通信インタフェースである。なお、図 1 に示す IF4c ~ 4e は、インタフェース部 4g の一例である。

【 0 1 4 6 】

入出力部 4h は、マウス、キーボード、タッチパネル、音声操作のためのマイク等の入力装置 (操作部)、並びにディスプレイ、スピーカ、及びプリンタ等の出力装置 (出力部、表示部) の少なくとも一方を含むことができる。例えば入力装置は、管理者等による各種操作やデータの入力等の作業に用いられてよく、出力装置は、各種通知等の出力に用いられてよい。

【 0 1 4 7 】

読取部 4i は、コンピュータ読取可能な記録媒体 4j に記録されたデータやプログラムを読み出す装置である。この記録媒体 4j にはリカバリプログラム 40 が格納されてもよい。

【 0 1 4 8 】

なお、記録媒体 4j としては、例えばフレキシブルディスク、CD、DVD、ブルーレイディスク等の光ディスクや、USBメモリやSDカード等のフラッシュメモリ等の非一時的な記録媒体が挙げられる。なお、CDとしては、CD-ROM、CD-R、CD-RW等が挙げられる。また、DVDとしては、DVD-ROM、DVD-RAM、DVD-R、DVD-RW、DVD+R、DVD+RW等が挙げられる。

【 0 1 4 9 】

上述したストレージ装置 10 のハードウェア構成は例示である。従って、ストレージ装置 10 内でのハードウェアの増減 (例えば任意のブロックの追加や省略)、分割、任意の組み合わせでの統合、バスの追加又は省略等は適宜行なわれてもよい。

【 0 1 5 0 】

〔 4 〕 その他

以上、本発明の好ましい実施形態について詳述したが、本発明は、かかる特定の実施形態に限定されるものではなく、本発明の趣旨を逸脱しない範囲内において、種々の変形、変更して実施することができる。

【 0 1 5 1 】

例えば、図 4 に示すストレージ装置 10 (CM4) の各機能ブロックは、任意の組み合わせで併合してもよく、分割してもよい。

【 0 1 5 2 】

また、ここまで、ストレージシステム 1 が複数のストレージ装置 10 (筐体) をそなえるクラスタ構成の分散ストレージシステムであるものとして説明したが、これに限定され

10

20

30

40

50

るものではない。例えばストレージシステム 1 は、図 20 に示すように、単一のストレージ装置 10 内に J B O D (Just a Brunck of Disks) 50 等の形態で実装されたディスクアレイをそなえてもよい。

【0153】

J B O D 50 は、複数の記憶装置 5 を搭載し、複数の記憶装置 5 に接続された I F 50 a を介して C M 4 の I F 4 d と通信を行なうことができる。この場合、C M 4 は、読出部 42 により複数の記憶装置 5 からシーケンシャルに読み出した情報ブロックを、書込部 44 によりキャッシュメモリ 43 に格納すればよい。そして、C M 4 は、キャッシュメモリ 43 の容量が逼迫した場合に、キャッシュメモリ 43 内の情報ブロックを復旧先の記憶装置 5 のリカバリ対象ブロックに反映すればよい。

10

【0154】

なお、図 20 に例示する構成において、C M 4 は図 4 に示す通信部 41、読出部 42、キャッシュメモリ 43、及び書込部 44 をそなえることができる。例えば読出部 42 は、記憶装置 5 から情報ブロックをシーケンシャルに読み出して通信部 41 に出力する。この情報ブロックは、通信部 41 を介して書込部 44 によりキャッシュメモリ 43 に格納され、リカバリ対象ブロックを有する記憶装置 5 へ書き込まれる。

【0155】

このように、図 20 及び図 4 に例示する構成において、キャッシュメモリ 43 は、通信部 41 から入力される第 1 ブロック情報を保持する保持部の一例であるといえる。また、書込部 44 は、通信部 41 から入力された第 1 ブロック情報に基づき、第 2 ブロック情報を段階的に復元する復元部の一例であるといえる。

20

【0156】

なお、ストレージシステム 1 は、複数のストレージ装置 10 (クラスタノード) 内にそれぞれ J B O D 50 等により複数の記憶装置 5 を搭載する、図 1 及び図 20 の双方の形態を持つ構成であってもよい。

【0157】

〔5〕付記

以上の第 1 及び第 2 実施形態に関し、更に以下の付記を開示する。

【0158】

(付記 1)

1 以上の第 1 記憶装置と、

前記 1 以上の第 1 記憶装置が記憶する複数の第 1 ブロック情報であって、故障した第 2 記憶装置が記憶する複数の第 2 ブロック情報の復元に用いられる前記複数の第 1 ブロック情報を、前記 1 以上の第 1 記憶装置から記憶領域のアドレス順に読み出す読出部と、

前記読出部により前記 1 以上の第 1 記憶装置から読み出し済の第 1 ブロック情報を、前記複数の第 2 ブロック情報を段階的に復元するために、前記複数の第 2 ブロック情報の復元先へ出力する出力部と、をそなえることを特徴とする、ストレージ装置。

30

【0159】

(付記 2)

前記ストレージ装置とは異なる第 1 ストレージ装置、又は、前記出力部、から入力された第 1 ブロック情報に基づき、前記第 2 ブロック情報を段階的に復元する復元部をさらにそなえることを特徴とする、付記 1 記載のストレージ装置。

40

【0160】

(付記 3)

前記復元部は、前記入力された第 1 ブロック情報及び前記第 2 ブロック情報の復元先が記憶する情報を用いて、消失訂正符号に基づく演算を行ない、演算結果を前記第 2 ブロック情報の復元先へ書き込むことを特徴とする、付記 2 記載のストレージ装置。

【0161】

(付記 4)

前記第 1 ストレージ装置又は前記出力部から入力される第 1 ブロック情報を保持する保

50

持部をさらにそなえ、

前記復元部は、所定のタイミングで、前記保持部が保持する 1 以上の第 1 ブロック情報に基づき、前記第 2 ブロック情報を段階的に復元することを特徴とする、付記 2 又は付記 3 記載のストレージ装置。

【 0 1 6 2 】

(付記 5)

前記第 2 ブロック情報の復元先は、前記 1 以上の第 1 記憶装置、又は、前記ストレージ装置とは異なる第 2 ストレージ装置にそなえられた第 1 記憶装置における空き記憶領域であり、

前記出力部は、前記読出部により前記第 1 記憶装置から読み出し済の第 1 ブロック情報を、前記復元先である第 1 記憶装置へ送信することを特徴とする、付記 1 ~ 4 のいずれか 1 項記載のストレージ装置。

【 0 1 6 3 】

(付記 6)

前記第 2 ブロック情報の復元先は、前記故障した第 2 記憶装置の代替となる第 3 記憶装置における空き記憶領域であり、

前記出力部は、前記読出部により前記第 1 記憶装置から読み出し済の第 1 ブロック情報を、前記第 3 記憶装置へ送信することを特徴とする、付記 1 ~ 4 のいずれか 1 項記載のストレージ装置。

【 0 1 6 4 】

(付記 7)

複数の記憶装置と、

前記複数の記憶装置のうちの複数の第 1 記憶装置の各々が記憶する複数の第 1 ブロック情報に基づき、前記複数の記憶装置のうちの故障した第 2 記憶装置が記憶する複数の第 2 ブロック情報を復元する 1 以上のストレージ装置と、をそなえ、

前記 1 以上のストレージ装置は、

複数の第 1 記憶装置の各々から、前記複数の第 1 ブロック情報を記憶領域のアドレス順に読み出し、

前記読み出しの処理において前記複数の第 1 記憶装置の各々から読み出し済の第 1 ブロック情報に基づき、前記複数の第 2 ブロック情報を段階的に復元することを特徴とする、ストレージシステム。

【 0 1 6 5 】

(付記 8)

コンピュータに、

1 以上の第 1 記憶装置が記憶する複数の第 1 ブロック情報であって、故障した第 2 記憶装置が記憶する複数の第 2 ブロック情報の復元に用いられる前記複数の第 1 ブロック情報を、前記 1 以上の第 1 記憶装置から記憶領域のアドレス順に読み出し、

前記読み出しの処理において前記 1 以上の第 1 記憶装置から読み出し済の第 1 ブロック情報を、前記複数の第 2 ブロック情報を段階的に復元するために、前記複数の第 2 ブロック情報の復元先へ出力する、

処理を実行させることを特徴とする、リカバリプログラム。

【 0 1 6 6 】

(付記 9)

複数の記憶装置と、前記複数の記憶装置に対する制御を行なう 1 以上のストレージ装置とをそなえるストレージシステムにおけるリカバリ方法であって、

前記 1 以上のストレージ装置は、

前記複数の記憶装置のうちの複数の第 1 記憶装置の各々が記憶する複数の第 1 ブロック情報に基づき、前記複数の記憶装置のうちの故障した第 2 記憶装置が記憶する複数の第 2 ブロック情報を復元し、

前記復元の処理において、

10

20

30

40

50

前記 1 以上のストレージ装置により、

複数の第 1 記憶装置の各々から、前記複数の第 1 ブロック情報を記憶領域のアドレス順に読み出し、

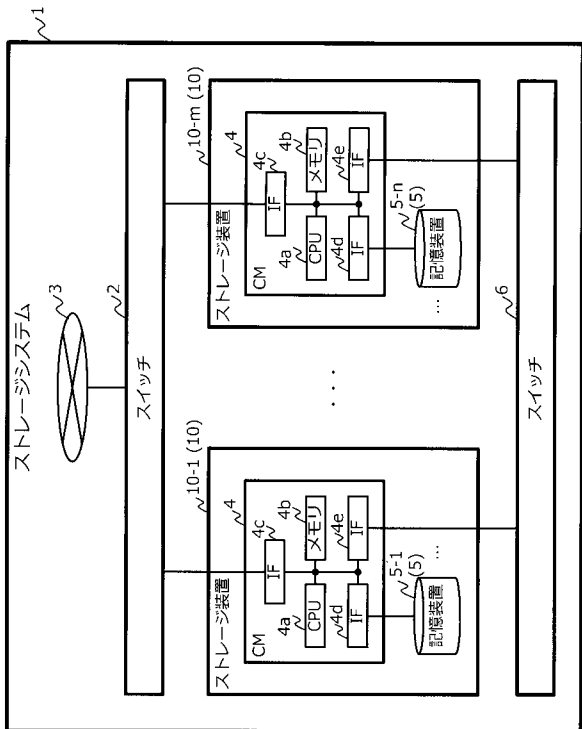
前記読み出しの処理において前記複数の第 1 記憶装置の各々から読み出し済の第 1 ブロック情報に基づき、前記複数の第 2 ブロック情報を段階的に復元する、ことを特徴とする、リカバリ方法。

【符号の説明】

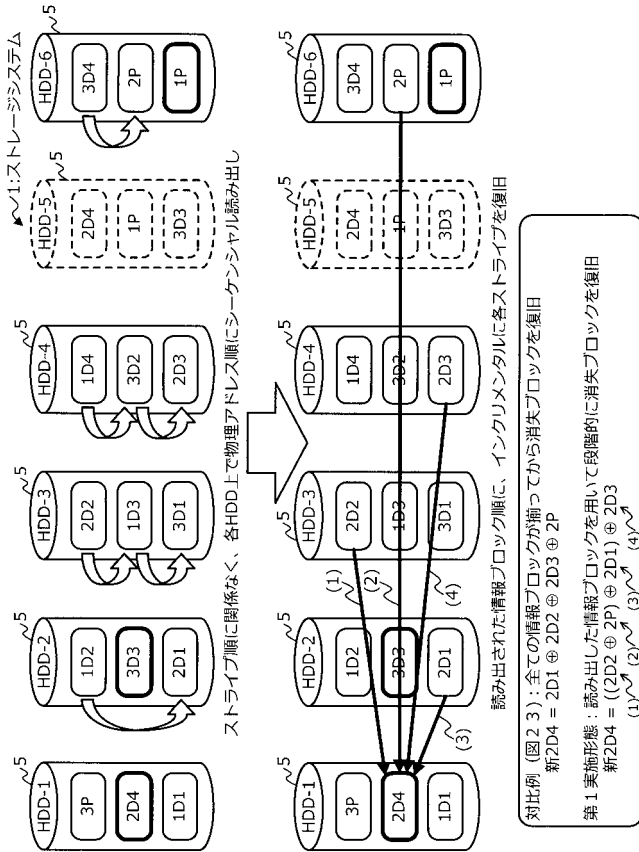
【 0 1 6 7 】

- | | | |
|------------------------|-------------|----|
| 1 | ストレージシステム | |
| 2, 6 | スイッチ | 10 |
| 3 | ネットワーク | |
| 4 | コントローラモジュール | |
| 4 a | C P U | |
| 4 b | メモリ | |
| 4 c ~ 4 e, 5 0 a | インタフェース | |
| 4 f | 記憶部 | |
| 4 g | インタフェース部 | |
| 4 h | 入出力部 | |
| 4 i | 読取部 | |
| 4 j | 記録媒体 | 20 |
| 5, 5 - 1 ~ 5 - n | 記憶装置 | |
| 1 0, 1 0 - 1 ~ 1 0 - m | ストレージ装置 | |
| 4 0 | リカバリプログラム | |
| 4 1 | 通信部 | |
| 4 2 | 読出部 | |
| 4 3 | キャッシュメモリ | |
| 4 4 | 書込部 | |
| 5 0 | J B O D | |

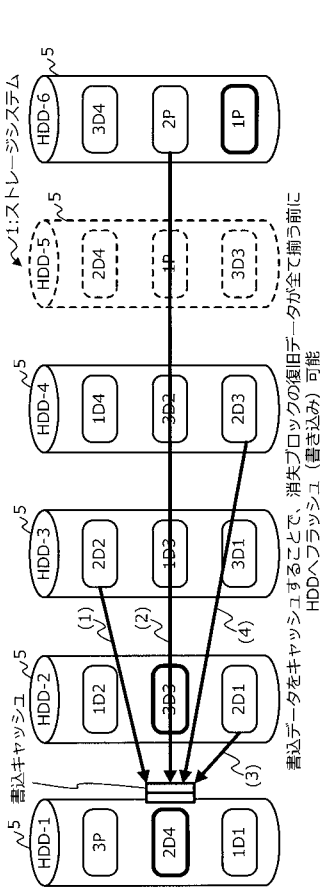
【 図 1 】



【 図 2 】

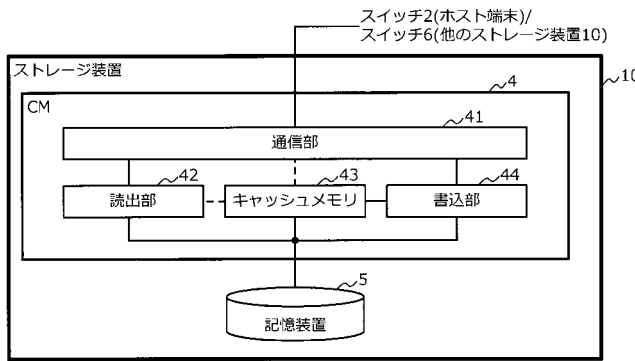


【 図 3 】

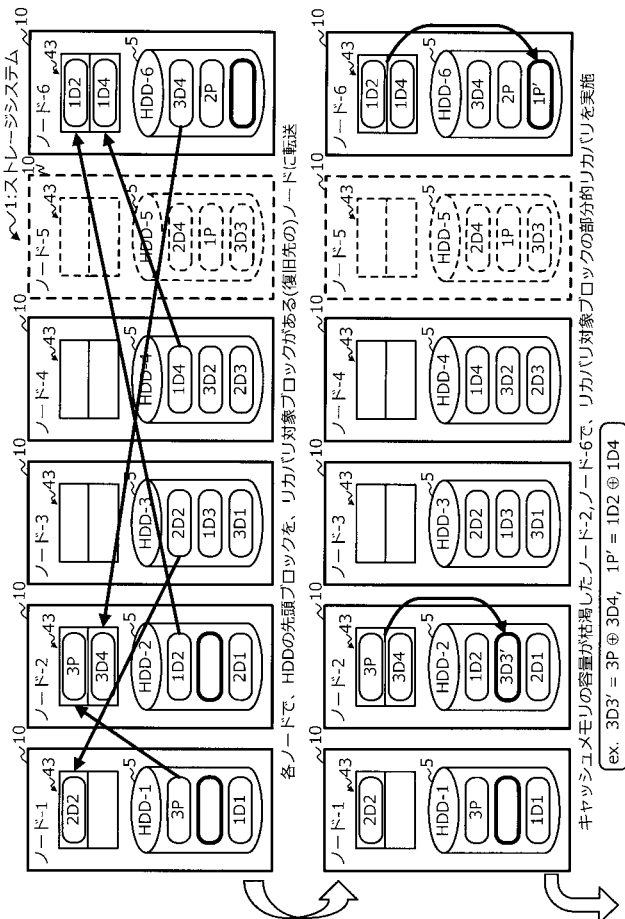


処理	ストライプ2の2D4復旧の進捗
(1) HDD-3から2D2を読み出し、HDD-1の新2D4の書き込キャッシュに反映	2D4 = 2D2
(2) HDD-6から2Pを読み出し、HDD-1の新2D4の書き込キャッシュに反映	2D4 = 2D4 ⊕ 2P
新2D4の書き込キャッシュをフラッシュ	
(3) HDD-2から2D1を読み出し、HDD-1の新2D4を書込キャッシュにリロードし、2D1をキャッシュに反映	2D4 = 2D4 ⊕ 2D1
(4) HDD-4から2D3を読み出し、HDD-1の新2D4の書き込キャッシュに反映	2D4 = 2D4 ⊕ 2D3
新2D4の書き込キャッシュをフラッシュ	

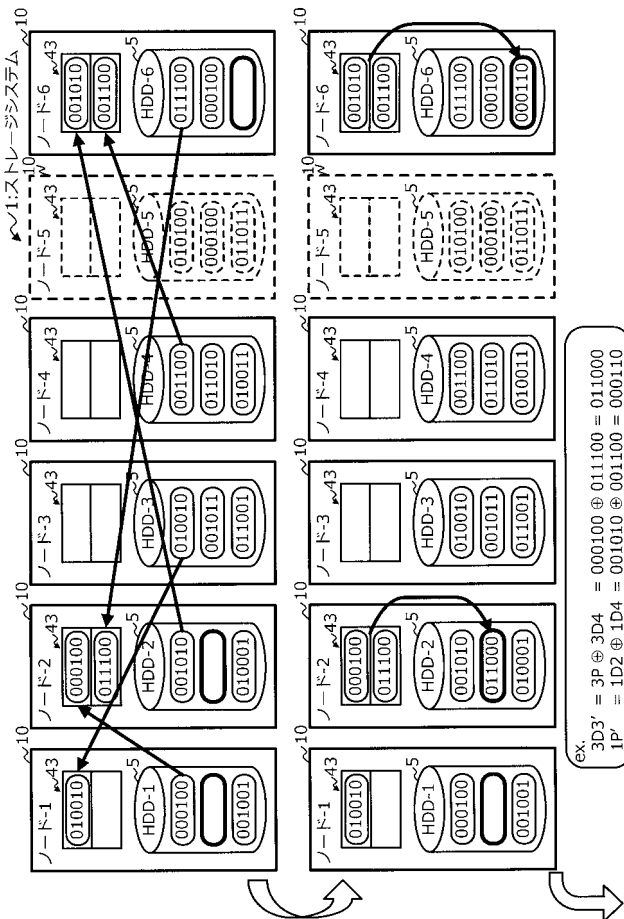
【 図 4 】



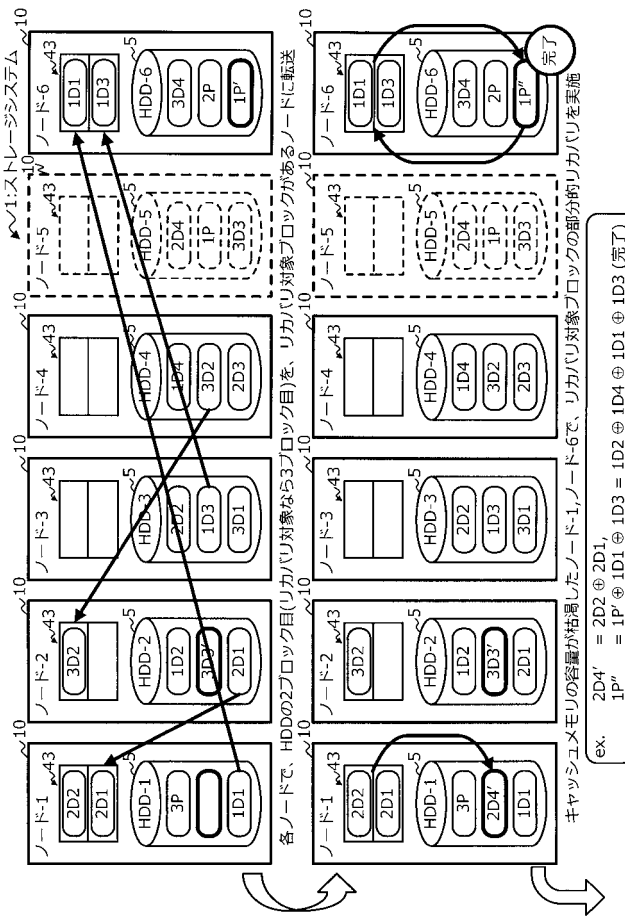
【 図 5 】



【 図 6 】



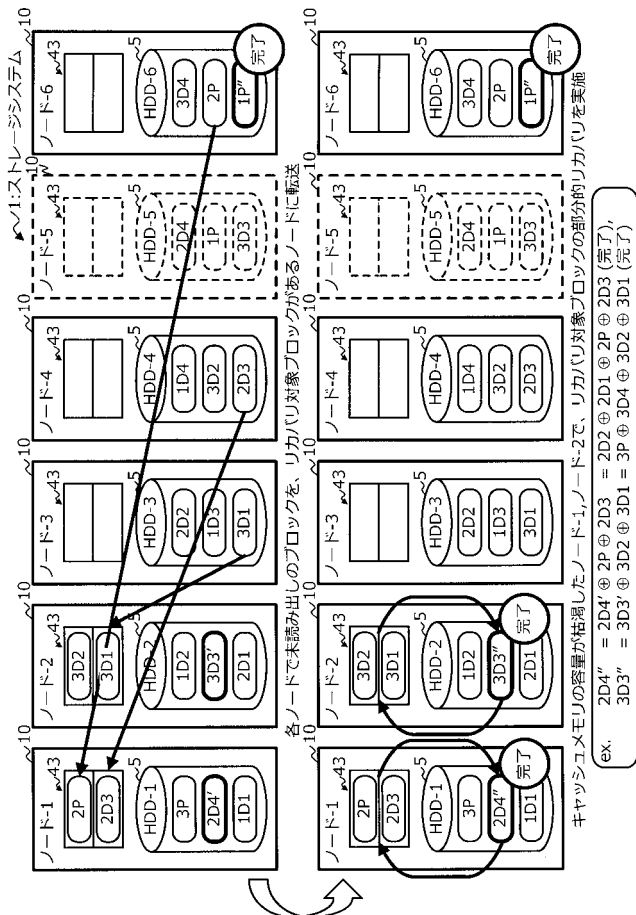
【 図 7 】



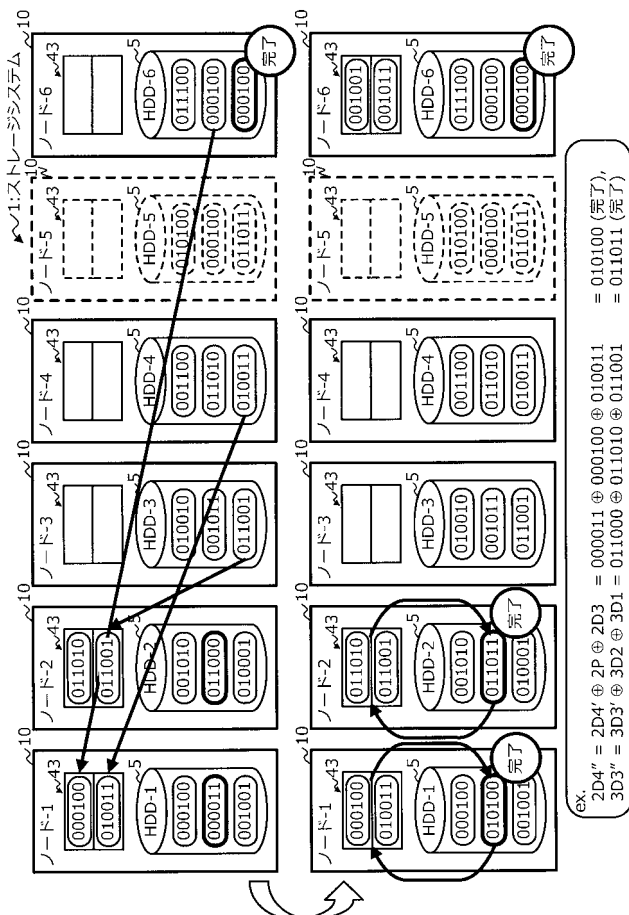
【 図 8 】



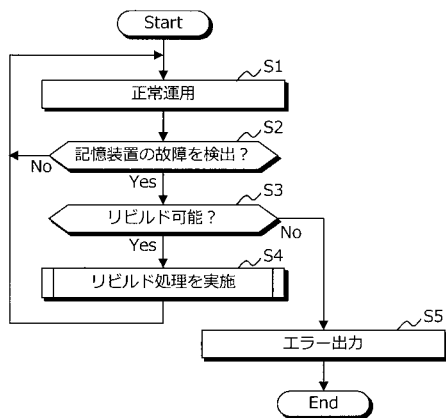
【図 9】



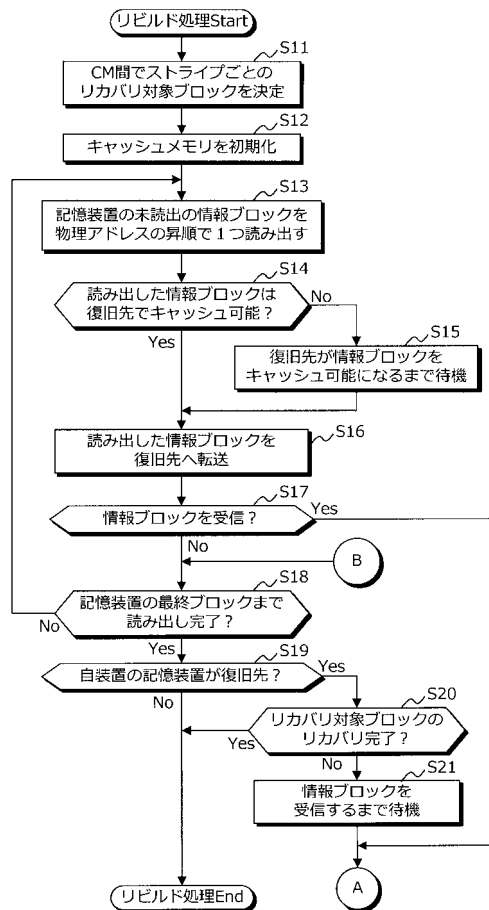
【図 10】



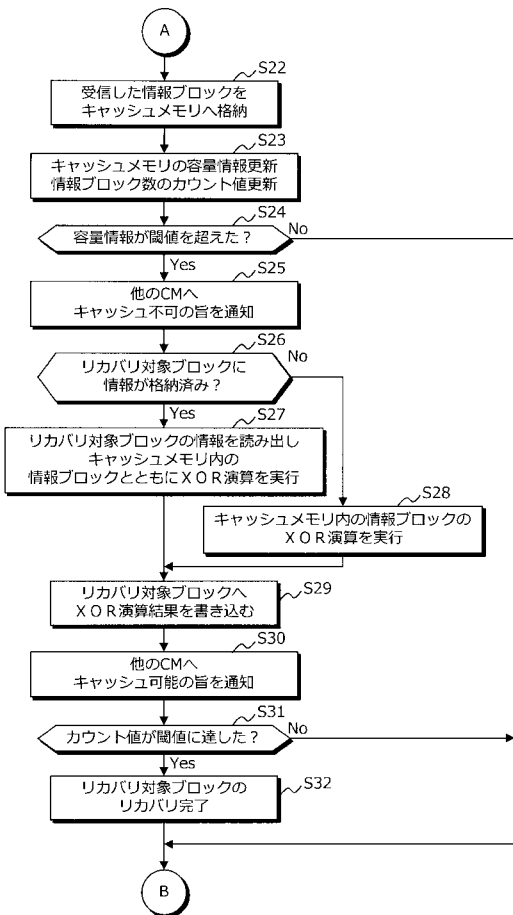
【図 11】



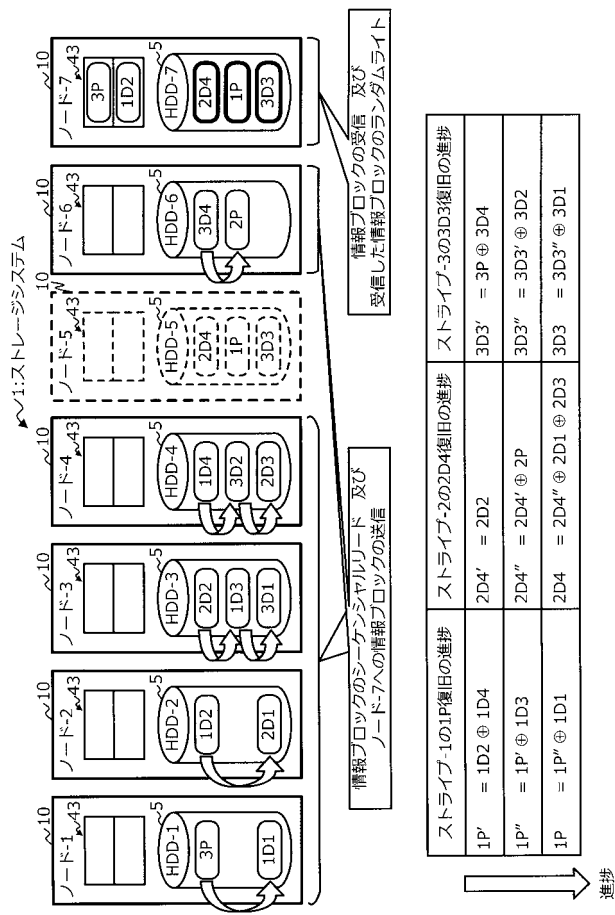
【図 12】



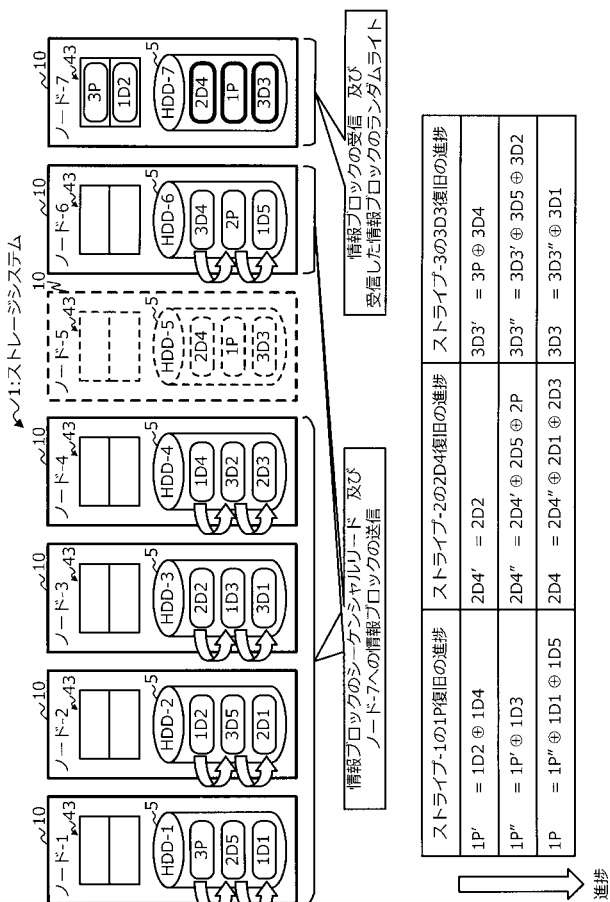
【図13】



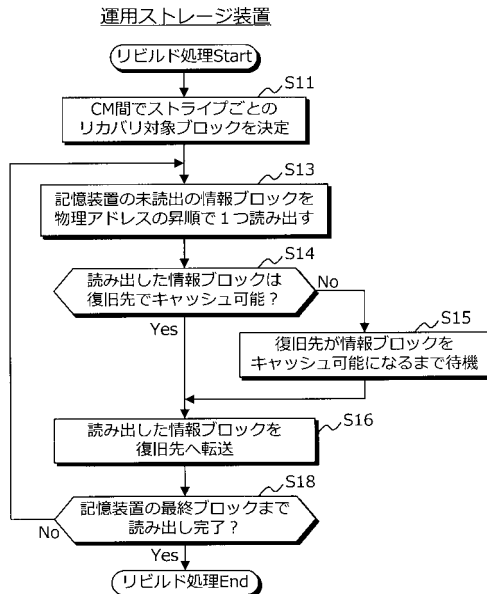
【図14】



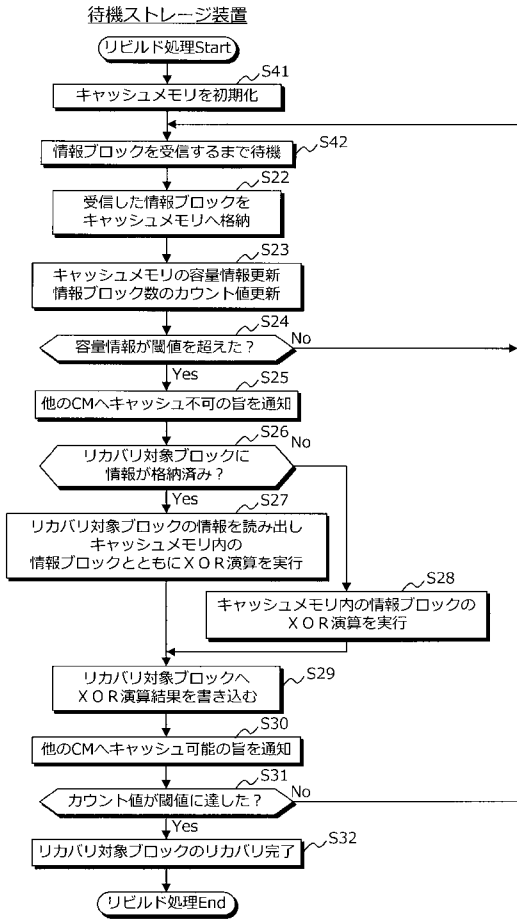
【図15】



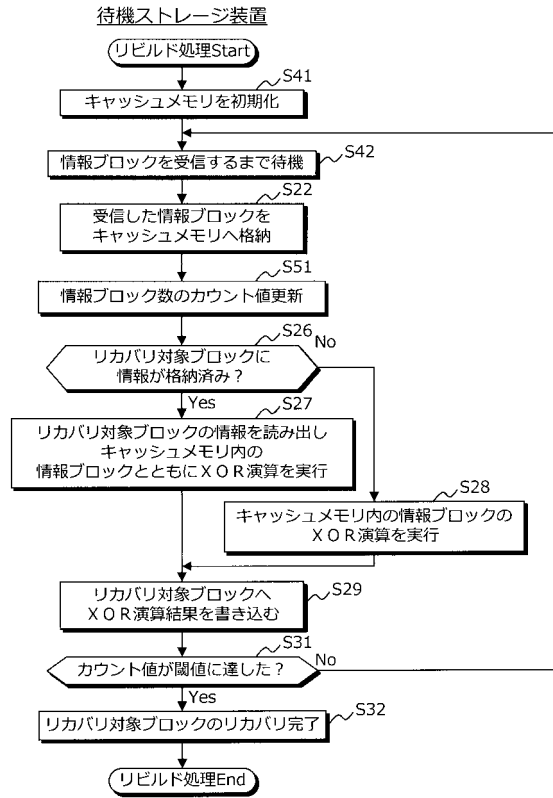
【図16】



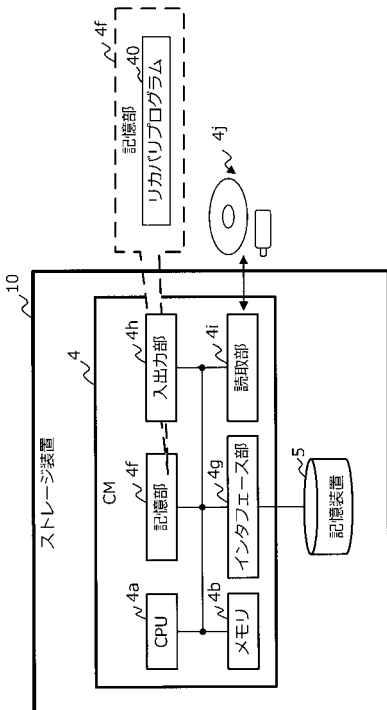
【 図 1 7 】



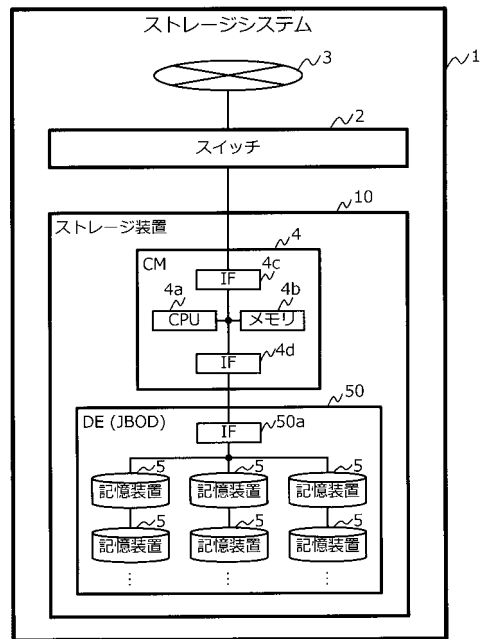
【 図 1 8 】



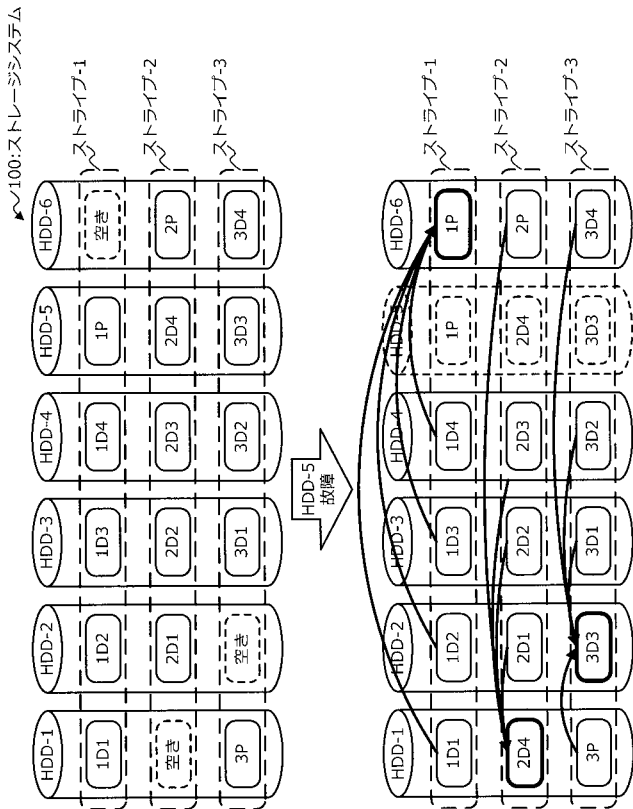
【 図 1 9 】



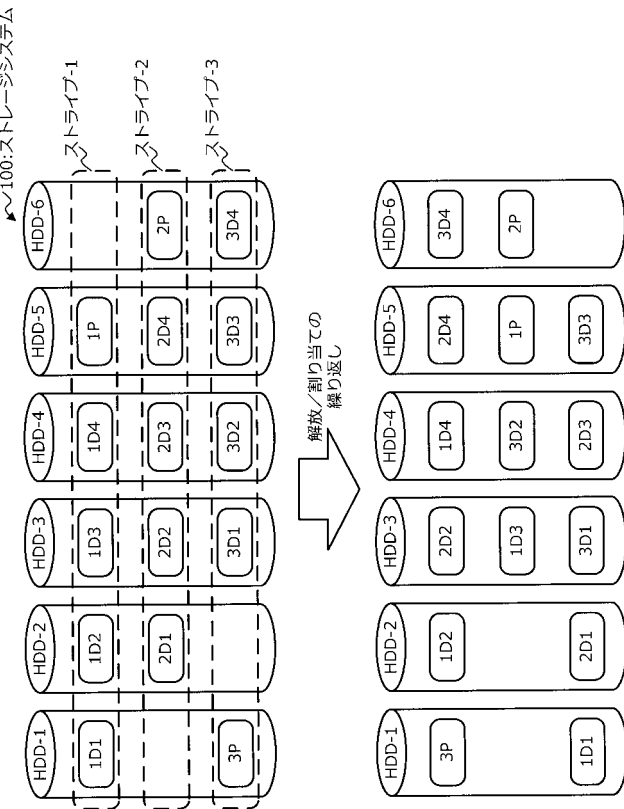
【 図 2 0 】



【 図 2 1 】

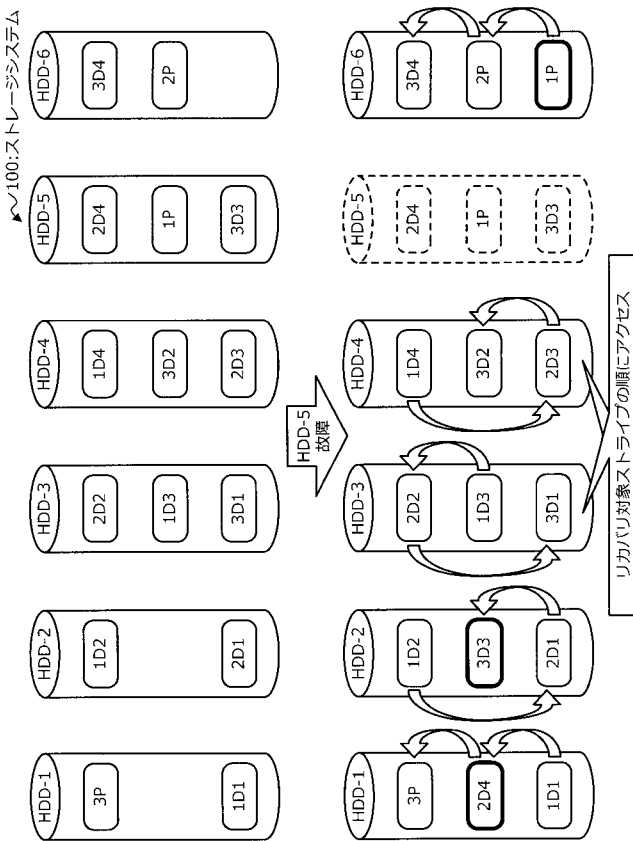


【 図 2 2 】



解放/割り当ての
繰り返し

【 図 2 3 】



リカバリ対象ストレージの順にアクセス