

(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2006-259894

(P2006-259894A)

(43) 公開日 平成18年9月28日(2006.9.28)

(51) Int. Cl.	F I	テーマコード (参考)
G06F 3/06 (2006.01)	G06F 3/06 304B	5B001
G06F 11/08 (2006.01)	G06F 11/08 310B	5B018
G06F 12/16 (2006.01)	G06F 12/16 320L	5B065

審査請求 有 請求項の数 11 O L (全 28 頁)

(21) 出願番号	特願2005-73669 (P2005-73669)	(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
(22) 出願日	平成17年3月15日 (2005.3.15)	(74) 代理人	100074099 弁理士 大菅 義之
		(74) 代理人	100067987 弁理士 久木元 彰
		(72) 発明者	望月 信哉 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
		(72) 発明者	伊藤 実希夫 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

最終頁に続く

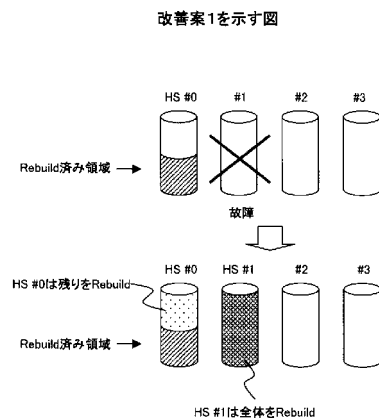
(54) 【発明の名称】 ストレージ制御装置および方法

(57) 【要約】

【課題】 パリティを用いてデータの冗長化を行うストレージ装置において、データ/パリティを格納する2つの記憶装置が故障したときの再構築処理を改善して、効率よく冗長性を回復する。

【解決手段】 故障ディスク#0のデータ/パリティの復元中にディスク#1が故障したとき、既に復元されて予備ディスクHS#0に書き込まれているデータ/パリティを破棄せずに、他の正常ディスクに格納されたデータ/パリティを用いて、ディスク#0の未復元領域とディスク#1のデータ/パリティを復元し、HS#0の対応する領域とHS#1にそれぞれ書き込む。

【選択図】 図4



【特許請求の範囲】

【請求項 1】

複数の記憶装置にデータおよびパリティを分散して格納する制御を行うことで、データの冗長化を実現するストレージ制御装置であって、

前記複数の記憶装置のうち第 1 の記憶装置が故障したとき、該第 1 の記憶装置以外の記憶装置に格納された情報を用いて該第 1 の記憶装置の情報を復元し、第 1 の予備記憶装置に書き込む第 1 の再構築手段と、

前記第 1 の記憶装置の情報を復元している間に第 2 の記憶装置が故障したとき、該第 1 および第 2 の記憶装置以外の記憶装置に格納された情報を用いて該第 1 の記憶装置の未復元領域の情報と第 2 の記憶装置の情報を復元し、前記第 1 の予備記憶装置の対応する領域と第 2 の予備記憶装置にそれぞれ書き込む第 2 の再構築手段とを備えることを特徴とするストレージ制御装置。

10

【請求項 2】

前記第 2 の再構築手段は、前記第 1 の記憶装置の未復元領域の情報を復元する処理と前記第 2 の記憶装置の情報を復元する処理をそれぞれ独立に並行して実行することを特徴とする請求項 1 記載のストレージ制御装置。

【請求項 3】

前記第 2 の再構築手段は、前記第 1 および第 2 の記憶装置以外の記憶装置に格納された、前記第 1 の記憶装置の復元済み領域に対応する領域の情報を用いて、該第 2 の記憶装置の対応する領域の情報を復元した後、該第 1 および第 2 の記憶装置以外の記憶装置に格納された、前記第 1 の記憶装置の未復元領域に対応する領域の情報を用いて、該第 1 の記憶装置の未復元領域および該第 2 の記憶装置の対応する領域の情報を復元することを特徴とする請求項 1 記載のストレージ制御装置。

20

【請求項 4】

前記第 2 の再構築手段は、前記第 1 の記憶装置の復元進捗位置と前記第 2 の記憶装置の復元進捗位置の差を閾値と比較し、該復元進捗位置の差が該閾値以上であれば、前記第 1 の記憶装置の未復元領域の情報を復元する処理と該第 2 の記憶装置の情報を復元する処理をそれぞれ独立に並行して実行し、該復元進捗位置の差が該閾値未満であれば、前記第 1 および第 2 の記憶装置以外の記憶装置に格納された、前記第 1 の記憶装置の復元済み領域に対応する領域の情報を用いて、該第 2 の記憶装置の対応する領域の情報を復元した後、該第 1 および第 2 の記憶装置以外の記憶装置に格納された、前記第 1 の記憶装置の未復元領域に対応する領域の情報を用いて、該第 1 の記憶装置の未復元領域および該第 2 の記憶装置の対応する領域の情報を復元することを特徴とする請求項 1 記載のストレージ制御装置。

30

【請求項 5】

前記第 2 の再構築手段は、前記第 1 および第 2 の記憶装置以外の記憶装置に格納された、前記第 1 の記憶装置の未復元領域に対応する領域の情報を用いて、該第 1 の記憶装置の未復元領域および該第 2 の記憶装置の対応する領域の情報を復元した後、該第 1 および第 2 の記憶装置以外の記憶装置に格納された、該第 1 の記憶装置の復元済み領域に対応する領域の情報を用いて、該第 2 の記憶装置の対応する領域の情報を復元することを特徴とする請求項 1 記載のストレージ制御装置。

40

【請求項 6】

前記第 2 の再構築手段は、前記第 1 および第 2 の記憶装置以外の記憶装置に格納された、前記第 1 の記憶装置の復元済み領域に対応する領域の情報を用いて、該第 2 の記憶装置の対応する領域の情報を復元する処理と、該第 1 および第 2 の記憶装置以外の記憶装置に格納された、前記第 1 の記憶装置の未復元領域に対応する領域の情報を用いて、該第 1 の記憶装置の未復元領域および該第 2 の記憶装置の対応する領域の情報を復元する処理を、並行して実行することを特徴とする請求項 1 記載のストレージ制御装置。

【請求項 7】

前記第 1 および第 2 の記憶装置の所定領域毎に復元済みか否かを表すビットマップ情報

50

を格納する格納手段をさらに備え、前記第2の再構築手段は、該ビットマップ情報を参照しながら復元済み以外の領域の情報を復元することを特徴とする請求項1乃至6記載のストレージ制御装置。

【請求項8】

前記第2の再構築手段は、前記第1および第2の記憶装置に対するアクセス要求が発生したとき、アクセス対象の情報を復元し、前記ビットマップ情報の該アクセス対象の情報に対応する位置に復元済みと記録することを特徴とする請求項7記載のストレージ制御装置。

【請求項9】

データの冗長化を実現するために、データおよびパリティを分散して格納する複数の記憶装置と、

前記複数の記憶装置のうち第1の記憶装置が故障したとき、該第1の記憶装置以外の記憶装置に格納された情報を用いて該第1の記憶装置の情報を復元し、第1の予備記憶装置に書き込む第1の再構築手段と、

前記第1の記憶装置の情報を復元している間に第2の記憶装置が故障したとき、該第1および第2の記憶装置以外の記憶装置に格納された情報を用いて該第1の記憶装置の未復元領域の情報と第2の記憶装置の情報を復元し、前記第1の予備記憶装置の対応する領域と第2の予備記憶装置にそれぞれ書き込む第2の再構築手段とを備えることを特徴とするストレージ装置。

【請求項10】

複数の記憶装置にデータおよびパリティを分散して格納する制御を行うことで、データの冗長化を実現するプロセッサのためのプログラムであって、

前記複数の記憶装置のうち第1の記憶装置が故障したとき、該第1の記憶装置以外の記憶装置に格納された情報を用いて該第1の記憶装置の情報を復元し、第1の予備記憶装置に書き込み、

前記第1の記憶装置の情報を復元している間に第2の記憶装置が故障したとき、該第1および第2の記憶装置以外の記憶装置に格納された情報を用いて該第1の記憶装置の未復元領域の情報と第2の記憶装置の情報を復元し、前記第1の予備記憶装置の対応する領域と第2の予備記憶装置にそれぞれ書き込む

処理を前記プロセッサに実行させることを特徴とするプログラム。

【請求項11】

複数の記憶装置にデータおよびパリティを分散して格納することで、データの冗長化を実現するストレージ制御方法であって、

前記複数の記憶装置のうち第1の記憶装置が故障したとき、該第1の記憶装置以外の記憶装置に格納された情報を用いて該第1の記憶装置の情報を復元し、第1の予備記憶装置に書き込み、

前記第1の記憶装置の情報を復元している間に第2の記憶装置が故障したとき、該第1および第2の記憶装置以外の記憶装置に格納された情報を用いて該第1の記憶装置の未復元領域の情報と第2の記憶装置の情報を復元し、前記第1の予備記憶装置の対応する領域と第2の予備記憶装置にそれぞれ書き込む

ことを特徴とするストレージ制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、RAID (Redundant Array of Inexpensive Disks) のような複数の記憶装置にデータおよびパリティを分散して格納し、記憶装置の故障時にデータおよびパリティの再構築 (Rebuild) 処理を行うストレージ制御装置および方法に関する。

【背景技術】

【0002】

RAIDは、複数のハードディスクを組み合わせ、冗長化された1台のハードディス

10

20

30

40

50

クとして管理する技術であり、ディスクへのデータ配置やデータの冗長化の方法により、RAID0～RAID6の7つのレベルに分類されている。このうち、RAID3～RAID6では、データから生成されたパリティをデータとは別に格納することで、冗長化を実現している。ディスク故障時には、このパリティを用いて故障したディスクのデータを復元する再構築処理が行われる（例えば、特許文献1参照）。

【0003】

RAID6は2つのディスクの故障に対応したRAIDレベルである。RAID6では、異なる2種類のパリティPとQがそれぞれが異なるディスクに分散して格納され、1つのディスクが故障したとき（1ディスク故障）の再構築処理と2つのディスクが故障したとき（2ディスク故障）の再構築処理では、復元方法が異なる。

10

【0004】

例えば、図23に示すように、ディスク10～14の5つのディスクからなるRAID装置において、ディスク10が故障してデータD0が失われた場合、予備ディスクであるホットスペア15を用いて1ディスク故障の再構築処理が行われる。そして、他のディスク11～13に格納されたデータD1、D2、およびパリティPからデータD0が復元される。

【0005】

これに対して、図24に示すように、ディスク10および11が故障してデータD0およびD1が失われた場合、ホットスペア15および16を用いて2ディスク故障の再構築処理が行われる。そして、他のディスク12～14に格納されたデータD2およびパリティP、QからデータD0およびD1が復元される。

20

【0006】

一般に、RAID6ではデータおよびパリティを担当するディスクがストライプ毎に異なるため、故障ディスクに格納された情報およびその復元に必要となる情報の種類もストライプ毎に異なってくる。そこで、以下の説明では、各ディスクに格納された情報をデータ/パリティと記すことにする。

【0007】

1ディスク故障から2ディスク故障になった場合、再構築処理も1ディスク故障の処理から2ディスク故障の処理へ切り替えて行う必要がある。例えば、図25に示すように、最初に故障したディスク#0をホットスペア(HS)と入れ替えて1ディスク故障の再構築処理を実行している間に、2番目のディスク#1が故障した場合、1ディスク故障の再構築処理ではデータ/パリティの復元ができなくなる。そこで、1ディスク故障の再構築処理を中止し、HS#0およびHS#1を対象に2ディスク故障の再構築処理を開始することが想定される。

30

【特許文献1】特開平03-240123号公報

【発明の開示】

【発明が解決しようとする課題】

【0008】

上述した2ディスク故障の再構築方法には、次のような問題がある。

図25に示した2ディスク故障の再構築処理では、HS#0の再構築済み領域へ既に格納された復元データ/パリティが破棄され、再び最初から再構築が実施されることとなるため、復元データ/パリティが有効活用されない。

40

【0009】

また、1ディスク故障の再構築処理に比べて、より大きな処理コストを要する2ディスク故障の再構築処理を、HS#0およびHS#1の全領域に対して実施するため、冗長性の回復までに多くの時間がかかる。

【0010】

本発明の課題は、パリティを用いてデータの冗長化を行うRAIDのようなストレージ装置において、データ/パリティを格納する2つの記憶装置が故障したときの再構築処理を改善して、効率よく冗長性を回復することである。

50

【課題を解決するための手段】

【0011】

図1は、本発明のストレージ制御装置の原理図である。図1のストレージ制御装置101は、再構築手段111および112を備え、複数の記憶装置102-1~102-Nにデータおよびパリティを分散して格納する制御を行うことで、データの冗長化を実現する。

【0012】

再構築手段111は、記憶装置102-1が故障したとき、それ以外の記憶装置に格納された情報を用いて記憶装置102-1の情報を復元し、予備記憶装置103-1に書き込む。再構築手段112は、記憶装置102-1の情報を復元している間に記憶装置102-2が故障したとき、記憶装置102-1および102-2以外の記憶装置に格納された情報を用いて記憶装置102-1の未復元領域の情報と記憶装置102-2の情報を復元し、予備記憶装置103-1の対応する領域と予備記憶装置103-2にそれぞれ書き込む。

10

【0013】

各記憶装置には、データまたはパリティが情報として格納されている。故障した記憶装置の情報は、その時点で故障していない正常な記憶装置に格納された情報を用いて復元され、対応する予備記憶装置に書き込まれる。記憶装置102-1の情報を復元している間に記憶装置102-2が故障したとき、記憶装置102-1の復元済み領域の情報はそのまま予備記憶装置103-1内に保存され、記憶装置102-1の未復元領域と記憶装置102-2の全領域を対象として復元処理が行われる。

20

【0014】

このような再構築制御によれば、1ディスク故障から2ディスク故障になったとき、既に復元されている情報が消去されずに有効利用される。また、記憶装置102-1の復元済み領域は復元対象から除外されるため、再構築処理に要する時間が削減される。さらに、記憶装置102-1の復元済み領域に対応する記憶装置102-2の領域については、処理コストの小さな1ディスク故障の再構築処理を適用できるため、さらなる効率化が可能である。

【0015】

ストレージ制御装置101は、例えば、後述する図2のコントローラ211、図21のホストバスアダプタ1911、または図22のホスト装置2001に対応する。

30

【発明の効果】

【0016】

本発明によれば、複数の記憶装置にデータおよびパリティを分散して格納するストレージ装置において、2つの記憶装置が故障したときの再構築処理が改善され、データの冗長性が効率よく回復する。

【発明を実施するための最良の形態】

【0017】

以下、図面を参照しながら、本発明を実施するための最良の形態を詳細に説明する。

図2は、実施形態のストレージシステムの構成例を示している。図2のストレージシステムは、ホスト装置201およびRAID装置202からなり、RAID装置202は、コントローラ211およびDisk#0~Disk#3の4つのディスクを備えるストレージ装置に対応する。

40

【0018】

各ディスクは、1つ以上の磁気ディスク装置からなり、ホスト装置201は、各ディスクを1つの記憶装置とみなしてデータのリード/ライトを行う。ただし、コントローラ211に接続されるディスクの数は4つに限られるわけではなく、一般にはより多くのディスクが接続される。

【0019】

コントローラ211は、プロセッサ221、メモリ222、およびキャッシュメモリ2

50

23を備え、Disk#0～Disk#3の故障時の再構築処理を行う。プロセッサ221は、メモリ222に格納されたプログラムを実行することにより、キャッシュメモリ223をデータバッファとして利用しながら再構築処理を行う。

【0020】

本実施形態においては、2ディスク故障の再構築処理を改善するために、図3に示す改善案1～6を採用する。まず、図4から図8までを参照しながら、各改善案の概要について説明する。

1. 個別実行案(改善案1)

コントローラ211は、1番目のディスクの再構築と2番目のディスクの再構築を独立に行う。特に、両方のディスクの故障部分に対応する二重故障部分の再構築処理では、2つのディスクのデータ/パリティをデータバッファ上で復元しながら、一方のディスクの復元データ/パリティのみをディスクにライトし、もう一方のディスクの復元データ/パリティはライトせずに破棄する。

10

【0021】

例えば、図4に示すように、最初にディスク#0が故障し、HS#0を用いてディスク#0のデータ/パリティを復元している間にディスク#1が故障した場合について説明する。この場合、ディスク#1が故障するまでの間は、正常なディスク#1～#3のデータ/パリティを用いて1ディスク故障の再構築処理によりディスク#0のデータ/パリティが復元される。

【0022】

ディスク#1が故障すると、HS#0については、再構築済み領域のデータ/パリティを破棄せずにそのまま保持し、残りの部分のデータ/パリティのみを2ディスク故障の再構築処理により復元する。この処理では、正常なディスク#2～#3のデータ/パリティを用いてHS#0の残りの部分のデータ/パリティが復元される。このとき、ディスク#1のデータ/パリティも同時に生成されるが、ディスク#1の再構築処理は別途独立に行われるため、生成されたデータ/パリティは破棄される。

20

【0023】

HS#1については、HS#0の再構築処理と並行して、正常なディスク#2～#3のデータ/パリティを用いて2ディスク故障の再構築処理により全体のデータ/パリティを復元する。このとき、2ディスク故障の再構築処理により同時に生成されるディスク#0のデータ/パリティは破棄される。

30

【0024】

このような再構築処理によれば、ディスク#1が故障したときに、ディスク#0の再構築処理が、復元対象のディスクのみが故障している1ディスク故障の再構築処理から、復元対象以外にも故障ディスクが存在する2ディスク故障の再構築処理に切り替えられ、復元済み領域のデータ/パリティはそのまま保持される。したがって、図25に示したように復元済み領域のデータ/パリティを再度復元し直す必要がなく、図25の再構築処理よりも早くディスク#0の復元が完了する。

2. 同じ進捗位置になるまで待つ案(改善案2)

コントローラ211は、2番目のディスクの故障時に、1番目のディスクの再構築処理を一旦停止し、2番目のディスクのみを1番目と同じ進捗位置になるまで再構築する。そして、同じ進捗位置になったら、それ以降は両方同時に再構築する。

40

【0025】

例えば、図5に示すように、最初にディスク#0が故障し、HS#0を用いてディスク#0のデータ/パリティを復元している間にディスク#1が故障した場合について説明する。ディスク#1が故障するまでの動作は、図4の場合と同様である。

【0026】

ディスク#1が故障すると、最初に、HS#0の再構築済み領域に対応するHS#1のデータ/パリティのみを、HS#0の再構築済み領域のデータ/パリティと正常なディスク#2～#3のデータ/パリティを用いて1ディスク故障の再構築処理により復元する。

50

そして、その部分のデータ/パリティが復元されたら、次に、HS#0およびHS#1の残りの部分のデータ/パリティを、正常なディスク#2～#3のデータ/パリティを用いて2ディスク故障の再構築処理により同時に復元する。

【0027】

このような再構築処理によれば、改善案1と同様の利点に加えて、復元中のディスク#0の復元済み領域のデータ/パリティを他のディスク#1の復元に有効利用することができ、ディスク#1の再構築処理が効率化される。

3. 混在案(改善案3)

改善案1と改善案2の混在案である。コントローラ211は、処理中に2つのディスクHS#0およびHS#1の再構築進捗位置をモニタしながら、復元方法を選択する。具体的には、進捗位置の差をチェックし、その差が所定の閾値以上であれば改善案1を適用して、HS#0の再構築を優先させ、差が閾値未満であれば改善案2を適用して、互いの進捗位置が同じとなるまでHS#1の再構築を行う。進捗位置の差のチェックは、所定単位のデータ/パリティを復元する度に毎回行うか、あるいは一定範囲のデータ/パリティを復元する度に行う。

10

【0028】

改善案2を適用した場合、HS#0およびHS#1の進捗位置の差は、冗長性が失われている二重故障部分の復元開始を待ち合わせる時間に対応し、進捗位置の差が大きいほど冗長性の回復は遅れることになる。そこで、待ち合わせ時間が一定時間以上となる場合は、それを避けるために改善案1を適用して、HS#0およびHS#1の再構築処理を並行して行う。

20

【0029】

ただし、改善案1とは異なり、HS#1の二重故障部分以外のデータ/パリティは、HS#0の再構築済み領域のデータ/パリティと正常なディスク#2～#3のデータ/パリティを用いて1ディスク故障の再構築処理により復元する。一般に、1ディスク故障の再構築処理に比べて、2ディスク故障の再構築処理の方が多くの計算量が必要となる。このため、先行しているHS#0の2ディスク故障の再構築処理より、HS#1の1ディスク故障の再構築処理の方が早く進行し、時間と共に進捗位置の差が縮まってくる可能性が高い。

【0030】

このような再構築処理によれば、改善案2と同様の利点に加えて、待ち合わせによる冗長性回復の遅れを回避することができる。また、並行処理によりプロセッサ資源等が有効利用されることも期待できる。

30

4. 二重故障部分を先に復元する案(改善案4)

コントローラ211は、2番目のディスクの故障時に、1番目のディスクの現在の進捗位置を復元完了位置(後述する停止位置)として保存し、その位置から両方同時の再構築を行う。そして、両方同時の再構築が最後まで終了したら、2番目のディスクの未済領域復旧のため、2番目のディスクの最初から復元完了位置までの再構築を行う。

【0031】

例えば、図6に示すように、最初にディスク#0が故障し、HS#0を用いてディスク#0のデータ/パリティを復元している間にディスク#1が故障した場合について説明する。ディスク#1が故障するまでの動作は、図4の場合と同様である。

40

【0032】

ディスク#1が故障すると、最初に、HS#0の再構築未済領域に対応するHS#0およびHS#1のデータ/パリティを、正常なディスク#2～#3のデータ/パリティを用いて2ディスク故障の再構築処理により同時に復元する。そして、その部分のデータ/パリティが復元されたら、次に、HS#1の残りの部分のデータ/パリティを、HS#0の再構築済み領域のデータ/パリティと正常なディスク#2～#3のデータ/パリティを用いて1ディスク故障の再構築処理により復元する。

【0033】

50

再構築処理中に R A I D 6 の通常のリード/ライトアクセスがあったときは、アクセス対象データの再構築済み/未済をチェックするか、あるいは、リード要求に対してデータを再構築しながら返すというように、通常のリード/ライトアクセスをすべて縮退動作にする。

【 0 0 3 4 】

このような再構築処理によれば、改善案 2 と同様の利点に加えて、二重故障部分のデータ/パリティを先に復元することで、R A I D グループとしての冗長性を短時間で回復することが可能になる。

5 . 二重故障復元と復元済み H S を用いた復元を並行して実施する案 (改善案 5)

コントローラ 2 1 1 は、改善案 4 において後で実施していた 2 番目のディスクの最初から復元完了位置までの再構築処理を、二重故障部分の再構築処理を待ち合わせずに並行して実施する。

10

【 0 0 3 5 】

例えば、図 7 に示すように、最初にディスク # 0 が故障し、H S # 0 を用いてディスク # 0 のデータ/パリティを復元している間にディスク # 1 が故障した場合について説明する。ディスク # 1 が故障するまでの動作は、図 4 の場合と同様である。

【 0 0 3 6 】

ディスク # 1 が故障すると、H S # 0 および H S # 1 の二重故障部分のデータ/パリティを、正常なディスク # 2 ~ # 3 のデータ/パリティを用いて 2 ディスク故障の再構築処理により復元する処理と、H S # 1 の残りの部分のデータ/パリティを、H S # 0 の再構築済み領域のデータ/パリティと正常なディスク # 2 ~ # 3 のデータ/パリティを用いて 1 ディスク故障の再構築処理により復元する処理とを、並行して実施する。

20

【 0 0 3 7 】

再構築処理中に通常のリード/ライトアクセスがあったときは、改善案 4 の場合と同様に、アクセス対象データの再構築済み/未済をチェックするか、あるいは縮退動作を行う。

【 0 0 3 8 】

このような再構築処理によれば、R A I D グループとしての冗長性の回復は改善案 4 より遅くなるが、冗長性が失われた二重故障部分と冗長性が残っている他の部分の復元処理が並行して行われるため、全体の復元に要する時間が短縮される。

30

6 . ランダム案 (改善案 6)

コントローラ 2 1 1 は、キャッシュメモリ 2 2 3 上のビットマップを用いて、ディスクの所定領域毎に再構築動作を行う。再構築順序としては、上記改善案 1 ~ 5 のいずれかを採用する。基本的には、データ/パリティの保全の観点から改善案 4 の順序を採用することが望ましいが、他の改善案の順序でも動作可能である。

【 0 0 3 9 】

また、通常のリード/ライトアクセスの延長で再構築動作が行われたときも、ビットマップには再構築済みと記録する。このため、リード/ライトアクセスとは無関係のシーケンシャルな再構築と、リード/ライトアクセスの延長で行われるワンポイント再構築の 2 種類の処理を併用して、再構築処理が行われる。ワンポイント再構築では、リード要求に対してデータを復元した場合、コントローラ 2 1 1 は復元データをディスクの対応位置にライトする。

40

【 0 0 4 0 】

コントローラ 2 1 1 は、再構築対象のディスク毎に、例えば、1 ストライプを 1 ビットに対応付けたビットマップを用意して、進捗状況を管理する。もし、ビットマップが失われた場合は、最初から再構築を開始する。1 つの論理ブロックを 1 ビットに対応付けたビットマップを用いてもよい。

【 0 0 4 1 】

また、最初にビットマップ用のメモリ領域が獲得できなかった場合、コントローラ 2 1 1 は以下のいずれかの動作を行う。

50

- ・上述の別の改善案で再構築を実行する。
- ・ビットマップのサイズに上限を設定しておき、これを超えたビットマップを必要とする場合、再構築は資源の獲得を待ってから行う。

【0042】

コントローラ211は、再構築動作を並行して行うことが可能であり、電源オフ/オンを考慮して、ビットマップのバックアップ/リストアの機能も有する。コントローラ211が冗長化(二重化)されている場合、ビットマップもコントローラ間で二重化することが基本であるが、二重化しなくてもデータロストに繋がることはない。上述したように、ビットマップが失われた場合は再構築が再始動される。

【0043】

例えば、図8に示すように、最初にディスク#0が故障し、HS#0を用いてディスク#0のデータ/パリティを復元している間にディスク#1が故障した場合について説明する。

【0044】

コントローラ211は、HS#0およびHS#1のそれぞれに対してビットマップを生成し、各ストライプを1ビットのデータで管理する。再構築未済のストライプについては対応するビットに“1”が記録され、再構築済みのストライプについては“0”が記録される。HS#0の再構築済み領域の全ストライプについては、ディスク#1が故障した時点で“0”が記録されており、残りのストライプについては再構築が行われたときに“0”が記録される。

【0045】

このような再構築処理によれば、改善案1~5と同様の利点に加えて、ワンポイント再構築により復元されたデータ/パリティをディスクに書き戻して復元済み領域として扱うことで、処理が効率化される。

【0046】

次に、図9から図19までを参照しながら、上述した各改善案の詳細について説明する。

各改善案の再構築処理は、ディスクの故障を契機として起動されるか、または他の再構築処理により起動される。コントローラ211は、契機となったディスクを担当Mainに指定し、必要に応じて他の故障ディスクを担当Subとして追加する。再構築処理では、図9に示すように、担当Mainと担当Subで共通となる現在位置について、復元処理が行われる。ここで、現在位置とは、現在、復元処理が実施されている場所を表す。

【0047】

また、コントローラ211は、RAID装置202を構成する全ディスクの情報を、すべての改善案に共通の制御情報としてキャッシュメモリ223に保持する。具体的には、図10に示すように、各ディスクについて、復元状態(復元済み、復元中、および未復元)、復元先頭位置、および停止位置(必要であれば)の情報が保持される。ただし、正常ディスクについては、すべての領域が復元済みに設定される。

【0048】

故障ディスクの復元先頭位置は、そのディスクを担当Mainとして実施されている復元処理の現在位置と一致し、復元処理が未実施の場合はそのディスクの端(図9の下端)に設定される。復元処理は、復元先頭位置から上へ向かう方向に進行する。停止位置は、復元処理を停止すべき場所を表す。

【0049】

復元状態は、ストリップ、論理ブロック等の単位領域毎に管理され、現在位置、復元先頭位置、および停止位置等の位置情報は、その単位領域のアドレスまたはその単位領域が属するストライプの識別子を用いて管理される。

【0050】

図11は、再構築処理のフローチャートである。コントローラ211は、まず、担当Mainの端(図9の下端)をそのディスクの現在位置として設定し(ステップ1101)

10

20

30

40

50

、復元ルーチンを実行して復元処理を行う（ステップ1102）。この復元処理では、復元対象ディスクの復元対象となるデータ/パリティが他のディスクのデータ/パリティを用いて生成され、対応するホットスペアにライトされる。1回の復元処理では、論理ブロック、ストライプ等の所定単位のデータ/パリティが復元されるが、一般的には、ストライプが所定単位として用いられる。

【0051】

次に、再構築制御のための復元後処理を行う（ステップ1103）。この復元後処理では、復元先頭位置の設定や復元処理を終了するか否かの判定等が行われる。その後、復元後処理において復元処理終了と判定されたか否かをチェックする（ステップ1104）。復元処理終了と判定されなければ、現在位置を1ストライプだけ進めて（ステップ1105）、ステップ1102以降の処理を繰り返し、復元処理終了と判定されれば、再構築処理を終了する。

10

【0052】

ステップ1102の復元ルーチンおよびステップ1103の復元後処理は、改善案毎に異なるため、以下、改善案1～6について順番に説明する。

1. 改善案1

改善案1では、図11の再構築処理は各ディスクの故障を契機として起動され、契機となったディスクが担当Mainに設定されるが、他の故障ディスクが担当Subとして追加されることはない。したがって、2番目のディスクが故障した後は、2つの再構築処理が並行して実行され、各再構築処理における復元対象ディスクは担当Mainのみとなる。

20

【0053】

図12は、改善案1の復元ルーチンのフローチャートである。コントローラ211は、まず、担当Mainを復元対象ディスクに設定し、その現在位置を設定して（ステップ1201）、故障ディスクの個数をチェックする（ステップ1202）。そして、故障ディスクが1個であれば、復元方法を1ディスク故障の再構築に決定する（ステップ1203）。

【0054】

次に、現在位置のストライプに属するデータ/パリティのうち、1ディスク故障の再構築に必要なものを正常ディスクからリードし（ステップ1205）、それらがすべてリードできたか否かをチェックする（ステップ1206）。データ/パリティがすべてリードできた場合は、それらを用いて同じストライプに属する復元対象ディスクのデータ/パリティを復元し、対応するホットスペアにライトする（ステップ1207）。

30

【0055】

ステップ1206においてリードエラーが発生した場合は、リード対象のディスクが故障したものと判断する。そこで、故障ディスクの個数をチェックし（ステップ1208）、それが2個であれば、ステップ1202以降の処理を行う。

【0056】

そして、復元方法を2ディスク故障の再構築に変更し（ステップ1204）、2ディスク故障の再構築に必要なデータ/パリティを正常ディスクからリードする（ステップ1205）。データ/パリティがすべてリードできた場合は、それらを用いて復元対象ディスクのデータ/パリティを復元し、ホットスペアにライトする（ステップ1207）。

40

【0057】

故障ディスクが2個のときにさらにリードエラーが発生した場合は、故障ディスクが3個となるため（ステップ1208）、復元不可能と判断し、エラー処理を行う（ステップ1209）。

【0058】

図13は、改善案1の復元後処理のフローチャートである。コントローラ211は、まず、復元ルーチンで用いた現在位置を担当Mainの復元先頭位置に設定し（ステップ1301）、担当Mainのすべての領域の復元が完了したか否かをチェックする（ステッ

50

ブ 1 3 0 2)。ここでは、復元ルーチンで用いた現在位置が担当 M a i n の端 (図 9 の上端) に達していれば、すべての領域の復元が完了したものと判定される。すべての領域の復元が完了していなければ、復元処理継続と判定し (ステップ 1 3 0 3)、すべての領域の復元が完了すれば、復元処理終了と判定する (ステップ 1 3 0 4)。

【 0 0 5 9 】

例えば、図 4 に示したように、最初にディスク # 0 が故障すると、ディスク # 0 を担当 M a i n として再構築処理が起動される。このとき、故障ディスクは 1 個であるから (図 1 2 のステップ 1 2 0 2)、復元方法は 1 ディスク故障の再構築となり (ステップ 1 2 0 3)、正常なディスク # 1 ~ # 3 のデータ / パリティのうち、ディスク # 0 のデータ / パリティを 1 ディスク故障の再構築により復元するために必要なものがリードされる (ステップ 1 2 0 5)。そして、リードされたデータ / パリティを用いてディスク # 0 のデータ / パリティが復元され、H S # 0 にライトされる (ステップ 1 2 0 7)。

10

【 0 0 6 0 】

ディスク # 0 の現在位置は、その復元先頭位置として設定され (図 1 3 のステップ 1 3 0 1)、復元処理継続と判定される (ステップ 1 3 0 3)。復元先頭位置は、ディスク # 0 に対するリード / ライトアクセス等の他の処理から参照される。このような復元ルーチンおよび復元後処理が 1 ストライプ毎に繰り返し実行される (図 1 1 のステップ 1 1 0 5)。

【 0 0 6 1 】

次に、ディスク # 1 が故障すると、故障ディスクは 2 個となるから (ステップ 1 2 0 2)、復元方法は 2 ディスク故障の再構築となり (ステップ 1 2 0 4)、正常なディスク # 2 ~ # 3 のデータ / パリティのうち、ディスク # 0 および # 1 のデータ / パリティを 2 ディスク故障の再構築により復元するために必要なものがリードされる (ステップ 1 2 0 5)。そして、リードされたデータ / パリティを用いてディスク # 0 および # 1 のデータ / パリティが復元され、そのうちディスク # 0 のデータ / パリティが H S # 0 にライトされる (ステップ 1 2 0 7)。

20

【 0 0 6 2 】

復元後処理については、ディスク # 1 の故障前と同様である。このような復元ルーチンおよび復元後処理が 1 ストライプ毎に繰り返し実行され (ステップ 1 1 0 5)、ディスク # 0 のすべての領域の復元が完了すれば (ステップ 1 3 0 4)、ディスク # 0 の再構築処理を終了する (ステップ 1 1 0 4)。

30

【 0 0 6 3 】

さらに、ディスク # 1 の故障時には、ディスク # 1 を担当 M a i n としてもう 1 つの再構築処理が起動される。このとき、故障ディスクは 2 個であるから (ステップ 1 2 0 2)、復元方法は 2 ディスク故障の再構築処理となり (ステップ 1 2 0 4)、正常なディスク # 2 ~ # 3 のデータ / パリティのうち、ディスク # 0 および # 1 のデータ / パリティを 2 ディスク故障の再構築処理により復元するために必要なものがリードされる (ステップ 1 2 0 5)。

【 0 0 6 4 】

そして、リードされたデータ / パリティを用いてディスク # 0 および # 1 のデータ / パリティが復元され、そのうちディスク # 1 のデータ / パリティが H S # 1 にライトされる (ステップ 1 2 0 7)。

40

【 0 0 6 5 】

復元後処理については、ディスク # 0 の場合と同様である。このような復元ルーチンおよび復元後処理が 1 ストライプ毎に繰り返し実行され (ステップ 1 1 0 5)、ディスク # 1 のすべての領域の復元が完了すれば (ステップ 1 3 0 4)、ディスク # 1 の再構築処理を終了する (ステップ 1 1 0 4)。

2 . 改善案 2

改善案 2 では、改善案 1 と同様に、図 1 1 の再構築処理は各ディスクの故障を契機として起動され、契機となったディスクが担当 M a i n に設定される。2 番目のディスクが故

50

障すると、1番目の故障ディスクを担当Mainとする再構築処理が中断され、2番目の故障ディスクを担当Mainとする再構築処理が開始される。そして、2番目の故障ディスクの現在位置が1番目と同じ進捗位置に達すると、1番目の故障ディスクが担当Subとして追加される。

【0066】

図14は、改善案2の復元ルーチンのフローチャートである。この場合、図12の復元ルーチンとは異なり、復元対象ディスクは担当Mainと担当Subで表され、2個まで設定できる。また、復元方法は、故障ディスクの個数ではなく、復元対象ディスクの個数に基づいて選択される。

【0067】

コントローラ211は、まず、担当Main/担当Subを復元対象ディスクに設定し、担当Mainの現在位置を設定する(ステップ1401)。このとき、担当Subが設定されていなければ、担当Mainのみが復元対象ディスクに設定される。

【0068】

次に、復元対象ディスクの個数をチェックし(ステップ1402)、復元対象ディスクが1個であれば、復元方法を1ディスク故障の再構築に決定する(ステップ1403)。そして、現在位置のストライプに属するデータ/パリティのうち、1ディスク故障の再構築に必要なものを正常ディスクからリードし(ステップ1405)、それらがすべてリードできたか否かをチェックする(ステップ1406)。データ/パリティがすべてリードできた場合は、それらを用いて同じストライプに属する復元対象ディスクのデータ/パリティを復元し、対応するホットスペアにライトする(ステップ1407)。

【0069】

ステップ1406においてリードエラーが発生した場合は、リード対象のディスクが故障したものと判断する。そこで、復元対象ディスクの個数をチェックし(ステップ1408)、それが1個であれば、故障ディスクを復元対象ディスクに追加して(ステップ1410)、ステップ1402以降の処理を行う。

【0070】

そして、復元方法を2ディスク故障の再構築に変更し(ステップ1404)、2ディスク故障の再構築に必要なデータ/パリティを正常ディスクからリードする(ステップ1405)。データ/パリティがすべてリードできた場合は、それらを用いて2個の復元対象ディスクのデータ/パリティを復元し、それぞれ対応するホットスペアにライトする(ステップ1407)。

【0071】

復元対象ディスクが2個のときにさらにリードエラーが発生した場合は、故障ディスクが3個となるため(ステップ1408)、復元不可能と判断し、エラー処理を行う(ステップ1409)。

【0072】

図15は、改善案2の復元後処理のフローチャートである。コントローラ211は、まず、復元ルーチン終了時の担当Mainの現在位置を担当Main/Subの復元先頭位置に設定し(ステップ1501)、以下の条件aが満たされるか否かをチェックする(ステップ1502)。

条件a: 担当Main以外に他の故障ディスクが存在し、担当Mainおよびその故障ディスクのいずれにも停止位置が設定されておらず、かつ、その故障ディスクの復元先頭位置が担当Mainのそれより後ろ(下方)である。

【0073】

他の故障ディスクの復元先頭位置が担当Mainより後ろの場合、担当Mainより復元処理が遅れていることを意味する。条件aが満たされれば、担当Mainの復元先頭位置を他の故障ディスクの停止位置に設定し(ステップ1506)、担当Mainの再構築処理を中断するために復元処理終了と判定する(ステップ1508)。

【0074】

10

20

30

40

50

条件 a が満たされなければ、次に、担当 Main のすべての領域の復元が完了したか否かをチェックする（ステップ 1503）。すべての領域の復元が完了していれば、復元処理終了と判定する（ステップ 1508）。

【0075】

すべての領域の復元が完了していなければ、次に、担当 Main に停止位置が設定されており、かつ、担当 Main の現在位置がその停止位置であるか否かをチェックする（ステップ 1504）。現在位置が停止位置であれば、他の故障ディスクを担当 Sub に追加し（ステップ 1507）、復元処理継続と判定する（ステップ 1505）。

【0076】

現在位置が停止位置でない場合、および、停止位置が設定されていない場合は、そのまま復元処理継続と判定する（ステップ 1505）。 10

例えば、図 5 に示したように、最初にディスク # 0 が故障すると、ディスク # 0 を担当 Main として再構築処理が起動される。このとき、復元対象ディスクは 1 個であるから（図 14 のステップ 1401）、復元方法は 1 ディスク故障の再構築となり（ステップ 1403）、正常なディスク # 1 ~ # 3 のデータ / パリティのうち、ディスク # 0 のデータ / パリティを 1 ディスク故障の再構築により復元するために必要なものがリードされる（ステップ 1405）。そして、リードされたデータ / パリティを用いてディスク # 0 のデータ / パリティが復元され、HS # 0 にライトされる（ステップ 1407）。

【0077】

ディスク # 0 の現在位置は、その復元先頭位置として設定され（図 15 のステップ 1501）、他の故障ディスクはないので（ステップ 1502）、復元処理継続と判定される（ステップ 1505）。このような復元ルーチンおよび復元後処理が 1 ストライプ毎に繰り返し実行される（図 11 のステップ 1105）。 20

【0078】

次に、ディスク # 1 が故障すると、担当 Main であるディスク # 0 の復元先頭位置はその現在位置に一致しており（ステップ 1501）、他の故障ディスクであるディスク # 1 の復元先頭位置はその下端に一致しているため、条件 a が満たされる（ステップ 1502）。そこで、ディスク # 0 の復元先頭位置がディスク # 1 の停止位置に設定され（ステップ 1506）、復元処理終了と判定される（ステップ 1508）。これにより、ディスク # 0 を担当 Main とする再構築処理が中断される（ステップ 1104）。 30

【0079】

このとき、ディスク # 1 を担当 Main として別の再構築処理が起動される。復元対象ディスクは 1 個であるから（ステップ 1401）、復元方法は 1 ディスク故障の再構築となり（ステップ 1403）、ディスク # 0、# 2、および # 3 のデータ / パリティのうち、ディスク # 1 のデータ / パリティを 1 ディスク故障の再構築により復元するために必要なものがリードされる（ステップ 1405）。ただし、ディスク # 0 については、HS # 0 にライトされた復元済みのデータ / パリティがリードされる。

【0080】

そして、リードされたデータ / パリティを用いてディスク # 1 のデータ / パリティが復元され、HS # 1 にライトされる（ステップ 1407）。 40

ディスク # 1 の現在位置は、その復元先頭位置として設定され（ステップ 1501）、ディスク # 1 には既に停止位置が設定されているため、条件 a は満たされない（ステップ 1502）。また、ディスク # 1 の現在位置はその停止位置に到達していないので（ステップ 1504）、復元処理継続と判定される（ステップ 1505）。このような復元ルーチンおよび復元後処理が 1 ストライプ毎に繰り返し実行される（ステップ 1105）。

【0081】

そして、ディスク # 1 の現在位置がその停止位置に到達すると（ステップ 1504）、処理が中断していたディスク # 0 が担当 Sub として追加され（ステップ 1507）、復元処理継続と判定される（ステップ 1505）。これにより、現在位置が更新される（ステップ 1105）。 50

【0082】

これにより、復元対象ディスクは2個となるから（ステップ1401）、復元方法は2ディスク故障の再構築処理となり（ステップ1404）、正常なディスク#2～#3のデータ/パリティのうち、ディスク#0および#1のデータ/パリティを2ディスク故障の再構築処理により復元するために必要なものがリードされる（ステップ1405）。そして、リードされたデータ/パリティを用いてディスク#0および#1のデータ/パリティが復元され、それぞれHS#0およびHS#1にライトされる（ステップ1407）。

【0083】

ディスク#1の現在位置は、ディスク#0および#1の復元先頭位置として設定され（ステップ1501）、ディスク#1の現在位置はその停止位置を越えているため（ステップ1504）、復元処理継続と判定される（ステップ1505）。

【0084】

このような復元ルーチンおよび復元後処理が1ストライプ毎に繰り返し実行され（ステップ1105）、ディスク#1のすべての領域の復元が完了すれば（ステップ1508）、ディスク#1を担当Mainとする再構築処理を終了する（ステップ1104）。この時点で、担当Subであるディスク#0の現在位置もその上端に達しているため、ディスク#0の復元も同時に完了する。

3. 改善案3

改善案3では、改善案2と同様に、図11の再構築処理は各ディスクの故障を契機として起動され、契機となったディスクが担当Mainに設定される。2番目のディスクが故障すると、2つの故障ディスクの進捗位置の差に応じて改善案1または2が選択される。

【0085】

進捗位置の差が閾値以上であれば改善案1が選択され、2つの再構築処理が並行して実行される。ただし、改善案1とは異なり、2番目の故障ディスクのデータ/パリティは、1番目の故障ディスクの復元済みのデータ/パリティと正常なディスクのデータ/パリティを用いて1ディスク故障の再構築処理により復元される。

【0086】

そして、進捗位置の差が閾値未満になると改善案2が選択され、1番目の故障ディスクを担当Mainとする再構築処理が中断され、2番目の故障ディスクを担当Mainとする再構築処理が開始される。そして、2番目の故障ディスクの現在位置が1番目と同じ進捗位置に達すると、1番目の故障ディスクが担当Subとして追加される。

【0087】

改善案3の復元ルーチンのフローチャートは改善案2と同様であり、復元後処理のフローチャートは図16のようになる。図16の復元後処理は、図15の復元後処理にステップ1603の判定を追加した構成を有する。

【0088】

ステップ1603において、コントローラ211は、他の故障ディスクと担当Mainの復元先頭位置の差を閾値と比較する。そして、復元先頭位置の差が閾値未満であれば、担当Mainの復元先頭位置を他の故障ディスクの停止位置に設定して（ステップ1607）、復元処理終了と判定する（ステップ1609）。また、復元先頭位置の差が閾値以上であれば、ステップ1604以降の処理を行う。

【0089】

したがって、改善案2で説明した条件aが満たされ、かつ、他の故障ディスクと担当Mainの復元先頭位置の差が閾値未満であれば、他の故障ディスクに停止位置が設定され（ステップ1607）、それ以外の場合は停止位置は設定されない。

【0090】

例えば、図5に示したように、最初にディスク#0が故障すると、ディスク#0を担当Mainとして再構築処理が起動され、ディスク#1が故障するまで改善案2と同様の処理が行われる。

【0091】

10

20

30

40

50

次に、ディスク# 1が故障すると、条件aが満たされるため(図16のステップ1602)、ディスク# 0および# 1の復元先頭位置の差が閾値と比較される(ステップ1603)。ここで、ディスク# 0が故障してから十分に時間が経過している場合は、ディスク# 0の復元処理がかなり進行しており、復元先頭位置の差は閾値を超えていると考えられる。この場合、ステップ1604以降の処理が行われ、ディスク# 0には停止位置が設定されていないので(ステップ1605)、復元処理継続と判定される(ステップ1606)。

【0092】

このとき、復元対象ディスクは1個であるから(図14のステップ1401)、復元方法は1ディスク故障の再構築となり(ステップ1403)、ディスク# 1~# 3のデータ/パリティのうち、ディスク# 0のデータ/パリティを1ディスク故障の再構築により復元するために必要なものがリードされる(ステップ1405)。

10

【0093】

ここで、ディスク# 1のデータ/パリティがリードされた場合、リードエラーが発生し、ディスク# 1が復元対象ディスクに追加されて(ステップ1410)、復元対象ディスクは2個となる。

【0094】

このため、復元方法は2ディスク故障の再構築処理に変更され(ステップ1404)、正常なディスク# 2~# 3のデータ/パリティのうち、ディスク# 0および# 1のデータ/パリティを2ディスク故障の再構築処理により復元するために必要なものがリードされる(ステップ1405)。そして、リードされたデータ/パリティを用いてディスク# 0および# 1のデータ/パリティが復元され、そのうちディスク# 0のデータ/パリティがHS# 0にライトされる(ステップ1407)。

20

【0095】

ディスク# 0の現在位置は、その復元先頭位置として設定され(ステップ1601)、条件aが満たされるが(ステップ1602)、復元先頭位置の差はまだ閾値より大きい(ステップ1603)。また、ディスク# 0には停止位置が設定されていないので(ステップ1605)、復元処理継続と判定される(ステップ1606)。このような復元ルーチンおよび復元後処理が1ストライプ毎に繰り返し実行される(図11のステップ1105)。

30

【0096】

さらに、ディスク# 1の故障時には、ディスク# 1を担当Mainとしてもう1つの再構築処理が起動される。復元対象ディスクは1個であるから(ステップ1401)、復元方法は1ディスク故障の再構築となり(ステップ1403)、ディスク# 0、# 2、および# 3のデータ/パリティのうち、ディスク# 1のデータ/パリティを1ディスク故障の再構築により復元するために必要なものがリードされる(ステップ1405)。ただし、ディスク# 0については、HS# 0にライトされた復元済みのデータ/パリティがリードされる。

【0097】

そして、リードされたデータ/パリティを用いてディスク# 1のデータ/パリティが復元され、HS# 1にライトされる(ステップ1407)。

40

ディスク# 1の現在位置は、その復元先頭位置として設定され(ステップ1601)、他の故障ディスクであるディスク# 0の復元先頭位置は担当Mainであるディスク# 1の復元先頭位置より前にあるため、条件aは満たされない(ステップ1602)。また、ディスク# 1には停止位置が設定されていないので(ステップ1605)、復元処理継続と判定される(ステップ1606)。このような復元ルーチンおよび復元後処理が1ストライプ毎に繰り返し実行される(ステップ1105)。

【0098】

こうして、ディスク# 0を担当Mainとする2ディスク故障の再構築処理と、ディスク# 1を担当Mainとする1ディスク故障の再構築処理とが並行して実行され、ディス

50

ク # 1 の復元先頭位置が徐々にディスク # 0 の復元先頭位置に接近してくる。

【 0 0 9 9 】

そして、ディスク # 0 を担当 M a i n とする再構築処理において、ディスク # 0 および # 1 の復元先頭位置の差が閾値未満になると (ステップ 1 6 0 3)、ディスク # 0 の復元先頭位置がディスク # 1 の停止位置に設定され (ステップ 1 6 0 7)、復元処理終了と判定される (ステップ 1 6 0 9)。これにより、ディスク # 0 を担当 M a i n とする再構築処理が中断される (ステップ 1 1 0 4)。

【 0 1 0 0 】

その後、ディスク # 1 を担当 M a i n とする再構築処理のみが継続されるが、ディスク # 1 には既に停止位置が設定されているため、やはり条件 a は満たされない (ステップ 1 6 0 2)。また、ディスク # 1 の現在位置はその停止位置に到達していないので (ステップ 1 6 0 5)、復元処理継続と判定される (ステップ 1 6 0 6)。

10

【 0 1 0 1 】

そして、ディスク # 1 の現在位置がその停止位置に到達すれば (ステップ 1 6 0 5)、処理が中断していたディスク # 0 が担当 S u b として追加され (ステップ 1 6 0 8)、復元処理継続と判定される (ステップ 1 6 0 6)。これにより、現在位置が更新される (ステップ 1 1 0 5)。

【 0 1 0 2 】

これにより、復元対象ディスクは 2 個となるから (ステップ 1 4 0 1)、復元方法は 2 ディスク故障の再構築処理となり (ステップ 1 4 0 4)、正常なディスク # 2 ~ # 3 のデータ / パリティのうち、ディスク # 0 および # 1 のデータ / パリティを 2 ディスク故障の再構築処理により復元するために必要なものがリードされる (ステップ 1 4 0 5)。そして、リードされたデータ / パリティを用いてディスク # 0 および # 1 のデータ / パリティが復元され、それぞれ H S # 0 および H S # 1 にライトされる (ステップ 1 4 0 7)。

20

【 0 1 0 3 】

ディスク # 1 の現在位置は、ディスク # 0 および # 1 の復元先頭位置として設定され (ステップ 1 6 0 1)、ディスク # 1 の現在位置はその停止位置を越えているため (ステップ 1 6 0 5)、復元処理継続と判定される (ステップ 1 6 0 6)。

【 0 1 0 4 】

このような復元ルーチンおよび復元後処理が 1 ストライプ毎に繰り返して実行され (ステップ 1 1 0 5)、ディスク # 1 のすべての領域の復元が完了すれば (ステップ 1 6 0 9)、ディスク # 1 を担当 M a i n とする再構築処理を終了する (ステップ 1 1 0 4)。この時点で、担当 S u b であるディスク # 0 の現在位置もその上端に達しているため、ディスク # 0 の復元も同時に完了する。

30

4 . 改善案 4

改善案 4 では、図 1 1 の再構築処理はディスクの故障を契機として起動されるか、あるいは他の再構築処理から起動され、契機となったディスクが担当 M a i n に設定される。前者の契機では、R A I D グループあたり 1 つの再構築処理のみが起動される。したがって、2 番目のディスクが故障しても、既に R A I D グループとして再構築処理が起動済みであれば、新たな再構築処理は起動されない。

40

【 0 1 0 5 】

2 番目のディスクが故障すると、1 番目の故障ディスクの現在位置が 2 番目の故障ディスクの停止位置として設定され、2 番目の故障ディスクが担当 S u b として追加されて、再構築処理が続行される。そして、1 番目の故障ディスクの復元が完了すると、2 番目の故障ディスクを担当 M a i n として、その下端から停止位置までの再構築処理が行われる。改善案 4 の復元ルーチンのフローチャートは改善案 2 と同様である。

【 0 1 0 6 】

図 1 7 は、改善案 4 の復元後処理のフローチャートである。コントローラ 2 1 1 は、まず、復元ルーチン終了時の担当 M a i n の現在位置を担当 M a i n / S u b の復元先頭位置に設定し (ステップ 1 7 0 1)、以下の条件 b が満たされるか否かをチェックする (ス

50

テップ1702)。

条件b：担当Main以外に他の故障ディスクが存在し、担当Mainおよびその故障ディスクのいずれにも停止位置が設定されていない。

【0107】

条件bが満たされれば、担当Mainの復元先頭位置を他の故障ディスクの停止位置に設定し、その故障ディスクを担当Subに追加する(ステップ1706)。そして、担当Mainのすべての領域の復元が完了したか否かをチェックする(ステップ1703)。条件bが満たされなければ、そのままステップ1703の処理を行う。

【0108】

すべての領域の復元が完了していれば、他の故障ディスクがあるか否かをチェックし、そのような故障ディスクがあれば、それを担当Mainとする別の再構築処理を起動する(ステップ1707)。そして、復元処理終了と判定する(ステップ1708)。他の故障ディスクがなければ、別の再構築処理を起動することなく復元処理終了と判定する(ステップ1708)。

10

【0109】

すべての領域の復元が完了していなければ、次に、担当Mainの現在位置がその停止位置である否かをチェックする(ステップ1704)。現在位置が停止位置であれば、復元処理終了と判定する(ステップ1708)。

【0110】

現在位置が停止位置でない場合、および、停止位置が設定されていない場合は、そのまま復元処理継続と判定する(ステップ1705)。

20

例えば、図6に示したように、最初にディスク#0が故障すると、ディスク#0を担当Mainとして再構築処理が起動される。このとき、復元対象ディスクは1個であるから(図14のステップ1401)、復元方法は1ディスク故障の再構築となり(ステップ1403)、正常なディスク#1~#3のデータ/パリティのうち、ディスク#0のデータ/パリティを1ディスク故障の再構築により復元するために必要なものがリードされる(ステップ1405)。そして、リードされたデータ/パリティを用いてディスク#0のデータ/パリティが復元され、HS#0にライトされる(ステップ1407)。

【0111】

ディスク#0の現在位置は、その復元先頭位置として設定され(図17のステップ1701)、他の故障ディスクはないので(ステップ1702)、復元処理継続と判定される(ステップ1705)。このような復元ルーチンおよび復元後処理が1ストライプ毎に繰り返し実行される(図11のステップ1105)。

30

【0112】

次に、ディスク#1が故障すると、条件bが満たされ(ステップ1702)、ディスク#0の復元先頭位置がディスク#1の停止位置に設定され、ディスク#1が担当Subとして追加される(ステップ1706)。しかし、ディスク#0には停止位置が設定されていないので(ステップ1704)、復元処理継続と判定される(ステップ1705)。

【0113】

このとき、復元対象ディスクは2個となるから(ステップ1401)、復元方法は2ディスク故障の再構築処理となり(ステップ1404)、正常なディスク#2~#3のデータ/パリティのうち、ディスク#0および#1のデータ/パリティを2ディスク故障の再構築処理により復元するために必要なものがリードされる(ステップ1405)。そして、リードされたデータ/パリティを用いてディスク#0および#1のデータ/パリティが復元され、それぞれHS#0およびHS#1にライトされる(ステップ1407)。

40

【0114】

ディスク#0の現在位置は、ディスク#0および#1の復元先頭位置として設定され(ステップ1701)、ディスク#1には既に停止位置が設定されているため、条件bは満たされない(ステップ1702)。また、ディスク#0には停止位置が設定されていないので(ステップ1704)、復元処理継続と判定される(ステップ1705)。

50

【0115】

このような復元ルーチンおよび復元後処理が1ストライプ毎に繰り返し実行され(ステップ1105)、ディスク#0のすべての領域の復元が完了すれば(ステップ1703)、ディスク#1を担当Mainとする別の再構築処理が起動され(ステップ1707)、復元処理終了と判定される(ステップ1708)。これにより、ディスク#0を担当Mainとする再構築処理が終了する(ステップ1104)。この時点で、担当Subであるディスク#1の現在位置は、その上端に達している。

【0116】

次に、ディスク#1を担当Mainとする再構築処理において、ディスク#1の下端が現在位置として設定される(ステップ1101)。このとき、復元対象ディスクは1個であるから(ステップ1401)、復元方法は1ディスク故障の再構築となり(ステップ1403)、ディスク#0、#2、および#3のデータ/パリティのうち、ディスク#1のデータ/パリティを1ディスク故障の再構築により復元するために必要なものがリードされる(ステップ1405)。ただし、ディスク#0については、HS#0にライトされた復元済みのデータ/パリティがリードされる。

10

【0117】

そして、リードされたデータ/パリティを用いてディスク#1のデータ/パリティが復元され、HS#1にライトされる(ステップ1407)。

ディスク#1の現在位置は、その復元先頭位置として設定され(ステップ1701)、ディスク#1には既に停止位置が設定されているため、条件bは満たされない(ステップ1702)。また、ディスク#1の現在位置はその停止位置に到達していないので(ステップ1704)、復元処理継続と判定される(ステップ1705)。

20

【0118】

このような復元ルーチンおよび復元後処理が1ストライプ毎に繰り返し実行され(ステップ1105)、ディスク#1の現在位置が停止位置に到達する。このとき、ディスク#1の現在位置はその上端に達していないため、すべての領域の復元が完了していないと判定される(ステップ1703)。しかし、現在位置が停止位置に一致するので(ステップ1704)、復元処理終了と判定される(ステップ1708)。これにより、ディスク#1を担当Mainとする再構築処理が終了し(ステップ1104)、ディスク#1の復元が完了する。

30

5. 改善案5

改善案5では、改善案1と同様に、図11の再構築処理は各ディスクの故障を契機として起動され、契機となったディスクが担当Mainに設定される。

【0119】

2番目のディスクが故障すると、改善案4と同様に、1番目の故障ディスクの現在位置が2番目の故障ディスクの停止位置として設定され、2番目の故障ディスクが担当Subとして追加されて、再構築処理が実行される。それと同時に、2番目の故障ディスクを担当Mainとする再構築処理が起動され、1番目の故障ディスクを担当Mainとする再構築処理と並行して実行される。

【0120】

改善案5の復元ルーチンのフローチャートは改善案2と同様であり、復元後処理のフローチャートは図18のようになる。図18の復元後処理は、図17の復元後処理からステップ1707の処理を除いた構成を有する。

40

【0121】

例えば、図7に示したように、最初にディスク#0が故障すると、ディスク#0を担当Mainとして再構築処理が起動され、ディスク#1が故障するまで改善案4と同様の処理が行われる。

【0122】

次に、ディスク#1が故障すると、条件bが満たされ(ステップ1802)、ディスク#0の復元先頭位置がディスク#1の停止位置に設定され、ディスク#1が担当Subと

50

して追加される(ステップ1806)。しかし、ディスク#0には停止位置が設定されていないので(ステップ1804)、復元処理継続と判定される(ステップ1805)。

【0123】

このとき、復元対象ディスクは2個となるから(ステップ1401)、復元方法は2ディスク故障の再構築処理となり(ステップ1404)、正常なディスク#2~#3のデータ/パリティのうち、ディスク#0および#1のデータ/パリティを2ディスク故障の再構築処理により復元するために必要なものがリードされる(ステップ1405)。そして、リードされたデータ/パリティを用いてディスク#0および#1のデータ/パリティが復元され、それぞれHS#0およびHS#1にライトされる(ステップ1407)。

【0124】

ディスク#0の現在位置は、ディスク#0および#1の復元先頭位置として設定され(ステップ1801)、ディスク#1には既に停止位置が設定されているため、条件bは満たされない(ステップ1802)。また、ディスク#0には停止位置が設定されていないので(ステップ1804)、復元処理継続と判定される(ステップ1805)。

【0125】

このような復元ルーチンおよび復元後処理が1ストライプ毎に繰り返し実行され(ステップ1105)、ディスク#0のすべての領域の復元が完了すれば(ステップ1803)、復元処理終了と判定される(ステップ1807)。これにより、ディスク#0を担当Mainとする再構築処理が終了する(ステップ1104)。この時点で、担当Subであるディスク#1の現在位置は、その上端に達している。

【0126】

さらに、ディスク#1の故障時には、ディスク#1を担当Mainとしてもう1つの再構築処理が起動され、ディスク#1の下端が現在位置として設定される(ステップ1101)。このとき、復元対象ディスクは1個であるから(ステップ1401)、復元方法は1ディスク故障の再構築となり(ステップ1403)、ディスク#0、#2、および#3のデータ/パリティのうち、ディスク#1のデータ/パリティを1ディスク故障の再構築により復元するために必要なものがリードされる(ステップ1405)。ただし、ディスク#0については、HS#0にライトされた復元済みのデータ/パリティがリードされる。

【0127】

そして、リードされたデータ/パリティを用いてディスク#1のデータ/パリティが復元され、HS#1にライトされる(ステップ1407)。

ディスク#1の現在位置は、その復元先頭位置として設定され(ステップ1801)、ディスク#1には既に停止位置が設定されているため、条件bは満たされない(ステップ1802)。また、ディスク#1の現在位置はその停止位置に到達していないので(ステップ1804)、復元処理継続と判定される(ステップ1805)。

【0128】

このような復元ルーチンおよび復元後処理が1ストライプ毎に繰り返し実行され(ステップ1105)、ディスク#1の現在位置が停止位置に到達する。このとき、ディスク#1の現在位置はその上端に達していないため、すべての領域の復元が完了していないと判定される(ステップ1803)。しかし、現在位置が停止位置に一致するので(ステップ1804)、復元処理終了と判定される(ステップ1807)。これにより、ディスク#1を担当Mainとする再構築処理が終了する(ステップ1104)。

【0129】

ディスク#0を担当Mainとする2ディスク故障の再構築処理と、ディスク#1を担当Mainとする1ディスク故障の再構築処理は並行して実行され、両方の再構築処理が終了した時点でディスク#1の復元が完了する。

6. 改善案6

改善案1~5の再構築処理において、各ディスクのストリップや論理ブロックのような所定領域毎にその復元状態を表すビットマップを、制御情報として追加する。ディスク全

10

20

30

40

50

体に対する進捗の制御は、改善案 1 ~ 5 により行われる。

【0130】

コントローラ 211 は、再構築処理または復元ルーチンの実行時に、ビットマップ中の復元位置に対応するビット情報を参照する。そして、図 19 に示すように、リード/ライトアクセスの延長等により既に復元済みであれば、復元処理をスキップする。これにより、余計な復元処理のコストを削減することができる。

【0131】

図 20 は、コントローラ 211 のプロセッサ 221 が処理に用いるプログラムおよびデータの提供方法を示している。情報処理装置等の外部装置 1801 や可搬記録媒体 1802 に格納されたプログラムおよびデータは、RAID 装置 202 のメモリ 222 にロード

10

【0132】

外部装置 1801 は、そのプログラムおよびデータを搬送する搬送信号を生成し、通信ネットワーク上の任意の伝送媒体を介して RAID 装置 202 に送信する。可搬記録媒体 1802 は、メモリカード、フレキシブルディスク、光ディスク、光磁気ディスク等の任意のコンピュータ読み取り可能な記録媒体である。プロセッサ 221 は、そのデータを用いてそのプログラムを実行し、必要な処理を行う。

【0133】

図 21 および 22 は、ストレージシステムの別の構成例を示している。図 21 は、ホスト装置に実装されたホストバスアダプタが再構築処理を行う例を示しており、図 22 は、ホスト装置に実装されたソフトウェアが再構築処理を行う例を示している。いずれの構成においても、必要なプログラムおよびデータは、RAID 装置 202 の場合と同様にして提供される。

20

【0134】

図 21 のストレージシステムは、ホスト装置 1901 および Disk # 0 ~ Disk # 3 からなり、ホスト装置 1901 は、ホストバスアダプタ 1911 を備える。ホストバスアダプタ 1911 は、プロセッサ 1921、メモリ 1922、およびキャッシュメモリ 1923 を備え、Disk # 0 ~ Disk # 3 の故障時の再構築処理を行う。このとき、プロセッサ 1921 は、メモリ 1922 に格納されたプログラムを実行することにより、上述した再構築処理を行う。

30

【0135】

図 22 のストレージシステムは、ホスト装置 2001 および Disk # 0 ~ Disk # 3 からなる。ホスト装置 2001 は、プロセッサ 2011 およびメモリ 2012、2013 を備え、Disk # 0 ~ Disk # 3 の故障時の再構築処理を行う。このとき、プロセッサ 2011 は、メモリ 2012 に格納されたプログラムを実行することにより、メモリ 2013 上で上述した再構築処理を行う。

【0136】

なお、以上の実施形態では、ディスク装置として磁気ディスク装置が用いられているが、本発明は、光ディスク装置、光磁気ディスク装置等の他のディスク装置や、テープ装置のような他の記憶装置を用いたストレージシステムに対しても、適用可能である。

40

【図面の簡単な説明】

【0137】

【図 1】本発明のストレージ制御装置の原理図である。

【図 2】第 1 のストレージシステムの構成図である。

【図 3】再構築処理の改善案を示す図である。

【図 4】改善案 1 を示す図である。

【図 5】改善案 2 を示す図である。

【図 6】改善案 4 を示す図である。

【図 7】改善案 5 を示す図である。

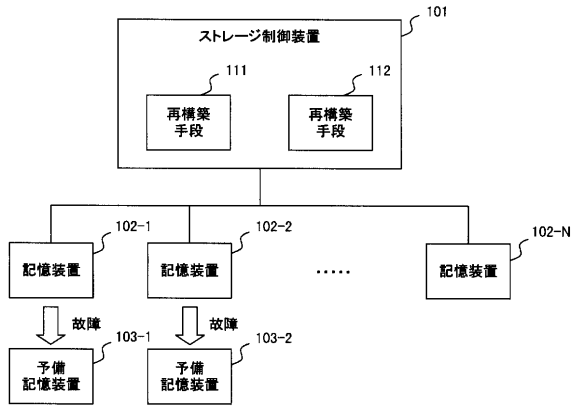
【図 8】改善案 6 を示す図である。

50

- 【図9】MainとSubの現在位置を示す図である。
- 【図10】各ディスクの復元状態を示す図である。
- 【図11】再構築処理のフローチャートである。
- 【図12】改善案1の復元ルーチンのフローチャートである。
- 【図13】改善案1の復元後処理のフローチャートである。
- 【図14】改善案2～5の復元ルーチンのフローチャートである。
- 【図15】改善案2の復元後処理のフローチャートである。
- 【図16】改善案3の復元後処理のフローチャートである。
- 【図17】改善案4の復元後処理のフローチャートである。
- 【図18】改善案5の復元後処理のフローチャートである。 10
- 【図19】改善案6の復元状態を示す図である。
- 【図20】プログラムおよびデータの提供方法を示す図である。
- 【図21】第2のストレージシステムの構成図である。
- 【図22】第3のストレージシステムの構成図である。
- 【図23】1ディスク故障のデータ復元を示す図である。
- 【図24】2ディスク故障のデータ復元を示す図である。
- 【図25】2ディスク故障の再構築処理を示す図である。
- 【符号の説明】
- 【0138】
- 101 ストレージ制御装置 20
- 102 - 1、102 - 2、102 - N 記憶装置
- 103 - 1、103 - 2 予備記憶装置
- 111、112 再構築手段
- 201、1901、2001 ホスト装置
- 202 RAID装置
- 211 コントローラ
- 211、1921、2011 プロセッサ
- 222、1922、2012、2013 メモリ
- 223、1923 キャッシュメモリ
- 1801 外部装置 30
- 1802 可搬記録媒体
- 1911 ホストバスアダプタ
- Disk0、Disk1、Disk2、Disk3 ディスク

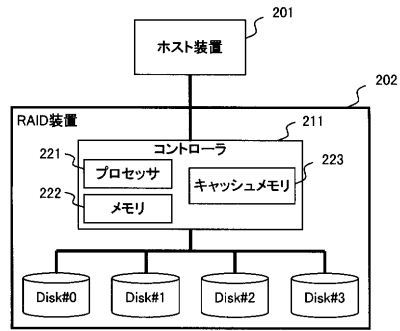
【 図 1 】

本発明のストレージ制御装置の原理図



【 図 2 】

第1のストレージシステムの構成図



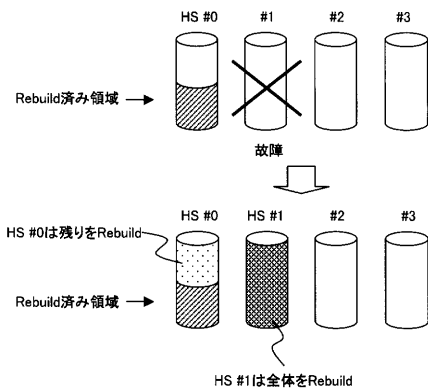
【 図 3 】

再構築処理の改善案を示す図

	方式	進捗管理方法	工数・難易度
1	個別実行案	現状通り	容易
2	同じ進捗位置になるまで待つ案	現状通り	容易
3	混在案	現状通り	中度
4	二重故障部分を先に復元する案	現状に加えて、復元完了位置を管理	通常処理に影響あり
5	二重故障復元と、復元済みHSを用いた復元を並行して実施する案	現状に加えて、復元完了位置を管理	
6	ランダム案	ビットマップ管理	工数がかかる 通常処理に影響あり

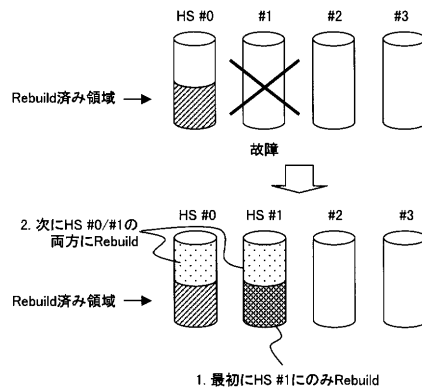
【 図 4 】

改善案1を示す図



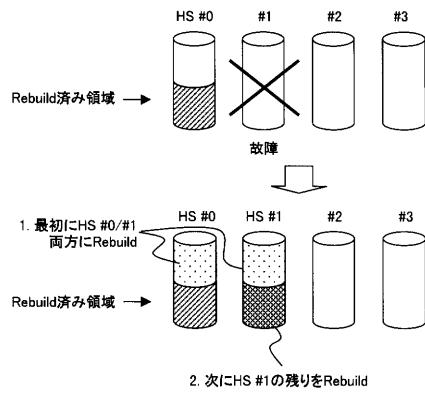
【 図 5 】

改善案2を示す図



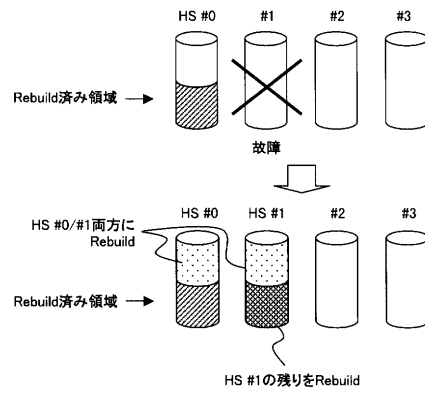
【 図 6 】

改善案4を示す図



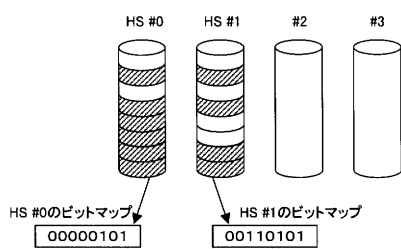
【 図 7 】

改善案5を示す図



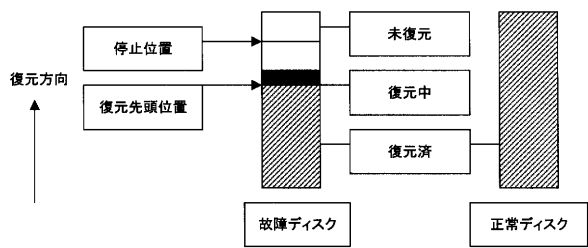
【 図 8 】

改善案6を示す図



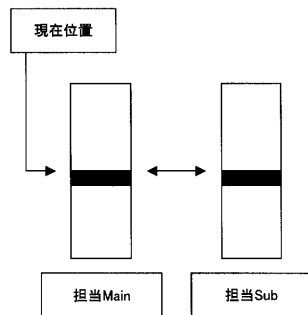
【 図 10 】

各ディスクの復元状態を示す図



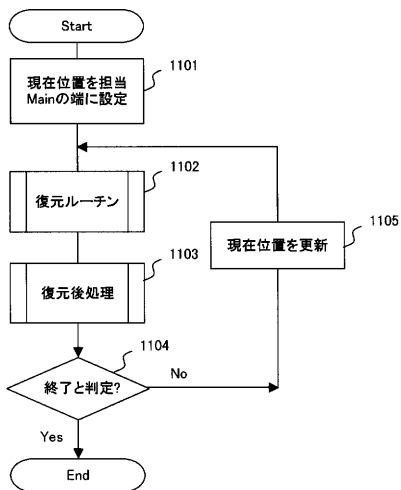
【 図 9 】

MainとSubの現在位置を示す図



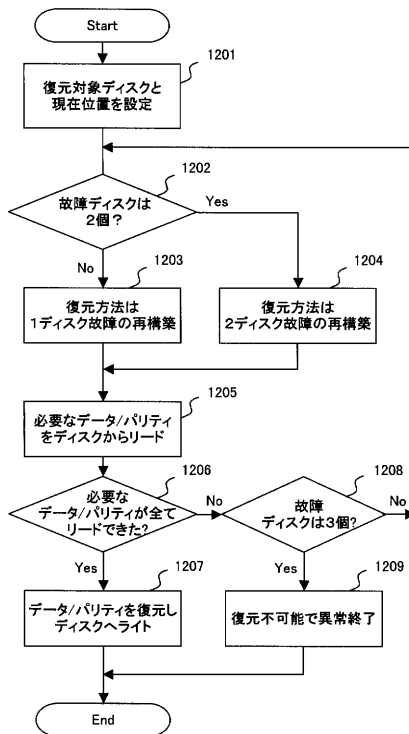
【図11】

再構築処理のフローチャート



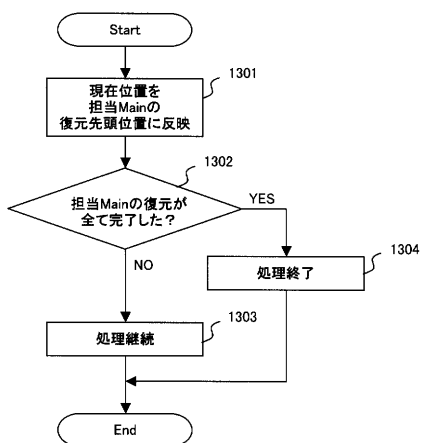
【図12】

改善案1の復元ルーチンのフローチャート



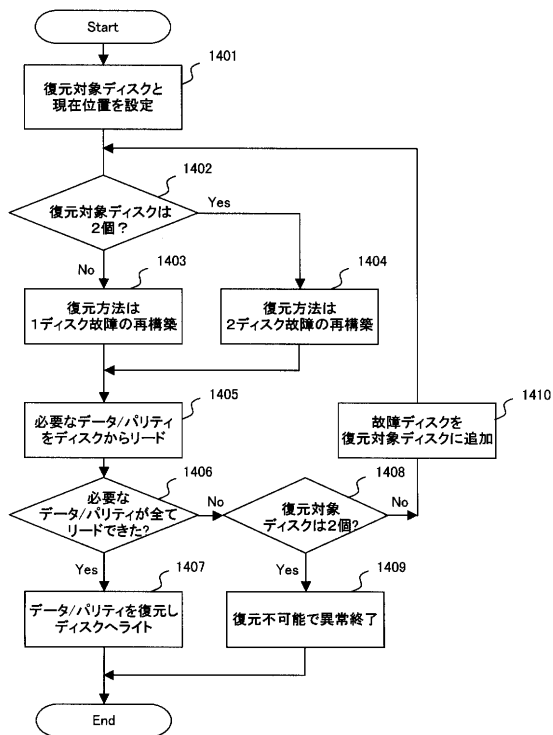
【図13】

改善案1の復元後処理のフローチャート



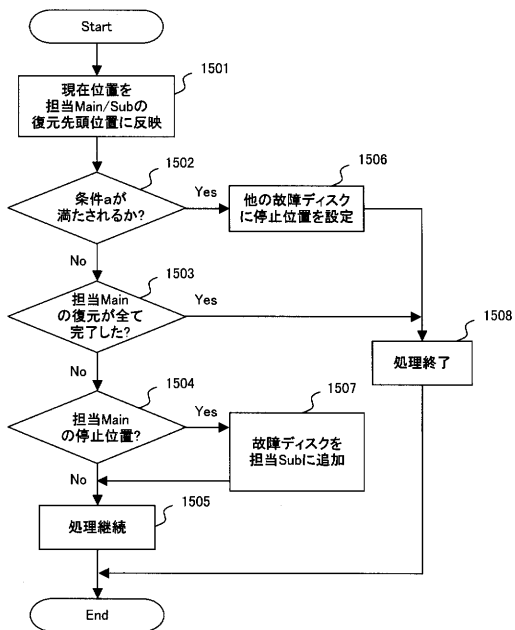
【図14】

改善案2~5の復元ルーチンのフローチャート



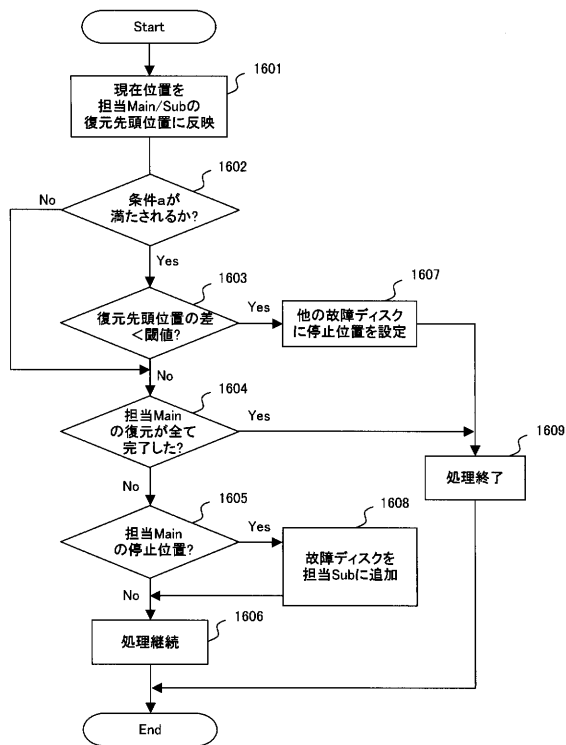
【 図 1 5 】

改善案2の復元後処理のフローチャート



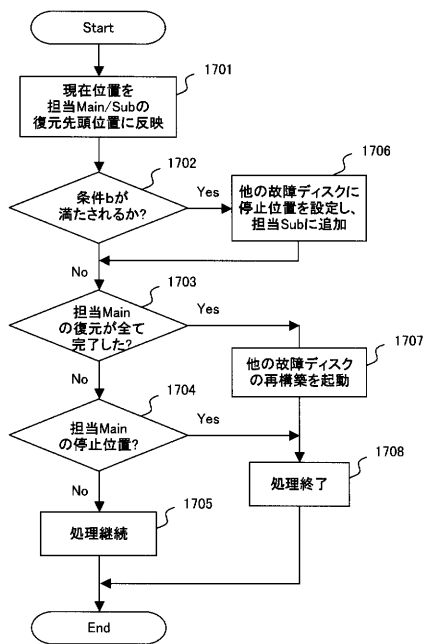
【 図 1 6 】

改善案3の復元後処理のフローチャート



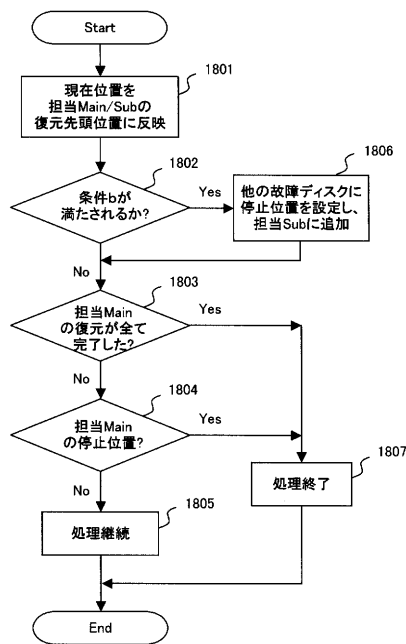
【 図 1 7 】

改善案4の復元後処理のフローチャート



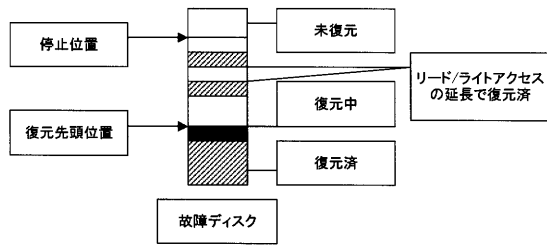
【 図 1 8 】

改善案5の復元後処理のフローチャート



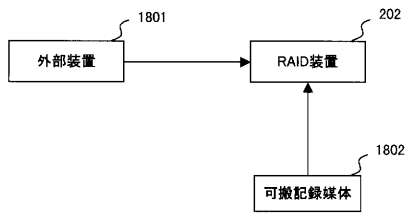
【図19】

改善案6の復元状態を示す図



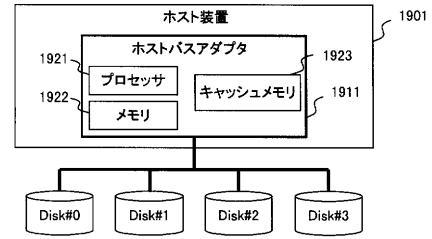
【図20】

プログラムおよびデータの提供方法を示す図



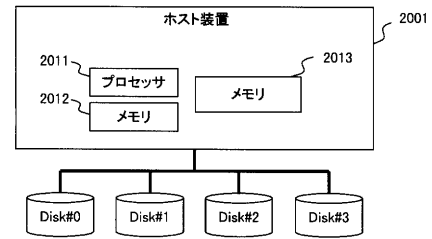
【図21】

第2のストレージシステムの構成図



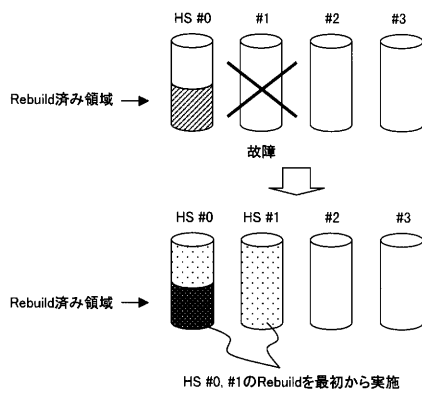
【図22】

第3のストレージシステムの構成図



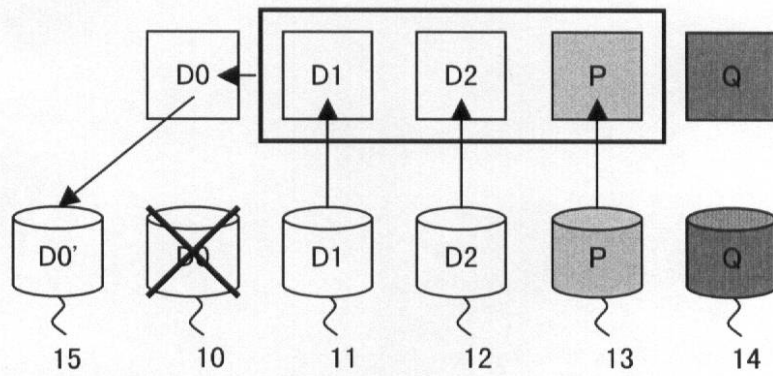
【図25】

2ディスク故障の再構築処理を示す図



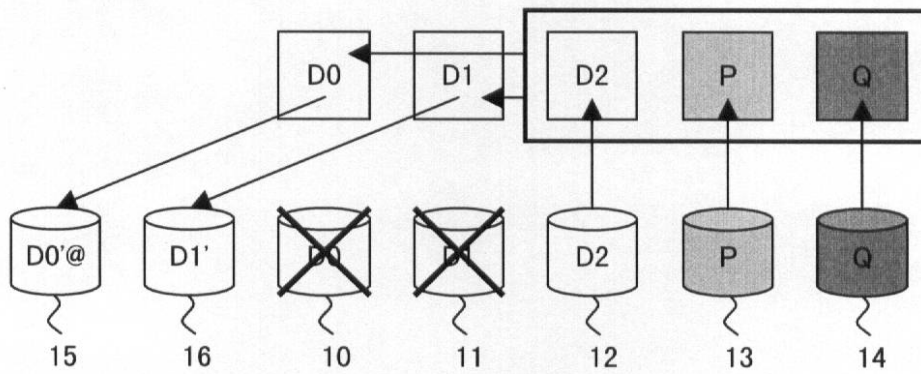
【 図 2 3 】

1ディスク故障のデータ復元を示す図



【 図 2 4 】

2ディスク故障のデータ復元を示す図



フロントページの続き

- (72)発明者 大黒谷 秀治郎
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
- (72)発明者 池内 和彦
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
- (72)発明者 高 橋 秀夫
神奈川県川崎市中原区上小田中4丁目1番1号 株式会社富士通コンピュータテクノロジーズ内
- (72)発明者 紺田 與志仁
神奈川県川崎市中原区上小田中4丁目1番1号 株式会社富士通コンピュータテクノロジーズ内
- (72)発明者 佐藤 靖丈
神奈川県川崎市中原区上小田中4丁目1番1号 株式会社富士通コンピュータテクノロジーズ内
- (72)発明者 越智 弘昭
神奈川県川崎市中原区上小田中4丁目1番1号 株式会社富士通コンピュータテクノロジーズ内
- (72)発明者 牧野 司
神奈川県川崎市中原区上小田中4丁目1番1号 株式会社富士通コンピュータテクノロジーズ内
- (72)発明者 久保田 典秀
神奈川県川崎市中原区上小田中4丁目1番1号 株式会社富士通コンピュータテクノロジーズ内
- Fターム(参考) 5B001 AB01 AB02 AD04
5B018 HA04 HA14 MA12
5B065 BA01 CA30 EA03 EA24