

(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2011-128917

(P2011-128917A)

(43) 公開日 平成23年6月30日 (2011.6.30)

(51) Int.Cl. F 1 テーマコード (参考)
G 0 6 F 3/06 (2006.01) G 0 6 F 3/06 3 0 4 B 5 B 0 6 5
 G 0 6 F 3/06 3 0 5 F

審査請求 未請求 請求項の数 7 O L (全 49 頁)

(21) 出願番号 特願2009-287068 (P2009-287068)
 (22) 出願日 平成21年12月18日 (2009.12.18)

(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 100092152
 弁理士 服部 毅巖
 (72) 発明者 野口 泰生
 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 櫻井 英樹
 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
 Fターム(参考) 5B065 BA01 EA02 EA24 EA31

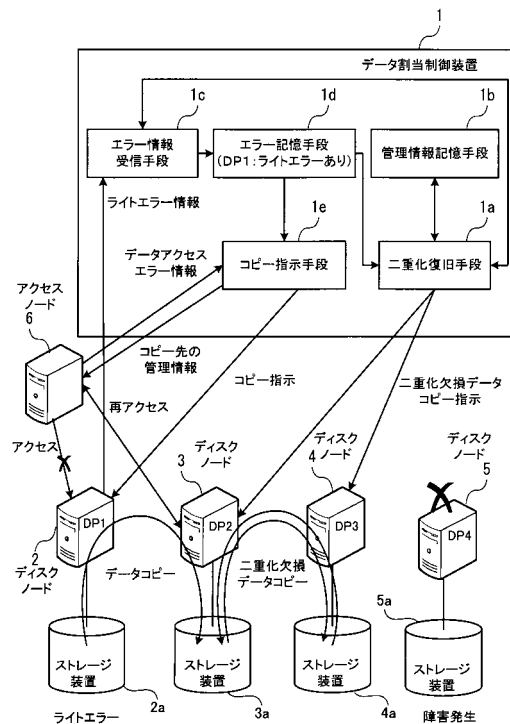
(54) 【発明の名称】 データ割当制御プログラム、データ割当制御方法、およびデータ割当制御装置

(57) 【要約】

【課題】 ディスクノードの障害発生時のデータの喪失を抑制する。

【解決手段】 二重化復旧手段 1 a は、二重化欠損データのコピーをディスクノードに指示し、二重化復旧処理を行う。エラー情報受信手段 1 c は、二重化復旧処理中に、ストレージ装置 2 a におけるライトエラーの発生を示すライトエラー情報を受け取ると、ライトエラーが発生したストレージ装置 2 a の識別情報をエラー記憶手段 1 d に格納する。コピー指示手段 1 e は、ライトエラーが発生したストレージ装置 2 a 内のアクセスが行われたデータに冗長データが存在しない場合、ディスクノード 2 に対してアクセスが行われたデータのコピーを指示する。

【選択図】 図 1



【特許請求の範囲】**【請求項 1】**

接続されたストレージ装置内のデータを管理する複数のディスクノードに対する管理対象データの割り当て指示を、コンピュータに実行させるデータ割当制御プログラムにおいて、

前記コンピュータに、

複数のストレージ装置に格納されたデータのうち同一内容の冗長データが存在しない二重化欠損データのコピーを、二重化欠損データを管理するディスクノードに指示する二重化復旧処理を行い、

前記二重化復旧処理中に、ストレージ装置に対するデータのライトエラーを示すライトエラー情報を受け取ると、ライトエラーが発生したストレージ装置の識別情報をエラー記憶手段に格納し、

前記エラー記憶手段を参照してライトエラーが発生したストレージ装置を判断し、ライトエラーが発生したストレージ装置に格納されているデータをコピー対象データとし、ライトエラーが発生していないストレージ装置内に前記コピー対象データの冗長データが存在しない場合、前記コピー対象データを管理するディスクノードに対して、ライトエラーが発生していないストレージ装置への前記コピー対象データのコピーを指示する、

処理を実行させることを特徴とするデータ割当制御プログラム。

【請求項 2】

前記二重化復旧処理では、二重化欠損データを順次選択し、選択した二重化欠損データを管理するディスクノードに対し、選択した二重化欠損データが格納されているストレージ装置と異なるストレージ装置であり、かつライトエラーが発生していないストレージ装置への、選択した二重化欠損データのコピーを指示することを特徴とする請求項 1 記載のデータ割当制御プログラム。

【請求項 3】

前記二重化復旧処理では、前記複数のストレージ装置内の二重化欠損データを調査し、二重化欠損データの格納場所を示す管理情報を管理情報記憶手段に格納し、前記管理情報記憶手段内の管理情報に示される二重化欠損データを順次選択して、選択した二重化欠損データを管理するディスクノードに対して選択した二重化欠損データのコピーを指示し、

前記コピー対象データのコピー指示の際には、前記コピー対象データのコピーが完了すると、コピー先のデータの管理情報を二重化欠損データとして前記管理情報記憶手段に格納する、

ことを特徴とする請求項 1 記載のデータ割当制御プログラム。

【請求項 4】

前記二重化復旧処理では、前記管理情報記憶手段内の管理情報で示されるすべての二重化欠損データのコピーが完了すると前記エラー記憶手段を参照してライトエラーが発生したストレージ装置の有無を判断し、ライトエラーが発生したストレージ装置がある場合、ライトエラーが発生したストレージ装置へのアクセスを停止することにより生じる二重化欠損データの管理情報を前記管理情報記憶手段に追加し、二重化復旧処理を続行する、

ことを特徴とする請求項 3 記載のデータ割当制御プログラム。

【請求項 5】

前記コピー対象データのコピー指示では、ライトエラーが発生しているストレージ装置内のデータに対するアクセスで発生したエラーを示すデータアクセスエラー情報を受け取ると、エラー発生時のアクセス対象となっていたデータをコピー対象データとし、前記コピー対象データのコピーが完了すると、コピー先のデータの管理情報を、新たなアクセス先として前記データアクセスエラー情報を送信した装置に応答することを特徴とする請求項 1 記載のデータ割当制御プログラム。

【請求項 6】

接続されたストレージ装置内のデータを管理する複数のディスクノードに対する管理対象データの割り当て指示をコンピュータで実行するデータ割当制御方法において、

10

20

30

40

50

前記コンピュータが、

複数のストレージ装置に格納されたデータのうち同一内容の冗長データが存在しない二重化欠損データのコピーを、二重化欠損データを管理するディスクノードに指示する二重化復旧処理を行い、

前記二重化復旧処理中に、ストレージ装置に対するデータのライトエラーを示すライトエラー情報を受け取ると、ライトエラーが発生したストレージ装置の識別情報をエラー記憶手段に格納し、

前記エラー記憶手段を参照してライトエラーが発生したストレージ装置を判断し、ライトエラーが発生したストレージ装置に格納されているデータをコピー対象データとし、ライトエラーが発生していないストレージ装置内に前記コピー対象データの冗長データが存在しない場合、前記コピー対象データを管理するディスクノードに対して、ライトエラーが発生していないストレージ装置への前記コピー対象データのコピーを指示する、

ことを特徴とするデータ割当制御方法。

【請求項 7】

接続されたストレージ装置内のデータを管理する複数のディスクノードに対する管理対象データの割り当て指示を行うデータ割当制御装置において、

複数のストレージ装置に格納されたデータのうち同一内容の冗長データが存在しない二重化欠損データのコピーを、二重化欠損データを管理するディスクノードに指示する二重化復旧処理を行う二重化復旧処理手段と、

前記二重化復旧処理中に、ストレージ装置に対するデータのライトエラーを示すライトエラー情報を受け取ると、ライトエラーが発生したストレージ装置の識別情報をエラー記憶手段に格納するエラー情報受信手段と、

前記エラー記憶手段を参照してライトエラーが発生したストレージ装置を判断し、ライトエラーが発生したストレージ装置に格納されているデータをコピー対象データとし、ライトエラーが発生していないストレージ装置内に前記コピー対象データの冗長データが存在しない場合、前記コピー対象データを管理するディスクノードに対して、ライトエラーが発生していないストレージ装置への前記コピー対象データのコピーを指示するコピー指示手段と、

を有することを特徴とするデータ割当制御装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は複数のディスクノードに対する管理対象データの割り当てを行うデータ割当制御プログラム、データ割当制御方法、およびデータ割当制御装置に関する。

【背景技術】

【0002】

データを複数のコンピュータで分散管理するシステムとして、マルチノードストレージシステムがある。マルチノードストレージシステムは、ネットワークに接続した複数のディスクノードと制御ノードとを有している。マルチノードストレージシステムでは、例えば制御ノードの管理の下、仮想的なディスク（論理ディスク）に格納するデータを複数のディスクノードに分散格納される。

【0003】

マルチノードストレージシステムでは、例えば論理ディスクがセグメント単位に分割されている。この場合、ディスクノードが有するストレージデバイスの記憶領域はスライス単位に分割される。このスライスは、例えば、セグメントと同サイズとされる。制御ノードによって、論理ディスクの各セグメントに対して、ストレージデバイスのスライスが割り当てられる。セグメントへのスライスの割り当て関係は、制御ノードからデータアクセスを行うコンピュータ（アクセスノード）に通知される。そして、アクセスノードから、セグメントを格納先として指定したデータが、そのセグメントに割り当てられたスライスを有するディスクノードに送られる。ディスクノードは、自己の管理するスライスが割り

10

20

30

40

50

当てられたセグメントのデータを受け取り、ストレージ装置内の該当スライスに格納する。

【0004】

このようなマルチノードストレージシステムによれば、ネットワークにディスクノードを追加することで管理可能なデータ容量を増やすことができる。そのため、システムの拡張が容易になる。

【0005】

また、マルチノードストレージシステムでは、セグメントに対して複数のスライスを割り当てることができる。1つのセグメントに2つのスライスを割り当てた場合、一方をプライマリスライス、他方をセカンダリスライスとする。プライマリスライスは、アクセスノードが直接リードまたはライトするスライスである。セカンダリスライスは、プライマリスライスへのライト時にライトデータのミラーリング先となるスライスである。このように、2つのスライスを用いたミラーリングによりデータを二重化することで、データの冗長性が確保される。

【0006】

なお、ストレージ装置の故障などでディスクノードの異常が検出されると、データの二重化復旧処理が実行される。例えば、1つのノードが異常となった場合、まず異常ノードがシステムから切り離される。次に異常ノードが持っていたデータと同じデータを持つ別のノードをコピー元とし、別のノードをコピー先として、データの冗長性が復元される。

【先行技術文献】

【特許文献】

【0007】

【特許文献1】国際公開第2004/104845号

【特許文献2】特開2005-301594号公報

【発明の概要】

【発明が解決しようとする課題】

【0008】

しかし、システムから切り離されたディスクノードへは、他のノードからアクセスできなくなる。そのためストレージ装置の故障が発生したディスクノードの切り離しを常に行うと、データの二重化復旧処理中に新たに別のディスクノードに接続されたストレージ装置で故障が発生した場合に、データを失う可能性がある。例えば、最初に故障したストレージ装置と、二重化復旧処理中に新たに故障したストレージ装置とのそれぞれに二重化して格納されていたデータは、新たに故障したストレージ装置が接続されたディスクノードの切り離しにより失われる可能性がある。

【0009】

本発明はこのような点に鑑みてなされたものであり、ディスクノードの障害発生時のデータの喪失を抑制することができるデータ割当制御プログラム、データ割当制御方法、およびデータ割当制御装置を提供することを目的とする。

【課題を解決するための手段】

【0010】

上記課題を解決するために、接続されたストレージ装置内のデータを管理する複数のディスクノードに対する管理対象データの割り当て指示を、コンピュータに実行させるデータ割当制御プログラムが提供される。このデータ割当制御プログラムは、コンピュータに以下の処理を実行させる。

【0011】

コンピュータは、複数のストレージ装置に格納されたデータのうち同一内容の冗長データが存在しない二重化欠損データのコピーを、二重化欠損データを管理するディスクノードに指示する二重化復旧処理を行う。コンピュータは、二重化復旧処理中に、ストレージ装置に対するデータのライトエラーを示すライトエラー情報を受け取ると、ライトエラーが発生したストレージ装置の識別情報をエラー記憶手段に格納する。コンピュータは、エ

10

20

30

40

50

ラー記憶手段を参照してライトエラーが発生したストレージ装置を判断し、ライトエラーが発生したストレージ装置に格納されているデータをコピー対象データとし、ライトエラーが発生していないストレージ装置内にコピー対象データの冗長データが存在しない場合、コピー対象データを管理するディスクノードに対して、ライトエラーが発生していないストレージ装置へのコピー対象データのコピーを指示する。

【0012】

また上記課題を解決するために、上記データ割当制御プログラムに基づいてコンピュータが実行する処理と同様の処理を行うデータ割当制御方法が提供される。さらに上記課題を解決するために、上記データ割当制御プログラムを実行するコンピュータが有する機能と同様の機能を有するデータ割当制御装置が提供される。

10

【発明の効果】

【0013】

二重化復旧処理中にライトエラーとなったストレージ装置内のデータの喪失を抑制することができる。

【図面の簡単な説明】

【0014】

【図1】実施の形態の概要を示す図である。

【図2】第2の形態のマルチノードストレージシステム構成の一例を示す図である。

【図3】第2の形態に用いる制御ノードのハードウェアの一構成例を示す図である。

【図4】論理ディスクのデータ構造の一例を示す図である。

20

【図5】第2の実施の形態に係るマルチノードストレージシステムの各装置の機能を示すブロック図である。

【図6】ストレージ装置のデータ構造の一例を示す図である。

【図7】メタデータ記憶部のデータ構造例を示す図である。

【図8】論理ディスクメタデータ記憶部のデータ構造の一例を示す図である。

【図9】エラー記憶部のデータ構造の一例を示す図である。

【図10】ディスクノードでのライトエラー検出時に装置間で送受信される情報の例を示す図である。

【図11】故障したディスクノードの切り離し処理の一例を示す図である。

【図12】リカバリ処理手順の一例を示すシーケンス図である。

30

【図13】リカバリ処理実行中のライトエラー発生時の処理手順の一例を示すシーケンス図である。

【図14】ディスクノードでのライトエラー検出および通知処理の手順の一例を示すフローチャートである。

【図15】エラーメッセージを受信したときのエラー処理手順の一例を示すフローチャートである。

【図16】リカバリ処理手順の一例を示すフローチャートである。

【図17】ライトエラーが発生したディスクノードへのアクセス要求時の処理手順を示すシーケンス図である。

【図18】アクセス処理手順の一例を示すフローチャートである。

40

【図19】ライトエラーが発生したディスクノードへのミラーライト処理の手順を示すシーケンス図である。

【図20】ライトエラーが発生したディスクノードへのスライスコピー処理時の処理手順を示すシーケンス図である。

【図21】スライスコピーのデータを受信したディスクノードの処理手順の一例を示すフローチャートである。

【図22】メタデータ照会要求時のスライス割当処理手順の一例を示すフローチャートである。

【図23】スライスコピーエラー時のスライス割り当て処理手順の一例を示すフローチャートである。

50

【図 2 4】プライマリスライスとセカンダリスライスとが割り当てられているセグメントの割り当て変更処理手順の一例を示すフローチャートである。

【図 2 5】プライマリスライスとセカンダリスライスとが割り当てられているセグメントの割り当て変更処理例を示す図である。

【図 2 6】プライマリスライスとリザーブスライスとが割り当てられているセグメントの割り当て変更処理手順の一例を示すフローチャートである。

【図 2 7】プライマリスライスとリザーブスライスとが割り当てられているセグメントの割り当て変更処理例を示す図である。

【図 2 8】シングルプライマリスライスが割り当てられているセグメントの割り当て変更処理手順の一例を示すフローチャートである。

【図 2 9】シングルプライマリスライスが割り当てられているセグメントの割り当て変更処理例を示す図である。

【図 3 0】スライスコピーを伴うスライス割当処理手順の一例を示すシーケンス図である。

【図 3 1】スライスコピーを伴わないスライス割当処理手順の一例を示すシーケンス図である。

【発明を実施するための形態】

【0015】

以下、本実施の形態について図面を参照して説明する。

〔第 1 の実施の形態〕

図 1 は、実施の形態の概要を示す図である。複数のディスクノード 2 ~ 5 には、それぞれストレージ装置 2 a , 3 a , 4 a , 5 a が接続されている。ディスクノード 2 ~ 5 は、接続されたストレージ装置 2 a , 3 a , 4 a , 5 a 内のデータを管理する。データ割当制御装置 1 は、ディスクノード 2 ~ 5 に対する管理対象データの割り当て指示を行う。例えば、データ割当制御装置 1 は、ネットワークを介してディスクノード 2 ~ 5 に接続されている。

【0016】

データ割当制御装置 1 は、既に割り当てられているデータの割り当て先を変更する場合、データのコピーをディスクノードに指示する。データ割り当て先の変更は、例えば、ストレージ装置に対するデータ書き込み時のライトエラーに応じて行われる。ライトエラー発生時のデータコピーを適切に指示するため、データ割当制御装置 1 は、二重化復旧手段 1 a、管理情報記憶手段 1 b、エラー情報受信手段 1 c、エラー記憶手段 1 d、およびコピー指示手段 1 e を有している。

【0017】

二重化復旧手段 1 a は、複数のストレージ装置 2 a , 3 a , 4 a , 5 a に格納されたデータのうち同一内容の冗長データが存在しない二重化欠損データのコピーを、二重化欠損データを管理するディスクノードに指示する二重化復旧処理を行う。例えば二重化復旧手段 1 a は、複数のストレージ装置 2 a , 3 a , 4 a , 5 a 内の二重化欠損データを調査し、二重化欠損データの格納場所を示す管理情報を管理情報記憶手段 1 b に格納する。二重化復旧手段 1 a は、管理情報記憶手段 1 b 内の管理情報に示される二重化欠損データを順次選択して、選択した二重化欠損データを管理するディスクノードに対して選択した二重化欠損データのコピーを指示する。選択した二重化欠損データのコピーの指示は、例えば、選択した二重化欠損データが格納されているストレージ装置と異なるストレージ装置であり、かつライトエラーが発生していないストレージ装置が、コピー先として指定される。

【0018】

なお二重化復旧手段 1 a は、管理情報記憶手段 1 b 内の管理情報で示されるすべての二重化欠損データのコピーが完了すると、ライトエラーが発生したストレージ装置を切り離し、その結果生じる二重化欠損データのコピーを行うことができる。例えば、二重化復旧手段 1 a は、管理情報記憶手段 1 b 内の管理情報で示されるすべての二重化欠損データの

10

20

30

40

50

コピーが完了すると、エラー記憶手段 1 d を参照してライトエラーが発生したストレージ装置の有無を判断する。ライトエラーが発生したストレージ装置がある場合、二重化復旧手段 1 a は、ライトエラーが発生したストレージ装置へのアクセスを停止することにより生じる二重化欠損データの管理情報を管理情報記憶手段 1 b に追加する。そして、二重化復旧手段 1 a は、管理情報記憶手段 1 b に追加された二重化復旧処理の管理情報に基づいて、二重化復旧処理を続行する。

【 0 0 1 9 】

管理情報記憶手段 1 b は、少なくとも二重化欠損データの格納場所を示す管理情報を記憶する。

エラー情報受信手段 1 c は、二重化復旧処理中にストレージ装置に対するデータのライトエラーを示すライトエラー情報を受け取ると、ライトエラーが発生したストレージ装置の識別情報をエラー記憶手段 1 d に格納する。またエラー情報受信手段 1 c は、二重化復旧処理中以外の期間に、ストレージ装置に対するデータのライトエラーを示すライトエラー情報を受け取る場合もある。この場合、エラー情報受信手段 1 c は、二重化復旧手段 1 a に、ライトエラーが発生したストレージ装置へのアクセスを停止することにより生じる二重化欠損データの二重化復旧処理を開始させる。

【 0 0 2 0 】

エラー記憶手段 1 d は、ライトエラーが発生したストレージ装置の識別情報を記憶する。なお、1台のディスクノードに1台のストレージ装置が接続されている場合、ディスクノードの識別情報を、ストレージ装置の識別情報として用いることもできる。また1台のディスクノードに複数のストレージ装置が接続されている場合、ディスクノードの識別情報と、ディスクノード内でのローカルなストレージ装置の識別情報とを組み合わせ、ストレージ装置の識別情報とすることもできる。

【 0 0 2 1 】

コピー指示手段 1 e は、エラー記憶手段 1 d を参照してライトエラーが発生したストレージ装置を判断する。そしてコピー指示手段 1 e は、ライトエラーが発生したストレージ装置に格納されているデータをコピー対象データとする。コピー指示手段 1 e は、ライトエラーが発生していないストレージ装置内にコピー対象データの冗長データが存在しない場合、コピー対象データを管理するディスクノードに対して、コピー対象データのコピーを指示する。コピー指示では、ライトエラーが発生していないストレージ装置がコピー先として指定される。例えば、コピー指示手段 1 e は、ライトエラーが発生したストレージ装置内のデータに対するアクセスで発生したエラーを示すデータアクセスエラー情報を他の装置から受け取ると、アクセス対象であったデータをコピー対象データとする。なお、コピー指示手段 1 e は、コピー対象データのコピーが完了すると、コピー先のデータの管理情報を二重化欠損データとして管理情報記憶手段 1 b に格納する。

【 0 0 2 2 】

なおコピー指示手段 1 e は、例えば、コピー対象データのコピーが完了すると、コピー先のデータの管理情報を、新たなアクセス先としてデータアクセスエラー情報を送信した他の装置に伝達する。図 1 の例では、アクセスノード 6 が設けられており、アクセスノード 6 からデータアクセスエラー情報が送信されている。そのため、コピー指示手段 1 e は、アクセスノード 6 に対してコピー先のデータの管理情報を送信する。

【 0 0 2 3 】

各ディスクノード 2 ~ 5 は、ライトエラー発生時に、データ割当制御装置 1 へライトエラー情報を送信する。また図 1 の例では、アクセスノード 6 が各ディスクノード 2 ~ 5 で管理されているデータにアクセスする。

【 0 0 2 4 】

ここで、ディスクノード 5 が故障すると、所定のタイミングで二重化復旧手段 1 a により二重化復旧処理が開始される。例えば、二重化復旧手段 1 a は、ディスクノード 2 ~ 4 で管理されているストレージ装置 2 a , 3 a , 4 a 内のデータのうち、冗長データが存在しない二重化欠損データを調べる。二重化復旧手段 1 a は、二重化欠損データの管理情報

10

20

30

40

50

を管理情報記憶手段 1 b に格納する。そして、二重化復旧手段 1 a は、二重化欠損データを順次選択し、二重化欠損データを管理しているディスクノードに対し、二重化欠損データが格納されているストレージ装置とは別のストレージ装置への二重化欠損データのコピーを指示する。二重化欠損データのコピー指示を受けたディスクノードでは、二重化欠損データのコピーが行われる。例えば、ディスクノード 3 が、ストレージ装置 3 a 内の二重化欠損データを、ディスクノード 4 に接続されたストレージ装置 4 a にコピーする。また、ディスクノード 4 が、ストレージ装置 4 a 内の二重化欠損データを、ディスクノード 3 に接続されたストレージ装置 3 a にコピーする。

【0025】

このような二重化復旧処理中に、ディスクノード 2 に接続されたストレージ装置 2 a でライトエラーが発生したものとす。例えば、ストレージ装置 2 a 内での定期的なデータの書き込みチェックが行われ、書き込みができず、かつストレージ装置 2 a 内でのリカバリができないような場合に、ストレージ装置 2 a からディスクノード 2 にライトエラーが出力される。ディスクノード 2 は、ストレージ装置 2 a でのライトエラーの発生を検出すると、データ割当制御装置 1 にライトエラー情報を送信する。ライトエラー情報には、例えば、システム内でストレージ装置 2 a を一意に識別する識別情報が含まれる。

10

【0026】

ディスクノード 2 が出力したライトエラー情報は、データ割当制御装置 1 のエラー情報受信手段 1 c で受信される。エラー情報受信手段 1 c は、二重化復旧処理中であることを確認し、ライトエラーが発生したストレージ装置の識別情報をエラー記憶手段 1 d に格納する。その結果、エラー記憶手段 1 d には、ストレージ装置 2 a にライトエラーが発生したことを示す情報が記憶される。

20

【0027】

ストレージ装置 2 a でライトエラーが発生すると、二重化復旧手段 1 a は、その後の二重化欠損データのコピー先からストレージ装置 2 a を除外する。

また、ストレージ装置 2 a でライトエラーが発生後、アクセスノード 6 がストレージ装置 2 a 内のデータへのアクセス要求をディスクノード 2 に送信すると、ディスクノード 2 からアクセスノード 6 にエラーが応答される。するとアクセスノード 6 は、データ割当制御装置 1 に、ストレージ装置 2 a へのアクセスでエラーが発生したことを示すデータアクセスエラー情報を送信する。データアクセスエラー情報には、エラー発生時のアクセス先であったアクセス対象データの識別情報が示される。

30

【0028】

データ割当制御装置 1 では、コピー指示手段 1 e がデータアクセスエラー情報を受信する。コピー指示手段 1 e は、データアクセスエラー情報に基づいて、エラー発生時のアクセス対象のデータを認識し、コピー対象データとする。そして、コピー指示手段 1 e は、コピー対象データの冗長データが、ライトエラーが発生していないストレージ装置内になければ、ディスクノード 2 に対して、コピー対象データの他のストレージ装置へのコピーを指示する。この際、コピー対象データのコピー先は、ライトエラーが発生していないいずれかのストレージ装置である。図 1 の例では、ストレージ装置 2 a 内のコピー対象データの、ストレージ装置 3 a へのコピーが指示されている。ディスクノード 2 は、コピー対象データのコピー指示に従って、コピー対象データをコピーする。

40

【0029】

コピー指示手段 1 e は、コピー対象データのコピーが完了すると、コピー先のデータの場所を示す管理情報を、新たなアクセス先としてアクセスノード 6 に通知する。アクセスノード 6 は、コピー指示手段 1 e からの管理情報の通知を受けて、コピー先となるストレージ装置 3 a を管理するディスクノード 3 へ再アクセスを行うことができる。

【0030】

このように、二重化復旧処理中にストレージ装置 2 a でライトエラーが発生しても、ストレージ装置 2 a のシステムからの切り離しは行われぬ。すなわち、ストレージ装置 2 a にライトエラーが発生しても、データを読み出せる可能性があるため、ストレージ装置

50

2 a の即時切り離しは行われぬ。そして、ライトエラーが発生したストレージ装置内のコピー対象データのうち、ライトエラーが発生していないストレージ装置に冗長データが存在しないコピー対象データについては、他のストレージ装置にコピーされる。コピー対象データのコピーが完了すると、そのコピー対象データについてはデータの喪失を免れたこととなる。すなわち、二重化復旧処理中にライトエラーとなったストレージ装置 2 a 内のデータの喪失が抑制されている。

【 0 0 3 1 】

またライトエラーが発生したストレージ装置 2 a が切り離されていないため、ストレージ装置 2 a 内に二重化欠損データがあれば、二重化復旧処理によって、ストレージ装置 2 a 内の二重化欠損データが他のストレージ装置にコピーされる。ストレージ装置 2 a から他のストレージ装置への二重化欠損データのコピーが完了すれば、その二重化欠損データについてはデータ喪失を免れたこととなる。すなわち、二重化復旧処理中にライトエラーとなったストレージ装置内のデータの喪失が抑制されている。

10

【 0 0 3 2 】

さらにストレージ装置 5 a が故障したことによる二重化復旧処理が完了した場合、すべてのデータに冗長データが存在することとなる。データとそのデータの冗長データとは、異なるストレージ装置に格納される。すると、二重化復旧処理が完了した場合、ライトエラーが発生したストレージ装置 2 a 内のすべてのデータは、同一内容のデータが他のストレージ装置にも格納されていることとなる。従って、ストレージ装置 5 a が故障したことによる二重化復旧処理が完了後は、ストレージ装置 2 a を切り離しても、データを喪失せずに済む。そこで、ストレージ装置 2 a でライトエラー発生時に管理情報記憶手段 1 b 内に登録されていたすべての二重化欠損データのコピー完了後は、ライトエラーが発生したストレージ装置 2 a が切り離される。すると、ストレージ装置 2 a 内のデータと同一内容の他のストレージ装置 3 a , 4 a 内のデータが二重化欠損データとなり、コピーされる。その結果、正常に動作するストレージ装置 3 a , 4 a 内ですべてのデータが二重化され、システムの信頼性を向上させることができる。

20

【 0 0 3 3 】

〔 第 2 の実施の形態 〕

第 2 の実施の形態は、論理ディスクを用いてデータアクセスを管理するマルチノードシステムの例である。

30

【 0 0 3 4 】

図 2 は、第 2 の形態のマルチノードストレージシステム構成の一例を示す図である。本実施の形態では、ネットワーク 1 0 を介して、複数のディスクノード 1 0 0 , 2 0 0 , 3 0 0 , 4 0 0 、制御ノード 5 0 0 、およびアクセスノード 6 0 0 が接続されている。ディスクノード 1 0 0 , 2 0 0 , 3 0 0 , 4 0 0 それぞれには、ストレージ装置 1 1 0 , 2 1 0 , 3 1 0 , 4 1 0 が接続されている。

【 0 0 3 5 】

ストレージ装置 1 1 0 には、複数のハードディスク装置 (H D D) 1 1 1 , 1 1 2 , 1 1 3 , 1 1 4 が実装されている。ストレージ装置 2 1 0 には、複数の H D D 2 1 1 , 2 1 2 , 2 1 3 , 2 1 4 が実装されている。ストレージ装置 3 1 0 には、複数の H D D 3 1 1 , 3 1 2 , 3 1 3 , 3 1 4 が実装されている。ストレージ装置 4 1 0 には、複数の H D D 4 1 1 , 4 1 2 , 4 1 3 , 4 1 4 が実装されている。各ストレージ装置 1 1 0 , 2 1 0 , 3 1 0 , 4 1 0 は、内蔵する H D D を用いた R A I D システムである。本実施の形態では、各ストレージ装置 1 1 0 , 2 1 0 , 3 1 0 , 4 1 0 の R A I D 5 のディスク管理サービスを提供する。

40

【 0 0 3 6 】

ディスクノード 1 0 0 , 2 0 0 , 3 0 0 , 4 0 0 は、接続されたストレージ装置 1 1 0 , 2 1 0 , 3 1 0 , 4 1 0 に格納されたデータを管理し、管理しているデータをネットワーク 1 0 経由で端末装置 2 1 , 2 2 , 2 3 に提供する。また、ディスクノード 1 0 0 , 2 0 0 , 3 0 0 , 4 0 0 は、冗長性を有するデータを管理している。すなわち、同一のデー

50

タが、少なくとも2つのディスクノードで管理されている。

【0037】

制御ノード500は、ディスクノード100, 200, 300, 400を管理する。例えば、制御ノード500は、ディスクノード100, 200, 300, 400から新たなストレージ装置の接続通知を受け取ると、論理ディスクを介して接続されたストレージ装置にアクセスできるようにする。

【0038】

アクセスノード600には、ネットワーク20を介して複数の端末装置21, 22, 23が接続されている。また、アクセスノード600には、論理ディスクが定義されている。そして、アクセスノード600は、端末装置21, 22, 23からの論理ディスクのデータのアクセス要求に応答して、ディスクノード100, 200, 300, 400内の対応するデータへアクセスする。

10

【0039】

なお、図2に示すように第2の実施の形態では、各ディスクノード100, 200, 300, 400に、1台ずつのストレージ装置110, 210, 310, 410が接続されている。そのため、ストレージ装置でのライトエラーなどの障害の発生は、そのストレージ装置が接続されたディスクノードの障害とみなすことができる。そして、障害が発生したストレージ装置を切り離す場合には、ディスクノードの切り離しが行われることとなる。

【0040】

図3は、第2の形態に用いる制御ノードのハードウェアの一構成例を示す図である。制御ノード500は、CPU (Central Processing Unit) 501によって装置全体が制御されている。CPU 501には、バス508を介してRAM (Random Access Memory) 502と複数の周辺機器が接続されている。

20

【0041】

RAM 502は、制御ノード500の主記憶装置として使用される。RAM 502には、CPU 501に実行させるOS (Operating System) のプログラムやアプリケーションプログラムの少なくとも一部が一時的に格納される。また、RAM 502には、CPU 501による処理に必要な各種データが格納される。

【0042】

バス508に接続されている周辺機器としては、ハードディスクドライブ (HDD: Hard Disk Drive) 503、グラフィック処理装置504、入力インタフェース505、光学ドライブ装置506、および通信インタフェース507がある。

30

【0043】

HDD 503は、内蔵したディスクに対して、磁氣的にデータの書き込みおよび読み出しを行う。HDD 503は、制御ノード500の二次記憶装置として使用される。HDD 503には、OSのプログラム、アプリケーションプログラム、および各種データが格納される。なお、二次記憶装置としては、フラッシュメモリなどの半導体記憶装置を使用することもできる。

【0044】

グラフィック処理装置504には、モニタ11が接続されている。グラフィック処理装置504は、CPU 501からの命令に従って、画像をモニタ11の画面に表示させる。モニタ11としては、CRT (Cathode Ray Tube) を用いた表示装置や液晶表示装置などがある。

40

【0045】

入力インタフェース505には、キーボード12とマウス13とが接続されている。入力インタフェース505は、キーボード12やマウス13から送られてくる信号をCPU 501に送信する。なお、マウス13は、ポインティングデバイスの一例であり、他のポインティングデバイスを使用することもできる。他のポインティングデバイスとしては、タッチパネル、タブレット、タッチパッド、トラックボールなどがある。

50

【 0 0 4 6 】

光学ドライブ装置 5 0 6 は、レーザ光などを利用して、光ディスク 1 4 に記録されたデータの読み取りを行う。光ディスク 1 4 は、光の反射によって読み取り可能なようにデータが記録された可搬型の記録媒体である。光ディスク 1 4 には、D V D (Digital Versatile Disc)、D V D - R A M、C D - R O M (Compact Disc Read Only Memory)、C D - R (Recordable) / R W (ReWritable) などがある。

【 0 0 4 7 】

通信インタフェース 5 0 7 は、ネットワーク 1 0 に接続されている。通信インタフェース 1 0 7 は、ネットワーク 1 0 を介して、他の装置との間でデータの送受信を行う。

以上のようなハードウェア構成によって、本実施の形態の処理機能を実現することができる。なお、図 3 では制御ノード 5 0 0 のハードウェア構成を示したが、ディスクノード 1 0 0, 2 0 0, 3 0 0, 4 0 0、およびアクセスノード 6 0 0 も同様のハードウェア構成で実現することができる。ただし、ディスクノード 1 0 0, 2 0 0, 3 0 0, 4 0 0 は、図 3 に示した機能に加え、ストレージ装置 1 1 0, 2 1 0, 3 1 0, 4 1 0 を外部接続するためのインタフェースを有している。

10

【 0 0 4 8 】

次に、マルチノードストレージシステムにおいて定義される論理ディスクのデータ構造について説明する。

図 4 は、論理ディスクのデータ構造の一例を示す図である。第 2 の実施の形態では、論理ディスク 6 0 には論理ディスク識別子「L V O L - X」が付与されている。ネットワーク経路で接続された 4 台のディスクノード 1 0 0, 2 0 0, 3 0 0, 4 0 0 には、個々のノードの識別用にそれぞれ「D P 1」、「D P 2」、「D P 3」、「D P 4」というディスクノード ID が付与されている。そして、各ディスクノード 1 0 0, 2 0 0, 3 0 0, 4 0 0 に接続されているストレージ装置 1 1 0, 2 1 0, 3 1 0, 4 1 0 は、ディスクノード ID と、各ディスクノード内のディスク ID との組によってネットワーク 1 0 で一意に識別される。

20

【 0 0 4 9 】

各ディスクノード 1 0 0, 2 0 0, 3 0 0, 4 0 0 が有するストレージ装置 1 1 0, 2 1 0, 3 1 0, 4 1 0 それぞれにおいて R A I D 5 のストレージシステムが構成されている。各ストレージ装置 1 1 0, 2 1 0, 3 1 0, 4 1 0 で提供される記憶機能は、複数のスライス 1 1 5 a ~ 1 1 5 c, 2 1 5 a ~ 2 1 5 c, 3 1 5 a ~ 3 1 5 c, 4 1 5 a ~ 4 1 5 c に分割されて管理されている。

30

【 0 0 5 0 】

論理ディスク 6 0 は、セグメント 6 1 ~ 6 4 という単位で構成される。セグメント 6 1 ~ 6 4 の記憶容量は、ストレージ装置 1 1 0, 2 1 0, 3 1 0, 4 1 0 における管理単位であるスライスの記憶容量と同じである。例えば、スライスの記憶容量が 1 ギガバイトとするとセグメントの記憶容量も 1 ギガバイトである。論理ディスク 6 0 の記憶容量はセグメント 1 つ当たりの記憶容量の整数倍である。セグメント 6 1 ~ 6 4 は、それぞれプライマリスライス 6 1 a, 6 2 a, 6 3 a, 6 4 a とセカンダリスライス 6 1 b, 6 2 b, 6 3 b, 6 4 b との組 (スライスペア) で構成される。

40

【 0 0 5 1 】

同一セグメントに属する 2 つのスライスは別々のディスクノードに属する。個々のスライスを管理する領域には論理ディスク識別子やセグメント情報や同じセグメントを構成するスライス情報の他にフラグがあり、そのフラグにはプライマリあるいはセカンダリなどを表す値が格納される。

【 0 0 5 2 】

図 4 の例では、論理ディスク 6 0 内のスライスの識別子を、「P」または「S」のアルファベットと数字との組合せで示している。「P」はプライマリスライスであることを示している。「S」はセカンダリスライスであることを示している。アルファベットに続く数字は、何番目のセグメントに属するのかを表している。例えば、1 番目のセグメント 6

50

1のプライマリスライスが「P1」で示され、セカンダリスライスが「S1」で示される。

【0053】

図5は、第2の実施の形態に係るマルチノードストレージシステムの各装置の機能を示すブロック図である。まずアクセスノード600の機能について説明する。アクセスノード600は、メタデータ照会部610、アクセス用メタデータ記憶部620、およびアクセス要求部630を有している。

【0054】

メタデータ照会部610は、論理ディスク60を定義するメタデータを、制御ノード500から取得する。具体的には、メタデータ照会部610は、アクセスノード600の起動時に制御ノード500に対して全メタデータの照会要求を送信する。すると、制御ノード500から論理ディスク60に関する全メタデータが送られてくる。また、メタデータ照会部610は、アクセス要求部630による任意のスライスへのデータアクセスがエラーとなったとき、アクセス対象のスライスが割り当てられたセグメントに関するメタデータの照会要求を制御ノード500に送信する。この場合、メタデータの照会要求が、データアクセスエラー情報の通知も兼ねることとなる。すると、制御ノード500から、該当セグメントの最新のメタデータが送られてくる。なお、メタデータ照会部610は、制御ノード500からメタデータを取得すると、そのメタデータをアクセス用メタデータ記憶部620に格納する。

【0055】

アクセス用メタデータ記憶部620は、論理ディスク60を定義するメタデータを記憶する。例えば、アクセスノード600のRAMの一部がアクセス用メタデータ記憶部620として使用される。なお、第2の実施の形態では、アクセスノード600は常にプライマリスライス（シングルプライマリスライスも含む）にアクセスを行う。そのため、アクセス用メタデータ記憶部620には、論理ディスク60のメタデータのうち、少なくともプライマリスライスに関するメタデータが格納されていればよい。

【0056】

アクセス要求部630は、端末装置21, 22, 23からの論理ディスク60上でのデータのアクセス要求に回答して、ストレージ装置110, 210, 310, 410のデータのアクセス要求（リード要求またはライト要求）を送信する。アクセス要求は、アクセス先のストレージ装置が接続されたディスクノードに対して送信される。具体的には、アクセス要求部630は、論理ディスク60のアドレスを指定したアクセス要求を受け取ると、まず、アクセス用メタデータ記憶部620を参照し、アクセス対象のデータが属するセグメントを判断する。次に、アクセス要求部630は、該当するセグメントにプライマリスライスとして割り当てられたスライスを判断する。そして、アクセス要求部630は、該当するスライスを管理するディスクノードに対して、そのスライス内のデータのアクセス要求を送信する。ディスクノードからアクセス結果が応答されると、アクセス要求部630は、端末装置21, 22, 23にアクセス結果を送信する。

【0057】

なお、アクセス要求部630は、アクセス先のディスクノードからエラーが返された場合、エラーが発生したセグメントをメタデータ照会部610に通知する。その後、アクセス要求部630は、データアクセスのリトライを行う。リトライでは、アクセス用メタデータ記憶部620を参照して、新たにプライマリスライスとした割り当てられたスライスを判断する処理から再度実行される。すなわち、前回のアクセス要求後にアクセス用メタデータ記憶部620内のメタデータが更新されていれば、更新後のメタデータに基づいてリトライ時のアクセス先となるディスクノードが判断される。

【0058】

次に、ディスクノード100の機能について説明する。ディスクノードは、アクセス処理部120、エラー通知部130、エラー記憶部140、メタデータ記憶部150、およびメタデータ管理部160を有する。

【 0 0 5 9 】

アクセス処理部 1 2 0 は、アクセスノード 6 0 0 からのアクセス要求に回答して、ストレージ装置 1 1 0 に対するデータアクセスを行う。具体的には、アクセス処理部 1 2 0 は、アクセスノード 6 0 0 からアクセス要求を受け取ると、メタデータ記憶部 1 5 0 を参照し、アクセス対象となるセグメントに対して、割り当てられているストレージ装置 1 1 0 内のスライスを判断する。

【 0 0 6 0 】

次に、アクセス処理部 1 2 0 は、アクセス要求で指定されているスライス内のデータに対してアクセスする。例えば、データリードのアクセス要求であれば、アクセス処理部 1 2 0 は、該当するデータをストレージ装置 1 1 0 から読み出す。また、データライトのアクセス要求であれば、アクセス処理部 1 2 0 は、ストレージ装置 1 1 0 内の該当する記憶領域にアクセス要求に含まれるデータを書き込む。そして、アクセス処理部 1 2 0 は、アクセス結果をアクセスノード 6 0 0 に送信する。データリードのアクセス要求の場合、ストレージ装置 1 1 0 から読み出したデータがアクセス結果に含まれる。

10

【 0 0 6 1 】

また、データライトのアクセス要求の場合、アクセス処理部 1 2 0 は、他のディスクノード 2 0 0 , 3 0 0 , 4 0 0 との間でミラーリング処理を実行する。ミラーリング処理では、例えばアクセス処理部 1 2 0 は、メタデータ記憶部 1 5 0 を参照し、ライトアクセスの対象となるストレージ装置 1 1 0 内のスライス（プライミスライス）と組となる、他のディスクノードのスライス（セカンダリスライス）を判断する。そして、アクセス処理部 1 2 0 は、セカンダリスライスを管理するディスクノードに対して、プライミスライスに書き込まれたデータと同じデータの書き込み要求を送信する。すると、データの書き込み要求を受信したディスクノードにおいてセカンダリスライスへのデータの書き込みが行われる。

20

【 0 0 6 2 】

ミラーリング処理では、他のディスクノードのプライミスライスに書き込まれたデータのセカンダリスライスへの書き込み要求が、他のディスクノードからディスクノード 1 0 0 に送信されることがある。ディスクノード 1 0 0 に送られたセカンダリスライスへのデータの書き込み要求は、アクセス処理部 1 2 0 で受け取られる。アクセス処理部 1 2 0 は、セカンダリスライスへのデータの書き込み要求を受け取ると、ストレージ装置 1 1 0 内の該当スライスにデータの書き込みを行う。

30

【 0 0 6 3 】

ストレージ装置 1 1 0 へのデータ書き込みの際には、ライトエラーが発生する可能性がある。ライトエラーが発生すると、アクセス処理部 1 2 0 は、エラー記憶部 1 4 0 にライトエラーの発生を示すエラー情報を格納すると共に、エラー通知部 1 3 0 に対してライトエラーの発生を通知する。さらにアクセス処理部 1 2 0 は、アクセスノード 6 0 0 からのアクセス要求に応じて実行したストレージ装置 1 1 0 へのライトアクセスがエラーとなった場合、アクセスノード 6 0 0 に対してエラーを返す（エラーメッセージを送信する）。

【 0 0 6 4 】

エラー通知部 1 3 0 は、アクセス処理部 1 2 0 からライトエラー発生の通知を受け取ると、制御ノード 5 0 0 へエラーの発生を示すエラーメッセージを送信する。例えばエラー通知部 1 3 0 は、ディスクノード 1 0 0 の識別子（ディスク ID）を含むエラーメッセージを制御ノード 5 0 0 に送信する。

40

【 0 0 6 5 】

エラー記憶部 1 4 0 は、ストレージ装置 1 1 0 のエラー情報を記憶する。エラー情報は、ライトエラーの有無を示す情報である。例えばエラー情報として、「ライトエラーあり」または「ライトエラーなし」の 2 つの状態を有するフラグを用いることができる。エラー記憶部 1 4 0 としては、例えば、ディスクノード 1 0 0 内の R A M の記憶領域の一部が用いられる。

【 0 0 6 6 】

50

なお、図5の例ではディスクノード100に対して1台のストレージ装置110が接続されている。そのため、エラー情報に、ストレージ装置の識別情報を含めなくても、エラー情報に対応するストレージ装置を一意に決定できる。ディスクノード100に複数のストレージ装置が接続されていた場合、ストレージ装置ごとのエラー情報がエラー記憶部140に格納される。その場合、各エラー情報には、ストレージ装置をディスクノード100内で一意に識別する識別子が含まれる。

【0067】

メタデータ記憶部150は、ディスクノード100が管理しているスライスのメタデータの記憶機能である。例えば、ディスクノード100のRAM内の一部の記憶領域がメタデータ記憶部150として使用される。

10

【0068】

メタデータ管理部160は、ストレージ装置110内の各スライスのメタデータを管理する。具体的には、メタデータ管理部160は、ディスクノード100起動時に、ストレージ装置110から各スライスのメタデータを読み出し、メタデータ記憶部150に格納する。また、メタデータ管理部160は、制御ノード500からメタデータの収集要求があれば、メタデータ記憶部150に格納されているメタデータを制御ノード500に送信する。さらにメタデータ管理部160は、制御ノード500からメタデータの変更要求を受け取ると、その変更要求で指定されたメタデータの内容を変更する。この際、メタデータ管理部160は、メタデータ記憶部150内のメタデータと、ストレージ装置110内のメタデータとを変更する。

20

【0069】

なお、図5には複数のディスクノード100, 200, 300, 400のうち、代表的にディスクノード100の機能を示したが、他のディスクノード200, 300, 400も同様の機能を有している。

【0070】

次に制御ノード500の機能について説明する。制御ノード500は、エラー受信部510、エラー記憶部520、論理ディスク管理部530および論理ディスクメタデータ記憶部540を有している。

【0071】

エラー受信部510は、ディスクノードからのライトエラーを示すエラーメッセージを受信する。エラー受信部510は、エラーメッセージを受信すると、二重化復旧のリカバリ処理中でなければ、エラーメッセージを送信したディスクノードでライトエラーが発生したことを示すエラー情報を、エラー記憶部520に格納する。また、エラー受信部510は、エラーメッセージを受信したときリカバリ処理中であれば、論理ディスク管理部530内のリカバリ処理部533にリカバリ処理の開始を指示する。

30

【0072】

エラー記憶部520は、ディスクノードでのライトエラーの有無を示すエラー情報を記憶する。例えば、RAM502またはHDD503の記憶領域の一部が、エラー記憶部520として使用される。

【0073】

論理ディスク管理部530は、マルチノードストレージシステム内の各ディスクノード100, 200, 300, 400のメタデータを管理する。論理ディスク管理部530は、メタデータ収集部531、スライス割当処理部532、およびリカバリ処理部533を有する。

40

【0074】

メタデータ収集部531は、所定のタイミングで、各ディスクノード100, 200, 300, 400からメタデータを収集する。例えば、メタデータ収集部531は、制御ノード500が起動した場合にメタデータを収集する。メタデータを収集するタイミングになると、メタデータ収集部531は、各ディスクノード100, 200, 300, 400にメタデータの収集要求を送信する。すると各ディスクノード100, 200, 300,

50

400からメタデータが応答される。メタデータ収集部531は、収集したメタデータを論理ディスクメタデータ記憶部540に格納する。

【0075】

スライス割当処理部532は、論理ディスク60の各セグメントへのスライスの割り当て処理を行う。例えば、スライス割当処理部532は、新たな論理ディスクが定義された場合などに、セグメントへのスライスの割り当て処理を行う。スライスの割り当て処理では、例えば、各セグメントに対して、異なるディスクノードで管理されている2つのスライスが割り当てられる。セグメントに割り当てるスライスとしては、他のセグメントに割り当てられておらず、かつ異常スライスではないスライス（フリースライス）が選択される。スライス割当処理部532は、新たにセグメントに割り当てられたスライスに関する論理ディスクメタデータ記憶部540内のメタデータを更新する。また、スライス割当処理部532は、更新されたメタデータを含むメタデータ変更要求を、そのメタデータに対応するスライスを管理しているディスクノードに送信する。

10

【0076】

さらにスライス割当処理部532は、アクセスノード600からのメタデータ照会要求を受け取ると、照会要求で指定されたセグメントに割り当てられたスライスのメタデータを、アクセスノード600に送信する。なお照会要求で指定されたセグメントに割り当てられたスライスが、ライトエラーが発生したディスクノードで管理されている場合がある。その場合、スライス割当処理部532は、指定されたセグメントに他のスライスを割り当て、その後、新たに割り当てられたスライスのメタデータをアクセスノード600に送信する。

20

【0077】

リカバリ処理部533は、所定のタイミングで、二重化復旧のリカバリ処理を実行する。例えば、リカバリ処理部533は、エラー受信部510からリカバリ指示を受け取ったときにリカバリ処理を開始する。また、リカバリ処理部533は、システム管理者の操作入力によるリカバリ処理開始指示があった場合に、リカバリ処理を開始することもできる。

【0078】

論理ディスクメタデータ記憶部540は、論理ディスク60を構成するセグメントへのスライスの割り当て関係を示すメタデータを記憶する記憶機能である。例えば、RAM502またはHDD503の記憶領域の一部が論理ディスクメタデータ記憶部540として使用される。

30

【0079】

次に、各ノードで管理されているデータの構造について説明する。

図6は、ストレージ装置のデータ構造の一例を示す図である。ストレージ装置110には、スライス115a, 115b, 115c, ...とは別に複数のメタデータ117a, 117b, 117c, ...が格納されている。

【0080】

ストレージ装置110に格納されたメタデータ117a, 117b, 117c, ...は、ディスクノード100の起動時にメタデータ管理部160によって読み出され、メタデータ記憶部150に格納される。

40

【0081】

なお、図6の例では、各スライスに隣接した記憶領域に、それぞれのスライスに関するメタデータが格納されているが、メタデータの格納領域は図6に示した場所には限らない。例えば、1つのメタデータ記憶領域内に、すべてのスライスのメタデータを記憶させることもできる。

【0082】

また、1つのスライスのデータを不連続の記憶領域に格納することもできる。その場合、例えば、各スライスが、所定のデータ長の複数のデータユニットに分割される。そして、各スライスのメタデータには、各データユニットに対応するデータの格納先となる記憶

50

領域を特定する情報が設定される。データの格納先となる記憶領域を特定する情報は、例えば、ストレージ装置の記憶領域の先頭からのオフセット値である。

【0083】

図7は、メタデータ記憶部のデータ構造例を示す図である。メタデータ記憶部150には、メタデータテーブル151が格納されている。メタデータテーブル151には、ディスクノードID、ディスクID、スライスID、状態、論理ディスクID、セグメントID、論理ディスクアドレス、ペアのディスクノードID、ペアのディスクID、およびペアのスライスIDの欄が設けられている。メタデータテーブル151内の横方向に並べられた情報同士が互いに関連付けられ、メタデータを示す1つのレコードを構成している。

【0084】

ディスクノードIDの欄は、ストレージ装置110を管理しているディスクノード100の識別情報(ディスクノードID)が設定される。

ディスクIDの欄には、ストレージ装置110の識別情報(ディスクID)が設定される。ディスクノード100に複数のストレージ装置が接続される場合、各ストレージ装置に異なるディスクIDが設定される。

【0085】

スライスIDの欄には、メタデータに対応するスライスのストレージ装置110内での識別情報(スライスID)が設定される。

状態の欄には、スライスの状態を示す状態フラグが設定される。スライスが論理ディスク60のセグメントに割り当てられていない場合、状態フラグ「F」が設定される。論理ディスク60のセグメントのプライマリストレージに割り当てられている場合、状態フラグ「P」または「SP」が設定される。種別フラグ「P」は、ペアとなるセカンダリセグメント(ミラースライス)が存在するプライマリスライスを示している。種別フラグ「SP」は、ミラースライスが障害などによって存在しなくなったプライマリスライス(シングルプライマリスライス)を示している。なお、種別フラグ「SP」は、そのスライスが割り当てられているセグメントの二重化状態が損なわれていることを示す欠損フラグでもある。論理ディスク60のセグメントのセカンダリストレージに割り当てられている場合、状態フラグ「S」が設定される。異常スライスと判定された場合、異常であることを示す状態フラグ「B」が設定される。異常スライスは、セグメントへの割り当て候補から除外される。

【0086】

なお、シングルプライマリスライスが割り当てられているセグメントは冗長性が確保されていない。以下、シングルプライマリスライスのみが割り当てられているセグメントを、欠損セグメントと呼ぶ。欠損セグメントに対しては、セカンダリスライスとして割り当てべきスライスが予約される。予約されたスライスは、リザーブスライスと呼ばれる。欠損セグメントのリザーブスライスに対してプライマリスライスからデータがコピーされる。これにより、欠損セグメントのデータの二重化状態が復旧する。このような二重化復旧処理中のリザーブスライスのメタデータに対しては、データのコピーが完了するまで種別フラグ「R」が設定される。リザーブスライスに対するプライマリスライスからのデータのコピーが完了すると、リザーブスライスはセカンダリスライスに変更される。

【0087】

論理ディスクIDの欄には、スライスに対応するセグメントが属する論理ディスク60を識別するための識別情報(論理ディスクID)が設定される。

セグメントIDの欄には、スライスが割り当てられたセグメントの識別情報(セグメントID)が設定される。

【0088】

論理ディスクアドレスの欄には、スライスが割り当てられているセグメントの先頭を示す論理ディスク60内でのアドレスが設定される。

ペアのディスクノードIDの欄には、ペアのスライス(同じセグメントに属する別のスライス)を有するストレージ装置を管理するディスクノードの識別情報(ディスクノード

10

20

30

40

50

ID)が設定される。

【0089】

ペアのディスクIDの欄には、ペアのスライスを有するストレージ装置の識別情報(ディスクID)が設定される。

ペアのスライスIDの欄には、ペアのスライスを、そのスライスが属するストレージ装置内で識別するための識別情報(スライスID)が設定される。

【0090】

図7にはディスクノード100のメタデータ記憶部150の内容を示しているが、他のディスクノード200,300,400も同様のメタデータ記憶部を有している。そして、各ディスクノード100,200,300,400のメタデータ記憶部に格納されたメタデータは、制御ノード500からの要求に応じて制御ノード500に送信される。制御ノード500では、ディスクノード100,200,300,400から収集したメタデータは、メタデータ収集部531により論理ディスクメタデータ記憶部540に格納される。

10

【0091】

図8は、論理ディスクメタデータ記憶部のデータ構造の一例を示す図である。論理ディスクメタデータ記憶部540には、論理ディスクメタデータテーブル541が格納されている。論理ディスクメタデータテーブル541には、ディスクノードID、ディスクID、スライスID、状態、論理ディスクID、セグメントID、論理ディスクアドレス、ペアのディスクノードID、ペアのディスクID、およびペアのスライスIDの欄が設けられている。論理ディスクメタデータテーブル541内の横方向に並べられた情報が互いに関連付けられ、メタデータを示す1つのレコードを構成している。論理ディスクメタデータテーブル541の各欄に設定される情報は、メタデータテーブル151の同名の欄と同種の情報である。

20

【0092】

論理ディスクメタデータテーブル541に格納されたメタデータは、アクセスノード600からの照会要求に回答して、アクセスノード600に送信される。アクセスノード600は、取得したメタデータを記憶する。具体的には、アクセスノード600のアクセス用メタデータ記憶部620にメタデータが格納される。

【0093】

アクセス用メタデータ記憶部620のデータ構造は、論理ディスクメタデータ記憶部540と同様である。なお、本実施の形態では、アクセスノード600は常にプライマリスライスまたはシングルプライマリスライスにアクセスする。そのため、アクセス用メタデータ記憶部620には、少なくともプライマリスライスおよびプライマリスライスに関するメタデータ(状態フラグが「P」または「SP」のメタデータ)が格納されていればよい。また、アクセス用メタデータ記憶部620には、各メタデータにおけるペアのディスクノードID、およびペアのスライスIDの各欄のデータは無くてもよい。

30

【0094】

図9は、エラー記憶部のデータ構造の一例を示す図である。エラー記憶部520には、各ディスクノード100,200,300,400のライトエラーの有無を示す情報が格納される。図9の例では、エラー記憶部520内にエラー管理テーブル521が設けられている。

40

【0095】

エラー管理テーブル521には、ディスクノードID、ディスクID、およびライトエラーの欄が設けられている。ディスクノードIDの欄には、各ディスクノードの識別子が設定される。ディスクIDの欄には、各ディスクノードに接続されているストレージ装置のディスクIDが設定される。ライトエラーの欄には、対応するディスクノードのストレージ装置でライトエラーが発生したか否かを示す情報が設定される。例えば、ライトエラーの欄には、ライトエラーに関して「あり」を示す状態と、「なし」を示す状態とを有するフラグが設定される。

50

【 0 0 9 6 】

以上のような構成のマルチノードストレージシステムにおいて、いずれかのディスクノードに故障が発生すると、リカバリ処理が実行される。ディスクノードの故障は、例えば、制御ノードで検知される。制御ノード500は、ディスクノードからのハートビートの途絶や、ディスクノードからのライトエラーの発生を示すエラーメッセージの受信によって、故障を検知することができる。

【 0 0 9 7 】

以下、ディスクノード400で最初にライトエラーが発生したことに起因するリカバリ処理中に、他のディスクノードでライトエラーが発生した場合を想定し、第2の実施の形態の処理を詳細に説明する。

10

【 0 0 9 8 】

図10は、ディスクノードでのライトエラー検出時に装置間で送受信される情報の例を示す図である。図10では、各装置の機能のうち、他の装置と通信を行う機能が示されている。

【 0 0 9 9 】

制御ノード500の論理ディスク管理部530から各ディスクノード100, 200, 300のメタデータ管理部160, 260, 360それぞれへは、メタデータ変更要求、スライスコピー要求などが送信される。ディスクノード100のエラー通知部130から制御ノード500のエラー受信部510へは、エラーメッセージが送信される。ディスクノード100, 200, 300のアクセス処理部120, 220, 320間では、ミラーリングによる書き込みデータや、スライスコピー時の書き込みデータなどが互いに送信される。

20

【 0 1 0 0 】

アクセスノード600のアクセス要求部630からディスクノード100のアクセス処理部120へは、リードまたはライトのアクセス要求が送信される。なお、図10には示していないが、アクセス要求部630からディスクノード200, 300のアクセス処理部220, 320へも、リードまたはライトのアクセス要求が送信される。アクセスノード600のメタデータ照会部610から制御ノード500の論理ディスク管理部530へは、メタデータ照会要求が送信される。

【 0 1 0 1 】

このような装置間の通信による装置の協働動作により、故障したディスクノードの切り離し処理、欠損セグメントの二重化を回復するリカバリ処理、およびリカバリ処理中のライトエラーに起因するエラー処理が進行する。そこで、まず故障したディスクノードの切り離し処理について説明する。

30

【 0 1 0 2 】

第2の実施の形態では、リカバリ処理の実行中を除き、故障したディスクノードはマルチノードストレージシステムから切り離される。故障したディスクノードの切り離しとは、アクセスノード600からのアクセス先、またはミラーリングによるコピー先から、故障したディスクノードを排除することである。

【 0 1 0 3 】

図11は、故障したディスクノードの切り離し処理の一例を示す図である。以下、図11に示す処理をステップ番号に沿って説明する。なお、図11の例には、リカバリ処理が実行されていない状況下でディスクノード400が故障した場合の処理手順が示されている。

40

【 0 1 0 4 】

[ステップS11] ディスクノード400は、ストレージ装置410へのライトアクセス時にライトエラーを検出すると、ライトエラーの発生を示すエラーメッセージを制御ノード500に送信する。

【 0 1 0 5 】

[ステップS12] 制御ノード500のエラー受信部510は、エラーメッセージを受

50

信すると、リカバリ処理が実行されていないことを確認後、論理ディスク管理部 5 3 0 にリカバリ指示を出す。

【 0 1 0 6 】

[ステップ S 1 3] 論理ディスク管理部 5 3 0 内のスライス割当処理部 5 3 2 は、ライトエラーメッセージを送信したディスクノード 4 0 0 を切り離すため、メタデータ更新処理を行う。メタデータ更新処理は、故障通知処理 (ステップ S 1 3 a) と論理ディスクメタデータ更新処理 (ステップ S 1 3 b) とに分かれる。

【 0 1 0 7 】

[ステップ S 1 3 a] 論理ディスク管理部 5 3 0 は、まずシングルプライマリ化を指示するメタデータ変更要求をディスクノード 1 0 0 , 2 0 0 , 3 0 0 に送信する。メタデータ変更要求は、ディスクノード 4 0 0 が管理しているスライスとペアとなるスライスについて、シングルプライマリスライスに変更すること (シングルプライマリ化) の要求である。メタデータ変更要求には、故障したディスクノード 4 0 0 のディスクノード ID が含まれる。

10

【 0 1 0 8 】

[ステップ S 1 4] シングルプライマリ化のメタデータ変更要求を受信したディスクノード 1 0 0 では、メタデータ管理部 1 6 0 がシングルプライマリ化の処理を実行する。例えば、メタデータ管理部 1 6 0 は、メタデータ記憶部 1 5 0 内のメタデータを参照し、ペアのディスクノード ID の欄に、ディスクノード 4 0 0 のディスクノード ID 「DP 4」が設定されているメタデータを検索する。そして、メタデータ管理部 1 6 0 は、検索で合致したメタデータテーブル 1 5 1 内のメタデータについて、状態をシングルプライマリスライスを示す「SP」に変更する。またメタデータ管理部 1 6 0 は、検索で合致したメタデータテーブル 1 5 1 内のメタデータについて、ペアのディスクノード ID、およびペアのスライス ID の欄に該当する領域の情報を削除する。情報が削除された領域の内容は、例えば「NULL」となる。

20

【 0 1 0 9 】

メタデータ管理部 1 6 0 は、メタデータ記憶部 1 5 0 内のメタデータの更新が完了すると、更新されたメタデータに対応するストレージ装置 1 1 0 内のメタデータの記憶領域に、更新されたメタデータを書き込む。これにより、ストレージ装置 1 1 0 内のメタデータが更新され、ストレージ装置 1 1 0 内のメタデータとメタデータ記憶部 1 5 0 内のメタデータとの同一性が保たれる。

30

【 0 1 1 0 】

メタデータ管理部 1 6 0 は、メタデータ記憶部 1 5 0 とストレージ装置 1 1 0 とのメタデータの更新が完了すると、制御ノード 5 0 0 に対してシングルプライマリ化完了応答を送信する。制御ノード 5 0 0 では、ディスクノード 1 0 0 からのシングルプライマリ化完了応答を受信することで、ディスクノード 1 0 0 でのシングルプライマリ化処理が完了したことを認識する。

【 0 1 1 1 】

なお、他のディスクノード 2 0 0 , 3 0 0 もディスクノード 1 0 0 と同様にシングルプライマリ化処理を実行する。そしてディスクノード 2 0 0 , 3 0 0 は、シングルプライマリ化処理が完了すると、制御ノード 5 0 0 に対してシングルプライマリ化完了応答を送信する。

40

【 0 1 1 2 】

[ステップ S 1 3 b] スライス割当処理部 5 3 2 は、各ディスクノード 1 0 0 , 2 0 0 , 3 0 0 からシングルプライマリ化完了応答を受け取ると、論理ディスクメタデータ記憶部 5 4 0 内のメタデータを更新する。具体的にはスライス割当処理部 5 3 2 は、ディスクノード 4 0 0 が管理しているスライスとペアとなるスライスのメタデータの状態を、シングルプライマリスライスに変更する。またスライス割当処理部 5 3 2 は、ディスクノード 4 0 0 が管理しているスライスとペアとなるスライスのメタデータから、ペアのディスクノード ID、およびペアのスライス ID の欄の情報を削除する。さらにスライス割当処理

50

部 5 3 2 は、ディスクノード 4 0 0 が管理しているスライスのメタデータを論理ディスクメタデータ記憶部 5 4 0 から削除するか、あるいは該当メタデータの状態をすべて異常であることを示す状態フラグ「B」に変更する。

【 0 1 1 3 】

スライス割当処理部 5 3 2 は、故障したディスクノード 4 0 0 を除くすべてのディスクノード 1 0 0 , 2 0 0 , 3 0 0 からメタデータを収集し、論理ディスクメタデータテーブル 5 4 1 を再作成することもできる。ディスクノード 1 0 0 , 2 0 0 , 3 0 0 内のメタデータは故障通知に回答して更新されているため、最新のメタデータを収集すれば、ディスクノード 4 0 0 がシステムに含まれない状態での論理ディスクメタデータテーブル 5 4 1 が作成できる。

10

【 0 1 1 4 】

制御ノード 5 0 0 で更新された論理ディスクメタデータテーブル 5 4 1 内のメタデータは、所定のタイミングでアクセスノード 6 0 0 に通知される。例えば、制御ノード 5 0 0 からの指示により、切り離されたディスクノード 4 0 0 のアクセス処理部の機能を停止させる。するとアクセスノード 6 0 0 がディスクノード 4 0 0 にアクセスすると、アクセスエラーとなる。この場合、アクセスノード 6 0 0 は、制御ノード 5 0 0 にメタデータ照会要求を送信する。制御ノード 5 0 0 は、メタデータ照会要求に回答して、最新のメタデータをアクセスノード 6 0 0 に通知する。また制御ノード 5 0 0 は、論理ディスクメタデータテーブル 5 4 1 内のメタデータを更新した場合、アクセスノード 6 0 0 からのメタデータ照会要求を待たずに、更新されたメタデータをアクセスノード 6 0 0 に送信してもよい。

20

【 0 1 1 5 】

ディスクノード 4 0 0 を切り離した後の論理ディスクメタデータテーブル 5 4 1 では、ディスクノード 4 0 0 が管理しているスライスとペアとなるスライスがすべてシングルプライマリスライスに変更されている。これにより、その後のアクセスノード 6 0 0 からのアクセスは、シングルプライマリ化されたスライスに対して実行されることとなる。また、ディスクノード 4 0 0 が管理しているスライスとペアとなるスライスがすべてシングルプライマリスライスに変更されていることで、ディスクノード 4 0 0 がミラーリング処理におけるコピー先となることもなくなる。さらに、ディスクノード 4 0 0 が管理しているスライスのメタデータが削除されているか、あるいは状態フラグ「B」に変更される。これにより、その後のセグメントへのスライス割り当てにおいて、ディスクノード 4 0 0 のスライスがセグメントに割り当てられることもなくなる。

30

【 0 1 1 6 】

このようにして、マルチノードストレージシステムの各装置からのディスクノード 4 0 0 へのアクセスを停止することができる。ディスクノード 4 0 0 をアクセス先とする処理がなくなることで、マルチノードストレージシステムからのディスクノード 4 0 0 の切り離しが完了する。

【 0 1 1 7 】

スライス割当処理部 5 3 2 は、ディスクノード 4 0 0 の切り離しが完了すると、リカバリ処理部 5 3 3 にリカバリ処理の開始を指示する。すると、リカバリ処理部 5 3 3 の制御によりリカバリ処理が実行される。

40

【 0 1 1 8 】

図 1 2 は、リカバリ処理手順の一例を示すシーケンス図である。以下、リカバリ処理の手順をステップ番号に沿って説明する。

[ステップ S 2 1] 制御ノード 5 0 0 のリカバリ処理部 5 3 3 は、故障したディスクノード 4 0 0 の切り離しが完了するとリカバリ処理を開始する。リカバリ処理は、リザーブスライスの割り当て処理 (ステップ S 2 1 a)、スライスコピー処理 (ステップ S 2 1 b)、およびスライス状態更新処理 (ステップ S 2 1 c) に分かれる。

【 0 1 1 9 】

[ステップ S 2 1 a] リカバリ処理部 5 3 3 は、論理ディスクメタデータテーブル 5 4

50

1を参照し、シングルプライマリスライス(状態フラグが「SP」のスライス)を選択する。図12の例では、ディスクノード200が管理するスライスが選択されている。シングルプライマリスライスが割り当てられているセグメントは、欠損セグメントである。そこでリカバリ処理部533は、論理ディスクメタデータテーブル541を参照し、欠損セグメントに対してセカンダリスライスとして割り当てるスライスを、フリースライス(状態が「F」のスライス)の中から選択する。選択されるフリースライスは、シングルプライマリスライスとは異なるディスクノードであり、かつライトエラーが発生していないディスクノードで管理されているスライスである。図12の例では、ディスクノード300で管理されているスライスが選択されている。

【0120】

リカバリ処理部533は、選択したシングルプライマリスライスを管理しているディスクノード200に対して、スライスのプライマリ化を指示するメタデータ変更要求を送信する。プライマリ化を指示するメタデータ変更要求には、欠損セグメントに割り当てられているシングルプライマリスライスのスライスIDが含まれる。

【0121】

またリカバリ処理部533は、選択したフリースライスを管理するディスクノード300に対して、リザーブスライスへの変更を指示するメタデータ変更要求を送信する。このメタデータ変更要求には、変更対象のメタデータに対応するスライスのスライスID、選択したシングルプライマリスライスを管理するディスクノードのディスクノードID、および選択したシングルプライマリスライスのスライスIDが含まれる。

【0122】

[ステップS22]ディスクノード200のメタデータ管理部260は、シングルプライマリ化を指示するメタデータ変更要求を受け取ると、メタデータ変更要求で指定されたスライスをシングルプライマリスライスに変更する。具体的には、メタデータ管理部260は、メタデータ変更要求で示される変更対象のスライスのメタデータを、メタデータ記憶部内から選択する。次にメタデータ管理部260は、選択したメタデータの状態の欄の種別フラグを、「SP」に変更する。またメタデータ管理部260は、ストレージ装置210内の変更対象のスライスのメタデータについても、状態がプライマリスライスになるように変更する。そして、メタデータ管理部260は、メタデータの更新が完了すると、更新完了応答を制御ノード500に送信する。

【0123】

[ステップS23]リザーブスライスへの変更を指示するメタデータ変更要求を受信したディスクノード300では、メタデータ管理部360が指定されたフリースライスの状態を、リザーブスライスに変更する。具体的には、メタデータ管理部360は、メタデータ変更要求で示される変更対象のスライスのメタデータを、メタデータ記憶部内から選択する。次にメタデータ管理部360は、選択したメタデータの状態の欄の種別フラグを、「R」に変更する。またメタデータ管理部360は、ストレージ装置310内の変更対象のスライスのメタデータについても、状態がリザーブスライスとなるように変更する。そして、メタデータ管理部360は、メタデータの更新が完了すると、更新完了応答を制御ノード500に送信する。

【0124】

なお、リカバリ処理部533は、ステップS22、S23のメタデータ更新が完了すると、論理ディスクメタデータ記憶部540内のメタデータを、ディスクノード200、300で更新されたメタデータと同じ内容に更新する。

【0125】

[ステップS21b]リカバリ処理部533は、ディスクノード200、300からメタデータ更新完了の応答を受け取ると、欠損セグメントのプライマリスライスを管理しているディスクノード200に対して、スライスコピー要求を送信する。スライスコピー要求には、欠損セグメントのプライマリスライスのスライスID、リザーブスライスを管理するディスクノードのディスクノードID、およびリザーブスライスのスライスIDが含

10

20

30

40

50

まれる。

【 0 1 2 6 】

[ステップ S 2 4] スライスコピー要求を受信したディスクノード 2 0 0 のアクセス処理部 2 2 0 は、指定されたスライスのペアのスライスを管理するディスクノード 3 0 0 に対して、指定されたスライス内のデータを送信する。データを送信する際には、データの格納先として、スライスコピー要求で指定されたスライスのペアのスライスのスライス ID が指定される。

【 0 1 2 7 】

[ステップ S 2 5] ディスクノード 3 0 0 のアクセス処理部 3 2 0 は、ディスクノード 2 0 0 から送られたデータを受信し、指定されたスライスにデータを格納する。アクセス処理部 3 2 0 は、データの格納が完了すると、書き込み完了応答をディスクノード 2 0 0 に送信する。

10

【 0 1 2 8 】

ディスクノード 2 0 0 のアクセス処理部 2 2 0 は、ディスクノード 3 0 0 からの書き込み完了応答を受け取ると、制御ノード 5 0 0 に対してスライスコピー完了の応答を送信する。

【 0 1 2 9 】

[ステップ S 2 1 c] 制御ノード 5 0 0 のリカバリ処理部 5 3 3 は、ディスクノード 2 0 0 からスライスコピー完了の応答を受け取ると、ディスクノード 2 0 0 に対してメタデータへのペアのスライスの設定を指示するメタデータ変更要求を送信する。このメタデータの変更要求には、欠損セグメントのプライミスライスのスライス ID、リザーブスライスを管理するディスクノードのディスクノード ID、およびリザーブスライスのスライス ID が含まれる。

20

【 0 1 3 0 】

また、リカバリ処理部 5 3 3 は、ディスクノード 2 0 0 からスライスコピー完了の応答を受け取ると、ディスクノード 3 0 0 に対してスライスのセカンダリ化を指示するメタデータ変更要求を送信する。セカンダリ化を指示するメタデータ変更要求には、欠損セグメントに割り当てられているリザーブスライスのスライス ID が含まれる。

【 0 1 3 1 】

[ステップ S 2 6] メタデータへのペアのスライスの設定を指示するメタデータ変更要求を受信したディスクノード 2 0 0 では、メタデータ管理部 2 6 0 がペア相手設定処理を行う。具体的には、メタデータ管理部 2 6 0 は、メタデータ変更要求で示される変更対象のスライスのメタデータを、メタデータ記憶部内から選択する。次にメタデータ管理部 2 6 0 は、選択したメタデータのペアのディスクノード ID、ペアのスライス ID の欄に、メタデータ変更要求に含まれているリザーブスライスのディスクノード ID とリザーブスライスのスライス ID とを設定する。またメタデータ管理部 2 6 0 は、ストレージ装置 2 1 0 内の変更対象のスライスのメタデータに対しても、ペアのディスクノード ID とペアのスライス ID とを設定する。そして、メタデータ管理部 2 6 0 は、メタデータの更新が完了すると、更新完了応答を制御ノード 5 0 0 に送信する。

30

【 0 1 3 2 】

[ステップ S 2 7] ディスクノード 3 0 0 のメタデータ管理部 3 6 0 は、セカンダリ化を指示するメタデータ変更要求を受け取ると、メタデータ変更要求で指定されたスライスをセカンダリスライスに変更する。具体的には、メタデータ管理部 3 6 0 は、メタデータ変更要求で示される変更対象のスライスのメタデータを、メタデータ記憶部内から選択する。次にメタデータ管理部 3 6 0 は、選択したメタデータの状態の欄の種別フラグを「S」に変更する。またメタデータ管理部 3 6 0 は、ストレージ装置 3 1 0 内の変更対象のスライスのメタデータについても、状態がセカンダリスライスになるように変更する。そして、メタデータ管理部 3 6 0 は、メタデータの更新が完了すると、更新完了応答を制御ノード 5 0 0 に送信する。

40

【 0 1 3 3 】

50

なお、リカバリ処理部 5 3 3 は、ステップ S 2 6 , S 2 7 のメタデータ更新が完了すると、論理ディスクメタデータ記憶部 5 4 0 内のメタデータを、ディスクノード 2 0 0 , 3 0 0 で更新されたメタデータと同じ内容に更新する。

【 0 1 3 4 】

以後、制御ノード 5 0 0 のリカバリ処理部 5 3 3 は、シングルプライマリスライスが無くなるまで、ステップ S 2 1 a , S 2 1 b , S 2 1 c の処理を繰り返し実行する。

ところで、図 1 1、図 1 2 の例では、リカバリ処理が実行されていない状況下においてディスクノード 4 0 0 でエラーが発生したため、ディスクノード 4 0 0 の切り離し処理と、リカバリ処理とが実行されている。一方、リカバリ処理実行中にいずれかのディスクノードでライトエラーが発生した場合、データの喪失を抑制するため、エラーが発生したディスクノードの即時の切り離しは行われない。以下、リカバリ処理実行中のライトエラーの発生によるエラー処理について説明する。

10

【 0 1 3 5 】

図 1 3 は、リカバリ処理実行中のライトエラー発生時の処理手順の一例を示すシーケンス図である。図 1 3 には、リカバリ処理部 5 3 3 によるリカバリ処理実行中に、ディスクノード 1 0 0 においてストレージ装置 1 1 0 へのライトエラーが検出された場合の処理が示されている。以下、図 1 3 に示す処理をステップ番号に沿って説明する。

【 0 1 3 6 】

[ステップ S 3 1] ディスクノード 1 0 0 のアクセス処理部 1 2 0 は、ストレージ装置 1 1 0 へのライトアクセス時にエラーを検出すると、ライトエラーの発生をエラー通知部 1 3 0 に伝える。するとエラー通知部 1 3 0 は、ライトエラーを示すエラーメッセージを制御ノード 5 0 0 に送信する。

20

【 0 1 3 7 】

[ステップ S 3 2] 制御ノード 5 0 0 では、エラー受信部 5 1 0 がディスクノード 1 0 0 からのエラーメッセージを受け取る。図 1 3 の例では、リカバリ処理部 5 3 3 によるリカバリ処理実行中であるため、エラー受信部 5 1 0 はエラー記憶部 5 2 0 に、ディスクノード 1 0 0 でライトエラーが発生したことを示すエラー情報を格納する。

【 0 1 3 8 】

このように、リカバリ処理が継続して実行されている間にディスクノード 1 0 0 でライトエラーが検出されると、ディスクノード 1 0 0 の切り離し処理に代えて、ディスクノード 1 0 0 でライトエラーが発生したことを示すエラー情報の記録が行われる。

30

【 0 1 3 9 】

次に、ディスクノード 1 0 0 におけるライトエラー検出および通知処理例について説明する。

図 1 4 は、ディスクノードでのライトエラー検出および通知処理の手順の一例を示すフローチャートである。以下、図 1 4 に示す処理をステップ番号に沿って説明する。

【 0 1 4 0 】

[ステップ S 4 1] ディスクノード 1 0 0 のアクセス処理部 1 2 0 は、ライトエラーを検出する。例えば、アクセス処理部 1 2 0 は、アクセスノード 6 0 0 からのデータライトのアクセス要求に応じてストレージ装置 1 1 0 に書き込みを行う際に、書き込みに失敗すると、ライトエラーが発生したものと判断する。また、アクセス処理部 1 2 0 は、他のディスクノード 2 0 0 , 3 0 0 からのミラーリング処理におけるライト要求に応じてストレージ装置 1 1 0 に書き込みを行う際に、書き込みに失敗したとき、ライトエラーが発生したものと判断することもできる。さらに、アクセス処理部 1 2 0 は、他のディスクノード 2 0 0 , 3 0 0 からのスライスデータのコピー要求に応じてストレージ装置 1 1 0 に書き込みを行う際に、書き込みに失敗すると、ライトエラーが発生したものと判断することもできる。

40

【 0 1 4 1 】

なお、ストレージ装置 1 1 0 内で定期的なライトエラーチェック処理が自動実行されている場合もある。この場合、エラーチェック処理においてライトエラーが検出されたこと

50

の通知をストレージ装置 110 から受け取ることで、アクセス処理部 120 がライトエラーの発生を認識できる。

【0142】

ところで、ストレージ装置は、多くの場合、データの書き込みができない単位記憶領域（ブロック）が発生すると、予備に用意されていたブロックにデータの書き込みを行う。このような書き込み先を変更する処理により、ストレージ装置内でライトエラーに対するリカバリが完了する。ストレージ装置内でリカバリが完了した場合、ストレージ装置は、ディスクノードからのデータのライト要求に対してエラー応答をせずにすむ。ただし、ストレージ装置における予備のブロックの数には限りがある。予備のブロックがすべて使用された後にデータの書き込みができないブロックが発生すると、そのブロックに書き込むべきデータの書き込み先が確保できない。この場合、ストレージ装置はディスクノードに対してライトエラーを応答する。アクセス処理部 120 は、ストレージ装置 110 内部でのリカバリが不能となった場合に出力されるライトエラーについてのみ、制御ノードへのエラー通知の対象のライトエラーと判断することができる。

10

【0143】

[ステップ S42] ライトエラーを検出したアクセス処理部 120 は、ライトエラーが発生したことを示すエラー情報を、ライトエラーが発生したストレージ装置のディスク ID に対応付けてエラー記憶部 140 に格納する。例えばアクセス処理部 120 は、エラー記憶部 140 に格納されているライトエラーの発生の有無を示すフラグを、「ライトエラー有り」の状態に変更する。

20

【0144】

[ステップ S43] アクセス処理部 120 は、エラー通知部 130 に対して、ライトエラーの発生を通知する。エラー通知部 130 は、アクセス処理部 120 からのエラー発生の通知を受けて、制御ノード 500 に対して、ライトエラーの発生を示すエラーメッセージを送信する。

【0145】

このようにして、ディスクノード 100 において、ストレージ装置 110 のライトエラー発生を検出すると、ライトエラーを示すエラー情報が記録されると共に、エラーメッセージが制御ノード 500 に送信される。エラーメッセージを受信した制御ノード 500 では、リカバリ処理を開始するか、あるいはエラー情報の記録が行われる。

30

【0146】

図 15 は、エラーメッセージを受信したときのエラー処理手順の一例を示すフローチャートである。以下、図 15 に示す処理をステップ番号に沿って説明する。

[ステップ S51] 制御ノード 500 のエラー受信部 510 は、ディスクノード 100 からのライトエラーを示すエラーメッセージを受信する。

【0147】

[ステップ S52] エラー受信部 510 は、リカバリ処理中か否かを判断する。例えば、リカバリ処理部 533 の実行プロセスが起動されている場合、リカバリ処理中と判断される。また、論理ディスクメタデータ記憶部 540 を参照し、シングルプライマリスライスが存在している間はリカバリ処理中であると判断することもできる。リカバリ処理中であれば、処理がステップ S53 に進められる。リカバリ処理中でなければ、処理がステップ S54 に進められる。

40

【0148】

[ステップ S53] エラー受信部 510 はリカバリ処理中であれば、エラー記憶部 520 に、ディスクノード 100 でライトエラーが発生したことを示すエラー情報を記録する。例えばエラー受信部 510 は、エラー記憶部 520 内のエラー管理テーブル 521 から、エラーメッセージを送信したディスクノードのディスクノード ID を検索する。そしてエラー受信部 510 は、検索で該当したディスクノード ID に対応付けられたライトエラーの欄に、ライトエラー有りを示すフラグを設定する。その後、エラー処理が終了する。

【0149】

50

[ステップ S 5 4] エラー受信部 5 1 0 はリカバリ処理中でなければ、リカバリ処理の開始を、リカバリ処理部 5 3 3 に指示する。これにより、リカバリ処理部 5 3 3 によるリカバリ処理が開始される。その後、エラー処理が終了する。

【 0 1 5 0 】

このようにして、ライトエラーのエラーメッセージ受信時には、リカバリ処理中か否かによって処理が分かれる。すなわち、リカバリ処理中であればエラー情報の記録が行われ、リカバリ処理中でなければリカバリ処理が開始される。

【 0 1 5 1 】

図 1 6 は、リカバリ処理手順の一例を示すフローチャートである。以下、図 1 6 に示す処理をステップ番号に沿って説明する。

[ステップ S 6 1] リカバリ処理部 5 3 3 は、エラーメッセージを送信したディスクノードを切り離し対象とする。このとき、リカバリ処理部 5 3 3 は、エラー管理テーブル 5 2 1 から、切り離し対象のディスクノードに対応するレコードを削除する。

【 0 1 5 2 】

[ステップ S 6 2] リカバリ処理部 5 3 3 は、切り離し対象のディスクノードが管理するスライスとペア相手のスライスを検出する。例えば、リカバリ処理部 5 3 3 は、論理ディスクメタデータテーブル 5 4 1 のペアのディスクノード ID の欄から、切り離し対象のディスクノードのディスクノード ID を検索する。検索により該当したメタデータに対応するスライスが、切り離し対象のディスクノードが管理するスライスとペア相手である。

【 0 1 5 3 】

[ステップ S 6 3] リカバリ処理部 5 3 3 は、ペア相手のスライスを管理するディスクノードに対して、ペア相手のスライスのシングルプライマリ化を指示する。例えば、リカバリ処理部 5 3 3 は、ペア相手のスライスのスライス ID を含むシングルプライマリ化を指示するメタデータ変更要求を、ペア相手のスライスを管理するディスクノードに送信する。

【 0 1 5 4 】

[ステップ S 6 4] リカバリ処理部 5 3 3 は、すべてのペア相手のスライスのシングルプライマリ化が完了したか否かを判断する。例えばリカバリ処理部 5 3 3 は、シングルプライマリ化を指示するメタデータ変更要求のすべてに対して、シングルプライマリ化完了の応答が返された場合に、シングルプライマリ化が完了したものと判断する。シングルプライマリ化が完了した場合、処理がステップ S 6 5 に進められる。シングルプライマリ化が未完了であれば、ステップ S 6 4 が繰り返される。

【 0 1 5 5 】

なお、ディスクノードでのシングルプライマリ化のメタデータ変更が完了すると、リカバリ処理部 5 3 3 は、論理ディスクメタデータテーブル 5 4 1 内のペア相手のスライスのメタデータについても、状態をシングルプライマリスライスに変更する。

【 0 1 5 6 】

[ステップ S 6 5] リカバリ処理部 5 3 3 は、シングルプライマリスライスが存在するか否かを判断する。例えばリカバリ処理部 5 3 3 は、論理ディスクメタデータテーブル 5 4 1 から、状態欄の種別フラグが「S P」のメタデータを検索する。少なくとも 1 つのメタデータが検索で見つかった場合、シングルプライマリスライスが存在すると判断される。種別フラグが「S P」のメタデータが検出されなければ、シングルプライマリスライスは存在しないと判断される。シングルプライマリスライスが存在する場合、処理がステップ S 6 6 に進められる。シングルプライマリスライスが存在しない場合、処理がステップ S 7 1 に進められる。

【 0 1 5 7 】

[ステップ S 6 6] リカバリ処理部 5 3 3 は、シングルプライマリスライスの 1 つを、リカバリ対象として選択する。例えばリカバリ処理部 5 3 3 は、論理ディスクメタデータテーブル 5 4 1 から種別フラグが「S P」のメタデータを 1 つ抽出し、抽出したメタデータに対応するスライスをリカバリ対象として選択する。

10

20

30

40

50

【 0 1 5 8 】

[ステップ S 6 7] リカバリ処理部 5 3 3 は、リザーブ化対象のスライスを選択する。例えばリカバリ処理部 5 3 3 は、論理ディスクメタデータテーブル 5 4 1 から種別フラグが「F」であり、かつリカバリ対象のスライスと異なるディスクノードで管理されているスライスのメタデータを 1 つ抽出する。そしてリカバリ処理部 5 3 3 は、抽出したメタデータに対応するスライスを、リザーブ化対象として選択する。

【 0 1 5 9 】

この際、リカバリ処理部 5 3 3 は、エラー管理テーブル 5 2 1 を参照してライトエラーが発生しているディスクノードを判断し、ライトエラーが発生しているディスクノードで管理されているスライスを、リザーブ化対象として選択しないようにしてもよい。すなわち、ライトエラーが発生しているディスクノードが管理するスライスは、リザーブ対象として選択しても、スライスコピー時にエラーとなり、リザーブ化対象の再選択を行うこととなる可能性がある。ライトエラーが発生しているディスクノードで管理されているスライスをリザーブ化対象から除外することで、リザーブ化対象の再選択処理の発生を抑制できる。

10

【 0 1 6 0 】

[ステップ S 6 8] リカバリ処理部 5 3 3 は、リカバリ対象のスライスとリザーブ化対象のスライスに関するメタデータを変更する。例えば、リカバリ対象スライスは、シングルプライマリスライスからプライマリスライスに変更される。リザーブ化対象のスライスは、フリースライスからリザーブスライスに変更される。具体的にはリカバリ処理部 5 3 3 は、各スライスを管理しているディスクノードに対してメタデータ変更要求を送信することで、ディスクノード内のメタデータを変更する。また、リカバリ処理部 5 3 3 は、論理ディスクメタデータテーブル 5 4 1 内のリカバリ対象のスライスとリザーブ化対象のスライスとのそれぞれのメタデータを変更する。

20

【 0 1 6 1 】

[ステップ S 6 9] リカバリ処理部 5 3 3 は、プライマリスライスのデータのリザーブスライスへのコピーを指示するスライスコピー要求を、シングルプライマリスライスを管理しているディスクノードに送信する。コピーが完了すると、ディスクノードからスライスコピー完了の応答が返される。

【 0 1 6 2 】

[ステップ S 7 0] リカバリ処理部 5 3 3 は、リカバリ対象のスライスにペア相手のスライスを設定するメタデータの変更を行う。またリカバリ処理部 5 3 3 は、リザーブ化したスライスをセカンダリ化するメタデータの変更を行う。具体的には、リカバリ処理部 5 3 3 は、プライマリスライスを管理しているディスクノードに対して、ペア相手のスライスとしてリザーブスライスを指定したメタデータ変更要求を送信する。また、リカバリ処理部 5 3 3 は、リザーブスライスを管理するディスクノードに対して、リザーブスライスからセカンダリスライスへの変更を指示するメタデータ変更要求を送信する。さらにリカバリ処理部 5 3 3 は、論理ディスクメタデータテーブル 5 4 1 内のプライマリスライスとリザーブスライスとのそれぞれのメタデータを変更する。その後、処理がステップ S 6 5 に進められる。

30

40

【 0 1 6 3 】

[ステップ S 7 1] リカバリ処理部 5 3 3 は、シングルプライマリスライスが無くなると、切り離されていないライトエラー発生ディスクノードの有無を判断する。例えばリカバリ処理部 5 3 3 は、エラー記憶部 5 2 0 を参照し、ライトエラー有りのディスクノードがある場合、切り離されていないライトエラー発生ディスクノードがあると判断する。切り離されていないライトエラー発生ディスクノードがある場合、処理がステップ S 7 2 に進められる。切り離されていないライトエラー発生ディスクノードがない場合、処理が終了する。

【 0 1 6 4 】

[ステップ S 7 2] リカバリ処理部 5 3 3 は、切り離されていないライトエラー発生デ

50

ディスクノードを1台選択し、切り離し対象とする。このとき、リカバリ処理部533は、エラー管理テーブル521から、切り離し対象のディスクノードに対応するレコードを削除する。その後、処理がステップS62に進められる。

【0165】

このように、リカバリ処理中にディスクノードでライトエラーが発生した場合、そのディスクノードの切り離しは、シングルプライマリスライスがなくなるまで延期される。シングルプライマリスライスがなくなるということは、すべてのセグメントのデータの二重化が復旧されている。シングルプライマリスライスがなくなると、ライトエラーが発生したディスクノードの切り離しと、その切り離しに伴って発生したシングルプライマリスライスの二重化復旧が行われる。

10

【0166】

ライトエラーが発生したディスクノードに対して外部からのアクセスがあった場合、アクセス対象のスライスが割り当てられたセグメントに関しては、ライトエラーが発生したディスクノードの切り離しを待たずに欠損セグメントに変更される。その結果、アクセス対象のスライスが割り当てられたセグメントに対する二重化復旧処理が実行される。ディスクノードに対する外部からのアクセスには、例えば、アクセスノード600からのリードまたはライトのアクセス要求、他のディスクノードからのミラーライト要求、リカバリ処理中のスライスデータコピー時のスライスコピー要求または書き込み要求などがある。ライトエラーが発生したディスクノードに対して外部からアクセスがあると、アクセス対象のスライスが制御ノード500に通知されることとなる。

20

【0167】

ただし、ライトエラーが発生したディスクノードが管理するデータに対する外部からのアクセスであっても、制御ノード500からのデータリードのアクセスに関しては、その要求に従ってディスクノードでデータリードのアクセス処理が実行される。例えば、リカバリ処理やスライス割当処理におけるスライスコピー要求が、ライトエラーが発生したディスクノードに出された場合、そのディスクノードはスライスコピーのコピー元となる。コピー元はコピー対象のスライスのデータを読み出せばよいので、ライトエラーが発生したディスクノードでも実行可能である。

【0168】

ディスクノードに対してどのようなアクセスが行われたのかによって、アクセス対象のスライスを制御ノード500に通知する手順が異なる。以下、アクセス対象のスライスを制御ノード500に通知する手順について説明する。

30

【0169】

図17は、ライトエラーが発生したディスクノードへのアクセス要求時の処理手順を示すシーケンス図である。以下、図17に示す処理をステップ番号に沿って説明する。なお図17の例では、ディスクノード100ではライトエラーの発生を示すエラー情報が記憶されているものとする。また、制御ノード500のリカバリ処理部533はリカバリ処理を実行中であるものとする。

【0170】

[ステップS81] アクセスノード600のアクセス要求部630は、端末装置からの要求に応じてディスクノード100にリードまたはライトのアクセス要求を送信する。

40

[ステップS82] アクセスノード600からのアクセス要求は、ディスクノード100のアクセス処理部120で受け取られる。アクセス処理部120は、エラー記憶部140を参照し、ライトエラーが発生していることを認識する。そして、アクセス処理部120は、アクセスノード600に対してエラー応答を送信する。この処理の詳細は後述する(図18参照)。

【0171】

[ステップS83] エラー応答を受信したアクセスノード600のアクセス要求部630は、メタデータ照会部610にメタデータの照会を依頼する。メタデータ照会部610は、アクセス要求のアクセス対象であったセグメントを指定して、該当セグメントのメタ

50

データ照会要求を制御ノード 500 に送信する。

【0172】

[ステップ S84] 制御ノード 500 では、スライス割当処理部 532 がメタデータ照会要求を受け取る。スライス割当処理部 532 は、メタデータ照会要求を受信したことで、受信したメタデータ照会要求で示されたセグメントのプライミスライスへのアクセスがエラーとなったことを認識する。またスライス割当処理部 532 はエラー記憶部 520 を参照し、エラーとなったスライスが、ライトエラーが検出されたディスクノード 100 で管理されているスライスであることを認識する。そこで、スライス割当処理部 532 は、スライス割当処理を実行する。なお図 17 の例では、ディスクノード 200 が管理しているスライスが、メタデータ照会要求で指定されたセグメントのプライミスライスとして割り当てられたものとする。スライス割当処理の詳細は後述する(図 22 参照)。

10

【0173】

スライス割当処理が完了すると、スライス割当処理部 532 はメタデータ照会要求への応答として、メタデータ照会要求で指定されたセグメントのプライミスライスのメタデータをアクセスノード 600 に送信する。

【0174】

[ステップ S85] アクセスノード 600 のアクセス要求部 630 は、取得したメタデータに従って、アクセスをリトライする。例えばアクセス要求部 630 は、取得したメタデータをアクセス用メタデータ記憶部 620 に格納する。またアクセス要求部 630 は、取得したメタデータに基づいて、アクセス対象のセグメントに割り当てられたスライスを判断し、該当スライスを管理しているディスクノードに対してアクセス要求を送信する。図 17 の例では、ディスクノード 200 に対してアクセス要求を送信する。

20

【0175】

[ステップ S86] ディスクノード 200 のアクセス処理部 220 は、アクセス要求に応じてストレージ装置 210 へデータアクセスを行う。そしてアクセス処理部 220 は、アクセス結果をアクセスノード 600 に応答する。

【0176】

次に、アクセスノード 600 からのアクセス要求に応じたディスクノードにおけるアクセス処理手順について詳細に説明する。

図 18 は、アクセス処理手順の一例を示すフローチャートである。以下、図 18 に示す処理をステップ番号に沿って説明する。

30

【0177】

[ステップ S91] ディスクノード 100 のアクセス処理部 120 は、アクセスノード 600 からのデータリードまたはデータライトのアクセス要求を受信する。

[ステップ S92] アクセス処理部 120 は、エラー記憶部 140 を参照し、ライトエラーの有無を判断する。ライトエラーがある場合、処理がステップ S93 に進められる。ライトエラーがなければ、処理がステップ S94 に進められる。

【0178】

[ステップ S93] アクセス処理部 120 は、ライトエラーがある場合、アクセスノード 600 に対してエラーを通知する。その後、アクセス処理が終了する。

40

[ステップ S94] アクセス処理部 120 は、ライトエラーがない場合、ストレージ装置 110 に対して、アクセス要求に応じたデータリードまたはデータライトのアクセスを行う。アクセス処理部 120 は、ストレージ装置 110 へのアクセスの結果を、アクセスノード 600 に応答する。例えばデータリードのアクセスであれば、アクセス処理部 120 は、読み出したデータをアクセスノード 600 に送信する。またデータライトのアクセスであれば、アクセス処理部 120 は、書き込み完了の応答メッセージをアクセスノード 600 に送信する。その後、アクセス処理が終了する。

【0179】

このようにして、アクセス要求に応じてディスクノード 100 でデータアクセスが実行される。

50

ディスクノード100へのアクセスは、アクセスノード600からのアクセス要求以外に、他のディスクノードからミラーライトのアクセスがある。以下、ミラーライト処理の手順について説明する。

【0180】

図19は、ライトエラーが発生したディスクノードへのミラーライト処理の手順を示すシーケンス図である。以下、図19に示す処理をステップ番号に沿って説明する。図19の例では、ディスクノード300がプライマリスライスを管理しており、ディスクノード100がセカンダリスライスを管理しているセグメントに関して、アクセスノード600がデータライトのアクセス要求を出力した場合を想定している。

【0181】

[ステップS101] アクセスノード600のアクセス要求部630は、ディスクノード300に対してデータライトのアクセス要求を送信する。

[ステップS102] ディスクノード300のアクセス処理部320は、データライトのアクセス要求を受け取ると、ストレージ装置310にデータを書き込む。またアクセス処理部320は、セカンダリスライスを管理しているディスクノード100に対してデータのミラーライト要求を送信する。

【0182】

[ステップS103] ディスクノード100のアクセス処理部120は、ミラーライト要求を受け取ると、ライトエラーの有無を判断する。図19の例では、すでにディスクノード100においてライトエラーが発生している。アクセス処理部120は、ライトエラーがある場合、ミラーライトエラーをディスクノード300に応答する。

【0183】

[ステップS104] ディスクノード300のアクセス処理部320は、ディスクノード100からのミラーライト応答を受け取ると、ライトエラーを示すエラーメッセージをアクセスノード600に送信する。

【0184】

[ステップS105] エラー応答を受信したアクセスノード600のアクセス要求部630は、メタデータ照会部610にメタデータの照会を依頼する。メタデータ照会部610は、アクセス要求のアクセス対象であったセグメントを指定して、該当セグメントのメタデータ照会要求を制御ノード500に送信する。

【0185】

[ステップS106] 制御ノード500のスライス割当処理部532は、メタデータ照会要求に応答し、スライス割り当て処理を行う。ステップS106およびその後の処理は、図17のステップS84～S86と同様である。

【0186】

このようにして、ライトエラーが発生したディスクノードへのミラーライトがエラーとなり、アクセスノード600からメタデータ照会要求が送信される。制御ノード500では、メタデータ照会要求に応答して、アクセス対象のセグメントに割り当てられたスライスが調査され、ライト要求のアクセス先であったスライスの割り当てが解除される。

【0187】

ライトエラーが発生したディスクノード100への他のディスクノードからのアクセスとしては、スライスコピーがある。例えばリカバリ処理によりディスクノード100が管理するスライスをリザーブスライスとしていた場合、プライマリスライスのデータがリザーブスライスにコピーされる。このようなスライス全体のデータのコピー処理が、スライスコピーである。ディスクノード100をコピー先としてスライスコピーを実行している最中にディスクノード100でライトエラーが発生する場合もある。そのような場合、スライスコピーがエラーとなる。

【0188】

図20は、ライトエラーが発生したディスクノードへのスライスコピー処理時の処理手順を示すシーケンス図である。以下、図20の処理をステップ番号に沿って説明する。

10

20

30

40

50

【ステップ S 1 1 1】ディスクノード 3 0 0 のアクセス処理部 3 2 0 は、制御ノード 5 0 0 からのスライスコピー要求に回答して、自己の管理するプライミスライスのデータをリザーブスライスにコピーする。図 2 0 の例では、ディスクノード 1 0 0 が管理するスライスがリザーブスライスである。

【0 1 8 9】

【ステップ S 1 1 2】ディスクノード 1 0 0 のアクセス処理部 1 2 0 は、ディスクノード 3 0 0 からスライスコピー対象のデータを取得する。このときアクセス処理部 1 2 0 は、ライトエラーの有無を判断する。ライトエラーが発生している場合、アクセス処理部 1 2 0 は、ディスクノード 3 0 0 に対してエラーを応答する。なお、スライスコピー対象のデータ受信時のディスクノード 1 0 0 の処理の詳細は後述する（図 2 1 参照）。

10

【0 1 9 0】

【ステップ S 1 1 3】ディスクノード 3 0 0 のアクセス処理部 3 2 0 は、ディスクノード 1 0 0 からのエラー応答を受け取ると、制御ノード 5 0 0 に対してエラーメッセージを送信する。エラーメッセージには、スライスコピーでエラーとなったスライスを管理しているディスク ID や、スライスコピーでエラーとなったスライスのスライス ID が含まれる。

【0 1 9 1】

【ステップ S 1 1 4】制御ノード 5 0 0 のスライス割当処理部 5 3 2 は、スライスコピーでエラーとなったスライスが割り当てられていたセグメントのスライス割り当て処理を行う。この処理の詳細は後述する（図 2 3 参照）。

20

【0 1 9 2】

このようにして、ライトエラーが発生したディスクノード 1 0 0 へのスライスコピーが行われると、スライスコピーがエラーとなり、制御ノード 5 0 0 でスライス割り当て処理が開始される。

【0 1 9 3】

図 2 1 は、スライスコピーのデータを受信したディスクノードの処理手順の一例を示すフローチャートである。以下、図 2 1 に示す処理をステップ番号に沿って説明する。

【ステップ S 1 2 1】ディスクノード 1 0 0 のアクセス処理部 1 2 0 は、スライスコピー対象のデータを受信する。

【0 1 9 4】

30

【ステップ S 1 2 2】アクセス処理部 1 2 0 は、ライトエラーの記録の有無を判断する。例えばアクセス処理部 1 2 0 は、エラー記憶部 1 4 0 を参照し、ライトエラーの有無を判断する。ライトエラーがある場合、処理がステップ S 1 2 3 に進められる。ライトエラーがない場合、処理がステップ S 1 2 4 に進められる。

【0 1 9 5】

【ステップ S 1 2 3】アクセス処理部 1 2 0 は、スライスコピーによるデータの送信元であるディスクノードに、エラーメッセージを送信する。その後、処理が終了する。

【ステップ S 1 2 4】アクセス処理部 1 2 0 は、受信したデータをストレージ装置 1 1 0 内のスライスコピー先のスライスに書き込む。その後、処理が終了する。

【0 1 9 6】

40

次に、スライス割当処理について詳細に説明する。スライス割当処理は、メタデータ照会要求に応じて実行される場合と、スライスコピー時のエラーメッセージに応じて実行される場合とがある。

【0 1 9 7】

図 2 2 は、メタデータ照会要求時のスライス割当処理手順の一例を示すフローチャートである。以下、図 2 2 に示す処理をステップ番号に沿って説明する。

【ステップ S 1 3 1】制御ノード 5 0 0 のスライス割当処理部 5 3 2 は、アクセスノード 6 0 0 からのメタデータ照会要求を受信する。

【0 1 9 8】

【ステップ S 1 3 2】スライス割当処理部 5 3 2 は、メタデータ照会要求で指定された

50

セグメントを変更対象セグメントとして、変更対象セグメントに割り当てられているスライスを検査する。例えばスライス割当処理部 5 3 2 は、論理ディスクメタデータ記憶部 5 4 0 を参照し、変更対象セグメントに割り当てられているスライスのメタデータを抽出する。そして、スライス割当処理部 5 3 2 は、抽出したメタデータの状態の欄を参照し、スライスの種別を判別する。プライマリスライスとセカンダリスライスとが割り当てられている場合、処理がステップ S 1 3 3 に進められる。プライマリスライスとリザーブスライスとが割り当てられている場合、処理がステップ S 1 3 4 に進められる。シングルプライマリスライスのみが割り当てられている場合、処理がステップ S 1 3 5 に進められる。

【 0 1 9 9 】

[ステップ S 1 3 3] スライス割当処理部 5 3 2 は、プライマリスライスとセカンダリスライスとの割り当て変更処理を実行する。この処理の詳細は後述する(図 2 4 参照)。

[ステップ S 1 3 4] スライス割当処理部 5 3 2 は、プライマリスライスとリザーブスライスとの割り当て変更処理を実行する。この処理の詳細は後述する(図 2 6 参照)。

【 0 2 0 0 】

[ステップ S 1 3 5] スライス割当処理部 5 3 2 は、シングルプライマリスライスの割り当て変更処理を実行する。この処理の詳細は後述する(図 2 8 参照)。

図 2 3 は、スライスコピーエラー時のスライス割り当て処理手順の一例を示すフローチャートである。以下、図 2 3 に示す処理をステップ番号に沿って説明する。

【 0 2 0 1 】

[ステップ S 1 4 1] 制御ノード 5 0 0 のスライス割当処理部 5 3 2 は、ディスクノードからスライスコピーエラーを受信する。

[ステップ S 1 4 2] スライス割当処理部 5 3 2 は、スライスコピー処理が行われたセグメントを変更対象セグメントとして、プライマリスライスとリザーブスライスとの割り当て変更処理を実行する。この処理の詳細は後述する(図 2 6 参照)。

【 0 2 0 2 】

次に、セグメントに割り当てられているスライスの種別に応じた割り当て変更処理の手順について詳細に説明する。

図 2 4 は、プライマリスライスとセカンダリスライスとが割り当てられているセグメントの割り当て変更処理手順の一例を示すフローチャートである。以下、図 2 4 に示す処理をステップ番号に沿って説明する。

【 0 2 0 3 】

[ステップ S 1 5 1] スライス割当処理部 5 3 2 は、変更対象セグメントに割り当てられている各スライスを管理しているディスクノードでのライトエラーの発生の有無を検査する。例えばスライス割当処理部 5 3 2 は、ステップ S 1 3 2 で抽出したメタデータのディスクノード ID の欄を参照し、変更対象セグメントに割り当てられているスライスを管理しているディスクノードを判断する。次に、スライス割当処理部 5 3 2 は、エラー記憶部 5 2 0 を参照し、変更対象セグメントに割り当てられているスライスを管理しているディスクノードにライトエラーが発生している否かを判断する。プライマリスライスとセカンダリスライスとのそれぞれを管理しているディスクノード共にライトエラーが発生していれば、処理がステップ S 1 5 2 に進められる。プライマリスライスを管理しているディスクノードのみにライトエラーが発生していれば、処理がステップ S 1 5 5 に進められる。セカンダリスライスを管理しているディスクノードのみにライトエラーが発生していれば、処理がステップ S 1 5 6 に進められる。

【 0 2 0 4 】

[ステップ S 1 5 2] スライス割当処理部 5 3 2 は、変更対象セグメントに割り当てられているプライマリスライスとセカンダリスライスとに関するメタデータを変更する。具体的には、スライス割当処理部 5 3 2 は、プライマリスライスとセカンダリスライスとのいずれか一方をフリースライスに変更すると共に、リザーブスライスの新規割り当てを行う。なおライトエラーが発生したディスクノードが管理するスライスは、セグメントへの割り当ては解除される。より詳細には、スライス割当処理部 5 3 2 は、セカンダリスライ

10

20

30

40

50

スを管理しているディスクノードへ、フリースライスへの変更を指示するメタデータ変更要求を送信する。またスライス割当処理部 5 3 2 は、フリースライスを管理しているディスクノードへ、リザーブスライスへの変更を指示するメタデータ変更要求を送信する。スライス割当処理部 5 3 2 は、各ディスクノードからメタデータ変更完了の応答を受け取ると、論理ディスクメタデータ記憶部 5 4 0 内のメタデータを、ディスクノードにおける変更内容と同様に変更する。なお第 2 の実施の形態では、セカンダリスライスをフリースライスに変更するものとする。また新規に割り当てられるリザーブスライスを選択する際には、スライス割当処理部 5 3 2 は、エラー記憶部 5 2 0 を参照し、ライトエラーが発生しているディスクノードを判断する。そしてスライス割当処理部 5 3 2 は、エラーが発生していないディスクノードが管理するフリースライスの中から、リザーブスライスとするスライスを選択する。

10

【 0 2 0 5 】

[ステップ S 1 5 3] スライス割当処理部 5 3 2 は、プライマリスライスからリザーブスライスへのスライスのデータのコピーを、プライマリスライスを管理しているディスクノードに対して指示する。

【 0 2 0 6 】

[ステップ S 1 5 4] スライス割当処理部 5 3 2 は、変更対象セグメントに割り当てられているプライマリスライスとリザーブスライスとに関するメタデータを変更する。具体的には、スライス割当処理部 5 3 2 は、プライマリスライスをフリースライスに変更し、リザーブスライスをシングルプライマリスライスに変更する。より詳細には、スライス割当処理部 5 3 2 は、プライマリスライスを管理しているディスクノードへ、プライマリスライスからフリースライスへの変更を指示するメタデータ変更要求を送信する。またスライス割当処理部 5 3 2 は、リザーブスライスを管理しているディスクノードへ、リザーブスライスからシングルプライマリスライスへの変更を指示するメタデータ変更要求を送信する。スライス割当処理部 5 3 2 は、各ディスクノードからメタデータ変更完了の応答を受け取ると、論理ディスクメタデータ記憶部 5 4 0 内のメタデータを、ディスクノードにおける変更内容と同様に変更する。その後、処理が終了する。

20

【 0 2 0 7 】

[ステップ S 1 5 5] スライス割当処理部 5 3 2 は、変更対象セグメントに割り当てられているプライマリスライスとセカンダリスライスとに関するメタデータを変更する。具体的には、スライス割当処理部 5 3 2 は、プライマリスライスをフリースライスに変更すると共に、セカンダリスライスをシングルプライマリスライスに変更する。より詳細には、スライス割当処理部 5 3 2 は、プライマリスライスを管理しているディスクノードへ、プライマリスライスからフリースライスへの変更を指示するメタデータ変更要求を送信する。またスライス割当処理部 5 3 2 は、セカンダリスライスを管理しているディスクノードへ、シングルプライマリスライスへの変更を指示するメタデータ変更要求を送信する。スライス割当処理部 5 3 2 は、各ディスクノードからメタデータ変更完了の応答を受け取ると、論理ディスクメタデータ記憶部 5 4 0 内のメタデータを、ディスクノードにおける変更内容と同様に変更する。その後、処理が終了する。

30

【 0 2 0 8 】

[ステップ S 1 5 6] スライス割当処理部 5 3 2 は、変更対象セグメントに割り当てられているプライマリスライスとセカンダリスライスとに関するメタデータを変更する。具体的には、スライス割当処理部 5 3 2 は、プライマリスライスをシングルプライマリスライスに変更すると共に、セカンダリスライスをフリースライスに変更する。より詳細には、スライス割当処理部 5 3 2 は、プライマリスライスを管理しているディスクノードへ、シングルプライマリスライスへの変更を指示するメタデータ変更要求を送信する。またスライス割当処理部 5 3 2 は、セカンダリスライスを管理しているディスクノードへ、セカンダリスライスからフリースライスへの変更を指示するメタデータ変更要求を送信する。スライス割当処理部 5 3 2 は、各ディスクノードからメタデータ変更完了の応答を受け取ると、論理ディスクメタデータ記憶部 5 4 0 内のメタデータを、ディスクノードにおける

40

50

変更内容と同様に変更する。その後、処理が終了する。

【0209】

図25は、プライマリスライスとセカンダリスライスとが割り当てられているセグメントの割り当て変更処理例を示す図である。図25の例では、変更対象セグメントのプライマリスライスをディスクノード100(DP1)が管理しており、変更対象セグメントのセカンダリスライスをディスクノード200(DP2)が管理しているものとする。またディスクノード300(DP3)には、ライトエラーが発生していないものとする。このような状況において、ディスクノード100が管理するプライマリスライス、ディスクノード200が管理するセカンダリスライス、およびディスクノード300が管理するフリースライスに着目する。

10

【0210】

図25では、いずれのディスクノードにおいてもライトエラーが検出されていない状態を初期状態としている。変更対象セグメントに割り当てられているスライスを管理しているディスクノード100, 200のライトエラーの有無を検査した結果に応じて、メタデータの変更内容が異なる。

【0211】

両方のディスクノード100, 200でライトエラーが発生している場合、各ディスクノード100, 200, 300の処理は以下の通りである。

ディスクノード100が管理しているプライマリスライス(状態フラグ「P」)は、最初のメタデータ変更では変更されず、スライスコピーのコピー元とされる。スライスコピーが完了後、ディスクノード100が管理しているプライマリスライスは、フリースライス(状態フラグ「F」)に変更される。ディスクノード200が管理しているセカンダリスライス(状態フラグ「S」)は、最初のメタデータ変更でフリースライス(状態フラグ「F」)に変更される。ディスクノード300が管理しているフリースライス(状態フラグ「F」)は、最初のメタデータ変更でリザーブスライス(状態フラグ「R」)に変更され、その後のスライスコピーのコピー先とされる。スライスコピーが完了後、ディスクノード300が管理しているリザーブスライスは、シングルプライマリスライス(状態フラグ「SP」)に変更される。

20

【0212】

ディスクノード100のみでライトエラーが発生している場合、各ディスクノード100, 200, 300の処理は以下の通りである。

30

ディスクノード100が管理しているプライマリスライス(状態フラグ「P」)は、フリースライス(状態フラグ「F」)に変更される。ディスクノード200が管理しているセカンダリスライス(状態フラグ「S」)は、シングルプライマリスライス(状態フラグ「SP」)に変更される。ディスクノード300が管理しているフリースライス(状態フラグ「F」)は変化しない。

【0213】

ディスクノード200のみでライトエラーが発生している場合、各ディスクノード100, 200, 300の処理は以下の通りである。

ディスクノード100が管理しているプライマリスライス(状態フラグ「P」)は、シングルプライマリスライス(状態フラグ「SP」)に変更される。ディスクノード200が管理しているセカンダリスライス(状態フラグ「S」)は、フリースライス(状態フラグ「F」)に変更される。ディスクノード300が管理しているフリースライス(状態フラグ「F」)は変化しない。

40

【0214】

図26は、プライマリスライスとリザーブスライスとが割り当てられているセグメントの割り当て変更処理手順の一例を示すフローチャートである。以下、図26に示す処理をステップ番号に沿って説明する。

【0215】

[ステップS161]スライス割当処理部532は、変更対象セグメントに割り当てら

50

れている各スライスを管理しているディスクノードでのライトエラーの発生の有無を検査する。プライミスライスとリザーブスライスとのそれぞれを管理しているディスクノード共にライトエラーが発生している場合、およびプライミスライスを管理しているディスクノードのみにライトエラーが発生している場合、処理がステップ S 1 6 2 に進められる。リザーブスライスを管理しているディスクノードのみにライトエラーが発生していれば、処理がステップ S 1 6 5 に進められる。

【 0 2 1 6 】

[ステップ S 1 6 2] スライス割当処理部 5 3 2 は、変更対象セグメントに割り当てられているプライミスライスとリザーブスライスとに関するメタデータを変更する。具体的には、スライス割当処理部 5 3 2 は、リザーブスライスをフリースライスに変更し、かつリザーブスライスの新規割り当てを行う。より詳細には、スライス割当処理部 5 3 2 は、リザーブスライスを管理しているディスクノードへ、リザーブスライスからフリースライスへの変更を指示するメタデータ変更要求を送信する。またスライス割当処理部 5 3 2 は、フリースライスを管理しているディスクノードへ、フリースライスからリザーブスライスへの変更を指示するメタデータ変更要求を送信する。スライス割当処理部 5 3 2 は、各ディスクノードからメタデータ変更完了の応答を受け取ると、論理ディスクメタデータ記憶部 5 4 0 内のメタデータを、ディスクノードにおける変更内容と同様に変更する。

10

【 0 2 1 7 】

[ステップ S 1 6 3] スライス割当処理部 5 3 2 は、プライミスライスからリザーブスライスへのスライスのデータのコピーを、プライミスライスを管理しているディスクノードに対して指示する。

20

【 0 2 1 8 】

[ステップ S 1 6 4] スライス割当処理部 5 3 2 は、変更対象セグメントに割り当てられているプライミスライスとリザーブスライスとに関するメタデータを変更する。具体的には、スライス割当処理部 5 3 2 は、プライミスライスをフリースライスに変更し、リザーブスライスをシングルプライミスライスに変更する。より詳細には、スライス割当処理部 5 3 2 は、プライミスライスを管理しているディスクノードへ、プライミスライスからフリースライスへの変更を指示するメタデータ変更要求を送信する。またスライス割当処理部 5 3 2 は、リザーブスライスを管理しているディスクノードへ、シングルプライミスライスへの変更を指示するメタデータ変更要求を送信する。スライス割当処理部 5 3 2 は、各ディスクノードからメタデータ変更完了の応答を受け取ると、論理ディスクメタデータ記憶部 5 4 0 内のメタデータを、ディスクノードにおける変更内容と同様に変更する。その後、処理が終了する。

30

【 0 2 1 9 】

[ステップ S 1 6 5] スライス割当処理部 5 3 2 は、変更対象セグメントに割り当てられているプライミスライスとリザーブスライスとに関するメタデータを変更する。具体的には、スライス割当処理部 5 3 2 は、プライミスライスをシングルプライミスライスに変更すると共に、リザーブスライスをフリースライスに変更する。より詳細には、スライス割当処理部 5 3 2 は、プライミスライスを管理しているディスクノードへ、シングルプライミスライスへの変更を指示するメタデータ変更要求を送信する。またスライス割当処理部 5 3 2 は、リザーブスライスを管理しているディスクノードへ、リザーブスライスからフリースライスへの変更を指示するメタデータ変更要求を送信する。スライス割当処理部 5 3 2 は、各ディスクノードからメタデータ変更完了の応答を受け取ると、論理ディスクメタデータ記憶部 5 4 0 内のメタデータを、ディスクノードにおける変更内容と同様に変更する。その後、処理が終了する。

40

【 0 2 2 0 】

図 2 7 は、プライミスライスとリザーブスライスとが割り当てられているセグメントの割り当て変更処理例を示す図である。図 2 7 の例では、変更対象セグメントのプライミスライスをディスクノード 1 0 0 (D P 1) が管理しており、変更対象セグメントのリザーブスライスをディスクノード 2 0 0 (D P 2) が管理しているものとする。またディ

50

スクノード 300 (DP3) には、ライトエラーが発生していないものとする。このような状況において、ディスクノード 100 が管理するプライミスライス、ディスクノード 200 が管理するリザーブスライス、およびディスクノード 300 が管理するフリースライスに着目する。

【0221】

図 27 では、いずれのディスクノードにおいてもライトエラーが検出されていない状態を初期状態としている。変更対象セグメントに割り当てられているスライスを管理しているディスクノード 100, 200 のライトエラーの有無を検査した結果に応じて、メタデータの変更内容が異なる。

【0222】

両方のディスクノード 100, 200 でライトエラーが発生している場合、各ディスクノード 100, 200, 300 の処理は以下の通りである。

ディスクノード 100 が管理しているプライミスライスは変更されず、スライスコピーのコピー元とされる。スライスコピーが完了後、ディスクノード 100 が管理しているプライミスライスは、フリースライス (状態フラグ「F」) に変更される。ディスクノード 200 が管理しているリザーブスライス (状態フラグ「R」) は、最初のメタデータ変更でフリースライス (状態フラグ「F」) に変更される。ディスクノード 300 が管理しているフリースライス (状態フラグ「F」) は、最初のメタデータ変更でリザーブスライス (状態フラグ「R」) に変更され、その後のスライスコピーのコピー先とされる。スライスコピーが完了後、ディスクノード 300 が管理しているリザーブスライスは、シングルプライミスライス (状態フラグ「SP」) に変更される。

【0223】

ディスクノード 100 のみでライトエラーが発生している場合の処理は、両方のディスクノード 100, 200 でライトエラーが発生している場合と同様である。なお、この場合、ディスクノード 200 にはライトエラーは発生していないため、ディスクノード 200 のリザーブスライスを、そのままスライスコピー先とすることもできる。

【0224】

ディスクノード 200 のみでライトエラーが発生している場合、各ディスクノード 100, 200, 300 の処理は以下の通りである。

ディスクノード 100 が管理しているプライミスライス (状態フラグ「P」) は、シングルプライミスライス (状態フラグ「SP」) に変更される。ディスクノード 200 が管理しているリザーブスライス (状態フラグ「R」) は、フリースライス (状態フラグ「F」) に変更される。ディスクノード 300 が管理しているフリースライス (状態フラグ「F」) は変化しない。

【0225】

図 28 は、シングルプライミスライスが割り当てられているセグメントの割り当て変更処理手順の一例を示すフローチャートである。以下、図 28 に示す処理をステップ番号に沿って説明する。

【0226】

[ステップ S171] スライス割当処理部 532 は、変更対象セグメントに割り当てられているプライミスライスに関するメタデータを変更する。具体的には、スライス割当処理部 532 は、シングルプライミスライスをプライミスライスに変更し、リザーブスライスの新規割り当てを行う。より詳細には、スライス割当処理部 532 は、シングルプライミスライスを管理しているディスクノードへ、プライミスライスへの変更を指示するメタデータ変更要求を送信する。またスライス割当処理部 532 は、フリースライスを管理しているディスクノードへ、リザーブスライスへの変更を指示するメタデータ変更要求を送信する。スライス割当処理部 532 は、各ディスクノードからメタデータ変更完了の応答を受け取ると、論理ディスクメタデータ記憶部 540 内のメタデータを、ディスクノードにおける変更内容と同様に変更する。

【0227】

【ステップ S 1 7 2】スライス割当処理部 5 3 2 は、プライマリスライスからリザーブスライスへのスライスのデータのコピーを、プライマリスライスを管理しているディスクノードに対して指示する。

【0 2 2 8】

【ステップ S 1 7 3】スライス割当処理部 5 3 2 は、変更対象セグメントに割り当てられているプライマリスライスとリザーブスライスとに関するメタデータを変更する。具体的には、スライス割当処理部 5 3 2 は、プライマリスライスをフリースライスに変更し、リザーブスライスをシングルプライマリスライスに変更する。より詳細には、スライス割当処理部 5 3 2 は、プライマリスライスを管理しているディスクノードへ、フリースライスへの変更を指示するメタデータ変更要求を送信する。またスライス割当処理部 5 3 2 は、リザーブスライスを管理しているディスクノードへ、シングルプライマリスライスへの変更を指示するメタデータ変更要求を送信する。スライス割当処理部 5 3 2 は、各ディスクノードからメタデータ変更完了の応答を受け取ると、論理ディスクメタデータ記憶部 5 4 0 内のメタデータを、ディスクノードにおける変更内容と同様に変更する。その後、処理が終了する。

10

【0 2 2 9】

図 2 9 は、シングルプライマリスライスが割り当てられているセグメントの割り当て変更処理例を示す図である。図 2 9 の例では、変更対象セグメントのシングルプライマリスライスをディスクノード 1 0 0 (D P 1) が管理しているものとする。またディスクノード 2 0 0 (D P 2) には、ライトエラーが発生していないものとする。このような状況において、ディスクノード 1 0 0 が管理するシングルプライマリスライス、ディスクノード 2 0 0 が管理するフリースライスに着目する。

20

【0 2 3 0】

図 2 9 では、いずれのディスクノードにおいてもライトエラーが検出されていない状態を初期状態としている。シングルプライマリスライスのみが割り当てられているセグメントに対してスライスの割り当て処理が実行されるのは、シングルプライマリスライスを管理しているディスクノード 1 0 0 でライトエラーが発生した場合である。その場合、各ディスクノード 1 0 0 , 2 0 0 において以下の処理が行われる。

【0 2 3 1】

ディスクノード 1 0 0 が管理しているシングルプライマリスライス (状態フラグ「S P」) は、最初のメタデータ変更ではプライマリスライスに変更され、その後、スライスコピーのコピー元とされる。スライスコピーが完了後、ディスクノード 1 0 0 が管理しているプライマリスライスは、フリースライス (状態フラグ「F」) に変更される。ディスクノード 2 0 0 が管理しているフリースライス (状態フラグ「F」) は、最初のメタデータ変更でリザーブスライス (状態フラグ「R」) に変更され、その後のスライスコピーのコピー先とされる。スライスコピーが完了後、ディスクノード 2 0 0 が管理しているリザーブスライスは、シングルプライマリスライス (状態フラグ「S P」) に変更される。

30

【0 2 3 2】

以上のようなメタデータ割当変更処理により、変更対象セグメントには、ライトエラーが発生していないディスクノードが管理するスライスが、シングルプライマリスライスとして割り当てられる。

40

【0 2 3 3】

次に、スライス割当処理のうち、スライスコピーを伴う処理における制御ノード 5 0 0 とディスクノードとの関係動作について説明する。スライスコピーを伴う処理は、ライトエラーが発生したディスクノードが管理するスライスのセグメントへの割り当てを解除すると、セグメントのデータが失われてしまう場合に行われる。

【0 2 3 4】

以下の例では、代表的に、処理対象セグメントに割り当てられているプライマリスライスとセカンダリスライスとのそれぞれを管理するディスクノードの両方でライトエラーが発生している場合の処理について説明する。

50

【 0 2 3 5 】

図 3 0 は、スライスコピーを伴うスライス割当処理手順の一例を示すシーケンス図である。この処理は、例えば図 2 4 のステップ S 1 5 1 からステップ S 1 5 2 に処理が遷移した場合に実行される。以下、図 3 0 に示す処理をステップ番号に沿って説明する。

【 0 2 3 6 】

〔ステップ S 1 8 1 〕スライスコピーを伴うスライス割当処理は、リザーブスライスの割り当て処理（ステップ S 1 8 1 a ）、スライスコピー処理（ステップ S 1 8 1 b ）、およびシングルプライマリ化処理（ステップ S 1 8 1 c ）に分かれる。

【 0 2 3 7 】

〔ステップ S 1 8 1 a 〕スライス割当処理部 5 3 2 は、各ディスクノード 1 0 0 , 2 0 0 , 3 0 0 に対してメタデータ変更要求を送信する。例えば、ディスクノード 1 0 0 へは、変更対象セグメントのプライミスライスのペアとなるセカンダリスライスの情報の削除を指示するメタデータ変更要求が送信される。またディスクノード 2 0 0 へは、変更対象セグメントのセカンダリスライスのフリースライスへの変更を指示するメタデータ変更要求が送信される。ディスクノード 3 0 0 へは、フリースライスを変更対象セグメントのリザーブスライスに変更することを指示するメタデータ変更要求が送信される。

【 0 2 3 8 】

〔ステップ S 1 8 2 〕ディスクノード 1 0 0 のメタデータ管理部 1 6 0 は、メタデータ変更要求に応じてメタデータ記憶部 1 5 0 とストレージ装置 1 1 0 とのメタデータを更新し、メタデータ変更完了の応答を制御ノード 5 0 0 に送信する。

【 0 2 3 9 】

〔ステップ S 1 8 3 〕ディスクノード 2 0 0 のメタデータ管理部 2 6 0 は、メタデータ変更要求に応じてメタデータ記憶部とストレージ装置 2 1 0 とのメタデータを更新し、メタデータ変更完了の応答を制御ノード 5 0 0 に送信する。

【 0 2 4 0 】

〔ステップ S 1 8 4 〕ディスクノード 3 0 0 のメタデータ管理部 3 6 0 は、メタデータ変更要求に応じてメタデータ記憶部とストレージ装置 3 1 0 とのメタデータを更新し、メタデータ更新完了の応答を制御ノード 5 0 0 に送信する。

【 0 2 4 1 】

スライス割当処理部 5 3 2 は、ステップ S 1 8 2 ~ S 1 8 3 のメタデータ更新が完了すると、論理ディスクメタデータ記憶部 5 4 0 内のメタデータを、ディスクノード 1 0 0 , 2 0 0 , 3 0 0 で更新されたメタデータと同じ内容に更新する。

【 0 2 4 2 】

〔ステップ S 1 8 1 b 〕スライス割当処理部 5 3 2 は、ディスクノード 1 0 0 , 2 0 0 , 3 0 0 からメタデータ更新完了の応答を受け取ると、プライミスライスを管理しているディスクノード 1 0 0 に対して、スライスコピー要求を送信する。スライスコピー要求には、変更対象セグメントのプライミスライスのスライス ID、リザーブスライスを管理するディスクノードのディスクノード ID、およびリザーブスライスのスライス ID が含まれる。

【 0 2 4 3 】

〔ステップ S 1 8 5 〕スライスコピー要求を受信したディスクノード 1 0 0 のアクセス処理部 1 2 0 は、指定されたスライスのペアのスライスを管理するディスクノード 3 0 0 に対して、指定されたスライス内のデータを送信する。データを送信する際には、データの格納先として、スライスコピー要求で指定されたスライスのペアのスライスのスライス ID が指定される。

【 0 2 4 4 】

〔ステップ S 1 8 6 〕ディスクノード 3 0 0 のアクセス処理部 3 2 0 は、ディスクノード 1 0 0 から送られたデータを受信し、指定されたスライスにデータを格納する。アクセス処理部 3 2 0 は、データの格納が完了すると、書き込み完了応答をディスクノード 1 0 0 に送信する。

10

20

30

40

50

【0245】

ディスクノード100のアクセス処理部120は、ディスクノード300からの書き込み完了応答を受け取ると、制御ノード500に対してスライスコピー完了の応答を送信する。

【0246】

[ステップS181c] 制御ノード500のスライス割当処理部532は、ディスクノード100からスライスコピー完了の応答を受け取ると、ディスクノード100に対してプライマリスライスのフリースライス化を指示するメタデータ変更要求を送信する。フリースライス化を指示するメタデータ変更要求には、例えば、変更対象セグメントに割り当てられているプライマリスライスのスライスIDが含まれる。

10

【0247】

また、スライス割当処理部532は、ディスクノード100からスライスコピー完了の応答を受け取ると、ディスクノード300に対してリザーブスライスのシングルプライマリ化を指示するメタデータ変更要求を送信する。シングルプライマリ化を指示するメタデータ変更要求には、例えば変更対象セグメントに割り当てられているリザーブスライスのスライスIDが含まれる。

【0248】

[ステップS187] ディスクノード100のメタデータ管理部160は、プライマリスライスのフリースライスへの変更を指示するメタデータ変更要求に応答して、メタデータ記憶部150とストレージ装置110とのメタデータを更新する。メタデータ管理部160は、メタデータの更新が完了すると、更新完了応答を制御ノード500に送信する。

20

【0249】

[ステップS188] ディスクノード300のメタデータ管理部360は、リザーブスライスのシングルプライマリスライスへの変更を指示するメタデータ変更要求に応答して、メタデータ記憶部とストレージ装置110とのメタデータを更新する。メタデータ管理部360は、メタデータの更新が完了すると、更新完了応答を制御ノード500に送信する。

【0250】

なお、スライス割当処理部532は、ステップS187, S188のメタデータ更新が完了すると、論理ディスクメタデータ記憶部540内のメタデータを、ディスクノード100, 200, 300で更新されたメタデータと同じ内容に更新する。

30

【0251】

このような手順によって、スライスコピーを伴うスライス割当処理が完了する。次に、スライスコピーを伴わないスライス割当処理における制御ノード500とディスクノードとの関係動作について説明する。

【0252】

図31は、スライスコピーを伴わないスライス割当処理手順の一例を示すシーケンス図である。この処理は、例えば図24のステップS151からステップS153に処理が遷移した場合に実行される。以下、図31に示す処理をステップ番号に沿って説明する。

【0253】

[ステップS191] スライス割当処理部532は、各ディスクノード100, 200に対してメタデータ変更要求を送信する。例えば、ディスクノード100へは、変更対象セグメントのプライマリスライスのフリースライスへの変更を指示するメタデータ変更要求が送信される。またディスクノード200へは、変更対象セグメントのセカンダリスライスのシングルプライマリスライスへの変更を指示するメタデータ変更要求が送信される。

40

【0254】

[ステップS192] ディスクノード100のメタデータ管理部160は、メタデータ変更要求に応じてメタデータ記憶部150とストレージ装置110とのメタデータを更新し、メタデータ変更完了の応答を制御ノード500に送信する。

50

【 0 2 5 5 】

[ステップ S 1 9 3] ディスクノード 2 0 0 のメタデータ管理部 2 6 0 は、メタデータ変更要求に応じてメタデータ記憶部とストレージ装置 2 1 0 とのメタデータを更新し、メタデータ変更完了の応答を制御ノード 5 0 0 に送信する。

【 0 2 5 6 】

スライス割当処理部 5 3 2 は、ステップ S 1 9 1 , S 1 9 2 のメタデータ更新が完了すると、論理ディスクメタデータ記憶部 5 4 0 内のメタデータを、ディスクノード 1 0 0 , 2 0 0 で更新されたメタデータと同じ内容に更新する。

【 0 2 5 7 】

以上説明したように、リカバリ処理中にディスクノードでライトエラーが発生した場合、そのディスクノードでのライトエラーの発生が、ディスクノードと制御ノードとに記録される。そしてリカバリ処理中に、ライトエラーが発生したディスクノードが管理するスライスにアクセスがあった場合、ライトエラーが発生したディスクノードが管理するスライスのセグメントへの割り当て解除によりセグメントのデータが喪失するか否かが判断される。そして、割り当て解除によりセグメントのデータを喪失する場合には、スライスコピー後にアクセスされたスライスの割り当てが解除される。その結果、データの喪失が抑制される。

10

【 0 2 5 8 】

しかも、ライトエラーが発生したディスクノードが管理するスライスのセグメントへの割り当て解除によりセグメントのデータが喪失しない場合であっても、ペアとなるスライスをシングルプライマリスライスにすることで、リカバリ処理の対象となる。従って、ライトエラーが発生したディスクノードが管理するスライスのうち、アクセスがあったスライスは優先的にリカバリ処理が実行されることとなる。

20

【 0 2 5 9 】

またリカバリ処理中にライトエラーが発生したディスクノードで管理されるスライスにアクセスがあると、アクセスされたスライスのセグメントへの割り当てが解除される。割り当てが解除されたセグメントは、シングルプライマリのみが割り当てられた欠損セグメントに変更される。一方、図 1 6 に示すように、リカバリ処理中はシングルプライマリスライスが順次選択され、二重化が復旧される。従って、欠損セグメントに変更されたセグメントは、リカバリ処理によって二重化状態が復旧される。その結果、マルチノードストレージシステムの信頼性が向上する。

30

【 0 2 6 0 】

ところで第 2 の実施の形態では、ライトエラーが発生したディスクノードのメタデータ割当処理の対象となったスライスは、メタデータ割当処理後にフリースライス(状態フラグ「F」)に変更している。フリースライスであっても、リカバリ処理およびスライス割当処理のいずれにおいても、ライトエラーが発生したディスクノード内のスライスがデータのコピー先(リザーブスライス)として選択されることはない。また、メタデータ割当処理後には、アクセスノード 6 0 0 に新たなメタデータが送信される(図 1 7 参照)。これにより、ライトエラーが発生したディスクノード内のメタデータ割当処理の対象となったスライスが、アクセスノード 6 0 0 からのアクセス対象となることもなくなる。その結果、ライトエラーが発生したディスクノード内のメタデータ割当処理の対象となったスライスは、マルチノードストレージシステム内でのアクセス対象から除外される。すなわち、メタデータ割当処理の対象となったスライスが切り離される。そして、アクセスノード 6 0 0 は、コピー先のスライスに対してアクセスを行うことができる。

40

【 0 2 6 1 】

なおライトエラーが発生したディスクノードのメタデータ割当処理の対象となったスライスは、メタデータ割当処理後に異常スライス(状態フラグ「B」)に変更してもよい。

また第 2 の実施の形態では、各ディスクノードに 1 台ずつのストレージ装置を接続しているため、ストレージ装置が故障するとディスクノードの切り離しが行われている。1 台のディスクノードに複数のストレージ装置が接続されている場合には、ストレージ装置ご

50

とに個別に切り離しを行うことができる。

【0262】

さらにライトエラーが発生したディスクノードにおいて、アクセスノード600からのデータリードのアクセス要求に対しては、要求に従ったストレージ装置へのアクセスを行い、リードしたデータをアクセスノード600に応答するようにしてもよい。すなわち第2の実施の形態では、ライトエラーが発生したディスクノードは、制御ノード500からのスライスコピー要求に対しては、ストレージ装置に対するスライス内のデータリードを行い、コピー先に転送する。その一方で、ライトエラーが発生したディスクノードは、アクセスノード600からのアクセス要求に対しては、リード要求であってもエラー応答を送信している。これは、アクセスがあったセグメントを欠損セグメントに変更し、リカバリ処理による二重化復旧を早期に行うためである。しかし、ライトエラーが発生していてもデータのリードは可能な場合がある。そこで、アクセスノード600への早期のデータ応答を優先する場合、ライトエラーが発生したディスクノードは、アクセス要求のうちリード要求に対してはストレージ装置へのリードアクセスを実行し、読み出したデータを応答してもよい。これにより、アクセスノード600は、スライス割当処理の完了を待たずにアクセス対象のデータを取得可能となる。

10

【0263】

〔その他の応用例〕

上記の処理機能は、コンピュータによって実現することができる。その場合、データ割当制御装置1、制御ノード500、ディスクノード100, 200, 300, 400が有すべき機能の処理内容を記述したプログラムが提供される。そのプログラムをコンピュータで実行することにより、上記処理機能がコンピュータ上で実現される。処理内容を記述したプログラムは、コンピュータで読み取り可能な記録媒体に記録しておくことができる。コンピュータで読み取り可能な記録媒体としては、磁気記憶装置、光ディスク、光磁気記録媒体、半導体メモリなどがある。磁気記憶装置には、ハードディスク装置(HDD)、フレキシブルディスク(FD)、磁気テープなどがある。光ディスクには、DVD、DVD-RAM、CD-ROM/RWなどがある。光磁気記録媒体には、MO(Magneto-Optical disc)などがある。

20

【0264】

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

30

【0265】

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

40

【0266】

また、上記の処理機能の少なくとも一部を、DSP(Digital Signal Processor)、ASIC(Application Specific Integrated Circuit)、PLD(Programmable Logic Device)などの電子回路で実現することもできる。

【0267】

以上、実施の形態を例示したが、実施の形態で示した各部の構成は同様の機能を有する他のものに置換することができる。また、他の任意の構成物や工程が付加されてもよい。さらに、前述した実施の形態のうちの任意の2以上の構成(特徴)を組み合わせたものであってもよい。

50

【0268】

以上の実施の形態に開示された技術には、以下の付記に示す技術が含まれる。

(付記1) 接続されたストレージ装置内のデータを管理する複数のディスクノードに対する管理対象データの割り当て指示を、コンピュータに実行させるデータ割当制御プログラムにおいて、

前記コンピュータに、

複数のストレージ装置に格納されたデータのうち同一内容の冗長データが存在しない二重化欠損データのコピーを、二重化欠損データを管理するディスクノードに指示する二重化復旧処理を行い、

前記二重化復旧処理中に、ストレージ装置に対するデータのライトエラーを示すライトエラー情報を受け取ると、ライトエラーが発生したストレージ装置の識別情報をエラー記憶手段に格納し、

前記エラー記憶手段を参照してライトエラーが発生したストレージ装置を判断し、ライトエラーが発生したストレージ装置に格納されているデータをコピー対象データとし、ライトエラーが発生していないストレージ装置内に前記コピー対象データの冗長データが存在しない場合、前記コピー対象データを管理するディスクノードに対して、ライトエラーが発生していないストレージ装置への前記コピー対象データのコピーを指示する、

処理を実行させることを特徴とするデータ割当制御プログラム。

【0269】

(付記2) 前記二重化復旧処理では、二重化欠損データを順次選択し、選択した二重化欠損データを管理するディスクノードに対し、選択した二重化欠損データが格納されているストレージ装置と異なるストレージ装置であり、かつライトエラーが発生していないストレージ装置への、選択した二重化欠損データのコピーを指示することを特徴とする付記1記載のデータ割当制御プログラム。

【0270】

(付記3) 前記二重化復旧処理では、前記複数のストレージ装置内の二重化欠損データを調査し、二重化欠損データの格納場所を示す管理情報を管理情報記憶手段に格納し、前記管理情報記憶手段内の管理情報に示される二重化欠損データを順次選択して、選択した二重化欠損データを管理するディスクノードに対して選択した二重化欠損データのコピーを指示し、

前記コピー対象データのコピー指示の際には、前記コピー対象データのコピーが完了すると、コピー先のデータの管理情報を二重化欠損データとして前記管理情報記憶手段に格納する、

ことを特徴とする付記1記載のデータ割当制御プログラム。

【0271】

(付記4) 前記二重化復旧処理では、前記管理情報記憶手段内の管理情報で示されるすべての二重化欠損データのコピーが完了すると前記エラー記憶手段を参照してライトエラーが発生したストレージ装置の有無を判断し、ライトエラーが発生したストレージ装置がある場合、ライトエラーが発生したストレージ装置へのアクセスを停止することにより生じる二重化欠損データの管理情報を前記管理情報記憶手段に追加し、二重化復旧処理を続行する、

ことを特徴とする付記3記載のデータ割当制御プログラム。

【0272】

(付記5) 前記コピー対象データのコピー指示では、ライトエラーが発生しているストレージ装置内のデータに対するアクセスで発生したエラーを示すデータアクセスエラー情報を受け取ると、エラー発生時のアクセス対象となっていたデータをコピー対象データとし、前記コピー対象データのコピーが完了すると、コピー先のデータの管理情報を、新たなアクセス先として前記データアクセスエラー情報を送信した装置に応答することを特徴とする付記1記載のデータ割当制御プログラム。

【0273】

10

20

30

40

50

(付記 6) 前記二重化復旧処理中以外の期間に、ストレージ装置に対するデータのライトエラーを示すライトエラー情報を受け取ると、ライトエラーが発生したストレージ装置へのアクセスを停止することにより生じる二重化欠損データの二重化復旧処理を開始することを特徴とする付記 1 記載のデータ割当制御プログラム。

【0274】

(付記 7) 接続されたストレージ装置内のデータを管理する複数のディスクノードに対する管理対象データの割り当て指示をコンピュータで実行するデータ割当制御方法において、

前記コンピュータが、

複数のストレージ装置に格納されたデータのうち同一内容の冗長データが存在しない二重化欠損データのコピーを、二重化欠損データを管理するディスクノードに指示する二重化復旧処理を行い、

前記二重化復旧処理中に、ストレージ装置に対するデータのライトエラーを示すライトエラー情報を受け取ると、ライトエラーが発生したストレージ装置の識別情報をエラー記憶手段に格納し、

前記エラー記憶手段を参照してライトエラーが発生したストレージ装置を判断し、ライトエラーが発生したストレージ装置に格納されているデータをコピー対象データとし、ライトエラーが発生していないストレージ装置内に前記コピー対象データの冗長データが存在しない場合、前記コピー対象データを管理するディスクノードに対して、ライトエラーが発生していないストレージ装置への前記コピー対象データのコピーを指示する、

ことを特徴とするデータ割当制御方法。

【0275】

(付記 8) 接続されたストレージ装置内のデータを管理する複数のディスクノードに対する管理対象データの割り当て指示を行うデータ割当制御装置において、

複数のストレージ装置に格納されたデータのうち同一内容の冗長データが存在しない二重化欠損データのコピーを、二重化欠損データを管理するディスクノードに指示する二重化復旧処理を行う二重化復旧処理手段と、

前記二重化復旧処理中に、ストレージ装置に対するデータのライトエラーを示すライトエラー情報を受け取ると、ライトエラーが発生したストレージ装置の識別情報をエラー記憶手段に格納するエラー情報受信手段と、

前記エラー記憶手段を参照してライトエラーが発生したストレージ装置を判断し、ライトエラーが発生したストレージ装置に格納されているデータをコピー対象データとし、ライトエラーが発生していないストレージ装置内に前記コピー対象データの冗長データが存在しない場合、前記コピー対象データを管理するディスクノードに対して、ライトエラーが発生していないストレージ装置への前記コピー対象データのコピーを指示するコピー指示手段と、

を有することを特徴とするデータ割当制御装置。

【符号の説明】

【0276】

- 1 データ割当制御装置
- 1 a 二重化復旧手段
- 1 b 管理情報記憶手段
- 1 c エラー情報受信手段
- 1 d エラー記憶手段
- 1 e コピー指示手段
- 2 ~ 5 ディスクノード
- 2 a , 3 a , 4 a , 5 a ストレージ装置
- 6 アクセスノード

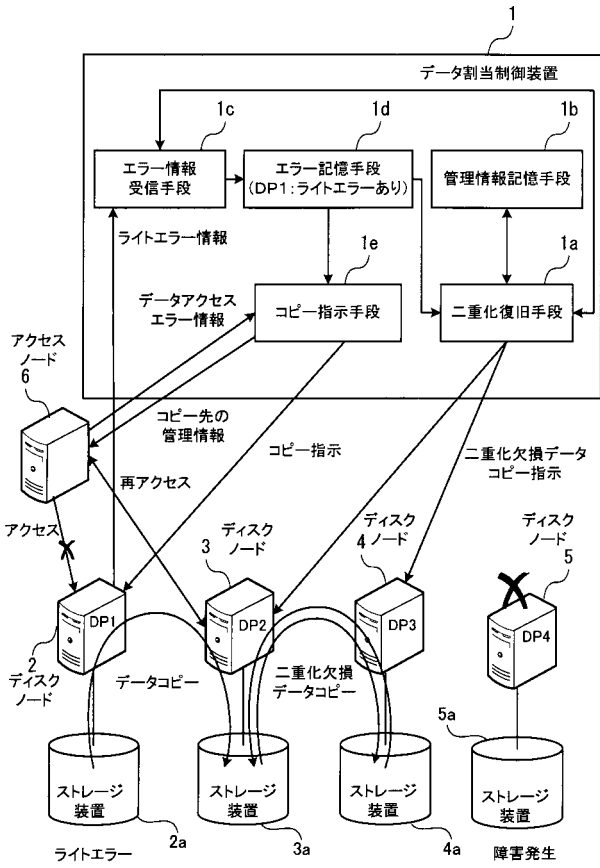
10

20

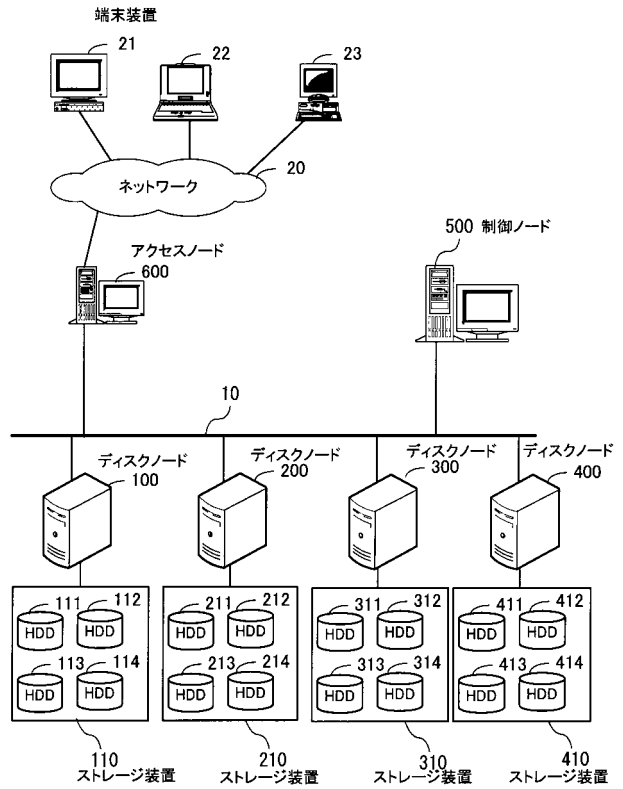
30

40

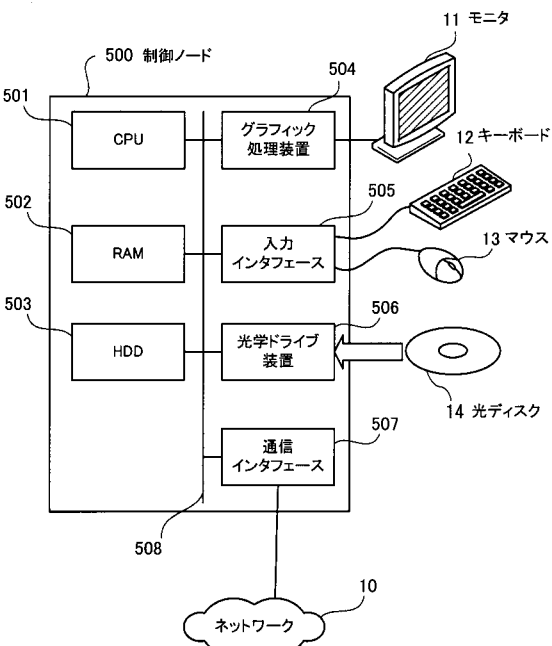
【 図 1 】



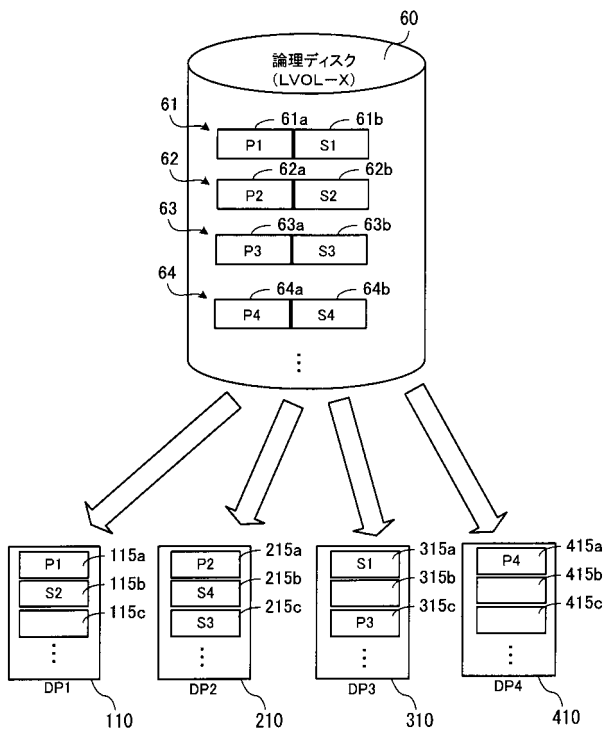
【 図 2 】



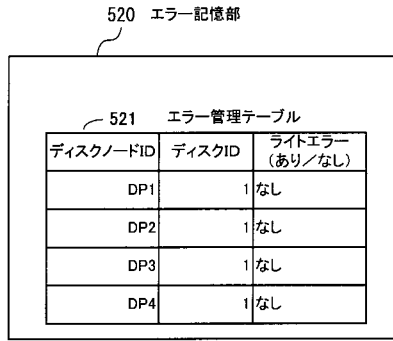
【 図 3 】



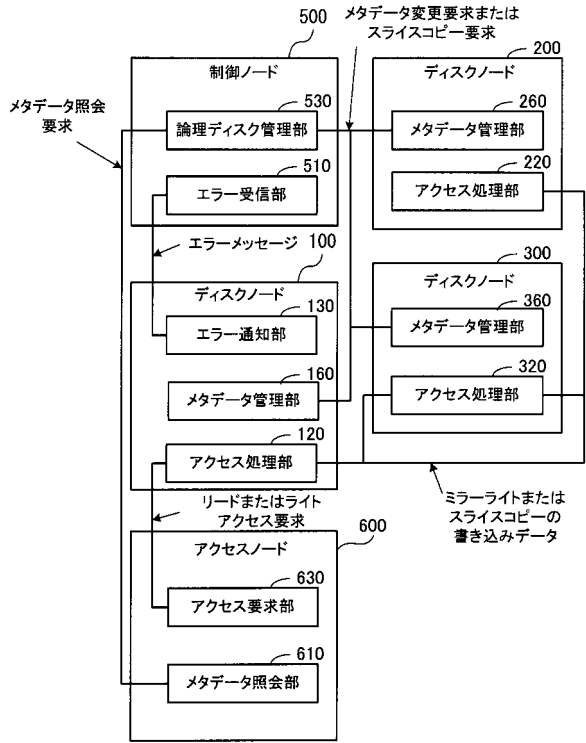
【 図 4 】



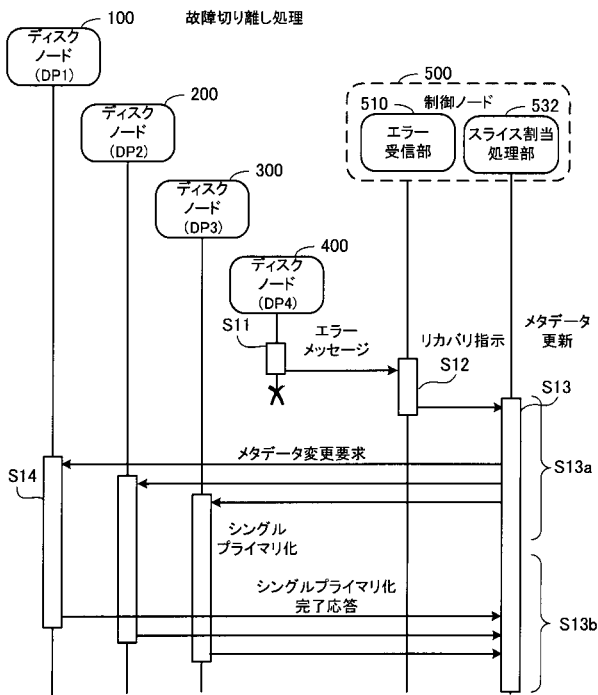
【 図 9 】



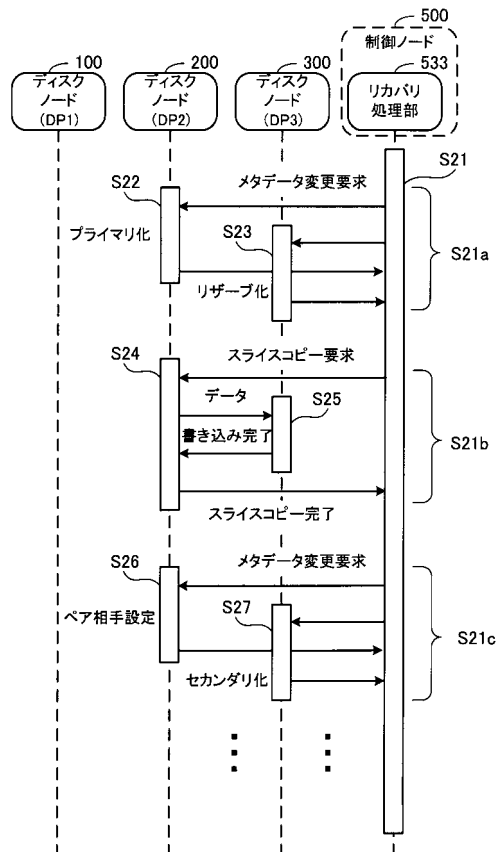
【 図 10 】



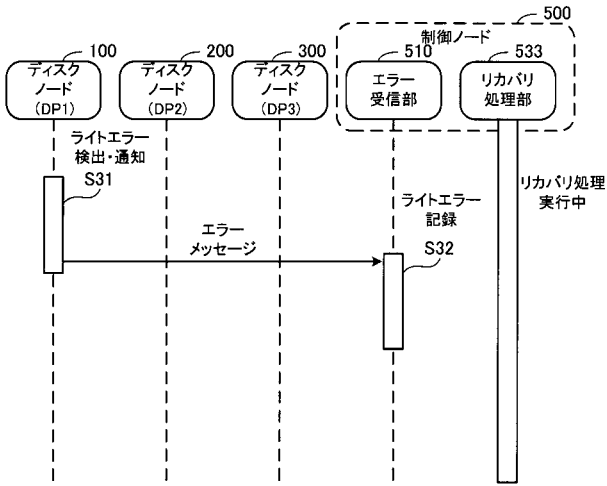
【 図 11 】



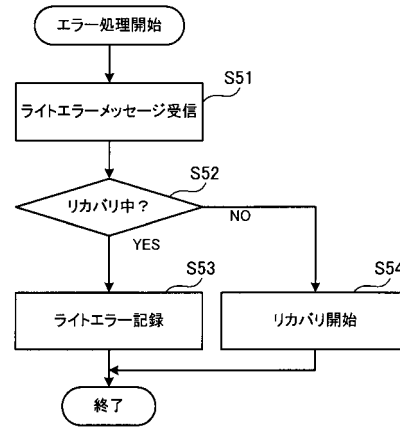
【 図 12 】



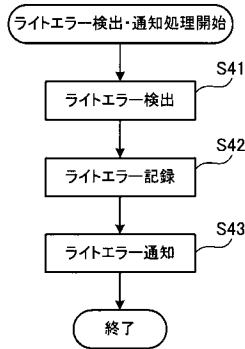
【図13】



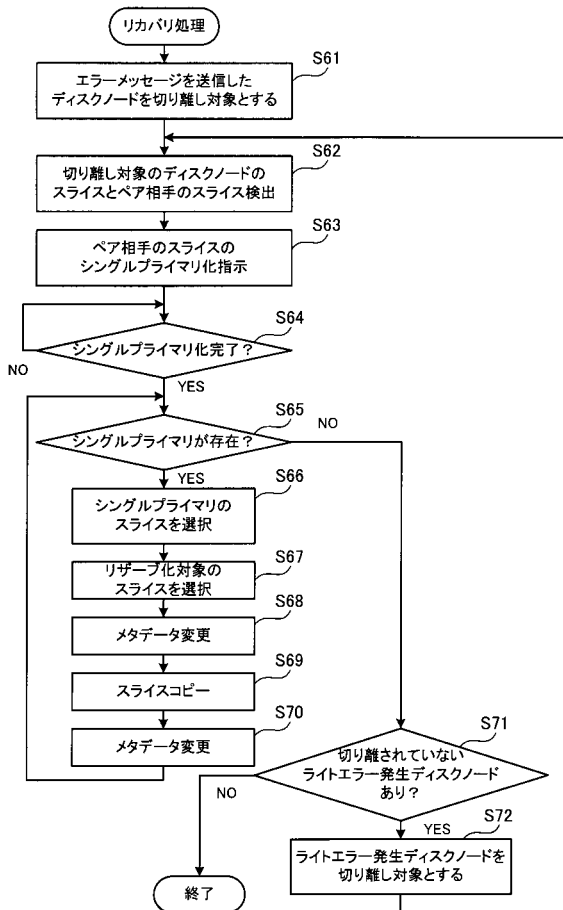
【図15】



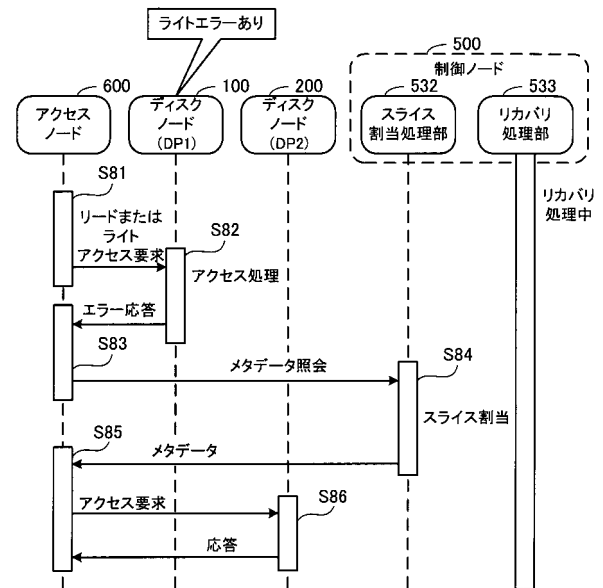
【図14】



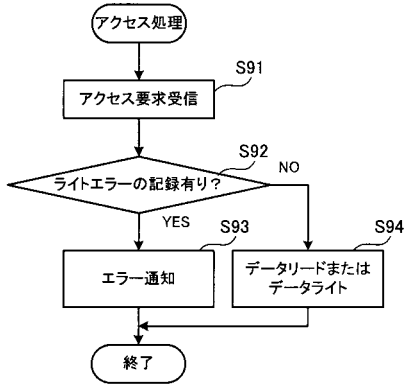
【図16】



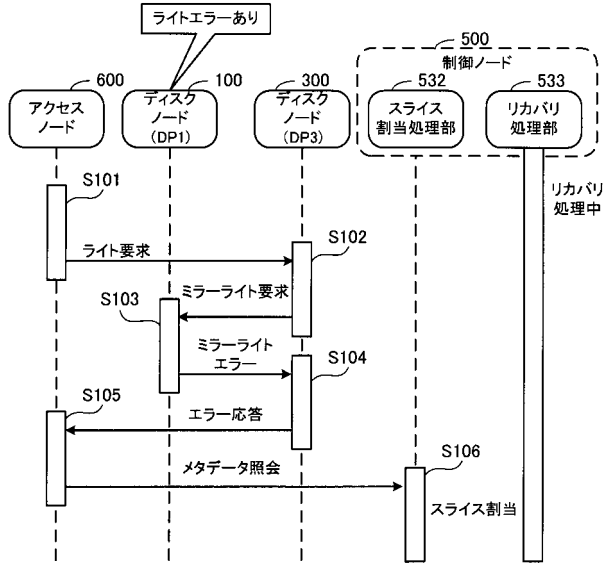
【図17】



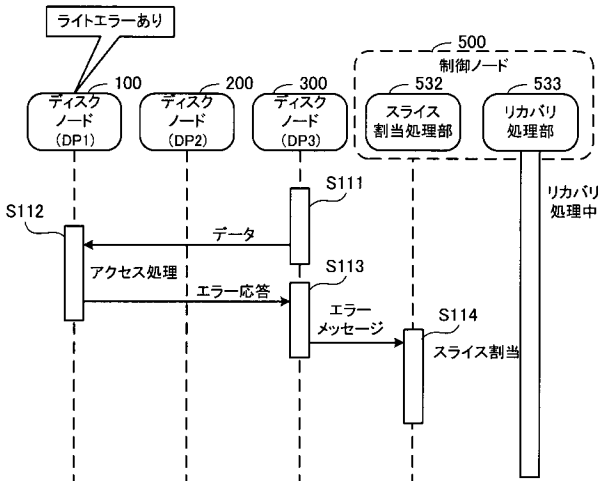
【 図 1 8 】



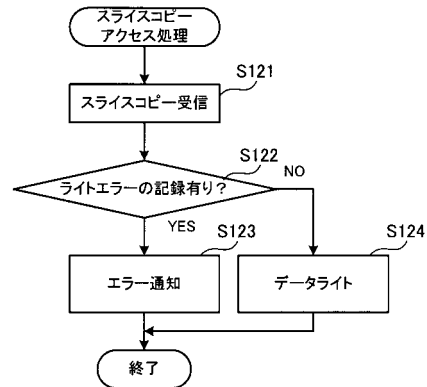
【 図 1 9 】



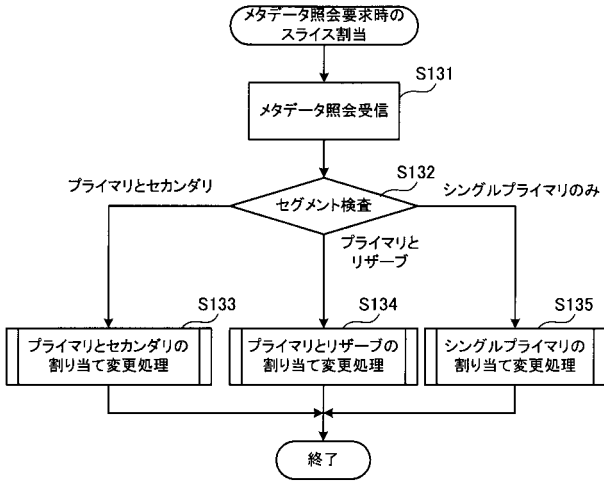
【 図 2 0 】



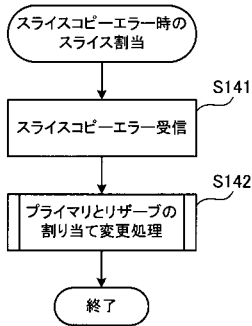
【 図 2 1 】



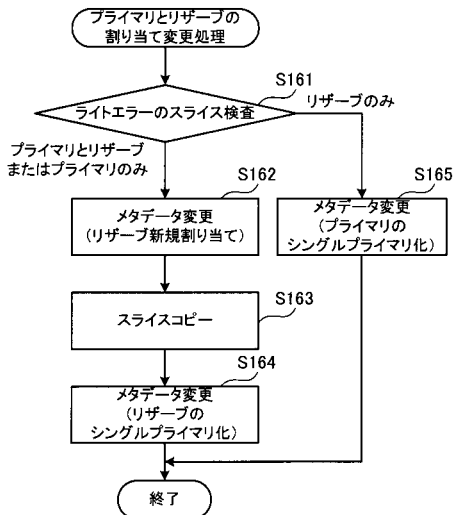
【図 2 2】



【図 2 3】



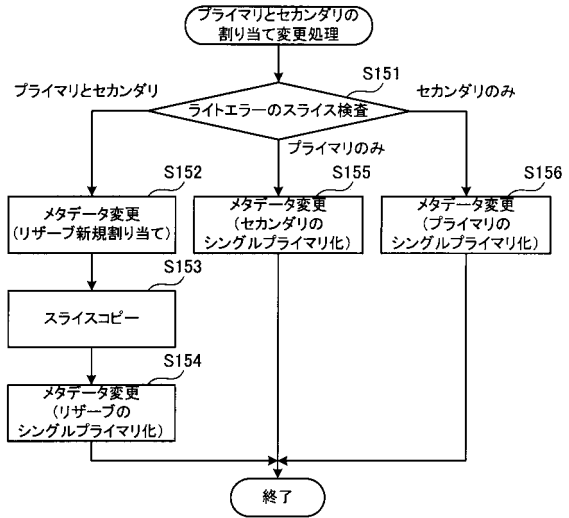
【図 2 6】



【図 2 7】

ライトエラー (DP1, DP2)	DP1	DP2	DP3	シーケンス
初期状態 (なし、なし)	P	R	F	
(あり、あり)	1.P→P 2.コピー元 3.P→F	1.R→F	1.F→R 2.コピー先 3.R→SP	1.メタデータ変更 2.スライスコピー 3.メタデータ変更
(あり、なし)	1.P→P 2.コピー元 3.P→F	1.R→F	1.F→R 2.コピー先 3.R→SP	1.メタデータ変更 2.スライスコピー 3.メタデータ変更
(なし、あり)	1.P→SP	1.R→F	変化なし	1.メタデータ変更

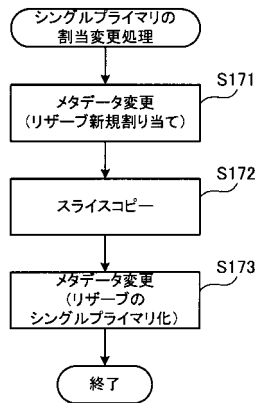
【図 2 4】



【図 2 5】

ライトエラー (DP1, DP2)	DP1	DP2	DP3	シーケンス
初期状態 (なし、なし)	P	S	F	
(あり、あり)	1.P→P 2.コピー元 3.P→F	1.S→F	1.F→R 2.コピー先 3.R→SP	1.メタデータ変更 2.スライスコピー 3.メタデータ変更
(あり、なし)	1.P→F	1.S→SP	変化なし	1.メタデータ変更
(なし、あり)	1.P→SP	1.S→F	変化なし	1.メタデータ変更

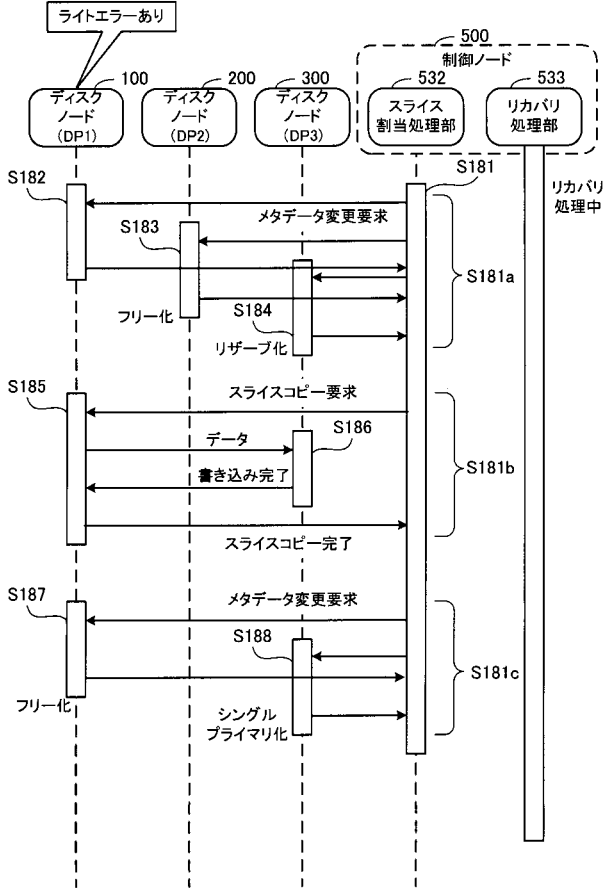
【図 2 8】



【図 2 9】

ライトエラー (DP1, DP2)	DP1	DP2	シーケンス
初期状態 (なし、なし)	SP	F	
(あり、なし)	1.SP→P 2.コピー元 3.P→F	1.F→R 2.コピー先 3.R→SP	1.メタデータ変更 2.スライスコピー 3.メタデータ変更

【図30】



【図31】

