

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4838878号  
(P4838878)

(45) 発行日 平成23年12月14日(2011.12.14)

(24) 登録日 平成23年10月7日(2011.10.7)

(51) Int.Cl. F I  
**G06F 12/00 (2006.01)**  
 G06F 12/00 501B  
 G06F 12/00 545A  
 G06F 12/00 520E

請求項の数 9 (全 54 頁)

(21) 出願番号	特願2009-276025 (P2009-276025)	(73) 特許権者	000005223 富士通株式会社
(22) 出願日	平成21年12月4日(2009.12.4)		神奈川県川崎市中原区上小田中4丁目1番1号
(65) 公開番号	特開2011-118712 (P2011-118712A)	(74) 代理人	100092152 弁理士 服部 毅巖
(43) 公開日	平成23年6月16日(2011.6.16)	(72) 発明者	田村 雅寿 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
審査請求日	平成23年2月17日(2011.2.17)	(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	荻原 一隆 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

最終頁に続く

(54) 【発明の名称】 データ管理プログラム、データ管理装置、およびデータ管理方法

(57) 【特許請求の範囲】

【請求項1】

複数のディスクノードそれぞれで管理されるストレージ装置にデータを分散格納するマルチノードストレージシステムを構成する1つの前記ディスクノードで実行すべきデータ管理処理をコンピュータに実行させるデータ管理プログラムにおいて、

前記コンピュータに、

仮想的な記憶領域を定義した論理ボリュームを複数に分割して得られるデータユニットを指定した書き込み要求があると、前記コンピュータに接続されたストレージ装置内の単位記憶領域の1つを指定された前記データユニットに割り当て、前記データユニットに割り当てられた前記単位記憶領域に対してデータの書き込みを行い、

書き込み対象となった前記データユニットに対応付けて、現在の時刻を最終書き込み時刻としてデータユニット情報記憶手段に格納し、

前記データユニット情報記憶手段を参照し、前記最終書き込み時刻から重複排除化猶予期間以上経過している前記データユニットを検出し、

仮想的な記憶領域を定義した重複排除ボリュームを複数に分割して得られる重複排除ユニットのうち重複排除化の対象とされた重複排除対象データの格納に使用されている重複排除ユニットの識別情報および前記重複排除ユニットを管理している前記ディスクノードの識別子を含む重複排除アドレスと、前記重複排除対象データに所定の演算を行うことで得られる前記重複排除対象データの固有値とを対応付けて管理するインデックスサーバから、検出された前記データユニットに割り当てられた前記単位記憶領域内のデータの

固有値に対応付けられた前記重複排除アドレスを取得し、

検出された前記データユニットに対応付けて、取得した前記重複排除アドレスを前記データユニット情報記憶手段に格納すると共に、検出された前記データユニットへの前記単位記憶領域の割り当てを解除する、

処理を実行させることを特徴とするデータ管理プログラム。

【請求項 2】

前記コンピュータに、さらに、

前記インデックスサーバから、前記重複排除ユニットの識別情報を指定した重複排除ユニット割当要求を受信すると、指定された前記重複排除ユニットに前記単位記憶領域の1つを割り当て、指定された前記重複排除ユニットをデータの読み出し先とする前記ディスクノードから前記重複排除対象データを取得し、指定された前記重複排除ユニットに割り当てた前記単位記憶領域に格納し、

前記重複排除ユニットの識別情報を指定したデータの読み出し要求を受信すると、指定された前記重複排除ユニットに割り当てられた前記単位記憶領域内の前記重複排除対象データを応答する、

処理を実行させることを特徴とする請求項 1 記載のデータ管理プログラム。

【請求項 3】

前記コンピュータに、さらに、

前記重複排除ユニットを前記重複排除対象データの読み出し先とする前記ディスクノードが存在しなくなった場合、前記重複排除ユニットへの前記単位記憶領域の割り当てを解除する、

処理を実行させることを特徴とする請求項 2 記載のデータ管理プログラム。

【請求項 4】

前記コンピュータにさらに、

前記重複排除ユニットに対応付けて、前記重複排除ユニットに割り当てられた前記単位記憶領域から最後にデータの読み出しが行われた時刻から保存期間経過後の保存期限満了時刻を重複排除ユニット情報記憶手段に格納し、

前記重複排除ユニット情報記憶手段を参照し、現在の時刻が前記保存期限満了時刻を過ぎている前記重複排除ユニットを検出し、検出した前記重複排除ユニットへの前記単位記憶領域の割り当てを解除する、

処理を実行させることを特徴とする請求項 3 記載のデータ管理プログラム。

【請求項 5】

前記コンピュータに、さらに、

前記データユニットを指定した読み出し要求があると、読み出し対象の前記データユニットに割り当てられた前記単位記憶領域がない場合、前記データユニット情報記憶手段を参照し、読み出し対象の前記データユニットに対応付けられた前記重複排除アドレスに示される前記ディスクノードから、前記重複排除アドレスに示される前記重複排除ユニットの識別情報を指定して前記重複排除対象データを取得する、

処理を実行させることを特徴とする請求項 1 記載のデータ管理プログラム。

【請求項 6】

複数のディスクノードそれぞれで管理されるストレージ装置にデータを分散格納するマルチノードストレージシステムの記憶領域管理処理をコンピュータに実行させるデータ管理プログラムにおいて、

前記コンピュータに、

仮想的な記憶領域を定義した重複排除ボリュームを複数に分割して得られる重複排除ユニットのうち前記ディスクノードで管理している前記重複排除ユニットの使用の有無、使用されている前記重複排除ユニットに割り当てられている単位記憶領域に格納された重複排除対象データに所定の演算を行うことで得られる固有値、および前記重複排除ユニットの識別情報を含む重複排除ユニット情報が、前記重複排除ユニットを管理している前記ディスクノードの識別子に対応付けて重複排除ユニット情報記憶手段に格納されており、

10

20

30

40

50

前記ディスクノードから、重複排除化データユニットに割り当てられた単位記憶領域内のデータに前記所定の演算を実行して得られた固有値を含む重複排除アドレス照会要求を受け取ると、前記重複排除ユニット情報記憶手段から、前記重複排除アドレス照会要求に示される固有値を含む前記重複排除ユニット情報を検索し、

検索により該当する前記重複排除ユニット情報が検出された場合、検出された前記重複排除ユニット情報に含まれる前記重複排除ユニットの識別情報と、検出された前記重複排除ユニット情報に対応する前記重複排除ユニットを管理している前記ディスクノードの識別子とを含む重複排除アドレスを、前記重複排除アドレス照会要求の送信元である前記ディスクノードに応答し、

検索により該当する前記重複排除ユニット情報が検出されなかった場合、前記重複排除ユニット情報記憶手段を参照し、未使用の前記重複排除ユニットを選択し、選択した前記重複排除ユニットを管理する前記ディスクノードに対して選択した重複排除ユニットへの前記単位記憶領域の割り当て要求を送信し、選択した前記重複排除ユニットの前記単位記憶領域割り当て後の前記重複排除ユニット情報を前記重複排除ユニット情報記憶手段に格納し、選択した前記重複排除ユニットの識別情報と、選択した前記重複排除ユニットを管理する前記ディスクノードの識別子とを含む重複排除アドレスを、前記重複排除アドレス照会要求の送信元である前記ディスクノードに応答する、

処理を実行させることを特徴とするデータ管理プログラム。

#### 【請求項 7】

前記コンピュータに、さらに、

前記重複排除ユニット情報には、前記重複排除ユニットに割り当てられた前記単位記憶領域から最後にデータの読み出しが行われた時刻から保存期間経過後の保存期限満了時刻が含まれており、前記重複排除ユニット情報記憶手段を参照し、現在の時刻が前記保存期限満了時刻を過ぎている前記重複排除ユニットを検出し、

検出された前記重複排除ユニット情報の取得元の前記ディスクノードから最新の前記重複排除ユニット情報を取得し、前記重複排除ユニット情報記憶手段に反映する、

処理を実行させることを特徴とする請求項 6 記載のデータ管理プログラム。

#### 【請求項 8】

複数のディスクノードそれぞれで管理されるストレージ装置にデータを分散格納するマルチノードストレージシステムを構成する 1 つの前記ディスクノードで実行すべきデータ管理処理をコンピュータによって実現するデータ管理装置において、

仮想的な記憶領域を定義した論理ボリュームを複数に分割して得られるデータユニットを指定した書き込み要求があると、前記コンピュータに接続されたストレージ装置内の単位記憶領域の 1 つを指定された前記データユニットに割り当て、前記データユニットに割り当てられた前記単位記憶領域に対してデータの書き込みを行う書き込みアクセス手段と、

書き込み対象となった前記データユニットに対応付けて、現在の時刻を最終書き込み時刻としてデータユニット情報記憶手段に格納する最終書き込み時刻更新手段と、

前記データユニット情報記憶手段を参照し、最終書き込み時刻から重複排除化猶予期間以上経過している前記データユニットを検出する重複排除化データユニット検出手段と、

仮想的な記憶領域を定義した重複排除ボリュームを複数に分割して得られる重複排除ユニットのうち重複排除化の対象とされた重複排除対象データの格納に使用されている重複排除ユニットの識別情報および前記重複排除ユニットを管理している前記ディスクノードの識別子を含む重複排除アドレスと、前記重複排除対象データに所定の演算を行うことで得られる前記重複排除対象データの固有値とを対応付けて管理するインデックスサーバから、検出された前記データユニットに割り当てられた前記単位記憶領域内のデータの固有値に対応付けられた前記重複排除アドレスを取得する重複排除アドレス取得手段と、

検出された前記データユニットに対応付けて、取得した前記重複排除アドレスを前記データユニット情報記憶手段に格納すると共に、検出された前記データユニットへの前記単位記憶領域の割り当てを解除する単位記憶領域割当解除手段と、

10

20

30

40

50

前記データユニットを指定した読み出し要求があると、読み出し対象の前記データユニットに割り当てられた前記単位記憶領域がある場合、割り当てられた前記単位記憶領域からデータを読み出し、読み出し対象の前記データユニットに割り当てられた前記単位記憶領域がない場合、前記データユニット情報記憶手段を参照し、読み出し対象の前記データユニットに対応付けられた前記重複排除アドレスに示される前記ディスクノードから、前記重複排除アドレスに示される前記重複排除ユニットの識別情報を指定して前記重複排除対象データを取得する読み出しアクセス手段と、

を有することを特徴とするデータ管理装置。

【請求項9】

複数のディスクノードそれぞれで管理されるストレージ装置にデータを分散格納するマルチノードストレージシステムを構成する1つの前記ディスクノードで実行すべきデータ管理処理をコンピュータで実行するデータ管理方法において、

前記コンピュータが、

仮想的な記憶領域を定義した論理ボリュームを複数に分割して得られるデータユニットを指定した書き込み要求があると、前記コンピュータに接続されたストレージ装置内の単位記憶領域の1つを指定された前記データユニットに割り当て、前記データユニットに割り当てられた前記単位記憶領域に対してデータの書き込みを行い、

書き込み対象となった前記データユニットに対応付けて、現在の時刻を最終書き込み時刻としてデータユニット情報記憶手段に格納し、

前記データユニット情報記憶手段を参照し、最終書き込み時刻から重複排除化猶予期間以上経過している前記データユニットを検出し、

仮想的な記憶領域を定義した重複排除ボリュームを複数に分割して得られる重複排除ユニットのうち重複排除化の対象とされた重複排除対象データの格納に使用されている重複排除ユニットの識別情報および前記重複排除ユニットを管理している前記ディスクノードの識別子を含む重複排除アドレスと、前記重複排除対象データに所定の演算を行うことで得られる前記重複排除対象データの固有値とを対応付けて管理するインデックスサーバから、検出された前記データユニットに割り当てられた前記単位記憶領域内のデータの前記固有値に対応付けられた前記重複排除アドレスを取得し、

検出された前記データユニットに対応付けて、取得した前記重複排除アドレスを前記データユニット情報記憶手段に格納すると共に、検出された前記データユニットへの前記単位記憶領域の割り当てを解除し、

前記データユニットを指定した読み出し要求があると、読み出し対象の前記データユニットに割り当てられた前記単位記憶領域がある場合、割り当てられた前記単位記憶領域からデータを読み出し、読み出し対象の前記データユニットに割り当てられた前記単位記憶領域がない場合、前記データユニット情報記憶手段を参照し、読み出し対象の前記データユニットに対応付けられた前記重複排除アドレスに示される前記ディスクノードから、前記重複排除アドレスに示される前記重複排除ユニットの識別情報を指定して前記重複排除対象データを取得する、

ことを特徴とするデータ管理方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明はデータ管理プログラム、データ管理装置、およびデータ管理方法に関する。

【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を有している。マルチノードストレージシステムでは、制御ノードの管理の下、仮想的なディスク（論理ボリューム）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。

10

20

30

40

50

## 【0003】

具体的には、マルチノードストレージシステムでは、例えば論理ボリューム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装置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、固定長であり、セグメントと同サイズとなる。制御ノードによって、論理ボリュームの各セグメントに対して、ストレージ装置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されたデータが、そのセグメントに割り当てられたスライスを管理するディスクノードに送られ、そのディスクノードが有するストレージ装置に格納される。

10

## 【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】

なお、コンピュータシステム内で同じ内容のデータを複数持つことは、ストレージの使用効率を悪化させることとなる。例えば、データのバックアップを定期的に行う場合、前回のバックアップ内容と比較してユニークなデータは、バックアップ対象の一部であることが多い。そこで、バックアップなどで移動させるデータの冗長性を低減する技術が考えられている。

20

## 【先行技術文献】

## 【特許文献】

## 【0006】

【特許文献1】国際公開第2004/104845号

【特許文献2】特開2007-234026号公報

## 【発明の概要】

## 【発明が解決しようとする課題】

## 【0007】

ところで、コンピュータが保持するデータを細分化して見てみると、運用中のシステム内でも同じ内容のデータが複数同時に存在する可能性がある。例えば同じファイルを添付した電子メールが、メールサーバを共有する複数のユーザに送信されると、そのメールサーバにおいて同じ内容のデータが異なる場所に格納される。

30

## 【0008】

特にマルチノードストレージシステムでは、論理ボリュームごとに異なるユーザにサービスを提供する場合がある。この場合、複数の論理ボリュームそれぞれに同じアプリケーションプログラムがインストールされることもあり得る。同じアプリケーションプログラムが異なる論理ボリュームにインストールされると、マルチノードストレージシステム全体としては、同じ内容のデータが複数の場所に格納されることとなる。

## 【0009】

しかし、従来は、マルチノードストレージシステムの異なるディスクノードに同じ内容のデータが重複して存在する場合、それらのデータの冗長性を低減させることができなかった。そのため、異なるディスクノードに同じ内容のデータが複数格納され、ストレージが効率的に利用されていなかった。

40

## 【0010】

本発明はこのように鑑みてなされたものであり、マルチノードストレージシステムにおけるデータの冗長性を低減することができるデータ管理プログラム、データ管理装置、およびデータ管理方法を提供することを目的とする。

## 【課題を解決するための手段】

## 【0011】

上記課題を解決するために、複数のディスクノードそれぞれで管理されるストレージ装

50

置にデータを分散格納するマルチノードストレージシステムを構成する1つのディスクノードで実行すべきデータ管理処理をコンピュータに実行させるデータ管理プログラムが提供される。このデータ管理プログラムは、コンピュータに以下の処理を実行させる。

【0012】

コンピュータは、仮想的な記憶領域を定義した論理ボリュームを複数に分割して得られるデータユニットを指定した書き込み要求があると、コンピュータに接続されたストレージ装置内の単位記憶領域の1つを指定されたデータユニットに割り当て、データユニットに割り当てられた単位記憶領域に対してデータの書き込みを行う。コンピュータは、書き込み対象となったデータユニットに対応付けて、現在の時刻を最終書き込み時刻としてデータユニット情報記憶手段に格納する。コンピュータは、データユニット情報記憶手段を参照し、最終書き込み時刻から重複排除化猶予期間以上経過しているデータユニットを検出する。コンピュータは、仮想的な記憶領域を定義した重複排除ボリュームを複数に分割して得られる重複排除ユニットのうち重複排除化の対象とされた重複排除対象データの格納に使用されている重複排除ユニットの識別情報および重複排除ユニットを管理しているディスクノードの識別子を含む重複排除アドレスと、重複排除対象データに所定の演算を行うことで得られる重複排除対象データの固有値とを対応付けて管理するインデックスサーバから、検出されたデータユニットに割り当てられた単位記憶領域内のデータの固有値に対応付けられた重複排除アドレスを取得する。コンピュータは、検出されたデータユニットに対応付けて、取得した重複排除アドレスをデータユニット情報記憶手段に格納すると共に、検出されたデータユニットへの単位記憶領域の割り当てを解除する。

10

20

【0013】

また、上記課題を解決するために、上記データ管理プログラムを実行するコンピュータと同様の機能を有するデータ管理装置が提供される。また、上記課題を解決するために、上記データ管理プログラムに基づいてコンピュータに実行する処理によるデータ管理方法が提供される。

【0014】

さらに、上記課題を解決するために、複数のディスクノードそれぞれで管理されるストレージ装置にデータを分散格納するマルチノードストレージシステムの記憶領域管理処理をコンピュータに実行させるデータ管理プログラムが提供される。このデータ管理プログラムは、コンピュータに以下の処理を実行させる。

30

【0015】

仮想的な記憶領域を定義した重複排除ボリュームを複数に分割して得られる重複排除ユニットのうちディスクノードで管理している重複排除ユニットの使用の有無、使用されている重複排除ユニットに割り当てられている単位記憶領域に格納された重複排除対象データに所定の演算を行うことで得られる固有値、および重複排除ユニットの識別情報を含む重複排除ユニット情報が、重複排除ユニットを管理しているディスクノードの識別子に対応付けて重複排除ユニット情報記憶手段に格納されている。コンピュータは、ディスクノードから、重複排除化データユニットに割り当てられた単位記憶領域内のデータに所定の演算を実行して得られた固有値を含む重複排除アドレス照会要求を受け取ると、重複排除ユニット情報記憶手段から、重複排除アドレス照会要求に示される固有値を含む重複排除ユニット情報を検索する。コンピュータは、検索により該当する重複排除ユニット情報が検出された場合、検出された重複排除ユニット情報に含まれる重複排除ユニットの識別情報と、検出された重複排除ユニット情報に対応する重複排除ユニットを管理しているディスクノードの識別子とを含む重複排除アドレスを、重複排除アドレス照会要求の送信元であるディスクノードに回答する。コンピュータは、検索により該当する重複排除ユニット情報が検出されなかった場合、重複排除ユニット情報記憶手段を参照し、未使用の重複排除ユニットを選択し、選択した重複排除ユニットを管理するディスクノードに対して選択した重複排除ユニットへの単位記憶領域の割り当て要求を送信し、選択した重複排除ユニットの単位記憶領域割り当て後の重複排除ユニット情報を重複排除ユニット情報記憶手段に格納し、選択した前記重複排除ユニットの識別情報と、選択した前記重複排除ユニット

40

50

を管理する前記ディスクノードの識別子とを含む重複排除アドレスを、重複排除アドレス照会要求の送信元であるディスクノードに応答する。

【発明の効果】

【0016】

マルチノードストレージシステムにおけるディスクノード間のデータの冗長性を低減することができる。

【図面の簡単な説明】

【0017】

【図1】第1の実施の形態のシステム構成を示すブロック図である。

【図2】第2の実施の形態のシステム構成例を示す図である。

10

【図3】本実施の形態に用いるディスクノードのハードウェア構成例を示す図である。

【図4】論理ボリュームと重複排除ボリュームへの記憶領域の割り当て関係を示す図である。

【図5】ストレージ装置に格納されている情報の例を示す図である。

【図6】論理ボリューム用スライスメタデータのデータ構造例を示す図である。

【図7】重複排除ボリューム用スライスメタデータのデータ構造例を示す図である。

【図8】制御ノードとアクセスノードとの機能を示すブロック図である。

【図9】制御ノードのスライスメタデータ記憶部のデータ構造例を示す図である。

【図10】制御ノードの論理ボリューム構成情報記憶部のデータ構造例を示す図である。

【図11】制御ノードの重複排除ボリューム構成情報記憶部のデータ構造例を示す図である。

20

【図12】アクセスノードの論理ボリューム構成情報記憶部のデータ構造例を示す図である。

【図13】ディスクノードとインデックスサーバとの機能を示すブロック図である。

【図14】ディスクノードが有する重複排除ボリューム構成情報記憶部のデータ構造例を示す図である。

【図15】重複排除ユニット情報収集処理を示すシーケンス図である。

【図16】重複排除ユニット情報記憶部のデータ構造例を示す図である。

【図17】重複排除ボリュームからのデータ読み出し処理の手順を示すシーケンス図である。

30

【図18】状態が「Blank」のデータユニットへの書き込み処理を示すシーケンス図である。

【図19】状態が「DeDup」のデータユニットの一部への書き込み処理を示すシーケンス図である。

【図20】ディスクノードにおける書き込み処理を示すフローチャートである。

【図21】状態が「Normal」のデータユニットへのパトロール処理を示すシーケンス図である。

【図22】状態が「DeDup」のデータユニットへのパトロール処理を示すシーケンス図である。

【図23】パトロール処理の手順を示すフローチャートである。

40

【図24】状態が「Normal」のデータユニットのパトロール処理の手順を示すフローチャートである。

【図25】重複排除ボリュームのパトロール処理の手順を示すフローチャートである。

【図26】インデックスサーバにおける重複排除アドレス検索処理の手順を示すフローチャートである。

【図27】未使用ユニット情報反映処理の手順を示すフローチャートである。

【図28】パトロールによるデータの保存先変更状況を示す図である。

【図29】重複排除ユニットへの関連付けの解消状況を示す図である。

【図30】重複排除ユニットの未使用状態への変更状況を示す図である。

【発明を実施するための形態】

50

## 【 0 0 1 8 】

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 〕

第 1 の実施の形態に係るマルチノードストレージシステムにおいて、データ実体を特別な仮想ボリューム（重複排除ボリューム）に保持し、ユーザの論理ボリュームには、そのデータ実体を指し示すインデックスを保持させる。このような処理を、重複排除（DeDup: Deduplication）処理と呼ぶ。これにより、同じ内容のデータを重複して保持することが抑止され、ストレージ容量の使用効率を向上する。

## 【 0 0 1 9 】

図 1 は、第 1 の実施の形態のシステム構成を示すブロック図である。第 1 の実施の形態にかかるとマルチノードストレージシステムでは、複数のディスクノード 1 ~ 3 それぞれで管理されるストレージ装置 5 ~ 7 にデータを分散格納する。このマルチノードストレージシステムでは、インデックスサーバ 4 によって、重複排除処理における重複排除化されたデータの格納場所が管理される。

10

## 【 0 0 2 0 】

ディスクノード 1 は、書き込みアクセス手段 1 a、最終書き込み時刻更新手段 1 b、データユニット情報記憶手段 1 c、重複排除化データユニット検出手段 1 d、重複排除アドレス取得手段 1 e、単位記憶領域割当解除手段 1 f、読み出しアクセス手段 1 g、重複排除ユニット割当手段 1 h、および重複排除データアクセス応答手段 1 i を有している。

## 【 0 0 2 1 】

書き込みアクセス手段 1 a は、仮想的な記憶領域を定義した論理ボリュームを複数に分割して得られるデータユニットを指定した書き込み要求を、図示していないコンピュータから取得する。すると書き込みアクセス手段 1 a は、ディスクノード 1 に接続されたストレージ装置 5 内の単位記憶領域の 1 つ（図 1 の例では単位記憶領域 5 a）を指定されたデータユニットに割り当てる。さらに書き込みアクセス手段 1 a は、データユニットに割り当てられた単位記憶領域 5 a に対してデータの書き込みを行う。

20

## 【 0 0 2 2 】

最終書き込み時刻更新手段 1 b は、書き込み対象となったデータユニットに対応付けて、現在の時刻を最終書き込み時刻としてデータユニット情報記憶手段 1 c に格納する。

データユニット情報記憶手段 1 c は、論理ボリュームに含まれるデータユニットごとに、最終書き込み時刻、または重複排除アドレスを記憶する記憶機能である。例えばデータユニット情報記憶手段 1 c は、ディスクノード 1 内のメモリの記憶領域の一部が使用される。また、ストレージ装置 5 の記憶領域の一部を、データユニット情報記憶手段 1 c として使用することもできる。

30

## 【 0 0 2 3 】

ここで重複排除アドレスとは、重複排除ユニットのうち重複排除化の対象とされた重複排除対象データの格納に使用されている重複排除ユニットの識別情報、および重複排除ユニットを管理しているディスクノードの識別子を含む情報である。重複排除ユニットは、仮想的な記憶領域を定義した重複排除ボリュームを複数に分割して得られる仮想的な記憶領域である。

40

## 【 0 0 2 4 】

重複排除化データユニット検出手段 1 d は、データユニット情報記憶手段 1 c を参照し、最終書き込み時刻から重複排除化猶予期間以上経過しているデータユニット（重複排除化データユニット）を検出する。重複排除化猶予期間は、重複排除化データユニット検出手段 1 d に予め設定されている。検出されたデータユニットに割り当てられた単位記憶領域 5 a に格納されたデータは、重複排除化の対象となる。

## 【 0 0 2 5 】

重複排除アドレス取得手段 1 e は、インデックスサーバ 4 から、検出されたデータユニットに割り当てられた単位記憶領域 5 a 内のデータの固有値に対応付けられた重複排除アドレスを取得する。ここで、データの固有値とは、データに所定の演算を行うことで得ら

50



れる値である。この場合の演算は、入力されるデータが異なれば異なる固有値を生成するような演算である。例えばデータに対してハッシュ関数による演算を行うことで得られるハッシュ値が、固有値として使用できる。この場合のハッシュ関数としては、例えば、入力されるデータが異なれば異なるハッシュ値を生成するようなハッシュ関数が用いられる。

【0026】

単位記憶領域割当解除手段1 fは、検出されたデータユニットに対応付けて、取得した重複排除アドレスをデータユニット情報記憶手段1 cに格納する。また単位記憶領域割当解除手段1 fは、検出されたデータユニットへの単位記憶領域5 aの割り当てを解除する。割当が解除された単位記憶領域5 aは、未使用の状態となり、他のデータユニットまたは重複排除ユニットへの割り当てが可能となる。なお重複排除アドレス取得時にインデックスサーバ4において重複排除ユニット割当処理が行われる場合がある。この場合、単位記憶領域割当解除手段1 fは、単位記憶領域5 aの割り当てを解除する前に、単位記憶領域5 aに格納されているデータを、重複排除アドレスで示されるディスクノードに転送する。

10

【0027】

読み出しアクセス手段1 gは、データユニットを指定した読み出し要求があると、読み出し対象のデータユニットに割り当てられた単位記憶領域5 aがある場合、割り当てられた単位記憶領域5 aからデータを読み出す。また読み出しアクセス手段1 gは、読み出し対象のデータユニットに割り当てられた単位記憶領域5 aがない場合、データユニット情報記憶手段1 cを参照する。そして読み出しアクセス手段1 gは、読み出し対象のデータユニットに対応付けられた重複排除アドレスに示されるディスクノードから、重複排除アドレスに示される重複排除ユニットの識別情報を指定して重複排除対象データを取得する。

20

【0028】

重複排除ユニット割当手段1 hは、インデックスサーバ4から重複排除対象データの重複排除ユニット割当要求を受信すると、未使用の重複排除ユニットに単位記憶領域の1つを割り当て、割り当てた単位記憶領域に重複排除対象データを格納する。

【0029】

重複排除データアクセス応答手段1 iは、重複排除ユニットの識別情報を指定したデータの読み出し要求を受信すると、指定された重複排除ユニットに割り当てられた単位記憶領域内の重複排除対象データを応答する。

30

【0030】

ディスクノード2, 3もディスクノード1と同様の機能の要素を有する。なお図1では、ディスクノード2の有する要素のうち、重複排除ユニット割当手段2 hと重複排除データアクセス応答手段2 iが示されている。また図1では、ディスクノード3の有する要素のうち、データユニット情報記憶手段3 c、重複排除化データユニット検出手段3 d、重複排除アドレス取得手段3 e、および単位記憶領域割当解除手段3 fが示されている。ディスクノード2, 3に示される各要素は、ディスクノード1内の同名の要素と同じ機能を有する。

40

【0031】

インデックスサーバ4は、重複排除ユニット情報記憶手段4 a、重複排除ユニットアドレス検索手段4 b、および重複排除ユニット割当要求手段4 cを有する。

重複排除ユニット情報記憶手段4 aは、重複排除ユニット情報を、その重複排除ユニット情報に対応する重複排除ユニットを管理するディスクノードの識別子に対応付けて記憶する。重複排除ユニット情報には、重複排除ボリュームの重複排除ユニットのうち、ディスクノード1~3で管理している重複排除ユニットの使用の有無が示される。また重複排除ユニット情報には、使用されている重複排除ユニットに割り当てられている単位記憶領域に格納された重複排除対象データの固有値、および重複排除ユニットの識別情報が含まれる。

50

## 【 0 0 3 2 】

重複排除ユニットアドレス検索手段 4 b は、ディスクノード 1 ~ 3 から、重複排除化データユニットに割り当てられた単位記憶領域内のデータに所定の演算を実行して得られた固有値を含む重複排除アドレス照会要求を受け取る。すると重複排除ユニットアドレス検索手段 4 b は、重複排除ユニット情報記憶手段 4 a から、重複排除アドレス照会要求に示される固有値を含む重複排除ユニット情報を検索する。重複排除ユニットアドレス検索手段 4 b は、検索により該当する重複排除ユニット情報が検出された場合、重複排除アドレスを、重複排除アドレス照会要求の送信元であるディスクノードに応答する。この重複排除アドレスには、検出された重複排除ユニット情報に含まれる重複排除ユニットの識別情報と、検出された重複排除ユニット情報に対応する重複排除ユニットを管理しているディスクノードの識別子とが含まれる。

10

## 【 0 0 3 3 】

また重複排除ユニットアドレス検索手段 4 b は、検索により該当する重複排除ユニット情報が検出されなかった場合、重複排除ユニット割当要求手段 4 c に、重複排除ユニットへの単位記憶領域の割り当て要求を依頼する。そして重複排除ユニットアドレス検索手段 4 b は、重複排除ユニット割当要求手段 4 c による割り当て結果に応じた重複排除アドレスを、重複排除アドレス照会要求の送信元であるディスクノードに応答する。この重複排除アドレスには、重複排除ユニット割当要求手段 4 c が選択した重複排除ユニットの識別情報と、重複排除ユニット割当要求手段 4 c が選択した重複排除ユニットを管理するディスクノードの識別子とが含まれる。このとき、重複排除ユニットアドレス検索手段 4 b は、応答内容に、重複排除ユニット割り当て処理が実行されたことを示す情報を含める。この情報を含めることで、重複排除アドレスを取得したディスクノードにおいて、重複排除化の対象となったデータユニットに割り当てられた記憶領域内のデータの転送が行われる。

20

## 【 0 0 3 4 】

重複排除ユニット割当要求手段 4 c は、重複検索により該当する重複排除ユニット情報が検出されなかった場合、重複排除ユニット情報記憶手段 4 a を参照し、未使用の重複排除ユニットを選択する。次に重複排除ユニット割当要求手段 4 c は、選択した重複排除ユニットを管理するディスクノードに対して選択した重複排除ユニットへの単位記憶領域の割り当て要求を送信する。さらに重複排除ユニット割当要求手段 4 c は、単位記憶領域割り当て後の重複排除ユニット情報を重複排除ユニット情報記憶手段 4 a に格納する。そして重複排除ユニット割当要求手段 4 c は、選択した重複排除ユニットの識別情報と、選択した重複排除ユニット情報を管理するディスクノードの識別子とを含む重複排除アドレスを、重複排除ユニットアドレス検索手段 4 b に渡す。

30

## 【 0 0 3 5 】

こ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、まずディスクノード 1 に対してデータ (data[a]) の書き込み要求が入力されたものとする。その場合、書き込みアクセス手段 1 a により、書き込み要求で指定されたデータユニットに単位記憶領域 5 a が割り当てられる。そして、書き込みアクセス手段 1 a により、単位記憶領域 5 a にデータ (data[a]) が書き込まれる。このとき、最終書き込み時刻更新手段 1 b により、書き込み対象となったデータユニットに対応付けて、現在の時刻が最終書き込み時刻としてデータユニット情報記憶手段 1 c に格納される。

40

## 【 0 0 3 6 】

その後、重複排除化データユニット検出手段 1 d により、最終書き込み時刻から重複排除化猶予期間以上経過しているデータユニット (重複排除化データユニット) が検出される。ここで、データ (data[a]) が重複排除猶予期間以上更新されなかったものとする。この場合、データ (data[a]) が格納された単位記憶領域 5 a が割り当てられたデータユニットについて、最終書き込み時刻から重複排除化猶予期間以上経過したことが検出される。

## 【 0 0 3 7 】

50

重複排除化データユニットが検出されると、重複排除アドレス取得手段 1 e により、インデックスサーバ 4 から、重複排除化データユニットに割り当てられた単位記憶領域 5 a 内のデータの固有値に対応付けられた重複排除アドレスが取得される。例えば、重複排除アドレス取得手段 1 e からインデックスサーバ 4 へ、重複排除化データユニットに割り当てられた単位記憶領域 5 a 内のデータ (data[a]) に所定の演算を実行して得られた固有値を含む重複排除アドレス照会要求が送信される。

【 0 0 3 8 】

インデックスサーバ 4 の重複排除ユニットアドレス検索手段 4 b では、重複排除アドレス照会要求に応答し、重複排除ユニット情報記憶手段 4 a から、重複排除アドレス照会要求に示される固有値を含む重複排除ユニット情報が検索される。この時点では、データ (data[a]) の固有値を有する重複排除ユニット情報は、重複排除ユニット情報記憶手段 4 a に存在しないものとする。この場合、検索により該当する重複排除ユニット情報が検出されない。

10

【 0 0 3 9 】

検索により該当する重複排除ユニット情報が検出されなかった場合、重複排除ユニット割当要求手段 4 c により、重複排除ユニット情報記憶手段 4 a から未使用の重複排除ユニットが選択される。この例では、ディスクノード 2 が管理する重複排除ユニットが選択されたものとする。次に重複排除ユニット割当要求手段 4 c により、選択した重複排除ユニットを管理するディスクノード 2 に対して選択した重複排除ユニットへの単位記憶領域の割り当て要求が送信される。すると、ディスクノード 2 の重複排除ユニット割当手段 2 h により、未使用の重複排除ユニットに単位記憶領域 6 a が割り当てられる。さらに重複排除ユニット割当要求手段 4 c により、選択した重複排除ユニットの単位記憶領域割り当て後の重複排除ユニット情報が重複排除ユニット情報記憶手段 4 a に格納される。そして重複排除ユニット割当要求手段 4 c により、選択した重複排除ユニットの識別情報と、選択した重複排除ユニット情報を管理するディスクノード 2 の識別子とを含む重複排除アドレスが重複排除ユニットアドレス検索手段 4 b に渡される。すると、重複排除ユニットアドレス検索手段 4 b により、重複排除ユニット割当要求手段 4 c から渡された重複排除アドレスが、重複排除アドレス照会要求の送信元であるディスクノード 1 に応答される。

20

【 0 0 4 0 】

ディスクノード 1 では、重複排除アドレスを重複排除アドレス取得手段 1 e が受け取る。すると単位記憶領域割当解除手段 1 f により、検出されたデータユニットに対応付けて、取得した重複排除アドレスがデータユニット情報記憶手段 1 c に格納される。また、単位記憶領域割当解除手段 1 f により、単位記憶領域 5 a に格納されているデータ (data[a]) が、重複排除アドレスで示されるディスクノード 2 に転送される。すると、ディスクノード 2 の重複排除データアクセス応答手段 2 i によって、割り当てた単位記憶領域 6 a にデータ (data[a]) が格納される。

30

【 0 0 4 1 】

その後、単位記憶領域割当解除手段 1 f により、検出されたデータユニットへの単位記憶領域 5 a の割り当てが解除される。割当が解除された単位記憶領域 5 a は、未使用の状態となり、他のデータユニットまたは重複排除ユニットへの割り当てが可能となる。

40

【 0 0 4 2 】

このようにして、ディスクノード 1 で管理されていたデータ (data[a]) が、重複排除対象データとしてディスクノード 2 で管理されるようになる。

ディスクノード 1 では、重複排除アドレスに基づいてデータ (data[a]) を読み出すことができる。例えば、ディスクノード 1 に対して読み出し要求があったとき、読み出しアクセス手段 1 g により、読み出し対象のデータユニットに割り当てられた単位記憶領域があれば、割り当てられた単位記憶領域 5 a からデータが読み出される。他方、単位記憶領域 5 a の割り当てが解除されたデータユニットへの読み出し要求が出された場合、ストレージ装置 5 内には既にデータが保存されていない。この場合、読み出しアクセス手段 1 g により、データユニット情報記憶手段 1 c に基づいて、読み出し対象のデータユニットに

50

対応付けられた重複排除アドレスに示されるディスクノード 2 からデータ (data[a]) が取得される。ディスクノード 2 からデータを取得する際には、重複排除アドレスに示される重複排除ユニットの識別情報が指定される。これにより、取得対象のデータが一意に決定される。

【 0 0 4 3 】

その後、ディスクノード 3 において、重複排除対象データと同じ内容のデータ (data[a]) がストレージ装置 7 に格納され、重複排除化猶予期間以上更新されなかったものとする。この場合、ディスクノード 3 の重複排除化データユニット検出手段 3 d により、データ (data[a]) が格納された単位記憶領域 7 a が割り当てられたデータユニットについて、最終書き込み時刻から重複排除化猶予期間以上経過したことが検出される。

10

【 0 0 4 4 】

重複排除化データユニットが検出されると、重複排除アドレス取得手段 3 e により、インデックスサーバ 4 から、重複排除化データユニットに割り当てられた単位記憶領域 7 a 内のデータの固有値に対応付けられた重複排除アドレスが取得される。例えば、重複排除アドレス取得手段 3 e からインデックスサーバ 4 へ、重複排除化データユニットに割り当てられた単位記憶領域 7 a 内のデータ (data[a]) の固有値を含む重複排除アドレス照会要求が送信される。

【 0 0 4 5 】

インデックスサーバ 4 の重複排除ユニットアドレス検索手段 4 b では、重複排除アドレス照会要求に回答し、重複排除ユニット情報記憶手段 4 a から、重複排除アドレス照会要求に示される固有値を含む重複排除ユニット情報が検索される。このとき単位記憶領域 7 a 内のデータと同じ内容のデータ (data[a]) が単位記憶領域 6 a に格納されている。そのため、単位記憶領域 6 a が割り当てられている重複排除ユニットの重複排除ユニット情報が、検索結果として抽出される。すると重複排除ユニットアドレス検索手段 4 b により、重複排除アドレスが、重複排除アドレス照会要求の送信元であるディスクノード 3 に回答される。

20

【 0 0 4 6 】

ディスクノード 3 では、重複排除アドレスを重複排除アドレス取得手段 3 e が受け取る。すると単位記憶領域割当解除手段 3 f により、検出されたデータユニットに対応付けて、取得した重複排除アドレスがデータユニット情報記憶手段 3 c に格納される。その後、単位記憶領域割当解除手段 3 f により、検出されたデータユニットへの単位記憶領域 7 a の割り当てが解除される。なお、ディスクノード 3 からディスクノード 1 と同様に、ストレージ装置 6 に格納されたデータ (data[a]) を読み出すことができる。

30

【 0 0 4 7 】

このようにして、マルチノードストレージシステムにおいて、複数のディスクノードで重複して格納されている同一内容の複数のデータが、1つのデータに統合される。すなわち、マルチノードストレージシステムにおいて、ノード間にまたがって重複するデータを除去して保持することで冗長性が低減され、ストレージ容量の使用効率を上げることができる。

【 0 0 4 8 】

また、上記実施の形態では、重複排除猶予期間以上更新されないデータのみが重複排除の対象となる。換言すると、書き込み頻度が高いデータに関しては重複排除が行われない。これにより、性能影響を低く抑えつつ重複するデータの除去が可能となる。

40

【 0 0 4 9 】

〔 第 2 の実施の形態 〕

次に、第 2 の実施の形態について説明する。第 2 の実施の形態は、論理ボリュームの記憶領域を一定の記憶容量のスライスに分割して管理するマルチノードストレージシステムに、第 1 の実施の形態で示したような重複排除技術を適用したものである。

【 0 0 5 0 】

第 2 の実施の形態では、論理ボリュームがセグメントに分けて管理されており、各セグ

50

メントにストレージ装置のスライスが割り当てられている。第2の実施の形態のシステムでは、このスライスをユニットに細分化し、ユニット単位でデータ実体を格納するための記憶領域（単位記憶領域）を割り当て、ユニット単位で重複データが存在するかどうかの管理を行う。

【0051】

また、第2の実施の形態に係るシステムでは定期的にパトロールが実行され、アクセス頻度が低くなったユーザの論理ボリュームのユニット（データユニット）を重複排除される。逆に重複排除されたデータユニットに書き込みが発生すると、重複排除状態が解除され、再度、論理ボリューム内に単位記憶領域が割り当てられる。また、論理ボリューム割り当て直後の書き込みがなされていない状態のユニットには単位記憶領域は割り当てられず、書き込みが発生して初めて単位記憶領域が割り当てられる。

10

【0052】

重複排除ボリュームにある重複排除ユニットの情報は、サーバ起動時にインデックスサーバに集められる。ディスクノードにおいてデータユニットを重複排除化するには、既に同じデータが存在するかが検索され、存在していなければ重複排除対象のデータを格納する新たな重複排除ユニットが設定される。

【0053】

重複排除ユニットには保存期限満了時刻が定められており、一定時間参照がなかった重複排除ユニットが保持するデータは廃棄される。システム内部では定期的にパトロールが実行され、すべてのデータユニットは指定された時間内に必ずパトロールが実施される。重複排除化されたデータユニットのパトロールの際には、そのデータユニットの重複排除化によるデータの参照先の重複排除ユニットに対して、参照が行われる。ユーザからの読み出し要求またはパトロールにより参照処理が行われた重複排除ユニットの保存期限満了時刻は、所定の保存期間だけ延長される。延長される保存期間より短い間隔でパトロールを実行することで、少なくとも1つのデータユニットからの参照先となっている重複排除ユニットについては、継続的にデータが保持される。他方、保存期限満了時刻を過ぎた重複排除ユニットは、単位記憶領域の割り当てが解除され、保持していたデータが消去される。

20

【0054】

図2は、第2の実施の形態のシステム構成例を示す図である。ネットワーク10には複数のディスクノード100、200、300、400、インデックスサーバ500、制御ノード600、およびアクセスノード700、800が接続されている。

30

【0055】

ディスクノード100、200、300、400は、それぞれ「DP-A」、「DP-B」、「DP-C」、「DP-D」のディスクノードIDが設定されており、このディスクノードIDによって一意に識別される。ディスクノード100、200、300、400は、それぞれストレージ装置110、210、310、410を有している。各ストレージ装置110、210、310、410は、例えばRAID（Redundant Arrays of Inexpensive Disks）5によりデータを管理するストレージシステムである。ディスクノード100、200、300、400は、制御ノード600からの指示に従って、論理ボリュームのデータを管理する。また、ディスクノード100、200、300、400は、インデックスサーバ500からの指示に従って、重複排除ボリューム内のデータを管理する。そして、ディスクノード100、200、300、400は、アクセスノード700、800からの要求に応じて、論理ボリュームにおける管理対象のデータを、ストレージ装置110、210、310、410へ入出力する。さらに、ディスクノード100、200、300、400は、他のディスクノードからの要求に応じて、重複排除ボリュームにおける管理対象のデータを、ストレージ装置110、210、310、410へ入出力する。

40

【0056】

インデックスサーバ500は、重複排除ボリュームに格納するデータのストレージ装置

50

110, 210, 310, 410への割り当てを行う。また、既に重複排除ボリュームに格納されているデータと同じ内容のデータが重複排除対象となった場合、該当するデータのアドレスを、重複排除対象となったデータを管理するディスクノードに通知する。さらに、インデックスサーバ500は、重複排除ボリューム内のデータのうち、いずれのディスクノードからも利用されていないデータを定期的に探索する。

**【0057】**

制御ノード600は、論理ボリュームや重複排除ボリュームを作成する。そして制御ノード600は、論理ボリュームや重複排除ボリューム内のセグメントへのストレージ装置110, 210, 310, 410の記憶領域の割り当てを行う。論理ボリュームと重複排除ボリュームへの記憶領域の割り当て結果は、制御ノード600からディスクノード100, 200, 300, 400に通知される。同様に論理ボリュームへの記憶領域の割り当て結果は、アクセスノード700, 800に通知される。

10

**【0058】**

アクセスノード700, 800は、ユーザからの論理ボリュームへのアクセス要求に応答し、制御ノード600から通知された論理ボリュームへの記憶領域の割り当て関係に基づいて、アクセス対象のデータを管理するディスクノードを判断する。さらにアクセスノード700, 800は、該当するディスクノードに対してデータのアクセス要求を送信する。そして、アクセスノード700, 800は、アクセス要求に対するディスクノードからの応答をユーザに返す。

**【0059】**

20

このようなシステムにより、論理ボリュームで管理されているデータに重複排除対象データがあれば、そのデータが、適宜、重複排除ボリュームに移動される。本実施の形態では、重複排除ボリューム内のデータは、ハッシュ値によって論理ボリューム内のデータと関連付けられる。重複排除ボリューム内のデータに対して、論理ボリューム内の複数のデータを関連付けることで、同じ内容のデータが重複して格納されることを排除できる。

**【0060】**

図3は、本実施の形態に用いるディスクノードのハードウェア構成例を示す図である。ディスクノード100は、CPU (Central Processing Unit) 101によって装置全体が制御されている。CPU 101には、バス109を介してRAM (Random Access Memory) 102と複数の周辺機器が接続されている。

30

**【0061】**

RAM 102は、ディスクノード100の主記憶装置として使用される。RAM 102には、CPU 101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 102には、CPU 101による処理に必要な各種データが格納される。

**【0062】**

バス109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) 103、グラフィック処理装置104、入力インタフェース105、光学ドライブ装置106、通信インタフェース107、およびストレージインタフェース108がある。

40

**【0063】**

HDD 103は、内蔵したディスクに対して、磁気的にデータの書き込みおよび読み出しを行う。HDD 103は、ディスクノード100の二次記憶装置として使用される。HDD 103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。

**【0064】**

グラフィック処理装置104には、モニタ11が接続されている。グラフィック処理装置104は、CPU 101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

50

がある。

【 0 0 6 5 】

入力インタフェース 1 0 5 には、キーボード 1 2 とマウス 1 3 とが接続されている。入力インタフェース 1 0 5 は、キーボード 1 2 やマウス 1 3 から送られてくる信号を CPU 1 0 1 に送信する。なお、マウス 1 3 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。

【 0 0 6 6 】

光学ドライブ装置 1 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、DVD (Digital Versatile Disc)、DVD-RAM、CD-ROM (Compact Disc Read Only Memory)、CD-R (Recordable) / RW (ReWritable) などがある。

10

【 0 0 6 7 】

通信インタフェース 1 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、他のコンピュータとの間でデータの送受信を行う。

【 0 0 6 8 】

ストレージインタフェース 1 0 8 にはストレージ装置 1 1 0 が接続されている。ストレージインタフェース 1 0 8 は、CPU 1 0 1 からの指示に従ってストレージ装置 1 1 0 へのデータの入出力を行う。

20

【 0 0 6 9 】

以上のようなハードウェア構成によって、第 2 の形態の処理機能を実現することができる。なお図 3 には、ディスクノード 1 0 0 のハードウェア構成例を示したが、他のディスクノード 2 0 0 , 3 0 0 , 4 0 0、インデックスサーバ 5 0 0、制御ノード 6 0 0、およびアクセスノード 7 0 0 , 8 0 0 も同様のハードウェアで実現することができる。

【 0 0 7 0 】

次に、論理ボリュームと重複排除ボリュームとに対するストレージ装置の記憶領域の割り当てについて説明する。

図 4 は、論理ボリュームと重複排除ボリュームへの記憶領域の割り当て関係を示す図である。図 4 の例では、2 つの論理ボリューム 2 0 , 3 0 と 1 つの重複排除ボリューム 4 0 とが設けられている。

30

【 0 0 7 1 】

論理ボリューム 2 0 には、識別子として論理ボリューム ID 「L V O L - X」が付与されている。論理ボリューム 2 0 は、複数のセグメント 2 1 , 2 2 , 2 3 , . . . を有している。論理ボリューム 3 0 には、識別子として論理ボリューム ID 「L V O L - Y」が付与されている。論理ボリューム 3 0 は、複数のセグメント 3 1 , 3 2 , 3 3 , . . . を有している。重複排除ボリューム 4 0 は、識別子として重複排除ボリューム ID 「D e D u p - U」が付与されている。重複排除ボリューム 4 0 は、複数のセグメント 4 1 , 4 2 , 4 3 , . . . を有している。

40

【 0 0 7 2 】

論理ボリューム 2 0 , 3 0 や重複排除ボリューム 4 0 の各セグメントは、所定の容量の仮想的なデータ記憶領域である。例えば 1 セグメント当たり 1 G B の記憶領域が設けられる。各セグメント 2 1 , 2 2 , 2 3 , . . . に対して、ストレージ装置 1 1 0 , 2 1 0 , 3 1 0 , 4 1 0 の記憶領域が割り当てられる。

【 0 0 7 3 】

ストレージ装置 1 1 0 には、複数のスライス 1 1 1 , 1 1 2 , 1 1 3 , 1 1 4 , . . . が定義されている。各スライス 1 1 1 , 1 1 2 , 1 1 3 , 1 1 4 , . . . は、論理ボリューム 2 0 , 3 0 のセグメントと同じ容量の記憶領域である。なおスライス 1 1 1 , 1 1 2 , 1 1 3 , 1 1 4 , . . . それぞれの記憶領域は、ストレージ装置 1 1 0 内の不連続 (連

50

続しても良い)の記憶領域である。例えば1スライスの記憶容量が1GBであれば、スライス111の記憶領域として、1GB分の記憶領域が所定データ長のデータユニット単位で、ストレージ装置110内に確保される。

【0074】

同様に、ストレージ装置210には、複数のスライス211, 212, 213, 214, …が定義されている。またストレージ装置310にも、複数のスライス311, 312, 313, 314, …が定義されている。さらにストレージ装置410にも、複数のスライス411, 412, 413, 414, …が定義されている。

【0075】

各ストレージ装置110, 210, 310, 410のスライスが、論理ボリューム20, 30または重複排除ボリューム40のセグメントに割り当てられる。図4の例では、論理ボリューム20のセグメント21~23のセグメントIDを、「X1」、「X2」、「X3」としている。また論理ボリューム30のセグメント31~33のセグメントIDを、「Y1」、「Y2」、「Y3」としている。さらに重複排除ボリューム40のセグメント41~43のセグメントIDを、「U1」、「U2」、「U3」としている。そして図4では、ストレージ装置110, 210, 310, 410内のスライスに対して、割り当て先のセグメントのセグメントIDを示している。

【0076】

ストレージ装置110, 210, 310, 410には、各スライスの定義情報であるメタデータが格納されている。

図5は、ストレージ装置に格納されている情報の例を示す図である。ストレージ装置110の記憶領域は、メタデータ領域110aとデータ実体領域110bとに分かれている。メタデータ領域110aは、各スライスを管理するためのメタデータが格納される。メタデータのうち、論理ボリュームに割り当てられたスライスのメタデータを、論理ボリューム用スライスメタデータ51, 52, 53, …とする。また、重複排除ボリュームに割り当てられたスライスのメタデータを、重複排除ボリューム用スライスメタデータ61, 62, 63, …とする。

【0077】

データ実体領域110bは、論理ボリュームまたは重複排除ボリュームに格納されるデータ実体を記憶する記憶領域である。データ実体領域110bは、ユニットと同じデータサイズの単位記憶領域71~76に分けられている。ここで、論理ボリュームに割り当てられたスライスの記憶領域となるユニット(データユニット)に対して、単位記憶領域71~73が割り当てられる。重複排除ボリュームに割り当てられたスライスのユニット(重複排除ユニット)に対して、単位記憶領域74~76が割り当てられる。

【0078】

メタデータ領域110a内の各メタデータには、そのメタデータに対応するスライス内の各ユニットに割り当てられた単位記憶領域の位置が示されている。また各メタデータには、割当先となるセグメントが属する論理ボリュームまたは重複排除ボリュームの識別子と、セグメントIDとが設定されている。これにより、メタデータを介して、各セグメントと実際のデータを格納する単位記憶領域とが関連付けられることとなる。すなわち、メタデータ領域110aは、第1の実施の形態におけるディスクノード1のデータユニット情報記憶手段1cの機能を含んでいる。

【0079】

メタデータの総数は、データ実体領域110bをスライスの記憶容量で除算したときの商で示される数(最少スライス数)以上となる。例えばデータ実体領域の記憶容量が300GBであり、メタデータのサイズが1GBであれば、少なくとも300個のスライスを定義することができ、そのスライスに対応するメタデータが作成できる。本実施の形態ではデータユニット単位でのシステム内での同一データの重複が排除される。そのため、最少スライス数よりも多くのスライスを、定義することが可能である。なお、メタデータは、ストレージ装置110内のスライスを論理ボリュームまたは重複排除ボリュームのセグ

10

20

30

40

50



メントに割り当てたときに制御ノード600により作成される。そして、制御ノード600からディスクノード100に対して新たに割り当てられたスライスのメタデータが送信され、ディスクノード100によってストレージ装置110のメタデータ領域110aにメタデータが格納される。

【0080】

次に、論理ボリューム用スライスメタデータと重複排除ボリューム用スライスメタデータとのデータ構造について説明する。

まず論理ボリューム用スライスメタデータのデータ構造について説明する。

【0081】

図6は、論理ボリューム用スライスメタデータのデータ構造例を示す図である。論理ボリューム用スライスメタデータ50には、論理ボリュームID、セグメントID、および最終パトロール時刻のフィールドが設けられている。論理ボリュームIDのフィールドには、論理ボリューム用スライスメタデータ50で定義されるスライスが割り当てられたセグメントが属する論理ボリュームの論理ボリュームIDが設定される。セグメントIDのフィールドには、論理ボリューム用スライスメタデータ50で定義されるスライスが割り当てられたセグメントのセグメントIDが設定される。最終パトロール時刻のフィールドには、論理ボリューム用スライスメタデータ50で定義されるスライスに対するパトロールを最後に実行した時刻が設定される。

10

【0082】

また論理ボリューム用スライスメタデータ50には、スライスに割り当てられたデータユニットごとに、データユニット情報50aが設けられている。データユニット管理情報は、スライスサイズをデータユニットサイズで除算したときの商で示される個数だけ設けられる。

20

【0083】

各データユニット管理情報には、状態情報のフィールドが設けられている。状態情報のフィールドには、「Blank」、「Normal」、「DeDup」の3つの状態のいずれかが設定される。「Blank」は、書き込みがされていないデータユニットであり、データ実体を持たないことを示す。なおスライスをセグメントに割り当てた直後は、すべてのデータユニットの状態が「Blank」となっている。

【0084】

「Normal」は、書き込みが行われたデータユニットであり、データユニットに割り当てられた単位記憶領域にデータ実体が格納されていることを示す。状態が「Normal」のデータユニット管理情報には、さらに最終書き込み時刻とデータ実体オフセットとのフィールドが設けられている。最終書き込み時刻のフィールドには、データユニットに対するデータの書き込みが最後に行われた時刻が設定される。データ実体オフセットのフィールドには、データ実体を格納している単位記憶領域の位置が、データ実体領域110bの先頭からのオフセットで示される。

30

【0085】

「DeDup」は、書き込みが行われたデータユニットであり、データ実体が重複排除対象となって重複排除ボリューム40で管理されていることを示す。状態が「DeDup」のデータユニット管理情報には、さらに重複排除ボリュームIDと重複排除オフセットとのフィールドが設けられている。重複排除ボリュームIDのフィールドには、データ実体を管理している重複排除ボリュームの重複排除ボリュームIDが設定される。重複排除オフセットのフィールドには、重複排除ボリューム内でのデータ実体を格納している単位記憶領域の位置が、重複排除ボリュームの先頭からのオフセットで示される。

40

【0086】

次に重複排除ボリューム用スライスメタデータのデータ構造について説明する。

図7は、重複排除ボリューム用スライスメタデータのデータ構造例を示す図である。重複排除ボリューム用スライスメタデータ60には、重複排除ボリュームID、セグメントID、および最終パトロール時刻のフィールドが設けられている。重複排除ボリュームI

50

Dのフィールドには、重複排除ボリューム用スライスマタデータ60で定義されるスライスが割り当てられたセグメントが属する重複排除ボリュームの重複排除ボリュームIDが設定される。セグメントIDのフィールドには、重複排除ボリューム用スライスマタデータ60で定義されるスライスが割り当てられたセグメントのセグメントIDが設定される。最終パトロール時刻のフィールドには、重複排除ボリューム用スライスマタデータ60で定義されるスライスに対するパトロールを最後に実行した時刻が設定される。

【0087】

また重複排除ボリューム用スライスマタデータ60には、スライスに割り当てられた重複排除ユニットごとに、重複排除ユニット情報60aが設けられている。重複排除ユニット情報60aは、スライスサイズを重複排除ユニットサイズで除算したときの商で示される個数だけ設けられる。セグメントIDで示されるセグメントの領域の開始位置に近いユニットの重複排除ユニット情報60aほど、重複排除ボリューム用スライスマタデータ内で上位に設定されている。すなわち、セグメントの領域の開始位置を示すオフセットと、重複排除ユニット情報60aの重複排除ボリューム用スライスマタデータ60内での順番により、その重複排除ユニット情報60aに対応するユニットの重複排除オフセットが求められる。

10

【0088】

重複排除ユニット情報60aには、ハッシュ値、保存期限満了時刻(expire時刻)、およびデータ実体オフセットのフィールドが設けられている。ハッシュ値のフィールドには、重複排除ユニットに書き込まれたデータ実体を所定のハッシュ関数で演算することで得られるハッシュ値が設定される。保存期限満了時刻のフィールドには、重複排除対象となったデータの保存期限が満了する時刻が設定される。重複排除ユニットの保存期限満了時刻が経過すると、該当する重複排除ユニット内が初期化される。

20

【0089】

なお、保存期限満了時刻は、重複排除ユニットへのデータ書き込み時に初期値が設定される。設定される保存期限満了時刻は、書き込み時刻を基準として、各スライスのパトロール間隔よりも長い時間経過後の時刻が設定される。

【0090】

重複排除ユニットに対して参照要求が出されると、この保存期限満了時刻は随時延長される。その際にも再設定する保存期限満了時刻は、参照時刻を基準として、各スライスのパトロール間隔よりも長い所定の保存期間経過後の時刻が設定される。

30

【0091】

なお、データの参照要求は、重複排除対象データに対してアクセスノードから参照要求が出された場合、および重複排除対象データと同じデータを有するスライスに対するパトロールが実行された場合に発生する。

【0092】

データ実体オフセットのフィールドには、重複排除ユニットのデータ実体が保存されている単位記憶領域の位置が、データ実体領域110bの先頭からのオフセットで示される。

【0093】

ストレージ装置110, 210, 310, 410に図5~図7に示した情報が格納されている状態で、システム内の各機器が起動されると、各機器において分散ストレージシステムの運用に必要な機能が動作開始する。

40

【0094】

図8は、制御ノードとアクセスノードとの機能を示すブロック図である。制御ノード600は、スライスマタデータ収集部610、スライスマタデータ記憶部620、論理ボリューム管理部630、論理ボリューム構成情報記憶部640、および重複排除ボリューム構成情報記憶部650を有している。

【0095】

スライスマタデータ収集部610は、制御ノード600の起動時に、各ディスクノード

50

100, 200, 300, 400からメタデータを収集する。すなわち、スライスメタデータ収集部610は、論理ボリューム用スライスメタデータと重複排除ボリューム用スライスメタデータとの取得要求を、各ディスクノード100, 200, 300, 400に対して送信する。すると各ディスクノード100, 200, 300, 400からストレージ装置110, 210, 310, 410に格納されている論理ボリューム用スライスメタデータと重複排除ボリューム用スライスメタデータとが制御ノード600に送信される。スライスメタデータ収集部610は、収集したメタデータを、スライスメタデータ記憶部620に格納する。

【0096】

スライスメタデータ記憶部620は、メタデータ（論理ボリューム用スライスメタデータと重複排除ボリューム用スライスメタデータ）を記憶する記憶領域である。具体的には、スライスメタデータ記憶部620には、図6や図7に示したデータ構造の論理ボリューム用スライスメタデータと重複排除ボリューム用スライスメタデータとが、送信元のディスクノードの識別子（ディスクノードID）に関連付けて記憶される。スライスメタデータ記憶部620としては、例えば制御ノード600内のRAMやHDDの記憶領域の一部が用いられる。

【0097】

論理ボリューム管理部630は、スライスメタデータ記憶部620に格納された論理ボリューム用スライスメタデータに基づいて、論理ボリューム構成情報を生成する。具体的には、論理ボリューム管理部630は、論理ボリューム用スライスメタデータのうち、論理ボリュームIDの値が同じものをグループ化する。次に論理ボリューム管理部630は、グループ化した各論理ボリューム用スライスメタデータから、セグメントIDと、その論理ボリューム用スライスメタデータの送信元であるディスクノードのディスクノードIDとの組を作成する。そして、論理ボリューム管理部630は、セグメントIDとディスクノードIDとの組をセグメントIDでソートし、グループに共通の論理ボリュームIDで示される論理ボリュームの構成情報とする。論理ボリューム管理部630は、生成した論理ボリューム構成情報を、論理ボリューム構成情報記憶部640に格納する。

【0098】

また論理ボリューム管理部630は、スライスメタデータ記憶部620に格納された重複排除ボリューム用スライスメタデータに基づいて、重複排除ボリューム構成情報を生成する。具体的には、論理ボリューム管理部630は、重複排除ボリューム用スライスメタデータのうち、重複排除ボリュームIDの値が同じものをグループ化する。次に論理ボリューム管理部630は、グループ化した各重複排除ボリューム用スライスメタデータから、セグメントIDと、その重複排除ボリューム用スライスメタデータの送信元であるディスクノードのディスクノードIDとの組を作成する。そして、論理ボリューム管理部630は、セグメントIDとディスクノードIDとの組をセグメントIDでソートし、グループに共通の重複排除ボリュームIDで示される重複排除ボリュームの構成情報とする。論理ボリューム管理部630は、生成した重複排除ボリューム構成情報を重複排除ボリューム構成情報記憶部650に格納する。

【0099】

さらに論理ボリューム管理部630は、論理ボリュームのセグメントを指定したメタデータ要求をアクセスノード700、800から受け取ると、該当セグメントのメタデータをアクセスノード700、800に通知する。同様に論理ボリューム管理部630は、重複排除ボリュームのセグメントを指定したメタデータ要求をディスクノード100, 200, 300, 400から受け取ると、該当セグメントのメタデータをディスクノード100, 200, 300, 400に通知する。

【0100】

アクセスノード700は、論理ボリューム構成情報記憶部710、データアクセス要求部720を有している。

論理ボリューム構成情報記憶部710は、論理ボリュームのセグメントごとに、そのセ

10

20

30

40

50

グメントのデータを管理しているディスクノードの情報を記憶する。具体的には、論理ボリュームIDに対応付けて、各セグメントの論理ボリューム内での占有域を示すオフセットの範囲とディスクノードIDとの組が格納されている。論理ボリューム構成情報記憶部710として、例えばアクセスノード700内のRAMまたはHDDの記憶領域の一部が使用される。

**【0101】**

データアクセス要求部720は、ユーザからのデータアクセス要求が発生すると、論理ボリューム構成情報記憶部710を参照して、アクセス対象のデータを管理しているディスクノードを判断する。そして、データアクセス要求部720は、アクセス対象のデータ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。

10

**【0102】**

なお、データアクセス要求部720は、論理ボリューム構成情報記憶部710内にアクセス対象のデータが属するセグメントの情報が格納されていない場合、制御ノード600に対して、該当セグメントのメタデータを要求する。データアクセス要求部720は、制御ノード600から応答されたメタデータによって、アクセス対象となるデータを管理しているディスクノードを判断する。そして、データアクセス要求部720は、取得したメタデータの内容に基づいて論理ボリューム構成情報記憶部710に新たなセグメントの情報を登録すると共に、アクセス対象のデータ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。

**【0103】**

データアクセス要求部720は、アクセス要求に対する応答をディスクノードから取得すると、応答されたアクセス結果をユーザからのデータアクセス要求に対する応答とする。

20

**【0104】**

なお、図8にはアクセスノード700の機能を示したが、アクセスノード800もアクセスノード700と同様の機能を有している。

次に、制御ノード600とアクセスノード700、800とが記憶する情報について詳細に説明する。

**【0105】**

図9は、制御ノードのスライスメタデータ記憶部のデータ構造例を示す図である。スライスメタデータ記憶部620には、メタデータの送信元であるディスクノードのディスクノードIDに対応付けて、そのディスクノードから送信された論理ボリューム用スライスメタデータと重複排除ボリューム用スライスメタデータとが格納されている。

30

**【0106】**

例えば、ディスクノード100からは、図5に示したストレージ装置110のメタデータ領域110a内の論理ボリューム用スライスメタデータ51、52、53、・・・と重複排除ボリューム用スライスメタデータ61、62、63、・・・とが送信される。これらのメタデータが、ディスクノード100のディスクノードID「DP-A」に対応付けてスライスメタデータ記憶部620に格納される。

**【0107】**

スライスメタデータ記憶部620に格納された情報に基づいて、論理ボリューム管理部630によって論理ボリューム構成情報や重複排除ボリューム構成情報が作成される。

図10は、制御ノードの論理ボリューム構成情報記憶部のデータ構造例を示す図である。論理ボリューム構成情報記憶部640には、論理ボリュームごとの論理ボリュームセグメント情報テーブル641、642が格納されている。論理ボリュームセグメント情報テーブル641には、論理ボリュームIDが設定されている。また論理ボリュームセグメント情報テーブル641には、セグメントID、オフセット、およびディスクノードIDの欄が設けられている。各欄の横方向に並べられた情報が関連付けられ、各セグメントのメタデータ(セグメントメタデータ641a)となる。

40

**【0108】**

50

セグメントIDの欄には、セグメントの識別番号が設定される。オフセットの欄には、セグメントに属する記憶領域の範囲が、論理ボリュームの記憶領域の先頭からのオフセットにより示されている。具体的には、セグメントに属する記憶領域の開始位置のオフセットと終了位置のオフセットとが、オフセットの欄に設定される。ディスクノードIDの欄には、セグメントに属するデータを管理するディスクノードのディスクノードIDが設定される。

#### 【0109】

図11は、制御ノードの重複排除ボリューム構成情報記憶部のデータ構造例を示す図である。重複排除ボリューム構成情報記憶部650には、重複排除ボリュームごとの重複排除ボリュームセグメント情報テーブル651が格納されている。重複排除ボリュームセグメント情報テーブル651には、重複排除ボリュームIDが設定されている。また重複排除ボリュームセグメント情報テーブル651には、セグメントID、オフセット、およびディスクノードIDの欄が設けられている。各欄の横方向に並べられた情報が関連付けられ、各セグメントのメタデータ(セグメントメタデータ651a)となる。各欄に設定される情報は、図10に示した論理ボリュームセグメント情報テーブル641の同名の欄と同種の情報である。

10

#### 【0110】

図12は、アクセスノードの論理ボリューム構成情報記憶部のデータ構造例を示す図である。論理ボリューム構成情報記憶部710には、アクセスノード700が使用する論理ボリュームの論理ボリュームセグメント情報テーブル711が格納されている。論理ボリュームセグメント情報テーブル711には、論理ボリュームIDが設定されている。また論理ボリュームセグメント情報テーブル711には、セグメントID、オフセット、およびディスクノードIDの欄が設けられている。各欄の横方向に並べられた情報が関連付けられ、各セグメントのメタデータ(セグメントメタデータ711a)となる。各欄に設定される情報は、図10に示した論理ボリュームセグメント情報テーブル641の同名の欄と同種の情報である。

20

#### 【0111】

なお、アクセスノード700が有する論理ボリュームセグメント情報テーブル711には、例えば、過去にアクセスしたセグメントのメタデータのみを登録しておくことができる。その場合、メタデータが未登録のセグメントにアクセスする場合、データアクセス要求部720によって、該当するメタデータが制御ノード600から取得される。また、アクセスノード700の起動時またはその他の所定のタイミングで、データアクセス要求部720が、論理ボリューム内のすべてのセグメントのメタデータを、制御ノード600から取得するようにしてもよい。その場合、論理ボリュームセグメント情報テーブル711には、過去にアクセスされたか否かに拘わらず、すべてのセグメントのメタデータが登録される。

30

#### 【0112】

次に、ディスクノード100, 200, 300, 400とインデックスサーバ500とが有する機能について説明する。

図13は、ディスクノードとインデックスサーバとの機能を示すブロック図である。ディスクノード100は、データアクセス部121、アクセスユニット識別部122、データ実体領域管理部123、重複排除ユニット情報応答部131、重複排除ユニット割当部132、重複排除アドレス照会部133、重複排除ボリューム構成情報記憶部141、重複排除データアクセス要求部142、重複排除データアクセス応答部143、およびパトリール部150を有する。

40

#### 【0113】

データアクセス部121は、アクセスノード700, 800から要求された論理ボリュームへのリード/ライトのアクセスを実行する。具体的には、データアクセス部121は、アクセスノード700, 800からのデータアクセス要求を受け取ると、アクセスユニット識別部122の機能を用いて、アクセス対象のユニットの状態を識別し、状態に応じ

50

た処理を行う。例えば状態が「Normal」であれば、データアクセス部 1 2 1 は、ストレージ装置 1 1 0 に格納されているデータユニットに対してアクセスする。また状態が「DeDup」であれば、データアクセス部 1 2 1 は、重複排除データアクセス要求部 1 4 2 の機能を用いて、アクセス対象のユニットに対応する重複排除ボリューム内の重複排除ユニットに対してアクセスする。

**【 0 1 1 4 】**

なお新たなユニットへのデータ書き込みのアクセスの場合、データアクセス部 1 2 1 は、データ実体領域管理部 1 2 3 に対してデータ実体を格納するデータユニットの割り当てを依頼する。そして、データアクセス部 1 2 1 は、データ実体領域管理部 1 2 3 が新たに割り当てたデータユニットに対して、データの書き込みを行う。

10

**【 0 1 1 5 】**

アクセスユニット識別部 1 2 2 は、データアクセス部 1 2 1 やパトロール部 1 5 0 からの要求に応じ、データユニットの状態を識別する。具体的には、アクセスユニット識別部 1 2 2 は、論理ボリューム用スライスメタデータを参照し、データアクセス部 1 2 1 やパトロール部 1 5 0 から指定された論理ボリューム内のアクセス位置に相当するデータユニットの状態を判断する。そして、アクセスユニット識別部 1 2 2 は、取得した状態をデータアクセス部 1 2 1 またはパトロール部 1 5 0 に応答する。

**【 0 1 1 6 】**

データ実体領域管理部 1 2 3 は、ユニットにおけるデータ実体用領域の割り当て / 解放を管理する。またデータ実体領域管理部 1 2 3 は、ディスクノード 1 0 0 管理下のストレージ装置 1 1 0 のデータ実体領域 1 1 0 b における未使用領域の管理を行う。未使用領域の管理とは、データ実体領域 1 1 0 b 中の各単位記憶領域について、いずれかのスライスメタデータのユニットに割り当てられているか否かを管理する処理である。スライスメタデータにおいてユニットの管理情報中のデータ実体オフセットで指定された位置から 1 ユニット分の記憶領域が、該当ユニットに割り当てられているデータ実体領域 1 1 0 b 中の単位記憶領域である。例えばデータ実体領域管理部 1 2 3 は、データ実体領域 1 1 0 b 内の各単位記憶領域に対応するフラグの集合（ビットマップ）を保持し、そのスライスメタデータに割り当てられた単位記憶領域にフラグを立てる（フラグの値を「1」に設定する）。そしてデータ実体領域管理部 1 2 3 は、新たにスライスメタデータのユニットに対して単位記憶領域を割り当てると、フラグが立っていない（フラグの値が「0」）単位記憶領域を 1 つ選択して、スライスメタデータ中のユニットに割り当てる。

20

30

**【 0 1 1 7 】**

重複排除ユニット情報応答部 1 3 1 は、インデックスサーバ 5 0 0 からの重複排除ユニット情報取得要求に対して応答する。具体的には、重複排除ユニット情報応答部 1 3 1 は、インデックスサーバ 5 0 0 から重複排除ユニット情報取得要求を受け取ると、メタデータ領域 1 1 0 a から重複排除ボリューム用スライスメタデータ 6 1 , 6 2 , 6 3 , . . . を取得する。そして重複排除ユニット情報応答部 1 3 1 は、取得した重複排除ボリューム用スライスメタデータ 6 1 , 6 2 , 6 3 , . . . を、インデックスサーバ 5 0 0 に送信する。

**【 0 1 1 8 】**

重複排除ユニット割当部 1 3 2 は、インデックスサーバ 5 0 0 からの指示に基づき、重複排除ユニットへのデータ実体領域 1 1 0 b 中の単位記憶領域の割り当て指示を行う。具体的には、重複排除ユニット割当部 1 3 2 は、重複排除ボリュームのセグメントに割り当てられたスライス内の重複排除ユニットを指定したデータ実体領域割り当て要求をインデックスサーバ 5 0 0 から受け取る。すると、重複排除ユニット割当部 1 3 2 は、データ実体領域管理部 1 2 3 に対して、該当重複排除ユニットへのデータ実体領域の割り当てを依頼する。するとデータ実体領域管理部 1 2 3 により、未使用の単位記憶領域が選択され、重複排除ユニットに対して選択された単位記憶領域が割り当てられる。

40

**【 0 1 1 9 】**

重複排除アドレス照会部 1 3 3 は、パトロール部 1 5 0 からの要求に応じて、インデッ

50

クスサーバ500に対し、指定されたデータの重複排除アドレス（重複排除ボリュームのID + 重複排除オフセット）の照会を行う。具体的には、重複排除アドレス照会部133は、パトロール部150から、データのハッシュ値を受け取ると、インデックスサーバ500に対して、ハッシュ値を含む重複排除アドレス照会要求を送信する。すると、インデックスサーバ500から、ハッシュ値に対応付けられている重複排除アドレスが応答される。重複排除アドレス照会部133は、インデックスサーバ500から送られた重複排除アドレスを、パトロール部150に渡す。

#### 【0120】

重複排除ボリューム構成情報記憶部141は、重複排除ボリューム構成情報を記憶する。例えば、RAM102またはHDD103の記憶領域の一部が、重複排除ボリューム構成情報記憶部141として使用される。重複排除ボリューム構成情報記憶部141には、図11に示した重複排除ボリュームセグメント情報テーブル651と同様の情報が格納される。ディスクノード100の重複排除ボリューム構成情報記憶部141には、例えば、重複排除ボリューム内のセグメントのうち、アクセスが行われたセグメントに関するメタデータのみが登録されている。なお、アクセスが行われたか否かに関係なく、すべてのセグメントに関するメタデータが重複排除ボリューム構成情報記憶部141に登録されている場合もある。

10

#### 【0121】

重複排除データアクセス要求部142は、データアクセス部121からの要求に応じて、重複排除ユニットを管理するディスクノードに対し、その重複排除ユニットのデータを要求する。なお重複排除データアクセス要求部142は、重複排除ボリューム構成情報記憶部141を参照して、重複排除ユニットを管理するディスクノードを判断する。例えばディスクノード200に対して重複排除ユニットのデータを要求した場合、ディスクノード200の重複排除データアクセス応答部243によって、該当するデータが応答される。

20

#### 【0122】

また、重複排除データアクセス要求部142は、重複排除ボリューム構成情報記憶部141内にアクセス対象の重複排除ユニットを含むセグメントのメタデータが格納されていない場合、該当するメタデータを制御ノード600から取得する。そして重複排除データアクセス要求部142は、取得したメタデータを重複排除ボリューム構成情報記憶部141に格納する。

30

#### 【0123】

重複排除データアクセス応答部143は、他のディスクノード200, 300, 400から重複排除ユニットのデータが要求されると、ストレージ装置110から該当データを取得する。そして、重複排除データアクセス応答部143は、データを要求したディスクノードに対して取得したデータを送信する。

#### 【0124】

例えば、ディスクノード200においてディスクノード100が管理する重複排除ユニットのデータへのアクセスが生じた場合、ディスクノード200内の重複排除データアクセス要求部242からデータが要求される。重複排除データアクセス応答部143は、その要求に応じてデータを取得し、取得したデータを重複排除データアクセス要求部242に送信する。

40

#### 【0125】

パトロール部150は、ディスクノード100内で定期的にパトロールを実行する。パトロール部150は、各スライスに設定されている最終パトロール時刻から指定した時刻が経過したスライスを随時パトロールする。またパトロール部150は、未使用な単位記憶領域も随時パトロールする。すなわちパトロールは、ディスクノード100管理下のストレージ装置110内の情報（各種メタデータ/データユニット/重複排除ユニット/未使用領域）が対象となる。具体的にはパトロール部150は、パトロール処理において、データユニットに関しては、ディスク保守的な読み書きを行う。ここで、ディスク保守的

50

な読み書きとは、対象データを読み出してデータのエラーの有無を判断し、エラーがなければ読み出したデータを元の位置に書き込む処理である。

【 0 1 2 6 】

またパトロール部 1 5 0 は、アクセスユニット識別部 1 2 2 を用いて、パトロールを実行するユニットの状態を識別し、それに応じた処理を行う。パトロール対象のユニットへの最後の書き込みから所定時間以上経過していれば、そのユニットを重複排除化対象と判断する。その場合、パトロール部 1 5 0 は、重複排除アドレス照会部 1 3 3 を用いて、重複対象ユニットの重複排除アドレスを入手する。

【 0 1 2 7 】

パトロール部 1 5 0 は、重複排除ユニットのパトロールに関しては、ディスク保守的な読み書きを行う。またパトロール部 1 5 0 は、保存期限満了時刻が過ぎている重複排除ユニットがある場合、重複排除ボリューム用スライスメタデータからそのユニットの情報をクリアする。さらにパトロール部 1 5 0 は、各種メタデータ、および未使用領域のパトロールでは、ディスク保守的な読み書きを行う。

【 0 1 2 8 】

インデックスサーバ 5 0 0 は、重複排除ユニット情報収集部 5 1 1、重複排除ユニット情報記憶部 5 1 2、重複排除ユニット割当要求部 5 1 3、および重複排除アドレス検索部 5 1 4 を有する。

【 0 1 2 9 】

重複排除ユニット情報収集部 5 1 1 は、各ディスクノード 1 0 0 , 2 0 0 , 3 0 0 , 4 0 0 から重複排除ユニットの情報を収集する。例えば、ディスクノード 1 0 0 から重複排除ユニットの情報を収集する場合、重複排除ユニット情報収集部 5 1 1 は、ディスクノード 1 0 0 の重複排除ユニット情報応答部 1 3 1 に対して、重複排除ユニット情報取得要求を送信する。すると重複排除ユニット情報応答部 1 3 1 から、ディスクノード 1 0 0 が管理するストレージ装置 1 1 0 内の重複排除ボリューム用スライスメタデータ 6 1 , 6 2 , 6 3 , . . . が応答される。重複排除ユニット情報収集部 5 1 1 は、収集した重複排除ユニットの情報（重複排除ボリューム用スライスメタデータ）を重複排除ユニット情報記憶部 5 1 2 に格納する。

【 0 1 3 0 】

重複排除ユニット情報記憶部 5 1 2 は、重複排除ユニット情報収集部 5 1 1 が収集した重複排除ユニットの情報（重複排除ボリューム用スライスメタデータ）を記憶する記憶機能である。例えば、インデックスサーバ 5 0 0 の R A M や H D D の記憶領域の一部が、重複排除ユニット情報記憶部 5 1 2 として使用される。なお重複排除ユニット情報記憶部 5 1 2 には、使用中の重複排除ユニットだけでなく未使用重複排除ユニットの情報も格納される。

【 0 1 3 1 】

重複排除ユニット割当要求部 5 1 3 は、重複排除アドレス検索部 5 1 4 からの要求に応じて、ディスクノード 1 0 0 , 2 0 0 , 3 0 0 , 4 0 0 に対して、指定した重複排除ユニットへの単位記憶領域の割り当てを要求する。具体的には重複排除ユニット割当要求部 5 1 3 は、重複排除アドレス検索部 5 1 4 が選択した重複排除ユニットが属するスライス管理するディスクノードに対して、重複排除ユニットへの単位記憶領域の割り当てを要求する。ディスクノードから割り当て完了の応答を受け取ると、重複排除ユニット割当要求部 5 1 3 は、割り当てが完了したことを重複排除アドレス検索部 5 1 4 に通知する。

【 0 1 3 2 】

重複排除アドレス検索部 5 1 4 は、ディスクノード 1 0 0 , 2 0 0 , 3 0 0 , 4 0 0 からの重複排除アドレス照会に対して応答する。例えば重複排除アドレス検索部 5 1 4 はディスクノード 1 0 0 の重複排除アドレス照会部 1 3 3 からハッシュ値を含む重複排除アドレス照会要求を取得すると、そのハッシュ値に対応する重複排除ユニット情報を、重複排除ユニット情報記憶部 5 1 2 から検索する。重複排除アドレス検索部 5 1 4 は、検索により重複排除ユニット情報がヒットすれば、ヒットした情報の重複排除アドレス（重複排除

10

20

30

40

50



ボリュームID + オフセット)を、重複排除アドレス照会部133に回答する。検索によりヒットする情報がなければ、重複排除アドレス検索部514は、重複排除ユニット情報記憶部512を参照し、未使用の重複排除ユニットを選択する。次に重複排除アドレス検索部514は、重複排除ユニット割当要求部513に対して、選択した重複排除ユニットへの単位記憶領域の割り当てを要求する。そして重複排除アドレス検索部514は、新たに割り当てられた重複排除アドレスを、ディスクノード100の重複排除アドレス照会部133に回答する。この際、重複排除アドレス検索部514は、重複排除ユニット情報記憶部512内の選択した重複排除ユニットの重複排除ユニットデータを更新する。

【0133】

なお、図13には複数のディスクノード100, 200, 300, 400のうち、ディスクノード100の機能のみを詳細に示しているが、他のディスクノード200, 300, 400もディスクノード100と同様の機能を有している。

10

【0134】

図14は、ディスクノードが有する重複排除ボリューム構成情報記憶部のデータ構造例を示す図である。重複排除ボリューム構成情報記憶部141には、重複排除ボリュームセグメント情報テーブル141aが格納されている。重複排除ボリュームセグメント情報テーブル141aのデータ構造は、図11に示した制御ノード600が有する重複排除ボリューム構成情報記憶部650内の重複排除ボリュームセグメント情報テーブル651と同様である。重複排除ボリュームセグメント情報テーブル141aには、例えば、ディスクノードが過去にアクセスしたセグメントのメタデータ(セグメントメタデータ141b)のみが登録されている。

20

【0135】

以上のような構成において、システム起動時には、重複排除ユニット情報収集部511によって、重複排除ユニット情報の収集が行われる。重複排除ユニット情報は、各ストレージ装置110, 210, 310, 410のメタデータ領域に保存されている。

【0136】

図15は、重複排除ユニット情報収集処理を示すシーケンス図である。以下、図15に示す処理をステップ番号に沿って説明する。なお、図15には2台のディスクノード100, 200とインデックスサーバ500との間の処理を示しているが、他のディスクノード300, 400とインデックスサーバ500の間でも同様の処理が実行される。

30

【0137】

[ステップS11]インデックスサーバ500が起動されると、重複排除ユニット情報収集部511は、すべてのディスクノード100, 200, 300, 400に対して、重複排除ユニット情報を要求する。

【0138】

[ステップS12]ディスクノード100の重複排除ユニット情報応答部131は、重複排除ユニット情報をインデックスサーバ500に回答する。具体的には重複排除ユニット情報応答部131は、ストレージ装置110のメタデータ領域110aから、重複排除ユニット情報を含む重複排除ボリューム用スライスメタデータ61, 62, 63, ...を取得する。そして重複排除ユニット情報応答部131は、取得した重複排除ボリューム用スライスメタデータ61, 62, 63, ...をインデックスサーバ500に送信する。

40

【0139】

[ステップS13]ディスクノード200もディスクノード100と同様に、重複排除ユニット情報を含む重複排除ボリューム用スライスメタデータをインデックスサーバ500に回答する。なお、図示していない他のディスクノード300, 400からも、インデックスサーバ500へ重複排除ボリューム用スライスメタデータが回答される。

【0140】

[ステップS14]インデックスサーバ500の重複排除ユニット情報収集部511は、ディスクノード100, 200, 300, 400から送られた重複排除ユニット情報を

50

、重複排除ユニット情報記憶部 5 1 2 に格納する。

【 0 1 4 1 】

このようにして、インデックスサーバ 5 0 0 において、重複排除ユニットの情報が収集され、重複排除ユニット情報記憶部 5 1 2 に格納される。

図 1 6 は、重複排除ユニット情報記憶部のデータ構造例を示す図である。重複排除ユニット情報記憶部 5 1 2 には、各ディスクノードから収集した重複排除ボリューム用スライスメタデータが、収集元のディスクノードのディスクノード ID に対応付けて格納されている。重複排除ボリューム用スライスメタデータの内容は、図 7 に示す通りである。すなわち、重複排除ユニット情報には、重複排除ボリューム ID が含まれると共に、重複排除ユニットごとの、データ実体オフセット、ハッシュ値、保存期限満了時刻が含まれる。

10

【 0 1 4 2 】

このような重複排除ユニット情報記憶部 5 1 2 内の重複排除ボリューム用スライスメタデータから、ハッシュ値を検索キーとして重複排除ユニット情報を検索することで、そのハッシュ値の生成元となったデータが重複排除対象となっているか否かを判断できる。検索を行う重複排除アドレス検索部 5 1 4 は、例えばハッシュテーブルとツリー構造を組み合わせた構成で、ハッシュ値の検索を行うことができる。

【 0 1 4 3 】

また、重複排除ユニット情報記憶部 5 1 2 を参照するところで、未使用の重複排除ユニットを検出できる。未使用の重複排除ユニットには、重複排除ボリューム用スライスメタデータにおいて「NULL」のデータが設定されている。重複排除ユニット割当要求部 5 1 3 は、重複排除ユニット情報記憶部 5 1 2 から未使用の重複排除ユニットを検索し、見つけた重複排除ユニットの新たな重複排除対象データへの割り当て要求処理を行う。なお、重複排除ユニット割当要求部 5 1 3 は、未使用の重複排除ユニットについて、例えばキュー管理することができる。すなわち、重複排除ユニット割当要求部 5 1 3 は、未使用の重複排除ユニットを示す重複排除ユニット情報へのポインタをキューにキューイングしておく。そして重複排除ユニット割当要求部 5 1 3 は、新たに重複排除ユニットの重複排除対象データへの割り当てが必要となると、キューの先頭のポインタで示される重複排除ユニットを割り当ての対象として選択する。また重複排除ユニット割当要求部 5 1 3 は、使用していた重複排除ユニットが未使用になると、その重複排除ユニットに関する重複排除ユニット情報へのポインタを、キューイングする。

20

30

【 0 1 4 4 】

なお、システム起動時には、制御ノード 6 0 0 のスライスメタデータ収集部 6 1 0 により、論理ボリューム用スライスメタデータと重複排除ボリューム用スライスメタデータとが収集される。そして、論理ボリューム管理部 6 3 0 によって、論理ボリューム構成情報と重複排除ボリューム構成情報とが作成される。

【 0 1 4 5 】

システム起動直後は、アクセスノード 7 0 0 , 8 0 0 が論理ボリュームにアクセスする際には、アクセス対象のデータを含むセグメントのセグメントメタデータを制御ノード 6 0 0 から取得する。そして、アクセスノード 7 0 0 , 8 0 0 は、取得したセグメントメタデータに基づいて、アクセス対象のデータを管理しているディスクノードのディスクノード ID を判断する。

40

【 0 1 4 6 】

同様にシステム起動直後は、ディスクノード 1 0 0 , 2 0 0 , 3 0 0 , 4 0 0 から他のディスクノードが管理する重複排除対象データにアクセスする場合、アクセス対象のデータを含むセグメントのセグメントメタデータを制御ノード 6 0 0 から取得する。そして、ディスクノード 1 0 0 , 2 0 0 , 3 0 0 , 4 0 0 は、取得したセグメントメタデータに基づいて、アクセス対象のデータを管理しているディスクノードのディスクノード ID を判断する。

【 0 1 4 7 】

< 読み出し処理 >

50

次に、読み出し要求時の処理について説明する。

図 17 は、重複排除ボリュームからのデータ読み出し処理の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。

【 0 1 4 8 】

[ ステップ S 2 1 ] アクセスノード 7 0 0 のデータアクセス要求部 7 2 0 から読み出し要求がディスクノード 1 0 0 に対して送信される。読み出し要求には、例えば、論理ボリューム ID、論理ボリュームの先頭からのアクセス対象データのオフセット、およびデータサイズが含まれる。

【 0 1 4 9 】

[ ステップ S 2 2 ] ディスクノード 1 0 0 では、データアクセス部 1 2 1 が読み出し要求を受信する。するとデータアクセス部 1 2 1 は、アクセスユニット識別部 1 2 2 に対してアクセス対象のユニットの状態の識別を依頼する。アクセスユニット識別部 1 2 2 は、論理ボリューム用スライスメタデータを参照し、アクセス対象データを含むユニットの状態を判断する。

【 0 1 5 0 】

例えばアクセスユニット識別部 1 2 2 は、読み出し要求に示される論理ボリューム ID で示される論理ボリュームに割り当てられたスライスの論理ボリューム用スライスメタデータを抽出する。論理ボリューム用スライスメタデータにはセグメント ID が含まれている。1 セグメント当たりの記憶容量は一定（例えば 1 GB）である。そのため、セグメント ID によって、該当スライスが割り当てられたセグメントの記憶領域の範囲（開始位置のオフセットと終了位置のオフセット）が判断できる。そこでアクセスユニット識別部 1 2 2 は、読み出し要求内のオフセットとデータサイズで示される領域を包含するセグメントに割り当てられたスライスの論理ボリューム用スライスメタデータを抽出する。さらにアクセスユニット識別部 1 2 2 は、抽出した論理ボリューム用スライスメタデータ内のデータユニット情報から、アクセス対象データが属するユニットのデータユニット情報を検索する。そして、アクセスユニット識別部 1 2 2 は、検索で該当したデータユニット情報を選択し、そのデータユニット情報の状態情報に基づきユニットの状態を判断する。判断結果は、アクセスユニット識別部 1 2 2 からデータアクセス部 1 2 1 に通知される。

【 0 1 5 1 】

図 17 の例では、状態が「DeDup」であったものとする。すなわち、アクセス対象データは、すでに重複排除対象データとなり、重複排除ボリュームで管理されている。状態が「DeDup」のデータユニット情報には、重複排除ボリューム ID と重複排除オフセットとが含まれている。状態が「DeDup」の場合、データアクセス部 1 2 1 から重複排除データアクセス要求部 1 4 2 に対して、重複排除ユニットへのアクセスが要求される。

【 0 1 5 2 】

[ ステップ S 2 3 ] 重複排除データアクセス要求部 1 4 2 は、読み出し要求に対応するデータユニット情報の重複排除ボリューム ID と重複排除オフセットとに基づいて、重複排除ボリュームセグメント情報テーブル 1 4 1 a（図 14 参照）からセグメントメタデータを検索する。具体的には、重複排除データアクセス要求部 1 4 2 は、選択したデータユニット情報の重複排除ボリューム ID が設定された重複排除ボリュームセグメント情報テーブル 1 4 1 a を、重複排除ボリューム構成情報記憶部 1 4 1 から抽出する。次に、重複排除データアクセス要求部 1 4 2 は、抽出した重複排除ボリュームセグメント情報テーブル 1 4 1 a から、選択したデータユニット情報の重複排除オフセットを包含するオフセットが設定されたセグメントメタデータを検索する。該当するセグメントメタデータが存在する場合、処理はステップ S 2 7 に進められる。該当するセグメントメタデータが存在しない場合、処理がステップ S 2 4 に進められる。

【 0 1 5 3 】

[ ステップ S 2 4 ] 重複排除データアクセス要求部 1 4 2 は、メタデータ要求を制御ノード 6 0 0 に送信する。メタデータ要求には、例えば、選択したデータユニット情報の重複排除オフセットが含まれる。

## 【 0 1 5 4 】

[ステップS 2 5] 制御ノード6 0 0の論理ボリューム管理部6 3 0は、セグメントメタデータを応答する。具体的には、制御ノード6 0 0は、重複排除ボリューム構成情報記憶部6 5 0内の重複排除ボリュームセグメント情報テーブル6 5 1から、メタデータ要求で示されるセグメントのセグメントメタデータを抽出する。例えば、メタデータ要求において重複排除オフセットが示されている場合、その値を包含するオフセットが設定されたセグメントメタデータが抽出される。そして、論理ボリューム管理部6 3 0は抽出したセグメントメタデータを、ディスクノード1 0 0に送信する。

## 【 0 1 5 5 】

[ステップS 2 6] ディスクノード1 0 0の重複排除データアクセス要求部1 4 2は、10  
重複排除ボリューム構成情報記憶部1 4 1内の重複排除ボリュームセグメント情報テーブル1 4 1 aに、取得したセグメントメタデータを追加登録する。

## 【 0 1 5 6 】

[ステップS 2 7] 重複排除データアクセス要求部1 4 2は、重複排除対象データの読み出し要求を、そのデータを管理しているディスクノード(図1 7の例では、ディスクノード2 0 0)に対して送信する。具体的には重複排除データアクセス要求部1 4 2は、ステップS 2 3~ステップS 2 6の処理で取得したセグメントメタデータに基づいて、重複排除対象データを管理するディスクノードを判断する。そして、重複排除データアクセス要求部1 4 2は、重複排除対象データが属する重複排除ボリュームの重複排除ボリュームID、および重複排除オフセットを指定した読み出し要求を、重複排除対象データを管理20  
するディスクノード2 0 0に送信する。なお、重複排除ボリュームID、および重複排除オフセットは、ステップS 2 2で選択したデータユニット情報から取得される。

## 【 0 1 5 7 】

[ステップS 2 8] ディスクノード2 0 0の重複排除データアクセス応答部2 4 3は、重複排除対象データの読み出し要求を受信すると、該当する重複排除対象データをディスクノード1 0 0に応答する。具体的には、重複排除データアクセス応答部2 4 3は、読み出し要求で指定された重複排除ボリュームIDが設定された重複排除ボリューム用スライスマタデータを、ストレージ装置2 1 0から取得する。次に重複排除データアクセス応答部2 4 3は、取得した重複排除ボリューム用スライスマタデータから、読み出し要求に示される重複排除オフセットに示されるユニットの重複排除ユニット情報を抽出する。さら30  
に、重複排除データアクセス応答部2 4 3は、抽出した重複排除ユニット情報に示されるデータ実体オフセットで示されるデータ実体領域内の単位記憶領域から、重複排除対象データを読み出す。そして、重複排除データアクセス応答部2 4 3は、読み出した重複排除対象データをディスクノード1 0 0に送信する。

## 【 0 1 5 8 】

[ステップS 2 9] ディスクノード1 0 0の重複排除データアクセス要求部1 4 2は、ディスクノード2 0 0から取得した重複排除対象データをデータアクセス部1 2 1に渡す。データアクセス部1 2 1は、受け取った重複排除対象データを、読み取り要求のアクセス対象データとして、アクセスノード7 0 0に送信する。

【 0 1 5 9 】40

このようにして、データの読み出し処理が実行される。データの読み出し処理において、読み出し対象のデータを含むユニットの状態が「DeDup」であれば、対応する重複排除対象データを管理するディスクノード2 0 0からデータが読み出される。なお図1 7の例では、アクセスノード7 0 0からの読み出し要求を受けたディスクノード1 0 0とは別のディスクノード2 0 0から重複排除対象データが読み出されているが、ディスクノード1 0 0において重複排除対象データを管理している場合もある。その場合、重複排除対象データアクセス要求部1 4 2は、ストレージ装置1 1 0から重複排除対象データを読み出す。

## 【 0 1 6 0 】

また、図1 7の例では、読み出し対象のデータを含むユニットの状態が「DeDup」の場50

合を示している。該当ユニットの状態が「Blank」の場合には、データアクセス部 1 2 1 はアクセスノード 7 0 0 に対して「0」の値を送信する。

【0 1 6 1】

また、該当ユニットの状態が「Normal」の場合には、データアクセス部 1 2 1 は、データユニット情報のデータ実体オフセットに基づいて、そのデータユニットに割り当てられている単位記憶領域を判断する。そしてデータアクセス部 1 2 1 は、その単位記憶領域からデータを読み出し、アクセスノード 7 0 0 に送信する。

【0 1 6 2】

<書き込み処理>

次にデータの書き込み処理について説明する。

図 1 8 は、状態が「Blank」のデータユニットへの書き込み処理を示すシーケンス図である。以下、図 1 8 に示す処理をステップ番号に沿って説明する。

【0 1 6 3】

[ステップ S 3 1] アクセスノード 7 0 0 のデータアクセス要求部 7 2 0 から書き込み要求がディスクノード 1 0 0 に対して送信される。書き込み要求には、例えば、論理ボリューム ID、論理ボリュームの先頭からのアクセス対象データのオフセット、および書き込むべきデータが含まれる。

【0 1 6 4】

[ステップ S 3 2] ディスクノード 1 0 0 では、データアクセス部 1 2 1 が書き込み要求を受信する。するとデータアクセス部 1 2 1 は、アクセスユニット識別部 1 2 2 にアクセス対象データを含むユニットの状態の判断を依頼し、判断結果を受け取る。状態判断処理の詳細は、図 1 7 のステップ S 2 2 の処理と同様である。

【0 1 6 5】

図 1 8 の例では、状態が「Blank」であったものとする。

[ステップ S 3 3] データアクセス部 1 2 1 は、データ実体領域管理部 1 2 3 に対してデータ実体領域内の単位記憶領域の割り当てを依頼する。するとデータ実体領域管理部 1 2 3 が、ストレージ装置 1 1 0 のデータ実体領域 1 1 0 b から未使用の単位記憶領域を選択し、アクセス対象のユニットに割り当てる。すなわち、論理ボリューム用スライスメタデータ内の該当ユニットのデータユニット情報に、選択した単位記憶領域のデータ実体オフセットを設定する。

【0 1 6 6】

[ステップ S 3 4] データアクセス部 1 2 1 は、アクセス対象のデータユニットのデータユニット情報の状態を「Normal」に変更する。

[ステップ S 3 5] データアクセス部 1 2 1 は、データ実体領域 1 1 0 b 内のアクセス対象のデータユニットに割り当てられた単位記憶領域内に、データを書き込む。

【0 1 6 7】

[ステップ S 3 6] データアクセス部 1 2 1 は、現在の時刻を、アクセス対象のデータユニットのデータユニット情報に最終書き込み時刻として設定する。

[ステップ S 3 7] データアクセス部 1 2 1 は、書き込み完了応答をアクセスノード 7 0 0 に送信する。

【0 1 6 8】

このようにして、状態が「Blank」のデータユニットへの書き込みが行われる。

図 1 9 は、状態が「DeDup」のデータユニットの一部への書き込み処理を示すシーケンス図である。以下、図 1 9 に示す処理をステップ番号に沿って説明する。

【0 1 6 9】

[ステップ S 4 1] アクセスノード 7 0 0 のデータアクセス要求部 7 2 0 から書き込み要求がディスクノード 1 0 0 に対して送信される。

[ステップ S 4 2] ディスクノード 1 0 0 では、データアクセス部 1 2 1 が書き込み要求を受信する。するとデータアクセス部 1 2 1 は、アクセスユニット識別部 1 2 2 にアクセス対象データを含むユニットの状態の判断を依頼し、判断結果を受け取る。

10

20

30

40

50

## 【 0 1 7 0 】

図 1 9 の例では、状態が「DeDup」であったものとする。また、データの書き込み領域は、ユニットの一部（全体でない）であるものとする。

【ステップ S 4 3】データアクセス部 1 2 1 は、重複排除データアクセス要求部 1 4 2 に対して重複排除対象データの読み出し要求を、そのデータを管理しているディスクノード（図 1 9 の例では、ディスクノード 2 0 0）に対して送信する。この処理の詳細は、図 1 7 のステップ S 2 7 の処理と同様である。

## 【 0 1 7 1 】

【ステップ S 4 4】ディスクノード 2 0 0 の重複排除データアクセス応答部 2 4 3 は、重複排除対象データの読み出し要求を受信すると、該当する重複排除対象データをディスクノード 1 0 0 に応答する。この処理の詳細は、図 1 7 のステップ S 2 8 の処理と同様である。

10

## 【 0 1 7 2 】

【ステップ S 4 5】ディスクノード 1 0 0 の重複排除データアクセス要求部 1 4 2 は、ディスクノード 2 0 0 から取得した重複排除対象データをデータアクセス部 1 2 1 に渡す。ディスクノード 1 0 0 のデータアクセス部 1 2 1 は、データ実体領域管理部 1 2 3 に対してデータ実体領域内の単位記憶領域の割り当てを依頼する。するとデータ実体領域管理部 1 2 3 が、ストレージ装置 1 1 0 のデータ実体領域 1 1 0 b から未使用の単位記憶領域を選択し、アクセス対象のユニットに割り当てる。

## 【 0 1 7 3 】

【ステップ S 4 6】データアクセス部 1 2 1 は、アクセス対象のデータユニットのデータユニット情報の状態を「Normal」に変更する。

【ステップ S 4 7】データアクセス部 1 2 1 は、取得した重複排除対象データを書き込み要求に応じて更新し、データ実体領域 1 1 0 b 内のアクセス対象のデータユニットに割り当てられた単位記憶領域内に、更新後のデータ実体を書き込む。

20

## 【 0 1 7 4 】

【ステップ S 4 8】データアクセス部 1 2 1 は、現在の時刻を、アクセス対象のデータユニットのデータユニット情報に最終書き込み時刻として設定する。

【ステップ S 4 9】データアクセス部 1 2 1 は、書き込み完了応答をアクセスノード 7 0 0 に送信する。

30

## 【 0 1 7 5 】

このようにして、状態が「DeDup」のデータユニットへの書き込みが行われる。

図 1 8 , 図 1 9 に示した状態が「Blank」または「DeDup」のデータユニットへの書き込み処理以外に、状態が「Normal」のデータユニットへの書き込み処理がある。状態が「Normal」のデータユニットへの書き込み処理では、既にデータユニットに割り当てられているデータ実体領域内の単位記憶領域に対して、データの書き込みが行われる。

## 【 0 1 7 6 】

また状態が「DeDup」のデータユニットへの書き込みの場合、データユニット内の全データの更新書き込みが行われる場合がある。その場合、既存の重複排除対象データの読み出しを行う必要はない。書き込み処理におけるこれらの条件判断を含むディスクノードでの書き込み処理を以下に示す。

40

## 【 0 1 7 7 】

図 2 0 は、ディスクノードにおける書き込み処理を示すフローチャートである。以下、図 2 0 に示す処理をステップ番号に沿って説明する。この処理は、アクセスノード 7 0 0 からの書き込み要求を受信した際に、ディスクノード 1 0 0 で実行される。

## 【 0 1 7 8 】

【ステップ S 5 1】データアクセス部 1 2 1 は、アクセスユニット識別部 1 2 2 の機能を用い、アクセス対象のデータユニットの状態を判断する。状態が「Normal」であれば、処理がステップ S 5 5 に進められる。状態が「Blank」であれば、処理がステップ S 5 3 に進められる。状態が「DeDup」であれば、処理がステップ S 5 2 に進められる。

50

## 【 0 1 7 9 】

[ステップS 5 2] データアクセス部 1 2 1 は、書き込みデータがユニット全体に対するデータか否かを判断する。ユニット全体に対するデータであれば、処理がステップS 5 3に進められる。ユニット全体に対するデータでなければ、処理がステップS 5 6に進められる。

## 【 0 1 8 0 】

[ステップS 5 3] データアクセス部 1 2 1 は、データ実体領域管理部 1 2 3 に対してアクセス対象のデータユニットに対する単位記憶領域の割り当てを依頼する。するとデータ実体領域管理部 1 2 3 がデータ実体領域 1 1 0 b から未使用の単位記憶領域を選択し、データユニットに割り当てる。

10

## 【 0 1 8 1 】

[ステップS 5 4] データアクセス部 1 2 1 は、アクセス対象のデータユニットのデータユニット情報の状態を「Normal」に変更する。

[ステップS 5 5] データアクセス部 1 2 1 は、書き込み要求に含まれるデータを、データ実体領域 1 1 0 b 内のアクセス対象のデータユニットに割り当てられた単位記憶領域内に書き込む。その後、処理がステップS 6 0に進められる。

## 【 0 1 8 2 】

[ステップS 5 6] データアクセス部 1 2 1 は、重複排除データアクセス要求部 1 4 2 を介してデータユニット全体の重複排除対象データを取得する。

[ステップS 5 7] データアクセス部 1 2 1 は、データ実体領域管理部 1 2 3 に対してアクセス対象のデータユニットに対する単位記憶領域の割り当てを依頼する。するとデータ実体領域管理部 1 2 3 がデータ実体領域 1 1 0 b から未使用の単位記憶領域を選択し、データユニットに割り当てる。

20

## 【 0 1 8 3 】

[ステップS 5 8] データアクセス部 1 2 1 は、アクセス対象のデータユニットのデータユニット情報の状態を「Normal」に変更する。

[ステップS 5 9] データアクセス部 1 2 1 は、取得した重複排除対象データを書き込み要求に応じて更新する。そしてデータアクセス部 1 2 1 は、データ実体領域 1 1 0 b 内のアクセス対象のデータユニットに割り当てられた単位記憶領域内に、更新後のデータ実体を書き込む。

30

## 【 0 1 8 4 】

[ステップS 6 0] データアクセス部 1 2 1 は、現在の時刻を、アクセス対象のデータユニットのデータユニット情報に最終書き込み時刻として設定する。

[ステップS 6 1] データアクセス部 1 2 1 は、書き込み完了応答をアクセスノード 7 0 0 に送信する。

## 【 0 1 8 5 】

以上のようにして、アクセス対象のデータユニットの状態に応じた適切な書き込み処理が実行される。

## &lt; パトロール処理 &gt;

次に、パトロール処理について詳細に説明する。パトロール処理では、論理ボリューム内のデータユニットから重複排除対象データを検出し、その重複排除対象データを重複排除ボリュームに移動する処理が行われる。さらにパトロール処理では、重複排除ボリューム内の重複排除ユニットから、論理ボリューム内のデータユニットとの関連性を失った重複排除ユニットを検出し、その重複排除対象データを削除する。なおパトロール処理は、スライス内のデータユニット、およびデータ実体領域 1 1 0 b 内の未使用の単位記憶領域に対して、所定のパトロール時間ごとに定期的に行われる。

40

## 【 0 1 8 6 】

図 2 1 は、状態が「Normal」のデータユニットへのパトロール処理を示すシーケンス図である。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。なお、図 2 1 のシーケンス図は、スライス内のデータユニットに対するパトロールの例である。

50

## 【 0 1 8 7 】

[ステップS 7 1] パトロール部 1 5 0 は、パトロール対象のデータユニットに割り当てられた単位記憶領域からデータ実体を読み出す。

[ステップS 7 2] パトロール部 1 5 0 は、最終書き込み時刻からの経過時間を判定する。具体的には、パトロール部 1 5 0 は、パトロール対象のユニットのデータユニット情報に設定されている最終書き込み時刻から、現在時刻までの経過時間を算出する。そしてパトロール部 1 5 0 は、経過時間が予め設定されている所定の重複排除化猶予期間を超えているか否かを判断する。最終書き込みから重複排除化猶予期間を経過している場合、パトロール対象のデータユニットへのデータの書き込み頻度が低いことを示す。その場合、パトロール部 1 5 0 は、パトロール対象のデータユニット内のデータを重複排除対象データとする。図 2 1 の例では、経過時間が重複排除化猶予期間以上経過しているものとする。

10

## 【 0 1 8 8 】

[ステップS 7 3] パトロール部 1 5 0 は、所定のハッシュ関数を用いて、ステップS 7 1 で読み出したデータのハッシュ値を算出する。

[ステップS 7 4] パトロール部 1 5 0 は、インデックスサーバ 5 0 0 に対して重複排除アドレス照会を行う。具体的には、パトロール部 1 5 0 は、重複排除アドレス照会部 1 3 3 に対して、ステップS 7 3 で算出したハッシュ値を指定した重複排除アドレスの照会を依頼する。すると、重複排除アドレス照会部 1 3 3 は、インデックスサーバ 5 0 0 に対して、ハッシュ値を含む重複排除アドレス照会要求を送信する。

20

## 【 0 1 8 9 】

[ステップS 7 5] インデックスサーバ 5 0 0 の重複排除アドレス検索部 5 1 4 は、重複排除アドレス照会要求に示されるハッシュ値が設定された重複排除ユニット情報を、重複排除ユニット情報記憶部 5 1 2 から検索する。図 2 1 の例では該当する重複排除ユニット情報が未登録であり、検索により何もヒットしないものとする。

## 【 0 1 9 0 】

[ステップS 7 6] インデックスサーバ 5 0 0 の重複排除アドレス検索部 5 1 4 は重複排除ユニット情報記憶部 5 1 2 から未使用の重複排除ユニットを選択し、重複排除ユニット割当要求部 5 1 3 に対して、選択した重複排除ユニットへの単位記憶領域の割当を依頼する。すると、重複排除ユニット割当要求部 5 1 3 は、選択された重複排除ユニットが属するスライスを管理するディスクノード(図 2 1 の例ではディスクノード 2 0 0)に対して、単位記憶領域の割当を要求する。

30

## 【 0 1 9 1 】

[ステップS 7 7] ディスクノード 2 0 0 は、インデックスサーバ 5 0 0 からの要求に応じて重複排除ユニットに対して単位記憶領域を割り当てる。

[ステップS 7 8] ディスクノード 2 0 0 は、重複排除ユニットへの単位記憶領域の割当完了を応答する。

## 【 0 1 9 2 】

[ステップS 7 9] インデックスサーバ 5 0 0 の重複排除ユニット割当要求部 5 1 3 は、割当が完了したことを重複排除アドレス検索部 5 1 4 に伝える。すると重複排除アドレス検索部 5 1 4 は、重複排除ユニット情報記憶部 5 1 2 内の新たに単位記憶領域を割り当てた重複排除ユニットの重複排除ユニット情報を更新する。更新された重複排除ユニット情報には、少なくともディスクノード 1 0 0 から送られたハッシュ値が設定される。

40

## 【 0 1 9 3 】

[ステップS 8 0] 重複排除アドレス検索部 5 1 4 は、重複排除アドレス(論理ボリューム ID と重複排除オフセット)をディスクノード 1 0 0 に応答する。この際、重複排除アドレス検索部 5 1 4 は、新たに重複排除ユニット割当処理(ステップS 7 6 ~ S 7 9)が行われたことを示す情報を、重複排除アドレスに付加する。

## 【 0 1 9 4 】

[ステップS 8 1] ディスクノード 1 0 0 の重複排除アドレス照会部 1 3 3 は、インデ

50



ックスサーバ500から送られた重複排除アドレスをパトロール部150に渡す。パトロール部150は、重複排除ユニット割当処理が行われたことを認識し、重複排除データアクセス要求部142に対して、取得した重複排除アドレスを指定して重複排除対象データの書き込みを要求する。すると、重複排除データアクセス要求部142が、ディスクノード200に対して、重複排除オフセットを書き込み位置として指定して、重複排除対象データを送信する。

【0195】

[ステップS82] ディスクノード200の重複排除データアクセス応答部243は、ディスクノード100から送られた重複排除対象データを、指定された重複排除オフセットに対応する重複排除ユニットに割り当てられた単位記憶領域に書き込む。この時、重複排除データアクセス応答部243は、重複排除ユニットの重複排除ユニット情報に、現在の時刻から所定の保存期間経過後の保存期限満了時刻を設定する。

10

【0196】

[ステップS83] ディスクノード200の重複排除データアクセス応答部243は、ディスクノード100に対して書き込み完了を応答する。

[ステップS84] ディスクノード100の重複排除データアクセス要求部142は、データ書き込みの完了をパトロール部150に伝える。すると、パトロール部150は、重複排除対象データが格納されていたデータユニットのデータユニット情報における状態を、「DeDup」に変更する。

【0197】

20

[ステップS85] パトロール部150は、重複排除対象データが格納されていたデータユニットに割り当てられていたデータ実体領域110b内の単位記憶領域を解放する。すなわち、そのデータユニットのデータユニット情報に設定されていた最終書き込み時刻とデータ実体オフセットが削除され、代わりに重複排除アドレス(重複排除ボリュームIDと重複排除オフセット)が設定される。

【0198】

[ステップS86] パトロール部150は、重複排除対象データが格納されていたデータユニットに割り当てられていたデータ実体領域110b内の単位記憶領域に対して、初期データを書き込む。

【0199】

30

このようにして、重複排除化猶予期間以上書き込みが行われていないデータユニット内のデータは、重複排除対象データとして重複排除ボリュームに移される。

なお図21の例では、ステップS75の処理で該当する重複排除ユニット情報が未登録であると判断されたため、ステップS76～ステップS79の重複排除ユニット割当処理が実行されている。そして重複排除ユニット割当処理が実行されたことにより、ステップS81～ステップS82のデータ転送処理が行われている。他方、ステップS75の検索で重複排除ユニット情報がヒットし、重複排除ユニット割当処理が行われない場合には、ディスクノード100からディスクノード200へのデータ転送も行われない。

【0200】

図22は、状態が「DeDup」のデータユニットへのパトロール処理を示すシーケンス図である。以下、図22に示す処理をステップ番号に沿って説明する。

40

[ステップS91] ディスクノード100のパトロール部150は、パトロール対象の単位記憶領域が割り当てられているデータユニットの状態が「DeDup」の場合、重複排除データアクセス要求部142に対して重複排除対象データの参照を要求する。すると重複排除データアクセス要求部142は、重複排除対象データを管理するディスクノード200に対して重複排除対象データの参照要求を送信する。

【0201】

[ステップS92] ディスクノード200の重複排除データアクセス応答部243は、参照要求に応じて重複排除対象データを取得する。この際、重複排除データアクセス応答部243は、重複排除対象データが格納された重複排除ユニットの保存期限満了時刻を、

50

現在の時刻から所定の保存期間経過後の保存期限満了時刻に更新する。

【 0 2 0 2 】

[ ステップ S 9 3 ] 重複排除データアクセス応答部 2 4 3 は、ディスクノード 1 0 0 に対して取得した重複排除対象データを応答する。

このようにして、状態が「DeDup」のデータユニットに対するパトロールが行われるごとに、対応する重複排除対象データが格納された重複排除ユニットの保存期限満了時刻が更新される。パトロール間隔は、重複排除対象データの保存期間よりも短い。そのため、重複排除ユニットの重複排除オフセットが設定されたデータユニットが少なくとも 1 つ存在する間は、重複排除対象データは削除されない。

【 0 2 0 3 】

以上のようなパトロールを実現するためのパトロール部 1 5 0 の処理を、以下に具体的に説明する。なお、同一のスライスに属するユニットに割り当てられた単位記憶領域については、そのスライスの最終パトロール時刻から所定のパトロール周期に相当する時間経過後に実行される。スライス内のパトロールが終了すると、パトロール部 1 5 0 によってそのスライスの最終パトロール時刻が更新される。

【 0 2 0 4 】

図 2 3 は、パトロール処理の手順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。なお、以下の処理は、各スライスのデータユニット、およびデータ実体領域 1 1 0 b 内の未使用の単位記憶領域それぞれに対して、所定の時間間隔で実行される。

【 0 2 0 5 】

[ ステップ S 9 4 ] パトロール部 1 5 0 は、パトロール対象がスライス内の使用領域として定義されたデータユニットか、未使用の単位記憶領域かを判断する。スライス内の領域として定義されたデータユニットに対するパトロールであれば、処理がステップ S 9 5 に進められる。未使用の単位記憶領域に対するパトロールであれば、処理がステップ S 1 0 0 に進められる。

【 0 2 0 6 】

[ ステップ S 9 5 ] パトロール部 1 5 0 は、パトロール対象のユニットが属するスライスが割り当てられたボリュームのタイプ（論理ボリュームか重複排除ボリュームか）を判断する。論理ボリュームに割り当てられている場合、処理がステップ S 9 6 に進められる。重複排除ボリュームに割り当てられている場合、処理がステップ S 9 9 に進められる。

【 0 2 0 7 】

[ ステップ S 9 6 ] パトロール対象のデータユニットが属するスライスが、論理ボリュームに割り当てられていれば、パトロール部 1 5 0 はデータユニットの状態を判断する。データユニットの状態は、そのデータユニットに対応するデータユニット情報内に設定されている。データユニットの状態が「Blank」であれば、処理が終了する。データユニットの状態が「DeDup」であれば、処理がステップ S 9 7 に進められる。データユニットの状態が「Normal」であれば、処理がステップ S 9 8 に進められる。

【 0 2 0 8 】

[ ステップ S 9 7 ] データユニットの状態が「DeDup」の場合、パトロール部 1 5 0 は、重複排除対象データを管理しているディスクノードに対して、重複排除対象データの参照要求を送信する。具体的にはパトロール部 1 5 0 は、重複排除データアクセス要求部 1 4 2 に対して、パトロール対象のデータユニットに対応する重複排除対象データの取得を依頼する。すると、重複排除データアクセス要求部 1 4 2 により、データユニット情報に示される重複排除ボリューム ID と重複排除オフセットとに基づいて、重複排除対象データを管理するディスクノードから重複排除対象データが取得される。そして、取得した重複排除対象データが、重複排除データアクセス要求部 1 4 2 からパトロール部 1 5 0 に渡される。その後、処理が終了する。

【 0 2 0 9 】

[ ステップ S 9 8 ] データユニットの状態が「Normal」の場合、パトロール部 1 5 0 は

10

20

30

40

50

、Normalユニットパトロール処理を実行する。この処理の詳細は後述する。Normalユニットパトロール処理の後、パトロール処理が終了する。

【0210】

【ステップS99】パトロール対象のデータユニットが重複排除ボリュームに割り当てられたスライスのユニット（重複排除ユニット）である場合、重複排除ユニットパトロール処理が実行される。重複排除ユニットパトロール処理の後、パトロール処理が終了する。

【0211】

【ステップS100】パトロール対象が未使用の単位記憶領域である場合、パトロール部150は、その単位記憶領域に初期データ（例えば「0」の連続）を書き込む。その後、パトロール処理が終了する。

10

【0212】

図24は、状態が「Normal」のデータユニットのパトロール処理の手順を示すフローチャートである。以下、図24に示す処理をステップ番号に沿って説明する。

【ステップS101】パトロール部150は、データユニットに割り当てられた単位記憶領域からデータを読み出す。

【0213】

【ステップS102】パトロール部150は、データユニットの最終書き込み時刻から重複排除化猶予期間以上経過したか否かを判断する。重複排除化猶予期間以上経過している場合、処理がステップS104に進められる。重複排除化猶予期間を経過していなければ、処理がステップS103に進められる。

20

【0214】

【ステップS103】パトロール部150は、重複排除化猶予期間を経過していなければ、読み出したデータを、データユニットに割り当てられた単位記憶領域に書き込む。その後、処理が終了する。

【0215】

【ステップS104】パトロール部150は、重複排除化猶予期間を経過していれば、ステップS101で読み出したデータのハッシュ値を算出する。

【ステップS105】パトロール部150は、重複排除アドレス照会部133を用いて重複排除アドレスの照会を行う。これにより、パトロール部150は、重複排除アドレスを取得する。

30

【0216】

【ステップS106】パトロール部150は、ステップS105の重複排除アドレスの照会に回答して、インデックスサーバ500において重複排除ユニット割当処理が実行されたか否かを判断する。重複排除ユニット割当処理の実行の有無は、重複排除アドレスの照会に対するインデックスサーバ500からの回答で示される。重複排除ユニット割当処理が新たに実行された場合、処理がステップS107に進められる。重複排除ユニット割当処理が実行されていない場合、処理がステップS108に進められる。

【0217】

【ステップS107】パトロール部150は、重複排除データアクセス要求部142を用いて、ステップS101で読み出したデータを、ステップS105で取得した重複排除アドレスで示される重複排除ユニットに転送する。

40

【0218】

【ステップS108】パトロール部150は、データユニットの状態を「DeDup」に変更する。

【ステップS109】パトロール部150は、データユニットに割り当てられた単位記憶領域を解放する。

【0219】

【ステップS110】パトロール部150は、解放した単位記憶領域に初期データを書き込む。

50

このようにして、論理ボリューム内の状態が「Normal」のデータユニットのパトロールが行われる。

【0220】

図25は、重複排除ボリュームのパトロール処理の手順を示すフローチャートである。以下、図25に示す処理をステップ番号に沿って説明する。

【ステップS111】パトロール部150は、重複排除ユニットに割り当てられた単位記憶領域からデータを読み出す。

【0221】

【ステップS112】パトロール部150は、現在の時刻が重複排除ユニットの保存期限満了時刻を過ぎているか否かを判断する。保存期限満了時刻を過ぎている場合、処理がステップS114に進められる。保存期限満了時刻を過ぎていない場合、処理がステップS113に進められる。

10

【0222】

【ステップS113】パトロール部150は、保存期限満了時刻を過ぎていなければ、ステップS111で読み出したデータを、重複排除ユニットに割り当てられた単位記憶領域に書き込む。その後、重複排除ボリュームパトロール処理が終了する。

【0223】

【ステップS114】パトロール部150は、保存期限満了時刻を過ぎている場合、重複排除ユニットを未使用(NULL)にする。すなわち、パトロール部150は、保存期限満了時刻を過ぎている重複排除ユニットの重複排除ユニット情報に対し、「NULL」を設定する。

20

【0224】

【ステップS115】パトロール部150は、重複排除ユニットに割り当てられている単位記憶領域を解放する(割り当てを解除する)。

【ステップS116】パトロール部150は、重複排除ユニットに割り当てられていた単位記憶領域に、初期データを書き込む。その後、重複排除ボリュームパトロール処理が終了する。

【0225】

このようにして、パトロール部150によるパトロール処理が実行される。

次に、パトロール処理の過程で重複排除アドレス照会要求が出された場合の、インデックスサーバ500の処理を詳細に説明する。

30

【0226】

図26は、インデックスサーバにおける重複排除アドレス検索処理の手順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。

【ステップS121】重複排除アドレス検索部514は、ハッシュ値を含む重複排除アドレス照会要求を受信する。

【0227】

【ステップS122】重複排除アドレス検索部514は、重複排除ユニット情報記憶部512から、重複排除アドレス照会要求に示されるハッシュ値と一致するハッシュ値を含む重複排除ユニットを検索する。

40

【0228】

【ステップS123】重複排除アドレス検索部514は、該当する重複排除ユニットがあるか否かを判断する。該当する重複排除ユニットがあれば、処理がステップS127に進められる。該当する重複排除ユニットがなければ、処理がステップS124に進められる。

【0229】

【ステップS124】重複排除アドレス検索部514は、現在定義されている重複排除ボリューム内に未使用の重複排除ユニットがあるか否かを判断する。未使用の重複排除ユニットがあれば、処理がステップS126に進められる。未使用の重複排除ユニットがなければ、処理がステップS125に進められる。

50

## 【 0 2 3 0 】

[ステップS 1 2 5] 重複排除アドレス検索部 5 1 4 は、重複排除ボリュームを拡張する。例えば、重複排除アドレス検索部 5 1 4 は、制御ノード 6 0 0 に対して重複排除ボリュームへのセグメントの追加を依頼する。制御ノード 6 0 0 では、論理ボリューム管理部 6 3 0 が重複排除ボリュームにセグメントを追加し、そのセグメントにいずれかのディスクボリュームで管理されているスライスを割り当て、重複排除ボリューム用スライスマタデータを生成する。生成された重複排除ボリューム用スライスマタデータは、割り当てられたスライスを管理するディスクノードに送信されると共に、インデックスサーバ 5 0 0 に送信される。インデックスサーバ 5 0 0 に送信される重複排除ボリューム用スライスマタデータには、その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が付与されている。インデックスサーバ 5 0 0 は、取得した重複排除ボリューム用スライスマタデータを重複排除ユニット情報記憶部 5 1 2 に格納する。新たに作成された重複排除ボリューム用スライスマタデータに含まれる重複排除ユニット情報は、すべて未使用である。すなわち、未使用の重複排除ユニットが生成されたこととなる。

10

## 【 0 2 3 1 】

[ステップS 1 2 6] 重複排除アドレス検索部 5 1 4 は、未使用の重複排除ユニットを、新たな重複排除対象データの格納先として割り当てる。

[ステップS 1 2 7] 重複排除アドレス検索部 5 1 4 は、重複排除アドレスをディスクノードに応答する。ステップS 1 2 3 で該当する重複排除ユニットが見つかった場合、その重複排除ユニットの重複排除アドレスが送信される。ステップS 1 2 3 で該当する重複排除ユニットが見つからなかった場合、ステップS 1 2 6 で割り当てた重複排除ユニットの重複排除アドレスが送信される。

20

## 【 0 2 3 2 】

このようにして、重複排除アドレス照会要求に応じた重複排除アドレス検索処理が行われる。ここで、未使用の重複排除ユニットが存在しなければ、適宜、重複排除ボリュームの拡張が行われる。すなわち、重複排除ボリュームの記憶容量は、初期状態では必要最小限の容量にとどめておき、重複排除ユニットの不足が発生した場合に拡張される。その結果、システム全体における資源の有効利用が可能となる。

## 【 0 2 3 3 】

以上のようにしてパトロールが実行される。その結果、ディスクノード 1 0 0 において最後の書き込みから重複排除化猶予期間を経過した論理ボリューム内のデータについては、データ実体が重複排除ボリュームに移される。論理ボリュームには、移行先の重複排除ユニットを示す重複排除アドレスが残される。

30

## 【 0 2 3 4 】

またパトロールにより重複排除ボリュームから、保存期限満了時刻を過ぎた重複排除ユニットが検出されると、その重複排除ユニットのデータは削除され、データ実体領域内の単位記憶領域が解放される。

## 【 0 2 3 5 】

< 未使用ユニット情報反映処理 >

パトロールの実行により、各ディスクノード 1 0 0 , 2 0 0 , 3 0 0 , 4 0 0 で管理されている重複排除ユニットについては、データユニットからの関連付けが無くなった場合に未使用に変更される。重複排除ユニットが未使用に変更されたことは、インデックスサーバ 5 0 0 の重複排除ユニット情報記憶部 5 1 2 にも反映される。

40

## 【 0 2 3 6 】

図 2 7 は、未使用ユニット情報反映処理の手順を示すフローチャートである。以下、図 2 7 に示す処理をステップ番号に沿って説明する。なお、以下の処理は、重複排除ユニット情報記憶部 5 1 2 に格納されている重複排除ユニット情報ごとに行われる。

## 【 0 2 3 7 】

[ステップS 1 3 1] インデックスサーバ 5 0 0 の重複排除ユニット情報収集部 5 1 1 は、現在の時刻が保存期限満了時刻を過ぎているか否かを判断する。保存期限満了時刻を

50

経過した場合、処理がステップS 1 3 2に進められる。保存期限満了時刻を経過していなければ、ステップS 1 3 1の処理が繰り返される。

【0 2 3 8】

【ステップS 1 3 2】重複排除ユニット情報収集部5 1 1は、処理対象の重複排除ユニット情報の取得元であるディスクノードに対して、その重複排除ユニット情報で示される重複排除ユニットの現在の状況を問い合わせる。なお問い合わせ対象の重複排除ユニットは、重複排除ボリュームIDと、重複排除ボリューム内でのオフセット（重複排除オフセット）によって一意に指定される。この問い合わせに対して、ディスクノードの重複排除ユニット情報応答部が、現在の状況を応答する。応答には、現在の重複排除ユニット情報が含まれる。

10

【0 2 3 9】

【ステップS 1 3 3】重複排除ユニット情報収集部5 1 1は、ディスクノードから受け取った応答内容により、重複排除ユニットの状態を判断する。具体的には、重複排除ユニット情報収集部5 1 1は、取得した重複排除ユニット情報が「NULL」であれば、該当する重複排除ユニットは未使用であると判断する。重複排除ユニット情報収集部5 1 1は、取得した重複排除ユニット情報が「NULL」でなければ、該当する重複排除ユニットは使用中であると判断する。重複排除ユニットが未使用であれば、処理がステップS 1 3 4に進められる。重複排除ユニットが使用中であれば、処理がステップS 1 3 6に進められる。

【0 2 4 0】

20

【ステップS 1 3 4】重複排除ユニット情報収集部5 1 1は、ディスクノードにおける重複排除ユニットの状態が未使用であれば、重複排除ユニット情報記憶部5 1 2内の対応する重複排除ユニットの状態も未使用にする。具体的には、重複排除ユニット情報収集部5 1 1は、重複排除ユニット情報をNULLに変更する。

【0 2 4 1】

【ステップS 1 3 5】重複排除ユニット情報収集部5 1 1は、ステップS 1 3 4で未使用に変更した重複排除ユニット情報へのポインタを、未使用ユニットキューに追加する。その後、処理が終了する。

【0 2 4 2】

【ステップS 1 3 6】重複排除ユニット情報収集部5 1 1は、ディスクノードにおける重複排除ユニットの状態が使用中であれば、ディスクノードから取得した重複排除ユニット情報から保存期限満了時刻を取得する。

30

【0 2 4 3】

【ステップS 1 3 7】重複排除ユニット情報収集部5 1 1は、重複排除ユニット情報記憶部5 1 2内の対応する重複排除ユニットの保存期限満了時刻を更新する。その後、処理が終了する。

【0 2 4 4】

このようにして、ディスクノード1 0 0, 2 0 0, 3 0 0, 4 0 0における重複排除ユニット情報の更新内容が、インデックスサーバ5 0 0に反映される。すなわちインデックスサーバにおいて重複排除ユニット情報の保存期限満了時刻が過ぎた場合、重複排除ユニットを管理しているディスクノードから最新の重複排除ユニット情報が取得される。そして、インデックスサーバ5 0 0では、取得した重複排除ユニット情報の内容が、重複排除ユニット情報記憶部5 1 2に反映される。

40

【0 2 4 5】

次に、パトロールによるデータの保存先の変更例について説明する。

図2 8は、パトロールによるデータの保存先変更状況を示す図である。図2 8の例では、論理ボリューム2 0のセグメントID「2」のセグメントにディスクノード1 0 0が管理するスライス1 1 1が割り当てられている。また論理ボリューム3 0のセグメントID「3」のセグメントにディスクノード2 0 0が管理するスライス2 1 1が割り当てられている。さらに重複排除ボリューム4 0のセグメントID「1」のセグメントにディスクノ

50

ード 3 0 0 が管理するスライス 3 1 1 が割り当てられている。

【 0 2 4 6 】

スライス 1 1 1 の先頭から 3 つのデータユニットに割り当てられた単位記憶領域には、それぞれ「data[a]」、「data[b]」、「data[c]」が格納されている。スライス 2 1 1 の先頭から 3 つのデータユニットに割り当てられた単位記憶領域には、それぞれ「data[u]」、「data[a]」、「data[v]」が格納されている。すなわち、スライス 1 1 1 の先頭のデータユニットと、スライス 2 1 1 の 2 番目のデータユニットとは、格納されているデータの内容が同じである。

【 0 2 4 7 】

また、スライス 1 1 1 の 2 番目のデータユニットとスライス 2 1 1 の 1 番目のデータユニットとは、頻繁に書き込みが行われ、最後の書き込みから重複排除化猶予期間が経過していないものとする。その他のデータユニットは、重複排除化猶予期間以上書き込みが行われていない。

10

【 0 2 4 8 】

このような状況下でパトロールが実行されると、重複排除化猶予期間以上書き込みが行われていないデータユニット内のデータが重複排除対象データとして選択され、重複排除ボリューム 4 0 に移動される。図 2 8 の例では、「data[a]」が重複排除ボリューム 4 0 の先頭の重複排除ユニットに移動され、「data[c]」が重複排除ボリューム 4 0 の 2 番目の重複排除ユニットに移動され、「data[v]」が重複排除ボリューム 4 0 の 3 番目の重複排除ユニットに移動されている。

20

【 0 2 4 9 】

データの移動が行われた移動元のデータユニットには、移動先の重複排除ユニットを一意に示す重複排除アドレスが設定される。図 2 8 では、重複排除ボリューム ID の後に、重複排除ユニットのオフセットを示している。図 2 8 の例では、重複排除ボリュームの最初のセグメント（セグメント ID 「 1 」）の重複排除ユニットが使用されているため、そのセグメント内での各重複排除ユニットの順番がオフセット値となる。

【 0 2 5 0 】

ハッシュ値は、データの内容が同一であれば同じ値となる。図 2 8 の例では、「data[a]」のハッシュ値は「hash[A]」、「data[c]」のハッシュ値は「hash[C]」、「data[v]」のハッシュ値は「hash[V]」である。

30

【 0 2 5 1 】

データを移動する際に算出したハッシュ値が既に重複排除ユニットに設定されていれば、その重複排除ユニットの重複排除アドレスがデータ移動元のデータユニットに設定される。図 2 8 の例では「data[a]」については、パトロール前は 2 つの論理ボリューム 2 0 , 3 0 の両方に重複して格納されていたが、パトロール後は、「data[a]」は重複排除ボリューム 4 0 に 1 つだけ格納されている。「data[a]」が格納されていたデータユニットには、同じ重複排除ユニットを示す重複排除アドレスが設定される。これにより、システム全体で保持すべきデータ量を削減し、資源の効率的利用が可能となる。

【 0 2 5 2 】

論理ボリュームにおいて重複排除対象データと判断され重複排除ボリュームに移行されたデータであっても、アクセスノードからのアクセスによりデータが更新されることがある。その場合には、更新後のデータがデータユニットに割り当てられた単位記憶領域に格納され、重複排除ユニットとの関連付けは解消される。

40

【 0 2 5 3 】

図 2 9 は、重複排除ユニットへの関連付けの解消状況を示す図である。図 2 9 の例では、論理ボリューム 3 0 のセグメント ID 「 3 」のセグメントに対する書き込みが行われた場合を想定している。具体的には、そのセグメントに割り当てられたスライス 2 1 1 の 2 番目と 3 番目のデータユニットのデータが、それぞれ「data[x]」、「data[y]」に更新されている。

【 0 2 5 4 】

50

この場合、更新されたデータユニットに割り当てられた単位記憶領域に「data[x]」、「data[y]」が格納され、各データユニットの重複排除アドレスは削除される。重複排除アドレスが削除されたことで、重複排除ユニットとの関連付けは解消される。

【0255】

さらに、すべてのデータユニットからの関連付けが解消された重複排除ユニットについては、保持していたデータが削除され未使用の状態に変更される。

図30は、重複排除ユニットの未使用状態への変更状況を示す図である。図30の例では、ディスクノード100でパトロールが行われた後に、ディスクノード300でパトロールが行われた状況を示している。ディスクノード100でパトロールが行われると、ディスクノード100が、データユニットとの関連付けを有する重複排除ユニットの参照（データ取得）を行う。すると参照された重複排除ユニットの保存期限満了時刻が更新される。

10

【0256】

その後、ディスクノード300でパトロールが行われると、データユニットとの関連付けを有していない重複排除ユニットについて、保存期限満了時刻を過ぎていることが検出される。該当する重複排除ユニットは、未使用に変更される。図30の例では、重複排除ボリューム40のセグメントID「1」のセグメントに割り当てられたスライス311の3番目の重複排除ユニットが未使用に変更されている。

【0257】

このようにして、重複排除ボリューム40から不要なデータを削除し、資源の効率的な利用が促進される。

20

〔その他の応用例〕

第2の実施の形態では、インデックスサーバ500と制御ノード600とが個別に設けられているが、1つの装置にインデックスサーバ500と制御ノード600との機能を搭載することもできる。またインデックスサーバ500の機能を、いずれかのディスクノード100, 200, 300, 400に搭載することもできる。

【0258】

第2の実施の形態では、重複排除ユニットに参照があれば、保存期限満了時刻が更新される。そこで、状態が「DeDup」のデータユニットのパトロールでは、関連付けられた重複排除ユニットの参照を実行することで、重複排除ユニットの保存期限満了時刻を更新させている。この処理では、重複排除ユニットの保存期限満了時刻を更新させることができれば、参照以外の機能を利用してよい。例えば、データユニットを管理するディスクノードから重複排除ユニットを管理するディスクノードへ、パトロールごとに保存期限満了時刻更新要求を送信するようにしてもよい。重複排除ユニットを管理するディスクノードでは、保存期限満了時刻更新要求に応答して、重複排除ユニットの保存期限満了時刻を更新する。

30

【0259】

また、ハッシュ関数によるハッシュ値の算出に代えて、例えば暗号化処理による暗号データを用いることもできる。

なお、上記の処理機能は、コンピュータによって実現することができる。その場合、ディスクノード100, 200, 300, 400、インデックスサーバ500、制御ノード600、およびアクセスノード700, 800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。

40

【0260】

50



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。

【0261】

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。

10

【0262】

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。

【0263】

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。

20

【0264】

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) 複数のディスクノードそれぞれで管理されるストレージ装置にデータを分散格納するマルチノードストレージシステムを構成する1つの前記ディスクノードで実行すべきデータ管理処理をコンピュータに実行させるデータ管理プログラムにおいて、

前記コンピュータに、

仮想的な記憶領域を定義した論理ボリュームを複数に分割して得られるデータユニットを指定した書き込み要求があると、前記コンピュータに接続されたストレージ装置内の単位記憶領域の1つを指定された前記データユニットに割り当て、前記データユニットに割り当てられた前記単位記憶領域に対してデータの書き込みを行い、

30

書き込み対象となった前記データユニットに対応付けて、現在の時刻を最終書き込み時刻としてデータユニット情報記憶手段に格納し、

前記データユニット情報記憶手段を参照し、前記最終書き込み時刻から重複排除化猶予期間以上経過している前記データユニットを検出し、

仮想的な記憶領域を定義した重複排除ボリュームを複数に分割して得られる重複排除ユニットのうち重複排除化の対象とされた重複排除対象データの格納に使用されている重複排除ユニットの識別情報および前記重複排除ユニットを管理している前記ディスクノードの識別子を含む重複排除アドレスと、前記重複排除対象データに所定の演算を行うことで得られる前記重複排除対象データの固有値とを対応付けて管理するインデックスサーバから、検出された前記データユニットに割り当てられた前記単位記憶領域内のデータの前記固有値に対応付けられた前記重複排除アドレスを取得し、

40

検出された前記データユニットに対応付けて、取得した前記重複排除アドレスを前記データユニット情報記憶手段に格納すると共に、検出された前記データユニットへの前記単位記憶領域の割り当てを解除する、

処理を実行させることを特徴とするデータ管理プログラム。

【0265】

(付記2) 前記コンピュータに、さらに、

前記インデックスサーバから、前記重複排除ユニットの識別情報を指定した重複排除ユニット割当要求を受信すると、指定された前記重複排除ユニットに前記単位記憶領域の1

50

つを割り当て、指定された前記重複排除ユニットをデータの読み出し先とする前記ディスクノードから前記重複排除対象データを取得し、指定された前記重複排除ユニットに割り当てた前記単位記憶領域に格納し、

前記重複排除ユニットの識別情報を指定したデータの読み出し要求を受信すると、指定された前記重複排除ユニットに割り当てられた前記単位記憶領域内の前記重複排除対象データを応答する、

処理を実行させることを特徴とする付記 1 記載のデータ管理プログラム。

【0266】

(付記 3) 前記コンピュータに、さらに、

前記重複排除ユニットを前記重複排除対象データの読み出し先とする前記ディスクノードが存在しなくなった場合、前記重複排除ユニットへの前記単位記憶領域の割り当てを解除する、

処理を実行させることを特徴とする付記 2 記載のデータ管理プログラム。

【0267】

(付記 4) 前記コンピュータにさらに、

前記重複排除ユニットに対応付けて、前記重複排除ユニットに割り当てられた前記単位記憶領域から最後にデータの読み出しが行われた時刻から保存期間経過後の保存期限満了時刻を重複排除ユニット情報記憶手段に格納し、

前記重複排除ユニット情報記憶手段を参照し、現在の時刻が前記保存期限満了時刻を過ぎている前記重複排除ユニットを検出し、検出した前記重複排除ユニットへの前記単位記憶領域の割り当てを解除する、

処理を実行させることを特徴とする付記 3 記載のデータ管理プログラム。

【0268】

(付記 5) 前記コンピュータに、さらに、

前記インデックスサーバから取得した前記重複排除アドレスが、新たに割り当てられた前記重複排除ユニットの前記重複排除アドレスの場合、前記データユニットに割り当てられた前記単位記憶領域内のデータを前記重複排除アドレスに示される前記ディスクノードへ送信後、前記データユニットへの前記単位記憶領域の割り当てを解除する、

処理を実行させることを特徴とする付記 2 記載のデータ管理プログラム。

【0269】

(付記 6) 前記コンピュータに、さらに、

定期的に前記データユニット情報記憶手段を参照し、前記重複排除アドレスが対応付けられた前記データユニットを検出し、前記重複排除アドレスに示される前記ディスクノードから、前記重複排除アドレスで示される前記重複排除ユニットの識別情報を指定してデータを取得する、

処理を実行させることを特徴とする付記 1 記載のデータ管理プログラム。

【0270】

(付記 7) 前記固有値は、前記重複排除対象データに対して所定のハッシュ関数による演算を行うことで得られるハッシュ値であることを特徴とする付記 1 記載のデータ管理プログラム。

【0271】

(付記 8) 前記重複排除ユニットの識別情報は、前記重複排除ボリューム内での前記重複排除ユニットの位置を示すオフセットであることを特徴とする付記 1 記載のデータ管理プログラム。

【0272】

(付記 9) 前記コンピュータに、さらに、

前記データユニットを指定した読み出し要求があると、読み出し対象の前記データユニットに割り当てられた前記単位記憶領域がない場合、前記データユニット情報記憶手段を参照し、読み出し対象の前記データユニットに対応付けられた前記重複排除アドレスに示される前記ディスクノードから、前記重複排除アドレスに示される前記重複排除ユニット

10

20

30

40

50

の識別情報を指定して前記重複排除対象データを取得する、  
処理を実行させることを特徴とする付記 1 記載のデータ管理プログラム。

【 0 2 7 3 】

(付記 1 0) 前記コンピュータに、さらに、  
前記データユニットを指定した読み出し要求があると、読み出し対象の前記データユニットに割り当てられた前記単位記憶領域がある場合、割り当てられた前記単位記憶領域からデータを読み出す、

処理を実行させることを特徴とする付記 1 記載のデータ管理プログラム。

【 0 2 7 4 】

(付記 1 1) 複数のディスクノードそれぞれで管理されるストレージ装置にデータを分散格納するマルチノードストレージシステムの記憶領域管理処理をコンピュータに実行させるデータ管理プログラムにおいて、

前記コンピュータに、

仮想的な記憶領域を定義した重複排除ボリュームを複数に分割して得られる重複排除ユニットのうち前記ディスクノードで管理している前記重複排除ユニットの使用の有無、使用されている前記重複排除ユニットに割り当てられている単位記憶領域に格納された重複排除対象データに所定の演算を行うことで得られる固有値、および前記重複排除ユニットの識別情報を含む重複排除ユニット情報が、前記重複排除ユニットを管理している前記ディスクノードの識別子に対応付けて重複排除ユニット情報記憶手段に格納されており

前記ディスクノードから、重複排除化データユニットに割り当てられた単位記憶領域内のデータに前記所定の演算を実行して得られた固有値を含む重複排除アドレス照会要求を受け取ると、前記重複排除ユニット情報記憶手段から、前記重複排除アドレス照会要求に示される固有値を含む前記重複排除ユニット情報を検索し、

検索により該当する前記重複排除ユニット情報が検出された場合、検出された前記重複排除ユニット情報に含まれる前記重複排除ユニットの識別情報と、検出された前記重複排除ユニット情報に対応する前記重複排除ユニットを管理している前記ディスクノードの識別子とを含む重複排除アドレスを、前記重複排除アドレス照会要求の送信元である前記ディスクノードに応答し、

検索により該当する前記重複排除ユニット情報が検出されなかった場合、前記重複排除ユニット情報記憶手段を参照し、未使用の前記重複排除ユニットを選択し、選択した前記重複排除ユニットを管理する前記ディスクノードに対して選択した重複排除ユニットへの前記単位記憶領域の割り当て要求を送信し、選択した前記重複排除ユニットの前記単位記憶領域割り当て後の前記重複排除ユニット情報を前記重複排除ユニット情報記憶手段に格納し、選択した前記重複排除ユニットの識別情報と、選択した前記重複排除ユニットを管理する前記ディスクノードの識別子とを含む重複排除アドレスを、前記重複排除アドレス照会要求の送信元である前記ディスクノードに応答する、

処理を実行させることを特徴とするデータ管理プログラム。

【 0 2 7 5 】

(付記 1 2) 前記コンピュータに、

前記重複排除ユニット情報を複数の前記ディスクノードそれぞれから取得し、取得元の前記ディスクノードの識別子に対応付けて重複排除ユニット情報記憶手段に格納する、

処理を実行させることを特徴とする付記 1 1 記載のデータ管理プログラム。

【 0 2 7 6 】

(付記 1 3) 前記コンピュータに、さらに、

前記重複排除ユニット情報には、前記重複排除ユニットに割り当てられた前記単位記憶領域から最後にデータの読み出しが行われた時刻から保存期間経過後の保存期限満了時刻が含まれており、前記重複排除ユニット情報記憶手段を参照し、現在の時刻が前記保存期限満了時刻を過ぎている前記重複排除ユニットを検出し、

検出された前記重複排除ユニット情報の取得元の前記ディスクノードから最新の前記重複排除ユニット情報を取得し、前記重複排除ユニット情報記憶手段に反映する、

10

20

30

40

50

処理を実行させることを特徴とする付記 1 1 記載のデータ管理プログラム。

【 0 2 7 7 】

(付記 1 4) 複数のディスクノードそれぞれで管理されるストレージ装置にデータを分散格納するマルチノードストレージシステムを構成する 1 つの前記ディスクノードで実行すべきデータ管理処理をコンピュータによって実現するデータ管理装置において、

仮想的な記憶領域を定義した論理ボリュームを複数に分割して得られるデータユニットを指定した書き込み要求があると、前記コンピュータに接続されたストレージ装置内の単位記憶領域の 1 つを指定された前記データユニットに割り当て、前記データユニットに割り当てられた前記単位記憶領域に対してデータの書き込みを行う書き込みアクセス手段と

10

書き込み対象となった前記データユニットに対応付けて、現在の時刻を最終書き込み時刻としてデータユニット情報記憶手段に格納する最終書き込み時刻更新手段と、

前記データユニット情報記憶手段を参照し、最終書き込み時刻から重複排除化猶予期間以上経過している前記データユニットを検出する重複排除化データユニット検出手段と、

仮想的な記憶領域を定義した重複排除ボリュームを複数に分割して得られる重複排除ユニットのうち重複排除化の対象とされた重複排除対象データの格納に使用されている重複排除ユニットの識別情報および前記重複排除ユニットを管理している前記ディスクノードの識別子を含む重複排除アドレスと、前記重複排除対象データに所定の演算を行うことで得られる前記重複排除対象データの固有値とを対応付けて管理するインデックスサーバから、検出された前記データユニットに割り当てられた前記単位記憶領域内のデータの前記固有値に対応付けられた前記重複排除アドレスを取得する重複排除アドレス取得手段と、

20

検出された前記データユニットに対応付けて、取得した前記重複排除アドレスを前記データユニット情報記憶手段に格納すると共に、検出された前記データユニットへの前記単位記憶領域の割り当てを解除する単位記憶領域割当解除手段と、

前記データユニットを指定した読み出し要求があると、読み出し対象の前記データユニットに割り当てられた前記単位記憶領域がある場合、割り当てられた前記単位記憶領域からデータを読み出し、読み出し対象の前記データユニットに割り当てられた前記単位記憶領域がない場合、前記データユニット情報記憶手段を参照し、読み出し対象の前記データユニットに対応付けられた前記重複排除アドレスに示される前記ディスクノードから、前記重複排除アドレスに示される前記重複排除ユニットの識別情報を指定して前記重複排除対象データを取得する読み出しアクセス手段と、

30

を有することを特徴とするデータ管理装置。

【 0 2 7 8 】

(付記 1 5) 複数のディスクノードそれぞれで管理されるストレージ装置にデータを分散格納するマルチノードストレージシステムを構成する 1 つの前記ディスクノードで実行すべきデータ管理処理をコンピュータで実行するデータ管理方法において、

前記コンピュータが、

仮想的な記憶領域を定義した論理ボリュームを複数に分割して得られるデータユニットを指定した書き込み要求があると、前記コンピュータに接続されたストレージ装置内の単位記憶領域の 1 つを指定された前記データユニットに割り当て、前記データユニットに割り当てられた前記単位記憶領域に対してデータの書き込みを行い、

40

書き込み対象となった前記データユニットに対応付けて、現在の時刻を最終書き込み時刻としてデータユニット情報記憶手段に格納し、

前記データユニット情報記憶手段を参照し、最終書き込み時刻から重複排除化猶予期間以上経過している前記データユニットを検出し、

仮想的な記憶領域を定義した重複排除ボリュームを複数に分割して得られる重複排除ユニットのうち重複排除化の対象とされた重複排除対象データの格納に使用されている重複排除ユニットの識別情報および前記重複排除ユニットを管理している前記ディスクノードの識別子を含む重複排除アドレスと、前記重複排除対象データに所定の演算を行うことで得られる前記重複排除対象データの固有値とを対応付けて管理するインデックスサーバか

50

ら、検出された前記データユニットに割り当てられた前記単位記憶領域内のデータの前記固有値に対応付けられた前記重複排除アドレスを取得し、

検出された前記データユニットに対応付けて、取得した前記重複排除アドレスを前記データユニット情報記憶手段に格納すると共に、検出された前記データユニットへの前記単位記憶領域の割り当てを解除し、

前記データユニットを指定した読み出し要求があると、読み出し対象の前記データユニットに割り当てられた前記単位記憶領域がある場合、割り当てられた前記単位記憶領域からデータを読み出し、読み出し対象の前記データユニットに割り当てられた前記単位記憶領域がない場合、前記データユニット情報記憶手段を参照し、読み出し対象の前記データユニットに対応付けられた前記重複排除アドレスに示される前記ディスクノードから、前記重複排除アドレスに示される前記重複排除ユニットの識別情報を指定して前記重複排除対象データを取得する、

10

ことを特徴とするデータ管理方法。

【符号の説明】

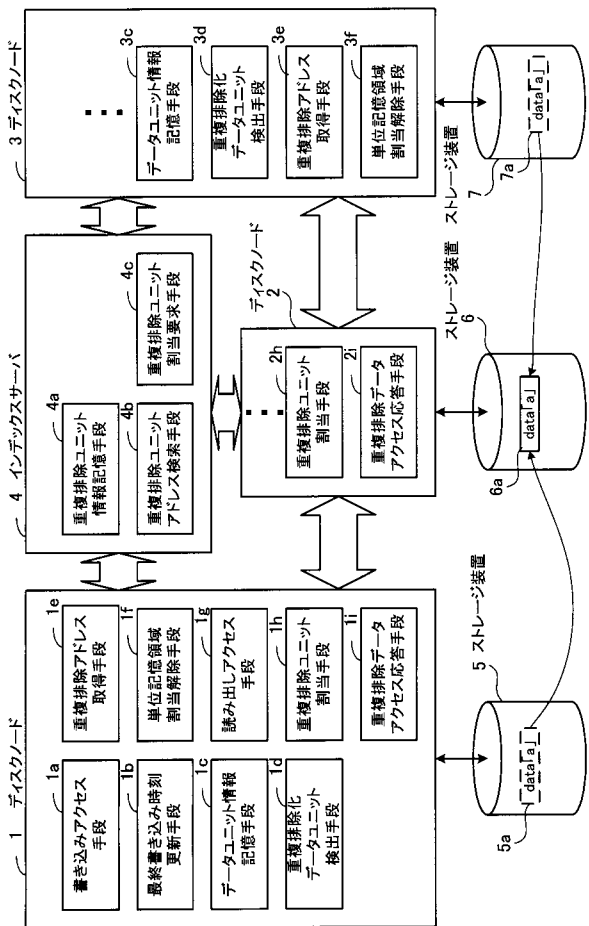
【0279】

- 1～3 ディスクノード
- 1 a 書き込みアクセス手段
- 1 b 最終書き込み時刻更新手段
- 1 c , 3 c データユニット情報記憶手段
- 1 d , 3 d 重複排除化データユニット検出手段
- 1 e , 3 e 重複排除アドレス取得手段
- 1 f , 3 f 単位記憶領域割当解除手段
- 1 g 読み出しアクセス手段
- 1 h , 2 h 重複排除ユニット割当手段
- 1 i , 2 i 重複排除データアクセス応答手段
- 4 インデックスサーバ
- 4 a 重複排除ユニット情報記憶手段
- 4 b 重複排除ユニットアドレス検索手段
- 4 c 重複排除ユニット割当要求手段
- 5～7 ストレージ装置
- 5 a , 6 a , 7 a 単位記憶領域

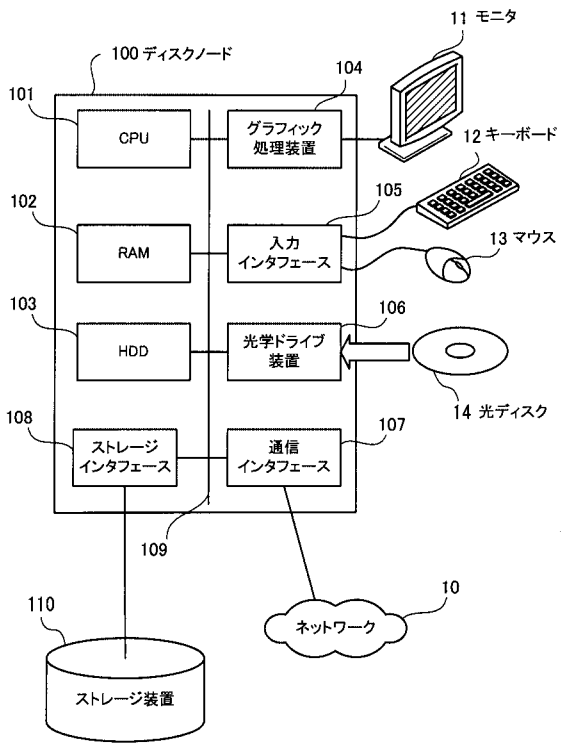
20

30

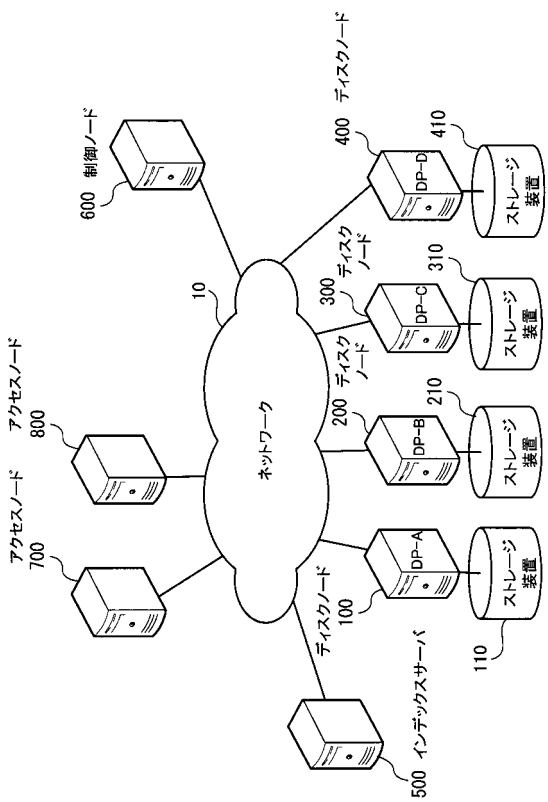
【図1】



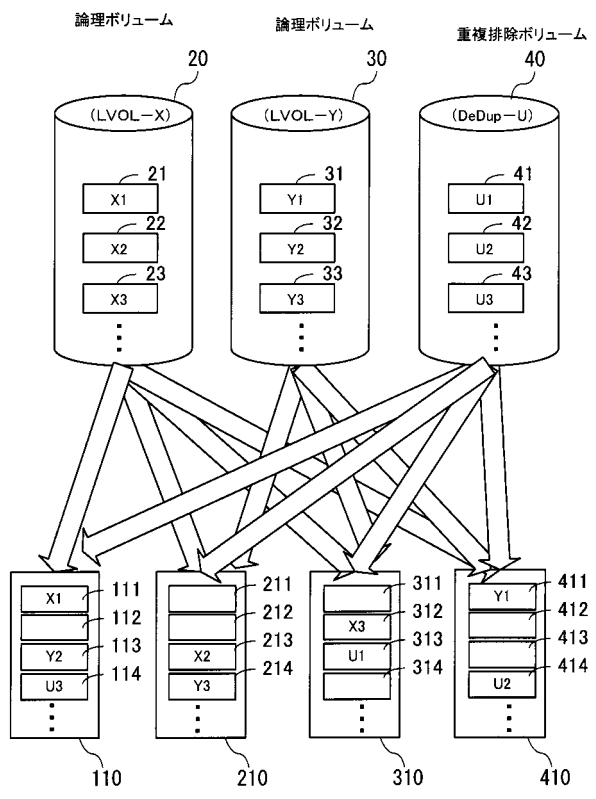
【図3】



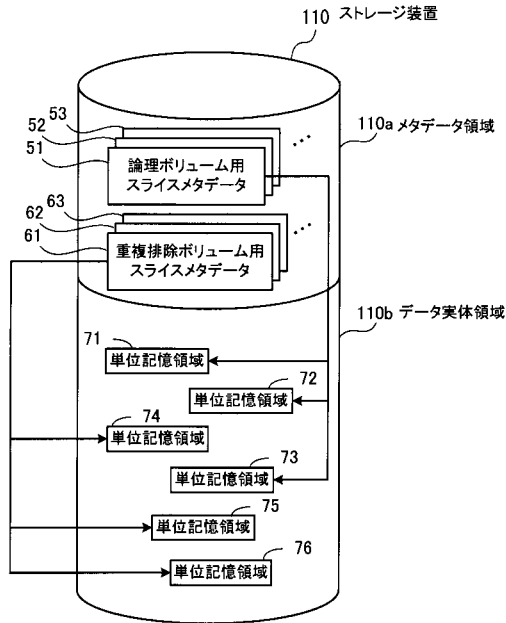
【図2】



【図4】



【図5】



【図6】

50 論理ボリューム用スライスメタデータ

論理ボリュームID		
セグメントID		
最終パトロール時刻		
Blank	-	
Normal	最終書き込み時刻	データ実体オフセット
Blank	-	
DeDup	重複排除ボリュームID	重複排除オフセット
Normal	最終書き込み時刻	データ実体オフセット
Blank	-	
⋮	⋮	⋮

データユニット情報 50a

スライスサイズ [個]  
ユニットサイズ

【図7】

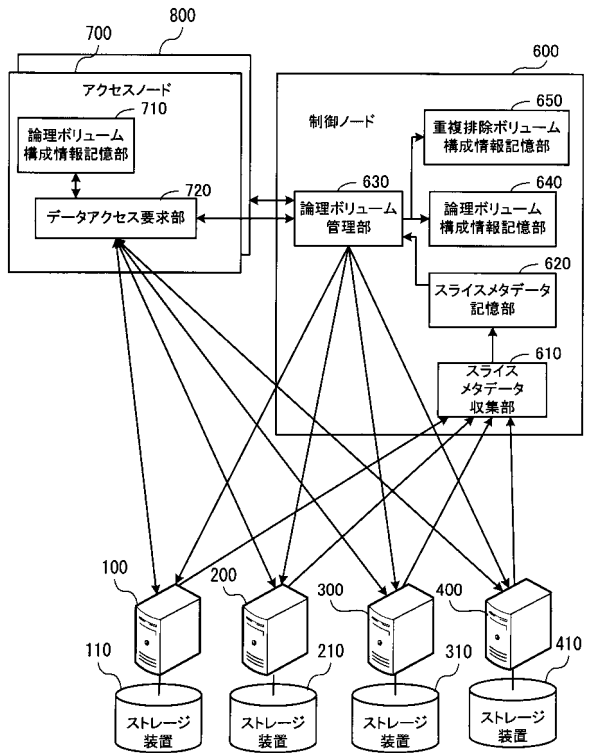
60 重複排除ボリューム用スライスメタデータ

重複排除ボリュームID		
セグメントID		
最終パトロール時刻		
ハッシュ値A	保存期限満了時刻A	データ実体オフセットA
ハッシュ値B	保存期限満了時刻B	データ実体オフセットB
NULL		
ハッシュ値C	保存期限満了時刻C	データ実体オフセットC
NULL		
ハッシュ値D	保存期限満了時刻D	データ実体オフセットD
⋮	⋮	⋮

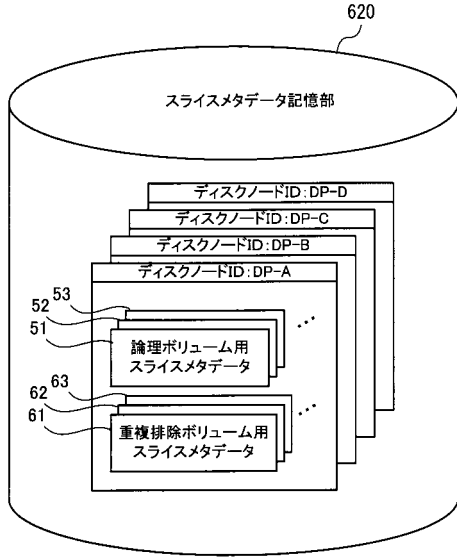
重複排除ユニット情報 60a

スライスサイズ [個]  
ユニットサイズ

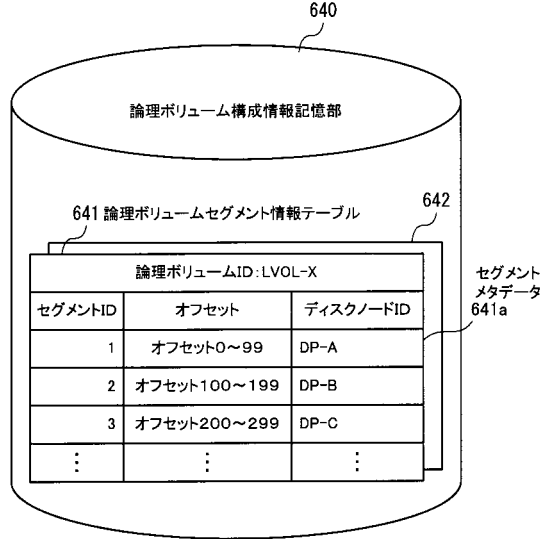
【図8】



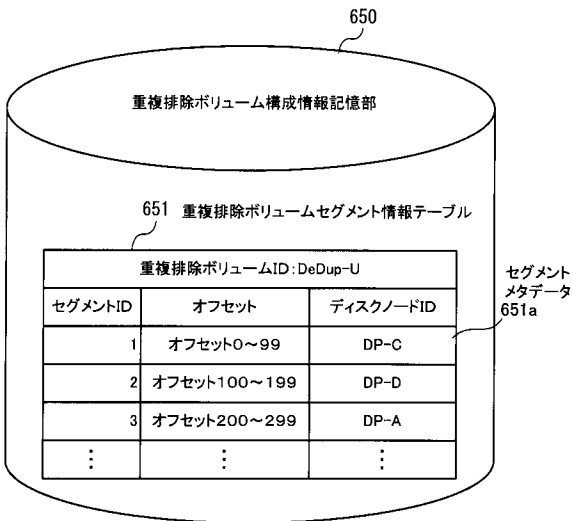
【図9】



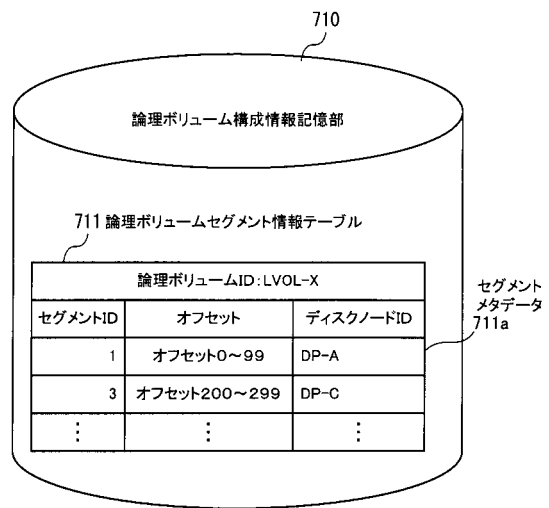
【図10】



【図11】



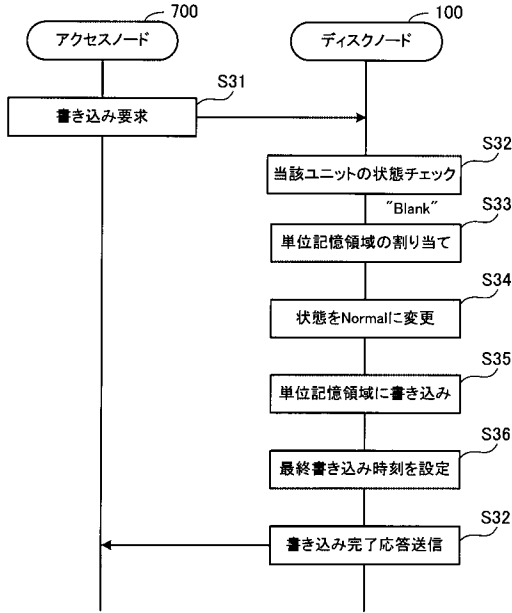
【図12】



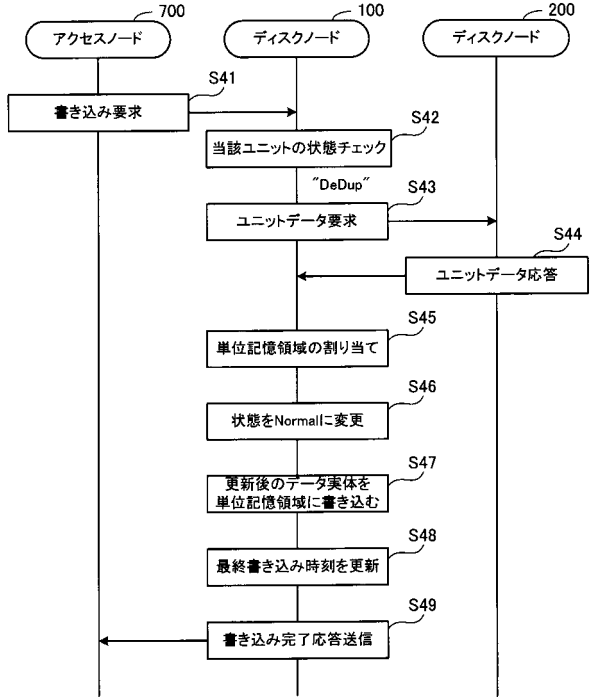




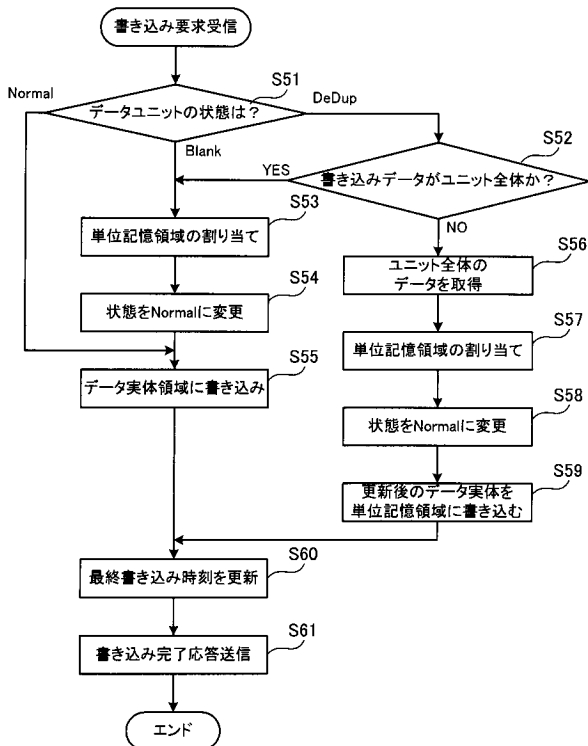
【図18】



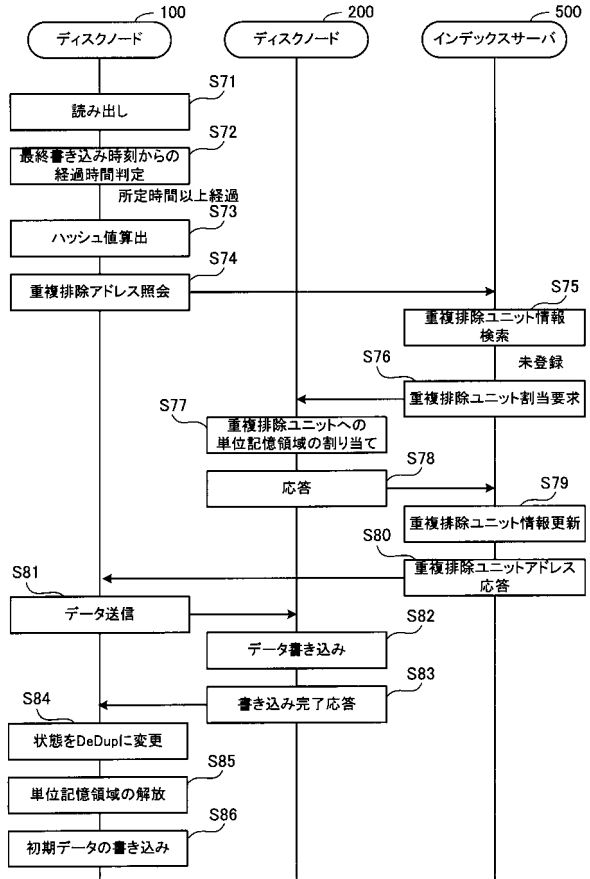
【図19】



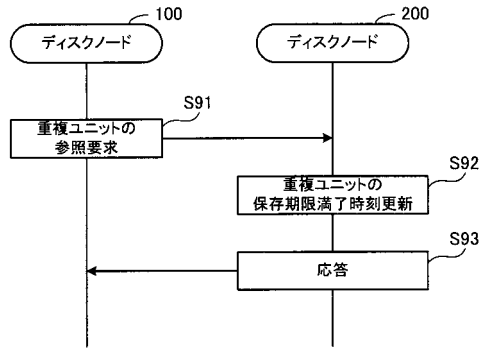
【図20】



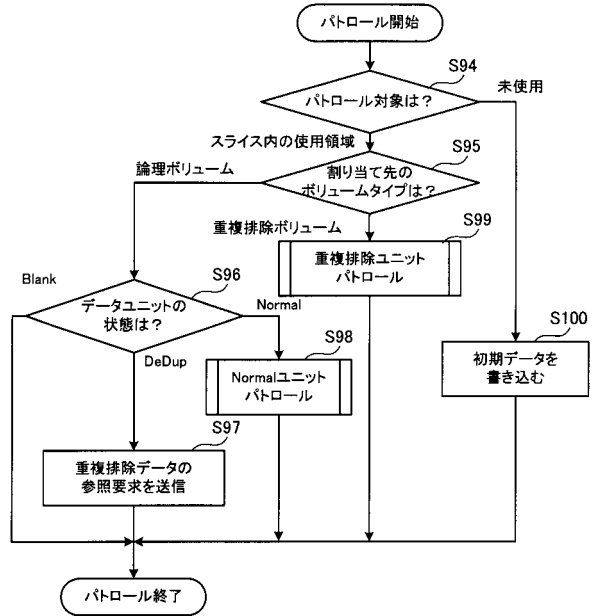
【図21】



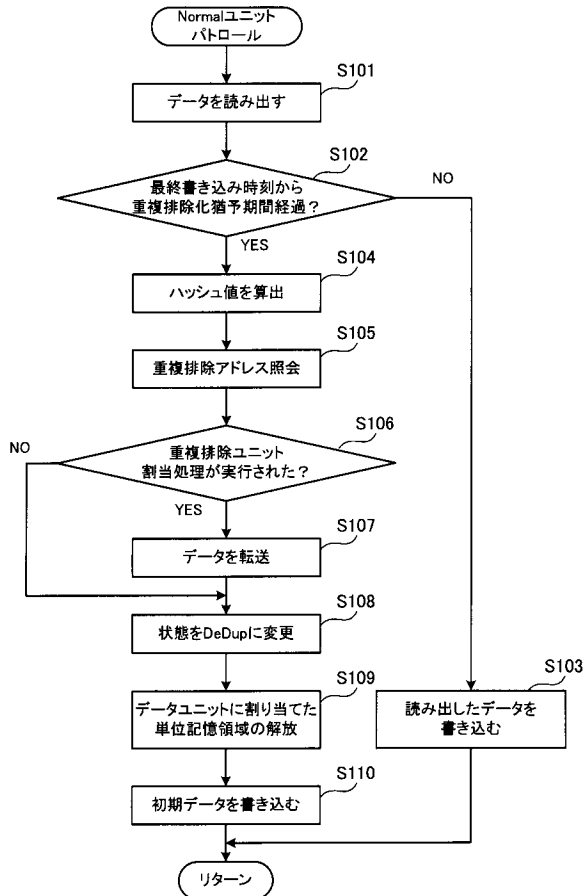
【図 2 2】



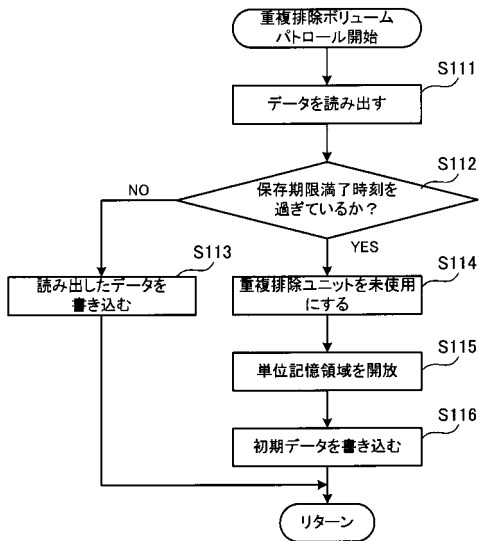
【図 2 3】



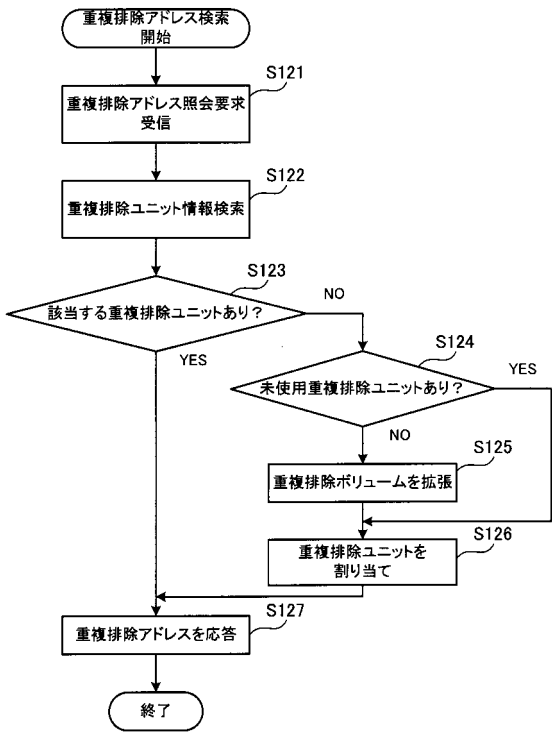
【図 2 4】



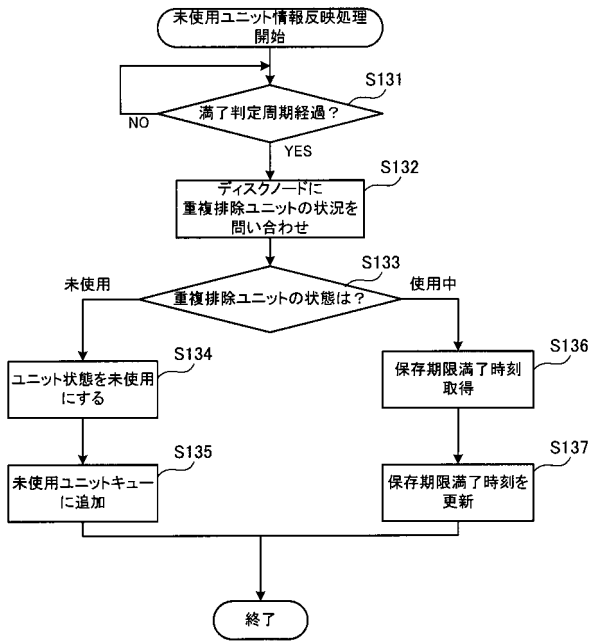
【図 2 5】



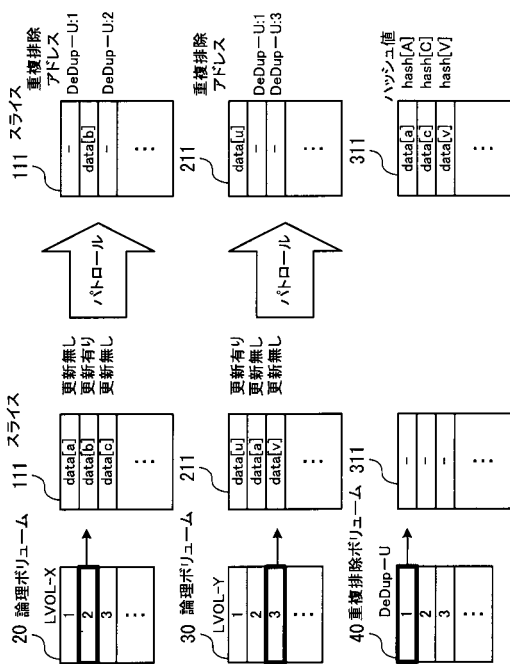
【図 26】



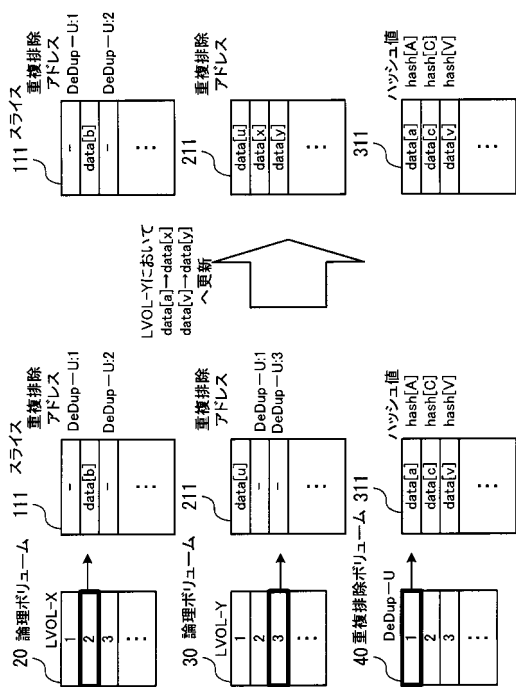
【図 27】



【図 28】



【図 29】





## フロントページの続き

- (72)発明者 丸山 哲太郎  
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 土屋 芳浩  
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 渡辺 高志  
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 熊野 達夫  
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
- (72)発明者 大江 和一  
神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

審査官 辻本 泰隆

- (56)参考文献 特開2009-237979(JP,A)  
特開2009-87021(JP,A)  
特開2009-146381(JP,A)  
特開2009-205201(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 12/00

JSTPlus/JMEDPlus/JST7580(JDreamII)

Cinii