

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5309043号
(P5309043)

(45) 発行日 平成25年10月9日(2013.10.9)

(24) 登録日 平成25年7月5日(2013.7.5)

(51) Int.Cl. F I
G 0 6 F 3 / 0 6 (2 0 0 6 . 0 1) G 0 6 F 3 / 0 6 3 0 1 A

請求項の数 14 外国語出願 (全 30 頁)

<p>(21) 出願番号 特願2010-16440 (P2010-16440) (22) 出願日 平成22年1月28日 (2010.1.28) (65) 公開番号 特開2010-182302 (P2010-182302A) (43) 公開日 平成22年8月19日 (2010.8.19) 審査請求日 平成24年7月24日 (2012.7.24) (31) 優先権主張番号 12/365,566 (32) 優先日 平成21年2月4日 (2009.2.4) (33) 優先権主張国 米国 (US)</p>	<p>(73) 特許権者 000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号 (74) 代理人 100093861 弁理士 大賀 真司 (74) 代理人 100129218 弁理士 百本 宏之 (72) 発明者 兼田 泰典 アメリカ合衆国 カリフォルニア州 95 129 サンノゼ コーデリア・アベニュー 1208 審査官 稲葉 崇</p>
---	---

最終頁に続く

(54) 【発明の名称】 ストレージシステム及びストレージシステムでの重複データ削除のための方法

(57) 【特許請求の範囲】

【請求項1】

ストレージ制御装置と、
チャンクプールを形成する複数のチャンクに分割されている複数のストレージ装置と、
前記ストレージ制御装置をホストコンピュータに接続するように構成したネットワーク
インターフェースと、
を備え、
前記ストレージ制御装置はメモリを備え、当該メモリのプログラムを実行してデータス
トレージボリュームを生成し、前記ネットワークインターフェースを経由して前記データ
ストレージボリュームを前記ホストコンピュータに利用可能にするストレージシステムで
あって、
前記メモリは、
前記ホストコンピュータから前記データストレージボリュームに向けられた書込コマン
ドに関連したデータに対応する識別子の複数を夫々有する複数の識別子テーブルと、
前記ホストコンピュータから前記ストレージ制御装置への書込みコマンドに伴うメタデ
ータの複数和、前記複数の識別子テーブルの夫々の情報とを有し、前記一つ又は複数のメ
タデータを前記複数の識別子テーブルの何れかに割り付けるメタデータテーブルと、
を有し、
前記複数の識別子テーブルの夫々は、前記識別子に対応するデータが格納された前記チ
ャンクの情報が当該識別子に割り付けられて構成され、

10

20

前記ストレージ制御装置は、前記メモリのプログラムを実行して、
 前記ホストコンピュータからの前記データストレージボリュームに向けられた書込コマンドを受信すると、前記識別子を計算し、
 前記メタデータテーブルを参照して、前記複数の識別子テーブルのうち、前記書込みコマンドに伴う前記メタデータに割り付けられた識別子テーブルを選択し、
 前記選択された識別子テーブルを参照して、前記計算された識別子の固有性を前記チャンク毎に確認し、
 前記計算された識別子が前記識別子テーブルにないものであれば、
 前記書込コマンドに関連した前記データに対応する識別子を生成し、
 前記チャンクプールから前記データストレージボリュームの前記書込みコマンドで指定されたデータ位置へ前記チャンクを割り付け、前記書込コマンドに関連した前記データを、前記割り付けられたチャンクに格納し、
 前記生成された識別子と前記割り付けられたチャンクの情報とに基づいて前記識別子テーブルを更新し、
 一方、前記計算された識別子が前記識別子テーブルに存在するものであれば、前記書込みコマンドに対応するデータを格納しない、
 ことを特徴とするストレージシステム。

10

【請求項 2】

前記識別子は、前記ホストコンピュータから受け取った前記データのSecure Hash Algorithm (S H A) 値を含む、
 ことを特徴とする請求項 1 によるストレージシステム。

20

【請求項 3】

前記識別子は、前記データのSecure Hash Algorithm (S H A) 値とハッシュ競合の数との組み合わせを含む、
 ことを特徴とする請求項 1 によるストレージシステム。

【請求項 4】

前記ストレージシステムは、前記ホストコンピュータから受け取った前記データが格納されているか否かを決定する前に、前記データを格納する、
 ことを特徴とする請求項 1 によるストレージシステム。

【請求項 5】

前記ストレージシステムは、前記データストレージボリュームに前記データを格納した後、前記識別子の計算および前記決定を、非同期に実行する、
 ことを特徴とする請求項 4 によるストレージシステム。

30

【請求項 6】

前記ストレージシステムを管理コンピュータに接続するように構成したマネージメントネットワークインターフェースをさらに備え、
 前記メタデータが前記管理コンピュータから登録される、
 ことを特徴とする請求項 1 によるストレージシステム。

【請求項 7】

前記識別子が、ハッシュ値に競合が検出された時に割り当てられる連続した番号をさらに含む、
 ことを特徴とする請求項 3 によるストレージシステム。

40

【請求項 8】

前記メタデータが前記ホストコンピュータのOSタイプである、
 ことを特徴とする請求項 1 によるストレージシステム。

【請求項 9】

前記ストレージ制御装置は、前記ホストコンピュータから受け取ったオブジェクト中のデータから識別子を計算する、
 ことを特徴とする請求項 1 によるストレージシステム。

【請求項 10】

50

前記識別子は、前記オブジェクト中の前記データのSecure Hash Algorithm (S H A) 値と、前記オブジェクトのデータサイズとを含む、
ことを特徴とする請求項 9 によるストレージシステム。

【請求項 1 1】

前記オブジェクトは、データ、ファイル名を含み、前記メタデータは当該ファイル名のファイルに対するものである、

ことを特徴とする請求項 9 によるストレージシステム。

【請求項 1 2】

前記メタデータは前記ファイル名のファイルに対する特徴を表すものであり、前記ストレージ制御装置は、一つのファイル名に対して複数の当該メタデータを設定することができる、

ことを特徴とする請求項 1 1 によるストレージシステム。

【請求項 1 3】

前記ストレージ制御装置は、複数の識別子テーブルの中から所定の 2 つ以上の識別子テーブルの相関性をチェックし、そのチェック結果に基づいて当該二つ以上の識別子テーブルを統合し、データの重複の削除に基づいて未使用チャンクを生成する、

ことを特徴とする請求項 1 によるストレージシステム。

【請求項 1 4】

ストレージ制御装置と、

チャンクプールを形成する複数のチャンクに分割されている複数のストレージ装置と、
前記ストレージ制御装置をホストコンピュータに接続するように構成したネットワーク
インターフェースと、

を備え、

前記ストレージ制御装置はメモリを備え、当該メモリのプログラムを実行してデータ
ストレージボリュームを生成し、前記ネットワークインターフェースを経由して前記データ
ストレージボリュームを前記ホストコンピュータに利用可能にし、

前記メモリは、

前記ホストコンピュータから前記データストレージボリュームに向けられた書込コマン
ドに関連したデータに対応する識別子の複数を夫々有する複数の識別子テーブルと、

前記ホストコンピュータから前記ストレージ制御装置への書込みコマンドに伴うメタ
データの複々と、前記複数の識別子テーブルの夫々の情報とを有し、前記一つ又は複数のメ
タデータを前記複数の識別子テーブルの何れかに割り付けるメタデータテーブルと、

を有し、

前記複数の識別子テーブルの夫々は、前記識別子に対応するデータが格納された前記チ
ャンクの情報が当該識別子に割り付けられて構成されているストレージシステムでの重複
データ削除のための方法であって、

前記ストレージ制御装置は、前記メモリのプログラムを実行して、

前記ホストコンピュータからの前記データストレージボリュームに向けられた書込コ
マンドを受信すると、前記識別子を計算し、

前記メタデータテーブルを参照して、前記複数の識別子テーブルのうち、前記書込みコ
マンドに伴う前記メタデータに割り付けられた識別子テーブルを選択し、

前記選択された識別子テーブルを参照して、前記計算された識別子の固有性を前記チャン
ク毎に確認し、

前記計算された識別子が前記識別子テーブルにないものであれば、

前記書込コマンドに関連した前記データに対応する識別子を生成し、

前記チャンクプールから前記データストレージボリュームの前記書込みコマンドで指定
されたデータ位置へ前記チャンクを割り付け、前記書込コマンドに関連した前記データを
、前記割り付けられたチャンクに格納し、

前記生成された識別子と前記割り付けられたチャンクの情報とに基づいて前記識別子テ
ーブルを更新し、

10

20

30

40

50

一方、前記計算された識別子が前記識別子テーブルに存在するものであれば、前記書込みコマンドに対応するデータを格納しない、

ことを特徴とするストレージシステムでの重複データ削除のための方法。

【発明の詳細な説明】

【技術分野】

【0001】

0001 本発明は、全般的には、データストレージシステムにおいて効率的に記憶容量を使用するために、重複したデータブロックあるいはファイルを削除する方法と装置に関するものである。具体的には、本発明は、メタデータグルーピングによって重複したデータブロックおよびファイルを速やかに発見することに関連するものである。

10

【背景技術】

【0002】

0002 重複削除は、重複したデータストリーム、データブロックあるいはファイルを削除するよう設計された機能であり、データストレージシステムにおいて、非常に効率的なデータ保存を提供するためにデータストレージシステムとバックアップ装置に実施される。1つの実施例では、重複したデータを削除するために、識別子が、データ自体から生成される。識別子はハッシュまたはMD5の、SHA（セキュアハッシュアルゴリズム）として生成することができる。データストレージシステムがデータを受け取ると、識別子が計算される。それから、データストレージシステムは、同じ識別子が識別子テーブルに既に格納されているか否かをチェックする。同じ識別子が識別子テーブルにある場合、受信データはデータストレージシステムに保存されない。識別子テーブルに同じ識別子がない場合には、受信データが保存されることとなる。例えば256ビット以上の十分なビット長を備えたハッシュの利用は、競合をめったに引き起こさないことには注目すべきである。大きなデータ量を管理するには、多数の識別子の値を計算し別に管理せねばならない。この理由で、非常に大量の記憶容量を備えたデータストレージシステムにおいては、識別子を確認するのに長い時間がかかることがありうる。

20

【0003】

0003 例えば、米国特許6,928,526号、題名“効率的なデータストレージシステム”は、重複したデータを削除する方法を開示している。データは、データ自体を使用して生成される識別子によって削除される。

30

【先行技術文献】

【特許文献】

【0004】

【特許文献1】米国特許6,928,526号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

しかしながら、従来技術は、メタデータに関連する複数のグループを定義することにより、データストレージシステムにおいて識別子を確認するためのCPU時間を縮小する方法及びシステムを提供していない。

40

【0006】

0004 本発明の実施形態は、先行技術における前述の1つ以上の不備に取り組み、データストレージシステムまたはバックアップの装置において、重複したデータを削除する機能によって、識別子を確認するためのCPU時間を縮小する方法と装置を提供する。

【課題を解決するための手段】

【0007】

0005 本発明の技法の一態様に従って、データストレージボリューム、データストレージボリュームに関連したメタデータを格納するメモリ、ストレージシステムをホストコンピュータに接続するように構成されたネットワークインターフェース、中央処理装

50

置、を備えたストレージシステムが準備される。ストレージシステムは、ホストコンピュータから受け取ったデータから識別子を計算し、データがデータストレージボリュームに格納されているか否かを識別子とメタデータによって決定する。

【0008】

0006 本発明の技法の別の態様に従って、データストレージボリューム、データストレージボリュームに関連したメタデータを格納するメモリ、ストレージシステムをホストコンピュータに接続するように構成されたネットワークインターフェース、中央処理装置、を備えたストレージシステムが準備される。ストレージシステムは、ホストコンピュータから受け取ったオブジェクト中のデータから識別子を計算し、識別子、メモリに格納されているメタデータ、およびオブジェクト中のメタデータによってデータストレージ

10

【0009】

0007 さらに、本発明の技法の別の態様に従って、チャンクプールを形成する多数のチャンクに分割されている多数のデータストレージ装置、ストレージシステムをホストコンピュータに接続するように構成されたネットワークインターフェース、中央処理装置とメモリを含むストレージコントローラ、を備えたストレージシステムにおいて実行される方法が提供される。発明の方法は、データストレージボリュームを供給し、ネットワークインターフェースによってホストコンピュータにデータストレージボリュームを利用可能にするステップと、ホストコンピュータからのデータストレージボリュームに向けられた書込コマンドを受信すると、書込コマンドに関連したデータに対応する識別子を計算

20

【0010】

0008 本発明に関連するさらなる態様は、一部分は以後の記述で述べられ、一部分は記述から明白になるか、もしくは発明の実施によって知るところとなるであろう。本発明の態様、特に、次の詳細な記述および最後尾に設けた請求項において指摘した態様は、要素と諸要素の組合せによって実現し達成することが出来る。

30

【0011】

0009 先の記述および以下の記述の両方ともに、典型的なものであって説明のためだけのものであり、主張している発明や応用を如何なる方法においても制限するようには意図されていないことは、当然理解すべきことである。

【0012】

0010 添付の図面は、この明細書に組み入れられ、この明細書の一部を構成して本発明の実施形態を例証し、記述と共に、本発明の技術の原理について説明し例示する役割を担っている。

具体的には：

【図面の簡単な説明】

40

【0013】

【図1a】0011 図1(a)は、本発明の概念が適用され得る情報システムの典型的な実施形態を示す。

【図1b】図1(b)は、本発明の概念が適用され得る情報システムの典型的な実施形態を示す。

【図2】0012 図2は、データボリューム管理テーブルの典型的な実施形態を示す。

【図3】0013 図3は、チャンクテーブルの典型的な実施形態を示す。

【図4】0014 図4は固有のチャンク識別番号を使用したチャンクの識別を示す。

【図5】0015 図5は、チャンク状態テーブルの典型的な実施形態を示す。

50

【図6】0016 図6(a)および6(b)は、メタデータマッピングテーブルの典型的な実施形態を示す。

【図7】0017 図7は、識別子テーブルの典型的な実施形態を示す。

【図8】0018 図8は、ボリューム生成プロセスの典型的な実施形態を示す。

【図9】0019 図9は、メタデータマッピングテーブルの初期化プロセスの典型的な実施形態を示す。

【図10】0020 図10は、書き込み動作の典型的な実施形態を示す。

【図11】0021 図11(a)は、書き込みコマンドの典型的な実施形態を示す。

0022 図11(b)は、読み出しコマンドの典型的な実施形態を示す。

【図12】0023 図12は、読み出し動作の典型的な実施形態を示す。

【図13a】0024 図13(a)は、本発明の方法が適用することができる情報システムの別の典型的な実施形態を示す。

【図13b】図13(b)は、本発明の方法が適用することができる情報システムの別の典型的な実施形態を示す。

【図14】0025 図14は、ファイル管理テーブルの別の典型的な実施形態を示す。

【図15a】0026 図15(a)は、メタデータマッピングテーブルの別の典型的な実施形態を示す。

【図15b】図15(b)は、メタデータマッピングテーブルの別の典型的な実施形態を示す。

【図15c】図15(c)は、メタデータマッピングテーブルの別の典型的な実施形態を示す。

【図15d】図15(d)は、メタデータマッピングテーブルの別の典型的な実施形態を示す。

【図16】0027 図16は、書き込み動作の別の典型的な実施形態を示す。

【図17】0028 図17(a)は、書き込みコマンドの別の典型的な実施形態を示す。0029 図17(b)は、読み出しコマンドの別の典型的な実施形態を示す。

【図18】0030 図18は、読み出し動作の別の典型的な実施形態を示す。

【図19】0031 図19は、書き込み動作の別の典型的な実施形態を示す。

【図20】0032 図20は、ファイル管理テーブルの別の典型的な実施形態を示す。

【図21】0033 図21は、識別子テーブル統合プログラムによって表示されたスクリーンの典型的な実施形態を示す。

【図22】0034 図22は、識別子テーブル統合プロセスの典型的な実施形態を示す。

【図23】0035 図23は、書き込み動作のための修正済のチャックテーブルを示す。

【図24】0036 図24は書き込みプロセスを示す。ここでの欄は、書き込みコマンドを受け取った後に非同期的に実行される。

【図25】0037 図25は、バックグラウンドでの重複削除のプロセスの典型的な実施形態を示す。

【図26】0038 図26は、本発明のシステムが実施され得るコンピュータプラットフォームの典型的な実施形態を示す。

【発明を実施するための形態】

【0014】

0039 以下の詳細な説明において、添付の図面が参照されるであろうが、その中において、同一の要素は類似の数字で示されている。前述の添付の図面は、本発明の原理と一致する特定の実施形態および実施例を図示の形で示しているのであって、限定するために示しているのではない。これらの実施例は、当分野の業者が発明を実施することを可能にするほどに十分に詳細に記述されている。また、これ以外の実施例も可能であり、構造

10

20

30

40

50

の変更および/または諸要素の置換が本発明の範囲および精神から外れることなく行なう事が可能なことも理解すべきである。したがって、次の詳細な説明は、限定する意味において解釈されるべきではない。

【0015】

0040 発明の第1実施形態のシステム構成を以下記述する。図1は、本発明の方法が適用できる情報システムの典型的な実施形態を示す。第1の実施形態の情報システムは、少なくともホストコンピュータ10、ストレージ装置100、管理コンピュータ500、データネットワーク50およびマネジメントネットワーク90からなる。

【0016】

0041 ホストコンピュータ10について以下記述する。少なくとも1台のホストコンピュータ10がデータネットワーク50によってストレージ装置100に接続される。特にこの実施形態では、6台のホストコンピュータ10a、10b、10c、10d、10eおよび10fが接続されている。少なくとも1つのOS13がホストコンピュータ上で実行される。アプリケーションプログラム14はOS13の上で実行することができる。OS13およびアプリケーションプログラム14のためのファイルおよびデータは、ストレージ装置100によって提供されるデータボリュームに格納される。OS13およびアプリケーションプログラム14は、ストレージ装置100へ書き込みおよび/または読み出しコマンドを出す。ホストコンピュータ10a、10bおよび10cはタイプAのOSを実行し、ホストコンピュータ10d、10eおよび10fはタイプBのOSを実行する。OSタイプはベンダーのOS名およびバージョン番号を使用して定義することが出来る。

【0017】

0042 ストレージ装置100について以下記述する。情報システムは、ストレージコントローラ150および1台以上のHDD101から構成される少なくとも1つのストレージ装置100を備えている。ストレージ装置100は、1つ以上のデータボリューム111をホストコンピュータ10に供給する。

【0018】

0043 管理コンピュータ500について以下記述する。情報システムは、マネジメントネットワーク90によってストレージ装置100を接続している、少なくとも1台の管理コンピュータ500を備えている。

【0019】

0044 データネットワーク50について以下記述する。ホストコンピュータ10およびストレージ装置100はデータネットワーク50によって接続している。この実施形態においては、データネットワーク50はファイバーチャンネルプロトコルを使用し実装されている。しかしながら、イーサネット(登録商標)やインフィニバンドのような他のネットワーク接続も、この目的に同様に使用することができる。ネットワークスイッチとハブはデータネットワーク50の構成要素を互いに接続するために使用することができる。

図1では、ファイバーチャンネルスイッチ55(FCSW55)が、データネットワーク50の構成要素を互いに接続するために使用されている。この目的のために、ホストコンピュータ10およびストレージ装置100はファイバーチャンネルデータネットワーク50に接続するために、1つ以上のファイバーチャンネルインターフェースボード(FCIF)を組込んでいる。

【0020】

0045 マネジメントネットワーク90について以下記述する。ストレージ装置100はマネジメントネットワーク90によって管理コンピュータ500に接続されている。この実施形態においてマネジメントネットワーク90はイーサネット(登録商標)プロトコルを使用し実装されている。しかしながら、他のネットワーク相互接続や、他の接続方法が、この目的に同様に使用することができる。ネットワークスイッチとハブはマネジメントネットワーク90の構成要素を互いに接続するために使用することができ、本発明

10

20

30

40

50

の本実施形態においては、ストレージ装置 100 および管理コンピュータ 500 は、イーサネット（登録商標）マネジメントネットワーク 90 に接続するために 1 つ以上のイーサネット（登録商標）インタフェースボード（イーサ I F）を持っている。

【0021】

0046 ホストコンピュータ 10 について以下詳細に記述する。ホストコンピュータ 10 は、メモリ 12 に格納したプログラムを実行するための CPU 11、プログラムとデータを格納するためのメモリ 12、データネットワーク 50 に接続するための F C I F 15、から構成されている。この実施形態では、CPU 11 は、メモリ 12 に格納された少なくとも 3 つのプログラムを実行する。

【0022】

0047 本発明のこの実施形態においては、メモリ 12 は、オペレーティングシステムプログラム 13（OS 13）、アプリケーションプログラム 14、OS 13 及び / またはアプリケーションプログラム 14 をインストールするためのインストーラプログラム 15、を格納している。

【0023】

0048 管理コンピュータ 500 について以下詳細に記述する。管理コンピュータ 500 は、メモリ 520 に格納したプログラムを実行するための CPU 510、プログラムとデータを格納するためのメモリ 520、マネジメントネットワーク 90 に接続するためのイーサ I F 590、から構成されている。

【0024】

0049 CPU 510 は、メモリ 520 に格納された少なくとも 3 つのプログラムを実行する。この実施形態において、メモリ 520 は、ストレージ装置 100 に、データボリュームを供給することを要求するための、データボリューム供給要求プログラム 521、ストレージ装置 100 への重複削除の範囲の定義の要求を出すための、重複削除範囲定義要求プログラム 522、管理スクリーンを表示し、2 台の識別子テーブルの違いの比率を計算し、識別子テーブル統合の要求を、ストレージ装置 100 に出すための、識別子テーブル統合プログラム 523、を格納している。

【0025】

0050 ストレージ装置 100 について以下詳細に記述する。ストレージ装置 100 は、データを格納するための 1 台以上の HDD 101、およびホストコンピュータへのデータボリュームを提供するための 1 台以上のストレージコントローラ 150 から構成されている。ストレージコントローラ 150 は、メモリ 152 に格納したプログラムを実行するための CPU 151、プログラムおよびデータを格納するためのメモリ 152、データネットワーク 50 に接続するための F C I F 155、HDD 101 に接続するための S A T A I F 156（HDD に F C、S C S I S A S のような別のインターフェースがある場合、適切なインターフェースが実装されねばならない。）、ホストコンピュータから受け取り、HDD から読み出されたデータを格納するためのキャッシュ 153、マネジメントネットワーク 90 に接続するためのイーサ I F 159、から構成されている。

【0026】

0051 CPU 151 は、メモリ 152 に格納された少なくとも 4 つのプログラムを実行する。この実施形態においては、メモリ 152 は、ホストコンピュータ 10 からの、少なくとも読み出し / 書込みコマンドに応答するための応答プログラム 161、ボリュームを生成し、ホストコンピュータ 10 にそれを割り付けるためのデータボリューム割付けプログラム 162、メタデータマッピングテーブルを更新するための重複削除範囲定義プログラム 163、ホストコンピュータ 10 から転送されるデータからメッセージダイジェスト 5（MD5）を計算し、同じ識別子が識別子テーブルに既に格納されているか否かを確認するための重複削除プログラム 165 を格納している。同じ識別子が見つからない場合には、重複削除プログラム 165 はデータを格納する。同じ識別子が見つかった場合には、重複削除プログラム 165 はデータを格納しない。この技術分野における同業者には、本発明のシステムは、MD5 アルゴリズムにのみ限定されるものではないことは理

10

20

30

40

50

解されよう。SHAのような、他の適切なハッシュあるいはダイジェスト機能も、本発明に同様に適用することができる。

【0027】

0052 重複削除プログラム165におけるMD5を計算し識別子機能を確認する機能は、計算と確認の加速のためにゲートアレイあるいはFPGA（フィールドプログラマブルゲートアレイ）上にハードウェアロジックとして実装することもできる。

【0028】

0053 データボリューム管理テーブルについて以下詳細に記述する。データボリューム管理テーブル166は生成されたボリュームを管理する。図2に示されるように、データボリューム管理テーブル166は、ボリューム番号を格納するための“ボリューム番号”列16601、データボリュームサイズ（ブロックの数）の格納のための“サイズ”列16602、ボリューム生成時に、管理コンピュータ500から付与するメタデータの格納のための“メタデータ”列16603、チャンクテーブル番号の格納のための“チャンクテーブル番号”列16604、から構成されている。

【0029】

0054 チャンクテーブルについて以下詳細に記述する。チャンクテーブル167は、データボリュームのLBA（論理的ブロックアドレス）と識別子の関係についての情報を格納する。識別子はMD5および別のシーケンシャル番号から構成される。シーケンシャル番号は、MD5に競合が生じた場合に使用される。各データボリュームはそれぞれ自分のチャンクテーブル167を持つ。LBAは、アクセスブロック位置を指定するために、ホストコンピュータによって読み/書き動作に使用される。図3に示されるように、チャンクテーブル167は、LBAを格納するための“LBA”列16701および識別子を格納するための“識別子”列16702から構成されている。

【0030】

0055 チャンク状態表について以下詳細に説明する。ストレージ装置100のHDDは多数のデータブロックで構成されている。データブロックのサイズは、通常、最近の製品では512バイトである。この実施形態では、チャンクは1つのデータブロックからなる。したがって、この実施形態では、チャンクのサイズは512バイトである。しかし、他のどのようなチャンクサイズも同様に使用することができる。チャンクはそれぞれ図4に示すように、一つずつ識別するために固有の番号を持っている。チャンク状態テーブル170は、チャンクの使用状況を管理する。図5は、この実施形態におけるチャンク状態表170の典型的な実施例を示す。チャンク状態表の典型的な実施形態はチャンク番号を格納するための“チャンク番号”列17001、チャンクが使用されているか否かを示す格納状態情報のための“状態”列17002、から構成されている。データが保存されている場合、重複削除プログラム165はチャンク状態表170を使用して、未使用のチャンクを探す。

【0031】

0056 メタデータマッピングテーブルについて以下詳細に記述する。メタデータマッピングテーブル180は、メタデータと識別子テーブル185の関係についての情報を格納する。データボリュームに書かれたデータは、メタデータマッピングテーブル180で指定された識別子テーブルによって確認される。図6に示されるように、メタデータマッピングテーブル180は、ボリュームに割り当てられたメタデータを保持するための“メタデータ”列18001と、識別子テーブル番号を保持するための“識別子テーブル番号”列18002と、から構成されている。

【0032】

0057 識別子テーブルについて以下詳細に記述する。識別子テーブル185は、チャンクの識別子、参照カウントおよびチャンク番号を管理する。識別子はMD5および別のシーケンシャル番号から構成される。MD5に競合が生じた場合に、シーケンシャル番号が使用される。図7に示すように、識別子テーブル185は、識別子を格納するための「識別子」列18501、チャンクテーブル167からの参照の数である参照カウント

10

20

30

40

50

の格納のための“参照カウント”列18502、データが保存されているチャンク番号を格納するための“チャンク番号”列18503、から構成されている。

【0033】

0058 ボリューム生成プロセスについて以下詳細に記述する。ボリューム生成プロセス800は、図8を参照して説明する。

【0034】

0059 ステップ810：データボリュームの供給要求プログラム521は、ストレージコントローラ150上のデータボリューム割付けプログラム162に、サイズ(ブロックの数)およびメタデータとともに、データボリューム供給要求を発行する。(メタデータは、“OSタイプA”、“OSタイプB”などのようなものである。)任意のタイプのメタデータを管理者が指定することができる。

10

【0035】

0060 ステップ820：データボリューム割付けプログラム162はデータボリューム管理テーブル166を更新する。図2は、6つのボリューム生成が完了した場合を示す。データボリューム111a、111bおよび111cは“OSタイプA”メタデータを伴って生成される。データボリューム111d、111eおよび111fは“OSタイプB”メタデータを伴って生成される。

【0036】

0061 メタデータマッピングテーブルを初期化するプロセスについて以下詳細に記述する。メタデータマッピングテーブル初期化プロセス900は、図9を参照して説明する。

20

【0037】

0062 ステップ910：重複削除範囲定義要求プログラム522は、ストレージコントローラ150上の重複削除範囲定義プログラム163へメタデータと共に重複削除範囲定義要求を出す。

【0038】

0063 ステップ920：重複削除範囲定義プログラム163は新しいメタデータを受け取り、重複削除範囲定義プログラム163は受信したメタデータのために新しい識別子テーブル185を割り付ける。

【0039】

30

0064 ステップ930：重複削除範囲定義プログラム163はメタデータマッピングテーブルを更新する。図6(a)は、2つのメタデータ(“OSタイプA”と“OSタイプB”)を受け取った場合を示す。識別子テーブル185aはメタデータ“OSタイプA”に割り付けられる。識別子テーブル185bはメタデータ“OSタイプB”に割り付けられる。

【0040】

0065 書き込み動作について以下詳細に記述する。書き込み動作1000を、図10を参照して説明する。書き込み動作は、OSとアプリケーションソフトウェアがインストールされ実行されている間に実行される。この実施形態では、OSタイプAがデータボリューム111a、111bおよび111cにインストールされている。OSタイプBが、データボリューム111d、111eおよび111fにインストールされている。非常に多数の書き込み動作が、実際のOSのインストール動作の間に実行されている事は注意すべきである。図10は、応答プログラム161および重複削除プログラム165における典型的な処理フローを示す。ホストコンピュータ10にあるインストーラプログラム15は、書き込みコマンドおよびデータをボリューム111へ発行する。この実施形態では、ホストコンピュータ10aがデータボリューム111aを使用し、ホストコンピュータ10bがデータボリューム111bを使用している、等。

40

【0041】

0066 ステップ1010：書き込みコマンドとデータを受け取る。書き込みコマンドはLBA、およびブロックの数についての情報を含んでいる。図11(a)は書き込

50

みコマンドを示す。書き込みコマンドはコマンド種別 (=書き込み)、L B A (=データの位置)、およびブロックの数 (=データのサイズ) についての情報を含んでいる。

【 0 0 4 2 】

0 0 6 7 ステップ 1 0 1 2 : ボリュームのメタデータはデータボリューム管理テーブル 1 6 6 から得られる。(書き込みコマンドを受け取る現在のデータボリュームが 1 1 1 a である場合、メタデータは " O S タイプ A " である (図 2 を参照)。)

【 0 0 4 3 】

0 0 6 8 ステップ 1 0 1 4 : 識別子テーブルはメタデータによりメタデータマッピングテーブル 1 8 0 から選ばれる。(書き込みコマンドを受け取る現在のデータボリュームが 1 1 1 a である場合、メタデータが " O S タイプ A " であるので、識別子テーブル 1 8 5 a が選択されている。) 下記ステップが各ブロックに対し実行される。

【 0 0 4 4 】

0 0 6 9 ステップ 1 0 1 6 : M D 5 値が受け取ったデータから計算される。

【 0 0 4 5 】

0 0 7 0 ステップ 1 0 1 8 : 書き込みコマンドによって指定された現在のデータ位置が割り付けられたチャンクを有するか否かのチェックが、チャンクテーブルを参照することにより行なわれる。チャンクが既に割り付けられている場合、プロセスはステップ 1 0 6 0 に進む。チャンクが割り付けられない場合、プロセスはステップ 1 0 2 0 に進む。

【 0 0 4 6 】

0 0 7 1 ステップ 1 0 2 0 : 選択された識別子テーブルからの同じ M D 5 の値を持っている欄が列挙される。同じ M D 5 の値が見つかった場合、プロセスはステップ 1 0 4 0 に進む。同じ M D 5 の値が見つからなかった場合、プロセスはステップ 1 0 2 2 に進む。

【 0 0 4 7 】

0 0 7 2 ステップ 1 0 2 2 : 識別子が生成される。識別子は M D 5 の値と 0 の組合せである。

【 0 0 4 8 】

0 0 7 3 ステップ 1 0 2 4 : チャンク状態表 1 7 0 によって未使用のチャンクを得る。

【 0 0 4 9 】

0 0 7 4 ステップ 1 0 2 6 : 使用状態を示すためのチャンク状態テーブル 1 7 0 を更新する。

【 0 0 5 0 】

0 0 7 5 ステップ 1 0 2 8 : チャンクヘデータを格納する。

【 0 0 5 1 】

0 0 7 6 ステップ 1 0 3 0 : 選択された識別子テーブル 1 8 5 を更新する。(ステップ 1 0 2 2 において) 生成された識別子が格納される。参照カウンタは 1 にセットされる。チャンク番号 (ステップ 1 0 2 2 で得られた) が格納される。

【 0 0 5 2 】

0 0 7 7 ステップ 1 0 3 2 : チャンクテーブル 1 6 7 を更新する。(書込コマンドを受け取る現在データボリュームが 1 1 1 a である場合、チャンクテーブル 1 6 7 a が更新される。) 生成された識別子は、現在の L B A 欄に格納される。

【 0 0 5 3 】

0 0 7 8 ステップ 1 0 4 0 : M D 5 の競合を回避するために、バイト単位データチェックを実行する。データが異なる場合、動作はステップ 1 0 5 0 に進む。データが異なる場合 (M D 5 の競合)、ステップ 1 0 4 2 に進む。(複数の同じ M D 5 が見つかった場合、バイト単位チェックを各々に実行する。)

【 0 0 5 4 】

0 0 7 9 ステップ 1 0 4 2 : 識別子を生成する。識別子は、M D 5 と現在の最大番号の次の番号となる新しいシーケンシャル番号の組合せである。(データが異なっても、

10

20

30

40

50

1つの同じMD5が見つかった場合、シーケンシャル番号は1である。データが異なっても、2つの同じMD5が見つかった場合、シーケンシャル番号は2である。)ステップ1024に進む。MD5はめったに競合しない。したがって、ステップ1042はまれな場合である。

【0055】

0080 ステップ1050：選択された識別子テーブル185を更新する。識別子に従い識別子テーブル185中の参照カウントを増加させる。

【0056】

0081 ステップ1052：チャンクテーブル167を更新する。識別子はLBAに対応して格納される。

10

【0057】

0082 ステップ1054：データを廃棄する。

【0058】

0083 ステップ1060：識別子に従い識別子テーブル185の参照カウントを減少させる。

【0059】

0084 ステップ1062：参照カウントが0であるか否かをチェックする。参照カウントが0である場合には、ステップ1064に進む。参照カウントが0でない場合には、ステップ1020に進む。

【0060】

0085 ステップ1064：チャンクを空ける。

20

【0061】

0086 ステップ1066：選択された識別子テーブル185を更新する。識別子の欄は選択された識別子テーブルから消去される。ステップ1020に進む。

【0062】

0087 上に言及したように、特定の識別子テーブルはデータボリュームに割り当てられるメタデータに従って使用される。1つの識別子テーブルは複数のデータボリュームに使用することができる。データボリュームは、種々のOSを保存することができる。しかしながら、種々のOSは種々のデータを含んでいるので、1つの識別子テーブルでは巨大になってしまう。1つの巨大な識別子テーブルで識別子確認を実行するには長い時間がかかる。この実施形態では、複数の識別子テーブルが使用される。1つの識別子テーブルはOSタイプAに使用され、もう一方はOSタイプBに使用される。各識別子テーブルは、2つのOSに対し1つの識別子テーブルの場合より小さくなる。メタデータによる分離の結果、識別子確認を実行するためのCPU時間はより短くなる。本発明を備えたデータストレージ装置の読み/書きアクセスの性能は改善される。識別子を確認するためのどのようなアルゴリズムでも、例えば二分木のようなものでも、使用することができる。

30

【0063】

0088 読み出し動作について以下詳細に記述する。読み出し動作1200を図12で説明する。OSとアプリケーションが実行されている間に読み出し動作は実行される。ホストコンピュータ10上の、OS13、およびアプリケーションプログラム14がボリューム111へ読み出しコマンドを発行する。この実施形態では、ホストコンピュータ10aはデータボリューム111aを、ホストコンピュータ10bはデータボリューム111bを、・・・使用する。

40

【0064】

0089 ステップ1210：読み出しコマンドを受け取る。読み出しコマンドは、LBAおよびブロックの数を含んでいる。図11(b)は読み出しコマンドを示す。読み出しコマンドは、コマンド種別(=読み出し)、LBA(=データの位置)、ブロックの数(=データのサイズ)を含んでいる。

【0065】

0090 ステップ1212：データボリューム管理テーブル166からボリューム

50

のメタデータを得る。(読み出しコマンドを受け取る現在データボリュームが111aである場合、メタデータは“OSタイプA”である(図2)。)

【0066】

0091 ステップ1214:メタデータに従いメタデータマッピングテーブル180から識別子テーブルを選択する。(読み出しコマンドを受け取る現在データボリュームが111aである場合、メタデータが“OSタイプA”であるので、識別子テーブル185aが選択される。)

【0067】

0092 ステップ1216:チャンクテーブル167を選択する。(読み出しコマンドを受け取る現在データボリュームが111aである場合、チャンクテーブル167aが選択される。)各ブロックに対し下記のステップが実行される。

10

【0068】

0093 ステップ1220:LBAに対応する選択されたチャンクテーブル167から識別子を得る。

【0069】

0094 ステップ1222:識別子に対応している識別子テーブル185からチャンク番号を得る。

【0070】

0095 ステップ1224:チャンク番号が指定されたチャンクのデータを、ホストコンピュータへ転送する。

20

【0071】

0096 本発明の第2の実施形態について以下詳細に説明する。図13は、本発明の方法が適用される情報システムの全体図の例を示す。第1の実施形態と第2の実施形態の違いは以下の通りである。

【0072】

0097 ストレージ装置100はファイル管理プログラム164を有する。

【0073】

0098 ホストコンピュータ10は、ファイルアクセスコマンド(図17に示される)によってファイルストレージ装置へのアクセスを行う。

【0074】

0099 データネットワークはイーサネット(登録商標)である。イーサネット(登録商標)スイッチ85(Ether SW 85)が相互接続のために使用されている。

30

【0075】

0100 ホストコンピュータ10がイーサネット(登録商標)データネットワーク80の接続のためにイーサネット(登録商標)インタフェースボード18(イーサIF18)を持っている。

【0076】

0101 ストレージ装置100がイーサネット(登録商標)データネットワーク80の接続のためにイーサネット(登録商標)インタフェースボード158(イーサIF158)を持っている。

40

【0077】

0102 ファイル管理プログラム164はファイル管理テーブル190を使用する。

ファイル管理テーブルはファイル名、ファイルサイズ、識別子およびメタデータの管理のために使用する。

【0078】

0103 識別子テーブル185中の各欄はそれぞれ複数のチャンク番号を保持することができる。

【0079】

0104 第2の実施形態のファイル管理テーブルについて以下詳細に説明する。図

50

14に示すように、ファイル管理テーブル190は、ファイル名を保持するための“ファイル名”列19001（この実施形態では、ファイル名はディレクトリー名とファイル名からなる）、ファイルサイズを保持するための“ファイルサイズ”列19002、MD5および別の連続番号からなる識別子を保持するための“識別子”列19003、メタデータの保持のための“メタデータ”列19004、とから構成されている。1つのファイル名に対応する複数のメタデータを保持することができる。

【0080】

0105 メタデータマッピングテーブルの初期化について以下詳細に説明する。メタデータマッピングテーブルの初期化のプロセスは第1の実施形態と同じである。しかしながら、生成者名、グループ名、組織名称、作成の時間スタンプ、生成されたアプリケーション名、ファイルタイプ、ホストコンピュータ名（物理的な及び/または仮想の）などなど、多様なメタデータを使用することができる。図15は、この実施形態におけるメタデータマッピングテーブル180の例を示す。図15(a)はファイルタイプによるメタデータマッピングテーブルを示す。図15(b)は組織名称によるメタデータマッピングテーブルを示す。図15(c)は作成時の時間スタンプによるメタデータマッピングテーブルを示す。

10

【0081】

0106 第2の実施形態の書き込み動作について以下詳細に記述する。

【0082】

0107 書き込み動作1600を図16により説明する。書き込み動作はOS13およびアプリケーションソフト14から実行される。

20

【0083】

0108 ステップ1610：書き込みコマンドとデータを受け取る。書き込みコマンドはファイル名、ファイルサイズおよびメタデータを含んでいる。図17(a)は書き込みコマンドを示す。

【0084】

0109 ステップ1616：データからのMD5を計算する。

【0085】

0110 ステップ1618：同じファイル名が既にファイル管理テーブル190に存在するか否かチェックする。同じファイル名が見つかった場合には、1660に進む。同じファイル名が見つからなかった場合には、ステップ1620に進む。

30

【0086】

0111 ステップ1620：書き込みコマンドに保持されているメタデータに従いメタデータマッピングテーブル180から識別子テーブルを選ぶ。

【0087】

0112 ステップ1621：選択された識別子テーブルから同じMD5の値を持っている欄を列挙する。同じMD5が見つかった場合には、ステップ1640に進む。同じMD5が見つからなかった場合には、ステップ1622に進む。

【0088】

0113 ステップ1622：識別子を生成する。識別子はMD5とゼロの組合せである。

40

【0089】

0114 ステップ1624：チャンク状態テーブルおよびファイルのサイズ（サイズはチャンクの数に変換されている）によって未使用のチャンクを得る。

【0090】

0115 ステップ1626：使用状態を示すためにチャンク状態表170を更新する。

【0091】

0116 ステップ1628：チャンクヘデータを格納する。

【0092】

50

0117 ステップ1630：選択された識別子テーブル185を更新する。生成された識別子を格納する。参照カウントは1にセットされる。チャンク番号（ステップ1624で得られた）を格納する。

【0093】

0118 ステップ1632：ファイル管理テーブル190を更新する。ファイル名、ファイルサイズ、生成された識別子（ステップ1622において）およびメタデータを格納する。

【0094】

0119 ステップ1640：MD5の競合を回避するためにバイト単位データチェックを行う。データが異なる場合、ステップ1650に進む。データが異なる場合には（MD5競合）、ステップ1642に進む。（複数の同じMD5が見つかった場合、各々に対しバイト単位チェックを行う。）

10

【0095】

0120 ステップ1642：識別子を生成する。識別子は、MD5と現在の最大番号の次のシーケンシャル番号との組合せである。（データが異なっても、1つの同じMD5が見つかった場合、シーケンシャル番号は1である。データが異なっても、2つの同じMD5が見つかった場合、シーケンシャル番号は2である。）

【0096】

0121 ステップ1650：選択された識別子テーブル185を更新する。識別子に従い識別子テーブル185の参照カウントを増加させる。

20

【0097】

0122 ステップ1652：ファイル管理テーブル190を更新する。ファイル名、ファイルサイズ、生成された識別子およびメタデータを格納する。

【0098】

0123 ステップ1654：データを廃棄する。

【0099】

0124 ステップ1660：ファイル管理テーブルに保持されたメタデータによって識別子テーブルを選択する。

【0100】

0125 ステップ1661：ファイル名に従いファイル管理テーブルから識別子を得る。

30

【0101】

0126 ステップ1662：識別子に従い識別子テーブル185の参照カウントを減少させる。

【0102】

0127 ステップ1663：参照カウントが0であるか否かチェックする。参照カウントが0である場合には、ステップ1664に進む。参照カウントが0でない場合は、ステップ1620に進む。

【0103】

0128 ステップ1664：識別子とともに欄に保持されているチャンク番号のチャンクを空ける。

40

【0104】

0129 ステップ1666：選択された識別子テーブル185を更新する。識別子を有する欄は選択された識別子テーブルから消去する。ステップ1620に進む。

【0105】

0130 上に述べたように、特定の識別子テーブルはファイルに割り当てられたメタデータに従って使用される。この実施形態では、複数の識別子テーブルが同様にメタデータマッピングテーブル180に従い使用されている。各識別子テーブルは、すべてのファイルに対し1つの識別子テーブルの場合よりもより小さくなる。メタデータによる分離の結果、識別子確認を実行するためのCPU時間はより短くなる。本発明を備えたデータ

50

ストレージ装置の読み / 書きアクセスの性能が改善される。

【 0 1 0 6 】

0 1 3 1 第 2 の実施形態の読み出し動作について以下詳細に記述する。読み出し動作 1 8 0 0 は図 1 8 で説明する。読み出し動作は OS 1 3 およびアプリケーションソフトウェア 1 4 から実行される。

【 0 1 0 7 】

0 1 3 2 ステップ 1 8 1 0 : 読み出しコマンドを受け取る。読み出しコマンドはファイル名を含んでいる。図 1 7 (b) は読み出しコマンドを示す。

【 0 1 0 8 】

0 1 3 3 ステップ 1 8 1 2 : ファイル名に従いファイル管理テーブル 1 9 0 からファイルのメタデータを得る。 10

【 0 1 0 9 】

0 1 3 4 ステップ 1 8 1 4 : メタデータに従いメタデータマッピングテーブル 1 8 0 から識別子テーブルを選ぶ。

【 0 1 1 0 】

0 1 3 5 ステップ 1 8 2 0 : ファイル名に従いファイル管理テーブル 1 9 0 から識別子を得る。

【 0 1 1 1 】

0 1 3 6 ステップ 1 8 2 2 : 識別子に従い識別子テーブル 1 8 5 からチャンク番号を得る。 20

【 0 1 1 2 】

0 1 3 7 ステップ 1 8 2 4 : ファイル名、ファイルサイズおよびメタデータをホストコンピュータへ転送する。

【 0 1 1 3 】

0 1 3 8 ステップ 1 8 2 6 : チャンク番号を指定したチャンクのデータをホストコンピュータへ転送する。

【 0 1 1 4 】

0 1 3 9 書き込み動作の変形について以下詳細に記述する。書き込み動作 1 9 0 0 の別の実施形態を図 1 9 を参照して説明する。書き込み動作は OS 1 3 およびアプリケーションソフトウェア 1 4 から実行される。この書き込み動作では、ファイルは MD 5 およびファイルのサイズによって重複削除される。識別子の一部であるシーケンシャル番号は使用しない。識別子は MD 5 およびファイルのサイズから構成される。図 2 0 に示すファイル管理テーブル 1 9 0 を書き込み動作に使用する。 30

【 0 1 1 5 】

0 1 4 0 ステップ 1 9 1 0 : 書き込みコマンドとデータを受け取る。書き込みコマンドはファイル名、ファイルサイズおよびメタデータを含んでいる。図 1 7 (a) は書き込みコマンドを示す。

【 0 1 1 6 】

0 1 4 1 ステップ 1 9 1 6 : データからの MD 5 を計算する。

【 0 1 1 7 】

0 1 4 2 ステップ 1 9 1 8 : 同じファイル名がファイル管理テーブル 1 9 0 に既に存在するか否かチェックする。同じファイル名が見つかった場合には、ステップ 1 9 6 0 に進む。同じファイル名が見つからなかった場合には、ステップ 1 9 2 0 に進む。 40

【 0 1 1 8 】

0 1 4 3 ステップ 1 9 2 0 : 書き込みコマンドに保持されているメタデータに従いメタデータマッピングテーブル 1 8 0 から識別子テーブルを選ぶ。

【 0 1 1 9 】

0 1 4 4 ステップ 1 9 2 1 : 同じ識別子 (MD 5 の値、ファイルのサイズ) が選択された識別子テーブルに存在しているか否かを確認する。同じ識別子が見つかった場合には、ステップ 1 9 5 0 に進む。同じ識別子が見つからなかった場合には、ステップ 1 9 2 50

2に進む。

【0120】

0145 ステップ1922：識別子を生成する。識別子はMD5とファイルサイズの組合せである。

【0121】

0146 ステップ1924：チャンク状態テーブルおよびファイルのサイズ（サイズはチャンクの数に変換されている）に従い未使用のチャンクを得る。

【0122】

0147 ステップ1926：使用状況を示すためにチャンク状態表170を更新する。

10

【0123】

0148 ステップ1928：チャンクヘータを格納する。

【0124】

0149 ステップ1930：選択した識別子テーブル185を更新する。識別子（MD5の値、ファイルのサイズ）を格納し、また、チャンク番号も格納する。参照カウントを1にセットする。

【0125】

0150 ステップ1932：ファイル管理テーブル190を更新する。ファイル名、ファイルサイズ、MD5およびメタデータを格納する。

【0126】

20

0151 ステップ1950：選択された識別子テーブル185を更新する。識別子に従い識別子テーブル185の参照カウントを増加させる。

【0127】

0152 ステップ1952：ファイル管理テーブル190を更新する。ファイル名、ファイルサイズ、MD5およびメタデータを格納する。

【0128】

0153 ステップ1954：データを廃棄する。

【0129】

0154 ステップ1960：ファイル管理テーブルに保持されたメタデータに従い識別子テーブルを選択する。

30

【0130】

0155 ステップ1961：ファイル名に従って識別子（MD5の値およびファイルのサイズ）を得る。

【0131】

0156 ステップ1962：識別子に従って識別子テーブル185参照カウントを減少させる。

【0132】

0157 ステップ1963：参照カウントが0であるか否かチェックする。参照カウントが0の場合には、ステップ1964に進む。参照カウントが0でない場合には、ステップ1920に進む。

40

【0133】

0158 ステップ1964：識別子と共に、欄に保持されているチャンク番号のチャンクを空ける。

【0134】

0159 ステップ1966：選択された識別子テーブル185を更新する。識別子を有する欄を選択した識別子テーブルから消去する。ステップ1920に進む。

【0135】

0160 識別子テーブル統合について以下詳細に記述する。管理コンピュータ500は識別子テーブル統合プログラム523を持つことが出来る。識別子テーブル統合プログラム523は図21に示すスクリーン2100を表示する。メタデータおよび識別子テ

50

ーブル番号はスクリーン 2 1 0 0 上のテーブル 2 1 0 1 に表示される。管理者はチェックボックスのチェックにより違いを比較するために 2 つのテーブルを選ぶことができる。図 2 1 では、メタデータ “アカウントティング” および “HR” が選択されている。識別子テーブル統合プログラムは 2 つの識別子テーブル間の相関性をチェックし、違いの比率を計算する。違いの比率は表示枠 2 1 0 2 に表示される。違いが大きい場合には、2 つのテーブルを統合すべきでない。統合はより長い待ち時間を引き起こすかも知れない。小さな違いの場合には、2 つのテーブルを統合してもよい。統合は重複削除により未使用チャンクを生み出す。管理者は 2 つのテーブルの統合のためにボタン 2 1 0 3 を押すことができる。識別子テーブル統合プログラム 5 2 3 は、ストレージ装置 1 0 0 へ識別子テーブル統合要求を出す。識別子テーブルの統合プロセスは、管理者からの要求なしの自動化にすることもでき得る。予め決定のしきい値を管理者がセットし、識別子テーブル統合プロセスの起動を決定するための比率と比較することができる。また、ストレージに十分なデータがあれば、比率が安定した量になるので、ストレージの使用状況のレベルを考慮すべきであろう。ストレージ装置 1 0 0 は、図 2 2 に示すように識別子テーブル統合プロセスをスタートする。

【 0 1 3 6 】

0 1 6 1 ステップ 2 2 1 0 : 識別子テーブル統合要求で指定された第 1 の識別子テーブルを選択する。

【 0 1 3 7 】

0 1 6 2 ステップ 2 2 1 2 : 識別子テーブル統合リクエストで指定された第 2 の識別子テーブルを選択する。

【 0 1 3 8 】

0 1 6 3 ステップ 2 2 1 4 : $i = 0$ と初期化する。

【 0 1 3 9 】

0 1 6 4 ステップ 2 2 2 0 : もし欄 [i] が最後の欄でなければ、2 2 5 0 へスキップする。

【 0 1 4 0 】

0 1 6 5 ステップ 2 2 2 2 : 欄 [i] の識別子 (MD5 の値およびファイルのサイズ) が第 1 の識別子テーブルに存在することを確認する。同じ識別子が存在する場合、ステップ 2 2 3 0 に進む。(第 1 の識別子テーブルの識別子を保持する欄の番号に j をセットする)。

【 0 1 4 1 】

0 1 6 6 ステップ 2 2 2 4 : 欄 [i] を第 1 の識別子テーブルに加える。

【 0 1 4 2 】

0 1 6 7 ステップ 2 2 3 0 : 欄 [i] に、チャンク番号が保持されているチャンクを空ける。

【 0 1 4 3 】

0 1 6 8 ステップ 2 2 3 2 : 次の方法で、第 1 の識別子テーブルの欄 [j] の参照カウントに第 2 の識別子テーブルの欄 [i] の参照カウント値を加える。(欄 [j] 参照カウント = 欄 [j] 参照カウント + 欄 [i] 参照カウント)。

【 0 1 4 4 】

0 1 6 9 ステップ 2 2 4 0 : i を増加させる。ステップ 2 2 2 0 に進む。

【 0 1 4 5 】

0 1 7 0 ステップ 2 2 5 0 : メタデータ管理テーブルを更新する (例えば、図 1 5 (a) は図 1 5 (d) として更新される)

【 0 1 4 6 】

0 1 7 1 ステップ 2 2 5 2 : 第 2 の識別子テーブルを削除する。

【 0 1 4 7 】

0 1 7 2 第 1 の実施形態の書き込み動作の代替実施形態について以下記述する。書き込み動作 2 3 0 0 の別の実施形態は図 2 3、2 4 および 2 5 を参照して説明する。書き

10

20

30

40

50

込み動作について、修正済のチャンクテーブルを、図 2 3 に示している。チャンクテーブルはチャンク番号情報を格納する。L B A 欄にチャンク番号がある場合、それは、チャンク番号で指定したチャンクに L B A に対するデータが一時的に格納されたことを意味する。この欄は、書込コマンドを受け取った後に非同期に実行されるであろう。対応する書き込みプロセスは図 2 4 に示されている。

【 0 1 4 8 】

0 1 7 3 ステップ 2 4 1 0 : 書込コマンドを受け取る。書込コマンドは、L B A およびブロックの番号を含んでいる。

【 0 1 4 9 】

0 1 7 4 ステップ 2 4 1 2 : チャンク状態テーブル 1 7 0 に従って未使用のチャンクを得る。 10

【 0 1 5 0 】

0 1 7 5 ステップ 2 4 1 4 : チャンクが使用されていることを示すためにチャンク状態テーブル 1 7 0 を更新する。

【 0 1 5 1 】

0 1 7 6 ステップ 2 4 1 6 : チャンクヘデータを格納する。

【 0 1 5 2 】

0 1 7 7 ステップ 2 4 1 8 : チャンクテーブルを更新する。この時に、ステップ 2 4 1 2 で受け取ったチャンク番号だけを、チャンクテーブルに格納する。

【 0 1 5 3 】

0 1 7 8 重複削除の評価は、書込コマンドを受け取った後に非同期に実行される。バックグラウンド重複削除プロセスを図 2 5 に示す。バックグラウンド重複削除は、各データボリュームで周期的に実行される。図 1 0 と図 2 5 との差はステップ 2 5 1 5 の中みにある。 20

【 0 1 5 4 】

0 1 7 9 ステップ 2 5 1 5 : テーブルから、一時的に格納したデータを列挙する。一時的に格納したデータは、チャンクテーブルのチャンク番号のチェックにより確認される。一時的に格納したデータは、チャンクテーブルにチャンク番号を有する。例えば、図 2 3 において、L B A 2 のデータはチャンク 1 0 0 0 3 に一時的に格納される。ステップ 2 5 1 5 に続くステップは、各ブロックに対し実行される。 30

【 0 1 5 5 】

典型的なコンピュータプラットフォーム

0 1 8 0 図 2 6 は、発明の技法の実施形態を実装することが出来るコンピュータ/サーバシステム 2 6 0 0 の実施形態を示すブロック図である。システム 2 6 0 0 は、コンピュータ/サーバプラットフォーム 2 6 0 1、周辺機器 2 6 0 2 およびネットワークリソース 2 6 0 3 を備えている。

【 0 1 5 6 】

0 1 8 1 コンピュータプラットフォーム 2 6 0 1 は、コンピュータプラットフォーム 2 6 0 1 の様々な部分の間で情報通信するための、データバス 2 6 0 4 あるいは他のコミュニケーション機構、ならびに、バス 2 6 0 4 と接続した、情報処理や、他の計算処理および制御タスクを行なうプロセッサ 2 6 0 5 を備えている。コンピュータプラットフォーム 2 6 0 1 は、また、プロセッサ 2 6 0 5 によって実行される命令とともに種々の情報を格納するための、バス 2 6 0 4 に接続した、ランダムアクセスメモリ (R A M) あるいは他の動的ストレージ装置のような、揮発性ストレージ装置 2 6 0 6 を備えている。揮発性ストレージ装置 2 6 0 6 は、また、プロセッサ 2 6 0 5 による命令の実行の間に一時変数あるいは他の中間情報を格納するために使用することも出来る。コンピュータプラットフォーム 2 6 0 1 は、さらに、基本入出力システム (B I O S) や様々なシステム構成パラメータの様な、プロセッサ 2 6 0 5 のための静的な情報や命令を格納するために、バス 2 6 0 4 と接続した読み取り専用メモリ (R O M または E P R O M) 2 6 0 7 あるいは他のスタティックストレージ装置を備えることがある。磁気ディスク、光ディス 40 50

クあるいはソリッドステートフラッシュメモリ装置のような持続形ストレージ装置 2608 は、情報と命令の格納のためにバス 2604 に接続し提供される。

【0157】

0182 コンピュータプラットフォーム 2601 は、コンピュータプラットフォーム 2601 のシステム管理者またはユーザへの情報の表示のために、陰極線管 (CRT)、プラズマディスプレイあるいは液晶ディスプレイ (LCD) のようなディスプレイ 2609 に、バス 2604 によってつなぐことができる。アルファニューメリックや他のキーを含む入力装置 2610 が、プロセッサ 2605 に情報やコマンドの選択を伝えるためにバス 2604 に接続されている。別のタイプのユーザ入力デバイスは、プロセッサ 2605 に命令情報とコマンドの選択を伝えるため、あるいはディスプレイ 2609 上でカーソル移動をコントロールするための、マウス、トラックボールあるいはカーソル指示キーのようなカーソル制御デバイス 2611 である。この入力装置は、2つの軸、第一軸 (例えば x) および第二軸 (例えば y) によって、典型的には2つの自由度を持っているが、これによって入力装置は平面上での位置を指定することが可能となる)。

10

【0158】

0183 外部ストレージ装置 2612 は、コンピュータプラットフォーム 2601 に追加のまたはリムーバブルストレージの容量を供給するために、バス 2604 によってコンピュータプラットフォーム 2601 につなぐ事が出来る。計算機システム 2600 の実施形態では、外部リムーバブルストレージ装置 2612 は、他の計算機システムとのデータの交換を容易にするために使用することが出来る。

20

【0159】

0184 本発明は、ここに記述された技術を実装するのに対し計算機システム 2600 を使用することに関係する。ある実施形態では、本発明のシステムはコンピュータプラットフォーム 2601 のようなマシン上で存在することができる。本発明の1つの実施形態によれば、ここに記述された技術は、揮発性メモリ 2606 に含まれている1つ以上の命令の1つ以上のシーケンスを実行するプロセッサ 2605 に応答して計算機システム 2600 によって行なわれる。そのような命令は、持続型ストレージ装置 2608 のような別のコンピュータ可読媒体から揮発性メモリ 2606 に読み込む事が出来る。揮発性メモリ 2606 に含まれている一連の命令の実行は、プロセッサ 2605 に、ここに記述されたプロセスのステップを実行させる。他の実施例においては、ハードワイヤードの回路類が、本発明を実装するためにソフトウェア命令の代わりに、あるいはソフトウェア命令と組み合わせて使用する事が出来る。このように、本発明の実施形態はハードウェア回路とソフトウェアのどんな特定の組合せにも制限されてはいない。

30

【0160】

0185 ここに使用されている用語“コンピュータ可読媒体”は、実行するプロセッサ 2605 に命令を提供することに関与するあらゆる媒体を表す。コンピュータ可読媒体は機械可読媒体の単に1つの例であり、ここに記述された方法及び/または技術のうちのどれであっても実装するために命令を運ぶことが出来る。そのような媒体は、不揮発性のメディアおよび揮発性のメディアを含み、これらに限定されることなく、多くの形式をとることが出来る。不揮発性のメディアは、例えば、ストレージ装置 2608 のように、光ディスクまたは磁気ディスクを含んでいる。揮発性のメディアは、揮発性ストレージ装置 2606 のように、ダイナミックメモリを含んでいる。

40

【0161】

0186 コンピュータ可読媒体の一般の形態は、例えば、フロッピー (登録商標) ディスク、フレキシブルディスク、ハードディスク、磁気テープあるいは何らかの他の磁気メディア、CD-ROM、何らかの他の光学媒体、パンチカード、紙テープ、何らかの他の物理的な孔部のパターンを備えたメディア、RAM、PROM、EPROM、フラッシュ EPROM、フラッシュドライブ、メモリカード、何らかの他のメモリチップあるいはカートリッジ、搬送波、など以下に記述されるように、あるいはコンピュータが読むことができる如何なる他のメディアをも含んでいる。

50

【 0 1 6 2 】

0 1 8 7 コンピュータ読取り可能なメディアの様々な形式は実行用プロセッサ 2 6 0 5 への 1 つ以上の命令の 1 つ以上のシーケンスを運ぶことに関係する。例えば、命令は、当初はリモートコンピュータから磁気ディスク上で運ばれるかも知れない。あるいは、リモートコンピュータはそのダイナミックメモリに命令をロードし、モデムを使用して、電話回線上で命令を送ることができる。計算機システム 2 6 0 0 に対して直接接続のモデムは、電話回線上のデータを受け取ることができ、データを赤外線信号に変換するために赤外線送信機を使用することができる。赤外線検知器は、赤外線信号で運ばれたデータを受け取ることができ、適切な回路類がデータバス 2 6 0 4 にデータを乗せることができる。バス 2 6 0 4 は、データを揮発性ストレージ装置 2 6 0 6 に運び、プロセッサ 2 6 0 5 が命令を検索し実行する。揮発性メモリ 2 6 0 6 によって受け取られた命令は、プロセッサ 2 6 0 5 によって実行する前に、あるいはその実行の後に持続型記憶装置 2 6 0 8 上に格納することも選択肢である。命令は、業界では良く知られた様々なネットワークデータ通信プロトコルを使用して、インターネットを通じてコンピュータプラットフォーム 2 6 0 1 へダウンロードすることも出来る。

10

【 0 1 6 3 】

0 1 8 8 コンピュータプラットフォーム 2 6 0 1 は、また、データバス 2 6 0 4 につながれたネットワークインターフェースカード 2 6 1 3 のような通信インターフェースを有する。通信インターフェース 2 6 1 3 は、ローカルネットワーク 2 6 1 5 につながっているネットワークリンク 2 6 1 4 に双方向のデータ通信接続を提供する。例えば、通信インターフェース 2 6 1 3 は、総合サービスデジタルネットワーク (I S D N) カードか、あるいは対応するタイプの電話回線にデータ通信接続を提供するモデムかもしれない。別の例として、通信インターフェース 2 6 1 3 は、互換性をもつ LAN にデータ通信接続を提供するローカルエリアネットワークインターフェースカード (L A N N I C) かかもしれない。良く知られている、8 0 2 . 1 1 a、8 0 2 . 1 1 b、8 0 2 . 1 1 g およびブルートゥースのような無線リンク、もまたネットワーク実装に使用することが出来る。すべてそのような実装においては、通信インターフェース 2 6 1 3 は、多様な形式の情報を表すデジタルデータストリームを運ぶ、電気的か、電磁氣的か、光学的な信号を送受信する。

20

【 0 1 6 4 】

0 1 8 9 ネットワークリンク 2 6 1 3 は、典型的には、1 つ以上のネットワークを通して他のネットワークリソースにデータ通信を提供する。例えば、ネットワークリンク 2 6 1 4 は、ローカルネットワーク 2 6 1 5 を通して、ホストコンピュータ 2 6 1 6、あるいはネットワークストレージ/サーバ 2 6 2 2 への接続を提供することが出来る。さらに、あるいは代わりに、ネットワークリンク 2 6 1 3 は、ゲートウェイ/ファイアウォール 2 6 1 7 を通してインターネットのような広域またはグローバルネットワーク 2 6 1 8 に接続することも出来る。このように、コンピュータプラットフォーム 2 6 0 1 は、遠隔ネットワークストレージ/サーバ 2 6 1 9 のように、インターネット 2 6 1 8 のいかなる場所に位置したネットワークリソースにもアクセス可能である。一方では、コンピュータプラットフォーム 2 6 0 1 はまた、ローカルエリアネットワーク 2 6 1 5 及び/またはインターネット 2 6 1 8 のいかなる場所に位置したクライアントによってもアクセスすることが出来る。ネットワーククライアント 2 6 2 0 および 2 6 2 1 は、プラットフォーム 2 6 0 1 に類似したコンピュータプラットフォームに基づいて、それら自身実装することが出来る。

30

40

【 0 1 6 5 】

0 1 9 0 ローカルネットワーク 2 6 1 5 およびインターネット 2 6 1 8 は両方とも、デジタルデータストリームを運ぶ電気的、電磁氣的、あるいは光学的信号を使用する。様々なネットワークを通る信号、およびコンピュータプラットフォーム 2 6 0 1 の間のデジタルデータを運ぶ、通信インターフェース 2 6 1 3 を通りネットワークリンク 2 6 1 4 上の信号は、情報を運ぶ搬送波の典型的な形態である。

50

【0166】

0191 コンピュータプラットフォーム2601は、インターネット2618、LAN2615、ネットワークリンク2614および通信インターフェース2613を含む多様なネットワーク、を通して、プログラムコードを含み、メッセージを送信しデータを受信することができる。インターネットの例において、システム2601がネットワークサーバとして働く場合、クライアント2620及び/または2621の上で走っているアプリケーションプログラムのために要求されるコードあるいはデータを、インターネット2618、ゲートウェイ/ファイアウォール2617、ローカルエリアネットワーク2615および通信インターフェース2613を通して送信することもありうる。同様に、他のネットワークリソースからコードを受け取ることもできる。

10

【0167】

0192 受け取られると、受信したコードはプロセッサ2605によって実行されるかも知れないし、及び/または、それぞれ持続型または揮発性ストレージ装置である2608および2606に各々格納され、あるいは後での実行のために他の不揮発性のストレージに格納されるかも知れない。

【0168】

0193 本発明が、どのような特定のファイアウォールシステムにも制限されていないことは注目すべきである。本発明の方針に基づく内容の処理システムは、3つのファイアウォールオペレーティングモード、のうちのどれでも使用できる、具体的には、NAT、ルーテッド、透過型である。

20

【0169】

0194 最後に、ここに記述されたプロセスと技術は、どのような特別の装置とも本質的には関係がなく、構成要素のいかなる適切な組合せによっても実行され得ることは理解すべきである。さらに、多様な型式の汎用目的の装置が、ここに記述された教えに従って使用することも出来る。あるいはまた、ここに記述された方法ステップを実行する専用の装置を構築することが有利であると判明するかも知れない。本発明は、特定の例に関して記述されているが、それは限定をするというよりむしろ全ての関連での例証となることを意図している。当分野での業者は、多くの違ったハードウェア、ソフトウェアおよびファームウェアの組合せが本発明を実施するのに適していることを認識するであろう。例えば、記述されるソフトウェアは、アセンブラー、C/C++、perl、シェル、PHP、Java（登録商標）などのような、種々様々なプログラミング言語あるいはスクリプト言語で実施されてもよい。

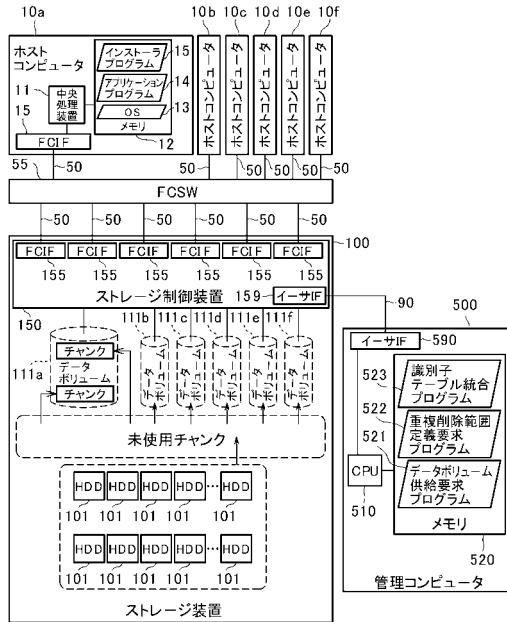
30

【0170】

0195 さらに、本発明のその他の実施も、ここに開示された本発明の明細書および実施の考察から、当分野の業者には明白になるであろう。記述された実施形態の種々の態様及び/または構成要素は、データ重複削除機能を備えたコンピュータ化されたストレージシステムにおいて、単独であるいは任意の組合せで使用することが出来る。明細書と実例は、典型的なものとしてのみ考慮されるよう意図されている。

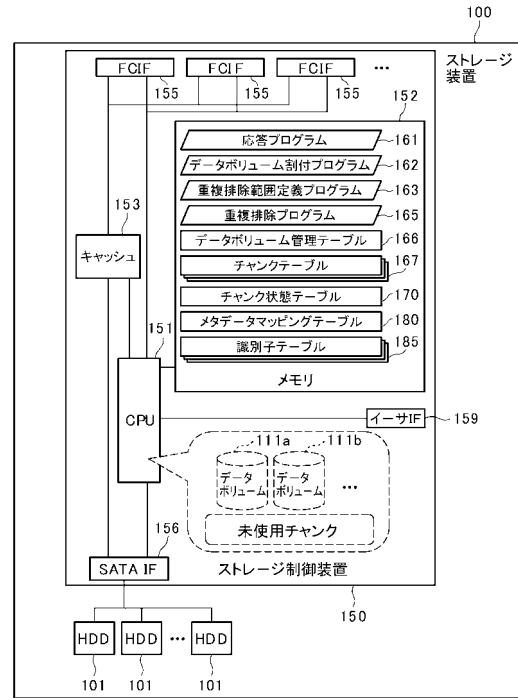
【図1a】

図1(a)



【図1b】

図1(b)



【図2】

図2

ボリューム番号	サイズ	メタデータ	チャンクテーブル番号
111a	2,000,000,000	OSタイプA	167a
111b	2,000,000,000	OSタイプA	167b
111c	2,000,000,000	OSタイプA	167c
111d	2,000,000,000	OSタイプB	167d
111e	2,000,000,000	OSタイプB	167e
111f	2,000,000,000	OSタイプB	167f
⋮			

データボリューム管理テーブル

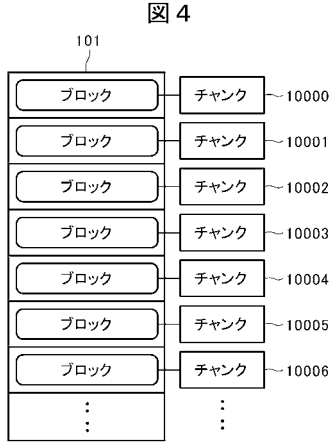
【図3】

図3

LBA	識別子	
	MD5	連続番号
0	0cc1cbcd165367614e3f49b1af433438	0
1	bdbc1ca29ba009a1d2ade3c39bc35acf	0
2		
3		
4		
⋮		

チャンクテーブル

【 図 4 】



【 図 5 】

図 5

チャンク番号	状態
10000	使用中
10001	使用中
10002	
10003	
10004	
10005	
...	

チャンク状態テーブル

【 図 6 】

図 6

(a)

メタデータ	識別子テーブル番号
OSタイプA	185a
OSタイプB	185b
...	

メタデータマッピングテーブル

(b)

メタデータ	識別子テーブル番号
OSタイプA、OSタイプB	185a
...	

メタデータマッピングテーブル

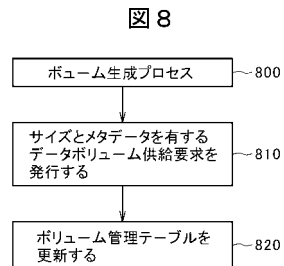
【 図 7 】

図 7

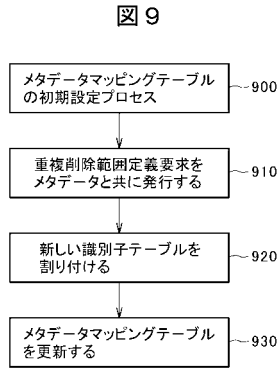
識別子	連続番号	参照数	チャンク番号	
			18501	18502
MDS	0	3	10000	10001
0ce1cbcd165367614c3f49b1af433438	0	3		
bdbc1ca29ba009a1d2adae3c39bc35acfc	0			
...				

識別子テーブル

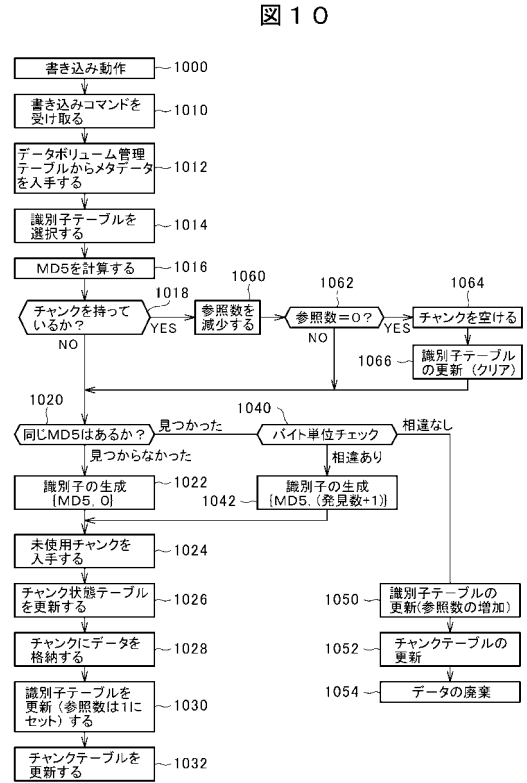
【 図 8 】



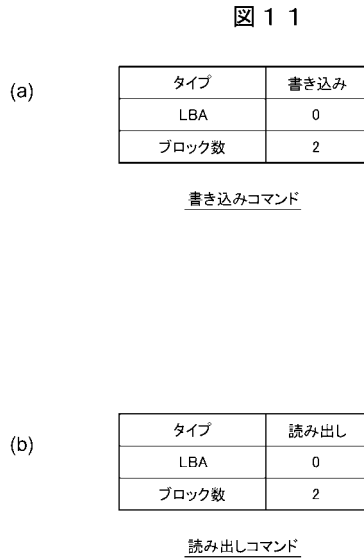
【 図 9 】



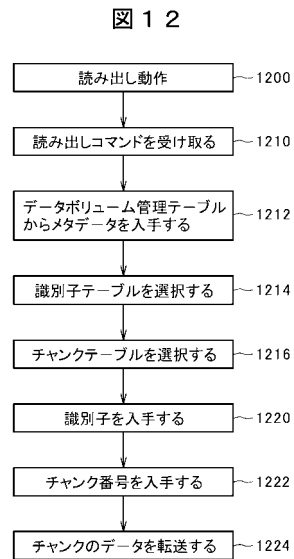
【 図 10 】



【 図 11 】

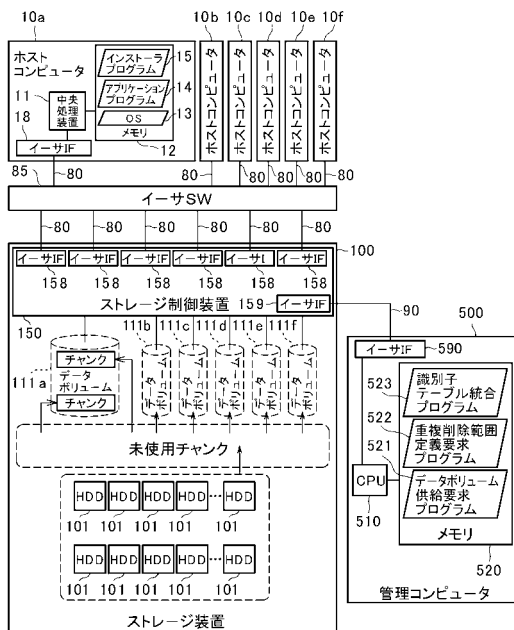


【 図 12 】



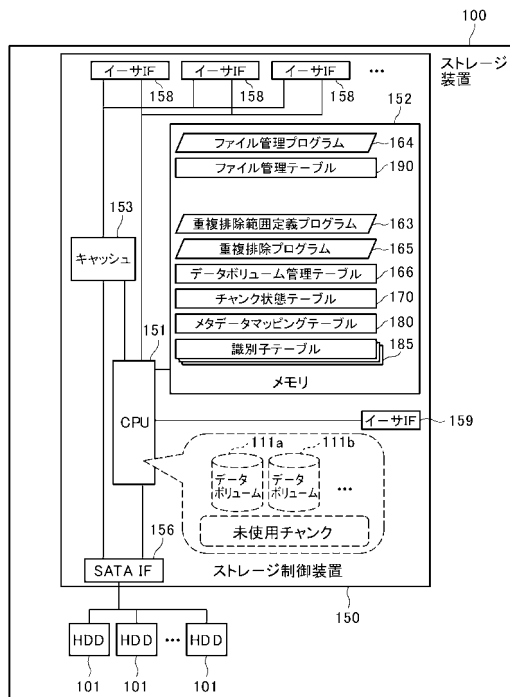
【図13a】

図13(a)



【図13b】

図13(b)



【図14】

図14

19001	19002	19003	19004
ファイル名	ファイルサイズ	識別子 MD5	メタデータ
/foo/sample.txt	2048	0cc1cbeecf165367614e349b1ef433438	Tom, アカウンティング, .txt, texteditor.Y, 10/31/2008
/foo/mail.txt	1024	bdac1ca29ba008a1d2a4e3c39bc39acf	Tom, mail, ma.IZ, 10/1/2008
/boo/mail.txt	1024	bdbc1ca29ba008a1d2a0c3c39bc39acf	Tom, アカウンティング, .txt, texteditor.Y, m.IZ, 10/1/2008
/boo/document.txt	4096	098f8bced4621d373cade4e82627b4f6	John, HR, .txt, texteditor.Y, 10/2/2008
		...	

ファイル管理テーブル

【図15a】

図15(a)

18001	18002
メタデータ	識別子テーブル 番号
.txt	185a
.doc	185b
.exe	185c
⋮	

メタデータマッピングテーブル

【図15b】

図15(b)

18001	18002
メタデータ	識別子テーブル 番号
アカウンティング	185a
HR	185b
R&D	185c
⋮	

メタデータマッピングテーブル

【図15c】

図15(c)

18001 メタデータ	18002 識別子テーブル番号
10/1/2008 ~ 10/31/2008	185a
11/1/2008 ~ 11/30/2008	185b
12/1/2008 ~ 12/31/2008	185c
⋮	

メタデータマッピングテーブル

【図15d】

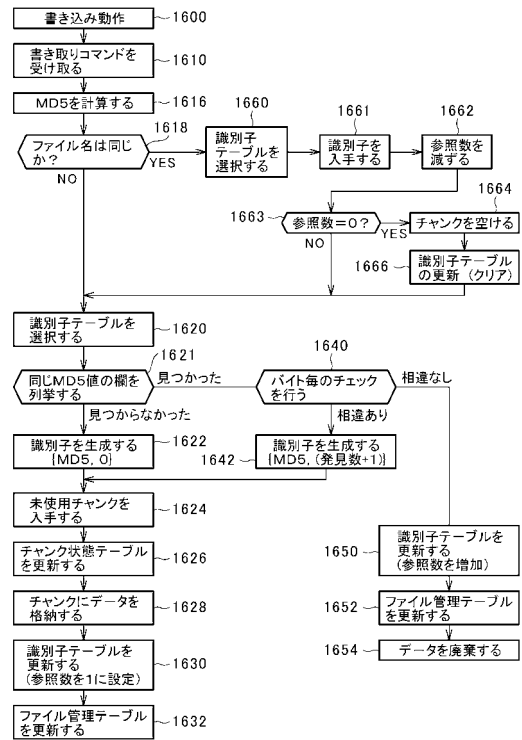
図15(d)

18001 メタデータ	18002 識別子テーブル番号
アカウントティング, HR	185a
R&D	185c
⋮	

メタデータマッピングテーブル

【図16】

図16



【図17】

図17

(a)

タイプ	書き込み
ファイル名	/foo/sample.txt
ファイルサイズ	2048 byte
メタデータ	Tom, アカウントティング, txt, texteditorX, 10/31/2008

書き込みコマンド

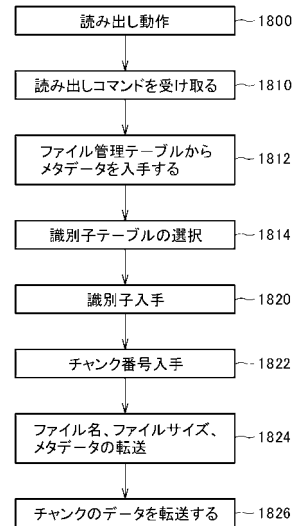
(b)

タイプ	読み出し
ファイル名	/foo/sample.txt

読み出しコマンド

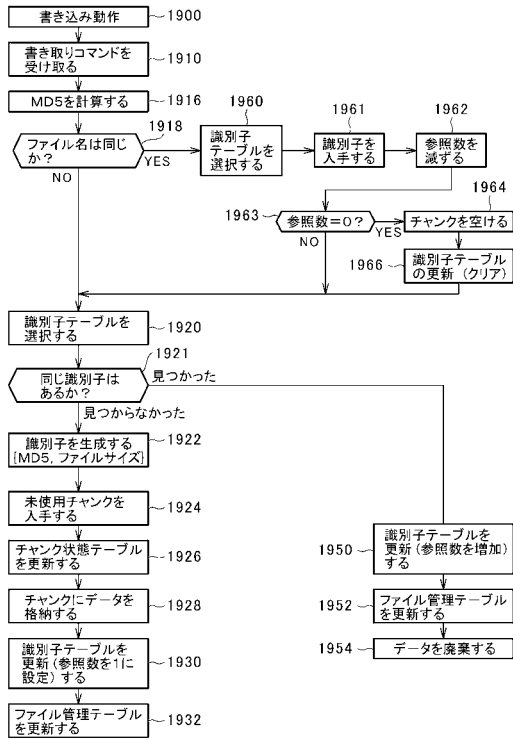
【図18】

図18



【図19】

図19



【図20】

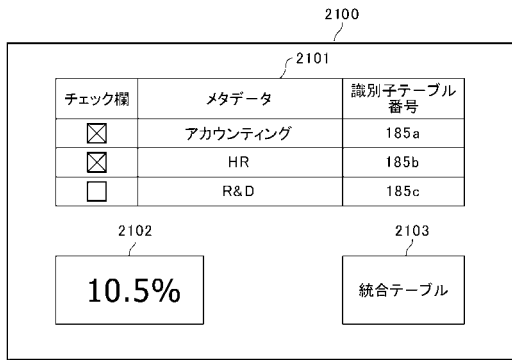
図20

ファイル名	ファイルサイズ	識別子	MDS	メタデータ
/foo/sample.txt	2048	0cc1bbcd165367614c3f49b1a433438	Tom, t.ct, 10/31/2008	アカウントティング
/foo/mail.txt	1024	bdbbc1ca29ba009a1d2ade3c39bc35acf	Tom, t.ct, 10/1/2008	アカウントティング
/foo/mail.txt	1024	bdbbc1ca29ba009a1d2ade3c39bc35acf	Tom, HR, t.ct, 10/1/2008	アカウントティング
/foo/document.txt	4096	098f6bcd4621d373cade4e832627b4f6	Tom, HR, t.ct, 10/27/2008	アカウントティング
				...

ファイル管理テーブル

【図21】

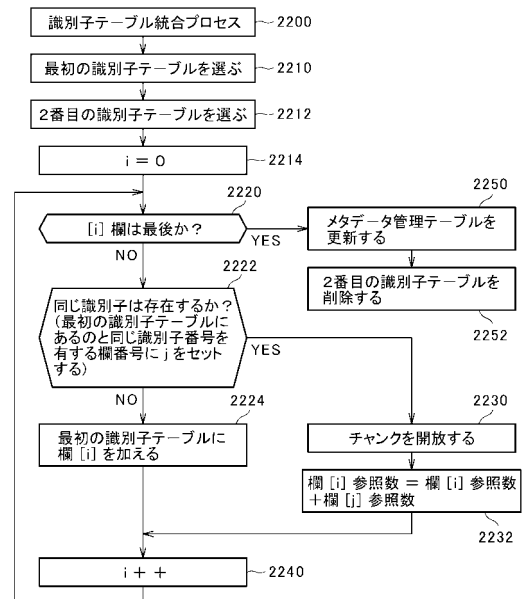
図21



スクリーン

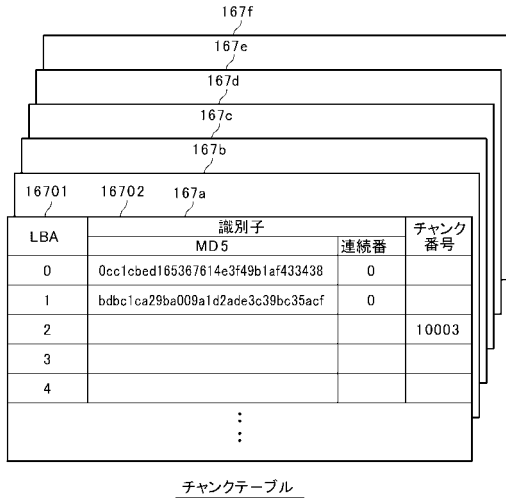
【図22】

図22



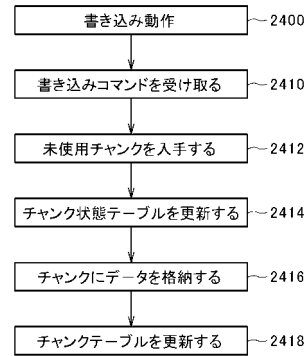
【図23】

図23



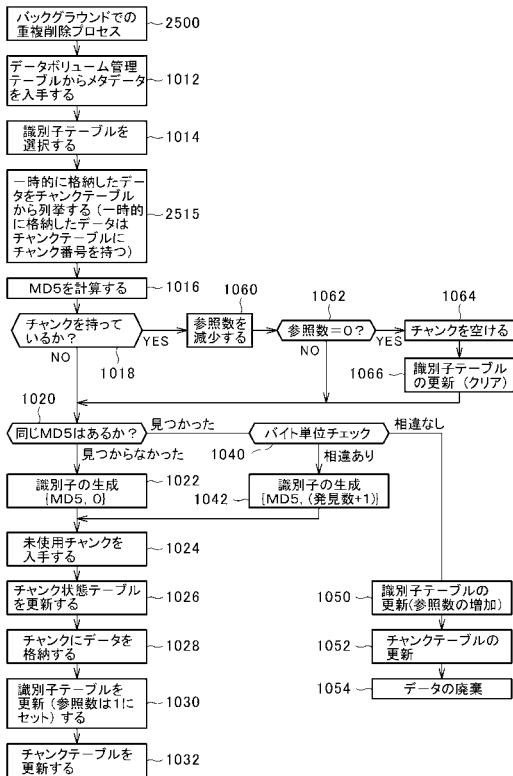
【図24】

図24



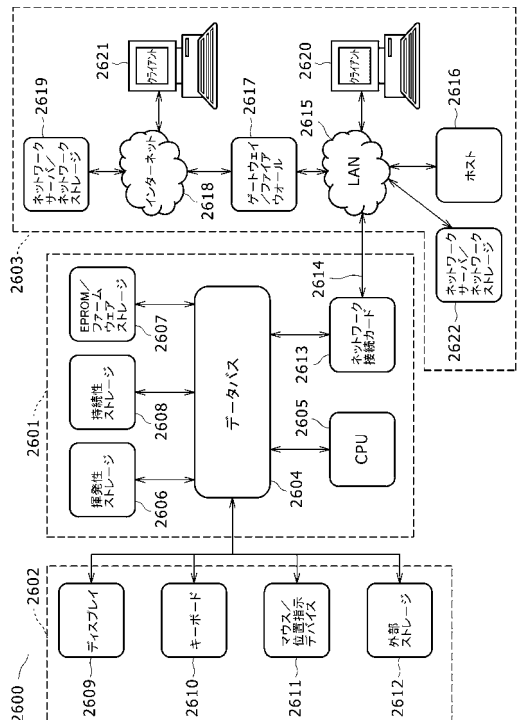
【図25】

図25



【図26】

図26



フロントページの続き

(56)参考文献 国際公開第2008/005211(WO, A1)
国際公開第2009/006278(WO, A1)

(58)調査した分野(Int.Cl., DB名)
G06F 3/06 - 3/08
G06F 12/00