

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2021-99611
(P2021-99611A)

(43) 公開日 令和3年7月1日(2021.7.1)

(51) Int. Cl. F I テーマコード (参考)
G06F 3/06 (2006.01)
 G06F 3/06 301W
 G06F 3/06 302A
 G06F 3/06 302J

審査請求 有 請求項の数 10 O L (全 29 頁)

(21) 出願番号	特願2019-230475 (P2019-230475)	(71) 出願人	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
(22) 出願日	令和1年12月20日 (2019.12.20)	(74) 代理人	110000279 特許業務法人ウィルフォート国際特許事務所
		(72) 発明者	水島 永雅 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
		(72) 発明者	吉原 朋宏 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
		(72) 発明者	島田 健太郎 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内

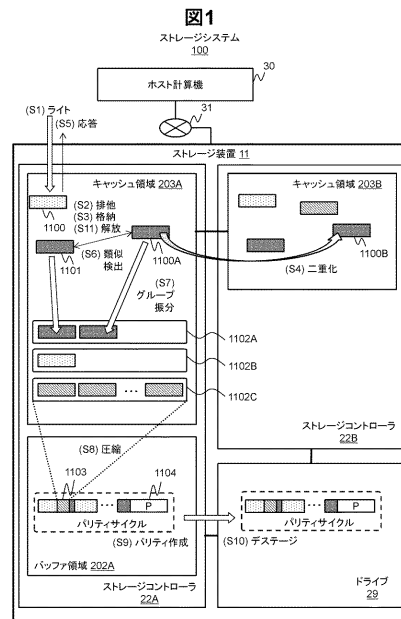
(54) 【発明の名称】 ストレージシステムおよびストレージシステムのデータ圧縮方法

(57) 【要約】

【課題】 ランダムに書き込まれたデータの圧縮率を高め、アクセス性能を向上させる。

【解決手段】 ストレージコントローラ22Aは、ドライブ29にリードまたはライトされるデータを格納するキャッシュ領域203Aを備え、キャッシュ領域203Aに格納されドライブ29に入力される複数のデータについて、そのデータ間の類似度に基づいてグループ化し、前記グループを選択し、選択したグループのデータをグループ単位で圧縮し、圧縮したデータをドライブ29に格納させる。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

物理記憶領域を有するドライブと、前記ドライブに入出力するデータを処理するコントローラを備え、

前記コントローラは、

前記ドライブにリードまたはライトされるデータを格納するキャッシュ領域を備え、

前記コントローラは、

前記キャッシュ領域に格納され前記ドライブに入力される複数のデータについて、そのデータ間の類似度に基づいてグループ化し、

前記グループを選択し、

選択したグループのデータをグループ単位で圧縮し、

前記圧縮したデータを前記ドライブに格納させるストレージシステム。

10

【請求項 2】

前記グループの選択は、前記グループ内のデータへのアクセス頻度および前記グループに含まれる複数のデータ間の類似度に基づいて行う請求項 1 に記載のストレージシステム。

【請求項 3】

前記コントローラは、前記グループのアクセス頻度が第 1 設定値以下かつ前記グループに含まれるデータ間の類似度が第 2 設定値以上の場合、前記キャッシュ領域に格納されたデータのグループを選択して圧縮する請求項 2 に記載のストレージシステム。

20

【請求項 4】

前記コントローラは、前記キャッシュ領域におけるダーティキャッシュ比率が所定値以上の場合、前記キャッシュ領域に格納されたデータのグループを選択して圧縮する請求項 3 に記載のストレージシステム。

【請求項 5】

前記グループ化するデータは、別個のライト要求に基づいて前記キャッシュ領域に格納された複数のデータを含む請求項 1 に記載のストレージシステム。

【請求項 6】

前記コントローラは、

前記ドライブに圧縮して格納された前記グループに含まれるデータ間の類似度を管理し

30

、前記ドライブに圧縮して格納されたデータを含むグループのうち、前記データ間の類似度が基準値以下のグループを選択し、

前記選択したグループに含まれるデータを前記ドライブから読み出して伸張し、前記伸張したデータを再グループ化及び再圧縮のために前記キャッシュ領域に格納する請求項 1 に記載のストレージシステム。

【請求項 7】

前記コントローラは、

圧縮したデータの伸長は、先頭側から完了し、

前記圧縮時には、同一グループに含まれるデータについて、前記アクセス頻度の高いデータを前記アクセス頻度の低いデータに比べて前記グループの先頭により近い位置に配置する請求項 2 に記載のストレージシステム。

40

【請求項 8】

前記コントローラは、前記データのハッシュ値に基づいて、前記類似度を評価し、

前記ハッシュ値の計算に用いる入力値は、前記データに含まれる同一長の複数の文字列であって、前記文字列は、前記データにおける出現頻度が所定数以上である請求項 1 に記載のストレージシステム。

【請求項 9】

前記コントローラは、

ライト要求にかかるデータについて、重複排除判定を行い、

50

重複排除要の場合には、重複排除処理を行い、
 重複排除不要の場合には、前記グループ化、圧縮及びドライブ格納を行う請求項 1 に記載のストレージシステム。

【請求項 10】

物理記憶領域を有するドライブと、前記ドライブに入出力するデータを処理するコントローラを備えるストレージシステムのデータ圧縮方法であって、

前記コントローラは、

前記ドライブにリードまたはライトされるデータを格納するキャッシュ領域を備え、

前記コントローラは、

前記キャッシュ領域に格納され前記ドライブに入力される複数のデータについて、そのデータ間の類似度に基づいてグループ化し、

10

前記グループを選択し、

選択したグループのデータをグループ単位で圧縮し、

前記圧縮したデータを前記ドライブに格納させるストレージシステムのデータ圧縮方法

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ストレージシステムおよびストレージシステムのデータ圧縮方法に関する。

20

【背景技術】

【0002】

ストレージシステムは、一般的に、1以上のストレージ装置を備える。1以上のストレージ装置の各々は、一般的に、記憶デバイスとして、例えば、HDD (Hard Disk Drive) 又はSSD (Solid State Drive) を備える。ストレージシステムは、SAN (Storage Area Network) 又はLAN (Local Area Network) などのネットワーク経由で、1又は複数の上位装置 (例えば、ホスト計算機30計算機) からアクセスされる。一般的に、ストレージ装置は、RAID (Redundant Array of Independent (or Inexpensive) Disks) 技術に従う高信頼化方法を用いることで信頼性を向上している。

30

【0003】

ストレージシステムには、物理的な記憶デバイスに保存するデータ量を削減してコストを低減するため、可逆圧縮機能を有するものがある。また、その圧縮率をできるだけ良くするために圧縮対象のデータに前処理を施すものもある。

【0004】

特許文献1には、その一つの技術が開示されている。すなわち、特許文献1に開示された技術では、ストレージシステム外部から書き込まれたファイルおよびシーケンスデータをチャンクという単位に分割し、チャンク間の類似性を評価し、類似するチャンク順に並び替えて新たなシーケンスデータを作り、それをグループ単位で圧縮して記憶デバイスに保存する。また、元のファイルやシーケンスデータを復元するためのチャンク並び替えリストも記憶デバイスに保存する。

40

【先行技術文献】

【特許文献】

【0005】

【特許文献1】米国特許第9367557号明細書

【発明の概要】

【発明が解決しようとする課題】

【0006】

50

特許文献 1 に記載のストレージシステムは、一度に書き込まれたファイルおよびシーケンスデータを圧縮して格納する時に、類似性に基づく並び替えを行ってから圧縮することによって圧縮率を改善する。ストレージシステムへの書き込みデータは、そのようなシーケンシャルな性質のものだけでなく、ランダムで互いに無関係なものもある。

【 0 0 0 7 】

後者のような書き込みデータを扱う場合においても、類似性に基づいて圧縮率を改善するストレージシステムを実現することが課題である。また、それを実現する際に、ストレージシステムのアクセス性能を改善することも課題である。

【 0 0 0 8 】

本発明は、上記事情に鑑みなされたものであり、その目的は、ランダムに書き込まれたデータの圧縮率を高め、アクセス性能を向上させることが可能なストレージシステムおよびストレージシステムのデータ圧縮方法を提供することにある。

【課題を解決するための手段】

【 0 0 0 9 】

上記目的を達成するため、第 1 の観点に係るストレージシステムは、物理記憶領域を有するドライブと、前記ドライブに入出力するデータを処理するコントローラを備え、前記コントローラは、前記ドライブにリードまたはライトされるデータを格納するキャッシュ領域を備え、前記コントローラは、前記キャッシュ領域に格納され前記ドライブに入力される複数のデータについて、そのデータ間の類似度に基づいてグループ化し、前記グループを選択し、選択したグループのデータをグループ単位で圧縮し、前記圧縮したデータを前記ドライブに格納させる。

【発明の効果】

【 0 0 1 0 】

本発明によれば、ランダムに書き込まれたデータの圧縮率を高め、アクセス性能を向上させることができる。

【図面の簡単な説明】

【 0 0 1 1 】

【図 1】図 1 は、実施形態に係るストレージシステムのデータ圧縮を伴うデータライト手順を示すブロック図である。

【図 2】図 2 は、実施形態に係るストレージシステムのハードウェア構成例を示すブロック図である。

【図 3】図 3 は、図 2 のボリューム管理テーブルの構成例を示す図である。

【図 4】図 4 は、図 2 のプール構成管理テーブルの構成例を示す図である。

【図 5】図 5 は、図 2 の RAID 構成管理テーブルの構成例を示す図である。

【図 6】図 6 は、図 2 のプール割当管理テーブルの構成例を示す図である。

【図 7】図 7 は、図 2 のドライブ割当管理テーブルの構成例を示す図である。

【図 8】図 8 は、図 2 のストレージ装置によって管理される論理記憶階層の構成例を示す図である。

【図 9】図 9 は、図 2 のグループ管理テーブルの構成例を示す図である。

【図 10】図 10 は、図 2 のメモリ割当管理テーブルの構成例を示す図である。

【図 11】図 11 は、図 2 の LRU リストテーブルとその遷移を示す図である。

【図 12】図 12 は、図 2 のストレージ装置におけるキャッシュメモリのリソース配分比率の変化を示す図である。

【図 13】図 13 は、図 2 のストレージ装置におけるスロットとグループの構成例を示す図である。

【図 14】図 14 は、図 2 のストレージ装置が実行する類似ハッシュ算出処理を示す図である。

【図 15】図 15 は、図 2 のストレージ装置が実行するデータリード処理を示すフローチャートである。

【図 16】図 16 は、図 2 のストレージ装置が実行するデータライト処理を示すフローチャートである。

10

20

30

40

50

ャートである。

【図17】図17は、図2のストレージ装置が実行する圧縮前データ構築処理を示すフローチャートである。

【図18】図18は、図2のストレージ装置が実行する圧縮・デステージ処理を示すフローチャートである。

【図19】図19は、図2のストレージ装置が実行するデータ削除リトライ処理を示すフローチャートである。

【図20】図20は、図2のストレージ装置におけるデータ伸張処理の過程を示す図である。

【発明を実施するための形態】

10

【0012】

実施形態について、図面を参照して説明する。なお、以下に説明する実施形態は特許請求の範囲に係る発明を限定するものではなく、また、実施形態の中で説明されている諸要素およびその組み合わせの全てが発明の解決手段に必須であるとは限らない。

【0013】

以下の説明では、「インターフェース部」は、ユーザインターフェース部と、通信インターフェース部とのうちの少なくとも1つを含んでよい。ユーザインターフェース部は、1以上のI/Oデバイス（例えば入力デバイス（例えばキーボード及びポインティングデバイス）と出力デバイス（例えば表示デバイス））と表示用計算機とのうちの少なくとも1つのI/Oデバイスを含んでよい。通信インターフェース部は、1以上の通信インターフェースデバイスを含んでよい。1以上の通信インターフェースデバイスは、1以上の同種の通信インターフェースデバイス（例えば1以上のNIC（Network Interface Card））であってもよいし、2以上の異種の通信インターフェースデバイス（例えばNICとHBA（Host Bus Adapter））であってもよい。

20

【0014】

また、以下の説明では、「メモリ部」は、1以上のメモリを含む。少なくとも1つのメモリは、揮発性メモリであってもよいし、不揮発性メモリであってもよい。メモリ部は、主に、プロセッサ部による処理の際に使用される。

【0015】

また、以下の説明では、「プロセッサ部」は、1以上のプロセッサを含む。少なくとも1つのプロセッサは、典型的には、CPU（Central Processing Unit）である。プロセッサは、処理の一部又は全部を行うハードウェア回路（例えば、CPUの処理を補助するオフロードエンジン）を含んでもよい。

30

【0016】

また、以下の説明では、「xxxテーブル」の表現にて情報を説明することがあるが、情報は、どのようなデータ構造で表現されていてもよい。すなわち、情報がデータ構造に依存しないことを示すために、「xxxテーブル」を「xxx情報」と言うことができる。また、以下の説明において、各テーブルの構成は一例であり、1つのテーブルは、2以上のテーブルに分割されてもよいし、2以上のテーブルの全部又は一部が1つのテーブルであってもよい。

40

【0017】

また、以下の説明では、同種の要素を区別しないで説明する場合には、参照符号のうちの共通符号を使用し、同種の要素を区別する場合は、参照符号（又は要素のID（例えば識別番号））を使用することがある。例えば、複数のストレージコントローラを区別しない場合には、「ストレージコントローラ22」と記載し、各ストレージコントローラを区別する場合には、「ストレージコントローラ22A」、「ストレージコントローラ22B」のように記載する。他の要素（例えばキャッシュ領域203、バッファ領域202、アドレス1100等）も同様である。

【0018】

また、以下の説明では、「ストレージシステム」は、1以上のストレージ装置を含む。

50

少なくとも1つのストレージ装置は、汎用的な物理計算機であってもよい。また、少なくとも1つのストレージ装置が、仮想的なストレージ装置であってもよいし、SDx (Software-Defined anything) を実行してもよい。SDxとしては、例えば、SDS (Software Defined Storage) (仮想的なストレージ装置の一例) 又はSDDC (Software-defined Data center) を採用することができる。

【0019】

図1は、実施形態に係るストレージシステムのデータ圧縮を伴うデータライト手順を示すブロック図である。

図1において、ストレージシステム100は、ホスト計算機30及びストレージ装置11を備える。ホスト計算機30は、ネットワーク31を介してストレージ装置11に接続され、管理計算機(図示せず)によって管理される。

10

【0020】

ストレージ装置11は、1以上のボリューム(仮想的な記憶領域)を有する。ストレージ装置11は、物理記憶領域を有するドライブ29と、ドライブ29を制御するストレージコントローラ22A、22Bを備える。ストレージコントローラ22Aは、ドライブにおいてリードまたはライトされるデータをキャッシュするキャッシュ領域203Aと、転送時のデータを一時的に保存するバッファ領域202を備える。ストレージコントローラ22Bは、キャッシュ領域203Aにキャッシュされたデータを二重化するキャッシュ領域203Bを備える。

20

【0021】

ホスト計算機30は、物理的な計算機でもよいし、物理的な計算機で実行される仮想的な計算機でもよい。ホスト計算機30は、ストレージシステム100において実行される仮想的な計算機でもよい。ホスト計算機30からは、ストレージ装置11のストレージコントローラ22A又はストレージコントローラ22Bに対してデータの書き込みが行われる。

【0022】

以下、このストレージシステム100において、ホスト計算機30からのデータライト命令に応じた処理手順について説明する。図1の例では、ホスト計算機30からのライト命令をストレージコントローラ22Aが受領した場合について示す。

30

【0023】

(S1)ストレージ装置11は、ホスト計算機30からネットワーク31を介してライト命令を受信する。ライト命令は、データとデータの割当先アドレス1100とを含んでいる。ストレージ装置11がライト命令を処理する際のデータ単位は、例えば8KBである。8KBに満たない場合は、周辺アドレスの不足データを補うため、当該装置内の記憶済みデータを読み出して補完し、8KBにしてから処理する。周辺アドレスに未書き込み部分がある場合は、その部分を不定値として管理し、暫定的にゼロデータで補完して処理する。8KBより大きい場合は、1つ以上の8KBデータと、8KBに満たない残りデータとに分け、上記補完により全て8KBデータにして個別に処理する。従って、以下の説明では、ストレージ装置11が扱うデータサイズは8KBであるものとする。ストレージ装置11は、ライト命令受信後に、S2以降のライト処理を開始する。

40

【0024】

(S2)ストレージ装置11は、ライト命令に応答して、割当先アドレス1100が示すスロットの排他を確保する。これによって、そのスロット内のデータが他のライト命令によって更新されることを防ぐ。「スロット」とは、ボリューム(VOL)における領域単位の種類である。具体的には、本実施形態のスロットは、ドライブ29への書き込みが行われたか否か及びバッファ領域202への転送が行われたか否か等の管理の単位となる。スロットのサイズは、例えば256KBである。本実施形態では、この領域を「スロット」と呼ぶが、他の名称で呼ばれてもよい。

【0025】

50

(S3) ストレージコントローラ22Aは、キャッシュ領域203Aにおいて、データの割当先アドレス1100に対応するアドレス1100Aにデータを格納する。

【0026】

(S4) ストレージコントローラ22Aは、キャッシュ領域203A内に格納されたデータをストレージコントローラ22Bに転送する。ストレージコントローラ22Bは、割当先アドレス1100に対応するキャッシュ領域203B内のアドレス1100Bに受領したデータを格納し、ストレージコントローラ22Aへ応答を返すことでストレージ装置11内でのデータの二重化を完了する。

【0027】

(S5) データの二重化を完了した後、ストレージ装置11は、ホスト計算機30に対してネットワーク31を介してライト完了を応答する。この時点で、ホスト計算機30は、ライトが完了したと認識する。

10

【0028】

(S6) ストレージコントローラ22Aは、キャッシュ領域203A内に格納されている複数のデータの中から、例えば、アドレス1100Aのデータの内容と類似関係にある1つ以上のデータを検出する。アドレス1100Aのデータの内容と類似関係にあるデータのアドレスを1101とする。

【0029】

類似関係にあるデータを検出するために、8KBのデータの内容から類似ハッシュと呼ばれる値(データ内容の特徴を表す小さいサイズ(例えば8B)の値)を算出する。例えば、類似ハッシュは、LSH(Locality Sensitive Hash)である。2つのデータの内容が類似しているほど、互いの類似ハッシュの値は近くなる。類似ハッシュの値が近いことは、ハミング距離が短いことを意味する。具体的な類似ハッシュ算出方法は後述する。

20

【0030】

(S7) ストレージコントローラ22Aは、キャッシュ領域203Aにキャッシュされたデータ間の類似度に基づいて、キャッシュ領域203Aにキャッシュされたデータをクレーピングしたグループを生成する。各グループにクレーピングされるデータは、ホスト計算機30から別個のライト要求に基づいてライトされたデータを含んでいてもよい。このとき、ストレージコントローラ22Aは、グループごとに複数のデータを格納するため、1つ以上の領域をキャッシュ領域203A内に確保する。グループごとに確保された領域のサイズは、例えば128KBである。そして、ストレージコントローラ22Aは、例えば、ある1つのグループ1102Aの領域に、アドレス1100Aのデータとアドレス1101のデータを転送する。転送するデータの数は2個以上であり、例えば1つのグループには、類似関係にある8KBデータを16個収納できる。別のグループ1102B、1102Cの領域には、(S6)の方法で検出した類似関係にあるデータをそれぞれ収納する。

30

【0031】

また、後述するように、ストレージコントローラ22Aは、選択したデータをグループリングする前に、重複排除処理を行って、その処理後の代表データのみをグループリングしてもよい。重複排除処理は、内容が完全に等しいN(Nは、2以上の整数)個のデータを1個の代表データに割り当てる処理である。この重複排除処理により、実際の格納データ量をN分の1に削減することができる。

40

【0032】

(S8) ストレージコントローラ22Aは、例えば、キャッシュ領域203Aからドライブ29へ書き出すグループ(例えば、グループ1102C)を選択し、選択したグループに含まれるデータをまとめて圧縮してバッファ領域202A内のアドレス1103に追記的に格納する。

【0033】

ここで、ストレージコントローラ22Aは、グループのアクセス頻度およびグループに

50

含まれるデータ間の類似度に基づいて、圧縮対象となるグループを選択する。例えば、ストレージコントローラ 22 A は、グループのアクセス頻度が第 1 設定値以下かつグループに含まれるデータ間の類似度が第 2 設定値以上のグループを圧縮対象として選択する。これにより、アクセス頻度が低く、データ間の類似度が高いグループを圧縮対象として選択し、キャッシュ領域 203 A から落とすことができる。このため、ドライブ容量を削減することが可能となるとともに、アクセス頻度が高いデータをキャッシュ領域 203 A に残すことができ、キャッシュヒット率を向上させることが可能となることから、アクセス性能を向上させることができる。

【0034】

さらに、ストレージコントローラ 22 A は、キャッシュ領域 203 A におけるダーティキャッシュ比率が上限値以上の場合、グループの圧縮を実行し、ダーティキャッシュ比率が上限値未満の場合、グループの圧縮を保留するようにしてもよい。これにより、キャッシュ領域 203 に新たにライトデータが受け入れられる Free 状態の領域を確保しつつ、より圧縮度を高めることが可能な類似したデータを可能な限りキャッシュ領域 203 に保持させることができ、データの圧縮率を高め、アクセス性能を向上させることが可能となる。

【0035】

なお、この圧縮処理は、バッファ領域 202 A 内に RAID パリティサイクル分以上の量が溜まるまでグループごとに実施される。このとき、ストレージコントローラ 22 A は、選択したグループの圧縮後に暗号化処理を行って、その暗号化処理後の圧縮データをアドレス 1103 に格納してもよい。

【0036】

(S9) ストレージコントローラ 22 A は、バッファ領域 202 A 内に追記的に格納した 1 つ以上のグループの圧縮データがパリティサイクル分以上の量に達すると、それらの内容を保証するためのパリティを生成し、バッファ領域 202 A 内のアドレス 1104 へ格納する。

【0037】

(S10) ストレージコントローラ 22 A は、バッファ領域 202 A 内の圧縮グループ (1 つ以上のグループのパリティサイクル分以上の圧縮データ) 及びそれに対応するパリティをドライブ 29 へ送信し、ドライブ 29 に書き出す (デステージ処理)。

【0038】

(S11) ストレージコントローラ 22 A は、デステージ処理が完了すると、(S2) において確保したスロットの排他を解放する。

【0039】

図 2 は、実施形態に係るストレージシステムのハードウェア構成例を示すブロック図である。

図 2 において、ストレージ装置 11 は、1 以上のストレージコントローラ 22 と、1 以上のストレージコントローラ 22 に接続された複数のドライブ 29 とを有する。

【0040】

ストレージコントローラ 22 は、FE_I/F (フロントエンドインターフェースデバイス) 23、プロセッサ 24、メモリ 25、215、BE_I/F (バックエンドインターフェースデバイス) 27、内部ネットワーク 26 およびオフロードエンジン 214 を備える。FE_I/F 23、プロセッサ 24、メモリ 25、BE_I/F 27 およびオフロードエンジン 214 は、内部ネットワーク 26 を介して接続されている。

【0041】

FE_I/F 23 は、ホスト計算機 30 との通信を行う。BE_I/F 27 は、ドライブ 29 との通信を行う。プロセッサ 24 は、ストレージ装置 11 全体を制御する。

【0042】

メモリ 25 は、プロセッサ 24 で使用されるプログラム及びデータを格納する。メモリ 25 は、プログラムを管理するプログラム領域 201、データの転送時の一時的な保存領

10

20

30

40

50

域であるバッファ領域 202、ホスト計算機 30からのライトデータ（ライト命令にตอบสนองして書き込まれるデータ）及びドライブ 29からのリードデータ（リード命令にตอบสนองして読み出されたデータ）を一時的に格納するキャッシュ領域 203及び種々のテーブルを格納するテーブル管理領域 206を有する。

【0043】

テーブル管理領域 206は、ボリュームに関する情報を保持するボリューム情報管理テーブル 207、プールに関する情報を保持するプール構成管理テーブル 208、RAID構成に関する情報を保持するRAID構成管理テーブル 209、プール割り当てに関する情報を保持するプール割当管理テーブル 210、ドライブ割り当てに関する情報を保持するドライブ割当管理テーブル 211、グループへのデータ割り当てに関する情報を保持するグループ管理テーブル 212、メモリ割り当てに関する情報を保持するメモリ割当管理テーブル 213及びスロットのアクセス頻度を管理するLRU（Least Recently Used）リストテーブル 217を格納する。

10

【0044】

ドライブ 29は、不揮発性のデータ記憶媒体を有する装置であり、例えばSSD（Solid State Drive）でもよいし、HDD（Hard Disk Drive）でもよい。複数のドライブ 29が、複数のRAIDグループ（パリティグループとも呼ばれる）を構成してよい。各RAIDグループは、1以上のドライブ 29から構成される。

【0045】

オフロードエンジン 214は、プロセッサ 24の行う処理を補助するハードウェア回路であり、図1を用いて説明したデータライトにおいて実施する、類似検出、圧縮、重複排除、暗号化およびパリティ生成などの処理の一部または全部をプロセッサ 24よりも高速に実行する。同様に、データリード（図による説明は省略）において実施する、伸張、復号およびパリティによる回復（ドライブ障害時のみ）などの処理の一部または全部をプロセッサ 24よりも高速に実行する。

20

【0046】

メモリ 215は、オフロードエンジン 214に直接接続された専用メモリである。メモリ 215は、データの転送時の一時的な保存領域であるバッファ領域 216を有する。メモリ 215は、オフロードエンジン 214で扱われるデータや、それらの処理に関する情報も保存する。メモリ 215を設けることで、内部ネットワーク 26とプロセッサ 24を経由するメモリ 25とオフロードエンジン 214間の転送量及び転送時間を削減し、性能向上を図ることができる。

30

【0047】

FE_I/F 23およびBE_I/F 27は、インターフェース部の一例である。メモリ 25、215は、メモリ部の一例である。プロセッサ 24とオフロードエンジン 214は、プロセッサ部の一例である。

【0048】

図3は、図2のボリューム管理テーブルの構成例を示す図である。

図3において、ボリューム管理テーブル 207は、ボリューム毎にエントリを有する。各エントリは、VOL_ID 41、VOL属性 42、VOL容量 43及びプールID 44の情報を格納する。

40

【0049】

VOL_ID 41は、ストレージ装置 11が管理している各ボリュームのIDである。VOL属性 42は、各ボリュームの属性を示す。各ボリュームの属性は、例えば、シンプロビジョニングのボリュームか、通常割り当てのボリュームか、データ削減が有効（ON）か無効（OFF）か、システムボリュームかなどである。システムボリュームは、類似関係にあるデータを含むグループをストレージ装置 11が管理するためのボリュームであり、ホスト計算機 30からは直接見えない領域である。ボリューム容量 43は、各ボリュームの容量を示す。プールID 44は、各ボリュームに関連付けられているプールのID

50

である。

【 0 0 5 0 】

プロセッサ 2 4 は、デステージ処理において、ボリューム管理テーブル 2 0 7 のボリューム属性 4 2 を参照することで、ホスト計算機 3 0 からデータのライト/リード命令を受けた時に、どのような処理を必要とするボリュームか判定できる。例えば、VOL_ID 4 1 が “ 1 0 ” のボリュームは、ボリューム属性 4 2 が “ シンプロビジョニング、削減 ON ” であるため、必要に応じてスロットに動的にプールを割り当てながら、ホスト計算機 3 0 のライト命令に対して図 1 で説明したデータライトを行う。

【 0 0 5 1 】

VOL_ID 4 1 が “ 0 ” のボリュームは、ボリューム属性 4 2 が “ シンプロビジョニング、削減 OFF ” であるため、必要に応じてスロットに動的にプールを割り当てながら、ホスト計算機 3 0 のライト命令に対してデータ削減処理（重複排除及び圧縮）を除いたデータライトを行う。

10

【 0 0 5 2 】

VOL_ID 4 1 が “ 3 0 ” のボリュームは、ボリューム属性 4 2 が “ 通常割当、削減 OFF ” であるため、ボリューム作成時に割り当てたプール容量範囲で、ホスト計算機 3 0 のライト命令に対してデータ削減処理（重複排除及び圧縮）を除いたデータライトを行う。

【 0 0 5 3 】

図 4 は、図 2 のプール構成管理テーブルの構成例を示す図である。

20

図 4 において、プールは、1 以上の RAID グループを基に構成された論理記憶領域である。プール構成管理テーブル 2 0 8 は、プール毎にエントリを有する。各エントリは、プール ID 5 1、RAID グループ ID 5 2、プール容量 5 3 及びプール使用容量 5 4 の情報を格納する。

【 0 0 5 4 】

プール ID 5 1 は、ストレージ装置 1 1 が管理している各プールの ID である。RAID グループ ID 5 2 は、プールの基になっている 1 以上の RAID グループの各々の ID である。プール容量 5 3 は、プールの容量を示す。プール使用容量 5 4 は、プールのプール容量のうち、ボリュームに割り当てられている領域の総量を示す。

【 0 0 5 5 】

図 5 は、図 2 の RAID 構成管理テーブルの構成例を示す図である。

30

図 5 において、RAID 構成管理テーブル 2 0 9 は、RAID 構成を管理する対象の RAID グループ毎にエントリを有する。各エントリは、RAID グループ ID 6 1、RAID レベル 6 2、ドライブ ID 6 3、ドライブ種別 6 4、容量 6 5 及び使用容量 6 6 の情報を格納する。

【 0 0 5 6 】

RAID グループ ID 6 1 は、ストレージ装置 1 1 が管理している各 RAID グループの ID である。RAID レベル 6 2 は、RAID グループに適用される RAID アルゴリズムの種別を示す。ドライブ ID 6 3 は、RAID グループを構成する 1 以上のドライブの各々の ID である。ドライブ種別 6 4 は、RAID グループを構成するドライブの種別（例えば HDD か SSD か）を示す。容量 6 5 は、RAID グループの容量を示す。使用容量 6 6 は、RAID グループの容量のうちの使用されている容量を示す。

40

【 0 0 5 7 】

図 6 は、図 2 のプール割当管理テーブルの構成例を示す図である。

図 6 において、プール割当管理テーブル 2 1 0 は、プールの割り当てを管理する対象のスロットの VOL アドレス（ボリューム内のスロットを示すアドレス）毎にエントリを有する。各エントリは、VOL_ID 7 1、VOL アドレス 7 2、プール ID 7 3、プールアドレス 7 4、圧縮前サイズ 7 5、圧縮後サイズ 7 6 及び類似度 7 7 の情報を格納する。

【 0 0 5 8 】

VOL_ID 7 1 は、ストレージ装置 1 1 が管理し、VOL アドレスによって識別され

50

るスロットが属するボリュームのIDである。VOLアドレス72は、それらスロットのVOLアドレスである。プールID73は、各スロット(256KB)に含まれる2つのグループ(各128KB)のデータをそれぞれ圧縮したデータを格納するために割り当てられたデータ領域を含むプールの各IDである。プールアドレス74は、2つのグループのデータをそれぞれ圧縮したデータを格納するために割り当てられたデータ領域のアドレス(プールに属するアドレス)である。圧縮前サイズ75は、各グループの圧縮前のデータサイズ(各128KB)を示す。圧縮後サイズ76は、各グループの圧縮後のデータサイズを示す。類似度77は、各グループを構成する16個の8KBのデータの相互の類似関係の強さを示す値である。類似度77が0%は、内容に全く類似関係がないことを示し、類似度77が100%は、内容が完全に一致している(可能性が高い)ことを示す。一般に、類似度77が高いグループほど圧縮率が高い(圧縮後サイズ76が小さい)ことが期待できる。

10

【0059】

ストレージ装置11で用いられる可逆圧縮アルゴリズムは、スライド辞書(Sliding Dictionary)圧縮をベースとしており、過去の文字列から一致する文字列を見つけて短い符号(発見位置までの距離と一致長)に置き換えてデータ量を削減する。例えば、30バイト一致する文字列を見つけて、2バイトの短い符号に置き換えられれば28バイト削減される。類似度が高いデータには共通する文字列が多く含まれており、圧縮するグループ内のデータの類似度が高いほど、この辞書圧縮が効果的に働くため、圧縮率が高くなる。

20

【0060】

図7は、図2のドライブ割当管理テーブルの構成例を示す図である。

図7において、ドライブ割当管理テーブル211は、ドライブ割当を管理するプールアドレス毎にエントリを有する。各エントリは、プールID81、プールアドレス82、RAIDグループID83、ドライブID84及びドライブアドレス85の情報を格納する。

【0061】

プールID81は、プールアドレスが属するプールのIDである。プールアドレス82は、ストレージ装置11が管理しているプールアドレスである。RAIDグループID83は、プールアドレスが示すデータ領域の基になっているRAIDグループのIDである。ドライブID84は、プールアドレスが示すデータ領域の基になっているドライブのIDである。ドライブアドレス85は、プールアドレスに対応したドライブアドレスである。

30

【0062】

図8は、図2のストレージ装置によって管理される論理記憶階層の構成例を示す図である。

図8において、ホストVOL1000は、シンプロビジョニングが適用され、データ削減が有効なボリュームであり、ホスト計算機30に提供される。システムVOL1001は、データ削減時の圧縮単位であるグループを管理するためのボリュームであり、ホスト計算機30からは直接見えない。

40

【0063】

ホストVOL1000内の各8KBデータアドレスに、システムVOL1001内のグループに含まれる各8KBデータアドレスが対応付けられる。図8の例では、VOL1000内の内容「A, B, D」を格納する3つのアドレスに、グループ1102D内の3つのアドレスが対応付けられ、VOL1000内の内容「C, E, F」を格納する3つのアドレスに、グループ1102E内の3つのアドレスが対応付けられている。内容「A, B, D」と「C, E, F」を格納するアドレスの対応先がそれぞれ同じグループ1102D、1102Eに含まれているのは、図1の類似検出(S6)により、それぞれの内容に類似関係があると判断されたためである。

【0064】

50

なお、重複排除処理が実施された場合は、ホストVOL1000内の複数の8KBデータアドレスに、システムVOL1001の1つの8KBデータアドレスが対応付けられることがある。図示していないが、異なるホストVOLの複数の8KBデータアドレスに、システムVOLの1つの8KBデータアドレスが対応付けられることもある。図8の例では、同一内容「A」を格納する異なる2つの8KBデータアドレス1105及び1106にグループ1102D内の1つの8KBデータアドレス1107が対応付けられ、同一内容「C」を格納する異なる2つの8KBデータアドレス1108及び1109にグループ1102E内の1つの8KBデータアドレス1110が対応付けられている。

【0065】

システムVOL1001内のスロット(256KB)に含まれる2つグループ(各128KB)には、圧縮された当該グループデータの格納先として、プール1002内のデータ領域がそれぞれ対応付けられる。図8の例では、グループ1102Dにデータ領域1103Dが対応付けられ、グループ1102Eにデータ領域1103Eが対応付けられている。なお、グループは圧縮されるため、各グループに対応するプール1002内のデータ領域のサイズは、グループサイズ(128KB)以下である。

10

【0066】

プール1002の空間は、チャンクと呼ばれる単位で分けられ、各チャンクはドライブアドレス空間1003へ対応付けられる。ドライブアドレス空間1003は、RAIDグループ1004を構成する複数のドライブ29(例えば4台)によって提供される物理的なデータ格納空間である。パリティ生成処理で生成するパリティ1104のサイズは、チャンクサイズである。ドライブアドレス空間1003は、複数のRAIDサイクルで構成される。ドライブ29の故障によるデータ消失をRAID技術で回復するため、各RAIDサイクルのパリティチャンクPは、1つのドライブに収まるように対応付けられる。

20

【0067】

ホストVOL1000からシステムVOL1001への割り当ては、図9のグループ管理テーブル212を基に管理される。そのテーブルの詳細は後述する。また、システムVOL1001からプール1002への割り当ては、図6のプール割当管理テーブル210を基に管理される。また、プール1002からドライブアドレス空間1003への割り当ては、図7のドライブ割当管理テーブル211を基に管理される。

【0068】

30

図9は、図2のグループ管理テーブルの構成例を示す図である。

図9において、グループ管理テーブル212は、図8のホストVOL1000の8KBデータに対応するシステムVOL1001のグループを管理するテーブルであり、管理対象の8KBデータのアドレス毎にエントリを有する。各エントリは、ホストVOL_ID901、ホストVOLアドレス902、位置番号903、システムVOL_ID904、システムVOLアドレス905、グループ番号906及び位置番号907の情報を格納する。

【0069】

ホストVOL_ID901は、ストレージ装置11が管理し、グループ対応の管理対象の8KBデータが属するスロットが属するボリューム(つまり、ホストVOL1000)のIDである。ホストVOLアドレス902は、管理対象の8KBデータが属するスロットのアドレスである。位置番号903は、管理対象の8KBデータが属するスロット(256KB)の中で、当該8KBデータが位置する場所を示す0~31の範囲の番号であり、先頭を0、末尾を31とする。

40

【0070】

システムVOL_ID904は、ストレージ装置11が管理し、ホストVOL1000の8KBデータに対応する、システムVOL1001のグループを含むスロットが属するボリュームのIDである。システムVOLアドレス905は、管理対象の8KBデータに対応する、システムVOL1001のグループを含むスロットのアドレスである。グループ番号906は、管理対象の8KBデータに対応する、システムVOL1001のグルー

50

プの番号であり、0または1である。位置番号907は、管理対象の8KBデータに対応する、システムVOL1001のグループ(128KB)の中で、当該8KBデータが対応する場所を示す0~15の範囲の番号であり、先頭を0、末尾を15とする。

【0071】

なお、システムVOLアドレス905、グループ番号906及び位置番号907において、いずれも“None”が設定されている第3行目のエントリの8KBデータについては、図1の(S7)グループ振分で示したグルーピングが未完了であり、システムVOL1001の対応先グループが未定の状態であることを示す。グルーピングが完了すると、これらの項目にアドレス及び番号が登録される。

【0072】

また、第4行目と第7行目のエントリで管理される2つの8KBデータは、システムVOL1001の同じグループに対応付けられており、両データの内容が類似関係にあることを示す。第2行目と第6行目のエントリで管理される2つの8KBデータも、システムVOL1001の同じグループに対応付けられており、両データの内容が類似関係にあることを示す。第5行目と第8行目のエントリで管理される2つの8KBデータは、システムVOL1001の同じグループに対応付けられており、かつ位置番号も同一であるため、両データの内容が重複関係にあることを示す。

【0073】

図10は、図2のメモリ割当管理テーブルの構成例を示す図である。

図10において、メモリ割当管理テーブル213は、キャッシュ領域203やバッファ領域202を利用するスロットについて、スロットアドレス毎にエントリを有する。各エントリは、VOL_ID908、VOLアドレス909、バッファ(BF)転送状態910、バッファ(BF)アドレス911、キャッシュ状態912及びキャッシュアドレス913の情報を格納する。

【0074】

VOL_ID908は、VOLアドレスによって識別されるスロットが属するボリュームのIDである。VOLアドレス909は、ホストVOL1000またはシステムVOL1001のスロットアドレスである。

【0075】

BF転送状態910は、システムVOL1001の管理対象スロットに含まれるグループのデータを圧縮したデータを一時的に保持するため、バッファ領域202に転送されたか否かの状態であり、“未”はまだ転送されていないことを示し、“済”はすでに転送されたことを示す。ホストVOL1000の管理対象スロットは、本項目を使わない。

【0076】

BFアドレス911は、圧縮されたシステムVOL1001のグループデータが転送されたバッファ領域202内のアドレスを示す。ホストVOL1000の管理対象スロットは、本項目を使わない。BF転送状態910が“未”の場合は、BFアドレス911に“None”が設定され、管理対象スロットに含まれるグループデータが圧縮されていないことを意味する。BF転送状態910が“済”であり、BFアドレス911が値を持つ場合、バッファ領域202に圧縮されたグループデータが保持されている状態を意味する。BF転送状態910が“済”であり、BFアドレス911が“None”である場合、バッファ領域202にあった圧縮されたグループデータがすでにドライブ29へ格納され、バッファ領域202の使用部分のアドレスは解放されたことを意味する。

【0077】

キャッシュ状態912は、ホストVOL1000またはシステムVOL1001の管理対象スロットのデータについて、ドライブ29への格納状態を示す。キャッシュ状態912が“Dirty”のデータは、ドライブ29への格納が未完了である状態を意味し、“Clean”のデータはドライブ29への格納が完了した状態を意味する。なお、ホストVOLスロットのキャッシュ状態912では、32個の8KBデータ毎にキャッシュ状態を管理し、システムVOLスロットのキャッシュ状態912では、2つの128KBグル

10

20

30

40

50

ープ毎に状態を管理する。システムVOLスロットのキャッシュ状態912が、2グループとも“Dirty”から“Clean”に変わる時、バッファ領域202の使用アドレスは解放され、BFアドレス911に“None”が設定される。ただし“Clean”になっても、キャッシュ領域203には、当該グループのデータは圧縮されていない状態で残る。当該データがキャッシュから落ちるまで、ホスト計算機30からストレージ装置11への当該データのリード命令はキャッシュヒットする。

【0078】

キャッシュアドレス913は、ホストVOL1000またはシステムVOL1001の管理対象スロットのデータを格納するために割り当てられたキャッシュ領域203の部分のアドレスである。キャッシュ状態912が全て“Clean”であるスロットは、キャッシュから落とすことができる。キャッシュから落とす場合、メモリ割当管理テーブル213からそのスロットのエントリを削除する。

10

【0079】

図11は、図2のLRUリストテーブルとその遷移を示す図である。

図11において、LRUリストテーブル217は、圧縮対象グループを選択するため、スロットのアクセス頻度を管理する。システムVOLアドレス2101は、アクセス頻度を管理するスロットの属するシステムVOLのIDを示す。システムVOLアドレス2102は、アクセス頻度を管理するスロットのアドレスを示す。本リストの上位にあるスロットほどアクセス頻度が低いことを意味する。

【0080】

あるデータがホスト計算機30からアクセスされると、アクセスデータを含むスロット（例えばアドレス=2600）は本リストの最下位に移動し、そのスロットより下位だった全てのスロット（例えばアドレス=3300、2400、4300など）はそれぞれ1段上へ移動する。このような移動処理を行うと、アクセス頻度の低いデータを含むスロットほどリスト上位に集まる。圧縮対象グループは、例えば、本リストの最上位2103に位置するスロット（例えばアドレス=4100）に含まれる2つのグループから選択することができる。

20

【0081】

図12は、図2のストレージ装置におけるキャッシュメモリのリソース配分比率の変化を示す図である。

30

図12において、図2のキャッシュ領域203は、ホストボリュームやシステムボリュームのスロットのデータを一時的に保持する。データの保持状態によって、キャッシュ領域203のメモリリソースは、以下の3つの状態に分けられる。

【0082】

データが保持されており、そのデータがすでにドライブ29に格納されている領域部分は、“Clean”である。データが保持されており、そのデータがまだドライブ29に格納されていない領域部分は、“Dirty”である。キャッシュ領域203のうちデータが保持されていない領域部分は、“Free”である。キャッシュ領域203は、それら3つの状態の配分比率1010で配分される。

【0083】

ホスト計算機30からデータのライト命令を受けて、キャッシュ領域203にデータが新たに保持される時、Free状態のリソースが使われて、その領域部分がDirty状態に変化する（1011）。このとき、キャッシュ領域203のFreeの比率が減り、Dirtyの比率が増える。

40

【0084】

キャッシュ領域203に保持したDirty状態のデータがドライブ29に格納されると、そのデータの保持領域は、Dirty状態からClean状態に変化する（1012）。このとき、キャッシュ領域203のDirtyの比率が減り、Cleanの比率が増える。

【0085】

50

キャッシュ領域 203 に新たにライトデータが受け入れられる Free 状態の領域を増やすためには、Clean 状態の保持領域を解放する必要がある。このとき、Clean 状態の保持領域のデータは、ドライブ 29 に格納済みなので、Clean 状態の保持領域を解放してもストレージ装置 11 としてデータは損失しない。これにより、その解放領域は、Clean 状態から Free 状態に変化する (1013)。このとき、キャッシュ領域 203 の Clean の比率が減り、Free の比率が増える。

【0086】

キャッシュ領域 203 のメモリリソースは、以上のようなライフサイクルで管理される。図 12 の 1012 で示す Dirty 状態から Clean 状態へ変化させるには、図 1 の (S8) 圧縮 ~ (S10) デステージの処理を行う必要がある。ホスト計算機 30 からデータのライト命令を受けてから、Free 状態のリソースを用意するため、(S8) 圧縮 ~ (S10) デステージの処理により Dirty 状態から Clean 状態へ変化させ、Clean 状態の領域の解放によって Free 状態へ変化させると、(S1) ライトから (S5) 応答までに長い時間がかかる。

10

【0087】

そこで、図 2 のストレージコントローラ 22 は、Dirty 状態のリソースの比率がある閾値に達したことを契機に、キャッシュ領域 203 にあるグループデータに対して (S8) 圧縮を開始して (S10) デステージまでの処理を行い、Dirty 状態から Clean 状態へ事前に変化させる。これにより、ストレージコントローラ 22 は、ライト命令を受けたら Clean 状態の領域の解放だけを行って Free 状態のリソースを直ぐに用意し、(S1) ライトから (S5) 応答までの時間を短縮することができる。

20

【0088】

Dirty 状態のリソースの比率がある閾値に達した時に、(S8) 圧縮の対象として選択するグループは、キャッシュ領域 203 において、その時点で最も長い時間アクセスされていないスロット (つまり、LRU スロット) に含まれるグループである。

【0089】

グループの選択方法は、この方法に限定されない。アクセスされていない時間が長い順に複数のスロット (M (M は 2 以上の整数) 個) を選択し、それらスロットに含まれるグループ (2M 個) の中で、図 6 のプール割当管理テーブルにおける類似度 77 が最も高いグループを、(S8) 圧縮の対象として選択してもよい。例えば、図 11 のリスト上位範囲 2104 は、アクセスされていない時間が長い複数のスロット (M = 4) を含む。これらのスロットに含まれる 8 つのグループで最も類似度 77 が高いグループを選択する。

30

【0090】

この方法の利点を以下に説明する。類似度が低いグループデータは、圧縮率が低い (データの削減量が少ない) ことが予想される。そのため、そのグループデータは、圧縮せずにキャッシュ領域 203 に残す方がよい。もし、新たに受け入れたライトデータの中に、より類似したデータがあれば、キャッシュ領域 203 の上でグループを作り直して類似度 77 を高めれば、圧縮率を改善できる可能性がある。一方、類似度が高いグループデータは、圧縮率が高い (データの削減量が多い) ことが期待できるため、キャッシュ領域 203 に残しても、新たに受け入れたライトデータの中に、より類似したデータがあることは期待できず、グループを作り直しても圧縮率が改善される可能性は低い。従って、最も類似度が高いグループを、(S8) 圧縮の対象として選択するのが望ましい。

40

【0091】

図 13 は、図 2 のストレージ装置におけるスロットとグループの構成例を示す図である。

図 13 において、ライト命令によって書かれる 8 KB データ 1201 は、ホスト計算機 30 から見えるボリューム 1210 で管理される。ボリューム 1210 のアドレス空間は、複数のスロット 1200 で構成され、ライト命令によって書かれる 8 KB データ 1201 の割当先アドレス 1100 は、スロット 1200 の 1 つが占めるアドレス空間に含まれる。

50

【 0 0 9 2 】

「スロットの排他を確保」とは、ホスト計算機 3 0 からのリード命令及びライト命令で指定されたアドレスが示すスロットに対するリード及びライトを防ぐ操作であり、排他を確保したことをホスト計算機 3 0 が認識するための情報が管理される。なお、この情報はビットマップ又は時間情報など識別できるものであれば種別は問わない。また、本実施形態において、「スロット」が、ボリューム（例えば、シンプロビジョニングに従うボリューム）における領域単位であるのに対し、「データ領域」は、スロットに割り当てられる領域（例えば、プール内の領域であるプール領域）である。例えば、ボリューム内の 1 個の 2 5 6 K B スロットに割り当てられたプール領域には、8 K B データを 3 2 個格納できる。

10

【 0 0 9 3 】

類似関係にある 1 6 個の 8 K B データを含むグループは、ホスト計算機 3 0 から見えない別のボリューム 1 2 2 0 で管理される。ボリューム 1 2 2 0 のアドレス空間も、複数のスロット 1 2 0 2 で構成される。ボリューム 1 2 2 0 のスロット 1 2 0 2 は、それぞれ 2 つのグループ 1 2 0 3 で構成される。それら 2 つのグループ 1 2 0 3 のアドレスは、それらを構成するスロット 1 2 0 2 が占めるアドレス空間に含まれる。また、ある 1 つのグループ 1 2 0 3 に 1 6 個含まれる、類似関係にある 8 K B データ 1 2 0 4 のアドレスは、そのグループ 1 2 0 3 が占めるアドレス空間に含まれる。

【 0 0 9 4 】

図 1 4 は、図 2 のストレージ装置が実行する類似ハッシュ算出処理を示す図である。なお、図 1 4 では、簡単のため、類似ハッシュの対象データの数は 2 つ（データ A、B）とし、各データのサイズは 3 7 B である場合を例にとる。データの数およびサイズが増えても、類似ハッシュの算出方法は同様である。

20

【 0 0 9 5 】

図 1 4 において、各データ A、B について、例えば、第 K 文字目から始まる連続 3 文字を 3 5 個抽出する。ここで、K は、1 ~ 3 5 である。以下、これらの 3 文字を単語と呼ぶ。2 つのデータが共通の単語を多く含むかを判断することで類似性の強さを評価する。単語は、各データ A、B 内のどの位置に存在してもよいので、開始点（すなわち、K）を 1 バイトステップで変えながら単語を抽出する。

【 0 0 9 6 】

次に、各データ A、B において、3 5 個の単語の出現頻度（以下、頻度と言う）を示す頻度順位表 1 4 1、1 4 2 を作成する。各頻度順位表 1 4 1、1 4 2 において、頻度の高いものから順に単語を並べる。ただし、頻度の少ない単語（例えば 1 回の単語）は、各頻度順位表 1 4 1、1 4 2 から除く。

30

【 0 0 9 7 】

図 1 4 では、2 つの頻度順位表 1 4 1、1 4 2 において共通する単語同士を線で結んでいる。共通する単語の数（積集合の要素数）を、両者の単語の和集合の要素数で割ることで、データ A、B 間の類似関係の強さ S がわかる。この例では、共通単語数が 6、和集合の要素数が 1 1 であるため、強さ $S = 6 / 1 1$ である。

【 0 0 9 8 】

なお、頻度順位表 1 4 1、1 4 2 から頻度の少ない単語を除いた理由は、スライド辞書圧縮において頻度の少ない単語は一致文字列になる可能性が低く、頻度の少ない単語よりも多い単語の共通性を評価する方が、圧縮率を高くする上で効果的だからである。

40

【 0 0 9 9 】

なお、図 1 4 の例では、単語のサイズを 3 バイトとしたが、それに限定しない。スライド辞書圧縮の多くが 3 ~ 4 バイト程度を一致文字列の最小検出長としているため、3 ~ 4 バイト程度が望ましい。

【 0 1 0 0 】

全てのデータについて頻度順位表をその付属情報として管理し、互いの類似関係の強さ S の計算に用いると、ストレージ装置 1 1 における管理情報の記憶量が多くなる。そこで

50

、全てのデータの頻度順位表を管理しなくても、その中のある2つのデータ間の類似関係の強さ S を近似的に評価できる `b-bit Min-wise Hash` という LSH アルゴリズムを適用する。そのアルゴリズムを以下に示す。

【0101】

n 個のハッシュ関数を用意する。次に、データ A の頻度順位表 141 に出現する 10 個の単語 144 をそれぞれ n (n は 2 以上の整数) 個のハッシュ関数に通してハッシュ値を計算する。この結果、データ A から $10 \times n$ 個のハッシュ値が生成される。

【0102】

そして、各ハッシュ関数 k (k は $1 \sim n$) からのハッシュ値 10 個の中の最小値 146 を求め、 n 個の最小値 146 を得る。それら n 個の最小値 146 を並べてデータ A についての類似ハッシュ 147 を構成する。

10

【0103】

データ B にも同様の処理を行う。つまり、 $7 \times n$ 個のハッシュ値が生成され、各ハッシュ関数 k からの 7 個のハッシュ値の中の最小値 146 を求め、 n 個の最小値 146 を得る。それら n 個の最小値 146 を並べてデータ B についての類似ハッシュ 147 を構成する。

【0104】

このようにして計算した各データ A、B についての 2 つの類似ハッシュ 147 のハミング距離を H 、ハッシュ関数のビット長を b とすると、類似関係の強さ S の近似値 J は、以下の式で与えられる。

20

【0105】

$$J = \{ (n - H) / n - (1/2)^b \} / \{ 1 - (1/2)^b \}$$

n と b が大きいほど、近似値 J の近似精度が向上する。

【0106】

以上の方法により、各データの類似ハッシュ 147 (各サイズは $n \times b$ ビット) を管理しておけば、ある 2 つのデータ間の類似関係の強さ S を近似的に求めることができる。例えば、 $n = 16$ 、 $b = 8$ の場合、8 KB データ当たり 16 B の付属情報 (0.2%) を持つだけでよい。

【0107】

ある 2 つのデータの類似ハッシュ 147 のハミング距離 $H = 0$ であるとき、 $J = 1$ となる。つまり、両データの類似関係の強さ S は 100% である。ハミング距離 $H = 0$ (類似ハッシュ 147 が一致) のとき、両データは全く同じ内容である可能性がある。このため、重複排除処理のための重複データ検出に類似ハッシュ 147 を用いることも可能である。実際に重複しているか否かは、8 KB 全体を比較して判定する必要があるが、類似ハッシュ 147 が一致するか否かで、比較候補を絞り込むことができる。以上、類似ハッシュ算出方法の一例を説明したが、類似ハッシュ算出方法は、この方法に限定されない。

30

【0108】

図 15 は、図 2 のストレージ装置が実行するデータリード処理を示すフローチャートである。

図 15 において、リード処理は、図 2 のホスト計算機 30 からネットワーク 31 を介してストレージ装置 11 がリード命令を受けた場合に開始される。リード命令では、例えば、ボリューム ID、アドレス及びデータサイズが指定される。

40

【0109】

S1801 において、プロセッサ 24 は、指定アドレスから特定されるスロットの排他を確保する。なお、スロット排他確保時に他の処理がスロットの排他を確保している場合、プロセッサ 24 は、一定の時間待ってから、S1801 を行う。

【0110】

次に、S1802 において、プロセッサ 24 は、リード対象データがキャッシュ領域 203 に存在するか否かを判定し、S1802 の判定結果が真の場合、S1807 に進む。

【0111】

50

一方、S 1 8 0 2 の判定結果が偽の場合、プロセッサ 2 4 は、S 1 8 0 3 において、R A I D グループを構成するドライブから対象データ含む圧縮されたグループデータをリードする。この際、プロセッサ 2 4 は、ホスト計算機 3 0 が指定したボリューム I D とアドレスから、図 6 のプール割当管理テーブル 2 1 0 のプール I D 7 3、プールアドレス 7 4 及び圧縮後サイズ 7 6 を特定し、図 7 のドライブ割当管理テーブル 2 1 1 からドライブ I D 8 4 及びドライブアドレス 8 5 を参照し、対象データを含む圧縮されたグループデータのドライブ上の格納場所及びサイズを特定する。

【 0 1 1 2 】

次に、S 1 8 0 4 において、プロセッサ 2 4 は、ドライブからリードした圧縮されたグループデータをバッファ領域 2 0 2 にライトする。

【 0 1 1 3 】

次に、S 1 8 0 5 において、プロセッサ 2 4 は、バッファ領域 2 0 2 上にライトした圧縮されたグループデータを伸張する。

【 0 1 1 4 】

次に、S 1 8 0 6 において、プロセッサ 2 4 は、バッファ領域 2 0 2 上の伸張されたグループデータから指定サイズの対象データを抽出する。

【 0 1 1 5 】

次に、S 1 8 0 7 において、プロセッサ 2 4 は、バッファ領域 2 0 2 上の対象データをホスト計算機 3 0 に転送する。ホスト計算機 3 0 は、このデータ転送が完了した時点でリード処理が終了したと認識する。

【 0 1 1 6 】

次に、S 1 8 0 8 において、プロセッサ 2 4 は、S 1 8 0 1 で確保していたスロット排他を解放する。

【 0 1 1 7 】

なお、S 1 8 0 5 の伸張処理は、プロセッサ 2 4 の負荷低減のため、図 2 のオフロードエンジン 2 1 4 が高性能なハードウェアを用いて代わりに行ってよい。その際に、オフロードエンジン 2 1 4 は、バッファ領域 2 0 2 の代わりにバッファ領域 2 1 6 を使うことにより、メモリ 2 5 の転送帯域の消費を抑えることができる。

【 0 1 1 8 】

図 1 6 は、図 2 のストレージ装置が実行するデータライト処理を示すフローチャートである。なお、以下の説明では、例えば、図 1 のストレージコントローラ 2 2 A に対応する図 2 のプロセッサ 2 4 をプロセッサ 2 4 A と記載するなど、図 1 の各ストレージコントローラ 2 2 A、2 2 B に属するものをそれぞれ参照符号に付した「A」及び「B」によって区別する。

【 0 1 1 9 】

図 1 6 において、データライト処理は、図 2 のホスト計算機 3 0 からストレージ装置 1 1 がライト命令を受信した場合に開始される。ライト命令では、例えば、ボリューム I D、アドレス及びデータサイズが指定される。このデータライト処理は、図 1 の (S 1) ~ (S 5) で実行される。

【 0 1 2 0 】

S 1 5 0 1 において、プロセッサ 2 4 A は、指定アドレスから特定されるスロットの排他を確保する。なお、スロット排他確保と同時に、プロセッサ 2 4 A は、データのライト先とするキャッシュ領域 2 0 3 A の部分を割り当てる。

【 0 1 2 1 】

次に、S 1 5 0 2 において、プロセッサ 2 4 A は、ホスト計算機 3 0 に対してライト処理の準備ができたことを示す「R e a d y」を応答する。プロセッサ 2 4 A は、「R e a d y」を受け取ったホスト計算機 3 0 から、ライトデータを受信する。

【 0 1 2 2 】

次に、S 1 5 0 3 において、プロセッサ 2 4 A は、ホスト計算機 3 0 から受信したライトデータをキャッシュ領域 2 0 3 A に割り当てたスロット領域にライトする。

10

20

30

40

50

【 0 1 2 3 】

次に、S 1 5 0 4において、プロセッサ 2 4 Aは、ストレージコントローラ 2 2 Aからストレージコントローラ 2 2 Bに対してキャッシュ領域 2 0 3 Aに格納したライトデータを転送し、キャッシュ領域 2 0 3 Bに格納することで二重化を行う。

【 0 1 2 4 】

次に、S 1 5 0 5において、プロセッサ 2 4 Aは、図 9のグループ管理テーブル 2 1 2と図 10のメモリ割当管理テーブル 2 1 3を更新する。ライトデータは、重複排除および類似グルーピング（両者を合わせて、圧縮前データ構築処理と呼ぶ）がまだ実施されていない。なお、類似グルーピングは、類似関係にあるデータのグルーピングである。そのため、プロセッサ 2 4 Aは、グループ管理テーブル 2 1 2において、ライトデータのホストVOLスロットに対応するシステムVOLスロットのアドレス 9 0 5、グループ番号 9 0 6および位置情報 9 0 7に“None”を設定する。また、プロセッサ 2 4 Aは、メモリ割当管理テーブル 2 1 3で、ライトデータのホストVOLスロットに対応するキャッシュ状態 9 1 2において、ライトデータのキャッシュ状態を“Dirty”に設定し、キャッシュアドレス 9 1 3には、キャッシュ領域 2 0 3 Aのライトデータ格納先アドレスを設定する。

10

【 0 1 2 5 】

次に、S 1 5 0 6において、プロセッサ 2 4 Aは、ネットワーク 3 1を介してホスト計算機 3 0に対してライト処理が完了したとして完了応答を返却する。

【 0 1 2 6 】

次に、S 1 5 0 7において、プロセッサ 2 4 Aは、S 1 5 0 1で確保していたスロットの排他を解放してライト処理を終了する。

20

【 0 1 2 7 】

図 1 7は、図 2のストレージ装置が実行する圧縮前データ構築処理を示すフローチャートである。

図 1 7において、圧縮前データ構築処理は、ライト処理が終了した後にプロセッサ 2 4が行う。圧縮前データ構築処理は、ライト処理終了直後に実行してもよいし、他に優先して行うべき処理があれば、その処理の後に実行してもよい。この圧縮前データ構築処理は、図 1の(S 6)～(S 7)で実行される。

【 0 1 2 8 】

S 1 6 0 1において、プロセッサ 2 4は、キャッシュ領域 2 0 3に格納したホストVOL上のライトデータを選択する。

30

【 0 1 2 9 】

次に、S 1 6 0 2において、プロセッサ 2 4は、選択したライトデータのアドレスから特定されるスロットの排他を確保する。

【 0 1 3 0 】

次に、S 1 6 0 3において、プロセッサ 2 4は、例えば、図 1 4で説明した方法に従って、選択したデータの類似ハッシュを算出する。

【 0 1 3 1 】

次に、S 1 6 0 4において、プロセッサ 2 4は、算出した類似ハッシュをハッシュテーブルに登録し、それまでに登録された1つ以上の類似ハッシュとの間でハミング距離Hを計算して、他のデータとの類似度を評価する。

40

【 0 1 3 2 】

次に、S 1 6 0 5において、類似ハッシュが一致し、かつ内容も一致するデータがある場合は、プロセッサ 2 4は、そのデータとライトデータとの間で重複排除処理 S 1 6 0 6を実行し、S 1 6 1 1に進む。

【 0 1 3 3 】

重複排除処理 S 1 6 0 6では、プロセッサ 2 4は、ライトデータと内容が一致するデータが含まれるシステムVOL上のスロットとグループを特定し、そのグループ内の一致するデータの位置番号を特定する。プロセッサ 2 4は、その特定結果に基づいて、S 1 6 1

50

1で、図9のグループ管理テーブル212を更新する。すなわち、プロセッサ24は、システムVOLアドレス905、グループ番号906及び位置番号907の値が等しくなるように設定する。このとき、グループ管理テーブル212は、例えば、第5行目と第8行目のエントリの対象データは重複排除されていることを示す。

【0134】

一方、S1605において、類似ハッシュが一致するデータがない場合、または、類似ハッシュが一致し、かつ内容も一致するデータがない場合は、S1607に進む。

【0135】

次に、S1607において、プロセッサ24は、ハミング距離Hの最小値があらかじめ規定した基準値以下であるかを判定する。判定結果が真であれば、プロセッサ24は、類似関係にあるデータが存在すると判断し、S1608において、その最小値をもたらずデータが含まれるシステムVOL上のグループを選択し、S1610に進む。

10

【0136】

一方、S1607の判定結果が偽であれば、プロセッサ24は、類似関係にあるデータは存在しないと判断し、S1609において、システムVOL上に新たにグループを作成する。この時、プロセッサ24は、グループデータ格納先としてキャッシュ領域203の部分を割り当て、S1610に進む。

【0137】

次に、S1610において、プロセッサ24は、システムVOL上の選択または作成されたグループに対応するキャッシュ領域部分に、ホストVOL上のライトデータを転送する。

20

【0138】

次に、S1610から遷移したS1611において、プロセッサ24は、グループ管理テーブル212を更新する。すなわち、プロセッサ24は、システムVOLアドレス905に選択または作成したシステムVOL上のグループを含むスロットのアドレスを設定し、グループ番号906に当該グループの番号を設定し、位置番号907にライトデータが転送されたグループ内での位置を設定する。なお、ライトデータを転送したグループが、新しく作成したグループであり(つまり、S1609を経ており)、当該グループを含むスロットにおける最初のグループである場合は、さらに、プロセッサ24は、メモリ割当管理テーブル213にそのスロットのエントリを追加する。すなわち、プロセッサ24は、VOL_ID908にシステムVOLのIDを設定し、VOLアドレス909に当該グループを含むスロットアドレスを設定する。また、BF転送状態910に“未”を、BFアドレス911に“None”を設定する。また、キャッシュ状態912の当該グループのキャッシュ状態に“Dirty”を設定する。また、キャッシュアドレス913にはキャッシュ領域203の割り当て部分のアドレスを設定する。

30

【0139】

次に、S1612において、プロセッサ24は、ライトデータのアドレスから特定されるスロットの排他を解放して圧縮前データ構築処理を終了する。

【0140】

なお、S1603の類似ハッシュ算出、S1604のハッシュ登録・評価は、プロセッサ24の負荷低減のため、図2のオフロードエンジン214が高性能なハードウェアを用いて代わりに行ってよい。その際、オフロードエンジン214は、メモリ215をワークメモリとして使うことにより、メモリ25の転送帯域の消費を抑えることができる。

40

【0141】

また、類似ハッシュを登録するためのハッシュテーブルは、キャッシュに乗っているデータの類似ハッシュを保持する。データがキャッシュから落とされるとき、そのデータの類似ハッシュをハッシュテーブルから削除する。

【0142】

また、S1610において、選択したグループのキャッシュ領域203にライトデータを格納する領域がない場合は、プロセッサ24は、そのグループの所属データで最も他の

50

所属データとの類似性が低いデータを選択し、そのグループから除外してライトデータの格納先を確保する。すなわち、プロセッサ24は、グループ管理テーブル212において、除外対象データのエントリのシステムVOLアドレス905、グループ番号906及び位置番号907に“None”を設定する。これは、除外されたデータが再び重複排除および類似グルーピング(図16参照)が施される準備ができたことを意味する。

【0143】

図18は、図2のストレージ装置が実行する圧縮・デステージ処理を示すフローチャートである。

図18において、圧縮・デステージ処理は、圧縮前データ構築処理が終了したグループデータについて、キャッシュ領域203の“Dirty”比率が閾値以上となったことを契機として、非同期的に行われる。ただし、プロセッサ24が低負荷な状態であれば、“Dirty”比率が閾値未満であっても、デステージ処理を開始してもよい。また、“Dirty”比率を問わず、タイマを用いて周期的にデステージ処理を開始してもよい。この圧縮・デステージ処理は、図1の(S8)~(S10)で実行される。

10

【0144】

S1701において、プロセッサ24は、キャッシュ領域203内で“Dirty”状態のスロットに含まれるグループの中から、デステージ対象のグループを選択する。具体的には、プロセッサ24は、上述の通り、LRUスロットに含まれるグループまたは最後のアクセスからの経過時間が長い複数スロットに含まれるグループのうち、最も類似度が高いグループを選択する。

20

【0145】

次に、S1702において、プロセッサ24は、デステージ対象のグループが属するシステムVOLスロットの排他を確保し、さらに、デステージするグループに転送されたデータが属するホストVOLスロットの排他を確保する。

【0146】

次に、S1703において、プロセッサ24は、デステージ対象のグループデータ(128KB)を読み出して、可逆圧縮処理を行う。

【0147】

次に、S1704において、プロセッサ24は、圧縮されたグループデータ(128KB未満)の格納先としてバッファ領域202の部分を割り当て、その格納先へ圧縮グループデータをライトする。なお、バッファ領域202のアドレス割り当てでは、パリティサイクル分の圧縮グループが集まるまで割り当てを繰り返すことが明らかなため、あらかじめパリティサイクル分のアドレスを一括して割り当ててもよい。

30

【0148】

次に、S1705において、プロセッサ24は、図10のメモリ割当管理テーブル213のBF転送状態910を“済”に更新する。また、BFアドレス911に、割り当てたバッファ部分のアドレスを設定する。

【0149】

次に、S1706において、プロセッサ24は、S1702で確保したスロットの排他を解放する。

40

【0150】

次に、S1707において、プロセッサ24は、バッファ内の圧縮グループデータの蓄積量を計算する。圧縮グループデータの蓄積量がパリティサイクル分よりも小さい場合、プロセッサ24は、S1701に戻ってデステージ対象のグループを追加で選択する。

【0151】

一方、パリティサイクル分の圧縮グループデータがバッファ領域202内に溜まった場合は、S1708に進む。なお、圧縮グループデータのサイズは可変であるため、バッファ領域202内の蓄積量が必ずしもパリティサイクル分に一致するとは限らない。パリティサイクル分に達する前にS1708へ処理を進めることもあり得る。

【0152】

50

次に、S 1 7 0 8において、プロセッサ 2 4は、バッファ領域 2 0 2内に蓄積された圧縮グループデータ列からパリティを生成する。

【 0 1 5 3 】

次に、S 1 7 0 9において、プロセッサ 2 4は、圧縮グループデータ列及び生成したパリティを、RAIDグループを構成するドライブ 2 9に書き出す。プロセッサ 2 4は、ドライブ 2 9への書き出し後に、バッファ領域の使用していた部分を解放する。

【 0 1 5 4 】

次に、S 1 7 1 0において、プロセッサ 2 4は、デステージ対象のグループを含むシステムVOLスロットの排他を確保し、さらに、デステージ対象のグループに転送されたデータが属するホストVOLスロットの排他を確保する。

10

【 0 1 5 5 】

次に、S 1 7 1 1において、プロセッサ 2 4は、メモリ割当管理テーブル 2 1 3で、デステージ対象のグループを含むシステムVOLスロットのキャッシュ状態 9 1 2において、当該グループのキャッシュ状態を“Clean”に更新する。その結果、プロセッサ 2 4は、2個のグループのキャッシュ状態がともに“Clean”となれば、BFアドレス 9 1 1を“None”に設定し、バッファ領域 2 0 2の使用部分を解放する。このシステムVOLスロットは“Clean”状態となり、いつでもキャッシュから解放して“Free”状態のキャッシュリソースを作ることができる。さらに、プロセッサ 2 4は、メモリ割当管理テーブル 2 1 3で、デステージ対象のグループに転送された各データが属するホストVOLスロットのキャッシュ状態 9 1 2において、当該データのキャッシュ状態を“Clean”に更新する。その結果、キャッシュ状態 9 1 2が全て“Clean”となったホストVOLスロットは、いつでもキャッシュから解放して“Free”状態のキャッシュリソースを作ることができる。

20

【 0 1 5 6 】

次に、S 1 7 1 2において、プロセッサ 2 4は、デステージ対象のグループを含むホストVOLスロットの排他を解放し、さらに、デステージ対象のグループに転送されたデータが属するホストVOLスロットの排他を解放し、処理を終了する。

【 0 1 5 7 】

以上の圧縮・デステージ処理において、キャッシュから解放されるグループは、アクセス頻度が低く、類似関係が強いデータが優先される。従って、ホスト計算機 3 0からのアクセスがキャッシュヒットする確率がより高くなるため、アクセス性能が向上する。また、ドライブ 2 9に格納されるデータの圧縮率がより高くなるため、データ削減効率が向上する。

30

【 0 1 5 8 】

なお、S 1 7 0 3の可逆圧縮は、プロセッサ 2 4の負荷低減のため、図 2のオフロードエンジン 2 1 4が高性能なハードウェアを用いて代わりに行ってよい。その際、オフロードエンジン 2 1 4は、メモリ 2 1 5のバッファ領域 2 1 6をバッファ領域 2 0 2の代わりとして使うことにより、メモリ 2 5の転送帯域の消費を抑えることができる。

【 0 1 5 9 】

図 1 9は、図 2のストレージ装置が実行するデータ削除リトライ処理を示すフローチャートである。

40

図 1 9において、データ削減リトライ処理は、図 2のドライブ 2 9にデステージした圧縮グループを構成するデータを再度キャッシュ領域 2 0 3に戻して、重複排除および類似グルーピングをリトライできる状態にする。デステージしたグループデータの類似度 7 7が悪かった場合、このデータ削減リトライ処理によって、前より類似関係の強いデータを含むグループが構築され得るため、ストレージ装置 1 1のデータ削減率を改善することができる。

【 0 1 6 0 】

S 1 9 0 1において、プロセッサ 2 4は、図 6のプール割当管理テーブル 2 1 0で管理されているシステムVOLスロットの1つを選択する。

50

【 0 1 6 1 】

次に、S 1 9 0 2において、プロセッサ 2 4は、プール割当管理テーブル 2 1 0で、そのシステムVOLスロットに含まれるグループの類似度 7 7が、あらかじめ規定した基準値以下であるかを判定する。

【 0 1 6 2 】

S 1 9 0 2の判定結果が偽であれば、プロセッサ 2 4は、そのシステムVOLスロットのデータ削減は不要として処理を終了する。判定結果が真であれば、S 1 9 0 3に進む。

【 0 1 6 3 】

次に、S 1 9 0 3において、プロセッサ 2 4は、そのシステムVOLスロットの排他を確保する。さらに、図 9のグループ管理テーブル 2 1 2を参照し、当該システムVOLスロット内の類似度が基準値以下のグループを構成するホストVOLスロットのデータを特定し、そのホストVOLスロットの排他も確保する。

10

【 0 1 6 4 】

次に、S 1 9 0 4において、プロセッサ 2 4は、図 1 0のメモリ割当管理テーブル 2 1 3を参照し、特定したデータがキャッシュ領域 2 0 3に存在するか否かを判定する。

【 0 1 6 5 】

S 1 9 0 4の判定結果が真の場合、プロセッサ 2 4は、メモリ割当管理テーブル 2 1 3において、そのデータのキャッシュ状態 9 1 2に“ D i r t y ”を設定する。また、グループ管理テーブル 2 1 2において、当該ホストVOLスロットのデータのエントリのシステムVOLアドレス 9 0 5、グループ番号 9 0 6及び位置番号 9 0 7に“ N o n e ”を設定し、既存のシステムVOLスロットのグループとの対応関係を切る。これは、当該データが重複排除および類似グルーピング（図 1 6参照）が施される準備ができたことを意味する。そして、S 1 9 0 9に進む。

20

【 0 1 6 6 】

一方、S 1 9 0 4の判定結果が偽の場合、プロセッサ 2 4は、S 1 9 0 5において、R A I Dグループのドライブから、S 1 9 0 2で類似度が基準値以下と判定した圧縮グループデータをリードする。この際、プロセッサ 2 4は、図 6のプール割当管理テーブル 2 1 0のプールID 7 3、プールアドレス 7 4及び圧縮後サイズ 7 6を特定し、ドライブ割当管理テーブル 2 1 1からドライブID 8 4及びドライブアドレス 8 5を参照し、その圧縮グループデータのドライブ上の格納場所及びサイズを特定する。

30

【 0 1 6 7 】

次に、S 1 9 0 6において、プロセッサ 2 4は、その圧縮グループデータをバッファ領域 2 0 2にライトする。

【 0 1 6 8 】

次に、S 1 9 0 7において、プロセッサ 2 4は、バッファ領域 2 0 2上の圧縮グループデータを伸張する。

【 0 1 6 9 】

次に、S 1 9 0 8において、プロセッサ 2 4は、キャッシュ領域 2 0 3内に当該データの格納用の領域を割り当て、バッファ領域 2 0 2上の伸張されたグループから抽出したデータを、その割り当て領域にライトする。この際、プロセッサ 2 4は、メモリ割当管理テーブル 2 1 3のVOL _ I D 9 0 8に当該データの属するホストVOLのIDを設定し、VOLアドレス 9 0 9に当該データの属するホストVOLスロットのアドレスを設定し、キャッシュ状態 9 1 2に“ D i r t y ”を設定し、キャッシュアドレス 9 1 3にその割り当て領域のアドレスを登録する。また、プロセッサ 2 4は、グループ管理テーブル 2 1 2において、当該ホストVOLスロットのデータのエントリを削除し、既存のシステムVOLスロットのグループとの対応関係を切る。これは、当該データは重複排除や類似グルーピング（図 1 6参照）が施される準備ができたことを意味する。そして、S 1 9 0 9に進む。

40

【 0 1 7 0 】

次に、S 1 9 0 9において、プロセッサ 2 4は、S 1 9 0 3で確保したシステムVOL

50

スロット及びホストVOLスロットの排他も解放し、データ削減リトライ処理を終了する。

【0171】

なお、S1907の伸張処理は、プロセッサ24の負荷低減のため、図2のオフロードエンジン214が高性能なハードウェアを用いて代わりに行ってよい。その際、オフロードエンジン214は、メモリ215のバッファ領域216をバッファ領域202の代わりとして使うことにより、メモリ25の転送帯域の消費を抑えることができる。

【0172】

図20は、図2のストレージ装置におけるデータ伸張処理の過程を示す図である。

図20において、プロセッサ24は、類似関係にある複数の8KBデータをグルーピングする処理において、データの並び順と伸張処理を最適化することで、データリード処理の性能を改善する。

10

【0173】

データリード処理では、ホスト計算機30からリード命令を受けたデータがキャッシュ領域203にはなく、ドライブ29に格納されていた場合、プロセッサ24は、圧縮されたグループデータ1103を伸張して、128KBサイズのグループ1102を復元し、その復元したデータからリード対象の8KBデータを抽出してホスト計算機30に応答する。

【0174】

このとき、プロセッサ24は、圧縮されたグループデータの伸張処理において、先頭の8KBデータから順に復元する。リード対象データが1100Cの場合、プロセッサ24は、復元までに範囲2001の部分を伸張すればよい。リード対象データが1100Dの場合、プロセッサ24は、復元までに範囲2002の部分を伸張すればよい。すなわち、リード対象データがグループデータの先頭に近いほど、プロセッサ24は、短い時間でデータを復元できる。

20

【0175】

プロセッサ24は、この特性を利用して、リード対象データが復元できた時点でグループデータの伸張処理を中断する。さらに、プロセッサ24は、ホストVOLのアドレス空間を複数の区間に分けて、所定時間範囲（例えば、過去10時間）におけるそれぞれの区間のアクセス頻度を調べる。

30

【0176】

プロセッサ24は、類似グルーピングの際に、ホスト計算機30からのアクセス頻度が高い区間にあるデータほど、グループの先頭に近い（位置番号が小さい）位置に配置する。アクセス頻度が低い区間にあるデータほど、グループの末尾に近い（位置番号が大きい）位置に配置する。これにより、ストレージ装置11は、ホスト計算機30からのリード命令に対して、平均応答時間を短縮することができる。

【0177】

なお、本発明は上記した実施形態に限定されるものではなく、様々な変形例が含まれる。例えば、上記した実施形態は、本発明のより良い理解のために詳細に説明したのであり、必ずしも説明の全ての構成を備えるものに限定されるものではない。

40

【0178】

また、上記の各構成、機能、処理部、処理手段等は、それらの一部又は全部を、例えば集積回路で設計する等によってハードウェアで実現してもよい。すなわち、類似ハッシュ算出、ハッシュ登録、ハッシュ評価及び可逆圧縮・伸張処理をオフロードエンジンの高速ハードウェアが代行して処理するという変形例以外の実施形態もある。

【0179】

また、上記の各構成、機能等は、プロセッサがそれぞれの機能を実現するプログラムを解釈し、実行することによってソフトウェアで実現してもよい。各機能を実現するプログラム、テーブル、ファイル等の情報は、不揮発性半導体メモリ、HDD、SSD等の記憶デバイス、または、ICカード、SDカード、DVD等の計算機読み取り可能な非一時的

50

データ記憶媒体に格納することができる。

【0180】

また、制御線及び情報線は説明上必要と考えられるものを示しており、製品上必ずしも全ての制御線及び情報線を示しているとは限らない。実際には、ほとんど全ての構成が相互に接続されていると考えてもよい。

【0181】

なお、本発明は上記した実施形態に限定されるものではなく、様々な変形例が含まれる。例えば、上記した実施形態は本発明を分かりやすく説明するために詳細に説明したものであり、必ずしも説明した全ての構成を備えるものに限定されるものではない。また、ある実施形態の構成の一部を他の実施形態の構成に置き換えることが可能であり、また、ある実施形態の構成に他の実施形態の構成を加えることも可能である。また、各実施形態の構成の一部について、他の構成の追加・削除・置換をすることが可能である。また、上記の各構成、機能、処理部、処理手段等は、それらの一部又は全部を、例えば集積回路で設計する等によりハードウェアで実現してもよい。

10

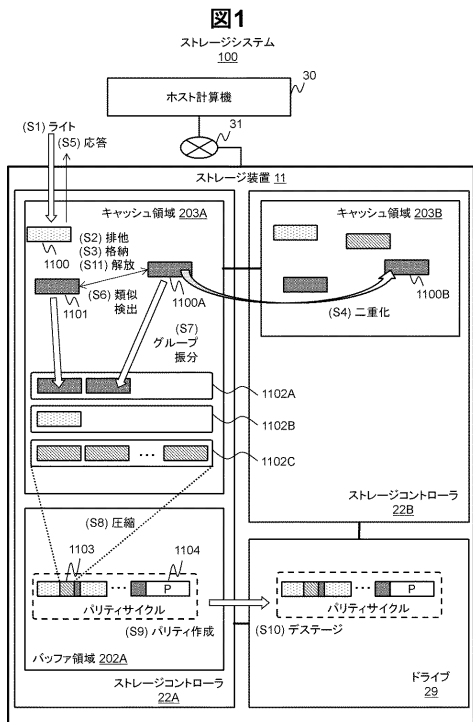
【符号の説明】

【0182】

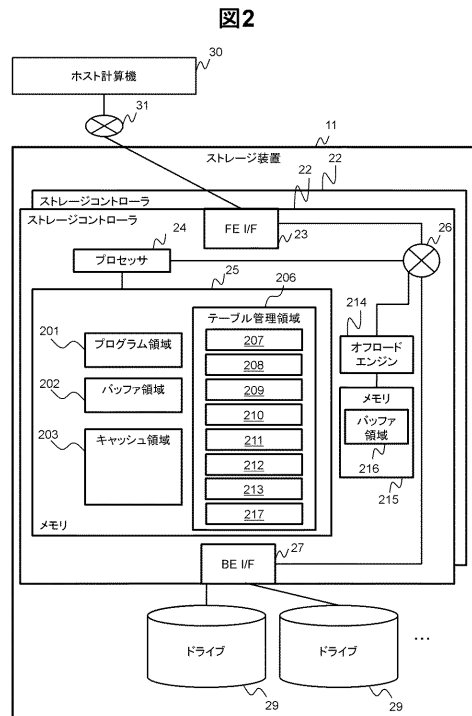
100 ストレージシステム、11 ストレージ装置、22、22A、22B ストレージコントローラ、24 プロセッサ、25、215 メモリ、202、216 バッファ領域、203、203A、203B キャッシュ領域、214 オフロードエンジン、29 ドライブ、30 ホスト計算機、31 ネットワーク

20

【図1】



【図2】



【 図 3 】

図3
ボリューム管理テーブル
207

VOL ID	VOL属性	VOL容量	プールID	ホスト可視
0	シンプロビジョニング、削減OFF	100GB	0	Yes
10	シンプロビジョニング、削減ON	200GB	0	Yes
20	通常割当、削減OFF	500GB	1	Yes
30	システム(グループ管理用)	200GB	0	No
...

【 図 4 】

図4
プール構成管理テーブル
208

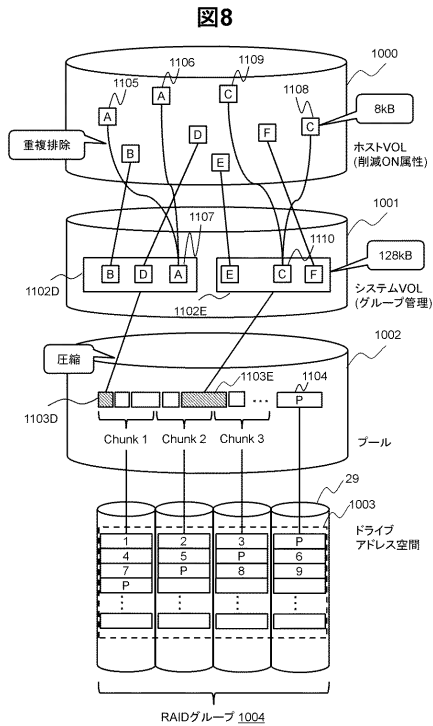
プールID	RAIDグループID	プール容量	プール使用容量
0	0	10TB	5TB
...

【 図 5 】

図5
RAID構成管理テーブル
209

RAIDグループID	RAIDレベル	ドライブID	ドライブ種別	容量	使用容量
0	RAID5	0 1 2 3	SSD	5TB	3TB
...

【 図 8 】



【 図 6 】

図6
プール割当管理テーブル
210

VOL ID	VOLアドレス	プールID	プールアドレス	圧縮前サイズ	圧縮後サイズ	類似度
30	1000	0	05610	128KB	54KB	84%
		0	17480	128KB	29KB	97%
30	1300	0	07140	128KB	73KB	71%
		0	14990	128KB	61KB	80%
30	1500	0	06340	128KB	92KB	54%
		0	04720	128KB	77KB	66%
30	1800	0	11390	128KB	42KB	90%
		0	10570	128KB	39KB	92%
...

【 図 7 】

図7
ドライブ割当管理テーブル
211

プールID	プールアドレス	RAIDグループID	ドライブID	ドライブアドレス
0	05610	0	0	0360
0	17480	1	2	1530
0	07140	1	3	0940
0	14990	1	3	0730
0	06340	0	1	1830
0	04720	0	0	1220
0	11390	1	2	0790
...

【 図 9 】

図9
グループ管理テーブル
212

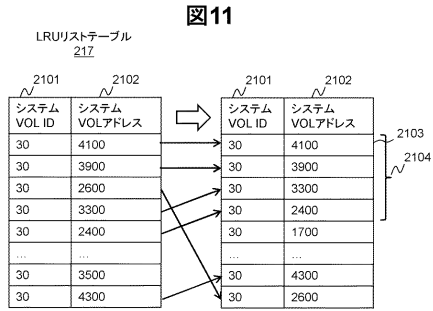
ホストVOL ID	ホストVOL アドレス	位置番号	システムVOL ID	システムVOL アドレス	グループ番号	位置番号
10	700	11	30	1300	1	7
10	500	8	30	1800	0	2
10	300	5	30	None	None	None
10	100	19	30	1500	1	10
10	200	10	30	1000	1	3
10	400	31	30	1800	0	5
10	600	17	30	1500	1	8
10	800	9	30	1000	1	3
...

【 図 10 】

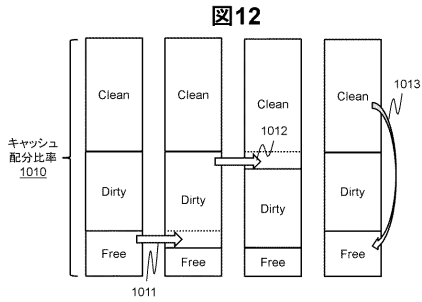
図10
メモリ割当管理テーブル
213

VOL ID	VOL アドレス	BF転送状態	BF アドレス	キャッシュ状態	キャッシュアドレス
10	300			[Dirty, Dirty, ..., Dirty]	2000
10	200			[Dirty, Dirty, ..., Dirty]	2200
10	500			[Dirty, Dirty, ..., Dirty]	2500
30	1300	未	None	[Dirty, Clean]	3000
30	1800	済	50	[Dirty, Dirty]	2100
30	1500	済	None	[Clean, Clean]	2300
30	1000	済	80	[Dirty, Clean]	3200
...

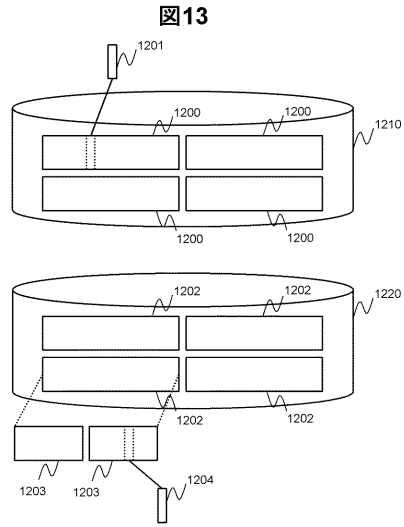
【図11】



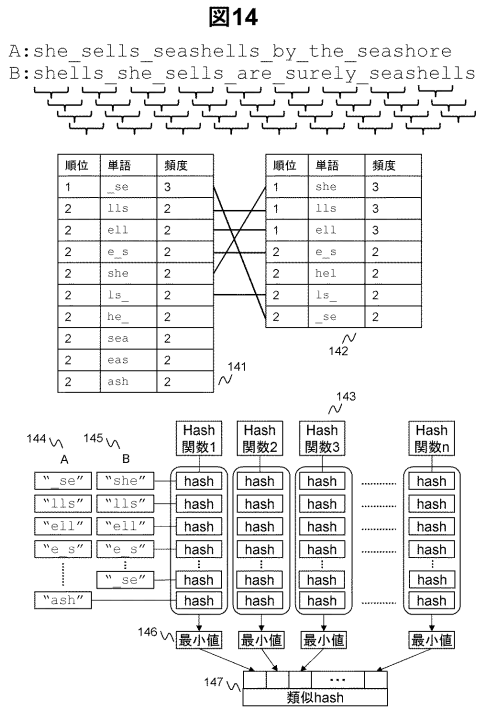
【図12】



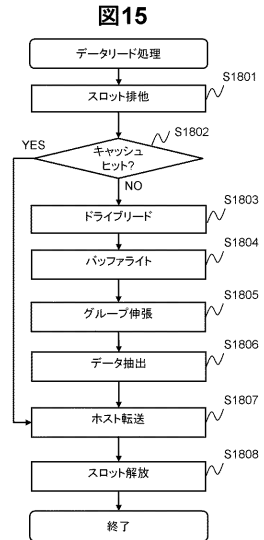
【図13】



【図14】



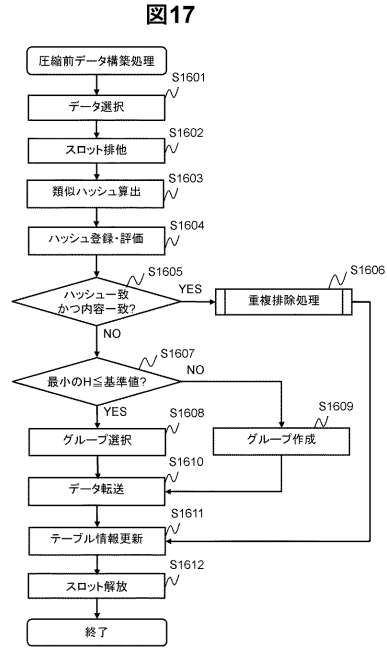
【図15】



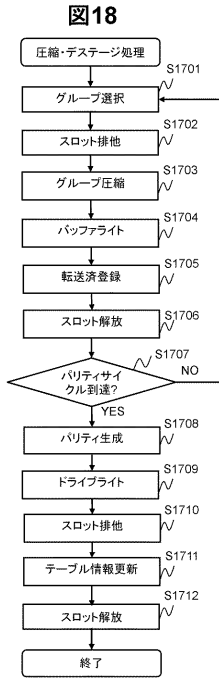
【 図 1 6 】



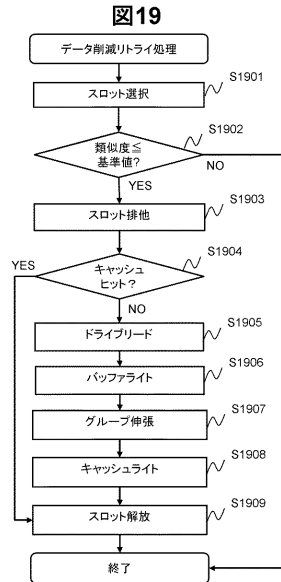
【 図 1 7 】



【 図 1 8 】



【 図 1 9 】



【 図 2 0 】

