

(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2008-9829

(P2008-9829A)

(43) 公開日 平成20年1月17日(2008.1.17)

(51) Int. Cl.

G06F 3/06 (2006.01)

F I

G06F 3/06 302A

テーマコード (参考)

5B065

審査請求 未請求 請求項の数 10 O L (全 23 頁)

(21) 出願番号 特願2006-181117 (P2006-181117)

(22) 出願日 平成18年6月30日 (2006.6.30)

(71) 出願人 000005223

富士通株式会社

神奈川県川崎市中原区上小田中4丁目1番
1号

(74) 代理人 100101856

弁理士 赤澤 日出夫

(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 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

Fターム(参考) 5B065 BA01 CA30 CC02 CH01

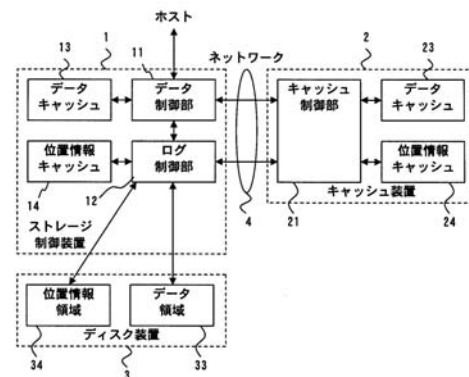
(54) 【発明の名称】 ストレージ制御プログラム、ストレージ制御装置、ストレージ制御方法

(57) 【要約】

【課題】 ログ書き込みに対するランダム read の性能を向上させるストレージ制御プログラム、ストレージ制御装置、ストレージ制御方法を提供する。

【解決手段】 ストレージ装置の制御をコンピュータに実行させるストレージ制御プログラムであって、外部からの write 要求に基づいて、前記ストレージ装置へデータのログ書き込みを行うデータ書き込みステップと、ネットワークを介して前記コンピュータに接続されたキャッシュ装置に、前記 write 要求により指定された論理位置と前記データ書き込みステップによるログ書き込みが行われた前記ストレージ装置内の物理位置との対応付けを位置情報として書き込むことができる位置情報制御ステップとをコンピュータに実行させる。

【選択図】 図1



【特許請求の範囲】

【請求項 1】

ストレージ装置の制御をコンピュータに実行させるストレージ制御プログラムであって、
外部からの `w r i t e` 要求に基づいて、前記ストレージ装置へデータのログ書き込みを行うデータ書き込みステップと、

ネットワークを介して前記コンピュータに接続されたキャッシュ装置に、前記 `w r i t e` 要求により指定された論理位置と前記データ書き込みステップによるログ書き込みが行われた前記ストレージ装置内の物理位置との対応付けを位置情報として書き込むことができる位置情報制御ステップと

をコンピュータに実行させるストレージ制御プログラム。

10

【請求項 2】

請求項 1 に記載のストレージ制御プログラムにおいて、

前記位置情報制御ステップは更に、前記キャッシュ装置に格納された前記位置情報から、外部からの `r e a d` 要求により指定された論理位置に対応する物理位置を読み出すことができ、

更に、前記位置情報制御ステップにより読み出された物理位置からデータを読み出すデータ読み出しステップと

をコンピュータに実行させることを特徴とするストレージ制御プログラム。

【請求項 3】

20

請求項 2 に記載のストレージ制御プログラムにおいて、

前記位置情報制御ステップは、前記コンピュータ内のキャッシュ、前記キャッシュ装置内のキャッシュを第 1 の階層的な記憶装置とし、該第 1 の階層的な記憶装置において前記位置情報を管理し、

前記キャッシュ装置内のキャッシュは前記コンピュータ内のキャッシュより大容量であることを特徴とするストレージ制御プログラム。

【請求項 4】

請求項 3 に記載のストレージ制御プログラムにおいて、

前記位置情報制御ステップは、前記コンピュータ内のキャッシュに最新の位置情報を書き込み、前記コンピュータ内のキャッシュに格納された位置情報のうち所定の使用頻度より低い位置情報を前記キャッシュ装置内のキャッシュへ移動させることを特徴とするストレージ制御プログラム。

30

【請求項 5】

請求項 3 または請求項 4 に記載のストレージ制御プログラムにおいて、

前記位置情報制御ステップは、前記第 1 の階層的な記憶装置に前記ストレージ装置を加えて第 1 の階層的な記憶装置とし、該第 1 の階層的な記憶装置において前記位置情報を管理し、

前記ストレージ装置は前記キャッシュ装置内のキャッシュより大容量であることを特徴とするストレージ制御プログラム。

【請求項 6】

40

請求項 5 に記載のストレージ制御プログラムにおいて、

前記位置情報制御ステップは、前記キャッシュ装置内のキャッシュに格納された位置情報のうち所定の使用頻度より低い位置情報を前記ストレージ装置へ移動させることを特徴とするストレージ制御プログラム。

【請求項 7】

請求項 2 乃至請求項 6 のいずれかに記載のストレージ制御プログラムにおいて、

前記データ書き込みステップは、前記コンピュータ内のキャッシュ、前記キャッシュ装置内のキャッシュ、前記ストレージ装置を前記第 2 の階層的な記録媒体とし、ライトスループにより該第 2 の階層的な記録媒体へデータを書き込み、

前記データ読み出しステップは、前記第 2 の階層的な記録媒体からデータを読み出し、

50

前記キャッシュ装置内のキャッシュは前記コンピュータ内のキャッシュより大容量であり、前記ストレージ装置は前記キャッシュ装置内のキャッシュより大容量であることを特徴とするストレージ制御プログラム。

【請求項 8】

ストレージ装置の制御を行うストレージ制御装置であって、

外部からの `w r i t e` 要求に基づいて、前記ストレージ装置へデータのログ書き込みを行うデータ制御部と、

ネットワークを介して前記ストレージ制御装置に接続されたキャッシュ装置に、前記 `w r i t e` 要求により指定された論理位置と前記データ制御部によるログ書き込みが行われた前記ストレージ装置内の物理位置との対応付けを位置情報として書き込むことができる位置情報制御部と

を備えるストレージ制御装置。

【請求項 9】

請求項 8 に記載のストレージ制御装置において、

前記位置情報制御部は更に、前記キャッシュ装置に格納された前記位置情報から、外部からの `r e a d` 要求により指定された論理位置に対応する物理位置を読み出すことができ

、
前記データ制御部は更に、前記位置情報制御部により読み出された物理位置からデータを読み出すことを特徴とするストレージ制御装置。

【請求項 10】

ストレージ装置の制御を行うストレージ制御方法であって、

外部からの `w r i t e` 要求に基づいて、前記ストレージ装置へデータのログ書き込みを行うデータ書き込みステップと、

ネットワークを介して前記コンピュータに接続されたキャッシュ装置に、前記 `w r i t e` 要求により指定された論理位置と前記データ書き込みステップによるログ書き込みが行われた前記ストレージ装置内の物理位置との対応付けを位置情報として書き込むことができる位置情報制御ステップと

を実行するストレージ制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ストレージ装置へのログ書き込みを行うストレージ制御プログラム、ストレージ制御装置、ストレージ制御方法に関するものである。

【背景技術】

【0002】

ハードディスクはシーケンシャルアクセスに比べてランダムアクセスが遅いという問題があり、ハードディスクを使ったストレージシステムは同様な問題を抱えている。この問題を解決するための方法がいくつか提案されている。

【0003】

第 1 の方法は、単純にディスクを大量に並べて並列に `I / O` を実行させることで性能を上げる方法である。この方法を用いて、大規模な `R A I D` (Redundant Arrays of Inexpensive Disks) システムを組むことにより、ランダムアクセスの `I O P S` (Input Output Per Second) と信頼性の双方を向上させることができる。

【0004】

第 2 の方法は、キャッシュを利用する方法である。この方法は、`r e a d` にキャッシュを利用するライトスルーキャッシュ、`r e a d` と `w r i t e` の双方でキャッシュを利用するライトバックキャッシュに分類できる。これらの方法を用いて、キャッシュの容量を大規模にしてヒット率を上げると、飛躍的にランダムアクセスの `I O P S` を向上させることができる。

【0005】

10

20

30

40

50

ただし、ライトバックキャッシュは、ディスクに書き込む前のデータをメモリに保持する必要があるため、信頼性が低下する。一方で、信頼性を確保しようとキャッシュを多重化すれば、キャッシュメモリのコストが倍増してしまう。

【0006】

次に、ランダムw r i t eについて説明する。

【0007】

ランダムw r i t eにおいて、ログ書き込みと呼ばれる方法がある。ログ書き込みは、ランダムw r i t eの論理位置（アドレス）の情報とデータをシーケンシャルに書き込んでいく。読み込みは、論理位置と物理位置の対応付けを示す位置情報テーブルの情報をすることにより、データを読み込んで再現させる。これにより、ランダムw r i t eは、シーケンシャルアクセス並のスループットが出せるようになる。ランダムw r i t eに強いログ書き込みと、ランダムr e a dに強いライトスルーキャッシュを組み合わせれば、ランダムアクセス性能・信頼性・コストについて良い結果を出すことができる。

10

【0008】

なお、本発明の関連ある従来技術として、キャッシュ装置、ディスク装置、制御装置をネットワーク上に分散配置することにより、キャッシュメモリ領域を増やすストレージシステムがある（例えば、特許文献1、特許文献2、特許文献3参照）。

【特許文献1】W O 2 0 0 3 / 0 6 5 1 9 5

【特許文献2】W O 2 0 0 3 / 0 7 5 1 4 7

【特許文献3】W O 2 0 0 4 / 0 2 7 6 2 5

20

【発明の開示】

【発明が解決しようとする課題】

【0009】

しかしながら、ログ書き込みとライトスルーキャッシュの組み合わせは効果が非常に高いが、ログ方式を用いた場合に特有の問題がある。特に、ログ書き込みのストレージシステムは、ホストが利用する書き込み位置（論理位置）と下位デバイス（ディスク装置）の書き込み位置（物理位置）との対応付けを示す位置情報テーブルを管理する必要があり、位置情報テーブルの検索に時間がかかることから、キャッシュミスの時の読み込みの性能が劣化するという問題がある。また、この位置情報テーブルはセクタとセクタを1対1で対応させる必要があり、おおむねボリュームサイズの1 / 6 4 ~ 1 / 5 1 2程度の空間が必要になる。

30

【0010】

本発明は上述した問題点を解決するためになされたものであり、ログ書き込みに対するランダムr e a dの性能を向上させるストレージ制御プログラム、ストレージ制御装置、ストレージ制御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】

上述した課題を解決するため、本発明は、ストレージ装置の制御をコンピュータに実行させるストレージ制御プログラムであって、外部からのw r i t e要求に基づいて、前記ストレージ装置へデータのログ書き込みを行うデータ書き込みステップと、ネットワークを介して前記コンピュータに接続されたキャッシュ装置に、前記w r i t e要求により指定された論理位置と前記データ書き込みステップによるログ書き込みが行われた前記ストレージ装置内の物理位置との対応付けを位置情報として書き込むことができる位置情報制御ステップとをコンピュータに実行させるものである。

40

【0012】

また、本発明に係るストレージ制御プログラムにおいて、前記位置情報制御ステップは更に、外部からのr e a d要求に基づいて、前記キャッシュ装置に記憶された前記位置情報から前記r e a d要求により指定された論理位置に対応する物理位置を読み出すことができ、更に、前記位置情報制御ステップにより読み出された物理位置からデータを読み出すデータ読み出しステップとをコンピュータに実行させることを特徴とする。

50

【 0 0 1 3 】

また、本発明に係るストレージ制御プログラムにおいて、前記位置情報制御ステップは、前記コンピュータ内のキャッシュ、前記キャッシュ装置内のキャッシュを第1の階層的な記憶装置とし、該第1の階層的な記憶装置において前記位置情報を管理し、前記キャッシュ装置内のキャッシュは前記コンピュータ内のキャッシュより大容量であることを特徴とする。

【 0 0 1 4 】

また、本発明に係るストレージ制御プログラムにおいて、前記位置情報制御ステップは、前記コンピュータ内のキャッシュに最新の位置情報を書き込み、前記コンピュータ内のキャッシュに格納された位置情報のうち所定の使用頻度より低い位置情報を前記キャッシュ装置内のキャッシュへ移動させることを特徴とする。

10

【 0 0 1 5 】

また、本発明に係るストレージ制御プログラムにおいて、前記位置情報制御ステップは、前記第1の階層的な記憶装置に前記ストレージ装置を加えて第1の階層的な記憶装置とし、該第1の階層的な記憶装置において前記位置情報を管理し、前記ストレージ装置は前記キャッシュ装置内のキャッシュより大容量であることを特徴とする。

【 0 0 1 6 】

また、本発明に係るストレージ制御プログラムにおいて、前記位置情報制御ステップは、前記キャッシュ装置内のキャッシュに格納された位置情報のうち所定の使用頻度より低い位置情報を前記ストレージ装置へ移動させることを特徴とする。

20

【 0 0 1 7 】

また、本発明に係るストレージ制御プログラムにおいて、前記データ書き込みステップは、前記コンピュータ内のキャッシュ、前記キャッシュ装置内のキャッシュ、前記ストレージ装置を前記第2の階層的な記録媒体とし、ライトスルーにより該第2の階層的な記録媒体へデータを書き込み、前記データ読み出しステップは、前記第2の階層的な記録媒体からデータを読み出し、前記キャッシュ装置内のキャッシュは前記コンピュータ内のキャッシュより大容量であり、前記ストレージ装置は前記キャッシュ装置内のキャッシュより大容量であることを特徴とする。

【 0 0 1 8 】

また、本発明は、ストレージ装置の制御を行うストレージ制御装置であって、外部からの `w r i t e` 要求に基づいて、前記ストレージ装置へデータのログ書き込みを行うデータ制御部と、ネットワークを介して前記ストレージ制御装置に接続されたキャッシュ装置に、前記 `w r i t e` 要求により指定された論理位置と前記データ制御部によるログ書き込みが行われた前記ストレージ装置内の物理位置との対応付けを位置情報として書き込むことができる位置情報制御部とを備えたものである。

30

【 0 0 1 9 】

また、本発明に係るストレージ制御装置において、前記位置情報制御部は更に、外部からの `r e a d` 要求に基づいて、前記キャッシュ装置に記憶された前記位置情報から前記 `r e a d` 要求により指定された論理位置に対応する物理位置を読み出すことができ、前記データ制御部は更に、前記位置情報制御部により読み出された物理位置からデータを読み出すことを特徴とする。

40

【 0 0 2 0 】

また、本発明に係るストレージ制御装置において、前記位置情報制御部は、前記ストレージ制御装置内のキャッシュ、前記キャッシュ装置内のキャッシュを第1の階層的な記憶装置とし、該第1の階層的な記憶装置において前記位置情報を管理し、前記キャッシュ装置内のキャッシュは前記ストレージ制御装置内のキャッシュより大容量であることを特徴とする。

【 0 0 2 1 】

また、本発明に係るストレージ制御装置において、前記位置情報制御部は、前記コンピュータ内のキャッシュに最新の位置情報を書き込み、前記ストレージ制御装置内のキャッ

50

シュに格納された位置情報のうち所定の使用頻度より低い位置情報を前記キャッシュ装置内のキャッシュへ移動させることを特徴とする。

【0022】

また、本発明に係るストレージ制御装置において、前記位置情報制御部は、前記第1の階層的な記憶装置に前記ストレージ装置を加えて第1の階層的な記憶装置とし、該第1の階層的な記憶装置において前記位置情報を管理し、前記ストレージ装置は前記キャッシュ装置内のキャッシュより大容量であることを特徴とする。

【0023】

また、本発明に係るストレージ制御装置において、前記位置情報制御部は、前記キャッシュ装置内のキャッシュに格納された位置情報のうち所定の使用頻度より低い位置情報を前記ストレージ装置へ移動させることを特徴とする。

10

【0024】

また、本発明に係るストレージ制御装置において、前記データ制御部は、前記ストレージ制御装置内のキャッシュ、前記キャッシュ装置内のキャッシュ、前記ストレージ装置を前記第2の階層的な記録媒体として管理し、ライトスルーにより該第2の階層的な記録媒体に対するデータの書き込み及び読み出しを行い、前記ストレージ制御装置内のキャッシュは前記キャッシュ装置内のキャッシュより高速であり、前記キャッシュ装置内のキャッシュは前記ストレージ装置より高速であることを特徴とする。

【0025】

また、本発明は、ストレージ装置の制御を行うストレージ制御方法であって、外部からのwrite要求に基づいて、前記ストレージ装置へデータのログ書き込みを行うデータ書き込みステップと、ネットワークを介して前記コンピュータに接続されたキャッシュ装置に、前記write要求により指定された論理位置と前記データ書き込みステップによるログ書き込みが行われた前記ストレージ装置内の物理位置との対応付けを位置情報として書き込むことができる位置情報制御ステップとを実行するものである。

20

【発明の効果】

【0026】

本発明によれば、ログ書き込みに対するランダムreadの性能、特にキャッシュミス時の性能を向上させることができる。

【発明を実施するための最良の形態】

30

【0027】

以下、本発明の実施の形態について図面を参照しつつ説明する。

【0028】

実施の形態1.

本実施の形態では、上述したログ書き込みとライトスルーキャッシュの組み合わせを用いるストレージ制御装置について説明する。

【0029】

まず、本実施の形態に係るストレージ制御装置を用いたストレージシステムの構成について説明する。

【0030】

40

図1は、本実施の形態に係るストレージシステムの構成の一例を示すブロック図である。このストレージシステムは、ストレージ制御装置1、キャッシュ装置2、ディスク装置3（ストレージ装置）、ネットワーク4を備える。ストレージ制御装置1は、外部のホストとディスク装置3にそれぞれ接続される。また、ストレージ制御装置1は、ネットワーク4を介してキャッシュ装置2に接続される。

【0031】

ストレージ制御装置1は、データ制御部11、ログ制御部12（位置情報制御部）、データキャッシュ13、位置情報キャッシュ14を備える。データキャッシュ13は、データ制御部11の内部に設けられても良い。位置情報キャッシュ14は、ログ制御部12の内部に設けられても良い。

50

【 0 0 3 2 】

キャッシュ装置 2 は、キャッシュ制御部 2 1、データキャッシュ 2 3、位置情報キャッシュ 2 4 を備える。キャッシュ制御部 2 1 は、データキャッシュ 2 3 及び位置情報キャッシュ 2 4 の制御を行う。データキャッシュ 2 3 は、データキャッシュ 1 3 より大容量のメモリで構成され、位置情報キャッシュ 2 4 は、位置情報キャッシュ 1 4 より大容量のメモリで構成される。なお、データキャッシュ 2 3 と位置情報キャッシュ 2 4 は、1 つのメモリとして構成されても良い。

【 0 0 3 3 】

ディスク装置 3 は、ログ書き込みによるデータを格納するデータ領域 3 3、位置情報を格納する位置情報領域 3 4 を備える。

10

【 0 0 3 4 】

また、データ制御部 1 1 は、ストレージ制御装置 1 におけるデータキャッシュ 1 3 及びキャッシュ装置 2 におけるデータキャッシュ 2 3 の制御を行う。また、データ制御部 1 1 は、ホストからの `read` 要求または `write` 要求に従ってログ制御部 1 2 への `read` 要求または `write` 要求を行い、結果をホストへ返す。ログ制御部 1 2 は、ストレージ制御装置 1 における位置情報キャッシュ 1 4 及びキャッシュ装置 2 における位置情報キャッシュ 2 4 の制御を行う。また、ログ制御部 1 2 は、データ制御部 1 1 からの `write` 要求のデータをログ形式に変換し、ディスク装置 3 のデータ領域 3 3 に書き込む。また、ログ制御部 1 2 は、データ制御部 1 1 からの `read` 要求で要求されたデータを、ディスク装置 3 のデータ領域 3 3 からデータ制御部 1 1 へ読み出す。

20

【 0 0 3 5 】

また、データ制御部 1 1 は、データキャッシュ 1 3、2 3、データ領域 3 3 を階層的な記憶装置として管理し、それぞれにデータを格納する。データへのアクセスにおいて、データキャッシュ 1 3 は、最も小容量、最も高速であり、最も使用頻度の高いデータが格納される。データキャッシュ 2 3 は、データキャッシュ 1 3 より大容量、低速であり、データキャッシュ 1 3 より使用頻度の低いデータが格納される。データ領域 3 3 は、最も大容量、最も低速であり、最も使用頻度の低いデータが格納される。

【 0 0 3 6 】

同様に、ログ制御部 1 2 は、位置情報キャッシュ 1 4、2 4、位置情報領域 3 4 を階層的な記憶装置として管理し、それぞれに位置情報テーブルを格納する。位置情報テーブルへのアクセスにおいて、位置情報キャッシュ 1 4 は、最も小容量、最も高速であり、最も使用頻度の高い位置情報が格納される。位置情報キャッシュ 2 4 は、位置情報キャッシュ 1 4 より大容量、低速であり、位置情報キャッシュ 1 4 より使用頻度の低い位置情報が格納される。位置情報領域 3 4 は、最も大容量、最も低速であり、最も使用頻度の低い位置情報が格納される。

30

【 0 0 3 7 】

次に、位置情報テーブルについて説明する。

【 0 0 3 8 】

位置情報テーブルの管理には、ファイルシステムのブロックマップなどに使われる方法と同じものを使う。図 2 は、本実施の形態に係る位置情報テーブルの構成の一例を示すブロック図である。この図のように、位置情報テーブルは、ホストが利用する論理位置とデータ領域 3 3 における物理位置との対応付けを示す位置情報を、2 ~ 3 レベルの間接テーブルとして表したものである。ログ制御部 1 2 は、この図の左端のテーブルに示された論理位置から、この図の右端に示された物理位置を検索する。このような位置情報テーブルを用いることにより、データ領域 3 3 におけるデータの部分的な更新や参照を高速にすることができ、大きなボリュームを管理することができる。

40

【 0 0 3 9 】

次に、`write` 処理について説明する。

【 0 0 4 0 】

ライトスルーキャッシュを用いるため、`write` 処理において、データ制御部 1 1 は

50

、データキャッシュ 13 , 23 への書き込みだけでなく、ディスク装置 3 への書き込みを行う。図 3 は、本実施の形態に係るストレージ制御装置による write 処理の動作の一例を示すフローチャートである。

【0041】

ホストから write 要求を受けたデータ制御部 11 は、ログ制御部 12 へ write 要求を送る (S11)。次に、ログ制御部 12 は、データ制御部 11 から受けた write 要求のデータに対するヘッダとトレイラを作成し (S12)、そのヘッダ、データ、トレイラをログ形式としてディスク装置 3 のデータ領域 33 へ書き込む (S13)。

【0042】

次に、ログ制御部 12 は、位置情報テーブル更新処理を行い (S21)、write 完了をデータ制御部 11 に通知する (S22)。次に、ログ制御部 12 から write 完了の通知を受け取ったデータ制御部 11 は、データキャッシュ 13 , 23 の更新を行い (S23)、このフローは終了する。

【0043】

次に、read 処理について説明する。

【0044】

read 処理において、ストレージ制御装置 1 は、データキャッシュ 13 , 23 においてキャッシュヒットすれば、そのキャッシュのデータをホストへ返し、キャッシュミスすれば、ディスク装置 3 のデータをホストへ返す。図 4 は、本実施の形態に係るストレージ制御装置による read 処理の動作の一例を示すフローチャートである。

【0045】

ホストから read 要求を受けたデータ制御部 11 は、ホストから要求されたデータがデータキャッシュ 13 , 23 においてヒットしたか否かの判断を行う (S31)。キャッシュヒットした場合 (S31, yes)、データ制御部 11 は、ヒットしたデータキャッシュからホストへ要求されたデータの読み出しを行い (S32)、ヒットしたデータキャッシュにおける LRU の更新を行い (S35)、このフローは終了する。データ制御部 11 は、データキャッシュ 13 における LRU を用いて、所定の使用頻度より低いデータを、データキャッシュ 13 からデータキャッシュ 23 へ移動させる。同様に、キャッシュ制御部 21 は、データキャッシュ 23 における LRU を用いて、所定の使用頻度より低いデータを、データキャッシュ 23 からデータ領域 33 へ移動させる。

【0046】

一方、キャッシュミスした場合 (S31, no)、データ制御部 11 は、ログ制御部 12 へ read 要求を送る (S41)。次に、ログ制御部 12 は、位置情報キャッシュ 24 の位置情報テーブルにおいて、指示されたデータの論理位置からディスク装置 3 上の物理位置を検索する検索処理を行い (S42)、検索された物理位置のデータをホストへ読み出し (S43)、このフローは終了する。

【0047】

次に、write 処理における位置情報テーブル検索処理 (S42) 及び read 処理における位置情報テーブル更新処理 (S21) の詳細について説明する。

【0048】

本実施の形態においては、ログ制御部 11 が、位置情報キャッシュ 14 , 24、位置情報領域 34 における全ての位置情報テーブルの管理を行う。また、ログ制御部 12 は、キャッシュ制御部 21 との通信を行うことにより、位置情報キャッシュ 24 の操作を行う。この通信において、ログ制御部 12 は、特許文献 3 のようなキャッシュプロトコルを用いる。

【0049】

図 5 は、本実施の形態に係るキャッシュプロトコルにおけるメッセージの構造の一例を示す構造図である。データ制御部 11 またはログ制御部 12 からキャッシュ制御部 21 へ送られるリクエストメッセージは、コマンド、ID、リクエスト (req) データからなる。リクエストメッセージの結果としてキャッシュ制御部 2 からデータ制御部 11 または

10

20

30

40

50

ログ制御部 12 へ送られるレスポンスメッセージは、コマンド、ID、レスポンス (res) データからなる。

【0050】

図 6 は、本実施の形態に係るキャッシュプロトコルにおけるデータの構造の一例を示す構造図である。割り当てリクエストのデータは、割り当てたいサイズからなる。これに対する割り当てレスポンスのデータは、割り当てられた位置のポインタとサイズからなる。開放リクエストのデータは、開放したい位置のポインタからなる。これに対する開放レスポンスのデータは、開放の結果が OK であるか NG であるかを示す。read リクエストのデータは、read したい位置のポインタとサイズからなる。これに対する read レスポンスのデータは、read したサイズとデータからなる。write リクエストのデータは、write したい位置のポインタとサイズとデータからなる。これに対する write レスポンスのデータは、write の結果が OK であるか NG であるかを示す。

10

【0051】

図 7 は、本実施の形態に係る位置情報テーブル検索処理及び位置情報テーブル更新処理の動作の一例を示すフローチャートである。まず、ログ制御部 12 は、データ制御部 11 からの write 要求 (更新処理時) または read 要求 (検索処理時) の対象が、最も高速なローカルのキャッシュである位置情報キャッシュ 14 の位置情報テーブルにおいてキャッシュヒットしたか否かの判断を行う (S51)。

【0052】

キャッシュヒットした場合 (S51, yes)、処理 S66 へ移行する。一方、キャッシュミスした場合 (S51, no)、ログ制御部 12 は、位置情報キャッシュ 24 においてキャッシュヒットしたか否かの判断を行う (S52)。キャッシュヒットした場合 (S52, yes)、ログ制御部 12 は、位置情報キャッシュ 24 から該当する位置情報を読み出し (S53)、処理 S61 へ移行する。ここで、ログ制御部 12 は、read リクエストをキャッシュ制御部 12 へ送信し、キャッシュ制御部 12 は、read レスポンスをログ制御部 12 へ送信する。一方、キャッシュミスした場合 (S52, no)、ログ制御部 12 は、位置情報領域 34 から該当する位置情報を読み出し (S54)、処理 61 へ移行する。

20

【0053】

次に、ログ制御部 12 は、処理 S53 または S54 で読み出した位置情報を位置情報キャッシュ 14 に追加で書き込み (S61)、位置情報キャッシュ 14 の LRU (Least Recently Used) を更新する (S62)。次に、ログ制御部 12 は、データ制御部 11 は、位置情報キャッシュ 14 の LRU を用いて所定の時刻より古い位置情報を位置情報キャッシュ 24 へ書き込む (S63)。ここで、ログ制御部 12 は、割り当てリクエストをキャッシュ制御部 12 へ送信し、キャッシュ制御部 12 は、割り当てレスポンスをログ制御部 12 へ送信する。

30

【0054】

次に、キャッシュ制御部 21 は、位置情報キャッシュ 24 の LRU を更新する (S64)。次に、ログ制御部 12 は、位置情報キャッシュ 24 の LRU を用いて所定の時刻より古い位置情報を位置情報領域 34 へ書き込む (S65)。ここで、ログ制御部 12 は、read リクエストをキャッシュ制御部 12 へ送信し、キャッシュ制御部 12 は、read レスポンスをログ制御部 12 へ送信し、更に、ログ制御部 12 は、開放リクエストをキャッシュ制御部 12 へ送信し、キャッシュ制御部 12 は、開放レスポンスをログ制御部 12 へ送信する。次に、ログ制御部 12 は、位置情報キャッシュ 14 において検索を行い、更に位置情報テーブル更新処理時には更新を行い (S66)、このフローは終了する。

40

【0055】

本実施の形態によれば、ストレージ制御装置 1 におけるログ制御部 12 が、ストレージ制御装置 1 における位置情報キャッシュ 14、キャッシュ装置 2 における位置情報キャッシュ 24、ディスク装置 3 における位置情報領域 34 の位置情報テーブルを階層的に管理することにより、ランダム read におけるディスク装置 3 上の物理位置の検索の性能を

50

向上させることができ、ログ書き込みに対するランダム read の性能を向上させることができる。

【 0 0 5 6 】

実施の形態 2 .

まず、本実施の形態に係るストレージ制御装置を用いたストレージシステムの構成について説明する。本実施の形態におけるストレージシステムの構成は、実施の形態 1 と同様であるが、データ制御部 1 1 とログ制御部 1 2 の動作が異なる。

【 0 0 5 7 】

次に、本実施の形態に係るストレージ制御装置の動作について説明する。write 処理及び read 処理の動作は、実施の形態 1 と同様であるが、write 処理における位置情報テーブル検索処理 (S 4 2) 及び read 処理における位置情報テーブル更新処理 (S 2 3) の内容が異なる。

10

【 0 0 5 8 】

次に、write 処理における位置情報テーブルの検索処理 (S 4 2) 及び read 処理における更新処理 (S 2 1) の詳細について説明する。

【 0 0 5 9 】

本実施の形態においては、ログ制御部 1 1 が位置情報キャッシュ 1 4 の管理を行い、キャッシュ制御部 2 1 が位置情報キャッシュ 2 4 及び位置情報領域 3 4 の管理を行う。また、ログ制御部 1 2 とキャッシュ制御部 2 1 は、互いに通信を行うことにより、位置情報キャッシュ 2 4 と位置情報領域 3 4 の操作を行う。この通信において、ログ制御部 1 2 とキャッシュ制御部 2 1 は、実施の形態 1 のキャッシュプロトコルにテーブル操作のためのプロトコルを追加したキャッシュプロトコルを用いる。

20

【 0 0 6 0 】

図 8 は、本実施の形態に係るキャッシュプロトコルにおけるデータの構造の一例を示す構造図である。検索リクエストのデータは、対象のデータの論理位置を示すオフセット、対象の長さからなる。検索レスポンスのデータは、検索結果のブロックの数を示すサイズ、ブロック毎の物理位置を表すオフセット、ブロック毎の長さからなる。更新リクエストのデータは、対象のデータのブロックの数を示すサイズ、ブロック毎の論理位置を表すオフセット、ブロック毎の長さからなる。更新レスポンスのデータは、更新リクエストに対する処理の結果が OK であるか NG であるかを示す。

30

【 0 0 6 1 】

テーブル書き出しリクエストのデータは、更新リクエストと同様、対象のデータのブロックの数を示すサイズ、ブロック毎の論理位置を表すオフセット、ブロック毎の長さからなる。テーブル書き出しレスポンスのデータは、更新レスポンスと同様、テーブル書き出しリクエストに対する処理の結果が OK であるか NG であるかを示す。テーブル読み込みリクエストのデータは、検索リクエストと同様、対象のデータの論理位置を示すオフセット、対象の長さからなる。テーブル読み込みレスポンスのデータは、検索レスポンスと同様、読み込み結果のブロックの数を示すサイズ、ブロック毎の物理位置を表すオフセット、ブロック毎の長さからなる。

【 0 0 6 2 】

40

図 9 は、本実施の形態に係る位置情報テーブル検索処理及び位置情報テーブル更新処理の動作の一例を示すフローチャートである。図中の点線の枠内はキャッシュ制御部 2 1 の動作を示し、点線の枠外はログ制御部 1 2 の動作を示す。まず、ログ制御部 1 2 は、データ制御部 1 1 により write 要求 (更新処理時) または read 要求 (検索処理時) の対象の位置情報が、位置情報キャッシュ 1 4 においてキャッシュヒットしたか否かの判断を行う (S 7 1) 。キャッシュヒットした場合 (S 7 1 , y e s) 、処理 S 9 4 へ移行する。一方、キャッシュミスした場合 (S 7 1 , y e s) 、ログ制御部 1 2 は、位置情報キャッシュ 2 4 において対象の位置情報を検索することを指示する検索リクエストをキャッシュ制御部 2 1 へ送信する (S 7 2) 。

【 0 0 6 3 】

50

次に、キャッシュ制御部 21 は、検索リクエストに従って、位置情報キャッシュ 24 において対象の位置情報を検索し (S 73)、位置情報キャッシュ 24 において対象の位置情報が見つかったか否かの判断を行う (S 74)。対象の位置情報が見つかった場合 (S 74, yes)、処理 S 94 へ移行する。一方、対象の位置情報が見つからなかった場合 (S 74, no)、位置情報領域 34 から対象の位置情報を読み出すことを指示するテーブル読み出しリクエストをログ制御部 12 へ送信する (S 75)。次に、ログ制御部 12 は、テーブル読み出しリクエストに従って、位置情報領域 34 から対象の位置情報を読み出し (S 76)、読み込んだ位置情報をテーブル読み込みレスポンスとしてキャッシュ制御部 21 へ送信する (S 77)。

【0064】

次に、キャッシュ制御部 21 は、位置情報キャッシュ 24 の LRU の更新を行い (S 81)、位置情報キャッシュ 24 の LRU を用いて、位置情報キャッシュ 24 のうち所定の時刻より古い位置情報を位置情報領域 34 へ書き出すことを指示するテーブル書き出しリクエストをログ制御部 12 へ送信する (S 82)。次に、ログ制御部 12 は、テーブル書き出しリクエストに従って、古い位置情報を位置情報領域 34 へ書き込み (S 83)、書き込み完了を示すテーブル書き出しレスポンスを送信する (S 84)。次に、キャッシュ制御部 21 は、得られた対象の位置情報を検索レスポンスとしてログ制御部 12 へ送信する (S 85)。

【0065】

次に、ログ制御部 12 は、検索レスポンスから得られた対象の位置情報を位置情報キャッシュ 14 に追加し (S 91)、位置情報キャッシュ 14 の LRU の更新を行い (S 92)、位置情報キャッシュ 14 の LRU を用いて位置情報キャッシュ 14 のうち所定の時刻より古い位置情報を削除する (S 93)。次に、ログ制御部 12 は、位置情報キャッシュ 14 において検索を行い、更に位置情報テーブル更新処理時には更新を行い (S 94)、このフローは終了する。

【0066】

本実施の形態によれば、ログ制御部 12 が、位置情報キャッシュ 14 及びキャッシュ制御部 21 を階層的に管理し、更に、キャッシュ制御部 21 が、位置情報キャッシュ 14 及び位置情報領域 34 を階層的に管理することにより、ディスク装置 3 における位置情報領域 34 を階層的に用いることにより、ランダム read におけるディスク装置 3 上の物理位置の検索の性能を向上させることができ、ログ書き込みに対するランダム read の性能を向上させることができる。

【0067】

なお、上述した実施の形態において、位置情報テーブルは、位置情報キャッシュ 14、24、位置情報領域 34 に書き込まれるとしたが、位置情報キャッシュ 24 が十分な容量を持つ場合、位置情報領域 24 を必要としない。

【0068】

なお、データ書き込みステップは、実施の形態における処理 S 11, S 12, S 13, S 22, S 23 に対応する。また、データ読み出しステップは、実施の形態における処理 S 31, S 41, S 32, S 33, S 43 に対応する。また、位置情報制御ステップは、実施の形態における処理 S 21, S 42 に対応する。

【0069】

また、本実施の形態に係るストレージ制御装置は、ストレージシステムに容易に適用することができ、ストレージシステムの性能をより高めることができる。

【0070】

更に、ストレージ制御装置を構成するコンピュータにおいて上述した各ステップを実行させるプログラムを、ストレージ制御プログラムとして提供することができる。上述したプログラムは、コンピュータにより読取り可能な記録媒体に記憶させることによって、ストレージ制御装置を構成するコンピュータに実行させることが可能となる。ここで、上記コンピュータにより読取り可能な記録媒体としては、ROM や RAM 等のコンピュータに

10

20

30

40

50

内部実装される内部記憶装置、ＣＤ－ＲＯＭやフレキシブルディスク、ＤＶＤディスク、光磁気ディスク、ＩＣカード等の可搬型記憶媒体や、コンピュータプログラムを保持するデータベース、或いは、他のコンピュータ並びにそのデータベースや、更に回線上の伝送媒体をも含むものである。

【００７１】

（付記１） ストレージ装置の制御をコンピュータに実行させるストレージ制御プログラムであって、

外部からの *w r i t e* 要求に基づいて、前記ストレージ装置へデータのログ書き込みを行うデータ書き込みステップと、

ネットワークを介して前記コンピュータに接続されたキャッシュ装置に、前記 *w r i t e* 要求により指定された論理位置と前記データ書き込みステップによるログ書き込みが行われた前記ストレージ装置内の物理位置との対応付けを位置情報として書き込むことができる位置情報制御ステップと

をコンピュータに実行させるストレージ制御プログラム。

（付記２） 付記１に記載のストレージ制御プログラムにおいて、

前記位置情報制御ステップは更に、前記キャッシュ装置に格納された前記位置情報から、外部からの *r e a d* 要求により指定された論理位置に対応する物理位置を読み出すことができ、

更に、前記位置情報制御ステップにより読み出された物理位置からデータを読み出すデータ読み出しステップと

をコンピュータに実行させることを特徴とするストレージ制御プログラム。

（付記３） 付記２に記載のストレージ制御プログラムにおいて、

前記位置情報制御ステップは、前記コンピュータ内のキャッシュ、前記キャッシュ装置内のキャッシュを第１の階層的な記憶装置とし、該第１の階層的な記憶装置において前記位置情報を管理し、

前記キャッシュ装置内のキャッシュは前記コンピュータ内のキャッシュより大容量であることを特徴とするストレージ制御プログラム。

（付記４） 付記３に記載のストレージ制御プログラムにおいて、

前記位置情報制御ステップは、前記コンピュータ内のキャッシュに最新の位置情報を書き込み、前記コンピュータ内のキャッシュに格納された位置情報のうち所定の使用頻度より低い位置情報を前記キャッシュ装置内のキャッシュへ移動させることを特徴とするストレージ制御プログラム。

（付記５） 付記３または付記４に記載のストレージ制御プログラムにおいて、

前記位置情報制御ステップは、前記第１の階層的な記憶装置に前記ストレージ装置を加えて第１の階層的な記憶装置とし、該第１の階層的な記憶装置において前記位置情報を管理し、

前記ストレージ装置は前記キャッシュ装置内のキャッシュより大容量であることを特徴とするストレージ制御プログラム。

（付記６） 付記５に記載のストレージ制御プログラムにおいて、

前記位置情報制御ステップは、前記キャッシュ装置内のキャッシュに格納された位置情報のうち所定の使用頻度より低い位置情報を前記ストレージ装置へ移動させることを特徴とするストレージ制御プログラム。

（付記７） 付記２乃至付記６のいずれかに記載のストレージ制御プログラムにおいて、

前記データ書き込みステップは、前記コンピュータ内のキャッシュ、前記キャッシュ装置内のキャッシュ、前記ストレージ装置を前記第２の階層的な記録媒体とし、ライトスルーにより該第２の階層的な記録媒体へデータを書き込み、

前記データ読み出しステップは、前記第２の階層的な記録媒体からデータを読み出し、

前記キャッシュ装置内のキャッシュは前記コンピュータ内のキャッシュより大容量であり、前記ストレージ装置は前記キャッシュ装置内のキャッシュより大容量であることを特徴とするストレージ制御プログラム。

10

20

30

40

50

(付記 8) ストレージ装置の制御を行うストレージ制御装置であって、

外部からの `w r i t e` 要求に基づいて、前記ストレージ装置へデータのログ書き込みを行うデータ制御部と、

ネットワークを介して前記ストレージ制御装置に接続されたキャッシュ装置に、前記 `w r i t e` 要求により指定された論理位置と前記データ制御部によるログ書き込みが行われた前記ストレージ装置内の物理位置との対応付けを位置情報として書き込むことができる位置情報制御部と

を備えるストレージ制御装置。

(付記 9) 付記 8 に記載のストレージ制御装置において、

前記位置情報制御部は更に、前記キャッシュ装置に格納された前記位置情報から、外部からの `r e a d` 要求により指定された論理位置に対応する物理位置を読み出すことができ

、

前記データ制御部は更に、前記位置情報制御部により読み出された物理位置からデータを読み出すことを特徴とするストレージ制御装置。

(付記 10) 付記 9 に記載のストレージ制御装置において、

前記位置情報制御部は、前記ストレージ制御装置内のキャッシュ、前記キャッシュ装置内のキャッシュを第 1 の階層的な記憶装置とし、該第 1 の階層的な記憶装置において前記位置情報を管理し、

前記キャッシュ装置内のキャッシュは前記ストレージ制御装置内のキャッシュより大容量であることを特徴とするストレージ制御装置。

(付記 11) 付記 10 に記載のストレージ制御装置において、

前記位置情報制御部は、前記コンピュータ内のキャッシュに最新の位置情報を書き込み、前記ストレージ制御装置内のキャッシュに格納された位置情報のうち所定の使用頻度より低い位置情報を前記キャッシュ装置内のキャッシュへ移動させることを特徴とするストレージ制御装置。

(付記 12) 付記 10 または付記 11 に記載のストレージ制御装置において、

前記位置情報制御部は、前記第 1 の階層的な記憶装置に前記ストレージ装置を加えて第 1 の階層的な記憶装置とし、該第 1 の階層的な記憶装置において前記位置情報を管理し、

前記ストレージ装置は前記キャッシュ装置内のキャッシュより大容量であることを特徴とするストレージ制御装置。

(付記 13) 付記 12 に記載のストレージ制御装置において、

前記位置情報制御部は、前記キャッシュ装置内のキャッシュに格納された位置情報のうち所定の使用頻度より低い位置情報を前記ストレージ装置へ移動させることを特徴とするストレージ制御装置。

(付記 14) 付記 9 乃至付記 13 のいずれかに記載のストレージ制御装置において、

前記データ制御部は、前記ストレージ制御装置内のキャッシュ、前記キャッシュ装置内のキャッシュ、前記ストレージ装置を前記第 2 の階層的な記録媒体として管理し、ライトスルーにより該第 2 の階層的な記録媒体に対するデータの書き込み及び読み出しを行い、

前記キャッシュ装置内のキャッシュは前記ストレージ制御装置内のキャッシュより大容量であり、前記ストレージ装置は前記キャッシュ装置内のキャッシュより大容量であることを特徴とするストレージ制御装置。

(付記 15) ストレージ装置の制御を行うストレージ制御方法であって、

外部からの `w r i t e` 要求に基づいて、前記ストレージ装置へデータのログ書き込みを行うデータ書き込みステップと、

ネットワークを介して前記コンピュータに接続されたキャッシュ装置に、前記 `w r i t e` 要求により指定された論理位置と前記データ書き込みステップによるログ書き込みが行われた前記ストレージ装置内の物理位置との対応付けを位置情報として書き込むことができる位置情報制御ステップと

を実行するストレージ制御方法。

(付記 16) 付記 15 に記載のストレージ制御方法において、

10

20

30

40

50

前記位置情報制御ステップは更に、前記キャッシュ装置に格納された前記位置情報から、外部からの read 要求により指定された論理位置に対応する物理位置を読み出すことができ、

更に、前記位置情報制御ステップにより読み出された物理位置からデータを読み出すデータ読み出しステップと

を実行することを特徴とするストレージ制御方法。

(付記 17) 付記 16 に記載のストレージ制御方法において、

前記位置情報制御ステップは、前記コンピュータ内のキャッシュ、前記キャッシュ装置内のキャッシュを第 1 の階層的な記憶装置とし、該第 1 の階層的な記憶装置において前記位置情報を管理し、

10

前記キャッシュ装置内のキャッシュは前記コンピュータ内のキャッシュより大容量であることを特徴とするストレージ制御方法。

(付記 18) 付記 17 に記載のストレージ制御方法において、

前記位置情報制御ステップは、前記コンピュータ内のキャッシュに最新の位置情報を書き込み、前記コンピュータ内のキャッシュに格納された位置情報のうち所定の使用頻度より低い位置情報を前記キャッシュ装置内のキャッシュへ移動させることを特徴とするストレージ制御方法。

(付記 19) 付記 17 または付記 18 に記載のストレージ制御方法において、

前記位置情報制御ステップは、前記第 1 の階層的な記憶装置に前記ストレージ装置を加えて第 1 の階層的な記憶装置とし、該第 1 の階層的な記憶装置において前記位置情報を管理し、

20

前記ストレージ装置は前記キャッシュ装置内のキャッシュより大容量であることを特徴とするストレージ制御方法。

(付記 20) 付記 19 に記載のストレージ制御方法において、

前記位置情報制御ステップは、前記キャッシュ装置内のキャッシュに格納された位置情報のうち所定の使用頻度より低い位置情報を前記ストレージ装置へ移動させることを特徴とするストレージ制御方法。

【図面の簡単な説明】

【0072】

【図 1】実施の形態 1 に係るストレージシステムの構成の一例を示すブロック図である。

30

【図 2】実施の形態 1 に係る位置情報テーブルの構成の一例を示すブロック図である。

【図 3】実施の形態 1 に係るストレージ制御装置による write 処理の動作の一例を示すフローチャートである。

【図 4】実施の形態 1 に係るストレージ制御装置による read 処理の動作の一例を示すフローチャートである。

【図 5】実施の形態 1 に係るキャッシュプロトコルにおけるメッセージの構造の一例を示す構造図である。

【図 6】実施の形態 1 に係るキャッシュプロトコルにおけるデータの構造の一例を示す構造図である。

【図 7】実施の形態 1 に係る位置情報テーブル検索処理及び位置情報テーブル更新処理の動作の一例を示すフローチャートである。

40

【図 8】実施の形態 2 に係るキャッシュプロトコルにおけるデータの構造の一例を示す構造図である。

【図 9】実施の形態 2 に係る位置情報テーブル検索処理及び位置情報テーブル更新処理の動作の一例を示すフローチャートである。

【符号の説明】

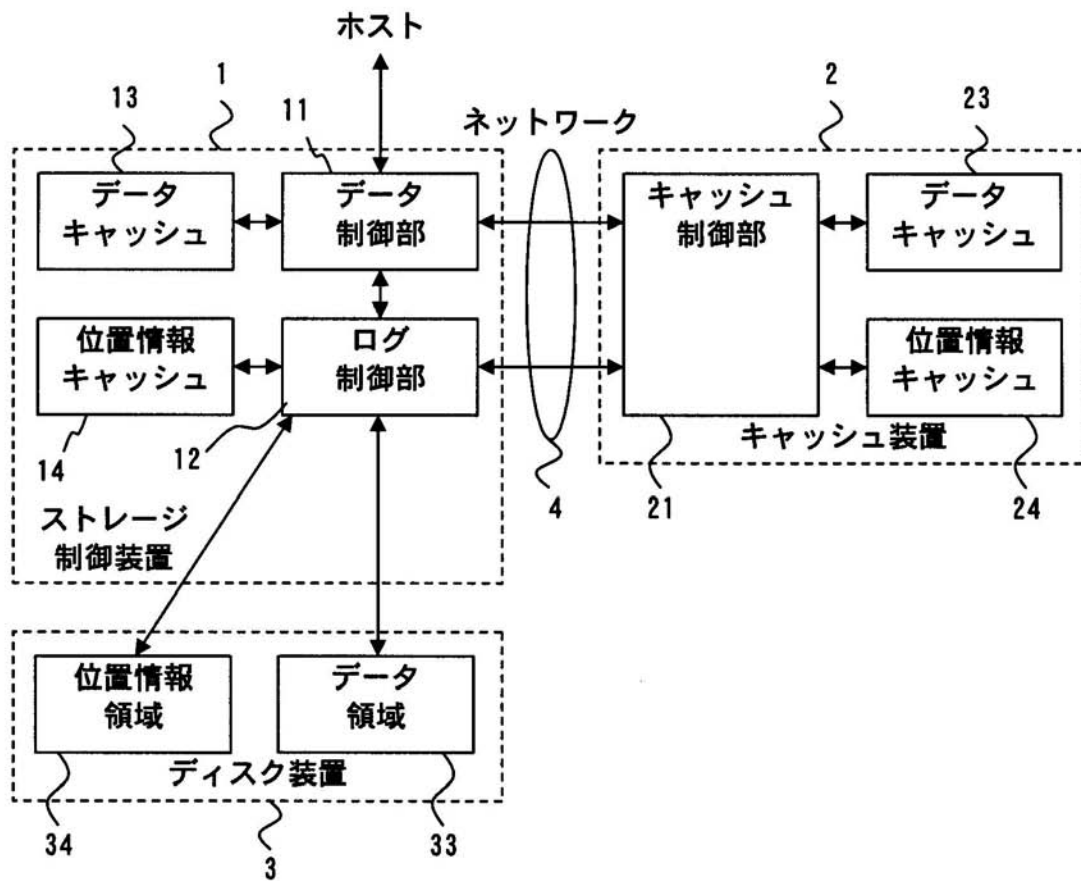
【0073】

1 ストレージ制御装置、2 キャッシュ装置、3 ディスク装置、4 ネットワーク、
11 データ制御部、12 ログ制御部、13, 23 データキャッシュ、14, 24
位置情報キャッシュ、21 キャッシュ制御部、33 データ領域、34 位置情報領域

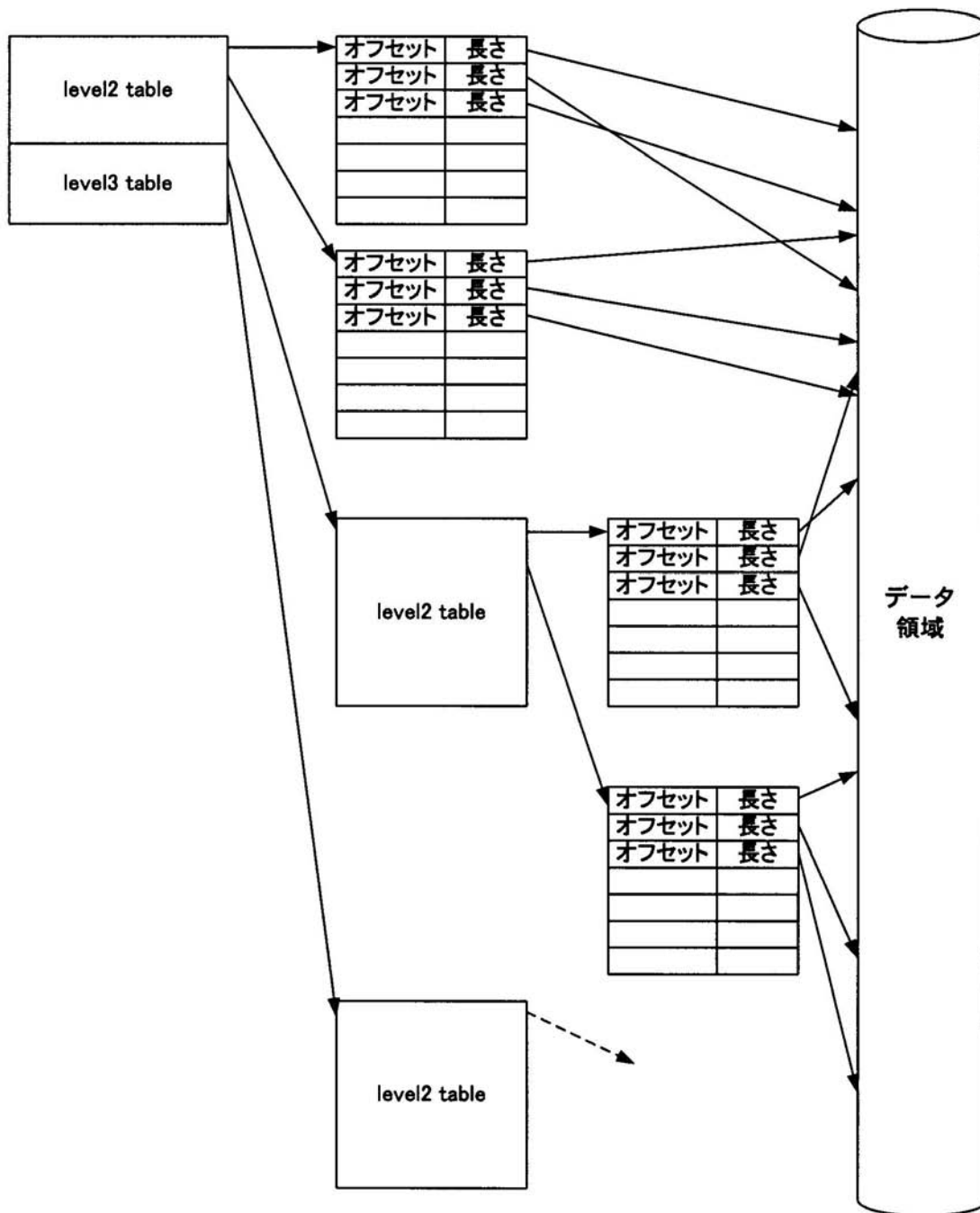
50

•

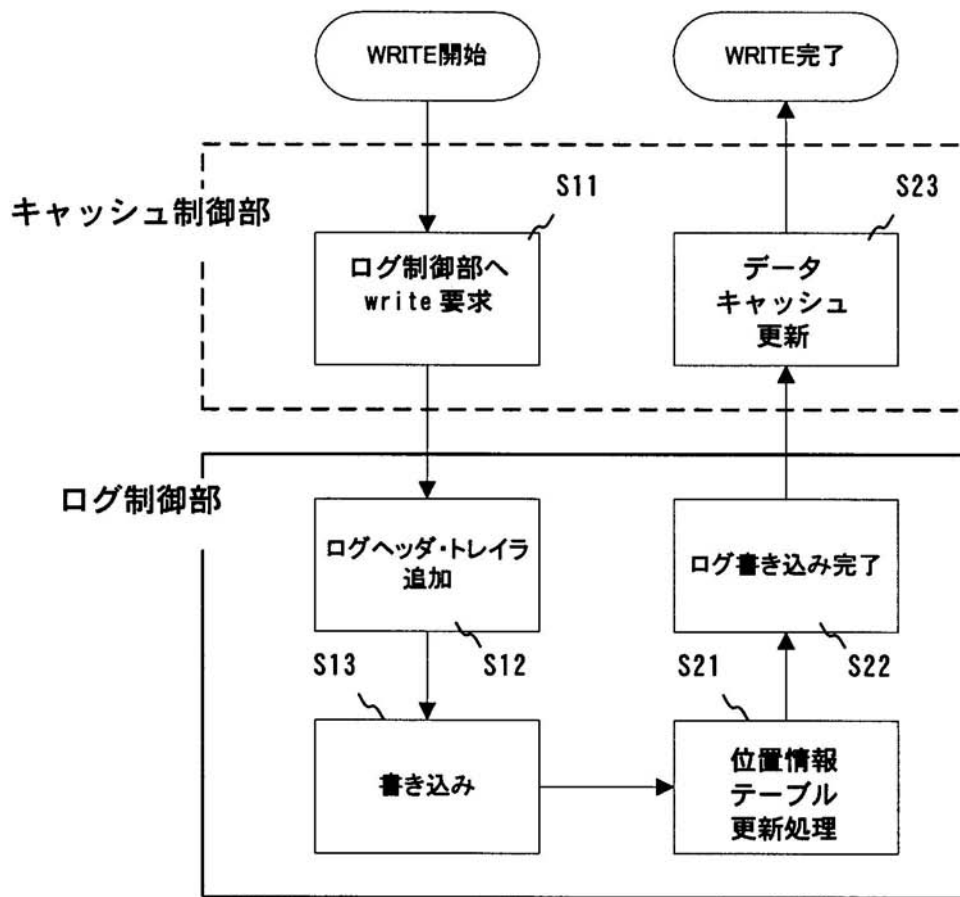
【図 1】



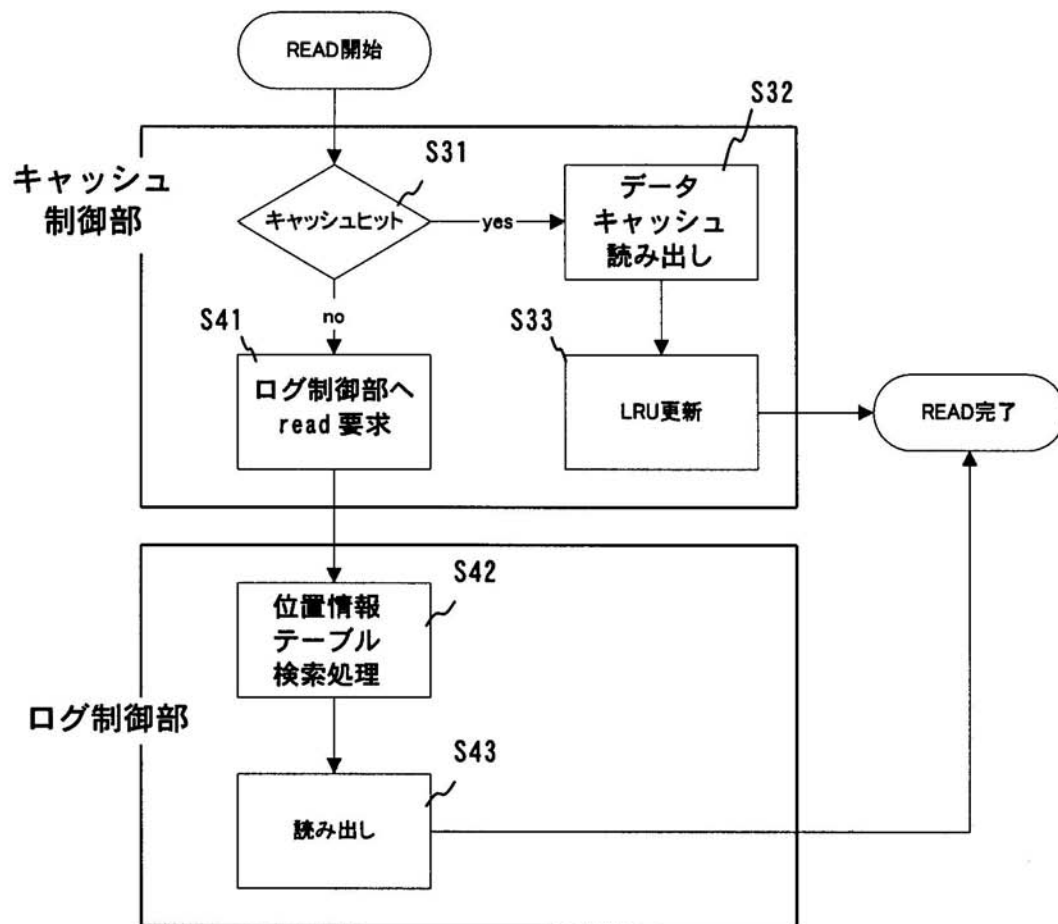
【 図 2 】



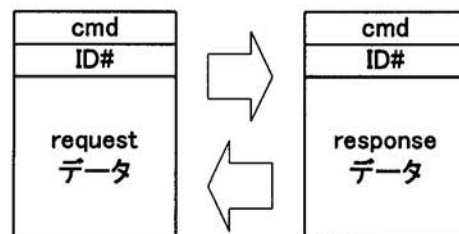
【 図 3 】



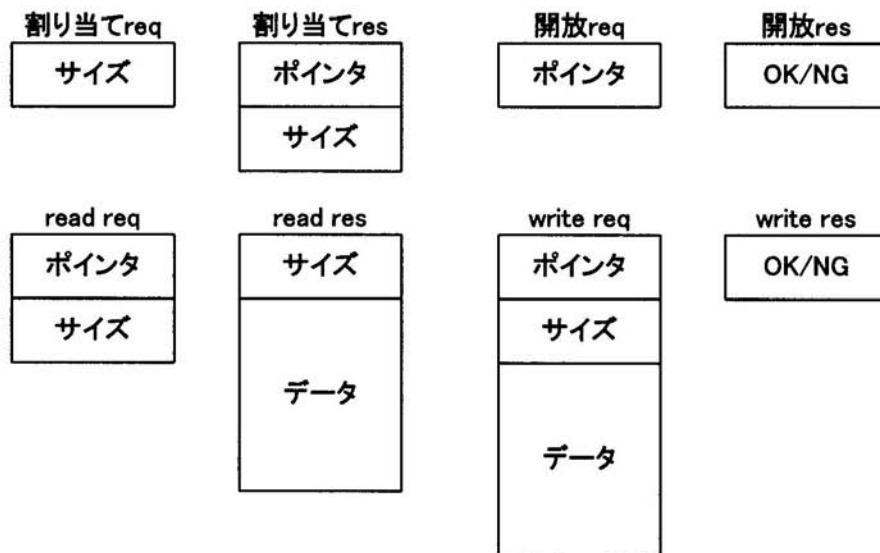
【 図 4 】



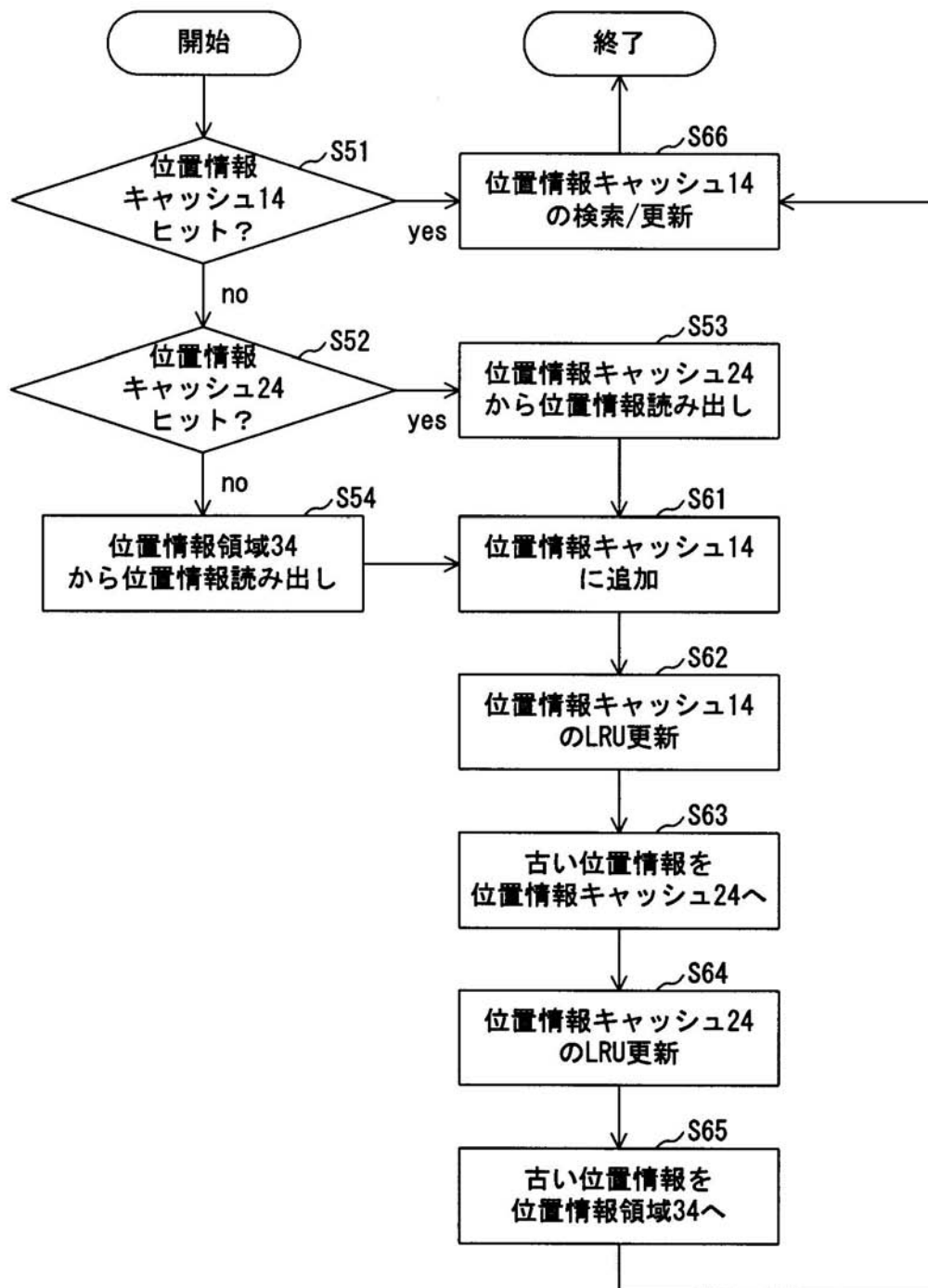
【 図 5 】



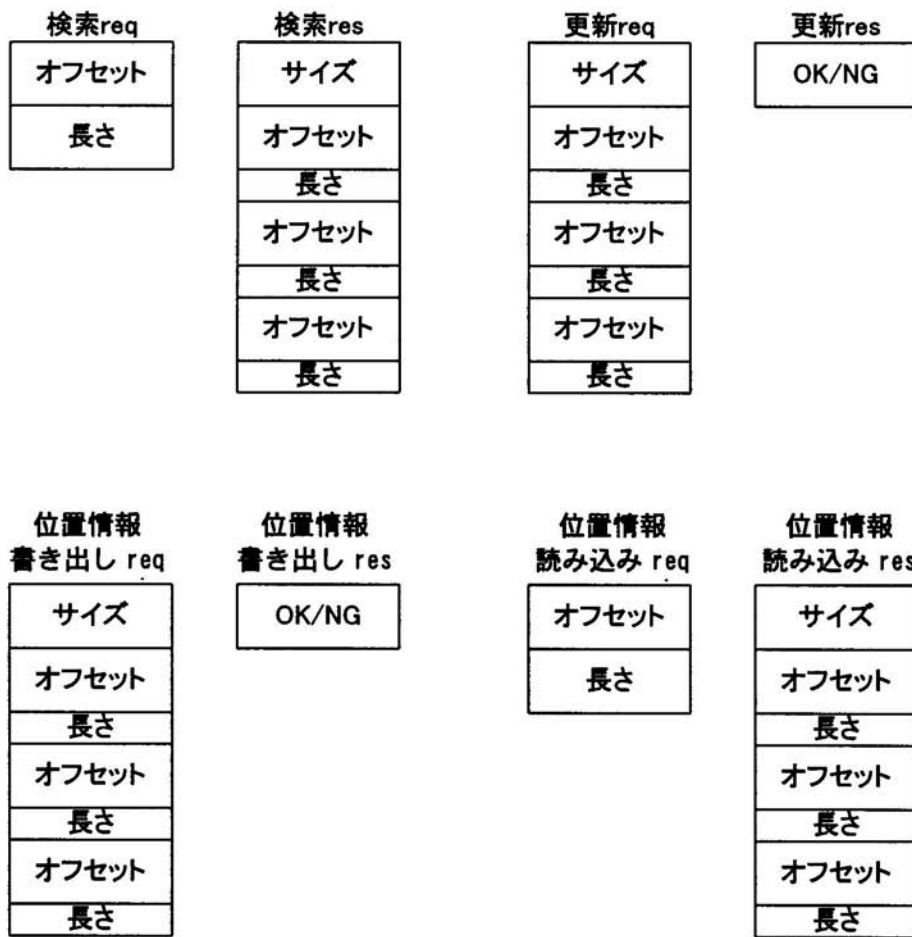
【 図 6 】



【図7】



【 図 8 】



【図9】

