

(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2018-181202

(P2018-181202A)

(43) 公開日 平成30年11月15日(2018.11.15)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 3/06 (2006.01)	G06F 3/06 301K	5B060
G06F 3/08 (2006.01)	G06F 3/08 H	5B160
G06F 12/00 (2006.01)	G06F 3/06 301R	
G06F 12/02 (2006.01)	G06F 12/00 597U	
	G06F 12/02 570A	
審査請求 未請求 請求項の数 8 O L (全 18 頁) 最終頁に続く		

(21) 出願番号 特願2017-83857 (P2017-83857)
 (22) 出願日 平成29年4月20日 (2017. 4. 20)

(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 110002147
 特許業務法人酒井国際特許事務所
 (72) 発明者 武田 直浩
 神奈川県川崎市中原区上小田中4丁目1番
 1号 株式会社富士通コンピュータテクノ
 ロジーズ内
 (72) 発明者 久保田 典秀
 神奈川県川崎市中原区上小田中4丁目1番
 1号 株式会社富士通コンピュータテクノ
 ロジーズ内

最終頁に続く

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

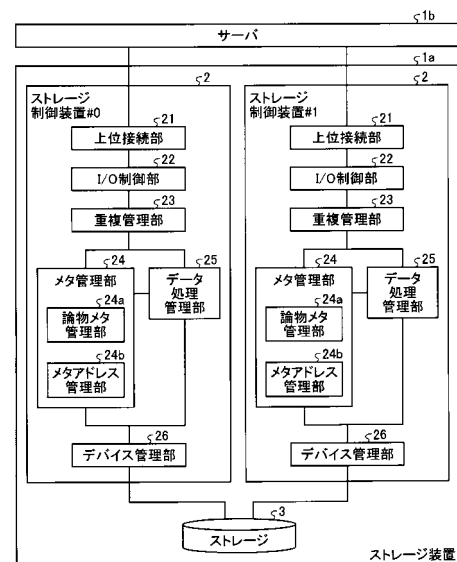
(57) 【要約】

【課題】SSDへの書き込み回数を減少させること。

【解決手段】メタ管理部24は、メタアドレス及び論物メタを用いて、仮想ボリュームの論理アドレスとSSDの物理アドレスの変換処理を行う。メタ管理部24は、論物メタ管理部24aと、メタアドレス管理部24bとを有する。論物メタ管理部24aが、論理アドレスと物理アドレスを対応付ける論物メタの情報を管理する。データ処理管理部25は、論物メタの情報をRAIDユニットの単位でSSDへ追記及びまとめ書きする。

【選択図】図8

実施例に係る情報処理システムの構成を示す図



【特許請求の範囲】**【請求項 1】**

書き込み回数に制限を有する記憶媒体を用いる記憶装置を制御するストレージ制御装置において、

前記記憶装置を使用する情報処理装置がデータの識別に用いる論理アドレスと前記記憶媒体上の該データの記憶位置を示す物理アドレスとを対応付けるアドレス変換情報を管理する変換情報管理部と、

前記変換情報管理部により管理されるアドレス変換情報を前記記憶媒体に追記及びまとめ書きを行う書込処理部と

を有することを特徴とするストレージ制御装置。

10

【請求項 2】

前記変換情報管理部により管理されるアドレス変換情報が前記記憶媒体に追記及びまとめ書きされた位置を示す物理アドレスをメタアドレスとして前記論理アドレスと対応付けて管理する変換情報位置管理部をさらに有することを特徴とする請求項 1 に記載のストレージ制御装置。

【請求項 3】

前記物理アドレスに記憶されたデータには、該データを参照する論理アドレスを示す参照情報が付加されることを特徴とする請求項 1 又は 2 に記載のストレージ制御装置。

【請求項 4】

前記書込処理部は、前記データの更新によって新たなデータが追記された場合に、更新前の前記データとともに前記参照情報を前記記憶媒体上に維持することを特徴とする請求項 3 に記載のストレージ制御装置。

20

【請求項 5】

前記メタアドレスは、前記記憶媒体の所定の位置に記憶されることを特徴とする請求項 2 に記載のストレージ制御装置。

【請求項 6】

前記記憶媒体は、SSDであることを特徴とする請求項 1 ～ 5 のいずれか 1 つに記載のストレージ制御装置。

【請求項 7】

書き込み回数に制限を有する記憶媒体を用いる記憶装置を制御するストレージ制御装置によるストレージ制御方法において、

30

前記記憶装置を使用する情報処理装置がデータの識別に用いる論理アドレスと前記記憶媒体上の該データの記憶位置を示す物理アドレスとを対応付けるアドレス変換情報を管理し、

前記アドレス変換情報を前記記憶媒体に追記及びまとめ書きを行うことを特徴とするストレージ制御方法。

【請求項 8】

書き込み回数に制限を有する記憶媒体を用いる記憶装置を制御するストレージ制御装置が有するコンピュータにより実行されるストレージ制御プログラムにおいて、

前記記憶装置を使用する情報処理装置がデータの識別に用いる論理アドレスと前記記憶媒体上の該データの記憶位置を示す物理アドレスとを対応付けるアドレス変換情報を管理し、

40

前記アドレス変換情報を前記記憶媒体に追記及びまとめ書きを行う

処理を前記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。

【発明の詳細な説明】**【技術分野】****【0001】**

本発明は、ストレージ制御装置、ストレージ制御方法及びストレージ制御プログラムに関する。

【背景技術】

50

【 0 0 0 2 】

昨今、ストレージ装置の記憶媒体は、HDD (Hard Disk Drive) からよりアクセススピードの速いSSD (Solid State Drive) 等のフラッシュメモリに移行している。SSDでは、メモリセルへの上書きを直接行うことはできず、例えば1MB (メガバイト) の大きさのブロックの単位でデータの消去が行われた後にデータの書き込みが行われる。

【 0 0 0 3 】

このため、ブロック内の一部のデータを更新する場合には、ブロック内の他のデータを退避し、ブロックを消去した後に、退避したデータと更新データを書き込むことが行われるため、ブロックの大きさに較べて小さいデータを更新する処理が遅い。また、SSDには書き込み回数の上限がある。このため、SSDでは、ブロックの大きさに較べて小さいデータの更新をできるだけ避けることが望ましい。そこで、ブロック内の一部のデータを更新する場合に、ブロック内の他のデータと更新データを新たなブロックに書き込むことが行われる。

10

【 0 0 0 4 】

しかし、新たなブロックを用いてデータの更新が行われると、データを記憶する物理アドレスが変更されるので、論理アドレスと物理アドレスを対応付ける管理データ (メタデータ) の更新が必要となる。また、ストレージ装置では、データの書き込み容量を削減するために、重複するデータブロックの排除が行われるが、重複排除 (Deduplication) のための管理データの更新も必要となる。

20

【 0 0 0 5 】

なお、ログ構造化ファイルシステムにおいて、ストレージ装置を第1領域と第2領域とに分けて、第1領域と第2領域を以下のように使用する技術がある。第2領域には、多数のデータと、多数のデータと関連する多数のノードとが格納される。第1領域には、多数のノード各々に対応する多数のノード識別子と、多数のノード識別子各々に対応する多数の物理アドレスとを含むノードアドレステーブルが格納される。この技術によれば、メタデータの修正のための付加的な書き込み動作を減らすことができる。

【 0 0 0 6 】

また、ランダムライトアクセスが行われた時に、不使用ページに基づいて選択されたブロックのページに記録されたデータをバッファに書き込み、ブロックを消去後にバッファに書き込まれたデータをブロックに書き込む技術がある。この技術によれば、ガベージコレクションを行わないので、IOPS (Input Output Per Second) 性能を向上させることができる。

30

【 0 0 0 7 】

また、N台のディスク装置から構成されるディスク記憶装置において、 $N \times K$ 個の論理ブロックに相当する容量を持つ書き込みバッファを備え、書き込みバッファに更新すべきデータの論理ブロックを蓄積し、制御装置が、以下の制御を行う技術がある。すなわち、制御装置は、蓄積した論理ブロックが $N \times K - 1$ 個に達するまで論理ブロックの更新を遅延させ、 $N \times K - 1$ 個に達すると当該論理ブロックの論理アドレスタグブロックを加えた $N \times K$ 個の論理ブロックを空領域に連続して順次書き込む。この技術によれば、論理アドレスと物理アドレスのマッピングを原理的に不要とすることで、安価で高速なディスク記憶装置を構築することができる。

40

【 先行技術文献 】

【 特許文献 】

【 0 0 0 8 】

【 特許文献 1 】 特開 2 0 1 4 - 7 1 9 0 6 号 公 報

【 特許文献 2 】 特開 2 0 1 0 - 2 3 7 9 0 7 号 公 報

【 特許文献 3 】 特開 平 1 1 - 5 3 2 3 5 号 公 報

【 発明の概要 】

【 発明が解決しようとする課題 】

50

【 0 0 0 9 】

論理アドレスと物理アドレスを対応付ける管理データ、重複排除のための管理データ等の更新では、ブロック内の一部のデータが更新されるので、管理データはメインメモリ上に配置することが望まれる。しかしながら、管理データが大きくなるとメインメモリ上に全ての管理データを保持することができない。このため、管理データをSSDに書き込むことになるが、管理データの更新によってSSDへの書き込み回数が増えるという問題がある。

【 0 0 1 0 】

本発明は、1つの側面では、管理データの更新によるSSDへの書き込み回数を減少させることを目的とする。

【課題を解決するための手段】

【 0 0 1 1 】

1つの態様では、ストレージ制御装置は、書き込み回数に制限を有する記憶媒体を用いる記憶装置を制御し、変換情報管理部と書込処理部とを有する。変換情報管理部は、記憶装置を使用する情報処理装置がデータの識別に用いる論理アドレスと記憶媒体上の該データの記憶位置を示す物理アドレスとを対応付けるアドレス変換情報を管理する。書込処理部は、変換情報管理部により管理されるアドレス変換情報を記憶媒体に追記及びまとめ書きを行う。

【発明の効果】

【 0 0 1 2 】

1つの側面では、本発明は、管理データの更新によるSSDへの書き込み回数を減少させることができる。

【図面の簡単な説明】

【 0 0 1 3 】

【図1】図1は、実施例に係るストレージ装置の記憶構成を示す図である。

【図2】図2は、RAIDユニットのフォーマットを示す図である。

【図3】図3は、参照メタのフォーマットを示す図である。

【図4】図4は、論物メタのフォーマットを示す図である。

【図5】図5は、実施例に係るメタメタ方式を説明するための図である。

【図6】図6は、メタアドレスのフォーマットを示す図である。

【図7】図7は、ドライブグループにおけるRAIDユニットの配置例を示す図である。

【図8】図8は、実施例に係る情報処理システムの構成を示す図である。

【図9】図9は、機能部間の関係を示す図である。

【図10A】図10Aは、重複のないデータの書き込み処理のシーケンスを示す図である。

【図10B】図10Bは、重複のあるデータの書き込み処理のシーケンスを示す図である。

【図11】図11は、読み出し処理のシーケンスを示す図である。

【図12A】図12Aは、メタメタ方式導入前のスモールライトの数を示す図である。

【図12B】図12Bは、メタメタ方式でのスモールライトの数を示す図である。

【図12C】図12Cは、メタメタ方式で旧データ無効化を行わない場合のスモールライトの数を示す図である。

【図13】図13は、実施例に係るストレージ制御プログラムを実行するストレージ制御装置のハードウェア構成を示す図である。

【発明を実施するための形態】

【 0 0 1 4 】

以下に、本願の開示するストレージ制御装置、ストレージ制御方法及びストレージ制御プログラムの実施例を図面に基づいて詳細に説明する。なお、この実施例は開示の技術を限定するものではない。

【実施例】

【 0 0 1 5 】

まず、実施例に係るストレージ装置のデータ管理方法について図 1 ~ 図 7 を用いて説明する。図 1 は、実施例に係るストレージ装置の記憶構成を示す図である。図 1 に示すように、実施例に係るストレージ装置は、R A I D (Redundant Arrays of Inexpensive Disks) 6 ベースのプール 3 a として複数の S S D 3 d を管理する。また、実施例に係るストレージ装置は、複数のプール 3 a を有する。

【 0 0 1 6 】

プール 3 a には、仮想化プールと階層化プールがある。仮想化プールは 1 つのティア 3 b を有し、階層化プールは 2 つ以上のティア 3 b を有する。ティア 3 b は、1 つ以上のドライブグループ 3 c を有する。ドライブグループ 3 c は、S S D 3 d のグループであり、6 ~ 2 4 台の S S D 3 d を有する。例えば、1 つのストライプを記憶する 6 台の S S D 3 d のうち、3 台はデータ記憶用に用いられ、2 台はパリティ記憶用に用いられ、1 台はホットスペア用に用いられる。なお、ドライブグループ 3 c は、2 5 台以上の S S D 3 d を有してよい。

10

【 0 0 1 7 】

実施例に係るストレージ装置は、R A I D ユニットの単位でデータを管理する。シン・プロビジョニングの物理割当の単位は、一般に固定サイズのチャンク単位で行われ、1 チャンクは 1 R A I D ユニットの該当する。以降の説明では、チャンクを R A I D ユニットと呼称する。R A I D ユニットは、プール 3 a から割り当てられる 2 4 M B の連続な物理領域である。実施例に係るストレージ装置は、データを R A I D ユニット単位でメインメモリ上にバッファリングし、追記型で S S D 3 d に書き込む。

20

【 0 0 1 8 】

図 2 は、R A I D ユニットのフォーマットを示す図である。図 2 に示すように、R A I D ユニットには、複数のユーザデータユニット(データログとも呼ばれる。)が含まれる。ユーザデータユニットには、参照メタと圧縮データが含まれる。参照メタは、S S D 3 d に書き込まれるデータの管理データである。

【 0 0 1 9 】

圧縮データは、S S D 3 d に書き込まれるデータが圧縮されたものである。データの大きさは最大 8 K B (キロバイト)である。圧縮率を 5 0 % とすると、実施例に係るストレージ装置は、1 つの R A I D ユニットの、例えば $24 \text{ MB} \div 4 = 5 \text{ KB}$ 5 4 6 1 個のユーザデータユニットが溜まると、R A I D ユニットの S S D 3 d に書き込む。

30

【 0 0 2 0 】

図 3 は、参照メタのフォーマットを示す図である。図 3 (a) に示すように、参照メタには、S B (Super Block) と最大 6 0 個の参照先の参照 L U N (Logical Unit Number: 論理ユニット番号) / L B A (Logical Block Address: 論理ブロックアドレス) 情報を書き込むことができる記憶容量の領域が確保されている。S B のサイズは 3 2 B (バイト)であり、参照メタのサイズは 5 1 2 B (バイト)である。各参照 L U N / L B A 情報のサイズは 8 B (バイト)である。参照メタは、重複排除により新しい参照先ができると、参照先が追加となり、参照メタが更新される。ただし、データの更新により参照先がなくなった場合にも参照 L U N / L B A 情報は削除されないで保持される。無効になった参照 L U N / L B A 情報はガベージコレクションにより回収される。

40

【 0 0 2 1 】

図 3 (b) に示すように、S B には、4 B の Header Length と、2 0 B の Hash Value と、2 B の Next Offset Block Count が含まれる。Header Length は、参照メタの長さである。Hash Value は、データのハッシュ値であり、重複排除のために用いられる。Next Offset Block Count は、次に格納する参照 L U N / L B A 情報の位置である。なお、Reserved は、将来の拡張用である。

【 0 0 2 2 】

図 3 (c) に示すように、参照 L U N / L B A 情報には、2 B の L U N と、6 B の L B

50

Aが含まれる。

【0023】

また、実施例に係るストレージ装置は、論物変換情報である論物メタを用いてデータの論理アドレスと物理アドレスの対応関係を管理する。図4は、論物メタのフォーマットを示す図である。実施例に係るストレージ装置は、8KBのデータ毎に、図4に示した情報を管理する。

【0024】

図4に示すように、論物メタの大きさは32Bである。論物メタには、2BのLUNと、6BのLBAがデータの論理アドレスとして含まれる。また、論物メタには、2BのCompression Byte Countが、圧縮されたデータのバイト数として含まれる。

10

【0025】

また、論物メタには、2BのNode Noと、1BのStorage Pool Noと、4BのRAID Unit Noと、2BのRAID Unit Offset LBAが物理アドレスとして含まれる。

【0026】

Node Noは、ユーザデータユニットを記憶するRAIDユニットが属するプール3aを担当するストレージ制御装置を識別するための番号である。なお、ストレージ制御装置については後述する。Storage Pool Noは、ユーザデータユニットを記憶するRAIDユニットが属するプール3aを識別するための番号である。RAID Unit Noは、ユーザデータユニットを記憶するRAIDユニットを識別するための番号である。RAID Unit Offset LBAは、ユーザデータユニットのRAIDユニット内でのアドレスである。

20

【0027】

実施例に係るストレージ装置は、RAIDユニットの単位で論物メタを管理する。実施例に係るストレージ装置は、論物メタをRAIDユニット単位でメインメモリ上にバッファリングし、バッファに例えば786432エントリ溜まると、論物メタを追記型でSSD3dにまとめ書きする。このため、実施例に係るストレージ装置は、論物メタがある場所を示す情報をメタメタ方式で管理する。

【0028】

30

図5は、実施例に係るメタメタ方式を説明するための図である。図5(d)に示すように、(1)、(2)、(3)、・・・で表されるユーザデータユニットは、RAIDユニットの単位でSSD3dにまとめ書きされる。そして、図5(c)に示すように、ユーザデータユニットの位置を示す論物メタも、RAIDユニットの単位でSSD3dにまとめ書きされる。

【0029】

そして、実施例に係るストレージ装置は、図5(a)に示すように、論物メタの位置をLUN/LBA毎にメタアドレスを用いてメインメモリ上で管理する。ただし、図5(b)に示すように、メインメモリから溢れたメタアドレス情報は、外部キャッシュ(2次キャッシュ)される。ここで、外部キャッシュとは、SSD3dでのキャッシュである。

40

【0030】

図6は、メタアドレスのフォーマットを示す図である。図6に示すように、メタアドレスの大きさは8Bである。メタアドレスには、Storage Pool Noと、RAID Unit Offset LBAと、RAID Unit Noとが含まれる。メタアドレスは、論物データのSSD3dでの格納位置を示す物理アドレスである。

【0031】

Storage Pool Noは、論物メタを記憶するRAIDユニットが属するプール3aを識別するための番号である。RAID Unit Offset LBAは、論物メタのRAIDユニット内のアドレスである。RAID Unit Noは、論物メタを記憶するRAIDユニットを識別するための番号である。

50

【 0 0 3 2 】

5 1 2 個のメタアドレスがメタアドレスページ (4 K B) として管理され、メタアドレスページの単位でメインメモリ上にキャッシングされる。また、メタアドレス情報は、R A I D ユニットの単位で例えば S S D 3 d の先頭から記憶される。

【 0 0 3 3 】

図 7 は、ドライブグループ 3 c における R A I D ユニットの配置例を示す図である。図 7 に示すように、メタアドレスを記憶する R A I D ユニットは、先頭に配置される。図 7 では、番号が「 0 」～「 1 2 」の R A I D ユニットが、メタアドレスを記憶する R A I D ユニットである。メタアドレスの更新があった場合は、メタアドレスを記憶する R A I D ユニットは上書き保存される。

10

【 0 0 3 4 】

論物メタを記憶する R A I D ユニット及びユーザデータユニットを記憶する R A I D ユニットは、それぞれのバッファがいっぱいになると順番にドライブグループに書き出される。図 7 では、ドライブグループにおいて、番号が「 1 3 」、「 1 7 」、「 2 7 」、「 4 0 」、「 5 1 」、「 6 3 」及び「 7 0 」の R A I D ユニットが、論物メタを記憶する R A I D ユニットであり、その他の R A I D ユニットが、ユーザデータユニットを記憶する R A I D ユニットである。

【 0 0 3 5 】

実施例に係るストレージ装置は、メタメタ方式によって最低限の情報をメインメモリに保持し、論物メタとユーザデータユニットを S S D 3 d に追記及びまとめ書きすることで S S D 3 d への書き込み回数を削減することができる。

20

【 0 0 3 6 】

次に、実施例に係る情報処理システムの構成について説明する。図 8 は、実施例に係る情報処理システムの構成を示す図である。図 8 に示すように、実施例に係る情報処理システム 1 は、ストレージ装置 1 a とサーバ 1 b とを有する。ストレージ装置 1 a は、サーバ 1 b が使用するデータを記憶する装置である。サーバ 1 b は、情報処理などの業務を行う情報処理装置である。ストレージ装置 1 a とサーバ 1 b との間は、F C (Fibre Channel) 及び i S C S I (Internet Small Computer System Interface) で接続される。

【 0 0 3 7 】

ストレージ装置 1 a は、ストレージ装置 1 a を制御するストレージ制御装置 2 とデータを記憶するストレージ (記憶装置) 3 とを有する。ここで、ストレージ 3 は、複数台の記憶装置 (S S D) 3 d の集まりである。

30

【 0 0 3 8 】

なお、図 8 では、ストレージ装置 1 a は、ストレージ制御装置 # 0 及びストレージ制御装置 # 1 で表される 2 台のストレージ制御装置 2 を有するが、ストレージ装置 1 a は、3 台以上のストレージ制御装置 2 を有してよい。また、図 8 では、情報処理システム 1 は、1 台のサーバ 1 b を有するが、情報処理システム 1 は、2 台以上のサーバ 1 b を有してよい。

【 0 0 3 9 】

ストレージ制御装置 2 は、ストレージ 3 を分担して管理し、1 つ以上のプール 3 a を担当する。ストレージ制御装置 2 は、上位接続部 2 1 と、I / O 制御部 2 2 と、重複管理部 2 3 と、メタ管理部 2 4 と、データ処理管理部 2 5 と、デバイス管理部 2 6 とを有する。

40

【 0 0 4 0 】

上位接続部 2 1 は、F C ドライバ及び i S C S I ドライバと I / O 制御部 2 2 との間の情報の受け渡しを行う。I / O 制御部 2 2 は、キャッシュメモリ上のデータを管理する。重複管理部 2 3 は、データ重複排除 / 復元の制御を行うことで、ストレージ装置 1 a 内に格納されているユニークなデータを管理する。

【 0 0 4 1 】

メタ管理部 2 4 は、メタアドレス及び論物メタを管理する。また、メタ管理部 2 4 は、メタアドレス及び論物メタを用いて、仮想ボリュームにおけるデータの識別に用いる論理

50

アドレスとSSD3dにおけるデータが記憶された位置を示す物理アドレスの変換処理を行う。

【0042】

メタ管理部24は、論物メタ管理部24aとメタアドレス管理部24bとを有する。論物メタ管理部24aは、論理アドレスと物理アドレスとを対応付けるアドレス変換情報に関連する論物メタを管理する。論物メタ管理部24aは、論物メタのSSD3dへの書き込み、及び、論物メタのSSD3dからの読み出しをデータ処理管理部25に依頼する。論物メタ管理部24aは、メタアドレスを用いて論物メタの記憶場所を特定する。

【0043】

メタアドレス管理部24bは、メタアドレスを管理する。メタアドレス管理部24bは、メタアドレスの外部キャッシュ(2次キャッシュ)への書き込み、及び、外部キャッシュからのメタアドレスの読み出しをデバイス管理部26に依頼する。

【0044】

データ処理管理部25は、ユーザデータを連続的なユーザデータユニットで管理し、RAIDユニットの単位でSSD3dに追記及びまとめ書きを行う。また、データ処理管理部25は、データの圧縮解凍、参照メタの生成を行う。ただし、データ処理管理部25は、データが更新された場合に、古いデータに対応するユーザデータユニットに含まれる参照メタの更新は行わない。

【0045】

また、データ処理管理部25は、論物メタをRAIDユニットの単位でSSD3dに追記及びまとめ書きを行う。論物メタの書き込みでは、1小ブロック(512B)に論物メタの16エントリが追記書きされるため、データ処理管理部25は、同一小ブロック内にLUNとLBAが同じものが存在しないように管理する。

【0046】

データ処理管理部25は、同一小ブロックにLUNとLBAが同じものが存在しないように管理することで、RAIDユニット番号とRAIDユニット内LBAにより、LUNとLBAを検索することができる。なお、データの消去単位である1MBのブロックと区別するため、ここでは512Bのブロックを小ブロックと呼ぶ。

【0047】

また、メタ管理部24から論物メタの読み出しを要求されると、データ処理管理部25は、メタ管理部24に指定された小ブロックから対象のLUNとLBAを検索して応答する。

【0048】

データ処理管理部25は、メインメモリ上のバッファであるライトバッファにライトデータを溜め、一定の閾値を超えるとSSD3dに書き出す。データ処理管理部25は、プール3aの物理スペースを管理し、RAIDユニットの配置を行う。デバイス管理部26は、RAIDユニットのストレージ3への書き込みを行う。

【0049】

図9は、機能部間の関係を示す図である。図9に示すように、重複管理部23とメタ管理部24との間では、論物メタの取得と更新が行われる。重複管理部23とデータ処理管理部25との間では、ユーザデータユニットのライトバック(writeback)とステージング(staging)が行われる。ここで、ライトバックとは、ストレージ3へのデータの書き込みであり、ステージングとは、ストレージ3からのデータの読み出しである。

【0050】

メタ管理部24とデータ処理管理部25との間では、論物メタのライト(Write)とリード(Read)が行われる。データ処理管理部25とデバイス管理部26との間では、追記データのストレージリードとストレージライトが行われる。メタ管理部24とデバイス管理部26との間では、外部キャッシュのストレージリードとストレージライトが行われる。デバイス管理部26とストレージ3との間では、ストレージリードとストレージライトが行われる。

10

20

30

40

50

【 0 0 5 1 】

次に、書き込み処理のシーケンスについて説明する。図 1 0 A は、重複のないデータの書き込み処理のシーケンスを示す図であり、図 1 0 B は、重複のあるデータの書き込み処理のシーケンスを示す図である。

【 0 0 5 2 】

重複のないデータの書き込み処理では、図 1 0 A に示すように、I / O 制御部 2 2 は、重複管理部 2 3 にデータのライトバックを要求する（ステップ S 1）。すると、重複管理部 2 3 は、データの重複判定を行って重複はないので（ステップ S 2）、新規のユーザデータユニットのライトをデータ処理管理部 2 5 に要求する（ステップ S 3）。

【 0 0 5 3 】

すると、データ処理管理部 2 5 は、ライトバッファを獲得し（ステップ S 4）、デバイス管理部 2 6 に R U（R A I D ユニット）の取得を要求する（ステップ S 5）。なお、ライトバッファを既に獲得している場合には、ライトバッファの獲得は不要である。そして、データ処理管理部 2 5 は、デバイス管理部 2 6 から D P #（S t o r a g e P o o l N o）と R U #（R A I D U n i t N o）を獲得する（ステップ S 6）。

【 0 0 5 4 】

そして、データ処理管理部 2 5 は、データを圧縮し（ステップ S 7）、参照メタを生成する（ステップ S 8）。そして、データ処理管理部 2 5 は、ライトバッファ内でユーザデータユニットの追記書きを行い（ステップ S 9）、まとめ書き判定を行う（ステップ S 1 0）。そして、データ処理管理部 2 5 は、まとめ書きが必要と判定した場合には、デバイス管理部 2 6 へライトバッファのまとめ書きを依頼する。そして、データ処理管理部 2 5 は、重複管理部 2 3 に D P # と R U # を応答する（ステップ S 1 1）。

【 0 0 5 5 】

すると、重複管理部 2 3 は、メタ管理部 2 4 に論物メタの更新を要求し（ステップ S 1 2）、メタ管理部 2 4 は、更新した論物メタのライトをデータ処理管理部 2 5 に要求する（ステップ S 1 3）。

【 0 0 5 6 】

すると、データ処理管理部 2 5 は、ライトバッファを獲得し（ステップ S 1 4）、デバイス管理部 2 6 に R U の取得を要求する（ステップ S 1 5）。なお、獲得するライトバッファは、ユーザデータユニット用のライトバッファとは異なるバッファである。また、ライトバッファを既に獲得している場合には、ライトバッファの獲得は不要である。そして、データ処理管理部 2 5 は、デバイス管理部 2 6 から D P # と R U # を獲得する（ステップ S 1 6）。

【 0 0 5 7 】

そして、データ処理管理部 2 5 は、ライトバッファ内で論物メタの追記書きを行い（ステップ S 1 7）、まとめ書き判定を行う（ステップ S 1 8）。そして、データ処理管理部 2 5 は、まとめ書きが必要と判定した場合には、デバイス管理部 2 6 へライトバッファのまとめ書きを依頼する。そして、データ処理管理部 2 5 は、メタ管理部 2 4 に D P # と R U # を応答する（ステップ S 1 9）。

【 0 0 5 8 】

すると、メタ管理部 2 4 は、メタアドレス更新のためにメタアドレスの追い出しが必要か否かを判定し（ステップ S 2 0）、追い出しが必要と判定した場合には、デバイス管理部 2 6 に追い出しを依頼する。そして、D P # と R U # に基づいてメタアドレスを更新する（ステップ S 2 1）。

【 0 0 5 9 】

そして、メタ管理部 2 4 は、重複管理部 2 3 に完了を通知し（ステップ S 2 2）、重複管理部 2 3 は、メタ管理部 2 4 から完了を通知されると、I / O 制御部 2 2 に完了を通知する（ステップ S 2 3）。

【 0 0 6 0 】

このように、データ処理管理部 2 5 が、ユーザデータユニットに加えて論物メタも追記

10

20

30

40

50

書き及びまとめ書きを行うことで、SSD3dへの書き込み回数を減らすことができる。

【0061】

また、重複のあるデータの書き込みの場合には、図10Bに示すように、I/O制御部22は、重複管理部23にデータのライトバックを要求する(ステップS31)。すると、重複管理部23は、データの重複判定を行って重複があるので(ステップS32)、重複するユーザデータユニットのライトをデータ処理管理部25に要求する(ステップS33)。

【0062】

すると、データ処理管理部25は、重複するユーザデータユニットを含むRAIDユニットについて、ストレージ3のリードをデバイス管理部26に要求する(ステップS34)。そして、デバイス管理部26が、重複するユーザデータユニットを含むRAIDユニットを読み出してデータ処理管理部25に応答する(ステップS35)。そして、データ処理管理部25は、ハッシュ値を比較し(ステップS36)、データの重複を確認する。

【0063】

そして、データ処理管理部25は、重複が確認されれば、重複するユーザデータユニットにおける参照メタに参照先を追加して参照メタを更新する(ステップS37)。データ処理管理部25は、参照メタが更新されたユーザデータユニットを含むRAIDユニットのストレージ3へのライトをデバイス管理部26に要求し(ステップS38)、デバイス管理部26から応答を受け取る(ステップS39)。そして、データ処理管理部25は、重複管理部23にDP#とRU#を応答する(ステップS40)。

【0064】

すると、重複管理部23は、メタ管理部24に論物メタの更新を要求し(ステップS41)、メタ管理部24は、更新した論物メタのライトをデータ処理管理部25に要求する(ステップS42)。

【0065】

すると、データ処理管理部25は、ライトバッファを獲得し(ステップS43)、デバイス管理部26にRUの取得を要求する(ステップS44)。そして、データ処理管理部25は、デバイス管理部26からDP#とRU#を獲得する(ステップS45)。

【0066】

そして、データ処理管理部25は、ライトバッファ内で論物メタの追記書きを行い(ステップS46)、まとめ書き判定を行う(ステップS47)。そして、データ処理管理部25は、まとめ書きが必要と判定した場合には、デバイス管理部26へライトバッファのまとめ書きを依頼する。そして、データ処理管理部25は、メタ管理部24にDP#とRU#を応答する(ステップS48)。

【0067】

すると、メタ管理部24は、メタアドレス更新のためにメタアドレスの追い出しが必要か否かを判定し(ステップS49)、追い出しが必要と判定した場合には、デバイス管理部26に追い出しを依頼する。そして、DP#とRU#に基づいてメタアドレスを更新する(ステップS50)。

【0068】

そして、メタ管理部24は、重複管理部23に完了を通知し(ステップS51)、重複管理部23は、メタ管理部24から完了を通知されると、I/O制御部22に完了を通知する(ステップS52)。

【0069】

このように、データ処理管理部25が、重複データについて、論物メタを追記書き及びまとめ書きを行うことで、SSD3dへの書き込み回数を減らすことができる。

【0070】

次に、読み出し処理のシーケンスについて説明する。図11は、読み出し処理のシーケンスを示す図である。図11に示すように、I/O制御部22は、データのステージングを重複管理部23に要求する(ステップS61)。すると、重複管理部23は、データの

10

20

30

40

50

論物メタの獲得をメタ管理部 2 4 に要求する (ステップ S 6 2)。

【0071】

すると、メタ管理部 2 4 は、データのメタアドレスがメインメモリ上にあることを確認し (ステップ S 6 3)、メタアドレスを指定して、論物メタのリードをデータ処理管理部 2 5 に要求する (ステップ S 6 4)。なお、データのメタアドレスがメインメモリ上にならない場合には、メタ管理部 2 4 は、デバイス管理部 2 6 にストレージ 3 からのリードを要求する。

【0072】

そして、データ処理管理部 2 5 は、論物メタを含む R A I D ユニットについて、ストレージ 3 からのリードをデバイス管理部 2 6 に要求し (ステップ S 6 5)、デバイス管理部 2 6 から R A I D ユニットを受け取る (ステップ S 6 6)。そして、データ処理管理部 2 5 は、R A I D ユニットから論物メタを検索し (ステップ S 6 7)、検索した論物メタをメタ管理部 2 4 に渡す (ステップ S 6 8)。

【0073】

すると、メタ管理部 2 4 は、論物メタを解析し (ステップ S 6 9)、ユーザデータユニットを含む R A I D ユニットの D P # と R U # と O f f s e t を重複管理部 2 3 に渡す (ステップ S 7 0)。ここで、O f f s e t は、R A I D ユニット内のユーザデータユニットのアドレスである。すると、重複管理部 2 3 は、D P # と R U # と O f f s e t を指定してユーザデータユニットのリードをデータ処理管理部 2 5 に要求する (ステップ S 7 1)。

【0074】

すると、データ処理管理部 2 5 は、ユーザデータユニットを含む R A I D ユニットのストレージ 3 からのリードをデバイス管理部 2 6 に要求し (ステップ S 7 2)、デバイス管理部 2 6 から R A I D ユニットを受け取る (ステップ S 7 3)。そして、データ処理管理部 2 5 は、O f f s e t を用いて R A I D ユニットから取り出したユーザデータユニットに含まれる圧縮データを伸張し (ステップ S 7 4)、ユーザデータユニットから参照メタを削除する (ステップ S 7 5)。

【0075】

そして、データ処理管理部 2 5 は、データを重複管理部 2 3 に渡し (ステップ S 7 6)、重複管理部 2 3 は、I / O 制御部 2 2 にデータを渡す (ステップ S 7 7)。

【0076】

このように、ストレージ制御装置 2 は、メタアドレスを用いて論物メタを取得し、論物メタを用いてユーザデータユニットを取得することで、ストレージ 3 からデータを読み出すことができる。

【0077】

次に、ストレージ制御装置 2 による書き込み処理の効果について図 1 2 A ~ 図 1 2 C を用いて説明する。図 1 2 A は、メタメタ方式導入前のスモールライトの数を示す図であり、図 1 2 B は、メタメタ方式でのスモールライトの数を示す図であり、図 1 2 C は、メタメタ方式で旧データ無効化を行わない場合のスモールライトの数を示す図である。ここで、スモールライトとは、ブロック (1 M B) に較べて小さな単位 (4 K B) のライトである。

【0078】

図 1 2 A に示すように、メタメタ方式導入前は、サーバ 1 b からの 8 K B のデータの書き込みに対して、データに関してはまとめ書きが行われるが、論物メタの更新、参照メタの更新に関してスモールライトが行われる。ここで、参照メタの更新としては、旧データ (参照 L U N / L B A 情報) の無効化がある。また、R A I D 6 の場合、データの書き込みに対応して 2 つのパリティ P と Q の書き込みが行われる。したがって、合計 6 回のスモールライトが行われる。

【0079】

これに対して、メタメタ方式では、図 1 2 B に示すように、論物メタの更新は追記書き

10

20

30

40

50

のためスモールライトではなくなり、3回のスモールライトで済む。さらに、旧データ無効化を行わない場合、図12Cに示すように、参照メタの更新もなくなり、スモールライトがなくなる。

【0080】

このように、メタメタ方式を用いることにより、ストレージ制御装置2は、スモールライトの数を減らし、書き込み処理を高速化することができる。また、ストレージ制御装置2は、旧データ無効化を行わないことで、さらにスモールライトの数を減らすことができる。

【0081】

上述してきたように、実施例では、論物メタ管理部24aが、データの論理アドレスと物理アドレスを対応付ける論物メタの情報を管理し、データ処理管理部25が、論物メタの情報をRAIDユニットの単位でSSD3dへ追記及びまとめ書きする。したがって、ストレージ制御装置2は、スモールライトの数を減らし、書き込み処理を高速化することができる。

【0082】

また、実施例では、メタアドレス管理部24bが、論理アドレスと論物メタのアドレスを対応付けるメタアドレスの情報を管理するので、論物メタ管理部24aはメタアドレスを用いて論物メタの位置を特定することができる。

【0083】

また、実施例では、データが更新された場合に、古いデータに対応するユーザデータユニットの参照メタを更新しない。したがって、ストレージ制御装置2は、スモールライトの数をさらに減らすことができる。

【0084】

また、実施例では、メタアドレスはメインメモリ上で管理され、溢れたメタアドレスの情報はSSD3dの所定の位置に記憶される。したがって、ストレージ制御装置2は、SSD3dの所定の位置から読み出すことでメタアドレスの情報を取得することができる。

【0085】

なお、実施例では、ストレージ制御装置2について説明したが、ストレージ制御装置2が有する構成をソフトウェアによって実現することで、同様の機能を有するストレージ制御プログラムを得ることができる。そこで、ストレージ制御プログラムを実行するストレージ制御装置2のハードウェア構成について説明する。

【0086】

図13は、実施例に係るストレージ制御プログラムを実行するストレージ制御装置2のハードウェア構成を示す図である。図13に示すように、ストレージ制御装置2は、メモリ41と、プロセッサ42と、ホストI/F43と、通信I/F44と、接続I/F45とを有する。

【0087】

メモリ41は、プログラムやプログラムの実行途中結果などを記憶するRAM(Random Access Memory)である。プロセッサ42は、メモリ41からプログラムを読み出して実行する処理装置である。

【0088】

ホストI/F43は、サーバ1bとのインタフェースである。通信I/F44は、他のストレージ制御装置2と通信するためのインタフェースである。接続I/F45は、ディスク3とのインタフェースである。

【0089】

そして、プロセッサ42において実行されるストレージ制御プログラムは、可搬記録媒体51に記憶され、メモリ41に読み込まれる。あるいは、ストレージ制御プログラムは、通信インタフェース44を介して接続されたコンピュータシステムのデータベースなどに記憶され、これらのデータベースから読み出されてメモリ41に読み込まれる。

【0090】

10

20

30

40

50

また、実施例では、SSD 3 d を不揮発性記憶媒体として用いる場合について説明したが、本発明はこれに限定されるものではなく、SSD 3 d と同様なデバイス特性を有する他の不揮発性記憶媒体を用いる場合にも同様に適用することができる。

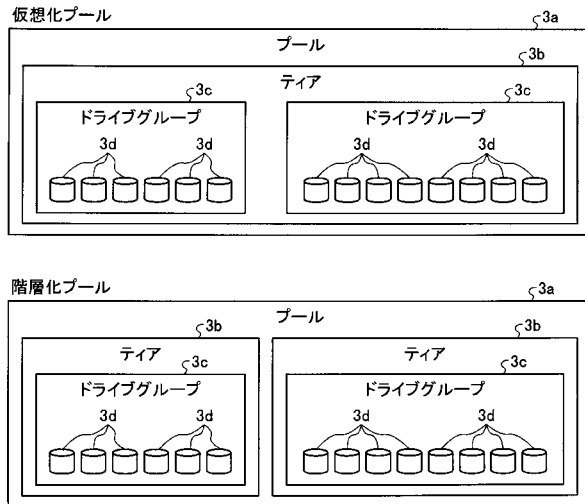
【符号の説明】

【0091】

1	情報処理システム	
1 a	ストレージ装置	
1 b	サーバ	
2	ストレージ制御装置	
3	ストレージ	10
3 a	プール	
3 b	ティア	
3 c	ドライブグループ	
3 d	SSD	
2 1	上位接続部	
2 2	I / O 制御部	
2 3	重複管理部	
2 4	メタ管理部	
2 4 a	論物メタ管理部	
2 4 b	メタアドレス管理部	20
2 5	データ処理管理部	
2 6	デバイス管理部	
4 1	メモリ	
4 2	プロセッサ	
4 3	ホスト I / F	
4 4	通信 I / F	
4 5	接続 I / F	
5 1	可搬記録媒体	

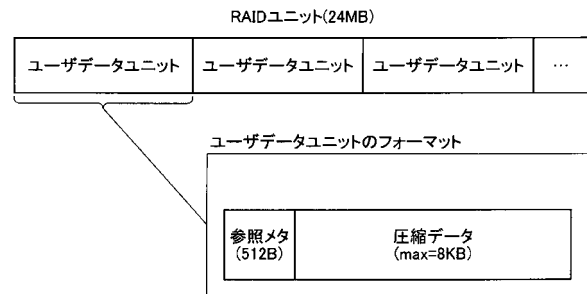
【図 1】

実施例に係るストレージ装置の記憶構成を示す図



【図 2】

RAIDユニットのフォーマットを示す図



【図 3】

参照メタのフォーマットを示す図



(b)SB(Super Block)

オフセット	バイト#0	バイト#1	バイト#2	バイト#3
0x00	Header Length			
0x04	Hash Value(20B)			
:				
0x14				
0x18	Next Offset Block Count		Reserved	
0x1c	Reserved			

(c)参照LUN/LBA情報

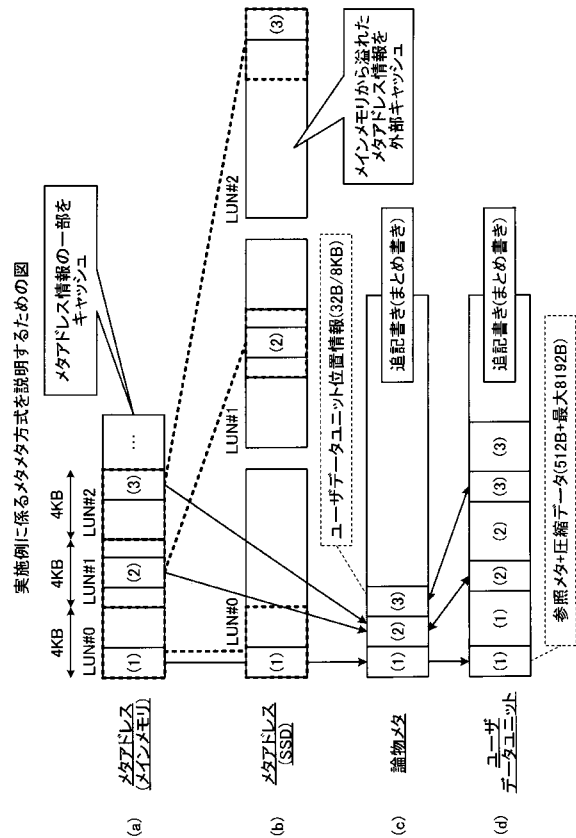
オフセット	バイト#0	バイト#1	バイト#2	バイト#3
0x00	LUN		LBA(上位バイト)	
0x04	LBA(下位バイト)			

【図 4】

論物メタのフォーマットを示す図

オフセット	バイト#0	バイト#1	バイト#2	バイト#3
0x00	Reserved	
0x04	LUN		LBA(上位バイト)	
0x08	LBA(下位バイト)			
0x0c	Compression Byte Count		Reserved	
0x10	Node No		Storage Pool No	Reserved
0x14	RAID Unit No			
0x18	RAID Unit Offset LBA		...	
0x1c	...			

【図 5】

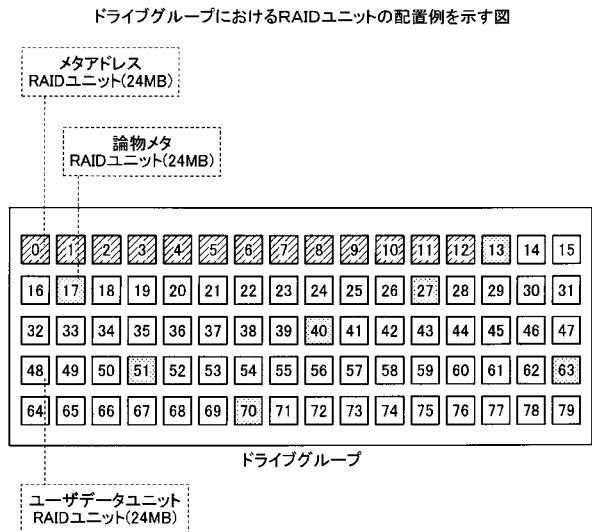


【図 6】

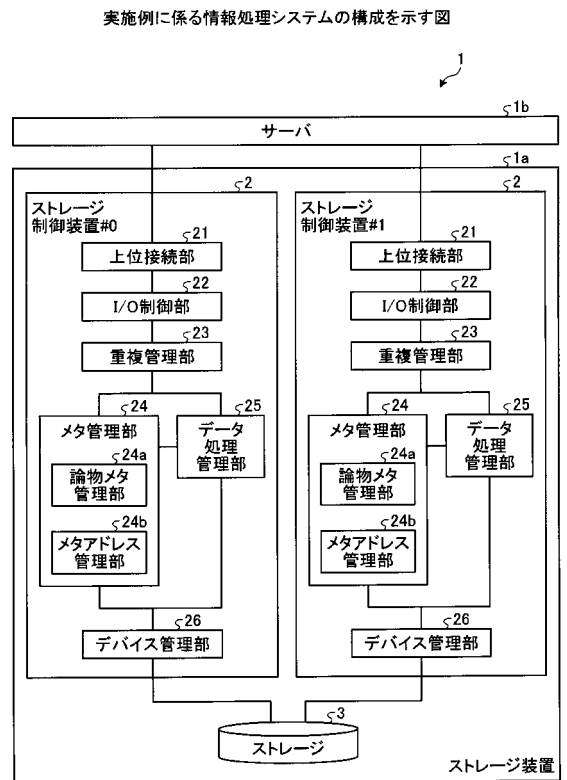
メタアドレスのフォーマットを示す図

オフセット	バイト#0	バイト#1	バイト#2	バイト#3
0x00	Storage Pool No	Reserved	RAID Unit Offset LBA	
0x04	RAID Unit No			

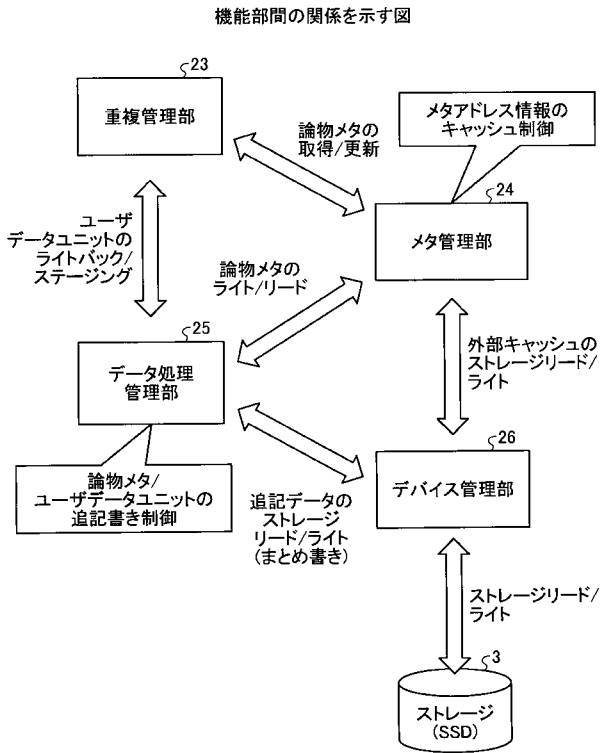
【図 7】



【図 8】

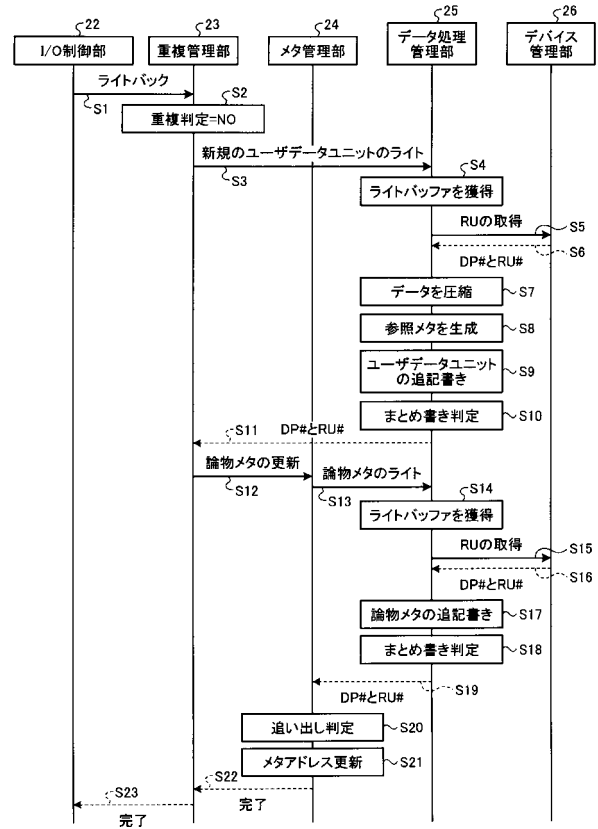


【図 9】



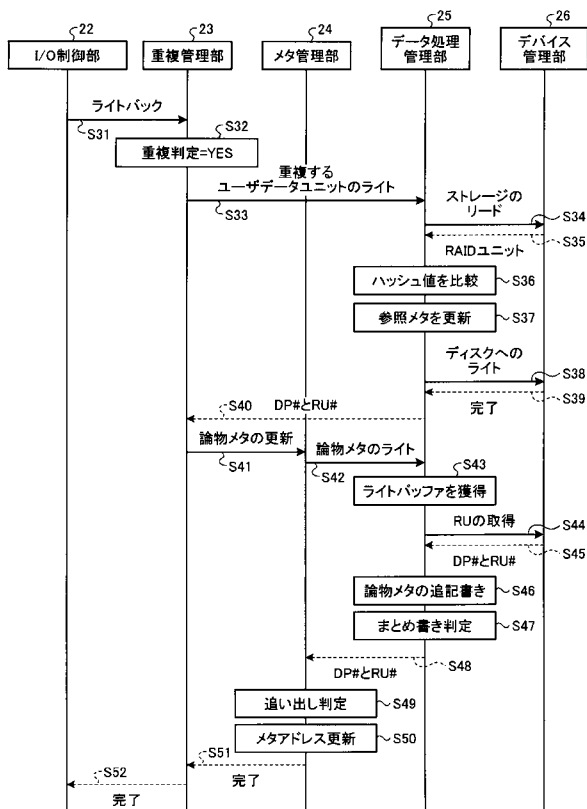
【図 10 A】

重複のないデータの書き込み処理のシーケンスを示す図



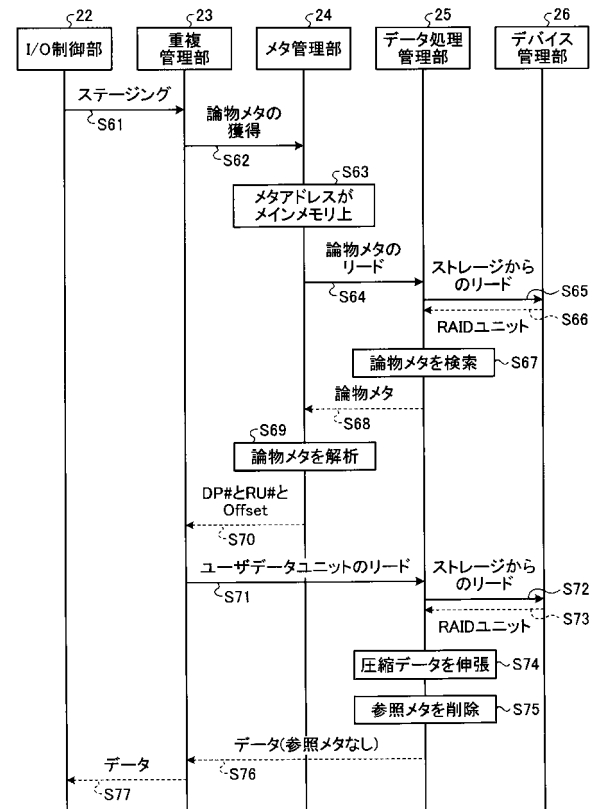
【図 10 B】

重複のあるデータの書き込み処理のシーケンスを示す図

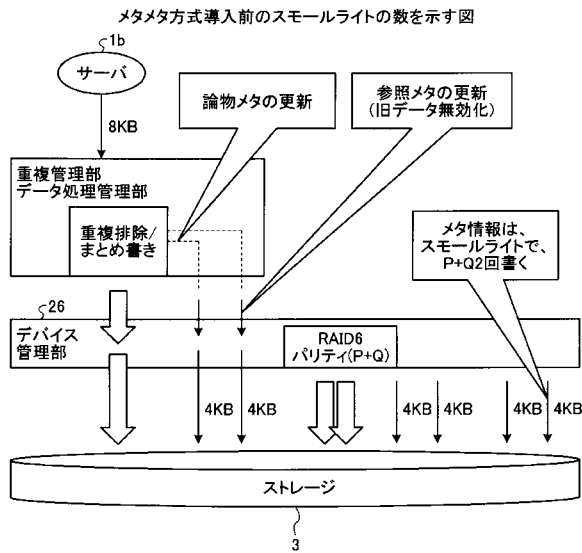


【図 11】

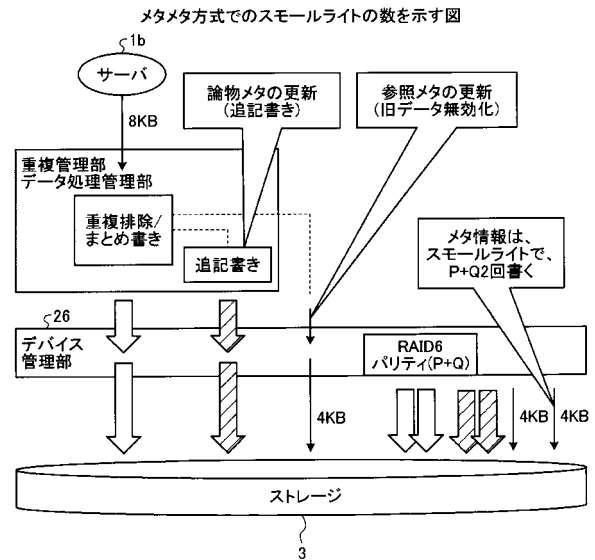
読み出し処理のシーケンスを示す図



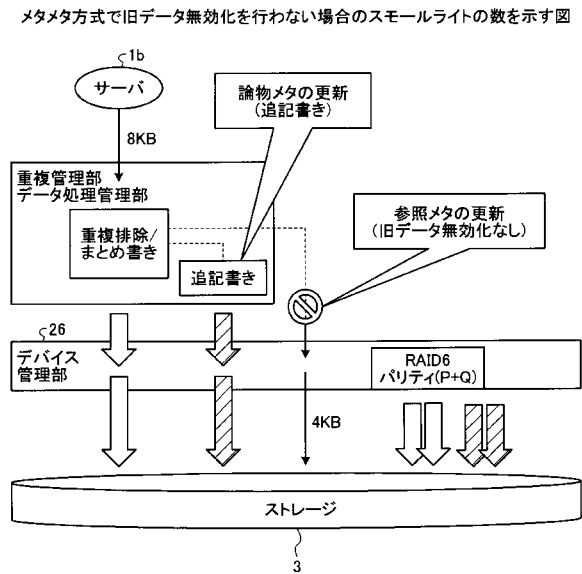
【図 1 2 A】



【図 1 2 B】

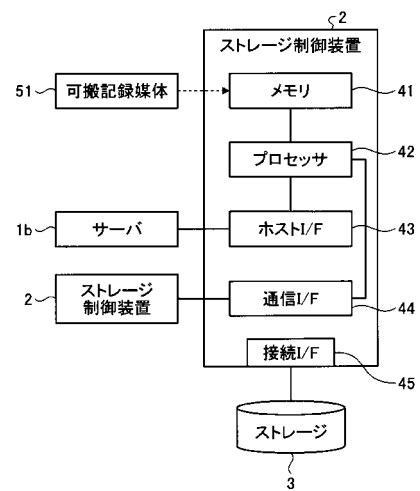


【図 1 2 C】



【図 1 3】

実施例に係るストレージ制御プログラムを実行するストレージ制御装置の
ハードウェア構成を示す図



フロントページの続き

(51)Int.Cl. F I テーマコード(参考)
G 0 6 F 12/02 5 1 0 A

- (72)発明者 紺田 與志仁
神奈川県川崎市中原区上小田中4丁目1番1号 株式会社富士通コンピュータテクノロジーズ内
- (72)発明者 倉澤 祐輔
神奈川県川崎市中原区上小田中4丁目1番1号 株式会社富士通コンピュータテクノロジーズ内
- (72)発明者 菊池 利夫
神奈川県川崎市中原区上小田中4丁目1番1号 株式会社富士通コンピュータテクノロジーズ内
- (72)発明者 田中 勇至
神奈川県川崎市中原区上小田中4丁目1番1号 株式会社富士通コンピュータテクノロジーズ内
- (72)発明者 梶山 真理乃
神奈川県川崎市中原区上小田中4丁目1番1号 株式会社富士通コンピュータテクノロジーズ内
- (72)発明者 鈴木 悠介
神奈川県川崎市中原区上小田中4丁目1番1号 株式会社富士通コンピュータテクノロジーズ内
- (72)発明者 篠 崎 祥成
神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
- (72)発明者 渡辺 岳志
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
- F ターム(参考) 5B060 AB26
5B160 AB26