

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6870246号
(P6870246)

(45) 発行日 令和3年5月12日(2021.5.12)

(24) 登録日 令和3年4月19日(2021.4.19)

(51) Int.Cl.		F I			
G06F 3/06	(2006.01)	G06F 3/06	301J		
G06F 3/08	(2006.01)	G06F 3/06	304Z		
G06F 13/10	(2006.01)	G06F 3/08	H		
		G06F 13/10	340A		

請求項の数 2 (全 31 頁)

(21) 出願番号	特願2016-174473 (P2016-174473)	(73) 特許権者	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
(22) 出願日	平成28年9月7日(2016.9.7)	(74) 代理人	100092978 弁理士 真田 有
(65) 公開番号	特開2018-41245 (P2018-41245A)	(72) 発明者	山本 和彦 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
(43) 公開日	平成30年3月15日(2018.3.15)	審査官	打出 義尚
審査請求日	令和1年6月11日(2019.6.11)		

最終頁に続く

(54) 【発明の名称】 ストレージ装置、及びストレージ制御装置

(57) 【特許請求の範囲】

【請求項1】

寿命の異なる複数の記憶部と、
前記複数の記憶部に格納されるデータブロックを管理する管理部と、
ホストからの、前記寿命の異なる複数の記憶部に対する書込み対象データブロックを、
前記データブロックとして一時的に格納する一時記憶部と、を備え、
前記管理部は、
前記データブロックの重複を排除する重複排除機能と、
前記重複排除機能によって管理される、前記一時記憶部に格納された前記データブロックの参照数に基づいて、前記データブロックに対するアクセス特性の判定または推定を行なう判定処理部と、
前記寿命の異なる複数の記憶部のうち、前記判定処理部によって判定された前記アクセス特性に応じた記憶部に、前記データブロックを格納する格納処理部と、を有し、
前記判定処理部は、前記参照数が2以上のデータブロックの前記アクセス特性は読み重視であると判定する一方、前記参照数が1で書込み後所定時間経過したデータブロックの前記アクセス特性は書込み重視であると推定し、
前記格納処理部は、
前記判定処理部によって前記アクセス特性が読み重視であると判定されたデータブロックを、前記一時記憶部から、前記寿命の異なる複数の記憶部のうち、他の記憶部よりも寿命の短い記憶部に格納する一方、

10

20

前記判定処理部によって前記アクセス特性が書込み重視であると推定されたデータブロックを、前記一時記憶部から、前記寿命の異なる複数の記憶部のうち、他の記憶部よりも寿命の長い記憶部に格納する、ストレージ装置。

【請求項 2】

寿命の異なる複数の記憶部に格納されるデータブロックを管理する管理部と、
ホストからの、前記寿命の異なる複数の記憶部に対する書込み対象データブロックを、
前記データブロックとして一時的に格納する一時記憶部と、を備え、

前記管理部は、

前記データブロックの重複を排除する重複排除機能と、

前記重複排除機能によって管理される、前記一時記憶部に格納された前記データブロックの参照数に基づいて、前記データブロックに対するアクセス特性の判定または推定を行なう判定処理部と、

前記寿命の異なる複数の記憶部のうち、前記判定処理部によって判定された前記アクセス特性に応じた記憶部に、前記データブロックを格納する格納処理部と、を有し、

前記判定処理部は、前記参照数が 2 以上のデータブロックの前記アクセス特性は読み出し重視であると判定する一方、前記参照数が 1 で書込み後所定時間経過したデータブロックの前記アクセス特性は書込み重視であると推定し、

前記格納処理部は、

前記判定処理部によって前記アクセス特性が読み出し重視であると判定されたデータブロックを、前記一時記憶部から、前記寿命の異なる複数の記憶部のうち、他の記憶部よりも寿命の短い記憶部に格納する一方、

前記判定処理部によって前記アクセス特性が書込み重視であると推定されたデータブロックを、前記一時記憶部から、前記寿命の異なる複数の記憶部のうち、他の記憶部よりも寿命の長い記憶部に格納する、ストレージ制御装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ストレージ装置、及びストレージ制御装置に関する。

【背景技術】

【0002】

SSD (Solid State Drive) は、フラッシュメモリを利用した記憶媒体である。フラッシュメモリには「書込み・消去動作による素子の劣化」という特性があることから、フラッシュメモリを素材とした SSD には書込み量の上限が存在する。書込み量の上限による SSD 寿命の指標は、例えば、DWPD によって表現される。DWPD は、Drive Write Per Day の略記であり、一日あたりにドライブ容量全体を何回書き換えることができるかを示す値である。

【0003】

書込み量の上限を持つ SSD の寿命を延ばすため、SSD を採用したストレージ装置では、SSD へ書き込むデータの量を小さくする「データ圧縮技術」や、データの書込みそのものを減らす「データ重複排除技術」が採用されている。

【先行技術文献】

【特許文献】

【0004】

【特許文献 1】特開 2005 - 108304 号公報

【特許文献 2】特開 2015 - 201050 号公報

【特許文献 3】特開 2015 - 204118 号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

今後、フラッシュメモリの技術として、高信頼・長寿命である SLC や MLC よりも、

10

20

30

40

50

書込み量の上限は下がるものの、より大容量かつより低価格を実現できるTLCといった技術が主流になると想定される。ドライブを大容量化することで、ドライブ全体の寿命指標であるDWPDの値が小さくても、一日あたりの書込み可能な容量そのものが大きくなるためである。SLCはSingle Level Cellの略記であり、MLCはMultiple Level Cellの略記であり、TLCはTriple Level Cellの略記である。

【0006】

しかし、業務のワークロードに依って、ストレージ装置に対するデータアクセス量やRead/Write比率は大きく異なる。したがって、複数の業務でストレージ装置を共用するミックスワークロード環境でSSDを安心して使うためには、SSDが大容量化しても長寿命化技術は依然として重要である。

10

【0007】

そのため、DWPDの小さいSSDを長寿命化させる手段として、書込み量を削減し、一度書き込んだデータについては基本的に読み出しで使うことを可能にするアクセス方法の開発が望まれる。しかしながら、データの圧縮処理や重複排除処理を行なうだけでは、そのデータが読み出し重視(Read Intensive)となるかどうかは分からない。このため、書込み重視(Write Intensive)のデータを、DWPDの小さなSSD、つまり寿命の短いSSDに格納すると、SSDの寿命を縮めてしまう。

【0008】

一つの側面で、本件明細書に開示の発明は、記憶部の寿命を延ばすことを目的とする。

【課題を解決するための手段】

20

【0009】

本件のストレージ装置は、寿命の異なる複数の記憶部と、前記複数の記憶部に格納されるデータブロックを管理する管理部と、ホストからの、前記寿命の異なる複数の記憶部に対する書込み対象データブロックを、前記データブロックとして一時的に格納する一時記憶部と、を備える。前記管理部は、前記データブロックの重複を排除する重複排除機能を有するとともに、判定処理部及び格納処理部を有する。前記判定処理部は、前記重複排除機能によって管理される、前記一時記憶部に格納された前記データブロックの参照数に基づいて、前記データブロックに対するアクセス特性の判定または推定を行なう。前記格納処理部は、前記寿命の異なる複数の記憶部のうち、前記判定処理部によって判定された前記アクセス特性に応じた記憶部に、前記データブロックを格納する。また、前記判定処理部は、前記参照数が2以上のデータブロックの前記アクセス特性は読み出し重視であると判定する一方、前記参照数が1で書込み後所定時間経過したデータブロックの前記アクセス特性は書込み重視であると推定する。さらに、前記格納処理部は、前記判定処理部によって前記アクセス特性が読み出し重視であると判定されたデータブロックを、前記一時記憶部から、前記寿命の異なる複数の記憶部のうち、他の記憶部よりも寿命の短い記憶部に格納する一方、前記判定処理部によって前記アクセス特性が書込み重視であると推定されたデータブロックを、前記一時記憶部から、前記寿命の異なる複数の記憶部のうち、他の記憶部よりも寿命の長い記憶部に格納する。

30

【発明の効果】

【0010】

記憶部の寿命を延ばすことができる。

40

【図面の簡単な説明】

【0011】

【図1】本発明の第1実施形態としてのストレージ制御装置を含むストレージ装置のハードウェア構成の一例を示すブロック図である。

【図2】図1に示すストレージ制御装置の機能構成の一例を示すブロック図である。

【図3】図2に示すストレージ制御装置を含むストレージ装置でのデータ書込み処理を説明するフローチャートである。

【図4】第1実施形態における論理ボリュームのデータ配置を説明する図である。

【図5】第1実施形態における論理ボリュームと不揮発性デバイスとを対応付けるマッピング

50

ングテーブルの構造の一例を示す図である。

【図6】図5に示すマッピングテーブルによるマッピングの具体例を示す図である。

【図7】第1実施形態におけるSCMフリーリストの構造の一例を示す図である。

【図8】第1実施形態におけるSCMアロケーションリストの構造の一例を示す図である。

【図9】第1実施形態のストレージ制御装置による初期化処理を説明するフローチャートである。

【図10】第1実施形態のストレージ制御装置による時間調整処理を説明するフローチャートである。

【図11】第1実施形態のストレージ制御装置による移動データ抽出処理を説明するフローチャートである。

10

【図12】第1実施形態のストレージ制御装置による移動データ抽出処理を説明するフローチャートである。

【図13】第1実施形態のストレージ制御装置による移動データ抽出処理を説明するフローチャートである。

【図14】本発明の第2実施形態としてのストレージ制御装置を含むストレージ装置のハードウェア構成の一例を示すブロック図である。

【図15】図14に示すストレージ制御装置の機能構成の一例を示すブロック図である。

【図16】図15に示すストレージ制御装置を含むストレージ装置でのデータ書込み処理を説明するフローチャートである。

20

【図17】第2実施形態における論理ボリュームと不揮発性デバイスとを対応付けるマッピングテーブルの構造の一例を示す図である。

【図18】図17に示すマッピングテーブルによるマッピングの具体例を示す図である。

【図19】第2実施形態のストレージ制御装置による移動データ抽出処理を説明するフローチャートである。

【図20】第2実施形態のストレージ制御装置による移動データ抽出処理を説明するフローチャートである。

【発明を実施するための形態】

【0012】

以下に、図面を参照し、本願の開示するストレージ装置、及びストレージ制御装置の実施形態について、詳細に説明する。ただし、以下に示す実施形態は、あくまでも例示に過ぎず、実施形態で明示しない種々の変形例や技術の適用を排除する意図はない。すなわち、本実施形態を、その趣旨を逸脱しない範囲で種々変形して実施することができる。また、各図は、図中に示す構成要素のみを備えるという趣旨ではなく、他の機能を含むことができる。そして、各実施形態は、処理内容を矛盾させない範囲で適宜組み合わせることが可能である。

30

【0013】

〔1〕本実施形態の概要

本実施形態（第1及び第2実施形態）では、ホストサーバからストレージ装置に対して対象データの保存要求を行なう際、対象データのアクセス特性が判断される。ここで、アクセス特性として、対象データが読出し重視であるか否かが判断される。対象データのアクセス特性が読出し重視でないと判断された場合、対象データのアクセス特性は、書込み重視であると推定してもよい。

40

【0014】

なお、読出し重視（Read Intensive）のデータブロックは、書込みアクセスよりも読出しアクセスの対象になる可能性の高いデータブロックである。逆に、書込み重視（Write Intensive）のデータブロックは、読出しアクセスよりも書込みアクセスの対象になる可能性の高いデータブロックである。以下において、読出し重視はリード重視と称し、書込み重視はライト重視と称する場合がある。

【0015】

50

そして、ストレージ装置内のDWPD（寿命指標）の異なるSSD（記憶部）の中から、アクセス特性の判断結果に応じた最適なSSDが選択される。対象データブロックは、選択されたSSDに書き込まれ格納保存される。例えば、読出し重視であると判定された対象データブロックは、DWPDの小さい（寿命の短い）SSDに格納されてよい。一方、読出し重視でないと判定された対象データブロック、つまり書込み重視であると推定された対象データブロックは、DWPDの大きい（寿命の長い）SSDに格納されてよい。また、DWPDの大きいSSDに格納された、書込み重視であると推定されたデータブロックのアクセス特性が、その推定後に読出し重視であると判定された場合、当該データブロックは、DWPDの大きいSSDからDWPDの小さいSSDへ移動されてもよい。

【0016】

10

このように、本実施形態（第1及び第2実施形態）によれば、アクセス特性が読出し重視であると判定されるデータブロックは、他のSSDよりもDWPDの小さいSSDに書き込まれる。また、アクセス特性が読出し重視でないデータブロック、つまりアクセス特性が書込み重視であると推定されるデータブロックは、他のSSDよりもDWPDの大きいSSDに書き込まれる。したがって、データブロックのアクセス特性に応じた、寿命指標の異なるSSDのうちの最適なSSDに、データブロックを格納することができ、特に、DWPDの小さいSSDの寿命を延ばすことができる。

【0017】

例えば、本実施形態（第1及び第2実施形態）のストレージ装置には、以下のように構成されていてもよい。

20

【0018】

本実施形態（第1及び第2実施形態）のストレージ装置は、寿命（DWPD）の異なる複数のSSD（記憶部）と、当該複数のSSDに格納されるデータブロックを管理する管理部（ストレージ制御装置）と、を備える。前記管理部は、判定処理部と格納処理部とを有する。前記判定処理部は、データブロックに対するアクセス特性の判定または推定を行なう。前記格納処理部は、寿命の異なる複数のSSDのうち、前記判定処理部によって判定または推定された前記アクセス特性に応じたSSDに、データブロックを格納する。

【0019】

このとき、前記管理部は、データブロックの重複を排除する重複排除機能を有し、前記判定処理部は、前記重複排除機能によって管理される、データブロックの参照数（重複回数）に基づいて、前記アクセス特性の判定または推定を行なってもよい。また、前記判定処理部は、前記参照数が2以上のデータブロックの前記アクセス特性は読出し重視であると判定し、前記格納処理部は、前記参照数が2以上のデータブロックを、他のSSDよりも寿命の短い（DWPDの小さい）SSDに格納してもよい。さらに、前記判定処理部は、前記アクセス特性が読出し重視でないデータブロックの前記アクセス特性は書込み重視であると推定し、前記格納処理部は、前記アクセス特性が読出し重視でないデータブロックを、他のSSDよりも寿命の長いSSDに格納してもよい。

30

【0020】

特に、図1～図13を参照しながら後述する第1実施形態のストレージ装置は、ホストサーバからの、寿命の異なる複数のSSDに対する書込み対象データブロックを、前記データブロックとして一時的に格納する一時記憶部を有してもよい。一時記憶部は、例えばSCM（Storage Class Memory）であってもよい。このとき、前記判定処理部は、SCMに格納されたデータブロックに対するアクセス特性の判定または推定を行なう。そして、前記格納処理部は、前記判定処理部によって前記アクセス特性が読出し重視であると判定されたデータブロックを、SCMから、他のSSDよりも寿命の短いSSDに格納する。一方、前記判定処理部は、前記参照数が1で書込み後所定時間経過したデータブロックの前記アクセス特性は書込み重視であると推定する。前記格納処理部は、前記判定処理部によって前記アクセス特性が書込み重視であると推定されたデータブロックを、SCMから、他のSSDよりも寿命の長いSSDに格納する。

40

【0021】

50

また、図14～図20を参照しながら後述する第2実施形態のストレージ装置では、他のSSDよりも寿命の長いSSDは、ホストサーバからの、寿命の異なる複数のSSDに対する書込み対象データブロックを、前記データブロックとして格納する。このとき、前記判定処理部は、他のSSDよりも寿命の長いSSDに格納されたデータブロックに対するアクセス特性の判定または推定を行なう。そして、前記格納処理部は、前記判定処理部によって前記アクセス特性が読出し重視であると判定されたデータブロックを、他のSSDよりも寿命の長いSSDから、他のSSDよりも寿命の短いSSDに格納する。

【0022】

(2) 第1実施形態

(2-1) ハードウェア構成

図1を参照しながら、本発明の第1実施形態としてのストレージ制御装置10を含むストレージ装置1のハードウェア構成の一例について説明する。図1は、そのハードウェア構成の一例を示すブロック図である。

【0023】

図1に示すように、第1実施形態のストレージ装置1は、ストレージ制御装置10とSSD部20とを有する。ストレージ制御装置10は、ホストサーバ30からの入出力要求に応じて、SSD部20に格納されるデータブロックを管理する。SSD部20には、ストレージ制御装置10によって制御される、DWPDの異なる(つまり寿命の異なる)複数のSSD21, 22が混在する。SSD21は、例えば、寿命の長いつまりDWPDの大きいSSDであり、SSD22は、例えば、SSD21に対し、寿命の短いつまりDWPDの小さいSSDである。

【0024】

第1実施形態のストレージ制御装置10は、CPU11、DRAM12、SCM13、SSDコントローラ14、及びホストインタフェース15を含む。ここで、CPUはCentral Processing Unitの略記であり、DRAMはDynamic Random Access Memoryの略記である。

【0025】

CPU11は、種々の制御や演算を行なうプロセッサ(処理部)の一例である。CPU11は、DRAM12、SCM13、SSDコントローラ14、及びホストインタフェース15と相互に通信可能に接続される。なお、プロセッサとしては、CPU11等の演算処理装置に代えて、電子回路、例えばMPU(Micro Processing Unit)、ASIC(Application Specific Integrated Circuit)、FPGA(Field Programmable Gate Array)等の集積回路(IC)が用いられてもよい。

【0026】

また、CPU11は、DRAM12やSCM13などの記憶部に格納されたOS(Operating System)やプログラムを実行することにより、種々の機能を果たす。特に、図2を参照しながら後述するごとく、第1実施形態のCPU11は、データキャッシング処理部111、データ圧縮・重複排除処理部112、SCMライトバッファ管理部113、データアクセス特性判定処理部114、データ移動制御部115、SSDプール管理部116、及びSSDドライバ117としての機能を果たしてもよい。

【0027】

DRAM12は、種々のデータやプログラムを一時的に格納する記憶装置であり、キャッシュ領域やメモリ領域を備える。キャッシュ領域は、ホストサーバ30から受信した書込み対象のデータや、ホストサーバ30に対して送信する読出し対象のデータを、一時的に格納する。メモリ領域は、CPU11がアプリケーションプログラムを実行する際に、データやプログラムを一時的に格納してもよい。アプリケーションプログラムは、例えば、本実施形態のストレージ制御機能を実現すべくCPU11が実行するプログラムであってもよい。

【0028】

SCM13は、ライトバッファ用として備えられ、DRAM12とSSD21, 22と

10

20

30

40

50

の中間のアクセス性能を有する。つまり、S C M 1 3 は、S S D 2 1 , 2 2 よりも高速のアクセス性能を有し、D R A M 1 2 は、S C M 1 3 よりも高速のアクセス性能を有する。S C M 1 3 は、ホストサーバ 3 0 からの、寿命の異なる複数の S S D 2 1 , 2 2 に対する書込み対象データブロックを、一時的に格納する一時記憶部の一例である。

【 0 0 2 9 】

S S D コントローラ 1 4 は、C P U 1 1 からの指示に従って S S D 部 2 0 における各 S S D 2 1 , 2 2 を制御する。S S D コントローラ 1 4 は、図 2 を参照しながら後述する、S S D プール管理部 1 1 6、及び S S D ドライバ 1 1 7 としての機能を果たしてもよい。

【 0 0 3 0 】

ホストインタフェース (Host I/F) 1 5 は、F C、S A S、E t h e r n e t (登録商標) などのハードウェアを用いて、ホストサーバ 3 0 と通信可能に接続される。このようなハードウェアを通信路として、本実施形態のストレージ装置 1 は、ホストサーバ 3 0 との間でデータの送受信を行なう。なお、F C は Fibre Channel の略記であり、S A S は Serial Attached S C S I の略記であり、S C S I は Small Computer System Interface の略記である。

【 0 0 3 1 】

ホストサーバ 3 0 からストレージ装置 1 (ストレージ制御装置 1 0) に対しデータ書込み要求が行なわれた場合、ホストサーバ 3 0 からの書込み対象のデータは、ホストインタフェース 1 5 を経由して D R A M 1 2 に一時的に格納される。このとき、ライトスルー (Write Through) 処理を採用した場合、D R A M 1 2 に格納されたデータが、装置電源オフ時でもデータ保持可能な不揮発性メディアである S C M 1 3 または S S D 2 1 , 2 2 に書き込まれるまで、ホストサーバ 3 0 へ書込み完了応答が返信されない。一方、ライトバック (Write Back) 処理を採用した場合、書込み対象のデータが D R A M 1 2 に格納された時点で、ホストサーバ 3 0 へ書込み完了応答が返され、その後、非同期に不揮発性メディアにデータ (データブロック) が移動される。

【 0 0 3 2 】

ライトバック処理の採用時には、ハードウェア故障などにより D R A M 1 2 上のデータが消失する場合を考慮し、ストレージ制御装置 1 0 をクラスタ構成とし、他のストレージ制御装置 1 0 とデータミラーリングすることが行なわれる場合がある。本実施形態において、ライトスルー処理とライトバック処理との違いはホストサーバ 3 0 への書込み完了応答を行なうタイミングの違いのみである。また、D R A M 1 2 上のデータのクラスタ間ミラーリング処理は本願技術の趣旨と関係ない。したがって、本実施形態では、ストレージ制御装置 1 0 が、クラスタ構成を採らないシングルノード構成であり、且つライトバック処理を採用する場合について説明する。なお、本実施形態において、読出し処理は、論理ボリュームのアドレスから読出し対象データの物理格納場所を特定し、読出し対象データをホストサーバ 3 0 に転送する処理であるので、本実施形態の説明では省略する。

【 0 0 3 3 】

〔 2 - 2 〕 機能構成

ついで、図 2 を参照しながら、第 1 実施形態のストレージ制御装置 (管理部) 1 0 の機能構成について説明する。なお、図 2 は、図 1 に示すストレージ制御装置 1 0 の機能構成の一例を示すブロック図である。

【 0 0 3 4 】

第 1 実施形態のストレージ制御装置 1 0 は、D W P D の異なる (つまり寿命の異なる) 複数の S S D 2 1 , 2 2 が混在するストレージ装置 1 において、複数の S S D 2 1 , 2 2 に格納されるデータブロックを管理する。ストレージ制御装置 1 0 において、C P U 1 1 は、プログラムを実行することで、図 2 に示すように、データキャッシング処理部 1 1 1、データ圧縮・重複排除処理部 1 1 2、S C M ライトバッファ管理部 1 1 3、データアクセス特性判定処理部 1 1 4、データ移動制御部 1 1 5、S S D プール管理部 1 1 6、及び S S D ドライバ 1 1 7 として機能してもよい。前述したように、S S D コントローラ 1 4 が、S S D プール管理部 1 1 6、及び S S D ドライバ 1 1 7 としての機能を果たしてもよ

10

20

30

40

50

い。

【0035】

なお、プログラムは、コンピュータ読取可能な記録媒体であって可搬型の非一時的な記録媒体に記録された形態で提供される。当該記録媒体としては、磁気ディスク、光ディスク、光磁気ディスクなどが挙げられる。また、光ディスクとしては、CD (Compact Disk)、DVD (Digital Versatile Disk)、ブルーレイディスクなどが挙げられる。CDは、CD-ROM (Read Only Memory)、CD-R (Recordable) / RW (ReWritable)などを含む。DVDは、DVD-RAM、DVD-ROM、DVD-R、DVD+R、DVD-RW、DVD+RW、HD (High Definition) DVDなどを含む。

【0036】

このとき、CPU 11は、上述のごとき記録媒体からプログラムを読み取って内部記憶装置または外付けの記憶装置に格納して用いる。CPU 11は、プログラムを、ネットワークを介して受信し内部記憶装置または外付けの記憶装置に格納して用いてもよい。

【0037】

ホストインタフェース15は、前述した通りホストサーバ30と通信可能に接続され、ホストサーバ30との間でデータの送受信を行なう。

【0038】

データキャッシング処理部111は、ホストサーバ30との間でデータの転送を高速に行なうためのデータキャッシュ機能を果たす。データキャッシング処理部111は、ホストサーバ30から転送された書込み対象データを、DRAM 12上で保持し、実格納場所であるSSD 21、22へ書き込む前にホストサーバ30に対して書込み完了応答を行なう。これにより、データキャッシング処理部111は、書込み処理性能を高速化させるバッファ処理を行なう。また、データキャッシング処理部111は、ホストサーバ30から読出し要求を受けたデータを、ホストサーバ30への転送後もDRAM 12で保持し、再度同じデータに対する読出し要求を受けた場合に高速応答を可能にするキャッシュ処理を行なう。さらに、データキャッシング処理部111は、DRAM 12にキャッシュされているデータブロックの有無を判定するキャッシュヒット判定処理と、長期間アクセスされないデータブロックをDRAM 12から破棄するLRU (Least Recently Used) 制御機能とも有する。なお、書込み要求に対して完了応答を返すタイミングは、DRAM 12へのキャッシング時点 (ライトバック方式) と、SCM 13またはSSD 21、22等の不揮発性メディアにデータブロックを格納した時点 (ライトスルー方式) とのいずれか一方を選択可能である。

【0039】

データ圧縮・重複排除処理部112は、ホストサーバ30等から受信するデータ量を削減する。データ圧縮・重複排除処理部112は、ホストサーバ30から転送される書込み対象データをブロック長単位 (通常4 Kbyteあるいは8 Kbyte) のデータブロックに分割し、分割後のデータブロックに対し、以下のような重複排除処理や圧縮処理を実行する。なお、書込み対象データのデータ長がブロック長単位に満たない場合、当該書込み対象データに対し、0データをパディングすることで、ブロック長単位のデータブロックが作成される。

【0040】

重複排除処理では、ストレージ装置1内に同じデータブロックが格納されているか否かが判定され、既に同じデータブロックが存在すれば、データブロックの格納は行なわれずにメタ情報が更新される。ここで、メタ情報には、論理ボリュームのLBA (Logical Block Address) 情報と物理SSDまたはSCMの実格納場所とをマッピングするマッピング情報が含まれるほか、後述する参照数が含まれる。圧縮処理では、データブロックが可逆形式で縮小・圧縮される。

【0041】

また、データ圧縮・重複排除処理部112は、ホストサーバ30からの読出し要求時、重複排除されたデータブロックを復元する機能や、圧縮されているデータブロックを伸長

10

20

30

40

50

する機能も有する。そして、データ圧縮・重複排除処理部 1 1 2 は、重複排除処理時に、後述する S S D プール管理部 1 1 6 と連携し、ホストサーバ 3 0 から L U N (Logical Unit Number) によって参照される論理ボリュームと物理 S S D 2 1 , 2 2 または S C M 1 3 との対応関係を記録するメタ情報や、後述する参照数を管理する。

【 0 0 4 2 】

データ圧縮・伸長論理や重複排除論理に関しては、一般に使用されており、本実施形態ではその説明は省略する。また、本実施形態において、重複排除処理と圧縮処理との実行順序、及び圧縮処理自体の有効・無効は関係ない。しかし、本実施形態では、重複排除処理で管理されるデータブロックの参照数は、データブロックのアクセス特性の判定・推定に利用されるため、データ圧縮・重複排除処理部 1 1 2 における重複排除機能は有効にする。参照数は、リファレンスカウンタ (Reference Counter) 値と表記されてもよいし、重複排除参照数あるいは重複回数と表記されてもよい。

10

【 0 0 4 3 】

S C M ライトバッファ管理部 1 1 3 は、D R A M 1 2 よりもアクセス性能は劣るが D R A M 1 2 よりも安価かつ大容量で不揮発性を有し、かつ S S D 2 1 , 2 2 よりも高価ではあるが S S D 2 1 , 2 2 よりも高速である S C M 1 3 を利用する。つまり、S C M ライトバッファ管理部 1 1 3 は、D R A M 1 2 と S S D 2 1 , 2 2 との中間のアクセス性能を有する S C M 1 3 で、ホストサーバ 3 0 からの書き込み対象データ (データブロック) を S S D 2 1 , 2 2 に格納する前にバッファリングする。これにより、ストレージ性能の向上が図られる。

20

【 0 0 4 4 】

また、第 1 実施形態では、実際に S S D 2 1 , 2 2 に格納する前に、S C M 1 3 上でデータブロックのアクセス特性の判定を行なうことで、S S D プール (S S D 群) 間のデータブロックの移動量を削減することが可能になる。S S D 部 2 0 には、D W P D の大きい S S D 2 1 が属する S S D プールと、D W P D の小さい S S D 2 2 が属する S S D プールとが含まれる (図 4 参照) 。

【 0 0 4 5 】

なお、S C M 1 3 を具備しないストレージ装置 1 0 A について、第 2 実施形態において図 1 4 ~ 図 2 0 を参照しながら説明する。つまり、第 1 実施形態では、S C M 1 3 を利用したデータ最適配置手法について説明し、第 2 実施形態では、S C M 1 3 を利用しないデータ最適配置手法について説明する。そのため、第 2 実施形態では、第 1 実施形態における S C M ライトバッファ管理部 1 1 3 が存在しない (図 1 5 参照) 。

30

【 0 0 4 6 】

データアクセス特性判定処理部 1 1 4 は、定期的にデータブロックの移動判定を行なうために、データブロックに対するアクセス特性の判定または推定を行なう判定処理部に相当する。データアクセス特性判定処理部 1 1 4 は、判定処理部 1 1 4 と略記する場合がある。特に、データアクセス特性判定処理部 1 1 4 は、ホストサーバ 3 0 から書き込まれるデータ (分割後のデータブロック) のアクセス特性が読出し重視 (Read Intensive) であるか否かを判断する。本実施形態では、読出し重視のデータブロックを D W P D の小さい (寿命の短い) S S D 2 2 に配置することで、寿命の短い S S D 2 2 の寿命を延ばすことを目的としている。このため、データアクセス特性判定処理部 1 1 4 によって、読出し重視のデータブロックが見い出される。ある一時点において、読出し重視ではないデータブロックは、全て書込み重視 (Write Intensive) として扱われる。その後の時間経過に伴い、先の判定タイミングにおいて書込み重視と判断されたデータブロックが、のちの判定タイミングで読出し重視であると判断されることは起き得る。

40

【 0 0 4 7 】

なお、第 1 実施形態において、判定処理部 1 1 4 は、前述した重複排除処理 (重複排除機能) によって管理される、データブロックの参照数 (重複回数) に基づいて、アクセス特性の判定または推定を行なってもよい。判定処理部 1 1 4 は、例えば、データブロックの参照数が 2 以上のデータブロックのアクセス特性は読出し重視であると判定する。また

50

、判定処理部 114 は、アクセス特性が読み重視でない（参照数が 2 以上でない）データブロックのアクセス特性は書き込み重視であると推定してもよい。このとき、判定処理部 114 は、参照数が 1 で書き込み後所定時間経過したデータブロックのアクセス特性は書き込み重視であると推定してもよい。

【0048】

データ移動制御部 115 は、データブロックの移動を実行するもので、寿命の異なる複数の SSD のうち、判定処理部 114 によって判定または推定されたアクセス特性に応じた SSD 21 または 22 に、データブロックを格納する格納制御部の一例である。特に、データ移動制御部 115 は、判定処理部 114 によるアクセス特性判定で読み重視と判断されたデータブロックを、SCM 13 から、DWPD の小さい SSD 22 へ移動させる処理を担う。SCM ライトバッファ管理部 113 の有無により、データ移動処理の動作が異なるが、その差異については後述する。

10

【0049】

なお、第 1 実施形態において、データ移動制御部 115 は、寿命の異なる複数の SSD 21, 22 のうち、判定処理部 114 によって判定または推定されたアクセス特性に応じた SSD 21 または 22 に、データブロックを格納する。このとき、データ移動制御部 115 は、参照数が 2 以上のデータブロックを、他の SSD 21 よりも寿命の短い（DWPD の小さい）SSD 22 に格納してもよい。また、データ移動制御部 115 は、アクセス特性が読み重視でないデータブロックを、他の SSD 22 よりも寿命の長い SSD 21 に格納してもよい。さらに、データ移動制御部 115 は、判定処理部 114 によってアクセス特性が読み重視であると判定されたデータブロックを、SCM 13 から、他の SSD 21 よりも寿命の短い SSD 22 に格納する。一方、データ移動制御部 115 は、判定処理部 114 によってアクセス特性が書き込み重視であると推定されたデータブロックを、SCM 13 から、他の SSD 22 よりも寿命の長い SSD 21 に格納する。

20

【0050】

SSD プール管理部 116 は、ホストサーバ 30 から LUN で参照される論理ボリュームと、物理 SSD 21, 22 または SCM 13 との対応関係を、データ圧縮・重複排除処理部 112 と連携して管理する。当該管理の手法は、物理的な SSD 21, 22 からチャンクと呼ばれるブロックを確保し、当該チャンクを論理ボリュームと対応付けるシンプロビジョニング（Thin Provisioning）機能において一般に使用される手法であり、本実施形態でもその手法が利用される。また、SSD プール管理部 116 は、寿命指標である DWPD が同一または近い SSD を SSD プールにまとめる機能を有していてもよい。本実施形態では、SSD 部 20 において、例えば、DWPD の大きい SSD 21 が属する SSD プールと、DWPD の小さい SSD 22 が属する SSD プールとにまとめられる（図 4 参照）。

30

【0051】

SSD ドライバ 117 は、SSD 21, 22 との間でデータ（データブロック）の送受信を行なう。

【0052】

〔2-3〕動作

40

次に、図 3 に示すフローチャート（ステップ S1 ~ S5）に従って、図 2 に示すストレージ制御装置 10 を含むストレージ装置 1 でのデータ書き込み処理について説明する。

【0053】

まず、ホストインタフェース 15 によってホストサーバ 30 から書き込み要求とともに書き込み対象のデータが受信されると、データキャッシング処理部 111 によって、書き込み対象のデータは、DRAM 12 上でバッファリングされる（ステップ S1）。そして、バッファリング時点で、データキャッシング処理部 111 からホストインタフェース 15 経由で、ホストサーバ 30 へ書き込み完了応答が返される。

【0054】

ホストサーバ 30 に対し書き込み完了応答を行なった後、DRAM 12 上のデータは、所

50

定のブロック長単位のデータブロックに分割される。そして、データ圧縮・重複排除処理部 1 1 2 によって、データブロック毎に、圧縮処理や重複排除処理といったデータ量削減処理が実行される（ステップ S 2）。本実施形態では、データブロックの参照数に基づきデータブロックのアクセス特性の判定を行なう。このため、圧縮処理の有効/無効や、圧縮処理と重複排除処理との処理順序は、本実施形態と関係ない。したがって、本実施形態では、重複排除処理に着目して説明を行なう。

【 0 0 5 5 】

重複排除処理では、D R A M 1 2 上の書込み対象のデータブロックが重複していると判断された場合、不揮発性デバイス（S C M 1 3 または S S D 2 1 , 2 2）に格納済の重複データと同じであることが分かるよう 2 種類のメタ情報が更新される。2 種類のメタ情報は、上述したマッピング情報及びリファレンスカウンタ値である。ただし、データブロックそのものの不揮発性デバイスへのライトバック（Write Back）処理は行なわれず、メタ情報更新後、D R A M 1 2 上のデータは破棄される。データブロックが重複していないと判断された場合、ライトバック処理が行なわれる。そして、D R A M 1 2 上のデータブロックは S C M 1 3 ヘコピーされた後、当該データブロックについての 2 種類のメタ情報が更新されてから、D R A M 1 2 上のデータブロックは削除される。

【 0 0 5 6 】

ここで、図 4 は、第 1 実施形態における論理ボリュームのデータ配置を説明する図であり、論理ボリュームとその実データが配置される不揮発性デバイスとの関係を示したものである。第 1 実施形態において、論理ボリュームの実データは、S C M 1 3 あるいは S S D 2 1 , 2 2 のどちらかに格納されている。D R A M 1 2 上の書込み対象データブロックがライトバックされた直後、当該書込み対象データは、S C M ライトバッファ管理部 1 1 3 によって、一旦、S C M 1 3 に保持される（ステップ S 3）。その後、当該書込み対象データは、判定処理部 1 1 4 によって、S C M 1 3 上のデータブロックのアクセス特性が判定される（ステップ S 4）。そして、当該書込み対象データブロックは、データ移動制御部 1 1 5 によって、アクセス特性の判定結果に応じて、S S D プールの S S D 2 1 または 2 2 に移動される（ステップ S 5）。

【 0 0 5 7 】

2 種類のメタ情報である、マッピング情報、及びリファレンスカウンタ値は、論理ボリュームのマッピングテーブル 1 2 0（図 5 参照）と S C M アロケーションリスト 1 2 2（図 8 参照）とを用いて、例えば D R A M 1 2 上で管理される。マッピング情報は、「論理ボリューム番号と論理ブロックアドレスとの組合せ」と、「不揮発性デバイスの番号とブロック位置との組合せ」とをマッピングする情報である。また、リファレンスカウンタ値は、一つの物理データブロックが複数の論理データブロックと対応している際の「論理データブロックの参照数」に相当する。なお、リファレンスカウンタ値は、シンプロビジョニングと呼ばれる「論理-物理アドレス変換処理」を行なう S S D プール管理部 1 1 6 においても保持される。シンプロビジョニング処理は、一般的な処理であるため、ここでは、リファレンスカウンタ値についての詳細な説明は省略する。

【 0 0 5 8 】

ここで、図 5 は、第 1 実施形態における論理ボリュームと不揮発性デバイスとを対応付けるマッピングテーブル 1 2 0 の構造の一例を示す図である。マッピングテーブル 1 2 0 の各エントリには、対応するデータブロックについて、次の情報が、図 5 に示すようなテーブル形式で、例えば D R A M 1 2 上において展開保持される。

【 0 0 5 9 】

・論理アドレス（ボリューム# LBA）： ホストサーバ 3 0 からアクセスされる L U N に対応する論理ボリュームの番号と、その論理ボリューム内でのデータブロックの位置を示す L B A との組合せ情報。

【 0 0 6 0 】

・フラグ（Flag）： 論理アドレスで示されるデータブロックが格納されている不揮発性デバイスが S C M 1 3 であるか S S D 2 1 , 2 2 であるかを示す情報。第 1 実施形態で

10

20

30

40

50

は、SCM13にデータブロックが格納されている時にフラグは1に設定される一方、SSD21または22にデータブロックが格納されている時にフラグは0に設定される。

【0061】

・不揮発性デバイスアドレス (SCM#-Page#またはSSD-Pool#-block#-Offset) : フラグの状態が1の時、SCM13上におけるデータブロックの位置の情報である一方、フラグの状態が0の時、データブロックのSSDプール内の論理位置の情報。

【0062】

なお、図6は、図5に示すマッピングテーブル120によるマッピングの具体例を示す図であり、図5に示すマッピングテーブル120によるマッピングを概念化したものである。

【0063】

図3のステップS3におけるSCMライトバッファ管理では、SCM13に確保されるライトバッファ領域の、SCMフリーリスト121 (図7参照)とSCMアロケーションリスト122 (図8参照)とが管理される。ライトバッファ領域は、ホストサーバ30からの書き込み対象データを、最終格納先であるSSD21, 22へ移動させるまで、一時的に格納する領域である。各々のリスト121, 122は、リンクリストとしてキュー管理される。DRAM12上からの書き込み対象データブロックをライトバックする際、SCMフリーリスト121で参照される空き領域 (フリー領域) が、ライトバック先の領域として選択されSCMフリーリスト121からデキューされ、SCMアロケーションリスト122にエンキューされる。なお、図7は、第1実施形態におけるSCMフリーリスト121の構造の一例を示す図である。図8は、第1実施形態におけるSCMアロケーションリスト122の構造の一例を示す図である。

【0064】

SCMフリーリスト121とSCMアロケーションリスト122との構造は基本的に同じで、SCMフリーリスト121は、SCM13における空き領域のリンクリストであり、SCMアロケーションリスト122は、SCM13において、論理ボリュームのデータが格納されている領域についてのリンクリストである。SCMフリーリスト121及びSCMアロケーションリスト122における各エントリには、次の情報が、それぞれ図7及び図8に示すようなテーブル形式で、例えばDRAM12上に展開保持される。

【0065】

・ヘッダ (Header) : エントリの識別情報とキュー操作時の排他処理用ロックフラグとを含む。図8に示す例では、SCM13のページ番号を利用したシーケンシャル番号が、エントリの識別情報として用いられている。

【0066】

・前エントリポインタ (Previous Entry) : リンクリストの直前のエントリをポイントする。エントリが先頭エントリの場合、アンカ (Free Top AnchorまたはAllocated Top Anchor) をポイントする。例えば、図7において先頭のPrevious EntryはFree Top Anchorをポイントし、図8において先頭のPrevious EntryとしてのTop AnchorはAllocated Top Anchorをポイントする。

【0067】

・次エントリポインタ (Next Entry) : リンクリストの直後のエントリをポイントする。エントリが最後尾エントリの場合、アンカ (Free Bottom AnchorまたはAllocated Bottom Anchor) をポイントする。例えば、図7において最後尾のNext EntryはFree Bottom Anchorをポイントし、図8において最後尾のNext EntryとしてのBottom AnchorはAllocated Bottom Anchorをポイントする。

【0068】

・SCMアドレス : 複数のSCM13が搭載された場合のSCM番号と、各SCM13内でページ単位に分割された領域のシーケンシャル番号との組合せであり、物理領域を特定している。なお、SCMフリーリスト121におけるSCMアドレスは、空き領域の位置を特定する上記組合せの情報である。

10

20

30

40

50

【 0 0 6 9 】

・リファレンスカウンタ： エントリに対応するデータブロックが複数の論理アドレスから参照されている場合に、その参照数（重複回数）を示す値。この値は、データ圧縮・重複排除処理部 1 1 2 の重複排除機能によって管理される。当該データブロックが重複排除されていないデータブロックつまり重複していないデータブロックである場合、リファレンスカウンタの値は 1 になる。図 8 に示す例では、シーケンシャル番号 # 2 0 のエントリの参照数は 3 であり、シーケンシャル番号 # 8 0 のエントリの参照数は 2 であり、シーケンシャル番号 # 1 5 のエントリの参照数は 1 である。

【 0 0 7 0 】

・T O D (Time Of Day)： エントリが S C M アロケーションリスト 1 2 2 にエンキューされた時刻。

10

【 0 0 7 1 】

・データ長： S C M 1 3 に格納されている、エントリに対応するデータブロックのサイズ（所定のブロック長）。データの圧縮処理を行わない場合、データブロックのサイズは、固定長で、S C M 1 3 を分割しているページサイズと同じサイズである。図 8 に示す例では、各エントリのデータ長には 8 1 9 2 が設定されている。

【 0 0 7 2 】

図 3 のステップ S 2 における重複排除処理では、ホストサーバ 3 0 から書き込まれた D R A M 1 2 上のデータブロックが、既に書き込み済みの他のデータブロックと重複すると判断され、且つ当該データブロックが S C M 1 3 上に存在している場合、以下の処理が行なわれる。つまり、当該データブロックに対応するマッピングテーブル 1 2 0 のフラグに 1 が設定されている場合、S C M アロケーションリスト 1 2 2 のリファレンスカウンタ値が 1 増加される。そして、D R A M 1 2 上のデータブロックのライトバック処理は行なわれることなく、D R A M 1 2 上のデータブロックは削除される。

20

【 0 0 7 3 】

なお、データ移動制御部 1 1 5 が後で実行する S S D 2 1 , 2 2 へのデータ移動処理によって、S C M 1 3 上のデータブロックの S S D 2 1 , 2 2 への移動を完了した時点で、S C M ライトバッファ管理部 1 1 3 は、当該データブロックに対応する S C M アロケーションリスト 1 2 2 のエントリを S C M フリーリスト 1 2 1 に戻す。

【 0 0 7 4 】

次に、図 3 のステップ S 4 におけるデータアクセス特性判定処理は、インターバル時間調整処理と移動データ抽出処理とを含む。インターバル時間調整処理は、データアクセス特性判定処理部 1 1 4 が動作する時間間隔を調整する処理である。移動データ抽出処理は、データアクセス特性判定処理部 1 1 4 によって実行され、S C M 1 3 から D W P D の大きい S S D 2 1 へ移動すべきデータブロックと、D W P D の小さい S S D 2 2 へ移動すべきデータブロックとを抽出する処理である。

30

【 0 0 7 5 】

インターバル時間調整処理においては、例えば図 9 に示す初期化処理と、例えば図 1 0 に示す時間調整処理とが含まれる。ここで、初期化処理は、装置稼働開始時及び動作モード設定変更時に実行される処理であり、時間調整処理は、インターバル時間毎に実行される処理である。また、移動データ抽出処理においては、例えば図 1 1 ~ 図 1 3 に示す処理が実行される。

40

【 0 0 7 6 】

ここで、図 9 に示すフローチャート（ステップ S 1 1 ~ S 1 5 ）に従って、第 1 実施形態のストレージ制御装置 1 0 による初期化処理について説明する。

【 0 0 7 7 】

初期化処理では、インターバル時間として固定時間が指定されている場合、つまり自動調整が選択されていない場合（ステップ S 1 1 の N O ルート）、指定される固定時間が設定され（ステップ S 1 2 ）、当該固定時間でインターバルタイマが起動され（ステップ S 1 5 ）、図 1 0 に示す処理が実行される。

50

【 0 0 7 8 】

一方、インターバル時間の自動調整が選択されている場合（ステップ S 1 1 の Y E S ルート）、インターバル時間のデフォルト値が設定される（ステップ S 1 3）。また、次回起動時に行なわれる自動調整で使用されるデータ移動数のデフォルト値が設定記憶される（ステップ S 1 4）。そして、インターバル時間のデフォルト値でインターバルタイムが起動され（ステップ S 1 5）、図 1 0 に示す処理が実行される。

【 0 0 7 9 】

ついで、図 1 0 に示すフローチャート（ステップ S 2 1 ~ S 2 8）に従って、第 1 実施形態のストレージ制御装置 1 0 による時間調整処理について説明する。

【 0 0 8 0 】

時間調整処理では、まず、図 1 1 ~ 図 1 3 を参照しながら後述する移動データ抽出処理が実行される（ステップ S 2 1）。そして、インターバル時間として固定時間が指定されている場合、つまり自動調整が選択されていない場合（ステップ S 2 2 の N O ルート）、当該固定時間でインターバルタイムが再起動され（ステップ S 2 8）、図 1 0 に示す処理が再実行される。これにより、固定時間間隔で移動データ抽出処理（ステップ S 2 1）が実行される。

【 0 0 8 1 】

一方、インターバル時間の自動調整が選択されている場合（ステップ S 2 2 の Y E S ルート）、移動データ抽出処理（ステップ S 2 1）で抽出される、小 D W P D の S S D 2 2 へ移動させるデータブロック数に関して、前回の数と今回の数とが比較される（ステップ S 2 3）。つまり、今回の移動データ抽出処理で抽出された、小 D W P D の S S D 2 2 へ移動させるデータブロック数と、前回の移動データ抽出処理で抽出された小 D W P D の S S D 2 2 へ移動させるデータブロック数とが比較される。なお、D W P D の小さい S S D 2 2 は、小 D W P D の S S D 2 2 と表記し、D W P D の大きい S S D 2 1 は、大 D W P D の S S D 2 1 と表記する場合がある。

【 0 0 8 2 】

比較の結果、前回のデータブロック数が今回のデータブロック数よりも大きい場合（ステップ S 2 4 の Y E S ルート）、次回のインターバル時間は縮小される（ステップ S 2 5）。例えば、今回のインターバル時間とインターバル時間の下限値との和を所定係数 2 で除算した値が、次回のインターバル時間として算出される。これにより、インターバル時間は、下限値を下回ることなく縮小される。

【 0 0 8 3 】

一方、前回のデータブロック数が今回のデータブロック数以下である場合（ステップ S 2 4 の N O ルート）、次回のインターバル時間は拡大される（ステップ S 2 6）。例えば、今回のインターバル時間とインターバル時間の上限値との和を所定係数 2 で除算した値が、次回のインターバル時間として算出される。これにより、インターバル時間は、上限値を上回ることなく拡大される。

【 0 0 8 4 】

次回のインターバル時間の算出後、今回の移動データ抽出処理（ステップ S 2 1）で抽出された、小 D W P D の S S D 2 2 へ移動させるデータブロック数は、次回の比較の際における前回のデータブロック数として記憶される（ステップ S 2 7）。この後、ステップ S 2 5 または S 2 6 において算出された次回のインターバル時間でインターバルタイムが再起動され（ステップ S 2 8）、図 1 0 に示す処理が再実行される。これにより、インターバル時間間隔で移動データ抽出処理（ステップ S 2 1）が実行される。

【 0 0 8 5 】

このように、第 1 実施形態においては、図 1 0 のステップ S 2 3 ~ S 2 6 の処理によってインターバル時間が自動調整される。これにより、S C M 1 3 に一時的に保存される書込み対象データ（データブロック）を S C M 1 3 から溢れさせることなく、S C M 1 3 の領域を効率的に利用することができる。

【 0 0 8 6 】

10

20

30

40

50

次に、図11～図13に示すフローチャートに従って、第1実施形態のストレージ制御装置10による移動データ抽出処理について説明する。

【0087】

まず、図11に示すフローチャート(ステップS31～S33)に従って、第1実施形態のストレージ制御装置10による移動データ抽出処理の全体的な流れについて説明する。データブロックの移動先のSSDプール(SSD群;図4参照)は、DWPDの値に応じてまとめられる複数のSSD21または22を含む。このとき、当該SSDプールへのアクセスは、チャンクと呼ばれるブロック単位で行なわれる。第1実施形態の移動データ抽出処理は、後述するように、SCM13上に格納されている移動すべきデータブロックを、チャンク単位毎にまとめる処理である。

10

【0088】

第1実施形態の移動データ抽出処理は、図11に示す3つの処理(ステップS31～S33)を含む。最初の処理(ステップS31)では、SCMアロケーションリスト122のエントリが、リファレンスカウンタ値をキーとして昇順にソートされる。そして、リファレンスカウンタ値が1であるエントリは、TODの古いものから順にソートされる。

【0089】

ステップS31でのソート処理後、リファレンスカウンタ値が2以上のデータブロックをDWPDの小さいSSDプール(SSD22)に移動させるための抽出処理が行なわれる(ステップS32)。ここで、リファレンスカウンタ値が2以上であるということは、複数の論理データブロックから参照されている重複データブロックであることを意味する。重複排除の仕組みから、その参照元の論理データが更新された場合でもリファレンスカウンタ値が減るだけで物理データブロックは変更されない。つまり、リファレンスカウンタ値が2以上のデータブロックは、その値が1になるまでは物理データブロックはリードオンのデータであることが保障される。このリファレンスカウンタ値が2以上のデータをまとめたチャンクは、読み重視(Read Intensive)のブロックであり、DWPDの小さいSSDプールに格納すべきデータである。図11のステップS32で実行される処理については、図12を参照しながら後述する。

20

【0090】

DWPDの小さいSSDプール(SSD22)へ移動するデータブロックの抽出完了後、DWPDの大きいSSDプール(SSD21)へ移動するデータを抽出する処理として、リファレンスカウンタ値が1のエントリを調べる処理が実行される(ステップS33)。図11のステップS33で実行される処理については、図13を参照しながら後述する。

30

【0091】

図12に示すフローチャート(ステップS41～S50)に従って、図11のステップS32で実行される、リファレンスカウンタ値が2以上のデータブロックの抽出処理について説明する。

【0092】

ここで、SCMアロケーションリスト122のエントリが前処理(ステップS31の処理)によってリファレンスカウンタ値をキーとして昇順にソートされている。そこで、図12に示す処理では、各エントリのデータ長に基づきチャンク残サイズを算出しながら、ある移動チャンク(ブロック)に入るだけ、エントリに対応するデータブロックが登録される。そして、当該移動チャンクがフルの状態になった場合、次の移動チャンクにデータブロックが登録される。この処理を、リファレンスカウンタ値が2以上のデータブロックが存在する限り、繰り返すことによって、DWPDの小さいSSDプールへ移動させるデータが抽出される。

40

【0093】

具体的には、図12に示すように、まず、ソート後のSCMアロケーションリスト122が参照され、リファレンスカウンタ値が2以上のエントリがあるか否かが判定される(

50

ステップS 4 1)。リファレンスカウンタ値が2以上のエントリがない場合(ステップS 4 1のNOルート)、リファレンスカウンタ値が2以上のデータブロックの抽出処理は終了する。

【0094】

一方、リファレンスカウンタ値が2以上のエントリがある場合(ステップS 4 1のYESルート)、リファレンスカウンタ値が2以上のデータブロックを登録する空の移動チャンクが設定される(ステップS 4 2)。また、チャンク残サイズとして、当該移動チャンクのチャンクサイズが設定される(ステップS 4 3)。

【0095】

この後、SCMアロケーションリスト122から、リファレンスカウンタ値が2以上のデータブロックが一つ選択され(ステップS 4 4)、選択したエントリのデータ長がチャンク残サイズよりも小さいか否かが判定される(ステップS 4 5)。選択したエントリのデータ長がチャンク残サイズよりも小さい場合(ステップS 4 5のYESルート)、選択したエントリに対応するデータブロックは、移動チャンクに登録される(ステップS 4 6)。この後、チャンク残サイズとして、チャンクサイズから選択したエントリのデータ長を減算した値が設定される(ステップS 4 7)。

10

【0096】

そして、リファレンスカウンタ値が2以上のエントリがあるか否かが判定される(ステップS 4 8)。リファレンスカウンタ値が2以上のエントリがない場合(ステップS 4 8のNOルート)、リファレンスカウンタ値が2以上のデータブロックの抽出処理は終了する。リファレンスカウンタ値が2以上のエントリがある場合(ステップS 4 8のYESルート)、特性判定処理部114はステップS 4 4の処理に戻る。

20

【0097】

また、選択したエントリのデータ長がチャンク残サイズ以上である場合(ステップS 4 5のNOルート)、新たに空の移動チャンクが設定される(ステップS 4 9)。また、チャンク残サイズとして、当該移動チャンクのチャンクサイズが設定されてから(ステップS 5 0)、特性判定処理部114はステップS 4 5の処理に戻る。

【0098】

図13に示すフローチャート(ステップS 5 1~S 6 0)に従って、図11のステップS 3 3で実行される、リファレンスカウンタ値が1のデータブロックの抽出処理について説明する。

30

【0099】

ここで、リファレンスカウンタ値が1であるということは、参照元のデータブロックが更新された場合、ポイントされている物理データブロックも更新されるということである。あるいは、別領域に更新データブロックが格納され、その物理データブロックが論理データブロックから新たにポイントされると同時に、ポイントされていた古い物理データブロックが無効化される場合もある。一方、時間の経過に伴い、重複データブロックが増え、リファレンスカウンタ値が1からもっと大きな値数に変わることもある。したがって、リファレンスカウンタ値が1のデータブロックは将来リードオンのデータブロックになる可能性もある。このため、ある程度の長期間に亘ってリファレンスカウンタ値を監視することが望ましい。しかし、SCM13の容量も有限であるため、第1実施形態では、ある規定時間が経過した、リファレンスカウンタ値が1のデータブロックは、DWP Dの大きいSSDプールへ移動させる。

40

【0100】

リファレンスカウンタ値が1のデータブロック(エントリ)は、前処理(ステップS 3 1)によってSCM13への登録時刻であるTODをキーとして古い順にソートされている。このソート結果が参照され、SCM13への登録後、規定時間以上経過しているデータブロック(エントリ)が抽出される。図13に示す処理でも、各エントリのデータ長に基づきチャンク残サイズを算出しながら、ある移動チャンク(ブロック)に入るだけ、エントリに対応するデータブロックが登録される。そして、当該移動チャンクがフルの状態

50

になった場合、次の移動チャンクにデータブロックが登録される。この処理を、規定時間以上の時間が経過しているデータブロックが存在する限り、繰り返すことによって、DWP Dの大きいSSDプールへ移動させるデータが抽出される。

【0101】

具体的には、図13に示すように、まず、ソート後のSCMアロケーションリスト122が参照され、リファレンスカウンタ値が1のエントリのうち登録後規定時間以上の時間が経過しているエントリがあるか否かが判定される(ステップS51)。該当するエントリがない場合(ステップS51のNORoot)、リファレンスカウンタ値が1のデータブロックの抽出処理は終了する。

【0102】

一方、該当するエントリがある場合(ステップS51のYESルート)、当該エントリに対応するデータブロックを登録する空の移動チャンクが設定される(ステップS52)。また、チャンク残サイズとして、当該移動チャンクのチャンクサイズが設定される(ステップS53)。

【0103】

この後、SCMアロケーションリスト122から、リファレンスカウンタ値が1で登録後規定時間以上の時間が経過しているデータブロックが一つ選択され(ステップS54)、選択したエントリのデータ長がチャンク残サイズよりも小さいか否かが判定される(ステップS55)。選択したエントリのデータ長がチャンク残サイズよりも小さい場合(ステップS55のYESルート)、選択したエントリに対応するデータブロックは、移動チャンクに登録される(ステップS56)。この後、チャンク残サイズとして、チャンクサイズから選択したエントリのデータ長を減算した値が設定される(ステップS57)。

【0104】

そして、リファレンスカウンタ値が1のエントリのうち登録後規定時間以上の時間が経過しているエントリがあるか否かが判定される(ステップS58)。該当するエントリがない場合(ステップS58のNORoot)、リファレンスカウンタ値が1のデータブロックの抽出処理は終了する。該当するエントリがある場合(ステップS58のYESルート)、特性判定処理部114はステップS54の処理に戻る。

【0105】

また、選択したエントリのデータ長がチャンク残サイズ以上である場合(ステップS55のNORoot)、新たに空の移動チャンクが設定される(ステップS59)。また、チャンク残サイズとして、当該移動チャンクのチャンクサイズが設定されてから(ステップS60)、特性判定処理部114はステップS55の処理に戻る。

【0106】

最後に、図3のステップS5におけるSSD21, 22へのデータ移動処理では、図3のステップS4におけるデータアクセス特性判定処理によって抽出された、SSD21または22へ移動すべきデータブロック(移動チャンク)を、SSD21または22へ移動させる処理が実行される。SCM13からSSDプールへのチャンク単位でのデータ移動は、シンプロビジョニングと呼ばれる一般的なデータ管理手法であるので、ここでは、データ移動処理そのものについては言及しない。

【0107】

SCM13上のデータブロックのSSD21, 22への移動が完了すると、図3のステップS5におけるSSD21, 22へのデータ移動処理では、メタデータが更新されるため、SCMライトバッファ管理部113へ次の二つの処理が依頼される。一つ目の処理は、移動させたデータブロックを管理していたSCMアロケーションリスト122のエントリを、SCMフリーリスト121へ戻す処理である。二つ目の処理は、移動させたデータブロックをポイントしているマッピングテーブル120の不揮発性デバイスアドレスの内容を、SCM13の位置情報からSSD21, 22の位置情報に置き換え、フラグの値を、SSDを示す値0に変更する処理である。

【0108】

10

20

30

40

50

以上のように、第1実施形態によれば、重複排除処理で管理されるリファレンスカウンタの値に着目してデータブロックを格納するSSD21, 22が振り分けられる。これにより、データブロックのアクセス特性（読出し重視か書込み重視か）に応じて寿命指標（DWPD）の異なるSSD21, 22へ最適なデータ格納が実現され、特に、DWPDの小さいSSD22の寿命を延ばすことが可能になる。

【0109】

〔3〕第2実施形態

〔3-1〕ハードウェア構成

前述したように、第1実施形態では、SCM13を利用したデータ最適配置手法について説明した。これに対し、第2実施形態では、SCM13を利用しないデータ最適配置手法について説明する。このため、図14に示すように、第2実施形態のストレージ装置1Aにおけるストレージ制御装置10Aのハードウェア構成は、SCMが省略されている点で第1実施形態のストレージ制御装置10と異なっている。図14は、本発明の第2実施形態としてのストレージ制御装置10Aを含むストレージ装置1Aのハードウェア構成の一例を示すブロック図である。なお、図14中、既述の符号と同一の符号は、同一もしくはほぼ同一の部分を示しているため、その説明は省略する。

10

【0110】

上述のように、第2実施形態のストレージ制御装置10Aでは、不揮発性デバイスとしてSSDのみを備えるハードウェア構成が採用される。第1実施形態で用いられるSCMは、SSDよりも高速なアクセスが可能であるが、SSDよりも高価であると考えられる。したがって、第2実施形態によれば、SCMを用いないため、第1実施形態と比べ、より安価な構成で第1実施形態と同様の作用効果を得ることができる。

20

【0111】

なお、第2実施形態のストレージ制御装置10Aも、第1実施形態と同様、クラスタ構成を採らないシングルノード構成であり、且つライトバック処理を採用する場合について説明する。

【0112】

〔3-2〕機能構成

ついで、図15を参照しながら、第2実施形態のストレージ制御装置（管理部）10Aの機能構成について説明する。なお、図15は、図14に示すストレージ制御装置10Aの機能構成の一例を示すブロック図である。図15中、既述の符号と同一の符号は、同一もしくはほぼ同一の部分を示しているため、その説明は省略する場合がある。

30

【0113】

第2実施形態のストレージ制御装置10Aも、第1実施形態と同様、DWPDの異なる複数のSSD21, 22が混在するストレージ装置1Aにおいて、複数のSSD21, 22に格納されるデータブロックを管理する。ストレージ制御装置10Aにおいて、CPU11は、プログラムを実行することで、図15に示すように、データキャッシング処理部111、データ圧縮・重複排除処理部112、データアクセス特性判定処理部114、データ移動制御部115、SSDプール管理部116、及びSSDドライバ117として機能してもよい。前述したように、SSDコントローラ14が、SSDプール管理部116、及びSSDドライバ117としての機能を果たしてもよい。なお、第2実施形態においても、プログラムは、第1実施形態と同様にして提供される。

40

【0114】

特に、第2実施形態では、DRAM12上のデータブロックのライトバック先は、DWPDの大きいSSDプール（寿命の長いSSD21）である。つまり、アクセス特性が不明なデータブロックは、DWPDの大きいSSDプールに格納しておく。そして、第2実施形態のデータアクセス特性判定処理部114は、SSD21上のデータブロックについて定期的に移動データ抽出処理を行なう。移動データ抽出処理によって抽出されたデータブロックは、第2実施形態のデータ移動制御部115によって、DWPDの大きいSSDプールからDWPDの小さいSSDプールへチャンク単位で移動される。

50

【 0 1 1 5 】

なお、第2実施形態で行なわれる、DWP Dの大きいSSD 2 1からDWP Dの小さいSSD 2 2へのデータブロックの移動処理は、第1実施形態のごとくSCMを具備する装置にも適用可能である。これにより、第1実施形態において、当初DWP Dの大きいSSD 2 1に格納され、書込み重視であると推定されたデータブロックが、推定後、読出し重視であると判定された場合に最適なデータ配置を行なうことができる。

【 0 1 1 6 】

第2実施形態のデータキャッシング処理部1 1 1及びデータ圧縮・重複排除処理部1 1 2は、それぞれ、第1実施形態のデータキャッシング処理部1 1 1及びデータ圧縮・重複排除処理部1 1 2と同様の機能を果たす。

10

【 0 1 1 7 】

第2実施形態のデータアクセス特性判定処理部1 1 4も、第1実施形態と同様、定期的にデータブロックの移動判定を行なうために、データブロックに対するアクセス特性の判定または推定を行なう判定処理部の一例である。特に、第2実施形態の判定処理部1 1 4は、DWP Dの大きいSSD 2 1に格納されたデータブロックに対するアクセス特性の判定または推定を行なう。

【 0 1 1 8 】

そして、第2実施形態のデータ移動制御部（格納処理部）1 1 5は、判定処理部1 1 4によってアクセス特性が読出し重視であると判定されたデータブロックを、DWP Dの大きいSSD 2 1から、DWP Dの小さいSSD 2 2に格納する。

20

【 0 1 1 9 】

なお、第2実施形態において、第1実施形態と同様、重複排除処理（重複排除機能）によって管理される、データブロックの参照数（重複回数）に基づいて、アクセス特性の判定または推定を行なってもよい。判定処理部1 1 4は、例えば、データブロックの参照数が2以上のデータブロックのアクセス特性は読出し重視であると判定する。また、判定処理部1 1 4は、アクセス特性が読出し重視でないデータブロックのアクセス特性は書込み重視であると推定してもよい。

【 0 1 2 0 】

また、第2実施形態のSSDプール管理部1 1 6及びSSDドライバ1 1 7は、それぞれ、第1実施形態のSSDプール管理部1 1 6及びSSDドライバ1 1 7と同様の機能を果たす。

30

【 0 1 2 1 】

〔 3 - 3 〕動作

次に、図1 6に示すフローチャート（ステップS 7 1～S 7 4）に従って、図1 5に示すストレージ制御装置1 0 Aを含むストレージ装置1 Aでのデータ書込み処理について説明する。

【 0 1 2 2 】

第1実施形態と同様、まず、ホストインタフェース1 5によってホストサーバ3 0から書込み要求とともに書込み対象のデータが受信されると、データキャッシング処理部1 1 1によって、書込み対象のデータは、DRAM 1 2上でバッファリングされる（ステップS 7 1）。そして、バッファリング時点で、データキャッシング処理部1 1 1からホストインタフェース1 5経由で、ホストサーバ3 0へ書込み完了応答が返される。

40

【 0 1 2 3 】

ホストサーバ3 0に対し書込み完了応答を行なった後、DRAM 1 2上のデータは、所定のブロック長単位のデータブロックに分割される。そして、データ圧縮・重複排除処理部1 1 2によって、データブロック毎に、圧縮処理や重複排除処理といったデータ量削減処理が実行される（ステップS 7 2）。第2実施形態でも、データブロックの参照数に基づきデータブロックのアクセス特性の判定を行なうため、重複排除処理に着目して説明を行なう。

【 0 1 2 4 】

50

重複排除処理では、DRAM 12上の書込み対象のデータブロックが重複していると判断された場合、不揮発性デバイス(SSD 21, 22)に格納済の重複データと同じであることが分かるよう2種類のメタ情報が更新される。2種類のメタ情報は、上述したマッピング情報及びリファレンスカウンタ値である。ただし、データブロックそのものの不揮発性デバイスへのライトバック処理は行なわれず、メタ情報更新後、DRAM 12上のデータは破棄される。データブロックが重複していないと判断された場合、ライトバック処理が行なわれる。そして、DRAM 12上のデータブロックはDWP Dの大きいSSD 21へコピーされた後、当該データブロックについての2種類のメタ情報が更新されてから、DRAM 12上のデータブロックは削除される。

【0125】

第2実施形態において、論理ボリュームの実データは、DWP Dの異なる複数のSSD プール(ここでは2種類のSSD 21, 22)のどこかに格納されている。DRAM 12上の書込み対象データブロックがライトバックされた直後、当該書込み対象データは、DWP Dの大きいSSD プールに保持される。その後、当該書込み対象データは、判定処理部114によって、DWP Dの大きいSSD プールにおけるデータブロックのアクセス特性が判定される(ステップS73)。そして、アクセス特性が読出し重視であると判定されたチャンク(データブロック)は、データ移動制御部115によって、DWP Dの大きいSSD プールからDWP Dの小さいSSD プールに移動される(ステップS74)。一方、アクセス特性が読出し重視でないデータブロックは、DWP Dの大きいSSD プールで保持される。

【0126】

また、第2実施形態において、2種類のメタ情報である、マッピング情報、及びリファレンスカウンタ値は、論理ボリュームのマッピングテーブル120A(図17参照)を用いて、例えばDRAM 12上で管理される。マッピング情報は、「論理ボリューム番号と論理ブロックアドレスとの組合せ」と、「不揮発性デバイスの番号とブロック位置との組合せ」とをマッピングする情報である。また、リファレンスカウンタ値は、一つの物理データブロックが複数の論理データブロックと対応している際の「論理データブロックの参照数」に相当する。なお、第1実施形態でも前述したように、リファレンスカウンタ値は、SSD プール管理部116においても保持されるが、ここでは言及しない。

【0127】

ここで、図17は、第2実施形態における論理ボリュームと不揮発性デバイスとを対応付けるマッピングテーブル120Aの構造の一例を示す図である。マッピングテーブル120Aの各エントリには、対応するデータブロックについて、次の情報が、図17に示すようなテーブル形式で、例えばDRAM 12上において展開保持される。

【0128】

・論理アドレス(ボリューム# LBA) : ホストサーバ30からアクセスされるLUNに対応する論理ボリュームの番号と、その論理ボリューム内でのデータブロックの位置を示すLBAとの組合せ情報。

【0129】

・フラグ(Flag) : 論理アドレスで示されるデータブロックが格納されているSSD プールがDWP Dの大きいSSD プール(寿命の長いSSD 21)であるかそれ以外であるかを示す情報。第2実施形態では、DWP Dの大きいSSD プールにデータブロックが格納されている時にフラグは1に設定される。一方、それ以外のSSD プール(第2実施形態ではDWP Dの小さいSSD プール/SSD 22)にデータブロックが格納されている時にフラグは0に設定される。

【0130】

・不揮発性デバイスアドレス(SSD-Pool#-block#-Offset) : 対応するデータブロックのSSD プール内の論理位置の情報。

【0131】

・リファレンスカウンタ(重複排除参照数) : エントリに対応するデータブロックが

10

20

30

40

50

複数の論理アドレスから参照されている場合に、その参照数（重複回数）を示す値。この値は、データ圧縮・重複排除処理部 112 の重複排除機能によって管理される。当該データブロックが重複排除されていないデータブロックつまり重複していないデータブロックである場合、リファレンスカウンタの値は 1 になる。図 17 に示す例では、不揮発性デバイスアドレスが SSD-Pool1-block15-Offset5 のデータブロックの参照数は 3 である。また、不揮発性デバイスアドレスが SSD-Pool1-block10-Offset7 のデータブロックの参照数は 2 であり、不揮発性デバイスアドレスが SSD-Pool2-block10-Offset3 のデータブロックの参照数は 2 である。また、不揮発性デバイスアドレスが SSD-Pool1-block10-Offset1 のデータブロックの参照数は 1 である。

【0132】

なお、図 18 は、図 17 に示すマッピングテーブル 120A によるマッピングの具体例を示す図であり、図 17 に示すマッピングテーブル 120A によるマッピングを概念化したものである。

【0133】

次に、図 16 のステップ S73 におけるデータアクセス特性判定処理は、第 1 実施形態と同様、インターバル時間調整処理と移動データ抽出処理とを含む。第 2 実施形態のインターバル時間調整処理は、第 1 実施形態と同様、データアクセス特性判定処理部 114 が動作する時間間隔を調整する処理である。第 2 実施形態のインターバル時間調整処理においても、図 9 及び図 10 を参照しながら前述した処理と同様の処理が実行されるので、その説明は省略する。

【0134】

ただし、第 2 実施形態では、図 10 の移動データ抽出処理として、第 1 実施形態の図 11 ~ 図 13 に示す処理の代わりに、図 19 及び図 20 に示す処理が実行される。つまり、第 2 実施形態の移動データ抽出処理では、データアクセス特性判定処理部 114 によって、DWPD の大きい SSD プールから DWPD の小さい SSD プールへ移動すべきデータブロックを抽出する処理が実行される。

【0135】

次に、図 19 及び図 20 に示すフローチャートに従って、第 2 実施形態のストレージ制御装置 10A による移動データ抽出処理について説明する。

【0136】

まず、図 19 に示すフローチャート（ステップ S81, S82）に従って、第 2 実施形態のストレージ制御装置 10A による移動データ抽出処理の全体的な流れについて説明する。第 2 実施形態において、データブロックの移動は、寿命の長い SSD 21 と寿命の短い SSD 22 との間での移動（図 16 のステップ S74 参照）になるので、チャンク単位で移動対象のデータブロックを抽出する。

【0137】

第 2 実施形態の移動データ抽出処理は、図 19 に示す 2 つの処理（ステップ S81, S82）を含む。最初の処理（ステップ S81）では、マッピングテーブル 120A においてフラグが 1 のエントリ、つまり寿命の長い SSD 21 に格納されるデータブロックが、不揮発性デバイスアドレスをキーとしてソートされ、チャンク順に並べられる。

【0138】

ステップ S81 でのソート処理後、判定処理部 114 は、各エントリのリファレンスカウンタ値を調査する。そして、判定処理部 114 は、リファレンスカウンタ値が 2 以上のエントリに対応するデータブロックの占める割合が設定値以上であるチャンクを、移動対象のチャンクとして抽出する（ステップ S82）。例えば設定値が 100% であれば、ステップ S82 で抽出される移動対象のチャンクに属する全てのデータブロックは、重複している、つまり重複排除の対象になっている。

【0139】

図 20 に示すフローチャート（ステップ S91 ~ S95）に従って、図 19 のステップ S82 で実行される、リファレンスカウンタ値が 2 以上のデータブロックの占める割合が

10

20

30

40

50

設定値以上のチャンクの抽出処理について説明する。

【0140】

まず、ステップS81でのソート結果が参照され、調査対象の最初のチャンクが選択される(ステップS91)。そして、選択されたチャンク内のデータブロックのうち、リファレンスカウンタ値が2以上のデータブロックの占める割合が、設定値以上であるか否かが判定される(ステップS92)。

【0141】

リファレンスカウンタ値が2以上のデータブロックの占める割合が設定値未満である場合(ステップS92のNORルート)、全てのチャンクの調査が完了したか否かが判定される(ステップS94)。全てのチャンクの調査が完了した場合(ステップS94のYESルート)、抽出処理は終了する。一方、全てのチャンクの調査が完了していない場合(ステップS94のNORルート)、未調査の次のチャンクが選択され(ステップS95)、特性判定処理部114はステップS92の処理に戻る。

10

【0142】

リファレンスカウンタ値が2以上のデータブロックの占める割合が設定値以上である場合(ステップS92のYESルート)、特性判定処理部114は、現在選択しているチャンクを、移動チャンク(移動対象のチャンク)として登録し(ステップS93)、ステップS94の処理に移行する。

【0143】

最後に、図16のステップS74におけるSSD21, 22間のデータ移動処理では、図16のステップS73におけるデータアクセス特性判定処理(図20の処理参照)によって抽出登録された移動チャンクを、DWPDの大きいSSD21からDWPDの小さいSSD22へ移動させる処理が実行される。

20

【0144】

チャンクを移動させる処理では、次の三つの処理が、SSDプール管理部116へ依頼される。一つ目の処理は、移動元チャンクのデータブロックを移動先チャンクへコピーする処理である。二つ目の処理は、移動元チャンク領域の無効化処理である。三つ目は、先の二つの処理が完了した後、マッピングテーブル120Aの不揮発性デバイスアドレスの内容を、移動先の新しいアドレスに置き換える処理である。

【0145】

上述のように、第2実施形態のストレージ制御装置10Aでは、重複排除処理で管理されるリファレンスカウンタの値に基づき、DWPDの大きいSSDプールからDWPDの小さいSSDプールへ、チャンク単位で読み重視のデータブロックの移動が行われる。これにより、データブロックのアクセス特性(読み重視か書き込み重視か)に応じて寿命指標(DWPD)の異なるSSD21, 22へ最適なデータ格納が実現され、特に、DWPDの小さいSSD22の寿命を延ばすことが可能になる。

30

【0146】

また、第2実施形態によれば、第1実施形態で備えられるSCMを用いないため、第1実施形態と比べ、より安価な構成で第1実施形態と同様の作用効果を得ることができる。

【0147】

{4} その他

以上、本発明の好ましい実施形態について詳述したが、本発明は、係る特定の実施形態に限定されるものではなく、本発明の趣旨を逸脱しない範囲内において、種々の変形、変更して実施することができる。

40

【0148】

例えば、上述した実施形態では、寿命の異なるSSD、つまり寿命指標であるDWPDのSSDが大小(長短)の二種類であ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なく、三種類以上で有る場合にも同様に適用され、上述した実施形態と同様の作用効果を得ることができる。

【0149】

50

〔 5 〕 付記

以上の各実施形態を含む実施形態に関し、さらに以下の付記を開示する。

【 0 1 5 0 】

(付記 1)

寿命の異なる複数の記憶部と、
前記複数の記憶部に格納されるデータブロックを管理する管理部と、を備え、
前記管理部は、
前記データブロックに対するアクセス特性の判定または推定を行なう判定処理部と、
前記寿命の異なる複数の記憶部のうち、前記判定処理部によって判定または推定された
前記アクセス特性に応じた記憶部に、前記データブロックを格納する格納処理部と、を有
する、ストレージ装置。

10

【 0 1 5 1 】

(付記 2)

前記管理部は、前記データブロックの重複を排除する重複排除機能を有し、
前記判定処理部は、前記重複排除機能によって管理される、前記データブロックの参照
数に基づいて、前記アクセス特性の判定または推定を行なう、付記 1 に記載のストレージ
装置。

【 0 1 5 2 】

(付記 3)

前記判定処理部は、前記参照数が 2 以上のデータブロックの前記アクセス特性は読出し
重視であると判定し、

20

前記格納処理部は、前記参照数が 2 以上のデータブロックを、前記寿命の異なる複数の
記憶部のうち、他の記憶部よりも寿命の短い記憶部に格納する、付記 2 に記載のストレ
ージ装置。

【 0 1 5 3 】

(付記 4)

前記判定処理部は、前記アクセス特性が読出し重視でないデータブロックの前記アクセ
ス特性は書込み重視であると推定し、

前記格納処理部は、前記アクセス特性が読出し重視でないデータブロックを、前記寿命
の異なる複数の記憶部のうち、他の記憶部よりも寿命の長い記憶部に格納する、付記 3 に
記載のストレージ装置。

30

【 0 1 5 4 】

(付記 5)

ホストからの、前記寿命の異なる複数の記憶部に対する書込み対象データブロックを、
前記データブロックとして一時的に格納する一時記憶部を有し、

前記判定処理部は、前記一時記憶部に格納された前記データブロックに対するアクセス
特性の判定または推定を行ない、

前記格納処理部は、前記判定処理部によって前記アクセス特性が読出し重視であると判
定されたデータブロックを、前記一時記憶部から、前記他の記憶部よりも寿命の短い記憶
部に格納する、付記 4 に記載のストレージ装置。

40

【 0 1 5 5 】

(付記 6)

前記判定処理部は、前記参照数が 1 で書込み後所定時間経過したデータブロックの前記
アクセス特性は書込み重視であると推定し、

前記格納処理部は、前記判定処理部によって前記アクセス特性が書込み重視であると推
定されたデータブロックを、前記一時記憶部から、前記他の記憶部よりも寿命の長い記憶
部に格納する、付記 5 に記載のストレージ装置。

【 0 1 5 6 】

(付記 7)

前記寿命の異なる複数の記憶部のうち、他の記憶部よりも寿命の長い記憶部は、ホスト

50

からの、前記寿命の異なる複数の記憶部に対する書込み対象データブロックを、前記データブロックとして格納し、

前記判定処理部は、前記他の記憶部よりも寿命の長い記憶部に格納された前記データブロックに対するアクセス特性の判定または推定を行ない、

前記格納処理部は、前記判定処理部によって前記アクセス特性が読み重視であると判定されたデータブロックを、前記他の記憶部よりも寿命の長い記憶部から、前記他の記憶部よりも寿命の短い記憶部に格納する、付記 2 ~ 付記 4 のいずれか一項に記載のストレージ装置。

【 0 1 5 7 】

(付記 8)

寿命の異なる複数の記憶部に格納されるデータブロックを管理する管理部を備え、

前記管理部は、

前記データブロックに対するアクセス特性を判定する判定処理部と、

前記寿命の異なる複数の記憶部のうち、前記判定処理部によって判定された前記アクセス特性に応じた記憶部に、前記データブロックを格納する格納処理部と、を有する、ストレージ制御装置。

【 0 1 5 8 】

(付記 9)

前記管理部は、前記データブロックの重複を排除する重複排除機能を有し、

前記判定処理部は、前記重複排除機能によって管理される、前記データブロックの参照数に基づいて、前記アクセス特性の判定または推定を行なう、付記 8 に記載のストレージ制御装置。

【 0 1 5 9 】

(付記 1 0)

前記判定処理部は、前記参照数が 2 以上のデータブロックの前記アクセス特性は読み重視であると判定し、

前記格納処理部は、前記参照数が 2 以上のデータブロックを、前記寿命の異なる複数の記憶部のうち、他の記憶部よりも寿命の短い記憶部に格納する、付記 9 に記載のストレージ制御装置。

【 0 1 6 0 】

(付記 1 1)

前記判定処理部は、前記アクセス特性が読み重視でないデータブロックの前記アクセス特性は書込み重視であると推定し、

前記格納処理部は、前記アクセス特性が読み重視でないデータブロックを、前記寿命の異なる複数の記憶部のうち、他の記憶部よりも寿命の長い記憶部に格納する、付記 1 0 に記載のストレージ制御装置。

【 0 1 6 1 】

(付記 1 2)

ホストからの、前記寿命の異なる複数の記憶部に対する書込み対象データブロックを、前記データブロックとして一時的に格納する一時記憶部を有し、

前記判定処理部は、前記一時記憶部に格納された前記データブロックに対するアクセス特性の判定または推定を行ない、

前記格納処理部は、前記判定処理部によって前記アクセス特性が読み重視であると判定されたデータブロックを、前記一時記憶部から、前記他の記憶部よりも寿命の短い記憶部に格納する、付記 1 1 に記載のストレージ制御装置。

【 0 1 6 2 】

(付記 1 3)

前記判定処理部は、前記参照数が 1 で書込み後所定時間経過したデータブロックの前記アクセス特性は書込み重視であると推定し、

前記格納処理部は、前記判定処理部によって前記アクセス特性が書込み重視であると推

10

20

30

40

50

定されたデータブロックを、前記一時記憶部から、前記他の記憶部よりも寿命の長い記憶部に格納する、付記 1 2 に記載のストレージ制御装置。

【 0 1 6 3 】

(付記 1 4)

前記寿命の異なる複数の記憶部のうち、他の記憶部よりも寿命の長い記憶部は、ホストからの、前記寿命の異なる複数の記憶部に対する書込み対象データブロックを、前記データブロックとして格納し、

前記判定処理部は、前記他の記憶部よりも寿命の長い記憶部に格納された前記データブロックに対するアクセス特性の判定または推定を行ない、

前記格納処理部は、前記判定処理部によって前記アクセス特性が読出し重視であると判定されたデータブロックを、前記他の記憶部よりも寿命の長い記憶部から、前記他の記憶部よりも寿命の短い記憶部に格納する、付記 9 ~ 付記 1 1 のいずれか一項に記載のストレージ制御装置。

10

【符号の説明】

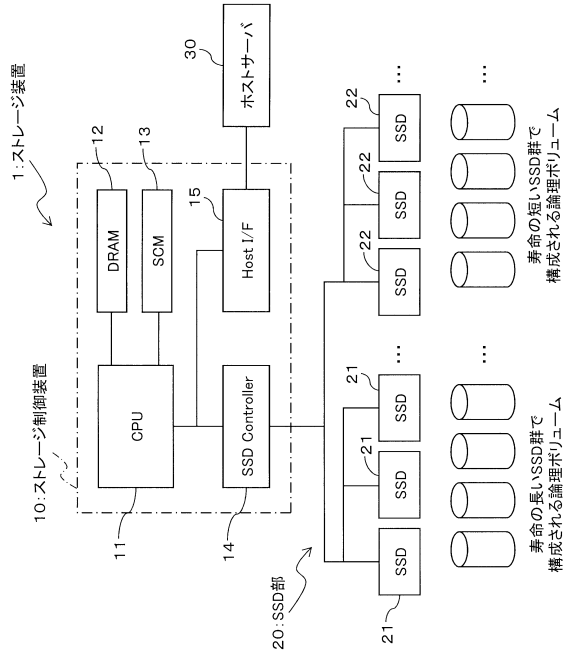
【 0 1 6 4 】

- 1 , 1 A ストレージ装置
- 1 0 , 1 0 A ストレージ制御装置 (管理部)
- 1 1 C P U (処理部)
- 1 1 1 データキャッシング処理部
- 1 1 2 データ圧縮・重複排除処理部
- 1 1 3 S C M ライトバッファ管理部
- 1 1 4 データアクセス特性判定処理部 (判定処理部)
- 1 1 5 データ移動制御部 (格納制御部)
- 1 1 6 S S D プール管理部
- 1 1 7 S S D ドライバ
- 1 2 0 , 1 2 0 A マッピングテーブル
- 1 2 1 S C M フリーリスト
- 1 2 2 S C M アロケーションリスト
- 1 2 D R A M
- 1 3 S C M (一時記憶部)
- 1 4 S S D コントローラ
- 1 5 ホストインタフェース (Host I/F)
- 2 0 S S D 部
- 2 1 寿命の長い S S D (D W P D の大きい S S D)
- 2 2 寿命の短い S S D (D W P D の小さい S S D)
- 3 0 ホストサーバ (サーバ)

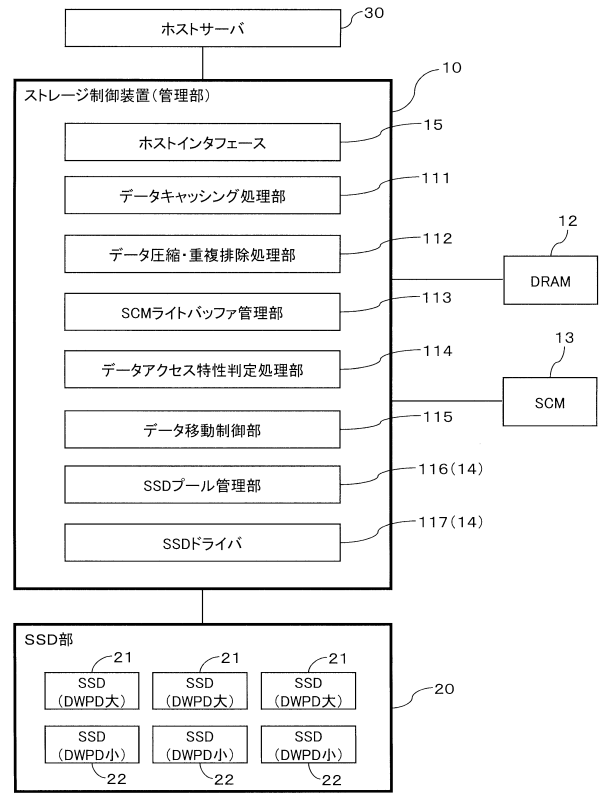
20

30

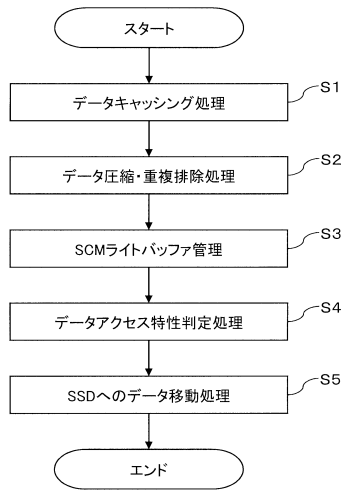
【図1】



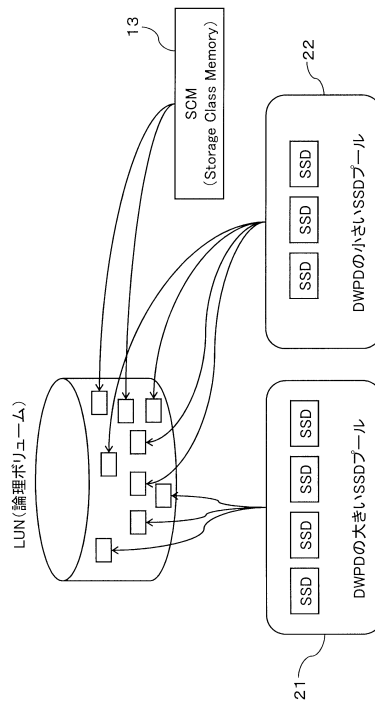
【図2】



【図3】



【図4】

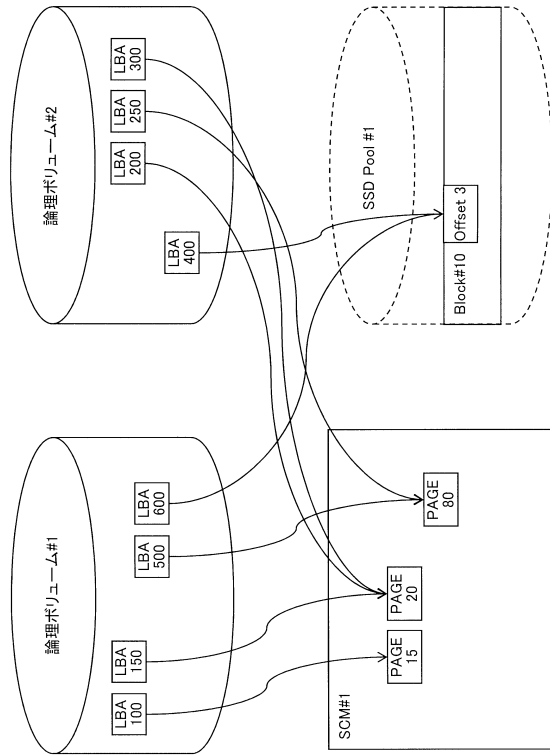


【図5】

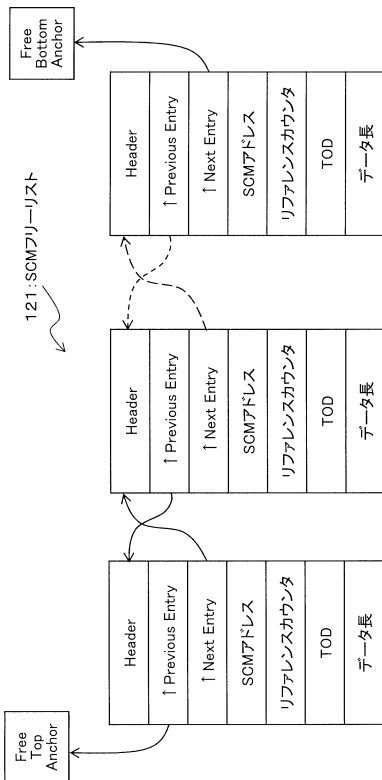
120: マッピングテーブル

論理アドレス	Flag	不揮発性デバイスアドレス
ボリューム# LBA	Flag	SCM#-Page# または SSD-Pool#-block#-Offset
ボリューム#1 LBA-100	1	SCM1-Page15
ボリューム#1 LBA-150	1	SCM1-Page20
ボリューム#1 LBA-500	1	SCM1-Page80
ボリューム#1 LBA-600	0	SSD-Pool1-block10-Offset3
ボリューム#2 LBA-200	1	SCM1-Page20
ボリューム#2 LBA-250	1	SCM1-Page80
ボリューム#2 LBA-300	1	SCM1-Page20
ボリューム#2 LBA-400	0	SSD-Pool1-block10-Offset3

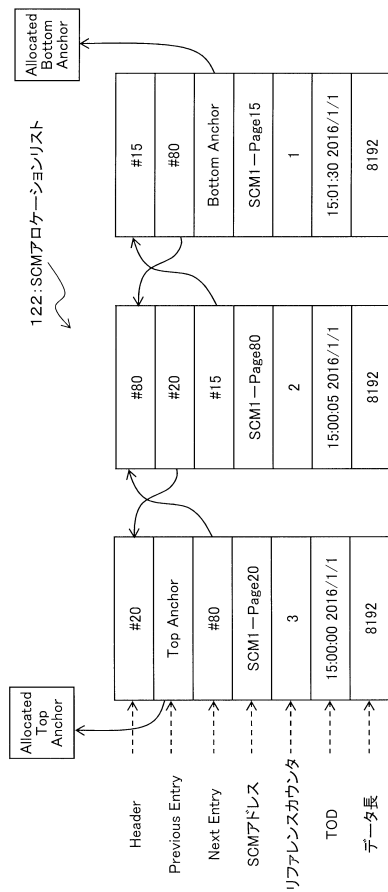
【図6】



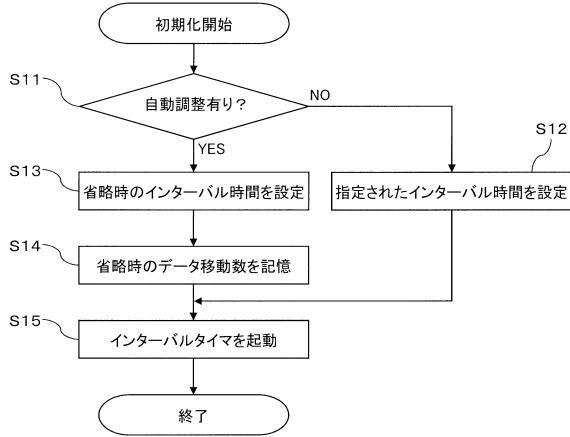
【図7】



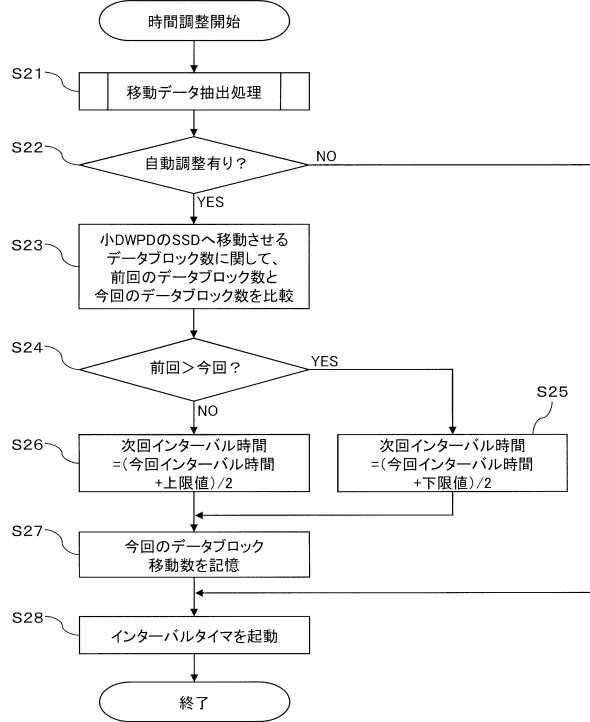
【図8】



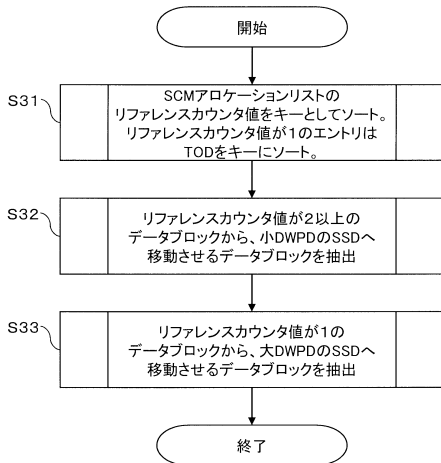
【図9】



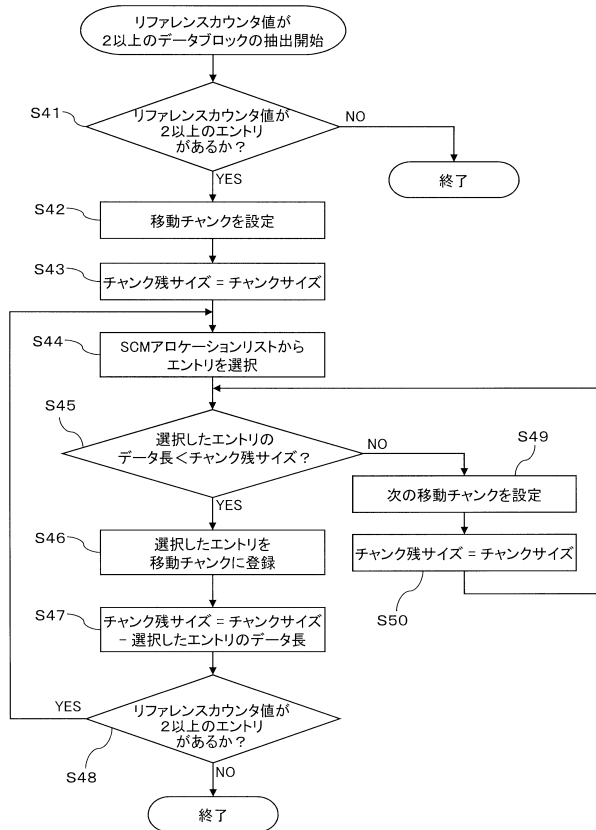
【図10】



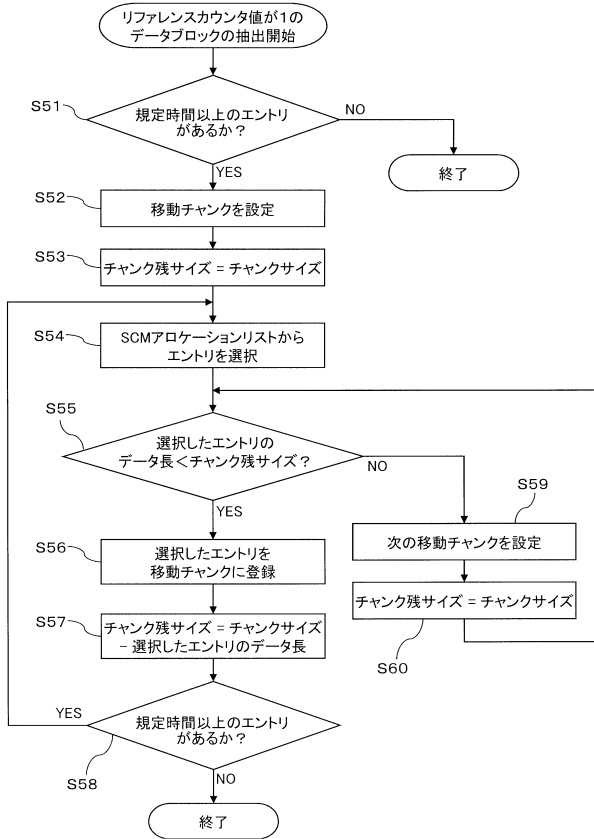
【図11】



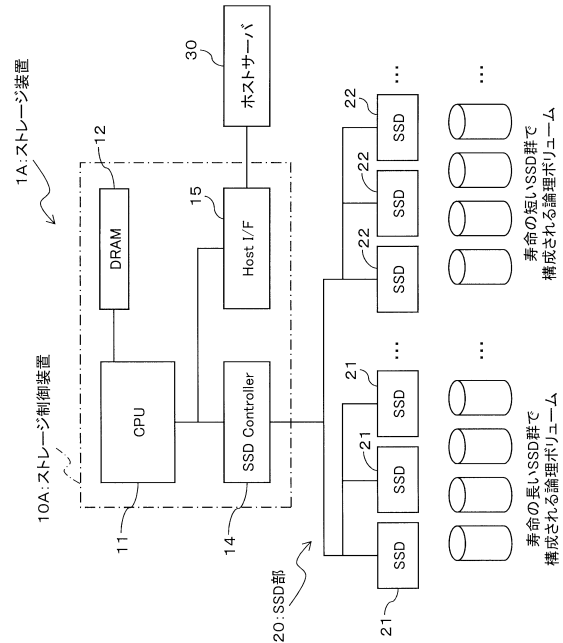
【図12】



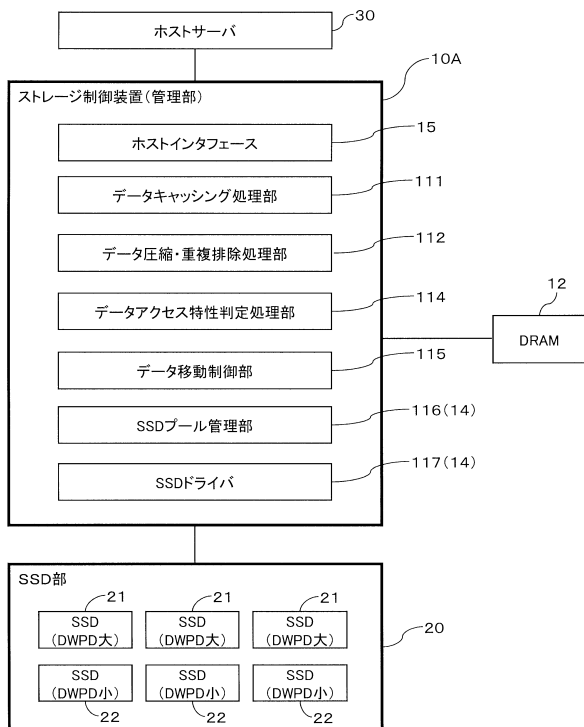
【図13】



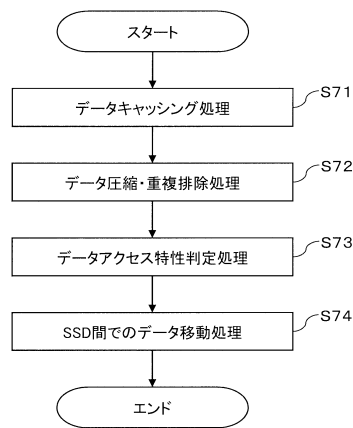
【図14】



【図15】



【図16】

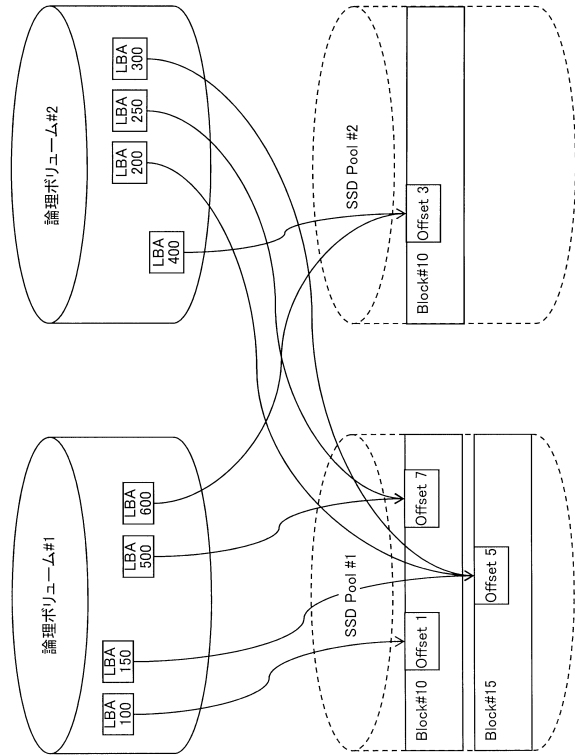


【図17】

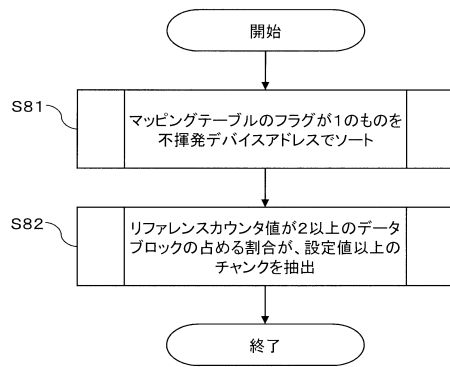
120A: マッピングテーブル

論理アドレス	Flag	不揮発性デバイスアドレス	リファレンスカウンタ
ボリューム# LBA	Flag	SSD-Pool#-block#-Offset	重複排除参照数
ボリューム#1 LBA-100	1	SSD-Pool1-block10-Offset1	1
ボリューム#1 LBA-150	1	SSD-Pool1-block15-Offset5	3
ボリューム#1 LBA-500	1	SSD-Pool1-block10-Offset7	2
ボリューム#1 LBA-600	0	SSD-Pool2-block10-Offset3	2
ボリューム#2 LBA-200	1	SSD-Pool1-block15-Offset5	3
ボリューム#2 LBA-250	1	SSD-Pool1-block10-Offset7	2
ボリューム#2 LBA-300	1	SSD-Pool1-block15-Offset5	3
ボリューム#2 LBA-400	0	SSD-Pool2-block10-Offset3	2

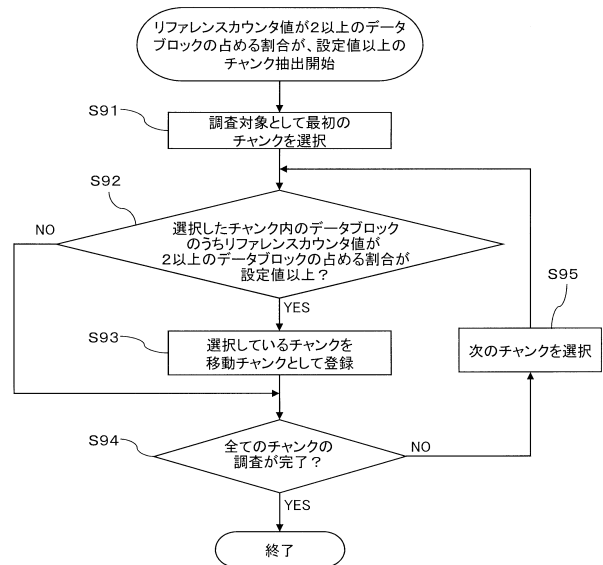
【図18】



【図19】



【図20】



フロントページの続き

- (56)参考文献 特開2015-204118(JP,A)
米国特許出願公開第2016/0179386(US,A1)
米国特許出願公開第2015/0199268(US,A1)

- (58)調査した分野(Int.Cl., DB名)
- | | |
|------|-------|
| G06F | 3/06 |
| G06F | 3/08 |
| G06F | 13/10 |