

(12) 特許協力条約に基づいて公開された国際出願

(19) 世界知的所有権機関
国際事務局



(10) 国際公開番号

WO 2015/173889 A1

(43) 国際公開日
2015年11月19日(19.11.2015)

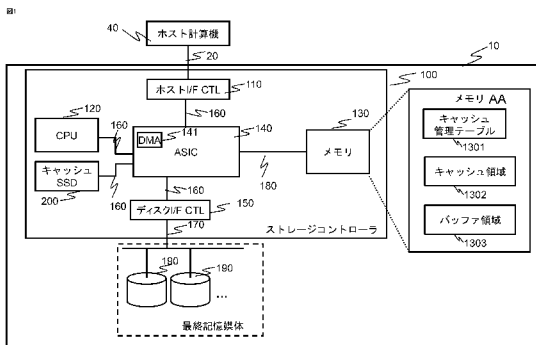
- (51) 国際特許分類:
G06F 12/08 (2006.01) G06F 12/12 (2006.01)
- (21) 国際出願番号: PCT/JP2014/062714
- (22) 国際出願日: 2014年5月13日(13.05.2014)
- (25) 国際出願の言語: 日本語
- (26) 国際公開の言語: 日本語
- (71) 出願人: 株式会社日立製作所 (HITACHI, LTD.) [JP/JP]; 〒1008280 東京都千代田区丸の内一丁目6番6号 Tokyo (JP).
- (72) 発明者: 伊藤 悠二 (ITO, Yuji); 〒1008280 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内 Tokyo (JP). 鈴木 彬史 (SUZUKI, Aki-fumi); 〒1008280 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内 Tokyo (JP). 山本 彰 (YAMAMOTO, Akira); 〒1008280 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内 Tokyo (JP).
- (74) 代理人: 特許業務法人第一国際特許事務所 (PATENT CORPORATE BODY DAI-ICHI KOKUSAI TOKKYO JIMUSHO); 〒1080014 東京都港区芝4丁目10番5号 Tokyo (JP).
- (81) 指定国 (表示のない限り、全ての種類の国内保護が可能): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JP, KE, KG, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.
- (84) 指定国 (表示のない限り、全ての種類の広域保護が可能): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, SZ, TZ, UG, ZM, ZW), ユーラシア (AM, AZ, BY, KG, KZ, RU, TJ, TM), ヨーロッパ (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG).

添付公開書類:

— 国際調査報告 (条約第 21 条(3))

(54) Title: STORAGE DEVICE

(54) 発明の名称: ストレージ装置



- 40 Host computer
- 100 Storage controller
- 110 Host interface CTL
- 120 CPU
- 130 AA Memory
- 150 Disk interface CTL
- 190 Final storage medium
- 200 Cache SSD
- 1301 Cache management table
- 1302 Cache region
- 1303 Buffer region

(57) Abstract: A storage device according to the present invention is provided with a cache device that caches data requested from an upper-level device such as a host computer. The cache device is provided with an internal controller that includes an internal cache configured from a memory capable of high speed access, such as a DRAM, and a storage unit configured from a storage medium having an erasure number limit, such as a flash memory (FM), and that has the function of controlling the internal cache and the storage unit and of performing a conversion and inverse-conversion process with respect to stored data. After the storage device stores write data in the cache device, the cache device, on the basis of an attribute or characteristics of the write data, determines the priority of an internal destage of the write data from the internal cache unit to the storage unit. The cache device then performs an internal destage from the internal cache unit to the storage unit preferentially from data with higher internal destage priority.

(57) 要約: 本発明に係るストレージ装置は、ホスト計算機などの上位装置から要求されたデータをキャッシュするキャッシュ装置を備える。キャッシュ装置は、DRAM等の高速アクセス可能なメモリで構成される内部キャッシュと、FM等の消去回数制限のある記憶媒体で構成される記憶部を有し、それらの制御および格納データに対する変換及び逆変換処理を行う機能を持つ内部コントローラを備える。ストレージ装置がキャッシュ装置にライトデータを格納した後、キャッシュ装置はライトデータの属性・特性に基づいて、ライトデータの内部キャッシュ部から記憶部への内部ディスティージの優先度を決定する。そしてキャッシュ装置は、内部ディスティージの優先度が高いデータを優先的に、内部キャッシュ部から記憶部に内部ディスティージする。

の優先度を決定する。そしてキャッシュ装置は、内部ディスティージの優先度が高いデータを優先的に、内部キャッシュ部から記憶部に内部ディスティージする。

WO 2015/173889 A1

明 細 書

発明の名称：ストレージ装置

技術分野

[0001] 本発明は、ストレージ装置に関する。

背景技術

[0002] 近年、サーバやストレージシステム（ストレージ装置）には、記憶媒体として半導体デバイスを用いたSSD（Solid State Drive）が盛んに利用されている。SSDは、従来のHDD（Hard Disk Drive）に比べてランダムアクセス性能に優れ、特にNANDフラッシュメモリを用いたSSDについては大容量化、低価格化が進み、HDDとの置き換えが進んでいる。また、不揮発記憶媒体としてSCM（Storage Class Memory）は、NANDフラッシュメモリに比べ小容量で高価格であるが、アクセス速度や書き換え回数に優れ、高性能・長寿命の記憶媒体として注目されている。SCMは、不揮発性半導体メモリの一種であり、例えば、MRAM（Magnetic Random Access Memory）、PRAM（Phase Change Random Access Memory）、ReRAM（Resistance Random Access Memory）である。

[0003] サーバやストレージシステムでは、データへのアクセス頻度の偏りを利用し、サーバやストレージコントローラに高速な記憶媒体を用いたキャッシュを設けることで、システム性能を向上させている。データの最終格納先であるHDDへのアクセスを逐一行うと、大きなオーバーヘッドとなってしまうためである。キャッシュに高頻度でアクセスされるデータを格納しておくことで、再度同データにアクセスする際にはHDDへアクセスを行わずに高速にアクセスすることができる。そのため、キャッシュ容量が大きいほど、高速にアクセスできるデータ量が増加する。つまりキャッシュでのデータのヒット率が向上して、システム性能の向上が期待できる。

[0004] サーバやストレージシステムにおけるキャッシュとしてDRAMが広く用いられている。DRAMはHDDより遥かに高速であるが、大容量化や価格に問題がある。近年では、DRAMに比べ安価かつ大容量な特性を持つSSDをストレージシステムのキャッシュとして適用する技術が公開されている（特許文献1）。DRAMよりも低価格で大容量なNANDフラッシュメモリをキャッシュとして用いることで、キャッシュヒット率を向上させ、システム性能、コストパフォーマンスの改善を可能としている。

先行技術文献

特許文献

[0005] 特許文献1：米国特許第8214595号明細書

発明の概要

発明が解決しようとする課題

[0006] フラッシュメモリでは、データ消去回数が限られている。さらにフラッシュメモリでは、同じ物理アドレスに対するデータの上書きはできず、一度リード／ライト単位（ページ）よりも大きい単位（ブロック）で消去してから書き込む必要がある。上書きのために逐一消去することはオーバーヘッドが大きく、限られた消去回数を浪費してしまう。そのため、フラッシュメモリを用いたシステムでは、ユーザが認識できる容量よりも大きな容量を持ち、上書きされるデータ用に更新領域を確保して、ユーザが認識できる論理的な領域における論理アドレスと、実際の物理アドレスとの対応付けを行って更新処理を行う。例えば、論理アドレスLに対して物理アドレスPを割り当てているデータが更新された場合、論理アドレスLに対して新たに物理アドレスP'を割り当て、その物理ページにデータを書き込み、物理アドレスPのデータを無効化する。この後、フラッシュメモリの容量が足りなくなった場合には、新たに書き込める領域を確保するために無効化したデータの消去を行う必要がある。つまり、書き込みが行われるということは、SSDの寿命を消費することにつながる。

[0007] サーバやストレージシステムのキャッシュにおいて、キャッシュへのデータの読み上げをステージング、キャッシュからデータの追い出しをディスステージングと呼ぶ。サーバやストレージシステムのキャッシュとして用いられる記憶媒体では、キャッシュデータのステージング、ディスステージングにより通常の最終記憶装置（データが最終的に管理される記憶装置）よりもデータ消去回数が非常に多くなる。なぜなら、システム全体の容量に対して小容量であるキャッシュでは、多数の最終記憶装置が処理するI/Oを少数のキャッシュ用の記憶媒体で処理することになるため、記憶媒体1台あたりの負荷が高くなって更新頻度が増加してしまうからである。また、リード処理時にデータのステージングを行う際には、キャッシュ用の記憶媒体に対してはリードデータを格納するためのライト処理が行われる。つまりリードのみのアクセスしかない場合でも、キャッシュ用記憶媒体の内容の書き換えが行われ、消去回数消費につながる。特にMLC (Multi Level Cell) タイプのNANDフラッシュメモリは、データ消去回数が限られているため、キャッシュとして用いると、最終記憶装置として用いるよりもさらに短寿命となってしまう。一方で、MLCタイプのNANDフラッシュメモリより長寿命なSLC (Single Level Cell) タイプのフラッシュメモリやSCMなどは、前述した通り、MLCタイプのフラッシュメモリよりも容量単価が高いため、キャッシュ容量を増やすことが難しく、システム性能の向上に効果的でない。

[0008] 寿命に課題のあるSSDをサーバやストレージシステムといった上位装置のキャッシュとして用いる場合、キャッシュする効果の小さいデータについては、高コスト、小容量だが長寿命なDRAMのみにキャッシュ、またはSSDにキャッシュしないことが考えられる。しかし、データをSSDにキャッシュしない場合、キャッシュ装置として用いるSSDの機能が利用不可となってしまう、上位装置のCPU、または別途搭載したハードウェアによって同等の機能を実現する必要がある。上位装置のCPUを用いる場合、処理オーバヘッドが大きくなり、システム性能が低下してしまう。別途ハードウ

エアを搭載する場合はシステムのコストが上昇してしまう。

[0009] 例えば、SSDが格納したデータを内部で暗号化し、暗号化後のデータをリードできる機能を持つとする。そのSSDをキャッシュ装置として用い、データをキャッシュする際にSSD内部で暗号化し、SSDからディステージされるデータを暗号化したまま最終記憶媒体に格納することで、上位装置のCPUなどのリソースを消費することなく全データの暗号化するシステムを実現できる。データのリード時は、SSDに一度格納することで、SSD内部で復号させ、さらにはキャッシュとしてヒットした場合にHDDよりも高速にデータを読み出せる。しかし、SSDの寿命を考慮して、前述したキャッシュする効果の小さいデータについてSSDへキャッシュしない場合、暗号化および復号化ができず、同等の処理を行う上位装置のCPUまたは別ハードウェアなどのリソースが必要となる。

課題を解決するための手段

[0010] 本発明の一実施形態に係るストレージ装置は、サーバやストレージコントローラといった上位装置から要求されたデータをキャッシュするキャッシュ装置を備え、キャッシュ装置は、高速・小容量・長寿命な内部キャッシュ部と、内部キャッシュ部より寿命の短い記憶部を有し、それらの制御を行う内部コントローラを備える。内部コントローラは、キャッシュ装置の内部キャッシュ部に格納された各データの属性・特性に応じて、内部キャッシュ部から記憶部への内部ディステージの優先度を制御する。内部コントローラは、内部ディステージ優先度の高いデータを優先的に、内部キャッシュ部から記憶部へとディステージし、内部ディステージ優先度が低いデータの記憶部へのライトを抑制する。

[0011] また本発明の一実施形態に係るストレージ装置が備えるキャッシュ装置は、キャッシュ装置に格納されたデータに対する、暗号化や圧縮などのデータ変換処理を行う機能を備え、上位装置から受信したライトデータにデータ変換処理を施してキャッシュ装置に格納し、上位装置に対してデータを送信する際には、データ逆変換を施して上位装置にデータを送信する。

発明の効果

[0012] 本発明によれば、フラッシュメモリなどの、データ消去回数の限られている記憶媒体を記憶部として用いたキャッシュ装置において、記憶部へのデータ書き込み回数を適切に抑制することができ、記憶部を長寿命化することができる。また、キャッシュ装置がデータ変換等の機能を備えている場合、短寿命な記憶部へのデータ書き込み回数を抑制しつつ、ライトデータに対するキャッシュ装置の機能の適用を可能にすることができる。

図面の簡単な説明

- [0013] [図1]図1は、実施例1に係るストレージシステムの一例の構成図である。
- [図2]図2は、実施例1に係るキャッシュ管理テーブルの一例を示す図である。
- [図3]図3は、実施例1に係るキャッシュSSDの一例の構成図である。
- [図4]図4は、実施例1に係る内部キャッシュ管理テーブルの一例を示す図である。
- [図5]図5は、実施例1に係る論理物理アドレス変換テーブルの一例を示す図である。
- [図6]図6は、実施例1に係るストレージコントローラとキャッシュSSDコントローラとの間のライト要求コマンドとその応答コマンドの一例を示す図である。
- [図7]図7は、実施例1に係るストレージコントローラとキャッシュSSDコントローラとの間のリードおよびデータ無効化要求コマンドとその応答コマンドの一例を示す図である。
- [図8]図8は、実施例1に係るストレージコントローラによるライト処理のフローチャートである。
- [図9]図9は、実施例1に係るストレージコントローラによるキャッシュSSDへのキャッシュ格納処理のフローチャートである。
- [図10]図10は、実施例1に係るストレージコントローラによるキャッシュSSDから最終記憶媒体へのディスページ処理を説明する図である。

[図11]図11は、実施例1に係るキャッシュSSDでのリード要求を受領した時の処理のフローチャートである。

[図12]図12は、実施例1に係るキャッシュSSDでのデータ無効化要求を受領した時の処理のフローチャートである。

[図13]図13は、実施例1に係るキャッシュSSDでのライト要求を受領した時の処理のフローチャートである。

[図14]図14は、実施例1に係るキャッシュSSDでの内部ディスティージ処理のフローチャートである。

[図15]図15は、実施例1に係るストレージコントローラによるリード処理のフローチャートである。

[図16]図16は、実施例2に係るキャッシュSSDでの内部ディスティージ処理のフローチャートである。

[図17]図17は、実施例2に係るストレージコントローラによるキャッシュSSDへのライト処理のフローチャートである。

[図18]図18は、実施例3に係るストレージコントローラによるライト処理のフローチャートである。

[図19]図19は、実施例3に係るストレージコントローラによるリード処理のフローチャートである。

[図20]図20は、実施例4に係るストレージコントローラによる、ライト処理のフローチャートである。

[図21]図21は、実施例5に係るキャッシュSSDで実行されるライト処理のフローチャートである。

発明を実施するための形態

[0014] 以下、図面を用いて、本発明の実施形態について述べる。

実施例 1

[0015] 図1は、本発明の実施例1に係る計算機システムの一例の構成図である。計算機システムは、キャッシュ制御装置の一例としてのストレージ装置10と、ホスト計算機40とを有する。ストレージ装置10とホスト計算機40

とは、SAN (Storage Area Network) やLAN (Local Area Network) などのネットワーク20を介して接続される。ホスト計算機40は、例えば、汎用的な計算機であり、ストレージ装置10に格納されたデータを読み出したり、ストレージ装置10にデータを格納させたりして、所定の業務処理を実行する。

[0016] ストレージ装置10は、ストレージコントローラ100と、1以上の最終記憶媒体190とを含む。最終記憶媒体190は、ホスト計算機40からストレージ装置10に対して書き込まれたデータを最終的に格納しておく記憶媒体であり、一例としてHDDが用いられる。ただしHDD以外の記憶媒体、たとえばSSDを用いる構成であってもよい。ストレージコントローラ100は、バス170を介して、最終記憶媒体190と接続される。

[0017] ストレージコントローラ100は、ホストインタフェースコントローラ (I/F CTL) 110、CPU120、例えばDRAMに代表されるメモリ130、ディスクインタフェースコントローラ (I/F CTL) 150、2次キャッシュとして用いられるキャッシュSSD200、およびASIC140を備える。ホストI/F CTL110、CPU120、メモリ130、およびディスクI/F CTL150、キャッシュSSD200は、PCIなどの専用内部バス160、メモリバス180を介してASIC140に接続される。ホストI/F CTL110は、ネットワーク20を介しての通信を行うためのデバイスである。CPU120は、ストレージ装置10全体の動作制御を行う。メモリ130は、ストレージコントローラ100のキャッシュ全体を管理するキャッシュ管理テーブル1301とデータをキャッシュするキャッシュ領域1302、データ転送のための一時格納領域として用いられるバッファ領域1303を含む。以降、キャッシュ領域1302のことをキャッシュDRAM1302と呼ぶこともある。ASIC140は、CPU120を介さずにメモリ130上のデータを転送するDMA (Direct Memory Access) 141を備える。ディスクI/F CTL150は、内部バス170を介して、最終記憶媒体190と通

信を行うためのデバイスである。

- [0018] ストレージコントローラ100は、ホスト計算機40がリードまたはライトするデータを、メモリ130上のキャッシュ領域1302またはキャッシュSSD200上にキャッシュする。キャッシュ領域1302またはキャッシュSSD200へのデータの格納及び読み出し処理については後述する。
- [0019] 図2は、実施例1に係るストレージコントローラ100のキャッシュを管理するキャッシュ管理テーブルの一例である。
- [0020] キャッシュ管理テーブル1301は、ホスト計算機40から書き込まれたデータが格納されている、キャッシュ領域1301またはキャッシュSSD200上の位置情報等を管理する情報で、ストレージコントローラ100のメモリ130上に格納されている。そしてキャッシュ管理テーブル1301の内容は、キャッシュへのデータの格納時等に、ストレージコントローラのCPU120によって更新される。
- [0021] ストレージコントローラ100（のCPU120）は、キャッシュDRAM1302及びキャッシュSSD200上の記憶領域を、64KBの固定領域に分割して管理しており、本発明の実施例ではこの64KBの固定領域のことをキャッシュセグメントと呼ぶ（または「セグメント」とも呼ばれる）。CPU120は各セグメントにストレージ装置10内で一意な識別番号（ID）を付して管理しており、この識別番号はキャッシュセグメントIDと呼ばれる（または「セグメントID」とも呼ばれる）。
- [0022] 図2の例では、キャッシュセグメントID13011の0x0~0x1FFFFFFFをキャッシュDRAM1302の各セグメント（計128GBの記憶領域）に付し、キャッシュセグメントID13011の0x200000~0x9FFFFFFFをキャッシュSSD200の各セグメント（計512GBの記憶領域）に付して管理している。ただしキャッシュ管理単位、キャッシュセグメントIDの割り当て方について、これに限定されない。
- [0023] また、ストレージコントローラ100は、キャッシュDRAM1302及びキャッシュSSD200上の記憶領域を、64KBのセグメント単位に管

理しているが、セグメントとは、ストレージコントローラ100がキャッシュDRAM1302及び／またはキャッシュSSD200内の記憶領域を確保、解放する時の最小単位であり、ストレージコントローラ100がキャッシュDRAM1302及びキャッシュSSD200に対してデータのアクセス（読み書き）を行う場合の、最小アクセス単位ではない。ストレージコントローラ100はホスト計算機40のストレージ装置10に対する最小アクセス単位であるディスクブロック（512バイト。セクタとも呼ばれる）単位で、キャッシュセグメントへのデータの読み書きを行うことが可能である。ただし、本発明の実施例に係るストレージ装置10において、後述するキャッシュSSD200が有するデータ処理機能により、データ圧縮や暗号化を行う場合、データ圧縮や暗号化を行うデータ単位である4KBが、最小アクセス単位になる。

[0024] キャッシュ管理テーブル1301には、キャッシュセグメントID13011で特定される各キャッシュセグメントについて、当該キャッシュセグメントに格納されているデータが本来格納されるべき最終記憶媒体190（あるいは複数の最終記憶媒体190から構成される論理ボリューム）上の格納位置（アドレス）であるLBA13012、当該キャッシュセグメント内に、実際にデータが格納されている位置を表すビットマップ13013、キャッシュ状態13014、最終記憶媒体190へのディスティージ時に参照する直近アクセス時刻13015（または最終アクセス時刻13015とも呼ばれる）が格納されている。

[0025] ビットマップ13013は、キャッシュセグメント上の、実際にデータが格納されている位置を表す128ビットの情報である。本発明の実施例では、キャッシュセグメントサイズは64KBだが、64KBの領域すべてにデータが格納されている（キャッシュされている）場合もあれば、必ずしもそうでない場合もある。ビットマップ13013の各ビットは、キャッシュセグメント内の各ディスクブロック（1ディスクブロックは512バイトとする。そのため、1キャッシュセグメントは128個のディスクブロックから

なる) にデータが格納されているか否かを表す。ビットマップ13013の先頭ビット(MSB)は、キャッシュセグメント内の先頭ディスクブロックにデータが格納されているか表しており、以下順に2、3、…番目のビットはそれぞれ、キャッシュセグメント内の2番目、3番目、…のディスクブロックにデータが格納されているか表している。各ビットの内容が1の場合、当該ビットに対応するキャッシュセグメント上ディスクブロックにデータが格納されていることを意味し、0の場合には、当該ビットに対応するキャッシュセグメント上ディスクブロックにデータが格納されていないことを意味する。

[0026] キャッシュ状態13014は、キャッシュされているデータが最終記憶媒体190に格納されているものに対して更新されているか否かを表すものである。キャッシュされているデータが最終記憶媒体190に格納されているものに対して更新されている場合にはキャッシュ状態13014には「dirty」が格納され、更新されていない場合(つまり、キャッシュセグメントに格納されているデータと最終記憶媒体190に格納されているデータと同じ内容である場合)には、キャッシュ状態13014には「clean」が格納される。cleanなデータである場合にはディスティージ時に最終記憶媒体190へ追い出す必要はない。一方、dirtyなデータは、最終記憶媒体190へ書き出す(ディスティージする)必要がある。なお、1キャッシュセグメント内にcleanなデータとdirtyなデータが混在している場合があるが、1キャッシュセグメント内にcleanなデータとdirtyなデータが混在している場合には、キャッシュ状態13014にはdirtyが格納される。

[0027] 直近アクセス時刻13015には、直近アクセス時刻13015に対応するキャッシュセグメントに最後にアクセスがあった時刻の情報が格納される。時刻情報としてたとえば、ある時点(2010年1月1日など)からの経過秒数が格納されるが、これ以外の形式で時刻情報を格納してもよい。直近アクセス時刻13015は、ストレージコントローラ100が最終記憶媒体

190へディスティージするデータを選択する際に、LRU (Least Recently Used) アルゴリズムを用いて選択する場合に必要な情報である。ただし、ディスティージ時のデータ選択ポリシーはLRUに限定されない。最終記憶媒体190へのディスティージするデータの選択を別ポリシーで行う場合は、キャッシュ管理テーブル1301には直近アクセス時刻13015に代えて、データ選択ポリシーに必要な情報を格納するようにする。

[0028] なお、本発明の実施例1に係るストレージ装置10では、キャッシュDRAM1302を1次キャッシュとし、キャッシュSSD200を2次キャッシュとして用いる。そのため、ホスト計算機40からのライトデータが、キャッシュDRAM1302、キャッシュSSD200の両方に格納されている場合もある。その場合、キャッシュセグメントID13011が0x0~0x1FFFFFF (キャッシュDRAM1302内セグメント) のエントリと、キャッシュセグメントID13011が0x200000~0x9FFFFFF (キャッシュSSD200内セグメント) のエントリの中に、LBA13012が同一であるエントリが存在する。

[0029] 図3は、本実施例におけるキャッシュSSD200の構成図の一例である。本実施例におけるキャッシュSSD200は、不揮発記憶媒体であるフラッシュメモリ (以下ではFMと略記することもある) を主たる記憶媒体とする。キャッシュSSD200は、内部に1または複数のSSDコントローラ (SSD CTL) 201と複数のFM2011を備える。

[0030] SSDコントローラ201は、その内部にインタフェースコントローラ (I/F CTL) 2001、プロセッサ2003、内部メモリ (DRAM) 2005、データ処理アシストユニット (Assist) 2006、FMコントローラ (FMC) 2007、およびデータ転送を相互に行う内部スイッチ (S/W) 2004を備えている。インタフェースコントローラ2001、プロセッサ2003、データ処理アシストユニット2006、FMコントローラ2007は内部バスを介して内部スイッチ2004に接続される。また内部メモリ2005は、メモリバス2009を介して内部スイッチ200

4に接続されている。

[0031] インタフェースコントローラ2001は、ストレージコントローラ100が備えるCPU120から、キャッシュSSD200に対する各種コマンドを受け付ける、あるいはデータ転送を行うためのものである。インタフェースコントローラ2001は、ストレージ装置10内のストレージコントローラ100が備えるASIC140と接続される。

[0032] CPU2003は、内部スイッチ2004を介してSSDコントローラ2001の各部位と接続され、内部メモリ2005に記録されたプログラム及び管理情報（後述する内部キャッシュ管理テーブル20052など）を用いて、SSDコントローラ2001全体を制御する。

[0033] 内部メモリ2005には具体的には、ビットコストはFM2011より高いものの、高速でかつFM2011より長寿命である、DRAMなどの揮発性メモリが用いられる。内部メモリ2005上には、SSDコントローラ2001でのデータ転送処理途中のデータを一時的に格納するために用いられる内部キャッシュ領域20051、内部キャッシュ領域20051に格納されたデータの情報を管理する内部キャッシュ管理テーブル20052、FM2011の格納情報を管理する論理物理変換テーブル20053、FM2011の物理ブロックに有効な情報が格納されているか否かなどの情報を格納するブロック情報20054が存在する。なお、これらの領域やテーブルについては、1つの内部メモリ2005上に格納してもよいし、複数の内部メモリ2005に分散して格納してもよい。以降、内部キャッシュ領域20051については、単に内部キャッシュ20051と呼ぶこともある。また、本発明の実施例に係るキャッシュSSDでは、内部キャッシュ20051のサイズ（容量）は、複数のFM2011の合計容量よりも小さいものとするが、内部キャッシュ20051のサイズがFM2011の合計容量よりも小さくない場合であっても、本発明は有効である。

[0034] FMコントローラ2007は、複数バス（例えば16）2010によってFM2011と接続する。各バス2010には複数（例えば2）のFM20

11が接続される。

[0035] また詳細は後述するが、本発明の実施例に係るキャッシュSSD200は、複数のFM2011の有する記憶領域から構成される論理的な記憶空間を、キャッシュSSD200が接続されるストレージコントローラ100に対して提供する。ストレージコントローラ100がキャッシュSSD200にアクセス（リード、ライト等）する際、たとえばキャッシュSSD200にデータの書き込み（ライト）を行う場合、論理的な記憶空間上の一次元の論理アドレスを指定したライトコマンドをキャッシュSSD200に対して発行することで、ライト対象データをキャッシュSSD200が提供する論理的記憶空間上に（仮想的に）書き込むことができる。なお、ライトコマンドとともにキャッシュSSD200に到来したデータは、一旦内部キャッシュ領域20051に格納され、その後、論理的な記憶空間に対応付けられたFM2011上の記憶領域へと移動（内部ディステージ）される。これらの制御は、CPU2003が内部メモリ2005に記録されたプログラム及び管理情報を用いて実行することで実現される。

[0036] キャッシュSSD200は、キャッシュSSD200内に格納するデータの圧縮／伸長やParity生成、暗号化／復号化などの、データ変換処理及び逆変換処理を行う機能を備える。データ処理アシストユニット2006は、これらのデータ変換処理／逆変換処理、たとえばデータ圧縮／伸長処理を行う機能を有するハードウェアである。データ処理アシストユニット2006は、CPU2003からの指示に従って、インタフェースコントローラ2001から到来してFM2011に書き込まれるデータに対するデータ変換処理（圧縮、暗号化など）を行い、またはFM2011からインタフェースコントローラ2001へ送出されるデータに対する、データ逆変換処理（圧縮データの伸長、暗号化データの復号化など）を行う。なお、これらのデータ変換／逆変換処理を行う機能をキャッシュSSD200に実装する方法として、データ処理アシストユニット（Assist）2006のようなハードウェアを設ける方法のほか、CPU2003にデータの圧縮／伸長やP

arity生成、暗号化／復号化などのデータ処理を行うプログラムを実行させることで、データ変換／逆変換を行う方法をとってもよい。

[0037] 以上説明した、SSDコントローラ201の各部位は、ASICやFPGA (Field Programmable Gate Array) として、一つの半導体素子内で構成してもよいし、複数の個別専用IC (Integrated Circuit) を相互に接続した構成であってもよい。また、キャッシュSSD200に内部メモリ2005用のキャパシタやバッテリーを備え、電断時などに揮発する内部メモリ2005に格納されたデータや管理情報をFM2011に退避する構成としてもよい。

[0038] なお、本実施例に係るキャッシュSSD200では、図3に示すように、主たる記憶媒体としてFMを使用しているが、キャッシュSSD200に搭載する記憶媒体はFMに限定されるものではなく、Phase Change RAMやResistance RAMなどの不揮発メモリであってもよい。また、FM2011の一部または全部を揮発性のRAM (DRAM等) とする構成であってもよい。

[0039] 続いて図4を用いて、内部キャッシュ管理テーブル20052について説明する。

[0040] 先に述べたとおり、ストレージコントローラ100はキャッシュSSD200が提供する論理的な記憶空間に対してデータの書き込み (ライト) を行う場合、ライト先として論理的記憶空間上のアドレス (論理アドレス) を指定したライトコマンドをキャッシュSSD200に対して発行する。ライトコマンドとともにキャッシュSSD200に到来したデータは、一旦内部キャッシュ領域20051に格納 (キャッシュ) される。キャッシュSSD200のCPU2003は、ライトデータのライト先である論理アドレスと、当該ライトデータが一旦格納される内部キャッシュ20051上の位置とを対応付けて、内部キャッシュ管理テーブル20052に記録して管理する。

[0041] ストレージコントローラ100のキャッシュ領域1302と同様、キャッシュSSD200のCPU2003は、内部キャッシュ領域20051を、

所定サイズ（たとえば4KBあるいは8KB）の固定領域（内部キャッシュセグメントと呼ぶ）に分割して管理している。これは、内部キャッシュ領域20051の割り当て、解放の単位であるとともに、データ処理アシストユニット2006による、データ処理（圧縮や暗号化）を行う際の単位である。またCPU2003は各内部キャッシュセグメントに、キャッシュSSD200内で一意な識別番号（内部キャッシュセグメントIDと呼ぶ。ただし「セグメントID」と略記されることもある）を付して管理する。

[0042] 内部キャッシュ管理テーブル20052には、内部キャッシュセグメントのセグメントID200521で特定される各内部キャッシュセグメントについて、当該内部キャッシュセグメントに格納されているデータの本来の書き込み先（先に述べた論理的な記憶空間上の位置（論理アドレス）であって、ストレージコントローラ100から指定された書き込み先アドレス）である論理アドレス200522、当該内部キャッシュセグメントのデータ処理（圧縮等）後のデータサイズを表すデータサイズ200523、内部キャッシュ状態200524、FMへのディステージ時に参照する直近アクセス時刻200525、FMへの内部ディステージ優先度200526（ディステージ優先度と略記されることもある）が格納されている。

[0043] データサイズ200523は、キャッシュSSD200にデータが格納される時に、データ処理アシストユニット2006によるデータ処理を行った結果、データサイズが変動する場合に用いられる。たとえばデータ処理として圧縮を行った場合、多くの場合データサイズが小さくなる。データサイズ200523には、データ処理（たとえば圧縮）後のデータサイズの情報が格納される。たとえば図4の例で、内部キャッシュセグメントID200521が0×1の行（20052-0）に格納されている情報を例にとると、内部キャッシュセグメントID200521が0×1である4KBの内部キャッシュ領域20051が確保されているが、当該内部キャッシュ領域にはデータ処理後のデータが格納されており、そのサイズは、データサイズ200523に格納されている内容を参照すると、1KBであることがわかる。

- [0044] なお、キャッシュSSD 200がデータ処理アシストユニット2006によるデータ処理として、たとえばデータ圧縮を行って圧縮後のデータを内部キャッシュ領域20051あるいはFM2011に格納する場合であっても、キャッシュSSD 200がストレージコントローラ100に提供している論理的な記憶空間に影響はなく、ストレージコントローラ100からは、あたかも圧縮前のデータ（非圧縮データ）を読み書きしているようにみえている。
- [0045] 内部キャッシュ状態200524は、先に説明したキャッシュ管理テーブル1301のキャッシュ状態13014と同様のもので、内部キャッシュ状態200524には「clean」または「dirty」のいずれかの情報が格納される。
- [0046] 直近アクセス時刻200525も、先に説明したキャッシュ管理テーブル1301の直近アクセス時刻13015と同様のもので、直近アクセス時刻200525に対応する内部キャッシュセグメントに最後にアクセスがあった時刻の情報が格納される。キャッシュ領域1302と同様、内部キャッシュ領域20051からFM2011へとデータを移動（ディステージ）する際に、移動対象データの選択ポリシーにはLRUアルゴリズムが用いられるが、ポリシーは必ずしもLRUアルゴリズムに限定されない。LRUアルゴリズムを用いない場合、直近アクセス時刻200525に代えて、そのディステージポリシーに応じた情報を内部キャッシュ管理テーブル20052に格納する。
- [0047] 続いて内部ディステージ優先度200526について説明する。本発明の実施例に係るキャッシュSSD 200では、ストレージコントローラ100から送信される、後述するライト要求コマンド3000に含まれるデータ情報3005を用いて、SSDコントローラ201が、そのライトデータがすぐにキャッシュSSD 200から追い出されるか（あるいはストレージコントローラ100が当該ライトデータに対して再アクセスする可能性が低い）、それとも保持される（ストレージコントローラ100から再アクセスされ

る)かを予測し、内部ディステージ優先度200526を決定し、内部キャッシュ管理テーブル20052の更新時に設定する。短寿命なFM2011のライトによる劣化を抑止するため、キャッシュSSD200からすぐ追いつ出されると予測されるデータについては、内部ディステージ優先度200526を低と設定して、できるだけFMへライトせずに内部キャッシュ20051で留める。予測通りであれば、ストレージコントローラ100による最終記憶媒体190へのディステージの際にはデータがキャッシュSSD200の内部キャッシュ20051上にある。キャッシュSSD200からのリードは内部キャッシュ20051でヒットとなり、続く無効化要求でキャッシュSSD200からデータが無効化されるため、FMへのライトが発生しない。この一連の処理の詳細については後述する。

[0048] FMへのディステージ優先度200526は、内部キャッシュ20051からデータを追いつ出す際に、直近アクセス時刻200525と合わせて、内部キャッシュ20051からFMへディステージすべきデータの選択に用いられる。ディステージ優先度200526には、高(High)または低(Low)のいずれかの情報が格納される(また、高(High)の情報を格納する代わりに「1」を格納する、あるいは低(Low)の情報を格納する代わりに「0」を格納するようにしてもよい)。なお、以降では、SSDコントローラ201による、内部キャッシュ20051からFMへのデータディステージのことを、内部ディステージと呼ぶ。内部ディステージする必要がある時、ディステージ優先度200526の高いデータの中から、LRUアルゴリズムにより内部ディステージするデータを選択する。ディステージ優先度200526の低いデータについては、ディステージ優先度200526の高いデータがない時にのみ、内部ディステージの対象として選ばれる。

[0049] ただし、内部ディステージポリシーは上記のポリシーに限定されない。たとえば直近アクセス時刻200525が一定時間以上前であった場合には、ディステージ優先度200526の高いデータに加えてディステージ優先度20

0526の低いデータも内部ディスティージ候補とする、などのポリシーを採用しても良い。また、ディスティージ優先度200526に、高(High)と低(Low)の2段階の情報を格納する態様のほか、複数段階の情報を格納するようにしてもよい。

[0050] 本発明の実施例1に係るストレージ装置10の場合、ストレージコントローラ100のCPU120が各ライトデータについて、シーケンシャルアクセスデータかランダムアクセスデータかを判定し(判定の方法には、ホスト計算機40から連続領域に対するアクセス要求が到来しているか否かを判定する等の、公知の手法を用いることができる)、キャッシュSSD200へライト要求する際に、ライトコマンド中に(具体的には後述するデータ情報3005に)、「ランダムアクセスデータ」または「シーケンシャルアクセスデータ」のいずれかの情報を格納してSSDコントローラ201へ送信する。シーケンシャルアクセスの場合、ライトデータに対する再アクセスが行われる可能性が低い(つまりキャッシュヒットする可能性が低い)。そのため、キャッシュDRAM1302およびキャッシュSSD200に格納されたシーケンシャルアクセスデータは、ヒットすることなく最終記憶媒体190にディスティージされる。そのため、SSDコントローラ201は、シーケンシャルアクセスデータであると通知されたデータはキャッシュヒットする可能性が低いと予測し、そのデータの内部キャッシュ20051への格納時に、内部ディスティージ優先度200526を「低」と設定する。一方、ランダムアクセスデータの場合には内部ディスティージ優先度200526を「高」と設定し、シーケンシャルアクセスのデータに比べて、内部ディスティージされやすくする(FMへライトされやすくする)。これにより、相対的にシーケンシャルアクセスデータを、できるだけ内部キャッシュ20051に留めて、内部ディスティージされにくくしている。

[0051] 続いて図5を用いて、キャッシュSSD200の論理物理アドレス変換/更新(論物変換/論物更新)について説明する。前述したとおり、FM2011は同物理アドレスに対して上書きができない。そのため、キャッシュS

SD200はライトコマンドを受け付けると、ライトコマンドで指定されている論理アドレスに対して新たな（未書き込みの）FM2011の物理アドレスを割り当て、当該割り当てられた物理アドレスに対してライトを行う。そのために、ストレージコントローラ100から指定される論理アドレス200531と、当該論理アドレス200531に対して割り当てられたFM2011上物理アドレス200532のマッピングを、内部メモリ2005上の論理物理アドレス変換テーブル20053に格納して、マッピングの管理を行う。SSDコントローラ201がリード／ライトを受けた際に、CPU2003が論理物理アドレス変換テーブル20053の参照、更新を行う。

[0052] SSDコントローラ201がリード要求を受けた場合、CPU2003は、論理物理アドレス変換テーブル20053を参照して、リード要求のあった論理アドレス200531に対応する物理アドレス200532を取得する。この処理を論物変換と呼ぶ。その後、CPU2003は、その物理アドレス200532を用いて、該当するFM2011に対してリード要求を行い、FM2011からデータをリードして、インタフェースコントローラ2001を介してストレージコントローラ100へリードデータを転送する。

[0053] SSDコントローラ201がライト要求を受けた場合、CPU2003はライトデータを格納するための、FM2011上の未書き込み領域を確保し、確保された領域のあるFM2011にデータをライトする。その後、CPU2003は、論理物理アドレス変換テーブル20053を参照して、ライト要求のあった論理アドレス200531に対応する物理アドレス200532を、新たにデータを格納した領域の物理アドレスへと更新する。この処理を論物更新と呼ぶ。

[0054] なお、図4の説明の際に述べたとおり、キャッシュSSD200では、データ処理アシストユニット2006によるデータ処理（圧縮や暗号化）を行う際に、4KB単位でデータ処理を行うため、図5の論理物理アドレス変換テーブル20053における論理アドレスの管理も4KB単位で行う。

- [0055] 次に、図6を用いてストレージコントローラ100のCPU120からキャッシュSSD200へ発行されるライトコマンド（ライト要求コマンド）について説明する。
- [0056] ライト要求コマンド3000には、コマンドを識別するコマンド番号（コマンドNo）3001、ライトであることを示す依頼内容3002、ライト先の論理空間上先頭アドレス3003、ライトデータのデータサイズ3004、SSDコントローラ201がFMへのディスページ優先度を設定するためのヒントとなるデータ情報3005が含まれる。依頼内容3002には、コマンドの種類を表す情報（ライトであることを示す情報）や、ライト対象データが現在格納されている領域のアドレス（たとえばキャッシュ領域1302のアドレス）のほか、データ処理アシストユニット2006によって、データ圧縮や暗号化などを行うか否かを指定する情報が含まれる。データ情報3005には先に述べたとおり、「ランダムアクセスデータ」または「シーケンシャルアクセスデータ」のいずれかの情報が格納される。
- [0057] キャッシュSSD200によるライト処理が終了（正常終了または失敗）した後、SSDコントローラ201はライト要求コマンド3000に対する応答コマンド4000をストレージコントローラ100のCPU120へ送信する。応答コマンド4000は、コマンドを識別するコマンド番号4001、完了または失敗などの情報を示す結果情報4002から成る。コマンド番号4001には、ライト要求コマンド3000に含まれているコマンド番号3001と同内容の情報が格納される。これにより、ストレージコントローラ100のCPU120は、自身が発行したライト要求コマンド3000によるライト要求が正常に終了したのか失敗したのかを判別できる。
- [0058] 続けて、図7を用いてストレージコントローラ100のCPU120からキャッシュSSD200へのリードまたはデータ無効化要求するためのコマンドについて説明する。ストレージコントローラ100のCPU120は、SSDコントローラ201にリードまたはデータ無効化要求コマンド5000を送信する。リードまたはデータ無効化要求コマンド5000は、コマン

ドを識別するコマンド番号5001、リードまたはデータ無効化であることを示す依頼内容5002、リードまたは無効化先の先頭アドレス5003、リードまたは無効化要求サイズ5004から成る。依頼内容5002には、コマンドの種類を表す情報（リード要求であること、あるいは無効化要求であることを示す情報）のほか、データ処理アシストユニット2006によって、リードデータを伸長あるいは復号化してストレージコントローラ100に返送するか否かを指定する情報が含まれる。

[0059] キャッシュSSD200が要求に対する処理を完了または失敗した後、応答コマンド6000がSSDコントローラ201からストレージコントローラ100のCPU120へ送信される。応答コマンド6000は、ライト要求コマンド3000に対する応答コマンド4000と同様、コマンドを識別するコマンド番号6001、完了または失敗などの情報を示す結果情報6002から成る。

[0060] 続いて図8を用いて、ホスト計算機40からのライト要求を受け取ったときにストレージ装置10で行われる処理の流れを説明する。ここでは、ストレージコントローラ100のCPU120によるキャッシュ先選択について、キャッシュDRAM1302を1次キャッシュとし、キャッシュSSD200を2次キャッシュとする例について説明する。また、以降では断りのない限り、ホスト計算機40からのリード、ライト要求で指定される領域が、ホスト計算機40がアクセスする最終記憶媒体190（あるいは1以上の最終記憶媒体190から構成されるボリューム）の記憶空間上の4KB境界に一致している場合の例（これはデータ処理アシストユニット2006によるデータ処理の最小単位が4KBのため）について説明する。また、図8の処理は特に断りのない限り、ストレージコントローラ100のCPU120が実行する処理である。

[0061] なお、ストレージコントローラ100のCPU120がキャッシュSSD200に発行するライト要求コマンド3000、リード要求コマンド5000と区別するため、以下ではホスト計算機40から受信するライト要求のこ

とを、ホストライト要求と呼ぶこともあり、またホスト計算機40から受信するリード要求のことを、ホストリード要求と呼ぶこともある。

[0062] ストレージコントローラ100のCPU120は、ホスト計算機40からホストライト要求を受けると(S10)、キャッシュ管理テーブル1301を用いてキャッシュDRAM1302のヒット判定を行う(S20)。キャッシュDRAM1302にヒットした、つまりキャッシュDRAM1302上にホストライト要求で指定されたデータを格納するためのキャッシュセグメントが確保済みの場合(S20:ヒット)、キャッシュDRAM1302上の該当セグメントに対してライトデータの上書きを行う(S60)。キャッシュDRAM1302でミスした場合(S20:ミス)、キャッシュDRAM1302に空き(未使用)のセグメントがあるかを確認する(S30)。空きがない場合(S30:No)、空きセグメントを確保するため、ホストライト要求で指定されたデータを格納するために必要となるサイズ以上のデータを、キャッシュDRAM1302から2次キャッシュであるキャッシュSSD200へとディステージする(S40)。このディステージ処理については後述する。

[0063] 空きがある場合(S30:Yes)またはデータディステージ後(S40の実行後)、キャッシュDRAM1302に該当する新たなセグメントを確保してキャッシュDRAM1302にライトデータを格納する(S50)。S50において、セグメントを確保すると、キャッシュ管理テーブル1301の内容を更新する。

[0064] キャッシュDRAM1302にデータ格納後、キャッシュ管理テーブル1301を用いてキャッシュSSD200にヒットしているか(ホストライト要求によってライトされるライト対象データを格納するためのキャッシュセグメントが確保済みか)を確認する(S70)。ヒットしている場合(S70:ヒット)には、ストレージコントローラ100はキャッシュSSD200に対して該当データの無効化を要求する(S80)。なぜなら、ヒットしている場合とは、キャッシュSSD200に、ホストライト要求によってラ

イトされるライト対象データの更新前データが格納されていることを意味するが、更新前データは不要なデータであるため無効化処理を行う。キャッシュSSD200でのデータ無効化処理については後述する。

[0065] 無効化要求が完了、またはミス（S70：ミス）で無効化の必要がない場合、ストレージコントローラ100はホスト計算機40へライト完了を通知する（S90）。また、キャッシュ管理テーブル1301を用いたキャッシュSSD200のヒット判定（S70）については、S20におけるキャッシュDRAM1302のヒット判定の際にキャッシュSSD200のヒット判定を行うようにしてもよい。

[0066] 図9を用いて、キャッシュDRAM1302の空きを確保するために行われる、キャッシュSSD200へのディステージ（S40）処理について説明する。この処理も、ストレージコントローラ100のCPU120によって実行される。

[0067] まずストレージコントローラ100のCPU120は、キャッシュDRAM1302からキャッシュSSD200へディステージすべきデータを選択する（S41）。この選択はすでに述べたように、キャッシュ管理テーブル1301を用いて、直近アクセス時刻（最終アクセス時刻）が最も古いデータから順に選択する。また、このディステージ処理は、ホストライト要求で指定されたデータを格納するために必要な領域を確保するための処理であるため、S41では、ディステージすべきデータを複数個（少なくともホストライト要求で指定されたデータと同量またはそれより大きい量）選択することもある。

[0068] 続いてキャッシュSSD200へディステージされるデータについて、キャッシュ管理テーブル1301を用いて、キャッシュSSD200のヒット判定、つまりディステージ対象のデータを格納するためのセグメントがキャッシュSSD200に確保済みか確認する処理を行う（S42）。ミスである場合（S42：ミス）、キャッシュSSD200に空きセグメントがあるかを確認する（S43）。空きがない場合（S43：No）、キャッシュS

SD200内のデータを最終記憶媒体190へとディステージする(S44)。このキャッシュSSD200から最終記憶媒体190へのディステージ処理については後述する。このディステージ処理後、またはキャッシュSSD200に空きがあった場合(S43:Yes)、キャッシュSSD200に該当するキャッシュセグメントを新たに確保し(S45)、キャッシュSSD200に対してライト要求コマンド3000を送信することで、キャッシュDRAM1302からキャッシュSSD200にディステージするデータのライトを行う(S46)。このライト要求コマンドを送信する際、データ情報3005に、「ランダムアクセスデータ」または「シーケンシャルアクセスデータ」のいずれかの情報を格納して送信する。S47で、ライト要求の完了結果である応答コマンド4000をキャッシュSSD200から受領し、キャッシュDRAM1302からのディステージ処理を完了する。

[0069] 図10を用いて、キャッシュSSD200の空きを確保するために行われる、最終記憶媒体190へのディステージ処理(S44)の流れについて説明する。

[0070] ストレージコントローラ100のCPU120は、キャッシュ管理テーブル1301を用いて、すでに述べたようにキャッシュSSD200からディステージするデータをLRUなどで選択する(S441)。続いてCPU120は、選択されたデータについて、キャッシュSSD200からデータをリードし(S442、S443、S444)、リードしたデータを最終記憶媒体190へライトする(S445、S446、S447)。その後、データは最終記憶媒体190に格納されたため、CPU120はキャッシュSSD200上のデータを無効化し(S448、S449、S4410)、最後に無効化したキャッシュSSD200上のセグメントについての情報をキャッシュ管理テーブル1301から削除することにより、セグメントの解放処理を行う(S4411)。キャッシュSSD200からのリード(S443)およびキャッシュSSD200上データの無効化要求(S449)については、図7に記載のコマンドを用いるが、この処理の詳細については後述す

る。

[0071] 図11を用いて、キャッシュSSD200で実行されるリード処理(S443)を説明する。

[0072] SSDコントローラ201のCPU2003は、リード要求コマンド5000を受領すると(S4431)、内部キャッシュ管理テーブル20052を用いて、内部キャッシュ20051のヒット判定を行う(S4432)。内部キャッシュ20051にヒットした場合(S4432:ヒット)、CPU2003は内部キャッシュ20051からデータをリードする(S4435)。内部キャッシュ20051でミスした場合(S4432:ミス)、CPU2003は前述の論物変換を行って(S4433)、FM2011からデータをリードする(S4434)。その後CPU2003は必要に応じてリードしたデータにキャッシュSSD200のデータ処理機能を適用し(S4436)、ストレージコントローラ100へデータを転送し(S4437)、CPU120へ完了通知する(S4438)。S4436の処理では、CPU2003はデータ処理アシストユニット2006を用いてデータ処理を行う。

[0073] 図11の例では、キャッシュSSD200の内部キャッシュ20051にリードデータをキャッシュしない。しかし、内部キャッシュ20051にリード時のデータをキャッシュしてもよい。その場合、FMからデータをリードする前に、内部キャッシュ20051の空きを確認し、空きがなければ内部ディスティージを行って空きを確保する。その後、空きセグメントを新たに確保し、FMからリードするデータをストレージコントローラ100へ転送するとともに、内部キャッシュ20051に格納する。また、本発明の実施例1に係るストレージ装置10では、ストレージコントローラ100が最終記憶媒体190にディスティージするために、キャッシュSSD200からデータを読み出す際には、データ変換/逆変換処理は行わせない。たとえば、キャッシュSSD200に圧縮データが格納されている場合、圧縮されたままのデータを読み出して最終記憶媒体190にディスティージするので、スト

レージコントローラ100が最終記憶媒体190にディステージするために、キャッシュSSD200に対してリード要求コマンド5000を発行する場合、リード要求コマンドの依頼内容5002には、リードデータの変換／逆変換は行わない旨が指定されている（そのため、上で説明したS4436の処理は行われぬ）。逆に、ストレージコントローラ100がホスト計算機40へ返送するためのリードデータをキャッシュSSD200から読み出す場合（かつ、キャッシュSSD200に、データ処理アシストユニット2006によって変換処理が施されたデータが格納されている場合）には、データ逆変換（伸長や復号化等）を行ったデータを読み出す必要がある。そのため、ストレージコントローラ100のCPU120がキャッシュSSD200に対して発行するリード要求コマンド5000の依頼内容5002には、リードデータの逆変換を行う旨が指定されている（そのため、上で説明したS4436の処理が行われる）。

[0074] 続いて、図12を用いて、キャッシュSSD200がデータ無効化要求を受信した時に行われる処理（S449）の流れを説明する。

[0075] SSDコントローラ201のCPU2003は、データ無効化要求コマンド5000を受領すると（S4491）、内部キャッシュ管理テーブル20052を用いて、内部キャッシュ20051のヒット判定を行う（S4492）。内部キャッシュにヒットした場合（S4492：ヒット）、CPU2003は内部キャッシュ20051の該当セグメントを解放する（S4493）。セグメント解放後または内部キャッシュ20051でミスした場合（S4492：ミス）、CPU2003は前述の論物変換を行って（S4494）、論理物理アドレス変換テーブル20053の該当テーブルエントリを確認する（S4495）。これは、内部キャッシュ20051およびFM2011の両方にデータが格納されていることがあり、その両方を無効化する必要があるため内部キャッシュヒット時にも論物変換を行う。FM2011にデータが格納されている場合は、該当テーブルエントリは有効であるため（S4495：有効）、CPU2003は該当物理アドレスのデータが無効

である旨をブロック情報20054に記録し(S4496)、さらに論理物理アドレス変換テーブル20053から、論理物理アドレス変換テーブル20053に記録されている該当物理アドレスの情報を削除する(S4497)。その後CPU2003は、ストレージコントローラ100のCPU120へ完了通知を行い(S4498)、処理を終了する。

[0076] 続いて図13を用いて、ストレージコントローラ100からライト要求コマンド(S46)を受信したときに、キャッシュSSD200が実行する処理について説明する。

[0077] SSDコントローラ201のCPU2003は、ライト要求コマンド3000を受領すると(S461)、コマンド3000の依頼内容3002(たとえば圧縮の指示、暗号化の指示など)に従って、必要に応じてデータ処理アシストユニット2006を用いて、ライトデータにデータ変換処理を施す(S462)。なお、本発明の実施例1に係るストレージ装置10では、ストレージコントローラ100がホスト計算機40から受信したライトデータをキャッシュSSD200に格納する場合には、キャッシュSSD200にデータ変換を行わせる。そのため、ストレージコントローラ100がホスト計算機40から受信したライトデータをキャッシュSSD200に格納するために、ライト要求コマンド3000を発行する場合には、当該ライト要求コマンドの依頼内容3002にはデータの変換を行う旨が指定されている。一方、最終記憶媒体190に格納されていたデータをキャッシュSSD200に格納(ステージング)する場合には、キャッシュSSD200にデータ変換/逆変換処理は行わせない。そのため、この場合にストレージコントローラ100のCPU120がキャッシュSSD200に発行するライト要求コマンドの依頼内容3002には、データの変換/逆変換は行わない旨が指定されている(そのため、上で説明したS462の処理は行われない)。

[0078] 続いてCPU2003は内部キャッシュ管理テーブル20052を用いて、内部キャッシュ20051のヒット判定を行う(S463)。内部キャッシュにヒットした場合(S463:ヒット)、CPU2003は内部キャッ

シュの該当セグメントにデータを上書きする（S467）。内部キャッシュ20051でミスした場合（S463：ミス）、CPU2003は内部キャッシュ20051に空きの（未使用の）セグメントがあるか確認する（S464）。空きがない場合（S464：No）、空きセグメントを確保するため、CPU2003は内部キャッシュ20051上のデータをFM2011に内部ディステージする（S465）。内部ディステージ処理については後述する。なお、S464の判定において、未使用のセグメントがなくなった場合に内部ディステージ（S465）処理を実施する態様に限定されるわけではない。たとえば、未使用の領域サイズ（セグメント数）が不足した場合（ストレージコントローラ100からライト要求のあったデータのサイズよりも未使用の領域のサイズが小さい場合、あるいは未使用の領域のサイズが所定の閾値未満になった場合など）に内部ディステージ（S465）処理を行うようにしてもよい。

[0079] 内部ディステージ後、または内部キャッシュ20051に空きがある場合（S464：Yes）、CPU2003は内部キャッシュ20051のセグメントを新たに確保し（S466）、内部キャッシュ20051にデータをライトする（S467）。なお、S467では、CPU2003はS461で受信したライト要求コマンド300のデータ情報3005に含まれている「ランダムアクセスデータ」または「シーケンシャルアクセスデータ」の情報に基づいて、ライトデータの格納されるキャッシュセグメントに対応する、内部キャッシュ管理テーブル20052のディステージ優先度200526の欄に、「高」または「低」を格納する。その後CPU2003は、ストレージコントローラ100のCPU120に応答コマンド4000を送信することでライト完了通知を行う（S468）。

[0080] 続いて図14を用いて、キャッシュSSD200で行われる内部ディステージ処理（S465）について説明する。

[0081] 内部キャッシュ20051からデータをFM2011へ内部ディステージする必要がある場合、SSDコントローラ201のCPU2003は、内部

ディスティージングするデータを選択する（S 4 6 4 1）。先に述べた通り、内部ディスティージング対象データの選択は、一例として、内部キャッシュ管理テーブル20052の直近アクセス時間200525とディスティージング優先度200526を用いて、ディスティージング優先度200526が「高」であるデータの中から直近アクセス時間200525の最も古いものを選択することによって行う。そしてディスティージング優先度200526が「高」であるデータが存在しない場合には、ディスティージング優先度200526が「低」のデータの中から直近アクセス時間200525の最も古いものを選択する。ただし、これ以外のデータ選択方法を採用してもよい。

[0082] 内部ディスティージングするデータを選択後、CPU2003は内部ディスティージング先となるFM2011の空き（未使用）領域（物理アドレス）を探す（S 4 6 4 2）。内部ディスティージング対象データを格納可能な空き領域（物理アドレス）が複数存在する場合、CPU2003はFM2011の寿命等を考慮して適切な物理アドレスを選択する。そしてCPU2003は、選択した書き先物理アドレスに対応するFM2011へのライト要求を行い（S 4 6 4 3）、書き先物理アドレスのデータが有効である旨をブロック情報20054に書き込む（S 4 6 4 4）。さらにCPU2003は、要求された論理アドレスに書き先物理アドレスに対応付けるために論理物理アドレス変換テーブル20053を更新する（S 4 6 4 5）。以上でデータをFM2011に格納したので、CPU2003は内部キャッシュ管理テーブル20052を更新して内部キャッシュ20051の該当セグメントを解放し（S 4 6 4 6）、内部ディスティージング処理を終了する。

[0083] 次に、図15を用いて、ホスト装置40からのリード要求（ホストリード要求）を受信した時に、ストレージ装置10が行う処理について説明する。

[0084] ストレージコントローラ100のCPU120は、ホスト計算機40からホストリード要求を受領すると（S 1 0 0）、キャッシュ管理テーブル1301を用いて、キャッシュヒット判定を行う（S 1 1 0、S 1 2 0）。キャッシュDRAM1302にヒットした場合（S 1 1 0：ヒット）、CPU1

20はキャッシュDRAM1302からデータをリードする(S190)。キャッシュDRAM1302にヒットせず(S110:ミス)、キャッシュSSD200にヒットした場合(S120:ヒット)、CPU120はキャッシュSSD200へリード要求コマンドを発行する(S140)。キャッシュSSD200で行われるリード処理は、S443で行われる処理と同様で、データが読み出される過程で、キャッシュSSD200はデータにキャッシュSSD200内部のデータ処理機能を適用し、データに対する処理(データ伸長やデータ復号化)を行う。そして処理が施されたデータがストレージコントローラ100のバッファ領域1303に転送される。

[0085] キャッシュSSD200でもミスした場合(S120:ミス)、CPU120は最終記憶媒体190からデータをリードし、メモリ130、ASIC140、あるいはIFCTL150内の一時記憶領域にデータを格納する(S130)。さらにCPU120は、キャッシュSSD200に空きセグメントが存在するか確認し(S150)、空きがないなら(S150:No)、キャッシュSSD200から最終記憶媒体190へディステージ処理を行う(S160)。なお、このキャッシュSSD200の空き確認とディステージ処理(S150、S160)については、最終記憶媒体190からのリード(S130)の前に行うようにしてもよい。ディステージ処理(S160)は、S44と同様である。ディステージ後、またはキャッシュSSD200に空きがある場合(S150:Yes)、CPU120はキャッシュSSD200に該当するキャッシュセグメントを新たに確保し、最終記憶媒体190からリードしたデータをライトする要求をキャッシュSSD200に発行する(S170)。キャッシュSSD200へのライト処理はS46と同様である。続いてS170でCPU120は、キャッシュSSD200に書き込んだデータをリードする要求を、キャッシュSSD200に発行する(S180)。このリード要求を受け付けたキャッシュSSD200ではS443と同様、データに対する処理(データ伸長やデータ復号化)を行い、処理が施されたデータがストレージコントローラ100のバッファ領域1

303に転送される。最後に、ストレージコントローラ100はいずれかの媒体からリードしたデータをホスト計算機40へ転送し(S200)、処理を終了する。

[0086] 上で述べたように、いずれのキャッシュでもミスし、最終記憶媒体190からリードした場合(S130)、CPU120は最終記憶媒体190からリードしたデータを一旦キャッシュSSD200にライトし、その後キャッシュSSD200にライトしたデータをリードする。最終記憶媒体190上のデータは、ライト時にキャッシュSSD200のデータ処理機能が適用されて、圧縮や暗号化等の処理がなされているので、一度キャッシュSSD200に格納し、再度データ処理機能を適用してデータ伸長あるいは復号化を行うことでホスト計算機40が扱えるデータにする必要があるためである。そのため、ホスト計算機40に送信するデータを、最終記憶媒体190から一旦キャッシュSSD200にライトする必要があるが、内部キャッシュ20051にヒットする、或いは内部キャッシュ20051の空き領域に書き込める可能性が高いため、FM2011へのライトが発生する可能性は低く、リードレイテンシも小さい。

[0087] なお、たとえばデータ処理機能によりデータ暗号化・復号化を行う場合、データ暗号化・復号化のために必要な暗号鍵などの情報については、ライト時からキャッシュSSD200が保持し続けてもよいし、ストレージコントローラ100が管理してもよい。ストレージコントローラ100が管理する場合、キャッシュSSD200へのライト要求コマンド3000、リード要求コマンド5000によって必要な情報をキャッシュSSD200に通知する。あるいはライト要求コマンド3000やリード要求コマンド5000とは異なる、別のコマンドを発行することで、必要な情報をキャッシュSSD200に通知するようにしてもよい。

[0088] また、上では、ホスト計算機40からのホストリード要求、ホストライト要求で指定される領域が、4KB境界に一致している場合の例について説明したが、ホスト計算機40からのリード、ライト要求で指定される領域が4

K B境界に一致しておらず、かつキャッシュDRAMあるいはキャッシュSSDにアクセス対象のデータが格納されていない（キャッシュミス）の場合には、最終記憶媒体190からアクセス対象の範囲を含む、4KB境界に一致したサイズのデータをキャッシュDRAMあるいはキャッシュSSDに読み出して、キャッシュされたデータに対して、リード、ライト処理を行えばよい。

実施例 2

[0089] 実施例2に係るストレージ装置10及びキャッシュSSD200のハードウェア構成は、実施例1で説明したものと同一であるため、ここでは説明しない。また、実施例2に係るストレージ装置10及びキャッシュSSD200で実施される大半の処理も、実施例1で説明したものと同一であるため、ここでは、実施例1と異なる点（具体的には、実施例1における、キャッシュSSD200の内部ディスティージ処理（図13：S465、図14）、そしてストレージコントローラ100のCPU120がキャッシュSSD200にライト要求を発行する際の処理（図9：S46）が、実施例2では異なる処理になる）を中心に説明する。

[0090] 実施例1に係るストレージ装置10及びキャッシュSSD200では、ディスティージ優先度200526を高く設定したデータが優先的にFM2011にライトされる。これにより、ディスティージ優先度200526が低いデータが極力FM2011にライトされないようにし、FM2011の寿命劣化を最小限に抑える。ただし、ディスティージ優先度200526が高いデータがキャッシュSSD200の内部キャッシュ20051に存在しない場合には、ディスティージ優先度200526が低いデータもFM2011にライトされる。一方、実施例2に係るストレージ装置10及びキャッシュSSD200では、ディスティージ優先度200526が高いデータのみをFM2011にライトするように制御する。

[0091] しかし、そのままでは内部ディスティージ優先度200526が低いデータが内部キャッシュ20051を占有してしまうため、ストレージコントロー

ラ100がキャッシュSSD200から最終記憶媒体190へデータをディスステージする必要がある。そのためにキャッシュSSD200からストレージコントローラ100のCPU120へディスステージが必要であることを通知する。この通知は、キャッシュSSD200がライト要求を受領したタイミングで行うが、それ以外の契機で通知するようにしてもよい。例えば、SSDコントローラ200が内部キャッシュ20051の空き容量を定期的に監視し、内部キャッシュ20051の空き容量が小さくなった際にストレージコントローラ100のCPU120にディスステージが必要であると通知してもよい。

[0092] 図16は、実施例2におけるキャッシュSSD200で行われる内部ディスステージ処理のフローの一例である。

[0093] SSDコントローラ201のCPU2003は、内部キャッシュ20051内のデータについて、内部ディスステージ優先度200526が高いデータの有無を確認し(S46411)、そのようなデータがない場合(S46411:No)、SSDコントローラ201から最終記憶媒体190へのディスステージが必要である旨、及び内部キャッシュ20051内のデータのうち、最終アクセス時刻(内部キャッシュ管理テーブル20052で管理しているデータの最終アクセス時刻200525)が最も古いデータ(キャッシュセグメントに格納されているデータ)の情報を、ライト要求コマンド3000に対する応答コマンド4000の結果情報4002に含めて、ストレージコントローラ100のCPU120へ通知し(S46412)、キャッシュSSD200の内部ディスステージ処理は終了する。なお、詳細は後述するが、ストレージコントローラ100のCPU120は、SSDコントローラ201から返送された応答コマンド4000の結果情報4002に、最終記憶媒体190へのディスステージが必要である旨の情報が含まれていた場合、キャッシュSSD200内のデータを最終記憶媒体190へディスステージする処理を実施する。

[0094] 一方、内部ディスステージ優先度200526が高いデータがある場合(S

46411: Yes)、CPU2003は実施例1における内部ディスク処理(S4641~S4646)と同様の処理を行う。

[0095] 続いて図17を用いて、ストレージコントローラ100のCPU120がキャッシュSSD200へライト要求を発行したときに、ストレージコントローラ100のCPU120の実行する処理の流れを説明する。

[0096] 実施例1と同様に、ストレージコントローラ100のCPU120がキャッシュSSD200へライト要求コマンド3000を送信することでライト要求(S46)した後、キャッシュSSD200からライト要求コマンド3000に対する応答コマンド4000でライト完了通知を受領する(S47)。本実施例では、CPU120はその結果情報4002からライト完了またはディステージが必要なのかを確認する(S48)。結果情報4002に、最終記憶媒体190へのディステージが必要である旨の情報が含まれていた場合(S48:ディステージ要)、CPU120はキャッシュSSD200からデータを最終記憶媒体190へディステージする(S49)。このディステージ処理は実施例1のS44(つまり図10に記載の処理)とほぼ同様であるが、CPU120はディステージするデータとして、キャッシュSSD200から通知された、内部キャッシュ20051内のデータのうち最終アクセス時刻が最も古いデータを選択するようにする。これは、内部キャッシュ20051のデータをディステージしなければ、キャッシュSSD200に新たにデータを書き込めないからである。ディステージ後、CPU120は再度キャッシュSSD200へライト要求を行い、ライトが完了するまで上記の処理を繰り返す。

[0097] 実施例2に係るストレージ装置10及びキャッシュSSD200では、ディステージ優先度が低いデータのFM2011への書き込みが完全に抑止される。そのため、キャッシュSSDの記憶媒体(FM)に書き込まれるデータ量を削減することができ、FM2011の寿命を延ばすことができる。

実施例 3

[0098] 実施例1または2に係るストレージ装置10では、ストレージコントロー

ラ100上のメモリ130（具体的にはメモリ130内のキャッシュDRAM1302）を1次キャッシュとし、キャッシュSSD200を2次キャッシュとして用いていたが、実施例3に係るストレージ装置10では、1次キャッシュを用いない。つまりキャッシュSSD200のみをストレージ装置10のキャッシュとして用いる点が、実施例1または2に係るストレージ装置10と異なる。

[0099] 実施例3に係るストレージ装置10の構成は、実施例1で説明したもの（図1に記載のもの）とほぼ同様であるため、図示は省略する。ただし実施例3に係るストレージ装置10の場合、メモリ130内にメモリキャッシュ領域1302が存在しない点が、実施例1に係るストレージ装置10と異なる。その他の点は実施例1で説明したものと同一である。

[0100] 続いて実施例3に係るストレージ装置10が管理する、キャッシュ管理テーブル1301について説明する。実施例3に係るキャッシュ管理テーブル1301に含まれる情報は、実施例1におけるものと同じであるので図示は省略する。ただし、実施例1におけるキャッシュ管理テーブル1301には、キャッシュDRAM1302のキャッシュセグメントについての情報（図2：1310）が格納されていたが、実施例3に係るストレージ装置10では、キャッシュとして用いる記憶領域は、キャッシュSSD200だけであるため、キャッシュSSDのキャッシュセグメントについての情報（図2：1311）だけが格納され、キャッシュDRAM1302のキャッシュセグメントについての情報（図2：1310）は格納されない。

[0101] なお、実施例3に係るキャッシュSSD200の構成、及びキャッシュSSD200が内部メモリ2005で管理する情報（内部キャッシュ管理テーブル、論理物理アドレス変換テーブル20053、ブロック情報20054）は、実施例1に係るキャッシュSSD200と同一であるため、ここでは説明を行わない。

[0102] 続いて図18を用いて、実施例3に係るストレージ装置10が、ホスト計算機40からライト要求を受信した時の処理の流れを説明する。この処理は

、一般的なディスクキャッシュ付きストレージ装置が行う処理と同様の処理である。

[0103] ストレージコントローラ100のCPU120は、ホスト計算機40からライト要求を受信する(S10)。この処理は図8を用いて説明した処理S10と同じである。続いてストレージコントローラ100は、キャッシュ管理テーブル1301を用いてキャッシュSSD200のヒット判定を行う(S20')。キャッシュSSD200にヒットした、つまりキャッシュSSD200上に、ライト要求で指定されたデータを格納するためのキャッシュセグメントが確保済みの場合(S20':ヒット)、キャッシュSSD200上の該当セグメントに対してライトデータの上書きを行う(S60')。キャッシュSSD200でミスした場合(S20':ミス)、キャッシュSSD200に空き(未使用)のセグメントがあるかを確認する(S30')。空きがない場合(S30':No)、空きセグメントを確保するため、キャッシュSSD200から最終記憶媒体190へデータのディステージを行う(S40')。このディステージ処理は、実施例1で説明した図10の処理(図9:S44の処理)と同じ処理である。

[0104] 空きがある場合(S30':Yes)またはデータディステージ後(S40'の実行後)、キャッシュSSD200に新たなセグメントを確保してキャッシュSSD200にライトデータを格納する(S50')。この処理は、実施例1で説明した、図9のS45、S46と同じ処理である。キャッシュSSD200にデータ格納後、ストレージコントローラ100はホスト計算機40へライト完了を通知し(S90)、処理を終了する。

[0105] 次に、図19を用いて、ホスト装置40からのリード要求を受信した時のストレージ装置10の行う処理について説明する。この処理は、図15を用いて説明した、実施例1に係るストレージ装置10が実施するリード処理から、キャッシュDRAM1302に対する処理を取り除いたものである。

[0106] 最初にストレージコントローラ100のCPU120は、ホスト計算機40からリード要求を受領すると(S100)、キャッシュ管理テーブル13

01を用いて、キャッシュヒット判定を行う（S120）。キャッシュSSD200にヒットした場合（S120：ヒット）、CPU120はキャッシュSSD200へリード要求する（S140）。キャッシュSSD200のリード処理は、実施例1で説明したS443の処理（つまり図11の処理）と同様である。キャッシュSSD200でキャッシュミスの場合（S120：ミス）、CPU120は最終記憶媒体190からデータをリードして、メモリ130、ASIC140、あるいはIFCTL150内の一時記憶領域に格納する（S130）。

[0107] さらにCPU120は、キャッシュSSD200に空きセグメントが存在するか確認し（S150）、空きがない場合（S150：No）、キャッシュSSD200から最終記憶媒体190へデータのディステージ処理を行う（S160）。このキャッシュSSD200の空き確認とディステージ処理（S150、S160）については、最終記憶媒体190からのリード（S130）の前に行うようにしてもよい。ディステージ処理（S160）は、S44と同様である。ディステージ後、またはキャッシュSSD200に空きがある場合（S150：Yes）、CPU120はキャッシュセグメントを新たに確保し、最終記憶媒体190からリードしたデータをキャッシュSSD200にライトする（S170）。キャッシュSSD200へのライト処理はS46と同様である。このライト処理の過程で、データにキャッシュSSD200内部のデータ処理機能が適用される。続いて、S170でCPU120はキャッシュSSD200に書き込んだデータをリードする要求を、キャッシュSSD200に発行する（S180）。このリード要求を受け付けたキャッシュSSD200ではS443と同様、データに対する処理（データ伸長やデータ復号化）を行い、処理が施されたデータがストレージコントローラ100のバッファ領域1303に転送される。最後にCPU120は、いずれかの媒体からリードしたデータをホスト計算機40へ転送し（S200）、処理を終了する。

[0108] なお、実施例3に係るキャッシュSSD200でも、実施例1と同様に、

内部キャッシュ領域20051に空きセグメントが存在しない場合、図14に記載された内部ディステージ処理が行われる。内部ディステージ処理では実施例1と同様、ディステージ優先度200526の高いデータの中から、LRUアルゴリズムにより内部ディステージするデータを選択することにより、シーケンシャルアクセスデータのように、再アクセスされる可能性が低いデータをFM2010に格納（内部ディステージ）されにくくしている。

[0109] また、実施例3に係るストレージ装置10においても、実施例2において説明した内部ディステージ処理（図16の処理）を採用してもよい。そうすると、実施例3に係るストレージ装置10は、実施例2に係るストレージ装置10と同様に、ディステージ優先度200526が高いデータのみをFM2011にライトするように動作する。

実施例 4

[0110] 続いて実施例4の説明を行う。実施例4に係るストレージ装置及びキャッシュSSDのハードウェア構成は、実施例1または2で説明したものと同一であるため、ここでは説明しない。実施例1または2に係るストレージ装置10では、ストレージコントローラ100上のメモリ130（具体的にはメモリ130内のキャッシュDRAM1302）を1次キャッシュとし、キャッシュSSD200を2次キャッシュとして用いていたが、実施例4では、ストレージコントローラ100のCPUがデータ特性（例えば、ライトデータの最終的な格納先）に応じて、キャッシュSSD200又はキャッシュDRAM1302のいずれか一方を選択してデータをキャッシュする。ライトデータの最終的な格納先については、ホストライト要求に含まれるライト対象データのアドレス情報（LBA等）を参照することで判断できる。

[0111] 最終記憶媒体190が、例えばSSDを始めとする高速な媒体である場合、アクセス速度が同等であるキャッシュSSD200にキャッシュしても、アクセス性能向上の効果が小さいため、ストレージコントローラ100のメモリ130上のキャッシュDRAM1302にのみ格納する。一方、最終記憶媒体190がHDDのような、SSDよりも低速な媒体である場合、キャ

ッシュSSD200にのみキャッシュする。この場合では、キャッシュDRAM1302に格納されたデータについては、キャッシュSSD200が備えるデータ処理機能を用いることができないため、ストレージコントローラ100のCPU120または別ハードウェアで同様の処理を行う。

- [0112] キャッシュSSD200にキャッシュされたデータの扱いは実施例3等と同等となるので詳細は省略する。
- [0113] 実施例4では、最終記憶媒体の種類などの第一のデータ特性に基づき、キャッシュDRAM1302又はキャッシュSSD200のいずれか一方にライトデータをキャッシュする。キャッシュSSD200に格納されたデータに関して、キャッシュSSD200が、ストレージコントローラ100のCPU120から受信した第二のデータ特性（例えば、ライトアクセスの種類（シーケンシャル／ランダム））に基づき、内部キャッシュからFMへのディスティージの優先度を決定する。
- [0114] 図20を用いて、ホスト装置40からのライト要求を受信した時のストレージ装置10の行う処理について説明する。
- [0115] ストレージコントローラ100のCPU120は、ホスト計算機40からライト要求を受信する（S10）。この処理は図8を用いて説明した処理S10と同じである。続いてCPU120は、ライト要求コマンドのLBAを確認して、ライトデータの最終格納先がSSDのような高速な媒体かHDDのような低速な媒体かを判定する（S11）。最終格納先がSSDであれば（S11：SSD）、ライトデータをキャッシュDRAM1302へ格納する（S13）。このキャッシュDRAM1302への格納処理については、図8のS70、S80を除いた処理となる。なぜなら、キャッシュDRAM1302とキャッシュSSD200は同じLBAのデータを保持しないため、キャッシュDRAM1302上のデータがキャッシュSSD200上にはないからである。
- [0116] 一方、最終格納先がHDDであれば（S11：HDD）、ライトデータをキャッシュSSD200に格納する（S12）。キャッシュSSD200へ

の格納処理は、図9のS42以降と同様である。

[0117] 実施例4に係るストレージ装置によれば、キャッシュSSD200には、HDDのように、SSDよりもアクセス性能の低い記憶媒体が用いられた最終記憶媒体190に格納されるデータのみがキャッシュされ、SSD等のようにキャッシュSSD200と同等のアクセス性能の記憶媒体が用いられた最終記憶媒体190に格納されるデータはキャッシュされない。そのため、キャッシュSSD200に書き込まれるデータ量を削減することができ、キャッシュSSD200の記憶媒体(FM2011)の寿命を延ばすことができる。

実施例 5

[0118] 続いて実施例5の説明を行う。実施例5に係るストレージ装置及びキャッシュSSDのハードウェア構成は、実施例1または2で説明したものと同一である。実施例5に係るストレージ装置及びキャッシュSSD200では、キャッシュSSD200のFM2011の劣化が進み、データを保持し続けることが難しくなった場合に、キャッシュSSD200が内部ディスティージを行わずに、実施例2と同様にストレージコントローラ100へディスティージを要求して、データを最終記憶媒体190へディスティージする。このとき、FM2011に格納されているデータのうち、キャッシュ状態13014が「dirty」のデータについても、ストレージコントローラ100へディスティージ要求する。これにより、キャッシュSSD200がキャッシュ可能な容量は小さくなるが、FM2011の劣化によってデータが消失してしまうことを防止しつつ、ライトデータまたは最終記憶媒体190からリードされたデータに対して、キャッシュSSD200の持つデータ変換等の機能を引き続き適用できる。

[0119] 図21を用いて、キャッシュSSD200のFM2011が寿命に到達した場合の動作を説明する。図21は、キャッシュSSD200がストレージコントローラ100からライト要求コマンドを受信したときに、キャッシュSSD200で実行される処理の流れを説明する図であり、多くの処理は実

施例1の図13と共通している。

- [0120] キャッシュSSD200は、実施例1と同様に、ストレージコントローラ100のCPU120からライト要求を受け取ると(S461)、データ処理を行い(S462)、内部キャッシュ管理テーブル20052を用いて、内部キャッシュ20051のヒット判定を行う(S463)。内部キャッシュにヒットした場合(S463:ヒット)、内部キャッシュの該当セグメントに上書きライトする(S467)。
- [0121] 内部キャッシュ20051でミスした場合(S463:ミス)、内部キャッシュ20051に空きの(未使用の)セグメントがあるか確認する(S464)。内部キャッシュ20051に空きがある場合(S464:Yes)、実施例1と同様に新たに内部キャッシュセグメントを確保し、データを内部キャッシュへライトする。一方、内部キャッシュ20051に空きのセグメントがない場合(S464:No)、FM2011が寿命に到達したかを判定する(S469)。ここでの判定方法として例えば、SSDコントローラ201のCPU2003が、上位装置からの要求とは無関係にデータを定期的にリードして、データに付与されているCRCなどの保証コードを用いてエラー率を確認し、エラー率が予め設定してある閾値を越えた場合に寿命に到達したと判定する、という方法があり得る。
- [0122] FM2011が寿命に到達していない場合(S469:No)、他の実施例で説明したライト処理と同様に内部ディスティージを行って(S465)、空きセグメントを確保する。空きセグメントを確保するための内部ディスティージを行うことができない場合、すなわちFM2011が寿命に到達した場合(S469:Yes)、内部ディスティージを行うとデータを消失する危険がある。そこで、実施例2と同様に、キャッシュSSD200からストレージコントローラ100のCPU120へディスティージが必要であることを要求する(S46412)。ただし、それ以外の契機でディスティージが必要であることをCPU120へ通知するようにしてもよい。実施例1と同様に、SSDコントローラ200が内部キャッシュ20051の空き容量を定期的

に監視し、内部キャッシュ20051の空き容量が小さくなった際にストレージコントローラ100のCPU120にディステージが必要であると通知してもよい。さらに、キャッシュSSD200はFM2011が寿命に到達したことをストレージコントローラ100に通知してもよい。

[0123] キャッシュSSD200のFM2011が寿命に到達した旨の通知を受けたストレージコントローラ100のCPU120は、すでにキャッシュSSD200に格納してあるデータについて、キャッシュ状態13014が「dirty」であるデータを最終記憶媒体190へディステージする。また、キャッシュ状態13014が「clean」であるデータについては、キャッシュSSD200に格納したままでは読み出せない可能性があるため、それらのデータについては、キャッシュ管理テーブル1301を更新し、キャッシュSSD200に格納されていないものとして扱う。具体的には、キャッシュの管理情報を、キャッシュSSD200に記録したデータが消失したものとして更新する。以降、キャッシュSSD200に格納されるデータは、内部キャッシュ20051に格納できないものについては、上述の通り、ストレージコントローラ100にディステージが必要であることを通知し、ストレージコントローラ100にディステージしてもらうため、寿命に到達したFM2011にデータが格納されることはない。

[0124] 本発明の実施例5に係るストレージ装置及びキャッシュSSD200によれば、ストレージコントローラ100のCPU120が、キャッシュSSD200からFM2011の寿命に基づく情報を取得し、キャッシュSSD200のDRAMなどからなる内部メモリ2005のデータを最終記憶媒体190へディステージすることで、信頼性を高めることを可能とする。

実施例 6

[0125] 続いて実施例6について説明する。実施例6に係るストレージ装置及びキャッシュSSDのハードウェア構成も、実施例1または2で説明したものと同様である。実施例6は、ストレージ装置が、複数の最終記憶媒体190をRAIDグループとして管理し、ホスト計算機40からのライトデータから

RAID Parity (以下、「Parity」と略記することもある)を生成し、ライトデータとともにParityを最終記憶媒体190に格納する構成を対象とする。

[0126] また実施例6に係るストレージ装置10では、キャッシュSSD200がストレージコントローラ100からの指示に基づいて、キャッシュSSD200に格納されたデータ(ホスト計算機40から書き込まれたライトデータ)からRAID Parityを生成する。この構成によれば、ストレージコントローラ100からRAID Parity生成に係る負荷をオフロードすることができ、ストレージ装置10全体の性能を向上させる事が可能となる。そのため、実施例6に係るキャッシュSSDは、Parityを生成する機能を備える。実施例1においても述べたとおり、この機能はデータ処理アシストユニット2006に備えられていてもよいし、あるいはデータ処理アシストユニット2006とは異なる、Parity生成用のハードウェアとしてキャッシュSSD200に実装されていてもよい。あるいは、CPU2003がParity生成を実行する構成であってもよい。

[0127] 実施例6に係るストレージ装置では、ストレージコントローラ100はホスト計算機40からのライトデータをキャッシュSSD200に格納した後、SSDに対しParity生成コマンドを発行する。

[0128] このParity生成コマンドは少なくとも、Parity生成対象とするデータの格納位置を指定するLBA(これは実施例1で説明した、キャッシュSSDがストレージコントローラに提供している論理的な記憶空間上のアドレスである)を複数個、生成されるParityの格納先を指定するLBAを1または2個(RAID5であれば1つ、RAID6であれば2つ)、そしてデータ長の情報を有している。

[0129] このParity生成コマンドをストレージコントローラ100のCPU120から受領したキャッシュSSD200は、コマンドにて指定されたLBAに対応付けられたデータを、FM2011または内部キャッシュ20051より取得し、Parity演算を実施する。そして取得されたデータを

用いて P a r i t y を生成し、内部キャッシュ 20051 に格納し、P a r i t y 生成コマンドにて指定された L B A に対応付けて管理する（具体的には、内部キャッシュ管理テーブル 20052 を用いて管理しておく）。

[0130] ストレージ装置は一般に、キャッシュの空き容量が低下し、新たなデータをキャッシュに載せる事ができない場合に、アクセス頻度の低いデータをキャッシュから HDD などの最終記憶媒体に転送するディステージ処理を行い、ディステージによって最終記憶媒体に記録されたデータをキャッシュから削除することで空き領域を生成する。またストレージ装置は一般に、ディステージされるべきと決定されたデータ（ディステージ対象データ）に対して P a r i t y を生成するため、生成された P a r i t y もディステージ対象データと同様に、生成後、短期間で最終記憶媒体に転送され、キャッシュから消去される。

[0131] このことから、本発明の実施例 6 に係るキャッシュ SSD でも、上で説明した P a r i t y 生成により生成された P a r i t y は、短期間でディステージされることが期待される。そのためキャッシュ SSD 200 は、生成された P a r i t y を内部キャッシュ 20051 に格納する際、当該 P a r i t y が内部ディステージ優先度の低いデータであるとして管理する（内部キャッシュ管理テーブル 20052 の、内部ディステージ優先度 200526 を「L o w」として記録する）。

[0132] これにより、ストレージコントローラが P a r i t y を最終記憶媒体 190 にディステージし、キャッシュ SSD 200 からのデータ（P a r i t y）消去が指示されるまでの期間、P a r i t y は優先的に D R A M 2005（内部キャッシュ 20051）に格納され、F M 2011 に書き込まれる（内部ディステージされる）確率を低下させることで、F M 2011 へのライト量を減らし、F M の劣化を軽減できる。

[0133] 以上で、本発明の実施例に係るストレージ装置及びキャッシュ SSD の説明を終了する。本発明の実施例に係るストレージ装置では、D R A M 等の高速の記憶媒体から成る内部キャッシュと、F M のような低価格で大容量の記

憶媒体を主たる記憶領域とするキャッシュ装置（キャッシュSSD）を備える。FMは、データ消去回数（書き換え回数）に制限があるという欠点があるため、高頻度の書き替えが行われると、寿命が短くなってしまふ。

[0134] 本発明の実施例に係るストレージ装置では、キャッシュ装置がデータの属性・特性に基づいて、キャッシュ対象データの内部ディステージの優先度を制御する。一例として、ストレージコントローラ等の上位装置がキャッシュ装置にデータを格納（キャッシュ）する際に、データの属性・特性に関する情報（たとえばシーケンシャルアクセス、ランダムアクセスである等）をキャッシュ装置に通知することにより、キャッシュ装置はデータの属性・特性を得る。これにより、たとえばシーケンシャルライトデータのように、上位装置から再アクセスされることなく追い出されてしまふデータ（キャッシュする効果が小さいデータ）が、内部キャッシュからFMにデータ移動（内部ディステージ）されることを抑制できるので、FMへのライト頻度を低下させ、キャッシュ装置の寿命（耐用年数）を延ばすことが可能である。

[0135] また、内部キャッシュからFMへのデータ移動（内部ディステージ）の優先度とは異なる基準（LRUなど）に基づき、データをキャッシュ装置から最終記憶媒体へディステージすることにより、内部キャッシュからFMを介さずに最終記憶媒体へディステージすることを可能とし、FMへのライト頻度を抑制している。

[0136] また、FMのような書き換え回数に制限のある記憶媒体を用いたキャッシュ装置を長寿命化させる場合、キャッシュする効果が小さいデータをキャッシュ装置に格納しないように制御する方法も考えられるが、そのように制御すると、本発明の実施例に係るキャッシュ装置（キャッシュSSD）のように、ライトデータの圧縮や暗号化などを行うデータ処理機能を有している場合、キャッシュする効果が小さいデータに対して、キャッシュ装置のデータ処理機能を適用できない。本発明の実施例に係るストレージ装置及びキャッシュ装置では、キャッシュする効果が小さいデータを極力FMに格納しないように制御できるため、キャッシュする効果が小さいデータに対してもキャ

ッシュ装置のデータ処理機能を適用でき、かつキャッシュ装置の長寿命化が可能である。

[0137] なお、実施例6で説明したとおり、本発明はストレージコントローラがデータの属性を指定し、その属性に基づいて、キャッシュ装置が内部ディスクの優先度を変更する例に限定されるものではない。実施例6に係るキャッシュSSDのように、キャッシュ装置がデータの種別によって自律的にディスク優先度を変更する態様であっても、FMへのライト頻度を低下させ、キャッシュ装置の寿命（耐用年数）を延ばすという効果を得ることが可能である。

符号の説明

[0138] 10 : ストレージ装置
20 : ネットワーク
40 : ホスト計算機
100 : ストレージコントローラ
110 : I/F CTL
120 : CPU
130 : メモリ
140 : ASIC
150 : I/F CTL
160 : 専用内部バス
170 : 内部バス
180 : メモリバス
190 : 最終記憶媒体
200 : キャッシュSSD
2001 : I/F CTL
2003 : CPU
2004 : 内部スイッチ
2005 : 内部メモリ

2006 : データ処理アシストユニット (A s s i s t)

2007 : FMC

2011 : FM

請求の範囲

[請求項1]

1 以上の最終記憶装置と、ホスト計算機及び前記最終記憶装置とが接続されるストレージコントローラを有するストレージ装置であって、

、

前記ストレージコントローラは、少なくともプロセッサと、前記ホスト計算機からのライトデータを一時的に格納するキャッシュ装置を有し、

前記キャッシュ装置は、書き替え回数に制限のある不揮発記憶媒体と、前記不揮発記憶媒体に比べて高速なアクセスが可能な揮発性記憶媒体から成る内部キャッシュを有し、

前記ストレージコントローラは、前記キャッシュ装置に前記ホスト計算機からのライトデータを格納する際、前記ライトデータの格納を指示するためのライト要求であって、前記ライトデータに関する情報が含まれている前記ライト要求を前記キャッシュ装置に発行し、

、

前記キャッシュ装置は、前記ライト要求を受信すると、前記ライト要求に伴うライトデータを前記内部キャッシュに格納するとともに、前記ライトデータに関する情報から前記ライトデータの内部ディステージ優先度を決定し、前記内部ディステージ優先度を前記ライトデータに対応付けて記憶し、

前記キャッシュ装置は前記内部キャッシュに格納された前記ライトデータのうち、前記ライトデータに対応付けられた前記内部ディステージ優先度に基づき、前記不揮発記憶媒体へと移動する、ことを特徴とする、ストレージ装置。

[請求項2]

前記ストレージコントローラは、前記ライトデータに関する情報として、ライトデータがシーケンシャルアクセスデータであるかランダムアクセスデータであるかを示す情報を、前記ライト要求に含ませて、前記キャッシュ装置に発行し、

前記キャッシュ装置は、前記ライト要求を受信すると、前記ライト要求に伴うライトデータを前記内部キャッシュに格納するとともに、前記ライトデータの特性に関する情報がシーケンシャルアクセスデータであることを示す情報であった場合、前記ライト要求に伴うライトデータの前記内部ディステージ優先度は低いと決定し、前記ライトデータの特性に関する情報がランダムアクセスデータであることを示す情報であった場合、前記ライト要求に伴うライトデータの前記内部ディステージ優先度は高いと決定することを特徴とする、請求項1に記載のストレージ装置。

[請求項3]

前記キャッシュ装置は、前記ライトデータに対して最後にアクセスがあった時刻を前記ライトデータに対応付けて記憶しており、

前記内部キャッシュの未使用領域の量が所定の閾値未満になった時、前記内部キャッシュに格納された前記ライトデータの中に、前記ライトデータに対応付けられた前記内部ディステージ優先度の高いライトデータが存在しなかった場合、前記ライトデータに対応付けられた前記内部ディステージ優先度が低い前記ライトデータの中から、前記ライトデータに最後にアクセスがあった時刻が最も古いデータを前記不揮発記憶媒体へと移動することを特徴とする、請求項2に記載のストレージ装置。

[請求項4]

前記キャッシュ装置は、前記ライトデータに対して最後にアクセスがあった時刻を前記ライトデータに対応付けて記憶しており、

前記キャッシュ装置は、前記内部キャッシュの未使用領域の量が所定の閾値未満になった時、前記内部キャッシュに格納された前記ライトデータの中に、前記ライトデータに対応付けられた前記内部ディステージ優先度の高いライトデータが存在しなかった場合、前記ストレージコントローラに、前記ライトデータの前記最終記憶装置へのディステージが必要である旨を通知し、

前記ストレージコントローラは、前記通知に応じて、前記キャッシ

装置内の内部キャッシュに格納された前記ライトデータの中から、前記ライトデータに最後にアクセスがあった時刻が最も古いデータを読み出して、前記最終記憶装置へ書き込むことを特徴とする、請求項2に記載のストレージ装置。

[請求項5] 前記キャッシュ装置は、前記前記ストレージコントローラからライトデータを受信すると、前記ライトデータに対して変換処理を施して、前記キャッシュ装置に格納し、

前記ストレージコントローラが前記キャッシュ装置に格納された前記変換処理が施されたデータを前記最終記憶装置へディステージする際には、前記変換処理が施されたデータを読み出して、前記最終記憶装置に格納することを特徴とする、請求項1に記載のストレージ装置。

[請求項6] 前記ストレージコントローラは、前記ホスト計算機からデータリード要求を受け付けると、

前記最終記憶装置から前記変換処理が施されたデータを読み出して、前記キャッシュ装置に格納し、

前記キャッシュ装置に対し、前記キャッシュ装置に格納された前記変換処理が施されたデータを読み出すリード要求であって、データ逆変換を行う旨が指定された前記リード要求を発行し、

前記キャッシュ装置は、前記リード要求を受信すると、前記変換処理が施されたデータを逆変換して前記ストレージコントローラに送信することを特徴とする、

請求項5に記載のストレージ装置。

[請求項7] 前記変換処理は、データの圧縮処理であって、前記変換処理が施されたデータを逆変換する処理は、圧縮データの伸長処理であることを特徴とする、

請求項6に記載のストレージ装置。

[請求項8] 前記変換処理は、データの暗号化処理であって、前記変換処理が施

されたデータを逆変換する処理は、暗号化データの復号化処理であることを特徴とする、

請求項6に記載のストレージ装置。

[請求項9]

前記ストレージ装置はさらに、前記ホスト計算機からのライトデータを一時的に格納するためのメモリを有し、

前記ストレージコントローラは、前記ホスト計算機からホストライト要求を受信すると、前記ホストライト要求で対象とされるライト対象データを格納するための空き領域が前記メモリに存在するか判定し、

前記空き領域が前記メモリに存在しない場合、前記メモリ内に格納されているデータのうち、最後にアクセスがあった時刻が最も古いデータを前記キャッシュ装置へ移動した後で空き領域を確保し、

前記確保された空き領域に前記ライト対象データを格納する、ことを特徴とする、請求項1に記載のストレージ装置。

[請求項10]

前記ストレージコントローラは、前記確保された空き領域に前記ライト対象データを格納した後、前記ライト対象データの更新前データが前記キャッシュ装置に格納されているか判定し、

前記ライト対象データの更新前データが前記キャッシュ装置に格納されている場合、前記キャッシュ装置に前記更新前データの無効化を要求することを特徴とする、

請求項9に記載のストレージ装置。

[請求項11]

1以上の最終記憶装置と、ホスト計算機及び前記最終記憶装置とが接続されるストレージコントローラを有するストレージ装置の制御方法であって、

前記ストレージコントローラは、少なくともプロセッサと、前記ホスト計算機からのライトデータを一時的に格納するキャッシュ装置を有し、

前記キャッシュ装置は、書き替え回数に制限のある不揮発記憶媒体

と、前記不揮発記憶媒体に比べて高速なアクセスが可能な揮発性記憶媒体から成る内部キャッシュを有するキャッシュ装置であって、

前記方法は、

前記ストレージコントローラが前記キャッシュ装置に前記ホスト計算機からのライトデータを格納する際、前記ライトデータの格納を指示するためのライト要求であって、前記ライトデータに関する情報が含まれている前記ライト要求を前記キャッシュ装置に発行し、

前記キャッシュ装置は前記ライト要求を受信すると、

前記内部キャッシュに、前記ライト要求に伴うライトデータを格納するための未使用の領域が存在するか判定し、

前記未使用の領域が存在しない場合、前記内部キャッシュに格納されたデータのうち、前記データに対応付けられた内部ディステージ優先度の高いデータを、前記不揮発記憶媒体へと移動した後、前記ライト要求に伴うライトデータを格納するための前記未使用の領域を確保して、

前記ライト要求に伴うライトデータを前記確保された前記未使用の領域に格納するとともに、前記ライトデータに関する情報から前記ライトデータの内部ディステージ優先度を決定し、前記内部ディステージ優先度を前記ライトデータに対応付けて記憶する、処理を実行することを特徴とする、ストレージ装置の制御方法。

[請求項12]

前記ストレージコントローラは、前記ライトデータに関する情報として、ライトデータがシーケンシャルアクセスデータであるかランダムアクセスデータであることを示す情報を、前記ライト要求に含ませて、前記キャッシュ装置に発行し、

前記キャッシュ装置は、前記ライト要求に伴うライトデータを前記確保された前記未使用の領域に格納する際、前記ライトデータに関する情報がシーケンシャルアクセスデータであることを示す情報であった場合、前記ライト要求に伴うライトデータの前記内部ディス

ページ優先度は低いと決定し、前記ライトデータの特性に関する情報がランダムアクセスデータであることを示す情報であった場合、前記ライト要求に伴うライトデータの前記内部ディスティージ優先度は高いと決定することを特徴とする、

請求項 1 1 に記載のストレージ装置の制御方法。

[請求項13]

前記キャッシュ装置は、前記ライトデータに対して最後にアクセスがあった時刻を前記ライトデータに対応付けて記憶しており、

前記内部キャッシュに、前記ライト要求に伴うライトデータを格納するための未使用の領域が存在せず、かつ前記内部キャッシュに格納された前記ライトデータの中に、前記ライトデータに対応付けられた前記内部ディスティージ優先度の高いライトデータが存在しなかった場合、

前記キャッシュ装置は、前記ライトデータに対応付けられた前記内部ディスティージ優先度が低い前記ライトデータの中から、前記ライトデータに最後にアクセスがあった時刻が最も古いデータを前記不揮発記憶媒体へと移動する処理を行うことを特徴とする、

請求項 1 2 に記載のストレージ装置の制御方法。

[請求項14]

前記キャッシュ装置は、前記ライトデータに対して最後にアクセスがあった時刻を前記ライトデータに対応付けて記憶しており、

前記キャッシュ装置は、前記内部キャッシュに、前記ライト要求に伴うライトデータを格納するための未使用の領域が存在せず、かつ前記内部キャッシュに格納された前記ライトデータの中に、前記ライトデータに対応付けられた前記内部ディスティージ優先度の高いライトデータが存在しなかった場合、

前記キャッシュ装置は、前記ストレージコントローラに、前記ライトデータの前記最終記憶装置へのディスティージが必要である旨を通知し、

前記ストレージコントローラは、前記通知に応じて、前記キャッシ

ユ装置内の内部キャッシュに格納された前記ライトデータの中から、前記ライトデータに最後にアクセスがあった時刻が最も古いデータを読み出して、前記最終記憶装置へ書き込む処理を行うことを特徴とする、
請求項 1 2 に記載のストレージ装置の制御方法。

[図1]

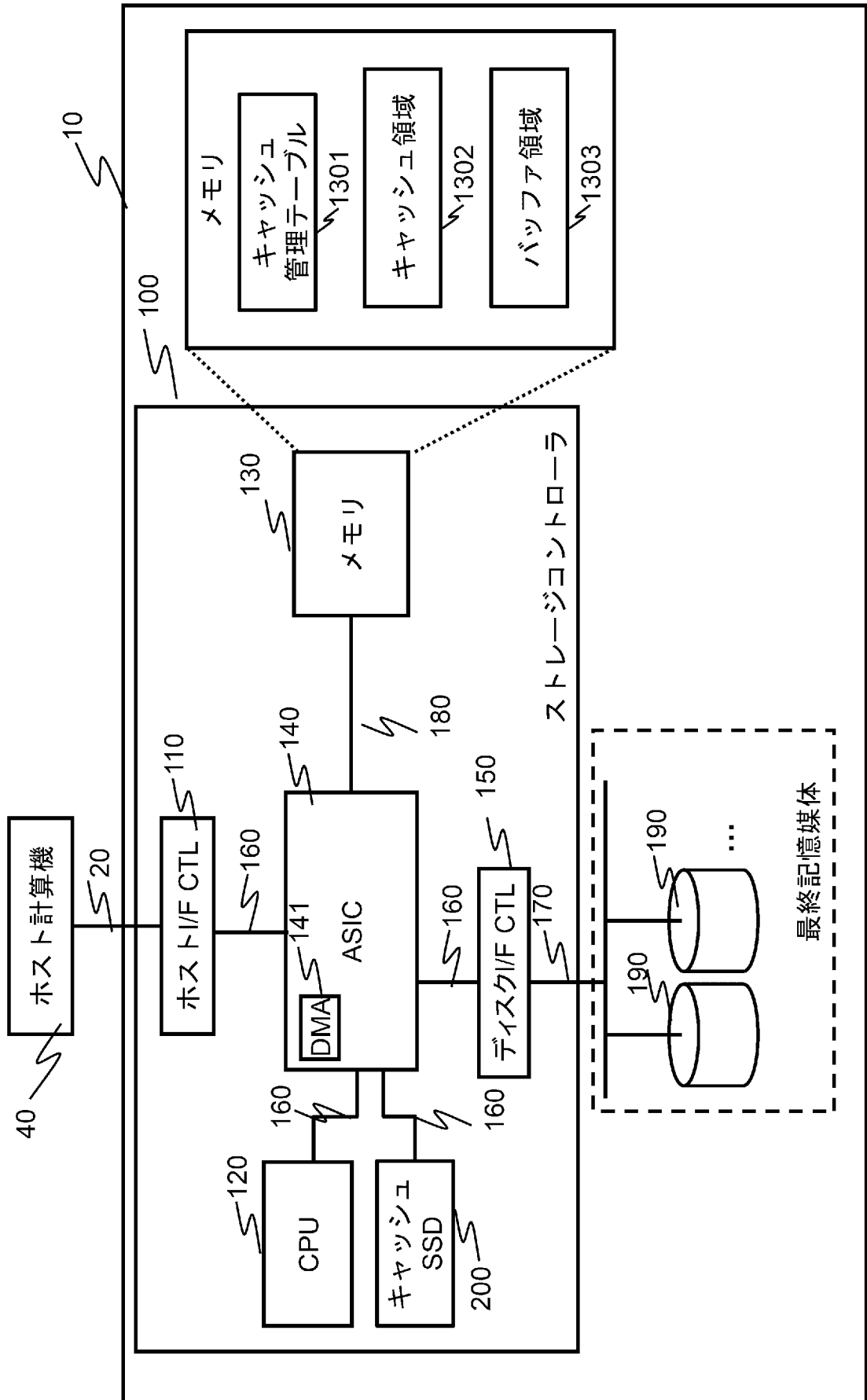


図1

[図2]

Cache segment ID	LBA	ビットマップ	キャッシュ状態	直近アクセス時刻
0x0	B45F24h	0xffff ffff ffff ffff	dirty	FEDCB
0x1	A4542Bh	0xffff ffff ffff ffff	clean	A987654
0x2	987810h	0x0fff ffff ffff ffff	clean	54321FE
...
0x1FFFFFF	ABCDEh	0xffff ffff ffff ffff	clean	123
0x200000	1234h	0xffff ffff ffff ffff	dirty	12345
...
0x9FFFFFF	9876h	0xffff ffff ffff ffff	clean	6789

13011 13012 13013 13014 13015

Cache DRAM 1310

Cache SSD 1311

図2

[図3]

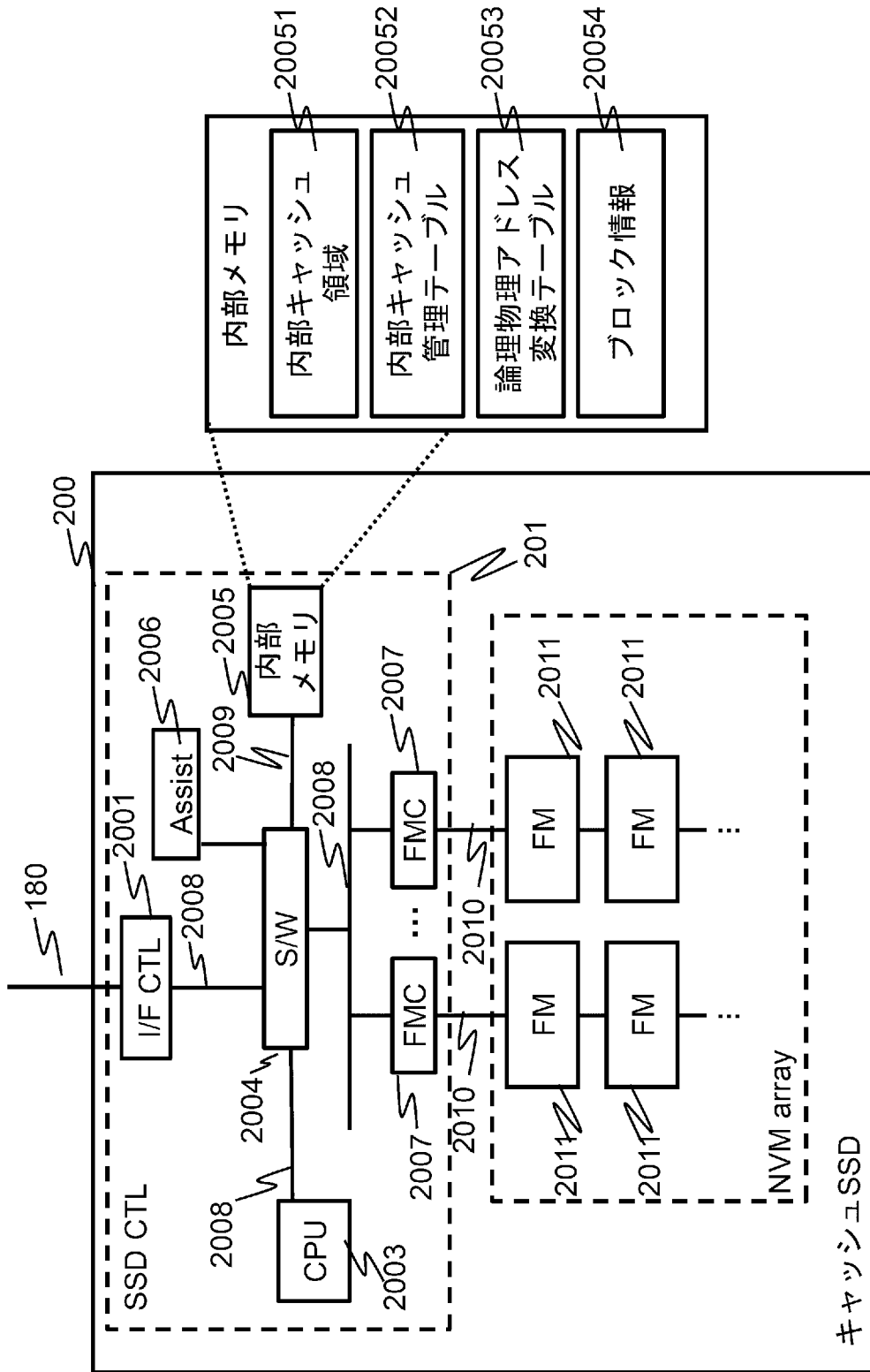


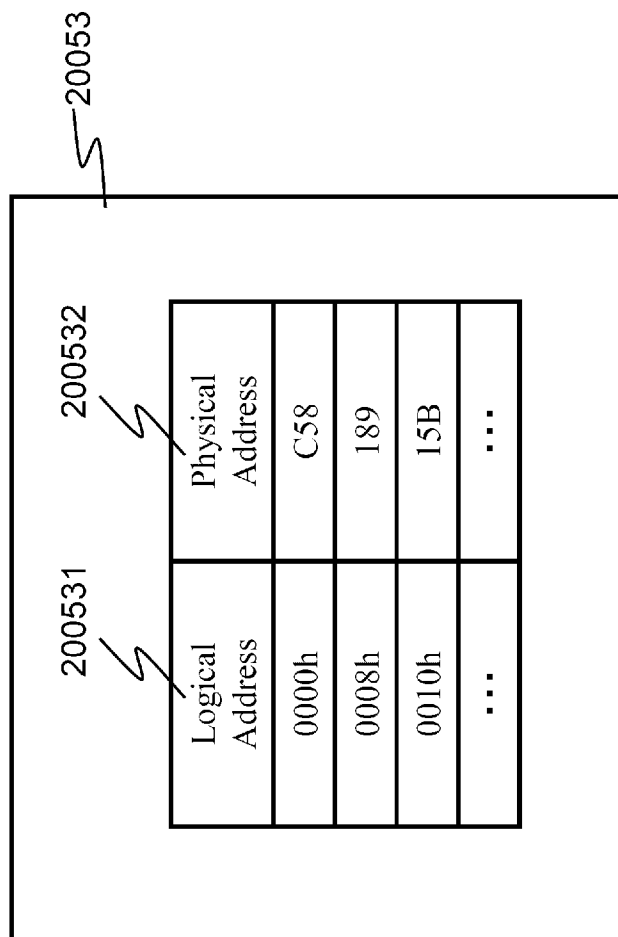
図3

[図4]

	200521	200522	200523	200524	200525	200526	
	⚡	⚡	⚡	⚡	⚡	⚡	20052
							20052-0
Internal Cache Segment ID	Logical Address (Cache)	データサイズ	内部キャッシュ状態	直近アクセス時刻	内部ディスク優先度		
0x0	12345h	4 KB	dirty	FEDCB	Low		
0x1	10000h	1 KB	clean	A987654	High		
0x2	2345h	4 KB	dirty	54321FE	Low		
...		

図4

[図5]



[図5]

[図6]

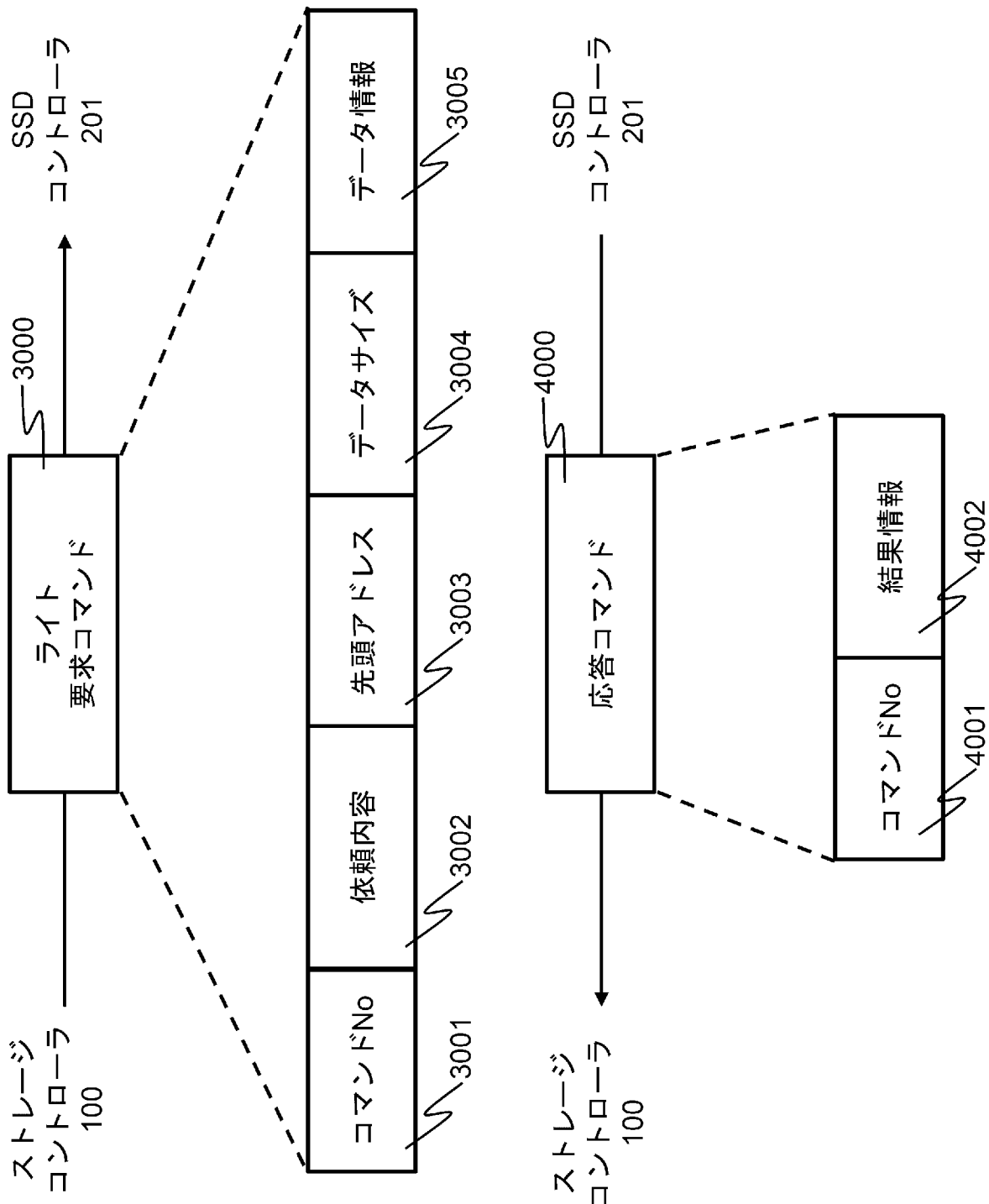


図6

[図7]

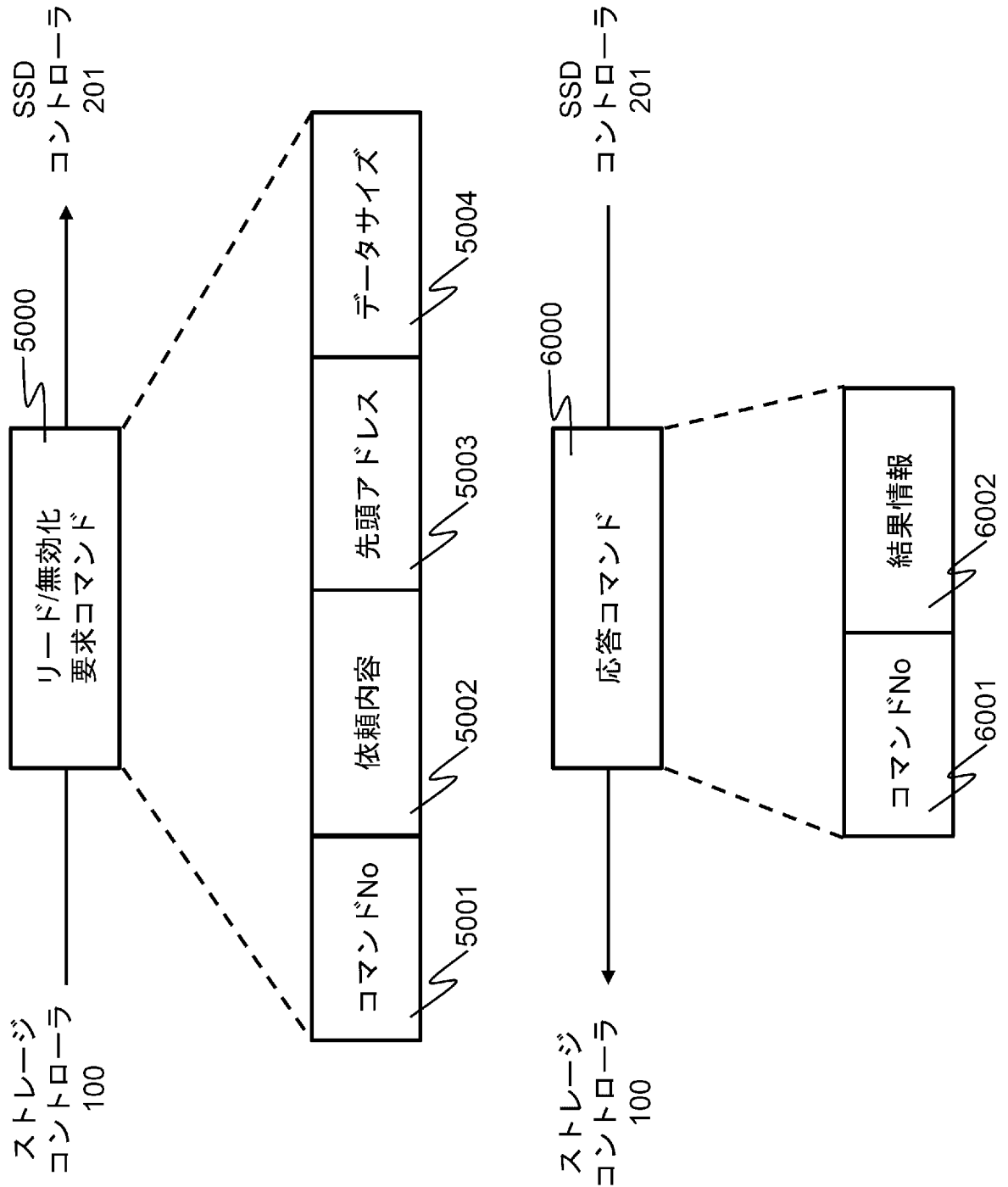


図7

[図8]

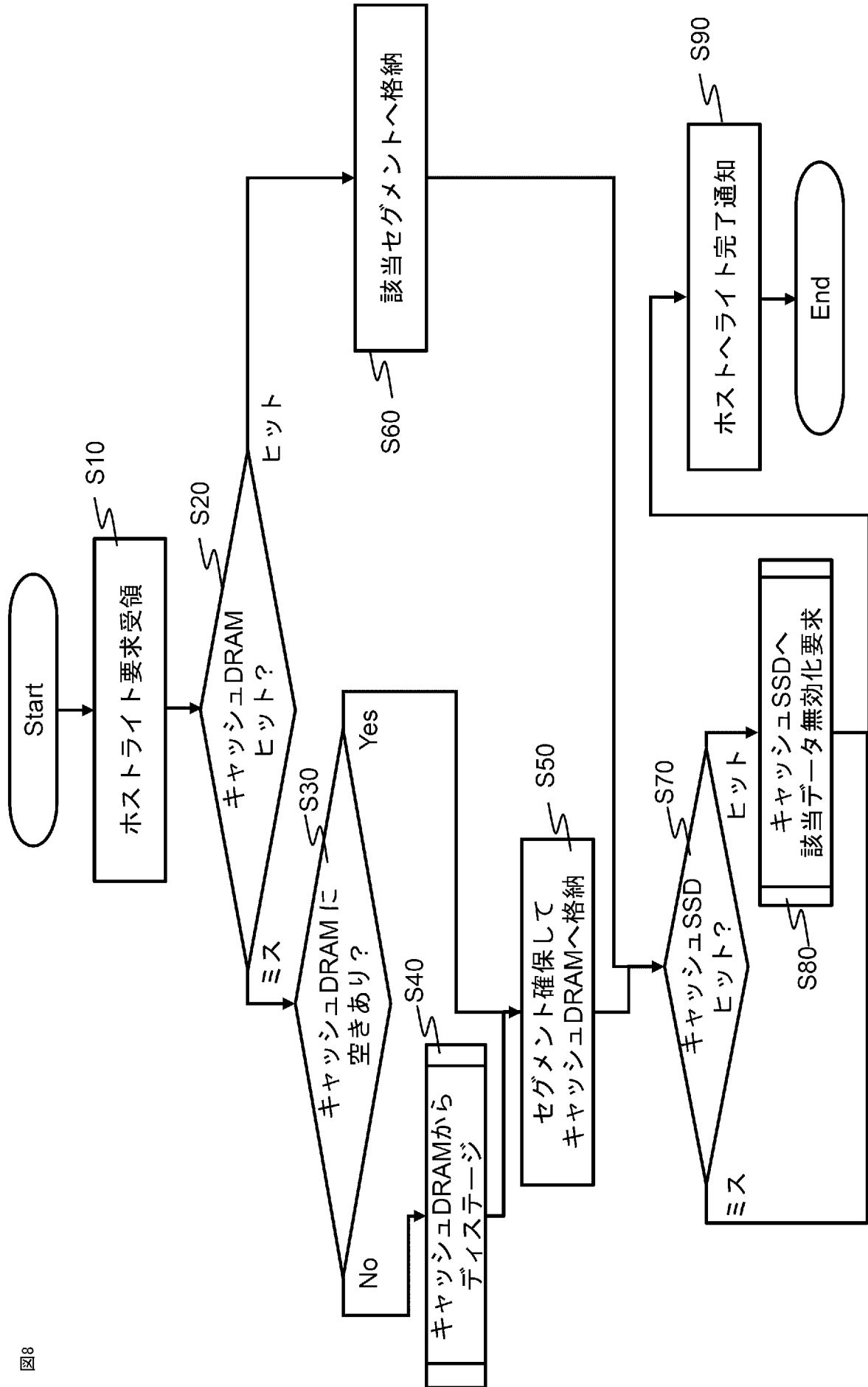


図8

図9

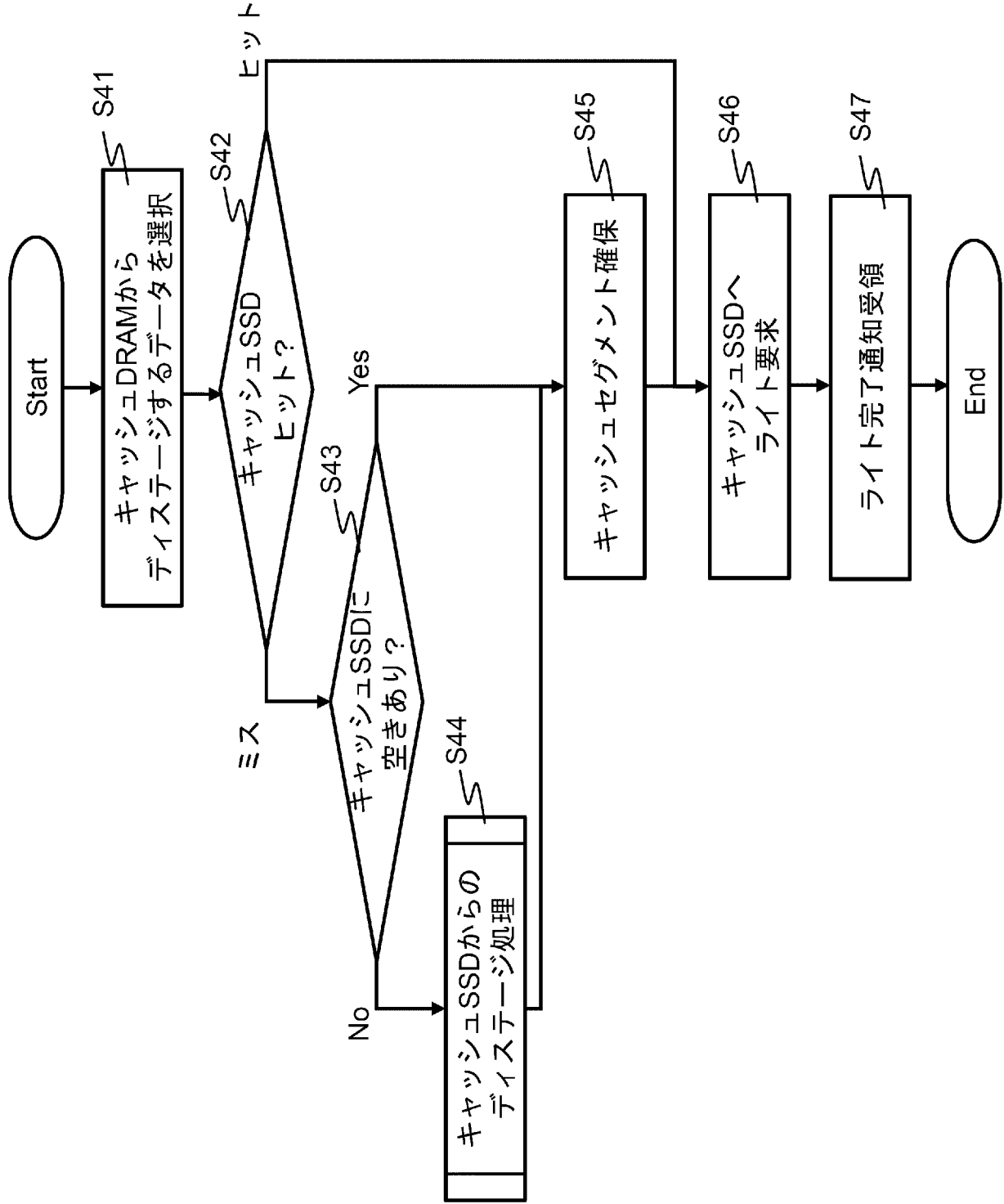


図9

[図10]

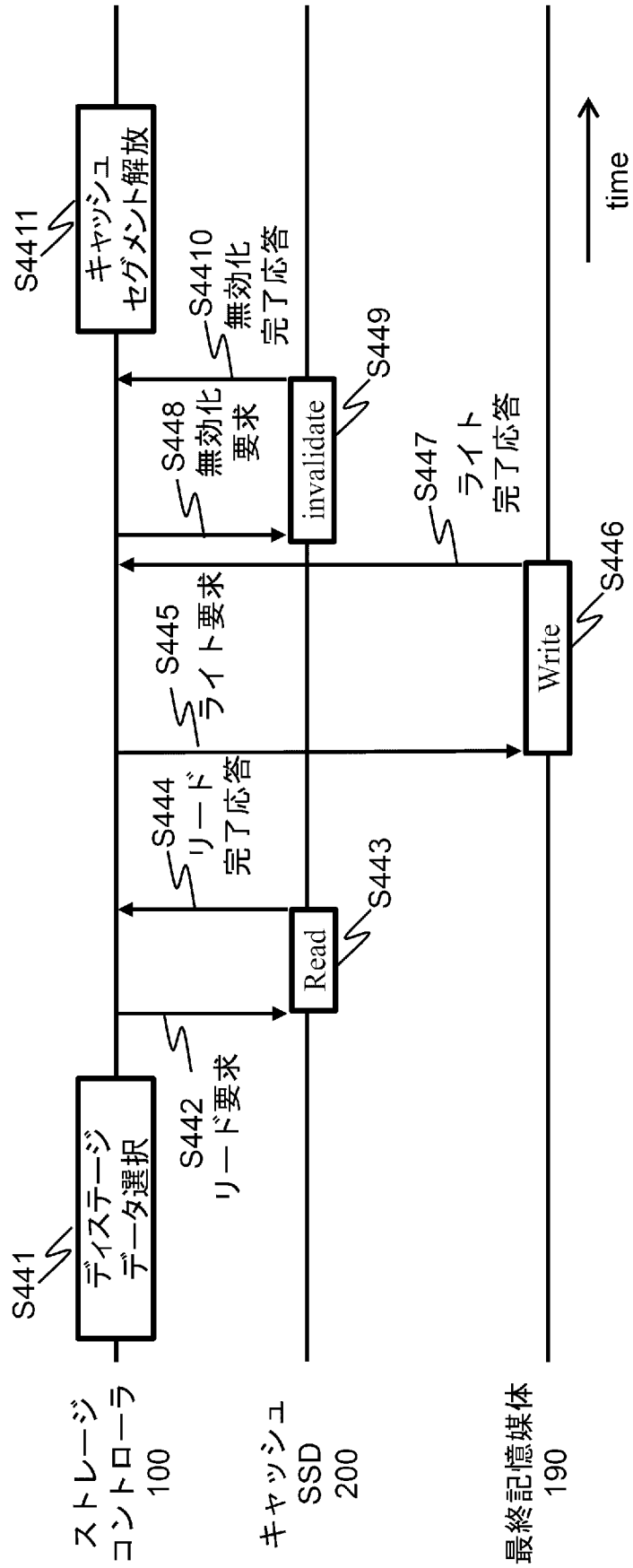


図10

図11

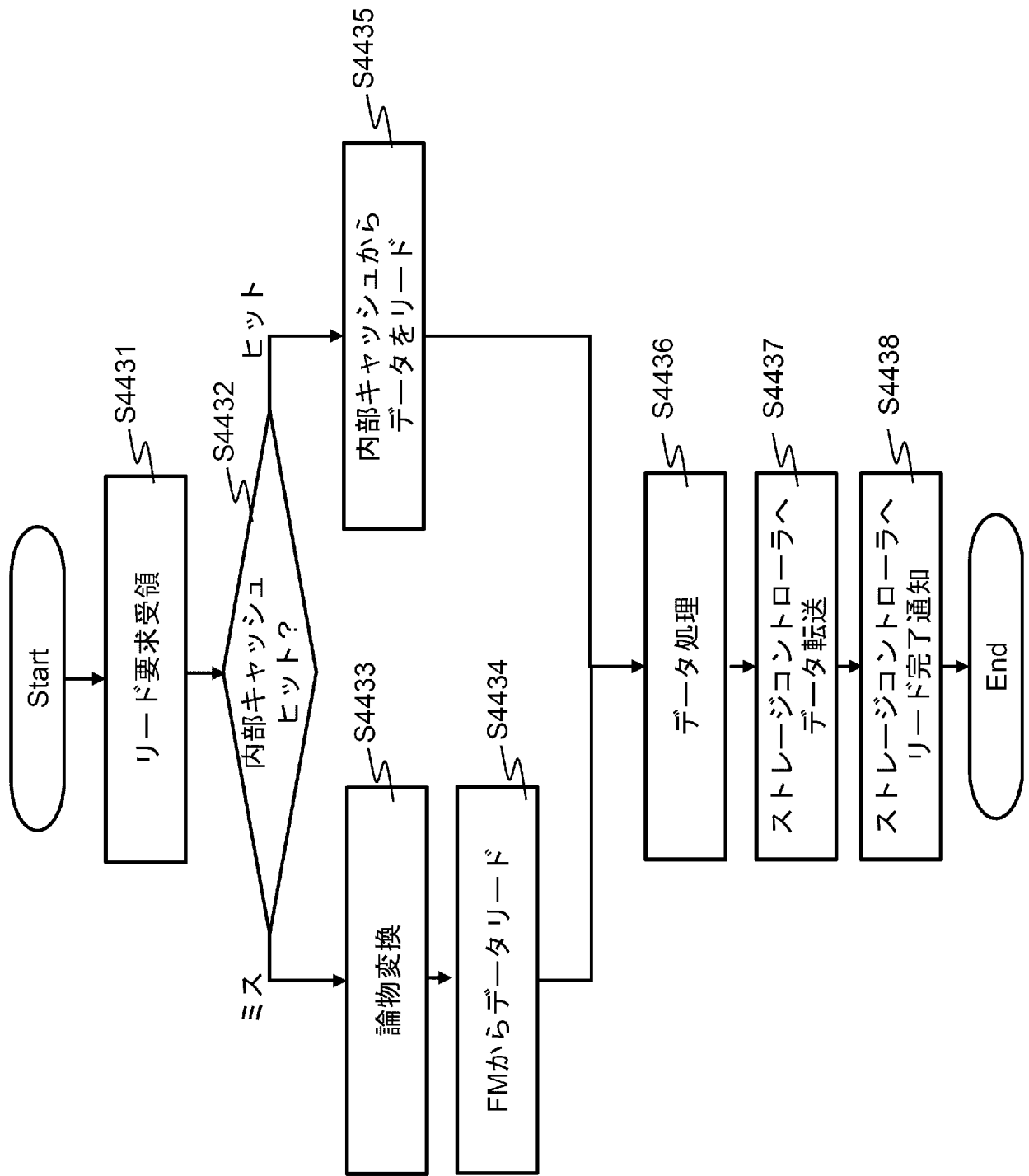


図11

[図12]

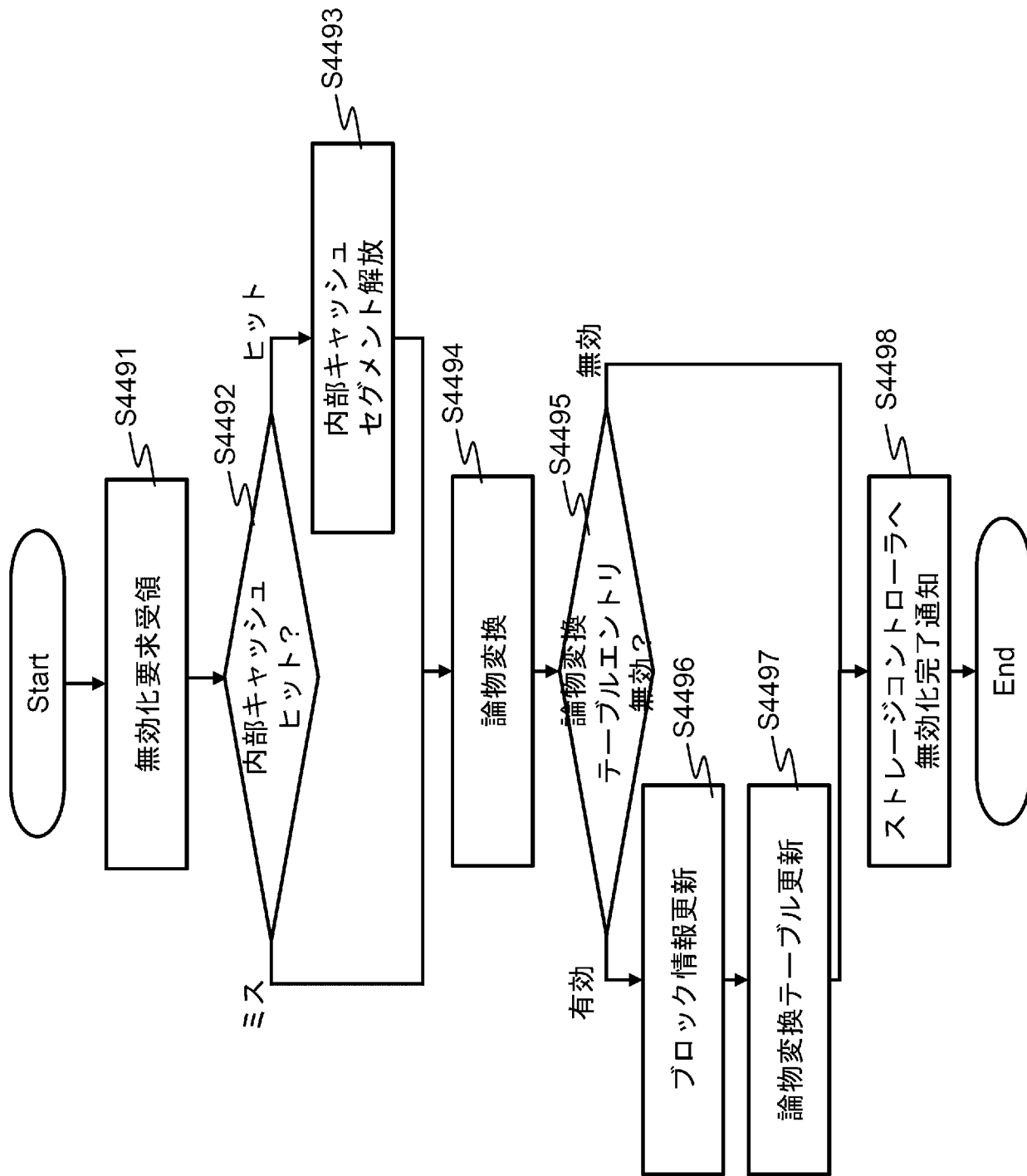


図12

[図13]

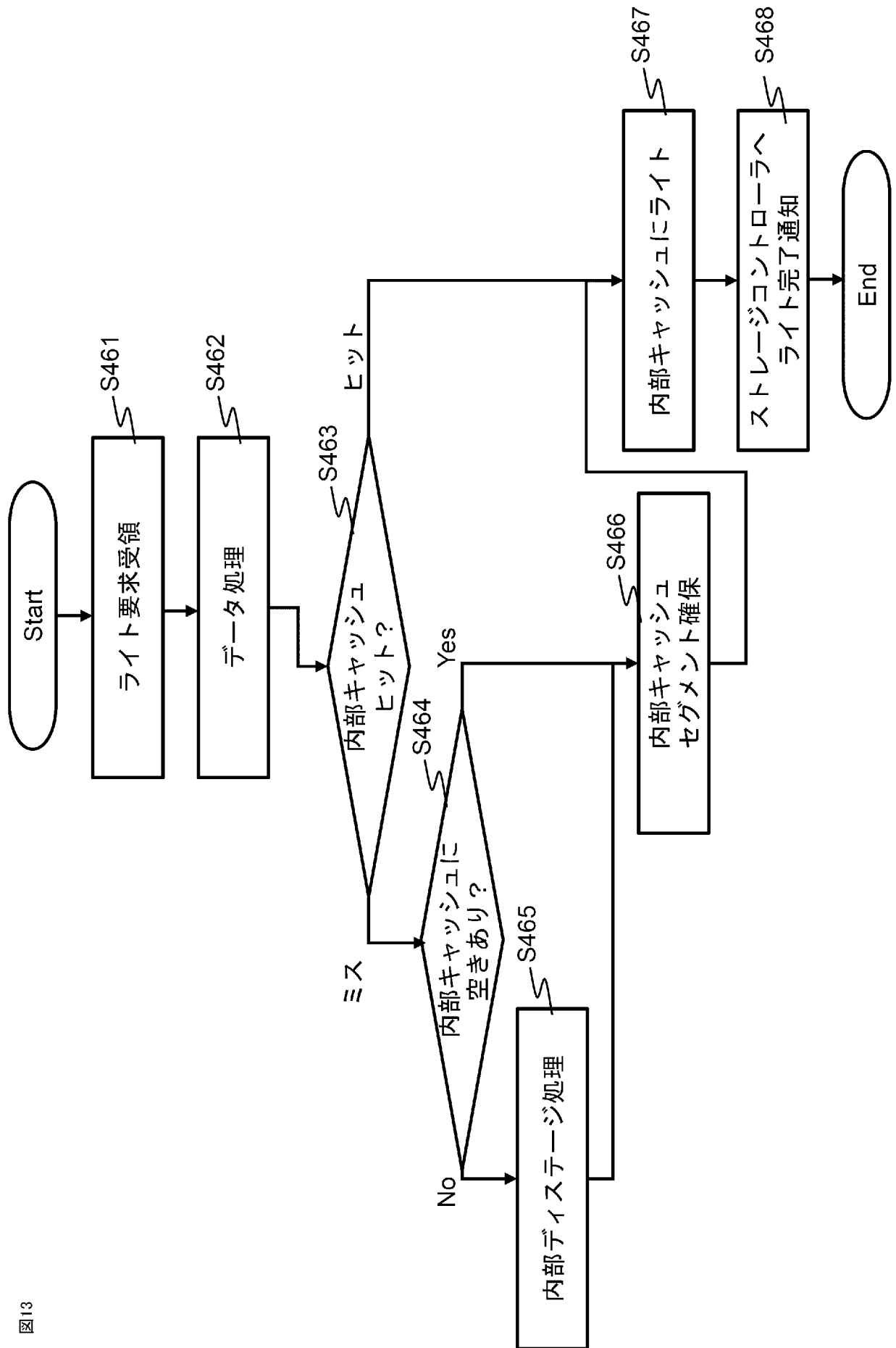
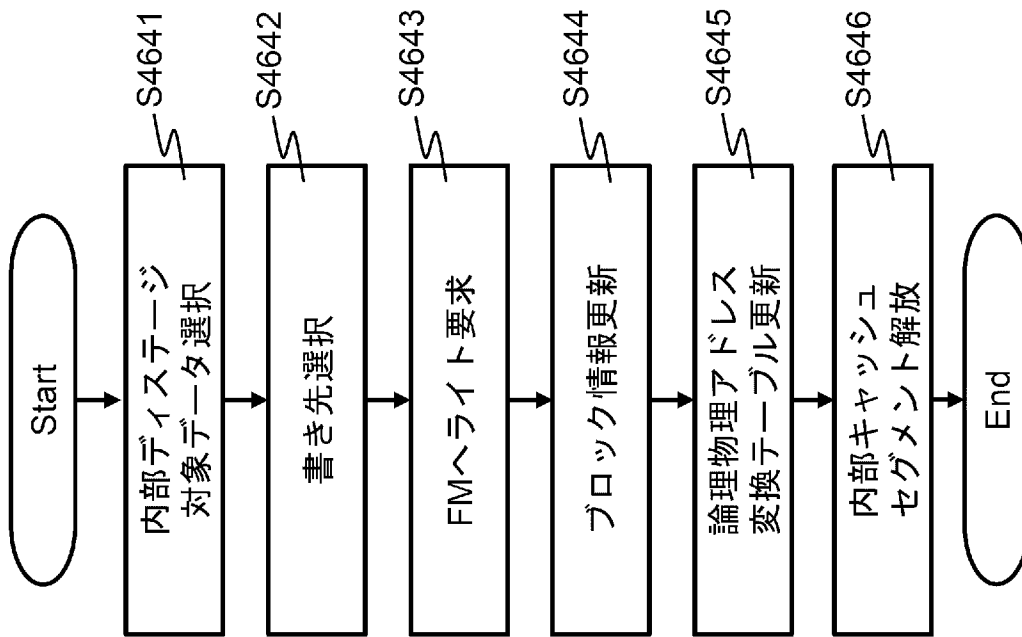


図13

[図14]



[図15]

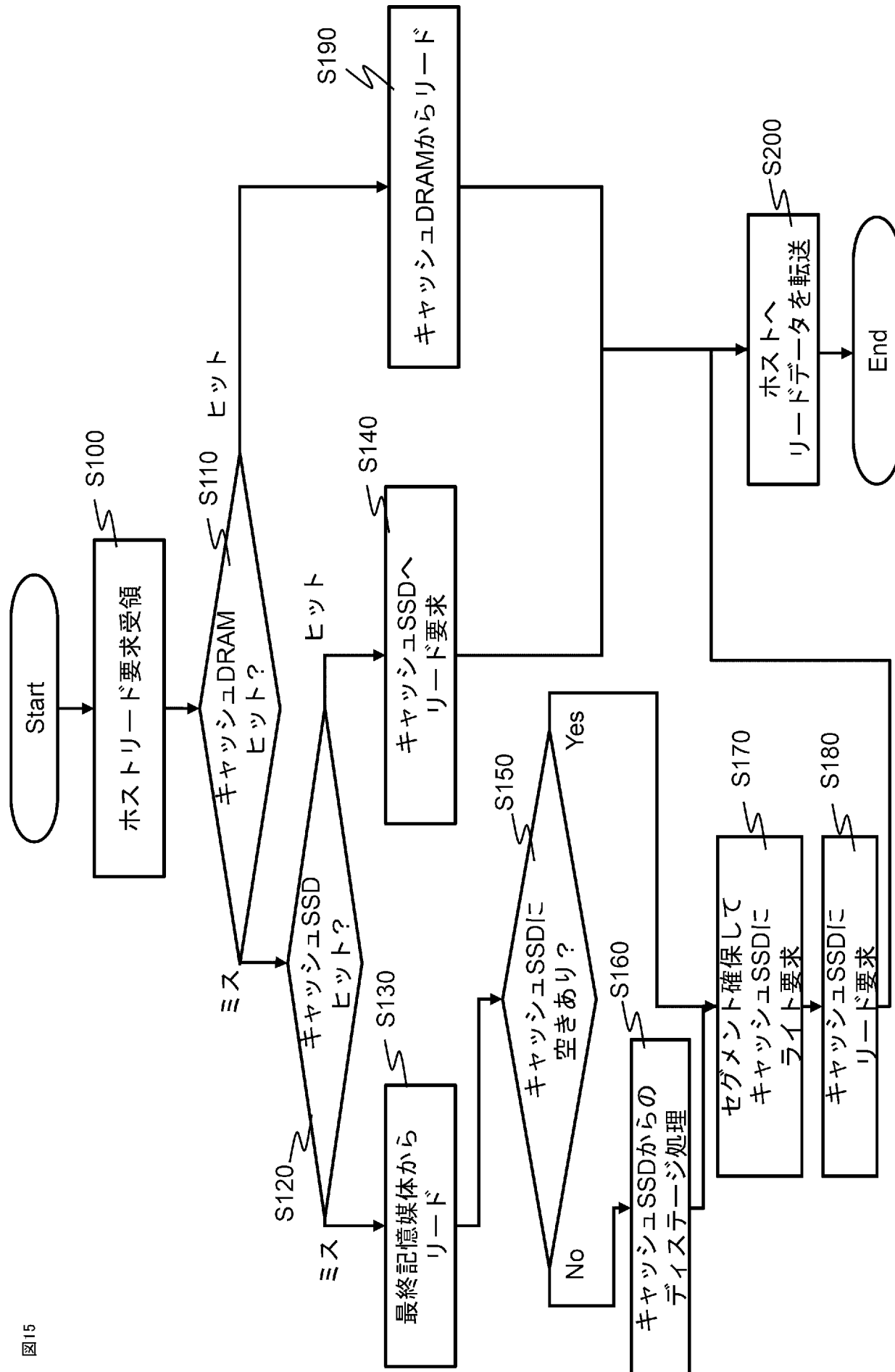


図15

[図16]

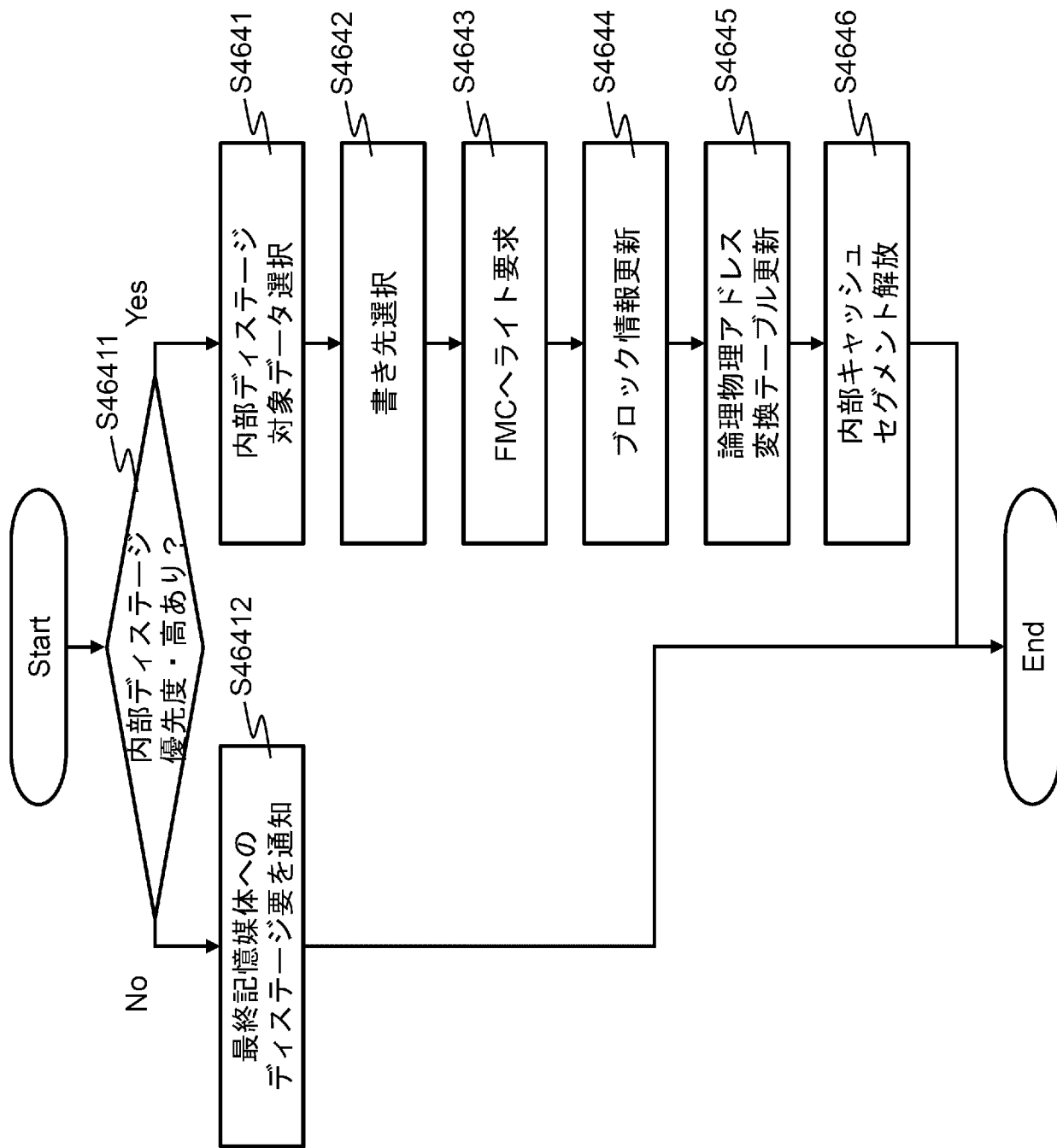


図16

[図17]

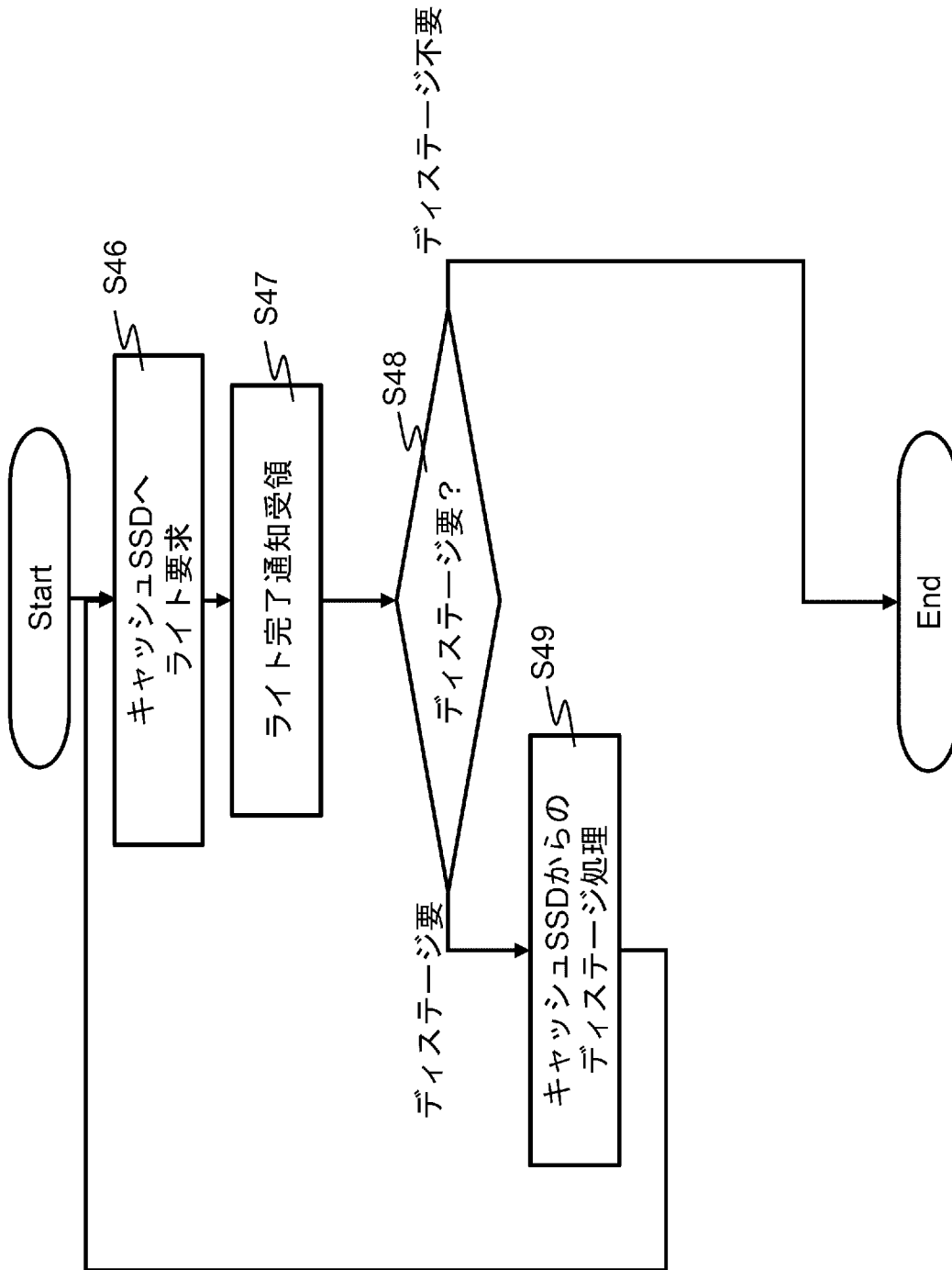


図17

[図18]

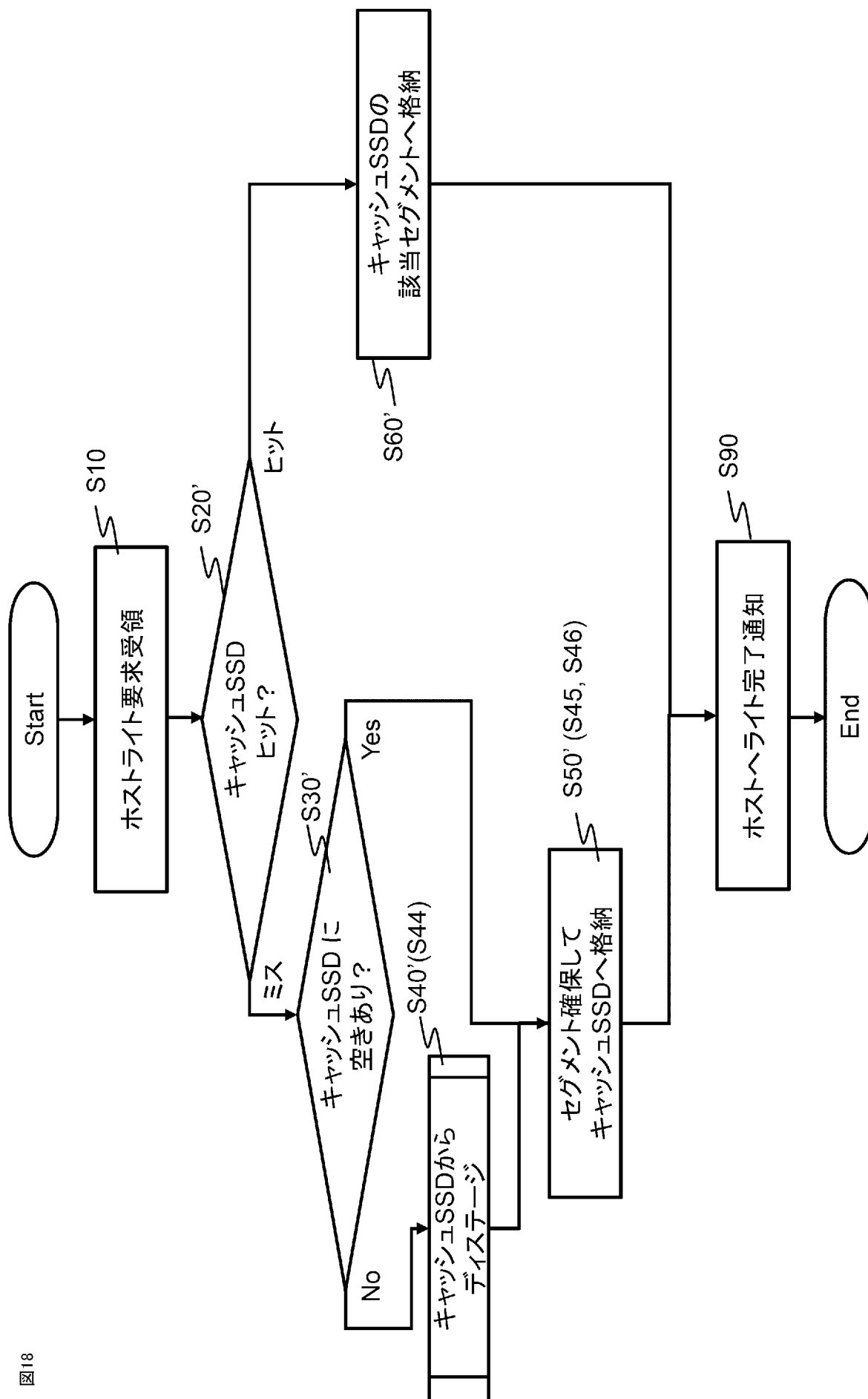


図18

[図19]

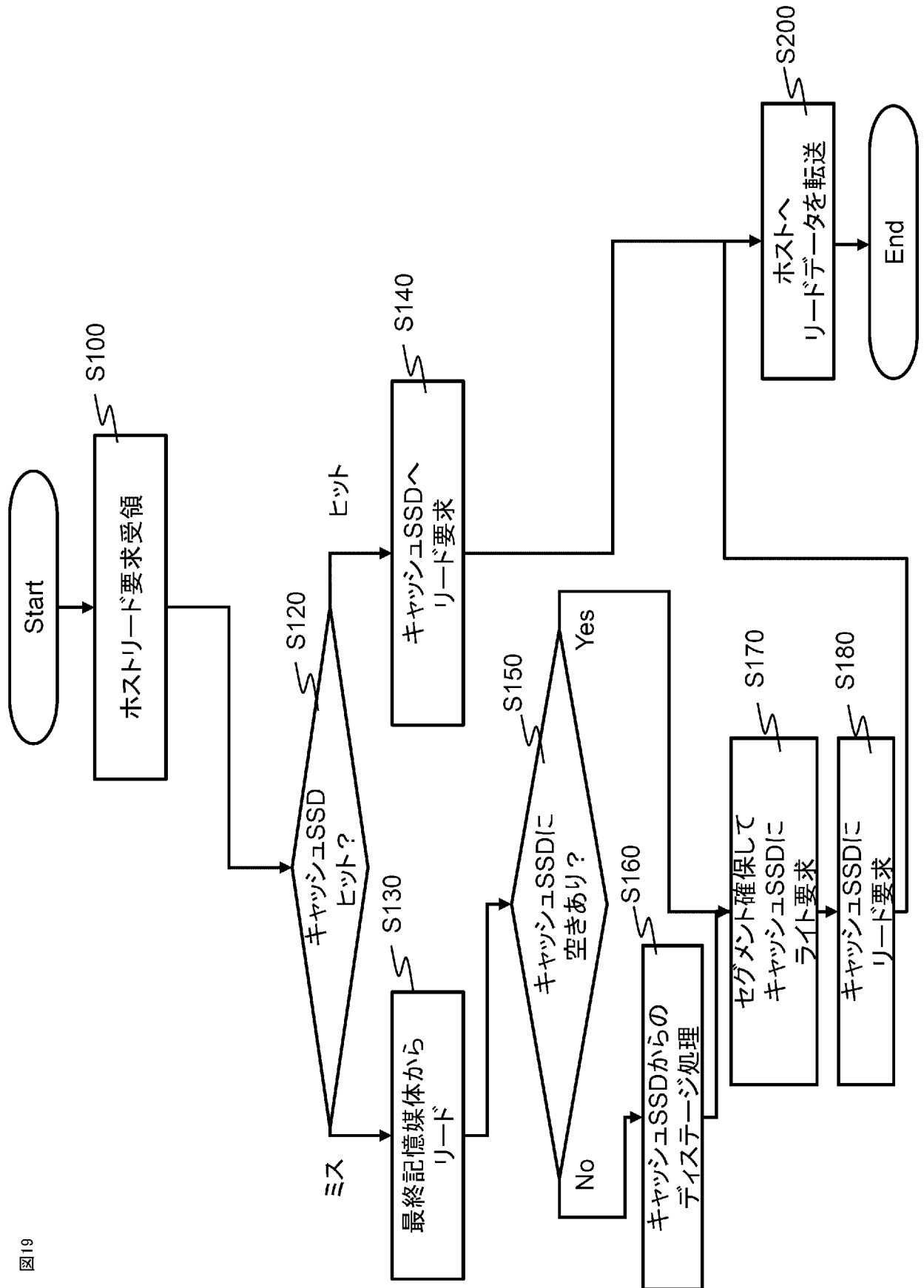


図19

[図20]

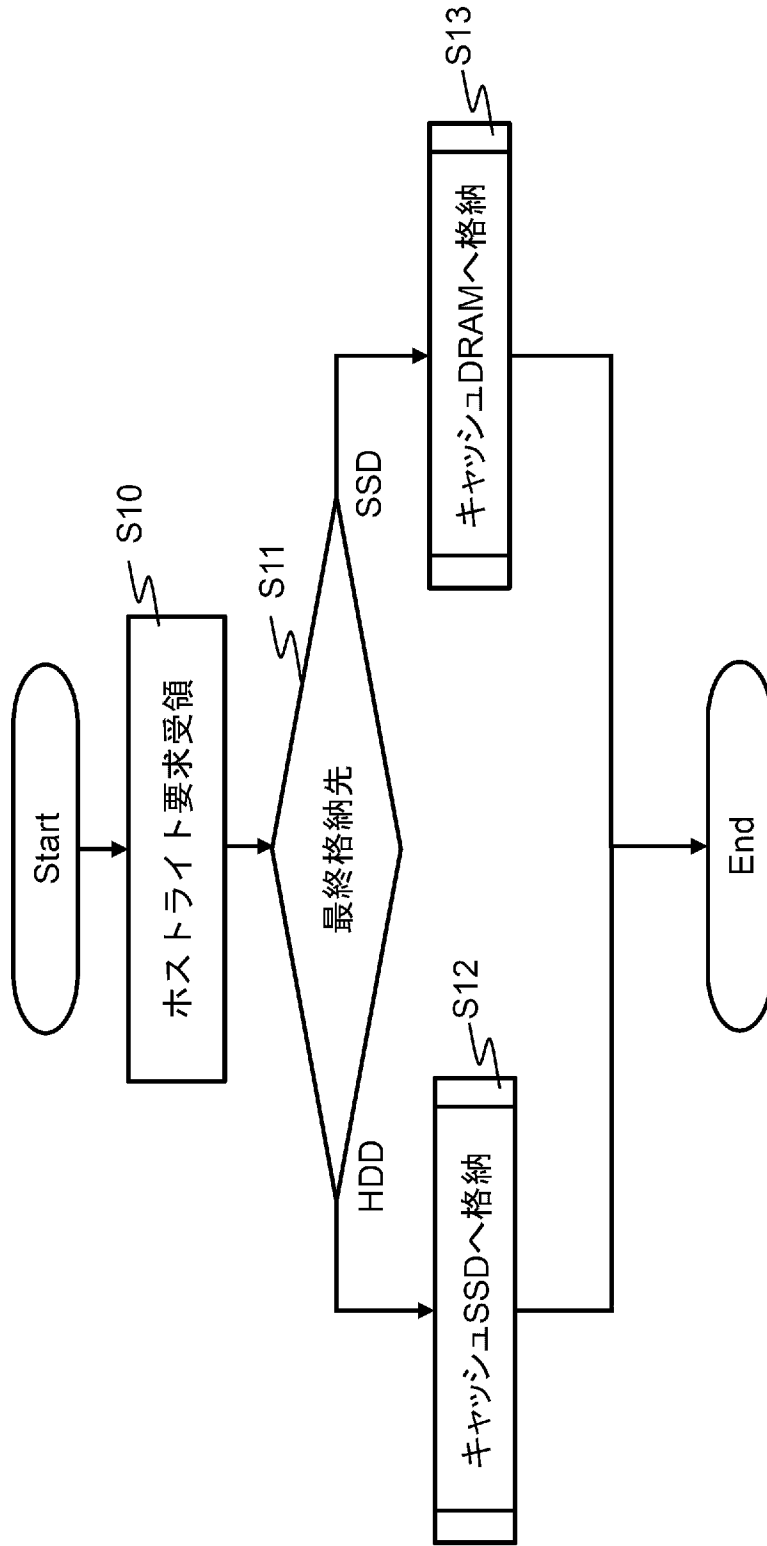


図20

[図21]

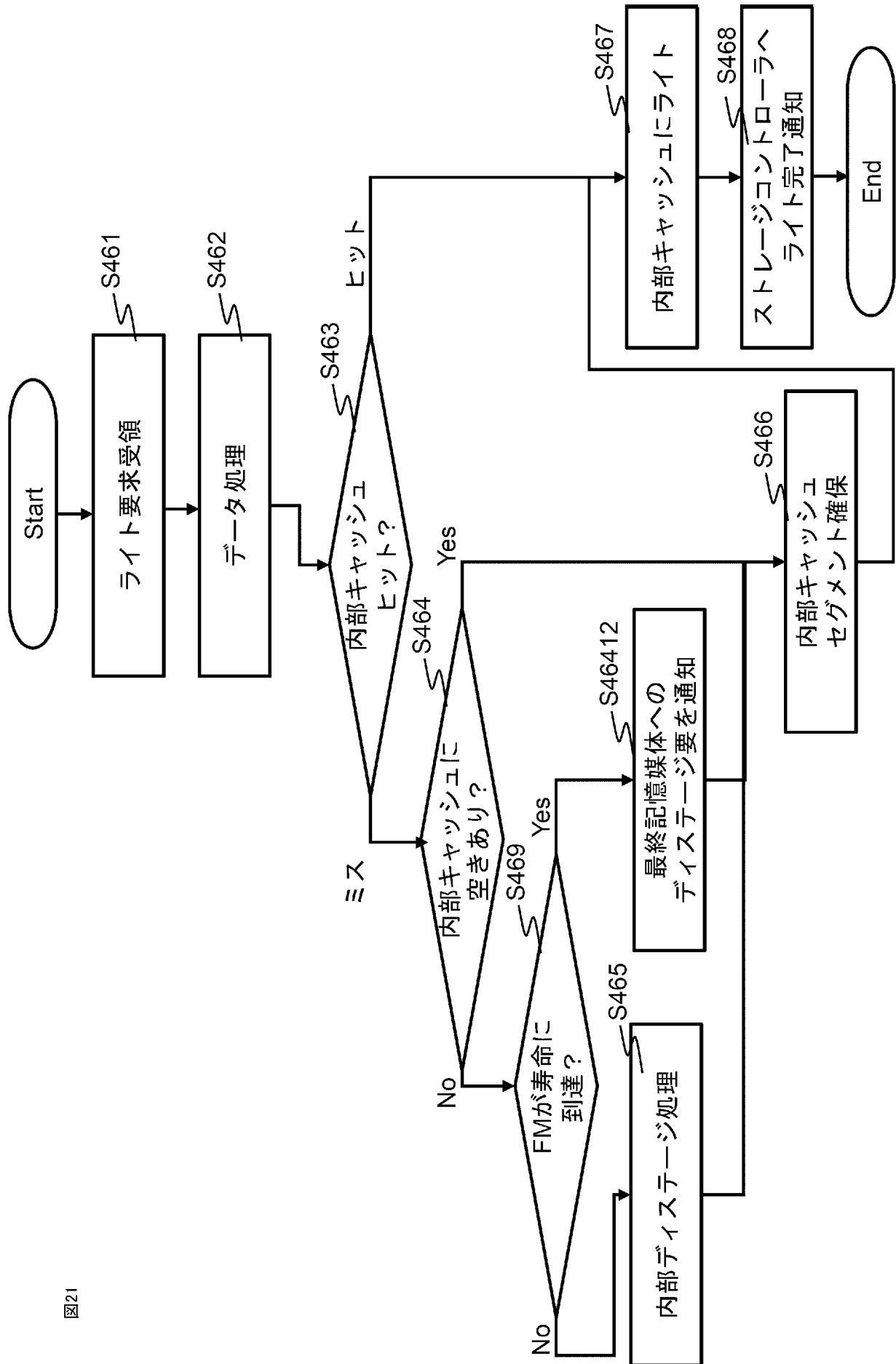


図21

INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2014/062714

A. CLASSIFICATION OF SUBJECT MATTER G06F12/08(2006.01)i, G06F12/12(2006.01)i		
According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED		
Minimum documentation searched (classification system followed by classification symbols) G06F12/08, G06F12/12		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Jitsuyo Shinan Koho 1922-1996 Jitsuyo Shinan Toroku Koho 1996-2014 Kokai Jitsuyo Shinan Koho 1971-2014 Toroku Jitsuyo Shinan Koho 1994-2014		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X Y A	JP 2008-217527 A (Hitachi, Ltd.), 18 September 2008 (18.09.2008), paragraphs [0010] to [0140]; fig. 1 to 2, 7 to 9, 12, 17 to 19 & US 2008/0222359 A1 & US 2011/0271066 A1 & US 2012/0233398 A1 & EP 1967953 A2	1-2, 4, 9, 11-12, 14 5-8 3, 10, 13
Y A	JP 2008-15918 A (Toshiba Corp.), 24 January 2008 (24.01.2008), paragraphs [0101] to [0105]; fig. 18 (Family: none)	5-8 1-4, 9-14
A	JP 2011-204060 A (NEC Corp.), 13 October 2011 (13.10.2011), entire text; all drawings & US 2011/0238908 A1	1-14
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/> See patent family annex.		
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search 30 June, 2014 (30.06.14)		Date of mailing of the international search report 08 July, 2014 (08.07.14)
Name and mailing address of the ISA/ Japanese Patent Office		Authorized officer
Facsimile No.		Telephone No.

INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2014/062714

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	JP 2009-205335 A (Hitachi, Ltd.), 10 September 2009 (10.09.2009), entire text; all drawings & US 2009/0216945 A1 & US 2011/0296091 A1 & US 2012/0254523 A1	1-14

A. 発明の属する分野の分類 (国際特許分類 (IPC)) Int.Cl. G06F12/08(2006.01)i, G06F12/12(2006.01)i		
B. 調査を行った分野 調査を行った最小限資料 (国際特許分類 (IPC)) Int.Cl. G06F12/08, G06F12/12		
最小限資料以外の資料で調査を行った分野に含まれるもの 日本国実用新案公報 1922-1996年 日本国公開実用新案公報 1971-2014年 日本国実用新案登録公報 1996-2014年 日本国登録実用新案公報 1994-2014年		
国際調査で使用した電子データベース (データベースの名称、調査に使用した用語)		
C. 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
X Y A	JP 2008-217527 A (株式会社日立製作所) 2008.09.18, 段落【0010】 -【0140】, 第1-2, 7-9, 12, 17-19 図 & US 2008/0222359 A1 & US 2011/0271066 A1 & US 2012/0233398 A1 & EP 1967953 A2	1-2, 4, 9, 11-12, 14 5-8 3, 10, 13
Y A	JP 2008-15918 A (株式会社東芝) 2008.01.24, 段落【0101】-【0105】, 第18 図 (ファミリーなし)	5-8 1-4, 9-14
<input checked="" type="checkbox"/> C欄の続きにも文献が列挙されている。 <input type="checkbox"/> パテントファミリーに関する別紙を参照。		
* 引用文献のカテゴリー 「A」 特に関連のある文献ではなく、一般的技術水準を示すもの 「E」 国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの 「L」 優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献 (理由を付す) 「O」 口頭による開示、使用、展示等に言及する文献 「P」 国際出願日前で、かつ優先権の主張の基礎となる出願日の後に公表された文献 「T」 国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの 「X」 特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの 「Y」 特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの 「&」 同一パテントファミリー文献		
国際調査を完了した日 30.06.2014	国際調査報告の発送日 08.07.2014	
国際調査機関の名称及びあて先 日本国特許庁 (ISA/J P) 郵便番号100-8915 東京都千代田区霞が関三丁目4番3号	特許庁審査官 (権限のある職員) 野田 佳邦 電話番号 03-3581-1101 内線 3565	5U 3450

C (続き) . 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
A	JP 2011-204060 A (日本電気株式会社) 2011. 10. 13, 全文, 全図 & US 2011/0238908 A1	1-14
A	JP 2009-205335 A (株式会社日立製作所) 2009. 09. 10, 全文, 全図 & US 2009/0216945 A1 & US 2011/0296091 A1 & US 2012/0254523 A1	1-14