

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2017-111476

(P2017-111476A)

(43) 公開日 平成29年6月22日(2017.6.22)

(51) Int.Cl.			F I	テーマコード (参考)		
G06F	12/16	(2006.01)	G06F 12/16	310A	5B018	
G06F	12/02	(2006.01)	G06F 12/02	530C	5B060	
G06F	12/00	(2006.01)	G06F 12/00	597U		

審査請求 未請求 請求項の数 13 O L (全 35 頁)

(21) 出願番号 特願2015-242997 (P2015-242997)
 (22) 出願日 平成27年12月14日 (2015.12.14)

(71) 出願人 000003078
 株式会社東芝
 東京都港区芝浦一丁目1番1号
 (74) 代理人 110001737
 特許業務法人スズエ国際特許事務所
 (72) 発明者 菅野 伸一
 東京都港区芝浦一丁目1番1号 株式会社東芝内
 Fターム(参考) 5B018 GA04 HA23 KA18 MA22 NA06
 5B060 AA10

(54) 【発明の名称】 メモリシステムおよび制御方法

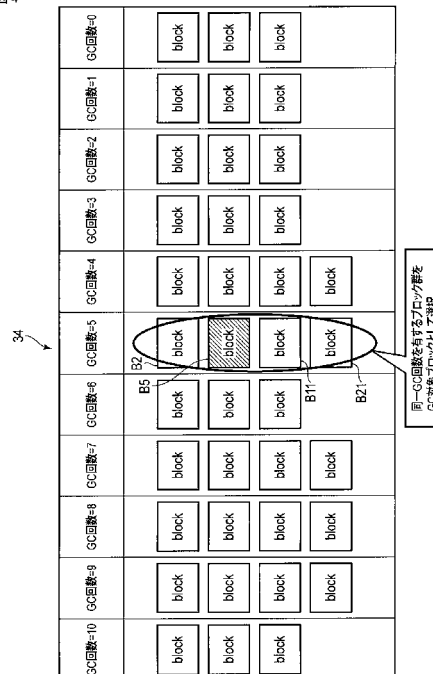
(57) 【要約】

【課題】データ局所性によるライトアンプリフィケーションを増加を抑制する。

【解決手段】実施形態によれば、メモリシステムは、不揮発性メモリと、コントローラとを具備する。前記コントローラは、前記複数のブロックの内の、ホストによって書き込まれたデータを含むブロック毎に、当該ブロック内のデータが前記ガベージコレクション動作によってコピーされた回数を示すガベージコレクション回数を管理する。前記コントローラは、同じガベージコレクション回数に関連づけられた複数の第1ブロックを、前記ガベージコレクション動作の対象ブロックとして選択する。前記コントローラは、前記複数の第1ブロック内の有効データをコピー先フリーブロックにコピーする。前記コントローラは、前記複数の第1ブロックのガベージコレクション回数に1を加えた値を、前記コピー先フリーブロックのガベージコレクション回数として設定する。

【選択図】 図4

図4



【特許請求の範囲】**【請求項 1】**

複数のブロックを含む不揮発性メモリと、
前記不揮発性メモリに電氣的に接続され、前記不揮発性メモリのガベージコレクション動作を実行するように構成されたコントローラとを具備し、

前記コントローラは、

前記複数のブロックの内の、ホストによって書き込まれたデータを含むブロック毎に、当該ブロック内のデータが前記ガベージコレクション動作によってコピーされた回数を示すガベージコレクション回数を管理し、

同じガベージコレクション回数に関連づけられた複数の第 1 ブロックを、前記ガベージコレクション動作の対象ブロックとして選択し、

前記複数の第 1 ブロック内の有効データをコピー先フリーブロックにコピーし、

前記複数の第 1 ブロックのガベージコレクション回数に 1 を加えた値を、前記コピー先フリーブロックのガベージコレクション回数として設定するように構成されているメモリシステム。

【請求項 2】

前記複数の第 1 ブロックは、無効データ量が最も多いブロックと、前記無効データ量が最も多いブロックのガベージコレクション回数と同じガベージコレクション回数に関連づけられた一以上の別のブロックとを含む請求項 1 記載のメモリシステム。

【請求項 3】

前記コントローラは、

同じガベージコレクション回数に関連づけられた第 1 ブロック群が前記ガベージコレクション動作の対象ブロックとして選択された場合、

前記第 1 ブロック群内の有効データの総量が第 1 閾値よりも少ないか否かを判定し、

前記第 1 ブロック群内の有効データの総量が前記第 1 閾値よりも少ない場合、前記第 1 ブロック群のガベージコレクション回数よりも 1 回以上少ないガベージコレクション回数に関連づけられた全てのブロック群の中で、最大のガベージコレクション回数に関連づけられている第 2 ブロック群を選択し、前記第 1 ブロック群の有効データと前記第 2 ブロック群内の有効データを前記コピー先フリーブロックにコピーするように構成されている請求項 1 記載のメモリシステム。

【請求項 4】

前記第 1 閾値は、一つのブロックに書き込み可能なデータの総量を示す値に設定されている請求項 3 記載のメモリシステム。

【請求項 5】

前記コントローラは、前記第 1 ブロック群のガベージコレクション回数に 1 を加えた値を、前記コピー先フリーブロックのガベージコレクション回数として設定するように構成されている請求項 3 記載のメモリシステム。

【請求項 6】

前記コントローラは、前記第 2 ブロック群のガベージコレクション回数に 1 を加えた値を、前記コピー先フリーブロックのガベージコレクション回数として設定するように構成されている請求項 3 記載のメモリシステム。

【請求項 7】

前記コントローラは、

同じガベージコレクション回数に関連づけられた第 1 ブロック群が前記ガベージコレクション動作の対象ブロックとして選択された場合、

前記第 1 ブロック群内の有効データの総量が第 1 閾値よりも少ないか否かを判定し、

前記第 1 ブロック群内の有効データの総量が前記第 1 閾値よりも少ない場合、前記第 1 ブロック群のガベージコレクション回数が第 2 閾値以上であるか否かを判定し、

前記第 1 ブロック群のガベージコレクション回数が前記第 2 閾値以上である場合、前記第 1 ブロック群のガベージコレクション回数よりも 1 回以上少ないガベージコレクション

10

20

30

40

50

回数に関連づけられた全てのブロック群の中で、最大のガベージコレクション回数に関連づけられている第2ブロック群を選択し、前記第1ブロック群の有効データと前記第2ブロック群内の有効データを前記コピー先フリーブロックにコピーし、

前記第1ブロック群のガベージコレクション回数が前記第2閾値よりも少ない場合、前記第2ブロック群の選択と前記第1ブロック群および前記第2ブロック群の有効データのコピーの代わりに、同じガベージコレクション回数を有する第3ブロック群を前記ガベージコレクション動作の対象ブロックとして選択するように構成されている請求項1記載のメモリシステム。

【請求項8】

前記第1閾値は、一つのブロックに書き込み可能なデータの総量を示す値に設定されている請求項7記載のメモリシステム。

10

【請求項9】

前記コントローラは、前記第1ブロック群のガベージコレクション回数に1を加えた値を、前記コピー先フリーブロックのガベージコレクション回数として設定するように構成されている請求項7記載のメモリシステム。

【請求項10】

前記コントローラは、前記第2ブロック群のガベージコレクション回数に1を加えた値を、前記コピー先フリーブロックのガベージコレクション回数として設定するように構成されている請求項7記載のメモリシステム。

【請求項11】

20

複数のブロックを含む不揮発性メモリと、
前記不揮発性メモリに電氣的に接続され、前記不揮発性メモリのガベージコレクション動作を実行するように構成されたコントローラとを具備し、

前記コントローラは、

前記複数のブロックの内の、ホストによって書き込まれたデータを含むブロック毎に、当該ブロック内のデータが前記ガベージコレクション動作によってコピーされた回数を示すガベージコレクション回数を管理し、

前記ホストによって書き込まれたデータを含むブロックから、無効データ量が最も多いブロックを選択し、

前記無効データ量が最も多いブロックと、前記無効データ量が最も多いブロックのガベージコレクション回数と同じガベージコレクション回数に関連づけられた一以上の第1ブロックとを、前記ガベージコレクション動作の対象ブロックとして選択し、

30

前記無効データ量が最も多いブロックの有効データと前記一以上の第1ブロック内の有効データとを、コピー先フリーブロックにコピーし、

前記ガベージコレクション動作の前記対象ブロックのガベージコレクション回数に1を加えた値を、前記コピー先フリーブロックのガベージコレクション回数として設定するように構成されているメモリシステム。

【請求項12】

複数のブロックを含む不揮発性メモリを制御し、前記不揮発性メモリのガベージコレクション動作を実行する制御方法であって、

40

前記複数のブロックの内の、ホストによって書き込まれたデータを含むブロック毎に、当該ブロック内のデータが前記ガベージコレクション動作によってコピーされた回数を示すガベージコレクション回数を管理することと、

同じガベージコレクション回数に関連づけられた複数の第1ブロックを、前記ガベージコレクション動作の対象ブロックとして選択することと、

前記複数の第1ブロック内の有効データをコピー先フリーブロックにコピーすることと、

前記複数の第1ブロックのガベージコレクション回数に1を加えた値を、前記コピー先フリーブロックのガベージコレクション回数として設定することとを具備する制御方法。

【請求項13】

50

前記複数の第1ブロックは、無効データ量が最も多いブロックと、前記無効データ量が最も多いブロックのガベージコレクション回数と同じガベージコレクション回数に関連づけられた一以上の別のブロックとを含む請求項12記載の制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明の実施形態は、不揮発性メモリを制御する技術に関する。

【背景技術】

【0002】

近年、不揮発性メモリを備えるメモリシステムが広く普及している。

10

【0003】

このようなメモリシステムの一つとして、NANDフラッシュ技術ベースのソリッドステートドライブ(SSD)が知られている。

【0004】

SSDは、その低電力消費、高性能という特徴により、様々なコンピュータのメインストレージとして使用されている。

【先行技術文献】

【特許文献】

【0005】

【特許文献1】国際公開第2012/020544号

20

【発明の概要】

【発明が解決しようとする課題】

【0006】

ところで、ホストによってSSDに書かれるデータには、そのデータの一部が頻繁に書き換えられ、残りの部分は頻繁に書き換えられない、というデータ局所性が存在する場合がある。

【0007】

このようなデータ局所性は、SSDのライトアンプリフィケーションを増加させ、結果としてSSDの性能および寿命に影響を及ぼす場合がある。

【0008】

30

本発明が解決しようとする課題は、データ局所性によるライトアンプリフィケーションを増加を抑制することができるメモリシステムおよび制御方法を提供することである。

【課題を解決するための手段】

【0009】

実施形態によれば、メモリシステムは、不揮発性メモリと、前記不揮発性メモリに電氣的に接続され、前記不揮発性メモリのガベージコレクション動作を実行するように構成されたコントローラとを具備する。前記コントローラは、前記複数のブロックの内の、ホストによって書き込まれたデータを含むブロック毎に、当該ブロック内のデータが前記ガベージコレクション動作によってコピーされた回数を示すガベージコレクション回数を管理する。前記コントローラは、同じガベージコレクション回数に関連づけられた複数の第1ブロックを、前記ガベージコレクション動作の対象ブロックとして選択する。前記コントローラは、前記複数の第1ブロック内の有効データをコピー先フリーブロックにコピーする。前記コントローラは、前記複数の第1ブロックのガベージコレクション回数に1を加えた値を、前記コピー先フリーブロックのガベージコレクション回数として設定する。

40

【図面の簡単な説明】

【0010】

【図1】実施形態に係るメモリシステムの構成例を説明するブロック図。

【図2】同実施形態のメモリシステムによって実行される、ガベージコレクション回数管理動作とガベージコレクション動作とを説明するための図。

【図3】同実施形態のメモリシステムにおいて用いられるガベージコレクション(GC)

50

回数管理リストの例を説明する図。

【図 4】ガベージコレクション回数管理リストに基づいて同実施形態のメモリシステムによって実行されるガベージコレクション対象ブロック選択動作を説明するための図。

【図 5】同実施形態のメモリシステムによって実行されるガベージコレクション動作を説明する図。

【図 6】同実施形態のメモリシステムに書き込まれる複数種のデータの例を説明する図。

【図 7】ガベージコレクション回数と複数種のデータ間のデータ量の割合との関係の例を説明する図。

【図 8】同実施形態のメモリシステムによって実行されるガベージコレクション動作の手順を説明するフローチャート。

【図 9】同実施形態のメモリシステムによって実行される、異なるガベージコレクション回数を有する 2 つのブロック群の有効データをマージする処理を含むガベージコレクション動作を説明する図。

【図 10】同実施形態のメモリシステムによって実行される、異なるガベージコレクション回数を有する 2 つのブロック群の有効データをマージする処理を含むガベージコレクション動作の手順を説明するフローチャート。

【図 11】マージ処理を特定のガベージコレクション回数以上のブロック群にのみ許可する動作を説明する図。

【図 12】マージ処理を特定のガベージコレクション回数以上のブロック群にのみ許可する動作を含むガベージコレクション動作の手順を説明するフローチャート。

【図 13】同実施形態のメモリシステムによって実行される、ホストからのデータの書き込み用にフリーブロックを順次割り当てる動作を説明するための図。

【図 14】同実施形態のメモリシステムによって使用されるブロック使用順序管理リストの例を説明するための図。

【図 15】同じ L B A へのライトが要求された時に同実施形態のメモリシステムによって実行される累積データ書き込み量算出動作を説明するための図。

【図 16】同実施形態のメモリシステムによって実行される累積データ書き込み量応答処理の処理シーケンスを説明するための図。

【図 17】同実施形態のメモリシステムによって実行される累積データ書き込み量応答処理の手順を説明するフローチャート。

【図 18】同実施形態のメモリシステムによって実行される累積データ書き込み量応答処理の別の処理シーケンスを説明するための図。

【図 19】同実施形態のメモリシステムによって実行される累積データ書き込み量応答処理の別の手順を説明するフローチャート。

【図 20】同実施形態のメモリシステムにおいて使用されるルックアップテーブルの例を説明するための図。

【図 21】同じ L B A へのライトが要求された時に同実施形態のメモリシステムによって実行される時間経過応答処理の手順を説明するためのフローチャート。

【図 22】同実施形態のメモリシステムから受信される累積データ書き込み量 / 時間経過情報に基づいてホストによって実行される処理の手順の例を説明するフローチャート。

【図 23】ホストの構成例を説明するブロック図。

【図 24】同実施形態のメモリシステムとホストとを含むコンピュータの構成例を示す図。

【発明を実施するための形態】

【0011】

以下、図面を参照して、実施形態を説明する。

まず、図 1 を参照して、一実施形態に係るメモリシステムを含む情報処理システム 1 の構成を説明する。

【0012】

このメモリシステムは、不揮発性メモリにデータをライトし、不揮発性メモリからデー

10

20

30

40

50

タをリードするように構成された半導体ストレージデバイスである。このメモリシステムは、例えば、NANDフラッシュ技術ベースのソリッドステートドライブ(SSD)3として実現されている。

【0013】

情報処理システム1は、ホスト(ホストデバイス)2と、SSD3とを含む。ホスト2は、サーバ、パーソナルコンピュータのような情報処理装置である。

【0014】

SSD3は、ホスト2として機能する情報処理装置のメインストレージとして使用され得る。SSD3は、情報処理装置に内蔵されてもよいし、情報処理装置にケーブルまたはネットワークを介して接続されてもよい。

10

【0015】

ホスト2とSSD3とを相互接続するためのインタフェースとしては、SCSI、Serial Attached SCSI(SAS)、ATA、Serial ATA(SATA)、PCI Express(PCIE)、Ethernet(登録商標)、Fibre channel等が使用し得る。

【0016】

SSD3は、コントローラ4、不揮発性メモリ(NANDメモリ)5、およびDRAM6を備える。NANDメモリ5は、限定されないが、複数のNANDフラッシュメモリチップを含んでいてもよい。

20

【0017】

NANDメモリ5は、多数のNANDブロック(ブロック)B0~Bm-1を含む。ブロックB0~Bm-1は、消去単位として機能する。ブロックは「物理ブロック」または「消去ブロック」と称されることもある。

【0018】

ブロックB0~Bm-1は多数のページ(物理ページ)を含む。つまり、ブロックB0~Bm-1の各々は、ページP0~Pn-1を含む。NANDメモリ5においては、データのリードおよびデータのライトはページ単位で実行される。データの消去はブロック単位で実行される。

【0019】

コントローラ4は、Toggle、ONFIのようなNANDインタフェース13を介して、不揮発性メモリであるNANDメモリ5に電氣的に接続されている。コントローラ4は、NANDメモリ5のデータ管理とNANDメモリ5のブロック管理とを実行するように構成されたフラッシュトランシエーション層(FTL)として機能し得る。

30

【0020】

データ管理には、(1)論理ブロックアドレス(LBA)と物理アドレスとの間の対応関係を示すマッピング情報の管理、(2)ページ単位のリード/ライトとブロック単位の消去動作とを隠蔽するための処理、等が含まれる。LBAと物理アドレスとの間のマッピングの管理は、論理物理アドレス変換テーブルとして機能するルックアップテーブル(LUT)33を用いて実行される。ルックアップテーブル(LUT)33は、所定の管理サイズ単位で、LBAと物理アドレスとの間のマッピングを管理する。ホスト2からのライトコマンドの多くは、4Kバイトのデータの書き込みを要求する。したがって、ルックアップテーブル(LUT)33は、例えば4Kバイト単位で、LBAと物理アドレスとの間のマッピングを管理してもよい。あるLBAに対応する物理アドレスは、このLBAのデータがライトされたNANDメモリ5内の物理記憶位置を示す。物理アドレスは、物理ブロックアドレスと物理ページアドレスとを含む。物理ページアドレスは全てのページに割り当てられており、また物理ブロックアドレスは全てのブロックに割り当てられている。

40

【0021】

ページへのデータ書き込みは、1消去サイクル当たり1回のみ可能である。

【0022】

このため、コントローラ4は、同じLBAへのライト(上書き)を、NANDメモリ5

50

上の別のページにマッピングする。つまり、コントローラ 4 は、この別のページにデータをライトする。そして、コントローラ 4 は、ルックアップテーブル (LUT) 33 を更新してこの LBA をこの別のページに関連付けると共に、元のページ (つまりこの LBA が関連付けられていた古いデータ) を無効化する。

【0023】

ブロック管理には、不良ブロックの管理と、ウェアレベリングと、ガベージコレクション動作等が含まれる。ウェアレベリングは、物理ブロックそれぞれのプログラム/イレース回数を平準化するための動作である。

【0024】

ガベージコレクション動作は、NANDメモリ 5 内のフリースペースを作り出すための動作である。このガベージコレクション動作は、NANDメモリ 5 のフリーブロックの個数を増やすため、有効データと無効データとが混在する幾つかのブロック内の全ての有効データを別のブロック (コピー先フリーブロック) にコピーする。そして、ガベージコレクション動作は、ルックアップテーブル (LUT) 33 を更新して、コピーされた有効データの LBA それぞれを正しい物理アドレスにマッピングする。有効データが別のブロックにコピーされることによって無効データのみになったブロックはフリーブロックとして開放される。これによって、このブロックは消去後に再利用することが可能となる。

10

【0025】

ホスト 2 は、ライトコマンドを SSD 3 に送出する。このライトコマンドは、ライトデータ (つまり書き込むべきデータ) の論理アドレス (開始論理アドレス) と、転送長とを含む。この実施形態においては、LBA が論理アドレスとして使用されるが、他の実施形態においてはオブジェクト ID が論理アドレスとして使用されても良い。LBA は、論理セクタ (論理ブロック) に付与されるシリアル番号によって表現される。シリアル番号はゼロから始まる。論理セクタのサイズは、例えば 512 バイトである。

20

【0026】

SSD 3 のコントローラ 4 は、ライトコマンド内の開始論理アドレスと転送長とによって指定されるライトデータを、NANDメモリ 5 内のブロックのページにライトする。さらに、コントローラ 4 は、ルックアップテーブル (LUT) 33 を更新することによって、ライトされたデータに対応する LBA を、このデータがライトされた物理記憶位置を示す物理アドレスにマッピングする。

30

【0027】

より詳しくは、コントローラ 4 は、NANDメモリ 5 内のフリーブロックの一つを、ホスト 2 からのデータの書き込みのために割り当てる。この割り当てられたブロックは、ホスト 2 からのデータが書き込まれるべき書き込み対象ブロックであり、「書き込み先ブロック」、または「入力ブロック」、等とも称される。コントローラ 4 は、ルックアップテーブル (LUT) 33 を更新しながら、ホスト 2 から受信されるライトデータを書き込み対象ブロック (書き込み先ブロック) の利用可能ページに順次書き込む。書き込み先ブロックに利用可能ページが無くなった場合に、コントローラ 4 は、新たなフリーブロックを書き込み先ブロックとして割り当てる。

【0028】

次に、コントローラ 4 の構成について説明する。

40

【0029】

コントローラ 4 は、ホストインタフェース 11、CPU 12、NAND インタフェース 13、DRAM インタフェース 14、SRAM 15 等を含む。これら CPU 12、NAND インタフェース 13、DRAM インタフェース 14、SRAM 15 は、バス 10 を介して相互接続される。

【0030】

ホストインタフェース 11 は、ホスト 2 から様々なコマンド (ライトコマンド、リードコマンド、アンマップ (UNMAP) コマンド、等) を受信する。

【0031】

50

ライトコマンドは、SSD3に対し、このライトコマンドによって指定されたデータをライトするように要求する。ライトコマンドは、ライトされるべき最初の論理ブロックのLBA（開始LBA）と、転送長（論理ブロックの数）とを含む。リードコマンドは、SSD3に対し、このリードコマンドによって指定されたデータをリードするように要求する。リードコマンドは、リードされるべき最初の論理ブロックのLBA（開始LBA）と、転送長（論理ブロックの数）とを含む。

【0032】

CPU12は、ホストインタフェース11、NANDインタフェース13、DRAMインタフェース14、SRAM15を制御するように構成されたプロセッサである。CPU12は、上述のFTLの処理に加え、ホスト2からの様々なコマンドを処理するためのコマンド処理等を実行する。

10

【0033】

例えば、コントローラ4がホスト2からライトコマンドを受信した時、CPU12の制御の下、コントローラ4はライトコマンドによって指定されるライトデータをNANDメモリ5に書き込む以下のライト動作を実行する。

【0034】

つまり、コントローラ4は、ライトデータを現在の書き込み先ブロックの物理記憶位置（利用可能ページ）に書き込み、そしてルックアップテーブル（LUT）33を更新して、ライトコマンドに含まれるLBA（開始LBA）にこの物理記憶位置の物理アドレスをマッピングする。

20

【0035】

これらFTL処理およびコマンド処理は、CPU12によって実行されるファームウェアによって制御されてもよい。このファームウェアは、CPU12を、ガベージコレクション（GC）回数管理部21、ガベージコレクション（GC）動作制御部22、および更新頻度情報応答部23として機能させる。

【0036】

ホスト2によってSSD3に書かれるデータには、そのデータの一部が頻繁に書き換えられ、残りの部分は頻繁には書き換えられない、というデータ局所性が存在する場合がある。この場合、例えば、無効データ量の多い上位幾つかのブロックをGC対象ブロックとして選択するという通常のGCアルゴリズムによってGC動作が実行されると、何度もGC動作が繰り返されるにつれて、更新頻度の高いデータと更新頻度の低いデータとが同じブロックに混在しやくなる。更新頻度の高いデータと更新頻度の低いデータとの混在は、SSD3のライトアンプリフィケーションを増加させる要因となり得る。

30

【0037】

なぜなら、更新頻度の高いデータ（Hotデータ）と更新頻度の低いデータ（Coldデータ）とが混在するブロックにおいては、Hotデータの更新によってブロック内の一部の領域だけが早いタイミングで無効化される一方、このブロック内の残りの領域（Coldデータ）は有効状態に長い間維持されるからである。

【0038】

もしHotデータのみによってブロックが満たされていたならば、このブロック内の全てのデータがそれらデータの更新（書き替え）によって、比較的速いタイミングで無効化される可能性が高い。したがって、このブロックは、ガベージコレクション動作を実行すること無しで、このブロックを消去することのみによって、再利用することが可能となる。

40

【0039】

一方、Coldデータのみによってブロックが満たされているならば、このブロック内の全てのデータは、長い間、有効状態に維持される。したがって、このブロックは、長い間、ガベージコレクション動作の対象とならない可能性が高い。

【0040】

ライトアンプリフィケーション（WA）は、以下のように定義される。

50

【 0 0 4 1 】

WA=「SSDにライトされたデータの総量」 / 「ホストからSSDにライトされたデータの総量」

「SSDにライトされたデータの総量」は、ホストからSSDにライトされたデータの総量とガベージコレクション動作等によって内部的にSSDにライトされたデータの総量との和に相当する。

【 0 0 4 2 】

ライトアンプリフィケーション(WA)の増加は、SSD3内のブロックそれぞれの書き換え回数(プログラム/イレース回数)の増加を引き起こす。つまり、ライトアンプリフィケーション(WA)が大きい程、ブロックのプログラム/イレース回数が、そのプログラム/イレース回数の上限値に速く達しやすくなる。この結果、SSD3の耐久性および寿命の劣化が引き起こされる。

【 0 0 4 3 】

本実施形態では、更新頻度の高いデータと更新頻度の低いデータとを分離できるようにするために、「ブロック内のデータのGC回数を考慮したGC機能」と、「LBAベースの更新頻度通知機能」とを有している。

【 0 0 4 4 】

ガベージコレクション(GC)回数管理部21およびガベージコレクション(GC)動作制御部22は、「ブロック内のデータのGC回数を考慮したGC機能」を実行する。「ブロック内のデータのGC回数を考慮したGC機能」は、データ局所性に起因するSSD3のライトアンプリフィケーションの増加を抑制可能な改善されたガベージコレクション(GC)動作を実行する。

【 0 0 4 5 】

ガベージコレクション(GC)回数管理部21は、ホスト2によって書き込まれたデータを含むブロック毎に、ガベージコレクション(GC)回数を管理する。あるブロックのGC回数は、このブロック内のデータがガベージコレクション(GC)動作によってコピーされた回数を示す。つまり、あるブロックのGC回数は、このブロック内のデータが有効データとして過去に何回コピーされたかを示す。

【 0 0 4 6 】

ホスト2によってデータがライトされた直後のブロック、つまりそのデータがGCによって一度も集められた(コピーされた)ことのないブロックについては、このブロックのGC回数はゼロに設定される。

【 0 0 4 7 】

GC回数がゼロである幾つかのブロックがGCの対象ブロック(コピー元ブロック)として選択され、これらブロックの有効データがコピー先フリーブロックにコピーされたならば、GC回数管理部21は、このコピー先フリーブロックのGC回数を1に設定する。コピー先フリーブロック内のデータは、有効データとしてGC対象ブロック(コピー元ブロック)から1回コピーされたデータであるからである。

【 0 0 4 8 】

有効データがコピー先フリーブロックにコピーされることによって無効データのみになった各ブロック(コピー元ブロック)は、フリーブロックとなる。フリーブロックはデータを含まないため、このフリーブロックのGC回数は管理する必要はない。

【 0 0 4 9 】

GC回数が1である幾つかのブロック(コピー元ブロック)がガベージコレクション(GC)の対象ブロックとして選択され、これらブロックの有効データがコピー先フリーブロックにコピーされたならば、GC回数管理部21は、このコピー先のフリーブロックのGC回数を2に設定する。このコピー先フリーブロック内のデータは、有効データとして過去に2回コピーされたデータであるからである。

【 0 0 5 0 】

このように、あるブロックに関連づけられたGC回数の値は、このブロック内のデータ

10

20

30

40

50

が過去のGC動作によって何回コピーされたか、つまりこのブロック内のデータに対して過去に何回のGC動作が実行されたかを示す。

【0051】

GC動作制御部22は、同じGC回数に関連づけられた幾つかのブロックをガベージコレクション(GC)動作の対象ブロックとして選択し、これら同じGC回数に関連づけられたブロックの有効データのみを同じコピー先ブロックにコピーする、という改良されたGC動作を実行する。

【0052】

例えば、GC動作制御部22は、同じGC回数に関連づけられたブロック群(つまり、同じGC回数を有するブロックの集合)の中から、幾つかのブロックを、GCのための対象ブロックとして選択する。GC動作制御部22は、このGC対象ブロックとして選択されたこれらブロック内の有効データをコピー先フリーブロックにコピーする。そして、GC対象ブロックとして選択されたこれらブロックのGC回数に1を加えた値が、GC回数管理部21によって、このコピー先フリーブロックのGC回数として設定される。

10

【0053】

「LBAベースの更新頻度通知機能」は、個々のライトコマンドに含まれる個々のLBAへのライトの頻度をホスト2に通知することによって、ホスト2がHotデータ/Coldデータを分離するのを効率良くアシストする機能である。

【0054】

「LBAベースの更新頻度通知機能」は、更新頻度情報応答部23によって実行される。

20

【0055】

更新頻度情報応答部23は、ホスト2からLBAを含むライトコマンドを受信した際に、このLBAへの前回のライトからこのLBAへの今回のライトまでの時間経過に関する値、またはこのLBAへの前回のライトからこのLBAへの今回のライトまでの累積データ書き込み量を、このライトコマンドに対する応答としてホスト2に通知する。これにより、ホスト2に対して、ユーザデータの実際の更新頻度(書き替え頻度)を知らせることができるので、ホスト2は、ユーザデータを互いに更新頻度の異なる複数種のデータ、例えば、頻繁に更新されるタイプのデータ(Hotデータ)、更新の頻度が低いタイプのデータ(Coldデータ)、HotデータとColdデータの間での更新頻度を有するタイプのデータ(Warmデータ)に、に分類できる。この結果、例えば、ホスト2は、必要に応じて、これら異なるタイプのデータを、異なるSSDに分散させるための処理等を実行することができる。

30

【0056】

次に、コントローラ4内の他のコンポーネントについて説明する。

【0057】

NANDインタフェース13は、CPU12の制御の下、NANDメモリ5を制御するように構成されたNANDコントローラである。

【0058】

DRAMインタフェース14は、CPU12の制御の下、DRAM6を制御するように構成されたDRAMコントローラである。

40

【0059】

DRAM6の記憶領域の一部は、NANDメモリ5にライトすべきデータを一時的に格納するためのライトバッファ(WB)31として利用されてもよい。また、DRAM6の記憶領域は、ガベージコレクション(GC)動作中に移動されるデータを一時的に格納するためのGCバッファ32として利用されてもよい。また、DRAM6の記憶領域は、上述のルックアップテーブル33の格納のために用いられてもよい。

【0060】

さらに、DRAM6の記憶領域は、GC回数管理リスト34、およびブロック使用順序管理リスト35として利用されてもよい。

50

【 0 0 6 1 】

G C回数管理リスト3 4は、ホスト2によって書き込まれたデータを含むブロック毎に、G C回数を保持するためのリストである。G C回数管理リスト3 4は、ブロックそれぞれのブロックID（例えば物理ブロックアドレス）とこれらブロック内のデータのG C回数との間の対応関係を示す表であってもよい。

【 0 0 6 2 】

あるいは、G C回数管理リスト3 4は、G C回数（例えば、G C回数 = 0 ~ G C回数 = n）別にブロックそれぞれを管理するための複数のG C回数リストから構成されてもよい。ここで、nは、管理すべきG C回数の上限値である。例えば、G C回数 = 0のG C回数リストは、0回のG C回数に関連づけられたブロックそれぞれのブロックID（例えば物理ブロックアドレス）のリストを保持する。G C回数 = 1のG C回数リストは、1回のG C回数に関連づけられたブロックそれぞれのブロックID（例えば物理ブロックアドレス）のリストを保持する。

10

【 0 0 6 3 】

ブロック使用順序管理リスト3 5は、書き込み先ブロック用に割り当てられたブロックそれぞれに付与される割り当て番号（シーケンシャル番号）を保持する。すなわち、コントローラ4は、書き込み対象ブロックとして割り当てられたブロックそれぞれに対してその割り当て順序を示す番号（割り当て番号）を付与する。番号は1から始まるシーケンシャル番号であってもよい。例えば、最初に書き込み先ブロック用に割り当てられたブロックには割り当て番号 = 1が付与され、2番目に書き込み先ブロック用に割り当てられたブロックには割り当て番号 = 2が付与され、3番目に書き込み先ブロック用に割り当てられたブロックには割り当て番号 = 3が付与される。これにより、どのブロックがどのような順序で書き込み先ブロックとして割り当てられたかを示すブロック使用履歴を管理することができる。割り当て番号としては、新たなフリーブロックが書き込み先ブロック用に割り当てられる度にインクリメントされるカウンタの値を使用できる。

20

【 0 0 6 4 】

S S D 3は、さらに他の様々な管理情報を保持していてもよい。このような管理情報の例には、物理アドレスそれぞれに対応する有効 / 無効フラグを保持するページ管理テーブルが含まれていても良い。各有効 / 無効フラグは、対応する物理アドレス（物理ページ）が有効であるか無効であることを示す。物理ページが有効であるとは、その物理ページ内のデータが有効データであることを意味する。物理ページが無効であるとは、その物理ページ内のデータが更新（書き替え）によって無効化されたデータであることを意味する。

30

【 0 0 6 5 】

次に、ホスト2の構成について説明する。

【 0 0 6 6 】

ホスト2は、様々なプログラムを実行する情報処理装置である。情報処理装置によって実行されるプログラムには、アプリケーションソフトウェアレイヤ4 1、オペレーティングシステム（OS）4 2、ファイルシステム4 3が含まれる。

【 0 0 6 7 】

一般に知られているように、オペレーティングシステム（OS）4 2は、ホスト2全体を管理し、ホスト2内のハードウェアを制御し、アプリケーションがハードウェアおよびS S D 3を使用することを可能にするための制御を実行するように構成されたソフトウェアである。

40

【 0 0 6 8 】

ファイルシステム4 3は、ファイルの操作（作成、保存、更新、削除等）のための制御を行うために使用される。例えば、Z F S、B t r f s、X F S、e x t 4、N T F Sなどがファイルシステム4 3として使用されても良い。あるいは、ファイルオブジェクトシステム（例えば、Ceph Object Storage Daemon）、Key Value Store System（例えば、Rocks DB）がファイルシステム4 3として使用されても良い。

【 0 0 6 9 】

50

様々なアプリケーションソフトウェアスレッドがアプリケーションソフトウェアレイヤ 4 1 上で走る。アプリケーションソフトウェアスレッドの例としては、クライアントソフトウェア、データベースソフトウェア、仮想マシン等がある。

【 0 0 7 0 】

アプリケーションソフトウェアレイヤ 4 1 がリードコマンドまたはライトコマンドのようなリクエストを S S D 3 に送付することが必要な時、アプリケーションソフトウェアレイヤ 4 1 は、O S 4 2 にそのリクエストを送付する。O S 4 2 はそのリクエストをファイルシステム 4 3 に送付する。ファイルシステム 4 3 は、そのリクエストを、コマンド（リードコマンド、ライトコマンド等）にトランスレートする。ファイルシステム 4 3 は、コマンドを、S S D 3 に送付する。S S D 3 からのレスポンスが受信された際、ファイルシステム 4 3 は、そのレスポンスを O S 4 2 に送付する。O S 4 2 は、そのレスポンスをアプリケーションソフトウェアレイヤ 4 1 に送付する。

10

【 0 0 7 1 】

次に、図 2 ~ 図 1 2 を参照して、「ブロック内のデータの G C 回数を考慮した G C 機能」の詳細を説明する。

【 0 0 7 2 】

図 2 は、S S D 3 によって実行される G C 回数管理動作と G C 動作とを示す。

【 0 0 7 3 】

S S D 3 のコントローラ 4 は、あるフリーブロックを、ホスト 2 からのデータ（ライトデータ）の書き込み用のブロック（書き込み先ブロック）として割り当て、ホスト 2 から受信されるライトデータをこの書き込み先ブロック内の利用可能ページに順次書き込む。現在の書き込み先ブロックの全てのページがデータで満たされた時、コントローラ 4 は、現在の書き込み先ブロックをアクティブブロック（データを含むブロック）として管理する。さらに、コントローラ 4 は、別のフリーブロックを新たな書き込み先ブロックとして割り当てる。このようにして、S S D 3 においては、ホスト 2 から受信されるデータ（ライトデータ）は、その到着順に、現在の書き込み先ブロックの最初のページから最後のページに向けて順次書き込まれる。

20

【 0 0 7 4 】

図 2 のブロック B 1 1 ~ B 1 7 は、ホスト 2 によってデータがライトされた直後のブロック、つまりそのブロック内のデータがガベージコレクション（G C）動作によって一度もコピーされたことのないブロックである。これらブロック B 1 1 ~ B 1 7 に対応する G C 回数は 0 である。

30

【 0 0 7 5 】

時間が経過するにつれ、ブロック B 1 1 ~ B 1 7 の各々のデータの一部は、その書き換えによって無効化されるかもしれない。これにより、ブロック B 1 1 ~ B 1 7 の各々においては、有効データと無効データとが混在される場合がある。

【 0 0 7 6 】

フリーブロックの数が閾個数以下に低下した場合、コントローラ 4 は、有効データと無効データとが混在される幾つかのブロックからフリーブロックを作り出す G C 動作を開始する。

40

【 0 0 7 7 】

コントローラ 4 は、まず、有効データと無効データとが混在する幾つかのブロックを G C 対象ブロックとして選択する。この G C 対象ブロックの選択においては、コントローラ 4 は、上述したように、同じ G C 回数に関連づけられたブロック群を G C 対象ブロックとして選択する。このブロック群は、例えば、最も無効データの量が多いブロックが属するブロック群、つまり最も無効データの量が多いブロックの G C 回数と同じ G C 回数を有するブロックの集合、であってよい。この場合、コントローラ 4 は、最初に、ホスト 2 によって書き込まれたデータを含むブロックから、無効データ量が最も多いブロックを選択してもよい。次いで、コントローラ 4 は、無効データ量が最も多いブロックと、この無効データ量が最も多いブロックの G C 回数と同じ G C 回数に関連づけられた一以上のブロック

50

とを、ガベージコレクション（GC）動作の対象ブロックとして選択してもよい。

【0078】

コントローラ4は、選択した幾つかのGC対象ブロック（同じGC回数に関連づけられた幾つかのブロック）内の有効データをコピー先フリーブロックにコピーし、これらGC対象ブロックのGC回数に1を加えた値を、コピー先フリーブロックのGC回数として設定する。これにより、GC対象ブロックのGC回数に1を加えた値がコピー先フリーブロックに引き継がれるので、コピー先フリーブロックのGC回数は、そのコピー先フリーブロック内のデータがGC動作によって過去に何回コピーされたかを正しく表すことができる。

【0079】

例えば、同じGC回数に関連づけられた2つのブロックB11、B12がGC対象ブロックとして選択され、これらブロックB11、B12の有効データがコピー先フリーブロックB21にコピーされたならば、このコピー先フリーブロックB21のGC回数は、ブロックB11、B12のGC回数（ここでは、0）に1を加えた値（ここでは1）に設定される。

【0080】

同様に、同じGC回数に関連づけられた3つのブロックB13、B14、B15がGC対象ブロックとして選択され、これらブロックB13、B14、B15の有効データがコピー先フリーブロックB22にコピーされたならば、このコピー先フリーブロックB22のGC回数は、ブロックB13、B14、B15のGC回数（ここでは、0）に1を加えた値（ここでは1）に設定される。

【0081】

同様に、同じGC回数に関連づけられた2つのブロックB16、B17がGC対象ブロックとして選択され、これらブロックB16、B17の有効データがコピー先フリーブロックB23にコピーされたならば、このコピー先フリーブロックB23のGC回数は、ブロックB16、B17のGC回数（ここでは、0）に1を加えた値（ここでは1）に設定される。

【0082】

時間が経過するに連れ、ブロックB21、B22、B23の各々のデータの一部は、その書き換えによって無効化されるかもしれない。これにより、ブロックB21、B22、B23の各々においては、有効データと無効データとが混在される場合がある。

【0083】

同じGC回数に関連づけられた2つのブロックB21、B22がGC対象ブロックとして選択され、これらブロックB21、B22の有効データがコピー先フリーブロックB31にコピーされたならば、このコピー先フリーブロックB31のGC回数は、ブロックB21、B22のGC回数（ここでは1）に1を加えた値（ここでは2）に設定される。

【0084】

このように、本実施形態では、ブロック毎に管理されるGC回数はそのブロック内のデータが過去のGC動作によってコピーされた回数を示す。このGC回数を正しく管理するために、GC対象ブロックのGC回数に1を加えた値が、コピー先フリーブロック内のデータに引き継がれる。

【0085】

図3は、GC回数管理リスト34の例を示す。

【0086】

管理すべきGC回数の上限値nが例えば10である場合、GC回数管理リスト34は、GC回数=0～GC回数=10にそれぞれ対応する11個のGC回数リストから構成されてもよい。

【0087】

GC回数=0のGC回数リストは、GC回数=0に関連づけられたブロックそれぞれのブロックID（例えば物理ブロックアドレス）のリストを示す。GC回数=1のGC回数

10

20

30

40

50

リストは、GC回数 = 1に関連づけられたブロックそれぞれのブロックID（例えば物理ブロックアドレス）のリストを示す。同様に、GC回数 = 10のGC回数リストは、GC回数 = 10に関連づけられたブロックそれぞれのブロックID（例えば物理ブロックアドレス）のリストを示す。各GC回数リストは、限定されないが、有効データと無効データとが混在するブロックだけを含んでもよい。

【0088】

図4は、コントローラ4によって実行されるGC対象ブロック選択動作を示す。

【0089】

GC対象ブロックを選択する処理においては、コントローラ4のガベージコレクション動作制御部22は、最初に、異なるGC回数にそれぞれ関連づけられた複数のブロック群（複数のGC回数リスト）から、GC対象のブロック群を選択してもよい。図4では、GC回数=5のブロック群（ブロックB2、ブロックB5、ブロックB11、ブロックB21）がGC対象のブロック群として選択され、さらに、このGC回数=5のブロック群から、幾つかのGC対象ブロックが選択される場合が例示されている。

【0090】

GC対象ブロックを選択する処理においては、例えば、まず、所定の条件に合致するブロックが最初のGC候補として選択されてもよい。所定の条件に合致するブロックは、アクティブブロック（ホスト2によって書き込まれたデータを含むブロック）の中で最も無効データ量が多いブロックであってもよい。他の実施形態では、所定の条件に合致するブロックは、アクティブブロックの中で最も古いブロックであってもよい。以下では、最も無効データ量が多いブロックが、最初のGC候補として選択される場合を想定する。

【0091】

最も無効データ量が多いブロックがブロックB5であるならば、コントローラ4は、ブロックB5を含むGC回数リスト（ここでは、GC回数 = 5のGC回数リスト）を特定し、このGC回数 = 5のGC回数リストによって示されるブロック群（ブロックB2、ブロックB5、ブロックB11、ブロックB21）をGC対象のブロック群として選択し、このGC対象のブロック群から、幾つかのGC対象ブロックを選択する。例えば、これらブロックB2、ブロックB5、ブロックB11、ブロックB21の中で無効データ量が多い上位幾つかのブロックがGC対象ブロックとして選択されてもよい。この場合、例えば、ブロックB5と、ブロックB2、ブロックB11、ブロックB21内の中で最も無効データ量が多い上位の一つ以上のブロックとが、GC対象ブロックとして選択されてもよい。

【0092】

図5は、コントローラ4によって実行されるGC動作を示す。

【0093】

コントローラ4は、全てのフリーブロックを含むフリーブロックプール（フリーブロックリスト）60を管理する。コントローラ4は、これらフリーブロックから一つのフリーブロックを選択する。コントローラ4は、選択されたフリーブロックを、コピー先フリーブロックB1000として割り当てる。コントローラ4は、同じGC回数を有するGC対象ブロック（ここでは、ブロックB2、B5、B11）からコピー先フリーブロックB1000に全ての有効データをコピーする。そして、コントローラ4は、ルックアップテーブル33を更新して有効データのLBAそれぞれをコピー先フリーブロックB1000の物理アドレスそれぞれにマッピングする。

【0094】

コピー先フリーブロックB1000のGC回数は、ブロックB2、B5、B11のGC回数（=5）+1に設定される。ブロックB1000は、GC回数 = 6のGC回数リストに追加される。ブロックB2、B5、B11は、有効データを含まないフリーブロックとなる。フリーブロックとなったブロックB2、B5、B11は、GC回数 = 5のGC回数リストから破棄される。

【0095】

図6は、SSD3に書き込まれる複数種のデータの例を示す。

10

20

30

40

50

【 0 0 9 6 】

図 6 では、互いに更新頻度の異なる 3 種類のデータ（データ A、データ B、データ C）が S S D 3 に書き込まれる場合が想定されている。S S D 3 のデータ記憶領域（L B A スペース）は、L B A グループ A、B、C に対応する 3 つのスペースを含む。

【 0 0 9 7 】

L B A グループ A に書き込まれるデータ A は更新頻度の低いデータであり、且つデータ A の量はデータ A、B、C の中で最も多い。つまり、L B A グループ A は最も大きい L B A 範囲を有する。

【 0 0 9 8 】

L B A グループ C に書き込まれるデータ C は、更新頻度の高いデータであり、且つデータ C の量はデータ A、B、C の中で最も少ない。つまり、L B A グループ C は最も小さい L B A 範囲を有する。

10

【 0 0 9 9 】

L B A グループ B に書き込まれるデータ B は、データ A とデータ C の中間の更新頻度を有するデータであり、且つデータ B の量はデータ A の量とデータ C の量の中間である。

【 0 1 0 0 】

S S D 3 の総ユーザ容量に対するデータ A の量の割合は、例えば、5 0 % であってもよい。S S D 3 の総ユーザ容量に対するデータ B の量の割合は、例えば、3 0 % であってもよい。S S D 3 の総ユーザ容量に対するデータ C の量の割合は、例えば、2 0 % であってもよい。

20

【 0 1 0 1 】

データ A の更新頻度つまり L B A グループ A へのライトの頻度は、例えば、2 0 % であってもよい。データ B の更新頻度つまり L B A グループ B へのライトの頻度は、例えば、3 0 % であってもよい。データ C の更新頻度つまり L B A グループ C へのライトの頻度は、例えば、5 0 % であってもよい。

【 0 1 0 2 】

この場合、例えば、S S D 3 がデータ A、データ B、データ C で満たされた後は、2 回のライトコマンドに 1 回の割合で、データ C（L B A グループ C）へのライトを要求するライトコマンドがホスト 2 から S S D 3 に発行され、また 5 回のライトコマンドに 1 回の割合で、データ A（L B A グループ A）へのライトを要求するライトコマンドがホスト 2 から S S D 3 に発行される。例えば、データ C は、2 回のライトコマンドに 1 回の割合（5 0 %）という高い頻度で更新される。

30

【 0 1 0 3 】

S S D 3 に書き込まれるデータが図 6 のようなデータ局所性を有する場合においては、図 6 の下部に示すように、各書き込み先ブロックにはデータ A、データ B、データ C が混在される。

【 0 1 0 4 】

一つの書き込み先ブロックにおいて、ブロックの容量に対するデータ C の量の割合は 5 0 %、ブロックの容量に対するデータ B の量の割合は 3 0 %、ブロックの容量に対するデータ A の量の割合は 2 0 % となる。

40

【 0 1 0 5 】

上述したように、データ C の量は、データ A、データ B よりも少なく、且つデータ C の更新頻度は、データ A、データ B よりも高いので、各ブロック内のデータ C のほとんどは速いタイミングで無効化される確率が高い。一方、データ A およびデータ B については、特にデータ A については、長い間、有効状態に維持される確率が高い。

【 0 1 0 6 】

データ C の更新（書き換え）によって無効データ量が増えたブロックそれぞれは、いずれ G C 対象ブロックとなり、これらブロックからコピー先フリーブロックに有効データがコピーされる。各 G C 対象ブロックにおいては、データ C の多くが無効化され且つデータ A、データ B の多くが有効データに維持されている確率が高い。このため、コピー先プロッ

50

クにおいては、GC対象ブロックに比べてデータAの量とデータBの量とが増え、代わりに、GC対象ブロックに比べてデータCの量が減る。

【0107】

本実施形態では、同じGC回数の幾つかのブロック内の有効データがコピー先フリーブロックにコピーされるので、GC回数の少ないブロック内の有効データとGC回数の多いブロック内の有効データとがGC動作によって同じコピー先フリーブロックにコピーされることはない。したがって、GC回数の多いブロックほど、そのブロックの容量に対するデータAの量の割合を増やすことができ、これによってデータA（Coldデータ）を、データC（ホットデータ）から分離することができる。

【0108】

図7は、GC回数と、各ブロック内のデータA、B、C間のデータ量の割合との関係の例を示す。

【0109】

GC回数 = 0の各ブロックにおいては、ブロックの容量に対するデータCの量の割合は50%、ブロックの容量に対するデータBの量の割合は30%、ブロックの容量に対するデータAの量の割合は20%である。

【0110】

ブロックの容量に対するデータCの量の割合は、1回または2回程度のGC動作によって速く低下される。GC回数が増えるにつれて、ブロックの容量に対するデータBの量の割合も徐々に低下される。

【0111】

上述したように、本実施形態では、GC回数の少ないブロック内の有効データとGC回数の多いブロック内の有効データとが同じコピー先フリーブロックにコピーされることはないので、データを含むブロックそれぞれを、(1)ほぼデータAのみを含むグループ（例えばGC回数7～10程度）、(2)データAとデータBとを含み、且つデータCをほとんど含まないグループ（例えばGC回数3～6程度）、(3)データAとデータBとデータCを含むグループ（例えばGC回数0～2程度）に分類できる。

【0112】

換言すれば、本実施形態では、同じGC回数のブロックについては、それらブロックに含まれるデータA、B、Cの量の割合を同じにすることができる。

【0113】

よって、同じGC回数の幾つかのブロック内の有効データを同じコピー先フリーブロックにコピーするという本実施形態の改良されたGC動作は、たとえSSD3に書かれるデータが高いデータ局所性を有する場合であっても、ほぼデータAのみを含むブロックのグループと、データAとデータBとを含み且つデータCをほとんど含まないブロックのグループと、データAとデータBとデータCを含むブロックのグループとを作ることができ、これによってHotデータとColdデータとを徐々に分離することができる。この結果、SSD3のライトアンプリフィケーションの増加を抑制することができる。

【0114】

図8のフローチャートは、コントローラ4によって実行されるGC動作の手順を示す。

【0115】

コントローラ4は、残りフリーブロックの数をチェックし（ステップS11）、残りフリーブロックの数が閾値 t_{h1} 以下であるか否かを判定する（ステップS12）。このチェックは、定期的に行われてもよい。例えば、新たなフリーブロックを書き込み先ブロックとして割り当てるべき時に残りフリーブロックの数をチェックしてもよい。

【0116】

残りフリーブロックの数が閾値 t_{h1} 以下であるならば（ステップS12のYES）、コントローラ4は、まず、全てのアクティブブロックから最初のGC候補を選択する。最初のGC候補は、最大無効データ量のブロックであってもよい。この場合、全てのアクティブブロックから最大無効データ量のブロックが最初のGC候補として選択される（ステ

10

20

30

40

50

ップ S 1 3)。コントローラ 4 は、G C 回数管理リスト 3 4 を参照して、最初の G C 候補 (ここでは、例えば、最大無効データ量のブロック) の G C 回数と同じ G C 回数に関連づけられたブロック群 (第 1 ブロック群) を選択し、さらに、この第 1 ブロック群から、幾つかの G C 対象ブロックを選択する (ステップ S 1 4)。ステップ S 1 4 では、最初の G C 候補 (例えば、最大無効データ量のブロック) が含まれている G C 回数リストによって示されるブロック群 (第 1 ブロック群) が選択され、そして第 1 ブロック群から、幾つかの G C 対象ブロックが選択される。この場合、最初の G C 候補 (例えば、最大無効データ量のブロック) と、この G C 回数リストに含まれる別の 1 以上のブロックとが、G C 対象ブロックとして選択されてもよい。

【 0 1 1 7 】

コントローラ 4 は、これら選択された G C 対象ブロック内の全ての有効データをコピー先フリーブロックにコピーする (ステップ S 1 5)。ステップ S 1 5 では、これら選択された G C 対象ブロック内の有効ページそれぞれから有効データがリードされ、リードされた有効データがコピー先フリーブロックの利用可能ページそれぞれに書き込まれる。ステップ S 1 5 では、さらに、コントローラ 4 は、ルックアップテーブル (L U T) 3 3 を更新して、コピーされた有効データの L B A をコピー先フリーブロックの物理アドレスに関連付けると共に、ページ管理テーブルを更新して、各 G C 対象ブロック内の元のページ (つまりこの L B A が関連付けられていた古いデータ) を無効化する。この場合、コントローラ 4 は、まず、ルックアップテーブル (L U T) 3 3 を参照することによって、コピーされた有効データが格納されている元のページの物理アドレスを取得してもよく、そして、ページ管理テーブルを更新して、この物理アドレスに対応する有効 / 無効フラグを無効を示す値に設定してもよい。

【 0 1 1 8 】

この後、コントローラ 4 は、これら選択された G C 対象ブロックの G C 回数 + 1、つまり第 1 ブロック群の G C 回数に 1 を加えた値を、コピー先フリーブロックの G C 回数として設定する (ステップ S 1 6)。

【 0 1 1 9 】

図 9 は、異なる G C 回数を有する 2 つのブロック群の有効データをマージする処理を含む G C 動作を示す。

【 0 1 2 0 】

例えば、最大無効データ量のブロックの G C 回数と同じ G C 回数に関連づけられたブロック群 (G C 対象ブロック群) に含まれる有効データの量が閾値よりも少ない場合、コントローラ 4 は、異なる G C 回数を有する 2 つのブロック群の有効データをマージする処理を実行する。この場合、コントローラ 4 は、G C 対象ブロック群の G C 回数と出来るだけ近い G C 回数を有する別の一つのブロック群を選択してもよい。

【 0 1 2 1 】

例えば、いま、最大無効データ量のブロックがブロック B 3 0 0 であり、ブロック B 3 0 0 の G C 回数が 1 0 である場合を想定する。この場合、コントローラ 4 は、G C 回数 = 1 0 の G C 回数管理リストに含まれるブロック群の総有効データ量をチェックする。例えば、G C 回数 = 1 0 の G C 回数管理リストに含まれるブロックがブロック B 3 0 0 のみである場合、あるいは G C 回数 = 1 0 の G C 回数管理リストに 2 つまたは 3 つ程度のブロックが含まれているがこれら各々の有効データ量が非常に少ない場合には、コントローラ 4 は、G C 回数 = 1 0 のブロック群と一緒に G C 動作が実行されるべきブロック群を選択する。

【 0 1 2 2 】

この場合、コントローラ 4 は、最大無効データ量のブロック B 3 0 0 の G C 回数よりも 1 回以上少ない G C 回数を有する全てのブロック群 (ここでは、G C 回数 9 のブロック群、G C 回数 8 のブロック群、G C 回数 7 のブロック群、... G C 回数 0 のブロック群) の中で、最大のガベージコレクション回数を有するブロック群を選択してもよい。

【 0 1 2 3 】

10

20

30

40

50

コントローラ 4 は、最初に GC 回数 = 9 の GC 回数管理リストを参照して、GC 回数 = 9 のブロックが存在するか否かを判定する。GC 回数 = 9 のブロックが存在しないならば、コントローラ 4 は、GC 回数 = 8 の GC 回数管理リストを参照して、GC 回数 = 8 のブロックが存在するか否かを判定する。

【0124】

GC 回数 = 9 のブロックが存在せず、GC 回数 = 8 のブロックが存在するならば、コントローラ 4 は、GC 回数 = 8 のブロック群（例えば、ブロック B 4 1、B 4 2、B 4 3）を、選択する。そして、コントローラ 4 は、ブロック B 3 0 0 の有効データと GC 回数 = 8 のブロック群の有効データとをコピー先フリーブロックにコピーする。この場合、ブロック B 4 1、B 4 2、B 4 3 の全ての有効データが必ずしも利用される必要は無く、ブロック B 4 1、B 4 2、B 4 3 内の少なくとも一つのブロック内の有効データが利用されればよい。

10

【0125】

図 10 のフローチャートは、異なる GC 回数を有する 2 つのブロック群の有効データをマージする処理を含む GC 動作の手順を示す。

【0126】

コントローラ 4 は、残りフリーブロックの数をチェックし（ステップ S 2 1）、残りフリーブロックの数が閾値 $t_h 1$ 以下であるか否かを判定する（ステップ S 2 2）。上述したように、このチェックは、定期的に行われてもよい。

【0127】

残りフリーブロックの数が閾値 $t_h 1$ 以下であるならば（ステップ S 2 2 の YES）、コントローラ 4 は、まず、全てのアクティブブロックから最初の GC 候補を選択する。最初の GC 候補は、最大無効データ量のブロックであってもよい。この場合、全てのアクティブブロックから最大無効データ量のブロックが最初の GC 候補として選択される（ステップ S 2 3）。コントローラ 4 は、GC 回数管理リスト 3 4 を参照して、最初の GC 候補（ここでは、例えば、最大無効データ量のブロック）の GC 回数と同じ GC 回数に関連づけられたブロック群（第 1 ブロック群）を選択し、このブロック群（第 1 ブロック群）の有効データの総量が閾値 $t_h 2$ 以下であるか否かを判定する（ステップ S 2 4）。

20

【0128】

閾値 $t_h 2$ の値は、固定であっても良いし、必要に応じて変更できる値であっても良い。閾値 $t_h 2$ の値が大きいほど、上述のマージ処理の実行が許可されやすくなる。

30

【0129】

例えば、閾値 $t_h 2$ は、SSD 3 内の一つのブロックの容量を示す値に予め設定されている。これにより、最初の GC 候補の GC 回数と同じ GC 回数に関連づけられたブロック群のみで GC 動作が実行できない場合にのみ、マージ処理の実行を許可することができる。あるいは、閾値 $t_h 2$ は、SSD 3 内の一つのブロックの容量の整数倍、例えば 2 倍の値に設定されている。これも良い。

【0130】

この第 1 ブロック群の有効データの総量が閾値 $t_h 2$ 以下でないならば（ステップ S 2 4 の NO）、コントローラ 4 は、この第 1 ブロック群から、幾つかの GC 対象ブロックを選択する（ステップ S 2 5）。ステップ S 2 5 では、最初の GC 候補（例えば、最大無効データ量のブロック）が含まれている GC 回数リストによって示される第 1 ブロック群から、これら GC 対象ブロックが選択される。この場合、最初の GC 候補（例えば、最大無効データ量のブロック）と、この GC 回数リストに含まれる別のブロックとが、GC 対象ブロックとして選択されてもよい。

40

【0131】

ステップ S 2 5 では、コントローラ 4 は、これら選択された GC 対象ブロック内の全ての有効データをコピー先フリーブロックにコピーする。ステップ S 2 5 では、さらに、コントローラ 4 は、ルックアップテーブル（LUT）3 3 を更新して、コピーされた有効データの LBA をコピー先フリーブロックの物理アドレスに関連付けると共に、各 GC 対象

50

ブロック内の元のページを無効化する。

【0132】

この後、コントローラ4は、これら選択されたGC対象ブロックのGC回数+1を、つまり第1ブロック群のGC回数に1を加えた値を、コピー先フリーブロックのGC回数として設定する(ステップS26)。

【0133】

一方、第1ブロック群の有効データの総量が閾値 t_{h2} 以下であるならば(ステップS24のYES)、コントローラ4は、この第1ブロック群のGC回数よりも1回以上少ないGC回数に関連づけられた全てのブロック群の中で、最大GC回数に関連づけられたブロック群(第2ブロック群)を選択する(ステップS27)。

10

【0134】

コントローラ4は、第1ブロック群の有効データと第2ブロック群の有効データとをコピー先フリーブロックにコピーする(ステップS28)。ステップS28では、さらに、コントローラ4は、ルックアップテーブル(LUT)33を更新して、コピーされた有効データのLBAをコピー先フリーブロックの物理アドレスに関連付けると共に、各GC対象ブロック内の元のページを無効化する。

【0135】

コントローラ4は、第2ブロック群のGC回数+1をコピー先フリーブロックのGC回数として設定するか、あるいは第1ブロック群のGC回数+1をコピー先フリーブロックのGC回数として設定する(ステップS29)。あるいは、第1ブロック群内のGC対象ブロックの数よりも第2ブロック群内のGC対象ブロックの数が多い場合には、第2ブロック群のGC回数+1をコピー先フリーブロックのGC回数として設定してもよく、第1ブロック群内のGC対象ブロックの数第2ブロック群内のGC対象ブロックの数よりも多い場合には、第1ブロック群のGC回数+1をコピー先フリーブロックのGC回数として設定してもよい。

20

【0136】

図11は、マージ処理を特定のGC回数以上のブロック群に対してのみ許可する動作を示す。

【0137】

GC回数の多いブロック内に含まれている有効データは、更新頻度の低いデータ(データA)である可能性が高い。しかし、データAも20%の割合で書き替えられるので、GC回数の多いブロック、例えばGC回数=10のブロック、についても、その無効データ量が多くなる場合がある。GC回数の多いブロック内の有効データは、これまで一度も更新(書き替え)されたことのないデータ、つまり、長い間、有効状態に維持されているデータである。このため、この有効データは、これからも更新されない確率が高い。

30

【0138】

一方、GC回数の少ないブロックにおいては、データBまたはデータCが含まれている可能性が高い。このようなブロックについては、そのブロックのGC動作をすぐ実行せずとも、時管理経過に伴ってブロック内の全てのデータが無効化される可能性がある。

40

【0139】

したがって、マージ処理を許可するブロック群をマージ許可閾値 t_{h3} 以上のGC回数を有するブロック群に対してのみに許可することにより、無駄なコピーの発生を防ぐことができ、GCの効率を高めることができる。

【0140】

図11では、マージ許可閾値 t_{h3} がGC回数=8に設定されている場合が例示されている。

【0141】

この場合、最初のGC候補のGC回数と同じGC回数に関連づけられたブロック群(第1ブロック群)のGC回数が8以上であるならば、第1ブロック群と他のブロック群とのマージ処理が許可される。

50

【 0 1 4 2 】

例えば、G C 回数 = 1 0 のブロック群と他のブロック群とのマージ処理、および G C 回数 = 9 のブロック群と他のブロック群とのマージ処理が、許可される。一方、例えば、G C 回数 = 7 のブロック群と他のブロック群とのマージ処理は禁止される。

【 0 1 4 3 】

図 1 2 のフローチャートは、マージ処理を特定の G C 回数以上のブロック群に対してのみ許可する動作を含む G C 動作の手順を示す。

【 0 1 4 4 】

この図 1 2 のフローチャートに示される G C 動作においては、図 1 0 で説明した処理に加え、ステップ S 3 0 ~ S 3 3 の処理が追加されている。以下では、ステップ S 3 0 ~ S 3 3 の処理を主に説明する。

10

【 0 1 4 5 】

第 1 ブロック群の有効データの総量が閾値 t_{h2} 以下であるならば (ステップ S 2 4 の Y E S)、コントローラ 4 の処理は、ステップ S 3 0 に進む。ステップ S 3 0 において、コントローラ 4 は、第 1 ブロック群の G C 回数がマージ許可閾値 t_{h3} 以上であるか否かを判定する。

【 0 1 4 6 】

第 1 ブロック群の G C 回数がマージ許可閾値 t_{h3} 以上であるならば (ステップ S 3 0 の Y E S)、コントローラ 4 は、図 1 0 で説明したステップ S 2 7 ~ S 2 9 のマージ処理を実行する。

20

【 0 1 4 7 】

一方、第 1 ブロック群の G C 回数 (最初の G C 候補のブロックの G C 回数) がマージ許可閾値 t_{h3} よりも少ないならば (ステップ S 3 0 の N O)、コントローラ 4 は、ステップ S 2 7 ~ S 2 9 のマージ処理の実行を禁止し、代わりに、ステップ S 3 1 ~ S 3 3 の処理を実行する。

【 0 1 4 8 】

ステップ S 3 1 において、コントローラ 4 は、第 1 ブロック群とは異なる別のブロック群を G C 対象ブロック群として選択する。例えば、コントローラ 4 は、最初の G C 候補のブロックの次に無効データ量が多いブロックを新たな G C 候補として選択し、この新たな G C 候補が含まれている G C 回数リストによって示されるブロック群を G C 対象ブロック群として選択してもよい。

30

【 0 1 4 9 】

次いで、コントローラ 4 は、選択された G C 対象ブロック群の有効データをコピー先フリーブロックにコピーし (ステップ S 3 2)、コピー先フリーブロックの G C 回数を、G C 対象ブロック群の G C 回数に 1 を加えた値に設定する (ステップ S 3 3)。

【 0 1 5 0 】

最初の G C 候補のブロックが、マージ許可閾値 t_{h3} よりも少ない G C 回数に関連付けられている場合には、この最初の G C 候補のブロックは、頻繁に更新されるデータを含んでいる可能性が高い。このため、コントローラ 4 は、最初の G C 候補のブロックに対する G C を実行せずに、このブロックの有効データが全て無効化されるまで待っても良い。

40

【 0 1 5 1 】

次に、図 1 3 ~ 図 2 2 を参照して、「L B A ベースの更新頻度通知機能」の詳細を説明する。

【 0 1 5 2 】

図 1 3 は、フリーブロックをホスト 2 からのデータの書き込み用に順次割り当てる動作を示す。

【 0 1 5 3 】

コントローラ 4 は、フリーブロックリスト 6 0 によって示されるフリーブロックの一つを書き込み先ブロック 6 2 として割り当てる。この場合、コントローラ 4 は、ブロック使用順序管理リスト 3 5 を更新して、書き込み先ブロック 6 2 として最初に割り当てられた

50

ブロックの割り当て番号（シーケンシャル番号）を 1 に設定する。ブロック使用順序管理リスト 3 5 は、図 1 4 に示されているように、ブロックアドレスそれぞれに対応する割り当て番号（シーケンシャル番号）を保持する。これら割り当て番号は、書き込み先ブロック 6 2 に割り当てられたブロックの順序関係を示す。つまり、コントローラ 4 は、書き込み対象ブロックとして割り当てられたブロックそれぞれに対してその割り当て順序を示す割り当て番号を付与し、これら割り当て番号をブロック使用順序管理リスト 3 5 を使用して管理する。

【 0 1 5 4 】

コントローラ 4 は、ホスト 2 から受信されるライトデータをライトバッファ 3 1 に書き込む。この後、コントローラ 4 は、ルックアップテーブル（LUT）3 3 を更新しながら、ライトバッファ 3 1 内のライトデータを書き込み先ブロック 6 2 の先頭ページから最終ページに向けて順次ライトする。

10

【 0 1 5 5 】

書き込み先ブロック 6 2 に利用可能ページが無くなったならば、コントローラ 4 は、書き込み先ブロック 6 2 をアクティブブロックリスト 6 1 に移動し、フリーブロックリスト 6 0 のフリーブロックを新たな書き込み先ブロック 6 2 として割り当てる。この場合、コントローラ 4 は、ブロック使用順序管理リスト 3 5 を更新して、この新たな書き込み先ブロック 6 2 として割り当てられたこのブロックの割り当て番号（シーケンシャル番号）を 2 に設定する。

【 0 1 5 6 】

アクティブブロックリスト 6 1 内の何れかのブロックの全てのデータがその更新によって無効化されたならば、このブロックはフリーブロックリスト 6 0 に移動される。

20

【 0 1 5 7 】

フリーブロックリスト 6 0 内のフリーブロックの数が閾値 $t h 1$ 以下に低下したならば、フリーブロックを作り出す上述の GC 動作が実行される。

【 0 1 5 8 】

図 1 5 は、同じ L B A へのライトが要求された時に実行される累積データ書き込み量算出動作を示す。

【 0 1 5 9 】

コントローラ 4 は、ホスト 2 からある L B A を含むライトコマンドを受信した際に、この L B A への前回のライトからの累積データ書き込み量を、このライトコマンドに対する応答としてホスト 2 に通知する。累積データ書き込み量は、受信されたライトコマンドの L B A と同じ L B A への前回のライトからライトコマンドの L B A への今回のライトまでの間にホスト 2 によって N A N D メモリ 5 に書き込まれたデータの総量を示す。

30

【 0 1 6 0 】

累積データ書き込み量は、例えば、次の値から算出することができる。

【 0 1 6 1 】

- (1) ブロック当たりの容量
- (2) ブロック内に含まれるページの数
- (3) 同じ L B A への前回のライトによってデータが書き込まれた N A N D メモリ 5 内の第 1 物理記憶位置（旧物理アドレス）
- (4) 今回のライトによってデータが書き込まれるべき N A N D メモリ 5 内の第 2 物理記憶位置（新物理アドレス）
- (5) 第 1 物理記憶位置（旧物理アドレス）を含むブロックの割り当てから第 2 物理記憶位置（新物理アドレス）を含むブロックの割り当てまでの間にホスト 2 からのデータの書き込みのために割り当てられたブロックの数

40

(1) ~ (4) の値は、S S D 3 内の通常の管理情報であり、累積データ書き込み量の算出のために専用に用意されたものではない。例えば、コントローラ 4 は、ルックアップテーブル（LUT）3 3 を参照することによって、受信されたライトコマンド内の L B A にマッピングされている物理アドレスを第 1 物理記憶位置として容易に取得することがで

50

きる。

【0162】

(5)の「ブロックの数」は、例えば、第1物理記憶位置を含むブロックに付与された割り当て番号と第2物理記憶位置を含むブロックに付与された割り当て番号とから容易に算出することができる。

【0163】

割り当て番号(シーケンシャル番号)は、図14のブロック使用順序管理リスト35によって管理されている。これら割り当て番号(シーケンシャル番号)の管理単位は、ブロック単位であるので、これら割り当て番号を保持するために必要な容量は少なく済む。したがって、累積データ書き込み量は、その算出のための専用の管理情報をほとんど使用すること無く、低コストで取得することができる。

10

【0164】

図15では、LBA10を含むライトコマンドが受信された時に実行される累積データ書き込み量算出動作を示している。

【0165】

ここでは、LBA10への前回のライトによってデータがブロックB51のページPxに既に書き込まれており、且つLBA10への今回のライトによってデータが現在の書き込み先ブロックB62のページPyに書き込まれるべき場合が想定されている。もしブロックB51の割り当て番号が10で、ブロックB51の割り当て番号が13であれば、ブロック51とブロックB62との間に2つの書き込み先ブロック(例えばブロックB52、B61)が割り当てられていたことが分かる。

20

【0166】

累積データ書き込み量は、 $d_1 + d_2 + d_2 + d_3$ で与えられる。

【0167】

ここで、 d_1 は、ページPxに後続するブロックB51内のページの数、またはこれらページの数に対応する容量を示す。 d_2 は、一つのブロック内のページの数、または一つのブロックの容量を示す。 d_3 は、ページPyに先行するブロックB62内のページの数、またはこのページの数に対応する容量を示す。

【0168】

LBA10を含む前回のライトコマンドの受信からLBA10を含む今回のライトコマンドの受信までの間にホスト2から受信されるライトコマンドの数が多いほど、累積データ書き込み量は増加する。したがって、上述の累積データ書き込み量は、LBA10によって指定されるデータの更新頻度、つまりLBA10へのライトの頻度を表すことができる。

30

【0169】

ライトコマンドの受信時に、コントローラ4は、以下の手順で累積データ書き込み量を取得(算出)してもよい。

【0170】

まず、コントローラ4は、ルックアップテーブル(LUT)33を参照してライトコマンドに含まれるLBA(ここではLBA10)にマッピングされている旧物理アドレス(ここではPA1)を取得する。そして、コントローラ4は、ブロック使用順序管理リスト35を参照して、旧物理アドレスによって指定されるブロックの割り当て番号(ここでは10)と、新物理アドレス(ここではPA2)によって指定されるブロックの割り当て番号(ここでは13)とを取得する。コントローラ4は、ブロック内に含まれるページの数と旧物理アドレス(PA1)とから d_1 を求め、ブロック内に含まれるページの数と新物理アドレス(PA2)とから d_3 を求める。さらに、コントローラ4は、割り当て番号(13)と割り当て番号(10)との間の差分から、旧物理アドレスによって指定されるブロックの割り当てから新物理アドレスによって指定されるブロックの割り当てまでの間に、書き込み先ブロックとして割り当てられたブロックの総数(ここでは、2)を求める。これにより、累積データ書き込み量($= d_1 + d_2 + d_2 + d_3$)を取得(算出)するこ

40

50

とができる。

【0171】

図16は、累積データ書き込み量応答処理の処理シーケンスを示す。

【0172】

ここでは、この処理シーケンスが、ライトコマンドとライトデータとが分割されているNCQ(Native Command Queuing)システムに適用される場合を想定する。

【0173】

ホスト2は、あるLBA(=LBAx)を示す開始LBAを含むライトコマンドをSSD3に送出する。このライトコマンドの受信に回答して、SSD3のコントローラ4は、LBAxへの前回のライトからLBAxへの今回のライトまでの累積データ書き込み量を算出し(ステップS41)、算出された累積データ書き込み量を含むコマンド許可応答をホスト2に送信する。コマンド許可応答は、受信されたライトコマンドに対するアクノリッジ(ライトコマンドの実行許可)を示す許可応答である。SSD3からホスト2に許可応答が送信されることにより、このライトコマンドによって指定されるライトデータの転送が開始される。許可応答は、実行を許可すべきライトコマンドを識別する値を含んでも良い。累積データ書き込み量は、例えば、バイトで表されてもよいし、論理ブロック(論理セクタ)の数によって表されてもよい。

10

【0174】

コマンド許可応答の受信に回答して、ホスト2は、ライトデータをSSD3に送出する。SSD3のコントローラ4は、ライトデータをライトバッファ31に書き込み、ライトバッファ31のライトデータを書き込み先ブロックに書き込み(ステップS42)、コマンド完了の応答(レスポンス)をホスト2に送信する。なお、ライトデータをライトバッファ31に書き込んだ時点でコマンド完了のレスポンスをホスト2に送信してもよい。

20

【0175】

ホスト2は、SSD3から受信されるコマンド許可応答に含まれる累積データ書き込み量に基づいて、LBAxのデータの実際の更新頻度(LBAxへのライトの頻度)を把握することができる。

【0176】

もしLBAxのデータの実際の更新頻度が、ホスト2によって予期されていたLBAxのデータの更新頻度と異なるならば、例えば、LBAxのデータの実際の更新頻度がホスト2によって予期されていたLBAxのデータの更新頻度よりも高いならば、ホスト2は、必要に応じて、送出したライトコマンドをアポートするためのアポートコマンドをSSD3に送出してもよい。この場合、ライトコマンドによって指定されたデータの書き込みは実行されない。

30

【0177】

図17のフローチャートは、コントローラ4によって実行される累積データ書き込み量応答処理の手順を示す。

【0178】

コントローラ4は、LBAxを開始LBAとして含むライトコマンドをホスト2から受信する(ステップS51)。コントローラ4は、LBAxにマッピングされている旧物理アドレスと、LBAxにマッピングされるべき新物理アドレスと、旧物理アドレスによって指定される物理記憶位置を含むブロックに付与された割り当て番号と、新物理アドレスによって指定される物理記憶位置を含むブロック(現在の書き込み先ブロック)に付与された割り当て番号、等に基づいて、LBAxへの前回の書き込みからLBAxへの今回の書き込みまでの累積データ書き込み量を算出する(ステップS52)。コントローラ4は、累積データ書き込み量を含む許可応答をホスト2へ返す(ステップS53)。

40

【0179】

コントローラ4は、このライトコマンドに対応するライトデータまたはこのライトコマンドをアポートするためのアポートコマンドのいずれがホスト2から受信されるかを判定

50

する（ステップ S 5 4）。

【 0 1 8 0 】

もしライトデータが受信されたならば、コントローラ 4 は、ステップ S 5 5 に進む。ステップ S 5 5 では、コントローラ 4 は、このライトデータをライトバッファ 3 1 に書き込み、ライトバッファ 3 1 内のライトデータを現在の書き込み先ブロックに書き込み、ルックアップテーブル（LUT）3 3 を更新して LBA x に新物理アドレスをマッピングし、そしてページ管理テーブルを更新して旧物理アドレス（旧データ）を無効化する。

【 0 1 8 1 】

この後、コントローラ 4 は、コマンド完了のレスポンスをホスト 2 へ返す（ステップ S 5 6）。

10

【 0 1 8 2 】

なお、上述したように、ライトデータをライトバッファ 3 1 に書き込んだ時点でコマンド完了のレスポンスをホスト 2 に送信してもよい。

【 0 1 8 3 】

一方、アポートコマンドが受信されたならば、コントローラ 4 は、このライトコマンドを破棄する（ステップ S 5 7）。

【 0 1 8 4 】

図 1 8 は、累積データ書き込み量応答処理の別の処理シーケンスを示す。

【 0 1 8 5 】

ホスト 2 は、ある LBA（= LBA x）を開始 LBA として含むライトコマンドを SSD 3 に送出する。このライトコマンドの受信に回答して、SSD 3 のコントローラ 4 は、コマンド許可応答をホスト 2 に送信する。コマンド許可応答の受信に回答して、ホスト 2 は、ライトデータを SSD 3 に送出する。ライトデータはライトバッファ 3 1 に書き込まれる。SSD 3 のコントローラ 4 は、累積データ書き込み量を算出する（ステップ S 5 8）。累積データ書き込み量を算出する処理は、ライトコマンドの受信に回答して開始してもよい。

20

【 0 1 8 6 】

この後、コントローラ 4 は、書き込み先ブロックへのライトデータの書き込みを実行し（ステップ S 5 9）、算出された累積データ書き込み量を含む、コマンド完了のレスポンスをホスト 2 に送信する。

30

【 0 1 8 7 】

なお、上述したように、ライトデータがライトバッファ 3 1 に書き込まれた時点で、累積データ書き込み量を含むコマンド完了のレスポンスをホスト 2 に送信してもよい。

【 0 1 8 8 】

図 1 9 のフローチャートは、累積データ書き込み量応答処理の別の手順を示す。

【 0 1 8 9 】

コントローラ 4 は、LBA x を開始 LBA として含むライトコマンドをホスト 2 から受信する（ステップ S 6 1）。コントローラ 4 は、許可応答をホスト 2 へ返す（ステップ S 6 2）。コントローラ 4 は、ライトデータをホスト 2 から受信する（ステップ S 6 3）。ライトデータはライトバッファ 3 1 に書き込まれる。

40

【 0 1 9 0 】

コントローラ 4 は、LBA x にマッピングされている旧物理アドレスと、LBA x にマッピングされるべき新物理アドレスと、旧物理アドレスによって指定される物理記憶位置を含むブロックに付与された割り当て番号と、新物理アドレスによって指定される物理記憶位置を含むブロック（現在の書き込み先ブロック）に付与された割り当て番号等に基づいて、LBA x への前回の書き込みから LBA x への今回の書き込みまでの累積データ書き込み量を算出する（ステップ S 6 4）。コントローラ 4 は、ステップ S 6 5 に進む。

【 0 1 9 1 】

ステップ S 6 5 では、コントローラ 4 は、ライトバッファ 3 1 内のライトデータを現在の書き込み先ブロックに書き込み、ルックアップテーブル（LUT）3 3 を更新して LBA

50

A x に新物理アドレスをマッピングし、そしてページ管理テーブルを更新して旧物理アドレス（旧データ）を無効化する。

【0192】

この後、コントローラ4は、累積データ書き込み量を含む、コマンド完了のレスポンスをホスト2へ返す（ステップS66）。

【0193】

なお、上述したように、ライトデータをライトバッファ31に書き込んだ時点でコマンド完了のレスポンスをホスト2に送信してもよい。

【0194】

次に、図20～図23を参照して、累積データ書き込み量の代わりに、同じLBAへの前回のライトからの時間経過値をホスト2に通知する処理について説明する。

10

【0195】

この時間経過値は同じLBAへの前回のライトからの時間経過に関する情報であり、時間経過値の例は、同じLBAへの前回のライトの時刻であってもよいし、同じLBAへの前回のライトの時刻とこの同じLBAへの今回のライトの時刻との間の時間間隔であってもよい。

【0196】

図20は、例えば4Kバイトのような所定の管理単位で、LBAと、物理アドレスと、前回ライトされた時刻との対応関係を管理するように構成されたルックアップテーブル（LUT）33の例を示す。

20

【0197】

ルックアップテーブル（LUT）33は、LBA毎に物理アドレス記憶領域33Aと時刻記憶領域33Bとを含む。各時刻記憶領域33Bは、対応するLBAへのライトが発生した時刻を示す値、つまり対応するLBAのデータがライトされた時刻を示す値、を保持するために使用される。各時刻記憶領域33Bに保持される時刻は、例えば、時分秒であってもよい。

【0198】

あるLBAを含むライトコマンドが受信された時、コントローラ4は、このLBAに対応する物理アドレス領域33Aに物理アドレスを登録すると共に、このLBAに対応する時刻領域33Bに、ライトコマンドによって指定されるデータ（ライトデータ）がライトされた時刻を登録する。物理アドレスは、ライトコマンドによって指定されたデータが書き込まれた物理記憶位置の物理アドレスを示す。ライトされた時刻は、ライトコマンドが受信された時刻であってもよいし、ライトコマンドによって指定されたデータがライトバッファ31に書き込まれた時刻であってもよいし、ライトコマンドによって指定されたデータがNANDメモリ5の書き込み先ブロックにライトされた時刻であってもよい。

30

【0199】

図21のフローチャートは、コントローラ4によって実行される時間経過応答処理の手順を示す。

【0200】

ここでは、時間経過値を含むコマンド許可応答をホスト2に送信する場合を想定する。

40

【0201】

コントローラ4は、LBAxを開始LBAとして含むライトコマンドをホスト2から受信する（ステップS71）。コントローラ4は、ルックアップテーブル（LUT）33を参照して、LBAxへの前回のライトの時刻、つまりLBAxを含む前回のライトコマンドによってデータがライトされた時刻を、取得する（ステップS72）。コントローラ4は、LBAxへの前回のライトの時刻を示す時間経過値を含む許可応答をホスト2へ返す（ステップS73）。上述したように、時間経過値は、LBAxへの前回のライトの時刻とLBAxの今回のライトの時刻との間の時間間隔、つまり現在時刻（LBAxへの今回のライトの時刻）からLBAxへの前回のライトの時刻を引いた値であってもよい。

【0202】

50

コントローラ 4 は、このライトコマンドに対応するライトデータまたはこのライトコマンドをアポートするためのアポートコマンドのどちらがホスト 2 から受信されるかを判定する（ステップ S 7 4）。

【0203】

ライトデータが受信されたならば、コントローラ 4 は、ステップ S 7 5 に進む。ステップ S 7 5 では、コントローラ 4 は、このライトデータをライトバッファ 3 1 に書き込み、ライトバッファ 3 1 内のライトデータを現在の書き込み先ブロックに書き込み、ルックアップテーブル（LUT）3 3 を更新して L B A x に新物理アドレスと新ライト時刻とをマッピングし、そしてページ管理テーブルを更新して旧物理アドレス（旧データ）を無効化する。

10

【0204】

この後、コントローラ 4 は、コマンド完了のレスポンスをホスト 2 へ返す（ステップ S 7 6）。

【0205】

なお、上述したように、ライトデータをライトバッファ 3 1 に書き込んだ時点でコマンド完了のレスポンスをホスト 2 に送信してもよい。

【0206】

一方、アポートコマンドが受信されたならば、コントローラ 4 は、このライトコマンドを破棄する（ステップ S 7 7）。

【0207】

図 2 1 のフローチャートでは、時間経過値を含むコマンド許可応答をホスト 2 に送信する場合を説明したが、時間経過値を含むコマンド完了のレスポンスをホスト 2 に送信してもよい。時間経過値を含むコマンド完了のレスポンスの送信は、図 1 8、図 1 9 と同様の手順によって実行することができる。

20

【0208】

図 2 2 のフローチャートは、SSD 3 から通知される累積データ書き込み量 / 時間経過値に基づいてホスト 2 によって実行される処理の手順を示す。

【0209】

ホスト 2 は、SSD 3 から通知される累積データ書き込み量 / 時間経過値に基づいて、データを更新頻度の異なる複数種のデータグループに分類してもよい。例えば、ホスト 2 のファイルシステム 4 3 がデータ管理部を含み、このデータ管理部が、データを複数種のデータグループに分類して、データを頻繁に更新されるデータグループ（Hot データ）と頻度には更新されないデータグループ（Cold データ）とに分離してもよい。SSD 3 に書き込んだデータの更新頻度がある閾値以上であるならば、データ管理部は、このデータが Hot データであると認識することができる。

30

【0210】

データ管理部は、同じ SSD 内における L B A 範囲それぞれの更新頻度をできるだけ同じ範囲の頻度に揃えるために、Hot データであると認識されたデータを SSD 3 から別のストレージデバイスに移動しても良い。

【0211】

あるいは、もし SSD 3 が高い耐久性を有する高価格 SSD として実現されているならば、Hot データを SSD 3 内に残し、Cold データを SSD 3 から別のストレージデバイスに移動しても良い。高い耐久性を有する高価格 SSD の例は、メモリセル当たり 1 ビットの情報を格納する S L C - S S D を含む。

40

【0212】

SSD の耐久性を示す指標の一つに、DWPD（Drive Write Per Day）がある。例えば、DWPD = 1 0 は、1 T バイトの総容量を有する SSD に関しては、1 日当たり 1 0 T バイト（= 1 0 × 1 T バイト）のデータのライトを 5 年間に渡って毎日実行することができることを意味する。

【0213】

50

以下では、前者のための処理の手順の例を説明する。

【0214】

ホスト2は、LBAxを含むライトコマンドをSSD3に送信し(ステップS81)、累積データ書き込み量または時間経過値を含む応答(許可応答、コマンド完了レスポンス)をSSD3から受信する(ステップS82)。

【0215】

ホスト2は、累積データ書き込み量または時間経過値に基づき、LBAxのデータの更新頻度(LBAxへのライトの頻度)が所定の上限頻度(閾値th4)以上であるか否かを判定する(ステップS83)。例えば、SSD3から累積データ書き込み量が通知されるケースにおいては、ホスト2は、累積データ書き込み量が閾値th4によって示される閾データ量以上であるかを判定してもよい。SSD3から時間経過値(同じLBAへの前回のライトの時刻)が通知されるケースにおいては、ホスト2は、現在時刻から前回のライトの時刻を引くことによって時間間隔を算出し、この時間間隔が、閾値th4によって示される閾時間間隔以上であるか否かを判定してもよい。あるいは、ホスト2は、累積データ書き込み量または時間経過値を、何回のライトアクセスに1回の割合でLBAxへのライトが発生するかを示す割合[パーセント]に換算し、この時間間隔が、閾値th4によって示される閾時間間隔以上であるか否かを判定してもよい。

10

【0216】

LBAxのデータの更新頻度(LBAxへのライトの頻度)が閾値th4以上であるならば(ステップS83のYES)、ホスト2は、LBAxのデータを高更新頻度データグループ(Hotデータ)に分類し(ステップS84)、LBAxのデータをSSD3から他のストレージデバイスに移動する(ステップS85)。

20

【0217】

ステップS84においては、もし累積データ書き込み量または時間経過値が含まれるレスポンスがライトコマンドに対する許可応答であったならば、ホスト2は、ライトコマンドをアボートする処理を実行してもよい。

【0218】

図23は、ホスト2として機能する情報処理装置のハードウェア構成例を示す。

【0219】

この情報処理装置は、サーバコンピュータ、またはパーソナルコンピュータとして実現される。この情報処理装置は、プロセッサ(CPU)101、メインメモリ102、BIOS-ROM103、ネットワークコントローラ105、周辺インタフェースコントローラ106、コントローラ107、およびエンベデッドコントローラ(EC)108等を含む。

30

【0220】

プロセッサ101は、この情報処理装置の各コンポーネントの動作を制御するように構成されたCPUである。このプロセッサ101は、複数のSSD3のいずれか1つからメインメモリ102にロードされる様々なプログラムを実行する。メインメモリ102は、DRAMのようなランダムアクセスメモリから構成される。プロセッサ101によって実行されるプログラムは、上述のアプリケーションソフトウェアレイヤ41、OS42およびファイルシステム43を含む。

40

【0221】

また、プロセッサ101は、不揮発性メモリであるBIOS-ROM103に格納された基本入出力システム(BIOS)も実行する。BIOSはハードウェア制御のためのシステムプログラムである。

【0222】

ネットワークコントローラ105は、有線LANコントローラ、無線LANコントローラのような通信デバイスである。周辺インタフェースコントローラ106は、USBデバイスのような周辺デバイスとの通信を実行するように構成されている。

【0223】

50

コントローラ107は、複数のコネクタ107Aにそれぞれ接続されるデバイスとの通信を実行するように構成されている。本実施形態では、複数のSSD3が複数のコネクタ107Aにそれぞれ接続される。コントローラ107は、SAS expander、PCIe Switch、PCIe expander、フラッシュアレイコントローラ、またはRAIDコントローラ等である。

【0224】

EC108は、情報処理装置の電力管理を実行するように構成されたシステムコントローラとして機能する。EC108は、ユーザによる電源スイッチの操作に応じて情報処理装置をパワーオンおよびパワーオフする。EC108はワンチップマイクロコントローラのような処理回路として実現されている。EC108は、キーボード(KB)などの入力デバイスを制御するキーボードコントローラを内蔵していてもよい。

10

【0225】

図22で説明した処理は、ファイルシステム43の制御の下、プロセッサ101によって実行される。

【0226】

図24は、複数のSSD3とホスト2とを含む情報処理装置の構成例を示す。

【0227】

この情報処理装置は、ラックに収容可能な薄い箱形の筐体201を備える。多数のSSD3は筐体201内に配置されても良い。この場合、各SSD3は筐体201の前面201Aに設けられたスロットに取り外し可能に挿入されてもよい。

20

【0228】

システムボード(マザーボード)202は筐体201内に配置される。システムボード(マザーボード)202上においては、CPU101、メモリ102、ネットワークコントローラ105、コントローラ107を含む様々な電子部品が実装されている。これら電子部品がホスト2として機能する。

【0229】

以上説明したように、本実施形態の「ブロック内のデータのGC回数を考慮したGC機能」によれば、ホスト2によって書き込まれたデータを含むブロック毎に、当該ブロック内のデータがガベージコレクション(GC)動作によってコピーされた回数を示すGC回数が管理され、且つ同じGC回数に関連づけられた複数のブロック(第1ブロック)が、ガベージコレクション(GC)動作の対象ブロックとして選択される。そして、これら第1ブロック内の有効データがコピー先フリーブロックにコピーされ、これら第1ブロックのGC回数に1を加えた値が、コピー先フリーブロックのGC回数として設定される。したがって、更新頻度の高いデータと更新頻度の低いデータとがGC動作によって一緒に同じブロックにコピーされてしまうことを防止できるようになる。これにより、GC回数の多いブロックほど、ブロックの容量に対する更新頻度の低いデータの量の割合を増やすことができるので、更新頻度の低いデータを、更新頻度の高いデータから分離することが可能とする。このことは、更新頻度の異なる複数種のデータが混在するブロックの数の増加を抑制できることを意味する。よって、たとえSSD3に書かれるデータが高いデータ局所性を有する場合であっても、更新頻度の高いデータと更新頻度の低いデータとが混在するブロックの数の増加を抑制でき、この結果、SSD3のライトアンプリフィケーションの増加を抑制できる。

30

40

【0230】

なお、本実施形態では、不揮発性メモリとしてNANDメモリを例示した。しかし、本実施形態の機能は、例えば、MRAM(Magnetoresistive Random Access Memory)、PRAM(Phase change Random Access Memory)、ReRAM(Resistive Random Access Memory)、又は、FeRAM(Ferroelectric Random Access Memory)のような他の様々な不揮発性メモリにも適用できる。

50

【0231】

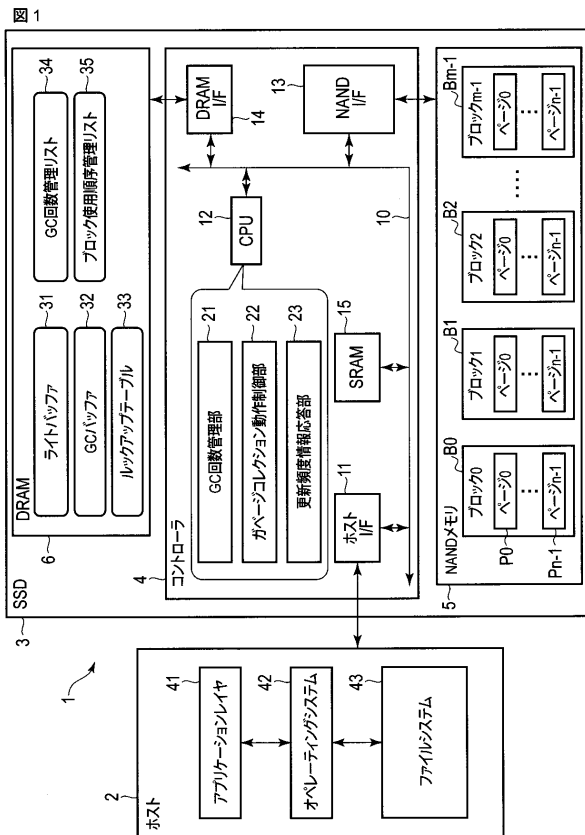
本発明のいくつかの実施形態を説明したが、これらの実施形態は、例として提示したものであり、発明の範囲を限定することは意図していない。これら新規な実施形態は、その他の様々な形態で実施されることが可能であり、発明の要旨を逸脱しない範囲で、種々の省略、置き換え、変更を行うことができる。これら実施形態やその変形は、発明の範囲や要旨に含まれるとともに、特許請求の範囲に記載された発明とその均等の範囲に含まれる。

【符号の説明】

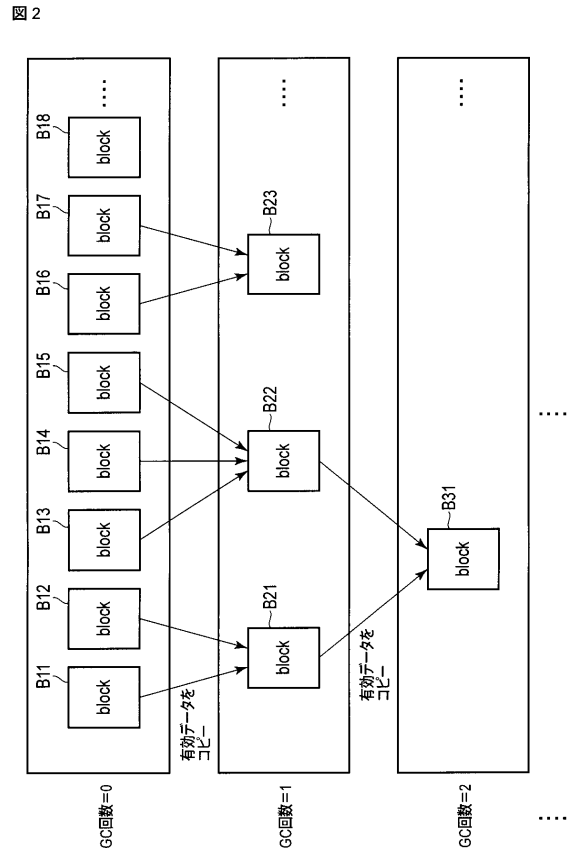
【0232】

2 ... ホスト、3 ... SSD、4 ... コントローラ、5 ... NANDメモリ、21 ... ガベージコレクション回数管理部、22 ... ガベージコレクション動作制御部、23 ... 更新頻度情報応答部。

【図1】



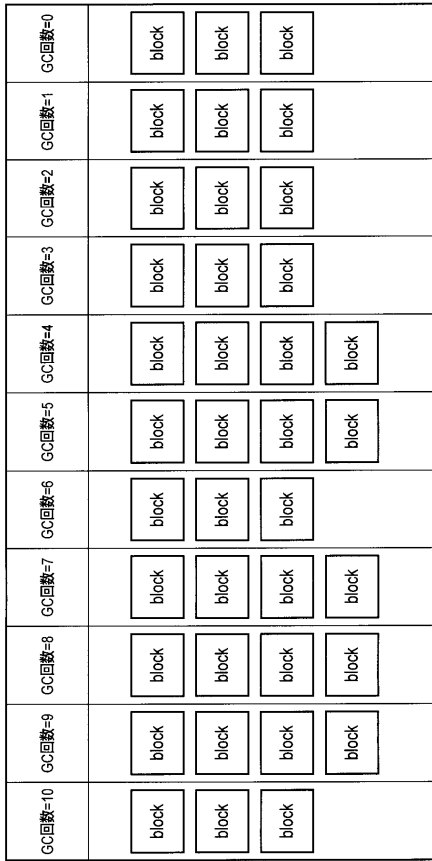
【図2】



【 図 3 】

図 3

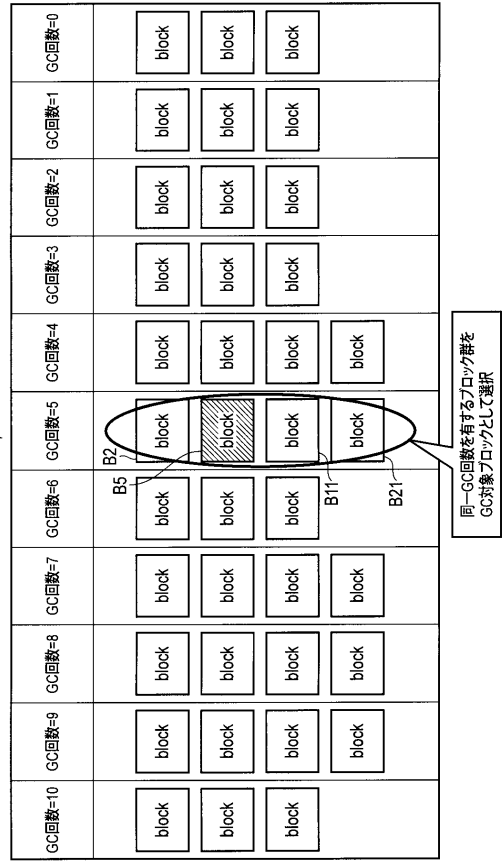
34



【 図 4 】

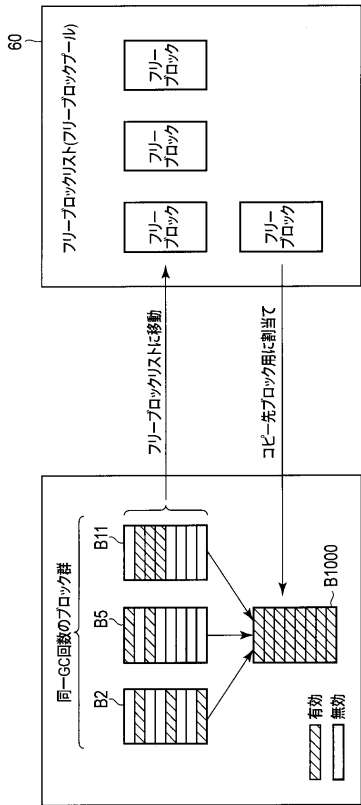
図 4

34



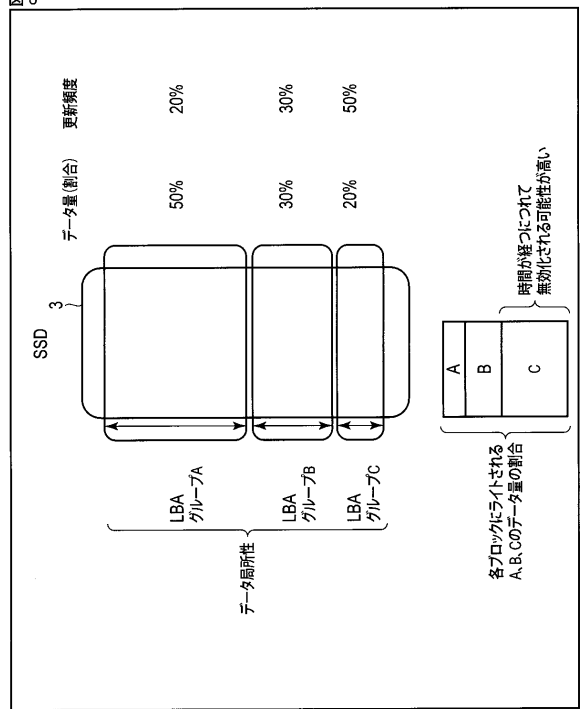
【 図 5 】

図 5



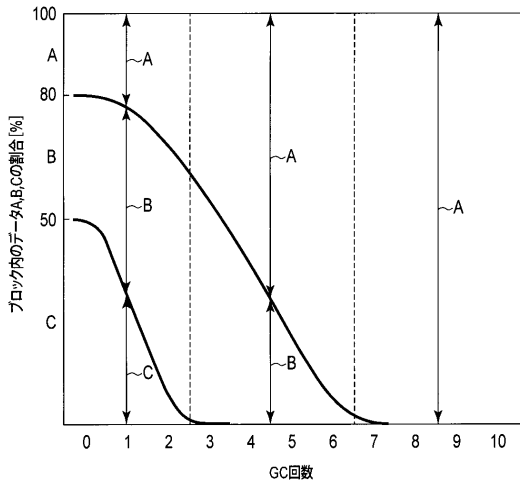
【 図 6 】

図 6



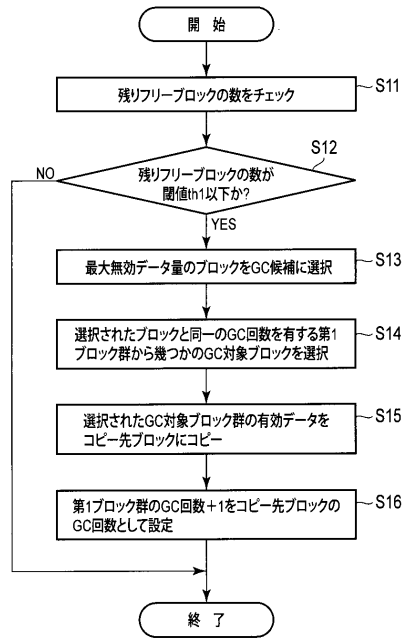
【 図 7 】

図 7



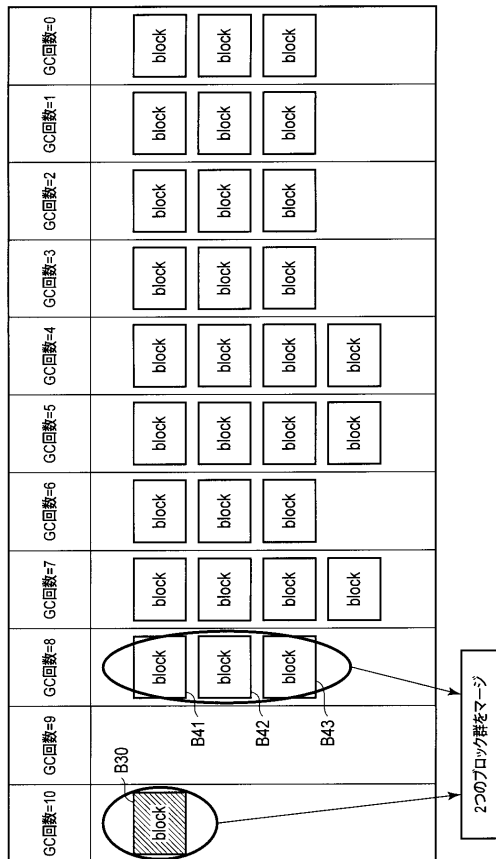
【 図 8 】

図 8



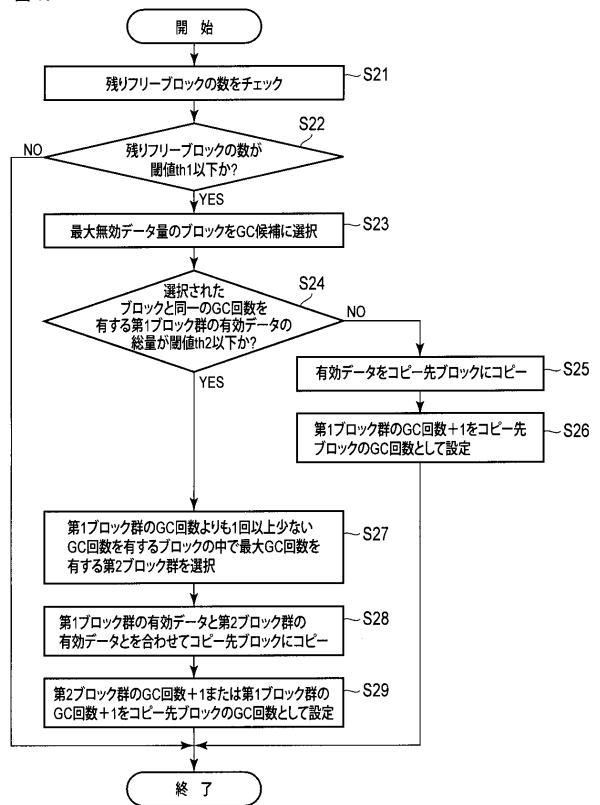
【 図 9 】

図 9

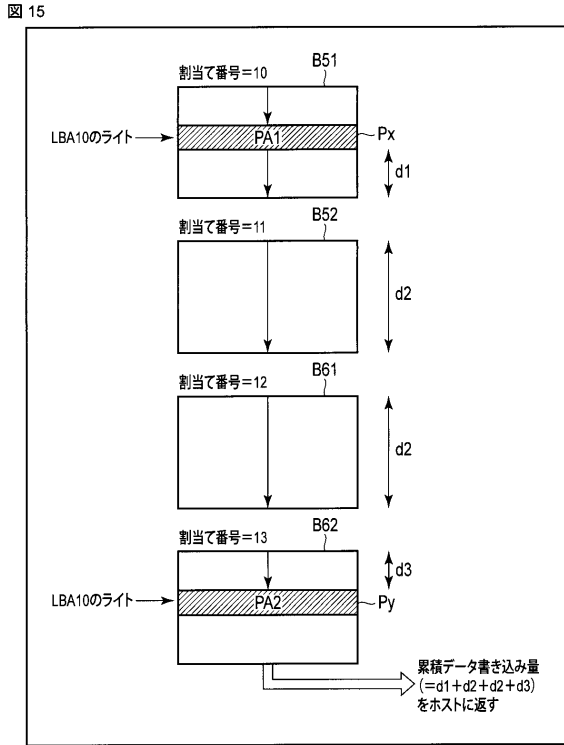


【 図 10 】

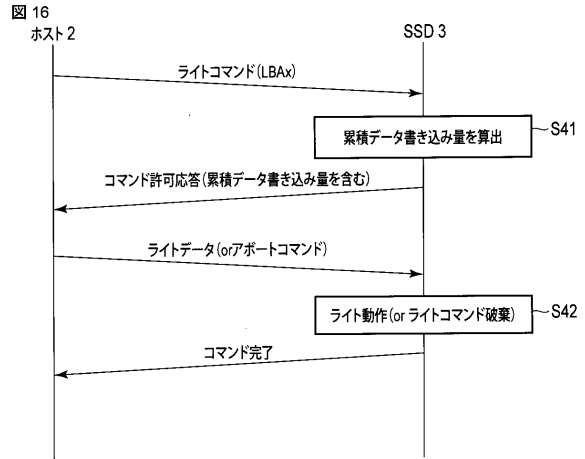
図 10



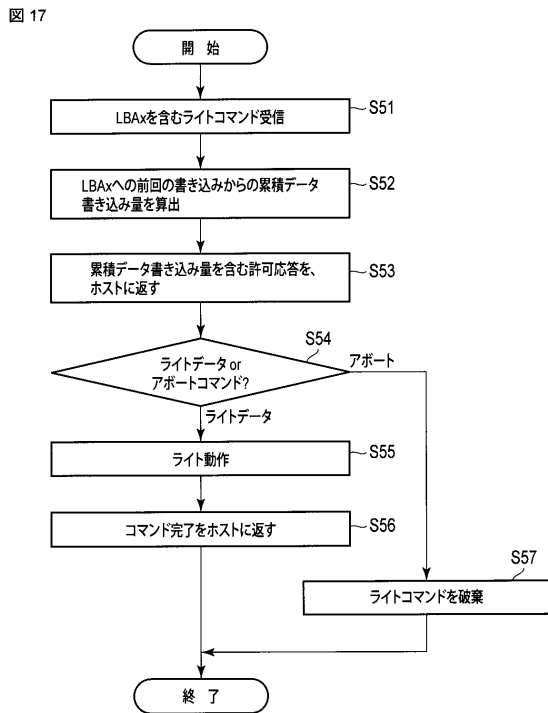
【 図 1 5 】



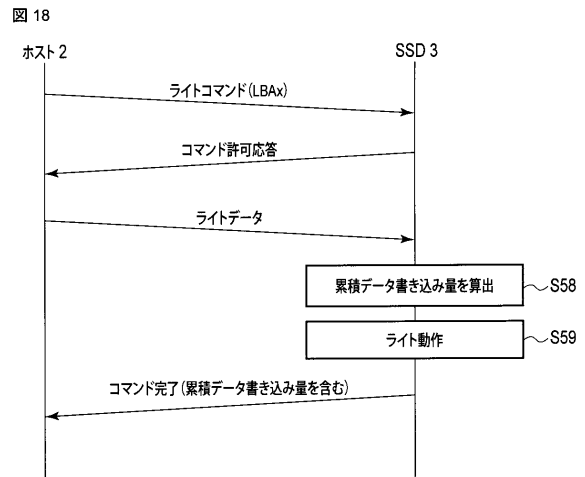
【 図 1 6 】



【 図 1 7 】

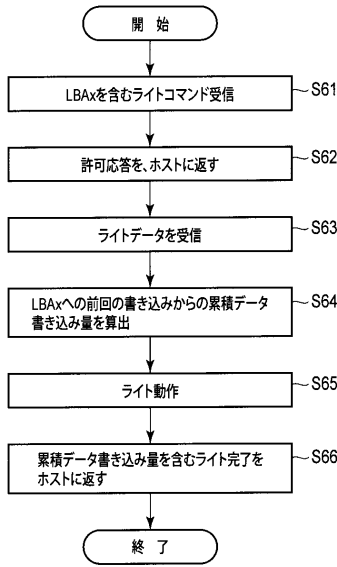


【 図 1 8 】



【 図 1 9 】

図 19



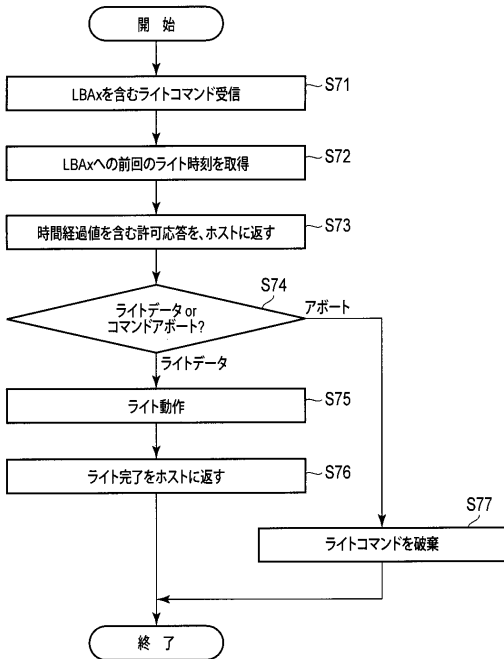
【 図 2 0 】

図 20

	33A 物理アドレス(PA)	33B ライトされた時刻
LBA #0		
#1		
#2		
#3		
#4		
#5		
#6		
⋮	⋮	⋮
#n		

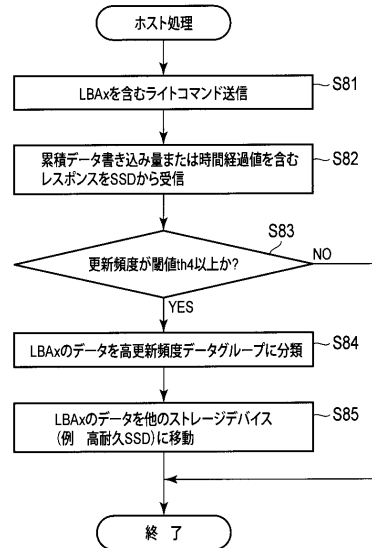
【 図 2 1 】

図 21



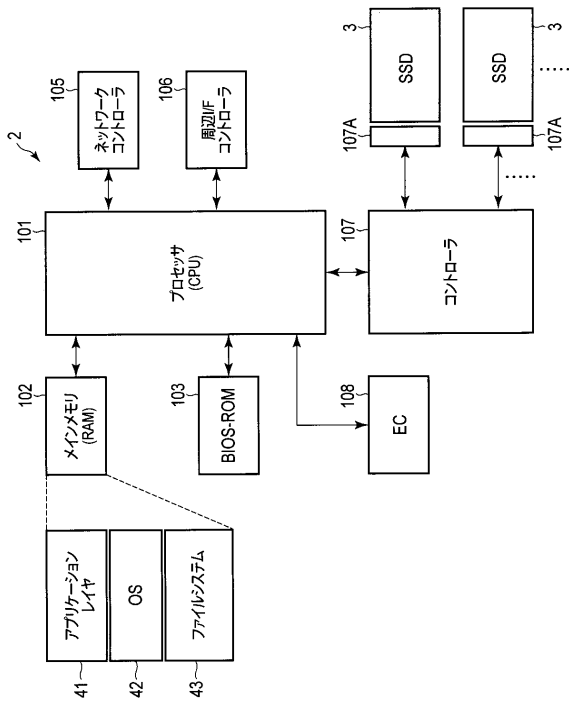
【 図 2 2 】

図 22



【 図 2 3 】

図 23



【 図 2 4 】

図 24

