

(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2009-199199

(P2009-199199A)

(43) 公開日 平成21年9月3日(2009.9.3)

(51) Int.Cl.	F I	テーマコード (参考)
<b>G 0 6 F 3/06 (2006.01)</b>	G 0 6 F 3/06 3 0 2 A	5 B 0 0 5
<b>G 0 6 F 12/08 (2006.01)</b>	G 0 6 F 12/08 5 5 1 B	5 B 0 6 5
	G 0 6 F 12/08 5 5 3 Z	
	G 0 6 F 12/08 5 4 1 Z	
	G 0 6 F 12/08 5 5 7	

審査請求 未請求 請求項の数 18 O L (全 17 頁)

(21) 出願番号 特願2008-38176 (P2008-38176)  
 (22) 出願日 平成20年2月20日 (2008.2.20)

(71) 出願人 000005108  
 株式会社日立製作所  
 東京都千代田区丸の内一丁目6番6号  
 (74) 代理人 100093861  
 弁理士 大賀 真司  
 (72) 発明者 水島 永雅  
 神奈川県川崎市麻生区王禅寺1099番地  
 株式会社日立製作所システム開発研究所  
 内  
 Fターム(参考) 5B005 JJ01 JJ12 MM01 UU23  
 5B065 BA05 CE01 CH02

(54) 【発明の名称】 ストレージシステム及びそのデータライト方法

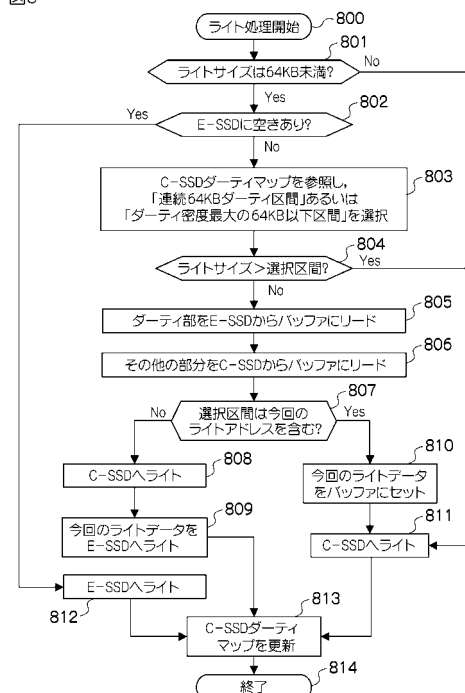
## (57) 【要約】

【解決課題】ライトデータを高性能側デバイスでキャッシュしてから、そのデータを高性能側デバイスから低性能側デバイスへコピーすると、低性能側デバイス内のフラッシュメモリの書き換え寿命が浪費される。

【解決手段】低性能側の不揮発性メモリデバイス内部のメモリ管理単位のサイズを保持し、ライトデータのサイズとメモリ管理単位のサイズとを比較する。ライトデータのほうが小さいときは、ライトデータを高性能側の不揮発性メモリデバイスへキャッシュし、そうでなければ低性能側デバイスへ書き込む。そして、高性能側デバイスにキャッシュされた複数のライトデータのアドレス値を参照して、キャッシュされたアドレス値がメモリ管理単位のサイズだけ連続しているアドレス区間を選択し、そのアドレス区間内に含まれるデータを、高性能側デバイスから低性能側デバイスへコピーする。

【選択図】 図 8

図8



**【特許請求の範囲】****【請求項 1】**

所定の性能である第 1 の不揮発性メモリデバイスと、前記所定の性能より高性能である第 2 の不揮発性メモリデバイスとを含んで構成されるストレージシステムであって、

前記第 1 の不揮発性メモリデバイス内のメモリを管理するメモリ管理単位のサイズを保持する保持部と、

上位装置からのライト要求に応答して、そのライト要求されたライトデータのサイズと前記メモリ管理単位のサイズを比較し、前記ライトデータのサイズが前記メモリ管理単位のサイズより小さいときは、前記ライトデータを前記第 2 の不揮発性メモリデバイスへ一時的に書き込み、前記ライトデータのサイズが前記メモリ管理単位のサイズ以上のときは、前記ライトデータを前記第 1 の不揮発性メモリデバイスへ書き込む制御部と、

を備えることを特徴とするストレージシステム。

**【請求項 2】**

前記制御部は、前記第 2 の不揮発性メモリデバイスに一時的に書き込まれた複数の前記ライトデータのアドレス値を参照し、その参照したアドレス値が前記メモリ管理単位のサイズだけ連続しているアドレス区間を選択し、その選択されたアドレス区間に含まれるライトデータを前記第 2 の不揮発性メモリデバイスから前記第 1 のメモリデバイスへコピーすることを特徴とする請求項 1 記載のストレージシステム。

**【請求項 3】**

前記制御部は、前記第 2 の不揮発性メモリデバイスに一時的に書き込まれた複数の前記ライトデータのアドレス値を参照し、前記メモリ管理単位サイズ以下で、前記アドレス値を最も多く含むアドレス区間を選択し、その選択されたアドレス区間に含まれるライトデータを前記第 2 の不揮発性メモリデバイスから読み出し、前記アドレス区間内に含むことができるライトデータを前記第 1 の不揮発性メモリデバイスから読み出し、これら読み出したライトデータから連続データを作成し、その作成した連続データを前記第 1 の不揮発性メモリデバイスへ書き込むことを特徴とする請求項 1 記載のストレージシステム。

**【請求項 4】**

前記第 1 の不揮発性メモリデバイスから読み出されるライトデータは、前記第 2 の不揮発性メモリデバイスから読み出されるライトデータと連続させたときに、前記アドレス区間に格納されるデータ量以下となるデータが読み出されることを特徴とする請求項 3 記載のストレージシステム。

**【請求項 5】**

前記第 2 の不揮発性メモリデバイスのメモリ管理単位のサイズは、前記第 1 の不揮発性メモリデバイスのメモリ管理単位のサイズより小さいことを特徴とする請求項 1 記載のストレージシステム。

**【請求項 6】**

前記第 1 の不揮発性メモリデバイスの性能と前記第 2 の不揮発性メモリデバイスの性能との差は、少なくともライト性能の差を含むことを特徴とする請求項 1 記載のストレージシステム。

**【請求項 7】**

前記第 1 の不揮発性デバイスはコンシューマ向けの半導体記憶装置であり、前記第 2 の不揮発性メモリデバイスは企業向けの半導体記憶装置であることを特徴とする請求項 1 記載のストレージシステム。

**【請求項 8】**

所定の性能である第 1 の不揮発性メモリデバイスと、前記所定の性能より高性能である第 2 の不揮発性メモリデバイスとを含んで構成されるストレージシステムのデータライト方法であって、

前記第 1 の不揮発性メモリデバイス内のメモリを管理するメモリ管理単位のサイズを保持するステップと、

上位装置からのライト要求に応答して、そのライト要求されたライトデータのサイズと

10

20

30

40

50

前記メモリ管理単位のサイズを比較するステップと、

前記ライトデータのサイズが前記メモリ管理単位のサイズより小さいときは、前記ライトデータを前記第2の不揮発性メモリデバイスへ一時的に書き込み、前記ライトデータのサイズが前記メモリ管理単位のサイズ以上のときは、前記ライトデータを前記第1の不揮発性メモリデバイスへ書き込むステップと、

を含むことを特徴とするストレージシステムのデータライト方法。

【請求項9】

前記第2の不揮発性メモリデバイスに一時的に書き込まれた複数の前記ライトデータのアドレス値を参照するステップと、

その参照したアドレス値が前記メモリ管理単位のサイズだけ連続しているアドレス区間を選択するステップと、

その選択されたアドレス区間に含まれるライトデータを前記第2の不揮発性メモリデバイスから前記第1のメモリデバイスへコピーするステップと、

を含むことを特徴とする請求項8記載のストレージシステムのデータライト方法。

【請求項10】

前記第2の不揮発性メモリデバイスに一時的に書き込まれた複数の前記ライトデータのアドレス値を参照するステップと、

前記メモリ管理単位サイズ以下で、前記アドレス値を最も多く含むアドレス区間を選択し、その選択されたアドレス区間に含まれるライトデータを前記第2の不揮発性メモリデバイスから読み出すステップと、

前記アドレス区間内に含むことができるライトデータを前記第1の不揮発性メモリデバイスから読み出すステップと、

これら読み出したライトデータから連続データを作成し、その作成した連続データを前記第1の不揮発性メモリデバイスへ書き込むステップと、

を含むことを特徴とする請求項8記載のストレージシステムのデータライト方法。

【請求項11】

前記第1の不揮発性メモリデバイスから読み出されるライトデータは、前記第2の不揮発性メモリデバイスから読み出されるライトデータと連続させたときに、前記アドレス区間に格納されるデータ量以下となるデータが読み出されることを特徴とする請求項10記載のストレージシステムのデータライト方法。

【請求項12】

前記第2の不揮発性メモリデバイスのメモリ管理単位のサイズは、前記第1の不揮発性メモリデバイスのメモリ管理単位のサイズより小さいことを特徴とする請求項8記載のストレージシステムのデータライト方法。

【請求項13】

前記第1の不揮発性メモリデバイスの性能と前記第2の不揮発性メモリデバイスの性能との差は、少なくともライト性能の差を含むことを特徴とする請求項8記載のストレージシステムのデータライト方法。

【請求項14】

前記第1の不揮発性デバイスはコンシューマ向けの半導体記憶装置であり、前記第2の不揮発性メモリデバイスは企業向けの半導体記憶装置であることを特徴とする請求項8記載のストレージシステムのデータライト方法。

【請求項15】

ストレージシステムに用いられるアダプタ装置であって、

所定の性能である第1の不揮発性メモリデバイスが有するインタフェースと接続する第1のインタフェースと、

前記所定の性能より高性能である第2の不揮発性メモリデバイスが有するインタフェースと接続する第2のインタフェースと、

前記第1の不揮発性メモリデバイス内のメモリを管理するメモリ管理単位のサイズを保持する保持部と、

10

20

30

40

50

上位装置からのライト要求に応答して、そのライト要求されたライトデータのサイズと前記メモリ管理単位のサイズを比較し、前記ライトデータのサイズが前記メモリ管理単位のサイズより小さいときは、前記ライトデータを前記第 2 の不揮発性メモリデバイスへ一時的に書き込み、前記ライトデータのサイズが前記メモリ管理単位のサイズ以上のときは、前記ライトデータを前記第 1 の不揮発性メモリデバイスへ書き込む制御部と、  
を備えることを特徴とするアダプタ装置。

【請求項 16】

前記第 1 のインタフェースと前記第 2 のインタフェースは、異なる仕様のインタフェースであり、前記第 1 のインタフェースが前記ストレージシステムとのインタフェースであり、

前記第 2 のインタフェースを前記ストレージシステムとのインタフェースと互換性を持たせるための第 3 のインタフェースを備えること、

を特徴とする請求項 15 記載のアダプタ装置。

【請求項 17】

前記第 1 のインタフェースはシリアル ATA インタフェースであり、前記第 2 のインタフェースは SAS インタフェースであることを特徴とする請求項 16 記載のアダプタ装置。

【請求項 18】

前記第 2 の不揮発性メモリデバイスのメモリ管理単位のサイズは、前記第 1 の不揮発性メモリデバイスのメモリ管理単位のサイズより小さいことを特徴とする請求項 15 記載のアダプタ装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、電氣的に書き換え可能な不揮発性メモリを使用した半導体記憶装置を搭載したストレージシステム及びそのデータライト方法に関し、特に半導体記憶装置のライト性能特性に応じてデータ格納をするストレージシステム及びそのデータライト方法に関する。

【背景技術】

【0002】

特許文献 1 には、性能差のある不揮発性デバイス 2 種からなる記憶装置においてライト性能を向上させる方法が開示されている。その方法では、低性能側の不揮発性デバイスがライト可能になるまで期間は、高性能側の不揮発性デバイスで所定量のライトデータをキャッシュし、後からそのデータを低性能側デバイスへコピーし、それ以降のライト先も低性能側デバイスへスイッチする。例えば、高性能側デバイスはフラッシュメモリであり、低性能側デバイスは磁気ディスクである。また、上記の「ライト可能になるまで期間」は磁気ヘッドのシーク時間に相当する。このような記憶装置はハイブリッドハードディスクと呼ばれている。

【特許文献 1】米国登録特許 7 1 3 6 9 7 3 号明細書

【発明の開示】

【発明が解決しようとする課題】

【0003】

いま、高性能側・低性能側デバイスともフラッシュメモリを用いた半導体記憶装置であるようなストレージシステムに、特許文献 1 のライト方法を適用した場合を考える。この方法では、高性能側デバイスから低性能側デバイスへデータコピーを行う制御において、そのデータサイズの最適性が考慮されていない。それゆえ、当該ストレージシステムに発生しうる問題として、低性能側デバイス内のフラッシュメモリの書き換え寿命（1メモリブロック当たり約 10 万回）が浪費される恐れがある。これは、一般に低性能のフラッシュメモリ記憶装置は、内部の制御ファームウェアで既定されたメモリ管理単位（例えば 64 KB）で内部フラッシュメモリの書き換えるように設計されていることに起因する。例

10

20

30

40

50

えば、上記のデバイス間コピーサイズが4KBである場合には、変更されない周辺データ60KBを加えた64KBのデータが内部フラッシュメモリの空き領域にプログラムされる。つまり、低性能側デバイス内部で無駄なデータプログラムが発生する。これが、低性能側デバイスの書き換え寿命の浪費となる。一方、高性能側デバイスは、一般に外部からのデータライトに対して必要最低限のデータを内部フラッシュメモリにプログラムするような制御を実施する。そのため、高性能側デバイス内部でこのような無駄なデータプログラムはほとんど発生しない。

【0004】

本発明は、以上の点を考慮してなされたもので、小さなサイズのライトデータを高性能側デバイスで複数回キャッシュしてから低性能側デバイスにライトし、平均のライト処理時間を削減し、ライト性能を向上させたストレージシステム及びそのデータライト方法を提案しようとするものである。

【課題を解決するための手段】

【0005】

本発明は、所定の性能である第1の不揮発性メモリデバイスと、前記所定の性能より高性能である第2の不揮発性メモリデバイスとを含んで構成されるストレージシステムであって、前記第1の不揮発性メモリデバイス内のメモリを管理するメモリ管理単位のサイズを保持する保持部と、上位装置からのライト要求に応答して、そのライト要求されたライトデータのサイズと前記メモリ管理単位のサイズを比較し、前記ライトデータのサイズが前記メモリ管理単位のサイズより小さいときは、前記ライトデータを前記第2の不揮発性メモリデバイスへ一時的に書き込み、前記ライトデータのサイズが前記メモリ管理単位のサイズ以上のときは、前記ライトデータを前記第1の不揮発性メモリデバイスへ書き込む制御部とを備えるストレージシステム及びそのデータライト方法である。

【0006】

さらに具体的には、性能差のある2種の不揮発性メモリデバイスから構成されるストレージシステムにおいて、以下を特徴とするデータライト方法を提供する。まず、低性能側の不揮発性メモリデバイス内部のメモリ管理単位のサイズを保持する。次に、前記ストレージシステムに対するライト要求に応答して以下を行う。(1)ライトデータのサイズと前記メモリ管理単位のサイズとを比較する。(2)当該ライトデータのほうが小さいときは、当該ライトデータを高性能側の不揮発性メモリデバイスへキャッシュし、そうでなければ前記低性能側デバイスへ書き込む。(3)次の(A)(B)のいずれか又は両方を行う。(A)前記高性能側デバイスにキャッシュされた複数の前記ライトデータのアドレス値を参照して、当該アドレス値が前記メモリ管理単位のサイズだけ連続しているアドレス区間を選択し、当該アドレス区間内に含まれる前記ライトデータを、前記高性能側デバイスから前記低性能側デバイスへコピーする。(B)前記高性能側デバイスにキャッシュされた複数の前記ライトデータのアドレス値を参照して、前記メモリ管理単位のサイズ以下で、前記アドレス値を最も多く含むアドレス区間を選択し、当該アドレス区間内に含むことができる前記ライトデータを前記高性能側デバイスから読み出し、ライトデータを前記低性能側デバイスから読み出し、当該アドレス区間の連続データを作成し、当該連続データを前記低性能側デバイスへライトする。

【発明の効果】

【0007】

本発明によれば、性能差のある2種の不揮発性メモリデバイスから構成されるストレージシステムにおいて、小さなサイズのライトデータを高性能側デバイスで複数回キャッシュしてから低性能側デバイスにライトするため、平均のライト処理時間が削減され、ストレージシステムのライト性能が向上するという効果を奏する。

【0008】

さらに、そのキャッシュデータを低性能側デバイスへライトする際には、低性能側デバイス内のフラッシュメモリの書き換え寿命が浪費されるようなデータサイズのライトが回避されるため、ストレージシステム内のフラッシュメモリの書き換え寿命が改善されると

10

20

30

40

50

いう効果を奏する。

【発明を実施するための最良の形態】

【0009】

以下、本発明の各実施形態について説明する。

まず、本発明を適用するストレージシステムにおいて、ユーザデータ記憶媒体として搭載される２種類の半導体記憶装置について、それぞれの内部ハードウェア構成とライト性能特性について図１～図４を用いて説明する。以下、半導体記憶装置を“Solid State Disk”と呼び、SSDと略す。

【0010】

１つ目のSSDは一般消費者向けに製品化されているコンシューマ向けSSD（以下、C-SSDと略す。）である。２つ目のSSDは企業向けに製品化されているエンタプライズ向けSSD（以下、E-SSDと略す。）である。C-SSDは原価や利益率をできるだけ安くして携帯型電子機器など向けの記憶装置市場に大量に流通させ、薄利多売で利益を得ることを意図した製品である。製造コストを安く抑えることを優先して安価なプロセッサを使用したり、メモリリソースを少なくしたりするため、性能はE-SSDより低い。

10

【0011】

一方、E-SSDは性能をできるだけ高くしてハイエンドの顧客要求を満足させることを意図した製品である。性能を高めることを優先して高価な部品を使用したり、高機能な制御ファームウェアを実装したりするため、C-SSDより製造コストが高い。業務用サーバ向け記憶装置などを主な応用先としており市場流通量があまり多くないため、利益率も高く設定されている。その結果、一般的なE-SSD価格は同容量のC-SSD価格の約５倍である。これは、コンシューマ向けハードディスクドライブとエンタプライズ向けハードディスクドライブに見られる価格差の存在と同様である。

20

【0012】

図１はC-SSD100のハードウェア構成を示す。C-SSD100は、メモリコントローラ110、およびフラッシュメモリ120を備える。フラッシュメモリ120はデータを不揮発に記憶する。メモリコントローラ110は、フラッシュメモリ120のデータの「リード」、「プログラム」、および「消去」を実行する。メモリコントローラ110は、プロセッサ112、SATA（シリアルATA）インタフェース111、データ転送部115、RAM113、およびROM114を備える。データ転送部115はバスロジックやフラッシュメモリ120の制御ロジックを含み、その他の構成要素111～114やフラッシュメモリ120と接続する。プロセッサ112はROM114に格納された制御ファームウェアに従ってデータ転送部115を制御する。RAM113は転送データ用バッファメモリや制御ファームウェア用ワークメモリとして機能する。また、フラッシュメモリ120は複数のフラッシュメモリチップ121で構成される。C-SSD100全体を動作させる電源はSATAインタフェース111から外部供給される。

30

【0013】

図２はE-SSD200のハードウェア構成を示す。E-SSD200は、メモリコントローラ210、フラッシュメモリ220、およびバックアップ電源230を備える。フラッシュメモリ220はデータを不揮発に記憶する。メモリコントローラ210は、フラッシュメモリ220のデータの「リード」、「プログラム」、および「消去」を実行する。メモリコントローラ210は、プロセッサ212、SAS（シリアル・アタッチド・SCSI）インタフェース211、データ転送部215、RAM213、およびROM214を備える。データ転送部215はバスロジックやフラッシュメモリ220の制御ロジックを含み、その他の構成要素211～214やフラッシュメモリ220と接続する。プロセッサ212はROM214に格納された制御ファームウェアに従ってデータ転送部215を制御する。RAM213は転送データ用バッファメモリや制御ファームウェア用ワークメモリとして機能する。また、フラッシュメモリ220は複数のフラッシュメモリチップ221で構成される。

40

50

## 【 0 0 1 4 】

なお、S A S インタフェース 2 1 1 は 2 つのポートを備え、2 つの独立したアクセスを非同期に受け入れることができる。一方のポートのアクセス経路に障害が置けても、もう一方を使ってアクセスを継続することが可能である。

## 【 0 0 1 5 】

E - S S D 2 0 0 全体を動作させる電源は、基本的に S A S インタフェース 2 1 1 から外部供給されるが、外部供給が絶たれた場合にはバックアップ電源 2 3 0 から供給される。外部供給が絶たれたときに R A M 2 1 3 内にフラッシュメモリ 2 2 0 に書き込むべきデータが残っていた場合は、バックアップ電源 2 3 0 からの電力を利用してそのデータをフラッシュメモリ 2 2 0 に書き込む。そして、切断された外部供給が再開するまでは外部アクセスは受け入れない。

10

## 【 0 0 1 6 】

図 3 ( a ) を用いて、C - S S D 1 0 0 のデータライト処理方式および性能特性を説明する。

## 【 0 0 1 7 】

各フラッシュメモリチップ 1 2 1 は複数 (例えば 4 0 9 6 個) のメモリブロック 3 0 1 で構成される。メモリブロックはフラッシュメモリの消去単位であり、そのサイズは例えば 2 5 6 K B である。1 つのメモリブロック 3 0 1 を消去する所要時間は 2 m s である。さらに、各メモリブロック 3 0 1 は複数 (例えば 6 4 個) のメモリページ 3 0 2 で構成される。メモリページはフラッシュメモリ 1 2 0 のプログラム単位であり、そのサイズは 4 K B である。1 つのメモリページ 3 0 2 をプログラムする所要時間は 5 0 0  $\mu$  s、リードする所要時間は 5 0  $\mu$  s である。C - S S D 1 0 0 内では連続する複数 (例えば 1 6 個) のメモリページ 3 0 2 をまとめた管理単位 3 0 3 を作る。C - S S D 1 0 0 外部からアクセスするときの論理アドレス空間を、そのサイズを単位として分割し、各分割要素をフラッシュメモリ 1 2 0 全体の物理アドレス (チップ番号、ブロック番号、管理単位番号) へ対応付けする。この対応付けしたテーブルをアドレス変換テーブルと呼ぶ。このアドレス変換テーブルは、C - S S D 1 0 0 外部からのライトアクセスによって更新される。なぜなら、フラッシュメモリ 1 2 0 は構造的に上書きできない記憶素子であるからである。つまり、プログラムすべきデータは元のデータとは違う未書き込みの領域に書き、元のデータがあったメモリブロック 3 0 1 は後で消去する必要があるため、各論理アドレスのデータの物理的所在は移動せざるを得ないからである。このアドレス変換テーブルは R A M 1 1 3 上に設置される。複数 (1 6 個) のメモリページをまとめたものを管理単位とするのは、対応付けの要素数を減らし、メモリリソース量を節約するためである。

20

30

## 【 0 0 1 8 】

いま、C - S S D 1 0 0 外部から 1 K B のデータライトがあったとすると、まず、プロセッサ 1 1 2 は、そのデータの論理アドレスを含むアドレス区間に対応する管理単位 3 0 4 を選択し、その領域内でライト対象外の 6 3 K B データ 3 0 7 を R A M 1 1 3 上にリードする (3 0 5)。そして 1 K B のライトデータ 3 0 6 を R A M 1 1 3 上に設定し、それら 6 4 K B のデータを未書き込みの管理単位 3 0 8 にプログラムする (3 0 9)。リード 3 0 5 の所要時間は 1 6 個のメモリページ 3 0 2 を読むため、1 6 回  $\times$  5 0  $\mu$  s = 0 . 8 m s である。プログラム 3 0 9 の所要時間は 1 6 個のメモリページを書くため、1 6 回  $\times$  5 0 0  $\mu$  s = 8 m s である。つまり、C - S S D 1 0 0 外部から 1 K B をデータライトする時にはデバイスレベルで 8 . 8 m s の時間を要する。なお、実効的な平均処理時間は、ライトデータの転送時間や、時々行うメモリブロック消去時間等がこれに加算されたものである。

40

## 【 0 0 1 9 】

図 3 ( b ) は、以上のライト方式に基づいて、C - S S D 1 0 0 におけるライトデータサイズと平均処理時間 (m s)、性能 (I O P S : 1 秒当たりの平均アクセス回数) との関係を表したものである。処理時間は左縦軸を用いた棒グラフ、性能は右縦軸を用いた実線グラフで示す。なお、平均処理時間はデバイスレベル処理時間とその他の処理時間 (デ

50

ータ転送等の時間)に分けて示す。

#### 【0020】

管理単位のサイズ64KBより小さいデータ(XKB)をライトした場合、C-SSD100内部では(64-X)KBのデータリードと64KBの管理単位データプログラムが発生するため、デバイスレベルで8~8.8msの時間を要する。また、128KB、256KBをライトした場合、C-SSD100内部ではプログラムすべき管理単位の個数に応じてそれぞれ16ms、32msの時間を要する。このように、C-SSD100に対する小単位のデータライトは、そのアドレス周辺データの移動を伴うため、内部のフラッシュメモリ120の有限な書き換え寿命は必要以上に浪費されることになる。

#### 【0021】

また、性能(IOPS)は平均処理時間の逆数として求めたものである。ライトデータのサイズが管理単位の64KBよりも大きい場合は、サイズが小さくなるにつれて性能は着実に増加するが、管理単位の64KBより小さくなると110 IOPS程度を漸近線として収束してしまう。

#### 【0022】

次に、図4(a)を用いて、E-SSD200のデータライト処理方式および性能特性を説明する。

#### 【0023】

各フラッシュメモリチップ221は、C-SSD100のフラッシュメモリチップ121と同じものであり、複数(例えば4096個)のメモリブロック301で構成される。各メモリブロック301は複数(例えば64個)のメモリページ302で構成される。E-SSD200内では1個のメモリページ302を管理単位とする。E-SSD200外部からアクセスするときの論理アドレス空間を、メモリページ302を単位として分割し、各分割要素をフラッシュメモリ220全体の物理アドレス(チップ番号、ブロック番号、ページ番号)へ対応付けする。この対応付けをアドレス変換テーブルと呼ぶ。このアドレス変換テーブルは、E-SSD200外部からのライトアクセスによって更新される。また、このアドレス変換テーブルはRAM213上に設置される。

#### 【0024】

いま、E-SSD200外部から1KBのデータライトが多数回あったとする。それらのデータはRAM213上に一旦バッファリングされる。そのうち、同一のページ論理アドレスに含まれる4個の1KBデータ310~313がバッファ上にあれば、それらを結合して4KBのページデータ314を作る。そして、そのデータを未書き込みの物理ページ315にプログラムする(316)。プログラム316の所要時間は1個のメモリページを書くため、1回×500μs=0.5msである。これは1KBライト4回分に相当するため、1回当たりの平均では約0.13msである。つまり、E-SSD200外部から1KBをデータライトする時にはデバイスレベルで0.13msの時間を要する。

#### 【0025】

上記ライト方式において、RAM213が枯渇するまでライトデータをバッファリングしても、4KBのページデータ314が作れない場合は、フラッシュメモリ220から不足データをリードし、ページデータ314を補完する。これはライト性能低下の要因となる。つまり、RAM213上にライトデータを多くバッファリングできる製品ほどライト性能が高くなる。そのため、高性能を追求したE-SSD200は大容量のRAM213を搭載する。

#### 【0026】

なお、実効的な平均処理時間は、ライトデータの転送時間や、時々行うメモリブロック消去時間等がこれに加算されたものである。

#### 【0027】

図4(b)は、以上のライト方式に基づいて、E-SSD200におけるライトデータサイズと平均処理時間(ms)、性能(IOPS)との関係を表したものである。処理時間は左縦軸を用いた棒グラフ、性能は右縦軸を用いた実線グラフで示す。なお、平均処理

10

20

30

40

50



時間はデバイスレベル処理時間とその他の処理時間（データ転送等の時間）に分けて示す。

【 0 0 2 8 】

E - S S D 2 0 0 内部ではライトすべきデータ以外のデータをできるだけフラッシュメモリ 2 2 0 にプログラムしないように制御されるため、内部のフラッシュメモリ 2 2 0 の有限な書き換え寿命は最も無駄なく消費されることになる。

【 0 0 2 9 】

また、性能（I O P S）は平均処理時間の逆数として求めたものである。性能は、ライトデータのサイズが小さくなるにつれて着実に増加し、ディスクドライブの最小ライト単位（1セクタ）である 0 . 5 K B では、1 0 K I O P S 程度に達する。これは C - S S D 1 0 0 の最大性能の約 1 0 0 倍に相当する。

【 0 0 3 0 】

さて、以上に示した C - S S D 1 0 0 と E - S S D 2 0 0 の特徴を踏まえて、本発明の実施形態を詳細に説明する。

【 0 0 3 1 】

図 5 は、本発明を適用したストレージシステム 5 0 0 の内部構成を示した図である。ストレージシステム 5 0 0 は、ホストパッケージ（以下、ホスト P K と略す。）5 1 1、5 2 1、M P U（マイクロプロセッサユニット）P K 5 1 3、5 2 3、キャッシュ P K 5 1 4、5 2 4、バックエンド P K 5 1 5、5 2 5 を備え、それぞれはスイッチ P K 5 1 2、5 2 2 に接続されている。ストレージシステム 5 0 0 の各 P K は冗長化（2 重化）構成を取っている。

【 0 0 3 2 】

ホスト P K 5 1 1、5 2 1 は、ホスト I / F として、F i b e r C h a n n e l や i S C S I 等の I / F コントローラを含むパッケージである。ストレージシステム 5 0 0 は、ホスト P K 5 1 1、5 2 1 と複数のホスト 5 0 1、5 0 2 との間を、S A N（Storage Area Network）5 0 3 を介して接続する。

【 0 0 3 3 】

M P U P K 5 1 3、5 2 3 は、ストレージシステム 5 0 0 を制御する M P U、制御ファームウェアやストレージシステムの構成情報を格納するためのメモリ、および、M P U やキャッシュ等をスイッチ P K 5 1 2、5 2 2 と接続するためのブリッジを含むパッケージである。

【 0 0 3 4 】

キャッシュ P K 5 1 4、5 2 4 は、ストレージシステム 5 0 0 に格納するユーザデータの一次記憶領域であるキャッシュメモリと、キャッシュとスイッチ P K とを接続するキャッシュコントローラを含むパッケージである。

【 0 0 3 5 】

バックエンド P K 5 1 5、5 2 5 は、ストレージシステム 5 0 0 内の複数の S S D ユニット（5 4 0 ~ 5 4 3、5 5 0 ~ 5 5 3 等）を制御する I / F コントローラを含むパッケージである。バックエンド P K 5 1 5、5 2 5 の I / F コントローラは、バックエンドスイッチ 5 1 6、5 2 6 を介して、複数の S S D ユニット 5 4 0 ~ 5 4 3、5 5 0 ~ 5 5 3 等と接続している。バックエンドスイッチ 5 1 6、5 2 6 は S A S 対応のホストバスアダプタとエキスパンダで構成されており、S A S インタフェースと S A T A インタフェースの両方をサポートする機能を持つ。

【 0 0 3 6 】

S S D ユニット（5 4 0 ~ 5 4 3、5 5 0 ~ 5 5 3 等）は、C - S S D 1 0 0、E - S S D 2 0 0、または両者をペアで内蔵する記憶装置ユニットである。各 S S D ユニット 5 4 0 ~ 5 4 3、5 5 0 ~ 5 5 3 等は、ポートが冗長化（2 重化）された S A T A または S A S インタフェースを持つ。そのため、各パッケージやバックエンドスイッチの 1 つに障害が発生した場合にも、冗長化されたポートのいずれかを經由して S S D ユニットのユーザデータへアクセス可能となっている。なお、S S D ユニット 5 4 0 ~ 5 4 3、5 5 0 ~

10

20

30

40

50

5 5 3 等に共通する内部構成については後述する。

【 0 0 3 7 】

ストレージシステム 5 0 0 は、S S D ユニットの故障によるユーザデータ損失を防ぐため、複数の S S D ユニットの R A I D グループを構成し、データの冗長性を図る。例えば、S S D ユニット 5 4 0 ~ 5 4 3 の 4 ユニットで、データ：パリティ比が 3 : 1 の R A I D 5 型のグループ 5 4 4 を構成したり、S S D ユニット 5 5 0 ~ 5 5 3 の 2 × 2 ユニットで R A I D 0 + 1 型のグループ 5 5 4 を構成したりすることができる。

【 0 0 3 8 】

ストレージシステム 5 0 0 は保守クライアント 5 0 4 と接続しており、ユーザは保守クライアント 5 0 4 を通じて、上記 R A I D グループの作成などのストレージ制御をおこなう。

10

【 0 0 3 9 】

以下、図 6、図 7 を用いて S S D ユニット 5 4 0 ~ 5 4 3、5 5 0 ~ 5 5 3 等に共通する内部構成について説明する。図 6 は本発明の第 1 の実施形態における S S D ユニットの内部構成、図 7 は本発明の第 2 の実施形態における S S D ユニットの内部構成である。

【 0 0 4 0 】

第 1 の実施形態では、多数（例えば全体の 9 5 %）の S S D ユニットの図 6（a）のような、S A T A マルチプレクサ 6 0 0 を接続した C - S S D 1 0 0 による構成（以下、構成 A と呼ぶ）とし、それ以外の少数（例えば全体の 5 %）の S S D ユニットの図 6（b）のような E - S S D 2 0 0 による構成（以下、構成 B と呼ぶ）とする。いずれの構成の S S D ユニットも必ず冗長型の R A I D グループに参加させ、各 S S D ユニットの格納データを保護する。なお、S A T A マルチプレクサ 6 0 0 とは、1 ポートの S A T A インタフェースを擬似的に 2 ポートの S A T A インタフェースに見せるアダプタ装置である。

20

【 0 0 4 1 】

第 1 の実施形態では、M P U P K 5 1 3（または 5 2 3）は、構成 A の S S D ユニットのライトすべきデータのうち、6 4 K B 未満のものは基本的に構成 B の S S D ユニットに代行ライトし、必要に応じてそのデータを複数個まとめて構成 B の S S D ユニットから構成 A の S S D ユニットへ移動するというライト処理を行う。

【 0 0 4 2 】

図 8 ~ 図 1 0 を用いて、第 1 の実施形態におけるストレージシステム 5 0 0 のライトアクセス処理手順を説明する。

30

【 0 0 4 3 】

M P U P K 5 1 3（または 5 2 3）は、キャッシュ P K 5 1 4（または 5 2 4）上にダーティなライトデータがあることを検出すると、それを S S D ユニットへライトする処理を開始する（8 0 0）。まず、ライトデータのサイズが 6 4 K B 未満であるかを判定する（8 0 1）。結果が偽（6 4 K B 以上）ならば、ライトデータを構成 A の S S D ユニットの一部へライトし（8 1 1）、ステップ 8 1 3 へ遷移する。一方、ステップ 8 0 1 の結果が真（6 4 K B 未満）ならば、構成 B の S S D ユニットに空き領域があるかを判定する（8 0 2）。結果が真（空きあり）ならば、ライトデータを構成 B の S S D ユニットの一部へライトし（8 1 2）、ステップ 8 1 3 へ遷移する。一方、ステップ 8 0 2 の結果が偽（空きなし）ならば、ステップ 8 0 3 へ遷移する。

40

【 0 0 4 4 】

M P U P K 5 1 3（または 5 2 3）は、ステップ 8 0 3 のために、構成 A の S S D ユニットのアドレス空間内で構成 B の S S D ユニットがライト代行している部分を管理するマップ（C - S S D ダーティマップ）を持つ。これは例えば、キャッシュ P K 5 1 4（または 5 2 4）の一部に設置する。図 9（図 1 0）のように、C - S S D ダーティマップは、構成 A の S S D ユニットのアドレス空間 9 0 0（1 0 0 0）において、構成 B の S S D ユニットがライト代行している部分 9 0 1（1 0 0 1）と、代行していない部分 9 0 2（1 0 0 2）をそれぞれ“1”と“0”で表記したビットマップである。このマップは、ストレージシステム 5 0 0 のシャットダウン時に、構成 B の S S D ユニットの一部にも不揮

50

発に保存しておく。ストレージシステム 500 の起動時には、これを読み出してキャッシュ P K 5 1 4 (または 5 2 4) の一部へ設置する。

【0045】

ステップ 803 では、M P U P K 5 1 3 (または 5 2 3) は、この C - S S D ダーティマップを参照し、図 9 の区間 910 のように、連続 64 K B がダーティとなっている (構成 B の S S D ユニットが記録代行している) アドレス区間を選択する。もし、そういう区間が存在しないならば、図 10 の区間 1010 のように、64 K B 以下の長さのアドレス区間のうち、ダーティとなっている部分の密度が最も高いものを選択する。

【0046】

次に、M P U P K 5 1 3 (または 5 2 3) は、図 10 の区間 1010 のようなものを選択した場合、ライトデータのサイズがその区間の長さよりも大きいかを判定する (804)。結果が真 (大きい) ならば、ライトデータを構成 A の S S D ユニットの一部へライトし (811)、ステップ 813 へ遷移する。一方、ステップ 804 の結果が偽 (大きくない) ならば、その選択区間においてダーティなアドレス部分のデータを構成 B の S S D ユニットからバッファ 920 (1020) にリードする (805、930、1030)。そのバッファ 920、1020 は例えば、キャッシュ P K 5 1 4 (または 5 2 4) の一部を利用する。そして、図 10 の区間 1010 のようなものを選択した場合、その選択区間においてクリーン (非ダーティ) なアドレス部分のデータを構成 A の S S D ユニットからバッファ 1020 にリードする (806、1040)。

【0047】

次に、M P U P K 5 1 3 (または 5 2 3) は、その選択区間が今回のライトデータのアドレスを含むかを判定する (807)。結果が真 (含む) ならば、そのライトデータをバッファ 920 (1020) に設定し (810)、バッファ 920 (1020) 上の当該選択区間データを構成 A の S S D ユニットの一部へライトし (811、940、1050)、ステップ 813 へ遷移する。一方、ステップ 807 の結果が偽 (含まない) ならば、バッファ 920 (1020) 上の当該選択区間データを構成 A の S S D ユニットの一部へライトし (808、940、1050)、さらに、構成 B の S S D ユニットで当該選択区間を代行していた部分 (つまり上書きしてもよい部分) へ今回のライトデータをライトし (809)、ステップ 813 へ遷移する。

【0048】

ステップ 813 では、M P U P K 5 1 3 (または 5 2 3) は、C - S S D ダーティマップにおいて、以上の手順を通して構成 A の S S D ユニットへライトしたアドレス部分をクリーン ("0") に設定し、また、以上の手順を通して構成 B の S S D ユニットへライトしたアドレス部分をダーティ ("1") に設定するように更新を行う。以上により、ライト処理を終了する (814)。

【0049】

第 2 の実施形態では、全ての S S D ユニット 540 ~ 543、550 ~ 553 等を図 7 のような、C - S S D 100 と、それよりも小容量 (例えば容量比で 5%) の E - S S D 200 と、それらと接続した S S D アダプタ 700 による構成とする。

【0050】

S S D アダプタ 700 は、C - S S D 100 および E - S S D 200 それぞれのユーザデータの「読み出し」、「書き込み」を実行する。S S D アダプタ 700 は、プロセッサ 704、S A T A インタフェース 701、702、データ転送部 7703、R A M 705、R O M 706、S A T A インタフェース 707、および S A S インタフェース 708 を備える。

【0051】

データ転送部 703 はバスロジックや S A S や S A T A の制御ロジックを含み、その他の構成要素 701, 702, 704 ~ 708 と接続する。プロセッサ 704 は R O M 706 に格納された制御ファームウェアに従ってデータ転送部 703 を制御する。R A M 705 は転送データ用バッファメモリや制御ファームウェア用ワークメモリとして機能する。

10

20

30

40

50

データ転送部 703 は、2 ポートの S A T A インタフェース 701、702 からの非同期なアクセスを受け入れることができる。C - S S D 100 は S A T A インタフェース 707 に 1 ポート接続され、E - S S D 200 は S A S インタフェース 708 に 2 ポート接続される。後者は 1 ポート接続でもよいが、耐障害性向上のためには冗長化が望ましい。

【0052】

なお、バックエンドスイッチ 516、526 は S A S インタフェースと S A T A インタフェースの両方をサポートするため、S S D アダプタ 700 に含まれる S A T A インタフェース 701、702 は、2 ポートの S A S インタフェースであってもよい。

【0053】

第 2 の実施形態では、S S D アダプタ 700 は、C - S S D 100 にライトすべきデータのうち、64KB 未満のものは基本的に E - S S D 200 に代行ライトし、必要に応じてそのデータを複数個まとめて E - S S D 200 から C - S S D 100 へ移動するというライト処理を行う。

【0054】

第 2 の実施形態におけるストレージシステム 500 のライトアクセス処理手順は、第 1 の実施形態で示した図 8 ~ 図 10 の手順と基本的に同じである。ただし、いくつかの点で異なるため、それを以下に示しておく。

【0055】

まず、処理の実行主体が M P U P K 513、523 ではなく、各 S S D ユニット内の S S D アダプタ 700 である。そして、ステップ 805、806、810 におけるバッファ、および C - S S D ダーティマップの所在はキャッシュ P K 514、524 の一部ではなく、各 S S D アダプタ 700 の R A M 704 内部である。

【0056】

第 2 の実施形態が第 1 の実施形態よりも優れている点は、本発明の適用範囲が S S D ユニットという小規模装置内に内包されるため、共通の S S D ユニットの既存の多様なストレージシステムへ換装するだけでよく、既存ストレージシステムそれぞれのライト制御ファームウェアに変更を行う必要がなく、導入障壁が低いことである。

【0057】

以上、2 つの実施形態を示しながら述べたライトアクセス処理手順によれば、ストレージシステム 500 のライト性能は向上し、フラッシュメモリ 120、220 の書き換え寿命は改善する効果をもたらす。以下、2 つのアクセスパターンを示し、C - S S D 100 のみで構成した従来のストレージシステムと、本発明のストレージシステム 500 とを比較しながら効果の大きさを説明する。

【0058】

1 つ目の例として、キャッシュ P K 514 (または 524) から 1KB のライトバックデータが断続的にあふれ出し、それらが最終的に 64KB の連続なアドレス区間を埋め尽くした場合を考える。従来は、C - S S D に 1KB のライトを少なくとも 64 回与えるため、 $64 \text{ 回} \times 8.8 \text{ ms} = 563.2 \text{ ms}$  のデバイス処理時間が必要だった。本発明では、E - S S D に 1KB のライトを少なくとも 64 回、その後で C - S S D に 64KB のライトを 1 回与えるため、 $64 \text{ 回} \times 0.13 \text{ ms} + 1 \text{ 回} \times 8 \text{ ms} = 16.32 \text{ ms}$  のデバイス処理時間で済む。両者とも、実効的な処理時間はこれにデータ転送等の時間も加算する必要があるが、それでも時間短縮効果の大きさは歴然である。また、フラッシュメモリ書き換え全体量は、従来は  $64 \text{ 回} \times 64 \text{ KB (C - S S D)} = 4096 \text{ KB}$  も発生していたが、本発明では  $64 \text{ 回} \times 1 \text{ KB (E - S S D)} + 1 \text{ 回} \times 64 \text{ KB (C - S S D)} = 128 \text{ KB}$  しか発生しない。

【0059】

2 つ目の例として、キャッシュ P K 514 (または 524) から 1KB のライトバックデータが断続的にあふれ出し、それらが最終的に 63KB の連続なアドレス区間を 1KB 間隔で 32 個埋めた場合を考える。従来は、C - S S D に 1KB のライトを少なくとも 32 回与えるため、 $32 \text{ 回} \times 8.8 \text{ ms} = 281.6 \text{ ms}$  のデバイス処理時間が必要だった

10

20

30

40

50

。本発明では、E - S S Dに1 K Bのライトを少なくとも3 2回、その後でC - S S Dから1 K Bのリードを3 1回、C - S S Dに6 3 K Bのライトを1回与えるため、 $3 2 \times 0 . 1 3 \text{ ms} + 3 1 \times 0 . 0 5 \text{ ms} + 1 \times 8 \text{ ms} = 1 3 . 7 1 \text{ ms}$ のデバイス処理時間で済む。両者とも、実効的な処理時間はこれにデータ転送等の時間も加算する必要があるが、それでも時間短縮効果の大きさは歴然である。また、フラッシュメモリ書き換え全体量は、従来は $3 2 \times 6 4 \text{ K B (C - S S D)} = 2 0 4 8 \text{ K B}$ も発生していたが、本発明では $3 2 \times 1 \text{ K B (E - S S D)} + 1 \times 6 4 \text{ K B (C - S S D)} = 9 6 \text{ K B}$ しか発生しない。

#### 【0060】

このように、本発明によれば、従来に比べて数10倍の性能向上および寿命改善がもたらされる。

#### 【0061】

第1の実施形態の効果の大きさは、図6 ( a ) 構成のS S Dユニットの総記憶容量と図6 ( b ) 構成のS S Dユニットの総記憶容量の比率に依存する。また、第2の実施形態の効果の大きさは、C - S S D 1 0 0の記憶容量とE - S S D 2 0 0の記憶容量の比率に依存する。いずれにおいても、E - S S Dの比率を大きくするほど、小単位のライトデータをE - S S Dで多く収納できるため、ストレージシステム全体としてのライト性能および書き換え寿命は向上する。しかし、ストレージシステムのコストパフォーマンス（費用対性能、費用対寿命）を考えるならば、E - S S Dの比率が大きければ大きいほどよいというものではない。

#### 【0062】

例えば、小単位のデータライトがストレージシステム500全体のユーザデータ容量の高々10%に集中する使用環境ならば、E - S S Dの比率が全体の10%程度となる構成をとるだけで十分な効果を楽しむことができる。しかし、その比率を10%以上にしても、E - S S D追加によるコストアップに見合うほどさらなる効果の増大はない。はじめに述べたように、E - S S D価格は同容量のC - S S D価格の約5倍であるから、10%分のE - S S D追加は50%のコストアップになり、ストレージシステムのドライブ原価は1.5倍となる。それでも上述の例のように数10倍の性能向上および寿命改善がもたらされる。すなわち、本発明の適用価値は、低価格なC - S S Dを主に構成した低価格帯ストレージシステムにおいて、使用環境に応じて適度な比率のE - S S Dを追加することで、コストパフォーマンスを最適化することにある。

#### 【0063】

なお、ストレージシステム500の使用環境が稼働時間の経過とともに変化し、全体のユーザデータ容量に対して小単位のデータライトが集中する領域が拡大した場合、ライトバックするアドレス区間のデータ密度は確率的に減少し、本発明による性能向上および寿命改善の効果は低下する。ここでさらにE - S S Dを追加すれば性能向上および寿命改善の効果を維持することができる。そこで、例えば、M P U P K 5 1 3、5 2 3がユーザデータアドレス空間内での小単位のデータライト回数分布を分析し、E - S S D追加により性能向上および寿命改善の効果を維持できると判断したならば、保守クライアント504を通じて、ユーザにE - S S D追加を促すメッセージを示すようにしてもよい。

#### 【0064】

なお、以上の説明においては、ライトデータをC - S S DまたはE - S S Dのどちらに書き込むかの判定基準や、ライトバックサイズの基準を64 K Bとしたが、この値はC - S S Dのメモリ管理方式により変わりうるメモリ管理単位である。よって、本発明はこの値を特定の数値に限定するものではない。C - S S Dのメモリ管理単位はそのC - S S Dの製造メーカーに問い合わせるなどして取得する。もし、メモリ管理単位を取得できない場合には、C - S S Dのライト性能テストを行い、図3 ( b ) のようなライトデータサイズと性能（または処理時間）の関係を示す特性グラフを描く。そして、性能曲線の傾きが大きく変化する点（またはライトデータサイズの減少に対して処理時間が下げ止まる点）を求めることで、C - S S Dのメモリ管理単位を推定する。この推定値をライト先判定基準

10

20

30

40

50

やライトバックサイズ基準値に適用してもよい。

【 0 0 6 5 】

なお、以上の説明においては、フラッシュメモリを記憶媒体とするストレージシステムへの実施形態を示したが、有限な書き換え寿命を持つその他の不揮発性メモリを記憶媒体とするストレージシステムについても、上述の発明を実施することが可能であり、本発明の効果を享受することは明白である。

【産業上の利用可能性】

【 0 0 6 6 】

本発明は、ストレージシステム及びそのデータライト方法に広く適用することができる。

10

【図面の簡単な説明】

【 0 0 6 7 】

【図 1】本発明の各実施形態に係わるコンシューマ向け半導体記憶装置の内部構成を示す図である。

【図 2】本発明の各実施形態に係わるエンタプライズ向け半導体記憶装置の内部構成を示す図である。

【図 3】本発明の各実施形態に係わるコンシューマ向け半導体記憶装置のデータライト処理方式とライト性能特性を示す図である。

【図 4】本発明の各実施形態に係わるエンタプライズ向け半導体記憶装置のデータライト処理方式とライト性能特性を示す図である。

20

【図 5】本発明の各実施形態に係わるストレージシステムの内部構成を示す図である

【図 6】本発明の第 1 の実施形態に係わる S S D ユニットの内部構成を示す図である。

【図 7】本発明の第 2 の実施形態に係わる S S D ユニットの内部構成を示す図である。

【図 8】本発明の実施形態に係わるストレージシステムのデータライト処理手順を示すフローチャートである。

【図 9】本発明の実施形態に係わるエンタプライズ向け半導体記憶装置にキャッシュされたデータをコンシューマ向け半導体記憶装置にライトする状況の一例を示す図である。

【図 1 0】本発明の実施形態に係わるエンタプライズ向け半導体記憶装置にキャッシュされたデータをコンシューマ向け半導体記憶装置にライトする状況のもう一例を示す図である。

30

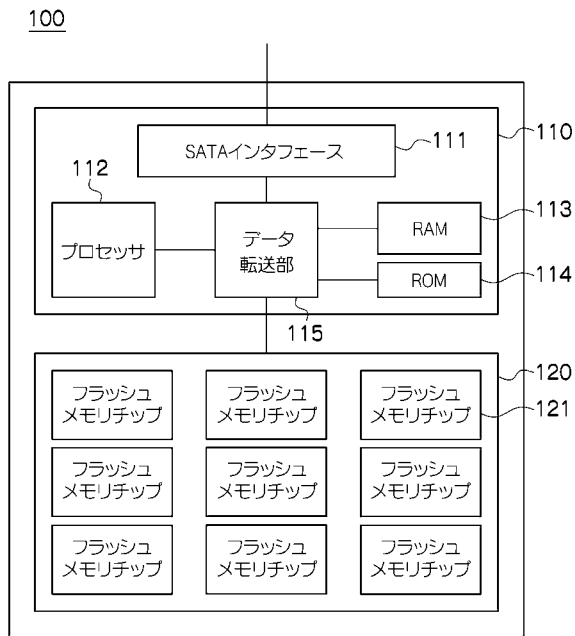
【符号の説明】

【 0 0 6 8 】

1 0 0 ... コンシューマ向け S S D、2 0 0 ... エンタプライズ向け S S D、5 0 0 ... ストレージシステム、5 4 0 ~ 5 4 3、5 5 0 ~ 5 5 3 ... S S D ユニット、7 0 0 ... S S D アダプタ

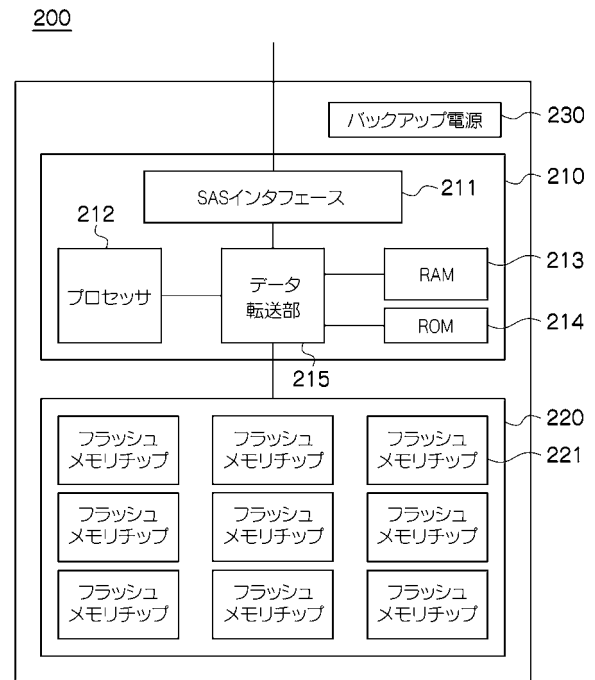
【図 1】

図1



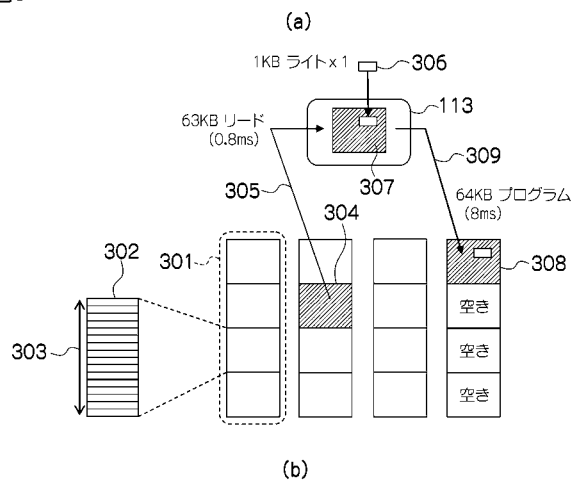
【図 2】

図2



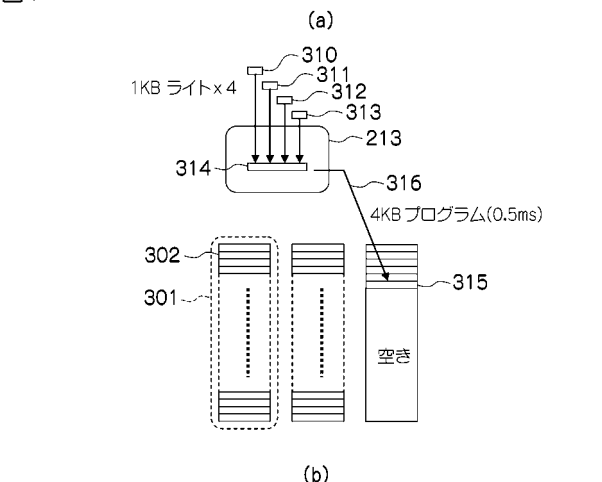
【図 3】

図3



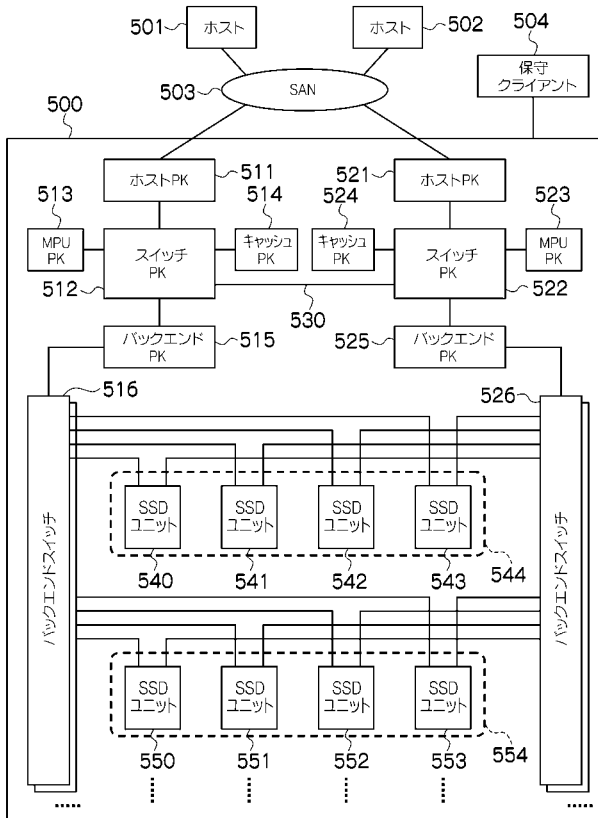
【図 4】

図4



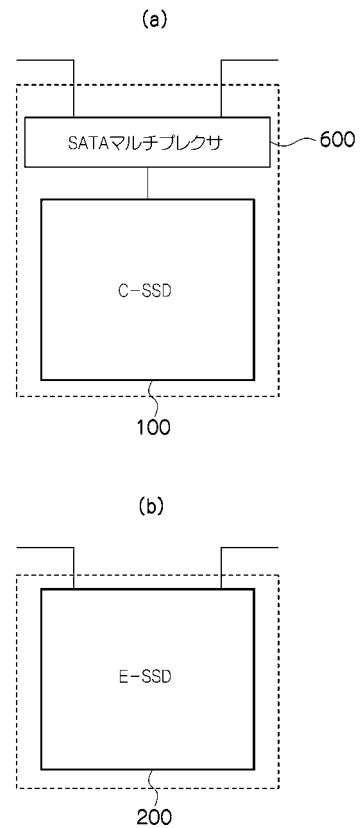
【図5】

図5



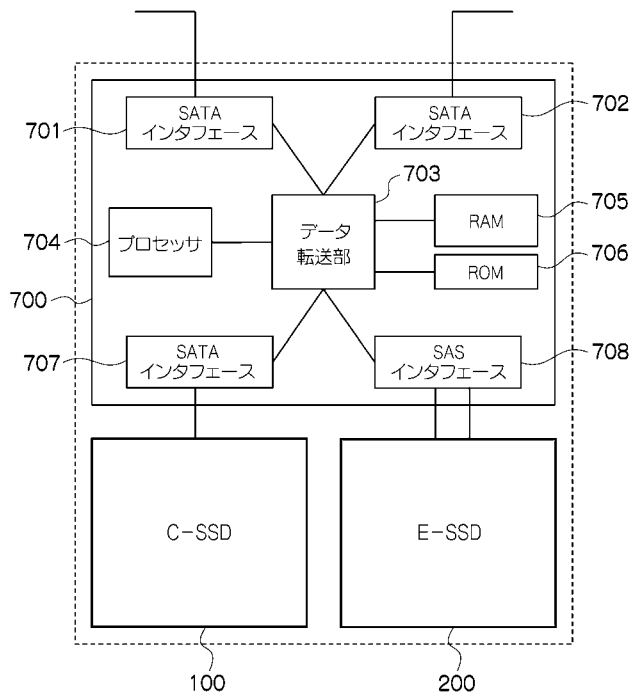
【図6】

図6



【図7】

図7



【図8】

図8

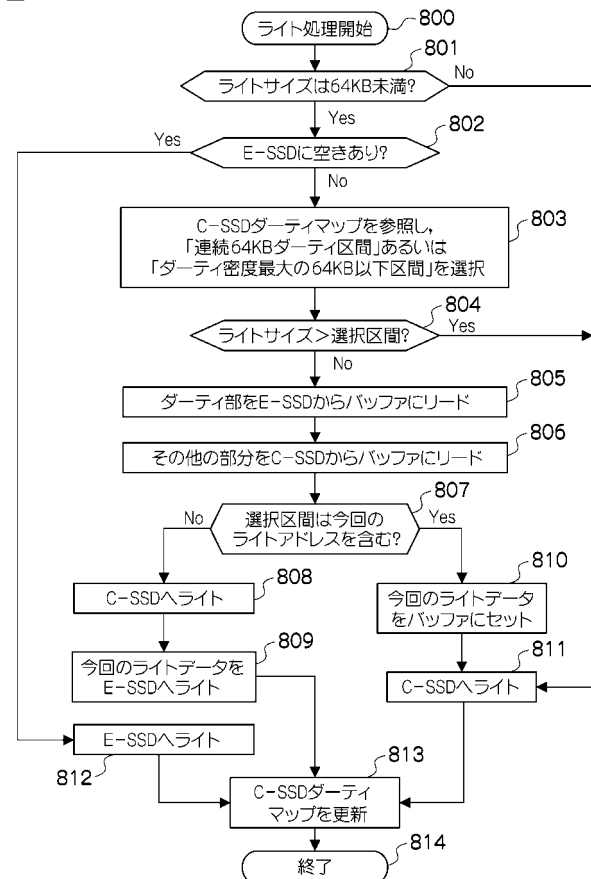




図10

