

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5221332号  
(P5221332)

(45) 発行日 平成25年6月26日 (2013. 6. 26)

(24) 登録日 平成25年3月15日 (2013. 3. 15)

|                   |                  |     |            |      |  |
|-------------------|------------------|-----|------------|------|--|
| (51) Int.Cl.      |                  | F I |            |      |  |
| <b>G06F 12/08</b> | <b>(2006.01)</b> |     | G06F 12/08 | 501F |  |
| <b>G06F 12/00</b> | <b>(2006.01)</b> |     | G06F 12/08 | 551B |  |
|                   |                  |     | G06F 12/08 | 557  |  |
|                   |                  |     | G06F 12/00 | 597U |  |

請求項の数 4 (全 32 頁)

|           |                               |           |                                       |
|-----------|-------------------------------|-----------|---------------------------------------|
| (21) 出願番号 | 特願2008-335568 (P2008-335568)  | (73) 特許権者 | 000003078<br>株式会社東芝<br>東京都港区芝浦一丁目1番1号 |
| (22) 出願日  | 平成20年12月27日 (2008. 12. 27)    | (74) 代理人  | 100089118<br>弁理士 酒井 宏明                |
| (65) 公開番号 | 特開2010-157142 (P2010-157142A) | (72) 発明者  | 矢野 浩邦<br>東京都港区芝浦一丁目1番1号 株式会社東芝内       |
| (43) 公開日  | 平成22年7月15日 (2010. 7. 15)      | (72) 発明者  | 加藤 亮一<br>東京都港区芝浦一丁目1番1号 株式会社東芝内       |
| 審査請求日     | 平成23年3月18日 (2011. 3. 18)      | (72) 発明者  | 檜田 敏克<br>東京都港区芝浦一丁目1番1号 株式会社東芝内       |

最終頁に続く

(54) 【発明の名称】 メモリシステム

(57) 【特許請求の範囲】

【請求項1】

キャッシュメモリと、

前記キャッシュメモリを介してデータが書き込まれる不揮発性半導体メモリと、

前記不揮発性半導体メモリのリソース使用量が所定値を越えている場合に、前記不揮発性半導体メモリのデータを整理してリソースを増加させる整理部と、

前記キャッシュメモリへのデータ書き込み処理の実行後、書き込み要求を待たせていない場合に、前記キャッシュメモリのリソース使用量が第1の閾値を越えかつ第1の閾値より大きい第2の閾値より小さい場合で、かつ前記整理部による整理が終了している場合は、

前記キャッシュメモリのリソース使用量が前記第1の閾値以下になるまで、前記キャッシュメモリのデータを前記不揮発性半導体メモリに追い出す第1の追い出し処理を実行し、

前記キャッシュメモリへのデータ書き込み処理の実行後、書き込み要求を待たせている場合に、前記キャッシュメモリのリソース使用量が第1の閾値を越えかつ第1の閾値より大きい第2の閾値より小さい場合は、前記第1の追い出し処理を実行せず、書き込み要求を受け付ける第1の追い出し制御部と、

前記キャッシュメモリへのデータ書き込み処理の実行後、書き込み要求を待たせていない場合に、前記キャッシュメモリのリソース使用量が前記第2の閾値を越える場合は、

前記整理部による整理が終了しているときは、前記キャッシュメモリのリソース使用量

が前記第2の閾値以下になるまで、前記キャッシュメモリのデータを前記不揮発性半導体メモリに追い出し、前記整理部による整理が終了していないときは、前記整理部による整理が終了した後、前記キャッシュメモリのリソース使用量が前記第2の閾値以下になるまで、前記キャッシュメモリのデータを前記不揮発性半導体メモリに追い出す第2の追い出し処理を実行し、

前記キャッシュメモリへのデータ書き込み処理の実行後、書き込み要求を待たせている場合に、前記キャッシュメモリのリソース使用量が前記第2の閾値を越えている場合は、書き込み要求を受け付けるよりも前に、前記第2の追い出し処理を実行する第2の追い出し制御部と、

を備えるメモリシステム。

10

【請求項2】

前記第1の追い出し制御部は、

前記第2の追い出し制御部による前記第2の追い出し処理の実行後、書き込み要求を待たせていない場合に、前記第1の追い出し処理を実行する

ことを特徴とする請求項1に記載のメモリシステム。

【請求項3】

キャッシュメモリと、

前記キャッシュメモリを介してデータが書き込まれる不揮発性半導体メモリと、

前記不揮発性半導体メモリのリソース使用量が所定値を越えている場合に、前記不揮発性半導体メモリのデータを整理してリソースを増加させる整理部と、

20

前記キャッシュメモリのリソース使用量が第1の閾値を越えかつ第1の閾値より大きい第2の閾値より小さい場合であって、かつ前記整理部による整理が終了している場合は、前記キャッシュメモリでのリソース使用量が前記第1の閾値以下になるまで、前記キャッシュメモリのデータを前記不揮発性半導体メモリに追い出す第1の追い出し制御部と、

前記キャッシュメモリのリソース使用量が前記第2の閾値を越え、かつ前記整理部による整理が終了している場合は、前記キャッシュメモリのリソース使用量が前記第2の閾値以下になるまで、前記キャッシュメモリのデータを前記不揮発性半導体メモリに追い出す第1の処理を実行し、

前記キャッシュメモリのリソース使用量が前記第2の閾値を越え、かつ前記整理部による整理が終了していない場合であって、さらに書き込み要求を処理しても前記キャッシュメモリのリソース使用量が前記第2の閾値より大きい許容限度としての最大値を越えない場合は、前記キャッシュメモリのデータを前記不揮発性半導体メモリに追い出すことなく前記キャッシュメモリへのデータ書き込みを行う第2の処理を実行し、

30

前記キャッシュメモリのリソース使用量が前記第2の閾値を越え、かつ前記整理部による整理が終了していない場合であって、さらに書き込み要求を処理すると前記最大値を越える場合は、前記整理部による整理が終了し、前記キャッシュメモリのデータを前記不揮発性半導体メモリに追い出した後、前記キャッシュメモリへのデータ書き込みを行う第3の処理を実行する、

第2の追い出し制御手段と、

を備えることを特徴とするメモリシステム。

40

【請求項4】

前記第1および第2の閾値が設定される、前記キャッシュメモリのデータを前記不揮発性半導体メモリに追い出すためのトリガとなるパラメータが、複数存在することを特徴とする請求項1乃至3の何れか一つに記載のメモリシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、不揮発性半導体メモリを備えたメモリシステムに関する。

【背景技術】

【0002】

50

コンピュータシステムに用いられる外部記憶装置として、NAND型フラッシュメモリなどのフラッシュメモリを搭載したSSD(Solid State Drive)が注目されている。フラッシュメモリは、磁気ディスク装置に比べ、高速、軽量などの利点を有している。SSD内には、複数のフラッシュメモリチップ、ホスト装置からの要求に応じて各フラッシュメモリチップのリード/ライト制御を行うコントローラ、各フラッシュメモリチップとホスト装置との間でデータ転送を行うためのバッファメモリ、電源回路、ホスト装置に対する接続インタフェースなどを備えている(例えば、特許文献1)。

【0003】

NAND型フラッシュメモリのように、不揮発性半導体記憶素子には、データを記憶させる場合に、ブロックと呼ばれる単位で一度データを消去してからその後に書き込みを行うものがあったり、ページと呼ばれる単位で読み出し/書き込みを行うものがあったり、消去/読み出し/書き込みの単位が固定されていたりするものがある。一方、パーソナルコンピュータなどのホスト機器がハードディスクをはじめとする二次記憶装置に対してデータの読み出し/書き込みを行う単位は、セクタと呼ばれる。セクタは、半導体記憶素子の消去/読み出し/書き込みの単位とは独立に定められており、通常、ブロック、ページ、セクタのサイズは、ブロック>ページ>セクタという関係にある。

10

【0004】

このように、半導体記憶素子の消去/読み出し/書き込みの単位は、ホスト機器の読み出し/書き込みの単位よりも大きい場合がある。このような半導体記憶素子を用いてハードディスクのようなパーソナルコンピュータの二次記憶装置を構成する場合、ホスト機器としてのパーソナルコンピュータからの小さなサイズのデータは、半導体記憶素子のブロックサイズ、ページサイズに適合させてアドレス変換を行う必要がある。

20

【0005】

また、このようなフラッシュメモリを用いて大容量の二次記憶装置を構成する場合においては、特許文献2に示されるように、フラッシュメモリとホスト装置との間に、キャッシュメモリを介在させて、フラッシュメモリでの書き込み回数(消去回数)を減らすように構成されていることが多い。キャッシュメモリにホスト装置からの書き込みが発生した際に、キャッシュメモリが満杯の場合は、キャッシュメモリからフラッシュメモリへのデータ追い出しを行ってから、キャッシュメモリにデータを書き込むことになる。しかしながら、キャッシュメモリがほぼ満杯になってから上記データ追い出しを行うのでは、データ追い出しを行う間、ホスト機器からの書き込み要求を待たせることになり、ホスト機器側から見て応答性のよい二次記憶装置を構成することができない。

30

【0006】

また、上述したように、データの消去単位(ブロック)と、データの管理単位が異なる場合、フラッシュメモリの書き換えが進むと、無効な(最新ではない)データによって、ブロックは穴あき状態になる。このような穴あき状態のブロックが増えると、実質的に使用可能なブロックが少なくなり、フラッシュメモリの記憶領域を有効利用できないので、有効な最新のデータを集めて、違うブロックに書き直すコンパクションと呼ばれるフラッシュメモリの整理処理が行われる。

【0007】

しかしながら、従来のキャッシュメモリの追い出し処理では、フラッシュメモリ側の整理状態を考慮していないので、フラッシュメモリ側の整理が進んでいない場合は、フラッシュメモリへの書き込みに時間がかかり、結果的にホスト側の書き込みコマンドに対する応答性が低下する。

40

【0008】

【特許文献1】特許第3688835号公報

【特許文献2】特表2007-528079号

【発明の開示】

【発明が解決しようとする課題】

【0009】

50

本発明は、ホストからの書き込み要求に対する応答性を全般的に向上させることが可能なメモリシステムを提供することを目的とする。

【課題を解決するための手段】

【0010】

本願発明の一態様によれば、揮発性のキャッシュメモリとしての第1の記憶部と、不揮発性の第2の記憶部と、ホスト装置からのデータを前記第1の記憶部を介して前記第2の記憶部に書き込むコントローラとを備えるメモリシステムにおいて、前記コントローラは、前記第2の記憶部でのリソース使用量が所定値を越えている場合に第2の記憶部のデータを整理してリソースを増加させる整理手段と、前記第1の記憶部のリソース使用量が第1の閾値を越えかつ第1の閾値より大きい第2の閾値より小さい場合であって、かつ前記整理手段による整理が終了している場合は、第1の記憶部でのリソース使用量が第1の閾値を越えなくなるまで、第1の記憶部のデータを第2の記憶部に追い出す第1の追い出し制御手段と、前記第1の記憶部のリソース使用量が第2の閾値を越え、かつ前記整理手段による整理が終了している場合は、第1の記憶部のリソース使用量が第2の閾値を越えなくなるまで、第1の記憶部のデータを第2の記憶部に追い出し、前記第1の記憶部のリソース使用量が第2の閾値を越え、かつ前記整理手段による整理が終了していない場合は、前記整理手段による整理が終了した後、第1の記憶部のデータを第2の記憶部に追い出す第2の制御手段とを備えることを特徴とする。

10

【発明の効果】

【0011】

本発明によれば、ホストからの書き込み要求に対する応答性を全般的に向上させることが可能なメモリシステムを提供できる。

20

【発明を実施するための最良の形態】

【0012】

以下に添付図面を参照して、この発明にかかるメモリシステムの実施の形態を詳細に説明する。

【0013】

以下、本発明の実施の形態について図面を参照して説明する。なお、以下の説明において、同一の機能および構成を有する要素については、同一符号を付し、重複説明は必要な場合にのみ行う。

30

【0014】

先ず、本明細書で用いる用語について定義しておく。

・物理ページ：NAND型フラッシュメモリ内部において一括して書き込み/読み出しが可能な単位のこと。

・論理ページ：SSD内部で設定される書き込み/読み出し単位であり、1以上の物理ページを組み合わせて構成される。

・物理ブロック：NAND型フラッシュメモリ内部において独立して消去可能な最小単位のことであり、複数の物理ページから構成される。

・論理ブロック：SSD内部で設定される消去単位であり、1以上の物理ブロックを組み合わせて構成される。論理ブロックは、複数の論理ページから構成される。

40

・セクタ：ホストからの最小アクセス単位のこと。セクタサイズは、例えば512B。

・クラスタ：SSD内部で「小さなデータ」を管理する管理単位。クラスタサイズはセクタサイズ以上であり、ホストのOSが採用するファイルシステムのデータ管理単位、または、論理ページサイズと等しくなるように定められる。例えば、クラスタサイズの2以上の自然数倍が論理ページサイズとなるように定められてもよい。

・トラック：SSD内部で「大きなデータ」を管理する管理単位。クラスタサイズの2以上の自然数倍がトラックサイズとなるように定められる。例えば、トラックサイズが論理ブロックサイズと等しくなるように定められてもよい。

・フリーブロック(FB)：内部に有効データを含まない、用途未割り当ての論理ブロックのこと。以下の、CFB、FFBの2種類がある。

50

・コンプリートフリーブロック (CFB) : 再利用のために消去動作を行う必要があるFBのこと。消去動作の実行後は、論理ブロックの先頭に位置する論理ページから書き込むことが可能である。

・フラグメントフリーブロック (FFB) : 未書き込みの論理ページが残っており、消去動作を実行することなく再利用が可能なFBのこと。残りの未書き込み状態のままの論理ページに書き込むことが可能である。

・バッドブロック (BB) : NAND型フラッシュメモリ上の、誤りが多いなど記憶領域として使用できない物理ブロックのこと。例えば、消去動作が正常に終了しなかった物理ブロックがバッドブロックBBとして登録される。

・書き込み効率 : 所定期間内における、ホストから書き込んだデータ量に対する、論理ブロックの消去量の統計値のこと。小さいほどNAND型フラッシュメモリの消耗度が小さい。

・有効クラスタ : 論理アドレスに対応するクラスタサイズの最新データ。

・無効クラスタ : 同一論理アドレスのデータが他の場所に書きこまれ、参照されることがなくなったクラスタサイズのデータ。

・有効トラック : 論理アドレスに対応するトラックサイズの最新データ。

・無効トラック : 同一論理アドレスのデータが他の場所に書きこまれ、参照されることがなくなったトラックサイズのデータ。

・MLC (Multi Level Cell) モード : 多値記憶が可能なNAND型フラッシュメモリにおいて、通常通り、上位ページおよび下位ページを使用して書き込みを行うモード。MLCモードで使用する1以上の物理ブロックを組み合わせ、MLCモードの論理ブロックが構成される。

・擬似SLC (Single Level Cell) モード : 多値記憶が可能なNAND型フラッシュメモリにおいて、下位ページのみを使用して書き込みを行うモード。擬似SLCモードで使用する1以上の物理ブロックを組み合わせ、擬似SLCモードの論理ブロックが構成される。一度擬似SLCモードで使用した物理ブロックであっても、消去動作後はMLCモードで使用する事が可能である。

[ 第1の実施形態 ]

【 0015 】

図1は、SSD (Solid State Drive) 100の構成例を示すブロック図である。SSD 100は、ATAインターフェース (ATA I/F) 2などのメモリ接続インターフェースを介してパーソナルコンピュータあるいはCPUコアなどのホスト装置 (以下、ホストと略す) 1と接続され、ホスト1の外部メモリとして機能する。また、SSD 100は、RS232Cインターフェース (RS232C I/F) などの通信インターフェース3を介して、デバッグ用 / 製造検査用機器200との間でデータを送受信することができる。SSD 100は、不揮発性半導体メモリとしてのNAND型フラッシュメモリ (以下、NANDメモリと略す) 10と、コントローラとしてのドライブ制御回路4と、揮発性半導体メモリとしてのDRAM20と、電源回路5と、状態表示用のLED6と、ドライブ内部の温度を検出する温度センサ7と、フューズ8とを備えている。

【 0016 】

電源回路5は、ホスト1側の電源回路から供給される外部直流電源から複数の異なる内部直流電源電圧を生成し、これら内部直流電源電圧をSSD 100内の各回路に供給する。また、電源回路5は、外部電源の立ち上がりを検知し、パワーオンリセット信号を生成して、ドライブ制御回路4に供給する。フューズ8は、ホスト1側の電源回路とSSD 100内部の電源回路5との間に設けられている。外部電源回路から過電流が供給された場合フューズ8が切断され、内部回路の誤動作を防止する。

【 0017 】

NANDメモリ10は、この場合、4並列動作を行う4つの並列動作要素10a~10dを有し、4つの並列動作要素10a~10dは、4つのチャンネル (ch0~ch3) によってドライブ制御回路4に接続されている。各並列動作要素10a~10dは、バンク

10

20

30

40

50

インターリーブが可能な複数のバンク（この場合、4バンク、Bank 0 ~ Bank 3）によって構成されており、各バンクは、複数のNANDメモリチップ（この場合、2メモリチップ、Chip 0、Chip 1）によって構成されている。各メモリチップは、例えば、それぞれ複数の物理ブロックを含むプレーン0、プレーン1の2つの領域（District）に分割されている。プレーン0およびプレーン1は、互いに独立した周辺回路（例えば、ロウデコーダ、カラムデコーダ、ページバッファ、データキャッシュ等）を備えており、倍速モードを使用することで、同時に消去/書き込み/読み出しを行うことが可能である。このように、NANDメモリ10の各NANDメモリチップは、複数のチャンネルによる並列動作、複数のバンクによるバンクインターリーブ動作、複数のプレーンを用いた倍速モードによる並列動作が可能である。なお、各メモリチップは、4つのプレーンに分割された構成であってもよいし、あるいは、全く分割されていなくてもよい。

10

#### 【0018】

DRAM 20は、ホスト1とNANDメモリ10間でのデータ転送用キャッシュおよび作業領域用メモリなどとして機能する。DRAM 20の作業領域用メモリに記憶されるものとしては、NANDメモリ10に記憶されている各種管理テーブル（後述する）が起動時などに展開されたマスターテーブル（スナップショット）、管理テーブルの変更差分であるログ情報などがある。DRAM 20の代わりに、FeRAM（Ferroelectric Random Access Memory）、MRAM（Magnetoresistive Random Access Memory）、PRAM（Phase change Random Access Memory）などを使用しても良い。ドライブ制御回路4は、ホスト1とNANDメモリ10との間でDRAM 20を介してデータ転送制御を行うとともに、SSD 100内の各構成要素を制御する。また、ドライブ制御回路4は、状態表示用LED 6にステータス表示用信号を供給するとともに、電源回路5からのパワーオンリセット信号を受けて、リセット信号およびクロック信号を自回路内およびSSD 100内の各部に供給する機能も有している。

20

#### 【0019】

各NANDメモリチップは、データ消去の単位である物理ブロックを複数配列して構成されている。図2(a)は、NANDメモリチップに含まれる1個の物理ブロックの構成例を示す等価回路図である。各物理ブロックは、X方向に沿って順に配列された $(p+1)$ 個のNANDストリングを備えている（ $p$ は、0以上の整数）。 $(p+1)$ 個のNANDストリングにそれぞれ含まれる選択トランジスタST 1は、ドレインがビット線BL 0 ~ BL  $p$ に接続され、ゲートが選択ゲート線SGDに共通接続されている。また、選択トランジスタST 2は、ソースがソース線SLに共通接続され、ゲートが選択ゲート線SGSに共通接続されている。

30

#### 【0020】

各メモリセルトランジスタMTは、半導体基板上に形成された積層ゲート構造を備えたMOSFET（Metal Oxide Semiconductor Field Effect Transistor）から構成される。積層ゲート構造は、半導体基板上にゲート絶縁膜を介して形成された電荷蓄積層（浮遊ゲート電極）、および電荷蓄積層上にゲート間絶縁膜を介して形成された制御ゲート電極を含んでいる。メモリセルトランジスタMTは、浮遊ゲート電極に蓄えられる電子の数に応じて閾値電圧が変化し、この閾値電圧の違いに応じてデータを記憶する。メモリセルトランジスタMTは、1ビットを記憶するように構成されていてもよいし、多値（2ビット以上のデータ）を記憶するように構成されていてもよい。

40

#### 【0021】

また、メモリセルトランジスタMTは、浮遊ゲート電極を有する構造に限らず、MONOS（Metal-Oxide-Nitride-Oxide-Silicon）型など、電荷蓄積層としての窒化膜界面に電子をトラップさせることでしきい値調整可能な構造であってもよい。MONOS構造のメモリセルトランジスタMTについても同様に、1ビットを記憶するように構成されていてもよいし、多値（2ビット以上のデータ）を記憶するように構成されていてもよい。

#### 【0022】

各NANDストリングにおいて、 $(q+1)$ 個のメモリセルトランジスタMTは、選択

50

トランジスタST1のソースと選択トランジスタST2のドレインとの間に、それぞれの電流経路が直列接続されるように配置されている。すなわち、複数のメモリセルトランジスタMTは、隣接するもの同士で拡散領域（ソース領域若しくはドレイン領域）を共有するような形でY方向に直列接続される。

【0023】

そして、最もドレイン側に位置するメモリセルトランジスタMTから順に、制御ゲート電極がワード線WL0~WLqにそれぞれ接続されている。従って、ワード線WL0に接続されたメモリセルトランジスタMTのドレインは選択トランジスタST1のソースに接続され、ワード線WLqに接続されたメモリセルトランジスタMTのソースは選択トランジスタST2のドレインに接続されている。

10

【0024】

ワード線WL0~WLqは、物理ブロック内のNANDストリング間で、メモリセルトランジスタMTの制御ゲート電極を共通に接続している。つまり、ブロック内において同一行にあるメモリセルトランジスタMTの制御ゲート電極は、同一のワード線WLに接続される。この同一のワード線WLに接続される(p+1)個のメモリセルトランジスタMTは1ページ（物理ページ）として取り扱われ、この物理ページごとにデータの書き込みおよびデータの読み出しが行われる。

【0025】

また、ビット線BL0~BLpは、ブロック間で、選択トランジスタST1のドレインを共通に接続している。つまり、複数のブロック内において同一列にあるNANDストリングは、同一のビット線BLに接続される。

20

【0026】

図2(b)は、例えば、1個のメモリセルトランジスタMTに2ビットの記憶を行う4値データ記憶方式でのしきい値分布を示す模式図である。4値データ記憶方式では、上位ページデータ“x”と下位ページデータ“y”で定義される4値データ“xy”の何れか1つをメモリセルトランジスタMTに保持可能である。

【0027】

この、4値データ“xy”は、メモリセルトランジスタMTのしきい値電圧の順に、例えば、データ“11”、“01”、“00”、“10”が割り当てられる。データ“11”は、メモリセルトランジスタMTのしきい値電圧が負の消去状態である。なお、データの割り当て規則はこれに限らない。また、1個のメモリセルトランジスタMTに3ビット以上の記憶を行う構成であってもよい。

30

【0028】

下位ページ書き込み動作においては、データ“11”（消去状態）のメモリセルトランジスタMTに対して選択的に、下位ビットデータ“y”の書き込みによって、データ“10”が書き込まれる。上位ページ書き込み前のデータ“10”のしきい値分布は、上位ページ書き込み後のデータ“01”とデータ“00”のしきい値分布の中間程度に位置しており、上位ページ書き込み後のしきい値分布よりブロードであってもよい。上位ページ書き込み動作においては、データ“11”のメモリセルと、データ“10”のメモリセルに対して、それぞれ選択的に上位ビットデータ“x”の書き込みが行われて、データ“01”およびデータ“00”が書き込まれる。擬似SLCモードでは、下位ページのみを使用して書き込みを行う。下位ページの書き込みは、上位ページの書き込みに比べて高速である。

40

【0029】

図3は、ドライブ制御回路4のハードウェア的な内部構成例を示すブロック図である。ドライブ制御回路4は、データアクセス用バス101、第1の回路制御用バス102、および第2の回路制御用バス103を備えている。第1の回路制御用バス102には、ドライブ制御回路4全体を制御するプロセッサ104が接続されている。第1の回路制御用バス102には、NANDメモリ10に記憶された各管理プログラム（FW：ファームウェア）をブートするブート用プログラムが格納されたブートROM105がROMコントロ

50

ーラ106を介して接続されている。また、第1の回路制御用バス102には、図1に示した電源回路5からのパワーオンリセット信号を受けて、リセット信号およびクロック信号を各部に供給するクロックコントローラ107が接続されている。

#### 【0030】

第2の回路制御用バス103は、第1の回路制御用バス102に接続されている。第2の回路制御用バス103には、図1に示した温度センサ7からのデータを受けるためのI<sup>2</sup>C回路108、状態表示用LED6にステータス表示用信号を供給するパラレルI/O (PIO)回路109、RS232C I/F3を制御するシリアルI/O (SIO)回路110が接続されている。

#### 【0031】

A T Aインタフェースコントローラ ( A T Aコントローラ ) 111、第1のE C C ( E r r o r C h e c k i n g a n d C o r r e c t i o n ) 回路112、N A N Dコントローラ113、およびD R A Mコントローラ114は、データアクセス用バス101と第1の回路制御用バス102との両方に接続されている。A T Aコントローラ111は、A T Aインタフェース2を介してホスト1との間でデータを送受信する。データアクセス用バス101には、データ作業領域およびファームウェア展開領域として使用されるS R A M 115がS R A Mコントローラ116を介して接続されている。N A N Dメモリ10に記憶されているファームウェアは起動時、ブートR O M 105に記憶されたブート用プログラムによってS R A M 115に転送される。

#### 【0032】

N A N Dコントローラ113は、N A N Dメモリ10とのインタフェース処理を行うN A N D I / F 117、第2のE C C回路118、およびN A N Dメモリ10 - D R A M 20間のアクセス制御を行うD M A転送制御用D M Aコントローラ119を備えている。第2のE C C回路118は第2の訂正符号のエンコードを行い、また、第1の誤り訂正符号のエンコードおよびデコードを行う。第1のE C C回路112は、第2の誤り訂正符号のデコードを行う。第1の誤り訂正符号、第2の誤り訂正符号は、例えば、ハミング符号、B C H ( B o s e C h a u d h u r i H o c q e n g h e m ) 符号、R S ( R e e d S o l o m o n ) 符号、或いはL D P C ( L o w D e n s i t y P a r i t y C h e c k ) 符号等であり、第2の誤り訂正符号の訂正能力は、第1の誤り訂正符号の訂正能力よりも高いとする。

#### 【0033】

図1に示したように、N A N Dメモリ10においては、4つの並列動作要素10a ~ 10dが各複数ビットの4チャンネル ( 4 c h ) を介して、ドライブ制御回路4内部のN A N Dコントローラ112に並列接続されており、4つの並列動作要素10a ~ 10dを並列動作させることが可能である。また、各チャンネルのN A N Dメモリ10は、バンクインターリーブが可能な4つのバンクに分割されており、各メモリチップのプレーン0およびプレーン1に対しても、同時にアクセスを行うことが可能である。したがって、1チャンネルにつき、最大8物理ブロック ( 4バンク x 2プレーン )、ほぼ同時に書き込みなどの処理を実行可能である。

#### 【0034】

図4は、プロセッサ104により実現されるファームウェアの機能構成例を示すブロック図である。プロセッサ104により実現されるファームウェアの各機能は、大きく、データ管理部120、A T Aコマンド処理部121、セキュリティ管理部122、ブートローダ123、初期化管理部124、デバッグサポート部125に分類される。

#### 【0035】

データ管理部120は、N A N Dコントローラ113、第1のE C C回路112を介して、N A N Dメモリ10 - D R A M 20間のデータ転送、N A N Dメモリ10に関する各種機能を制御する。A T Aコマンド処理部121は、A T Aコントローラ111、およびD R A Mコントローラ114を介して、データ管理部120と協働してD R A M 20 - ホスト1間のデータ転送処理を行う。セキュリティ管理部122は、データ管理部120およびA T Aコマンド処理部121と協働して各種のセキュリティ情報を管理する。

10

20

30

40

50

## 【 0 0 3 6 】

ブートローダ 1 2 3 は、パワーオン時、各管理プログラム（ファームウェア）を N A N D メモリ 1 0 から S R A M 1 1 5 にロードする。初期化管理部 1 2 4 は、ドライブ制御回路 4 内の各コントローラ / 回路の初期化を行う。デバッグサポート部 1 2 5 は、外部から R S 2 3 2 C インタフェースを介して供給されたデバッグ用データを処理する。主に、データ管理部 1 2 0、A T A コマンド処理部 1 2 1、およびセキュリティ管理部 1 2 2 が、S R A M 1 1 5 に記憶される各管理プログラムをプロセッサ 1 0 4 が実行することによって実現される機能部である。

## 【 0 0 3 7 】

本実施形態では、主としてデータ管理部 1 2 0 が実現する機能について説明する。データ管理部 1 2 0 は、A T A コマンド処理部 1 2 1 が記憶デバイスである N A N D メモリ 1 0 や D R A M 2 0 に対して要求する機能の提供（ホストからの Write 要求、Cache Flush 要求、Read 要求等の各種コマンドへの応答）と、ホスト 1 から与えられる論理アドレスと N A N D メモリ 1 0 との対応関係の管理と、スナップショット、ログによる管理情報の保護と、D R A M 1 0 および N A N D メモリ 1 0 を利用した高速で効率の良いデータ読み出し / 書き込み機能の提供と、N A N D メモリ 1 0 の信頼性の確保などを行う。

## 【 0 0 3 8 】

図 5 は、N A N D メモリ 1 0 および D R A M 2 0 内に形成された機能ブロックを示すものである。ホスト 1 と N A N D メモリ 1 0 との間には、D R A M 2 0 上に構成されたライトキャッシュ（W C）2 1 およびリードキャッシュ（R C）2 2 が介在している。W C 2 1 はホスト 1 からの Write データを一時保存し、R C 2 2 は N A N D メモリ 1 0 からの Read データを一時保存する。N A N D メモリ 1 0 内のブロックは、書き込み時の N A N D メモリ 1 0 に対する消去の量を減らすために、データ管理部 1 2 0 により、前段ストレージ領域（F S : Front Storage）1 2、中段ストレージ領域（I S : Intermediate Storage）1 3 およびメインストレージ領域（M S : Main Storage）1 1 という各管理領域に割り当てられている。F S 1 2 は、W C 2 1 からのデータを「小さな単位」であるクラスタ単位に管理するものであり、小データを短期間保存する。I S 1 3 は、F S 1 2 から溢れたデータを「小さな単位」であるクラスタ単位に管理するものであり、小データを長期間保存する。M S 1 1 は、W C 2 1、F S 1 2、I S 1 3 からのデータを「大きな単位」であるトラック単位で管理する。

## 【 0 0 3 9 】

つぎに、図 5 の各構成要素の具体的な機能構成について詳述する。ホスト 1 は S S D 1 0 0 対し、Read または Write する際には、A T A インタフェースを介して論理アドレスとしての L B A（Logical Block Addressing）を入力する。L B A は、図 6 に示すように、セクタ（サイズ：5 1 2 B）に対して 0 からの通し番号をつけた論理アドレスである。本実施の形態においては、図 5 の各構成要素である W C 2 1、R C 2 2、F S 1 2、I S 1 3、M S 1 1 の管理単位として、L B A の下位（s + 1）ビット目から上位のビット列で構成されるクラスタアドレスと、L B A の下位（s + t + 1）ビットから上位のビット列で構成されるトラックアドレスとを定義する。この実施の形態では、トラックと論理ブロックのサイズは同じとする。論理ブロックとは、N A N D メモリ 1 0 のチップ上の物理ブロックを複数組み合わせることであり、この実施の形態では、論理ブロックは 1 つの物理ブロックを並列チャネル数分（この場合、図 1 に示すように 4 c h）まとめた単位のことをいう。論理ページも同様であり、物理ページを 4 c h 分まとめた単位のことをいう。また、論理ブロックは、バンクインターリーブを有効利用するため、同じバンクに属する物理ブロックから選択される。

## 【 0 0 4 0 】

・リードキャッシュ（R C）2 2

R C 2 2 は、ホスト 1 からの Read 要求に対して、N A N D メモリ 1 0（F S 1 2、I S 1 3、M S 1 1）からの Read データを一時的に保存するための領域である。ホスト 1 へのデータ転送は、基本的に、R C 2 2 から行う。なお、W C 2 1 から N A N D メモリ 1 0 へ

10

20

30

40

50

のデータの書き込みを行う際には、同一論理アドレスの R 2 2 上のデータを無効にする。

【 0 0 4 1 】

・ライトキャッシュ ( W C ) 2 1

W C 2 1 は、ホスト 1 からの Write 要求に対して、ホスト 1 からの Write データを一時的に保存するための領域である。W C 2 1 上のデータは、クラスタ単位で管理し、書き込みと有効データの管理はセクタ単位で行う。W C 2 1 のリソースが不足した場合、W C 2 1 の記憶データを N A N D 1 0 に追い出す。ホスト 1 から R C 2 2 上のデータと同一の論理アドレスに対する書き込みが行われた場合、その最新データは W C 2 1 上に保存される。そのため、同一の論理アドレスに対応するデータが、W C 2 1、R C 2 2、N A N D メモリ 1 0 上にある場合には、データの新しさは、W C 2 1、R C 2 2、N A N D メモリ 1 0 10  
の順となるため、ホスト 1 に返すデータも W C 2 1 上のデータを優先する。

【 0 0 4 2 】

・メインストレージ領域 ( M S ) 1 1

M S 1 1 はトラック単位でデータの管理が行われ、ほとんどのユーザデータが格納される。W C 2 1 上で有効クラスタの多いトラック ( 高密度トラック ) は、W C 1 2 から直接 M S 1 1 に書き込まれる。その他、M S 1 1 には、F S 1 2、I S 1 3 で管理しきれなくなったデータが入力される。M S 1 1 に入力されたトラックと同一 L B A のトラックについては、論理ブロック内で無効化し、この論理ブロックを解放する。M S 1 1 に入力されたトラックと同一 L B A のトラックに属するクラスタについては、論理ブロック内で無効化し、論理ブロック内の全クラスタが無効になった論理ブロックは解放する。M S 1 1 は、M L C モードの複数の論理ブロックで構成される。この実施の形態では、トラックと論理ブロックのサイズは同じとしているので、F S 1 2 や I S 1 3 で行われる追記処理や、I S 1 3 で行われるコンパクション ( 有効クラスタのみを集めて新しい論理ブロックを作り、無効なクラスタ部分を解放する処理 ) は不要となる。もしトラックサイズが論理ブロックサイズよりも小さい場合は、F S 1 2 や I S 1 3 で行われる追記処理や、I S 1 3 で行われるコンパクションを適用してもよい。20

【 0 0 4 3 】

・前段ストレージ領域 ( F S ) 1 2

F S 1 2 はクラスタ単位でデータを管理される F I F O 構造のバッファであり、入力は複数のクラスタをまとめた論理ページ単位で行われる。F S 1 2 には、W C 2 1 上で有効クラスタ数の少ないトラック ( 低密度トラック ) が最初に書き込まれる。データの書き込み順序で論理ブロックが並んだ F I F O 構造となっている。F S 1 2 に存在するクラスタと同一 L B A のクラスタが F S 1 2 に入力された場合、F S 1 2 内のクラスタを無効化するだけでよく、書き換え動作を伴わない。F S 1 2 に入力されたクラスタと同一 L B A のクラスタについては、論理ブロック内で無効化し、論理ブロック内の全クラスタが無効になった論理ブロックは解放する。F S 1 1 の F I F O 構造の最後まで到達した論理ブロックに格納されたクラスタは、ホスト 1 から再書き込みされる可能性の低いクラスタとみなし、論理ブロックごと I S 1 3 の管理下に移動する。F S 1 2 は、この実施の形態では、書き込みの高速化を図るため擬似 S L C モードの複数の論理ブロックで構成される。なお、F S 1 2 は、M L C モードの複数の論理ブロックで構成されてもよい。更新頻度の高いデータは F S 1 2 を通過している最中に無効化され、更新頻度の低いデータだけが F S 1 2 から溢れていくため、更新頻度の高いデータと低いデータとを F S 1 2 で選り分けることができる。これにより、後段の I S 1 3 でコンパクションが頻繁に発生する可能性を低減させることが可能である。30

【 0 0 4 4 】

・中段ストレージ領域 ( I S ) 1 3

I S 1 3 は、再書き込みされる可能性の低いクラスタを格納するためのバッファであり、F S 1 3 と同様にクラスタ単位でデータの管理が行われる。I S 1 3 に存在するクラスタと同一 L B A のクラスタが F S 1 2、I S 1 3 に入力された場合、I S 1 3 内のクラスタを無効化するだけでよく、書き換え動作を伴わない。I S 1 3 においては、F S 1 2 と40  
50

同様、データの書き込まれた順序（F S 1 2 から移動された順序）が古い論理ブロックから並んだリスト構造をとるが、コンパクションを行う点がF S 1 2 と異なる。I S 1 3 の容量や管理テーブルの都合で飽和した場合は、コンパクション（I S 1 3 から有効クラスタを集めてI S 1 3 へ書き戻すこと）やデフラグ（F S 1 2 およびI S 1 3 のクラスタをトラックに統合して、M S 1 1 へ追い出すこと）を行う。I S 1 3 は、この実施の形態では、M L C モードの論理ブロックと擬似S L C モードの論理ブロックの混在で構成される。すなわち、F S 1 2 からI S 1 3 に移動されるブロックは擬似S L C モードの論理ブロックであるが、I S 1 3 内でコンパクションする際に、M L C モードの論理ブロックに書き直す。なお、F S 1 2 がM L C モードの論理ブロックで構成される場合は、I S 1 3 もM L C モードの論理ブロックのみで構成されることになる。

10

## 【 0 0 4 5 】

図7は、データ管理部120が図5に示した各構成要素を制御管理するための管理テーブルを示すものである。D R A M 2 0 を管理するためのテーブルとしては、R C 管理テーブル23、W C トラックテーブル24、W C トラック情報テーブル25、W C 高密度トラック情報テーブル26、W C 低密度トラック情報テーブル27などがある。N A N D メモリ10を管理するためのテーブルとしては、トラックテーブル30、クラスタディレクトリテーブル31、クラスタテーブル32、クラスタブロック情報テーブル33、論物変換テーブル40などがある。N A N D メモリ10を管理するためのテーブルは、正引きアドレス変換で参照するテーブル、逆引きアドレス変換で参照するテーブルに分けられる。正引きアドレス変換とは、データのL B A から実際にデータが記憶されている論理ブロック

20

## 【 0 0 4 6 】

## ・ R C 管理テーブル23

R C 管理テーブル23は、N A N D メモリ10からR C 2 2 に転送されたデータを管理するためのものである。

## 【 0 0 4 7 】

## ・ W C トラックテーブル24

W C 2 1 上に記憶されたデータに関するW C トラック情報をL B A からルックアップするためのハッシュテーブルであり、L B A のトラックアドレスのL S B 数ビットをインデックスとし、インデックス毎に複数のエントリ（タグ）を有する。各タグには、L B A トラックアドレスと該トラックアドレスに対応するW C トラック情報へのポインタが記憶されている。

30

## 【 0 0 4 8 】

## ・ W C トラック情報テーブル25

W C トラック情報テーブル25には、アクセスのあったW C トラック情報の新旧の順序をL R U (Least Recently used) で双方向リストで管理するためのW C トラックL R U 情報テーブル25 a と、空いているW C トラック情報の番号を管理するW C トラック空き情報テーブル25 b とがある。W C 2 1 からN A N D にデータを追い出すときに、W C トラックL R U 情報テーブル25 a を用いて最も古くにアクセスされたトラックを取り出す。

40

## 【 0 0 4 9 】

W C トラック情報は、W C 2 1 内に存在する複数のトラックの1つに対応する。

W C トラック情報には、

(1) W C 2 1 内に存在するトラックアドレス、トラック内のW C 2 1 上の有効クラスタの個数、各クラスタが有効であるかどうかの情報、各クラスタがW C 2 1 のどこに存在するかを示すW C 内クラスタ位置情報、

(2) 1 クラスタに含まれる複数のセクタのうちどのセクタに有効なデータを保持しているかを示す情報（セクタビットマップ）、

(3) トラックの状態情報（有効、無効、A T A からのデータ転送中、N A N D に書き込

50

み中など)

などが含まれている。なお、上記のWCトラック情報では、有効クラスタが存在する記憶位置で自トラック内に存在するクラスタアドレスのLSB(t)ビットを管理するようにしたが、クラスタアドレスの管理方法は任意であり、例えば、自トラック内に存在するクラスタアドレスのLSB(t)ビット自体を管理するようにしてもよい(図6参照)。

【0050】

・WC高密度トラック情報テーブル26

MS11に書き込むことになる高密度(トラック内で有効クラスタ数が所定パーセント以上)のトラック情報を管理するためのもので、高密度トラックに関するWCトラック情報とその個数を管理している。

10

【0051】

・WC低密度トラック情報テーブル27

FS12に書き込むことになる低密度(トラック内で有効クラスタ数が所定パーセント未満)のトラック情報を管理するためのもので、低密度トラックのクラスタ数の合計を管理している。

【0052】

・トラックテーブル30(正引き)

LBAのトラックアドレスからトラック情報を取得するためのテーブルである。トラック情報としては、

(1)論理ブロックアドレス(トラックのデータが記憶されている論理ブロックを示す情報である)

20

(2)クラスタディレクトリ番号(トラック内のデータの少なくとも一部がFS12またはIS13に記憶されている場合に有効となる情報であり、トラック内のデータがFS12またはIS13に記憶されている場合に、トラック毎に存在するクラスタディレクトリテーブルのテーブル番号を示す情報である)

(3)FS/ISクラスタ数(このトラック内のクラスタが、いくつFS12またはIS13に記憶されているかを示す情報であり、デフラグするかどうかを決めるために使用する)。

【0053】

・クラスタディレクトリテーブル31(正引き)

30

トラック内のデータがFS12またはIS13に記憶されている場合に、その論理ブロックまでたどるための中間的なテーブルであり、トラック別に備えられている。各クラスタディレクトリテーブル31に登録されるクラスタディレクトリ情報は、クラスタテーブル32のテーブル番号を示す情報(クラスタテーブル番号情報)の配列からなる。LBAのクラスタアドレスのLSB(t)ビット中の上位数ビットで、1つのクラスタディレクトリテーブル31中に配列されている複数のクラスタテーブル番号情報からひとつの情報を選択する。

【0054】

このクラスタディレクトリテーブル31としては、書き込み時刻を基準として、クラスタディレクトリ情報(クラスタテーブル番号情報の配列)の新旧の順序を、対応するトラックアドレスとともに、LRU(Least Recently used)で双方向リストで管理するためのクラスタディレクトリLRUテーブル31aと、空いているクラスタディレクトリを、対応するトラックアドレスとともに、双方向リストで管理するクラスタディレクトリ空き情報テーブル31bとがある。

40

【0055】

・クラスタテーブル32(正引き)

クラスタディレクトリテーブル31と関連し、トラック内のデータがFS12またはIS13に記憶されている場合に、どの論理ブロックのどのクラスタ位置にデータが記憶されているかを管理するテーブルである。トラックテーブル30からクラスタディレクトリテーブル31を経由して間接参照される。実体は、複数クラスタ分の論理ブロックアドレ

50

ス+クラスタ位置の配列である。LBAのクラスタアドレスのLSB(t)ビット中の下位数ビットで、1つのクラスタテーブル32中に配列されている複数の(論理ブロックアドレス+クラスタ位置)からひとつの情報を選択する。後述のクラスタブロック情報の番号とその中のクラスタ位置の情報も配列としてもつ。

#### 【0056】

・クラスタブロック情報テーブル33(逆引き)

クラスタブロックとは、論理ブロックのうちクラスタ単位でデータを記憶するものをいう。クラスタブロック情報は、FS12、IS13の論理ブロックを管理するための情報であり、論理ブロック内にどのようなクラスタが入っているかを示す情報である。クラスタブロック情報同士を双方向リストとしてFS12、IS13内のFIFOの順序で連結される。

10

クラスタブロック情報は、

(1)論理ブロックアドレス

(2)有効クラスタ数

(3)当該論理ブロックに含まれるクラスタのLBA

を有する。

クラスタブロック情報テーブル33は、使われていないクラスタブロック情報を管理する空き情報管理用のクラスタブロック情報テーブル33a、FS12に含まれるクラスタブロック情報を管理するFS用のクラスタブロック情報テーブル33b、IS13に含まれるクラスタブロック情報を管理するIS用のクラスタブロック情報テーブル33cを有し、各テーブル33a~33cは、双方向リストとして管理されている。逆引きアドレス変換の主な用途はIS13のコンパクションであり、コンパクション対象の論理ブロックにどのようなクラスタが記憶されているかを調べ、データを他の場所へ書き直すために使用する。よって、本実施の形態では、逆引きアドレス変換はクラスタ単位でデータを記憶しているFS12、IS13のみを対象としている。

20

#### 【0057】

・論物変換テーブル40(正引き)

論物変換テーブル40は、論理ブロックアドレスと物理ブロックアドレスとの変換、寿命に関する情報を管理するためのテーブルである。論理ブロックアドレス毎に、当該論理ブロックに所属する複数の物理ブロックアドレスを示す情報、当該論理ブロックアドレスの消去回数を示す消去回数情報、クラスタブロック情報の番号などの情報を有している。あるLBAのデータを他の場所へ書き直すには、元のクラスタブロック内のLBAを無効にする必要があり、LBAからクラスタブロックをたどる必要がある。そのために、論物変換テーブル40で管理する論理ブロックの管理情報に、クラスタブロック情報の識別子を記憶している。

30

#### 【0058】

(スナップショット、ログ)

上記各管理テーブルで管理される管理情報によって、ホスト1で使用されるLBAと、SSD100で使用される論理NANDアドレス(論理ブロックアドレス+オフセット)と、NANDメモリ10で使用される物理NANDアドレス(物理ブロックアドレス+オフセット)との間を対応付けることができ、ホスト1とNANDメモリ10との間のデータのやり取りを行うことが可能となる。

40

#### 【0059】

上記各管理テーブルのうちNAND管理用のテーブル(図7のトラックテーブル30、クラスタディレクトリテーブル31、クラスタテーブル32、クラスタブロック情報テーブル33、論物変換テーブル40など)は、不揮発性のNANDメモリ10の所定の領域に記憶されており、起動時に、NANDメモリ10に記憶されていた各管理テーブルを揮発性のDRAM20の作業領域に展開して、この展開された管理テーブルをデータ管理部120が使用することで、各管理テーブルは更新されていく。DRAM20上に展開された各管理テーブルをマスターテーブルと呼ぶ。このマスターテーブルは、電源が切れても

50

、電源が切れる以前の状態に復元する必要があるため、このためマスターテーブルを不揮発性のNANDメモリ10に保存する仕組みが必要となる。スナップショットは、NANDメモリ10上の不揮発性の管理テーブルの全体を指し、DRAM20に展開されたマスターテーブルをそのままNANDメモリ10に保存することを、スナップショットをとるとも表現する。ログは、管理テーブルの変更差分のことである。マスターテーブルの更新の度に、スナップショットをとっていたのでは、速度も遅く、NANDメモリ10への書き込み数が増えるために、通常は変更差分としてのログだけをNANDメモリ10に記録していく。ログをマスターテーブルに反映し、NANDメモリ10に保存することを、コミットすると表現する。

#### 【0060】

図8に、データ更新時に、スナップショットとログがどのように更新されるかを示す。データ管理部120がデータ更新する際に、マスターテーブルに加えた変更内容をDRAM20上のログ(DRAMログと呼ぶ)に蓄積する。管理テーブルの種類によっては、マスターテーブルを直接更新し、更新内容をDRAMログに蓄積したり、マスターテーブルには直接変更を加えず、変更領域をDRAMログ上に確保して、その領域に更新内容を記録したりする。データの読み書き処理の際には、マスターテーブルの他に蓄積されたDRAMログも参照する。

#### 【0061】

データの更新が安定したら、ログのコミットを行う。コミット処理では、DRAMログの内容を必要に応じてマスターテーブルに反映させ、さらにDRAMログの内容をNANDメモリ10に保存して不揮発化する。スナップショットをNANDメモリ10に保存するのは、正常な電源断シーケンスの際、ログの保存領域が不足した場合などとする。ログまたはスナップショットがNANDメモリ10に書き終わった時点で、管理テーブルの不揮発化が完了する。

#### 【0062】

##### ・Read処理

つぎに、読み出し処理の概要について説明する。ATAコマンド処理部121から、Readコマンドおよび読み出しアドレスとしてのLBAが入力されると、データ管理部120は、RC管理テーブル23とWCトラックテーブル24を検索することで、WC21またはRC22にLBAに対応するデータが存在しているか否かを探索し、キャッシュヒットの場合は、該当LBAに対応するWC21またはRC22のデータを読み出して、ATAコマンド処理部121に送る。

#### 【0063】

データ管理部120は、RC22またはWC21でヒットしなかった場合は、検索対象のデータがNANDメモリ10のどこに格納されているかを検索する。データがMS11に記憶されている場合は、データ管理部120は、LBAトラックテーブル30論物変換テーブル40と辿ることで、MS11上のデータを取得する。一方、データがFS12, IS13に記憶されている場合は、データ管理部120は、LBAトラックテーブル30クラスタディレクトリテーブル31クラスタテーブル32論物変換テーブル40と辿ることで、FS12, IS13上のデータを取得する。

#### 【0064】

##### ・Write処理

##### (WC21での処理)

つぎに、書き込み処理の概要について説明する。書き込み処理では、ATAコマンド処理部121からWriteコマンドおよび書き込みアドレスとしてのLBAが入力されると、LBAで指定されたデータをWC21に書き込む。WC21に空き領域がない場合は、DRAM管理用の各種管理テーブルを参照してWC21からデータを追い出して、NANDメモリ10に書き込み、空き領域を作成する。トラック内の有効クラスタ数が所定パーセント未満のトラックは低密度トラックとし、クラスタサイズデータとしてFS12を追い出し先とする。FS12が追い出し先の場合は、トラック内の有効クラスタを論理ページ

10

20

30

40

50

単位で書き込む。

【 0 0 6 5 】

トラック内の有効クラスタ数が所定パーセント以上のトラックは高密度トラックとし、トラックサイズのデータとしてMS 1 1を追い出し先とする。MS 1 1が追い出し先の場合は、トラックサイズのデータのまま論理ブロック全体に書き込む。書き込み対象の論理ブロック数が複数の場合は、倍速モードやバンクインターリーブを利用して転送効率を上げる。WC 2 1に書き込まれたデータに応じて、またNANDメモリ10へのデータ追い出しに応じて、DRAM管理用の各種管理テーブルを更新する。

【 0 0 6 6 】

( MS 1 1 への書き込み )

MS 1 1 への書き込みは、図 9 に示すように、次の手順で実行される。

1 . DRAM 2 0 上にトラックのデータイメージを作成 ( 穴埋め処理 ) する。すなわち、WC 2 1 に存在しないクラスタ、WC 2 1 に全セクタを保持していないクラスタに関しては、NANDメモリ10から読み出して、WC 2 1 のデータと統合する。

2 . MS 1 1 用に、論理ブロック ( トラックブロック ) をCFBから確保する。トラックブロックとは、論理ブロックのうちトラック単位でデータを記憶するものをいう。

3 . 作成したトラックのデータイメージを確保した論理ブロックに書き込む。

4 . トラックのLBAからトラック情報を調べ、トラック情報と書き込んだ論理ブロックに対応する論理ブロックアドレスとを関連付け、NAND管理用の所要のテーブルに登録する。

5 . WC 2 1 , NANDメモリ10の古いデータを無効化する。

【 0 0 6 7 】

( FS 1 2 への書き込み )

FS 1 2 への書き込みは、DRAM 2 0 上にクラスタのデータイメージを作成 ( 穴埋め処理 ) し、新たに確保する論理ブロック ( クラスタブロック ) に対し論理ページ単位の書き込みを、擬似SLCモードを使用して行う。確保する論理ブロックは、書き込むデータイメージ以上の書き込み可能な論理ページをもつフラグメントフリーブロック ( FFB ) を優先し、ない場合はコンプリートフリーブロック ( CFB ) を使用する。FS 1 2 への書き込みは、図 1 0 に示すように、以下の手順で実行する。

【 0 0 6 8 】

WC 2 1 からFS 1 2 に低密度トラックのデータを書き込むための論理ブロック ( クラスタブロック ) のことをFS Input Buffer ( 以下、FSIB ) と呼ぶ。

1 . WC 2 1 から入力された低密度トラック内の総データ量が小さい場合、すなわち有効クラスタ数が所定の閾値よりも少ない場合には、それを書き込めるFFBを確保し、FSIBとする。

2 . WC 2 1 から渡された低密度トラック内の総データ量が大きい場合、すなわち有効クラスタ数が所定の閾値以上の場合には、CFBを確保し、FSIBとする。このとき、並列で書き込むことが出来る複数の論理ブロックを確保し、FSIBとする。

3 . DRAM 2 0 上で、書き込むクラスタのデータイメージを作成する。すなわち、WC 2 1 に全セクタを保持していないクラスタに関しては、WC 2 1 上に存在しないセクタのデータをNANDメモリ10から読み出し、WC 2 1 上のセクタのデータと統合する。

4 . WC 2 1 上のクラスタと、作業領域上に作ったクラスタイメージをFSIBに書き込む。

5 . FSIBをFS 1 2 のリストに追加する。

6 . 書き込んだトラックを、クラスタディレクトリLRUテーブル3 1 aの末尾に挿入しなおす。

【 0 0 6 9 】

( FS 1 2 からIS 1 3 への移動 )

FS 1 2 管理下の論理ブロック数が所定の最大論理ブロック数を越えている場合は、図 1 1 に示すように、FS 1 2 から溢れた論理ブロックをそのままIS 1 3 に移動する。一

10

20

30

40

50

度の処理単位で移動する論理ブロック数は、溢れた論理ブロック内の有効クラスタ数に応じて、以下のルールで決定する。

- ・溢れた論理ブロック内のクラスタ数がMLCモードの1論理ブロック分の境界に近くなるように、FS12の最も古い論理ブロックから移動する論理ブロックを追加する。MLCモードの1論理ブロック分の境界に近くするのは、コンパクション後の論理ブロックに、なるべく多くの有効クラスタを収容することを目的とする。

- ・クラスタ数がIS13で同時にコンパクションできるクラスタ数を超える場合は、IS13で同時にコンパクションできるクラスタ数以下になるようなブロック数とする。

- ・移動ブロック数には、上限値を設ける。

#### 【0070】

(IS13でのコンパクションとデフラグ)

IS13では、IS管理下の論理ブロック数が最大論理ブロック数を越えた場合に、MS11へのデータ移動(デフラグ処理)と、コンパクション処理によって、管理下の論理ブロック数を最大数以下に抑える。データの消去単位(論理ブロック)と、データの管理単位(クラスタ)が異なる場合、NANDメモリ10の書き換えが進むと、無効なデータによって、論理ブロックは穴あき状態になる。このような穴あき状態の論理ブロックが増えると、実質的に使用可能な論理ブロックが少なくなり、NANDメモリ10の記憶領域を有効利用できないので、有効クラスタを集めて、違う論理ブロックに書き直すことをコンパクションという。デフラグ処理とは、FS12, IS13のクラスタをトラックに統合して、MS11に追い出す処理をいう。

#### 【0071】

つぎに、本実施の形態の要部について説明する。上述したSSDにおいては、NANDメモリ10の書き込みに時間がかかる、書き込み回数に制限がある、書き込みの大きさの単位が固定であるなどの理由で、ランダムアクセス可能な高速なメモリとしてのDRAM20上にWC21を設け、WC21に一時的にデータを蓄えてからNANDメモリ10にデータを書き込むようにしている。また、前述したように、NANDメモリ10への書き込み回数(消去回数)を減らすために、WC21のデータに関し、大きなデータ(高密度トラック)はMS11へ、小さなデータ(低密度トラック)はFS11へ書き込むような記憶部の切替え制御も行われている。

#### 【0072】

また、NANDメモリ10にデータが書き込まれて、各記憶部のリソース(容量や管理テーブルのエントリ数)が不足した場合には、コンパクションやデフラグなどのNANDの整理を行うことで、記憶部のリソースを確保している。同様に、WC21に関しても、そのリソース(領域や管理テーブルのエントリ数)が限界を超える場合には、前述したように、NANDメモリ10にデータを追い出して、WC21のリソースを確保する。このときの追い出し条件としてよく用いられるのは、次のホスト1からの書き込みに耐えられるだけの空きリソースをWC21で確保しておくという条件である。この条件では、ホスト1からWC21への書き込みを常に受け付けられるようにすることによって、単一の書き込みコマンドへのレスポンスは向上するが、NANDメモリ10の整理が進んでおらず、NANDメモリ10への書き込みに時間がかかると、後続の書き込みコマンドへのレスポンスは低下する。

#### 【0073】

そこで、本実施の形態では、WC21からNANDメモリ10への追い出しを早めに行うための閾値(AF閾値:Auto Flush閾値)を設ける。NANDメモリ10での整理(コンパクションおよびデフラグなど)が充分に進んでいる場合には、早めにWC21からNANDメモリ10にデータを追い出すことによって、NANDメモリ10での整理を早く始め、それによって、後続のWC21からの追い出しも高速に行えるようになる。その結果、WC21のリソース(メモリ領域、管理テーブルのエントリ)を多く確保できるため、後続の書き込みコマンドのレスポンスが向上する。

#### 【0074】

10

20

30

40

50

以下、WC 21での追い出し処理について詳細に説明する。まず、図12を用いてWC 21の管理構造をより詳細に説明する。この実施の形態では、図1に示したように、NANDメモリ10の各並列動作要素10a~10dは、バンクインターリーブ可能な4つのバンク(Bank0~3)を有し、各メモリチップは、並列動作可能なプレーン0およびプレーン1の2つのプレーンを有するものとする。そして、この実施の形態では、図12に示すような、LBAとNANDメモリ10のバンク/プレーンの割り当てを行うようにしている。すなわち、LBAのトラックアドレスのLSB1ビット目にプレーン(P)を割り当て、LBAのトラックアドレスのLSB2ビット、3ビット目にバンク(B)を割り当てている。なお、このようなLBAに対するバンク割り当ては、NANDメモリ10におけるMS11に対する書き込みの際にのみ使用される。

10

## 【0075】

WCトラックテーブル24は、前述したように、WC21上に記憶されたデータに関するWCトラック情報をLBAからトラックアップするための例えばハッシュテーブルであり、LBAのトラックアドレスのプレーン/バンク割り当てビット(P, B)を含むLSB数ビット(g個)をインデックスとし、g個のインデックス毎にn個(way)のエントリ(タグ)を有する。各タグには、LBAトラックアドレスと該トラックアドレスに対応するWCトラック情報へのポインタが記憶されている。したがって、WC21には、(g×n)個の異なるトラックをキャッシュすることができる。WCトラックテーブル24の更新に応じて、WCトラックテーブル24のインデックス毎の空きエントリ数(または使用エントリ数)1が計数されており、これらの空きエントリ数1は、WC21の追い出しをトリガするための1つのパラメータ(WCリソース使用量)となる。

20

## 【0076】

WCトラック情報テーブル25は、図13にも示すように、WCトラック情報をLRUで例えば双方向リストで管理するためのWCトラックLRU情報テーブル25aと、空いているWCトラック情報の番号を例えば双方向リストとして管理するWCトラック空き情報テーブル25bとを有する。

## 【0077】

WCトラック情報には、前述したように、

- (1)WC21内に存在するトラック内の有効クラスタ数を示す情報、フルのクラスタ数(セクタデータが満杯のクラスタの個数)を示す情報、
- (2)LBAのLSB側の数ビット部分であるクラスタ内オフセットに基づいて作成される情報であり、1クラスタに含まれる複数のセクタのうちのどのセクタに有効なデータを保持しているかを示す情報(セクタビットマップ)、
- (3)トラックの状態情報(有効、無効、ATAからのデータ転送中、NANDに書き込み中など)、
- (4)トラック単位にオール0のデータが含まれるか否かを識別する情報
- (5)クラスタ位置情報:(図13に示すように、LBAのトラック内クラスタインデックス(tビット)に対応するクラスタ領域番号をインデックスとした(( $2^t - 1$ )個)のクラスタ領域が確保され、各クラスタ領域には、当該クラスタデータがWC21のどこに存在するかを示すクラスタ位置情報が格納される。空きのクラスタに対応するクラスタ領域番号のクラスタ領域には、無効値が格納される。)

30

などが含まれている。

## 【0078】

WC21においては、前述したように、最大(g×n)個の異なるトラックをキャッシュすることが可能であるが、WCトラックLRU情報テーブル25aによって、WC21で使用されているトラックに関する情報を管理している。一方、WCトラック空き情報テーブル25bは、WC21にキャッシュすることが可能な最大トラック数(g×n)に対する空きのWCトラック情報を管理している。WCトラックLRU情報テーブル25aに登録されているWCトラック情報の個数をd個とした場合、WCトラック空き情報テーブル25bでは、(g×n)-d個の空きのWCトラック情報を管理している。WCトラッ

40

50

ク空き情報テーブル 25 b に用意されている WCトラック情報は、WCトラック LRU 情報テーブル 25 a で使用される WCトラック情報の領域を WC 21 に確保するためのもので、新しいトラックを管理し始めるときに、図 13 に示す各種情報を格納する。すなわち、WC 21 で新たな WCトラック情報を管理する必要が生じた場合、WCトラック空き情報テーブル 25 b から 1 つ WCトラック情報を確保し、確保した WCトラック情報に所要の情報を格納してから WCトラック LRU 情報テーブル 25 a のリンクに接続し直すようにする。なお、新たに確保された WCトラック情報に対応するトラックアドレスなどの情報は WCトラックテーブル 24 に登録され、高密度トラックの場合は高密度トラック情報テーブル 26 に登録される。

【 0079 】

WCトラック空き情報テーブル 25 b から新たな WCトラック情報を確保する度に ( 図 13 のフリーリストから WCトラック情報のリストが 1 個抜かれる度に )、図 13 に示す WC 21 での WCトラック情報の空き数を示す WCトラック情報の空き数 2 が - 1 され、WC 21 から NANDメモリ 10 に対する追い出しなどの発生によって、WCトラック LRU 情報テーブル 25 a に登録された WCトラック情報が解放されて WCトラック空き情報テーブル 25 b に戻される度に、WCトラック情報の空き数 2 が + 1 される。勿論、WCトラック情報空き数情報 2 の代わりに WC 21 での WCトラック情報の使用数を管理するようにしてもよい。WCトラック情報空き数情報 ( または使用数 ) 2 は、WC 21 の追い出しをトリガするための 1 つのパラメータ ( WCリソース使用量 ) となる。

【 0080 】

図 12 に示した WCクラスタ領域管理テーブル 29 は、各トラック中で空いているクラスタ領域を管理するためのもので、空いているクラスタ領域番号を FIFO 構造や双方向リンクリストなどによって管理している。また、WCクラスタ領域管理テーブル 29 によってクラスタ領域の総空き数が管理されている。クラスタ領域の使用数の最大値とは、WC 21 のキャッシュ容量に対応し、例えば、32 MB の WC 21 の場合は、クラスタ領域の使用数の最大値とは、32 MB に対応するクラスタ数になる。クラスタ領域の総空き数が 0 のときが、クラスタ領域の使用数の最大値に対応する。クラスタ領域の総空き数 ( または使用数 ) 3 は、WC 21 の追い出しをトリガするための 1 つのパラメータ ( WCリソース使用量 ) となる。

【 0081 】

図 14 は、WC高密度トラック情報テーブル 26 を示すものである。この高密度トラック情報テーブル 26 は、MS 11 に書き込むことになる、有効クラスタ数が多い高密度トラックに関するトラック情報を、LBA のトラックアドレスのバンク番号 ( B ) をインデックスとして管理するためのハッシュテーブルであり、各インデックス毎に m 個のエントリ ( way ) を有する。確実に並列に書き込めるバンク別に、WC 21 上の高密度トラックを管理することによって、WC 21 から MS 11 への追い出しの際のトラック検索に要する時間を高速化する。また、この WC 高密度トラック情報テーブル 26 によって、MS 11 に書き込むことになるトラック数を、各バンクで、最大 m 個となるように規制することで、WC 21 から NANDメモリ 10 への最大フラッシュ ( Flush ) 時間を抑えるようにする。WC 高密度トラック情報テーブル 26 における同じインデックスには、倍速モードで書けるトラック ( プレーン 0、1 ) が有る可能性もあるし、ない可能性もある。この WC 高密度トラック情報テーブル 26 におけるエントリ数を計数することによって、バンク番号毎に、高密度トラック情報の個数を管理している。このバンク毎の高密度トラック情報の個数 4 は、WC 21 の追い出しをトリガするための 1 つのパラメータ ( WCリソース使用量 ) となる。

【 0082 】

WC低密度トラック情報テーブル 27 ( 図 7 参照 )

FS 12 に書き込むことになる低密度のトラック情報を管理するためのもので、低密度トラックのクラスタ数の合計値を管理している。この低密度トラックのクラスタ数の合計値 5 は、WC 21 の追い出しをトリガするための 1 つのパラメータ ( WCリソース使用

10

20

30

40

50

量)となる。

【0083】

図15は、WC21の追い出しをトリガするための複数のパラメータ(WCリソース使用量)と、2つの閾値(AF閾値Caf、上限値Clmt)との関係の一例を示す図である。図15に示すように、WC21の追い出しをトリガするための複数のパラメータとしては、前述したように、

- ・WCクラスタ領域(クラスタ領域の総空き数) 3
- ・バンク毎の高密度トラック情報の個数(MS行きトラック数) 4
- ・低密度トラックのクラスタ数の合計値(FS行きクラスタ数あるいはFS行きクラスタデータ量) 5
- ・WCトラック情報の個数(WCトラック情報空き数) 2
- ・WCトラックテーブルのインデックス毎の使用エントリ数(または空きエントリ数) 1
- ・フルのトラック数 6

がある。これらWC21の追い出しをトリガするための複数のパラメータ1~6をWCリソース使用量とも呼ぶ。

【0084】

フルのトラック数6とは、セクタおよびクラスタが満杯のトラックの個数である。フルのトラックとは、ホスト1からデータが書き込まれ、WCトラック情報内のフルのクラスタ数が、クラスタ領域番号の個数である2<sup>1</sup>個になったトラックであり、フルのトラック数を計数する専用のカウンタ(図示せず)を有している。

【0085】

これら6個のパラメータ1~6には、WC21の追い出し処理のために、2つの閾値(AF閾値Caf、上限値Clmt)が設定されている。図15において、各パラメータ1~6に設定されている最大値maxは各パラメータ1~6が取り得る実質的な最大値を示すもので、基本的にWC21の追い出しをトリガする閾値としての意味はない。

【0086】

上限値Clmtは、パラメータ1~6がこれ以上になると次のwrite要求を受け付けられない可能性のある閾値であり、上限値Clmtを越えた場合は、ホスト1からのwrite要求を待たせる可能性がある。よって、上限値Clmtは、ホスト1からの次のwrite要求を待たせるための閾値であるともいうことができる。いずれかのパラメータ1~6が上限値Clmtを越えて、次のwrite要求を待たせている間、NANDメモリ10の整理が終わっている場合は、WC21からNANDメモリ10へデータを追い出して、該当するパラメータが上限値Clmt以下になるようにする。全てのパラメータ1~6が上限値Clmt以下になると、ホスト1からの次のwrite要求を受け付ける。いずれかのパラメータ1~6が上限値Clmtを越えて、次のwrite要求を待たせている間、NANDメモリの整理が終わっていない場合は、NANDメモリ10の整理を優先して、NANDメモリ10へのデータ追い出しを実行しない。NANDメモリ10の整理が終了したら、WC21からNANDメモリ10へデータを追い出して、該当するパラメータを上限値Clmt以下とし、その後ホスト1からの次のwrite要求を受け付ける。

【0087】

AF閾値Cafは、オートフラッシュ(Auto Flush)処理を行わせるための閾値である。Auto Flush処理とは、ホスト1からのFlushコマンドに関係なく行う処理であり、状況に応じてデータ管理部120の判断によってWC21のデータの一部または全てをNANDメモリ10に追い出す処理である。Auto Flush処理は、ホスト1からのWriteコマンド終了後に実行される処理であり、先行してWC21に一定の空き領域を作っておくことで、書き込み性能をトータル的に向上させるための処理であって、いずれかのパラメータ1~6がAF閾値Caf以上になると、Auto Flush処理を実行して早めにWC21からデータをNANDメモリ10に追い出す。したがって、通常は、AF閾値Caf<上限値Clmt<最大値maxの関係にある。いずれかのパラメータ1~6がAF閾値Cafを越

えていた場合、NANDメモリ10の整理の状況を確認し、NANDメモリ10の整理が終了していた場合に、Auto Flush処理を実行する。いずれかのパラメータ 1 ~ 6がAF閾値Cafを越えていた場合でも、NANDメモリ10の整理が終了していない場合は、WC21はまだホスト1からのwrite要求を待たせるほど余裕のない状況ではないので、NANDメモリ1の整理を優先させる。

#### 【0088】

つぎに、各パラメータ 1 ~ 6毎に、最大値max、上限値Climt、AF閾値Cafについて説明する。WCクラスタ領域(クラスタ領域の総空き数) 3の最大値maxであるZとは、WC21の容量であり、WC21が32MBの容量を有する場合は、Z = 32MBである。WCクラスタ領域(クラスタ領域の総空き数) 3のAF閾値Cafとしては、例えば最大値maxの半分であるZ/2に設定する。上限値Climtとしては、ホスト1からの1回のデータ転送サイズを考慮して決定する。例えば、(7/8 ~ 15/16)Z程度の値に設定する。

10

#### 【0089】

WCトラック情報の個数 2の最大値maxは、WCトラックテーブル24の総エントリ数であり、この場合(g × n)である。AF閾値Cafとしては、例えば最大値maxの半分である(g × n) / 2に設定する。WCトラック情報の個数 2に関する上限値Climtは、図15では設定されていないが、(g × n)より小さく、(g × n) / 2より大きな適当な値を設定してもよい。

20

#### 【0090】

WCトラックテーブルのインデックス毎の使用エントリ数 1の最大値maxは、nである。1に関するAF閾値Cafは、図15では設定されていないが、例えば最大値maxの半分であるn / 2程度の値に設定するようにしてもよい。上限値Climtとしては、残り1(WCトラックテーブル24の最後のエントリ(way)しか残っていない状態)とする。

#### 【0091】

フルのトラック数 6に関しては、AF閾値Caf(= y)のみが設定されている。フルのトラック数 6が増えても、他のパラメータが空いていれば、次のホストからのwrite要求に応えることができるので、フルのトラック数 6については、上限値Climtを設定していない。

30

#### 【0092】

つぎに、MS行きトラック数(バンク毎) 4の最大値maxとは、図14に示す高密度トラック情報テーブル26のエントリ(way)数であるmである。このmという数値は、WC21の全データをNANDメモリ10に追い出す指令であるフラッシュ(flush)コマンドを処理するのに要する時間を考慮して決定されている。また、FS行きクラスタ数(FS行きクラスタデータ量) 5の最大値Qとは、FS12に追い出す低密度トラックのクラスタ数(あるいはクラスタデータ量)の合計値の最大値であり、このQという数値も、flushコマンドを処理するのに要する時間、NANDメモリ10側のリソース使用量(FS12への書き込みを抑えることで、FS12、IS13用のブロック数や、FS12、IS13用の管理テーブル量の増大を抑える)などを考慮して決定されている。flushコマンドを実行する際には、高密度トラックはMS11に追い出し、低密度トラックはFS12に追い出す必要がある。

40

#### 【0093】

前述したように、この実施の形態ではトラック = 論理ブロックとしている。論理ブロックとは、NANDメモリ10のチップ上の物理ブロックを複数組み合わせる仮想的ブロックのことであり、この実施の形態では、図1に示す4つの並列動作要素10a ~ 10d内の各物理ブロックに1回4ch並列動作させた単位を論理ブロックという。図14に示したWC高密度トラック情報テーブル26においては、1つのバンクに関する1つのエントリには1つのトラックアドレスが登録されており、1つのバンクの1つのトラックの追い出しを行うことは、1つの論理ブロックをNANDメモリ10に(高密度トラ

50

ックであるので、正確にはMS11に)1回書き込むことに相当する。

【0094】

一方、FS12は、この実施の形態では、擬似SLCモードで動作しており、FS12での論理ブロックサイズは、4値のMLCモードで動作しているMS11での論理ブロックサイズの半分である。したがって、FS12への1回の4ch並列書き込みによる書き込みサイズも、MS11への書き込みの場合の半分になる。その反面、擬似SLCモードでの書き込みは、MLCモードより数倍速い。

【0095】

flushコマンド処理に要する時間Tflは、高密度トラックをMS11に追い出す所要時間Taと、低密度トラックをFS12に追い出す所要時間Tbと、ログの書き込み処理など他の処理に要する時間Tc(固定値)の合計になる。則ち、 $Tfl = Ta + Tb + Tc$ となる。高密度トラックをMS11に追い出す所要時間Taは、4バンクでのバンクインタリーブを用いた場合を想定して、 $4 \times 1$ 回の書き込み所要時間(固定値)  $\times$  書き込み回数( $u1$ )となる。低密度トラックをFS12に追い出す所要時間Tbは、1回の書き込み所要時間(固定値)  $\times$  書き込み回数( $u2$ )となる。flushコマンド処理に要する時間Tflを所定の最悪時間Tflmax(固定値)時間以内に抑えようとした場合、例えば $u1 = u2$ として、下記の式

$$Tflmax(\text{固定値}) = Ta(=4 \times 1 \text{ 回の書き込み所要時間(固定値)} \times \text{書き込み回数}(u1)) + Tb(=1 \text{ 回の書き込み所要時間(固定値)} \times \text{書き込み回数}(u2)) + Tc(\text{固定値})$$

から、書き込み回数 $u1$ 、 $u2$ を求めることができる。

【0096】

このようにして求めた $u1$ が、MS行きトラック数(バンク毎)4の最大値maxとしてのmである。WC高密度トラック情報テーブル26のway数mはこのようにして決定する。また、上記によりflushコマンド処理に対する最悪時間Tflmaxを満足するFS12への書き込み回数 $u2$ も決定することができたので、(FS12への1回の論理ブロックの書き込みサイズ)  $\times$  (書き込み回数 $u2$ )を求めることで、FS行きクラスタデータ量5の最大値maxとしてのQ(MB)を求めることができる。

【0097】

MS行きトラック数(バンク毎)4に関するAF閾値Cafは、例えば最大値mの半分 $m/2$ に設定する。 $m/2$ に設定したことで、Auto Flush処理では、4トラックをMS11に書き込んでいる間に、残りの4トラックをホスト1からWC21に書き込むようなダブルバッファ的な並列処理を実行することができる。4に関する上限値Clmtとしては、例えば $(6/8 \sim 7/8)m$ 程度の値に設定する。上限値Clmtは、最大値maxからの残り分だけ空けておけば、次のwrite要求をNANDメモリ10に追い出さなくとも受け付けることが可能な量を考慮して設定する。

【0098】

FS行きクラスタ数(FS行きクラスタデータ量)5のAF閾値Cafは、例えば最大値Qの $1/4$ である $Q/4$ 程度に設定する。 $Q/4$ は、例えば、FS12に対して並列に書けば、ホスト1からWC21への書き込み速度と同等の速度が得られる値とする。5に対する上限値Clmtとしては、例えば $(6/8 \sim 7/8)Q$ 程度の値に設定する。上限値Clmtは、4と同様、最大値からの残り分だけ空けておけば、次のwrite要求をNANDメモリ10に追い出さなくとも受け付けることが可能な量を考慮して設定する。

【0099】

図16は、本実施の形態の要部の機能構成を示すブロック図である。前述したように、NANDメモリ10は、ユーザデータを記憶するMS11、FS12、IS13からなるユーザデータ記憶部を備えている。DRAM20には、WC21が備えられている。データ管理部であるコントローラ120は、ホスト1からのデータをWC21に対して書き込む制御を行うWC書き込み制御部210と、WC21からNANDメモリ10にデータを追い出す追い出し制御部211と、WC21から追い出され

10

20

30

40

50

たデータをNANDメモリ10に書き込む制御を実行するNAND書き込み制御部213と、NANDメモリ10での論理ブロック整理（コンパクション、デフラグなど）を実行するNAND整理部212を備えている。NAND整理部212は、NANDメモリ10での整理処理の状態を示すNAND整理状態信号（現在整理処理を実行中であるか終了したかを示す信号）を、WC書き込み制御部210およびWC追い出し制御部211に逐次送っている。WC書き込み制御部210およびWC追い出し制御部211は、NAND整理状態信号に基づきNANDでのブロック整理状態を判断する。

【0100】

（NANDメモリの整理）

NAND整理部212が行うNANDメモリ10の整理について説明する。NANDの整理とは、

- ・FS12, IS13管理下の論理ブロック数を所定の閾値以下にするためのコンパクション/デフラグ処理

- ・NAND管理用テーブル（クラスタディレクトリテーブル31、クラスタテーブル32など）のエントリ数を所定の閾値以下にするためのデフラグ処理などを含む。

これらFS12, IS13管理下の論理ブロック数、NAND管理用テーブルのエントリ数などNANDの整理の際に考慮するパラメータを総称してNANDリソース使用量と呼ぶ。NAND整理部212は、各NANDリソース使用量が閾値を越えているときには、NANDメモリ10の整理を実行し、現在処理を実行中である旨を示すNAND整理状態信号をWC書き込み制御部210およびWC追い出し制御部211に送る。

【0101】

（コンパクション）

コンパクション処理とは、IS13で行われる処理であり、IS管理下の論理ブロックが所定の閾値を越えた際に、有効な最新のクラスタデータを集めて、違う論理ブロックに書き直し、無効クラスタを解放することをいう。なお、この実施の形態では、トラックサイズと論理ブロックサイズを同じにしているために、MS11ではコンパクションは発生しないが、トラックサイズと論理ブロックサイズが異なる場合は、MS11でもコンパクションが発生する。このような場合は、MS11でのコンパクションを含めてNANDメモリ10の整理を行う必要がある。また、FS11においても、コンパクションを実行するようにしてもよく、その場合は、FS12でのコンパクションを含めてNANDメモリ10の整理を行う必要がある。

【0102】

（デフラグ）

デフラグ処理とは、FS12, IS13管理下の論理ブロック数を所定の閾値を越えた際、あるいはクラスタディレクトリテーブル31、クラスタテーブル32などのFS12, IS13管理用のテーブルのエントリ数が所定の閾値を越えた際に、FS12あるいはIS13のクラスタデータをトラックに統合して、MS11に追い出す処理をいう。

【0103】

つぎに、NANDメモリ10への追い出しを含むWC21での書き込み処理について説明する。WC21では、Writeコマンドで一時的に超過する場合を除いて、すべてのWC21のリソース使用量が常に上限値C<sub>lmt</sub>以下になるよう、WC21上のデータをNANDメモリ10に追い出す。また、ホスト1からの要求を待たせている場合を除いて、WCリソース使用量がAF閾値C<sub>af</sub>以下になるまで追い出しを継続する。具体的には、以下のような制御が実行される。

【0104】

（WCリソース使用量 > 上限値C<sub>lmt</sub>の場合）

まず、ホスト1からのWC21へのデータ書き込み終了後、WCリソース使用量が上限値C<sub>lmt</sub>を越えていた場合の処理について説明する。WC追い出し制御部211は、ホスト1からのデータ書き込み終了後、WCリソース使用量 1 ~ 6の状態を判断し、WC

10

20

30

40

50

リソース使用量 1 ~ 6 の何れかが上限値 C lmt を越えていた場合であって、ホスト 1 からつぎの Write 要求が来ていない場合は、つぎのような処理を実行する。WC 追い出し制御部 2 1 1 は、NAND 整理状態信号によって NAND メモリ 1 0 の整理状態を確認し、NAND メモリ 1 0 の整理が終了している場合は、上限値 C lmt を越えている要因となっている WC 2 1 のデータを MS 1 1 または FS 1 2 の何れかまたは両方に追い出して、全ての WC リソース使用量 1 ~ 6 を上限値 C lmt 以下とする。追い出し対象のトラックは、LRU で古いトラックの順から優先的に選択するとか、上限値 C lmt 以下に WC リソース使用量を抑えるための処理速度を優先して追い出しトラックを選択するとかして決定する。

【 0 1 0 5 】

WC 追い出し制御部 2 1 1 は、全ての WC リソース使用量 1 ~ 6 が上限値 C lmt 以下となった時点で、ホスト 1 からのつぎの Write 要求が来ていない場合は、WC リソース使用量が AF 閾値 Caf 以下になるまで追い出しを継続する。すなわち、WC 追い出し制御部 2 1 1 は、NAND 整理状態信号によって NAND メモリ 1 0 の整理状態を確認し、NAND メモリ 1 0 の整理が終了している場合は、AF 閾値 Caf を越えている要因となっている WC 2 1 のデータを MS 1 1 または FS 1 2 の何れかまたは両方に追い出して、全ての WC リソース使用量 1 ~ 6 を AF 閾値 Caf 以下とする。追い出し対象のトラックは、例えば、書き込み効率のよいフルのトラックを優先的に追い出す、あるいは LRU で古いトラックの順から優先的に選択するとかして決定する。一方、WC 追い出し制御部 2 1 1 は、全ての WC リソース使用量 1 ~ 6 が上限値 C lmt 以下となった時点で、ホスト 1 からのつぎの Write 要求が来ている場合、あるいは NAND 整理状態信号によって NAND メモリ 1 0 の整理が終了していないと判断した場合は、ホスト 1 からの要求あるいは NAND メモリ 1 0 の整理を優先して、Auto Flush 処理を実行しない。

【 0 1 0 6 】

また、WC 追い出し制御部 2 1 1 は、ホスト 1 からのデータ書き込み終了後、WC リソース使用量 1 ~ 6 の何れかが上限値 C lmt を越えていた場合であって、ホスト 1 からつぎの Write 要求が来ている場合は、つぎのような処理を実行する。WC 追い出し制御部 2 1 1 は、ホスト 1 からのつぎの Write 要求を待たせるとともに、NAND 整理状態信号によって NAND メモリ 1 0 の整理状態を確認し、NAND メモリ 1 0 の整理が終了している場合は、上限値 C lmt を越えている要因となっている WC 2 1 のデータを MS 1 1 または FS 1 2 の何れかまたは両方に追い出して、全ての WC リソース使用量 1 ~ 6 を上限値 C lmt 以下とする。そして、WC 追い出し制御部 2 1 1 は、全ての WC リソース使用量 1 ~ 6 が上限値 C lmt 以下となった時点で、ホスト 1 からの Write 要求を受け付ける。しかし、WC 追い出し制御部 2 1 1 は、ホスト 1 からのつぎの Write 要求を待たせた状態中に、NAND 整理状態信号によって NAND メモリ 1 0 の整理が終了していないと判断した場合は、NAND メモリ 1 0 の整理を優先して、NAND メモリ 1 0 の整理が終了するまで待機し、NAND メモリ 1 0 の整理の終了を確認した後に、上限値 C lmt を越えている要因となっている WC 2 1 のデータを MS 1 1 または FS 1 2 の何れかまたは両方に追い出して、全ての WC リソース使用量 1 ~ 6 を上限値 C lmt 以下とする。そして、全ての WC リソース使用量 1 ~ 6 が上限値 C lmt 以下となった時点で、ホスト 1 からの Write 要求を受け付ける。

【 0 1 0 7 】

( AF 閾値 Caf < WC リソース使用量 < 上限値 C lmt の場合 )

つぎに、ホスト 1 からの WC 2 1 へのデータ書き込み終了後、WC リソース使用量が上限値 C lmt を越えていないが、AF 閾値 Caf を越えていた場合の処理について説明する。WC 追い出し制御部 2 1 1 は、ホスト 1 からのデータ書き込み終了後、WC リソース使用量 1 ~ 6 の何れかが AF 閾値 Caf を越えていた場合に、ホスト 1 からのつぎの Write 要求が来ている場合は、そのまま Auto Flush 処理を実行せず、ホスト 1 からのつぎの Write 要求を受け付ける。

【 0 1 0 8 】

10

20

30

40

50

しかし、ホスト1からのつぎのWrite要求が来ていない場合は、つぎのような処理を実行する。WC追出し制御部211は、NAND整理状態信号によってNANDメモリ10の整理状態を確認し、NANDメモリ10の整理が終了している場合は、AF閾値を越えている要因となっているWC21のデータをMS11またはFS12の何れかまたは両方に追い出して、全てのWCリソース使用量1~6をAF閾値Caf以下とする。追い出し対象のトラックは、例えば、書き込み効率のよいフルのトラックを優先的に追い出す、あるいはLRUで古いトラックの順から優先的に選択するとかして決定する。

**【0109】**

一方、WC追出し制御部211は、WCリソース使用量1~6の何れかがAF閾値Cafを越えていた場合でも、NAND整理状態信号によってNANDメモリ10の整理が終了していないと判断した場合は、NANDメモリ10の整理を優先して、Auto Flush処理を実行しない。つぎのホスト1からのWriteコマンドが受信される前に、NAND整理状態信号によってNANDメモリ10の整理終了が確認された場合は、Auto Flush処理を実行する。このように、書き込み終了後にNANDメモリ10の整理が終わっている場合は、Auto Flush処理を実行することで、WC21を早めに余裕をもって空けることができる。

**【0110】**

(他の実施の形態)

上記では、まず、ホスト1からのWC21へのデータ書き込み終了後、WCリソース使用量1~6が上限値Climtを越えていた場合であって、かつホスト1からつぎのWrite要求が来ている場合であって、さらにNANDメモリ10の整理が終了していない場合は、ホスト1からのつぎのWrite要求を待たせてNANDメモリの整理を行い、NANDメモリ10の整理の終了後、追い出し処理を実行するようにしたが、つぎのような制御も可能である。すなわち、上記の条件が成立した場合、ホスト1のつぎのWrite要求を受け付けてもWCリソース使用量1~6が最大値maxを越えない場合は、ホスト1のつぎのWrite要求を受け付けて、WC21にデータを書き込ませ、リソース使用量1~6が最大値maxを越える場合は、ホスト1のつぎのWrite要求を待たせて、追い出し処理を行って全てのWCリソース使用量1~6が上限値Climt以下になるようにする。

**【0111】**

なお、上記実施の形態では、ホスト1からのWC21へのデータ書き込み終了後、すなわちデータの書き込み前に、WC21からNANDメモリ10への追い出しを行うようにしているが、WC21へのデータ書き込みと、NANDメモリ10への追い出しとを並列処理するようにしてもよい。

**【0112】**

このように本実施の形態によれば、WCリソース使用量を上限値Cafより値が小さいAF閾値Cafと比較し、WCリソース使用量がAF閾値Cafを越えているときは、NANDメモリ10での整理の状態を確認し、NANDメモリ10での整理が充分に進んでいる場合には、その時間を有効活用して早めにWC21からNANDメモリ10にデータを追い出すようにしているので、NANDメモリ10での整理を早く始めることができ、それによって、後続のWC21からの追い出しも高速に行えるようになる。その結果、WC21のリソース(メモリ領域、管理テーブルのエントリ)を多く確保できるため、後続の書き込みコマンドのレスポンスが向上し、ホストからの書き込み要求に対する応答性を全般的に向上させることができる。また、この実施の形態では、複数のWCリソース使用量1~6を用意し、これら複数のWCリソース使用量1~6毎に、上限値CafとAF閾値Cafを用意し、いずれかのWCリソース使用量1~6が対応する上限値CafまたはAF閾値Cafを越えた場合、上限値CafまたはAF閾値Cafを越えた原因となるデータを追い出すようにしているので、WCリソース使用量を上限値CafまたはAF閾値Caf以下にする際の追い出し処理を効率良くかく高速化することが可能となる。

**【0113】**

[第2の実施の形態]

図17は、SSD100を搭載したパーソナルコンピュータ1200の一例を示す斜視図である。パーソナルコンピュータ1200は、本体1201、及び表示ユニット1202を備えている。表示ユニット1202は、ディスプレイハウジング1203と、このディスプレイハウジング1203に収容された表示装置1204とを備えている。

【0114】

本体1201は、筐体1205と、キーボード1206と、ポインティングデバイスであるタッチパッド1207とを備えている。筐体1205内部には、メイン回路基板、ODD(optical disk device)ユニット、カードスロット、及びSSD100等が収容されている。

【0115】

カードスロットは、筐体1205の周壁に隣接して設けられている。周壁には、カードスロットに対向する開口部1208が設けられている。ユーザは、この開口部1208を通じて筐体1205の外部から追加デバイスをカードスロットに挿抜することが可能である。

【0116】

SSD100は、従来のHDDの置き換えとして、パーソナルコンピュータ1200内部に実装された状態として使用してもよいし、パーソナルコンピュータ1200が備えるカードスロットに挿入した状態で、追加デバイスとして使用してもよい。

【0117】

図18は、SSDを搭載したパーソナルコンピュータのシステム構成例を示している。パーソナルコンピュータ1200は、CPU1301、ノースブリッジ1302、主メモリ1303、ビデオコントローラ1304、オーディオコントローラ1305、サウスブリッジ1309、BIOS-ROM1310、SSD100、ODDユニット1311、エンベデッドコントローラ/キーボードコントローラIC(EC/KBC)1311、及びネットワークコントローラ1312等を備えている。

【0118】

CPU1301は、パーソナルコンピュータ1200の動作を制御するために設けられたプロセッサであり、SSD100から主メモリ1303にロードされるオペレーティングシステム(OS)を実行する。更に、ODDユニット1311が、装填された光ディスクに対して読み出し処理及び書き込み処理の少なくとも1つの処理の実行を可能にした場合に、CPU1301は、それらの処理の実行をする。

【0119】

また、CPU1301は、BIOS-ROM1310に格納されたシステムBIOS(Basic Input Output System)も実行する。尚、システムBIOSは、パーソナルコンピュータ1200内のハードウェア制御のためのプログラムである。

【0120】

ノースブリッジ1302は、CPU1301のローカルバスとサウスブリッジ1309との間を接続するブリッジデバイスである。ノースブリッジ1302には、主メモリ1303をアクセス制御するメモリコントローラも内蔵されている。

【0121】

また、ノースブリッジ1302は、AGP(Accelerated Graphics Port)バス1314等を介してビデオコントローラ1304との通信、及びオーディオコントローラ1305との通信を実行する機能も有している。

【0122】

主メモリ1303は、プログラムやデータを一時的に記憶し、CPU1301のワークエリアとして機能する。主メモリ1303は、例えばDRAMから構成される。

【0123】

ビデオコントローラ1304は、パーソナルコンピュータ1200のディスプレイモータとして使用される表示ユニット1202を制御するビデオ再生コントローラである。

【0124】

10

20

30

40

50

オーディオコントローラ 1305 は、パーソナルコンピュータ 1200 のスピーカ 1306 を制御するオーディオ再生コントローラである。

【0125】

サウスブリッジ 1309 は、LPC (Low Pin Count) バス上の各デバイス、及び PCI (Peripheral Component Interconnect) バス 1315 上の各デバイスを制御する。また、サウスブリッジ 1309 は、各種ソフトウェア及びデータを格納する記憶装置である SSD 100 を、ATA インタフェースを介して制御する。

【0126】

パーソナルコンピュータ 1200 は、セクタ単位で SSD 100 へのアクセスを行う。ATA インタフェースを介して、書き込みコマンド、読出しコマンド、フラッシュコマンド等が SSD 100 に入力される。

10

【0127】

また、サウスブリッジ 1309 は、BIOS-ROM 1310、及び ODD ユニット 1311 をアクセス制御するための機能も有している。

【0128】

EC/KBC 1311 は、電力管理のためのエンベデッドコントローラと、キーボード (KB) 1206 及びタッチパッド 1207 を制御するためのキーボードコントローラとが集積された 1 チップマイクロコンピュータである。

【0129】

この EC/KBC 1311 は、ユーザによるパワーボタンの操作に応じてパーソナルコンピュータ 1200 の電源を ON/OFF する機能を有している。ネットワークコントローラ 1312 は、例えばインターネット等の外部ネットワークとの通信を実行する通信装置である。

20

【図面の簡単な説明】

【0130】

【図1】SSDの構成例を示すブロック図。

【図2】NANDメモリチップに含まれる1個のブロックの構成例と、4値データ記憶方式でのしきい値分布を示す図。

【図3】ドライブ制御回路のハードウェア的な内部構成例を示すブロック図。

【図4】プロセッサの機能構成例を示すブロック図。

30

【図5】NANDメモリおよびDRAM内に形成された機能構成を示すブロック図。

【図6】LBA論理アドレスを示す図。

【図7】データ管理部内の管理テーブルの構成例を示す図。

【図8】スナップショットとログの生成形態を概念的に示す図。

【図9】MSへの書き込み手順を示す図。

【図10】FSへの書き込みを示す図。

【図11】FSからISへのブロック移動を示す図。

【図12】WCでの管理構造を示す図。

【図13】WCトラック情報テーブルを示す図。

【図14】WC高密度トラック情報テーブルを示す図。

40

【図15】WCリソース名(パラメータ名)とAF閾値、上限値との関係を示す図。

【図16】本実施の形態の要部構成を示す機能ブロック図。

【図17】SSDを搭載したパーソナルコンピュータの全体図。

【図18】SSDを搭載したパーソナルコンピュータのシステム構成例を示す図。

【符号の説明】

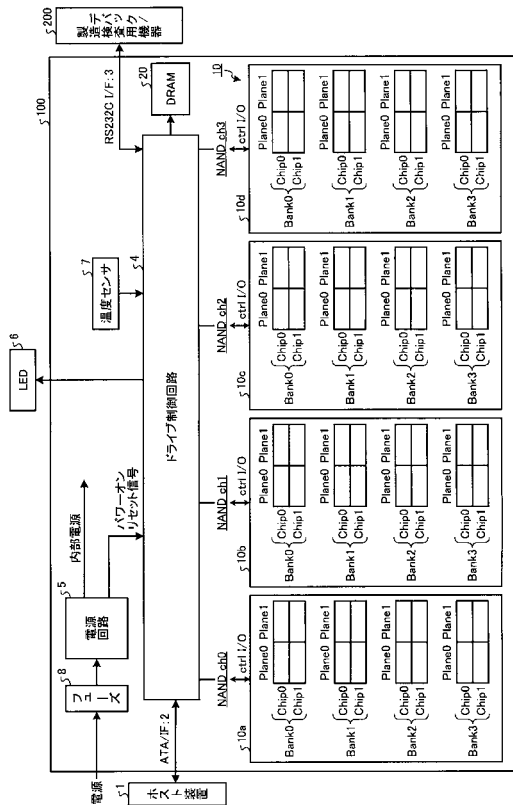
【0131】

1 ホスト装置、2 ATAインタフェース、4 ドライブ制御回路、5 電源回路、7 温度センサ、10 NANDメモリ、11 MS、12 FS、13 IS、20 DRAM、21 WC、22 RC、24 WCトラックテーブル、25 WCトラック情報テーブル、26 高密度トラック情報テーブル、27 低密度トラック情報テーブ

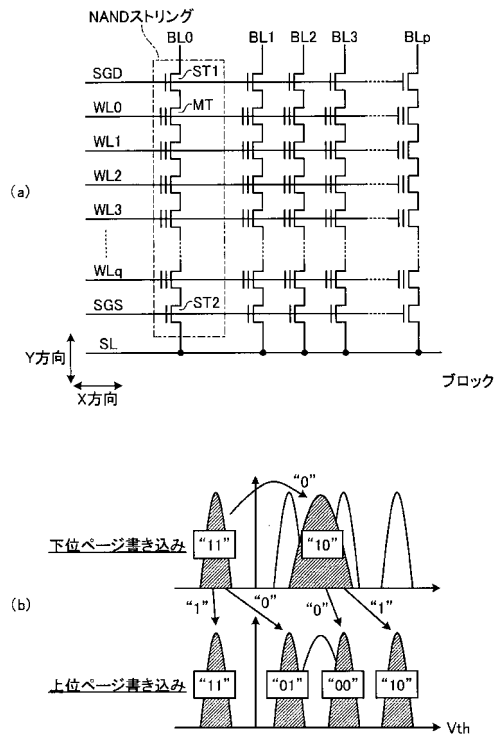
50

ル、30 トラックテーブル、31 クラスタディレクトリテーブル、32 クラスタテーブル、33 クラスタブロック情報テーブル、40 論物変換テーブル、120 データ管理部、210 WC書き込み制御部、211 WC追い出し制御部、212 NAND整理部、213 NAND書き込み制御部。

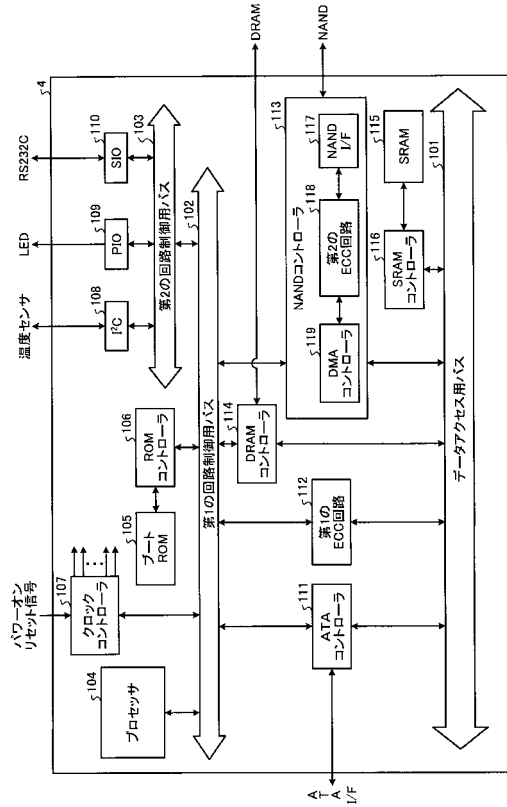
【図1】



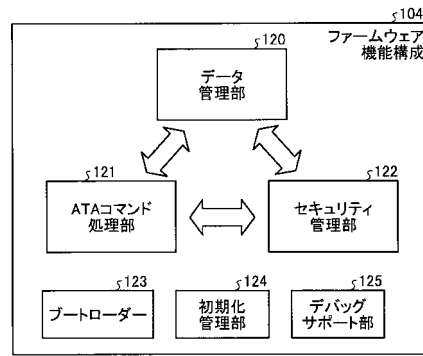
【図2】



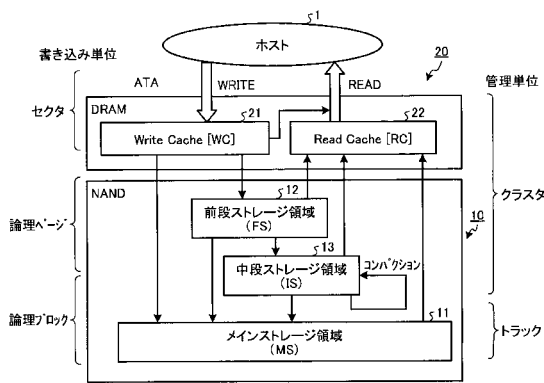
【図3】



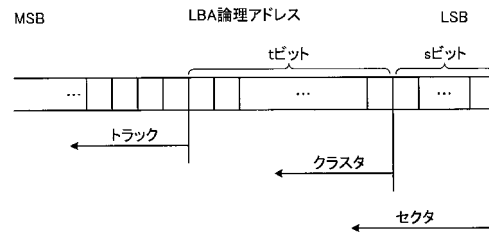
【図4】



【図5】

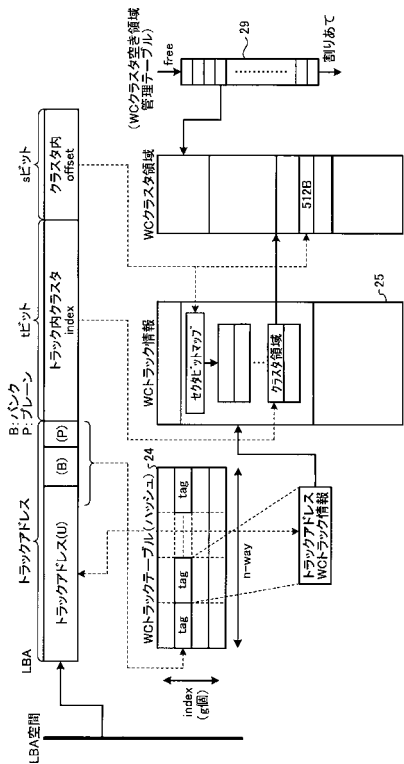


【図6】

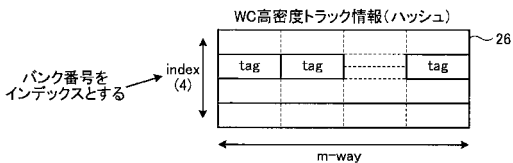




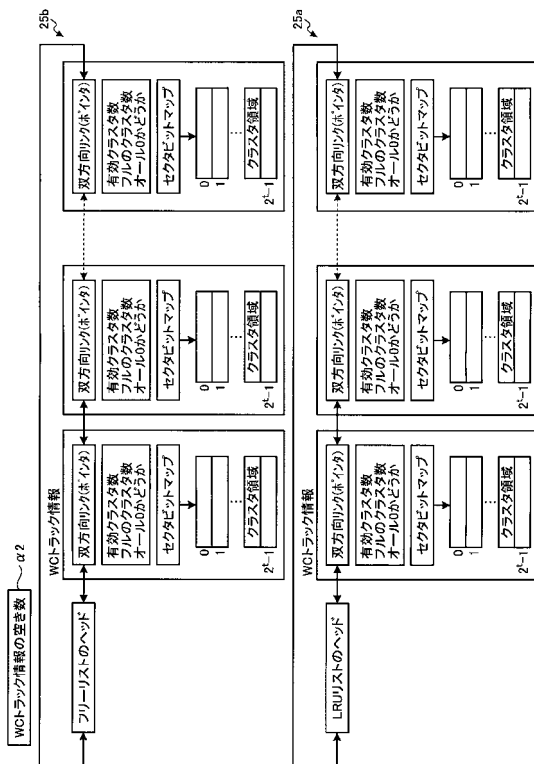
【 図 1 2 】



【 図 1 4 】



【 図 1 3 】



【 図 1 5 】

| パラメータ名                       | 最大値max  | AF閾値CAF   | 上限値Omt       |
|------------------------------|---------|-----------|--------------|
| WCクラスタ領域(MB)                 | Z       | Z/2       | (7/8~15/16)Z |
| MS行きトラック数(バンク毎)              | m       | m/2       | (6/8~7/8)m   |
| FS行きクラスタ数(MB)                | Q       | Q/4       | (6/8~7/8)Q   |
| WCトラッキング情報の数                 | (g × n) | (g × n)/2 |              |
| WCトラッキングアドレスのインデックス毎の使用エントリ数 | n       |           | 残1           |
| フルのトラック数                     |         | y         |              |

α3

α4

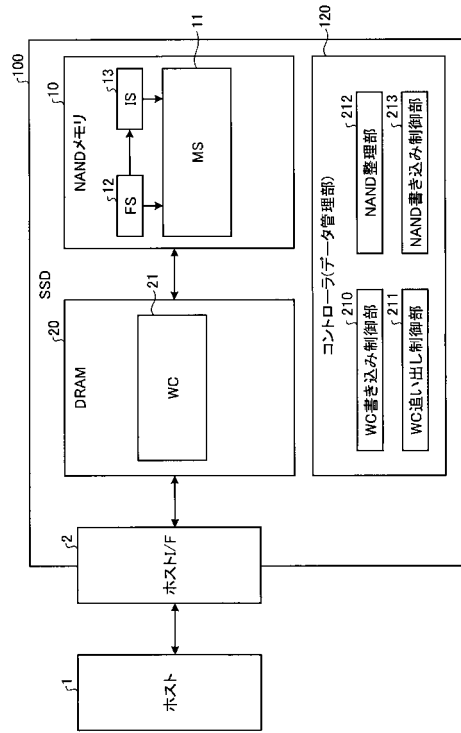
α5

α2

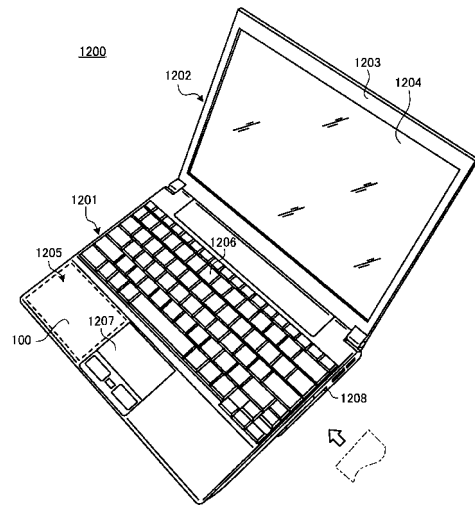
α1

α6

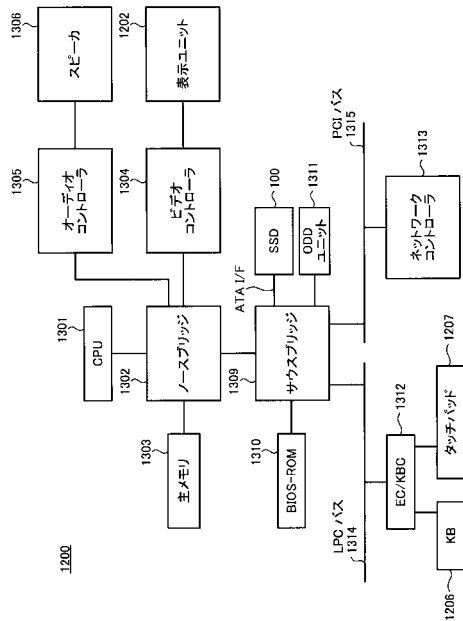
【図16】



【図17】



【図18】



---

フロントページの続き

審査官 渡部 博樹

(56)参考文献 特開平11-288387(JP,A)  
特開2000-089983(JP,A)

(58)調査した分野(Int.Cl., DB名)  
G06F12/00 - G06F12/06  
G06F12/08 - G06F12/16