

(19)日本国特許庁(JP)

## (12)特許公報(B2)

(11)特許番号  
特許第7091203号  
(P7091203)

(45)発行日 令和4年6月27日(2022.6.27)

(24)登録日 令和4年6月17日(2022.6.17)

(51)国際特許分類

F I

G 0 6 F	12/02	(2006.01)	G 0 6 F	12/02	5 3 0 E
G 0 6 F	12/00	(2006.01)	G 0 6 F	12/00	5 9 7 U
G 0 6 F	3/08	(2006.01)	G 0 6 F	12/00	5 7 1 A
G 0 6 F	3/06	(2006.01)	G 0 6 F	3/08	H
			G 0 6 F	3/06	3 0 4 Z

請求項の数 8 (全77頁) 最終頁に続く

(21)出願番号 特願2018-175148(P2018-175148)  
 (22)出願日 平成30年9月19日(2018.9.19)  
 (65)公開番号 特開2020-46963(P2020-46963A)  
 (43)公開日 令和2年3月26日(2020.3.26)  
 審査請求日 令和3年3月18日(2021.3.18)

(73)特許権者 318010018  
キオクシア株式会社  
東京都港区芝浦三丁目1番21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 東芝メ  
モリ株式会社内  
 審査官 酒井 恭信

最終頁に続く

(54)【発明の名称】 メモリシステムおよび制御方法

(57)【特許請求の範囲】

【請求項1】

ホストに接続可能なメモリシステムであって、  
 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を制御するコントローラと  
 を具備し、  
 前記コントローラは、  
前記不揮発性メモリを論理的に分割することによって得られる複数の領域の中のいずれか  
 の領域に割り当て済みのブロック毎に、各ブロックからのデータの読み出しを伴う未実行  
 または実行中の命令の数を示すカウンタを管理し、  
あるブロックを対象とするデータの読み出し処理を実行する場合、そのブロックに対応す  
 る前記カウンタの値を1つ加算し、あるブロックを対象とするデータの読み出し処理を終  
 了した場合、そのブロックに対応する前記カウンタの値を1つ減算し、  
前記複数の領域の中のいずれかの領域に割り当て済みの第1ブロックを前記複数の領域の  
 中の任意の領域に再割り当て可能な状態に遷移させる命令を前記ホストから受信した際、  
 前記第1ブロックに対応する前記カウンタの値が1以上である場合、前記カウンタの値が  
 0になってから前記第1ブロックを再割り当て可能な状態に遷移させる、  
 ーように構成されているメモリシステム。

【請求項2】

ホストに接続可能なメモリシステムであって、

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を制御するコントローラとを具備し、  
前記コントローラは、  
前記不揮発性メモリを論理的に分割することによって得られる複数の領域の中のいずれかの領域に割り当て済みのブロック毎に、各ブロックからのデータの読み出しを伴う未実行または実行中の命令の数を示すカウンタを管理し、  
あるブロックを対象とするデータの読み出し処理を実行する場合、そのブロックに対応する前記カウンタの値を1つ加算し、あるブロックを対象とするデータの読み出し処理を終了した場合、そのブロックに対応する前記カウンタの値を1つ減算し、  
前記複数の領域の中のいずれかの領域に割り当て済みの第2ブロックを前記複数の領域の中の任意の領域に再割り当て可能な状態に遷移させる命令を前記ホストから受信した際、前記第2ブロックに対応する前記カウンタの値が0でない場合、前記第2ブロックを再割り当て可能な状態に遷移させる命令に対する応答としてエラーを前記ホストに通知する、  
ように構成されているメモリシステム。

10

【請求項3】

前記コントローラは、前記複数の領域の中のいずれかの領域に割り当て済みの第3ブロックを前記複数の領域の中の任意の領域に再割り当て可能な状態に遷移させる命令が前記ホストから受信されている状況下において、前記第3ブロックからのデータの読み出しを伴う命令が前記ホストから受信された場合、前記第3ブロックからのデータの読み出しを伴う命令に対する応答としてエラーを前記ホストに通知するように構成されている請求項1または2に記載のメモリシステム。

20

【請求項4】

前記複数の領域は、各々が複数のブロックを含むQoSドメインである請求項1～3のいずれか1項に記載のメモリシステム。

【請求項5】

あるQoSドメインに割り当て済みのブロックを任意のQoSドメインに再割り当て可能な状態に遷移させる命令は、QoSドメインIDと、前記QoSドメインIDで示されるQoSドメインの中のブロックを示すブロックアドレスと、を含む請求項4に記載のメモリシステム。

30

【請求項6】

前記コントローラは、  
前記不揮発性メモリに含まれる複数のブロックの各々が1つのQoSドメインのみに属するように、前記複数のブロックを複数のQoSドメインに分類し、  
前記QoSドメイン毎に、未割り当てのブロックのリストと割り当て済みのブロックのリストとを管理する、  
ように構成されている請求項4または5に記載のメモリシステム。

【請求項7】

複数のブロックを含む不揮発性メモリを制御するコントローラによって実行される制御方法であって、  
前記不揮発性メモリを論理的に分割することによって得られる複数の領域の中のいずれかの領域に割り当て済みのブロック毎に、各ブロックからのデータの読み出しを伴う未実行または実行中の命令の数を示すカウンタを管理し、  
あるブロックを対象とするデータの読み出し処理を実行する場合、そのブロックに対応する前記カウンタの値を1つ加算し、あるブロックを対象とするデータの読み出し処理を終了した場合、そのブロックに対応する前記カウンタの値を1つ減算し、  
前記複数の領域の中のいずれかの領域に割り当て済みの第1ブロックを前記複数の領域の中の任意の領域に再割り当て可能な状態に遷移させる命令をホストから受信した際、前記第1ブロックに対応する前記カウンタの値が1以上である場合、前記カウンタの値が0になってから前記第1ブロックを再割り当て可能な状態に遷移させる、

40

50

制御方法。

【請求項 8】

複数のブロックを含む不揮発性メモリを制御するコントローラによって実行される制御方法であって、

前記不揮発性メモリを論理的に分割することによって得られる複数の領域の中のいずれかの領域に割り当て済みのブロック毎に、各ブロックからのデータの読み出しを伴う未実行または実行中の命令の数を示すカウンタを管理し、

あるブロックを対象とするデータの読み出し処理を実行する場合、そのブロックに対応する前記カウンタの値を1つ加算し、あるブロックを対象とするデータの読み出し処理を終了した場合、そのブロックに対応する前記カウンタの値を1つ減算し、

前記複数の領域の中のいずれかの領域に割り当て済みの第2ブロックを前記複数の領域の中の任意の領域に再割り当て可能な状態に遷移させる命令をホストから受信した際、前記第2ブロックに対応する前記カウンタの値が0でない場合、前記第2ブロックを再割り当て可能な状態に遷移させる命令に対する応答としてエラーを前記ホストに通知する、

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。

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。

最近では、ホストとストレージとの間の新たなインタフェースが提案され始めている。

【先行技術文献】

【非特許文献】

【0003】

【文献】Yiying Zhang, 外, "De-indirection for flash-based SSDs with nameless writes." FAST. 2012, [online], [平成29年9月13日検索], インターネット URL: <http://www.usenix.org/system/files/conference/fast12/zhang.pdf>

【発明の概要】

【発明が解決しようとする課題】

【0004】

しかし、一般に、NAND型フラッシュメモリの制御は複雑であるため、I/O性能を改善するための新たなインタフェースの実現に際しては、ホストとストレージ(メモリシステム)との間の適切な役割分担を考慮することが必要とされる。

本発明が解決しようとする課題は、I/O性能の改善を図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】

【0005】

実施形態によれば、ホストに接続可能なメモリシステムは、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を制御するコントローラとを含む。前記コントローラは、前記不揮発性メモリを論理的に分割することによって得られる複数の領域の中のいずれかの領域に割り当て済みのブロック毎に、各ブロックからのデータの読み出しを伴う未実行または実行中の命令の数を示すカウンタを管理する。前記コントローラは、あるブロックを対象とするデータの読み出し処理を実行する場合、そのブロックに対応する前記カウンタの値を1つ加算し、あるブロックを対象とするデータの読み出し処理を終了した場合、そのブロックに対応する前記カウンタの値を1つ減算する。前記コントローラは、前記複数の領域の中のいずれかの領域に割り当て済

10

20

30

40

50

みの第1ブロックを前記複数の領域の中の任意の領域に再割り当て可能な状態に遷移させる命令を前記ホストから受信した際、前記第1ブロックに対応する前記カウンタの値が1以上である場合、前記カウンタの値が0になってから前記第1ブロックを再割り当て可能な状態に遷移させる。

【図面の簡単な説明】

【0006】

【図1】ホストと実施形態のメモリシステム（フラッシュストレージデバイス）との関係を示すブロック図。

【図2】従来型SSDとホストとの間の役割分担と、同実施形態のフラッシュストレージデバイスとホストとの間の役割分担とを説明するための図。

【図3】複数のホストと複数のフラッシュストレージデバイスとの間のデータ転送がネットワーク機器を介して実行される計算機システムの構成例を示すブロック図。

【図4】同実施形態のメモリシステムの構成例を示すブロック図。

【図5】同実施形態のメモリシステムに設けられたNANDインタフェースと複数のNAND型フラッシュメモリダイとの関係を示すブロック図。

【図6】複数のブロックの集合によって構築されるスーパーブロックの構成例を示す図。

【図7】ホストが論理アドレスとブロック番号とを指定し且つ同実施形態のメモリシステムがブロック内物理アドレス（ブロック内オフセット）を決定するデータ書き込み動作と、ホストがブロック番号とブロック内物理アドレス（ブロック内オフセット）とを指定するデータ読み出し動作とを説明するための図。

【図8】同実施形態のメモリシステムに適用されるライトコマンドを説明するための図。

【図9】図8のライトコマンドに対するレスポンスを説明するための図。

【図10】同実施形態のメモリシステムに適用されるTrimコマンドを説明するための図。

【図11】物理アドレスを表す、ブロック番号およびオフセットを説明するための図。

【図12】ライトコマンドに応じて実行される書き込み動作を説明するための図。

【図13】不良ページをスキップする書き込み動作を説明するための図。

【図14】不良ページをスキップする書き込み動作の別の例を説明するための図。

【図15】論理アドレスとデータのペアをブロック内のページに書き込む動作を説明するための図。

【図16】データをブロック内のページのユーザデータ領域に書き込み、このデータの論理アドレスをこのページの冗長領域に書き込む動作を説明するための図。

【図17】スーパーブロックが使用される場合におけるブロック番号とオフセットとの関係を説明するための図。

【図18】同実施形態のメモリシステムに適用される最大ブロック番号ゲットコマンドを説明するための図。

【図19】最大ブロック番号ゲットコマンドに対するレスポンスを説明するための図。

【図20】同実施形態のメモリシステムに適用されるブロックサイズゲットコマンドを説明するための図。

【0007】

【図21】ブロックサイズゲットコマンドに対するレスポンスを説明するための図。

【図22】同実施形態のメモリシステムに適用されるブロックアロケートコマンド（ブロック割り当て要求）を説明するための図。

【図23】ブロックアロケートコマンドに対するレスポンスを説明するための図。

【図24】ホストと同実施形態のメモリシステムとによって実行されるブロック情報取得処理を示すシーケンスチャート。

【図25】ホストと同実施形態のメモリシステムとによって実行される書き込み処理のシーケンスを示すシーケンスチャート。

【図26】すでに書き込まれているデータに対する更新データを書き込むデータ更新動作を示す図。

10

20

30

40

50

【図 2 7】同実施形態のメモリシステムによって管理されるブロック管理テーブルを更新する動作を説明するための図。

【図 2 8】ホストによって管理されるルックアップテーブル（論理物理アドレス変換テーブル）を更新する動作を説明するための図。

【図 2 9】無効化すべきデータに対応するブロック番号および物理アドレスを示すホストからの通知に応じてブロック管理テーブルを更新する動作を説明するための図。

【図 3 0】同実施形態のメモリシステムに適用されるリードコマンドを説明するための図。

【図 3 1】同実施形態のメモリシステムによって実行されるリード動作を説明するための図。

【図 3 2】ホストからのリードコマンドに応じて、異なる物理記憶位置にそれぞれ格納されているデータ部をリードする動作を説明するための図。

10

【図 3 3】ホストと同実施形態のメモリシステムとによって実行されるリード処理のシーケンスを示すシーケンスチャート。

【図 3 4】同実施形態のメモリシステムに適用されるガベージコレクション（GC）制御コマンドを説明するための図。

【図 3 5】同実施形態のメモリシステムに適用される GC 用コールバックコマンドを説明するための図。

【図 3 6】ホストと同実施形態のメモリシステムとによって実行されるガベージコレクション（GC）動作の手順を示すシーケンスチャート。

【図 3 7】ガベージコレクション（GC）のために実行されるデータコピー動作の例を説明するための図。

20

【図 3 8】図 3 7 のデータコピー動作の結果に基づいて更新されるホストのルックアップテーブルの内容を説明するための図。

【図 3 9】ホストと同実施形態のメモリシステムとのシステムアーキテクチャを示す図。

【図 4 0】同実施形態のメモリシステム上における仮想ストレージデバイスの定義例を示す図。

【0 0 0 8】

【図 4 1】同実施形態のメモリシステム上において仮想ストレージデバイス毎に QoS ドメインが管理される例を示す図。

【図 4 2】同実施形態のメモリシステムのブロックリユースコマンド受信時におけるフラッシュストレージデバイス 3 の動作手順（第 1 ケース）を示すフローチャート。

30

【図 4 3】同実施形態のメモリシステムのブロックリユースコマンド受信時におけるフラッシュストレージデバイス 3 の動作手順（第 2 ケース）を示すフローチャート。

【図 4 4】同実施形態のメモリシステムによって実行される I/O コマンド処理を説明するためのブロック図。

【図 4 5】同実施形態のメモリシステムによって実行される複数段階の書き込み動作を説明するための図。

【図 4 6】同実施形態のメモリシステム内のある書き込み先ブロックへのデータの書き込み順序を説明するための図。

【図 4 7】不揮発性メモリのデータ書き込み単位と同じサイズの単位でライトデータをホストから同実施形態のメモリシステムに転送する動作を説明するための図。

40

【図 4 8】同実施形態のメモリシステムによって実行されるデータ書き込み処理の手順を示すフローチャート。

【図 4 9】同実施形態のメモリシステムによって実行されるデータ書き込み処理の別の手順を示すフローチャート。

【図 5 0】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への解放可能通知の送信処理の手順を示すフローチャート。

【図 5 1】ホストによって実行されるライトデータ破棄処理の手順を示すフローチャート。

【図 5 2】最後のライトコマンドが受信されてから闕期間、次のライトコマンドが受信されない場合に、同実施形態のメモリシステムによって実行されるダミーデータ書き込み処

50

理を説明するための図。

【図 5 3】同実施形態のメモリシステムによって実行されるダミーデータ書き込み処理の手順を示すフローチャート。

【図 5 4】内部バッファを使用して同実施形態のメモリシステムによって実行されるデータ転送動作を示すブロック図。

【図 5 5】内部バッファを使用して同実施形態のメモリシステムによって実行されるデータ書き込み処理の手順を示すフローチャート。

【図 5 6】同実施形態のメモリシステムによって実行されるデータ読み出し処理の手順を示すフローチャート。

【図 5 7】同実施形態のメモリシステムに適用されるブロックリユースコマンドコマンドを説明するための図。

10

【図 5 8】同実施形態のメモリシステムに適用されるライトコマンドの別の例を説明するための図。

【図 5 9】図 5 8 のライトコマンドに対するレスポンスを説明するための図。

【図 6 0】ホストと同実施形態のメモリシステムとによって実行される書き込み動作処理のシーケンスの別の例を示すシーケンスチャート。

【0 0 0 9】

【図 6 1】同実施形態のメモリシステムに適用されるガベージコレクション ( G C ) 制御コマンドの別の例を説明するための図。

【図 6 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 ( G C ) 動作の手順の別の例を示すシーケンスチャート。

20

【図 6 3】同実施形態のメモリシステムにおける書き込み先ブロックの割り当ての手順を示すフローチャート。

【図 6 4】同実施形態のメモリシステムにおける G C デスティネーションブロックの割り当ての手順を示すフローチャート。

【発明を実施するための形態】

【0 0 1 0】

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計算機システムの構成を説明する。

30

このメモリシステムは、不揮発性メモリにデータを書き込み、不揮発性メモリからデータを読み出す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、NANDフラッシュ技術ベースのフラッシュストレージデバイス 3 として実現されている。

【0 0 1 1】

この計算機システムは、ホスト ( ホストデバイス ) 2 と、複数のフラッシュストレージデバイス 3 とを含んでいてもよい。ホスト 2 は、複数のフラッシュストレージデバイス 3 によって構成されるフラッシュアレイをストレージとして使用するよう構成されたサーバであってもよい。ホスト ( サーバ ) 2 と複数のフラッシュストレージデバイス 3 は、インタフェース 5 0 を介して相互接続される ( 内部相互接続 ) 。この内部相互接続のためのインタフェース 5 0 としては、これに限定されないが、P C I E x p r e s s ( P C I e ) ( 登録商標 ) 、 N V M E x p r e s s ( N V M e ) ( 登録商標 ) 、 E t h e r n e t ( 登録商標 ) 、 N V M e o v e r F a b r i c s ( N V M e O F ) 等を使用し得る。

40

【0 0 1 2】

ホスト 2 として機能するサーバの典型例としては、データセンター内のサーバが挙げられる。

ホスト 2 がデータセンター内のサーバによって実現されるケースにおいては、このホスト ( サーバ ) 2 は、ネットワーク 5 1 を介して複数のエンドユーザ端末 ( クライアント ) 6 1 に接続されてもよい。ホスト 2 は、これらエンドユーザ端末 6 1 に対して様々なサービスを提供することができる。

【0 0 1 3】

50

ホスト（サーバ）2によって提供可能なサービスの例には、（1）システム稼働プラットフォームを各クライアント（各エンドユーザ端末61）に提供するプラットフォーム・アズ・ア・サービス（PaaS）、（2）仮想サーバのようなインフラストラクチャを各クライアント（各エンドユーザ端末61）に提供するインフラストラクチャ・アズ・ア・サービス（IaaS）、等がある。

【0014】

複数の仮想マシンが、このホスト（サーバ）2として機能する物理サーバ上で実行されてもよい。ホスト（サーバ）2上で走るこれら仮想マシンの各々は、対応する幾つかのクライアント（エンドユーザ端末61）に各種サービスを提供するように構成された仮想サーバとして機能することができる。

10

【0015】

ホスト（サーバ）2は、フラッシュアレイを構成する複数のフラッシュストレージデバイス3を管理するストレージ管理機能と、エンドユーザ端末61それぞれに対してストレージアクセスを含む様々なサービスを提供するフロントエンド機能とを含む。

従来型SSDにおいては、NAND型フラッシュメモリのブロック/ページの階層構造はSSD内のフラッシュトランスレーション層（FTL）によって隠蔽されている。つまり、従来型SSDのFTLは、（1）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を使用して、論理アドレスそれぞれとNAND型フラッシュメモリの物理アドレスそれぞれとの間のマッピングを管理する機能、（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機能、（3）NAND型フラッシュメモリのガベージコレクション（GC）を実行する機能、等を有している。論理アドレスそれぞれとNAND型フラッシュメモリの物理アドレスとの間のマッピングは、ホストからは見えない。NAND型フラッシュメモリのブロック/ページ構造もホストからは見えない。

20

【0016】

一方、ホストにおいても、一種のアドレス変換（アプリケーションレベルアドレス変換）が実行されることがある。このアドレス変換は、アプリケーションレベルアドレス変換テーブルを使用して、アプリケーションレベルの論理アドレスそれぞれとSSD用の論理アドレスそれぞれとの間のマッピングを管理する。また、ホストにおいても、SSD用の論理アドレス空間上に生じるフラグメントの解消のために、この論理アドレス空間上のデータ配置を変更する一種のGC（アプリケーションレベルGC）が実行される。

30

【0017】

しかし、ホストおよびSSDがそれぞれアドレス変換テーブルを有するという冗長な構成（SSD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を有し、ホストはアプリケーションレベルアドレス変換テーブルを有する）においては、これらアドレス変換テーブルを保持するために膨大なメモリリソースが消費される。さらに、ホスト側のアドレス変換とSSD側のアドレス変換とを含む2重のアドレス変換は、I/O性能を低下させる要因にもなる。

【0018】

さらに、ホスト側のアプリケーションレベルGCは、SSDへのデータ書き込み量を実際のユーザデータ量の数倍（例えば2倍）程度に増やす要因となる。このようなデータ書き込み量の増加は、SSDのライトアンプリフィケーションとあいまってシステム全体のストレージ性能を低下させ、またSSDの寿命も短くする。

40

【0019】

このような問題点を解消するために、従来型SSDのFTLの機能の全てをホストに移すという対策も考えられる。

しかし、この対策を実装するためには、NAND型フラッシュメモリのブロックおよびページをホストが直接的にハンドリングすることが必要となる。NAND型フラッシュメモリにおいては、ページ書き込み順序制約があるため、ホストがページを直接ハンドリングすることは困難である。また、NAND型フラッシュメモリにおいては、ブロックが不良ページ（バッドページ）を含む場合がある。バッドページをハンドリングすることはホス

50

トにとってはなおさら困難である。

#### 【 0 0 2 0 】

そこで、本実施形態では、F T Lの役割はホスト2とフラッシュストレージデバイス3との間で分担される。概していえば、ホスト2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を管理するが、ホスト2はデータが書き込まれるべきブロックのブロック番号とこのデータに対応する論理アドレスだけを指定し、このデータが書き込まれるべきこのブロック内の位置（書き込み先位置）はフラッシュストレージデバイス3によって決定される。決定されたこのブロック内の位置（書き込み先位置）を示すブロック内物理アドレスは、フラッシュストレージデバイス3からホスト2に通知される。

#### 【 0 0 2 1 】

このように、ホスト2はブロックのみをハンドリングし、ブロック内の位置（例えば、ページ、ページ内の位置）はフラッシュストレージデバイス3によってハンドリングされる。フラッシュストレージデバイス3にデータを書き込む必要がある時、ホスト2は、ブロック番号を選択（またはフラッシュストレージデバイス3にフリーブロックを割り当てるように要求）し、論理アドレスと、選択したブロックのブロック番号（またはフラッシュストレージデバイス3によって通知される割り当てられたブロックのブロック番号）とを指定するライト要求（ライトコマンド）をフラッシュストレージデバイス3に送信する。フラッシュストレージデバイス3は、指定されたブロック番号を有するブロックにホスト2からのデータを書き込む。この場合、フラッシュストレージデバイス3は、このブロック内の位置（書き込み先位置）を決定し、ホスト2からのデータをこのブロック内の位置（書き込み先位置）に書き込む。そして、フラッシュストレージデバイス3は、このブロック内の位置（書き込み先位置）を示すブロック内物理アドレスを、ライト要求に対するレスポンス（返り値）としてホスト2に通知する。以下では、ホスト2に移されたF T L機能をグローバルF T Lと称する。

#### 【 0 0 2 2 】

ホスト2のグローバルF T Lは、ストレージサービスを実行する機能、ウェア制御機能、高可用性を実現するための機能、同じ内容を有する複数の重複データ部がストレージに格納されることを防止する重複排除（D e - d u p l i c a t i o n）機能、ガベージコレクション（G C）ブロック選択機能、Q o S制御機能等を有する。Q o S制御機能には、Q o Sドメイン毎（またはブロック毎）にアクセス単位を決める機能が含まれる。アクセス単位は、ホスト2がライト/リードすることが可能な最小データサイズ（G r a i n）を示す。フラッシュストレージデバイス3は単一、あるいは複数のアクセス単位（G r a i n）をサポートしており、ホスト2は、フラッシュストレージデバイス3が複数のアクセス単位をサポートしている場合にはQ o Sドメイン毎（またはブロック毎）に、使用すべきアクセス単位をフラッシュストレージデバイス3に指示することができる。

#### 【 0 0 2 3 】

また、Q o S制御機能には、Q o Sドメイン間の性能干渉をできるだけ防ぐための機能が含まれている。この機能は、基本的には、安定したレイテンシを保つための機能である。一方、フラッシュストレージデバイス3は、ローレベルアブストラクション（L L A）を実行することができる。L L AはN A N D型フラッシュメモリのアブストラクションのための機能である。L L Aは、不良ページ（バッドページ）を隠蔽する機能、ページ書き込み順序制約を守る機能を含む。L L Aは、G C実行機能も含む。G C実行機能は、ホスト2によって指定されたコピー元ブロック（G Cソースブロック）内の有効データを、ホスト2によって指定されたコピー先ブロック（G Cデスティネーションブロック）にコピーする。フラッシュストレージデバイス3のG C実行機能は、有効データを書き込むべきG Cデスティネーションブロック内の位置（コピー先位置）を決定し、G Cソースブロック内の有効データを、G Cデスティネーションブロック内のコピー先位置にコピーする。

#### 【 0 0 2 4 】

図2は、従来型S S Dとホストとの間の役割分担と、本実施形態のフラッシュストレージデバイス3とホスト2との間の役割分担とを示す。

10

20

30

40

50



図 2 の左部は、従来型 S S D と仮想ディスクサービスを実行するホストとを含む計算機システム全体の階層構造を表している。

【 0 0 2 5 】

ホスト（サーバ）においては、複数のエンドユーザに複数の仮想マシンを提供するための仮想マシンサービス 1 0 1 が実行される。仮想マシンサービス 1 0 1 上の各仮想マシンにおいては、対応するエンドユーザによって使用されるオペレーティングシステムおよびユーザアプリケーション 1 0 2 が実行される。

【 0 0 2 6 】

また、ホスト（サーバ）においては、複数のユーザアプリケーション 1 0 2 に対応する複数の仮想ディスクサービス 1 0 3 が実行される。各仮想ディスクサービス 1 0 3 は、従来型 S S D 内のストレージリソースの容量の一部を、対応するユーザアプリケーション 1 0 2 用のストレージリソース（仮想ディスク）として割り当てる。各仮想ディスクサービス 1 0 3 においては、アプリケーションレベルアドレス変換テーブルを使用して、アプリケーションレベルの論理アドレスを S S D 用の論理アドレスに変換するアプリケーションレベルアドレス変換も実行される。さらに、ホストにおいては、アプリケーションレベル G C 1 0 4 も実行される。

【 0 0 2 7 】

ホスト（サーバ）から従来型 S S D へのコマンドの送信および従来型 S S D からホスト（サーバ）へのコマンド完了のレスポンスの返送は、ホスト（サーバ）および従来型 S S D の各々に存在する I / O キュー 2 0 0 を介して実行される。

従来型 S S D は、ライトバッファ（W B ） 3 0 1、ルックアップテーブル（L U T ） 3 0 2、ガベージコレクション機能 3 0 3、N A N D 型フラッシュメモリ（N A N D フラッシュユアレイ） 3 0 4 を含む。従来型 S S D は、一つのルックアップテーブル（L U T ） 3 0 2 のみを管理しており、N A N D 型フラッシュメモリ（N A N D フラッシュユアレイ） 3 0 4 のリソースは複数の仮想ディスクサービス 1 0 3 によって共有される。

【 0 0 2 8 】

この構成においては、仮想ディスクサービス 1 0 3 下のアプリケーションレベル G C 1 0 4 と従来型 S S D 内のガベージコレクション機能 3 0 3（L U T レベル G C）とを含む重複した G C により、ライトアンプリフィケーションが大きくなる。また、従来型 S S D においては、あるエンドユーザまたはある仮想ディスクサービス 1 0 3 からのデータ書き込み量の増加によって G C の頻度が増加し、これによって他のエンドユーザまたは他の仮想ディスクサービス 1 0 3 に対する I / O 性能が劣化するというノイジーネイバー問題が生じうる。

【 0 0 2 9 】

また、各仮想ディスクサービス内のアプリケーションレベルアドレス変換テーブルと従来型 S S D 内の L U T 3 0 2 とを含む重複したリソースの存在により、多くのメモリリソースが消費される。

図 2 の右部は、本実施形態のフラッシュストレージデバイス 3 とホスト 2 とを含む計算機システム全体の階層構造を表している。

【 0 0 3 0 】

ホスト（サーバ） 2 においては、複数のエンドユーザに複数の仮想マシンを提供するための仮想マシンサービス 4 0 1 が実行される。仮想マシンサービス 4 0 1 上の各仮想マシンにおいては、対応するエンドユーザによって使用されるオペレーティングシステムおよびユーザアプリケーション 4 0 2 が実行される。

【 0 0 3 1 】

また、ホスト（サーバ） 2 においては、複数のユーザアプリケーション 4 0 2 に対応する複数の I / O サービス 4 0 3 が実行される。これら I / O サービス 4 0 3 には、L B A ベースのブロック I / O サービス、キー・バリュー・ストアサービスなどが含まれてもよい。各 I / O サービス 4 0 3 は、論理アドレスそれぞれとフラッシュストレージデバイス 3 の物理アドレスそれぞれとの間のマッピングを管理するルックアップテーブル（L U T）

10

20

30

40

50

4 1 1を含む。ここで、論理アドレスとは、アクセス対象のデータを識別可能な識別子を意味する。この論理アドレスは、論理アドレス空間上の位置を指定する論理ブロックアドレス（LBA）であってもよいし、あるいは、キー・バリュー・ストアのキー（タグ）であってもよいし、キーのハッシュ値であってもよい。

【0032】

LBAベースのブロックI/Oサービスにおいては、論理アドレス（LBA）それぞれとフラッシュストレージデバイス3の物理アドレスそれぞれとの間のマッピングを管理するLUT411が使用されてもよい。

キー・バリュー・ストアサービスにおいては、論理アドレス（つまり、キーのようなタグ）それぞれとこれら論理アドレス（つまり、キーのようなタグ）に対応するデータが格納されているフラッシュストレージデバイス3内の物理記憶位置を示す物理アドレスそれぞれとの間のマッピングを管理するLUT411が使用されてもよい。このLUT411においては、タグと、このタグによって識別されるデータが格納されている物理アドレスと、このデータのデータ長との対応関係が管理されてもよい。

10

【0033】

各エンドユーザは、使用すべきアドレッシング方法（LBA、キー・バリュー・ストアのキー、等）を選択することができる。

これら各LUT411は、ユーザアプリケーション402からの論理アドレスそれぞれをフラッシュストレージデバイス3用の論理アドレスそれぞれに変換するのではなく、ユーザアプリケーション402からの論理アドレスそれぞれをフラッシュストレージデバイス3の物理アドレスそれぞれに変換する。つまり、これら各LUT411は、フラッシュストレージデバイス3用の論理アドレスを物理アドレスに変換するテーブルとアプリケーションレベルアドレス変換テーブルとが統合（マージ）されたテーブルである。

20

【0034】

また、各I/Oサービス403は、GCブロック選択機能を含む。GCブロック選択機能は、対応するLUTを使用して各ブロックの有効データ量を管理することができ、これによってGCソースブロックを選択することができる。

ホスト（サーバ）2においては、上述のQoSドメイン毎にI/Oサービス403が存在してもよい。あるQoSドメインに属するI/Oサービス403は、対応するQoSドメイン内のユーザアプリケーション402によって使用される論理アドレスそれぞれと対応するQoSドメインに割り当てられたリソースグループに属するブロック群のブロック番号それぞれとの間のマッピングを管理してもよい。

30

【0035】

ホスト（サーバ）2からフラッシュストレージデバイス3へのコマンドの送信およびフラッシュストレージデバイス3からホスト（サーバ）2へのコマンド完了のレスポンス等の返送は、ホスト（サーバ）2およびフラッシュストレージデバイス3の各々に存在するI/Oキュー500を介して実行される。これらI/Oキュー500も、複数のQoSドメインに対応する複数のキューグループに分類されていてもよい。

【0036】

フラッシュストレージデバイス3は、複数のQoSドメインに対応する複数のライトバッファ（WB）601、複数のQoSドメインに対応する複数のガベージコレクション（GC）機能602、NAND型フラッシュメモリ（NANDフラッシュアレイ）603を含む。

40

【0037】

この図2の右部に示す構成においては、上位階層（ホスト2）はブロック境界を認識することができるので、ブロック境界/ブロックサイズを考慮してユーザデータを各ブロックに書き込むことができる。つまり、ホスト2はNAND型フラッシュメモリ（NANDフラッシュアレイ）603の個々のブロックを認識することができ、これにより、例えば、一つのブロック全体に一齐にデータを書き込む、一つのブロック内のデータ全体を削除または更新によって無効化する、といった制御を行うことが可能となる。この結果、一つの

50

ブロックに有効データと無効データが混在されるという状況を起こりにくくすることが可能となる。したがって、GCを実行することが必要となる頻度を低減することができる。GCの頻度を低減することにより、ライトアンプリフィケーションが低下され、フラッシュストレージデバイス3の性能の向上、フラッシュストレージデバイス3の寿命の最大化を実現できる。このように、上位階層（ホスト2）がブロック番号を認識可能な構成は有用である。

**【0038】**

一方、データが書き込まれるべきブロック内の位置は、上位階層（ホスト2）ではなく、フラッシュストレージデバイス3によって決定される。したがって、不良ページ（バッドページ）を隠蔽することができ、またページ書き込み順序制約を守ることができる。

10

**【0039】**

図3は、図1のシステム構成の変形例を示す。

図3においては、複数のホスト2Aと複数のフラッシュストレージデバイス3との間のデータ転送がネットワーク機器（ここでは、ネットワークスイッチ1）を介して実行される。すなわち、図3の計算機システムにおいては、図1のホスト（サーバ）2のストレージ管理機能がマネージャ2Bに移され、且つホスト（サーバ）2のフロントエンド機能が複数のホスト（エンドユーザサービス用ホスト）2Aに移されている。

**【0040】**

マネージャ2Bは、複数のフラッシュストレージデバイス3を管理し、各ホスト（エンドユーザサービス用ホスト）2Aからの要求に応じて、これらフラッシュストレージデバイス3のストレージリソースを各ホスト（エンドユーザサービス用ホスト）2Aに割り当てる。

20

**【0041】**

各ホスト（エンドユーザサービス用ホスト）2Aは、ネットワークを介して一つ以上のエンドユーザ端末61に接続される。各ホスト（エンドユーザサービス用ホスト）2Aは、上述の統合（マージ）された論理物理アドレス変換テーブルであるルックアップテーブル（LUT）を管理する。各ホスト（エンドユーザサービス用ホスト）2Aは、自身のLUTを使用して、対応するエンドユーザによって使用される論理アドレスそれぞれと自身に割り当てられたリソースの物理アドレスそれぞれとの間のマッピングのみを管理する。したがって、この構成は、システムを容易にスケールアウトすることを可能にする。

30

**【0042】**

各ホスト2AのグローバルFTLは、ルックアップテーブル（LUT）を管理する機能、高可用性を実現するための機能、QoS制御機能、GCブロック選択機能等を有する。マネージャ2Bは、複数のフラッシュストレージデバイス3を管理するための専用のデバイス（計算機）である。マネージャ2Bは、各ホスト2Aから要求された容量分のストレージリソースを予約するグローバルリソース予約機能を有する。さらに、マネージャ2Bは、各フラッシュストレージデバイス3の消耗度を監視するためのウェア監視機能、予約されたストレージリソース（NANDリソース）を各ホスト2Aに割り当てるNANDリソース割り当て機能、QoS制御機能、グローバルロック管理機能、等を有する。

**【0043】**

各フラッシュストレージデバイス3のローレベルアブストラクション（LLA）は、不良ページ（バッドページ）を隠蔽する機能、ページ書き込み順序制約を守る機能、ライトバッファを管理する機能、GC実行機能等を有する。

40

図3のシステム構成によれば、各フラッシュストレージデバイス3の管理はマネージャ2Bによって実行されるので、各ホスト2Aは、自身に割り当てられた一つ以上のフラッシュストレージデバイス3にI/O要求を送信する動作と、フラッシュストレージデバイス3からのレスポンスを受信するという動作とのみを実行すればよい。つまり、複数のホスト2Aと複数のフラッシュストレージデバイス3との間のデータ転送はネットワークスイッチ1のみを介して実行され、マネージャ2Bはこのデータ転送には関与しない。また、上述したように、ホスト2Aそれぞれによって管理されるLUTの内容は互いに独立して

50

いる。よって、容易にホスト 2 A の数を増やすことができるので、スケールアウト型のシステム構成を実現することができる。

【 0 0 4 4 】

図 4 は、フラッシュストレージデバイス 3 の構成例を示す。

フラッシュストレージデバイス 3 は、コントローラ 4 および不揮発性メモリ ( N A N D 型フラッシュメモリ ) 5 を備える。フラッシュストレージデバイス 3 は、ランダムアクセスメモリ、例えば、D R A M 6 も備えていてもよい。

【 0 0 4 5 】

N A N D 型フラッシュメモリ 5 は、マトリクス状に配置された複数のメモリセルを含むメモリセルアレイを含む。N A N D 型フラッシュメモリ 5 は、2 次元構造の N A N D 型フラッシュメモリであってもよいし、3 次元構造の N A N D 型フラッシュメモリであってもよい。

10

【 0 0 4 6 】

N A N D 型フラッシュメモリ 5 のメモリセルアレイは、複数のブロック B L K 0 ~ B L K m - 1 を含む。ブロック B L K 0 ~ B L K m - 1 の各々は多数のページ ( ここではページ P 0 ~ P n - 1 ) によって編成される。ブロック B L K 0 ~ B L K m - 1 は、消去単位として機能する。ブロックは、「消去ブロック」、「物理ブロック」、または「物理消去ブロック」と称されることもある。ページ P 0 ~ P n - 1 の各々は、同一ワード線に接続された複数のメモリセルを含む。ページ P 0 ~ P n - 1 は、データ書き込み動作およびデータ読み込み動作の単位である。

20

コントローラ 4 は、T o g g l e、オープン N A N D フラッシュインタフェース ( O N F I ) のような N A N D インタフェース 1 3 を介して、不揮発性メモリである N A N D 型フラッシュメモリ 5 に電氣的に接続されている。コントローラ 4 は、N A N D 型フラッシュメモリ 5 を制御するように構成されたメモリコントローラ ( 制御回路 ) である。

【 0 0 4 7 】

N A N D 型フラッシュメモリ 5 は、図 5 に示すように、複数の N A N D 型フラッシュメモリダイを含む。各 N A N D 型フラッシュメモリダイは、複数のブロック B L K を含むメモリセルアレイとこのメモリセルアレイを制御する周辺回路とを含む不揮発性メモリダイである。個々の N A N D 型フラッシュメモリダイは独立して動作可能である。このため、N A N D 型フラッシュメモリダイは、並列動作単位として機能する。N A N D 型フラッシュメモリダイは、「N A N D 型フラッシュメモリチップ」または「不揮発性メモリチップ」とも称される。図 5 においては、N A N D インタフェース 1 3 に 1 6 個のチャンネル C h 1、C h 2、... C h 1 6 が接続されており、これらチャンネル C h 1、C h 2、... C h 1 6 の各々に、同数 ( 例えばチャンネル当たり 2 個のダイ ) の N A N D 型フラッシュメモリダイそれぞれが接続されている場合が例示されている。各チャンネルは、対応する N A N D 型フラッシュメモリダイと通信するための通信線 ( メモリバス ) を含む。

30

【 0 0 4 8 】

コントローラ 4 は、チャンネル C h 1、C h 2、... C h 1 6 を介して N A N D 型フラッシュメモリダイ # 1 ~ # 3 2 を制御する。コントローラ 4 は、チャンネル C h 1、C h 2、... C h 1 6 を同時に駆動することができる。

40

チャンネル C h 1 ~ C h 1 6 に接続された 1 6 個の N A N D 型フラッシュメモリダイ # 1 ~ # 1 6 は第 1 のバンクとして編成されてもよく、またチャンネル C h 1 ~ C h 1 6 に接続された残りの 1 6 個の N A N D 型フラッシュメモリダイ # 1 7 ~ # 3 2 は第 2 のバンクとして編成されてもよい。バンクは、複数のメモリモジュールをバンクインタリーブによって並列動作させるための単位として機能する。図 5 の構成例においては、1 6 チャンネルと、2 つのバンクを使用したバンクインタリーブとによって、最大 3 2 個の N A N D 型フラッシュメモリダイを並列動作させることができる。

【 0 0 4 9 】

本実施形態では、コントローラ 4 は、各々が複数のブロック B L K から構成される複数のブロック ( 以下、スーパーブロックと称する ) を管理してもよく、スーパーブロックの単

50

位で消去動作を実行してもよい。

スーパーブロックは、これに限定されないが、NAND型フラッシュメモリダイ#1～#32から一つずつ選択される計32個のブロックBLKを含んでいてもよい。なお、NAND型フラッシュメモリダイ#1～#32の各々はマルチプレーン構成を有していてもよい。例えば、NAND型フラッシュメモリダイ#1～#32の各々が、2つのプレーンを含むマルチプレーン構成を有する場合には、一つのスーパーブロックは、NAND型フラッシュメモリダイ#1～#32に対応する64個のプレーンから一つずつ選択される計64個のブロックBLKを含んでいてもよい。図6には、一つのスーパーブロックSBが、NAND型フラッシュメモリダイ#1～#32から一つずつ選択される計32個のブロックBLK(図5においては太枠で囲まれているブロックBLK)から構成される場合が例示されている。

10

#### 【0050】

図4に示されているように、コントローラ4は、ホストインタフェース11、CPU12、NANDインタフェース13、およびDRAMインタフェース14等を含む。これらCPU12、NANDインタフェース13、DRAMインタフェース14は、バス10を介して相互接続される。

#### 【0051】

このホストインタフェース11は、ホスト2との通信を実行するように構成されたホストインタフェース回路である。このホストインタフェース11は、例えば、PCIeコントローラ(NVMeコントローラ)であってもよい。ホストインタフェース11は、ホスト2から様々な要求(コマンド)を受信する。これら要求(コマンド)には、ライト要求(ライトコマンド)、リード要求(リードコマンド)、他の様々な要求(コマンド)が含まれる。

20

#### 【0052】

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を制御するように構成されたプロセッサである。CPU12は、フラッシュストレージデバイス3の電源オンにตอบสนองしてNAND型フラッシュメモリ5または図示しないROMから制御プログラム(ファームウェア)をDRAM6にロードし、そしてこのファームウェアを実行することによって様々な処理を行う。なお、ファームウェアはコントローラ4内の図示しないSRAM上にロードされてもよい。このCPU12は、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ことができる。CPU12の動作は、CPU12によって実行される上述のファームウェアによって制御される。なお、コマンド処理の一部または全部は、コントローラ4内の専用ハードウェアによって実行してもよい。

30

#### 【0053】

CPU12は、ライト動作制御部21、リード動作制御部22、およびGC動作制御部23として機能することができる。これらライト動作制御部21、リード動作制御部22、およびGC動作制御部23においては、図2の右部に示すシステム構成を実現するためのアプリケーションプログラムインタフェース(API)が実装されている。

#### 【0054】

ライト動作制御部21は、ブロック番号と論理アドレスを指定するライト要求(ライトコマンド)をホスト2から受信する。論理アドレスは、書き込むべきデータ(ユーザデータ)を識別可能な識別子であり、例えば、LBAであってもよいし、あるいはキー・バリュー・ストアのキーのようなタグであってもよいし、キーのハッシュ値であってもよい。ブロック番号は、このデータが書き込まれるべきブロックを指定する識別子である。ブロック番号としては、複数のブロック内の任意の一つを一意に識別可能な様々な値を使用し得る。ブロック番号によって指定されるブロックは、物理ブロックであってもよいし、上述のスーパーブロックであってもよい。ライトコマンドを受信した場合、ライト動作制御部21は、まず、ホスト2からのデータを書き込むべき、この指定されたブロック番号を有するブロック(書き込み先ブロック)内の位置(書き込み先位置)を決定する。次いで、

40

50

ライト動作制御部 2 1 は、ホスト 2 からのデータ（ライトデータ）を、この書き込み先ブロックの書き込み先位置に書き込む。この場合、ライト動作制御部 2 1 は、ホスト 2 からのデータのみならず、このデータとこのデータの論理アドレスの双方を書き込み先ブロックに書き込むことができる。

【 0 0 5 5 】

そして、ライト動作制御部 2 1 は、この書き込み先ブロックの上述の書き込み先位置を示すブロック内物理アドレスをホスト 2 に通知する。このブロック内物理アドレスは、この書き込み先ブロック内の書き込み先位置を示すブロック内オフセットによって表される。

【 0 0 5 6 】

この場合、このブロック内オフセットは、書き込み先ブロックの先頭から書き込み先位置までのオフセット、つまり書き込み先ブロックの先頭に対する書き込み先位置のオフセットを示す。書き込み先ブロックの先頭から書き込み先位置までのオフセットのサイズは、ページサイズとは異なるサイズを有する粒度（Grain）の倍数で示される。粒度（Grain）は、上述のアクセス単位である。粒度（Grain）のサイズの最大値は、ブロックサイズまでに制限される。換言すれば、ブロック内オフセットは、書き込み先ブロックの先頭から書き込み先位置までのオフセットをページサイズとは異なるサイズを有する粒度の倍数で示す。

10

【 0 0 5 7 】

粒度（Grain）は、ページサイズよりも小さいサイズを有していてもよい。例えば、ページサイズが 16 K バイトである場合、粒度（Grain）は、そのサイズが 4 K バイトであってもよい。この場合、ある一つのブロックにおいては、各々サイズが 4 K バイトである複数のオフセット位置が規定される。ブロック内の最初のオフセット位置に対応するブロック内オフセットは、例えば 0 であり、ブロック内の次のオフセット位置に対応するブロック内オフセットは、例えば 1 である、ブロック内のさらに次のオフセット位置に対応するブロック内オフセットは、例えば 2 である。

20

【 0 0 5 8 】

あるいは、粒度（Grain）は、ページサイズよりも大きなサイズを有していてもよい。例えば、粒度（Grain）は、ページサイズの数倍のサイズであってもよい。ページサイズが 16 K バイトである場合、粒度は、32 K バイトのサイズであってもよい。

【 0 0 5 9 】

このように、ライト動作制御部 2 1 は、ホスト 2 からのブロック番号を有するブロック内の書き込み先位置を自身で決定し、そしてホスト 2 からのライトデータをこのブロック内のこの書き込み先位置に書き込む。そして、ライト動作制御部 2 1 は、この書き込み先位置を示すブロック内物理アドレス（ブロック内オフセット）をライト要求に対応するレスポンス（返り値）としてホスト 2 に通知する。あるいは、ライト動作制御部 2 1 は、ブロック内物理アドレス（ブロック内オフセット）のみをホスト 2 に通知するのではなく、論理アドレスとブロック番号とブロック内物理アドレス（ブロック内オフセット）との組をホスト 2 に通知してもよい。

30

【 0 0 6 0 】

したがって、フラッシュストレージデバイス 3 は、ブロック番号をホスト 2 にハンドリングさせつつ、ページ書き込み順序制約、バッドページ、ページサイズ等を隠蔽することができる。

40

この結果、ホスト 2 は、ブロック境界は認識できるが、ページ書き込み順序制約、バッドページ、ページサイズについては意識することなく、どのユーザデータがどのブロック番号に存在するかを管理することができる。

【 0 0 6 1 】

リード動作制御部 2 2 は、物理アドレス（すなわち、ブロック番号およびブロック内オフセット）を指定するリード要求（リードコマンド）をホスト 2 から受信した場合、これらブロック番号およびブロック内オフセットに基づいて、リード対象のブロック内のリード対象の物理記憶位置からデータをリードする。リード対象のブロックは、ブロック番号に

50

よって特定される。このブロック内のリード対象の物理記憶位置は、ブロック内オフセットによって特定される。このブロック内オフセットを使用することにより、ホスト2は、NAND型フラッシュメモリの世代毎の異なるページサイズをハンドリングする必要がない。

#### 【0062】

リード対象の物理記憶位置を得るために、リード動作制御部22は、まず、このブロック内オフセットを、ページサイズを表す粒度の数（ページサイズが16Kバイトで粒度（Grain）が4Kバイトである場合には、ページサイズを表す粒度の数は4）で除算し、そしてこの除算によって得られる商および余りを、リード対象のページ番号およびリード対象のページ内オフセットとしてそれぞれ決定してもよい。

10

#### 【0063】

GC動作制御部23は、NAND型フラッシュメモリ5のガベージコレクションのためのコピー元ブロック番号（GCソースブロック番号）およびコピー先ブロック番号（GCデスティネーションブロック番号）を指定するGC制御コマンドをホスト2から受信した場合、NAND型フラッシュメモリ5の複数のブロックから、指定されたコピー元ブロック番号を有するブロックと指定されたコピー先ブロック番号を有するブロックとをコピー元ブロック（GCソースブロック）およびコピー先ブロック（GCデスティネーションブロック）として選択する。GC動作制御部23は、選択されたGCソースブロックに格納されている有効データを書き込むべきGCデスティネーションブロック内のコピー先位置を決定し、有効データをGCデスティネーションブロック内のコピー先位置にコピーする。

20

#### 【0064】

そして、GC動作制御部23は、有効データの論理アドレスと、コピー先ブロック番号と、GCデスティネーションブロック内のコピー先位置を示すブロック内物理アドレス（ブロック内オフセット）とを、ホスト2に通知する。

有効データ/無効データの管理は、ブロック管理テーブル32を使用して実行されてもよい。このブロック管理テーブル32は、例えば、ブロック毎に存在してもよい。あるブロックに対応するブロック管理テーブル32においては、このブロック内のデータそれぞれの有効/無効を示すビットマップフラグが格納されている。ここで、有効データとは、論理アドレスから最新のデータとして紐付けられているデータであって、後にホスト2からリードされる可能性があるデータを意味する。無効データとは、もはやホスト2からリードされる可能性が無いデータを意味する。例えば、ある論理アドレスに関連付けられているデータは有効データであり、どの論理アドレスにも関連付けられていないデータは無効データである。

30

#### 【0065】

上述したように、GC動作制御部23は、コピー元ブロック（GCソースブロック）内に格納されている有効データを書き込むべきコピー先ブロック（GCデスティネーションブロック）内の位置（コピー先位置）を決定し、有効データをコピー先ブロック（GCデスティネーションブロック）のこの決定された位置（コピー先位置）にコピーする。この場合、GC動作制御部23は、有効データとこの有効データの論理アドレスの双方を、コピー先ブロック（GCデスティネーションブロック）にコピーしてもよい。

40

#### 【0066】

本実施形態では、上述したように、ライト動作制御部21は、ホスト2からのデータ（ライトデータ）とホスト2からの論理アドレスの双方を書き込み先ブロックに書き込むことができる。このため、GC動作制御部23は、コピー元ブロック（GCソースブロック）内の各データの論理アドレスをこのコピー元ブロック（GCソースブロック）から容易に取得することができるので、コピーされた有効データの論理アドレスをホスト2に容易に通知することができる。

#### 【0067】

NANDインタフェース13は、CPU12の制御の下、NAND型フラッシュメモリ5を制御するように構成されたメモリ制御回路である。DRAMインタフェース14は、C

50

P U 1 2 の制御の下、D R A M 6 を制御するように構成された D R A M 制御回路である。D R A M 6 の記憶領域の一部は、内部バッファ（共有キャッシュ）3 1 の格納のために使用される。また、D R A M 6 の記憶領域の他の一部は、ブロック管理テーブル 3 2 の格納のために使用される。なお、これら内部バッファ（共有キャッシュ）3 1、およびブロック管理テーブル 3 2 は、コントローラ 4 内の図示しない S R A M に格納されてもよい。

【 0 0 6 8 】

図 7 は、ホスト 2 が論理アドレスとブロック番号とを指定し且つフラッシュストレージデバイス 3 がブロック内物理アドレス（ブロック内オフセット）を決定するデータ書き込み動作と、ホスト 2 がブロック番号とブロック内物理アドレス（ブロック内オフセット）とを指定するデータ読み出し動作とを示す。

10

【 0 0 6 9 】

データ書き込み動作は以下の手順で実行される。

( 1 ) ホスト 2 のライト処理部 4 1 2 がフラッシュストレージデバイス 3 にデータ（ライトデータ）を書き込むことが必要な時、ライト処理部 4 1 2 は、フリーブロックを割り当てるようにフラッシュストレージデバイス 3 に要求してもよい。フラッシュストレージデバイス 3 のコントローラ 4 は、N A N D 型フラッシュメモリ 5 のフリーブロック群を管理するブロック割り当て部 7 0 1 を含む。ブロック割り当て部 7 0 1 がライト処理部 4 1 2 からこの要求（ブロック割り当て要求）を受信した時、ブロック割り当て部 7 0 1 は、フリーブロック群の一つのフリーブロックをホスト 2 に割り当て、割り当てられたブロックのブロック番号（B L K #）をホスト 2 に通知する。

20

【 0 0 7 0 】

あるいは、ライト処理部 4 1 2 がフリーブロック群を管理する構成においては、ライト処理部 4 1 2 が自身で書き込み先ブロックを選択してもよい。

( 2 ) ライト処理部 4 1 2 は、ライトデータに対応する論理アドレス（例えば L B A）と書き込み先ブロックのブロック番号（B L K #）とを指定するライト要求をフラッシュストレージデバイス 3 に送信する。

【 0 0 7 1 】

( 3 ) フラッシュストレージデバイス 3 のコントローラ 4 は、データ書き込み用のページを割り当てるページ割り当て部 7 0 2 を含む。ページ割り当て部 7 0 2 がライト要求を受信した時、ページ割り当て部 7 0 2 は、ライト要求によって指定されたブロック番号を有するブロック（書き込み先ブロック）内の書き込み先位置を示すブロック内物理アドレス（ブロック内 P B A）を決定する。ブロック内物理アドレス（ブロック内 P B A）は、上述のブロック内オフセット（単にオフセットとしても参照される）によって表すことができる。コントローラ 4 は、ライト要求によって指定されたブロック番号と、ブロック内物理アドレス（ブロック内 P B A）とに基づいて、ホスト 2 からのライトデータを、書き込み先ブロック内の書き込み先位置に書き込む。

30

【 0 0 7 2 】

( 4 ) コントローラ 4 は、書き込み先位置を示すブロック内物理アドレス（ブロック内 P B A）をライト要求に対するレスポンスとしてホスト 2 に通知する。あるいは、コントローラ 4 は、ライトデータに対応する論理アドレス（L B A）と、書き込み先位置を示すブロック内 P B A（オフセット）との組を、ライト要求に対するレスポンスとしてホスト 2 に通知してもよい。換言すれば、コントローラは、ブロック内物理アドレス、または論理アドレスとブロック番号とブロック内物理アドレスとの組のいずれかを、ホスト 2 に通知する。ホスト 2 においては、ライトデータが書き込まれた物理記憶位置を示す物理アドレス（ブロック番号、ブロック内物理アドレス（ブロック内オフセット））が、このライトデータの論理アドレスにマッピングされるように、L U T 4 1 1 が更新される。

40

【 0 0 7 3 】

データリード動作は以下の手順で実行される。

( 1 ) ' ホスト 2 がフラッシュストレージデバイス 3 からデータをリードすることが必要な

50



時、ホスト 2 は、L U T 4 1 1 を参照して、リードすべきデータの論理アドレスに対応する物理アドレス（ブロック番号、ブロック内物理アドレス（ブロック内オフセット））を L U T 4 1 1 から取得する。

【 0 0 7 4 】

（ 2 ）' ホスト 2 は、取得されたブロック番号およびブロック内物理アドレス（ブロック内オフセット）を指定するリード要求をフラッシュストレージデバイス 3 に送出する。フラッシュストレージデバイス 3 のコントローラ 4 がこのリード要求をホスト 2 から受信した時、コントローラ 4 は、ブロック番号およびブロック内物理アドレスに基づいて、リード対象のブロックおよびリード対象の物理記憶位置を特定し、このリード対象のブロック内のリード対象の物理記憶位置からデータをリードする。

10

【 0 0 7 5 】

図 8 は、フラッシュストレージデバイス 3 に適用されるライトコマンドを示す。ライトコマンドは、フラッシュストレージデバイス 3 にデータの書き込みを要求するコマンドである。このライトコマンドは、コマンド I D、ブロック番号 B L K #、論理アドレス、長さ、等を含んでもよい。

【 0 0 7 6 】

コマンド I D はこのコマンドがライトコマンドであることを示す I D（コマンドコード）であり、ライトコマンドにはライトコマンド用のコマンド I D が含まれる。ブロック番号 B L K # は、データが書き込まれるべきブロックを一意に識別可能な識別子（ブロックアドレス）である。

20

【 0 0 7 7 】

論理アドレスは、書き込まれるべきライトデータを識別するための識別子である。この論理アドレスは、上述したように、L B A であってもよいし、キー・バリュー・ストアのキーであってもよいし、キーのハッシュ値であってもよい。論理アドレスが L B A である場合には、このライトコマンドに含まれる論理アドレス（開始 L B A）は、ライトデータが書き込まれるべき論理位置（最初の論理位置）を示す。

【 0 0 7 8 】

長さは、書き込まれるべきライトデータの長さを示す。この長さ（データ長）は、粒度（G r a i n）の数によって指定されてもよいし、L B A の数によって指定されてもよいし、あるいはそのサイズがバイトによって指定されてもよい。

30

ホスト 2 からライトコマンドを受信した時、コントローラ 4 は、ライトコマンドによって指定されたブロック番号を有するブロック内の書き込み先位置を決定する。この書き込み先位置は、ページ書き込み順序の制約およびパッドページ等を考慮して決定される。そして、コントローラ 4 は、ホスト 2 からのデータを、ライトコマンドによって指定されたブロック番号を有するこのブロック内のこの書き込み先位置に書き込む。

【 0 0 7 9 】

図 9 は、図 8 のライトコマンドに対するレスポンスを示す。このレスポンスは、ブロック内物理アドレス、長さを含む。ブロック内物理アドレスは、データが書き込まれたブロック内の位置（物理記憶位置）を示す。ブロック内物理アドレスは、上述したように、ブロック内オフセットによって指定可能である。長さは、書き込まれたデータの長さを示す。この長さ（データ長）は、粒度（G r a i n）の数によって指定されてもよいし、L B A の数によって指定されてもよいし、あるいはそのサイズがバイトによって指定されてもよい。

40

【 0 0 8 0 】

あるいは、このレスポンスは、ブロック内物理アドレスおよび長さだけでなく、論理アドレスおよびブロック番号をさらに含んでもよい。論理アドレスは、図 8 のライトコマンドに含まれていた論理アドレスである。ブロック番号は、図 8 のライトコマンドに含まれていた論理アドレスである。

【 0 0 8 1 】

図 1 0 は、フラッシュストレージデバイス 3 に適用される T r i m コマンドを示す。

50

このTrimコマンドは、無効にすべきデータが格納されている物理記憶位置を示すブロック番号およびブロック内物理アドレス（ブロック内オフセット）を含むコマンドである。つまり、このTrimコマンドは、LBAのような論理アドレスではなく、物理アドレスを指定可能である。このTrimコマンドは、コマンドID、物理アドレス、長さを含む。

#### 【0082】

コマンドIDはこのコマンドがTrimコマンドであることを示すID（コマンドコード）であり、TrimコマンドにはTrimコマンド用のコマンドIDが含まれる。物理アドレスは、無効化すべきデータが格納されている最初の物理記憶位置を示す。本実施形態では、この物理アドレスは、ブロック番号とオフセット（ブロック内オフセット）との組み合わせによって指定される。

10

#### 【0083】

長さは、無効化すべきデータの長さを示す。この長さ（データ長）は、粒度（Grain）の数によって指定されてもよいし、バイトによって指定されてもよい。コントローラ4は、複数のブロックの各々に含まれるデータそれぞれの有効/無効を示すフラグ（ビットマップフラグ）をブロック管理テーブル32を使用して管理する。無効にすべきデータが格納されている物理記憶位置を示すブロック番号およびオフセット（ブロック内オフセット）を含むTrimコマンドをホスト2から受信した場合、コントローラ4は、ブロック管理テーブル32を更新して、Trimコマンドに含まれるブロック番号およびブロック内オフセットに対応する物理記憶位置のデータに対応するフラグ（ビットマップフラグ）を無効を示す値に変更する。

20

#### 【0084】

図11は、ブロック内物理アドレスを規定するブロック内オフセットを示す。ブロック番号はある一つのブロックBLKを指定する。各ブロックBLKは、図11に示されているように、複数のページ（ここでは、ページ0～ページn）を含む。ページサイズ（各ページのユーザデータ格納領域）が16Kバイトであり、粒度（Grain）が4KBのサイズであるケースにおいては、このブロックBLKは、 $4 \times (n + 1)$ 個の領域に論理的に分割される。

#### 【0085】

オフセット+0はページ0の最初の4KB領域を示し、オフセット+1はページ0の2番目の4KB領域を示し、オフセット+2はページ0の3番目の4KB領域を示し、オフセット+3はページ0の4番目の4KB領域を示す。オフセット+4はページ1の最初の4KB領域を示し、オフセット+5はページ1の2番目の4KB領域を示し、オフセット+6はページ1の3番目の4KB領域を示し、オフセット+7はページ1の4番目の4KB領域を示す。

30

#### 【0086】

図12は、ライトコマンドに応じて実行される書き込み動作を示す。いま、ブロックBLK#1が書き込み先ブロックとして割り当てられている場合を想定する。コントローラ4は、ページ0、ページ1、ページ2、...ページnという順序で、データをページ単位でブロックBLK#1に書き込む。

40

#### 【0087】

図12においては、ブロックBLK#1のページ0に16Kバイト分のデータがすでに書き込まれている状態で、ブロック番号（=BLK#1）、論理アドレス（LBAx）および長さ（=4）を指定するライトコマンドがホスト2から受信された場合が想定されている。コントローラ4は、ブロックBLK#1のページ1を書き込み先位置として決定し、ホスト2から受信される16Kバイト分のライトデータをブロックBLK#1のページ1に書き込む。そして、コントローラ4は、このライトコマンドに対するレスポンスとして、オフセット（ブロック内オフセット）、長さをホスト2に返す。このケースにおいては、オフセット（ブロック内オフセット）は+5であり、長さは4である。あるいは、コントローラ4は、このライトコマンドに対するレスポンスとして、論理アドレス、ブロック

50

番号、オフセット（ブロック内オフセット）、長さをホスト 2 に返してもよい。このケースにおいては、論理アドレスは  $LBA \times$  であり、ブロック番号は  $BLK \# 1$  であり、オフセット（ブロック内オフセット）は  $+ 5$  であり、長さは 4 である。

【 0 0 8 8 】

図 1 3 は、不良ページ（パッドページ）をスキップする書き込み動作を示す。

図 1 3 においては、ブロック  $BLK \# 1$  のページ 0、ページ 1 にデータがすでに書き込まれている状態で、ブロック番号（ $= BLK \# 1$ ）、論理アドレス（ $LBA \times + 1$ ）および長さ（ $= 4$ ）を指定するライトコマンドがホスト 2 から受信された場合が想定されている。もしブロック  $BLK \# 1$  のページ 2 が不良ページであるならば、コントローラ 4 は、ブロック  $BLK \# 1$  のページ 3 を書き込み先位置として決定し、ホスト 2 から受信される 16 K バイト分のライトデータをブロック  $BLK \# 1$  のページ 3 に書き込む。そして、コントローラ 4 は、このライトコマンドに対するレスポンスとして、オフセット（ブロック内オフセット）、長さをホスト 2 に返す。このケースにおいては、オフセット（ブロック内オフセット）は  $+ 1 2$  であり、長さは 4 である。あるいは、コントローラ 4 は、このライトコマンドに対するレスポンスとして、論理アドレス、ブロック番号、オフセット（ブロック内オフセット）、長さをホスト 2 に返してもよい。このケースにおいては、論理アドレスは  $LBA \times + 1$  であり、ブロック番号は  $BLK \# 1$  であり、オフセット（ブロック内オフセット）は  $+ 1 2$  であり、長さは 4 である。

10

【 0 0 8 9 】

図 1 4 は、不良ページをスキップする書き込み動作の別の例を示す。

図 1 4 においては、不良ページを挟む 2 つのページに跨がってデータが書き込まれる場合が想定されている。いま、ブロック  $BLK \# 2$  のページ 0、ページ 1 にデータがすでに書き込まれており、且つ内部バッファ（共有キャッシュ）3 1 に未書き込みの 8 K バイト分のライトデータが残っている場合を想定する。この状態で、ブロック番号（ $= BLK \# 2$ ）、論理アドレス（ $LBA y$ ）および長さ（ $= 6$ ）を指定するライトコマンドが受信されたならば、コントローラ 4 は、未書き込みの 8 K バイトライトデータと、ホスト 2 から新たに受信される 24 K バイトライトデータ内の最初の 8 K バイトライトデータとを使用して、ページサイズに対応する 16 K バイトライトデータを準備する。そして、コントローラ 4 は、この準備した 16 K バイトライトデータをブロック  $BLK \# 2$  のページ 2 に書き込む。

20

30

【 0 0 9 0 】

もしブロック  $BLK \# 2$  の次のページ 3 が不良ページであるならば、コントローラ 4 は、ブロック  $BLK \# 2$  のページ 4 を次の書き込み先位置として決定し、ホスト 2 から受信された 24 K バイトライトデータ内の残りの 16 K バイトライトデータを、ブロック  $BLK \# 2$  のページ 4 に書き込む。

【 0 0 9 1 】

そして、コントローラ 4 は、このライトコマンドに対するレスポンスとして、2 つのオフセット（ブロック内オフセット）と、2 つの長さをホスト 2 に返す。このケースにおいては、このレスポンスは、オフセット（ $= + 1 0$ ）、長さ（ $= 2$ ）、オフセット（ $= + 1 6$ ）、長さ（ $= 4$ ）を含んでもよい。あるいは、コントローラ 4 は、このライトコマンドに対するレスポンスとして、 $LBA y$ 、ブロック番号（ $= BLK \# 2$ ）、オフセット（ $= + 1 0$ ）、長さ（ $= 2$ ）、ブロック番号（ $= BLK \# 2$ ）、オフセット（ $= + 1 6$ ）、長さ（ $= 4$ ）をホスト 2 に返してもよい。

40

【 0 0 9 2 】

図 1 5、図 1 6 は、論理アドレスとデータのペアをブロック内のページに書き込む動作を示す。

各ブロックにおいて、各ページは、ユーザデータを格納するためのユーザデータ領域と管理データを格納するための冗長領域とを含んでもよい。ページサイズは  $16 K B + \text{アルファ}$  である。

【 0 0 9 3 】

50

コントローラ 4 は、4 K B ユーザデータとこの 4 K B ユーザデータに対応する論理アドレス（例えば L B A）との双方を書き込み先ブロック B L K に書き込む。この場合、図 1 5 に示すように、各々が L B A と 4 K B ユーザデータとを含む 4 つのデータセットが同じページに書き込まれてもよい。ブロック内オフセットは、セット境界を示してもよい。

【 0 0 9 4 】

あるいは、図 1 6 に示されているように、4 つの 4 K B ユーザデータがページ内のユーザデータ領域に書き込まれ、これら 4 つの 4 K B ユーザデータに対応する 4 つの L B A がこのページ内の冗長領域に書き込まれてもよい。

図 1 7 は、スーパーブロックが使用されるケースにおけるブロック番号とオフセット（ブロック内オフセット）との関係を示す。以下では、ブロック内オフセットは単にオフセットとしても参照される。

10

【 0 0 9 5 】

ここでは、図示を簡単化するために、ある一つのスーパーブロック S B # 1 が 4 つのブロック B L K # 1 1、B L K # 2 1、B L K # 3 1、B L K # 4 1 から構成されている場合が想定されている。コントローラ 4 は、ブロック B L K # 1 1 のページ 0、ブロック B L K # 2 1 のページ 0、ブロック B L K # 3 1 のページ 0、ブロック B L K # 4 1 のページ 0、ブロック B L K # 1 1 のページ 1、ブロック B L K # 2 1 のページ 1、ブロック B L K # 3 1 のページ 1、ブロック B L K # 4 1 のページ 1、... という順序でデータを書き込む。

【 0 0 9 6 】

オフセット + 0 はブロック B L K # 1 1 のページ 0 の最初の 4 K B 領域を示し、オフセット + 1 はブロック B L K # 1 1 のページ 0 の 2 番目の 4 K B 領域を示し、オフセット + 2 はブロック B L K # 1 1 のページ 0 の 3 番目の 4 K B 領域を示し、オフセット + 3 はブロック B L K # 1 1 のページ 0 の 4 番目の 4 K B 領域を示す。

20

【 0 0 9 7 】

オフセット + 4 はブロック B L K # 2 1 のページ 0 の最初の 4 K B 領域を示し、オフセット + 5 はブロック B L K # 2 1 のページ 0 の 2 番目の 4 K B 領域を示し、オフセット + 6 はブロック B L K # 2 1 のページ 0 の 3 番目の 4 K B 領域を示し、オフセット + 7 はブロック B L K # 2 1 のページ 0 の 4 番目の 4 K B 領域を示す。

【 0 0 9 8 】

同様に、オフセット + 1 2 はブロック B L K # 4 1 のページ 0 の最初の 4 K B 領域を示し、オフセット + 1 3 はブロック B L K # 4 1 のページ 0 の 2 番目の 4 K B 領域を示し、オフセット + 1 4 はブロック B L K # 4 1 のページ 0 の 3 番目の 4 K B 領域を示し、オフセット + 1 5 はブロック B L K # 4 1 のページ 0 の 4 番目の 4 K B 領域を示す。

30

【 0 0 9 9 】

オフセット + 1 6 はブロック B L K # 1 1 のページ 1 の最初の 4 K B 領域を示し、オフセット + 1 7 はブロック B L K # 1 1 のページ 1 の 2 番目の 4 K B 領域を示し、オフセット + 1 8 はブロック B L K # 1 1 のページ 1 の 3 番目の 4 K B 領域を示し、オフセット + 1 9 はブロック B L K # 1 1 のページ 1 の 4 番目の 4 K B 領域を示す。

【 0 1 0 0 】

オフセット + 2 0 はブロック B L K # 2 1 のページ 1 の最初の 4 K B 領域を示し、オフセット + 2 1 はブロック B L K # 2 1 のページ 1 の 2 番目の 4 K B 領域を示し、オフセット + 2 2 はブロック B L K # 2 1 のページ 1 の 3 番目の 4 K B 領域を示し、オフセット + 2 3 はブロック B L K # 2 1 のページ 1 の 4 番目の 4 K B 領域を示す。

40

【 0 1 0 1 】

同様に、オフセット + 2 8 はブロック B L K # 4 1 のページ 1 の最初の 4 K B 領域を示し、オフセット + 2 9 はブロック B L K # 4 1 のページ 1 の 2 番目の 4 K B 領域を示し、オフセット + 3 0 はブロック B L K # 4 1 のページ 1 の 3 番目の 4 K B 領域を示し、オフセット + 3 1 はブロック B L K # 4 1 のページ 1 の 4 番目の 4 K B 領域を示す。

【 0 1 0 2 】

50

図 18 は、フラッシュストレージデバイス 3 に適用される最大ブロック番号ゲットコマンドを示す。

最大ブロック番号ゲットコマンドは、フラッシュストレージデバイス 3 から最大ブロック番号を取得するためのコマンドである。ホスト 2 は、フラッシュストレージデバイス 3 に最大ブロック番号ゲットコマンドに送信することにより、フラッシュストレージデバイス 3 に含まれるブロックの数を示す最大ブロック番号を認識することができる。最大ブロック番号ゲットコマンドは、最大ブロック番号ゲットコマンド用のコマンド ID を含み、パラメータは含まない。

【 0 1 0 3 】

図 19 は、最大ブロック番号ゲットコマンドに対するレスポンスを示す。

10

最大ブロック番号ゲットコマンドをホスト 2 から受信した時、フラッシュストレージデバイス 3 は、図 19 に示すレスポンスをホスト 2 に返す。このレスポンスは、最大ブロック番号（つまり、フラッシュストレージデバイス 3 に含まれる利用可能なブロックの総数）を示すパラメータを含む。

【 0 1 0 4 】

図 20 は、フラッシュストレージデバイス 3 に適用されるブロックサイズゲットコマンドを示す。

ブロックサイズゲットコマンドは、フラッシュストレージデバイス 3 からブロックサイズを取得するためのコマンドである。ホスト 2 は、フラッシュストレージデバイス 3 にブロックサイズゲットコマンドに送信することにより、フラッシュストレージデバイス 3 に含まれる NAND 型フラッシュメモリ 5 のブロックサイズを認識することができる。

20

【 0 1 0 5 】

なお、別の実施形態では、ブロックサイズゲットコマンドは、ブロック番号を指定するパラメータを含んでいてもよい。あるブロック番号を指定するブロックサイズゲットコマンドをホスト 2 から受信した場合、フラッシュストレージデバイス 3 は、このブロック番号を有するブロックのブロックサイズをホスト 2 に返す。これにより、たとえ NAND 型フラッシュメモリ 5 に含まれるブロックそれぞれのブロックサイズが不均一である場合であっても、ホスト 2 は、個々のブロックそれぞれのブロックサイズを認識することができる。

【 0 1 0 6 】

図 21 は、ブロックサイズゲットコマンドに対するレスポンスを示す。

30

ブロックサイズゲットコマンドをホスト 2 から受信した時、フラッシュストレージデバイス 3 は、ブロックサイズ（NAND 型フラッシュメモリ 5 に含まれるブロックそれぞれの共通のブロックサイズ）をホスト 2 に返す。この場合、もしブロック番号がブロックサイズゲットコマンドによって指定されていたならば、フラッシュストレージデバイス 3 は、上述したように、このブロック番号を有するブロックのブロックサイズをホスト 2 に返す。

【 0 1 0 7 】

図 22 は、フラッシュストレージデバイス 3 に適用されるブロックアロケートコマンドを示す。

ブロックアロケートコマンドは、フラッシュストレージデバイス 3 にブロック（フリーブロック）の割り当てを要求するコマンド（ブロック割り当て要求）である。ホスト 2 は、ブロックアロケートコマンドをフラッシュストレージデバイス 3 に送信することによって、フリーブロックを割り当てるようにフラッシュストレージデバイス 3 に要求し、これによってブロック番号（割り当てられたフリーブロックのブロック番号）を取得することができる。

40

【 0 1 0 8 】

フラッシュストレージデバイス 3 がフリーブロック群をフリーブロックリストによって管理し、ホスト 2 はフリーブロック群を管理しないケースにおいては、ホスト 2 は、フリーブロックを割り当てるようにフラッシュストレージデバイス 3 に要求し、これによってブロック番号を取得する。一方、ホスト 2 がフリーブロック群を管理するケースにおいては、ホスト 2 は、フリーブロック群の一つを自身で選択することができるので、ブロックア

50

ロケートコマンドをフラッシュストレージデバイス 3 に送信する必要は無い。

【 0 1 0 9 】

図 2 3 は、ブロックアロケートコマンドに対するレスポンスを示す。

ブロックアロケートコマンドをホスト 2 から受信した時、フラッシュストレージデバイス 3 は、フリーブロックリストから、ホスト 2 に割り当てるべきフリーブロックを選択し、選択したフリーブロックのブロック番号を含むレスポンスをホスト 2 に返す。

【 0 1 1 0 】

図 2 4 は、ホスト 2 とフラッシュストレージデバイス 3 とによって実行されるブロック情報取得処理を示す。

ホスト 2 がフラッシュストレージデバイス 3 の使用を開始する時、ホスト 2 は、まず、最大ブロック番号ゲットコマンドをフラッシュストレージデバイス 3 に送信する。フラッシュストレージデバイス 3 のコントローラは、最大ブロック番号をホスト 2 に返す。最大ブロック番号は、利用可能なブロックの総数を示す。なお、上述のスーパーブロックが使用されるケースにおいては、最大ブロック番号は、利用可能なスーパーブロックの総数を示してもよい。

10

【 0 1 1 1 】

次いで、ホスト 2 は、ブロックサイズゲットコマンドをフラッシュストレージデバイス 3 に送信して、ブロックサイズを取得する。この場合、ホスト 2 は、ブロック番号 1 を指定するブロックサイズゲットコマンド、ブロック番号 2 を指定するブロックサイズゲットコマンド、ブロック番号 3 を指定するブロックサイズゲットコマンド、... をフラッシュストレージデバイス 3 にそれぞれ送信して、全てのブロックそれぞれのブロックサイズを個別に取得してもよい。

20

【 0 1 1 2 】

このブロック情報取得処理により、ホスト 2 は、利用可能ブロック数、個々のブロックのブロックサイズを認識することができる。

図 2 5 は、ホスト 2 とフラッシュストレージデバイス 3 とによって実行される書き込み処理のシーケンスを示す。

ホスト 2 は、まず、書き込みのために使用すべきブロック（フリーブロック）を自身で選択するか、またはブロックアロケートコマンドをフラッシュストレージデバイス 3 に送信することによってフリーブロックを割り当てるようにフラッシュストレージデバイス 3 に要求する。そして、ホスト 2 は、自身で選択したブロックのブロック番号 B L K #（またはフラッシュストレージデバイス 3 によって割り当てられたフリーブロックのブロック番号 B L K #）と、論理アドレス（L B A）と、長さを含むライトコマンドをフラッシュストレージデバイス 3 に送信する（ステップ S 2 0）。

30

【 0 1 1 3 】

フラッシュストレージデバイス 3 のコントローラ 4 がこのライトコマンドを受信した時、コントローラ 4 は、ホスト 2 からのライトデータを書き込むべき、このブロック番号 B L K # を有するブロック（書き込み先ブロック B L K #）内の書き込み先位置を決定し、この書き込み先ブロック B L K # の書き込み先位置にライトデータを書き込む（ステップ S 1 1）。ステップ S 1 1 では、コントローラ 4 は、論理アドレス（ここでは L B A）とライトデータの双方を書き込み先ブロックに書き込んでよい。

40

【 0 1 1 4 】

コントローラ 4 は、書き込み先ブロック B L K # に対応するブロック管理テーブル 3 2 を更新して、書き込まれたデータに対応するビットマップフラグ（つまり、このデータが書き込まれたオフセット（ブロック内オフセット）に対応するビットマップフラグ）を 0 から 1 に変更する（ステップ S 1 2）。

【 0 1 1 5 】

例えば、図 2 6 に示されているように、開始 L B A が L B A x である 1 6 K バイト更新データがブロック B L K # 1 のオフセット + 4 ~ + 7 に対応する物理記憶位置に書き込まれた場合を想定する。この場合、図 2 7 に示されているように、ブロック B L K # 1 用のブ

50

ロック管理テーブルにおいては、オフセット + 4 ~ + 7 に対応するビットマップフラグそれぞれが 0 から 1 に変更される。

【 0 1 1 6 】

そして、図 2 5 に示すように、コントローラ 4 は、このライトコマンドに対するレスポンスをホスト 2 に返す (ステップ S 1 3 )。このレスポンスは、このデータが書き込まれたオフセット (ブロック内オフセット) を少なくとも含む。

ホスト 2 がこのレスポンスを受信した時、ホスト 2 は、ホスト 2 によって管理されている L U T 4 1 1 を更新して、書き込まれたライトデータに対応する論理アドレスそれぞれに物理アドレスをマッピングする。図 2 8 に示されているように、L U T 4 1 1 は、複数の論理アドレス (例えば L B A ) それぞれに対応する複数のエントリを含む。ある論理アドレス (例えばある L B A ) に対応するエントリには、この L B A に対応するデータが格納されている N A N D 型フラッシュメモリ 5 内の位置 (物理記憶位置) を示す物理アドレス P B A、つまりブロック番号、オフセット (ブロック内オフセット) が格納される。図 2 6 に示されているように、開始 L B A が L B A x である 1 6 K バイト更新データがブロック B L K # 1 のオフセット + 4 ~ + 7 に対応する物理記憶位置に書き込まれたならば、図 2 8 に示されているように、L U T 4 1 1 が更新されて、L B A x に対応するエントリに B L K # 1、オフセット + 4 が格納され、L B A x + 1 に対応するエントリに B L K # 1、オフセット + 5 が格納され、L B A x + 2 に対応するエントリに B L K # 1、オフセット + 6 が格納され、L B A x + 3 に対応するエントリに B L K # 1、オフセット + 7 が格納される。

【 0 1 1 7 】

図 2 5 に示すように、この後、ホスト 2 は、上述の更新データの書き込みによって不要になった以前のデータを無効化するための T r i m コマンドをフラッシュストレージデバイス 3 に送信する。図 2 6 に示されているように、以前のデータがブロック B L K # 0 のオフセット + 0、オフセット + 1、オフセット + 2、オフセット + 3 に対応する位置に格納されている場合には、図 2 9 に示すように、ブロック番号 (= B L K # 0)、オフセット (= + 0)、長さ (= 4) を指定する T r i m コマンドがホスト 2 からフラッシュストレージデバイス 3 に送信される。フラッシュストレージデバイス 3 のコントローラ 4 は、この T r i m コマンドに応じて、ブロック管理テーブル 3 2 を更新する (図 2 5、ステップ S 1 4)。ステップ S 1 5 においては、図 2 9 に示すように、ブロック B L K # 0 用のブロック管理テーブルにおいて、オフセット + 0 ~ + 3 に対応するビットマップフラグそれぞれが 1 から 0 に変更される。

【 0 1 1 8 】

図 3 0 は、フラッシュストレージデバイス 3 に適用されるリードコマンドを示す。リードコマンドは、フラッシュストレージデバイス 3 にデータの読み出しを要求するコマンドである。このリードコマンドは、コマンド I D、物理アドレス P B A、長さ、転送先ポインタを含む。

【 0 1 1 9 】

コマンド I D はこのコマンドがリードコマンドであることを示す I D (コマンドコード) であり、リードコマンドにはリードコマンド用のコマンド I D が含まれる。物理アドレス P B A は、データが読み出されるべき最初の物理記憶位置を示す。物理アドレス P B A は、ブロック番号、オフセット (ブロック内オフセット) によって指定される。

【 0 1 2 0 】

長さは、リードすべきデータの長さを示す。このデータ長は、G r a i n の数によって指定可能である。

転送先ポインタは、読み出されたデータが転送されるべきホスト 2 内のメモリ上の位置を示す。

一つのリードコマンドは、物理アドレス P B A (ブロック番号、オフセット) と長さの組を複数指定することができる。

【 0 1 2 1 】

図 3 1 は、リード動作を示す。

ここでは、ブロック番号 (= B L K # 2 )、オフセット (= + 5 )、長さ (= 3 ) を指定するリードコマンドがホスト 2 から受信された場合が想定されている。フラッシュストレージデバイス 3 のコントローラ 4 は、ブロック番号 (= B L K # 2 )、オフセット (= + 5 )、長さ (= 3 ) に基づいて、B L K # 2 からデータ d 1 ~ d 3 をリードする。この場合、コントローラ 4 は、B L K # 2 のページ 1 から 1 ページサイズ分のデータをリードし、このリードデータからデータ d 1 ~ データ d 3 を抽出する。次いで、コントローラ 4 は、データ d 1 ~ データ d 3 を、転送先ポインタによって指定されるホストメモリ上に転送する。

【 0 1 2 2 】

図 3 2 は、ホスト 2 からのリードコマンドに応じて、異なる物理記憶位置にそれぞれ格納されているデータ部をリードする動作を示す。

ここでは、ブロック番号 (= B L K # 2 )、オフセット (= + 1 0 )、長さ (= 2 )、ブロック番号 (= B L K # 2 )、オフセット (= + 1 6 )、長さ (= 4 ) を指定するリードコマンドがホスト 2 から受信された場合が想定されている。フラッシュストレージデバイス 3 のコントローラ 4 は、ブロック番号 (= B L K # 2 )、オフセット (= + 1 0 )、長さ (= 2 ) に基づいて、B L K # 2 のページ 2 から 1 ページサイズ分のデータをリードし、このリードデータからデータ d 1 ~ データ d 2 を抽出する。次いで、コントローラ 4 は、ブロック番号 (= B L K # 2 )、オフセット (= + 1 6 )、長さ (= 4 ) に基づいて、B L K # 2 のページ 4 から 1 ページサイズ分のデータ (データ d 3 ~ データ d 6 ) をリードする。そして、コントローラ 4 は、データ d 1 ~ データ d 2 とデータ d 3 ~ データ d 6 とを結合することによって得られる長さ (= 6 ) のリードデータを、リードコマンド内の転送先ポインタによって指定されるホストメモリ上に転送する。

【 0 1 2 3 】

これにより、たとえブロック内に不良ページが存在する場合であっても、リードエラーを引き起こすことなく、別個の物理記憶位置からデータ部をリードすることができる。また、たとえデータが 2 つのブロックに跨がって書き込まれている場合であっても、このデータを一つのリードコマンドの発行によってリードすることができる。

【 0 1 2 4 】

図 3 3 は、ホスト 2 とフラッシュストレージデバイス 3 とによって実行されるリード処理のシーケンスを示す。

ホスト 2 は、ホスト 2 によって管理されている L U T 4 1 1 を参照して、ユーザアプリケーションからのリード要求に含まれる論理アドレスをブロック番号、オフセットに変換する。そして、ホスト 2 は、このブロック番号、オフセット、長さを指定するリードコマンドをフラッシュストレージデバイス 3 に送信する。

【 0 1 2 5 】

フラッシュストレージデバイス 3 のコントローラ 4 がリードコマンドをホスト 2 から受信した時、コントローラ 4 は、このリードコマンドによって指定されたブロック番号に対応するブロックをリード対象のブロックとして決定するとともに、このリードコマンドによって指定されたオフセットに基づいてリード対象のページを決定する (ステップ S 3 1 )。ステップ S 3 1 では、コントローラ 4 は、まず、リードコマンドによって指定されたオフセットを、ページサイズを表す粒度の数 (ここでは、4 ) で除算してもよい。そして、コントローラ 4 は、この除算によって得られる商および余りを、リード対象のページ番号およびリード対象のページ内オフセット位置としてそれぞれ決定してもよい。

【 0 1 2 6 】

コントローラ 4 は、ブロック番号、オフセット、長さによって規定されるデータを N A N D 型フラッシュメモリ 5 からリードし (ステップ S 3 2 )、このリードデータをホスト 2 に送信する。

図 3 4 は、フラッシュストレージデバイス 3 に適用される G C 制御コマンドを示す。

【 0 1 2 7 】

10

20

30

40

50



GC制御コマンドは、GCソースブロック番号およびGCデスティネーションブロック番号をフラッシュストレージデバイス3に通知するために使用される。ホスト2は、各ブロックの有効データ量/無効データ量を管理しており、有効データ量がより少ない幾つかのブロックをGCソースブロックとして選択することができる。また、ホスト2は、フリーブロックリストを管理しており、幾つかのフリーブロックをGCデスティネーションブロックとして選択することができる。このGC制御コマンドは、コマンドID、GCソースブロック番号、GCデスティネーションブロック番号、等を含んでもよい。

【0128】

コマンドIDはこのコマンドがGC制御コマンドであることを示すID（コマンドコード）であり、GC制御コマンドにはGC制御コマンド用のコマンドIDが含まれる。

10

GCソースブロック番号は、GCソースブロックを示すブロック番号である。ホスト2は、どのブロックをGCソースブロックとすべきかを指定することができる。ホスト2は、複数のGCソースブロック番号を一つのGC制御コマンドに設定してもよい。

【0129】

GCデスティネーションブロック番号は、GCデスティネーションブロックを示すブロック番号である。ホスト2は、どのブロックをGCデスティネーションブロックとすべきかを指定することができる。ホスト2は、複数のGCデスティネーションブロック番号を一つのGC制御コマンドに設定してもよい。

【0130】

図35は、GC用コールバックコマンドを示す。

20

GC用コールバックコマンドは、GCによってコピーされた有効データの論理アドレスとこの有効データのコピー先位置を示すブロック番号およびオフセットとをホスト2に通知するために使用される。

GC用コールバックコマンドは、コマンドID、論理アドレス、長さ、デスティネーション物理アドレスを含んでよい。

【0131】

コマンドIDはこのコマンドがGC用コールバックコマンドであることを示すID（コマンドコード）であり、GC用コールバックコマンドにはGC用コールバックコマンド用のコマンドIDが含まれる。

論理アドレスは、GCによってGCソースブロックからGCデスティネーションブロックにコピーされた有効データの論理アドレスを示す。

30

【0132】

長さは、このコピーされたデータの長さを示す。このデータ長は、粒度（Grain）の数によって指定されてもよい。

デスティネーション物理アドレスは、有効データがコピーされたGCデスティネーションブロック内の位置を示す。デスティネーション物理アドレスは、ブロック番号、オフセット（ブロック内オフセット）によって指定される。

【0133】

図36は、ガベージコレクション（GC）動作の手順を示す。

例えば、ホスト2は、ホスト2によって管理されているフリーブロックリストに含まれている残りフリーブロックの数が閾値以下に低下した場合、GCソースブロックおよびGCデスティネーションブロックを選択し、選択されたGCソースブロックおよび選択されたGCデスティネーションブロックを指定するGC制御コマンドをフラッシュストレージデバイス3に送信する（ステップS41）。あるいは、ライト処理部412がフリーブロック群を管理する構成においては、残りフリーブロックの数が閾値以下に低下した際にライト処理部412がホスト2にその旨通知を行ない、通知を受信したホスト2がブロック選択およびGC制御コマンドの送信を行なってもよい。

40

【0134】

このGC制御コマンドを受信すると、フラッシュストレージデバイス3のコントローラ4は、GCソースブロック内の有効データを書き込むべきGCデスティネーションブロック

50

内の位置（コピー先位置）を決定する動作と、GCソースブロック内の有効データをGCデスティネーションブロック内のコピー先位置にコピーする動作とを含むデータコピー動作を実行する（ステップS51）。ステップS51では、コントローラ4は、GCソースブロック（コピー元ブロック）内の有効データのみならず、この有効データとこの有効データに対応する論理アドレスの双方を、GCソースブロック（コピー元ブロック）からGCデスティネーションブロック（コピー先ブロック）にコピーする。これにより、GCデスティネーションブロック（コピー先ブロック）内にデータと論理アドレスとのペアを保持することができる。

【0135】

また、ステップS51では、GCソースブロック内の全ての有効データのコピーが完了するまでデータコピー動作が繰り返し実行される。複数のGCソースブロックがGC制御コマンドによって指定された場合には、全てのGCソースブロック内の全ての有効データのコピーが完了するまでデータコピー動作が繰り返し実行される。

10

【0136】

そして、コントローラ4は、コピーされた有効データ毎に、その有効データの論理アドレス（LBA）と、その有効データのコピー先位置を示すデスティネーション物理アドレス等を、GC用コールバックコマンドを使用してホスト2に通知する（ステップS52）。ある有効データに対応するデスティネーション物理アドレスは、この有効データがコピーされたコピー先ブロック（GCデスティネーションブロック）のブロック番号と、この有効データがコピーされたコピー先ブロック内の物理記憶位置を示すブロック内物理アドレス（ブロック内オフセット）とによって表される。

20

【0137】

ホスト2がこのGC用コールバックコマンドを受信した時、ホスト2は、ホスト2によって管理されているLUT411を更新して、コピーされた各有効データに対応する論理アドレスにデスティネーション物理アドレス（ブロック番号、ブロック内オフセット）をマッピングする（ステップS42）。

【0138】

図37は、GCのために実行されるデータコピー動作の例を示す。

図37では、GCソースブロック（ここではブロックBLK#50）のオフセット+4に対応する位置に格納されている有効データ（LBA=10）が、GCデスティネーションブロック（ここではブロックBLK#100）のオフセット+0に対応する位置にコピーされ、GCソースブロック（ここではブロックBLK#50）のオフセット+10に対応する位置に格納されている有効データ（LBA=20）が、GCデスティネーションブロック（ここではブロックBLK#100）のオフセット+1に対応する位置にコピーされた場合が想定されている。この場合、コントローラ4は、{LBA10、BLK#100、オフセット(=+0)、LBA20、BLK#100、オフセット(=+1)}をホストに通知する（GC用コールバック処理）。

30

【0139】

図38は、図37のデータコピー動作の結果に基づいて更新されるホスト2のLUT411の内容を示す。

40

このLUT411においては、LBA10に対応するブロック番号およびオフセットは、BLK#50、オフセット(=+4)から、BLK#100、オフセット(=+0)に更新される。同様に、LBA20に対応するブロック番号およびオフセットは、BLK#50、オフセット(=+10)から、BLK#100、オフセット(=+1)に更新される。

【0140】

LUT411が更新された後、ホスト2は、BLK#50およびオフセット(=+4)を指定するTrimコマンドをフラッシュストレージデバイス3に送信して、BLK#50のオフセット(=+4)に対応する位置に格納されているデータを無効化してもよい。さらに、ホスト2は、BLK#50およびオフセット(=+10)を指定するTrimコマンドをフラッシュストレージデバイス3に送信して、BLK#50のオフセット(=+1

50

0) に対応する位置に格納されているデータを無効化してもよい。

あるいは、ホスト2からTrimコマンドを送信せず、GC処理の一環としてコントローラ4がブロック管理テーブル32を更新してこれらのデータを無効化してもよい。

【0141】

以上説明したように、本実施形態によれば、第1の論理アドレスと第1のブロック番号とを指定するライト要求をホスト2から受信した場合、フラッシュストレージデバイス3のコントローラ4は、ホスト2からのデータを書き込むべき、第1のブロック番号を有するブロック(書き込み先ブロック)内の位置(書き込み先位置)を決定し、ホスト2からのデータを書き込み先ブロックの書き込み先位置に書き込み、第1の位置を示す第1のブロック内物理アドレス、または第1の論理アドレスと第1のブロック番号と第1のブロック内物理アドレスとの組のいずれかを、ホスト2に通知する。

10

【0142】

したがって、ホスト2がブロック番号をハンドリングし、フラッシュストレージデバイス3がページ書き込み順序制約/パッドページ等を考慮して、ホスト2によって指定されるブロック番号を有するブロック内の書き込み先位置(ブロック内オフセット)を決定するという構成を実現できる。ホスト2がブロック番号をハンドリングすることにより、上位階層(ホスト2)のアプリケーションレベルアドレス変換テーブルと従来型SSDのLUTレベルアドレス変換テーブルとのマージを実現できる。また、フラッシュストレージデバイス3は、NAND型フラッシュメモリ5の特徴/制約を考慮してNAND型フラッシュメモリ5を制御することができる。さらに、ホスト2はブロック境界を認識することができるので、ブロック境界/ブロックサイズを考慮してユーザデータを各ブロックに書き込むことができる。これにより、ホスト2が同一ブロック内のデータをデータ更新等によって一斉に無効化する等の制御を行うことが可能となるので、GCが実行される頻度を下げることが可能となる。この結果、ライトアンプリフィケーションが低下され、フラッシュストレージデバイス3の性能の向上、フラッシュストレージデバイス3の寿命の最大化を実現できる。

20

【0143】

したがって、ホスト2とフラッシュストレージデバイス3との間の適切な役割分担を実現でき、これによってホスト2とフラッシュストレージデバイス3とを含むシステム全体のI/O性能の向上を図ることができる。

30

また、ガベージコレクションのためのコピー元ブロック番号およびコピー先ブロック番号を指定する制御コマンドをホスト2から受信した場合、フラッシュストレージデバイス3のコントローラ4は、複数のブロックから、コピー元ブロック番号を有する第2のブロックとコピー先ブロック番号を有する第3のブロックとを選択し、第2のブロックに格納されている有効データを書き込むべき第3のブロック内のコピー先位置を決定し、有効データを第3のブロックのコピー先位置にコピーする。そして、コントローラは、有効データの論理アドレスと、コピー先ブロック番号と、第3のブロック内のコピー先位置を示す第2のブロック内物理アドレスとを、ホスト2に通知する。これにより、GCにおいても、ホスト2がブロック番号(コピー元ブロック番号、コピー先ブロック番号)のみをハンドリングし、フラッシュストレージデバイス3がコピー先ブロック内のコピー先位置を決定する、という構成を実現できる。

40

【0144】

なお、フラッシュストレージデバイス3は、ストレージレイ内に設けられる複数のフラッシュストレージデバイス3の一つとして利用されてもよい。ストレージレイは、サーバ計算機のような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。ストレージレイは、このストレージレイ内の複数のフラッシュストレージデバイス3を制御するコントローラを含む。フラッシュストレージデバイス3がストレージレイに適用された場合には、このストレージレイのコントローラが、フラッシュストレージデバイス3のホスト2として機能してもよい。

【0145】

50

また、本実施形態では、不揮発性メモリとしてNAND型フラッシュメモリを例示した。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。

【0146】

図39は、ホスト2とフラッシュストレージデバイス3とのシステムアーキテクチャを示す。具体的には、図39は、ホスト2に含まれるライトデータバッファ51およびフラッシュトランザクション部52と、フラッシュストレージデバイス3に含まれるライト動作制御部21、リード動作制御部22および(GC動作制御部23を含む)最適化処理部53との関係を示す。

10

【0147】

ホスト2は、ライトデータをホストメモリ上のライトデータバッファ51に格納し、そしてライトコマンドをフラッシュストレージデバイス3に発行する。このライトコマンドは、このライトデータが存在するライトデータバッファ51上の位置を示すデータポイントと、このライトデータを識別するタグ(例えばLBA)と、このライトデータの長さ、ライトデータが書き込まれるべきブロックを示す識別子(ブロックアドレス、またはストリームID)とを含んでいてもよい。

20

【0148】

フラッシュストレージデバイス3は、以下のタイプ#1 - ストレージデバイス、タイプ#2 - ストレージデバイス、タイプ#3 - ストレージデバイスのうちの任意のストレージデバイスとして実現される。

タイプ#1 - ストレージデバイスは、ホスト2が、データが書き込まれるべきブロックとこのデータが書き込まれるべきページアドレスの双方を指定するタイプのストレージデバイスである。タイプ#1 - ストレージデバイスに適用されるライトコマンドは、ブロックアドレス、ページアドレス、データポイント、長さを含む。ブロックアドレスは、ホスト2から受信されるライトデータが書き込まれるべきブロックを指定する。ページアドレスは、このライトデータが書き込まれるべきこのブロック内のページを指定する。データポイントは、このライトデータが存在するホスト2内のメモリ上の位置を示す。長さは、このライトデータの長さを示す。

30

【0149】

タイプ#2 - ストレージデバイスは、ホスト2がデータが書き込まれるべきブロックを指定し、ストレージデバイスがこのデータが書き込まれるべきこのブロック内の位置(ページ)を指定するタイプのストレージデバイスである。タイプ#2 - ストレージデバイスに適用されるライトコマンドは、書き込まれるべきライトデータを識別するためのタグ(例えば、LBA、キー)、ブロックアドレス、データポイント、長さを含む。さらに、ライトコマンドは、QoSドメインIDを含んでもよい。QoSドメインIDは、NAND型フラッシュメモリを論理的に分割することによって得られる複数の領域の一つを指定する。複数の領域の各々は、複数のブロックを含む。タイプ#2 - ストレージデバイスは、不良ページ、ページ書き込み順序の制約を考慮して、データが書き込まれるべきページを決定することができる。

40

【0150】

つまり、フラッシュストレージデバイス3がタイプ#2 - ストレージデバイスとして実現されているケースにおいては、フラッシュストレージデバイス3は、ブロックをホスト2にハンドリングさせつつ、ページ書き込み順序制約、バッドページ、ページサイズ等を隠蔽することができる。この結果、ホスト2は、ブロック境界を認識でき、且つページ書き込み順序制約、バッドページ、ページサイズについては意識することなく、どのユーザデータがどのブロックに存在するかを管理することができる。

50

## 【 0 1 5 1 】

タイプ# 3 - ストレージデバイスは、ホスト 2 がデータを識別するタグ（例えば L B A ）を指定し、ストレージデバイスがこのデータが書き込まれるべきブロックおよびページの双方を決定するタイプのストレージデバイスである。タイプ# 3 - ストレージデバイスに適用されるライトコマンドは、書き込まれるべきライトデータを識別するためのタグ（例えば、L B A、キー）、ストリーム I D、データポインタ、長さを含む。ストリーム I D は、このライトデータに関連付けられたストリームの識別子である。フラッシュストレージデバイス 3 がタイプ# 3 - ストレージデバイスとして実現されているケースにおいては、フラッシュストレージデバイス 3 は、ストリーム I D それぞれとブロックアドレスそれぞれとの間のマッピングを管理する管理テーブルを参照して、このライトデータが書き込まれるべきブロックを決定する。さらに、フラッシュストレージデバイス 3 は、論理物理アドレス変換テーブルと称されるアドレス変換テーブルを使用して、タグ（L B A ）それぞれと N A N D 型フラッシュメモリの物理アドレスそれぞれとの間のマッピングを管理する。

10

## 【 0 1 5 2 】

フラッシュストレージデバイス 3 がタイプ# 1 - ストレージデバイスとして実現されている場合、フラッシュストレージデバイス 3 においては、ライト動作制御部 2 1 の制御の下、このブロックの識別子によって指定される書き込み先ブロックの書き込み動作の進行に合わせて、ライトデータバッファ 5 1 から内部バッファ（共有キャッシュ）3 1 へのデータ転送が D M A C によって実行される。このデータ転送は、N A N D 型フラッシュメモリ 5 のデータ書き込み単位と同じデータサイズの単位で実行される。ライト動作制御部 2 1 の制御の下、書き込むべきライトデータが、内部バッファ（共有キャッシュ）3 1 から、この書き込み先ブロックを含む N A N D 型フラッシュメモリチップ 1 5 に転送され、そして、ライト動作制御部 2 1 からこの N A N D 型フラッシュメモリチップ 1 5 に書き込み指示用の N A N D コマンドが送出される。

20

## 【 0 1 5 3 】

フラッシュストレージデバイス 3 がタイプ# 2 - ストレージデバイスとして実現されている場合には、ライト動作制御部 2 1 は、ホスト 2 から受信されるブロック割り当て要求に回答して、フリーブロックの一つを書き込み先ブロックとしてホスト 2 に割り当てる処理も実行する。ブロック割り当て要求は、Q o S ドメイン I D を含んでいてもよい。ライト制御部 2 1 は、この Q o S ドメイン I D に属するフリーブロックの一つを書き込み先ブロックとして決定し、この書き込み先ブロックのブロックアドレスをホスト 2 に通知する。これにより、ホスト 2 は、このブロックアドレスと、データポインタと、タグ（例えば L B A ）と、長さとを指定するライトコマンドを発行することができる。このライトデータがこの書き込み先ブロックに書き込まれた後、ライト動作制御部 2 1 は、このライトデータが書き込まれた書き込み先ブロックを示すブロックアドレスと、このライトデータが書き込まれたこの書き込み先ブロック内のページを示すページアドレスと、このライトデータのタグ（例えば L B A ）とをホスト 2 に通知する。ホスト 2 のフラッシュトランслーション部 5 2 は、タグ（例えば L B A ）それぞれと N A N D 型フラッシュメモリ 5 の物理アドレス（ブロックアドレス、ページアドレス、等）それぞれとの間のマッピングを管理するアドレス変換テーブルである L U T 4 1 1 を含む。ブロックアドレス、ページアドレス、およびタグ（例えば L B A ）がフラッシュストレージデバイス 3 から通知された場合、フラッシュトランслーション部 5 2 は、L U T 4 1 1 を更新し、通知されたタグ（例えば L B A ）に、通知された物理アドレス（ブロックアドレス、ページアドレス）をマッピングする。フラッシュトランслーション部 5 2 は、L U T 4 1 1 を参照することによって、リード要求に含まれるタグ（例えば L B A ）を物理アドレス（ブロックアドレス、ページアドレス）に変換することができ、これによって物理アドレスを含むリードコマンドをフラッシュストレージデバイス 3 に発行することができる。

30

40

## 【 0 1 5 4 】

フラッシュストレージデバイス 3 がタイプ# 1 - ストレージデバイスまたはタイプ# 2 -

50

ストレージデバイスとして実現されている場合、リード動作制御部 2 2 は、リードコマンドに含まれる物理アドレスに基づき、NAND型フラッシュメモリチップ 1 5 に読み出し指示用の NAND コマンドを送出する。フラッシュストレージデバイス 3 がタイプ # 3 - ストレージデバイスとして実現されている場合には、リード動作制御部 2 2 は、アドレス変換テーブルを参照して、リードコマンドに含まれるタグ (LBA) に対応する物理アドレスを取得し、取得した物理アドレスに基づき、NAND型フラッシュメモリチップ 1 5 に読み出し指示用の NAND コマンドを送出する。

【 0 1 5 5 】

リード動作制御部 2 2 の制御の下、NAND型フラッシュメモリチップ 1 5 から読み出されたデータは、内部バッファ (共有キャッシュ) 3 1 に転送される。そして、リード動作制御部 2 2 の制御の下、内部バッファ (共有キャッシュ) 3 1 からホスト 2 へのデータ転送が DMA C によって実行される。また、リード動作制御部 2 2 は、読み出すべきリードデータがホスト 2 のライトデータバッファ 5 1 内に存在する場合、ライトデータバッファ 5 1 からリードデータを取得することができる。あるいは、ホスト 2 に対して、リードデータをライトデータバッファ 5 1 から取得することを指示してもよい。なお、ライトデータバッファ 5 1 上のライトデータが格納される領域は、ライト動作制御部 2 1 による NAND型フラッシュメモリ 5 への書き込みが完了した場合にライト動作制御部 2 1 からホスト 2 へ送信される解放可能通知により、ホスト 2 側において解放される。たとえば、ライト動作制御部 2 1 による NAND型フラッシュメモリ 5 への書き込みが失敗し、ライトデータが別の場所 (異なるページやブロック) に書き込まれる場合、その書き込みに必要なデータについて、まだ解放されていない、ホスト 2 のライトデータバッファ 5 1 の領域からフラッシュストレージデバイス 3 の内部バッファ (共有キャッシュ) 3 1 へのデータ転送が再実行される。データの再書き込みは、エラーが検出された範囲で実行されてもよいし、ライトコマンドの範囲すべてで実行されてもよい。解放可能通知は、ライトコマンド単位でホスト 2 に通知されてもよいし、ホスト 2 のデータ使用単位でホスト 2 に通知されてもよい。

【 0 1 5 6 】

(GC動作制御部 2 3 を含む)最適化処理部 5 3 は、たとえば、ホスト 2 から受信されるブロック解放要求に応答して、割り当て済みのブロックをフリーブロックに戻す処理などを実行する。ホスト 2 は、ブロック解放要求を、ブロックリユースコマンドとしてフラッシュストレージデバイス 3 に送信する。ブロックリユースコマンドで指定され得る割り当て済みのブロックは、フラッシュストレージデバイス 3 がタイプ # 1 - ストレージデバイスとして実現され、かつ、ホスト 2 がフリーブロック群を管理しない場合、または、フラッシュストレージデバイス 3 がタイプ # 2 - ストレージデバイスとして実現されている場合において、ブロックアロケートコマンドとしてホスト 2 から受信したブロック割り当て要求に応答して、フリーブロックの中から割り当てたブロックである。また、最適化処理部 5 3 は、たとえば、ホスト 2 から受信される GC 制御コマンドに応答して、あるブロックのデータを別のブロックにコピーする処理などを実行する。

【 0 1 5 7 】

また、フラッシュストレージデバイス 3 がホスト 2 から受信する各種コマンドには、優先度が含まれ得る。つまり、フラッシュストレージデバイス 3 は、先にホスト 2 から受信したコマンドに優先して、後からホスト 2 から受信したコマンドを実行し得る。コマンドの実行順の制御は、たとえば、ホスト 2 から受信される各種コマンドが一時的に格納される I/O コマンドキューからコマンドを取り出す際にコマンド間の優先度を比較することによって実現できる。I/O コマンドキューは、QoS ドメイン毎に設けられるものであってもよいし、後述する仮想ストレージデバイス (VD: Virtual Device) 毎に設けられるものであってもよいし、各フラッシュストレージデバイス 3 に 1 つずつ設けられるものであってもよい。

【 0 1 5 8 】

NAND型フラッシュメモリ 5 が複数の NAND型フラッシュメモリチップ 1 5 を含むフ

10

20

30

40

50

ラッシュストレージデバイス3においては、1以上の仮想ストレージデバイスを定義することができる。図40は、フラッシュストレージデバイス3上における仮想ストレージデバイスの定義例を示す。

【0159】

図40中、(A)は、仮想ストレージデバイス間でNANDインタフェース13に接続されるチャンネルが共有される、複数の仮想ストレージデバイスの定義例を示す。(B)は、仮想ストレージデバイス間でNANDインタフェース13に接続されるチャンネルが共有されない、複数の仮想ストレージデバイスの定義例を示す。(C)は、NAND型フラッシュメモリ5が含む複数のNAND型フラッシュメモリチップ15すべてを用いた、1つの仮想ストレージデバイスの定義例を示す。(D)は、NAND型フラッシュメモリ5が含む複数のNAND型フラッシュメモリチップ15それぞれを個別に用いた、NAND型フラッシュメモリチップ15と同数、つまり最大数の仮想ストレージデバイスの定義例を示す。

10

【0160】

このように、フラッシュストレージデバイス3上には、様々な形態で、1以上の仮想ストレージデバイスを定義することができる。仮想ストレージデバイスを定義すると、たとえばNAND型フラッシュメモリチップ15の消耗度を監視するためのウェア監視などを、この仮想ストレージデバイス毎に実行することができる。

【0161】

また、1以上の仮想ストレージデバイスを定義し得るフラッシュストレージデバイス3においては、仮想ストレージデバイス毎にQoSドメインを管理することができる。図41は、仮想ストレージデバイス毎にQoSドメインが管理される例を示す。

20

フラッシュストレージデバイス3のブロックは、同一の仮想ストレージデバイス上に定義されるQoSドメイン間で共有される。ブロックを取り扱う単位は、複数のブロックで構成されるスーパーブロックの単位であってもよい。つまり、スーパーブロックが、QoSドメイン間で共有されてもよい。たとえばエンドユーザ毎にQoSドメインが割り当てられている場合において、QoSドメインを示すQoSドメインIDを含むブロックアロケートコマンドがホスト2から受信されると、仮想ストレージデバイス内で共有されるフリーブロック群の中の1つのフリーブロックが、QoSドメインIDで示されるQoSドメインに割り当てられる。

30

【0162】

一方、QoSドメインIDとブロックアドレスとを含むブロックリユースコマンドがホスト2から受信されると、QoSドメインIDで示されるQoSドメインに割り当てられているブロックの中のブロックアドレスで示されるブロックが、フリーブロックとしてフリーブロック群に戻される。QoSドメインに割り当てられているブロックをフリーブロックとしてフリーブロック群に戻すことは、ブロックを解放するとも称される。解放されたブロックは、その後、たとえばホスト2からのブロックアロケートコマンドなどによって、その仮想ストレージデバイス内のいずれのQoSドメインにも割り当てられ得る。

【0163】

ところで、あるQoSドメイン内のあるブロックを対象とするブロックリユースコマンドがホスト2から受信された際、フラッシュストレージデバイス3内において、そのブロックを対象とするリードコマンドが実行中または未実行の状態であった場合、そのリードコマンドよりも先にブロックリユースコマンドが実行されてしまうと、たとえば値が不定のデータがホスト2に返却されるなどのおそれがある。前述したように、フラッシュストレージデバイス3がホスト2から受信する各種コマンドには、優先度が含まれ得るので、フラッシュストレージデバイス3においては、先にホスト2から受信したリードコマンドに優先して、後からホスト2から受信したブロックリユースコマンドが実行され得る。また、リードコマンドのほか、たとえば、GC制御コマンドに回答して、そのブロック内のデータを他のブロックにコピーする場合にも、同様の事態が発生し得る。つまり、そのブロックを対象とするデータの読み出し処理を実行中または未実行の状態にある場合、この実

40

50

行中または未実行のデータの読み出し処理において意図しないデータが読み出されてしまうおそれがある。

【0164】

このような事態をホスト2側の制御によって防止するためには、ホスト2において、たとえば仕掛り中のデータの読み出し処理の有無をブロック毎に管理することなどが必要となる。そこで、このフラッシュストレージデバイス3は、このような事態を防止するための仕組みを備えて、ホスト2の負担を軽減するようにしてもよい。

【0165】

フラッシュストレージデバイス3は、ホスト2からブロックリユースコマンドが受信された場合において、そのブロックリユースコマンドで指定されるブロックを対象とするデータの読み出し処理が実行中または未実行の状態であるならば、ホスト2に対してエラーを通知し、あるいは、実行中または未実行の処理が終了するまでブロックリユースコマンドの実行を保留し、実行中または未実行の処理が終了したら、ブロックリユースコマンドを実行する。

【0166】

この仕組みをフラッシュストレージデバイス3が備えることで、ホスト2は、たとえば解放しようとするブロックを対象とする仕掛り中のデータの読み出し処理の有無などを気にせず、ブロックリユースコマンドをフラッシュストレージデバイス3に送信することができるようになる。つまり、ホスト2の負担を軽減することが実現される。

【0167】

この仕組みは、たとえば、最適化処理部53が、ブロックリユースコマンドの受信時またはブロックリユースコマンドの実行時、そのブロックリユースコマンドで指定されるブロックを対象とするリードコマンドやGC制御コマンドがI/Oコマンドキュー42に格納されていないかを検索することによって実現できる。図41においては、I/Oコマンドキュー42がQoSドメイン毎に設けられる例が示されており、この場合は、最適化処理部53は、ブロックリユースコマンドに含まれるQoSドメインIDで示されるQoSドメインに対して設けられるI/Oコマンドキュー42について、ブロックリユースコマンドに含まれるブロックアドレスで示されるブロックを対象とするリードコマンドまたはGC制御コマンドが存在していないかどうかを調べる。I/Oコマンドキュー42が仮想ストレージデバイス毎または各フラッシュストレージデバイスに1つずつ設けられる場合には、最適化処理部53は、そのI/Oコマンドキュー42について、ブロックリユースコマンドに含まれるQoSドメインIDで示されるQoSドメインの中のブロックリユースコマンドに含まれるブロックアドレスで示されるブロックを対象とするリードコマンドまたはGC制御コマンドが存在していないかどうかを調べる。存在する場合、最適化処理部53は、ホスト2に対してエラーを通知し、あるいは、I/Oコマンドキュー42に存在する、ブロックリユースコマンドで指定されるブロックを対象とするリードコマンドまたはGC制御コマンドが終了するまでブロックリユースコマンドの実行を保留し、それらが終了したら、ブロックリユースコマンドを実行する。

【0168】

あるいは、この仕組みは、たとえば、フリーブロック群から選ばれてQoSドメインに割り当てられたブロックそれぞれについて、そのブロックを対象として実行中のリードコマンドの数、および、そのブロックをコピー元として実行中のGC制御コマンドの数を示すカウンタをメタデータなどとして設けることによって実現できる。たとえば、リード動作制御部21や(GC動作制御部23を含む)最適化処理部53は、あるブロックを対象とするデータの読み出し処理を実行する場合、そのブロックのカウンタの値を1つ加算する。また、リード動作制御部21や最適化処理部53は、データの読み出し処理を終了した場合、そのブロックのカウンタの値を1つ減算する。最適化処理部53は、ブロックリユースコマンドの受信時またはブロックリユースコマンドの実行時、そのブロックリユースコマンドで指定されるブロックのカウンタの値が0でない場合、ホスト2に対してエラーを通知し、あるいは、カウンタの値が0になるまでブロックリユースコマンドの実行を保

10

20

30

40

50



留し、カウンタの値が 0 になったら、ブロックリユースコマンドを実行する。

【 0 1 6 9 】

図 4 2 は、ブロックリユースコマンド受信時におけるフラッシュストレージデバイス 3 の動作手順（第 1 ケース）を示すフローチャートである。なお、ここでは、ブロックリユースコマンド受信時を想定するが、以下に説明する動作は、ブロックリユースコマンド実行時に行われるものであってもよい。

【 0 1 7 0 】

最適化処理部 2 3 は、ブロックリユースコマンドがホスト 2 から受信された場合（ステップ A 1 ）、そのブロックリユースコマンドで指定されるブロックを対象とする実行中または実行待ちの読み出し処理が存在するか否かを判定する（ステップ A 2 ）。存在しない場合（ステップ A 2 : NO ）、最適化処理部 2 3 は、指定されたブロックをフリーブロック化（解放）し、リユース完了を示すレスポンスをホスト 2 に返す（ステップ A 3 ）。 10

【 0 1 7 1 】

一方、存在する場合（ステップ A 2 : YES ）、最適化処理部 2 3 は、エラーをホスト 2 に通知する（ステップ A 4 ）。 10

図 4 3 は、ブロックリユースコマンド受信時におけるフラッシュストレージデバイス 3 の動作手順（第 2 ケース）を示すフローチャートである。ここでも、ブロックリユースコマンド受信時を想定するが、以下に説明する動作は、ブロックリユースコマンド実行時に行われるものであってもよい。

【 0 1 7 2 】

最適化処理部 2 3 は、ブロックリユースコマンドがホスト 2 から受信された場合（ステップ A 1 1 ）、そのブロックリユースコマンドで指定されるブロックを対象とする実行中または実行待ちの読み出し処理が存在するか否かを判定する（ステップ A 1 2 ）。 20

【 0 1 7 3 】

存在しない場合（ステップ A 1 2 : NO ）、最適化処理部 2 3 は、即時的に、指定されたブロックをフリーブロック化（解放）し、リユース完了を示すレスポンスをホスト 2 に返す（ステップ A 1 4 ）。一方、存在する場合（ステップ A 1 2 : YES ）、続いて、最適化処理部 2 3 は、該当する読み出し処理がすべて完了したか否かの判定を行い（ステップ A 1 3 ）、すべて完了したら（ステップ A 1 3 : YES ）、指定されたブロックをフリーブロック化（解放）し、リユース完了を示すレスポンスをホスト 2 に返す（ステップ A 1 4 ）。 30

【 0 1 7 4 】

なお、以上では、ブロックリユースコマンドがホスト 2 から受信された際、そのブロックリユースコマンドで指定されるブロックを対象とする実行中または実行待ちのリードコマンドや GC 制御コマンドを受信済みである場合におけるブロックリユースコマンドの取り扱いを説明した。このフラッシュストレージデバイス 3 は、さらに、ブロックリユースコマンドがホスト 2 から受信された後、そのブロックリユースコマンドで指定されるブロックを対象とするリードコマンドや GC 制御コマンドがホスト 2 から受信された場合、エラーをホスト 2 に返却するようにしてもよい。

【 0 1 7 5 】

また、図 4 1 を参照して説明したように、フラッシュストレージデバイス 3 のブロックは、たとえば仮想ストレージデバイス毎に管理される QoS ドメイン間で共有される。つまり、たとえば仮想ストレージデバイス毎にフリーブロック群が管理され、そのフリーブロック群の中から各 QoS ドメインへフリーブロックが割り当てられていく。

【 0 1 7 6 】

ブロックにデータを書き込むケースは、大きく分けて、ホスト 2 から受信されるライトコマンドに応じて、ホスト 2 のライトデータバッファ 5 1 に格納されているデータを書き込むケースと、ホスト 2 から受信される GC 制御コマンドに応じて、フラッシュストレージデバイス 3 の別のブロックに格納されているデータを書き込むケースとが存在する。ホスト 2 のライトデータバッファ 5 1 に格納されているデータは新しく、フラッシュストレ 50

ジデバイス3の別のブロックに格納されているデータは古い。したがって、これらのデータを同一のブロックに混在させると、ライトアンプリフィケーションが悪化するおそれがある。そこで、このフラッシュストレージデバイス3は、データが書き込まれるべきブロックおよびページの双方をストレージデバイスが決定するタイプ#3 - ストレージデバイスとして実現されている場合、QoSドメイン毎に、ホスト2からのデータを書き込むためのブロックと、フラッシュストレージデバイス3内のデータをコピーするためのブロックとに分離する仕組みを備えてもよい。ブロックがスーパーブロックの単位で取り扱われる場合には、スーパーブロックを、ホスト2からのデータを書き込むためのスーパーブロックと、フラッシュストレージデバイス3内のデータをコピーするためのスーパーブロックとに分離する。つまり、QoSドメイン毎に、空きページを含むブロックとして、ホスト2からのデータを書き込むためのブロックと、フラッシュストレージデバイス3内のデータをコピーするためのブロックとを各々確保する。

10

## 【0177】

このブロックの分離は、たとえば、フリーブロック群から選ばれてQoSドメインに割り当てられるブロックそれぞれについて、その用途を示す属性情報をメタデータなどとして保持することによって実現できる。QoSドメインが利用を開始された時、ホスト2からのデータを書き込むためのブロックと、フラッシュストレージデバイス3内のデータをコピーするためのブロックとのいずれも確保されていない。なお、ブロックが確保されているとは、空きページを含むブロックが割り当てられていることである。

## 【0178】

たとえば、ライト動作制御部21は、あるQoSドメインに関してホスト2からのデータの書き込みを実行する場合、そのQoSドメインにおいて属性情報がホスト2からのデータを書き込むためのブロックであることを示すブロックが確保されていなければ、フリーブロック群の中の1つのフリーブロックをそのQoSドメイン用に取得して、取得したブロックにデータを書き込む。この取得時、ライト動作制御部21は、そのブロックがホスト2からのデータを書き込むためのブロックであることを示す属性情報をメタデータとして記録する。一方、属性情報がホスト2からのデータを書き込むためのブロックであることを示す、ホスト2からのデータを書き込むためのブロックが確保されているならば、ライト動作制御部21は、そのブロック中の最後に書き込みが行われたページに続くページからデータの書き込みを実行する。データの書き込み途中で、そのブロックの最後のページまでデータが書き込まれると、ブロック未確保の状態に戻るので、ライト動作制御部21は、フリーブロック群の中の1つのフリーブロックをそのQoSドメイン用に取得して、取得したブロックに続きのデータを書き込む。この取得時にも、ライト動作制御部21は、そのブロックがホスト2からのデータを書き込むためのブロックであることを示す属性情報をメタデータとして記録する。

20

30

## 【0179】

また、たとえば、(GC動作制御部23を含む)最適化処理部53は、あるQoSドメインに関してデータのコピーを実行する場合、そのQoSドメインにおいて属性情報がフラッシュストレージデバイス3内のデータをコピーするためのブロックが確保されていなければ、フリーブロック群の中の1つのフリーブロックをそのQoSドメイン用に取得して、取得したブロックにデータを書き込む(コピーする)。この取得時、最適化処理部53は、そのブロックがフラッシュストレージデバイス3内のデータをコピーするためのブロックであることを示す属性情報をメタデータとして記録する。一方、属性情報がフラッシュストレージデバイス3内のデータをコピーするためのブロックであることを示す、フラッシュストレージデバイス3内のデータをコピーするためのブロックが確保されているならば、最適化処理部53は、そのブロック中の最後に書き込みが行われたページに続くページからデータの書き込みを実行する。データの書き込み途中で、そのブロックの最後のページまでデータが書き込まれると、ブロック未確保の状態に戻るので、最適化処理部53は、フリーブロック群の中の1つのフリーブロックをそのQoSドメイン用に取得して、取得したブロックにデータを書き込む。この取得時にも、最適化処理部53は、そのブ

40

50

ロックがフラッシュストレージデバイス 3 内のデータをコピーするためのブロックであることを示す属性情報をメタデータとして記録する。

【 0 1 8 0 】

このように、このフラッシュストレージデバイス 3 は、ホスト 2 からの新しいデータを書き込むためのブロックと、フラッシュストレージデバイス 3 内の古いデータをコピーするためのブロックとを分離することによって、ライトアンプリフィケーションの悪化を防ぐことができる。

【 0 1 8 1 】

また、図 3 9 を参照して説明したように、リード動作制御部 2 2 は、読み出すべきリードデータがホスト 2 のライトデータバッファ 5 1 内に存在する場合、ライトデータバッファ 5 1 からリードデータを取得することができる。一方、ライトデータバッファ 5 1 上のライトデータが格納される領域は、ライト動作制御部 2 1 からホスト 2 へ解放可能通知が送信されると、ホスト 2 側において解放される。そこで、このフラッシュストレージデバイス 3 は、ライトデータバッファ 5 1 内に存在するライトデータ中のデータが対象となるリードコマンドがホスト 2 から受信された場合、そのリードコマンドが終了するまで、そのデータが格納される領域についての解放可能通知をホスト 2 へ送信しないようにする仕組みを備えてもよい。

【 0 1 8 2 】

この仕組みは、たとえば、ホスト 2 のライトデータバッファ 5 1 に格納されるライトデータに対し、ホスト 2 から受信されるライトコマンド単位またはホスト 2 のデータ使用単位で、データの書き込み処理数および読み出し処理数の残数を示すカウンタをメタデータなどとして設けることによって実現できる。カウンタは、たとえば、解放可能通知がホスト 2 に通知される単位と一致させて設けられる。解放可能通知がライトコマンド単位でホスト 2 に通知される場合において、カウンタがホスト 2 のデータ使用単位で設けられてもよい。

【 0 1 8 3 】

カウンタがホスト 2 のデータ使用単位で設けられるものと想定すると、ライト動作制御部 2 1 は、各カウンタの初期値として、NAND 型フラッシュメモリ 5 へのデータの書き込みに必要なデータの転送回数 + 1 をセットする。+ 1 は、エラーが検出された場合の再書き込み処理のために加算しておくものである。

【 0 1 8 4 】

ライト動作制御部 2 1 は、NAND 型フラッシュメモリ 5 へデータを転送する都度、対応するカウンタの値を 1 ずつ減算する。あるデータ使用単位についてデータの転送が終了したとすると、その時点で、一般的には、カウンタの値は 1 となる。ライト動作制御部 2 1 は、転送したデータすべてが NAND 型フラッシュメモリ 5 に書き込まれ、エラーが検出された場合の再書き込み処理が不要となったことが確定すると、対応するカウンタの値をさらに 1 減算する。この時点で、一般的には、カウンタの値は 0 となる。仮に、解放可能通知がホスト 2 のデータ使用単位でホスト 2 に通知されるものと想定すると、ライト動作制御部 2 1 は、カウンタの値が 0 になったことを検知した場合、対応する領域についての解放可能通知をホスト 2 へ通知する。なお、エラーが検出された場合、ライト動作制御部 2 1 は、再書き込み処理に必要なデータの転送回数をカウンタに再加算する。NAND 型フラッシュメモリ 5 にデータを転送し終えた後にエラーが検出されても、カウンタの値が 0 になっていないことから、解放可能通知はホスト 2 へ通知されず、再書き込み処理に必要なデータはホスト 2 のライトデータバッファ 5 1 に存在する。したがって、ホスト 2 のライトデータバッファ 5 1 からフラッシュストレージデバイス 3 の内部バッファ（共有キャッシュ）3 1 へのデータ転送を再実行することができる。

【 0 1 8 5 】

リード動作制御部 2 1 も、ホスト 2 のライトデータバッファ 5 1 に存在するライトデータ中のデータが対象となるリードコマンドがホスト 2 から受信された場合、そのデータに対応するカウンタの値を 1 加算する。そして、その読み出し処理が終了したら、リード動作

10

20

30

40

50

制御部 2 1 は、対応するカウンタの値を 1 減算する。

【 0 1 8 6 】

ライトデータバッファ 5 1 上に存在するライトデータ中の読み出し処理の対象となるデータについては、NAND型フラッシュメモリ 5 への書き込みが終了しても、対応するカウンタの値は 0 にはならず、解放可能通知がホスト 2 へ通知されない。つまり、カウンタの値を 1 加算することで、リード動作制御部 2 1 は、ライトデータバッファ 5 1 上の目的の領域を解放禁止状態とする。したがって、ホスト 2 のライトデータバッファ 5 1 に存在するライトデータ中のデータが対象となるリードコマンドがホスト 2 から受信されている状況下において、そのデータを含むライトデータバッファ 5 1 上の領域がホスト 2 側において解放されてしまうことがない。

10

【 0 1 8 7 】

なお、ホスト 2 のライトデータバッファ 5 1 に存在するライトデータ中のデータが対象となるリードコマンドがホスト 2 から受信された場合、リード動作制御部 2 1 は、そのデータを必ずしもライトデータバッファ 5 1 から読み出さなくてもよく、ライトデータの NAND型フラッシュメモリ 5 への書き込みが終了し、読み出し可能な状態になっているならば、NAND型フラッシュメモリ 5 から読み出してもよい。この場合、ライトデータバッファ 5 1 上のライトデータは、たとえば、予備のデータなどとして活用し得る。

【 0 1 8 8 】

解放可能通知がライトコマンド単位でホスト 2 に通知され、カウンタがホスト 2 のデータ使用単位で設けられる場合、ライト動作制御部 2 1 は、ライトコマンドで書き込まれるべきライトデータに対応する複数のカウンタすべての値が 0 になった時点で、解放可能通知をホスト 2 に通知する。

20

【 0 1 8 9 】

また、リード動作制御部 2 2 が、ライトデータバッファ 5 1 からリードデータを取得することができる点に着目し、このフラッシュストレージデバイス 3 は、タイプ # 2 - ストレージデバイスとして実現されている場合、ホスト 2 からのライトデータを書き込み予定のページアドレスを、NAND型フラッシュメモリ 5 への書き込みの終了を待機することなく、ホスト 2 に通知する仕組みを備えてもよい。この仕組みをフラッシュストレージデバイス 3 が備える場合、ホスト 2 は、たとえば、ライトコマンドで書き込んだデータがフラッシュストレージデバイス 3 において読み出し可能な状態となるまで待たされることなく、ライトコマンドで書き込んだデータ中のデータを対象とするリードコマンドを速やかに発行することが可能となる。

30

【 0 1 9 0 】

この仕組みは、たとえば、リード動作制御部 2 1 が、ライトコマンドの受信時にホスト 2 から通知されるライトデータバッファ 5 1 上のライトデータに関する情報を、たとえばメタデータなどとして設ける書き込み先ブロック毎のライトバッファリストに登録し、ホスト 2 のデータ使用単位毎に、ライトデータを書き込み予定のページアドレスをホスト 2 に通知することによって実現できる。ライトバッファリストに登録されたライトデータのサイズが、書き込み先ブロックの残り書き込み領域のサイズより大きくても構わない。この場合、リード動作制御部 2 1 は、まず、その書き込み先ブロックに書き込み可能な分について書き込み予定のページアドレスをホスト 2 に通知し、その書き込み先ブロックへの書き込みが終了して、新たな書き込み先ブロックが確保された後、残りの分について書き込み予定のページアドレスをホスト 2 に通知する。ライトデータはすべてライトデータバッファ 5 1 上に存在し、書き込み先ブロックの確保は極めて短時間で済むので、ライトデータの書き込みがブロックを跨ぐ場合であっても実用上の問題はない。

40

【 0 1 9 1 】

NAND型フラッシュメモリ 5 への書き込み時にエラーが検出された場合、リード動作制御部 2 1 は、新たに決定される書き込み予定のページアドレスをホスト 2 に改めて通知する。なお、ページアドレスのホスト 2 への通知は、前述したように、書き込み予定のページアドレスが決定した時点で行ってもよいし、ホスト 2 のデータ使用単位で書き込みが終

50

了する毎に行ってもよい。前者の場合、エラーが検出された場合にホスト 2 への通知が複数回発生し得るが、ホスト 2 への通知が速い。後者の場合、ホスト 2 への通知が前者の場合よりも遅くはなるが、NAND型フラッシュメモリ 5 への書き込み時にエラーが検出された場合であっても、その回数に関係なく、1 度の通知で済むことになる。

【0192】

次に、ホスト 2 のライトデータバッファ 5 1 を用いるライトコマンド処理を含む、このフラッシュストレージデバイス 3 によって実行される各種 I/O コマンド処理について詳細に説明する。

図 4 4 は、フラッシュストレージデバイス 3 によって実行される I/O コマンド処理を示す。

10

【0193】

上述したように、本実施形態では、フラッシュストレージデバイス 3 はタイプ # 1 - ストレージデバイス、タイプ # 2 - ストレージデバイス、タイプ # 3 - ストレージデバイスのいずれであってもよいが、図 4 4 では、フラッシュストレージデバイス 3 がタイプ # 1 - ストレージデバイスである場合を例示する。

【0194】

ホスト 2 によって発行される各ライトコマンドは、ブロックアドレス、ページアドレス、データポインタ、長さを含む。発行された各ライトコマンドは I/O コマンドキュー 4 2 に入れられる。ホスト 2 によって発行される各リードコマンドも、ブロックアドレス、ページアドレス、データポインタ、長さを含む。発行された各リードコマンドも I/O コマンドキュー 4 2 に入れられる。

20

【0195】

ホスト 2 がライトデータの書き込みをフラッシュストレージデバイス 3 に要求することを望む場合、ホスト 2 は、まず、このライトデータをホストメモリ上のライトデータバッファ 5 1 に格納し、そして、ライトコマンドをフラッシュストレージデバイス 3 に発行する。このライトコマンドは、このライトデータが書き込まれるべき書き込み先ブロックを示すブロックアドレスと、このライトデータが書き込まれるべきこの書き込み先ブロック内のページを示すページアドレスと、このライトデータが存在するライトデータバッファ 5 1 内の位置を示すデータポインタと、このライトデータの長さを含む。

【0196】

フラッシュストレージデバイス 3 は、プログラム/リードシーケンサ 4 1 を含む。このプログラム/リードシーケンサ 4 1 は、上述のライト制御部 2 1 およびリード制御部 2 2 によって実現される。プログラム/リードシーケンサ 4 1 は、I/O コマンドキュー 4 2 に入れられたコマンドそれぞれを任意の順序で実行することができる。

30

【0197】

プログラム/リードシーケンサ 4 1 が、ある同じ書き込み先ブロックを指定する 1 以上のライトコマンドを I/O コマンドキュー 4 2 から取得した後、プログラム/リードシーケンサ 4 1 は、この書き込み先ブロックの書き込み動作の進行に合わせて、この書き込み先ブロックに書き込むべき次のライトデータ（例えば、1 ページサイズ分のライトデータ）を内部バッファ（共有キャッシュ）3 1 またはライトデータバッファ 5 1 から取得するための転送要求を内部バッファ（共有キャッシュ）3 1 に送出する。この転送要求は、データポインタと長さを含んでもよい。この転送要求に含まれるデータポインタは、1 つのライトコマンドに関連付けられたライトデータを分割、または同じ書き込み先ブロックを指定する 2 以上のライトコマンドに関連付けられた 2 以上のライトデータを結合する処理によって算出される。つまり、プログラム/リードシーケンサ 4 1 は、同じ書き込み先ブロックを示す識別子を有する一つ以上のライトコマンドに関連付けられたライトデータの集合を、その先頭から NAND 型フラッシュメモリ 5 のデータ書き込み単位と同じサイズを有する境界で区切り、各境界に対応するホストメモリ内の位置を特定する。これによって、プログラム/リードシーケンサ 4 1 は、書き込み単位と同じサイズの単位でライトデータをホスト 2 から取得することができる。

40

50

## 【 0 1 9 8 】

この転送要求に含まれるデータポインタは、この1ページサイズ分のライトデータが存在するライトデータバッファ51上の位置を示す。この1ページサイズ分のライトデータは、この書き込み先ブロックを指定する複数のライトコマンドに関連付けられた複数の小さなサイズのライトデータの集合であってもよいし、この書き込み先ブロックを指定する一つのライトコマンドに関連付けられた大きなサイズのライトデータの一部であってもよい。

## 【 0 1 9 9 】

さらに、プログラムノードシーケンサ41は、この1ページサイズ分のライトデータが書き込まれるべき書き込み先ブロックのブロックアドレスと、この1ページサイズ分のライトデータが書き込まれるべきページのページアドレスとを内部バッファ（共有キャッシュ）31に送出する。

10

## 【 0 2 0 0 】

フラッシュストレージデバイス3のコントローラ4は内部バッファ（共有キャッシュ）31を制御するキャッシュコントローラを含んでいてもよい。この場合、このキャッシュコントローラは内部バッファ（共有キャッシュ）31をあたかも制御ロジックであるかのように動作させることができる。内部バッファ（共有キャッシュ）31と複数の書き込み先ブロック#0、#1、#2、...、#nとの間には、複数のフラッシュコマンドキュー43が存在する。これらフラッシュコマンドキュー43は、複数のNAND型フラッシュメモリチップにそれぞれ対応づけられている。

20

## 【 0 2 0 1 】

内部バッファ（共有キャッシュ）31、つまりキャッシュコントローラは、転送要求によって指定されたこの1ページサイズ分のライトデータが内部バッファ（共有キャッシュ）31に存在するか否かを判定する。

この転送要求によって指定されたこの1ページサイズ分のライトデータが内部バッファ（共有キャッシュ）31に存在するならば、内部バッファ（共有キャッシュ）31、つまりキャッシュコントローラは、この1ページサイズ分のライトデータを、このライトデータが書き込まれるべき書き込み先ブロックを含むNAND型フラッシュメモリチップに転送する。さらに、内部バッファ（共有キャッシュ）31、つまりキャッシュコントローラは、このライトデータが書き込まれるべき書き込み先ブロックを含むNAND型フラッシュメモリチップに、この書き込み先ブロックのブロックアドレス、このライトデータが書き込まれるべきページアドレス、書き込み指示用のNANDコマンド（フラッシュライトコマンド）を、フラッシュコマンドキュー43を介して送出する。フラッシュコマンドキュー43は、NAND型フラッシュメモリチップ毎に設けられている。このため、内部バッファ（共有キャッシュ）31、つまりキャッシュコントローラは、このライトデータが書き込まれるべき書き込み先ブロックを含むNAND型フラッシュメモリチップに対応するフラッシュコマンドキュー43に、この書き込み先ブロックのブロックアドレス、このライトデータが書き込まれるべきページアドレス、書き込み指示用のNANDコマンド（フラッシュライトコマンド）を入れる。

30

## 【 0 2 0 2 】

なお、内部バッファ（共有キャッシュ）31からNAND型フラッシュメモリチップへのこの1ページサイズ分のライトデータの転送が、このライトデータをNAND型フラッシュメモリチップに書き込むために必要な最終回のデータ転送であるならば、内部バッファ（共有キャッシュ）31、つまりキャッシュコントローラは、このライトデータを内部バッファ（共有キャッシュ）31から破棄し、このライトデータが格納されていた領域を空き領域として確保する。NAND型フラッシュメモリチップにデータを一回転送することを伴う書き込み動作（例えば、フル・シーケンス書き込み動作、等）によってライトデータを書き込み先ブロックに書き込むケースにおいては、NAND型フラッシュメモリチップへの初回のデータ転送が最終回のデータ転送となる。一方、NAND型フラッシュメモリチップにデータを複数回転送することを伴う書き込み動作（例えば、フォギー・ファイ

40

50

ン書き込み動作)によってライトデータを書き込み先ブロックに書き込むケースにおいては、最後のファイン書き込みのために必要なNAND型フラッシュメモリチップへのデータ転送が最終回のデータ転送となる。

#### 【0203】

次に、転送要求によって指定されたこの1ページサイズ分のライトデータが内部バッファ(共有キャッシュ)31に存在しない場合について説明する。

この転送要求によって指定されたこの1ページサイズ分のライトデータが内部バッファ(共有キャッシュ)31に存在しないならば、内部バッファ(共有キャッシュ)31、つまりキャッシュコントローラは、この転送要求(データポインタ、長さ)をDMAC15に送出する。この転送要求(データポインタ、長さ)に基づいて、DMAC15は、この1ページサイズ分のライトデータをホストメモリ上のライトデータバッファ51から内部バッファ(共有キャッシュ)31に転送する。このデータ転送が終了すると、DMAC15は、転送完了(Done)と、このデータポインタ、この長さとを、内部バッファ(共有キャッシュ)31、つまりキャッシュコントローラに通知する。

10

#### 【0204】

内部バッファ(共有キャッシュ)31に空き領域が存在するならば、内部バッファ(共有キャッシュ)31、つまりキャッシュコントローラは、DMA転送によってライトデータバッファ51から取得されたライトデータを、この空き領域に格納する。

内部バッファ(共有キャッシュ)31に空き領域が存在しないならば、内部バッファ(共有キャッシュ)31、つまりキャッシュコントローラは、内部バッファ(共有キャッシュ)31内の最も古いライトデータを内部バッファ(共有キャッシュ)31から破棄し、最も古いライトデータが格納されていた領域を空き領域として確保する。そして、内部バッファ(共有キャッシュ)31、つまりキャッシュコントローラは、DMA転送によってライトデータバッファ51から取得されたライトデータを、この空き領域に格納する。

20

#### 【0205】

フォギー・ファイン書き込み動作のような多段階の書き込み動作が使用されるケースにおいては、キャッシュコントローラは、フォギー書き込み動作のような第1段階の書き込み動作が終了している内部バッファ(共有キャッシュ)31内のライトデータのうちで、最も古いライトデータを破棄する。

#### 【0206】

データ書き込み量の少ない書き込み先ブロックへのデータ書き込み動作の進行速度に比べ、データ書き込み量の多い書き込み先ブロックへのデータ書き込み動作の進行速度は速くなる傾向がある。このため、データ書き込み量の多い書き込み先ブロックに書き込まれるべきライトデータは頻繁にライトデータバッファ51から内部バッファ(共有キャッシュ)31に転送される。この結果、この最も古いライトデータは、ホスト2から書き込まれるデータ量が比較的少ない書き込み先ブロックへのライトデータである可能性が高い。したがって、フォギー書き込み動作のような第1段階の書き込み動作が終了している内部バッファ(共有キャッシュ)31内のライトデータのうちで最も古いライトデータを破棄するという方法を使用することにより、ホスト2とフラッシュストレージデバイス3との間のデータトラフィックを効率よく低減することが可能となる。

30

40

#### 【0207】

なお、フォギー書き込み動作のような第1段階の書き込み動作が終了している内部バッファ(共有キャッシュ)31内のライトデータの中から破棄すべきライトデータを選択するためのアルゴリズムは、最も古いデータを選択するファースト・イン・ファースト・アウトに限定されず、LRU、ランダムのような他のアルゴリズムを使用してもよい。

#### 【0208】

プログラム/リードシーケンサ41は、各NAND型フラッシュメモリチップからステータス、つまり、書き込み完了(Done)、書き込み失敗(Error)、ブロックアドレス、ページアドレス、を受信する。そして、これらステータスに基づいて、プログラム/リードシーケンサ41は、ライトコマンド毎に、このライトコマンドに関連付けられた

50

ライトデータ全体に対する書き込み動作（NAND型フラッシュメモリチップに同じデータを1回または複数回転送する書き込み動作）の全てが終了したか否かを判定する。あるライトコマンドに関連付けられたライトデータ全体に対する書き込み動作の全てが終了したならば、プログラム/リードシーケンサ41は、このライトコマンドのコマンド完了を示すレスポンス（Done）をホスト2に送信する。このコマンド完了を示すレスポンス（Done）は、このライトコマンドを一意に識別するコマンドIDを含む。

【0209】

次に、リードコマンドの処理について説明する。

リードコマンドは、リードすべきデータが格納されているブロックを示すブロックアドレスと、このデータが格納されているページを示すページアドレスと、このデータが転送されるべきホストメモリ上のリードデータバッファ53内の位置を示すデータポイントと、このデータの長さを含む。

10

【0210】

プログラム/リードシーケンサ41は、リードコマンドによって指定されたブロックアドレスおよびページアドレスを内部バッファ（共有キャッシュ）31に送出し、リードコマンドによって指定されたデータの読み出しを内部バッファ（共有キャッシュ）31に要求する。

【0211】

内部バッファ（共有キャッシュ）31、つまりキャッシュコントローラは、NAND型フラッシュメモリチップに、このブロックアドレスと、このページアドレスと、リード指示用のNANDコマンド（フラッシュリードコマンド）を、フラッシュコマンドキュー43を介して送出する。NAND型フラッシュメモリチップから読み出されたデータは、DMAC15によってリードデータバッファ53に転送される。

20

【0212】

なお、リードコマンドによって指定されたデータが、書き込み動作が終了していないデータ、または書き込み動作の全てが終了しているがNAND型フラッシュメモリ5からまだ読み出し可能となっていないデータである場合、内部バッファ（共有キャッシュ）31、つまりキャッシュコントローラは、内部バッファ（共有キャッシュ）31にこのデータが存在するか否かを判定してもよい。内部バッファ（共有キャッシュ）31にこのデータが存在するならば、このデータが内部バッファ（共有キャッシュ）31から読み出され、そしてDMAC15によってリードデータバッファ53に転送される。

30

【0213】

一方、内部バッファ（共有キャッシュ）31にこのデータが存在しないならば、このデータは、まず、DMAC15によってライトデータバッファ51から内部バッファ（共有キャッシュ）31に転送される。そして、このデータが内部バッファ（共有キャッシュ）31から読み出され、そしてDMAC15によってリードデータバッファ53に転送される。

【0214】

図45は、フラッシュストレージデバイス3によって実行される複数段階の書き込み動作を示す。

ここでは、4つのワード線を往復する場合のフォギー・ファイン書き込み動作を例示する。また、ここでは、NAND型フラッシュメモリ5が、メモリセル当たり4ビットのデータを格納するQLC-フラッシュである場合を想定する。NAND型フラッシュメモリ5内の一つの特定の書き込み先ブロック（ここでは、書き込み先ブロックBLK#1）に対するフォギー・ファイン書き込み動作は以下のように実行される。

40

【0215】

(1) まず、4ページ（P0～P3）分のライトデータがページ単位でNAND型フラッシュメモリ5に転送され、この書き込み先ブロックBLK#1内のワード線WL0に接続された複数のメモリセルに、これら4ページ（P0～P3）分のライトデータを書き込むためのフォギー書き込み動作が実行される。

【0216】

50



(2) 次いで、次の4ページ(P4~P7)分のライトデータがこのNAND型フラッシュメモリ5にページ単位で転送され、この書き込み先ブロックBLK#1内のワード線WL1に接続された複数のメモリセルに、これら4ページ(P4~P7)分のライトデータを書き込むためのフォギー書き込み動作が実行される。

【0217】

(3) 次いで、次の4ページ(P8~P11)分のライトデータがこのNAND型フラッシュメモリ5にページ単位で転送され、この書き込み先ブロックBLK#1内のワード線WL2に接続された複数のメモリセルに、これら4ページ(P8~P11)分のライトデータを書き込むためのフォギー書き込み動作が実行される。

【0218】

(4) 次いで、次の4ページ(P12~P15)分のライトデータがこのNAND型フラッシュメモリ5にページ単位で転送され、この書き込み先ブロックBLK#1内のワード線WL3に接続された複数のメモリセルに、これら4ページ(P12~P15)分のライトデータを書き込むためのフォギー書き込み動作が実行される。

【0219】

(5) ワード線WL3に接続された複数のメモリセルに対するフォギー書き込み動作が終了すると、書き込み対象のワード線はワード線WL0に戻り、ワード線WL0に接続された複数のメモリセルに対するファイン書き込み動作の実行が可能となる。そして、ワード線WL0に対するフォギー書き込み動作で使用された4ページ(P0~P3)分のライトデータと同じ4ページ(P0~P3)分のライトデータがページ単位でNAND型フラッシュメモリ5に再び転送され、この書き込み先ブロックBLK#1内のワード線WL0に接続された複数のメモリセルに、これら4ページ(P0~P3)分のライトデータを書き込むためのファイン書き込み動作が実行される。これにより、ページP0~P3に対するフォギー・ファイン書き込み動作が終了する。

【0220】

(6) 次いで、次の4ページ(P16~P19)分のライトデータがこのNAND型フラッシュメモリ5にページ単位で転送され、この書き込み先ブロックBLK#1内のワード線WL4に接続された複数のメモリセルに、これら4ページ(P16~P19)分のライトデータを書き込むためのフォギー書き込み動作が実行される。

【0221】

(7) ワード線WL4に接続された複数のメモリセルに対するフォギー書き込み動作が終了すると、書き込み対象のワード線はワード線WL1に戻り、ワード線WL1に接続された複数のメモリセルに対するファイン書き込み動作の実行が可能となる。そして、ワード線WL1に対するフォギー書き込み動作で使用された4ページ(P4~P7)分のライトデータと同じ4ページ(P4~P7)分のライトデータがページ単位でNAND型フラッシュメモリ5に再び転送され、この書き込み先ブロックBLK#1内のワード線WL1に接続された複数のメモリセルに、これら4ページ(P4~P7)分のライトデータを書き込むためのファイン書き込み動作が実行される。これにより、ページP4~P7に対するフォギー・ファイン書き込み動作が終了する。

【0222】

(8) 次いで、次の4ページ(P20~P23)分のライトデータがこのNAND型フラッシュメモリ5にページ単位で転送され、この書き込み先ブロックBLK#1内のワード線WL5に接続された複数のメモリセルに、これら4ページ(P20~P23)分のライトデータを書き込むためのフォギー書き込み動作が実行される。

【0223】

(9) ワード線WL5に接続された複数のメモリセルに対するフォギー書き込み動作が終了すると、書き込み対象のワード線はワード線WL2に戻り、ワード線WL2に接続された複数のメモリセルに対するファイン書き込み動作の実行が可能となる。そして、ワード線WL2に対するフォギー書き込み動作で使用された4ページ(P8~P11)分のライトデータと同じ4ページ(P8~P11)分のライトデータがページ単位でNAND型フ

10

20

30

40

50

ラッシュメモリ 5 に再び転送され、この書き込み先ブロック B L K # 1 内のワード線 W L 2 に接続された複数のメモリセルに、これら 4 ページ ( P 8 ~ P 1 1 ) 分のライトデータを書き込むためのファイン書き込み動作が実行される。これにより、ページ P 8 ~ P 1 1 に対するフォギー・ファイン書き込み動作が終了する。

【 0 2 2 4 】

図 4 6 は、書き込み先ブロック B L K # 1 へのデータの書き込み順序を示す。

ここでは、図 7 と同様に、4 つのワード線を往復する場合のフォギー・ファイン書き込み動作が実行される場合を想定する。

図 4 6 の左部に示されるデータ d 0、データ d 1、データ d 2、データ d 3、データ d 4、データ d 5、データ d 6、データ d 7、...、データ d 2 5 2、データ d 2 5 3、データ d 2 5 4、データ d 2 5 5 は、書き込み先ブロック B L K # 1 を指定する複数のライトコマンドそれぞれに対応する複数のライトデータを示している。ここでは、図示の簡単化のために、全てのライトデータが同じサイズを有している場合が想定されている。

10

【 0 2 2 5 】

図 4 6 の右部は、書き込み先ブロック B L K # 1 へのデータの書き込み順序を示している。書き込み動作は、ワード線 W L 0 に接続された複数のメモリセルへのデータ d 0 の書き込み ( フォギー書き込み )、ワード線 W L 1 に接続された複数のメモリセルへのデータ d 1 の書き込み ( フォギー書き込み )、ワード線 W L 2 に接続された複数のメモリセルへのデータ d 2 の書き込み ( フォギー書き込み )、ワード線 W L 3 に接続された複数のメモリセルへのデータ d 3 の書き込み ( フォギー書き込み )、ワード線 W L 0 に接続された複数のメモリセルへのデータ d 0 の書き込み ( ファイン書き込み )、ワード線 W L 4 に接続された複数のメモリセルへのデータ d 4 の書き込み ( フォギー書き込み )、ワード線 W L 1 に接続された複数のメモリセルへのデータ d 1 の書き込み ( ファイン書き込み )、ワード線 W L 5 に接続された複数のメモリセルへのデータ d 5 の書き込み ( フォギー書き込み )、ワード線 W L 2 に接続された複数のメモリセルへのデータ d 2 の書き込み ( ファイン書き込み )、... という順序で実行される。

20

【 0 2 2 6 】

図 4 7 は、N A N D 型フラッシュメモリ 5 のデータ書き込み単位と同じサイズの単位でライトデータをホスト 2 からフラッシュストレージデバイス 3 に転送する動作を示す。

図 4 7 の左部に示されるデータ d 1、データ d 2、データ d 3、データ d 4、データ d 5、データ d 6、データ d 7、データ d 8、データ d 9、データ d 1 0、... は、書き込み先ブロック B L K # 1 を指定する 1 0 個のライトコマンドにそれぞれ対応する 1 0 個のライトデータを示している。ライトデータの長さ ( サイズ ) は、個々のライトコマンド毎に異なる。図 4 7 では、データ d 1、データ d 2、データ d 3、データ d 4 の各々が 4 K バイトのサイズを有し、データ d 5 が 8 K バイトのサイズを有し、データ d 6 が 4 0 K バイトのサイズを有し、データ d 7 が 1 6 K バイトのサイズを有し、データ d 8、データ d 9 の各々が 8 K バイトのサイズを有し、データ d 1 0 が 1 M バイトのサイズを有する場合が想定されている。

30

【 0 2 2 7 】

ホスト 2 から受信される各ライトコマンドはデータポインタ、長さ、ブロック識別子 ( 例えばブロックアドレス ) を含むので、フラッシュストレージデバイス 3 のコントローラ 4 は、ホスト 2 から受信されるライトコマンドを、複数の書き込み先ブロックにそれぞれ対応する複数のグループに分類することができる。上述のデータ d 1、データ d 2、データ d 3、データ d 4、データ d 5、データ d 6、データ d 7、データ d 8、データ d 9、データ d 1 0、... は、書き込み先ブロック B L K # 1 に対応するグループに分類された 1 0 個のライトコマンドにそれぞれ対応している。これら 1 0 個のライトコマンドは、書き込み先ブロック B L K # 1 を示すブロック識別子 ( 例えばブロックアドレス ) を含むライトコマンドである。

40

【 0 2 2 8 】

フラッシュストレージデバイス 3 のコントローラ 4 は、書き込み先ブロック B L K # 1 を

50

指定するこれらライトコマンド内のデータポインタおよび長さに基づいて、データ d 1、データ d 2、データ d 3、データ d 4、データ d 5、データ d 6、データ d 7、データ d 8、データ d 9、データ d 10 がそれぞれ存在するライトデータバッファ 5 1 上の位置、およびデータ d 1、データ d 2、データ d 3、データ d 4、データ d 5、データ d 6、データ d 7、データ d 8、データ d 9、データ d 10 それぞれの長さを管理する。そして、コントローラ 4 は、一つのライトコマンドに関連付けられた大きなサイズのライトデータを複数のライトデータ（複数のデータ部）に分割、または 2 以上のライトコマンドにそれぞれ関連付けられた小さなサイズの 2 以上のライトデータを互いに結合することによって得られる、NAND 型フラッシュメモリ 5 のデータ書き込み単位と同じサイズを有するライトデータを、ホスト 2 から取得する。

10

#### 【0229】

図 4 7 では、コントローラ 4 は、最初に、各々が 4 K バイトのサイズを有するデータ d 1、データ d 2、データ d 3、データ d 4 を互いに結合することによって得られる 16 K バイトライトデータをホスト 2 のライトデータバッファ 5 1 から取得する。この場合、コントローラ 4 は、これに限定されないが、例えば、4 回の DMA 転送によって、この 16 K バイトライトデータをホスト 2 のライトデータバッファ 5 1 から内部バッファ 3 1 に転送してもよい。最初の DMA 転送では、データ d 1 の先頭位置を指定する転送元アドレスと、データ長 = 4 K B が DMAC 1 5 にセットされてもよい。データ d 1 の先頭位置を指定する転送元アドレスは、データ d 1 に対応するライトコマンド内のデータポインタによって表される。2 回目の DMA 転送では、データ d 2 の先頭位置を指定する転送元アドレスと、データ長 = 4 K B が DMAC 1 5 にセットされてもよい。データ d 2 の先頭位置を指定する転送元アドレスは、データ d 2 に対応するライトコマンド内のデータポインタによって表される。3 回目の DMA 転送では、データ d 3 の先頭位置を指定する転送元アドレスと、データ長 = 4 K B が DMAC 1 5 にセットされてもよい。データ d 3 の先頭位置を指定する転送元アドレスは、データ d 3 に対応するライトコマンド内のデータポインタによって表される。4 回目の DMA 転送では、データ d 4 の先頭位置を指定する転送元アドレスと、データ長 = 4 K B が DMAC 1 5 にセットされてもよい。データ d 4 の先頭位置を指定する転送元アドレスは、データ d 4 に対応するライトコマンド内のデータポインタによって表される。

20

#### 【0230】

そして、コントローラ 4 は、DMA 転送によって取得されるこの 16 K バイトライトデータ（d 1、d 2、d 3、d 4）を、書き込み先ブロック BLK # 1 のページ P 0 に書き込まれるべきデータとして NAND 型フラッシュメモリ 5 に転送する。

30

コントローラ 4 は、書き込み先ブロック BLK # 1 の次の書き込み先ページをページ P 1 に変更し、8 K バイトのサイズを有するデータ d 3 と、データ d 6 内の先頭の 8 K バイトデータ d 6 - 1 とを互いに結合することによって得られる 16 K バイトライトデータをホスト 2 のライトデータバッファ 5 1 から取得する。この場合、コントローラ 4 は、これに限定されないが、例えば、2 回の DMA 転送によって、この 16 K バイトライトデータをホスト 2 のライトデータバッファ 5 1 から内部バッファ 3 1 に転送してもよい。最初の DMA 転送では、データ d 5 の先頭位置を指定する転送元アドレスと、データ長 = 8 K B が DMAC 1 5 にセットされてもよい。データ d 5 の先頭位置を指定する転送元アドレスは、データ d 5 に対応するライトコマンド内のデータポインタによって表される。2 回目の DMA 転送では、データ d 6 - 1 の先頭位置を指定する転送元アドレスと、データ長 = 8 K B が DMAC 1 5 にセットされてもよい。データ d 6 - 1 の先頭位置を指定する転送元アドレスは、データ d 6 に対応するライトコマンド内のデータポインタによって表される。

40

#### 【0231】

そして、コントローラ 4 は、この 16 K バイトライトデータ（d 5、d 6 - 1）を、書き込み先ブロック BLK # 1 のページ P 1 に書き込まれるべきデータとして NAND 型フラッシュメモリ 5 に転送する。

コントローラ 4 は、書き込み先ブロック BLK # 1 の次の書き込み先ページをページ P 2

50

に変更し、データ d 6 の残りの 3 2 K バイトデータのうちの最初の 1 6 K バイトデータ d 6 - 2 をホスト 2 のライトデータバッファ 5 1 から取得する。この場合、コントローラ 4 は、これに限定されないが、例えば、1 回の D M A 転送によって、この 1 6 K バイトライトデータをホスト 2 のライトデータバッファ 5 1 から内部バッファ 3 1 に転送してもよい。この D M A 転送では、データ d 6 - 2 の先頭位置を指定する転送元アドレスと、データ長 = 1 6 K B が D M A C 1 5 にセットされてもよい。データ d 6 - 2 の先頭位置を指定する転送元アドレスは、データ d 6 に対応するライトコマンド内のデータポインタの値に 8 K B 分のオフセットを加算することによって求めることができる。

【 0 2 3 2 】

そして、コントローラ 4 は、この 1 6 K バイトライトデータ ( d 6 - 2 ) を、書き込み先ブロック B L K # 1 のページ P 2 に書き込まれるべきデータとして N A N D 型フラッシュメモリ 5 に転送する。

10

コントローラ 4 は、書き込み先ブロック B L K # 1 の次の書き込み先ページをページ P 3 に変更し、データ d 6 の残りの 1 6 K バイトデータ d 6 - 3 をホスト 2 のライトデータバッファ 5 1 から取得する。この場合、コントローラ 4 は、これに限定されないが、例えば、1 回の D M A 転送によって、この 1 6 K バイトライトデータをホスト 2 のライトデータバッファ 5 1 から内部バッファ 3 1 に転送してもよい。この D M A 転送では、データ d 6 - 3 の先頭位置を指定する転送元アドレスと、データ長 = 1 6 K B が D M A C 1 5 にセットされてもよい。データ d 6 - 3 の先頭位置を指定する転送元アドレスは、データ d 6 に対応するライトコマンド内のデータポインタの値に 2 4 K B 分のオフセットを加算することによって求めることができる。

20

【 0 2 3 3 】

そして、コントローラ 4 は、この 1 6 K バイトライトデータ ( d 6 - 3 ) を、書き込み先ブロック B L K # 1 のページ P 3 に書き込まれるべきデータとして N A N D 型フラッシュメモリ 5 に転送する。

そして、コントローラ 4 は、フォギー書き込み動作によって 4 ページ分のデータ ( P 0 ~ P 3 ) を書き込み先ブロック B L K # 1 のワード線 W L 0 に接続された複数のメモリセルに書き込む。

【 0 2 3 4 】

コントローラ 4 は、書き込み先ブロック B L K # 1 の次の書き込み先ページをページ P 4 に変更し、1 6 K バイトのサイズを有するデータ d 7 をホスト 2 のライトデータバッファ 5 1 から取得する。この場合、コントローラ 4 は、これに限定されないが、例えば、1 回の D M A 転送によって、この 1 6 K バイトライトデータをホスト 2 のライトデータバッファ 5 1 から内部バッファ 3 1 に転送してもよい。この D M A 転送では、データ d 7 の先頭位置を指定する転送元アドレスと、データ長 = 1 6 K B が D M A C 1 5 にセットされてもよい。データ d 7 の先頭位置を指定する転送元アドレスは、データ d 7 に対応するライトコマンド内のデータポインタによって表される。

30

【 0 2 3 5 】

そして、コントローラ 4 は、この 1 6 K バイトライトデータ ( d 7 ) を、書き込み先ブロック B L K # 1 のページ P 4 に書き込まれるべきデータとして N A N D 型フラッシュメモリ 5 に転送する。

40

コントローラ 4 は、書き込み先ブロック B L K # 1 の次の書き込み先ページをページ P 5 に変更し、8 K バイトのサイズを有するデータ d 8 と、8 K バイトのサイズを有するデータ d 9 とを互いに結合することによって得られる 1 6 K バイトライトデータをホスト 2 のライトデータバッファ 5 1 から取得する。この場合、コントローラ 4 は、これに限定されないが、例えば、2 回の D M A 転送によって、この 1 6 K バイトライトデータをホスト 2 のライトデータバッファ 5 1 から内部バッファ 3 1 に転送してもよい。最初の D M A 転送では、データ d 8 の先頭位置を指定する転送元アドレスと、データ長 = 8 K B が D M A C 1 5 にセットされてもよい。データ d 8 の先頭位置を指定する転送元アドレスは、データ d 8 に対応するライトコマンド内のデータポインタによって表される。2 回目の D M A 転

50

送では、データ d 9 の先頭位置を指定する転送元アドレスと、データ長 = 8 K B が D M A C 1 5 にセットされてもよい。データ d 9 の先頭位置を指定する転送元アドレスは、データ d 9 に対応するライトコマンド内のデータポインタによって表される。

【 0 2 3 6 】

そして、コントローラ 4 は、この 1 6 K バイトライトデータ ( d 8 、 d 9 ) を、書き込み先ブロック B L K # 1 のページ P 5 に書き込まれるべきデータとして N A N D 型フラッシュメモリ 5 に転送する。

コントローラ 4 は、書き込み先ブロック B L K # 1 の次の書き込み先ページをページ P 6 に変更し、データ d 1 0 内の先頭の 1 6 K バイトデータ d 1 0 - 1 をホスト 2 のライトデータバッファ 5 1 から取得する。この場合、コントローラ 4 は、これに限定されないが、例えば、1 回の D M A 転送によって、この 1 6 K バイトライトデータをホスト 2 のライトデータバッファ 5 1 から内部バッファ 3 1 に転送してもよい。この D M A 転送では、データ d 1 0 - 1 の先頭位置を指定する転送元アドレスと、データ長 = 1 6 K B が D M A C 1 5 にセットされてもよい。データ d 1 0 - 1 の先頭位置を指定する転送元アドレスは、データ d 1 0 に対応するライトコマンド内のデータポインタによって表される。

【 0 2 3 7 】

そして、コントローラ 4 は、この 1 6 K バイトライトデータ ( d 1 0 - 1 ) を、書き込み先ブロック B L K # 1 のページ P 6 に書き込まれるべきデータとして N A N D 型フラッシュメモリ 5 に転送する。

コントローラ 4 は、書き込み先ブロック B L K # 1 の次の書き込み先ページをページ P 7 に変更し、データ d 1 0 内の次の 1 6 K バイトデータ d 1 0 - 2 をホスト 2 のライトデータバッファ 5 1 から取得する。この場合、コントローラ 4 は、これに限定されないが、例えば、1 回の D M A 転送によって、この 1 6 K バイトライトデータをホスト 2 のライトデータバッファ 5 1 から内部バッファ 3 1 に転送してもよい。この D M A 転送では、データ d 1 0 - 2 の先頭位置を指定する転送元アドレスと、データ長 = 1 6 K B が D M A C 1 5 にセットされてもよい。データ d 1 0 - 2 の先頭位置を指定する転送元アドレスは、データ d 1 0 に対応するライトコマンド内のデータポインタの値に 1 6 K B 分のオフセットを加算することによって求めることができる。

【 0 2 3 8 】

そして、コントローラ 4 は、この 1 6 K バイトライトデータ ( d 1 0 - 2 ) を、書き込み先ブロック B L K # 1 のページ P 7 に書き込まれるべきデータとして N A N D 型フラッシュメモリ 5 に転送する。

そして、コントローラ 4 は、フォギー書き込み動作によって 4 ページ分のデータ ( P 4 ~ P 7 ) を書き込み先ブロック B L K # 1 のワード線 W L 1 に接続された複数のメモリセルに書き込む。

【 0 2 3 9 】

このように、コントローラ 4 は、書き込み先ブロック B L K # 1 の書き込み動作の進行に合わせて、書き込み先ブロック B L K # 1 の書き込み先ページに転送すべき 1 6 K バイトデータをホスト 2 から取得する。

そして、ワード線 W L 3 に接続された複数のメモリセルに対するフォギー書き込み動作が終了すると、ワード線 W L 0 に接続された複数のメモリセルに対するファイン書き込み動作の実行が可能となる。コントローラ 4 は、書き込み先ブロック B L K # 1 の次の書き込み先ページをページ P 1 に変更し、上述と同様の手順で、ライトデータ ( P 0 ~ P 3 ) をページ単位で N A N D 型フラッシュメモリ 5 に再び転送し、そしてファイン書き込み動作によってこれら 4 ページ分のライトデータ ( P 0 ~ P 3 ) を書き込み先ブロック B L K # 1 のワード線 W L 0 に接続された複数のメモリセルに書き込む。

【 0 2 4 0 】

これにより、最初の 6 つのライトコマンド、つまりデータ d 1 に対応するライトコマンド、データ d 2 に対応するライトコマンド、データ d 3 に対応するライトコマンド、データ d 4 に対応するライトコマンド、データ d 5 に対応するライトコマンド、データ d 6 に対

10

20

30

40

50

応するライトコマンドの各々に関しては、各ライトコマンドに関連付けさせたライトデータ全体に対するフォギー・ファイン書き込み動作の全てが終了し、且つデータ d 1 ~ d 6 の各々は N A N D 型フラッシュメモリ 5 から読み出し可能となる。このため、コントローラ 4 は、最初の 6 つのライトコマンドにそれぞれに対応する 6 つのコマンド完了レスポンスをホスト 2 に返す。

【 0 2 4 1 】

なお、図 4 7 では、書き込み先ブロック B L K # 1 を指定するライトコマンドそれぞれに関連づけられたライトデータを、書き込み先ブロック B L K # 1 の書き込み動作の進行に合わせて 1 6 K バイトの単位でホスト 2 からフラッシュストレージデバイス 3 に転送する動作を説明したが、他の各書き込み先ブロック B L K # に関して、図 9 で説明した動作と同様の動作が実行される。

10

【 0 2 4 2 】

図 4 8 のフローチャートは、フラッシュストレージデバイス 3 によって実行されるデータ書き込み処理の手順を示す。ここでは、フラッシュストレージデバイス 3 が、タイプ # 2 - ストレージデバイスとして実現されており、かつ、ホスト 2 からのライトデータを書き込み予定のページアドレスを、N A N D 型フラッシュメモリ 5 への書き込みの終了を待機することなく、ホスト 2 に通知する仕組みを備えていることを想定する。

【 0 2 4 3 】

フラッシュストレージデバイス 3 のコントローラ 4 は、データポインタ、長さ、ブロック識別子（例えばブロックアドレス）を各々が含むライトコマンドそれぞれをホスト 2 から受信する（ステップ B 1）。

20

次いで、コントローラ 4 は、特定の書き込み先ブロックを指定する一つのライトコマンドに対応する大きなサイズのライトデータを 2 以上のデータ部に分割、またはこの特定の書き込み先ブロックを指定する 2 以上のライトコマンドに対応する 2 以上のライトデータを結合し、これによって N A N D 型フラッシュメモリ 5 の書き込み単位（データ転送サイズ）と同じサイズの単位でデータをホスト 2 からフラッシュストレージデバイス 3 に転送する（ステップ B 2）。ステップ B 2 では、図 4 7 で説明したように、例えば、小さなサイズを有する幾つかのライトデータ部を互いに結合することによって得られる一つの 1 6 K バイトデータ、あるいは大きなサイズを有するライトデータを分割することによって得られる幾つかの 1 6 K バイトデータの 하나가、ホスト 2 からフラッシュストレージデバイス 3 に転送される。フラッシュストレージデバイス 3 が内部バッファ 3 1 を含む構成であるケースにおいては、ホスト 2 からフラッシュストレージデバイス 3 に転送される各 1 6 バイトライトデータは内部バッファ 3 1 に格納される。また、ステップ B 2 では、小さなサイズを有する幾つかのライトデータ部を互いに結合するために、コントローラ 4 は、ある書き込み先ブロックを指定する識別子を有する先行するライトコマンドに関連付けられたライトデータのサイズが書き込み単位（例えば 1 6 K バイト）よりも小さい場合、この書き込み先ブロックを指定する識別子を有する後続のライトコマンドの受信を待つ。

30

【 0 2 4 4 】

コントローラ 4 は、ホスト 2 から転送された 1 6 K バイトデータについて、この 1 6 K バイトデータを特定の書き込み先ブロックに書き込む前に、この 1 6 K バイトデータの書き込み先に割り当てた特定の書き込み先ブロック内のアドレスをホスト 2 に通知する（ステップ B 3）。それから、コントローラ 4 は、ホスト 2 から転送された 1 6 K バイトデータを N A N D 型フラッシュメモリ 5 に転送し、この 1 6 K バイトデータを、この特定の書き込み先ブロック内の書き込み先に割り当てたアドレスに書き込む（ステップ B 4）。コントローラ 4 は、この書き込みが成功したか否かを判定し（ステップ B 5）、エラーの場合（ステップ B 5 : N O）、ステップ B 3 から処理を繰り返す。つまり、同一のデータについてのホスト 2 へのアドレスの通知は複数回発生し得る。成功の場合は（ステップ B 5 : Y E S）、ステップ B 6 へと進む。

40

【 0 2 4 5 】

そして、コントローラ 4 は、この特定の書き込み先ブロックを指定する一つのライトコマ

50

ンドに関連付けられたライトデータ全体に対する書き込み動作（NAND型フラッシュメモリ5に同じデータを1回または複数回転送することを伴う書き込み動作）の全てが終了したか否かを判定する（ステップB6）。

【0246】

この特定の書き込み先ブロックを指定する一つのライトコマンドに関連付けられたライトデータ全体に対する書き込み動作の全てが終了したならば、コントローラ4は、このライトコマンドのコマンド完了を示すレスポンスをホスト2に返す（ステップB7）。ライトデータが格納されるライトデータバッファ51上の領域に関する解放可能通知のホスト2への送信タイミングについては後述する。

【0247】

図49のフローチャートは、フラッシュストレージデバイス3によって実行されるデータ書き込み処理の別の手順を示す。ここでも、フラッシュストレージデバイス3が、タイプ#2-ストレージデバイスとして実現されており、かつ、ホスト2からのライトデータを書き込み予定のページアドレスを、NAND型フラッシュメモリ5への書き込みの終了を待機することなく、ホスト2に通知する仕組みを備えていることを想定する。

【0248】

フラッシュストレージデバイス3のコントローラ4は、データポインタ、長さ、ブロック識別子（例えばブロックアドレス）を各々が含むライトコマンドそれぞれをホスト2から受信する（ステップB11）。

次いで、コントローラ4は、特定の書き込み先ブロックを指定する一つのライトコマンドに対応する大きなサイズのライトデータを2以上のデータ部に分割、またはこの特定の書き込み先ブロックを指定する2以上のライトコマンドに対応する2以上のライトデータを結合し、これによってNAND型フラッシュメモリ5の書き込み単位（データ転送サイズ）と同じサイズの単位でデータをホスト2からフラッシュストレージデバイス3に転送する（ステップB12）。ステップB12では、図47で説明したように、例えば、小さなサイズを有する幾つかのライトデータ部を互いに結合することによって得られる一つの16Kバイトデータ、あるいは大きなサイズを有するライトデータを分割することによって得られる幾つかの16Kバイトデータの 하나가、ホスト2からフラッシュストレージデバイス3に転送される。フラッシュストレージデバイス3が内部バッファ31を含む構成であるケースにおいては、ホスト2からフラッシュストレージデバイス3に転送される各16Kバイトライトデータは内部バッファ31に格納される。また、ステップB2では、小さなサイズを有する幾つかのライトデータ部を互いに結合するために、コントローラ4は、ある書き込み先ブロックを指定する識別子を有する先行するライトコマンドに関連付けられたライトデータのサイズが書き込み単位（例えば16Kバイト）よりも小さい場合、この書き込み先ブロックを指定する識別子を有する後続のライトコマンドの受信を待つ。

【0249】

コントローラ4は、ホスト2から転送された16KバイトデータをNAND型フラッシュメモリ5に転送し、この16Kバイトデータを、この特定の書き込み先ブロックに書き込む（ステップB13）。

コントローラ4は、この書き込みが成功したか否かを判定し（ステップB14）、エラーの場合（ステップB14：NO）、ステップB13から処理を繰り返す。成功の場合は（ステップB14：YES）、ステップB15へと進む。

【0250】

コントローラ4は、ホスト2から転送された16Kバイトデータの書き込み先に割り当てた特定の書き込み先ブロック内のアドレスをホスト2に通知する（ステップB15）。

そして、コントローラ4は、この特定の書き込み先ブロックを指定する一つのライトコマンドに関連付けられたライトデータ全体に対する書き込み動作（NAND型フラッシュメモリ5に同じデータを1回または複数回転送することを伴う書き込み動作）の全てが終了したか否かを判定する（ステップB16）。

【0251】

10

20

30

40

50

この特定の書き込み先ブロックを指定する一つのライトコマンドに関連付けられたライトデータ全体に対する書き込み動作の全てが終了したならば、コントローラ 4 は、このライトコマンドのコマンド完了を示すレスポンスをホスト 2 に返す（ステップ B 17）。ライトデータが格納されるライトデータバッファ 5 1 上の領域に関する解放可能通知のホスト 2 への送信タイミングについては後述する。

【0252】

図 5 0 のフローチャートは、フラッシュストレージデバイス 3 によって実行されるホスト 2 への解放可能通知の送信処理の手順を示す。

コントローラ 4 は、第 1 に、特定の書き込み先ブロックを指定する一つのライトコマンドに関連付けられたライトデータ全体に対する書き込み動作の全てが終了したか否かを判定する（ステップ C 1）。また、コントローラ 4 は、第 2 に、このライトデータを対象とするリードコマンドが存在するか否かを判定する（ステップ C 2）。なお、ステップ C 1 の処理とステップ C 2 の処理とは、並列的に実行される。そして、コントローラ 4 は、特定の書き込み先ブロックを指定する一つのライトコマンドに関連付けられたライトデータ全体に対する書き込み動作の全てが終了し（ステップ C 1：YES）、かつ、このライトデータを対象とするリードコマンドが存在しない場合に（ステップ C 2：NO）、このライトデータが格納されるライトデータバッファ 5 1 上の領域に関する解放可能通知をホスト 2 に送信する（ステップ C 3）。

【0253】

図 5 1 のフローチャートは、ホスト 2 によって実行されるライトデータ破棄処理の手順を示す。

ホスト 2 は、ライトコマンドのコマンド完了を示すレスポンスをフラッシュストレージデバイス 3 から受信したか否かを判定する（ステップ D 1）。あるライトコマンドのコマンド完了を示すレスポンスをフラッシュストレージデバイス 3 から受信した場合（ステップ D 1：YES）、ホスト 2 は、さらに、このライトコマンドに関連付けられたライトデータに関する解放可能通知をフラッシュストレージデバイス 3 から受信したか否かを判定する（ステップ D 2）。このライトデータに関する解放可能通知をフラッシュストレージデバイス 3 から受信した場合（ステップ D 2：YES）、ホスト 2 は、このライトコマンドに関連付けられたライトデータをライトデータバッファ 5 1 から破棄する（ステップ D 3）。

【0254】

図 5 2 は、ある書き込み先ブロックを指定する最後のライトコマンドが受信されてから閾期間、この書き込み先ブロックを指定する次のライトコマンドが受信されない場合に、フラッシュストレージデバイス 3 によって実行されるダミーデータ書き込み処理を示す。

【0255】

図 5 2 の左部に示されるデータ d 1、データ d 2、データ d 3、データ d 4 は、書き込み先ブロック B L K # 1 を指定する 4 個のライトコマンドそれぞれに対応する 4 個のライトデータを示している。図 5 2 では、データ d 1、データ d 2、データ d 3、データ d 4 の各々が 4 K バイトのサイズを有する場合が想定されている。

【0256】

(1) コントローラ 4 は、データ d 1、データ d 2、データ d 3、データ d 4 を互いに結合することによって得られる 16 K バイトライトデータをホスト 2 のライトデータバッファ 5 1 から取得する。そして、コントローラ 4 は、この 16 K バイトライトデータを、書き込み先ブロック B L K # 1 のページ P 0 に書き込まれるべきデータとして NAND 型フラッシュメモリ 5 に転送する。書き込み先ブロック B L K # 1 を指定する最後のライトコマンド、つまりデータ d 4 の書き込みを要求したライトコマンド、が受信されてから閾期間、書き込み先ブロック B L K # 1 を指定する後続のライトコマンドが受信されない場合、コントローラ 4 は、最後のライトコマンドのコマンド完了を示すレスポンスを所定時間内にホスト 2 に返すことを可能にするために、書き込み先ブロック B L K # 1 内の 1 以上のページにダミーデータを書き込み、次のライトデータが書き込まれるべき書き込み先ブ

10

20

30

40

50



ロック B L K # 1 内の書き込み先ページの位置を進める。例えば、コントローラ 4 は、ページ P 1 ~ P 3 に対応する 3 ページ分のダミーデータをページ単位で N A N D 型フラッシュメモリ 5 に転送し、フォギー書き込み動作によって 4 ページ分のデータ ( P 0 ~ P 3 ) を書き込み先ブロック B L K # 1 のワード線 W L 0 に接続された複数のメモリセルに書き込む。

【 0 2 5 7 】

( 2 ) 次いで、コントローラ 4 は、ページ P 4 ~ P 7 に対応する 4 ページ分のダミーデータをページ単位で N A N D 型フラッシュメモリ 5 に転送し、フォギー書き込み動作によって 4 ページ分のデータ ( P 4 ~ P 7 ) を書き込み先ブロック B L K # 1 のワード線 W L 1 に接続された複数のメモリセルに書き込む。

10

【 0 2 5 8 】

( 3 ) 次いで、コントローラ 4 は、ページ P 8 ~ P 1 1 に対応する 4 ページ分のダミーデータをページ単位で N A N D 型フラッシュメモリ 5 に転送し、フォギー書き込み動作によって 4 ページ分のデータ ( P 8 ~ P 1 1 ) を書き込み先ブロック B L K # 1 のワード線 W L 2 に接続された複数のメモリセルに書き込む。

【 0 2 5 9 】

( 4 ) 次いで、コントローラ 4 は、ページ P 1 2 ~ P 1 5 に対応する 4 ページ分のダミーデータをページ単位で N A N D 型フラッシュメモリ 5 に転送し、フォギー書き込み動作によって 4 ページ分のデータ ( P 1 2 ~ P 1 5 ) を書き込み先ブロック B L K # 1 のワード線 W L 3 に接続された複数のメモリセルに書き込む。

20

【 0 2 6 0 】

( 5 ) 次いで、コントローラ 4 は、データ d 1、データ d 2、データ d 3、データ d 4 を互いに結合することによって得られる 1 6 K バイトライトデータをライトデータバッファ 5 1 または内部バッファ 3 1 から N A N D 型フラッシュメモリ 5 に転送し、さらに、W L 0 のフォギー書き込み動作で使用した 3 ページ分のダミーデータ ( P 0 ~ P 3 ) と同じ 3 ページ分のダミーデータ ( P 0 ~ P 3 ) をページ単位で N A N D 型フラッシュメモリ 5 に転送する。そして、コントローラ 4 は、ファイン書き込み動作によって 4 ページ分のデータ ( P 0 ~ P 3 ) を書き込み先ブロック B L K # 1 のワード線 W L 0 に接続された複数のメモリセルに書き込む。これにより、データ d 1、データ d 2、データ d 3、データ d 4 の複数段階の書き込み動作が全て完了し、データ d 1、データ d 2、データ d 3、データ d 4 が N A N D 型フラッシュメモリ 5 読み出し可能となる。コントローラ 4 は、データ d 1 の書き込みを要求した最初のライトコマンドのコマンド完了を示すレスポンスと、データ d 2 の書き込みを要求した 2 番目のライトコマンドのコマンド完了を示すレスポンスと、データ d 3 の書き込みを要求した 3 番目のライトコマンドのコマンド完了を示すレスポンスと、データ d 4 の書き込みを要求した 4 番目のライトコマンドのコマンド完了を示すレスポンスとをホスト 2 に返す。

30

【 0 2 6 1 】

本実施形態では、N A N D 型フラッシュメモリ 5 のデータ書き込み単位と同じデータサイズの単位でライトデータがホスト 2 からフラッシュストレージデバイス 3 に転送され、あるライトコマンドのライトデータ全体の書き込み動作の全てが終了した時点で、またはこのライトデータ全体の書き込み動作の全てが終了し且つこのライトデータ全体が読み出し可能となった時点で、このライトコマンドのコマンド完了を示すレスポンスがホスト 2 に返される。このため、例えば、小さなライトデータのある書き込み先ブロックに書き込むことを要求するライトコマンドがホスト 2 からフラッシュストレージデバイス 3 に発行された後にしばらくの間、この書き込み先ブロックを指定する後続のライトコマンドがホスト 2 から発行されない場合には、このライトコマンドのタイムアウトエラーが起こる可能性がある。本実施形態では、コントローラ 4 は、あるブロック識別子を有する最後のライトコマンドがホスト 2 から受信されてから闕期間このブロック識別子を有する次のライトコマンドが受信されない場合、ダミーデータを、このブロック識別子に対応する書き込み先ブロック内の次の 1 以上の未書き込みページに書き込む。したがって、必要に応じて、

40

50

この書き込み先ブロックの書き込み動作を進行させることができるので、ライトコマンドのタイムアウトエラーが起こることを防止することができる。

【0262】

図53のフローチャートは、フラッシュストレージデバイス3によって実行されるダミーデータ書き込み処理の手順を示す。ここでは、フォギー・ファイン書き込み動作のような複数段階の書き込み動作によってデータが書き込み先ブロックに書き込まれる場合を想定する。

【0263】

フラッシュストレージデバイス3のコントローラ4は、ある書き込み先ブロックを指定する最後のライトコマンドに関連付けられたライトデータをフォギー書き込み動作のような第1段階の書き込み動作によってこの書き込み先ブロックに書き込む。この最後のライトコマンドの受信から闕期間(Th)、この書き込み先ブロックを指定する次のライトコマンドが受信されない場合(ステップS31のYES)、コントローラ4は、最後のライトコマンドに関連付けられたライトデータが書き込まれた書き込み先ブロック内のページに後続する1以上のページにダミーデータを書き込み、これによって、次のライトデータが書き込まれるべきこの書き込み先ブロック内の書き込み先ページの位置を進める(ステップS32)。この書き込み先ブロックへのダミーデータの書き込みによって書き込み先ページの位置が進み、これによって最後のライトコマンドに関連付けられたライトデータのファイン書き込み動作(第2段階の書き込み動作)が実行可能となると、コントローラ4は、最後のライトコマンドに関連付けられたライトデータをライトデータバッファ51または内部バッファ(共有キャッシュ)31からNAND型フラッシュメモリ5に再び転送し、このライトデータのファイン書き込み動作を実行する(ステップS33)。

【0264】

最後のライトコマンドに関連付けられたライトデータのファイン書き込み動作が終了すると、つまりこのライトデータ全体の複数段階の書き込み動作の全てが終了すると、コントローラ4は、この最後のライトコマンドのコマンド完了を示すレスポンスをホスト2に返す(ステップS34)。

【0265】

このように、複数段階の書き込み動作によってライトデータを書き込み先ブロックに書き込むケースにおいては、コントローラ4は、最後のライトコマンドに関連付けられたライトデータの第2段階の書き込み動作が実行可能になるように、ダミーデータをこの書き込み先ブロック内の1以上のページに書き込み、次のライトデータが書き込まれるべきこの書き込み先ブロック内の書き込み先ページの位置を進める。

【0266】

図54は、内部バッファ(共有キャッシュ)31を使用してコントローラ4によって実行されるデータ転送動作を示す。

内部バッファ(共有キャッシュ)31は、複数の書き込み先ブロックBLK#1、BLK#2、...、BLK#nによって共有される。フラッシュストレージデバイス3のコントローラ4は、書き込み先ブロックBLK#1、BLK#2、...、BLK#nの各々について以下の処理を実行する。

【0267】

以下では、書き込み先ブロックBLK#1を例示して説明する。

コントローラ4が、書き込み先ブロックBLK#1を指定する1以上のライトコマンドを受信した後、コントローラ4は、書き込み先ブロックBLK#1を指定する一つのライトコマンドに関連付けられたライトデータを複数のライトデータに分割、または書き込み先ブロックBLK#1を指定する2以上のライトコマンドにそれぞれ関連付けられたライトデータを互いに結合することによって得られる、NAND型フラッシュメモリ5の書き込み単位と同じサイズを有するライトデータを、ライトデータバッファ51から取得する。そして、コントローラ4は、ライトデータバッファ51から取得される、各々がNAND型フラッシュメモリ5の書き込み単位と同じサイズを有する複数のライトデータを内部バ

10

20

30

40

50

ッファ（共有キャッシュ）31に格納する。

【0268】

ライトデータバッファ51は、必ずしも、ホストメモリ上の連続する一つの領域から構成される必要は無く、図54に示されているように、複数のライトデータバッファ51-1、51-2、...、51-nによって実現されてもよい。

コントローラ4は、書き込み先ブロックBLK#1に次に書き込むべきライトデータ（第1のライトデータ）を内部バッファ（共有キャッシュ）31から取得し、第1のライトデータをNAND型フラッシュメモリ5に転送し、フォギー書き込み動作のような第1段階の書き込み動作によってこのライトデータを書き込み先ブロックBLK#1に書き込む。

【0269】

内部バッファ（共有キャッシュ）31にホスト2からのライトデータを効率よく蓄積できるようにするため、ホスト2から取得されるライトデータを格納するための空き領域が内部バッファ（共有キャッシュ）31に無い場合には、コントローラ4は、フォギー書き込み動作のような第1段階の書き込み動作が終了している内部バッファ（共有キャッシュ）31内のライトデータ（フォギーステートのライトデータ）を破棄して、空き領域を内部バッファ（共有キャッシュ）31に確保する。

【0270】

例えば、内部バッファ（共有キャッシュ）31に空き領域がない状態でホスト2から任意の書き込み先ブロックを指定する新たなライトコマンドを受信した場合には、コントローラ4は、フォギー書き込み動作のような第1段階の書き込み動作が終了している内部バッファ（共有キャッシュ）31内のライトデータ（フォギーステートのライトデータ）を破棄して、新たなライトコマンドに対応するライトデータを格納可能な空き領域を内部バッファ（共有キャッシュ）31に確保してもよい。

【0271】

例えば、内部バッファ（共有キャッシュ）31全体が多数のフォギーステートのライトデータで満たされている状態でホスト2から新たなライトコマンドを受信した場合には、コントローラ4は、これらフォギーステートのライトデータの中から破棄すべき特定のライトデータを選択してもよく、この選択したライトデータを破棄してもよい。これにより、制限された容量を有する内部バッファ（共有キャッシュ）31を、複数の書き込み先ブロック間で効率よく共有することができる。

【0272】

コントローラ4は、第1のライトデータのファイン書き込み動作のような第2段階の書き込み動作を実行すべき時点において第1のライトデータが内部バッファ（共有キャッシュ）31に存在しない場合には、第1のライトデータを取得するための要求（転送要求：DMA転送要求）をホスト2に送信することによってホスト2のライトデータバッファ51から再び取得する。この取得された第1のライトデータは内部バッファ（共有キャッシュ）31に格納されてもよい。そして、コントローラ4は、取得された第1のライトデータをNAND型フラッシュメモリ5に転送し、ファイン書き込み動作のような第2段階の書き込み動作によってこの第1のライトデータを書き込み先ブロックBLK#1に書き込む。

【0273】

第1のライトデータのファイン書き込み動作のような第2段階の書き込み動作を実行すべき時点において第1のライトデータが内部バッファ（共有キャッシュ）31に存在している場合には、コントローラ4は、この内部バッファ（共有キャッシュ）31から第1のライトデータを取得し、取得された第1のライトデータをNAND型フラッシュメモリ5に転送し、ファイン書き込み動作のような第2段階の書き込み動作によってこの第1のライトデータを書き込み先ブロックBLK#1に書き込む。

【0274】

NAND型フラッシュメモリ5への第1のライトデータの最終回のデータ転送（ここでは、ファイン書き込み動作のためのデータ転送）を行った後、コントローラ4は、この第1のライトデータを内部バッファ（共有キャッシュ）31から破棄することによって内部バ

10

20

30

40

50

バッファ（共有キャッシュ）31に空き領域を確保する。あるいは、コントローラ4はこの第1のライトデータのファイン書き込み動作が終了した場合に、この第1のライトデータを内部バッファ（共有キャッシュ）31から破棄してもよい。

【0275】

さらに、コントローラ4は、あるライトコマンドに関連付けられたライトデータ全体のファイン書き込み動作が終了した場合、あるいはこのライトデータ全体のファイン書き込み動作が終了し且つこのライトデータ全体がNAND型フラッシュメモリ5から読み出し可能となった場合に、このライトコマンドのコマンド完了を示すレスポンスをホスト2に返す。

【0276】

内部バッファ（共有キャッシュ）31はある限られた容量を有しているが、書き込み先ブロックの数がある一定数以下であるならば、第2段階の書き込み動作を実行すべき時点において第1のライトデータが内部バッファ（共有キャッシュ）31に存在する確率（ヒット率）は比較的高い。したがって、同じライトデータをホスト2からフラッシュストレージデバイス3に複数回転送すること無く、フォギー・ファイン書き込み動作のような複数段階の書き込み動作を実行することができる。これにより、ホスト2とフラッシュストレージデバイス3との間のデータトラフィックを削減できるので、データ書き込みの度に同じライトデータをホスト2からフラッシュストレージデバイス3に複数回転送する場合に比し、フラッシュストレージデバイス3のI/O性能を向上することができる。

【0277】

書き込み先ブロックの数は、ホスト2を利用するクライアントの数と同数であってよい。この場合、あるクライアントに対応するデータはこのクライアントに対応する書き込み先ブロックに書き込まれ、他のクライアントに対応するデータは別の書き込み先ブロックに書き込まれる。したがって、ホスト2を利用するクライアントの数が増えるにつれて、内部バッファ（共有キャッシュ）31のヒット率は低下する。しかし、第1のライトデータが内部バッファ（共有キャッシュ）31に存在しない場合（ミス）には、コントローラ4はこの第1のライトデータをホスト2から取得する。したがって、クライアントの数が増加しても、フォギー・ファイン書き込み動作のような複数段階の書き込み動作を正常に行うことができる。

【0278】

よって、フラッシュストレージデバイス3は、フラッシュストレージデバイス3を共有するクライアントの数の増加（つまり、同時に利用可能な書き込み先ブロックの数の増加）に柔軟に対応でき且つホスト2とフラッシュストレージデバイス3との間のデータトラフィックを削減することができる。

【0279】

ここでは、書き込み先ブロックBLK#1にデータを書き込むための書き込み処理について説明したが、他の全ての書き込み先ブロックの各々に対しても同様の書き込み処理が実行される。

図55のフローチャートは、内部バッファ（共有キャッシュ）31を使用してコントローラ4によって実行されるデータ書き込み処理の手順を示す。

【0280】

コントローラ4は、データポインタと、ライトデータの長さ、複数の書き込み先ブロックのいずれか一つを指定する識別子（例えばブロックアドレス）とを各々が含む1つ以上のライトコマンドをホスト2から受信する（ステップS101）。同じ書き込み先ブロックを示す識別子を有する一つ以上のライトコマンドを受信した後、コントローラ4は、これらライトコマンド内の一つのライトコマンドに関連付けられたライトデータを複数のライトデータに分割、または同じ書き込み先ブロックを示す識別子を有する2以上のライトコマンドにそれぞれ関連付けられたライトデータを互いに結合することによって得られる、NAND型フラッシュメモリ5の書き込み単位と同じサイズを有するライトデータを、ライトデータバッファ51から内部バッファ（共有キャッシュ）31に転送する（ステッ

10

20

30

40

50

プ S 1 0 2 )。

【 0 2 8 1 】

コントローラ 4 は、この書き込み先ブロックに次に書き込むべきライトデータを内部バッファ (共有キャッシュ) 3 1 から取得し、このライトデータを NAND 型フラッシュメモリ 5 に転送し、フォギー書き込み動作によってこの第 1 のライトデータをこの書き込み先ブロックに書き込む (ステップ S 1 0 3、S 1 0 4)。NAND 型フラッシュメモリ 5 が QLC - フラッシュとして実現されている場合には、ステップ S 1 0 3 では、4 ページ分のライトデータがページ単位で NAND 型フラッシュメモリ 5 に転送され、ステップ S 1 0 4 では、4 ページ分のライトデータがフォギー書き込み動作によってこの書き込み先ブロック内の書き込み対象の一つのワード線に接続された複数のメモリセルに書き込まれる。

10

【 0 2 8 2 】

なお、ライトデータバッファ 5 1 から内部バッファ (共有キャッシュ) 3 1 へのライトデータの転送は、各書き込み先ブロックの書き込み動作の進行に合わせて実行される。例えば、ある書き込み先ブロックのあるページに書き込むべきライトデータを NAND 型フラッシュメモリチップに転送する動作が終了した場合に、この書き込み先ブロックの次のページに書き込むべきライトデータがライトデータバッファ 5 1 から内部バッファ (共有キャッシュ) 3 1 へ転送されてもよい。あるいは、ある書き込み先ブロックのあるページに書き込むべきライトデータをこの書き込み先ブロックを含む NAND 型フラッシュメモリチップに転送する動作が終了し、且つこのライトデータをこの書き込み先ブロックに書き込む動作が終了した場合に、この書き込み先ブロックの次のページに書き込むべきライトデータがライトデータバッファ 5 1 から内部バッファ (共有キャッシュ) 3 1 へ転送されてもよい。

20

【 0 2 8 3 】

フォギー書き込み動作が行われたこのライトデータのファイン書き込み動作を開始すべき時点において、コントローラ 4 は、このライトデータが、内部バッファ (共有キャッシュ) 3 1 に存在するか否かを判定する。

このライトデータが内部バッファ (共有キャッシュ) 3 1 に存在するならば (ステップ S 1 0 6 の YES)、コントローラ 4 は、このライトデータを内部バッファ (共有キャッシュ) 3 1 から取得し、このライトデータを NAND 型フラッシュメモリ 5 に転送し、ファイン書き込み動作によってこのライトデータをこの書き込み先ブロックに書き込む (ステップ S 1 0 7、S 1 0 8、S 1 0 9)。これにより、このライトデータは NAND 型フラッシュメモリ 5 から読み出し可能となる。

30

【 0 2 8 4 】

コントローラ 4 は、ライトコマンド毎に、そのライトデータ全体のフォギー・ファイン書き込み動作が終了し且つそのライトデータ全体が NAND 型フラッシュメモリ 5 から読み出し可能となったか否かを判定する。そして、コントローラ 4 は、フォギー・ファイン書き込み動作が終了し且つ NAND 型フラッシュメモリ 5 から読み出し可能となったライトデータに対応するライトコマンドのコマンド完了を示すレスポンスをホスト 2 に返す (ステップ S 1 1 0)。もしステップ S 1 0 9 の処理によってあるライトコマンドに関連付けられたライトデータ全体のファイン書き込み動作が終了したならば、ステップ S 1 1 0 では、このライトコマンドのコマンド完了を示すレスポンスがホスト 2 に返されてもよい。

40

【 0 2 8 5 】

図 5 6 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データ読み出し処理の手順を示す。ここでは、フラッシュストレージデバイス 3 が、ホスト 2 のライトデータバッファ 5 1 内に存在するライトデータ中のデータが対象となるリードコマンドがホスト 2 から受信された場合、そのリードコマンドが終了するまで、そのデータが格納される領域についての解放可能通知をホスト 2 へ送信しないようにする仕組みを備えていることを想定する。

【 0 2 8 6 】

上述したように、コントローラ 4 は、ホスト 2 から受信されるリードコマンドによって指

50

定されたデータが、その書き込み動作（NAND型フラッシュメモリ5に同じデータを1回または複数回転送する書き込み動作）の全てが終了していないデータ、またはその書き込み動作の全てが終了しているがNAND型フラッシュメモリ5からまだ読み出し可能となっていないデータである場合、このデータが内部バッファ（共有キャッシュ）31に存在するか否かを判定する。このデータが内部バッファ（共有キャッシュ）31に存在しない場合、コントローラ4は、このデータをライトデータバッファ51から取得し、このデータを内部バッファ（共有キャッシュ）31に格納し、このデータを内部バッファ（共有キャッシュ）31からホスト2に返す。

【0287】

具体的には、以下のデータ読み出し処理が実行される。

10

コントローラ4がホスト2からリードコマンドを受信した場合（ステップE1：YES）、コントローラ4は、このリードコマンドによって指定されたデータが、その書き込み動作の全てが終了し、且つNAND型フラッシュメモリ5から読み出し可能なデータであるか否かを判定する（ステップE2）。

【0288】

このデータがNAND型フラッシュメモリ5から読み出し可能であるならば（ステップE2：YES）、コントローラ4は、このデータをNAND型フラッシュメモリ5から読み出し、読み出されたデータをホスト2に返す（ステップE3）。ステップE3では、コントローラ4は、この読み出されたデータを、リードコマンドに含まれるデータポイントによって指定されるリードデータバッファ53内の位置に転送する。

20

【0289】

このデータがNAND型フラッシュメモリ5から読み出し可能でないならば（ステップE2：NO）、コントローラ4は、まず、ライトデータバッファ51上のデータが破棄されないように、解放可能通知のホスト2への送信が禁止される状態に設定する（ステップE5）。そして、コントローラ4は、このデータが内部バッファ（共有キャッシュ）31に存在するか否かを判定する（ステップE5）。

【0290】

このデータが内部バッファ（共有キャッシュ）31に存在するならば（ステップE5：YES）、コントローラ4は、このデータを内部バッファ（共有キャッシュ）31から読み出し、この読み出したデータをホスト2に返す（ステップE6）。

30

ステップE6では、コントローラ4は、この読み出されたデータを、リードコマンドに含まれるデータポイントによって指定されるリードデータバッファ53内の位置に転送する。

【0291】

このデータが内部バッファ（共有キャッシュ）31に存在しないならば（ステップE5：NO）、コントローラ4は、このデータをライトデータバッファ51から取得し、内部バッファ（共有キャッシュ）31に格納する（ステップE7）。ステップE7では、DMA C15によってこのデータがライトデータバッファ51から内部バッファ（共有キャッシュ）31の空き領域に転送される。内部バッファ（共有キャッシュ）31の空き領域が無い場合には、内部バッファ（共有キャッシュ）31の空き領域を確保する処理が実行される。そして、コントローラ4は、このデータを内部バッファ（共有キャッシュ）31から読み出し、この読み出したデータをホスト2に返す（ステップE6）。ステップE6では、コントローラ4は、この読み出されたデータを、リードコマンドに含まれるデータポイントによって指定されるリードデータバッファ53内の位置に転送する。そして、コントローラ4は、ステップE4で設定した、解放可能通知のホスト2への送信が禁止される状態を解除する（ステップE8）。

40

【0292】

図57は、フラッシュストレージデバイス3に適用されるブロックリユースコマンドを示す。

ブロックリユースコマンドは、たとえば無効データや不要データのみが格納されているなどの理由で不要となった割り当て済みのブロックをフリーブロックに戻すことをフラッシュ

50

ュストレージデバイス3に対して要求するコマンド(ブロック解放要求)である。ブロックリユースコマンドは、QoSドメインを指定するQoSドメインIDと、フリーブロック化(解放)するブロックを指定するブロックアドレスとを含む。

【0293】

また、図58は、フラッシュストレージデバイス3に適用されるライトコマンドの別の例を示す。具体的には、図8に示したライトコマンドが、フラッシュストレージデバイス3がタイプ#1-ストレージデバイスとして実現されている場合において適用されるものであるのに対して、図58に示すライトコマンドは、フラッシュストレージデバイス3がタイプ#2-ストレージデバイスとして実現されている場合において適用されるものである。図58中に施されているハッチングは、図8との相違点を示している。

10

【0294】

ライトコマンドは、フラッシュストレージデバイス3にデータの書き込みを要求するコマンドである。このライトコマンドは、コマンドID、QoSドメインID、論理アドレス、長さ、等を含んでもよい。

コマンドIDはこのコマンドがライトコマンドであることを示すID(コマンドコード)であり、ライトコマンドにはライトコマンド用のコマンドIDが含まれる。

【0295】

QoSドメインIDは、データが書き込まれるべきQoSドメインを一意に識別可能な識別子である。あるエンドユーザからのライト要求に応じてホスト2から送信されるライトコマンドは、このエンドユーザに対応するQoSドメインを指定するQoSドメインIDを含んでもよい。ネームスペースIDがQoSドメインIDとして扱われてもよい。

20

【0296】

論理アドレスは、書き込まれるべきライトデータを識別するための識別子である。この論理アドレスは、上述したように、LBAであってもよいし、キー・バリュー・ストアのキーであってもよい。論理アドレスがLBAである場合には、このライトコマンドに含まれる論理アドレス(開始LBA)は、ライトデータが書き込まれるべき論理位置(最初の論理位置)を示す。

【0297】

長さは、書き込まれるべきライトデータの長さを示す。この長さ(データ長)は、粒度(Grain)の数によって指定されてもよいし、LBAの数によって指定されてもよいし、あるいはそのサイズがバイトによって指定されてもよい。

30

上述したように、コントローラ4は、NAND型フラッシュメモリ5内の多数のブロックの各々が一つのグループのみに属するようにNAND型フラッシュメモリ5内の多数のブロックを複数のグループ(複数のQoSドメイン)に分類することができる。そして、コントローラ4は、グループ(QoSドメイン)毎に、フリーブロックリスト(フリーブロックプール)とアクティブブロックリスト(アクティブブロックプール)とを管理することができる。

【0298】

各ブロックの状態は、有効データを格納しているアクティブブロックと、有効データを格納していないフリーブロックとに大別される。アクティブブロックである各ブロックは、アクティブブロックリストによって管理される。一方、フリーブロックである各ブロックは、フリーブロックリストによって管理される。

40

【0299】

ホスト2からライトコマンドを受信した時、コントローラ4は、ホスト2からのデータが書き込まれるべきブロック(書き込み先ブロック)およびこの書き込み先ブロック内の位置(書き込み先位置)を決定する。コントローラ4は、QoSドメインIDに対応するQoSドメインに属するフリーブロック群の一つを書き込み先ブロックとして決定してもよい。書き込み先位置は、ページ書き込み順序の制約およびパッドページ等を考慮して決定される。そして、コントローラ4は、ホスト2からのデータを、書き込み先ブロック内の書き込み先位置に書き込む。

50

## 【0300】

なお、この書き込み先ブロック全体がユーザデータで満たされたならば、コントローラ4は、この書き込み先ブロックをアクティブブロックリスト（アクティブブロックプール）に移動する。そして、コントローラ4は、このQoSドメインに対応するフリーブロックリストからフリーブロックを再び選択し、この選択した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。

## 【0301】

もしフリーブロックリストによって管理されている残りフリーブロックの数が所定のポリシーによって定められる閾値以下に低下した場合あるいはホスト2からガベージコレクションを実施する指示があった場合、コントローラ4は、このQoSドメインのガベージコレクションを開始してもよい。

10

## 【0302】

このQoSドメインのガベージコレクションでは、コントローラ4は、このQoSドメインに対応するアクティブブロック群からコピー元ブロック（GCソースブロック）とコピー先ブロック（GCデスティネーションブロック）を選択する。どのブロックをGC候補（コピー元ブロック）として選択するかは、ホスト2によって指定される上述のポリシーに従って決定されてもよいし、ホスト2から指定されても良い。ポリシーも基づく場合には例えば、有効データ量が最も少ないブロックがGC候補（コピー元ブロック）として選択されてもよい。

## 【0303】

図59は、図58のライトコマンドに対するレスポンスを示す。なお、図59中に施されているハッチングも、図9との相違点を示している。

20

このレスポンスは、論理アドレス、物理アドレス、長さを含む。

論理アドレスは、図7のライトコマンドに含まれていた論理アドレスである。

## 【0304】

物理アドレスは、図7のライトコマンドに応じてデータが書き込まれたNAND型フラッシュメモリ5内の物理記憶位置を示す。本実施形態では、この物理アドレスは、ブロック番号とページ番号との組み合わせではなく、上述したように、ブロック番号とオフセット（ブロック内オフセット）との組み合わせによって指定される。ブロック番号は、フラッシュストレージデバイス3内の全てのブロックの任意の一つを一意に識別可能な識別子である。全てのブロックに異なるブロック番号が付与されている場合には、これらブロック番号を直接使用してもよい。あるいは、ブロック番号は、ダイ番号と、ダイ内ブロック番号との組み合わせによって表現されてもよい。長さは、書き込まれるべきライトデータの長さを示す。この長さ（データ長）は、粒度（Grain）の数によって指定されてもよいし、LBAの数によって指定されてもよいし、あるいはそのサイズがバイトによって指定されてもよい。

30

## 【0305】

図60は、ホスト2とフラッシュストレージデバイス3とによって実行される書き込み動作処理のシーケンスの別の例を示す。具体的には、図25に示したシーケンスが、フラッシュストレージデバイス3がタイプ#1-ストレージデバイスとして実現されている場合におけるものであるのに対して、図60に示すシーケンスは、フラッシュストレージデバイス3がタイプ#2-ストレージデバイスとして実現されている場合におけるものである。

40

## 【0306】

ホスト2は、QoSドメインID、LBA、長さを含むライトコマンドをフラッシュストレージデバイス3に送信する。フラッシュストレージデバイス3のコントローラ4がこのライトコマンドを受信した時、コントローラ4は、ホスト2からのライトデータを書き込むべき書き込み先ブロックおよびこの書き込み先ブロック内の位置を決定する。より詳しくは、コントローラ4は、フリーブロックリストから一つのフリーブロックを選択し、選択した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（ステップS11）。つまり、この選択されたフリーブロックおよびこの選択されたフリーブロック内の利用可能な

50



最初のページが、ホスト 2 からのライトデータを書き込むべき書き込み先ブロックおよびこの書き込み先ブロック内の位置として決定される。もし書き込み先ブロックがすでに割り当てられている場合には、このステップ 1 2 における書き込み先ブロック割り当て処理を実行する必要は無い。すでに割り当てられている書き込み先ブロック内の利用可能な次のページが、ホスト 2 からのライトデータを書き込むべき書き込み先ブロック内の位置として決定される。

【0307】

コントローラ 4 は、複数の QoS ドメインに対応する複数のフリーブロックリストを管理してもよい。ある QoS ドメインに対応するフリーブロックリストにおいては、この QoS ドメインに対して予約されたブロック群のみが登録されてもよい。この場合、ステップ S 1 2 では、コントローラ 4 は、ライトコマンドの QoS ドメイン ID によって指定される QoS ドメインに対応するフリーブロックリストを選択し、この選択したフリーブロックリストから一つのフリーブロックを選択し、この選択したフリーブロックを書き込み先ブロックとして割り当ててもよい。これにより、異なる QoS ドメインに対応するデータが同じブロックに混在されてしまうことを防止することができる。

10

【0308】

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書き込み先ブロックに書き込む（ステップ S 1 2）。ステップ S 1 2 では、コントローラ 4 は、論理アドレス（ここでは LBA）とライトデータの双方を書き込み先ブロックに書き込む。

コントローラ 4 は、ブロック管理テーブル 3 2 を更新して、書き込まれたデータに対応するビットマップフラグ（つまり、このデータが書き込まれた物理記憶位置の物理アドレスに対応するビットマップフラグ）を 0 から 1 に変更する（ステップ S 1 3）。例えば、図 2 6 に示されているように、開始 LBA が LBA<sub>x</sub> である 16 K バイト更新データがブロック BLK # 1 のオフセット + 4 ~ + 7 に対応する物理記憶位置に書き込まれた場合を想定する。この場合、図 2 7 に示されているように、ブロック BLK # 1 用のブロック管理テーブルにおいては、オフセット + 4 ~ + 7 に対応するビットマップフラグそれぞれが 0 から 1 に変更される。

20

【0309】

コントローラ 4 は、このライトコマンドに対するレスポンスをホスト 2 に返す（ステップ S 1 4）。例えば、図 2 6 に示されているように、開始 LBA が LBA<sub>x</sub> である 16 K バイト更新データがブロック BLK # 1 のオフセット + 4 ~ + 7 に対応する物理記憶位置に書き込まれたならば、LBA<sub>x</sub>、ブロック番号（= BLK 1）、オフセット（= + 4）、長さ（= 4）を含むレスポンスがコントローラ 4 からホスト 2 に送信される。

30

【0310】

ホスト 2 がこのレスポンスを受信した時、ホスト 2 は、ホスト 2 によって管理されている LUT を更新して、書き込まれたライトデータに対応する論理アドレスそれぞれに物理アドレスをマッピングする。図 2 8 に示されているように、LUT は、複数の論理アドレス（例えば LBA）それぞれに対応する複数のエントリを含む。ある論理アドレス（例えばある LBA）に対応するエントリには、この LBA に対応するデータが格納されている NAND 型フラッシュメモリ 5 内の位置（物理記憶位置）を示す物理アドレス PBA、つまりブロック番号、オフセット（ブロック内オフセット）が格納される。図 2 6 に示されているように、開始 LBA が LBA<sub>x</sub> である 16 K バイト更新データがブロック BLK # 1 のオフセット + 4 ~ + 7 に対応する物理記憶位置に書き込まれたならば、図 2 8 に示されているように、LUT が更新されて、LBA<sub>x</sub> に対応するエントリに BLK # 1、オフセット + 4 が格納され、LBA<sub>x</sub> + 1 に対応するエントリに BLK # 1、オフセット + 5 が格納され、LBA<sub>x</sub> + 2 に対応するエントリに BLK # 1、オフセット + 6 が格納され、LBA<sub>x</sub> + 3 に対応するエントリに BLK # 1、オフセット + 7 が格納される。

40

【0311】

この後、ホスト 2 は、上述の更新データの書き込みによって不要になった以前のデータを無効化するための Trim コマンドをフラッシュストレージデバイス 3 に送信する（ステ

50

ップS 2 1)。図 2 6 に示されているように、以前のデータがブロック B L K # 0 のオフセット + 0、オフセット + 1、オフセット + 2、オフセット + 3 に対応する位置に格納されている場合には、図 2 9 に示すように、ブロック番号 (= B L K # 0)、オフセット (= + 0)、長さ (= 4) を指定する T r i m コマンドがホスト 2 からフラッシュストレージデバイス 3 に送信される。フラッシュストレージデバイス 3 のコントローラ 4 は、この T r i m コマンドに応じて、ブロック管理テーブル 3 2 を更新する (ステップ S 1 5)。ステップ S 1 5 においては、図 2 9 に示すように、ブロック B L K # 0 用のブロック管理テーブルにおいて、オフセット + 0 ~ + 3 に対応するビットマップフラグそれぞれが 1 から 0 に変更される。

#### 【 0 3 1 2 】

図 6 1 は、フラッシュストレージデバイス 3 に適用される G C 制御コマンドの別の例を示す。具体的には、図 3 4 に示した G C 制御コマンドが、フラッシュストレージデバイス 3 がタイプ # 1 - ストレージデバイスとして実現されている場合において適用されるものであるのに対して、図 6 1 に示す G C 制御コマンドは、フラッシュストレージデバイス 3 がタイプ # 2 - ストレージデバイスとして実現されている場合において適用されるものである。図 6 1 中に施されているハッチングは、図 3 4 との相違点を示している。

#### 【 0 3 1 3 】

G C 制御コマンドは、コマンド I D、ポリシー、ソース Q o S ドメイン I D、デスティネーション Q o S ドメイン I D、等を含んでもよい。

コマンド I D はこのコマンドが G C 制御コマンドであることを示す I D (コマンドコード) であり、G C 制御コマンドには G C 制御コマンド用のコマンド I D が含まれる。

#### 【 0 3 1 4 】

ポリシーは、G C 候補ブロック (G C ソースブロック) を選択するための条件 (G C ポリシー) を指定するパラメータである。フラッシュストレージデバイス 3 のコントローラ 4 は、複数の G C ポリシーをサポートしている。

コントローラ 4 によってサポートされている G C ポリシーには、有効データ量が少ないブロックを優先的に G C 候補ブロック (G C ソースブロック) として選択するというポリシー (G r e e d y) が含まれてもよい。

#### 【 0 3 1 5 】

また、コントローラ 4 によってサポートされている G C ポリシーには、低い更新頻度を有するデータ (コールドデータ) が集められているブロックを、高い更新頻度を有するデータ (ホットデータ) が集められているブロックよりも優先的に G C 候補ブロック (G C ソースブロック) として選択するというポリシーが含まれていてもよい。

#### 【 0 3 1 6 】

さらに、G C ポリシーは、G C 開始条件を指定してもよい。G C 開始条件は、例えば、残りフリーブロックの個数を示してもよい。

コントローラ 4 は、有効データを含むブロック群をアクティブブロックリストによって管理しており、G C を実行する場合には、G C 制御コマンドによって指定された G C ポリシーに基づいて、アクティブブロックリストによって管理されているブロック群から一つ以上の G C 候補ブロック (G C ソースブロック) を選択する。

#### 【 0 3 1 7 】

Q o S ドメイン I D は、G C を実行すべき Q o S ドメインを指定するパラメータである。コントローラ 4 は、Q o S ドメイン I D によって指定される Q o S ドメインに属するブロック群、つまりこの Q o S ドメインに対応するアクティブブロックリストから、一つ以上の G C 候補ブロック (G C ソースブロック) を選択する。また、コントローラ 4 は、Q o S ドメイン I D によって指定される Q o S ドメインに属するフリーブロック群内の一つ以上のフリーブロックを G C デスティネーションブロックとして選択する。

#### 【 0 3 1 8 】

コントローラ 4 は、Q o S ドメインに対応する残りフリーブロックの数がポリシーによって指定される閾値以下になった場合に、G C を開始してもよい。もし G C の強制実行を指

10

20

30

40

50

定するポリシーを含むGC制御コマンドを受信したならば、コントローラ4は、ホスト2からこのGC制御コマンドを受信した時にGCを即座に開始してもよい。

【0319】

図62は、ガベージコレクション(GC)動作の手順の別の例を示す。具体的には、図36に示した手順が、フラッシュストレージデバイス3がタイプ#1-ストレージデバイスとして実現されている場合におけるものであるのに対して、図62に示す手順は、フラッシュストレージデバイス3がタイプ#2-ストレージデバイスとして実現されている場合におけるものである。図62中に施されているハッチングは、図36との相違点を示している。

【0320】

フラッシュストレージデバイス3のコントローラ4は、ホスト2によって指定されたポリシーに基づいて、QoSドメインIDによって指定されるQoSドメインに属するブロック群から、有効データと無効データとが混在する一つ以上のGCソースブロック(コピー元ブロック)を選択する(ステップS41)。次いで、コントローラ4は、QoSドメインIDによって指定されるQoSドメインに属するフリーブロック群から一つ以上のフリーブロックを選択し、選択したフリーブロックをGCデスティネーションブロック(コピー先ブロック)として割り当てる(ステップS42)。

【0321】

コントローラ4は、GCソースブロック(コピー元ブロック)内の全ての有効データをGCデスティネーションブロック(コピー先ブロック)にコピーする(ステップS44)。ステップS44では、コントローラ4は、GCソースブロック(コピー元ブロック)内の有効データのみならず、この有効データとこの有効データに対応する論理アドレスの双方を、GCソースブロック(コピー元ブロック)からGCデスティネーションブロック(コピー先ブロック)にコピーする。これにより、GCデスティネーションブロック(コピー先ブロック)内にデータと論理アドレスとのペアを保持することができる。

【0322】

そして、コントローラ4は、コピーされた有効データの論理アドレスと、この有効データがコピーされたGCデスティネーションブロック(コピー先ブロック)内の位置を示すデスティネーション物理アドレス(ブロック番号、オフセット(ブロック内オフセット))を、GC用コールバックコマンドを使用してホスト2に通知する(ステップS44)。なお、ステップS44では、コントローラ4は、コピーされた有効データの論理アドレスとデスティネーション物理アドレスとみならず、ソース物理アドレスもホスト2に通知してもよい。

【0323】

ホスト2がこのGC用コールバックコマンドを受信した時、ホスト2は、ホスト2によって管理されているLUTを更新して、コピーされた有効データに対応する論理アドレスそれぞれにデスティネーション物理アドレスをマッピングする(ステップS51)。

【0324】

図63および図64のフローチャートは、フラッシュストレージデバイス3が、ホスト2からのデータを書き込むためのブロックと、フラッシュストレージデバイス3内のデータをコピーするためのブロックとに分離する仕組みを備えている場合における、フリーブロックの割り当ての手順を示す。

【0325】

図63に示す手順は、図60に示す書き込み動作処理のシーケンス中のステップS11(書き込み先ブロックの割り当て)におけるフリーブロックの割り当ての手順である。フラッシュストレージデバイス3のコントローラ4は、ホスト2からのライトデータを書き込むためのブロック(書き込み先ブロック)が割り当てられているか否かを判定する(ステップF1)。仮に、空きページを含むブロックが割り当てられていても、そのブロックが、フラッシュストレージデバイス3内のデータをコピーするためのブロックであったならば、コントローラ4は、ホスト2からのライトデータを書き込むためのブロックは割

10

20

30

40

50

り当てられていないと判定する。この判定は、たとえば、ブロックのメタデータなどとして保持される、そのブロックの用途を示す属性情報に基づいて実行される。

【0326】

ホスト2からのライトデータを書き込むためのブロックは割り当てられていないと判定した場合(ステップF1:NO)、コントローラ4は、たとえば同一の仮想ストレージデバイス内のQoSドメイン間で共有されるフリーブロック群の中の1つのフリーブロックを、ホスト2からのライトデータを書き込むためのブロックとして割り当てる(ステップF2)。このとき、コントローラ4は、このブロックのメタデータなどとして、ホスト2からのライトデータを書き込むためのブロックであることを示す属性情報を記録する。

【0327】

一方、図64に示す手順は、図62に示すガベージコレクション(GC)動作の手順中のステップS42(GCデスティネーションブロックの割り当て)におけるフリーブロックの割り当ての手順である。

フラッシュストレージデバイス3のコントローラ4は、GCソースブロック(コピー元ブロック)内の有効データをコピーするためのブロック(GCデスティネーションブロック)、つまり、フラッシュストレージデバイス3内のデータをコピーするためのブロックが割り当てられているか否かを判定する(ステップF11)。仮に、空きページを含むブロックが割り当てられていても、そのブロックが、ホスト2からのライトデータを書き込むためのブロックであったならば、コントローラ4は、フラッシュストレージデバイス3内のデータをコピーするためのブロックは割り当てられていないと判定する。

【0328】

フラッシュストレージデバイス3内のデータをコピーするためのブロックは割り当てられていないと判定した場合(ステップF11:NO)、コントローラ4は、たとえば同一の仮想ストレージデバイス内のQoSドメイン間で共有されるフリーブロック群の中の1つのフリーブロックを、フラッシュストレージデバイス3内のデータをコピーするためのブロックとして割り当てる(ステップF12)。このとき、コントローラ4は、このブロックのメタデータなどとして、フラッシュストレージデバイス3内のデータをコピーするためのブロックであることを示す属性情報を記録する。

【0329】

以上のように、本実施形態のフラッシュストレージデバイス3によれば、I/O性能の改善を図ることができる。

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】

【0330】

2...ホスト、3...フラッシュストレージデバイス、4...コントローラ、5...NAND型フラッシュメモリ、21...ライト動作制御部、22...リード動作制御部、23...GC動作制御部。

10

20

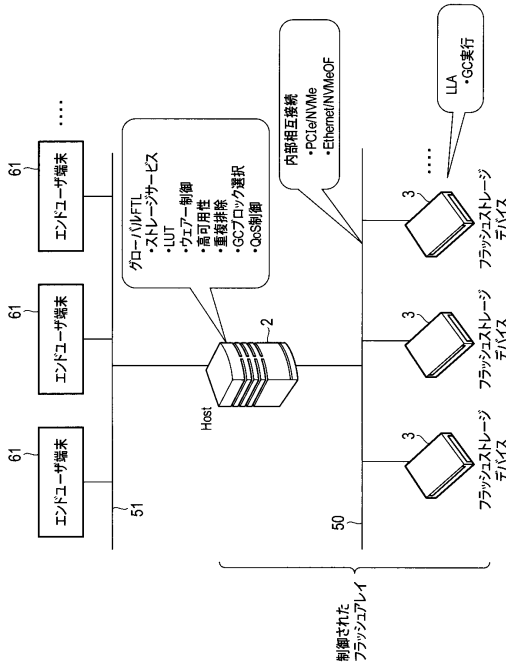
30

40

50

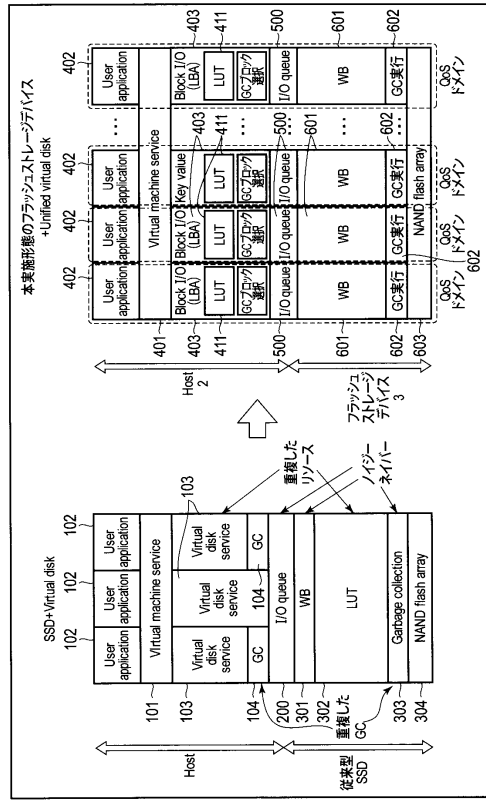
【図面】  
【図 1】

図 1



【図 2】

図 2

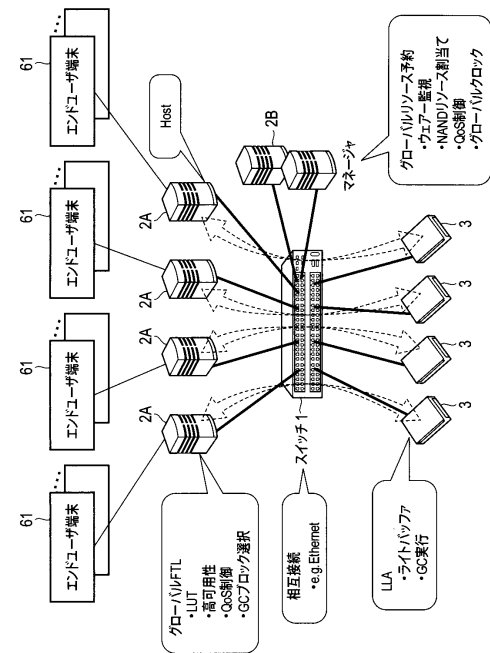


10

20

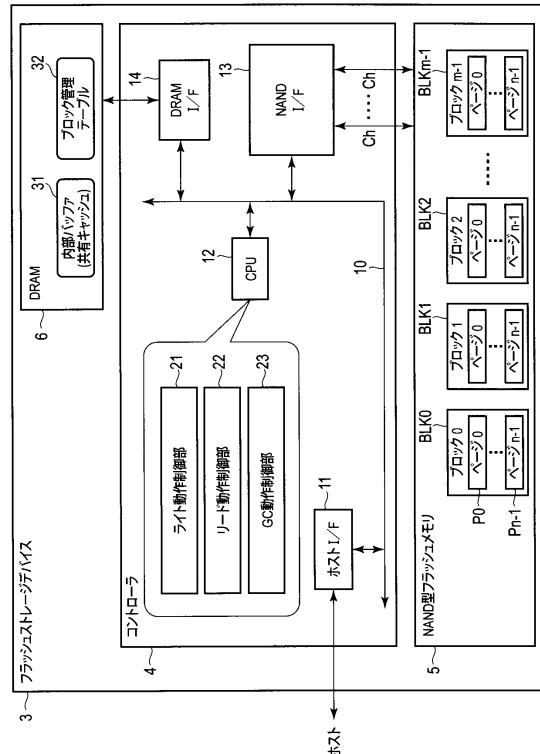
【図 3】

図 3



【図 4】

図 4

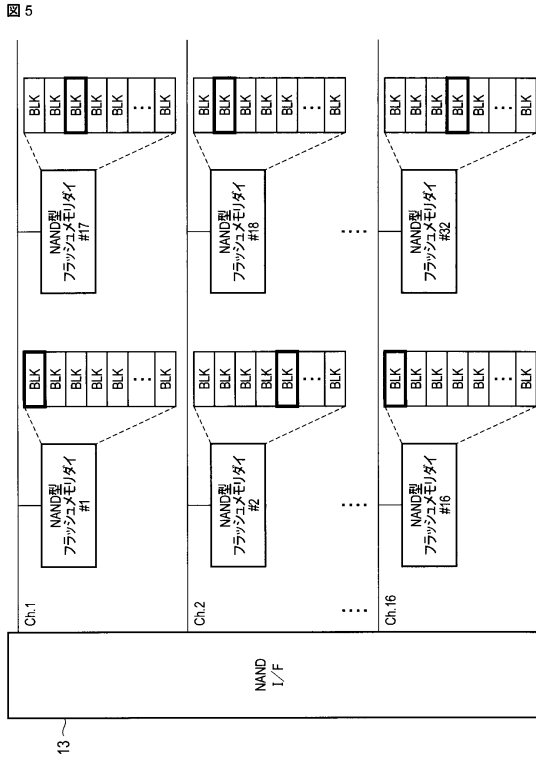


30

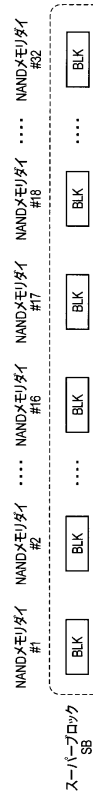
40

50

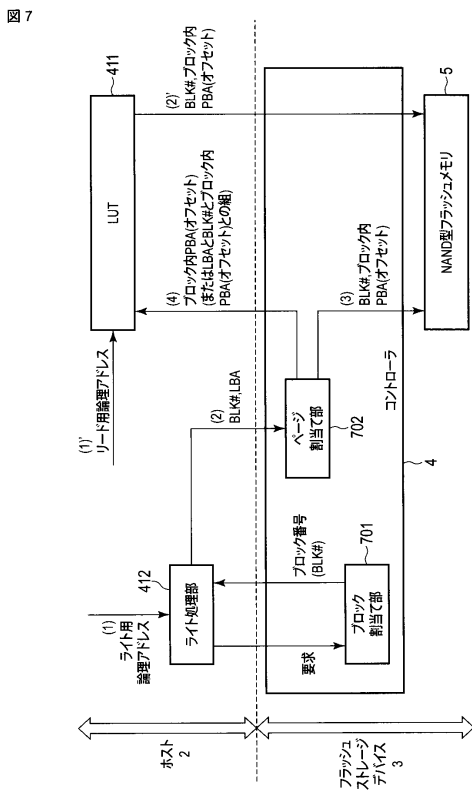
【図 5】



【図 6】



【図 7】



【図 8】

図 8  
ライトコマンド

コマンドID	説明
コマンドID	ライトコマンド用のコマンドID
ブロック番号 (BLK#)	データが書き込まれるべきブロックを指定する物理アドレス
論理アドレス	論理アドレスはデータが書き込まれるべき最初の論理位置を示す
長さ	書き込むべきデータのサイズ

10

20

30

40

50

【 図 9 】

図 9

ライトコマンドに対するレスポンス

	説明
ブロック内物理アドレス	ブロック内物理アドレスは、データが書き込まれたブロック内の位置を示す。ブロック内物理アドレスは、ブロック内オフセットによって指定可能。
長さ	書き込まれたデータの長さ(データ長は、Grainの数によって指定可能)

【 図 1 0 】

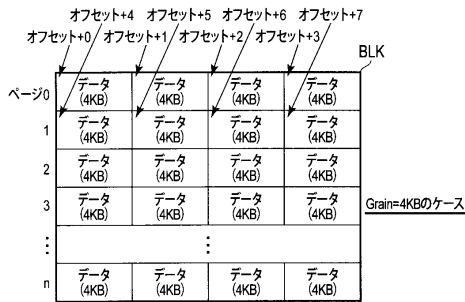
図 10

Trimコマンド(物理アドレスを指定可能)

	説明
コマンドID	Trimコマンド用のコマンドID
物理アドレス	物理アドレスは、無効化すべきデータが格納されている最初の物理記憶位置を示す。物理アドレスは、ブロック番号およびオフセットによって指定可能。
長さ	無効化すべきデータの長さ(データ長は、Grainの数によって指定可能)

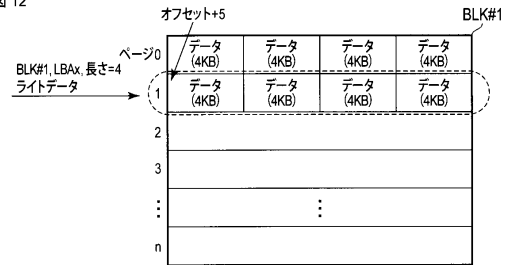
【 図 1 1 】

図 11



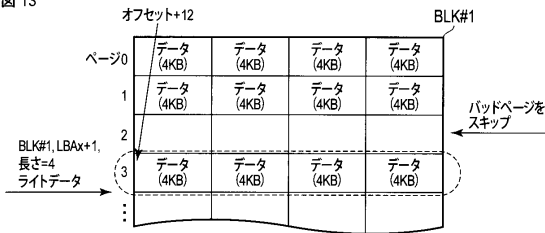
【 図 1 2 】

図 12



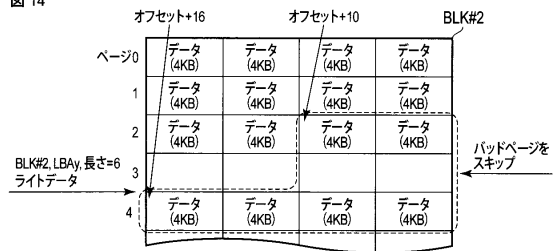
【 図 1 3 】

図 13



【 図 1 4 】

図 14



10

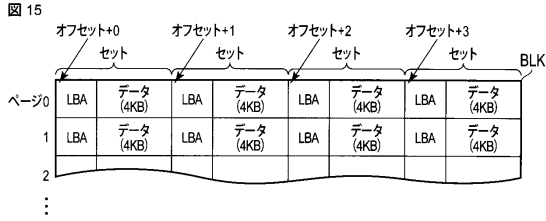
20

30

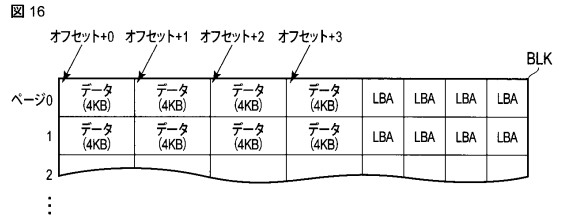
40

50

【 図 15 】

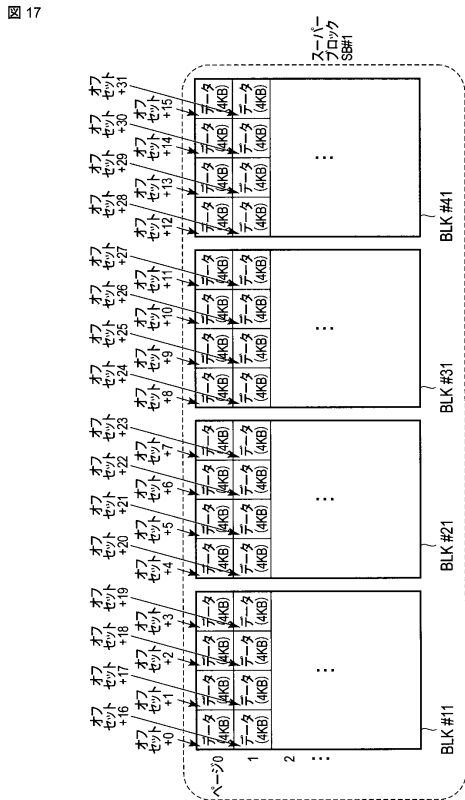


【 図 16 】



10

【 図 17 】



【 図 18 】

図 18

最大ブロック番号ゲットコマンド

コマンドID	説明
コマンドID	最大ブロック番号ゲット用のコマンドID
パラメータなし	

20

30

40

50



【図 19】

図 19

最大ブロック番号取得コマンドに対するレスポンス

	説明
最大ブロック番号	最大ブロック番号はデバイス内ブロック番号の最大値(つまりデバイス内に存在する利用可能なブロックの数)を示す

【図 20】

図 20

ブロックサイズ取得コマンド

	説明
コマンドID	ブロックサイズ取得用のコマンドID
(オプション) ブロック番号	

10

【図 21】

図 21

ブロックサイズ取得コマンドに対するレスポンス

	説明
ブロックサイズ	ブロック番号が指定されたならば、デバイスは、指定されたブロック番号に対応するブロックのサイズをホストに返す

【図 22】

図 22

ブロックアローケートコマンド

	説明
コマンドID	ブロックアローケート用のコマンドID
	ホストは、フリーブロックを割当てるようにデバイスに要求し、これによってブロック番号をデバイスから取得することができる

20

30

40

50

【 図 2 3 】

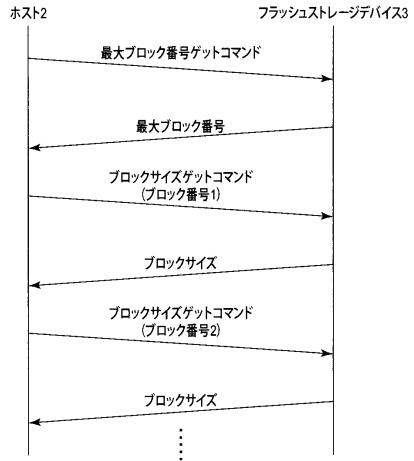
図 23

ブロックアローケートコマンドに対するレスポンス

	説明
ブロック番号	デバイスはフリーブロックリストから、ホストに割当てべきブロックを選択し、選択されたブロックのブロック番号をホストに返す

【 図 2 4 】

図 24

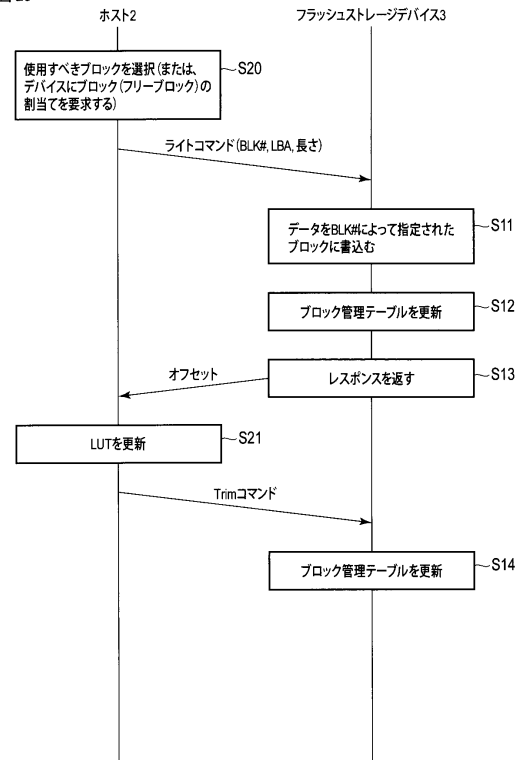


10

20

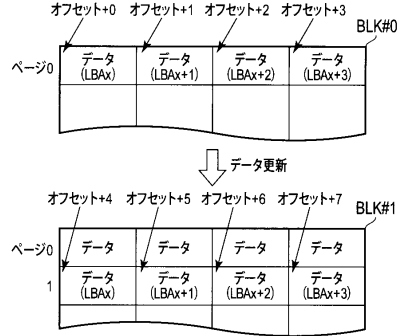
【 図 2 5 】

図 25



【 図 2 6 】

図 26



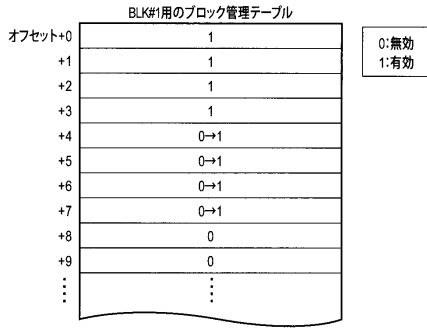
30

40

50

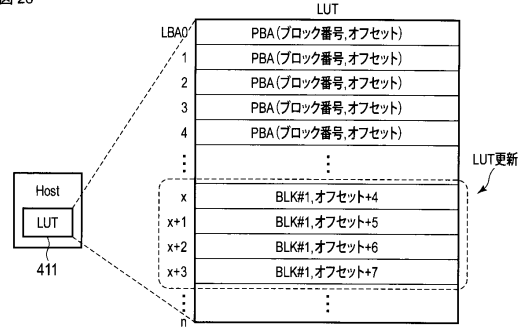
【図 27】

図 27



【図 28】

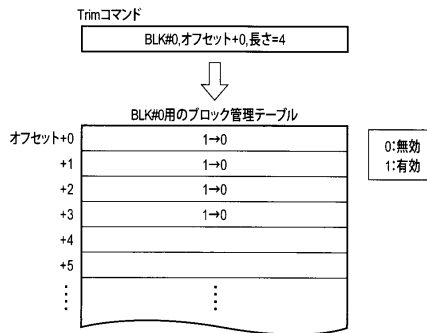
図 28



10

【図 29】

図 29



【図 30】

図 30

リードコマンド

リードコマンド	説明
コマンドID	リードコマンド用のコマンドID
物理アドレス (PBA) #1	物理アドレスはデータが読み出されるべき最初の物理記憶位置を示す。物理アドレスはブロック番号およびオフセットによって指定可能。
長さ#1	リードすべきデータの長さ(データ長はGrainの数によって指定可能)
物理アドレス (PBA) #2	物理アドレスはデータが読み出されるべき最初の物理記憶位置を示す。物理アドレスはブロック番号およびオフセットによって指定可能。
長さ#2	リードすべきデータの長さ(データ長はGrainの数によって指定可能)
転送先ポインタ	リードされたデータが転送されるべきホストメモリ内の位置

20

30

40

50



【 図 3 5 】

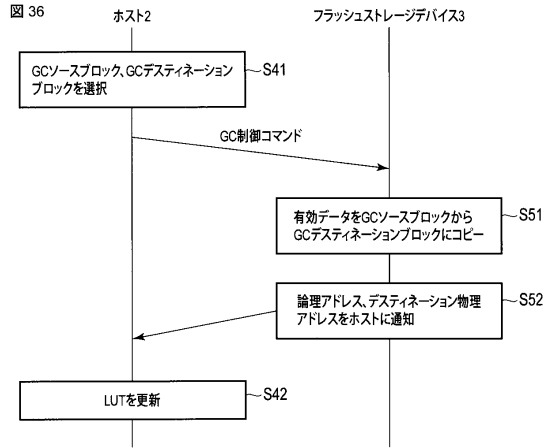
図 35

GC用コールバックコマンド

	説明
コマンドID	GC用コールバック用のコマンドID
論理アドレス	コピーされた有効データの論理アドレス
長さ	コピーされた有効データの長さ
デスティネーション物理アドレス	デスティネーション物理アドレスは、有効データがコピーされたGCデスティネーションブロック内の位置を示す。デスティネーション物理アドレスは、ブロック番号およびオフセットによって指定可能。

【 図 3 6 】

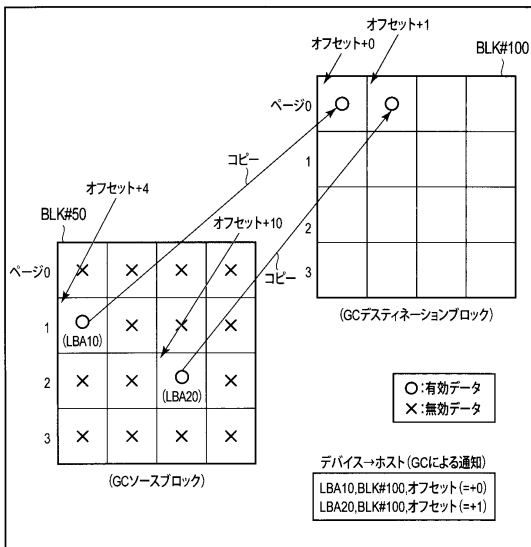
図 36



10

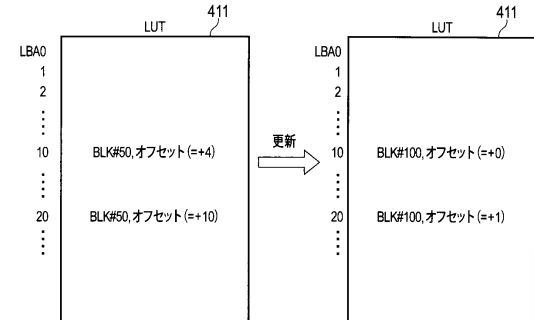
【 図 3 7 】

図 37



【 図 3 8 】

図 38



20

30

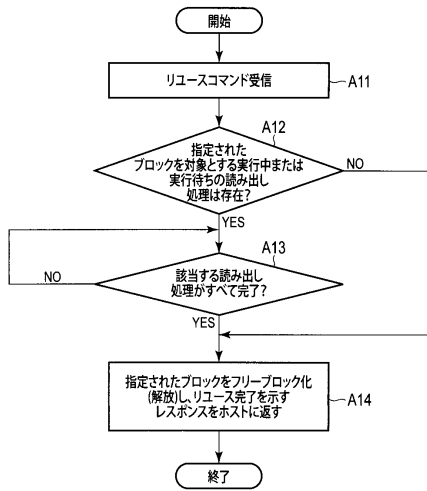
40

50



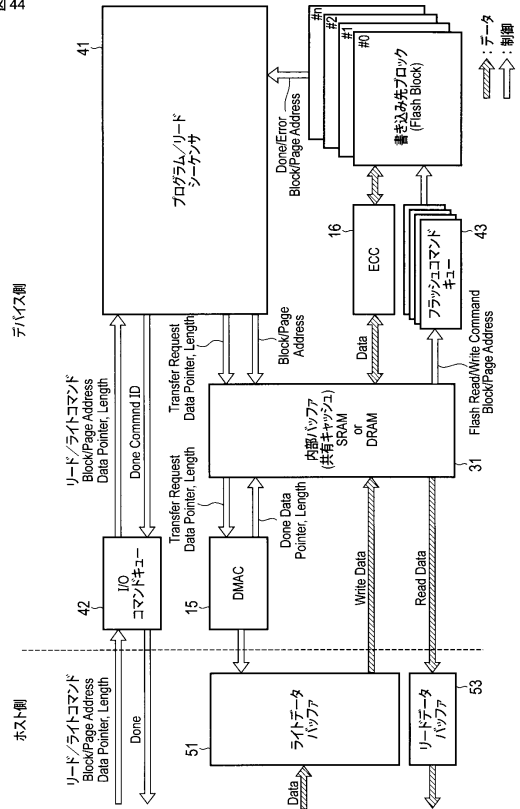
【図 4 3】

図 43



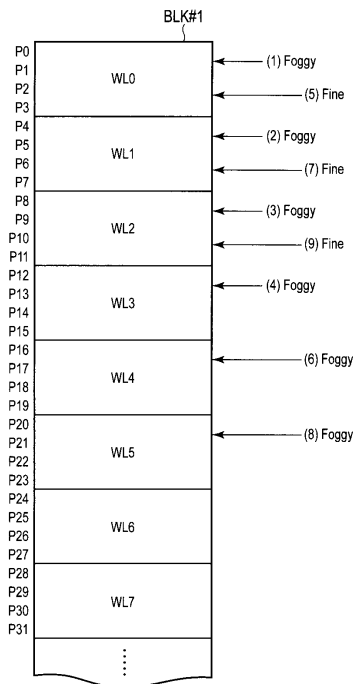
【図 4 4】

図 44



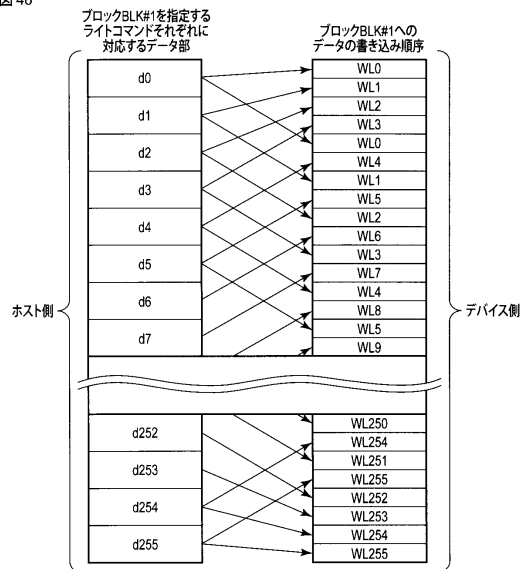
【図 4 5】

図 45



【図 4 6】

図 46



10

20

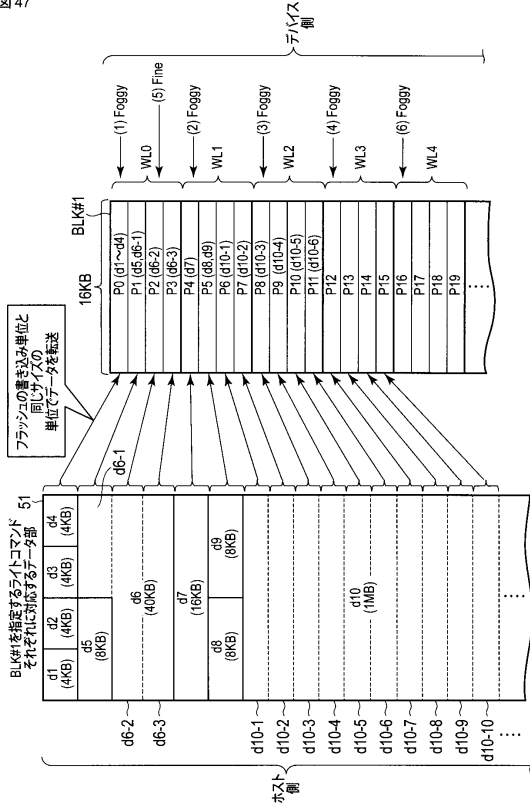
30

40

50

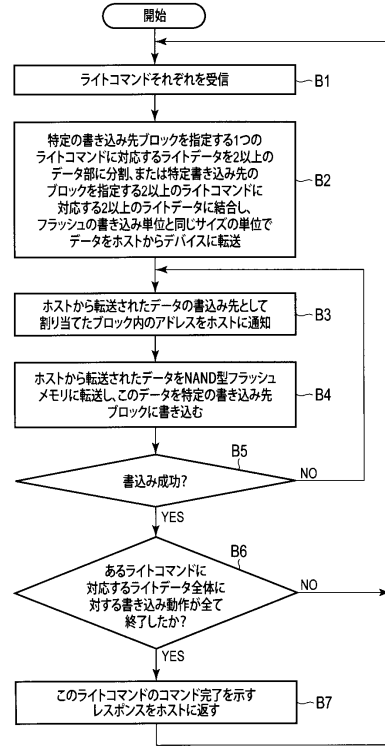
【 図 4 7 】

図 47



【 図 4 8 】

図 48

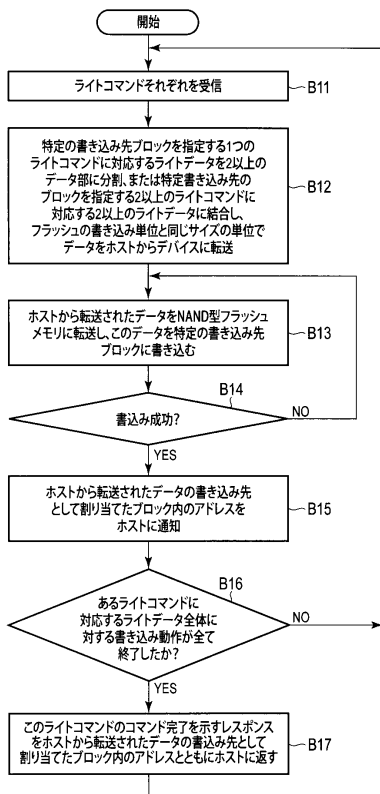


10

20

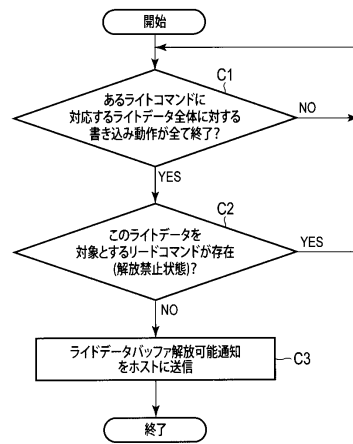
【 図 4 9 】

図 49



【 図 5 0 】

図 50



30

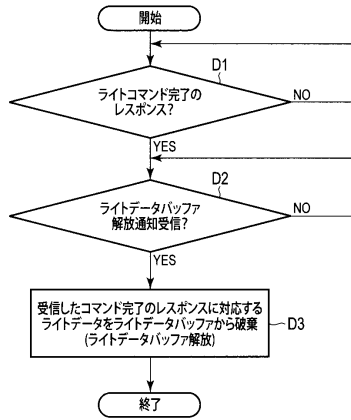
40

50



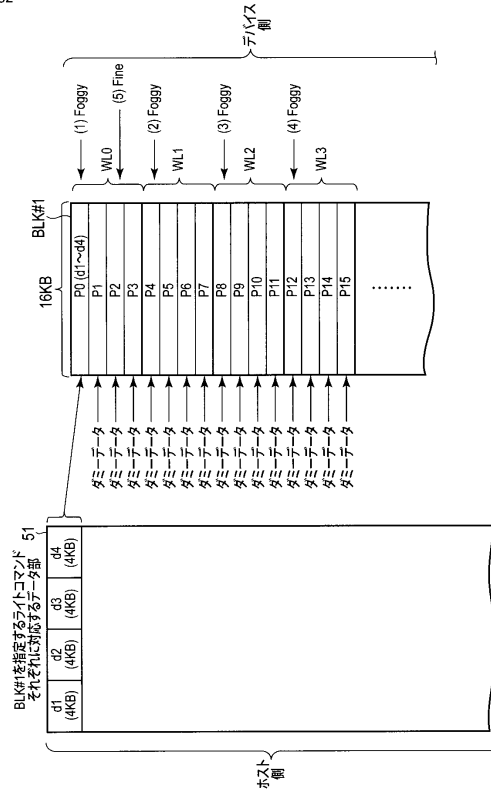
【 図 5 1 】

図 51



【 図 5 2 】

図 52

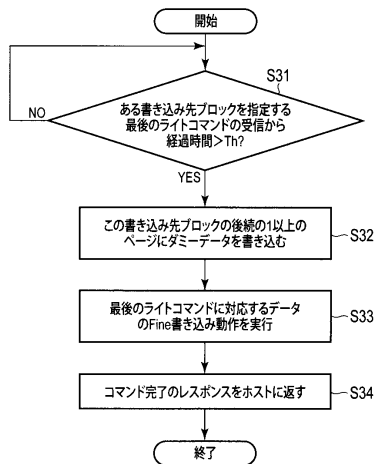


10

20

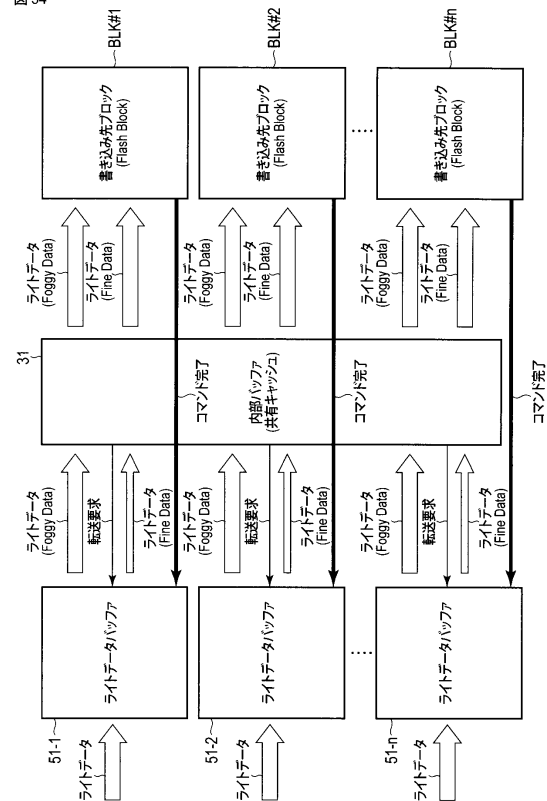
【 図 5 3 】

図 53



【 図 5 4 】

図 54



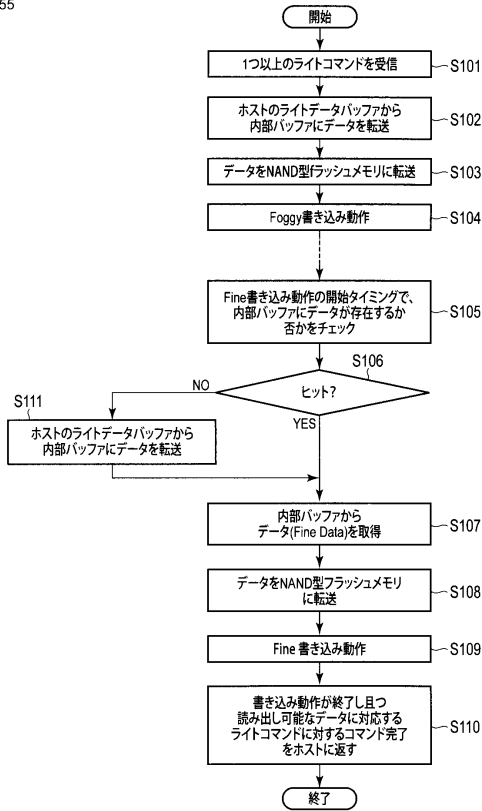
30

40

50

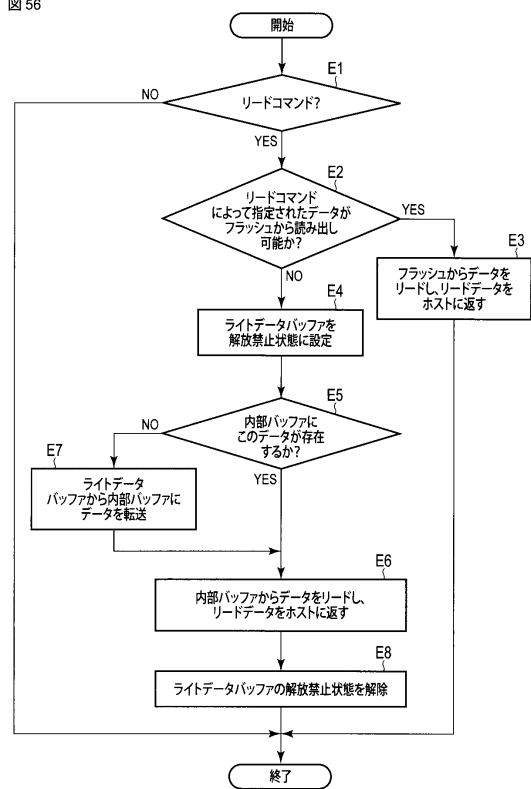
【 図 5 5 】

図 55



【 図 5 6 】

図 56



【 図 5 7 】

図 57

ブロックリユースコマンド

コマンドID	説明
コマンドID	ブロックリユース用のコマンドID
QoSドメインID	QoSドメインを指定する識別子
ブロックアドレス	フリーブロック化(解放)するブロックを指定する物理アドレス

【 図 5 8 】

図 58

ライトコマンド

コマンドID	説明
コマンドID	ライトコマンド用のコマンドID
QoSドメインID	データが書き込まれるべきQoSドメインを一貫に識別するための識別子
論理アドレス	論理アドレスはデータが書き込まれるべき論理位置を示す
長さ	書き込むべきデータの長さ(データ長は、Grainの数によって指定可能)

10

20

30

40

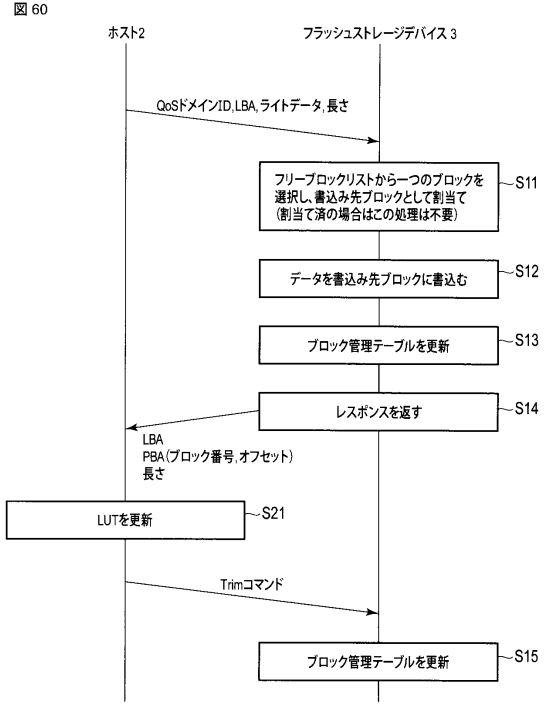
50

【 図 5 9 】

図 59  
ライトコマンドに対するレスポンス

	説明
論理アドレス	ライトコマンドに含まれていた論理アドレス
物理アドレス	物理アドレスはデータが書き込まれた物理記憶位置を示す。物理アドレスは、ブロック番号およびオフセットによって指定される
長さ	書き込まれたデータの長さ(データ長は、Grainの数によって指定可能)

【 図 6 0 】



10

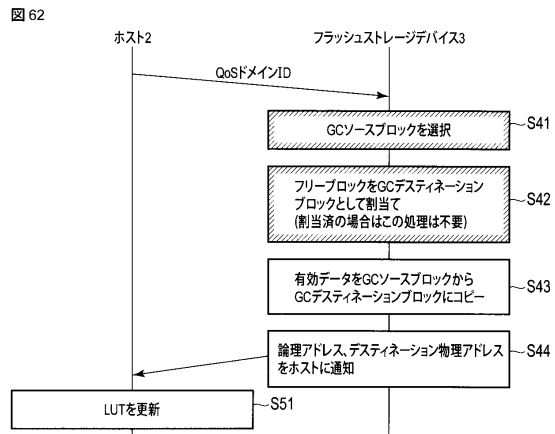
20

【 図 6 1 】

図 61  
GC制御コマンド

	説明
コマンドID	GC制御用のコマンドID
ポリシー	ホストは、GC候補ブロックを選択するためのポリシーを指定することができる。デバイスは、複数のポリシーをサポートしており、ホストによって指定されたポリシーを使用してGC候補ブロックを選択する
QoSドメインID	GCを実行すべきQoSドメインを一意に識別するための識別子

【 図 6 2 】



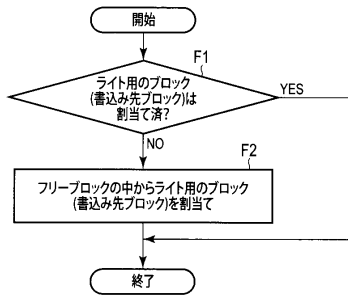
30

40

50

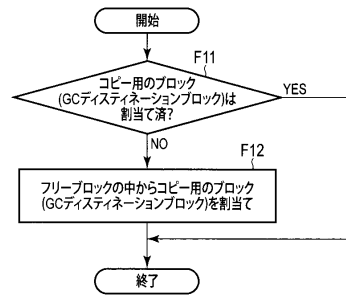
【 図 6 3 】

図 63



【 図 6 4 】

図 64



10

20

30

40

50

## フロントページの続き

(51)国際特許分類

F I  
G 0 6 F 3/06 3 0 2 E

(56)参考文献

特開 2 0 1 7 - 1 6 2 0 6 5 ( J P , A )  
米国特許出願公開第 2 0 1 7 / 0 2 6 2 1 7 5 ( U S , A 1 )  
特開 2 0 1 8 - 1 4 2 2 3 7 ( J P , A )  
特開 2 0 1 3 - 1 0 9 4 1 9 ( J P , A )

(58)調査した分野 (Int.Cl. , D B 名)

G 0 6 F 1 2 / 0 0 - 1 2 / 0 2  
G 0 6 F 3 / 0 8  
G 0 6 F 3 / 0 6