

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2008-16024

(P2008-16024A)

(43) 公開日 平成20年1月24日(2008.1.24)

(51) Int. Cl.	F I	テーマコード (参考)
<b>G06F 3/06 (2006.01)</b>	G06F 3/06 302A	5B005
<b>G06F 12/08 (2006.01)</b>	G06F 3/06 540	5B065
	G06F 12/08 501F	
	G06F 12/08 557	
	G06F 12/08 543B	
審査請求 未請求 請求項の数 19 O L 外国語出願 (全 15 頁)		

(21) 出願番号 特願2007-172478 (P2007-172478)  
 (22) 出願日 平成19年6月29日 (2007.6.29)  
 (31) 優先権主張番号 11/480, 128  
 (32) 優先日 平成18年6月30日 (2006.6.30)  
 (33) 優先権主張国 米国 (US)

(71) 出願人 500373758  
 シーゲイト テクノロジー エルエルシー  
 アメリカ合衆国, カリフォルニア, スコッ  
 ツ バレイ, ピー. オー. ボックス 66  
 360, ディスク ドライブ 920  
 (74) 代理人 100066692  
 弁理士 浅村 皓  
 (74) 代理人 100072040  
 弁理士 浅村 肇  
 (74) 代理人 100091339  
 弁理士 清水 邦明  
 (74) 代理人 100094673  
 弁理士 林 拓三

最終頁に続く

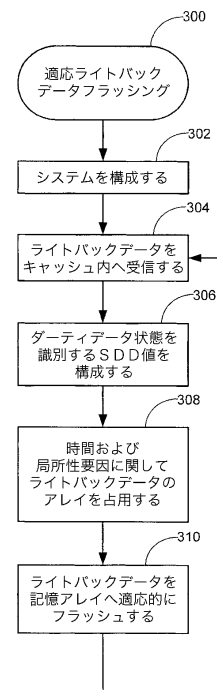
(54) 【発明の名称】 キャッシュされたデータのダイナミック適応フラッシング

## (57) 【要約】

【課題】 キャッシュされたライトバックデータを記憶アレイにフラッシングする方法および装置を得る。

【解決手段】 ライトバックデータのセットが記憶アレイの異なる場所にわたってデータの実質的に均一な分布を維持する目的でキャッシュメモリ内に整列して蓄積される。その後、データの整列されたセットは実質的にライトバックデータの付加セットがホストによりキャッシュメモリへ提供されるレートでキャッシュメモリから記憶アレイへ転送される。好ましくは、各ライトバックデータは複数の隣接データブロックを含み、好ましくは記憶アレイに関してデータの選択された近接範囲内の別個のアクセスコマンドの操作と共に記憶装置に書き込まれる（フラッシュされる）。好ましくは、アレイ内のライトバックデータの各セットに対してストライプデータ記述子（SDD）が維持される。

【選択図】 図8



## 【特許請求の範囲】

## 【請求項 1】

記憶アレイの場所に関して記憶アレイ内のライトバックデータのセットのアレイを形成するステップと、実質的にライトバックデータの付加セットがホストによりキャッシュメモリへ提供されるレートでライトバックデータの前記セットをキャッシュメモリから記憶アレイへ転送するステップと、を含む方法。

## 【請求項 2】

請求項 1 に記載の方法であって、ライトバックデータの各セットは複数の隣接データブロックを含む方法。

## 【請求項 3】

請求項 1 に記載の方法であって、転送ステップは記憶アレイのトランスデューサを記録媒体に隣接する選択された場所へ移動させて選択されたアクセスコマンドをサービスするステップと、選択されたアクセスコマンドの選択された場所とライトバックデータの選択されたセットに対する関連場所との近接に関してライトバックデータの選択されたセットをキャッシュメモリから媒体へ転送するステップと、を含む方法。

## 【請求項 4】

請求項 3 に記載の方法であって、前記近接はトランスデューサの選択されたシーク長を含む方法。

## 【請求項 5】

請求項 1 に記載の方法であって、さらに、アレイ内のライトバックデータの各セットに対してストライプデータ記述子 (SDD) を与えるステップを含み、SDD はライトバックデータを記憶アレイへフラッシングするための準備完了として識別する値を含む方法。

## 【請求項 6】

請求項 1 に記載の方法であって、キャッシュメモリはライトバックデータの m セットを格納し、蓄積ステップは記憶アレイへのフラッシングに備えてアレイの異なる非隣接場所に関連する前記ライトバックデータの n セットを配列するステップを含み、n は m よりも小さい方法。

## 【請求項 7】

請求項 1 に記載の方法であって、さらに、ライトバックデータの前記付加セットがキャッシュメモリへ移されるレートを監視するステップを含み、それに関連して転送ステップを実施して前記キャッシュメモリ内に前記ライトバックデータの実質的に均一な分布を維持する方法。

## 【請求項 8】

記憶アレイの関連する場所へ転送するのに利用できるライトバックデータのセットを格納するキャッシュメモリと、ライトバックデータの前記セットをキャッシュメモリ内に配列して実質的に均一な分布を記憶アレイを横切って提供し、かつ実質的にライトバックデータの付加セットがホストによりキャッシュメモリへ提供されるレートでライトバックデータの前記セットをキャッシュメモリから記憶アレイへ転送するように構成されたプロセッサと、を含む装置。

## 【請求項 9】

請求項 8 に記載の装置であって、プロセッサはライトバックデータのセットを選択的にフラッシングリストへ移すキャッシュ・マネージャを含む装置。

## 【請求項 10】

請求項 8 に記載の装置であって、記憶アレイは記録媒体に隣接する選択された場所へ移されて選択されたアクセスコマンドをサービスするように構成されたトランスデューサを含み、プロセッサは選択されたアクセスコマンドの選択された場所とライトバックデータの選択されたセットに対する関連場所との間の近接に関してライトバックデータの選択されたセットをキャッシュメモリから媒体へ転送する装置。

## 【請求項 11】

請求項 10 に記載の装置であって、前記近接はトランスデューサの選択されたシーク長

10

20

30

40

50

を含む装置。

【請求項 1 2】

請求項 8 に記載の装置であって、プロセッサはアレイ内のライトバックデータの各セットに対するストライプデータ記述子 (SDD) を発生し、SDD はライトバックデータを記憶アレイへのフラッシング準備完了として識別する値を含む装置。

【請求項 1 3】

請求項 8 に記載の装置であって、キャッシュメモリはライトバックデータの m セットを格納し、プロセッサは前記ライトバックデータの n セットを記憶アレイフラッシングするために配置し、n は m よりも小さい装置。

【請求項 1 4】

請求項 8 に記載の装置であって、プロセッサはさらにライトバックデータの前記付加セットがキャッシュメモリへ移されるレートを監視し、それに関して前記データを転送してキャッシュメモリ内に前記ライトバックデータの実質的に均一な分布を維持する装置。

【請求項 1 5】

記憶アレイの関連する場所へ転送するのに利用できるライトバックデータの複数のセットを配列するキャッシュメモリと、実質的にライトバックデータの付加セットがホストによりキャッシュメモリへ提供されるレートでライトバックデータの前記セットを前記アレイから記憶アレイへ転送する第 1 の手段と、を含む装置。

【請求項 1 6】

請求項 1 5 に記載の装置であって、第 1 の手段はキャッシュ・マネージャを含む装置。

【請求項 1 7】

請求項 1 5 に記載の装置であって、記憶アレイは記録媒体に隣接する選択された場所へ移されて選択されたアクセスコマンドをサービスするように構成されたトランスデューサを含み、第 1 の手段は選択されたアクセスコマンドの選択された場所とライトバックデータの選択されたセットに対する関連場所との間の近接に関してライトバックデータの選択されたセットをキャッシュメモリから媒体へ転送する装置。

【請求項 1 8】

請求項 1 7 に記載の装置であって、前記近接はトランスデューサの選択されたシーク長を含む装置。

【請求項 1 9】

請求項 1 5 に記載の装置であって、第 1 の手段はさらにライトバックデータの前記付加セットがキャッシュメモリへ移されるレートを監視し、前記データを転送してプール内にライトバックデータの実質的に均一な分布を維持する装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は一般的にデータ記憶システムの分野に関し、制約はしないが、特に、記憶アレイへのキャッシュされたデータのダイナミック適応フラッシング方法および装置に向けられている。

【背景技術】

【0002】

高速かつ効率的な方法でデータにアクセスするために記憶装置が使用される。ある種の記憶装置は、媒体表面上に規定されたトラックにデータを書込みその後データを読み出す 1 つ以上のデータトランスデューサと共に、回転可能な記録媒体を使用する。

【0003】

マルチデバイス・アレイ (MDA) は多数の記録装置を利用して統合されたメモリスペースを形成する。MDA に対する 1 つの広く利用されるフォーマットは RAID (redundant array of independent discs) 構成を利用し、入力データはアレイ内の多数の記憶装置にわたって格納される。RAID レベルに応じて、ミラーリング、ストライピングおよびパリティコード発生を含むさまざまな技術を利用

10

20

30

40

50

用して格納されたデータの完全性を高めることができる。

【発明の開示】

【発明が解決しようとする課題】

【0004】

記憶容量および性能のレベルは絶えず高まることが要求され続けているため、このようなアレイ内の記憶装置の動作が管理される方法の改善が継続的に必要とされている。本発明の好ましい実施例は一般的にこれらおよびその他の改善に向けられる。

【課題を解決するための手段】

【0005】

本発明の好ましい実施例は一般的に記憶アレイへキャッシュされたライトバックデータをフラッシングする装置および方法に向けられている。 10

【0006】

好ましい実施例では、記憶アレイの異なる場所にわたってデータの実質的に均一な分布を維持する目的で、ライトバックデータのセットがキャッシュメモリ内に整列して蓄積される。その後、データの整列されたセットは実質的にライトバックデータの付加セットがホストによりキャッシュメモリへ提供されるレートでキャッシュメモリから記憶アレイへ転送される。

【0007】

ライトバックデータの各セットは好ましくは複数の隣接データブロックを含み、好ましくは記憶アレイに関してデータの選択された近接範囲内の別個のアクセスコマンドの操作と共に記憶装置に書き込まれる（フラッシュされる）。好ましくは、アレイ内のライトバックデータの各セットに対してストライプデータ記述子（SDD）が維持される。 20

【0008】

このようにして、キャッシュされたライトバックデータのフラッシングは全体ホストI/Oアクセスレートに著しい変動を生じることはない。

【0009】

下記の詳細な説明を読みかつ添付図を検討すれば、特許請求された本発明を特徴づけるこれらおよびさまざまな他の特徴および利点が明白となる。

【実施例】

【0010】

図1はユーザデータを格納かつ検索するように構成された典型的な記憶装置100を示す。好ましくは、装置100はハードディスクドライブとして特徴づけられるが、所望により、他の装置構成も容易に利用することができる。 30

【0011】

ベースデッキ102がトップカバー（図示せず）と嵌合して密閉筐体を形成する。スピンドルモータ104が筐体内に搭載されて、好ましくは、磁気記録ディスクとして特徴づけられる媒体106を制御可能に回転させる。

【0012】

制御可能に移動可能なアクチュエータ108はボイスコイルモータ（VCM）112へ電流を印加することにより、媒体表面上に規定されたトラックに隣接するリード/ライト・トランスデューサ110のアレイを移動させる。フレックス回路アセンブリ114はアクチュエータ108および外部搭載印刷回路板（PCB）116上の装置制御エレクトロニクス間の電気通信パスを提供する。 40

【0013】

図2はn個の記憶装置（SD）100を有利に内蔵して統合記憶アレイ122を形成する典型的なネットワークシステム120を一般的に示している。冗長コントローラ124、126が好ましくは記憶アレイ122およびサーバ128間でデータを転送するように動作する。サーバ128はローカルエリアネットワーク（LAN）、インターネット、等のファブリック130に接続されている。

【0014】

遠隔ユーザは、それぞれ、パーソナルコンピュータ（PC）132, 134, 136を介してファブリック130にアクセスする。このようにして、選択されたユーザは記憶スペース122にアクセスして、所望により、データを書込みまたは検索することができる。

#### 【0015】

好ましくは、装置100およびコントローラ124, 126はマルチデバイス・アレイ（MDA）138に内蔵される。好ましくは、MDAは1つ以上の選択されたRAID（redundant array of independent discs）構成を使用して装置100に亘ってデータを格納する。図2には1つだけのMDAと3人の遠隔ユーザしか例示されていないが、これは単なる説明が目的であって制約的意味合いは無く、所望により、ネットワークシステム120は任意数およびタイプのMDA、サーバ、クライアントおよびホスト装置、ファブリック構成およびプロトコル、等を利用することができる。図3は図2のネットワーク内で有用なアレイコントローラ構成140を示している。

10

#### 【0016】

図3は中間バス146（「Eバス」と呼ばれる）により接続された2つのインテリジェント記憶プロセッサ（ISP）142, 144を説明する。好ましくは、各ISP142, 144は共通コントローラボード上の別個の集積回路パッケージ内に配置される。好ましくは、各ISP142, 144はファイバチャネル・サーバリンク148, 150を介してアップストリーム・アプリケーション・サーバと通信し、かつファイバチャネル・記憶リンク152, 154を介して記憶装置100と通信する。

20

#### 【0017】

ポリシープロセッサ156, 158はコントローラ140に対するリアルタイム・オペレーティングシステム（RTOS）を実行し、PCIバス160, 162を介して各ISP142, 144と通信する。ポリシープロセッサ156, 158は、さらに、カスタマイズされた論理を実行して定められた記憶アプリケーションに対してISP142, 144と共に精巧な処理タスクを実施する。所望により、ISP142, 144およびポリシープロセッサ156, 158は動作中に必要とされるメモリモジュール164, 166にアクセスする。

#### 【0018】

図4は図3の選択されたISPに対する好ましい構造を提供する。纏めて168に示すいくつかのファンクション・コントローラはホスト交換、直接メモリアクセス（DMA）、排他的or（XOR）、コマンドルーティング、メタデータ制御、およびディスク交換等のいくつかのコントローラ動作に対するファンクション・コントローラ・コア（FCC）として機能する。好ましくは、各FCCはメモリ交換その他のスケジューリングタスクを容易にする非常に柔軟な機能セットおよびインターフェイスを含んでいる。

30

#### 【0019】

一般的に170に示すいくつかのリストマネージャは、キャッシュテーブル管理、メタデータ保守、およびバッファ管理等のコントローラ動作中にさまざまなデータおよびメモリ管理タスクに使用される。好ましくは、リストマネージャ170はメモリ上の単純ではあるが明確に定義された操作を実施してFCC168により指令されるタスクを遂行する。好ましくは、各リストマネージャはFCCによるメモリアクセスに対するメッセージプロセッサとして作動して、好ましくは、規定されたプロトコルに従って受信メッセージにより規定される操作を実行する。

40

#### 【0020】

リストマネージャ170は、それぞれ、交換メモリブロック172、キャッシュテーブル・ブロック174、バッファメモリ・ブロック176およびSRAM178を含むいくつかのメモリモジュールと通信しかつ制御する。ファンクション・コントローラ168およびリストマネージャ170は、それぞれ、クロスポイントスイッチ（CPS）モジュール180を介して通信する。このようにして、コントローラ168の選択されたファンク

50

ション・コアはC P S 1 8 0を介した対応するリストマネージャ1 7 0への通信パスウェイを確立して状態を通信し、メモリモジュールをアクセスし、あるいは所望のI S P動作を呼び出すことができる。

【0 0 2 1】

同様に、選択されたリストマネージャ1 7 0はC P S 1 8 0を介してファンクション・コントローラ1 6 8へ応答を通信し返すことができる。図示されてはいないが、好ましくは、図4の各エレメント間にその間のデータ転送を調整するために別個のデータバス接続が確立される。所望により、他の構成を容易に利用することがお判りであろう。

【0 0 2 2】

P C Iインターフェイス(I / F)モジュール1 8 2はポリシープロセッサ1 5 6およびI S P 1 4 2間でトランザクションを確立して指示する。E - B U S I / Fモジュール1 8 4はF C Cおよび各I S P 1 4 2, 1 4 4のリストマネージャ間のE - B U S 1 4 6を介した通信を容易にする。所望により、ポリシープロセッサ1 5 6, 1 5 8はE - B U S 1 4 6を介したシステムの他の部分との通信を開始し受信することもできる。

【0 0 2 3】

図3および4のコントローラ・アーキテクチャはアレイに対するスケラブルで、非常に機能的なデータ管理および制御を有利に提供する。好ましくは、ストライプ・バッファリスト(S B L)および他のメタデータ構造は記録媒体上のストライプ境界および記憶トランザクション中にディスクストライプに関連付けられるデータを専用に格納するキャッシュ内の基準データバッファに揃えられる。

【0 0 2 4】

処理効率をさらに高めるために、好ましくは、コントローラ・アーキテクチャは新しいライトバックデータ・キャッシング方法論を利用する。一般的に、これはメモリ内の記憶装置1 0 0に書き込まれるデータのキャッシング、および後でこのようなライトバックデータの記憶装置1 0 0への転送をスケジューリングする(フラッシング)ことを含んでいる。

【0 0 2 5】

一般的に、データの時間および局所性の両方を考慮する2次元方法を使用してライトバックデータの隣接ブロックのセットがキャッシュメモリ内に配列される。キャッシュされたライトバックデータの実質的に均一な分布が維持されて、他のアクセス操作と共にデータを書き込むための最適機会を提供する。

【0 0 2 6】

好ましくは、データの隣接ブロックのセットは付加ライトバックデータセットがホストによりキャッシュメモリへ提供されるレートと実質的に一致するレートでキャッシュメモリから記憶アレイへ書き込まれる。このようにして、観察されたホストI / O転送レートの大きな変動は実質的に解消される。

【0 0 2 7】

図5に示すように、好ましくは、キャッシュされたデータはストライプ・データ記述子(S D D) 1 9 2と呼ばれるデータ構造を使用してキャッシュ・マネージャ(C M) 1 9 0によりノードベースで管理される。各S D Dは関連するデータへの最近および現在のアクセスに関するデータを保持する。各S D Dは好ましくは対応するR A I Dストライプ1 9 4(すなわち、特定のパリティセットに関連する選択された装置1 0 0上の全データ)に揃えられ、特定のS B L 1 9 6に従う。

【0 0 2 8】

C M 1 9 0により管理される各キャッシュノードは、好ましくは、ある特定のS D Dを参照し、論理ディスクの定められたセット(装置1 0 0のサブセット)に対するアクティブS D Dは、好ましくは、標準フォワードおよびバックワード・リンクリストを使用して仮想ブロックアドレス(V B A)を介して昇ベキの順でリンクされる。

【0 0 2 9】

好ましくは、V B A値はしばしばR A I D A l l o c a t i o n G r i d S y s t 50

e m ( R A G S ) と呼ばれるグリッドシステムを使用して R A I D データ編成と揃えられる。一般的に、同じ R A I D ストリップ 1 9 8 に属するブロックの任意特定の集り（たとえば、特定のパリティセットに寄与する全データ）が特定のシート上の特定の信頼できる記憶ユニット（ R S U ）に割当てられる。

【 0 0 3 0 】

数枚のシートからなるブックが異なる装置 1 0 0 からのブロックの多数の隣接セットから構成される。実際のシートおよび V B A に基づいて、ブックは特定の装置または装置セットを示す（冗長度が利用される場合）ゾーンへ細分割される。

【 0 0 3 1 】

好ましくは、各 S D D はアクセス歴、ロック状態、最終オフセット、最終ブロック、タイムスタンプデータ（一日の時間、 T O D ）、データがどのゾーン（ブック）に属するかの識別子、および利用される R A I D レベルを含むデータのさまざまな状態を示す変数を含んでいる。好ましくは、S D D に関連するデータのライトバック（「ダーティ」データ）状態がダーティデータ、ダーティバッファ、ダーティ L R U およびフラッシング L R U 値に関して管理される。

【 0 0 3 2 】

好ましくは、C M 1 9 0 は同時に作動していくつかの異なるレベルのライトバックデータ・プロセスをシステム要求条件に応じて管理する。第 1 のレベルは一般的にフル R A I D ストリップ 1 9 8 が検出される時のフル S D D 構造の周期的フラッシングを含んでいる。これは S D D が関連するデータをダーティとして識別する時に R A I D レベル変数に基づいて定められた S D D 1 9 2 に対して容易に実施することができる。好ましくは、これは十分連続的な隣接 S D D 構造がダーティデータで十分満たされているかを決定するバックワード検査を含んでいる。そうであれば、S D D 構造はフラッシングリスト（1 9 9 に示す）上に置かれ、データのフラッシングを開始する要求がなされる。フラッシングリスト状態は S D D 1 9 2 のフラッシング L R U 値を使用して設定することができる。

【 0 0 3 3 】

データの小さなセットのフラッシングは好ましくは S D D ベースで処理される。ダーティブロックを有しロックブロックの無い任意の S D D は好ましくはダーティ L R U として設定され経過期間（ a g e ）（たとえば、データがキャッシュ待機フラッシングにおいて過ごした時間）によりソートされる。特定の経過期間に達すると、好ましくはフラッシング L R U 変数が設定されフラッシングリスト 1 9 9 が更新される。

【 0 0 3 4 】

好ましくは、フラッシングリストからのデータのフラッシングの攻撃性は実質的に付加ダーティデータがキャッシュ内に入るレートでダーティデータを押出すように適応的に調節される。特定範囲の連続的ダーティブロックをフラッシングする予定の時は、C M 1 9 0 は好ましくは近似局所性を有する R A I D レベルに基づいて他の範囲のダーティブロック、すなわち、シーク時間に関して「すぐ近い」または同じ R A I D パリティストリップ 1 9 9 へのアクセスを伴うブロックを捜し出す。

【 0 0 3 5 】

コントローラーアーキテクチャがこれを実施する好ましい方法を図 6 について見ることができ、それはキャッシュされたライトバックデータのアレイ 2 0 0 を表している。アレイ 2 0 0 は C M 1 9 0 またはコントローラーの他の処理ブロックにより維持される。

【 0 0 3 6 】

アレイ 2 0 0 内の各セル 2 0 2 は一般的に記憶装置 1 0 0 内の選択された局所性に対応し、ブック、シートおよび / またはアレイ内のゾーンとして編成することができる。装置内の境界は、たとえば、各コラムが異なる記憶装置 1 0 0 を表しコラム内の各セルはその装置を横切る異なる半径方向バンドを表すように選択することができる。

【 0 0 3 7 】

セルはフラッシングリスト 1 9 9 にフラグを付されている隣接ライトバックデータのセットで「占有されている（ p o p u l a t e d ）」。特に、各占有ブロック 2 0 4（図 6

10

20

30

40

50

に「X」で示す)は記憶装置100内の異なる論理的または物理的場所に対応する変動するサイズのデータブロックの1つの(またはいくつかの)グループを表す。このようにして、キャッシュ内のライトバックデータのセットは記憶アレイ122へ転送されることを予測してプールすることができる。

#### 【0038】

アレイ200はさまざまなデータ装置100にわたってデータのライトバックをスケジューリングする有用なフォーマットを提供する。好ましい実施例では、記憶装置100内の選択された場所にわたってアクセスする(たとえば、リードコマンドを実施する)ように特定のアクセスコマンドがスケジュールされると、アレイ200はアクセスコマンドと共に効率的にサービスされるライトバックデータの利用可能なブロックを識別するように参照される。

10

#### 【0039】

アレイ200を横切って占有セル204の実質的に均一な分布を維持し、着信パーティデータのレートをキャッシュと整合させるためにパーティセットが選択的にアレイ200に加えられる。あるシステム負荷の元で、CM190は比較的多数のフラッシング操作をロードアップして、論理的セットまたは物理的媒体を横切ってIDからODへ進行するライト等の、ショートシークのクラスタを生成するように構成することができる。

#### 【0040】

関連する実施例では、図6のアレイ200は個別のRAIDストライプ(対応するSD192を有する図5の194等)を表すようにセル202を配置するように構成することができ、アレイ200のコラムは前記したRAGSグリッドのコラムに対応することができる。このような場合、定められたロー内の特定の占有セル204のフラッシングを使用して同じロー内の他の占有セルにより使用されるパリティ(グリッド内には示されていない)への参照を示すことができる。

20

#### 【0041】

このようなフラッシング操作の同時スケジューリングにより、特に、RAID-5およびRAID-6環境において性能改善を導くことができ、それはそのロー内の全占有セルに対する4つの(RAID-6の場合は6つの)I/Oアクセスの内の2つが同じパリティRAIDストライプ194にアクセスするためである。

#### 【0042】

もう1つの好ましい実施例では、アレイ200は各コラム(またはロー)が別個の記憶装置100を表し、各セル202は一般的に装置の媒体106の半径方向幅にわたって異なる隣接ゾーン、すなわち領域、に対応するように配置される。

30

#### 【0043】

図7の「W」表記は一般的にこれらのさまざまな場所内のデータの未決ライトバックセットに対応し、したがって各セル202内の装置にわたって分布されるライト機会を表す。各セル202内のW表記の空間的場所は一般的に関連する領域内のそのデータの論理的または物理的場所を表す。W表記は媒体106へ書き込むためにフラッシングリスト199から利用可能な全ライトバックデータセットを必ずしも表すものではない。

#### 【0044】

リード優先環境では、リードコマンドが優先権を有し、したがって一般的にライトコマンドのサービスよりも優先的にサービスされる。しかしながら、ライト優勢環境では、一般的にリードに比べて比較的多数のライトコマンドがある。1つのこのようなリードコマンドは図7において「R」表記で識別され、一般的に関連するデータが検索される媒体106の領域に対応する。

40

#### 【0045】

この実施例では、キャッシュ・マネージャ190は好ましくは関連するリードコマンドを実施してR表記のデータを検索するよう装置100に指令する。このリード操作の終りに、キャッシュ・マネージャ190は、さらに、好ましくは一般的にリードコマンドの近くの(たとえば、同じセル202内の)1つ以上ライトバック操作を装置100に実施さ

50



せるように進行する。

【0046】

図7は「丸W」表記を使用してこのようにサービスされる2つのこのようなライトバックデータセットを識別する、すなわち、2つの丸W表記ライトバックフラッシュが関連するリードコマンドの終りに生じる(R表記)。リードコマンド近くの(たとえば、セル202内の)全ライトバックデータがフラッシュされる必要はなく、望ましくないことさえあることに注目願いたい。しかしながら、近くのデータの少なくともいくつかはフラッシュされ、装置100のトランスデューサ110は一般的にこの近くにあるため、これらのライトバックデータ・フラッシュ操作は低減されたシークレーテンシ(seek latencies)で実施することができる。

10

【0047】

選択されたライトバックデータがフラッシュされると、キャッシュ・マネージャ190は、フラッシュリストから得られる、この同じ領域への付加ライトバックデータ・セットでアレイ200を「埋め戻す」ように進行する。このようにして、新しいライトコマンドは媒体106のさまざまな半径方向幅にわたってライトバックデータ機会の均一な分布を実質的に維持するようにアレイ200に対して計られる(metered)。

【0048】

ディスクアクセスを要する未決リードコマンドが無い限り、キャッシュ・マネージャ190は一般的に前と同様にライトバックデータをフラッシュするように動作する。しかしながら、各新しいリードコマンドが発せられると、リードコマンドに優先権が与えられ1つ以上の付加ライトバックセットがリードコマンドの一般的に近くからフラッシュされる。

20

【0049】

好ましくは、これは次のリードコマンドがどこへ向けられても効率的な方法でフラッシュすることができる1つ以上のライトバックデータセットがその近くにある環境を提供する。好ましい実施例では、キャッシュ・マネージャ190は、任意の定められた時間に未決の30を超えるWおよび2を超えるRが無いように(待ち行列コマンドに対する合計32の「スロット」に対して)、各装置100に対するアレイ200内の「R」に対する「W」の選択された比率を維持するように動作する。しかしながら、他の比率も容易に使用することができる。比率はキャッシュ・マネージャ190により経験されるリード/ライト・コマンドミックス内のバースト変化に関して時間と共に調節することもできる。

30

【0050】

図8はライトバックデータ・フラッシング・ルーチン300を説明し、一般的に本発明の好ましい実施例に従って実施されるステップを表す。

【0051】

システムはステップ302において初期構成される。好ましくは、これはフラッシングリスト199に対するさまざまな境界および記憶装置100の物理的領域をカバーする対応するアレイ200の初期識別を含む。所望により、記憶スペースの適切なサブセットに対する異なるアレイおよびリストを維持することができ、あるいは単一の統合リスト/アレイを維持することができる。

40

【0052】

正規のシステム操作が次に開始され、ステップ304に示すように、これはキャッシュメモリへのライトバック(ダーティ)データの周期的提供を含む。このようなライトバックデータは主として図1のPC132, 134, 136等のホストからのデータライト操作から生じるものと考えられ、その場合、コントローラ124は好ましくはライトバックデータを選択されたキャッシュ場所(図4の176等)に格納し、ライト完了信号を開始装置へ返す。しかしながら、ライトバックデータはシステム状態データ、選択されたメモリバックアップ、メタデータ、等の内部発生ライトとすることもできる。

【0053】

ステップ306に示すように、SDD192は好ましくは関連するライトバックデータ

50

に対して更新される。ダーティデータおよびダーティバッファ値はデータをダーティとして識別するように初期設定することができる。その後、データセットはフルストリップ198とのデータの関係、経過期間、および付加データのキャッシュ内への進入レートを含む前記したいくつかの要因に関してフラッシングリスト199へ移される。アレイ200は前記ブロックのフラッシングリスト199への移動に回答してフラッシングに利用できる隣接データブロックのセットを識別するようにステップ308において対応的に占用される。

#### 【0054】

ステップ310において、ライトバックデータの選択されたセットは記憶装置100へフラッシュされる。好ましくは、これは他の近接アクセス操作と共に生じるが、装置100にわたってより大きい逐次フラッシング操作をスケジュールすることもできる。前記したように、CM190または他のプロセスは好ましくはキャッシュされたライトバックデータがさらにキャッシュメモリへ導入されるレートに関して利用できるライトバックデータブロックの実質的に均一な分布を維持するように動作する。

#### 【0055】

ここで検討されたさまざまな好ましい実施例は従来技術に優る利点を提供する。開示された方法論は時間と局所性の両方がフラッシング・アルゴリズムの要因として記憶装置100へのデータの効率的フラッシングを提供する点においてダイナミックである。さらに、方法論は付加ダーティデータがキャッシュメモリへ導入されるレートに実質的に揃えるのに適応的である。好ましい実施例では、キャッシュメモリ内の付加ダーティデータのセットはアレイ200へかつアレイ200から記憶ディスクへ選択的に計られて実質的に平坦な負荷量を維持する。このようにして、ホストI/O内の著しい変動が回避される。

#### 【0056】

ここに提示された好ましい実施例は複数のディスクドライブ記憶装置を利用するマルチデバイスアレイに向けられているが、それは単なる説明の目的であって制約的意味合いはない。むしろ、特許請求される発明は任意数のさまざまな環境で利用して効率的なデータ処理を促進させることができる。

#### 【0057】

本発明のさまざまな実施例の非常に多くの特徴および利点を、本発明のさまざまな実施例の構造および機能の詳細と共に、前記した明細書に記載してきたが、この詳細な説明は説明用にすぎず詳細、特に、部品の構造および配置に関して本発明の原理内で添付特許請求の範囲が表現される用語の広範な一般的意味により示される限界まで変更を行うことができる。たとえば、特定のエレメントは特定の応用に応じて本発明の精神および範囲を逸脱することなく変動することができる。

#### 【図面の簡単な説明】

#### 【0058】

【図1】本発明の好ましい実施例に従って構成かつ作動される記憶装置を一般的に示す図である。

【図2】図1に示すようないくつかの記憶装置を利用するネットワークシステムの機能的ブロック図である。

【図3】図2のコントローラの好ましいアーキテクチャの一般的表現を示す図である。

【図4】図3の選択されたインテリジェント記憶プロセッサの機能的ブロック図である。

【図5】好ましい実施例に従って記憶アレイにデータをフラッシュするように動作するキャッシュ・マネージャを一般的に示す図である。

【図6】好ましい実施例に従って記憶アレイのいくつかの異なる場所にわたってライトバックデータ機会の分布を提供する図5のキャッシュ・マネージャにより維持されるライトバックデータセットのアレイを表す図である。

【図7】もう1つの好ましい実施例に従った図6のアレイの一部を示す図である。

【図8】好ましい実施例に従って実施されるステップを示すライトバックデータ・フラッシング・ルーチンに対するフロー図である。

10

20

30

40

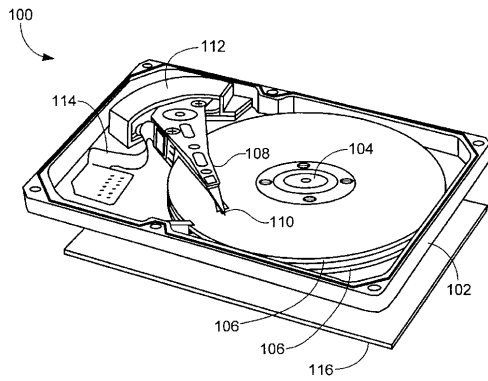
50

## 【符号の説明】

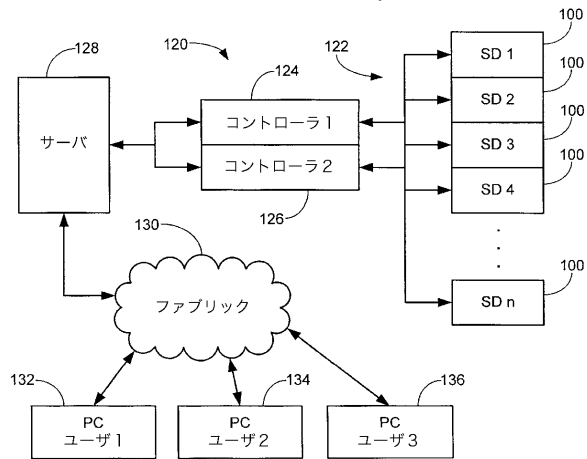
## 【0059】

100	記憶装置	
102	ベースデッキ	
104	スピンドルモータ	
106	回転媒体	
108	アクチュエータ	
110	リード/ライト・トランスデューサ	
112	ボイスコイルモータ	
114	フレックス回路アセンブリ	10
116	印刷回路板	
120	ネットワーク・システム	
122	記録アレイ	
124, 126	冗長コントローラ	
128	サーバ	
130	ファブリック	
132, 134, 136	パーソナルコンピュータ	
140	アレイコントローラ構成	
142, 144	インテリジェント記憶プロセッサ	
146	中間バス	20
148, 150	ファイバチャネル・サーバリンク	
152, 154	ファイバチャネル記憶リンク	
156, 158	ポリシー・プロセッサ	
160, 162	PCIバス	
164, 166	アクセスメモリ・モジュール	
168	ファンクション・コントローラ	
170	リストマネージャ	
172	交換メモリブロック	
174	キャッシュテーブル・ブロック	
176	バッファメモリ・ブロック	30
178	SRAM	
180	クロスポイントスイッチ・モジュール	
182	PCIインターフェイス・モジュール	
184	E-BUS IFモジュール	
190	キャッシュ・マネージャ	
192	ストライプ・データ記述子	
194, 198	RAIDストライプ	
196	SBL	
199	フラッシングリスト	
200	アレイ	40
202, 204	セル	

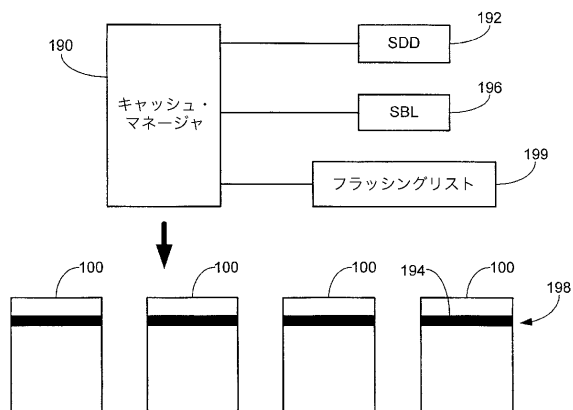
【図 1】



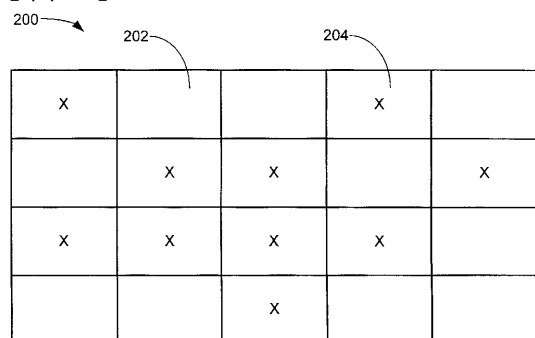
【図 2】



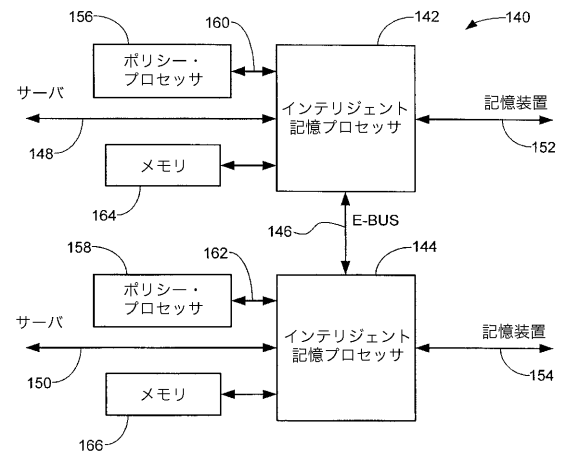
【図 5】



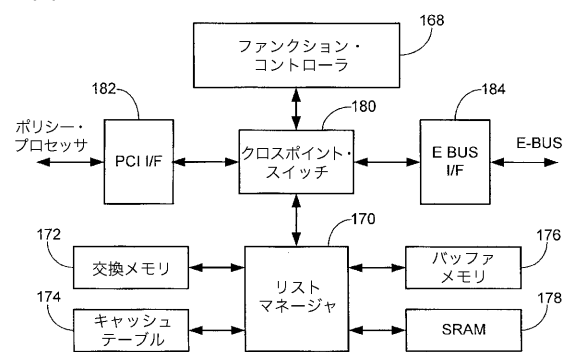
【図 6】



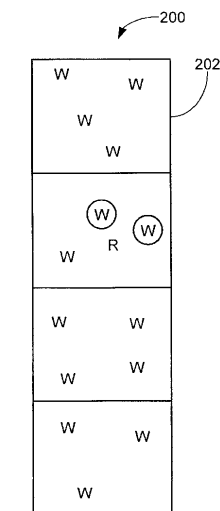
【図 3】



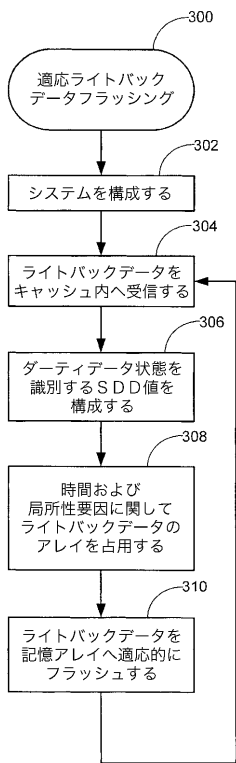
【図 4】



【図 7】



【図 8】



---

フロントページの続き

- (72)発明者 クラーク イー・レッベルス  
アメリカ合衆国、コロラド、コロラドスプリングス、 ピニョン バレー ロード 5301
- (72)発明者 マイケル ディー・ウォーカー  
アメリカ合衆国、コロラド、コロラドスプリングス、 イースト ウィラメッド 219
- (72)発明者 デーヴィッド ビー・デセンゾ  
アメリカ合衆国、コロラド、プエブロ、 ミッドナイト アベニュー 319
- F ターム(参考) 5B005 JJ12 MM11 NN02 VV01  
5B065 CA30 CH01

【外国語明細書】

2008016024000001.pdf