

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5594942号  
(P5594942)

(45) 発行日 平成26年9月24日(2014.9.24)

(24) 登録日 平成26年8月15日(2014.8.15)

(51) Int. Cl. F I  
**G06F 3/06 (2006.01)** G06F 3/06 301G  
 G06F 3/06 540

請求項の数 14 (全 15 頁)

<p>(21) 出願番号 特願2008-169143 (P2008-169143)</p> <p>(22) 出願日 平成20年6月27日 (2008.6.27)</p> <p>(65) 公開番号 特開2009-15845 (P2009-15845A)</p> <p>(43) 公開日 平成21年1月22日 (2009.1.22)</p> <p>審査請求日 平成23年6月16日 (2011.6.16)</p> <p>(31) 優先権主張番号 11/771, 980</p> <p>(32) 優先日 平成19年6月29日 (2007.6.29)</p> <p>(33) 優先権主張国 米国 (US)</p> <p>前置審査</p>	<p>(73) 特許権者 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、95014 カリフォル                  ニア州、クパチーノ、サウス・デ・アンザ                  ・ブールバード、10200</p> <p>(74) 代理人 100064746                  弁理士 深見 久郎</p> <p>(74) 代理人 100085132                  弁理士 森田 俊雄</p> <p>(74) 代理人 100083703                  弁理士 仲村 義平</p> <p>(74) 代理人 100096781                  弁理士 堀井 豊</p> <p>(74) 代理人 100111246                  弁理士 荒川 伸夫</p> <p style="text-align: right;">最終頁に続く</p>
--	---

(54) 【発明の名称】 好適なゾーン・スケジューリング

(57) 【特許請求の範囲】

【請求項1】

データ・ストレージ・システムであって、  
 ポリシーエンジンを備え、

前記ポリシー・エンジンは、前記データ・ストレージ・システムに対するネットワーク  
 負荷についての定性データを連続的に収集して前記負荷を動的に特徴付け、該特徴づけに  
 より得られた負荷の特徴に応じて、ライトバック・データに対する転送要求についてのコ  
 マンド・キューの内容を記憶媒体上にデータ記憶領域として設けられた複数の物理ゾーン  
 において物理ゾーン単位で選択されたライトバック・データに対する転送要求だけに選択  
 的に限定し、

前記ネットワーク負荷は、該ネットワークを介して前記データ・ストレージ・システム  
 に該ネットワークに結合されるホストから与えられる一連のアクセス・コマンドの少なく  
 とも種別および転送速度により表され、

前記定性データは前記負荷に含まれるアクセス・コマンドが読取りコマンドであるか書  
 込みコマンドであるかを特定するデータであり、前記負荷の特徴は少なくとも当該特定さ  
 れたアクセス・コマンドの分布を表す、データ・ストレージ・システム。

【請求項2】

前記ポリシー・エンジンにตอบสนองして、前記負荷の特徴が、前記負荷が、前記ポリシー  
 エンジンが前記ネットワーク負荷の状態の判定のための基準として用いる予想された定常状  
 態と同じか、またはそれ以下であることを示している場合に、前記複数の物理ゾーンが均

一な割合で選択されてアクセスされる均一な分布にもとづいて前記コマンド・キューの内容を設定して前記複数の物理ゾーンが均一な割合でアクセスされるよう管理し、前記負荷の特徴が、均一でない分布のアクセス・コマンド活動のバーストを示している場合に、物理ゾーンにもとづいて前記コマンド・キューの内容を設定して前記物理ゾーン単位で前記複数の物理ゾーンがアクセスされるように管理するキャッシュ・マネージャを備える、請求項 1 に記載のデータ・ストレージ・システム。

【請求項 3】

前記負荷の特徴が、前記負荷に含まれるアクセス・コマンドのレイテンシに敏感なアクセス・コマンドに対する速度に敏感なアクセス・コマンドの比率で表され、

前記レイテンシに敏感なコマンドは、コマンド処理においてレイテンシの影響が速度による影響よりも大きなコマンドであり、

前記速度に敏感なコマンドはコマンド処理において速度の影響がレイテンシによる影響よりも大きなコマンドである、請求項 1 または 2 に記載のデータ・ストレージ・システム。

【請求項 4】

前記速度に敏感なコマンドが、ライトバック・キャッシュ・コマンドであり、前記レイテンシに敏感なコマンドが、読取りコマンドおよびライトスルー・キャッシュ・コマンドのうちの少なくとも一方である、請求項 3 に記載のデータ・ストレージ・システム。

【請求項 5】

前記負荷の特徴が、各アクセス・コマンドに関連するサイズで表される、請求項 1 に記載のデータ・ストレージ・システム。

【請求項 6】

各物理ゾーンが、前記記憶媒体の全周にわたる環状部分である、請求項 1 に記載のデータ・ストレージ・システム。

【請求項 7】

独立のドライブ ( R A I D ) の冗長アレイを定義するために割り当てられた媒体を有するデータ・ストレージのアレイを備え、各物理ゾーンが、2 つ以上のデータ・ストレージにわたってストライプ化されて割り当てられたストレージ空間である、請求項 1 から 6 のいずれかに記載のデータ・ストレージ・システム。

【請求項 8】

前記コマンド・キューがシーク・マネージャに提示され、前記シーク・マネージャは、コマンド・プロファイルを定義するコマンド・キュー内の転送要求を選択的に発行する、請求項 1 から 7 のいずれかに記載のデータ・ストレージ・システム。

【請求項 9】

前記ポリシー・エンジンが、発行時のコマンドの分布を示すコマンド・プロファイルを決定する際に、シーク・マネージャを管理する規則として負荷の特徴に対応する規則および前記コマンド・キューからの転送要求の処理時に前記システムにおいて実現されるべき目標に対応する規則を決定する、請求項 1 から 8 のいずれかに記載のデータ・ストレージ・システム。

【請求項 10】

前記ポリシー・エンジンが、前記負荷に含まれるアクセス・コマンドのレイテンシに敏感なコマンドに対する速度に敏感なコマンドの比率でにおいて、該発行時のコマンドの分布を示すコマンド・プロファイルを前記負荷の特徴に選択的に対応させ、

前記レイテンシに敏感なコマンドは、コマンド処理においてレイテンシの影響が速度による影響よりも大きなコマンドであり、

前記速度に敏感なコマンドはコマンド処理において速度の影響がレイテンシによる影響よりも大きなコマンドである、請求項 1 から 9 のいずれかに記載のデータ・ストレージ・システム。

【請求項 11】

前記ポリシー・エンジンが、前記負荷に含まれるアクセス・コマンドのレイテンシに敏

10

20

30

40

50

感なコマンドに対する速度に敏感なコマンドの比率において、前記コマンド・プロファイルを前記負荷の特徴とマッチさせ、

前記レイテンシに敏感なコマンドは、コマンド処理においてレイテンシの影響が速度による影響よりも大きなコマンドであり、

前記速度に敏感なコマンドはコマンド処理において速度の影響がレイテンシによる影響よりも大きなコマンドである、請求項 9 に記載のデータ・ストレージ・システム。

【請求項 1 2】

前記ポリシー・エンジンが、前記コマンド・プロファイルを前記アクセス・コマンドの所望の最大レイテンシと選択的に対応させる、請求項 9 に記載のデータ・ストレージ・システム。

10

【請求項 1 3】

前記ポリシー・エンジンが、前記コマンド・プロファイルを異なる L U N に割り当てられた優先度を選択的に対応させる、請求項 9 に記載のデータ・ストレージ・システム。

【請求項 1 4】

前記ポリシー・エンジンが、有限状態機械である、請求項 1 に記載のデータ・ストレージ・システム。

【発明の詳細な説明】

【技術分野】

【0 0 0 1】

本発明の実施形態は、概して、データ・ストレージ・システムの分野に関し、特に、分散アレイ・ストレージ・システムにおけるシーク・コマンド・プロファイルを適合可能に管理するための装置および方法に関するが、これに限定されない。

20

【背景技術】

【0 0 0 2】

コンピュータ・ネットワークは、工業規格アーキテクチャのデータ転送速度が、インテル社 ( Intel Corporation ) の 8 0 3 8 6 プロセッサのデータ・アクセス速度に追いつくことができなくなった時に急激に増大し始めた。ローカル・エリア・ネットワーク ( L A N ) は、ネットワーク内のデータ・ストレージ容量を強化することにより、ストレージ・エリア・ネットワーク ( S A N ) に進化した。ユーザは、装置を結合し S A N 内の装置で扱われる関連データにより、直接取付ストレージにより可能となるより 1 桁上の処理能力、そして扱いやすいコストで有意の利点を実現している。

30

【0 0 0 3】

さらに最近は、データ・ストレージ・サブシステムを制御するためのネットワーク・セントリック・アプローチの方向への動きがある。すなわち、ストレージを強化したのと同じ方法で、サーバから取り出され、ネットワーク自身に送られるストレージの機能を制御するシステムにも同じ動きがある。例えば、ホスト・ベースのソフトウェアは、インテリジェント・スイッチまたは特殊化したネットワーク・ストレージ・サービス・プラットフォームに保守および管理タスクを委託することができる。アプライアンス・ベースの解決方法を使用すれば、ホストで稼働するソフトウェアが必要なくなるし、企業内にノードとして設置されているコンピュータで動作することができる。いずれにせよ、インテリジェント・ネットワーク解決方法は、これらのものをストレージ割当ルーチン、バックアップ・ルーチン、およびホストによらない障害許容スキームとして中央に集めることができる。

40

【発明の開示】

【発明が解決しようとする課題】

【0 0 0 4】

インテリジェンスをホストからネットワークに移動すればこのようないくつかの問題を解決することはできるが、仮想ストレージのプレゼンテーションをホストに変更する際の柔軟性の一般的な不足に関連する固有の問題は解決しない。例えば、データを格納する方法を、通常でないホスト負荷活動のバーストを収容するように適合させる必要がある場合

50

がある。その各データ・ストレージ容量の自己決定による割当て、管理、および保護、およびグローバルなストレージ要件に適應するように、ネットワークへその容量を仮想ストレージ空間として提示するインテリジェント・データ・ストレージ・サブシステムが求められている。この仮想ストレージ空間は、複数のストレージ・ボリューム内に提供することができる。本発明の目指しているのはこのための解法である。

【課題を解決するための手段】

【0005】

本発明の実施形態は、概して、バースト負荷条件下でコマンド・キューの内容を選択的に入手することによるダーティ・データのフラッシング性能の最適化に関する。

【0006】

ある実施形態においては、ポリシー・エンジンが、負荷を動的に特徴付けるために、データ・ストレージ・システムへのネットワーク負荷に関する定性情報を連続的に収集し、キャッシュされたライトバック・データおよびホスト読取りコマンドに対するデータ転送要求から入手するコマンド・キューの内容に負荷の特徴を連続的に相関付け、内容を記憶媒体の複数の予め定義したゾーンのゾーン・ベースで選択されるライトバック・データに対するこれらの転送要求だけに制限するデータ・ストレージ・システムおよび関連する方法が提供される。

【0007】

特許請求の範囲に記載する本発明を特徴付けるこれらおよび種々の他の機能および利点は、下記の詳細な説明を読み、関連する図面を見れば理解することができるだろう。

【発明を実施するための最良の形態】

【0008】

図1は、本発明の実施形態を含む例示としてのコンピュータ・システム100である。1つまたは複数のホスト102は、ローカル・エリア・ネットワーク(LAN)および/またはワイド・エリア・ネットワーク(WAN)106により、1つまたは複数のネットワークに取り付けられているサーバ104にネットワークで接続している。好適には、LAN/WAN106は、ワールド・ワイド・ウェブを通して通信するために、インターネット・プロトコル・ネットワークング・インフラストラクチャを使用することが好ましい。ホスト102は、多数のインテリジェント記憶素子(ISE)108のうちの1つまたは複数上に格納しているデータをルーチン的に必要とするサーバ104内に常駐しているアプリケーションにアクセスする。それ故、SAN110は、格納しているデータにアクセスするために、サーバ104をISE108に接続する。ISE108は、その内部の企業またはデスクトップ・クラスの記憶媒体により、直列ATAおよびファイバ・チャネルのような種々の選択した通信プロトコルによりデータを格納するために、データ・ストレージ容量109を提供する。

【0009】

図2は、図1のコンピュータ・システム100の一部の簡単な図面である。3つのホスト・バス・アダプタ(HBA)103は、ネットワークまたはファブリック110を介して1対のISE108(それぞれAおよびBで示す)と相互に作用する。各ISE108は、好適には、独立ドライブの冗長アレイ(RAID)として特徴付けられている一組のストレージとして、データ・ストレージ容量109上で動作することが好ましい二重化冗長制御装置112(A1、A2およびB1、B2で示す)を含む。すなわち、好適には、制御装置112およびデータ・ストレージ容量109は、種々の制御装置112が並列の冗長リンクを使用し、システム100が格納しているユーザ・データのうちの少なくともいくつかは、少なくとも一組のデータ・ストレージ容量109内の冗長フォーマットに格納されるように、障害許容配置を使用することが好ましい。

【0010】

図3は、本発明の例示としての実施形態により組み立てたISE108である。シェルフ114は、ミッドプレーン116と電氣的に接続している制御装置112に収容する形で係合するための空洞を定める。シェルフ114は、キャビネット(図示せず)内に支持

10

20

30

40

50

される。1対の複数ドライブ・アセンブリ(MDA)118は、ミッドプレーン116の同じ側面上のシェルフ114内に収容される形で係合している。ミッドプレーン116の対向側面には、非常電力供給を行うデュアル・バッテリー122、デュアル交流電源124およびデュアル・インタフェース・モジュール126が接続している。好適には、デュアル構成要素は、一方のあるいは両方のMDA118を同時に動作し、それにより構成要素が故障した場合にバックアップ保護を行うように構成することが好ましい。

#### 【0011】

図4は、それぞれが5つのデータ・ストレージ128を支持している上部隔壁130および下部隔壁132を有するMDA118の拡大分解等角図である。隔壁130、132は、ミッドプレーン116(図3)と係合するコネクタ136を有する共通の回路基板134と接続するためにデータ・ストレージ128を整合する。ラッパー138は、電磁妨害シールドを行う。MDA118のこの例示としての実施形態は、参照により本明細書に組み込むものとする譲受人に譲渡される「複数のディスク・アレイのためのキャリア装置および方法(Carrier Device and Method for a Multiple Disc Array)」という名称の米国特許第7,133,291号の主題である。MDA118のもう1つの例示としての実施形態は、本発明の譲受人に譲渡される、参照により本明細書に組み込むものとする同じ名称の米国特許第7,177,145号の主題である。他の等価の実施形態の場合には、MDA118は、密封されたエンクロージャ内に設置することができる。

#### 【0012】

図5は、本発明の実施形態と一緒に使用するのに適して、回転媒体ディスク・ドライブの形をしているデータ・ストレージ128の等角図である。動体データ記憶媒体と回転スピンドルを下記の説明のために使用するが、他の等価の実施形態の場合には、固体メモリ素子のような非回転媒体デバイスが使用される。図5の例示としての実施形態の場合には、データ記憶ディスク138は、読取り/書込みヘッド(「ヘッド」)142にディスク138のデータ記憶位置を示すためにモータ140により回転する。ヘッド142は、ディスク138の内部トラックと外部トラックとの間をヘッド142が半径方向に移動している間に、ボイス・コイル・モータ(VCM)146に応じる回転アクチュエータ144の遠い方の端部のところに支持されている。ヘッド142は、フレックス回路150を通して回路基板148に電氣的に接続している。回路基板148は、データ・ストレージ128の機能を制御する制御信号を受信し、送信することができる。コネクタ152は、回路基板148に電氣的に接続していて、データ・ストレージ128をMDA118の回路基板134(図4)と接続することができる。

#### 【0013】

図6は、制御装置112のうちの1つの図面である。制御装置112は、1つの集積回路で具体化することもできるし、必要に応じて多数の個々の回路間で分散することもできる。好適には、プログラマブル・コンピュータ・プロセッサであることを特徴とするプロセッサ154は、プログラミング・ステップ、および好適には不揮発性メモリ156(フラッシュ・メモリまたは類似物など)およびダイナミック・ランダム・アクセス・メモリ(DRAM)158内に格納している処理データにより制御を行う。

#### 【0014】

ファブリック・インタフェース(I/F)回路160は、ファブリック110を介して他の制御装置112およびHBA103と通信し、デバイスI/F回路162は、ストレージ128と通信する。I/F回路160、162および経路制御装置164は、キャッシュ166を使用するなどして、HBA103を介してネットワーク・デバイスとISE108との間でコマンドおよびデータを送るために通信経路を形成する。別々に図示してあるが、経路制御装置164およびI/F回路160、162は一体に形成することができることを理解することができるだろう。

#### 【0015】

好適には、ホスト処理機能を増大するために、ストレージ128への仮想ブロックをフ

10

20

30

40

50

ラッシュするように R A I D コンテナ・サービス ( R C S ) に要求することにより、キャッシュ・マネージャが、書込みコマンドの特定のサブセットに対してフラッシング活動を作動させるまで、仮想ブロックに対する書込みコマンドはキャッシュ 1 6 6 内にライトバック・キャッシュされ、その内部に懸案として保持される。確実に媒体を更新する R A I D アルゴリズムにより媒体の更新を行う目的で、R C S は、シーク・マネージャに特定のデータ転送を行うために要求を送るアルゴリズムを実行する。シーク・マネージャは、キャッシュされたライトバック・コマンド、およびもっと優先度の高いホスト読取りコマンドからのデータ転送要求を発行する許可を実際に与えるために、特定のストレージ 1 2 8 に対するコマンド・キューを管理する。シーク・マネージャは、実際に転送要求を発行する許可を与える関連するデータ転送を行うためのリソースを割り当てる。

10

## 【 0 0 1 6 】

I S E 1 0 8 のデータ・ストレージ容量は、データをストレージ 1 2 8 に格納する場合に、およびデータをストレージ 1 2 8 から検索する場合に、参照される論理装置の形に組織される。システム構成情報は、ユーザ・データおよび関連するパリティと、ミラー・データおよび各記憶位置間の関係を定義する。システム構成情報は、さらに、論理ブロック・アドレス ( L B A ) の用語のようなもので、データに割り当てられたストレージ容量のブロックと関連するメモリ記憶位置との間の関係を識別する。システム構成情報は、さらに、論理ブロック・アドレスにマッピングされる仮想ブロック・アドレスを定義することによる仮想化を含むことができる。

## 【 0 0 1 7 】

20

制御装置 1 1 2 アーキテクチャは、有利にスケールリングすることができる非常に機能的なデータ管理を行い、ストレージ容量の制御を行う。好適には、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を、記憶媒体上のストライプ境界、および記憶処理中ディスク・ストライプと関連するデータを格納するための専用のキャッシュ 1 6 6 内の参照データ・バッファと整合することが好ましい。

## 【 0 0 1 8 】

動作中、キャッシュ 1 6 6 は、S A N 1 1 0 により H B A 1 0 3 を通して、ユーザ・データおよび I / O 転送に関連する他の情報を格納する。要求されなかった不確かなデータを含むストレージ 1 2 8 から検索したリードバック・データを、ストレージ 1 2 8 宛のアクセス・コマンドのスケジューリングを要求する代わりに、以降の要求したデータがキャッシュ 1 6 6 から直接転送されるように、以降の「キャッシュ・ヒット」をあてにして、キャッシュ 1 6 6 内に暫くの間保持することができる。同様に、ストレージ 1 2 8 に書き込むデータが、キャッシュされるようにライトバック・キャッシュ・ポリシーが使用され、完了肯定応答が H B A 1 0 3 を介して開始ネットワーク・デバイスに返送されるが、ストレージ 1 2 8 へのデータの実際の書込みは、後の都合のよい時間にスケジューリングされる。

30

## 【 0 0 1 9 】

それ故、通常、制御装置 1 1 2 は、各エントリの状態を含むキャッシュ 1 6 6 の内容の正確な制御を維持しなければならない。このような制御は、テーブル構造に関連するアドレスを使用するスキップ・リスト配置により実行することが好ましい。スキップ・リストは、キャッシュ 1 6 6 の一部内に維持することが好ましいが、必要に応じて他のメモリ・スペースを使用することもできる。

40

## 【 0 0 2 0 】

図 7 は、経路制御装置 1 6 4 ( 図 6 ) 内に常駐するキャッシュ・マネージャ 1 7 0、R A I D コンテナ・サービス 1 7 2、ポリシー・エンジン 1 7 4、およびシーク・マネージャ 1 7 6 を示す機能ブロック図である。シーク・マネージャ 1 7 6 は 1 つしか図示していないが、ストレージ 1 2 8 に対して専用のシーク・マネージャ 1 7 6 が存在する。そのため、これらのシーク・マネージャは、ポリシー・エンジン 1 7 4 からのシーク規則に個々に応答する。

## 【 0 0 2 1 】

50

これらの機能ブロックは、ソフトウェアまたはハードウェアで実施することができる。ハードウェアで実施する場合には、ポリシー・エンジン 174 は有限状態機械であるが、これに限定されない。いずれにせよ、ポリシー・エンジン 174 は、経路 178 を介して、I/O 単位ベースでファブリック I/F 160 経由で受信したアクセス・コマンドについての定性データを連続的に収集する。ポリシー・エンジン 174 は、動的にネットワーク負荷を特徴付け、それに続けてシーク・マネージャ 176 を管理する経路 179 を介してシーク規則をその後で発行する。シーク・マネージャは、経路 180 を介してライトバック・データおよびホスト読取り要求を転送するために、データ転送要求のコマンド・キューに問い合わせ、コマンド・プロファイルを定義するために経路 182 を通してデータ転送要求を発行する許可を選択的に与える。

10

**【0022】**

ポリシー・エンジン 174 は、シーク・マネージャ 176 に対する規則を作成する際に性能の目標 188 に応じることができる。目標 188 は、速度に敏感なコマンドに対するレイテンシに敏感なコマンドの比率（ライトバック・キャッシングに対する書込みコマンドに対する読取りコマンドの比率）でのネットワーク負荷のある要因である所望のコマンド・プロファイルの強化、異なる LUN クラスに割り当てた優先度の強化、所望の読取りコマンドのレイテンシの強化のような量的なものであっても質的なものであってもよいが、これらに限定されない。ポリシー・エンジン 174 は、また、サイズ（帯域幅）のような他のものでホスト負荷を特徴付ける定性データを収集することもできる。

**【0023】**

さらに、ポリシー・エンジン 174 は、シーク・マネージャ 176 を管理する規則を作成する際にシステム状態情報 190 に応じることができる。例えば、制限なしで、電源インジケータは、ポリシー・エンジン 174 に ISE 108 がバッテリーのバックアップ電源に切り替わったことを知らせることができる。この状態の場合、ポリシー・エンジン 174 は、計画した制限付きの電力利用度に関してキャッシュ 166 を積極的にフラッシングするために付随の事態を実施する可能性がある。ポリシー・エンジン 174 は、コマンド・プロファイルをストレージ 128 に調整する場合に、シーク・マネージャ 176 を管理するシーク規則を作成する際に、アクセス・コマンド・データ転送に直接関連しない懸念中のバックグラウンド I/O 192 または I/O の状態に応じることができる。

20

**【0024】**

それ故、ポリシー・エンジン 174 は、キャッシュされたライトバック・コマンドおよびより優先度の高いホスト読取りコマンドから入手したデータ転送のコマンド・キュー内の複数のデータ転送から選択したデータ転送要求を発行する目的で、シーク・マネージャ 176 を管理するシーク規則を定義するために、負荷の特徴、目標 188、システム状態 190、およびバックグラウンド I/O の任意の組み合わせを使用することができる。

30

**【0025】**

例示としての例の場合には、ポリシー・エンジンは、レイテンシに敏感なコマンドに対する速度に敏感なコマンドの比率で、ネットワーク負荷を特徴付ける。この説明のために、ライトバック・キャッシング・スキームを仮定する。それ故、ライトバック・キャッシュ・コマンドは、速度に敏感なコマンドであると見なされる。何故なら、任意の時点でデータ・ストレージ 128 にどの要求をフラッシングするかは大した問題ではないからである。実際には、ダーティ・データとしてキャッシュ 166 内で未決状態である場合に、速度に敏感な要求を上書きすることすらできるからである。問題は、速度に敏感なコマンドを、キャッシュ 166 が飽和状態になることを防止する速度でフラッシングすることである。

40

**【0026】**

一方、1つまたは複数のデータ・ストレージ 128 内に格納しているデータを読み出すためのアクセス・コマンドは、同様に、ネットワーク・アプリケーションが、アクセス・コマンドが満足するまでそれ以上の処理を阻止する恐れがある。アクセス・コマンドを満足させる時間、すなわち、レイテンシ期間は、アプリケーションの性能にとって非常に重

50

要なものである。そのため、このようなコマンドは、レイテンシに敏感なコマンドと呼ばれる。ある状況の場合には、ホストは、ライトバック・キャッシングを許可しないことを選択することができる。この場合、ライトスルー・キャッシュ・コマンドと呼ばれる書込みコマンドは、同様にレイテンシに敏感なコマンドとして分類される。

#### 【0027】

定性データを収集する際に、ポリシー・エンジン174は、1秒の各間隔のような、しかしこれに限定されない所定の各サンプリング期間中にカウントを照合することが好ましい。書込みコマンドに対する読取りコマンドの比率でデータを収集するために、例えば、フリーランニング・カウンタは、連続的に上記比率を追跡するために1秒刻みで指針を移動させるポイントにより設定することができる。カウンタは、9番目のスロットで現在の1秒の比を計算するため前の8回の1秒サンプル比のような所望の数の前に観察した比率を保持する。1秒刻みの目盛の上では、指針が回転し、指し示した履歴値を減算し、最近のサンプル値を加算し、次に、比率の最近の移動平均を計算するために8で除算を行う。

10

#### 【0028】

さらに、例示としての例について説明すると、ポリシー・エンジン174が、ホスト負荷バースト活動が発生していることを観察している場合には、ポリシー・エンジンは、キャッシュ166内での飽和状態を防止する際に、ネットワーク負荷に関連してコマンド・プロファイルを修正するために、シーク・マネージャ176を管理する規則を発行することができる。この規則は、(レイテンシに敏感な)読取りコマンドに対する(速度に敏感な)書込みコマンドの比率でコマンド・プロファイルをネットワーク負荷にマッチさせることができる。この規則は、また飽和状態からできるだけ速く滑らかに回復するために、最大レイテンシおよびLUNクラス優先度のような他の目標を修正することもできるし、または一時的に延期することさえできる。

20

#### 【0029】

このような状態の場合、シーク・マネージャ176に書込みコマンドを積極的にフラッシングすることが強く求められた場合には、特にストレージ128内のディスク138(図5)の複数の物理ゾーンの選択した物理ゾーン内に位置するものだけへのコマンド・キュー内の転送要求を入手するために、負荷特徴に関連するキャッシュ・マネージャ170の機能を変更することにより、フラッシング性能を有意に改善するという決定が行われる。

30

#### 【0030】

図8は、ストライプ・データ記述子(SDD)とも呼ばれるデータ構造を使用するノード・ベースでのキャッシュ166(図6)を管理するキャッシュ・マネージャ170の略図である。各SDDは、それが関連するデータへの最近および現在のアクセスに関連するデータを保持する。各SDDは、アクセス履歴、ロックした状態、最後のオフセット、最後のブロック、タイムスタンプ・データ(時刻、TOD)、データがどのゾーン(ブック)に属するのかわを示す識別子、および使用するRAIDレベルを含むデータの種々の状態を示す変数を含むことが好ましい。好適には、SDDに関連するデータのライトバック(「ダーティ」データ)状態は、ダーティ・データ、ダーティ・バッファ、ダーティLRUおよびフラッシングLRU値と関連して管理することが好ましい。

40

#### 【0031】

各SDDは、対応するRAIDストライプ194(すなわち、特定のパリティ・セットと関連する選択したストレージ128上のすべてのデータ)と整合し、特定のストライプ・バッファ・リスト(SBL)に適合することが好ましい。各キャッシュ・ノードは、順方向および逆方向にリンクしているリストを使用して、仮想ブロック・アドレス(VBA)を介して昇順にリンクしている所与の組の論理ディスクに対する能動SDD構造により、いくつかの特定のSDDを参照することが好ましい。

#### 【0032】

好適には、VBA値は、RAID割当てグリッド・システム(RAGS)と呼ばれるグリッド・システムを使用して、RAIDデータ組織と整合される。通常、同一のRAID

50

ストリップ 196 (例えば、特定のパリティ・セットに貢献するすべてのデータなど) に属するブロックの任意の特定の集合体が、特定のシート上の特定の信頼性のあるストレージ・ユニット (RSU) に割り当てられる。ブックはいくつかのシートからできていて、異なるストレージ 128 からのブロックの複数の隣接する組から作られる。実際のシートおよび VBA に基づいて、このブックを、さらに、(冗長性を使用する場合) 特定のデバイスまたはデバイスの組を示すゾーンに再分割することができる。

#### 【0033】

好適には、キャッシュ・マネージャ 170 は、ライトバック・コマンドおよびホスト読取りコマンドをできるだけ効率的に処理するために適合できるように管理することが好ましい。ダーティ・ブロックおよびロックしていないブロックを含む任意の SDD は、ダーティとしてセットし、古さ (例えば、キャッシュ待機フラッシング中にデータが消費した時間など) で分類することが好ましい。特定の古さになると、フラッシング LRU 変数をセットすることができ、フラッシング・リストが更新される。

10

#### 【0034】

図 9 は、コマンド・キュー内の転送要求のレイ 198 である。これらの要求は、記憶容量およびブック 200 が存在するストレージ 128 のブック 200 の交点により定義されるセルによりマッピングされる。セルは、X で示すこれらのセルが 1 つまたは複数の転送要求が懸案中であるこれらの位置を示すように、転送要求により配置される。レイ 198 は、読取りコマンドを実行するために、アクセス・コマンドが特定の位置にトランスジューサ 142 (図 5) を送った場合に使用されるシーク管理ツールであり、ライトバック・データの他の使用可能なブロックは、アクチュエータ 144 (図 5) をショート・ストロークする際に現在の位置に関連してサービスを行うために容易に識別することができる。

20

#### 【0035】

ダーティ・データに対する転送要求は、シーク・マネージャ 176 に均一な分布が提示されるように、レイ 198 にランダムに追加される。コマンド・キューに対するこの均一に分布しているベースは、ネットワーク負荷が予想した定常状態である場合に、またはそれ以下である場合に使用することが好ましい CSCAN およびエレベータ・ルーチンのような従来のシーク・スケジューリング・ルーチンにうまく適している。図 10 は、例えば、ヘッド 142 が、どのように移動の予め定義した経路に沿って均一に分布している転送要求を発行しているシーク・マネージャ 176 と一緒に、内部の半径 202 から半径方向に外側に向かってディスク 138 を横断し、次に、外側の半径 204 から半径方向に内側に横断するのを示す。

30

#### 【0036】

しかし、本実施形態の場合には、コマンド・キューからのデータのフラッシングの積極性は、アクセス・コマンドのホスト負荷に結びついている。すなわち、比較的大きなホスト負荷がかかっている間に十分積極的にフラッシングしないと、キャッシュ 126 が飽和する恐れがある。それ故、均一でない分布による負荷バースト活動の間にホスト負荷の特徴にコマンド・プロファイルをマッチさせるためにシーク・マネージャ 174 に重きを置くポリシー・エンジン 174 と一緒に、ポリシー・エンジン 174 は、また、それにより自身が均一な分布によるのではなく、物理ゾーンによりコマンド・キューを定義するモードへのキャッシュ・マネージャ 170 の切替も管理する。

40

#### 【0037】

図 11 は、図 9 のレイ 198 を形成しているストレージ 128 のディスク 138 のうちの 1 つである。これらのストレージ 128 のこのディスク 138 およびすべての他のディスク 138 は、記憶容量のブック 200 により定義される環状ゾーン内に区分される。図 12 は、各ブック 200 内で現在懸案中のライトバック・データに対する転送要求の数の一例である。この実施形態によれば、ポリシー・エンジン 174 は、この場合はブック 2004 である、最大数の懸案中のライトバック要求を含むブック 200 を選択するためにキャッシュ・マネージャ 170 を選択的に管理し、コマンド・キューの内容をブック 2

50

004内のこれらのライトバック要求だけに制限する。ある実施形態の場合には、キャッシュ・マネージャ170は、すべてのライトバック要求がそのブック200から発行されるまで、該最多のブック200に制限することができる。他の実施形態の場合には、キャッシュ・マネージャ170は、周期的にライトバック要求の分布を再評価することができ、異なるブック200がかなり大きい数のライトバック要求を有している場合には、キャッシュ・マネージャは、コマンド・キューの内容を該最多のブック200内のライトバック要求にシフトする。

【0038】

図13は、本発明の実施形態による好適なゾーン・スケジューリングのための方法202のステップを示すフローチャートである。この方法202は、ライトバック・コマンドおよびホスト読取りコマンドに対する転送要求の均一な分布によるコマンド・キューをローディングするデフォルト・モードでブロック204からスタートする。デフォルト・モードは、ポリシー・エンジンが、ネットワーク負荷に関するデータを収集する1秒の間隔のような、しかしこれに限定されない予め定義した間隔中に実施される。最新のデータは、例えば、書込みに対する読取りの比率およびI/O速度で、ホスト負荷を動的に特徴付けるためにブロック206で使用される。

【0039】

ブロック208においては、ポリシー・エンジンは、均一でない分布を含むI/Oコマンドのバーストがネットワーク負荷を監視することによりはつきりわかるか否かを判定する。ブロック208における判定が「いいえ」である場合には、制御はブロック204に戻り、そのためデフォルト状態が継続する。しかし、ブロック208における判定が「はい」である場合には、ブロック210において、ポリシー・エンジンは、コマンド・プロファイルの入手の際に連続的に調整する目的で、ホスト負荷の特徴、およびおそらく目標188、システム状態190およびバックグラウンドI/O192を使用して好適なゾーン・スケジューリング規則を呼び出す。例えば、制限無しで、飽和状態でコマンドを読み出すための高い書込みコマンドが発生している場合には、ポリシー・エンジンは、飽和状態から回復するまで、読取りに対する書込みの比率でコマンド・キューの内容をホスト負荷にマッチさせるためにキャッシュ・マネージャを管理することができる。ポリシー・エンジンは、飽和状態からできるだけ速くまた滑らかに回復するために、読取りレイテンシおよびLUNクラス優先度のような他の規則を修正することさえできるし、または一時的に中止することさえできる。好適なゾーン・スケジューリング規則は、負荷データの次のバッチが収集される1秒の間隔のような所定の間隔中に呼び出され、制御は、ブロック206に戻る。

【0040】

通常、この実施形態は、ネットワーク・アクセス・コマンドに応じてデータを転送するためにネットワークに接続するように構成されているストレージ・アレイ、およびコマンド・キュー内の転送要求の分布をアクセス・コマンドのネットワーク負荷の観察した特徴に相関付けることによりコマンド・キューを入手するための手段を予想している。この説明および添付の特許請求の範囲の意味のために、「入手するための手段」という用語は、明らかに、本明細書に記載する構造、および制御装置112がネットワーク負荷を特徴付け、特徴によりコマンド・キューの内容を直接調整することができるようにするその等価物を含む。フラッシング・リストを「直接」調整することにより、「入手するための手段」は、明らかに、ネットワーク負荷の特徴およびポリシー・エンジンからの規則に関連して、均一な分布からゾーン分布へのようなこれに応じて内容を切り替えるキャッシュ・マネージャを予測している。この説明および添付の特許請求の範囲の場合、「入手するための手段」は、ライトバック・データおよびコマンド・キューの内容に間接的に影響を与える恐れがある高度に局所化したバースト活動の周期となるPOS (point of sale) 取引のようなその発生の一致する位置を単に予想していない。

【0041】

上記説明内で、本発明の種々の実施形態の構造および機能の詳細と一緒に、本発明の種

10

20

30

40

50

々の実施形態の多くの特徴および利点を説明してきたが、この詳細な説明は、例示としてのためだけのものであって、添付の特許請求の範囲を説明している用語の広い一般的な意味により示す全範囲に、特に本発明の原理の部材の構造および配置を変えることができることを理解されたい。例えば、本発明の精神および範囲から逸脱することなしに、特定の処理環境により特定の要素を変えることができる。

【0042】

さらに、本明細書に記載する実施形態は、データ・ストレージ・アレイに関するものであるが、当業者であれば、特許請求の範囲に記載の主題は、それに限定されるものではなく、本発明の精神および範囲から逸脱することなしに、種々の他の処理システムも使用することができることを理解することができるだろう。

10

【図面の簡単な説明】

【0043】

【図1】本発明の実施形態を組み込むコンピュータ・システムの図面である。

【図2】図1のコンピュータ・システムの一部の簡単な図面である。

【図3】本発明の実施形態によるインテリジェント記憶素子の分解等角図である。

【図4】図3のインテリジェント記憶素子の複数のドライブ・アレイの分解等角図である。

【図5】図4の複数のドライブ・アレイで使用する例示としてのデータ・ストレージである。

【図6】インテリジェント記憶素子内のアレイ制御装置の機能ブロック図である。

20

【図7】図6のアレイ制御装置の一部の機能ブロック図である。

【図8】インテリジェント・ストレージ素子内のアレイ制御装置およびデータ・ストレージのアレイの一部を示す略図である。

【図9】ライトバック・データ要求をランダムにマッピングする際にキャッシュ・マネージャが使用するアレイである。

【図10】CSCANシーク・スケジューリング技術と一緒に、均一に分布しているライトバック・データ要求を発行するストレージの略図である。

【図11】データ・ストレージ容量のブックにより定義される環状のゾーンを示すストレージ・ディスクの簡単な平面図である。

【図12】図11の各ブック内のライトバック要求の観察したカウントを示す。

30

【図13】本発明の実施形態による好適なゾーン・スケジューリングのための方法のステップを示すフローチャートである。

【符号の説明】

【0044】

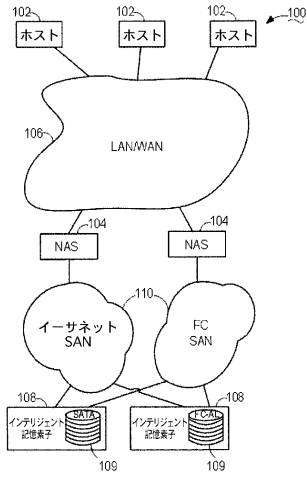
- 100 コンピュータ・システム
- 102 ホスト
- 103 ホスト・バス・アダプタ(HBA)
- 104 サーバ
- 106 LAN/WAN
- 108 インテリジェント記憶素子(ISE)
- 109 データ・ストレージ容量
- 110 SAN
- 112 二重化冗長制御装置
- 114 シェルフ
- 116 ミッドプレーン
- 118 ドライブ・アセンブリ(MDA)
- 122 デュアル・バッテリー
- 124 デュアル交流電源
- 126 デュアル・インタフェース・モジュール
- 128 データ・ストレージ

40

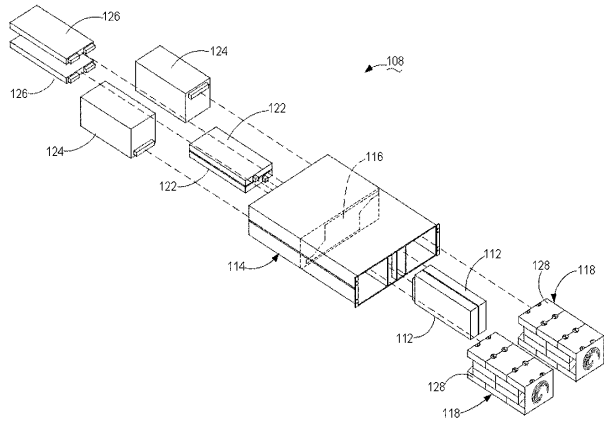
50

1 3 0 , 1 3 2	隔壁	
1 3 4	回路基板	
1 3 6	コネクタ	
1 3 8	ラッパ	
1 4 0	モータ	
1 4 2	読取り / 書込みヘッド	
1 4 4	回転アクチュエータ	
1 4 6	ボイス・コイル・モータ ( V C M )	
1 4 8	回路基板	
1 5 0	フレックス回路	10
1 5 2	コネクタ	
1 5 4	プロセッサ	
1 5 6	不揮発性メモリ	
1 5 8	ダイナミック・ランダム・アクセス・メモリ ( D R A M )	
1 6 0 , 1 6 2	I / F 回路	
1 6 4	経路制御装置	
1 6 6	キャッシュ	
1 7 0	キャッシュ・マネージャ	
1 7 2	R A I D コンテナ・サービス	
1 7 4	ポリシー・エンジン	20
1 7 6	シーク・マネージャ	
1 7 8 , 1 8 0 , 1 8 2	経路	
1 8 8	目標	
1 9 0	システム状態	
1 9 8	アレイ	
2 0 0	ブック	

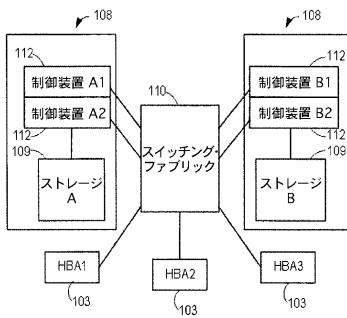
【図1】



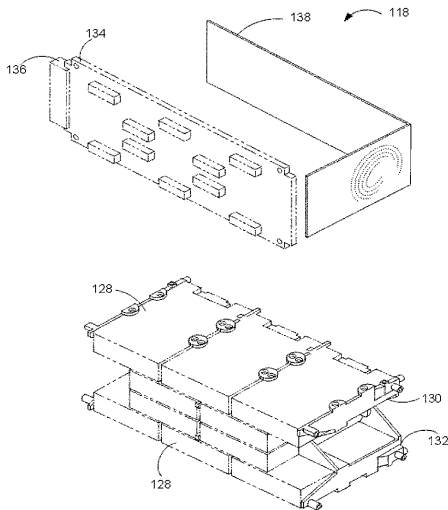
【図3】



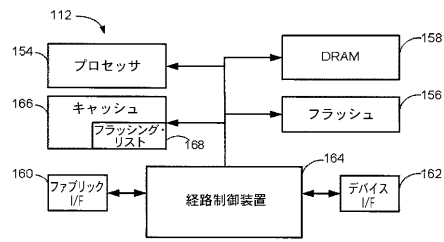
【図2】



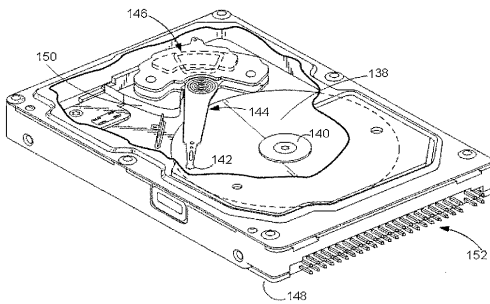
【図4】



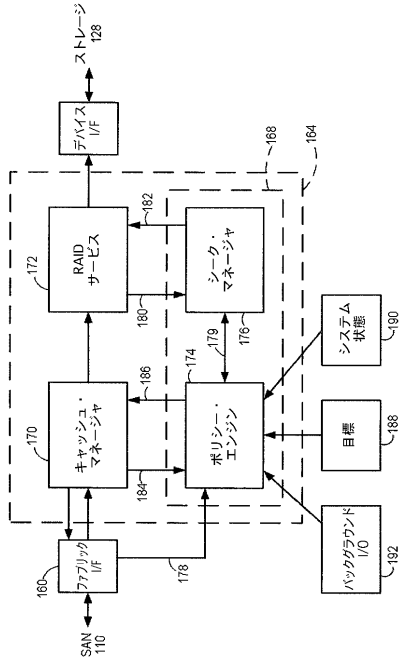
【図6】



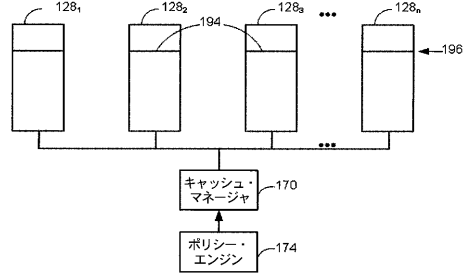
【図5】



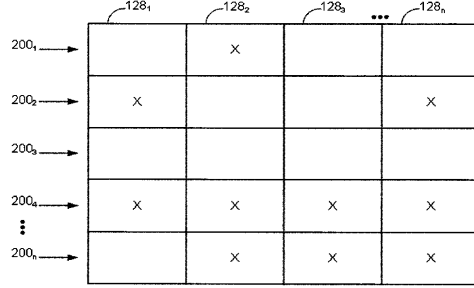
【図 7】



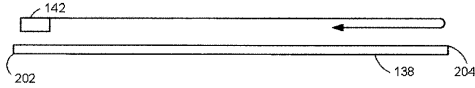
【図 8】



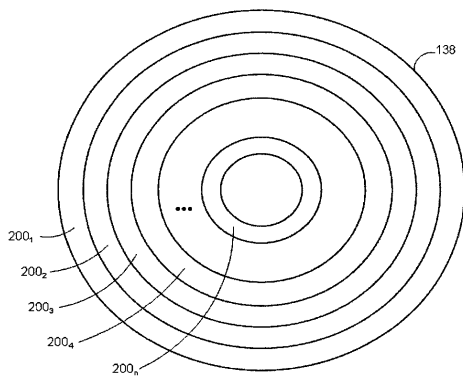
【図 9】



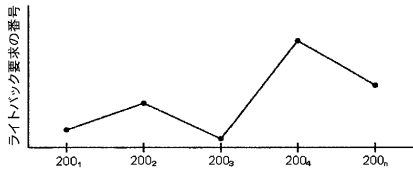
【図 10】



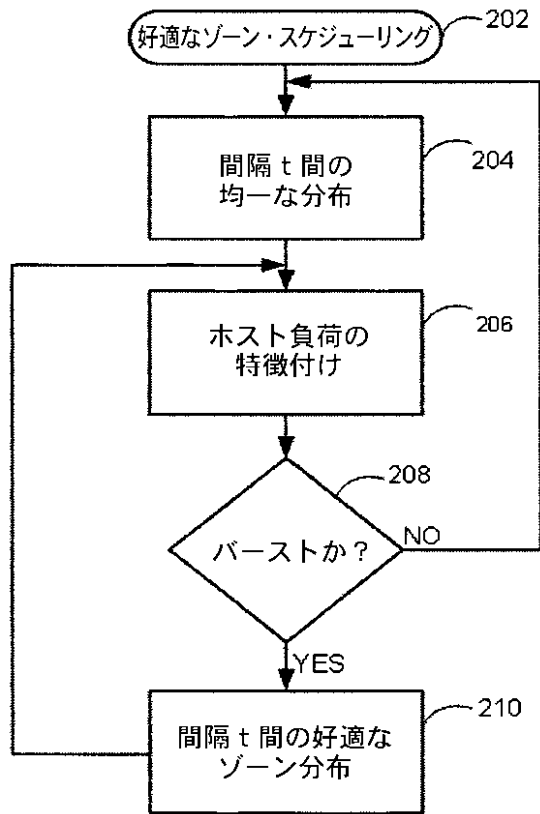
【図 11】



【図 12】



【図 13】



## フロントページの続き

(74)代理人 100124523

弁理士 佐々木 真人

(72)発明者 クラーク エドワード ルッベルス

アメリカ合衆国、コロラド、コロラド スプリングス、 ピニオン バレイ ロード 5301

(72)発明者 ロバート マイケル レスター

アメリカ合衆国、コロラド、コロラド スプリングス、 ロシヨルト ループ 14710

審査官 坂東 博司

(56)参考文献 特開平08-115169(JP,A)

特開2004-295860(JP,A)

特開2006-120118(JP,A)

米国特許出願公開第2002/0156972(US,A1)

米国特許出願公開第2005/0076115(US,A1)

米国特許出願公開第2006/0253621(US,A1)

特開平06-110772(JP,A)

特開平05-257614(JP,A)

特開2007-115233(JP,A)

Binny S.Gill and Dharmendra S. Modha, Proceedings of the 4th conference on USENIX Conference on File and Storage Technologies, WOW: wise ordering of writes - combining spatial and temporal locality in non-volatile caches, USENIX Association, 2005年12月13日, 129~142, URL, [http://static.usenix.org/event/fast05/tech/full\\_papers/gill/gill.pdf#search='WOW Wise Ordering for Writes Combining Spatial and Temporal Locality in NonVolatile Caches'](http://static.usenix.org/event/fast05/tech/full_papers/gill/gill.pdf#search='WOW Wise Ordering for Writes Combining Spatial and Temporal Locality in NonVolatile Caches')

(58)調査した分野(Int.Cl., DB名)

G06F 3/06