

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5376624号  
(P5376624)

(45) 発行日 平成25年12月25日 (2013. 12. 25)

(24) 登録日 平成25年10月4日 (2013.10. 4)

(51) Int. Cl.		F I			
<b>G06F 13/10</b>	<b>(2006.01)</b>	G06F 13/10	340A		
<b>G06F 3/06</b>	<b>(2006.01)</b>	G06F 3/06	301A		
		G06F 3/06	301F		

請求項の数 20 (全 16 頁)

(21) 出願番号	特願2008-166092 (P2008-166092)	(73) 特許権者	500373758
(22) 出願日	平成20年6月25日 (2008. 6. 25)		シーゲイト テクノロジー エルエルシー
(65) 公開番号	特開2009-9572 (P2009-9572A)		アメリカ合衆国、95014 カリフォル
(43) 公開日	平成21年1月15日 (2009. 1. 15)		ニア州、クパチーノ、サウス・デ・アンザ
審査請求日	平成23年6月7日 (2011. 6. 7)		・ブルバード、10200
(31) 優先権主張番号	11/768, 850	(74) 代理人	100064746
(32) 優先日	平成19年6月26日 (2007. 6. 26)		弁理士 深見 久郎
(33) 優先権主張国	米国 (US)	(74) 代理人	100085132
			弁理士 森田 俊雄
		(74) 代理人	100083703
			弁理士 仲村 義平
		(74) 代理人	100096781
			弁理士 堀井 豊
		(74) 代理人	100109162
			弁理士 酒井 将行

最終頁に続く

(54) 【発明の名称】 ホスト適応シーク技術環境

(57) 【特許請求の範囲】

【請求項 1】

データ・ストレージ・システムにおけるコントローラであって、前記コントローラは、前記コントローラへのネットワーク負荷についての該ネットワーク負荷の性質を表すデータである定性情報を連続的に収集し該収集された定性情報に従って前記負荷を動的に特徴づけ該特徴づけにより得られる前記負荷の動的特徴を策定するポリシー・エンジン論理を含むホスト適応シーク技術環境 (H A S T E) モジュールを含み、前記ポリシー・エンジン論理は、キャッシュされたライトバック・データを含むフラッシング・リスト内の複数の入出力 (I / O) 要求から選択した I / O 要求を発行するシーク・マネージャを支配しているホスト適応シーク技術環境規則である H A S T E 規則を定義する際に、前記負荷の前記動的特徴および予め定義した前記システムにおいて実現すべきシーク性能の目標値を示す性能目標を使用し、前記シーク・マネージャは、前記 I / O 要求の発行により、前記データ・ストレージ・システム内のデータ・ストレージ・デバイスへ前記リスト内の I / O 要求をフラッシュする、コントローラ。

【請求項 2】

前記動的特徴は、該ネットワーク上のコマンドの転送速度よりも前記データ・ストレージ・システムのレイテンシの影響がより大きいレイテンシに敏感なコマンドに対する前記ネットワーク上のコマンドの転送速度の影響が前記データ・ストレージ・システムのレイテンシの影響よりも大きい速度に敏感なコマンドの比率に関するものであり、前記レイテンシは、アクセスコマンドの指示が満たされるまでに要する期間である、請求項 1 に記載

のコントローラ。

【請求項 3】

前記速度に敏感なコマンドは、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は、読取りコマンドおよびライトスルー・キャッシュ・コマンドのうち少なくとも 1 つである、請求項 2 に記載のコントローラ。

【請求項 4】

前記動的特徴は、各コマンドと関連する帯域幅に関するものである、請求項 1 に記載のコントローラ。

【請求項 5】

前記 H A S T E 規則は、該ネットワーク上のコマンドの転送速度よりも前記データ・ストレージ・システムのレイテンシの影響がよりレイテンシに敏感なコマンドに対する前記ネットワーク上のコマンドの転送速度の影響が前記データ・ストレージ・システムのレイテンシの影響よりも大きい速度に敏感なコマンドの比率に関して、前記負荷にコマンドの分布または配列を示すコマンド・プロファイルを選択的にマッチさせ、前記レイテンシは、アクセスコマンドの指示が満たされるまでに要する期間である、請求項 1 に記載のコントローラ。

10

【請求項 6】

前記 H A S T E 規則は、前記コマンドに割り当てられた論理ユニット番号 ( L U N ) 優先度の強制的な設定に関連する、請求項 2 に記載のコントローラ。

【請求項 7】

前記 H A S T E 規則は、前記コマンドの読取りコマンドのレイテンシの強制設定に関連する、請求項 2 に記載のコントローラ。

20

【請求項 8】

個別データストレージデバイスのアレイに結合されており、前記ポリシー・エンジン論理からの前記 H A S T E 規則に個々に応じる前記アレイ内の各データ・ストレージデバイス専用のシーク・マネージャを含む、請求項 1 に記載のコントローラ。

【請求項 9】

ストレージ・システムへのネットワーク負荷のコマンド・ストリームを監視し、該コマンド・ストリーム内のコマンドについての該コマンドの性質を表すデータである定性情報を収集するステップと、

30

前記収集された定性情報に従って前記負荷を動的に特徴づけ、該特徴づけにより得られる前記負荷の動的特徴を策定するステップと、

前記動的特徴を使用して、前記ネットワークと前記ストレージ・システム間のキャッシュされたライトバック・データの通信を含むフラッシング・リスト内の複数の I / O 要求から選択した I / O 要求を発行するための、前記動的特徴に関連するホスト適応シーク技術環境規則である H A S T E 規則を生成するステップとを含み、前記 I / O 要求の発行により前記ストレージ・システム内のストレージデバイスへ前記リスト内の I / O 要求がフラッシュされる、方法。

【請求項 10】

前記使用するステップは、前記ストレージ・システムに対する予め定義した前記システムにおいて実現すべきシーク性能の目標値を示す性能目標を考慮した規則を使用するステップを含む、請求項 9 に記載の方法。

40

【請求項 11】

前記監視するステップは、前記負荷を、該ネットワーク上のコマンドの転送速度よりも前記ストレージ・システムのレイテンシの影響がより大きいレイテンシに敏感なコマンドに対する前記ネットワーク上のコマンドの転送速度の影響が前記データ・ストレージ・システムのレイテンシの影響よりも大きい速度に敏感なコマンドの比率の点から特徴づけるステップを含み、前記レイテンシは、アクセスコマンドの指示が満たされるまでに要する期間である、請求項 9 に記載の方法。

【請求項 12】

50

前記監視するステップは、ライトバック・キャッシュ・コマンドが前記速度に敏感なコマンドであり、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が、前記レイテンシに敏感なコマンドであると特徴づけるステップを含む、請求項 11 に記載の方法。

【請求項 13】

前記監視するステップは、前記負荷を、各アクセスコマンドに関連する帯域幅の点から特徴づけるステップを含む、請求項 9 に記載の方法。

【請求項 14】

前記使用するステップは、前記 H A S T E 規則を、該ネットワーク上のコマンドの転送速度よりも前記ストレージ・システムのレイテンシの影響がより大きいレイテンシに敏感なコマンドに対する前記ネットワーク上のコマンドの転送速度の影響が前記データ・ストレージ・システムのレイテンシの影響よりも大きい速度に敏感なコマンドの比率に関連づけるステップを含み、前記レイテンシは、アクセスコマンドの指示が満たされるまでに要する期間である、請求項 9 に記載の方法。

10

【請求項 15】

前記使用するステップは、該ネットワーク上のコマンドの転送速度よりも前記ストレージ・システムのレイテンシの影響がより大きいレイテンシに敏感なコマンドに対する前記ネットワーク上のコマンドの転送速度の影響が前記ストレージ・システムのレイテンシの影響よりも大きい速度に敏感なコマンドの比率に関して、前記負荷にコマンドの分布または配列を示すコマンド・プロファイルを選択的にマッチさせるステップを含み、前記レイテンシは、アクセスコマンドの指示が満たされるまでに要する期間である、請求項 9 に記載の方法。

20

【請求項 16】

前記使用するステップは、前記規則を、前記コマンドの論理ユニット番号 ( L U N ) クラス優先度の強制的な設定に関連づけるステップを含む、請求項 11 に記載の方法。

【請求項 17】

前記使用するステップは、前記規則を、前記コマンドの読取りコマンドに対するレイテンシの強制設定に関連づけるステップを含む、請求項 11 に記載の方法。

【請求項 18】

ストレージ・システムへのネットワーク負荷のコマンド・ストリーム内のコマンドについての該コマンドの性質を表すデータである定性情報を収集するステップと、

30

前記収集された定性情報を使用して、読取りコマンドに対する書込みコマンドの比率に関して前記負荷を動的に特徴付けるステップと、

前記ネットワークと前記ストレージ・システム間のキャッシュされたライトバック・データの通信を含むフラッシング・リスト内の複数の I / O 要求から選択した I / O 要求を発行するための、前記動的に特徴付けるステップにより特徴づけられた動的特徴に関連するホスト適応シーク技術環境規則である H A S T E 規則を生成するステップとを含み、前記 I / O 要求の発行により前記ストレージ・システム内のストレージデバイスへ前記リスト内の I / O 要求がフラッシュされる、方法。

【請求項 19】

40

前記使用するステップは、前記負荷の前記動的な特徴付けにより得られる動的特徴が予め定められたしきい値比率未満の場合に第 1 の規則を生成し、前記負荷の前記動的特徴が前記予め定められたしきい値よりも大きい場合に異なる第 2 の規則を生成する、請求項 18 に記載の方法。

【請求項 20】

前記使用するステップは、前記負荷の前記動的な特徴に、複数の選択された I / O 要求をマッチさせるステップを含む、請求項 18 に記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

50

本発明の実施形態は、概して、分散データ・ストレージ・システムの分野に関し、特に、分散レイ・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。

【背景技術】

【0002】

コンピュータ・ネットワーキングは、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 80386 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 (LAN) 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 (SAN) に進化した。ユーザは、装置を結合し、SAN 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以上の桁の処理能力、そして扱いやすいコストで有意の利点を実現している。

10

【0003】

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。

20

【発明の開示】

【発明が解決しようとする課題】

【0004】

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の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。

30

【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、分散データ・ストレージ・システムのホスト適応シーク技術環境 (HASTE) に関する。

40

【0006】

ある実施形態の場合には、データ・ストレージ・システムおよび関連する方法は、負荷を動的に特徴付けるために、データ・ストレージ・システムにネットワーク負荷についての定性情報を連続的に収集し、コマンド・プロファイルを、特徴に関連するデータ・ストレージ・システムのデータ・ストレージに連続的に相関付けるポリシー・エンジンにより HASTE を実行する。

【0007】

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】

50

## 【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の必要とするサーバ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。

10

## 【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、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。

20

## 【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。

30

## 【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。

40

## 【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

50

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。

10

**【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】**

20

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ラッシュするようにRAIDコンテナ・サービス(RCS)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166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RAIDアルゴリズムにより媒体の更新を行う目的で、RCS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128に対するコマンド・キューを管理する。特定の時点でいずれのコマンドを発行するのかの選択は、経路制御装置164内に常駐するホスト適応シーク技術環境(HASTE: Host Adaptive Seek Technique Environment)モジュール168と協働してシーク・マネージャが行う。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。

30

40

**【0016】**

ISE108のデータ・ストレージ容量は、データをストレージ128に格納する場合に、およびデータをストレージ128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(LBA)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

50

による仮想化を含むことができる。

【 0 0 1 7 】

制御装置 1 1 2 アーキテクチャは、有利にスケー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。

【 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。

【 0 0 2 0 】

キャッシュ 1 6 6 は、ストライプ・データ記述子 ( S D D ) と呼ばれるデータ構造を使用して、制御装置 1 1 2 によりノード・ベースで管理される。各 S D D 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 S D D は、対応する R A I D ストライプ ( すなわち、特定のパリティ・セットに関連する選択したストレージ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 ( S B L ) に適合することが好ましい。

【 0 0 2 1 】

制御装置 1 1 2 により管理される各キャッシュ・ノードは、好適には、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 ( V B A ) を通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 S D D 構造によりいくつかの特定の S D D を参照することが好ましい。

【 0 0 2 2 】

好適には、 V B A の値は、 R A I D 割当てグリッド・システム ( R A G S ) とも呼ばれるグリッド・システムを使用して、 R A I D データ組織と整合される。通常、同じ R A I D ストリップ ( 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できるストレージ・ユニットに割り当てられる。ブックは多数のシートからできていて、異なるストレージ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 V B A 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分割することができる。

【 0 0 2 3 】

各 S D D 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 ( 時刻、 T O D ) 、データがいずれのゾーン ( ブック ) に属するのかわを示す識別子、および使用する R A I D 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。 S D D に関連するライトバック ( 「ダーティ」データ状態は、ダーティ・デ

10

20

30

40

50

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の値と関連して管理することが好ましい。

#### 【0024】

制御装置112は、システム要件により、多数の異なるレベルのところでのライトバック・データ・プロセスを管理するために同時に動作することが好ましい。第1のレベルは、通常、全RAIDストリップが検出された場合に、全SDD構造の周期的フラッシングを含む。このことは、SDDが関連するデータをダーティと識別した場合に、RAIDレベル変数に基づいて所与のSDDに対して容易に行うことができる。好適には、このことは、十分な連続している隣接SDDが、十分ダーティなデータで満たされているか否かを判定するために逆方向のチェックを含む。そうである場合には、これらのSDD構造は、

10

#### 【0025】

もっと小さな組のデータのフラッシングもSDDをベースとして処理することが好ましい。ダーティ・ブロックおよびロックされていないブロックを含む任意のSDDは、ダーティLRUとしてセットし、古さの程度（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）により区分けすることが好ましい。特定のエイジングに達した場合には、フラッシングLRU変数を設定し、コマンド・キューを更新することが好ましい。

#### 【0026】

連続しているダーティ・ブロックの特定の範囲がフラッシングに対してスケジューリングされると、制御装置112は、最も近い位置を有するRAIDレベルに基づいて、ダーティ・ブロックの他の範囲、すなわち、シーク時間の点で「近い」ブロック、または同じRAIDパリティ・ストリップへのアクセスを含むブロックを配置することが好ましい。

20

#### 【0027】

この実施形態によれば、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、I/Oコマンドのホスト負荷と結びついている。すなわち、比較的大きな負荷がかかっている時に十分積極的にフラッシングを行わないと、キャッシュ126が飽和する恐れがある。逆に、ホストの負荷が比較的低い時にあまり積極的にフラッシングすると、キャッシュが不足してポテンシャル・キャッシュ・ヒットを満足できないままになる恐れがある。

30

両方のシナリオともISE108システムの性能に悪影響を及ぼす。

#### 【0028】

図7は、キャッシュ・マネージャ170および経路制御装置164（図6）内に常駐するRAIDコンテナ・サービス172を示す機能ブロック図である。この図は、またHASTEモジュール168のポリシー・エンジン174およびシーク・マネージャ176も示す。シーク・マネージャ176は1つしか図示してないが、ストレージ128に対して専用のシーク・マネージャ176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174からHASTE規則に個々に応答する。

#### 【0029】

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内に位置することができる。ハードウェア内に位置する場合には、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/F160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、動的にホスト負荷を特徴づけそれに続けてシーク・マネージャ176を支配する経路179を介してHASTE規則を発行する。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フラッシングするために、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。ポリシー・エンジン174は、経路184を通してキャッシュ166の状態を引き続き知らされ、経路186を通してキャッシュ・マ

40

50

ネージャにH A S T E 規則を同様に発行することができる。

【 0 0 3 0 】

ポリシー・エンジン 1 7 4 は、1 つまたは複数のネットワーク要求デバイスからの I / O コマンドの現在の速度のようなリアルタイムの負荷についての定量データを収集することができる。ポリシー・エンジン 1 7 4 は、負荷を動的に特徴付け、コマンド・プロファイルの特徴付けとの関係においてストレージ 1 2 8 に連続的に調整するために、負荷に関する定性データを収集する。例えば、好適には、ポリシー・エンジン 1 7 4 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、ホストの負荷を特徴付ける連続データをリアルタイムで収集することが好ましい。

【 0 0 3 1 】

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 1 2 8 にいずれの要求をフラッシングするかはたいした問題ではないからである。実際には、速度に敏感な要求は、ダーティ・データとしてキャッシュ 1 6 6 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さえできるからである。問題は、速度に敏感なコマンドを、キャッシュ 1 6 6 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。

【 0 0 3 2 】

一方、1 つまたは複数のストレージ 1 2 8 内に格納しているデータを読み出すためのホスト・アクセス・コマンドは、同様に、ホスト・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。

【 0 0 3 3 】

ポリシー・エンジン 1 7 4 は、また、関連するデータ・ファイル（帯域幅）のサイズ、アクセス・コマンドを開始する H B A 1 0 3 および / またはネットワーク・デバイス、ストレージ 1 2 8 のアクセス履歴、またはブック・アクセス履歴の期間のようなその任意の一部、タイムスタンプ・データ、R A I D クラス、およびアクセス・コマンドが送られる L U N クラスのような、しかしこれに限定されない、ホストの負荷を特徴付ける定性データを収集することができる。

【 0 0 3 4 】

定性データを収集する際に、ポリシー・エンジン 1 7 4 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。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。

【 0 0 3 5 】

ポリシー・エンジン 1 7 4 は、シーク・マネージャ 1 7 6 に対する規則を作成する際に性能の目標 1 8 8 に応じることができる。目標 1 8 8 は、量的（定量的）なものであっても質的（定性的）なものであってもよいが、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）という点でネットワーク負荷のある種の要因である所望のコマンド・プロファイルの強化（強制設定）、異なる L U N クラスへの割り当てられた優先度の強化（強制設定）、所望の読取りコマンド・レイテンシの強化（強制設定）等に限定されな

10

20

30

40

50

い。それ故、ポリシー・エンジン174は、キャッシュに格納しているライトバック・コマンドおよびもっと高い優先度ホスト読取りコマンドからの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を発行する許可を与える目的で、シーク・マネージャ176を支配しているHASTE規則を定義するために、負荷特性および予め定義した性能目標188の両方を使用することができる。

#### 【0036】

さらに、ポリシー・エンジン174は、シーク・マネージャ176を支配してい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ISE108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。この状態において、ポリシー・エンジン174は、投影され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不測の事態を実行する可能性が高い。ポリシー・エンジン174は、また、ストレージ128へのコマンド・プロファイルを調整する時に、シーク・マネージャ176を支配しているHASTE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与しない懸案のバックグラウンドI/O192またはI/Oの状態に応じることができる。

10

#### 【0037】

図8は、本発明の例示としての実施形態によりHASTEを実施するための方法200のステップを示すフローチャートである。この方法200は、レイテンシおよびアドレス要因と関連するように、ダーティ・データの均一の分布によりランダムにフラッシングするために、HASTE規則を実行するデフォルト・モードでブロック202から開始する。デフォルト・モードは、その間に定性HASTEデータが収集される、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行される。最近のHASTEデータは、書込みに対する読取りの比率で、ホストの負荷を動的に特徴付けるために、ブロック204で使用される。

20

#### 【0038】

ブロック206においては、ポリシー・エンジンは、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っきり分かるか否かを判定する。ブロック206における判定が「いいえ」である場合には、制御は、ブロック202に戻り、デフォルト状態のままである。しかし、ブロック206の判定が「はい」である場合には、ブロック208において、ポリシー・エンジンは、ストレージへのコマンド・プロファイル内で引き続き調整を行う目的でHASTE規則を呼び出すために、ホストの負荷、そして恐らく、目標188、システム状態190およびバックグラウンドI/O192を使用する。例えば、制限無しに、飽和状態で読取りコマンドに対する書込みコマンドの比率が高くなった場合には、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の点でホストの負荷にコマンド・プロファイルを一致させるために、シーク・マネージャを支配することができる。ポリシー・エンジンは、できるだけ迅速に滑らかに飽和状態から回復するために、読取りレイテンシおよびLUNクラス優先度ののような他の規則を修正することさえできるし、一時的に中止させることもできる。HASTE規則は、その間にHASTEデータの次のバッチが収集され、制御がブロック204に戻る1秒間隔のような所定の間隔中に呼び出される。

30

40

#### 【0039】

例示としての実施形態の場合には、HASTEデータは、最近の1秒の間隔中に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の決定、およびその比率の移動平均に対する比較の実施などで、ホストの負荷を質的に特徴付ける。下記の例示としての例のために、すべてのLUNクラスは、すべてのLUNクラスへのシステム・リソースの利用度がバランスのとれたものになるように、同じ優先度を有するものとして処理される。しかし、他の等価の実施形態の場合には、LUNクラスには、シーク・マネージャ176を支配するHASTE規則に分かれる優先度レベルを割り当てることができる。また、この例示としての例のために、HASTEデータは、RAID-1ストレージ・アレイ

50

に割り当てられるストレージ 128 の 1 つはプールに対して入手される。

【 0 0 4 0 】

この例の場合には、HASTEモードの目標 188 は、各ストレージ 128 が、要求しているネットワーク・デバイスからの到着のその平均速度に比例して書込みコマンドおよび読取りコマンドを混合することである。分析は、適当な RAID レベルで分かれる書込みに対する読取りの平均比率を観察することにより開始する。RAID - 1 の場合には、例えば、各ホスト書込みコマンドに対して 2 つの書込みコマンドが発生する。この情報は、(計算のために) 16384 の仮定した分母値を含む整数の分子であると見なすことができる「速度に敏感な要求の要因」(FRSR) にとり入れられる。例示としての例の場合には、FRSR は 5461 である。それ故、「レイテンシに敏感な要求の要因」は下式により計算することができる。

10

【数 1】

$$FLSR = 16384 - 5461$$

$$FLSR = 10923$$

【 0 0 4 1 】

ここで、ストレージ 128 は、200 IOPS でアクセス・コマンドの所与の混合物を処理することができるものと仮定する。それ故、この前の 8 回の 1 秒サンプル間隔の移動平均は、1600 のアクセス・コマンドに跨る。この移動平均は、観察した 1067 の読取りコマンドおよび 533 の書込みコマンドに基づくものと仮定し、最後の 1 秒サンプル間隔中には、66 の読取りコマンドおよび 34 の書込みコマンドが存在したと仮定する。正規化 FRSR を計算すると、下式のようなになる。

20

【数 2】

$$FRSR = \frac{(533 + 34 = 567) \ll 14}{(1600 + 100 = 1700)}$$

$$FRSR = 5464$$

30

【 0 0 4 2 】

FRSR の目標は 5461 であり、実際の FRSR は目標を超える。それ故、ポリシー・エンジン 174 は、現在のコマンド・プロファイルを調整しない。

【 0 0 4 3 】

しかし、上記と同じ状況において、最新の 1 秒サンプルは、70 の読取りコマンドおよび 30 の書込みコマンドを生成したと仮定する。この正規化 FRSR を計算すると、下式のようなになる。

【数 3】

$$FRSR = \frac{(533 + 30 = 563) \ll 14}{(1600 + 100 = 1700)}$$

$$FRSR = 5425$$

40

【 0 0 4 4 】

目標の FRSR と実際の FRSR とを比較する場合のマイナスのデルタは、速度検出 (ライトバック) コマンドの数をコマンド・プロファイル内で増大する必要があることを示す。速度検出コマンドに対する適当な修正を決定するために、下記の関係を使用する。

【数4】

$$FRSR = \frac{((RS.IOPS + X) \ll 14)}{TOTAL.IOPS + X}$$

【0045】

Xについてこの式を解くと、下式のようになる。

【数5】

$$X = \frac{((TOTAL.IOPS * FRSR) - (RS.IOPS \ll 14))}{FLSR}$$

10

【0046】

上記例の場合には、下式のようになる。

【数6】

$$X = \frac{((1700 * 5461) - (563 \ll 14))}{10923}$$

$$X = 5$$

20

【0047】

それ故、シーク・モード176は、所望の速度に観察した速度を一致させるために余分の書込みコマンドを実行させる。このことは、例えば、ドライブ自身を、高い優先度（またはレイテンシ検出）コマンドになるように、内部に待ち行列を含む5つの書込みコマンドを促進させることにより行うことができる。

【0048】

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HASTE内のストレージにコマンド・プロファイルを制御するための手段を予想する。この説明および添付の特許請求の範囲の意味のために、「制御するための手段」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、特徴により制御装置・プロファイル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。コマンド・プロファイルを「直接」調整することにより、「制御するための手段」は、制御装置112が、特徴に応じてダーティ・データおよび未決のホスト読取り要求からの複数のI/Oコマンドから選択したI/Oコマンドの発行を実際に調整することをはっきりと意味する。この説明および添付の特許請求の範囲の意味のために、「制御するための手段」という用語は、それによりキャッシュ・マネージャが、間接的にコマンド・プロファイルに影響を与える恐れがあるフラッシュ・リストを決定する機構の単なる調整を予想していない。

30

40

【0049】

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。

【0050】

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い。

50

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。

【図面の簡単な説明】

【0051】

【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。 10

【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。

【図7】インテリジェント記憶素子内のアレイ制御装置の一部の機能ブロック図である。

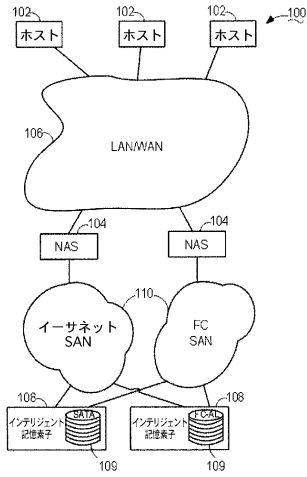
【図8】本発明の実施形態によるHASTEを呼び出すための方法のステップを示すフローチャートである。

【符号の説明】

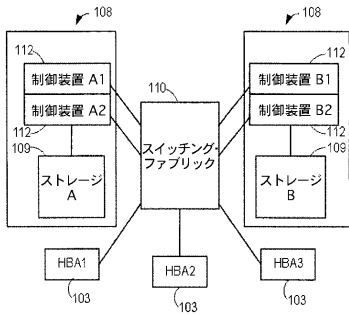
【0052】

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)	20
104	サーバ	
106	LAN/WAN	
108	インテリジェント記憶素子(ISE)	
110	SAN	
109	データ・ストレージ容量	
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	30
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	
130, 132	隔壁	
134	回路基板	
136	コネクタ	
138	ラッパ	
140	モータ	
142	読取り/書込みヘッド	
148	回路基板	40
150	フレックス回路	
154	プロセッサ	
160, 162	I/F回路	
164	経路制御装置	
166	キャッシュ	
168	ホスト適応シーク技術環境(HASTE)モジュール	

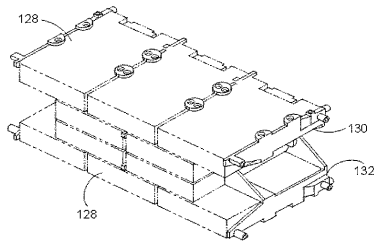
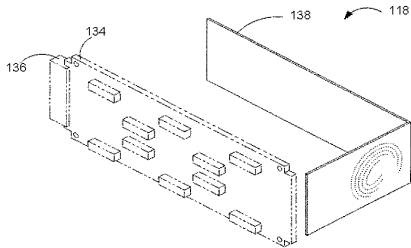
【図1】



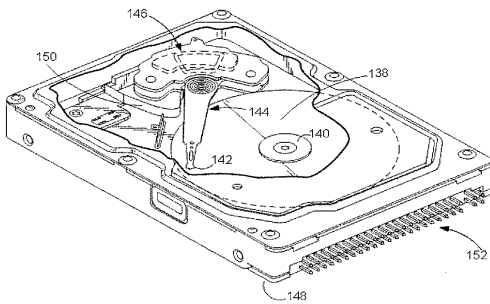
【図2】



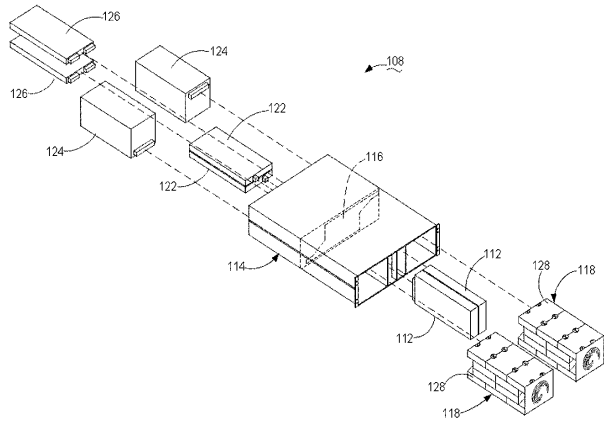
【図4】



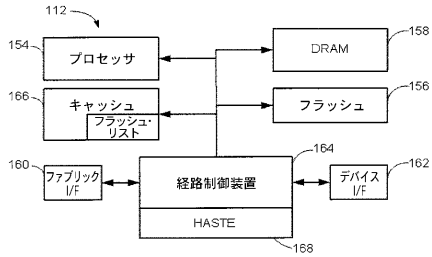
【図5】



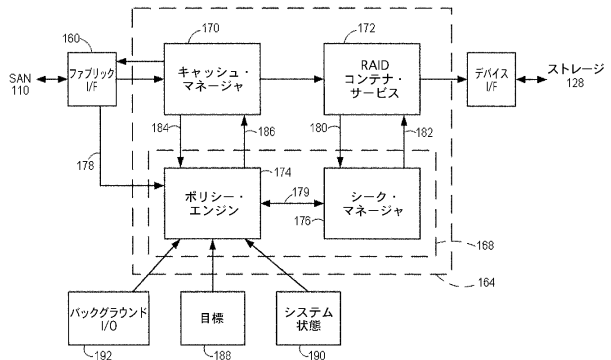
【図3】



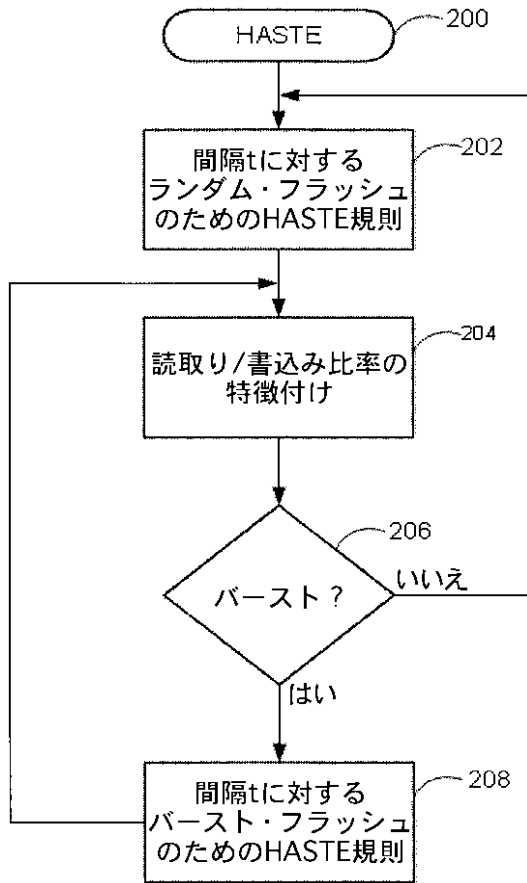
【図6】



【図7】



【図8】



## フロントページの続き

- (74)代理人 100111246  
弁理士 荒川 伸夫
- (74)代理人 100124523  
弁理士 佐々木 真人
- (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 銘三
- (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 パレイ 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06-243042(JP, A)  
特開平05-108274(JP, A)  
特開2004-295860(JP, A)  
特表2006-521640(JP, A)  
米国特許出願公開第2002/0015972(US, A1)  
特表2002-542568(JP, A)  
米国特許出願公開第2004/0162901(US, A1)  
米国特許出願公開第2003/0149838(US, A1)  
米国特許出願公開第2002/0156972(US, A1)  
米国特許第06339811(US, B1)  
米国特許第06170042(US, B1)

(58)調査した分野(Int.Cl., DB名)

G06F 13/10

G06F 3/06