

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5638744号
(P5638744)

(45) 発行日 平成26年12月10日 (2014. 12. 10)

(24) 登録日 平成26年10月31日 (2014. 10. 31)

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

請求項の数 14 (全 14 頁)

(21) 出願番号	特願2008-166109 (P2008-166109)	(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2009-9573 (P2009-9573A)		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22日 (2011. 6. 22)		・ブルバード、10200
審査番号	不服2013-17342 (P2013-17342/J1)	(74) 代理人	100064746
審査請求日	平成25年9月9日 (2013. 9. 9)		弁理士 深見 久郎
(31) 優先権主張番号	11/768, 849	(74) 代理人	100085132
(32) 優先日	平成19年6月26日 (2007. 6. 26)		弁理士 森田 俊雄
(33) 優先権主張国	米国 (US)	(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	100111246
			弁理士 荒川 伸夫

最終頁に続く

(54) 【発明の名称】 コマンド・キュー・ローディング

(57) 【特許請求の範囲】

【請求項1】

データ・ストレージ・システムにおける制御装置であって、前記制御装置は、
ネットワーク・デバイスから前記制御装置へのコマンドを含む負荷の動的特徴を作成するために、および前記動的特徴を、前記コマンドに含まれるデータ転送要求のコマンド・キューの深さに連続的に関連させるように対応のコマンドキューの深さを選択するために、前記ネットワーク・デバイスから前記制御装置への前記負荷に関する定性データを連続的に前記負荷に含まれるコマンドから収集するポリシー・エンジンを含む経路制御装置を備え、前記定性データは、前記負荷内のコマンドを速度に敏感なコマンドとレイテンシに敏感なコマンドとを識別して分類し、前記動的特徴は、前記収集された定性データに従って前記負荷を動的に特徴づけた結果得られたものであり、

10

前記速度に敏感なコマンドは、前記ネットワークデバイスからの該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コマンドであり、前記レイテンシに敏感なコマンドは、前記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コマンドである、制御装置。

【請求項2】

前記経路制御装置は、さらに、前記ポリシーエンジンにตอบสนองして、前記負荷が、デフォルトモードとして期待される定常状態と同じか、またはそれ以下であることを前記動的特徴が示している場合に、コマンド・キューの深さとして第1のコマンド・キューの深さを

20

設定し、前記動的特徴がコマンド活動のバーストを示している場合に、前記第1のコマンド・キューの深さより深い第2のコマンド・キューの深さを設定するキャッシュ・マネージャを備える、請求項1に記載の制御装置。

【請求項3】

前記第2のコマンド・キューの深さが、前記第1のコマンド・キューの深さより少なくとも約1桁深い、請求項2に記載の制御装置。

【請求項4】

前記選択したコマンド・キューの深さが、前記選択したコマンド・キューの深さを有する前記コマンド・キューからデータ転送コマンドを選択的に発行してコマンドの分布または配列を示すコマンド・プロファイルを規定するシーク・マネージャに提示される、請求項3に記載の制御装置。

10

【請求項5】

前記ポリシー・エンジンが、前記動的特徴に関連して、また前記コマンド・プロファイルを定義する際に、前記シーク・マネージャが現在実現する目標について前記シーク・マネージャのシーク動作を規定する規則を設定する、請求項4に記載の制御装置。

【請求項6】

前記ポリシー・エンジンの規則が、前記レイテンシに敏感なコマンドに対する前記速度に敏感なコマンドの数の比率の点において前記コマンド・プロファイルを前記動的特徴に相関させる、請求項5に記載の制御装置。

【請求項7】

20

前記ポリシー・エンジンの規則が、前記レイテンシに敏感なコマンドの数に対する前記速度に敏感なコマンドの数の比率の点において前記コマンド・プロファイルを前記動的特徴とマッチさせる、請求項6に記載の制御装置。

【請求項8】

前記ポリシー・エンジンの規則が、前記コマンド・プロファイルを前記コマンドの所望の最大レイテンシと相関させる、請求項5に記載の制御装置。

【請求項9】

前記ポリシー・エンジンの規則が、前記コマンド・プロファイルを異なるLUNのグループに割り当てられた優先度と相関させる、請求項5に記載の制御装置。

【請求項10】

30

前記動的特徴が、前記レイテンシに敏感なコマンドに対する前記速度に敏感なコマンドの比率である、請求項3から9のいずれかに記載の制御装置。

【請求項11】

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10に記載の制御装置。

【請求項12】

前記負荷の動的特徴が、各コマンドに関連するファイル・サイズである、請求項11に記載の制御装置。

【請求項13】

40

前記ポリシー・エンジンが、有限状態機械を備える、請求項12に記載の制御装置。

【請求項14】

制御装置によって行われる方法であって、

ネットワーク・デバイスからストレージ・システムへのコマンド・ストリーム負荷を監視して、前記コマンド・ストリーム内のコマンドについての定性データを前記コマンドから収集して、前記定性データに従って前記負荷を動的に特徴づけ、該特徴づけにより前記負荷の動的特徴を作成するステップを含み、前記定性データは、前記コマンドを前記ネットワークデバイスからの該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と前記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

50

影響がより大きいレイテンシに敏感なコマンドとを識別して分類し、

前記動的特徴を使用して、前記コマンドに含まれるデータ転送要求の選択すべきコマンド・キューの深さを、前記動的特徴が前記コマンドに含まれるデータ転送要求の対応のコマンド・キューの深さに連続的に関連させるように規定するステップをさらに含む、方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明の実施形態は、概して、ストレージ・システム分野に関し、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。

10

【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。

20

【0003】

さらに最近では、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。

30

【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。

40

【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、バースト負荷の状況下でシーク・コマンド・キューを選択的にローディング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。

50

【 0 0 0 6 】

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、負荷の特徴をライトバック・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コマンド・キューの深さに連続的に相関付けるために、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。

【 0 0 0 7 】

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 】

10

【 0 0 0 8 】

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 1 0 0 である。1 つまたは複数のホスト 1 0 2 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) 1 0 6 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 1 0 4 にネットワークで接続している。好適には、LAN/WAN 1 0 6 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・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 1 0 2 は、多数のインテリジェント・ストレージ素子 (ISE) 1 0 8 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的に必要とするサーバ 1 0 4 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、SAN 1 1 0 は、格納しているデータにアクセスするために、サーバ 1 0 4 を ISE 1 0 8 に接続する。ISE 1 0 8 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・チャン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 1 0 9 を提供する。

20

【 0 0 0 9 】

図 2 は、図 1 のコンピュータ・システム 1 0 0 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) 1 0 3 は、ネットワークまたはファブリック 1 1 0 を介して1 対の ISE 1 0 8 (それぞれ A および B で示す) と相互に作用する。各 ISE 1 0 8 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 1 0 9 上で動作することが好ましい二重化冗長制御装置 1 1 2 (A 1、A 2 および B 1、B 2 で示す) を含む。すなわち、好適には、制御装置 1 1 2 およびデータ・ストレージ容量 1 0 9 は、種々の制御装置 1 1 2 が並列の冗長リンクを使用し、システム 1 0 0 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 1 0 9 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。

30

【 0 0 1 0 】

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 ISE 1 0 8 である。シェルフ 1 1 4 は、ミッドプレーン 1 1 6 と電氣的に接続している制御装置 1 1 2 に収容する形で係合するための空洞を定める。シェルフ 1 1 4 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) 1 1 8 は、ミッドプレーン 1 1 6 の同じ側面上のシェルフ 1 1 4 内に収容される形で係合している。ミッドプレーン 1 1 6 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 1 2 2、デュアル交流電源 1 2 4 およびデュアル・インタフェース・モジュール 1 2 6 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 MDA 1 1 8 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。

40

【 0 0 1 1 】

図 4 は、それぞれが 5 つのデータ・ストレージ 1 2 8 を支持している上部隔壁 1 3 0 および下部隔壁 1 3 2 を有する MDA 1 1 8 の拡大分解等角図である。隔壁 1 3 0、1 3 2 は、ミッドプレーン 1 1 6 (図 3) と係合するコネクタ 1 3 6 を有する共通の回路基板 1 3 4 と接続するためにデータ・ストレージ 1 2 8 を整合する。ラッパー 1 3 8 は、電磁妨

50

害シールドを行う。MDA 118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。MDA 118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、MDA 118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に電気的に接続している。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に電気的に接続していて、データ・ストレージ128をMDA 118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。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。

【0016】

I S E 1 0 8 のデータ・ストレージ容量は、データをドライブ 1 2 8 に格納する場合に、およびデータをドライブ 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 】

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。

10

【 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。

20

【 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。

30

【 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。

40

【 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 (R S U) 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

50

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 の値と関連して管理することが好ましい。

【 0 0 2 4 】

制御装置 1 1 2 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 R A I D ストリップが検出された場合に、全 S D D 構造の周期的フラッシングを含む。このことは、S D D が関連するデータをダーティと識別した場合に、R A I D レベル変数に基づいて所与の S D D 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 S D D 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 S D D 構造は、コマンド・キュー内に入れられ、データのフラッシングを開始するようにとの要求が行われる。

【 0 0 2 5 】

もっと小さな組のデータのフラッシングも S D D 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 S D D は、ダーティ L R U 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 L R U 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。

【 0 0 2 6 】

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 1 1 2 は、最も近い位置を有する R A I D 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 R A I D 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。

【 0 0 2 7 】

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 1 2 6 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、ホスト読取り要求のレイテンシに悪影響を与える恐れがある。両方のシナリオとも I S E 1 0 8 システムの性能に悪影響を及ぼす。

【 0 0 2 8 】

図 7 は、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、および経路制御装置 1 6 4（図 6）内に常駐するシーク・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 9 】

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装することができる。ハードウェアで実施する場合には、ポリシー・エンジン 1 7 4 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 1 7 4 は、経路 1 7 8 を介して、I / O 単位ベースでファブリック I / F 1 6 0 経由で受信したアクセス・コマンドにつ

10

20

30

40

50

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コマンド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。

【0030】

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、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、I/O コマンドに関連するファイル・サイズのような、しかしこれに限定されない他のものでホストの負荷を特徴付ける定性データを収集することもできる。

【0031】

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 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施する。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与しない懸案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。

【0032】

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 の任意の組合わせを使用することができる。

【0033】

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。

【0034】

一方、1つまたは複数のストレージ 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原因となり得る。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。この場合には、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

10

20

30

40

50

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。

【0035】

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ル間隔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、自立式カウンタは、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、現在の1秒比率を照合するために9番目のスロットを含む、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、最新のサンプル値を加算し、次に、比率の最新の移動平均を計算するために8で除算を行う。

10

【0036】

例示としての例について引き続き説明すると、ポリシー・エンジン174が、負荷バーストが発生していることを観察した場合には、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ために、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、読取りコマンドに対する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、一時的に延期することもできる。

20

【0037】

このような状態の場合、シーク・マネージャ176にキャッシュ・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、転送要求のコマンド・キューの深さを増大することにより、フラッシング性能を有意に改善することができることが分かっている。すなわち、この実施形態は、また、シーク・マネージャ176に提示されたコマンド・キューを「ロードアップ」するキャッシュ・マネージャ170を管理する規則を定義するために、負荷の特徴を使用するポリシー・エンジン174を予想している。

【0038】

図8および図9は、それぞれ本発明の実施形態による、定常状態の負荷の状態およびI/Oバースト状態中のアレイ制御装置112の関連部分を示す略図である。194および196で示すキャッシュ166内のライトバック・キャッシュされたデータのリストは、通常、両方の状態の場合、数千ノード程度である。ドライブ・キュー198は、比較的短く、両方の状態の場合一定である。しかし、キャッシュ・マネージャ170は、定常状態の場合、第1のコマンド・キューの深さ200を定義するために、ポリシー・エンジン174からの規則に応じ、バースト状態において第1のコマンド・キューの深さ200より有意に深い第2のコマンド・キューの深さ202を定義することに留意されたい。通常、コマンド・キューの深さ200は、約20~40転送要求の範囲内である。飽和状態でバースト状態が発生している場合には、調整したコマンド・キューの深さ202は、約200~500の要求の範囲内にあるか、または少なくとも定常状態の状態より大きいことが好ましい。

30

【0039】

図10は、本発明の実施形態によるコマンド・キュー・ローディングのための方法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、20~40程度の転送要求のようなデフォルト・モード・コマンド・キューの深さで、ブロック20から開始する。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、書込みに対する読取りの比率およびI/O速度等で、ホスト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。

40

【0040】

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わかるか否かを判定する。ブロック206に

50

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、そのためデフォルト状態が継続する。しかし、ブロック206における判定が「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、デフォルト深さより約1レベル深いような有意に深い深さにコマンド・キューのローディングを呼び出す目的で、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を使用する。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・プロファイルをホストにマッチさせるために、シーク・マネージャを管理することができる。それを支持して、ポリシー・エンジンは、この実施形態によりコマンド・キューの深さを有意に増大する。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、一時的に中止することさえできる。ローディングしたコマンド・キュー規則は、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る。次に、制御はブロック204に戻る。

10

【0041】

通常、この実施形態は、ネットワーク・アクセス・コマンドに応じて転送データへネットワークを接続するように構成されているストレージ・アレイ、およびデータ転送要求のサイズをアクセス・コマンドのネットワーク負荷の観察した特徴と関連付けることにより、コマンド・キューをローディング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、「ローディングのための手段」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。

20

【0042】

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。

30

【0043】

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。

【図面の簡単な説明】**【0044】**

【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・ストレージ素子の分解等角図である

40

。【図4】図3のインテリジェント・ストレージ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】図3のインテリジェント・ストレージ素子内のアレイ制御装置の機能ブロック図である。

【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】予想した定常状態のネットワーク負荷状態中のコマンド・キューのサイズを示す略図である。

50

【図9】図8類似の図面であるが、I/Oバースト状態中にシーク・マネージャに提示されるローディングしたコマンド・キューを示す。

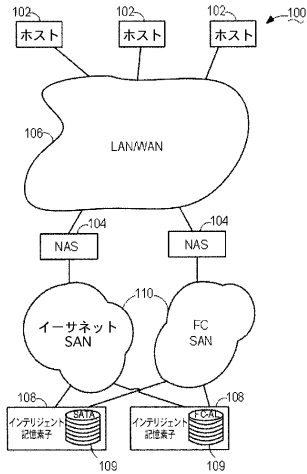
【図10】本発明の実施形態によるコマンド・キュー・ローディングのための方法のステップを示すフローチャートである。

【符号の説明】

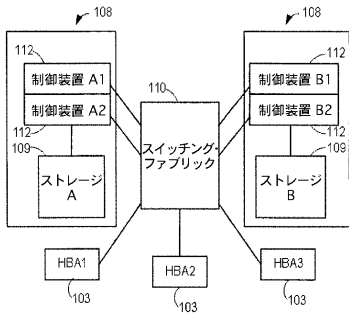
【0045】

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	10
106	LAN/WAN	
108	インテリジェント・ストレージ素子 (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	20
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	
144	回転アクチュエータ	
146	ボイス・コイル・モータ (VCM)	30
148	回路基板	
150	フレックス回路	
152	コネクタ	
154	プロセッサ	
156	不揮発性メモリ	
158	ダイナミック・ランダム・アクセス・メモリ (DRAM)	
160, 162	I/F回路	
164	経路制御装置	
166	キャッシュ	

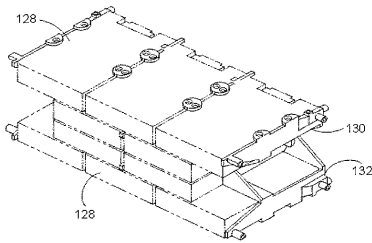
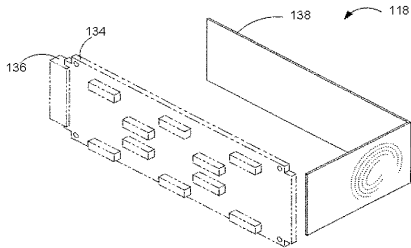
【図1】



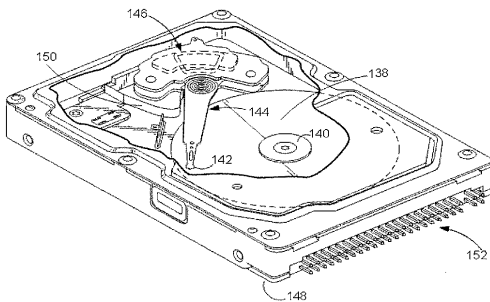
【図2】



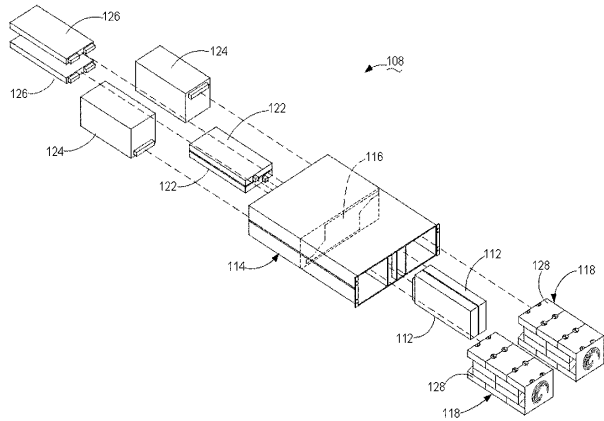
【図4】



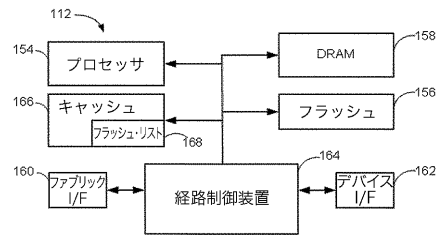
【図5】



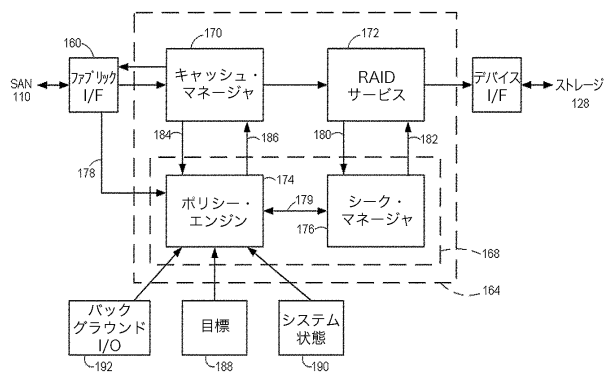
【図3】



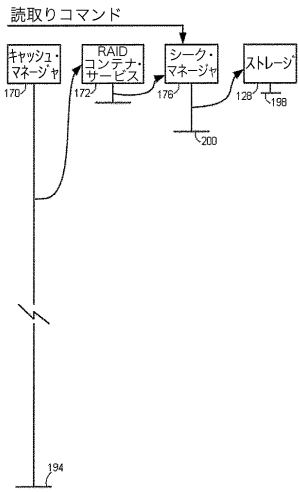
【図6】



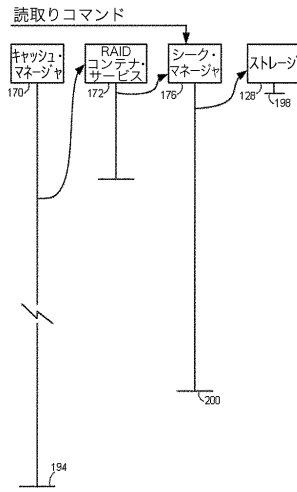
【図7】



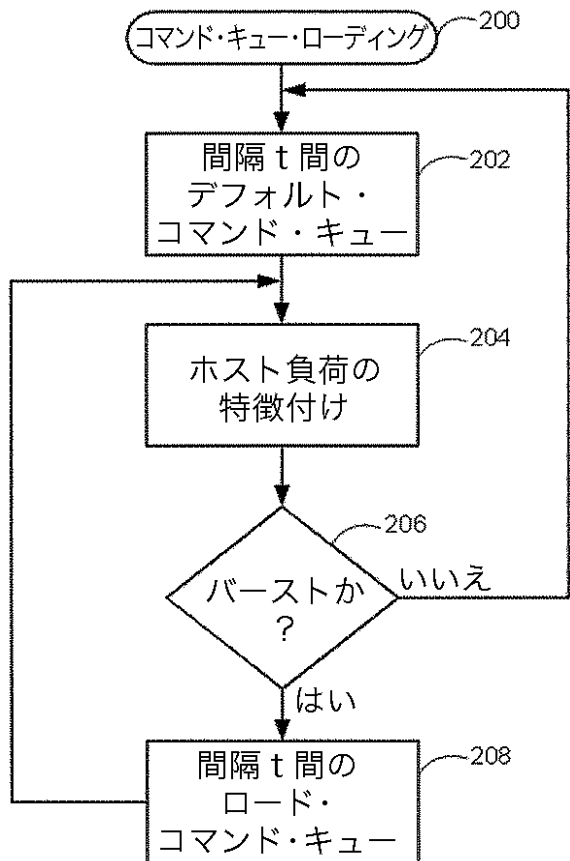
【図 8】



【図 9】



【図 10】



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

合議体

審判長 乾 雅浩

審判官 千葉 輝久

審判官 和田 志郎

(56)参考文献 米国特許第5426736 (US, A)

特表2006-521640 (JP, A)

特開2004-295860 (JP, A)

特開平6-243042 (JP, A)

特開2006-18689 (JP, A)

特開2007-115233 (JP, A)

特開2004-171411 (JP, A)

(58)調査した分野(Int.Cl., DB名)

G06F3/06-3/08