

(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2004-86914

(P2004-86914A)

(43) 公開日 平成16年3月18日(2004.3.18)

(51) Int. Cl.<sup>7</sup>

G06F 3/06

F I

G06F 3/06 3 O 1 F  
 G06F 3/06 3 O 2 Z  
 G06F 3/06 5 4 O

テーマコード (参考)

5 B O 6 5

審査請求 未請求 請求項の数 10 O L (全 12 頁)

(21) 出願番号 特願2003-303987 (P2003-303987)  
 (22) 出願日 平成15年8月28日 (2003. 8. 28)  
 (31) 優先権主張番号 10/233107  
 (32) 優先日 平成14年8月28日 (2002. 8. 28)  
 (33) 優先権主張国 米国 (US)

(71) 出願人 503003854  
 ヒューレット・パカード デベロップメント カンパニー エル. ピー.  
 アメリカ合衆国 テキサス州 77070  
 ヒューストン 20555 ステイト  
 ハイウェイ 249  
 (74) 代理人 100087642  
 弁理士 古谷 聡  
 (74) 代理人 100076680  
 弁理士 溝部 孝彦  
 (74) 代理人 100121061  
 弁理士 西山 清春

最終頁に続く

(54) 【発明の名称】 コンピュータシステム内の記憶装置のパフォーマンスの最適化

(57) 【要約】

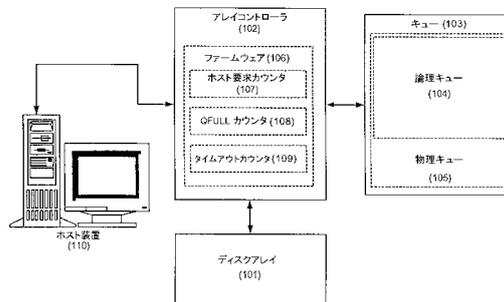
【課題】

多数のホストを接続してもホスト側タイムアウトイベントや記憶装置側の Q F U L L イベントの発生を最小限に抑える記憶装置を提供すること。

【解決手段】

データ記憶装置(101,102)は、ディスクアレイ(101)、該アレイを制御するためのアレイコントローラ(102)、及び、ホストシステムからディスクアレイ(101)へのコマンドを入れるためのキュー(103)で構成される。アレイコントローラ(102)にインストールされたプログラム(106)がキューの論理サイズ(104)を調節してパフォーマンスを最適化する。

【選択図】 図1



## 【特許請求の範囲】

## 【請求項 1】

ディスクアレイ(101)と、  
前記アレイ(101)を制御するためのアレイコントローラ(102)と、  
ホストシステムから前記ディスクアレイ(101)へのコマンドを入れるためのキュー(103)とからなり、  
前記アレイコントローラ(102)にインストールされたプログラム(106)が前記キュー(104)の論理サイズを調節してパフォーマンスを最適化する、データ記憶装置(101,102)。

## 【請求項 2】

少なくとも1つのホスト装置(110)と少なくとも1つのデータ記憶装置(101,102)とを含むコンピュータシステム内のパフォーマンスを最適化する方法であって、前記データ記憶装置(101,102)のキューの論理サイズを調節するステップ(207,209)を含む方法。

## 【請求項 3】

前記データ記憶装置(101,102)にさらなる記憶容量を追加するステップをさらに含む、請求項2の方法。

## 【請求項 4】

ホスト要求をカウントするステップ(201)と、  
Q F U L L イベントをカウントするステップ(202-2)と、  
タイムアウトイベントをカウントするステップ(203-2)と、  
をさらに含む、請求項2の方法。

## 【請求項 5】

発生したQ F U L L イベントの総数を受信したホスト要求の現在の数で割ることによりQ F U L L 率を計算するステップ(204)と、  
前記Q F U L L 率が所定の閾値を超えた場合、前記キューの論理サイズを増大させるステップ(207)と、  
をさらに含む、請求項4の方法。

## 【請求項 6】

発生したタイムアウトイベントの総数を受信したホスト要求の現在の数で割ることによりタイムアウト率を計算するステップ(206)と、  
前記タイムアウト率が所定の閾値を超えた場合、前記キューの論理サイズを減少させるステップ(209)と、  
をさらに含む、請求項4の方法。

## 【請求項 7】

少なくとも1つのホスト装置(110)と、  
少なくとも1つのデータ記憶装置(101,102)と、  
前記少なくとも1つのホスト装置(110)と前記少なくとも1つのデータ記憶装置(101,102)とを接続するネットワークとからなるコンピュータシステムであって、  
前記少なくとも1つのデータ記憶装置(101,102)が、  
ディスクアレイ(101)と、  
前記アレイ(101)を制御するためのアレイコントローラ(102)と、  
ホストシステムから前記ディスクアレイへのコマンドを入れるキュー(103)とからなり、  
前記アレイコントローラ(102)にインストールされたプログラム(106)が前記キューの論理サイズ(104)を調節してパフォーマンスを最適化する、コンピュータシステム。

## 【請求項 8】

前記少なくとも1つのホスト装置(110)にインストールされたドライバ(301)と、  
前記少なくとも1つのホスト装置(110)上のコマンドキュー(302)とをさらに含み、  
前記ドライバ(301)が、前記少なくとも1つのデータ記憶装置(101,102)上のQ F U L L イベント(108)又はタイムアウトイベント(109)にตอบสนองして、前記ホスト装置(110)上の前記コマンドキュー(302)のサイズを調節する、請求項7のコンピュータシステム。

## 【請求項 9】

少なくとも1つのホスト装置(110)と、少なくとも1つのデータ記憶装置(101,102)と、前記少なくとも1つのホスト装置(110)と前記少なくとも1つのデータ記憶装置(101,102)とを接続するネットワークとを含み、前記少なくとも1つのデータ記憶装置(101,102)がディスクアレイ(101)と、該アレイ(101)を制御するためのアレイコントローラ(102)と、ホストシステムから前記ディスクアレイへのコマンドを入れるためのキュー(103)とを含む、コンピュータシステムにおいてパフォーマンスを最適化する方法であって、

前記少なくとも1つのホスト装置(110)によって課せられた前記少なくとも1つのデータ記憶装置(101,102)に対する要求に応じて前記キューの論理サイズを調節するステップ(207,209)を含む方法。

10

## 【請求項 10】

前記少なくとも1つのデータ記憶装置(101,102)上のQ F U L L イベント(108)又はタイムアウトイベント(109)に応じて前記少なくとも1つのホスト装置(110)上のコマンドキューのサイズを調節するステップをさらに含む、請求項9の方法。

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

本発明は、コンピュータシステムの分野に関し、特に、ホスト装置の要求とデータ記憶装置の容量とのバランスをとる手段および方法を提供する。

## 【背景技術】

20

## 【0002】

コンピュータシステムには、通常、ホストと呼ばれる1つまたは複数のコンピュータが含まれる。複数のコンピュータが使用される場合、それらのコンピュータは、データの共有を可能にするネットワークによって相互接続される。通常、そのようなネットワークは、ネットワーク接続されたそれらのコンピュータに対して追加のデータ記憶容量を提供するための1つまたは複数のデータ記憶装置も含む。一般的なデータ記憶装置はディスクアレイであり、R A I D (Redundant Array of Independent (またはInexpensive) Disks)と呼ばれることもある。ディスクアレイは、接続されたホストに対してデータ記憶装置を提供する2つ以上のハードディスクまたは同種のディスクである。

## 【0003】

30

コンピュータシステムは、システム中の全コンポーネントを均等なバランスのパフォーマンスにした時に最適に動作する。ホストシステムおよび接続された記憶装置は、記憶装置のホスト要求を満たす能力がホストによって生成された作業負荷とおおよそ等価になるように、通常バランスを取らなければならない。さらに複雑な構成の場合、多数のホストが1つまたは複数の記憶装置との間でデータを転送することが可能になっている場合がある。複数のホストが1つの記憶装置にアクセスする場合、その記憶装置は、接続されたすべてのホストのパフォーマンス要求に対応できることが重要である。パフォーマンスの観点から記憶装置の能力がホストの能力よりも大幅に低い場合、記憶装置はシステムパフォーマンス全体における制限要素になる可能性がある。

## 【0004】

40

ファイバーチャネル(F C)テクノロジーは、ホスト装置と記憶装置とを接続またはネットワーク化するテクノロジーである。ファイバーチャネルテクノロジーは、1つの記憶装置に接続された何千ものホストを含むことが可能な構成を可能にする。ファイバーチャネルは、S C S I (小型コンピュータシステムインタフェース)プロトコル、たとえばS C S I - 3通信プロトコルを使用する。したがって、何千に及ぶ可能性があるホストの各々は、ディスクアレイ内で最大65,536までの論理記憶装置(L U N)をアドレス指定することができる。今日では、何千ものL U Nを作成する能力を有するファイバーチャネルテクノロジーを用いて通信するディスクアレイが出荷されている。

## 【0005】

ディスクアレイ構成は、4つ以下のドライブを有するアレイから100を超えるドライ

50

ブを有するアレイまで、さまざまである。最終的に、ディスクアレイのパフォーマンスは、アレイ中のディスク数によって制限される。これは、すべての作業が最終的にディスクに渡されることになるからである。

【0006】

明らかに、ファイバーチャネルテクノロジーは、1つの記憶装置に多数のホストを接続することをユーザに勧めるものである。多数のホストは、ディスクアレイ等の記憶装置に対して莫大な量の作業を命令する可能性がある。したがって、そのような環境にいるユーザは、何百ものディスクを有するエンタープライズディスクアレイを用いた作業に慣れているので、多数のホストをサポートするのに必要なパフォーマンスを受け入れることができる場合がある。

10

【0007】

高要求環境では、各ホストが複数のLUNに対する複数の入出力(I/O)要求を常に有する場合がある。そうした要求の各々は、処理されるまでキューに入れる必要があり、キューが一杯の場合は拒絶する必要がある。

【0008】

システムパフォーマンスは、広範囲に研究されてきた複雑な課題である。多数のパフォーマンス要因および変形を考慮すると、ホスト/アプリケーション構成によって生成される作業負荷が記憶システムの能力を超えししまうか否かをユーザが設計の時点で知ることは、不可能ではないが非常に困難なものとなりうる。したがって、記憶パフォーマンスとのバランスがとれていないシステムをユーザが設定してしまうことは、決してめずらしいことではない。

20

【0009】

1つのアンバランスなシステムにおいて、ユーザは、接続されたホストによって生成される作業負荷を大幅に超える記憶装置のパフォーマンス能力を有する場合がある。ユーザは、そのソリューション全体のパフォーマンスに満足できず、所望のパフォーマンス結果を達成するために、さらにホスト容量を追加しなければならない場合がある。

【0010】

本発明は、ホスト/アプリケーションの作業負荷が接続された記憶装置の能力を超過する相補的なシナリオを扱う。本明細書で使用する「過剰構成(over-configuration)」という用語は、ホスト/アプリケーションのパフォーマンス要求が記憶装置のパフォーマンス能力を超過するシナリオを説明するのに用いられる。別の言い方をすると、システム設計者がソリューション全体の要求に対して記憶装置を十分に供給しなかったということである。

30

【0011】

過剰構成されたシステムは、不満よりも、もっと深刻な問題を引き起こす場合がある。データ記憶アレイがシステムパフォーマンス要求を満たすことができない場合、ユーザアプリケーションの応答回数が許容できないほど多くなる。もっと極端な場合、アプリケーションエラー、サーバクラッシュおよびホストクラスタ障害などの悲惨な事態の兆候が現れる場合もある。それらの悲惨な事態が発生する理由は、ホストのソフトウェア(ドライバ、アプリケーション、ボリュームマネージャなど)が、記憶装置によって要求が完了されるのを待つ時間を制限しているからである。ディスクアレイ等の記憶装置の応答時間がホストソフトウェアによって課せられた制限時間を超えると、悲惨な事態が発生する可能性がある。明らかに、過剰構成を検出する方法および過剰構成を完全に回避する方法が必要とされている。

40

【0012】

現在、過剰構成問題は、ユーザが試行錯誤によって発見している。この発見は、大抵の場合サポート要求によって行なわれ、その場合、サポート要員はシステムがアンバランスな状態で動作していることを推測できるにすぎない。これは、パフォーマンス問題が発生した際に、それらのアンバランスな構成を回避する方法をユーザに知らせることに重点をおく受動的な解決方法であり、それらのアンバランスな構成を検出および軽減する方法を

50

ユーザに知らせるものではない。

【0013】

パフォーマンスの複雑性および動的性質が原因で、この発見は、非常に誤りを引き起こしがちな処理になる可能性がある。アンバランスな構成は極めて一般的である。アンバランスな構成が設定されると、ユーザはまず、システム全体のパフォーマンスが許容できないものであることを確認しなければならない。ユーザは、手がかりを提供する様々なツールを用いて、十分に実行されていないシステムを確認する。しかしながら、この処置には、構成、アプリケーション、及び、様々な情報源から取得したその他多数の情報の断片を手作業で検査することが含まれる。この処置は専門家による自学自習の分析で終わるのが常であり、専門家はそのシステムがアンバランスであると結論する。この結論が得られると、次のステップとして、さらに記憶装置を追加し、パフォーマンス問題が解決されることを期待することになる。

10

【0014】

一般に、I/O要求が記憶装置に送信されると、そのI/O要求は、許容または拒絶される。許容された要求は、サービスキューに配置され、さらなる処理を待つ。新たな要求を取り込んでキューに入れる記憶装置の能力を超えると、その要求は拒絶され、キューが満杯であるというステータス(QFULL)とともにホストシステムに返される。QFULLはこの状態についての適当な戻りステータスであるが、記憶装置によっては、QFULLではなくBUSYを返す場合もある。これは、ホストシステムドライバがBUSY状態およびQFULL状態について一般に異なる再試行挙動を示すという点で重要である。

20

【0015】

記憶装置はQFULLイベントが発生し得ないほどそのサービスキューを大きくすることもできるが、この解決方法にも限界がある。サービスキューのサイズは、実際には、ホストシステムドライバおよびアプリケーションによって課せられたタイムアウト時間によって制限される。記憶装置によりアプリケーションが永久待機させられるということが絶対に起こらないようにするため、ホストシステムドライバ及びいくつかのソフトウェアアプリケーションは、I/O要求に対してタイムアウトを課している。要求がタイムアウトすると、ホストは、記憶装置に何らかの問題があることを推定する。この状態から復帰するため、ホストシステムは、そのI/Oを中止すべきことを記憶装置に命令し、その後その要求を再試行する(何らかの有限数の再試行回数まで)。

30

【0016】

キューに入るI/O要求の数が極めて多い場合、それらの要求を処理する際の待ち時間は、ホストシステムのドライバ及び/又はホストシステムのアプリケーションによって課せられたタイムアウト制限を超えてしまう場合がある。そうしたタイムアウトが発生する理由は、単に、記憶装置がサービスキュー中にある要求の数を許容時間内に処理することができないからである。そのような理由から、記憶装置のサービスキューのサイズについては実際上制限がある。

【0017】

QFULLイベントおよびタイムアウトの両方に対するホストの反応動作は、パフォーマンスの点で費用がかかる。実際に、両イベントはシステム内を伝搬し、実際の、すなわち知覚される障害を生じさせる可能性がある。両イベントは、ホストプログラム(すなわち、ファイバーチャネルドライバ)においてエラー事象として扱われる。それらのイベントは、ホストシステムのログに現れ、重大な障害として扱われる場合が多い。この解釈は、他の挙動次第で正しい場合もあればそうでない場合もある。

40

【発明の開示】

【発明が解決しようとする課題】

【0018】

QFULLイベントおよびタイムアウトイベントの処理は、記憶装置の設計者にジレンマを残す。一方、記憶サービスキューのサイズは、最大まで大きくすることが望ましい。キューを大きくすることは、より多数のホストが接続できるようになることを意味する。

50

これにより、不要な Q F U L L イベントを生じさせることなくアレイの接続性を向上させることができる。ホストおよびアプリケーションの I / O 要求特性が適当なものであれば、比較的低速な記憶装置についても、大量のホストを備えたバランスのとれたシステムを構成することが可能であるかもしれない。一方、記憶装置の設計者は、最悪の場合のアクセスパターンを有する大型の構成であってもタイムアウトイベントを発生させることなく対処することができるように、キューサイズを十分小さく保つことを望んでいる。明らかに、システム設計者がこれらの競合する考慮事項のバランスをとろうとすることは、従来は極めて困難であった。

【課題を解決するための手段】

【0019】

10

多数の考え得る実施形態のうちの一つにおいて、本発明は、ディスクアレイと、ディスクアレイを制御するアレイコントローラと、ホストシステムからディスクアレイへのコマンドを入れるためのキューとを有するデータ記憶装置を提供する。アレイコントローラにインストールされた例えばファームウェアなどのプログラムは、パフォーマンスを最適化するためにキューの論理サイズを調節する。

【0020】

本発明の他の実施形態は、少なくとも一つのホスト装置と少なくとも一つのデータ記憶装置とを有するコンピュータシステム内のパフォーマンスを最適化する方法であって、データ記憶装置のキューの論理サイズを調節することを含む方法も提供する。

【0021】

20

本発明の他の実施形態は、少なくとも一つのホスト装置と、少なくとも一つのデータ記憶装置と、少なくとも一つのホスト装置と少なくとも一つのデータ記憶装置とを接続するネットワークとを有するコンピュータシステムも提供する。データ記憶装置は、ディスクアレイと、ディスクアレイを制御するアレイコントローラと、ホストシステムからディスクアレイへのコマンドを入れるキューとを有する。アレイコントローラにインストールされたプログラムは、パフォーマンスを最適化するためにキューの論理サイズを調節する。

【0022】

本発明の他の実施形態は、少なくとも一つのホスト装置と、少なくとも一つのデータ記憶装置と、少なくとも一つのホスト装置と少なくとも一つのデータ記憶装置とを接続するネットワークとを有するコンピュータシステムであって、前記少なくとも一つのデータ記憶装置がディスクアレイと、ディスクアレイを制御するアレイコントローラと、ホストシステムからディスクアレイへのコマンドを入れるキューとを有するコンピュータシステム内のパフォーマンスを最適化する方法も提供する。この方法は、少なくとも一つのホスト装置によって課せられた少なくとも一つのデータ記憶装置への要求に応じて、キューの論理サイズを調節することを含む。

30

【0023】

添付の図面は、本発明の様々な実施形態を例示する。図面は、下記の説明とあわせて本発明の原理を例示および説明するものである。例示の実施形態は本発明の例であり、本発明の特許請求の範囲を限定するものではない。

【0024】

40

図面を通して同一の符号は類似要素を指すものであり、必ずしも同一の要素を指すものではない。

【発明を実施するための最良の形態】

【0025】

本発明は、ユーザおよびサポート要員がキューサイズおよび応答時間に関してシステムについて最適な選択を行うことを可能にする新たなアレイキュー手段を導入している。本発明の原理に従う記憶装置のファームウェア又は他のプログラムによって、サポート要員は、システムの動作中にアレイのサービスキューのサイズを動的に変更することができる。変更はほぼ瞬時に反映される。

【0026】

50

図1は、本発明の一実施形態による、ホスト装置とデータ記憶装置とを備えた改良されたコンピュータシステムを示すブロック図である。本明細書および特許請求の範囲で使用する「記憶装置」という用語は、ディスクアレイ、ディスクドライブ、テープドライブ、任意のSCSIデバイスまたはファイバチャネル装置を含む任意のデータ記憶装置を指すものとして使用するが、それらに限定はしない。図1に示すように、例示のシステムは、1つまたは複数のディスクアレイ(101)を含むデータ記憶装置に接続された1つまたは複数のホスト装置(110)を含む。

**【0027】**

ディスクアレイ(101)はアレイコントローラ(102)によって制御される。アレイコントローラ(102)は、ホスト装置(110)からI/O要求を受信し、それらの要求に応じてディスクアレイ(101)を制御する。アレイコントローラ(102)は、ホスト要求を処理することができるようになるまでそれらをキュー(103)に入れておく。様々なホスト装置(110)がこの1つの共通キューに対して要求を提出する。

10

**【0028】**

本発明の原理によると、キューの物理サイズ(105)はキューの論理サイズ(104)と区別される。キューの物理サイズ(105)は、ファームウェア、ソフトウェアまたは他のプログラム(以下「ファームウェア」と呼ぶ)(106)がアレイコントローラ(102)にインストールされた時点でサポートされていた最大論理キューサイズ(104)を収容できるように設定される。論理キューサイズ(104)は、キュー(103)がホスト(110)にQFULLイベントを返す前に格納しているホストI/O要求の数である。論理キューサイズ(104)は、システム動作中に、アレイコントローラ(102)のファームウェア(106)により、物理キューサイズ(105)の最大値まで動的に調節することができる。ファームウェア(106)のアルゴリズムは、論理キューサイズ(104)を物理キューサイズ(105)よりも増大させようとする試みを防止するように構成することが好ましい。

20

**【0029】**

また、本発明によるディスクアレイファームウェア(106)は、本記憶装置システムの診断および調節に利用可能なキューパフォーマンス情報を追跡してログ記録する3つのイベントカウンタも導入している。それらのカウンタは、記憶装置が遭遇したホスト要求(107)、QFULLイベント(108)およびタイムアウトイベント(109)の数を追跡する。実際には、記憶装置は、要求の処理の中止を通知された回数をカウントすることにより、タイムアウトイベントを間接的にカウントすることができる。中止の大半がタイムアウトの結果生じるものと仮定した場合。

30

**【0030】**

ユーザおよびサポートチームメンバは、調節可能な論理キューサイズ(104)とともに、記憶装置のそれらのカウンタ(107~109)を用いて取り込んだ情報を用いて、自分達のシステムの作業負荷・パフォーマンスを調節することができる。目的は、QFULLイベントを実質的になくすのに十分な程度大きいと同時に、タイムアウトイベントを最小限にするのに十分な程度小さい論理キューサイズ(104)を見つけることである。QFULLイベントやタイムアウトを完全になくすことは、実際的でない。良くバランスのとれた正常なコンピュータシステムであっても、QFULLイベントおよびタイムアウトは、稀なイベント又は異常なイベント(作業負荷の予期せぬスパイク等)に起因して発生する場合がある。

40

**【0031】**

目的は、より正確には、論理キューサイズ(104)を操作して2つの割合、すなわちホスト要求の数に対するQFULLイベントの数の割合(QFULL率)と、ホスト要求の数に対するタイムアウトイベントの数の割合(タイムアウト率)とを最小化することであるべきである。目的は、QFULL率とタイムアウト率との両方をゼロに近づける論理キューサイズ(104)を見付けることになる。

**【0032】**

50

図2は、本発明の一実施形態による、ホスト装置およびデータ記憶装置を備えたコンピュータシステムを動作させる方法を示すフロー図である。図2に示すように、論理キューサイズを最適化するプロセスは、QFULL率およびタイムアウト率を最小化するなわちゼロに近づけようと試みる。

【0033】

まず、コントローラファームウェアのカウンタがホスト要求をカウントする。ホスト要求を受信(200)する毎に、ホスト要求カウンタをインクリメントする(201)。ホスト要求によってQFULLイベントが生じた場合(202-1)、QFULLイベントカウンタをインクリメントする(202-2)。同様に、ホスト要求が処理されなかったためにタイムアウトイベントが生じた場合(203-1)、タイムアウトカウンタをインクリメントする(203-2)。

10

【0034】

次に、QFULLイベントの数を受信したホスト要求の現在の数で除算することにより、QFULL率を計算する(204)。そして、このQFULL率を閾値と比較する(205)。理想的なシステムではこの閾値がゼロであるが、実際のシステムではそれよりも高く設定される。QFULL率が許容可能な閾値を超過すると(205)、アルゴリズムは論理キューサイズを増大させる(207)。このようにして追加のキュー容量を与えることにより、QFULLイベントの数が減少する。

【0035】

次に、タイムアウトイベントの数を受信したホスト要求の現在の数で除算することにより、タイムアウト率を計算する(206)。ここでも、理想的なシステムではこの閾値がゼロであるが、実際のシステムではそれよりも高くすることが多い。タイムアウト率が許容可能な閾値を超過すると(208)、アルゴリズムは論理キューサイズを減少させる(209)。このようにして実行待ちの要求の数を制限することにより、タイムアウトイベントが防止される。

20

【0036】

両方の割合がそれぞれの閾値を超過した場合(210)、記憶装置は、生成されているホスト作業負荷に対処することができなくなる。これは、異常に高いホスト作業負荷に起因した一時的な状態である可能性もある。しかしながら、この状態が持続するならば、システムは恐らく過剰構成されている。そのような場合、ファームウェアは、システム管理者に対し、システムが過剰構成されている可能性があるという通知を生成することができる。そして、システム管理者は正しい処置を考えることができる。正しい対応は、ディスクアレイの能力を増大させる(キャッシュ、アレイコントローラまたはディスクを追加することにより)ことや、単に別のディスクアレイを追加することになるであろう。

30

【0037】

上記の説明では、ディスクアレイその他の記憶装置のキューに重点を置いている。しかしながら、ホスト要求はホスト自体のキューも通過する。また、ホストドライバおよびホストバスアダプタ(HBA)も、それら自体の内部キューで要求を処理する。それらのキューも本発明の原理に基づいて動的に調節することができるので、ホスト内で使用されるアルゴリズムを理解することは重要である。

40

【0038】

記憶装置では、I/O要求を管理する手段としてコマンドサービスキューが使用される。コマンドサービスキューは、低頻度、短期間の高い要求も受け入れるとともに、通常の負荷も処理できる程度に十分大きいものである必要がある。ディスクアレイのI/O処理速度がホストのI/O要求速度と同じかそれよりも速い場合、サービスキューの利用率は、長期にわたって変化しないかまたは減少してゆく。ホストのI/O要求速度がそれらの要求をディスクアレイが処理できる速度を超過した場合、キュー容量を使い果たすまでサービスキューレベルが増大してゆく。

【0039】

複雑なストレージエリアネットワーク(SAN)環境におけるディスクアレイには、複

50

数のホストによって要求がポストされる。それらのホストの各々は、バースト要求を送信することができる。本明細書で使用する「バースト」という用語は、要求の数が一定であること、及び、それらの要求がディスクアレイの要求を処理する能力を超える速度で伝達されることを意味する。従って、バースト要は、アレイのサービスキューの需要量を増大させる結果となる。各ホストは、自分が送信するバーストのサイズをホストバスアダプタ（HBA）ドライバによって課されたコマンドキュー制限に従って制限する。

**【0040】**

ビジー環境においてホストが示す通常の挙動は、ホストがバースト要求（ホストのキューで制限されたもの）を送信した後、それらの要求がアレイによって処理されるのを待つというものになる。ディスクアレイがそれらの要求を処理するのに従って、ホスト側のコマンドキューが減少し、ホストはより多くの要求を送信できるようになる。ホスト側のコマンドキューサイズを低減することを含む策により、QFULLの戻りステータスで拒絶されるコマンド数を最低限に抑えることができる。

10

**【0041】**

実際には、ホストの要求速度が一定でないように、アレイのI/O処理速度も一定ではない。上記のように、本発明の原理に基づくと、記憶装置は、システムの必要性に合わせてサービスキュー制限値を調節することが可能になる。さらに、ホスト装置（単数の場合も複数の場合もある）のコマンドキューも、システムの必要性に合わせて調節することが可能である。

**【0042】**

図3は、本発明の他の実施形態によるシステムを示す。このシステムは、少なくとも1つのホスト装置と少なくとも1つの記憶装置とを有し、ホスト装置（単数の場合も複数の場合もある）のコマンドキューのサイズを調節して、システムを均等にすることができる。

20

**【0043】**

たとえば、キュー（103）を有する特定のディスクアレイ（101）が、750コマンドの論理サイズ（104）を有するものと仮定する。この容量にも関わらず、この特定のディスクアレイのキュー深さ、すなわち接続されたホスト装置（110）によりキュー（103）に配置される要求は、通常の動作中に約100コマンドであり、適度な数のコマンドを完了するのに0.2秒よりも長くかかることが分かっている。

30

**【0044】**

この実施例では、論理キューサイズ（104）をデフォルトの750コマンドから実際の使用数である100コマンドに変更することにより、ディスクアレイ（101）のコマンドレイテンシ曲線がさらに広がるのを防止することができる。そして、ホスト（110）がもっとアクティブになると、アレイコントローラ（102）は、この活力の増加によって生じるコマンドを、QFULLステータスの指示を用いて拒絶することにより応答する。

**【0045】**

本発明の原理に基づくと、ホストドライバ（301）は、ホスト側のコマンドキュー（302）を減少させることにより応答することが好ましい。これにより、ホスト（110）から記憶装置（101～103）に送信されるコマンドの速度が減少する。

40

**【0046】**

平衡状態に達した後、アレイ（101）は、作業負荷が増大する前と同じコマンドレイテンシレベルでコマンドを処理することになる。ホスト側では、システムコールインタフェースのブロッキングの性質により、アプリケーションに要求速度を「キュー満杯閾値」で管理されたレベルに強制的に維持させる。このシナリオは、応答時間を所望のレベルに維持するが、要求間の最小時間は増大する。この要求間の時間の増大は、長い目で見ればほとんどのユーザが満足しないことになるといふことの兆候である。もっと良い解決方法は、構成にさらなる記憶処理能力を追加することである。

**【0047】**

50

また、この「キュー満杯閾値」の調節を利用して、長いドライバ関連タイムアウトが発生する割合に影響を与えることもできる。特定の作業負荷について、オペレーティングシステム（O/S）ドライバ（301）が5秒毎に1回の割合でタイムアウトイベントを報告している間、記憶装置（101～103）が400コマンドの論理キュー深さを維持しているものと仮定する。その場合、コマンドレイテンシ分布の末尾を短くするためには、キューの論理サイズ（104）を300まで縮小することが好ましい。

【0048】

すると、アレイドコントローラ（102）は、QFULLステータスによって大量のコマンドを拒絶することになる。ホスト側のドライバ（301）は、QFULL戻りステータスがなくなるまでホスト側のキュー（302）の深さを縮小することにより、QFULLステータスに応じることが好ましい。今度は記憶装置（101～103）がもっと短いコマンドキュー（104）で動作しているので、著しい待ち時間（例えば30秒など）をもつコマンドの割合が大幅に減少するはずである。

10

【0049】

ホスト側のドライバ（301）は、ホスト側のコマンドキュー（302）の深さを再度増大させることにより、記憶装置（101～103）を時々検査する。これを行なう場合、ホスト側のドライバ（301）がホスト側のキュー（302）の深さを再度縮小させることを判断するまで、QFULL応答が生じることになる。

【0050】

先の説明は、本発明の例示及び説明のために提供したものに過ぎない。その説明は、網羅的なものにする 것도、本発明を開示したいずれかの形態そのままに制限することも意図していない。上記の示唆を考慮して、多数の変更および変形が可能である。

20

【0051】

好ましい実施形態は、本発明の原理とその実際の使用を最もよく例示する目的で選択され、説明されている。先の説明は、他の当業者が本発明を様々な実施形態で意図する特定の用途に適するように様々な変更を加えて最もよく利用できるようにすることを意図している。本発明の範囲は、特許請求の範囲で規定されることを意図している。

【図面の簡単な説明】

【0052】

【図1】本発明の一実施形態による改良されたコンピュータシステムを示すブロック図であり、該システムは、少なくとも1つのホスト装置と、調節可能なキューを有する少なくとも1つのデータ記憶装置とを有する。

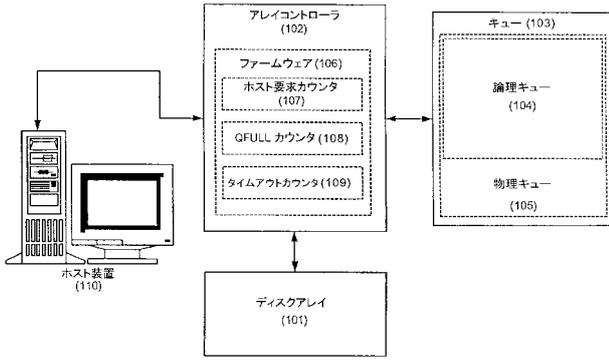
30

【図2】本発明の一実施形態による、ホスト装置とデータ記憶装置とを備えたコンピュータシステムを動作させる方法を示すフロー図である。

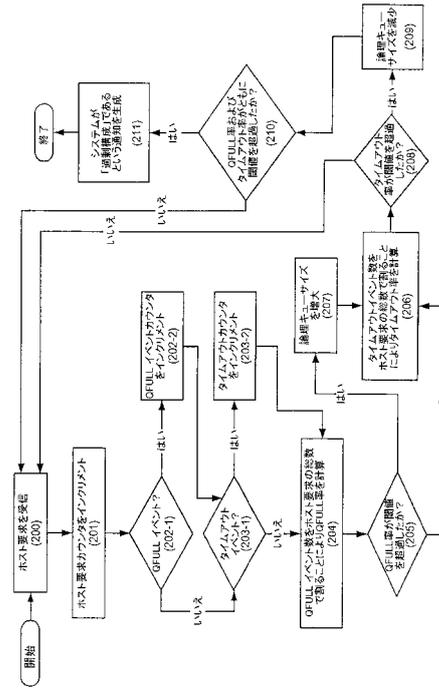
【図3】本発明の他の実施形態によるコンピュータシステムを示す図であり、該システムは、少なくとも1つのホスト装置と少なくとも1つの記憶装置とを有し、ホスト装置のコマンドキューのサイズを調節して、システムを均等にするようにできている。

。

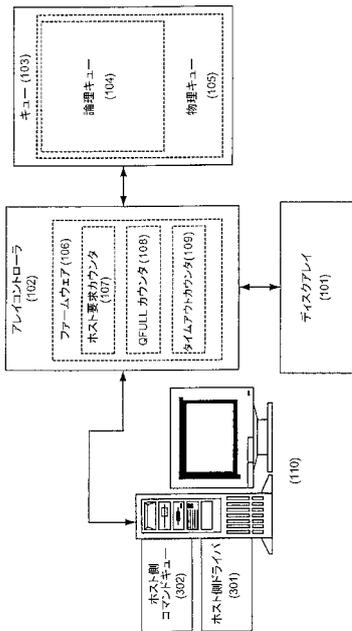
【 図 1 】



【 図 2 】



【 図 3 】



---

フロントページの続き

- (72)発明者 ランディ・ジェイ・マシューズ  
アメリカ合衆国アイダホ州 8 3 7 1 3 , ボイス, ヒッコリー・ナット・ストリート・1 1 2 4 5
- (72)発明者 マーク・イー・レフェブル  
アメリカ合衆国アイダホ州 8 3 7 1 3 , ボイス, ウエスト・アナブルック・ドライブ・1 3 4 0 1
- (72)発明者 リachel・エル・アールバース  
アメリカ合衆国アイダホ州 8 3 6 1 6 , イーグル, ウエスト・キャンポ・レーン・1 9 5 1
- (72)発明者 ウェイド・エイ・ドルフィン  
アメリカ合衆国アイダホ州 8 3 6 1 9 , イーグル, ウエスト・キャンポ・1 9 5 1
- (72)発明者 ダグラス・エル・ボイト  
アメリカ合衆国アイダホ州 8 3 7 0 2 , ボイス, ノース・トゥエンティフォース・3 0 3 0
- Fターム(参考) 5B065 BA01 CA15 CA30 CC08 CH18 EK06