

(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2006-106975

(P2006-106975A)

(43) 公開日 平成18年4月20日(2006.4.20)

(51) Int. Cl. F I テーマコード (参考)
G06F 3/06 (2006.01) G06F 3/06 305F 5B065
 G06F 3/06 540

審査請求 未請求 請求項の数 6 O L (全 9 頁)

<p>(21) 出願番号 特願2004-290501 (P2004-290501) (22) 出願日 平成16年10月1日 (2004.10.1)</p>	<p>(71) 出願人 000001007 キヤノン株式会社 東京都大田区下丸子3丁目30番2号 (74) 代理人 100090273 弁理士 國分 孝悦 (72) 発明者 外山 猛 東京都大田区下丸子3丁目30番2号 キヤノン株式会社内 (72) 発明者 伊藤 博康 東京都大田区下丸子3丁目30番2号 キヤノン株式会社内 Fターム(参考) 5B065 BA01 CA11 CA30 CH05 EA01 EA15</p>
--	--

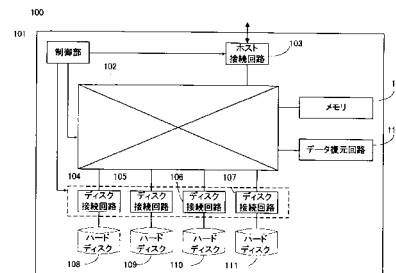
(54) 【発明の名称】 ディスクアレイ装置及びその制御方法

(57) 【要約】

【課題】 ディスクアレイ装置の読み出し応答時間の短縮を図る。

【解決手段】 ハードディスク装置の所定領域ごとの障害情報を記憶し、ホストからの要求データが記録されている領域が読み出し障害を発生する領域であるか否かを該障害情報に基づいて判断し、その判断結果に基づいて、当該領域からデータを読み出す第1の読み出し方式、又は、他のハードディスク装置から復元用データを読み出して該要求データを復元する第2の読み出し方式のいずれかを選択して実行するようにした。

【選択図】 図1



【特許請求の範囲】

【請求項 1】

ハードディスク装置の所定領域ごとの障害情報を記憶した記憶手段と、
上位装置からの第 1 のデータが記録されている対応領域が読み出し障害を発生する領域
であるか否かを前記障害情報に基づいて判断する判断手段と、

前記判断手段の判断結果に基づいて、前記対応領域から前記第 1 のデータを読み出す第
1 の読み出し方式、又は、前記要求データを復元するための第 2 のデータを読み出した後
、該第 2 のデータから前記第 1 のデータを復元する第 2 の読み出し方式のいずれかを選択
して実行する制御手段とを有することを特徴とするディスクアレイ装置。

【請求項 2】

前記記憶手段に記憶された障害情報は、読み出しの実行結果に応じて更新されることを
特徴とする請求項 1 に記載のディスクアレイ装置。

【請求項 3】

前記ハードディスク装置は、自ら読み出しの再試行を実行しないように設定されている
ことを特徴とする請求項 1 又は 2 に記載のディスクアレイ装置。

【請求項 4】

ハードディスク装置の所定領域ごとの障害情報を記憶し、
上位装置からの第 1 のデータが記録されている対応領域が読み出し障害を発生する領域
であるか否かを前記障害情報に基づいて判断し、

該判断結果に基づいて、前記対応領域から前記第 1 のデータを読み出す第 1 の読み出し
方式、又は、前記要求データを復元するための第 2 のデータを読み出した後、該第 2 のデ
ータから前記第 1 のデータを復元する第 2 の読み出し方式のいずれかを選択して実行する
こと特徴とするディスクアレイ装置の制御方法。

【請求項 5】

ハードディスク装置の所定領域ごとの障害情報を記憶し、
上位装置からの第 1 のデータが記録されている対応領域が読み出し障害を発生する領域
であるか否かを前記障害情報に基づいて判断し、

該判断結果に基づいて、前記対応領域から前記第 1 のデータを読み出す第 1 の読み出し
方式、又は、前記要求データを復元するための第 2 のデータを読み出した後、該第 2 のデ
ータから前記第 1 のデータを復元する第 2 の読み出し方式のいずれかを選択して実行する
処理をコンピュータにて実行させることを特徴とするコンピュータプログラム。

【請求項 6】

請求項 5 に記載のコンピュータプログラムを記録したことを特徴とするコンピュータ読
み取り可能な記録媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ディスクアレイ装置、その制御方法、コンピュータプログラム、及びコンピ
ュータ読み取り可能な記録媒体に関する。

【背景技術】

【0002】

近年、複数のハードディスク装置を並列に動作させて読み書き速度を高速化し、冗長構
成により信頼性を高めたディスクアレイ装置が普及している。

【0003】

ディスクアレイ装置には、R A I D (Redundant Arrays of Inexpensive Disks) と呼
ばれる技術が採用されている。R A I D には、用途に応じて様々な構成が存在する。図 5
に R A I D 5 の構成例を示す。5 0 1 ~ 5 0 4 は、それぞれ別個のハードディスク装置で
ある。D 0 ~ D 8 は、それぞれストライプと呼ばれる所定量のデータ格納領域であり、デ
ィスクアレイ装置におけるデータ管理制御の単位となっている。

【0004】

10

20

30

40

50

ディスクアレイ装置を介してハードディスク装置をアクセスする装置（以下、「ホスト」と称する）は、各ストライプデータがいずれのハードディスク装置に記録されているかの情報は知る必要はなく、ホストからの論理アドレスを各ハードディスク装置上のアドレスに変換する作業はディスクアレイ装置が担う。

【0005】

P02、P35、P68は、それぞれ冗長データの格納領域である。P02にはストライプD0～D2内のデータの排他的論理和が格納されている。P35には、ストライプD3～D5内のデータの排他的論理和が格納されている。P68には、ストライプD6～D8内のデータの排他的論理和が格納されている。

【0006】

このようにRAID5では、いずれかのハードディスク装置が壊れた場合でも、残りのハードディスク装置から壊れたハードディスク装置のデータを復元することが可能な構成となっている。

【0007】

一方で、ハードディスク装置に対しては、信頼性ととも、読み出し速度向上の要求が年々高まっており、この要求に答えるべく各種の改良案が提案されている。例えば、特許文献1には、ストライプデータを読み出す際、当該ストライプデータの復元に必要なデータ（以下、復元用データと称する）も並列に読み出しおき、前記ストライプデータの読み出しエラーが発生した場合にはリトライをすることなく前記他のデータにより前記ストライプデータを復元する技術が提案されている。

【0008】

【特許文献1】特開平10-247133号公報

【発明の開示】

【発明が解決しようとする課題】

【0009】

上記従来例は、ストライプデータの読み出しと当該ストライプデータの復元用データの読み出しとを並行して行うことによりディスクアレイ装置の読み出し応答時間を短縮しようとしている。

【0010】

しかしながら、並列に読み出されたデータは、共にディスクアレイ装置内の共有のデータ経路（例えば、上記公報における図1の103）を流れることになるので、読み出しエラーが発生しない場合の読み出し応答時間はむしろ長くなってしまいうという問題がある。

【0011】

そのため、読み出しエラーの発生率が低い場合には、読み出しのスループットはむしろ悪くなってしまいう。

【0012】

本発明は上記のような点に鑑みてなされたものであり、読み出し障害の発生により読み出し応答時間が長くなってしまいうのを防止することを目的とする。

【課題を解決するための手段】

【0013】

本発明のディスクアレイ装置は、ハードディスク装置の所定領域ごとの障害情報を記憶した記憶手段と、上位装置からの第1のデータが記録されている対応領域が読み出し障害を発生する領域であるか否かを前記障害情報に基づいて判断する判断手段と、前記判断手段の判断結果に基づいて、前記対応領域から前記第1のデータを読み出す第1の読み出し方式、又は、前記要求データを復元するための第2のデータを読み出した後、該第2のデータから前記第1のデータを復元する第2の読み出し方式のいずれかを選択して実行する制御手段とを有する点に特徴を有する。

本発明のディスクアレイ装置の制御方法は、ハードディスク装置の所定領域ごとの障害情報を記憶し、上位装置からの第1のデータが記録されている対応領域が読み出し障害を発生する領域であるか否かを前記障害情報に基づいて判断し、該判断結果に基づいて、前

10

20

30

40

50

記対応領域から前記第1のデータを読み出す第1の読み出し方式、又は、前記要求データを復元するための第2のデータを読み出した後、該第2のデータから前記第1のデータを復元する第2の読み出し方式のいずれかを選択して実行する点に特徴を有する。

本発明のコンピュータプログラムは、ハードディスク装置の所定領域ごとの障害情報を記憶し、上位装置からの第1のデータが記録されている対応領域が読み出し障害を発生する領域であるか否かを前記障害情報に基づいて判断し、該判断結果に基づいて、前記対応領域から前記第1のデータを読み出す第1の読み出し方式、又は、前記要求データを復元するための第2のデータを読み出した後、該第2のデータから前記第1のデータを復元する第2の読み出し方式のいずれかを選択して実行する処理をコンピュータにて実行させる点に特徴を有する。

10

本発明のコンピュータ読み取り可能な記録媒体は、上記本発明のコンピュータプログラムを記録した点に特徴を有する。

【発明の効果】

【0014】

本発明によれば、読み出し障害の発生により読み出し応答時間が長くなってしまふのを防止することができる。それと同時に、従来例に対して、ディスクアレイ装置内のデータ経路の帯域に必要以上の負荷がかからないので、その空いている帯域を他の並列動作に割り当てることが可能となる。

【発明を実施するための最良の形態】

【0015】

以下、添付図面を参照して、本発明の好適な実施形態について説明する。図1は、本実施形態におけるディスクアレイ装置の構成を示した図である。101は、装置全体の制御を司る制御部であり、CPU、制御プログラムを格納するROM、制御に必要なワークメモリを格納するRAM等から構成される。

20

【0016】

102は、スイッチ回路であり、制御部101による設定により、ホスト接続回路103、ディスク接続回路104～107、メモリ112、データ復元回路の各ユニット間のデータ経路を形成する。

【0017】

103は、不図示のホストと接続するためのホスト接続回路であり、例えば、PCI Expressのような高速シリアルインターフェースによりホストとの間のインターフェース制御を行う。

30

【0018】

104～107はディスク接続回路であり、例えば、Serial ATAのような高速シリアルインターフェースによりハードディスク装置108～111との間のインターフェース制御を行う。

【0019】

108～111は、ハードディスク装置である。制御部101はハードディスク装置108～111によりRAID5のディスクアレイを構成する。

【0020】

112は、メモリであり、不図示のメモリコントローラによる制御によりメモリへのデータの書き込み、及び、メモリからのデータの読み出しが行われる。

40

【0021】

113は、データ復元回路であり、復元用データを用いてストライプデータの復元を行う。データ復元回路113は内部に所定サイズのメモリを有しており、該メモリ内のデータと入力されたデータとの排他的論理和を計算した結果を該メモリに出力する。

【0022】

次に、図2を参照して、本実施形態におけるハードディスク装置からのデータの読み出し動作を説明する。図2において、メモリ112、ハードディスク装置108～111、及びデータ復元回路113は、図1と同じ符号を用いている。

50

【0023】

まず、ハードディスク装置108内の領域201のデータを読み出す動作（以下、「直接読み出し」と称する）を説明する。制御部101はスイッチ回路102を制御して、ディスク接続回路104とメモリ112との間のデータ経路を形成し、領域201のデータをメモリ112内の領域207に転送する。

【0024】

その後、制御部101はスイッチ回路102を制御して、メモリ112とホスト接続回路103との間のデータ経路を形成し、領域202のデータをホストへ転送する。

【0025】

次に、領域201のデータの復元用データ、すなわち、領域203、領域204、及び領域205のデータを読み出して領域201のデータを復元する動作（以下、「復元読み出し」と称する）を説明する。 10

【0026】

制御部101はスイッチ回路102を制御してハードディスク装置109とデータ復元回路113との間のデータ経路を形成し、領域203のデータをデータ復元回路113内の領域206に転送する。

【0027】

この転送処理の際、領域203からのデータは領域206内のデータとの排他的論理和をとられてから領域206に記録される。

【0028】

領域206内は予めゼロに初期化されているので、この転送処理により領域206には領域203のデータと同じデータが記録される。 20

【0029】

次に、制御部101は、スイッチ回路102を制御してハードディスク装置110とデータ復元回路113内の領域206に転送する。

【0030】

この転送処理の際、領域204からのデータは領域206内のデータと排他的論理和をとられた後、領域206に記録される。

【0031】

次に、制御部101は、スイッチ回路102を制御してハードディスク装置111とデータ復元回路113内の領域206に転送する。 30

【0032】

この転送処理の際、領域205からのデータは領域206内のデータと排他的論理和を取られた後、領域206に記録される。

【0033】

領域206には領域201のデータが復元されているので、制御部101は領域206のデータを領域207に転送する。

【0034】

その後、制御部101はスイッチ回路102を制御して、メモリ112とホスト接続回路103との間のデータ経路を形成し、領域207のデータをホストへ転送する。 40

【0035】

このように、ディスクアレイ装置100は、ホストからの読み出し要求に対して、要求データが記録されているハードディスク装置から読み出す直接読み出し方式と、該要求データの復元用データを他のハードディスク装置から読み出した後、要求データを復元する復元読み出し方式とを有する。

【0036】

いずれの方式を選択するかは、ホストからの要求データの記録されている領域が読み出し障害を発生する領域であるか否かの判断結果に基づく。すなわち、ホストからの要求データが記録されている領域が読み出し障害を発生する領域でないと判断した場合は、直接読み出し方式を選択する。ホストからの要求データが記録されている領域が読み出し障害 50

を発生する領域であると判断した場合は、復元読み出し方式を選択する。ホストからの要求データの記録されている領域が読み出し障害を発生する領域であるか否かの判断は、当該領域において過去読み出し障害が発生したか否かの情報を記録した障害領域情報を参照して行われる。

【0037】

次に、図3を参照して、本実施形態における障害領域情報のフォーマットを説明する。301～304は、ディスクアレイ装置100内の不図示の不揮発性メモリ内にハードディスク装置108～111ごとに設けられた障害領域情報である。

【0038】

ディスクアレイ装置100では、各ハードディスク装置のデータ領域を所定サイズのデータブロックごとに管理している。例えば、305に示すように各データブロックが障害領域であるか否かをビット情報で表したビットマップとして管理する。そして、ある領域で、直接読み出し方式で読み出し障害が発生したら、該領域を含むデータブロックのビットを1にセットする。以後、該データブロックは読み出し障害の発生する領域として判断される。

10

【0039】

次に、図4のフローチャートを参照して、本実施形態におけるディスクアレイ装置の読み出し動作を説明する。ハードディスク装置108～111は、セクタの読み出しがエラーした場合に自らリトライを実行する機能を備えている。ディスクアレイ装置100は、ハードディスク装置108～111の各々について、この自らリトライを実行する機能が動作しないモードに予め設定しておく。これにより、ディスクアレイ装置100がハードディスク装置からデータの読み出しを開始してから読み取り障害の有無を判断するまでの時間を短縮することができる。

20

【0040】

まず、ステップS401では、ホストからの読み出し要求を、ホスト接続回路103を介して受信する。

【0041】

ステップS402では、要求データがいずれのハードディスク装置に記録されているかを調べた後、該要求データが記録された領域が読み出し障害の発生する領域であるか否かを障害領域情報301～304を参照して判断する。

30

【0042】

例えば、要求データがハードディスク装置108に記録されている場合は、障害領域情報301を参照して判断がなされる。

【0043】

以下の説明では、要求データがハードディスク装置108に記録されている場合の例を適宜挙げて補足する。

【0044】

ステップS402の判断が肯定判断の場合はステップS403に進み、直接読み出し方式により要求データがハードディスク装置108から読み出される。読み出されたデータは、メモリ112上に展開される。

40

【0045】

ステップS404では、ステップS403において読み出し障害が発生したか否かを判断する。この判断が否定判断の場合はステップS405に進み、読み出したデータをホストへ転送する。また、ステップS404の判断が否定判断の場合はステップS406に進み、ホストからの要求データが記憶されている領域が読み出し障害を発生する領域であることを示すべく障害領域情報301を変更する。

【0046】

ステップS407の判断では、再読み出しを実行すべきか否かを判断する。再読み出しを実行すべきか否かは、予めディスクアレイ装置100において設定されているものとする。ステップS407の判断が否定判断であればステップS408に進み、ホストへエラ

50

ーを報告する。ステップS 4 0 7の判断が肯定判断の場合、及び、ステップS 4 0 2の判断が肯定判断の場合は、ステップS 4 0 9に進む。

【0047】

ステップS 4 0 7では、復元読み出し方式によりデータの読み出しが実行される。すなわち、ホストからの要求データの復元用データをハードディスク装置109～111から読み出し、該要求データを復元する。復元されたデータは、メモリ112上に展開され、その後、ステップS 4 0 5でホストに転送される。

【0048】

このように、本実施形態によれば、ハードディスク装置の所定領域ごとに障害の有無を示す障害情報を記憶し、ホストからの要求データが記録されている領域が読み出し障害を発生する領域であるか否かを該障害情報に基づいて判断し、その判断結果に基づいて、当該領域からデータを読み出す直接読み出し方式、又は、他のハードディスク装置から復元用データを読み出して該要求データを復元する復元読み出し方式のいずれかを選択して実行するようにした。

10

【0049】

これにより、読み出し障害の発生により、読み出し応答時間が長くなってしまふのを防止することができる。それと同時に、直接読み出しと復元用読み出しのいずれか一方のみを実行するので、すなわち、直接読み出しと復元読み出しを同時並列で実行しないので、ディスクアレイ装置内のデータ経路の帯域に空きが生じ、その空いている帯域を他の並列動作に割り当てることが可能となる。

20

【0050】

なお、本発明の目的は、上述した実施形態の機能を実現するソフトウェアのプログラムコードを記録した記憶媒体を、システム或いは装置に供給し、そのシステム或いは装置のコンピュータ（又はCPUやMPU）が記憶媒体に格納されたプログラムコードを読み出し実行することによっても、達成されることは言うまでもない。

【0051】

この場合、記憶媒体から読み出されたプログラムコード自体が上述した実施形態の機能を実現することになり、プログラムコード自体及びそのプログラムコードを記憶した記憶媒体は本発明を構成することになる。

【0052】

プログラムコードを供給するための記憶媒体としては、例えば、フレキシブルディスク、ハードディスク、光ディスク、光磁気ディスク、CD-ROM、CD-R、磁気テープ、不揮発性のメモリカード、ROM等を用いることができる。

30

【0053】

また、コンピュータが読み出したプログラムコードを実行することにより、上述した実施形態の機能が実現されるだけでなく、そのプログラムコードの指示に基づき、コンピュータ上で稼動しているOS（基本システム或いはオペレーティングシステム）などが実際の処理の一部又は全部を行い、その処理によって上述した実施形態の機能が実現される場合も含まれることは言うまでもない。

【0054】

さらに、記憶媒体から読み出されたプログラムコードが、コンピュータに挿入された機能拡張ボードやコンピュータに接続された機能拡張ユニットに備わるメモリに書き込まれた後、そのプログラムコードの指示に基づき、その機能拡張ボードや機能拡張ユニットに備わるCPU等が実際の処理の一部又は全部を行い、その処理によって上述した実施形態の機能が実現される場合も含まれることは言うまでもない。

40

【図面の簡単な説明】

【0055】

【図1】本実施形態におけるディスクアレイ装置の構成を示した図である。

【図2】本実施形態におけるハードディスク装置からのデータの読み出し動作を説明する図である。

50

【図3】本実施形態における障害領域情報のフォーマットを説明する図である。

【図4】本実施形態におけるディスクアレイ装置の読み出し動作のフローチャートである。

【図5】RAID5の動作を説明する図である。

【符号の説明】

【0056】

100 ディスクアレイ装置

101 制御部

102 スイッチ回路

103 ホスト接続回路

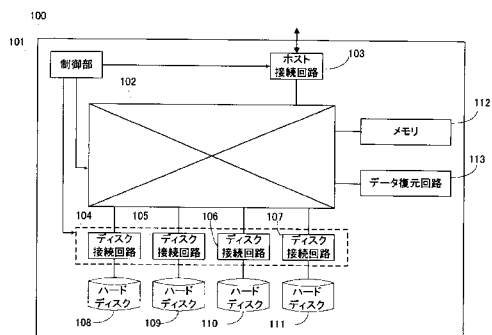
104 ~ 107 ディスク接続回路

108 ~ 111 ハードディスク装置

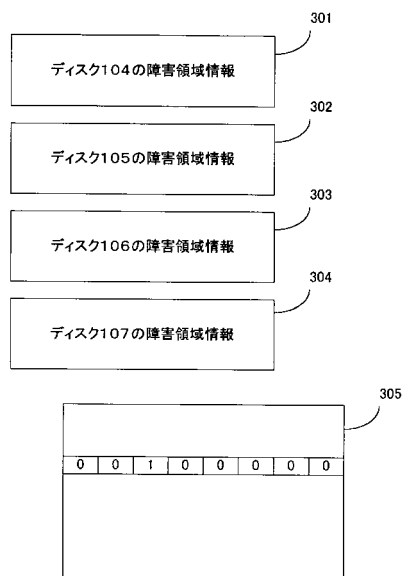
112 メモリ

113 データ復元回路

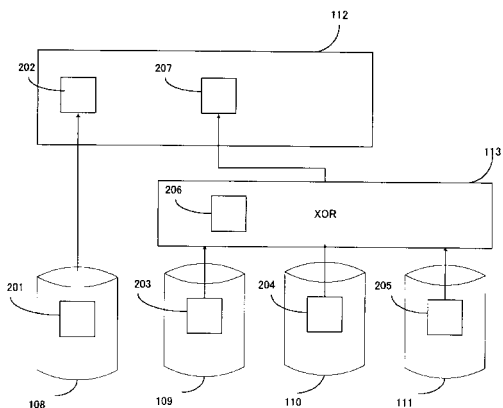
【図1】



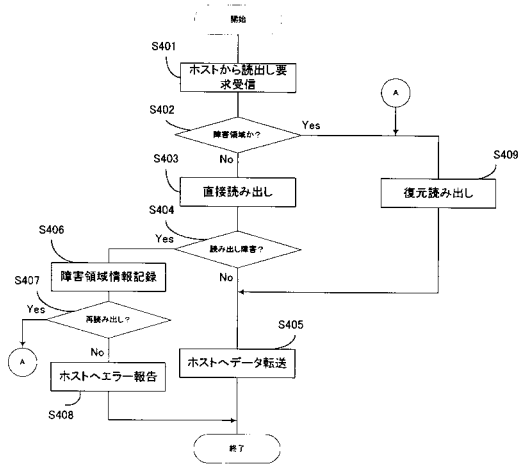
【図3】



【図2】



【 図 4 】



【 図 5 】

