

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4156499号  
(P4156499)

(45) 発行日 平成20年9月24日(2008.9.24)

(24) 登録日 平成20年7月18日(2008.7.18)

(51) Int.Cl. F I  
**G 0 6 F 3/06 (2006.01)**  
 G 0 6 F 3/06 3 0 4 R  
 G 0 6 F 3/06 3 0 5 C  
 G 0 6 F 3/06 5 4 0

請求項の数 31 (全 27 頁)

(21) 出願番号	特願2003-400517 (P2003-400517)	(73) 特許権者	000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号
(22) 出願日	平成15年11月28日(2003.11.28)	(74) 代理人	100079108 弁理士 稲葉 良幸
(65) 公開番号	特開2005-165443 (P2005-165443A)	(74) 代理人	100093861 弁理士 大賀 眞司
(43) 公開日	平成17年6月23日(2005.6.23)	(72) 発明者	加納 東 神奈川県小田原市中里322番2号 株式会社日立製作所 RAIDシステム事業部内
審査請求日	平成18年11月16日(2006.11.16)	(72) 発明者	小河 卓二 神奈川県小田原市中里322番2号 株式会社日立製作所 RAIDシステム事業部内
早期審査対象出願			
前置審査			

最終頁に続く

(54) 【発明の名称】 ディスクアレイ装置

(57) 【特許請求の範囲】

【請求項1】

第一の情報処理装置から書き込み要求又は/及び読み出し要求を受信して、第一の F i b r e Channel - Arbitrated Loop ( F C - A L ) を介して R A I D ( Redundant Array of Inexpensive Disks ) 方式に従ってデータの書き込み又は/及び読み出しを制御する第一のディスクコントローラと、

前記第一の情報処理装置又は/及び第二の情報処理装置から書き込み要求又は/及び読み出し要求を受信して、第二の F C - A L を介して R A I D 方式に従ってデータの書き込み又は/及び読み出しを制御する第二のディスクコントローラと、

前記第一の F C - A L の第一のサーキット及び前記第二の F C - A L の第二のサーキットに接続され且つ前記第一のディスクコントローラ及び前記第二のディスクコントローラにおける書き込み又は読み出し制御に基づいて送受信される F i b r e Channel ( F C ) プロトコルに従った信号とシリアル A T Attachment ( シリアル A T A ) プロトコルに従った信号との間の信号変換を行う変換器と、前記変換器によって変換済みのシリアル A T A プロトコルに従った信号によって受信されるデータを記憶するシリアル A T A ディスクドライブと、を有するシリアル A T A ディスクドライブユニットと、

複数の前記シリアル A T A ディスクドライブユニットと、前記複数のシリアル A T A ディスクドライブユニットへの電力の供給を行う第一の電源モジュールと、を有する第一の筐体と

を備え、

前記第一のサーキットは、前記第一のFC - AL上の経路をバイパスする機能を有し、前記第二のサーキットは、前記第二のFC - AL上の経路をバイパスする機能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】

請求項1に記載のディスクアレイ装置において、

前記第一のサーキットは、ポートバイパスサーキットであり、

前記第一の筐体は、前記第一のFC - ALによる信号の通信に利用され且つ複数の前記第一のサーキットを有する第一の回路基板を有しており、

前記複数のシリアルATAディスクドライブユニットは、前記第一の筐体内において前記第一の回路基板に接続されるものである

ことを特徴とするディスクアレイ装置。

10

【請求項3】

請求項1又は2に記載のディスクアレイ装置において、

前記第二のサーキットは、ポートバイパスサーキットであり、

前記第一の筐体は、前記第二のFC - ALによる信号の通信に利用され且つ複数の前記第二のサーキットを有する第二の回路基板を有しており、

前記複数のシリアルATAディスクドライブユニットは、前記第一の筐体内において前記第二の回路基板に接続されるものである

ことを特徴とするディスクアレイ装置。

20

【請求項4】

請求項1乃至3のいずれかに記載のディスクアレイ装置において、

前記複数のシリアルATAディスクドライブユニットは、少なくとも一つのRAIDグループを形成しており、

前記RAIDグループは、前記第一の情報処理装置又は/及び前記第二の情報処理装置からの書き込み要求又は読み出し要求の宛先となる複数の論理ボリューム識別子に対応付けられるものである

ことを特徴とするディスクアレイ装置。

【請求項5】

請求項1乃至3のいずれかに記載のディスクアレイ装置において、

FCディスクドライブを有するFCディスクドライブユニットと、

複数の前記FCディスクドライブユニットと、前記複数のFCディスクドライブユニットへの電力の供給を行う第二の電源モジュールと、を有する第二の筐体と、を有しており、

前記複数のシリアルATAディスクドライブは、少なくとも一つの第一のRAIDグループを形成しており、

前記複数のFCディスクドライブは、少なくとも一つの第二のRAIDグループを形成しているものである

ことを特徴とするディスクアレイ装置。

30

【請求項6】

請求項5に記載のディスクアレイ装置において、

前記第一のディスクコントローラは、前記第一のRAIDグループ内のシリアルATAディスクドライブと、前記第二のRAIDグループ内のFCディスクドライブと、のいずれのディスクドライブへのデータの書き込み又は/及び読み出しも可能である

ことを特徴とするディスクアレイ装置。

40

【請求項7】

請求項1乃至6のいずれかに記載のディスクアレイ装置において、

前記複数のシリアルATAディスクドライブユニットは、前記第一のFC - ALに接続される第一の接続部と、前記第二のFC - ALに接続される第二の接続部と、を有するものである

50

ことを特徴とするディスクアレイ装置。

【請求項 8】

請求項 5 又は 6 のいずれかに記載のディスクアレイ装置において、

前記第一のディスクコントローラ又は / 及び前記第二のディスクコントローラは、前記第二の RAID グループ内の FC ディスクドライブに格納されるデータには行わないデータの信頼性のチェックを、前記第一の RAID グループ内のシリアル ATA ディスクドライブ内に格納されるデータへ行うものである

ことを特徴とするディスクアレイ装置。

【請求項 9】

請求項 5 又は 6 のいずれかに記載のディスクアレイ装置において、

前記第一のディスクコントローラ又は / 及び前記第二のディスクコントローラは、前記第一の RAID グループ内のシリアル ATA ディスクドライブへのデータの書き込み又は読み出しの際に、前記第二の RAID グループ内の FC ディスクドライブへのデータの書き込み又は読み出し制御と異なる制御を行うものである

ことを特徴とするディスクアレイ装置。

【請求項 10】

第一の情報処理装置から書き込み要求を受信して、第一の Fibre Channel - Arbitrated Loop (FC - AL) を介して RAID (Redundant Array of Inexpensive Disks) 方式に従ってデータの書き込みを制御する第一のディスクコントローラと、

前記第一の情報処理装置又は / 及び第二の情報処理装置から書き込み要求を受信して、第二の FC - AL を介して RAID 方式に従ってデータの書き込みを制御する第二のディスクコントローラと、

前記第一の FC - AL の第一のサーキットに接続される第一の接続部と、前記第二の FC - AL の第二のサーキットに接続される第二の接続部と、前記第一のディスクコントローラから前記第一の接続部を介して受信される Fibre Channel (FC) プロトコルに従った信号及び前記第二のディスクコントローラから前記第二の接続部を介して受信される FC プロトコルに従った信号の両方をシリアル ATA Attachment (シリアル ATA) プロトコルに従った信号へ変換を行う変換器と、前記変換器によって変換済みのシリアル ATA プロトコルに従った信号によって受信されるデータを記憶するシリアル ATA ディスクドライブと、を有する第一のディスクドライブユニットと、

複数の前記第一のディスクドライブユニットと、を有する第一の筐体と

を備え、

前記第一のサーキットは、前記第一の FC - AL 上の経路をバイパスする機能を有し、

前記第二のサーキットは、前記第二の FC - AL 上の経路をバイパスする機能を有することを特徴とするディスクアレイ装置。

【請求項 11】

第一の情報処理装置から書き込み要求を受信して、第一の Fibre Channel - Arbitrated Loop (FC - AL) を介して RAID (Redundant Array of Inexpensive Disks) 方式に従ってデータの書き込みを制御する第一のディスクコントローラと、

前記第一の情報処理装置又は / 及び第二の情報処理装置から書き込み要求を受信して、第二の FC - AL を介して RAID 方式に従ってデータの書き込みを制御する第二のディスクコントローラと、

前記第一の FC - AL の第一のサーキット及び前記第二の FC - AL の第二のサーキットに接続され前記第一のディスクコントローラから第一の接続部を介して受信される Fibre Channel (FC) プロトコルに従った信号をシリアル ATA Attachment (シリアル ATA) プロトコルに従った信号へ変換しかつ前記第二のディスクコントローラから第二の接続部を介して受信される FC プロトコルに従った信号をシリアル ATA プロトコルに従った信号へ変換する変換器と、前記変換器によって変換済みのシリア

10

20

30

40

50

ル A T A プロトコルに従った信号によって受信されるデータを記憶するシリアル A T A ディスクドライブと、を有する第一のディスクドライブユニットと、  
 複数の前記第一のディスクドライブユニットを有する第一の筐体と、  
 を備え、  
前記第一のサーキットは、前記第一の F C - A L 上の経路をバイパスする機能を有し、  
前記第二のサーキットは、前記第二の F C - A L 上の経路をバイパスする機能を有する  
 ことを特徴とするディスクアレイ装置。

【請求項 1 2】

請求項 1 0 又は 1 1 に記載のディスクアレイ装置において、  
 前記第一の筐体は、前記複数の第一のディスクドライブユニットへの電力の供給を行う  
 第一の電源モジュールを備える 10  
 ことを特徴とするディスクアレイ装置。

【請求項 1 3】

請求項 1 0 乃至 1 2 のいずれかに記載のディスクアレイ装置において、  
前記第一のサーキットは、ポートバイパスサーキットであり、  
 前記第一の筐体は、前記第一の F C - A L による信号の通信に利用され、複数の前記第  
 一のサーキットを有する第一の回路基板を有しており、  
 前記複数の第一のディスクドライブユニットは、各々の前記第一の接続部を介して前記  
 第一の回路基板に接続されるものである 20  
 ことを特徴とするディスクアレイ装置。

【請求項 1 4】

請求項 1 0 乃至 1 3 のいずれかに記載のディスクアレイ装置において、  
前記第二のサーキットは、ポートバイパスサーキットであり、  
 前記第一の筐体は、前記第二の F C - A L による信号の通信に利用され、複数の前記第  
 二のサーキットを有する第二の回路基板を有しており、  
 前記複数の第一のディスクドライブユニットは、各々の前記第二の接続部を介して前記  
 第二の回路基板に接続されるものである  
 ことを特徴とするディスクアレイ装置。

【請求項 1 5】

請求項 1 0 乃至 1 4 のいずれかに記載のディスクアレイ装置において、 30  
 前記複数の第一のディスクドライブユニットは、少なくとも一つの R A I D グループを  
 形成しており、  
 前記 R A I D グループは、前記第一の情報処理装置又は / 及び前記第二の情報処理装置  
 からの書き込み要求の宛先となる複数の論理ボリューム識別子に対応付けられるものであ  
 る  
 ことを特徴とするディスクアレイ装置。

【請求項 1 6】

請求項 1 0 乃至 1 4 のいずれかに記載のディスクアレイ装置において、  
 F C ディスクドライブを有する第二のディスクドライブユニットと、  
 複数の前記第二のディスクドライブユニットを有する第二の筐体と、を有しており、 40  
 前記複数のシリアル A T A ディスクドライブは、少なくとも一つの第一の R A I D グル  
 ープを形成しており、  
 前記複数の F C ディスクドライブは、少なくとも一つの第二の R A I D グループを形成  
 しているものである  
 ことを特徴とするディスクアレイ装置。

【請求項 1 7】

請求項 1 6 に記載のディスクアレイ装置において、  
 前記第一のディスクコントローラは、前記第一の R A I D グループ内のシリアル A T A  
 ディスクドライブと、前記第二の R A I D グループ内の F C ディスクドライブと、のいず  
 れのディスクドライブへのデータの書き込みも可能である 50

ことを特徴とするディスクアレイ装置。

【請求項 18】

請求項 10 乃至 17 のいずれかに記載のディスクアレイ装置において、  
前記第一の筐体は、さらに、  
前記第一のディスクコントローラに前記第一の FC - AL を介して接続される第一の  
コネクタと、  
前記第二のディスクコントローラに前記第二の FC - AL を介して接続される第一の  
コネクタと、を有する  
ことを特徴とするディスクアレイ装置。

【請求項 19】

第一の情報処理装置から受信される書き込み要求及び読み出し要求に応じて、第一の Fibre Channel - Arbitrated Loop (FC - AL) を介して RAID (Redundant Array of Inexpensive Disks) 方式に従ってデータの書き込み及び読み出しを制御する第一のディスクコントローラと

、  
前記第一の情報処理装置及び第二の情報処理装置の少なくとも 1 つから受信される書き込み要求及び読み出し要求に応じて、第二の FC - AL を介して RAID 方式に従って データの書き込み及び読み出しを制御する第二のディスクコントローラと、

複数のディスクドライブユニットを含む増設筐体と、  
を備え、

前記複数のディスクドライブユニットの各々は、変換器とシリアル ATA Attachment (シリアル ATA) ディスクドライブとを備えており、

前記変換器は、前記第一の FC - AL の第一のサーキットと前記第二の FC - AL の第二のサーキットとに接続され、前記第一の FC - AL 及び前記第二の FC - AL のファイバチャネルプロトコルと前記シリアル ATA ディスクドライブのシリアル ATA プロトコルとの間のプロトコルの変換を行うものであり、

複数のディスクドライブユニットの各々は、前記第一のディスクコントローラと前記第二のディスクコントローラとの制御に応じてデータが書込まれ又は/及び読み出されるものであり、

前記第一のサーキットは、前記第一の FC - AL 上の経路をバイパスする機能を有し、  
前記第二のサーキットは、前記第二の FC - AL 上の経路をバイパスする機能を有する  
ことを特徴とするディスクアレイシステム。

【請求項 20】

請求項 19 に記載のディスクアレイシステムにおいて、

前記増設筐体は複数の電力供給モジュールを含み、前記複数の電力供給モジュールの各々は、前記複数のディスクドライブユニットに電力を供給するように制御するものであり、

前記第一のディスクコントローラ及び前記第二のディスクコントローラは、前記増設筐体の外に備えられるものである、

ことを特徴とするディスクアレイシステム。

【請求項 21】

請求項 20 に記載のディスクアレイシステムにおいて、

前記増設筐体は、第一のコネクタと、第二のコネクタとを備えており、

前記第一のコネクタは、前記第一のディスクコントローラと前記複数のディスクドライブユニットとの間のデータの中継に利用される前記第一の FC - AL に接続されるものであり、

前記第二のコネクタは、前記第二のディスクコントローラと前記複数のディスクドライブユニットとの間のデータの中継に利用される前記第二の FC - AL に接続されるものである、

ことを特徴とするディスクアレイシステム。

10

20

30

40

50

## 【請求項 2 2】

請求項 1 9 又は 2 0 に記載のディスクアレイシステムにおいて、  
前記ディスクアレイシステムは、前記第一のディスクコントローラと前記第二のディスクコントローラとを含む基本筐体を備えており、

前記第一及び第二のサーキットは、ポートバイパスサーキットであり、

前記増設筐体は、複数の前記第一及び第二のサーキットを備えており、

前記複数の前記第一及び第二のサーキットのうちの一つは、前記複数のディスクドライブユニットのうち少なくとも一つのディスクドライブユニットに対応する前記変換器に接続されるものである、

ことを特徴とするディスクアレイシステム。

10

## 【請求項 2 3】

請求項 1 9 乃至 2 2 のいずれかに記載のディスクアレイシステムにおいて、

前記第一及び第二のサーキットは、ポートバイパスサーキットであり、

前記増設筐体は、前記第一のディスクコントローラに接続される前記第一のサーキットと、前記第二のディスクコントローラに接続される前記第二のサーキットと、を備えている、

ことを特徴とするディスクアレイシステム。

## 【請求項 2 4】

請求項 1 9 乃至 2 2 のいずれかに記載のディスクアレイシステムにおいて、

前記第一及び第二のサーキットは、ポートバイパスサーキットであり、

前記増設筐体は、回路基板を備えており、  
前記回路基板は、複数の前記第一及び第二のサーキットを備えており、前記複数のディスクドライブユニットの各々に対応する前記変換器に接続されるものである、

ことを特徴とするディスクアレイシステム。

20

## 【請求項 2 5】

請求項 1 9 乃至 2 4 のいずれかに記載のディスクアレイシステムにおいて、

前記第一のディスクコントローラと前記第二のディスクコントローラとを有する基本筐体を備える、

ことを特徴とするディスクアレイシステム。

## 【請求項 2 6】

請求項 1 9 乃至 2 5 のいずれかに記載のディスクアレイシステムにおいて、

前記第一のディスクコントローラによって実行され、前記複数のディスクドライブユニットのうち 2 つ以上のディスクドライブユニットを用いて第一の RAID グループを形成するように制御する第一のプログラムと、

前記第二のディスクコントローラによって実行され、前記複数のディスクドライブユニットのうち 2 つ以上のディスクドライブユニットを用いて第二の RAID グループを形成するように制御する第二のプログラムとを備える、

ことを特徴とするディスクアレイシステム。

30

## 【請求項 2 7】

請求項 1 9 乃至 2 6 のいずれかに記載のディスクアレイシステムにおいて、

前記ディスクアレイシステムは、ファイバチャネルインタフェースを有するファイバチャネルディスクドライブを各々に有する複数の他のディスクドライブユニットと、前記複数の他のディスクドライブユニットに対する電力の供給を制御する複数の第二の電力供給モジュールと、を備える他の増設筐体を備える、

ことを特徴とするディスクアレイシステム。

40

## 【請求項 2 8】

請求項 2 7 に記載のディスクアレイシステムにおいて、

前記第一のディスクコントローラは、前記複数のディスクドライブユニット内の複数の前記シリアル ATA ディスクドライブ及び前記複数の他のディスクドライブユニット内の複数の前記ファイバチャネルディスクドライブに対して、データの読み出し及び書き込み

50

を制御するものである、

ことを特徴とするディスクアレイシステム。

【請求項 29】

請求項 28 に記載のディスクアレイシステムにおいて、

前記第二のディスクコントローラは、前記複数のディスクドライブユニット内の複数の前記シリアル ATA ディスクドライブ及び前記複数の他のディスクドライブユニット内の複数の前記ファイバチャネルディスクドライブに対して、データの読み出し及び書き込みを制御することが可能なものである、

ことを特徴とするディスクアレイシステム。

【請求項 30】

請求項 19 乃至 29 のいずれかに記載のディスクアレイシステムにおいて、

前記第一のディスクコントローラによって実行され、第一の論理ユニット番号 (LUN) に対応する第一の論理ボリュームを形成するように制御する第一のコンピュータプログラムと、

前記第二のディスクコントローラによって実行され、第二の LUN に対応する第二の論理ボリュームを形成するように制御する第二のコンピュータプログラムと、

を備えることを特徴とするディスクアレイシステム。

【請求項 31】

請求項 19 乃至 30 のいずれかに記載のディスクアレイシステムにおいて、

前記複数のディスクドライブユニットの各々は、前記増設筐体に対して挿入されるものであり、かつ前記増設筐体から取り外すことが可能なものである、

ことを特徴とするディスクアレイシステム。

【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及びディスクアレイ装置の制御方法に関する。

【背景技術】

【0002】

近年、ディスクアレイ装置における記憶容量の増大に伴い、情報処理システムにおける重要性は益々高まってきている。そこで、情報処理装置等からのデータ入出力要求に対して、要求された位置に正しくデータの書き込みを行うこと、読み出したデータが不正である場合にはそれを検知することが重要である。

【0003】

特許文献 1 においては、磁気ディスク装置に 2 つのヘッドを持たせ、同一のデータを 2 つのヘッドから読み出して比較することにより、磁気ディスク装置における書き込み及び読み出しの信頼性を高める方法が開示されている。

【特許文献 1】特開平 5 - 150909 号公報

【発明の開示】

【発明が解決しようとする課題】

【0004】

特許文献 1 の方法をディスクアレイ装置に適用する場合、各磁気ディスクにヘッドを 2 つ持たせる必要があるため、ハードディスクドライブの製造単価が高くなってしまふ。そこで、ヘッドの追加等の物理的な構造の変更を行うことなく、ハードディスクドライブにおける信頼性を高める方法が求められている。

【0005】

また、ディスクアレイ装置においては、ファイバチャネルのハードディスクドライブに加えて、シリアル ATA やパラレル ATA 等のハードディスクドライブ等も利用されはじめている。これは、シリアル ATA やパラレル ATA 等のハードディスクドライブは、ファイバチャネルのハードディスクドライブと比較して信頼性は劣るが価格が低いためである。そこで、このようにファイバチャネルとシリアル ATA 等の規格のハードディスクド

10

20

30

40

50

ライブを組み合わせて構成するディスクアレイ装置において、ファイバチャネル以外のハードディスクドライブにおける信頼性を高める方法が求められている。

【0006】

本発明は上記課題を鑑みてなされたものであり、ディスクアレイ装置及びディスクアレイ装置の制御方法を提供することを目的とする。

【課題を解決するための手段】

【0007】

上記目的を達成するため本発明のうち主たる発明に係るディスクアレイ装置は、第一の情報処理装置から書き込み要求又は/及び読み出し要求を受信して、第一のFibre Channel - Arbitrated Loop (FC - AL)を介してRAID (Redundant Array of Inexpensive Disks)方式に従ってデータの書き込み又は/及び読み出しを制御する第一のディスクコントローラと、前記第一の情報処理装置又は/及び第二の情報処理装置から書き込み要求又は/及び読み出し要求を受信して、第二のFC - ALを介してRAID方式に従ってデータの書き込み又は/及び読み出しを制御する第二のディスクコントローラと、前記第一のFC - ALの第一のサーキット及び前記第二のFC - ALの第二のサーキットに接続され前記第一のディスクコントローラ及び前記第二のディスクコントローラにおける書き込み又は読み出し制御に基づいて送受信されるFibre Channel (FC)プロトコルに従った信号とシリアルATA Attachment (シリアルATA)プロトコルに従った信号との間の信号変換を行う変換器と、前記変換器によって変換済みのシリアルATAプロトコルに従った信号によって受信されるデータを記憶するシリアルATAディスクドライブと、を有するシリアルATAディスクドライブユニットと、複数の前記シリアルATAディスクドライブユニットと、前記複数のシリアルATAディスクドライブユニットへの電力の供給を行う第一の電源モジュールと、を有する第一の筐体とを備え、前記第一のサーキットは、前記第一のFC - AL上の経路をバイパスする機能を有し、前記第二のサーキットは、前記第二のFC - AL上の経路をバイパスする機能を有することを特徴とする。

【0008】

また本発明においては、第一の情報処理装置から書き込み要求を受信して、第一のFibre Channel - Arbitrated Loop (FC - AL)を介してRAID (Redundant Array of Inexpensive Disks)方式に従ってデータの書き込みを制御する第一のディスクコントローラと、前記第一の情報処理装置又は/及び第二の情報処理装置から書き込み要求を受信して、第二のFC - ALを介してRAID方式に従ってデータの書き込みを制御する第二のディスクコントローラと、前記第一のFC - ALの第一のサーキットに接続される第一の接続部と、前記第二のFC - ALの第二のサーキットに接続される第二の接続部と、前記第一のディスクコントローラから前記第一の接続部を介して受信されるFibre Channel (FC)プロトコルに従った信号及び前記第二のディスクコントローラから前記第二の接続部を介して受信されるFCプロトコルに従った信号の両方をシリアルATA Attachment (シリアルATA)プロトコルに従った信号へ変換を行う変換器と、前記変換器によって変換済みのシリアルATAプロトコルに従った信号によって受信されるデータを記憶するシリアルATAディスクドライブと、を有する第一のディスクドライブユニットと、複数の前記第一のディスクドライブユニットと、を有する第一の筐体とを備え、前記第一のサーキットは、前記第一のFC - AL上の経路をバイパスする機能を有し、前記第二のサーキットは、前記第二のFC - AL上の経路をバイパスする機能を有することを特徴とする。

【0009】

さらに本発明においては、第一の情報処理装置から書き込み要求を受信して、第一のFibre Channel - Arbitrated Loop (FC - AL)を介してRAID (Redundant Array of Inexpensive Disks)方式に従ってデータの書き込みを制御する第一のディスクコントローラと、前記第一の

10

20

30

40

50



情報処理装置又は/及び第二の情報処理装置から書き込み要求を受信して、第二のFC - ALを介してRAID方式に従ってデータの書き込みを制御する第二のディスクコントローラと、前記第一のFC - ALの第一のサーキット及び前記第二のFC - ALの第二のサーキットに接続され前記第一のディスクコントローラから第一の接続部を介して受信されるFibre Channel (FC) プロトコルに従った信号をシリアルATA Attachment (シリアルATA) プロトコルに従った信号へ変換し、かつ前記第二のディスクコントローラから第二の接続部を介して受信されるFC プロトコルに従った信号をシリアルATA プロトコルに従った信号へ変換する変換器と、前記変換器によって変換済みのシリアルATA プロトコルに従った信号によって受信されるデータを記憶するシリアルATA ディスクドライブと、を有する第一のディスクドライブユニットと、複数の前記第一のディスクドライブユニットを有する第一の筐体と、を備え、前記第一のサーキットは、前記第一のFC - AL上の経路をバイパスする機能を有し、前記第二のサーキットは、前記第二のFC - AL上の経路をバイパスする機能を有することを特徴とする。

10

## 【0010】

さらに本発明においては、第一の情報処理装置から受信される書き込み要求及び読み出し要求に応じて、第一のFibre Channel - Arbitrated Loop (FC - AL) を介してRAID (Redundant Array of Inexpensive Disks) 方式に従ってデータの書き込み及び読み出しを制御する第一のディスクコントローラと、前記第一の情報処理装置及び第二の情報処理装置の少なくとも1つから受信される書き込み要求及び読み出し要求に応じて、第二のFC - ALを介してRAID方式に従ってデータの書き込み及び読み出しを制御する第二のディスクコントローラと、複数のディスクドライブユニットを含む増設筐体と、を備え、前記複数のディスクドライブユニットの各々は、変換器とシリアルATA Attachment (シリアルATA) ディスクドライブとを備えており、前記変換器は、前記第一のFC - ALの第一のサーキットと前記第二のFC - ALの第二のサーキットとに接続され、前記第一のFC - AL及び前記第二のFC - ALのファイバチャネルプロトコルと前記シリアルATA ディスクドライブのシリアルATA プロトコルとの間のプロトコルの変換を行うものであり、複数のディスクドライブユニットの各々は、前記第一のディスクコントローラと前記第二のディスクコントローラとの制御に応じてデータが書込まれ又は/及び読み出されるものであり、前記第一のサーキットは、前記第一のFC - AL上の経路をバイパスする機能を有し、前記第二のサーキットは、前記第二のFC - AL上の経路をバイパスする機能を有することを特徴とする。

20

30

## 【0011】

その他、本願が開示する課題、及びその解決方法は、発明を実施するための最良の形態の欄、及び図面により明らかにされる。

## 【発明の効果】

## 【0012】

ディスクアレイ装置及びディスクアレイ装置の制御方法を提供することができる。

## 【発明を実施するための最良の形態】

## 【0013】

= = 装置構成 = =

図1(a)は本発明の一実施例として説明するディスクアレイ装置10の正面図であり、図1(b)はディスクアレイ装置10の背面図である。図2(a)は、ディスクアレイ装置10に装着される基本筐体20を正面側から見た斜視図であり、図2(b)は基本筐体20を背面側から見た斜視図である。図3(a)は、ディスクアレイ装置10に装着される増設筐体30を正面側から見た斜視図であり、図3(b)は増設筐体30を背面側から見た斜視図である。

40

## 【0014】

図1(a), (b)に示すように、ディスクアレイ装置10は、ラックフレーム11をベースとして構成される。ラックフレーム11の内側左右側面の上下方向には、複数段に

50

わたって前後方向にマウントフレーム 1 2 が形成され、このマウントフレーム 1 2 に沿って基本筐体 2 0 および増設筐体 3 0 が引き出し式に装着される。図 2 ( a ) , ( b ) に示すように、基本筐体 2 0 および増設筐体 3 0 には、ディスクアレイ装置 1 0 の各種機能を提供するボードやユニットが装着されている。

【 0 0 1 5 】

図 2 ( a ) に示すように、基本筐体 2 0 の正面上段側には、ハードディスクドライブ 5 1 が装填された複数のディスクドライブユニット 5 2 が並べて装着されている。ハードディスクドライブ 5 1 は、例えば、FC - AL 規格、SCSI 1 ( Small Computer System Interface 1 ) 規格、SCSI 2 規格、SCSI 3 規格、パラレル ATA ( AT Attachment ) 規格、シリアル ATA 規格などの通信機能を提供する通信インタフェースを有するハードディスクドライブである。

10

【 0 0 1 6 】

また、基本筐体 2 0 の正面下段側には、バッテリーユニット 5 3、ハードディスクドライブ 5 1 の稼働状態などが表示される表示パネル 5 4、フレキシブルディスクドライブ 5 5 が装着されている。バッテリーユニット 5 3 には二次電池が内蔵されている。バッテリーユニット 5 3 は、停電などにより AC / DC 電源 5 7 からの電力供給が途絶えた場合に、ボードやユニットに電力を供給するバックアップ電源として機能する。表示パネル 5 4 には、ハードディスクドライブ 5 1 の稼働状態などを表示する LED ランプなどの表示デバイスが設けられている。フレキシブルディスクドライブ 5 5 は、メンテナンス用プログラムをロードする場合などに用いられる。

20

【 0 0 1 7 】

図 2 ( b ) に示すように、基本筐体 2 0 の背面上段側の両側面側には、1 枚ずつ電源コントローラボード 5 6 が装着されている。電源コントローラボード 5 6 は、複数のハードディスクドライブ 5 1 と通信可能に接続されている。電源コントローラボード 5 6 と複数のハードディスクドライブ 5 1 は、ループ状の通信経路、例えば、FC - AL の方式 ( トポロジー ) で通信を行う通信経路によって通信可能に接続されている。

【 0 0 1 8 】

電源コントローラボード 5 6 は、AC / DC 電源 5 7 の状態監視やハードディスクドライブ 5 1 の状態監視、ハードディスクドライブ 5 1 の電源供給の制御、冷却装置の冷却能力の制御、表示パネル 5 4 上の表示デバイスの制御、筐体各部の温度監視などを行う回路が実装されている。なお、冷却装置は、ディスクアレイ装置 1 0 内や筐体 2 0 , 3 0 内を冷却する装置であり、例えば、インタークーラー、ヒートシンク、空冷式の冷却ファンなどである。電源コントローラボード 5 6 にはファイバチャネルケーブルのコネクタ 6 7 が設けられ、このコネクタにはファイバチャネルケーブル 9 1 が接続される。

30

【 0 0 1 9 】

図 2 ( b ) に示すように、基本筐体 2 0 の背面上段側の前記 2 枚の電源コントローラボード 5 6 に挟まれた空間には、AC / DC 電源 5 7 が 2 台並べて装着されている。AC / DC 電源 5 7 は、ハードディスクドライブ 5 1、ボード、ユニットなどに電源を供給する。AC / DC 電源 5 7 は、電源コントローラボード 5 6 と接続されており、電源コントローラボード 5 6 からの信号により各ハードディスクドライブ 5 1 に電源を供給できるように設定されている。

40

【 0 0 2 0 】

なお、本実施の形態においては、各筐体 2 0 , 3 0 の電源供給に関するセキュリティを確保するために、基本筐体 2 0 および増設筐体 3 0 に電源コントローラボード 5 6 と AC / DC 電源 5 7 とを各 2 台ずつ冗長に装着させることとしているが、電源コントローラボード 5 6 と AC / DC 電源 5 7 とを各 1 台ずつ装着させることとしてもよい。

AC / DC 電源 5 7 には、AC / DC 電源 5 7 の出力をオン・オフするためのブレーカスイッチ 6 4 が設けられている。

【 0 0 2 1 】

図 2 ( b ) に示すように、AC / DC 電源 5 7 の下方には、2 台の空冷式の冷却ファン

50

ユニット 5 8 が並べて装着されている。冷却ファンユニット 5 8 には、1 台以上の冷却ファン 6 6 が実装されている。冷却ファン 6 6 は、筐体内に空気を流入・流出させることでハードディスクドライブ 5 1 や A C / D C 電源 5 7 などから発生する熱を筐体外部に排出する。なお、基本筐体 2 0 や増設筐体 3 0、およびこれらに装着されているボードやユニットには、筐体 2 0、3 0 内に空気を循環させる通気路や通気口が形成され、冷却ファン 6 6 により筐体 2 0 内の熱が外部に効率よく排出される仕組みになっている。冷却ファン 6 6 は、ハードディスクドライブ 5 1 ごとに設けることとしてもよいが、チップやユニットの数を削減できることから筐体ごとに大きな冷却ファン 6 6 を設けることが好ましい。

【 0 0 2 2 】

冷却ファンユニット 5 8 は、コントローラボード 5 9 もしくは電源コントローラボード 5 6 と制御ライン 4 8 で接続されており、冷却ファンユニット 5 8 の冷却ファン 6 6 の回転数は、この制御ライン 4 8 を通じてコントローラボード 5 9 もしくは電源コントローラボード 5 6 により制御される。

【 0 0 2 3 】

図 2 ( b ) に示すように、基本筐体 2 0 の背面下段側には、1 枚のコントローラボード 5 9 が装着されている。コントローラボード 5 9 には、基本筐体 2 0 および増設筐体 3 0 に装着されているハードディスクドライブ 5 1 との間の通信インタフェースと、ハードディスクドライブ 5 1 の動作の制御 ( 例えば、R A I D 方式による制御 ) やハードディスクドライブ 5 1 の状態監視を行う回路などが実装されている。

【 0 0 2 4 】

なお、本実施の形態において、電源コントローラボード 5 6 がハードディスクドライブ 5 1 の電源供給の制御や冷却装置の冷却能力の制御を行うこととしているが、これらの制御をコントローラボード 5 9 が行うこととしてもよい。

【 0 0 2 5 】

また、本実施例においては、コントローラボード 5 9 は、情報処理装置 3 0 0 との間の通信インタフェースの機能、例えば、S C S I 規格やファイバチャネル規格の通信機能を提供する通信インタフェースボード 6 1 や、ハードディスクドライブ 5 1 への書き込みデータや読み出しデータが記憶されるキャッシュメモリ 6 2 などを実装する形態としているが、これらを別のポートが実装する形態としてもよい。

【 0 0 2 6 】

コントローラボード 5 9 に装着される通信インタフェースボード 6 1 には、情報処理装置 3 0 0 と接続するための、ファイバチャネル、Ethernet ( 登録商標 ) などのプロトコルで構築された S A N ( Storage Area Network )、L A N ( Local Area Network )、もしくは、S C S I などの所定のインタフェース規格に準拠した外部コネクタ 6 3 が設けられ、ディスクアレイ装置 1 0 は、このコネクタ 6 3 に接続される通信ケーブル 9 2 を介して情報処理装置 3 0 0 と接続される。

【 0 0 2 7 】

なお、基本筐体 2 0 のハードディスクドライブ 5 1 の制御に関するセキュリティを確保するために、2 枚のコントローラボード 5 9 を冗長に装着させることとしてもよい。

【 0 0 2 8 】

図 3 ( a ) に示すように、増設筐体 3 0 の正面側には、ハードディスクドライブ 5 1 が収容された複数のディスクドライブユニット 5 2 が並べて装着されている。図 3 ( b ) に示すように、増設筐体 3 0 の背面両側面側には、それぞれ一枚ずつ電源コントローラボード 5 6 が装着されている。また、2 枚の電源コントローラボード 5 6 に挟まれた空間には、A C / D C 電源 5 7 が 2 台並べて装着されている。また、A C / D C 電源 5 7 の下方には、2 台の冷却ファンユニット 5 8 が並べて装着されている。A C / D C 電源 5 7 には、A C / D C 電源 5 7 の出力をオン・オフするためのブレーカスイッチ 6 4 が設けられている。

【 0 0 2 9 】

本実施の形態においては、上述したように増設筐体 3 0 の電源供給に関するセキュリテ

10

20

30

40

50

ィを確保するために、増設筐体30に電源コントローラボード56とAC/DC電源57とを各2台ずつ冗長に装着させることとしているが、電源コントローラボード56とAC/DC電源57とを各1台ずつ装着させることとしてもよい。なお、ハードディスクドライブ51の電源供給の制御や冷却装置の冷却能力の制御などの電源コントローラボード56の機能をコントローラボード59に実装することとしてもよい。

【0030】

図4にディスクドライブユニット52に収容されているハードディスクドライブ51の構成の一例を示す。ハードディスクドライブ51は、その筐体70内に、磁気ディスク73、アクチュエータ71、スピンドルモータ72、データの読み書きを行うヘッド74、ヘッド74等の機構部分を制御する機構制御回路75、磁気ディスク73へのデータの読み書き信号を制御する信号処理回路76、通信インタフェース回路77、各種コマンドやデータが入出力されるインタフェースコネクタ79、電源コネクタ80等を備えて構成される。なお、通信インタフェース回路77には、データを一時的に格納するためのキャッシュメモリも含まれている。なお、後述するコントローラ500におけるキャッシュメモリ62と区別するため、ハードディスクドライブ51が備えるキャッシュメモリをディスクキャッシュと称する。

【0031】

ハードディスクドライブ51は、例えば、コンタクトスタートストップ(CSS: Contact Start Stop)方式の3.5インチサイズの磁気ディスクや、ロード/アンロード方式の2.5インチサイズの磁気ディスクなどを備える記憶装置である。3.5インチサイズの磁気ディスクは、例えば、SCSI1、SCSI2、SCSI3、FC-ALなどの通信インタフェースを有している。一方、2.5インチサイズの磁気ディスクは、例えば、パラレルATA、シリアルATAなどの通信インタフェースを有している。

【0032】

2.5インチサイズの磁気ディスクをディスクアレイ装置10の筐体20,30に収容する場合には、3.5インチの形状をした容器に収めるようにしてもよい。これにより、磁気ディスクの衝撃耐力性能を向上させることが可能となる。なお、2.5インチサイズの磁気ディスクと3.5インチサイズの磁気ディスクとは、通信インタフェースが異なるだけでなく、I/O性能、消費電力、寿命の点などで異なっている。2.5インチサイズの磁気ディスクは、3.5インチサイズの磁気ディスクに比べ、I/O性能が優れておらず、寿命が短い。しかし、3.5インチサイズの磁気ディスクに比べ、消費電力が少ないという点で優れている。

【0033】

==ディスクアレイ装置のハードウェア構成==

図5は、本発明の一実施例として説明するディスクアレイ装置10のハードウェア構成を示すブロック図である。

【0034】

図5に示すように、ディスクアレイ装置10には、SANを介して情報処理装置300が接続されている。情報処理装置300は、例えば、パーソナルコンピュータ、ワークステーション、メインフレームコンピュータなどである。

【0035】

ディスクアレイ装置10は、前述したように、基本筐体20と1つ又は複数の増設筐体30を備えている。本実施の形態においては、基本筐体20は、コントローラ500、ハードディスクドライブ51などを備えている。コントローラ500は、チャンネル制御部501、ディスク制御部502、CPU503、メモリ504、キャッシュメモリ62、及びデータコントローラ505などを備えており、前述のコントローラボード59に実装されている。また、増設筐体30は、ハードディスクドライブ51などを備えている。基本筐体及び増設筐体のハードディスクドライブ51は、FC-AL506により、ディスク制御部502と通信可能に接続されている。なお、ディスク制御部502とハードディスクドライブ51との接続形態の詳細については後述する。

## 【 0 0 3 6 】

チャンネル制御部 5 0 1 は情報処理装置 3 0 0 との間で通信を行うインタフェースである。チャンネル制御部 5 0 1 は、ファイバチャンネルプロトコルに従ってブロックアクセス要求を受け付ける機能を有する。

## 【 0 0 3 7 】

ディスク制御部 5 0 2 は、CPU 5 0 3 からの指示により、ハードディスクドライブ 5 1 との間でデータのやりとりを行うインタフェースである。ディスク制御部 5 0 2 は、ハードディスクドライブ 5 1 を制御するコマンドなどを規定するプロトコルに従ってハードディスクドライブ 5 1 に対するデータ入出力要求を送信する機能を備える。

## 【 0 0 3 8 】

CPU 5 0 3 は、ディスクアレイ装置 1 0 の全体の制御を司るもので、メモリ 5 0 4 に格納されたマイクロプログラムを実行することにより、チャンネル制御部 5 0 1、ディスク制御部 5 0 2、及びデータコントローラ 5 0 6 等の制御を行う。マイクロプログラムとは、図 6 に示すデータ READ 処理 6 0 1 やデータ WRITE 処理 6 0 2 などである。

## 【 0 0 3 9 】

キャッシュメモリ 6 2 は、チャンネル制御部 5 0 1 とディスク制御部 5 0 2 との間で授受されるデータを一時的に記憶するために用いられる。

## 【 0 0 4 0 】

データコントローラ 5 0 5 は、CPU 5 0 3 の制御によりチャンネル制御部 5 0 1 とキャッシュメモリ 6 2 との間又はキャッシュメモリ 6 2 とディスク制御部 5 0 2 との間のデータ転送を行うものである。

## 【 0 0 4 1 】

コントローラ 5 0 0 は、ハードディスクドライブ 5 1 をいわゆる R A I D (Redundant Array of Inexpensive Disks) 方式に規定される R A I D レベル (例えば、0, 1, 5) で制御する機能を備えている。R A I D 方式においては、複数のハードディスクドライブ 5 1 が 1 つのグループ (以後、R A I D グループと称する) として管理されている。R A I D グループ上には、情報処理装置 3 0 0 からのアクセス単位である論理ボリュームが形成されており、各論理ボリュームには L U N (Logical Unit Number) と呼ばれる識別子が付与されている。R A I D の構成情報は、図 6 に示すようにメモリ 5 0 4 に R A I D 構成テーブル 6 0 3 として記憶されており、データ READ 処理 6 0 1 やデータ WRITE 処理 6 0 2 の実行時に CPU 5 0 3 により参照される。

## 【 0 0 4 2 】

なお、ディスクアレイ装置は、以上に説明した構成のもの以外にも、例えば、N F S (Network File System) などのプロトコルにより情報処理装置 3 0 0 からファイル名指定によるデータ入出力要求を受け付けるように構成された N A S (Network Attached Storage) として機能するものなどであってもよい。

## 【 0 0 4 3 】

== ハードディスクドライブの接続形態 ==

次に、コントローラ 5 0 0 とハードディスクドライブ 5 1 との接続形態について説明する。

図 7 は、基本筐体 2 0 にファイバチャンネルのハードディスクドライブ 5 1 が収容されている場合における、ディスク制御部 5 0 2 と各ハードディスクドライブ 5 1 との接続形態を示している。

## 【 0 0 4 4 】

ディスク制御部 5 0 2 は、F C - A L 5 0 6 で複数のハードディスクドライブ 5 1 と接続されている。F C - A L 5 0 6 は、複数の P B C (Port Bypass Circuit) 7 0 1 を備えている。ファイバチャンネルのハードディスクドライブ 5 1 は、この P B C 7 0 1 を介して F C - A L 5 0 6 に接続されている。P B C 7 0 1 は、チップ化された電子スイッチであり、ディスク制御部 5 0 2 やハードディスクドライブ 5 1 などバイパスし電氣的に F C - A L 5 0 6 から除外する機能も有している。具体的には、P B C 7 0 1 は、障害が発

10

20

30

40

50

生したハードディスクドライブ51をFC - AL506から切り離して、他のハードディスクドライブ51とディスク制御部502との間の通信を可能にする。

【0045】

また、PBC701は、FC - AL506の動作を維持したままでハードディスクドライブ51の抜き差しを可能にする。例えば、ハードディスクドライブ51が新たに装着された場合にはそのハードディスクドライブ51をFC - AL506に取り込み、ディスク制御部502との間の通信を可能にする。なお、PBC701の回路基板は、ディスクアレイ装置10のラックフレーム11に設けられているか、もしくは、その一部または全部がコントローラボード59や電源コントローラボード56に実装されていることとしてもよい。

10

【0046】

図8は、基本筐体20にシリアルATAのハードディスクドライブ51が収容されている場合における、ディスク制御部502と各ハードディスクドライブ51との接続形態を示している。

各ハードディスクドライブ51は、コンバータ801を介してFC - AL506のPBC701に接続されている。コンバータ801はファイバチャネルプロトコルとシリアルATAプロトコルとを変換する回路である。コンバータ801は、プロトコル変換機能が組み込まれた1つのチップであり、各ディスクドライブユニット52内に設けられている。

【0047】

20

図9は、基本筐体20にシリアルATAのハードディスクドライブ51が収容されている場合における、もう一つの接続形態を示している。

コンバータ901は、図7におけるコンバータ801と同様にファイバチャネルプロトコルとシリアルATAプロトコルとを変換する回路である。コンバータ901はFC - AL506のPBC701に接続されており、1つのコンバータ901には、複数のハードディスクドライブ51がスイッチ902を介して接続されている。スイッチ902は、ハードディスクドライブ51が複数のコンバータ901に接続されている場合において、どのハードディスクドライブ51と通信を行うかを選択する回路である。スイッチ902は、各ディスクドライブユニット52内に設けられている。コンバータ901は、プロトコル変換機能が組み込まれた1つのチップであるか、複数の回路により構成されている。コンバータ901は、例えば「米国特許出願公開第2003/0135577号明細書」にて開示されているSATAマスタデバイスの構成により実現することが可能である。コンバータ901は、コントローラボード59や電源コントローラボード56等を実装されている。

30

【0048】

== 信頼性を高めるための制御 ==

以上に説明したディスクアレイ装置10において、ハードディスクドライブ51からの読み出し又はハードディスクドライブ51への書き込みの信頼性を高める方法について説明する。

【0049】

40

== RAID構成でのパリティチェック ==

まず、RAID構成においてハードディスクドライブ51に記憶されているデータが不正な状態となっていないか検査する方法について説明する。ここで、不正な状態とは、データがディスク制御部502から指定された場所に指定された内容で書き込まれていない状態のことである。

【0050】

図10は、RAID5においてハードディスクドライブ51にデータが記憶されている様子を表している。RAID5においては、複数のハードディスクドライブ51によりRAIDグループ1001が形成されている。図10の例では、ハードディスクドライブ51には、データA～DとデータA～Dに対する誤りを検出するためのパリティデータP(

50

A - D) が記憶されている。同様に、データ E ~ H とデータ E ~ H に対するパリティデータ P ( E - H ) が記憶されている。このような、データとパリティデータとの組合せのことを、ストライプグループ 1 0 0 2 と呼ぶこととする。このようなストライプグループ 1 0 0 2 が形成されている RAID 構成において、コントローラ 5 0 0 はストライプグループ 1 0 0 2 の全てのデータ及びパリティデータを読み出すことにより、データが不正な状態となっていないか検査することができる。まず、CPU 5 0 3 からの指示により、ディスク制御部 5 0 2 がデータ A ~ D とパリティデータ P ( A - D ) とを読み出す。次に、CPU 5 0 3 は、データ A ~ D とパリティデータ P ( A - D ) とでパリティチェックを行うことにより、データ A ~ D のいずれかが不正な状態となっていないか検査することができる。

10

#### 【 0 0 5 1 】

コントローラ 5 0 0 は、情報処理装置 3 0 0 からデータの読み出し要求を受信した際に、当該データを含むストライプグループの全てのデータとパリティデータとを読み出すようにすることもできる。これにより、コントローラ 5 0 0 がハードディスクドライブ 5 1 から不正なデータを読み出して情報処理装置 3 0 0 に送信することを防止することが可能となる。なお、不正なデータの検査はデータの読み出し要求を受信した際にかかわらず、任意のタイミングで行うこととしてもよい。これにより、データの読み出し性能に影響を与えずに、不正なデータの検出を行うことが可能となる。

#### 【 0 0 5 2 】

また、図 1 1 に示す更新管理テーブル 1 1 0 1 を用いて、ハードディスクドライブ 5 1 に書き込まれたデータが不正な状態となっていないか検査することができる。更新管理テーブル 1 1 0 1 はドライブ番号とセクタ番号とで構成され、メモリ 5 0 4 に記憶されている。本実施の形態においてはセクタ番号は L B A ( Logical Block Address ) で定義されており、L B A # 1 - 1 2 8 のように 1 2 8 L B A 単位で管理されている。なお、セクタ番号をまとめる単位は 1 2 8 に限らず任意の単位でよい。CPU 5 0 3 は、ディスク制御部 5 0 2 を介してデータをハードディスクドライブ 5 1 に書き込むと、更新管理テーブル 1 1 0 1 の当該ハードディスクドライブ 5 1 の書き込みを行ったセクタの値を「 1 」に変更する。CPU 5 0 3 は、更新管理テーブル 1 1 0 1 において「 1 」が記憶されているハードディスクドライブ 5 1 の対象セクタのストライプグループの全てのデータとパリティデータとをディスク制御部 5 0 2 を介して読み出し、パリティチェックを行う。CPU 5 0 3 は、読み出したデータが不正でない場合は、更新管理テーブル 1 1 0 1 において当該セクタの値を「 0 」に変更する。CPU 5 0 3 は、チャンネル制御部 5 0 1 を介して情報処理装置 3 0 0 からデータの読み出し要求を受信すると、更新管理テーブル 1 1 0 1 を参照して当該データが記憶されているセクタが検査済みであるかどうかを確認する。当該データが記憶されているセクタが検査済みでない場合は、CPU 5 0 3 は前述の手順に従い当該データの属するストライプグループのデータを検査する。このように、ハードディスクドライブ 5 1 に書き込まれたデータについて、当該データに対する読み出し要求を受信する前に当該データの検査を実施しておくことにより、データの読み出し性能の低下を抑えることが可能となる。また、検査の未済を更新管理テーブル 1 1 0 1 に記憶し、検査が行われていないデータを読み出す際にはパリティチェックを行うため、不正データの読み出しを防止することが可能となる。

20

30

40

#### 【 0 0 5 3 】

== W R I T E データに対する検査 ==

次に、ハードディスクドライブ 5 1 にデータを書き込んだ際に、当該データが正しく書き込まれているか検査する方法について説明する。

#### 【 0 0 5 4 】

図 1 2 は、コントローラ 5 0 0 がハードディスクドライブ 5 1 にデータを書き込む際の CPU 5 0 3 での制御を表すフローチャートである。CPU 5 0 3 は、チャンネル制御部 5 0 1 を介して情報処理装置 3 0 0 からデータの書き込み要求を受信すると、当該データのハードディスクドライブ 5 1 への書き込み指示をディスク制御部 5 0 2 に送信する ( S 1

50

201)。そして、CPU503は当該データが書き込まれた磁気ディスクのヘッドの位置を移動させるシーク処理の実行指示をディスク制御部に送信する(S1202)。次に、CPU503は、キャッシュメモリ62から当該データを読み出し(S1203)、磁気ディスクから当該データを読み出す(S1204)。CPU503は、キャッシュメモリ62のデータと磁気ディスクのデータとが一致しているか比較する(S1205)。2つのデータが一致していない場合、CPU503は、書き込みが正常に行われていないことを情報処理装置300に通知する(S1206)。

このように、磁気ディスクに記憶されているデータとキャッシュメモリ62に記憶されているデータとを比較することにより、磁気ディスクに正しくデータが書き込まれているか確認することが可能である。また、書き込まれているデータが不正な状態となっている場合においても、データがキャッシュメモリ62に残っているため、データを失うことがない。なお、比較するデータを磁気ディスクとキャッシュメモリ62とから読み出す前に、シーク処理等により当該ハードディスクドライブが備えるヘッドを移動することにより、書き込み時にヘッドの位置が不正であった場合に、同じ位置から再度読み出すことを防止することが可能となる。

#### 【0055】

図12の処理においては、書き込まれたデータの全部をキャッシュメモリ62と磁気ディスクとから読み出し比較することによりデータの検査を実施したが、データの全部ではなく、先頭と末尾の1セグメント等、そのデータの一部を読み出して比較することとしてもよい。例えば、シリアルATAのハードディスクドライブは、データのバックアップ等の用途に用いられるため、サイズの大きいデータ(シーケンシャルデータ)が書き込まれることが多い。このような場合、書き込みを行ったデータの全部について、磁気ディスクに記憶されているデータとキャッシュメモリ62に記憶されているデータとを比較することは、書き込み処理の性能を著しく低下させる要因となる。また、シーケンシャルデータの書き込み時に書き込み位置の誤り等が発生した場合は、そのデータの全部が不正となっている可能性が高い。そのため、データの一部を検査することでデータが不正となっているか判断可能である場合が多い。つまり、書き込みを行ったデータの一部、例えば先頭と末尾の1セグメント等について比較することにより、書き込み処理の性能低下を抑えたうえで、不正データのチェックを行うことが可能である。

#### 【0056】

また、ハードディスクドライブ51に書き込まれたデータのサイズに応じて、データの検査方法を変更することとしてもよい。図13は、書き込まれたデータがシーケンシャルデータであるかどうかに応じて、検査方法を変更する処理を示すフローチャートである。CPU503は、ハードディスクドライブ51にデータを書き込む指示をディスク制御部502に送信する(S1301)。そして、CPU503は当該データが書き込まれた磁気ディスクのヘッドの位置を移動させるシーク処理の実行指示をディスク制御部に送信する(S1302)。CPU503は、当該データがシーケンシャルデータであるかどうか判断する(S1303)。なお、シーケンシャルデータであるかどうかの判断は、書き込まれたデータのサイズが既定のサイズ以上であるかどうかにより行う。

CPU503は、当該データがシーケンシャルデータである場合は、キャッシュメモリ62と磁気ディスクとから当該データの先頭と末尾の1セグメントを読み出す。また、CPU503は、当該データがシーケンシャルデータでない場合は、キャッシュメモリ62と磁気ディスクとから当該データの全部を読み出す(S1306, S1307)。その後、CPU503は読み出した2つのデータが一致しているか比較し(S1308)、一致していない場合は書き込みが正常に行われていないことを情報処理装置300に通知する(S1309)。

このように、書き込んだデータがシーケンシャルデータである場合は、当該データの一部について、磁気ディスクとキャッシュメモリ62とに記憶されているデータを比較することにより、書き込み処理の性能低下を抑えた上で、データの不正を検出することが可能である。また、書き込んだデータがシーケンシャルデータでない場合は、書き込んだデー

10

20

30

40

50



データの全部について、磁気ディスクとキャッシュメモリ 6 2 とに記憶されているデータを比較することにより、シーケンシャルデータの場合ほど書き込み処理の性能を著しく低下させることなく、データの不正を完全に検出することが可能である。

**【 0 0 5 7 】**

ハードディスクドライブ 5 1 は、データの書き込み性能を向上させるため、コントローラ 5 0 0 からデータの書き込み要求を受信すると、当該データをディスクキャッシュのみに書き込み、コントローラ 5 0 0 に書き込み完了を通知する機能を備えている場合がある。この場合、図 1 2 および図 1 3 に説明した方法では、書き込まれたデータの検査を行うことができない。図 1 4 は、ハードディスクドライブ 5 1 がこのような機能を備えている場合において、書き込まれたデータの検査を行う処理のフローチャートを示す図である。CPU 5 0 3 は、ハードディスクドライブ 5 1 への書き込み回数が所定の回数を超過していないか監視している ( S 1 4 0 1 )。所定の回数を超過すると、CPU 5 0 3 はディスクキャッシュに記憶されているデータを磁気ディスクに書き込むよう、ディスク制御部 4 5 0 2 を介してハードディスクドライブに通知する ( S 1 4 0 2 )。そして、CPU 5 0 3 は、当該データをキャッシュメモリ 6 2 と磁気ディスクとから読み出す ( S 1 4 0 3 , S 1 4 0 4 )。CPU 5 0 3 は、キャッシュメモリ 6 2 のデータと磁気ディスクのデータとが一致しているか確認し ( S 1 4 0 5 )、一致していない場合は書き込みが正常に行われていないことを情報処理措置 3 0 0 に通知する ( S 1 4 0 6 )。これにより、前述した書き込み処理の性能を高める機能を用いた上で、データの不正を検出することが可能となる。なお、図 1 4 の処理では、書き込み回数が所定の回数を超過した場合に磁気ディスクへのデータの書き込みと書き込まれたデータの検査を行うこととしたが、所定の時間が経過した場合や、ディスクキャッシュの空領域が無くなった場合などを契機としてもよい。

**【 0 0 5 8 】**

また、シリアル A T A のハードディスクドライブ 5 1 においては、ヘッドの障害により、データの書き込みが正しく行われていないことが多い。そこで、ハードディスクドライブ 5 1 からデータを読み出す際に、ヘッドの障害を検出する方法を説明する。

**【 0 0 5 9 】**

図 1 5 はヘッドチェック管理テーブル 1 5 0 1 を示す図である。ヘッドチェック管理テーブル 1 5 0 1 はドライブ番号とヘッド番号とセクタ番号とで構成され、メモリ 5 0 4 に記憶されている。セクタ番号は更新管理テーブル 1 1 0 1 と同様に L B A により定義されている。CPU 5 0 3 は、ディスク制御部 5 0 2 を介してデータをハードディスクドライブ 5 1 に書き込むと、ヘッドチェック管理テーブル 1 5 0 1 の当該データの書き込みを行ったヘッドの当該セクタの「更新有無」の値を「 1 」に変更する。

**【 0 0 6 0 】**

図 1 6 は、CPU 5 0 3 が実行するヘッドチェック処理のフローチャートを示す図である。CPU 5 0 3 は、検査ヘッド番号に初期値として 1 を設定する ( S 1 6 0 1 )。CPU 5 0 3 は、一定時間経過するのを待ち ( S 1 6 0 2 )、検査ヘッド番号で指定されるヘッドを用いて磁気ディスクの管理領域に検査用のデータを書き込む ( S 1 6 0 3 )。なお管理領域は磁気ディスク上の予め定められている記憶領域である。次に CPU 5 0 3 は、管理領域に書き込まれているデータを読み出し ( S 1 6 0 4 )、読み出したデータと検査用のデータとが一致しているか確認する ( S 1 6 0 5 )。

**【 0 0 6 1 】**

データが一致している場合、CPU 5 0 3 は当該ヘッドに異常が無いと判断し、ヘッドチェック管理テーブル 1 5 0 1 の当該ヘッドの「更新有無」を「 0 」に変更する ( S 1 6 0 6 )。CPU 5 0 3 は、検査ヘッド番号に 1 を加算する ( S 1 6 0 7 )。検査ヘッド番号がヘッド番号の最大値より大きい確認し ( S 1 6 0 8 )、大きい場合は検査ヘッド番号を 1 に設定する。CPU 5 0 3 は、設定されたヘッド番号について、ヘッドチェック処理を繰り返し実行する。

**【 0 0 6 2 】**

管理領域から読み出したデータと検査用のデータとが一致していない場合、CPU 5 0

10

20

30

40

50

3は当該ハードディスクドライブ51に異常が発生していることを情報処理装置300に通知し、処理を終了する。

【0063】

図17は、CPU503が情報処理装置300からデータの読み出し要求を受信した際の処理のフローチャートを示す図である。CPU503は、チャンネル制御部501を介して情報処理装置300からデータの読み出し要求を受信する(S1701)。CPU503は、ヘッドチェック管理テーブル1501から、当該データが記憶されているハードディスクドライブ51の対象セクタの「更新有無」を確認する(S1702, S1703)。「更新有無」が「1」である場合は、当該ハードディスクドライブ51の当該LBAについてデータの書き込みが行われているが、前述したヘッドチェック処理が行われていない状態を示している。「更新有無」が「0」である場合は、CPU503は当該データをハードディスクドライブ51から読み出す(S1708)。

10

【0064】

「更新有無」が「1」である場合、CPU503は前述したヘッドチェック処理と同様に、当該ヘッドを用いて磁気ディスクの管理領域に検査用のデータを書き込む(S1704)。なお管理領域は磁気ディスク上の予め定められている記憶領域である。次にCPU503は、管理領域に書き込まれているデータを読み出し(S1705)、読み出したデータと検査用のデータとが一致しているか確認する(S1706)。

データが一致している場合、CPU503は当該ヘッドに異常が無いと判断し、ヘッドチェック管理テーブル1501の当該ヘッドの「更新有無」を「0」に変更する(S1707)。そして、CPU503は当該読み出し要求に従いハードディスクドライブ51からデータを読み出す(S1708)。

20

管理領域から読み出したデータと検査用のデータとが一致していない場合、CPU503は当該ハードディスクドライブ51に異常が発生していることを情報処理装置300に通知し(S1709)、当該ハードディスクドライブ51からデータを読み出さずに処理を終了する。

【0065】

これにより、ハードディスクドライブ51に書き込まれているデータを読み出す際に、当該データの書き込みを行ったヘッドが正常であるかどうかを確認することができる。ヘッドが異常である場合、データが正しく書き込まれていない可能性や、データの読み出しを正しく行うことができない可能性がある。データの読み出し時にヘッドの異常を検知することにより、不正なデータを読み出すことを防止することが可能となる。

30

【0066】

== パリティ付与による検査 ==

前述したRAID構成でのストライプグループの全てのデータを読み出してパリティチェックする方法では、ストライプグループの中のどのデータが不正な状態となっているか判断することができなかった。そのため、不正なデータの読み出しを防止することは可能であるが、不正なデータを復旧させることができず、データを損失してしまう可能性がある。そこで、ストライプグループにおけるパリティデータとは別に、各データにパリティデータを付与する方法について説明する。

40

【0067】

CPU503は、データをハードディスクドライブ51に書き込む際の最小単位を必ず複数セクタとし、これら複数セクタに対する誤りを検出するためのパリティデータを生成する。本実施の形態において、この複数セクタのデータとパリティデータとの組合せをデータユニットと称することとする。CPU503は、チャンネル制御部501を介して情報処理装置300からデータの書き込み要求を受信すると、当該データからデータユニットを形成する。CPU503は当該データユニットをディスク制御部502を介してハードディスクドライブ51に書き込む。

図18は、ハードディスクドライブに1つのデータ1801が書き込まれている様子を示す図である。データ1801は複数のセクタS#1~S#4で構成され、これら複数セ

50

クタのデータ1801に対するパリティデータ1802とでデータユニット1803が形成されている。CPU503はチャンネル制御部501を介して情報処理装置300からデータの読み出し要求を受信すると、ディスク制御部502を介して当該データのデータユニット1803を読み出し、当該データのパリティチェックを行うことで当該データが不正な状態となっていないか検査する。このように、読み出し要求の対象となっているデータのみを読み出して、当該データが不正な状態となっていないか判断することが可能となる。また、ハードディスクドライブ51がRAID5のように冗長性のあるRAID構成である場合には、ストライプグループにおける他のデータ及びパリティデータとを用いて、当該データを復元することが可能であるため、データを損失することがない。

#### 【0068】

1つのハードディスクドライブ51においてヘッドの障害等が発生している場合は、不正な状態となっているセクタが複数発生する可能性が高い。データユニット1803が1つのハードディスクドライブ51に書き込まれている場合にデータユニット1803のうちの複数のセクタが不正な状態となると、パリティチェックにより不正を検出できない場合がある。

#### 【0069】

そこで、図19に示すように、CPU503は、前述したデータユニット1803をディスク制御部502を介してRAIDグループ内の複数のハードディスクドライブ51に分散して書き込むこととしてもよい。図20は、データユニット管理テーブル2001を示している図である。データユニット管理テーブル2001は、複数セクタで構成されるデータユニット1803がどのハードディスクドライブ51のどのLBAに対応しているかを示している。図20の例では、000~129までの130セクタで1つのデータユニット1803が形成され、このデータユニットはドライブ番号#0とドライブ番号#1のハードディスクドライブ51の000~064までのLBAにより構成されていることが示されている。CPU503は、情報処理装置300からデータの書き込み要求を受信すると、データユニット管理テーブル2001を参照し、データユニット1803ごとに複数のハードディスクドライブ51に分散してデータを書き込む。

#### 【0070】

これにより、1つのハードディスクドライブに障害が発生している場合においても、データの不正を検出することができる可能性が高くなる。

#### 【0071】

＝ファイバチャネルとシリアルATAとが混在する環境＝

次に、ファイバチャネルのハードディスクドライブ51とシリアルATAのハードディスクドライブ51とが混在しているディスクアレイ装置10における説明を行う。

#### 【0072】

図21は、第一の筐体2101にファイバチャネルのハードディスクドライブ51、第二の筐体2102にシリアルATAのハードディスクドライブ51が収容されているディスクアレイ装置を示すブロック図である。なお、第一の筐体2101及び第二の筐体2102とは、基本筐体20または増設筐体30のことである。各ハードディスクドライブ51のディスク制御部502との接続形態については、前述した通りである。また、図19においては、1つのコンバータ901に複数のシリアルATAのハードディスクドライブが接続される形態を示しているが、前述したように、コンバータ801が各ディスクドライブユニットに設けられて接続されているものとしてもよい。

#### 【0073】

このようなディスクアレイ装置10においては、ファイバチャネルのハードディスクドライブ51と比較して信頼性の低いシリアルATAのハードディスクドライブ51の信頼性を高めることが求められている。そこで、コントローラ500は、シリアルATAのハードディスクドライブ51のみに対して、前述した信頼性を高める方法を適用する。これにより、基幹業務等の高いアクセス性能が要求される処理に用いられるファイバチャネルのハードディスクドライブ51に対するデータの読み書き性能を落とさずに、シリアルA

10

20

30

40

50

T Aのハードディスクドライブ5 1に対するデータの読み書きの信頼性を高めることができる。また、シリアルA T Aのハードディスクドライブ5 1の各磁気ディスクにヘッドを2つずつ設ける等の物理的な構造の変更が必要でないため、シリアルA T Aのハードディスクドライブ5 1の製造コストを抑えることが可能である。

【0074】

なお、本実施の形態においては、ファイバチャネルのハードディスクドライブ5 1とシリアルA T Aのハードディスクドライブ5 1とが混在しているとしたが、信頼性の異なるインタフェース規格のハードディスクドライブ5 1であれば他のものでもよい。例えば、シリアルA T Aのハードディスクドライブ5 1の代わりにパラレルA T Aのハードディスクドライブ5 1であるとしてもよい。

10

【0075】

以上、本実施の形態について説明したが、上記実施例は本発明の理解を容易にするためのものであり、本発明を限定して解釈するためのものではない。本発明は、その趣旨を逸脱することなく、変更、改良され得ると共に、本発明にはその等価物も含まれる。

【図面の簡単な説明】

【0076】

【図1】本実施の形態における、ディスクアレイ装置の外観を示す図である。

【図2】本実施の形態における、ディスクアレイ装置の基本筐体の構成を示す図である。

【図3】本実施の形態における、ディスクアレイ装置の増設筐体の構成を示す図である。

【図4】本実施の形態における、ハードディスクドライブの構成を示す図である。

20

【図5】本実施の形態における、ディスクアレイ装置の構成を示す図である。

【図6】本実施の形態における、コントローラのCPUが実行するマイクロプログラムがメモリに記憶されている状態を示す図である。

【図7】本実施の形態における、ファイバチャネルのハードディスクドライブをコントローラのディスク制御部と接続する形態を示す図である。

【図8】本実施の形態における、シリアルA T Aのハードディスクドライブをコントローラのディスク制御部と接続する第一の形態を示す図である。

【図9】本実施の形態における、シリアルA T Aのハードディスクドライブをコントローラのディスク制御部と接続する第二の形態を示す図である。

【図10】本実施の形態における、RAIDグループを構成するハードディスクドライブにデータが書き込まれている例を示す図である。

30

【図11】本実施の形態における、更新管理テーブルを示す図である。

【図12】本実施の形態における、データ書き込み時にキャッシュメモリと磁気ディスクとに記憶されているデータを比較するフローチャートを示す図である。

【図13】本実施の形態における、データ書き込み時にデータサイズを考慮してキャッシュメモリと磁気ディスクとに記憶されているデータを比較するフローチャートを示す図である。

【図14】本実施の形態における、ディスクキャッシュに記憶されているデータを磁気ディスクに書き込む際にキャッシュメモリと磁気ディスクとに記憶されているデータを比較するフローチャートを示す図である。

40

【図15】本実施の形態における、ヘッドチェック管理テーブルを示す図である。

【図16】本実施の形態における、定期的実施するヘッドチェックのフローチャートを示す図である。

【図17】本実施の形態における、データの読み出し時にヘッドチェックを実施するフローチャートを示す図である。

【図18】本実施の形態における、データユニットが1つのハードディスクドライブに書き込まれている例を示す図である。

【図19】本実施の形態における、データユニットが複数のハードディスクドライブに分散して書き込まれている例を示す図である。

【図20】本実施の形態における、データユニット管理テーブルを示す図である。

50

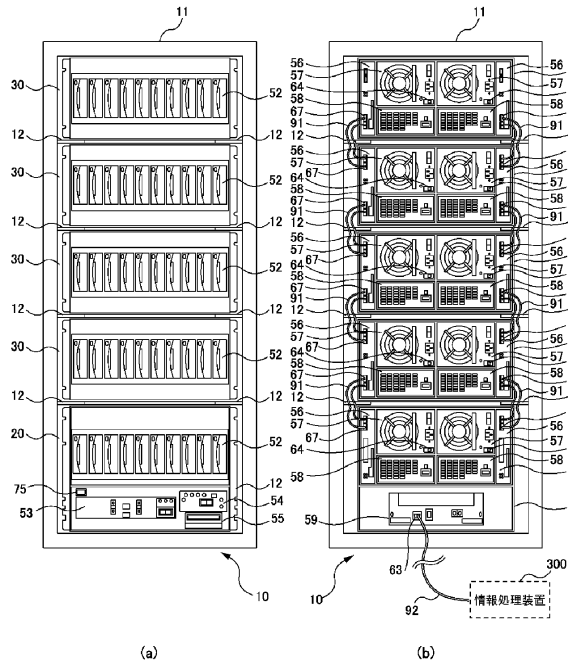
【図 2 1】本実施の形態における、第一の筐体にファイバチャネルのハードディスクドライブが収容され、第二の筐体にシリアル A T A のハードディスクドライブが収容されているディスクアレイ装置の構成を示す図である。

【符号の説明】

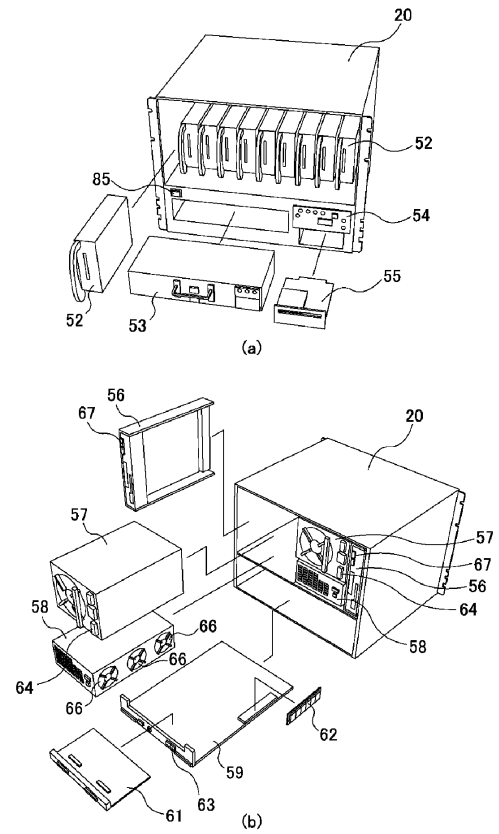
【 0 0 7 7 】

1 0	ディスクアレイ装置	1 1	ラックフレーム	
1 2	マウントフレーム	2 0	基本筐体	
3 0	増設筐体	4 8	制御ライン	
4 9	電源供給ライン	5 1	ハードディスクドライブ	
5 2	ディスクドライブユニット	5 6	電源コントローラボード	10
5 7	A C / D C 電源	5 8	冷却ファンユニット	
5 9	コントローラボード	6 1	通信インターフェースボード	
6 2	キャッシュメモリ	6 4	ブレーカスイッチ	
6 6	冷却ファン	6 7	コネクタ	
7 0	ディスクドライブの筐体	7 3	磁気ディスク	
8 5	メインスイッチ	8 1	電源コントローラ	
9 1	ファイバチャネルケーブル	3 0 0	情報処理装置	
5 0 0	コントローラ	5 0 1	チャネル制御部	
5 0 2	ディスク制御部	5 0 3	C P U	
5 0 4	メモリ	5 0 5	データコントローラ	20
5 0 6	F C - A L	7 0 1	P B C	
8 0 1	コンバータ	9 0 1	コンバータ	
9 0 2	スイッチ	1 0 0 1	R A I D グループ	
1 0 0 2	ストライプグループ	1 1 0 1	更新管理テーブル	
1 5 0 1	ヘッドチェック管理テーブル			
1 8 0 1	複数セクタのデータ	1 8 0 2	パリティデータ	
1 8 0 3	データユニット	2 0 0 1	データユニット管理テーブル	

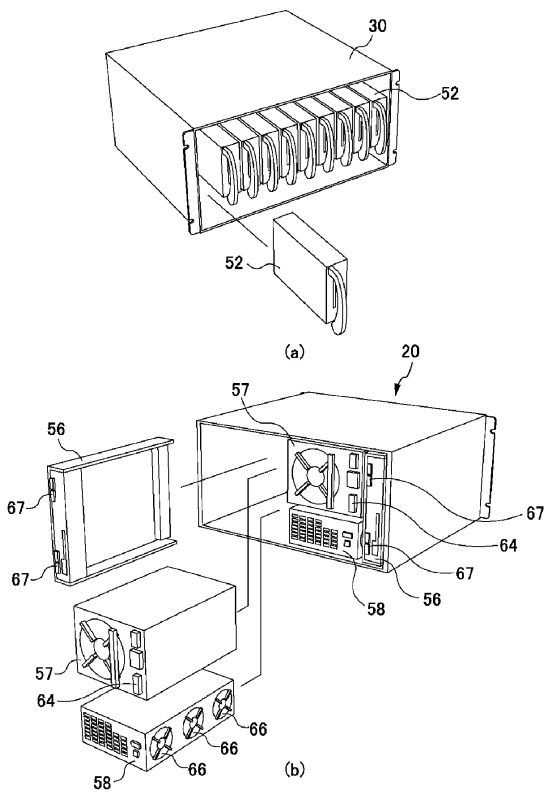
【 図 1 】



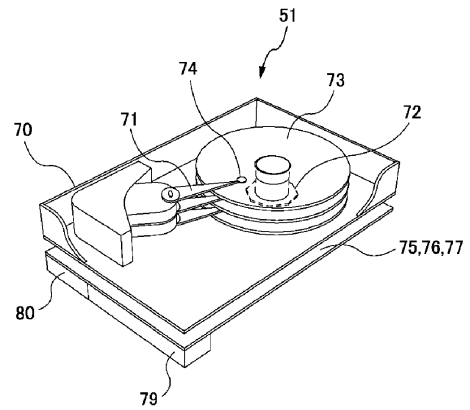
【 図 2 】



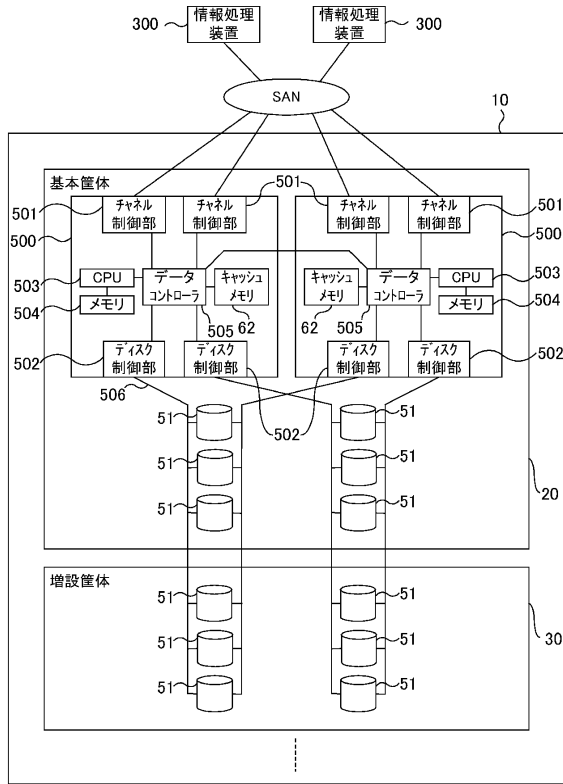
【 図 3 】



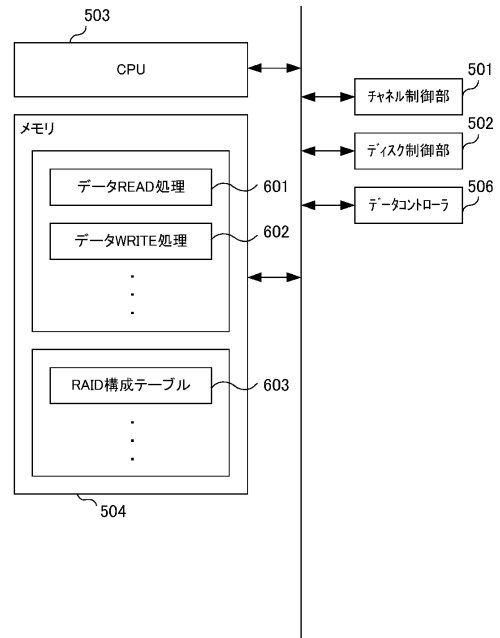
【 図 4 】



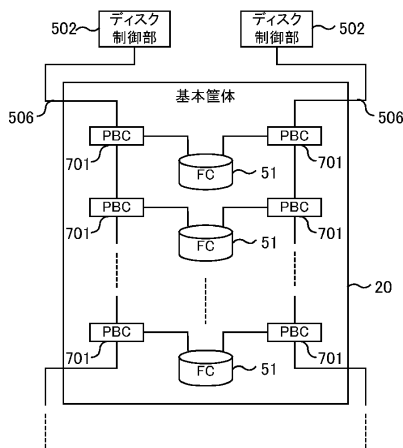
【図5】



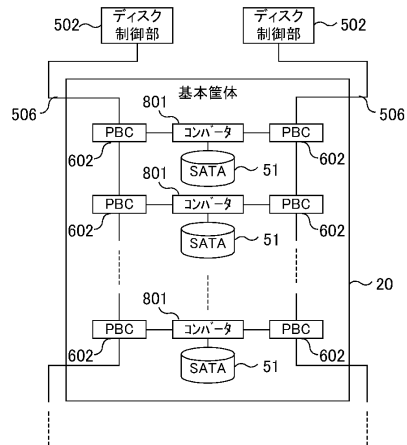
【図6】



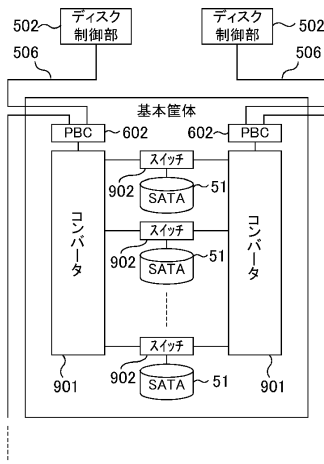
【図7】



【図8】



【図9】

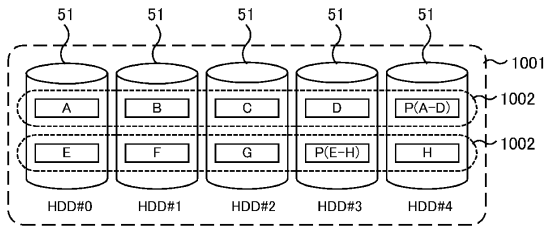


【図11】

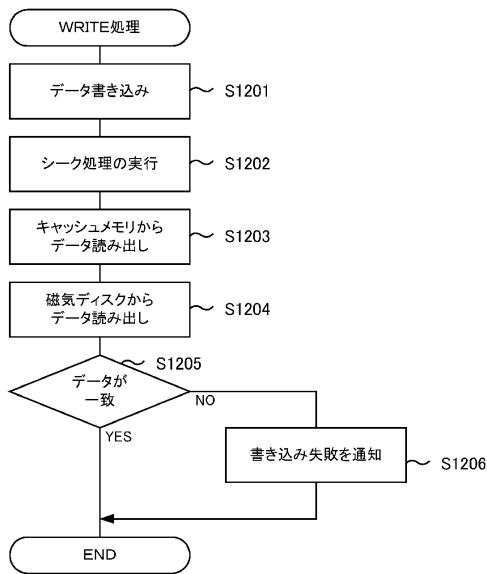
セクタ番号 ドライブ番号	LBA #1-128	LBA #129-256	LBA #257-384	...
HDD#0	0	0	1	...
HDD#1	1	0	0	...
HDD#2	0	1	0	...
...	...	...	...	...

1101

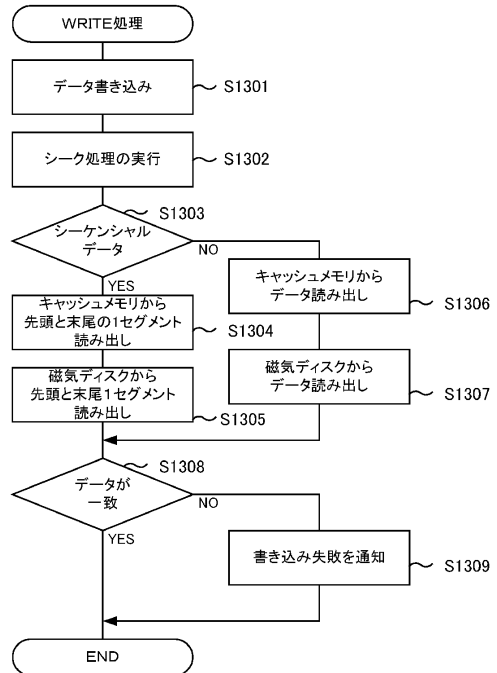
【図10】



【図12】

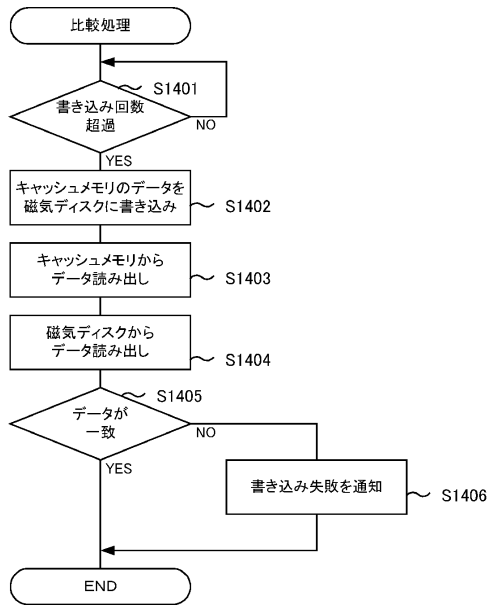


【図13】





【図14】

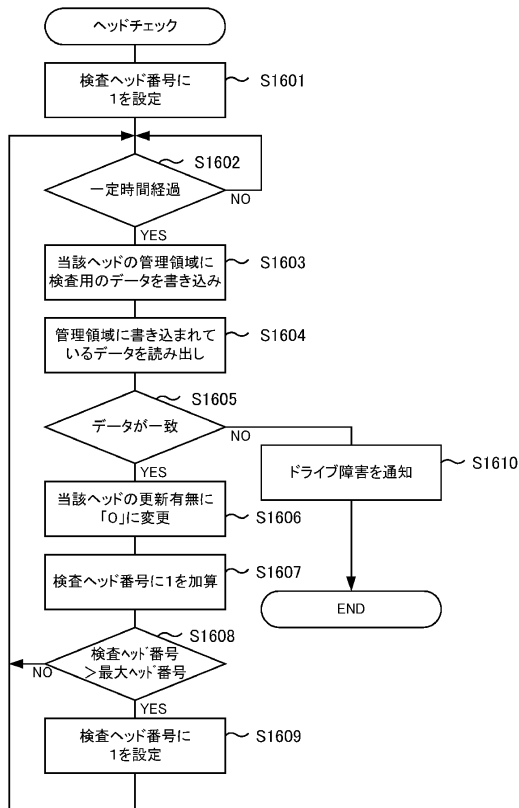


【図15】

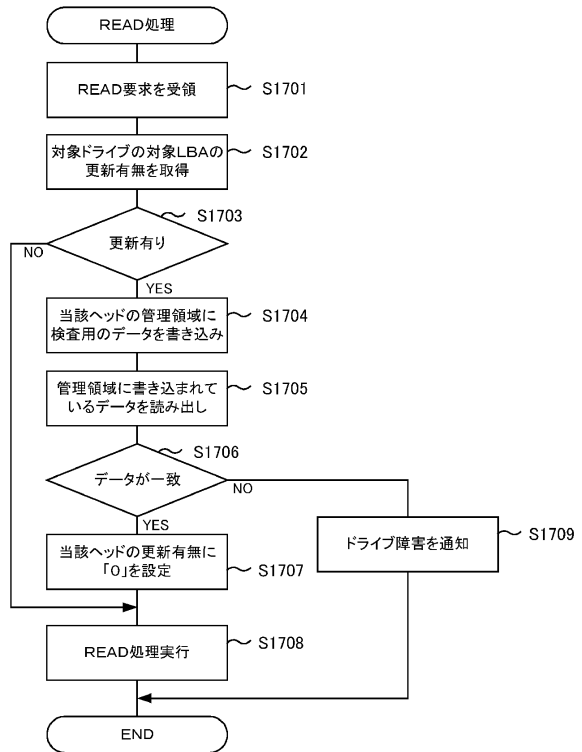
		LBA #1-128	LBA #129-256	LBA #257-384	...
HDD#0	ヘッド番号	#0	#0	#1	...
	更新有無	1	0	0	...
HDD#1	ヘッド番号	#0	#0	#1	...
	更新有無	0	1	0	...
...	...	...	...	...	...

1501

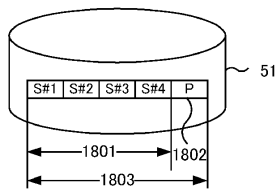
【図16】



【図17】



【図18】

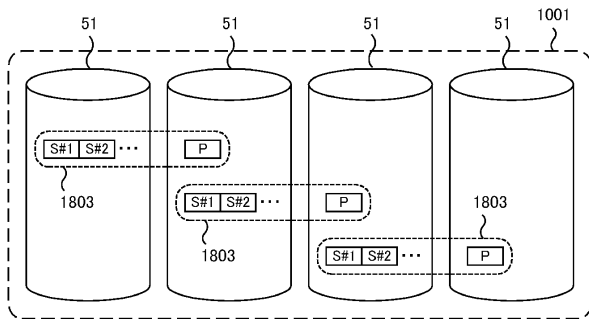


【図20】

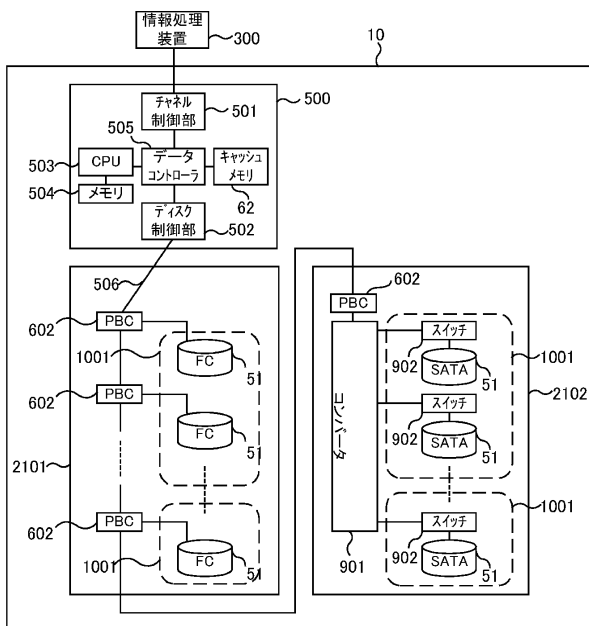
データ ユニット	ドライブ	ドライブLBA
000-129	#0	000-064
	#1	000-064
130-259	#2	000-064
	#0	065-129
⋮	⋮	⋮

2001

【図19】



【図21】



---

フロントページの続き

(72)発明者 八木沢 育哉

神奈川県川崎市麻生区王禅寺1099番地 株式会社日立製作所 システム開発研究所内

審査官 藤井 浩

(56)参考文献 特開2002-007077(JP,A)  
国際公開第03/083636(WO,A1)  
特表2005-520258(JP,A)  
特開2003-108508(JP,A)  
特開2003-108510(JP,A)  
特開2001-222385(JP,A)  
特開2003-303055(JP,A)  
米国特許出願公開第2003/0135577(US,A1)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06 - 3/08