

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4969791号
(P4969791)

(45) 発行日 平成24年7月4日(2012.7.4)

(24) 登録日 平成24年4月13日(2012.4.13)

(51) Int.Cl. F I
G 0 6 F 3 / 0 6 (2006.01) G O 6 F 3 / 0 6 3 O 4 N
 G O 6 F 3 / 0 6 5 4 O

請求項の数 10 (全 21 頁)

(21) 出願番号	特願2005-97505 (P2005-97505)	(73) 特許権者	000005108
(22) 出願日	平成17年3月30日 (2005.3.30)		株式会社日立製作所
(65) 公開番号	特開2006-277487 (P2006-277487A)		東京都千代田区丸の内一丁目6番6号
(43) 公開日	平成18年10月12日 (2006.10.12)	(74) 代理人	110000062
審査請求日	平成19年11月26日 (2007.11.26)		特許業務法人第一国際特許事務所
		(72) 発明者	新井 政弘
			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
		(72) 発明者	松並 直人
			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	菅原 浩二

最終頁に続く

(54) 【発明の名称】 ディスクアレイ装置およびその制御方法

(57) 【特許請求の範囲】

【請求項1】

ディスクアレイコントローラと前記コントローラを制御する制御プログラムを備えるディスクアレイ装置であって、

前記ディスクアレイコントローラは、物理的に分離不可能な複数個のプロセッサコアを内蔵したCPUを搭載しており、

前記制御プログラムは、各々処理完結動作するプログラムモジュールと、前記プロセッサコアを管理する情報を有し、該管理情報を用いて、該プロセッサの状態を管理する管理部とを含み、

前記各プロセッサコアはそれぞれ独立した単位プロセッサであり、

前記単位プロセッサには、前記各々完結したプログラムモジュールが動的または固定、半固定の何れかになるように割り当てられ、

該割り当てられたプログラムモジュールの1つが、前記管理部の管理情報を用いて自己および他の単位プロセッサの動作負荷または動作状態を監視し、該単位プロセッサの負荷分散が必要であると判断したとき、前記プログラムモジュール群動作中に該負荷分散が必要である単位プロセッサのプログラムモジュールを未割当単位プロセッサに動的に割り当て、負荷分散処理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、

前記割り当てられたプログラムモジュールの1つが、前記単位プロセッサの動作負荷ま

10

20

たは動作状態を監視する監視制御プログラムモジュールであり、該監視制御プログラムモジュールは、固定または半固定になるように割り当てた単位プロセッサの動作負荷または動作状態が、前記管理部の管理情報に基づいて負荷分散開始を決定する規定値に達すると、該単位プロセッサへのプログラムモジュール割り当てを維持しつつ、該単位プロセッサとは別の単位プロセッサに前記プログラムモジュールを動的となるよう追加して割り当てることを特徴とするディスクアレイ装置。

【請求項3】

請求項1に記載のディスクアレイ装置において、

前記単位プロセッサは、固定または半固定になるように前記プログラムモジュールを割り当てた単位プロセッサの動作負荷が相対的に高いとき、動的になるように割り当てたプログラムモジュールの動作負荷が相対的に低い単位プロセッサの処理を終了させ、終了させた単位プロセッサに、前記相対的に動作負荷が高い単位プロセッサに割り当てたプログラムモジュールを、動的になるように割り当てることを特徴とするディスクアレイ装置。

10

【請求項4】

請求項3に記載のディスクアレイ装置において、

前記単位プロセッサは、前記プログラムモジュールの動作負荷が相対的に高い単位プロセッサがグループ化され、該グループ化された同一のグループ内に属する単位プロセッサで、かつ、単位プロセッサが動的になるように割り当てたプログラムモジュールの動作負荷が相対的に低い単位プロセッサの処理を終了させ、該終了した単位プロセッサに、前記相対的に動作負荷が高い単位プロセッサに割り当てたプログラムモジュールを、動的になるように割り当てることを特徴とするディスクアレイ装置。

20

【請求項5】

請求項3に記載のディスクアレイ装置において、

前記単位プロセッサは、どのグループにも属さない単位プロセッサに、前記相対的に動作負荷が高い単位プロセッサに割り当てた制御プログラムモジュールを、動的になるよう割り当てるとともに、割り当てた単位プロセッサを、前記相対的に動作負荷が高い単位プロセッサが属するグループに属するようにすることを特徴とするディスクアレイ装置。

【請求項6】

ディスクアレイコントローラと前記コントローラを制御する制御プログラムを備えるディスクアレイ装置であって、

30

前記ディスクアレイコントローラは、物理的に分離不可能な複数個のプロセッサコアを内蔵したCPUを搭載しており、各プロセッサコアはそれぞれ独立した単位プロセッサとなり、前記制御プログラムは、各々処理完結動作するプログラムモジュール群と、前記プロセッサコアを管理する情報を有し、該管理情報を用いて、該プロセッサの状態を管理する管理部を含み、

前記単位プロセッサには、前記各々処理完結したプログラムモジュールが動的または固定、半固定の何れかになるように割り当てられ、該制御プログラムの1つである監視制御プログラムモジュールは、前記管理部の管理情報を用いて自己および他の単位プロセッサの動作負荷または動作状態を監視し、特定のプロセッサコアに障害が発生した場合に、前記プログラムモジュール群間での制御、非制御の関係を表す特権レベルの低い制御プログラムモジュールから順にすべて終了させ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6に記載のディスクアレイ装置において、

前記単位プロセッサは、前記障害が発生したプロセッサコアの単位プロセッサに割り当てたプログラムモジュールの割り当て形式が半固定であり、かつ、前記プログラムモジュールの割り当て形式を動的とした別の単位プロセッサがあるとき、該単位プロセッサの割り当て形式を動的から半固定に変更することを特徴とするディスクアレイ装置。

【請求項8】

請求項6に記載のディスクアレイ装置において、

前記単位プロセッサは、前記障害が発生したプロセッサコアの単位プロセッサに前記プ

50

プログラムモジュールの割り当て形式が固定であるとき、他系のディスクアレイコントローラに処理引き継ぎを依頼し、他のプログラムモジュールの終了処理を省略する障害処理を行うことを特徴とするディスクアレイ装置。

【請求項 9】

請求項 1 に記載のディスクアレイ装置において、

前記ディスクアレイコントローラは、物理的に分離不可能な複数個のプロセッサコアを内蔵した CPU を搭載しており、内部スイッチ、不揮発性メモリコントローラ、揮発性メモリコントローラ、ディスクアレイコントローラ間転送コントローラ、パリティ演算器および CPU 内部キャッシュを有し、前記制御プログラムは、監視制御プログラムモジュール、運用管理プログラムモジュール、RAID 制御プログラムモジュール、NAS 制御プログラムモジュール、ホスト I/O 制御プログラムモジュール、ドライブ I/O 制御プログラムモジュール、初期 MPU コア割当管理テーブル、初期閾値管理テーブル、RAID 設定管理テーブル、LU 設定管理テーブル、および NAS ボリューム管理テーブルを備えており、前記単位プロセッサは、自己および他の単位プロセッサの制御プログラムの負荷状態を監視し、前記管理部の管理情報である負荷閾値および閾値オーバ回数が規定回数を超えると、同一グループ内の未割当コアを検索し、検出した未割当コアに、前記プログラムモジュールを割り当てて、

10

【請求項 10】

物理的に分離不可能な複数個のプロセッサコアを内蔵した CPU を搭載し、各プロセッサコアはそれぞれ単位プロセッサとなるディスクアレイコントローラと、前記コントローラを制御し、各々処理完結動作するプログラムモジュールと、前記プロセッサコアを管理する情報を有し、該管理情報を用いて、該プロセッサコアの状態を管理する管理部を含む制御プログラムを備えるディスクアレイ装置の制御方法であって、

20

前記単位プロセッサに、前記各々完結したプログラムモジュールが動的または固定、半固定の何れかになるように割り当て、

前記割り当てた前記プログラムモジュールの 1 つが、前記管理部の管理情報を用いて自己および他の単位プロセッサの動作負荷または動作状態を監視し、

該単位プロセッサの負荷分散が必要であると判断したとき、前記プログラムモジュール群動作中に該負荷分散が必要である単位プロセッサのプログラムモジュールを未割当単位プロセッサに動的に割り当て、負荷分散処理することを特徴とするディスクアレイ装置の制御方法。

30

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、複数のディスク装置を備えるディスクアレイ装置、およびディスクアレイ装置の制御方法に関する。

【背景技術】

【0002】

複数のディスク装置を備えることによってデータに対するアクセス速度、ディスク装置に対する信頼性を向上させる技術として、RAID (Redundant Array of Independent Disks) が知られている。一般にディスクアレイ装置では、この RAID 制御を装置内に備えるディスクアレイコントローラが担っており、従来より処理性能向上のためディスクアレイコントローラの高性能化が図られてきた。

40

【0003】

ディスクアレイコントローラにはアレイコントローラの処理プログラムを動作させるための CPU が搭載されている。高性能化の一手法として、CPU を複数乗せることによって処理を並列化し、高速化する手法が考えられる。

【特許文献 1】特開平 9 - 160889 号公報

【特許文献 2】特開平 6 - 35871 号公報

50

【発明の開示】**【発明が解決しようとする課題】****【0004】**

しかしながら、複数CPUを搭載したディスクアレイ装置では、CPU間での情報共有のために高価な共有メモリが必要となる。このため性能の向上分に対する相対的なコストが高くなってしまいコストパフォーマンスが悪化してしまうという問題がある。また、CPU個別にクロックや電源供給回路が必要となるため、実装面でも肥大化し、結果として製造コストも増加してしまうといった問題がある。

【0005】

一方、制御面で見ると、特許文献1では処理をジョブ単位で制御プロセッサないしプロセッサ群に割り当てる方法が取られているが、この場合、複数の分割された制御（ジョブ）が1つのプロセッサないしプロセッサ群内に混在して存在するため、例えばディスクアレイ装置の運用管理機能がセキュリティクラック等の攻撃を受けてダウンしてしまうような場合、それにひきずられて他の制御もダウンしてしまう可能性がある。

【0006】

本発明は、上記に示した課題を解決するためになされたものであり、高コストパフォーマンスの実現、他制御への影響抑止、リソースの有効活用を可能とするディスクアレイ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】**

上記課題を解決するために、本発明の実施形態では、複数プロセッサコア(MPUコア)を内蔵したCPUを搭載するディスクアレイコントローラを備え、処理完結しているプログラムモジュールを各コアに静的または動的に割り当てて処理動作するディスクアレイ装置を提供する。ここで、処理完結しているプログラムモジュールとは、同一プログラム内で起動・停止・障害復旧処理及び外部への通知のための論理インターフェースを備え、他のプログラムと連携非同期で動作可能なプログラムを指す。

【0008】

本発明の実施形態に係るディスクアレイ装置は、MPUコアを固定資源または動的資源として管理する手段と、各コアの負荷状況を監視し、必要に応じて外部からの指示なしに各制御へのコアの割当て数を変動させる監視制御手段と、コアに障害が発生した際に動作中プログラム間の依存関係を考慮して終了させ、必要に応じて他系のディスクコントローラにフェイルオーバーする障害処理手段とを備えることを特徴とする。

【0009】

すなわち、本発明は、ディスクアレイコントローラを備えるディスクアレイ装置であって、前記ディスクアレイコントローラは、物理的に分離不可能な複数個のプロセッサコアを内蔵したCPUを搭載しており、各プロセッサコアはそれぞれ単位プロセッサとなり、一の単位プロセッサは、自己及び他の単位プロセッサを個別に管理し、各々完結した制御プログラムを、自己又は他の単位プロセッサにプロセッサ単位で割り当てる際に、適宜動作終了可能とする動的になるように、又は前記CPU全体の動作終了まで動作可能とする非動的になるように割り当て、割り当てた制御プログラムの動作負荷又は動作状態を単位プロセッサごとに管理するディスクアレイ装置である。

【発明を実施するための最良の形態】**【0010】**

本発明を実施するための最良の形態を説明する。

以下、本発明に係るディスクアレイ装置およびディスクアレイ装置の制御方法について図面を参照しつつ、実施例に基づいて説明する。

【0011】

実施例1を説明する。本実施例のディスクアレイ装置の構成について、図1～図8を参照して説明する。図1は本実施例に係るディスクアレイ装置の概略構成を示す説明図である。図2は本実施例に係るディスクアレイ装置の概観図である。図3は本実施例に係るデ

10

20

30

40

50

ディスクアレイ装置におけるディスクアレイコントローラの機能的構成を示すブロック図であり、図4はその概観図を表している。図5は本実施例に係るディスクアレイ装置においてディスクアレイコントローラに搭載される複数プロセッサコア(MPUコア)を内蔵したCPU(以下、「マルチコアCPU」と記す)の機能的構成を示すブロック図である。

【0012】

また、図6は制御プログラムに含まれる各プログラムモジュール及び管理テーブルを示す説明図であり、図7、8は前記管理テーブルのうちの2つを具体的に示した一例である。以下、順に説明を行う。

【0013】

本実施例におけるディスクアレイ装置1は、ディスクアレイコントローラ11、12、
接続インターフェース130、131、132、電源105、106、および、複数のディスク装置D00~D2Nを備えている。複数のディスク装置D00~D2Nは、例えば、図2に示すようにしてディスク装置1に備えられると共に、RAIDシステムを構成している。

10

【0014】

ディスクアレイコントローラ11、12は、制御プログラムを実行することによって、ディスクアレイ装置1における各種制御処理を実行する制御回路である。本実施例においては、2つのディスクアレイコントローラ11、12が備えられているが、1つまたは3つ以上のディスクアレイコントローラが備えられていてもよい。ディスクアレイコントローラ11、12は信号線101を介して相互に通信可能に接続されている。ディスクアレイコントローラ11、12はまた、ストレージネットワーク40を介して各ホスト20、21、22と接続され、管理用ネットワーク30とを介して管理用端末装置31と接続されている。ストレージネットワークは、例えば、ファイバチャネルによるFC-SAN(Storage Area Network)やTCP/IPネットワークを利用したIP-SANなどであり、管理用ネットワークはTCP/IPネットワークを利用したLANや、シリアルケーブルによるPoint to Pointネットワークである。

20

【0015】

ディスクアレイコントローラ11、12は、接続インターフェース130、131、132を介して複数のディスク装置D00~D2Nと接続されている。より具体的には、接続インターフェース130は、ディスクアレイコントローラ11、12と信号線102を介して直接接続されており、定期的な通信を行っている。また、各接続インターフェース130、131、132は互いに信号線103を介して接続されている。従って、接続インターフェース131は接続インターフェース130を介して、接続インターフェース132は接続インターフェース130、131を介してディスクアレイコントローラ11、12と接続されている。

30

【0016】

接続インターフェース130は、複数のディスク装置D00~D0Nと接続され、接続インターフェース131は複数のディスク装置D10~D1Nと接続され、接続インターフェース132は複数のディスク装置D20~D2Nと接続されている。

【0017】

ディスクアレイコントローラ11、12を含む接続インターフェース130及び複数のディスク装置D00~D0Nのグループは、例えば、基本筐体と呼ばれ、接続インターフェース131及び複数のディスク装置D10~D1Nのグループ、及び接続インターフェース132及び複数のディスク装置D20~D2Nのグループは、例えば、増設筐体と呼ばれる。なお、図1からも明らかなように、増設筐体は0ないし1つであってもよく、あるいは、3つ以上であってもよい。

40

【0018】

なお、本実施例では基本筐体をディスクアレイコントローラ11、12および接続インターフェース130、複数のディスク装置D00~D0Nから成るグループとして記載しているが、基本筐体に複数のディスク装置D00~D0Nを含まない形態でも良い。

50

【 0 0 1 9 】

ホスト 2 0、2 1、2 2 は、例えば、各種データを入力する端末装置であり、ホスト 2 0、2 1、2 2 において処理されたデータは、逐次、ディスクアレイ装置 1 に対し送られ、ディスクアレイ装置 1 に格納される。なお、ホスト 2 0 ~ 2 1 は、1 つであっても良く、あるいは、4 つ以上備えられてもよい。

【 0 0 2 0 】

電源 1 0 5 は、電力線 1 0 7 及び接続インターフェース 1 3 0 を介して複数のディスク D 0 0 ~ D 0 N に動作のための電力を供給し、また、電力線 1 0 7 及び接続インターフェース 1 3 0、及び電力線 1 0 4 を介してディスクアレイコントローラ 1 1、1 2 の動作のための電力を供給する。同様に、電源 1 0 6 は、電力線 1 0 7 及び接続インターフェース 1 3 1、1 3 2 を介し、複数のディスク装置 D 1 0 ~ 1 N、D 2 0 ~ D 2 N に動作のための電力を供給する。

10

【 0 0 2 1 】

各ディスク装置 D 0 0 ~ D 2 N は、ハードディスクドライブであり、例えば、A T A 規格のハードディスクドライブや S A S 規格のハードディスクドライブが用いられる。

【 0 0 2 2 】

管理用端末 3 1 は、ディスクアレイ装置 1 に対する保守管理を実行するために用いられる端末装置である。管理用端末装置 3 1 には、管理画面 3 2 が備えられており、管理者は管理画面 3 2 を通じて、ディスクアレイ装置 1 の状態を管理する。

20

【 0 0 2 3 】

図 3 を参照して、ディスクアレイコントローラ 1 1 の内部構成について説明する。なお、ディスクアレイコントローラ 1 2 に関しても同様の内部構成を有している。ディスクアレイコントローラ 1 1 は、マルチコア C P U 1 1 0、不揮発性メモリ 1 1 1、揮発メモリ 1 1 2、ホスト側物理ポート 1 1 3、ドライブ側物理ポート 1 1 4、管理ネットワーク用物理ポート 1 1 5、ブート R O M 1 1 6 を備えている。

【 0 0 2 4 】

不揮発性メモリ 1 1 1 には、ディスクアレイ装置 1 を制御するための制御プログラムが格納されている。不揮発性メモリ 1 1 1 は、電力の供給が停止しても継続的にデータを記憶しておくことができ、例えば、F l a s h メモリが用いられる。制御プログラム 1 1 9 については、図 6 を用いて後述する。

30

【 0 0 2 5 】

揮発性メモリ 1 1 2 は、ディスク装置 D から読み出したデータ、ディスク装置 D に書き込むデータ、及び、マルチコア C P U 1 1 0 による演算結果を一時的に格納するためのデータバッファ領域 1 1 2 と、マルチコア 1 1 0 によって制御プログラム 1 1 9 を実行するために、制御プログラム 1 1 9 を読み出して格納しておく制御プログラム配置領域 1 1 7 とを備える。揮発性メモリは停電などによりメモリへの電力供給が停止するとデータを記憶しておくことができなくなるメモリであり、例えば D R A M (D y n a m i c R a n d o m A c c e s s M e m o r y) などが用いられる。

【 0 0 2 6 】

ホスト側物理ポート 1 1 3 は、ストレージネットワーク 4 0 へ物理的に接続し、ホスト 2 0 ~ 2 2 と電気的信号を送受信するための伝送路の受け口である。

40

【 0 0 2 7 】

ドライブ側物理ポート 1 1 4 は、接続インターフェースと電気的信号を送受信するための伝送路の受け口である。

【 0 0 2 8 】

管理ネットワーク用物理ポート 1 1 5 は、管理用ネットワークへ接続し、管理用端末装置 3 1 と電気的信号を送受信するための受け口である。

【 0 0 2 9 】

ブート R O M 1 1 6 は、ディスクアレイ装置 1 を起動した際、不揮発性メモリ 1 1 1 に格納された制御プログラム 1 1 9 を、揮発性メモリ 1 1 2 の制御プログラム配置領域 1 1

50

7に読み出すためのイニシャルプログラムローダを格納した読み出し専用メモリ (Read Only Memory) である。

【0030】

マルチコアCPU110は、複数のプロセッサコア(MPUコア)を内蔵する演算処理装置であり、信号線120を介して、不揮発性メモリ111、揮発性メモリ112、ホスト側物理ポート113、ドライブ側物理ポート114、管理ネットワーク用物理ポート115、ブートROM116と相互に接続されている。マルチコアCPU110は、信号線116を介して、不揮発性メモリ111からのデータの読み出しと書き込み、揮発性メモリ112からのデータの読み出しと書き込みを実行するほか、信号線120とホスト側物理ポート113を介してホスト20~22との間でコマンドおよびデータの送受信を実行し、信号線120とドライブ側物理ポート114、及び、図1に示した信号線102、103、接続インターフェース130~132を介して、ディスク装置Dとの間でコマンドおよびデータの送受信を実行する。

10

【0031】

また、マルチコアCPU110は、信号線101を介して他系コントローラと相互に接続されているおり、他系コントローラとの間でデータおよびコマンドの送受信を実行する。図3はディスクアレイコントローラ11を示しているので、他系コントローラとは、具体的には、ディスクアレイコントローラ12を指す。なお、ディスクアレイコントローラ12に対する他系コントローラとはディスクアレイコントローラ11を意味する。

20

【0032】

図4は、ディスクアレイコントローラ11の概観図を示したものである。なお、ディスクアレイコントローラ12の概観図も同様である。

【0033】

ディスクアレイコントローラ11には、マルチコアCPU110をはじめとし、図3の内部構成で示した部品が基板回路上に配置され、接続されているほか、接続コネクタ122、ブラケット121、ホスト接続用ポート123、管理ネットワーク接続用ポート124、エラー表示LED125が備えられている。

【0034】

また、図4では記載を省略するが、マルチコアCPU110や周辺回路上の半導体の上には、発熱による破壊を防ぐ為の放熱板やファンが取り付けられることもある。

30

【0035】

接続コネクタ122は接続インターフェース130上に設けられたコネクタと嵌合することによって、信号線101、102、及び電力線104と物理的な接続を行う。接続コネクタ上にはこのほかにディスクアレイ装置上で必要となるいくつかの信号線を備えている。たとえば、ディスクアレイコントローラ挿抜検出用信号線などである。

【0036】

ホスト接続用ポート123はストレージネットワーク40へ接続するためのケーブルのコネクタを挿入する接続口であり、管理ネットワーク接続用ポート124は、管理用ネットワーク30へ接続するためのケーブルのコネクタを挿入する接続口である。

40

【0037】

エラー表示LED125はディスクアレイコントローラ11に障害が発生し保守交換が必要となった際に、エラーが発生していることを視覚的にディスクアレイ装置1の外部に通知するための表示灯である。

【0038】

図5を参照して、ディスクアレイコントローラ11に搭載されているマルチコアCPU110の機能的構成の一例について説明する。マルチコアCPU110は、複数のMPUCOA(プロセッサコア)1110~111Nと、内部スイッチ1120、不揮発メモリコントローラ1130、ディスクアレイコントローラ間転送コントローラ1140、パリティ演算器1150、CPU内部キャッシュ1160、揮発メモリコントローラ1170を

50

備える。複数のMPUコア1110～111Nはそれぞれした独立したプロセッサであるが、1つの半導体素子上に作成されており物理的に個別に切り離して利用することはできない点が、従来のマルチプロセッサとは異なる。

【0039】

内部スイッチ1120は複数のMPUコア1110～111N、不揮発メモリコントローラ1130、ディスクアレイコントローラ間転送コントローラ1140、パリティ演算器1150、CPU内部キャッシュ1160、揮発メモリコントローラをスイッチ機構によって高速に相互接続している。

【0040】

不揮発性メモリコントローラ1130はMPUコア1110～111Nの指示に基づいて不揮発性メモリ111との間でデータの転送を実行するI/Oコントローラである。

10

【0041】

ディスクアレイコントローラ間転送コントローラ1140は、他系のディスクアレイコントローラに搭載されたマルチコアCPU110とコマンドとデータの送受信を行う。

【0042】

パリティ演算器1150は、各MPUコア1110～111Nの指示に基づいて与えられたデータのパリティ生成や整合性の確認を高速に行うのに用いられる。たとえばXOR演算器が含まれる。

【0043】

CPU内部キャッシュ1160は、複数のMPUコア1110～111NからマルチコアCPU110外部にある揮発性メモリ112よりも高速にアクセスできる揮発性のメモリである。CPU内部キャッシュ1160は複数のMPUコア1110～111Nが演算結果を一時的に格納するほか、MPUコアの動作状況を格納するMPUコア管理テーブルなどを格納するのに用いられる。

20

【0044】

揮発性メモリコントローラ1170はMPUコア1110～111Nの指示に基づいて揮発性メモリ112との間でデータの転送を実行するI/Oコントローラである。

【0045】

以上の様に、本発明のマルチコアCPU110は複数のMPUコアのほかにI/Oコントローラなどの周辺制御回路を内蔵しているが、少なくとも複数のMPUコアを内蔵していれば他の構成を取っていても良い。例えば不揮発メモリコントローラは内蔵されずに外部回路として設けられる場合や、マルチコアCPUにTCP/IPコントローラが内蔵されるような場合である。

30

【0046】

図6を参照して、制御プログラム119の詳細について説明する。制御プログラム119は、監視制御プログラムモジュールPr1、運用管理プログラムモジュールPr2、RAID制御プログラムモジュールPr3、NAS制御プログラムモジュールPr4、ホストI/O制御プログラムモジュールPr5、ドライブI/O制御プログラムモジュールPr6、初期MPUコア割当管理テーブルTb1、初期閾値管理テーブルTb2、RAID設定管理テーブルTb3、LU設定管理テーブルTb4、NASボリューム管理テーブルTb5を備えている。

40

【0047】

監視制御プログラムモジュールPr1は、プログラムモジュールPr2～Pr6の動作管理及びMPUコアの資源管理を行う、特権レベルの最も高いプログラムモジュールであり、他のプログラムモジュールの処理負荷状況に応じて自律的にMPUコアの割当て数の変更を実行する。

【0048】

運用管理プログラムモジュールPr2は、監視制御プログラムモジュールPr1に次いで特権の高いプログラムモジュールであり、管理用端末装置31からのディスクアレイ装置1に対する運用・保守管理に関する設定を受け、他のプログラムモジュールを通じて

50

実際に設定を実行するプログラムモジュールである。運用・保守管理には、例えば、RAIDグループの設定、LU(Logical Unit)の設定、NAS(Network Attached Storage)の設定、ホストマッピングの設定、監視機能の設定、制御プログラムのアップデート、稼働状況の確認、動作ログの閲覧・取得、ディスク表面検査の設定、ディスクアレイ装置の起動・停止などがある。

【0049】

RAID制御プログラムモジュールPr3は、運用管理プログラムモジュールPr2を介して、管理用端末装置31から管理用端末画面32を通じた指定内容に基づき、RAIDグループの作成やLUの作成、LUの初期化、ホスト20~22に対するLUとLU番号(LUN)との対応付けなどを行うほか、ホスト20~22から受領したコマンドの解釈を行い、必要であれば演算を行いながら、適切なディスク装置Dに対しコマンドを発行してデータの読み書きを実行し、ホスト20~22との間でコマンド処理の結果通知やデータの送受信を実行する。RAID制御プログラムはPr3は、また、次に説明するNAS制御プログラムとの間でも同様のコマンド送受信、演算、ディスク装置Dへのアクセスを実行する。

10

【0050】

NAS制御プログラムモジュールPr4は、ホスト20~22からのファイルレベルアクセス処理を実行するプログラムモジュールであり、その他にRAID制御プログラムモジュールPr3によって作成されたLU上へのファイルシステムの構築や構築されたファイルシステムサイズの変更、ファイルやディレクトリへのアクセス権の設定、アクセス認証方法の設定などを行う。

20

【0051】

ホストI/O制御プログラムモジュールPr5は、RAID制御プログラムモジュールPr3が作成した転送リストの指示に基づき、ホスト20~22と揮発性メモリ112にあるデータバッファ領域118との間で、パリティ演算器1150を介してデータ整合性のチェックを行いながら、データ転送を実行するプログラムモジュールである。ドライブI/O制御プログラムモジュールPr6は、同様にRAID制御プログラムモジュールPr4が指定する転送リストに基づき、ディスク装置Dと揮発性メモリ112にあるデータバッファ領域118との間で、パリティ演算器1150を介してデータ整合性のチェックを行いながら、データ転送を実行するプログラムモジュールである。

30

【0052】

初期MPUCOA割当管理テーブルTb1は、ディスクアレイ装置1が起動した際の、MPUCOA資源管理内容と、各制御プログラムモジュールPr1~Pr6の初期割当を指定した管理表である。初期MPUCOA割当て管理テーブルTb1の詳細は、図7を用いて後述する。

【0053】

初期閾値管理テーブルTb2は、各プログラムモジュールPr1~Pr6の負荷の上限閾値の初期値を管理しているテーブルである。初期閾値管理テーブルTb2の詳細は、図8を用いて後述する。

【0054】

RAID設定管理テーブルTb3は、RAIDグループを構成するディスク装置についての各種情報を管理するために用いられるテーブルであり、例えば、RAIDグループNo、RAIDグループの総記憶容量、RAIDレベル、RAIDグループを構成するディスク装置、正常・異常の状態などの情報を管理する。

40

【0055】

LU設定管理テーブルTb4は、RAIDグループ上に作成される論理ユニット(Logical Unit)を管理するためのテーブルであり、例えば、LUN、所属するRAIDグループNo、論理ユニットに設定されている記憶容量、正常・異常の状態などの情報を管理する。

【0056】

50

N A S ボリューム設定管理テーブル T b 5 は、N A S が作成するファイルシステムを管理するテーブルであり、例えば、ファイルシステム名、ファイルシステムを構成する L U の番号、ファイルシステムのフォーマット、ファイルシステムの総容量、ファイルシステムの使用容量、ファイルシステムの差管理の有無、差管理容量、差管理容量の使用率、ファイルシステムの状態などの情報を管理する。

【 0 0 5 7 】

図 7 に初期 M P U コア割当管理テーブル T b 1 の一例を示す。図 7 に示す初期 M P U コア割当管理テーブル T b 1 は、管理上振付けた M P U コア番号、同 M P U を初期の段階で使用するプログラムモジュール名、モジュールの割当形式、特権レベル、グループ N o についての情報を保持している。

10

【 0 0 5 8 】

モジュールの割当形式は、固定・半固定・動的の 3 種類がある。固定は初期 M P U コア割当管理テーブルで指定した M P U コア以外を割当てることができないことを示し、半固定は初期 M P U コア割当管理テーブルで指定した M P U コアに加え、他の M P U コアを負荷分散用に割当可能であることを示す。動的は具体的なプログラムモジュールが初期の段階では割当てられず、未使用で開始することを示している。図 7 を例にとると、監視制御プログラムモジュールと運用管理プログラムモジュールはそれぞれ M P U コア 0、M P U コア 1 に固定的に割当てられ変更不可であり、R A I D 制御は当初 M P U コア 2 が割当てられるが、必要に応じ他のコアも割当可能であることを示している。

【 0 0 5 9 】

20

特権レベルは、プログラムモジュール間での制御、非制御の関係を表す整数値であり、小さいほど特権性が高いことを示す。図 7 を例にとると、監視制御プログラムモジュールは最も高い特権レベル 0 であり、他のいずれのプログラムモジュールからも制御指示を受けず、逆にいずれのプログラムモジュールにも制御指示が可能であることが分かる。運用管理モジュールの特権レベルは 1 であり、これより高い監視制御プログラムモジュールからの制御指示は受けるが、それ以外からは制御指示を受けない。逆に、例えば、特権レベルが 1 つ低い、特権レベル 2 の R A I D 制御プログラムモジュールに対しては制御指示が可能である。なお、特権レベルが同じとは、お互いに独立している場合などであり、制御指示をお互いに与えることはない。

【 0 0 6 0 】

30

グループ N o は、負荷分散で融通できる M P U コアのグループを示している。図 7 を例に取れば、M P U コア N o 6 及び 7 の M P U コアは、グループ N o が 0 0 3 なので、同じグループ N o を持つ M P U コア N o 2、3 に割当てられているプログラムモジュールの負荷分散用に利用が可能である。

【 0 0 6 1 】

図 8 は、初期閾値管理テーブル T b 2 の一例を示したものである。図 8 に示す初期閾値管理テーブル T b 2 は、使用モジュールと、負荷分散要と判断するための負荷閾値、負荷分散開始を決定する閾値オーバ回数についての情報を保持している。図 8 を例にとると、監視制御プログラムモジュールは負荷閾値に 1 0 0 が指定されており、負荷は 1 0 0 % を超えることはないから、負荷が高くなっても分散しないことを示している。一方、R A I D 制御プログラムモジュールの負荷閾値は 9 0 % であり、割当 M P U コアの負荷が 9 0 % を超えると閾値オーバ回数 1 回としてカウントされる。負荷は定期的に監視され、閾値オーバ回数が閾値オーバ規定回数欄に記載されている 5 回を超えると負荷分散を開始する。

40

【 0 0 6 2 】

図 9 は C P U 内部キャッシュ内で監視制御プログラムモジュールが管理する M P U コア割当テーブルの一例を示す図である。図 9 に示す M P U コア割当テーブルには、M P U コア N o、M P U コアの使用負荷率、M P U コアのグループ N o、状態、M P U コアを使用しているプログラムモジュール名、モジュールの割当形式、特権レベル、及び、負荷閾値のオーバ回数について情報を保持している。

【 0 0 6 3 】

50

図10は管理用端末31の管理端末画面32においてMPUCOA利用状況の情報を確認した際の画面を示している。画面にはMPUCOA NoとそのMPUCOAで動作しているプログラムモジュール名、状態、負荷状況が表示され、情報が更新されるたびに書き換えられる。

【0064】

図11～13を参照して本実施例に係るディスクアレイ装置の制御方法について説明する。

【0065】

図11を参照して、本実施例に係るディスクアレイ装置1の装置起動時において制御プログラムをマルチコアCPUの各MPUCOAに割り当てる処理手順について説明する。図11は装置起動時の制御プログラムのロードとマルチコアMPU110の各MPUCOA1110～111Nへのプログラムモジュール割り当ての実行手順を示すフローチャートである。

10

【0066】

図11に示すフローチャートは、ディスクアレイ装置1の電源スイッチがONにされ、ディスクアレイコントローラ11、12への通電が開始された際に実行される。

【0067】

ディスクアレイコントローラ11、12への通電が開始されると、マルチコアCPU110は、ブートROM上にあるイニシャルプログラムローダを実行する(ステップS1000)。イニシャルプログラムローダは、マルチコアCPU110の動作チェックと揮発性メモリ112のチェックを実施する(ステップS1010)。マルチコアCPUの動作チェックとは、一例としては、各MPUCOA1110～111Nが正常であるか否か確認させるCPUのセルフチェック機能が挙げられる。また揮発性メモリのチェックとはメモリへ特定パターンの値をライトした後、正しくリードできるか否か確認するチェック方法がある。

20

【0068】

イニシャルプログラムローダによってチェックした結果が正常であった場合には(ステップS1020: Yes)、次の処理に進み、制御プログラム119を揮発性メモリ112の制御プログラム配置領域117にロードし、イニシャルプログラムローダから制御プログラムへと処理を移す(ステップS1030)。

30

【0069】

制御プログラムモジュールでは、最初に監視制御プログラムモジュールPr1が起動される。そしてCPU内部キャッシュ1160上にMPUCOA管理テーブルTb6を作成し、監視制御プログラムモジュールPr1の情報を登録する(ステップS1040)。その後、監視制御プログラムモジュールPr1は初期MPUCOA割り当て管理テーブルTb1を参照し、他のモジュールに使用MPUCOAを割り当て、MPUCOA管理テーブルTb6の情報を更新する(ステップS1050)。次に、初期閾値管理テーブルTb2に基づき、負荷閾値と閾値オーバ回数規定値を設定する(ステップS1060)。MPUCOAを割り当てられた各プログラムモジュールは、起動して各制御処理を開始する(ステップS1070)。

40

【0070】

これらの一連の作業は、管理用端末31およびその管理画面32を通じて、管理者に報告される(ステップS1080)。

【0071】

一方、起動後にイニシャルプログラムロードがマルチコアCPU、揮発性メモリをチェックした結果が正常でなかった場合には(ステップS1020: No)、本ディスクアレイコントローラは動作不可能であると判断し、エラーを通知する。エラー通知とは例えば、イニシャルプログラムローダによる管理用端末31への通知や装置に搭載されたスピーカーによるアラーム音などの発生による。

【0072】

50

図12を参照して、本実施例に係るディスクアレイ装置1の制御プログラム内監視制御プログラムモジュールP b 1による自律負荷分散処理について説明する。図12は、監視制御プログラムモジュールP b 1が過負荷を検出し、当該過負荷プログラムモジュールの処理を負荷分散させるために追加のM P Uコアを割り当てる処理手順を示すフローチャートである。

【0073】

図12のフローチャートは、ディスクアレイ装置1の電源スイッチがONにされ、ディスクアレイコントローラ11、12を含む基本筐体、および、増設筐体に通電が開始され、図11に示した手順によって制御プログラムの起動処理が完了した後に、装置が停止されるまで監視制御プログラムモジュールによって繰り返し実行される。

10

【0074】

監視制御プログラムモジュールは、動作中の各プログラムモジュールへを使用しているM P Uコアごとの使用率(負荷状態)の問い合わせ情報を取得するとともに、M P Uコア管理テーブルT b 6の情報を更新する(ステップS 2 0 0 0)。次に、監視制御プログラムモジュールP r 1は、初期閾値管理テーブルT b 2に基づいて設定された負荷閾値を照らし合わせ、各負荷状態が閾値を超えているモジュールがないかチェックする。もし負荷閾値を超えているモジュールがあれば(ステップS 2 0 1 0 : Y e s)、閾値を超えている当該M P Uコアの閾値オーバ回数を1増加させ、M P Uコア管理テーブルT b 6上に反映する。次に、当該閾値オーバ回数が、規定回数を超えていないかチェックする。閾値オーバ回数が規定回数Nを越える値となった場合(ステップS 2 0 3 0 : Y e s)、監視制御プログラムモジュールP r 1は、当該プログラムモジュール処理の過負荷状態が続いており、負荷分散が必要であると判断する。監視制御プログラムモジュールP r 1は、M P Uコア管理テーブルT b 6を参照し、当該過負荷プログラムモジュールが割り当てられているM P UコアのグループN oを取得後、同一グループN oを持つ未割当のM P Uコアを検索する(ステップS 2 0 4 0)。未割当M P Uコアが存在した場合(ステップS 2 0 5 0 : Y e s)は、当該過負荷モジュールに未割当M P Uコアを割当て、M P U管理テーブルT b 6の情報を更新し、当該未割当コア欄に動作プログラムモジュール名の登録と、特権レベルの変更を行う。また、過負荷判定されたM P Uコアの閾値オーバ回数をゼロにクリアする。そして、これらの変化状況を管理端末に通知して(ステップS 2 2 3 0)処理を終了する。

20

30

【0075】

一方、当該過負荷モジュールに現在割り当てられているM P Uコアと同一のグループN oを持つ未割当M P Uコアが存在しない場合(ステップS 2 0 5 0 : N o)は、グループN oが-1、すなわちどのグループにも属していない未割当M P Uコアを検索する(ステップS 2 0 6 0)。

【0076】

どのグループにも属していない未割当M P Uコアが存在した場合(ステップS 2 0 7 0 : Y e s)は、当該過負荷モジュールに未割当M P Uコアを割当て、M P U管理テーブルT b 6の情報を更新し、当該未割当コア欄に動作プログラムモジュール名の登録と、特権レベルの変更を行う。また、過負荷判定されたM P Uコアの閾値オーバ回数をゼロにクリアする。そして、これらの変化状況を管理端末に通知して(ステップS 2 2 3 0)処理を終了する。

40

【0077】

どのグループにも属していない未割当M P Uコアが存在しなかった場合(ステップS 2 0 7 0 : Y e s)は、未割当M P Uコアの利用はできないので、すでに利用されている動的割当のM P Uコアのうち、もっとも利用率の低いM P Uコアの割当を変えることを考え、当該過負荷モジュールに現在割り当てられているM P Uコアと同一のグループN oを持ち、割当方式が「動的」であるM P Uコアを検索する(ステップ2 0 9 0)。検索して見つかった動的割当のM P Uコアの負荷状態が、当該過負荷モジュールの使用しているM P Uコアの負荷値の2倍よりも大きい場合(ステップS 2 0 9 0 : Y e s)、すなわち、過

50

負荷プログラムモジュールの処理を2つのMPUコアに負荷分散した際の1つあたりのMPUコア負荷が現在割り当てられている処理の負荷よりも大きい場合、現在割り当てられている処理を終了し、当該動的割当MPUコアの特権レベルを255に変更し、未割当状態にする。(ステップS2210)。その後、当該過負荷モジュールに未割当MPUコアを割当て、MPU管理テーブルTb6の情報を更新し、当該未割当コア欄に動作プログラムモジュール名の登録と、特権レベルの変更を行う。また、過負荷判定されたMPUコアの閾値オーバ回数をゼロにクリアする。そして、これらの変化状況を管理端末に通知して(ステップS2230)処理を終了する。

【0078】

検索して見つかった動的割当のMPUコアの負荷状態が、当該過負荷モジュールの使用しているMPUコアの負荷値の2倍と同等かそれよりも大きい場合(ステップS2090:No)、当該プログラムモジュールの負荷分散を断念し、管理端末31に高負荷状態であることを通知し、処理を終了する。

【0079】

なお、図12の説明の冒頭に述べた、各MPUコアの負荷状態監視の際に、負荷閾値を超えたモジュールが存在しない場合(ステップS2010:No)や、存在しても閾値オーバ回数が規定回数未満(ステップ2030:No)だった場合には、負荷分散のためのMPU割当変更処理は不要であると判断し、一定時間待ち(ステップS2200)、間隔を置いた後、負荷状態の監視を繰り返す。

【0080】

図13を参照して、本実施例に係るディスクアレイ装置のディスクアレイコントローラにおいて、マルチコアCPUの一部MPUコアに障害が発生した際の処理について説明する。マルチコアCPUの内部には、複数のMPUコアが存在するが、物理的に切り離すことができないため、あるMPUコアに物理障害が生じた場合には、CPU全体を交換する必要がある。図13は、この障害対応のため、動作中のプログラムモジュールをフェイルオーバさせるために停止する処理手順を示したフローチャートである。

【0081】

図13に示すフローチャートは、マルチコアCPU110のうち、ある特定のMPUコア111x(xは整数)に障害が発生した際に実行される。

【0082】

監視制御プログラムPr1は、あるMPUコア111xに障害が発生したのを検出すると(ステップ3000)、MPUコア割当テーブルTb6を参照してMPUコア111xの割当状況を確認するとともに、状態を「障害」に変更する(ステップS3010)。具体的には、割当済みであるか否かと、割当済みである場合の割当形式について確認する。

【0083】

障害が発生したMPUコア111xが固定割当てで利用されている場合(ステップS3020:Yes)、MPUコア111xは他の処理の監視制御や管理用端末31による運用管理制御を行っているコアであるため、正常なフェイルオーバ処理が行えなくなる可能性がある。そこで、監視制御プログラムPr1は動作可能であれば、管理端末31への障害通知と保守交換用のエラー表示LED125の点灯を試み(ステップS3200)、信号線101を通じて他系のディスクアレイコントローラに障害情報を送信して、処理の引継ぎ依頼を通知する(ステップS3070)。なお、ステップS3200、S3070に示す処理が行えないような重大な障害が発生している場合には、信号線101を介したディスクアレイコントローラ間の定期的な通信が断絶するため、他系のディスクアレイコントローラによってこの障害を検出される。

【0084】

他系のディスクアレイコントローラは、ステップS3070に示すように処理引継ぎの依頼を受けた場合や、もしくはディスクコントローラ間にある信号線101を介した通信できなくなった場合には、もう一方のディスクアレイコントローラに障害が発生したものと判断し、信号線を介して、障害が発生しているディスクアレイコントローラのマルチコ

10

20

30

40

50

アCPU110にリセット信号を常時印加する。常時リセット信号が印加されていることにより、障害マルチコアCPU110はリセットされ続けるので、それ以上動作することができなくなる。このため、MPUCOA111xの誤動作によってデータを破壊してしまうような恐れがなくなる。なお、ここでは常時リセットにより動作を停止させているが、給電の停止など他の手段によってもよい。

【0085】

障害を発生したMPUCOA111xが半固定割り当てで利用されている場合、監視制御プログラムモジュールPr1は、MPUCOA割当テーブルTb6を確認し、障害を発生したMPUCOA111xと同一のプログラムモジュールが存在しないか検索する(ステップS3300)。

10

【0086】

同一のプログラムモジュールが見つかった場合には、監視制御プログラムモジュールPr1は当該プログラムが動作しているMPUCOAの割当形式を動的から半固定に変更し(ステップS3310)、その後管理端末に障害を通知して、保守交換用のLEDを点灯させる(ステップS3040)。

【0087】

次に監視制御プログラムは、MPUCOA割当テーブルTb6を参照して、特権レベルの低いプログラムモジュールから順に終了していく(ステップ3050)。

【0088】

以降ステップS3070~S3080に処理を進めるが、ステップS3070~S3080については、既に説明しているため省略する。

20

【0089】

一方、障害コアに割り当てられていたプログラムモジュールと同一のモジュールが割り当てられたMPUCOAが存在しない場合には、ステップS3040に処理を進める。

【0090】

障害を発生したMPUCOA111が、固定割当でもなく(ステップS3020:No)、また半固定でもない場合(ステップS3030:No)、すなわち動的で割り当てられている場合には、監視制御プログラムモジュールPr1はステップS3040~S3080に示す処理を行う。なお、ステップS3040~S3080の処理については先に説明しているため省略する。

30

【0091】

以上の実施例によれば、本発明にかかるディスクアレイ装置は、複数個あるMPUCOAを管理テーブルを用いて固定的または半固定的または動的な資源として管理するために、他の障害ですべての処理に影響が生じるような問題を回避することができる。また各MPUCOAの負荷状況を監視することで、CPU資源をコア単位で効率よく活用することができる。

【0092】

また、障害発生時に特権レベルの確認によって動作中のプログラムモジュールの依存関係を考慮しつつフェイルオーバーすることが可能であり、マルチコアCPUで生じる特定コアの障害時の対応を行うことが可能である。

40

【0093】

以上実施例で説明したが、本発明の他の実施形態1は、前記一の単位プロセッサは、非動的になるように割り当てる制御プログラムのうちの1以上を、非動的になるように割り当てる単位プロセッサとは別の単位プロセッサに動的になるように割り当てるディスクアレイ装置である。

【0094】

本発明の他の実施形態2は、前記一の単位プロセッサは、非動的になるように割り当てる制御プログラムのうちの1以上を、非動的になるように割り当てる単位プロセッサのみに割り当てるディスクアレイ装置である。

【0095】

50

本発明の他の実施形態3は、前記一の単位プロセッサは、非動的になるように制御プログラムを割り当てた単位プロセッサの動作負荷が相対的に高いとき、動的になるように割り当てた制御プログラムの動作負荷が相対的に低い単位プロセッサの処理を終了させ、終了させた単位プロセッサに、前記相対的に動作負荷が高い単位プロセッサに割り当てた制御プログラムを、動的になるように割り当てるディスクアレイ装置である。

【0096】

本発明の他の実施形態4は、前記一の単位プロセッサは、制御プログラムの動作負荷が相対的に高い単位プロセッサが属するグループ内で、かつ、動的になるように割り当てた制御プログラムの動作負荷が相対的に低い単位プロセッサの処理を終了させ、終了させ単位プロセッサに、前記相対的に動作負荷が高い単位プロセッサに割り当てた制御プログラムを、動的になるように割り当てるディスクアレイ装置である。

10

【0097】

本発明の他の実施形態5は、前記一の単位プロセッサは、どのグループにも属さない単位プロセッサに、前記相対的に動作負荷が高い単位プロセッサに割り当てた制御プログラムを、動的になるよう割り当てるとともに、割り当てた単位プロセッサを、前記相対的に動作負荷が高い単位プロセッサが属するグループに属するとするディスクアレイ装置である。

【0098】

本発明の他の実施形態6は、前記ディスクアレイコントローラは、物理的に分離不可能な複数個のプロセッサコアを内蔵したCPUを搭載しており、各プロセッサコアはそれぞれ単位プロセッサとなり、一の単位プロセッサは、すべての単位プロセッサを個別に管理し、各々完結した制御プログラムを、すべての単位プロセッサにプロセッサ単位で割り当てる際に、適宜動作終了可能とする動的になるように、又は前記CPU全体の動作終了まで動作可能とする非動的になるように割り当て、割り当てた制御プログラムの動作負荷又は動作状態を単位プロセッサごとに管理し、特定のプロセッサコアに障害が発生した場合に、該プロセッサコアの単位プロセッサに制御プログラムを動的に又は非動的になるように割り当てたか及び制御プログラムの終了順を考慮し、すべての制御プログラムを終了させるディスクアレイ装置である。

20

【0099】

本発明の他の実施形態7は、前記一の単位プロセッサは、前記障害が発生したプロセッサコアの単位プロセッサに割り当てた制御プログラムが非動的になるよう割り当てたものであり、かつ、前記制御プログラムを動的になるよう割り当てた別の単位プロセッサがあるとき、前記制御プログラムを、前記別の単位プロセッサに非動的になるよう割り当てたとするディスクアレイ装置である。

30

【0100】

本発明の他の実施形態8は、前記一の単位プロセッサは、前記障害が発生したプロセッサコアの単位プロセッサに割り当てた制御プログラムが該単位プロセッサのみに非動的になるように割り当てた制御プログラムであるとき、前記CPU全体の動作を終了させるディスクアレイ装置である。

【0101】

40

本発明の他の実施形態9は、前記ディスクアレイコントローラは、物理的に分離不可能な複数個のプロセッサコアを内蔵したCPUを搭載しており、内部スイッチ、不揮発性メモリコントローラ、揮発性メモリコントローラ、ディスクアレイコントローラ間転送コントローラ、パリティ演算器及びCPU内部キャッシュを有し、前記制御プログラムは、監視制御プログラムモジュール、運用管理プログラムモジュール、RAID制御プログラムモジュール、NAS制御プログラムモジュール、ホストI/O制御プログラムモジュール、ドライブI/O制御プログラムモジュール、初期MPUCOA割当管理テーブル、初期閾値管理テーブル、RAID設定管理テーブル、LU設定管理テーブル、及びNASボリューム管理テーブルを備えており、前記一の単位プロセッサは、自己及び他の単位プロセッサの制御プログラムの負荷状態を監視し、負荷閾値及び閾値オーバ回数が規定回数を超え

50

ると、同一グループ内の未割当コアを検索し、検出した未割当コアに、前記制御プログラムを割り当てるディスクアレイ装置である。

【0102】

本発明の他の実施形態10は、物理的に分離不可能な複数個のプロセッサコアを内蔵したCPUを搭載し、各プロセッサコアはそれぞれ単位プロセッサとなるディスクアレイコントローラを備えるディスクアレイ装置の制御方法であって、すべての単位プロセッサを個別に管理し、各々完結した制御プログラムを、すべての単位プロセッサにプロセッサ単位で割当ての際に、適宜動作終了可能とする動的になるように、又は前記CPU全体の動作終了まで動作可能とする非動的になるように割り当て、割り当てた制御プログラムの動作負荷又は動作状態を単位プロセッサごとに管理するディスクアレイ装置の制御方法である。

10

【図面の簡単な説明】

【0103】

【図1】本実施例に係るディスクアレイ装置の概略構成を示す説明図。

【図2】本実施例に係るディスクアレイ装置の概観図。

【図3】本実施例に係るディスクアレイ装置におけるディスクアレイコントローラの機能的構成を示すブロック図。

【図4】本実施例に係るディスクアレイ装置におけるディスクアレイコントローラの概観図。

【図5】ディスクアレイコントローラに搭載される複数プロセッサコア(MPUコア)を内蔵したCPU(マルチコアCPU)の機能的構成の一例を示すブロック図。

20

【図6】制御プログラムに含まれる各プログラムモジュール及び管理テーブルを示す説明図。

【図7】制御プログラムの初期MPUコア割当管理テーブルの一例を示す図。

【図8】制御プログラムの初期閾値管理テーブルの一例を示す図。

【図9】MPUコア割当テーブルの一例を示す図。

【図10】管理用端末装置31の管理画面32におけるMPUコア利用状況に関する表示画面。

【図11】制御プログラムをマルチコアCPUの各MPUコアに割当てる処理手順を示すフローチャート。

30

【図12】監視制御プログラムモジュールが負荷分散制御指示行う際の動作を示すフローチャート。

【図13】マルチコアCPUの一部MPUコアに障害が発生した際の障害対応処理を表すフローチャート。

【符号の説明】

【0104】

1 ... ディスクアレイ装置

11、12 ... ディスクアレイコントローラ

101、102、103 ... 信号線

104、107 ... 電力線

40

105、106 ... 電源

110 ... 複数MPUコア内蔵CPU(マルチコアCPU)

111 ... 不揮発性メモリ

112 ... 揮発性メモリ

113 ... ホスト側接続ポート

114 ... ドライブ側接続ポート

115 ... 管理用ネットワーク用物理ポート

116 ... ブートROM

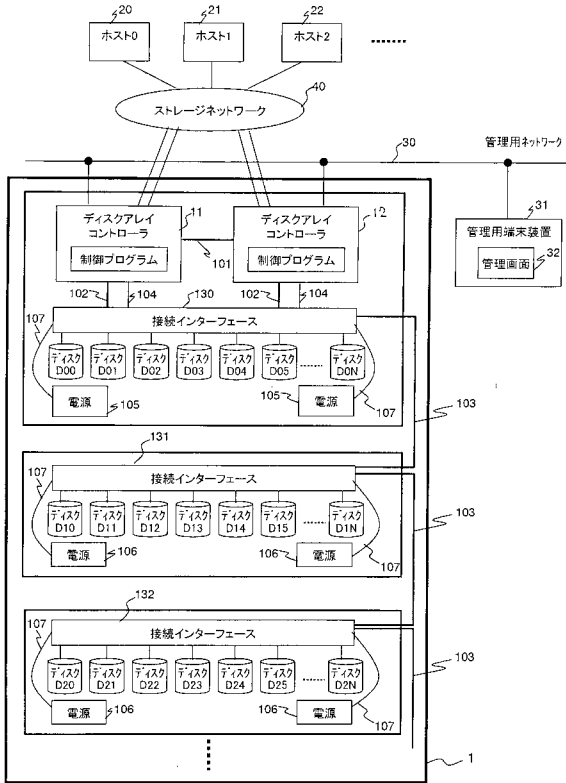
117 ... 制御プログラム配置領域

118 ... データバッファ領域

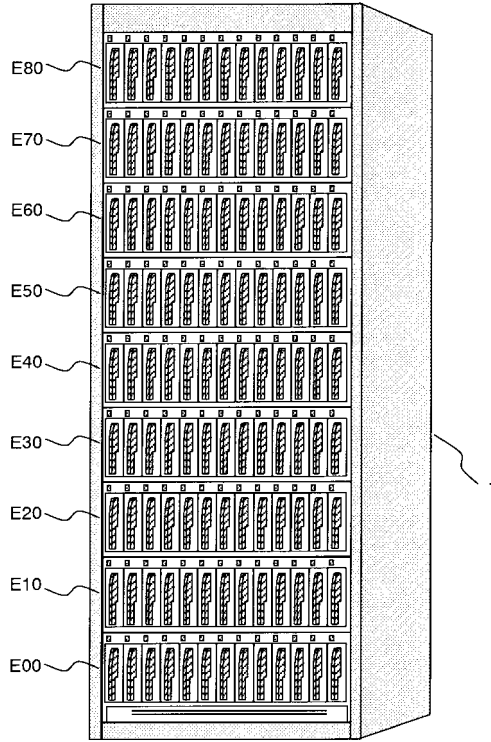
50

1 1 9 ... 制御プログラム	
1 1 1 0 ~ 1 1 1 N ... M P U コア	
1 1 2 0 ... C P U 内部スイッチ	
1 1 3 0 ... 不揮発性メモリコントローラ	
1 1 4 0 ... ディスクアレイコントローラ間転送コントローラ	
1 1 5 0 ... パリティ演算器	
1 1 6 0 ... C P U 内部キャッシュ	
1 1 7 0 ... 揮発性メモリコントローラ	
1 2 0 ... 信号線	
1 2 2 ... 接続コネクタ	10
1 2 3 ... ホスト側接続口	
1 2 4 ... 管理用ネットワーク接続口	
1 2 5 ... エラー L E D	
1 3 0、1 3 1、1 3 2 ... 接続インターフェース	
D 0 0 ~ D 2 N ... ディスク装置	
E 0 0 ~ E 8 0 ... ディスク筐体	
2 0、2 1、2 2 ... ホスト	
3 0 ... 管理用ネットワーク	
3 1 ... 管理用端末装置	
3 2 ... 管理画面	20
4 0 ... ストレージネットワーク	
P r 1 ... 監視制御プログラムモジュール	
P r 2 ... 運用管理プログラムモジュール	
P r 3 ... R A I D 制御プログラムモジュール	
P r 4 ... N A S 制御プログラムモジュール	
P r 5 ... ホスト I / O 制御プログラムモジュール	
P r 6 ... ドライブ I / O 制御プログラムモジュール	
T b 1 ... 初期 M P U コア 割当管理テーブル	
T b 2 ... 初期 閾値管理テーブル	
T b 3 ... R A I D 設定管理テーブル	30
T b 4 ... L U 設定管理テーブル	
T b 5 ... N A S ボリューム管理テーブル	
T b 6 ... M P U コア 管理テーブル	

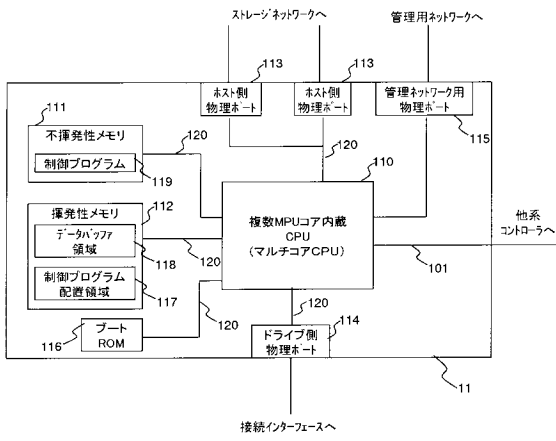
【図1】



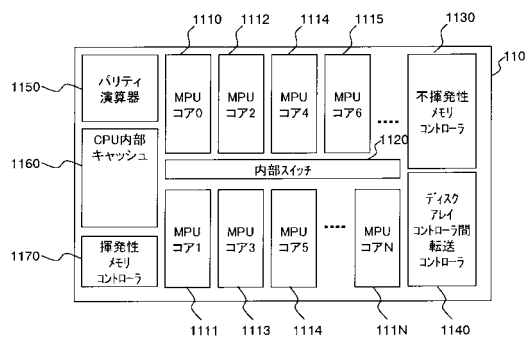
【図2】



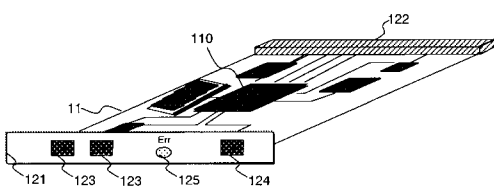
【図3】



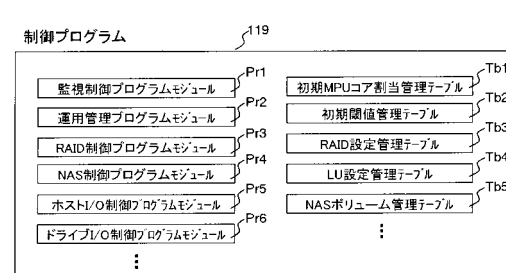
【図5】



【図4】



【図6】



【図7】

Tb1

MPUコアNo.	使用モジュール	割当形式	特権レベル	グループNo
0	監視制御	固定	0	001
1	運用管理	固定	1	002
2	RAID制御	半固定	2	003
3	NAS制御	半固定	3	003
4	ホストI/O制御	半固定	4	004
5	ドライブI/O制御	半固定	4	004
6	(未割当)	動的	255	003
7	(未割当)	動的	255	003
8	(未割当)	動的	255	004
9	(未割当)	動的	255	-1

【図8】

Tb2

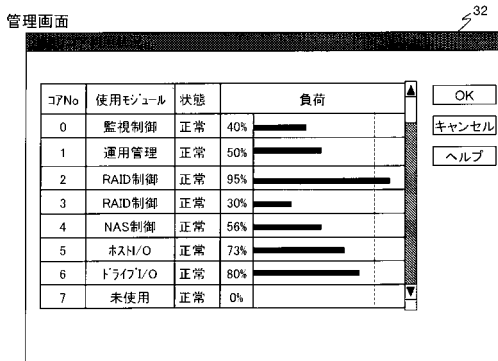
使用モジュール	負荷閾値(%)	閾値オーバー規定回数
監視制御	100	1
運用管理	100	1
RAID制御	90	5
NAS制御	90	5
ホストI/O制御	90	5
ドライブI/O制御	90	5

【図9】

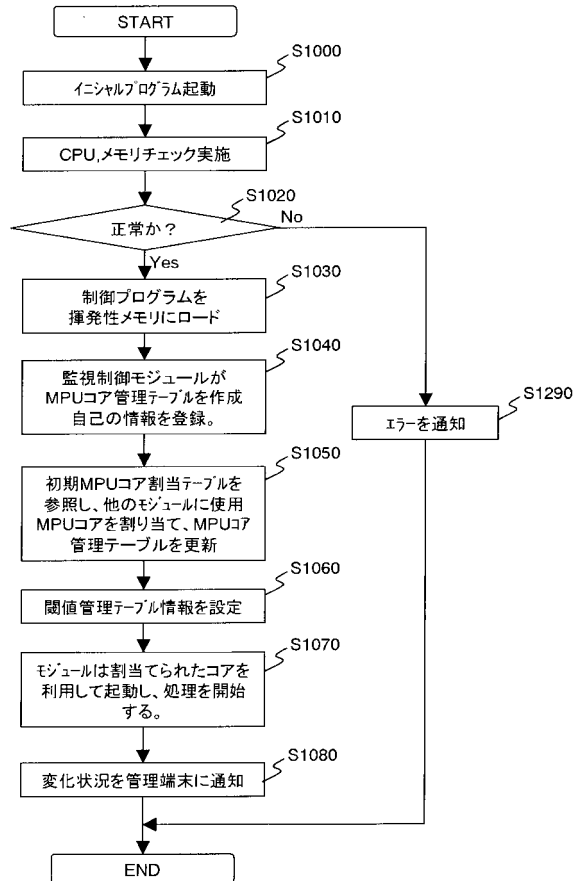
Tb6

MPUコアNo.	負荷(%)	グループNo.	状態	使用モジュール	割当形式	特権	負荷閾値オーバー回数
0	50%	001	正常	監視制御	固定	0	0
1	60%	002	正常	運用管理	固定	1	0
2	70%	003	正常	RAID制御	半固定	2	0
3	95%	003	正常	NAS制御	半固定	3	0
4	70%	003	正常	ホストI/O制御	半固定	4	0
5	40%	004	正常	ドライブI/O制御	半固定	4	0
6	0%	004	正常	(未割当)	動的	255	0
7	0%	004	正常	(未割当)	動的	255	0
8	0%	000	正常	(未割当)	動的	255	0

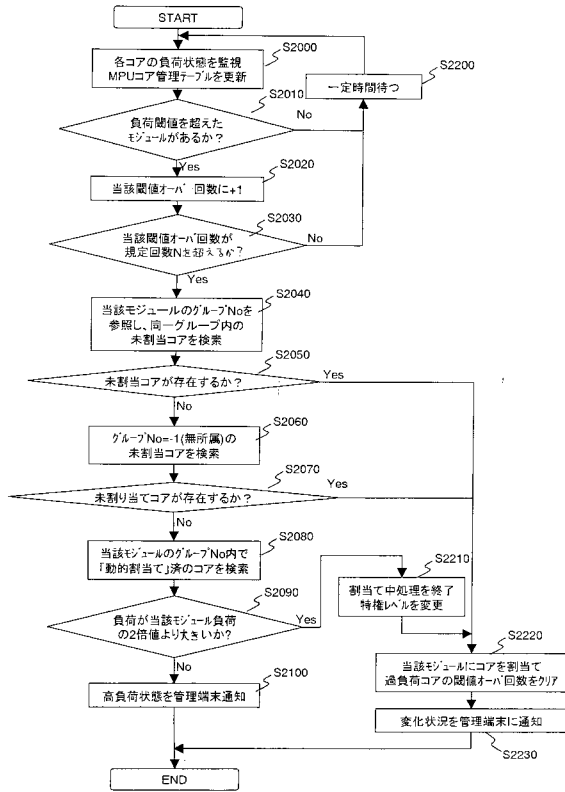
【図10】



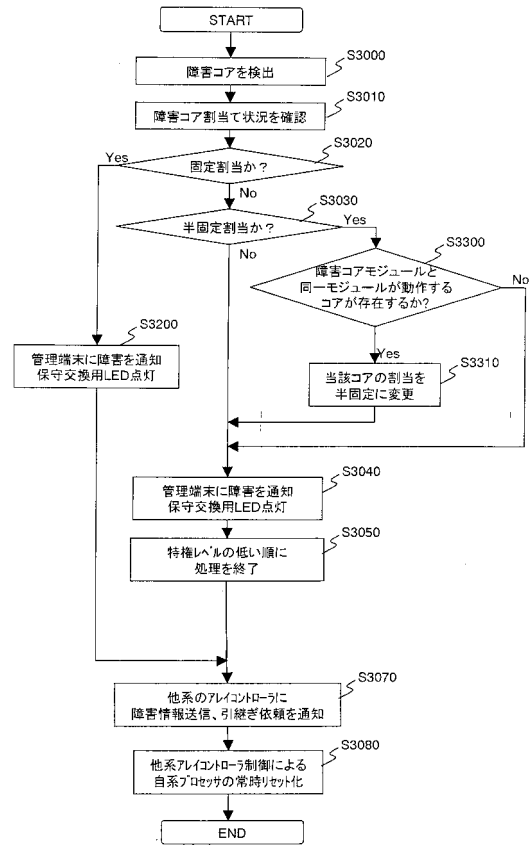
【図11】



【図12】



【図13】



フロントページの続き

- (56)参考文献 特開平02-238556(JP,A)
特開2001-167040(JP,A)
特開平04-235662(JP,A)
特開平09-265435(JP,A)
特開2000-148709(JP,A)
特開2002-353960(JP,A)
特開2005-099984(JP,A)
特開平09-274608(JP,A)
特開平11-053327(JP,A)
特開2003-208362(JP,A)
特開平11-316726(JP,A)
特開2005-266841(JP,A)
枝廣 正人,マルチコア向けソフトウェア・プラットフォームを開発し携帯電話機に適用,日経エレクトロニクス,日本,日経BP社,2005年 3月28日,第896号,125~136頁
進藤 智則,シングルコアよりマルチコア 第1部<発想の転換> 迫るマイクロプロセッサ危機 性能はコアの数で稼ぐ,日経エレクトロニクス,日本,日経BP社 Nikkei Business Publications,Inc.,2004年 8月30日,第881号,98~105頁
進藤 智則,シングルコアよりマルチコア 第3部<ハードウェア> オンチップの恩恵生かす 高速バスや記憶階層がカギ,日経エレクトロニクス,日本,日経BP社 Nikkei Business Publications,Inc.,2004年 8月30日,第881号,116~121頁
高務 祐哲 他,エネルギー最小周波数を利用したタスク再配置によるマルチプロセッサ向け消費エネルギー削減手法,電子情報通信学会技術研究報告 ICD2004-233~246 集積回路,日本,社団法人電子情報通信学会,2005年 3月 4日,第104巻,第711号,37~42頁

(58)調査した分野(Int.Cl.,DB名)

G06F 3/06