

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第3641872号  
(P3641872)

(45) 発行日 平成17年4月27日(2005.4.27)

(24) 登録日 平成17年2月4日(2005.2.4)

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

F I

G06F 3/06  
G06F 12/08

G06F 3/06 540  
G06F 3/06 301X  
G06F 3/06 302E  
G06F 3/06 304N  
G06F 12/08 320

請求項の数 12 (全 16 頁)

(21) 出願番号	特願平8-85370	(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8年4月8日(1996.4.8)	(74) 代理人	100095511 弁理士 有近 紳志郎
(65) 公開番号	特開平9-274544	(72) 発明者	山本 康友 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
(43) 公開日	平成9年10月21日(1997.10.21)	(72) 発明者	山本 彰 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
審査請求日	平成12年3月10日(2000.3.10)		
審判番号	不服2002-4389(P2002-4389/J1)		
審判請求日	平成14年3月14日(2002.3.14)		

最終頁に続く

(54) 【発明の名称】 記憶装置システム

(57) 【特許請求の範囲】

【請求項1】

データ処理装置が直接アクセスを行う論理的記憶装置が配置された物理的記憶装置を含む複数の物理的記憶装置と、前記データ処理装置と前記複数の物理的記憶装置の間のデータ転送を制御する記憶制御装置とを有する記憶装置システムにおいて、  
予め定めた指標に基づいて、前記論理的記憶装置を、前記複数の物理的記憶装置のうち、該論理的記憶装置が配置された物理的記憶装置とは異なる物理的記憶装置に再配置する際に、再配置先の物理的記憶装置に、前記論理的記憶装置に格納されたデータ全体を連続的に格納する論理的記憶装置再配置手段を有することを特徴とする記憶装置システム。

【請求項2】

データ処理装置が直接アクセスを行う論理的記憶装置が配置された物理的記憶装置を含む複数の物理的記憶装置と、前記データ処理装置と前記複数の物理的記憶装置の間のデータ転送を制御する記憶制御装置とを有する記憶装置システムにおいて、  
前記データ転送の制御の運用中にデータ処理装置の論理的記憶装置へのアクセス情報を指標として採取するアクセス情報採取手段と、前記指標に基づいて前記論理的記憶装置を、前記複数の物理的記憶装置のうち、該論理的記憶装置が配置された物理的記憶装置とは異なる物理的記憶装置に再配置する際に、再配置先の物理的記憶装置に、前記論理的記憶装置に格納されたデータ全体を連続的に格納する論理的記憶装置再配置手段とを有することを特徴とする記憶装置システム。

【請求項3】

10

20

請求項 2 に記載の記憶装置システムにおいて、前記アクセス情報が、前記データ処理装置から前記論理的記憶装置へのアクセス頻度情報を含むことを特徴とする記憶装置システム。

【請求項 4】

請求項 2 または請求項 3 に記載の記憶装置システムにおいて、前記アクセス情報が、前記データ処理装置から前記論理的記憶装置へのアクセスパターン情報を含むことを特徴とする記憶装置システム。

【請求項 5】

請求項 1 に記載の記憶装置システムにおいて、前記指標が、前記論理的記憶装置に求められる信頼性であることを特徴とする記憶装置システム。

10

【請求項 6】

請求項 1 から請求項 5 のいずれかに記載の記憶装置システムにおいて、前記指標を保守員に提示する指標提示手段と、保守員からの再配置指示を受け付ける再配置指示受付手段とを具備したことを特徴とする記憶装置システム。

【請求項 7】

請求項 1 から請求項 5 のいずれかに記載の記憶装置システムにおいて、データ処理装置からの再配置指示を受け付ける再配置指示受付手段を具備したことを特徴とする記憶装置システム。

【請求項 8】

複数の物理的記憶装置と、データ処理装置と前記複数の物理的記憶装置の間のデータ転送を制御する記憶制御装置とを有する記憶装置システムにおいて、

20

前記複数の物理的記憶装置のうちの第一の物理的記憶装置群にまたがって設定され、前記データ処理装置が直接アクセスを行う論理的記憶装置を有し、

前記記憶制御装置は、予め定めた指標に基づいて、前記第一の物理的記憶装置群とは異なる第二の物理的記憶装置群を選択すると共に、前記第二の物理的記憶装置群に、前記論理的記憶装置に格納されたデータ全体を連続的に格納する論理的記憶装置再配置手段を有することを特徴とする記憶装置システム。

【請求項 9】

請求項 1 から請求項 8 のいずれかに記載の記憶装置システムにおいて、再配置中の論理的記憶装置にデータ処理装置からのアクセスがあったとき、再配置中の論理的記憶装置の再配置完了領域と再配置未完領域とを識別し、前記アクセス位置が前記再配置完了領域ならば再配置先の物理的記憶装置にアクセスさせ、前記アクセス位置が前記再配置未完領域ならば再配置前の物理的記憶装置にアクセスさせるアクセス位置切替手段をさらに具備したことを特徴とする記憶装置システム。

30

【請求項 10】

データ処理装置が直接アクセスを行う論理的記憶装置が配置された物理的記憶装置を含む複数の物理的記憶装置と、前記データ処理装置と前記複数の物理的記憶装置の間のデータ転送を制御する記憶制御装置とを有する記憶装置システムにおいて、

前記記憶制御装置は、前記データ処理装置によるアクセス状況を取得し、前記アクセス状況に基づいて前記論理的記憶装置のデータを、前記複数の物理的記憶装置のうち、該論理的記憶装置に配置されている物理的記憶装置から、該論理的記憶装置が配置された物理的記憶装置とは異なる物理的記憶装置に移動させることを特徴とする記憶装置システム。

40

【請求項 11】

キャッシュメモリを有し、

前記記憶制御装置は、前記第一の物理的記憶装置群から前記キャッシュメモリへデータを読み出す手段と、前記読み出されたデータを前記第二の物理的記憶装置群に書き込む手段とを有することを特徴とする請求項 8 記載の記憶装置システム。

【請求項 12】

前記第二の物理的記憶装置群は、前記第一の物理的記憶装置群よりもデータの読み出し速度が高速であることを特徴とする請求項 8 又は請求項 11 記載の記憶装置システム。

50

## 【発明の詳細な説明】

## 【0001】

## 【発明の属する技術分野】

本発明は、記憶装置システムに関し、さらに詳しくは、シーケンシャルアクセスの場合やランダムアクセスでヒット率が低い場合でもアクセス性能を向上することが出来る記憶装置システムおよびデータの信頼性を向上することが出来る記憶装置システムに関する。

特に、本発明は、ディスクアレイ向きの高機能ディスク装置、その高機能ディスク装置とディスク制御装置とにより構成される記憶装置サブシステム、およびその記憶装置サブシステムとデータ処理装置とにより構成される情報処理システムに有用である。

## 【0002】

## 【従来の技術】

シカゴのイリノイ大学で開かれた「ACM SIGMOD」会議において発表された論文「D.Patterson,G.gibson,and R.H.Kartz;A Case for Redundant Arrays of Inexpensive Disks (RAID),ACM SIGMOD Conference,Chicago,IL,(June 1988),pp.109-116」は、ディスクアレイ上のデータ配置に関する技術を開示している。

## 【0003】

また、特開平7-84732号公報では、ディスク装置の一部をディスクキャッシュの如く用いる技術が開示されている。具体的には、ディスク装置を一時的にデータを格納するテンポラリ領域と最終的にデータを書き込む領域とに分け、更新データはパリティを生成せず一旦テンポラリ領域に二重書きし、非同期にパリティ生成し、最終領域に書き込む。

## 【0004】

一方、電気情報通信学会技術研究報告「DE95-68(茂木他:Hot Mirroringを用いたディスクアレイのディスク故障時の性能評価、1995年12月、電気情報通信学会技報 Vol.95-No.407, pp.19-24)」には、アクセス頻度の違いにより、データを保持するRAIDレベルを動的に変更する技術が開示されている。具体的には、ディスク装置をRAID1構成の部分とRAID5構成の部分に分け、ライトアクセスのあったデータを優先的にRAID1構成の部分に格納するようにデータの格納位置を動的に変更することにより、アクセス頻度の高いデータはRAID1構成の部分に格納し、アクセス頻度の低いものはRAID5構成の部分に格納するように出来る。

この技術によれば、記憶容量の異なる物理ディスク装置やRAIDレベルの異なる物理ディスク装置を記憶装置サブシステム内で混在させることが可能であり、論理ディスク装置内のデータを、そのアクセス頻度やアクセスパターンなどの指標に基づいて、任意の物理ディスク装置に格納することが出来る。また、アクセス頻度の高いデータを、より高速な物理ディスク装置に格納するように、動的に格納位置を変更することも出来る。

なお、RAID1のディスクアレイは、データ処理装置からの書き込みデータに対して、その複製をミラーと呼ばれる副ディスク装置に書き込み、データの信頼性を確保する。冗長データが元のデータの複製であるため、冗長データ作成のオーバーヘッドが小さく、アクセス性能が良い。但し、物理的記憶装置の使用効率は、50%と低い。一方、RAID5のディスクアレイは、データ処理装置からの複数の書き込みデータに対して、パリティと呼ばれる冗長データを作成する。パリティ作成時に更新前データと更新前パリティのリードが必要なため、冗長データ作成のオーバーヘッドが大きく、アクセス性能は悪い。但し、複数のデータに対して1つのパリティを作成するため、記憶装置の使用効率はRAID1に比べ高い。

## 【0005】

## 【発明が解決しようとする課題】

上記従来技術では、アクセスするデータ単位でデータの格納位置の変更を行うため、データ処理装置が直接アクセスを行う論理ディスク装置上では連続なデータが、実際にデータを記憶する物理ディスク装置上では非連続となってしまう。このため、一連のデータをリード/ライトするシーケンシャルアクセスの場合、実際には複数データをまとめてリード

10

20

30

40

50

ノライトできなくなり、アクセス性能の低下を招く問題点がある。

【0006】

一方、上記報告「DE95-68」の従来技術では、ライトの度に、アクセス頻度が低いと判断したデータをRAID1構成の部分からRAID5構成の部分に移し、空いたRAID1構成の部分にライトデータを書き込むため、アクセスパターンがランダムアクセスでヒット率が低い場合には、RAID1構成の部分に移したデータの多くは再びRAID5構成の部分に戻されることになる。このため、ヒット率が低い場合、アクセス性能の向上は期待できず、逆にデータを移す処理のオーバーヘッドがアクセス性能の低下を引き起こす問題点がある。

【0007】

また、上記の従来技術では、データの信頼性の向上については全く考慮されていない問題点がある。

【0008】

そこで、本発明の第1の目的は、シーケンシャルアクセスの場合やランダムアクセスでヒット率が低い場合でも、アクセス性能を向上することが出来る記憶装置システムを提供することにある。

また、本発明の第2の目的は、データの信頼性を向上することが出来る記憶装置システムを提供することにある。

【0009】

【課題を解決するための手段】

第1の観点では、本発明は、データ処理装置が直接アクセスを行う論理的記憶装置が配置された物理的記憶装置を含む複数の物理的記憶装置と、前記データ処理装置と前記複数の物理的記憶装置の間のデータ転送を制御する記憶制御装置とを有する記憶装置システムにおいて、予め定めた指標に基づいて、前記論理的記憶装置を、前記複数の物理的記憶装置のうち、該論理的記憶装置が配置された物理的記憶装置とは異なる物理的記憶装置に再配置する際に、再配置先の物理的記憶装置に、前記論理的記憶装置に格納されたデータ全体を連続的に格納する論理的記憶装置再配置手段を有することを特徴とする記憶装置システムを提供する。上記第1の観点による記憶装置システムでは、アクセスするデータ単位でデータの格納位置の変更を行うのではなく、論理的記憶装置を単位として物理的記憶装置への再配置を行い、且つ、再配置先の物理的記憶装置にデータを連続的に格納する。従って、シーケンシャルアクセスの場合でも、アクセス性能を向上することが出来る。また、ライトの度にデータの格納位置の変更を行うのではなく、予め定めた指標に基づいて前記再配置を行うから、ランダムアクセスでヒット率が低い場合でも、アクセス性能を向上することが出来る。

【0010】

第2の観点では、本発明は、データ処理装置が直接アクセスを行う論理的記憶装置が配置された物理的記憶装置を含む複数の物理的記憶装置と、前記データ処理装置と前記複数の物理的記憶装置の間のデータ転送を制御する記憶制御装置とを有する記憶装置システムにおいて、前記データ転送の制御の運用中にデータ処理装置の論理的記憶装置へのアクセス情報を指標として採取するアクセス情報採取手段と、前記指標に基づいて前記論理的記憶装置を、前記複数の物理的記憶装置のうち、該論理的記憶装置が配置された物理的記憶装置とは異なる物理的記憶装置に再配置する際に、再配置先の物理的記憶装置に、前記論理的記憶装置に格納されたデータ全体を連続的に格納する論理的記憶装置再配置手段とを有することを特徴とする記憶装置システムを提供する。

上記第2の観点による記憶装置システムでは、アクセスするデータ単位でデータの格納位置の変更を行うのではなく、論理的記憶装置を単位として物理的記憶装置への再配置を行い、且つ、再配置先の物理的記憶装置にデータを連続的に格納する。従って、シーケンシャルアクセスの場合でも、アクセス性能を向上することが出来る。また、ライトの度にデータの格納位置の変更を行うのではなく、アクセス情報を採取し、それを統計的に利用して前記再配置を行うから、ランダムアクセスでヒット率が低い場合でも、アクセス性能を

10

20

30

40

50

向上することが出来る。

【0011】

第3の観点では、本発明は、上記構成の記憶装置システムにおいて、前記アクセス情報が、前記データ処理装置から前記論理的記憶装置へのアクセス頻度情報を含むことを特徴とする記憶装置システムを提供する。

上記第3の観点による記憶装置システムでは、アクセス頻度の高い論理的記憶装置をより高速な物理的記憶装置へ再配置することが出来る。従って、アクセス性能を向上することが出来る。

【0012】

第4の観点では、本発明は、上記構成の記憶装置システムにおいて、前記アクセス情報が、前記データ処理装置から前記論理的記憶装置へのアクセスパターン情報を含むことを特徴とする記憶装置システムを提供する。

上記第4の観点による記憶装置システムでは、シーケンシャルアクセスの比率の高い論理的記憶装置をよりシーケンシャルアクセス性能の高い物理的記憶装置へ再配置することが出来る。従って、アクセス性能を向上することが出来る。

【0013】

第5の観点では、本発明は、上記構成の記憶装置システムにおいて、前記指標が、前記論理的記憶装置に求められる信頼性であることを特徴とする記憶装置システムを提供する。

上記第5の観点による記憶装置システムでは、信頼性が高いことが求められる論理的記憶装置をより信頼性の高い物理的記憶装置へ再配置することが出来る。従って、データの信頼性を向上することが出来る。

【0014】

第6の観点では、本発明は、上記構成の記憶装置システムにおいて、前記指標を保守員に提示する指標提示手段と、保守員からの再配置指示を受け付ける再配置指示受付手段とを具備したことを特徴とする記憶装置システムを提供する。

上記第6の観点による記憶装置システムでは、保守員が再配置指示を入力できるため、非常に柔軟に前記再配置を行うことが出来る。

【0015】

第7の観点では、本発明は、上記構成の記憶装置システムにおいて、データ処理装置からの再配置指示を受け付ける再配置指示受付手段を具備したことを特徴とする記憶装置システムを提供する。

上記第7の観点による記憶装置システムでは、データ処理装置が再配置指示を入力できるため、保守員では判断不可能な高度の条件下で前記再配置を行うことが出来る。

【0016】

第8の観点では、本発明は、複数の物理的記憶装置と、データ処理装置と前記複数の物理的記憶装置の間のデータ転送を制御する記憶制御装置とを有する記憶装置システムにおいて、前記複数の物理的記憶装置のうちの第一の物理的記憶装置群にまたがって設定され、前記データ処理装置が直接アクセスを行う論理的記憶装置を有し、前記記憶制御装置は、予め定めた指標に基づいて、前記第一の物理的記憶装置群とは異なる第二の物理的記憶装置群を選択すると共に、前記第二の物理的記憶装置群に、前記論理的記憶装置に格納されたデータ全体を連続的に格納する論理的記憶装置再配置手段を有することを特徴とする記憶装置システムを提供する。

上記第8の観点による記憶装置システムでは、再配置により、アクセス性能を向上することが出来る。

【0017】

第9の観点では、本発明は、上記構成の記憶装置システムにおいて、再配置中の論理的記憶装置にデータ処理装置からのアクセスがあったとき、再配置中の論理的記憶装置の再配置完了領域と再配置未完領域とを識別し、前記アクセス位置が前記再配置完了領域ならば再配置先の物理的記憶装置にアクセスさせ、前記アクセス位置が前記再配置未完領域ならば再配置前の物理的記憶装置にアクセスさせるアクセス位置切替手段をさらに具備したこ

10

20

30

40

50

とを特徴とする記憶装置システムを提供する。

上記第9の観点による記憶装置システムでは、再配置中の論理的記憶装置の再配置完了領域と再配置未完領域とを識別し、データ処理装置からのアクセス位置を切り替えるから、データ処理装置と物理的記憶装置の間のデータ転送を運用中に再配置を行うことが出来る。

第10の観点では、本発明は、データ処理装置が直接アクセスを行う論理的記憶装置が配置された物理的記憶装置を含む複数の物理的記憶装置と、前記データ処理装置と前記複数の物理的記憶装置の間のデータ転送を制御する記憶制御装置とを有する記憶装置システムにおいて、前記記憶制御装置は、前記データ処理装置によるアクセス状況を取得し、前記アクセス状況に基づいて前記論理的記憶装置のデータを、前記複数の物理的記憶装置のうち、該論理的記憶装置に配置されている物理的記憶装置から、該論理的記憶装置が配置された物理的記憶装置とは異なる物理的記憶装置に移動させることを特徴とする記憶装置システムを提供する。

10

上記第10の観点による記憶装置システムでは、データ処理装置によるアクセス状況に応じて論理的記憶装置のデータを、前記複数の物理的記憶装置のうち、該論理的記憶装置に配置されている物理的記憶装置から、該論理的記憶装置が配置された物理的記憶装置とは異なる物理的記憶装置へと移動させるから、アクセス性能を向上することが出来る。

第11の観点では、本発明は、キャッシュメモリを有し、前記記憶制御装置は、前記第一の物理的記憶装置群から前記キャッシュメモリへデータを読み出す手段と、前記読み出されたデータを前記第二の物理的記憶装置群に書き込む手段とを有することを特徴とする上記第8の観点の記憶装置システムを提供する。

20

第12の観点では、本発明は、前記第二の物理的記憶装置群は、前記第一の物理的記憶装置群よりもデータの読み出し速度が高速であることを特徴とする上記第8又は第11の観点の記憶装置システムを提供する。

【0018】

【発明の実施の形態】

以下、本発明の実施形態を説明する。なお、これにより本発明が限定されるものではない。

【0019】

- 第1の実施形態 -

30

第1の実施形態は、各論理ディスク装置のアクセス情報を記憶制御装置で採取し、SVP（サービスプロセッサ）を通じて保守員に提示し、このアクセス情報に基づく保守員の再配置指示により、論理ディスク装置の物理ディスク装置への再配置を行うものである。

【0020】

図1は、本発明の第1の実施形態にかかる記憶制御装置を含む情報処理システムのブロック図である。

この情報処理システム1は、データ処理装置100と、記憶制御装置104と、1台以上の物理ディスク装置105と、SVP111とを接続してなっている。

【0021】

前記データ処理装置100は、CPU101と、主記憶102と、チャンネル103とを有している。

40

【0022】

前記記憶制御装置104は、1つ以上のディレクタ106と、キャッシュメモリ107と、ディレクトリ108と、不揮発性メモリ109と、不揮発性メモリ管理情報110と、論理物理対応情報300と、論理ディスク装置情報400と、アクセス情報500を有している。

前記ディレクタ106は、データ処理装置100のチャンネル103と物理ディスク装置105の間のデータ転送、データ処理装置100のチャンネル103と前記キャッシュメモリ107の間のデータ転送および前記キャッシュメモリ107と物理ディスク装置105の間のデータ転送を行う。

50

前記キャッシュメモリ 107 には、物理ディスク装置 105 の中のアクセス頻度の高いデータをロードしておく。このロード処理は、前記ディレクタ 106 が実行する。ロードするデータ的具体例は、データ処理装置 100 の CPU 101 のアクセス対象データや、このアクセス対象データと物理ディスク装置 105 上の格納位置に近いデータ等である。

前記ディレクトリ 108 は、前記キャッシュメモリ 107 の管理情報を格納する。

前記不揮発性メモリ 109 は、前記キャッシュメモリ 107 と同様に、物理ディスク装置 105 の中のアクセス頻度の高いデータをロードしておく。

前記不揮発性メモリ管理情報 110 は、前記不揮発性メモリ 109 の管理情報を格納する。

前記論理物理対応情報 300 は、各論理ディスク装置 (図 2 の 200) が配置されている物理ディスク装置 105 上の位置および各物理ディスク装置 105 に配置されている論理ディスク装置 (図 2 の 200) を示す情報である。この情報を用いて、データ処理装置 100 の CPU 101 のアクセス対象データの物理ディスク装置 105 上の格納領域の算出などを行う。 10

前記論理ディスク装置情報 400 は、各論理ディスク装置 (図 2 の 200) のアクセス可否等の状態を示す。

前記アクセス情報 500 は、各論理ディスク装置 (図 2 の 200) のアクセス頻度やアクセスパターンなどの情報である。

#### 【0023】

論理物理対応情報 300 と論理ディスク情報 400 は、電源断などによる消失を防ぐために不揮発の媒体に記録する。 20

#### 【0024】

前記物理ディスク装置 105 は、データを記録する媒体と、記録されたデータを読み書きする装置とから構成される。

#### 【0025】

前記 SVP 111 は、アクセス情報 500 の保守員への提示や、保守員からの再配置指示 620 の入力を受け付けを行う。また、保守員からの情報処理システム 1 への指示の発信や、情報処理システム 1 の障害状態等の保守員への提示を行う。

#### 【0026】

図 2 は、論理ディスク装置 200 と物理ディスク装置 105 の関連を表わした図である。 30  
論理ディスク装置 200 は、データ処理装置 100 の CPU 101 が直接アクセスする見掛け上のディスク装置で、アクセス対象データが実際に格納される物理ディスク装置 105 と対応している。論理ディスク装置 200 上のデータは、シーケンシャルアクセスを考慮して、物理ディスク装置 105 上に連続的に配置されている。論理ディスク装置 200 のデータが配置されている物理ディスク装置 105 がディスクアレイ構成の場合、該論理ディスク装置 200 は複数の物理ディスク装置 105 と対応する。また、物理ディスク装置 105 の容量が論理ディスク装置 200 より大きく、複数の論理ディスク装置のデータを 1 台の物理ディスク装置 105 に格納できる場合には、該物理ディスク装置 105 は複数の論理ディスク装置 200 と対応する。この論理ディスク装置 200 と物理ディスク装置 105 の対応は前記論理物理対応情報 300 で管理される。例えば、データ処理装置 100 の CPU 101 が論理ディスク装置 200 のデータ 201 をリードする時、記憶制御装置 104 で論理物理対応情報 300 に基づき論理ディスク装置 200 に対応する物理ディスク装置 105 を求め、その物理ディスク装置 105 の領域内のデータ格納位置 202 を求め、データ転送を行う。 40

#### 【0027】

図 3 は、論理物理対応情報 300 を表わした図である。

論理物理対応情報 300 は、論理ディスク構成情報 310 と、物理ディスク構成情報 320 とから構成される。前記論理ディスク構成情報 310 は、各論理ディスク装置 200 が配置されている物理ディスク装置 105 上の領域に関する情報であり、論理ディスク装置 200 から対応する物理ディスク装置 105 を求める時に用いる。一方、前記物理ディス 50

ク構成情報 320 は、各物理ディスク装置 105 に配置されている論理ディスク装置 200 に関する情報で、物理ディスク装置 105 から対応する論理ディスク装置 200 を求める時に用いる。

#### 【0028】

前記論理ディスク構成情報 310 は、物理ディスク装置グループ 311、RAID 構成 312 および開始位置 313 の組を、論理ディスク装置 200 の数だけ有している。

前記物理ディスク装置グループ 311 は、当該論理ディスク装置 200 が配置されている物理ディスク装置 105 を示す情報である。

前記 RAID 構成 312 は、前記物理ディスク装置グループ 311 の RAID レベルを示す。

10

前記開始位置 313 は、当該論理ディスク装置 200 が物理ディスク装置 105 上で配置されている先頭位置を示す。

#### 【0029】

前記物理ディスク構成情報 320 は、論理ディスク装置グループ 321 を、物理ディスク装置 105 の数だけ有している。

前記論理ディスク装置グループ 321 は、当該物理ディスク装置 105 に配置されている論理ディスク装置 200 を示す。

#### 【0030】

図 4 は、論理ディスク情報 400 を表わした図である。

論理ディスク情報 400 は、論理ディスク状態 401 と再配置完了ポインタ 402 とを、論理ディスク装置 200 の数だけ有している。

20

前記論理ディスク状態 401 は、「正常」「閉塞」「フォーマット中」「再配置中」などの論理ディスク装置 200 の状態を表わす。

前記再配置完了ポインタ 402 は、前記論理ディスク状態 401 が「再配置中」の時のみ有効な情報で、当該論理ディスク装置 200 の再配置処理を完了している領域の次の位置すなわち当該論理ディスク装置 200 が未だ再配置処理を終えていない領域の先頭位置を示す。「再配置中」におけるデータアクセス時、再配置完了ポインタ 402 よりも前の領域へのアクセスの場合には、再配置後の物理ディスク装置 105 へアクセスしなければならない。一方、再配置完了ポインタ 402 以後の領域へのアクセスの場合には、再配置前の物理ディスク装置 105 へアクセスしなければならない。

30

#### 【0031】

図 5 は、アクセス情報 500 を表わしている。

アクセス情報 500 は、アクセス頻度情報 501 とアクセスパターン情報 502 とを、論理ディスク装置 200 の数だけ有している。このアクセス情報 500 は、記憶制御装置 104、データ処理装置 100、SVP 111 のいずれからも参照することが出来る。

前記アクセス頻度情報 501 は、単位時間あたりの当該論理ディスク装置 200 へのアクセス回数を管理する。このアクセス頻度情報 501 は、各論理ディスク装置 200 の中でアクセス頻度の高いもの又は低いものを求める指標として用いる。

前記アクセスパターン情報 502 は、当該論理ディスク装置 200 へのシーケンシャルアクセスとランダムアクセスの割合を管理する。このアクセスパターン情報 502 は、シーケンシャルアクセスが多く、よりシーケンシャル性能の高い物理ディスク装置 105 に再配置するのが望ましい論理ディスク装置 200 を求める指標として用いる。

40

#### 【0032】

次に、記憶制御装置 104 の動作を説明する。

図 6 は、記憶制御装置 104 の動作を詳細に表わした図である。

まず、リード/ライト処理時の動作について説明する。

ディレクタ 106 は、通常リード/ライト処理を実行する際、CPU 101 からチャンネル 103 を経由して CPU からの指示 600 を受け取る。この CPU からの指示 600 は、リード(またはライト)対象のレコードが記憶されている論理ディスク装置 200 を指定する指定情報 1 と、リード(またはライト)対象のレコードが記憶されている論理ディス

50



ク装置 200 内の位置 (トラック, セクタ, レコード) を指定する指定情報 2 とを含んでいる。

ディレクタ 106 は、物理ディスク装置上のアクセス位置算出処理 (610) で、前記 CPU からの指示 600 と論理物理対応情報 300 とを用いて、物理ディスク装置 105 上でのアクセス位置を算出する。この物理ディスク装置アクセス位置算出処理 (610) については図 8 を参照して後で詳述する。

その後、たとえばリード処理では、算出した物理ディスク装置 105 上のデータ格納位置 202 のデータをキャッシュメモリ 107 上に読み上げてデータ 201 とし、その読み上げたデータ 201 をチャンネル 103 を通じて主記憶 102 に転送する。

#### 【0033】

次に、アクセス情報 500 の採取処理について説明する。

CPU 101 からのリード/ライト処理のアクセス時に、ディレクタ 106 は、アクセス対象論理ディスク装置 200 のアクセス情報 500 を更新する。

アクセス頻度情報 501 の採取は、例えば、アクセスの度に内部カウンタをカウントアップしていき、一定時間または一定回数のアクセス経過後のアクセス時に、前記内部カウンタからアクセス頻度を判定する。

アクセスパターン情報 502 の採取は、例えば、アクセスの度に内部カウンタにシーケンシャルアクセス回数をカウントアップしていき、一定時間または一定回数のアクセス経過後のアクセス時に、前記内部カウンタからアクセスパターンを判定する。

#### 【0034】

次に、再配置指示 620 を説明する。

保守員は、SVP 111 を通じて提示されたアクセス情報 500 を参照して、各論理ディスク装置 200 の再配置の必要性を検討する。この検討の結果、再配置を決定した論理ディスク装置 200 があれば、SVP 111 を通じて記憶制御装置 104 に対して再配置指示 620 を出す。

この再配置指示 620 は、再配置対象の論理ディスク装置 200 を 2 つ指定する指示情報 1-2 からなる。

保守員が行う検討の内容は、後述する第 3 の実施形態で図 10 を参照して説明する論理ディスク装置再配置要否決定処理 (910) と同様である。

#### 【0035】

次に、論理ディスク装置再配置処理 (630) を説明する。

ディレクタ 106 は、前記再配置指示 620 を受けて、指定された 2 つの論理ディスク装置 200 の間で論理ディスク装置再配置処理 (630) を行う。

図 7 は、論理ディスク装置再配置処理部 630 の処理フロー図である。

ステップ 700 では、論理ディスク情報 400 のうちの指定された 2 つの論理ディスク装置 200 の論理ディスク状態 401 を「再配置中」に設定する。

ステップ 701 では、論理ディスク情報 400 のうちの指定された 2 つの論理ディスク装置 200 の再配置完了ポインタ 402 を各論理ディスク装置 200 の先頭位置に初期化する。

ステップ 702 では、論理ディスク情報 400 のうちの指定された 2 つの論理ディスク装置 200 の再配置完了ポインタ 402 をチェックし、全領域の再配置が完了していなければステップ 703 へ進み、完了していればステップ 707 へ進む。

#### 【0036】

ステップ 703 では、再配置完了ポインタ 402 が示すデータ位置から再配置処理の 1 回の処理単位分のデータに対して物理ディスク装置 105 からキャッシュメモリ 107 上へのデータ転送を行う。ここで、1 回の処理単位分のデータ量は、再配置対象の 2 つの論理ディスク装置 200 の冗長データ 1 つに対応する各データ量の最小公倍数に決定される。たとえば、再配置を RAID 5 の論理ディスク装置 200 と RAID 1 の論理ディスク装置 200 の間で行うならば、RAID 1 の論理ディスク装置 200 の冗長データ 1 つに対応するデータ量は“1”であるから、1 回の処理単位分のデータ量は、RAID 5 の論理

10

20

30

40

50

ディスク装置 200 の冗長データ 1 つに対応するデータ量すなわちパリティ 1 つに対応するデータ量に決定される。

【0037】

ステップ 704 では、再配置対象の各論理ディスク装置 200 の再配置先論理ディスク装置 200 がパリティを有する RAID レベルのものである場合、キャッシュメモリ 107 上の再配置対象の 1 回の処理単位分のデータ 201 に対してパリティを生成する。

ステップ 705 では、キャッシュメモリ 107 上の再配置対象の 1 回の処理単位分のデータ 201 および前記ステップ 704 で作成したパリティを、再配置先の物理ディスク装置 105 へ書き込む。

ステップ 706 では、1 回の処理単位分だけ再配置完了ポインタ 402 を進める。そして、前記ステップ 702 に戻る。 10

【0038】

なお、上記ステップ 703, 704 において、データおよびパリティは、不揮発性メモリ 109 にも転送して二重化し、キャッシュ障害によるデータ消失を防ぐ。この理由は、上記ステップ 705 での書き込み時に、例えば、第 1 の論理ディスク装置 200 と第 2 の論理ディスク装置 200 のデータのうち、第 1 の論理ディスク装置 200 のデータを物理ディスク装置 105 (元は第 2 の論理ディスク装置 200 に配置されていた物理ディスク装置 105) へ書き込んだ段階で障害によりキャッシュメモリ 107 上のデータがアクセス不能になったとすると、書き込みが終了していない第 2 の論理ディスク装置 200 のデータが消失するからである (元は第 2 の論理ディスク装置 200 に配置されていた物理ディスク装置 105 には、上記のように第 1 の論理ディスク装置 200 のデータが上書きされてしまっている)。 20

【0039】

ステップ 707 では、論理物理対応情報 300 を更新する。すなわち、論理ディスク構成情報 310 と物理ディスク構成情報 321 を変更する。

ステップ 708 では、論理ディスク情報 400 の論理ディスク状態 401 を元の状態に戻し、再配置処理 (630) を終了する。

【0040】

次に、物理ディスク装置アクセス位置算出処理 (610) を説明する。

図 8 は、物理ディスク装置アクセス位置算出処理部 610 の処理フロー図である。 30

ステップ 800 では、論理ディスク情報 400 のうちのアクセス対象論理ディスク装置 200 の論理ディスク状態 401 が「再配置中」であるか否かをチェックし、「再配置中」ならばステップ 801 に進み、「再配置中で」なければステップ 803 に進む。

【0041】

ステップ 801 では、論理ディスク情報 400 のうちのアクセス対象論理ディスク装置 200 の再配置完了ポインタ 402 とアクセスデータ位置とを比較し、アクセスデータ位置が再配置完了ポインタ 402 の指す位置以後ならばステップ 802 に進み、アクセスデータ位置が再配置完了ポインタ 402 の指す位置より前ならばステップ 803 に進む。

【0042】

ステップ 802 では、当該論理ディスク装置 200 の再配置先の論理ディスク装置 200 をアクセス対象にする。そして、ステップ 804 へ進む。 40

【0043】

ステップ 803 では、当該論理ディスク装置 200 をアクセス対象とする。

【0044】

ステップ 804 では、アクセス対象の論理ディスク装置 200 に対応した物理ディスク装置 105 上でのアクセス位置を、論理物理対応情報 300 を用いて算出する。

【0045】

以上の第 1 の実施形態にかかる情報処理システム 1 および記憶制御装置 104 によれば、アクセス情報 500 に基づく保守員の判断により、アクセス頻度の高い論理ディスク装置をより高速な物理ディスク装置へ再配置することが出来る。また、シーケンシャルアクセ 50

スの比率の高い論理ディスク装置をよりシーケンシャルアクセス性能の高い物理ディスク装置へ再配置することが出来る。従って、アクセス性能を向上することが出来る。

【 0 0 4 6 】

- 第 2 の実施形態 -

上記第 1 の実施形態を変形して、記憶制御装置 1 0 4 からアクセス情報 5 0 0 をデータ処理装置 1 0 0 に提示し、データ処理装置 1 0 0 が再配置要否を決定し記憶制御装置 1 0 4 に再配置指示 ( 6 2 0 相当 ) を出すようにしてもよい。

【 0 0 4 7 】

- 第 3 の実施形態 -

第 3 の実施形態は、再配置指示を S V P 1 1 1 やデータ処理装置 1 0 0 から受けるのではなく、記憶制御装置 1 0 4 が自己決定するものである。 10

【 0 0 4 8 】

図 9 は、記憶制御装置 1 0 4 の動作を詳細に表わした図である。

第 1 の実施形態 ( 図 6 ) との違いは、論理ディスク再配置要否決定処理部 9 1 0 が再配置指示 6 2 0 を出すことである。

【 0 0 4 9 】

図 1 0 は、上記論理ディスク再配置要否決定処理部 9 1 0 の処理フロー図である。

この論理ディスク再配置要否決定処理 ( 9 1 0 ) は、ディレクタ 1 0 6 が一定周期で各論理ディスク装置 2 0 0 のアクセス情報 5 0 0 を検査して行う。

ステップ 1 0 0 0 では、アクセス情報 5 0 0 のアクセス頻度情報 5 0 1 を参照し、アクセス頻度が規定値を超え且つ配置されている物理ディスク装置 1 0 5 が比較的低速なものである論理ディスク装置 ( 以下、これを第 1 候補論理ディスク装置という ) 2 0 0 があるか否かをチェックし、該当する論理ディスク装置 2 0 0 があればステップ 1 0 0 1 へ進み、なければステップ 1 0 0 5 へ進む。 20

【 0 0 5 0 】

ステップ 1 0 0 1 では、前記第 1 候補論理ディスク装置 2 0 0 のアクセスパターン情報 5 0 2 を参照し、シーケンシャルアクセスの比率が規定値以上であるか否かをチェックし、規定値以上でなければステップ 1 0 0 2 へ進み、規定値以上であればステップ 1 0 0 4 へ進む。

【 0 0 5 1 】

ステップ 1 0 0 2 では、前記第 1 候補論理ディスク装置 2 0 0 より高速な物理ディスク装置 1 0 5 に配置されている論理ディスク装置 2 0 0 のアクセス頻度情報 5 0 1 を参照し、アクセス頻度が規定値以下の論理ディスク装置 ( 以下、これを第 2 候補論理ディスク装置という ) 2 0 0 があるか否かをチェックし、あればステップ 1 0 0 3 へ進み、なければステップ 1 0 0 5 へ進む。 30

【 0 0 5 2 】

ステップ 1 0 0 3 では、前記第 1 候補論理ディスク装置 2 0 0 と前記第 2 候補論理ディスク装置 2 0 0 の間で再配置処理 ( 6 3 0 ) が必要であると決定し、再配置指示 6 2 0 を出す。そして、処理を終了する。

【 0 0 5 3 】

ステップ 1 0 0 4 では、前記第 1 候補論理ディスク装置 2 0 0 よりシーケンシャル性能の高い物理ディスク装置 1 0 5 に配置されている論理ディスク装置 2 0 0 のアクセスパターン情報 5 0 2 を参照し、シーケンシャルアクセスの比率が規定値以下の論理ディスク装置 ( 以下、これを第 2 候補論理ディスク装置という ) 2 0 0 があるか否かをチェックし、あれば前記ステップ 1 0 0 3 へ進み、なければ前記ステップ 1 0 0 2 へ進む。 40

【 0 0 5 4 】

ステップ 1 0 0 5 では、論理ディスク装置 2 0 0 の再配置処理 ( 6 3 0 ) は不要であると決定する。そして、処理を終了する。

【 0 0 5 5 】

以上の第 3 の実施形態にかかる情報処理システム 1 および記憶制御装置 1 0 4 によれば、 50

アクセス情報 500 に基づいて自動的に、アクセス頻度の高い論理ディスク装置をより高速な物理ディスク装置へ再配置することが出来る。また、シーケンシャルアクセスの比率の高い論理ディスク装置をよりシーケンシャルアクセス性能の高い物理ディスク装置へ再配置することが出来る。従って、アクセス性能を向上することが出来る。

【0056】

- 第4の実施形態 -

上記第1～第3の実施形態を变形して、アクセス情報 500 に代えて又は加えて、論理ディスク装置 200 に要求される信頼性を再配置処理要否決定の指標に用いてもよい。信頼性を指標に用いれば、論理ディスク装置 200 上のデータの信頼性を向上させることが出来る。

10

【0057】

【発明の効果】

本発明の記憶装置システムによれば、シーケンシャルアクセスの場合やランダムアクセスでヒット率が低い場合でも、アクセス性能を向上することが出来る。また、本発明の記憶装置システムによれば、データの信頼性を向上することが出来る。

【図面の簡単な説明】

【図1】本発明の第1の実施形態にかかる記憶制御装置を含む情報処理システムのブロック図である。

【図2】論理ディスク装置と物理ディスク装置との対応関係の説明図である。

【図3】論理物理対応情報の構成例示図である。

20

【図4】論理ディスク情報の構成例示図である。

【図5】アクセス情報の構成例示図である。

【図6】本発明の第1の実施形態における記憶制御装置の動作を示すブロック図である。

【図7】論理ディスク装置再配置処理部の処理フロー図である。

【図8】物理ディスク装置アクセス位置算出処理部の処理フロー図である。

【図9】本発明の第3の実施形態における記憶制御装置の動作を示すブロック図である。

【図10】論理ディスク装置再配置要否決定処理部の処理フロー図である。

【符号の説明】

1 ... 情報処理システム

100 ... データ処理装置

30

101 ... CPU

102 ... 主記憶

103 ... チャンネル

104 ... 記憶制御装置

105 ... 物理ディスク装置

106 ... ディレクタ

107 ... キャッシュメモリ

108 ... キャッシュディレクトリ

109 ... 不揮発性メモリ

110 ... 不揮発性メモリ管理情報

40

111 ... SVP

200 ... 論理ディスク装置

201 ... データ

202 ... データ格納位置

300 ... 論理物理対応情報

400 ... 論理ディスク情報

500 ... アクセス情報

600 ... CPUからの指示

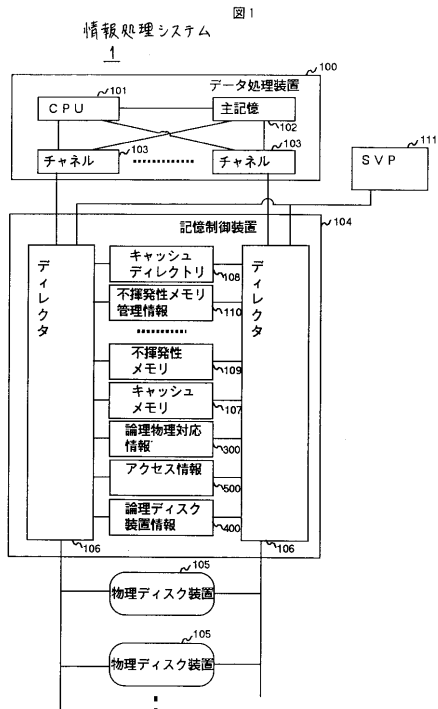
610 ... 物理ディスク装置上のアクセス位置算出処理部

620 ... 指示情報

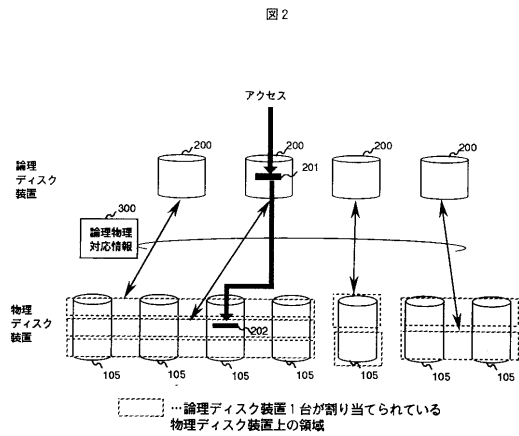
50

630 ... 論理ディスク装置再配置処理部  
 910 ... 論理ディスク再配置要否決定処理部

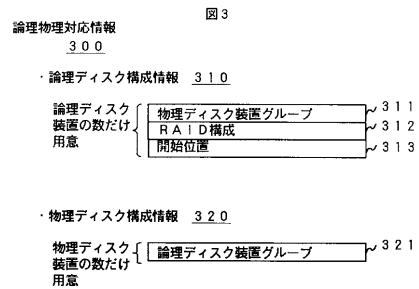
【 図 1 】



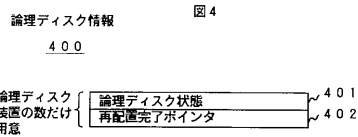
【 図 2 】



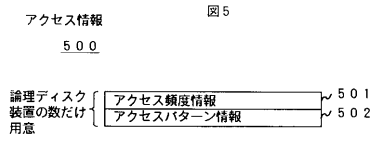
【 図 3 】



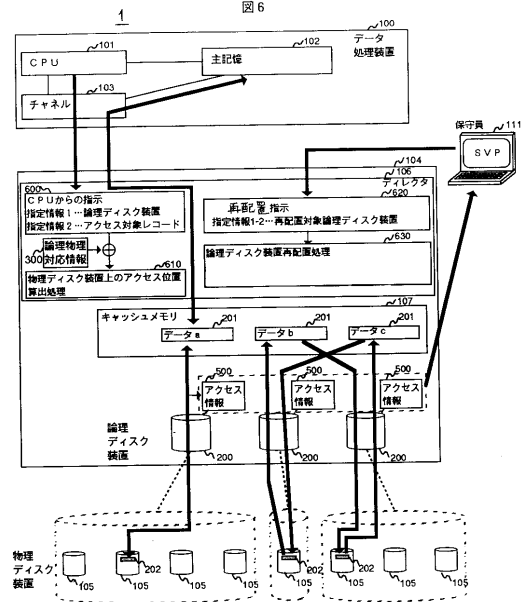
【 図 4 】



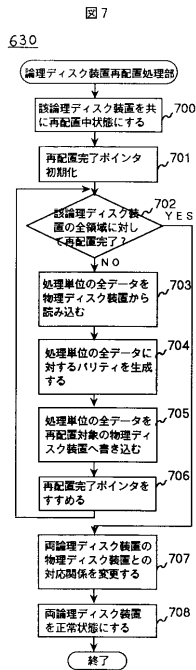
【 図 5 】



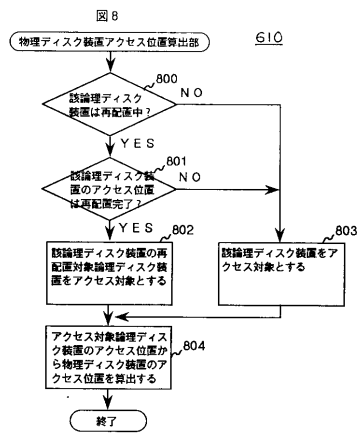
【 図 6 】



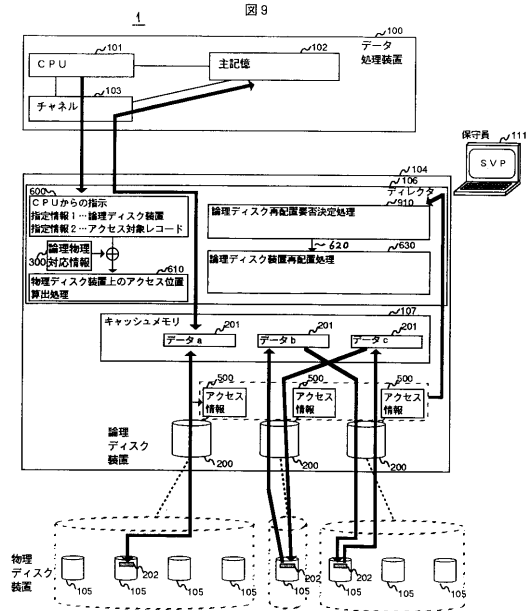
【 図 7 】



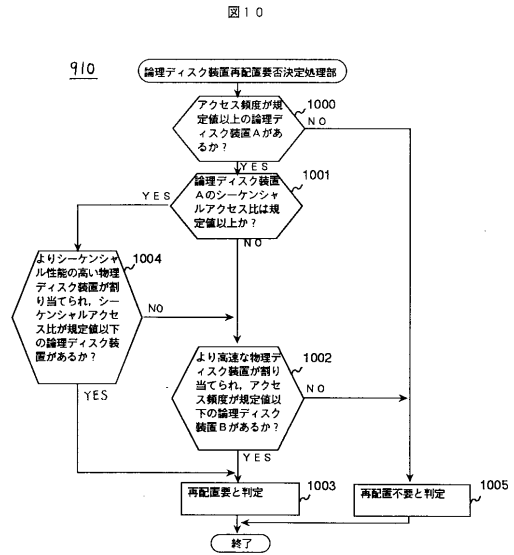
【 図 8 】



【 図 9 】



【 図 10 】



---

フロントページの続き

(72)発明者 佐藤 孝夫  
神奈川県小田原市国府津2880番地  
事業部内  
株式会社日立製作所 ストレージシステム

合議体

審判長 大日方 和幸

審判官 大野 克人

審判官 矢島 伸一

(56)参考文献 特開平5 - 100801 (JP, A)  
特開平6 - 139027 (JP, A)  
特開平7 - 146757 (JP, A)  
特開平3 - 102418 (JP, A)  
特開平8 - 63298 (JP, A)  
特開昭62 - 67629 (JP, A)  
特開平7 - 141121 (JP, A)  
月刊アスキー 第16巻9号通巻183号p166 (株式会社アスキー)