

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5438827号
(P5438827)

(45) 発行日 平成26年3月12日(2014.3.12)

(24) 登録日 平成25年12月20日(2013.12.20)

(51) Int. Cl.	F I
G06F 3/06 (2006.01)	G06F 3/06 301J
	G06F 3/06 301E
	G06F 3/06 301Z
	G06F 3/06 540
	G06F 3/06 301M

請求項の数 13 (全 40 頁)

(21) 出願番号	特願2012-516979 (P2012-516979)	(73) 特許権者	000005108
(86) (22) 出願日	平成21年10月9日(2009.10.9)		株式会社日立製作所
(65) 公表番号	特表2012-531653 (P2012-531653A)		東京都千代田区丸の内一丁目6番6号
(43) 公表日	平成24年12月10日(2012.12.10)	(74) 代理人	110000279
(86) 国際出願番号	PCT/JP2009/005297		特許業務法人ウィルフォート国際特許事務所
(87) 国際公開番号	W02011/042940	(72) 発明者	川口 裕太郎
(87) 国際公開日	平成23年4月14日(2011.4.14)		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23年12月27日(2011.12.27)	(72) 発明者	石川 篤
			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

最終頁に続く

(54) 【発明の名称】 記憶制御装置及び仮想ボリュームの制御方法

(57) 【特許請求の範囲】

【請求項1】

仮想的に形成される仮想ボリュームと、
一つまたは複数の記憶装置を含む、複数のRAIDグループと、
前記各RAIDグループにそれぞれストライプ状に設けられる第1実記憶領域であって、それぞれ複数の第2実記憶領域を有する複数の第1実記憶領域を管理するためのプール部と、

上位装置から前記仮想ボリュームに関するライトコマンドが発行された場合、前記各第1実記憶領域のうち所定の第1実記憶領域に含まれる前記各第2実記憶領域の中から所定の第2の実記憶領域を選択し、その所定の第2実記憶領域を前記ライトコマンドに対応する前記仮想ボリューム内の所定領域に対応付ける制御部であって、一つの前記第1実記憶領域に一つの前記仮想ボリュームを対応付ける制御部と、

前記仮想ボリュームに対応付けられている前記各第2実記憶領域の中から移動対象の第2実記憶領域を予め選択し、その移動対象の第2実記憶領域に記憶されているデータの移動先となる移動先第1実記憶領域を、前記各第1実記憶領域のうち前記移動対象の第2実記憶領域が設けられているRAIDグループ以外の他のRAIDグループ内の各第1実記憶領域の中から予め選択する移動先決定部と、

予め選択した前記移動対象の第2実記憶領域および前記移動先第1実記憶領域を対応付けて記憶する移動先記憶部と、

前記上位装置が前記移動対象として予め選択した第2実記憶領域に対応するコマンドを

発行した場合は、前記コマンドの処理の中で、前記移動対象の第2実記憶領域に記憶されているデータを、前記移動先記憶部により記憶されている前記移動先第1実記憶領域内に移動させるコマンド処理部と、
を備える記憶制御装置。

【請求項2】

(1) 前記移動先記憶部は、
 (1-1) 前記移動対象の第2実記憶領域を特定する移動対象情報と、
 (1-2) 前記移動先第1実記憶領域を特定する移動先情報と、
 (1-3) 前記各RAIDグループの負荷に関する負荷情報と、
 (1-4) 前記各RAIDグループの使用容量に関する使用容量情報と、
 を記憶しており、
 (2) 移動先決定部は、
 (2-1) 所定時刻が到来した場合またはユーザから指示された場合のいずれかの場合に、前記移動先記憶部に記憶されている前記移動先情報を消去させ、さらに、
 (2-2) 前記各第2実記憶領域のうちヌルデータのみが記憶されている第2実記憶領域を解放して、未使用の第2実記憶領域に変更させ、さらに、
 (2-3) 前記負荷情報に基づく負荷分散処理であって、相対的に高負荷の第2実記憶領域を相対的に低負荷の第1実記憶領域内に移動させるための、新たな移動先情報を作成して、前記移動先記憶部に記憶させる負荷分散処理と、
 (2-4) 前記使用容量情報に基づく使用容量平均化処理であって、相対的に使用容量の大きい第1実記憶領域内の第2実記憶領域を相対的に使用容量の小さい第1実記憶領域内に移動させるための、新たな移動先情報を作成して、前記移動先記憶部に記憶させる使用容量平均化処理とを、
 それぞれ実行し、
 (3) 前記コマンド処理部は、
 (3-1) 前記上位装置から前記移動対象の第2実記憶領域についてのライトコマンドが発行された場合、前記移動対象の第2実記憶領域に記憶されているデータを読み出し、その読み出されたデータと前記ライトコマンドに係るライトデータとをマージして、前記移動先第1実記憶領域内の第2実記憶領域に書き込み、さらに、前記上位装置に前記ライトコマンドの処理が完了した旨を通知し、
 (3-2) 前記上位装置から前記移動対象の第2実記憶領域についてのリードコマンドが発行された場合、前記移動対象の第2実記憶領域からデータを読み出して前記上位装置に送信し、前記上位装置に前記リードコマンドの処理が完了した旨を通知した後で、前記移動対象の第2実記憶領域から読み出されたデータを前記移動先第1実記憶領域内の第2実記憶領域に書き込む、
 請求項1に記載の記憶制御装置。

【請求項3】

前記コマンド処理部は、前記上位装置から前記移動対象の第2実記憶領域についてのライトコマンドが発行された場合、前記移動対象の第2実記憶領域に記憶されているデータを読み出し、その読み出されたデータと前記ライトコマンドに係るライトデータとをマージして、前記移動先第1実記憶領域内の第2実記憶領域に書き込む、
 請求項1に記載の記憶制御装置。

【請求項4】

前記コマンド処理部は、前記上位装置から前記移動対象の第2実記憶領域についてのリードコマンドが発行された場合、前記移動対象の第2実記憶領域からデータを読み出して前記上位装置に送信し、前記上位装置に前記リードコマンドの処理が完了した旨を通知した後で、前記移動対象の第2実記憶領域から読み出されたデータを前記移動先第1実記憶

領域内の第 2 実記憶領域に書き込む、
請求項 3 に記載の記憶制御装置。

【請求項 5】

前記移動先記憶部は、
前記移動対象の第 2 実記憶領域を特定する移動対象情報と、
前記移動先第 1 実記憶領域を特定する移動先情報と、
前記各 R A I D グループの負荷に関する負荷情報と、
を記憶しており、
移動先決定部は、
前記移動先記憶部に記憶されている前記移動先情報を消去させ、
相対的に高負荷の第 2 実記憶領域を相対的に低負荷の第 1 実記憶領域内に移動させる
ための、新たな移動先情報を前記負荷情報に基づいて作成し、前記移動先記憶部に記憶さ
せる、
請求項 3 に記載の記憶制御装置。

10

【請求項 6】

前記移動先記憶部は、
前記移動対象の第 2 実記憶領域を特定する移動対象情報と、
前記移動先第 1 実記憶領域を特定する移動先情報と、
前記各 R A I D グループの負荷に関する負荷情報と、
前記各 R A I D グループの使用容量に関する使用容量情報と、
を記憶しており、
移動先決定部は、
前記移動先記憶部に記憶されている前記移動先情報を消去させ、
相対的に高負荷の第 2 実記憶領域を相対的に低負荷の第 1 実記憶領域内に移動させる
ための、新たな移動先情報を前記負荷情報に基づいて作成し、前記移動先記憶部に記憶さ
せ、さらに、
相対的に使用容量の大きい第 1 実記憶領域内の第 2 実記憶領域を相対的に使用容量の
小さい第 1 実記憶領域内に移動させるための、新たな移動先情報を前記使用容量情報に基
づいて作成し、前記移動先記憶部に記憶させる、
請求項 3 に記載の記憶制御装置。

20

30

【請求項 7】

前記移動先決定部は、前記移動先記憶部に記憶されている前記移動先情報を消去させた
後で、前記各第 2 実記憶領域のうちヌルデータのみが記憶されている第 2 実記憶領域を解
放して、未使用の第 2 実記憶領域に変更させる、
請求項 5 に記載の記憶制御装置。

【請求項 8】

前記移動先決定部は、前記移動対象の第 2 実記憶領域を前記第 1 実記憶領域の単位で複
数選択する、請求項 1 に記載の記憶制御装置。

40

【請求項 9】

前記制御部は、
前記仮想ボリュームを生成する場合に、前記仮想ボリューム内の各仮想的記憶領域を
、初期データが記憶されている初期化用の第 2 実記憶領域に対応付け、
前記上位装置から前記仮想ボリュームに関する前記ライトコマンドが発行された場合
、前記ライトコマンドに対応する前記仮想的記憶領域の対応付け先を、前記初期化用の第
2 実記憶領域から、選択された前記所定の第 2 実記憶領域に切り替える、

50

請求項 1 に記載の記憶制御装置。

【請求項 1 0】

前記制御部は、前回のライト要求に対応して前記仮想ボリュームに対応付けられた前記第 2 実記憶領域に連続する、未使用の前記第 2 実記憶領域を、前記所定の第 2 実記憶領域として前記仮想ボリュームに対応付ける、
請求項 9 に記載の記憶制御装置。

【請求項 1 1】

前記仮想ボリュームに記憶されるデータには、所定サイズ毎に保証コードが設定されており、前記保証コードは、前記 R A I D グループを識別するためのデータと、前記第 1 実記憶領域を識別するためのデータと、前記第 1 実記憶領域内における前記第 2 実記憶領域を識別するためのデータとを含んでいる、
請求項 1 0 に記載の記憶制御装置。

10

【請求項 1 2】

ライトコマンドに応じて実記憶領域が割り当てられる仮想ボリュームを制御するための方法であって、

前記仮想ボリュームは、複数の仮想的記憶領域を備えており、

複数の R A I D グループを管理するためのプール部を作成し、

20

前記各 R A I D グループは、複数の記憶装置を跨るようしてストライプ状に形成される第 1 実記憶領域であって、前記仮想的記憶領域に対応する第 2 実記憶領域をそれぞれ複数ずつ有する第 1 実記憶領域を複数備えており、

前記上位装置から前記仮想ボリュームに関するライトコマンド要求が発行された場合、一つの前記第 1 実記憶領域に複数の前記仮想ボリュームが対応付けられないようにして、前記ライトコマンドに対応する前記仮想的記憶領域に、前記各第 1 実記憶領域のうち所定の第 1 実記憶領域に含まれる所定の第 2 実記憶領域を対応付け、

前記仮想的記憶領域に対応付けられる前記所定の第 2 実記憶領域に、前記上位装置から受領したライトデータを記憶させ、

前記仮想的記憶領域に対応付けられている前記各第 2 実記憶領域の中から、相対的に高負荷の第 2 実記憶領域を移動対象の第 2 実記憶領域として予め選択し、

30

前記各第 1 実記憶領域のうち前記移動対象の第 2 実記憶領域が設けられている R A I D グループ以外の他の R A I D グループ内の各第 1 実記憶領域であって、相対的に低負荷の第 1 実記憶領域を、前記移動対象の第 2 実記憶領域に記憶されているデータの移動先となる移動先第 1 実記憶領域として予め選択し、

予め選択した前記移動対象の第 2 実記憶領域および前記移動先第 1 実記憶領域を対応付けて移動先記憶部に記憶させ、

前記上位装置が前記移動対象として予め選択した第 2 実記憶領域に対応するライトコマンドを発行した場合は、前記移動対象の第 2 実記憶領域に記憶されているデータを読み出し、その読み出されたデータと前記ライトコマンドに係るライトデータとをマージして、
前記移動先第 1 実記憶領域内の第 2 実記憶領域に書き込む、
仮想ボリュームの制御方法。

40

【請求項 1 3】

前記上位装置から前記移動対象の第 2 実記憶領域についてのリードコマンドが発行された場合には、前記移動対象の第 2 実記憶領域からデータを読み出して前記上位装置に送信し、前記上位装置に前記リードコマンドの処理が完了した旨を通知した後で、前記移動対象の第 2 実記憶領域から読み出されたデータを前記移動先第 1 実記憶領域内の第 2 実記憶領域に書き込む、

請求項 1 2 に記載の仮想ボリュームの制御方法。

50

【発明の詳細な説明】**【技術分野】****【0001】**

本発明は、記憶制御装置及び仮想ボリュームの制御方法に関する。

【背景技術】**【0002】**

企業等のユーザは、記憶制御装置を用いてデータを管理する。記憶制御装置は、RAID (Redundant Array of Independent Disks) に基づく記憶領域上に、論理ボリュームを形成する。その論理ボリュームは、ホストコンピュータ(以下、ホスト)に提供される。

【0003】

ユーザの使用するデータ量は日々増大するため、現状に合わせて設定されたボリュームサイズでは、いずれ容量が不足する。これに対し、データ量の増加を見越して、ボリュームサイズを現在必要なサイズよりも過大に設定すると、不要不急のディスクドライブが多くなり、コストが増加する。

【0004】

そこで、仮想的な論理ボリュームを用意し、実際の使用に応じて、仮想的な論理ボリュームに実記憶領域を割り当てる技術が提案されている(特許文献1、特許文献2)。なお、仮想的な論理ボリュームに割り当てられた実記憶領域が特定の实ボリュームに偏在するのを防止するために、データを再配置させる技術も知られている(特許文献3)。

【先行技術文献】**【特許文献】****【0005】**

【特許文献1】 米国特許第6823442号明細書

【特許文献2】 特開2007-310861号公報

【特許文献3】 特開2008-234158号公報

【発明の概要】**【発明が解決しようとする課題】****【0006】**

前記文献(US6,823,442B1)には、ストレージサーバシステムが、仮想ボリューム上のブロックアドレスに関するライト要求を受信した場合に、そのブロックアドレスに対応する仮想ボリュームページアドレスに対して、論理的なデータページを割り当てる。そして、その論理的なデータページにデータが書き込まれる。

【0007】

前記文献には、複数の物理ディスク上のエリアから得られるチャンクレット(chunklet)という概念のエリアに基づいて、特定のRAIDレベルを有する論理ディスクを構成するための管理方法が記載されている。

【0008】

しかしながら、その管理方法は、物理ディスクドライブ単位にRAIDグループを構成する、記憶制御装置の物理エリア管理方法と全く異なる。従って、前記文献に記載の管理方法を、物理ディスクドライブ単位にRAIDグループを構成する記憶制御装置に、そのまま適用することはできない。

【0009】

仮に、前記文献に記載の技術を前記管理方法に適用する場合は、通常の論理ボリュームと仮想的な論理ボリュームとの両方を提供可能な記憶制御装置において、通常の論理ボリュームと仮想的な論理ボリュームとで、それぞれ物理エリアの管理方法が異なってしまう、記憶制御装置の構造が複雑化するという問題を生じる。ここで、通常の論理ボリュームとは、ボリューム生成時に、そのボリュームサイズと同容量の物理エリア(物理的記憶領域)が予め割り当てられる論理ボリュームを意味する。仮想的な論理ボリュームとは、ボリュームサイズが仮想化された論理ボリュームであって、ライト要求に応じて、物理エリアが割り当てられる論理ボリュームを意味する。

10

20

30

40

50

【0010】

つまり、前記文献に記載の技術を、物理ディスクドライブ単位でRAIDグループを構成する記憶制御装置にもしも適用したとすると、異なる複数の管理方法で物理エリアを管理しなければならず、構成が複雑化し、開発コストも増大する。

【0011】

さらに、前記文献では、ライト要求の受領時に、仮想ボリュームページアドレスに対応するテーブルページが割り当てられていない場合、ストレージサーバシステムは、まず最初にテーブルページを割当て、その次に、論理的なデータページを割り当てる。従って、前記文献に記載の技術では、テーブルページを割り当てた後で、データページを割り当てる必要があり、ライト処理の性能が低下するという問題がある。

10

【0012】

さらに、前記文献では、上述のような割当て処理を行うため、データページ専用のプールとテーブルページ専用のプールとをそれぞれ別々に設ける必要があり、システム構造が複雑化する。

【0013】

そこで、本発明の目的は、RAIDグループの物理的構成を考慮して仮想ボリュームに効率的に記憶領域を対応付けることのできる記憶制御装置及び仮想ボリュームの制御方法を提供することにある。本発明の他の目的は、各RAIDグループを均等に使用して仮想ボリュームを構成することのできる記憶制御装置及び仮想ボリュームの制御方法を提供することにある。本発明のさらに別の目的は、仮想ボリュームに効率的に記憶領域を対応付けることができ、かつ、仮想ボリュームの応答性能を向上できるようにした記憶制御装置及び仮想ボリュームの制御方法を提供することにある。本発明の更なる目的は、後述する実施形態の記載から明らかになるであろう。

20

【課題を解決するための手段】

【0014】

上記課題を解決すべく、本発明の第1観点に従う記憶制御装置は、仮想的に形成される仮想ボリュームと、一つまたは複数の記憶装置を含む、複数のRAIDグループと、各RAIDグループにそれぞれストライプ状に設けられる第1実記憶領域であって、それぞれ複数の第2実記憶領域を有する複数の第1実記憶領域を管理するためのプール部と、上位装置から仮想ボリュームに関するライトコマンドが発行された場合、各第1実記憶領域のうち所定の第1実記憶領域に含まれる各第2実記憶領域の中から所定の第2の実記憶領域を選択し、その所定の第2実記憶領域をライトコマンドに対応する仮想ボリューム内の所定領域に対応付ける制御部であって、一つの第1実記憶領域に一つの仮想ボリュームを対応付ける制御部と、仮想ボリュームに対応付けられている各第2実記憶領域の中から移動対象の第2実記憶領域を選択し、その移動対象の第2実記憶領域に記憶されているデータの移動先となる移動先第1実記憶領域を、各第1実記憶領域のうち移動対象の第2実記憶領域が設けられているRAIDグループ以外の他のRAIDグループ内の他の各第1実記憶領域の中から選択する移動先決定部と、移動対象の第2実記憶領域と、移動先第1実記憶領域とを対応付けて記憶する移動先記憶部と、上位装置が移動対象の第2実記憶領域に対応するコマンドを発行した場合は、コマンドの処理の中で、移動対象の第2実記憶領域に記憶されているデータを、移動先記憶部により記憶されている移動先第1実記憶領域内に移動させるコマンド処理部と、を備える。

30

40

【0015】

第2観点では、第1観点において、移動先記憶部は、移動対象の第2実記憶領域を特定する移動対象情報と、移動先第1実記憶領域を特定する移動先情報と、各RAIDグループの負荷に関する負荷情報と、各RAIDグループの使用容量に関する使用容量情報と、を記憶しており、移動先決定部は、所定時刻が到来した場合またはユーザから指示された場合のいずれかの場合に、移動先記憶部に記憶されている移動先情報を消去させ、さらに、各第2実記憶領域のうちヌルデータのみが記憶されている第2実記憶領域を解放して、未使用の第2実記憶領域に変更させ、さらに、負荷情報に基づく負荷分散処理であって、

50

相対的に高負荷の第2実記憶領域を相対的に低負荷の第1実記憶領域内に移動させるための、新たな移動先情報を作成して、移動先記憶部に記憶させる負荷分散処理と、使用容量情報に基づく使用容量平均化処理であって、相対的に使用容量の大きい第1実記憶領域内の第2実記憶領域を相対的に使用容量の小さい第1実記憶領域内に移動させるための、新たな移動先情報を作成して、移動先記憶部に記憶させる使用容量平均化処理とを、それぞれ実行し、コマンド処理部は、上位装置から移動対象の第2実記憶領域についてのライトコマンドが発行された場合、移動対象の第2実記憶領域に記憶されているデータを読み出し、その読み出されたデータとライトコマンドに係るライトデータとをマージして、移動先第1実記憶領域内の第2実記憶領域に書き込み、さらに、上位装置にライトコマンドの処理が完了した旨を通知し、上位装置から移動対象の第2実記憶領域についてのリードコマンドが発行された場合、移動対象の第2実記憶領域からデータを読み出して上位装置に送信し、上位装置にリードコマンドの処理が完了した旨を通知した後で、移動対象の第2実記憶領域から読み出されたデータを移動先第1実記憶領域内の第2実記憶領域に書き込むようになっている。

10

【0016】

第3観点では、第1観点において、コマンド処理部は、上位装置から移動対象の第2実記憶領域についてのライトコマンドが発行された場合、移動対象の第2実記憶領域に記憶されているデータを読み出し、その読み出されたデータとライトコマンドに係るライトデータとをマージして、移動先第1実記憶領域内の第2実記憶領域に書き込む。

【0017】

20

第4観点では、コマンド処理部は、上位装置から移動対象の第2実記憶領域についてのリードコマンドが発行された場合、移動対象の第2実記憶領域からデータを読み出して上位装置に送信し、上位装置にリードコマンドの処理が完了した旨を通知した後で、移動対象の第2実記憶領域から読み出されたデータを移動先第1実記憶領域内の第2実記憶領域に書き込む。

【0018】

第5観点では、第3観点において、移動先記憶部は、移動対象の第2実記憶領域を特定する移動対象情報と、移動先第1実記憶領域を特定する移動先情報と、各RAIDグループの負荷に関する負荷情報と、を記憶しており、移動先決定部は、移動先記憶部に記憶されている移動先情報を消去させ、相対的に高負荷の第2実記憶領域を相対的に低負荷の第1実記憶領域内に移動させるための、新たな移動先情報を負荷情報に基づいて作成し、移動先記憶部に記憶させる。

30

【0019】

第6観点では、第3観点において、移動先記憶部は、移動対象の第2実記憶領域を特定する移動対象情報と、移動先第1実記憶領域を特定する移動先情報と、各RAIDグループの負荷に関する負荷情報と、各RAIDグループの使用容量に関する使用容量情報と、を記憶しており、移動先決定部は、移動先記憶部に記憶されている移動先情報を消去させ、相対的に高負荷の第2実記憶領域を相対的に低負荷の第1実記憶領域内に移動させるための、新たな移動先情報を負荷情報に基づいて作成し、移動先記憶部に記憶させ、さらに、相対的に使用容量の大きい第1実記憶領域内の第2実記憶領域を相対的に使用容量の小さい第1実記憶領域内に移動させるための、新たな移動先情報を使用容量情報に基づいて作成し、移動先記憶部に記憶させる。

40

【0020】

第7観点では、第5観点において、移動先決定部は、移動先記憶部に記憶されている移動先情報を消去させた後で、各第2実記憶領域のうちヌルデータのみが記憶されている第2実記憶領域を解放して、未使用の第2実記憶領域に変更させる。

【0021】

第8観点では、移動先決定部は、移動対象の第2実記憶領域を第1実記憶領域の単位で複数選択する。

【0022】

50

第9観点では、制御部は、仮想ボリュームを生成する場合に、仮想ボリューム内の各仮想的記憶領域を、初期データが記憶されている初期化用の第2実記憶領域に対応付け、上位装置から仮想ボリュームに関するライトコマンドが発行された場合、ライトコマンドに対応する仮想的記憶領域の対応付け先を、初期化用の第2実記憶領域から、選択された所定の第2実記憶領域に切り替える。

【0023】

第10観点では、制御部は、前回のライト要求に対応して仮想ボリュームに対応付けられた第2実記憶領域に連続する、未使用の第2実記憶領域を、所定の第2実記憶領域として仮想ボリュームに対応付ける。

【0024】

第11観点では、仮想ボリュームに記憶されるデータには、所定サイズ毎に保証コードが設定されており、保証コードは、RAIDグループを識別するためのデータと、第1実記憶領域を識別するためのデータと、第1実記憶領域内における第2実記憶領域を識別するためのデータとを含んでいる。

【0025】

第12観点で従う仮想ボリュームの制御方法は、ライトコマンドに応じて実記憶領域が割り当てられる仮想ボリュームを制御するための方法であって、仮想ボリュームは、複数の仮想的記憶領域を備えており、複数のRAIDグループを管理するためのプール部を作成し、各RAIDグループは、複数の記憶装置を跨るようにしてストライプ状に形成される第1実記憶領域であって、仮想的記憶領域に対応する第2実記憶領域をそれぞれ複数ずつ有する第1実記憶領域を複数備えており、上位装置から仮想ボリュームに関するライトコマンド要求が発行された場合、一つの第1実記憶領域に複数の仮想ボリュームが対応付けられないようにして、ライトコマンドに対応する仮想的記憶領域に、各第1実記憶領域のうち所定の第1実記憶領域に含まれる所定の第2実記憶領域を対応付け、仮想的記憶領域に対応付けられる所定の第2実記憶領域に、上位装置から受領したライトデータを記憶させ、仮想的記憶領域に対応付けられている各第2実記憶領域の中から、相対的に高負荷の第2実記憶領域を移動対象の第2実記憶領域として選択し、各第1実記憶領域のうち移動対象の第2実記憶領域が設けられているRAIDグループ以外の他のRAIDグループ内の他の各第1実記憶領域であって、相対的に低負荷の第1実記憶領域を、移動対象の第2実記憶領域に記憶されているデータの移動先となる移動先第1実記憶領域として選択し、移動対象の第2実記憶領域と、移動先第1実記憶領域とを対応付けて移動先記憶部に記憶させ、上位装置が移動対象の第2実記憶領域に対応するライトコマンドを発行した場合は、移動対象の第2実記憶領域に記憶されているデータを読み出し、その読み出されたデータとライトコマンドに係るライトデータとをマージして、移動先第1実記憶領域内の第2実記憶領域に書き込む。

【0026】

第13観点では、第12観点において、上位装置から移動対象の第2実記憶領域についてのリードコマンドが発行された場合には、移動対象の第2実記憶領域からデータを読み出して上位装置に送信し、上位装置にリードコマンドの処理が完了した旨を通知した後で、移動対象の第2実記憶領域から読み出されたデータを移動先第1実記憶領域内の第2実記憶領域に書き込む。

【0027】

第14観点で従う記憶制御装置は、仮想的に形成される仮想ボリュームと、一つまたは複数の記憶装置を含む、複数のRAIDグループと、各RAIDグループにそれぞれストライプ状に設けられる第1実記憶領域であって、それぞれ複数の第2実記憶領域を有する複数の第1実記憶領域を管理するためのプール部と、上位装置から仮想ボリュームに関するライトコマンドが発行された場合、各第1実記憶領域のうち所定の第1実記憶領域に含まれる各第2実記憶領域の中から所定の第2の実記憶領域を選択し、その所定の第2実記憶領域をライトコマンドに対応する仮想ボリューム内の所定領域に対応付ける制御部であって、一つの第1実記憶領域に一つの仮想ボリュームを対応付ける制御部と、仮想ボリ

10

20

30

40

50

ームに対応付けられている各第1実記憶領域の中から移動対象の第1実記憶領域を選択し、その移動対象の第1実記憶領域に記憶されているデータの移動先となる移動先RAIDグループを、移動対象の第1実記憶領域が設けられているRAIDグループ以外の他のRAIDグループの中から選択する移動先決定部と、移動対象の第1実記憶領域と、移動先RAIDグループとを対応付けて記憶する移動先記憶部と、上位装置が移動対象の第1実記憶領域に対応するコマンドを発行した場合は、コマンドの処理の中で、移動対象の第1実記憶領域に記憶されているデータを、移動先記憶部により記憶されている移動先RAIDグループ内に移動させるコマンド処理部と、を備える。

【0028】

本発明の構成の少なくとも一部は、コンピュータプログラムとして構成できる。このコンピュータプログラムは、記録媒体に固定して配布したり、通信ネットワークを介して配信することができる。さらに、前記観点の組合せ以外の他の組合せも本発明の範囲に含まれる。

10

【図面の簡単な説明】

【0029】

【図1】図1は、本発明の実施形態の全体概念を示す説明図である。

【図2】図2は、記憶制御装置を含むシステムの全体構成を示す説明図である。

【図3】図3は、記憶制御装置のブロック図である。

【図4】図4は、プログラム及びテーブルを示す図である。

【図5】図5は、仮想ボリュームとチャンク及びページの関係を示す図である。

20

【図6】図6は、チャンク及びページを管理するテーブル群を示す図である。

【図7】図7は、仮想ボリュームの生成時におけるテーブル群の接続の様子を示す説明図である。

【図8】図8は、ライトデータを書き込む場合に、仮想ボリュームに割り当てるページを初期設定用ページから所定のページに切り替える様子を示す図である。

【図9】図9は、移動先管理テーブルを示す図である。

【図10】図10は、記憶制御装置の全体動作の流れを示す図である。

【図11】図11は、プール作成処理を示すフローチャートである。

【図12】図12は、データに付加される保証コードの構成を示す図である。

【図13】図13は、チャンクの状態遷移を示す図である。

30

【図14】図14は、チャンクを管理するためのキューを示す図である。

【図15】図15は、フォーマット処理を示すフローチャートである。

【図16】図16は、フォーマット処理の一部を示すフローチャートである。

【図17】図17は、仮想ボリューム生成処理を示すフローチャートである。

【図18】図18は、ライト処理を示すフローチャートである。

【図19】図19は、図18に続くフローチャートである。

【図20】図20は、図19に続くフローチャートである。

【図21】図21は、チャンクを仮想ボリュームに割り当てる処理を示すフローチャートである。

【図22】図22は、ページを仮想ボリュームに割り当てる処理を示すフローチャートである。

40

【図23】図23は、ページ状態を変更する処理を示すフローチャートである。

【図24】図24は、リード処理を示すフローチャートである。

【図25】図25は、図24に続くフローチャートである。

【図26】図26は、図25に続くフローチャートである。

【図27】図27は、移動先を決定する処理のフローチャートである。

【図28】図28は、負荷分散処理を示すフローチャートである。

【図29】図29は、使用容量を平均化させる処理のフローチャートである。

【図30】図30は、本発明の効果の一つを示す図である。

【図31】図31は、第2実施例に係り、移動先を決定する処理を示すフローチャートで

50

ある。

【図32】図32は、0データのみを記憶するページを解放する処理を示すフローチャートである。

【発明を実施するための形態】

【0030】

以下、図面に基づいて、本発明の実施の形態を説明する。最初に、本発明の概要を説明し、次に、実施例について説明する。本発明は、後述のように、仮想ボリューム5への実記憶領域の割当てを、チャンク7単位で行う。チャンク7は、複数のページ8から構成される。一つのチャンク7には、一つの仮想ボリューム5が対応付けられる。つまり、一つのチャンク7が異なる複数の仮想ボリューム5に対応付けられることはない。このため、チャンク7の記憶領域を効率的に使用することができる。

10

【0031】

仮想ボリューム5の生成時に、各仮想的記憶領域5Aと初期設定用のページ8とが予め対応付けられる。ホスト2から仮想ボリューム5へのライト要求が発行されると、チャンク7内のページ8が順番に使用されて、ライト要求に関わる仮想的記憶領域5Aに割り当てられる。その割り当てられたページ8にライトデータが書き込まれる。データの書込み時には、ライト要求に係る仮想的記憶領域5Aの接続先が、初期設定用のページ8から、チャンク7内の所定のページ8に切り替えられる。所定のページ8とは、前回のライト処理時に使用されたページに連続するページである。つまり、ライトデータを書き込む際には、仮想的記憶領域5Aに割り当てられるページを、初期設定用のページ8から所定のページ8に切り替えるだけで済むため、仮想ボリューム5の応答性能を向上できる。

20

【0032】

さらに、チャンク7内に空きページ8が無くなると、新たなチャンク7が選択されて仮想ボリューム5に割り当てられる。新たなチャンク7は、別のRAIDグループ6b内のチャンク群から選択される。これにより、各RAIDグループ6a、6b間で負荷を分散させることができる。

【0033】

さらに、本発明では、負荷等に基づいて移動対象のページ8またはチャンク7を選択して、その移動先を決定する。そして、移動対象のページ8またはチャンク7にホスト2がアクセスした場合、ホスト2から発行されるコマンド（ライトコマンドまたはリードコマンド）を処理しながら、データを移動させる。

30

【0034】

図1は、本発明の実施形態の概要を示す説明図である。図1に関する以下の記載は、本発明の理解及び実施に必要な程度で本発明の概要を示しており、本発明の範囲は図1に示す構成に限定されない。

【0035】

図1に示すシステムは、例えば、記憶制御装置1と、「上位装置」としてのホスト2とを備える。ホスト2は、例えば、サーバコンピュータまたはメインフレームコンピュータのようなコンピュータ装置として構成される。ホスト2がホストコンピュータの場合、例えば、FICON（Fibre Connection：登録商標）、ESCON（Enterprise System Connection：登録商標）、ACONARC（Advanced Connection Architecture：登録商標）、FIBARC（Fibre Connection Architecture：登録商標）等の通信プロトコルに従って、データ通信が行われる。ホスト2がサーバコンピュータ等の場合、例えば、FCP（Fibre Channel Protocol）またはiSCSI（internet Small Computer System Interface）等の通信プロトコルに従って、データ通信が行われる。

40

【0036】

記憶制御装置1は、ホスト2に通信ネットワークを介して接続される。記憶制御装置1は、例えば、コントローラ3と、記憶装置4と、仮想ボリューム5（1）、5（2）とを備える。特に区別する必要がない場合、仮想ボリューム5と呼ぶ。各RAIDグループ6

50

a, 6 bは、それぞれ複数ずつの記憶装置4から構成される。特に区別しない場合、RAIDグループ6と呼ぶ。

【0037】

記憶装置4としては、例えば、ハードディスクデバイス、半導体メモリデバイス、光ディスクデバイス、光磁気ディスクデバイス、磁気テープデバイス、フレキシブルディスクデバイス等のデータを読み書き可能な種々のデバイスを利用可能である。

【0038】

記憶装置4としてハードディスクデバイスを用いる場合、例えば、FC (Fibre Channel) ディスク、SCSI (Small Computer System Interface) ディスク、SATA ディスク、ATA (AT Attachment) ディスク、SAS (Serial Attached SCSI) ディスク等を用いることができる。また、例えば、フラッシュメモリ、FeRAM (Ferroelectric Random Access Memory)、MRAM (Magnetoresistive Random Access Memory)、相変化メモリ (Ovonic Unified Memory)、RRAM (Resistance RAM) 等の種々の記憶装置4を用いることもできる。さらに、例えば、フラッシュメモリデバイスとハードディスクドライブのように、種類の異なる記憶装置4を混在させる構成でもよい。

【0039】

各RAIDグループ6 a, 6 bの有する物理的記憶領域は、ストライプ状の複数のチャンク7に区切られる。各チャンク7は、連続する複数のページ8から構成される。チャンク7は「第1実記憶領域」に該当し、ページ8は「第2実記憶領域」に該当する。理解のために、一方のRAIDグループ6 aに属する第1チャンク7に符号 " a 1 " を与え、第1チャンク7 (a 1) に属する各ページに連番を添える。他方のRAIDグループ6 bについても同様である。従って、例えば、 " a 2 - 3 " は、RAIDグループ6 aの第2チャンク内の3番目のページであることを意味し、 " b 1 - 1 " は、RAIDグループ6 bの第1チャンク内の1番目のページであることを意味する。

【0040】

仮想ボリューム5は、複数の仮想的記憶領域5 Aから構成される。仮想的記憶領域5 Aとページ8のサイズは同一である。一つの例では、1枚のページ8のサイズはSZ1バイト (例えば、32MB)、一つのチャンク7のサイズはSZ2バイト (例えば、1GB)、仮想ボリューム5のサイズはSZ3バイト (例えば、10GB) である。この場合、一つの仮想ボリューム5はN1個 (例えば、10個) のチャンク7から構成され、一つのチャンク7はN2枚 (例えば、32枚) のページ8から構成される。上記括弧内の数値は、理解のための一例に過ぎず、本発明の範囲は上記数値に限定されない。上述のページサイズ、チャンクサイズ、仮想ボリュームサイズ等は、可変に設定できる。

【0041】

コントローラ3は、記憶制御装置1の動作を制御する。例えば、コントローラ3は、ユーザからの指示に基づいて、RAIDグループ6 a, 6 b及び仮想ボリューム5を生成させる。また、コントローラ3は、ホスト2から発行されるコマンド (リードコマンド、ライトコマンド) に応じて処理を実行し、その処理結果をホスト2に送信する。

【0042】

さらに、コントローラ3は、ライトコマンドを受領した場合、ライトコマンドにより指定される仮想的記憶領域5 Aに、ページ8が割り当てられているか否かを判断する。指定される仮想的記憶領域5 Aにページ8が割り当てられていない場合、コントローラ3は、チャンク7内の所定のページ8を、指定される仮想的記憶領域5 Aに割り当てる。コントローラ3は、割り当てられた所定ページ8にライトデータを書き込む。

【0043】

コントローラ3は、一つのチャンク7に一つの仮想ボリューム5だけが対応付けられるように、ページ割当てを制御する。一つのチャンク7に含まれる各ページ8は、一つの仮想ボリューム5にのみ割り当てられる。一つのチャンク7内にそれぞれ異なる複数の仮想ボリューム5に割り当てられるページ8が混在することはない。一つのチャンク7内では、論理アドレスの値を問わずに、連続するページ8が使用される。

【 0 0 4 4 】

一方の仮想ボリューム 5 (1) を例に挙げて説明すると、最初のライトコマンドについて、チャンク 7 (a 1) 内の先頭ページ 8 (a 1 - 1) が使用され、次のライトコマンドについては、その先頭ページ 8 (a 1 - 1) に続く次のページ 8 (a 1 - 2) が使用され、さらに別のライトコマンドについては、次のページ 8 (a 1 - 3) が使用される。そして、最後のページ 8 (a 1 - 4) を使用した後で、さらにライトコマンドを受領した場合、新たなチャンク 7 (b 1) が仮想ボリューム 5 (1) に割り当てられる。

【 0 0 4 5 】

他方の仮想ボリューム 5 (2) には、RAIDグループ 6 a の第 2 チャンク 7 内の先頭ページ 8 (a 2 - 1) が割り当てられている。もしも、仮想ボリューム 5 (2) を対象とする新たなライトコマンドが発行された場合、次のページ 8 が仮想ボリューム 5 (2) に割り当てられる。そのページ 8 には、" a 2 - 2 " の符号が添えられるはずであるが、図 1 では省略する。

10

【 0 0 4 6 】

このように、仮想ボリューム 5 には、ホスト 2 に見せかけているボリュームサイズよりも小さいサイズの実記憶領域 (ページ 7、チャンク 8) が割り当てられた状態で、ホスト 2 に提供される。そして、ホスト 2 からのライトコマンドに応じて、必要量の実記憶領域が動的に割り当てられる。

【 0 0 4 7 】

さらに、コントローラ 3 は、以下に述べるように、各 RAIDグループ 6 a , 6 b 間でデータを移動させる。コントローラ 3 は、例えば、データへのアクセス頻度 (負荷) 等に基づいて、移動対象のデータを選択する。コントローラ 3 は、移動対象のデータを移動させる移動先として、低負荷の RAIDグループ 6 を選択する。または、コントローラ 3 は、低負荷の RAIDグループ内のチャンクを移動先として選択することもできる。

20

【 0 0 4 8 】

上述のように、コントローラ 3 は、移動対象のデータ及びその移動先を決定し (S 1)、記憶させる (S 2)。ホスト 2 が移動対象のデータにアクセスすると、コントローラ 3 は、移動対象データを読み出してコマンド処理を行うとともに、その移動対象データを移動先の記憶領域に移動させる (S 3)。

【 0 0 4 9 】

このように構成される本実施形態では、複数のページ 8 を有するチャンク 7 の単位で、仮想ボリューム 5 に実記憶領域 (物理的記憶領域) を割り当て、かつ、一つのチャンク 7 を一つの仮想ボリューム 5 にのみ割り当てる。従って、後述のように、チャンク 7 内の記憶領域を有効に使用できる。また、通常の論理ボリュームと同様にして、仮想ボリューム 5 を管理でき、制御構造を簡素化できる。

30

【 0 0 5 0 】

本実施形態では、複数の RAIDグループ 6 a , 6 b を均等に使用するべく、各 RAIDグループ 6 a , 6 b からそれぞれチャンク 7 (a 1) , 7 (b 1) を選択して、仮想ボリューム 5 (1) に割り当てる。これにより、各 RAIDグループ 6 a , 6 b の負荷を均一化することができる。

40

【 0 0 5 1 】

本実施形態では、例えば、RAIDグループ 6 a , 6 B 間の負荷を分散させべくデータ移動計画を予め作成しておき、ホスト 2 が移動対象のデータにアクセスした場合に、ホスト 2 の要求する処理を実行しながら、移動対象のデータを予め設定された移動先に移動させる。従って、ホスト 2 からのアクセスを契機として、通常のコマンド処理の中でデータ移動を行うことができる。リード処理の場合は、読み出したデータを移動先に書き込む作業が発生するが、移動先は低負荷の記憶領域の中から選択されるため、記憶制御装置の応答性能に与える影響は少ない。

【 0 0 5 2 】

さらに、本実施形態では、格納先の移動が予定されたものの、ホスト 2 にアクセスされ

50

なかったデータは、移動契機が発生しないので、移動されることはない。従って、実記憶領域へのアクセスを低減し、ホスト2により使用されるデータのみを効率的に移動させることができる。

【0053】

さらに、本実施形態では、ホスト2によるアクセスを移動契機とするため、移動先を決定するタイミングと、実際の移動タイミングとの間に時間遅れが生じる。しかし、ホスト2により頻繁にアクセスされるデータは、直ちに低負荷の記憶領域に移動されるため、タイミングのずれによる影響は小さい。以下、本実施形態を詳細に説明する。

【実施例1】

【0054】

図2は、本実施例に係る記憶制御装置10を含む情報処理システムの全体構成を示す説明図である。この情報処理システムは、例えば、少なくとも一つの記憶制御装置10と、一つまたは複数のホスト20と、少なくとも一つの管理サーバ70とを含んで構成することができる。

【0055】

先に図1で述べた実施形態との対応関係を説明する。記憶制御装置10は記憶制御装置1に、ホスト20はホスト2に、コントローラ30はコントローラ3に、記憶装置40は記憶装置4に、仮想ボリューム50Vは仮想ボリューム5に、RAIDグループ90はRAIDグループ6a, 6bに、それぞれ対応する。図1で述べた説明と重複する説明は、できるだけ省略する。

【0056】

ホスト20と記憶装置10とは、第1通信ネットワーク80を介して接続される。第1通信ネットワーク80は、例えば、FC-SAN (Fibre Channel-Storage Area Network) やIP-SAN (Internet Protocol_SAN) のように構成される。

【0057】

管理サーバ70は、記憶制御装置10の設定を変更等するための装置である。管理サーバ70は、例えば、LAN (Local Area Network) のような第2通信ネットワーク81を介して、記憶制御装置10に接続されている。なお、ホスト20にストレージ管理機能を設けて、ホスト20側から記憶制御装置10の設定変更等を行うように構成することもできる。

【0058】

記憶制御装置10の詳細は後述するが、記憶制御装置10は、仮想ボリューム50Vと通常ボリューム50Nとを備える。なお、図中では、論理ボリュームを"LU"と表記している。LUとは、Logical Unitの略である。

【0059】

仮想ボリューム50Vは、図1で述べたように、仮想的に生成される論理ボリュームであって、ホスト20からのライトコマンドに応じて記憶領域が割り当てられるボリュームである。つまり、仮想ボリューム50Vは、ホスト20に提供されるボリュームサイズと、実際に有する記憶領域のサイズとが一致しない。通常ボリューム50Nは、RAIDグループ90の有する記憶領域に基づいて生成されるボリュームである。

【0060】

プール部60は、複数のRAIDグループ90の有する記憶領域を管理する。プール部60で管理されている記憶領域は、チャンク91 (図5参照) 単位で仮想ボリューム50Vに割り当てられる。

【0061】

コントローラ30は、各ボリューム50V, 50Nの論理アドレスをRAIDグループ90の物理アドレスに変換等して、データを記憶装置40に書き込んだり、あるいは、記憶装置40から読み出したデータの物理アドレスを論理アドレスに変換等して、データをホスト20に送信する。

【0062】

10

20

30

40

50

さらに、後述のように、コントローラ 30 は、データ移動計画を作成し、ホスト 20 からのアクセスを契機としてデータ移動計画を実行する。

【0063】

図 3 は、コントローラ 30 の構成を示すブロック図である。記憶制御装置 10 には、複数の増設筐体 43 を接続することができる。増設筐体 43 は、複数の記憶装置 40 を収容している。ユーザは、必要に応じて、増設筐体 43 を接続することにより、システムの総記憶容量を増大させることができる。

【0064】

記憶制御装置 10 は、複数のコントローラ 30 (# 0) , 30 (# 1) を備える。いずれか一方のコントローラ 30 が障害等によって停止した場合でも、他方のコントローラ 30 により動作を継続することができる。以下、特に区別する必要が無い場合、コントローラ 30 と称する。

10

【0065】

コントローラ 30 は、例えば、第 1 通信回路 310 (図中、 F E I / F) と、第 2 通信回路 320 (図中、 S A S) と、データ転送制御回路 330 (図中、 D C T L) と、キャッシュメモリ 340 (図中、 C M) と、ローカルメモリ 350 (図中、 L M) と、マイクロプロセッサ 360 (図中、 M P U) と、メモリコントローラ 370 (図中、 M C) と、エキスパンダ 380 (図中、 E X P) とを備える。

【0066】

第 1 通信回路 310 は、ホスト 20 と通信を行うための制御回路である。第 2 通信回路 320 は、各記憶装置 40 と通信を行うための制御回路である。データ転送制御回路 330 は、記憶制御装置 10 内のデータの流れを制御するための回路である。各データ転送制御回路 330 は、互いに接続されている。キャッシュメモリ 340 は、例えば、ホスト 20 から受領したライトデータと、記憶装置 40 から読み出されるデータとを記憶する。さらに、キャッシュメモリ 340 には、記憶制御装置 10 の構成または動作を管理するための管理用データ等が記憶される場合もある。

20

【0067】

ローカルメモリ 350 は、例えば、マイクロプロセッサ 360 により使用される各種データを記憶する。マイクロプロセッサ 360 は、記憶装置 40 またはローカルメモリ 350 からコンピュータプログラムを読み込んで実行することにより、後述のように、記憶制御装置 10 の動作を制御する。

30

【0068】

メモリコントローラ 370 は、マイクロプロセッサ 360 をローカルメモリ 350 及びデータ転送制御回路 330 に接続させるための制御回路である。エキスパンダ 380 は、第 2 通信回路 320 の通信ポートを拡張するための回路である。なお、図 3 に示す構成は一例であって、他の構成の記憶制御装置にも本発明を適用できる。

【0069】

図 4 は、記憶制御装置 10 の有するプログラム及びテーブルの例を示す。例えば、メモリ 360 にはプログラム 110 ~ 130 が記憶され、キャッシュメモリ 340 にはテーブル 140 ~ 160 が記憶される。または、特定の記憶装置 40 内に、プログラム及びテーブルを格納し、必要に応じてメモリにロードして使用する構成でもよい。

40

【0070】

R A I D 制御プログラム 110 は、R A I D を制御するためのプログラムである。R A I D 制御プログラム 110 は、ドライブ管理テーブル 140 等を用いて、R A I D 構成を管理する。

【0071】

コマンド処理プログラム 120 は、ホスト 20 からのコマンドを処理するためのプログラムである。コマンドとしては、例えば、ライトコマンド、リードコマンド、その他のコマンドを挙げることができる。その他のコマンドとしては、例えば、ボリュームの作成、ボリュームの削除、ボリュームコピーの指示、仕様等を問い合わせるためのコマンドを挙

50

げることができる。コマンド処理プログラム 120 は、ライトコマンド及びリードコマンドを処理する場合、その処理対象のデータに「データの移動先」が関連づけられているか否かを確認する。もしも、処理対象データにデータの移動先が関連づけられている場合、コマンド処理プログラム 120 は、コマンド処理の中で、処理対象データを予め設定されている移動先に移動させる。

【0072】

移動先決定プログラム 130 は、仮想ボリューム 50V に記憶されているデータのうち、第1条件を満たすデータを移動対象データとして選択し、第2条件を満たす移動先に移動させる。

【0073】

例えば、第1条件は、予め設定されている負荷の値を超える RAID グループ 90 のうち、最も高い負荷値を有すること、である。例えば、第2条件は、予め設定される他の負荷値を下回る RAID グループ 90 のうち、最も低い負荷値を有すること、である。

【0074】

他の第1条件は、使用容量が予め設定される使用容量の値よりも大きい RAID グループ 90 のうち最も使用容量が大きいこと、である。他の第2条件は、予め設定される他の使用容量値を下回る RAID グループ 90 のうち、最も小さい使用容量を有すること、である。使用容量は、百分率形式で表すこともできる。

【0075】

ホストアクセスを契機としてデータを移動させると、移動元 RAID グループ 90 は、その負荷及び使用容量がともに減少する。移動先 RAID グループ 90 は、その負荷及び使用容量がともに増加する。

【0076】

ドライブ管理テーブル 140 は、各記憶装置 40 を管理するための情報である。ドライブ管理テーブル 140 は、例えば、各記憶装置 40 の種類、容量、所属する RAID グループの番号等を管理する。

【0077】

ページ管理テーブル 150 は、各ページ 92 及び各チャンク 91 等を管理する。ページ管理テーブル 150 の詳細は、図 6 ~ 図 8 で後述する。

【0078】

移動先管理テーブル 160 は、移動対象データの移動先に関する情報と、移動先を決定するために使用される情報とを管理する。移動先管理テーブル 160 の詳細は、図 9 で後述する。

【0079】

図 5 は、仮想ボリューム 50V とチャンク 91 との関係を示す説明図である。仮想ボリューム 50V は、複数の仮想的記憶領域 500 を有する。チャンク 91 は、複数のページ 92 を有する。図 5 では、便宜上、各ページ 92 を、ストライプ状のチャンク 91 を横方向に区切って形成するかのよう示すが、実際には、各ページ 92 は、ストライプ列に沿うようにして形成される。

【0080】

図 5 では、最初に、第1のチャンク 91 (#0) が仮想ボリューム 50V に割り当てられたとする。仮想ボリューム 50V へのライトコマンド(ライト要求)が受領される度に、第1チャンク 91 (#0) 内のページ 92 が順番に選択されて、ライトコマンドに対応する仮想的記憶領域 500 に対応付けられる。ライトデータは、仮想的記憶領域 500 に対応付けられたページ 92 に書き込まれる。つまり、そのページ 92 を構成する各記憶装置 40 の各記憶領域にライトデータが書き込まれる。

【0081】

第1チャンク 91 (#0) の最初のページ 92 (1-0) が使用された後、次のページ 92 (1-1) が使用され、さらに次のページ 92 (1-2) が使用される。そして、第1チャンク 91 (#0) の最終ページ 92 (1-4) まで使用されたとする。これにより

10

20

30

40

50

、第1チャンク91(#0)の全ページ92が使用されたことになる。

【0082】

新たなライトコマンドがホスト20から発行されると、第2のチャンク91(#1)が他のRAIDグループ90内から選択され、第2チャンク91(#1)の先頭ページ92(2-0)が使用される。以下、第1チャンク91(#0)で述べたと同様に、第2チャンク91(#1)内の各ページ92が順番に使用される。

【0083】

このように、仮想ボリューム50Vには、複数のチャンク91(#0)、91(#1)が対応付けられる。それらのチャンク91(#0)、91(#1)は、それぞれ別々のRAIDグループ90の中から選択される。つまり、仮想ボリューム50Vには、複数のRAIDグループ90の中から選択される複数のチャンク91が対応付けられる。ホスト20による仮想ボリューム50Vへのデータ書込みに応じて、対応付けられたチャンク91内のページ92が順番に使用される。

10

【0084】

なお、ホスト20が仮想ボリューム50V内のデータ消去を要求する場合、消去対象のデータの論理アドレスに対応するページ92は開放され、未使用ページに戻る。未使用ページには、別のライトデータが書き込まれる。

【0085】

ホスト20が仮想ボリューム50Vを使用すればするほど、多くのページ92が仮想ボリューム50Vに次第に割り当てられていく。これにより、仮想ボリューム50Vが実際に有する物理的記憶領域の量と、仮想ボリューム50Vがホスト20に見せかけている容量との差は縮まっていく。

20

【0086】

プール部60内に、新たなRAIDグループ90が追加された場合を検討する。図示は省略するが、その新たなRAIDグループの番号を#2とする。チャンク91(#1)を使い切った後であれば、新たなRAIDグループ90(#2)からチャンク91(#2)が選択されるかも知れない。

【0087】

しかし、せっかく新たなRAIDグループ90(#2)が追加されても、前のチャンク91(#1)を使い終わるまで、新たなRAIDグループ90(#2)を使用できないならば、それは無駄である。

30

【0088】

そこで、本実施例では、各RAIDグループ90(#0)~90(#2)の負荷または/及び使用容量を均等にするためのデータ移動計画を作成し、ホストアクセスを契機として、データを移動させる。これにより、プール部60に新たに追加されるRAIDグループ90(#2)を比較的速やかに使い始めることができる。

【0089】

図6は、仮想ボリューム50Vを管理するためのページ管理テーブル150を示す説明図である。図6に示すように、ページ管理テーブル150は、複数のテーブル151~155等を含んで構成される。ページ管理テーブル150は、キャッシュメモリ340上に設けられ、その後、所定の記憶装置40内の所定の管理領域に格納される。また、ページ管理テーブル150は、ローカルメモリ350にもコピーされて、マイクロプロセッサ360により使用される。

40

【0090】

プールインデックス151は、例えば、プール部60で管理されているRAIDグループ90の識別番号等の、プール部60の構成(状態及び属性を含む)に関する情報を管理するための情報である。

【0091】

仮想ボリュームインデックス152は、仮想ボリューム50Vの構成に関する情報を管理するための情報である。仮想ボリュームインデックス152は、例えば、仮想ボリュー

50

ム50Vに割り当てられているチャンク91の識別番号、及び、リンクされる仮想アドレスインデックス153の情報等を管理する。

【0092】

仮想アドレスインデックス153は、仮想アドレスブロック154へのポインタを管理するための情報である。仮想アドレスブロック154は、ページアドレス情報155へのポインタを管理するための情報である。

【0093】

例えば、仮想ボリューム50Vのボリュームサイズを10GBとすると、仮想アドレスインデックス153は、仮想ボリューム50Vの仮想アドレス領域を4GBずつの領域に分けて管理する(最初の2つの領域は4GB、最後の一つの領域は2GBである。)。仮想アドレスブロック154は、それぞれ4GBの範囲をカバー可能である。このように、本実施例では、仮想ボリューム50Vの有する仮想アドレス領域を、仮想アドレスインデックス153と仮想アドレスブロック154との2段階に分けて階層管理する。これにより、範囲を絞り込んで検索することができ、該当ページ92に速やかにアクセスすることができる。なお、上記の数値(4GB、10GB、2GB)は、説明のための一例に過ぎず、本発明はそれらの数値に限定されない。前記各数値は可変に設定できる。

10

【0094】

ページアドレス情報155は、仮想ボリューム50Vを構成する各仮想的記憶領域500(つまり、仮想ページ500)の構成情報を管理するための情報である。ページアドレス情報155には、例えば、仮想ページに対応付けられる物理ページ92を示す物理アドレス及びページ状態が含まれる。

20

【0095】

プール内RAIDグループインデックス110は、プール部60で管理されている各RAIDグループ90の構成情報を管理するための情報である。プール内RAIDグループインデックス110は、例えば、プール部60内の各RAIDグループ90が有する各チャンク91の情報等を管理する。また、プール内RAIDグループインデックス110は、未割当てチャンクキューの先頭及び末尾をそれぞれ示すためのポインタを含む。

【0096】

チャンクインデックス120は、各仮想ボリューム50Vにどこまでチャンク91が割り当てられているかを管理するためのポインタを含む。つまり、チャンクインデックス120は、各仮想ボリューム50Vに割り当てられているチャンク数などを管理する。

30

【0097】

図7は、図6に示すページ管理テーブル150が初期化された状態を示す。仮想ボリューム50Vの生成時に、ページ管理テーブル150は、図7に示すように初期化される。仮想ボリューム50Vに含まれる各仮想的記憶領域500(仮想ページ500)は、特定のチャンク91内の特定のページ92にマッピングされる。

【0098】

例えば、図7に示すように、RAIDグループ90内の先頭チャンク91の最終ページ92が、初期化用の特定ページとして使用される。先頭チャンク91には、上述のテーブル群などの管理情報が記憶される。先頭チャンク91内の先頭ページ92から所定数のページ92までは、管理情報の退避領域として使用される。設定可能な仮想ボリュームの数等によって異なるが、管理情報の合計サイズはチャンクサイズ(例えば、1GB)未満となる。従って、少なくとも、先頭チャンク91の最終ページ92に管理情報が格納されることはない。つまり、先頭チャンク91の最終ページ92は、管理情報の格納先として使用されることがない。

40

【0099】

そこで、先頭チャンク91の最終ページである、初期化用のページ92に、予めゼロデータのみを記憶させる。そして、ページ管理テーブル150を初期化する場合には、仮想ボリューム50V内の全ての仮想的記憶領域500を、初期化用のページ92に対応付けておく。

50

【 0 1 0 0 】

これにより、仮想ボリューム 5 0 V の定義時（仮想ボリューム 5 0 V の生成時）に、テーブル群のマッピングに異常が生じていないことを予め確認することができる。さらに、図 8 に太線で示すように、仮想ボリューム 5 0 V にライトコマンドが発行された場合には、そのライトコマンドで指定される論理アドレスに対応する仮想的記憶領域 5 0 0 を、初期設定用のページ 9 2 からライトデータを書き込むべき所定のページに、接続し直すだけでよい。従って、仮想的記憶領域 5 0 0 の対応付け先を切り替えるだけで、ライトデータを格納させることができ、仮想ボリューム 5 0 V の応答性能を高めることができる。

【 0 1 0 1 】

図 9 は、移動先管理テーブル 1 6 0 の一例を示す。移動先管理テーブル 1 6 0 も、前記ページ管理テーブル 1 5 0 と同様に、例えば、キャッシュメモリ等に記憶させることができる。

10

【 0 1 0 2 】

移動先管理テーブル 1 6 0 は、例えば、チャンクインデックス 1 6 1 と、チャンク情報管理テーブル 1 6 2 とを対応付けることにより構成される。チャンク情報管理テーブル 1 6 2 は、各チャンク毎に管理される管理情報として、移動先情報と、移動先決定用情報とを記憶する。

【 0 1 0 3 】

移動先情報とは、移動対象のデータを格納するための格納先を示す情報である。移動先情報としては、例えば、移動先として決定された R A I D グループ 9 0 の番号等が用いられる。なお、R A I D グループ番号に限らず、R A I D グループ番号とチャンク番号の組合せ等を用いてもよい。

20

【 0 1 0 4 】

移動先決定用情報とは、移動対象のデータの移動先を決定するために使用される情報である。移動先決定用情報としては、例えば、そのデータへのアクセス頻度等を用いることができる。後述のように、R A I D グループ 9 0 の使用容量（使用率）を移動先決定用情報として用いることもできる。

【 0 1 0 5 】

さらに、移動先を決定するために用いられる負荷情報としては、アクセス頻度に限らず、例えば、アクセスサイズ（ライトデータのサイズ等）、アクセスパターン（シーケンシャルアクセスかランダムアクセスか）等を用いてもよい。

30

【 0 1 0 6 】

チャンクインデックス 1 6 1 は、仮想ボリューム 5 0 V と、チャンク情報管理テーブル 1 6 2 の所定エントリとを対応付ける。所定エントリとは、チャンクインデックス 1 6 1 の示す仮想ボリューム 5 0 V に対応付けられているチャンク 9 0 の管理情報を記憶しているエントリである。

【 0 1 0 7 】

図 1 0 は、記憶制御装置 1 0 の全体動作を理解するためのフローチャートである。このフローチャートには、ユーザの手順も含まれている。まず最初に、ユーザは、管理サーバ 7 0 を介して記憶制御装置 1 0 に所定の指示を与えることにより、プール部 6 0 に R A I D グループ 9 0 を生成させ、さらに、その R A I D グループ 9 0 をフォーマットさせ、未割当てチャンクキュー等を作成させる（S 1 0）。ここで、プール部 6 0 内で管理される、仮想ボリューム 5 0 V 用の R A I D グループ 9 0 と、通常ボリューム 5 0 N 用の R A I D グループ 9 0 とは、連番で管理される。

40

【 0 1 0 8 】

続いて、ユーザは、管理サーバ 7 0 から記憶制御装置 1 0 に別の所定の指示を与えることにより、仮想ボリューム 5 0 V を作成させる（S 1 1）。上述の通り、仮想ボリューム 5 0 V の生成時に、各仮想的記憶領域 5 0 0 はそれぞれ初期設定用のページ 9 2 に対応付けられる。ここで、仮想ボリューム 5 0 V と通常ボリューム 5 0 N とは、連番で管理される。これにより、仮想ボリューム 5 0 V と通常ボリューム 5 0 N とを共通の管理方式で管

50

理することができ、記憶制御装置 10 内に仮想ボリューム 50 V と通常ボリューム 50 N とを混在させることができる。

【0109】

続いて、ユーザは、ホスト 20 と仮想ボリューム 50 V とを接続させる (S12)。ユーザは、ホスト 20 に繋がる LUN (Logical Unit Number) に、仮想ボリューム 50 V を接続させ、さらに、WWN (Logical Unit Number) の登録等を行わせる。

【0110】

ホスト 20 は、仮想ボリューム 50 V を認識し (S13)、仮想ボリューム 50 V に向けてライトコマンド等のコマンドを発行する。記憶制御装置 10 は、ホスト 20 からのコマンドに応じた処理を行い、その結果をホスト 20 に送信する (S14)。

10

【0111】

記憶制御装置 10 は、所定の周期により、または、管理サーバ 70 からのユーザ指示により、仮想ボリューム 50 V に対応付けられている各 RAID グループ 90 内のデータについて、移動先を決定して記憶する (S14)。

【0112】

記憶制御装置 10 は、移動対象として選択されたデータにホスト 20 がアクセスすると、そのホストアクセスに関するコマンドを処理しながら、予定されたデータ移動を実施する (S15)。以下、各処理の詳細を図を改めて説明する。

【0113】

図 11 は、プール作成処理を示すフローチャートである。以下に述べる各フローチャートは、各処理の概要を示す。いわゆる当業者であれば、図示されたステップの入れ替え、変更、削除、あるいは新たなステップの追加を行うことができるであろう。

20

【0114】

図 11 に示す処理は、プール作成要求が与えられると開始される。以下、動作の主体をコントローラ 30 とする。コントローラ 30 は、作成対象のプール部 60 のプールインデックス 151 について、そのプール状態を更新し (S20)、RAID グループ 90 を作成する (S21)。コントローラ 30 は、プール内 RAID グループインデックス 110 の状態を "処理中" に変化させる (S22)。

【0115】

コントローラ 30 は、管理情報を退避させるための領域を RAID グループ 90 内に設定し (S23)、さらに、チャンクインデックス 120 を作成する (S24)。プール部 60 内の全チャンク 91 について、以下のステップ S26 ~ S28 がそれぞれ実行される (S25)。

30

【0116】

コントローラ 30 は、対象チャンクに対応するページアドレス情報 155 を初期化し (S26)、対象チャンクの状態を "フォーマット待ち" に変更する (S27)。コントローラ 30 は、管理情報を退避させるための退避要求キューに、管理情報の対比要求をエンキューする (S28)。

【0117】

各チャンク 91 について S26 ~ S28 を実行した後、コントローラ 30 は、プール内 RAID グループインデックス 110 の状態を "有効" に変更する (S29)。そして、コントローラ 30 は、キャッシュメモリ 340 にヌルデータをステージングさせて (S30)、本処理を終了する。ライトデータの書き込まれていない仮想的記憶領域 500 からのデータ読み出しが要求された場合に、物理的記憶領域であるページ 92 にアクセスすることなく、ホスト 20 にヌルデータを返すためである。

40

【0118】

図 12 は、データ及び保証コードについて説明する図である。図 12 (a) に示すように、本実施例では、例えば、512 バイトのデータ D10 毎に 8 バイトの保証コード D11 を付加して、記憶装置 40 に記憶させる。保証コード D11 は、論理アドレスを検証するための部分と、ビットエラーを検証するための部分とを含むことができる。以下では、

50

論理アドレスを検証するための部分に着目して説明する。

【 0 1 1 9 】

図 1 2 (b) は、通常ボリューム 5 0 N に記憶されるデータに付加される、保証コードの構成を示す。通常ボリューム 5 0 N に関する保証コードのうち論理アドレスの検証に使用される部分は、4 ビットの予約領域 D 1 1 0 と、1 2 ビットの L U N 領域 D 1 1 1 と、1 6 ビットの L B A 領域 D 1 1 2 とを含んでいる。L B A とは、Logical Block Address の略である。L U N 領域 D 1 1 1 には、通常ボリューム 5 0 N に対応付けられる L U N が格納される。L B A 領域 D 1 1 2 には、データ D 1 0 の論理アドレスが格納される。

【 0 1 2 0 】

図 1 2 (c) は、仮想ボリューム 5 0 V に記憶されるデータに付加される、保証コードの構成を示す。仮想ボリューム 5 0 V に関する保証コードのうち論理アドレスの検証に使用される部分は、4 ビットの予約領域 D 1 1 0 と、8 ビットの R A I D グループ識別領域 D 1 1 3 と、4 ビットのチャンク識別領域 D 1 1 4 と、1 6 ビットのチャンク内 L B A オフセット領域 D 1 1 5 とを含んでいる。R A I D グループ識別領域 D 1 1 3 には、R A I D グループ 9 0 を識別するための情報が格納される。チャンク識別領域 D 1 1 4 には、チャンク 9 1 を識別するための情報のうち下位の 4 ビットが格納される。チャンク内 L B A オフセット領域 D 1 1 5 には、データ D 1 0 が格納されているチャンク 9 1 内において、そのチャンク 9 1 の先頭論理アドレスからのオフセット値が格納される。

【 0 1 2 1 】

図 1 3 , 図 1 4 に基づいて、チャンク 9 1 を管理するためのキューを説明する。図 1 3 は、キューの使用法を示す説明図である。プール部 6 0 が作成されると、プール部 6 0 で管理される各 R A I D グループ 9 0 の有する各チャンク 9 1 は、フォーマット待ちチャンクキュー Q 1 0 に登録されて管理される。フォーマットが開始されると、フォーマット中のチャンク 9 1 は、処理完了待ちチャンクキュー Q 2 0 に移される。そして、フォーマットが完了すると、フォーマット済チャンク 9 1 は、未割当てチャンクキュー Q 3 0 に移される。

【 0 1 2 2 】

図 1 4 は、各キュー Q 1 0 , Q 2 0 , Q 3 0 を模式的に示す説明図である。フォーマット待ちチャンクキュー Q 1 0 は、各 R A I D グループ 9 0 毎に用意される。そして、フォーマット処理を開始する場合、フォーマット待ちチャンクキュー Q 1 0 内の各 R A I D グループ 9 0 から所定の順番でチャンク 9 1 が取り出されて、処理完了待ちキュー Q 2 0 に接続される。フォーマットが完了したチャンク 9 1 は、上述の通り、未割当てチャンクキュー Q 3 0 に接続される。処理完了待ちチャンクキュー Q 2 0 に繋がれた順番で、フォーマットは完了するため、処理完了待ちキュー Q 2 0 内の各チャンク 9 1 の順番と、未割当てチャンクキュー Q 3 0 内の各チャンク 9 1 の順番とは、通常の場合、一致する。

【 0 1 2 3 】

図 1 5 , 図 1 6 に基づいてチャンク 9 1 をフォーマットする処理を説明する。図 1 5 は、フォーマット処理の全体を示し、図 1 6 は、フォーマット処理の一部を示す。

【 0 1 2 4 】

コントローラ 3 0 は、フォーマット待ちチャンクキュー Q 1 0 を確認することにより、フォーマット待ちのチャンク 9 1 が有るか否かを判定する (S 5 0) 。フォーマット待ちチャンクキュー Q 1 0 にチャンク 9 1 が登録されている場合 (S 5 0 : YES) 、コントローラ 3 0 は、プール部 6 0 内の各 R A I D グループ 9 0 毎に、ラウンドロビン方式でチャンクを選択し、以下のステップ S 5 1 ~ S 5 4 を実行する。

【 0 1 2 5 】

図 1 4 に示す例で選択方法を説明すると、第 1 R A I D グループ 9 0 (# 1) から一つのチャンク 9 1 (1 - 1) を選択した後、第 2 R A I D グループ 9 0 (# 2) から別の一つのチャンク 9 1 (2 - 1) を選択し、さらに、第 3 R A I D グループ 9 0 (# 3) からさらに別の一つのチャンク 9 1 (3 - 1) を選択する。選択されたチャンク 9 1 (1 - 1) 、 9 1 (2 - 1) , 9 1 (3 - 1) について、後述する S 5 2 ~ S 5 4 をそれぞれ実行

10

20

30

40

50

する。チャンク 9 1 (1 - 2) , 9 1 (2 - 2) , 9 1 (3 - 2) のセットについても、S 5 2 ~ S 5 4 を実行する。さらに、チャンク 9 1 (1 - 3) , 9 1 (2 - 3) , 9 1 (3 - 3) のセットについても、S 5 2 ~ S 5 4 を実行する。以下同様である。

【 0 1 2 6 】

コントローラ 3 0 は、フォーマット待ちチャンクキュー Q 1 0 から、対象 R A I D グループ 9 0 のチャンク 9 1 を一つ選択し、フォーマット待ちチャンクキュー Q 1 0 からデキューさせる (S 5 2) 。コントローラ 3 0 は、選択されたチャンク 9 1 を処理完了待ちチャンクキュー Q 2 0 にエンキューし (S 5 3) 、選択されたチャンク 9 1 についてのフォーマットジョブを実行する (S 5 4) 。フォーマットジョブの詳細については、図 1 6 と共に後述する。

10

【 0 1 2 7 】

フォーマットジョブが完了すると、コントローラ 3 0 は、フォーマット済のチャンク 9 1 を、処理完了待ちチャンクキュー Q 2 0 からデキューし (S 5 5) 、そのフォーマット済のチャンク 9 1 を未割当てチャンクキュー Q 3 0 にエンキューする (S 5 6) 。

【 0 1 2 8 】

図 1 6 は、図 1 4 中に S 5 4 で示されるフォーマットジョブの詳細を示すフローチャートである。コントローラ 3 0 は、処理対象のチャンク 9 1 の識別番号を取得し (S 6 0) 、対象チャンク 9 1 をフォーマットする範囲を決定する (S 6 1) 。そして、コントローラ 3 0 は、対象チャンク 9 1 についてのエクステンロックを取得する (S 6 2) 。これにより、対象チャンク 9 1 が別のプログラムによって使用されるのを防止する。

20

【 0 1 2 9 】

コントローラ 3 0 は、対象チャンク 9 1 について所定サイズ毎に、後述する S 6 4 ~ S 7 0 を実行する (S 6 3) 。つまり、コントローラ 3 0 は、対象チャンク 9 1 を、所定サイズの単位領域毎にフォーマットするようになっている。

【 0 1 3 0 】

コントローラ 3 0 は、データ用のキャッシュセグメントを確保し (S 6 4) 、続いて、パリティ用のキャッシュセグメントを確保する (S 6 5) 。キャッシュセグメントとは、キャッシュメモリ 3 4 0 の記憶領域を管理する単位である。

【 0 1 3 1 】

コントローラ 3 0 は、論理アドレスを算出し (S 6 6) 、ゼロデータの生成を要求し (S 6 7) 、さらに、パリティの生成を要求する (S 6 8) 。ゼロデータは、S 6 4 で確保されたキャッシュセグメントを用いて生成される。パリティは、S 6 5 で確保されたキャッシュセグメントを用いて生成される。コントローラ 3 0 は、データ用に確保されたキャッシュセグメントを開放させ (S 6 9) 、さらに、パリティ用に確保されたキャッシュセグメントも開放させる (S 7 0) 。

30

【 0 1 3 2 】

コントローラ 3 0 は、対象チャンク 9 1 のエクステンロックを開放し (S 7 1) 、対象チャンク 9 1 についてのフォーマット完了を確認してから (S 7 2 : YES) 、図 1 5 の処理に戻る。

【 0 1 3 3 】

図 1 7 は、仮想ボリューム 5 0 V を作成する処理を示すフローチャートである。コントローラ 3 0 は、プール部 6 0 の状態が正常であるか否かを判定する (S 8 0) 。プール部 6 0 に障害が発生している場合 (S 8 0 : NO) 、仮想ボリューム 5 0 V を生成することはできないため、コントローラ 3 0 は、エラー処理を実行する (S 8 1) 。エラー処理では、例えば、プール部 6 0 に異常が生じているために仮想ボリューム 5 0 V を生成できない旨を、ユーザに通知する。

40

【 0 1 3 4 】

プール部 6 0 が正常な場合 (S 8 0 : YES) コントローラ 3 0 は、仮想ボリュームインデックス 1 5 2 の状態を " 無効 " から " 処理中 " に変更させる (S 8 2) 。予め用意されている仮想ボリュームの識別番号には、仮想ボリュームの状態の初期値として " 無効 " が予め

50

設定されている。仮想ボリューム50Vの生成中では、その状態は”無効”から”処理中”に変化する。仮想ボリューム50Vの生成が完了すると、その状態は”処理中”から”有効”に変化する。

【0135】

コントローラ30は、状態が”処理中”に変更された仮想ボリューム50Vについて、仮想アドレスインデックス153を作成し(S83)、さらに、その仮想アドレスインデックス153に対応付けられる仮想アドレスブロック154を作成する(S84)。さらに、コントローラ30は、仮想アドレスブロック154に対応付けられるページアドレス情報155を作成する(S85)。

【0136】

コントローラ30は、プール内RAIDグループインデックス110を更新し(S86)、仮想ボリュームインデックス152の状態を”処理中”から”有効”に変更する(S87)。コントローラ30は、仮想ボリューム50Vを正常に作成することができたか否かを確認し(S88)、仮想ボリューム50Vを正常に作成できた場合には(S88:YES)、本処理を終了させる。仮想ボリュームを正常に作成できなかった場合(S88:NO)、エラー処理が行われる(S89)。エラー処理では、例えば、仮想ボリューム50Vを正常に作成できなかった旨をユーザに通知する。

【0137】

なお、便宜上、仮想ボリューム50Vを正常に作成できたか否かを最後に判定するかのよう説明したが、実際には、各テーブル101-104の作成時にそれぞれ正常に作成できたか否かが判定される。そして、正常に作成できなかった場合には、エラー処理が行われる。

【0138】

図17に示す処理を行うことにより、図5に示すテーブル群が作成され、仮想ボリューム50V内の各仮想的記憶領域500は、初期設定用のページ92にそれぞれ対応付けられる。従って、仮想ボリューム50Vを正常に作成できた時点で、各仮想的記憶領域500と実ページ92との対応付けを正常に行えることが確認される。

【0139】

図16は、ライトコマンドの処理を示すフローチャートである。コントローラ30は、ホスト20からコマンドを受領すると、そのコマンドが仮想ボリューム50Vを対象とするライトコマンドであるか否かを判定する(S100)。

【0140】

そのコマンドが、通常ボリューム50Nを対象とするライトコマンドである場合(S100:NO)、コントローラ30は、通常のライト処理を実行する(S101)。通常のライト処理では、例えば、ホスト20から受領したライトデータをキャッシュメモリ340に書き込み、キャッシュメモリ340へのライトデータ書き込みが完了した時に、ホスト20に処理完了を通知する。その後、適切なタイミングを見計らって、キャッシュメモリ340に記憶されたライトデータを、記憶装置40に書き込む。

【0141】

仮想ボリューム50Vを対象とするライトコマンドを受領した場合(S100:YES)、コントローラ30は、移動先管理テーブル160を参照することにより、書き込み対象の仮想ページ500に割り当てられている実ページ92に移動先が設定されているか否かを判定する(S102)。

【0142】

ライトデータの書き込まれる実ページ92に移動先が設定されている場合(S102:YES)、コントローラ30は、移動先のRAIDグループ90に利用可能なチャック91が有るか否かを判定する(S103)。利用可能なチャック91とは、既に作成されており、仮想ボリューム50Vに対応付けられているチャックである。このように、既に作成されているチャック91を利用して、ページを移動させることにより、割り当てられるチャックの数を少なくできる。つまり、記憶容量を効率的に使用できる。

10

20

30

40

50

【 0 1 4 3 】

利用可能なチャンク 9 1 が無い場合 (S103:NO)、コントローラ 3 0 は、初期状態になっている新規チャンク 9 1 が有るか否かを判定する (S 1 0 4)。つまり、コントローラ 3 0 は、直ちに使用できる新規チャンク 9 1 が、移動先 R A I D グループ 9 0 内に用意されているか否かを判定する。

【 0 1 4 4 】

使用可能な新規チャンク 9 1 が移動先 R A I D グループ 9 0 内に無い場合 (S104:NO)、コントローラ 3 0 は、エラー処理を行う (S 1 0 7)。エラー処理では、例えば、「記憶容量不足です。」「R A I D グループを作成してください。」等のエラーメッセージを、管理サーバ 7 0 を介してユーザに通知する。

10

【 0 1 4 5 】

移動先 R A I D グループ 9 0 内に利用可能なチャンク 9 1 が存在する場合 (S103:YES)、コントローラ 3 0 は、S 1 0 4 をスキップして S 1 0 5 に移る。

【 0 1 4 6 】

コントローラ 3 0 は、書き込み対象の実ページ 9 2 に記憶されるデータを、移動先 R A I D グループ 9 0 内に設けられる新ページ 9 2 に移動させるべく、チャンク割当て変更処理 (S 1 0 5) と、ページ割当て変更処理 (S 1 0 6) とを実行する。チャンク割当て変更処理及びページ割当て変更処理については、後述する。その後、コントローラ 3 0 は、図 1 9 の S 1 1 4 に移る。

【 0 1 4 7 】

20

つまり、S 1 0 2 ~ S 1 0 6 は、移動先の設定されている実ページ 9 2 についてライトコマンドが発行された場合に、移動先チャンク 9 1 内の移動先ページ 9 2 を、ライトコマンドの対象である仮想ページ 5 0 0 に対応付ける。換言すれば、書き込み対象の仮想ページ 5 0 0 の割当て先を、移動元ページ 9 2 から移動先ページ 9 2 に切り替える。

【 0 1 4 8 】

図 1 9 は、図 1 8 に続くフローチャートである。書き込み対象の実ページ 9 2 に移動先 R A I D グループ 9 0 が設定されていない場合 (S102:NO)、コントローラ 3 0 は、ライトコマンドで指定された仮想ボリューム 5 0 V について、現在使用中のチャンク 9 1 が有るか否かを判定する (S 1 0 8)。使用中のチャンク 9 1 が有る場合 (S108:YES)、後述の S 1 1 2 に移行する。使用中のチャンク 9 1 が無い場合 (S108:NO)、コントローラ 3 0

30

【 0 1 4 9 】

通常の場合、図 1 5 に示すフォーマット処理によって、新規チャンクは初期状態に設定されている。それにもかかわらず、新規チャンク 9 1 が初期状態になっていない場合 (S109:NO)、エラー処理が行われる (S 1 1 0)。エラー処理では、例えば、初期状態のチャンク 9 1 が存在しない旨を、管理サーバ 7 0 を介してユーザに通知する。

【 0 1 5 0 】

新規チャンク 9 1 が初期状態になっている場合 (S109:YES)、コントローラ 3 0 は、チャンク割当て変更処理を行う (S 1 1 1)。チャンク割当て変更処理の詳細は、図 2 1 で述べるが、先に簡単に説明すると、コントローラ 3 0 は、未割当てチャンクキュー Q 3 0 から一つのチャンク 9 1 を選択して仮想ボリューム 5 0 V に対応付け、そのチャンク 9 1 の状態を " 割当済み (使用中) " に変更等する。

40

【 0 1 5 1 】

コントローラ 3 0 は、チャンク 9 1 内の各ページ 9 2 のうち、使用しようとするページ 9 2 が初期状態になっているか否かを判定する (S 1 1 2)。使用しようとするページ 9 2 が初期状態の場合 (S112:YES)、ページ割当て変更処理が行われる (S 1 1 3)。

【 0 1 5 2 】

ページ割当て変更処理の詳細は、図 2 2 で後述する。簡単に説明すると、ページ割当て変更処理では、仮想ボリューム 5 0 V に割り当てられるページ 9 2 の状態を " 割当済み (

50

使用中) ”に変更し、仮想アドレスブロック 1 5 4 及びページアドレス情報 1 5 5 を更新させる。S 1 1 3 の後、S 1 1 4 に移る。

【 0 1 5 3 】

使用しようとするページ 9 2 が初期状態ではない場合 (S 1 1 2 : N 0)、つまり、使用しようとするページ 9 2 が初期設定用のページではない場合、S 1 1 3 はスキップされて、S 1 1 4 に移る。

【 0 1 5 4 】

コントローラ 3 0 は、ライトデータを記憶させるためのキャッシュセグメントを確保し (S 1 1 4)、さらに、ライトデータを転送するための D M A (Direct Memory Access) 転送リストを作成する (S 1 1 5)。そして、コントローラ 3 0 は、保証コードのアドレス部分 (L A) を算出する (S 1 1 6)。

10

【 0 1 5 5 】

図 2 0 を参照する。コントローラ 3 0 は、ホスト 2 0 から受領したライトデータをキャッシュメモリ 3 4 0 に D M A 転送させる (S 1 1 7)。キャッシュメモリ 3 4 0 にライトデータを記憶させた後で、コントローラ 3 0 は、ライトコマンドの処理が完了した旨をホスト 2 0 に通知する (S 1 1 8)。

【 0 1 5 6 】

キャッシュメモリ 3 4 0 へのライトデータ書き込み完了後に、ホスト 2 0 に処理完了を通知する方式を非同期方式と呼ぶ。これに対し、記憶装置 4 0 にライトデータが書き込まれるのを待ってから、ホスト 2 0 に処理完了を通知する方式を同期方式と呼ぶ。非同期方式または同期方式のいずれを用いても良い。

20

【 0 1 5 7 】

コントローラ 3 0 は、ライトデータに設定されている論理アドレスを、記憶装置 4 0 に記憶させるための物理アドレスに変換する (S 1 1 9)。コントローラ 3 0 は、キャッシュセグメントを確保する (S 1 2 0)。コントローラ 3 0 は、パリティを生成するために必要な旧データを記憶装置 4 0 から読み出して (S 1 2 1)、S 1 2 0 で確保したキャッシュセグメントに格納させる。

【 0 1 5 8 】

コントローラ 3 0 は、ホスト 2 0 から受領したライトデータと記憶装置 4 0 から読み出した旧データとに基づいて、新たなパリティを算出する (S 1 2 2)。コントローラ 3 0 は、キャッシュメモリ 3 4 0 に記憶されたライトデータを、記憶装置 4 0 (図 2 0 中、ディスクと表示) に転送して記憶させる (S 1 2 3)。

30

【 0 1 5 9 】

コントローラ 3 0 は、キャッシュメモリ 3 4 0 に記憶されているライトデータの状態を ” ダーティ ” から ” クリーン ” に変更する (S 1 2 4)。 ” ダーティ ” 状態とは、キャッシュメモリ 3 4 0 だけに記憶されている状態を示す。 ” クリーン状態 ” とは、記憶装置 4 0 に書き込まれた状態を示す。最後に、コントローラ 3 0 は、S 1 1 4 及び S 1 2 0 で確保されたキャッシュセグメントを開放し、本処理を終了する (S 1 2 5)。

【 0 1 6 0 】

上述の通り、書き込み対象の実ページ 9 2 (移動元ページ) に移動先が設定されている場合、S 1 0 2 ~ S 1 0 6 によって、その移動元ページ 9 2 に対応付けられていた仮想ページ 5 0 0 は、移動先ページ 9 2 に新たに対応付けられる。従って、S 1 1 9 ~ S 1 2 5 では、移動先ページ 9 2 にライトデータが書き込まれる。これにより、ライトコマンドを処理する中で、移動対象データを移動元ページ 9 2 から移動先ページ 9 2 に移動させることができる。

40

【 0 1 6 1 】

図 2 1 は、図 1 8 の S 1 0 5 及び図 1 9 の S 1 1 1 に示されるチャンク割当て変更処理の詳細を示すフローチャートである。コントローラ 3 0 は、処理対象チャンク 9 1 の状態を所定の状態に変更させる (S 1 3 0)。

【 0 1 6 2 】

50

例えば、新規チャック 9 1 を仮想ボリューム 5 0 V に割り当てる場合、新規チャック 9 1 の状態は " 未割当て (未使用) " から " 割当済み (使用中) " に変更される。また、例えば、仮想ボリューム 5 0 V に割り当てられているチャック 9 1 を開放する場合、そのチャック 9 1 の状態は " 割当済み (使用中) " から " フォーマット待ち " に変更される。

【 0 1 6 3 】

コントローラ 3 0 は、チャック 9 1 の状態を変更した後で、そのチャック 9 1 に対応する仮想アドレスインデックス 1 5 3 を更新させる (S 1 3 1) 。さらに、コントローラ 3 0 は、チャックインデックス 1 2 0 を更新させ (S 1 3 2) 、続いて、プール内 R A I D グループインデックス 1 1 0 を更新させる (S 1 3 3) 。

【 0 1 6 4 】

図 2 2 は、図 1 8 中の S 1 0 6 及び図 1 9 の S 1 1 3 に示されるページ割当て変更処理の詳細を示すフローチャートである。コントローラ 3 0 は、処理対象ページ 9 2 の状態を所定の状態に変更させる (S 1 4 0) 。例えば、新規ページ 9 2 を仮想ボリューム 5 0 V 内の仮想的記憶領域 5 0 0 に割り当てる場合、その新規ページ 9 2 の状態は " 未割当て (未使用) " から " 割当済み (使用中) " に変更される。また、例えば、仮想ボリューム 5 0 V に割当済みのページ 9 2 を開放する場合、その割当済みページ 9 2 の状態は " 割当済み (使用中) " から " フォーマット待ち " に変更される。仮想ボリューム 5 0 V に新たなページ 9 2 を割り当てる方法の詳細については、図 2 3 で後述する。

【 0 1 6 5 】

ページ状態を変更した後、コントローラ 3 0 は、処理対象ページ 9 2 に対応する仮想アドレスブロック 1 5 4 を更新させ (S 1 4 1) 、さらに、処理対象ページ 9 2 に対応するページアドレス情報 1 5 5 を更新させる (S 1 4 2) 。そして、コントローラ 3 0 は、管理情報 (図 6 に示すテーブル群) を、管理情報退避領域 (図 1 1 の S 2 3 参照) に退避させる (S 1 2 3) 。

【 0 1 6 6 】

図 2 3 は、図 2 2 の S 1 4 0 で示される処理の一部を示すフローチャートである。上述の通り、ホスト 2 0 からのライトコマンドに応じて、仮想ボリューム 5 0 V にチャック単位で実記憶領域が割り当てられ、その割り当てられたチャック 9 1 内の実ページ 9 2 が順番に使用されていく。割り当てられたチャック 9 1 を使い切ると、新たなチャック 9 1 が仮想ボリューム 5 0 V に割り当てられる。以下、図 2 3 を参照しながら、ページ 9 2 の使用方法を説明する。

【 0 1 6 7 】

コントローラ 3 0 は、各ページ 9 2 の状態を管理するためのページ状態管理テーブル 1 5 8 を備える。そのテーブル 1 5 8 は、例えば、チャック識別番号の欄 1 5 8 0 と、ページ識別番号の欄 1 5 8 1 と、ページの状態を示す欄 1 5 8 2 とを備える。

【 0 1 6 8 】

ページ状態欄 1 5 8 2 には、例えば、 " 使用中 (あるいは割当済み) " 、 " 未使用 (あるいは未割当て) " 、 " 開放 (あるいはフォーマット待ち) " 等のような予め用意されている状態のうちのいずれか一つの状態が設定される。なお、説明の便宜上、ページ状態管理テーブル 1 5 8 は、ページアドレス情報 1 5 5 と別体の情報であるかのように示すが、実際には、ページアドレス情報 1 5 5 だけで各ページの状態を管理できる。

【 0 1 6 9 】

コントローラ 3 0 は、テーブル 1 5 8 を参照し、現在使用中のチャック 9 1 内に未使用ページ 9 2 が有る場合、その未使用ページ 9 2 を使用する (S 1 4 0 0) 。現在使用中のチャック 9 1 内に未使用ページ 9 2 が無い場合、コントローラ 3 0 は、テーブル 1 5 8 を参照し、現在使用中のチャック 9 1 内に開放ページ 9 2 が有る場合、その開放されたページ 9 2 を使用する (S 1 4 0 1) 。

【 0 1 7 0 】

現在使用中のチャック 9 1 内に未使用ページ 9 2 も開放ページ 9 2 もいずれも存在しない場合、コントローラ 3 0 は、テーブル 1 5 8 を参照し、使用済チャック 9 1 内の開放ペ

10

20

30

40

50

ージを使用する (S 1 4 0 2)。つまり、コントローラ 3 0 は、対象の仮想ボリューム 5 0 V について既に使用されたチャンク 9 1 から、開放されたページ 9 2 を検出して再使用する。

【 0 1 7 1 】

なお、使用済チャンク 9 1 内に開放ページ 9 2 が無い場合、コントローラ 3 0 は、図 1 9 で述べたように、未使用のチャンク 9 1 を仮想ボリューム 5 0 V に対応付け、そのチャンク 9 1 の先頭ページ 9 2 を仮想ボリューム 5 0 V に割り当てる。

【 0 1 7 2 】

図 2 4 は、リードコマンドの処理を示すフローチャートである。コントローラ 3 0 は、ホスト 2 0 から受領したコマンドが、仮想ボリューム 5 0 V を対象とするリードコマンドであるか否かを判別する (S 1 5 0)。

10

【 0 1 7 3 】

ホスト 2 0 から受領したコマンドが通常ボリューム 5 0 N へのリードコマンドである場合 (S150:NO)、コントローラ 3 0 は、通常のリード処理を実行する (S 1 5 1)。例えば、コントローラ 3 0 は、ホスト 2 0 から要求されたデータがキャッシュメモリ 3 4 0 に記憶されているか否かを判定する。要求されたデータがキャッシュメモリ 3 4 0 上に存在する場合、コントローラ 3 0 は、キャッシュメモリ 3 4 0 からデータを読み出してホスト 2 0 に送信する。ホスト 2 0 の要求するデータがキャッシュメモリ 3 4 0 上に存在しない場合、コントローラ 3 0 は、記憶装置 4 0 からデータを読み出してキャッシュメモリ 3 4 0 に記憶させ、そのデータをホスト 2 0 に送信する。

20

【 0 1 7 4 】

ホスト 2 0 から発行されたコマンドが、仮想ボリューム 5 0 V からデータを読み出すためのリードコマンドである場合 (S150:YES)、コントローラ 3 0 は、リード対象の仮想ボリューム 5 0 V についてエクステンロックを取得する (S 1 5 2)。

【 0 1 7 5 】

リードコマンドは、データの読出し先の論理アドレスを指定する。コントローラ 3 0 は、指定された論理アドレスに対応する仮想的記憶領域 5 0 0 を検出し、図 6 に示すテーブル群を参照して、その仮想的記憶領域 5 0 0 に割り当てられているページ 9 2 の状態を取得する。コントローラ 3 0 は、リード対象のページ 9 2 の状態が " 初期状態 " であるか否かを判定する (S 1 5 3)。

30

【 0 1 7 6 】

リード対象のページ 9 2 が初期状態の場合 (S153:YES)、コントローラ 3 0 は、ホスト 2 0 に送信すべきヌルデータがキャッシュメモリ 3 4 0 に記憶されているか否かを判定する (S 1 6 0)。図 1 1 の S 3 0 で述べたように、プール部 6 0 の作成時に、キャッシュメモリ 3 4 0 の所定のキャッシュセグメントにヌルデータが予め記憶される。従って、通常の場合、S 1 6 0 では " Y E S " と判定されて、後述の S 1 6 1 に移る。キャッシュメモリ 3 4 0 にヌルデータが記憶されていない場合 (S160:NO)、後述の S 1 5 5 に移り、所定チャンクの所定ページ (例えば、初期設定用ページ) から、ヌルデータが読み出されて、ホスト 2 0 に送信される (S 1 5 5 ~ S 1 5 9 , S 1 7 1 ~ S 1 7 5)。

【 0 1 7 7 】

40

S 1 5 3 に戻る。リード対象のページ 9 2 が初期状態では無い場合 (S153:NO)、つまり、リード対象ページにライトデータが書き込まれている場合、コントローラ 3 0 は、リード対象データに関するパリティを算出する (S 1 5 4)。そして、コントローラ 3 0 は、キャッシュセグメントを確保し (S 1 5 5)、第 2 通信回路 3 2 0 に向けてリード要求を発行する (S 1 5 6)。コントローラ 3 0 は、論理アドレスを物理アドレスに変換し (S 1 5 7)、保証コードのアドレス部分 (L A) を算出する (S 1 5 8)。コントローラ 3 0 は、第 2 通信回路 3 2 0 を介して、記憶装置 4 0 からキャッシュメモリ 3 4 0 にリード対象データを転送させる (S 1 5 9)。

【 0 1 7 8 】

図 2 5 に移る。コントローラ 3 0 は、キャッシュメモリ 3 4 0 から第 1 通信回路 3 1 0

50

にDMA転送させるための、DMA転送リストを設定する(S161)。続いて、コントローラ30は、キャッシュメモリ340上に記憶されているデータを、第1通信回路310を介してホスト20に送信させる(S162)。そして、コントローラ30は、リードコマンドに係る仮想ページ500に対応付けられている実ページ92に、移動先RAIDグループ90が設定されているか否かを判定する(S163)。リード対象の実ページ92に移動先RAIDグループ90が設定されていない場合(S163:NO)、図26に示すフローチャートに移る。

【0179】

リード対象の実ページ92に移動先RAIDグループ90が設定されている場合(S163:YES)、コントローラ30は、移動先RAIDグループ90に利用可能なチャンク91が有るか否かを判定する(S164)。利用可能なチャンク91が無い場合(S164:NO)、移動先RAIDグループ90に初期状態の新規チャンク91が有るか否かを判定する(S1653)。

10

【0180】

移動先RAIDグループ90内に初期状態の新規チャンク91が無い場合(S165:NO)、コントローラ30は、図26に示すフローチャートに移る。移動対象のデータを移動先RAIDグループ90に移動させることができないためである。

【0181】

一方、移動先RAIDグループ90内に利用可能なチャンク91が有る場合(S164:YES)、または、移動先RAIDグループ90が初期状態の新規チャンク91を有する場合(S165:YES)、のいずれかの場合には、コントローラ30は、リードコマンドの処理が完了した旨をホスト20に通知する(S166)。

20

【0182】

さらに、コントローラ30は、読み出したデータを移動先RAIDグループ90内に移動させるべく、以下のステップを実行する。コントローラ30は、図21で述べたチャンク割当て処理(S167)と、図22で述べたページ割当て変更処理(S168)とを実行する。これにより、リード対象の仮想ページ500は、移動元ページ92から、移動先RAIDグループ90内の新たなページ92に対応付けられる。

【0183】

コントローラ30は、キャッシュメモリ340上のデータを、記憶装置40に転送させて書き込ませる(S169)。コントローラ30は、キャッシュセグメントを開放し(S170)、さらに、エクステンロックを開放する(S171)。

30

【0184】

一方、リード対象の実ページ92に移動先RAIDグループ90が設定されていない場合(S163:NO)、図26に移る。

【0185】

コントローラ30は、S155で確保したキャッシュセグメントを開放し(S172)、エクステンロックも開放する(S173)。最後に、コントローラ30は、リードコマンドの処理が完了した旨をホスト20に通知し(S174)、本処理を終了する。

【0186】

40

このように、コントローラ30は、リードコマンドを処理する中で、移動対象のデータを予め設定される移動先に移動させる。そこで、移動先を決定する方法を説明する。

【0187】

図27は、移動先を決定するための処理を示すフローチャートである。データを再配置させるための処理と呼び変えてもよい。本処理は、例えば、予め設定される所定の周期毎に実行される。または、ユーザからの指示に応じて本処理を実行することもできる。

【0188】

コントローラ30は、移動先管理テーブル160に設定されている移動先RAIDグループ90に関する情報をクリアさせる(S180)。つまり、新たなデータ移動計画を作成する前に、前回作成されたデータ移動計画をリセットする。以下、各プール部60毎に

50

(S181)、負荷を分散させる処理(S182)と、使用容量を平均化させる処理(S183)とが行われる。

【0189】

図28は、図27のS182に示される負荷分散処理のフローチャートである。コントローラ30は、予め設定される条件に従って、移動元RAIDグループ90と、移動先RAIDグループ90とを選択する(S190)。

【0190】

例えば、コントローラ30は、使用容量が50%以下であるRAIDグループの中から、負荷が最も低いRAIDグループを移動先RAIDグループとして選択する。さらに、コントローラ30は、RAIDグループの中から、負荷が最も高いRAIDグループを移動元RAIDグループとする。

10

【0191】

もしも、上記条件を満たすRAIDグループが存在せず、移動元RAIDグループまたは移動先RAIDグループのいずれか一つでも決定することができなかつた場合、コントローラ30は、図28に示すループを抜ける。

【0192】

コントローラ30は、移動元チャンクを決定する(S191)。今までの説明では、ページ単位でデータを移動させる場合を述べたが、チャンク単位でデータを移動させることもできる。

【0193】

コントローラ30は、移動元RAIDグループに所属する各チャンクの中で負荷が高いチャンクから順番に移動元チャンクとして決定する。移動元チャンクに指定されるチャンクの負荷の合計は、RAIDグループ間の負荷差の50%以内とする。

20

【0194】

データの移動によって、移動元RAIDグループの負荷が低下し、移動先RAIDグループの負荷は増大する。RAIDグループ間の負荷差を考慮して移動元チャンクを選定することにより、データ移動の結果、移動元RAIDグループの負荷が移動先RAIDグループの負荷よりも小さくなるのを防止できる。

【0195】

コントローラ30は、上記方法で選ばれた移動元RAIDグループ、移動先グループ及び移動元チャンクを、移動先管理テーブル160に記憶させる(S192)。コントローラ30は、S190で選択された移動元RAIDグループ及び移動先RAIDグループを、処理対象から除外する(S193)。

30

【0196】

コントローラ30は、プール部60内のRAIDグループ数が1以下になるか、または、上記条件を満たすRAIDグループが見つからなくなるまで、本処理を繰り返す。

【0197】

図29は、図27のS183に示される負荷分散処理のフローチャートである。コントローラ30は、予め設定される条件に従って、移動元RAIDグループ90と、移動先RAIDグループ90とを選択する(S200)。

40

【0198】

例えば、コントローラ30は、使用容量が50%以下のRAIDグループの中から、使用容量が最も小さいRAIDグループを移動先RAIDグループとして選択する。さらに、コントローラ30は、RAIDグループの中で使用容量が最も大きいRAIDグループを移動元RAIDグループとして選択する。

【0199】

もしも、上記条件を満たすRAIDグループが存在せず、移動元RAIDグループまたは移動先RAIDグループのいずれか一つでも決定することができなかつた場合、コントローラ30は、図29に示すループを抜ける。

【0200】

50

コントローラ30は、移動元チャンクを決定する(S191)。例えば、コントローラ30は、移動元RAIDグループに所属するチャンクのうち未だ移動先が決定されておらず、かつ、負荷が低いチャンクから順番に、移動元チャンクに決定する。負荷が0のチャンクは選択対象から外される。

【0201】

コントローラ30は、移動元チャンクを決定する(S201)。移動元チャンクとして選択されるチャンクの数、RAIDグループ間の未使用チャンク数差の50%以内に保たれる。

【0202】

さらに、移動先RAIDグループの負荷が、移動元RAIDグループの負荷よりも高い場合、移動元チャンクとして選定されるチャンクの数、移動先RAIDグループにおいて移動先が指定されているチャンクの数以下に保持される。

10

【0203】

このように構成される本実施例では、ホストアクセスが発生する前に、仮想ボリューム50Vに割り当てられている実ページ92のデータを移動させるためのデータ移動計画を作成して保存する。本実施例では、移動対象データについてホストアクセスが発生すると、そのホストアクセスに係わるコマンド処理の中で、移動対象のデータを移動させることができる。つまり、データマイグレーション用の特別なプログラムを、コマンド処理のためのプログラムと別に実行するのではなく、コマンドを処理するための一連の流れの中でデータ移動を行う。

20

【0204】

従って、本実施例では、比較的簡易な構成で、比較的効率的にデータを移動させることができ、プール部60内の実記憶領域を略均等に使用することができ、特定のRAIDグループにアクセスが偏ったりするのを抑制できる。これにより、記憶制御装置10の応答性能の低下を防止できる。

【0205】

上述のデータ移動に関する本実施例の効果は、本実施例の基本的構成と結合することにより発揮される。本実施例の基本的構成とは、仮想ボリューム50Vにチャンク単位で実記憶領域を割り当て、かつ、一つのチャンク91は一つの仮想ボリューム50Vに専属させる構成である。

30

【0206】

本実施例では、図30に示すように、物理的記憶領域であるRAIDグループ90を効率的に使用することができる。図30は、本発明の効果を模式的に示す説明図である。図30(a)は、本発明を適用しない場合を示し、図30(b)は、本発明を適用した場合を示す。

【0207】

通常ボリューム50Nの場合を先に説明する。通常ボリューム50Nの場合は、RAIDグループ90内の連続した記憶領域を使用することができる。従って、一つのストライプ列に、複数の通常ボリューム50Nに関するデータが混在することは無い。

【0208】

仮想ボリューム50Vの場合、必要に応じて実記憶領域が割り当てられ、データは離散的に管理される。もしも、仮想ボリュームにページ単位で実記憶領域を割り当てる場合は、同一のストライプ列に複数の仮想ボリュームに関するデータが混在しないように制御する必要がある。何故なら、一つのストライプ列に複数のボリュームが混在すると、パリティ生成等の処理が複雑化し、データ入出力時のオーバーヘッドが増大して、記憶制御装置の性能が低下するためである。

40

【0209】

そこで、図30(a)に示すように、各ページの先頭をストライプ列の先頭に一致させるようにして、データを記憶する方法が考えられる。この方法によれば、横一列のストライプ内に、異なるボリュームのデータが混在するという事態は生じない。

50

【0210】

しかし、本実施例の記憶制御装置10は、RAIDグループ90を構成する記憶装置40の台数を自由に設定することができるため、ページサイズとストライプサイズとは必ずしも一致しない。ページサイズとストライプサイズとが不一致の場合において、ページの先頭をストライプ列の先頭に一致させるようにしてデータを配置すると、図30(a)に空白領域として示すように無駄な領域が発生する。従って、図30(a)に示す方法では、RAIDグループ90の有する実記憶領域を有効に利用することができず、その利用効率が低いという問題がある。

【0211】

そこで、本発明では、図30(b)に示すように、実記憶領域をページ単位で仮想ボリュームに割り当てるのではなく、複数ページを有するチャンク91単位で実記憶領域を仮想ボリュームに割り当てる。そして、上述のように、仮想ボリューム50Vに関するライトコマンドを受領するたびに、チャンク91内のページ92を連続的に使用する。チャンク91内の各ページ92は、同一の仮想ボリューム50Vに対応付けられる。異なる仮想ボリュームに関するページ92が同一のチャンク91内に混在することはない。このように、本実施例によれば、RAIDグループ90の有する実記憶領域を効率的に使用することができる。

10

【0212】

本実施例では、上述の通り、仮想ボリューム50Vの識別番号と通常ボリューム50Nの識別番号とを特に区別することなく、各ボリューム50V、50Nを連続番号で管理している。また、本実施例では、仮想ボリューム50Vのために使用されるRAIDグループ90と、通常ボリューム50Nが設けられるRAIDグループ90とを特に区別せずに、連続番号で管理している。従って、本実施例の記憶制御装置10は、仮想ボリューム50Vと通常ボリューム50Nとを比較的簡易な制御構造で共通に管理し、両方のボリューム50V、50Nを混在させることができる。

20

【0213】

本実施例では、複数のRAIDグループ90からチャンク91を順番に選択して、仮想ボリューム50Vに割り当てる。従って、プール部60内の各RAIDグループ90の負荷を均等にすることができる。

【0214】

本実施例では、図7、図17で述べたように、仮想ボリューム50Vを作成するときに各テーブル151-155間を関連づけ、仮想ボリューム50V内の全ての仮想的記憶領域500を、初期設定用のページ92に割り当てる。

30

【0215】

従って、本実施例では、仮想ボリューム50Vを正常に作成できた時点で、各仮想的記憶領域500と実ページ92との対応付けが正常に行われることを確認できる。つまり、本実施例では、ライトコマンドを受領する前に、仮想ボリューム50Vへのチャンク91及び初期設定用ページ92の仮の割り当てが完了している。これにより、仮想ボリューム50Vが正常に動作するか否かを、ライトコマンド受領前に事前に確認することができ、信頼性及び使い勝手が向上する。

40

【0216】

さらに、本実施例では、ライトコマンドを受領した場合に、ライトコマンドで指定される論理アドレスに対応する仮想的記憶領域500の対応付け先を、仮割り当てされた初期設定用ページ92から、所定チャンク91内の所定ページ92に切り替えるだけで済む。これにより、比較的速やかにライトコマンドを処理することができ、記憶制御装置10の応答性能を高めることができる。

【実施例2】

【0217】

図31、図32に基づいて第2実施例を説明する。本実施例は、第1実施例の変形例に該当する。従って、第1実施例との相違を中心に説明する。本実施例は、データ移動計画

50

を作成する場合に、ヌルデータのみが記憶されている実ページ 92 を仮想ボリューム 50V から開放し、未使用ページに戻す。

【0218】

図31は、本実施例による移動先決定処理を示すフローチャートである。本フローチャートは、図27に示すフローチャートに比べて、ゼロデータ（ヌルデータ）を削除するための処理（S210）をさらに備える。移動先決定処理では、最初に、ゼロデータのみを記憶する実ページ92を未使用ページに戻し、その後で、負荷分散処理及び使用容量平均化処理を実行する。

【0219】

図32は、ゼロデータ削除処理のフローチャートである。コントローラ30は、各RAIDグループ毎に、以下の処理を行う（S2101）。コントローラ30は、処理対象のRAIDグループ90の各ページ92のうち、ゼロデータのみが記憶されているページ92を検出する（S2102）。

【0220】

コントローラ30は、ゼロデータのみを記憶するページ92に対応する仮想ページ500を、図7で述べた初期化用の特定ページに対応付け（S2103）、ページ管理テーブル160を更新させる（S2104）。

【0221】

上述の構成を有する本実施例も第1実施例と同様の効果を奏する。さらに、本実施例では、データ移動計画を作成する場合に、ヌルデータのみが記憶されている実ページ92を開放して、未使用の実ページ92に戻すため、プール部60内の実記憶領域を有効に利用できる。つまり、本実施例では、データ移動計画を作成するたびに、プール部60内のページ92のうち無駄に使用されているページ92を開放できる。

【0222】

なお、本発明は、上述した実施形態に限定されない。当業者であれば、例えば、上記各実施例を適宜組み合わせる等のように、本発明の範囲内で、種々の追加や変更等を行うことができる。

【符号の説明】

【0223】

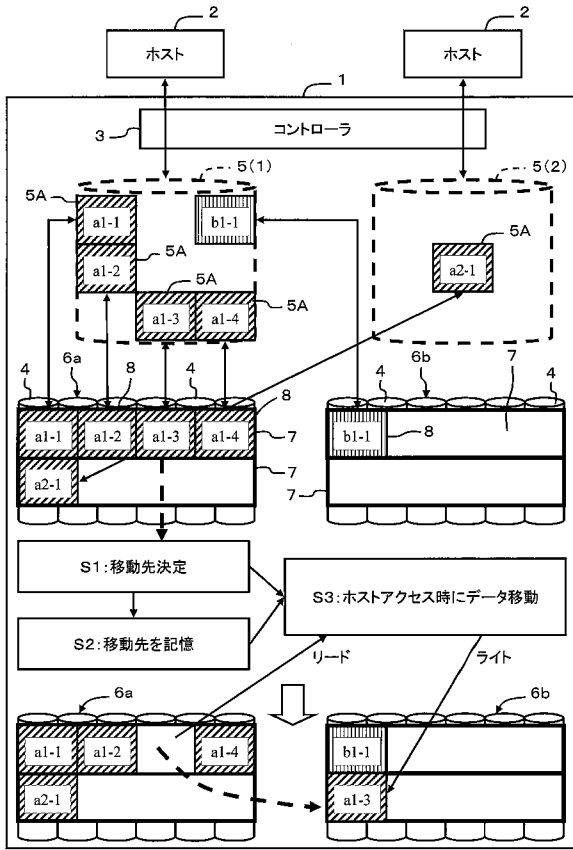
1：記憶制御装置、2：ホスト、3：コントローラ、4：記憶装置、5（1）、5（2）：仮想ボリューム、5A：仮想的記憶領域、6a、6b：RAIDグループ、7：チャンク、8：ページ、10記憶制御装置、20：ホスト、30：コントローラ、40：記憶装置、50V：仮想ボリューム、60：プール部、70：管理サーバ、90：RAIDグループ、91：チャンク、92：ページ、500：仮想ページ。

10

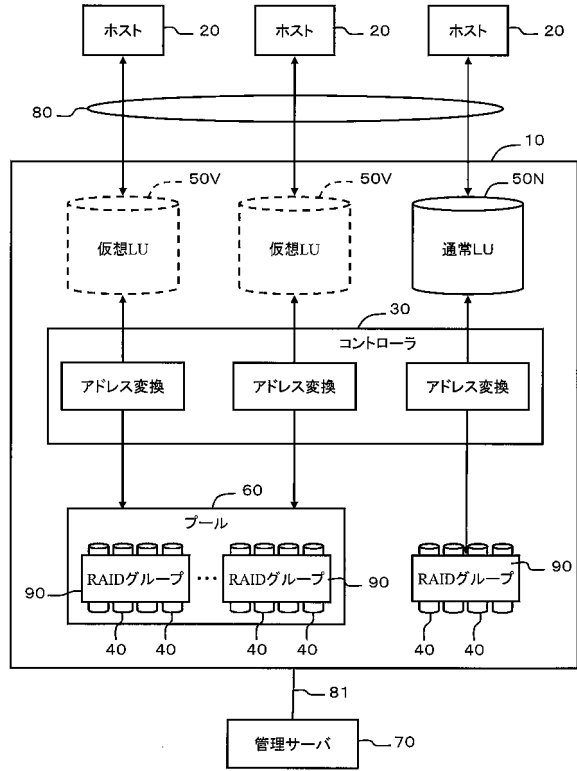
20

30

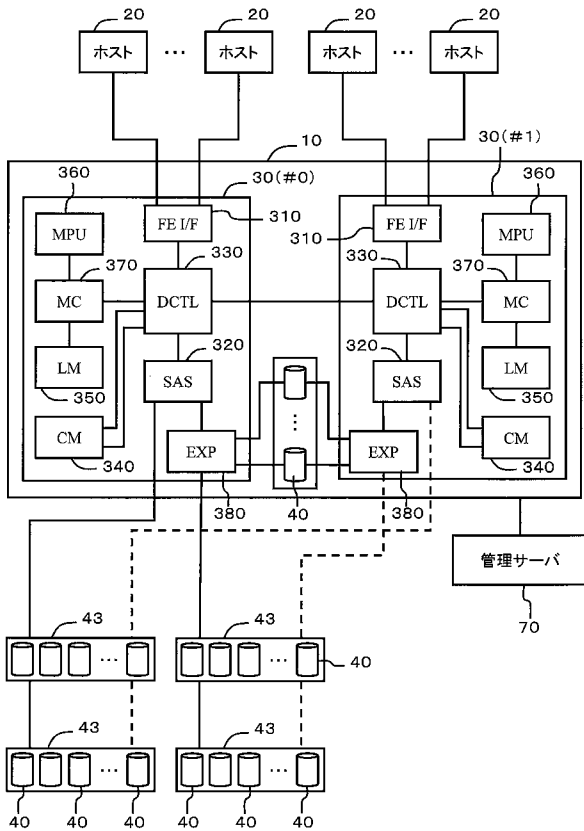
【図1】



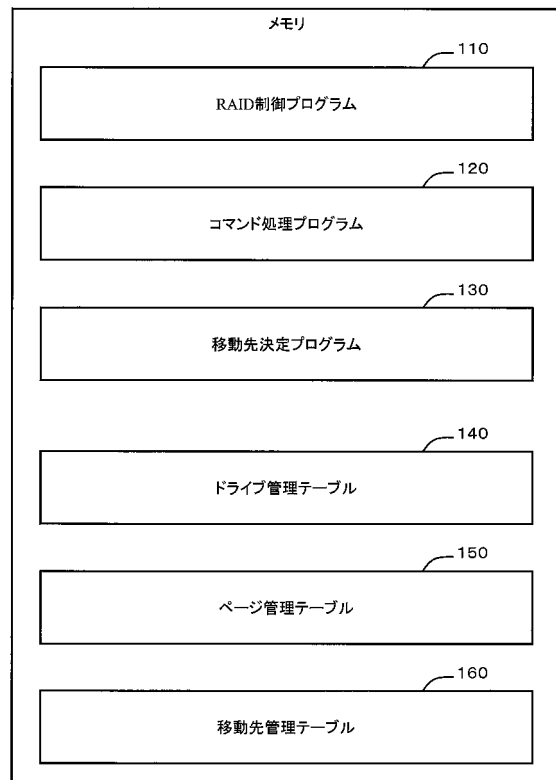
【図2】



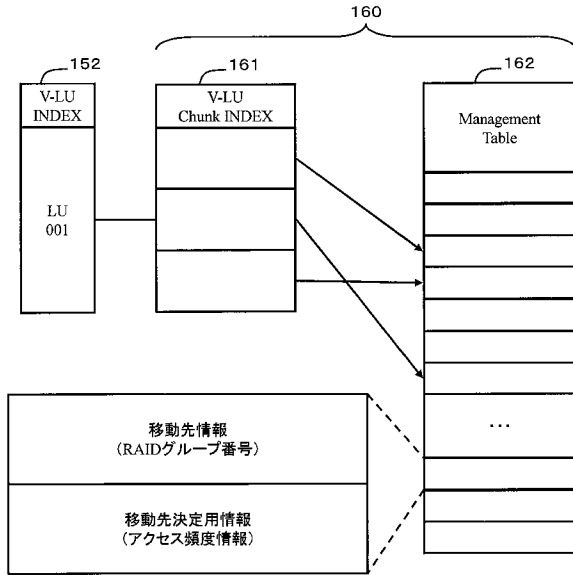
【図3】



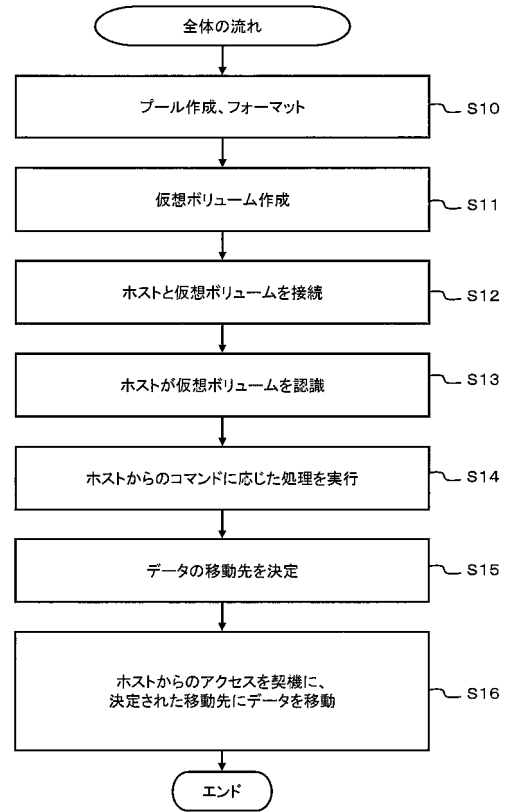
【図4】



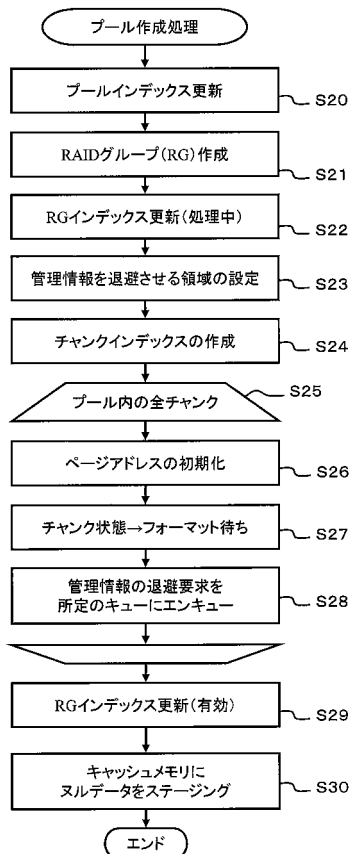
【図9】



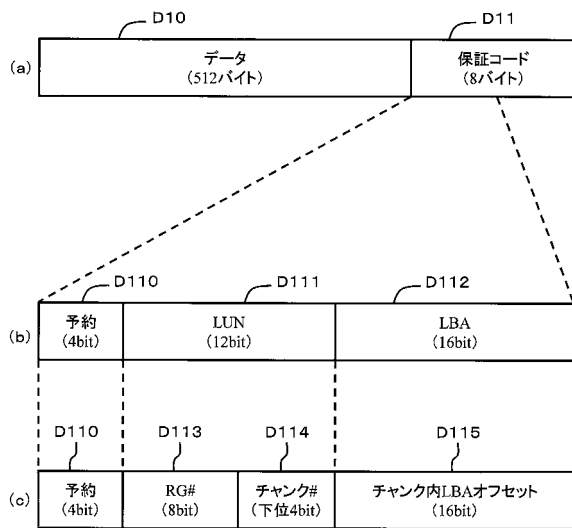
【図10】



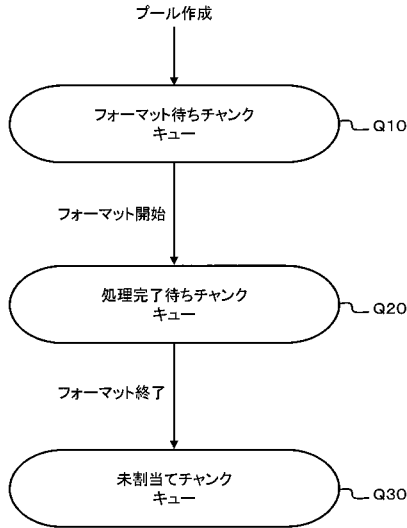
【図11】



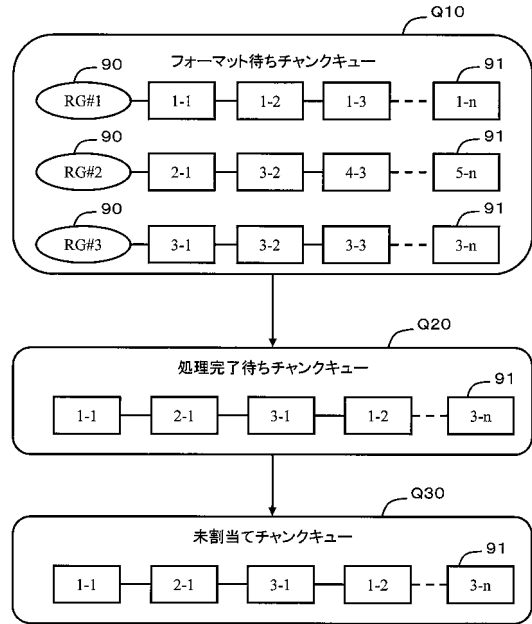
【図12】



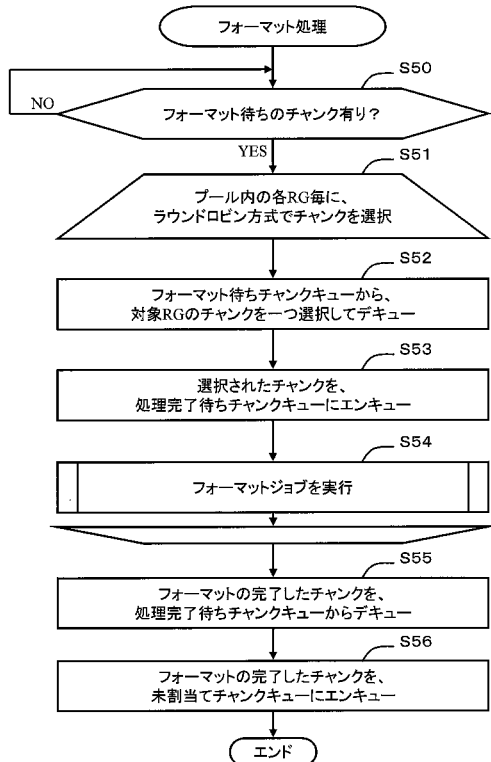
【図13】



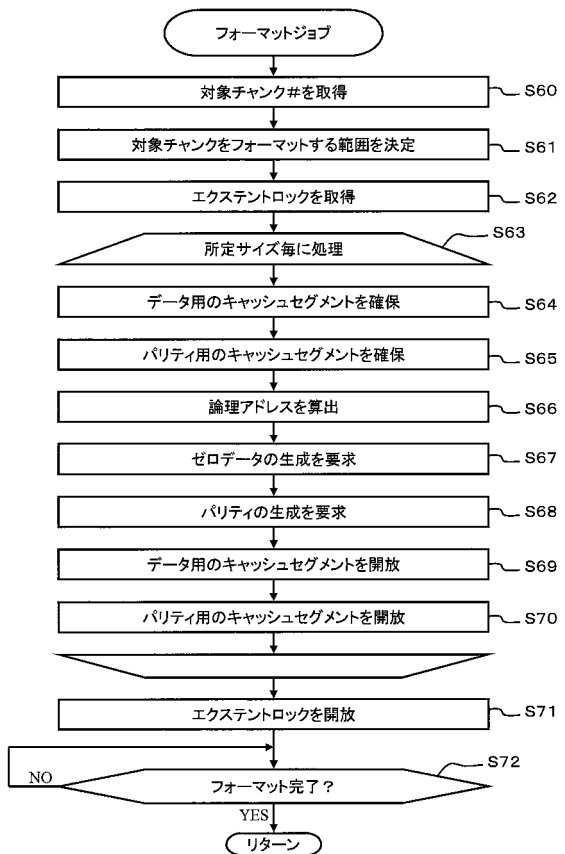
【図14】



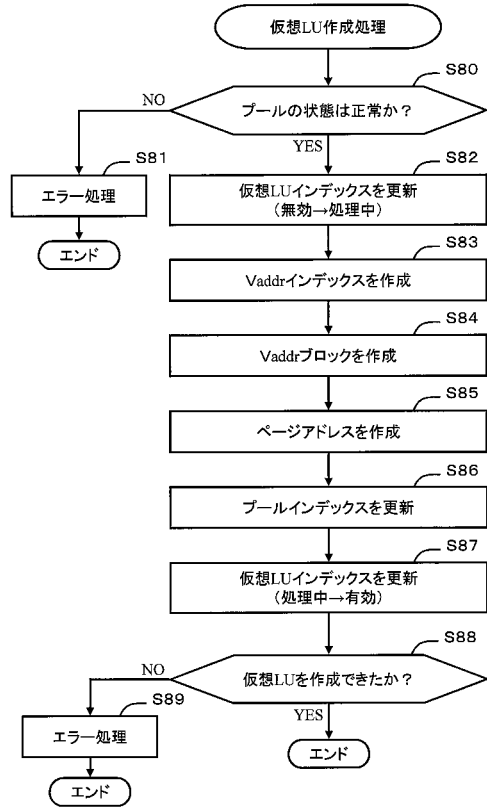
【図15】



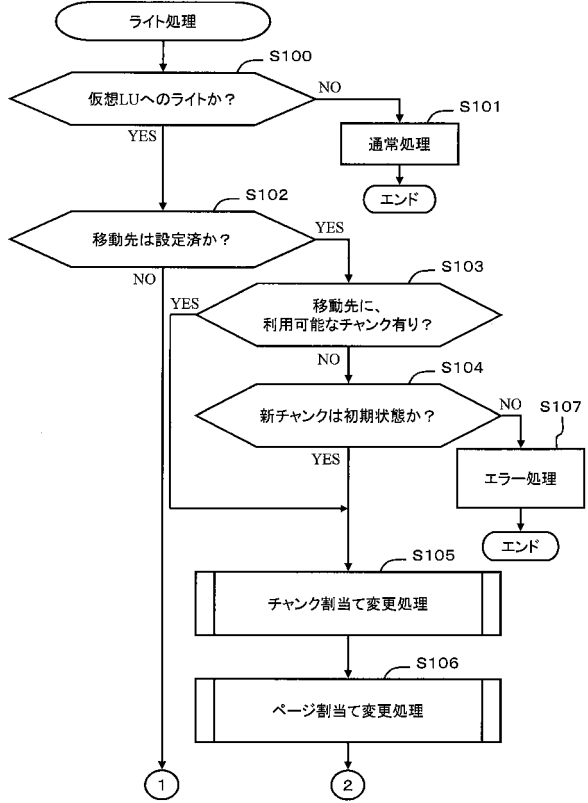
【図16】



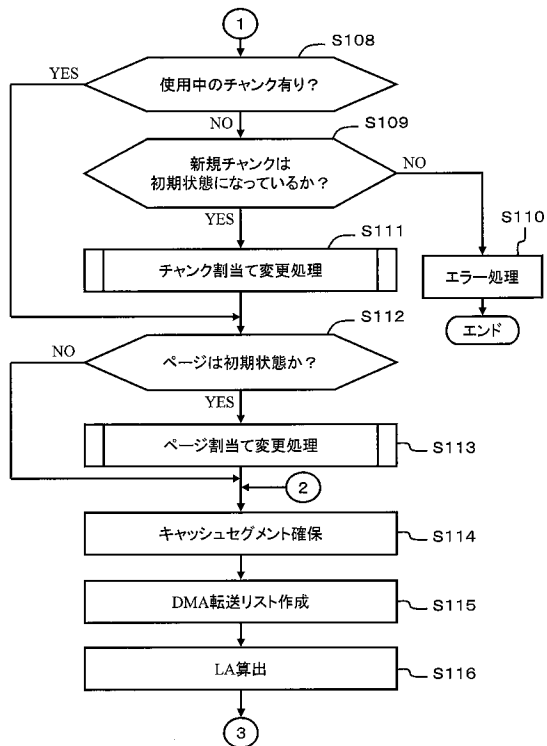
【図17】



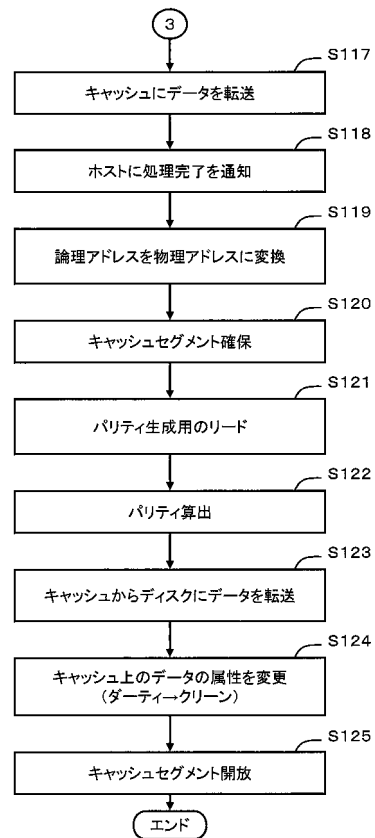
【図18】



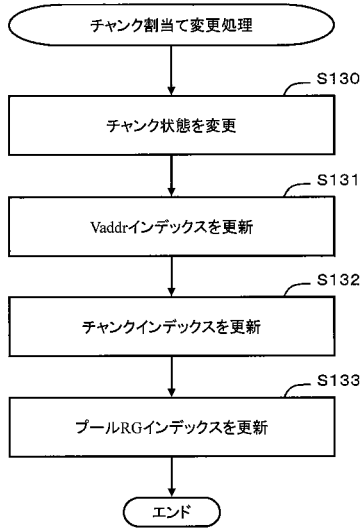
【図19】



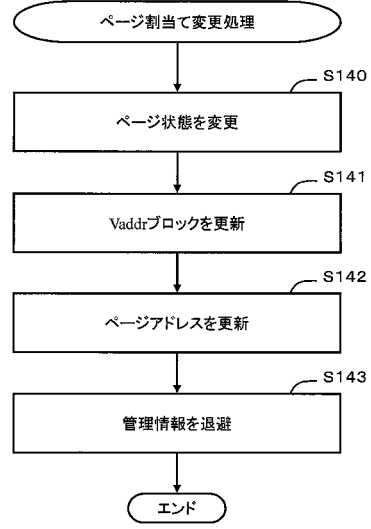
【図20】



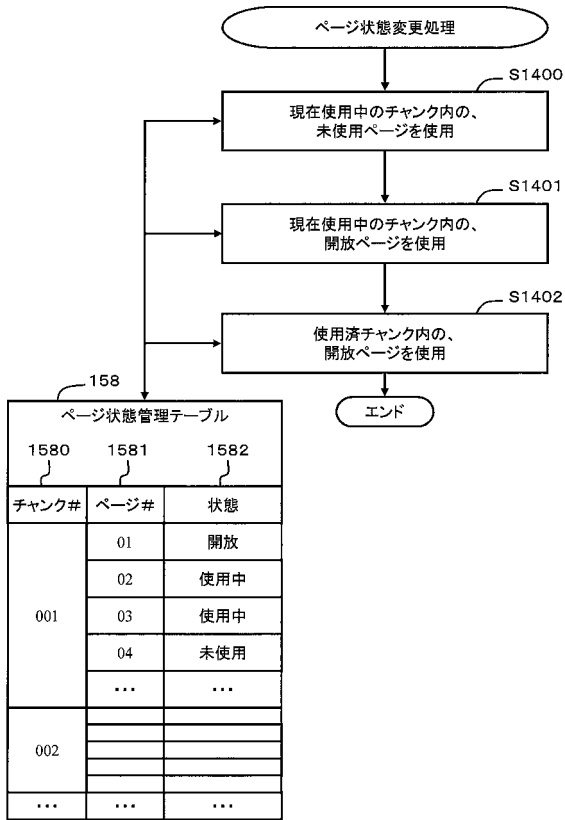
【図 2 1】



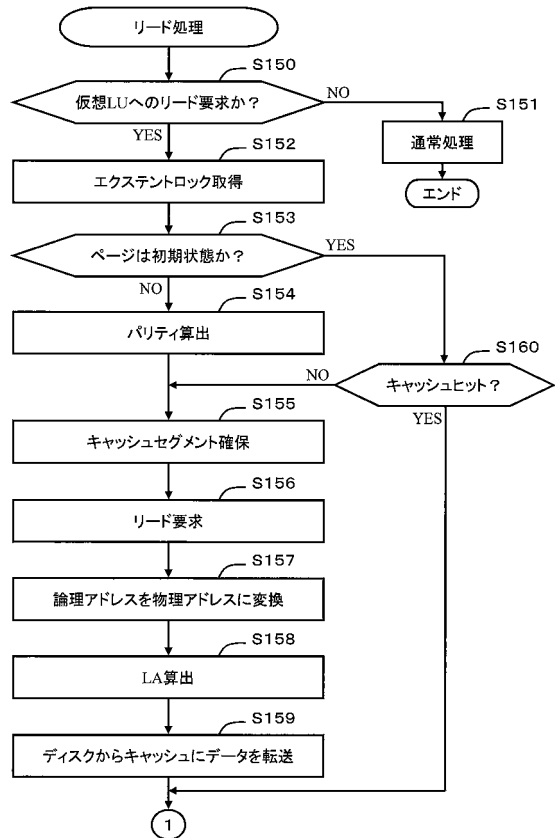
【図 2 2】



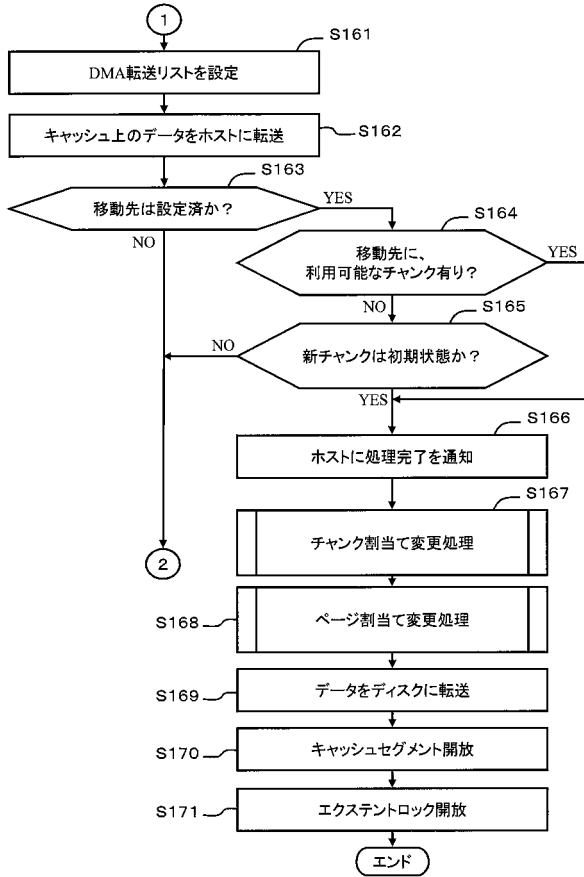
【図 2 3】



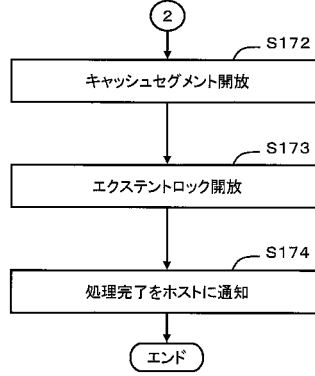
【図 2 4】



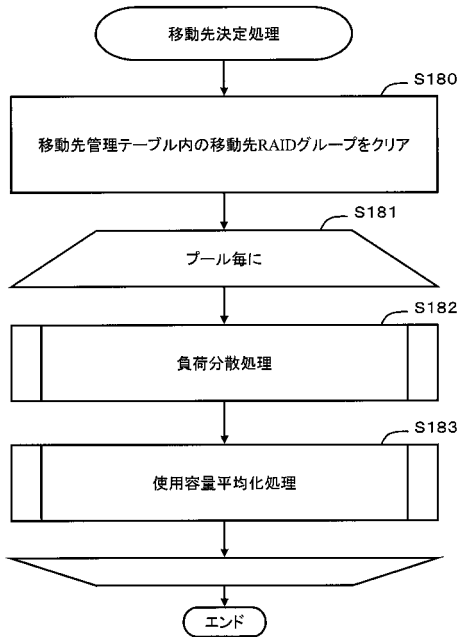
【図 25】



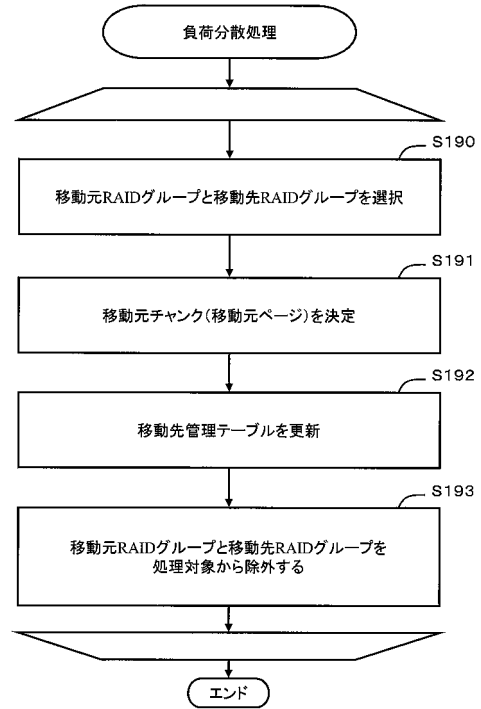
【図 26】



【図 27】



【図 28】



フロントページの続き

(72)発明者 内海 勝広

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

審査官 木村 貴俊

(56)参考文献 特開2008-059353(JP,A)

特開2008-090741(JP,A)

特開2008-046986(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06 - 3/08

G06F 12/00 - 12/16

G06F 13/00 - 13/42