

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5365236号  
(P5365236)

(45) 発行日 平成25年12月11日(2013.12.11)

(24) 登録日 平成25年9月20日(2013.9.20)

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

請求項の数 9 (全 19 頁)

<p>(21) 出願番号 特願2009-32335 (P2009-32335)                  (22) 出願日 平成21年2月16日 (2009.2.16)                  (65) 公開番号 特開2010-191499 (P2010-191499A)                  (43) 公開日 平成22年9月2日 (2010.9.2)                  審査請求日 平成23年12月8日 (2011.12.8)</p>	<p>(73) 特許権者 000004237                  日本電気株式会社                  東京都港区芝五丁目7番1号                  (74) 代理人 100124811                  弁理士 馬場 資博                  (74) 代理人 100088959                  弁理士 境 廣巳                  (72) 発明者 デマタビティヤ スムドゥ                  東京都港区芝五丁目7番1号 日本電気株                  式会社内                  審査官 木村 貴俊</p>
--	---

最終頁に続く

(54) 【発明の名称】 ストレージシステム

(57) 【特許請求の範囲】

【請求項1】

複数の記憶手段と、これら複数の記憶手段に対してデータを記憶するデータ処理手段と、を備え、

前記データ処理手段は、

前記複数の記憶手段のうち、ストレージシステムにおけるデータ処理状況を表すジャーナルを記憶するジャーナル用記憶手段を設定すると共に、当該設定されたジャーナル用記憶手段とは異なる他の複数の前記記憶手段を、記憶対象データを形成する複数のフラグメントデータをそれぞれ分散して記憶するフラグメント用記憶手段としてそれぞれ設定する記憶先設定手段と、

前記記憶先設定手段にて前記ジャーナル用記憶手段として設定した前記記憶手段に前記ジャーナルを記憶すると共に、前記フラグメント用記憶手段として設定した複数の前記記憶手段に前記複数のフラグメントデータを分散してそれぞれ記憶する分散記憶制御手段と、を備え、

前記記憶先設定手段は、

前記ジャーナル用記憶手段を順次異なる前記記憶手段に変更して設定すると共に、  
 $\{ ( \text{前記ジャーナル用記憶手段として設定された記憶手段の空き容量} ) - ( \text{他の記憶手段のうち最も空き容量の少ない記憶手段の空き容量} ) \} > ( \text{ : 予め設定された値 } )$   
 を満たすタイミングで、前記ジャーナル用記憶手段を変更して設定する、  
 ストレージシステム。

## 【請求項 2】

請求項 1 に記載のストレージシステムであって、  
前記記憶先設定手段は、全ての前記記憶手段を変更対象として、前記ジャーナル用記憶手段を順次変更して設定する、  
ストレージシステム。

## 【請求項 3】

請求項 1 又は 2 に記載のストレージシステムであって、  
前記記憶先設定手段は、空き容量が最も少ない前記記憶手段を、前記ジャーナル用記憶手段として設定する、  
ストレージシステム。

10

## 【請求項 4】

請求項 1 乃至 3 のいずれか一項に記載のストレージシステムであって、  
前記分散記憶制御手段にて前記ジャーナルを前記記憶手段に記憶すると共に、当該ジャーナルをストレージシステムに装備された揮発性メモリにも記憶するジャーナル記憶制御手段を備えた、  
ストレージシステム。

## 【請求項 5】

請求項 4 に記載のストレージシステムであって、  
前記ジャーナル記憶制御手段は、前記記憶先設定手段にて前記ジャーナル用記憶手段を変更設定したときに、前記揮発性メモリに記憶されているジャーナルを、前記ジャーナル用記憶手段として新たに変更設定された前記記憶手段に記憶する、  
ストレージシステム。

20

## 【請求項 6】

請求項 4 又は 5 に記載のストレージシステムであって、  
前記ジャーナル記憶制御手段は、ストレージシステムの障害発生時に、前記揮発性メモリに記憶しているジャーナルを、前記フラグメント用記憶手段として設定されている前記記憶手段に書き出して記憶する、  
ストレージシステム。

## 【請求項 7】

請求項 6 に記載のストレージシステムであって、  
全ての前記記憶手段に、ジャーナルを記憶するジャーナル記憶領域を形成し、  
前記ジャーナル記憶制御手段は、前記ストレージシステムの障害発生時に、前記揮発性メモリに記憶しているジャーナルを前記フラグメント用記憶手段として設定されている全ての前記記憶手段にそれぞれ記憶する、  
ストレージシステム。

30

## 【請求項 8】

複数の記憶手段を備えた情報処理装置に、  
前記複数の記憶手段に対してデータを記憶するデータ処理手段を実現させると共に、  
前記データ処理手段は、  
前記複数の記憶手段のうち、情報処理装置におけるデータ処理状況を表すジャーナルを記憶するジャーナル用記憶手段を設定すると共に、当該設定されたジャーナル用記憶手段とは異なる他の複数の前記記憶手段を、記憶対象データを形成する複数のフラグメントデータをそれぞれ分散して記憶するフラグメント用記憶手段としてそれぞれ設定する記憶先設定手段と、

40

前記記憶先設定手段にて前記ジャーナル用記憶手段として設定した前記記憶手段に前記ジャーナルを記憶すると共に、前記フラグメント用記憶手段として設定した複数の前記記憶手段に前記複数のフラグメントデータを分散してそれぞれ記憶する分散記憶制御手段と、  
を備え、

前記記憶先設定手段は、

前記ジャーナル用記憶手段を順次異なる前記記憶手段に変更して設定すると共に、

50

{ (前記ジャーナル用記憶手段として設定された記憶手段の空き容量) - (他の記憶手段のうち最も空き容量の少ない記憶手段の空き容量) } > ( : 予め設定された値 )  
を満たすタイミングで、前記ジャーナル用記憶手段を変更して設定する、  
 プログラム。

【請求項 9】

複数の記憶手段を備えたストレージシステムにて、  
 前記複数の記憶手段に対してデータを記憶するデータ処理工程を有し、  
 前記データ処理工程は、  
 前記複数の記憶手段のうち、ストレージシステムにおけるデータ処理状況を表すジャーナルを記憶するジャーナル用記憶手段を設定すると共に、当該設定されたジャーナル用記憶手段とは異なる他の複数の前記記憶手段を、記憶対象データを形成する複数のフラグメントデータをそれぞれ分散して記憶するフラグメント用記憶手段としてそれぞれ設定する記憶先設定工程と、

前記記憶先設定工程にて前記ジャーナル用記憶手段として設定した前記記憶手段に前記ジャーナルを記憶すると共に、前記フラグメント用記憶手段として設定した複数の前記記憶手段に前記複数のフラグメントデータを分散してそれぞれ記憶する分散記憶制御工程と、を有し、

前記記憶先設定工程は、

前記ジャーナル用記憶手段を順次異なる前記記憶手段に変更して設定すると共に、

{ (前記ジャーナル用記憶手段として設定された記憶手段の空き容量) - (他の記憶手段のうち最も空き容量の少ない記憶手段の空き容量) } > ( : 予め設定された値 )  
を満たすタイミングで、前記ジャーナル用記憶手段を変更して設定する、  
 データ処理方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ストレージシステムにかかり、特に、データを分散して複数の記憶装置に記憶するストレージシステムに関する。

【背景技術】

【0002】

近年、コンピュータの発達及び普及に伴い、種々の情報がデジタルデータ化されている。このようなデジタルデータを保存しておく装置として、磁気テープや磁気ディスクなどの記憶装置がある。そして、保存すべきデータは日々増大し、膨大な量となるため、大容量なストレージシステムが必要となっている。また、記憶装置に費やすコストを削減しつつ、信頼性も必要とされる。これに加えて、後にデータを容易に取り出すことが可能であることも必要である。その結果、自動的に記憶容量や性能の増大を実現できると共に、重複記憶を排除して記憶コストを削減し、さらには、冗長性の高いストレージシステムが望まれている。

【0003】

このような状況に応じて、近年では、特許文献 1 に示すように、コンテンツアドレスストレージシステムが開発されている。このコンテンツアドレスストレージシステムは、データを分散して複数の記憶装置に記憶すると共に、このデータの内容に応じて特定される固有のコンテンツアドレスによって、当該データを格納した格納位置が特定される。

【0004】

具体的に、コンテンツアドレスストレージシステムでは、図 1 に示すように、ストレージサーバ 100 の分割部 101 にて、記憶対象データであるデータブロックを複数のフラグメントに分割すると共に、冗長データとなるフラグメントをさらに付加する。そして、これら複数のフラグメントを、ストレージサーバ 100 内の記憶装置 102 や、他のストレージサーバに装備された記憶装置 103 , 104 など、複数の記憶装置にそれぞれ格納

10

20

30

40

50

する。そして、後に、コンテンツアドレスを指定することにより、当該コンテンツアドレスにて特定される格納位置に格納されているデータつまりフラグメントを読み出し、複数のフラグメントから分割前の所定のデータを復元することができる。

【0005】

また、上記コンテンツアドレスは、データの内容に応じて固有となるよう生成される。このため、重複データであれば同じ格納位置のデータを参照することで、同一内容のデータを取得することができる。従って、重複データを別々に格納する必要がなく、重複記録を排除し、データ容量の削減を図ることができる。

【0006】

そして、上述したようなストレージシステムにおいて、同一のデータブロックから分割された2個以上のフラグメントを同一のサーバに書き込む場合は、耐ディスク障害性を保つために、同じディスクに2個以上のフラグメントを書き込まないように各ディスクに分散して記憶する。一方、同じディスクに異なるデータブロックに属した2個以上のフラグメントの書き込みを行う必要がある場合には、書き込み待ちのフラグメントをキューに溜めてから書き込む。そして、キューから1または複数のフラグメントをディスクに書き込む処理は、トランザクション方式で処理され、トランザクションをジャーナルに記録する。

10

【0007】

ここで、図2に、上述したトランザクションをジャーナルに記憶する様子を示す。この図に示すように、ジャーナル(J)は、データ(D)の書込み開始直前に書かれるスタートエントリー、データ書込み終了直後に書かれるコミットエントリー、用済みのスタート・コミットエントリーに対する削除エントリー、で構成される。

20

【0008】

そして、上述したジャーナル110は、冗長性を保つために、図3に示すように、ハードウェアRAID (Redundant Arrays of Inexpensive Disks) を構成するディスク領域のOS (Operating System) 領域120内に、固定的に配置される。つまり、ジャーナルは、2つのディスクに2重化されたOS領域120内に配置される。なお、関連する技術として、特許文献2には、ジャーナルを記憶する空きスペースがなくなった場合には、当該ジャーナルを記憶するファイルを切り替える、という技術が開示されている。

【0009】

一方で、ディスク容量を最大に活用するためには、OS領域120があるディスク(ディスク0、ディスク1)の残容量もデータ領域130として使う、ことが必要である。これは、ディスク単体が大容量化している今日においては、OS領域が相対的に小さいため、極めて重要である。従って、OS領域120が設けられたディスク及び他のディスク、つまり、全てのディスク(ディスク0~ディスクn)に、データそのものを記憶するデータ領域130を形成する必要がある。

30

【先行技術文献】

【特許文献】

【0010】

【特許文献1】特開2005-235171号公報

40

【特許文献2】特開平2-54346号公報

【発明の概要】

【発明が解決しようとする課題】

【0011】

しかしながら、OS領域120を有するディスクにも、フラグメントデータといった記憶対象データを分散記憶すると、ストレージシステムにおける書き込み性能の低下、という問題が生じる。具体的には、図4に示すように、OS領域が形成されたディスクにフラグメントデータを格納すると、当該フラグメントデータとジャーナルとが同一ディスク(ディスク0, 1)に格納されることとなる。すると、かかるディスクにおいて書き込む処理の競合が発生し、書き込みが遅れ、フラグメント化された元のデータブロック全体の書

50

き込み完了も遅れる。その結果、書き込み性能が低下する、という問題が生じる。さらに、ディスク障害が生じた場合には、データとジャーナルの双方が消失する場合が生じるため、信頼性も低下する、という問題が生じる。

【0012】

このため、本発明の目的は、上述した課題である、冗長性の向上を図りつつ、書き込み性能の向上を図ることができるストレージシステムを提供する、ことにある。

【課題を解決するための手段】

【0013】

かかる目的を達成するため本発明の一形態であるストレージシステムは、  
複数の記憶手段と、これら複数の記憶手段に対してデータを記憶するデータ処理手段と、  
を備えている。

そして、上記データ処理手段は、

上記複数の記憶手段のうち、ストレージシステムにおけるデータ処理状況を表すジャーナルを記憶するジャーナル用記憶手段を設定すると共に、当該設定されたジャーナル用記憶手段とは異なる他の複数の上記記憶手段を、記憶対象データを形成する複数のフラグメントデータをそれぞれ分散して記憶するフラグメント用記憶手段としてそれぞれ設定する記憶先設定手段と、

上記記憶先設定手段にて上記ジャーナル用記憶手段として設定した上記記憶手段に上記ジャーナルを記憶すると共に、上記フラグメント用記憶手段として設定した複数の上記記憶手段に上記複数のフラグメントデータを分散してそれぞれ記憶する分散記憶制御手段と、  
を備えた、という構成を採る。

【0014】

また、本発明の他の形態であるプログラムは、  
複数の記憶手段を備えた情報処理装置に、  
上記複数の記憶手段に対してデータを記憶するデータ処理手段を実現させるプログラム  
である。

そして、上記データ処理手段は、

上記複数の記憶手段のうち、情報処理装置におけるデータ処理状況を表すジャーナルを記憶するジャーナル用記憶手段を設定すると共に、当該設定されたジャーナル用記憶手段とは異なる他の複数の上記記憶手段を、記憶対象データを形成する複数のフラグメントデータをそれぞれ分散して記憶するフラグメント用記憶手段としてそれぞれ設定する記憶先設定手段と、

上記記憶先設定手段にて上記ジャーナル用記憶手段として設定した上記記憶手段に上記ジャーナルを記憶すると共に、上記フラグメント用記憶手段として設定した複数の上記記憶手段に上記複数のフラグメントデータを分散してそれぞれ記憶する分散記憶制御手段と、  
を備えた、という構成を採る。

【0015】

また、本発明の他の形態であるデータ処理方法は、  
複数の記憶手段を備えたストレージシステムにて、  
上記複数の記憶手段に対してデータを記憶するデータ処理工程を有し、  
上記データ処理工程は、

上記複数の記憶手段のうち、ストレージシステムにおけるデータ処理状況を表すジャーナルを記憶するジャーナル用記憶手段を設定すると共に、当該設定されたジャーナル用記憶手段とは異なる他の複数の上記記憶手段を、記憶対象データを形成する複数のフラグメントデータをそれぞれ分散して記憶するフラグメント用記憶手段としてそれぞれ設定する記憶先設定工程と、

上記記憶先設定工程にて上記ジャーナル用記憶手段として設定した上記記憶手段に上記ジャーナルを記憶すると共に、上記フラグメント用記憶手段として設定した複数の上記記憶手段に上記複数のフラグメントデータを分散してそれぞれ記憶する分散記憶制御工程と、  
を有する、という構成を採る。

10

20

30

40

50

## 【発明の効果】

## 【0016】

本発明は、以上のように構成されることにより、ストレージシステムにおいて、記憶するデータの冗長性の向上を図りつつ、書き込み性能の向上を図ることができる。

## 【図面の簡単な説明】

## 【0017】

【図1】本発明に関連するストレージシステムにおけるデータの記録時の様子を示す図である。

【図2】本発明に関連するストレージシステムにおけるジャーナルの記憶の様子を示す図である。

【図3】本発明に関連するストレージシステムにおけるディスクの構成を示す図である。

【図4】本発明に関連するストレージシステムにおけるフラグメントデータとジャーナルの記録の様子を示す図である。

【図5】本発明の実施形態1におけるストレージシステムの構成を示す機能ブロック図である。

【図6】本発明の実施形態2におけるシステム全体の構成を示すブロック図である。

【図7】本発明の実施形態2におけるストレージシステムの概略を示すブロック図である。

【図8】図7に開示したストレージシステムの構成を示す機能ブロック図である。

【図9】図8に開示したストレージシステムによるデータの記憶動作を説明するための説明図である。

【図10】ストレージシステムに装備されたディスクの記憶領域を示す図である。

【図11】ストレージシステムに記憶されたデータの一例を示す図である。

【図12】ストレージシステムにおけるデータの記憶動作を示す説明図である。

【図13】ストレージシステムにおけるデータの記憶動作を示す説明図である。

【図14】ストレージシステムの動作を示すフローチャートである。

【図15】ストレージシステムの動作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】

## 【0018】

## &lt;実施形態1&gt;

本発明の第1の実施形態を、図5を参照して説明する。図5は、ストレージシステムの構成を示す機能ブロック図である。なお、本実施形態では、ストレージシステムの概略を説明する。

## 【0019】

図5に示すように、本実施形態におけるストレージシステム1は、複数の記憶手段5と、これら複数の記憶手段5に対してデータを記憶するデータ処理手段2と、を備えている。

そして、上記データ処理手段2は、

上記複数の記憶手段5のうち、ストレージシステムにおけるデータ処理状況を表すジャーナルを記憶するジャーナル用記憶手段を設定すると共に、当該設定されたジャーナル用記憶手段とは異なる他の複数の上記記憶手段を、記憶対象データを形成する複数のフラグメントデータをそれぞれ分散して記憶するフラグメント用記憶手段としてそれぞれ設定する記憶先設定手段3と、

上記記憶先設定手段3にて上記ジャーナル用記憶手段として設定した上記記憶手段5に上記ジャーナルを記憶すると共に、上記フラグメント用記憶手段として設定した複数の上記記憶手段5に上記複数のフラグメントデータを分散してそれぞれ記憶する分散記憶制御手段4と、を備えている。

## 【0020】

上記構成のストレージシステムによると、ジャーナルを記憶するジャーナル用記憶手段と、記憶対象データを形成する各フラグメントデータを記憶するフラグメント用記憶手段

10

20

30

40

50

とが、重複しないよう設定される。従って、ジャーナルが特定の記憶手段に固定的に記憶されることなく、また、記憶手段に障害が発生した場合であってもフラグメントデータとジャーナルとが同時に消失することを抑制できる。さらに、ジャーナルとフラグメントデータが1つの記憶手段に記憶されることがないため、記憶手段に対するデータの書き込み処理の競合を抑制することができる。その結果、記憶するデータの冗長性の向上を図りつつ、書き込み性能の向上を図ることができる。

【0021】

また、上記ストレージシステムでは、上記記憶先設定手段は、上記ジャーナル用記憶手段を順次異なる上記記憶手段に変更して設定する、という構成を採る。また、上記ストレージシステムでは、上記記憶先設定手段は、全ての上記記憶手段を変更対象として、上記ジャーナル用記憶手段を順次変更して設定する、という構成を採る。

10

【0022】

また、上記ストレージシステムでは、上記記憶先設定手段は、上記ジャーナル用記憶手段として設定された上記記憶手段の空き容量と、他の上記記憶手段の空き容量と、に基づくタイミングで、上記ジャーナル用記憶手段を変更して設定する、という構成を採る。例えば、上記記憶先設定手段は、 $\{ ( \text{上記ジャーナル用記憶手段として設定された記憶手段の空き容量} ) - ( \text{他の記憶手段のうち最も空き容量の少ない記憶手段の空き容量} ) \} > ( \quad : \text{予め設定された値} )$ 、を満たすタイミングで、上記ジャーナル用記憶手段を変更して設定する。

【0023】

また、上記ストレージシステムでは、上記記憶先設定手段は、空き容量が最も少ない上記記憶手段を、上記ジャーナル用記憶手段として設定する、という構成を採る。

20

【0024】

このように、ジャーナルを格納する領域を順次変更設定することで、比較的、小さいデータ容量であるジャーナルを効率的に記憶することができる。

【0025】

また、上記ストレージシステムでは、上記分散記憶制御手段にて上記ジャーナルを上記記憶手段に記憶すると共に、当該ジャーナルをストレージシステムに装備された揮発性メモリにも記憶するジャーナル記憶制御手段を備えた、という構成とる。

【0026】

そして、上記ストレージシステムでは、上記ジャーナル記憶制御手段は、上記記憶先設定手段にて上記ジャーナル用記憶手段を変更設定したときに、上記揮発性メモリに記憶されているジャーナルを、上記ジャーナル用記憶手段として新たに変更設定された上記記憶手段に記憶する、という構成を採る。

30

ストレージシステム

【0027】

さらに、上記ストレージシステムでは、上記ジャーナル記憶制御手段は、上記揮発性メモリに記憶しているジャーナルを、ストレージシステムの作動状況に応じて、上記フラグメント用記憶手段として設定されている上記記憶手段に記憶する、という構成を採る。このとき、例えば、上記ジャーナル記憶制御手段は、上記ストレージシステムの障害発生時に、上記揮発性メモリに記憶しているジャーナルを上記フラグメント用記憶手段として設定されている全ての上記記憶手段にそれぞれ記憶する、という構成を採る。

40

ストレージシステム

【0028】

これにより、ストレージシステムの障害時などに、ジャーナルをいずれかの記憶手段に記憶させておくことができる。従って、仮にジャーナル用記憶手段に障害が発生した場合であっても、ジャーナルが他の記憶手段に保持され、これを参照して後に復旧が可能となる。その結果、さらなる冗長性の向上を図ることができる。

【0029】

また、上述したストレージシステムは、情報処理装置に、プログラムが組み込まれるこ

50

とで実現できる。具体的に、本発明の他の形態であるプログラムは、

複数の記憶手段を備えた情報処理装置に、

上記複数の記憶手段に対してデータを記憶するデータ処理手段を実現させるプログラムであり、

上記データ処理手段は、

上記複数の記憶手段のうち、情報処理装置におけるデータ処理状況を表すジャーナルを記憶するジャーナル用記憶手段を設定すると共に、当該設定されたジャーナル用記憶手段とは異なる他の複数の上記記憶手段を、記憶対象データを形成する複数のフラグメントデータをそれぞれ分散して記憶するフラグメント用記憶手段としてそれぞれ設定する記憶先設定手段と、

上記記憶先設定手段にて上記ジャーナル用記憶手段として設定した上記記憶手段に上記ジャーナルを記憶すると共に、上記フラグメント用記憶手段として設定した複数の上記記憶手段に上記複数のフラグメントデータを分散してそれぞれ記憶する分散記憶制御手段と、を備えた、という構成を採る。

【0030】

そして、上記プログラムでは、上記記憶先設定手段は、上記ジャーナル用記憶手段を順次異なる上記記憶手段に変更して設定する、という構成を採る。

【0031】

また、上述したストレージシステムが作動することにより実行される、本発明の他の形態であるデータ処理方法は、

複数の記憶手段を備えたストレージシステムにて、

上記複数の記憶手段に対してデータを記憶するデータ処理工程を有し、

上記データ処理工程は、

上記複数の記憶手段のうち、ストレージシステムにおけるデータ処理状況を表すジャーナルを記憶するジャーナル用記憶手段を設定すると共に、当該設定されたジャーナル用記憶手段とは異なる他の複数の上記記憶手段を、記憶対象データを形成する複数のフラグメントデータをそれぞれ分散して記憶するフラグメント用記憶手段としてそれぞれ設定する記憶先設定工程と、

上記記憶先設定工程にて上記ジャーナル用記憶手段として設定した上記記憶手段に上記ジャーナルを記憶すると共に、上記フラグメント用記憶手段として設定した複数の上記記憶手段に上記複数のフラグメントデータを分散してそれぞれ記憶する分散記憶制御工程と、を有する、という構成を採る。

【0032】

そして、上記データ処理方法では、上記記憶先設定工程は、上記ジャーナル用記憶手段を順次異なる上記記憶手段に変更して設定する、という構成を採る。

【0033】

上述した構成を有する、プログラム、又は、データ処理方法、の発明であっても、上記ストレージシステムと同様の作用を有するために、上述した本発明の目的を達成することができる。

【0034】

<実施形態2>

本発明の第2の実施形態を、図6乃至図15を参照して説明する。図6は、システム全体の構成を示すブロック図である。図7は、ストレージシステムの概略を示すブロック図であり、図8は、構成を示す機能ブロック図である。図9は、ストレージシステムによるデータの記憶動作を説明するための説明図である。図10は、ストレージシステムに装備されたディスクの記憶領域を示す図である。図11は、ストレージシステムに記憶されたデータの一例を示す図である。図12乃至図13は、ストレージシステムにおけるデータの記憶動作を示す説明図である。図14乃至図15は、ストレージシステムの動作を示すフローチャートである。

【0035】

ここで、本実施形態は、上述した実施形態 1 にて開示したストレージシステムの具体的な一例を示すものである。そして、以下では、ストレージシステムが、複数台のサーバコンピュータが接続されて構成されている場合を説明する。但し、本発明におけるストレージシステムは、複数台のコンピュータにて構成されることに限定されず、1台のコンピュータで構成されていてもよい。

**【0036】****[構成]**

図 6 に示すように、対象本発明におけるストレージシステム 10 は、ネットワーク N を介してバックアップ処理を制御するバックアップシステム 11 に接続している。そして、バックアップシステム 11 は、ネットワーク N を介して接続されたバックアップ対象装置 12 に格納されているバックアップ対象データ（記憶対象データ）を取得し、ストレージシステム 10 に対して記憶するよう要求する。これにより、ストレージシステム 10 は、記憶要求されたバックアップ対象データをバックアップ用に記憶する。

10

**【0037】**

そして、図 7 に示すように、本実施形態におけるストレージシステム 10 は、複数のサーバコンピュータが接続されて構成を採っている。具体的に、ストレージシステム 10 は、ストレージシステム 10 自体における記憶再生動作を制御するサーバコンピュータであるアクセラレータノード 10A と、データを格納する記憶装置を備えたサーバコンピュータであるストレージノード 10B と、を備えている。なお、アクセラレータノード 10A の数とストレージノード 10B の数は、図 7 に示したものに限定されず、さらに多くの各ノード 10A, 10B が接続されて構成されていてもよい。

20

**【0038】**

さらに、本実施形態におけるストレージシステム 10 は、データを分割及び冗長化し、分散して複数の記憶装置に記憶すると共に、このデータの内容に応じて特定される固有のコンテンツアドレスによって、当該データを格納した格納位置を特定するコンテンツアドレスストレージシステムである。具体的な構成については詳述する。

**【0039】**

また、以下では、ストレージシステム 10 が 1 つのシステムであるとして、当該ストレージシステム 10 が備えている構成及び機能を説明する。つまり、以下に説明するストレージシステム 10 が有する構成及び機能は、アクセラレータノード 10A あるいはストレージノード 10B のいずれに備えられていてもよい。なお、ストレージシステム 10 は、図 7 に示すように、必ずしもアクセラレータノード 10A とストレージノード 10B とを備えていることに限定されず、いかなる構成であってもよい。また、ストレージシステム 10 は、コンテンツアドレスストレージシステムであることにも限定されない。

30

**【0040】**

図 8 に、ストレージシステム 10 の構成を示す。この図に示すように、ストレージシステム 10 は、データを記憶するデータ記憶装置 30 と、当該データ記憶装置 30 に対するデータの記憶及び読み出し動作を制御するデータ処理装置 20（データ処理手段）と、を備えている。なお、実際には、データ処理装置 20 は、図 7 に示したアクセラレータノード 10A 及びストレージノード 10B が備えている CPU（Central Processing Unit）などの複数の演算装置にて構成されている。また、記憶装置 30 は、図 7 に示したアクセラレータノード 10A 及びストレージノード 10B が備えているハードディスクなどの記憶装置にて構成されている。

40

**【0041】**

そして、図 8 に示すように、上記データ処理装置 20 は、プログラムが組み込まれることにより構築された、フラグメント生成部 21 と、分散記憶制御部 22 と、ジャーナル記

50

憶制御部 2 3 と、記憶先設定部 2 4 と、を備えている。また、データ記憶装置 3 0 には、ハードディスクドライブなど複数の記憶装置 3 1 (記憶手段)の他に、ジャーナルを格納する揮発性のメモリ 3 2 を備えている。以下、各構成について説明してください。

【 0 0 4 2 】

まず、上記データ処理装置 2 0 の基本的な構成、つまり、フラグメント生成部 2 1 と分散記憶制御部 2 2 とによる分散記憶機能について、図 9 を参照して説明する。

【 0 0 4 3 】

フラグメント生成部 2 1 は、まず、バックアップ対象データ A の入力を受けると、当該バックアップ対象データ A を、所定容量 (例えば、6 4 K B ) のブロックデータ D に分割する。そして、さらに、ブロックデータ D を複数の所定の容量のフラグメントデータに分割する。例えば、符号 D 1 ~ D 9 に示すように、9 つのフラグメントデータ (分割データ 4 1 ) に分割する。さらに、分割したフラグメントデータのうちいくつかは欠けた場合であっても、元となるブロックデータを復元可能なよう冗長データを生成し、上記分割したフラグメントデータ 4 1 に追加する。例えば、符号 D 1 0 ~ D 1 2 に示すように、3 つのフラグメントデータ (冗長データ 4 2 ) を付加する。これにより、9 つの分割データ 4 1 と、3 つの冗長データとにより構成される 1 2 個のフラグメントデータからなるデータセット 4 0 を生成する。

【 0 0 4 4 】

そして、分散記憶制御部 2 2 (分散記憶制御手段)は、上述したように生成されたデータセットを構成する各フラグメントデータを、各記憶装置 3 1 に形成された各記憶領域に、それぞれ分散して格納する。例えば、1 2 個のフラグメントデータ D 1 ~ D 1 2 を生成した場合には、1 2 個の記憶装置 3 1 にそれぞれ形成したデータ格納ファイル F 1 ~ F 1 2 に、各フラグメントデータ D 1 ~ D 1 2 を 1 つずつそれぞれ格納する。

【 0 0 4 5 】

また、本実施形態における分散記憶制御部 2 2 は、さらに、ストレージシステム内におけるデータ処理状況、例えば、データ格納処理状況やデータ削除処理状況などの記録を表すジャーナルを、記憶装置 3 1 に格納する機能を有する。このとき、分散記憶制御部 2 2 は、記憶先設定部 2 4 によって設定された各記憶装置 3 1 に対して、それぞれジャーナルとフラグメントデータとを記憶する。つまり、ジャーナル用記憶装置として設定された記憶装置 3 1 にジャーナルを記憶し、フラグメント用記憶装置として設定された記憶装置に各フラグメントデータを記憶する。

【 0 0 4 6 】

次に、記憶先設定部 2 4 (記憶先設定手段)について説明する。記憶先設定部 2 4 は、上述したように、分散記憶制御部 2 2 にて記憶されるジャーナルやフラグメントデータの記憶先となる記憶装置 3 1 を設定する機能を有する。

【 0 0 4 7 】

ここで、本実施形態では、図 1 0 に示すように、記憶装置 3 1 として、n + 1 台のディスク (0 ~ n) を備えていることとする。そして、これらディスクのうち、ディスク 0 とディスク 1 には、二重化された OS (Operating System) 領域 5 2 を備えている。また、本実施形態では、全てのディスク 0 ~ n に、ジャーナルを記憶するジャーナル記憶領域 5 1 と、フラグメントデータを記憶するデータ領域 5 3 と、を形成している。つまり、OS 領域 5 2 を備えるディスク 0 , 1 及び他のディスク 2 ~ n の全てに、ジャーナル記憶領域 5 1 とデータ領域 5 3 とが形成されている。

【 0 0 4 8 】

上述した状況において、記憶先設定部 2 4 は、1 つのディスクをジャーナル用記憶装置として設定すると共に、このジャーナル用記憶装置として設定したディスクとは異なる他の全てのディスクを、フラグメント用記憶装置として設定する。具体的に、記憶先設定部 2 4 は、まず、初期状態では、データ領域 5 3 の空き容量が最も少ないディスクを、ジャーナル用記憶装置として設定する。このとき、空き容量が最も少ないディスクが複数存在する場合には、そのうち任意の 1 つのディスクをジャーナル用記憶装置として記憶する。

なお、空き容量が最も少ないディスクをジャーナル用記憶装置として設定する理由は、ジャーナルのデータ量が比較的少ないためである。そして、その他のディスクを、フラグメント用記憶装置として設定する。

【 0 0 4 9 】

なお、記憶先設定部 2 4 は、上述したように各ディスクを設定すると、データ記憶装置 3 0 内に形成された図 1 1 に示すようなディスク利用情報テーブルに、ディスク番号毎に、記憶するデータと、データ領域の空き容量と、を記憶する。例えば、ディスク 2 をジャーナル用記憶装置として設定した場合には、ディスク 2 に、「ジャーナル書き込み用」という情報を記憶し、その他のフラグメント用記憶装置として設定されたディスクに対しては、「データ書き込み用」という情報を記憶する。また、データ処理装置 2 0 は、一定の時間間隔にて、常に各ディスクのデータ領域 5 3 の空き容量を調べ、当該各ディスクに対応して上記テーブルに記憶しておく ( D 0 ~ D n ) 。

10

【 0 0 5 0 】

そして、上述したように、分散記憶制御部 2 2 にて、各ディスクにジャーナルやフラグメントデータが分散記憶されるが、当該各ディスクの空き容量が変化する。この空き容量の変化に応じて、上記記憶先設定部 2 4 は、ジャーナル用記憶装置としてのディスクを変更して設定する。具体的には、上述したように常に更新されているディスク利用情報テーブル内のデータ領域の空き容量に基づいて、

{ ( ジャーナル用記憶装置として設定されたディスクのデータ領域の空き容量 ) - ( 他のディスクのうち最も空き容量の少ないディスクのデータ領域の空き容量 ) } > ( : 予め設定された値 )

20

を満たすかどうかを、一定の時間間隔で常に調べる。そして、上記条件を満たしたタイミングで、そのときにデータ領域 5 3 の空き容量が最も少ないディスクを、新たなジャーナル用記憶装置として変更設定する。つまり、フラグメント用記憶装置のディスクの空き容量が小さくなると、上記条件式を満たすこととなり、そのタイミングで、当該空き容量の最も小さいディスクを、データ量の少ないジャーナルを記憶するジャーナル用記憶装置として変更設定する。

【 0 0 5 1 】

以上のようにして、記憶先設定部 2 4 は、全てのディスクが、順次、ジャーナル用記憶装置として設定されるよう、変更設定を行う。なお、記憶先設定部 2 4 は、必ずしも全てのディスクがジャーナル用記憶装置として設定されるよう変更する必要はない。また、ジャーナル用記憶装置を変更するタイミングは、上述した条件を満たすときに限定されない。例えば、ジャーナル用記憶装置として設定されたディスクのデータ領域の空き容量と、他のディスクの空き容量と、に基づく他の条件式を満たしたタイミングで、ジャーナル用記憶装置の変更設定を行ってもよい。

30

【 0 0 5 2 】

また、上記ジャーナル記憶制御部 2 3 ( ジャーナル記憶制御手段 ) は、上述したように、分散記憶制御部 2 2 にてジャーナル用記憶装置として設定されたディスクに記憶するジャーナルを、別途、メモリ 3 2 にも記憶保持するよう制御する。そして、ジャーナル記憶制御部 2 3 は、上述したように記憶先設定部 2 4 にて新たにジャーナル用記憶装置としてのディスクが設定されると、当該新たなジャーナル用記憶装置としてのディスクに、メモリ 3 2 内に保持されているジャーナルを書き出す。

40

【 0 0 5 3 】

さらに、ジャーナル記憶制御部 2 3 は、ストレージシステム 1 0 の作動状況に応じて、例えば、記憶装置 3 1 に障害が発生するなど、ストレージシステム 1 0 の障害発生時に、メモリ 3 2 に記憶しているジャーナルを、ディスクに書き出す。このとき、特に、ジャーナル用記憶装置として設定されているディスクとは異なり、フラグメント用記憶装置として設定されている他のディスクにジャーナルを記憶する。

【 0 0 5 4 】

[ 動作 ]

50

次に、上述した構成のストレージシステム 1 の動作（データ処理工程）を、図 1 2 乃至図 1 5 を参照して説明する。まず、初期状態では、データ領域の空き容量が最も少ないディスクを、ジャーナル用記憶装置に設定する（図 1 4 のステップ S 1）。そして、設定状況を、ディスク利用情報テーブルに記憶する。また、このテーブルを参照して、ジャーナル用記憶装置以外のフラグメント用記憶装置として設定されたディスク上に、データファイルを作成し、フラグメントデータを記憶する領域を確保する（図 1 4 のステップ S 2）。

【 0 0 5 5 】

その後、上述したディスクの設定に基づいて、記憶対象データの分割データ及び冗長データであるフラグメントデータと、ジャーナルと、をそれぞれ各ディスクに記憶していく。例えば、ディスク 2 がジャーナル用記憶装置として設定されている場合には、図 1 2（A）に示すように、ジャーナルをメモリ 3 2 を介してディスク 2 に記憶し、また、フラグメントデータを、その他のディスクに分散して記憶する。なお、メモリ 3 2 は、ジャーナルを保持したままである。

【 0 0 5 6 】

ここで、ジャーナルの記憶処理について、さらに図 1 3 を参照して説明する。なお、図 1 3 は、データ処理装置 2 0 内に構築されたジャーナルを記憶する処理に関連する機能を図示したものである。この図に示すように、ジャーナルは、まず、メモリ 3 2 に保持され、当該メモリ 3 2 には、常に保持される。そして、ジャーナル書き込み器 4 1 が、ジャーナル用ディスク切り替え器 4 2 にてジャーナル用記憶装置として切り替えられたディスクに、メモリ 3 2 に記憶されたジャーナルを記憶する。

【 0 0 5 7 】

続いて、上述したように、ジャーナルとフラグメントデータが各ディスクに分散記憶されている間は、常に各ディスクのデータ領域 5 3 の空き容量を更新する（図 1 4 のステップ S 3、図 1 3 の空き容量監視器）。そして、上述した条件、つまり、 $\{ (\text{ジャーナル用記憶装置として設定されたディスクのデータ領域の空き容量}) - (\text{他のディスクのうち最も空き容量の少ないディスクのデータ領域の空き容量}) \} > (\text{予め設定された値})$

を満たすと（図 1 4 のステップ S 4 で Yes）、そのときにデータ領域 5 3 の空き容量が最も少ないディスクを、新たなジャーナル用記憶装置として変更設定し（図 1 4 のステップ 5）、ディスクの設定状況をディスク利用情報テーブルに記憶する（記憶先設定工程）。なお、空き容量が最も少ないディスクが複数存在する場合には、そのうち任意の 1 つのディスクをジャーナル用記憶装置として記憶する。

【 0 0 5 8 】

続いて、新たにジャーナル用記憶装置として設定されたディスク上のデータファイル、つまり、フラグメント用の領域が全て閉じられる（クローズされる）のを待つ（図 1 4 のステップ S 6）。そして、新たにジャーナル用記憶装置として設定されたディスクのデータファイルが全て閉じられた直後に、当該ディスクのジャーナル用記憶装置としての利用を開始する（図 1 4 のステップ S 7）。このとき、まず、メモリ 3 2 に保持されているジャーナルを、新たなジャーナル記憶装置のディスクに書き出す（図 1 4 のステップ S 8）。なお、このとき、メモリ 3 2 内にもジャーナルを保持しておく。これにより、ジャーナルの移動中にクラッシュなどの障害が発生した場合であっても、ディスク内のジャーナルとメモリ内のジャーナルとを用いて、リカバリすることができる。

【 0 0 5 9 】

その後は、上述同様に、ジャーナル用記憶装置以外のフラグメント用記憶装置として設定されたディスク上に、データファイルを作成し、フラグメントデータを記憶する領域を確保する（図 1 4 のステップ S 2）。そして、上記テーブルの設定に基づいて、記憶対象データの分割データ及び冗長データであるフラグメントデータと、ジャーナルと、をそれぞれ各ディスクに記憶していく（分散記憶制御工程）。例えば、図 1 2（A）の状態から、ディスク 0 が新たなジャーナル用記憶装置として設定された場合には、図 1 2（B）に

示すように、ジャーナルをメモリ32を介してディスク0に記憶し、また、フラグメントデータを、その他のディスクに分散して記憶する。なお、メモリ32は、ジャーナルを保持したままである。

#### 【0060】

続いて、ストレージシステムに障害が発生した時の動作を、図15を参照して説明する。まず、サーバのシャットダウンやディスク障害を含む「Expected Error」などのイベントの発生を契機に、ストレージシステム1に障害が発生したことを検出する。なお、かかる障害の検出は、例えば、図13に示す、シャットダウン検出器46、ディスク障害検出器47、エラー検出器48にて行う(図15のステップS11)。

#### 【0061】

そして、上述したようなストレージシステム1の障害を検出すると、ジャーナル書き込み処理及びジャーナル用ディスクの切り替え処理を停止する(図15のステップS12)。その後、メモリ32内のジャーナルを、ディスクに書き出す処理を行う。具体的には、ジャーナル用記憶装置として設定されたディスク以外の他の全てのディスクであって、かつ、正常動作中のディスクに対して、メモリ32に保持されているジャーナルを書き出す(図15のステップS13でNo、ステップS14でYes、ステップS15)。なお、図13において、ジャーナルレプリケート器45にて、上記ジャーナルの書き出し処理が実行される。

#### 【0062】

そして、書き込みの対象となった全てのディスクに対して、ジャーナルの書き出し処理が完了すると(図15のステップ13でYes)、ストレージ・サービスの停止あるいはストレージシステムのシャットダウン継続命令を出す。

#### 【0063】

以上のように、本実施形態におけるストレージシステム1によると、ジャーナルが特定のディスクに固定的に記憶されることなく、また、ディスクに障害が発生した場合であっても、フラグメントデータとジャーナルとが同時に消失することを抑制できる。特に、ジャーナルとフラグメントデータが1つのディスクに同時に記憶されることがないため、ディスクに対するデータの書き込み処理の競合を抑制することができる。その結果、記憶するデータの冗長性の向上を図りつつ、書き込み性能の向上を図ることができる。そして、ジャーナルを格納するディスクを順次変更設定することで、比較的、小さいデータ容量であるジャーナルを効率的に記憶することができる。

#### 【0064】

また、ジャーナルをディスクに記憶すると共に、メモリにも保持しているため、冗長性が確保される。さらに、ストレージシステムの障害時などに、ジャーナルをいずれかのディスクに記憶させることで、さらなる冗長性の向上を図ることができる。

#### 【産業上の利用可能性】

#### 【0065】

本発明は、複数のコンピュータを接続して構成されるストレージシステムに利用することができ、産業上の利用可能性を有する。

#### 【符号の説明】

#### 【0066】

- 1 ストレージシステム
- 2 データ処理手段
- 3 記憶先設定手段
- 4 分散記憶制御手段
- 5 記憶手段
- 10 ストレージシステム
- 10A アクセラレータノード
- 10B ストレージノード
- 11 バックアップシステム

10

20

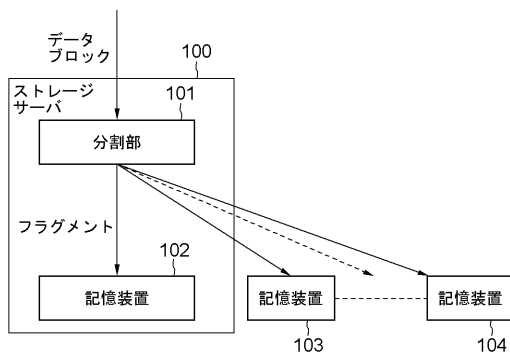
30

40

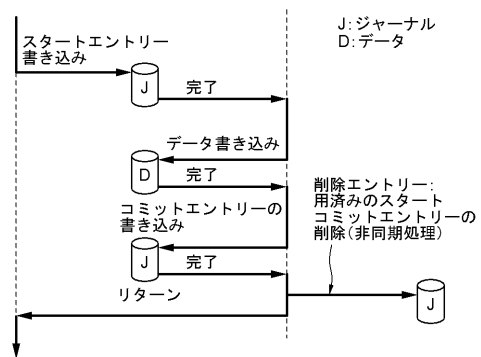
50

- 1 2 バックアップ装置
- 2 1 フラグメント生成部
- 2 2 分散記憶制御部
- 2 3 ジャーナル記憶制御部
- 2 4 記憶先設定部
- 3 1 記憶装置
- 3 2 メモリ

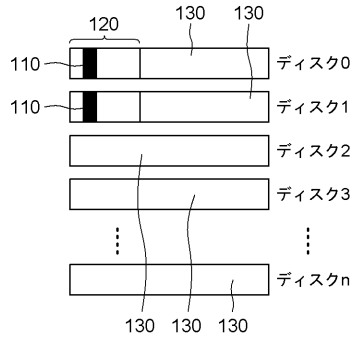
【図 1】



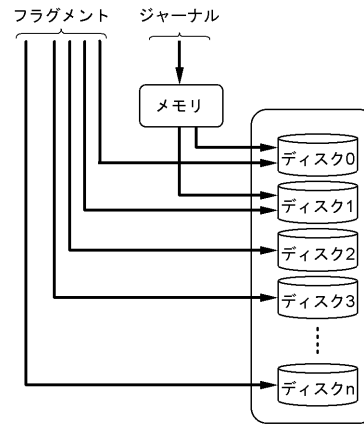
【図 2】



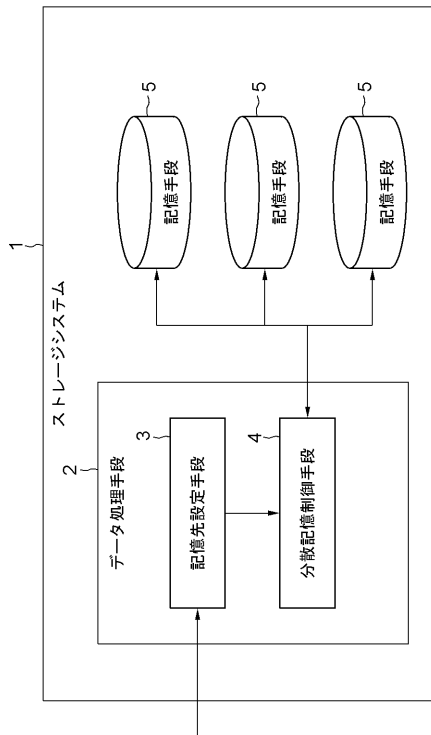
【図3】



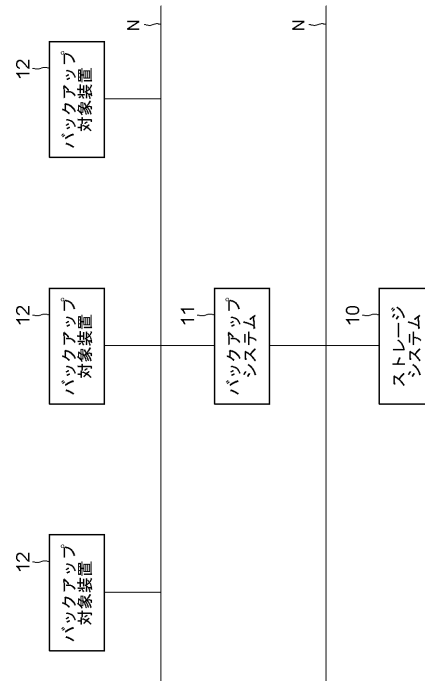
【図4】



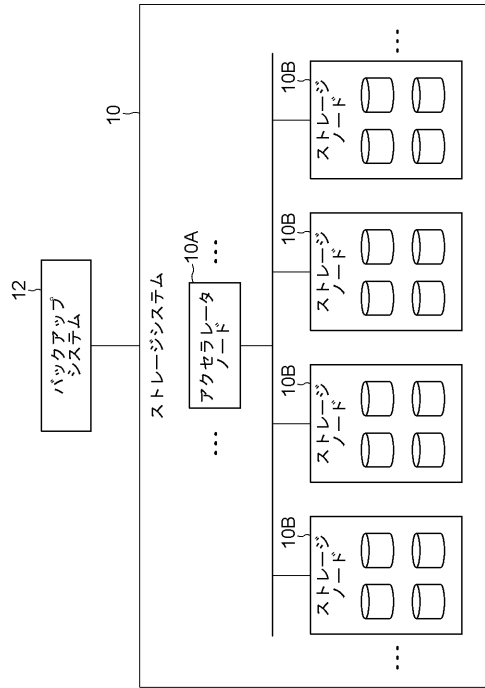
【図5】



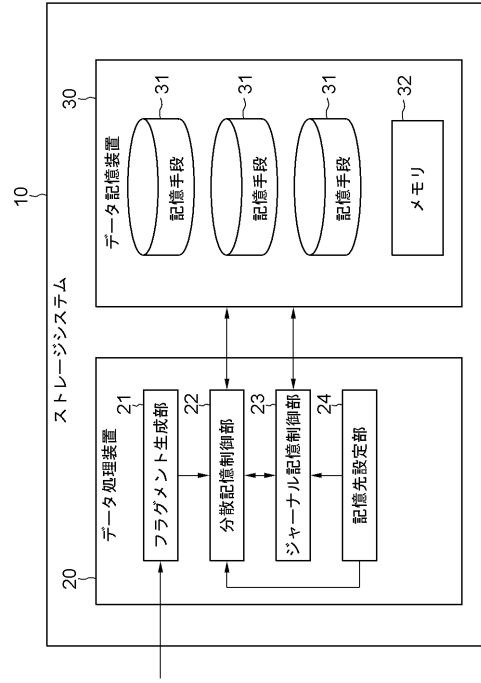
【図6】



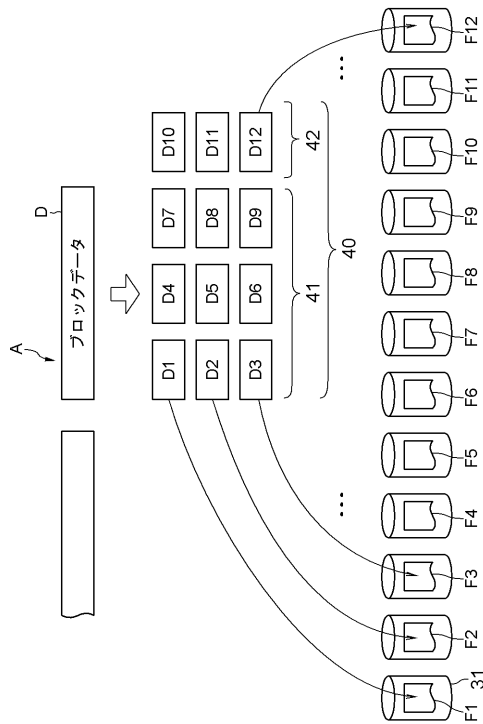
【図7】



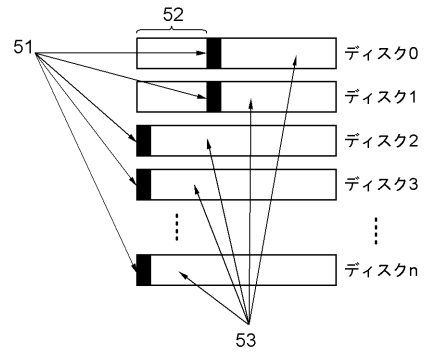
【図8】



【図9】



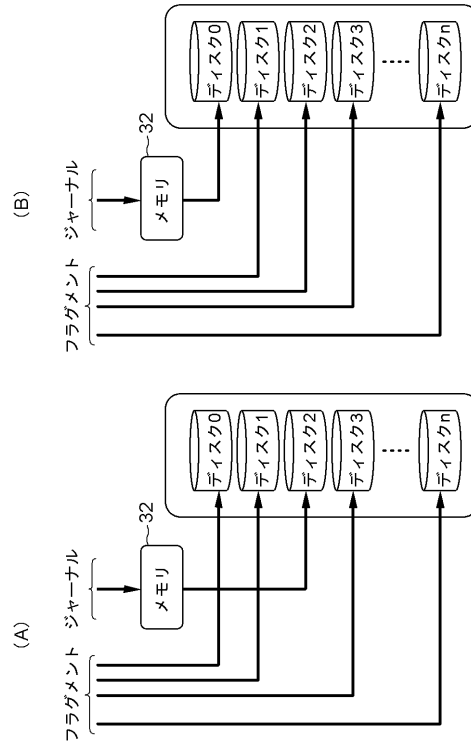
【図10】



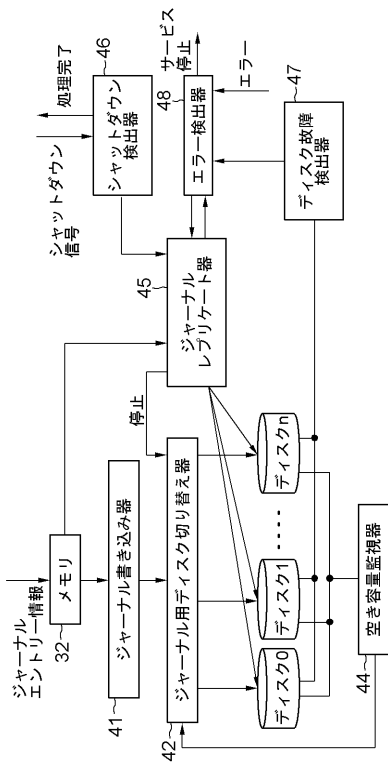
【図 1 1】

ディスク番号	データ/ジャーナル	データ領域の空き容量
ディスク0	データ書き込み用	D0
ディスク1	データ書き込み用	D1
ディスク2	ジャーナル書き込み用	D2
ディスク3	データ書き込み用	D3
ディスク4	データ書き込み用	D4
⋮	⋮	⋮
ディスクn	データ書き込み用	Dn

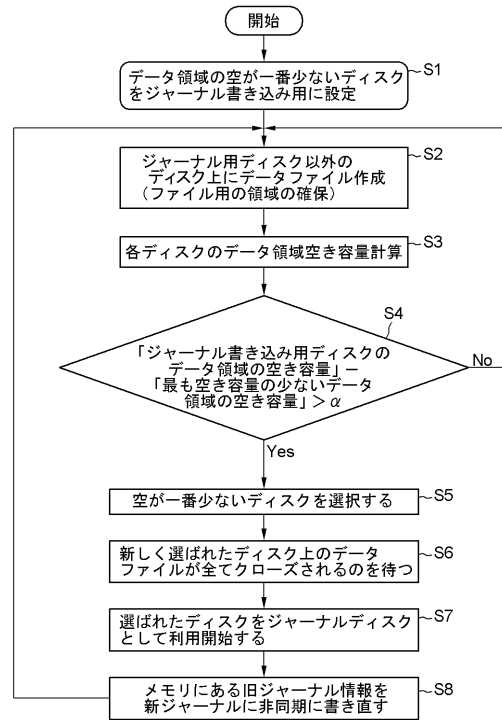
【図 1 2】



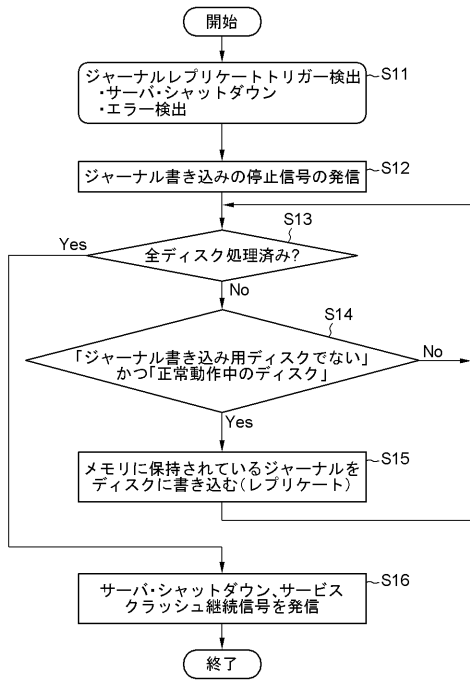
【図 1 3】



【図 1 4】



【図15】



---

フロントページの続き

(56)参考文献 特開2001-184176(JP,A)  
特開平09-282211(JP,A)  
特開2004-287648(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06 - 3/08  
G06F 13/00 - 13/42