

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5412902号
(P5412902)

(45) 発行日 平成26年2月12日(2014.2.12)

(24) 登録日 平成25年11月22日(2013.11.22)

(51) Int. Cl. F I
G06F 3/06 (2006.01) G06F 3/06
 G06F 3/06 301Z

請求項の数 7 (全 18 頁)

<p>(21) 出願番号 特願2009-63903 (P2009-63903) (22) 出願日 平成21年3月17日 (2009.3.17) (65) 公開番号 特開2010-218193 (P2010-218193A) (43) 公開日 平成22年9月30日 (2010.9.30) 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 (74) 代理人 100131428 弁理士 若山 剛 (72) 発明者 鈴木 秀明 東京都港区芝五丁目7番1号 日本電気株式会社社内 審査官 古河 雅輝</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】

ディスク装置毎に当該ディスク装置に対する単位時間あたりの稼働状況を計測して、稼働した前記ディスク装置を稼働ディスクと判定し、稼働していない前記ディスク装置を非稼働ディスクと判定すると共に、前記ディスク装置内の記憶領域が複数に分割された各区画における単位時間当たりの負荷を計測して、計測結果を記憶する性能監視部と、

前記性能監視部による計測結果に基づいて、稼働ディスクと判定された前記ディスク装置に形成された前記区画のうち、負荷が他の前記区画と比較して高い前記区画における負荷を軽減するよう、負荷が他の前記区画と比較して高い前記区画に格納されたデータの一部又は全部を、他の前記ディスク装置に格納して、データの再配置を行い、その後、前記性能監視部による計測結果に基づいて、稼働ディスクと判定された前記ディスク装置に形成された前記区画のうち、負荷が他の前記区画と比較して低い前記区画に格納されているデータを、稼働ディスクと判定された他の前記ディスク装置に格納するデータ再配置部と

を備えた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、

前記データ再配置部は、前記性能監視部による計測結果に基づいて、負荷が低い順から予め設定された順番までの前記区画に格納されているデータを、稼働ディスクと判定された他の前記ディスク装置に格納する、

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、

前記データ再配置部は、非稼働ディスクと判定された前記ディスク装置を含めて、前記データの再配置を行う、

ストレージシステム。

【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、

前記データ再配置部は、前記性能監視部による計測結果に基づいて、負荷が高い順から予め設定された順番までの前記区画における負荷を軽減するよう、前記各区画に格納されたデータの再配置を行う、

ストレージシステム。

【請求項 5】

請求項 1 乃至 4 のいずれか一項に記載のストレージシステムであって、

前記性能監視部は、前記各区画における単位時間あたりの入力開始から当該入力に対する出力完了までに要する時間であるターンアラウンドタイムの累積時間を、前記各区画における前記負荷として計測する、

ストレージシステム。

【請求項 6】

コンピュータに、

ディスク装置毎に当該ディスク装置に対する単位時間あたりの稼働状況を計測して、稼働した前記ディスク装置を稼働ディスクと判定し、稼働していない前記ディスク装置を非稼働ディスクと判定すると共に、前記ディスク装置内の記憶領域が複数に分割された各区画における単位時間当たりの負荷を計測して、計測結果を記憶する性能監視部と、

前記性能監視部による計測結果に基づいて、稼働ディスクと判定された前記ディスク装置に形成された前記区画のうち、負荷が他の前記区画と比較して高い前記区画における負荷を軽減するよう、負荷が他の前記区画と比較して高い前記区画に格納されたデータの一部又は全部を、他の前記ディスク装置に格納して、データの再配置を行い、その後、前記性能監視部による計測結果に基づいて、稼働ディスクと判定された前記ディスク装置に形成された前記区画のうち、負荷が他の前記区画と比較して低い前記区画に格納されているデータを、稼働ディスクと判定された他の前記ディスク装置に格納するデータ再配置部と

、
を実現させるためのプログラム。

【請求項 7】

ディスク装置毎に当該ディスク装置に対する単位時間あたりの稼働状況を計測して、稼働した前記ディスク装置を稼働ディスクと判定し、稼働していない前記ディスク装置を非稼働ディスクと判定すると共に、前記ディスク装置内の記憶領域が複数に分割された各区画に単位時間当たりの負荷を計測して、計測結果を記憶して、性能監視を行い、

前記計測結果に基づいて、稼働ディスクと判定された前記ディスク装置に形成された前記区画のうち、負荷が他の前記区画と比較して高い前記区画における負荷を軽減するよう、負荷が他の前記区画と比較して高い前記区画に格納されたデータの一部又は全部を、他の前記ディスク装置に格納し、その後、前記計測結果に基づいて、稼働ディスクと判定された前記ディスク装置に形成された前記区画のうち、負荷が他の前記区画と比較して低い前記区画に格納されているデータを、稼働ディスクと判定された他の前記ディスク装置に格納して、データの再配置を行う、

データ格納方法。

【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、複数の記憶装置を装備したストレージシステムに関する。

【背景技術】

【0002】

近年、データ記憶容量、信頼性の向上を図るべく、多数のディスク装置を備えたストレージシステムが開発されている。例えば、N A S (Network Attached Storage) 装置や、グリッドストレージシステムといったものがある。

【0003】

ここで、仮に、数百台のクライアントが起動時の特定処理で読み込むファイルが、ストレージシステム上に格納されているとする。すると、このファイルには、数百台のクライアントが同時に起動する際にアクセスするため、当該ファイルに負荷が集中し、ホットスポットとなる。ところが、このようなホットスポットが出現する時間は短く、1日平均のディスクビジー率でみると、システム全体のわずかなりソースしか占有しない。また、ストレージシステムでは、一般的に、容量がフルに近づくとき性能低下が起きやすくなるという問題がある。これに対応すべく、また、突発的な容量増加に耐えられるよう、70%以下など余裕をもって運用されるのが普通である。なお、I L M (情報ライフサイクル管理: Information Lifecycle

Management) の考え方では、大部分のデータはアクセス頻度が低いデータ(例えば1日~数ヶ月に1回以下)だと言われている。

【0004】

しかしながら、上述したように、ストレージシステム全体の稼働率(負荷・容量)が低い場合であっても、多数のディスク装置を備えているストレージシステムでは、これら多数のディスク装置を絶えず駆動する必要がある。その結果、ストレージシステムの消費電力が下がりにくい、という問題がある。

【0005】

一方で、特許文献1では、N A S 装置などのストレージシステムにおいて、省エネルギー化を図ることを目的とした技術を開示している。具体的には、ストレージシステムを構成する仮想ストレージ毎の負荷を計測し、負荷が高い仮想ストレージの処理を、他のストレージ処理部にて分割し、負荷が低い仮想ストレージの処理を他のストレージ処理部にて統合する。これにより、仮想ストレージシステムを制御しているストレージ処理部への電力供給を制御している。

【先行技術文献】

【特許文献】

【0006】

【特許文献1】特開2003-296153号公報

【発明の概要】

【発明が解決しようとする課題】

【0007】

ところが、上述した特許文献1に開示の技術では、仮想ストレージ単位で負荷を計測しており、当該仮想ストレージを制御するストレージ処理部に対する電力供給を制御しているだけであるため、物理的なディスク装置に対する電力供給を適切に制御していない。従って、ディスク装置の消費電力の低下を図ることができず、依然としてシステム全体の消費電力が下がりにくい、という問題が生じる。

【0008】

このため、本発明の目的は、上述した課題である、消費電力の低下が困難であることを解決す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】

【0009】

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、ディスク装置毎に当該ディスク装置に対する単位時間あたりの稼働状況を計測して、稼

10

20

30

40

50

働した上記ディスク装置を稼働ディスクと判定し、稼働していない上記ディスク装置を非稼働ディスクと判定すると共に、単位時間当たりの上記ディスク装置における負荷を計測して、計測結果を記憶する性能監視部と、

上記性能監視部による計測結果に基づいて、稼働ディスクと判定された上記ディスク装置に格納されているデータを、当該ディスク装置における負荷よりも高い負荷の他の上記ディスク装置に格納するデータ再配置部と、
を備える。

【0010】

また、本発明の他の形態であるプログラムは、
コンピュータに、

ディスク装置毎に当該ディスク装置に対する単位時間あたりの稼働状況を計測して、稼働した上記ディスク装置を稼働ディスクと判定し、稼働していない上記ディスク装置を非稼働ディスクと判定すると共に、単位時間当たりの上記ディスク装置における負荷を計測して、計測結果を記憶する性能監視部と、

上記性能監視部による計測結果に基づいて、稼働ディスクと判定された上記ディスク装置に格納されているデータを、当該ディスク装置における負荷よりも高い負荷の他の上記ディスク装置に格納するデータ再配置部と、
を実現させるためのプログラムである。

【0011】

また、本発明の他の形態であるデータ格納方法は、

ディスク装置毎に当該ディスク装置に対する単位時間あたりの稼働状況を計測して、稼働した上記ディスク装置を稼働ディスクと判定し、稼働していない上記ディスク装置を非稼働ディスクと判定すると共に、単位時間当たりの上記ディスク装置における負荷を計測して、計測結果を記憶して、性能監視を行い、

上記計測結果に基づいて、稼働ディスクと判定された上記ディスク装置に格納されているデータを、当該ディスク装置における負荷よりも高い負荷の他の上記ディスク装置に格納して、データの再配置を行う、
という構成を採る。

【発明の効果】

【0012】

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全体のさらなる省電力化を図ることができる。

【図面の簡単な説明】

【0013】

【図1】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2】図1に開示したストレージシステムにおけるディスク装置の負荷の計測の一例を示す図である。

【図3】図1に開示した再配置候補リスト記憶部に記憶されているデータの一例を示す図である。

【図4】図1に開示したストレージシステムに設定されているデータの一例を示す図である。

【図5】図1に開示した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6】図1に開示した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7】図1に開示したストレージシステムにおけるデータ再配置時の様子を示す図である。

【図8】図1に開示した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9】図1に開示したストレージシステムにおけるデータ再配置時の様子を示す図である。

【図10】図1に開示した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11】図1に開示したストレージシステムにおけるデータ再配置時の様子を示す図で

10

20

30

40

50

ある。

【図 1 2】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。

【発明を実施するための形態】

【0014】

<実施形態 1>

本発明の第 1 の実施形態を、図 1 乃至図 1 1 を参照して説明する。図 1 は、本実施形態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 2 は、ディスク装置の負荷の計測の一例を示す図である。図 3 乃至図 4 は、ストレージシステムに記憶されているデータの一例を示す図である。図 5 乃至図 1 1 は、ストレージシステムの動作を示す図である。

10

【0015】

[構成]

本実施形態におけるストレージシステム 10 は、図 1 に示すように、複数のディスク装置 6, 7, …, n を備えており、例えば、NAS 装置やグリッドストレージシステムによって実現されている。なお、図 1 では、ストレージシステムが、ディスク装置 6 等と一体的に構成されているよう図示しているが、本発明におけるストレージシステムは、ディスク装置 6 等と分離されており、ディスク装置 6 等に対するデータの記録再生制御を行うコンピュータにて構成されていてもよい。

【0016】

そして、ストレージシステム 10 は、装備された演算装置にプログラムが組み込まれることによって構築された、ファイル管理部 1 と、性能監視部 2 と、再配置候補リスト記憶部 3 と、ファイル再配置制御部 4 と、ファイル群選別制御部 5 と、を備えている。なお、上記各部 1 ~ 5 による処理に付随して記憶することが必要なデータを記憶するフラッシュメモリなどの記憶装置も備えている。以下、各構成について詳述する。

20

【0017】

上記ファイル管理部 1 は、ディスク装置 6 等内に格納されたファイルデータの格納位置、つまり、どのディスクのどの区画に格納されているかということを管理している。なお、ディスク装置 6 等は、図 1 の符号 6 1, 7 1, n 1 等に示すように、その記憶容量が複数の区画に分割されている。そして、各区画あるいは区画間をまたがって、ファイルデータ 6 2, 7 2, n 2 等が格納されている。例えば、一つのディスク装置の記憶容量は、1

30

000 程度の区画に分割されている。

【0018】

また、上記性能監視部 2 は、上述した各区画をユニークな区画番号により識別する。そして、性能監視部 2 は、一定期間（例えば、10 分）における区画毎の I/O（Input/Output）度数分布と、ディスク全体のビジー率を観測する。具体的に、性能監視部 2 は、該当するディスク装置の単位時間当たりの稼働率、つまり、一定期間あたりの稼働している時間の割合を、上記ディスクビジー率として算出する。そして、ディスクビジー率が「0」ではない、つまり、一定期間（単に時間）内に稼働したディスク装置を、稼働ディスクとして判定し、ディスクビジー率が「0」、つまり、一定期間（単位時間）内に稼働しなかったディスク装置を、非稼働ディスクとして判定する。さらに、性能監視部 2 あるいは

40

ストレージシステム 10 に装備された他の機能により、非稼働ディスクとして判定したディスク装置は停止される。つまり、非稼働ディスクであるディスク装置への電源供給は停止される。なお、非稼働ディスクは未使用という意味ではない。上述した単位時間あたりに稼働していないだけであって、データが格納されている可能性はある。

【0019】

なお、上述したように計測対象となるディスク装置は、ディスク装置単体であることに限定されず、複数のディスク装置であってもよい。つまり、グループ化された複数のディスク装置を一つのディスク装置として、そのディスクビジー率を計測してもよい。

【0020】

また、性能監視部 2 は、上述したように計測したディスクビジー率と、各区画毎の I/O

50

0度数分布とから、各区画毎のビジー率を算出して、この値を各区画の負荷とする。このとき、各区画に対するRead回数とWrite回数とに基づいて、各区画に対するRead処理時におけるReadビジー率、Write処理時におけるWriteビジー率、をそれぞれ算出する。

【0021】

ここで、図2を参照して、所定のディスク装置6の各区画におけるビジー率算出の一例を説明する。まず、図2(A)に示すように、ディスク装置6が区画#1~#7に分割されているとする。そして、このときのディスク装置6全体のビジー率が80%であると計測され、また、各区画#1~#7に対するRead回数とWrite回数とがそれぞれ図2(B)に示す回数であったとする。つまり、ディスク装置6を7つの区画#1~#7に分けた時の各区画におけるI/O度数が、図2(B)の上側の図に示すように、区画#1からReadが、4回、1回、8回、0回、0回、0回、4回であり、図2(B)の下側の図に示すように、区画#1からWriteが、0回、2回、0回、0回、2回、0回、0回であるとする。この場合に、ディスクビジー率(80%)を、各区画のI/O度数分布で比例分配することにより、各区画のビジー率を算出する。すると、上記の場合には、全I/O数を合計すると20回であるため、I/O一回で全体のビジー率80%うちの4%に貢献することとなる。従って、図2(C)に示すように、各区画のReadビジー率は、それぞれ12%、4%、32%、0%、0%、0%、16%となり、各区画のWriteビジー率は、それぞれ0%、8%、0%、0%、8%、0%、0%となる。このようにして、各区画のReadビジー率、Writeビジー率を算出し、当該各区画の負荷とする。

【0022】

なお、上述した各区画のビジー率の算出は、ディスクビジー率が「0」ではなく稼働ディスクと判定されたディスク装置に対して行う。そして、算出した区画ビジー率が「0」でなければ、性能監視部2は、当該算出した区画ビジー率(Readビジー率、Writeビジー率)を再配置候補リスト記憶部3に入力して記憶する。

【0023】

そして、上記再配置候補リスト記憶部3は、上述したように性能監視部2から受け取った各区画のビジー率、つまり、Readビジー率、Writeビジー率を、それぞれビジー率が高い順と低い順にソートして、それらの優先度を管理する。具体的には、図3に示すように、Read・Writeの二種類の処理において、それぞれ高ビジー率・低ビジー率の二種類の、計4種類の優先度リストを管理する。そして、各優先度リストにおいては、ストレージシステムにある全区画の中から、その値が高い順、及び、低い順に、予め設定された個数(例えば、5つ)ずつ保持している。このリストに登録される区画は、後述のファイル再配置制御部4にてファイルが再配置されるか、さらにより高い、もしくは、低い区画がリストに登録され、各リストの上位に位置する設定個数内から追い出されるまで、当該リスト中に存在し続ける。ただし、一つのリストに同じ区画番号は2度出現しない。例えば、リストをヒープ構造でもち、定期的に各エリアのビジー率をリストの最下位の値と比較して、リストの圏内であれば最下位の値と交換後にソートすることで、このような優先度リストを実現できる。

【0024】

ここで、再配置候補リスト記憶部3における入力された各区画におけるRead・Writeビジー率のリストのソート方法について説明する。まず、すでにリストの最大個数を越えて区画が記憶されている場合には、新たに入力された区画のビジー率と、記憶しているリストの最下位の区画のビジー率と比較し、これよりも高い値(区画が高い順に並んでいる場合)あるいは低い値(区画が低い順に並んでいる場合)であれば、区画番号とビジー率の値を交換する。つまり、区画番号をリストに追加する。その後、交換が発生したリストをソートし、重複する区画が既に登録されていれば、その区画の登録が古い方を削除する。

【0025】

例えば、最大リスト個数が「3」の高Readビジー率リストがあり、現在空である場合、前述のディスクの区画Readビジー率12%(#1)、4%(#2)、32%(#3)、16%(#7)を挿入するケースを考える。まず、リストが空なので、#1、#2、#3が

10

20

30

40

50

リストに登録されソートされる。リストの状況は、#3 #1 #2の順番となる。次に#7(16%)を挿入すると、リストは既に最大個数(3)なので、最下位の区画#2(4%)と追加区画#7(16%)のビジー率を比較する。ここでは、#7のビジー率が高いため、#2と#7とを交換する。そして、再度ソートを行い、新たなリスト#3 #7 #1を得る。なお、#7の前後に他に#7はないため、挿入した区画#7は重複しておらず、重複削除の必要は無い。

【0026】

なお、区画ビジー率のリストは、図3に示すように、ReadとWriteを別々に形成されていることに限定されない。例えば、Readだけ、あるいは、Writeだけに対応するリストを記憶していてもよく、あるいは、Read/Writeの区別無しで、上述した性能監視部2で区画毎のビジー率を計測して、再配置候補リスト記憶部3で管理してもよい。

10

【0027】

さらに、上記では、各区画の負荷を表す値として、区画毎のビジー率を利用したが、他の計測した値や算出した値を、各区画の負荷とし、これに基づいて、図3に示すような各区画のリストを管理してもよい。例えば、上記性能監視部2は、各区画における単位時間あたりの入力開始から当該入力に対する出力完了までに要する時間であるターンアラウンドタイム(TAT: Turn Around Time)を計測し、これらの累積時間を、各区画における負荷としてもよい。そして、この負荷の値が高い順や低い順に、ファイルの再配置を行う区画を選択するリストの並び替えを行ってもよい。このように、区画の負荷を表す値として、TATを用いた場合には、当該TATの値は故障気味のディスク装置で高くなること

20

【0028】

また、上記ファイル郡選別制御部4は、ある区画に含まれるファイルを特定する機能を有する。具体的には、ファイルシステムのレイアウト情報を利用し、区画番号から、その区画を利用しているファイル群やメタデータ群を選別できる。このような逆マップ情報は、昨今のファイルシステムでは内部情報として普通に存在している。

【0029】

また、上記ファイル再配置制御部5は、定期的に(例えば1日1回)、上述した再配置候補リスト記憶部3に記憶されたリストに基づいて、各区画に格納されているファイルの再配置を計画して、実行する機能を有する。具体的には、大きく分けて、以下の3つの処理A, B, Cを実行する。

30

【0030】

(処理A)

まず、再配置候補リスト記憶部3に基づいて、ビジー率の高い区画(一番高い区画、あるいは、高い順から所定の順番までの区画)を選択し、現在稼働中のディスクに負荷を分散した場合の効果を予測して計算する。具体的な負荷の分散方法としては、例えば、ファイル群の一部を別ディスクに移動する、高いWrite要求のある単一ファイルをRAID5(Redundant Arrays of Inexpensive Disks 5)などの分散方式で複数ディスクに分散する、高いRead要求のあるファイルの複製を複数ディスクに作成する、などの再配置処理が考えられる。また、予測計算では、高いRead要求の区画のみのファイルを分散、高いWrite要求の区画のみのファイルを分散、高いRead/Write両方の区画のファイルを分散した場合のそれぞれのケースで、予め設定された評価関数により効果を数値化し、この数値に基づいて適したファイル再配置処理を決定する。

40

【0031】

そして、データ再配置後の予測した評価等の値(例えば、評価関数による評価値や各区画のビジー率)が、予め設定された範囲内となる場合には、ファイルの再配置処理を実行する。その後、再配置に伴い、ファイル管理部1の情報を適切に変更し、再配置候補リスト記憶部3から当該区画を削除して、処理を終了する。

50

【 0 0 3 2 】

(処理 B)

上記処理 A で、ファイル再配置後の負荷が軽減されない場合には、この処理 B を行う。この処理 B では、まず、再配置候補リスト記憶部 3 から、ビジー率の高い区画を選択し、現在、非稼働であるディスクを含めて負荷を分散した場合の効果と、ビジー率と消費電力に関する評価関数 $E V(N)$ を用いて、予測計算する。なお、この評価関数 $E V(N)$ における N は非稼働ディスクを含めた使用するディスク台数であるため、すでに現在の稼働ディスク数が下限値となる。

【 0 0 3 3 】

そして、上記評価関数 (N) は、ディスク台数が増えることでビジー率が下がる効果と、消費電力が増える効果がバランスした最適のディスク稼働台数が得られるものを使う。例えば、ビジー率とディスク数の定数倍を足す、図 4 に示すような評価関数 $E V(N)$ が用いられ、この評価関数に基づく評価値は、値が低い方が評価がよいこととなる。従って、図 4 の例では、 N_{max} が最も効率のよいディスク台数を示している。なお、性能や消費電力が予め決めた閾値を外れる場合は、 $E V(N)$ を 0 とすることで、システムの動作をコントロールできる。

【 0 0 3 4 】

そして、上述した評価関数 (N) を利用した評価の結果、一台以上の非稼働ディスクを利用してファイルの再配置を行うべきと判断された場合、ファイルの再配置を行う。その後、再配置に伴い、ファイル管理部 1 の情報を適切に変更し、再配置候補リスト記憶部 3 から当該区画を削除して、処理を終了する。

【 0 0 3 5 】

なお、上述した処理 A、B による高ビジー率の区画に対して、負荷を軽減するファイルの再配置処理は、上述した処理に限定されない。

【 0 0 3 6 】

(処理 C)

上記処理 A、B を実行した後に、この処理 C を実行する。この処理 C では、再配置候補リスト記憶部 3 から、ビジー率の低い区画 (一番低い区画、あるいは、低い順から所定の順番までの区画) を選択し、その区画上のファイル群を別の稼働ディスクに移動した場合の効果を表す評価の値を、予測計算する。そして、計算後のビジー率などの値が予め決められた範囲に入る場合には、ファイルの再配置を行い、再配置候補リスト記憶部 3 から当該区画を削除する。このとき、選択された区画のファイルの再配置は、稼働ディスクであって、選択された区画のビジー率よりも高いビジー率の区画に移動することで実行し、特に、ファイルが移動された後の区画が形成されたディスク装置を停止できる可能性が高いディスク装置であるとよい。

【 0 0 3 7 】

なお、上述したストレージシステムは、図 1 に示すように、各処理を行う機能が 1 台のコンピュータに装備されている場合を説明したが、各機能が複数台のコンピュータに装備されていてもよい。例えば、上記ファイル管理部 1 が独立したサーバコンピュータに存在し、クライアントからファイルの物理位置を隠蔽するメタデータサーバの役割を果たす、分散ファイルシステムとして構成されていてもよい。

【 0 0 3 8 】

また、上述したストレージシステムでは、各ディスク装置を、独立したファイルシステムとし、ファイル管理部 1 が全体を統括構成する構成も考えられる。例えば、画像検索システムを考えた場合に、ファイル管理部 1 を検索処理部と考え、アクセスパターンに最適な画像データのディスク配置を行う事ができる。

【 0 0 3 9 】

[動作]

次に、上述したストレージシステムの動作を説明する。はじめに、性能監視部 2 にて一定時間ごと (例えば、10 分毎) に実行される動作について、図 5 のフローチャートを参

10

20

30

40

50

照して説明する。

【 0 0 4 0 】

まず、各ディスク装置のビジー率と、各区画のRead・WriteのI/O度数分布を取得して記録する(ステップS1)。続いて、ディスク装置のビジー率が「0」ではなく、稼働ディスクと判定された場合には(ステップS2でYes)、そのディスク装置に形成された区画毎のビジー率を、ディスクのビジー率を各区画のI/O度数分布で比例分配することにより算出する(ステップS3)。そして、算出した区画ビジー率が、「0」でなければ、再配置候補リスト記憶部3に入力する(ステップS4)。このとき、再配置候補リスト記憶部3では、上述したように、各リストを高い順、あるいは、低い順に、ソートする。以上の処理が、一定時間ごとに実行される。

10

【 0 0 4 1 】

続いて、ファイル再配置制御部4にて定期的に(例えば1日1回)実行される動作について、図6乃至図11に示すフローチャート及びファイル再配置の様子を示す図を参照して、説明する。

【 0 0 4 2 】

はじめに、上述した処理Aを実行するが、このときの動作を、図6及び図7を参照して説明する。まず、再配置候補リスト記憶部3から、ビジー率の高い区画を選択する(ステップS11)。このとき、例えば、図7(A)の符号6に示すディスク装置内の区画#11が選択されたとする。そして、この区画における負荷を、現在稼働中の他のディスク装置(区画)に分散した場合の効果を、ビジー率にて予測計算する(ステップS12)。

20

【 0 0 4 3 】

ここでは、図7(A)の矢印に示すように、選択された区画#11内のファイルを、稼働中の他のディスク装置7内の区画#12に移動して、負荷を分散する場合を考える。このような負荷分散処理を行った場合の効果を、予め記憶された評価関数などの評価情報に基づいて評価値として算出する。そして、この評価値が、再配置により効果が上がると判定できる予め設定された基準を満たすか否かを判断する(ステップS13)。基準を満たす場合には(ステップS13でYes)、予測した負荷分散処理を実際に行う(ステップS14)。すると、図7(B)に示すように、区画#11内に格納された1つのファイルが、稼働中の他のディスク装置7内の区画#12に移動して格納された状態となる。

【 0 0 4 4 】

その後は、ファイルの再配置に伴い、ファイル管理部1の情報を適切に変更する。また、再配置候補リスト記憶部3内の高ビジー率リストから当該区画を削除して(ステップS15)、処理を終了する。

30

【 0 0 4 5 】

ここで、上述したステップS12で予測した分散処理の効果が基準を満たさなかった場合には(ステップS13でNo)、上記処理Bを実行する。このときの動作を、図8及び図9を参照して説明する。

【 0 0 4 6 】

まず、再配置候補リスト記憶部3から、ビジー率の高い区画を選択する(ステップS21)。このとき、例えば、図9(A)の符号6に示すディスク装置内の区画#11が選択されたとする。そして、この区画における負荷を、上記処理Aとは異なり、今度は、現在、非稼働であるディスク装置も含めて、負荷を分散した場合の効果を予測計算する(ステップS22)。

40

【 0 0 4 7 】

ここでは、図9(A)の矢印に示すように、選択された区画#11内のファイルを、非稼働である他のディスク装置9内の区画#14に移動して、負荷を分散する場合を考える。このような負荷分散処理を行った場合の効果を、予め記憶された評価関数などの評価情報に基づいて評価値として算出する。そして、この評価値が、再配置により効果が上がると判定できる予め設定された基準を満たすか否かを判断する(ステップS23)。基準を満たす場合には(ステップS23でYes)、予測した負荷分散処理を実際に行う(

50

ステップS24)。すると、図9(B)に示すように、区画#11内に格納された1つのファイルが、非稼働の他のディスク装置9内の区画#14に移動して格納された状態となり、また、非稼働であったディスク装置9は稼働状態となる。

【0048】

その後は、ファイルの再配置に伴い、ファイル管理部1の情報を適切に変更する。また、再配置候補リスト記憶部3内の高ビジー率リストから当該区画を削除して(ステップS25)、処理を終了する。

【0049】

以上のように、本実施形態におけるストレージシステムでは、負荷の高い区画に格納されたファイルを分散したり、あるいは複製して、他のディスク装置に格納するなどして、データの再配置を行っている。これにより、該当する区画が形成されたディスク装置の負荷が軽減されることとなる。その結果、ストレージシステム全体において電力が効率的に利用され、省電力化を図ることができる。特に、非稼働ディスクを増やして負荷分散するさいには、実際の分散処理の前に評価関数などを用いて性能予測を行うことで、性能及び省電力化のバランスを維持することができる。

【0050】

続いて、上記処理Aあるいは処理Bが実行された後に、上記処理Cを実行する。このときの動作を、図10及び図11を参照して説明する。

【0051】

まず、再配置候補リスト記憶部3から、ビジー率の低い区画を選択する(ステップS31)。このとき、例えば、図11(A)の符号8に示すディスク装置内の区画#13が選択されたとする。そして、この区画に配置されているファイルを、稼働中の他のディスク装置(区画)に移動する場合の効果を実測計算する(ステップS32)。

【0052】

ここでは、図11(A)の矢印に示すように、選択された区画#13内のファイルを、稼働中の他のディスク装置7内の区画#12に移動する場合を考える。このような負荷分散処理を行った場合の効果を実測計算された評価関数などの評価情報に基づいて評価値として算出して、この評価値が、再配置により効果が上がると判定できる予め設定された基準を満たすか否かを判断する(ステップS33)。基準を満たす場合には(ステップS33でYes)、予測したファイルの移動を実際に行う(ステップS34)。すると、図11(B)に示すように、区画#13内に格納されたファイルが、稼働中の他のディスク装置7内の区画#12に移動して格納された状態となる。

【0053】

その後は、ファイルの再配置に伴い、ファイル管理部1の情報を適切に変更する。また、再配置候補リスト記憶部3内の低ビジー率リストから当該区画を削除して(ステップS35)、処理を終了する。

【0054】

そして、その後は、上述したようにもともと負荷の低かった区画#13内のファイルが他の区画に移動されることから、当該区画#13にはその後の負荷がさらに低減される。従って、後に性能監視部2にてディスク装置のビジー率が計測されたときに、区画#13が位置するディスク装置8のビジー率が「0」であった場合には、当該ディスク装置8を非稼働ディスクと判定でき、その後は、電源供給が停止される。

【0055】

以上により、本発明によると、負荷の低い区画内のファイルを他のディスク装置に移動することで、非稼働ディスクを効率的に増加させることができる。従って、非稼働となったディスク装置に対する電源供給を停止することができ、ストレージシステム全体の省電力化を図ることができる。

【0056】

ここで、上述したように、区画ビジー率を計算し、ソート済みリストを管理するための処理量を見積もる。

10

20

30

40

50

単位時間あたりの処理：1000台の1TBディスク装置を1000個のエリアに分け、全ディスクが稼働ディスクで、10分に一回ビジー率を集計した場合

ディスクビジー率取得回数：1K個/10分 = 1.67個/秒

各エリアのサイズ：1GB

エリア数：1M個

ビジー率の計算：1M個/10分 = 1.67K個/秒

上記のうち、10%がリストの圏内に入ったとすると、リストへのエリア挿入と再ソートが167回/秒である。そして、リストの長さが短ければ、昨今のCPUにとって大きな負荷ではない。また、稼働ディスク数が1/10になれば、管理対象の区画も1/10になり管理処理は大幅に削減される。このようにディスクが1000台つながっていても、少ない計算と統計情報で、再配置処理を実行することができる。

10

【0057】

<実施形態2>

本発明の第1の実施形態を、図12を参照して説明する。図12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。

【0058】

図1に示すように、本実施形態におけるストレージシステム100は、

ディスク装置111, 112毎に当該ディスク装置に対する単位時間あたりの稼働状況を計測して、稼働した上記ディスク装置を稼働ディスクと判定し、稼働していない上記ディスク装置を非稼働ディスクと判定すると共に、単位時間当たりの上記ディスク装置における負荷を計測して、計測結果を記憶する性能監視部101と、

20

上記性能監視部による計測結果に基づいて、稼働ディスクと判定された上記ディスク装置に格納されているデータを、当該ディスク装置における負荷よりも高い負荷の他の上記ディスク装置に格納するデータ再配置部102と、を備える。

【0059】

そして、上記ストレージシステムでは、

上記性能監視部は、上記ディスク装置内の記憶領域が複数に分割された各区画における負荷を計測して、計測結果を記憶し、

30

上記データ再配置部は、上記性能監視部による計測結果に基づいて、稼働ディスクと判定された上記ディスク装置に形成された上記区画のうち、負荷が他の上記区画と比較して低い上記区画に格納されているデータを、稼働ディスクと判定された他の上記ディスク装置に格納する、という構成を採る。

【0060】

また、上記ストレージシステムでは、

上記データ再配置部は、上記性能監視部による計測結果に基づいて、負荷が低い順から予め設定された順番までの上記区画に格納されているデータを、稼働ディスクと判定された他の上記ディスク装置に格納する、という構成を採る。

40

【0061】

上記発明によると、まず、ストレージシステムでは、ディスク装置毎に単位時間当たりの稼働状況を計測する。そして、単位時間あたりに稼働したディスク装置を稼働ディスクと判定し、単位時間あたりに稼働していないディスク装置を非稼働ディスクと判定する。これにより、例えば、非稼働ディスクと判定されたディスク装置を停止することで、ストレージシステム全体の省電力化を図ることができる。

【0062】

また、ストレージシステムでは、上述したようにディスク装置毎の稼働状況を計測すると共に、ディスク装置における負荷、例えば、ディスク装置内に形成された各区画毎にお

50

ける負荷を計測して記憶している。そして、稼働ディスクであるディスク装置内に形成された区画のうち、負荷が他の区画と比較して低い区画に格納されているデータを、稼働ディスクである他のディスク装置に格納する。なお、上述したデータの移動は、負荷が低い区画順に実行してもよい。

【0063】

これにより、データを稼働中の他のディスク装置に移動させたディスク装置は、もともと負荷が低かったため、その後、アクセスが減少して、非稼働ディスクとなる可能性が高くなる。そして、非稼働ディスクとなった場合には、上述したように、そのディスク装置を停止する。これにより、ストレージシステム全体のさらなる省電力化を図ることができる。

10

【0064】

また、上記ストレージシステムでは、

上記データ再配置部は、上記性能監視部による計測結果に基づいて、稼働ディスクと判定された上記ディスク装置に形成された上記区画のうち、負荷が他の上記区画と比較して高い上記区画における負荷を軽減するよう、上記各区画に格納されたデータの再配置を行う、

という構成を採る。

【0065】

また、上記ストレージシステムの他の態様において、

上記データ再配置部は、稼働ディスクと判定された上記ディスク装置に形成された上記区画のうち、負荷が他の上記区画と比較して高い上記区画に格納されたデータの一部又は全部を、他の上記ディスク装置に格納して、データの再配置を行う、

という構成を採る。

20

【0066】

また、上記ストレージシステムでは、

上記データ再配置部は、非稼働ディスクと判定された上記ディスク装置を含めて、上記データの再配置を行う、

という構成を採る。

【0067】

また、上記ストレージシステムでは、

上記データ再配置部は、上記性能監視部による計測結果に基づいて、負荷が高い順から予め設定された順番までの上記区画における負荷を軽減するよう、上記各区画に格納されたデータの再配置を行う、

という構成を採る。

30

【0068】

また、上記ストレージシステムでは、

上記性能監視部は、上記各区画における単位時間あたりの入力開始から当該入力に対する出力完了までに要する時間であるターンアラウンドタイムの累積時間を、上記各区画における上記負荷として計測する、

という構成を採る。

40

【0069】

上記発明によると、ストレージシステムでは、負荷の高い区画に格納されたデータを分散あるいは複製して他のディスク装置に格納するなどして、データの再配置を行う。これにより、該当する区画が形成されたディスク装置の負荷が軽減されるため、ストレージシステム全体において電力が効率的に利用され、省電力化を図ることができる。

【0070】

そして、上述したストレージシステムは、コンピュータにプログラムが組み込まれることで実現できる。具体的に、本発明の他の形態であるプログラムは、コンピュータに、

ディスク装置毎に当該ディスク装置に対する単位時間あたりの稼働状況を計測して、稼働した上記ディスク装置を稼働ディスクと判定し、稼働していない上記ディスク装置を非

50

稼働ディスクと判定すると共に、単位時間当たりの上記ディスク装置における負荷を計測して、計測結果を記憶する性能監視部と、

上記性能監視部による計測結果に基づいて、稼働ディスクと判定された上記ディスク装置に格納されているデータを、当該ディスク装置における負荷よりも高い負荷の他の上記ディスク装置に格納するデータ再配置部と、
を実現させるためのプログラムである。

【0071】

そして、上記プログラムでは、

上記性能監視部は、上記ディスク装置内の記憶領域が複数に分割された各区画における負荷を計測して、計測結果を記憶し、

上記データ再配置部は、上記性能監視部による計測結果に基づいて、稼働ディスクと判定された上記ディスク装置に形成された上記区画のうち、負荷が他の上記区画と比較して低い上記区画に格納されているデータを、稼働ディスクと判定された他の上記ディスク装置に格納する、
という構成を採る。

【0072】

また、上記プログラムでは、

上記データ再配置部は、上記性能監視部による計測結果に基づいて、稼働ディスクと判定された上記ディスク装置に形成された上記区画のうち、負荷が他の上記区画と比較して高い上記区画における負荷を軽減するよう、上記各区画に格納されたデータの再配置を行う、
という構成を採る。

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格納方法は、

ディスク装置毎に当該ディスク装置に対する単位時間あたりの稼働状況を計測して、稼働した上記ディスク装置を稼働ディスクと判定し、稼働していない上記ディスク装置を非稼働ディスクと判定すると共に、単位時間当たりの上記ディスク装置における負荷を計測して、計測結果を記憶して、性能監視を行い、

上記計測結果に基づいて、稼働ディスクと判定された上記ディスク装置に格納されているデータを、当該ディスク装置における負荷よりも高い負荷の他の上記ディスク装置に格納して、データの再配置を行う。

【0073】

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格納方法は、

上記性能監視時に、上記ディスク装置内の記憶領域が複数に分割された各区画における負荷を計測して、計測結果を記憶し、

上記データの再配置時に、上記性能監視による計測結果に基づいて、稼働ディスクと判定された上記ディスク装置に形成された上記区画のうち、負荷が他の上記区画と比較して低い上記区画に格納されているデータを、稼働ディスクと判定された他の上記ディスク装置に格納する、
という構成を採る。

【0074】

そして、上記データ格納方法の他の態様において、

上記データ再配置時は、上記性能監視による計測結果に基づいて、稼働ディスクと判定された上記ディスク装置に形成された上記区画のうち、負荷が他の上記区画と比較して高い上記区画における負荷を軽減するよう、上記各区画に格納されたデータの再配置を行う、
という構成を採る。

【0075】

上述した構成を有する、プログラム、又は、データ格納方法、の発明であっても、上記

10

20

30

40

50

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。

【産業上の利用可能性】

【0076】

本発明は、ディスクアレイ装置、NAS装置、アグリゲートNAS装置、グリッドストレージなど、複数のディスク装置を搭載したストレージシステムに適用することができ、産業上の利用可能性を有する。

【符号の説明】

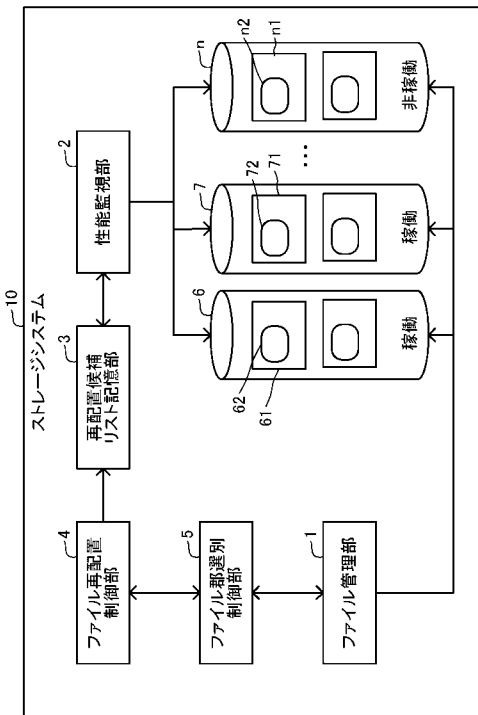
【0077】

- 1 ファイル管理部
- 2 性能監視部
- 3 再配置候補リスト記憶部
- 4 ファイル再配置制御部
- 5 ファイル郡選別制御部
- 6, 7, n ディスク装置
- 10, 100 ストレージシステム
- 101 性能監視部
- 102 データ再配置制御部
- 111, 112 ディスク装置

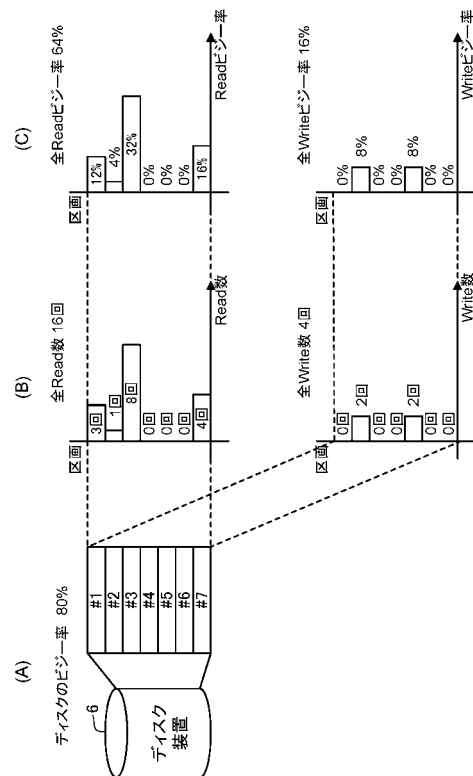
10

20

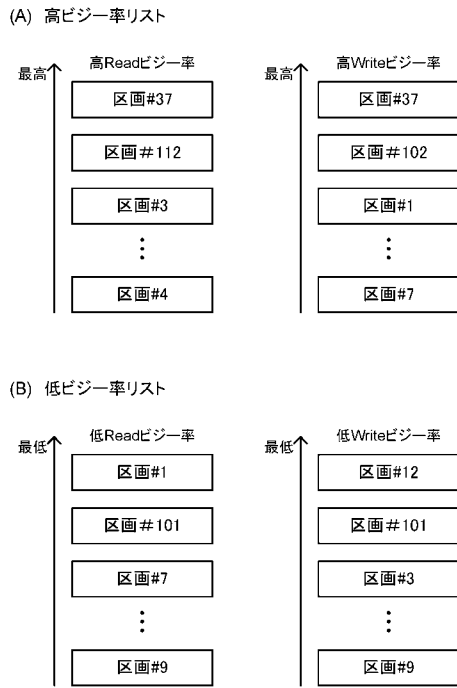
【図1】



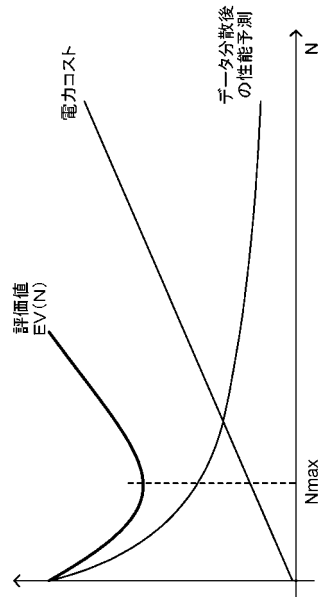
【図2】



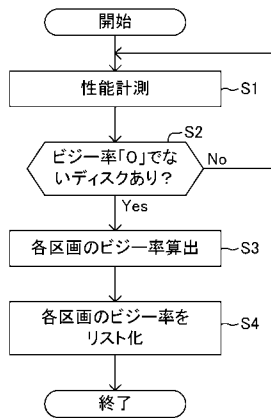
【図3】



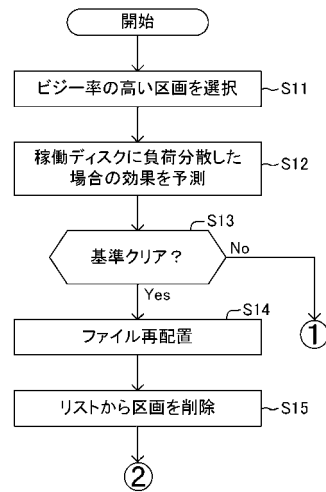
【図4】



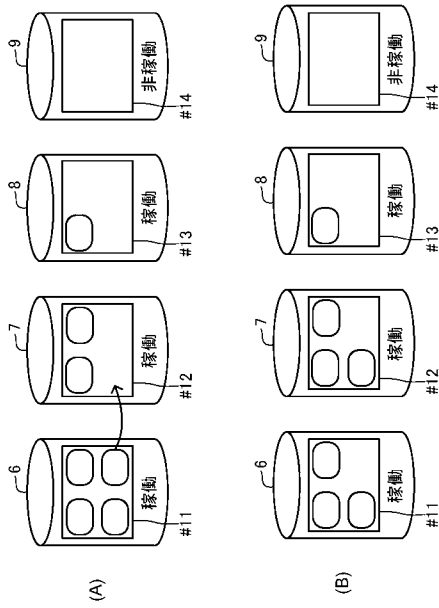
【図5】



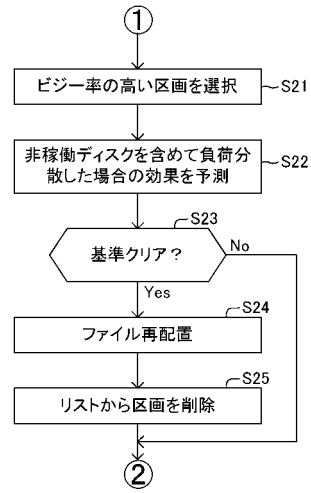
【図6】



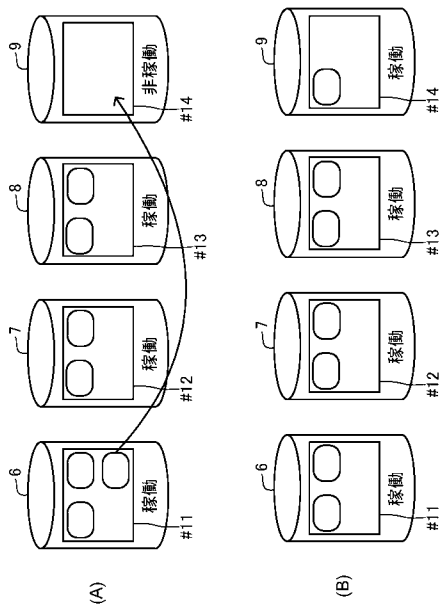
【図7】



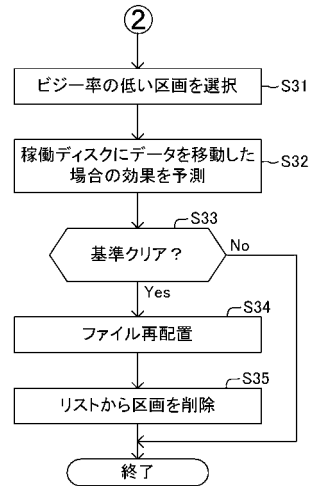
【図8】



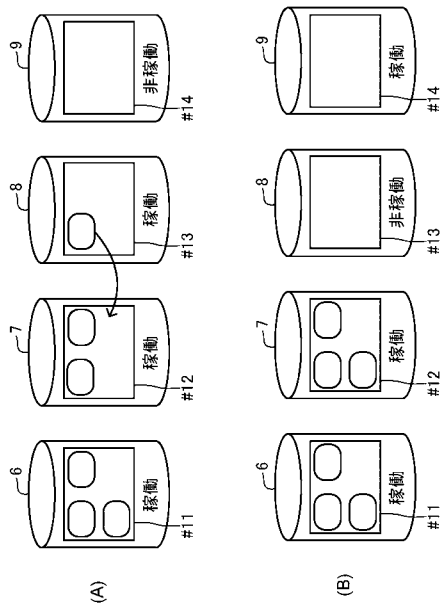
【図9】



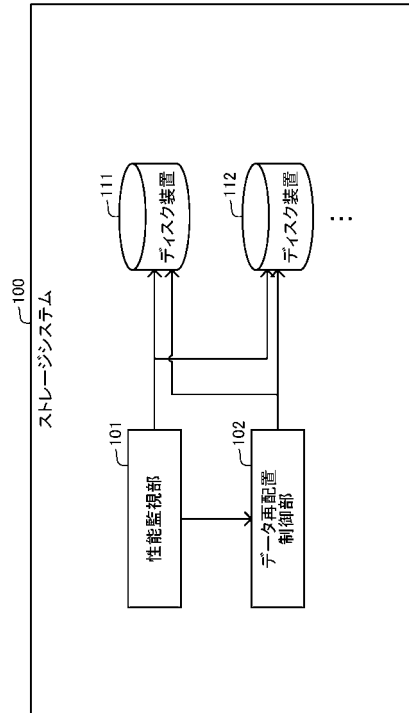
【図10】



【図 11】



【図 12】



フロントページの続き

- (56)参考文献 特開2001-093220(JP,A)
特開2000-231454(JP,A)
特開2008-217575(JP,A)
特開2008-146141(JP,A)
特開2008-003719(JP,A)
特開2007-293486(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06 - 3/08
G06F 12/00
G06F 13/10 - 13/14