

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6136460号
(P6136460)

(45) 発行日 平成29年5月31日(2017.5.31)

(24) 登録日 平成29年5月12日(2017.5.12)

(51) Int.Cl.		F I			
G06F	3/06	(2006.01)	G06F	3/06	304P
G06F	3/08	(2006.01)	G06F	3/08	H
			G06F	3/06	54O
			G06F	3/06	304Z

請求項の数 9 (全 24 頁)

(21) 出願番号	特願2013-70675 (P2013-70675)	(73) 特許権者	000005223
(22) 出願日	平成25年3月28日 (2013.3.28)		富士通株式会社
(65) 公開番号	特開2014-194667 (P2014-194667A)		神奈川県川崎市中原区上小田中4丁目1番1号
(43) 公開日	平成26年10月9日 (2014.10.9)	(74) 代理人	100089118
審査請求日	平成27年12月4日 (2015.12.4)		弁理士 酒井 宏明
		(72) 発明者	小野 貴継
			神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		審査官	田上 隆一

最終頁に続く

(54) 【発明の名称】 情報処理装置、情報処理装置の制御プログラムおよび情報処理装置の制御方法

(57) 【特許請求の範囲】

【請求項 1】

データが読み書きされる半導体記憶装置と、前記半導体記憶装置の予備である予備半導体記憶装置とを有する情報処理装置において、

前記半導体記憶装置について、書き込まれたデータを消去できる残余回数に基づいた余命指標値、および、前記半導体記憶装置に対する読出しと書き込みとの読み書きの比率を収集する収集部と、

前記収集部が収集した前記比率においてデータの読出しの比率が多い場合、第1の閾値を前記半導体記憶装置の交換基準となる閾値として決定し、データの書き込みの比率が多い場合、前記第1の閾値よりも値が大きい第2の閾値を前記半導体記憶装置の交換基準となる閾値として決定する決定部と、

前記収集部が収集した前記余命指標値が、前記決定部が決定した閾値よりも短い場合は、前記半導体記憶装置を前記予備半導体記憶装置と交換する交換部と

を有することを特徴とする情報処理装置。

【請求項 2】

前記情報処理装置はさらに、

前記予備半導体記憶装置を複数有し、

前記決定部は、

収集した前記比率においてデータの読出しの比率が多い場合、前記複数の予備半導体記憶装置のうち、余命指標値が前記第1の閾値よりも大きく、かつ、前記第2の閾値よりも

小さい予備半導体記憶装置を選択し、

前記交換部は、

前記半導体記憶装置を前記決定部が選択した予備半導体記憶装置と交換することを特徴とする請求項 1 に記載の情報処理装置。

【請求項 3】

前記決定部は、

収集した前記比率においてデータの書き込みの比率が多い場合、前記複数の予備半導体記憶装置のうち、余命が前記第 2 の閾値よりも大きい予備半導体記憶装置を選択することを特徴とする請求項 2 に記載の情報処理装置。

【請求項 4】

前記収集部は、

前記比率として、単位時間あたりに前記半導体記憶装置から読出されたデータ量と前記半導体記憶装置に書き込まれたデータ量との比率を収集し、

前記決定部は、

単位時間あたりの前記半導体記憶装置から読み出したデータ量が所定の閾値よりも多い場合、データの読出しの比率が多いと判定することを特徴とする請求項 1 ~ 3 のいずれか 1 つに記載の情報処理装置。

【請求項 5】

前記収集部は、

前記比率として、単位時間あたりに前記半導体記憶装置からデータを読出した回数とデータを書き込んだ回数との比率を収集し、

前記決定部は、

単位時間あたりの前記データの読出し回数が所定の閾値よりも多い場合、データの読出しの比率が多いと判定することを特徴とする請求項 1 ~ 3 のいずれか 1 つに記載の情報処理装置。

【請求項 6】

前記情報処理装置はさらに、

データが読み書きされる複数の半導体記憶装置と、

前記半導体記憶装置が、二重化されている場合、前記交換部が、前記半導体記憶装置を前記予備の半導体記憶装置と交換する間、前記半導体記憶装置に対する書き込みデータを保持する保持部と、

前記交換部が、交換対象の半導体記憶装置が記憶するデータと、前記保持部が保持したデータとを、前記予備半導体記憶装置と、前記複数の半導体記憶装置のうち前記予備半導体記憶装置と二重化される半導体記憶装置とに複製する複製部と

を有することを特徴とする請求項 1 ~ 5 のいずれか 1 つに記載の情報処理装置。

【請求項 7】

前記情報処理装置はさらに、

データが読み書きされる複数の半導体記憶装置を有し、

前記交換部は、

前記複数の半導体記憶装置に、データとデータを復元するためのパリティとが分散して格納されている場合、前記複数の半導体記憶装置のうち、いずれかの半導体記憶装置を前記予備半導体記憶装置と交換したとき、前記複数の半導体記憶装置のうち、前記予備半導体記憶装置以外の半導体記憶装置が記憶するデータとパリティとを用いて、前記予備半導体記憶装置が記憶するデータを生成することを特徴とする請求項 1 ~ 5 のいずれか 1 つに記載の情報処理装置。

【請求項 8】

データが読み書きされる半導体記憶装置と、前記半導体記憶装置の予備である予備半導体記憶装置とを有する情報処理装置の制御プログラムにおいて、

前記情報処理装置に、

前記半導体記憶装置について、書き込まれたデータを消去できる残余回数に基づいた余

10

20

30

40

50

命指標値、および、前記半導体記憶装置に対する読出しと書き込みとの読み書きの比率を収集させ、

収集された前記比率においてデータの読出しの比率が多い場合、第１の閾値を前記半導体記憶装置の交換基準となる閾値として決定し、データの書き込みの比率が多い場合、前記第１の閾値よりも値が大きい第２の閾値を前記半導体記憶装置の交換基準となる閾値として決定させ、

収集された前記余命指標値が、決定された前記閾値よりも短い場合は、前記半導体記憶装置を前記予備半導体記憶装置と交換させる

ことを特徴とする情報処理装置の制御プログラム。

【請求項９】

データが読み書きされる半導体記憶装置と、前記半導体記憶装置の予備である予備半導体記憶装置とを有する情報処理装置の制御方法において、

前記情報処理装置が、

前記半導体記憶装置について、書き込まれたデータを消去できる残余回数に基づいた余命指標値、および、前記半導体記憶装置に対する読出しと書き込みとの読み書きの比率を収集し、

収集された前記比率においてデータの読出しの比率が多い場合、第１の閾値を前記半導体記憶装置の交換基準となる閾値として決定し、データの書き込みの比率が多い場合、前記第１の閾値よりも値が大きい第２の閾値を前記半導体記憶装置の交換基準となる閾値として決定し、

収集された前記余命指標値が、決定された前記閾値よりも短い場合は、前記半導体記憶装置を前記予備半導体記憶装置と交換する

ことを特徴とする情報処理装置の制御方法。

【発明の詳細な説明】

【技術分野】

【０００１】

本発明は、情報処理装置、情報処理装置の制御プログラムおよび情報処理装置の制御方法に関する。

【背景技術】

【０００２】

従来、不揮発性の半導体を含む記憶装置をストレージとして用いる情報処理装置が知られている。このような情報処理装置の一例として、ＮＡＮＤ型のフラッシュメモリを有するＳＳＤ（Solid State Drive）をストレージとして用いる情報処理装置が知られている。

【０００３】

ここで、ＮＡＮＤ型のフラッシュメモリでは、データを保持する複数のメモリセルを含むページ単位で書き込み及び読出しが行われ、複数のページを含むブロック単位でデータの消去が行われる。しかし、メモリセルは、データの書き換えを行う度に劣化するので、データの書き換えを何度も実行すると、正常に情報を記録できなくなる。このため、同じメモリセルに対するデータの書き換えが頻繁に発生すると、メモリセルが正常に情報を記録できなくなり、正常に情報を記録できなくなったメモリセルを含むブロックが不良ブロックとなる。

【０００４】

そこで、メモリセルへの書き込み回数や消去回数を平坦化することで、不良ブロックの発生を防ぎ、ＳＳＤの寿命を改善するウェアレベリングの技術が知られている。例えば、ウェアレベリングの技術が適用されたＳＳＤは、頻繁に更新されるブロックに格納されたデータを更新回数が少ないブロックに移動することで、ＳＳＤが有するメモリセル全体の更新回数を平坦化する。

【先行技術文献】

【特許文献】

10

20

30

40

50

【 0 0 0 5 】

【特許文献 1】特開 2 0 0 7 - 3 2 3 2 2 4 号公報

【発明の概要】

【発明が解決しようとする課題】

【 0 0 0 6 】

しかしながら、上述したウェアレベリングの技術は、1つのSSD内における書き込み回数や消去回数の制限を根本的に解決するものではなく、SSDの寿命を延命するに過ぎない。このため、情報処理装置は、フラッシュメモリの書き込み回数や消去回数の制限を考慮せずに1つのSSDを使用し続けることができない。

【 0 0 0 7 】

本願は、上述した問題に鑑みてなされたものであり、1つの側面では、寿命を考慮せずにSSDをストレージとして利用することを目的とする。

【課題を解決するための手段】

【 0 0 0 8 】

1つの側面では、データが読み書きされる半導体記憶装置と、半導体記憶装置の予備である予備半導体記憶装置とを有する情報処理装置である。また、情報処理装置は、半導体記憶装置について、書き込まれたデータを消去できる残余回数に基づいた余命指標値、および、半導体記憶装置に対する読出しと書き込みとの読み書きの比率を収集する。また、情報処理装置は、収集した比率においてデータの読出しの比率が多い場合、第1の閾値を前記半導体記憶装置の交換基準となる閾値として決定し、データの書き込みの比率が多い場合、前記第1の閾値よりも値が大きい第2の閾値を前記半導体記憶装置の交換基準となる閾値として決定する。そして、情報処理装置は、収集した余命指標値が、決定した閾値よりも短い場合は、半導体記憶装置を予備半導体記憶装置と交換する。

【発明の効果】

【 0 0 0 9 】

1つの側面では、寿命を考慮せずにSSDをストレージとして利用することができる。

【図面の簡単な説明】

【 0 0 1 0 】

【図 1】図 1 は、実施例 1 に係る情報処理装置を説明する図である。

【図 2】図 2 は、実施例 1 に係る情報処理装置が実行する処理を説明する図である。

【図 3】図 3 は、実施例 1 に係る管理サーバの機能構成を説明する図である。

【図 4】図 4 は、余命管理テーブルの一例を説明する図である。

【図 5】図 5 は、ワークロード管理テーブルの一例を説明する図である。

【図 6】図 6 は、実施例 1 に係る情報処理装置が実行する処理の流れを説明するフローチャートである。

【図 7】図 7 は、選択部が実行する処理の流れを説明するための第 1 のフローチャートである。

【図 8】図 8 は、選択部が実行する処理の流れを説明するための第 2 のフローチャートである。

【図 9】図 9 は、実施例 2 に係る情報処理装置が実行する処理の流れを説明するフローチャートである。

【図 1 0】図 1 0 は、制御プログラムを実行するコンピュータを説明する図である。

【発明を実施するための形態】

【 0 0 1 1 】

以下に添付図面を参照して本願に係る情報処理装置、情報処理装置の制御プログラムおよび情報処理装置の制御方法の実施例を図面に基づいて詳細に説明する。なお、この実施例により開示技術が限定されるものではない。また、各実施例は、矛盾しない範囲で適宜組みあわせてもよい。

【実施例 1】

【 0 0 1 2 】

10

20

30

40

50

以下の実施例 1 では、図 1 を用いて、本願に係る情報処理装置を説明する。図 1 は、実施例 1 に係る情報処理装置を説明する図である。図 1 に示すように、情報処理装置 1 は、CPU (Central Processing Unit) プール 2、ディスクエリアネットワーク 3、ディスクプール 4、管理ネットワーク 7、管理サーバ 10 を有する。

【0013】

また、CPU プール 2 は、複数のノード 5 a ~ 5 e を有する。ここで、ノード 5 a は、それぞれ CPU とメモリとを有し、各種アプリケーションプログラムを独立、または、連携して実行可能な装置であり、例えば、CPU とメモリとを搭載したサーバボードである。なお、ノード 5 b ~ 5 e も、ノード 5 a と同様の機能を発揮するものとして、以下の説明を省略する。

10

【0014】

一方、ディスクプール 4 は、複数の SSD 6 a ~ 6 f を有する。なお、図 1 に示すように、ディスクプール 4 は、SSD 6 a ~ 6 f 以外にも複数の SSD を有するものとする。各 SSD 6 a ~ 6 f は、ノード 5 a ~ 5 e が実行するアプリケーションプログラムが使用する半導体記憶装置であり、データの消去回数に基づく余命を有する。例えば、各 SSD 6 a ~ 6 f のうち、データの消去が行われていない半導体記憶装置については、余命の指標値である余命指標値が「100」となり、データの消去が何度も行われ、寿命となった SSD については、余命指標値が「0」となる。

【0015】

ディスクエリアネットワーク 3 は、各ノード 5 a ~ 5 e と、各ノード 5 a ~ 5 e がアプリケーションプログラムを実行する際に使用する SSD とを接続するネットワークである。例えば、ディスクエリアネットワーク 3 は、管理サーバ 10 からの指定に従って、ノード 5 a と、ディスクプール 4 が有する 1 つ以上の任意の SSD とを接続、もしくは接続の解除を行う。

20

【0016】

このような構成の元、ノード 5 a は、ディスクプール 4 が有する各 SSD のうち、任意の SSD を用いてアプリケーションプログラムを実行する。例えば、ノード 5 a は、SSD 6 a と SSD 6 b とを用いて、アプリケーションプログラムを実行する。

【0017】

管理サーバ 10 は、管理ネットワーク 7 を介して、各ノード 5 a ~ 5 e と接続されている。そして、管理サーバ 10 は、例えば、ノード 5 b が、新たなサービスを提供するため、アプリケーションプログラムを実行する場合は、SSD 6 d とノード 5 b とを接続するようにディスクエリアネットワーク 3 に指示する。この結果、ノード 5 b は、SSD 6 d を用いて、アプリケーションプログラムを実行し、各種サービスを提供することができる。

30

【0018】

このように、情報処理装置 1 は、複数のノード 5 a ~ 5 e とディスクプール 4 が有する各 SSD 6 a ~ 6 f とを用いて、任意の数のアプリケーションプログラムを実行し、各種サービスを提供することができる。例えば、情報処理装置 1 は、ノード 5 a、5 b を SSD 6 a、6 b と組み合わせて Web サービスを提供させ、ノード 5 c ~ 5 e と SSD 6 c、6 d とを組み合わせて大規模データ処理のサービスを提供させることができる。

40

【0019】

ここで、SSD 6 a は、NAND 型のフラッシュメモリや相変化メモリ (PCM: Phase Change Memory) 等、不揮発性の半導体メモリを有する。ここで、不揮発性の半導体メモリは、データの消去回数やデータの書き込み回数に制限があるデータセルを有する。例えば、1 つのデータセルが 1 ビットのデータを記憶する SLC (Single Level Cell) 方式では、データの消去回数が約 1 万回に制限される。また、1 つのデータセルが複数ビットのデータを記憶する MLC (Multi Level Cell) 方式では、データの消去回数が約 10 万回に制限される。

【0020】

このように、SSD 6 a が有するメモリセルには、データの消去回数に制限が存在する

50

ため、SSDには寿命が存在する。このため、情報処理装置1は、各SSD 6a~6fがデータを消去することができる残余回数に基づく余命指標値を収集する。例えば、情報処理装置1は、smartmontoolsを利用してMedia_Wearout_Indicatorの値を余命として収集する。そして、情報処理装置1は、各ノード5a~5eが利用するSSDの余命が所定の閾値を下回った場合は、利用中のSSDを他のSSDに変更する。すなわち、情報処理装置1は、ディスクプール4が有する各SSD 6a~6fについてのウェアレベリングを実行する。このため、情報処理装置1は、寿命を考慮せずにSSDをストレージとして利用することができる。

【0021】

例えば、図2は、実施例1に係る情報処理装置が実行する処理を説明する図である。例えば、図2中(A)に示すSSDは、ノード5a、5bがアプリケーションプログラムの実行に利用するSSDである。すなわち、情報処理装置1は、ノード5a、5bと図2中(A)に示すSSDとを組み合わせ、1つのシステムとして動作させている。また、図2中(B)に示すSSDは、利用されていないSSDである。

【0022】

ここで、図2中(C)に示すSSDの余命が所定の閾値を下回った場合は、情報処理装置1は、図2中(C)に示すSSDを図2中(D)に示すSSDと論理的な交換を行う。すなわち、情報処理装置1は、ディスクエリアネットワーク3の接続を変更し、ノード5a、5bと図2中(C)に示すSSDとの接続を切断し、ノード5a、5bと図2中(D)に示すSSDとを接続する。この結果、ノード5a、5bは、図2中(E)に示すSSDを用いて、アプリケーションプログラムを継続して実行する。

【0023】

また、例えば、SSD 6aの余命が短い際に、データの書き込み等、メモリセルにおけるデータの消去を伴うアクセスを行うのは、不具合が生じる可能性があるものの、データの読出し等、データの消去を伴わないアクセスについては、許容してもよい。そこで、情報処理装置1は、各ノード5a~5eが実行するアプリケーションプログラムが、単位時間あたりに読出したデータ量と書き込んだデータ量との比率を収集する。

【0024】

そして、情報処理装置1は、収集した比率に基づいて、アプリケーションプログラムの性質を判定する。その後、情報処理装置1は、判定したアプリケーションプログラムの性質に応じて、SSDを交換するかを判定するための閾値を決定し、SSD 6a~6fの余命が決定した閾値よりも短い場合は、SSDの交換を行う。

【0025】

例えば、情報処理装置1は、ノード5aが実行するアプリケーションプログラムがデータの書き込みよりもデータの読出しを頻繁に実行するリードインテンシブな性質である場合は、閾値の値「10」を決定する。そして、情報処理装置1は、ノード5aが使用するSSD 6aの余命が「10」よりも短い場合は、ノード5aが利用するSSD 6aを他のSSDと交換する。

【0026】

若しくは、情報処理装置1は、収集した比率においてデータの読み出しの比率が多い場合、直接、SSDを交換するかを判定するための閾値「10」を決定し、SSD 6a~6fの余命が決定した閾値よりも短い場合は、SSDの交換を行うようにしても良い。

【0027】

一方、情報処理装置1は、ノード5aが実行するアプリケーションプログラムがデータの読出しよりもデータの書き込みを頻繁に実行するライトインテンシブな性質である場合は、閾値の値「30」を決定する。すなわち、情報処理装置1は、アプリケーションプログラムの性質がライトインテンシブである場合は、アプリケーションプログラムの性質がリードインテンシブである場合の閾値よりも、大きな値を閾値とする。そして、情報処理装置1は、ノード5aが使用するSSD 6aの余命が「30」よりも短い場合は、ノード5aが利用するSSD 6aを他のSSDと交換する。

【 0 0 2 8 】

若しくは、情報処理装置 1 は、収集した比率においてデータの書き込みの比率が多い場合、直接、SSD を交換するかを判定するための閾値「30」を決定し、SSD 6 a ~ 6 f の余命が決定した閾値よりも短い場合は、SSD の交換を行うようにしても良い。

【 0 0 2 9 】

また、例えば、ノード 5 a が SSD 6 a のみを使用する際に、SSD 6 a の寿命が尽きてしまうと、SSD 6 a を他の SSD と交換する間、ノード 5 a が処理を実行できなくなる。また、ノード 5 a が使用する SSD 6 a を単に SSD 6 d と交換しただけでは、ノード 5 a が使用していたデータがなくなるので、ノード 5 a が継続して処理を実行できない。

10

【 0 0 3 0 】

そこで、情報処理装置 1 においては、各ノード 5 a ~ 5 e は、2 つの SSD を用いて、RAID (Redundant Arrays of Inexpensive Disks) 1 を構成し、データをミラーリングする。例えば、ノード 5 a は、RAID 1 によりデータがミラーリングされた SSD 6 a と SSD 6 b とを使用する。そして、情報処理装置 1 は、ノード 5 a が使用する SSD 6 a を SSD 6 d と交換する場合は、SSD 6 a とノード 5 a との接続を解除し、SSD 6 a のデータを SSD 6 d に複製する。そして、情報処理装置 1 は、データを複製した SSD 6 d とノード 5 a とを接続する。

【 0 0 3 1 】

また、ノード 5 a は、SSD 6 a との接続が解除されてから SSD 6 d が接続されるまでの間、以下の処理を実行する。すなわち、ノード 5 a は、データの書き込みを行う場合は、書き込み対象のデータをメモリ上にバッファリングし、データの読出しを行う場合は、接続されている SSD 6 b からデータの読出しを行う。そして、ノード 5 a は、SSD 6 d が接続されると、メモリ上にバッファリングしたデータを SSD 6 b、および SSD 6 d に反映させる。このため、情報処理装置 1 は、各ノード 5 a ~ 5 e にアプリケーションプログラムの実行を一時停止させずとも、SSD の交換を行うことができる。

20

【 0 0 3 2 】

次に、図 3 を用いて、管理サーバ 10 が有する機能構成について説明する。図 3 は、実施例 1 に係る管理サーバの機能構成を説明する図である。なお、図 3 に示す例では、情報処理装置 1 が有するノード 5 a ~ 5 e のうち、ノード 5 a、5 b を記載し、情報処理装置 1 が有する SSD のうち、SSD 6 a ~ 6 d を記載した。

30

【 0 0 3 3 】

図 3 に示すように、管理サーバ 10 は、記憶部 11、通信部 14、更新部 15、判定部 16、決定部 17、選択部 18、設定部 19 を有する。また、記憶部 11 は、余命管理テーブル 12、ワークロード管理テーブル 13 を記憶する。

【 0 0 3 4 】

以下、図 4、5 を用いて、記憶部 11 が記憶する余命管理テーブル 12、ワークロード管理テーブル 13 に格納された情報について説明する。まず、図 4 を用いて、余命管理テーブル 12 に格納された情報の一例を説明する。図 4 は、余命管理テーブルの一例を説明する図である。図 4 に示すように、余命管理テーブル 12 には、各 SSD 6 a ~ 6 f の余命が格納されている。詳細には、余命管理テーブル 12 には、SSD ID (Identifier)、接続先、余命が対応付けて格納されている。ここで、SSD ID とは、各 SSD 6 a ~ 6 f を識別するための識別子である。また、接続先とは、対応付けられた SSD ID が示す SSD が接続されているノードを示す情報である。また、余命とは、対応付けられた SSD ID の余命である。

40

【 0 0 3 5 】

例えば、図 4 に示す例では、余命管理テーブル 12 は、SSD ID が「SSD # 0」の SSD について、ノード 5 a が接続されており、余命が「100」であることを記憶する。また、余命管理テーブル 12 は、SSD ID が「SSD # 1」の SSD について、ノード 5 b が接続されており、余命が「80」であることを記憶する。また、余命管理テ

50

ーブル12は、SSD IDが「SSD#2」のSSDについては、どのノード5a~5eも接続されておらず、余命が「50」であることを記憶する。

【0036】

次に、図5を用いて、ワークロード管理テーブル13に格納された情報について説明する。図5は、ワークロード管理テーブルの一例を説明する図である。図5に示すように、ワークロード管理テーブル13には、各ノード5a~5eが実行するアプリケーションプログラムが、使用中のSSDから読みだしたデータと書き込んだデータとの比率が格納されている。

【0037】

詳細には、ワークロード管理テーブル13には、ノードIDとプログラム実行フラグと比率とが対応付けて格納されている。ここで、ノードIDとは、各ノード5a~5eを識別する識別子である。またプログラム実行フラグとは、対応付けられたノードIDが示すノードがプログラムを実行中か否かを示すフラグである。また、比率とは、対応付けられたノードIDが示すノードが、1秒間に読出したデータ量(Byte)と書き込んだデータ量(Byte)との比率を示す情報である。

【0038】

例えば、ワークロード管理テーブル13は、ノード5aがアプリケーションプログラムを実行しており、アプリケーションプログラムが読出したデータ量と書き込んだデータ量との比率が「50:50」である旨を記憶する。また、ワークロード管理テーブル13は、ノード5bがアプリケーションプログラムが読出したデータ量と書き込んだデータ量との比率が「90:10」である旨を記憶する。また、ワークロード管理テーブル13は、ノード5cがアプリケーションプログラムを実行していないので、比率が「null」である旨を記憶する。

【0039】

図3に戻り、通信部14は、管理ネットワーク7を介して、管理サーバ10と各ノード5a~5eとの通信、およびディスクエリアネットワーク3に対する指示を制御する。詳細には、通信部14は、更新部15および設定部19と各ノード5a~5eとの間の通信を制御する。また、通信部14は、設定部19からディスクエリアネットワーク3に対する指示を制御する。

【0040】

更新部15は、所定の時間間隔で、記憶部11が記憶する余命管理テーブル12とワークロード管理テーブル13とを更新する。例えば、各ノード5a~5eは、所定の時間間隔で、smartmontoolsを利用してMedia_Wearout_Indicatorの値を取得し、取得したMedia_Wearout_Indicatorの値を余命指標値とし、SSD IDと対応付けて管理サーバ10へ送信する。また、各ノード5a~5eは、読出したデータの量と書き込んだデータの量との比率をLinux(登録商標)のsarコマンド等、任意の手法で取得し、取得した比率を管理サーバ10へ送信する。また、各ノード5a~5eは、アプリケーションプログラムを起動した場合は、アプリケーションプログラムの起動を管理サーバ10に通知し、アプリケーションプログラムの実行を終了した場合は、アプリケーションプログラムの終了を管理サーバ10へ送信する。

【0041】

更新部15は、各ノード5a~5eからSSD ID、余命指標値、比率を受信する。そして、更新部15は、受信したSSD IDと余命指標値とに基づいて、余命管理テーブル12を更新し、受信した比率を用いて、ワークロード管理テーブル13を更新する。

【0042】

例えば、更新部15は、SSD ID「SSD#0」と余命「50」とをノード5aから受信すると、余命管理テーブル12からSSD ID「SSD#0」が格納されたエントリを抽出する。そして、更新部15は、抽出したエントリの接続先を「ノード5a」に更新し、余命を「50」に更新する。また、更新部15は、ノード5aから比率「40:60」を受信した場合には、ワークロード管理テーブル13からノードID「ノード5a

10

20

30

40

50

」が格納されたエントリを抽出する。そして、更新部 15 は、抽出したエントリの比率を「40:60」に更新する。

【0043】

また、更新部 15 は、ノード 5 a からアプリケーションの起動を通知された場合は、ワークロード管理テーブル 13 からノード ID「ノード 5 a」が格納されたエントリを抽出し、抽出したエントリのプログラム実行フラグを「1」に更新する。また、更新部 15 は、ノード 5 a からアプリケーションの終了を通知された場合は、ワークロード管理テーブル 13 からノード ID「ノード 5 a」が格納されたエントリを抽出し、抽出したエントリのプログラム実行フラグを「0」に更新する。

【0044】

判定部 16 は、ワークロード管理テーブル 13 に格納された読出しと書き込みの比率に応じて、各ノード 5 a ~ 5 e が実行するアプリケーションプログラムの性質を判定する。例えば、判定部 16 は、ワークロード管理テーブル 13 を参照し、各 SSD 6 a ~ 6 f に対する読出しと書き込みの比率を取得する。

【0045】

そして、判定部 16 は、読出したデータ量を算出し、算出したデータ量が所定の閾値よりも大きい場合は、アプリケーションプログラムがリードインテンシブであると判定する。また、判定部 16 は、算出したデータ量が所定の閾値よりも小さい場合は、アプリケーションプログラムがライトインテンシブであると判定する。なお、判定部 16 は、アプリケーションプログラムによる読出しの比率が書き込みの比率よりも多い場合は、アプリケーションプログラムがリードインテンシブであると判定してもよい。

【0046】

例えば、ワークロード管理テーブル 13 に図 5 に示す情報が格納されている場合には、判定部 16 は、以下の処理を実行する。例えば、判定部 16 は、ノード 5 b が実行するアプリケーションプログラムの読出し比率が書き込み比率よりも大きいので、ノード 5 b が実行するアプリケーションプログラムがリードインテンシブであると判定する。また、判定部 16 は、ノード 5 c が実行するアプリケーションプログラムの書き込み比率が読出し比率よりも大きいので、ノード 5 c が実行するアプリケーションプログラムがライトインテンシブであると判定する。

【0047】

ここで、ノード 5 a が実行するアプリケーションプログラムについては、読出しと書き込みの比率が同じである。このような場合には、判定部 16 は、システムの安定性を考慮し、ノード 5 a が実行するアプリケーションプログラムがライトインテンシブであると判定してもよい。なお、判定部 16 は、アプリケーションプログラムの性質を判定した場合は、判定結果を決定部 17 に通知する。

【0048】

なお、判定部 16 は、アプリケーションプログラムの性質を判定せず、各 SSD 6 a ~ 6 f に対する読出しと書き込みのどちらかが多いかを判定し、判定結果を決定部 17 に通知してもよい。

【0049】

決定部 17 は、各ノード 5 a ~ 5 e が実行するアプリケーションプログラムの性質に基づいて、SSD の交換基準となる閾値を決定する。例えば、決定部 17 は、ノード 5 a が実行するアプリケーションプログラムがライトインテンシブである旨の通知を判定部 16 から受信する。このような場合は、決定部 17 は、ノード 5 a が使用する SSD の交換基準となる閾値を「30」とし、ノード 5 a が使用する SSD の交換基準となる閾値が「30」である旨を選択部 18 に通知する。

【0050】

一方、決定部 17 は、ノード 5 c が実行するアプリケーションプログラムがライトインテンシブである旨の通知を判定部 16 から受信する。このような場合は、決定部 17 は、ノード 5 c が使用する SSD の交換基準となる閾値を「10」とし、ノード 5 c が使用する

10

20

30

40

50

るSSDの交換基準となる閾値が「10」である旨を選択部18に通知する。すなわち、決定部17は、各ノード5a～5eが実行するアプリケーションプログラムがリードインテンシブである場合は、アプリケーションプログラムがライトインテンシブである際よりも小さい値の閾値を決定する。

【0051】

なお、上述した例では、決定部17は、各ノード5a～5eが実行するアプリケーションプログラムがリードインテンシブである場合は、閾値を「10」とし、ライトインテンシブである場合は、閾値を「30」とした。しかし、実施例はこれに限定されるものではない。例えば、決定部17は、ワークロード管理テーブル13に格納された比率に応じて、閾値を決定してもよい。例えば、決定部17は、ライトインテンシブなアプリケーションプログラムについて、ワークロード管理テーブル13に格納された書き込みの比率が所定の値よりも大きい場合は、閾値を「30」よりも大きい値にしてもよい。

10

【0052】

なお、決定部17は、各SSD6a～6fに対する読出しと書き込みのどちらが多いかを示す判定結果を受信した場合は、受信した判定結果に基づいて、SSD6a～6fの交換基準となる閾値を決定しても良い。例えば、決定部17は、SSD6aに対して読出しが多い旨の判定結果を受信した場合は、SSD6aの交換基準となる閾値を「10」に決定する。また、決定部17は、SSD6aに対して書き込みが多い旨の判定結果を受信した場合は、SSD6aの交換基準となる閾値を「30」に決定する。

【0053】

20

選択部18は、アプリケーションプログラムの性質に応じて決定した閾値に応じて、各ノード5a～5eが使用するSSDを交換するか否かを判定する。例えば、選択部18は、決定部17からノード5aが使用するSSDの交換基準となる閾値が「30」である旨を受信する。すると、選択部18は、余命管理テーブル12を参照し、ノード5aが使用する1つ以上のSSDを識別する。そして、選択部18は、識別したSSDの余命指標値が閾値「30」よりも小さいか否かを判定し、余命指標値が閾値「30」よりも小さい場合は、SSDを交換すると判定する。

【0054】

また、選択部18は、SSDを交換すると判定した場合は、以下の処理を実行する。例えば、選択部18は、余命管理テーブル12を参照し、接続先が未接続のSSDを抽出する。そして、選択部18は、ノード5aが実行するアプリケーションがリードインテンシブである場合は、抽出したSSDのうち、余命指標値が「10」より大きく「30」以下であるSSDを検索する。その後、選択部18は、余命指標値が「10」より大きく「30」以下であるSSDを検出した場合は、ノード5aが使用するSSDのうち交換対象となるSSDのSSD IDと余命管理テーブル12から検出したSSDのSSD IDとを設定部19に通知する。

30

【0055】

また、選択部18は、余命指標値が「10」より大きく「30」以下であるSSDを検出できなかった場合は、抽出したSSDのうち、余命指標値が「30」より大きなSSDを検索する。そして、選択部18は、余命指標値が「30」より大きなSSDを検出した場合は、ノード5aが使用するSSDのうち交換対象となるSSDのSSD IDと余命管理テーブル12から検出したSSDのSSD IDとを設定部19に通知する。なお、選択部18は、余命指標値が「30」より大きなSSDを検出できなかった場合は、交換可能なSSDがない旨を管理者等に通知する。

40

【0056】

一方、選択部18は、ノード5aが実行するアプリケーションがライトインテンシブである場合は、余命指標値が「30」より大きなSSDを検索する。そして、選択部18は、余命指標値が「30」より大きなSSDを検出した場合は、ノード5aが使用するSSDのうち交換対象となるSSDのSSD IDと余命管理テーブル12から検出したSSDのSSD IDとを設定部19に通知する。

50

【 0 0 5 7 】

なお、選択部 18 は、判定部 16 がアプリケーションの性質を判定しない場合であっても、同様の処理を行うことで、交換する S S D を選択できる。例えば、選択部 18 は、S S D 6 a について閾値が「10」であり、S S D 6 a の余命指標値が「10」よりも小さい場合は、S S D 6 a と交換する S S D として、余命指標値が「10」より大きく「30」以下となる S S D を選択する。また、選択部 18 は、S S D 6 a について閾値が「30」であり、S S D 6 a の余命指標値が「30」より小さい場合は、S S D 6 a と交換する S S D として、余命指標値が「30」より大きい S S D を選択する。

【 0 0 5 8 】

設定部 19 は、各ノード 5 a ~ 5 e が使用する S S D の設定変更を行う。例えば、設定部 19 は、選択部 18 から、S S D 6 a の S S D I D と S S D 6 d の S S D I D とを受信する。すると、設定部 19 は、ワークロード管理テーブル 13 を参照し、S S D 6 a を使用中のノード 5 a を識別する。そして、設定部 19 は、識別したノード 5 a に対して S S D 6 a の接続解除手続きを依頼する。このような場合は、ノード 5 a は、S S D 6 a の接続解除手続きを実行し、S S D 6 a の接続を解除する。

【 0 0 5 9 】

また、設定部 19 は、ノード 5 a と S S D 6 a との接続を解除するようディスクエリアネットワーク 3 に指示する。この結果、ノード 5 a と S S D 6 a との接続は解除される。また、設定部 19 は、ワークロード管理テーブル 13 を参照し、アプリケーションプログラムを実行していないノードを識別する。例えば、設定部 19 は、ワークロード管理テーブル 13 に図 5 に示す情報が格納されている場合は、ノード 5 c を識別する。そして、設定部 19 は、ディスクエリアネットワーク 3 にノード 5 c と S S D 6 a 、6 d とを接続するよう指示する。そして、設定部 19 は、ノード 5 c に S S D 6 a のデータを S S D 6 d に複製するよう依頼する。

【 0 0 6 0 】

すると、ノード 5 c は、S S D 6 a のデータを S S D 6 d に複製し、複製完了を設定部 19 に通知する。すると、設定部 19 は、ノード 5 c と S S D 6 d との接続を解除し、ノード 5 a と S S D 6 d とを接続するようにディスクエリアネットワーク 3 に指示する。そして、設定部 19 は、ノード 5 a に S S D 6 d の接続手続きを依頼する。この結果、ノード 5 a は、S S D 6 a と同じデータを記憶する S S D 6 d を用いて、アプリケーションプログラムの実行を継続する。

【 0 0 6 1 】

なお、ノード 5 a は、S S D 6 a との接続が解除されてから S S D 6 d が接続されるまでの間、S S D 6 b からデータの読出しを行う。また、ノード 5 a は、データの書き込みを行う場合は、メモリ上に書き込み対象のデータをバッファし、S S D 6 d が接続された後に、バッファしたデータを S S D 6 b 、6 d に反映させる。なお、ノード 5 a は、バッファオーバーフローが発生した場合は、S S D 6 b に対してデータの書き込みを行い、S S D 6 d が接続された後に、S S D 6 b から S S D 6 d へデータを複製してもよい。

【 0 0 6 2 】

次に、図 6 を用いて、実施例 1 に係る情報処理装置 1 が実行する処理の流れについて説明する。図 6 は、実施例 1 に係る情報処理装置 1 が実行する処理の流れを説明するフローチャートである。なお、以下の説明では、ノード 5 a が S S D 6 a 、6 b を使用しており、S S D 6 a を S S D 6 d に交換する処理の流れについて説明する。また、以下の説明では、ノード 5 c がアプリケーションプログラムを実行していないものとする。

【 0 0 6 3 】

例えば、ノード 5 a は、S S D へのアクセスを監視し、読出しと書き込みの比率である Read / Write 比率を管理サーバ 10 に通知する (ステップ S 101)。また、ノード 5 a は、使用中の S S D 6 a 、6 b の余命を確認し、余命指標値を管理サーバ 10 へ送信する (ステップ S 102)。すると、管理サーバ 10 は、受信した Read / Write 比率、および余命指標値に従って、ワークロード管理テーブル 13 と余命管理テーブ

10

20

30

40

50

ル 1 2 とを更新する (ステップ S 1 0 3)。

【 0 0 6 4 】

また、管理サーバ 1 0 は、ワークロード管理テーブル 1 3 に格納された R e a d / W r i t e 比率に基づいて、ノード 5 a が実行するアプリケーションプログラムがリードインテンシブであるかライトインテンシブであるかを判定する (ステップ S 1 0 4)。そして、管理サーバ 1 0 は、ノード 5 a に接続されている S S D 6 a、6 b の余命が、ノード 5 a が実行するアプリケーションプログラムの性質に応じた閾値を超えていないか確認する (ステップ S 1 0 5)。

【 0 0 6 5 】

また、管理サーバ 1 0 は、アプリケーションプログラムを実行中の全てのノードについて、アプリケーションプログラムの性質に応じた閾値と S S D の余命とのチェックを行ったか判定する (ステップ S 1 0 6)。そして、管理サーバ 1 0 は、アプリケーションプログラムを実行中の全てのノードについて、アプリケーションプログラムの性質に応じた閾値と S S D の余命とのチェックを行っていない場合は (ステップ S 1 0 6 N o)、以下の処理を実行する。すなわち、管理サーバ 1 0 は、アプリケーションプログラムを実行中の各ノードについて、ステップ S 1 0 5 を実行する。

10

【 0 0 6 6 】

また、管理サーバ 1 0 は、アプリケーションプログラムを実行中の全てのノードについて、アプリケーションプログラムの性質に応じた閾値と S S D の余命とのチェックを行った場合は (ステップ S 1 0 6 Y e s)、以下の処理を実行する。すなわち、管理サーバ 1 0 は、余命指標値が閾値を超えた S S D、すなわち余命指標値が閾値よりも小さい S S D があるか否かを判定する (ステップ S 1 0 7)。また、管理サーバ 1 0 は、余命指標値が閾値よりも小さい S S D がある場合は (ステップ S 1 0 7 Y e s)、使用されていない S S D にアプリケーションプログラムの性質に応じた閾値を満たす S S D が存在するか否かを判定する (ステップ S 1 0 8)。

20

【 0 0 6 7 】

また、管理サーバ 1 0 は、使用されていない S S D にアプリケーションプログラムの性質に応じた閾値を満たす S S D、例えば S S D 6 d が存在する場合は (ステップ S 1 0 8 Y e s)、以下の処理を実行する。すなわち、管理サーバ 1 0 は、余命が閾値を超えた S S D を使用するノード、例えばノード 5 a に対して、余命が閾値を超えた S S D、例えば S S D 6 a の接続解除手続きを依頼する (ステップ S 1 0 9)。すると、ノード 5 a は、S S D 6 a の接続解除手続きを実行し、手続き完了を管理サーバ 1 0 に通知する (ステップ S 1 1 0)。

30

【 0 0 6 8 】

また、管理サーバ 1 0 は、ノード 5 a から S S D 6 a の接続解除手続きの完了通知を受信すると、以下の S S D 交換処理を実行する (ステップ S 1 1 1)。すなわち、管理サーバ 1 0 は、ディスクエリアネットワーク 3 に指示し、ノード 5 a から S S D 6 a の接続を解除し、ノード 5 c と S S D 6 a、6 d とを接続し、接続完了をノード 5 b に通知する。すると、ノード 5 c は、S S D 6 a から S S D 6 d へデータのコピーを実行し、実行完了を管理サーバ 1 0 に通知する (ステップ S 1 1 2)。

40

【 0 0 6 9 】

すると、管理サーバ 1 0 は、ノード 5 c から S S D 6 d の接続を解除し、S S D 6 d をノード 5 a に接続してノード 5 a に S S D 6 d を接続した旨を通知する (ステップ S 1 1 3)。すると、ノード 5 a は、S S D 6 d の接続手続きを実行する (ステップ S 1 1 4)。また、管理サーバ 1 0 は、ノード 5 c から S S D 6 a の接続を解除する (ステップ S 1 1 5)。

【 0 0 7 0 】

また、管理サーバ 1 0 は、閾値を超えた S S D で未交換の S S D が存在するか判定し (ステップ S 1 1 6)、閾値を超えた S S D で未交換の S S D が存在しない場合は (ステップ S 1 1 6 N o)、処理を終了する。一方、管理サーバ 1 0 は、閾値を超えた S S D で未

50

交換のSSDが存在する場合は(ステップS116Yes)、ステップS108を実行する。なお、管理サーバ10は、使用していないSSDでアプリケーションプログラムの性質に応じた閾値を満たすSSDが存在しない場合は(ステップS108No)、そのまま処理を終了する。

【0071】

次に、図7、図8を用いて、図6中ステップS108にて使用されていないSSDからアプリケーションプログラムの性質に応じた閾値を満たすSSDがあるか否かを判定する処理の流れについて説明する。まず、図7を用いて、アプリケーションプログラムの性質がリードインテンシブである際に、選択部18が実行する処理の流れについて説明する。

【0072】

図7は、選択部が実行する処理の流れを説明するための第1のフローチャートである。例えば、選択部18は、使用されていないSSDに、余命指標値が「10」よりも大きく、「30」以下のSSDが存在するか否かを判定する(ステップS201)。ここで、選択部18は、余命指標値が「10」よりも大きく、「30」以下のSSDが存在しない場合は(ステップS201No)、余命指標値が「30」より大きく「100」以下のSSDが存在するか否かを判定する(ステップS202)。

【0073】

そして、選択部18は、余命指標値が「30」より大きく「100」以下のSSDが存在する場合は(ステップS202Yes)、閾値を満たすSSDを選択し(ステップS203)、処理を終了する。また、選択部18は、使用されていないSSDに、余命指標値が「10」よりも大きく、「30」以下のSSDが存在する場合は(ステップS201Yes)、閾値を満たすSSDを選択し(ステップS203)、処理を終了する。また、選択部18は、余命指標値が「30」より大きく「100」以下のSSDが存在しない場合は(ステップS202No)、SSDの選択を行わずに、処理を終了する。

【0074】

次に、図8を用いて、アプリケーションプログラムの性質がライトインテンシブである際に、選択部18が実行する処理の流れについて説明する。図8は、選択部が実行する処理の流れを説明するための第2のフローチャートである。例えば、選択部18は、アプリケーションプログラムの性質がライトインテンシブである場合は、使用されていないSSDに余命指標値が「30」より大きく「100」以下のSSDが存在するか判定する(ステップS301)。

【0075】

そして、選択部18は、使用されていないSSDに余命指標値が「30」より大きく「100」以下のSSDが存在する場合は(ステップS301Yes)、閾値を満たすSSDを選択し(ステップS302)、処理を終了する。一方、選択部18は、使用されていないSSDに余命指標値が「30」より大きく「100」以下のSSDが存在しない場合は(ステップS301No)、SSDを選択せずに処理を終了する。

【0076】

[情報処理装置1の効果]

上述したように、情報処理装置1は、SSD6aと、例えばSSD6aの予備であるSSD6fとを有する。また、情報処理装置1は、SSD6aについて、余命指標値と、SSD6aに対する書込みと読み出しとの読み書きの比率を収集する。そして、情報処理装置1は、収集した比率に基づいて、SSD6aの交換基準となる閾値を決定する。その後、情報処理装置1は、収集した余命指標値が決定した閾値よりも短い場合は、SSD6aをSSD6fと交換する。このため、情報処理装置1は、ノード5aにSSD6a~6fをストレージ装置として継続して利用させることができる。

【0077】

また、情報処理装置1は、SSD6aに対する読み書きの比率に応じた閾値を用いて、SSD6aの交換を判定するので、各SSD6a~6fの余命を使いこなすことができる。この結果、情報処理装置1は、各SSD6a~6fをストレージとして利用した際の寿命

10

20

30

40

50

を最大限利用することができる。

【0078】

なお、情報処理装置1は、収集した比率に応じて、ノード5aが実行するアプリケーションプログラムの性質を判定し、判定したアプリケーションプログラムの性質に応じて、ノード5aが使用するSSD6aの交換基準となる閾値を決定してもよい。その後、情報処理装置1は、収集した余命指標値が、決定した閾値よりも短い場合は、SSD6aを他のSSDと交換する。このような場合にも、情報処理装置1は、ノード5aにSSD6a～6fをストレージ装置として継続して利用させることができる。

【0079】

また、情報処理装置1は、ノード5aが実行するアプリケーションプログラムの性質に応じた閾値を用いて、SSDの交換を行ってもよい。このような場合にも、情報処理装置1は、各SSD6a～6fの余命を使いきり、各SSD6a～6fをストレージとして利用した際の寿命を最大限利用することができる。

10

【0080】

また、情報処理装置1は、収集した比率においてデータの書き込みの比率が多い場合は、データの読み出しの比率が多い場合に設定する第1の閾値よりも値が大きい第2の閾値をSSD6aの交換基準となる閾値とする。このため、情報処理装置1は、書き込みの頻度が多いアプリケーションプログラムを実行するノードに対して余命が長いSSDを割り当て、読み出しの頻度が多いアプリケーションプログラムを実行するノードに対して余命が短いSSDを割り当てる。この結果、情報処理装置1は、各SSD6a～6fの寿命を最大限利用することができる。

20

【0081】

なお、情報処理装置1は、ノード5aが実行するアプリケーションプログラムの性質がリードインテンシブであるかライトインテンシブであるかを判定してもよい。そして、情報処理装置1は、アプリケーションプログラムがライトインテンシブである場合は、リードインテンシブである際に用いる閾値よりも値が大きい閾値を用いることとしてもよい。

【0082】

また、情報処理装置1は、使用されていないSSD、すなわち予備のSSDを複数有し、SSD6aを交換する場合は、使用されていないSSDのうち、以下の条件を満たすSSDを選択する。すなわち、情報処理装置1は、SSD6aから収集した比率についてデータ読み出しの比率が多い場合、予備のSSDのうち、余命指標値が第1の閾値よりも大きく、かつ、余命指標値が第2の閾値よりも小さいSSDを選択する。そして、情報処理装置1は、SSD6aを選択したSSDと交換する。このため、情報処理装置1は、使用されていないSSDのうち、余命が短いSSDを優先して使用することができる。

30

【0083】

なお、情報処理装置1は、ノード5aが実行するアプリケーションプログラムの性質がリードインテンシブである場合は、使用されていないSSDのうち、以下の条件を満たすSSDを選択してもよい。すなわち、情報処理装置1は、余命指標値が、アプリケーションプログラムの性質がリードインテンシブである際の閾値よりも大きく、かつ、アプリケーションプログラムの性質がライトインテンシブである際の閾値よりも短いSSDを選択してもよい。

40

【0084】

また、情報処理装置1は、SSD6aから収集した比率についてデータ書き込みの比率が多い場合、予備のSSDのうち、余命指標値が第2の閾値よりも大きいSSDを選択する。そして、情報処理装置1は、SSD6aを選択したSSDと交換する。このため、情報処理装置1は、書き込みの頻度が多いアプリケーションプログラムを実行するノードに対し、余命に余裕があるSSDを割り当てることができる。この結果、情報処理装置1は、SSD6a～6fの余命が尽き、各ノード5a～5eが使用するSSDが急に使用できなくなるといった事態を防ぐことができるので、アプリケーションプログラムを安定して実行することができる。

50

【 0 0 8 5 】

なお、情報処理装置 1 は、ノード 5 a が実行するアプリケーションプログラムの性質がライトインテンシブである場合は、使用されていない S S D のうち、以下の条件を満たす S S D を選択してもよい。すなわち、情報処理装置 1 は、余命指標値が、アプリケーションプログラムの性質がライトインテンシブである際の閾値よりも大きい S S D を選択する。そして、情報処理装置 1 は、ノード 5 a が使用する S S D を選択した S S D と交換してもよい。

【 0 0 8 6 】

また、情報処理装置 1 は、単位時間あたりに S S D 6 a から読み出されたデータ量が所定の閾値よりも多い場合は、データの読み出しの比率が多いと判定する。このため、情報処理装置 1 は、多くのデータを読み出すアプリケーションプログラムに対し、余命指標値が小さい S S D を割り当てることができる。

10

【 0 0 8 7 】

なお、情報処理装置 1 は、ノード 5 a が実行するアプリケーションプログラムが読み出したデータ量が所定の閾値よりも多い場合は、ノード 5 a が実行するアプリケーションプログラムがリードインテンシブであると判定する。そして、情報処理装置 1 は、リードインテンシブなアプリケーションプログラムに余命指標値が小さい S S D を割り当ててもよい。

【 0 0 8 8 】

また、情報処理装置 1 は、例えば、ノード 5 a が S S D 6 a と S S D 6 b とを R A I D 1 によりミラーリングした際に、S S D 6 a を S S D 6 d と交換する場合は、以下の処理を実行する。すなわち、情報処理装置 1 は、S S D 6 a を S S D 6 d とを交換するまでの間、ノード 5 a から S S D 6 a、6 b に書き込もうとしたデータをメモリ上に保持する。そして、情報処理装置 1 は、S S D 6 a のデータを S S D 6 d に複製し、メモリ上に保持したデータを S S D 6 b および 6 d に複製する。このため、情報処理装置 1 は、ノード 5 a が実行するアプリケーションプログラムを動作させたままで S S D の交換を行うことができる。

20

【 0 0 8 9 】

このように、情報処理装置 1 は、各 S S D 6 a ~ 6 f の余命が尽きるまで、各 S S D 6 a ~ 6 f を交換するか否かを判定する。ここで、情報処理装置 1 には、複数のノード 5 a ~ 5 e が設置され、各ノード 5 a ~ 5 e が実行するアプリケーションには、それぞれ異なる性質を有する。このため、各ノード 5 a ~ 5 e が R A I D 1 によるミラーリングを構築した場合にも、各 S S D 6 a ~ 6 f の余命に偏りが発生することとなる。この結果、情報処理装置 1 は、ノード 5 a が R A I D 1 によりミラーリングされた S S D を用いる場合であっても、各 S S D 6 a ~ 6 f を長期に渡り効率よく利用することができる。

30

【実施例 2】

【 0 0 9 0 】

これまで本発明の実施例について説明したが実施例は、上述した実施例以外にも様々な異なる形態にて実施されてよいものである。そこで、以下では実施例 2 として本発明に含まれる他の実施例を説明する。

【 0 0 9 1 】

(1) 他の R A I D について

上述した情報処理装置 1 では、ノード 5 a が R A I D 1 により、S S D 6 a、6 b のデータをミラーリングしていた。しかし、実施例はこれに限定されるものではない。例えば、ノード 5 a は、5 台の S S D に対し、データとデータを復元するためのパリティとを分散して格納する R A I D 5 を構築してもよい。また、ノード 5 a は、6 台の S S D に対し、データとデータを復元するためのパリティとを分散して格納する R A I D 6 を構築してもよい。

40

【 0 0 9 2 】

このように、ノード 5 a が R A I D 5、または R A I D 6 を構築中には、情報処理装置 1 は、交換対象となる S S D が記憶するデータの複製を行わなくともよい。例えば、情報

50

処理装置 1 は、ノード 5 a が S S D 6 a ~ 6 e を用いて R A I D 5 を構築中に、S S D 6 a を S S D 6 f と交換する場合は、S S D 6 b ~ 6 e から S S D 6 f が記憶すべきデータを復元することができる。

【 0 0 9 3 】

このため、情報処理装置 1 は、ノード 5 a が使用中の S S D 6 a を S S D 6 f と交換する場合は、S S D 6 a とノード 5 a との接続を解除する。そして、情報処理装置 1 は、S S D 6 a から S S D 6 f にデータを複製せずに S S D 6 f とノード 5 b との接続を行う。この結果、ノード 5 a は、R A I D 5 の機能により、S S D 6 b ~ 6 e から S S D 6 f が記憶すべきデータを復元する。また、R A I D 5 の機能により、ノード 5 a は、S S D 6 a と S S D 6 f とを交換する間、データの読出しを継続することができ、データの書き込みもバッファする必要はない。このため、ノード 5 a は、S S D 6 a と S S D 6 f とを交換する間、アプリケーションプログラムを継続して実行することができる。

10

【 0 0 9 4 】

次に、図 9 を用いて、R A I D 5 を構築中のノード 5 a に対して情報処理装置 1 が実行する処理の流れについて説明する。図 9 は、実施例 2 に係る情報処理装置が実行する処理の流れを説明するフローチャートである。なお、図 9 に示す処理のうちステップ S 4 0 1 ~ S 4 0 8 については、図 6 に示すステップ S 1 0 1 ~ S 1 0 8 と同様の処理であるものとして、説明を省略する。また、以下の説明では、ノード 5 a が S S D 6 a ~ 6 e を用いて R A I D 5 を構築しており、ノード 5 a が実行するアプリケーションプログラムの性質に応じた閾値を満たす S S D 6 f と S S D 6 a とを交換する処理の流れについて説明する。

20

【 0 0 9 5 】

例えば、管理サーバ 1 0 は、ノード 5 a から S S D 6 a の接続を解除し、S S D 6 f をノード 5 a に接続し、余命管理テーブル 1 2 を更新してノード 5 a に S S D 6 f を接続した旨を通知する（ステップ S 4 0 9 ）。すると、ノード 5 a は、S S D 6 b ~ 6 e を用いて S S D 6 f が記憶すべきデータを復元し、S S D 6 b ~ 6 f を用いて R A I D 5 を再構築する。そして、ノード 5 a は、R A I D 5 の再構築が終了した旨を管理サーバ 1 0 に通知する（ステップ S 4 1 0 ）。また、管理サーバ 1 0 は、閾値を超えた S S D で未交換の S S D が存在するか判定し（ステップ S 4 1 1 ）、閾値を超えた S S D で未交換の S S D が存在しない場合は（ステップ S 4 1 1 N o ）、処理を終了する。一方、管理サーバ 1 0 は、閾値を超えた S S D で未交換の S S D が存在する場合は（ステップ S 4 1 1 Y e s ）、ステップ S 4 0 8 を実行する。

30

【 0 0 9 6 】

このように、情報処理装置 1 は、ノード 5 a が使用する S S D 6 a ~ 6 e にデータとデータを復元するためのパリティとが格納されている場合に、S S D 6 a ~ 6 e のいずれかを交換する時は、以下の処理を実行する。例えば、情報処理装置 1 は、S S D 6 a と S S D 6 f とを交換し、S S D 6 b ~ 6 e が記憶するデータとパリティとを用いて、S S D 6 f が記憶するデータを復元する。このため、情報処理装置 1 は、交換対象となる S S D のデータを複製せずとも、ノード 5 a がアプリケーションプログラムを実行したままで、S S D の交換を行うことができる。

40

【 0 0 9 7 】

（ 2 ）情報処理装置 1 の構成について

上述した実施例 1 では、ディスクエリアネットワーク 3 を介して、C P U プール 2 が有する各ノード 5 a ~ 5 e に対し、ディスクプール 4 が有する各 S S D 6 a ~ 6 e のうち、任意の S S D を接続する例について説明した。しかし、実施例はこれに限定されるものではない。例えば、情報処理装置 1 は、ディスクエリアネットワーク 3 の代わりに、S A N（Storage Area Netowrk）を用いて、各ノード 5 a ~ 5 e と各 S S D 6 a ~ 6 e とを接続してもよい。また、情報処理装置 1 は、S S D だけではなく、H D D（Hard Disk Drive）等任意の装置をストレージとしてディスクプール 4 に含めてもよい。

【 0 0 9 8 】

50

(3) アプリケーションプログラムの性質について

上述した実施例 1 では、管理サーバ 10 は、アプリケーションプログラムによる読出しと書き込みの比率に基づいて、アプリケーションプログラムの性質を判定した。しかし、実施例はこれに限定されるものではない。例えば、管理サーバ 10 は、アプリケーションプログラムによる読出しと書き込みの比率だけではなく、例えば、平均的な書き込みデータ量を考慮に加えてアプリケーションプログラムの性質を判定してもよい。

【 0 0 9 9 】

また、管理サーバ 10 は、アプリケーションプログラムを実行するノードの数を考慮に加えてもよい。また、管理サーバ 10 は、必ずしも、アプリケーションプログラムの性質を判定する必要はない。すなわち、管理サーバ 10 は、各 SSD 6 a ~ 6 e に対する読出しと書き込みの比率に応じて、交換の指標となる閾値を直接設定し、設定した閾値と余命指標値とに応じて、各 SSD 6 a ~ 6 e の交換を行っても良い。

10

【 0 1 0 0 】

(4) RAID について

上述した実施例 1 では、ノード 5 a が RAID 1 を構築する例について説明し、実施例 2 では、ノード 5 a が RAID 6 を構築する例について説明した。しかしながら、実施例はこれに限定されるものではない。例えば、情報処理装置 1 は、ノード 5 a が RAID 1 を構築し、ノード 5 b が RAID 5、または RAID 6 を構築中には、ノード 5 a に対して図 6 に示した処理を実行し、ノード 5 b に対して図 9 に示した処理を実行してもよい。

【 0 1 0 1 】

20

(5) 収集する情報について

上述した実施例 1 では、管理サーバ 10 は、単位時間あたりに読出したデータ量と書き込んだデータ量との比率を収集した。しかし、実施例はこれに限定されるものではなく、例えば、管理サーバ 10 は、単位時間あたりにデータを読出した回数とデータを書き込んだ回数とを各ノード 5 a ~ 5 e から収集してもよい。また、管理サーバ 10 は、データを読出した回数が所定の閾値よりも多い場合は、リードインテンシブであると判断しても良い。すなわち、管理サーバ 10 は、データを読出した回数とデータを書き込んだ回数との比率に応じて、アプリケーションプログラムがリードインテンシブであるかライトインテンシブであるかを判断して良い。

【 0 1 0 2 】

30

このように、情報処理装置 1 は、ノード 5 a が実行するアプリケーションプログラムが SSD 6 a からデータを読出した回数が所定の閾値よりも多い場合は、データの書き込みの比率が多いと判定する。例えば、情報処理装置 1 は、ノード 5 a が実行するアプリケーションプログラムがリードインテンシブであると判定する。このため、情報処理装置 1 は、データを何度も読出すアプリケーションプログラムに対し、余命指標値が小さい SSD を割り当てることができる。

【 0 1 0 3 】

(6) プログラム

ところで、実施例 1 に係る管理サーバ 10 は、ハードウェアを利用して各種の処理を実現する場合を説明した。しかし、実施例はこれに限定されるものではなく、あらかじめ用意されたプログラムをコンピュータで実行することによって実現するようにしてもよい。そこで、以下では、図 10 を用いて、実施例 1 に示した管理サーバ 10 と同様の機能を有するプログラムを実行するコンピュータの一例を説明する。図 10 は、制御プログラムを実行するコンピュータを説明する図である。

40

【 0 1 0 4 】

図 10 に例示されたコンピュータ 100 は、ROM (Read Only Memory) 110、HDD (Hard Disk Drive) 120、RAM (Random Access Memory) 130、CPU 140 がバス 160 で接続される。また、図 10 に例示されたコンピュータ 100 は、パケットを送受信するための I/O (Input Output) 150 を有する。

【 0 1 0 5 】

50

HDD 120は、図3に示した記憶部11が記憶する余命管理テーブル12と同様の情報である余命管理テーブル121、および、ワークロード管理テーブル13と同様の情報であるワークロード管理テーブル122を記憶する。また、RAM 130には、制御プログラム131があらかじめ保持される。CPU 140が制御プログラム131をRAM 130から読出して実行することによって、図10に示す例では、制御プログラム131は、制御プロセス141として機能するようになる。なお、制御プロセス141は、図3に示した更新部15、判定部16、決定部17、選択部18、設定部19と同様の機能を発揮する。

【0106】

なお、本実施例で説明した制御プログラムは、あらかじめ用意されたプログラムをパーソナルコンピュータやワークステーションなどのコンピュータで実行することによって実現することができる。このプログラムは、インターネットなどのネットワークを介して配布することができる。また、このプログラムは、ハードディスク、フレキシブルディスク(FD)、CD-ROM(Compact Disc Read Only Memory)、MO(Magneto Optical Disc)、DVD(Digital Versatile Disc)などのコンピュータで読取可能な記録媒体に記録される。また、このプログラムは、コンピュータによって記録媒体から読出されることによって実行することもできる。

【0107】

以上の各実施例を含む実施形態に関し、さらに以下の付記を開示する。

【0108】

(付記1) データが読み書きされる半導体記憶装置と、前記半導体記憶装置の予備である予備半導体記憶装置とを有する情報処理装置において、

前記半導体記憶装置について、書き込まれたデータを消去できる残余回数に基づいた余命指標値、および、前記半導体記憶装置に対する読出しと書き込みとの読み書きの比率を収集する収集部と、

前記収集部が収集した前記比率に基づいて、前記半導体記憶装置の交換基準となる閾値を決定する決定部と、

前記収集部が収集した前記余命指標値が、前記決定部が決定した閾値よりも短い場合は、前記半導体記憶装置を前記予備半導体記憶装置と交換する交換部と

を有することを特徴とする情報処理装置。

【0109】

(付記2) 前記決定部は、

収集した前記比率においてデータの読出しの比率が多い場合、第1の閾値を前記半導体記憶装置の交換基準となる閾値とし、データの書き込みの比率が多い場合、前記第1の閾値よりも値が大きい第2の閾値を前記半導体記憶装置の交換基準となる閾値とすることを特徴とする付記1に記載の情報処理装置。

【0110】

(付記3) 前記情報処理装置はさらに、

前記予備半導体記憶装置を複数有し、

前記決定部は、

収集した前記比率においてデータの読出しの比率が多い場合、前記複数の予備半導体記憶装置のうち、余命指標値が前記第1の閾値よりも大きく、かつ、前記第2の閾値よりも小さい予備半導体記憶装置を選択し、

前記交換部は、

前記半導体記憶装置を前記選択部が選択した予備半導体記憶装置と交換することを特徴とする付記2に記載の情報処理装置。

【0111】

(付記4) 前記選択部は、

収集した前記比率においてデータの書き込みの比率が多い場合、前記複数の予備半導体記憶装置のうち、余命が前記第2の閾値よりも大きい予備半導体記憶装置を選択すること

10

20

30

40

50

を特徴とする付記 3 に記載の情報処理装置。

【 0 1 1 2 】

(付記 5) 前記収集部は、

前記比率として、単位時間あたりに前記半導体記憶装置から読出されたデータ量と前記半導体記憶装置に書き込まれたデータ量との比率を収集し、

前記判定部は、

単位時間当たりの前記半導体記憶装置から読み出したデータ量が所定の閾値よりも多い場合、データの読出しの比率が多いと判定することを特徴とする付記 1 ~ 4 のいずれか 1 つに記載の情報処理装置。

【 0 1 1 3 】

(付記 6) 前記収集部は、

前記比率として、単位時間あたりに前記半導体記憶装置からデータを読出した回数とデータを書き込んだ回数との比率を収集し、

前記判定部は、

単位時間当たりの前記データの読出し回数が所定の閾値よりも多い場合、データの読出しの比率が多いと判定することを特徴とする付記 1 ~ 4 のいずれか 1 つに記載の情報処理装置。

【 0 1 1 4 】

(付記 7) 前記情報処理装置はさらに、

データが読み書きされる複数の半導体記憶装置と、

前記半導体記憶装置が、二重化されている場合、前記交換部が、前記半導体記憶装置を前記予備の半導体記憶装置と交換する間、前記半導体記憶装置に対する書き込みデータを保持する保持部と、

前記交換部が、交換対象の半導体記憶装置が記憶するデータと、前記保持部が保持したデータとを、前記予備半導体記憶装置と、前記複数の半導体記憶装置のうち前記予備半導体記憶装置と二重化される半導体記憶装置とに複製する複製部と

を有することを特徴とする付記 1 ~ 6 のいずれか 1 つに記載の情報処理装置。

【 0 1 1 5 】

(付記 8) 前記情報処理装置はさらに、

データが読み書きされる複数の半導体記憶装置を有し、

前記交換部は、

前記複数の半導体記憶装置に、データとデータを復元するためのパリティとが分散して格納されている場合、前記複数の半導体記憶装置のうち、いずれかの半導体記憶装置を前記予備半導体記憶装置と交換したとき、前記複数の半導体記憶装置のうち、前記予備半導体記憶装置以外の半導体記憶装置が記憶するデータとパリティとを用いて、前記予備半導体記憶装置が記憶するデータを生成することを特徴とする付記 1 ~ 6 のいずれか 1 つに記載の情報処理装置。

【 0 1 1 6 】

(付記 9) データが読み書きされる半導体記憶装置と、前記半導体記憶装置の予備である予備半導体記憶装置とを有する情報処理装置の制御プログラムにおいて、

前記情報処理装置に、

前記半導体記憶装置について、書き込まれたデータを消去できる残余回数に基づいた余命指標値、および、前記半導体記憶装置に対する読出しと書き込みとの読み書きの比率を収集させ、

収集された前記比率に基づいて、前記半導体記憶装置の交換基準となる閾値を決定させ

、
収集された前記余命指標値が、決定された前記閾値よりも短い場合は、前記半導体記憶装置を前記予備半導体記憶装置と交換させる

ことを特徴とする情報処理装置の制御プログラム。

【 0 1 1 7 】

(付記 10) データが読み書きされる半導体記憶装置と、前記半導体記憶装置の予備である予備半導体記憶装置とを有する情報処理装置の制御方法において、

前記情報処理装置が、

前記半導体記憶装置について、書き込まれたデータを消去できる残余回数に基づいた余命指標値、および、前記半導体記憶装置に対する読出しと書き込みとの読み書きの比率を収集し、

収集された前記比率に基づいて、前記半導体記憶装置の交換基準となる閾値を決定し、

収集された前記余命指標値が、決定された前記閾値よりも短い場合は、前記半導体記憶装置を前記予備半導体記憶装置と交換する

ことを特徴とする情報処理装置の制御方法。

10

【符号の説明】

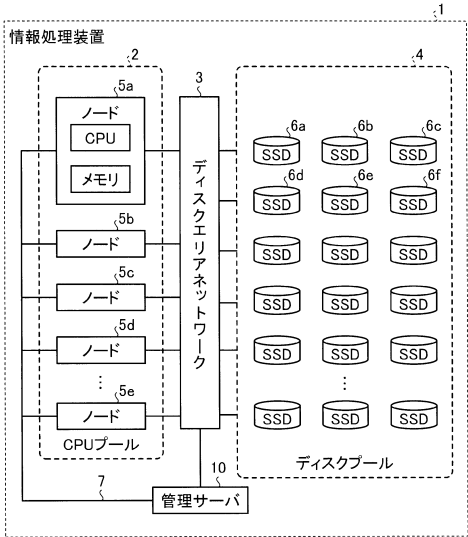
【0118】

- 1 情報処理装置
- 2 CPUプール
- 3 ディスクエリアネットワーク
- 4 ディスクプール
- 5 a ~ 5 e ノード
- 6 a ~ 6 f SSD
- 7 管理ネットワーク
- 10 管理サーバ
- 11 記憶部
- 12 余命管理テーブル
- 13 ワークロード管理テーブル
- 14 通信部
- 15 更新部
- 16 判定部
- 17 決定部
- 18 選択部
- 19 設定部

20

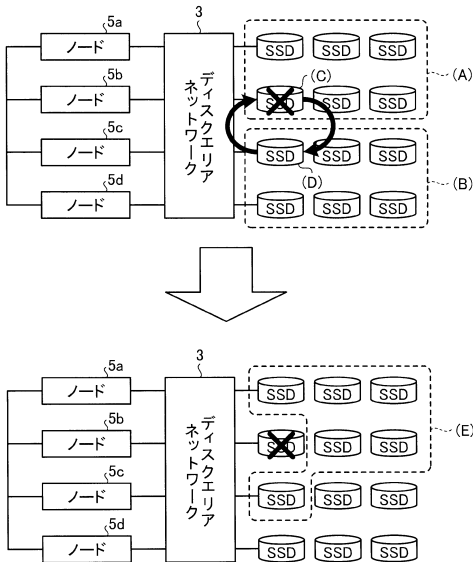
【図 1】

実施例1に係る情報処理装置を説明する図



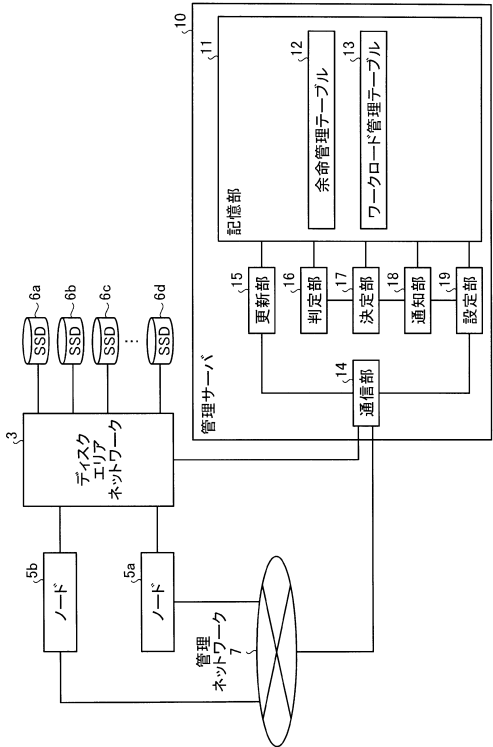
【図 2】

実施例1に係る情報処理装置が実行する処理を説明する図



【図 3】

実施例1に係る管理サーバの機能構成を説明する図



【図 4】

余命管理テーブルの一例を説明する図

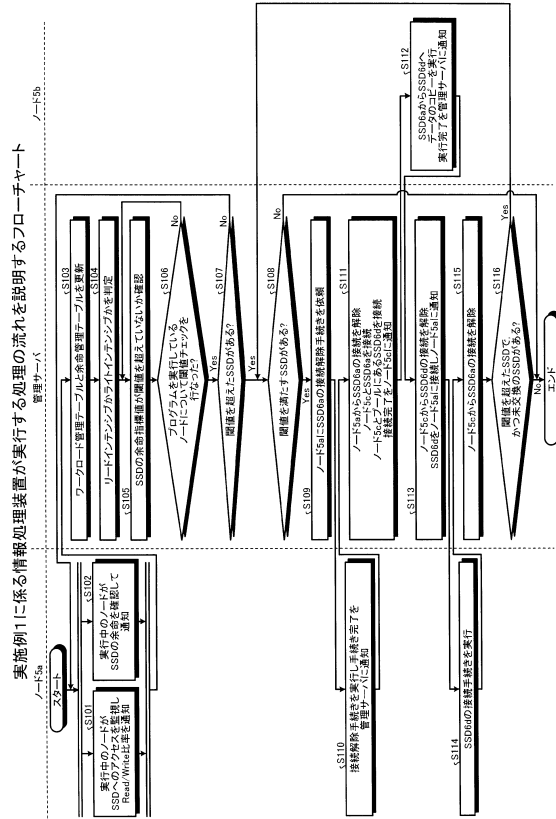
SSD ID	接続先	余命
SSD#0	ノード5a	100
SSD#1	ノード5b	80
SSD#2	-(未接続)	50
...		

【図 5】

ワークロード管理テーブルの一例を説明する図

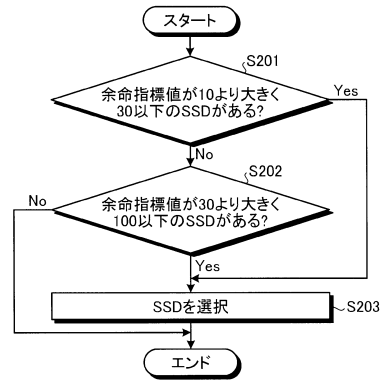
ノードID	プログラム実行フラグ	比率 (Byte/s)
ノード5a	1(実行している)	50:50
ノード5b	1(実行している)	90:10
ノード5c	0(実行していない)	null
...		

【 図 6 】



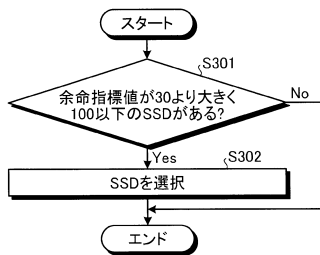
【圖 7】

選択部が実行する処理の流れを説明するための第1のフローチャート



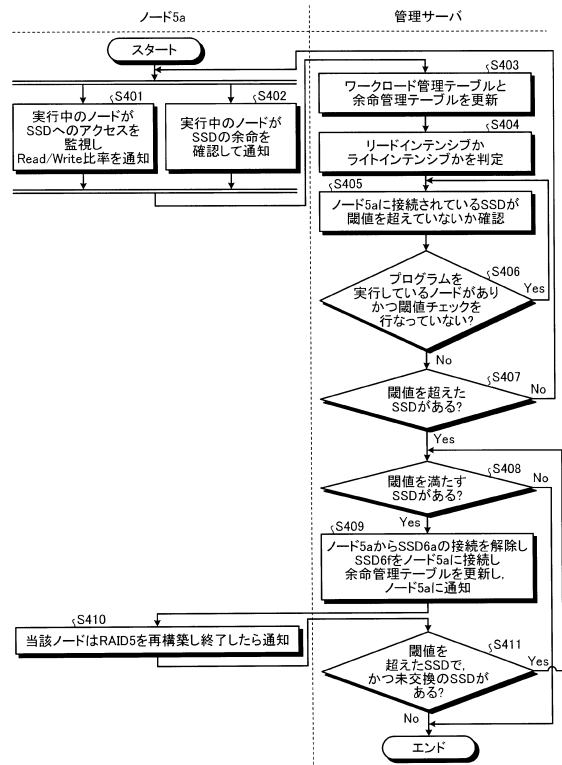
【 図 8 】

選択部が実行する処理の流れを説明するための第2のフローチャート



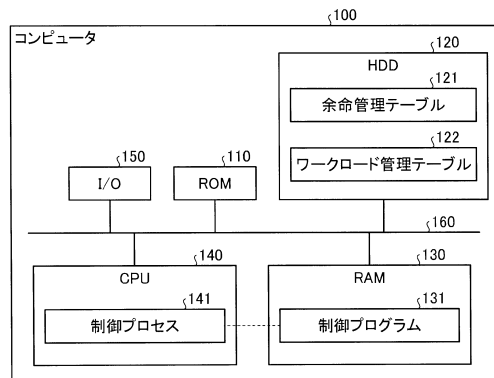
【 図 9 】

実施例2に係る情報処理装置が実行する処理の流れを説明するフローチャート



【図 10】

制御プログラムを実行するコンピュータを説明する図



フロントページの続き

(56)参考文献 特開 2 0 1 3 - 0 2 0 5 4 4 (J P , A)
特開 2 0 0 7 - 1 1 5 2 3 2 (J P , A)
特開 2 0 0 7 - 3 2 3 2 2 4 (J P , A)
特開 2 0 0 8 - 0 9 7 2 3 7 (J P , A)

(58)調査した分野(Int.Cl. , D B 名)
G 0 6 F 3 / 0 6
G 0 6 F 3 / 0 8