

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6094267号
(P6094267)

(45) 発行日 平成29年3月15日 (2017.3.15)

(24) 登録日 平成29年2月24日 (2017.2.24)

(51) Int. Cl.		F I			
G 0 6 F	12/00	(2006.01)	G 0 6 F	12/00	5 0 1 B
G 0 6 F	3/06	(2006.01)	G 0 6 F	3/06	3 0 2 J
			G 0 6 F	3/06	3 0 1 W

請求項の数 10 (全 26 頁)

<p>(21) 出願番号 特願2013-40842 (P2013-40842)</p> <p>(22) 出願日 平成25年3月1日 (2013.3.1)</p> <p>(65) 公開番号 特開2014-170304 (P2014-170304A)</p> <p>(43) 公開日 平成26年9月18日 (2014.9.18)</p> <p>審査請求日 平成28年2月5日 (2016.2.5)</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</p> <p>(74) 代理人 100124811 弁理士 馬場 資博</p> <p>(74) 代理人 100088959 弁理士 境 廣巳</p> <p>(72) 発明者 山本 拓明 東京都港区芝五丁目7番1号 日本電気株式会社内</p> <p>審査官 漆原 孝治</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】

記憶対象データを記憶装置に格納すると共に、当該記憶装置に既に記憶されている前記記憶対象データと同一のデータ内容の他の記憶対象データを前記記憶装置に格納する場合に、当該記憶装置に既に記憶されている前記記憶対象データを前記他の記憶対象データとして参照させるデータ格納制御部と、

前記記憶装置の所定の領域内においてデフラグ範囲とされた箇所に格納された前記記憶対象データを、前記記憶装置の他の領域内に格納し直すデフラグ処理部と、を備え、

前記データ格納制御部は、前記記憶装置に記憶されている前記記憶対象データ毎に、当該記憶対象データが他の記憶対象データとして参照されている数である参照数を記憶し、

前記デフラグ処理部は、前記記憶対象データの参照数に応じて当該記憶対象データを前記記憶装置の他の領域内において後にデフラグ範囲となる箇所に格納する、
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、

前記デフラグ処理部は、前記参照数が予め設定された閾値未満である前記記憶対象データを、前記記憶装置の他の領域内において後にデフラグ範囲となる箇所に格納し、前記参照数が予め設定された閾値以上である前記記憶対象データを、前記記憶装置の他の領域内において後に非デフラグ範囲となる箇所に格納する、
ストレージシステム。

【請求項 3】

請求項 2 に記載のストレージシステムであって、

前記記憶装置の他の領域は、予め設定されたデフラグ範囲と、データが未格納な範囲である未使用範囲と、が隣接して形成されており、

前記デフラグ処理部は、前記参照数が閾値未満である前記記憶対象データを、前記記憶装置の他の領域における前記未使用範囲内であり前記デフラグ範囲との隣接箇所からデータが連続して位置する箇所に格納し、前記未使用範囲内における前記デフラグ範囲との隣接箇所から連続して前記記憶対象データが格納された範囲と前記デフラグ範囲とを連結した範囲を新たなデフラグ範囲とする、
ストレージシステム。

10

【請求項 4】

請求項 3 に記載のストレージシステムであって、

前記記憶装置の他の領域は、予め設定された非デフラグ範囲と、データが未格納な範囲である未使用範囲と、が隣接して形成されており、

前記デフラグ処理部は、前記参照数が閾値以上である前記記憶対象データを、前記記憶装置の他の領域における前記未使用範囲であり前記非デフラグ範囲との隣接箇所からデータが連続して位置する箇所に格納し、前記未使用範囲内における前記非デフラグ範囲との隣接箇所から連続して前記記憶対象データが格納された範囲と前記非デフラグ範囲とを連結した範囲を新たな非デフラグ範囲とする、
ストレージシステム。

20

【請求項 5】

請求項 2 に記載のストレージシステムであって、

前記記憶装置の他の領域は、予め設定されたデフラグ範囲にデータが未格納な範囲である未使用範囲の一端側が隣接すると共に、当該未使用範囲の他端側が予め設定された非デフラグ範囲に隣接し、当該未使用範囲がデフラグ範囲と非デフラグ範囲とに挟まれて形成されており、

前記デフラグ処理部は、前記参照数が閾値未満である前記記憶対象データを、前記記憶装置の他の領域における前記未使用範囲内であり前記デフラグ範囲との隣接箇所からデータが連続して位置する箇所に格納し、前記未使用範囲内における前記デフラグ範囲との隣接箇所から連続して前記記憶対象データが格納された範囲と前記デフラグ範囲とを連結した範囲を新たなデフラグ範囲とし、前記参照数が閾値以上である前記記憶対象データを、前記記憶装置の他の領域における前記未使用範囲であり前記非デフラグ範囲との隣接箇所からデータが連続して位置する箇所に格納し、前記未使用範囲内における前記非デフラグ範囲との隣接箇所から連続して前記記憶対象データが格納された範囲と前記非デフラグ範囲とを連結した範囲を新たな非デフラグ範囲とする、
ストレージシステム。

30

【請求項 6】

請求項 1 乃至 5 のいずれかに記載のストレージシステムであって、

前記記憶装置の所定の領域は、予め設定されたデフラグ範囲と、データが未格納な範囲である未使用範囲と、が隣接して形成されており、

前記データ格納制御部は、前記記憶対象データを新たに前記記憶装置に格納する際に、前記記憶装置の所定の領域における前記未使用範囲内であり前記デフラグ範囲との隣接箇所からデータが連続して位置する箇所に格納し、前記未使用範囲内における前記デフラグ範囲との隣接箇所から連続して前記記憶対象データが格納された範囲と前記デフラグ範囲とを連結した範囲を新たなデフラグ範囲とする、
ストレージシステム。

40

【請求項 7】

請求項 4 に記載のストレージシステムであって、

前記記憶装置の他の領域は、前記デフラグ範囲と前記未使用範囲とが隣接して形成された第一領域と、当該第一領域とは異なり前記非デフラグ範囲と前記未使用範囲とが隣接し

50

て形成された第二領域と、を有する、
ストレージシステム。

【請求項 8】

請求項 7 に記載のストレージシステムであって、
前記第一領域は、所定の記憶装置に形成されており、
前記第二領域は、前記所定の記憶装置よりもデータ読み出し処理が高速な他の記憶装置
に形成されている、
ストレージシステム。

【請求項 9】

情報処理装置に、
記憶対象データを記憶装置に格納すると共に、当該記憶装置に既に記憶されている前記
記憶対象データと同一のデータ内容の他の記憶対象データを前記記憶装置に格納する場合
に、当該記憶装置に既に記憶されている前記記憶対象データを前記他の記憶対象データと
して参照させるデータ格納制御部と、

前記記憶装置の所定の領域内においてデフラグ範囲とされた箇所に格納された前記記憶
対象データを、前記記憶装置の他の領域内に格納し直すデフラグ処理部と、を実現させる
と共に、

前記データ格納制御部は、前記記憶装置に記憶されている前記記憶対象データ毎に、当
該記憶対象データが他の記憶対象データとして参照されている数である参照数を記憶し、

前記デフラグ処理部は、前記記憶対象データの参照数に応じて当該記憶対象データを前
記記憶装置の他の領域内において後にデフラグ範囲となる箇所に格納する、
ことを実現させるためのプログラム。

【請求項 10】

記憶対象データを記憶装置に格納すると共に、当該記憶装置に既に記憶されている前記
記憶対象データと同一のデータ内容の他の記憶対象データを前記記憶装置に格納する場合
に、当該記憶装置に既に記憶されている前記記憶対象データを前記他の記憶対象データと
して参照させてデータ格納制御を行うと共に、前記記憶装置に記憶されている前記記憶対
象データ毎に、当該記憶対象データが他の記憶対象データとして参照されている数である
参照数を記憶し、

前記記憶装置の所定の領域内においてデフラグ範囲とされた箇所に格納された前記記憶
対象データを、前記記憶装置の他の領域内に格納し直すデフラグ処理を実行し、

前記デフラグ処理時に、前記記憶対象データの参照数に応じて当該記憶対象データを前
記記憶装置の他の領域内において後にデフラグ範囲となる箇所に格納する、
デフラグ方法。

【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、同一内容のデータの重複記憶を排除す
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。

【 0 0 0 3 】

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。また、コンテンツアドレスストレージシステムの中には、所定のデータ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して、これら複数のフラグメントをそれぞれ複数の記憶装置にそれぞれ格納する、というものもある。

【 0 0 0 4 】

そして、上述したようなコンテンツアドレスストレージシステムでは、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。

10

【 0 0 0 5 】

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値、例えばデータのハッシュ値、に基づいて生成される。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、重複記録を排除して、データ容量の削減を図ることができる。

【 0 0 0 6 】

20

特に、上述したような重複記憶を排除する機能を有するストレージシステムでは、ファイルなど書き込み対象となるデータを所定容量の複数のブロックデータに分割して圧縮し、記憶装置に書き込む。このように、ファイルを分割したブロックデータ単位で重複記憶を排除することで、重複率が增大し、データ容量の削減を図っている。そして、バックアップを行うストレージシステムに提供することで、バックアップの容量を節約したり、レプリケーション時の帯域の節約を図っている。

【 0 0 0 7 】

一方で、上述した重複排除を行うストレージシステムでは、格納したデータが上書きされることはない。つまり、新しい非重複データを追記するか、古く参照されていないデータを削除するか、といったどちらかしか行われぬ。このため、重複排除ストレージシステムを長期間運用すると、ストレージ内部で空き領域が断片化し、入出力性能（I/O性能）が低下する。そのため、重複排除ストレージシステムでは、定期的なデフラグ処理が必要になる。

30

【 先行技術文献 】

【 特許文献 】

【 0 0 0 8 】

【 特許文献 1 】 特開 2 0 0 5 - 2 3 5 1 7 1 号 公 報

【 特許文献 2 】 特開 2 0 1 0 - 2 8 7 0 4 9 号 公 報

【 発明の概要 】

【 発明が解決しようとする課題 】

40

【 0 0 0 9 】

しかしながら、デフラグ処理は、ユーザによる運用とは関係のないストレージシステムにおける I/O 負荷を発生させるため、デフラグ期間中の I/O 性能を低下させてしまう、という問題が生じる。ここで、特許文献 2 には、メモリ内におけるページ分割数が閾値を超えるブロックに対してデフラグ処理を行う、という技術が開示されている。かかる技術によると、デフラグ処理を行う領域が制限され、I/O 性能の低下を抑制できる可能性があるが、データの内容を考慮してデフラグを行っていない。従って、重複排除ストレージシステムに格納される重複排除を伴うデータに対して、デフラグを効率よく行うことができるか不明であり、断片化を解消するためには結局のところ全ての領域に対してデフラグを行う必要がある。その結果、重複排除ストレージシステムにおける I/O 性能の低下

50

を解消することができない。

【 0 0 1 0 】

このため、本発明の目的は、重複排除ストレージシステムにおけるデフラグ処理による I / O 性能の低下、という問題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】

【 0 0 1 1 】

本発明の一形態であるストレージシステムは、

記憶対象データを記憶装置に格納すると共に、当該記憶装置に既に記憶されている前記記憶対象データと同一のデータ内容の他の記憶対象データを前記記憶装置に格納する場合に、当該記憶装置に既に記憶されている前記記憶対象データを前記他の記憶対象データとして参照させるデータ格納制御部と、

前記記憶装置の所定の領域内においてデフラグ範囲とされた箇所に格納された前記記憶対象データを、前記記憶装置の他の領域内に格納し直すデフラグ処理部と、を備え、

前記データ格納制御部は、前記記憶装置に記憶されている前記記憶対象データ毎に、当該記憶対象データが他の記憶対象データとして参照されている数である参照数を記憶し、

前記デフラグ処理部は、前記記憶対象データの参照数に応じて当該記憶対象データを前記記憶装置の他の領域内において後にデフラグ範囲となる箇所に格納する、という構成をとる。

【 0 0 1 2 】

本発明の他の形態であるプログラムは、

情報処理装置に、

記憶対象データを記憶装置に格納すると共に、当該記憶装置に既に記憶されている前記記憶対象データと同一のデータ内容の他の記憶対象データを前記記憶装置に格納する場合に、当該記憶装置に既に記憶されている前記記憶対象データを前記他の記憶対象データとして参照させるデータ格納制御部と、

前記記憶装置の所定の領域内においてデフラグ範囲とされた箇所に格納された前記記憶対象データを、前記記憶装置の他の領域内に格納し直すデフラグ処理部と、を実現させると共に、

前記データ格納制御部は、前記記憶装置に記憶されている前記記憶対象データ毎に、当該記憶対象データが他の記憶対象データとして参照されている数である参照数を記憶し、

前記デフラグ処理部は、前記記憶対象データの参照数に応じて当該記憶対象データを前記記憶装置の他の領域内において後にデフラグ範囲となる箇所に格納する、ことを実現させるためのプログラムである。

【 0 0 1 3 】

本発明の他の形態であるデフラグ方法は、

記憶対象データを記憶装置に格納すると共に、当該記憶装置に既に記憶されている前記記憶対象データと同一のデータ内容の他の記憶対象データを前記記憶装置に格納する場合に、当該記憶装置に既に記憶されている前記記憶対象データを前記他の記憶対象データとして参照させてデータ格納制御を行うと共に、前記記憶装置に記憶されている前記記憶対象データ毎に、当該記憶対象データが他の記憶対象データとして参照されている数である参照数を記憶し、

前記記憶装置の所定の領域内においてデフラグ範囲とされた箇所に格納された前記記憶対象データを、前記記憶装置の他の領域内に格納し直すデフラグ処理を実行し、

前記デフラグ処理時に、前記記憶対象データの参照数に応じて当該記憶対象データを前記記憶装置の他の領域内において後にデフラグ範囲となる箇所に格納する、という構成をとる。

【発明の効果】

【 0 0 1 4 】

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けるデフラグ

10

20

30

40

50

による I / O 性能の低下を抑制することができる。

【図面の簡単な説明】

【0015】

【図1】本発明の実施形態1におけるストレージシステムを含むシステム全体の構成を示すブロック図である。

【図2】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の概略を示すブロック図である。

【図3】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4】図3に開示したストレージシステムにおけるデータ書き込み処理の様子を説明するための説明図である。 10

【図5】図3に開示したストレージシステムにおけるデータ書き込み処理の様子を説明する説明図である。

【図6】図3に開示したストレージシステムにおけるデータ構造を説明する説明図である。

【図7】図3に開示したストレージシステムにおけるデータ書き込み処理の様子を説明する説明図である。

【図8】図3に開示したストレージシステムにおけるデータ書き込み処理の様子を説明する説明図である。

【図9】図3に開示したストレージシステムにおけるデータ書き込み処理の様子を説明する説明図である。 20

【図10】図3に開示したストレージシステムにおけるデフラグ処理の様子を説明する説明図である。

【図11】図3に開示したストレージシステムにおけるデフラグ処理の様子を説明する説明図である。

【図12】図3に開示したストレージシステムにおけるデフラグ処理の様子を説明する説明図である。

【図13】図3に開示したストレージシステムにおけるデフラグ処理の様子を説明する説明図である。

【図14】図3に開示したストレージシステムにおけるデータ書き込み処理の動作を示すフローチャートである。 30

【図15】図3に開示したストレージシステムにおけるデフラグ処理の動作を示すフローチャートである。

【図16】本発明の実施形態2におけるストレージシステムに記憶されたデータの様子を示す図である。

【図17】実施形態2のストレージシステムにおけるデフラグ処理の様子を説明する説明図である。

【図18】実施形態2のストレージシステムにおけるデフラグ処理の様子を説明する説明図である。

【図19】実施形態2のストレージシステムにおけるデフラグ処理の様子を説明する説明図である。 40

【図20】実施形態2のストレージシステムにおけるデフラグ処理の動作を示すフローチャートである。

【図21】実施形態2のストレージシステムの他の構成を示すブロック図である。

【図22】本発明の付記1におけるストレージシステムの構成を示すブロック図である。

【発明を実施するための形態】

【0016】

<実施形態1>

本発明の第1の実施形態を、図1乃至図15を参照して説明する。図1は、システム全体の構成を示すブロック図である。図2は、ストレージシステムの概略を示すブロック図 50

であり、図3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。図4乃至図13は、ストレージシステムにおけるデータ書き込み処理及びデフラグ処理の動作を説明するための説明図である。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。

【0017】

ここで、本実施形態は、後述する付記に記載のストレージシステム等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。

10

【0018】

[構成]

図1に示すように、本発明におけるストレージシステム1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4に接続している。そして、バックアップシステム4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5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1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1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。

【0019】

そして、図2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、複数のサーバコンピュータが接続された構成を採っている。具体的に、ストレージシステム1は、ストレージシステム1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2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3と、を備えている。なお、アクセラレータノード2の数とストレージノード3の数は、図2に示したものに限定されず、さらに多くの各ノード2,3が接続されて構成されていてもよい。

20

【0020】

さらに、本実施形態におけるストレージシステム1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、記憶するデータの内容に応じて設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。このコンテンツアドレスストレージシステムについては、後に詳述する。

30

【0021】

なお、以下では、ストレージシステム1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1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1が有する構成及び機能は、アクセラレータノード2あるいはストレージノード3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1は、図2に示すように、必ずしもアクセラレータノード2とストレージノード3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く、例えば、1台のコンピュータにて構成されていてもよい。さらには、ストレージシステム1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限定されず、重複排除機能を有しているストレージシステムであればよい。

40

【0022】

図3に、本実施形態におけるストレージシステム1の構成を示す。この図に示すように、ストレージシステム1は、サーバコンピュータにて構成され、所定の演算処理を行う演算装置(図示せず)と、記憶装置20と、を備える。そして、ストレージシステム1は、上記演算装置にプログラムが組み込まれることで構築された、書き込み部11(データ格納制御部)、読み出し部12、削除部13、デフラグ部14(デフラグ処理部)、を備えている。

【0023】

なお、実際には、上述したストレージシステム1が備える構成は、図2に示したアクセ

50

ラレータノード2及びストレージノード3がそれぞれ備えているCPU (Central Processing Unit) などの演算装置やハードディスクドライブなどの記憶装置にて構成されている。

【0024】

ここで、上述したように、本実施形態におけるストレージシステム1は、コンテンツアドレスストレージシステムである。このため、ストレージシステム1は、コンテンツアドレスを利用してデータを記憶装置20に格納する機能を有しており、以下に説明するように、データを分割及び分散し、かつ、コンテンツアドレスにて格納位置を特定して、データを格納する。以下、ストレージシステム1にてコンテンツアドレスを利用した書き込み部11によるデータ書き込み処理の一例、及び、読み出し部12による読み出し処理の一例について、図4及び図5を参照して説明する。但し、以下に説明する処理は、重複排除型のストレージシステムにおけるデータの書き込み処理及び読み出し処理の一例である。従って、本発明におけるストレージシステム1は、以下の方法でデータを書き込んだり読み出すことに限定されず、他の方法で重複記憶を排除してデータの書き込み処理及び読み出し処理を行ってもよい。

10

【0025】

まず、図4及び図5の矢印Y1に示すように、ストレージシステム1が書き込み要求されたファイルAの入力を受ける。すると、図4及び図5の矢印Y2に示すように、ファイルAを所定容量 (例えば、64KB) のブロックデータDに分割する。

【0026】

続いて、分割されたブロックデータDのデータ内容に基づいて、当該データ内容を代表する固有のハッシュ値Hを算出する (図5の矢印Y3)。例えば、ハッシュ値Hは、予め設定されたハッシュ関数を用いて、ブロックデータDのデータ内容から算出する。

20

【0027】

続いて、ファイルAのブロックデータDのハッシュ値Hを用いて、当該ブロックデータDが既に格納されているか否かを調べる。具体的には、まず、既に格納されているブロックデータDは、そのハッシュ値Hと格納位置を表すコンテンツアドレスCAとが、関連付けられてMFI (Main Fragment Index) ファイルに登録されている。従って、格納前に算出したブロックデータDのハッシュ値HがMFIファイル内に存在している場合には、既に同一内容のブロックデータDが格納されていると判断できる (図5の矢印Y4)。この場合には、格納前のブロックデータDのハッシュ値Hと一致したMFI内のハッシュ値Hに関連付けられているコンテンツアドレスCAを、当該MFIファイルから取得する。そして、このコンテンツアドレスCAを、書き込み要求されたブロックデータDのコンテンツアドレスCAとして返却する。

30

【0028】

そして、返却されたコンテンツアドレスCAが参照する既に格納されているデータを、書き込み要求されたブロックデータDとして使用する。つまり、書き込み要求されたブロックデータDの格納先として、返却されたコンテンツアドレスCAが参照する領域を指定することで、当該書き込み要求されたブロックデータDを記憶したこととする。これにより、書き込み要求にかかるブロックデータDを、実際に記憶装置20内に記憶する必要がなくなる。

40

【0029】

また、書き込み要求にかかるブロックデータDがまだ記憶されていないと判断された場合には、以下のようにして書き込み要求にかかるブロックデータDの書き込みを行う。まず、書き込み要求にかかるブロックデータDを圧縮して、図5の矢印Y5に示すように、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、図4の符号D1~D9に示すように、9つのフラグメントデータ (分割データ91) に分割する。そしてさらに、分割したフラグメントデータのうちのいくつかが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91に追加する。例えば、図4の符号D10~D12に示すように、3つのフラグメントデータ (冗

50

長データ92)を追加する。これにより、9つの分割データ91と、3つの冗長データ92とにより構成される12個のフラグメントデータからなるデータセット90を生成する。

【0030】

続いて、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、記憶装置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、図4に示すように、12個のフラグメントデータD1~D12を生成した場合には、複数の記憶装置内にそれぞれ形成したデータ格納ファイルに、各フラグメントデータD1~D12を1つずつそれぞれ格納する(図5の矢印Y6参照)。

【0031】

続いて、ストレージシステム1は、上述したように格納したフラグメントデータD1~D12の格納位置、つまり、当該フラグメントデータD1~D12にて復元されるブロックデータDの格納位置を表すコンテンツアドレスCAを生成して管理する。具体的には、格納したブロックデータDの内容に基づいて算出したハッシュ値Hの一部(ショートハッシュ)(例えば、ハッシュ値Hの先頭8B(バイト))と、論理格納位置を表す情報と、を組み合わせて、コンテンツアドレスCAを生成する。そして、このコンテンツアドレスCAを、ストレージシステム1内のファイルシステムに返却する(図5の矢印Y7)。すると、ストレージシステム1は、バックアップ対象データのファイル名などの識別情報と、コンテンツアドレスCAとを関連付けてファイルシステムで管理する。

【0032】

また、ブロックデータDのコンテンツアドレスCAと、当該ブロックデータDのハッシュ値Hと、を関連付けて、各ストレージノード3がMFIファイルにて管理する。このように、上記コンテンツアドレスCAは、ファイルを特定する情報やハッシュ値Hなどと関連付けられて、アクセラレータノード2やストレージノード3の記憶装置に格納される。

【0033】

さらに、ストレージシステム1は、上述したように格納したファイルを読み出す制御を行う。例えば、ストレージシステム1に対して、特定のファイルを指定して読み出し要求があると、まず、ファイルシステムに基づいて、読み出し要求にかかるファイルに対応するハッシュ値の一部であるショートハッシュと論理位置の情報からなるコンテンツアドレスCAを指定する。そして、コンテンツアドレスCAがMFIファイルに登録されているか否かを調べる。登録されていない場合は、要求されたデータは格納されていないため、エラーを返却する。

【0034】

一方、読み出し要求にかかるコンテンツアドレスCAが登録されている場合には、上記コンテンツアドレスCAにて指定される格納位置を特定し、この特定された格納位置に格納されている各フラグメントデータを、読み出し要求されたデータとして読み出す。このとき、各フラグメントが格納されているデータ格納ファイルと、当該データ格納ファイルのうち1つのフラグメントデータの格納位置が分かれば、同一の格納位置から他のフラグメントデータの格納位置を特定することができる。

【0035】

そして、読み出し要求に応じて読み出した各フラグメントデータからブロックデータDを復元する。さらに、復元したブロックデータDを複数連結し、ファイルAなどの一群のデータに復元して返却する。

【0036】

次に、上述したように記憶装置20内に書き込まれたデータに対するデフラグ処理を行うための構成を、さらに説明する。図3に示すように、ストレージシステム1は、領域管理表15、ブロック管理表16、ファイル管理表17を備える。これらの表15, 16, 17は、図示しない主記憶装置に形成されている。また、補助記憶装置にて構成されている記憶装置20の記憶領域は、複数の領域21に分割されて管理されている。

【0037】

10

20

30

40

50

図6は、領域管理表15、ブロック管理表16、ファイル管理表17の詳細、及び、これら各表と記憶装置20及び各領域21との関係を示している。

【0038】

領域管理表15の「領域番号」の列は、各領域21を識別するための一意な番号を表す。そして、記憶装置20内の領域21は、書き込み時において同時に1つだけ書き込み可能となり、その書き込み可能な領域21のみを「使用中フラグ」を「yes」とすることで特定する。「領域サイズ」は、各領域21に格納可能なデータの合計サイズを表す。「デフラグ範囲」は、後述するようにデフラグ部14がデフラグを行う際の対象とすべき範囲を表す。「未使用範囲」は、一切のデータが格納されていない範囲であり、後述するように書き込み部11やデフラグ部14によってデータが格納される範囲を表す。「長寿命範囲」は、後述するようにデフラグ部14によって重複率が高いブロックが格納される範囲を表す。「最終書込時刻」は、領域21に最後にデータが書き込まれたときの時刻である。

10

【0039】

ブロック管理表16の「ハッシュ値」の列は、ファイルを分割した各ブロックデータのハッシュ値を表している。このハッシュ値は、ブロックデータのデータ内容を代表する値であり、書き込み部11は、このハッシュ値を用いて、上述したように書き込み対象であるブロックデータが既に記憶装置20に格納されているか否かを調べる。「重複率」の列は、各ブロックデータの重複率を表している。重複率は、ブロックデータが他のブロックデータとして参照されている数(参照数)の値である。「領域番号」、「物理アドレス」及び「ブロックサイズ」は、各ブロックデータが記憶装置20において格納されている物理的な位置を指し示している。

20

【0040】

ファイル管理表17は、ファイルを構成するブロックデータのリストを、ハッシュ値のリストとして保持している。

【0041】

そして、本実施形態における書き込み部11は、ファイルを記憶装置20の領域21に格納する際には、当該ファイルを分割したブロックデータのハッシュ値を、上記ファイル管理表17及びブロック管理表16に登録する。このとき、ブロックデータのハッシュ値から、当該ブロックデータが既に記憶装置20内の領域21に記憶されているか否かを判断する。そして、記憶対象となるブロックデータが既に記憶装置20内に存在している場合には、当該存在しているブロックデータを参照し、ブロック管理表16の該当ハッシュ値の行における重複率(参照する)の値をインクリメントする。

30

【0042】

一方、記憶対象となるブロックデータのハッシュ値が未知のものである場合は、領域管理表15を参照し、「使用中フラグ」が「yes」となっている領域21を見つけ出し、該当領域21の未使用範囲の先頭に該当ブロックデータを書き込む。そして、ブロックデータのサイズ分だけ未使用範囲の先頭を縮小し、デフラグ範囲の末尾を拡大し、領域管理表15の該当領域21の最終書き込み時刻を更新する。加えて、ブロック管理表16に新規エントリを追加し、書き込みの際に使用した領域番号、物理アドレス、およびブロックサイズを格納する。最後に、ファイル管理表17の該当ファイル名のリストに書き込んだブロックのハッシュ値を追記する。

40

【0043】

また、読み出し部12は、ファイル管理表17を参照し、読み出し対象となるファイルを構成するブロックデータのハッシュ値を列挙する。そして、ブロック管理表16を参照し、該当ブロックが格納されている領域、物理アドレスおよびブロックサイズを割り出し、該当する格納場所からブロックデータの読み出しを行う。

【0044】

また、削除部13は、削除対象のファイルを構成するブロックデータに該当するブロック管理表16のエントリにおける「重複率」をデクリメントする。重複率が0になった場

50

合、該当するエントリをブロック管理表 1 6 から削除する。ただし、重複率は、削除時に即座にデクリメントされる必要はなく、ファイルつまりブロックデータ削除後の適切なタイミングでデクリメントすることも可能である。

【 0 0 4 5 】

次に、本実施形態におけるデフラグ部 1 4 について説明する。デフラグ部 1 4 は、所定のタイミングで、記憶装置 2 0 の所定の領域 2 1 であるデフラグ元領域内に格納されたブロックデータを、他の領域 2 1 であるデフラグ先領域内に格納し直すデフラグ処理を実行する。このとき、デフラグ部 1 4 は、デフラグ元領域のデフラグ範囲にあるブロックデータをデフラグ対象とし、当該ブロックデータの重複率に応じて、デフラグ先領域に格納する箇所を決定する。

10

【 0 0 4 6 】

例えば、予め設定された一定値（閾値）未満の重複率であるブロックデータは、デフラグ先領域の未使用範囲のうち、デフラグ範囲に隣接する箇所から格納データが連続するよう格納する。そして、ブロックデータ格納後やデフラグ終了時には、デフラグ先領域のデフラグ範囲から未使用範囲に格納したブロックデータまでの連続する範囲を新たなデフラグ範囲とし、その分、未使用範囲を縮小する。

【 0 0 4 7 】

一方、予め設定された一定値（閾値）以上の高い重複率であるブロックデータは、デフラグ先領域の未使用範囲のうち、長寿命範囲（非デフラグ範囲）に隣接する箇所から格納データが連続するよう格納する。そして、ブロックデータ格納後やデフラグ終了時には、

20

【 0 0 4 8 】

[動作]

次に、上述した構成のストレージシステムの動作を、図 7 乃至図 1 5 を参照して説明する。まず、ストレージシステムにファイルを書き込むときの動作を、図 7 乃至図 9 と、図 1 4 のフローチャートを参照して説明する。

【 0 0 4 9 】

まず、書き込み部 1 1 は、書き込み対象のファイル（データストリーム）をフィンガープリントに基づいて所定容量のブロックデータに分割する（ステップ S 1 ）。次に、分割

30

【 0 0 5 0 】

その後、書き込み部 1 1 は、ブロック管理表 1 6 を確認し、該当するハッシュ値を持つブロックデータが記憶装置 2 0 内に存在しているか否かを確認する（ステップ S 3 ）。ブロックデータが既に記憶装置 2 0 内に存在する場合は（ステップ S 3 : Y e s ）、書き込み対象であるブロックデータとして既に存在するブロックデータを参照し、ブロック管理表 1 6 内の参照されたブロックデータに該当するハッシュ値の行における重複率の値をインクリメントする（ステップ S 4 ）。

【 0 0 5 1 】

一方、書き込み部 1 1 は、書き込み対象である該当するブロックデータのハッシュ値が未知のもの、つまり、記憶装置 2 0 に格納されていない場合は（ステップ S 3 : N o ）、領域管理表 1 5 を参照し、使用中フラグが「yes」となっている領域 2 1（所定の領域）を見つけ出す。ここでは、該当する領域として、図 7 の書き込み領域 3 0 が見つかったとする。なお、書き込み領域 3 0 は、それぞれ所定の容量を有するデフラグ範囲 3 1、未使用範囲 3 2、長寿命範囲 3 3 が、この順に連なって設定された領域である。つまり、未使用範囲 3 2 がデフラグ範囲 3 1 と長寿命範囲 3 3 とに挟まれており、未使用範囲 3 2 の先頭（一端側）にデフラグ範囲 3 1 が隣接しており、未使用範囲 3 2 の末尾（他端側）に長寿命範囲 3 3 が隣接している。なお、書き込み領域 3 0 には、必ずしも長寿命範囲 3 3 が設定されていなくてもよい。

40

【 0 0 5 2 】

50

そして、書き込み部 11 は、図 8 に示すように、書き込み領域 30 の未使用範囲 32 の先頭、つまり、デフラグ範囲 31 との隣接箇所に、ブロックデータを書き込む（ステップ 55、図 8 の網掛け部分参照）。なお、未使用範囲 32 のデフラグ範囲 31 との隣接箇所側に既に他のブロックデータが格納されている場合には、当該デフラグ範囲 31 からデータが連続して位置する箇所に、ブロックデータを書き込む。

【 0053 】

その後、書き込み部 11 は、図 9 に示すように、書き込んだブロックデータのサイズ分だけデフラグ範囲 31 の末尾を拡大すると共に、当該書き込んだブロックデータのサイズ分だけ未使用範囲 32 の先頭を縮小する。換言すると、未使用範囲 31 の先頭側に位置するデフラグ範囲 31 との隣接箇所から連続してブロックデータが格納された範囲（図 9 の網掛け部分）と、デフラグ範囲（図 8 の符号 31）と、を連結した範囲を、新たなデフラグ範囲（図 9 の符号 31）とする。そして、領域管理表 15 における書き込み領域 30 の最終書込時刻を更新する。

10

【 0054 】

加えて、書き込み部 11 は、ブロック管理表 16 に新規エントリを追加し、ブロックデータを書き込んだ領域の領域番号、物理アドレス、および、ブロックデータのサイズを格納する。最後に、ファイル管理表 17 の該当ファイル名のリストに、書き込んだブロックデータのハッシュ値を追記する。

【 0055 】

次に、読み出し部 12 による動作を説明する。ファイルの読み出し要求があると、読み出し部 12 は、ファイル管理表 17 を参照し、ファイルを構成する各ブロックデータのハッシュ値を読み出して列挙する。そして、ブロック管理表 16 を参照し、該当ブロックデータが格納されている領域および物理アドレスを割り出し、かかる格納場所からブロックデータの読み出し、及び、ファイルの復元を行うことで、要求されたファイルを読み出す。

20

【 0056 】

次に、削除部 13 による動作を説明する。ファイルの削除要求があると、削除部 13 は、削除対象のブロックデータに該当するブロック管理表 16 のエントリにおいて、重複率をデクリメントする。重複率が 0 になった場合、該当するエントリをブロック管理表 16 から削除する。ただし、削除時に即座に重複率をデクリメントする必要はなく、データ削除後の適切なタイミングでデクリメントすることも可能である。

30

【 0057 】

次に、デフラグ部 14 によるデフラグ処理の動作を、図 10 乃至図 13 と、図 15 のフローチャートを参照して説明する。書き込み実施後、使用中フラグが「yes」となっている領域 21 について未使用領域が空になった場合、デフラグを実行する。デフラグには、2つの領域を用いる。このとき、図 13 に示すように領域管理表 15 を用いて、デフラグ元となる領域 21 は、デフラグ範囲が空でない領域のうち最終書き込み時刻が最も古い領域を選択する。また、デフラグ先となる領域 21 は、デフラグ範囲が空（0～0）である領域を選択する。これによって、ここでは図 10 に示すように、デフラグ元となる領域 21 としてデフラグ元領域 40（図 13 の領域 2）が選択され、デフラグ先となる領域 21 としてデフラグ先領域 50（図 13 の領域 1）が選択されたこととする。そして、デフラグ元領域 40 からデフラグ先領域 50 へブロックデータを移動することでデフラグを行なう。

40

【 0058 】

なお、デフラグ元領域 40 及びデフラグ先領域 50 は、それぞれ所定の容量を有するデフラグ範囲 41, 51、未使用範囲 42, 52、長寿命範囲 43, 53 が、この順に連なって設定された領域である。つまり、未使用範囲 42, 52 がデフラグ範囲 41, 51 と長寿命範囲 43, 53 とに挟まれており、未使用範囲 42, 52 の先頭（一端側）にデフラグ範囲 41, 51 が隣接しており、未使用範囲 42, 52 の末尾（他端側）に長寿命範囲 43, 53 が隣接している。

50

【 0 0 5 9 】

まず、デフラグ部 1 4 は、デフラグ元領域 4 0 内のデフラグ範囲 4 1 に含まれるブロックデータを列挙する（ステップ S 1 1）。次に、デフラグ部 1 4 は、各ブロックデータについて、ブロック管理表 1 6 を参照し、該当ブロックデータの重複率を確認する（ステップ S 1 2）。

【 0 0 6 0 】

そして、デフラグ部 1 4 は、該当ブロックデータの重複率が一定値未満である場合は（ステップ S 1 3：Yes）、図 1 1 の矢印 M 2，M 4 に示すように、デフラグ先領域 5 0 の未使用領域 5 2 の先頭にブロックデータを移動する（ステップ S 1 4）。つまり、デフラグ先領域 5 0 の未使用範囲 5 2 のうち、デフラグ範囲 5 1 との隣接箇所にブロックデータを
10
書き込む（図 1 1 の網掛け部分参照）。なお、未使用範囲 5 2 のデフラグ範囲 5 1 との隣接箇所側に既に他のブロックデータが格納されている場合には、当該デフラグ範囲 5 1 からデータが連続して位置する箇所に、ブロックデータを書き込む。

【 0 0 6 1 】

その後、デフラグ部 1 4 は、図 1 2 に示すように、デフラグ先領域 5 0 において、書き込んだブロックデータのサイズ分だけデフラグ範囲 5 1 の末尾を拡大すると共に、当該書き込んだブロックデータのサイズ分だけ未使用範囲 5 2 の先頭を縮小する。換言すると、未使用範囲 5 2 の先頭側に位置するデフラグ範囲 5 1 との隣接箇所から連続してブロック
20
データが格納された範囲（図 1 1 の網掛け部分）と、デフラグ範囲（図 1 1 の符号 5 1）と、を連結した範囲を、新たなデフラグ範囲（図 1 2 の符号 5 1）とする。このデフラグ先領域 5 0 の新たなデフラグ範囲 5 1 は、後にデフラグ対象となる。そして、領域管理表 1 5 におけるデフラグ先領域 5 0 の最終書込時刻を更新する。

【 0 0 6 2 】

また、デフラグ部 1 4 は、該当ブロックデータの重複率が一定値以上である場合は（ステップ S 1 3：No）、図 1 1 の矢印 M 1，M 3 に示すように、デフラグ先領域 5 0 の未使用領域 5 2 の末尾にブロックデータを移動する（ステップ S 1 5）。つまり、デフラグ先領域 5 0 の未使用範囲 5 2 のうち、長寿命範囲 5 3 との隣接箇所にブロックデータを書き込む（図 1 1 の網掛け部分参照）。なお、未使用範囲 5 2 の長寿命範囲 5 3 との隣接箇所側に既に他のブロックデータが格納されている場合には、当該長寿命範囲 5 1 からデータ
30
が連続して位置する箇所に、ブロックデータを書き込む。

【 0 0 6 3 】

その後、デフラグ部 1 4 は、図 1 2 に示すように、デフラグ先領域 5 0 において、書き込んだブロックデータのサイズ分だけ長寿命範囲 5 3 の先頭を拡大すると共に、当該書き込んだブロックデータのサイズ分だけ未使用範囲 5 2 の末尾を縮小する。換言すると、未使用範囲 5 2 の末尾側に位置する長寿命範囲 5 3 との隣接箇所から連続してブロックデータが格納された範囲（図 1 1 の網掛け部分）と、長寿命範囲（図 1 1 の符号 5 3）と、を
40
連結した範囲を、新たな長寿命範囲（図 1 2 の符号 5 3）とする。

【 0 0 6 4 】

上述したデフラグ処理により、デフラグ元領域 4 0 のデフラグ範囲 4 1 内のブロックデータの重複率が「0」となった場合、ブロック管理表 1 6 の該当する行を削除することで、
40
該当ブロックデータを破棄する。この時点で、デフラグ先領域 5 0 の未使用範囲 5 2 が空になった場合、デフラグを中断する。また、ブロック管理表 1 6 の該当するハッシュ値をもつ行について、領域番号と物理アドレスを、移動先を指し示すように更新する。これらの処理がすべて完了した時点で、デフラグ元領域 4 0 のデフラグ範囲 4 1 を未使用範囲 4 2 に統合し、当該デフラグ範囲 4 1 を「0-0」とする。

【 0 0 6 5 】

デフラグ後、デフラグ部 1 4 あるいは書き込み部 1 1 は、デフラグ先領域 5 0 の使用中フラグを「yes」に変更し、それ以外の領域 2 1 すべてについて使用中フラグを「no」にする。

【 0 0 6 6 】

10

20

30

40

50

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、デフラグ時に、重複率の高いブロックデータと重複率の低いブロックデータとを、別々の場所（デフラグ範囲と長寿命範囲）に格納している。これにより、次回のデフラグ時には、重複率の高いブロックを格納している長寿命範囲をデフラグの対象から除外し、重複率の低いブロックデータを格納しているデフラグ範囲のみをデフラグ対象としている。このため、領域 2 1 全体をデフラグする場合と比較して、デフラグ処理に要する I/O 負荷を低減させることができる。

【 0 0 6 7 】

また、データ書き込み時には、領域管理表 1 5 における使用中フラグを用いることで、書き込み可能な領域を限定している。これにより、他の書き込み可能でない領域についてはデータの書き込みが発生せず、ブロックの削除のみが行なわれる。そのため、時間経過とともに、書き込みが行なわれていない領域は空き領域に変わることが期待できる。

10

【 0 0 6 8 】

さらに、重複率の低いデータは、時間経過とともに大半が削除され、連続した空き領域に変わることが期待できる。そのため、重複率の低いデータをまとめて一箇所に格納し、しばらく経ってから空き領域として利用することができる。一方で、重複率の高いデータは、以後削除される可能性が低い。そのため、一箇所にまとめて格納しておき、空き領域の断片化の発生を抑制することができる。

【 0 0 6 9 】

< 実施形態 2 >

20

次に、本発明の第 2 の実施形態を、図 1 6 乃至図 2 0 を参照して説明する。図 1 6 は、領域管理表の一例を示す図である。図 1 7 乃至図 1 9 は、ストレージシステムにおけるデフラグ処理の動作を説明するための説明図である。図 2 0 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。

【 0 0 7 0 】

本実施形態におけるストレージシステムは、上述した実施形態 1 とほぼ同様の構成をとっているが、デフラグ時におけるブロックデータを格納するデフラグ先領域の構成が異なる。以下、実施形態 1 と異なる構成について主に説明する。

【 0 0 7 1 】

まず、本実施形態における記憶装置 2 0 内の領域 2 1 は、図 1 6 の領域管理表 1 5 に示すように、主に重複率の低いブロックデータを格納する「短寿命領域」と、主に重複率の高いブロックデータを格納するための「長寿命範囲」といった専用領域が設定されている。

30

【 0 0 7 2 】

そして、図 1 7 に示すように、「短寿命領域」は、デフラグ時に、一次デフラグ先領域 6 0（第一領域）として選択される。この「短寿命領域」である一次デフラグ先領域 6 0 は、所定の容量を有するデフラグ範囲 6 1、未使用範囲 6 2、長寿命範囲 6 3 が、この順に連なって設定された領域である。つまり、未使用範囲 6 2 がデフラグ範囲 6 1 と長寿命範囲 6 3 とに挟まれており、未使用範囲 6 2 の先頭（一端側）にデフラグ範囲 6 1 が隣接しており、未使用範囲 6 2 の末尾（他端側）に長寿命範囲 6 3 が隣接している。

40

【 0 0 7 3 】

また、図 1 7 に示すように、「長寿命領域」は、デフラグ時に、二次デフラグ先領域 7 0（第二領域）として選択される。この「長寿命領域」である二次デフラグ先領域 7 0 は、所定の容量を有するデフラグ範囲 7 1、未使用範囲 7 2、長寿命範囲 7 3 が、この順に連なって設定された領域である。つまり、未使用範囲 7 2 がデフラグ範囲 7 1 と長寿命範囲 7 3 とに挟まれており、未使用範囲 7 2 の先頭（一端側）にデフラグ範囲 7 1 が隣接しており、未使用範囲 7 2 の末尾（他端側）に長寿命範囲 7 3 が隣接している。

【 0 0 7 4 】

次に、デフラグ時の動作を説明する。まず、図 1 6 に示す領域管理表 1 5 より、デフラグ元となるデフラグ元領域 4 0 として、デフラグ範囲が空でない領域のうち、最終書き込

50

み時刻が最も古い領域を選択する。また、デフラグ先となるデフラグ先領域 60, 70 は、短寿命用と長寿命用の 2 種類を選択する。短寿命用の一次デフラグ領域 60 として、特性が短寿命の領域のうち未使用範囲が最大である領域を選択する。長寿命用の二次デフラグ領域 70 として、特性が長寿命の領域のうちデフラグ範囲が空である領域を選択する。ここでは、図 17 に示すように、短寿命用の一次デフラグ領域 60 と、長寿命用の二次デフラグ領域 70 と、が選択されたこととする。そして、デフラグ元領域 40 から一次デフラグ先領域 60、二次デフラグ領域 70 へブロックデータを移動することでデフラグを行う。

【0075】

まず、デフラグ部 14 は、デフラグ元領域 40 内のデフラグ範囲 41 に含まれるブロックデータを列挙する (ステップ S21)。次に、デフラグ部 14 は、各ブロックデータについて、ブロック管理表 16 を参照し、該当ブロックデータの重複率を確認する (ステップ S22)。

10

【0076】

そして、デフラグ部 14 は、該当ブロックデータの重複率が一定値未満である場合は (ステップ S23: Yes)、図 18 に示すように、一次デフラグ先領域 60 の未使用領域 62 の先頭にブロックデータを移動する (ステップ S24)。つまり、デフラグ先領域 60 の未使用範囲 62 のうち、デフラグ範囲 61 との隣接箇所にはブロックデータを書き込む (図 18 の網掛け部分参照)。

【0077】

20

その後、デフラグ部 14 は、図 19 に示すように、一次デフラグ先領域 60 において、書き込んだブロックデータのサイズ分だけデフラグ範囲 61 の末尾を拡大すると共に、当該書き込んだブロックデータのサイズ分だけ未使用範囲 62 の先頭を縮小する。換言すると、未使用範囲 62 の先頭側に位置するデフラグ範囲 61 との隣接箇所から連続してブロックデータが格納された範囲 (図 18 の網掛け部分) と、デフラグ範囲 (図 18 の符号 61) と、を連結した範囲を、新たなデフラグ範囲 (図 19 の符号 61) とする。これにより、一次デフラグ先領域 60 の新たなデフラグ範囲 61 は、後にデフラグ対象となる。そして、領域管理表 15 における一次デフラグ先領域 60 の最終書込時刻を更新する。

【0078】

また、デフラグ部 14 は、該当ブロックデータの重複率が一定値以上である場合は (ステップ S23: No)、図 18 に示すように、二次デフラグ先領域 70 の未使用領域 72 の末尾にブロックデータを移動する (ステップ S25)。つまり、二次デフラグ先領域 70 の未使用範囲 72 のうち、長寿命範囲 73 との隣接箇所にはブロックデータを書き込む (図 18 の網掛け部分参照)。

30

【0079】

その後、デフラグ部 14 は、図 19 に示すように、二次デフラグ先領域 70 において、書き込んだブロックデータのサイズ分だけ長寿命範囲 73 の先頭を拡大すると共に、当該書き込んだブロックデータのサイズ分だけ未使用範囲 72 の末尾を縮小する。換言すると、未使用範囲 72 の末尾側に位置する長寿命範囲 73 との隣接箇所から連続してブロックデータが格納された範囲 (図 18 の網掛け部分) と、長寿命範囲 (図 18 の符号 73) と、を連結した範囲を、新たな長寿命範囲 (図 19 の符号 73) とする。

40

【0080】

上述したデフラグ処理により、デフラグ元領域 40 のデフラグ範囲 41 内のブロックデータの重複率が「0」となった場合、ブロック管理表 16 の該当する行を削除することで、該当ブロックデータを破棄する。この時点で、一次デフラグ先領域 60 の未使用範囲 62 が空になった場合、デフラグを中断し、二次デフラグ先領域 70 の未使用範囲 72 が空になった場合、他の特性が超寿命となっている領域を選びなおし続行する。また、ブロック管理表 16 の該当するハッシュ値をもつ行について、領域番号と物理アドレスを、移動先を指し示すように更新する。これらの処理がすべて完了した時点で、デフラグ元領域 40 のデフラグ範囲 41 を未使用範囲 42 に統合し、当該デフラグ範囲 41 を「0-0」とす

50

る。

【0081】

デフラグ後、デフラグ部14あるいは書き込み部11は、一次デフラグ先領域60の使用フラグを「yes」に変更し、それ以外の領域21すべてについて使用中フラグを「no」にする。

【0082】

以上のように、重複率に応じてブロックデータを格納するデフラグ先として、「短寿命領域」と「長寿命範囲」といった専用領域を設定することで、デフラグ時におけるデフラグ先領域の確保が容易になり、デフラグ処理を実行する手段の実装が簡潔になる。

【0083】

ここで、上記記憶装置20として、データのアクセス速度、特に、読み出し速度が異なる2種類の装置を用いてもよい。例えば、記憶装置20は、図21に示すように、ハードディスクドライブといった一次記憶装置20Aと、当該一次記憶装置20Aよりもデータのアクセスが高速なSSD(Solid State Drive)といった二次記憶装置20Bと、を備える。そして、一次記憶装置20A内の領域21は、主に重複率の低いブロックデータを格納する上記一次デフラグ先領域60として使用し、二次記憶装置20B内の領域21は、主に重複率の高いブロックデータを格納する上記二次デフラグ先領域70として使用する。

【0084】

これにより、重複率の高いブロックデータは高頻度で参照されることとなるが、かかるブロックデータがデフラグ処理により集約される上記二次デフラグ先領域70を、読み出し速度が高速な二次記憶装置20Bに設けることで、読み込み性能を改善することができる。

【0085】

<付記>

上記実施形態の一部又は全部は、以下の付記のようにも記載されうる。以下、本発明におけるストレージシステム(図22参照)、プログラム、デフラグ方法の構成の概略を説明する。但し、本発明は、以下の構成に限定されない。

【0086】

(付記1)

記憶対象データを記憶装置120に格納すると共に、当該記憶装置に既に記憶されている前記記憶対象データと同一のデータ内容の他の記憶対象データを前記記憶装置に格納する場合に、当該記憶装置に既に記憶されている前記記憶対象データを前記他の記憶対象データとして参照させるデータ格納制御部111と、

前記記憶装置の所定の領域内においてデフラグ範囲とされた箇所に格納された前記記憶対象データを、前記記憶装置の他の領域内に格納し直すデフラグ処理部112と、を備え、

前記データ格納制御部111は、前記記憶装置に記憶されている前記記憶対象データ毎に、当該記憶対象データが他の記憶対象データとして参照されている数である参照数を記憶し、

前記デフラグ処理部112は、前記記憶対象データの参照数に応じて当該記憶対象データを前記記憶装置の他の領域内において後にデフラグ範囲となる箇所に格納する、ストレージシステム100。

【0087】

(付記2)

付記1に記載のストレージシステムであって、

前記デフラグ処理部は、前記参照数が予め設定された閾値未満である前記記憶対象データを、前記記憶装置の他の領域内において後にデフラグ範囲となる箇所に格納し、前記参照数が予め設定された閾値以上である前記記憶対象データを、前記記憶装置の他の領域内において後に非デフラグ範囲となる箇所に格納する、

10

20

30

40

50

ストレージシステム。

【0088】

(付記3)

付記2に記載のストレージシステムであって、

前記記憶装置の他の領域は、予め設定されたデフラグ範囲と、データが未格納な範囲である未使用範囲と、が隣接して形成されており、

前記デフラグ処理部は、前記参照数が閾値未満である前記記憶対象データを、前記記憶装置の他の領域における前記未使用範囲内であり前記デフラグ範囲との隣接箇所からデータが連続して位置する箇所に格納し、前記未使用範囲内における前記デフラグ範囲との隣接箇所から連続して前記記憶対象データが格納された範囲と前記デフラグ範囲とを連結した範囲を新たなデフラグ範囲とする、

10

ストレージシステム。

【0089】

(付記4)

付記3に記載のストレージシステムであって、

前記記憶装置の他の領域は、予め設定された非デフラグ範囲と、データが未格納な範囲である未使用範囲と、が隣接して形成されており、

前記デフラグ処理部は、前記参照数が閾値以上である前記記憶対象データを、前記記憶装置の他の領域における前記未使用範囲であり前記非デフラグ範囲との隣接箇所からデータが連続して位置する箇所に格納し、前記未使用範囲内における前記非デフラグ範囲との隣接箇所から連続して前記記憶対象データが格納された範囲と前記非デフラグ範囲とを連結した範囲を新たな非デフラグ範囲とする、

20

ストレージシステム。

【0090】

(付記5)

付記2に記載のストレージシステムであって、

前記記憶装置の他の領域は、予め設定されたデフラグ範囲にデータが未格納な範囲である未使用範囲の一端側が隣接すると共に、当該未使用範囲の他端側が予め設定された非デフラグ範囲に隣接し、当該未使用範囲がデフラグ範囲と非デフラグ範囲とに挟まれて形成されており、

30

前記デフラグ処理部は、前記参照数が閾値未満である前記記憶対象データを、前記記憶装置の他の領域における前記未使用範囲内であり前記デフラグ範囲との隣接箇所からデータが連続して位置する箇所に格納し、前記未使用範囲内における前記デフラグ範囲との隣接箇所から連続して前記記憶対象データが格納された範囲と前記デフラグ範囲とを連結した範囲を新たなデフラグ範囲とし、前記参照数が閾値以上である前記記憶対象データを、前記記憶装置の他の領域における前記未使用範囲であり前記非デフラグ範囲との隣接箇所からデータが連続して位置する箇所に格納し、前記未使用範囲内における前記非デフラグ範囲との隣接箇所から連続して前記記憶対象データが格納された範囲と前記非デフラグ範囲とを連結した範囲を新たな非デフラグ範囲とする、

ストレージシステム。

40

【0091】

(付記6)

付記1乃至5のいずれかに記載のストレージシステムであって、

前記記憶装置の所定の領域は、予め設定されたデフラグ範囲と、データが未格納な範囲である未使用範囲と、が隣接して形成されており、

前記データ格納制御部は、前記記憶対象データを新たに前記記憶装置に格納する際に、前記記憶装置の所定の領域における前記未使用範囲内であり前記デフラグ範囲との隣接箇所からデータが連続して位置する箇所に格納し、前記未使用範囲内における前記デフラグ範囲との隣接箇所から連続して前記記憶対象データが格納された範囲と前記デフラグ範囲とを連結した範囲を新たなデフラグ範囲とする、

50

ストレージシステム。

【 0 0 9 2 】

(付記 7)

付記 4 に記載のストレージシステムであって、

前記記憶装置の他の領域は、前記デフラグ範囲と前記未使用範囲とが隣接して形成された第一領域と、当該第一領域とは異なり前記非デフラグ範囲と前記未使用範囲とが隣接して形成された第二領域と、を有する、

ストレージシステム。

【 0 0 9 3 】

(付記 8)

付記 7 に記載のストレージシステムであって、

前記第一領域は、所定の記憶装置に形成されており、

前記第二領域は、前記所定の記憶装置よりもデータ読み出し処理が高速な他の記憶装置に形成されている、

ストレージシステム。

【 0 0 9 4 】

(付記 9)

情報処理装置に、

記憶対象データを記憶装置に格納すると共に、当該記憶装置に既に記憶されている前記記憶対象データと同一のデータ内容の他の記憶対象データを前記記憶装置に格納する場合に、当該記憶装置に既に記憶されている前記記憶対象データを前記他の記憶対象データとして参照させるデータ格納制御部と、

前記記憶装置の所定の領域内においてデフラグ範囲とされた箇所に格納された前記記憶対象データを、前記記憶装置の他の領域内に格納し直すデフラグ処理部と、を実現させると共に、

前記データ格納制御部は、前記記憶装置に記憶されている前記記憶対象データ毎に、当該記憶対象データが他の記憶対象データとして参照されている数である参照数を記憶し、

前記デフラグ処理部は、前記記憶対象データの参照数に応じて当該記憶対象データを前記記憶装置の他の領域内において後にデフラグ範囲となる箇所に格納する、ことを実現させるためのプログラム。

【 0 0 9 5 】

(付記 9 - 2)

付記 9 に記載のプログラムであって、

前記デフラグ処理部は、前記参照数が予め設定された閾値未満である前記記憶対象データを、前記記憶装置の他の領域内において後にデフラグ範囲となる箇所に格納し、前記参照数が予め設定された閾値以上である前記記憶対象データを、前記記憶装置の他の領域内において後に非デフラグ範囲となる箇所に格納する、

プログラム。

【 0 0 9 6 】

(付記 1 0)

記憶対象データを記憶装置に格納すると共に、当該記憶装置に既に記憶されている前記記憶対象データと同一のデータ内容の他の記憶対象データを前記記憶装置に格納する場合に、当該記憶装置に既に記憶されている前記記憶対象データを前記他の記憶対象データとして参照させてデータ格納制御を行うと共に、前記記憶装置に記憶されている前記記憶対象データ毎に、当該記憶対象データが他の記憶対象データとして参照されている数である参照数を記憶し、

前記記憶装置の所定の領域内においてデフラグ範囲とされた箇所に格納された前記記憶対象データを、前記記憶装置の他の領域内に格納し直すデフラグ処理を実行し、

前記デフラグ処理時に、前記記憶対象データの参照数に応じて当該記憶対象データを前記記憶装置の他の領域内において後にデフラグ範囲となる箇所に格納する、

10

20

30

40

50

デフラグ方法。

【 0 0 9 7 】

(付記 1 0 - 2)

付記 1 0 に記載のデフラグ方法であって、

前記デフラグ処理時に、前記参照数が予め設定された閾値未満である前記記憶対象データを、前記記憶装置の他の領域内において後にデフラグ範囲となる箇所に格納し、前記参照数が予め設定された閾値以上である前記記憶対象データを、前記記憶装置の他の領域内において後に非デフラグ範囲となる箇所に格納する、

デフラグ方法。

【 0 0 9 8 】

なお、上述したプログラムは、記憶装置に記憶されていたり、コンピュータが読み取り可能な記録媒体に記録されている。例えば、記録媒体は、フレキシブルディスク、光ディスク、光磁気ディスク、及び、半導体メモリ等の可搬性を有する媒体である。

【 0 0 9 9 】

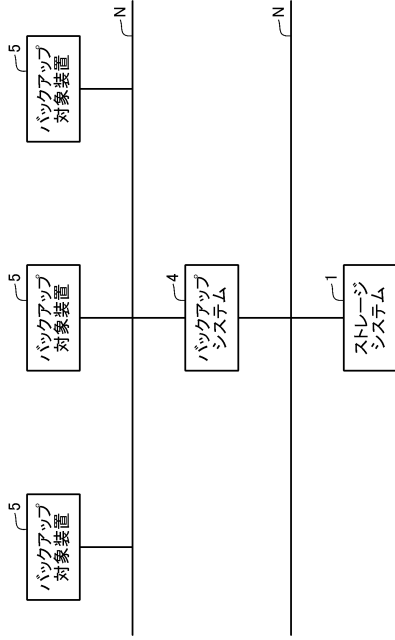
以上、上記実施形態等を参照して本願発明を説明したが、本願発明は、上述した実施形態に限定されるものではない。本願発明の構成や詳細には、本願発明の範囲内で当業者が理解しうる様々な変更をすることができる。

【 符号の説明 】

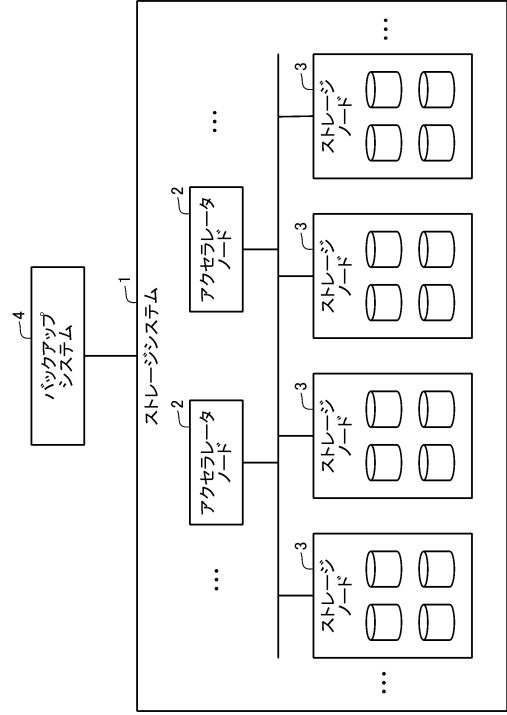
【 0 1 0 0 】

- | | | |
|-----------------------------|------------|----|
| 1 | ストレージシステム | 20 |
| 2 | アクセラレータノード | |
| 3 | ストレージノード | |
| 4 | バックアップシステム | |
| 5 | バックアップ対象装置 | |
| 1 1 | 書き込み部 | |
| 1 2 | 読み出し部 | |
| 1 3 | 削除部 | |
| 1 4 | デフラグ部 | |
| 1 5 | 領域管理表 | |
| 1 6 | ブロック管理表 | 30 |
| 1 7 | ファイル管理表 | |
| 2 0 | 記憶装置 | |
| 2 0 A | 一次記憶装置 | |
| 2 0 B | 二次記憶装置 | |
| 2 1 | 領域 | |
| 3 0 | 書き込み領域 | |
| 4 0 | デフラグ元領域 | |
| 5 0 | デフラグ先領域 | |
| 6 0 | 一次デフラグ領域 | |
| 7 0 | 二次デフラグ領域 | 40 |
| 3 1 , 4 1 , 5 1 , 6 1 , 7 1 | デフラグ範囲 | |
| 3 2 , 4 2 , 5 2 , 6 2 , 7 2 | 未使用範囲 | |
| 3 3 , 4 3 , 5 3 , 6 3 , 7 3 | 長寿命範囲 | |
| 1 0 0 | ストレージシステム | |
| 1 1 1 | データ格納制御部 | |
| 1 1 2 | デフラグ処理部 | |
| 1 2 0 | 記憶装置 | |

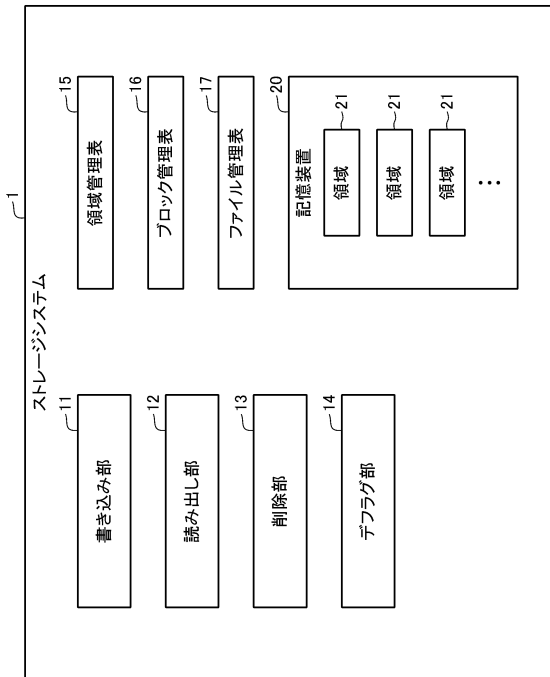
【図1】



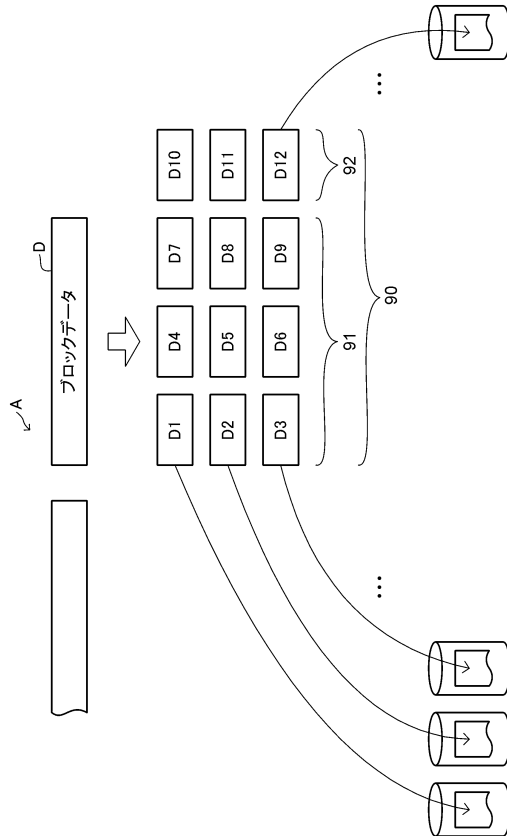
【図2】



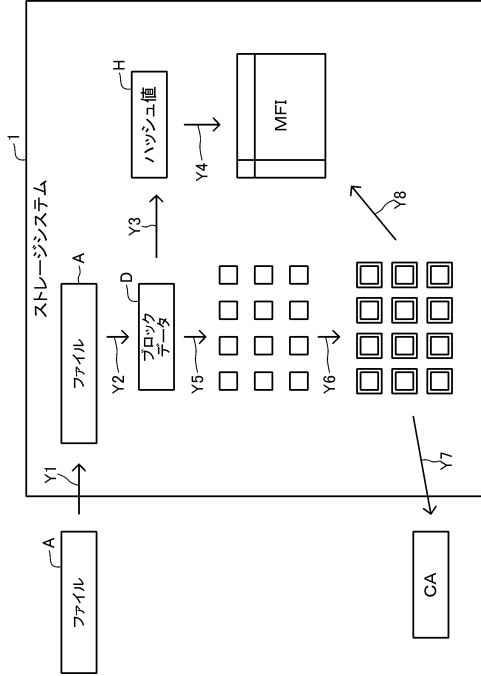
【図3】



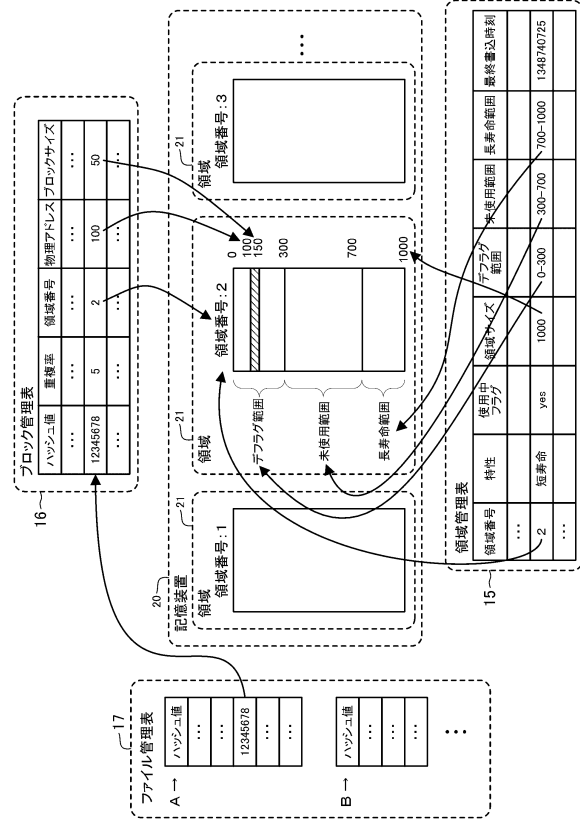
【図4】



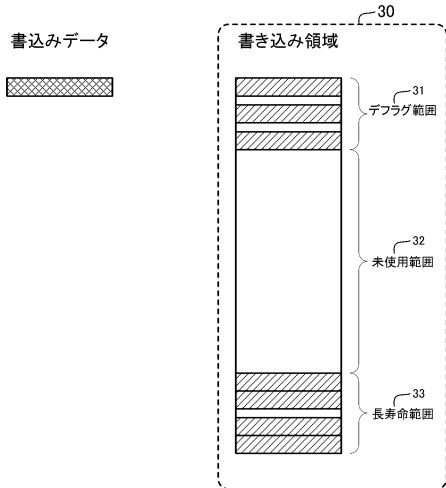
【図5】



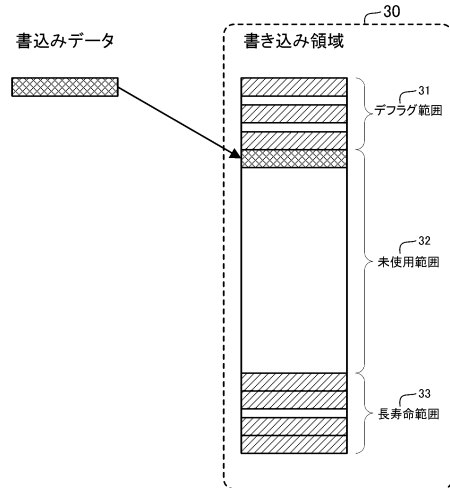
【図6】



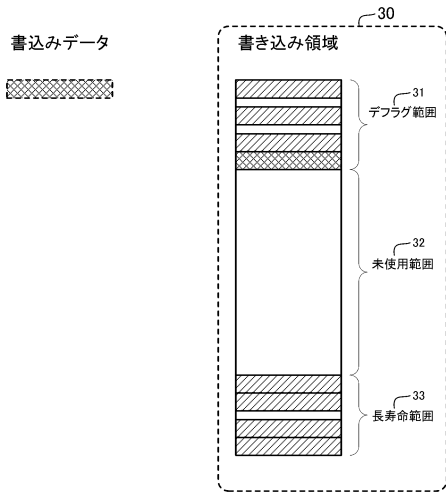
【図7】



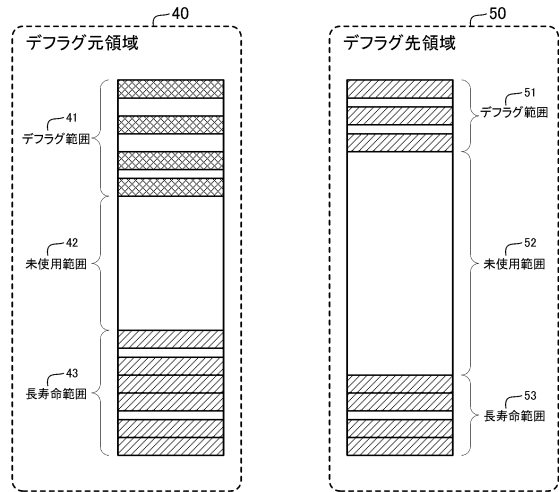
【図8】



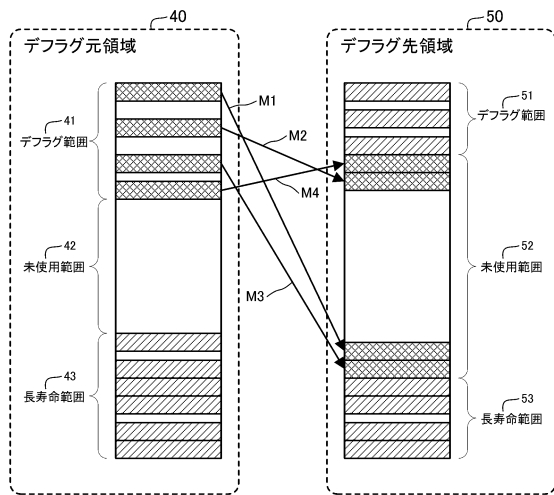
【図9】



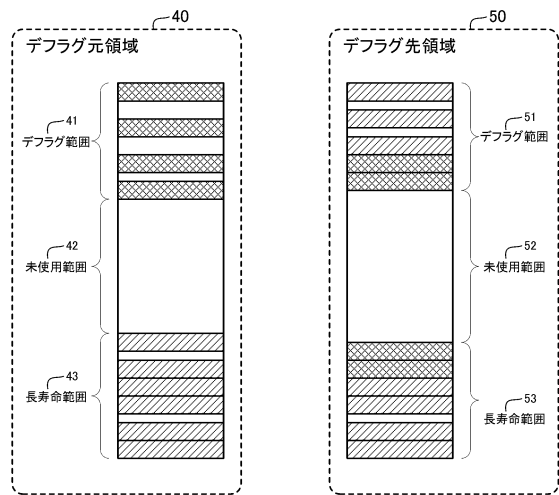
【図10】



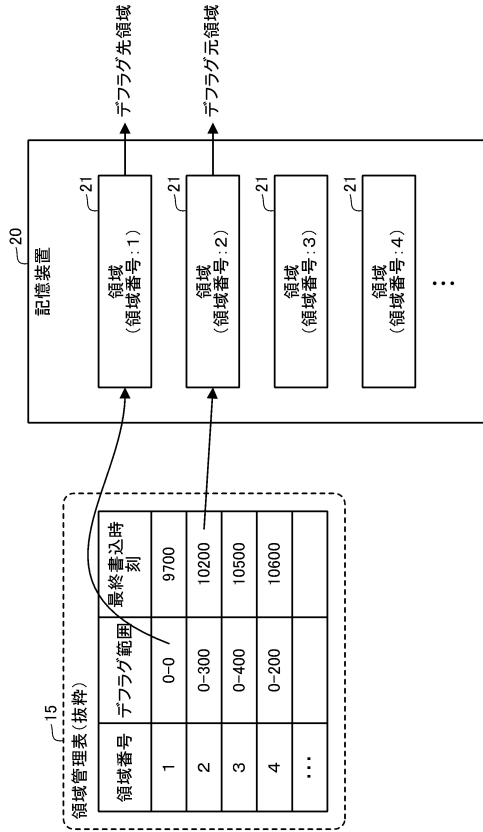
【図11】



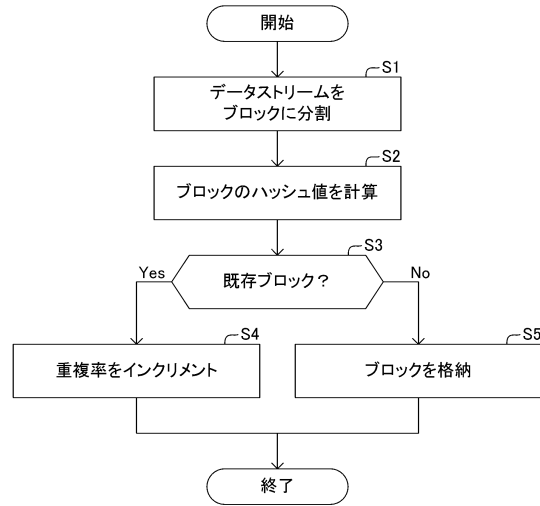
【図12】



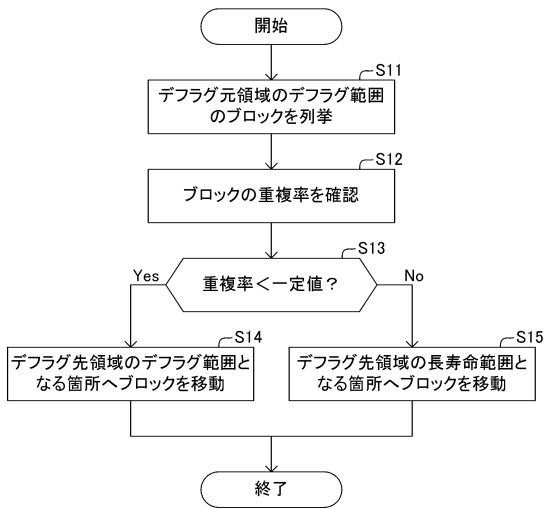
【図13】



【図14】



【図15】

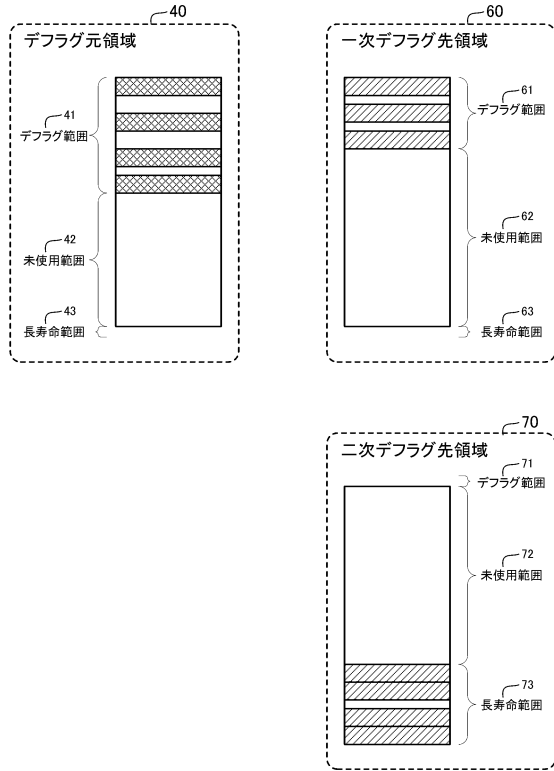


【図16】

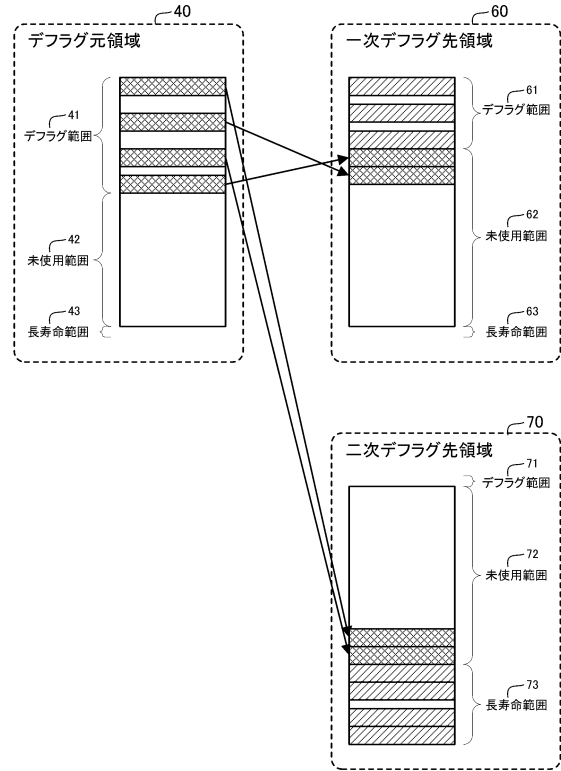
領域管理表 (抜粋) (15)

領域番号	特性	使用中フラグ	領域サイズ	デフラグ範囲	未使用範囲	長寿命範囲	最終書込時刻
...							
2	短寿命	yes	1000	0-300	300-700	700-1000	1348740725
...							
5	長寿命	-	1000	0-0	0-600	600-1000	-
...							

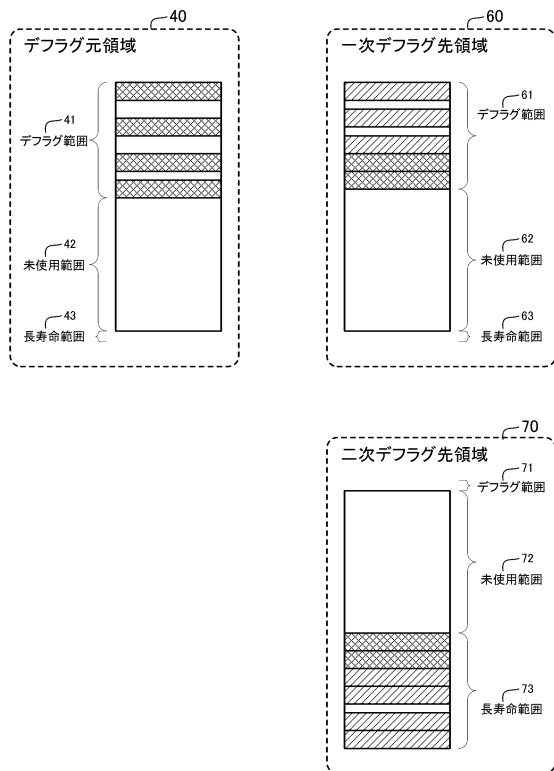
【図17】



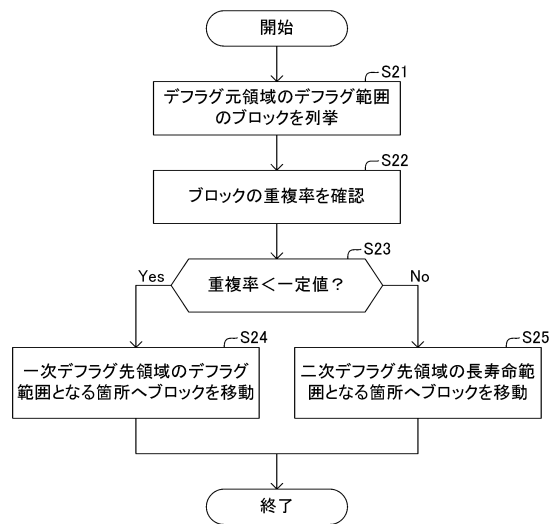
【図18】



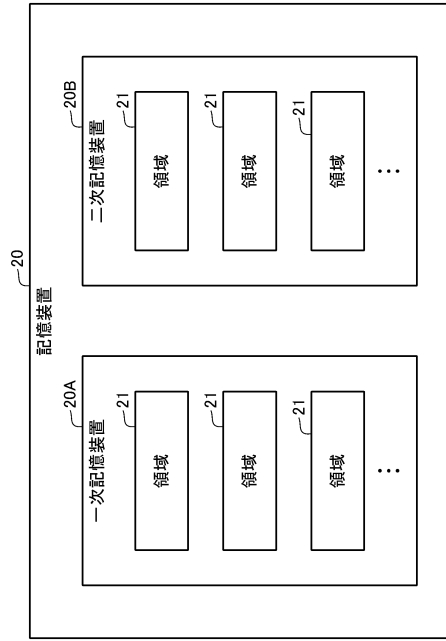
【図19】



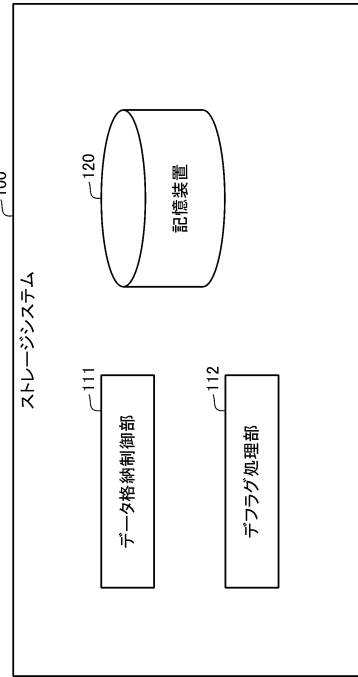
【図20】



【図 2 1】



【図 2 2】



フロントページの続き

- (56)参考文献 特開2011-170665(JP,A)
特開2009-205201(JP,A)
特開2010-218194(JP,A)
米国特許出願公開第2013/0226881(US,A1)

- (58)調査した分野(Int.Cl., DB名)
G06F 12/00
G06F 3/06